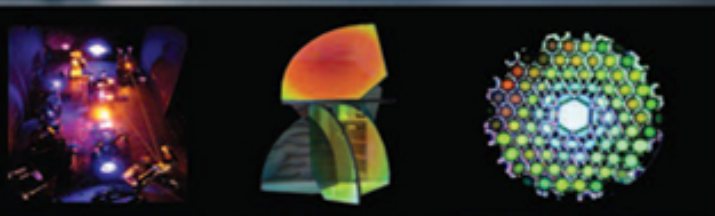


Third Edition

Sponsored by the Optical Society of America

HANDBOOK OF OPTICS

Volume I-V



Editor-in-Chief:
Michael Bass

Associate Editors:
Casimer M. DeCusatis
Jay M. Enoch
Vasudevan Lakshminarayanan
Guifang Li
Carolyn MacDonald
Virendra N. Mahajan
Eric Van Stryland

OSA[®]

HANDBOOK OF OPTICS

DO NOT DUPLICATE

ABOUT THE EDITORS

Editor-in-Chief: Dr. Michael Bass is professor emeritus at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Associate Editors:

Dr. Casimer M. DeCusatis is a distinguished engineer and technical executive with IBM Corporation.

Dr. Jay M. Enoch is dean emeritus and professor at the School of Optometry at the University of California, Berkeley.

Dr. Vasudevan Lakshminarayanan is professor of Optometry, Physics, and Electrical Engineering at the University of Waterloo, Ontario, Canada.

Dr. Guifang Li is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Dr. Carolyn MacDonald is a professor at the University at Albany, and director of the Center for X-Ray Optics.

Dr. Virendra N. Mahajan is a distinguished scientist at The Aerospace Corporation.

Dr. Eric Van Stryland is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

HANDBOOK OF OPTICS

THIRD EDITION

Sponsored by the
OPTICAL SOCIETY OF AMERICA

Michael Bass Editor-in-Chief

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*



New York Chicago San Francisco Lisbon London Madrid
Mexico City Milan New Delhi San Juan Seoul
Singapore Sydney Toronto

Copyright © 2010 by The McGraw-Hill Companies, Inc. (with the exception of Chapters I:14 and I:15, copyright © Russell A. Chipman). All rights reserved. Except as permitted under the United States Copyright Act of 1976, no part of this publication may be reproduced or distributed in any form or by any means, or stored in a database or retrieval system, without the prior written permission of the publisher.

ISBN 978-0-07-175342-5, MHID 0-07-175342-7

The material in this eBook also appears in the print version of this title. ISBN: 978-0-07-170160-0, MHID: 0-07-170160-5.

All trademarks are trademarks of their respective owners. Rather than put a trademark symbol after every occurrence of a trademarked name, we use names in an editorial fashion only, and to the benefit of the trademark owner, with no intention of infringement of the trademark. Where such designations appear in this book, they have been printed with initial caps.

McGraw-Hill eBooks are available at special quantity discounts to use as premiums and sales promotions, or for use in corporate training programs. To contact a representative please e-mail us at bulksales@mcgraw-hill.com.

TERMS OF USE

This is a copyrighted work and The McGraw-Hill Companies, Inc. (“McGraw-Hill”) and its licensors reserve all rights in and to the work. Use of this work is subject to these terms. Except as permitted under the Copyright Act of 1976 and the right to store and retrieve one copy of the work, you may not decompile, disassemble, reverse engineer, reproduce, modify, create derivative works based upon, transmit, distribute, disseminate, sell, publish or sublicense the work or any part of it without McGraw-Hill’s prior consent. You may use the work for your own noncommercial and personal use; any other use of the work is strictly prohibited. Your right to use the work may be terminated if you fail to comply with these terms.

THE WORK IS PROVIDED “AS IS.” MCGRAW-HILL AND ITS LICENSORS MAKE NO GUARANTEES OR WARRANTIES AS TO THE ACCURACY, ADEQUACY OR COMPLETENESS OF OR RESULTS TO BE OBTAINED FROM USING THE WORK, INCLUDING ANY INFORMATION THAT CAN BE ACCESSED THROUGH THE WORK VIA HYPERLINK OR OTHERWISE, AND EXPRESSLY DISCLAIM ANY WARRANTY, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. McGraw-Hill and its licensors do not warrant or guarantee that the functions contained in the work will meet your requirements or that its operation will be uninterrupted or error free. Neither McGraw-Hill nor its licensors shall be liable to you or anyone else for any inaccuracy, error or omission, regardless of cause, in the work or for any damages resulting therefrom. McGraw-Hill has no responsibility for the content of any information accessed through the work. Under no circumstances shall McGraw-Hill and/or its licensors be liable for any indirect, incidental, special, punitive, consequential or similar damages that result from the use of or inability to use the work, even if any of them has been advised of the possibility of such damages. This limitation of liability shall apply to any claim or cause whatsoever whether such claim or cause arises in contract, tort or otherwise.

Information contained in this work has been obtained by The McGraw-Hill Companies, Inc. (“McGraw-Hill”) from sources believed to be reliable. However, neither McGraw-Hill nor its authors guarantee the accuracy or completeness of any information published herein, and neither McGraw-Hill nor its authors shall be responsible for any errors, omissions, or damages arising out of use of this information. This work is published with the understanding that McGraw-Hill and its authors are supplying information but are not attempting to render engineering or other professional services. If such services are required, the assistance of an appropriate professional should be sought.

COVER ILLUSTRATIONS

Left: Poincaré sphere describing light's polarization states is shown floating in front of a depolarized field of polarization ellipses, with linearly and circularly polarized fields propagating on its left and right, respectively. See Chaps. 12 and 15.

Middle: Triplet lens developed for photographic applications that can zero out the primary aberrations by splitting the positive lens of a doublet into two and placing one on each side of the negative lens. See Chap. 17.

Right: Micrographs of different optical storage media showing the straight and narrow tracks with 1.6- μm spacing between adjacent tracks. The recorded information bits appear as short marks along each track. See Chap. 35.

This page intentionally left blank.

DO NOT DUPLICATE

CONTENTS OF VOLUME I

Part 1. Geometrical Optics

Chapter 1. General Principles of Geometrical Optics <i>Douglas S. Goodman</i>	1.3
<hr/>	
1.1 Glossary / 1.3	
1.2 Introduction / 1.7	
1.3 Fundamentals / 1.8	
1.4 Characteristic Functions / 1.13	
1.5 Rays in Heterogeneous Media / 1.18	
1.6 Conservation of Étendue / 1.22	
1.7 Skew Invariant / 1.23	
1.8 Refraction and Reflection at Interfaces between Homogeneous Media / 1.23	
1.9 Imaging / 1.26	
1.10 Description of Systems of Revolution / 1.32	
1.11 Tracing Rays in Centered Systems of Spherical Surfaces / 1.35	
1.12 Paraxial Optics of Systems of Revolution / 1.37	
1.13 Images About Known Rays / 1.43	
1.14 Gaussian Lens Properties / 1.44	
1.15 Collineation / 1.56	
1.16 System Combinations: Gaussian Properties / 1.63	
1.17 Paraxial Matrix Methods / 1.65	
1.18 Apertures, Pupils, Stops, Fields, and Related Matters / 1.74	
1.19 Geometrical Aberrations of Point Images: Description / 1.85	
1.20 References / 1.92	

Part 2. Physical Optics

Chapter 2. Interference <i>John E. Greivenkamp</i>	2.3
<hr/>	
2.1 Glossary / 2.3	
2.2 Introduction / 2.3	
2.3 Waves and Wavefronts / 2.3	
2.4 Interference / 2.5	
2.5 Interference by Wavefront Division / 2.14	
2.6 Interference by Amplitude Division / 2.19	
2.7 Multiple Beam Interference / 2.28	
2.8 Coherence and Interference / 2.36	
2.9 Applications of Interference / 2.42	
2.10 References / 2.42	
Chapter 3. Diffraction <i>Arvind S. Marathay and John F. McCalmont</i>	3.1
<hr/>	
3.1 Glossary / 3.1	
3.2 Introduction / 3.1	
3.3 Light Waves / 3.2	
3.4 Huygens-Fresnel Construction / 3.4	
3.5 Cylindrical Wavefront / 3.13	

- 3.6 Mathematical Theory of Diffraction / 3.21
- 3.7 Stationary Phase Approximation / 3.29
- 3.8 Vector Diffraction / 3.32
- 3.9 Acknowledgments / 3.38
- 3.10 References / 3.38

Chapter 4. Transfer Function Techniques *Glenn D. Boreman* **4.1**

- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.1
- 4.3 Definitions / 4.2
- 4.4 MTF Calculations / 4.3
- 4.5 MTF Measurements / 4.6
- 4.6 References / 4.8

Chapter 5. Coherence Theory *William H. Carter* **5.1**

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.1
- 5.3 Some Elementary Classical Concepts / 5.2
- 5.4 Definitions of Coherence Functions / 5.4
- 5.5 Model Sources / 5.9
- 5.6 Propagation / 5.13
- 5.7 Spectrum of Light / 5.19
- 5.8 Polarization Effects / 5.22
- 5.9 Applications / 5.22
- 5.10 References / 5.23
- 5.11 Additional Reading / 5.26

Chapter 6. Coherence Theory: Tools and Applications
*Gisele Bennett, William T. Rhodes,
and J. Christopher James* **6.1**

- 6.1 Glossary / 6.1
- 6.2 Introduction / 6.2
- 6.3 Key Definitions and Relationships / 6.2
- 6.4 Propagation, Diffraction, and Scattering: Enhanced Backscatter and the Lau Effect / 6.5
- 6.5 Image Formation: Lukosz-Type Super-Resolving System / 6.9
- 6.6 Efficient Sampling of Coherence Functions / 6.10
- 6.7 An Example of When Not to Use Coherence Theory / 6.12
- 6.8 Concluding Remarks / 6.13
- 6.9 References / 6.13

Chapter 7. Scattering by Particles *Craig F. Bohren* **7.1**

- 7.1 Glossary / 7.1
- 7.2 Introduction / 7.2
- 7.3 Scattering: An Overview / 7.3
- 7.4 Scattering by Particles: Basic Concepts and Terminology / 7.4
- 7.5 Scattering by an Isotropic, Homogeneous Sphere: The Archetype / 7.11
- 7.6 Scattering by Regular Particles / 7.14
- 7.7 Computational Methods for Nonspherical Particles / 7.15
- 7.8 References / 7.17

Chapter 8. Surface Scattering *Eugene L. Church
and Peter Z. Takacs* **8.1**

- 8.1 Glossary of Principal Symbols / 8.1
- 8.2 Introduction / 8.2

- 8.3 Notation / 8.2
- 8.4 The Fresnel-Kirchhoff Approximation / 8.5
- 8.5 The Rayleigh-Rice (RR) or Small-Perturbation Approximation / 8.9
- 8.6 Effects of Finite Illumination Area / 8.12
- 8.7 Surface Statistics / 8.12
- 8.8 Surface Finish Specification / 8.16
- 8.9 Retrospect and Prospect / 8.17
- 8.10 References and Endnotes / 8.18

Chapter 9. Volume Scattering in Random Media 9.1
Aristide Dogariu and Jeremy Ellis

- 9.1 Glossary / 9.1
- 9.2 Introduction / 9.2
- 9.3 General Theory of Scattering / 9.3
- 9.4 Single Scattering / 9.4
- 9.5 Multiple Scattering / 9.8
- 9.6 References / 9.18

Chapter 10. Optical Spectroscopy and Spectroscopic Lineshapes 10.1
Brian Henderson

- 10.1 Glossary / 10.1
- 10.2 Introductory Comments / 10.2
- 10.3 Theoretical Preliminaries / 10.3
- 10.4 Rates of Spectroscopic Transition / 10.4
- 10.5 Lineshapes of Spectral Transitions / 10.6
- 10.6 Spectroscopy of One-Electron Atoms / 10.7
- 10.7 Multielectron Atoms / 10.10
- 10.8 Optical Spectra and the Outer Electronic Structure / 10.12
- 10.9 Spectra of Tri-Positive Rare Earth Ions / 10.16
- 10.10 Vibrational and Rotational Spectra of Molecules / 10.18
- 10.11 Lineshapes in Solid State Spectroscopy / 10.22
- 10.12 References / 10.27

Chapter 11. Analog Optical Signal and Image Processing 11.1
Joseph W. Goodman

- 11.1 Glossary / 11.1
- 11.2 Introduction / 11.1
- 11.3 Fundamental Analog Operations / 11.2
- 11.4 Analog Optical Fourier Transforms / 11.3
- 11.5 Spatial Filtering / 11.5
- 11.6 Coherent Optical Processing of Synthetic Aperture Radar Data / 11.6
- 11.7 Coherent Optical Processing of Temporal Signals / 11.8
- 11.8 Optical Processing of Two-Dimensional Images / 11.12
- 11.9 Incoherent Processing of Discrete Signals / 11.17
- 11.10 Concluding Remarks / 11.20
- 11.11 References / 11.20

Part 3. Polarized Light

Chapter 12. Polarization *Jean M. Bennett* 12.3

- 12.1 Glossary / 12.3
- 12.2 Basic Concepts and Conventions / 12.4
- 12.3 Fresnel Equations / 12.6

- 12.4 Basic Relations for Polarizers / 12.14
- 12.5 Polarization by Nonnormal-Incidence Reflection (Pile of Plates) / 12.15
- 12.6 Polarization by Nonnormal-Incidence Transmission (Pile of Plates) / 12.18
- 12.7 Quarter-Wave Plates and Other Phase Retardation Plates / 12.24
- 12.8 Matrix Methods for Computing Polarization / 12.27
- 12.9 References / 12.30

Chapter 13. Polarizers *Jean M. Bennett* 13.1

- 13.1 Glossary / 13.1
- 13.2 Prism Polarizers / 13.2
- 13.3 Glan-Type Prisms / 13.8
- 13.4 Nicol-Type Prisms / 13.15
- 13.5 Polarizing Beam-Splitter Prisms / 13.18
- 13.6 Feussner Prisms / 13.22
- 13.7 Noncalcite Polarizing Prisms / 13.23
- 13.8 Dichroic and Diffraction-Type Polarizers / 13.24
- 13.9 Non-Normal-Incidence Reflection and Transmission Polarizers / 13.33
- 13.10 Retardation Plates / 13.43
- 13.11 Variable Retardation Plates and Compensators / 13.53
- 13.12 Half-Shade Devices / 13.56
- 13.13 Miniature Polarization Devices / 13.57
- 13.14 References / 13.58

Chapter 14. Mueller Matrices *Russell A. Chipman* 14.1

- 14.1 Glossary / 14.1
- 14.2 Conventions / 14.3
- 14.3 Objectives / 14.3
- 14.4 Stokes Parameters and Mueller Matrices / 14.4
- 14.5 The Stokes Parameters and the Poincaré Sphere / 14.4
- 14.6 Mueller Matrices / 14.6
- 14.7 Sequences of Polarization Elements / 14.7
- 14.8 Polarization Elements' Properties in the Mueller Calculus / 14.7
- 14.9 Rotation of an Element About the Optical Axis / 14.8
- 14.10 Nonpolarizing Mueller Matrices / 14.8
- 14.11 Mueller Matrices of Ideal Polarizers / 14.8
- 14.12 Retarder Mueller Matrices / 14.11
- 14.13 Retarder Mueller Matrices Ambiguities and Retarder Space / 14.14
- 14.14 Transmittance and Diattenuation / 14.16
- 14.15 Polarizance / 14.18
- 14.16 Mueller Matrices of Diattenuators / 14.18
- 14.17 Normalizing a Mueller Matrix / 14.19
- 14.18 Coordinate System for the Mueller Matrix / 14.19
- 14.19 Mueller Matrices for Refraction / 14.20
- 14.20 Mueller Matrices for Reflection / 14.21
- 14.21 Conversion between Mueller Matrices and Jones Matrices / 14.22
- 14.22 Nondepolarizing Mueller Matrices and Mueller-Jones Matrices / 14.24
- 14.23 Homogeneous and Inhomogeneous Polarization Elements / 14.25
- 14.24 Mueller Matrices Near the Identity Matrix, Weak Polarization Elements / 14.26
- 14.25 Matrix Roots of Nondepolarizing Mueller Matrices / 14.27
- 14.26 Depolarization and the Depolarization Index / 14.30
- 14.27 Degree of Polarization Surfaces and Maps / 14.31
- 14.28 The Depolarization Index / 14.32
- 14.29 The Average Degree of Polarization / 14.32
- 14.30 Determining Mueller Matrix Properties / 14.33
- 14.31 Generators for Depolarization / 14.33
- 14.32 Interpretation of Arbitrary Mueller Matrices, the Polar Decomposition of Mueller Matrices / 14.39
- 14.33 Physically Realizable Mueller Matrices / 14.40
- 14.34 Acknowledgments / 14.42

14.35 References / 14.43

Chapter 15. Polarimetry *Russell A. Chipman*

15.1

- 15.1 Glossary / 15.1
- 15.2 Objectives / 15.3
- 15.3 Polarimeters / 15.3
- 15.4 Light-Measuring Polarimeters / 15.3
- 15.5 Sample-Measuring Polarimeters / 15.4
- 15.6 Complete and Incomplete Polarimeters / 15.4
- 15.7 Polarization Generators and Analyzers / 15.4
- 15.8 Classes of Polarimeters / 15.5
- 15.9 Time-Sequential Measurements / 15.5
- 15.10 Polarization Modulation / 15.5
- 15.11 Division of Aperture / 15.5
- 15.12 Division of Amplitude / 15.5
- 15.13 Spectropolarimeters / 15.6
- 15.14 Imaging Polarimeters / 15.6
- 15.15 Definitions / 15.6
- 15.16 Stokes Vectors and Mueller Matrices / 15.8
- 15.17 Phenomenological Definition of the Stokes Vector / 15.9
- 15.18 Polarization Properties of Light Beams / 15.9
- 15.19 Mueller Matrices / 15.11
- 15.20 Data Reduction for Light-Measuring Polarimeters / 15.11
- 15.21 Sample-Measuring Polarimeters for Measuring Mueller Matrix Elements / 15.13
- 15.22 Polarimetric Measurement Equation and Polarimetric Data-Reduction Equation / 15.14
- 15.23 Dual Rotating Retarder Polarimeter / 15.16
- 15.24 Incomplete Sample-Measuring Polarimeters / 15.16
- 15.25 Nonideal Polarization Elements / 15.17
- 15.26 Elliptical and Circular Polarizers and Analyzers / 15.17
- 15.27 Common Defects of Polarization Elements / 15.19
- 15.28 Polarization Modulators, Retardance Modulators / 15.20
- 15.29 Rotating Retarders / 15.20
- 15.30 Photo-Elastic Modulators / 15.21
- 15.31 Liquid Crystal Retarders / 15.21
- 15.32 Electro-Optical Modulators / 15.23
- 15.33 Magneto-Optical Modulators / 15.23
- 15.34 Fiber Squeezers / 15.24
- 15.35 Polarimeter Design Metrics / 15.24
- 15.36 Singular Value Decomposition Examples / 15.26
- 15.37 Polarimeter Error Analysis / 15.27
- 15.38 The Mueller Matrix for Polarization Component Characterization / 15.28
- 15.39 Retro-Reflection Testing and Correction for Supplemental Optics / 15.28
- 15.40 Applications of Polarimetry / 15.29
- 15.41 Ellipsometry and Generalized Ellipsometry / 15.30
- 15.42 Liquid Crystal Cell and System Testing / 15.32
- 15.43 Polarization Aberrations / 15.35
- 15.44 Remote Sensing / 15.37
- 15.45 Polarization Light Scattering / 15.38
- 15.46 Ophthalmic Polarimetry / 15.39
- 15.47 Acknowledgments / 15.41
- 15.48 References / 15.41

Chapter 16. Ellipsometry *Rasheed M. A. Azzam*

16.1

- 16.1 Glossary / 16.1
- 16.2 Introduction / 16.2
- 16.3 Conventions / 16.3
- 16.4 Modeling and Inversion / 16.4
- 16.5 Transmission Ellipsometry / 16.10

- 16.6 Instrumentation / 16.10
- 16.7 Jones-Matrix Generalized Ellipsometry / 16.19
- 16.8 Mueller-Matrix Generalized Ellipsometry / 16.19
- 16.9 Applications / 16.21
- 16.10 References / 16.21

Part 4. Components

Chapter 17. Lenses *R. Barry Johnson* 17.3

- 17.1 Glossary / 17.3
- 17.2 Introduction / 17.4
- 17.3 Basics / 17.5
- 17.4 Stops and Pupils / 17.8
- 17.5 F-Number and Numerical Aperture / 17.9
- 17.6 Magnifier or Eye Loupe / 17.9
- 17.7 Compound Microscopes / 17.10
- 17.8 Field and Relay Lenses / 17.10
- 17.9 Aplanatic Surfaces and Immersion Lenses / 17.10
- 17.10 Single Element Lens / 17.12
- 17.11 Landscape Lenses and the Influence of Stop Position / 17.17
- 17.12 Two-Lens Systems / 17.20
- 17.13 Achromatic Doublets / 17.22
- 17.14 Triplet Lenses / 17.26
- 17.15 Symmetrical Lenses / 17.26
- 17.16 Double-Gauss Lenses / 17.27
- 17.17 Petzval Lenses / 17.28
- 17.18 Telephoto Lenses / 17.29
- 17.19 Inverted or Reverse Telephoto Lenses / 17.29
- 17.20 Performance of Representative Lenses / 17.29
- 17.21 Rapid Estimation of Lens Performance / 17.36
- 17.22 Bibliography / 17.40

Chapter 18. Afocal Systems *William B. Wetherell* 18.1

- 18.1 Glossary / 18.1
- 18.2 Introduction / 18.2
- 18.3 Gaussian Analysis of Afocal Lenses / 18.2
- 18.4 Keplerian Afocal Lenses / 18.7
- 18.5 Galilean and Inverse Galilean Afocal Lenses / 18.15
- 18.6 Relay Trains and Periscopes / 18.17
- 18.7 Reflecting and Catadioptric Afocal Lenses / 18.19
- 18.8 References / 18.23

Chapter 19. Nondispersive Prisms *William L. Wolfe* 19.1

- 19.1 Glossary / 19.1
- 19.2 Introduction / 19.1
- 19.3 Inversion, Reversion / 19.2
- 19.4 Deviation, Displacement / 19.2
- 19.5 Summary of Prism Properties / 19.2
- 19.6 Prism Descriptions / 19.2
- 19.7 References / 19.29

Chapter 20. Dispersive Prisms and Gratings *George J. Zissis* 20.1

- 20.1 Glossary / 20.1
- 20.2 Introduction / 20.1

- 20.3 Prisms / 20.2
- 20.4 Gratings / 20.3
- 20.5 Prism and Grating Configurations and Instruments / 20.4
- 20.6 References / 20.15

Chapter 21. Integrated Optics *Thomas L. Koch,
Frederick J. Leonberger, and Paul G. Suchoski* **21.1**

- 21.1 Glossary / 21.1
- 21.2 Introduction / 21.2
- 21.3 Device Physics / 21.3
- 21.4 Integrated Optics Materials and Fabrication Technology / 21.13
- 21.5 Circuit Elements / 21.21
- 21.6 Applications of Integrated Optics / 21.31
- 21.7 Future Trends / 21.39
- 21.8 References / 21.41

Chapter 22. Miniature and Micro-Optics *Tom D. Milster and Tomasz S. Tkaczyk* **22.1**

- 22.1 Glossary / 22.1
- 22.2 Introduction / 22.2
- 22.3 Uses of Micro-Optics / 22.2
- 22.4 Micro-Optics Design Considerations / 22.2
- 22.5 Molded Microlenses / 22.8
- 22.6 Diamond Turning / 22.15
- 22.7 Lithography for Making Refractive Components / 22.18
- 22.8 Monolithic Lenslet Modules / 22.25
- 22.9 Distributed-Index Planar Microlenses / 22.26
- 22.10 Micro-Fresnel Lenses / 22.31
- 22.11 Liquid Lenses / 22.37
- 22.12 Other Technologies / 22.42
- 22.13 References / 22.47

Chapter 23. Binary Optics *Michael W. Farn
and Wilfrid B. Veldkamp* **23.1**

- 23.1 Glossary / 23.1
- 23.2 Introduction / 23.2
- 23.3 Design—Geometrical Optics / 23.2
- 23.4 Design—Scalar Diffraction Theory / 23.10
- 23.5 Design—Vector Diffraction Theory / 23.13
- 23.6 Fabrication / 23.14
- 23.7 References / 23.17

Chapter 24. Gradient Index Optics *Duncan T. Moore* **24.1**

- 24.1 Glossary / 24.1
- 24.2 Introduction / 24.1
- 24.3 Analytic Solutions / 24.2
- 24.4 Mathematical Representation / 24.2
- 24.5 Axial Gradient Lenses / 24.3
- 24.6 Radial Gradients / 24.5
- 24.7 Radial Gradients with Curved Surfaces / 24.7
- 24.8 Shallow Radial Gradients / 24.7
- 24.9 Materials / 24.8
- 24.10 References / 24.9

Part 5. Instruments

Chapter 25. Cameras <i>Norman Goldberg</i>	25.3
<hr/>	
25.1 Glossary / 25.3	
25.2 Introduction / 25.3	
25.3 Background / 25.4	
25.4 Properties of the Final Image / 25.5	
25.5 Film Choice / 25.5	
25.6 Resolving Fine Detail / 25.5	
25.7 Film Sizes / 25.6	
25.8 Display / 25.6	
25.9 Distributing the Image / 25.7	
25.10 Video Cameras / 25.7	
25.11 Instant Pictures / 25.8	
25.12 Critical Features / 25.8	
25.13 Time Lag / 25.8	
25.14 Automation / 25.10	
25.15 Flash / 25.16	
25.16 Flexibility through Features and Accessories / 25.16	
25.17 Advantages of Various Formats / 25.17	
25.18 Large Format: A Different World / 25.18	
25.19 Special Cameras / 25.20	
25.20 Further Reading / 25.26	
<hr/>	
Chapter 26. Solid-State Cameras <i>Gerald C. Holst</i>	26.1
<hr/>	
26.1 Glossary / 26.1	
26.2 Introduction / 26.2	
26.3 Imaging System Applications / 26.3	
26.4 Charge-Coupled Device Array Architecture / 26.3	
26.5 Charge Injection Device / 26.6	
26.6 Complementary Metal-Oxide Semiconductor / 26.8	
26.7 Array Performance / 26.9	
26.8 Camera Performance / 26.12	
26.9 Modulation Transfer Function / 26.14	
26.10 Resolution / 26.15	
26.11 Sampling / 26.16	
26.12 Storage, Analysis, and Display / 26.19	
26.13 References / 26.20	
<hr/>	
Chapter 27. Camera Lenses <i>Ellis Betensky, Melvin H. Kreitzer, and Jacob Moskovich</i>	27.1
<hr/>	
27.1 Introduction / 27.1	
27.2 Imposed Design Limitations / 27.1	
27.3 Modern Lens Types / 27.2	
27.4 Classification System / 27.17	
27.5 Lens Performance Data / 27.24	
27.6 Acknowledgments / 27.25	
27.7 Further Reading / 27.25	
<hr/>	
Chapter 28. Microscopes <i>Rudolf Oldenbourg and Michael Shribak</i>	28.1
<hr/>	
28.1 Glossary / 28.1	
28.2 Introduction / 28.1	
28.3 Optical Arrangements, Lenses, and Resolution / 28.3	
28.4 Contrast and Imaging Modes / 28.24	

- 28.5 Manipulation of Specimen / 28.54
28.6 Acknowledgment / 28.55
28.7 References / 28.56

Chapter 29. Reflective and Catadioptric Objectives **29.1**
Lloyd Jones

- 29.1 Glossary / 29.1
29.2 Introduction / 29.2
29.3 Glass Varieties / 29.2
29.4 Introduction to Catadioptric and Reflective Objectives / 29.2
29.5 Field-of-View Plots / 29.34
29.6 Definitions / 29.36
29.7 References / 29.38

Chapter 30. Scanners *Leo Beiser and R. Barry Johnson* **30.1**

- 30.1 Glossary / 30.1
30.2 Introduction / 30.2
30.3 Scanned Resolution / 30.6
30.4 Scanners for Remote Sensing / 30.14
30.5 Scanning for Input/Output Imaging / 30.25
30.6 Scanner Devices and Techniques / 30.34
30.7 Scan-Error Reduction / 30.48
30.8 Agile Beam Steering / 30.51
30.9 References / 30.64
30.10 Further Reading / 30.68

Chapter 31. Optical Spectrometers *Brian Henderson* **31.1**

- 31.1 Glossary / 31.1
31.2 Introduction / 31.2
31.3 Optical Absorption Spectrometers / 31.2
31.4 Luminescence Spectrometers / 31.5
31.5 Photoluminescence Decay Time / 31.12
31.6 Polarization Spectrometers / 31.15
31.7 High-Resolution Techniques / 31.23
31.8 Light Scattering / 31.30
31.9 References / 31.31

Chapter 32. Interferometers *Parameswaran Hariharan* **32.1**

- 32.1 Glossary / 32.1
32.2 Introduction / 32.2
32.3 Basic Types of Interferometers / 32.2
32.4 Three-Beam and Double-Passed Two-Beam Interferometers / 32.7
32.5 Fringe-Counting Interferometers / 32.8
32.6 Two-Wavelength Interferometry / 32.9
32.7 Frequency-Modulation Interferometers / 32.9
32.8 Heterodyne Interferometers / 32.10
32.9 Phase-Shifting Interferometers / 32.10
32.10 Phase-Locked Interferometers / 32.11
32.11 Laser-Doppler Interferometers / 32.12
32.12 Laser-Feedback Interferometers / 32.13
32.13 Fiber Interferometers / 32.14
32.14 Interferometric Wave Meters / 32.16
32.15 Second-Harmonic and Phase-Conjugate Interferometers / 32.17
32.16 Stellar Interferometers / 32.19
32.17 Gravitational-Wave Interferometers / 32.21
32.18 References / 32.22

Chapter 33. Holography and Holographic Instruments <i>Lloyd Huff</i>	33.1
<hr/>	
33.1 Glossary / 33.1	
33.2 Introduction / 33.2	
33.3 Background and Basic Principles / 33.2	
33.4 Holographic Interferometry / 33.4	
33.5 Holographic Optical Elements / 33.13	
33.6 Holographic Inspection / 33.16	
33.7 Holographic Lithography / 33.22	
33.8 Holographic Memory / 33.24	
33.9 Conclusion / 33.25	
33.10 References / 33.25	
Chapter 34. Xerographic Systems <i>Howard Stark</i>	34.1
<hr/>	
34.1 Introduction and Overview / 34.1	
34.2 Creation of the Latent Image / 34.1	
34.3 Development / 34.5	
34.4 Transfer / 34.10	
34.5 Fusing / 34.10	
34.6 Cleaning and Erasing / 34.10	
34.7 Control Systems / 34.11	
34.8 Color / 34.11	
34.9 References / 34.13	
Chapter 35. Principles of Optical Disk Data Storage <i>Masud Mansuripur</i>	35.1
<hr/>	
35.1 Introduction / 35.1	
35.2 Preliminaries and Basic Definitions / 35.2	
35.3 The Optical Path / 35.7	
35.4 Automatic Focusing / 35.12	
35.5 Automatic Tracking / 35.14	
35.6 Thermomagnetic Recording Process / 35.17	
35.7 Magneto-Optical Readout / 35.21	
35.8 Materials of Magneto-Optical Recording / 35.25	
35.9 Concluding Remarks / 35.28	
35.10 Further Information / 35.30	
35.11 Bibliography / 35.31	
Index I.1	

CONTENTS OF VOLUME II

Part 1. Design

Chapter 1. Techniques of First-Order Layout *Warren J. Smith* 1.3

- 1.1 Glossary / 1.3
- 1.2 First-Order Layout / 1.4
- 1.3 Ray-Tracing / 1.4
- 1.4 Two-Component Systems / 1.5
- 1.5 Afocal Systems / 1.7
- 1.6 Magnifiers and Microscopes / 1.8
- 1.7 Afocal Attachments / 1.8
- 1.8 Field Lenses / 1.8
- 1.9 Condensers / 1.10
- 1.10 Zoom or Varifocal Systems / 1.11
- 1.11 Additional Rays / 1.12
- 1.12 Minimizing Component Power / 1.13
- 1.13 Is It a Reasonable Layout? / 1.13
- 1.14 Achromatism / 1.14
- 1.15 Athermalization / 1.15

Chapter 2. Aberration Curves in Lens Design *Donald C. O'Shea and Michael E. Harrigan* 2.1

- 2.1 Glossary / 2.1
- 2.2 Introduction / 2.1
- 2.3 Transverse Ray Plots / 2.2
- 2.4 Field Plots / 2.4
- 2.5 Additional Considerations / 2.5
- 2.6 Summary / 2.6
- 2.7 References / 2.6

Chapter 3. Optical Design Software *Douglas C. Sinclair* 3.1

- 3.1 Glossary / 3.1
- 3.2 Introduction / 3.2
- 3.3 Lens Entry / 3.2
- 3.4 Evaluation / 3.8
- 3.5 Optimization / 3.16
- 3.6 Other Topics / 3.21
- 3.7 Buying Optical Design Software / 3.22
- 3.8 Summary / 3.24
- 3.9 References / 3.24

Chapter 4. Optical Specifications *Robert R. Shannon* 4.1

- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.1
- 4.3 Preparation of Optical Specifications / 4.5
- 4.4 Image Specifications / 4.6
- 4.5 Element Description / 4.8

- 4.6 Environmental Specifications / 4.10
- 4.7 Presentation of Specifications / 4.10
- 4.8 Problems with Specification Writing / 4.11
- 4.9 References / 4.12

Chapter 5. Tolerancing Techniques *Robert R. Shannon* **5.1**

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.1
- 5.3 Wavefront Tolerances / 5.3
- 5.4 Other Tolerances / 5.7
- 5.5 Starting Points / 5.8
- 5.6 Material Properties / 5.9
- 5.7 Tolerancing Procedures / 5.9
- 5.8 Problems in Tolerancing / 5.11
- 5.9 References / 5.11

Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.* **6.1**

- 6.1 Glossary / 6.1
- 6.2 Introduction and Summary / 6.1
- 6.3 Mounting Individual Rotationally Symmetric Optics / 6.2
- 6.4 Multicomponent Lens Assemblies / 6.5
- 6.5 Mounting Windows and Domes / 6.11
- 6.6 Mounting Small Mirrors and Prisms / 6.11
- 6.7 Mounting Moderate-Sized Mirrors / 6.17
- 6.8 Contact Stresses in Optics / 6.21
- 6.9 Temperature Effects on Mounted Optics / 6.21
- 6.10 References / 6.25

Chapter 7. Control of Stray Light *Robert P. Breault* **7.1**

- 7.1 Glossary / 7.1
- 7.2 Introduction / 7.1
- 7.3 Concepts / 7.2
- 7.4 Optical Software for Stray Light Analysis / 7.24
- 7.5 Methods / 7.27
- 7.6 Conclusion / 7.30
- 7.7 Sources of Information on Stray Light and Scattered Light / 7.31
- 7.8 References / 7.32

Chapter 8. Thermal Compensation Techniques
Philip J. Rogers and Michael Roberts **8.1**

- 8.1 Glossary / 8.1
- 8.2 Introduction / 8.2
- 8.3 Homogeneous Thermal Effects / 8.2
- 8.4 Tolerable Homogeneous Temperature Change (No Compensation) / 8.5
- 8.5 Effect of Thermal Gradients / 8.6
- 8.6 Intrinsic Athermalization / 8.7
- 8.7 Mechanical Athermalization / 8.8
- 8.8 Optical Athermalization / 8.12
- 8.9 References / 8.15

Part 2. Fabrication

Chapter 9. Optical Fabrication *Michael P. Mandina* **9.3**

- 9.1 Introduction / 9.3
- 9.2 Material Forms of Supply / 9.3

9.3	Basic Steps in Spherical Optics Fabrication /	9.4
9.4	Plano Optics Fabrication /	9.7
9.5	Asphere Optics Fabrication /	9.7
9.6	Crystalline Optics /	9.8
9.7	Purchasing Optics /	9.9
9.8	Conclusion /	9.9
9.9	References /	9.9

Chapter 10.	Fabrication of Optics by Diamond Turning	10.1
	<i>Richard L. Rhorer and Chris J. Evans</i>	

10.1	Glossary /	10.1
10.2	Introduction /	10.1
10.3	The Diamond-Turning Process /	10.2
10.4	The Advantages of Diamond Turning /	10.2
10.5	Diamond-Turnable Materials /	10.4
10.6	Comparison of Diamond Turning and Traditional Optical Fabrication /	10.6
10.7	Machine Tools for Diamond Turning /	10.6
10.8	Basic Steps in Diamond Turning /	10.8
10.9	Surface Finish of Diamond-Turned Optics /	10.9
10.10	Metrology of Diamond-Turned Optics /	10.12
10.11	Conclusions /	10.13
10.12	References /	10.14

Part 3. Testing

Chapter 11.	Orthonormal Polynomials in Wavefront Analysis	11.3
	<i>Virendra N. Mahajan</i>	

	Abstract /	11.3
11.1	Glossary /	11.3
11.2	Introduction /	11.4
11.3	Orthonormal Polynomials /	11.5
11.4	Zernike Circle Polynomials /	11.6
11.5	Zernike Annular Polynomials /	11.13
11.6	Hexagonal Polynomials /	11.21
11.7	Elliptical Polynomials /	11.21
11.8	Rectangular Polynomials /	11.27
11.9	Square Polynomials /	11.30
11.10	Slit Polynomials /	11.30
11.11	Aberration Balancing and Tolerancing, and Diffraction Focus /	11.30
11.12	Isometric, Interferometric, and PSF Plots for Orthonormal Aberrations /	11.36
11.13	Use of Circle Polynomials for Noncircular Pupils /	11.37
11.14	Discussion and Conclusions /	11.39
11.15	References /	11.40

Chapter 12.	Optical Metrology	12.1
	<i>Zacarias Malacara and Daniel Malacara-Hernández</i>	

12.1	Glossary /	12.1
12.2	Introduction and Definitions /	12.2
12.3	Length and Straightness Measurements /	12.2
12.4	Angle Measurements /	12.10
12.5	Curvature and Focal Length Measurements /	12.17
12.6	References /	12.25

Chapter 13.	Optical Testing	13.1
	<i>Daniel Malacara-Hernández</i>	

13.1	Glossary /	13.1
13.2	Introduction /	13.1

- 13.3 Classical Noninterferometric Tests / 13.1
- 13.4 Interferometric Tests / 13.7
- 13.5 Increasing the Sensitivity of Interferometers / 13.13
- 13.6 Interferogram Evaluation / 13.14
- 13.7 Phase-Shifting Interferometry / 13.18
- 13.8 Measuring Aspherical Wavefronts / 13.23
- 13.9 References / 13.28

Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant* 14.1

- 14.1 Glossary / 14.1
- 14.2 Introduction / 14.1
- 14.3 Plotting CGHs / 14.3
- 14.4 Interferometers Using Computer-Generated Holograms / 14.4
- 14.5 Accuracy Limitations / 14.6
- 14.6 Experimental Results / 14.7
- 14.7 Discussion / 14.9
- 14.8 References / 14.9

Part 4. Sources

Chapter 15. Artificial Sources *Anthony LaRocca* 15.3

- 15.1 Glossary / 15.3
- 15.2 Introduction / 15.3
- 15.3 Radiation Law / 15.4
- 15.4 Laboratory Sources / 15.7
- 15.5 Commercial Sources / 15.13
- 15.6 References / 15.53

Chapter 16. Lasers *William T. Silfvast* 16.1

- 16.1 Glossary / 16.1
- 16.2 Introduction / 16.2
- 16.3 Laser Properties Associated with the Laser Gain Medium / 16.4
- 16.4 Laser Properties Associated with Optical Cavities or Resonators / 16.19
- 16.5 Special Laser Cavities / 16.25
- 16.6 Specific Types of Lasers / 16.29
- 16.7 References / 16.37

Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman* 17.1

- 17.1 Glossary / 17.1
- 17.2 Introduction / 17.2
- 17.3 Light-Generation Processes / 17.2
- 17.4 Light Extraction / 17.6
- 17.5 Device Structures / 17.8
- 17.6 Material Systems / 17.15
- 17.7 Substrate Technology / 17.20
- 17.8 Epitaxial Technology / 17.21
- 17.9 Wafer Processing / 17.23
- 17.10 Led Quality and Reliability / 17.25
- 17.11 Led-Based Products / 17.29
- 17.12 References / 17.35

Chapter 18. High-Brightness Visible LEDs <i>Winston V. Schoenfeld</i>	18.1
<hr/>	
18.1 The Materials Systems / 18.1	
18.2 Substrates and Epitaxial Growth / 18.2	
18.3 Processing / 18.3	
18.4 Solid-State Lighting / 18.4	
18.5 Packaging / 18.5	
Chapter 19. Semiconductor Lasers <i>Pamela L. Derry, Luis Figueroa, and Chi-Shain Hong</i>	19.1
<hr/>	
19.1 Glossary / 19.1	
19.2 Introduction / 19.3	
19.3 Applications for Semiconductor Lasers / 19.3	
19.4 Basic Operation / 19.4	
19.5 Fabrication and Configurations / 19.6	
19.6 Quantum Well Lasers / 19.9	
19.7 High-Power Semiconductor Lasers / 19.18	
19.8 High-Speed Modulation / 19.30	
19.9 Spectral Properties / 19.36	
19.10 Surface-Emitting Lasers / 19.39	
19.11 Conclusion / 19.41	
19.12 References / 19.43	
Chapter 20. Ultrashort Optical Sources and Applications <i>Jean-Claude Diels and Ladan Arissian</i>	20.1
<hr/>	
20.1 Introduction / 20.1	
20.2 Description of Optical Pulses and Pulse Trains / 20.2	
20.3 Pulse Evolution toward Steady State / 20.9	
20.4 Coupling Circulating Pulses Inside a Cavity / 20.12	
20.5 Designs of Cavities with Two Circulating Pulses / 20.15	
20.6 Analogy of a Two-Level System / 20.22	
20.7 Conclusion / 20.28	
20.8 References / 20.28	
Chapter 21. Attosecond Optics <i>Zenghu Chang</i>	21.1
<hr/>	
21.1 Glossary / 21.1	
21.2 Introduction / 21.2	
21.3 The Driving Laser / 21.4	
21.4 Attosecond Pulse Generation / 21.6	
21.5 Attosecond Pulse Characterization / 21.8	
21.6 Acknowledgments / 21.10	
21.7 References / 21.10	
Chapter 22. Laser Stabilization <i>John L. Hall, Matthew S. Taubman, and Jun Ye</i>	22.1
<hr/>	
22.1 Introduction and Overview / 22.1	
22.2 Servo Principles and Issues / 22.5	
22.3 Practical Issues / 22.12	
22.4 Summary and Outlook / 22.23	
22.5 Conclusions and Recommendations / 22.24	
22.6 Acknowledgments / 22.24	
22.7 References / 22.24	

Chapter 23. Quantum Theory of the Laser *János A. Bergou,
Berthold-Georg Englert, Melvin Lax, Marian O. Scully,
Herbert Walther, and M. Suhail Zubairy* 23.1

- 23.1 Glossary / 23.1
23.2 Introduction / 23.5
23.3 Some History of the Photon Concept / 23.6
23.4 Quantum Theory of the Laser / 23.14
23.5 The Laser Phase-Transition Analogy / 23.35
23.6 Exotic Masers and Lasers / 23.40
23.7 Acknowledgments / 23.45
23.8 References / 23.46

Part 5. Detectors

Chapter 24. Photodetectors *Paul R. Norton* 24.3

- 24.1 Scope / 24.3
24.2 Thermal Detectors / 24.4
24.3 Quantum Detectors / 24.6
24.4 Definitions / 24.10
24.5 Detector Performance and Sensitivity / 24.13
24.6 Other Performance Parameters / 24.18
24.7 Detector Performance / 24.21
24.8 References / 24.101
24.9 Suggested Readings / 24.102

Chapter 25. Photodetection *Abhay M. Joshi
and Gregory H. Olsen* 25.1

- 25.1 Glossary / 25.1
25.2 Introduction / 25.2
25.3 Principle of Operation / 25.3
25.4 Applications / 25.11
25.5 Reliability / 25.13
25.6 Future Photodetectors / 25.15
25.7 Acknowledgment / 25.17
25.8 References / 25.18
25.9 Additional Reading / 25.19

Chapter 26. High-Speed Photodetectors
J. E. Bowers and Y. G. Wey 26.1

- 26.1 Glossary / 26.1
26.2 Introduction / 26.3
26.3 Photodetector Structures / 26.3
26.4 Speed Limitations / 26.5
26.5 *p-i-n* Photodetectors / 26.10
26.6 Schottky Photodiode / 26.16
26.7 Avalanche Photodetectors / 26.17
26.8 Photoconductors / 26.20
26.9 Summary / 26.24
26.10 References / 26.24

Chapter 27. Signal Detection and Analysis
John R. Willison 27.1

- 27.1 Glossary / 27.1
27.2 Introduction / 27.1

27.3	Prototype Experiment / 27.2
27.4	Noise Sources / 27.3
27.5	Applications Using Photomultipliers / 27.6
27.6	Amplifiers / 27.10
27.7	Signal Analysis / 27.12
27.8	References / 27.15

Chapter 28.	Thermal Detectors	<i>William L. Wolfe and Paul W. Kruse</i>	28.1
--------------------	--------------------------	---	-------------

28.1	Glossary / 28.1
28.2	Thermal Detector Elements / 28.1
28.3	Arrays / 28.7
28.4	References / 28.13

Part 6. Imaging Detectors

Chapter 29.	Photographic Films	<i>Joseph H. Altman</i>	29.3
--------------------	---------------------------	-------------------------	-------------

29.1	Glossary / 29.3
29.2	Structure of Silver Halide Photographic Layers / 29.4
29.3	Grains / 29.5
29.4	Processing / 29.5
29.5	Exposure / 29.5
29.6	Optical Density / 29.6
29.7	The D-Log H Curve / 29.8
29.8	Spectral Sensitivity / 29.11
29.9	Reciprocity Failure / 29.11
29.10	Development Effects / 29.12
29.11	Color Photography / 29.12
29.12	Microdensitometers / 29.15
29.13	Performance of Photographic Systems / 29.16
29.14	Image Structure / 29.17
29.15	Acutance / 29.17
29.16	Graininess / 29.19
29.17	Sharpness and Graininess Considered Together / 29.22
29.18	Signal-to-Noise Ratio and Detective Quantum Efficiency / 29.22
29.19	Resolving Power / 29.24
29.20	Information Capacity / 29.24
29.21	List of Photographic Manufacturers / 29.25
29.22	References / 29.25

Chapter 30.	Photographic Materials	<i>John D. Baloga</i>	30.1
--------------------	-------------------------------	-----------------------	-------------

30.1	Introduction / 30.1
30.2	The Optics of Photographic Films and Papers / 30.2
30.3	The Photophysics of Silver Halide Light Detectors / 30.7
30.4	The Stability of Photographic Image Dyes toward Light Fade / 30.10
30.5	Photographic Spectral Sensitizers / 30.13
30.6	General Characteristics of Photographic Films / 30.18
30.7	References / 30.28

Chapter 31.	Image Tube Intensified Electronic Imaging	<i>C. Bruce Johnson and Larry D. Owen</i>	31.1
--------------------	--	---	-------------

31.1	Glossary / 31.1
31.2	Introduction / 31.2
31.3	The Optical Interface / 31.3

- 31.4 Image Intensifiers / 31.7
- 31.5 Image Intensified Self-Scanned Arrays / 31.19
- 31.6 Applications / 31.27
- 31.7 References / 31.30

Chapter 32. Visible Array Detectors *Timothy J. Tredwell* **32.1**

- 32.1 Glossary / 32.1
- 32.2 Introduction / 32.2
- 32.3 Image Sensing Elements / 32.2
- 32.4 Readout Elements / 32.12
- 32.5 Sensor Architectures / 32.21
- 32.6 References / 32.35

Chapter 33. Infrared Detector Arrays *Lester J. Kozlowski and Walter F. Kosonocky* **33.1**

- 33.1 Glossary / 33.1
- 33.2 Introduction / 33.3
- 33.3 Monolithic FPAs / 33.10
- 33.4 Hybrid FPAs / 33.14
- 33.5 Performance: Figures of Merit / 33.23
- 33.6 Current Status and Future Trends / 33.28
- 33.7 References / 33.31

Part 7. Radiometry and Photometry

Chapter 34. Radiometry and Photometry *Edward F. Zalewski* **34.3**

- 34.1 Glossary / 34.3
- 34.2 Introduction / 34.5
- 34.3 Radiometric Definitions and Basic Concepts / 34.7
- 34.4 Radiant Transfer Approximations / 34.13
- 34.5 Absolute Measurements / 34.20
- 34.6 Photometry / 34.37
- 34.7 References / 34.44

Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection *James M. Palmer* **35.1**

- 35.1 Glossary / 35.1
- 35.2 Introduction and Terminology / 35.2
- 35.3 Transmittance / 35.3
- 35.4 Absorptance / 35.4
- 35.5 Reflectance / 35.4
- 35.6 Emittance / 35.7
- 35.7 Kirchhoff's Law / 35.7
- 35.8 Relationship between Transmittance, Reflectance, and Absorptance / 35.7
- 35.9 Measurement of Transmittance / 35.8
- 35.10 Measurement of Absorptance / 35.10
- 35.11 Measurement of Reflectance / 35.10
- 35.12 Measurement of Emittance / 35.14
- 35.13 References / 35.16
- 35.14 Further Reading / 35.23

Chapter 36. Radiometry and Photometry: Units and Conversions	<i>James M. Palmer</i>	36.1
<hr/>		
36.1	Glossary / 36.1	
36.2	Introduction and Background / 36.2	
36.3	Symbols, Units, and Nomenclature in Radiometry / 36.4	
36.4	Symbols, Units, and Nomenclature in Photometry / 36.5	
36.5	Conversion of Radiometric Quantities to Photometric Quantities / 36.11	
36.6	Conversion of Photometric Quantities to Radiometric Quantities / 36.12	
36.7	Radiometric/Photometric Normalization / 36.14	
36.8	Other Weighting Functions and Conversions / 36.17	
36.9	References / 36.17	
36.10	Further Reading / 36.18	
Chapter 37. Radiometry and Photometry for Vision Optics	<i>Yoshi Ohno</i>	37.1
<hr/>		
37.1	Introduction / 37.1	
37.2	Basis of Physical Photometry / 37.1	
37.3	Photometric Base Unit—the Candela / 37.3	
37.4	Quantities and Units in Photometry and Radiometry / 37.3	
37.5	Principles in Photometry and Radiometry / 37.8	
37.6	Practice in Photometry and Radiometry / 37.11	
37.7	References / 37.12	
Chapter 38. Spectroradiometry	<i>Carolyn J. Sher DeCusatis</i>	38.1
<hr/>		
38.1	Introduction / 38.1	
38.2	Definitions, Calculations, and Figures of Merit / 38.1	
38.3	General Features of Spectroradiometry Systems / 38.7	
38.4	Typical Spectroradiometry System Designs / 38.13	
38.5	References / 38.19	
Chapter 39. Nonimaging Optics: Concentration and Illumination	<i>William Cassarly</i>	39.1
<hr/>		
39.1	Introduction / 39.1	
39.2	Basic Calculations / 39.2	
39.3	Software Modeling of Nonimaging Systems / 39.6	
39.4	Basic Building Blocks / 39.8	
39.5	Concentration / 39.12	
39.6	Uniformity and Illumination / 39.22	
39.7	Acknowledgments / 39.41	
39.8	References / 39.41	
Chapter 40. Lighting and Applications	<i>Anurag Gupta and R. John Koschel</i>	40.1
<hr/>		
40.1	Glossary / 40.1	
40.2	Introduction / 40.1	
40.3	Vision Biology and Perception / 40.3	
40.4	The Science of Lighting Design / 40.6	
40.5	Luminaires / 40.24	
40.6	Lighting Measurements / 40.51	
40.7	Lighting Application Areas / 40.54	
40.8	Acknowledgments / 40.71	
40.9	References / 40.72	

This page intentionally left blank.

DO NOT DUPLICATE

CONTENTS OF VOLUME III

Chapter 1. Optics of the Eye	<i>Neil Charman</i>	1.1
<hr/>		
1.1	Glossary / 1.1	
1.2	Introduction / 1.3	
1.3	Ocular Parameters and Ametropia / 1.4	
1.4	Ocular Transmittance and Retinal Illuminance / 1.8	
1.5	Factors Affecting In-Focus Retinal Image Quality / 1.12	
1.6	Final Retinal Image Quality / 1.21	
1.7	Depth-of-Focus and Accommodation / 1.28	
1.8	Eye Models / 1.36	
1.9	Two Eyes and Stereopsis / 1.38	
1.10	Movements of the Eyes / 1.42	
1.11	Conclusion / 1.45	
1.12	References / 1.45	
Chapter 2. Visual Performance	<i>Wilson S. Geisler and Martin S. Banks</i>	2.1
<hr/>		
2.1	Glossary / 2.1	
2.2	Introduction / 2.2	
2.3	Optics, Anatomy, Physiology of the Visual System / 2.2	
2.4	Visual Performance / 2.14	
2.5	Acknowledgments / 2.41	
2.6	References / 2.42	
Chapter 3. Psychophysical Methods	<i>Denis G. Pelli and Bart Farell</i>	3.1
<hr/>		
3.1	Introduction / 3.1	
3.2	Definitions / 3.2	
3.3	Visual Stimuli / 3.3	
3.4	Adjustments / 3.4	
3.5	Judgments / 3.6	
	Magnitude Estimation / 3.8	
3.6	Stimulus Sequencing / 3.9	
3.7	Conclusion / 3.9	
3.8	Tips from the Pros / 3.10	
3.9	Acknowledgments / 3.10	
3.10	References / 3.10	
Chapter 4. Visual Acuity and Hyperacuity	<i>Gerald Westheimer</i>	4.1
<hr/>		
4.1	Glossary / 4.1	
4.2	Introduction / 4.2	
4.3	Stimulus Specifications / 4.2	
4.4	Optics of the Eye's Resolving Capacity / 4.4	
4.5	Retinal Limitations—Receptor Mosaic and Tiling of Neuronal Receptive Fields / 4.5	
4.6	Determination of Visual Resolution Thresholds / 4.6	
4.7	Kinds of Visual Acuity Tests / 4.7	
4.8	Factors Affecting Visual Acuity / 4.9	
4.9	Hyperacuity / 4.14	
4.10	Resolution, Superresolution, and Information Theory / 4.15	

- 4.11 Summary / 4.16
4.12 References / 4.16

Chapter 5. Optical Generation of the Visual Stimulus *Stephen A. Burns and Robert H. Webb* **5.1**

- 5.1 Glossary / 5.1
5.2 Introduction / 5.1
5.3 The Size of the Visual Stimulus / 5.2
5.4 Free or Newtonian Viewing / 5.2
5.5 Maxwellian Viewing / 5.4
5.6 Building an Optical System / 5.8
5.7 Light Exposure and Ocular Safety / 5.18
5.8 Light Sources / 5.19
5.9 Coherent Radiation / 5.19
5.10 Detectors / 5.21
5.11 Putting It Together / 5.21
5.12 Conclusions / 5.24
5.13 Acknowledgments / 5.24
5.14 General References / 5.25
5.15 References / 5.26

Chapter 6. The Maxwellian View: with an Addendum on Apodization *Gerald Westheimer* **6.1**

- 6.1 Glossary / 6.1
6.2 Introduction / 6.2
6.3 Postscript (2008) / 6.13

Chapter 7. Ocular Radiation Hazards *David H. Sliney* **7.1**

- 7.1 Glossary / 7.1
7.2 Introduction / 7.2
7.3 Injury Mechanisms / 7.2
7.4 Types of Injury / 7.3
7.5 Retinal Irradiance Calculations / 7.7
7.6 Examples / 7.8
7.7 Exposure Limits / 7.9
7.8 Discussion / 7.11
7.9 References / 7.15

Chapter 8. Biological Waveguides *Vasudevan Lakshminarayanan and Jay M. Enoch* **8.1**

- 8.1 Glossary / 8.1
8.2 Introduction / 8.2
8.3 Waveguiding in Retinal Photoreceptors and the Stiles-Crawford Effect / 8.3
8.4 Waveguides and Photoreceptors / 8.3
8.5 Photoreceptor Orientation and Alignment / 8.5
8.6 Introduction to the Models and Theoretical Implications / 8.8
8.7 Quantitative Observations of Single Receptors / 8.15
8.8 Waveguide Modal Patterns Found in Monkey/Human Retinal Receptors / 8.19
8.9 Light Guide Effect in Cochlear Hair Cells and Human Hair / 8.24
8.10 Fiber-Optic Plant Tissues / 8.26
8.11 Sponges / 8.28
8.12 Summary / 8.29
8.13 References / 8.29

Chapter 9. The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution		9.1
<i>Jay M. Enoch and Vasudevan Lakshminarayanan</i>		
9.1	Glossary / 9.1	
9.2	Introduction / 9.2	
9.3	The Problem and an Approach to Its Solution / 9.3	
9.4	Sample Point-by-Point Estimates of SCE-1 and Integrated SCE-1 Data / 9.6	
9.5	Discussion / 9.13	
9.6	Teleological and Developmental Factors / 9.14	
9.7	Conclusions / 9.14	
9.8	References / 9.15	
Chapter 10. Colorimetry		10.1
<i>David H. Brainard and Andrew Stockman</i>		
10.1	Glossary / 10.1	
10.2	Introduction / 10.2	
10.3	Fundamentals of Colorimetry / 10.3	
10.4	Color Coordinate Systems / 10.11	
10.5	Matrix Representations and Calculations / 10.24	
10.6	Topics / 10.32	
10.7	Appendix—Matrix Algebra / 10.45	
10.8	References / 10.49	
Chapter 11. Color Vision Mechanisms		11.1
<i>Andrew Stockman and David H. Brainard</i>		
11.1	Glossary / 11.1	
11.2	Introduction / 11.3	
11.3	Basics of Color-Discrimination Mechanisms / 11.9	
11.4	Basics of Color-Appearance Mechanisms / 11.26	
11.5	Details and Limits of the Basic Model / 11.31	
11.6	Conclusions / 11.79	
11.7	Acknowledgments / 11.85	
11.8	References / 11.86	
Chapter 12. Assessment of Refraction and Refractive Errors and Their Influence on Optical Design		12.1
<i>B. Ralph Chou</i>		
12.1	Glossary / 12.1	
12.2	Introduction / 12.3	
12.3	Refractive Errors / 12.3	
12.4	Assessment of Refractive Error / 12.5	
12.5	Correction of Refractive Error / 12.8	
12.6	Binocular Factors / 12.15	
12.7	Consequences for Optical Design / 12.17	
12.8	References / 12.17	
Chapter 13. Binocular Vision Factors That Influence Optical Design		13.1
<i>Clifton Schor</i>		
13.1	Glossary / 13.1	
13.2	Combining the Images in the Two Eyes into One Perception of the Visual Field / 13.3	
13.3	Distortion of Space by Monocular Magnification / 13.13	
13.4	Distortion of Space Perception from Interocular Aniso-Magnification (Unequal Binocular Magnification) / 13.16	
13.5	Distortions of Space from Convergence Responses to Prism / 13.19	
13.6	Eye Movements / 13.19	

- 13.7 Coordination and Alignment of the Two Eyes / 13.20
13.8 Effects of Lenses and Prism on Vergence and Phoria / 13.25
13.9 Prism-Induced Errors of Eye Alignment / 13.27
13.10 Head and Eye Responses to Direction (Gaze Control) / 13.29
13.11 Focus and Responses to Distance / 13.30
13.12 Video Head Sets, Head's Up Displays and Virtual Reality: Impact on Binocular Vision / 13.31
13.13 References / 13.35

Chapter 14. Optics and Vision of the Aging Eye *John S. Werner, Brooke E. Scheffrin, and Arthur Bradley* 14.1

- 14.1 Glossary / 14.1
14.2 Introduction / 14.2
14.3 The Graying of the Planet / 14.2
14.4 Senescence of the Eye's Optics / 14.4
14.5 Senescent Changes in Vision / 14.14
14.6 Age-Related Ocular Diseases Affecting Visual Function / 14.22
14.7 The Aging World from the Optical Point of View: Presbyopic Corrections / 14.27
14.8 Conclusions / 14.30
14.9 Acknowledgments / 14.30
14.10 References / 14.30

Chapter 15. Adaptive Optics in Retinal Microscopy and Vision *Donald T. Miller and Austin Roorda* 15.1

- 15.1 Glossary / 15.1
15.2 Introduction / 15.2
15.3 Properties of Ocular Aberrations / 15.4
15.4 Implementation of AO / 15.7
15.5 Application of AO to the Eye / 15.15
15.6 Acknowledgments / 15.24
15.7 References / 15.24

Chapter 16. Refractive Surgery, Correction of Vision, PRK and LASIK *L. Diaz-Santana and Harilaos Ginis* 16.1

- 16.1 Glossary / 16.1
16.2 Introduction / 16.2
16.3 Refractive Surgery Modalities / 16.9
16.4 Laser Ablation / 16.15
16.5 Acknowledgments / 16.19
16.6 References / 16.19

Chapter 17. Three-Dimensional Confocal Microscopy of the Living Human Cornea *Barry R. Masters* 17.1

- 17.1 Glossary / 17.1
17.2 Introduction / 17.3
17.3 Theory of Confocal Microscopy / 17.3
17.4 The Development of Confocal Instruments / 17.3
17.5 The Scanning Slit and Laser Scanning Clinical Confocal Microscopes / 17.6
17.6 Clinical Applications of Confocal Microscopy / 17.8
17.7 Perspectives / 17.9
17.8 Summary / 17.10
17.9 Acknowledgments / 17.10
17.10 References / 17.10

Chapter 18. Diagnostic Use of Optical Coherence Tomography in the Eye *Johannes F. de Boer* 18.1

- 18.1 Glossary / 18.1
18.2 Introduction / 18.2

18.3	Principle of OCT: Time Domain OCT /	18.3
18.4	Principle of OCT: Spectral Domain OCT /	18.5
18.5	Principle of OCT: Optical Frequency Domain Imaging /	18.7
18.6	SD-OCT Versus OFDI /	18.9
18.7	Sensitivity Advantage of SD-OCT Over TD-OCT /	18.9
18.8	Noise Analysis of SD-OCT Using Charge Coupled Devices (CCDs) /	18.9
18.9	Signal to Noise Ratio and Autocorrelation Noise /	18.11
18.10	Shot-Noise-Limited Detection /	18.12
18.11	Depth Dependent Sensitivity /	18.13
18.12	Motion Artifacts and Fringe Washout /	18.15
18.13	OFDI at 1050 NM /	18.15
18.14	Functional Extensions: Doppler OCT and Polarization Sensitive OCT /	18.18
18.15	Doppler OCT and Phase Stability /	18.18
18.16	Polarization Sensitive OCT (PS-OCT) /	18.20
18.17	PS-OCT in Ophthalmology /	18.24
18.18	Retinal Imaging with SD-OCT /	18.27
18.19	Conclusion /	18.29
18.20	Acknowledgment /	18.30
18.21	References /	18.30

Chapter 19. Gradient Index Optics in the Eye

Barbara K. Pierscionek

19.1

19.1	Glossary /	19.1
19.2	Introduction /	19.2
19.3	The Nature of an Index Gradient /	19.2
19.4	Spherical Gradients /	19.2
19.5	Radial Gradients /	19.3
19.6	Axial Gradients /	19.5
19.7	The Eye Lens /	19.5
19.8	Fish /	19.6
19.9	Octopus /	19.7
19.10	Rat /	19.7
19.11	Guinea Pig /	19.8
19.12	Rabbit /	19.8
19.13	Cat /	19.9
19.14	Bovine /	19.9
19.15	Pig /	19.11
19.16	Human/primate /	19.12
19.17	Functional Considerations /	19.14
19.18	Summary /	19.15
19.19	References /	19.15

Chapter 20. Optics of Contact Lenses *Edward S. Bennett*

20.1

20.1	Glossary /	20.1
20.2	Introduction /	20.2
20.3	Contact Lens Material, Composition, and Design Parameters /	20.3
20.4	Contact Lens Power /	20.6
20.5	Other Design Considerations /	20.20
20.6	Convergence and Accommodation Effects /	20.25
20.7	Prismatic Effects /	20.30
20.8	Magnification /	20.31
20.9	Summary /	20.34
20.10	Acknowledgments /	20.34
20.11	References /	20.34

Chapter 21. Intraocular Lenses *Jim Schwiegerling*

21.1

21.1	Glossary /	21.1
21.2	Introduction /	21.2
21.3	Cataract Surgery /	21.4

- 21.4 Intraocular Lens Design / 21.5
- 21.5 Intraocular Lens Side Effects / 21.20
- 21.6 Summary / 21.22
- 21.7 References / 21.22

Chapter 22. Displays for Vision Research *William Cowan* 22.1

- 22.1 Glossary / 22.1
- 22.2 Introduction / 22.2
- 22.3 Operational Characteristics of Color Monitors / 22.3
- 22.4 Colorimetric Calibration of Video Monitors / 22.20
- 22.5 An Introduction to Liquid Crystal Displays / 22.34
- 22.6 Acknowledgments / 22.40
- 22.7 References / 22.40

Chapter 23. Vision Problems at Computers *Jeffrey Anshel and James E. Sheedy* 23.1

- 23.1 Glossary / 23.1
- 23.2 Introduction / 23.4
- 23.3 Work Environment / 23.4
- 23.4 Vision and Eye Conditions / 23.9
- 23.5 References / 23.12

Chapter 24. Human Vision and Electronic Imaging
Bernice E. Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allebach 24.1

- 24.1 Introduction / 24.1
- 24.2 Early Vision Approaches: The Perception of Imaging Artifacts / 24.2
- 24.3 Higher-Level Approaches: The Analysis of Image Features / 24.6
- 24.4 Very High-Level Approaches: The Representation of Aesthetic and Emotional Characteristics / 24.9
- 24.5 Conclusions / 24.10
- 24.6 Additional Information on Human Vision and Electronic Imaging / 24.11
- 24.7 References / 24.11

Chapter 25. Visual Factors Associated with Head-mounted Displays *Brian H. Tsou and Martin Shenker* 25.1

- 25.1 Glossary / 25.1
- 25.2 Introduction / 25.1
- 25.3 Common Design Considerations among All HMDs / 25.2
- 25.4 Characterizing HMD / 25.7
- 25.5 Summary / 25.10
- 25.6 Appendix / 25.10
- 25.7 Acknowledgments / 25.12
- 25.8 References / 25.12

Index I.1

CONTENTS OF VOLUME IV

Part 1. Properties

Chapter 1. Optical Properties of Water *Curtis D. Mobley* 1.3

- 1.1 Introduction / 1.3
- 1.2 Terminology, Notation, and Definitions / 1.3
- 1.3 Radiometric Quantities Useful in Hydrologic Optics / 1.4
- 1.4 Inherent Optical Properties / 1.9
- 1.5 Apparent Optical Properties / 1.12
- 1.6 The Optically Significant Constituents of Natural Waters / 1.13
- 1.7 Particle Size Distributions / 1.15
- 1.8 Electromagnetic Properties of Water / 1.16
- 1.9 Index of Refraction / 1.18
- 1.10 Measurement of Absorption / 1.20
- 1.11 Absorption by Pure Sea Water / 1.21
- 1.12 Absorption by Dissolved Organic Matter / 1.22
- 1.13 Absorption by Phytoplankton / 1.23
- 1.14 Absorption by Organic Detritus / 1.25
- 1.15 Bio-Optical Models for Absorption / 1.27
- 1.16 Measurement of Scattering / 1.29
- 1.17 Scattering by Pure Water and by Pure Sea Water / 1.30
- 1.18 Scattering by Particles / 1.30
- 1.19 Wavelength Dependence of Scattering: Bio-Optical Models / 1.35
- 1.20 Beam Attenuation / 1.40
- 1.21 Diffuse Attenuation and Jerlov Water Types / 1.42
- 1.22 Irradiance Reflectance and Remote Sensing / 1.46
- 1.23 Inelastic Scattering and Polarization / 1.47
- 1.24 Acknowledgments / 1.50
- 1.25 References / 1.50

Chapter 2. Properties of Crystals and Glasses *William J. Tropf, Michael E. Thomas, and Eric W. Rogala* 2.1

- 2.1 Glossary / 2.1
- 2.2 Introduction / 2.3
- 2.3 Optical Materials / 2.4
- 2.4 Properties of Materials / 2.5
- 2.5 Properties Tables / 2.36
- 2.6 References / 2.77

Chapter 3. Polymeric Optics *John D. Lytle* 3.1

- 3.1 Glossary / 3.1
- 3.2 Introduction / 3.1
- 3.3 Forms / 3.2
- 3.4 Physical Properties / 3.2
- 3.5 Optical Properties / 3.5
- 3.6 Optical Design / 3.7
- 3.7 Processing / 3.11
- 3.8 Coatings / 3.17
- 3.9 References / 3.18

Chapter 4. Properties of Metals *Roger A. Paquin* 4.1

- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.2
- 4.3 Summary Data / 4.11
- 4.4 References / 4.70

Chapter 5. Optical Properties of Semiconductors
*David G. Seiler, Stefan Zollner, Alain C. Diebold,
and Paul M. Amirtharaj* 5.1

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.3
- 5.3 Optical Properties / 5.8
- 5.4 Measurement Techniques / 5.56
- 5.5 Acknowledgments / 5.83
- 5.6 Summary and Conclusions / 5.83
- 5.7 References / 5.91

Chapter 6. Characterization and Use of Black Surfaces for Optical Systems *Stephen M. Pompea and Robert P. Breault* 6.1

- 6.1 Introduction / 6.1
- 6.2 Selection Process for Black Baffle Surfaces in Optical Systems / 6.10
- 6.3 The Creation of Black Surfaces for Specific Applications / 6.13
- 6.4 Environmental Degradation of Black Surfaces / 6.16
- 6.5 Optical Characterization of Black Surfaces / 6.18
- 6.6 Surfaces for Ultraviolet and Far-Infrared Applications / 6.21
- 6.7 Survey of Surfaces with Optical Data / 6.34
- 6.8 Paints / 6.35
- 6.9 Conclusions / 6.59
- 6.10 Acknowledgments / 6.59
- 6.11 References / 6.60
- 6.12 Further Readings / 6.67

Chapter 7. Optical Properties of Films and Coatings
Jerzy A. Dobrowolski 7.1

- 7.1 Introduction / 7.1
- 7.2 Theory and Design of Optical Thin-Film Coatings / 7.5
- 7.3 Thin-Film Manufacturing Considerations / 7.10
- 7.4 Measurements on Optical Coatings / 7.12
- 7.5 Antireflection Coatings / 7.15
- 7.6 Two-Material Periodic Multilayers Theory / 7.32
- 7.7 Multilayer Reflectors—Experimental Results / 7.39
- 7.8 Cutoff, Heat-Control, and Solar-Cell Cover Filters / 7.53
- 7.9 Beam Splitters and Neutral Filters / 7.61
- 7.10 Interference Polarizers and Polarizing Beam Splitters / 7.69
- 7.11 Bandpass Filters / 7.73
- 7.12 High Performance Optical Multilayer Coatings / 7.96
- 7.13 Multilayers for Two or Three Spectral Regions / 7.98
- 7.14 Phase Coatings / 7.101
- 7.15 Interference Filters with Low Reflection / 7.104
- 7.16 Reflection Filters and Coatings / 7.106
- 7.17 Special Purpose Coatings / 7.113
- 7.18 References / 7.114

Chapter 8. Fundamental Optical Properties of Solids *Alan Miller* **8.1**

- 8.1 Glossary / 8.1
- 8.2 Introduction / 8.3
- 8.3 Propagation of Light in Solids / 8.4
- 8.4 Dispersion Relations / 8.14
- 8.5 Lattice Interactions / 8.16
- 8.6 Free Electron Properties / 8.21
- 8.7 Band Structures and Interband Transitions / 8.24
- 8.8 References / 8.32

Chapter 9. Photonic Bandgap Materials *Pierre R. Villeneuve* **9.1**

- 9.1 Glossary / 9.1
- 9.2 Introduction / 9.2
- 9.3 Maxwell's Equations / 9.2
- 9.4 Three-Dimensional Photonic Crystals / 9.4
- 9.5 Microcavities in Three-Dimensional Photonic Crystals / 9.6
- 9.6 Microcavities in Photonic Crystals with Two-Dimensional Periodicity / 9.8
- 9.7 Waveguides / 9.12
- 9.8 Conclusion / 9.17
- 9.9 References / 9.18

Part 2. Nonlinear Optics
Chapter 10. Nonlinear Optics *Chung L. Tang* **10.3**

- 10.1 Glossary / 10.3
- 10.2 Introduction / 10.4
- 10.3 Basic Concepts / 10.5
- 10.4 Material Considerations / 10.19
- 10.5 Appendix / 10.21
- 10.6 References / 10.23

Chapter 11. Coherent Optical Transients *Paul R. Berman and Duncan G. Steel* **11.1**

- 11.1 Glossary / 11.1
- 11.2 Introduction / 11.2
- 11.3 Optical Bloch Equations / 11.3
- 11.4 Maxwell-Bloch Equations / 11.6
- 11.5 Free Polarization Decay / 11.7
- 11.6 Photon Echo / 11.11
- 11.7 Stimulated Photon Echo / 11.15
- 11.8 Phase Conjugate Geometry and Optical Ramsey Fringes / 11.19
- 11.9 Two-Photon Transitions and Atom Interferometry / 11.22
- 11.10 Chirped Pulse Excitation / 11.25
- 11.11 Experimental Considerations / 11.26
- 11.12 Conclusion / 11.28
- 11.13 References / 11.28

Chapter 12. Photorefractive Materials and Devices *Mark Cronin-Golomb and Marvin Klein* **12.1**

- 12.1 Introduction / 12.1
- 12.2 Materials / 12.10
- 12.3 Devices / 12.28
- 12.4 References / 12.38
- 12.5 Further Reading / 12.45

Chapter 13. Optical Limiting *David J. Hagan* 13.1

- 13.1 Introduction / 13.1
 13.2 Basic Principles of Passive Optical Limiting / 13.4
 13.3 Examples of Passive Optical Limiting in Specific Materials / 13.9
 13.4 References / 13.13

Chapter 14. Electromagnetically Induced Transparency *Jonathan P. Marangos and Thomas Halfmann* 14.1

- 14.1 Glossary / 14.1
 14.2 Introduction / 14.2
 14.3 Coherence in Two- and Three-Level Atomic Systems / 14.4
 14.4 The Basic Physical Concept of Electromagnetically Induced Transparency / 14.5
 14.5 Manipulation of Optical Properties by Electromagnetically Induced Transparency / 14.10
 14.6 Electromagnetically Induced Transparency, Driven by Pulsed Lasers / 14.15
 14.7 Steady State Electromagnetically Induced Transparency, Driven by CW Lasers / 14.16
 14.8 Gain without Inversion and Lasing without Inversion / 14.18
 14.9 Manipulation of the Index of Refraction in Dressed Atoms / 14.19
 14.10 Pulse Propagation Effects / 14.20
 14.11 Ultraslow Light Pulses / 14.22
 14.12 Nonlinear Optical Frequency Conversion / 14.24
 14.13 Nonlinear Optics at Maximal Atomic Coherence / 14.28
 14.14 Nonlinear Optics at the Few Photon Level / 14.32
 14.15 Electromagnetically Induced Transparency in Solids / 14.33
 14.16 Conclusion / 14.36
 14.17 Further Reading / 14.36
 14.18 References / 14.37

Chapter 15. Stimulated Raman and Brillouin Scattering *John Reintjes and Mark Bashkansky* 15.1

- 15.1 Introduction / 15.1
 15.2 Raman Scattering / 15.1
 15.3 Stimulated Brillouin Scattering / 15.43
 15.4 References / 15.54
 15.5 Additional References / 15.60

Chapter 16. Third-Order Optical Nonlinearities *Mansoor Sheik-Bahae and Michael P. Hasselbeck* 16.1

- 16.1 Introduction / 16.1
 16.2 Quantum Mechanical Picture / 16.4
 16.3 Nonlinear Absorption and Nonlinear Refraction / 16.7
 16.4 Kramers-Kronig Dispersion Relations / 16.9
 16.5 Optical Kerr Effect / 16.11
 16.6 Third-Harmonic Generation / 16.14
 16.7 Stimulated Scattering / 16.14
 16.8 Two-Photon Absorption / 16.19
 16.9 Effective Third-Order Nonlinearities; Cascaded $\chi^1: \chi^1$ Processes / 16.20
 16.10 Effective Third-Order Nonlinearities; Cascaded $\chi^{(2)}: \chi^{(2)}$ Processes / 16.22
 16.11 Propagation Effects / 16.24
 16.12 Common Experimental Techniques and Applications / 16.26
 16.13 References / 16.31

Chapter 17. Continuous-Wave Optical Parametric Oscillators *Majid Ebrahim-Zadeh* 17.1

- 17.1 Introduction / 17.1
 17.2 Continuous-Wave Optical Parametric Oscillators / 17.2

- 17.3 Applications / 17.21
 17.4 Summary / 17.29
 17.5 References / 17.31

Chapter 18. Nonlinear Optical Processes for Ultrashort Pulse Generation *Uwe Siegner and Ursula Keller* 18.1

- 18.1 Glossary / 18.1
 18.2 Abbreviations / 18.3
 18.3 Introduction / 18.3
 18.4 Saturable Absorbers: Macroscopic Description / 18.5
 18.5 Kerr Effect / 18.11
 18.6 Semiconductor Ultrafast Nonlinearities: Microscopic Processes / 18.15
 18.7 References / 18.23

Chapter 19. Laser-Induced Damage to Optical Materials *Marion J. Soileau* 19.1

- 19.1 Introduction / 19.1
 19.2 Practical Estimates / 19.2
 19.3 Surface Damage / 19.2
 19.4 Package-Induced Damage / 19.4
 19.5 Nonlinear Optical Effects / 19.5
 19.6 Avoidance of Damage / 19.5
 19.7 Fundamental Mechanisms / 19.6
 19.8 Progress in Measurements of Critical NLO Parameters / 19.9
 19.9 References / 19.11

Part 3. Quantum and Molecular Optics

Chapter 20. Laser Cooling and Trapping of Atoms *Harold J. Metcalf and Peter van der Straten* 20.3

- 20.1 Introduction / 20.3
 20.2 General Properties Concerning Laser Cooling / 20.4
 20.3 Theoretical Description / 20.6
 20.4 Slowing Atomic Beams / 20.11
 20.5 Optical Molasses / 20.13
 20.6 Cooling Below the Doppler Limit / 20.17
 20.7 Trapping of Neutral Atoms / 20.21
 20.8 Applications / 20.26
 20.9 References / 20.39

Chapter 21. Strong Field Physics *Todd Ditmire* 21.1

- 21.1 Glossary / 21.1
 21.2 Introduction and History / 21.2
 21.3 Laser Technology Used in Strong Field Physics / 21.4
 21.4 Strong Field Interactions with Single Electrons / 21.5
 21.5 Strong Field Interactions with Atoms / 21.10
 21.6 Strong Field Interactions with Molecules / 21.22
 21.7 Strong Field Nonlinear Optics in Gases / 21.27
 21.8 Strong Field Interactions with Clusters / 21.31
 21.9 Strong Field Physics in Underdense Plasmas / 21.36
 21.10 Strong Field Physics at Surfaces of Overdense Plasmas / 21.46
 21.11 Applications of Strong Field Interactions with Plasmas / 21.52
 21.12 References / 21.55

**Chapter 22. Slow Light Propagation in Atomic
and Photonic Media** *Jacob B. Khurgin* 22.1

- 22.1 Glossary / 22.1
- 22.2 Introduction / 22.2
- 22.3 Atomic Resonance / 22.2
- 22.4 Bandwidth Limitations in Atomic Schemes / 22.9
- 22.5 Photonic Resonance / 22.9
- 22.6 Slow Light in Optical Fibers / 22.13
- 22.7 Conclusion / 22.15
- 22.8 References / 22.16

Chapter 23. Quantum Entanglement in Optical Interferometry
*Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver,
and Jonathan P. Dowling* 23.1

- 23.1 Introduction / 23.1
- 23.2 Shot-Noise Limit / 23.4
- 23.3 Heisenberg Limit / 23.6
- 23.4 “Digital” Approaches / 23.7
- 23.5 N00n State / 23.9
- 23.6 Quantum Imaging / 23.13
- 23.7 Toward Quantum Remote Sensing / 23.14
- 23.8 References / 23.15

Index 1.1

DO NOT DUPLICATE

CONTENTS OF VOLUME V

Part 1. Measurements

Chapter 1. Scatterometers *John C. Stover* 1.3

- 1.1 Glossary / 1.3
- 1.2 Introduction / 1.3
- 1.3 Definitions and Specifications / 1.5
- 1.4 Instrument Configurations and Component Descriptions / 1.7
- 1.5 Instrumentation Issues / 1.11
- 1.6 Measurement Issues / 1.13
- 1.7 Incident Power Measurement, System Calibration, and Error Analysis / 1.14
- 1.8 Summary / 1.16
- 1.9 References / 1.16

Chapter 2. Spectroscopic Measurements *Brian Henderson* 2.1

- 2.1 Glossary / 2.1
- 2.2 Introductory Comments / 2.2
- 2.3 Optical Absorption Measurements of Energy Levels / 2.2
- 2.4 The Homogeneous Lineshape of Spectra / 2.13
- 2.5 Absorption, Photoluminescence, and Radiative Decay Measurements / 2.19
- 2.6 References / 2.24

Part 2. Atmospheric Optics

Chapter 3. Atmospheric Optics *Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman* 3.3

- 3.1 Glossary / 3.3
- 3.2 Introduction / 3.4
- 3.3 Physical and Chemical Composition of the Standard Atmosphere / 3.6
- 3.4 Fundamental Theory of Interaction of Light with the Atmosphere / 3.11
- 3.5 Prediction of Atmospheric Optical Transmission: Computer Programs and Databases / 3.22
- 3.6 Atmospheric Optical Turbulence / 3.26
- 3.7 Examples of Atmospheric Optical Remote Sensing / 3.36
- 3.8 Meteorological Optics / 3.40
- 3.9 Atmospheric Optics and Global Climate Change / 3.43
- 3.10 Acknowledgments / 3.45
- 3.11 References / 3.45

Chapter 4. Imaging through Atmospheric Turbulence *Virendra N. Mahajan and Guang-ming Dai* 4.1

- Abstract / 4.1
- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.2
- 4.3 Long-Exposure Image / 4.3
- 4.4 Kolmogorov Turbulence and Atmospheric Coherence Length / 4.7

4.5	Application to Systems with Annular Pupils /	4.10
4.6	Modal Expansion of Aberration Function /	4.17
4.7	Covariance and Variance of Expansion Coefficients /	4.20
4.8	Angle of Arrival Fluctuations /	4.23
4.9	Aberration Variance and Approximate Strehl Ratio /	4.27
4.10	Modal Correction of Atmospheric Turbulence /	4.28
4.11	Short-Exposure Image /	4.31
4.12	Adaptive Optics /	4.35
4.13	Summary /	4.36
4.14	Acknowledgments /	4.37
4.15	References /	4.37

Chapter 5. Adaptive Optics *Robert Q. Fugate* **5.1**

5.1	Glossary /	5.1
5.2	Introduction /	5.2
5.3	The Adaptive Optics Concept /	5.2
5.4	The Nature of Turbulence and Adaptive Optics Requirements /	5.5
5.5	AO Hardware and Software Implementation /	5.21
5.6	How to Design an Adaptive Optical System /	5.38
5.7	Acknowledgments /	5.46
5.8	References /	5.47

Part 3. Modulators

Chapter 6. Acousto-Optic Devices *I-Cheng Chang* **6.1**

6.1	Glossary /	6.3
6.2	Introduction /	6.4
6.3	Theory of Acousto-Optic Interaction /	6.5
6.4	Acousto-Optic Materials /	6.16
6.5	Acousto-Optic Deflector /	6.22
6.6	Acousto-Optic Modulator /	6.31
6.7	Acousto-Optic Tunable Filter /	6.35
6.8	References /	6.45

Chapter 7. Electro-Optic Modulators *Georgeanne M. Purvinis and Theresa A. Maldonado* **7.1**

7.1	Glossary /	7.1
7.2	Introduction /	7.3
7.3	Crystal Optics and the Index Ellipsoid /	7.3
7.4	The Electro-Optic Effect /	7.6
7.5	Modulator Devices /	7.16
7.6	Applications /	7.36
7.7	Appendix: Euler Angles /	7.39
7.8	References /	7.40

Chapter 8. Liquid Crystals *Sebastian Gauza and Shin-Tson Wu* **8.1**

	Abstract /	8.1
8.1	Glossary /	8.1
8.2	Introduction to Liquid Crystals /	8.2
8.3	Types of Liquid Crystals /	8.4
8.4	Liquid Crystals Phases /	8.8
8.5	Physical Properties /	8.13
8.6	Liquid Crystal Cells /	8.25
8.7	Liquid Crystals Displays /	8.29
8.8	Polymer/Liquid Crystal Composites /	8.36

- 8.9 Summary / 8.37
 8.10 References / 8.38
 8.11 Bibliography / 8.39

Part 4. Fiber Optics

Chapter 9. Optical Fiber Communication Technology and System Overview *Ira Jacobs* 9.3

- 9.1 Introduction / 9.3
 9.2 Basic Technology / 9.4
 9.3 Receiver Sensitivity / 9.8
 9.4 Bit Rate and Distance Limits / 9.12
 9.5 Optical Amplifiers / 9.13
 9.6 Fiber-Optic Networks / 9.14
 9.7 Analog Transmission on Fiber / 9.15
 9.8 Technology and Applications Directions / 9.17
 9.9 References / 9.17

Chapter 10. Nonlinear Effects in Optical Fibers *John A. Buck* 10.1

- 10.1 Key Issues in Nonlinear Optics in Fibers / 10.1
 10.2 Self- and Cross-Phase Modulation / 10.3
 10.3 Stimulated Raman Scattering / 10.4
 10.4 Stimulated Brillouin Scattering / 10.7
 10.5 Four-Wave Mixing / 10.9
 10.6 Conclusion / 10.11
 10.7 References / 10.12

Chapter 11. Photonic Crystal Fibers *Philip St. J. Russell and Greg J. Pearce* 11.1

- 11.1 Glossary / 11.1
 11.2 Introduction / 11.2
 11.3 Brief History / 11.2
 11.4 Fabrication Techniques / 11.4
 11.5 Modeling and Analysis / 11.6
 11.6 Characteristics of Photonic Crystal Cladding / 11.7
 11.7 Linear Characteristics of Guidance / 11.11
 11.8 Nonlinear Characteristics of Guidance / 11.22
 11.9 Intrafiber Devices, Cutting, and Joining / 11.26
 11.10 Conclusions / 11.28
 11.11 Appendix / 11.28
 11.12 References / 11.28

Chapter 12. Infrared Fibers *James A. Harrington* 12.1

- 12.1 Introduction / 12.1
 12.2 Nonoxide and Heavy-Metal Oxide Glass IR Fibers / 12.3
 12.3 Crystalline Fibers / 12.7
 12.4 Hollow Waveguides / 12.10
 12.5 Summary and Conclusions / 12.13
 12.6 References / 12.13

Chapter 13. Sources, Modulators, and Detectors for Fiber Optic Communication Systems *Elsa Garmire* 13.1

- 13.1 Introduction / 13.1
 13.2 Double Heterostructure Laser Diodes / 13.3

13.3	Operating Characteristics of Laser Diodes / 13.8	
13.4	Transient Response of Laser Diodes / 13.13	
13.5	Noise Characteristics of Laser Diodes / 13.18	
13.6	Quantum Well and Strained Lasers / 13.24	
13.7	Distributed Feedback and Distributed Bragg Reflector Lasers / 13.28	
13.8	Tunable Lasers / 13.32	
13.9	Light-Emitting Diodes / 13.36	
13.10	Vertical Cavity Surface-Emitting Lasers / 13.42	
13.11	Lithium Niobate Modulators / 13.48	
13.12	Electroabsorption Modulators / 13.55	
13.13	Electro-Optic and Electrorefractive Modulators / 13.61	
13.14	<i>PIN</i> Diodes / 13.63	
13.15	Avalanche Photodiodes, MSM Detectors, and Schottky Diodes / 13.71	
13.16	References / 13.74	
Chapter 14. Optical Fiber Amplifiers <i>John A. Buck</i>		14.1
<hr/>		
14.1	Introduction / 14.1	
14.2	Rare-Earth-Doped Amplifier Configuration and Operation / 14.2	
14.3	EDFA Physical Structure and Light Interactions / 14.4	
14.4	Other Rare-Earth Systems / 14.7	
14.5	Raman Fiber Amplifiers / 14.8	
14.6	Parametric Amplifiers / 14.10	
14.7	References / 14.11	
Chapter 15. Fiber Optic Communication Links (Telecom, Datacom, and Analog) <i>Casimer DeCusatis and Guifang Li</i>		15.1
<hr/>		
15.1	Figures of Merit / 15.2	
15.2	Link Budget Analysis: Installation Loss / 15.6	
15.3	Link Budget Analysis: Optical Power Penalties / 15.8	
15.4	References / 15.18	
Chapter 16. Fiber-Based Couplers <i>Daniel Nolan</i>		16.1
<hr/>		
16.1	Introduction / 16.1	
16.2	Achromaticity / 16.3	
16.3	Wavelength Division Multiplexing / 16.4	
16.4	$1 \times N$ Power Splitters / 16.4	
16.5	Switches and Attenuators / 16.4	
16.6	Mach-Zehnder Devices / 16.4	
16.7	Polarization Devices / 16.5	
16.8	Summary / 16.6	
16.9	References / 16.6	
Chapter 17. Fiber Bragg Gratings <i>Kenneth O. Hill</i>		17.1
<hr/>		
17.1	Glossary / 17.1	
17.2	Introduction / 17.1	
17.3	Photosensitivity / 17.2	
17.4	Properties of Bragg Gratings / 17.3	
17.5	Fabrication of Fiber Gratings / 17.4	
17.6	The Application of Fiber Gratings / 17.8	
17.7	References / 17.9	
Chapter 18. Micro-Optics-Based Components for Networking <i>Joseph C. Palais</i>		18.1
<hr/>		
18.1	Introduction / 18.1	
18.2	Generalized Components / 18.1	

- 18.3 Network Functions / 18.2
 18.4 Subcomponents / 18.5
 18.5 Components / 18.9
 18.6 References / 18.12

Chapter 19. Semiconductor Optical Amplifiers *Jay M. Wiesenfeld and Leo H. Spiekman* **19.1**

- 19.1 Introduction / 19.1
 19.2 Device Basics / 19.2
 19.3 Fabrication / 19.15
 19.4 Device Characterization / 19.17
 19.5 Applications / 19.22
 19.6 Amplification of Signals / 19.22
 19.7 Switching and Modulation / 19.28
 19.8 Nonlinear Applications / 19.29
 19.9 Final Remarks / 19.36
 19.10 References / 19.36

Chapter 20. Optical Time-Division Multiplexed Communication Networks *Peter J. Delfyett* **20.1**

- 20.1 Glossary / 20.1
 20.2 Introduction / 20.3
 20.3 Multiplexing and Demultiplexing / 20.3
 20.4 Introduction to Device Technology / 20.12
 20.5 Summary and Future Outlook / 20.24
 20.6 Bibliography / 20.25

Chapter 21. WDM Fiber-Optic Communication Networks *Alan E. Willner, Changyuan Yu, Zhongqi Pan, and Yong Xie* **21.1**

- 21.1 Introduction / 21.1
 21.2 Basic Architecture of WDM Networks / 21.4
 21.3 Fiber System Impairments / 21.13
 21.4 Optical Modulation Formats for WDM Systems / 21.27
 21.5 Optical Amplifiers in WDM Networks / 21.37
 21.6 Summary / 21.44
 21.7 Acknowledgments / 21.44
 21.8 References / 21.44

Chapter 22. Solitons in Optical Fiber Communication Systems *Pavel V. Mamyshev* **22.1**

- 22.1 Introduction / 22.1
 22.2 Nature of the Classical Soliton / 22.2
 22.3 Properties of Solitons / 22.4
 22.4 Classical Soliton Transmission Systems / 22.5
 22.5 Frequency-Guiding Filters / 22.7
 22.6 Sliding Frequency-Guiding Filters / 22.8
 22.7 Wavelength Division Multiplexing / 22.9
 22.8 Dispersion-Managed Solitons / 22.12
 22.9 Wavelength-Division Multiplexed Dispersionmanaged Soliton Transmission / 22.15
 22.10 Conclusion / 22.17
 22.11 References / 22.17

Chapter 23. Fiber-Optic Communication Standards 23.1
Casimer DeCusatis

- 23.1 Introduction / 23.1
- 23.2 ESCON / 23.1
- 23.3 FDDI / 23.2
- 23.4 Fibre Channel Standard / 23.4
- 23.5 ATM/SONET / 23.6
- 23.6 Ethernet / 23.7
- 23.7 Infiniband / 23.8
- 23.8 References / 23.8

Chapter 24. Optical Fiber Sensors 24.1
Richard O. Claus, Ignacio Matias, and Francisco Arregui

- 24.1 Introduction / 24.1
- 24.2 Extrinsic Fabry-Perot Interferometric Sensors / 24.2
- 24.3 Intrinsic Fabry-Perot Interferometric Sensors / 24.4
- 24.4 Fiber Bragg Grating Sensors / 24.5
- 24.5 Long-Period Grating Sensors / 24.8
- 24.6 Comparison of Sensing Schemes / 24.13
- 24.7 Conclusion / 24.13
- 24.8 References / 24.13
- 24.9 Further Reading / 24.14

Chapter 25. High-Power Fiber Lasers and Amplifiers 25.1
Timothy S. McComb, Martin C. Richardson, and Michael Bass

- 25.1 Glossary / 25.1
- 25.2 Introduction / 25.3
- 25.3 Fiber Laser Limitations / 25.6
- 25.4 Fiber Laser Fundamentals / 25.7
- 25.5 Fiber Laser Architectures / 25.9
- 25.6 LMA Fiber Designs / 25.18
- 25.7 Active Fiber Dopants / 25.22
- 25.8 Fiber Fabrication and Materials / 25.26
- 25.9 Spectral and Temporal Modalities / 25.29
- 25.10 Conclusions / 25.33
- 25.11 References / 25.33

PART 5. X-Ray and Neutron Optics

SUBPART 5.1. INTRODUCTION AND APPLICATIONS

Chapter 26. An Introduction to X-Ray and Neutron Optics 26.5
Carolyn MacDonald

- 26.1 History / 26.5
- 26.2 X-Ray Interaction with Matter / 26.6
- 26.3 Optics Choices / 26.7
- 26.4 Focusing and Collimation / 26.9
- 26.5 References / 26.11

Chapter 27. Coherent X-Ray Optics and Microscopy 27.1
Qun Shen

- 27.1 Glossary / 27.1
- 27.2 Introduction / 27.2

- 27.3 Fresnel Wave Propagation / 27.2
- 27.4 Unified Approach for Near- and Far-Field Diffraction / 27.2
- 27.5 Coherent Diffraction Microscopy / 27.4
- 27.6 Coherence Preservation in X-Ray Optics / 27.5
- 27.7 References / 27.5

Chapter 28. Requirements for X-Ray Diffraction *Scott T. Misture* **28.1**

- 28.1 Introduction / 28.1
- 28.2 Slits / 28.1
- 28.3 Crystal Optics / 28.3
- 28.4 Multilayer Optics / 28.5
- 28.5 Capillary and Polycapillary Optics / 28.5
- 28.6 Diffraction and Fluorescence Systems / 28.5
- 28.7 X-Ray Sources and Microsources / 28.7
- 28.8 References / 28.7

Chapter 29. Requirements for X-Ray Fluorescence *Walter Gibson and George Havrilla* **29.1**

- 29.1 Introduction / 29.1
- 29.2 Wavelength-Dispersive X-Ray Fluorescence (WDXRF) / 29.2
- 29.3 Energy-Dispersive X-Ray Fluorescence (EDXRF) / 29.3
- 29.4 References / 29.12

Chapter 30. Requirements for X-Ray Spectroscopy *Dirk Lützenkirchen-Hecht and Ronald Frahm* **30.1**

- 30.1 References / 30.5

Chapter 31. Requirements for Medical Imaging and X-Ray Inspection *Douglas Pfeiffer* **31.1**

- 31.1 Introduction to Radiography and Tomography / 31.1
- 31.2 X-Ray Attenuation and Image Formation / 31.1
- 31.3 X-Ray Detectors and Image Receptors / 31.4
- 31.4 Tomography / 31.5
- 31.5 Computed Tomography / 31.5
- 31.6 Digital Tomosynthesis / 31.7
- 31.7 Digital Displays / 31.8
- 31.8 Conclusion / 31.9
- 31.9 References / 31.10

Chapter 32. Requirements for Nuclear Medicine *Lars R. Furenlid* **32.1**

- 32.1 Introduction / 32.1
- 32.2 Projection Image Acquisition / 32.2
- 32.3 Information Content in SPECT / 32.3
- 32.4 Requirements for Optics For SPECT / 32.4
- 32.5 References / 32.4

Chapter 33. Requirements for X-Ray Astronomy *Scott O. Rohrbach* **33.1**

- 33.1 Introduction / 33.1
- 33.2 Trade-Offs / 33.2
- 33.3 Summary / 33.4

Chapter 34. Extreme Ultraviolet Lithography *Franco Cerrina and Fan Jiang* 34.1

- 34.1 Introduction / 34.1
- 34.2 Technology / 34.2
- 34.3 Outlook / 34.5
- 34.4 Acknowledgments / 34.6
- 34.5 References / 34.7

Chapter 35. Ray Tracing of X-Ray Optical Systems *Franco Cerrina and Manuel Sanchez del Rio* 35.1

- 35.1 Introduction / 35.1
- 35.2 The Conceptual Basis of SHADOW / 35.2
- 35.3 Interfaces and Extensions of SHADOW / 35.3
- 35.4 Examples / 35.4
- 35.5 Conclusions and Future / 35.5
- 35.6 References / 35.6

Chapter 36. X-Ray Properties of Materials *Eric M. Gullikson* 36.1

- 36.1 X-Ray and Neutron Optics / 36.2
- 36.2 Electron Binding Energies, Principal K- and L-Shell Emission Lines, and Auger Electron Energies / 36.3
- 36.3 References / 36.10

SUBPART 5.2. REFRACTIVE AND INTERFERENCE OPTICS

Chapter 37. Refractive X-Ray Lenses *Bruno Lengeler and Christian G. Schroer* 37.3

- 37.1 Introduction / 37.3
- 37.2 Refractive X-Ray Lenses with Rotationally Parabolic Profile / 37.4
- 37.3 Imaging with Parabolic Refractive X-Ray Lenses / 37.6
- 37.4 Microfocusing with Parabolic Refractive X-Ray Lenses / 37.7
- 37.5 Prefocusing and Collimation with Parabolic Refractive X-Ray Lenses / 37.8
- 37.6 Nanofocusing Refractive X-Ray Lenses / 37.8
- 37.7 Conclusion / 37.11
- 37.8 References / 37.11

Chapter 38. Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region *Malcolm R. Howells* 38.1

- 38.1 Introduction / 38.1
- 38.2 Diffraction Properties / 38.1
- 38.3 Focusing Properties / 38.3
- 38.4 Dispersion Properties / 38.6
- 38.5 Resolution Properties / 38.7
- 38.6 Efficiency / 38.8
- 38.7 References / 38.8

Chapter 39. Crystal Monochromators and Bent Crystals *Peter Siddons* 39.1

- 39.1 Crystal Monochromators / 39.1
- 39.2 Bent Crystals / 39.5
- 39.3 References / 39.6

Chapter 40. Zone Plates *Alan Michette* **40.1**

- 40.1 Introduction / 40.1
- 40.2 Geometry of a Zone Plate / 40.1
- 40.3 Zone Plates as Thin Lenses / 40.3
- 40.4 Diffraction Efficiencies of Zone Plates / 40.4
- 40.5 Manufacture of Zone Plates / 40.8
- 40.6 Bragg-Fresnel Lenses / 40.9
- 40.7 References / 40.10

Chapter 41. Multilayers *Eberhard Spiller* **41.1**

- 41.1 Glossary / 41.1
- 41.2 Introduction / 41.1
- 41.3 Calculation of Multilayer Properties / 41.3
- 41.4 Fabrication Methods and Performance / 41.4
- 41.5 Multilayers for Diffractive Imaging / 41.9
- 41.6 References / 41.10

Chapter 42. Nanofocusing of Hard X-Rays with Multilayer Laue Lenses *Albert T. Macrander, Hanfei Yan, Hyon Chol Kang, Jörg Maser, Chian Liu, Ray Conley, and G. Brian Stephenson* **42.1**

- Abstract / 42.1
- 42.1 Introduction / 42.2
- 42.2 MLL Concept and Volume Diffraction Calculations / 42.4
- 42.3 Magnetron-Sputtered MLLs / 42.5
- 42.4 Instrumental Beamline Arrangement and Measurements / 42.9
- 42.5 Takagi-Taupin Calculations / 42.12
- 42.6 Wedged MLLs / 42.12
- 42.7 MMLs with Curved Interfaces / 42.14
- 42.8 MLL Prospects / 42.15
- 42.9 Summary / 42.17
- 42.10 Acknowledgments / 42.17
- 42.11 References / 42.18

Chapter 43. Polarizing Crystal Optics *Qun Shen* **43.1**

- 43.1 Introduction / 43.1
- 43.2 Linear Polarizers / 43.2
- 43.3 Linear Polarization Analyzers / 43.4
- 43.4 Phase Plates for Circular Polarization / 43.5
- 43.5 Circular Polarization Analyzers / 43.6
- 43.6 Acknowledgments / 43.8
- 43.7 References / 43.8

SUBPART 5.3. REFLECTIVE OPTICS

Chapter 44. Image Formation with Grazing Incidence Optics *James E. Harvey* **44.3**

- 44.1 Glossary / 44.3
- 44.2 Introduction to X-Ray Mirrors / 44.3
- 44.3 Optical Design and Residual Aberrations of Grazing Incidence Telescopes / 44.6
- 44.4 Image Analysis for Grazing Incidence X-Ray Optics / 44.12
- 44.5 Validation of Image Analysis for Grazing Incidence X-Ray Optics / 44.16
- 44.6 References / 44.18

Chapter 45. Aberrations for Grazing Incidence Optics	45.1
<i>Timo T. Saha</i>	
45.1 Grazing Incidence Telescopes / 45.1	
45.2 Surface Equations / 45.1	
45.3 Transverse Ray Aberration Expansions / 45.3	
45.4 Curvature of the Best Focal Surface / 45.5	
45.5 Aberration Balancing / 45.5	
45.6 On-Axis Aberrations / 45.6	
45.7 References / 45.8	
Chapter 46. X-Ray Mirror Metrology	46.1
<i>Peter Z. Takacs</i>	
46.1 Glossary / 46.1	
46.2 Introduction / 46.1	
46.3 Surface Finish Metrology / 46.2	
46.4 Surface Figure Metrology / 46.3	
46.5 Practical Profile Analysis Considerations / 46.6	
46.6 References / 46.12	
Chapter 47. Astronomical X-Ray Optics	47.1
<i>Marshall K. Joy and Brian D. Ramsey</i>	
47.1 Introduction / 47.1	
47.2 Wolter X-Ray Optics / 47.2	
47.3 Kirkpatrick-Baez Optics / 47.7	
47.4 Hard X-Ray Optics / 47.9	
47.5 Toward Higher Angular Resolution / 47.10	
47.6 References / 47.11	
Chapter 48. Multifoil X-Ray Optics	48.1
<i>Ladislav Pina</i>	
48.1 Introduction / 48.1	
48.2 Grazing Incidence Optics / 48.1	
48.3 Multifoil Lobster-Eye Optics / 48.2	
48.4 Multifoil Kirkpatrick-Baez Optics / 48.3	
48.5 Summary / 48.4	
48.6 References / 48.4	
Chapter 49. Pore Optics	49.1
<i>Marco W. Beijersbergen</i>	
49.1 Introduction / 49.1	
49.2 Glass Micropore Optics / 49.1	
49.3 Silicon Pore Optics / 49.6	
49.4 Micromachined Silicon / 49.7	
49.5 References / 49.7	
Chapter 50. Adaptive X-Ray Optics	50.1
<i>Ali Khounsary</i>	
50.1 Introduction / 50.1	
50.2 Adaptive Optics in X-Ray Astronomy / 50.2	
50.3 Active and Adaptive Optics for Synchrotron- and Lab-Based X-Ray Sources / 50.2	
50.4 Conclusions / 50.8	
50.5 References / 50.8	
Chapter 51. The Schwarzschild Objective	51.1
<i>Franco Cerrina</i>	
51.1 Introduction / 51.1	
51.2 Applications to X-Ray Domain / 51.3	
51.3 References / 51.5	

Chapter 52. Single Capillaries	Donald H. Bilderback and Sterling W. Cornaby	52.1
<hr/>		
52.1	Background / 52.1	
52.2	Design Parameters / 52.1	
52.3	Fabrication / 52.4	
52.4	Applications of Single-Bounce Capillary Optics / 52.5	
52.5	Applications of Condensing Capillary Optics / 52.6	
52.6	Conclusions / 52.6	
52.7	Acknowledgments / 52.6	
52.8	References / 52.6	
Chapter 53. Polycapillary X-Ray Optics	Carolyn MacDonald and Walter Gibson	53.1
<hr/>		
53.1	Introduction / 53.1	
53.2	Simulations and Defect Analysis / 53.3	
53.3	Radiation Resistance / 53.5	
53.4	Alignment and Measurement / 53.5	
53.5	Collimation / 53.8	
53.6	Focusing / 53.9	
53.7	Applications / 53.10	
53.8	Summary / 53.19	
53.9	Acknowledgments / 53.19	
53.10	References / 53.19	
SUBPART 5.4. X-RAY SOURCES		
Chapter 54. X-Ray Tube Sources	Susanne M. Lee and Carolyn MacDonald	54.3
<hr/>		
54.1	Introduction / 54.3	
54.2	Spectra / 54.4	
54.3	Cathode Design and Geometry / 54.10	
54.4	Effect of Anode Material, Geometry, and Source Size on Intensity and Brightness / 54.11	
54.5	General Optimization / 54.15	
54.6	References / 54.17	
Chapter 55. Synchrotron Sources	Steven L. Hulbert and Gwyn P. Williams	55.1
<hr/>		
55.1	Introduction / 55.1	
55.2	Theory of Synchrotron Radiation Emission / 55.2	
55.3	Insertion Devices (Undulators and Wigglers) / 55.9	
55.4	Coherence of Synchrotron Radiation Emission in the Long Wavelength Limit / 55.17	
55.5	Conclusion / 55.20	
55.6	References / 55.20	
Chapter 56. Laser-Generated Plasmas	Alan Michette	56.1
<hr/>		
56.1	Introduction / 56.1	
56.2	Characteristic Radiation / 56.2	
56.3	Bremsstrahlung / 56.8	
56.4	Recombination Radiation / 56.10	
56.5	References / 56.10	
Chapter 57. Pinch Plasma Sources	Victor Kantsyrev	57.1
<hr/>		
57.1	Introduction / 57.1	
57.2	Types of Z-Pinch Radiation Sources / 57.2	

- 57.3 Choice of Optics for Z-Pinch Sources / 57.4
 57.4 References / 57.5

Chapter 58. X-Ray Lasers *Greg Tallents* 58.1

- 58.1 Free-Electron Lasers / 58.1
 58.2 High Harmonic Production / 58.2
 58.3 Plasma-Based EUV Lasers / 58.2
 58.4 References / 58.4

Chapter 59. Inverse Compton X-Ray Sources *Frank Carroll* 59.1

- 59.1 Introduction / 59.1
 59.2 Inverse Compton Calculations / 59.2
 59.3 Practical Devices / 59.2
 59.4 Applications / 59.3
 59.5 Industrial/Military/Crystallographic Uses / 59.4
 59.6 References / 59.4

SUBPART 5.5. X-RAY DETECTORS

Chapter 60. Introduction to X-Ray Detectors *Walter Gibson and Peter Siddons* 60.3

- 60.1 Introduction / 60.3
 60.2 Detector Type / 60.3
 60.3 Summary / 60.9
 60.4 References / 60.10

Chapter 61. Advances in Imaging Detectors *Aaron Couture* 61.1

- 61.1 Introduction / 61.1
 61.2 Flat-Panel Detectors / 61.3
 61.3 CCD Detectors / 61.7
 61.4 Conclusion / 61.8
 61.5 References / 61.8

Chapter 62. X-Ray Spectral Detection and Imaging *Eric Lifshin* 62.1

- 62.1 References / 62.6

SUBPART 5.6. NEUTRON OPTICS AND APPLICATIONS

Chapter 63. Neutron Optics *David Mildner* 63.3

- 63.1 Neutron Physics / 63.3
 63.2 Scattering Lengths and Cross Sections / 63.5
 63.3 Neutron Sources / 63.12
 63.4 Neutron Optical Devices / 63.15
 63.5 Refraction and Reflection / 63.19
 63.6 Diffraction and Interference / 63.23
 63.7 Polarization Techniques / 63.27
 63.8 Neutron Detection / 63.31
 63.9 References / 63.35

Chapter 64. Grazing-Incidence Neutron Optics *Mikhail Gubarev
and Brian Ramsey* **64.1**

- 64.1 Introduction / 64.1
64.2 Total External Reflection / 64.1
64.3 Diffractive Scattering and Mirror Surface Roughness Requirements / 64.2
64.4 Imaging Focusing Optics / 64.3
64.5 References / 64.7

Index **I.1**

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

CONTRIBUTORS TO HANDBOOK OF OPTICS

Jan P. Allebach *Electronic Imaging Systems Laboratory, School of Electrical and Computer Engineering, Purdue University, West Lafayette, Indiana* (VOL. III, CHAP. 24, “Human Vision and Electronic Imaging”)

Joseph H. Altman *Institute of Optics, University of Rochester, Rochester, New York* (VOL. II, CHAP. 29, “Photographic Films”)

Paul M. Amirtharaj *Sensors and Electron Devices Directorate, U.S. Army Research Laboratory, Adelphi, Maryland* (VOL. IV, CHAP. 5, “Optical Properties of Semiconductors”)

Jeffrey Anshel *Corporate Vision Consulting, Encinitas, California* (VOL. III, CHAP. 23, “Vision Problems at Computers”)

Ladan Arissian *Texas A&M University, College Station, Texas, and National Research Council of Canada, Ottawa, Ontario, Canada* (VOL. II, CHAP. 20, “Ultrashort Optical Sources and Applications”)

Francisco Arregui *Public University Navarra, Pamplona, Spain* (VOL. V, CHAP. 24, “Optical Fiber Sensors”)

Rasheed M. A. Azzam *Department of Electrical Engineering, University of New Orleans, New Orleans, Louisiana* (VOL. I, CHAP. 16, “Ellipsometry”)

John D. Baloga *Imaging Materials and Media, Eastman Kodak Company, Rochester, New York* (VOL. II, CHAP. 30, “Photographic Materials”)

Martin S. Banks *School of Optometry, University of California, Berkeley, California* (VOL. III, CHAP. 2, “Visual Performance”)

Mark Bashkansky *Optical Sciences Division, Naval Research Laboratory, Washington, D.C.* (VOL. IV, CHAP. 15, “Stimulated Raman and Brillouin Scattering”)

Michael Bass *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. V, CHAP. 25, “High-Power Fiber Lasers and Amplifiers”)

Marco W. Beijersbergen *Cosine Research B.V./Cosine Science & Computing B.V., Leiden University, Leiden, Netherlands* (VOL. V, CHAP. 49, “Pore Optics”)

Leo Beiser *Consultant, Flushing, New York* (VOL. I, CHAP. 30, “Scanners”)

Edward S. Bennett *College of Optometry, University of Missouri, St. Louis, Missouri* (VOL. III, CHAP. 20, “Optics of Contact Lenses”)

Gisele Bennett *Electro-Optical Systems Laboratory and School of Electrical and Computer Engineering, Georgia Tech Research Institute, Georgia Institute of Technology, Atlanta, Georgia* (VOL. I, CHAP. 6, “Coherence Theory: Tools and Applications”)

Jean M. Bennett* *Research Department, Michelson Laboratory, Naval Air Warfare Center, China Lake, California* (VOL. I, CHAPS. 12, 13, “Polarization,” “Polarizers”)

János A. Bergou *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Department of Physics and Astronomy, Hunter College of the City University of New York, New York, New York* (VOL. II, CHAP. 23, “Quantum Theory of the Laser”)

Paul R. Berman *Physics Department, University of Michigan, Ann Arbor, Michigan* (VOL. IV, CHAP. 11, “Coherent Optical Transients”)

*Deceased.

- Ellis Betensky** *Opcon Associates, Inc., West Redding, Connecticut* (VOL. I, CHAP. 27, “Camera Lenses”)
- Donald H. Bilderback** *Cornell High Energy Synchrotron Source, School of Applied and Engineering Physics, Cornell University, Ithaca, New York* (VOL. V, CHAP. 52, “Single Capillaries”)
- Craig F. Bohren** *Pennsylvania State University, University Park, Pennsylvania* (VOL. I, CHAP. 7, “Scattering by Particles”)
- Glenn D. Boreman** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. I, CHAP. 4, “Transfer Function Techniques”)
- John E. Bowers** *Department of Electrical and Computer Engineering, University of California, Santa Barbara, California* (VOL. II, CHAP. 26, “High-Speed Photodetectors”)
- Arthur Bradley** *School of Optometry, Indiana University, Bloomington, Indiana* (VOL. III, CHAP. 14, “Optics and Vision of the Aging Eye”)
- David H. Brainard** *Department of Psychology, University of Pennsylvania, Philadelphia, Pennsylvania* (VOL. III, CHAPS. 10, 11, “Colorimetry”, “Color Vision Mechanisms”)
- Robert P. Breault** *Breault Research Organization, Tucson, Arizona* (VOL. II, CHAP. 7, “Control of Stray Light”; VOL. IV, CHAP. 6, “Characterization and Use of Black Surfaces for Optical Systems”)
- John A. Buck** *Georgia Institute of Technology, School of Electrical and Computer Engineering, Atlanta, Georgia* (VOL. V, CHAPS. 10, 14, “Nonlinear Effects in Optical Fibers”, “Optical Fiber Amplifiers”)
- Stephen A. Burns** *School of Optometry, Indiana University, Bloomington, Indiana* (VOL. III, CHAP. 5, “Optical Generation of the Visual Stimulus”)
- Frank Carroll** *MXISystems, Nashville, Tennessee* (VOL. V, CHAP. 59, “Inverse Compton X-Ray Sources”)
- William H. Carter*** *Naval Research Laboratory, Washington, D.C.* (VOL. I, CHAP. 5, “Coherence Theory”)
- William Cassarly** *Optical Research Associates, Pasadena, California* (VOL. II, CHAP. 39, “Nonimaging Optics: Concentration and Illumination”)
- Franco Cerrina** *Department of Electrical and Computer Engineering, University of Wisconsin, Madison, Wisconsin* (VOL. V, CHAPS. 34, 35, 51, “Extreme Ultraviolet Lithography”, “Ray Tracing of X-Ray Optical Systems”, “The Schwarzschild Objective”)
- I-Cheng Chang** *Accord Optics, Sunnyvale, California* (VOL. V, CHAP. 6, “Acousto-Optic Devices”)
- Zenghu Chang** *Department of Physics, Kansas State University, Cardwell Hall, Manhattan, Kansas* (VOL. II, CHAP. 21, “Attosecond Optics”)
- Neil Charman** *Department of Optometry and Vision Sciences, University of Manchester, Manchester, United Kingdom* (VOL. III, CHAP. 1, “Optics of the Eye”)
- Russell A. Chipman** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. I, CHAPS. 14, 15, “Mueller Matrices”, “Polarimetry”)
- B. Ralph Chou** *School of Optometry, University of Waterloo, Waterloo, Ontario, Canada* (VOL. III, CHAP. 12, “Assessment of Refraction and Refractive Errors and Their Influence on Optical Design”)
- Eugene L. Church** *Instrumentation Division, Brookhaven National Laboratory, Upton, New York* (VOL. I, CHAP. 8, “Surface Scattering”)
- James H. Churnside** *National Oceanic and Atmospheric Administration, Earth System Research Laboratory, Boulder, Colorado* (VOL. V, CHAP. 3, “Atmospheric Optics”)
- Richard O. Claus** *Virginia Tech, Blacksburg, Virginia* (VOL. V, CHAP. 24, “Optical Fiber Sensors”)
- Ray Conley** *X-Ray Science Division, Argonne National Laboratory, Argonne, Illinois, and National Synchrotron Light Source II, Brookhaven National Laboratory, Upton, New York* (VOL. V, CHAP. 42, “Nanofocusing of Hard X-Rays with Multilayer Laue Lenses”)
- Sterling W. Cornaby** *Cornell High Energy Synchrotron Source, School of Applied and Engineering Physics, Cornell University Ithaca, New York* (VOL. V, CHAP. 52, “Single Capillaries”)

*Deceased.

- Aaron Couture** *GE Global Research Center, Niskayuna, New York* (VOL. V, CHAP. 61, “Advances in Imaging Detectors”)
- William Cowan** *Department of Computer Science, University of Waterloo, Waterloo, Ontario, Canada* (VOL. III, CHAP. 22, “Displays for Vision Research”)
- M. George Craford** *Hewlett-Packard Co., San Jose, California* (VOL. II, CHAP. 17, “Light-Emitting Diodes”)
- Katherine Creath** *Optineering, Tucson, Arizona, and College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. II, CHAP. 14, “Use of Computer-Generated Holograms in Optical Testing”)
- Mark Cronin-Golomb** *Department of Biomedical Engineering, Tufts University, Medford, Massachusetts* (VOL. IV, CHAP. 12, “Photorefractive Materials and Devices”)
- Guang-ming Dai** *Laser Vision Correction Group, Advanced Medical Optics, Milpitas, California* (VOL. V, CHAP. 4, “Imaging through Atmospheric Turbulence”)
- Johannes F. de Boer** *Department of Physics, VU University, Amsterdam, and Rotterdam Ophthalmic Institute, Rotterdam, The Netherlands* (VOL. III, CHAP. 18, “Diagnostic Use of Optical Coherence Tomography in the Eye”)
- Casimer DeCusatis** *IBM Corporation, Poughkeepsie, New York* (VOL. V, CHAPS. 15, 23, “Fiber Optic Communication Links (Telecom, Datacom, and Analog)”, “Fiber-Optic Communication Standards”)
- Peter J. Delfyett** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. V, CHAP. 20, “Optical Time-Division Multiplexed Communication Networks”)
- Pamela L. Derry** *Boeing Defense & Space Group, Seattle, Washington* (VOL. II, CHAP. 19, “Semiconductor Lasers”)
- L. Diaz-Santana** *Department of Optometry and Visual Science, City University, London, United Kingdom* (VOL. III, CHAP. 16, “Refractive Surgery, Correction of Vision, PRK, and Lasik”)
- Alain C. Diebold** *College of Nanoscale Science and Engineering, University at Albany, Albany, New York* (VOL. IV, CHAP. 5, “Optical Properties of Semiconductors”)
- Jean-Claude Diels** *Departments of Physics and Electrical Engineering, University of New Mexico, Albuquerque, New Mexico* (VOL. II, CHAP. 20, “Ultrashort Optical Sources and Applications”)
- Todd Ditmire** *Texas Center for High Intensity Laser Science, Department of Physics, The University of Texas at Austin, Austin, Texas* (VOL. IV, CHAP. 21, “Strong Field Physics”)
- Jerzy A. Dobrowolski** *Institute for Microstructural Sciences, National Research Council of Canada, Ottawa, Ontario, Canada* (VOL. IV, CHAP. 7, “Optical Properties of Films and Coatings”)
- Aristide Dogariu** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. I, CHAP. 9, “Volume Scattering in Random Media”)
- Jonathan P. Dowling** *Hearne Institute for Theoretical Physics, Department of Physics and Astronomy, Louisiana State University, Baton Rouge, Louisiana* (VOL. IV, CHAP. 23, “Quantum Entanglement in Optical Interferometry”)
- Majid Ebrahim-Zadeh** *ICFO—Institut de Ciències Fotoniques, Mediterranean Technology Park, Barcelona, Spain, and Institutio Catalana de Recerca i Estudis Avancats (ICREA), Passeig Lluis Companys, Barcelona, Spain* (VOL. IV, CHAP. 17, “Continuous-Wave Optical Parametric Oscillators”)
- Jeremy Ellis** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. I, CHAP. 9, “Volume Scattering in Random Media”)
- Berthold-Georg Englert** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Max-Planck-Institut für Quantenoptik, Garching bei München, Germany, and Abteilung Quantenphysik der Universität Ulm, Ulm, Germany* (VOL. II, CHAP. 23, “Quantum Theory of the Laser”)
- Jay M. Enoch** *School of Optometry, University of California at Berkeley, Berkeley, California* (VOL. III, CHAPS. 8, 9, “Biological Waveguides”, “The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution”)
- Chris J. Evans** *Zygo Corporation, Middlefield, Connecticut* (VOL. II, CHAP. 10, “Fabrication of Optics by Diamond Turning”)

- Bart Farell** *Institute for Sensory Research, Syracuse University, Syracuse, New York* (VOL. III, CHAP. 3, "Psychophysical Methods")
- Michael W. Farn** *MIT/Lincoln Laboratory, Lexington, Massachusetts* (VOL. I, CHAP. 23, "Binary Optics")
- Luis Figueroa** *Boeing Defense & Space Group, Seattle, Washington* (VOL. II, CHAP. 19, "Semiconductor Lasers")
- Ronald Frahm** *Bergische Universität Wuppertal, Wuppertal, Germany* (VOL. V, CHAP. 30, "Requirements for X-Ray Spectroscopy")
- Robert Q. Fugate** *Starfire Optical Range, Directed Energy Directorate, Air Force Research Laboratory, Kirtland Air Force Base, New Mexico* (VOL. V, CHAP. 5, "Adaptive Optics")
- Lars R. Furenlid** *University of Arizona, Tucson, Arizona* (VOL. V, CHAP. 32, "Requirements for Nuclear Medicine")
- Elsa Garmire** *Dartmouth College, Hanover, New Hampshire* (VOL. V, CHAP. 13, "Sources, Modulators, and Detectors for Fiber Optic Communication Systems")
- Sebastian Gauza** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. V, CHAP. 8, "Liquid Crystals")
- Wilson S. Geisler** *Department of Psychology, University of Texas, Austin, Texas* (VOL. III, CHAP. 2, "Visual Performance")
- Walter Gibson** *X-Ray Optical Systems, Inc., East Greenbush, New York* (VOL. V, CHAPS. 29, 53, 60, "Requirements for X-Ray Fluorescence", "Polycapillary X-Ray Optics", "Introduction to X-Ray Detectors")
- Harilaos Ginis** *Institute of Vision and Optics, University of Crete, Greece* (VOL. III, CHAP. 16, "Refractive Surgery, Correction of Vision, PRK, and Lasik")
- Norman Goldberg** *Madison, Wisconsin* (VOL. I, CHAP. 25, "Cameras")
- Douglas S. Goodman** *Corning Tropel Corporation, Fairport, New York* (VOL. I, CHAP. 1, "General Principles of Geometrical Optics")
- Joseph W. Goodman** *Department of Electrical Engineering, Stanford University, Stanford, California* (VOL. I, CHAP. 11, "Analog Optical Signal and Image Processing")
- John E. Greivenkamp** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. I, CHAP. 2, "Interference")
- Mikhail Gubarev** *NASA/Marshall Space Flight Center, Huntsville, Alabama* (VOL. V, CHAP. 64, "Grazing-Incidence Neutron Optics")
- Eric M. Gullikson** *Center for X-Ray Optics, Lawrence Berkeley National Laboratory, Berkeley, California* (VOL. V, CHAP. 36, "X-Ray Properties of Materials")
- Anurag Gupta** *Optical Research Associates, Tucson, Arizona* (VOL. II, CHAP. 40, "Lighting and Applications")
- David J. Hagan** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. IV, CHAP. 13, "Optical Limiting")
- Roland H. Haitz** *Hewlett-Packard Co., San Jose, California* (VOL. II, CHAP. 17, "Light-Emitting Diodes")
- Thomas Halfmann** *Institute of Applied Physics, Technical University of Darmstadt, Darmstadt, Germany* (VOL. IV, CHAP. 14, "Electromagnetically Induced Transparency")
- John L. Hall** *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (VOL. II, CHAP. 22, "Laser Stabilization")
- Parameswaran Hariharan** *School of Physics, University of Sydney, Sydney, Australia* (VOL. I, CHAP. 32, "Interferometers")
- Michael E. Harrigan** *Harrigan Optical Design, Victor, New York* (VOL. II, CHAP. 2, "Aberration Curves in Lens Design")
- James A. Harrington** *Rutgers University, Piscataway, New Jersey* (VOL. V, CHAP. 12, "Infrared Fibers")
- James E. Harvey** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. V, CHAP. 44, "Reflective Optics")

- Michael P. Hasselbeck** *Department of Physics and Astronomy, University of New Mexico, Albuquerque, New Mexico* (VOL. IV, CHAP. 16, “Third-Order Optical Nonlinearities”)
- George Havrilla** *Los Alamos National Laboratory, Los Alamos, New Mexico* (VOL. V, CHAP. 29, “Requirements for X-Ray Fluorescence”)
- Brian Henderson** *Department of Physics and Applied Physics, University of Strathclyde, Glasgow, United Kingdom* (VOL. I, CHAPS. 10, 31, “Optical Spectroscopy and Spectroscopic Lineshapes”; “Optical Spectrometers”; VOL. V, CHAP. 2, “Spectroscopic Measurements”)
- Kenneth O. Hill** *Communications Research Centre, Ottawa, Ontario, Canada, and Nu-Wave Photonics, Ottawa, Ontario, Canada* (VOL. V, CHAP. 17, “Fiber Bragg Gratings”)
- Gerald C. Holst** *JCD Publishing, Winter Park, Florida* (VOL. I, CHAP. 26, “Solid-State Cameras”)
- Chi-Shain Hong** *Boeing Defense & Space Group, Seattle, Washington* (VOL. II, CHAP. 19, “Semiconductor Lasers”)
- Malcolm R. Howells** *Advanced Light Source, Lawrence Berkeley National Laboratory, Berkeley, California* (VOL. V, CHAP. 38, “Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region”)
- Lloyd Huff** *Research Institute, University of Dayton, Dayton, Ohio* (VOL. I, CHAP. 33, “Holography and Holographic Instruments”)
- Steven L. Hulbert** *National Synchrotron Light Source, Brookhaven National Laboratory, Upton, New York* (VOL. V, CHAP. 55, “Synchrotron Sources”)
- Sean D. Huver** *Hearne Institute for Theoretical Physics, Department of Physics and Astronomy, Louisiana State University, Baton Rouge, Louisiana* (VOL. IV, CHAP. 23, “Quantum Entanglement in Optical Interferometry”)
- Ira Jacobs** *The Bradley Department of Electrical and Computer Engineering, Virginia Polytechnic Institute and State University, Blacksburg, Virginia* (VOL. V, CHAP. 9, “Optical Fiber Communication Technology and System Overview”)
- J. Christopher James** *Electro-Optical Systems Laboratory, Georgia Tech Research Institute, Georgia Institute of Technology, Atlanta, Georgia* (VOL. I, CHAP. 6, “Coherence Theory: Tools and Applications”)
- Fan Jiang** *Electrical and Computer Engineering & Center for Nano Technology, University of Wisconsin, Madison* (VOL. V, CHAP. 34, “Extreme Ultraviolet Lithography”)
- C. Bruce Johnson** *Johnson Scientific Group, Inc., Phoenix, Arizona* (VOL. II, CHAP. 31, “Image Tube Intensified Electronic Imaging”)
- R. Barry Johnson** *Consultant, Huntsville, Alabama* (VOL. I, CHAPS. 17, 30, “Lenses”; “Scanners”)
- Lloyd Jones** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. I, CHAP. 29, “Reflective and Catadioptric Objectives”)
- Abhay M. Joshi** *Discovery Semiconductors, Inc., Cranbury, New Jersey* (VOL. II, CHAP. 25, “Photodetection”)
- Marshall K. Joy** *National Aeronautics and Space Administration, Marshall Space Flight Center, Huntsville, Alabama* (VOL. V, CHAP. 47, “Astronomical X-Ray Optics”)
- Hyon Chol Kang** *Materials Science Division, Argonne National Laboratory, Argonne, Illinois, and Advanced Materials Engineering Department, Chosun University, Gwangju, Republic of Korea* (VOL. V, CHAP. 42, “Nanofocusing of Hard X-Rays with Multilayer Laue Lenses”)
- Victor Kantsyrev** *Physics Department, University of Nevada, Reno, Nevada* (VOL. V, CHAP. 57, “Pinch Plasma Sources”)
- Ursula Keller** *Institute of Quantum Electronics, Physics Department, Swiss Federal Institute of Technology (ETH), Zurich, Switzerland* (VOL. IV, CHAP. 18, “Nonlinear Optical Processes for Ultrashort Pulse Generation”)
- Ali Khounsary** *Argonne National Laboratory, Argonne, Illinois* (VOL. V, CHAP. 50, “Adaptive X-Ray Optics”)
- Jacob B. Khurgin** *Department of Electrical and Computer Engineering, Johns Hopkins University, Baltimore, Maryland* (VOL. IV, CHAP. 22, “Slow Light Propagation in Atomic and Photonic Media”)

- Dennis K. Killinger** *Center for Laser Atmospheric Sensing, Department of Physics, University of South Florida, Tampa, Florida* (VOL. V, CHAP. 3, “Atmospheric Optics”)
- Marvin Klein** *Intelligent Optical Systems, Inc., Torrance, California* (VOL. IV, CHAP. 12, “Photorefractive Materials and Devices”)
- Thomas L. Koch** *Lehigh University, Bethlehem, Pennsylvania* (VOL. I, CHAP. 21, “Integrated Optics”)
- R. John Koshel** *Photon Engineering LLC, and College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. II, CHAP. 40, “Lighting and Applications”)
- Walter F. Kosonocky*** *New Jersey Institute of Technology, University Heights, Newark, New Jersey* (VOL. II, CHAP. 33, “Infrared Detector Arrays”)
- Lester J. Kozlowski** *Altasens, Inc., Westlake Village, California* (VOL. II, CHAP. 33, “Infrared Detector Arrays”)
- Melvin H. Kreitzer** *Opcon Associates, Inc., Cincinnati, Ohio* (VOL. I, CHAP. 27, “Camera Lenses”)
- Paul W. Kruse** *Consultant, Edina, Minnesota* (VOL. II, CHAP. 28, “Thermal Detectors”)
- Vasudevan Lakshminarayanan** *School of Optometry and Departments of Physics and Electrical Engineering, University of Waterloo, Waterloo, Ontario, Canada* (VOL. III, CHAPS. 8, 9, “Biological Waveguides”, “The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution”)
- Anthony LaRocca†** *General Dynamics, Advanced Information Systems, Ypsilanti, Michigan* (VOL. II, CHAP. 15, “Artificial Sources”)
- Melvin Lax*** *Department of Physics, City College of the City University of New York, New York, New York* (VOL. II, CHAP. 23, “Quantum Theory of the Laser”)
- Hwang Lee** *Hearne Institute for Theoretical Physics, Department of Physics and Astronomy, Louisiana State University, Baton Rouge, Louisiana* (VOL. IV, CHAP. 23, “Quantum Entanglement in Optical Interferometry”)
- Susanne M. Lee** *GE Global Research, Nikayuna, New York* (VOL. V, CHAP. 54, “X-Ray Tube Sources”)
- Bruno Lengeler** *Physikalisches Institut, RWTH Aachen University, Aachen, Germany* (VOL. V, CHAP. 37, “Refractive X-Ray Lenses”)
- Frederick J. Leonberger** *MIT Center for Integrated Photonic Systems, Cambridge, Massachusetts* (VOL. I, CHAP. 21, “Integrated Optics”)
- Guifang Li** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. V, CHAP. 15, “Fiber Optic Communication Links (Telecom, Datacom, and Analog)”)
- Eric Lifshin** *College of Nanoscale Science and Engineering, University at Albany, Albany, New York* (VOL. V, CHAP. 62, “X-Ray Spectral Detection and Imaging”)
- Chian Liu** *X-Ray Science Division, Argonne National Laboratory, Argonne, Illinois* (VOL. V, CHAP. 42, “Nanofocusing of Hard X-Rays with Multilayer Laue Lenses”)
- Dirk Lützenkirchen-Hecht** *Bergische Universität Wuppertal, Wuppertal, Germany* (VOL. V, CHAP. 30, “Requirements for X-Ray Spectroscopy”)
- John D. Lytle** *Advanced Optical Concepts, Santa Cruz, California* (VOL. IV, CHAP. 3, “Polymeric Optics”)
- Carolyn MacDonald** *University at Albany, Albany, New York* (VOL. V, CHAPS. 26, 53, 54, “An Introduction to X-Ray and Neutron Optics”, “Polycapillary X-Ray Optics”, “X-Ray Tube Sources”)
- Albert T. Macrander** *X-Ray Science Division, Argonne National Laboratory, Argonne, Illinois* (VOL. V, CHAP. 42, “Nanofocusing of Hard X-Rays with Multilayer Laue Lenses”)
- Virendra N. Mahajan** *The Aerospace Corporation, El Segundo, California* (VOL. II, CHAP. 11, “Orthonormal Polynomials in Wavefront Analysis”; VOL. V, CHAP. 4, “Imaging through Atmospheric Turbulence”)

*Deceased.

†Retired.

- Daniel Malacara-Hernández** *Centro de Investigaciones en Óptica, A. C., León, Gto., México* (VOL. II, CHAPS. 12, 13, “Optical Metrology”, “Optical Testing”)
- Zacarias Malacara** *Centro de Investigaciones en Óptica, A. C., León, Gto., México* (VOL. II, CHAP. 12, “Optical Metrology”)
- Theresa A. Maldonado** *Department of Electrical and Computer Engineering, Texas A&M University, College Station, Texas* (VOL. V, CHAP. 7, “Electro-Optic Modulators”)
- Pavel V. Mamyshev** *Bell Laboratories—Lucent Technologies, Holmdel, New Jersey* (VOL. V, CHAP. 22, “Solitons in Optical Fiber Communication Systems”)
- Michael P. Mandina** *Brandon Light, Optimax Systems, Inc., Ontario, New York* (VOL. II, CHAP. 9, “Optical Fabrication”)
- Masud Mansuripur** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. I, CHAP. 35, “Principles of Optical Disk Data Storage”)
- Jonathan P. Marangos** *Quantum Optics and Laser Science Group, Blackett Laboratory, Imperial College, London, United Kingdom* (VOL. IV, CHAP. 14, “Electromagnetically Induced Transparency”)
- Arvind S. Marathay** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. I, CHAP. 3, “Diffraction”)
- Jörg Maser** *X-Ray Science Division, Argonne National Laboratory, Argonne, Illinois, and Center for Nanoscale Materials, Argonne National Laboratory, Argonne, Illinois* (VOL. V, CHAP. 42, “Nanofocusing of Hard X-Rays with Multilayer Laue Lenses”)
- Barry R. Masters** *Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts* (VOL. III, CHAP. 17, “Three-Dimensional Confocal Microscopy of the Living Human Cornea”)
- Ignacio Matias** *Public University Navarra, Pamplona, Spain* (VOL. V, CHAP. 24, “Optical Fiber Sensors”)
- John F. McCalmont** *Air Force Research Laboratory, Sensors Directorate, Wright-Patterson AFB, Ohio* (VOL. I, CHAP. 3, “Diffraction”)
- Timothy S. McComb** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. V, CHAP. 25, “High-Power Fiber Lasers and Amplifiers”)
- Harold J. Metcalf** *Department of Physics, State University of New York, Stony Brook, New York* (VOL. IV, CHAP. 20, “Laser Cooling and Trapping of Atoms”)
- Alan Michette** *King’s College, London, United Kingdom* (VOL. V, CHAPS. 40, 56, “Zone Plates”, “Laser Generated Plasmas”)
- David Mildner** *NIST Center for Neutron Research, National Institute of Standards and Technology, Gaithersburg, Maryland* (VOL. V, CHAP. 63, “Neutron Optics”)
- Alan Miller** *Scottish Universities Physics Alliance, School of Engineering and Physical Sciences, Heriot-Watt University, Edinburgh, Scotland* (VOL. IV, CHAP. 8, “Fundamental Optical Properties of Solids”)
- Donald T. Miller** *School of Optometry, Indiana University, Bloomington, Indiana* (VOL. III, CHAP. 15, “Adaptive Optics in Retinal Microscopy and Vision”)
- Tom D. Milster** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. I, CHAP. 22, “Miniature and Micro-Optics”)
- Scott T. Mixture** *Kazuo Inamori School of Engineering, Alfred University, Alfred, New York* (VOL. V, CHAP. 28, “Requirements for X-Ray Diffraction”)
- Curtis D. Mobley** *Applied Electromagnetics and Optics Laboratory, SRI International, Menlo Park, California* (VOL. IV, CHAP. 1, “Optical Properties of Water”)
- Duncan T. Moore** *The Institute of Optics, and Gradient Lens Corporation, Rochester, New York* (VOL. I, CHAP. 24, “Gradient Index Optics”)
- Jacob Moskovich** *Opcon Associates, Inc., Cincinnati, Ohio* (VOL. I, CHAP. 27, “Camera Lenses”)
- Daniel Nolan** *Corning Inc., Corning, New York* (VOL. V, CHAP. 16, “Fiber-Based Couplers”)

- Paul R. Norton** *U.S. Army Night Vision and Electronics Directorate, Fort Belvoir, Virginia* (VOL. II, CHAP. 24, "Photodetectors")
- Donald C. O'Shea** *Georgia Institute of Technology, School of Physics, Atlanta, Georgia* (VOL. II, CHAP. 2, "Aberration Curves in Lens Design")
- Yoshi Ohno** *Optical Technology Division, National Institute of Standards and Technology, Gaithersburg, Maryland* (VOL. II, CHAP. 37, "Radiometry and Photometry for Vision Optics")
- Rudolf Oldenbourg** *Marine Biological Laboratory, Woods Hole, Massachusetts, and Physics Department, Brown University, Providence, Rhode Island* (VOL. I, CHAP. 28, "Microscopes")
- Gregory H. Olsen** *Sensors Unlimited, Inc., Princeton, New Jersey* (VOL. II, CHAP. 25, "Photodetection")
- Larry D. Owen** *NuOptics International, Phoenix, Arizona* (VOL. II, CHAP. 31, "Image Tube Intensified Electronic Imaging")
- Joseph C. Palais** *Ira A. Fulton School of Engineering, Arizona State University, Tempe, Arizona* (VOL. V, CHAP. 18, "Micro-Optics-Based Components for Networking")
- James M. Palmer*** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. II, CHAPS. 35, 36, "Measurement of Transmission, Absorption, Emission, and Reflection"; "Radiometry and Photometry: Units and Conversions")
- Zhongqi Pan** *University of Louisiana at Lafayette, Lafayette, Louisiana* (VOL. V, CHAP. 21, "WDM Fiber-Optic Communication Networks")
- Thrasylvoulos N. Pappas** *Department of Electrical and Computer Engineering, Northwestern University, Evanston, Illinois* (VOL. III, CHAP. 24, "Human Vision and Electronic Imaging")
- Roger A. Paquin** *Advanced Materials Consultant, Tucson, Arizona, and Optical Sciences Center, University of Arizona, Tucson* (VOL. IV, CHAP. 4, "Properties of Metals")
- Greg J. Pearce** *Max-Planck Institute for the Science of Light, Erlangen, Germany* (VOL. V, CHAP. 11, "Photonic Crystal Fibers")
- Denis G. Pelli** *Psychology Department and Center for Neural Science, New York University, New York* (VOL. III, CHAP. 3, "Psychophysical Methods")
- Douglas Pfeiffer** *Boulder Community Hospital, Boulder, Colorado* (VOL. V, CHAP. 31, "Requirements for Medical Imaging and X-Ray Inspection")
- Barbara K. Pierscionek** *Department of Biomedical Sciences, University of Ulster, Coleraine, United Kingdom* (VOL. III, CHAP. 19, "Gradient Index Optics in the Eye")
- Ladislav Pina** *Faculty of Nuclear Sciences and Physical Engineering, Czech Technical University, Prague, Holesovickach* (VOL. V, CHAP. 48, "Multifoil X-Ray Optics")
- Stephen M. Pompea** *National Optical Astronomy Observatory, Tucson, Arizona* (VOL. IV, CHAP. 6, "Characterization and Use of Black Surfaces for Optical Systems")
- Georgeanne M. Purvinis** *The Battelle Memorial Institute, Columbus, Ohio* (VOL. V, CHAP. 7, "Electro-Optic Modulators")
- Brian D. Ramsey** *National Aeronautics and Space Administration, Marshall Space Flight Center, Huntsville, Alabama* (VOL. V, CHAPS. 47, 64, "Astronomical X-Ray Optics"; "Grazing-Incidence Neutron Optics")
- John Reintjes** *Optical Sciences Division, Naval Research Laboratory, Washington, D.C.* (VOL. IV, CHAP. 15, "Stimulated Raman and Brillouin Scattering")
- Lloyd Huff Research** *Institute, University of Dayton, Dayton, Ohio* (VOL. I, CHAP. 33, "Holography and Holographic Instruments")
- William T. Rhodes** *School of Electrical and Computer Engineering, Georgia Institute of Technology, and Department of Electrical Engineering and Imaging Technology Center, Florida Atlantic University, Boca Raton, Florida* (VOL. I, CHAP. 6, "Coherence Theory: Tools and Applications")

*Deceased.

- Richard L. Rhorer** *National Institute of Standards and Technology, Gaithersburg, Maryland* (VOL. II, CHAP. 10, "Fabrication of Optics by Diamond Turning")
- Martin C. Richardson** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. V, CHAP. 25, "High-Power Fiber Lasers and Amplifiers")
- Michael Roberts** *Pilkington Optronics, Wales, United Kingdom* (VOL. II, CHAP. 8, "Thermal Compensation Techniques")
- Eric W. Rogala** *Raytheon Missile Systems, Tucson, Arizona* (VOL. IV, CHAP. 2, "Properties of Crystals and Glasses")
- Philip J. Rogers** *Pilkington Optronics, Wales, United Kingdom* (VOL. II, CHAP. 8, "Thermal Compensation Techniques")
- Bernice E. Rogowitz** *IBM T. J. Watson Research Center, Hawthorne, New York* (VOL. III, CHAP. 24, "Human Vision and Electronic Imaging")
- Scott O. Rohrbach** *Optics Branch, Goddard Space Flight Center, NASA, Greenbelt, Maryland* (VOL. V, CHAP. 33, "Requirements for X-Ray Astronomy")
- Austin Roorda** *School of Optometry, University of California, Berkeley, California* (VOL. III, CHAP. 15, "Adaptive Optics in Retinal Microscopy and Vision")
- Laurence S. Rothman** *Harvard-Smithsonian Center for Astrophysics, Atomic and Molecular Physics Division, Cambridge, Massachusetts* (VOL. V, CHAP. 3, "Atmospheric Optics")
- Philip St. J. Russell** *Max-Planck Institute for the Science of Light, Erlangen, Germany* (VOL. V, CHAP. 11, "Photonic Crystal Fibers")
- Timo T. Saha** *NASA/Goddard Space Flight Center, Greenbelt, Maryland* (VOL. V, CHAP. 45, "Aberrations for Grazing Incidence Optics")
- Manuel Sanchez del Rio** *European Synchrotron Radiation Facility, Grenoble, France* (VOL. V, CHAP. 35, "Ray Tracing of X-Ray Optical Systems")
- Brooke E. Schefrin** *Department of Psychology, University of Colorado, Boulder, Colorado* (VOL. III, CHAP. 14, "Optics and Vision of the Aging Eye")
- Winston V. Schoenfeld** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. II, CHAP. 18, "High-Brightness Visible LEDs")
- Clifton Schor** *School of Optometry, University of California, Berkeley, California* (VOL. III, CHAP. 13, "Binocular Vision Factors That Influence Optical Design")
- Christian G. Schroer** *Institute of Structural Physics, TU Dresden, Dresden, Germany* (VOL. V, CHAP. 37, "Refractive X-Ray Lenses")
- Jim Schwiegerling** *Department of Ophthalmology, University of Arizona, Tucson, Arizona* (VOL. III, CHAP. 21, "Intraocular Lenses")
- Marian O. Scully** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Max-Planck-Institut für Quantenoptik, Garching bei München, Germany* (VOL. II, CHAP. 23, "Quantum Theory of the Laser")
- David G. Seiler** *Semiconductor Electronics Division, National Institute of Standards and Technology, Gaithersburg, Maryland* (VOL. IV, CHAP. 5, "Optical Properties of Semiconductors")
- Robert R. Shannon[†]** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. II, CHAPS. 4, 5, "Optical Specifications", "Tolerancing Techniques")
- James E. Sheedy** *College of Optometry, Pacific University, Forest Grove, Oregon* (VOL. III, CHAP. 23, "Vision Problems at Computers")
- Mansoor Sheik-Bahae** *Department of Physics and Astronomy, University of New Mexico, Albuquerque, New Mexico* (VOL. IV, CHAP. 16, "Third-Order Optical Nonlinearities")

[†]Retired.

- Qun Shen** *National Synchrotron Light Source II, Brookhaven National Laboratory, Upton, New York* (VOL. V, CHAPS. 27, 43, “Coherent X-Ray Optics and Microscopy”, “Polarizing Crystal Optics”)
- Martin Shenker** *Martin Shenker Optical Design, Inc., White Plains, New York* (VOL. III, CHAP. 25, “Visual Factors Associated with Head-Mounted Displays”)
- Carolyn J. Sher DeCusatis** *Pace University, White Plains, New York* (VOL. II, CHAP. 38, “Spectroradiometry”)
- Michael Shribak** *Marine Biological Laboratory, Woods Hole, Massachusetts* (VOL. I, CHAP. 28, “Microscopes”)
- Peter Siddons** *National Synchrotron Light Source, Brookhaven National Laboratory, Upton, New York* (VOL. V, CHAPS. 39, 60, “Crystal Monochromators and Bent Crystals”, “Introduction to X-Ray Detectors”)
- Uwe Siegner** *Institute of Quantum Electronics, Physics Department, Swiss Federal Institute of Technology (ETH), Zurich, Switzerland* (VOL. IV, CHAP. 18, “Nonlinear Optical Processes for Ultrashort Pulse Generation”)
- William T. Silfvast** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. II, CHAP. 16, “Lasers”)
- Douglas C. Sinclair** *Sinclair Optics, Inc., Fairport, New York* (VOL. II, CHAP. 3, “Optical Design Software”)
- David H. Sliney** *Consulting Medical Physicist, Fallston, Maryland, and Retired, U.S. Army Center for Health Promotion and Preventive Medicine, Laser/Optical Radiation Program, Aberdeen Proving Ground, Maryland* (VOL. III, CHAP. 7, “Ocular Radiation Hazards”)
- Warren J. Smith*** *Kaiser Electro-Optics, Inc., Carlsbad, California* (VOL. II, CHAP. 1, “Techniques of First-Order Layout”)
- Marion J. Soileau** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. IV, CHAP. 19, “Laser-Induced Damage to Optical Materials”)
- Leo H. Spiekman** *Alphion Corp., Princeton Junction, New Jersey* (VOL. V, CHAP. 19, “Semiconductor Optical Amplifiers”)
- Eberhard Spiller** *Spiller X-Ray Optics, Livermore, California* (VOL. V, CHAP. 41, “Multilayers”)
- Howard Stark†** *Xerox Corporation, Corporate Research and Technology, Rochester, New York* (VOL. I, CHAP. 34, “Xerographic Systems”)
- Duncan G. Steel** *Physics Department, University of Michigan, Ann Arbor, Michigan* (VOL. IV, CHAP. 11, “Coherent Optical Transients”)
- G. Brian Stephenson** *Center for Nanoscale Materials, Argonne National Laboratory, Argonne, Illinois, Materials Science Division, Argonne National Laboratory, Argonne, Illinois* (VOL. V, CHAP. 42, “Nanofocusing of Hard X-Rays with Multilayer Laue Lenses”)
- Andrew Stockman** *Department of Visual Neuroscience, UCL Institute of Ophthalmology, London, United Kingdom* (VOL. III, CHAPS. 10, 11, “Colorimetry”, “Color Vision Mechanisms”)
- John C. Stover** *The Scatter Works, Inc., Tucson, Arizona* (VOL. V, CHAP. 1, “Scatterometers”)
- Paul G. Suchoski** *Audigence Inc., Melbourne, Florida* (VOL. I, CHAP. 21, “Integrated Optics”)
- Peter Z. Takacs** *Instrumentation Division, Brookhaven National Laboratory, Upton, New York* (VOL. I, CHAP. 8, “Surface Scattering”; VOL. V, CHAP. 46, “X-Ray Mirror Metrology”)
- Greg Tallents** *University of York, York, United Kingdom* (VOL. V, CHAP. 58, “X-Ray Lasers”)
- Chung L. Tang** *School of Electrical and Computer Engineering, Cornell University, Ithaca, New York* (VOL. IV, CHAP. 10, “Nonlinear Optics”)
- Matthew S. Taubman** *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (VOL. II, CHAP. 22, “Laser Stabilization”)
- Michael E. Thomas** *Applied Physics Laboratory, Johns Hopkins University, Laurel, Maryland* (VOL. IV, CHAP. 2, “Properties of Crystals and Glasses”)

*Deceased.

†Retired.

- Tomasz S. Tkaczyk** *Department of Bioengineering, Rice University, Houston, Texas* (VOL. I, CHAP. 22, "Miniature and Micro-Optics")
- Timothy J. Tredwell** *Sensor Systems Division, Imager Systems Development Laboratory, Eastman Kodak Company, Rochester, New York* (VOL. II, CHAP. 32, "Visible Array Detectors")
- William J. Tropf** *Applied Physics Laboratory, Johns Hopkins University, Laurel, Maryland* (VOL. IV, CHAP. 2, "Properties of Crystals and Glasses")
- Brian H. Tsou** *Air Force Research Laboratory, Wright Patterson AFB, Ohio* (VOL. III, CHAP. 25, "Visual Factors Associated with Head-Mounted Displays")
- Peter van der Straten** *Debye Institute, Department of Atomic and Interface Physics, Utrecht University, Utrecht, The Netherlands* (VOL. IV, CHAP. 20, "Laser Cooling and Trapping of Atoms")
- Wilfrid B. Veldkamp** *MIT/Lincoln Laboratory, Lexington, Massachusetts* (VOL. I, CHAP. 23, "Binary Optics")
- Pierre R. Villeneuve** *Department of Physics, Massachusetts Institute of Technology, Cambridge, Massachusetts* (VOL. IV, CHAP. 9, "Photonic Bandgap Materials")
- Herbert Walther*** *Max-Planck-Institut für Quantenoptik, Garching bei München, Germany, and Sektion Physik der Universität München, Garching bei München, Germany* (VOL. II, CHAP. 23, "Quantum Theory of the Laser")
- Robert H. Webb** *The Schepens Eye Research Institute, Boston, Massachusetts* (VOL. III, CHAP. 5, "Optical Generation of the Visual Stimulus")
- Robert H. Weissman** *Hewlett-Packard Co., San Jose, California* (VOL. II, CHAP. 17, "Light-Emitting Diodes")
- John S. Werner** *Department of Ophthalmology & Vision Science, University of California, Davis, Sacramento, California* (VOL. III, CHAP. 14, "Optics and Vision of the Aging Eye")
- Gerald Westheimer** *Division of Neurobiology, University of California, Berkeley, California* (VOL. III, CHAPS. 4, 6, "Visual Acuity and Hyperacuity", "The Maxwellian View with an Addendum on Apodization")
- William B. Wetherell†** *Optical Research Associates, Framingham, Massachusetts* (VOL. I, CHAP. 18, "Afocal Systems")
- Yih G. Wey** *Department of Electrical and Computer Engineering, University of California, Santa Barbara, California* (VOL. II, CHAP. 26, "High-Speed Photodetectors")
- Jay M. Wiesenfeld** *Bell Laboratories, Alcatel-Lucent, Murray Hill, New Jersey* (VOL. V, CHAP. 19, "Semiconductor Optical Amplifiers")
- Christoph F. Wildfeuer** *Hearne Institute for Theoretical Physics, Department of Physics and Astronomy, Louisiana State University, Baton Rouge, Louisiana* (VOL. IV, CHAP. 23, "Quantum Entanglement in Optical Interferometry")
- Gwyn P. Williams** *Free Electron Laser, Thomas Jefferson National Accelerator Facility, Newport News, Virginia* (VOL. V, CHAP. 55, "Synchrotron Sources")
- John R. Willison** *Stanford Research Systems, Inc., Sunnyvale, California* (VOL. II, CHAP. 27, "Signal Detection and Analysis")
- Alan E. Willner** *University of Southern California, Los Angeles, California* (VOL. V, CHAP. 21, "WDM Fiber-Optic Communication Networks")
- William L. Wolfe** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. I, CHAP. 19, "Nondispersive Prisms"; VOL. II, CHAP. 28, "Thermal Detectors")
- Shin-Tson Wu** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (VOL. V, CHAP. 8, "Liquid Crystals")
- James C. Wyant** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. II, CHAP. 14, "Use of Computer-Generated Holograms in Optical Testing")
- Yong Xie** *Texas Instruments Inc., Dallas, Texas* (VOL. V, CHAP. 21, "WDM Fiber-Optic Communication Networks")

*Deceased.

†Retired.

Hanfei Yan *Center for Nanoscale Materials, Argonne National Laboratory, Argonne, Illinois, and National Synchrotron Light Source II, Brookhaven National Laboratory, Upton, New York* (VOL. V, CHAP. 42, “Nanofocusing of Hard X-Rays with Multilayer Laue Lenses”)

Jun Ye *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (VOL. II, CHAP. 22, “Laser Stabilization”)

Paul R. Yoder, Jr. *Consultant in Optical Engineering, Norwalk, Connecticut* (VOL. II, CHAP. 6, “Mounting Optical Components”)

Changyuan Yu *National University of Singapore, and A *STAR Institute for Infocomm Research, Singapore* (VOL. V, CHAP. 21, “WDM Fiber-Optic Communication Networks”)

Edward F. Zalewski *College of Optical Sciences, University of Arizona, Tucson, Arizona* (VOL. II, CHAP. 34, “Radiometry and Photometry”)

George J. Zissis *Environmental Research Institute of Michigan, Ann Arbor, Michigan* (VOL. I, CHAP. 20, “Dispersive Prisms and Gratings”)

Stefan Zollner *Freescall Semiconductor, Inc., Hopewell Junction, New York* (VOL. IV, CHAP. 5, “Optical Properties of Semiconductors”)

M. Suhail Zubairy *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Department of Electronics, Quaid-i-Azam University, Islamabad, Pakistan* (VOL. II, CHAP. 23, “Quantum Theory of the Laser”)

DO NOT DUPLICATE

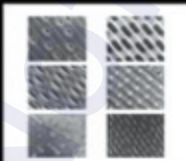
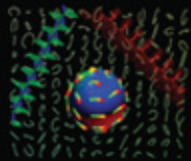
Third Edition

Sponsored by the Optical Society of America

HANDBOOK OF OPTICS

Volume I

*Geometrical and Physical Optics, Polarized Light,
Components and Instruments*



Editor-in-Chief:
Michael Bass

Associate Editors:
Casimer M. DeCusatis
Jay M. Enoch
Vasudevan Lakshminarayanan
Guifang Li
Carolyn MacDonald
Virendra N. Mahajan
Eric Van Stryland

OSA[®]

HANDBOOK OF OPTICS

DO NOT DUPLICATE

ABOUT THE EDITORS

Editor-in-Chief: Dr. Michael Bass is professor emeritus at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Associate Editors:

Dr. Casimer M. DeCusatis is a distinguished engineer and technical executive with IBM Corporation.

Dr. Jay M. Enoch is dean emeritus and professor at the School of Optometry at the University of California, Berkeley.

Dr. Vasudevan Lakshminarayanan is professor of Optometry, Physics, and Electrical Engineering at the University of Waterloo, Ontario, Canada.

Dr. Guifang Li is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Dr. Carolyn MacDonald is a professor at the University at Albany, and director of the Center for X-Ray Optics.

Dr. Virendra N. Mahajan is a distinguished scientist at The Aerospace Corporation.

Dr. Eric Van Stryland is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

HANDBOOK OF OPTICS

Volume I
Geometrical and Physical Optics,
Polarized Light,
Components and Instruments

THIRD EDITION

Sponsored by the
OPTICAL SOCIETY OF AMERICA

Michael Bass Editor-in-Chief
*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

Virendra N. Mahajan Associate Editor
*The Aerospace Corporation
El Segundo, California*



New York Chicago San Francisco Lisbon London Madrid
Mexico City Milan New Delhi San Juan Seoul
Singapore Sydney Toronto

This page intentionally left blank.

DO NOT DUPLICATE

COVER ILLUSTRATIONS

Left: Poincaré sphere describing light's polarization states is shown floating in front of a depolarized field of polarization ellipses, with linearly and circularly polarized fields propagating on its left and right, respectively. See Chaps. 12 and 15.

Middle: Triplet lens developed for photographic applications that can zero out the primary aberrations by splitting the positive lens of a doublet into two and placing one on each side of the negative lens. See Chap. 17.

Right: Micrographs of different optical storage media showing the straight and narrow tracks with 1.6- μm spacing between adjacent tracks. The recorded information bits appear as short marks along each track. See Chap. 35.

This page intentionally left blank.

DO NOT DUPLICATE

CONTENTS

Contributors	xvii
Brief Contents of All Volumes	xix
Editors' Preface	xxv
Preface to Volume I	xxvii
Glossary and Fundamental Constants	xxix

Part 1. Geometrical Optics

Chapter 1. General Principles of Geometrical Optics	1.3
<i>Douglas S. Goodman</i>	
<hr/>	
1.1	Glossary / 1.3
1.2	Introduction / 1.7
1.3	Fundamentals / 1.8
1.4	Characteristic Functions / 1.13
1.5	Rays in Heterogeneous Media / 1.18
1.6	Conservation of Étendue / 1.22
1.7	Skew Invariant / 1.23
1.8	Refraction and Reflection at Interfaces between Homogeneous Media / 1.23
1.9	Imaging / 1.26
1.10	Description of Systems of Revolution / 1.32
1.11	Tracing Rays in Centered Systems of Spherical Surfaces / 1.35
1.12	Paraxial Optics of Systems of Revolution / 1.37
1.13	Images About Known Rays / 1.43
1.14	Gaussian Lens Properties / 1.44
1.15	Collineation / 1.56
1.16	System Combinations: Gaussian Properties / 1.63
1.17	Paraxial Matrix Methods / 1.65
1.18	Apertures, Pupils, Stops, Fields, and Related Matters / 1.74
1.19	Geometrical Aberrations of Point Images: Description / 1.85
1.20	References / 1.92

Part 2. Physical Optics

Chapter 2. Interference	2.3
<i>John E. Greivenkamp</i>	
<hr/>	
2.1	Glossary / 2.3
2.2	Introduction / 2.3
2.3	Waves and Wavefronts / 2.3
2.4	Interference / 2.5
2.5	Interference by Wavefront Division / 2.14
2.6	Interference by Amplitude Division / 2.19
2.7	Multiple Beam Interference / 2.28
2.8	Coherence and Interference / 2.36
2.9	Applications of Interference / 2.42
2.10	References / 2.42

Chapter 3. Diffraction *Arvind S. Marathay and John F. McCalmont* 3.1

- 3.1 Glossary / 3.1
- 3.2 Introduction / 3.1
- 3.3 Light Waves / 3.2
- 3.4 Huygens-Fresnel Construction / 3.4
- 3.5 Cylindrical Wavefront / 3.13
- 3.6 Mathematical Theory of Diffraction / 3.21
- 3.7 Stationary Phase Approximation / 3.29
- 3.8 Vector Diffraction / 3.32
- 3.9 Acknowledgments / 3.38
- 3.10 References / 3.38

Chapter 4. Transfer Function Techniques *Glenn D. Boreman* 4.1

- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.1
- 4.3 Definitions / 4.2
- 4.4 MTF Calculations / 4.3
- 4.5 MTF Measurements / 4.6
- 4.6 References / 4.8

Chapter 5. Coherence Theory *William H. Carter* 5.1

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.1
- 5.3 Some Elementary Classical Concepts / 5.2
- 5.4 Definitions of Coherence Functions / 5.4
- 5.5 Model Sources / 5.9
- 5.6 Propagation / 5.13
- 5.7 Spectrum of Light / 5.19
- 5.8 Polarization Effects / 5.22
- 5.9 Applications / 5.22
- 5.10 References / 5.23
- 5.11 Additional Reading / 5.26

Chapter 6. Coherence Theory: Tools and Applications
*Gisele Bennett, William T. Rhodes,
and J. Christopher James* 6.1

- 6.1 Glossary / 6.1
- 6.2 Introduction / 6.2
- 6.3 Key Definitions and Relationships / 6.2
- 6.4 Propagation, Diffraction, and Scattering: Enhanced Backscatter and the Lau Effect / 6.5
- 6.5 Image Formation: Lukosz-Type Super-Resolving System / 6.9
- 6.6 Efficient Sampling of Coherence Functions / 6.10
- 6.7 An Example of When Not to Use Coherence Theory / 6.12
- 6.8 Concluding Remarks / 6.13
- 6.9 References / 6.13

Chapter 7. Scattering by Particles *Craig F. Bohren* 7.1

- 7.1 Glossary / 7.1
- 7.2 Introduction / 7.2
- 7.3 Scattering: An Overview / 7.3
- 7.4 Scattering by Particles: Basic Concepts and Terminology / 7.4
- 7.5 Scattering by an Isotropic, Homogeneous Sphere: The Archetype / 7.11
- 7.6 Scattering by Regular Particles / 7.14

- 7.7 Computational Methods for Nonspherical Particles / 7.15
 7.8 References / 7.17

Chapter 8. Surface Scattering *Eugene L. Church
and Peter Z. Takacs* **8.1**

- 8.1 Glossary of Principal Symbols / 8.1
 8.2 Introduction / 8.2
 8.3 Notation / 8.2
 8.4 The Fresnel-Kirchhoff Approximation / 8.5
 8.5 The Rayleigh-Rice (RR) or Small-Perturbation Approximation / 8.9
 8.6 Effects of Finite Illumination Area / 8.12
 8.7 Surface Statistics / 8.12
 8.8 Surface Finish Specification / 8.16
 8.9 Retrospect and Prospect / 8.17
 8.10 References and Endnotes / 8.18

Chapter 9. Volume Scattering in Random Media *Aristide Dogariu and Jeremy Ellis* **9.1**

- 9.1 Glossary / 9.1
 9.2 Introduction / 9.2
 9.3 General Theory of Scattering / 9.3
 9.4 Single Scattering / 9.4
 9.5 Multiple Scattering / 9.8
 9.6 References / 9.18

**Chapter 10. Optical Spectroscopy and Spectroscopic
Lineshapes** *Brian Henderson* **10.1**

- 10.1 Glossary / 10.1
 10.2 Introductory Comments / 10.2
 10.3 Theoretical Preliminaries / 10.3
 10.4 Rates of Spectroscopic Transition / 10.4
 10.5 Lineshapes of Spectral Transitions / 10.6
 10.6 Spectroscopy of One-Electron Atoms / 10.7
 10.7 Multielectron Atoms / 10.10
 10.8 Optical Spectra and the Outer Electronic Structure / 10.12
 10.9 Spectra of Tri-Positive Rare Earth Ions / 10.16
 10.10 Vibrational and Rotational Spectra of Molecules / 10.18
 10.11 Lineshapes in Solid State Spectroscopy / 10.22
 10.12 References / 10.27

Chapter 11. Analog Optical Signal and Image Processing *Joseph W. Goodman* **11.1**

- 11.1 Glossary / 11.1
 11.2 Introduction / 11.1
 11.3 Fundamental Analog Operations / 11.2
 11.4 Analog Optical Fourier Transforms / 11.3
 11.5 Spatial Filtering / 11.5
 11.6 Coherent Optical Processing of Synthetic Aperture Radar Data / 11.6
 11.7 Coherent Optical Processing of Temporal Signals / 11.8
 11.8 Optical Processing of Two-Dimensional Images / 11.12
 11.9 Incoherent Processing of Discrete Signals / 11.17
 11.10 Concluding Remarks / 11.20
 11.11 References / 11.20

Part 3. Polarized Light

Chapter 12. Polarization *Jean M. Bennett* **12.3**

- 12.1 Glossary / 12.3
- 12.2 Basic Concepts and Conventions / 12.4
- 12.3 Fresnel Equations / 12.6
- 12.4 Basic Relations for Polarizers / 12.14
- 12.5 Polarization by Nonnormal-Incidence Reflection (Pile of Plates) / 12.15
- 12.6 Polarization by Nonnormal-Incidence Transmission (Pile of Plates) / 12.18
- 12.7 Quarter-Wave Plates and Other Phase Retardation Plates / 12.24
- 12.8 Matrix Methods for Computing Polarization / 12.27
- 12.9 References / 12.30

Chapter 13. Polarizers *Jean M. Bennett* **13.1**

- 13.1 Glossary / 13.1
- 13.2 Prism Polarizers / 13.2
- 13.3 Glan-Type Prisms / 13.8
- 13.4 Nicol-Type Prisms / 13.15
- 13.5 Polarizing Beam-Splitter Prisms / 13.18
- 13.6 Feussner Prisms / 13.22
- 13.7 Noncalcite Polarizing Prisms / 13.23
- 13.8 Dichroic and Diffraction-Type Polarizers / 13.24
- 13.9 Non-Normal-Incidence Reflection and Transmission Polarizers / 13.33
- 13.10 Retardation Plates / 13.43
- 13.11 Variable Retardation Plates and Compensators / 13.53
- 13.12 Half-Shade Devices / 13.56
- 13.13 Miniature Polarization Devices / 13.57
- 13.14 References / 13.58

Chapter 14. Mueller Matrices *Russell A. Chipman* **14.1**

- 14.1 Glossary / 14.1
- 14.2 Conventions / 14.3
- 14.3 Objectives / 14.3
- 14.4 Stokes Parameters and Mueller Matrices / 14.4
- 14.5 The Stokes Parameters and the Poincaré Sphere / 14.4
- 14.6 Mueller Matrices / 14.6
- 14.7 Sequences of Polarization Elements / 14.7
- 14.8 Polarization Elements' Properties in the Mueller Calculus / 14.7
- 14.9 Rotation of an Element About the Optical Axis / 14.8
- 14.10 Nonpolarizing Mueller Matrices / 14.8
- 14.11 Mueller Matrices of Ideal Polarizers / 14.8
- 14.12 Retarder Mueller Matrices / 14.11
- 14.13 Retarder Mueller Matrices Ambiguities and Retarder Space / 14.14
- 14.14 Transmittance and Diattenuation / 14.16
- 14.15 Polarizance / 14.18
- 14.16 Mueller Matrices of Diattenuators / 14.18
- 14.17 Normalizing a Mueller Matrix / 14.19
- 14.18 Coordinate System for the Mueller Matrix / 14.19
- 14.19 Mueller Matrices for Refraction / 14.20
- 14.20 Mueller Matrices for Reflection / 14.21
- 14.21 Conversion between Mueller Matrices and Jones Matrices / 14.22
- 14.22 Nondepolarizing Mueller Matrices and Mueller-Jones Matrices / 14.24
- 14.23 Homogeneous and Inhomogeneous Polarization Elements / 14.25
- 14.24 Mueller Matrices Near the Identity Matrix, Weak Polarization Elements / 14.26
- 14.25 Matrix Roots of Nondepolarizing Mueller Matrices / 14.27
- 14.26 Depolarization and the Depolarization Index / 14.30

14.27	Degree of Polarization Surfaces and Maps / 14.31
14.28	The Depolarization Index / 14.32
14.29	The Average Degree of Polarization / 14.32
14.30	Determining Mueller Matrix Properties / 14.33
14.31	Generators for Depolarization / 14.33
14.32	Interpretation of Arbitrary Mueller Matrices, the Polar Decomposition of Mueller Matrices / 14.39
14.33	Physically Realizable Mueller Matrices / 14.40
14.34	Acknowledgments / 14.42
14.35	References / 14.43

Chapter 15. Polarimetry *Russell A. Chipman*

15.1

15.1	Glossary / 15.1
15.2	Objectives / 15.3
15.3	Polarimeters / 15.3
15.4	Light-Measuring Polarimeters / 15.3
15.5	Sample-Measuring Polarimeters / 15.4
15.6	Complete and Incomplete Polarimeters / 15.4
15.7	Polarization Generators and Analyzers / 15.4
15.8	Classes of Polarimeters / 15.5
15.9	Time-Sequential Measurements / 15.5
15.10	Polarization Modulation / 15.5
15.11	Division of Aperture / 15.5
15.12	Division of Amplitude / 15.5
15.13	Spectropolarimeters / 15.6
15.14	Imaging Polarimeters / 15.6
15.15	Definitions / 15.6
15.16	Stokes Vectors and Mueller Matrices / 15.8
15.17	Phenomenological Definition of the Stokes Vector / 15.9
15.18	Polarization Properties of Light Beams / 15.9
15.19	Mueller Matrices / 15.11
15.20	Data Reduction for Light-Measuring Polarimeters / 15.11
15.21	Sample-Measuring Polarimeters for Measuring Mueller Matrix Elements / 15.13
15.22	Polarimetric Measurement Equation and Polarimetric Data-Reduction Equation / 15.14
15.23	Dual Rotating Retarder Polarimeter / 15.16
15.24	Incomplete Sample-Measuring Polarimeters / 15.16
15.25	Nonideal Polarization Elements / 15.17
15.26	Elliptical and Circular Polarizers and Analyzers / 15.17
15.27	Common Defects of Polarization Elements / 15.19
15.28	Polarization Modulators, Retardance Modulators / 15.20
15.29	Rotating Retarders / 15.20
15.30	Photo-Elastic Modulators / 15.21
15.31	Liquid Crystal Retarders / 15.21
15.32	Electro-Optical Modulators / 15.23
15.33	Magneto-Optical Modulators / 15.23
15.34	Fiber Squeezers / 15.24
15.35	Polarimeter Design Metrics / 15.24
15.36	Singular Value Decomposition Examples / 15.26
15.37	Polarimeter Error Analysis / 15.27
15.38	The Mueller Matrix for Polarization Component Characterization / 15.28
15.39	Retro-Reflection Testing and Correction for Supplemental Optics / 15.28
15.40	Applications of Polarimetry / 15.29
15.41	Ellipsometry and Generalized Ellipsometry / 15.30
15.42	Liquid Crystal Cell and System Testing / 15.32
15.43	Polarization Aberrations / 15.35
15.44	Remote Sensing / 15.37
15.45	Polarization Light Scattering / 15.38
15.46	Ophthalmic Polarimetry / 15.39
15.47	Acknowledgments / 15.41
15.48	References / 15.41

Chapter 16. Ellipsometry *Rasheed M. A. Azzam* **16.1**

- 16.1 Glossary / 16.1
- 16.2 Introduction / 16.2
- 16.3 Conventions / 16.3
- 16.4 Modeling and Inversion / 16.4
- 16.5 Transmission Ellipsometry / 16.10
- 16.6 Instrumentation / 16.10
- 16.7 Jones-Matrix Generalized Ellipsometry / 16.19
- 16.8 Mueller-Matrix Generalized Ellipsometry / 16.19
- 16.9 Applications / 16.21
- 16.10 References / 16.21

Part 4. Components

Chapter 17. Lenses *R. Barry Johnson* **17.3**

- 17.1 Glossary / 17.3
- 17.2 Introduction / 17.4
- 17.3 Basics / 17.5
- 17.4 Stops and Pupils / 17.8
- 17.5 F-Number and Numerical Aperture / 17.9
- 17.6 Magnifier or Eye Loupe / 17.9
- 17.7 Compound Microscopes / 17.10
- 17.8 Field and Relay Lenses / 17.10
- 17.9 Aplanatic Surfaces and Immersion Lenses / 17.10
- 17.10 Single Element Lens / 17.12
- 17.11 Landscape Lenses and the Influence of Stop Position / 17.17
- 17.12 Two-Lens Systems / 17.20
- 17.13 Achromatic Doublets / 17.22
- 17.14 Triplet Lenses / 17.26
- 17.15 Symmetrical Lenses / 17.26
- 17.16 Double-Gauss Lenses / 17.27
- 17.17 Petzval Lenses / 17.28
- 17.18 Telephoto Lenses / 17.29
- 17.19 Inverted or Reverse Telephoto Lenses / 17.29
- 17.20 Performance of Representative Lenses / 17.29
- 17.21 Rapid Estimation of Lens Performance / 17.36
- 17.22 Bibliography / 17.40

Chapter 18. Afocal Systems *William B. Wetherell* **18.1**

- 18.1 Glossary / 18.1
- 18.2 Introduction / 18.2
- 18.3 Gaussian Analysis of Afocal Lenses / 18.2
- 18.4 Keplerian Afocal Lenses / 18.7
- 18.5 Galilean and Inverse Galilean Afocal Lenses / 18.15
- 18.6 Relay Trains and Periscopes / 18.17
- 18.7 Reflecting and Catadioptric Afocal Lenses / 18.19
- 18.8 References / 18.23

Chapter 19. Nondispersive Prisms *William L. Wolfe* **19.1**

- 19.1 Glossary / 19.1
- 19.2 Introduction / 19.1
- 19.3 Inversion, Reversion / 19.2
- 19.4 Deviation, Displacement / 19.2
- 19.5 Summary of Prism Properties / 19.2

- 19.6 Prism Descriptions / 19.2
19.7 References / 19.29

Chapter 20. Dispersive Prisms and Gratings *George J. Zissis* **20.1**

- 20.1 Glossary / 20.1
20.2 Introduction / 20.1
20.3 Prisms / 20.2
20.4 Gratings / 20.3
20.5 Prism and Grating Configurations and Instruments / 20.4
20.6 References / 20.15

Chapter 21. Integrated Optics *Thomas L. Koch,
Frederick J. Leonberger, and Paul G. Suchoski* **21.1**

- 21.1 Glossary / 21.1
21.2 Introduction / 21.2
21.3 Device Physics / 21.3
21.4 Integrated Optics Materials and Fabrication Technology / 21.13
21.5 Circuit Elements / 21.21
21.6 Applications of Integrated Optics / 21.31
21.7 Future Trends / 21.39
21.8 References / 21.41

Chapter 22. Miniature and Micro-Optics *Tom D. Milster and Tomasz S. Tkaczyk* **22.1**

- 22.1 Glossary / 22.1
22.2 Introduction / 22.2
22.3 Uses of Micro-Optics / 22.2
22.4 Micro-Optics Design Considerations / 22.2
22.5 Molded Microlenses / 22.8
22.6 Diamond Turning / 22.15
22.7 Lithography for Making Refractive Components / 22.18
22.8 Monolithic Lenslet Modules / 22.25
22.9 Distributed-Index Planar Microlenses / 22.26
22.10 Micro-Fresnel Lenses / 22.31
22.11 Liquid Lenses / 22.37
22.12 Other Technologies / 22.42
22.13 References / 22.47

Chapter 23. Binary Optics *Michael W. Farn
and Wilfrid B. Veldkamp* **23.1**

- 23.1 Glossary / 23.1
23.2 Introduction / 23.2
23.3 Design—Geometrical Optics / 23.2
23.4 Design—Scalar Diffraction Theory / 23.10
23.5 Design—Vector Diffraction Theory / 23.13
23.6 Fabrication / 23.14
23.7 References / 23.17

Chapter 24. Gradient Index Optics *Duncan T. Moore* **24.1**

- 24.1 Glossary / 24.1
24.2 Introduction / 24.1
24.3 Analytic Solutions / 24.2
24.4 Mathematical Representation / 24.2
24.5 Axial Gradient Lenses / 24.3

- 24.6 Radial Gradients / 24.5
- 24.7 Radial Gradients with Curved Surfaces / 24.7
- 24.8 Shallow Radial Gradients / 24.7
- 24.9 Materials / 24.8
- 24.10 References / 24.9

Part 5. Instruments

Chapter 25. Cameras *Norman Goldberg* 25.3

- 25.1 Glossary / 25.3
- 25.2 Introduction / 25.3
- 25.3 Background / 25.4
- 25.4 Properties of the Final Image / 25.5
- 25.5 Film Choice / 25.5
- 25.6 Resolving Fine Detail / 25.5
- 25.7 Film Sizes / 25.6
- 25.8 Display / 25.6
- 25.9 Distributing the Image / 25.7
- 25.10 Video Cameras / 25.7
- 25.11 Instant Pictures / 25.8
- 25.12 Critical Features / 25.8
- 25.13 Time Lag / 25.8
- 25.14 Automation / 25.10
- 25.15 Flash / 25.16
- 25.16 Flexibility through Features and Accessories / 25.16
- 25.17 Advantages of Various Formats / 25.17
- 25.18 Large Format: A Different World / 25.18
- 25.19 Special Cameras / 25.20
- 25.20 Further Reading / 25.26

Chapter 26. Solid-State Cameras *Gerald C. Holst* 26.1

- 26.1 Glossary / 26.1
- 26.2 Introduction / 26.2
- 26.3 Imaging System Applications / 26.3
- 26.4 Charge-Coupled Device Array Architecture / 26.3
- 26.5 Charge Injection Device / 26.6
- 26.6 Complementary Metal-Oxide Semiconductor / 26.8
- 26.7 Array Performance / 26.9
- 26.8 Camera Performance / 26.12
- 26.9 Modulation Transfer Function / 26.14
- 26.10 Resolution / 26.15
- 26.11 Sampling / 26.16
- 26.12 Storage, Analysis, and Display / 26.19
- 26.13 References / 26.20

Chapter 27. Camera Lenses *Ellis Betensky, Melvin H. Kreitzer, and Jacob Moskovich* 27.1

- 27.1 Introduction / 27.1
- 27.2 Imposed Design Limitations / 27.1
- 27.3 Modern Lens Types / 27.2
- 27.4 Classification System / 27.17
- 27.5 Lens Performance Data / 27.24
- 27.6 Acknowledgments / 27.25
- 27.7 Further Reading / 27.25

Chapter 28. Microscopes <i>Rudolf Oldenbourg and Michael Shribak</i>	28.1
<hr/>	
28.1 Glossary / 28.1	
28.2 Introduction / 28.1	
28.3 Optical Arrangements, Lenses, and Resolution / 28.3	
28.4 Contrast and Imaging Modes / 28.24	
28.5 Manipulation of Specimen / 28.54	
28.6 Acknowledgment / 28.55	
28.7 References / 28.56	
Chapter 29. Reflective and Catadioptric Objectives <i>Lloyd Jones</i>	29.1
<hr/>	
29.1 Glossary / 29.1	
29.2 Introduction / 29.2	
29.3 Glass Varieties / 29.2	
29.4 Introduction to Catadioptric and Reflective Objectives / 29.2	
29.5 Field-of-View Plots / 29.34	
29.6 Definitions / 29.36	
29.7 References / 29.38	
Chapter 30. Scanners <i>Leo Beiser and R. Barry Johnson</i>	30.1
<hr/>	
30.1 Glossary / 30.1	
30.2 Introduction / 30.2	
30.3 Scanned Resolution / 30.6	
30.4 Scanners for Remote Sensing / 30.14	
30.5 Scanning for Input/Output Imaging / 30.25	
30.6 Scanner Devices and Techniques / 30.34	
30.7 Scan-Error Reduction / 30.48	
30.8 Agile Beam Steering / 30.51	
30.9 References / 30.64	
30.10 Further Reading / 30.68	
Chapter 31. Optical Spectrometers <i>Brian Henderson</i>	31.1
<hr/>	
31.1 Glossary / 31.1	
31.2 Introduction / 31.2	
31.3 Optical Absorption Spectrometers / 31.2	
31.4 Luminescence Spectrometers / 31.5	
31.5 Photoluminescence Decay Time / 31.12	
31.6 Polarization Spectrometers / 31.15	
31.7 High-Resolution Techniques / 31.23	
31.8 Light Scattering / 31.30	
31.9 References / 31.31	
Chapter 32. Interferometers <i>Parameswaran Hariharan</i>	32.1
<hr/>	
32.1 Glossary / 32.1	
32.2 Introduction / 32.2	
32.3 Basic Types of Interferometers / 32.2	
32.4 Three-Beam and Double-Passed Two-Beam Interferometers / 32.7	
32.5 Fringe-Counting Interferometers / 32.8	
32.6 Two-Wavelength Interferometry / 32.9	
32.7 Frequency-Modulation Interferometers / 32.9	
32.8 Heterodyne Interferometers / 32.10	
32.9 Phase-Shifting Interferometers / 32.10	
32.10 Phase-Locked Interferometers / 32.11	

- 32.11 Laser-Doppler Interferometers / 32.12
- 32.12 Laser-Feedback Interferometers / 32.13
- 32.13 Fiber Interferometers / 32.14
- 32.14 Interferometric Wave Meters / 32.16
- 32.15 Second-Harmonic and Phase-Conjugate Interferometers / 32.17
- 32.16 Stellar Interferometers / 32.19
- 32.17 Gravitational-Wave Interferometers / 32.21
- 32.18 References / 32.22

Chapter 33. Holography and Holographic Instruments 33.1
Lloyd Huff

- 33.1 Glossary / 33.1
- 33.2 Introduction / 33.2
- 33.3 Background and Basic Principles / 33.2
- 33.4 Holographic Interferometry / 33.4
- 33.5 Holographic Optical Elements / 33.13
- 33.6 Holographic Inspection / 33.16
- 33.7 Holographic Lithography / 33.22
- 33.8 Holographic Memory / 33.24
- 33.9 Conclusion / 33.25
- 33.10 References / 33.25

Chapter 34. Xerographic Systems 34.1
Howard Stark

- 34.1 Introduction and Overview / 34.1
- 34.2 Creation of the Latent Image / 34.1
- 34.3 Development / 34.5
- 34.4 Transfer / 34.10
- 34.5 Fusing / 34.10
- 34.6 Cleaning and Erasing / 34.10
- 34.7 Control Systems / 34.11
- 34.8 Color / 34.11
- 34.9 References / 34.13

Chapter 35. Principles of Optical Disk Data Storage 35.1
Masud Mansuripur

- 35.1 Introduction / 35.1
- 35.2 Preliminaries and Basic Definitions / 35.2
- 35.3 The Optical Path / 35.7
- 35.4 Automatic Focusing / 35.12
- 35.5 Automatic Tracking / 35.14
- 35.6 Thermomagnetic Recording Process / 35.17
- 35.7 Magneto-Optical Readout / 35.21
- 35.8 Materials of Magneto-Optical Recording / 35.25
- 35.9 Concluding Remarks / 35.28
- 35.10 Further Information / 35.30
- 35.11 Bibliography / 35.31

CONTRIBUTORS

- Rasheed M. A. Azzam** *Department of Electrical Engineering, University of New Orleans, New Orleans, Louisiana (CHAP. 16)*
- Leo Beiser** *Consultant, Flushing, New York (CHAP. 30)*
- Gisele Bennett** *Electro-Optical Systems Laboratory and School of Electrical and Computer Engineering, Georgia Tech Research Institute, Georgia Institute of Technology, Atlanta, Georgia (CHAP. 6)*
- Jean M. Bennett*** *Research Department, Michelson Laboratory, Naval Air Warfare Center, China Lake, California (CHAPS. 12, 13)*
- Ellis Betensky** *Opcon Associates, Inc., West Redding, Connecticut (CHAP. 27)*
- Craig F. Bohren** *Pennsylvania State University, University Park, Pennsylvania (CHAP. 7)*
- Glenn D. Boreman** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida (CHAP. 4)*
- William H. Carter*** *Naval Research Laboratory, Washington, D.C. (CHAP. 5)*
- Russell A. Chipman** *College of Optical Sciences, University of Arizona, Tucson, Arizona (CHAPS. 14, 15)*
- Eugene L. Church** *Instrumentation Division, Brookhaven National Laboratory, Upton, New York (CHAP. 8)*
- Aristide Dogariu** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida (CHAP. 9)*
- Jeremy Ellis** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida (CHAP. 9)*
- Michael W. Farn** *MIT/Lincoln Laboratory, Lexington, Massachusetts (CHAP. 23)*
- Norman Goldberg** *Madison, Wisconsin (CHAP. 25)*
- Douglas S. Goodman** *Corning Tropel Corporation, Fairport, New York (CHAP. 1)*
- Joseph W. Goodman** *Department of Electrical Engineering, Stanford University, Stanford, California (CHAP. 11)*
- John E. Greivenkamp** *College of Optical Sciences, University of Arizona, Tucson, Arizona (CHAP. 2)*
- Parameswaran Hariharan** *School of Physics, University of Sydney, Sydney, Australia (CHAP. 32)*
- Brian Henderson** *Department of Physics and Applied Physics, University of Strathclyde, Glasgow, United Kingdom (CHAPS. 10, 31)*
- Gerald C. Holst** *JCD Publishing, Winter Park, Florida (CHAP. 26)*
- Lloyd Huff** *Research Institute, University of Dayton, Dayton, Ohio (CHAP. 33)*
- J. Christopher James** *Electro-Optical Systems Laboratory, Georgia Tech Research Institute, Georgia Institute of Technology, Atlanta, Georgia (CHAP. 6)*
- R. Barry Johnson** *Consultant, Huntsville, Alabama (CHAPS. 17, 30)*
- Lloyd Jones** *College of Optical Sciences, University of Arizona, Tucson, Arizona (CHAP. 29)*
- Thomas L. Koch** *Lehigh University, Bethlehem, Pennsylvania (CHAP. 21)*
- Melvin H. Kreitzer** *Opcon Associates, Inc., Cincinnati, Ohio (CHAP. 27)*

*Deceased.

- Frederick J. Leonberger** *MIT Center for Integrated Photonic Systems, Cambridge, Massachusetts* (CHAP. 21)
- Masud Mansuripur** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 35)
- Arvind S. Marathay** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 3)
- John F. McCalmont** *Air Force Research Laboratory, Sensors Directorate, Wright-Patterson AFB, Ohio* (CHAP. 3)
- Tom D. Milster** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 22)
- Duncan T. Moore** *The Institute of Optics, and Gradient Lens Corporation, Rochester, New York* (CHAP. 24)
- Jacob Moskovich** *Opcon Associates, Inc., Cincinnati, Ohio* (CHAP. 27)
- Rudolf Oldenbourg** *Marine Biological Laboratory, Woods Hole, Massachusetts, and Physics Department, Brown University, Providence, Rhode Island* (CHAP. 28)
- William T. Rhodes** *School of Electrical and Computer Engineering, Georgia Institute of Technology, and Department of Electrical Engineering and Imaging Technology Center, Florida Atlantic University, Boca Raton, Florida* (CHAP. 6)
- Michael Shribak** *Marine Biological Laboratory, Woods Hole, Massachusetts* (CHAP. 28)
- Howard Stark**[†] *Xerox Corporation, Corporate Research and Technology, Rochester, New York* (CHAP. 34)
- Paul G. Suchoski** *Audigence Inc., Melbourne, Florida* (CHAP. 21)
- Peter Z. Takacs** *Instrumentation Division, Brookhaven National Laboratory, Upton, New York* (CHAP. 8)
- Tomasz S. Tkaczyk** *Department of Bioengineering, Rice University, Houston, Texas* (CHAP. 22)
- Wilfrid B. Veldkamp** *MIT/Lincoln Laboratory, Lexington, Massachusetts* (CHAP. 23)
- William B. Wetherell**[†] *Optical Research Associates, Framingham, Massachusetts* (CHAP. 18)
- William L. Wolfe** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 19)
- George J. Zissis** *Environmental Research Institute of Michigan, Ann Arbor, Michigan* (CHAP. 20)

[†]Retired.

BRIEF CONTENTS OF ALL VOLUMES

VOLUME I. GEOMETRICAL AND PHYSICAL OPTICS, POLARIZED LIGHT, COMPONENTS AND INSTRUMENTS

PART 1. GEOMETRICAL OPTICS

Chapter 1. General Principles of Geometrical Optics *Douglas S. Goodman*

PART 2. PHYSICAL OPTICS

Chapter 2. Interference *John E. Greivenkamp*

Chapter 3. Diffraction *Arvind S. Marathay and John F. McCalmont*

Chapter 4. Transfer Function Techniques *Glenn D. Boreman*

Chapter 5. Coherence Theory *William H. Carter*

Chapter 6. Coherence Theory: Tools and Applications *Gisele Bennett, William T. Rhodes,
and J. Christopher James*

Chapter 7. Scattering by Particles *Craig F. Bohren*

Chapter 8. Surface Scattering *Eugene L. Church and Peter Z. Takacs*

Chapter 9. Volume Scattering in Random Media *Aristide Dogariu and Jeremy Ellis*

Chapter 10. Optical Spectroscopy and Spectroscopic Lineshapes *Brian Henderson*

Chapter 11. Analog Optical Signal and Image Processing *Joseph W. Goodman*

PART 3. POLARIZED LIGHT

Chapter 12. Polarization *Jean M. Bennett*

Chapter 13. Polarizers *Jean M. Bennett*

Chapter 14. Mueller Matrices *Russell A. Chipman*

Chapter 15. Polarimetry *Russell A. Chipman*

Chapter 16. Ellipsometry *Rasheed M. A. Azzam*

PART 4. COMPONENTS

Chapter 17. Lenses *R. Barry Johnson*

Chapter 18. Afocal Systems *William B. Wetherell*

Chapter 19. Nondispersive Prisms *William L. Wolfe*

Chapter 20. Dispersive Prisms and Gratings *George J. Zissis*

Chapter 21. Integrated Optics *Thomas L. Koch, Frederick J. Leonberger, and Paul G. Suchoski*

Chapter 22. Miniature and Micro-Optics *Tom D. Milster and Tomasz S. Tkaczyk*

Chapter 23. Binary Optics *Michael W. Farn and Wilfrid B. Veldkamp*

Chapter 24. Gradient Index Optics *Duncan T. Moore*

PART 5. INSTRUMENTS

Chapter 25. Cameras *Norman Goldberg*

Chapter 26. Solid-State Cameras *Gerald C. Holst*

Chapter 27. Camera Lenses *Ellis Betensky, Melvin H. Kreitzer, and Jacob Moskovich*

Chapter 28. Microscopes *Rudolf Oldenbourg and Michael Shribak*

- Chapter 29. Reflective and Catadioptric Objectives *Lloyd Jones*
Chapter 30. Scanners *Leo Beiser and R. Barry Johnson*
Chapter 31. Optical Spectrometers *Brian Henderson*
Chapter 32. Interferometers *Parameswaran Hariharan*
Chapter 33. Holography and Holographic Instruments *Lloyd Huff*
Chapter 34. Xerographic Systems *Howard Stark*
Chapter 35. Principles of Optical Disk Data Storage *Masud Mansuripur*

VOLUME II. DESIGN, FABRICATION, AND TESTING; SOURCES AND DETECTORS; RADIOMETRY AND PHOTOMETRY

PART 1. DESIGN

- Chapter 1. Techniques of First-Order Layout *Warren J. Smith*
Chapter 2. Aberration Curves in Lens Design *Donald C. O'Shea and Michael E. Harrigan*
Chapter 3. Optical Design Software *Douglas C. Sinclair*
Chapter 4. Optical Specifications *Robert R. Shannon*
Chapter 5. Tolerancing Techniques *Robert R. Shannon*
Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.*
Chapter 7. Control of Stray Light *Robert P. Breault*
Chapter 8. Thermal Compensation Techniques *Philip J. Rogers and Michael Roberts*

PART 2. FABRICATION

- Chapter 9. Optical Fabrication *Michael P. Mandina*
Chapter 10. Fabrication of Optics by Diamond Turning *Richard L. Rhorer and Chris J. Evans*

PART 3. TESTING

- Chapter 11. Orthonormal Polynomials in Wavefront Analysis *Virendra N. Mahajan*
Chapter 12. Optical Metrology *Zacarias Malacara and Daniel Malacara-Hernández*
Chapter 13. Optical Testing *Daniel Malacara-Hernández*
Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant*

PART 4. SOURCES

- Chapter 15. Artificial Sources *Anthony LaRocca*
Chapter 16. Lasers *William T. Silfvast*
Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman*
Chapter 18. High-Brightness Visible LEDs *Winston V. Schoenfeld*
Chapter 19. Semiconductor Lasers *Pamela L. Derry, Luis Figueroa, and Chi-shain Hong*
Chapter 20. Ultrashort Optical Sources and Applications *Jean-Claude Diels and Ladan Arissian*
Chapter 21. Attosecond Optics *Zenghu Chang*
Chapter 22. Laser Stabilization *John L. Hall, Matthew S. Taubman, and Jun Ye*
Chapter 23. Quantum Theory of the Laser *János A. Bergou, Berthold-Georg Englert, Melvin Lax, Marian O. Scully, Herbert Walther, and M. Suhail Zubairy*

PART 5. DETECTORS

- Chapter 24. Photodetectors *Paul R. Norton*
Chapter 25. Photodetection *Abhay M. Joshi and Gregory H. Olsen*
Chapter 26. High-Speed Photodetectors *John E. Bowers and Yih G. Wey*
Chapter 27. Signal Detection and Analysis *John R. Willison*
Chapter 28. Thermal Detectors *William L. Wolfe and Paul W. Kruse*

PART 6. IMAGING DETECTORS

- Chapter 29. Photographic Films *Joseph H. Altman*
Chapter 30. Photographic Materials *John D. Baloga*

- Chapter 31. Image Tube Intensified Electronic Imaging *C. Bruce Johnson and Larry D. Owen*
 Chapter 32. Visible Array Detectors *Timothy J. Tredwell*
 Chapter 33. Infrared Detector Arrays *Lester J. Kozlowski and Walter F. Kosonocky*

PART 7. RADIOMETRY AND PHOTOMETRY

- Chapter 34. Radiometry and Photometry *Edward F. Zalewski*
 Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection *James M. Palmer*
 Chapter 36. Radiometry and Photometry: Units and Conversions *James M. Palmer*
 Chapter 37. Radiometry and Photometry for Vision Optics *Yoshi Ohno*
 Chapter 38. Spectroradiometry *Carolyn J. Sher DeCusatis*
 Chapter 39. Nonimaging Optics: Concentration and Illumination *William Cassarly*
 Chapter 40. Lighting and Applications *Anurag Gupta and R. John Koshel*

VOLUME III. VISION AND VISION OPTICS

- Chapter 1. Optics of the Eye *Neil Charman*
 Chapter 2. Visual Performance *Wilson S. Geisler and Martin S. Banks*
 Chapter 3. Psychophysical Methods *Denis G. Pelli and Bart Farell*
 Chapter 4. Visual Acuity and Hyperacuity *Gerald Westheimer*
 Chapter 5. Optical Generation of the Visual Stimulus *Stephen A. Burns and Robert H. Webb*
 Chapter 6. The Maxwellian View with an Addendum on Apodization *Gerald Westheimer*
 Chapter 7. Ocular Radiation Hazards *David H. Sliney*
 Chapter 8. Biological Waveguides *Vasudevan Lakshminarayanan and Jay M. Enoch*
 Chapter 9. The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution *Jay M. Enoch and Vasudevan Lakshminarayanan*
 Chapter 10. Colorimetry *David H. Brainard and Andrew Stockman*
 Chapter 11. Color Vision Mechanisms *Andrew Stockman and David H. Brainard*
 Chapter 12. Assessment of Refraction and Refractive Errors and Their Influence on Optical Design
B. Ralph Chou
 Chapter 13. Binocular Vision Factors That Influence Optical Design *Clifton Schor*
 Chapter 14. Optics and Vision of the Aging Eye *John S. Werner, Brooke E. Scheffrin, and Arthur Bradley*
 Chapter 15. Adaptive Optics in Retinal Microscopy and Vision *Donald T. Müller and Austin Roorda*
 Chapter 16. Refractive Surgery, Correction of Vision, PRK, and LASIK *L. Diaz-Santana and Harilaos Ginis*
 Chapter 17. Three-Dimensional Confocal Microscopy of the Living Human Cornea *Barry R. Masters*
 Chapter 18. Diagnostic Use of Optical Coherence Tomography in the Eye *Johannes F. de Boer*
 Chapter 19. Gradient Index Optics in the Eye *Barbara K. Pierscionek*
 Chapter 20. Optics of Contact Lenses *Edward S. Bennett*
 Chapter 21. Intraocular Lenses *Jim Schwiegerling*
 Chapter 22. Displays for Vision Research *William Cowan*
 Chapter 23. Vision Problems at Computers *Jeffrey Anshel and James E. Sheedy*
 Chapter 24. Human Vision and Electronic Imaging *Bernice E. Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allebach*
 Chapter 25. Visual Factors Associated with Head-Mounted Displays *Brian H. Tsou and Martin Shenker*

VOLUME IV. OPTICAL PROPERTIES OF MATERIALS, NONLINEAR OPTICS, QUANTUM OPTICS

PART 1. PROPERTIES

- Chapter 1. Optical Properties of Water *Curtis D. Mobley*
 Chapter 2. Properties of Crystals and Glasses *William J. Tropf, Michael E. Thomas, and Eric W. Rogala*
 Chapter 3. Polymeric Optics *John D. Lytle*
 Chapter 4. Properties of Metals *Roger A. Paquin*

- Chapter 5. Optical Properties of Semiconductors *David G. Seiler, Stefan Zollner, Alain C. Diebold, and Paul M. Amirtharaj*
- Chapter 6. Characterization and Use of Black Surfaces for Optical Systems *Stephen M. Pompea and Robert P. Breault*
- Chapter 7. Optical Properties of Films and Coatings *Jerzy A. Dobrowolski*
- Chapter 8. Fundamental Optical Properties of Solids *Alan Miller*
- Chapter 9. Photonic Bandgap Materials *Pierre R. Villeneuve*

PART 2. NONLINEAR OPTICS

- Chapter 10. Nonlinear Optics *Chung L. Tang*
- Chapter 11. Coherent Optical Transients *Paul R. Berman and D. G. Steel*
- Chapter 12. Photorefractive Materials and Devices *Mark Cronin-Golomb and Marvin Klein*
- Chapter 13. Optical Limiting *David J. Hagan*
- Chapter 14. Electromagnetically Induced Transparency *Jonathan P. Marangos and Thomas Halfmann*
- Chapter 15. Stimulated Raman and Brillouin Scattering *John Reintjes and M. Bashkansky*
- Chapter 16. Third-Order Optical Nonlinearities *Mansoor Sheik-Bahae and Michael P. Hasselbeck*
- Chapter 17. Continuous-Wave Optical Parametric Oscillators *M. Ebrahim-Zadeh*
- Chapter 18. Nonlinear Optical Processes for Ultrashort Pulse Generation *Uwe Siegner and Ursula Keller*
- Chapter 19. Laser-Induced Damage to Optical Materials *Marion J. Soileau*

PART 3. QUANTUM AND MOLECULAR OPTICS

- Chapter 20. Laser Cooling and Trapping of Atoms *Harold J. Metcalf and Peter van der Straten*
- Chapter 21. Strong Field Physics *Todd Ditmire*
- Chapter 22. Slow Light Propagation in Atomic and Photonic Media *Jacob B. Khurgin*
- Chapter 23. Quantum Entanglement in Optical Interferometry *Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver, and Jonathan P. Dowling*

VOLUME V. ATMOSPHERIC OPTICS, MODULATORS, FIBER OPTICS, X-RAY AND NEUTRON OPTICS

PART 1. MEASUREMENTS

- Chapter 1. Scatterometers *John C. Stover*
- Chapter 2. Spectroscopic Measurements *Brian Henderson*

PART 2. ATMOSPHERIC OPTICS

- Chapter 3. Atmospheric Optics *Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman*
- Chapter 4. Imaging through Atmospheric Turbulence *Virendra N. Mahajan and Guang-ming Dai*
- Chapter 5. Adaptive Optics *Robert Q. Fugate*

PART 3. MODULATORS

- Chapter 6. Acousto-Optic Devices *I-Cheng Chang*
- Chapter 7. Electro-Optic Modulators *Georgianne M. Purvinis and Theresa A. Maldonado*
- Chapter 8. Liquid Crystals *Sebastian Gauza and Shin-Tson Wu*

PART 4. FIBER OPTICS

- Chapter 9. Optical Fiber Communication Technology and System Overview *Ira Jacobs*
- Chapter 10. Nonlinear Effects in Optical Fibers *John A. Buck*
- Chapter 11. Photonic Crystal Fibers *Philip St. J. Russell and G. J. Pearce*
- Chapter 12. Infrared Fibers *James A. Harrington*
- Chapter 13. Sources, Modulators, and Detectors for Fiber Optic Communication Systems *Elsa Garmire*
- Chapter 14. Optical Fiber Amplifiers *John A. Buck*

- Chapter 15. Fiber Optic Communication Links (Telecom, Datacom, and Analog) *Casimer DeCusatis and Guifang Li*
- Chapter 16. Fiber-Based Couplers *Daniel Nolan*
- Chapter 17. Fiber Bragg Gratings *Kenneth O. Hill*
- Chapter 18. Micro-Optics-Based Components for Networking *Joseph C. Palais*
- Chapter 19. Semiconductor Optical Amplifiers *Jay M. Wiesenfeld and Leo H. Spiekman*
- Chapter 20. Optical Time-Division Multiplexed Communication Networks *Peter J. Delfyett*
- Chapter 21. WDM Fiber-Optic Communication Networks *Alan E. Willner, Changyuan Yu, Zhongqi Pan, and Yong Xie*
- Chapter 22. Solitons in Optical Fiber Communication Systems *Pavel V. Mamyshev*
- Chapter 23. Fiber-Optic Communication Standards *Casimer DeCusatis*
- Chapter 24. Optical Fiber Sensors *Richard O. Claus, Ignacio Matias, and Francisco Arregui*
- Chapter 25. High-Power Fiber Lasers and Amplifiers *Timothy S. McComb, Martin C. Richardson, and Michael Bass*

PART 5. X-RAY AND NEUTRON OPTICS

Subpart 5.1. Introduction and Applications

- Chapter 26. An Introduction to X-Ray and Neutron Optics *Carolyn A. MacDonald*
- Chapter 27. Coherent X-Ray Optics and Microscopy *Qun Shen*
- Chapter 28. Requirements for X-Ray diffraction *Scott T. Mixture*
- Chapter 29. Requirements for X-Ray Fluorescence *George J. Havrilla*
- Chapter 30. Requirements for X-Ray Spectroscopy *Dirk Lützenkirchen-Hecht and Ronald Frahm*
- Chapter 31. Requirements for Medical Imaging and X-Ray Inspection *Douglas Pfeiffer*
- Chapter 32. Requirements for Nuclear Medicine *Lars R. Furenlid*
- Chapter 33. Requirements for X-Ray Astronomy *Scott O. Rohrbach*
- Chapter 34. Extreme Ultraviolet Lithography *Franco Cerrina and Fan Jiang*
- Chapter 35. Ray Tracing of X-Ray Optical Systems *Franco Cerrina and M. Sanchez del Rio*
- Chapter 36. X-Ray Properties of Materials *Eric M. Gullikson*

Subpart 5.2. Refractive and Interference Optics

- Chapter 37. Refractive X-Ray Lenses *Bruno Lengeler and Christian G. Schroer*
- Chapter 38. Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region *Malcolm R. Howells*
- Chapter 39. Crystal Monochromators and Bent Crystals *Peter Siddons*
- Chapter 40. Zone Plates *Alan Michette*
- Chapter 41. Multilayers *Eberhard Spiller*
- Chapter 42. Nanofocusing of Hard X-Rays with Multilayer Laue Lenses *Albert T. Macrander, Hanfei Yan, Hyon Chol Kang, Jörg Maser, Chian Liu, Ray Conley, and G. Brian Stephenson*
- Chapter 43. Polarizing Crystal Optics *Qun Shen*

Subpart 5.3. Reflective Optics

- Chapter 44. Reflective Optics *James Harvey*
- Chapter 45. Aberrations for Grazing Incidence Optics *Timo T. Saha*
- Chapter 46. X-Ray Mirror Metrology *Peter Z. Takacs*
- Chapter 47. Astronomical X-Ray Optics *Marshall K. Joy and Brian D. Ramsey*
- Chapter 48. Multifoil X-Ray Optics *Ladislav Pina*
- Chapter 49. Pore Optics *Marco Beijersbergen*
- Chapter 50. Adaptive X-Ray Optics *Ali Khounsary*
- Chapter 51. The Schwarzschild Objective *Franco Cerrina*
- Chapter 52. Single Capillaries *Donald H. Bilderback and Sterling W. Cornaby*
- Chapter 53. Polycapillary X-Ray Optics *Carolyn MacDonald and Walter M. Gibson*

Subpart 5.4. X-Ray Sources

- Chapter 54. X-Ray Tube Sources *Susanne M. Lee and Carolyn MacDonald*
Chapter 55. Synchrotron Sources *Steven L. Hulbert and Gwyn P. Williams*
Chapter 56. Laser Generated Plasmas *Alan Michette*
Chapter 57. Pinch Plasma Sources *Victor Kantsyrev*
Chapter 58. X-Ray Lasers *Greg Tallents*
Chapter 59. Inverse Compton X-Ray Sources *Frank Carroll*

Subpart 5.5. X-Ray Detectors

- Chapter 60. Introduction to X-Ray Detectors *Walter M. Gibson and Peter Siddons*
Chapter 61. Advances in Imaging Detectors *Aaron Couture*
Chapter 62. X-Ray Spectral Detection and Imaging *Eric Lifshin*

Subpart 5.6. Neutron Optics and Applications

- Chapter 63. Neutron Optics *David Mildner*
Chapter 64. Grazing-Incidence Neutron Optics *Mikhail Gubarev and Brian Ramsey*

DO NOT DUPLICATE

EDITORS' PREFACE

The third edition of the *Handbook of Optics* is designed to pull together the dramatic developments in both the basic and applied aspects of the field while retaining the archival, reference book value of a handbook. This means that it is much more extensive than either the first edition, published in 1978, or the second edition, with Volumes I and II appearing in 1995 and Volumes III and IV in 2001. To cover the greatly expanded field of optics, the *Handbook* now appears in five volumes. Over 100 authors or author teams have contributed to this work.

Volume I is devoted to the fundamentals, components, and instruments that make optics possible. Volume II contains chapters on design, fabrication, testing, sources of light, detection, and a new section devoted to radiometry and photometry. Volume III concerns vision optics only and is printed entirely in color. In Volume IV there are chapters on the optical properties of materials, non-linear, quantum and molecular optics. Volume V has extensive sections on fiber optics and x ray and neutron optics, along with shorter sections on measurements, modulators, and atmospheric optical properties and turbulence. Several pages of color inserts are provided where appropriate to aid the reader. A purchaser of the print version of any volume of the *Handbook* will be able to download a digital version containing all of the material in that volume in PDF format to one computer (see download instructions on bound-in card). The combined index for all five volumes can be downloaded from www.HandbookofOpticsOnline.com.

It is possible by careful selection of what and how to present that the third edition of the *Handbook* could serve as a text for a comprehensive course in optics. In addition, students who take such a course would have the *Handbook* as a career-long reference.

Topics were selected by the editors so that the *Handbook* could be a desktop (bookshelf) general reference for the parts of optics that had matured enough to warrant archival presentation. New chapters were included on topics that had reached this stage since the second edition, and existing chapters from the second edition were updated where necessary to provide this compendium. In selecting subjects to include, we also had to select which subjects to leave out. The criteria we applied were: (1) was it a specific application of optics rather than a core science or technology and (2) was it a subject in which the role of optics was peripheral to the central issue addressed. Thus, such topics as medical optics, laser surgery, and laser materials processing were not included. While applications of optics are mentioned in the chapters there is no space in the *Handbook* to include separate chapters devoted to all of the myriad uses of optics in today's world. If we had, the third edition would be much longer than it is and much of it would soon be outdated. We designed the third edition of the *Handbook of Optics* so that it concentrates on the principles of optics that make applications possible.

Authors were asked to try to achieve the dual purpose of preparing a chapter that was a worthwhile reference for someone working in the field and that could be used as a starting point to become acquainted with that aspect of optics. They did that and we thank them for the outstanding results seen throughout the *Handbook*. We also thank Mr. Taisuke Soda of McGraw-Hill for his help in putting this complex project together and Mr. Alan Tourtlotte and Ms. Susannah Lehman of the Optical Society of America for logistical help that made this effort possible.

We dedicate the third edition of the *Handbook of Optics* to all of the OSA volunteers who, since OSA's founding in 1916, give their time and energy to promoting the generation, application, archiving, and worldwide dissemination of knowledge in optics and photonics.

Michael Bass, Editor-in-Chief

Associate Editors:

Casimer M. DeCusatis

Jay M. Enoch

Vasudevan Lakshminarayanan

Guifang Li

Carolyn MacDonald

Virendra N. Mahajan

Eric Van Stryland

This page intentionally left blank.

DO NOT DUPLICATE

PREFACE TO VOLUME I

The third edition of the *Handbook of Optics* has been completely reorganized, expanded, and updated. The four volumes of the second edition grew to five in the current edition. Each volume is divided into parts, where each part, sometimes referred to as a section, consists of several chapters related to a certain topic. Volumes I and II are devoted primarily to the basic concepts of optics and optical phenomena, sometimes called classical optics. Volume I starts with geometrical optics and continues with physical optics. This includes interference, diffraction, coherence theory, and scattering. A new chapter on tools and applications of coherence theory has been added. A several-chapter section follows devoted to issues of polarized light. The chapter on polarimetry has been updated and its content on the Mueller matrices now appears in a separate chapter by that title. Next there are chapters on components such as lenses, afocal systems, nondispersive and dispersive prisms, and special optics that include integrated, miniature and micro-, binary, and gradient index optics. Finally, there are several chapters on instruments. They include cameras and camera lenses, microscopes, reflective and catadioptric objectives, scanners, spectrometers, interferometers, xerographic systems, and optical disc data storage.

There are many other chapters in this edition of the *Handbook* that could have been included in Volumes I and II. However, page limitations prevented that. For example, in Volume V there is a section on Atmospheric Optics. It consists of three chapters, one on transmission through the atmosphere, another on imaging through atmospheric turbulence, and a third on adaptive optics to overcome some of the deleterious effects of turbulence.

The chapters are generally aimed at the graduate students, though practicing scientists and engineers will find them equally suitable as references on the topics discussed. Each chapter has sufficient references for additional and/or further study.

The whole *Handbook* has been retyped and the figures redrawn. The reader will find that the figures in the new edition are crisp. Ms. Arushi Chawla and her team from Glyph International have done an outstanding job in accomplishing this monumental task. Many of the authors updated and proofread their chapters. However, some authors have passed away since the second edition and others couldn't be located. Every effort has been made to ensure that such chapters have been correctly reproduced.

Virendra N. Mahajan
The Aerospace Corporation
Associate Editor

This page intentionally left blank.

DO NOT DUPLICATE

GLOSSARY AND FUNDAMENTAL CONSTANTS

Introduction

This glossary of the terms used in the *Handbook* represents to a large extent the language of optics. The symbols are representations of numbers, variables, and concepts. Although the basic list was compiled by the author of this section, all the editors have contributed and agreed to this set of symbols and definitions. Every attempt has been made to use the same symbols for the same concepts throughout the entire *Handbook*, although there are exceptions. Some symbols seem to be used for many concepts. The symbol α is a prime example, as it is used for absorptivity, absorption coefficient, coefficient of linear thermal expansion, and more. Although we have tried to limit this kind of redundancy, we have also bowed deeply to custom.

Units

The abbreviations for the most common units are given first. They are consistent with most of the established lists of symbols, such as given by the International Standards Organization ISO¹ and the International Union of Pure and Applied Physics, IUPAP.²

Prefixes

Similarly, a list of the numerical prefixes¹ that are most frequently used is given, along with both the common names (where they exist) and the multiples of ten that they represent.

Fundamental Constants

The values of the fundamental constants³ are listed following the sections on SI units.

Symbols

The most commonly used symbols are then given. Most chapters of the *Handbook* also have a glossary of the terms and symbols specific to them for the convenience of the reader. In the following list, the symbol is given, its meaning is next, and the most customary unit of measure for the quantity is presented in brackets. A bracket with a dash in it indicates that the quantity is unitless. Note that there is a difference between units and dimensions. An angle has units of degrees or radians and a solid angle square degrees or steradians, but both are pure ratios and are dimensionless. The unit symbols as recommended in the SI system are used, but decimal multiples of some of the dimensions are sometimes given. The symbols chosen, with some cited exceptions, are also those of the first two references.

RATIONALE FOR SOME DISPUTED SYMBOLS

The choice of symbols is a personal decision, but commonality improves communication. This section explains why the editors have chosen the preferred symbols for the *Handbook*. We hope that this will encourage more agreement.

Fundamental Constants

It is encouraging that there is almost universal agreement for the symbols for the fundamental constants. We have taken one small exception by adding a subscript B to the k for Boltzmann's constant.

Mathematics

We have chosen i as the imaginary almost arbitrarily. IUPAP lists both i and j , while ISO does not report on these.

Spectral Variables

These include expressions for the wavelength λ , frequency ν , wave number σ , ω for circular or radian frequency, k for circular or radian wave number and dimensionless frequency x . Although some use f for frequency, it can be easily confused with electronic or spatial frequency. Some use $\tilde{\nu}$ for wave number, but, because of typography problems and agreement with ISO and IUPAP, we have chosen σ ; it should not be confused with the Stefan-Boltzmann constant. For spatial frequencies we have chosen ξ and η , although f_x and f_y are sometimes used. ISO and IUPAP do not report on these.

Radiometry

Radiometric terms are contentious. The most recent set of recommendations by ISO and IUPAP are L for radiance [$\text{Wcm}^{-2}\text{sr}^{-1}$], M for radiant emittance or exitance [Wcm^{-2}], E for irradiance or incidence [Wcm^{-2}], and I for intensity [Wsr^{-2}]. The previous terms, W , H , N , and J , respectively, are still in many texts, notably Smith⁴ and Lloyd⁵ but we have used the revised set, although there are still shortcomings. We have tried to deal with the vexatious term *intensity* by using *specific intensity* when the units are $\text{Wcm}^{-2}\text{sr}^{-1}$, *field intensity* when they are Wcm^{-2} , and *radiometric intensity* when they are Wsr^{-1} .

There are two sets to terms for these radiometric quantities, which arise in part from the terms for different types of reflection, transmission, absorption, and emission. It has been proposed that the *ion* ending indicate a process, that the *ance* ending indicate a value associated with a particular sample, and that the *ivity* ending indicate a generic value for a "pure" substance. Then one also has reflectance, transmittance, absorptance, and emittance as well as reflectivity, transmissivity, absorptivity, and emissivity. There are now two different uses of the word emissivity. Thus the words *exitance*, *incidence*, and *sterance* were coined to be used in place of emittance, irradiance, and radiance. It is interesting that ISO uses radiance, exitance, and irradiance whereas IUPAP uses radiance, exitance [*sic*], and irradiance. We have chosen to use them both, i.e., emittance, irradiance, and radiance will be followed in square brackets by exitance, incidence, and sterance (or vice versa). Individual authors will use the different endings for transmission, reflection, absorption, and emission as they see fit.

We are still troubled by the use of the symbol E for irradiance, as it is so close in meaning to electric field, but we have maintained that accepted use. The spectral concentrations of these quantities, indicated by a wavelength, wave number, or frequency subscript (e.g., L_λ) represent partial differentiations; a subscript q represents a photon quantity; and a subscript ν indicates a quantity normalized to the response of the eye. Thereby, L_ν is luminance, E_ν illuminance, and M_ν and I_ν luminous emittance and luminous intensity. The symbols we have chosen are consistent with ISO and IUPAP.

The refractive index may be considered a radiometric quantity. It is generally complex and is indicated by $\tilde{n} = n - ik$. The real part is the relative refractive index and k is the extinction coefficient. These are consistent with ISO and IUPAP, but they do not address the complex index or extinction coefficient.

Optical Design

For the most part ISO and IUPAP do not address the symbols that are important in this area.

There were at least 20 different ways to indicate focal ratio; we have chosen FN as symmetrical with NA; we chose f and efl to indicate the effective focal length. Object and image distance, although given many different symbols, were finally called s_o and s_i since s is an almost universal symbol for distance. Field angles are θ and ϕ ; angles that measure the slope of a ray to the optical axis are u ; u can also be $\sin u$. Wave aberrations are indicated by W_{ijk} , while third-order ray aberrations are indicated by σ_i and more mnemonic symbols.

Electromagnetic Fields

There is no argument about \mathbf{E} and \mathbf{H} for the electric and magnetic field strengths, Q for quantity of charge, ρ for volume charge density, σ for surface charge density, etc. There is no guidance from Refs. 1 and 2 on polarization indication. We chose \perp and \parallel rather than p and s , partly because s is sometimes also used to indicate scattered light.

There are several sets of symbols used for reflection transmission, and (sometimes) absorption, each with good logic. The versions of these quantities dealing with field amplitudes are usually specified with lower case symbols: r , t , and a . The versions dealing with power are alternately given by the uppercase symbols or the corresponding Greek symbols: R and T versus ρ and τ . We have chosen to use the Greek, mainly because these quantities are also closely associated with Kirchhoff's law that is usually stated symbolically as $\alpha = \epsilon$. The law of conservation of energy for light on a surface is also usually written as $\alpha + \rho + \tau = 1$.

Base SI Quantities

length	m	meter
time	s	second
mass	kg	kilogram
electric current	A	ampere
temperature	K	kelvin
amount of substance	mol	mole
luminous intensity	cd	candela

Derived SI Quantities

energy	J	joule
electric charge	C	coulomb
electric potential	V	volt
electric capacitance	F	farad
electric resistance	Ω	ohm
electric conductance	S	siemens
magnetic flux	Wb	weber
inductance	H	henry
pressure	Pa	pascal
magnetic flux density	T	tesla
frequency	Hz	hertz
power	W	watt
force	N	newton
angle	rad	radian
angle	sr	steradian

Prefixes

Symbol	Name	Common name	Exponent of ten
F	exa		18
P	peta		15
T	tera	trillion	12
G	giga	billion	9
M	mega	million	6
k	kilo	thousand	3
h	hecto	hundred	2
da	deca	ten	1
d	deci	tenth	-1
c	centi	hundredth	-2
m	milli	thousandth	-3
μ	micro	millionth	-6
n	nano	billionth	-9
p	pico	trillionth	-12
f	femto		-15
a	atto		-18

Constants

c	speed of light vacuo [299792458 ms ⁻¹]
c_1	first radiation constant = $2\pi^2 h = 3.7417749 \times 10^{-16}$ [Wm ²]
c_2	second radiation constant = $hc/k = 0.014838769$ [mK]
e	elementary charge [$1.60217733 \times 10^{-19}$ C]
g_n	free fall constant [9.80665 ms ⁻²]
h	Planck's constant [$6.6260755 \times 10^{-34}$ Ws]
k_B	Boltzmann constant [1.380658×10^{-23} JK ⁻¹]
m_e	mass of the electron [$9.1093897 \times 10^{-31}$ kg]
N_A	Avogadro constant [6.0221367×10^{23} mol ⁻¹]
R_∞	Rydberg constant [10973731.534 m ⁻¹]
ϵ_0	vacuum permittivity [$\mu_0^{-1}c^{-2}$]
σ	Stefan-Boltzmann constant [5.67051×10^{-8} Wm ⁻¹ K ⁻⁴]
μ_0	vacuum permeability [$4\pi \times 10^{-7}$ NA ⁻²]
μ_B	Bohr magneton [$9.2740154 \times 10^{-24}$ JT ⁻¹]

General

B	magnetic induction [Wbm ⁻² , kgs ⁻¹ C ⁻¹]
C	capacitance [f, C ² s ² m ⁻² kg ⁻¹]
C	curvature [m ⁻¹]
c	speed of light in vacuo [ms ⁻¹]
c_1	first radiation constant [Wm ²]
c_2	second radiation constant [mK]
D	electric displacement [Cm ⁻²]
E	incidence [irradiance] [Wm ⁻²]
e	electronic charge [coulomb]
E_v	illuminance [lux, lmm ⁻²]
E	electrical field strength [Vm ⁻¹]
E	transition energy [J]
E_g	band-gap energy [eV]
f^g	focal length [m]
f_f	Fermi occupation function, conduction band
f_v	Fermi occupation function, valence band

FN	focal ratio (<i>f</i> /number) [—]
<i>g</i>	gain per unit length [m^{-1}]
g_{th}	gain threshold per unit length [m^{-1}]
H	magnetic field strength [Am^{-1} , $\text{Cs}^{-1} \text{m}^{-1}$]
<i>h</i>	height [m]
<i>I</i>	irradiance (see also <i>E</i>) [Wm^{-2}]
<i>I</i>	radiant intensity [Wsr^{-1}]
<i>I</i>	nuclear spin quantum number [—]
<i>I</i>	current [A]
<i>i</i>	$\sqrt{-1}$
Im()	imaginary part of
<i>J</i>	current density [Am^{-2}]
j	total angular momentum [$\text{kg m}^2 \text{s}^{-1}$]
$J_1()$	Bessel function of the first kind [—]
<i>k</i>	radian wave number $=2\pi/\lambda$ [rad cm^{-1}]
k	wave vector [rad cm^{-1}]
<i>k</i>	extinction coefficient [—]
<i>L</i>	sterance [radiance] [$\text{Wm}^{-2} \text{sr}^{-1}$]
L_v	luminance [cdm^{-2}]
<i>L</i>	inductance [h, $\text{m}^2 \text{kg C}^2$]
<i>L</i>	laser cavity length
<i>L, M, N</i>	direction cosines [—]
<i>M</i>	angular magnification [—]
<i>M</i>	radiant exitance [radiant emittance] [Wm^{-2}]
<i>m</i>	linear magnification [—]
<i>m</i>	effective mass [kg]
MTF	modulation transfer function [—]
<i>N</i>	photon flux [s^{-1}]
<i>N</i>	carrier (number) density [m^{-3}]
<i>n</i>	real part of the relative refractive index [—]
\tilde{n}	complex index of refraction [—]
NA	numerical aperture [—]
OPD	optical path difference [m]
<i>P</i>	macroscopic polarization [C m^{-2}]
Re()	real part of [—]
<i>R</i>	resistance [Ω]
r	position vector [m]
<i>S</i>	Seebeck coefficient [VK^{-1}]
<i>s</i>	spin quantum number [—]
<i>s</i>	path length [m]
S_o	object distance [m]
S_i	image distance [m]
T	temperature [K, C]
<i>t</i>	time [s]
<i>t</i>	thickness [m]
<i>u</i>	slope of ray with the optical axis [rad]
<i>V</i>	Abbe reciprocal dispersion [—]
<i>V</i>	voltage [V , $\text{m}^2 \text{kg s}^{-2} \text{C}^{-1}$]
<i>x, y, z</i>	rectangular coordinates [m]
<i>Z</i>	atomic number [—]

Greek Symbols

α	absorption coefficient [cm^{-1}]
α	(power) absorptance (absorptivity)

ϵ	dielectric coefficient (constant) [—]
ϵ	emittance (emissivity) [—]
ϵ	eccentricity [—]
ϵ_1	Re (ϵ)
ϵ_2	Im (ϵ)
τ	(power) transmittance (transmissivity) [—]
ν	radiation frequency [Hz]
ω	circular frequency = $2\pi\nu$ [rads ⁻¹]
ω	plasma frequency [Hz]
λ	wavelength [μm , nm]
σ	wave number = $1/\lambda$ [cm ⁻¹]
σ	Stefan Boltzmann constant [Wm ⁻² K ⁻¹]
ρ	reflectance (reflectivity) [—]
θ, ϕ	angular coordinates [rad, °]
ξ, η	rectangular spatial frequencies [m ⁻¹ , r ⁻¹]
ϕ	phase [rad, °]
ϕ	lens power [m ⁻²]
Φ	flux [W]
χ	electric susceptibility tensor [—]
Ω	solid angle [sr]

Other

\Re	responsivity
$\exp(x)$	e^x
$\log_a(x)$	log to the base a of x
$\ln(x)$	natural log of x
$\log(x)$	standard log of x : $\log_{10}(x)$
Σ	summation
Π	product
Δ	finite difference
δx	variation in x
dx	total differential
∂x	partial derivative of x
$\delta(x)$	Dirac delta function of x
δ_{ij}	Kronecker delta

REFERENCES

1. Anonymous, *ISO Standards Handbook 2: Units of Measurement*, 2nd ed., International Organization for Standardization, 1982.
2. Anonymous, *Symbols, Units and Nomenclature in Physics*, Document U.I.P. 20, International Union of Pure and Applied Physics, 1978.
3. E. Cohen and B. Taylor, "The Fundamental Physical Constants," *Physics Today*, 9 August 1990.
4. W. J. Smith, *Modern Optical Engineering*, 2nd ed., McGraw-Hill, 1990.
5. J. M. Lloyd, *Thermal Imaging Systems*, Plenum Press, 1972.

William L. Wolfe
 College of Optical Sciences
 University of Arizona
 Tucson, Arizona

PART

1

GEOMETRICAL
OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

GENERAL PRINCIPLES OF GEOMETRICAL OPTICS

Douglas S. Goodman

Corning Tropel Corporation
Fairport, New York

1.1 GLOSSARY

(NS)	indicates nonstandard terminology
<i>italics</i>	definition or first usage
∇	gradient ($\partial/\partial x, \partial/\partial y, \partial/\partial z$)
prime, unprime	before and after, object and image space (not derivatives)
A	auxiliary function for ray tracing
A, A'	area, total field areas, object and image points
AB	directed distance from A to B
\mathbf{a}	unit axis vector, vectors
a_o, a_b, a_l	coefficients in characteristic function expansion
B	matrix element for symmetrical systems
B	auxiliary function for ray tracing
B, B'	arbitrary object and image points
\mathbf{b}	binormal unit vector of a ray path
\mathcal{B}	interspace (between) term in expansion
C	matrix element for conjugacy
$C(\mathcal{C}, \mathcal{B}, \mathcal{F})$	characteristic function
c	speed of light in vacuum
c	surface vertex curvature, spherical surface curvature
c_s	sagittal curvature
c_t	tangential curvature
D	auxiliary distance function for ray tracing
d	distance from origin to mirror
d	nominal focal distance
d, d'	arbitrary point to conjugate object, image points $d = AO, d' = A'O'$
d, d'	axial distances, distances along rays
d_H	hyperfocal distance
d_N	near focal distance
d_F	far focal distance

dA	differential area
ds	differential geometrical path length
E	image irradiance
E_0	axial image irradiance
E, E'	entrance and exit pupil locations
e	eccentricity
e_x, e_y, e_z	coefficients for collineation
F	matrix element for front side
F, F'	front and rear focal points
FN	F-number
FN_m	F-number for magnification m
$F()$	general function
$F(x, y, z)$	general surface function
f, f'	front and rear focal lengths $f = PF, f' = P'F'$
G	diffraction order
g, g'	focal lengths in tilted planes
h, h'	ray heights at objects and images, field heights
\mathcal{H}	hamiltonian
I, I'	incidence angles
\mathbf{I}	unit matrix
i, i'	paraxial incidence angles
\mathcal{J}	image space term in characteristic function expansion
L	surface x -direction cosine
L	paraxial invariant
l, l'	principal points to object and image axial points $l = PO, l' = P'O'$ axial distances from vertices of refracting surface $l = VO, l' = V'O'$
\mathcal{L}	lagrangian for heterogeneous media
M	lambertian emittance
M	surface z -direction cosine
m	transverse magnification
m_L	longitudinal magnification
m_α	angular magnification
m_E	paraxial pupil magnification
m_N	nodal point magnification = n/n'
m_p	pupil magnification in direction cosines
m_O	magnification at axial point
m_x, m_y, m_z	magnifications in the $x, y,$ and z directions
N	surface z -direction cosine
N, N'	nodal points
NA, NA'	numerical aperture
n	refractive index
\mathbf{n}	normal unit vector of a ray path
NS	nonstandard
O, O'	axial object and image points
\mathcal{O}	object space term in expansion

P	power (radiometric)
P, P'	principal points
$P(\alpha, \beta; x, y)$	pupil shape functions
$P'(\alpha', \beta'; x', y')$	
p	period of grating
\mathbf{p}	ray vector, optical direction cosine $\mathbf{p} = n \mathbf{r} = (p_x, p_y, p_z)$
p	pupil radius
p_x, p_y, p_z	optical direction cosines
$Q(\alpha, \beta; x, y)$	pupil shape functions relative to principal direction cosines
$Q'(\alpha', \beta'; x', y')$	
q	resolution parameter
q_i	coordinate for Lagrange equations
\dot{q}_i	derivative with respect to a parameter
q, q'	auxiliary functions for collineation
\mathbf{q}	unit vector along grating lines
R	matrix element for rear side
r	radius of curvature, vertex radius of curvature
\mathbf{r}	ray unit direction vector $\mathbf{r} = (\alpha, \beta, \gamma)$
\mathbf{S}	surface normal $\mathbf{S} = (L, M, N)$
$S(x, y, x', y')$	point eikonal $V(x, y, z_0; x', y', z_0')$
s	geometrical length
s	axial length
s, s'	distances associated with sagittal foci
\mathcal{S}	skew invariant
$T(\alpha, \beta; \alpha', \beta')$	angle characteristic function
t	thickness, vertex-to-vertex distance
t, t'	distances associated with tangential foci
t	time
\mathbf{t}	tangent unit vector of a ray path
U, U'	meridional ray angles relative to axis
u, u'	paraxial ray angles relative to axis
u_M	paraxial marginal ray angle
u_C	paraxial chief ray angle
u_1, u_2, u_3, u_4	homogeneous coordinates for collineation
V	optical path length
$V(\mathbf{x}; \mathbf{x}')$	point characteristic function
V, V'	vertex points
v	speed of light in medium
W_{LMN}	wavefront aberration term
W_x, W_y, W_z	wavefront aberration terms for reference shift
$W(\xi, \eta; x, y, z)$	wavefront aberration function
$W'(\alpha, \beta; x', y')$	angle-point characteristic function
$W(x, y; \alpha', \beta')$	point-angle characteristic function
$\mathbf{x} = (x, y, z)$	position vector
$\mathbf{x}(\sigma)$	parametric description of ray path

$\dot{\mathbf{x}}(\sigma)$	derivative with respect to a parameter
$\ddot{\mathbf{x}}(\sigma)$	second derivative with respect to a parameter
y	meridional ray height, paraxial ray height
y_M	paraxial marginal ray height
y_C	paraxial chief ray height
y_P, y'_P	paraxial ray height at the principal planes
z	axis of revolution
$z(\rho)$	surface sag
z_{sphere}	sag of a sphere
z_{conic}	sag of a conic
z, z'	focal point to object and image distances $z = FO, z' = F'O'$
α, β, γ	ray direction cosines
α, β, γ	entrance pupil directions
α', β', γ'	exit pupil direction cosines
α_0, β_0	principal direction of entrance pupil
α'_0, β'_0	principal direction of exit pupil
$\alpha_{\text{max}}, \alpha_{\text{min}}$	extreme pupil directions
$\beta_{\text{max}}, \beta_{\text{min}}$	extreme pupil directions
Γ	$n' \cos I' - n \cos I$
$\delta x, \delta y, \delta z$	reference point shifts
$\Delta\alpha, \Delta\beta$	angular ray aberrations
$\Delta x, \Delta y, \Delta z$	shifts
ε	surface shape parameter
$\varepsilon_x, \varepsilon_y$	transverse ray aberrations
ξ, η	pupil coordinates—not specific
θ	ray angle to surface normal
	marginal ray angle
	plane tilt angle
κ	conic parameter
κ	curvature of a ray path
λ	wavelength
ψ	azimuth angle
	field angle
ϕ	power, surface power
	azimuth
ρ	radius of curvature of a ray path
	distance from axis
	radial pupil coordinate
σ	ray path parameter
	general parameter for a curve
τ	reduced axial distances
	torsion of a ray path
$\tau(\alpha', \beta'; x', y')$	pupil transmittance function
ω, ω'	reduced angle $\omega = nu, \omega' = n'u'$
$d\omega$	differential solid angle

1.2 INTRODUCTION

The Subject

Geometrical optics is both the object of abstract study and a body of knowledge necessary for design and engineering. The subject of geometrical optics is small, since so much can be derived from a single principle, that of Fermat, and large since the consequences are infinite and far from obvious. Geometrical optics is deceptive in that much that seems simple is loaded with content and implications, as might be suggested by the fact that some of the most basic results required the likes of Newton and Gauss to discover them. Most of what appears complicated seems so because of obscuration with mathematical terminology and excessive abstraction. Since it is so old, geometrical optics tends to be taken for granted and treated too casually by those who consider it to be “understood.” One consequence is that what has been long known can be lost if it is not recirculated by successive generations of textbook authors, who are pressed to fit newer material in a fairly constant number of pages.

The Contents

The material in this chapter is intended to be that which is most fundamental, most general, and most useful to the greatest number of people. Some of this material is often thought to be more esoteric than practical, but this opinion is less related to its essence than to its typical presentation. There are no applications per se here, but everything is applicable, at least to understanding. An effort has been made to compensate here for what is lacking elsewhere and to correct some common errors. Many basic ideas and useful results have not found their way into textbooks, so are little known. Moreover, some basic principles are rarely stated explicitly. The contents are weighted toward the most common type of optical system, that with rotational symmetry consisting of mirrors and/or lens elements of homogeneous materials. There is a section “Rays in Heterogeneous Media,” an application of which is gradient index optics discussed in Chap. 24. The treatment here is mostly monochromatic. The topics of caustics and anisotropic media are omitted, and there is little specifically about systems that are not figures of revolution. The section on aberrations is short and mostly descriptive, with no discussion of lens design, a vast field concerned with the practice of aberration control. Because of space limitations, there are too few diagrams.

Terminology

Because of the complicated history of geometrical optics, its terminology is far from standardized. Geometrical optics developed over centuries in many countries, and much of it has been rediscovered and renamed. Moreover, concepts have come into use without being named, and important terms are often used without formal definitions. This lack of standardization complicates communication between workers at different organizations, each of which tends to develop its own optical dialect. Accordingly, an attempt has been made here to provide precise definitions. Terms are italicized where defined or first used. Some needed nonstandard terms have been introduced, and these are likewise italicized, as well as indicated by “NS” for “nonstandard.”

Notation

As with terminology, there is little standardization. And, as usual, the alphabet has too few letters to represent all the needed quantities. The choice here has been to use some of the same symbols more than once, rather than to encumber them with superscripts and subscripts. No symbol is used in a given section with more than one meaning. As a general practice nonprimed and primed quantities are used to indicate before and after, input and output, and object and image space.

References

No effort has been made to provide complete references, either technical or historical. (Such a list would fill the entire chapter.) The references were not chosen for priority, but for elucidation or interest, or because of their own references. Newer papers can be found by computer searches, so the older ones have been emphasized, especially since older work is receding from view beneath the current flood of papers. In geometrical optics, nothing goes out of date, and much of what is included here has been known for a century or so—even if it has been subsequently forgotten.

Communication

Because of the confusion in terminology and notation, it is recommended that communication involving geometrical optics be augmented with diagrams, graphs, equations, and numeric results, as appropriate. It also helps to provide diagrams showing both first-order properties of systems, with object and image positions, pupil positions, and principal planes, as well as direction cosine space diagrams, as required, to show angular subtenses of pupils.

1.3 FUNDAMENTALS

What Is a Ray?

Geometrical optics, which might better be called *ray optics*, is concerned with the light ray, an entity that does not exist. It is customary, therefore, to begin discussions of geometrical optics with a theoretical justification for the use of the ray. The real justification is that, like other successful models in physics, rays are indispensable to our thinking, notwithstanding their shortcomings. The ray is a model that works well in some cases and not at all in others, and light is necessarily thought about in terms of rays, scalar waves, electromagnetic waves, and with quantum physics—depending on the class of phenomena under consideration.

Rays have been defined with both corpuscular and wave theory. In corpuscular theory, some definitions are (1) the path of a corpuscle and (2) the path of a photon. A difficulty here is that energy densities can become infinite. Other efforts have been made to define rays as quantities related to the wave theory, both scalar and electromagnetic. Some are (1) wavefront normals, (2) the Poynting vector, (3) a discontinuity in the electromagnetic field,^{1,2} (4) a descriptor of wave behavior in short wavelength or high frequency limit,³ and (5) quantum mechanically.⁴ One problem with these definitions is that there are many ordinary and simple cases where wavefronts and Poynting vectors become complicated and/or meaningless. For example, in the simple case of two coherent plane waves interfering, there is no well-defined wavefront in the overlap region. In addition, rays defined in what seems to be a reasonable way can have undesirable properties. For example, if rays are defined as normals to wavefronts, then, in the case of gaussian beams, rays bend in a vacuum.

An approach that avoids the difficulties of a physical definition is that of treating rays as mathematical entities. From definitions and postulates, a variety of results is found, which may be more or less useful and valid for light. Even with this approach, it is virtually impossible to think “purely geometrically”—unless rays are treated as objects of geometry, rather than optics. In fact, we often switch between ray thinking and wave thinking without noticing it, for instance in considering the dependence of refractive index on wavelength. Moreover, geometrical optics makes use of quantities that must be calculated from other models, for example, the index of refraction. As usual, Rayleigh⁵ has put it well: “We shall, however, find it advisable not to exclude altogether the conceptions of the wave theory, for on certain most important and practical questions no conclusion can be drawn without the use of facts which are scarcely otherwise interpretable. Indeed it is not to be denied that the too rigid separation of optics into geometrical and physical has done a good deal of harm, much that is essential to a proper comprehension of the subject having fallen between the two schools.”

The ray is inherently ill-defined, and attempts to refine a definition always break down. A definition that seems better in some ways is worse in others. Each definition provides some insight into the behavior of light, but does not give the full picture. There seems to be a problem associated with the uncertainty principle involved with attempts at definition, since what is really wanted from a ray is a specification of both position and direction, which is impossible by virtue of both classical wave properties and quantum behavior. So the approach taken here is to treat rays without precisely defining them, and there are few reminders hereafter that the predictions of ray optics are imperfect.

Refractive Index

For the purposes of this chapter, the optical characteristics of matter are completely specified by its refractive index. The *index of refraction* of a medium is defined in geometrical optics as

$$n = \frac{\text{speed of light in vacuum}}{\text{speed of light in medium}} = \frac{c}{v} \quad (1)$$

A *homogeneous medium* is one in which n is the same everywhere. In an *inhomogeneous* or *heterogeneous medium* the index varies with position. In an *isotropic medium* n is the same at each point for light traveling in all directions and with all polarizations, so the index is described by a scalar function of position. Anisotropic media are not treated here.

Care must be taken with equations using the symbol n , since it sometimes denotes the ratio of indices, sometimes with the implication that one of the two is unity. In many cases, the difference from unity of the index of air (≈ 1.0003) is important. Index varies with wavelength, but this dependence is not made explicit in this chapter, most of which is implicitly limited to monochromatic light. The output of a system in polychromatic light is the sum of outputs at the constituent wavelengths.

Systems Considered

The optical systems considered here are those in which spatial variations of surface features or refractive indices are large compared to the wavelength. In such systems ray identity is preserved; there is no “splitting” of one ray into many as occurs at a grating or scattering surface.

The term *lens* is used here to include a variety of systems. *Dioptric* or *refractive* systems employ only refraction. *Catoptric* or *reflective* systems employ only reflection. *Catadioptric* systems employ both refraction and reflection. No distinction is made here insofar as refraction and reflection can be treated in a common way. And the term lens may refer here to anything from a single surface to a system of arbitrary complexity.

Summary of the Behavior and Attributes of Rays

Rays propagate in straight lines in homogeneous media and have curved paths in heterogeneous media. Rays have positions, directions, and speeds. Between any pair of points on a given ray there is a geometrical path length and an optical path length. At smooth interfaces between media with different indices rays refract and reflect. Ray paths are reversible. Rays carry energy, and power per area is approximated by ray density.

Reversibility

Rays are reversible; a path can be taken in either direction, and reflection and refraction angles are the same in either direction. However, it is usually easier to think of light as traveling along rays in a particular direction, and, of course, in cases of real instruments there usually is such a direction. The solutions to some equations may have directional ambiguity.

Groups of Rays

Certain types of groups of rays are of particular importance. Rays that originate at a single point are called a *normal congruence* or *orthotomic system*, since as they propagate in isotropic media they are associated with perpendicular wavefronts. Such groups are also of interest in image formation, where their reconvergence to a point is important, as is the path length of the rays to a reference surface used for diffraction calculations. Important in radiometric considerations are groups of rays emanating from regions of a source over a range of angles. The changes of such groups as they propagate are constrained by conservation of brightness. Another group is that of two meridional paraxial rays, related by the two-ray invariant.

Invariance Properties

Individual rays and groups of rays may have *invariance properties*—relationships between the positions, directions, and path lengths—that remain constant as a ray or a group of rays passes through an optical system.⁶ Some of these properties are completely general, e.g., the conservation of étendue and the perpendicularity of rays to wavefronts in isotropic media. Others arise from symmetries of the system, e.g., the skew invariant for rotationally symmetric systems. Other invariances hold in the paraxial limit. There are also differential invariance properties.^{7,8} Some ray properties not ordinarily thought of in this way can be thought of as invariances. For example, Snell's law can be thought of as a refraction invariant $n \sin I$.

Description of Ray Paths

A *ray path* can be described parametrically as a locus of points $\mathbf{x}(\sigma)$, where σ is any monotonic parameter that labels points along the ray. The description of curved rays is elaborated in the section on heterogeneous media.

Real Rays and Virtual Rays

Since rays in homogeneous media are straight, they can be extrapolated infinitely from a given region. The term *real* refers to the portion of the ray that “really” exists, or the accessible part, and the term *virtual* refers to the extrapolated, or inaccessible, part.

Direction

At each position where the refractive index is continuous a ray has a unique direction. The direction is given by its unit *direction vector* \mathbf{r} , whose cartesian components are direction cosines (α, β, γ) , i.e.,

$$\mathbf{r} = (\alpha, \beta, \gamma) \quad (2)$$

where $|\mathbf{r}|^2 = \alpha^2 + \beta^2 + \gamma^2 = 1$. The three direction cosines are not independent, and one is often taken to depend implicitly on the other two. In this chapter it is usually γ , which is

$$\gamma(\alpha, \beta) = \sqrt{1 - \alpha^2 - \beta^2} \quad (3)$$

Another vector with the same direction as \mathbf{r} is

$$\mathbf{p} = n\mathbf{r} = n(\alpha, \beta, \gamma) = (p_x, p_y, p_z) \quad (4)$$

where $|\mathbf{p}|^2 = n^2$. Several names are used for this vector, including the *optical direction cosine* and the *ray vector*.

Geometrical Path Length

Geometrical path length is geometrical distance measured along a ray between any two points. The differential unit of length is

$$ds = \sqrt{dx^2 + dy^2 + dz^2} \quad (5)$$

The path length between points \mathbf{x}_1 and \mathbf{x}_2 on a ray described parametrically by $\mathbf{x}(\sigma)$, with derivative $\dot{\mathbf{x}}(\sigma) = d\mathbf{x}(\sigma)/d\sigma$ is

$$s(\mathbf{x}_1; \mathbf{x}_2) = \int_{\mathbf{x}_1}^{\mathbf{x}_2} ds = \int_{\mathbf{x}_1}^{\mathbf{x}_2} \frac{ds}{d\sigma} d\sigma = \int_{\mathbf{x}_1}^{\mathbf{x}_2} \sqrt{|\dot{\mathbf{x}}(\sigma)|^2} d\sigma \quad (6)$$

Optical Path Length

The *optical path length* between two points \mathbf{x}_1 and \mathbf{x}_2 through which a ray passes is

$$\text{Optical path length} = V(\mathbf{x}_1; \mathbf{x}_2) = \int_{\mathbf{x}_1}^{\mathbf{x}_2} n(\mathbf{x}) ds = c \int_{\mathbf{x}_1}^{\mathbf{x}_2} \frac{ds}{v} = c \int dt \quad (7)$$

The integral is taken along the ray path, which may traverse homogeneous and inhomogeneous media, and include any number of reflections and refractions. Path length can be defined for virtual rays. In some cases, path length should be considered positive definite, but in others it can be either positive or negative, depending on direction.⁹ If \mathbf{x}_0 , \mathbf{x}_1 , and \mathbf{x}_2 are three points on the same ray, then

$$V(\mathbf{x}_0; \mathbf{x}_2) = V(\mathbf{x}_0; \mathbf{x}_1) + V(\mathbf{x}_1; \mathbf{x}_2) \quad (8)$$

Equivalently, the time required for light to travel between the two points is

$$\text{Time} = \frac{\text{optical path length}}{c} = \frac{V}{c} = \frac{1}{c} \int_{\mathbf{x}_1}^{\mathbf{x}_2} n(\mathbf{x}) ds = \int_{\mathbf{x}_1}^{\mathbf{x}_2} \frac{ds}{v} \quad (9)$$

In homogeneous media, rays are straight lines, and the optical path length is $V = n \int ds = (\text{index}) \times (\text{distance between the points})$.

The optical path length integral has several interpretations, and much of geometrical optics involves the examination of its meanings. (1) With both points fixed, it is simply a scalar, the optical path length from one point to another. (2) With one point fixed, say \mathbf{x}_0 , then treated as a function of \mathbf{x} , the surfaces $V(\mathbf{x}_0; \mathbf{x}) = \text{constant}$ are geometrical wavefronts for light originating at \mathbf{x}_0 . (3) Most generally, as a function of both arguments $V(\mathbf{x}_1; \mathbf{x}_2)$ is the *point characteristic function*, which contains all the information about the rays between the region containing \mathbf{x}_1 and that containing \mathbf{x}_2 . There may not be a ray between all pairs of points.

Fermat's Principle

According to Fermat's principle¹⁰⁻¹⁵ the optical path between two points through which a ray passes is an extremum. Light passing through these points along any other nearby path would take either more or less time. The principle applies to different *neighboring* paths. The optical path length of a ray may not be a global extremum. For example, the path lengths of rays through different facets of a Fresnel lens have no particular relationship. Fermat's principle applies to entire systems, as well as to any portion of a system, for example, to any section of a ray. In a homogeneous medium, the extremum is a straight line or, if there are reflections, a series of straight line segments.

The extremum principle can be described mathematically as follows.¹⁶ With the end points fixed, if a nonphysical path differs from a physical one by an amount proportional to δ , the nonphysical optical path length differs from the actual one by a quantity proportional to δ^2 or to a higher order. If the order is three or higher, the first point is imaged at the second-to-first order. Roughly speaking,

the higher the order, the better the image. A point is imaged stigmatically when a continuum of neighboring paths have the same length, so the equality holds to all orders. If they are sufficiently close, but vary slightly, the deviation from equality is a measure of the aberration of the imaging. An extension of Fermat's principle is given by Hopkins.¹⁷

Ray and wave optics are related by the importance of path length in both.^{18,19} In wave optics, optical path length is proportional to phase change, and the extremum principle is associated with constructive interference. The more alike the path lengths are from an object point to its image, the less the differences in phase of the wave contributions, and the greater the magnitude of the net field. In imaging this connection is manifested in the relationship of the wavefront aberration and the eikonal.

Fermat's principle is a unifying principle of geometrical optics that can be used to derive laws of reflection and refraction, and to find the equations that describe ray paths and geometrical wavefronts in heterogeneous and homogeneous media. It is one of a number of variational principles based historically on the idea that nature is economical, a unifying principle of physics. The idea that the path length is an extremum could be used mathematically without interpreting the refractive index in terms of the speed of light.

Geometrical Wavefronts

For rays originating at a single point, a *geometrical wavefront* is a surface that is a locus of constant optical path length from the source. If the source point is located at \mathbf{x}_0 and light leaves at time t_0 , then the wavefront at time t is given by

$$V(\mathbf{x}_0; \mathbf{x}) = c(t - t_0) \quad (10)$$

The function $V(\mathbf{x}; \mathbf{x}_0)$, as a function of \mathbf{x} , satisfies the *eikonal equation*

$$\begin{aligned} n(\mathbf{x})^2 &= \left(\frac{\partial V}{\partial x}\right)^2 + \left(\frac{\partial V}{\partial y}\right)^2 + \left(\frac{\partial V}{\partial z}\right)^2 \\ &= |\Delta V(\mathbf{x}; \mathbf{x}_0)|^2 \end{aligned} \quad (11)$$

This equation can also be written in relativistic form, with a four-dimensional gradient as $0 = \Sigma(\partial V / \partial x_i)$.²⁰

For constant refractive index, the eikonal equation has some simple solutions, one of which is $V = n[\alpha(x - x_0) + \beta(y - y_0) + \gamma(z - z_0)]$, corresponding to a parallel bundle of rays with directions (α, β, γ) . Another is $V = n[(x - x_0)^2 + (y - y_0)^2 + (z - z_0)^2]^{1/2}$, describing rays traveling radially from a point (x_0, y_0, z_0) .

In isotropic media, the rays and wavefronts are everywhere perpendicular to each other, a condition referred to as *orthotomic*. According to the *Malus-Dupin principle*, if a group of rays emanating from a single point is reflected and/or refracted any number of times, the perpendicularity of rays to wavefronts is maintained. The direction of a ray from \mathbf{x}_0 at \mathbf{x} is that of the gradient of $V(\mathbf{x}_0; \mathbf{x})$

$$\mathbf{p} = n\mathbf{r} = \nabla V$$

or

$$n\alpha = \frac{\partial V}{\partial x} \quad n\beta = \frac{\partial V}{\partial y} \quad n\gamma = \frac{\partial V}{\partial z} \quad (12)$$

In a homogeneous medium, all wavefronts can be found from any one wavefront by a construction. Wavefront normals, i.e., rays, are projected from the known wavefront, and loci of points equidistant therefrom are other wavefronts. This gives wavefronts in both directions, that is, both subsequent and previous wavefronts. (A single wavefront contains no directional information.) The construction also gives virtual wavefronts, those which would occur or would have occurred if the medium extended infinitely. This construction is related to that of Huygens for wave optics. At each point on a wavefront there are two principal curvatures, so there are two foci along each ray and two caustic surfaces.^{8,21}

The geometrical wavefront is analogous to the surface of constant phase in wave optics, and the eikonal equation can be obtained from the wave equation in the limit of small wavelength.^{3,4} A way in which wave optics differs from ray optics is that the phase fronts can be modified by phase changes that occur on reflection, transmission, or in passing through foci.

Fields of Rays

In many cases the optical direction cosine vectors \mathbf{p} form a field, where the optical path length is the potential, and the geometrical wavefronts are equipotential surfaces. The potential changes with position according to

$$dV = n\alpha dx + n\beta dy + n\gamma dz = n\mathbf{r} \cdot d\mathbf{x} = \mathbf{p} \cdot d\mathbf{x} \quad (13)$$

If $d\mathbf{x}$ is in the direction of a ray, then $dV/dx = n$, the maximum rate of change. If $d\mathbf{x}$ is perpendicular to a ray, then $dV/dx = 0$. The potential difference between any two wavefronts is

$$V_2 - V_1 = \int_{\mathbf{x}_1}^{\mathbf{x}_2} dV \quad (14)$$

where \mathbf{x}_1 and \mathbf{x}_2 are any two points on the respective wavefronts, and the integrand is independent of the path. Other relationships for rays originating at a single point are

$$0 = \nabla \times \mathbf{p} = \nabla \times (n\mathbf{r}) \quad \text{and} \quad 0 = \oint \mathbf{p} \cdot d\mathbf{x} \quad (15)$$

where the integral is about a closed path.³ These follow since \mathbf{p} is a gradient, Eq. (13). In regions where the rays are folded onto themselves by refraction or reflections, \mathbf{p} and V are not single-valued, so there is not a field.

1.4 CHARACTERISTIC FUNCTIONS

Introduction

Characteristic functions contain all the information about the path lengths between pairs of points, which may either be in a contiguous region or physically separated, e.g., on the two sides of a lens. These functions were first considered by Hamilton,²² so their study is referred to as *hamiltonian optics*. They were rediscovered in somewhat different form by Bruns^{23,24} and referred to as eikonals, leading to a confusing set of names for the various functions. The subject is discussed in a number of books.²⁵⁻³⁶

Four parameters are required to specify a ray. For example, an input ray is defined in the $z = 0$ plane by coordinates (x, y) and direction (α, β) . So four functions of four variables specify how an incident ray emerges from a system. In an output plane $z' = 0$, the ray has coordinates $x' = x'(x, y, \alpha, \beta)$, $y' = y'(x, y, \alpha, \beta)$, and directions $\alpha' = \alpha'(x, y, \alpha, \beta)$, $\beta' = \beta'(x, y, \alpha, \beta)$. Because of Fermat's principle, these four functions are not independent, and the geometrical optics properties of a system can be fully characterized by a single function.³²

For any given system, there is a variety of characteristic functions related by Legendre transformations, with different combinations of spatial and angular variables.³⁴ The different functions are suited for different types of analysis. *Mixed* characteristic functions have both spatial and angular arguments. Those functions that are of most general use are discussed next. The others may be useful in special circumstances. If the regions have constant refractive indices, the volumes over which the characteristic functions are defined can be extended virtually from physically accessible to inaccessible regions.

From any of its characteristic functions, all the properties of a system involving ray paths can be found, for example, ray positions, directions, and geometrical wavefronts. An important use of

the characteristic functions is demonstrating general principles and fundamental limitations. Much of this can be done by using the general properties, e.g., symmetry under rotation. (Unfortunately, it is not always known how closely the impossible can be approached.)

Point Characteristic Function

The *point characteristic function* is the optical path integral $V(\mathbf{x}; \mathbf{x}') = V(x, y, z; x', y', z')$ taken as a function of both points \mathbf{x} and \mathbf{x}' . At point \mathbf{x} where the index is n ,

$$-n\alpha = \frac{\partial V}{\partial x} \quad -n\beta = \frac{\partial V}{\partial y} \quad -n\gamma = \frac{\partial V}{\partial z} \quad \text{or} \quad -\mathbf{p} = \nabla V \quad (16)$$

Similarly, at \mathbf{x}' , where the index is n' ,

$$n'\alpha' = \frac{\partial V}{\partial x'} \quad n'\beta' = \frac{\partial V}{\partial y'} \quad n'\gamma' = \frac{\partial V}{\partial z'} \quad \text{or} \quad \mathbf{p}' = \nabla' V \quad (17)$$

It follows from the above equations and Eq. (4) that the point characteristic satisfies two conditions:

$$n^2 = |\nabla V|^2 \quad \text{and} \quad n'^2 = |\nabla' V|^2 \quad (18)$$

Therefore, the point characteristic is not an arbitrary function of six variables. The total differential of V is

$$dV(\mathbf{x}; \mathbf{x}') = \mathbf{p}' \cdot d\mathbf{x}' - \mathbf{p} \cdot d\mathbf{x} \quad (19)$$

“This expression can be said to contain all the basic laws of optics”.³⁶

Point Eikonal

If reference planes in object and image spaces are fixed, for which we use z_0 and z'_0 , then the *point eikonal* is $S(x, y; x', y') = V(x, y, z_0; x', y', z'_0)$. This is the optical path length between pairs of points on the two planes. The function is not useful if the planes are conjugate, since more than one ray through a pair of points can have the same path length. The function is arbitrary, except for the requirement³⁷ that

$$\frac{\partial^2 S}{\partial x \partial x'} \frac{\partial^2 S}{\partial y \partial y'} - \frac{\partial^2 S}{\partial x \partial y'} \frac{\partial^2 S}{\partial x' \partial y} \neq 0 \quad (20)$$

The partial derivatives of the point eikonal are

$$-n\alpha = \frac{\partial S}{\partial x} \quad -n\beta = \frac{\partial S}{\partial y} \quad \text{and} \quad n'\alpha' = \frac{\partial S}{\partial x'} \quad n'\beta' = \frac{\partial S}{\partial y'} \quad (21)$$

The relative merits of the point characteristic function and point eikonal have been debated.³⁷⁻³⁹

Angle Characteristic

The *angle characteristic function* $T(\alpha, \beta; \alpha', \beta')$, also called the *eikonal*, is related to the point characteristic by

$$T(\alpha, \beta; \alpha', \beta') = V(x, y, z; x', y', z') + n(\alpha x + \beta y + \gamma z) - n'(\alpha' x' + \beta' y' + \gamma' z') \quad (22)$$

Here the input plane z and output plane z' are fixed and are implicit parameters of T .



FIGURE 1 Geometrical interpretation of the angle characteristic function for constant object and image space indices. There is, in general, a single ray with directions (α, β, γ) in object space and $(\alpha', \beta', \gamma')$ in image space. Point O is the coordinate origin in object space, and O' is that in image space. From the origins, perpendiculars to the ray are constructed, which intersect the ray at Q and Q' . The angle characteristic function $T(\alpha, \beta; \alpha', \beta')$ is the path length from Q to Q' .

This equation is really shorthand for a Legendre transformation to coordinates $p_x = \partial V / \partial x$, etc. In principle, the expressions of Eq. (16) are used to solve for x and y in terms of α and β , and likewise Eq. (17) gives x' and y' in terms of α' and β' , so

$$\begin{aligned} T(\alpha, \beta; \alpha', \beta') = & V[x(\alpha, \beta), y(\alpha, \beta), z; x'(\alpha', \beta'), y'(\alpha', \beta'), z'] \\ & + n[\alpha x(\alpha, \beta) + \beta y(\alpha, \beta) + \sqrt{1 - \alpha^2 - \beta^2} z] \\ & - n'[\alpha' x'(\alpha', \beta') + \beta' y'(\alpha', \beta') + \sqrt{1 - \alpha'^2 - \beta'^2} z'] \end{aligned} \quad (23)$$

The angle characteristic is an arbitrary function of four variables that completely specify the directions of rays in two regions. This function is not useful if parallel incoming rays give rise to parallel outgoing rays, as is the case with afocal systems, since the relationship between incoming and outgoing directions is not unique. The partial derivatives of the angular characteristic function are

$$\frac{\partial T}{\partial \alpha} = n \left(x - \frac{\alpha}{\gamma} z \right), \quad \frac{\partial T}{\partial \beta} = n \left(y - \frac{\beta}{\gamma} z \right) \quad (24)$$

$$\frac{\partial T}{\partial \alpha'} = -n' \left(x' - \frac{\alpha'}{\gamma'} z' \right), \quad \frac{\partial T}{\partial \beta'} = -n' \left(y' - \frac{\beta'}{\gamma'} z' \right) \quad (25)$$

These expressions are simplified if the reference planes are taken to be $z = 0$ and $z' = 0$. The geometrical interpretation of T is that it is the path length between the intersection point of rays with perpendicular planes through the coordinate origins in the two spaces, as shown in Fig. 1 for the case of constant n and n' . If the indices are heterogeneous, the construction applies to the tangents to the rays. Of all the characteristic functions, T is most easily found for single surfaces and most easily concatenated for series of surfaces.

Point-Angle Characteristic

The *point-angle characteristic function* is a mixed function defined by

$$\begin{aligned} W(x, y, z; \alpha', \beta') = & V(x, y, z; x', y', z') - n'(\alpha' x' + \beta' y' + \gamma' z') \\ = & T(\alpha, \beta; \alpha', \beta') - n(\alpha x + \beta y + \gamma z) \end{aligned} \quad (26)$$

As with Eq. (22), this equation is to be understood as shorthand for a Legendre transformation. The partial derivatives with respect to the spatial variables are related by equations like those of Eq. (16), so $n^2 = |\nabla W|^2$, and the derivatives with respect to the angular variables are like those of Eq. (25). This function is useful for examining transverse ray aberrations for a given object point, since $\partial W/\partial \alpha'$, $\partial W/\partial \beta'$ give the intersection points (x', y') in plane z for rays originating at (x, y) in plane z .

Angle-Point Characteristic

The *angle-point characteristic function* is

$$\begin{aligned} W'(\alpha, \beta; x', y', z') &= V(x, y, z; x', y', z') + n(\alpha x + \beta y + \gamma z) \\ &= T(\alpha, \beta; \alpha', \beta') - n'(\alpha' x' + \beta' y' + \gamma' z) \end{aligned} \quad (27)$$

Again, this is shorthand for the Legendre transformation. This function satisfies relationships like those of Eq. (17) and satisfies $n'^2 = |\nabla' W'|^2$. Derivatives with respect to spatial variables are like those of Eq. (21). It is useful when input angles are given and output angles are to be found.

Expansions About an Arbitrary Ray

If two points on a ray that are not conjugate are taken as coordinate origins, and the z axes of the coordinate systems are taken to lie along the rays, then the expansion to second order of the point eikonal about these points is

$$\begin{aligned} S(x_1, y_1; x_2, y_2) &= v + a_1 x_1^2 + b_1 x_1 y_1 + c_1 y_1^2 + a_2 x_2^2 + b_2 x_2 y_2 + c_2 y_2^2 \\ &\quad + dx_1 x_2 + ey_1 y_2 + fx_1 y_2 + gy_1 x_2 \end{aligned} \quad (28)$$

The other characteristic functions have similar expansions. These expansions have three types of terms, those associated with the input space, the output space, and “interspace” terms. From the coefficients, information about imaging along a known ray is obtained. This subject is treated in the references for the section “Images About Known Rays.”

Expansions About the Axis

For rotationally symmetric systems, the building blocks for an expansion about the axis are

$$\text{Object space term: } \mathcal{O} = x^2 + y^2 \quad \text{or} \quad \alpha^2 + \beta^2 \quad (29)$$

$$\text{Image space term: } \mathcal{I} = x'^2 + y'^2 \quad \text{or} \quad \alpha'^2 + \beta'^2 \quad (30)$$

$$\begin{aligned} \text{Interspace term: } \mathcal{B} = xx' + yy' \quad \text{or} \quad \alpha\alpha' + \beta\beta' \quad \text{or} \quad x\alpha' + y\beta' \\ \text{or} \quad \alpha x' + \beta y' \end{aligned} \quad (31)$$

(Here $\mathcal{B} \equiv$ “between.”) The interspace term combines the variables included in \mathcal{O} and \mathcal{I} . The general form can be written as a series

$$C(\mathcal{O}, \mathcal{B}, \mathcal{I}) = \sum_{L,M,N} a_{LMN} \mathcal{O}^L \mathcal{B}^M \mathcal{I}^N \quad (32)$$

To second order, the expansion is

$$\begin{aligned} C(\mathcal{O}, \mathcal{B}, \mathcal{I}) &= a_0 + a_{100} \mathcal{O} + a_{010} \mathcal{B} + a_{001} \mathcal{I} + a_{200} \mathcal{O}^2 + a_{020} \mathcal{B}^2 + a_{002} \mathcal{I}^2 \\ &\quad + a_{110} \mathcal{O} \mathcal{B} + a_{101} \mathcal{O} \mathcal{I} + a_{011} \mathcal{B} \mathcal{I} \end{aligned} \quad (33)$$

The constant term is the optical path length between coordinate origins in the two spaces. It is often unimportant, but it does matter if two systems are used in parallel, as in an interferometer. The three first-order terms give the paraxial approximation. For imaging systems, the second-order terms are associated with third-order ray aberrations, and so on.³⁰ It is also possible to expand the characteristic functions in terms of three linear combinations of \mathcal{O} , \mathcal{B} , and \mathcal{I} . These combinations can be chosen so that the characteristic function of an aberration-free system depends on only one of the three terms, and the other two describe the aberrations.^{26,31,40}

Paraxial Forms for Rotationally Symmetric Systems

These functions contain one each of the object space, image space, and interspace terms, with coefficients a_O , a_I , and a_B . The coefficients of the object and image space terms depend on the input and output plane locations. That of the interspace term depends on the system power. Point eikonal:

$$S(x', y'; x, y) = a + a_O(x^2 + y^2) + a_B(xx' + yy') + a_I(x'^2 + y'^2) \quad (34)$$

Angle characteristic:

$$T(\alpha', \beta'; \alpha, \beta) = a + a_O(\alpha^2 + \beta^2) + a_B(\alpha\alpha' + \beta\beta') + a_I(\alpha'^2 + \beta'^2) \quad (35)$$

Point-angle characteristic:

$$W(x, y; \alpha', \beta') = a + a_O(x^2 + y^2) + a_B(x\alpha' + y\beta') + a_I(\alpha'^2 + \beta'^2) \quad (36)$$

Angle-point characteristic:

$$W'(\alpha, \beta, x', y') = a + a_O(\alpha^2 + \beta^2) + a_B(\alpha x' + \beta y') + a_I(x'^2 + y'^2) \quad (37)$$

The corresponding coefficients in these expressions are different from each other. The familiar properties of paraxial and gaussian optics can be found from these functions by taking the appropriate partial derivatives.

Some Ideal Characteristic Functions

For a system that satisfies certain conditions, the form of a characteristic function can sometimes be found. Thereafter, some of its properties can be determined. Some examples of characteristic functions follow, in each of which expression the function F is arbitrary.

For maxwellian perfect imaging (defined below) by a rotationally symmetric system between planes at $z = 0$ and $z' = 0$ related by transverse magnification m , the point characteristic function, defined for $z' \neq 0$, is

$$V(x', y', z'; x, y) = F(x^2 + y^2) + [(x' - mx)^2 + (y' - my)^2 + z'^2]^{1/2} \quad (38)$$

Expanding the expression above for small x, x', y, y' give the paraxial form, Eq. (34). The form of the point-angle characteristic is

$$W(x, y; \alpha', \beta') = F(x^2 + y^2) - m(n'\alpha'x + n\beta'y) \quad (39)$$

The form of the angle-point characteristic is

$$W'(\alpha, \beta; x', y') = F(x'^2 + y'^2) + \frac{1}{m}(n\alpha x' + n\beta y') \quad (40)$$

The functions F are determined if the imaging is also stigmatic at one additional point, for example, at the center of the pupil.^{26,34,40,41} The angular characteristic function has the form

$$T(\alpha, \beta; \alpha', \beta') = F[(n\alpha - mn'\alpha')^2 + (n\beta - mn'\beta')^2] \quad (41)$$

where F is any function.

For a lens of power ϕ that stigmatically images objects at infinity in a plane, and does so in either direction,

$$S(x, y; x', y') = -\phi(xx' + yy') \quad \text{and} \quad T(\alpha, \beta; \alpha', \beta') = \frac{nn'}{\phi}(\alpha\alpha' + \beta\beta') \quad (42)$$

Partially differentiating with respect to the appropriate variables shows that for such a system, the heights of point images in the rear focal plane are proportional to the sines of the incident angles, rather than the tangents.

1.5 RAYS IN HETEROGENEOUS MEDIA

Introduction

This section provides equations for describing and determining the curved ray paths in a heterogeneous or inhomogeneous medium, one whose refractive index varies with position. It is assumed here that $n(\mathbf{x})$ and the other relevant functions are continuous and have continuous derivatives to whatever order is needed. Various aspects of this subject are discussed in a number of books and papers.^{42–49} This material is often discussed in the literature on gradient index lenses^{50–54} and in discussions of microwave lenses.^{55–58}

Differential Geometry of Space Curves

A curved ray path is a space curve, which can be described by a standard parametric description, $\mathbf{x}(\sigma) = [x(\sigma), y(\sigma), z(\sigma)]$, where σ is an arbitrary parameter.^{46,59–62}

Different parameters may be used according to the situation. The path length s along the ray is sometimes used, as is the axial position z . Some equations change form according to the parameter, and those involving derivatives are simplest when the parameter is s . Derivatives with respect to the parameter are denoted by dots, so $\dot{\mathbf{x}}(\sigma) = d\mathbf{x}(\sigma)/d\sigma = [\dot{x}(\sigma), \dot{y}(\sigma), \dot{z}(\sigma)]$. A parameter other than s is a function of s , so $d\mathbf{x}(\sigma)/ds = (d\mathbf{x}/d\sigma)(d\sigma/ds)$.

Associated with space curves are three mutually perpendicular unit vectors, the tangent vector \mathbf{t} , the principal normal \mathbf{n} , and the binormal \mathbf{b} , as well as two scalars, the curvature and the torsion. The direction of a ray is that of its unit *tangent vector*

$$\mathbf{t} = \frac{\dot{\mathbf{x}}(\sigma)}{|\dot{\mathbf{x}}(\sigma)|} = \dot{\mathbf{x}}(s) = (\alpha, \beta, \gamma) \quad (43)$$

The tangent vector \mathbf{t} is the same as the direction vector \mathbf{r} used elsewhere in this chapter. The rate of change of the tangent vector with respect to path length is

$$\kappa \mathbf{n} = \dot{\mathbf{t}}(s) = \ddot{\mathbf{x}}(s) = \left(\frac{d\alpha}{ds}, \frac{d\beta}{ds}, \frac{d\gamma}{ds} \right) \quad (44)$$

The *normal vector* is the unit vector in this direction

$$\mathbf{n} = \frac{\ddot{\mathbf{x}}(s)}{|\ddot{\mathbf{x}}(s)|} \quad (45)$$

The vectors \mathbf{t} and \mathbf{n} define the *osculating plane*. The *curvature* $\kappa = |\ddot{\mathbf{x}}(s)|$ is the rate of change of direction of \mathbf{t} in the osculating plane.

$$\kappa^2 = \frac{|\dot{\mathbf{x}}(\sigma) \times \ddot{\mathbf{x}}(\sigma)|^2}{|\dot{\mathbf{x}}(\sigma)|^6} = |\ddot{\mathbf{x}}(s)|^2 = \left(\frac{d\alpha}{ds} \right)^2 + \left(\frac{d\beta}{ds} \right)^2 + \left(\frac{d\gamma}{ds} \right)^2 \quad (46)$$

The radius of curvature is $\rho = 1/\kappa$. Perpendicular to the osculating plane is the unit *binormal vector*

$$\mathbf{b} = \mathbf{t} \times \mathbf{n} = \frac{\dot{\mathbf{x}}(s) \times \ddot{\mathbf{x}}(s)}{|\ddot{\mathbf{x}}(s)|} \quad (47)$$

The *torsion* is the rate of change of the normal to the osculating plane

$$\tau = \mathbf{b}(s) \cdot \frac{d\mathbf{n}(s)}{ds} = \frac{[\dot{\mathbf{x}}(\sigma) \times \ddot{\mathbf{x}}(\sigma)] \cdot \ddot{\mathbf{x}}(\sigma)}{|\dot{\mathbf{x}}(\sigma) \times \ddot{\mathbf{x}}(\sigma)|^2} = \frac{[\dot{\mathbf{x}}(s) \times \ddot{\mathbf{x}}(s)] \cdot \ddot{\mathbf{x}}(s)}{|\dot{\mathbf{x}}(s)|^2} \quad (48)$$

The quantity $1/\tau$ is the *radius of torsion*. For a plane curve, $\tau = 0$ and \mathbf{b} is constant. The rates of change of \mathbf{t} , \mathbf{n} , and \mathbf{b} are given by the Frenet equations:

$$\dot{\mathbf{t}}(s) = \kappa \mathbf{n}, \quad \dot{\mathbf{n}}(s) = -\kappa \mathbf{t} + \tau \mathbf{b}, \quad \dot{\mathbf{b}}(s) = -\tau \mathbf{n} \quad (49)$$

In some books, $1/\kappa$ and $1/\tau$ are used for what are denoted here by κ and τ .

Differential Geometry Equations Specific to Rays

From the general space curve equations above and the differential equations below specific to rays, the following equations for rays are obtained. Note that n here is the refractive index, unrelated to \mathbf{n} . The tangent and normal vectors are related by Eq. (59), which can be written

$$\nabla \log n = \kappa \mathbf{n} + (\nabla \log n \cdot \mathbf{t}) \mathbf{t} \quad (50)$$

The osculating plane always contains the vector ∇n . Taking the dot product with \mathbf{n} in the above equation gives

$$\kappa = \frac{\partial \log n}{\partial N} = \mathbf{n} \cdot \nabla \log n = \mathbf{b} \cdot (\dot{\mathbf{x}} \times \nabla \log n) \quad (51)$$

The partial derivative $\partial/\partial N$ is in the direction of the principal normal, so rays bend toward regions of higher refractive index. Other relations⁴⁶ are

$$\mathbf{n} = \rho \dot{\mathbf{x}}(s) \times [\nabla \log n \times \dot{\mathbf{x}}(s)] \quad (52)$$

$$\mathbf{b} = \rho \dot{\mathbf{x}}(s) \times \nabla \log n \quad \text{and} \quad 0 = \mathbf{b} \cdot \nabla n \quad (53)$$

$$\tau = \frac{(\dot{\mathbf{x}}(s) \times \nabla n) \cdot \nabla \dot{n}}{|\nabla n \times \dot{\mathbf{x}}(s)|^2} \quad (54)$$

Variational Integral

Written in terms of parameter σ , the optical path length integral, Eq. (7), is

$$V = \int n ds = \int \left(n \frac{ds}{d\sigma} \right) d\sigma = \int \mathcal{L} d\sigma \quad (55)$$

The solution for ray paths involves the calculus of variations in a way analogous to that used in classical mechanics, where the time integral of the lagrangian \mathcal{L} is an extremum.⁶³ If \mathcal{L} has no explicit dependence on σ , the mechanical analogue to the optics case is of no explicit time dependence.

Differential Equations for Rays

General Differential Equations Because the optical path length integral is an extremum, the integrand \mathcal{L} satisfies the Euler equations.⁴⁶ For an arbitrary coordinate system, with coordinates q_1, q_2, q_3 and the derivatives with respect to the parameter $\dot{q}_i = dq_i/d\sigma$, the differential equations for the path are

$$0 = \frac{d}{d\sigma} \frac{\partial \mathcal{L}}{\partial \dot{q}_i} - \frac{\partial \mathcal{L}}{\partial q_i} = \frac{d}{d\sigma} \left(n \frac{\partial}{\partial \dot{q}_i} \frac{ds}{d\sigma} \right) - \frac{\partial}{\partial q_i} \left(n \frac{ds}{d\sigma} \right); \quad i=1, 2, 3 \quad (56)$$

Cartesian Coordinates with Unspecified Parameter In cartesian coordinates $ds/d\sigma = (\dot{x}^2 + \dot{y}^2 + \dot{z}^2)^{1/2}$, so the x equation is

$$0 = \frac{d}{d\sigma} \left(n \frac{\partial}{\partial \dot{x}} \frac{ds}{d\sigma} \right) - \frac{ds}{d\sigma} \frac{\partial n}{\partial x} = \frac{d}{d\sigma} \left[\frac{n\dot{x}}{(\dot{x}^2 + \dot{y}^2 + \dot{z}^2)^{1/2}} \right] - (\dot{x}^2 + \dot{y}^2 + \dot{z}^2)^{1/2} \frac{\partial n}{\partial x} \quad (57)$$

Similar equations hold for y and z .

Cartesian Coordinates with Parameter $\sigma = s$ With $\sigma = s$, so $ds/d\sigma = 1$, an expression, sometimes called the *ray equation*, is obtained.²⁸

$$\nabla n = \frac{d}{ds} \left(n \frac{d\mathbf{x}(s)}{ds} \right) = n \frac{d^2\mathbf{x}(s)}{ds^2} + \frac{dn[\mathbf{x}(s)]}{ds} \frac{d\mathbf{x}(s)}{ds} \quad (58)$$

Using $dn/ds = \nabla n \cdot \dot{\mathbf{x}}$, the ray equation can also be written

$$\nabla n = n\ddot{\mathbf{x}} + (\nabla n \cdot \dot{\mathbf{x}})\dot{\mathbf{x}} \quad \text{or} \quad \nabla \log n = \ddot{\mathbf{x}} + (\nabla \log n \cdot \dot{\mathbf{x}})\dot{\mathbf{x}} \quad (59)$$

Only two of the component equations are independent, since $|\dot{\mathbf{x}}| = 1$.

Cartesian Coordinates with Parameter $\sigma = \int ds/n$ The parameter $\sigma = \int ds/n$, for which $ds/d\sigma = n$ and $n^2 = \dot{x}^2 + \dot{y}^2 + \dot{z}^2$, gives⁴⁴

$$\frac{d^2\mathbf{x}}{d\sigma^2} = \nabla \left(\frac{1}{2}n^2 \right) \quad (60)$$

This equation is analogous to Newton's law of motion for a particle, $\mathbf{F} = m d^2\mathbf{x}/dt^2$, so the ray paths are like the paths of particles in a field with a potential proportional to $n^2(\mathbf{x})$. This analogy describes paths, but not speeds, since light travels slower where n is greater, whereas the particles would have greater speeds.^{64,65}

Euler Equations for Parameter $\sigma = z$ If $\sigma = z$, then $ds/d\sigma = (\dot{x}^2 + \dot{y}^2 + 1)^{1/2}$ and $\mathcal{L} = \mathcal{L}(x, y; \dot{x}, \dot{y}; z)$. This gives^{45,49}

$$0 = \frac{d}{dz} \left(n \frac{\partial}{\partial \dot{x}} \frac{ds}{dz} \right) - \frac{ds}{dz} \frac{\partial n}{\partial x} = \frac{d}{dz} \left[\frac{n\dot{x}}{(1 + \dot{x}^2 + \dot{y}^2)^{1/2}} \right] - (1 + \dot{x}^2 + \dot{y}^2)^{1/2} \frac{\partial n}{\partial x} \quad (61)$$

with a similar equation for y . The equations can also be written (Refs. 51, app. A, and 66) as

$$n\ddot{x} = (1 + \dot{x}^2 + \dot{y}^2) \left(\frac{\partial n}{\partial x} - \frac{\partial n}{\partial z} \dot{x} \right) \quad n\ddot{y} = (1 + \dot{x}^2 + \dot{y}^2) \left(\frac{\partial n}{\partial y} - \frac{\partial n}{\partial z} \dot{y} \right) \quad (62)$$

This parameter is particularly useful when n is rotationally symmetric about the z axis.

Hamilton's Equations with Cartesian Coordinates for Parameter $\sigma = z$ A set of Hamilton's equations can also be written in cartesian coordinates using z as the parameter.^{45,49} The canonical momenta in cartesian coordinates are the optical direction cosines

$$p_x = \frac{\partial \mathcal{L}}{\partial \dot{x}} = n\alpha \quad p_y = \frac{\partial \mathcal{L}}{\partial \dot{y}} = n\beta \quad (63)$$

The hamiltonian is

$$\mathcal{H}(x, y; p_x, p_y; z) = \dot{x}p_x + \dot{y}p_y - \mathcal{L} = -\sqrt{n^2(x, y, z) - (p_x^2 + p_y^2)} \quad (64)$$

Hamilton's equations are

$$\frac{dx}{dz} = \frac{\partial \mathcal{H}}{\partial p_x}, \quad \frac{dy}{dz} = \frac{\partial \mathcal{H}}{\partial p_y}, \quad \frac{dp_x}{dz} = -\frac{\partial \mathcal{H}}{\partial x}, \quad \frac{dp_y}{dz} = -\frac{\partial \mathcal{H}}{\partial y} \quad (65)$$

It is not possible to write a set of Hamilton's equations using an arbitrary parameter and three canonical momenta, since they are not independent.⁶⁷ Another equation is

$$\frac{\partial \mathcal{H}}{\partial z} = \frac{d\mathcal{H}}{dz} = \frac{1}{\gamma} \frac{\partial n}{\partial z} \quad (66)$$

Paraxial Form of Hamilton's Equations for $\sigma = z$ In the paraxial limit, if n_0 is the average index, the above set of equations gives⁴⁹

$$\frac{d^2x(z)}{dz^2} = \frac{1}{n_0} \frac{\partial n}{\partial x}, \quad \frac{d^2y(z)}{dz^2} = \frac{1}{n_0} \frac{\partial n}{\partial y} \quad (67)$$

Other Forms A variety of additional differential equations can be obtained with various parameters.⁶⁷ Time cannot be used as a parameter.⁶⁸ The equations can also be expressed in a variety of coordinate systems.^{56,58,69-71}

Refractive Index Symmetries

When the refractive index has symmetry or does not vary with one or more of the spatial variables, the above equations may simplify and take special forms. If, in some coordinate system, n does not vary with a coordinate q_i , so $\partial n / \partial q_i = 0$, and if, in addition, $\partial / \partial q_i (ds/d\sigma) = 0$, then

$$\frac{\partial \mathcal{L}}{\partial q_i} = 0 \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial \dot{q}_i} = n \frac{\partial}{\partial \dot{q}_i} \left(\frac{ds}{d\sigma} \right) = \text{constant} \quad (68)$$

There is an associated invariance of the ray path.^{44,49,56,58} (This is analogous to the case in mechanics where a potential does not vary with some coordinate.) A more esoteric approach to symmetries involves Noether's theorem.^{72,73} There are a number of special cases.

If the index is rotationally symmetric about the z axis, $n = n(x^2 + y^2, z)$, then $\partial \mathcal{L} / \partial \phi = 0$, where ϕ is the azimuth angle, and the constant of motion is analogous to that of the z component of angular momentum in mechanics for a potential with rotational symmetry. The constant quantity is the *skew invariant*, discussed in the section "Skew Invariant."

If the refractive index is a function of radius, $n = n(r)$, there are two constants of motion. The ray paths lie in planes through the center ($r = 0$) and have constant angular motion about an axis through the center that is perpendicular to this plane, so $\mathbf{x} \times \mathbf{p}$ is constant. If the plane is in the x - y plane, then $n(\alpha y - \beta x)$ is constant. This is analogous to motion of a particle in a central force field. Two of the best-known examples are the Maxwell fisheye^{48,74} for which $n(r) \propto (1 + r^2)^{-1}$, and the Luneburg lens,^{45,75} for which $n(r) = \sqrt{2 - r^2}$ for $r \leq 1$ and $n = 1$ for $r > 1$.

If n does not vary with z , then $\mathcal{H} = n\gamma$ is constant for a ray as a function of z , according to Eq. (66).

If the medium is layered, so the index varies in only the z direction, then $n\alpha$ and $n\beta$ are constant. If θ is the angle relative to the z axis, then $n(z)\sin\theta(z)$ is constant, giving Snell's law as a special case.

The homogeneous medium, where $\partial n/\partial x = \partial n/\partial y = \partial n/\partial z = 0$, is a special case in which there are three constants of motion, $n\alpha$, $n\beta$, and $n\gamma$, so rays travel in straight lines.

1.6 CONSERVATION OF ÉTENDUE

If a bundle of rays intersects a constant z plane in a small region of size $dx dy$ and has a small range of angles $d\alpha d\beta$, then as the light propagates through a lossless system, the following quantity remains constant:

$$n^2 dx dy d\alpha d\beta = n^2 dA d\alpha d\beta = n^2 dA \cos\theta d\omega = dx dy dp_x dp_y, \quad (69)$$

Here $dA = dx dy$ is the differential area, $d\omega$ is the solid angle, and θ is measured relative to the normal to the plane. The integral of this quantity

$$\int n^2 dx dy d\alpha d\beta = \int n^2 dA d\alpha d\beta = \int n^2 dA \cos\theta d\omega = \int dx dy dp_x dp_y, \quad (70)$$

is the *étendue*, and is also conserved. For lambertian radiation of radiance L_e , the total power transferred is $P = \int L_e n^2 d\alpha d\beta dx dy$. The *étendue* and related quantities are known by a variety of names,⁷⁶ including *generalized Lagrange invariant*, *luminosity*, *light-gathering power*, *light grasp*, *throughput*, *acceptance*, *optical extent*, and *area-solid-angle-product*. The angle term is not actually a solid angle, but is weighted. It does approach a solid angle in the limit of small extent. In addition, the integrations can be over area, giving $n^2 d\alpha d\beta \int dA$, or over angle, giving $n^2 dA \int d\alpha d\beta$. A related quantity is the geometrical vector flux,⁷⁷ with components $(\int dp_y dp_z, \int dp_x dp_z, \int dp_x dp_y)$. In some cases these quantities include a brightness factor, and in others they are purely geometrical. The *étendue* is related to the information capacity of a system.⁷⁸

As special case, if the initial and final planes are conjugate with transverse magnification $m = dx'/dx = dy'/dy$, then

$$n^2 d\alpha d\beta = n'^2 m^2 d\alpha' d\beta' \quad (71)$$

Consequently, the angular extents of the entrance and exit pupil in direction cosine coordinates are related by

$$n^2 \int_{\text{entrance pupil}} d\alpha d\beta = n'^2 m^2 \int_{\text{exit pupil}} d\alpha' d\beta' \quad (72)$$

See also the discussion of image irradiance in the section on apertures and pupils.

This conservation law is general; it does not depend on index homogeneity or on axial symmetry. It can be proven in a variety of ways, one of which is with characteristic functions.^{79–81} Phase space arguments involving Liouville's theorem can also be applied.^{82–85} Another type of proof involves thermodynamics, using conservation of radiance (or brightness) or the principle of detailed balance.^{86–89} Conversely, the thermodynamic principle can be proven from the geometrical optics one.^{90–92} In the paraxial limit for systems of revolution the conservation of *étendue* between object and image planes is related to the two-ray paraxial invariant, Eq. (152). Some historical aspects are discussed by Rayleigh.^{93,94}

1.7 SKEW INVARIANT

In a rotationally symmetric system, whose indices may be constant or varying, a *skew ray* is one that does not lie in a plane containing the axis. The *skewness* of such a ray is

$$\mathcal{S} = n(\alpha y - \beta x) = p_x y - p_y x \quad (73)$$

As a skew ray propagates through the system, this quantity, known as the *skew invariant*, does not change.⁹⁵⁻¹⁰⁴ For a meridional ray, one lying in a plane containing the axis, $\mathcal{S} = 0$. The skewness can be written in vector form as

$$\mathcal{S} = \mathbf{a} \cdot (\mathbf{x} \times \mathbf{p}) \quad (74)$$

where \mathbf{a} is a unit vector along the axis, \mathbf{x} is the position on a ray, and \mathbf{p} is the optical cosine vector at that position.

This invariance is analogous to the conservation of the axial component of angular momentum in a cylindrical force field, and it can be proven in several ways. One is by performing the rotation operations on α , β , x , and y (as discussed in the section on heterogeneous media). Another is by means of characteristic functions. It can also be demonstrated that \mathcal{S} is not changed by refraction or reflection by surfaces with radial gradients. The invariance holds also for diffractive optics that are figures of rotation.

A special case of the invariant relates the intersection points of a skew ray with a given meridian. If a ray with directions (α, β) in a space of index n intersects the $x = 0$ meridian with height y , then at another intersection with this meridian in a space with index n' , its height y' , and direction cosine α' are related by

$$n\alpha y = n'\alpha' y' \quad (75)$$

The points where rays intersect the same meridian are known as *diapoints* and the ratio y'/y as the *diamagnification*.⁹⁹

1.8 REFRACTION AND REFLECTION AT INTERFACES BETWEEN HOMOGENEOUS MEDIA

Introduction

The initial ray direction is specified by the unit vector $\mathbf{r} = (\alpha, \beta, \gamma)$. After refraction or reflection the direction is $\mathbf{r}' = (\alpha', \beta', \gamma')$. At the point where the ray intersects the surface, its normal has direction $\mathbf{S} = (L, M, N)$.

The *angle of incidence* I is the angle between a ray and the surface normal at the intersection point. This angle and the corresponding outgoing angle I' are given by

$$|\cos I| = |\mathbf{r} \cdot \mathbf{S}| = |\alpha L + \beta M + \gamma N| \quad (76)$$

$$|\cos I'| = |\mathbf{r}' \cdot \mathbf{S}| = |\alpha' L + \beta' M + \gamma' N|$$

In addition

$$|\sin I| = |\mathbf{r} \times \mathbf{S}| \quad \text{and} \quad |\sin I'| = |\mathbf{r}' \times \mathbf{S}| \quad (77)$$

The signs of these expressions depend on which way the surface normal vector is directed. The surface normal and the ray direction define the *plane of incidence*, which is perpendicular to the vector cross product $\mathbf{S} \times \mathbf{r} = (M\gamma - N\beta, N\alpha - L\gamma, L\beta - M\alpha)$. After refraction or reflection, the outgoing ray is in the same plane. This symmetry is related to the fact that optical path length is an extremum.

The laws of reflection and refraction can be derived from Fermat's principle, as is done in many books. At a planar interface, the reflection and refraction directions are derived from Maxwell's equations using the boundary conditions. For scalar waves at a plane interface, the directions are related to the fact that the number of oscillation cycles is the same for incident and outgoing waves.

Refraction

At an interface between two homogeneous and isotropic media, described by indices n and n' , the incidence angle I and the outgoing angle I' are related by *Snell's law*:¹⁰⁵

$$n' \sin I' = n \sin I \quad (78)$$

If $\sin I' > 1$, there is total internal reflection. Another relationship is

$$n' \cos I' = \sqrt{n'^2 - n^2 \sin^2 I} = \sqrt{n'^2 - n^2 + n^2 \cos^2 I} \quad (79)$$

Snell's law can be expressed in a number of ways, one of which is

$$n[\mathbf{r} + (\mathbf{r} \cdot \mathbf{S})\mathbf{S}] = n'[\mathbf{r}' + (\mathbf{r}' \cdot \mathbf{S})\mathbf{S}] \quad (80)$$

Taking the cross product of both sides with \mathbf{S} gives another form

$$n'(\mathbf{S} \times \mathbf{r}') = n(\mathbf{S} \times \mathbf{r}) \quad (81)$$

A quantity that appears frequently in geometrical optics (but which has no common name or symbol) is

$$\Gamma = n' \cos I' - n \cos I \quad (82)$$

It can be written in several ways

$$\Gamma = (n\mathbf{r} - n'\mathbf{r}') \cdot \mathbf{S} = -n \cos I + \sqrt{n'^2 - n^2 \sin^2 I} = -n \cos I + \sqrt{n'^2 - n^2 + n^2 \cos^2 I} \quad (83)$$

In terms of Γ , Snell's law is

$$n'\mathbf{r}' = n\mathbf{r} + \Gamma\mathbf{S} \quad (84)$$

or

$$n'\alpha' = n\alpha + L\Gamma, \quad n'\beta' = n\beta + M\Gamma, \quad n'\gamma' = n\gamma + N\Gamma \quad (85)$$

The outgoing direction is expressed explicitly as a function of incident direction by

$$n'\mathbf{r}' = n\mathbf{r} + \mathbf{S}[n\mathbf{r} \cdot \mathbf{S} - \sqrt{n'^2 - n^2 + (n\mathbf{r} \cdot \mathbf{S})^2}] \quad (86)$$

If the surface normal is in the z direction, the equations simplify to

$$n'\alpha' = n\alpha, \quad n'\beta' = n\beta, \quad n'\gamma' = \sqrt{n'^2 - n^2 + n^2\gamma^2} \quad (87)$$

If $\beta = 0$, this reduces to $n'\alpha' = n\alpha$, the familiar form of Snell's law, written with direction cosines, with $n'\gamma' = (n'^2 - n^2\alpha^2)^{1/2}$, corresponding to Eq. (79). Another relation from Eq. (85) is

$$\frac{n'\alpha' - n\alpha}{L} = \frac{n'\beta' - n\beta}{M} = \frac{n'\gamma' - n\gamma}{N} = \Gamma \quad (88)$$

All of the above expressions can be formally simplified by using $\mathbf{p} = n\mathbf{r}$ and $\mathbf{p}' = n'\mathbf{r}'$. For a succession of refractions by parallel surfaces,

$$n_1 \sin I_1 = n_2 \sin I_2 = n_3 \sin I_3 = \dots \quad (89)$$

so the angles within any two media are related by their indices alone, regardless of the intervening layers. Refractive indices vary with wavelength, so the angles of refraction do likewise.

Reflection

The reflection equations can be derived from those for refraction by setting the index of the final medium equal to the negative of that of the incident medium, i.e., $n' = -n$, which gives $\Gamma = -2n \cos I$. The angles of incidence and reflection are equal

$$I' = I \quad (90)$$

The incident and reflected ray directions are related by

$$\mathbf{S} \times \mathbf{r}' = -\mathbf{S} \times \mathbf{r} \quad (91)$$

Another expression is

$$\mathbf{r}' = \mathbf{r} - (2\mathbf{S} \cdot \mathbf{r})\mathbf{S} = \mathbf{r} - (2 \cos I)\mathbf{S} \quad (92)$$

The components are

$$\alpha' = \alpha - 2L \cos I \quad \beta' = \beta - 2M \cos I \quad \gamma' = \gamma - 2N \cos I \quad (93)$$

This relationship can be written in terms of dyadics¹⁰⁶ as $\mathbf{r}' = (\mathbf{I} - \mathbf{SS}) \cdot \mathbf{r}$. This is equivalent to the matrix form¹⁰⁷⁻¹¹¹

$$\begin{pmatrix} \alpha' \\ \beta' \\ \gamma' \end{pmatrix} = \begin{pmatrix} 1-2L^2 & -2LM & -2LN \\ -2LM & 1-2M^2 & -2MN \\ -2LN & -2MN & 1-2N^2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \\ \gamma \end{pmatrix} \quad (94)$$

Each column of this matrix is a set of direction cosines and is orthogonal to the others, and likewise for the rows. The matrix associated with a sequence of reflections from plane surfaces is calculated by multiplying the matrices for each reflection. Another relationship is

$$\frac{\alpha' - \alpha}{L} = \frac{\beta' - \beta}{M} = \frac{\gamma' - \gamma}{N} \quad (95)$$

If the surface normal is in the z direction, so $(L, M, N) = (0, 0, 1)$, then

$$\alpha' = \alpha \quad \beta' = \beta \quad \gamma' = -\gamma \quad (96)$$

Reflection by a Plane Mirror: Positions of Image Points

If light from a point (x, y, z) is reflected by a plane mirror whose distance from the coordinate origin is d , and whose surface normal has direction (L, M, N) , the image point coordinates (x', y', z') are given by

$$\begin{pmatrix} x' \\ y' \\ z' \\ 1 \end{pmatrix} = \begin{pmatrix} 1-2L^2 & -2LM & -2LN & 2dL \\ -2LM & 1-2M^2 & -2MN & 2dM \\ -2LN & -2MN & 1-2N^2 & 2dN \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix} \quad (97)$$

This transformation involves both rotation and translation, with only the former effect applying if $d = 0$. It is an affine type of collinear transformation, discussed in the section on collineation. The effect of a series of reflections by plane mirrors is found by a product of such matrices. The transformation can also be formulated in terms of quaternions.^{110,112}

Diffractive Elements

The changes in directions produced by gratings or diffractive elements can be handled in an ad hoc geometrical way^{113,114}

$$n' \mathbf{r}'_G \times \mathbf{S} = n \mathbf{r} \times \mathbf{S} + G \frac{\lambda}{p} \mathbf{q} \quad (98)$$

Here λ is the vacuum wavelength, p is the grating period, \mathbf{q} is a unit vector tangent to the surface and parallel to the rulings, and G is the diffraction order. Equations (81) and (91) are special cases of this equation for the 0th order.

1.9 IMAGING

Introduction

Image formation is the principal use of lenses. Moreover, lenses form images even if this is not their intended purpose. This section provides definitions, and discusses basic concepts and limitations. The purposes of the geometrical analysis of imaging include the following: (1) discovering the nominal relationship between an object and its image, principally the size, shape, and location of the image, which is usually done with paraxial optics; (2) determining the deviations from the nominal image, i.e., the aberrations; (3) estimating image irradiance; (4) understanding fundamental limitations—what is inherently possible and impossible; and (5) supplying information for diffraction calculations, usually the optical path lengths.

Images and Types of Images

A definition of image (Webster 1934¹¹⁵) is: “The optical counterpart of an object produced by a lens, mirror, or other optical system. It is a geometrical system made up of foci corresponding to the parts of the object.” The point-by-point correspondence is the key, since a given object can have a variety of different images.

Image irradiance can be found only approximately from geometrical optics, the degree of accuracy of the predictions varying from case to case. In many instances wave optics is required, and for objects that are not self-luminous, an analysis involving partial coherence is also needed.

The term *image* is used in a variety of ways, so clarification is useful. The light from an object produces a three-dimensional distribution in image space. The *aerial image* is the distribution on a mathematical surface, often that of best focus, the locus of points of the images of object points. An aerial image is never the final goal; ultimately, the light is to be captured. The *receiving surface* (NS) is that on which the light falls, the distribution of which there can be called the *received image* (NS). This distinction is important in considerations of defocus, which is a relationship, not an absolute. The record thereby produced is the *recorded image* (NS). The recorded image varies with the position of the receiving surface, which is usually intended to correspond with the aerial image surface. In this section, “image” means aerial image, unless otherwise stated.

Object Space and Image Space

The object is said to exist in *object space*; the image, in *image space*. Each space is infinite, with a physically accessible region called *real*, and an inaccessible region, referred to as *virtual*. The two spaces may overlap physically, as with reflective systems. Corresponding quantities and locations associated with the object and image spaces are typically denoted by the same symbol, with a prime indicating image space. Positions are specified by a coordinate system (x, y, z) in object space and (x', y', z') in image space. The refractive indices of the object and image spaces are n and n' .

Image of a Point

An *object point* is thought of as emitting rays in all directions, some of which are captured by the lens, whose internal action converges the rays, more or less, to an *image point*, the term “point” being used even if the ray convergence is imperfect. Object and image points are said to be *conjugate*. Since geometrical optics is reversible, if A' is the image of A , then A is the image of A' .

Mapping Object Space to Image Space

If every point were imaged stigmatically, then the entire object space would be mapped into the image space according to a transformation

$$x' = x'(x, y, z) \quad y' = y'(x, y, z) \quad z' = z'(x, y, z) \quad (99)$$

The mapping is reciprocal, so the equations can be inverted. If n and n' are constant, then the mapping is a collinear transformation, discussed below.

Images of Extended Objects

An *extended object* can be thought of as a collection of points, a subset of the entire space, and its stigmatic image is the set of conjugate image points. A surface described by $0 = F(x, y, z)$ has an image surface

$$0 = F'(x', y', z') = F[x(x', y', z'), y(x', y', z'), z(x', y', z')] \quad (100)$$

A curve described parametrically by $\mathbf{x}(\sigma) = [x(\sigma), y(\sigma), z(\sigma)]$ has an image curve

$$\mathbf{x}'(\sigma) = \{x'[x(\sigma), y(\sigma), z(\sigma)], y'[x(\sigma), y(\sigma), z(\sigma)], z'[x(\sigma), y(\sigma), z(\sigma)]\} \quad (101)$$

Rotationally Symmetric Lenses

Rotationally symmetric lenses have an *axis*, which is a ray path (unless there is an obstruction). All planes through the axis, the *meridians* or *meridional planes*, are planes with respect to which there is bilateral symmetry. An *axial object point* is conjugate to an *axial image point*. An axial image point is located with a single ray in addition to the axial one. Off-axis object and image points are in the same meridian, and may be on the same or opposite sides of the axis. The *object height* is the distance of a point from the axis, and the *image height* is that for its image. It is possible to have rotational symmetry without bilateral symmetry, as in a system made of crystalline quartz,¹¹⁶ but such systems are not discussed here. For stigmatic imaging by a lens rotationally symmetric about the z axis, Eq. (99) gives

$$x' = x'(x, z) \quad y' = y'(y, z) \quad z' = z'(z) \quad (102)$$

Planes Perpendicular to the Axis

The arrangement most often of interest is that of planar object and receiving surfaces, both perpendicular to the lens axis. When the terms *object plane* and *image plane* are used here without further elaboration, this is the meaning. (This arrangement is more common for manufactured systems with flat detectors, than for natural systems, for instance, eyes, with their curved retinas.)

Magnifications

The term *magnification* is used in a general way to denote the ratio of conjugate object and image dimensions, for example, heights, distances, areas, volumes, and angles. A single number is inadequate when object and image shapes are not geometrically similar. The term magnification implies an increase, but this is not the case in general.

Transverse Magnification

With object and image planes perpendicular to the axis, the relative scale factor of length is the *transverse magnification* or *lateral magnification*, denoted by m , and usually referred to simply as “the magnification.” The transverse magnification is the ratio of image height to object height, $m = h'/h$. It can also be written in differential form, e.g., $m = dx'/dx$ or $m = \Delta x'/\Delta x$. The transverse magnification has a sign, and it can have any value from $-\infty$ to $+\infty$. Areas in such planes are scaled by m^2 . A lens may contain plane mirrors that affect the image parity or it may be accompanied by external plane mirrors that reorient images and change their parity, but these changes are independent of the magnification at which the lens works.

Longitudinal Magnification

Along the rotational axis, the *longitudinal magnification*, m_L , also called *axial magnification*, is the ratio of image length to object length in the limit of small lengths, i.e., $m_L = dz'/dz$.

Visual Magnification

With visual instruments, the perceived size of the image depends on its angular subtense. *Visual magnification* is the ratio of the angular subtense of an image relative to that of the object viewed directly. Other terms are used for this quantity, including “magnification,” “power,” and “magnifying power.” For objects whose positions can be controlled, there is arbitrariness in the subtense without the instrument, which is greatest when the object is located at the near-point of the observer. This distance varies from person to person, but for purposes of standardization the distance is taken to be 250 mm. For instruments that view distant objects there is no arbitrariness of subtense with direct viewing.

Ideal Imaging and Disappointments in Imaging

Terms such as *perfect imaging* and *ideal imaging* are used in various ways. The ideal varies with the circumstances, and there are applications in which imaging is not the purpose, for instance, energy collection and Fourier transformation. The term *desired imaging* might be more appropriate in cases where that which is desired is fundamentally impossible. Some deviations from what is desired are called *aberrations*, whether their avoidance is possible or not. Any ideal that can be approximated must agree in its paraxial limit ideal with what a lens actually does in its paraxial limit.

Maxwellian Ideal for Single-Plane Imaging

The most common meaning of perfect imaging is that elucidated by Maxwell,¹¹⁷ and referred to here as *maxwellian ideal* or *maxwellian perfection*. This ideal is fundamentally possible. The three conditions for such imaging of points in a plane perpendicular to the lens axis are: (1) Each point is imaged stigmatically. (2) The images of all the points in a plane lie on a plane that is likewise perpendicular to the axis, so the field is flat, or free from field curvature. (3) The ratio of image heights to object heights is the same for all points in the plane. That is, transverse magnification is constant, or there is no distortion.

The Volume Imaging Ideal

A more demanding ideal is that points everywhere in regions of constant index be imaged stigmatically and that the imaging of every plane be flat and free from distortion. For planes perpendicular to the lens axis, such imaging is described mathematically by the collinear transformation, discussed later. It is inherently impossible for a lens to function in this way, but the mathematical apparatus of collineation is useful in obtaining approximate results.

Paraxial, First-Order, and Gaussian Optics

The terms “paraxial,” “first order,” and “gaussian” are often used interchangeably, and their consideration is merged with that of collineation. The distinction is often not made, probably because these descriptions agree in result, although differing in approach. One of the few discussions is that of Southall.¹¹⁸ A *paraxial* analysis has to do with the limiting case in which the distances of rays from the axis approach zero, as do the angles of the rays relative to the axis. The term *first order* refers to the associated mathematical approximation in which the positions and directions of such rays are computed with terms to the first order only in height and angle. *Gaussian* refers to certain results of the paraxial optics, where lenses are black boxes whose properties are summarized by the existence and locations of cardinal points. In the limit of small heights and angles, the equations of collineation are identical to those of paraxial optics. Each of these is discussed in greater detail below.

Fundamental Limitations

There are fundamental geometrical limitations on optical systems, relating to the fact that a given ray passes through many points and a given point lies on many rays. So the images of points on the same line or plane, or on different planes, are not independent. A set of rays intersecting at several points in object space cannot be made to satisfy arbitrary requirements in image space. Such limitations are best studied by the methods of hamiltonian optics.

Stigmatic Imaging

If all rays from an object point converge precisely to an image point, the imaging of this point is said to be *stigmatic*. The optical path lengths of all rays between two such points are identical. A stigmatic image point is located by the intersection of any two rays that pass through the points. An *absolute instrument* is one which images all points stigmatically.¹¹⁹ For such imaging

$$n\delta x = n'\delta x', \quad n\delta y = n'\delta y', \quad n\delta z = n'\delta z' \quad (103)$$

where conjugate length elements are δx and $\delta x'$, δy and $\delta y'$, δz and $\delta z'$.

Path Lengths and Conjugate Points

All the rays from an object point to its stigmatic image point have the same optical path length. For focal lenses, the paths lengths for different pairs of conjugate points in a plane perpendicular to the axis are different, except for points on circles centered on the axis. For afocal lenses path lengths are nominally the same for all points in planes perpendicular to the axis. For afocal lenses with transverse magnification $\pm n/n'$, path lengths can be the same for all points. In general, the path lengths between different points on an object and image surface are equal only if the shape of the image surface is that of a wavefront that has propagated from a wavefront whose shape is that of the object surface.

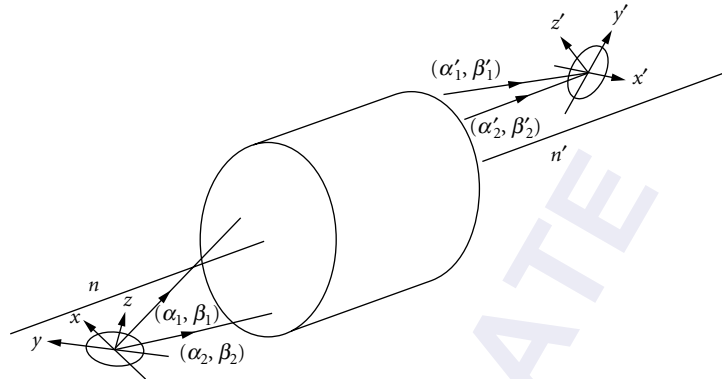


FIGURE 2 The cosine condition. A small area in object space about the origin in the x - y plane is imaged to the region around the origin of the x' - y' plane in image space. A pair of rays from the origin with direction cosines (α_1, β_1) and (α_2, β_2) arrive with (α'_1, β'_1) and (α'_2, β'_2) . The direction cosines and the transverse magnification in the planes are related by Eq. (104).

The Cosine Condition

The *cosine condition* relates object space and image space ray angles, if the imaging is stigmatic over some area.^{116,120,121} Let the x - y plane lie in the object surface and the x' - y' plane in the conjugate surface (Fig. 2). Two rays leaving a point in the object region have direction cosines (α_1, β_1) and (α_2, β_2) , and the rays on the image side have (α'_1, β'_1) and (α'_2, β'_2) . If the imaging is stigmatic, with local transverse magnification m on the surface, then

$$m = \frac{n(\alpha_1 - \alpha_2)}{n'(\alpha'_1 - \alpha'_2)} = \frac{n(\beta_1 - \beta_2)}{n'(\beta'_1 - \beta'_2)} \quad (104)$$

In the limit as $\alpha_1 \rightarrow \alpha_2$ and $\beta_1 \rightarrow \beta_2$, the cosine condition gives

$$m = \frac{n d\alpha}{n' d\alpha'} = \frac{nd\beta}{n' d\beta'} \quad (105)$$

This condition also holds in the more general case of *isoplanatic* imaging, where there is aberration that is locally constant across the region in question.^{122,123}

The Abbe Sine Condition

The *sine condition* or *Abbe sine condition*^{119,124} is a special case of the cosine condition for object and image planes perpendicular to the axis in regions about the axis. For a plane with transverse magnification m , let θ be the angle relative to the axis made by a ray from an axial object point, and θ' be that in image space. If the lens is free of coma

$$m = \frac{n \sin \theta}{n' \sin \theta'} = \frac{n\alpha}{n'\alpha'} = \frac{n\beta}{n'\beta'} \quad (106)$$

for all θ and θ' . There are signs associated with θ and θ' , so that $m > 0$ if they have the same sign, and $m < 0$ if the signs are opposite. This equation is sometimes written with m replaced by the ratio of paraxial angles. There is sometimes the implication that θ and θ' refer only to the extreme angles passing through the lens, when in fact the sine condition dictates that the ratio of the sines is the constant for all angles. For an object at infinity, the sine condition is

$$\sin\theta' = -\frac{y}{f'} \quad \text{or} \quad n'\beta' = -y\phi \quad (107)$$

where y is the height of a ray parallel to the axis, ϕ is the power of the lens, and f' is the rear focal length. These relationships hold to a good approximation in most lenses, since small deviations are associated with large aberrations. A deviation from this relationship is called *offense against the sine condition*, and is associated with coma.^{123,125–128} The sine condition does not apply where there are discontinuities in ray behavior, for example, in devices like Fresnel lenses, or to diffraction-based devices like zone plates.

The Herschel Condition

The *Herschel condition* is a relationship that holds if the imaging is stigmatic for nearby points along the axis.^{119,129,130} The two equivalent relations are

$$m = \frac{n \sin(\frac{1}{2}\theta)}{n' \sin(\frac{1}{2}\theta')} \quad \text{and} \quad m_L = \frac{n \sin^2(\frac{1}{2}\theta)}{n' \sin^2(\frac{1}{2}\theta')} = \frac{n(1-\gamma)}{n'(1-\gamma')} \quad (108)$$

The Herschel condition is inconsistent with the sine condition unless $m \pm n/n'$. So, in general, stigmatic imaging in one plane precludes that in others.

Sine and Herschel Conditions for Afocal Systems

For afocal systems the sine condition and Herschel condition are identical. For rays entering parallel to the axis at y and leaving at y' , they are

$$m = \frac{y'}{y} \quad (109)$$

That is, the ratio of incoming and outgoing heights is independent of the incoming height. (Ref. 128, chap. 3, "The Sine Condition and Herschel's Condition").

Stigmatic Imaging Possibilities

For object and image spaces with constant refractive indices, stigmatic imaging is only possible for the entire spaces for afocal lenses with identical transverse and longitudinal magnifications $m = \pm n/n'$ and $|m_L| = |m|$. Such lenses re-create not only the intersection points, but the wavefronts, since the corresponding optical path lengths are the same in both spaces, Eq. (103). For other lenses with constant object and image space indices, the maxwellian ideal can be met for only a single surface. In addition, a single point elsewhere can be imaged stigmatically.^{127,131} Nonplanar surfaces can be imaged stigmatically, a well-known example being the imaging of spherical surfaces by a spherical refracting surface, for a particular magnification.¹¹⁹ For systems with spherical symmetry, it is possible that two nonplanar surfaces be stigmatically imaged.¹³² In addition, various systems with heterogeneous indices can image stigmatically over a volume.

1.10 DESCRIPTION OF SYSTEMS OF REVOLUTION

Introduction

This section is concerned with the optical description of lens and mirror systems that are figures of revolution.^{133–145} From a mechanical viewpoint, optical systems are comprised of lenses and mirrors. From the point of view of the light, the system is regions of media with different indices, separated by interfaces of various shapes. This section is limited to homogeneous isotropic media. It is further restricted to reflecting and refracting surfaces that are nominally smooth, and to surfaces that are figures of revolution arranged so their axes are collinear, so the entire system is a figure of revolution about the *lens axis*. (The often-used term “optical axis” is also used in crystallography. Moreover, the axis is often mechanical as well as “optical.”) The lens axis is the z axis of an orthogonal coordinate system, with the x - y plane perpendicular. The distance from a point to the axis is $\rho = \sqrt{x^2 + y^2}$. Along the axis, the positive direction is from left to right.

Terminology

A *meridian* or *meridional plane* contains the axis, all such planes being equivalent. Meridional planes are planes of bilateral symmetry if the indices are homogeneous and isotropic. Some optical systems are comprised of pieces of surfaces of revolution, in which case it is still useful to discuss the behavior about the axis.

Reflection, Unfolded Diagrams

Light passes through refractive systems more or less in the same direction relative to the axis. In reflective and catadioptric systems, the light may change directions. (It may not, in the case of grazing incidence optics.) In order to consider all types of systems in a unified way, pictorially and algebraically, reflections can often be “unfolded,” i.e., represented pictorially as transmission, with mirrors replaced by hypothetical lenses with the same properties, Figs. 3 and 18. Some structures must be taken into account several times in unfolding. For example, a hole may block light at one point along a ray and transmit it at another. (In some considerations, unfolding can be misleading—for instance, those involving stray light.)

Description of Surfaces

A *surface* is an interface between media with different refractive indices—either refracting or reflecting. The surface is the optical element, produced by a lens, which is a mechanical element. Surfaces can be described mathematically in many ways. (For example, conics can be described as loci of points with certain relationships.) In optical instruments, the entire surface is rarely used, and the axial region is usually special, so the description usually begins there and works out. The *vertex* of a figure of revolution intersects with the axis, and is a local extremum. The plane perpendicular to the axis and tangent to the vertex will be referred to as the *vertex plane* (NS). A surface can be described by its *sag*, the directed distance $z(\rho)$ from the vertex plane to the surface, Fig. 4. The vertex is usually taken to have $z(0) = 0$. The *vertex curvature* or *paraxial curvature* c and radius of curvature r are given by

$$c = \frac{1}{r} = \left. \frac{\partial^2 z(\rho)}{\partial \rho^2} \right|_{\rho=0} \quad (110)$$

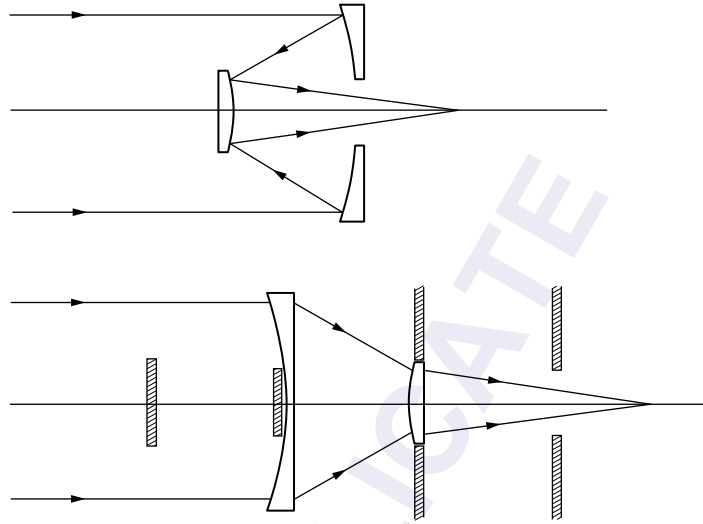


FIGURE 3 Example of an unfolded diagram. The two-mirror system above has an unfolded representation below. The reflective surfaces are replaced by thin lens equivalents. Their obstructions and the finite openings are accounted for by dummy elements.

For an arbitrary surface, this curvature is identical to that of the sphere which is a best fit on axis. The sign convention for curvature and radius is that c and r are positive if the center of curvature is to the right of the vertex, as in the case shown in Fig. 4. In general, the curvature is mathematically more foolproof than radius, since curvature can be zero, but it is never infinite, whereas radius is never zero, but may be infinite.

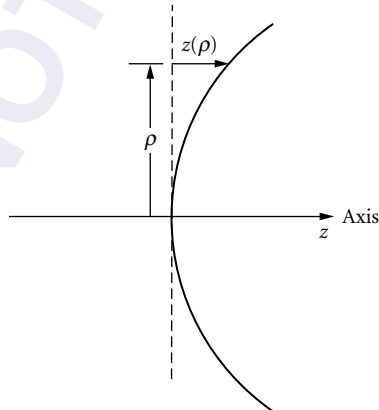


FIGURE 4 Description of a surface of revolution. The distance from the axis is ρ , and the sag $z(\rho)$ is the distance from the vertex tangent plane to the surface.

Spherical Surfaces

The spherical surface is the most common in optics, since it is most naturally produced. Spherical is the default, and is assumed when no other mention is made. Aspheres are thought of as deviating from spheres, rather than spheres as special cases of more general forms. The equation for a sphere with radius r , curvature c , and a vertex at $z = 0$ is

$$\rho^2 + (z-r)^2 = r^2 \quad (111)$$

The sag is given by

$$z(\rho) = r - \sqrt{r^2 - \rho^2} = r(1 - \sqrt{1 - c^2 \rho^2}) = \frac{c\rho^2}{1 + \sqrt{1 - c^2 \rho^2}} \quad (112)$$

The Taylor expansion is

$$z(\rho) = \frac{1}{2}c\rho^2 + \frac{1}{8}c^3\rho^4 + \frac{1}{16}c^5\rho^6 + \frac{5}{128}c^7\rho^8 + \frac{7}{256}c^9\rho^{10} + \dots \quad (113)$$

At the point (x, y, z) on the surface of the sphere, the surface normal has direction cosines

$$(L, M, N) = \left(\frac{x}{r}, \frac{y}{r}, \frac{z-r}{r} \right) = (cx, cy, cz-1) \quad (114)$$

Conics of Rotation

The general form of a conic of rotation about the z axis is

$$z(\rho) = \frac{r}{\varepsilon} (1 - \sqrt{1 - \varepsilon c^2 \rho^2}) = \frac{c\rho^2}{1 + \sqrt{1 - \varepsilon c^2 \rho^2}} \quad (115)$$

The value of ε determines the type of conic, as given in the table below. It is common in optics to use κ , the *conic parameter* or *conic constant*, related by

$$\kappa = \varepsilon - 1 \quad \text{or} \quad \varepsilon = 1 + \kappa \quad (116)$$

Another parameter used to describe conics is the *eccentricity* e , used in the polar coordinate form for conics about their focus: $r(\theta) = a / (1 + e \cos \theta)$ where $e^2 = -\kappa$. In the case of paraboloids, the first form of Eq. (115) breaks down. A cone can be approximated by a hyperbola with $\kappa = -\sec^2 \theta$, where θ is the cone half angle.

Conic Type and Value of Parameter

Parameter	ε	κ	e
Oblate ellipsoid	$\varepsilon > 1$	$\kappa > 0$	—
Sphere	$\varepsilon = 1$	$\kappa = 0$	0
Prolate ellipsoid	$0 < \varepsilon < 1$	$-1 < \kappa < 0$	$0 < e < 1$
Paraboloid	$\varepsilon = 0$	$\kappa = -1$	$e = 1$
Hyperboloid	$\varepsilon < 0$	$\kappa < -1$	$e > 1$

The Taylor expansion for the sag is

$$z(\rho) = \frac{1}{2}c\rho^2 + \frac{1}{8}\varepsilon c^3\rho^4 + \frac{1}{16}\varepsilon^2 c^5\rho^6 + \frac{5}{128}\varepsilon^3 c^7\rho^8 + \frac{7}{256}\varepsilon^4 c^9\rho^{10} + \dots \quad (117)$$

The surface normals are

$$(L, M, N) = [1 - 2c(\varepsilon - 1)z + c^2\varepsilon(\varepsilon - 1)z^2]^{-1/2} (cx, cy, cz - 1) \quad (118)$$

The sagittal and tangential curvatures are

$$c_s = \frac{c}{[1+(1-\varepsilon)c^2\rho^2]^{1/2}}, \quad c_t = \frac{c}{[1+(1-\varepsilon)c^2\rho^2]^{3/2}} \quad (119)$$

General Asphere of Revolution

For an arbitrary figure of revolution all of whose derivatives are continuous, the Taylor expansion is

$$z(\rho) = \frac{1}{2}c\rho^2 + q_4\rho^4 + q_6\rho^6 + \dots \quad (120)$$

An asphere is often treated as a sphere that matches at the vertex and a deviation therefrom:

$$z(\rho) = z_{\text{sphere}}(\rho) + a_4\rho^4 + a_6\rho^6 + \dots \quad (121)$$

Alternatively, nonconic aspheres can be treated as conics and a deviation therefrom:

$$z(\rho) = z_{\text{conic}}(\rho) + b_4\rho^4 + b_6\rho^6 + \dots \quad (122)$$

The expansion coefficients are different in each case. Additional information on the coefficients is given by Malacara¹⁴⁴ and Brueggemann.¹³⁵ The sagittal and tangential curvatures are given in general by

$$c_s = \frac{\dot{z}(\rho)}{\rho[1+\dot{z}(\rho)^2]^{1/2}}, \quad c_t = \frac{\ddot{z}(\rho)}{[1+\dot{z}(\rho)^2]^{3/2}} \quad (123)$$

Here $\dot{z}(\rho) = dz(\rho)/d\rho$ and $\ddot{z}(\rho) = d^2z(\rho)/d\rho^2$.

1.11 TRACING RAYS IN CENTERED SYSTEMS OF SPHERICAL SURFACES

Introduction

Ray tracing is the process of calculating the paths of rays through optical systems. Two operations are involved, propagation from one surface to the next and refraction or reflection at the surfaces. Exact equations can be written for spherical surfaces and conics of revolution with homogeneous media.¹⁴⁶⁻¹⁵³ Conics are discussed by Welford.¹⁵² For general aspheres, the intersection position is found by iterating.^{153,154} Nonsymmetric systems are discussed by Welford.¹⁵²

Description and Classification of Rays in a Lens

For optical systems with rotational symmetry, rays are typically described in terms of the axial parameter z . A ray crosses each constant z plane at a point (x, y) with direction cosines (α, β, γ) , where γ is not independent. Thus a ray is described by $[x(z), y(z)]$ and $[\alpha(z), \beta(z)]$.

For systems that are figures of revolution, *meridional rays* are those lying in a meridional plane, a plane that includes the axis, and other rays are *skew rays*. The *axial ray* corresponds to the axis of revolution. Rays are also classified according to their proximity to the axis. *Paraxial rays* are those in the limit of having small angles and small distances from the axis. Nonparaxial rays are sometimes referred to as *finite rays* or *real rays*. *Ray fans* are groups of rays lying in a plane. A *tangential fan* lies in a meridian, and intersects at a *tangential focus*. A *sagittal fan* lies in a plane perpendicular to a meridian, and intersects at a *sagittal focus*.

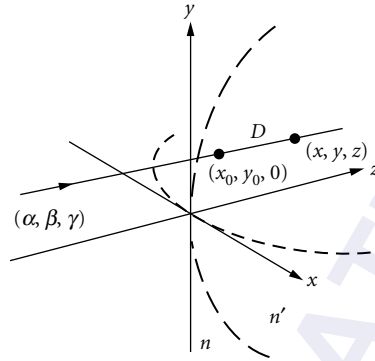


FIGURE 5 Intersection points. A ray with direction cosines (α, β, γ) intersects the vertex tangent plane at $(x_0, y_0, 0)$ and the optical surface at (x, y, z) . The distance between these points is D , given by Eq. (126).

Transfer

In propagating through a homogeneous medium, a ray originating at (x_1, y_1, z_1) with directions (α, β, γ) intersects a z_2 plane at

$$x_2 = x_1 + \frac{\alpha}{\gamma}(z_2 - z_1) \quad \text{and} \quad y_2 = y_1 + \frac{\beta}{\gamma}(z_2 - z_1) \quad (124)$$

Intersection Points

Let the intersection points of the ray with the vertex plane at $z = 0$ be $(x_0, y_0, 0)$, Fig. 5. Define auxiliary functions

$$A(x_0, y_0; \alpha, \beta; c) = \gamma - c(\alpha x_0 + \beta y_0) \quad \text{and} \quad B(x_0, y_0, c) = c^2(x_0^2 + y_0^2) \quad (125)$$

The distance D along the ray from this point to the surface is given by

$$cD = A - \sqrt{A^2 - B} = \frac{B}{A + \sqrt{A^2 - B}} \quad (126)$$

The intersection point has the coordinates

$$x = x_0 + \alpha D, \quad y = y_0 + \beta D, \quad z = \gamma D \quad (127)$$

The incidence angle I at the intersection point is given by

$$\cos I = \sqrt{A^2 - B} \quad (128)$$

so

$$\Gamma = -n\sqrt{A^2 - B} + \sqrt{n'^2 - n^2 + n^2(A^2 - B)} \quad (129)$$

Mathematically, double intersection points are possible, so they must be checked for. If the ray misses the surface, then $A^2 < B$. If there is total internal reflection, the second square root in Eq. (129) is imaginary.

Refraction and Reflection by Spherical Surfaces

Rays refract or reflect at surfaces with reference to the local normal at the intersection point. The surface normal given by Eq. (114) is substituted in the general form for refraction, Eq. (85), to give

$$n'\alpha' = n\alpha - \Gamma cx, \quad n'\beta' = n\beta - \Gamma cy, \quad n'\gamma' = n\gamma - \Gamma(1 - cz) \quad (130)$$

For reflection, the above equations are used, with $n' = -n$, so $\Gamma = 2n \cos I = 2n\sqrt{A^2 - B}$.

Meridional Rays

The meridian is customarily taken to be that for which $x = 0$, so the direction cosines are $(0, \beta, \gamma)$. Let U be the angle between the axis and the ray, so $\beta = \sin U$ and $\gamma = \cos U$. The transfer equation, Eq. (124), becomes

$$y_2 = y_1 + \tan U(z_2 - z_1) \quad (131)$$

The second equation of Eq. (130) can be written

$$n' \sin U' - n \sin U = -cy(n' \cos I' - n \cos I) \quad (132)$$

If the directed distance from the vertex to the intersection point of the incident ray with the axis is l , the outgoing angle is

$$U' = U + \arcsin[(l-1)\sin U] - \arcsin\left[\frac{n}{n'}(l-1)\sin U\right] \quad (133)$$

The directed distance l' from the vertex to the axial intersection of the refracted ray is given by

$$cl' = 1 + (cl-1)\frac{n \sin U}{n' \sin U'} \quad (134)$$

For reflection, setting $n' = -n$ gives

$$U' = U + 2\arcsin[(l-1)\sin U] \quad (135)$$

1.12 PARAXIAL OPTICS OF SYSTEMS OF REVOLUTION

Introduction

The term *paraxial* is used in different ways. In one, paraxial rays are those whose distances from the axis and whose angles relative to the axis are small. This leaves questions of how small is small enough and how this varies from system to system. The other interpretation of the term, which is used here, is that paraxial rays represent a limiting case in which the distances from the axis and angles relative to the axis vanish. *Paraxial optics* then describes the behavior of systems in this limit. The ray-tracing equations in the paraxial limit are linear in angle and in distance from the axis, hence the term *first-order optics*, which is often considered equivalent to paraxial. (There are no 0th-order terms since the expansion is taken about the axis, so a ray with an initial height and angle of zero, i.e., a ray along the axis, has the same outgoing height and angle.) The linearity of the paraxial equations makes them simple and understandable, as well as expressible in matrix form. Paraxial ray tracing is discussed to some extent by almost every book that treats geometrical optics.

Paraxial ray tracing is done to determine the gaussian properties of lenses, to locate image positions and magnifications, and to locate pupils and determine their sizes. Another use of paraxial ray tracing, not discussed here, is the computation of third-order aberrations.¹⁵⁵

Paraxial imaging is perfect in the sense that it agrees with the Maxwell ideal and with that of collimation. Point images everywhere are stigmatic, fields are flat, and there is no distortion. Aberration is often thought of as the difference between the behavior of finite rays and that of paraxial rays. If this approach is taken, then in the process of lens design, finite rays are made to agree, insofar as possible, with the paraxial ones, which cannot be similarly changed. In the paraxial limit, surfaces are described by their vertex curvatures, so conics, aspheres, and spheres are indistinguishable, the difference being in the fourth power and above. Consequently, aberrations can be altered by changing the surface asphericity without changing paraxial properties. A paraxial treatment can be done even if a system is missing the axial region, as in the case with central obstructions and off-axis sections of figures of revolution.

This section is concerned with systems of mirrors and lenses with rotational symmetry and homogeneous refractive indices. In this case, it suffices to work in a single meridian. Generalizations are found in the sections in this chapter on images about known rays and rays in heterogeneous media. Other generalizations involve expansions about given rays in systems that are not rotationally symmetric.

The Paraxial Limit

The lens axis is the z axis, and rays in the $x = 0$ meridian are considered. Ray heights are y , and angles relative to the axis are u . In the paraxial limit, the quantities u , $\tan u$, and $\sin u = \beta$ are indistinguishable. The z -direction cosine is $\gamma = \cos u \approx 1$. Since the ray angles and heights are small, incidence angles are likewise, so $i \approx \sin i$, $\cos I \approx 1$, $\cos I' \approx 1$, and $\Gamma = n' \cos I' - n \cos I \approx n' - n$.

Transfer

In traversing a distance t between two planes, the height of a meridional ray changes from y to y' according to Eq. (124), $y' = y + t\beta/\gamma$. In the paraxial limit, this equation becomes

$$y' = y + tu \quad (136)$$

If a ray travels from one curved surface to the next, the distance t equals the vertex separation to first order, since the correction for the surface sag is of second order in height and angle. This term is given above in Eq. (127).

Refraction

The paraxial form of Snell's law, Eq. (78), is

$$n' i' = n i \quad (137)$$

Reflection

The law of reflection is the same for paraxial as for finite rays,

$$i' = -i \quad (138)$$

Angle of Incidence at a Surface

A ray with an angle u , which intersects a surface of curvature c at height y , makes an angle i with the local surface normal of the surface given by

$$i = u + yc \quad (139)$$

This equation is easily remembered from two special cases. When $y = 0$, the intersection is at the vertex, so $i = u$. When $u = -cy$, the ray is directed through the center of curvature, so $i = 0$.

Refraction at a Surface

The above equation combined with that for Snell's law gives

$$n' u' = nu - yc(n' - n) \quad (140)$$

This equation can also be obtained from the exact equation, $n' \beta' = n \beta - \Gamma cy$, Eq. (125). In the paraxial limit, $\Gamma = n' - n$, and the intersection height y is that in the vertex plane.

Reflection at a Surface

The relationship between incident and outgoing angles at a reflecting surface is found by combining Eqs. (138) and (139), to be

$$u' = -u - 2cy \quad (141)$$

Refraction and Reflection United: Surface Power

Reflection and refraction can be treated the same way mathematically by thinking of reflection as refraction with $n' = -n$, in which case Eq. (140) gives Eq. (141). A reflecting surface can be represented graphically as a thin convex-plano or concave-plano thin lens with index $-n$, where n is the index of the medium, Fig. 18. For both refraction and reflection,

$$n' u' = nu - y\phi \quad (142)$$

where the surface power ϕ is

$$\phi = c(n' - n) \quad (143)$$

If the surface is approached from the opposite direction, then n and n' are switched, as is the sign of c , so ϕ is the same in both directions. Thus ϕ is a scalar property of the interface, which can be positive, negative, or zero. The power is zero if $n' = n$ or $c = 0$. If $n' = n$, the surface is "invisible," and the rays are not bent. If $c = 0$, the rays are bent. For a planar refracting surface $n' u' = nu$, and a planar reflecting surface gives $u' = -u$.

Principal Focal Lengths of a Surface

A ray approaching a surface parallel to the axis ($u = 0$) with a height y has an outgoing angle given by

$$n' u' = -y\phi \quad (144)$$

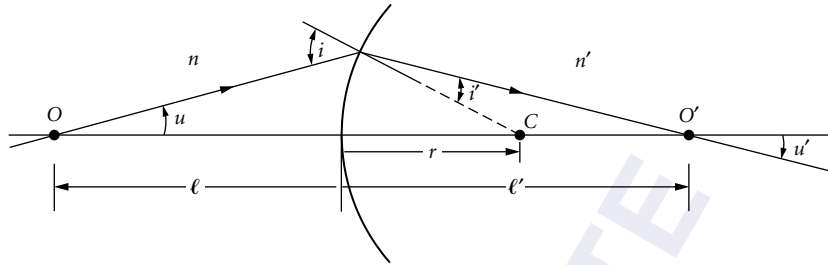


FIGURE 6 Refraction at a single spherical surface with center C and radius r . Axial object point O is imaged at O' .

This ray intercepts the axis at the *rear focal point*, whose directed distance from the vertex is $f' = y/u' = n'/\phi$. This directed distance is the *rear focal length*. Similarly, a ray entering from the right with $u' = 0$ intercepts the axis at the *front focal point*, a directed distance from the vertex of $f = y/u = -n/\phi$, the *front focal length*. Thus, a surface has a single power and two focal lengths, among which the following relationships hold:

$$f' = \frac{n}{\phi}, \quad f = -\frac{n}{\phi}, \quad \phi = -\frac{n}{f} = \frac{n'}{f'}, \quad \frac{f'}{f} = -\frac{n'}{n} \quad (145)$$

For a refracting surface, the signs of f' and f are opposite. For a reflecting surface $f' = f$.

Axial Object and Image Locations for a Single Surface

A ray from an axial point a directed distance l from the vertex of a surface that makes an angle u with the axis intersects the surface at height $y = -l/u$. After refraction or reflection, the ray angle is u' , and the ray intersects the axis at a distance $l' = -y/u'$ from the vertex, Fig. 6. Substituting for u and u' in Eq. (142), the relationship between axial object and image distances is

$$\frac{n'}{l'} = \frac{n}{l} + \phi \quad (146)$$

This can also be written

$$n \left(\frac{1}{r} - \frac{1}{l} \right) = n' \left(\frac{1}{r} - \frac{1}{l'} \right) \quad (147)$$

This is a special case of the equations below for imaging about a given ray. The transverse magnification is $m = l'/l$.

Paraxial Ray Tracing

Paraxial rays are traced through an arbitrary lens by a sequence of transfers between surfaces and power operations at surfaces. Each transfer changes height but not angle, and each power operation changes angle but not height. An image can be found by applying Eq. (136) and Eq. (142) successively. Alternatively, matrix methods described later or in Sec. 1.17 can be used.

Linearity of Paraxial Optics

For both the transfer and power operations, the outgoing heights and angles depend linearly on the incoming heights and angles. So a system described by a sequence of such operations is also linear. Therefore, a ray that enters with height y and angle u leaves with $y'(y, u)$ and $u'(y, u)$ given by

$$y' = \left(\frac{\partial y'}{\partial y} \right) y + \left(\frac{\partial y'}{\partial u} \right) u \quad \text{and} \quad u' = \left(\frac{\partial u'}{\partial y} \right) y + \left(\frac{\partial u'}{\partial u} \right) u \quad (148)$$

These equations can also be thought of as the first terms of Taylor expansions of exact expressions for $y'(y, u)$ and $u'(y, u)$. These partial derivatives depend on the structure of the system, and they can be determined by tracing two rays through the system. The partial derivatives, other than $\partial u'/\partial y$, also depend on the axial locations of the input and output surfaces. The changes with respect to these locations are treated easily by matrix methods.

The Two-Ray Paraxial Invariant

The various rays that pass through a lens are not acted upon independently, so there are several invariants that involve groups of rays. Consider two meridional paraxial rays that pass through a lens. At a given plane, where the medium has an index n , one ray has height y_1 and angle u_1 , and the other has y_2 and u_2 . The quantity

$$L = n(y_1 u_2 - y_2 u_1) \quad (149)$$

which we refer to as the *paraxial invariant* (NS), is unchanged as the rays pass through the system. Applying Eq. (136) and Eq. (142) to the above expression shows that this quantity does not change upon transfer or upon refraction and reflection. The invariant is also related to the general skew invariant, Eq. (73), since a paraxial skew ray can be decomposed into two meridional rays.

Another version of the invariance relationship is as follows. Two objects with heights y_1 and y_2 are separated axially by d_{12} . If their image heights are y'_1 and y'_2 , and the image separation is d'_{12} , then

$$n \frac{y_1 y_2}{d_{12}} = n' \frac{y'_1 y'_2}{d'_{12}} \quad (150)$$

An additional version of the invariance relationship is

$$\left(\frac{\partial y'}{\partial y} \right) \left(\frac{\partial u'}{\partial u} \right) - \left(\frac{\partial y'}{\partial u} \right) \left(\frac{\partial u'}{\partial y} \right) = \frac{n}{n'} \quad (151)$$

where the partial derivatives, Eq. (148), describe the action of any system.

The invariant applies *regardless* of the system. Thus, for example, if the lens changes, as with a zoom system, so that both of the outgoing rays change, their invariant remains. The invariant arises from basic physical principles that are manifested in a variety of ways, for example, as conservation of brightness and Liouville's theorem, discussed earlier in the section on conservation of étendue. This invariance shows that there are fundamental limits on what optical systems can do. Given the paraxial heights and angles of two input rays, only three of the four output heights and angles can be chosen arbitrarily. Likewise, only three of the four partial derivatives above can be freely chosen. The invariant is not useful if it vanishes identically. This occurs if the two rays are scaled versions of one another, which happens if both $u_1 = 0$ and $u_2 = 0$ for some z , or if both rays pass through the same axial object point, in which case $y_1 = 0$ and $y_2 = 0$. The invariant also vanishes if one of the rays lies along the axis, so that $y_1 = 0$ and $u_1 = 0$.

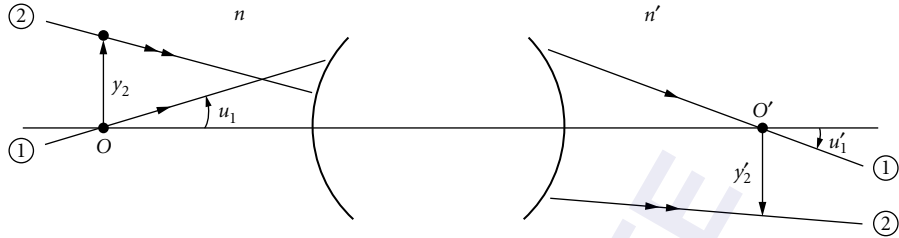


FIGURE 7 An object and image plane with ray 1 through the axial points and ray 2 through off-axis points. The location and magnification of an image plane can be found by tracing a ray from the axial object point O to axial image point O' . The magnification is given by Eq. (153). In the figure, u_1 and u'_1 have opposite signs, so the transverse magnification is negative.

Image Location and Magnification

To locate an image plane, any ray originating at the axial object point can be traced through the system to determine where it again intersects the axis, Fig. 7. The magnification for these conjugates can be found in two ways. One is to trace an arbitrary ray from any off-axis point in the object plane. The ratio of its height in the image plane to that in the object plane is the transverse magnification.

Alternately, the magnification can be found from the initial and final angles of the ray through the axial points. Let ray 1 leave the axial object point, so $y_1 = 0$. Let ray 2 originate in the object plane some distance from the axis. At the object plane $L = ny_2u_1$, and at the image plane $y'_2 = 0$, so $L = n'y'_2u'_1$. Therefore,

$$L = ny_2u_1 = n'y'_2u'_1 \quad (152)$$

So the magnification is

$$m = \frac{y'_2}{y_2} = \frac{nu_1}{n'u'_1} \quad (153)$$

The relative signs of u and u' determine that of the magnification. Equation (153) is a paraxial form of the sine condition Eq. (106). Squaring this equation gives $L^2 = n^2y_2^2u_1^2$, which is proportional to a paraxial form of the étendue. These matters are discussed further in the sections on conservation of étendue and on apertures. The quantity ny_2u_1 is sometimes referred to as *the invariant*, but it is not the most general form.

Three-Ray Rule

A further consequence of the paraxial invariant and of the linearity of paraxial optics is that once the paths of two paraxial meridional rays has been found, that of any third ray is determined. Its heights and angles are a linear combination of those of the other two rays. Given three rays, each pair has an invariant: $L_{12} = n(y_1u_2 - y_2u_1)$, $L_{23} = n(y_2u_3 - y_3u_2)$, and $L_{31} = n(y_3u_1 - y_1u_3)$. Therefore, in every plane

$$y_3 = -\frac{L_{23}}{L_{12}}y_1 + \frac{L_{31}}{L_{12}}y_2 \quad \text{and} \quad u_3 = -\frac{L_{23}}{L_{12}}u_1 + \frac{L_{31}}{L_{12}}u_2 \quad (154)$$

This assumes that no pair of the three rays are simply scaled versions of one another, i.e., that both $L_{23} \neq 0$ and $L_{31} \neq 0$.

Switching Axial Object and Viewing Positions

If an axial object and axial viewing position are switched, the apparent size of the image is unchanged. Put more precisely, let an object lie in a plane intersecting the axial point A and let its image be viewed from an axial point B' in image space that is not conjugate to A . If the object and viewing positions are switched, so the eye is at A and the object plane is at B' , the subtense of the object as seen by the eye is unchanged.¹⁵⁶⁻¹⁵⁹

1.13 IMAGES ABOUT KNOWN RAYS

Given a ray, referred to here as the *central ray* (also “base ray”), other rays from a point on the central ray making a small angle with respect to it are focused at or near other points on the central ray. These foci can be determined if the path of a central ray is known, as well as the indices of the media through which it passes, and the principal curvatures at the surfaces where it intersects. Here indices are constant. At each intersection point with an optical surface, the wavefront has two principal curvatures, as does the surface. After refraction or reflection, the wavefront has two different principal curvatures. Accordingly, if a single point is imaged, there are two astigmatic focal lines at some orientation. These foci are perpendicular, but they do not necessarily lie in planes perpendicular to that of the central ray. The imaging of a small extended region is generally skewed, so, for example, a small square in a plane perpendicular to the central ray can be imaged as a rectangle, parallelogram, or trapezoid.

This is a generalization of paraxial optics, in which the central ray is the axis of a system of revolution. While not difficult conceptually, the general case of an arbitrary central ray and an arbitrary optical system is algebraically complicated. This case can also be analyzed with a hamiltonian optics approach, using an expansion of a characteristic function about the central ray, like that of Eq. (28). The subject is sometimes referred to as *parabasal optics*, and the central ray as the *base ray*. This subject has been discussed by numerous authors¹⁶⁰⁻¹⁸⁶ under various names, e.g., “narrow beams,” “narrow pencils,” and “first order.”

The following is limited to the case of meridional central rays and surfaces that are figures of revolution. The surface, at the point of intersection, has two principal curvatures c_s and c_t . [See Eqs. (119) and (123).] For spherical surfaces, $c_s = c_t = c$, and for planar surfaces $c = 0$. There is a focus for the sagittal fan and one for the tangential one, Fig. 8, the two foci coinciding if the imaging is stigmatic. After one or more surfaces are encountered, the separated foci are the sources for subsequent imaging. Let s and t be the directed distances from the intersection point of the central ray and the surface to the object point, and s' and t' be the distances from intersection point to the foci. The separation $|s' - t'|$ is known as the *astigmatic difference*.

For refraction

$$\frac{n'}{s'} = \frac{n}{s} + c_s \Gamma \quad \text{and} \quad \frac{n' \cos^2 I'}{t'} = \frac{n \cos^2 I}{t} + c_t \Gamma \quad (155)$$

where $\Gamma = n' \cos I' - n \cos I$, Eq. (82). The sagittal equation is simpler, providing a mnemonic for remembering which equation is which: “S” = sagittal = simpler. If the surface is spherical, and the ray fan makes an arbitrary angle of ψ with the meridian, then¹⁷⁵

$$\frac{n'}{d'} (1 - \cos^2 \psi \sin^2 I') = \frac{n}{d} (1 - \cos^2 \psi \sin^2 I) + c \Gamma \quad (156)$$

where d and d' are the distances along the central ray from the surface to the object and image points. For normal incidence at a spherical surface $\Gamma = n' - n$, so both equations become

$$\frac{n'}{d'} = \frac{n}{d} + c(n' - n) \quad (157)$$

This also applies to surfaces of revolution if the central ray lies along the axis. This equation is identical to the paraxial equation, Eq. (146).

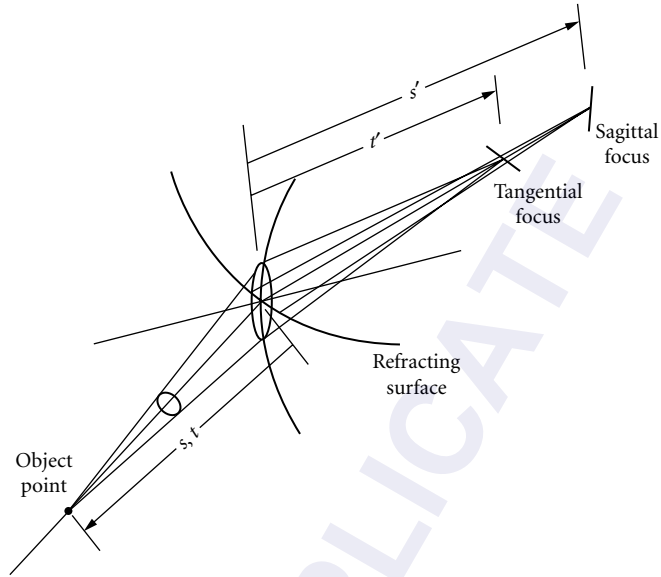


FIGURE 8 Astigmatic imaging of a point by a single refracting surface. The distance from the surface to the object point along the central ray of the bundle is $s = t$. The distances from the surface to the sagittal focus is s' , and that to the tangential focus is t' , as given by Eq. (155).

The corresponding relations for reflection are obtained by setting $n' = -n$ and $I' = I$ in the refraction equations, giving

$$\frac{1}{s'} = -\frac{1}{s} + 2c_s \cos I \quad \text{and} \quad \frac{1}{t'} = -\frac{1}{t} + \frac{2c_t}{\cos I} \quad (158)$$

For stigmatic imaging between the foci of reflective conics, $s = t$ is the distance from one focus to a point on the surface, and $s' = t'$ is that from the surface to the other focus. Therefore, $c_t = c_s \cos^2 I$. The reflection analogue to Eq. (156), for a spherical surface is

$$\frac{1}{d'} = -\frac{1}{d} + \frac{2c \cos I}{1 - \cos^2 \psi \sin^2 I} \quad (159)$$

These equations are known by several names, including *Coddington's equations*, *Young's astigmatic formulae*, and the *s- and t-trace formulae*.

1.14 GAUSSIAN LENS PROPERTIES

Introduction

The meaning of the term *gaussian optics* is not universally agreed upon, and it is often taken to be indistinguishable from paraxial optics or first-order optics, as well as collineation. Here the term is considered to apply to those aspects of paraxial optics discovered by Gauss,¹⁸⁷ who recognized that all rotationally symmetric systems of lens elements can be described paraxially by certain system properties. In particular, lenses can be treated as black boxes described by two axial length parameters and the locations of special points, called *cardinal points*, also called *Gauss points*. Once a lens is so

TABLE 1 Gaussian Notation and Definitions

By convention, in the diagrams the object space is to the left of the lens, image space is to the right, and rays go left to right. Object space quantities are unprimed, and image space quantities are primed, and quantities or positions that correspond in some way have same symbol, primed and unprimed. This correspondence can have several forms, e.g., the same type of thing or conjugate. The term *front* refers to object space, or left side, and *rear* to image space, or right side. A “front” entity may actually be behind a “rear” one. For example, a negative singlet has its object space focal point behind lens.

Scalars
n and n' object and image space refractive indices ϕ power m transverse magnification m_N nodal plane magnification = n/n' m_l longitudinal magnification m_α angular magnification u and u' paraxial ray angles (the positive direction is counterclockwise from the axis) y and y' paraxial ray heights y_p paraxial ray height at the principal planes = y'_p
Axial points
Cardinal points: Focal points F and F' , not conjugate Principal points P and P' , conjugate $m = +1$ Nodal points N and N' , conjugate $m_N = n/n'$ Other points: Axial object and image points O and O' , conjugate Arbitrary object and image points A and A' , B and B' Vertices V and V' , not conjugate, in general
Directed axial distances
These distances here are between axial points and are directed. Their signs are positive if from left to right and vice versa. Types of distances: entirely in object or image space, between spaces Principal focal lengths: $f = PF$ and $f' = P'F'$ Principal points to object and image axial points: $l = PO$ and $l' = P'O'$ Front and rear focal points to object and image axial points: $z = FO$ and $z' = F'O'$ Relations: $l = f + z$ and $l' = f' + z'$ Arbitrary point to conjugate object and image points: $d = AO$ and $d' = A'O'$
Distances between object space and image space points involve distances within both spaces, as well as a distance <i>between</i> the spaces, e.g., PP' , FF' , VV' , and OO' . The distances between spaces depend on the particular structure of the lens. They can be found by paraxial ray tracing.

characterized, knowledge of its actual makeup is unnecessary for many purposes, and repeated ray traces need not be performed. For example, given the object location, the image location and magnification are determined from the gaussian parameters. From the gaussian description of two or more lenses, that of a coaxial combination can be found. Another consequence of Gauss's discovery is that there is an infinity of specific embodiments for any external prescription.

The lenses considered in this section are figures of revolution with uniform object space and image space indices n and n' . All quantities discussed in this section are paraxial, so the prefix “paraxial” is not repeated. For the purposes of this section, no distinction is made between real and virtual rays. Those in each space are considered to extend infinitely, and intersection points may be either accessible or inaccessible. The quantities used in this section are found in Table 1.

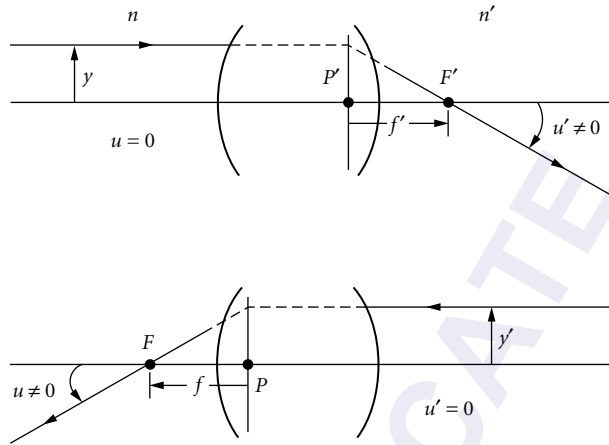


FIGURE 9 Diagrams for determining power, focal points, and focal lengths. Rays parallel to the axis in one space cross the axis in the other space at the focal points. The principal planes are at the intersections of entering and leaving rays. The power is given by Eq. (159). The lens in this diagram has positive power, a positive rear focal length, and a negative front focal length.

Power, Focal Lenses, and Afocal Lenses

A paraxial ray entering a lens parallel to the axis at a height y leaves with some angle u' , Fig. 9. Likewise, a ray entering from the opposite side with height y' leaves with angle u . The *power* of the lens is defined by

$$\phi = -n' \frac{u'}{y} = n \frac{u}{y'} \quad (160)$$

The outgoing ray can have any angle, and the power can be positive, negative, or zero. If $u' = 0$, then $\phi = 0$ and the lens is *afocal* or *telescopic*. Lenses for which $\phi \neq 0$ are referred to here as *focal*, although the term “nonafocal” is more common. Afocal lenses are fundamentally different from focal ones, and are treated separately next. Power is the same in both directions, i.e., whether the ray enters from left to right or from right to left. The lens in Fig. 9 has $\phi > 0$, and that in Fig. 10

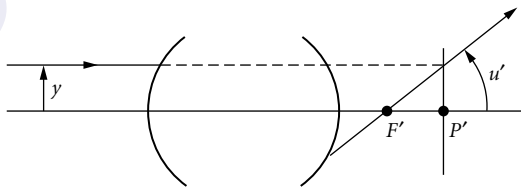


FIGURE 10 A lens with negative power and negative rear focal length. An incoming ray parallel to the axis with a positive height leaves the lens with a positive angle. The rear focal plane precedes the rear principal plane.

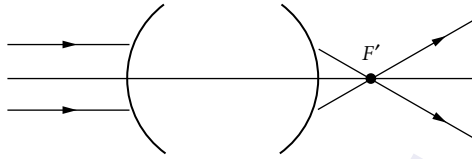


FIGURE 11 An ambiguous diagram. Two rays that enter a lens parallel to its axis converge at the rear focal point F' . Without specifying which ray is which, the sign of the power is not known.

has $\phi < 0$. Diagrams such as Fig. 11 show the location of the principal focal point, but not the sign of the power; two rays enter and two leave, but there is no indication of which is which. (Note that some negative lenses have accessible rear focal points.) Another expression for power involves two rays at arbitrary angles and heights. If two incident rays have (y_1, u_1) and (y_2, u_2) , and a nonzero invariant $L = n(y_1 u_2 - y_2 u_1)$, and the outgoing ray angles are u'_1 and u'_2 , then

$$\phi = -\frac{nn'}{L}(u'_1 u_2 - u'_2 u_1) \quad (161)$$

Focal Lenses

Focal lenses are those for which $\phi \neq 0$. Their cardinal points are the principal focal points, the principal points, and the nodal points. These points may be located anywhere on axis relative to the physical lens system. If they are inside a lens, then the intersection points referred to below are virtual. The cardinal points are pairs consisting of a member in object space and one in image space. The one in object space is often referred to as *front*, and the one in image space as *rear*, but this terminology may be misleading, since the points can be in any sequence along the axis.

Principal Focal Points Rays entering a lens parallel to its axis cross the axis at the *principal focal points* or *focal points*. Rays parallel to the axis in object space intersect the axis at the *rear focal point* F' in image space and those parallel in image space intersect at the *front focal point* F in object space, Fig. 9. The *principal focal planes* or *focal planes* are the planes perpendicular to the axis at the focal points. The terms *focal point* and *focal plane* are often used to refer to the images of any point or plane. In this chapter, *image point* is used for other points where rays are focused and *image plane* for other planes.

Principal Planes The *principal planes* are the conjugate planes for which the transverse magnification is unity, Fig. 12. The intersections of the principal planes and the axis are the *principal points*,

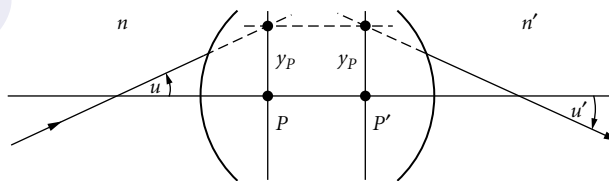


FIGURE 12 Principal planes as effective ray-bending surfaces. Incoming and outgoing paraxial rays intersect the object and image space principal planes at the same height y_P . The angles are related by Eq. (161).

denoted by P and P' . The *rear principal plane* is the locus of intersections between $u = 0$ rays incident from the left and their outgoing portions, Fig. 9. Likewise, the *front principal plane* is the intersection so formed with the rays for which $u' = 0$. A ray intersecting the first principal plane with height y_p and angle u leaves the second principal plane with height $y' = y_p$ and an angle given by

$$n' u' = nu - y_p \phi \quad (162)$$

The lens behaves as if the incoming ray intercepts the front principal plane, is transferred to the second with its height unchanged, and is bent at the second by an amount proportional to its height and to the power of lens. The power of the lens determines the amount of bending. For rays passing through the principal points, $y_p = 0$, so $u'/u = n/n'$.

Principal Focal Lengths The *focal lengths*, also called *effective focal lengths*, are the directed distances from the principal points to the focal points. The front and rear focal lengths are

$$PF = f = -\frac{n}{\phi} \quad \text{and} \quad P'F' = f' = \frac{n'}{\phi} \quad (163)$$

The two focal lengths are related by

$$\phi = -\frac{n}{f} = \frac{n'}{f'} \quad \text{and} \quad \frac{f}{f'} = -\frac{n}{n'} \quad (164)$$

This ratio is required by the paraxial invariant.¹⁸⁸ If $n = n'$, then $f' = -f$. If $n = n' = 1$, then

$$f' = -f = \frac{1}{\phi} \quad (165)$$

The focal lengths are the axial scaling factors for the lens, so axial distances in all equations can be scaled to them.

Nodal Points The *nodal points* are points of unit angular magnification. A paraxial ray entering the *object space nodal point* N leaves the *image space nodal point* N' at the same angle, Fig. 13. The planes containing the nodal points are called *nodal planes*. A *nodal ray* is one that passes through the nodal points. Such a ray must cross the axis, and the point where it does so physically is sometimes called the *lens center*. In general, this point has no special properties. (Gauss suggested an alternate “lens center,” the point midway between the principal points. Rotating a lens front to rear about this point would leave object and image positions and magnifications unchanged.)

If the refractive indices of the object space and image space are the same, the nodal points correspond to the principal points. If not, both nodal points are shifted according to

$$PN = P'N' = \frac{n' - n}{\phi} = f + f' \quad (166)$$

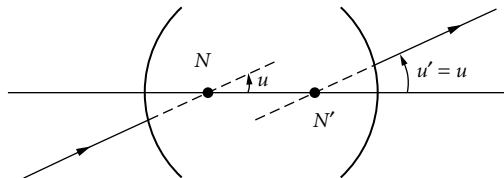


FIGURE 13 Nodal points. A paraxial ray through the object space nodal point N passes through image space nodal point N' with the same angle.

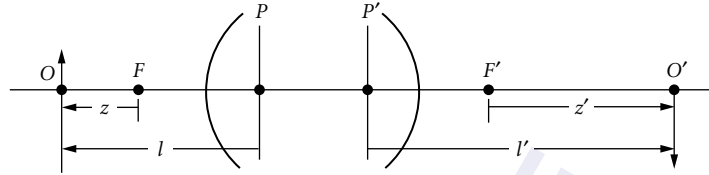


FIGURE 14 Directed distances used in common conjugate equations for focal lenses. Distances z and z' are from the focal points to axial conjugate points. Distances l and l' are from the principal points to axial conjugate points.

The distances from the nodal points to the focal points are

$$N'F' = -f \quad \text{and} \quad NF = -f' \quad (167)$$

The nodal points are conjugate, and the transverse magnification of the nodal planes is

$$m_N = \frac{n}{n'} \quad (168)$$

These equations can be recalled by the simple example of the single refracting surface, for which both nodal points correspond to the center of curvature.

Conjugate Equations For an object plane perpendicular to the axis at point O , there is an image plane perpendicular to the axis at O' , in which the transverse magnification is m . Note that specifying magnification implies both object and image positions. There is a variety of *conjugate equations* (NS) that relate their positions and magnifications. The equations differ in which object space and image space reference points are used from which to measure the directed distances to the object and image. These equations can be written in several ways, as given below, and with reference to Fig. 14. Axial distances can be scaled to the focal lengths, or the distances can be scaled to the indices, with a common power term remaining.

The simplest conjugate equation is *Newton's equation*, for which the reference points are the focal points and the lengths therefrom are $z = FO$ and $z' = F'O'$. The equation can be written in several forms:

$$zz' = ff' \quad \text{or} \quad \frac{z'}{f'} \frac{z}{f} = 1 \quad \text{or} \quad \frac{z'}{n'} \frac{z}{n} = \frac{1}{\phi^2} \quad (169)$$

More generally, if A and A' are any pair of axial conjugate points, as are B and B' , then

$$FA \times F'A' = FB \times F'B' \quad (170)$$

Another form is that for which the reference points are the principal points and the directed distances are $l = PO$ and $l' = P'O'$:

$$1 = \frac{f'}{l'} + \frac{f}{l} \quad \text{or} \quad \frac{n'}{l'} = \frac{n}{l} + \phi \quad (171)$$

If the reference points are arbitrary conjugates with magnification m_A and the axial distances are $d = AO$ and $d' = A'O'$, then

$$m_A \frac{n'}{d'} = \frac{1}{m_A} \frac{n}{d} + \phi \quad \text{or} \quad \frac{d'}{f'} = \frac{m_A^2 \frac{d}{f}}{1 - m_A \frac{d}{f}} \quad (172)$$

This equation also relates the curvatures of a wavefront at conjugate points. For a point source at A the radius of the wavefront at O is d , so at O' the radius is d' .

If the reference points are the nodal points, $m_A = m_N = n/n'$, and the axial distances are $d = NO$ and $d' = N'O'$, then

$$1 = \frac{f}{d'} + \frac{f'}{d} \quad \text{or} \quad \frac{n}{d'} = \frac{n'}{d} + \phi \quad (173)$$

The most general equation relating conjugate points is obtained when both reference points are arbitrary. Let the reference point in object space be a point A , at which the magnification is m_A , and that in image space be B' , associated with magnification m'_B . If $d = AO$ and $d' = B'O'$, then

$$\begin{aligned} \frac{1}{\phi} \left(1 - \frac{m'_B}{m_A} \right) &= \frac{1}{m_A} d - m'_B d' + \phi d d' \quad \text{or} \\ d' &= \frac{\frac{1}{m_A} d + \left(\frac{m'_B}{m_A} - 1 \right) \frac{1}{\phi}}{\phi d - m'_B} \end{aligned} \quad (174)$$

All the other conjugate equations are special cases of this one with the appropriate choice of m_A and m'_B .

If the reference point in object space is the focal point, and that in image space is the principal plane, then $m_A = \infty$ and $m'_B = 1$, giving

$$\frac{n'}{z'\phi} = \frac{l\phi}{n} + 1 \quad \text{or} \quad \frac{f'}{z'} = \frac{l}{f} + 1 \quad (175)$$

Likewise, if the object space reference point is P and the image space reference is F' , then

$$\frac{n'}{l'\phi} = \frac{z\phi}{n} + 1 \quad \text{or} \quad \frac{f'}{l'} = \frac{z}{f} + 1 \quad (176)$$

A relationship between distances to the object and image from the principal points and those from the focal points is

$$1 = \frac{z'}{l'} + \frac{z}{l} = \frac{F'O'}{P'O'} + \frac{FO}{PO} \quad (177)$$

Transverse Magnification In planes perpendicular to the axis, the *transverse magnification*, usually referred to simply as the *magnification*, is

$$m = \frac{x'}{x} = \frac{y'}{y} = \frac{dx'}{dx} = \frac{dy'}{dy} \quad (178)$$

There are several equations for magnification as a function of object position or image position, or as a relationship between the two. Newton's equations are

$$m = -\frac{f}{z} = -\frac{z'}{f'} = \frac{f}{f-l} = \frac{f'-l'}{f'} \quad (179)$$

Other relationships are

$$m = \frac{n}{n'} \frac{l'}{l} = -\frac{f}{f'} \frac{l'}{l} = \frac{z'}{l'} \frac{l}{z} \quad (180)$$

If $n = n'$, then $m = l'/l$. Another form, with respect to conjugate planes of magnification m_A is

$$mm_A = \frac{n d'}{d n'} = \frac{f d'}{d f'} \quad (181)$$

If d and d' are distances from the nodal points, $m = d'/d$. The change of magnification with respect to object or image positions with conjugacy maintained is

$$\frac{dm}{dz'} = -\frac{1}{f'} = \frac{m}{z'} \quad \text{and} \quad \frac{dm}{dz} = \frac{f}{z^2} = \frac{m^2}{f} = \frac{m}{z} \quad (182)$$

Images of Distant Objects If an object at a great distance from the lens subtends an angle ψ from the axis at the lens, then its paraxial linear extent is $y = z\psi$. The image height is

$$y' = my = \frac{-f}{z} y = -f\psi = \frac{n}{n'} f' \psi \quad \text{and} \quad \frac{dy'}{d\psi} = \frac{n}{n'} f' \quad (183)$$

If a distant object moves perpendicularly to the axis, then its image moves in the opposite direction if $f' > 0$ and in the same direction if $f' < 0$, so long as n and n' have the same sign.

Distance Between Object and Image The directed distance from an axial object point to its image contains three terms, one in object space, one in image space, and one relating the two spaces. The first two depend on the magnification and focal lengths. The interspace term depends on the particular structure of the lens, and is found by paraxial ray tracing. The most commonly used interspace distance is PP' , since it equals zero for a thin lens, but the equations using FF' are simpler. Newton's equations give $z = -f/m$ and $z' = -mf'$, so the object-to-image distance is

$$OO' = FF' - z + z' = FF' - f'm + \frac{f}{m} = FF' - \frac{1}{\phi} \left(n'm + \frac{n}{m} \right) \quad (184)$$

This is the basic equation from which all others are derived. If the reference points are the principal points, then

$$OO' = PP' + f'(1-m) - f \left(1 - \frac{1}{m} \right) = PP' + \frac{1}{\phi} \left[n'(1-m) + n \left(1 - \frac{1}{m} \right) \right] \quad (185)$$

If the object-to-image distance is given, the magnification is

$$m = \frac{1}{2n'} (-q \pm \sqrt{q^2 - 4nn'})$$

where $q = \phi(OO' - PP') - n - n'$. (186)

There are two magnifications for which OO' is the same. The magnitude of their product is n/n' . The derivative of the object-to-image distance with respect to the magnification is

$$\frac{d}{dm} OO' = -f' - \frac{f}{m^2} = -f' - \frac{z^2}{f} = \frac{1}{\phi} \left(\frac{n}{m^2} - n' \right) \quad (187)$$

Extrema occur at $m \pm \sqrt{n/n'}$, giving $m = \pm 1$ if $n = n'$. The extrema are

$$OO' - FF' = \pm \frac{2}{\phi} \sqrt{nn'} = \pm 2\sqrt{-ff'} \quad (188)$$

or

$$OO' - PP' = \frac{1}{\phi} (n + n' \pm 2\sqrt{nn'}) = f' - f \pm 2\sqrt{-ff'} \quad (189)$$

For the common case of $n' = n$, the object-to-image distance is

$$OO' = PP' + f' \left(2 - m - \frac{1}{m} \right) \quad (190)$$

OO' is the same for magnifications m and $1/m$. For a lens with $f' > 0$, the extremum object-to-image distances are $OO' - PP' = 4f'$ with $m = -1$ and $OO' - PP' = 0$ for $m = +1$. If the object-to-image distance and the focal length are given, then the magnification is

$$m = -\frac{1}{2}s \pm \sqrt{\frac{1}{4}s^2 - 1}$$

$$\text{where } s = \frac{1}{f'}(OO' - PP') - 2. \quad (191)$$

The two values of m are reciprocal of each other.

Axial Separations and Longitudinal Magnification Two axial points A and B are imaged at A' and B' with magnifications m_A and m_B . Newton's equations give the object separation

$$\Delta z = z_A - z_B = \frac{m_A m_B}{m_B - m_A} f \quad (192)$$

The separation of their images is

$$\Delta z' = z'_A - z'_B = (m_B - m_A) f' \quad (193)$$

The ratio of the image and object separations is

$$\frac{\Delta z'}{\Delta z} = \frac{z'_A - z'_B}{z_A - z_B} = \frac{A'B'}{AB} = \frac{n'}{n} m_A m_B = -\frac{f'}{f} m_A m_B \quad (194)$$

If m_A and m_B have different signs, then the direction of $A'B'$ is opposite to that of AB . This occurs when A and B are on opposite sides of the front focal point. In the limit as the separation between A and B vanishes, m_A and m_B both approach the same magnification m . The *longitudinal magnification* m_L is the ratio of axial separations in the limit of small separations

$$m_L = \lim_{A \rightarrow B} \frac{A'B'}{AB} = \frac{dz'}{dz} = \frac{n'}{n} m^2 = -\frac{z'}{z} \quad (195)$$

This quantity is also called the *axial magnification*. Since m^2 is always positive, as an object moves axially in a given direction, its image moves in a constant direction. There is a discontinuity in image position when the object crosses the focal point, but the direction of motion stays the same. At the nodal points, the transverse and longitudinal magnifications are equal.

Angular Magnification The ratio of the outgoing to incoming ray angles, u'/u , is sometimes called the *angular magnification* m_α . If the ray passes through conjugate axial points with magnification m , then the angular magnification is

$$m_\alpha = \frac{u'}{u} = \frac{n}{n'} \frac{1}{m} \quad (196)$$

If the ray leaves an object point with height y in a plane for which the magnification is m , the outgoing ray angle is given by

$$n' u' = \frac{1}{m} n u - y \phi = \frac{1}{m} (n u - y' \phi) \quad (197)$$

The ratio u'/u is not constant unless $y = 0$ or $\phi = 0$.

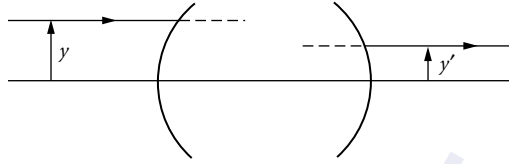


FIGURE 15 Afocal lens. Paraxial rays entering parallel to the axis leave parallel, in general at a different height. The ratio of the heights is the transverse magnification, which is constant.

Relationship Between Magnifications The transverse, angular, and longitudinal magnifications are related by

$$m_{\alpha} m_L = m \quad (198)$$

This relationship is connected to the paraxial invariant and also holds for afocal lenses.

Reduced Coordinates Many relationships are formally simplified by using reduced axial distances $\tau = z/n$ and $\tau' = z'/n'$ and reduced angles $\omega = nu$, $\omega' = n'u'$, which are paraxial optical direction cosines. For example, the angular magnification is $\omega'/\omega = 1/m$, and the longitudinal magnification is $d\tau'/d\tau = m^2$.

Mechanical Distances The cardinal points can be located anywhere on axis relative to the physical structure of the lens. The *vertex* of a lens is its extreme physical point on axis. The object space vertex is denoted by V and the image space vertex by V' . The two vertices are not, in general, conjugate. The *front focal distance* FV is that from the vertex to the front focal point, and the *rear focal distance* $V'F'$ is that from the rear vertex to the rear focal point. Likewise, the *front working distance* OV is the distance from the object to the vertex, and the *rear working distance* $V'O'$ is that from the vertex to the image. These lengths have no significance to the gaussian description of a lens. For example, a lens of a given focal length can have any focal distance and vice versa. For a *telephoto* lens the focal length is greater than the focal distance, and for a *retrofocus* lens the focal distance is greater than the focal length.

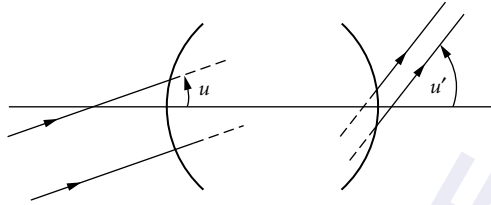
Afocal Lenses

An *afocal* or *telescopic* lens^{189–191} is one for which $\phi = 0$. A ray entering with $u = 0$ leaves with $u' = 0$, Fig. 15. There are no principal focal points or focal lengths. In general, the $u = 0$ ray leaves at a different height than that at which it enters. The ratio y'/y is the same for all such rays, so the transverse magnification m is constant. Likewise, the longitudinal magnification is constant, equaling $m_L = (n'/n)m^2$, as is the angular magnification $u'/u = m_{\alpha} = n/(n'm)$. A parallel bundle of rays entering at angle u leaves as a parallel bundle at $u' = m_{\alpha}u$, Fig. 16. Summarizing:

$$m = \text{const}, \quad m_L = \frac{n'}{n} m^2 = \text{const}, \quad m_{\alpha} = \frac{n}{n' m} = \text{const}, \quad m = m_L m_{\alpha} \quad (199)$$

Any two of these magnifications provide the two scaling factors that describe the system. If $m = n/n'$, then $m_L = m$ and $m_{\alpha} = 1$, so image space is a scaled version of object space.

Afocal lenses differ fundamentally from focal lenses. Objects at infinity are imaged by afocal lenses at infinity, and objects at finite distances are imaged at finite distances. An afocal lens has no cardinal points and the focal length is undefined. Afocal lenses have no principal planes. If $m \neq 1$ there are no unit magnification conjugates, and if $m = 1$ there is only unit magnification. Likewise, there are no nodal points; the angular magnification is either always unity or always differs from unity.



It is sometimes stated or implied that an afocal lens is a focal one with an infinite focal length, but this description is dubious. For example, the above equations relating magnification and conjugate positions to focal length are meaningless for afocal lenses, and they cannot be made useful by substituting $f = \infty$. The equations for the afocal lenses can be obtained from those for focal lenses with a limiting process, but for most purposes this approach is not helpful.

If the positions for a single axial conjugate pair A and A' are known, other pairs are located from the property of constant longitudinal magnification. If O and O' are another pair of conjugates, then

$$A'O' = m_L AO \quad (200)$$

As a function of distance AO , the object-to-image distance OO' is

$$OO' = AA' + (m_L - 1)AO \quad (201)$$

where AA' is the separation between the initially known conjugates. If $m_L = 1$, the object-to-image distance is constant. Otherwise, it can take any value. For all afocal lenses, except those for which $m_L = 1$, there is a position, sometimes called the *center*, at which $OO' = 0$, so the object and image planes coincide.

A principal use of afocal lenses is in viewing distant objects, as with binoculars. An object of height h at a great distance d from the lens subtends an angle $\psi \approx h/d$. The image height is $h' = mh$, and the image distance is approximately $d' \approx m^2 d$. So the image subtends an angle $\psi' \approx m\psi = \psi/m_\alpha$. Thus a telescope used visually produces an image which is actually smaller, but which is closer by a greater factor, so the subtense increases.

Determination of Gaussian Parameters

If a lens prescription is given, its gaussian properties can be obtained by paraxially tracing any two meridional rays whose invariant is not zero. A common choice for focal lenses is the rays with $u = 0$ and $u' = 0$, which give F , P , F' , and P' . If a lens is known to be afocal, a single ray not parallel to the axis suffices, since such a ray gives a pair of conjugates and the angular magnification. If it is not known that the lens is afocal, two rays show that it is, as well as giving the required information about conjugates. Alternately, a matrix representation of the lens can be determined, from which the cardinal points are found, as described in the matrix section. The gaussian properties can also be determined experimentally in a number of ways.

Basic Systems

Single Refracting Surface Media of indices n and n' are separated by a surface of curvature c and radius r . The power is $\phi = (n' - n)c$. The principal points coincide at the vertex. The nodal points coincide at the center of curvature. The distance from principal points to nodal points is r .

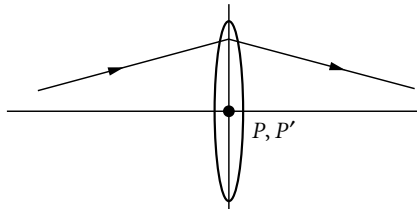


FIGURE 17 The thin lens approximation. The thickness of the lens is negligible, and the principal planes are coincident, so rays bend at the common plane.

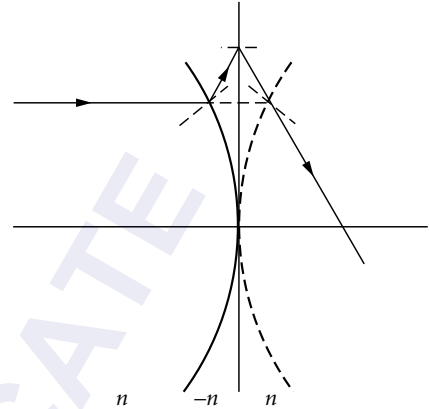


FIGURE 18 Reflecting surface represented unfolded. A convex plano thin lens with index $n' = -n$, where n is the index of the medium.

Thick Lens The term *thick lens* usually denotes a singlet whose vertex-to-vertex distance is not negligible, where negligibility depends on the application. For a singlet of index n in vacuum with curvatures c_1 and c_2 and thickness t , measured from vertex to vertex

$$\phi = \frac{1}{f'} = (n-1) \left[c_1 - c_2 - \frac{n-1}{n} t c_1 c_2 \right] \quad (202)$$

A given power may be obtained with a variety of curvatures and indices. For a given power, higher refractive index gives lower curvatures. The principal planes are located relative to the vertices by

$$VP = \frac{n-1}{n} \frac{t c_2}{\phi} \quad \text{and} \quad V'P' = -\frac{n-1}{n} \frac{t c_1}{\phi} \quad (203)$$

These equations can be derived by treating the lens as the combination of two refracting surfaces. Two additional relationships are

$$PP' = VV' - \frac{n-1}{n} \frac{t(c_1 + c_2)}{\phi} \quad \text{and} \quad \frac{V'P'}{VP} = \frac{r_2}{r_1} = \frac{c_1}{c_2} \quad (204)$$

Thin Lens A *thin lens* is the limiting case of a refracting element whose thickness is negligible, so the principal planes coincide, and the ray bending occurs at a single surface, Fig. 17. In the limit as $t \rightarrow 0$, for a lens in vacuum the thick lens expressions give

$$\phi = \frac{1}{f'} = (n-1)(c_1 - c_2), \quad VP = V'P' = 0, \quad PP' = 0 \quad (205)$$

Single Reflecting Surface A reflecting surface has power $\phi = 2n/r = 2nc$. The principal points are located at the vertex. The nodal points are at the center of curvature.

Mirror as a Thin Lens In unfolding systems, a mirror can be thought of as a convex or concave plano thin lens, with an index $-n$, where n is the index of the medium in which it works, Fig. 18. All the thin lens equations apply, as well as those for third-order aberration equations, which are not discussed here.

1.15 COLLINEATION

Introduction

Collineation is a mathematical transformation that approximates the imaging action of a lens with homogeneous refractive indices in both spaces. This transformation takes points to points, lines to lines, and planes to planes. With an actual lens, incoming rays become outgoing rays, so lines go exactly to lines. In general, however, rays that intersect in object space do not intersect in image space, so points do not go to points, nor planes to planes. The collinear transformation is an approximate description of image geometry with the intervening optical system treated as a black box, not a theory that describes the process of image formation. Collineation is also referred to as *projective transformation*. The historical development of this approach, which was first applied to optics by Möbius,¹⁹² is discussed by Southall.¹⁹³ Several authors give extensive discussions.^{193–197} Projective transformation is used in computer graphics, and is discussed in this context in a number of recent books and papers.

The imaging described by collineation is, by definition, stigmatic everywhere, and planes are imaged without curvature. And for rotationally symmetric lenses, planes perpendicular to the axis are imaged without distortion. So the three conditions of maxwellian perfection are satisfied for all conjugates. Consequently, collineation is often taken as describing ideal imaging of the entire object space. However, it is physically impossible for a lens to image as described by collineation, except for the special case of an afocal lens with $m = m_L = n/n'$. The putative ray intersections of collineation violate the equality of optical path lengths for the rays involved in the imaging of each point. The intrinsic impossibility manifests itself in a variety of ways. As an example, for axial points in a plane with transverse magnification m and ray angles θ and θ' relative to the axis, collineation gives $m \propto \tan \theta / \tan \theta'$, but optical path length considerations require that $m \propto \sin \theta / \sin \theta'$. Another violation is that of the skew invariant $\mathcal{S} = n(\alpha\gamma - \beta x)$. The ratio of this quantity before and after collineation is not unity, but $\mathcal{S}'/\mathcal{S} = \gamma'/\gamma$, where γ is the axial direction cosine in object space and γ' is that in image space.

The expressions for collineation do not contain refractive indices, another manifestation of their not accounting for optical path length. Rather than the refractive index ratio n'/n , which occurs in many imaging equations, the expressions of collineation involve ratios of focal lengths. For afocal lenses there are ratios of transverse and longitudinal magnifications or ratios of the focal lengths of the lenses making up the afocal system.

The expressions for actual ray behavior take the form of collineation in the paraxial, and, more generally, paraxial limits. So paraxial calculations provide the coefficients of the transformation for any particular system.

Collineation is most often treated by starting with the general form, and then reducing its complexity by applying the symmetries of a rotationally symmetric system, to give familiar simple equations such as Newton's.¹⁹⁸ Alternatively, it is possible to begin with the simple forms and to derive the general ones therefrom with a succession of images, along with translations and rotations. However, the more important use of these properties is in treating lenses lacking rotational symmetry. This includes those comprising elements that are arbitrarily oriented, that is, tilted or decentered—either intentionally or unintentionally. Other examples are nonplanar objects, tilted object planes, and arbitrary three-dimensional object surfaces.

Lenses, along with plane mirror systems, can form a succession of images and can produce translations and rotations. Correspondingly, a succession of collinear transformations is a collinear transformation, and these transformations form a group. It is associative, corresponding to the fact that a series of imaging operations can be associated pairwise in any way. There is a unit transformation, corresponding physically to nothing or to a unit magnification afocal lens. There is an inverse, so an image distorted as a result of object or lens tilt can be rectified by an appropriately designed system—to the extent that collineation validly describes the effects.

General Equations

The general form of the collinear transformation is

$$x' = \frac{a_1x + b_1y + c_1z + d_1}{ax + by + cz + d}, \quad y' = \frac{a_2x + b_2y + c_2z + d_2}{ax + by + cz + d}, \quad z' = \frac{a_3x + b_3y + c_3z + d_3}{ax + by + cz + d} \quad (206)$$

At least one of the denominator coefficients, a , b , c , d , is not zero. The equations can be inverted, so there is a one-to-one correspondence between a point (x, y, z) in object space and a point (x', y', z') in image space. The inverted equations are formally identical, and can be written by replacing unprimed quantities with primed ones and vice versa in the above equation. It is seen that a plane is transformed to a plane, since $a'x' + b'y' + c'z' + d' = 0$ has the same form as a function of (x, y, z) . An intersection of two planes gives a line. It can also be shown that a line transforms to a line by writing the equation for a line in parametric form, with parameter σ , $x(\sigma) = x_0 + \alpha\sigma$, $y(\sigma) = y_0 + \beta\sigma$, $z(\sigma) = z_0 + \gamma\sigma$. Substituting in the transformation equations, it is found that $dx'/dy' = (dx'/d\sigma)/(dy'/d\sigma)$ is constant, as are other such ratios.

These equations contain 16 coefficients, but it is possible to divide all three equations through by one of the coefficients, so there are 15 independent coefficients in general. Since the location of an image point is described by three coordinates, five points that are not coplanar determine the transformation.

The ratios of the coefficient dimensions are determined by the fact that x , y , z and x' , y' , z' are lengths. A variety of schemes can be used and, in the expressions below, a given symbol may have different dimensions.

There are two major categories of the transformation, according to whether the denominator varies or is constant. That with a varying denominator corresponds to focal lenses. For afocal lenses, the denominator is constant, and the general form of the transformation is

$$x' = a_1x + b_1y + c_1z + d_1, \quad y' = a_2x + b_2y + c_2z + d_2, \quad z' = a_3x + b_3y + c_3z + d_3 \quad (207)$$

Here coefficient d has been normalized to unity. Such a transformation is called *affine* or *telescopic*.

Coordinate Systems and Degrees of Freedom

The transformation involves two coordinate systems. The origin of each is located by three parameters, as is the orientation of each. This leaves three parameters that describe the other aspects of the transformation for the most general case of no symmetry. The number is reduced to two if there is rotational symmetry.

In addition to considering the transformation of the entire space, there are other cases, especially the imaging of planes. In each situation, there are specific coordinate systems in which the aspects of the relationship, other than position and orientation, are most simply expressed. Accordingly, different coordinate systems are used in the following sections. Thus, for example, the z axis in one expression may not be the same as that for another.

Simplest Form of the General Transformation

For focal lenses, the denominators are constant for a set of parallel planes

$$ax + by + cz + d = \text{constant} \quad (208)$$

Each such plane is conjugate to one of a set of parallel planes in the other space. Within each of these planes, the quantities $\partial x'/\partial x$, $\partial x'/\partial y$, $\partial x'/\partial z$ are constant, as are the other such derivatives.

Therefore, magnifications do not vary with position over these planes, although they do vary with direction. There is one line that is perpendicular to these planes in one space whose conjugate is perpendicular to the conjugate planes in the other space. It can be taken to be the z axis in one space and the z' axis in the other. The aximuths of the x - y and x' - y' axes are found by imaging a circle in each space, which gives an ellipse in the other. The directions of the major and minor axes determine the orientations of these coordinate axes. The *principal focal planes* are the members of this family of planes for which

$$0 = ax + by + cz + d \quad (209)$$

Lines that are parallel in one space have conjugates that intersect at the principal focal plane in the other. The *principal focal points* are the intersection of the axes with the focal planes.

Using these simplifying coordinate systems, the general transformation is

$$x' = \frac{a_1 x}{cz + d}, \quad y' = \frac{b_1 y}{cz + d}, \quad z' = \frac{c_3 z + d_3}{cz + d} \quad (210)$$

One of the six coefficients can be eliminated, and two of the others are determined by the choice of origins for the z axis and z' axis. If the origins are taken to be at the principal focal points, the transformation equations are

$$x' = \frac{e_x x}{z}, \quad y' = \frac{e_y y}{z}, \quad z' = \frac{e_z}{z} \quad (211)$$

where e_x , e_y , e_z are constants. Unless $e_x = e_y$ the images of shapes in constant z planes vary with their orientations. Squares in one orientation are imaged as rectangles, and in others as parallelograms. Squares in planes not perpendicular to the axes are imaged, in general, with four unequal sides.

For afocal lenses, the simplest form is

$$x' = m_x x, \quad y' = m_y y, \quad z' = m_z z \quad (212)$$

Spheres in one space are imaged as ellipsoids in the other. The principal axes of the ellipsoids give the directions of the axes for which the imaging equations are simplest.

Conjugate Planes

A pair of conjugate planes can be taken to have $x = 0$ and $x' = 0$, so the general transformation between such planes is

$$y' = \frac{b_2 y + c_2 z + d_2}{by + cz + d}, \quad z' = \frac{b_3 y + c_3 z + d_3}{by + cz + d} \quad (213)$$

There are eight independent coefficients, so four points that are not in a line define the transformation. In each space, two parameters specify the coordinate origins and one the orientation. Two parameters describe the other aspects of the transformation.

The simplest set of coordinates is found by a process like that described above. For focal lenses, constant denominators define a line set of parallel lines

$$by + cz + d = \text{constant} \quad (214)$$

with similar conjugate lines in the other space. There is a line that is perpendicular to this family in one space, whose conjugate is perpendicular in the other, which can be taken as the z axis on one side and the z' axis on the other. There is a *principal focal line* in the plane in each space, and a *principal focal point*, at its intersection with the axis. In this coordinate system the transformation is

$$y' = \frac{b_2 y}{cz + d}, \quad z' = \frac{c_3 z + d_3}{cz + d} \quad (215)$$

Of the six coefficients, four are independent and two are fixed by the choice of origins. If $z = 0$ and $z' = 0$ are at the principal focal points, then

$$y' = \frac{e_y y}{z}, \quad z' = \frac{e_z}{z} \quad (216)$$

where e_y and e_z are constants.

For afocal lenses, the general transformation between conjugate planes is

$$y' = b_2 y + c_2 z + d_2, \quad z' = b_3 y + c_3 z + d_3 \quad (217)$$

The simplest form of the transformation is

$$y' = m_y y, \quad z' = m_z z \quad (218)$$

where m_y and m_z are constants.

Conjugate Lines

A line can be taken to have $x = 0, y = 0, x' = 0, y' = 0$, so its transformation is

$$z' = \frac{c_3 z + d_3}{c z + d} \quad (219)$$

There are three independent coefficients, so three points determine them. The origins in the two spaces account for two of the parameters, leaving one to describe the relative scaling. The simplest forms are

$$\text{Focal: } z' = \frac{e_z}{z}; \quad \text{Afocal: } z' = m_z z \quad (220)$$

There is a relationship between distances along a line (or ray) that is unchanged in collineation.^{193,199} If four points on a line A, B, C, D have images A', B', C', D' , the *double ratio* or *cross ratio* is invariant under projective transformation, that is,

$$\frac{AC}{BC} \frac{BD}{AD} = \frac{A'C'}{B'C'} \frac{B'D'}{A'D'} \quad (221)$$

where AC is the distance from A to C , and likewise for other pairs.

Matrix Representation of the Transformation

The transformation can be expressed in linear form by using the variables (u_1, u_2, u_3, u_4) and (u'_1, u'_2, u'_3, u'_4) , where $x = u_1/u_4, y = u_2/u_4, z = u_3/u_4$ and $x' = u'_1/u'_4, y' = u'_2/u'_4, z' = u'_3/u'_4$. These are referred to as *homogeneous coordinates*. The transformation can be written

$$\begin{pmatrix} u'_1 \\ u'_2 \\ u'_3 \\ u'_4 \end{pmatrix} = \begin{pmatrix} a_1 & b_1 & c_1 & d_1 \\ a_2 & b_2 & c_2 & d_2 \\ a_3 & b_3 & c_3 & d_3 \\ a & b & c & d \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} \quad (222)$$

In terms of the cartesian coordinates and an additional pair of terms q and q' , the transformation can be expressed as

$$\begin{pmatrix} q'x' \\ q'y' \\ q'z' \\ q' \end{pmatrix} = \begin{pmatrix} a_1 & b_1 & c_1 & d_1 \\ a_2 & b_2 & c_2 & d_2 \\ a_3 & b_3 & c_3 & d_3 \\ a & b & c & d \end{pmatrix} \begin{pmatrix} qx \\ qy \\ qz \\ q \end{pmatrix} \quad (223)$$

The dimensions of q and q' depend on the choice of coefficient dimensions. Here $q'/q = ax + by + cz + d$, the equation for the special set of planes.

Certain sections of the matrix are associated with various aspects of the transformation.²⁰⁰ The first three elements in the rightmost column have to do with translation. This is shown by setting $(x, y, z) = (0, 0, 0)$ to locate the conjugate in the other space. The first three elements in the bottom row are related to perspective transformation. The upper left-hand 3×3 array expresses rotation, skew, and local magnification variation.

For the simple form of the transformation expressed in Eq. (211), $a_1 = e_x$, $b_2 = e_y$, $d_3 = e_z$, $c = 1$, and the rest of the coefficients vanish. The general matrix representation for the afocal transformation is

$$\begin{pmatrix} x' \\ y' \\ z' \\ 1 \end{pmatrix} = \begin{pmatrix} a_1 & b_1 & c_1 & d_1 \\ a_2 & b_2 & c_2 & d_2 \\ a_3 & b_3 & c_3 & d_3 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix} \quad (224)$$

The quantities q and q' can also be included, in which case $q' = q$. In the simplest afocal form, Eq. (212), the matrix is diagonal with $a_1 = m_x$, $b_2 = m_y$, $d_3 = m_z$, and the rest of the nondiagonal coefficients vanishing. A succession of collineations can be treated by multiplying the matrices that describe them.²⁰¹ To combine lenses with arbitrary orientations and to change coordinate systems, compatible rotation and translation matrices are required. The transformation for a pure rotation with direction cosines (L, M, N) is

$$\begin{pmatrix} x' \\ y' \\ z' \\ 1 \end{pmatrix} = \begin{pmatrix} 1-2L^2 & -2LM & -2LN & 0 \\ -2LM & 1-2M^2 & -2MN & 0 \\ -2LN & -2MN & 1-2N^2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix} \quad (225)$$

The transformation for translation by $(\Delta x, \Delta y, \Delta z)$ is

$$\begin{pmatrix} x' \\ y' \\ z' \\ 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & \Delta x \\ 0 & 1 & 0 & \Delta y \\ 0 & 0 & 1 & \Delta z \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix} \quad (226)$$

The quantities q and q' can be included if necessary. The transformations associated with conjugate planes can likewise be expressed with 3×3 matrices, and the transformations of lines with 2×2 matrices.

Rotationally Symmetric Lenses

For rotationally symmetric lenses, the simplest forms are obtained with the z and z' axes corresponding to the lens axis in the two spaces. There is one less degree of freedom than in the general case, and $a_1 = b_2$ in Eq. (210). The general transformation is thus

$$x' = \frac{a_1 x}{cz + d}, \quad y' = \frac{a_1 y}{cz + d}, \quad z' = \frac{c_3 z + d_3}{cz + d} \quad (227)$$

There are four degrees of freedom, two associated with the lens and two with the choice of coordinate origins. For focal lenses, the two axial length parameters are f and f' . If the coordinate origins are at the focal points,

$$x' = -\frac{fx}{z}, \quad y' = -\frac{fy}{z}, \quad z' = \frac{ff'}{z} \quad (228)$$

If the coordinate origins are conjugate and related by magnification m_0 , then

$$x' = \frac{m_0 x}{1 + z/f}, \quad y' = \frac{m_0 y}{1 + z/f}, \quad z' = \frac{(f'/f)m_0^2 z}{1 + z/f} \quad (229)$$

The constant term in the numerator of the z' expression is the longitudinal magnification for $z = 0$, for which point $dz'/dz = (f'/f)m_0^2$. A special case of these equations is that for which the principal points are the origins, so $m_0 = 1$.

For rotationally symmetric afocal lenses, the two degrees of freedom are the transverse magnification $m_x = m_y = m$, and the longitudinal magnification $m_z = m_L$. The simplest set of transformation equations is

$$x' = mx, \quad y' = my, \quad z' = m_L z \quad (230)$$

where $z = 0$ and $z' = 0$ are conjugate. If $m = \pm 1$ and $m_L = \pm 1$ the image space replicates object space, except possibly for orientation. If $m_L = m$, the spaces are identical except for overall scaling and orientation. The m and m_L appear as functions of ratios of focal lengths of the lenses that make up the afocal system.

Rays for Rotationally Symmetric Lenses

A skew ray with direction cosines (α, β, γ) in object space is described in parametric form with parameter z as follows

$$x(z) = x_0 + \frac{\alpha}{\gamma} z, \quad y(z) = y_0 + \frac{\beta}{\gamma} z \quad (231)$$

For a focal lens, if $z = 0$ is taken to be the front focal plane, and $z' = 0$ is the rear focal plane, the parametric form of the ray in image space is

$$x'(z') = \left(-f \frac{\alpha}{\gamma}\right) + \left(-\frac{x_0}{f'}\right) z', \quad y'(z') = \left(-f \frac{\beta}{\gamma}\right) + \left(-\frac{y_0}{f'}\right) z' \quad (232)$$

So $x'_0 = -f\alpha/\gamma$, $y'_0 = -f\beta/\gamma$, $\alpha'/\gamma' = -x_0/f'$, $\beta'/\gamma' = -y_0/f'$. For meridional rays with $x = 0$, if θ and θ' are the ray angles in the two spaces, then $\tan \theta = \beta/\gamma$, $\tan \theta' = -y_0/f'$, and

$$\frac{\tan \theta}{\tan \theta'} = \frac{f'}{f} m \quad (233)$$

where m is the transverse magnification in a plane where the meridional ray crosses the axis.

For afocal lenses, if $z = 0$ and $z' = 0$ are conjugate planes, the ray in image space is given by

$$x'(z') = mx_0 + \left(\frac{m}{m_L} \frac{\alpha}{\gamma}\right) z', \quad y'(z') = my_0 + \left(\frac{m}{m_L} \frac{\beta}{\gamma}\right) z' \quad (234)$$

For meridional rays with $x = 0$,

$$\frac{\tan \theta}{\tan \theta'} = \frac{m_L}{m} \quad (235)$$

Tilted Planes with Rotationally Symmetric Lenses

A plane making an angle θ with the lens axis in object space has an image plane that makes an angle θ' , given by Eq. (233), the so-called *Scheimpflug condition*.^{202,203} A tilted plane and its image are perpendicular to a meridian of the lens, Fig. 19. There is bilateral symmetry on these planes about the

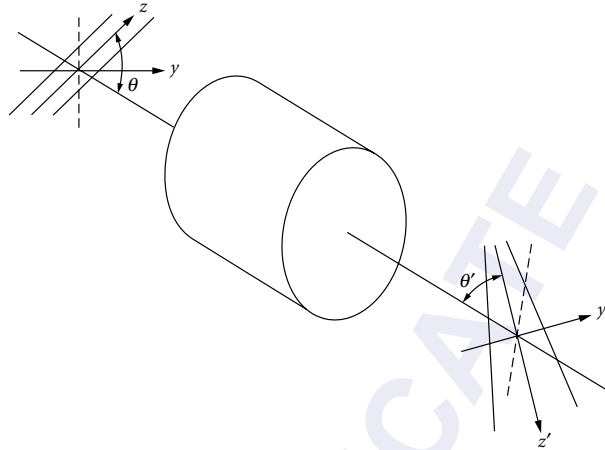


FIGURE 19 The image plane for a tilted object plane. The y - z plane is the object plane and the y' - z' plane is the image plane. The angles between the planes and the lens axis are θ and θ' , which are related by Eq. (232). The conjugate points in these planes are related by Eq. (235).

intersection line with the meridian, which is taken to be the z axis in object space and the z' axis in image space. The perpendicular coordinates are y and y' . Letting m_0 be the transverse magnification for the axial point crossed by the planes, the transform equations are

$$y' = \frac{m_0 y}{1 + z/g}, \quad z' = \frac{(g'/g)m_0^2 z}{1 + z/g} \quad (236)$$

Here g and g' are the focal lengths in the tilted planes, the distances from the principal planes to the focal planes of the lens, measured along the symmetry line, so

$$g = \frac{f}{\cos \theta}, \quad g' = \frac{f'}{\cos \theta'}, \quad \text{and} \quad \frac{g'}{g} = \sqrt{\left(\frac{f'}{f}\right)^2 \cos^2 \theta + \frac{1}{m_0^2} \sin^2 \theta} \quad (237)$$

As $\theta \rightarrow 90^\circ$, g and g' become infinite, and $(g'/g)m_0 \rightarrow 1$, giving $y' \rightarrow m_0 y$ and $z' \rightarrow m_0 z$. (Forms like Newton's equations may be less convenient here, since the distances from the axes to the focal points may be large.)

For an afocal lens with transverse magnification m and longitudinal magnification m_L , the object and image plane angles are related by Eq. (235). The conjugate equations for points in the planes are

$$y' = my, \quad z' = (m_L^2 \cos^2 \theta + m^2 \sin^2 \theta)^{1/2} z \quad (238)$$

Here the origins may be the axial intersection point, or any other conjugate points.

Some General Properties

For all collinear transformations, points go to points, lines to lines, and planes to planes. In general, angles at intersections, areas, and volumes are changed. The degree of a curve is unchanged, so, for example, a conic is transformed into a conic. For focal systems, a "closed" conic, an ellipse or circle, may be imaged as either a closed or an "open" one, a parabola or hyperbola. For afocal systems, the

closedness and openness are preserved. With focal systems, the imaging of a shape varies with its location, but for afocal systems it does not. For afocal systems parallelness of lines is maintained, but for focal systems the images of parallel lines intersect. For afocal systems, equal distances along lines are imaged as equal distances, but are different unless the magnification is unity.

1.16 SYSTEM COMBINATIONS: GAUSSIAN PROPERTIES

Introduction

This section deals with combinations of systems, each of which is of arbitrary complexity. From a gaussian description of each lens and the geometry of the combination, the gaussian description of the net system can be found. If two rotationally symmetric lenses are put in series with a common axis, the resultant system is also rotationally symmetric. Its gaussian description is found from that of the two constituent lenses and their separations. The net magnification is the product of the two contributions, i.e., $m = m_1 \times m_2$. Matrix methods are particularly convenient for handling such combinations, and the results below can be demonstrated easily thereby. If two rotationally symmetric lenses are combined so their axes do not coincide, the combination can be handled with appropriate coordinate translations and rotations in the intermediate space, or by means of collineation. In the most general case, where subsystems without rotational symmetry are combined, the general machinery of collineation can be applied. There are three classes of combinations: focal-focal, focal-afocal, and afocal-afocal.

Focal-Focal Combination: Coaxial

The first lens has power ϕ_1 and principal points at P_1 and P'_1 , Fig. 20. The index preceding the lens is n and that following it is n_{12} . The second lens has power ϕ_2 and principal points at P_2 and P'_2 , with preceding index n_{12} and following index n' . The directed distance from the rear principal point of the first lens to the first principal point of the second lens is $d = P'_1P_2$, which may be positive or negative, since the lenses may have external principal planes. The power of the combination is

$$\phi = \phi_1 + \phi_2 - \frac{1}{n_{12}} d \phi_1 \phi_2 \quad (239)$$

The two principal planes of the combination are located relative to those of the contributing lenses by directed distances

$$P_1P = +\frac{n}{n_{12}} d \frac{\phi_2}{\phi} \quad P'_2P' = -\frac{n'}{n_{12}} d \frac{\phi_1}{\phi} \quad (240)$$

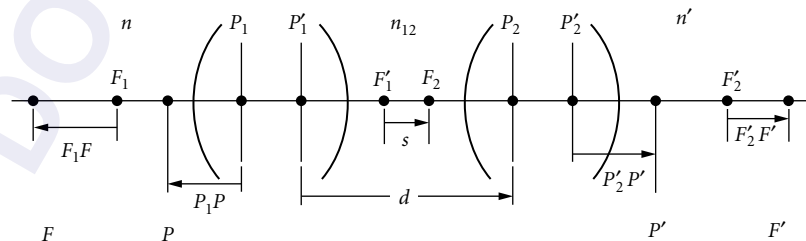


FIGURE 20 Coaxial combination of two focal lenses. The cardinal points of the two lenses are shown above the axis and those of the system below. The directions in this drawing are only one possible case.

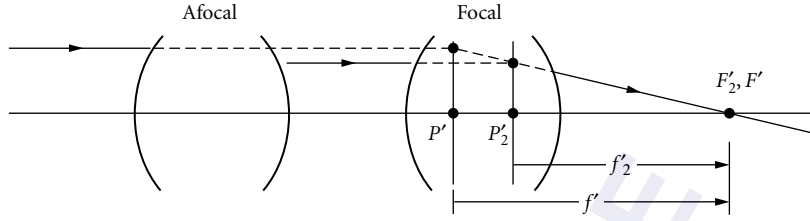


FIGURE 21 Coaxial combination of a focal lens and an afocal lens. In this drawing the afocal lens has a transverse magnification $0 < m_1 < 1$ and the focal lens has a positive power. The combination is a focal lens with focal length $f' = f'_2/m_1$. The focal point on the side of the focal lens is at the focal point of that lens alone.

If $\phi = 0$, the combination is afocal and there are no principal planes. In applying these equations, the inner-space index n_{12} must be the same as that for which the two lenses are characterized. For example, if two thick lenses are characterized in air and combined with water between them, these equations cannot be used by simply changing n_{12} . It would be necessary to characterize the first lens with water following it and the second lens with water preceding it.

Another set of equations involves the directed distance from the rear focal point of the first lens to the front focal point of the second, $s = F_1'F_2$. The power and focal lengths of the combination are

$$\phi = -\frac{1}{n_{12}}s\phi_1\phi_2, \quad f = +\frac{f_1f_2}{s}, \quad f' = -\frac{f_1'f_2'}{s} \quad (241)$$

The focal points are located with respect to those of the contributing lenses by

$$F_1F = +\frac{nn_{12}}{s\phi_1^2} = \frac{n_{12}}{n}\frac{f_1^2}{s}, \quad F_2'F' = -\frac{n'n_{12}}{s\phi_2^2} = -\frac{n'_{12}}{n'}\frac{f_2'^2}{s} \quad (242)$$

Another relationship is $(F_1F)(F_2'F') = ff'$. The system is afocal if $s = 0$. There are many special cases of such combinations. Another case is that when the first principal point of the second lens is at the rear focal point of the first, in which case the system focal length is that of the first. These relationships are proven by Welford.²⁰⁴

Focal-Afocal: Coaxial

A focal lens combined with an afocal lens is focal, Fig. 21. Here we take the afocal lens to be to the left, with magnification m_1 . The focal lens to the right has power ϕ_2 and rear focal length f_2' . The power of the combination is ϕ_2m_1 , and the rear focal length of the combination is $f' = f_2'/m_1$. On the side of the focal lens, the location of the principal focal point is unchanged. On the side of the afocal lens, the system focal point is located at the image of the focal point of the focal lens in the space between the two. Changing the separation between the lenses does not change the power or the position of the principal focal point relative to that of the focal lens. The principal focal point on the afocal lens side does move.

Afocal-Afocal: Coaxial

The combination of two afocal lenses is itself afocal, Fig. 22. If the two lenses have transverse magnifications m_1 and m_2 , the combination has $m = m_1m_2$. A pair of conjugate reference positions is found from the conjugates in the outer regions to any axial point in the inner space. If the separation

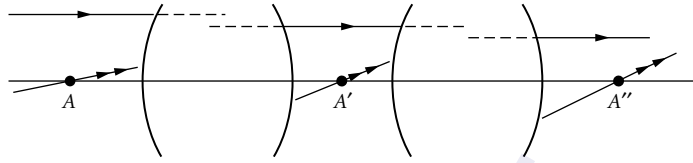


FIGURE 22 Coaxial combination of two afocal lenses. An internal point A' has an object space conjugate A and an image space conjugate A'' . These two points can be used for position references in the overall object and image spaces.

between the two lenses changes, the combination remains afocal and the magnification is fixed, but the conjugate positions change. This result extends to a combination of any number of afocal lenses.

Noncoaxial Combinations: General

The most general combinations can be handled by the machinery of collineation. The net collineation can be found by multiplying the matrices that describe the constituents, with additional rotation and translation matrices to account for their relative positions. After obtaining the overall matrix, object and image space coordinate systems can be found in which the transformation is simplest. This approach can also be used to demonstrate general properties of system combinations. For example, by multiplying matrices for afocal systems, it is seen that a succession of afocal lenses with any orientation is afocal.

1.17 PARAXIAL MATRIX METHODS

Introduction

Matrix methods provide a simple way of representing and calculating the paraxial properties of lenses and their actions on rays. These methods contain no physics beyond that contained in the paraxial power and transfer equations, Eq. (136) and Eq. (142), but they permit many useful results to be derived mechanically, and are especially useful for lens combinations. The matrix description of systems is also useful in elucidating fundamental paraxial properties. With the symbolic manipulation programs now available, matrix methods also provide a means of obtaining useful expressions.

The optical system is treated as a black box represented by a matrix. The axial positions of the input and output planes are arbitrary. The matrix describes the relationship between what enters and what leaves, but contains no information about the specifics of the system within, and there is an infinity of systems with the same matrix representation.

The origin of matrix methods in optics is not clear. Matrices were used by Samson²⁰⁵ who referred to them as “schemes.” Matrices appear without comment in a 1908 book.²⁰⁶ Matrix methods are treated in papers^{207,208} and in many books.^{209–218} Notation is not standardized, and many treatments are complicated by notation that conceals the basic structures.

This section is limited to rotationally symmetric lenses with homogeneous media. References are provided for systems with cylindrical elements. This treatment is monochromatic, with the wavelength dependence of index not made explicit.

The matrices are simplified by using *reduced axial distances* $\tau = t/n$ and *reduced angles* $\omega = nu$. The paraxial angles u are equivalent to direction cosines, and the reduced angles are optical direction cosines in the paraxial limit. For brevity, ω and τ are usually referred to in this section simply as “angle” and “distance.”

Basic Idea: Linearity

Paraxial optics is concerned with the paraxial heights and paraxial angles of rays. A meridional ray entering a system has a given height y and angle ω and leaves with another height y' and angle ω' . Paraxial optics is linear, as discussed above, in the sense that both the outgoing height and angle depend linearly on the incoming height and angle. Writing Eq. (148) in terms of ω 's gives

$$y' = \left(\frac{\partial y'}{\partial y} \right) y + \left(\frac{\partial y'}{\partial \omega} \right) \omega \quad \text{and} \quad \omega' = \left(\frac{\partial \omega'}{\partial y} \right) y + \left(\frac{\partial \omega'}{\partial \omega} \right) \omega \quad (243)$$

The partial derivatives are constant for a given system. This linearity is the basis of the matrix treatment, since these equations can be written in matrix form:

$$\begin{pmatrix} y' \\ \omega' \end{pmatrix} = \begin{pmatrix} \frac{\partial y'}{\partial y} & \frac{\partial y'}{\partial \omega} \\ \frac{\partial \omega'}{\partial y} & \frac{\partial \omega'}{\partial \omega} \end{pmatrix} \begin{pmatrix} y \\ \omega \end{pmatrix} \quad (244)$$

Basic Operations

The basic operations in paraxial ray tracing are transfer, Eq. (136), between surfaces and refraction or reflection at surfaces, Eq. (142).

Transfer Matrix

Transfer changes the height of a ray, in general, leaving the angle unchanged. In terms of reduced quantities, the relationships are

$$y' = y + tu = y + \frac{t}{n} un = y + \tau \omega \quad \text{and} \quad \omega' = \omega \quad (245)$$

The transfer matrix is

$$\begin{pmatrix} 1 & \tau \\ 0 & 1 \end{pmatrix} \quad (246)$$

For left-to-right transfer, $\tau > 0$. This gives a difference in signs between some of the terms in expressions here and those in the gaussian section, where directed distances are measured from a reference point related to the lens to the object.

Power Matrix

Refraction or reflection changes the angle of a ray, but not its height. The equations for reduced quantities are

$$n' u' = nu - y\phi = \omega' = \omega - y\phi \quad \text{and} \quad y' = y \quad (247)$$

Here $\phi = c(n' - n)$ for refraction and $\phi = -2nc$ for reflection, where c is the surface curvature, Eq. (143). The power matrix is

$$\begin{pmatrix} 1 & 0 \\ -\phi & 1 \end{pmatrix} \quad (248)$$

A planar reflecting or refracting surface has $\phi = 0$, so it is represented by the unit matrix.

Arbitrary System

A general system consists of a series of surfaces with powers ϕ_1, ϕ_2, \dots that are separated from one another by distances τ_1, τ_2, \dots . Its matrix is the product

$$\begin{pmatrix} 1 & \tau_N \\ 0 & 1 \end{pmatrix} \dots \begin{pmatrix} 1 & 0 \\ -\phi_2 & 1 \end{pmatrix} \begin{pmatrix} 1 & \tau_2 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -\phi_1 & 1 \end{pmatrix} \begin{pmatrix} 1 & \tau_1 \\ 0 & 1 \end{pmatrix} \quad (249)$$

By convention, the successive matrices are concatenated from right to left, whereas ray tracing is done left to right.

A special case is a succession of transfers, itself a transfer.

$$\text{Succession of transfers: } \begin{pmatrix} 1 & \tau_1 + \tau_2 + \dots \\ 0 & 1 \end{pmatrix} \quad (250)$$

Another is a series of refractions with no intervening transfer, itself a power operation.

$$\text{Succession of powers: } \begin{pmatrix} 1 & 0 \\ -(\phi_1 + \phi_2 + \dots) & 1 \end{pmatrix} \quad (251)$$

Matrix Elements

Each matrix element has a physical significance, and the terms can be given mnemonic symbols associated with the conditions under which they are zero. (This practice is not standard.) If the initial ray angle is zero, the outgoing angles depend on the incident ray heights and the power of the system, according to $\omega' = -\phi y$, so $\partial\omega'/\partial y = -\phi$. If the initial surface is at the front focal plane, the outgoing ray angles depend only on the incident height, so $\partial\omega'/\partial\omega = 0$. This term is denoted by F for "front." Similarly, if the final surface is at the real focal plane, the outgoing ray heights depend only on the incoming angles, so $\partial y'/\partial y = R$ for "rear." If the initial and final planes are conjugate, then all incoming rays at a given height y have the outgoing height $y' = my$, regardless of their angle, so $\partial y'/\partial\omega = 0$ for conjugate planes. Since this term is related to the condition of conjugacy, $\partial y'/\partial\omega = C$ for "conjugate." With this notation, the general matrix is

$$\begin{pmatrix} R & C \\ -\phi & F \end{pmatrix} \quad (252)$$

Dimensions

The terms R and F are dimensionless. C has the dimensions of length, and those of ϕ are inverse length. Dimensional analysis, as well as the consideration of Eq. (248), shows that the F and R terms will always contain products of equal numbers of ϕ_i 's and τ_k 's, such as $\phi_k \tau_l$. The ϕ expression contains terms like ϕ_k and $\tau_k \phi_l \phi_m$, with one more power term than distance terms. Similarly, C has terms like τ_k and $\tau_k \tau_l \phi_m$.

Determinant

Both the transfer and power matrices have unit determinants. Therefore, any product of such matrices has a unit determinant, a fact that is related to the two-ray paraxial invariant.

$$\begin{vmatrix} R & C \\ -\phi & F \end{vmatrix} = FR + C\phi = 1 \quad (253)$$

This provides an algebraic check. For afocal lenses and conjugate arrangements, $FR = 1$.

Possible Zeros

The possible arrangements of zeros in a system matrix is limited by the unit determinant restriction. There can be a single zero anywhere. In this case, either $C = 1/\phi$ or $F = 1/R$, and the remaining nonzero term can have any value. There can be two zeros on either diagonal. No row or column can contain two zeros, since a system represented by such a matrix would violate conservation of brightness. A matrix with double zeros in the bottom row would collimate all rays, regardless of their incoming position and direction. A matrix with all zeros in the top row represents a system that would bring all incoming light to a single point. A system whose matrix has double zeros in the first column would bring all incoming light to a focus on the axis. For double zeros in the second row, the system would concentrate all light diverging from an input point in a single output point with a single direction.

Operation on Two Rays

Instead of considering a single input and output ray, the matrix formalism can be used to treat a pair of rays, represented by a 2×2 matrix. In this case

$$\begin{pmatrix} y'_1 & y'_2 \\ \omega'_1 & \omega'_2 \end{pmatrix} = \begin{pmatrix} R & C \\ -\phi & F \end{pmatrix} \begin{pmatrix} y_1 & y_2 \\ \omega_1 & \omega_2 \end{pmatrix} \quad (254)$$

Since the system matrix has a unit determinant, the determinants of the incoming and outgoing ray matrices are identical:

$$L_{12} = y'_1 \omega'_2 - y'_2 \omega'_1 = y_1 \omega_2 - y_2 \omega_1 \quad (255)$$

This is the paraxial invariant, Eq. (149). It is possible to operate on more than two rays, but never necessary, since any third ray is a linear combination of two, Eq. (154). Operations on two rays can also be handled with a complex notation in which two ray heights and two angles are each represented by a complex number.^{219,220}

Conjugate Matrix

For conjugate planes, $y' = my$, so $C = 0$, $R = m$, and $F = 1/m$, giving

$$\begin{pmatrix} m & 0 \\ -\phi & 1/m \end{pmatrix} \quad (256)$$

The $1/m$ term gives the angular magnification, $u'/u = n/n'm$, Eq. (196). This matrix also holds for afocal lenses, in which case $\phi = 0$.

Translated Input and Output Planes

For a given system, the locations of the input and output planes are arbitrary. If the input plane is translated by τ and the output plane by τ' , the resultant matrix is

$$\begin{pmatrix} R - \tau'\phi & C + \tau R + \tau'F - \tau\tau'\phi \\ -\phi & F - \tau\phi \end{pmatrix} \quad (257)$$

Note that the object-space translation term τ is grouped with F and the image-space term τ' with R . The equation $C = 0 = \tau R - \tau'F - \tau\tau'\phi$ gives all pairs of τ and τ' for which the input and output surfaces are conjugate.

Principal Plane-to-Principal Plane

If the input and output planes are the principal planes, then the matrix is a conjugate one, for which $m = +1$.

$$\begin{pmatrix} 1 & 0 \\ -\phi & 1 \end{pmatrix} \quad (258)$$

This is also the matrix representing a thin lens.

Nodal Plane-to-Nodal Plane

The nodal points are conjugate, with unit angular magnification, so $u' = u$ and $\omega' = n' \omega/n$. Thus

$$\begin{pmatrix} n/n' & 0 \\ -\phi & n'/n \end{pmatrix} \quad (259)$$

The transverse magnification $m_N = n/n'$ equals unity when $n = n'$. This matrix has no meaning for afocal lenses.

Focal Plane-to-Focal Plane

If the initial surface is at the front principal focal plane and the final surface is at the rear focal plane, the matrix is

$$\begin{pmatrix} 0 & 1/\phi \\ -\phi & 0 \end{pmatrix} \quad (260)$$

This is the “Fourier transform” arrangement, in which incident heights are mapped as angles and vice versa.

Translation from Conjugate Positions

If the input plane is translated τ from a plane associated with magnification m and the output plane is translated a distance τ' from the conjugate plane, the matrix is

$$\begin{pmatrix} m - \tau'\phi & m\tau + \tau'/m - \tau\tau'\phi \\ -\phi & 1/m - \tau\phi \end{pmatrix} \quad (261)$$

Setting $C = 0$ gives an equation that locates all other pairs of conjugate planes relative to the first one, Eq. (172).

Translation from Principal Planes

If the initial conjugate planes are the principal planes, then

$$\begin{pmatrix} 1 - \tau'\phi & \tau + \tau' - \tau\tau'\phi \\ -\phi & 1 - \tau\phi \end{pmatrix} \quad (262)$$

The equation for other conjugates is $C = 0 = \tau + \tau' - \tau\tau'\phi$, corresponding to Eq. (170). It follows that the distance from the input surface to the first principal plane is $\tau = (1 - F)/\phi$ and the distance from the output surface to the second principal plane is $\tau' = (1 - R)/\phi$.

Translation from Focal Planes

If the input plane is a distance τ from the front focal plane and the output plane a distance τ' from the rear focal plane, the matrix is

$$\begin{pmatrix} -\phi\tau' & \frac{1}{\phi}(1-\phi^2\tau\tau') \\ -\phi & -\phi\tau \end{pmatrix} \quad (263)$$

Thus F and R are proportional to the distances of the input and output surfaces from the object space and image space focal planes. Using Newton's formulas, this can also be written

$$\begin{pmatrix} m' & \frac{1}{\phi}\left(1-\frac{m'}{m}\right) \\ -\phi & \frac{1}{m} \end{pmatrix} \quad (264)$$

Here m' is the magnification that would obtain if the image point were as located by R , and m is that if the object point were located by F . The conjugate term vanishes when $m = m'$.

Conjugate Relative to Principal Focal Planes

If Eq. (263) is a conjugate matrix, it becomes

$$\begin{pmatrix} -\phi\tau' & 0 \\ -\phi & -\phi\tau \end{pmatrix} \quad (265)$$

The vanishing C term gives $0 = 1/\phi - \phi\tau\tau'$, which is the Newton equation usually written as $zz' = ff'$. The magnification terms are the other Newton's equations, $m = -\phi\tau'$ and $1/m = -\phi\tau$, which are usually written as $m = -z'/f' = -f/z$.

Afocal Lens

For afocal lenses $\phi = 0$. Since the determinant is unity, $F = 1/R$. And since the transverse magnification is constant, $R = m$, giving

$$\begin{pmatrix} m & C \\ 0 & 1/m \end{pmatrix} \quad (266)$$

A ray with $\omega = 0$ has $y' = my$, and $\omega' = \omega/m$ for all y . At conjugate positions, an afocal lens has the matrix

$$\begin{pmatrix} m & 0 \\ 0 & 1/m \end{pmatrix} \quad (267)$$

Performing a translation in both object and images spaces from the conjugate position gives

$$\begin{pmatrix} m & m\tau + \frac{\tau'}{m} \\ 0 & 1/m \end{pmatrix} \quad (268)$$

Setting $C = 0$ gives $\tau' = -m^2\tau$, which relates the location of a single conjugate pair to all others, Eq. (200).

Symmetrical Lenses

For lenses with symmetry about a central plane and symmetrically located input and output surfaces, $F = R$, so the matrix has the form

$$\begin{pmatrix} B & C \\ -\phi & B \end{pmatrix} \quad (269)$$

where $B^2 = 1 - \phi C$. The conjugate matrix has $m = \pm 1$.

Reversing Lenses

When a lens is flipped left to right along with their media, the matrix of the reversed system is obtained from that of the original one by switching the F and R terms.

$$\begin{pmatrix} F & C \\ -\phi & R \end{pmatrix} \quad (270)$$

This reversal maintains the exterior reference planes, that is, the input surface for the initial system becomes the output surface for the flipped one and vice versa.

Inverse Systems

By the “inverse” of a lens is meant a second system that undoes the effect of a given one. That is, the rays at the output surface of the second system have the same height and angle as those at the input of the first system. The combination of a system and its inverse is afocal with unit magnification. The matrix representing the inverse system is the inverse of that representing the system.

$$\begin{pmatrix} F & -C \\ \phi & R \end{pmatrix} \quad (271)$$

The matrix provides no instruction as to how such a lens is made up. Alternatively, the inverse matrix can be interpreted as that whose input is y' and ω' , with outputs y and ω .

Series of Arbitrary Lenses

The matrix for two successive lenses is

$$\begin{pmatrix} R_1 R_2 - C_2 \phi_1 & C_1 R_2 + C_2 F_1 \\ -\phi_1 F_2 - \phi_2 R_1 & F_1 F_2 - C_1 \phi_2 \end{pmatrix} = \begin{pmatrix} R_2 & C_2 \\ -\phi_2 & F_2 \end{pmatrix} \begin{pmatrix} R_1 & C_1 \\ -\phi_1 & F_1 \end{pmatrix} \quad (272)$$

For example, two given lenses separated by some distance have the matrix

$$\begin{pmatrix} R_2 & C_2 \\ -\phi_2 & F_2 \end{pmatrix} \begin{pmatrix} 1 & \tau \\ 0 & 1 \end{pmatrix} \begin{pmatrix} R_1 & C_1 \\ -\phi_1 & F_1 \end{pmatrix} \quad (273)$$

Multiplying from right to left gives a running product or “cumulative matrix,” that shows the effect of the system up to a given plane.

Decomposition

Matrix multiplication is associative, so the system representation can be broken up in a number of ways. For example, the portion of a lens before and after the aperture stop can be used to find the

pupil locations and magnifications. An arbitrary lens matrix can be written as a product of three matrices:²²¹

$$\begin{pmatrix} R & C \\ -\phi & F \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ -\phi/R & 1 \end{pmatrix} \begin{pmatrix} R & 0 \\ 0 & 1/R \end{pmatrix} \begin{pmatrix} 1 & C/R \\ 0 & 1 \end{pmatrix} \quad (274)$$

or

$$\begin{pmatrix} R & C \\ -\phi & F \end{pmatrix} = \begin{pmatrix} 1 & C/F \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1/F & 0 \\ 0 & F \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -\phi/F & 1 \end{pmatrix} \quad (275)$$

Thus a general lens is equivalent to a succession of three systems. One has power and works at unit magnification. The second is a conjugate afocal matrix. The third is a translation. Each of these systems is defined by one of the three terms, either, R , ϕ/R , C/R or F , ϕ/F , C/F . This is another manifestation of the three degrees of freedom of paraxial systems.

Matrix Determination by Two-Ray Specification

If a two-ray input matrix is given along with the desired output, or the two input and output rays are measured to determine the matrix of an unknown lens, Eq. (254) gives

$$\begin{pmatrix} R & C \\ -\phi & F \end{pmatrix} = \begin{pmatrix} y'_1 & y'_2 \\ \omega'_1 & \omega'_2 \end{pmatrix} \begin{pmatrix} y_1 & y_2 \\ \omega_1 & \omega_2 \end{pmatrix}^{-1} \quad (276)$$

so

$$\begin{pmatrix} R & C \\ -\phi & F \end{pmatrix} = \frac{1}{y_1\omega_2 - y_2\omega_1} \begin{pmatrix} y'_1\omega_2 - y'_2\omega_1 & y'_2y_1 - y'_1y_2 \\ \omega'_1\omega_2 - \omega'_2\omega_1 & \omega'_2y_1 - y_2\omega'_1 \end{pmatrix} \quad (277)$$

The denominator of the multiplicative factor is the paraxial invariant associated with the two rays, Eq. (149). As a special case, the two rays could be the marginal and chief rays. The input and output pairs must have the same invariant, or the matrix thus found will not have a unit determinant.

Experimental Determination of Matrix Elements

The matrix elements for an unknown lens can, in principle, be determined experimentally. One method, as mentioned in the preceding section, is to measure the heights and angles of an arbitrary pair of rays. Another method is as follows. The power term is found in the usual way by sending a ray into the lens parallel to the axis and measuring its outgoing angle. To find $C = \partial y'/\partial \omega$, the input ray angle is varied, while its height is unchanged. If the output height is graphed, its slope is C . Likewise, the other partial derivatives in Eq. (243) can be found by changing one of the input parameters while the other is fixed. The four measurements are redundant, the unit determinant providing a check of consistency.

Angle Instead of Reduced Angle

The matrices above can be modified to use the angles u and u' , instead of the reduced angles. In terms of matrix theory, this amounts to a change in basis vectors, which is accomplished by multiplying by diagonal vectors with elements 1 and n or 1 and n' . The result is

$$\begin{pmatrix} y \\ u' \end{pmatrix} = \begin{pmatrix} R & nC \\ -\frac{1}{n'}\phi & \frac{n}{n'}F \end{pmatrix} \begin{pmatrix} y \\ u \end{pmatrix} \quad (278)$$

This matrix has a constant determinant n/n' . The form of Eq. (252) is simpler.

Other Input-Output Combinations

Referring to Eq. (244), any pair of the four quantities y , ω , y' , and ω' can be taken as inputs, with the other two as outputs, and the relationships can be expressed in matrix form. The four matrices in this section cannot be multiplied to account for the concatenation of lenses. If the angles are given, the heights are

$$\begin{pmatrix} y \\ y' \end{pmatrix} = \frac{1}{\phi} \begin{pmatrix} F & -1 \\ 1 & -R \end{pmatrix} \begin{pmatrix} \omega \\ \omega' \end{pmatrix} \quad (279)$$

The matrix is undefined for afocal lenses, for which the relationship of ω and ω' is independent of heights. Similarly, the angles can be expressed as functions of the heights by

$$\begin{pmatrix} \omega \\ \omega' \end{pmatrix} = \frac{1}{C} \begin{pmatrix} -R & 1 \\ -1 & F \end{pmatrix} \begin{pmatrix} y \\ y' \end{pmatrix} \quad (280)$$

For conjugates the expression breaks down, since there is no fixed relationship between heights and angles. If the input is a height on one side and an angle on the other, then

$$\begin{pmatrix} y' \\ \omega \end{pmatrix} = \frac{1}{F} \begin{pmatrix} 1 & C \\ \phi & 1 \end{pmatrix} \begin{pmatrix} y \\ \omega' \end{pmatrix} \quad (281)$$

For the inverse situation,

$$\begin{pmatrix} y \\ \omega' \end{pmatrix} = \frac{1}{R} \begin{pmatrix} 1 & -C \\ -\phi & 1 \end{pmatrix} \begin{pmatrix} y' \\ \omega \end{pmatrix} \quad (282)$$

The determinants of these matrices are, respectively, C , ϕ , R , and F .

Derivative Matrices

If the axial position of the input surface changes, the rate of change of the output quantities is

$$\begin{pmatrix} dy'/dz \\ d\omega'/dz \end{pmatrix} = \begin{pmatrix} 0 & R \\ 0 & -\phi \end{pmatrix} \begin{pmatrix} y \\ \omega \end{pmatrix} \quad (283)$$

If the axial position of the output surface can change, the rate of change of output quantities is

$$\begin{pmatrix} dy'/dz' \\ d\omega'/dz' \end{pmatrix} = \begin{pmatrix} -\phi & F \\ 0 & 0 \end{pmatrix} \begin{pmatrix} y \\ \omega \end{pmatrix} \quad (284)$$

Higher derivatives vanish.

Skew rays

The matrix formalism can be used to treat a paraxial skew ray, represented by a 2×2 matrix of x and y positions and directions α and β . In this case

$$\begin{pmatrix} x' & y' \\ n'\alpha' & n'\beta' \end{pmatrix} = \begin{pmatrix} R & C \\ -\phi & F \end{pmatrix} \begin{pmatrix} x & y \\ n\alpha & n\beta \end{pmatrix} \quad (285)$$

Since the lens matrix has a unit determinant, the determinants of the incoming and outgoing ray matrices are identical:

$$n'(y'\alpha' - x'\beta') = n(y\alpha - x\beta) \quad (286)$$

From Eq. (73), this is the skew invariant.

Relationship to Characteristic Functions

A lens matrix can be related to any one of the four paraxial characteristic functions, Eqs. (34) through (37), each of which has three first coefficients, associated with the three degrees of freedom of the matrix. Brouwer and Walther²²² derive the paraxial matrices from more general matrices based on the point angle characteristic function.

Nonrotationally Symmetric Systems

Systems comprising cylindrical lenses can also be treated paraxially by matrices.^{223–228,221} The more general case of a treatment around an arbitrary ray is also represented by a 4×4 matrix.²²⁹ This is treated by several of the references to the section “Images About Known Rays.”

1.18 APERTURES, PUPILS, STOPS, FIELDS, AND RELATED MATTERS

Introduction

This section is concerned with the finite size of lenses and their fields, as expressed in various limitations of linear dimensions and angles, and with some of the consequences of these limits. (Other consequences, for example, resolution limitations, are in the domain of wave optics.) Terminology in this area is not well defined, and the terms typically used are insufficient for all the aspects of the subject, so this section deals considerably with definitions.

Field Size and Field Stop

The *field* or *field of view* of a lens is the region of object space from which light is captured or the region of image space that is used. The field size may be described in angular, linear, or area units, depending on the circumstances. (It can be described in still other ways, e.g., the number of pixels.) In and of itself, a lens does not have a definite field size, but beyond a certain size, image quality diminishes, both with respect to aberration correction and to light collection. A *field stop* is a physical delimiter of the field, which may be in either object or image space. A detector may be the delimiter.

Aperture Stop

Each object point can be thought of as emitting rays in all directions. Since lenses are finite in size, only some of the rays pass through them. The rays that do pass are referred to as *image-forming rays*, the ensemble of which is the *image-forming bundle*, also called the *image-forming cone*, although the bundle may not be conical. The bundle associated with each object point is delimited by one or more physical structures of the lens. For axial object points, the delimiting structure is called the *aperture*, the *stop*, or the *aperture stop*. The aperture may be either within the lens or outside of it on either side, Fig. 23. The aperture may be a structure whose sole purpose is delimiting the bundle, or it may be the edge of an optical element or a lens mount. The aperture stop may be fixed or adjustable, for instance, an iris. Which structure acts as the aperture can change with object position, Fig. 24. The size and position of the aperture do not effect the gaussian properties of the lens, i.e., the cardinal points and the conjugate locations and magnifications. They do affect the image irradiance, the aberrations, and the effects of defocus. The aperture is most commonly centered on axis, but this is not always so. With visual instruments, the aperture stop for the entire system may be either an aperture in the optics or the iris of the observer's eye.

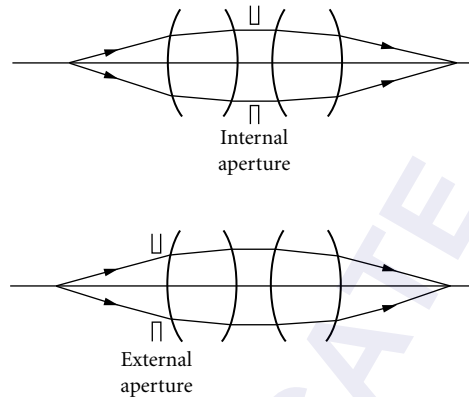


FIGURE 23 Axial ray cone and aperture stop. The upper lens has an internal aperture, and the lower one has an external aperture on the object side.

Marginal Rays and Chief Rays

Ray bundles are described to a considerable extent by specifying their central and extreme rays. For object planes perpendicular to the lens axis, there are two meridional rays of particular importance, defining the extremities of field and aperture, Fig. 25. These rays are reciprocal in that one is to the pupil what the other is to the field.

The *marginal ray* originates at the axial object point, intersects the conjugate image point, and passes through the edge of the aperture. This term is also used for rays from other field points that pass through the extremes of the aperture. The *paraxial marginal ray* is the marginal ray in the paraxial limit.

The *chief ray* or *principal ray* originates at the edge of the object field, intersects the edge of the image field, and passes approximately through the center of the aperture, and hence approximately through the center of the pupils. (Here we use “chief ray,” since the prefix “principal” is so commonly used for other entities.) The term is also used for the central ray of other bundles. The *paraxial chief ray* passes exactly through the centers of the aperture and both paraxial pupils.

Field Angle

The *field angle* is that subtended by the field of view at the lens. This term is ambiguous, since several angles can be used, as well as angles in both object and image space. A nodal ray angle is the same in both spaces. If the nodal points are not at the pupils, the chief ray angle differs on the two sides. The ratio

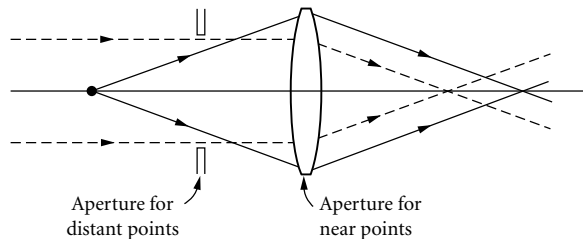


FIGURE 24 An example of change of aperture with axial object position. For distant points the aperture is the nominal stop. For near points the aperture is the rim of the lens.

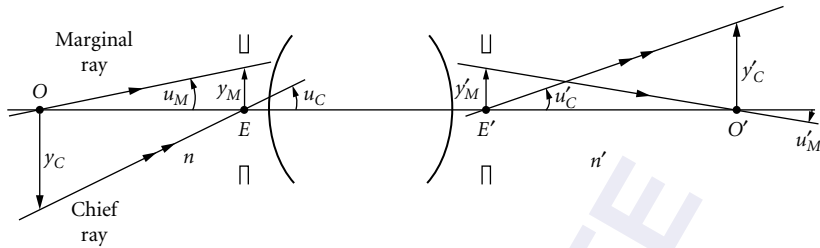


FIGURE 25 Schematic diagram of a lens with object and image planes, entrance and exit pupils, and marginal and chief rays. The entrance pupil is located at E and the exit pupil at E' . The chief ray passes through the edges of the fields and the centers of the pupils. The marginal ray passes through the axial object and image points and the edges of the pupils.

of paraxial chief ray angles is proportional to the paraxial pupil magnification, as discussed later, Eq. (289). If the lens is telecentric, the chief ray angle is zero. An afocal lens has no nodal points, and the paraxial ratio of output angles to input angles is constant. The concept of field angle is most useful with objects and/or images at large distances, in which case on the long conjugate side the various ray angles are nearly identical. On the short conjugate side, ambiguity is removed by giving the focal length, the linear size of the detector, and the principal plane and exit pupil positions. For finite conjugates, such information should be provided for both spaces.

Pupils

The term *pupil* is used in several ways, and care should be taken to distinguish between them. There are paraxial pupils, “real” pupils, pupils defined as ranges of angles, and pupil reference spheres used for aberration definition and diffraction calculations. The *entrance pupil* is the aperture as seen from object space—more precisely, as seen from a particular point in object space. If the aperture is physically located in object space, the entrance pupil is identical to the aperture. Otherwise, the entrance pupil is the image of the aperture in object space formed by the portion of the lens on the object side of the aperture. If the aperture is in image space, the entrance pupil is its image formed by the entire lens. Similarly, the *exit pupil* is the aperture as seen from image space. A *real pupil* is a physically accessible image of the aperture or the aperture itself, and a *virtual pupil* is an inaccessible image. Visual instruments often have external pupils, where the user’s eye is located. The axial entrance pupil point is denoted here by E and the exit pupil by E' .

The pupils can be located anywhere on axis, except that they cannot coincide with the object or image. It is common to draw pupils as shown in Fig. 25, but they can also be on the side of the object or image away from the lens. The pupils are usually centered on axis, but not necessarily. Aberrations may shift pupils from nominal axial centration.

Both pupils are conjugate to the aperture, so they are conjugate to each other. The term *pupil imaging* refers to the relationship of the pupils with respect to each other and to the aperture. In pupil imaging, the chief ray of the lens is the marginal ray and vice versa. The *pupil magnification* m_p denotes the ratio of exit pupil size to entrance pupil size. The size may be specified as linear or an angular extent, and the pupil magnification may be a transverse magnification, finite or paraxial, or a ratio of angular subtenses. In general, there is *pupil aberration*, so the image of the aperture in each space is aberrated, as is that of the imaging of one pupil to the other. Pupil imaging is subject to chromatic aberration, so positions, sizes, and shapes of pupils may vary with wavelength.

There is ambiguity about pupil centers and chief rays for several reasons. The center can be taken with respect to linear, angular, or direction cosine dimensions. Because of spherical pupil aberration, a ray through the center of the pupil may not also pass through the center of the aperture, and vice versa. The angular dimensions of pupils may change with field position. Pupil aberrations cause the actual pupil shape to be different from that of the paraxial pupil.

Pupils that are not apertures can have any linear size, since the aperture can be imaged at any magnification. If the aperture is within the lens, there is no particular relationship between the positions and linear sizes of the entrance and exit pupils, since the portions of the lens that precede and follow the aperture have no specific relationship. There is a relationship between the angular subtense of the pupils, as discussed below.

The angular size and shape of the pupils can vary with field position, and the pupils can change position if the aperture changes with object position. If the lens changes internally, as with a zoom, the sizes and positions of the pupils change.

Paraxial Description

The *paraxial pupils* are the paraxial images of the aperture. They are usually planar and perpendicular to the axis and are implicitly free from aberration. The paraxial chief ray passes through the center of both pupils and the aperture, and the paraxial marginal ray through the edges. The object and pupil magnifications and the distances from object to entrance pupil and from exit pupil to image are related by Eq. (194). If the object at O is imaged at O' with magnification m , and the pupil magnification from entrance pupil at E to exit pupil at E' is m_E , then from Eq. (194)

$$O'E' = \frac{n'}{n} m m_E OE \quad (287)$$

Paraxial Invariant for Full Field and Full Aperture

Let the height of the paraxial marginal ray be y_M at the entrance pupil and y'_M at the exit pupil, and that of the paraxial chief ray by y_C , at the object plane and y'_C at the image plane, Fig. 25. Let the angles of these rays be u_M, u_C, u'_M, u'_C . The two-ray paraxial invariant, Eq. (149), is

$$L = n y_C u_M = n y_M u_C = n' y'_M u'_C = n' y'_C u'_M \quad (288)$$

This relationship was rediscovered several times, so the conserved quantity is referred to by a variety of names, including the *Lagrange invariant*, the *Helmholtz invariant*, the *Smith invariant*, and with various hyphenated combinations of the proper names.^{230, 231} Further discussions are found in the sections on paraxial optics and on the étendue. The paraxial transverse magnification and paraxial pupil magnifications are related to the paraxial marginal and chief ray angles by

$$m = \frac{y'_C}{y_C} = \frac{n u_M}{n' u'_M} \quad \text{and} \quad m_p = \frac{y'_M}{y_M} = \frac{n u_C}{n' u'_C} \quad (289)$$

Pupil Directions

For some purposes, pupils are best described as ranges of directions, specified in direction cosines, rather than by linear extents of aperture images. Here the term *pupil directions* (NS) is used. This is particularly the case when dealing with a given region of the object. The construction for this description is shown in Fig. 26. The x and y axes of the object-space coordinate system lie in the object surface, and the x' and y' axes of the image-space coordinate system lie in the image surface. From a point on the object plane, the extreme set of rays that passes through the lens is found. Its intersection with a unit sphere about the object point is found, and perpendiculars are dropped to the unit circle on (or tangent to) the object plane, giving the extent in direction cosines.

The entrance pupil is delimited by a closed curve described by a relationship $0 = P(\alpha, \beta; x, y)$, and the exit pupil is likewise defined by $0 = P'(\alpha', \beta'; x', y')$. The spatial argument is included to indicate that the shape varies, in general, with the field position. There may be multiple regions, as in the case of central obstructions. It is usually preferable to define the pupils relative to *principal directions* (NS) (α_0, β_0) in object space and (α'_0, β'_0) in image space, where the two directions are

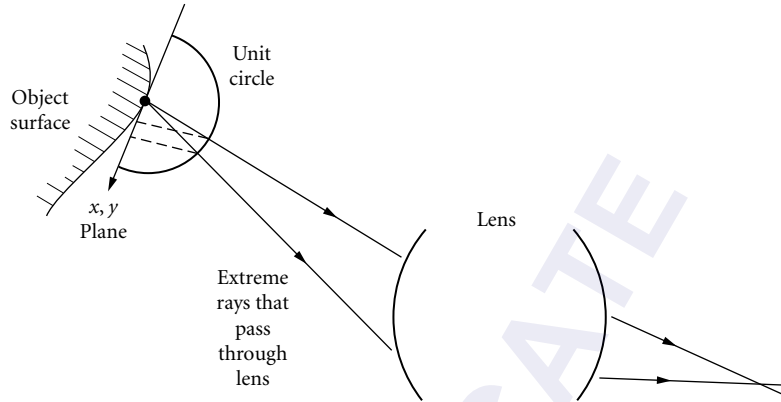


FIGURE 26 Construction for the description of the pupils with direction cosines. An x - y plane is tangent to the object surface at the object point, and a unit sphere is centered on the point. The intersections with the unit sphere of the rays are projected to the tangent plane to give the pupil direction cosines.

those of the same ray in the two spaces, usually a meridional ray. The principal directions are analogous to the chief rays. The entrance pupil is then given by $0 = Q(\alpha - \alpha_0, \beta - \beta_0; x, y)$ and the exit pupil by $0 = Q'(\alpha' - \alpha'_0, \beta' - \beta'_0; x', y')$. For example, for a field point on the $x = 0$ meridian, the expression for the pupil might be well approximated by an ellipse, $0 = a\alpha^2 + b(\beta - \beta_0)^2$, where $(0, \beta_0)$ is the chief ray direction. If the imaging is stigmatic, the relationship between entrance and exit pupil angular shapes is provided by the cosine condition, Eq. (104).

$$Q'(\alpha', \beta'; x', y') = Q(m_p \alpha' - \alpha'_0, m_p \beta' - \beta'_0; x, y) \quad (290)$$

The entrance and exit pupils have the same shapes when described in direction cosine space. They are scaled according to the *pupil angular magnification* (NS) $m_p = n/n'm$. The orientations may be the same or rotated 180° . There is no particular relationship between (α_0, β_0) and (α'_0, β'_0) which can, for example, be changed by field lenses. The principal directions are, however, usually in the same meridian as the object and image points, in which case $\alpha_0/\beta_0 = \alpha'_0/\beta'_0$. If the field point is in the x meridian, and the central ray is in this meridian, then $\alpha_0 = 0$ and $\alpha'_0 = 0$. Even with aberrations, Eq. (290) usually holds to a good approximation. The aberration *pupil distortion* refers to a deviation from this shape constancy.

Pupil Directional Extent: Numerical Aperture and Its Generalizations

The angular extent of a pupil extent is limited by some extreme directions. In the example above of the elliptical shape, for instance, there are two half widths

$$\frac{1}{2}(\alpha_{\max} - \alpha_{\min}) \quad \text{and} \quad \frac{1}{2}(\beta_{\max} - \beta_{\min}) \quad (291)$$

For a rotationally symmetric lens with a circular aperture, the light from an axial object point in a medium of index n is accepted over a cone whose vertex angle is θ_{\max} . The *object space numerical aperture* is defined as

$$NA = n \sin \theta_{\max} = n \sqrt{(\alpha^2 + \beta^2)_{\max}} = n \alpha_{\max} = n \beta_{\max} \quad (292)$$

Likewise, on the image side, where the index is n' and the maximum angle is θ'_{\max} , the *image space numerical aperture* is

$$NA' = n' \sin \theta'_{\max} = n' \sqrt{(\alpha'^2 + \beta'^2)_{\max}} = n' \alpha'_{\max} = n' \beta'_{\max} \quad (293)$$

If the lens is free of coma, the sine condition, Eq. (106), gives for finite conjugates

$$m = \frac{n \sin \theta_{\max}}{n' \sin \theta'_{\max}} = \frac{NA}{NA'} \quad (294)$$

For infinite conjugates

$$\sin \theta'_{\max} = -\frac{y_{\max}}{f'} \quad \text{or} \quad n' \sin \theta'_{\max} = NA' = n \beta_{\max} = -y_{\max} \phi \quad (295)$$

If there is coma, these relationships are still good approximations. For a given lens and a given aperture size, the numerical aperture varies with the axial object position.

F-Number and Its Problems

The *F-number* is written in a variety of ways, including F/no. and F/#. It is denoted here by FN . The F-number is not a natural physical quantity, is not defined and used consistently in the literature, and is often used in ways that are both wrong and confusing.^{232,233} Moreover, there is no need to use the F-number, since everything that it purports to describe or approximately describes is treated properly with direction cosines. The most common definition for F-number, applied to the case of an object at infinity, is

$$FN = \frac{\text{focal length}}{\text{entrance pupil diameter}} = \frac{1}{2 \tan \theta'} \quad (296)$$

where θ' is the outgoing angle of the axial imaging cone. In general, the F-number is associated with the tangents of collinear transformations, rather than the sines (or direction cosines) that are physically appropriate. It presumes that a nonparaxial ray entering parallel to the axis at height y leaves the rear principal plane at the same height and intersects the rear focal point, so that $\tan \theta' = y/f'$. However, this particular presumption contradicts Eq. (294), and in general, collineation does not accurately describe lens behavior, as discussed above.

Other problems with F-number, as it is used in the literature, include the following: (1) It is not defined consistently. For example, the literature also contains the definition F-number = (focal length)/(exit pupil diameter). (2) For lenses used at finite conjugates, the F-number is often stated for an object at infinity. In fact, given only the numerical aperture for an object at infinity, that for other conjugates cannot be determined. (3) There are confusing descriptions of variation of F-number with conjugates, for example, the equation $FN_m = (1 + m)FN_{\infty}$, where FN_m is the F-number for magnification m and FN_{∞} is that for an object at infinity. In fact, numerical apertures for various magnifications are not so related. (4) The object and image space numerical apertures are related by Eq. (294), but there is no such relationship for tangents of angles, except that predicted by collineation, Eq. (232), which is approximate. (5) With off-axis field points and noncircular pupils, the interpretation of F-number is more ambiguous. (6) Afocal systems have finite numerical apertures when used at finite conjugates, but they have no analogue to Eq. (295). (7) Object and image space refractive indices are not accounted for by the F-number, whereas they are by the numerical aperture. (8) The F-number is often used as a descriptor of radiometric throughput, rather than of ray angles per se.

A related quantity is the *T-number*,²³⁴ which accounts for both the convergence angle of the imaging cone and the fraction of power transmitted by the lens. This is useful as a single-number descriptor, but it is subject to all the confusion associated with the F-number.

Image Irradiance for Lambertian Objects

If the light from a region of an object is lambertian with a power/area M , then the emitted power per angle with angle according to $(M/\pi) \cos \theta d\omega = (M/\pi) d\alpha d\beta$. The power captured by the entrance pupil from a small object area dA is

$$dP = \frac{1}{\pi} M dA \int_{\text{entrance pupil}} d\alpha d\beta \quad (297)$$

(For a full hemisphere $\int d\alpha d\beta = \pi$, giving $dP = M dA$.) If there are no losses within the lens, the power reaching the conjugate image region dA' is the same. Using the conservation of étendue equation, Eq. (72), the image irradiance is

$$E = \frac{dP}{dA'} = \frac{1}{\pi} M \frac{n'^2}{n^2} \int_{\text{exit pupil}} d\alpha' d\beta' \quad (298)$$

The image irradiance does not depend explicitly on the magnification, but magnification is included implicitly, since, for a given lens, the subtense of the exit pupil varies with conjugates.

This equation applies everywhere in the field, and it applies to arbitrary object surface positions and orientations, so long as the direction cosines are defined with respect to the local object and image surface normals. These equations apply regardless of the chief ray angles, so they are applicable, for example, with telecentricity. In general, the pupil shape and principal direction vary with field position, so there is a gradation of irradiance in the image of a uniform lambertian object.

These equations do not account for all that influences image irradiance, for example, lens absorption and reflection. These effects can be included in the above expressions by adding an appropriate weighting function of angle and field in the above integrals, giving

$$E(x', y') = \frac{dP}{dA'} = \frac{1}{\pi} M(x, y) \frac{n'^2}{n^2} \int \tau(\alpha', \beta'; x', y') d\alpha' d\beta' \quad (299)$$

where $\tau(\alpha', \beta'; x', y')$ is the lens transmittance as a function of the direction cosines for the image point (x', y') . With externally illuminated objects that are not lambertian scatterers, these relationships do not hold. For example, in optical projectors the illumination is matched to the object and imaging lens to give nominally uniform image irradiance.

Axial Image Irradiance for Lambertian Objects

In the special case of circular pupils and axial object surfaces perpendicular to the axis, the collected power and image irradiance given above are

$$dP = M dA \sin^2 \theta \quad \text{and} \quad E = M \frac{n'^2}{n^2} \sin^2 \theta' \quad (300)$$

Power/Pixel

From wave optics, a lens working at the “resolution limit” has an image pixel size $q\lambda/n' \sin \theta'$, where λ is the vacuum wavelength and q is a dimensionless factor, typically of the order of unity. Applying Eq. (300) gives

$$\text{Power/pixel} = q^2 M \left(\frac{\lambda}{n} \right)^2 \quad (301)$$

$M(\lambda/n)^2$ is the energy emitted per square wavelength of object area. This is a fundamental radiometric quantity. Increasing q gives a greater numerical aperture than is nominally required for resolution, but in practice the aberration correction may be such that the actual resolution is not greater.

Cosine-to-the-Fourth Approximation

For distant, planar, uniform lambertian objects perpendicular to the lens axis, if the entrance pupil is well approximated by a circle, then the image irradiance varies approximately with the object space field angle ψ according to the *cosine-to-the-fourth* relationship

$$E(\psi) = E_0 \cos^4 \psi \quad (302)$$

where E_0 is the axial irradiance. There are three contributions to this dependence. (1) The angular distribution of a lambertian emitter varies as $\cos \psi$. (2) The distance from the field point to the entrance pupil varies as $1/d^2 \propto \cos^2 \psi$. (3) Insofar as the pupil behaves as a rigid circle, its projected solid angle varies approximately as $\cos \psi$. The cosine-to-the-fourth relationship should be used only as a guideline, since ray tracing permits more accurate calculations, and because of the ambiguities in the meaning of the field angle, as discussed above, and elsewhere.^{235–239} For example, field angle is meaningless with telecentricity. Some lenses, especially wide-angle ones, are specifically designed so the pupil subtense increases with the field angle in order to compensate for effects (1) and (2) above, to produce a sufficiently uniform image.²⁴⁰

Total Lens Étendue

The total amount of power from a lambertian object that can be transferred through a lens is

$$\frac{1}{\pi} M \int_{\text{field}} dx dy \int_{\text{pupil}} d\alpha d\beta \quad (303)$$

The pupil integral may vary over the field. If the pupil is round and constant over the field, the étendue is proportional to $A(\text{NA})^2$, where A is the area of the field. This quantity is also related to the total number of pixels in the field, and the ability of the lens to transfer information.²⁴¹ The term “area-solid angle product” is sometimes used, but this is an approximation. The total étendue is proportional paraxially to $\sim L^2$, where L is given by Eq. (288).

Vignetting

Vignetting occurs when an image-forming bundle is truncated by two or more physical structures in different planes, Fig. 27. Typically, one is the nominal aperture and another is the edge of a lens. Another case is that of central obstructions away from the aperture. When vignetting occurs, the image irradiance is changed, and its diminution with field height is faster than it otherwise would be. Aberration properties are also changed, so vignetting is sometimes used to eliminate light that would unacceptably blur the image.

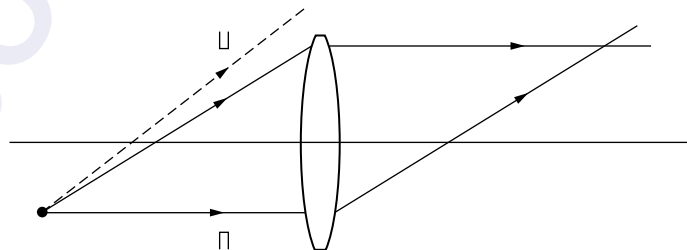


FIGURE 27 Example of vignetting. The dashed ray passes through the aperture, but misses the lens.

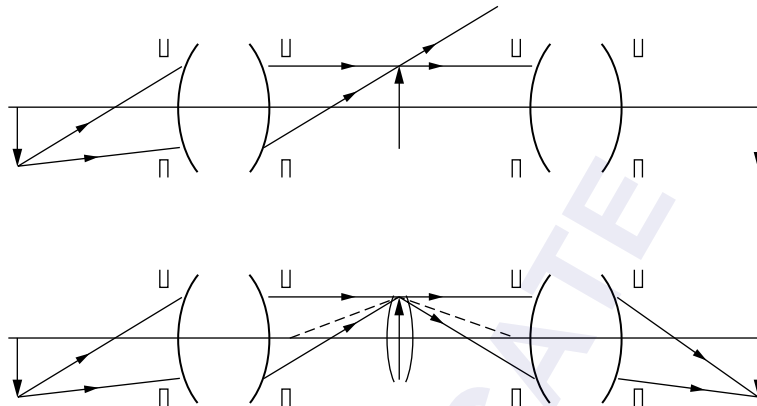


FIGURE 28 A pair of lenses relaying an image with and without a field lens. In the top figure, there is no field lens, and some of the light forming the intermediate image does not pass through the second lens. The amount lost depends on the two numerical apertures and increases with distance from the axis. In the lower figure, a field lens at the intermediate image forms an image of the exit pupil of the first lens into the entrance pupil of the next. No light is lost unless the numerical aperture of the second lens is less than that of the first.

Lens Combinations and Field Lenses

When lenses are used to relay images, the light is transferred without loss only if the exit pupil of one corresponds with the entrance pupil of the next. An example of the failure to meet this requirement is shown in Fig. 28. The axial point is reimaged satisfactorily, but off-axis bundles are vignetted. To transfer the light properly, a *field lens* in the vicinity of the intermediate image is used to image the exit pupil of the preceding lens into the entrance pupil of the next one. If the field lens is a thin lens in the image plane, then its magnification with respect to the image is unity. In practice, the field lens is usually shifted axially, so scratches or dust on its surface are out of focus. Its magnification then differs from unity. The focal length of a thin field lens in air is given by $1/f' = 1/a + 1/b$, where a is the distance from exit pupil of first lens to the field lens, and b is that from field lens to the entrance pupil of the second lens. The exit pupil is reimaged with a magnification b/a . If the sizes of the various pupils and their images are not matched, then the aperture of the combination is determined by the smallest. Field lenses affect aberrations.

Defocus

When the object and image-receiving surface are not conjugate there is *defocus*. If either the object or the receiving surface is considered to be correctly positioned, the defocus is associated with the other. Another situation is that in which the object and receiving surfaces are conjugate, but both are wrongly located, so that the image is sharp but the magnification is not what is desired.

Defocus has two basic geometrical effects, if there are no aberrations, Fig. 29. One is blurring, since the rays from an object point do not converge to a single point on the receiving surface. The blur size varies linearly with the axial defocus in image space and with the cone angle of the image-forming bundle. The shape of the blur is that of the exit pupil, projected on the receiving surface. The other effect of defocus is a lateral shift in position of the blur's centroid relative to that of the correctly focused point. The shift depends on the chief ray angle on the side of the lens where the defocus occurs. In the simplest case, the shift is approximately linear with field height, so acts as a

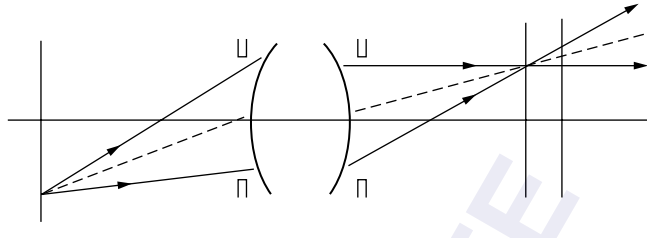


FIGURE 29 Defocus of the receiving surface. A receiving surface is shown in focus and shifted axially. The image of a point on the shifted surface is blurred, and its centroid is translated radially.

change of magnification. If the object is tilted or is not flat, the effects of defocus vary across the field in a more complicated way. Aberrations affect the nature of the blur. With some aberrations, the blur is different on the two sides of focus. With spherical aberration, the blur changes in quality, and with astigmatism the orientation of the blur changes.

In considering the geometrical imaging of a small region of a lambertian object, there is an implicit assumption that the pupil is filled uniformly with light. In imaging an extended object that is externally illuminated, the light from a given region may not fill the pupil uniformly, so the character of the blurring is affected by the angular properties of the illumination and scattering properties of the object.

The amount of defocus can be described in either object or image space, and it can be measured in a variety of ways, for example, axial displacement, displacement along a chief ray, geometrical blur size, and wavefront aberration. The axial displacements in object and image space differ, in general, and are related by the longitudinal magnification. As expressed in wavefront aberration, i.e., optical path length, defocus is the same in both spaces. There are also various functional measurements of defocus, for example, the sizes of recorded images through focus.

Telecentricity

A lens is *telecentric* if the chief rays are parallel to one another. Most commonly, they are also parallel to the lens axis and perpendicular to the object and/or image planes that are perpendicular to the axis, Fig. 30. Telecentricity is often described by speaking of pupils at infinity, but the consideration of ray angles is more concrete and more directly relevant. A lens is *telecentric in object space* if the chief rays in object space are parallel to the axis, $\alpha_0 = 0$ and $\beta_0 = 0$. In this case the image of the aperture formed by the portion of the lens preceding it is at infinity and the aperture is at the rear focal plane of the portion preceding it. Similarly, a lens is *telecentric in image space* if the aperture is at the

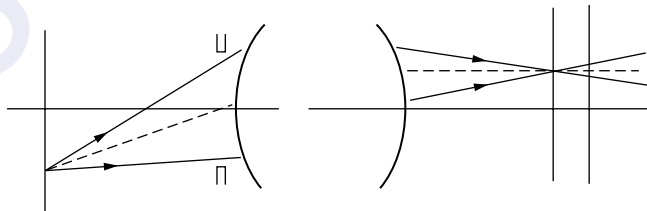


FIGURE 30 Example of telecentricity. The lens shown is telecentric in image space, in which chief rays are parallel to the axis. An axial shift in the receiving surface results in blurring, but does not translate the centroid, so there is no change in image scale.

front focal point of the subsequent optics, so $\alpha'_0 = 0$ and $\beta'_0 = 0$. More generally, but less commonly, the chief rays can be parallel to each other, but not necessarily to the axis, and not necessarily perpendicular to a (possibly tilted) object or image plane.

With tilted object and image surfaces and nonaxial pupils, the chief rays are not perpendicular to the object and/or image surfaces, but their angles are everywhere the same, so defocus can result in a rigid shift of the entire image.

A focal lens can be nontelecentric or telecentric on either side, but it cannot be doubly telecentric. An afocal lens can be nontelecentric, or doubly telecentric, but it cannot be telecentric on one side. A doubly telecentric lens must be afocal, and a singly telecentric lens cannot be afocal.

For a lens that is telecentric in image space, if the receiving surface is defocused, the image of a point is blurred, but its centroid stays fixed. However, if it is not telecentric in object space, then the scale changes if the object is defocused. The converse holds for object-space telecentricity without image-space telecentricity. For a doubly telecentric lens, an axial shift of either the object or the receiving plane produces blurring without a centroid shift. Although the magnification of an afocal lens does not change with conjugates, there can be an effective change with defocus if it is not telecentric. If the pupil is not on the axis or if the object and image planes are tilted, there can be telecentricity without the chief rays being perpendicular to the object and/or image planes. In these cases, defocus results in a rigid shift of the entire image.

Nominal telecentricity can be negated in several ways. Pupil aberrations may change the chief ray angles across the field. For an extended object that is externally illuminated the pupil may not be filled uniformly by light from a given region, so defocus can product a lateral image shift.

Depth of Focus and Depth of Field

The *depth of focus* and *depth of field* are the amounts of defocus that the receiving surface or object may undergo before the recorded image becomes unacceptable. The criterion depends on the application—the nature of the object, the method of image detection, and so on, and there are both ray and wave optics criteria for goodness of focus. For example, a field of separated point objects differs from that of extended objects. Depth of focus is usually discussed in terms of blurring, but there are cases where lateral shifts become unacceptable before blurring. For example, in nature photography blurring is more critical than geometrical deformation, while the opposite may be true in metrology.

Range of Focus and Hyperfocal Distance

In some cases, a geometrical description of defocus is applicable, and the allowable blur is specified as an angle.^{242,243,234} The *hyperfocal distance* is

$$\text{Hyperfocal distance} = \frac{\text{diameter of the entrance pupil}}{\text{maximum acceptable angular blur}} = d_H \quad (304)$$

Let the object distance at which the lens is focused be d , the nearest distance at which the image is acceptable be d_N , and the furthest distance be d_F . All of these quantities are positive definite. The following relations are obtained:

$$d_F = \frac{d_H d}{d_H - d} \quad \text{and} \quad d_N = \frac{d_H d}{d_H + d} \quad (305)$$

The distances to either side of best focus are

$$d_F - d = \frac{d^2}{d_H - d} \quad \text{and} \quad d - d_N = \frac{d^2}{d_H + d} \quad (306)$$

The total range of focus is

$$d_F - d_N = \frac{2d^2 d_H}{d_H^2 - d^2} = \frac{2d}{(d_H/d)^2 - 1} \quad (307)$$

For $d > d_H$ the above quantities involving d_F are infinite (not negative). If the lens is focused at the hyperfocal distance or beyond, then everything more distant is adequately focused. If the lens is focused at the hyperfocal distance, i.e., $d = d_H$, the focus is adequate everywhere beyond half this distance, and this setting gives the greatest total range. If the lens is focused at infinity, then objects beyond hyperfocal distance are adequately focused. The hyperfocal distance decreases as the lens is stopped down.

1.19 GEOMETRICAL ABERRATIONS OF POINT IMAGES: DESCRIPTION

Introduction

In instrumental optics, the term *aberration* refers to a departure from what is desired, whether or not it is physically possible. Terms such as “perfect system” and “ideal system” indicate what the actual is compared to, and these terms themselves are not absolute, but depend on what is wished for. The ideal may be intrinsically impossible, in which case a deviation therefrom is not a defect. A further distinction is between aberrations inherent in a design and those that result from shortcomings in fabrication.

This section considers only the description of aberrations of point images, with the lens treated as a black box, whose action with respect to aberrations is accounted for by what leaves the exit pupil. A full consideration of aberrations involves, among other things, their causes, their correction, their various manifestations, and their evaluation. Aberrated images of extended objects are formed by overlapping blurs from the individual points. The analysis of such images is object- and application-dependent, and is beyond the scope of this section. Aberrations do vary with wavelength, but most of this discussion involves monochromatic aberrations, those at a single wavelength. In addition, aberrations vary with magnification. Aberrations are discussed to some extent in many books that treat geometrical optics.^{244–253}

Descriptions

Aberration has many manifestations, and can be described in a variety of ways. For example, geometrical wavefronts, path lengths, ray angles, and ray intersection points can all differ from the nominal (and in wave optics there are additional manifestations). Terms such as “wavefront aberration” and “ray aberration” do not refer to fundamentally different things, but to different aspects of the same thing. Often, a single manifestation of the aberration is considered, according to what is measurable, what best describes the degradation in a particular application, or what a lens designer prefers to use for optimization during the design process.

Classification

Aberrations are classified and categorized in a variety of ways. These include pupil dependence, field dependence, order, evenness and oddness, pupil and field symmetry, and the nature of change through focus—symmetrical and unsymmetrical. In addition, there are natural groupings, e.g., astigmatism and field curvature. The classification systems overlap, and the decompositions are not unique. The complete aberration is often described as a series of terms, several schemes being used, as discussed below. The names of aberrations, such as “spherical,” “coma,” and “astigmatism,” are not standardized, and a given name may have different meanings with respect to different expansions. Furthermore, the effects of aberrations are not simply separated. For example, “pure coma” can have effects usually associated with distortion. Defocus is sometimes taken to be a type of aberration, and

it is useful to think of it in this way, since it is represented by a term in the same expansion and since the effects of aberrations vary with focus. The number of terms in an expansion is infinite, and familiar names are sometimes associated with unfamiliar terms. To improve clarity, it is recommended that all the terms in an expansion be made explicit up to agreed-upon values too small to matter, and that, in addition, the net effect be shown graphically. Further, it is often helpful to show more than one of an aberration's manifestations.

Pupil and Field Coordinates

In this section, all the quantities in the equation are in image space, so primes are omitted. Field coordinates are x and y , with $h^2 = x^2 + y^2$, and (x, y) is the nominal image point in a plane $z = 0$. Direction cosines equally spaced on the exit pupil should be used for pupil coordinates but, in practice, different types of coordinates are used, including linear positions, spatial frequencies, and direction cosines. Here the pupil coordinates are ξ and η , which are dimensionless, with $\rho^2 = \xi^2 + \eta^2$. The overall direction of the pupil may vary with field. Here the $(\xi, \eta) = (0, 0)$ is always taken at the pupil center, the meaning of which may not be simple, as discussed in the section on pupils above. The angle of a meridian in the pupil is ψ . Entrance and exit pupil coordinates must be distinguished. For diffraction calculations, the exit pupil should be sampled at equal intervals in direction cosines, but a set of rays from an object point that is equally spaced in direction cosines may leave with uneven spacing, as a result of aberrations.

Wavefront Aberration

If an object point is imaged stigmatically, then the optical path lengths of all rays from the object point to its image are identical, and the geometrical wavefronts leaving the exit pupil are spherical. In the presence of aberrations, the wavefront is no longer spherical. Rather than describing the wavefront shape, it is usually preferable to consider the difference between the actual wavefront, and a nominal wavefront, often called the *reference sphere*, centered at a *reference point* that is usually the nominal image point. This reference sphere is usually taken to intersect the center of the pupil, since this gives the most accurate diffraction calculations. The *wavefront aberration* W is the optical path length from reference sphere to wavefront, or vice versa, according to the convention used, Fig. 31. Two sign conventions are in use; a positive wavefront aberration may correspond either to a wavefront which lags or leads the reference sphere. For each nominal image point (x, y, z) , the wavefront aberration is a function of the pupil coordinates (ξ, η) , so the functional form is $W(\xi, \eta; x, y, z)$,

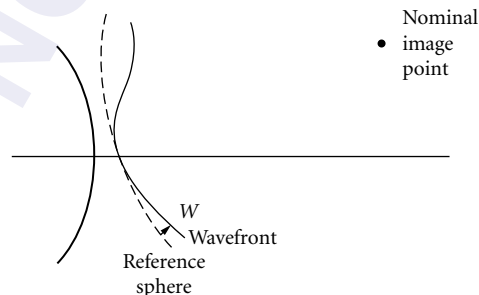


FIGURE 31 Wavefront aberration. The reference sphere is concentric with the nominal image point. The wavefront is taken that is tangent to the reference sphere in the center of the pupil. The wavefront aberration function is the distance from the reference sphere to the wavefront as a function of pupil coordinates.

with the z usually suppressed, since the image plane is usually taken to be fixed. For a given lens prescription, W is found by tracing a set of rays from each object point to the reference sphere and calculating their path lengths. If the absolute path length is unimportant, the choice of the reference sphere's radius is not critical. Considered from the point of view of wave optics, the image of a point is degraded by phase differences across the reference sphere, so absolute phase is of no consequence, and the zero of the wavefront aberration can be chosen arbitrarily. By convention and convenience, the zero is usually taken at the center of the pupil, so $W(0, 0, x, y) = 0$. Absolute optical path lengths are significant for imaging systems with paths that separate between object and image in cases where there is coherence between the various image contributions. An error in absolute optical path length is called *piston error*. This results in no ray aberrations, so it is omitted from some discussions.

Ray Aberrations

In the presence of aberrations, the rays intersect any surface at different points than they would otherwise. The intersection of the rays with the receiving surface, usually a plane perpendicular to the axis, is most often of interest. The *transverse ray aberration* is the vectorial displacement (ϵ_x, ϵ_y) between a nominal intersection point and the actual one. The displacement is a function of the position of the nominal image point (x, y) and the position in the pupil through which the ray passes (ξ, η) . A complete description of transverse ray aberrations is given by

$$\epsilon_x(\xi, \eta; x, y) \quad \text{and} \quad \epsilon_y(\xi, \eta; x, y) \quad (308)$$

The *longitudinal aberration* is the axial displacement from nominal of an axial intersection point. This description is useful for points on the axis of rotationally symmetrical systems, in which case all rays intersect the axis. Such aberrations have both transverse and longitudinal aspects. The intersection with a meridian can also be used. The *diapoint* is the point where a ray intersects the same meridian as that containing the object point.²⁴⁷ For an image nominally located at infinity, aberrations can be described by the slope of the wavefront relative to that of the nominal, that is, by ray angles rather than intersection points. A hypothetical ideal focusing lens can also be imagined to convert to transverse aberrations.

A *ray intercept diagram* shows the intersection points of a group of rays with the receiving surface.²⁵⁴ The rays are usually taken to arise from a single object point and to uniformly sample the pupil, with square or hexagonal arrays commonly used. The ray intercept diagrams can suffer from artifacts of the sampling array, which can be checked for by using more than one type of array. Other pupil loci, for instance, principal meridians and annuli, can be employed to show particular aspects of the aberration. Intercept diagrams can also be produced for a series of surfaces through focus. Image quality may be better than ray diagrams suggest, since destructive interference can reduce the irradiance in a region relative to that predicted by the ray density.

Relationship of Wavefront and Ray Aberrations

Since rays are normal to geometrical wavefronts, Fig. 32, transverse ray aberrations are proportional to the slope of the wavefront aberration function. For systems of rotation with image space index n and marginal ray angle θ , the transverse aberrations are to a good approximation²⁵¹

$$\epsilon_x = \frac{1}{n \sin \theta} \frac{\partial W}{\partial \xi} \quad \epsilon_y = \frac{1}{n \sin \theta} \frac{\partial W}{\partial \eta} \quad (309)$$

The refractive index appears since W is an optical path length. If the rays are nominally parallel, then the partial derivatives give the angular ray errors

$$\Delta\alpha = \frac{1}{np} \frac{\partial W}{\partial \xi} \quad \Delta\beta = \frac{1}{np} \frac{\partial W}{\partial \eta} \quad (310)$$

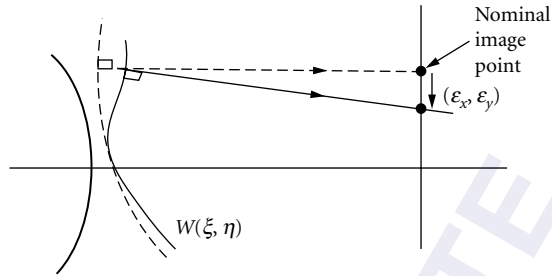


FIGURE 32 Ray aberration. Rays intersect the receiving plane at positions shifted from the nominal.

where p is the linear radius of the exit pupil, which cannot be infinite if the image is at infinity. These expressions may also have a multiplicative factor of -1 , depending on the sign conventions. A sum of wavefront aberrations gives a transverse aberration that is the sum of the contributing ones.

Ray Densities

The density of rays near the nominal image point is²⁵¹

$$\frac{1}{\text{Density}} \propto \left(\frac{\partial^2 W}{\partial \xi^2} \right) \left(\frac{\partial^2 W}{\partial \eta^2} \right) - \left(\frac{\partial^2 W}{\partial \xi \partial \eta} \right)^2 \quad (311)$$

Caustics are the surfaces where ray densities are infinite. Here, geometrical optics predicts infinite power/area, so the ray model is quantitatively inaccurate in this case.

Change of Reference Points

The center of the reference sphere may be displaced from the nominal image point. If the reference point is changed by linear displacement $(\delta x, \delta y, \delta z)$, then the wavefront aberration function changes from W to W' according to

$$W'(\xi, \eta; x, y; \delta x, \delta y, \delta z) = W(\xi, \eta; x, y) + W_x \xi + W_y \eta + W_z (\xi^2 + \eta^2) \quad (312)$$

where $W_x = n \sin \theta \delta x$

$$W_y = n \sin \theta \delta y \quad (313)$$

$$W_z = \frac{1}{2} n \sin^2 \theta \delta z$$

The transverse ray aberration ϵ'_x and ϵ'_y with respect to the new reference points are

$$\epsilon'_x = \epsilon_x + \delta x + \sin \theta \delta z \quad \epsilon'_y = \epsilon_y + \delta y + \sin \theta \delta z \quad (314)$$

The change through focus is accounted for by varying δz . Setting $\epsilon'_x = \epsilon'_y = 0$ gives the parametric equations $x(\delta z)$ and $y(\delta z)$ for a ray with pupil coordinates (ξ, η) , relative to the nominal ray near the nominal image point.

Aberration Symmetries for Systems with Rotational Symmetry

If the lens, including the aperture, is a figure of rotation, only certain aberration forms are possible. For object points on axis, the wavefront aberration and the image blur are figures of revolution. For off-axis points, both wavefront aberration and blur are bilaterally symmetrical about the meridional plane containing the object point. For object points on a circle centered on the axis, the wavefront and ray aberrations are independent of azimuth, relative to the local meridian. In practice, there is always some imperfection, so the symmetries are imperfect and additional aberration forms arise.

Wavefront Aberration Forms for Systems with Rotational Symmetry

Here the pupil is taken to be circular, with the coordinate origin taken at the center. The field coordinates are normalized so $x^2 + y^2 = h^2 = 1$ at the edge of the field. The pupil coordinates are normalized, so that $\xi^2 + \eta^2 = \rho^2 = 1$ on the rim of the pupil. The algebra is simplified by using dimensionless coordinates. To add dimensions and actual sizes, replace the ξ by ξ/ξ_{\max} and likewise for other variables. The simplest combinations of pupil and field coordinates with rotational symmetry are

$$x^2 + y^2 = h^2, \quad \xi^2 + \eta^2 = \rho^2, \quad \xi x + \eta y \quad (315)$$

The general wavefront aberration function can be expressed as a series of such terms raised to integral powers,

$$W(x, y; \xi, \eta) = \sum_{L, M, N=0} W_{LMN} (x^2 + y^2)^L (\xi^2 + \eta^2)^M (x\xi + y\eta)^N \quad (316)$$

where L, M, N are positive integers. The terms can be grouped in *orders* according to the sum $L + M + N$, where, by convention, the order equals $2(L + M + N) - 1$. The order number refers more directly to ray aberration forms than to wavefront forms, and it is always odd. The first-order terms are those for which $L + M + N = 1$, for the third-order terms the sum is two, and so on. The number of terms in the Q th order is $1 + (Q + 1)(Q + 7)/8$. For orders 1, 3, 5, 7, 9 the number of terms is 3, 6, 10, 15, 21. For each order, one contribution is a piston error, which is sometimes excluded from the count.

The expression of Eq. (316) is related to the characteristic function for a rotationally symmetrical system, Eq. (32). If the spatial coordinates are taken to be those of the object point, this is the point-angle characteristic function. In the hamiltonian optics viewpoint, the characteristic function is a sum of two parts. The first-order terms specify the nominal properties, and those of higher orders the deviation therefrom. This is discussed in the references given in that section. The term for which $L = M = N = 0$ has to do with absolute optical path length.

Since there is bilateral symmetry about all meridians, the expansion can be simplified by considering object points in a single meridian, customarily taken to be that for which $x = 0$. Doing so and letting the fractional field height be $y = h$ gives the wavefront aberration function

$$W(h; \rho, \eta) = \sum_{L, M, N=0} W_{LMN} h^{2L+N} \rho^{2M} \eta^N = \sum_{A, B, C} W'_{ABC} h^A \rho^B \eta^C \quad (317)$$

where $A = 2L + N$, $B = 2M$, $C = N$, and the order equals $(A + B + C) - 1$. Another form is obtained with the fractional pupil radius ρ and the pupil azimuth ψ , the angle from the $x = 0$ meridian, so $\eta = \rho \cos \psi$. With these pupil variables the wavefront aberration function is

$$W(h; \rho, \psi) = \sum_{L, M, N=0} W_{LMN} h^{2L+N} \rho^{2M+N} \cos^N \psi = \sum_{A, B, C} W''_{ABC} h^A \rho^B \cos^C \psi \quad (318)$$

where $A = 2L + N$, $B = 2M + N$, $C = N$, and the order is $A + B - 1$. For orders above the first, the W_{LMN} , W'_{ABC} , and W''_{ABC} are the *wavefront aberration coefficients*.

For a given field position, the wavefront aberration function for circular pupils can also be decomposed into the *Zernike polynomials*, also called *circle polynomials*, a set of functions complete and orthonormal on a circle.^{249,255–257}

Third-Order Aberrations and Their Near Relatives

There are six third-order terms. The *Seidel aberrations* are spherical, coma, astigmatism, field curvature, distortion, and there is also a piston-error term. Expressions for these aberrations are given below, along with some higher-order ones that fall in the same classification. The terminology of higher-order aberrations is not standardized, and there are forms that do not have third-order analogues. This section uses the notation of the second expression of Eq. (318), without the primes on the coefficients.

It is useful to include *defocus* as a term in aberration expansions. Its wavefront aberration and transverse ray aberrations are

$$W = W_{020}\rho^2 \quad \varepsilon_x \propto 2W_{020}\xi \quad \varepsilon_y \propto 2W_{020}\eta \quad (319)$$

Coefficient W_{020} is similar to W_z , Eq. (313).

In *spherical aberration* the wavefront error is a figure of revolution in the pupil. The individual terms of the expansion have the form ρ^{2N} . The form that appears on axis, and which is independent of field position is

$$W = W_{020}\rho^2 + W_{040}\rho^4 + W_{060}\rho^6 + \dots \quad (320)$$

where defocus has been included. The W_{040} term is the third-order term, the W_{060} is the fifth-order term, etc. The ray aberrations are

$$\begin{aligned} \varepsilon_x &\propto 2W_{020}\xi + 4W_{040}\rho^2\xi + 6W_{060}\rho^4\xi + \dots \\ \varepsilon_y &\propto 2W_{020}\eta + 4W_{040}\rho^2\eta + 6W_{060}\rho^4\eta + \dots \end{aligned} \quad (321)$$

There are also higher-order off-axis terms, called *oblique spherical aberration*, with forms $h^{2L}\rho^{2M}$. Spherical is an even aberration.

In *coma*, the wavefront aberration varies linearly with field height, so the general form is $h\rho^{2M}\eta = h\rho^{2M+1}\cos\psi$. Coma is an odd aberration. The wavefront expansion is

$$W = (W_{131}\rho^2 + W_{151}\rho^4 + \dots)\eta h = (W_{131}\rho^3 + W_{151}\rho^5 + \dots)\cos\psi h \quad (322)$$

The ray aberrations are

$$\begin{aligned} \varepsilon_x &\propto [W_{131}(2\xi\eta) + 4W_{151}(\xi^2 + \eta^2)\xi\eta + \dots]h \\ \varepsilon_y &\propto [W_{131}(\xi^2 + 3\eta^2) + W_{151}(\xi^4 + 5\xi^2\eta^2 + 6\eta^4) + \dots]h \end{aligned} \quad (323)$$

In *astigmatism* the wavefront aberration is cylindrical. The third-order term is

$$W = W_{222}h^2\eta^2 = W_{222}h^2\rho^2\cos^2\psi \quad (324)$$

with ray aberration

$$\varepsilon_x = 0 \quad \varepsilon_y \propto 2W_{222}h^2\eta \quad (325)$$

Field curvature, also known as *Petzval curvature*, is a variation of focal position in the axial direction with field height. In its presence, the best image of a planar object lies on a nonplanar surface. Its absence is called *field flatness*. The wavefront aberration form is

$$W = (W_{220}h^2 + W_{420}h^4 + W_{620}h^6 + \dots)\rho^2 \quad (326)$$

with symmetrical blurs given by

$$\begin{aligned} \varepsilon_x &\propto (W_{220}h^2 + W_{420}h^4 + W_{620}h^6 + \dots)\xi \\ \varepsilon_y &\propto (W_{220}h^2 + W_{420}h^4 + W_{620}h^6 + \dots)\eta \end{aligned} \quad (327)$$

The curvature of the best focus surface may have the same sign across the field, or there may be curvatures of both signs.

Astigmatism and field curvature are often grouped together. Combining defocus, third-order astigmatism, and third-order field curvature, the wavefront aberration can be written

$$W = W_{020}(\xi^2 + \eta^2) + [W_{220}\xi^2 + (W_{220} + W_{222})\eta^2]h^2 \quad (328)$$

The resultant ray aberration is

$$\varepsilon_x \propto [W_{020} + W_{220}h^2]\xi, \quad \varepsilon_y \propto [W_{020} + (W_{222} + W_{220})h^2]\eta \quad (329)$$

A tangential fan of rays, one that lies in the $x = 0$ meridian, has $\xi = 0$, so $\varepsilon_x = 0$. The *tangential focus* occurs where $\varepsilon_y = 0$, which occurs for a defocus of $W_{020} = -(W_{220} + W_{222})h^2$. Combining this result with Eq. (314) gives $\delta z \propto h^2$, the equation for the tangential focal surface. A sagittal fan of rays crosses the pupil in the $\eta = 0$ meridian, so $\varepsilon_y = 0$. The *sagittal focus* occurs where $\varepsilon_x = 0$, i.e., on the surface given by $W_{020} = -W_{220}h^2$.

In general, *distortion* is a deviation from geometrical similarity between object and image. For rotationally symmetrical lenses and object and image planes perpendicular to the axis, the error is purely radial, and can be thought of as a variation of magnification with field height. The aberration forms are

$$W = (W_{111}h + W_{311}h^3 + W_{511}h^5 + \dots)\eta \quad (330)$$

with

$$\varepsilon_x = 0, \quad \varepsilon_y \propto W_{111}h + W_{311}h^3 + W_{511}h^5 + \dots \quad (331)$$

In *pincushion* distortion the magnitude of magnification increases monotonically with field height, so the image is stretched radially. In *barrel distortion* the magnitude decreases, so the image is squeezed. In general, the aberration coefficients can be both positive and negative, so the direction of distortion can change as a function of field height and the distortion may vanish for one or more field heights.

For *piston error* the wavefront differs uniformly across the pupil from its nominal in a way that varies with field height.

$$W = W_{000} + W_{200}h^2 + W_{400}h^4 + W_{600}h^6 + \dots, \quad \varepsilon_x = \varepsilon_y = 0 \quad (332)$$

There are no transverse ray aberrations.

Chromatic Aberrations

In general, the properties of optical systems vary with wavelength. The term *chromatic aberration* often refers to the variation in paraxial properties as a function of wavelength. Thus, *axial color* is related to differences of focal length and principal plane location with wavelength, and *lateral color*

is related to variations of magnification with wavelength. Also, the monochromatic aberrations vary in magnitude with wavelength. Reflective systems have identical ray properties at all wavelengths, but their wave properties vary with color, since a given variation in path length has an effect on phase that varies with wavelength.

Stop Size and Aberration Variation

For a given system, if the size of the aperture is changed, the marginal ray is changed, but not the chief ray. If the aperture is reduced, depth of focus and depth of field increase and image irradiance decreases. The rays from axial object points are more nearly paraxial, so the imaging tends to be better corrected. For off-axis points, some aberrations are changed and others are not. Distortion, as defined with respect to the chief ray, is not changed. Field curvature per se does not change, since the aperture size does not change the location of the best image surface (if there are no other aberrations), but the depth of focus does change, so a flat detector can cover a larger field.

Stop Position and Aberration Variation

For a given system, if the aperture is moved axially, the image-forming bundle passes through different portions of the lens elements. Accordingly, some aberrations vary with the position of the stop. Lens design involves an operation called the *stop shift*, in which the aperture is moved axially while its size is adjusted to keep the numerical apertures constant. In this operation, the marginal ray is fixed, while the chief ray is changed. This does not change the aberrations on axis. Most of those for off-axis points are changed, but third-order field curvature is unchanged.

1.20 REFERENCES

1. R. K. Luneburg, *Mathematical Theory of Optics*, University of California Press, Berkeley, 1964, chap. 1.
2. M. Kline and I. W. Kay, *Electromagnetic Theory and Geometrical Optics*, Interscience Publishers, New York, 1965.
3. M. Born and E. Wolf, *Principles of Optics*, Pergamon, New York, 1980.
4. D. Marcuse, *Light Transmission Optics*, Van Nostrand Reinhold, New York, 1989, chap. 3.6.
5. L. Rayleigh, "Optics," *Encyclopedia Britannica* vol. XVII, 1884, and *Scientific Papers*, Dover, New York, 1964, sec. 119, p. 385.
6. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986, chap. 6.
7. M. Herzberger, "On the Fundamental Optical Invariant, the Optical Tetrality Principle, and on the New Development of Gaussian Optics Based on This Law," *J. Opt. Soc. Am.*, vol. 25, 1935, pp. 295–304.
8. O. N. Stavroudis, *The Optics of Rays, Wavefronts, and Caustics*, Academic, New York, 1972, chap. 13, "The Fundamental Optical Invariant."
9. G. Forbes and B. D. Stone, "Restricted Characteristic Functions for General Optical Configurations," *J. Opt. Soc. Am.*, vol. 10, 1993, pp. 1263–1269.
10. P. Fermat, "Fermat's principle," *A Source Book in Physics*, W. F. Magie, Harvard, Cambridge, Mass., 1963, p. 278.
11. P. de Fermat, *Oeuvres*, vol. II, P. Tomnery et C. Henry (eds.), Paris, 1891.
12. P. de Fermat, *Oeuvres de Fermat*, C. Henry, P. Tannery (eds.), vols. i–iv and supplement Gauthier-Villars, Paris, 1854–1922 vols. I–IV and supplement vol. III, French translations of Latin papers.
13. R. P. Feynman, R. B. Leighton, and M. Sands, *The Feynman Lectures on Physics*, Addison Wesley, Reading, Mass., 1963 vol 1, chap. 6, "Optics: The Principle of Least Time."
14. B. B. Rossi, *Optics*, Addison-Wesley, Reading, Mass., 1957.
15. E. Hecht, *Optics*, Addison-Wesley, Reading, Mass., 1987.
16. M. V. Klein and T. E. Furtak, *Optics*, 2d ed. Wiley, New York, 1986.

17. H. H. Hopkins, "An Extension of Fermat's Theorem," *Optica Acta*, vol. 17, 1970, pp. 223–225.
18. A. Walther, "Systematic Approach to the Teaching of Lens Theory," *Am. J. Phys.*, vol. 35, 1967, pp. 808–816.
19. A. Walther, "Lenses, Wave Optics, and Eikonal Functions," *J. Opt. Soc. Am.*, vol. 59, 1969, pp. 1325–1333.
20. L. Landau and E. Lifshitz, *Classical Theory of Fields*, Addison-Wesley, Reading, Mass., 1951, sec. 7.1, "Geometrical Optics."
21. J. A. Kneisly, "Local Curvature of Wavefronts in an Optical System," *J. Opt. Soc. Am.*, vol. 54, 1965, pp. 229–235.
22. W. R. Hamilton, *Geometrical Optics. Mathematical Papers of Sir William Rowan Hamilton*, vol. 1, A. W. Conway and J. L. Synge (eds.), Cambridge University Press, 1931.
23. H. Bruns, "Das Eikonal," *Abh. Königl. sächs. Abhand. Ges. Wiss., Leipzig, Math Phys K1*, vol. 21, 1895, pp. 325–436.
24. K. Schwarzschild, *Untersuchungen zur geometrischen Optik*, I–II, *Abh. Königl. Ges. Wiss. Goettingen, Math Phys. K1*, Neue Folge 4, 1905, pp. 1–54.
25. Czapski-Eppenstein, *Grundzuge der Theorie der Optische Instrumente*, Leipzig, 1924.
26. C. G. Steward, *The Symmetrical Optical System*, Cambridge University Press, 1928.
27. M. Herzberger, *Strahlenoptik*, Springer, Berlin, 1931.
28. J. L. Synge, *Geometrical Optics. An Introduction to Hamilton's Method*, Cambridge University Press, 1937.
29. C. Caratheodory, *Geometrische Optik*, Springer-Verlag, Berlin, 1937.
30. L. Rayleigh, "Hamilton's Principle and the Five Aberrations of Von Seidel," *Phil. Mag.*, vol. XV, 1908, pp. 677–687, or *Scientific Papers*, Dover, New York, 1964.
31. R. J. Pegis, "The Modern Development of Hamiltonian Optics" *Progress in Optics*, E. Wolf (ed.), vol. 1, North Holland, 1961, pp. 1–29.
32. R. K. Luneburg, *Mathematical Theory of Optics*, University of California Press, Berkeley, 1964, sec. 19.
33. W. Brouwer and A. Walther, "Geometrical Optics," *Advanced Optical Techniques*, A. C. S. Van Heel (ed.), North-Holland, Amsterdam, 1967, pp. 503–570.
34. H. A. Buchdahl, *An Introduction to Hamiltonian Optics*, Cambridge U.P., 1970.
35. M. Born and E. Wolf, *Principles of Optics*, Pergamon, New York, 1980, chap. 4.
36. M. Herzberger, "Geometrical Optics," *Handbook of Physics*, E. U. Condon and H. Odishaw (ed.), McGraw-Hill, 1958, chap. 6-2.
37. M. Herzberger, "On the Characteristic Function of Hamilton, the Eiconal of Bruns, and Their Use in Optics," *J. Opt. Soc. Am.*, vol. 26, 1936, pp. 177–180.
38. M. Herzberger, "Hamilton's Characteristic Function and Bruns Eiconal," *J. Opt. Soc. Am.*, vol. 27, 1937, pp. 133–317.
39. J. L. Synge, "Hamilton's Method in Geometrical Optics," *J. Opt. Soc. Am.*, vol. 27, 1937, pp. 75–82.
40. T. Smith, "On Perfect Optical Instruments," *Proc. Roy. Soc.*, vol. 55, 1945, pp. 293–304.
41. C. H. F. Velzel, "Breaking the Boundaries of Optical System Design and Construction," *Trends in Modern Optics*, J. W. Goodman (ed.), Academic, Boston, 1991, pp. 325–338.
42. R. S. Heath, *A Treatise on Geometrical Optics*, Cambridge University Press, 1895, chap. 13, "Refraction through Media of Varying Density."
43. R. A. Herman, *A Treatise on Geometrical Optics*, Cambridge University Press, 1900, chap. 13 "Heterogeneous Media"
44. J. L. Synge, *Geometrical Optics. An Introduction to Hamilton's Method*, Cambridge University Press, 1937, chap. 5, "Heterogeneous Isotropic Media."
45. R. K. Luneburg, *Mathematical Theory of Optics*, University of California Press, Berkeley, 1964, chaps. 2, 3.
46. O. N. Stavroudis, *The Optics of Rays, Wavefronts, and Caustics*, Academic, New York, 1972, chap. 11 "The Inhomogeneous Medium."
47. A. K. Ghatak and K. Thyagarajan, *Contemporary Optics*, Plenum, New York, 1978.
48. M. Born and E. Wolf, *Principles of Optics*, Pergamon, New York, 1980, chap. 3.2.
49. D. Marcuse, *Light Transmission Optics*, Van Nostrand Reinhold, New York, 1989, chap. 3.

50. E. W. Marchand, "Gradient Index Lenses," *Progress in Optics*, vol. 11, North Holland, Amsterdam, 1973, pp. 305–337.
51. E. W. Marchand, *Gradient Index Optics*, Academic, New York, 1978.
52. A. Sharma, D. V. Kumar, and A. K. Ghatak, "Tracing Rays Through Graded-Index Media: A New Method," *Appl. Opt.* vol. 21, 1982, pp. 984–987.
53. D. T. Moore, ed., *Selected Papers on Gradient-Index Optics*, SPIE, 1992.
54. D. Moore, "Gradient Index Optics," *Handbook of Optics*, vol. II, chap. 9, McGraw-Hill, New York, 2d ed., 1994.
55. J. Brown, *Microwave Lenses*, Methuen, London, 1953.
56. S. Cornbleet, *Microwave Optics*, Academic, New York, 1976, chap. 2, "Non-Uniform Media."
57. S. Cornbleet, "Geometrical Optics Reviewed: A New Light on an Old Subject," *Proc. IEEE*, vol. 71, 1983, pp. 471–502.
58. S. Cornbleet, *Microwave and Optical Ray Geometry*, Wiley, New York, 1984.
59. W. Blaschke, *Vorlesungen über Differentialgeometrie*, Springer, 1930, Dover, 1945.
60. E. Kreyszig, *Differential Geometry*, Univ. Toronto, 1959, Dover, New York, 1991.
61. J. J. Stoker, *Differential Geometry*, Wiley, New York, 1969.
62. D. J. Struik, *Lectures in Classical Differential Geometry*, Addison-Wesley, 1961, Dover, 1990.
63. H. Goldstein, *Classical Mechanics* Addison-Wesley, Reading, Mass., 1980, 2d edition. Variational principles in physics are discussed in many books on classical mechanics, a good reference list is supplied by Goldstein.
64. J. A. Arnaud, "Analogy Between Optical Rays and Nonrelativistic Particle Trajectory: A Comment," *Am. J. Phys.*, vol. 44, 1976, pp. 1067–1069.
65. J. Evans, M. Rosenquist, "' $F = ma$ ' Optics," *Am. J. Phys.*, vol. 54, 1986, pp. 876–883.
66. D. T. Moore, "Ray Tracing in Gradient-Index Media," *J. Opt. Soc. Am.*, vol. 65, 1975, pp. 451–455.
67. G. W. Forbes, "On Variational Problems in Parametric Form," *Am. J. Phys.*, vol. 59, 1991, pp. 1130–1140.
68. L. Landau and E. Lifshitz, *Classical Theory of Fields*, Addison-Wesley, Reading, Mass., 1951, sec. 7.1, "Geometrical Optics."
69. H. A. Buchdahl, "Rays in Gradient Index Media: Separable Systems," *J. Opt. Soc. Am.*, vol. 63, 1973, pp. 46–49.
70. S. Cornbleet, "Ray Paths in Nonuniform Axially Symmetric Medium," *IEE Journal of Microwaves Optics and Acoustics*, vol. 2, 1978, pp. 194–200.
71. S. Cornbleet, "Ray Tracing in a General Inhomogeneous Cylindrical Medium," *IEE Journal of Microwaves Optics and Acoustics*, vol. 3, 1979, pp. 181–185.
72. J. W. Blaker and M. A. Tavel, "The Application of Noether's Theorem to Optical Systems," *Amer. Jour. Phys.*, vol. 42, 1974, pp. 857–861.
73. W. B. Joyce, "Comments on 'The Application of Noether's Theorem to Optical Systems,'" *Amer. Jour. Phys.*, vol. 43, 1975, p. 455.
74. J. C. Maxwell, "Solutions of Problems," problem no. 2, *Cambridge and Dublin Math. Journal*, vol. 8, 1854, p. 188, and *The Scientific Papers of James Clerk Maxwell*, Cambridge University Press, 1890, Dover, New York, 1965, pp. 74–79.
75. S. P. Morgan, "General Solution of the Luneberg Lens Problem," *J. Appl. Phys.*, vol. 29, 1958, pp. 1358–1368.
76. W. H. Steel, "Luminosity, Throughput, or Étendue?" *Appl. Opt.*, vol. 13, 1974, pp. 704–705.
77. R. Winston, W. T. Welford, "Geometrical Vector Flux and Some Nonimaging Concentrators," *J. Opt. Soc. Am.*, vol. 69, 1979, pp. 532–536.
78. D. Gabor, "Light and Information," *Progress in Optics*, vol. 1, North-Holland, Amsterdam, 1961, chap. 4, pp. 111–153.
79. W. T. Welford, R. Winston, *The Optics of Nonimaging Concentrators*, Academic, New York, 1978, app. A.
80. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986.
81. W. T. Welford and R. Winston, *High Collection Nonimaging Optics*, Academic, New York, 1989, app. A.
82. G. T. di Francia, "Parageometrical Optics," *J. Opt. Soc. Am.*, vol. 40, 1950, pp. 600–602.
83. R. Winston, "Light Collection within the Framework of Geometrical Optics," *J. Opt. Soc. Am.*, vol. 60, 1970, pp. 245–247.

84. T. Jansson and R. Winston, "Liouville's Theorem and Concentrator Optics," *J. Opt. Soc. Am. A.*, vol. 3, 1986, pp. 7–8.
85. D. Marcuse, *Light Transmission Optics*, Van Nostrand Reinhold, New York, 1989, chap. 3.7, "Liouville's Theorem," pp. 112–124.
86. R. Clausius, "Die Concentration von Wärme und Lichtstrahlen und die Grenzen ihrer Wirkung," *Pogg. Ann.*, cxxi, S. 1, 1864.
87. R. Clausius, *The Mechanical Theory of Heat*, Macmillan, London, 1879, "On the Concentration of Rays of Heat and Light and on the Limits of its Action." (English translation.)
88. H. Helmholtz, "Die theoretische Grenze für die Leitungsfähigkeit der Mikroskope," *Pogg. Ann. Jubelband*, 1874, pp. 556–584.
89. S. Liebes, "Brightness—On the Ray Invariance of B/n^2 ," *Am. J. Phys.*, vol. 37, 1969, pp. 932–934.
90. Nicodemus, "Radiance," *Am. J. Phys.*, vol. 31, 1963, pp. 368–377.
91. R. W. Boyd, *Radiometry and the Detection of Optical Radiation*, Wiley, New York, 1983, chap. 5.
92. M. V. Klein, T. E. Furtak, *Optics*, 2d ed., Wiley, New York, 1986, sec. 4.2, "Radiometry and Photometry."
93. L. Rayleigh, "Notes, Chiefly Historical, on Some Fundamental Propositions in Optics," *Phil. Mag.*, vol. 21, 1886, pp. 466–476. Collected works, vol. II, Dover, New York, 1964, article 137, pp. 513–521.
94. J. P. C. Southall, *Principles and Methods of Geometrical Optics*, Macmillan, London, 1910.
95. T. Smith, "On Tracing Rays through an Optical System," *Proc. Phys. Soc. London*, vol. 33, 1921, pp. 174–178.
96. H. H. Hopkins, "An Optical Magnification for Finite Objects," *Nature*, vol. 159, 1947, pp. 573–574.
97. J. Marshall, "Particle Counting by Cerenkov Radiation," *Phys. Rev.*, vol. 86, 1952, pp. 685–693.
98. H. A. Buchdahl, *Optical Aberration Coefficients*, Oxford, London, 1954, sec. 4.
99. M. Herzberger, *Modern Geometrical Optics*, Interscience Publishers, New York, 1958.
100. W. T. Welford, "A Note on the Skew Invariant of Optical Systems," *Optica Acta*, vol. 15, 1968, pp. 621–623.
101. O. N. Stavroudis, *The Optics of Rays, Wavefronts, and Caustics*, Academic, New York, 1972, p. 208.
102. W. T. Welford, *Aberrations of the Symmetrical Optical System*, Academic Press, London, 1974, sec. 5.4, p. 66.
103. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986, sec. 6.4, p. 84.
104. W. T. Welford and R. Winston, *High Collection Nonimaging Optics*, Academic, New York, 1989, p. 228.
105. R. Descartes, "Snell's Law," *A Source Book in Physics*, W. F. Magie, ed., Harvard, Cambridge, Mass., 1963, p. 265.
106. L. Silberstein, *Simplified Methods of Tracing Rays through any Optical System of Lenses, Prisms, and Mirrors*, Longmans, Green, and Co., London, 1918, sec. 8, "Dyadic Representing the Most General Reflector."
107. T. Smith, "On Systems of Plane Reflecting Surfaces," *Trans. Opt. Society*, vol. 30, 1928, pp. 68–78.
108. R. K. Luneburg, *Mathematical Theory of Optics*, University of California Press, Berkeley, 1964, app. II.
109. R. E. Hopkins, *Applied Optics and Optical Engineering*, R. Kingslake (ed.), vol. 3, Academic, New York, 1965, chap. 7, pp. 269–308. "Mirror and Prism Systems."
110. L. Levi, *Applied Optics: A Guide to Modern Optical System Design*, vol. 1, Wiley, New York, 1968, pp. 346–354.
111. Lian Tongshu, *Theory of Conjugation for Reflecting Prisms: Adjustment and Image Stabilization of Optical Instruments*, International Academic Publishers, Pergamon, Oxford, 1991.
112. H. Wagner, "Zur mathematischen Behandlung von Spiegelungen," *Optik*, vol. 8, 1951, pp. 456–472.
113. G. H. Spencer, M. V. R. K. Murty, "General Ray-Tracing Procedure," *J. Opt. Soc. Am.*, vol. 52, 1962, pp. 672–678.
114. G. T. di Francia, "Parageometrical Optics," *J. Opt. Soc. Am.*, vol. 40, 1950, pp. 600–602.
115. *Webster's International Unabridged Dictionary*, 2d ed., 1934.
116. H. A. Buchdahl, *An Introduction to Hamiltonian Optics*, Cambridge University Press, 1970.
117. J. C. Maxwell, "On the General Laws of Optical Instruments," *Quarterly Journal of Pure and Applied Mathematics*, Feb. 1858, and *The Scientific Papers of James Clerk Maxwell*, Cambridge UP, 1890, Dover, New York, 1965.
118. J. P. C. Southall, *The Principles and Methods of Geometrical Optics*, Macmillan, London, 1910.
119. M. Born and E. Wolf, *Principles of Optics*, Pergamon, New York, 1980, sec. 4.2.3.
120. T. Smith, "The Optical Cosine Law," *Trans. Opt. Soc. London*, vol. 24, 1922–23, pp. 31–39.

121. G. C. Steward, *The Symmetrical Optical Systems*, Cambridge University Press., 1928.
122. W. T. Welford, "Aplanatism and Isoplanatism," *Progress in Optics*, vol. 13, E. Wolf (ed.), North-Holland, Amsterdam, 1976, pp. 267–292.
123. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986, sec. 94.
124. E. Abbe, "Über die Bedingungen des Aplanatismus der Linsensysteme, Sitzungsber. der Jenaischen Gesellschaft für Med. u. Naturw.," 1879, pp. 129–142; *Gesammelte Abhandlungen*, Bd. I, pp. 213–226; *Carl. Reper. Phys.*, vol. 16, 1880.
125. A. E. Conrady, *Applied Optics and Optical Design*, Oxford University Press, 1929, Dover, New York, 1992.
126. H. H. Hopkins, "The Optical Sine-Condition," *Proc. Phys. Soc. Lond.*, vol. 58, 1946, pp. 92–99.
127. T. Smith, "On Perfect Optical Instruments," *Proc. Phys. Soc.*, vol. LX 1948, pp. 293–304.
128. H. H. Hopkins, *Wave Theory of Aberrations*, Clarendon Press, Oxford, 1950.
129. J. F. W. Herschel, "On the Aberration of Compound Lenses and Object Glasses," *Phil. Trans. Roy. Soc.*, vol. 8, 1821, pp. 222–267.
130. H. H. Hopkins, "Herschel's Condition," *Proc. Phys. Soc. Lond.*, vol. 58, 1946, pp. 100–105.
131. C. H. F. Velzel, "Breaking the Boundaries of Optical System Design and Construction," *Trends in Modern Optics*, J. W. Goodman (ed.), Academic, Boston, 1991, pp. 325–338.
132. T. Smith, "The Theory of Aplanatic Surfaces," *Trans. Opt. Soc.*, vol. 29, 1927–28, pp. 179–186.
133. Military Standardization Handbook MIL-HDBK-141, *Optical Design*, Defense Supply Agency, 1962.
134. R. Kingslake, "Basic Geometrical Optics," *Applied Optics and Optical Engineering*, vol. I, R. Kingslake (ed.), Academic, New York, 1965, chap. 6.
135. H. P. Brueggemann, *Conic Mirrors*, Focal, New York, 1968.
136. R. Kingslake, *Lens Design Fundamentals*, Academic, New York, 1978.
137. R. Kingslake, *Optical System Design*, Academic, New York, 1983.
138. G. G. Slyusarev, *Aberration and Optical Design Theory*, Hilger, Bristol, 1984.
139. M. V. Klein and T. E. Furtak, *Optics*, 2d ed., Wiley, New York, 1986.
140. R. E. Hopkins, "Geometrical Optics," *Geometrical and Instrumental Optics*, D. Malacara (ed.), vol. 25 of *Methods of Experimental Physics*, Academic, New York, 1988, pp. 7–58.
141. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1990.
142. R. Shannon, "Aspheric Surfaces," *Applied Optics and Optical Engineering*, vol. 8, Academic, New York, 1980, chap. 3, pp. 55–85.
143. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986.
144. D. Malacara, "An Optical Surface and Its Characteristics," *Optical Shop Testing*, 2d ed., D. Malacara (ed.), Wiley, New York, 1992, app. I.
145. *ISO Fundamental Standards Optics and Optical Instruments*, "Indications in Optical Drawings Part 12: Aspheric Surfaces," April, ISO/TC 172/SC 1, 1992.
146. G. H. Spencer, M. V. R. K. Murty, "General Ray-Tracing Procedure," *J. Opt. Soc. Am.*, vol. 52, 1962, pp. 672–678.
147. W. T. Welford, *Aberrations of the Symmetrical Optical System*, Academic Press, London, 1974, chap. 4.
148. R. Kingslake, *Lens Design Fundamentals*, Academic, New York, 1978.
149. R. Kingslake, *Optical System Design*, Academic, New York, 1983.
150. G. G. Slyusarev, *Aberration and Optical Design Theory*, Hilger, Bristol, 1984.
151. M. V. Klein and T. E. Furtak, *Optics*, 2d ed., Wiley, New York, 1986, sec. 3.1.
152. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986, chaps. 4, 5, sec. 4.7.
153. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1992, chap. 10, sec. 10.5.
154. A. Nussbaum and R. A. Phillips, *Contemporary Optics for Scientists and Engineers*, Prentice-Hall, Englewood Cliffs, N. J., 1976, p. 95.
155. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1992.
156. L. Rayleigh, "Notes, Chiefly Historical, on some Fundamental Propositions in Optics," *Phil. Mag.*, vol. 21, 1886, pp. 466–476. Collected works, Dover, New York, 1964, vol. II, article 137, pp. 513–521.

157. J. P. C. Southall, *Principles and Methods of Geometrical Optics*, Macmillan, London, 1910.
158. M. Herzberger, "On the Fundamental Optical Invariant, the Optical Tetrality Principle, and on the New Development of Gaussian Optics Based on This Law," *J. Opt. Soc. Am.*, vol. 25, 1935, pp. 295–304.
159. W. Brouwer and A. Walther, "Geometrical Optics," *Advanced Optical Techniques*, A. C. S. Van Heel (ed.), North-Holland, Amsterdam, 1967, p. 576.
160. T. Young, *Phil. Trans. Roy. Soc.*, vol. 102, 1801, pp. 123–128.
161. H. Coddington, *A Treatise on the Reflexion and Refraction of Light*, Cambridge, 1829, pp. 66ff.
162. J. C. Maxwell, "On the Application of Hamilton's Characteristic Function to the Theory of an Optical Instrument Symmetrical about Its Axis," *Proc London Math Society*, vol. VI, 1875. *The Scientific Papers of James Clerk Maxwell*, Cambridge University Press, 1890, Dover, 1965.
163. J. C. Maxwell, "On Hamilton's Characteristic Function for a Narrow Beam of Light," *Proc. London Math. Society*, vol. VI, 1847–75. *The Scientific Papers of James Clerk Maxwell*, Cambridge University Press, 1890, Dover, 1965.
164. J. C. Maxwell, "On the Focal Lines of a Refracted Pencil," *Proc. London Math. Society*, vol. VI, 1874–75. *The Scientific Papers of James Clerk Maxwell*, Cambridge University Press, 1890, Dover, 1965.
165. J. Larmor, "The Characteristics of an Asymmetric Optical Computation," *Proc. London Math. Soc.*, vol. 20, 1889, pp. 181–194.
166. J. Larmor, "The Simplest Specification of a Given Optical Path, and the Observations Required to Determine It," *Proc. London Math. Soc.*, vol. 23, 1892, pp. 165–173.
167. R. S. Heath, *A Treatise on Geometrical Optics*, Cambridge University Press, 1895, chap. 8 "On the General Form and Properties of a Thin Pencil. General Refraction of Thin Pencils."
168. R. W. Sampson, "A Continuation of Gauss's 'Dioptrische Untersuchungen,'" *Proc London. Math. Soc.*, vol. 24, 1897, pp. 33–83.
169. T. Smith, "Note on Skew Pencils Traversing a Symmetrical Instrument," *Trans. Opt. Soc.*, vol. 30, 1928–29, pp. 130–133.
170. T. Smith, "Imagery Around a Skew Ray," *Trans. Opt. Soc.*, vol. 31, 1929–30, pp. 131–156.
171. T. Smith, "Canonical Forms in the Theory of Asymmetric Optical Systems," *Trans. Opt. Soc.*, vol. 29, 1928–1929, pp. 88–98.
172. M. Herzberger, "First-Order Laws in Asymmetrical Optical Systems. Part I. The Image of a Given Congruence: Fundamental Conceptions," *J. Opt. Soc. Am.*, vol. 26, 1936, pp. 345–359.
173. J. L. Synge, *Geometrical Optics. An Introduction to Hamilton's Method*, Cambridge University Press, 1937, chap. 3, "Thin Bundles of Rays."
174. H. H. Hopkins, "A Transformation of Known Astigmatism Formulae," *Proc. Phys. Soc.*, vol. 55, 1945, pp. 663–668.
175. H. H. Hopkins, *Wave Theory of Aberrations*, Clarendon Press, Oxford, 1950, chap. 5, "Astigmatism."
176. J. A. Kneisly, "Local Curvature of Wavefronts in an Optical System," *J. Opt. Soc. Am.*, vol. 54, 1963, pp. 229–235.
177. P. J. Sands, "When is First-Order Optics Meaningful," *J. Opt. Soc. Am.*, vol. 58, 1968, pp. 1365–1368.
178. H. A. Buchdahl, *An Introduction to Hamiltonian Optics*, Cambridge University Press, 1970, sec. 10, "Parabasal Optics in General."
179. O. N. Stavroudis, *The Optics of Rays, Wavefronts, and Caustics*, Academic, New York, 1972, chap. 10, "Generalized Ray Tracing."
180. P. J. Sands, "First-Order Optics of the General Optical System," *J. Opt. Soc. Am.*, vol. 62, 1972, pp. 369–372.
181. H. H. Hopkins, "Image Formation by a General Optical System. 1: General Theory," *Appl. Opt.*, vol. 24, 1985, pp. 2491–2505.
182. H. H. Hopkins, "Image Formation by a General Optical System. 2: Computing Methods," *Appl. Opt.*, vol. 24, 1985, pp. 2506–2519.
183. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986, sec. 9.5, "Optics Round a Finite Principal Ray," p. 185.
184. C. H. F. Velzel, "Image Formation by a General Optical System, Using Hamilton's Method," *J. Opt. Soc. Am. A*, vol. 4, 1987, pp. 1342–1348.

185. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1992, chap. 10.6, “Coddington’s Equations.”
186. B. D. Stone and G. W. Forbes, “Characterization of First-Order Optical Properties for Asymmetric Systems,” *J. Opt. Soc. Am. A*, vol. 9, 1992, pp. 478–489.
187. C. F. Gauss, *Dioptrische Untersuchungen*, delivered to the Royal Society, Dec 10, 1840. *Transactions of the Royal Society of Sciences at Göttingen*, vol. I, 1843, “Optical Investigations” (translated by Bernard Rosett), July 1941. (Photocopies of a typed translation.)
188. R. Kingslake, “Basic Geometrical Optics,” *Applied Optics and Optical Engineering*, vol. I, R. Kingslake (ed.), Academic, New York, 1965, chap. 6, p. 214.
189. W. B. Wetherell, “Afocal Lenses,” *Applied Optics and Optical Engineering*, vol. 10, R. R. Shannon and J. C. Wyant (eds.), Academic, New York, 1987, chap. 3, pp. 110–192.
190. W. B. Wetherell, “Afocal Lenses and Mirror Systems,” *Handbook of Optics*, vol. II, chap. 2, McGraw-Hill, New York, 2d ed., 1994.
191. D. S. Goodman, “Afocal Systems,” *Geometrical and Instrumental Optics*, vol. 25 of *Methods of Experimental Physics*, D. Malacara, ed., Academic, New York, 1988, pp. 132–142.
192. A. F. Möbius, “Entwicklung der Lehre von dioptrischen Bildern mit Huelfe der Collineations-Verwandschaft,” *Leipziger Berichte*, vol. 7, 1855 pp. 8–32. (Translation: “Development of the teaching (or science) of dioptric images by means of the collineation relationship.”)
193. J. P. C. Southall, *The Principles and Methods of Geometrical Optics*, Macmillan, London, 1910.
194. S. Czapski, *Theorie der Optischen Instrumente nach Abbe*, Verlag von Eduard Trewendt, Breslau, 1893.
195. P. K. L. Drude, *The Theory of Optics*, Longmans, Green, and Co., New York, 1901, Dover, New York, 1959. (1900 English translation, C. R. Mann, R. A. Millikan.)
196. E. Wandersleb, “Abbe’s Geometrical Image Formation of Optical Images,” *The Formation of Images in Optical Instruments*, Von Rohr, ed., London, 1920, pp. 83–124. (English translation, H. M. Stationary Office.)
197. H. Chrétien, *Calcul des Combinaison Optique*, Mason, Paris, 1980, pp. 25ff.
198. M. Born and E. Wolf, *Principles of Optics*, Pergamon, New York, 1980, sec. 4.3, “Projective Transformations (Collineation) with Axial Symmetry.”
199. J. P. C. Southall, *Mirrors, Prisms and Lenses*, Macmillan, New York, 1933.
200. M. A. Penna, R. R. Patterson, *Projective Geometry and Its Applications to Computer Graphics*, Prentice-Hall, Englewood Cliffs, N.J., 1986.
201. J. C. Chastang, “Oblique Imaging in the Paraxial Approximation,” 1990 OSA annual meeting, Boston.
202. T. Scheimpflug, “Die Herstellung von Karten und Platen auf photographischem Wege,” *Journal of the Royal Academy of Sciences*, Vienna, vol. 116, 1907, pp. 235–266.
203. J. Sasian, “Image Plane Tilt in Optical Systems,” *Optical Engineering*, vol. 31, 1992, pp. 527–532.
204. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986.
205. R. A. Samson, “A Continuation of Gauss’s *Dioptrische Untersuchungen*,” *Proc. London Mathematical Society*, vol. 29, 1897, pp. 33–83.
206. G. G. Leathern, *The Elementary Theory of the Symmetrical Optical Instrument*, Hafner, New York, 1908.
207. K. Halbach, “Matrix Representation of Gaussian Optics,” *Am. J. Phys.*, vol. 32, 1964, pp. 90–108.
208. D. C. Sinclair, “The Specification of Optical Systems by Paraxial Matrices,” *Applications of Geometrical Optics*, W. J. Smith (ed.), *SPIE Proceedings*, vol. 39, SPIE, Bellingham, Wash., 1973, pp. 141–149.
209. E. L. O’Neil, *Introduction to Statistical Optics*, Addison-Wesley, Reading, Mass., 1963.
210. W. Brouwer, *Matrix Methods in Optical Instrument Design*, W. A. Benjamin, New York, 1964.
211. J. W. Blaker, *Geometrical Optics: The Matrix Theory*, Marcel Dekker, New York, 1971.
212. R. A. Longhurst, *Geometrical and Physical Optics*, Longman, London, 1973.
213. A. Gerrard and J. M. Burch, *Introduction of Matrix Methods in Optics*, Wiley, London, New York, 1974.
214. A. Nussbaum and R. A. Phillips, *Contemporary Optics for Scientists and Engineers*, Prentice-Hall, Englewood Cliffs, N.J., 1976.
215. H. Kogelnik, “Propagation of Laser Beams,” *Applied Optics and Optical Engineering*, vol. 7, R. R. Shannon and J. C. Wyant (eds.), Academic, New York, 1979, chap. 6, pp. 156–190.

216. M. V. Klein, T. E. Furtak, *Optics*, 2d ed., Wiley, New York, 1986.
217. K. D. Moller, *Optics*, University Science, 1988.
218. R. Guenther, *Modern Optics*, Wiley, New York, 1990.
219. A. Marechal, "Optique Geometrique General," *Fundamentals of Optics*, S. Fluegge (ed.), *Handbuch der Physik*, vol. 24, Springer-Verlag, Berlin, 1956, pp. 44–170.
220. A. Marechal, "Optique Geometrique et Algebre des Matrices," *Application de L'Algebre Moderne a Quelque Probleme de Physique Classique*, M. Parodi, Gauthier-Villars, Paris, 1961, chap. 15, pp. 306–346.
221. B. Macukow, H. H. Arsenault, "Matrix Decompositions for Nonsymmetrical Optical Systems," *J. Opt. Soc. Am.*, vol. 73, 1983, pp. 1360–1366.
222. W. Brouwer and A. Walther, "Geometrical Optics," *Advanced Optical Techniques*, A. C. S. Van Heel, (ed.), North-Holland, Amsterdam, 1967, chap. 16, pp. 505–570.
223. H. H. Arsenault, "The Rotation of Light Fans by Cylindrical Lenses," *Opt. Commun.*, vol. 31, 1979, pp. 275–278.
224. H. H. Arsenault, "Generalization of the Principal Plane Concept in Matrix Optics," *Am. J. Phys.*, vol. 48, 1980, pp. 397–399.
225. H. H. Arsenault, "A Matrix Representation for Non-Symmetrical Optical Systems," *J. Opt.*, vol. 11, Paris, 1980, pp. 87–91.
226. M. P. Keating, "A System Matrix for Astigmatic Optical Systems. Introduction and Dioptric Power Relations," *Am. J. Optometry and Physiological Optics*, vol. 58, 1981, pp. 810–819.
227. H. H. Arsenault and B. Macukow, "Factorization of the Transfer Matrix for Symmetrical Optical Systems," *J. Opt. Soc. Am.*, vol. 73, 1983, pp. 1350–1359.
228. A. E. Attard, "Matrix Optical Analysis of Skew Rays in Mixed Systems of Spherical and Orthogonal Cylindrical Lenses," *Appl. Opt.*, vol. 23, 1984, pp. 2706–2709.
229. B. D. Stone and G. W. Forbes, "Characterization of First-Order Optical Properties for Asymmetric Systems," *J. Opt. Soc. Am. A*, vol. 9, 1992, pp. 478–489.
230. L. Rayleigh, "Notes, Chiefly Historical, on Some Fundamental Propositions in Optics," *Phil. Mag.*, vol. 21, 1886, pp. 466–476. Collected works, vol. II, Dover, New York, 1964, article 137, pp. 513–521.
231. J. P. C. Southall, *Principles and Methods of Geometrical Optics*, Macmillan, London, 1910.
232. M. R. Hatch and D. E. Stoltzman, "The F-Stops Here," *Optical Spectra*, June 1980, pp. 88–91.
233. D. S. Goodman, "The F-Word in Optics," *Optics and Photonics News*, vol. 4, April 1993, pp. 38–39.
234. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1992.
235. R. Kingslake, "The Effective Aperture of a Photographic Objective," *J. Opt. Soc. Am.*, vol. 32, 1945, pp. 518–520.
236. M. Reiss, "The Cos⁴ Law of Illumination," *J. Opt. Soc. Am.*, vol. 35, 1945, pp. 283–288.
237. J. C. Gardner, "Validity of the Cosine-Fourth-Power Law of Illumination," *Journal of Research of the National Bureau of Standards*, vol. 39, 1947, pp. 213–219.
238. M. Reiss, "Notes on the Cos⁴ Law of Illumination," *J. Opt. Soc. Am.*, vol. 38, 1948, pp. 980–986.
239. R. Kingslake, "Illumination in Optical Instruments," *Applied Optics and Optical Engineering*, vol. 2, Academic, New York, 1965, chap. 5, pp. 195–228.
240. G. Slyusarev, "L'Eclaircissement de L'Image Formée par les Objectifs Photographiques Grand-Angulaires," *Journal of Physics (USSR)*, vol. 4, 1941, pp. 537–545.
241. D. Gabor, "Light and Information," *Progress in Optics*, E. Wolf, ed., vol. 1, 1961.
242. S. F. Ray, *Applied Photographic Optics*, Focal Press, London, 1988, chap. 22, "Depth of Field and Depth of Focus."
243. R. Kingslake, *Lenses in Photography*, Barnes, New York, 1963, SPIE, Bellingham WA, 1992.
244. A. E. Conrady, *Applied Optics and Optical Design*, vol. 1, 2, A. E. Conrady, Dover, New York, 1929.
245. H. H. Hopkins, *Wave Theory of Aberrations*, Clarendon Press, Oxford, 1950.
246. H. A. Buchdahl, *Optical Aberration Coefficients*, Oxford, London, 1954.
247. M. Herzberger, *Modern Geometrical Optics*, Interscience Publishers, New York, 1958.
248. R. Kingslake, *Lens Design Fundamentals*, Academic, New York, 1978.
249. M. Born and E. Wolf, *Principles of Optics*, Pergamon, New York, 1980.

250. G. C. Slyusarev, *Aberration and Optical Design Theory*, Hilger, Bristol, 1984.
251. W. T. Welford, *Aberrations of the Symmetrical Optical Systems*, Academic Press, London, 1974.
252. W. T. Welford, *Aberrations of Optical Systems*, Hilger, Bristol, 1986.
253. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1992.
254. D. O'Shea, "Aberration Curves in Lens Design," *Handbook of Optics*, vol. I, chap. 33, McGraw-Hill, New York, 2d ed., 1994.
255. F. Zernike, "Beugungstheorie des Schneidenver-Eahrens und Seiner Verbesserten Form, der Phasenkontrastmethode," *Physica*, vol. 1, 1934, p. 689.
256. C.-J. Kim, R. R. Shannon, "Catalogue of Zernike Polynomials," *Applied Optics and Optical Engineering*, vol. 10, R. R. Shannon, J. C. Wyant (eds.), Academic, New York, 1987, chap. 4, pp. 193–233.
257. V. N. Mahajan, "Zernike Polynomials and Wavefront Fitting," D. Malacara (ed.), 3d ed. *Optical Shop Testing*, Wiley, New York, 2007, pp. 498–546.

PART

2

PHYSICAL OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

INTERFERENCE

John E. Greivenkamp

College of Optical Sciences
University of Arizona
Tucson, Arizona

2.1 GLOSSARY

A	amplitude
E	electric field vector
\mathbf{r}	position vector
x, y, z	rectangular coordinates
ϕ	phase

2.2 INTRODUCTION

Interference results from the superposition of two or more electromagnetic waves. From a classical optics perspective, interference is the mechanism by which light interacts with light. Other phenomena, such as refraction, scattering, and diffraction, describe how light interacts with its physical environment. Historically, interference was instrumental in establishing the wave nature of light. The earliest observations were of colored fringe patterns in thin films. Using the wavelength of light as a scale, interference continues to be of great practical importance in areas such as spectroscopy and metrology.

2.3 WAVES AND WAVEFRONTS

The *electric field vector* due to an electromagnetic field at a point in space is composed of an amplitude and a phase

$$\mathbf{E}(x, y, z, t) = \mathbf{A}(x, y, z, t)e^{i\phi(x, y, z, t)} \quad (1)$$

or

$$\mathbf{E}(\mathbf{r}, t) = \mathbf{A}(\mathbf{r}, t)e^{i\phi(\mathbf{r}, t)} \quad (2)$$

where \mathbf{r} is the position vector and both the amplitude \mathbf{A} and phase ϕ are functions of the spatial coordinate and time. As described in Chap. 12, "Polarization," the polarization state of the field is contained in the temporal variations in the amplitude vector.

This expression can be simplified if a linearly polarized monochromatic wave is assumed:

$$\mathbf{E}(x, y, z, t) = \mathbf{A}(x, y, z, t) e^{i[\omega t - \phi(x, y, z)]} \quad (3)$$

where ω is the angular frequency in radians per second and is related to the frequency ν by

$$\omega = 2\pi\nu \quad (4)$$

Some typical values for the optical frequency are 5×10^{14} Hz for the visible, 10^{13} Hz for the infrared, and 10^{16} Hz for the ultraviolet. Note that in the expression for the electric field vector, the time dependence has been eliminated from the amplitude term to indicate a constant linear polarization. The phase term has been split into spatial and temporal terms. At all locations in space, the field varies harmonically at the frequency ω .

Plane Wave

The simplest example of an electromagnetic wave is the *plane wave*. The plane wave is produced by a monochromatic point source at infinity and is approximated by a collimated light source. The complex amplitude of a linearly polarized plane wave is

$$\mathbf{E}(x, y, z, t) = \mathbf{E}(\mathbf{r}, t) = \mathbf{A} e^{i[\omega t - \mathbf{k} \cdot \mathbf{r}]} \quad (5)$$

where \mathbf{k} is the wave vector. The wave vector points in the direction of propagation, and its magnitude is the wave number $k = 2\pi/\lambda$, where λ is the wavelength. The wavelength is related to the temporal frequency by the speed of light v in the medium:

$$\lambda = \frac{v}{\nu} = 2\pi \frac{v}{\omega} = \frac{c}{n\nu} = 2\pi \frac{c}{n\omega} \quad (6)$$

where n is the index of refraction, and c is the speed of light in a vacuum. The amplitude \mathbf{A} of a plane wave is a constant over all space, and the plane wave is clearly an idealization.

If the direction of propagation is parallel to the z axis, the expression for the complex amplitude of the plane wave simplifies to

$$\mathbf{E}(x, y, z, t) = \mathbf{A} e^{i[\omega t - kz]} \quad (7)$$

We see that the plane wave is periodic in both space and time. The spatial period equals the wavelength in the medium, and the temporal period equals $1/\nu$. Note that the wavelength changes with index of refraction, and the frequency is independent of the medium.

Spherical Wave

The second special case of an electromagnetic wave is the spherical wave which radiates from an isotropic point source. If the source is located at the origin, the complex amplitude is

$$E(r, t) = (A/r) e^{i[\omega t - kr]} \quad (8)$$

where $r = (x^2 + y^2 + z^2)^{1/2}$. The field is spherically symmetric and varies harmonically with time and the radial distance. The radial period is the wavelength in the medium. The amplitude of the field decreases as $1/r$ for energy conservation. At a large distance from the source, the spherical wave can be approximated by a plane wave. Note that the vector characteristics of the field (its polarization) are not considered here as it is not possible to describe a linear polarization pattern

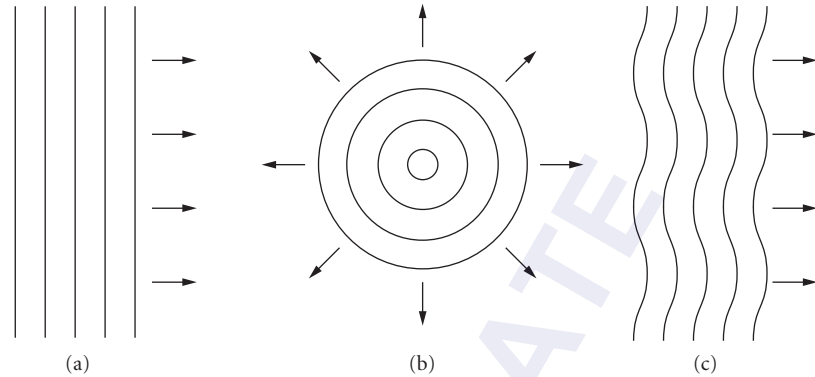


FIGURE 1 Examples of wavefronts: (a) plane wave; (b) spherical wave; and (c) aberrated plane wave.

of constant amplitude that is consistent over the entire surface of a sphere. In practice, we only need to consider an angular segment of a spherical wave, in which case this polarization concern disappears.

Wavefronts

Wavefronts represent surfaces of constant phase for the electromagnetic field. Since they are normally used to show the spatial variations of the field, they are drawn or computed at a fixed time. Wavefronts for plane and spherical waves are shown in Fig. 1a and b. The field is periodic, and a given value of phase will result in multiple surfaces. These surfaces are separated by the wavelength. A given wavefront also represents a surface of constant optical path length (OPL) from the source. The OPL is defined by the following path integral:

$$\text{OPL} = \int_S^P n(s) ds \quad (9)$$

where the integral goes from the source S to the observation point P , and $n(s)$ is the index of refraction along the path. Variations in the index or path can result in irregularities or aberrations in the wavefront. An aberrated plane wavefront is shown in Fig. 1c. Note that the wavefronts are still separated by the wavelength.

The local normal to the wavefront defines the propagation direction of the field. This fact provides the connection between wave optics and ray or geometrical optics. For a given wavefront, a set of rays can be defined using the local surface normals. In a similar manner, a set of rays can be used to construct the equivalent wavefront.

2.4 INTERFERENCE

The net complex amplitude is the sum of all of the component fields,

$$\mathbf{E}(x, y, z, t) = \sum_i \mathbf{E}_i(x, y, z, t) \quad (10)$$

and the resulting field intensity is the time average of the modulus squared of the total complex amplitude

$$I(x, y, z, t) = \langle |\mathbf{E}(x, y, z, t)|^2 \rangle \quad (11)$$

where $\langle \rangle$ indicates a time average over a period much longer than $1/\nu$. If we restrict ourselves to two interfering waves \mathbf{E}_1 and \mathbf{E}_2 , this result simplifies to

$$I(x, y, z, t) = \langle |\mathbf{E}_1|^2 \rangle + \langle |\mathbf{E}_2|^2 \rangle + \langle \mathbf{E}_1 \cdot \mathbf{E}_2^* \rangle + \langle \mathbf{E}_1^* \cdot \mathbf{E}_2 \rangle \quad (12)$$

or

$$I(x, y, z, t) = I_1 + I_2 + \langle \mathbf{E}_1 \cdot \mathbf{E}_2^* \rangle + \langle \mathbf{E}_1^* \cdot \mathbf{E}_2 \rangle \quad (13)$$

where I_1 and I_2 are the intensities due to the two beams individually, and the (x, y, z, t) dependence is now implied for the various terms.

This general result can be greatly simplified if we assume linearly polarized monochromatic waves of the form in Eq. (3):

$$\mathbf{E}_i(x, y, z, t) = \mathbf{A}_i(x, y, z, t) e^{i[\omega_i t - \phi_i(x, y, z)]} \quad (14)$$

The resulting field intensity is

$$I(x, y, z, t) = I_1 + I_2 + 2(\mathbf{A}_1 \cdot \mathbf{A}_2) \cos[(\omega_1 - \omega_2)t - (\phi_1(x, y, z) - \phi_2(x, y, z))] \quad (15)$$

The interference effects are contained in the third term, and we can draw two important conclusions from this result. First, if the two interfering waves are orthogonally polarized, there will be no visible interference effects, as the dot product will produce a zero coefficient. Second, if the frequencies of the two waves are different, the interference effects will be modulated at a temporal beat frequency equal to the difference frequency.

Interference Fringes

We will now add the additional restrictions that the two linear polarizations are parallel and that the two waves are at the same optical frequency. The expression for the intensity pattern now becomes

$$I(x, y, z) = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos[\Delta\phi(x, y, z)] \quad (16)$$

where $\Delta\phi = \phi_1 - \phi_2$ is the phase difference. This is the basic equation describing interference. The detected intensity varies cosinusoidally with the phase difference between the two waves as shown in Fig. 2. These alternating bright and dark bands in the intensity pattern

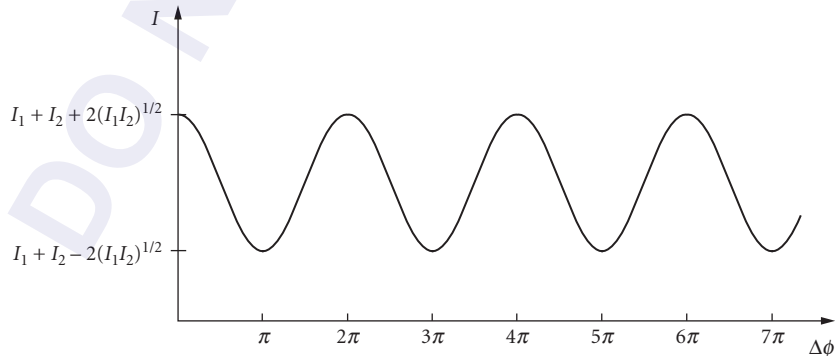


FIGURE 2 The variation in intensity as a function of the phase difference between two interfering waves.

TABLE 1 The Phase Difference and OPD for Bright and Dark Fringes (m an Integer)

	$\Delta\phi$	OPD
Bright fringe	$2m\pi$	$m\lambda$
Dark fringe	$2(m+1)\pi$	$(m+1/2)\lambda$

are referred to as *interference fringes*, and along a particular fringe, the phase difference is constant.

The phase difference is related to the difference in the optical path lengths between the source and the observation point for the two waves. This is the *optical path difference* (OPD):

$$\text{OPD} = \text{OPL}_1 - \text{OPL}_2 = \left(\frac{\lambda}{2\pi}\right) \Delta\phi \quad (17)$$

or

$$\Delta\phi = \left(\frac{2\pi}{\lambda}\right) \text{OPD} \quad (18)$$

The phase difference changes by 2π every time the OPD increases by a wavelength. The OPD is therefore constant along a fringe.

Constructive interference occurs when the two waves are in phase, and a bright fringe or maximum in the intensity pattern results. This corresponds to a phase difference of an integral number of 2π 's or an OPD that is a multiple of the wavelength. A dark fringe or minimum in the intensity pattern results from *destructive interference* when the two waves are out of phase by π or the OPD is an odd number of half wavelengths. These results are summarized in Table 1. For conditions between these values, an intermediate value of the intensity results. Since both the OPD and the phase difference increase with the integer m , the absolute value of m is called the *order of interference*.

As we move from one bright fringe to an adjacent bright fringe, the phase difference changes by 2π . Each fringe period corresponds to a change in the OPD of a single wavelength. It is this inherent precision that makes interferometry such a valuable metrology tool. The wavelength of light is used as the unit of measurement. Interferometers can be configured to measure small variations in distance, index, or wavelength.

When two monochromatic waves are interfered, the interference fringes exist not only in the plane of observation, but throughout all space. This can easily be seen from Eq. (16) where the phase difference can be evaluated at any z position. In many cases, the observation of interference is confined to a plane, and this plane is usually assumed to be perpendicular to the z axis. The z dependence in Eq. (16) is therefore often not stated explicitly, but it is important to remember that interference effects will exist in other planes.

Fringe Visibility

It is often more convenient to rewrite Eq. (16) as

$$I(x, y) = I_0(x, y) \{1 + \gamma(x, y) \cos[\Delta\phi(x, y, z)]\} \quad (19)$$

or

$$I(x, y) = I_0(x, y) \{1 + \gamma(x, y) \cos[2\pi \text{OPD}(x, y)/\lambda]\} \quad (20)$$

where $I_0(x, y) = I_1(x, y) + I_2(x, y)$, and

$$\gamma(x, y) = \frac{2[I_1(x, y)I_2(x, y)]^{1/2}}{I_1(x, y) + I_2(x, y)} \quad (21)$$

Since the cosine averages to zero, $I_0(x, y)$ represents the average intensity, and $\gamma(x, y)$ is the local *fringe contrast* or *visibility*. The fringe visibility can also be equivalently calculated using the standard formula for modulation:

$$\gamma(x, y) = \frac{I_{\max}(x, y) - I_{\min}(x, y)}{I_{\max}(x, y) + I_{\min}(x, y)} \quad (22)$$

where I_{\max} and I_{\min} are the maximum and minimum intensities in the fringe pattern.

The fringe visibility will have a value between 0 and 1. The maximum visibility will occur when the two waves have equal intensity. Not surprisingly, the visibility will drop to zero when one of the waves has zero intensity. In general, the intensities of the two waves can vary with position, so that the average intensity and fringe visibility can also vary across the fringe pattern. The average intensity in the observation plane equals the sum of the individual intensities of the two interfering waves. The interference term redistributes this energy into bright and dark fringes.

Two Plane Waves

The first special case to consider is the interference of two plane waves of equal intensity, polarization and frequency. They are incident at angles θ_1 and θ_2 on the observation plane, as shown in Fig. 3. The plane of incidence is the x - z plane (the two \mathbf{k} -vectors are contained in this plane). According to Eq. (5), the complex amplitude for each of these plane waves is

$$\mathbf{E}_i(x, y, z, t) = \mathbf{A}e^{i(\omega t - kz \cos(\theta_i) - kx \sin(\theta_i))} \quad (23)$$

where the dot product has been evaluated. For simplicity we will place the observation plane at $z = 0$, and the phase difference between the two waves is

$$\Delta\phi(x, y) = kx(\sin\theta_1 - \sin\theta_2) = (2\pi x/\lambda)(\sin\theta_1 - \sin\theta_2) \quad (24)$$

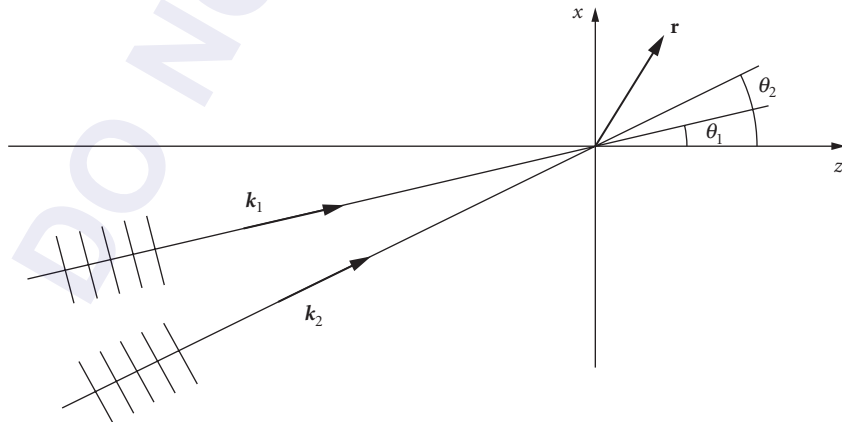


FIGURE 3 The geometry for the interference of two plane waves.

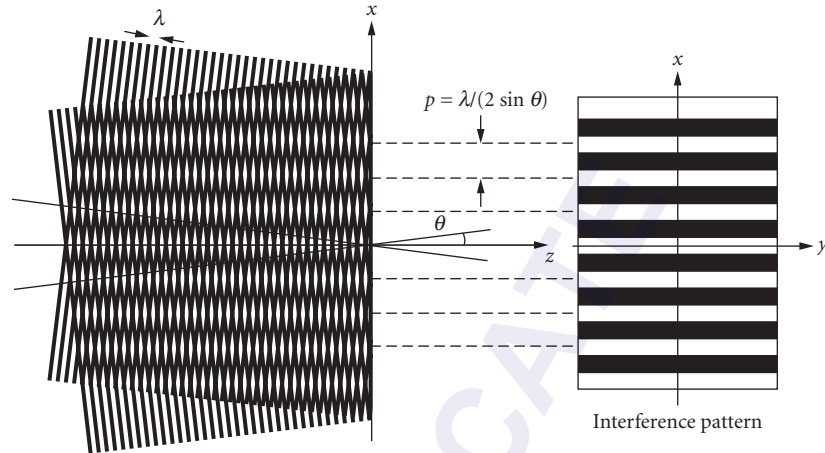


FIGURE 4 The interference of plane waves incident at $\pm\theta$ resulting in straight fringes.

The resulting intensity from Eq. (19) is

$$I(x, y) = I_0 \{1 + \cos[(2\pi x/\lambda)(\sin\theta_1 - \sin\theta_2)]\} \quad (25)$$

where $I_0 = 2A^2$ is twice the intensity of each of the individual waves. Straight equispaced fringes are produced. The fringes are parallel to the y axis, and the fringe period depends on the angle between the two interfering beams.

The fringe period p is

$$p = \frac{\lambda}{\sin\theta_1 - \sin\theta_2} \quad (26)$$

and this result can also be obtained by noting that a bright fringe will occur whenever the phase difference equals a multiple of 2π . A typical situation for interference is that the two angles of incidence are equal and opposite, $\theta_1 = -\theta_2 = \theta$. The angle between the two beams is 2θ . Under this condition, the period is

$$p = \frac{\lambda}{2\sin\theta} \approx \frac{\lambda}{2\theta} \quad (27)$$

and the small-angle approximation is given. As the angle between the beams gets larger, the period decreases. For example, the period is 3.8λ at 15° (full angle of 30°) and is λ at 30° (full angle of 60°). The interference of two plane waves can be visualized by looking at the overlap or moiré of two wavefront patterns (Fig. 4). Whenever the lines representing the wavefronts overlap, a fringe will result. This description also clearly shows that the fringes extend parallel to the z axis and exist everywhere the two beams overlap.

Plane Wave and Spherical Wave

A second useful example to consider is the interference of a plane wave and a spherical wave. Once again the two waves have the same frequency. The plane wave is at normal incidence, the spherical wave is due to a source at the origin, and the observation plane is located at $z = R$. The wavefront shape at the observation plane will be a spherical shell of radius R .

Starting with Eq. (8), the complex amplitude of the spherical wave in the observation plane is

$$E(\rho, t) = (A/R)e^{i[\omega t - k(R^2 + \rho^2)^{1/2}]} \approx (A/R)e^{i[\omega t - k(R + \rho^2/2R)]} \quad (28)$$

where $\rho = (x^2 + y^2)^{1/2}$, and the square root has been expanded in the second expression. This expansion approximates the spherical wave by a parabolic wave with the same vertex radius. An additional assumption is that the amplitude of the field A/R is constant over the region of interest. The field for the plane wave is found by evaluating Eq. (23) at $z = R$ and $\theta = 0$. The phase difference between the plane and the sphere is then

$$\Delta\phi(\rho) \approx \frac{\pi\rho^2}{\lambda R} \quad (29)$$

and the resulting intensity pattern is

$$I(\rho) = I_0 \left[1 + \cos\left(\frac{\pi\rho^2}{\lambda R}\right) \right] \quad (30)$$

The fringe pattern comprises concentric circles, and the radial fringe spacing decreases as the radius ρ increases. The intensities of the two waves have been assumed to be equal at the observation plane. This result is valid only when ρ is much smaller than R .

The radius of the m th bright fringe can be found by setting $\Delta\phi = 2\pi m$:

$$\rho_m = \sqrt{2m\lambda R} \quad (31)$$

where m is an integer. The order of interference m increases with radius. Figure 5 shows a visualization of this situation using wavefronts. This fringe pattern is the Newton's ring pattern and is discussed in more detail later, under "Fizeau Interferometer." This picture also shows that the radii of the fringes increase as the square root of R .

The analysis of the spherical wave could also have been done by using the sag of a spherical wavefront to produce an OPD and then converting this value to a phase difference. The quadratic approximation for the sag of a spherical surface is $\rho^2/2R$. This corresponds to the OPD

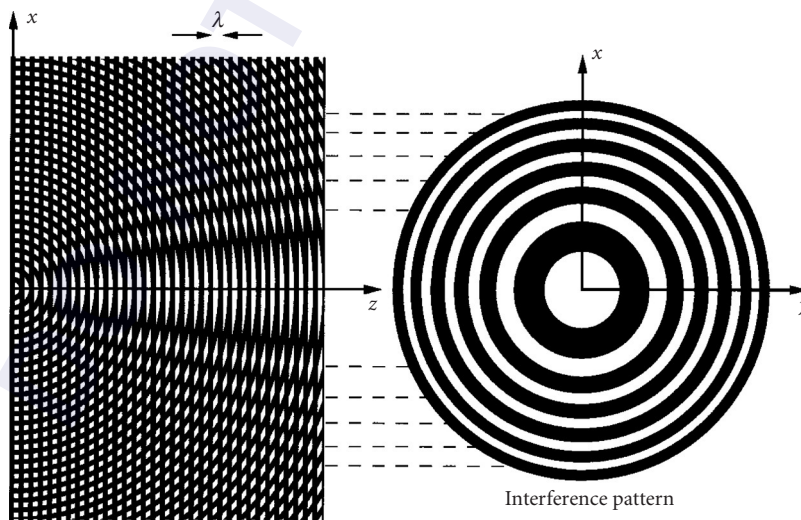


FIGURE 5 The interference of a plane wave and a spherical wave.

between the spherical and planar wavefronts. The equivalent phase difference [Eq. (18)] is then $\pi\rho^2/\lambda R$, as before.

Two Spherical Waves

When considering two spherical waves, there are two particular geometries that we want to examine. The first places the observation plane perpendicular to a line connecting the two sources, and the second has the observation plane parallel to this line. Once again, the sources are at the same frequency.

When the observations are made on a plane perpendicular to a line connecting the two sources, we can use Eq. (28) to determine the complex amplitude of the two waves:

$$E_i(\rho, t) \approx (A/R)e^{i[\omega t - k(R_i + \rho^2/2R_i)]} \quad (32)$$

Let $d = R_1 - R_2$ be the separation of the two sources. For simplicity, we have also assumed that the amplitudes of the two waves are equal (R is an average distance). The phase difference between the two waves is

$$\Delta\phi = \left(\frac{\pi\rho^2}{\lambda}\right)\left(\frac{1}{R_1} - \frac{1}{R_2}\right) + \frac{2\pi d}{\lambda} \approx \frac{2\pi d}{\lambda} - \left(\frac{\pi\rho^2}{\lambda}\right)\left(\frac{d}{R^2}\right) \quad (33)$$

where the approximation $R_1 R_2 \approx R^2$ has been made. There are two terms to this phase difference. The second is a quadratic phase term identical in form to the result obtained from spherical and plane waves. The pattern will be symmetric around the line connecting the two sources, and its appearance will be similar to Newton's rings. The equivalent radius of the spherical wave in Eq. (29) is R^2/d . The first term is a constant phase shift related to the separation of the two sources. If this term is not a multiple of 2π , the center of the fringe pattern will not be a bright fringe; if the term is π , the center of the pattern will be dark. Except for the additional phase shift, this intensity pattern is not distinguishable from the result in the previous section. It should be noted, however, that a relative phase shift can be introduced between a spherical wave and a plane wave to obtain this same result.

An important difference between this pattern and the Newton's ring pattern is that the order of interference ($|m|$ defined by $\Delta\phi = 2\pi m$) or phase difference is a maximum at the center of the pattern and decreases with radius. The Newton's ring pattern formed between a plane and a spherical wave has a minimum order of interference at the center of the pattern. This distinction is important when using polychromatic sources.

There are several ways to analyze the pattern that is produced on a plane that is parallel to a line connecting the two sources. We could evaluate the complex amplitudes by using Eq. (28) and moving the center of the spherical waves to $\pm d/2$ for the two sources. An equivalent method is to compare the wavefronts at the observation plane. This is shown in Fig. 6. The OPD between the two wavefronts is

$$\text{OPD}(x, y) = \frac{[(x+d/2)^2 + y^2]}{2L} - \frac{[(x-d/2)^2 + y^2]}{2L} \quad (34)$$

where the quadratic approximation for the wavefront sag has been assumed, and L is the distance between the sources and the observation plane. After simplification, the OPD and phase differences are

$$\text{OPD}(x, y) = \frac{xd}{L} \quad (35)$$

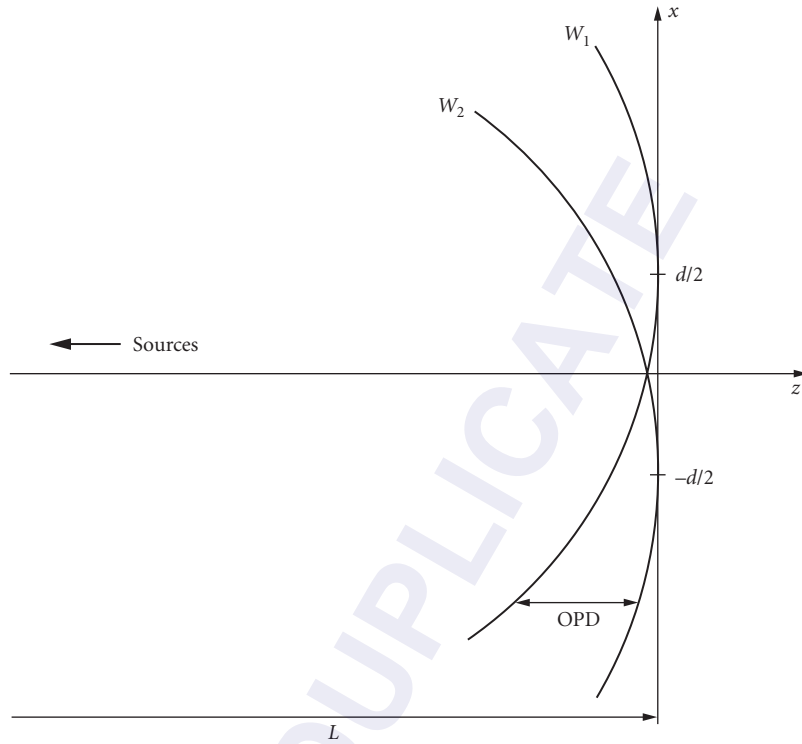


FIGURE 6 The interference of two spherical waves on a plane parallel to the sources.

and

$$\Delta\phi(x, y) = \frac{2\pi xd}{\lambda L} \tag{36}$$

Straight equispaced fringes parallel to the y axis are produced. The period of the fringes is $\lambda L/d$. This fringe pattern is the same as that produced by two plane waves. Note that these fringes increase in spacing as the distance from the sources increases. The approximations used require that L be much larger than ρ and d .

Figure 7 shows the creation of the fringe patterns for two point sources. The full three-dimensional pattern is a series of nested hyperboloids symmetric about the line connecting the sources. Above the two sources, circular fringes approximating Newton's rings are produced, and perpendicular to the sources, the fringes appear to be straight and equispaced. The actual appearance of these patterns is modified by the approximations used in the derivations, and as a result, these two specific patterns have limited lateral extent.

Aberrated Wavefronts

When an aberrated or irregularly shaped wavefront is interfered with a reference wavefront, an irregularly shaped fringe pattern is produced. However, the rules for analyzing this pattern are the same as with any two wavefronts. A given fringe represents a contour of constant OPD or phase difference between the two wavefronts. Adjacent fringes differ in OPD by one wavelength

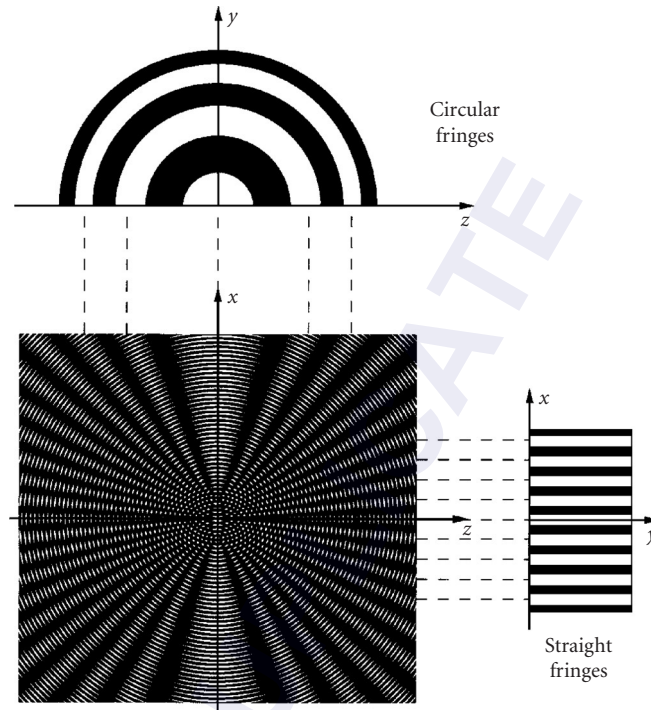


FIGURE 7 The interference of two spherical waves.

or equivalently correspond to a phase difference of 2π . If the reference is a plane wave, the absolute shape of the irregular wavefront is obtained. If the reference is a spherical wave, or another aberrated wave, the measured OPD or phase difference map represents the difference between the two wavefronts.

Temporal Beats

In Eq. (15) it was noted that if the waves are at different frequencies, the interference effects are modulated by a beat frequency. Rewriting this expression assuming equal-intensity parallel-polarized beams produces

$$I(x, y, t) = I_0 \{1 + \cos[2\pi \Delta \nu t - \Delta \phi(x, y)]\} \quad (37)$$

where $\Delta \nu = \nu_1 - \nu_2$. The intensity at a given location will now vary sinusoidally with time at the beat frequency $\Delta \nu$. The phase difference $\Delta \phi$ appears as a spatially varying phase shift of the beat frequency. This is the basis of the heterodyne technique used in a number of interferometers. It is commonly used in distance-measuring interferometers.

In order for a heterodyne system to work, there must be a phase relationship between the two sources even though they are at different frequencies. One common method for obtaining this is accomplished by starting with a single source, splitting it into two beams, and frequency-shifting one beam with a known Doppler shift. The system will also work in reverse; measure the interferometric beat frequency to determine the velocity of the object producing the Doppler shift.

Coherence

Throughout this discussion of fringe patterns, we have assumed that the two sources producing the two waves have the same frequency. In practice, this requires that both sources be derived from a single source. Even when two different frequencies are used [Eq. (37)] there must be an absolute phase relation between the two sources. If the source has finite size, it is considered to be composed of a number of spatially separated, independently radiating point sources. If the source has a finite spectral bandwidth, it is considered to be composed of a number of spatially coincident point sources with different frequencies. These reductions in the spatial or temporal coherence of the source will decrease the visibility of the fringes at different locations in space. This is referred to as *fringe localization*. These effects will be discussed later in this chapter and also in Chap. 5, “Coherence Theory.”

There are two general methods to produce mutually coherent waves for interference. The first is called *wavefront division*, where different points on a wavefront are sampled to produce two new wavefronts. The second is *amplitude division*, where some sort of beamsplitter is used to divide the wavefront at a given location into two separate wavefronts. These methods are discussed in the next sections.

2.5 INTERFERENCE BY WAVEFRONT DIVISION

Along a given wavefront produced by a monochromatic point source, the wavefront phase is constant. If two parts of this wavefront are selected and then redirected to a common volume in space, interference will result. This is the basis for *interference by wavefront division*.

Young’s Double-Slit Experiment

In 1801, Thomas Young performed a fundamental experiment for demonstrating interference and the wave nature of light. Monochromatic light from a single pinhole illuminates an opaque screen with two additional pinholes or slits. The light diffracts from these pinholes and illuminates a viewing screen at a distance large compared to the pinhole separation. Since the light illuminating the two pinholes comes from a single source, the two diffracted wavefronts are coherent and interference fringes form where the beams overlap.

In the area where the two diffracted beams overlap, they can be modeled as two spherical waves from two point sources, and we already know the form of the solution for the interference from our earlier discussion. Equispaced straight fringes are produced, and the period of the fringes is $\lambda L/d$, where L is the distance to the screen and d is the separation of the pinholes. The fringes are oriented perpendicular to the line connecting the two pinholes.

Even though we already know the answer, there is a classic geometric construction we should consider that easily gives the OPD between the two wavefronts at the viewing screen. This is shown in Fig. 8. S_0 illuminates both S_1 and S_2 and is equidistant from both slits. The OPD for an observation point P at an angle θ or position x is

$$\text{OPD} = \overline{S_2P} - \overline{S_1P} \quad (38)$$

We now draw a line from S_1 to B that is perpendicular to the second ray. Since L is much larger than d , the distances from B to P and S_1 to P are approximately equal. The OPD is then

$$\text{OPD} \approx \overline{S_2B} = d \sin \theta \approx d \theta \approx \frac{dx}{L} \quad (39)$$

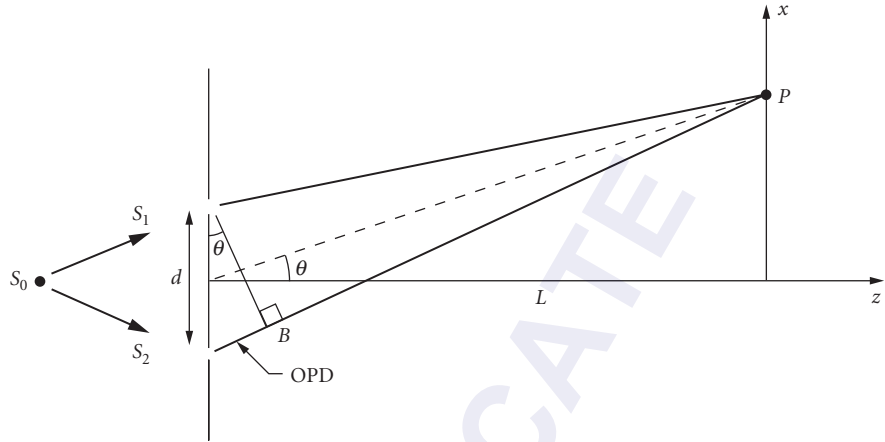


FIGURE 8 Young's double-slit experiment.

and constructive interference or a bright fringe occurs when the OPD is a multiple of the wavelength: $OPD \approx m\lambda$, where m is an integer. The condition for the m th order bright fringe is

$$\text{Bright fringe: } \sin(\theta) \approx \theta = \frac{m\lambda}{d} \quad \text{or} \quad x = \frac{m\lambda L}{d} \quad (40)$$

This construction is useful not only for interference situations, but also for diffraction analysis.

Effect of Slit Width

The light used to produce the interference pattern is diffracted by the pinholes or slits. Interference is possible only if light is directed in that direction. The overall interference intensity pattern is therefore modulated by the single-slit diffraction pattern (assuming slit apertures):

$$I(x) = I_0 \operatorname{sinc}^2\left(\frac{Dx}{\lambda L}\right) \left[1 + \gamma(x) \cos\left(\frac{2\pi xd}{\lambda L}\right) \right] \quad (41)$$

where D is the slit width, and a one-dimensional expression is shown. The definition of a sinc function is

$$\operatorname{sinc}(\alpha) = \frac{\sin(\pi\alpha)}{\pi\alpha} \quad (42)$$

where the zeros of the function occur when the argument α is an integer. The intensity variation in the y direction is due to diffraction only and is not shown. Since the two slits are assumed to be illuminated by a single source, there are no coherence effects introduced by using a pinhole or slit of finite size.

The term $\gamma(x)$ is included in Eq. (41) to account for variations in the fringe visibility. These could be due to unequal illumination of the two slits, a phase difference of the light reaching the slits, or a lack of temporal or spatial coherence of the source S_0 .

Other Arrangements

Several other arrangements for producing interference by division of wavefront are shown in Fig. 9. They all use a single source and additional optical elements to produce two separate and mutually coherent sources. Fresnel's biprism and mirror produce the two virtual source images, Billet's split lens produces two real source images, and Lloyd's mirror produces a single virtual source image as a companion to the original source. Interference fringes form wherever the two resulting waves overlap (shaded regions). One significant difference between these arrangements and Young's two slits is that a large section of the initial wavefront is used instead of just two points. All of these systems are much more light efficient, and they do not rely on diffraction to produce the secondary wavefronts.

In the first three of these systems, a bright fringe is formed at the zero OPD point between the two sources as in the double-slit experiment. With Lloyd's mirror, however, the zero OPD point has a dark fringe. This is due to the π phase shift that is introduced into one of the beams on reflection from the mirror.

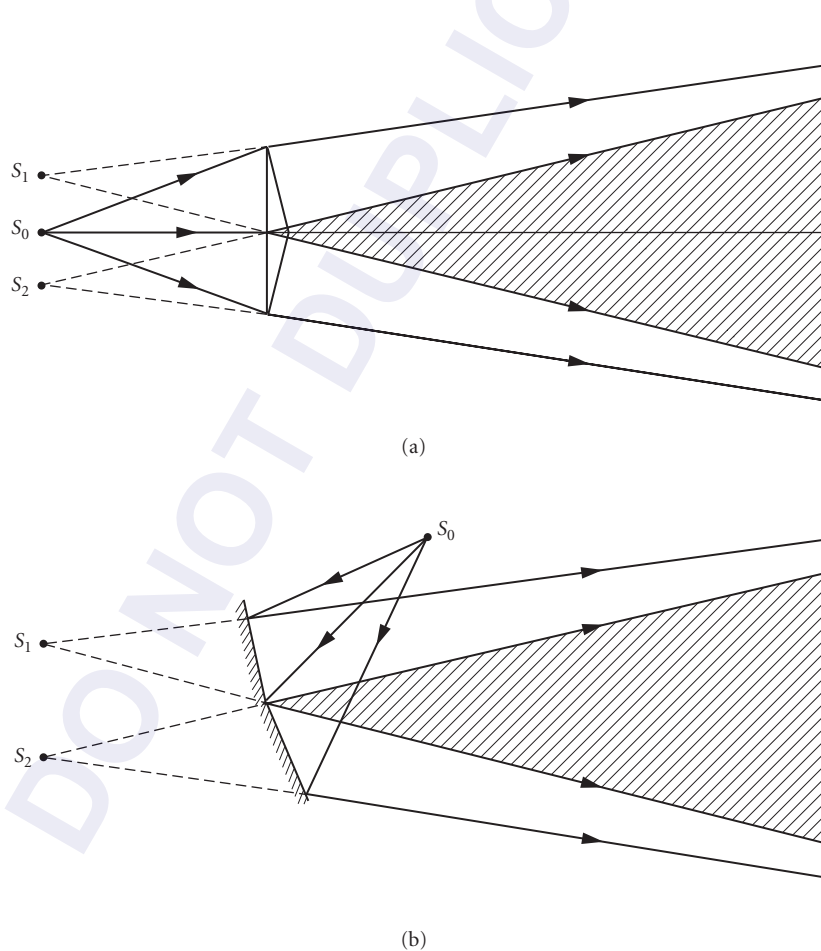


FIGURE 9 Arrangements for interference by division of wavefront: (a) Fresnel's biprism; (b) Fresnel's mirror; (c) Billet's split lens; and (d) Lloyd's mirror.

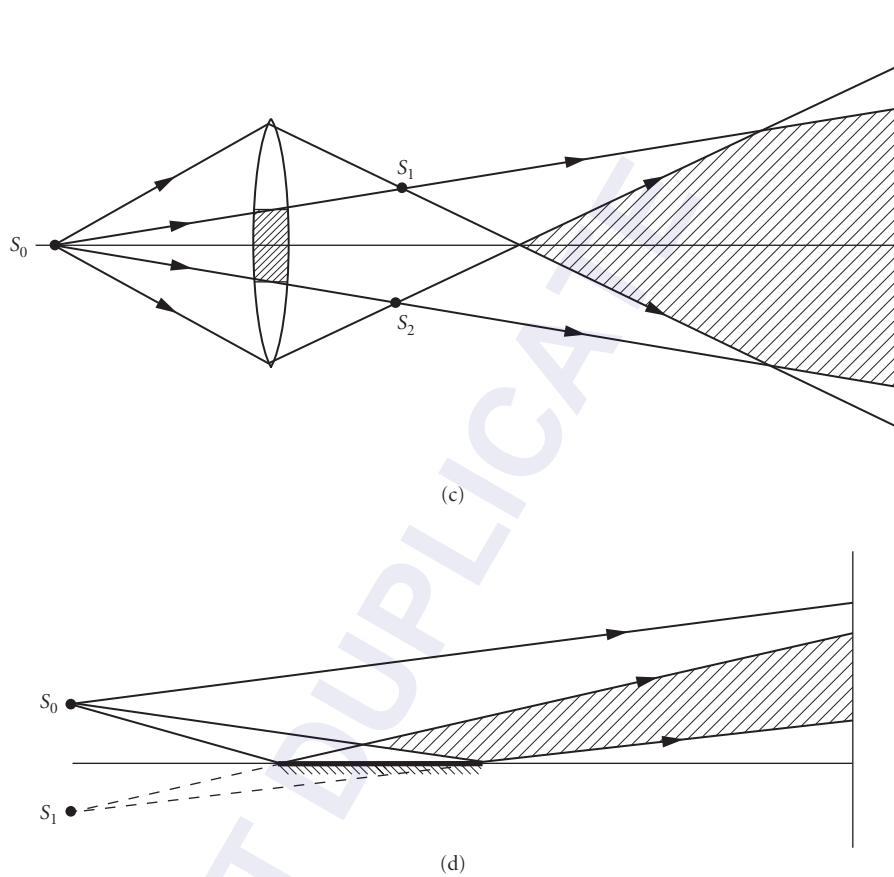


FIGURE 9 (Continued)

Source Spectrum

The simple fringe pattern produced by the two-slit experiment provides a good example to examine the effects of a source with a finite spectrum. In this model, the source can be considered to be a collection of sources, each radiating independently and at a different wavelength. All of these sources are collocated to produce a point source. (Note that this is an approximation, as a true point source must be monochromatic.) At each wavelength, an independent intensity pattern is produced:

$$I(x, \lambda) = I_0 \left[1 + \cos \left(\frac{2\pi x d}{\lambda L} \right) \right] = I_0 \left[1 + \cos \left(\frac{2\pi \text{OPD}}{\lambda} \right) \right] \quad (43)$$

where the period of the fringes is $\lambda L/d$, and a fringe visibility of one is assumed. The total intensity pattern is the sum of the individual fringe patterns:

$$I(x) = \int_0^\infty S(\lambda) I(x, \lambda) d\lambda = \int_0^\infty S(\nu) I(x, \nu) d\nu \quad (44)$$

where $S(\lambda)$ or $S(\nu)$ is the source intensity spectrum which serves as a weighting function.

The effect of this integration can be seen by looking at a simple example where the source is composed of three different wavelengths of equal intensity. To further aid in visualization, let's use Blue (400 nm), Green (500 nm), and Red (600 nm). The result is shown in Fig. 10a. There are three cosine patterns, each with a period proportional to the wavelength. The total intensity is the sum of these curves. All three curves line up when the OPD is zero ($x=0$), and the central bright fringe is now surrounded by two-colored dark fringes. These first dark fringes have a red to blue coloration with increasing OPD. As we get further away from the zero OPD condition, the three patterns get out of phase, the pattern washes out, and the color saturation decreases. This is especially true when the source is composed of more than three wavelengths.

It is common in white light interference situations for one of the two beams to undergo an additional π phase shift. This is the situation in Lloyd's mirror. In this case, there is a central dark fringe at zero OPD with colored bright fringes on both sides. This is shown in Fig. 10b, and the pattern is complementary to the previous pattern. In this case the first bright fringe shows a blue to red color smear. The dark central fringe is useful in determining the location of zero OPD between the two beams.

The overall intensity pattern and resulting fringe visibility can be computed for a source with a uniform frequency distribution over a frequency range of $\Delta\nu$:

$$I(x) = \frac{1}{\Delta\nu} \int_{\nu_0 - \Delta\nu/2}^{\nu_0 + \Delta\nu/2} I(x, \nu) d\nu = \frac{1}{\Delta\nu} \int_{\nu_0 - \Delta\nu/2}^{\nu_0 + \Delta\nu/2} I_0 \left[1 + \cos\left(\frac{2\pi\nu xd}{cL}\right) \right] d\nu \quad (45)$$

where ν_0 is the central frequency, and the $1/\Delta\nu$ term is a normalization factor to assure that the average intensity is I_0 . After integration and simplification, the result is

$$I(x) = I_0 \left[1 + \operatorname{sinc}\left(\frac{xd\Delta\nu}{cL}\right) \cos\left(\frac{2\pi\nu_0 xd}{cL}\right) \right] \quad (46)$$

where the sinc function is defined in Eq. (42). A fringe pattern due to the average optical frequency results, but it is modulated by a sinc function that depends on $\Delta\nu$ and x . The absolute value of the

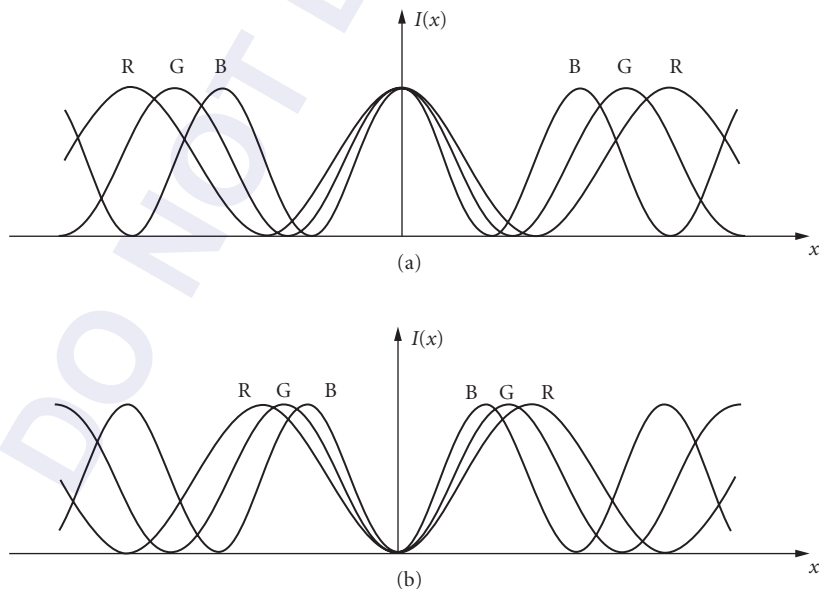


FIGURE 10 The interference pattern produced by a source with three separate wavelengths: (a) zero OPD produces a bright fringe and (b) zero OPD produces a dark fringe.

sinc function is the fringe visibility $\gamma(x)$, and it depends on both the spectral width and position of observation. The negative portions of the sinc function correspond to a π phase shift of the fringes.

It is informative to rewrite this expression in terms of the OPD:

$$I(x) = I_0 \left[1 + \operatorname{sinc} \left(\frac{\text{OPD} \Delta \nu}{c} \right) \cos \left(\frac{2\pi \text{OPD}}{\lambda_0} \right) \right] \quad (47)$$

where λ_0 is the wavelength corresponding to ν_0 . Good fringe visibility is obtained only when either the spectral width is small (the source is quasi-monochromatic) or the OPD is small. The fringes are *localized* in certain areas of space. This result is consistent with the earlier graphical representations. In the area where the OPD is small, the fringes are in phase for all wavelengths. As the OPD increases, the fringes go out of phase since they all have different periods, and the intensity pattern washes out.

This result turns out to be very general: for an incoherent source, the fringes will be localized in the vicinity of zero OPD. There are two other things we should notice about this result. The first is that the first zero of the visibility function occurs when the OPD equals $c/\Delta\nu$. This distance is known as the *coherence length* as it is the path difference over which we can obtain interference. The second item is that the visibility function is a scaled version of the Fourier transform of the source frequency spectrum. It is evaluated for the OPD at the measurement location. The Fourier transform of a uniform distribution is a sinc function. We will discuss this under “Coherence and Interference” later in the chapter.

2.6 INTERFERENCE BY AMPLITUDE DIVISION

The second general method for producing interference is to use the same section of a wavefront from a single source for both resulting wavefronts. The original wavefront amplitude is split into two or more parts, and each fraction is directed along a different optical path. These waves are then recombined to produce interference. This method is called *interference by amplitude division*. There are a great many interferometer designs based on this method. A few will be examined here, and many more will be discussed in Chap. 32, “Interferometers.”

Plane-Parallel Plate

A first example of interference by amplitude division is a plane-parallel plate illuminated by a monochromatic point source. Two virtual images of the point source are formed by the Fresnel reflections at the two surfaces, as shown in Fig. 11. Associated with each of the virtual images is a spherical wave, and interference fringes form wherever these two waves overlap. In this case, this is the volume of space on the source side of the plate. The pattern produced is the same as that found for the interference of two spherical waves (discussed earlier under “Two Spherical Waves”), and nonlocalized fringes are produced. The pattern is symmetric around the line perpendicular to the plate through the source. If a screen is placed along this axis, a pattern similar to circular Newton’s ring fringes are produced as described by Eq. (33), where $d = 2t/n$ is now the separation of the virtual sources. The thickness of the plate is t , its index is n , and the distance R is approximately the screen-plate separation plus the source-plate separation. We have ignored multiple reflections in the plate. As with the interference of two spherical waves, the order of interference is a maximum at the center of the pattern.

The interference of two plane waves can be obtained by illuminating a wedged glass plate with a plane wavefront. If the angle of incidence on the first surface is θ and the wedge angle is α , two plane waves are produced at angles θ and $\theta + 2n\alpha$ due to reflections at the front and rear surfaces. Straight equispaced fringes will result in the volume of space where the two reflected waves overlap. The period of these fringes on a screen parallel to the plate is given by Eq. (26), where the two reflected angles are used.

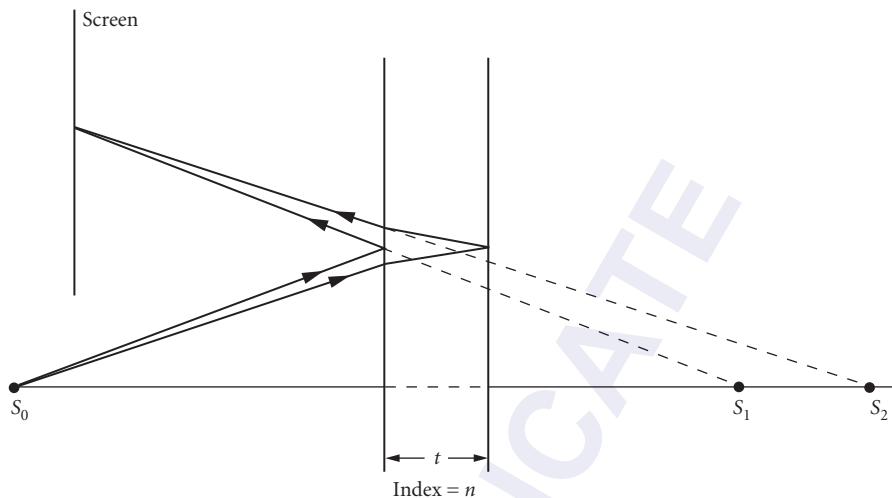


FIGURE 11 Interference from a plane-parallel plate and a point source.

Extended Source

An *extended source* is modeled as a collection of independent point sources. If the source is quasi-monochromatic, all of the point sources radiate at the same nominal frequency, but without a phase relationship. Each point source will produce its own interference pattern, and the net intensity pattern is the sum or integral of all the individual intensity patterns. This is the spatial analogy to the temporal average examined earlier under “Source Spectrum.”

With an extended source, the fringes will be localized where the individual fringe position or spacing is not affected by the location of the point sources that comprise the extended source. We know from our previous examples that a bright fringe (or a dark fringe, depending on phase shifts) will occur when the OPD is zero. If there is a location where the OPD is zero independent of source location, all of the individual interference patterns will be in phase, and the net pattern will show good visibility. In fact, the three-dimensional fringe pattern due to a point source will tend to shift or pivot around this zero-OPD location as the point source location is changed. The individual patterns will therefore be out of phase in areas where the OPD is large, and the average intensity pattern will tend to wash out in these regions as the source size increases.

The general rule for fringe visibility with an extended quasi-monochromatic source is that the fringes will be localized in the region where the OPD between the two interfering wavefronts is small. For a wedged glass plate, the fringes are localized in or near the wedge, and the best visibility occurs as the wedge thickness approaches zero and is perhaps just a few wavelengths. The allowable OPD will depend on the source size and the method of viewing the fringes. This result explains why, under natural light, interference effects are seen in thin soap bubbles but not with other thicker glass objects. An important exception to this rule is the plane-parallel plate where the fringes are localized at infinity.

Fringes of Equal Inclination

There is no section of a plane-parallel plate that produces two reflected wavefronts with zero OPD. The OPD is constant, and we would expect, based on the previous section, that no high-visibility fringes would result with an extended source. If, however, a lens is used

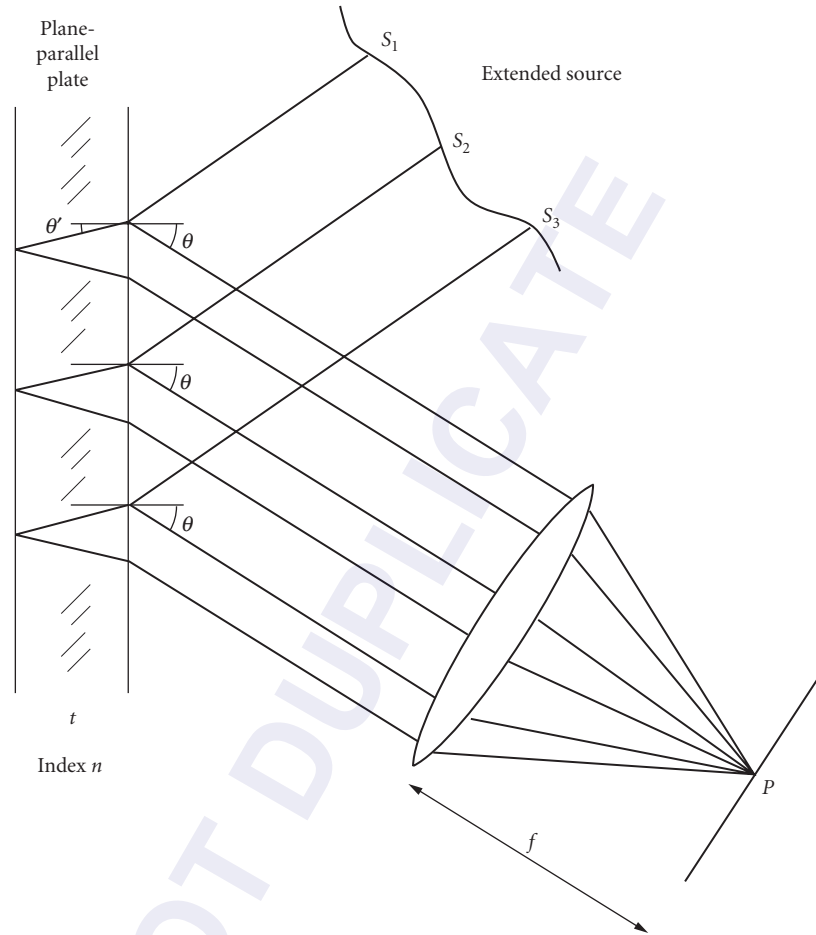


FIGURE 12 The formation of fringes of equal inclination.

to collect the light reflected from the plate, fringes are formed in the back focal plane of the lens. This situation is shown in Fig. 12, and any ray leaving the surface at a particular angle θ is focused to the same point P . For each incident ray at this angle, there are two parallel reflected rays: one from the front surface and one from the back surface. The reflections from different locations on the plate at this angle are due to light from different points in the extended source. The OPD for any pair of these reflected rays is the same regardless of the source location. These rays will interfere at P and will all have the same phase difference. High-visibility fringes result. Different points in the image plane correspond to different angles. The formation of these fringes localized at infinity depends on the two surfaces of the plate being parallel.

The OPD between the reflected rays is a function of the angle of incidence θ , the plate index n , and thickness t :

$$\text{OPD} = 2nt \cos \theta' \quad (48)$$

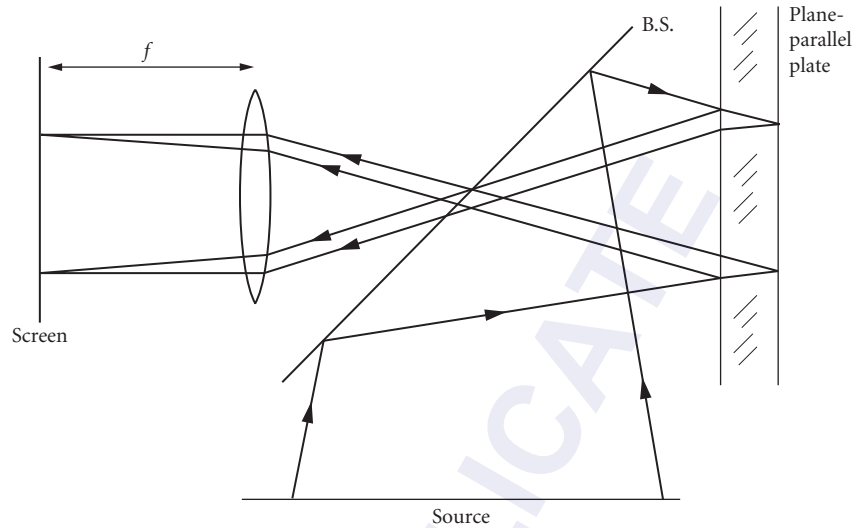


FIGURE 13 The formation of Haidinger fringes.

where θ' is the internal angle. Taking into account the half-wave shift due to the phase change difference of π between an internal and an external reflection, a dark fringe will result for angles satisfying

$$2nt \cos \theta' = m\lambda \quad \text{or} \quad \cos \theta' = \frac{m\lambda}{2nt} \quad (49)$$

where m is an integer. Since only the angle of incidence determines the properties of the interference (everything else is constant), these fringes are called *fringes of equal inclination*. They appear in the back focal plane of the lens and are therefore localized at infinity since infinity is conjugate to the focal plane. As the observation plane is moved away from the focal plane, the visibility of the fringes will quickly decrease.

When the axis of the lens is normal to the surfaces of the plate, a beamsplitter arrangement is required to allow light from the extended source to be reflected into the lens as shown in Fig. 13. Along the axis, $\theta = \theta' = 90^\circ$, and symmetry requires that the fringes are concentric about the axis. In this special case, these fringes are called *Haidinger fringes*, and they are identical in appearance to Newton's rings [Eq. (30)]. If there is an intensity maximum at the center, the radii of the other bright fringes are proportional to the square roots of integers. As with other fringes formed by a plane-parallel plate (discussed earlier), the order of interference decreases with the observation radius on the screen. As θ' increases, the value of m decreases.

Fringes of Equal Thickness

The existence of fringes of equal inclination depends on the incident light being reflected by two parallel surfaces, and the angle of incidence is the mechanism which generates changes in the OPD. There are many arrangements with an extended source where the reflections are not parallel, and the resulting changes in OPD dominate the angle-of-incidence considerations. The fringes produced in this situation are called *fringes of equal thickness*, and we have stated earlier that they will be localized in regions where the OPD between the two reflections is small.

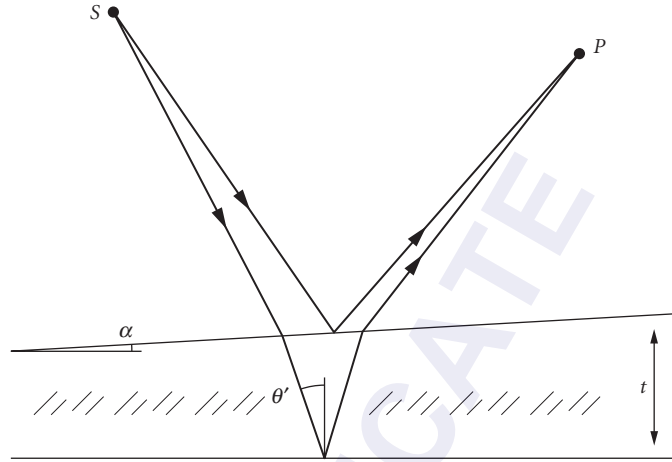


FIGURE 14 The ray path between a point source and an observation point for a wedged plate.

An example of fringes of equal thickness occurs with a wedged glass plate illuminated by a quasi-monochromatic extended source. We know that for each point in the source, a pattern comprising equispaced parallel fringes results, and the net pattern is the sum of all of these individual patterns. However, it is easier to examine this summation by looking at the OPD between the two reflected rays reaching an observation point P from a source point S . This is shown in Fig. 14. The wedge angle is α , the thickness of the plate at this location is t , its index is n , and the internal ray angle is θ' . The exact OPD is difficult to calculate, but under the assumption that α is small and the wedge is sufficiently thin, the following result for the OPD is obtained:

$$\text{OPD} \approx 2nt \cos \theta' \quad (50)$$

As other points on the source are examined, the reflection needed to get light to the observation point will move to a different location on the plate, and different values of both t and θ' will result. Different source points may have greatly different OPDs, and in general the fringe pattern will wash out in the vicinity of P .

This reduction in visibility can be avoided if the observation point is placed in or near the wedge. In this case, all of the paths between S and P must reflect from approximately the same location on the wedge, and the variations in the thickness t are essentially eliminated. The point P where the two reflected rays cross may be virtual. The remaining variations in the OPD are from the different θ' 's associated with different source points. This variation may be limited by observing the fringe pattern with an optical system having a small entrance pupil. This essentially limits the amount of the source that is used to examine any area on the surface. A microscope or the eye focused on the wedge can be used to limit the angles. If the range of values of θ' is small, high-visibility fringes will appear to be localized at the wedge. The visibility of the fringes will decrease as the wedge thickness increases.

It is common to arrange the system so that the fringes are observed in a direction approximately normal to the surface. Taking into account the additional phase shift introduced at the reflection from one of the surfaces, the conditions for bright and dark fringes are then

$$\text{Bright: } 2nt - \frac{\lambda}{2} = m\lambda \quad (51)$$

and

$$\text{Dark: } 2nt = m\lambda \quad (52)$$

where m is an integer greater than or equal to zero. Since t increases linearly across the wedge, the observed pattern will be straight equispaced fringes.

These same conditions hold for any plate where the two surfaces are not parallel. The surfaces may have any shape, and as long as the surface angles are small and the plate is relatively thin, high-visibility fringes localized in the plate are observed. Along a given fringe the value of m is constant, so that a fringe represents a contour of constant optical path length nt . If the index is constant, we have fringes of equal thickness. The fringes provide a contour map of the plate thickness, and adjacent fringes correspond to a change of thickness of $\lambda/2n$. An irregularly shaped pattern will result from the examination of a plate of irregular thickness.

Thin Films

With the preceding background, we can easily explain the interference characteristics of *thin films*. There are two distinct types of films to consider. The first is a thin film of nonuniform thickness, and examples are soap bubbles and oil films on water. The second type is a uniform film, such as would be obtained by vacuum deposition and perhaps used as an antireflection coating. Both of these films share the characteristic of being extremely thin—usually not more than a few wavelengths thick and often just a fraction of a wavelength thick.

With a nonuniform film, fringes of equal thickness localized in the film are produced. There will be a dark fringe in regions of the film where it is substantially thinner than a half wave. We are assuming that the film is surrounded by a lower-index medium such as air so that there is an extra π phase shift. If white light is used for illumination, colored bands will be produced similar to those diagramed in Fig. 10*b* (the curves would need to be modified to rescale the x axis to OPD or film thickness). Each color will produce its first maximum in intensity when the optical thickness of the film is a quarter of that wavelength. As the film thickness increases, the apparent fringe color will first be blue, then green, and finally red. These colored fringes are possible because the film is very thin, and the order of interference m is often zero or one [Eqs. (51) and (52)]. The interference patterns in the various colors are just starting to get out of phase, and interference colors are visible. As the film thickness increases, the various wavelength fringes become jumbled, and distinct fringe patterns are no longer visible.

When a uniform thin film is examined with an extended source, fringes of equal inclination localized at infinity are produced. These fringes will be very broad since the thickness of the film is very small, and large angles will be required to obtain the necessary OPD for a fringe [Eq. (49)]. A common use of this type of film is as an antireflection coating. In this application, a uniform coating that has an optical thickness of a quarter wavelength is applied to a substrate. The coating index is lower than the substrate index, so an extra phase shift is not introduced. A wave at normal incidence is reflected by both surfaces of the coating, and these reflected waves are interfered. If the incident wavelength matches the design of the film, the two reflected waves are out of phase and interfere destructively. The reflected intensity will depend on the Fresnel reflection coefficients at the two surfaces, but will be less than that of the uncoated surface. When a different wavelength is used or the angle of incidence is changed, the effectiveness of the antireflection coating is reduced. More complicated film structures comprising many layers can be produced to modify the reflection or transmission characteristics of the film.

Fizeau Interferometer

The *Fizeau interferometer* compares one optical surface to another by placing them in close proximity. A typical arrangement is shown in Fig. 15, where the extended source is filtered to be quasi-monochromatic. A small air gap is formed between the two optical surfaces, and fringes of equal thickness are observed between the two surfaces. Equations (51) and (52) describe the location of

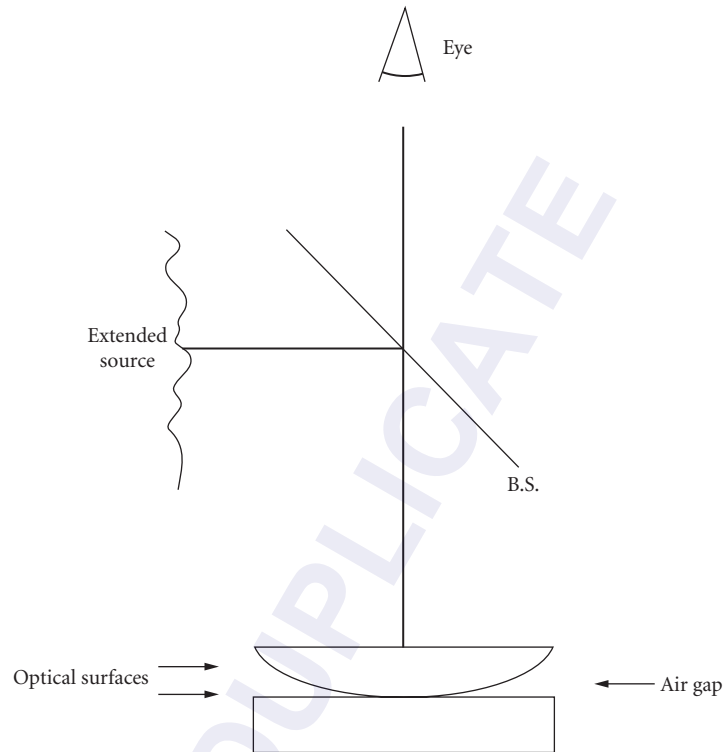


FIGURE 15 Fizeau interferometer.

the fringes, and the index of the thin wedge is now that of air. Along a fringe, the gap is of constant thickness, and adjacent fringes correspond to a change of thickness of a half wavelength. This interferometer is sometimes referred to as a *Newton interferometer*.

This type of interferometer is the standard test instrument in an optical fabrication shop. One of the two surfaces is a reference or known surface, and the interferometric comparison of this reference surface and the test surface shows imperfections in the test part. Differences in radii of the two surfaces are also apparent. The fringes are easy to interpret, and differences of as little as a twentieth of a wavelength can be visually measured. These patterns and this interferometer are further discussed in Chap. 13, "Optical Testing," in Vol. II. The interferometer is often used without the beamsplitter, and the fringes are observed in the direct reflection of the source from the parts.

The classic fringe pattern produced by a Fizeau interferometer is *Newton's rings*. These are obtained by comparing a convex sphere to a flat surface. The parabolic approximation for the sag of a sphere of radius R is

$$\text{sag}(\rho) = \frac{\rho^2}{2R} \quad (53)$$

and ρ is the radial distance from the vertex of the sphere. If we assume the two surfaces are in contact at $\rho=0$, the OPD between the reflected waves is twice the gap, and the condition for a dark fringe is

$$\rho = \sqrt{m\lambda R} \quad (54)$$

Circular fringes that increase in radius as the square root of ρ are observed. Note that a dark fringe occurs at the center of the pattern. In reflection, this point must be dark, as there is no interface at the contact point to produce a reflection.

Michelson Interferometer

There are many two-beam interferometers which allow the surfaces producing the two wavefronts to be physically separated by a large distance. These instruments allow the two wavefronts to travel along different optical paths. One of these is the *Michelson interferometer* diagrammed in Fig. 16a. The two interfering wavefronts are produced by the reflections from the two mirrors. A plate beamsplitter with one face partially silvered is used, and an identical block of glass is placed in one of the arms of the interferometer to provide the same amount of glass path in each arm. This cancels the effects of the dispersion of the glass beamsplitter and allows the system to be used with white light since the optical path difference is the same for all wavelengths.

Figure 16b provides a folded view of this interferometer and shows the relative optical position of the two mirrors as seen by the viewing screen. It should be obvious that the two

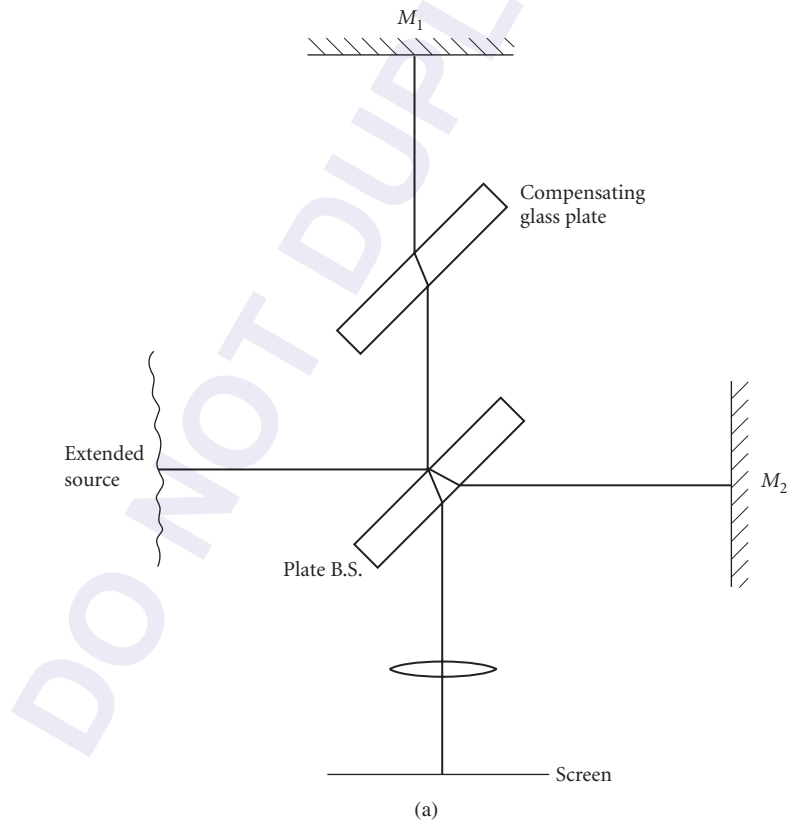


FIGURE 16 Michelson interferometer: (a) schematic view and (b) folded view showing the relative optical position of the two mirrors.

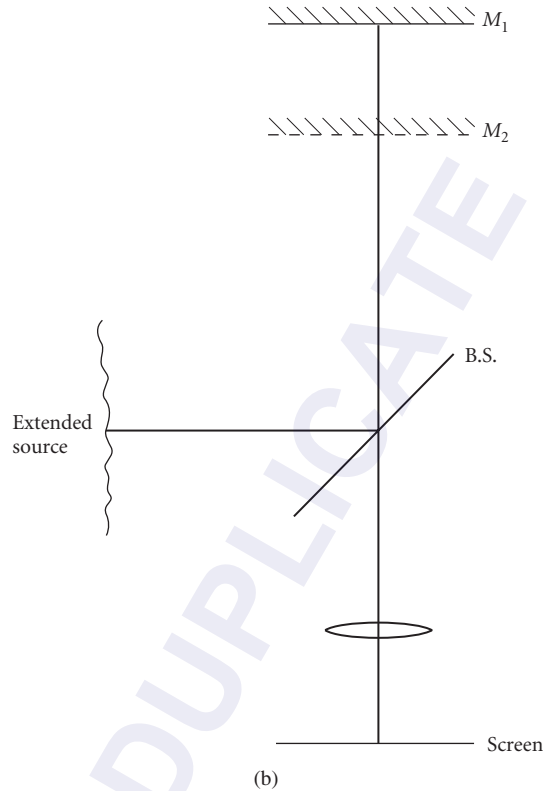


FIGURE 16 (Continued)

mirrors can be thought of as the two surfaces of a “glass” plate that is illuminated by the source. In this case, the index of the fictitious plate is one, and the reflectivity at the two surfaces is that of the mirrors. Depending on the mirror orientations and shapes, the interferometer either mimics a plane-parallel plate of adjustable thickness, a wedge of arbitrary angle and thickness, or the comparison of a reference surface with an irregular or curved surface. The type of fringes that are produced will depend on this configuration, as well as on the source used for illumination.

When a monochromatic point source is used, nonlocalized fringes are produced, and the imaging lens is not needed. Two virtual-source images are produced, and the resulting fringes can be described by the interference of two spherical waves (discussed earlier). If the mirrors are parallel, circular fringes centered on the line normal to the mirrors result as with a plane-parallel plate. The source separation is given by twice the apparent mirror separation. If the mirrors have a relative tilt, the two source images appear to be laterally displaced, and hyperbolic fringes result. Along a plane bisecting the source images, straight equispaced fringes are observed.

When an extended monochromatic source is used, the interference fringes are localized. If the mirrors are parallel, fringes of equal inclination or Haidinger fringes (as described earlier) are produced. The fringes are localized at infinity and are observed in the rear focal plane of the imaging lens. Fringes of equal thickness localized at the mirrors are generated when the mirrors are tilted.

The apparent mirror separation should be kept small, and the imaging lens should focus on the mirror surface.

If the extended source is polychromatic, colored fringes localized at the mirrors result. They are straight for tilted mirrors. The fringes will have high visibility only if the apparent mirror separation or OPD is smaller than the coherence length of the source. Another way of stating this is that the order of interference m must be small to view the colored fringes. As m increases, the fringes will wash out. The direct analogy here is a thin film. As the mirror separation is varied, the fringe visibility will vary. The fringe visibility as a function of mirror separation is related to the source frequency spectrum (see under “Source Spectrum” and “Coherence and Interference”), and this interferometer forms the basis of a number of spectrometers. When the source spectrum is broad, chromatic fringes cannot be viewed with the mirrors parallel. This is because the order of interference for fringes of equal inclination is a maximum at the center of the pattern.

An important variation of the Michelson interferometer occurs when monochromatic collimated light is used. This is the *Twyman-Green interferometer*, and is a special case of point-source illumination with the source at infinity. Plane waves fall on both mirrors, and if the mirrors are flat, nonlocalized equispaced fringes are produced. Fringes of equal thickness can be viewed by imaging the mirrors onto the observation screen. If one of the mirrors is not flat, the fringes represent changes in the surface height. The two surfaces are compared as in the Fizeau interferometer. This interferometer is an invaluable tool for optical testing.

2.7 MULTIPLE BEAM INTERFERENCE

Throughout the preceding discussions, we have assumed that only two waves were being interfered. There are many situations where multiple beams are involved. Two examples are the diffraction grating and a plane-parallel plate. We have been ignoring multiple reflections, and in some instances these extra beams are very important. The net electric field is the sum of all of the component fields. The two examples noted above present different physical situations: all of the interfering beams have a constant intensity with a diffraction grating, and the intensity of the beams from a plane-parallel plate decreases with multiple reflections.

Diffraction Grating

A *diffraction grating* can be modeled as a series of equispaced slits, and the analysis bears a strong similarity to the Young’s double slit (discussed earlier). It operates by division of wavefront, and the geometry is shown in Fig. 17. The slit separation is d , the OPD between successive beams for a given observation angle θ is $d \sin(\theta)$, and the corresponding phase difference $\Delta\phi = 2\pi d \sin(\theta)/\lambda$. The field due to the n th slit at a distant observation point is

$$E_j(\theta) = A e^{i(j-1)\Delta\phi} \quad j = 1, 2, \dots, N \quad (55)$$

where all of the beams have been referenced to the first slit, and there are N total slits. The net field is

$$E(\theta) = \sum_{j=1}^N E_j(\theta) = A \sum_{j=1}^N (e^{i\Delta\phi})^{j-1} \quad (56)$$

which simplifies to

$$E(\theta) = A \left(\frac{1 - e^{iN\Delta\phi}}{1 - e^{i\Delta\phi}} \right) \quad (57)$$

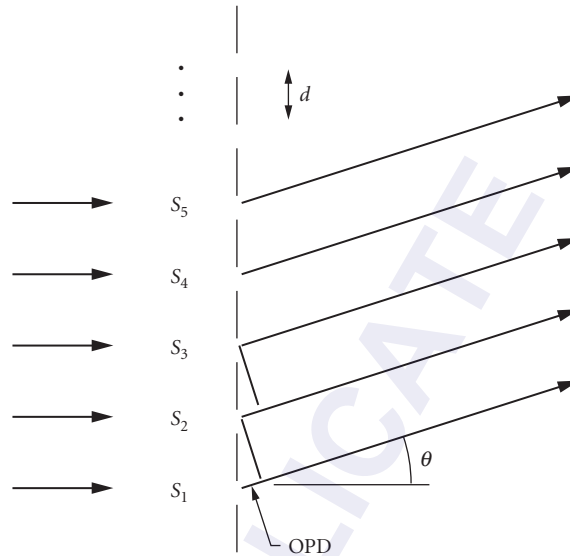


FIGURE 17 Diffraction grating: multiple-beam interference by division of wavefront.

The resulting intensity is

$$I(\theta) = I_0 \left[\frac{\sin^2\left(\frac{N\Delta\phi}{2}\right)}{\sin^2\left(\frac{\Delta\phi}{2}\right)} \right] = I_0 \left[\frac{\sin^2\left(\frac{N\pi d \sin(\theta)}{\lambda}\right)}{\sin^2\left(\frac{\pi d \sin(\theta)}{\lambda}\right)} \right] \quad (58)$$

where I_0 is the intensity due to an individual slit.

This intensity pattern is plotted in Fig. 18 for $N = 5$. The result for $N = 2$, which is the double-slit experiment, is also shown. The first thing to notice is that the locations of the maxima are the same, independent of the number of slits. A maximum of intensity is obtained whenever the phase difference between adjacent slits is a multiple of 2π . These maxima occur at the diffraction angles given by

$$\sin(\theta) = \frac{m\lambda}{d} \quad (59)$$

where m is an integer. The primary difference between the two patterns is that with multiple slits, the intensity at the maximum increases to N^2 times that due to a single slit, and this energy is concentrated into a much narrower range of angles. The full width of a diffraction peak between intensity zero corresponds to a phase difference $\Delta\phi$ of $4\pi/N$.

The number of intensity zeros between peaks is $N - 1$. As the number of slits increases, the angular resolution or resolving power of the grating greatly increases. The effects of a finite slit width can be added by replacing I_0 in Eq. (58) by the single-slit diffraction pattern. This intensity variation forms an envelope for the curve in Fig. 18.

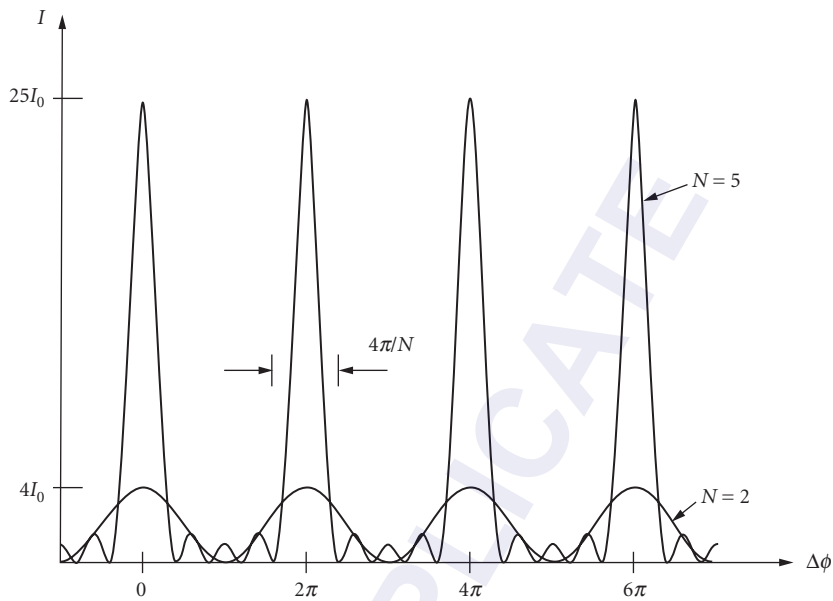


FIGURE 18 The interference patterns produced by gratings with 2 and 5 slits.

Plane-Parallel Plate

The plane-parallel plate serves as a model to study the interference of multiple waves obtained by division of amplitude. As we shall see, the incremental phase difference between the interfering beams is constant but, in this case, the beams have different intensities. A plate of thickness t and index n with all of the reflected and transmitted beams is shown in Fig. 19. The amplitude reflection and transmission coefficients are ρ and ρ' , and τ and τ' , where the primes indicate

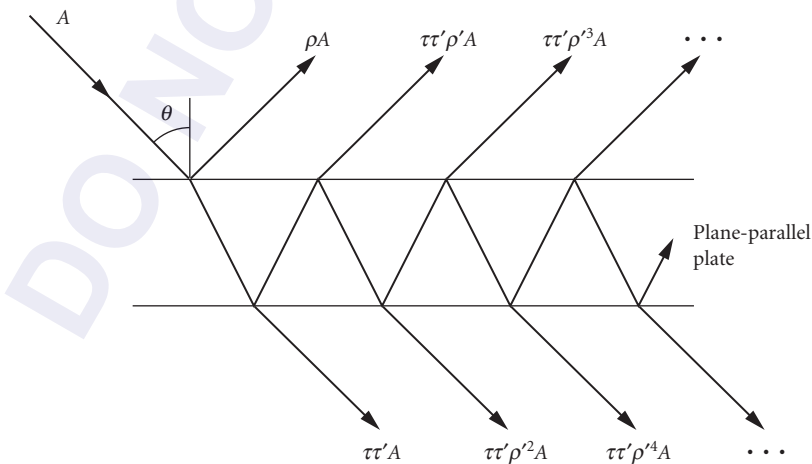


FIGURE 19 Plane-parallel plate: multiple-beam interference by division of amplitude.

reflection or transmission from within the plate. The first reflected beam is 180° out of phase with the other reflected beams since it is the only beam to undergo an external reflection, and $\rho = -\rho'$. Note that ρ' occurs only in odd powers for the reflected beams. Each successive reflected or transmitted beam is reduced in amplitude by ρ^2 . The phase difference between successive reflected or transmitted beams is the same as we found when studying fringes of equal inclination from a plane-parallel plate:

$$\Delta\phi = \left[\frac{4\pi nt \cos(\theta')}{\lambda} \right] \quad (60)$$

where θ' is the angle inside the plate.

The transmitted intensity can be determined by first summing all of the transmitted amplitudes:

$$E(\Delta\phi) = \sum_{j=1}^{\infty} E_j = A\tau\tau' \sum_{j=1}^{\infty} (\rho^2 e^{i\Delta\phi})^{j-1} \quad (61)$$

where the phase is referenced to the first transmitted beam. The result of the summation is

$$E(\Delta\phi) = \left(\frac{A\tau\tau'}{1 - \rho^2 e^{i\Delta\phi}} \right) \quad (62)$$

The transmitted intensity I_t is the squared modulus of the amplitude which, after simplification, becomes

$$\frac{I_t}{I_0} = \frac{1}{1 + \left(\frac{2\rho}{1 - \rho^2} \right)^2 \sin^2(\Delta\phi/2)} \quad (63)$$

where I_0 is the incident intensity. We have also assumed that there is no absorption in the plate, and therefore $\tau\tau' + \rho^2 = 1$. Under this condition of no absorption, the sum of the reflected and transmitted light must equal the incident light: $I_t + I_r = I_0$. The expressions for the transmitted and reflected intensities are then

$$\frac{I_t}{I_0} = \frac{1}{1 + F \sin^2(\Delta\phi/2)} \quad (64)$$

and

$$\frac{I_r}{I_0} = \frac{F \sin^2(\Delta\phi/2)}{1 + F \sin^2(\Delta\phi/2)} \quad (65)$$

and F is defined as

$$F \equiv \left(\frac{2\rho}{1 - \rho^2} \right)^2 \quad (66)$$

F is the *coefficient of finesse* of the system and is a function of the surface reflectivity only. The value of F will have a large impact on the shape of the intensity pattern. Note that the reflected intensity could also have been computed by summing the reflected beams.

A maximum of transmitted intensity, or a minimum of reflected intensity, will occur when $\Delta\phi/2 = m\pi$, where m is an integer. Referring back to Eq. (60), we find that this corresponds to the angles

$$\cos\theta' = \frac{m\lambda}{2nt} \quad (67)$$

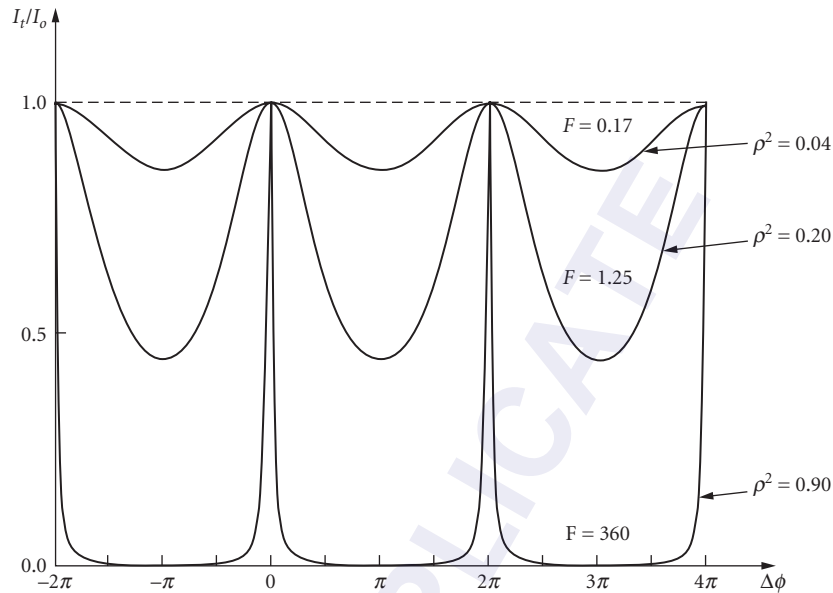


FIGURE 20 The transmitted intensity of a multiple-beam interference pattern produced by a plane-parallel plate.

This is exactly the same condition that was found for a plane-parallel plate with two beams [Eq. (49)]. With an extended source, fringes of equal inclination are formed, and they are localized at infinity. They must be at infinity since all of the reflected or transmitted beams are parallel for a given input angle. The fringes are observed in the rear focal plane of a viewing lens. If the optical axis of this lens is normal to the surface, circular fringes about the axis are produced. The locations of the maxima and minima of the fringes are the same as were obtained with two-beam interference.

The shape of the intensity profile of these multiple beam fringes is not sinusoidal, as it was with two beams. A plot of the transmitted fringe intensity [Eq. (64)] as a function of $\Delta\phi$ is shown in Fig. 20 for several values of F . When the phase difference is a multiple of 2π , we obtain a bright fringe independent of F or ρ . When F is small, low-visibility fringes are produced. When F is large, however, the transmitted intensity is essentially zero unless the phase has the correct value. It drops off rapidly for even small changes in $\Delta\phi$. The transmitted fringes will be very narrow bright circles on an essentially black background. The reflected intensity pattern is one minus this result, and the fringe pattern will be very dark bands on a uniform bright background. The reflected intensity profile is plotted in Fig. 21 for several values of F .

The value of F is a strong function of the surface reflectivity $R = \rho^2$. We do not obtain appreciable values of F until the reflectivity is approximately one. For example, $R = 0.8$ produces $F = 80$, while $R = 0.04$ gives $F = 0.17$. This latter case is typical for uncoated glass, and dim broad fringes in reflection result, as in Fig. 21. The pattern is approximately sinusoidal, and it is clear that our earlier assumptions about ignoring multiple reflections when analyzing a plane-parallel plate are valid for many low-reflectivity situations.

The multiple beam interference causes an energy redistribution much like that obtained from a diffraction grating. A strong response is obtained only when all of the reflected beams at a given angle add up in phase. The difference between this pattern and that of a diffraction pattern is that there are no oscillations or zeros between the transmitted intensity maxima. This is a result of the

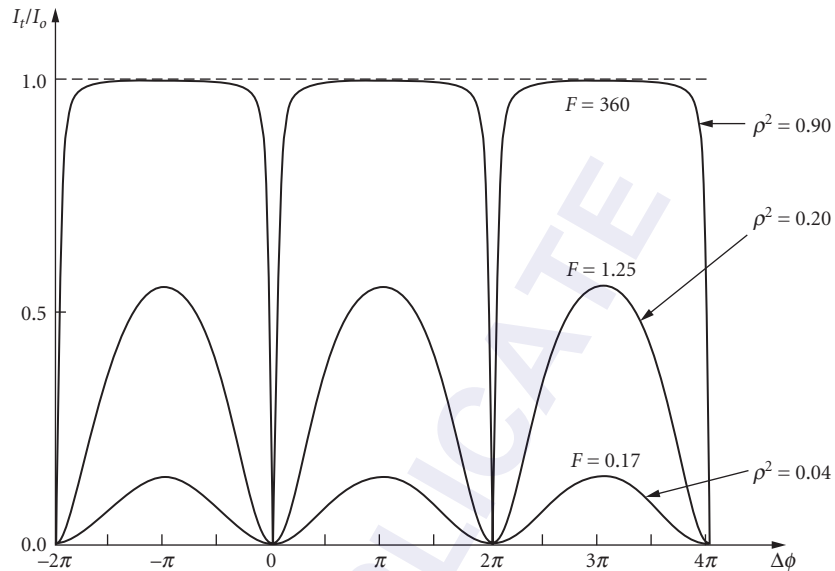


FIGURE 21 The reflected intensity of a multiple-beam interference pattern produced by a plane-parallel plate.

unequal amplitudes of the interfering beams. With a diffraction grating, all of the beams have equal amplitude, and the resultant intensity oscillates as more beams are added.

Multiple-beam fringes of equal thickness can be produced by two high-reflectivity surfaces in close proximity in a Fizeau interferometer configuration. The dark fringes will narrow to sharp lines, and each fringe will represent a contour of constant OPD between the surfaces. As before, a dark fringe corresponds to a gap of an integer number of half wavelengths. The area between the fringes will be bright. The best fringes will occur when the angle and the separation between the surfaces is kept small. This will prevent the multiple reflections from walking off or reflecting out of the gap.

Fabry-Perot Interferometer

The *Fabry-Perot interferometer* is an important example of a system which makes use of multiple-beam interference. This interferometer serves as a high-resolution spectrometer and also as an optical resonator. In this latter use, it is an essential component of a laser. The system is diagrammed in Fig. 22, and it consists of two highly reflective parallel surfaces separated by a distance t . These two separated reflective plates are referred to as a *Fabry-Perot etalon* or cavity, and an alternate arrangement has the reflected coatings applied to the two surfaces of a single glass plate. The two lenses serve to collimate the light from a point on the extended source in the region of the cavity and to then image this point onto the screen. The screen is located in the focal plane of the lens so that fringes of equal inclination localized at infinity are viewed. As we have seen, light of a fixed wavelength will traverse the etalon only at certain well-defined angles. Extremely sharp multiple-beam circular fringes in transmission are produced on the screen, and their profile is the same as that shown in Fig. 20.

If the source is not monochromatic, a separate independent circular pattern is formed for each wavelength. Equation (67) tells us that the location or scale of the fringes is dependent on the wavelength. If the source is composed of two closely spaced wavelengths, the ring structure is doubled,

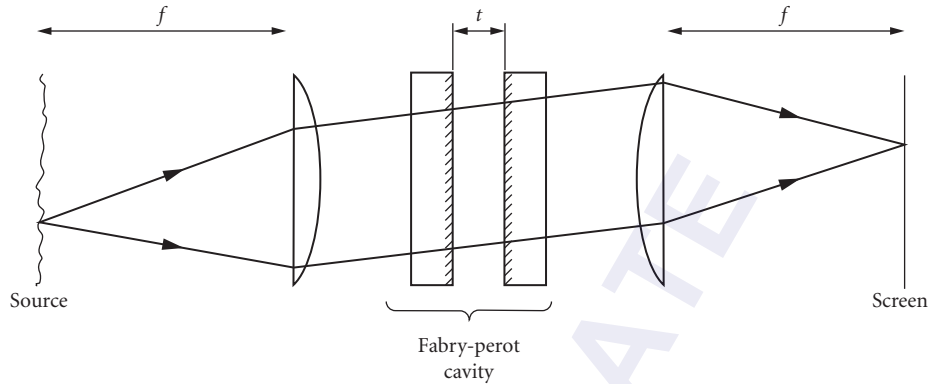


FIGURE 22 Fabry-Perot interferometer.

and the separation of the two sets of rings allows the hyperfine structure of the spectral lines to be evaluated directly. More complicated spectra, usually composed of discrete spectral lines, can also be measured. This analysis is possible even though the order of interference is highest in the center of the pattern. If the phase change $\Delta\phi$ due to the discrete wavelengths is less than the phase change between adjacent fringes, nonoverlapping sharp fringes are seen.

A quantity that is often used to describe the performance of a Fabry-Perot cavity is the *fineness* \mathcal{F} . It is a measure of the number of resolvable spectral lines, and is defined as the ratio of the phase difference between adjacent fringes to the full width-half maximum FWHM of a single fringe. Since the fringe width is a function of the coefficient of fineness, the fineness itself is also a strong function of reflectivity. The phase difference between adjacent fringes is 2π , and the half width-half maximum can be found by setting Eq. (64) equal to $\frac{1}{2}$ and solving for $\Delta\phi$. The FWHM is twice this value, and under the assumption that F is large,

$$\text{FWHM} = \frac{4}{\sqrt{F}} \quad (68)$$

and the fineness is

$$\mathcal{F} = \frac{2\pi}{\text{FWHM}} = \frac{\pi\sqrt{F}}{2} = \frac{\pi\rho}{1-\rho^2} = \frac{\pi\sqrt{R}}{1-R} \quad (69)$$

where ρ is the amplitude reflectivity, and R is the intensity reflectivity. Typical values for the fineness of a cavity with flat mirrors is about 30 and is limited by the flatness and parallelism of the mirrors. There are variations in $\Delta\phi$ across the cavity. Etalons consisting of two curved mirrors can be constructed with a much higher fineness, and values in excess of 10,000 are available.

Another way of using the Fabry-Perot interferometer as a spectrometer is suggested by rewriting the transmission [Eq. (64)] in terms of the frequency ν :

$$T = \frac{I_t}{I_0} = \frac{1}{1 + F \sin^2(2\pi\nu t/c)} \quad (70)$$

where Eq. (60) relates the phase difference to the wavelength, t is the mirror separation, and an index of one and normal incidence ($\theta' = 0$) have been assumed. This function is plotted in Fig. 23, and a series of transmission spikes separated in frequency by $c/2t$ are seen. A maximum occurs whenever the value of the sine is zero. The separation of these maxima is known as the *free spectral range*, FSR. If the separation of the mirrors is changed slightly, these transmission peaks will scan the frequency axis. Since the order of interference m is usually very large, it takes

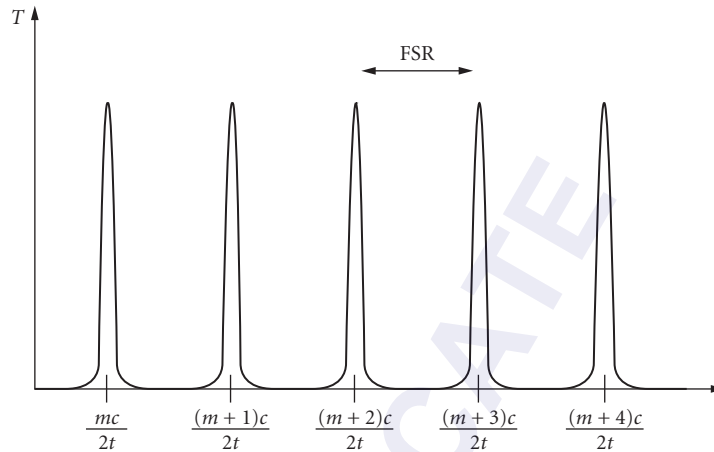


FIGURE 23 The transmission of a Fabry-Perot cavity as a function of frequency.

only a small change in the separation to move the peaks by one FSR. In fact, to scan one FSR, the required change in separation is approximately t/m . If the on-axis transmitted intensity is monitored while the mirror separation is varied, a high-resolution spectrum of the source is obtained. The source spectrum must be contained within one free spectral range so that the spectrum is probed by a single transmission peak at a time. If this were not the case, the temporal signal would contain simultaneous contributions from two or more frequencies resulting from different transmission peaks. Under this condition there are overlapping orders, and it is often prevented by using an auxiliary monochromator with the scanning Fabry-Perot cavity to preselect or limit the frequency range of the input spectrum. The resolution $\Delta\nu$ of the trace is limited by the finesse of the cavity.

For a specific cavity, the value of m at a particular transmission peak, and some physical insight into the operation of this spectrometer, is obtained by converting the frequency of a particular transmission mode $mc/2t$ into wavelength:

$$\lambda = \frac{2t}{m} \quad \text{or} \quad t = m \frac{\lambda}{2} \quad (71)$$

For the m th transmission maximum, exactly m half waves fit across the cavity. This also implies that the round-trip path within the cavity is an integer number of wavelengths. Under this condition, all of the multiply-reflected beams are in phase everywhere in the cavity, and therefore all constructively interfere. A maximum in the transmission occurs. Other maxima occur at different wavelengths, but these specific wavelengths must also satisfy the condition that the cavity spacing is an integer number of half wavelengths.

These results also allow us to determine the value of m . If a 1-cm cavity is used and the nominal wavelength is 500 nm, $m = 40,000$ and $\text{FSR} = 1.5 \times 10^{10}$ Hz. The wavelength interval corresponding to this FSR is 0.0125 nm. If a 1-mm cavity is used instead, the results are $m = 4000$ and $\text{FSR} = 1.5 \times 10^{11}$ Hz = 0.125 nm. We see now that to avoid overlapping orders, the spectrum must be limited to a very narrow range, and this range is a function of the spacing. Cavities with spacings of a few tens of μm 's are available to increase the FSR. Increasing the FSR does have a penalty. The finesse of a cavity depends only on the reflectivities, so as the FSR is increased by decreasing t , the FWHM of the transmission modes increases to maintain a constant ratio. The number of resolvable spectrum lines remains constant, and the absolute spectral resolution decreases.

A mirror translation of a half wavelength is sufficient to cover the FSR of the cavity. The usual scanning method is to separate the two mirrors with a piezoelectric spacer. As the applied voltage is changed, the cavity length will also change. An alternate method is to change the index of the air in the cavity by changing the pressure.

2.8 COHERENCE AND INTERFERENCE

The observed fringe visibility is a function of the spatial and temporal coherence of the source. The classical assumption for the analysis is that every point on an extended source radiates independently and therefore produces its own interference pattern. The net intensity is the sum of all of the individual intensity patterns. In a similar manner, each wavelength or frequency of a nonmonochromatic source radiates independently, and the temporal average is the sum of the individual temporal averages. *Coherence theory* allows the interference between the light from two point sources to be analyzed, and a good visual model is an extended source illuminating the two pinholes in Young's double slit. We need to determine the relationship between the light transmitted through the two pinholes. Coherence theory also accounts for the effects of the spectral bandwidth of the source.

With interference by division of amplitude using an extended source, the light from many point sources is combined at the observation point, and the geometry of the interferometer determines where the fringes are localized. Coherence theory will, however, predict the spectral bandwidth effects for division of amplitude interference. Each point on the source is interfered with an image of that same point. The temporal coherence function relates the interference of these two points independently of other points on the source. The visibility function for the individual interference pattern due to these two points is computed, and the net pattern is the sum of these patterns for the entire source. The temporal coherence effects in division of amplitude interference are handled on a point-by-point basis across the source.

In this section, the fundamentals of coherence theory as it relates to interference are introduced. Much more detail on this subject can be found in Chap. 5, "Coherence Theory."

Mutual Coherence Function

We will consider the interference of light from two point sources or pinholes. This light is derived from a common origin so that there may be some relationship between the complex fields at the two sources. We will represent these amplitudes at the pinholes as $E_1(t)$ and $E_2(t)$, as shown in Fig. 24. The propagation times between the two sources and the observation point are t_1 and t_2 ,

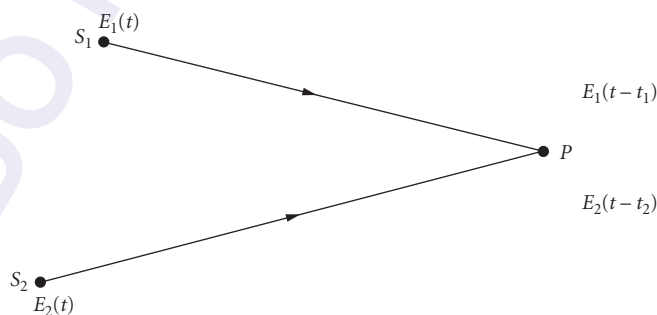


FIGURE 24 Geometry for examining the mutual coherence of two sources.

where the times are related to the optical path lengths by $t_i = \text{OPL}_i/c$. The two complex amplitudes at the observation point are then $E_1(t - t_1)$ and $E_2(t - t_2)$, where the amplitudes have been scaled to the observation plane. The time-average intensity at the observation point can be found by returning to Eq. (13), which is repeated here with the time dependence:

$$I = I_1 + I_2 + \langle E_1(t - t_1)E_2^*(t - t_2) \rangle + \langle E_1^*(t - t_1)E_2(t - t_2) \rangle \quad (72)$$

where I_1 and I_2 are the intensities due to the individual sources. If we now shift our time origin by t_2 , we obtain

$$I = I_1 + I_2 + \langle E_1(t + \tau)E_2^*(t) \rangle + \langle E_1^*(t + \tau)E_2(t) \rangle \quad (73)$$

where

$$\tau = t_2 - t_1 = \frac{\text{OPL}_2 - \text{OPL}_1}{c} = \frac{\text{OPD}}{c} \quad (74)$$

The difference in transit times for the two paths is τ . The last two terms in the expression for the intensity are complex conjugates, and they contain the interference terms.

We will now define the *mutual coherence function* $\Gamma_{12}(\tau)$:

$$\Gamma_{12}(\tau) = \langle E_1(t + \tau)E_2^*(t) \rangle \quad (75)$$

which is the cross correlation of the two complex amplitudes. With this identification, the intensity of the interference pattern is

$$I = I_1 + I_2 + \Gamma_{12}(\tau) + \Gamma_{12}^*(\tau) \quad (76)$$

or, recognizing that a quantity plus its complex conjugate is twice the real part,

$$I = I_1 + I_2 + 2 \text{Re}\{\Gamma_{12}(\tau)\} \quad (77)$$

It is convenient to normalize the mutual coherence function by dividing by the square root of the product of the two self-coherence functions. The result is the *complex degree of coherence*:

$$\gamma_{12}(\tau) = \frac{\Gamma_{12}(\tau)}{\sqrt{\Gamma_{11}(0)\Gamma_{22}(0)}} = \frac{\Gamma_{12}(\tau)}{\sqrt{\langle |E_1(t)|^2 \rangle \langle |E_2(t)|^2 \rangle}} = \frac{\Gamma_{12}(\tau)}{\sqrt{I_1 I_2}} \quad (78)$$

and the intensity can be rewritten:

$$I = I_1 + I_2 + 2\sqrt{I_1 I_2} \text{Re}\{\gamma_{12}(\tau)\} \quad (79)$$

We can further simplify the result by writing $\gamma_{12}(\tau)$ as a magnitude and a phase:

$$\gamma_{12}(\tau) = |\gamma_{12}(\tau)| e^{i\phi_{12}(\tau)} = |\gamma_{12}(\tau)| e^{i[\alpha_{12}(\tau) - \Delta\phi(\tau)]} \quad (80)$$

where $\alpha_{12}(\tau)$ is associated with the source, and $\Delta\phi(\tau)$ is the phase difference due to the OPD between the two sources and the observation point [Eq. (18)]. The quantity $|\gamma_{12}(\tau)|$ is known as the *degree of coherence*. The observed intensity is therefore

$$I = I_1 + I_2 + 2\sqrt{I_1 I_2} |\gamma_{12}(\tau)| \cos[\alpha_{12}(\tau) - \Delta\phi(\tau)] \quad (81)$$

The effect of $\alpha_{12}(\tau)$ is to add a phase shift to the intensity pattern. The fringes will be shifted. A simple example of this situation is Young's double-slit experiment illuminated by a tilted plane wave or a decentered source. With quasi-monochromatic light, the variations of both $|\gamma_{12}(\tau)|$ and $\alpha_{12}(\tau)$ with τ are slow with respect to changes of $\Delta\phi(\tau)$, so that the variations in the interference pattern in the observation plane are due primarily to changes in $\Delta\phi$ with position.

A final rewrite of Eq. (81) leads us to the intensity pattern at the observation point:

$$I = I_0 \left\{ 1 + \frac{2\sqrt{I_1 I_2}}{I_1 + I_2} |\gamma_{12}(\tau)| \cos[\alpha_{12}(\tau) - \Delta\phi(\tau)] \right\} \quad (82)$$

where $I_0 = I_1 + I_2$. The fringe visibility is therefore

$$\gamma(\tau) = \frac{2\sqrt{I_1 I_2}}{I_1 + I_2} |\gamma_{12}(\tau)| \quad (83)$$

and is a function of the degree of coherence and τ . Remember that τ is just the temporal measure of the OPD between the two sources and the observation point. If the two intensities are equal, the fringe visibility is simply the degree of coherence: $\gamma(\tau) = |\gamma_{12}(\tau)|$. The degree of coherence will take on values between 0 and 1. The source is *coherent* when $|\gamma_{12}(\tau)| = 1$, and completely incoherent when $|\gamma_{12}(\tau)| = 0$. The source is said to be *partially coherent* for other values. No fringes are observed with an incoherent source, and the visibility is reduced with a partially coherent source.

Spatial Coherence

The spatial extent of the source and its distance from the pinholes will determine the visibility of the fringes produced by the two pinhole sources (see Fig. 25). Each point on the source will produce a set of Young's fringes, and the position of this pattern in the observation plane will shift with source position. The value of $\alpha_{12}(\tau)$ changes with source position. The existence of multiple shifted patterns will reduce the overall visibility. As an example, consider a quasi-monochromatic source that consists of a several point sources arranged in a line. Each produces a high modulation fringe pattern in the observation plane (Fig. 26a), but there is a lateral shift between each pattern. The net pattern shows a fringe with the same period as the individual patterns, but it has a reduced modulation due to the shifts (Fig. 26b). This reduction in visibility can be predicted by calculating the degree of coherence $|\gamma_{12}(\tau)|$ at the two pinholes.

Over the range of time delays between the interfering beams that are usually of interest, the degree of coherence is a slowly varying function and is approximately equal to the value at $\tau = 0$: $|\gamma_{12}(\tau)| = |\gamma_{12}(0)| = |\gamma_{12}|$. The *van Cittert–Zernike theorem* allows the degree of coherence in the geometry of Fig. 25 to be calculated. Let θ be the angular separation of the two pinholes as seen

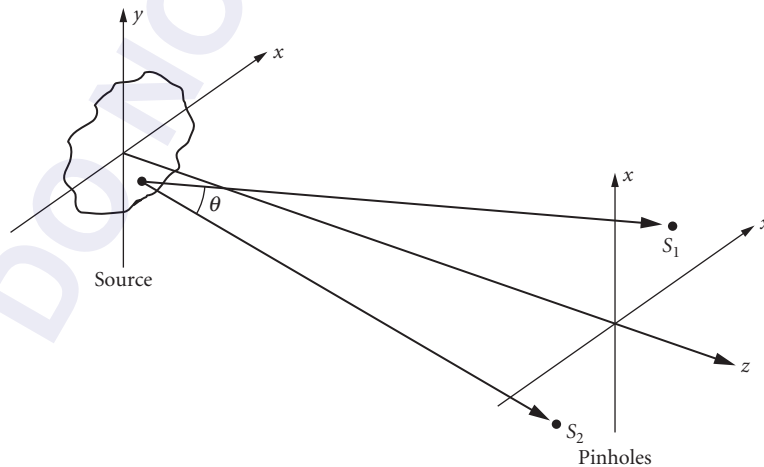


FIGURE 25 An extended source illuminating two pinholes.

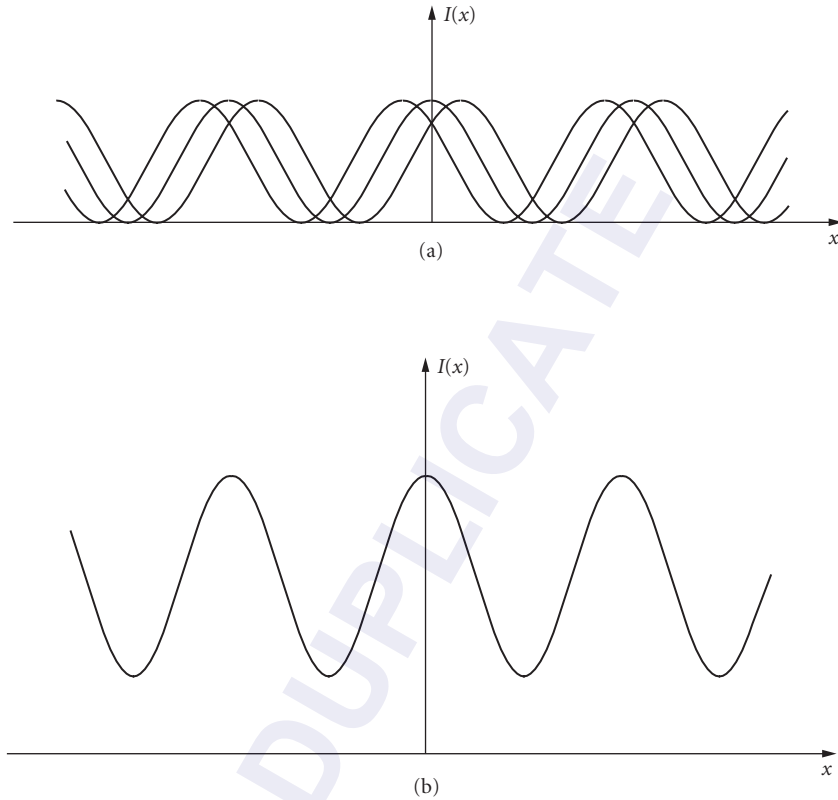


FIGURE 26 The interference pattern produced by a linear source: (a) the individual fringe patterns and (b) the net fringe pattern with reduced visibility.

from the source. This theorem states that degree of coherence between two points is the modulus of the scaled and normalized Fourier transform of the source intensity distribution:

$$|\gamma_{12}| = \left| \frac{\iint_S I(\xi, \eta) e^{i(2\pi/\lambda)(\xi\theta_x + \eta\theta_y)} d\xi d\eta}{\iint_S I(\xi, \eta) d\xi d\eta} \right| \quad (84)$$

where θ_x and θ_y are the x and y components of the pinhole separation θ , and the integral is over the source.

Two cases that are of particular interest are a slit source and a circular source. The application of the van Cittert–Zernike theorem yields the two coherence functions:

$$\text{Slit source of width } w: \quad |\gamma_{12}| = \left| \text{sinc} \left(\frac{w\theta_x}{\lambda} \right) \right| = \left| \text{sinc} \left(\frac{wa}{\lambda z} \right) \right| \quad (85)$$

$$\text{Circular source of diameter } d: \quad |\gamma_{12}| = \left| \frac{2J_1 \left(\frac{\pi d\theta_x}{\lambda} \right)}{\frac{\pi d\theta_x}{\lambda}} \right| = \left| \frac{2J_1 \left(\frac{\pi da}{\lambda z} \right)}{\frac{\pi da}{\lambda z}} \right| \quad (86)$$

where a is the separation of the pinholes, z is the distance from the source to the pinholes, the sinc function is defined by Eq. (42), and J_1 is a first-order Bessel function. The pinholes are assumed to be located on the x axis. These two functions share the common characteristic of a central core surrounded by low-amplitude side lobes. We can imagine these functions of pinhole spacing mapped onto the aperture plane. The coherence function is centered on one of the pinholes. If the other pinhole is then within the central core, high-visibility fringes are produced. If the pinhole spacing places the second pinhole outside the central core, low-visibility fringes result.

Michelson Stellar Interferometer

The *Michelson stellar interferometer* measures the diameter of stars by plotting out the degree of coherence due to the light from the star. The system is shown in Fig. 27. Two small mirrors separated by the distance a sample the light and serve as the pinholes. The spacing between these mirrors can be varied. This light is then directed along equal path lengths into a telescope, and the two beams interfere in the image plane. To minimize chromatic effects, the input light should be filtered to a small range of wavelengths. The modulation of the fringes is measured as a function of the mirror spacing to measure the degree of coherence in the plane of the mirrors. This result will follow

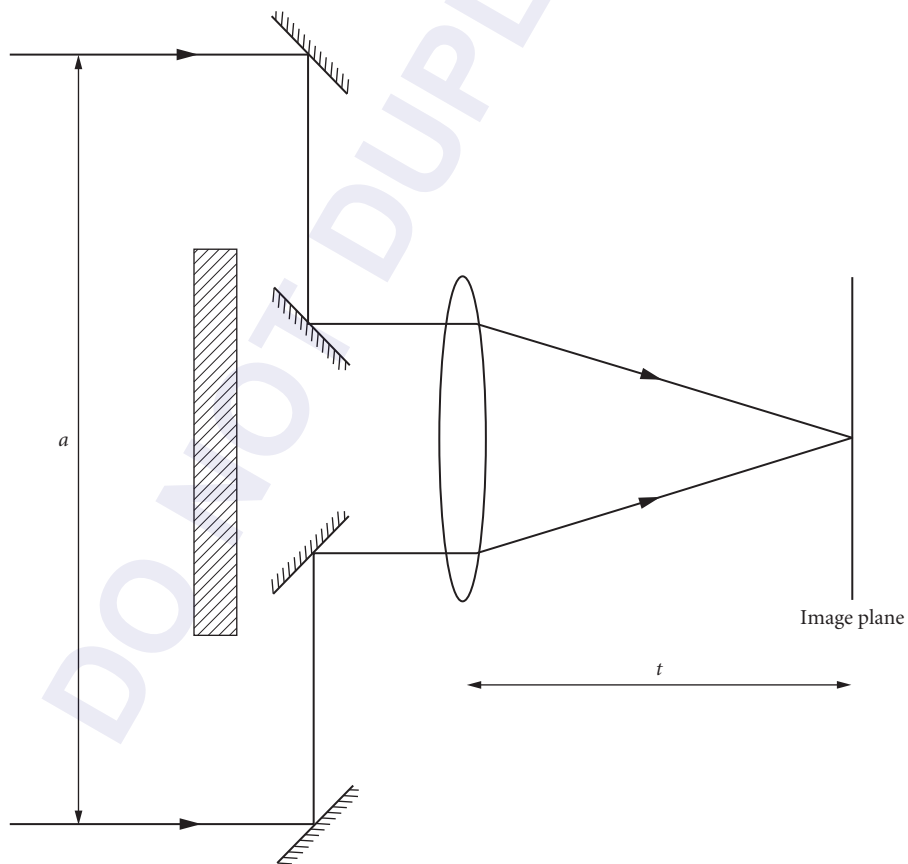


FIGURE 27 Michelson stellar interferometer.

Eq. (86) for a circular star, and the fringe visibility will go to zero when $a = 1.22\lambda\alpha$, where $\alpha = d/z$ is the angular diameter of the star. We measure the mirror separation that produces zero visibility to determine α . In a similar manner, this interferometer can be used to measure the spacing of two closely spaced stars.

Temporal Coherence

When examining temporal coherence effects, we use a source of small dimensions (a point source) that radiates over a range of wavelengths. The light from this source is split into two beams and allowed to interfere. One method to do this is to use an amplitude-splitting interferometer. Since the two sources are identical, the mutual coherence function becomes the *self-coherence function* $\Gamma_{11}(\tau)$. Equal-intensity beams are assumed. The complex degree of temporal coherence becomes

$$\gamma_{11}(\tau) = \frac{\Gamma_{11}(\tau)}{\Gamma_{11}(0)} = \frac{\langle E_1(t+\tau)E_1^*(t) \rangle}{\langle |E_1(t)|^2 \rangle} \quad (87)$$

After manipulation, it follows from this result that $\gamma_{11}(\tau)$ is the normalized Fourier transform of the source intensity spectrum $S(\nu)$:

$$\gamma_{11}(\tau) = \frac{FT\{S(\nu)\}}{\int_0^\infty S(\nu)d\nu} = \frac{\int_0^\infty S(\nu)e^{i2\pi\nu\tau}d\nu}{\int_0^\infty S(\nu)d\nu} \quad (88)$$

The fringe visibility is the modulus of this result. Since $\gamma_{11}(\tau)$ has a maximum at $\tau=0$, the maximum fringe visibility will occur when the time delay between the two beams is zero. This is consistent with our earlier observation under "Source Spectrum" that the fringes will be localized in the vicinity of zero OPD.

As an example, we will repeat the earlier problem of a uniform source spectrum:

$$S(\nu) = \text{rect}\left(\frac{\nu - \nu_0}{\Delta\nu}\right) \quad (89)$$

where ν_0 is the average frequency and $\Delta\nu$ is the bandwidth. The resulting intensity pattern is

$$I = I_0\{1 + \text{Re}\{\gamma_{12}(\tau)\}\} = I_0[1 + \text{sinc}(\tau\Delta\nu)\cos(2\pi\tau\nu_0)] \quad (90)$$

where the sinc function is the Fourier transform of the rect function. Using $\tau = \text{OPD}/c$ from Eq. (74), we can rewrite this equation in terms of the OPD to obtain the same result expressed in Eq. (47).

Laser Sources

The laser is an important source for interferometry, as it is a bright source of coherent radiation. Lasers are not necessarily monochromatic, as they may have more than one longitudinal mode, and it is important to understand the unique temporal coherence properties of a laser in order to get good fringes. The laser is a Fabry-Perot cavity that contains a gain medium. Its output spectrum is therefore a series of discrete frequencies separated by $c/2nL$, where L is the cavity length. For gas lasers, the index is approximately equal to one, and we will use this value for the analysis. If $G(\nu)$ is the gain bandwidth, the frequency spectrum is

$$S(\nu) = G(\nu) \text{comb}\left(\frac{2L\nu}{c}\right) \quad (91)$$

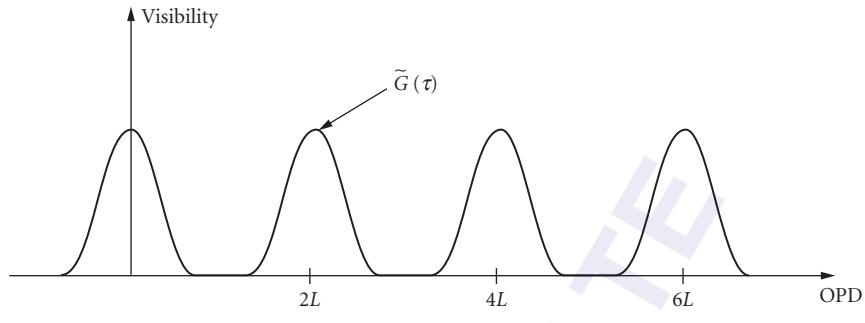


FIGURE 28 The fringe visibility versus OPD for a laser source.

where a comb function is a series of equally spaced delta functions. The number of modes contained under the gain bandwidth can vary from 1 or 2 up to several dozen. The resulting visibility function can be found by using Eq. (88):

$$\gamma(\tau) = |\gamma_{11}(\tau)| = \left| \tilde{G}(\tau) * \text{comb} \left(\frac{c\tau}{2L} \right) \right| = \left| \tilde{G}(\tau) * \text{comb} \left(\frac{\text{OPD}}{2L} \right) \right| \quad (92)$$

where $\tilde{G}(\tau)$ is the normalized Fourier transform of the gain bandwidth, and * indicates convolution. This result is plotted in Fig. 28, where $\tilde{G}(\tau)$ is replicated at multiples of $2L$. The width of these replicas is inversely proportional to the gain bandwidth. We see that as long as the OPD between the two optical paths is a multiple of twice the cavity length, high-visibility fringes will result. This condition is independent of the number of longitudinal modes of the laser. If the laser emits a single frequency, it is a coherent source and good visibility results for any OPD.

2.9 APPLICATIONS OF INTERFERENCE

The fundamental measurement unit associated with interference is the wavelength of light. Every time the OPD in the system changes by one wave, an additional fringe is produced. Because of this sensitivity, interferometers find widespread use in many metrology and optical testing applications. Many of these applications are detailed in subsequent chapters of this *Handbook*, including Chap. 32, “Interferometers,” in this volume, and Chap. 12, “Optical Metrology,” and Chap. 13, “Optical Testing,” in Vol. II. The applications of interferometry include distance and angle measurement, surface figure and finish metrology, profilometry, and spectroscopy. Techniques such as phase-shifting interferometry, heterodyne interferometry, and stitching interferometry have enabled the analysis of the interference patterns associated with the many interferometric measurement techniques in use.

The use of lasers in interferometers has greatly increased their utility. Because of their long coherence length, interference fringes can be produced even when there is a large OPD between the two interfering beams. Instruments such as the Tywman-Green interferometer and the laser-Fizeau interferometer can be used in a compact form to test very large optical surfaces.

2.10 REFERENCES

1. M. Born and E. Wolf, *Principles of Optics*, Pergamon, New York, 1975.
2. R. W. Ditchburn, *Light*, Academic, New York, 1976.

3. M. Francon, *Optical Interferometry*, Academic, New York, 1966.
4. M. H. Freeman, *Optics*, Butterworths, London, 1990.
5. J. D. Gaskill, *Linear Systems, Fourier Transforms and Optics*, Wiley, New York, 1978.
6. E. P. Goodwin and J. C. Wyant, *Field Guide to Interferometric Optical Testing*, SPIE, Bellingham, Wash., 2006.
7. P. Hariharan, *Optical Interferometry*, Academic Press, San Diego, Calif., 1985.
8. P. Hariharan (ed.), *Selected Papers on Interferometry*, SPIE, Bellingham, Wash., 1990.
9. P. Hariharan, *Basics of Interferometry*, Academic Press, San Diego, Calif., 1992.
10. E. Hecht, *Optics*, Addison-Wesley, Reading, Mass., 1989.
11. F. A. Jenkins and H. E. White, *Fundamentals of Optics*, McGraw-Hill, New York, 1976.
12. R. A. Longhurst, *Geometrical and Physical Optics*, Longman, London, 1973.
13. D. Malacara (ed.), *Selected Papers on Optical Shop Metrology*, SPIE, Bellingham, Wash., 1990.
14. D. Malacara (ed.), *Optical Shop Testing*, 2d and 3d editions, Wiley, New York, 1992 and 2007.
15. A. S. Marathay, *Elements of Optical Coherence Theory*, Wiley, New York, 1982.
16. J. R. Meyer-Arendt, *Introduction to Classical and Modern Optics*, Prentice-Hall, Englewood Cliffs, N.J., 1989.
17. G. O. Reynolds, J. B. De Velis, G. B. Parrent, and B. J. Thompson, *The New Physical Optics Notebook: Tutorials in Fourier Optics*, SPIE, Bellingham, Wash., 1989.
18. B. E. A. Saleh and M. C. Teich, *Fundamentals of Photonics*, Wiley, New York, 1991.
19. W. H. Steel, *Interferometry*, Cambridge, London, 1967.
20. F. G. Smith and J. H. Thompson, *Optics*, Wiley, New York, 1971.
21. R. W. Wood, *Physical Optics*, Optical Society of America, Washington D.C., 1988.

This page intentionally left blank.

DO NOT DUPLICATE

DIFFRACTION

Arvind S. Marathay

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

John F. McCalmont

*Air Force Research Laboratory
Sensors Directorate
Wright-Patterson AFB, Ohio*

3.1 GLOSSARY

A	amplitude
E	electric field
f	focal length
G	Green function
E	irradiance
p, q, m	direction cosines
\mathbf{r}	spatial vector
\mathbf{S}	Poynting vector
t	time
ϵ	dielectric constant
μ	permeability
ν	frequency
ψ	wave function
\wedge	Fourier transform

3.2 INTRODUCTION

Starting with waves as solutions to the wave equation obtained from Maxwell's equations, the basics of diffraction of light are covered in this chapter. The discussion includes those applications where the geometry permits analytical solutions. At appropriate locations references are given to the literature and/or textbooks for further reading. The discussion is limited to an explanation of diffraction, and how it may be found in some simple cases with plots of fringe structure.

3.3 LIGHT WAVES

Light waves propagate through free space or a vacuum. They exhibit the phenomenon of diffraction with every obstacle they encounter. Maxwell's equations form a theoretical basis for describing light in propagation, diffraction, scattering and, in general, its interaction with material media. Experience has shown that the electric field \mathbf{E} plays a central role in detection of light and interaction of light with matter. We begin with some mathematical preliminaries.

The electric field \mathbf{E} obeys the wave equation in free space or vacuum:

$$\nabla^2 \mathbf{E} - \frac{1}{c^2} \frac{\partial^2 \mathbf{E}}{\partial t^2} = 0 \quad (1)$$

where c is the velocity of light in vacuum. Each Cartesian component E_j ($j = x, y, z$) obeys the equation and, as such, we use a scalar function $\psi(\mathbf{r}, t)$ to denote its solutions, where the radius vector \mathbf{r} has components, $\mathbf{r} = i\hat{x} + j\hat{y} + k\hat{z}$. The wave equation is a linear second-order partial differential equation. Linear superposition of its linearly independent solutions offers the most general solution. It has traveling plane waves, spherical waves, and cylindrical waves as examples of its solutions. These solutions represent optical wave forms. A frequently used special case of these solutions is the time harmonic version of these waves. We start with the Fourier transform on time,

$$\psi(\mathbf{r}, t) = \int \hat{\psi}(\mathbf{r}, \nu) \exp(-i2\pi\nu t) d\nu \quad (2)$$

where ν is a temporal (linear) frequency in hertz. The spectrum $\hat{\psi}(\vec{r}, \nu)$ obeys the Helmholtz equation,

$$\nabla^2 \hat{\psi} + k^2 \hat{\psi} = 0 \quad (3)$$

with the propagation constant $k = 2\pi/\lambda = 2\pi\nu/c = \omega/c$, where λ is the wavelength and ω is the circular frequency. A Fourier component traveling in a medium of refractive index $n = \sqrt{\epsilon}$, where ϵ is the dielectric constant, is described by the Helmholtz equation with k^2 replaced by $n^2 k^2$. As a further special case, a plane wave may be harmonic in time as well as in space.

$$\psi(\mathbf{r}, t) = A \cos(\mathbf{k} \cdot \mathbf{r} - \omega t) \quad (4)$$

where $\mathbf{k} \equiv k\hat{s}$, \hat{s} is a unit vector in the direction of propagation, and A is a constant. An expanding spherical wave may be written in the form

$$\psi(r, t) = \frac{A}{r} \cos(kr - \omega t) \quad (5)$$

For convenience of operations, a complex function is frequently used. For example, we write

$$\psi(\mathbf{r}, t) = A \exp[i(\mathbf{k} \cdot \mathbf{r} - \omega t)] \quad (6)$$

in place of Eq. (4) bearing in mind that only its real part corresponds to the optical wave form. The function $\psi(\mathbf{r}, t)$ is called the optical "disturbance" while the coefficient A is the amplitude.

In the general discussion of diffraction phenomenon throughout this chapter several classic source books have been used.¹⁻¹⁰ This discussion is a blend of ideas contained in these sources.

The mathematical solutions described heretofore, although ideal, are nevertheless often approximated in practice. A suitable experimental arrangement with a self-luminous source and a condensing lens to feed light to a small enough pinhole fitted with a narrowband spectral filter serves as a quasi-monochromatic, or almost monochromatic, point source. In Fig. 1, light behind the pinhole S is in the form of ever-expanding spherical waves. These waves are of limited spatial extent; all are

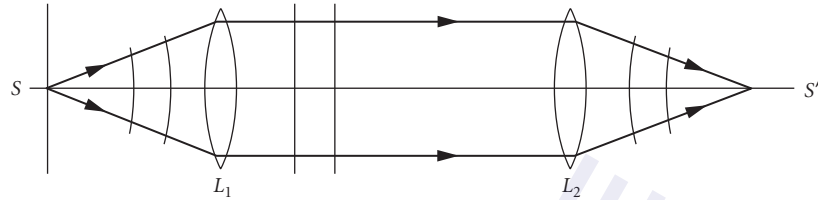


FIGURE 1 Experimental layout to describe the notation used for spherical and plane waves. S : pinhole source, L_1, L_2 : lenses, S' : image.

approximately contained in a cone with its apex at the pinhole. When intercepted by a converging lens L_1 , with the pinhole on its axis and at the front focal point, these spherical waves are converted to plane waves behind L_1 . These plane waves also are limited spatially to the extent dictated by the aperture of the converging lens. A second converging lens, L_2 , is behind the first converging lens and is oriented so that both lenses have a common optical axis and can form an image of the pinhole. The image S' is on the axis at the focal point behind the second lens and is formed by converging spherical waves. These waves, which converge toward the image, are limited spatially to the extent dictated by the aperture of the second lens and are approximately contained in a cone with its apex at the image of the pinhole.

It is necessary to clarify that “a small enough pinhole” means that the optics behind the pinhole are not able to resolve its structure.¹ A “narrowband filter” means that its pass band $\Delta\nu$ is much smaller than the mean frequency $\bar{\nu}$, that is, $\Delta\nu \ll \bar{\nu}$. In this situation, the experimental arrangement may be described by a quasi-monochromatic theory, provided that the path differences Δl of concern in the optics that follow the pinhole are small enough, as given by, $\Delta l \leq c/\Delta\nu$. If the path differences Δl involved are unable to obey this condition, then a full polychromatic treatment of the separate Fourier components contained within $\Delta\nu$ is necessary, even if $\Delta\nu \ll \bar{\nu}$. See, for example, Beran and Parrent¹¹ and Marathay.¹²

Limiting the extent of plane waves and spherical waves, as discussed before, causes diffraction, a topic of primary concern in this chapter. The simplifying conditions stated above are assumed to hold throughout the chapter, unless stated otherwise.

As remarked earlier, the electric field \mathbf{E} [V/m] plays a central role in optical detection. There are detectors that attain a steady state for constant incident beam power [W], and there are those like the *photographic plate* that integrate the incident power over a certain time. For a constant beam power, the darkening of the photographic plate depends on the product of the power and exposure time. Since detectors inherently take the time average, the quantity of importance is the average radiant power [W]. Furthermore, light beams have a finite cross-sectional area, so it is meaningful to talk about the average power in the beam per unit area of its cross-section measured in square meters or square centimeters. In the standard radiometric nomenclature, this sort of measurement is called *irradiance*, E [Wcm^{-2}]. For a plane wave propagating in free space, the irradiance may be expressed in terms of the Poynting vector \mathbf{S} by

$$E = \langle \mathbf{S} \rangle = \left(\frac{1}{2} \right) \left(\frac{\epsilon_0}{\mu_0} \right)^{1/2} \langle \mathbf{E} \cdot \mathbf{E} \rangle \quad (7)$$

The constants given in Eq. (7) may not be displayed with every theoretical result. The Poynting vector and irradiance are discussed further in Ref. 12 (pp. 280–285).

Light is properly described by a transverse vector field. Nevertheless, a scalar field is a convenient artifice to use in understanding the wave nature of light without the added complication of the vector components. The transverse nature of the field will be accounted for when the situation calls for it.

3.4 HUYGENS-FRESNEL CONSTRUCTION

Without the benefit of a fundamental theory based on Maxwell's equations and the subsequent mathematical development, Huygens sought to describe wave propagation in the days before Maxwell. Waves are characterized by constant-phase surfaces called *wavefronts*. If the initial shape at time t of such a wavefront is known in a vacuum or in any medium, Huygens proposed a geometrical construction to obtain its shape at a later time, $t + \Delta t$. He regarded each point of the initial wavefront as the origin of a new disturbance that propagates in the form of secondary wavelets in all directions with the same speed as the speed of propagation of the initial wave in the medium. These secondary wavelets of radii $c\Delta t$ are constructed at each point of the initial wavefront. A surface tangential to all these secondary wavelets, called the *envelope* of all these wavelets, is then the shape and position of the wavefront at time $t + \Delta t$. With this construct, Huygens explained the phenomena of reflection and refraction of the wavefront. To explain the phenomenon of diffraction, Fresnel modified Huygens' construction by attributing the property of mutual interference to the secondary wavelets. The modified Huygens' construction is called the Huygens-Fresnel construction. With further minor modifications it helps explain the phenomenon of diffraction and its various aspects, including those that are not so intuitively obvious.

Fresnel Zones

Let P_0 be a point source of light that produces monochromatic spherical waves. A typical spherical wave, $A/r_0 \exp[-i(\omega t - kr_0)]$, of radius r_0 at time t is shown in Fig. 2. The coefficient A stands for the amplitude of the wave at unit distance from the source P_0 . At a later time this wave will have progressed to assume a position passing through a point of observation P with radius, $r_0 + b$. Fresnel zone construction on the initial wave offers a way to obtain the wave in the future by applying the Huygens-Fresnel construction. The zone construction forms a simple basis for studying and understanding diffraction of light.

From the point of observation P , we draw spheres of radii $b, b + \lambda/2, b + 2\lambda/2, b + 3\lambda/2, \dots, b + j\lambda/2, \dots$, to mark zones on the wave in its initial position, as shown in Fig. 2. The zones are

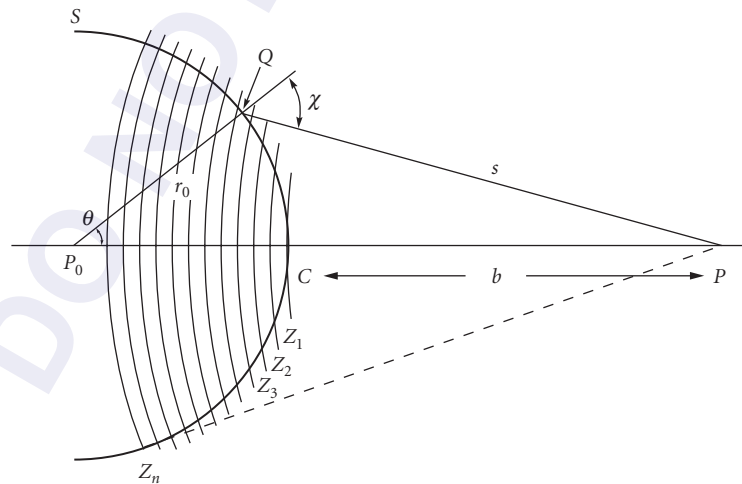


FIGURE 2 Fresnel zone construction. P_0 : point source. S : wavefront. r_0 : radius of the wavefront. b : distance CP . s : distance QP . (After Born and Wolf.¹)

labeled z_1, z_2, \dots, z_j . The zone boundaries are successively half a wavelength away from the point of observation P . By the Huygens-Fresnel construction, each point of the wave forms a source of a secondary disturbance. Each secondary source produces wavelets that are propagated to the point P . A linear superposition of the contribution of all such wavelets yields the resulting amplitude at the point P . It is reasonable to expect that the contribution of the secondary wavelets is not uniform in all directions. For example, a wavelet at C is in line with the source P_0 and the point of observation P , while a wavelet at Q sees the point P at an angle χ with respect to the radius vector from the source P_0 .

To account for this variation, an obliquity or inclination factor $K(\chi)$ is introduced. In the phenomenological approach developed by Fresnel, no special form of $K(\chi)$ is used. It is assumed to have the value unity at C where $\chi = 0$, and it is assumed to decrease at first slowly and then rapidly as χ increases. The obliquity factors for any two adjacent zones are nearly equal and it is assumed that it becomes negligible for zones with high enough index j .

The total contribution to the disturbance at P is expressed as an area integral over the primary wavefront,

$$\psi(P) = A \frac{\exp[-i(\omega t - kr_0)]}{r_0} \iint_S \frac{\exp(iks)}{s} K(\chi) dS \quad (8)$$

where dS is the area element at Q . The subscript S on the integrals denotes the region of integration on the wave surface. The integrand describes the contribution of the secondary wavelets. Fresnel-zone construction provides a convenient means of expressing the area integral as a sum over the contribution of the zones.

For optical problems, the distances involved, such as r_0 and b , are much larger than the wavelength λ . This fact is used very effectively in approximating the integral. The phases of the wavelets within a zone will not differ by more than π . The zone boundaries are successively $\lambda/2$ further away from the point of observation P . The average distance of successive zones from P differs by $\lambda/2$; the zones, therefore, are called half-period zones. Thus, the contributions of the zones to the disturbance at P alternate in sign,

$$\psi(P) = \psi_1 - \psi_2 + \psi_3 - \psi_4 + \psi_5 - \psi_6 + \dots \quad (9)$$

where j stands for the contribution of the j th zone, $j = 1, 2, 3, \dots$. The contribution of each annular zone is directly proportional to the zone area and is inversely proportional to the average distance of the zone to the point of observation P . The ratio of the zone area to its average distance from P is independent of the zone index j . Thus, in summing the contributions of the zones we are left with only the variation of the obliquity factor $K(\chi)$. To a good approximation, the obliquity factors for any two adjacent zones are nearly equal and for a large enough zone index j the obliquity factor becomes negligible. The total disturbance at the point of observation P may be approximated by

$$\psi(P) = 1/2(\psi_1 \pm \psi_n) \quad (10)$$

where the index n stands for the last zone contributing to P . The \pm sign is taken according to whether n is odd or even. For an unobstructed wave, the integration is carried out over the whole spherical wave. In this case, the last term ψ_n is taken to be zero. Thus, the resulting disturbance at the point of observation P equals one-half of the contribution of the first Fresnel zone,

$$\psi(P) = 1/2\psi_1 \quad (11)$$

The contribution ψ_1 is found by performing the area integral of Eq. (8) over the area of the first zone. The procedure results in

$$\psi(P) = \frac{A}{r_0 + b} \lambda \exp\{-i[\omega t - k(r_0 + b) - \pi/2]\} \quad (12)$$

whereas a freely propagating spherical wave from the source P_0 that arrives at point P is known to have the form

$$\psi(P) = \frac{A}{r_0 + b} \exp\{-i[\omega t - k(r_0 + b)]\} \quad (12')$$

The synthesized wave of Eq. (12) can be made to agree with this fact, if one assumes that the complex amplitude of the secondary waves, $\exp(iks)/s$ of Eq. (8) is $[1/\lambda \exp(-i\pi/2)]$ times the primary wave of unit amplitude and zero phase. With the time dependence $\exp(-i\omega t)$, the secondary wavelets are required to oscillate a quarter of a period ahead of the primary.

The synthesis of propagation of light presented above has far-reaching consequences. The phenomenon of light diffraction may be viewed as follows. Opaque objects that interrupt the free propagation of the wave block some or parts of zones. The zones, or their portions that are unobstructed, contribute to the diffraction amplitude (disturbance) at the point of observation P . The obstructed zones do not contribute.

Diffraction of Light from Circular Apertures and Disks

Some examples of unobstructed zones are shown in Fig. 3. Suppose a planar opaque screen with a circular aperture blocks the free propagation of the wave. The center C of the aperture is on the axis joining the source point S and the observation point P , as shown in Fig. 4. The distance and the size of the aperture are such that, with respect to point P , only the first two zones are uncovered as in Fig. 3a. To obtain the diffraction amplitude for an off-axis point such as P , one has to redraw the zone structure as in Fig. 4. Figure 3b shows the zones and parts of zones uncovered by the circular aperture in this case. Figure 3c shows the uncovered zones for an irregularly shaped aperture.

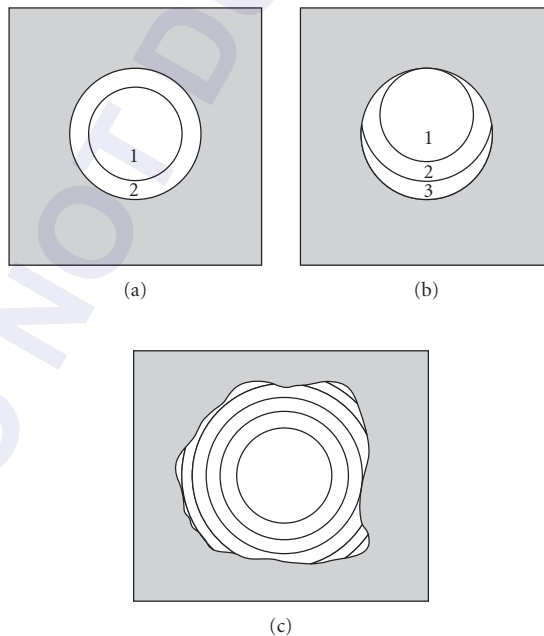


FIGURE 3 Some examples of unobstructed Fresnel zones that contribute to the amplitude at the observation point P . (After Andrews.⁹)

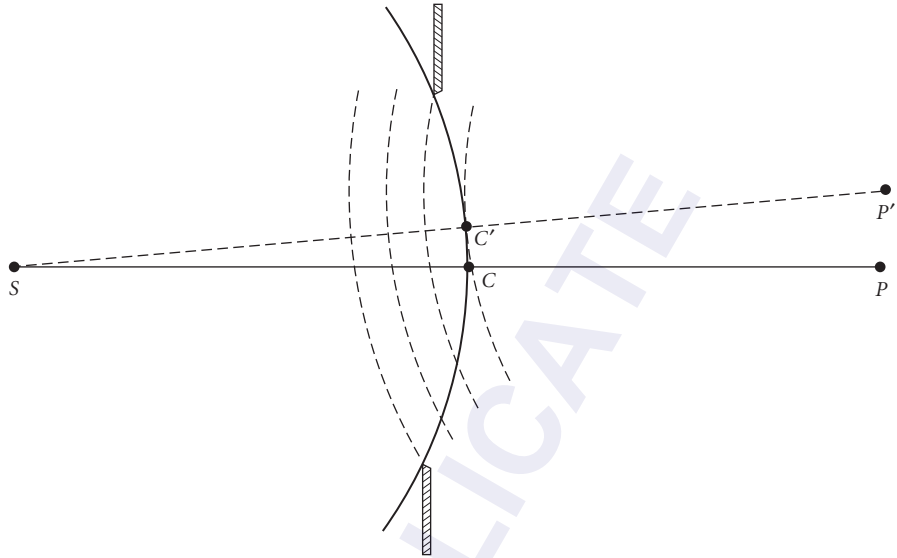


FIGURE 4 The redrawn zone structure for use with an off-axis point P' . (After Andrews.⁹)

In Fig. 3a the first two zones are uncovered. Following Eq. (9), the resulting diffraction amplitude at P for this case is

$$\psi(P) = \psi_1 - \psi_2 \quad (13)$$

but, since these two contributions are nearly equal, the resulting amplitude is $\psi(P) = 0!$

Relocating point P necessitates redrawing the zone structure. The first zone may just fill the aperture if point P is placed farther away from it. In this case the resulting amplitude is

$$\psi(P) = \psi_1 \quad (14)$$

which is twice what it was for the unobstructed wave! Therefore the irradiance is four times as large!

On the other hand, if the entire aperture screen is replaced by a small opaque disk, the irradiance at the center of the geometrical shadow is the same as that of the unobstructed wave! To verify this, suppose that the disk diameter and the distance allows only one Fresnel zone to be covered by the disk. The rest of the zones are free to contribute and *do* contribute. Per Eq. (9) we have

$$\psi(P) = -\psi_2 + \psi_3 - \psi_4 + \psi_5 - \psi_6 + \dots$$

The discussion after Eq. (9) also applies here and the resulting amplitude on the axis behind the center of the disk is

$$\psi(P) = -\frac{1}{2}\psi_2 \quad (15)$$

which is the same as the amplitude of the unobstructed wave. Thus, the irradiance is the same at point P as though the wave were unobstructed. As the point P moves farther away from the disk, the radius of the first zone increases and becomes larger than the disk.

Alternatively one may redraw the zone structure starting from the edge of the disk. The analysis shows that the point P continues to be a bright spot of light. As the point P moves closer to the disk, more and more Fresnel zones get covered by the disk, but the analysis continues to predict a bright spot at P . There comes a point where the unblocked zone at the edge of the disk is significantly weak; the point P continues to be bright but has reduced irradiance. Still closer to the disk, the analysis ceases to apply because P enters the near-field region, where the distances are comparable to

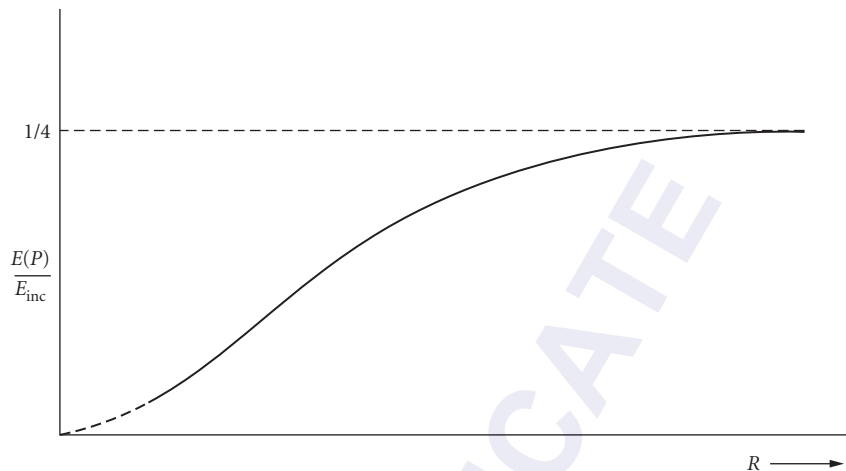


FIGURE 5 Variation of on-axis irradiance behind an opaque disk. R : distance along the axis behind the disk. (From Marion and Heald.⁵)

the size of the wavelength. In Fig. 5, the variation of irradiance on the axial region behind the disk is shown. It is remarkable that the axial region is nowhere dark! For an interesting historical note, see Refs. 1 and 5.

For comparison, we show the variation of on-axis irradiance behind a circular opening in Fig. 6. It shows several on-axis locations where the irradiance goes to zero. These correspond to the situation where an even number of zones are exposed through the circular aperture. Only the first few zeros are shown, since the number of zeros per unit length (linear density) increases as the point P is moved closer to the aperture. The linear density increases as the square of the index j when P moves

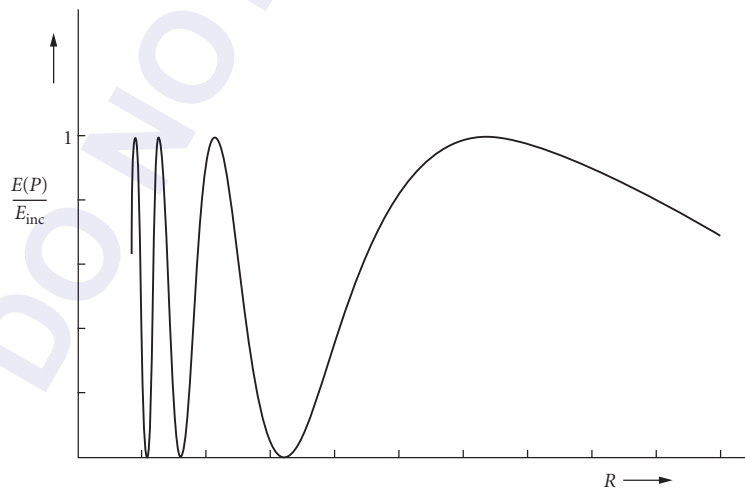


FIGURE 6 Variation of on-axis irradiance behind a circular opening. R : distance along the axis behind the opening. (From Marion and Heald.⁵)

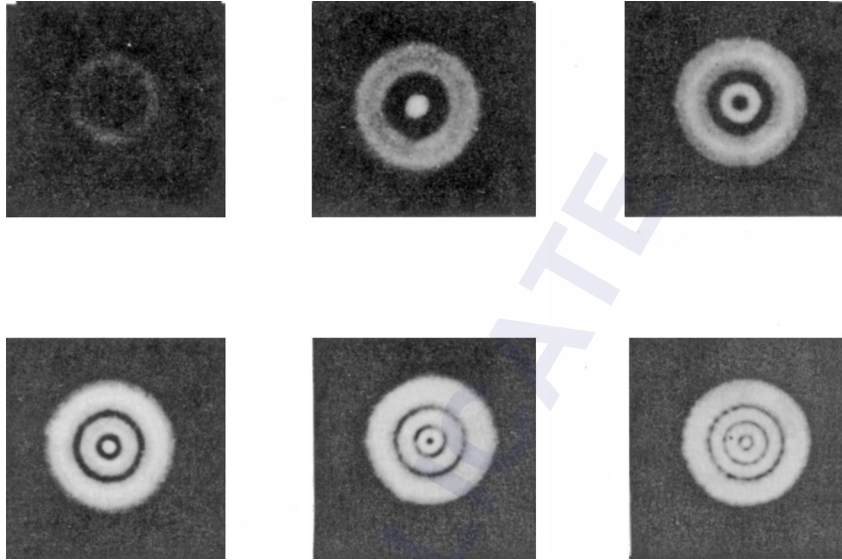


FIGURE 7 A series of pictures of diffraction patterns from circular apertures. (After Andrews.⁹)

closer to the aperture. While far enough away, there comes a point where the first zone fills the aperture and, thereafter, there are no more zeros as the distance increases.

Figure 7 shows a series of diffraction patterns from a circular aperture. The pictures are taken at different distances from the aperture to expose one, two, three, etc., zones. Each time an odd number of zones is uncovered the center spot becomes bright. As we approach the pictures at the bottom right, more zones are exposed.

Babinet Principle

Irradiances for the on-axis points are quite different for the circular disk than for the screen with a circular opening. The disk and the screen with a hole form of a pair of complementary screens, that is, the open areas of one are the opaque areas of the other and vice versa. Examples of pairs of such complementary screens are shown in Fig. 8. Observe that the open areas of screen S_a taken with the open areas of the complementary screen S_b add up to no screen at all.

The Babinet principle states that the wave disturbance $\psi_s(P)$ at any point of observation P due to a diffracting screen S_a added to the disturbance $\psi_{CS}(P)$ due to the complementary screen S_b at the same point P equals the disturbance at P due to the unobstructed wave, that is,

$$\psi_s(P) + \psi_{CS}(P) = \psi_{UN}(P) \quad (16)$$

Recall that the wave disturbance at any point of observation P behind the screen is a linear superposition of the contributions of the unobstructed zones or portions thereof. This fact, with the observation that the open areas of screen S_a taken with the open areas of the complementary screen S_b add up to no screen at all, implies the equality indicated by the Babinet principle.

The application of Babinet's principle to diffraction problems can reduce the complexity of the analysis considerably. For an example of this, we once again return to diffraction of light due to a circular aperture and an opaque disk. For on-axis amplitude, the Rayleigh-Sommerfeld diffraction integral [see Eq. (43) later] can be evaluated in closed form.⁶

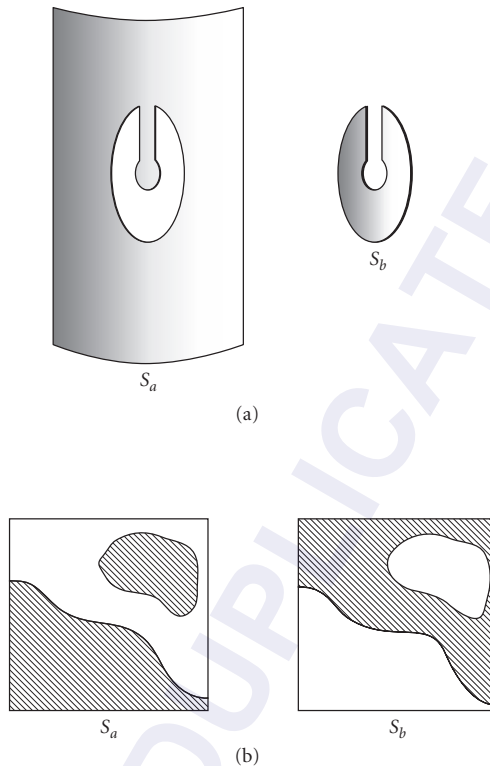


FIGURE 8 Examples of complementary screens, labeled S_a and S_b . (After Jackson¹³ and Andrews.⁹)

A circular aperture of radius a in a dark screen is illuminated by a normally incident plane wave, $A \exp(ikz)$, where z is the axis perpendicular to the plane of the aperture. The on-axis diffracted field ψ_{cir} is

$$\psi_{\text{cir}}(z) = A \exp(ikz) - \frac{Az}{\sqrt{a^2 + z^2}} \exp\left(ik\sqrt{a^2 + z^2}\right) \quad (17)$$

The corresponding irradiance at z is

$$E_{\text{cir}}(z) = |\psi_{\text{cir}}|^2 = A^2 \left[1 + \frac{z^2}{a^2 + z^2} - \frac{2z}{\sqrt{a^2 + z^2}} \cos\left(k\left(\sqrt{a^2 + z^2} - z\right)\right) \right] \quad (18)$$

See Fig. 9 for the on-axis irradiance plotted against the distance z from the circular aperture.

Next consider the dark screen and circular aperture is replaced by an opaque disk of the same radius a . Applying Babinet's principle, the on-axis diffracted field for the opaque disk illuminated by a normally incident plane wave is

$$\psi_{\text{disk}} = A \exp(ikz) - \psi_{\text{cir}} \quad (19)$$

$$\psi_{\text{disk}}(z) = \frac{Az}{\sqrt{a^2 + z^2}} \exp\left(ik\sqrt{a^2 + z^2}\right) \quad (20)$$

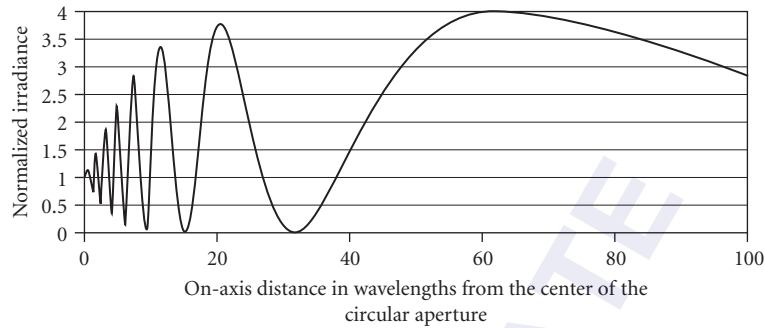


FIGURE 9 Normalized on-axis irradiance behind a circular aperture of radius 8λ , plotted as a function of z , the distance from the aperture in terms of wavelengths. At $z = 0$, the assumed boundary condition of unity is obtained. The maxima of the oscillations increases monotonically to a maximum value of 4 at $z = d^2/4\lambda$ when one Fresnel zone fills the aperture. For $z > d^2/4\lambda$ the irradiance decreases monotonically to zero.

The corresponding irradiance at z is

$$E_{\text{disk}}(z) = |\psi_{\text{disk}}|^2 = A^2 \frac{z^2}{a^2 + z^2} \quad (21)$$

See Fig. 10 for the variation of the on-axis irradiance with distance z from the opaque disk of radius a .

The behavior of the on-axis irradiance for the case of the opaque disk is quite different from that of the complementary circular aperture. There is no simple relationship between the irradiances of the two cases because they involve a squaring operation that brings in cross-terms.

It is important to note that the closed form expressions of Eq. (17) through Eq. (21) are valid only within the approximations of the Rayleigh-Sommerfeld theory. The value of unity obtained by Eq. (18) reproduces the assumed boundary conditions of the theory.

Zone Plate

If alternate zones are blocked the contribution of the unblocked zones will add in phase to yield a large irradiance at the point of observation. An optical device that blocks alternate zones is called a *zone plate*. Figure 11 shows two zone plates made up of concentric circles with opaque alternate

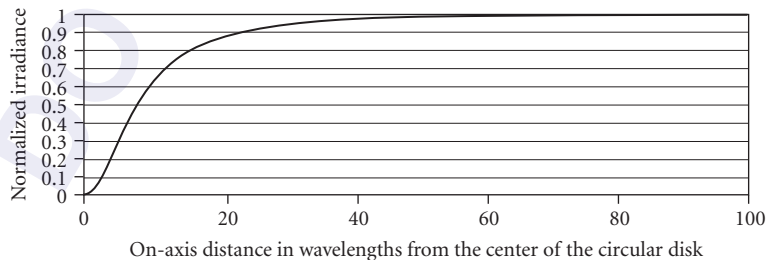


FIGURE 10 Normalized on-axis irradiance behind a disk of radius 8λ , plotted as a function of z , the distance from the aperture in terms of wavelengths. The irradiance is zero at $z = 0$ and increases monotonically to unity for large z .

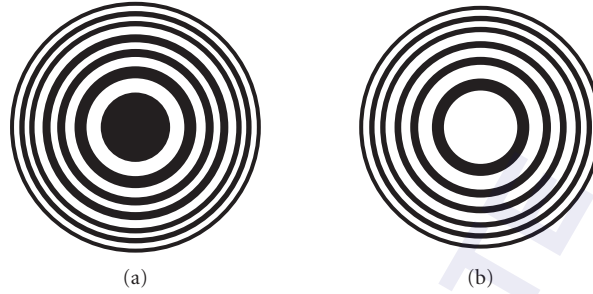


FIGURE 11 Two zone plates made up of concentric circles with alternate zones made opaque. They block odd-indexed or even-indexed zones, respectively. The radii of the zone boundaries are proportional to the square root of natural numbers. (From Hecht and Zajac.¹⁰)

zones. They block odd-indexed or even-indexed zones, respectively. The radii of the zone boundaries are proportional to the square root of natural numbers.

We place a point source at a distance r_0 in front of the zone plate. If R_m is the radius of the m th zone, a bright image of this source is observed at a distance b behind the plate, so that

$$\frac{1}{r_0} + \frac{1}{b} = \frac{m\lambda}{R_m^2} \quad (22)$$

where λ is the wavelength of light from the source. This equation for the condition on the distance b is like the paraxial lens formula from which the focal length of the zone plate may be identified or may be obtained by setting the source distance $r_0 \rightarrow \infty$.

The focal length f_1 so obtained is

$$f_1 = \frac{R_m^2}{m\lambda} \quad (23)$$

and is called the *primary focal length*. For unlike, the case of the lens, the zone plate has several secondary focal lengths. These are given by

$$f_{2n-1} = \frac{R_m^2}{(2n-1)m\lambda} \quad (24)$$

where $n = 1, 2, 3, \dots$. In the case of the primary focal length, each opaque zone of the zone plate covers exactly one Fresnel zone. The secondary focal length f_3 is obtained when each opaque zone covers three Fresnel zones. It is a matter of regrouping the right-hand side of Eq. (9) in the form

$$\begin{aligned} \psi(P) = & (\psi_1 - \psi_2 + \psi_3) - \langle \psi_4 - \psi_5 + \psi_6 \rangle + (\psi_7 - \psi_8 + \psi_9) \\ & - \langle \psi_{10} - \psi_{11} + \psi_{12} \rangle + (\psi_{13} - \psi_{14} + \psi_{15}) - \dots \end{aligned} \quad (25)$$

The zone plate in Fig. 11b, for example, blocks all even-indexed zones. It corresponds to omitting the terms enclosed in the angular brackets, $\langle \dots \rangle$ in Eq. (25). The remaining terms grouped in parentheses add in phase to form a secondary image of weaker irradiance. The higher-order images are formed successively closer to the zone plate and are successively weaker in irradiance.

Further discussion may be found in several books listed in the references, for example, Ref. 10 (p. 375). The radii of the concentric circles in a zone plate are proportional to the square root of natural numbers. For equidistant source and image locations, say 10 cm at a wavelength of 500 nm,

$R_m = \sqrt{m} \times 0.16$ mm. Due to the smallness of the radii, a photographic reduction of a large-scale drawing is used.

Incidentally, the pair of zone plates of Fig. 11 form a pair of complementary screens. Per the Babinet principle, the groupings are

$$\begin{aligned}\psi_{\text{UN}}(P) &= \psi_s(P) + \psi_{\text{CS}}(P) \\ &= (\psi_1 + \psi_3 + \psi_5 + \psi_7 + \psi_9 + \psi_{11} + \dots) \\ &\quad - (\psi_2 + \psi_4 + \psi_6 + \psi_8 + \psi_{10} + \psi_{12} + \dots)\end{aligned}\quad (26)$$

The first group of terms corresponds to the zone plate of Fig. 11*b* and the second group of items corresponds to Fig. 11*a*.

3.5 CYLINDRICAL WAVEFRONT

A line source generates cylindrical wavefronts. It is frequently approximated in practice by a slit source, which, in turn, can illuminate straight edges and rectangular or slit apertures (see Fig. 12*a*). In this case, as we shall see, the phenomena of diffraction can be essentially reduced to a one-dimensional analysis for this source and aperture geometry.

Fresnel zones for cylindrical wavefronts take the form of rectangular strips, as shown in Fig. 12*a*. The edges of these strip zones are $\lambda/2$ farther away from the point of observation P . The treatment for the cylindrical wave parallels the treatment used for the spherical wave in Sec. 3.4. The line M_0 on the wavefront intersects at right angles to the line joining the source S and the point of observation P . Refer to M_0 as the axis line of the wavefront with respect to the point P . Let a be the radius of the wavefront with respect to the source slit and let b be the distance of P from M_0 . Fresnel zones are now in the form of strips above and below M_0 and are parallel to it. The line pairs $M_1M'_1$, $M_2M'_2$, etc., are marked $\lambda/2$ farther away from the point of observation P . Fresnel zones are now half-period strips. Thus $PM_m = b + m\lambda/2$ and, to a good approximation, the arc length $(M_mM_{m+1}) = \sqrt{mab\lambda/(a+b)(\sqrt{m+1} - \sqrt{m})}$. For small values of m such as 1, 2, etc., the arc widths decrease rapidly while, for large values of m the neighboring strips have nearly equal widths. The lower-order strips have much larger areas compared to the ones further up from M_0 . This effect is much more dominant than the variation of the obliquity factor $K(\chi)$ which has been neglected in this analysis.

Consider one single strip as marked in Fig. 12*b*. Imagine that this strip is divided into half-period sections as shown. The wavefront is almost planar over the width of this strip. All the sections on either side of the arc M_1M_2 contribute to the disturbance at P . The boundaries of these are marked N_1, N_2 , etc. The area of these sections are proportional to $\sqrt{c\lambda}(\sqrt{n+1} - \sqrt{n})$. The areas of those half-period sections decrease rapidly at first and then slowly. The contribution to the disturbance at P from the higher-order sections is alternately positive and negative with respect to the first section near M_1M_2 . Consequently, their contribution to the total disturbance at P is nearly zero.

The disturbance at P due to a single strip consists of the dominant contribution of the two sections from N_1 to N_1' . This conclusion holds for all the strips of the cylindrical wave. Following the procedure of Eq. (9) employed for the spherical wave,

$$\psi(P) = \psi_1 - \psi_2 + \psi_3 - \psi_4 + \psi_5 - \psi_6 + \dots \quad (9')$$

where $\psi(P)$ is the disturbance at P and m denotes the secondary wavelet contributions from strip zones of either side of the axis line M_0 of Fig. 12*a*. As in Eq. (11) the series can be summed, but here we need to account for the strip zone contribution from both sides of the axis line M_0 ; therefore, we have

$$\psi(P) = \psi_1 \quad (11')$$

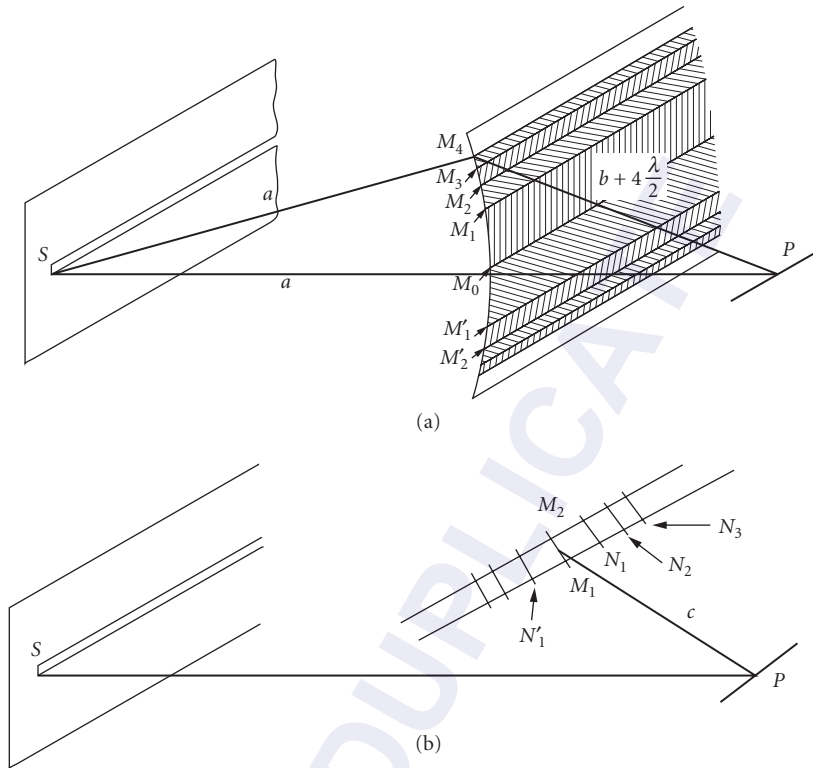


FIGURE 12 Fresnel zones for cylindrical wavefront. *S*: slit source. *a*: radius of the cylindrical wavefront. *b*: distance M_0P . M_1, M'_1 , etc.: zone boundaries. (Adapted from Jenkins and White.²)

The first zone contributions can be computed and compared with the freely propagating cylindrical wave.

Fresnel Diffraction from Apertures with Rectangular Symmetry

Straight Edge A cylindrical wave from a slit source *S* illuminates an opaque screen with a straight edge *AB* oriented parallel to the slit, as shown in Fig. 13. It shows three special positions of the point of observation *P*. In Fig. 13*a*, *P* is such that all the strip zones above the axis line M_0 are exposed, while those below are blocked. The point *P* in Fig. 13*b* is such that the strip zones above M_0 and one zone below, marked by the edge M_1 , are exposed. In Fig. 13*c*, *P* has moved into the geometrical shadow region. The strip M_1M_0 and all those below M_0 are blocked.

Following the discussion in Sec. 3.4, we discuss the disturbance at *P* for the three cases of Fig. 13. At the edge of the geometrical shadow

$$\psi_a(P) = 1/2\psi_1 \tag{27}$$

For Fig. 13*b* we have

$$\psi_b(P) = 1/2\psi_1 + \psi_1 = 3/2\psi_1 \tag{28}$$

As *P* is moved further up, there comes a point for which two strip zones below M_0 are exposed, resulting in $\psi(P) = 3/2\psi_1 - \psi_2$. As *P* explores the upper half of the observation plane, the amplitude

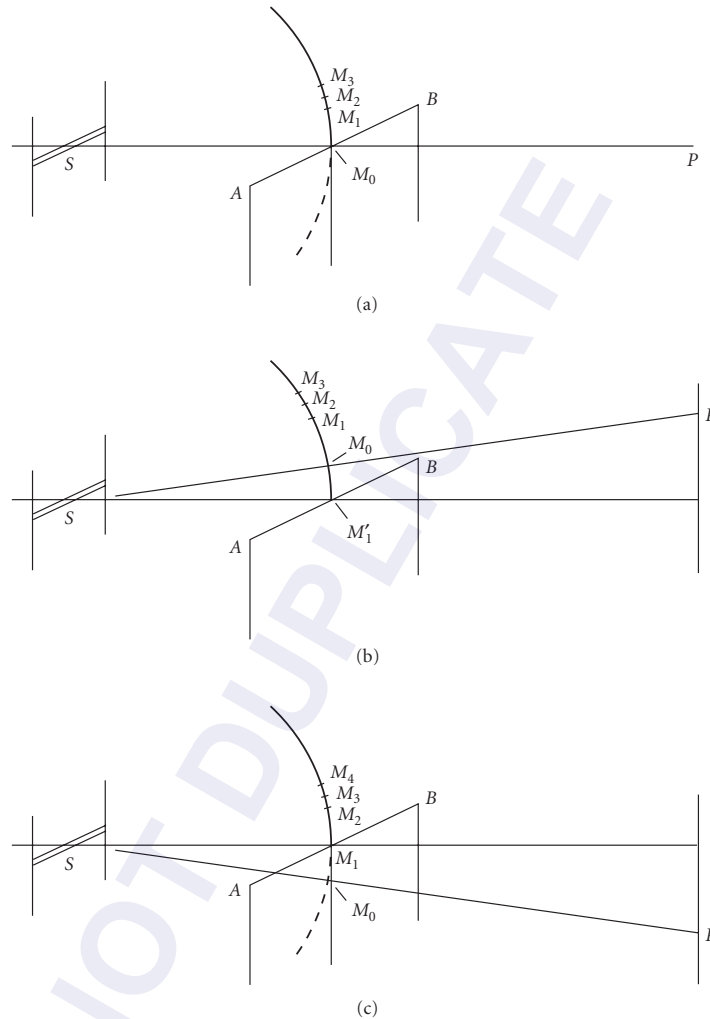


FIGURE 13 Fresnel zones for a cylindrical wavefront. The edges of these strip zones are $\lambda/2$ farther away from the point of observation P . S : slit source. P : point of observation. M_0 : axis line of the cylindrical wave. AB : straight-edge opaque obstruction. (After Jenkins and White.²)

and, hence, the irradiance goes through maxima and minima according to whether an odd or even number of (lower) strip zones is exposed. Furthermore, the maxima decrease gradually while the minima increase gradually until the fringes merge into a uniform illumination that corresponds to the unobstructed wave.

In the geometrical shadow (see Fig. 13c)

$$\psi_c(P) = -1/2\psi_2 \quad (29)$$

As P goes further down, we get $\psi(P) = 1/2\psi_3$; in general, the number of exposed zones decreases and the irradiance falls off monotonically.

A mathematical analysis of Fresnel diffraction from apertures with rectangular symmetry is possible with the use of Fresnel integrals and the Cornu's spiral (vibration curve). Irradiance of the

diffraction pattern from a straight edge illuminated by a cylindrical wave is found by relating the spiral to the cylindrical wavefront and to the plane of observation. Figure 13 shows the diffraction geometry. The line or slit source is aligned parallel to the straight edge. The on-axis point of observation P is as shown in Fig. 13 and is defined as $P(0)$. An off-axis point of observation (see Fig. 13b and c) is defined as $P(x)$ at x . The plane passing through $P(0)$ and $P(x)$ is perpendicular to the edge. As shown in Fig. 12a, the radius of the cylindrical wave is a and the distance between M_0 and $P(0)$ is b . The distance of a point on the wavefront above M_0 labeled M_m to the point $P(0)$ is $b' = b + m\lambda/2$. The path difference $b' - b = \Delta$. To a good approximation,

$$\Delta \cong s^2 \frac{a+b}{2ab} \quad (30)$$

where s is the arc length between M_0 and M_m but is approximated by the corresponding chord length. For computational purposes the arc length s is converted to a dimensionless parameter v by defining the phase difference as

$$\delta = \frac{2\pi}{\lambda} \Delta = \frac{2\pi}{\lambda} s^2 \frac{a+b}{2ab} \cong \frac{\pi}{2} v^2 \Rightarrow v = s \sqrt{\frac{2(a+b)}{\lambda ab}} \quad (31)$$

With Fig. 12 we can relate the arc length s with the coordinate x in the plane of observation,

$$x = s \left(\frac{a+b}{a} \right) \quad (32)$$

and therefore,

$$v = s \sqrt{\frac{2(a+b)}{\lambda ab}} = x \sqrt{\frac{2a}{\lambda b(a+b)}} \quad (33)$$

Cornu's Spiral The diffraction amplitude [see Eq. (8)] is given by

$$\psi(P) = A \frac{\exp[-i(\omega t - kr_0)]}{r_0} \iint_s \frac{\exp(iks)}{s} K(\chi) dS \quad (34)$$

For a cylindrical wave illumination of the straight edge, the double integral reduces to a single integral. The single integral describes integration along a line on the cylindrical wavefront parallel to the line source. For $r_0 = a$ and $s = b$ much larger than the width of the strip of Fresnel zone on the cylindrical wavefront, we use the approximation,

$$\psi(P) = C \frac{\exp(ik(a+b))}{ab} K(\chi) \int_{s_1}^{s_2} \exp(ik\Delta) ds \quad (35)$$

In this expression, $C = A \exp(-i\omega t)$. The formal limits of integration s_1 and s_2 designate the limit of integration appropriate for a slit aperture oriented parallel to the line source. For the straight edge problem, the upper limit is ∞ and the lower limit s_1 describes the location of the straight-edge with respect to the position of the point of observation P along the x axis. Thus, the amplitude is proportional to

$$\psi(P) = C' \left[\int_{-\infty}^{s_1} \exp(ik\Delta) dS \right] = C' \left[\int_{-\infty}^0 \exp(ik\Delta) dS + \int_0^{s_1} \exp(ik\Delta) dS \right] \quad (36)$$

$$\psi(P) = C' \left\{ \left[\frac{1}{2} + C(v_1) \right] + i \left[\frac{1}{2} + S(v_1) \right] \right\} \quad (37)$$

$$E(P) = |\psi(P)|^2 = |C'|^2 \left\{ \left[\frac{1}{2} + C(v_1) \right]^2 + \left[\frac{1}{2} + S(v_1) \right]^2 \right\} \quad (38)$$

We have defined

$$C' = C \frac{\exp(ik(a+b))}{ab} K(\chi) \quad \text{and} \quad |C'|^2 = \frac{|C|^2 K(\chi)^2}{a^2 b^2}$$

For the case of a cylindrical wave incident, Cornu's spiral helps in the calculation of the diffraction amplitude (Plane wave incidence is a special case in the limit where radius a of the cylindrical wave becomes very large compared to the distance b of the observation screen.) Cornu's spiral is defined in terms of the dimensionless parameter v . The description of this spiral begins with the definitions

$$C(v) = \int_0^v \cos\left(\frac{\pi}{2}v^2\right) dv \quad S(v) = \int_0^v \sin\left(\frac{\pi}{2}v^2\right) dv \quad (39)$$

1. Cornu's spiral is a plot of $S(v)$ along the vertical axis and $C(v)$ along the horizontal axis with v as a parameter, $-\infty \leq v \leq \infty$.
2. The arc length measured along the spiral from the origin is v ; $\delta v = \sqrt{[\delta C(v)]^2 + [\delta S(v)]^2}$; $\delta C(v)$ and $\delta S(v)$ are the projections of the arc length on the horizontal and vertical axes, respectively.
3. The vector length from the origin to any point v on the spiral is proportional to diffracted amplitude.
4. The angle ϕ made by the vector measured from the horizontal equals the phase of the diffracted light;

$$\tan \phi = \frac{\delta S}{\delta C} = \tan \frac{\pi}{2} v^2 \Rightarrow \phi = \frac{\pi}{2} v^2.$$

5. The radius of curvature = $\frac{dv}{d\phi} = \frac{1}{\pi v}$.
6. For large v , the spiral winds about two limit points

$$v \rightarrow +\infty \quad C(+\infty) = +\frac{1}{2} \quad S(+\infty) = +\frac{1}{2}$$

$$v \rightarrow -\infty \quad C(-\infty) = -\frac{1}{2} \quad S(-\infty) = -\frac{1}{2}$$

7. The magnitude of the diffraction integral has its maximum value when

$$\phi = \frac{3\pi}{4} \Rightarrow v = \sqrt{\frac{3}{2}}.$$

8. Subsidiary maxima at $\phi = \frac{3\pi}{4} + 2\pi n \Rightarrow v = \sqrt{\frac{3}{2} + 4n}$, where $n = 1, 2, 3, \dots$

9. Minimum values at $\phi = \frac{7\pi}{4} + 2\pi m \Rightarrow v = \sqrt{\frac{7}{2} + 4m}$, where $m = 0, 1, 2, 3, \dots$

Cornu's spiral is plotted in Fig. 14 and tabulated in Table 1.

Figure 15 shows a plot of the irradiance in the diffraction pattern of light from a straight edge. The maxima and minima follow the description given in 8 and 9 in the above list. It is interesting to observe that the edge is neither at the maximum of the first fringe nor at the halfway point. It appears at one-fourth of the irradiance of the unobstructed wave. For a complete vectorial solution of diffraction of light by a straight edge, see Sommerfeld.¹⁴

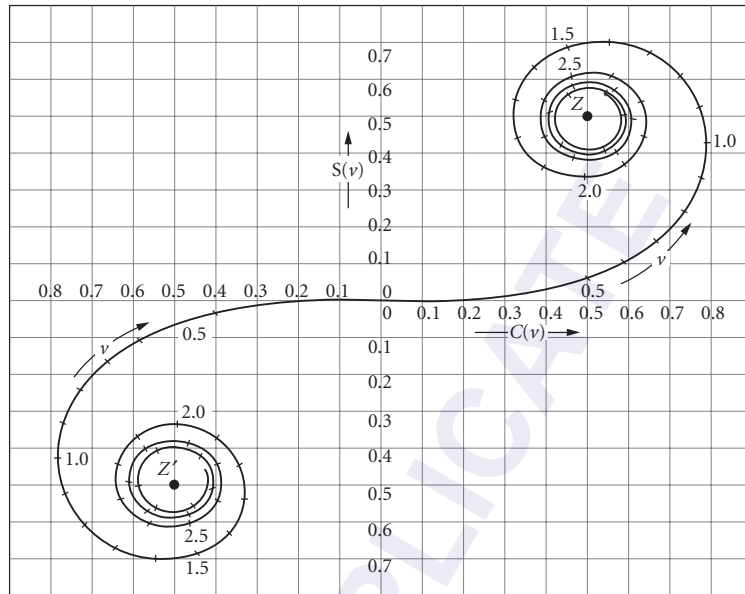


FIGURE 14 Cornu's spiral (vibration curve) for use with cylindrical waves and apertures with rectangular symmetry. (Adapted from Jenkins and White.²)

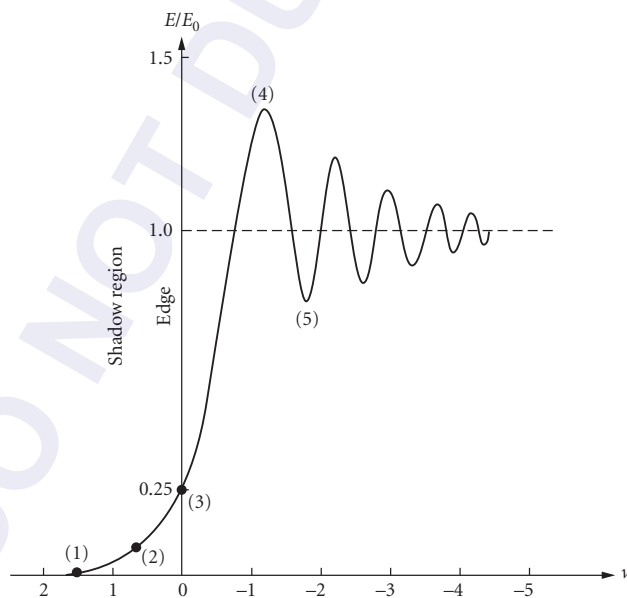


FIGURE 15 The plot of irradiance in the diffraction pattern of a straight edge AB . The plot is normalized to unity for the irradiance of the unobstructed wave. Labels 1 and 2 show points P in the geometrical shadow. Label 3 is at the edge of the geometrical shadow, while labels 4 and 5 are in the illuminated region. v is a unitless variable to label distances along the plane of observation. (From Hecht and Zajac.¹⁰)

TABLE 1 Table of Fresnel Integrals

ν	$C(\nu)$	$S(\nu)$	ν	$C(\nu)$	$S(\nu)$
0.00	0.0000	0.0000	4.50	0.5261	0.4342
0.10	0.1000	0.0005	4.60	0.5673	0.5162
0.20	0.1999	0.0042	4.70	0.4914	0.5672
0.30	0.2994	0.0141	4.80	0.4338	0.4968
0.40	0.3975	0.0334	4.90	0.5002	0.4350
0.50	0.4923	0.0647	5.00	0.5637	0.4992
0.60	0.5811	0.1105	5.05	0.5450	0.5442
0.70	0.6597	0.1721	5.10	0.4998	0.5624
0.80	0.7230	0.2493	5.15	0.4553	0.5427
0.90	0.7648	0.3398	5.20	0.4389	0.4969
1.00	0.7799	0.4383	5.25	0.4610	0.4536
1.10	0.7638	0.5365	5.30	0.5078	0.4405
1.20	0.7154	0.6234	5.35	0.5490	0.4662
1.30	0.6386	0.6863	5.40	0.5573	0.5140
1.40	0.5431	0.7135	5.45	0.5269	0.5519
1.50	0.4453	0.6975	5.50	0.4784	0.5537
1.60	0.3655	0.6389	5.55	0.4456	0.5181
1.70	0.3238	0.5492	5.60	0.4517	0.4700
1.80	0.3336	0.4508	5.65	0.4926	0.4441
1.90	0.3944	0.3734	5.70	0.5385	0.4595
2.00	0.4882	0.3434	5.75	0.5551	0.5049
2.10	0.5815	0.3743	5.80	0.5298	0.5461
2.20	0.6363	0.4557	5.85	0.4819	0.5513
2.30	0.6266	0.5531	5.90	0.4486	0.5163
2.40	0.5550	0.6197	5.95	0.4566	0.4688
2.50	0.4574	0.6192	6.00	0.4995	0.4470
2.60	0.3890	0.5500	6.05	0.5424	0.4689
2.70	0.3925	0.4529	6.10	0.5495	0.5165
2.80	0.4675	0.3915	6.15	0.5146	0.5496
2.90	0.5624	0.4101	6.20	0.4676	0.5398
3.00	0.6058	0.4963	6.25	0.4493	0.4954
3.10	0.5616	0.5818	6.30	0.4760	0.4555
3.20	0.4664	0.5933	6.35	0.5240	0.4560
3.30	0.4058	0.5192	6.40	0.5496	0.4965
3.40	0.4385	0.4296	6.45	0.5292	0.5398
3.50	0.5326	0.4152	6.50	0.4816	0.5454
3.60	0.5880	0.4923	6.55	0.4520	0.5078
3.70	0.5420	0.5750	6.60	0.4690	0.4631
3.80	0.4481	0.5656	6.65	0.5161	0.4549
3.90	0.4223	0.4752	6.70	0.5467	0.4915
4.00	0.4984	0.4204	6.75	0.5302	0.5362
4.10	0.5738	0.4758	6.80	0.4831	0.5436
4.20	0.5418	0.5633	6.85	0.4539	0.5060
4.30	0.4494	0.5540	6.90	0.4732	0.4624
4.40	0.4383	0.4622	6.95	0.5207	0.4591

Note: This table is adapted from Jenkins and White.²

Rectangular Aperture Figure 16 is a series of diagrams of irradiance distributions for light diffracted by single-slit apertures. A pair of marks on the horizontal axis indicate the edges of the geometrical shadow of the slit relative to the diffraction pattern. In all cases, relatively little light falls in the geometrical shadow region. The last diagram corresponds to a rather wide slit. It appears as two opposing straightedge diffraction patterns corresponding to the two edges of the slit.

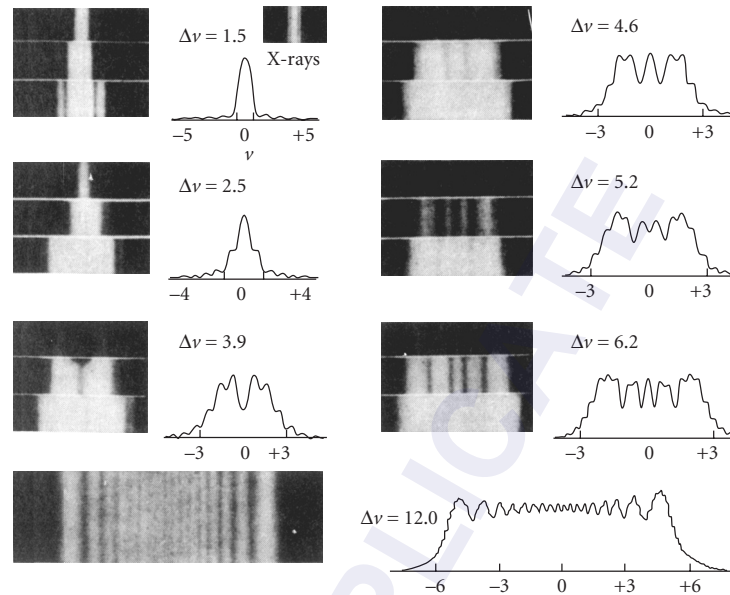


FIGURE 16 A series of diagrams of irradiance distributions for light diffracted by single-slit apertures of different widths. (From Jenkins and White.²)

These patterns may be interpreted as obtained with the plane of observation fixed for different-size slits. Alternately, the slit size may be held fixed but move the plane of observation. For the first diagram the plane is far away. For the successive diagrams the plane is moved closer to the slit. The plane of observation is the closest for the last diagram of Fig. 16. The important parameter is the angular subtense of the slit to the observation plane. A similar comment applies to the case of the circular aperture^{2,9} as shown in Fig. 7.

Opaque Strip Obstruction A slit aperture and an opaque strip or a straight wire form a pair of complementary screens. In Fig. 17 photographs of Fresnel diffraction patterns produced by narrow wires are shown with the corresponding theoretical curves. These theoretical curves show some more detail. Generally, the figures show the characteristic unequally spaced diffraction fringes of a straight edge on either side of the geometrical shadow. These fringes get closer and closer together, independent of the width of the opaque obstruction, and finally merge into a uniform illumination. Figure 17 also shows the maximum in the center and equally spaced narrow fringes within the shadow. The width of these fringes is inversely proportional to the width of the obstruction. We shall now discuss this detail.

Figure 18 shows the arrangement of the source S , opaque strip AB , and the plane of observation. A point x in the geometrical shadow receives light from Fresnel zones of both sides of the opaque strip. At each edge of the opaque strip the exposed zones add up effectively to one-half of the contribution of a single zone adjacent to that edge. Owing to the symmetry, the resulting disturbance from each edge starts out in phase. Light from the two edges adds constructively or destructively according to whether the path difference to the point x in the shadow region is an even or an odd multiple of $\lambda/2$. The situation is similar to two coherent sources separated by the width of the opaque strip. Young examined these fringes inside the geometrical shadow. In particular, he showed that if an opaque screen is introduced on one side of the opaque strip to block that part of the wave, then the straightedge diffraction fringes due to that edge, as well as the interference fringes in the shadow region, vanished.

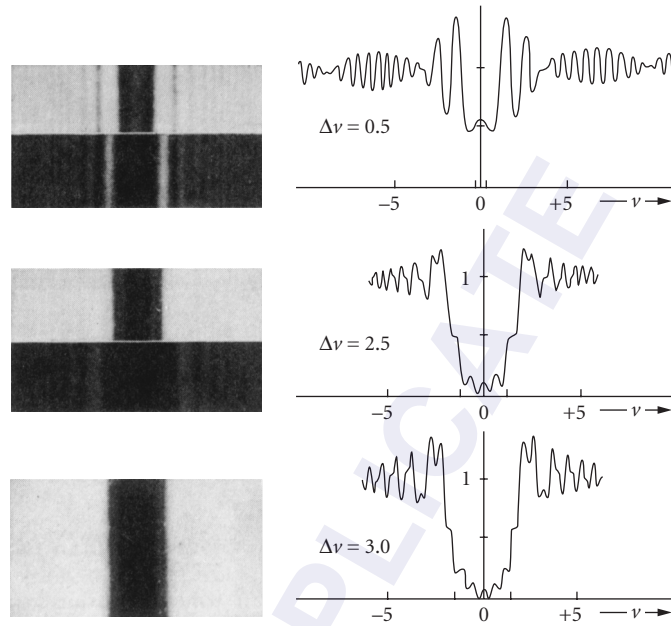


FIGURE 17 Fresnel diffraction patterns produced by narrow wires are shown with the corresponding theoretical curves. (From Jenkins and White.²)

3.6 MATHEMATICAL THEORY OF DIFFRACTION

Kirchhoff showed that the Huygens-Fresnel construction follows from an integral theorem starting from the wave equation. The resulting mathematical expression is called the Fresnel-Kirchhoff diffraction formula.¹ This theory was further refined by Rayleigh and Sommerfeld.^{7,12}

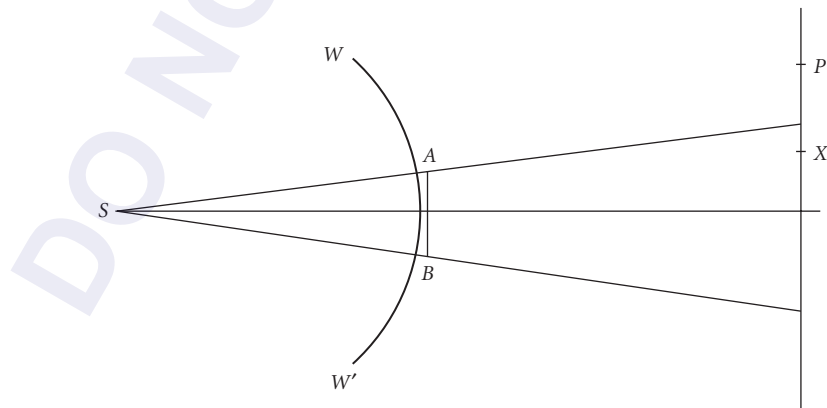


FIGURE 18 Arrangement of the source S , opaque strip AB , and the plane of observation. Point x is in the geometrical shadow region.

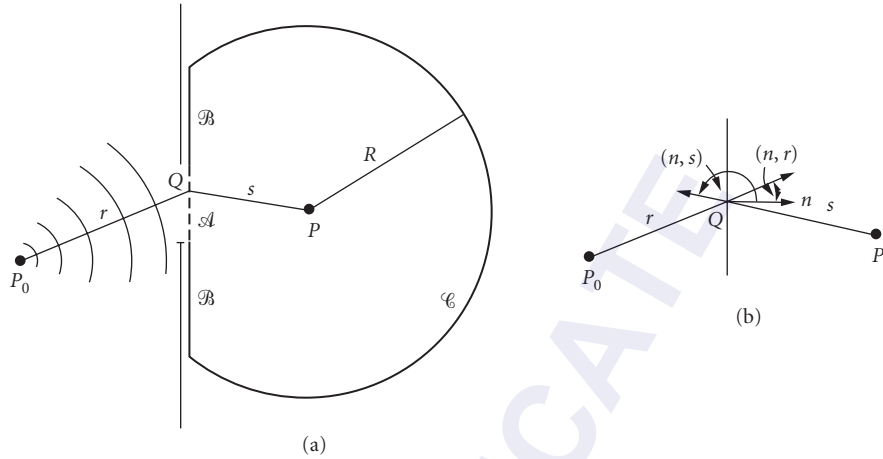


FIGURE 19 The closed surface for the diffraction problem. It is made up of the aperture plane and a large partial sphere centered at the point of observation P . (From Born and Wolf.¹)

It is well known in wave theory that the field values inside a volume enclosed by a bounding surface are determined by the values of the field and/or its normal derivative on this bounding surface. The solution is expressed in terms of the Green function of the problem, as in

$$\psi(P) = \left(\frac{1}{4\pi}\right) \iint_S \left\{ \psi \left(\frac{\partial G}{\partial n} \right) - G \left(\frac{\partial \psi}{\partial n} \right) \right\} dS \quad (40)$$

where G is the Green function of the problem. The integral is over the arbitrary closed surface S . The symbol $\partial/\partial n$ stands for the normal derivative with the normal pointing into the volume enclosed by the surface.¹ A convenient Green function is the expanding spherical wave, $\exp(iks)/s$ from the point of observation P . The closed surface for the diffraction problem is made up of the aperture plane and a large partial sphere centered at the point of observation P , as shown in Fig. 19.

This is the starting point of Kirchhoff theory. It requires specifying the field values and its normal derivative on the bounding surface to obtain the field $\psi(P)$ at P in the volume enclosed by the surface. It is possible to show that the contribution of the surface integral on the partial sphere is zero. Kirchhoff assumed that the field and its normal derivative are zero on the opaque portion of the aperture plane. On the open areas of the aperture plane he assumed the values to be the same as incident (unperturbed) values. If the incident field is an expanding spherical wave $(a/r) \exp(ikr)$, then the field $\psi(P)$ is given by

$$\psi(P) = -\frac{ia}{2\lambda} \iint_A \left[\frac{\exp(ikr)}{r} \right] \left[\frac{\exp(iks)}{s} \right] [\cos(n, r) - \cos(n, s)] dS \quad (41)$$

The area integral is over the open areas A of the aperture. As shown in Fig. 19, (n, s) and (n, r) are the angles made by s and r , respectively, with the normal to the aperture plane. The above equation is referred to as the Fresnel-Kirchhoff diffraction formula. From a strictly mathematical point of view the specification of field and its normal derivative over-specifies the boundary conditions. It is possible to modify the Green function so that only the field or its normal derivative $\partial\psi/\partial n$ needs to be specified. With this modification one obtains

$$\psi(P) = -\left(\frac{ia}{\lambda}\right) \iint_A \left[\frac{\exp(ikr)}{r} \right] \left[\frac{\exp(iks)}{s} \right] \cos(n, s) dS \quad (42)$$

This is referred to as the Rayleigh-Sommerfeld diffraction formula. Other than mathematical consistency, both formulas yield essentially similar results when applied to practical optical situations. They both use the approximate boundary conditions, namely, that the field is undisturbed in the open areas of the aperture and zero on the opaque regions of the aperture plane. The cosine factors in the above formulas play the role of the obliquity factor of the Huygens wave used in Eq. (8), more generally, the field (for a single temporal frequency) at the point of observation $P(x, y, z)$ may be expressed by¹²

$$\psi(x, y, z) = \iint_A \psi(x_s, y_s, 0) \left[\frac{1}{2\pi} \frac{z}{\rho} (1 - ik\rho) \frac{\exp(ik\rho)}{\rho^2} \right] dx_s dy_s \quad (43)$$

where $\psi(x_s, y_s, 0)$ are the values of the field in the aperture A , at $z = 0$. The expression in the square brackets is the normal derivative of the modified Green function. In this expression $\rho = [(x - x_s)^2 + (y - y_s)^2 + z^2]^{1/2}$ is the distance between a point in the aperture and the point of observation P , and the ratio z/ρ is the direction cosine of the difference vector. In the far zone where $k\rho \gg 1$, Eq. (43) reduces to Eq. (42) for the case of spherical wave illumination. Since the expression in the square brackets depends on the coordinate difference, $(x - x_s)$ and $(y - y_s)$, Eq. (43) has the form of a convolution integral. It is well known that the convolution integral has a corresponding product relationship in the Fourier-spatial-frequency domain. The two-dimensional Fourier decomposition of the field is

$$\psi(x, y, z) = \iint \hat{\psi}(p/\lambda, q/\lambda, z) \exp[+i2\pi(px + qy)/\lambda] d(p/\lambda) d(q/\lambda) \quad (44)$$

where p and q are the two-direction cosines. The third-direction cosine m is defined by

$$\begin{aligned} m &= +(1 - p^2 - q^2)^{1/2} & \text{for } p^2 + q^2 \leq 1 \\ m &= +i(p^2 + q^2 - 1)^{1/2} & \text{for } p^2 + q^2 > 1 \end{aligned} \quad (45)$$

A similar decomposition as in Eq. (44) is used for the field in the aperture at $z = 0$, wherein the finite area of the aperture is included in the description of the incident field. With the help of Weyl's plane-wave decomposition of a spherical wave,

$$\frac{\exp(ikr)}{r} = \frac{i}{\lambda} \iint \frac{1}{m} \exp(ikmz) \exp\left[\frac{i2\pi}{\lambda}(px + qy)\right] dp dq \quad (46)$$

the Fourier transform of the expression in square brackets in Eq. (43) can be found. The relationship in the Fourier domain has the form

$$\hat{\psi}(p/\lambda, q/\lambda, z) = \hat{\psi}(p/\lambda, q/\lambda, 0) \exp(ikmz) \quad (47)$$

The inverse Fourier transform yields the disturbance in x, y, z space at point P . At $z = 0$ it reproduces the assumed boundary conditions, a property not shared by the Fresnel-Kirchhoff formula.

A plane-wave decomposition describes a function in (x, y, z) space in terms of the weighted sum of plane waves, each propagating in a direction given by the direction cosines (p, q, m) . Equation (47) may be referred to as the angular spectrum formulation of diffraction. For application of this formulation see Ref. 15.

Fresnel and Fraunhofer Approximations

In practical optical situations, diffraction is mainly studied in the forward direction, that is, for small angles from the direction of propagation of the incident field. Furthermore, the distances involved

are much larger than the wavelength λ , $r \gg \lambda$. In this situation the distance ρ of Eq. (43) may be approximated by the low-order terms in the binomial expansion of the square root

$$\rho = [(x-x_s)^2 + (y-y_s)^2 + z^2]^{1/2} \cong \left(r - \frac{xx_s + yy_s}{r} + \frac{x_s^2 + y_s^2}{2r} \right) \quad (48)$$

where r is the radial distance of the observation point P , $r = \sqrt{x^2 + y^2 + z^2}$. When the terms quadratic in the aperture variables are retained, namely $(x_s^2 + y_s^2)$, we have a description of Fresnel diffraction. Let d be the maximum dimension of the aperture. If the plane of observation is moved to distance $z \gg d^2/\lambda$, the quadratic terms are negligible and Eq. (43) is approximated by

$$\psi(x, y, z) = \left(-\frac{i}{\lambda r} \right) \exp(ikr) \iint_A \psi(x_s, y_s, 0) \exp \left[-\frac{ik(xx_s + yy_s)}{r} \right] dx_s dy_s \quad (49)$$

This is the formula for Fraunhofer diffraction.

Fraunhofer Diffraction

Far enough away from the aperture, $z \gg d^2/\lambda$, Fraunhofer-type diffraction is found. Equation (49) shows that it has the form of a Fourier transform of the light distribution in the aperture. For more general conditions on the distance and angles to obtain Fraunhofer diffraction, see Born and Wolf.¹ Thus, instead of moving the observation plane to the far field, parallel light incident on the aperture can be brought to a focus by a converging lens as in Fig. 20, thus producing a Fraunhofer pattern of the aperture in the focal plane.

In an imaging situation (see Fig. 21), a diverging spherical wave is brought to a focus in the image plane. This is also an example of Fraunhofer diffraction pattern of the light distribution in the aperture A by a converging spherical wave. To realize Fraunhofer diffraction, a similar situation is obtained when a narrow diffracting aperture is held next to the eye focused on a distant point source. The diffraction pattern is observed in the plane of the source.

An optical processing setup is shown in Fig. 22 where collimated or parallel light is incident normally on plane 1. In this arrangement an inverted image of plane 1 is formed in plane 3. The imaging process may be thought of as a Fourier transform (Fraunhofer diffraction) of the light distribution in plane 1 onto plane 2, followed by another Fourier transform of the light distribution in plane 2 onto plane 3.

Recall our earlier discussion in relation to Eqs. (48) and (49). When the quadratic phase factor, $\exp[i\pi(x_s^2 + y_s^2)/\lambda r]$, may be approximated by unity, we are in the domain of Fraunhofer diffraction.

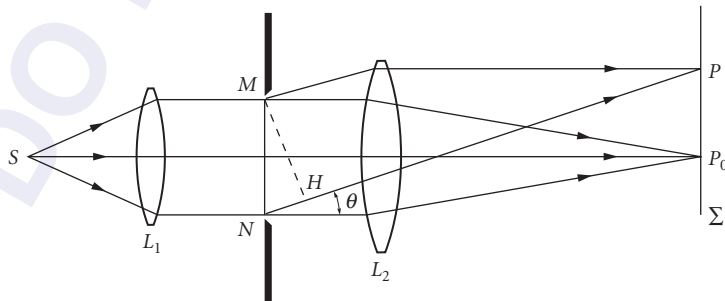


FIGURE 20 Arrangement to observe a Fraunhofer diffraction by a slit aperture. (After Rossi.⁸)

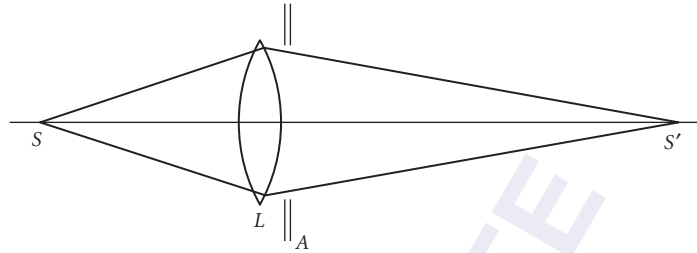


FIGURE 21 Fraunhofer diffraction of an aperture A , with a converging spherical wave. S : point source. L : converging lens. S' : image.

From the point of view of Fresnel zone construction, the far-field condition, $z \gg d^2/\lambda$, means that for these distances z the first Fresnel zone overfills the aperture. The entire aperture contributes to the disturbance at any point in the Fraunhofer pattern. In Fresnel diffraction only relatively small portions of the aperture contribute to any one point in the pattern.

In this context, the term Fresnel number is frequently used. It is defined in terms of the product of two ratios. The radius r of the aperture to the wavelength λ times the radius of the aperture to the distance b measured from the aperture to the plane of observation:

$$\text{Fresnel number} \equiv N = \frac{r}{\lambda} \cdot \frac{r}{b} = \frac{1}{4b} \frac{d^2}{\lambda} \quad (50)$$

Thus, the Fresnel number can also be expressed as the ratio of the far-field distance, d^2/λ , to the distance b from the aperture. With the definition of the Fresnel zones in Sec. 3.4, these ratios indicate that the Fresnel number equals the number of Fresnel zones that may be drawn within the aperture from a point P at a distance b from the aperture.

Thus, well within the Fresnel region, $b \ll d^2/\lambda$, the Fresnel number is large. There are many zones in the aperture. As seen in Figs. 7 and 16, very little light falls within the geometrical shadow region; most of the light is in the confines of the aperture boundary dictated by geometrical optics. In the study of cavity resonators^{16,17} and modes it is found that diffraction losses are small for large Fresnel numbers, $N \gg 1$. In the Fraunhofer region $b > d^2/\lambda$, $N < 1$ where the first Fresnel zone overfills the aperture as pointed out before.

In Figs. 23 and 24 the theoretical plots of Fraunhofer patterns of a rectangular aperture and a circular aperture, respectively, are shown. In the rectangular case the central maximum has equally spaced zeros on either side, while in the circular case the central maximum is surrounded by unequally spaced concentric dark rings. In both cases the central maximum occurs at the geometrical image of the point source that produced parallel light illumination on the aperture.

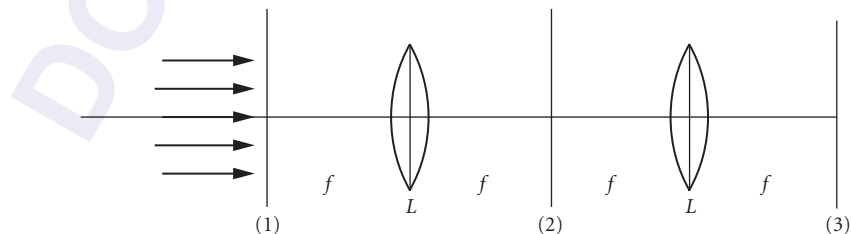


FIGURE 22 Optical processing arrangement. Collimated or parallel light is incident normally on plane 1. The system images plane 1 onto plane 3 with a Fourier transform in plane 2.

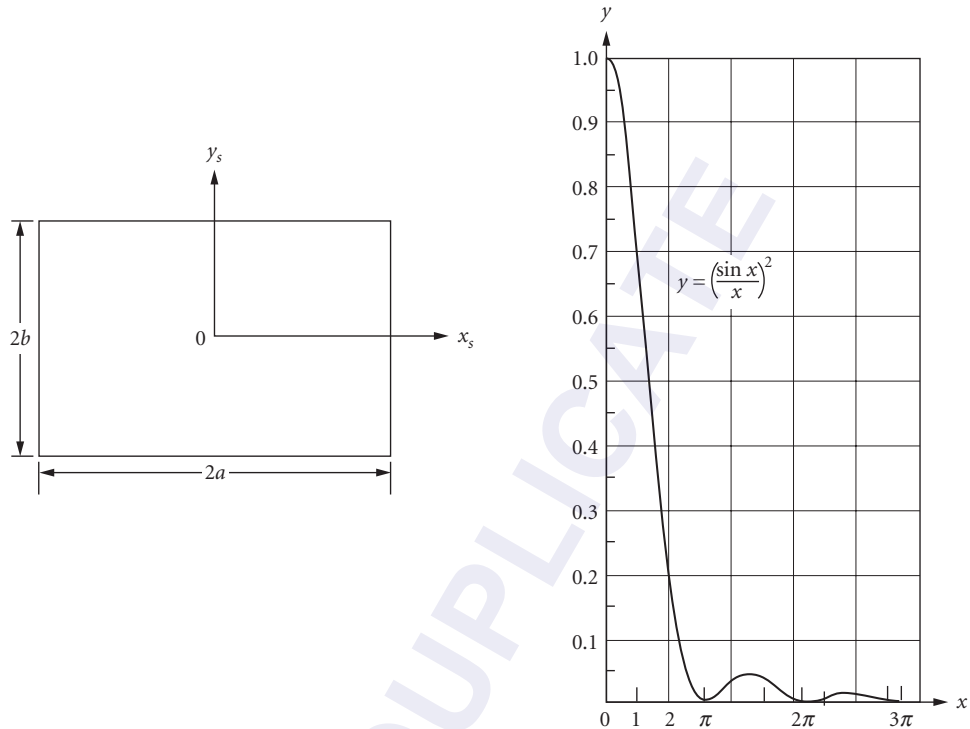


FIGURE 23 Rectangular aperture (in coordinates x_s and y_s) and one section of the diffraction pattern. Normalized irradiance y plotted against a unitless variable x as discussed in the text. (From Born and Wolf.¹)

The unitless variable x shown in the plots is defined as follows. (1) In the case of rectangular aperture, $x = 2\pi a p / \lambda$, where $2a$ is the width of the aperture in the x_s direction. In the other dimension $y = 2\pi b q / \lambda$, and $2b$ is the dimension in the y_s direction. As before, p and q are the direction cosines of the vector joining the center of the aperture to the point of observation. (2) In the case of circular aperture, the unitless radial variable $x = 2\pi a w / \lambda$, where $2a$ is the diameter of the aperture in the x_s, y_s plane and $w = \sqrt{p^2 + q^2}$.

In the far field the size of the diffraction pattern is very large compared to the aperture that produced it. In the focal plane of the lens, $z = f$ and the size of the diffraction pattern is much smaller than the aperture. In both cases the patterns are in a reciprocal width relationship, that is, if the aperture is narrow in the x_s direction compared to y_s , the pattern is broader in the x direction compared to the y . A converging spherical lens illuminated by a plane wave produces in the focal plane a Fraunhofer diffraction pattern of the amplitude and phase of the aperture of a circular lens. When the lens has negligible phase errors, the diffraction pattern has a bright disk in the center surrounded by concentric dark rings. This is called an *Airy disk* and it plays an important role in the Rayleigh criterion of resolving power.

Fraunhofer Diffraction Pattern of a Double Slit

The diffraction pattern of two slits may be observed by using the optical arrangement of Fig. 25. The center-to-center separation of the two slits is h . The off-axis point P is in the direction θ from the axis as shown in the figure. The maxima and minima are determined according to whether the path

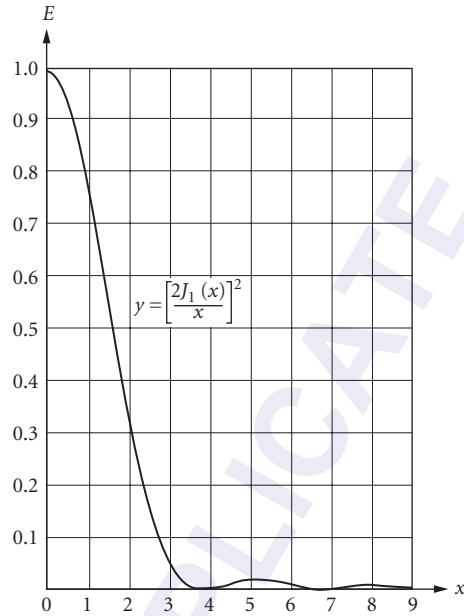


FIGURE 24 A section of the diffraction pattern of a circular aperture. The normalized irradiance E is plotted against a unitless variable x in the plane of observation as discussed in the text. (From Born and Wolf.¹)

difference O_1H is an even or odd multiple of a half-wave. Let E_0 be the irradiance at the center of the single-slit diffraction pattern. The irradiance distribution in the plane of observation is given by

$$E = 4E_0 \left(\frac{\sin \alpha}{\alpha} \right)^2 (\cos \delta)^2 \quad (51)$$

where $\delta = \pi h(\sin \theta)/\lambda$. The irradiance at the center of the double-slit pattern is $4E_0$. The second term, $(\sin \alpha/\alpha)^2$, describes the diffraction pattern of a single slit of width $2a$. Here $\alpha = 2\pi a(\sin \theta)/\lambda$.

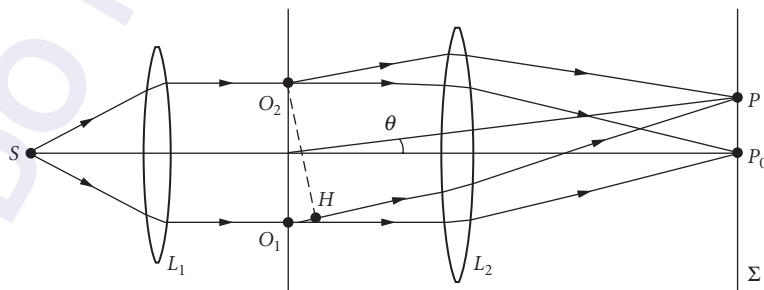


FIGURE 25 Arrangement to observe the Fraunhofer diffraction of an aperture consisting of two slits. S : point source. P : point of observation. O_1, O_2 : two slit apertures with center-to-center separation h . (From Rossi.⁸)

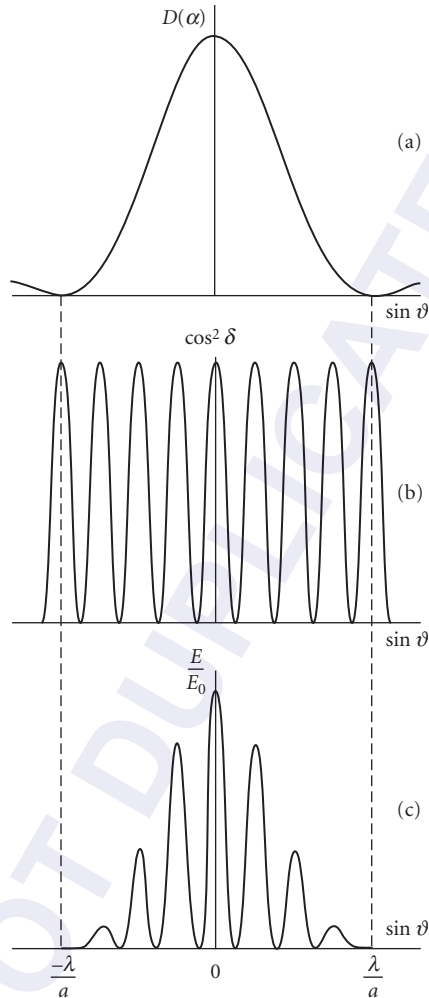


FIGURE 26 (a) Plot of a single-slit diffraction pattern $D(\alpha)$; (b) plot of a two-slit interference pattern; and (c) their product E/E_0 . (From Rossi.⁸)

The term $(\cos \delta)^2$ is the interference pattern of two slits. These two patterns as well as their product are sketched in Fig. 26.

Diffraction Grating

In Fig. 27, an arrangement similar to Fig. 25 permits observation of the Fraunhofer diffraction pattern of a grating, of N parallel and equidistant slits. The center-to-center separation between neighboring slits is h . As in the two-slit case, the Fraunhofer pattern consists of the diffraction due to one slit times the interference pattern of N slits. The irradiance distribution in the plane of observation is given by

$$E = N^2 E_0 \left(\frac{\sin \alpha}{\alpha} \right)^2 \left(\frac{\sin N\gamma}{N \sin \gamma} \right)^2 \quad (52)$$

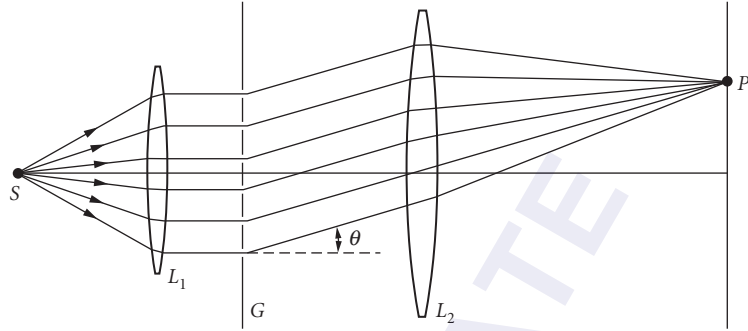


FIGURE 27 Arrangement to observe the Fraunhofer diffraction of an aperture consisting of N slits. S : slit source. P : point of observation. G : grating. (From Rossi.⁸)

where $\gamma = \pi h(\sin \theta)/\lambda$ and $N^2 E_0$ is proportional to the irradiance at the center of the N -slit pattern. The term $(\sin \alpha/\alpha)^2$ is the single-slit pattern as used with Eq. (51). In the case of multiple slits each slit is very narrow; hence, this pattern is very broad, a characteristic of Fraunhofer diffraction. The interference term $(\sin N\gamma/N\sin \gamma)^2$ shows prominent maxima when both the numerator and denominator are simultaneously zero; this happens when $\gamma = \pi h(\sin \theta)/\lambda = m\pi$, where m is an integer. It leads to the grating equation, namely,

$$h \sin \theta = m\lambda \quad (53)$$

There are several, $(N-1)$, subsidiary minima in between the principal maxima. This happens when the numerator is zero but the denominator is not, $\gamma = m\pi/N$. For the case of $N = 10$, these effects are sketched in Fig. 28, which shows the effect of the product of the diffraction and interference terms.

In general, as N increases the subsidiary maxima become more nearly negligible, while the principal maxima become narrower, being proportional to $(1/N)$. The location of the principal maxima other than the zeroth order ($m = 0$) are proportional to the wavelength λ . The diffraction grating thus forms an important spectroscopic tool. Further discussion of gratings is given by Petit¹⁸ and Gaylord and Moharam.¹⁹

3.7 STATIONARY PHASE APPROXIMATION

The diffracted field in the Rayleigh-Sommerfeld diffraction theory is given by

$$\psi(x, y, z) = \iint_A \psi(x_s, y_s, 0) \left[\frac{1}{2\pi} \frac{z}{\rho} (1 - ik\rho) \frac{\exp(ik\rho)}{\rho^2} \right] dx_s dy_s \quad (54)$$

where $\psi(x_s, y_s, 0)$ is the field in aperture A . The diffracted field can also be represented by

$$\psi(x, y, z) = \iint_A \tilde{\psi}(L/\lambda, M/\lambda, 0) \exp[i2\pi(Lx + My + Nz)/\lambda] d(L/\lambda) d(M/\lambda) \quad (55)$$

The phase term in this integral is

$$\phi(L, M) \equiv 2\pi(Lx + My + Nz)/\lambda = \frac{2\pi}{\lambda} \left\{ Lx + My + \sqrt{[1 - (L^2 + M^2)]} z \right\} \quad (56)$$

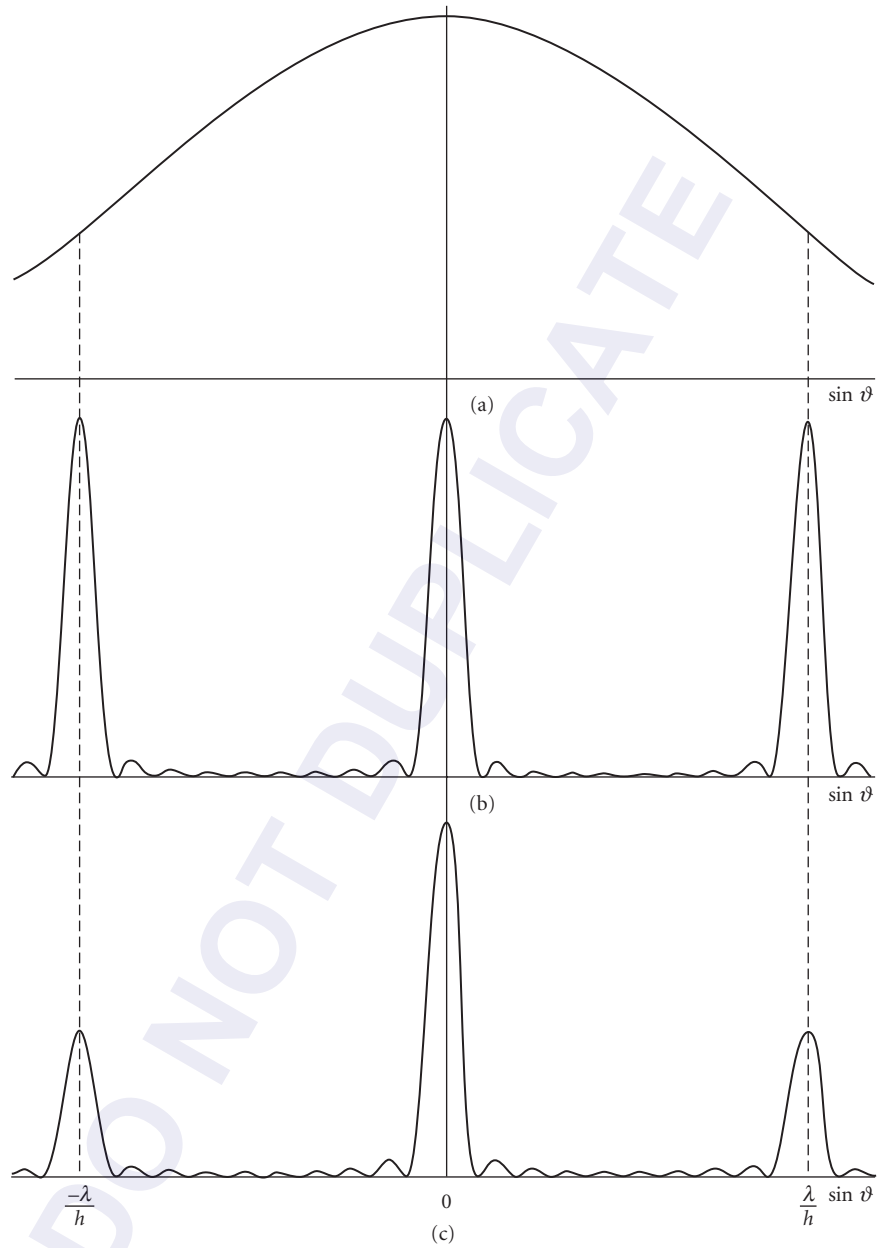


FIGURE 28 (a) Irradiance plot of a single-slit diffraction pattern; (b) partial plot of an $N = 10$ slit interference pattern; and (c) their product. (From Rossi.⁸)

The special values of L and M that make the first derivatives of the phase zero,

$$\frac{\partial\phi}{\partial L}=0=\frac{\partial\phi}{\partial M} \quad (57)$$

are

$$L_0 = \pm \frac{x}{r} \quad \text{and} \quad M_0 = \pm \frac{y}{r} \quad (58)$$

where $r = \sqrt{x^2 + y^2 + z^2}$. The negative sign is omitted for forward propagation, $z > 0$. The phase is approximated by

$$\phi(L, M) \approx \phi(L_0, M_0) + \frac{1}{2} [\alpha(L - L_0)^2 + \beta(M - M_0)^2 + 2\gamma(L - L_0)(M - M_0)]$$

where the higher-order terms are neglected and α , β , and γ are the second derivatives evaluated at $L = L_0$ and $M = M_0$. These constant coefficients are given by

$$\phi(L_0, M_0) = kr \quad (59)$$

$$\alpha \equiv \frac{\partial^2\phi}{\partial L^2} = -kr \frac{x^2 + z^2}{z^2} \quad (60a)$$

$$\beta \equiv \frac{\partial^2\phi}{\partial M^2} = -kr \frac{y^2 + z^2}{z^2} \quad (60b)$$

and

$$\gamma \equiv \frac{\partial^2\phi}{\partial L \partial M} = \frac{\partial^2\phi}{\partial M \partial L} = -kr \frac{xy}{z^2} \quad (60c)$$

The resulting phase function is used in the double integral to obtain the diffracted field, $\psi(x, y, z)$. The reader may also refer to Ref. 1, app. III, Eq. (20). The above procedure yields the stationary phase approximation for the diffracted field given by

$$\psi(x, y, z) \approx -\frac{i}{\lambda} \frac{z}{r} \frac{\exp(ikr)}{r} \psi_{\text{SP}}\left(\frac{x}{r}, \frac{y}{r}\right) \quad (61)$$

where

$$\psi_{\text{SP}}\left(\frac{x}{r}, \frac{y}{r}\right) \equiv \iint_A \psi(x_s, y_s, 0) \exp\left[\frac{-i2\pi}{\lambda} \left(\frac{xx_s + yy_s}{r}\right)\right] dx_s dy_s$$

The diffracted field on a hemisphere is simply the spatial Fourier transform of the field distribution in the aperture as long as the distance r to the observation point satisfies the far-field condition

$$r \gg \frac{N^2 a^2}{\lambda} \quad (62)$$

where $N = z/r = \cos(\theta)$, the third direction cosine, a is the radius of the aperture, and λ is the wavelength of the light incident on the aperture. (θ is measured from the z axis, and as θ increases the far-field condition is weakened.) For observation points not satisfying the far-field condition, the higher-order terms of the stationary phase approximation cannot be neglected. Harvey²⁰ and

Harvey and Shack²¹ have shown that these terms can be considered as aberrations of the spatial Fourier transform of the aperture field on the hemisphere of observation. In the stationary phase approximation, there is no restriction on the direction cosines, L , M , and N . Hence the diffracted field amplitude in Eq. (61) is valid over the entire hemisphere.

3.8 VECTOR DIFFRACTION

The popularity of the Fresnel-Kirchhoff diffraction formula in the scalar case stems from the fact that it is widely applicable and relatively easy to use. In the study of electromagnetic diffraction,^{13,22} a similar formula can be obtained [see Ref. 13, Eq. (9.156)] but it has limited applicability because of the boundary conditions that must be satisfied.

These conditions are the ones related to perfectly conducting screens. They are not adequately approximated at optical frequencies. The study with finite conductivity makes for complicated mathematical procedures. From the point of view of instrumental optics the applicability of the theory then is severely limited.

In the optical literature, periodic structures such as gratings (both shallow and deep compared to the wavelength) have been studied. Boundary conditions are applied to perfectly conducting grating profiles.^{18,19} The equation of the grating dictating the angular positions of the diffraction orders such as Eq. (53) continues to apply; the amount of power found in the different orders is significantly different in the vector theory compared to the scalar theory.

A special case of interest is discussed in detail by Jackson.¹³ Consider a plane wave incident at an angle α on a thin, perfectly conducting screen with a circular hole of radius a in the x - y plane. The polarization vector (E field) of the incident wave lies in the x - z plane, which is taken to be the plane of incidence. The arrangement is shown in Fig. 29 where \mathbf{k}_0 stands for the wave vector of the incident wave and \mathbf{k} is used for the diffracted field.

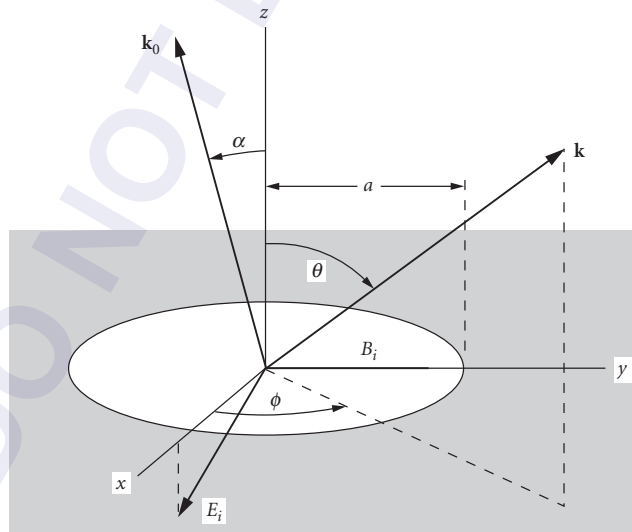


FIGURE 29 Coordinate system and aperture geometry for vector diffraction. α : angle of incidence. The E field is in the xz plane. \mathbf{k}_0 : the wave vector of the incident wave. \mathbf{k} : the diffracted field. (From Jackson.¹³)

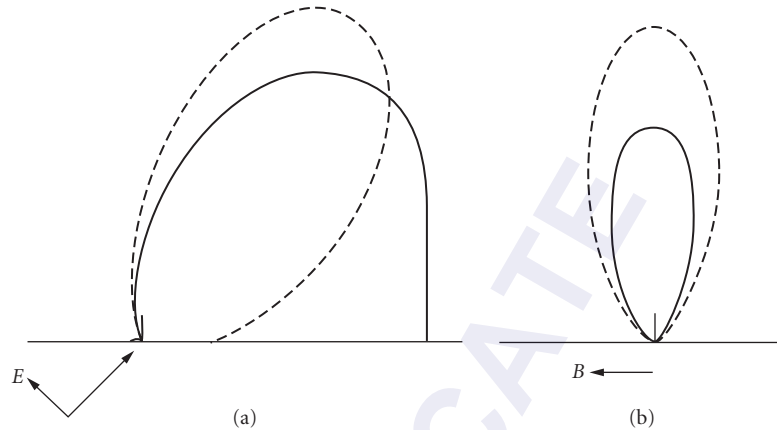


FIGURE 30 Fraunhofer diffraction pattern for a circular opening one wavelength in diameter in a thin-plane conducting screen. The angle of incidence is 45° . (a) Power-per-unit solid angle (radiant intensity) in the plane of incidence and (b) perpendicular to it. The solid (dotted) curve gives the vector (scalar) approximation to it. (From Jackson.¹³)

The vector and scalar approximations are compared in Fig. 30. The angle of incidence is equal to 45° and the aperture is one wavelength in diameter, $ka = \pi$. The angular distribution is shown in Fig. 30 for two cases. Figure 30a shows the distribution of the power per unit solid angle in the plane of incidence which contains the E field and Fig. 30b the distribution for the plane perpendicular to it. Both vector and scalar theories contain the Airy-disk-type distribution; the differences show in the angular distribution.

For normal incidence $\alpha = 0$ and $ka \gg 1$ the polarization dependence is unimportant and the diffraction is confined to very small angles in the forward direction (Airy-disk-type distribution) as we found before in Fig. 24 under Fraunhofer diffraction.

The Vector Huygens-Fresnel Secondary Source

Ideally, the fundamental model of any diffraction theory of light would retain the simplicity of Huygens' scalar secondary source and wavefront construction theory but also account for its vector nature. It has been shown that an electromagnetic wavefront can be modeled as a set of fictitious oscillating electric and magnetic surface charge and current densities existing at all points on the wavefront. The vector Huygens secondary source is a unit composed of two fictitious coincident dipoles; one electric and the other magnetic their magnitudes and orientation dictated by the wavefront boundary conditions. The fields of the vector Huygens secondary source are composed of the linear, vector superposition of the fields of these electric and magnetic dipoles. The electric dipole's axis lies in the plane of the page, is oriented in the vertical direction, and is located at the origin. The magnetic dipole's axis is perpendicular to the plane of the page and is also located at the origin. The vector of (a) the radiated electric field (Fig. 31a) is tangent to the spherical wavefront (represented by the outer circle) and lies in the plane of the page and (b) the radiated magnetic field (Fig. 31b) is tangent to the spherical wavefront (represented by the outer circle) and is perpendicular to the plane of the page. The Poynting vector (Fig. 31c) points radially outward. The magnitude of the irradiance is proportional to the length of the chord from the origin to the irradiance plot along the radius to the point of tangency. The strength of these vectors is proportional to the length of the chord from the origin to the field plot along the radius to the point of tangency.

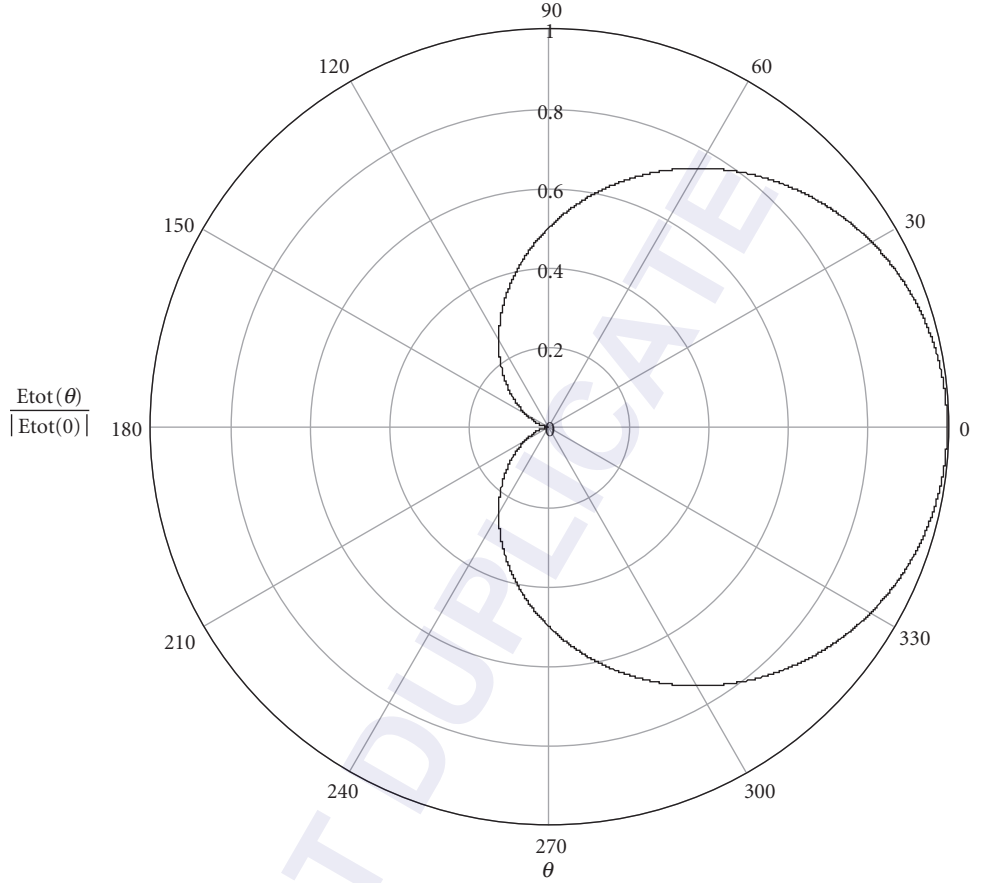


FIGURE 31a This figure is the mapping of the EM dipole's normalized electric field strength on the radiated spherical wavefront in the far zone of the EM dipole as a function of the angle between the direction of observation and \hat{k}_{EM} which lies in the plane of the page and is in the direction of the vector from the origin to 0 degrees.

The plot is rotationally symmetric about \hat{k}_{EM} which lies in the plane of the page and is in the direction of the vector from the origin to 0 degrees.

The diffracted field at any observation point is the summation of the fields radiated from the electromagnetic dipoles in the aperture visible to the observation point as given by

$$\vec{E}(\vec{r}) = \frac{-i}{2\lambda^2} \iiint_{EPW} \left\{ \left[\hat{R}_s \times \vec{E}(\vec{r}_s) \right] \times \hat{R}_s - \sqrt{\frac{\mu_0}{\epsilon_0}} \left[\hat{R}_s \times \vec{H}(\vec{r}_s) \right] \right\} \frac{\exp(ikR_s)}{R_s} d^3\vec{r}_s \quad (63)$$

$$\vec{H}(\vec{r}) = \frac{-i}{2\lambda^2} \iiint_{EPW} \left\{ \left[\hat{R}_s \times \vec{H}(\vec{r}_s) \right] \times \hat{R}_s + \sqrt{\frac{\epsilon_0}{\mu_0}} \left[\hat{R}_s \times \vec{E}(\vec{r}_s) \right] \right\} \frac{\exp(ikR_s)}{R_s} d^3\vec{r}_s \quad (64)$$

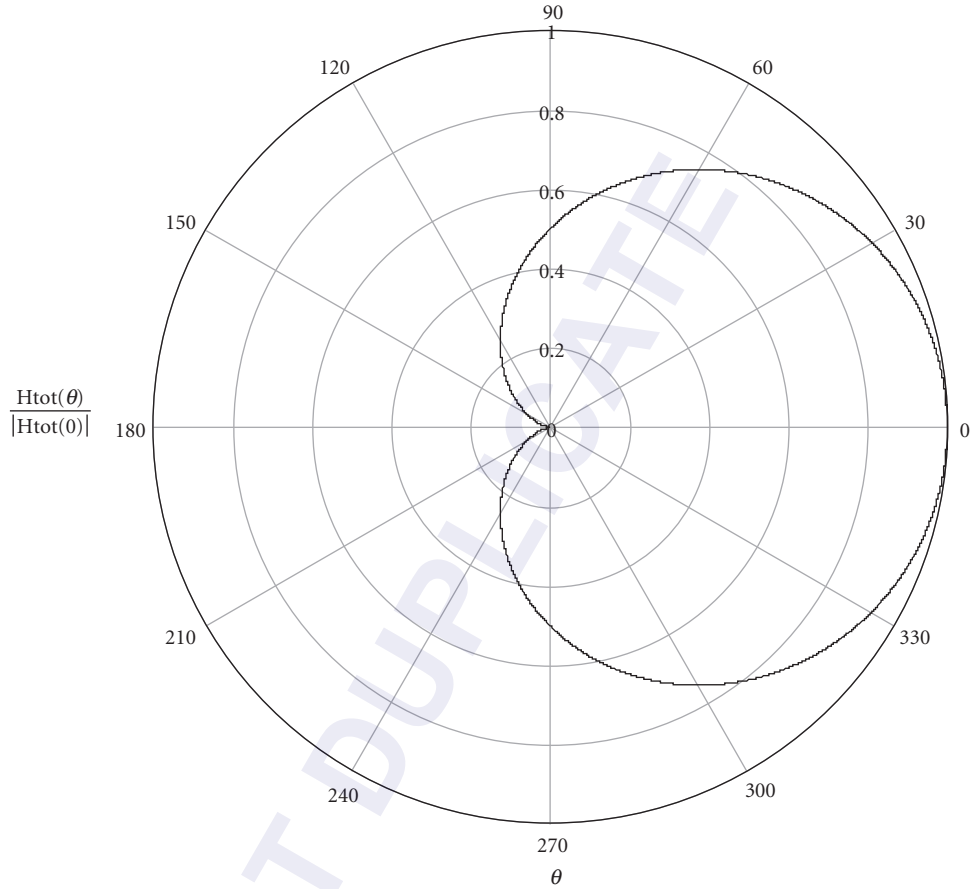


FIGURE 31b This figure is the mapping of the EM dipole's normalized magnetic field strength on the radiated spherical wavefront in the far zone of the EM dipole as a function of the angle between the direction of observation and \hat{k}_{EM} which lies in the plane of the page and is in the direction of the vector from the origin to 0 degrees.

In these expressions, $\vec{E}(\vec{r})$ and $\vec{H}(\vec{r})$ are, respectively, the diffracted electric and magnetic fields at the point of observation \vec{r} . In the integrands, $\vec{E}(\vec{r}_s)$ and $\vec{H}(\vec{r}_s)$ are the fields of the wavefront incident on the aperture at point \vec{r}_s and are related by

$$\vec{H}(\vec{r}_s) = \sqrt{\frac{\epsilon_0}{\mu_0}} \hat{n}(\vec{r}_s) \times \vec{E}(\vec{r}_s) \quad (65)$$

Here $\hat{n}(\vec{r}_s)$ is the normal to the wavefront. We have defined $\vec{R}_s = \vec{r} - \vec{r}_s$ and use R_s for the magnitude and \hat{R}_s for the unit vector in the direction of \vec{R}_s . The letters *EPW* under the volume integral stand for "exposed parts of the primary wavefront." The integration is restricted to the open areas of the aperture.

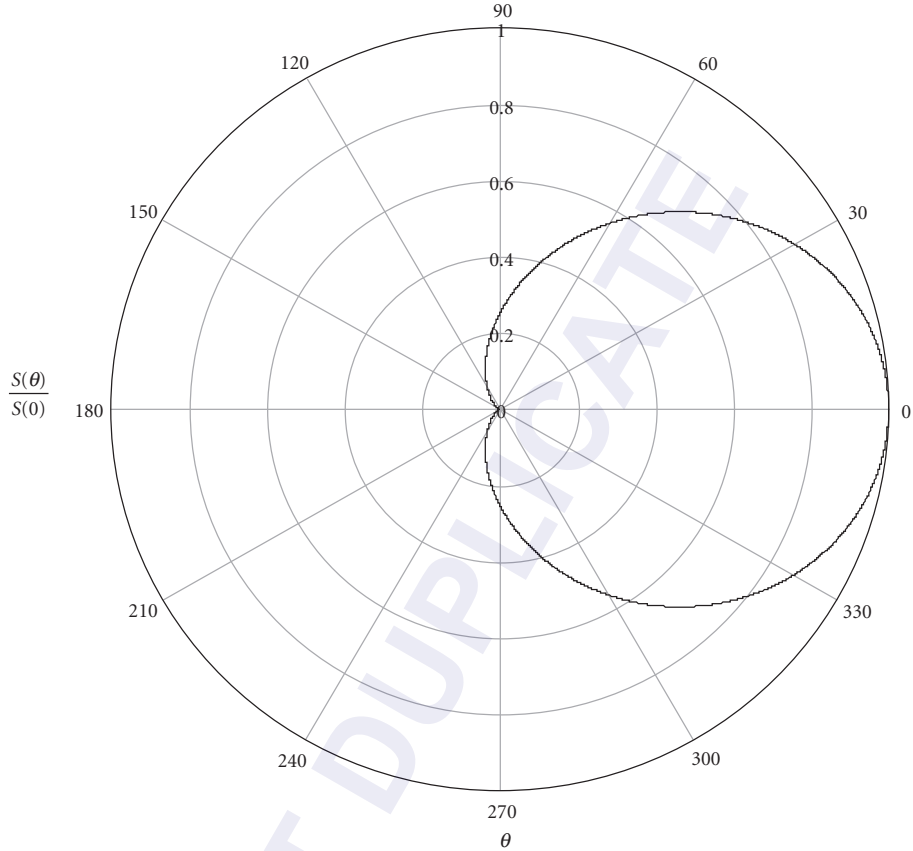


FIGURE 31c This figure is the mapping of the EM dipole's normalized irradiance on the radiated spherical wavefront in the far zone of the EM dipole as a function of the angle between the direction of observation and \hat{k}_{EM} , which lies in the plane of the page and is in the direction of the vector from the origin to 0 degrees.

By use of vector identities, the curly bracket in Eq. (63) may be rewritten in the form

$$\left\{ \left[\hat{R}_s \times \vec{E}(\vec{r}_s) \right] \times \hat{R}_s - \sqrt{\frac{\mu_o}{\epsilon_o}} \left[\hat{R}_s \times \vec{H}(\vec{r}_s) \right] \right\} \quad (66)$$

$$= \vec{E}_{\perp}(\vec{r}_s) + \vec{E}(\vec{r}_s)(\cos\chi) - \hat{n}(\vec{r}_s)(E_{\parallel}(\vec{r}_s))$$

In this expression $\vec{E}_{\perp}(\vec{r}_s)$ is the transverse component of $\vec{E}(\vec{r}_s)$ perpendicular to the direction \hat{R}_s and $E_{\parallel}(\vec{r}_s)$ is the longitudinal component of $\vec{E}(\vec{r}_s)$ parallel to the direction \hat{R}_s . The symbol χ stands for the angle between the unit vectors \hat{R}_s and $\hat{n}(\vec{r}_s)$. In the special case where the angle χ is zero, the curly bracket reduces to $2\vec{E}_{\perp}(\vec{r}_s)$.

For further details and additional references, we refer to McCalmont²³ and Marathay and McCalmont.²⁴ Figure 32 is a sequence of irradiance profiles due to the diffraction of light by a narrow slit based on Eqs. (63) and (64). The sequence describes diffraction from deep within the Fresnel region into the Fraunhofer zone.

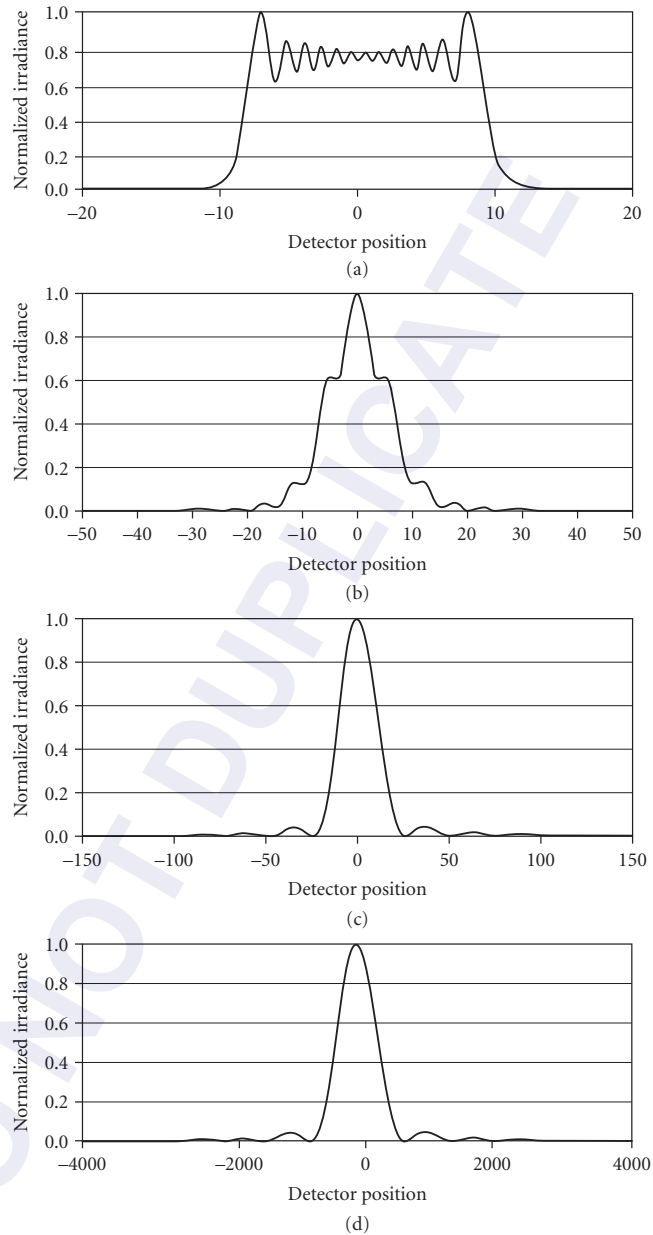


FIGURE 32 Flux density profiles along the horizontal x axis on an observation screen at the following distances: (a) 5λ , (b) 100λ , (c) 500λ , and (d) $15,000\lambda$. These profiles are due to diffraction by a rectangular slit and are based on Eqs. (63) and (64). The slit is of width 20λ . The first zero of the Fraunhofer pattern is at 2.81° from the optical axis. The incident field is a plane wave of unit amplitude, at normal incidence on the plane of the aperture and polarized in the vertical y -direction. The position on the x axis is in terms of wavelengths from the origin of the observation plan.

3.9 ACKNOWLEDGMENTS

The chapter is dedicated by Arvind Marathay to his wife Sunita and family and by John McCalmont to his wife Ingrid and family.

3.10 REFERENCES

Exhaustive and/or the latest listing of references is not the intent but the following source books were used throughout the chapter.

1. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon, Oxford, 1980.
2. F. A. Jenkins and H. E. White, *Fundamentals of Optics*, McGraw-Hill, New York, 1957.
3. R. S. Longhurst, *Geometrical and Physical Optics*, Wiley, New York, 1967.
4. B. K. Mathur and T. P. Pandya, *Geometrical and Physical Optics*, Gopal Printing Press, Kanpur, India, 1972.
5. J. B. Marion and M. A. Heald, *Classical Electromagnetic Radiation*, 2d ed., Academic Press, New York, 1980.
6. A. S. Marathay and J. F. McCalmont, "On the Usual Approximation Used in the Rayleigh-Sommerfeld Diffraction Theory," *J. Opt. Soc. Am.* **21**(4):510–516, 2004.
7. J. W. Goodman, *Introduction to Fourier Optics*, McGraw-Hill, New York, 1968.
8. B. Rossi, *Optics*, Addison-Wesley, Reading, Massachusetts, 1957.
9. C. L. Andrews, *Optics of the Electromagnetic Spectrum*, Prentice-Hall, Englewood Cliffs, New Jersey, 1960.
10. E. Hecht and A. Zajac, *Optics*, Addison-Wesley, Reading, Massachusetts, 1976.
11. M. J. Beran and G. B. Parrent Jr., *Theory of Partial Coherence*, Prentice-Hall, Englewood Cliffs, New Jersey, 1964.
12. A. S. Marathay, *Elements of Optical Coherence Theory*, Wiley, New York, 1985.
13. J. D. Jackson, *Classical Electrodynamics*, Wiley, New York, 1962.
14. A. E. Sommerfeld, *Electrodynamics*, Academic Press, New York, 1964 (Translation of *Über Theoretische Volesung Physik*, Vol. 3, *Electrodynamik*).
15. W. H. Carter, "Wave Theory of a Simple Lens," *Optica Acta*, **20**(10):805–826, 1973.
16. K. D. Möller, *Optics*, University Science Books, Mill Valley, California, 1988.
17. H. Kogelnik and T. Li, "Laser Beams and Resonators," *Appl. Opt.* **5**:1550–1565, 1966.
18. R. Petit, *Electromagnetic Theory of Gratings*, Springer-Verlag, New York, 1980.
19. T. K. Gaylord and M. G. Moharam, "Analysis and Applications of Optical Diffraction by Gratings," *Proc. IEEE* **73**:894–937, 1985.
20. J. E. Harvey, "Fourier Treatment of Near-Field Scalar Diffraction Theory," *Am. J. Phys.* **47**(11), Nov. 1979.
21. J. E. Harvey and R. V. Shack, "Aberrations of Diffracted Wave Fields," *Appl. Opt.* **17**:3003–3009, 1978.
22. J. A. Stratton, *Electromagnetic Theory*, McGraw-Hill, New York, 1941.
23. J. F. McCalmont, *A Vector Huygens-Fresnel Model of the Diffraction of Electromagnetic Waves*, Ph.D. dissertation, University of Arizona, 1999.
24. A. S. Marathay and J. F. McCalmont, "Vector Diffraction Theory for Electromagnetic Waves," *J. Opt. Soc. Am.* **18**:2585–2593, 2001.

4

TRANSFER FUNCTION TECHNIQUES

Glenn D. Boreman

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

4.1 GLOSSARY

B	spot full width
CTF	contrast transfer function (square wave response)
$e(x)$	edge response
FN	focal ratio
$F(\xi, \eta)$	Fourier transform of $f(x, y)$
$f(x, y)$	object function
$G(\xi, \eta)$	Fourier transform of $g(x, y)$
$g(x, y)$	image function
$H(\xi, \eta)$	Fourier transform of $h(x, y)$
$h(x, y)$	impulse response
$\ell(x)$	line response
$S(\xi, \eta)$	power spectrum
W	detector dimension
$\delta(x)$	delta function
$\theta(\xi, \eta)$	phase transfer function
**	two-dimensional convolution

4.2 INTRODUCTION

Transfer functions are a powerful tool for analyzing optical and electro-optical systems. The interpretation of objects and images in the frequency domain makes available the whole range of linear-systems analysis techniques. This approach can facilitate insight, particularly in the treatment of complex optical problems. For example, when several optical subsystems are combined, the overall transfer function is the multiplication of the individual transfer functions. The corresponding analysis, without the use of transfer functions, requires convolution of the corresponding impulse responses.

4.3 DEFINITIONS

The image quality of an optical or electro-optical system can be characterized by either the system's impulse response or its Fourier transform, the transfer function. The impulse response $h(x, y)$ is the two-dimensional image formed in response to a delta-function object. Because of the limitations imposed by diffraction and aberrations, the image quality produced depends on the following: the wavelength distribution of the source; the F-number (FN) at which the system operates; the field angle at which the point source is located; and the choice of focus position.

A continuous object $f(x, y)$ can be decomposed, using the sifting property of delta functions, into a set of point sources, each with a strength proportional to the brightness of the object at that location. The final image $g(x, y)$ obtained is the superposition of the individually weighted impulse responses. This result is equivalent to the convolution of the object with the impulse response:

$$f(x, y) ** h(x, y) = g(x, y) \quad (1)$$

where the double asterisk denotes a two-dimensional convolution.

The validity of Eq. (1) requires shift invariance and linearity. Shift invariance is necessary for the definition of a single impulse response and linearity is necessary for the superposition of impulse responses. These assumptions are often violated in practice, but the convenience of a transfer-function analysis dictates that we preserve this approach if possible. While most optical systems are linear, electro-optical systems that include a receiver (such as photographic film, detector arrays, and xerographic media) are often nonlinear. A different impulse response (and hence transfer function) is obtained for inputs of different strengths. In optical systems with aberrations that depend on field angle, separate impulse responses are defined for different regions of the image plane.

Although $h(x, y)$ is a complete specification of image quality (given a set of optical parameters), additional insight is gained by use of the transfer function. A transfer-function analysis considers the imaging of sinusoidal objects, rather than point objects. It is more convenient than an impulse-response analysis because the combined effect of two or more subsystems can be calculated by a point-by-point multiplication of the transfer functions, rather than by convolving the individual impulse responses. Using the convolution theorem of Fourier transforms, we can rewrite the convolution of Eq. (1) as a multiplication of the corresponding spectra:

$$F(\xi, \eta) \times H(\xi, \eta) = G(\xi, \eta) \quad (2)$$

where the uppercase variables denote the Fourier transforms of the corresponding lowercase variables: $F(\xi, \eta)$ is the object spectrum; $G(\xi, \eta)$ is the image spectrum; $H(\xi, \eta)$ is the spectrum of the impulse response. As a transfer function, $H(\xi, \eta)$ multiplies the object spectrum to yield the image spectrum. The variables ξ and η are spatial frequencies in the x and y directions. Spatial frequency is the reciprocal of the crest-to-crest distance of a sinusoidal waveform used as a basis function in the Fourier analysis of an object or image. In two dimensions, a sinusoid of arbitrary orientation has a spatial period along both the x and y axes. The reciprocals of these spatial periods are the spatial frequencies ξ and η . Typical units of spatial frequency are cycles/millimeter when describing an image, and cycles/milliradian when describing an object at a large distance. For an object located at infinity, these two representations are related through the focal length of the image-forming optical system:

$$\xi_{\text{angular}} [\text{cycles/mrad}] = 0.001 \times \xi [\text{cycles/mm}] \times f [\text{mm}] \quad (3)$$

The function $H(\xi, \eta)$ in Eq. (2) is usually normalized to have unit value at zero frequency. This yields a transfer function relative to the response at low frequency, and ignores frequency-independent attenuations, such as losses caused by Fresnel reflection or by obscurations. This normalization is appropriate for most optical systems, because¹ the transfer function of an incoherent optical system is proportional to the two-dimensional autocorrelation of the exit pupil, which is maximum at zero frequency. For more general imaging systems (for example, the human eye, photographic film, and electronic imaging systems), the transfer function is not necessarily maximum at the origin, and may be more useful in an unnormalized form.

With the above normalization, $H(\xi, \eta)$ is called the optical transfer function (OTF). In general, OTF is a complex function, having both a magnitude and a phase portion:

$$\text{OTF}(\xi, \eta) = H(\xi, \eta) = |H(\xi, \eta)| \exp\{-j\theta(\xi, \eta)\} \quad (4)$$

The magnitude of the OTF, $|H(\xi, \eta)|$, is referred to as the modulation transfer function (MTF), while the phase portion of the OTF, $\theta(\xi, \eta)$, is referred to as the phase transfer function (PTF).

MTF is the magnitude response of the imaging system to sinusoids of different spatial frequencies. This response is described in terms of the modulation depth, a measure of visibility or contrast:

$$M = \frac{A_{\max} - A_{\min}}{A_{\max} + A_{\min}} \quad (5)$$

where A refers to a value of the waveform (typically W/cm^2 vs position) that describes the object or image. These quantities are nonnegative, so the sinusoids always have a dc bias. Modulation depth is thus a number between 0 and 1. The effect of the finite-size impulse response is that the modulation depth in the image is less than that in the object. This attenuation is usually more severe at high frequencies. MTF is the ratio of image modulation to object modulation, as a function of spatial frequency:

$$\text{MTF}(\xi, \eta) = \frac{M_{\text{image}}(\xi, \eta)}{M_{\text{object}}(\xi, \eta)} \quad (6)$$

PTF describes the relative phases with which the various sinusoidal components recombine in the image. A linear phase such as $\text{PTF} = x_0 \xi$ corresponds to a shift of the image by an amount x_0 , each frequency component being shifted the amount required to reproduce the original waveform at the displaced location. For impulse responses that are symmetric about the ideal image point, the PTF exhibits phase reversals, with a value of either 0 or π radians as a function of spatial frequency. A general impulse response that is real but not even yields a PTF that is a nonlinear function of frequency, resulting in image degradation. Linearity of PTF is a sensitive test for aberrations (such as coma) which produces asymmetric impulse responses, and is often a design criterion.

4.4 MTF CALCULATIONS

OTF can be calculated from wave-optics considerations. For an incoherent optical system, the OTF is proportional to the two-dimensional autocorrelation of the exit pupil. This calculation can account for any phase factors across the pupil, such as those arising from aberrations or defocus. A change of variables is required for the identification of an autocorrelation (a function of position in the pupil) as a transfer function (a function of image-plane spatial frequency). The change of variables is

$$\xi = \frac{x}{\lambda d_i} \quad (7)$$

where x is the autocorrelation shift distance in the pupil, λ is the wavelength, and d_i is the distance from the exit pupil to the image. A system with an exit pupil of full width D has an image-space cut-off frequency consistent with Eq. (7):

$$\xi_{\text{cutoff}} = \frac{1}{(\lambda \text{ FN})} \quad (8)$$

where FN equals (focal length)/ D for a system with the object at infinity, and d_i/D for a system operating at finite conjugates.

A diffraction-limited system has a purely real OTF. Diffraction-limited MTFs represent the best performance that a system can achieve, for a given FN and λ , and accurately describe systems with

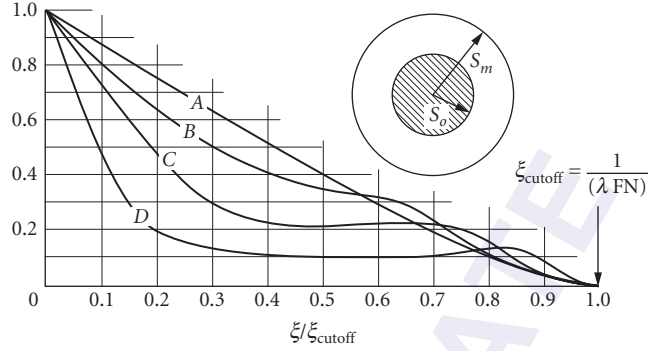


FIGURE 1 (A) Diffraction-limited MTF for system with circular pupil (no obscuration: $S_o/S_m = 0$). (B) through (D) are diffraction-limited MTF for a system with an annular pupil: (B) $S_o/S_m = 0.25$; (C) $S_o/S_m = 0.5$; (D) $S_o/S_m = 0.75$. (Adapted from Ref. 4, p. 322.)

negligible aberrations, whose impulse-response size is dominated by diffraction effects. A diffraction-limited system with a square exit pupil of dimensions $D \times D$ has a linear MTF along ξ or η :

$$\text{MTF}\left(\frac{\xi}{\xi_{\text{cutoff}}}\right) = 1 - \frac{\xi}{\xi_{\text{cutoff}}} \quad (9)$$

For a system with a circular exit pupil of diameter D , the MTF is circularly symmetric, with ξ profile:²

$$\text{MTF}\left(\frac{\xi}{\xi_{\text{cutoff}}}\right) = \frac{2}{\pi} \left\{ \cos^{-1}\left(\frac{\xi}{\xi_{\text{cutoff}}}\right) - \frac{\xi}{\xi_{\text{cutoff}}} \left[1 - \left(\frac{\xi}{\xi_{\text{cutoff}}}\right)^2 \right]^{1/2} \right\} \quad \text{if } \xi \leq \xi_{\text{cutoff}} \quad (10)$$

$$= 0 \quad \text{if } \xi > \xi_{\text{cutoff}}$$

Equation (10) is plotted in Fig. 1, along with MTF curves obtained for annular pupils, which arise in obscured systems such as Cassegrain telescopes. The plots are functions of the obscuration ratio, and the emphasis at high frequencies has been obtained by an overall decrease in flux reaching the image, proportional to the obscured area. If the curves in Fig. 1 were plotted without normalization to 1 at $\xi = 0$, they would all be contained under the envelope of the unobscured diffraction-limited curve.

A system exhibiting effects of both diffraction and aberrations has an MTF curve bounded by the diffraction-limited MTF curve as the upper envelope. Aberrations broaden the impulse response, resulting in a narrower and lower MTF, with less integrated area.

The effect of defocus on the MTF is shown in Fig. 2. The MTF curves resulting from third-order spherical aberration are shown in Fig 3. MTF results for specific cases of other aberrations are contained in Ref. 3.

A geometrical-aberration OTF can be calculated from ray-trace data, without regard for diffraction effects. Optical-design computer programs typically yield a diagram of ray-intersection density in the image plane, a geometrical-optics spot diagram. A geometrical-aberration OTF is calculated by Fourier transforming the spot-density distribution. The OTF thus obtained is accurate if the impulse-response size is dominated by aberration effects. A one-dimensional uniform blur spot of full width B has the following OTF in the ξ direction:

$$\text{OTF}(\xi) = \frac{\sin(\pi\xi B)}{\pi\xi B} \quad (11)$$

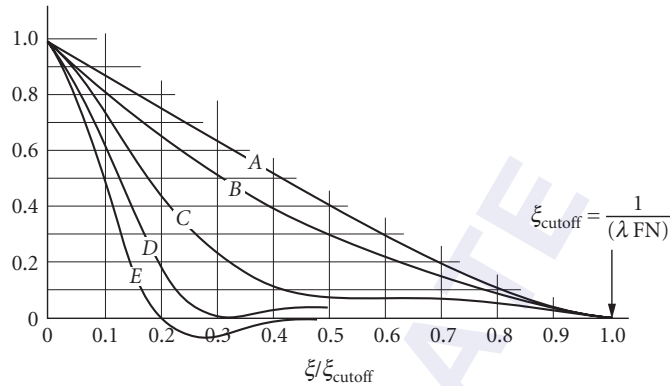


FIGURE 2 Diffraction MTF for a defocused system: (A) in focus, $OPD = 0.0$; (B) defocus $= \lambda/2N \sin^2 u$, $OPD = \lambda/4$; (C) defocus $= \lambda/N \sin^2 u$, $OPD = \lambda/2$; (D) defocus $= 3\lambda/2N \sin^2 u$, $OPD = 3\lambda/4$; and (E) defocus $= 2\lambda/N \sin^2 u$, $OPD = \lambda$. (Adapted from Ref. 4, p. 320.)

which has a zero at $\xi = 1/B$, and also exhibits the phase reversals mentioned above. When an MTF has been calculated from ray trace data, an approximation to the total system MTF may be made⁴ by multiplying the diffraction-limited MTF of the proper FN and λ with the ray-trace data MTF. This is equivalent to a convolution of the spot profiles from diffraction and geometrical aberrations.

In electronic imaging systems, an electronics subsystem performs signal-handling and signal-processing functions. The performance characterization of electronic networks by transfer-function techniques is well established. The usual independent variable for these time-domain transfer functions is the temporal frequency f (Hz). To interpret the electronics transfer function in the same units as the image-plane spatial frequency (cycles/mm), the temporal frequencies are divided by the scan velocity (mm/s). For a scanning system, this is the velocity of the instantaneous field of view, referred to as image coordinates. For a staring system, an effective scan velocity is the

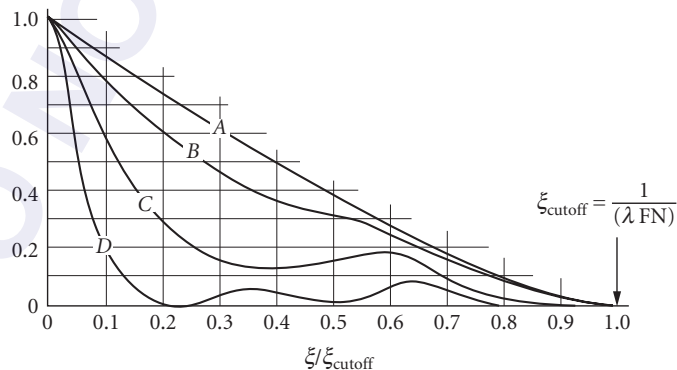


FIGURE 3 Diffraction MTF for system with third-order spherical aberration (image plane midway between marginal and paraxial foci): (A) $LA_m = 0.0$, $OPD = 0$; (B) $LA_m = 4\lambda/N \sin^2 u$, $OPD = \lambda/4$; (C) $LA_m = 8\lambda/N \sin^2 u$, $OPD = \lambda/2$; and (D) $LA_m = 16\lambda/N \sin^2 u$, $OPD = \lambda$. (Adapted from Ref. 4, p. 322.)

horizontal dimension of the image plane divided by the video line time. With this change of variables from temporal frequencies to spatial frequencies, the electronics can be analyzed as simply an additional subsystem, with its own transfer function that will multiply the transfer functions of the other subsystems. It should be noted that an electronics transfer function is not bounded by a pupil autocorrelation the way an optical transfer function is. Thus, it need not be maximum at the origin, and can amplify certain frequencies and have sharp cutoffs at others. Thus, the usual normalization of MTF may not be appropriate for analysis of the electronics subsystems, or for the entire imaging system including the electronics.

An unavoidable impact of the electronics subsystem is the contribution of noise to the image. This limits the amount of electronic amplification that is useful in recovering modulation depth lost in other subsystems. A useful figure of merit, which has been validated to correlate with image visibility,⁵ is the area between two curves: the MTF and the noise power spectrum. To facilitate comparison on the same graph, the noise power spectrum is expressed in modulation depth units, and is interpreted as a noise-equivalent modulation depth (the modulation needed for unit signal-to-noise ratio) as a function of spatial frequency.

The detector photosensitive area has finite size, rather than being a true point. It thus performs some spatial averaging⁶ on any irradiance distribution that falls on it. Large detectors exhibit more attenuation of high spatial frequencies than do small detectors. For a detector of dimension W in the x direction, the MTF is

$$\text{MTF}(\xi) = \left| \frac{\sin(\pi\xi W)}{\pi\xi W} \right| \quad (12)$$

which has a zero at $\xi = 1/W$. This MTF component applies to any system with detectors, and will multiply the MTFs of other subsystems.

In electronic imaging systems, the image is typically sampled in both directions. The distance between samples will determine the image-plane spatial frequency at which aliasing artifacts will occur. Care must be taken in the calculation of MTF, because different impulse responses are possible depending on the location of the impulse response with respect to the sampling positions. This violates the assumption of shift-invariance needed for a transfer-function analysis.⁷ One approach for defining a generalized MTF is to average over all possible positions of the impulse response with respect to the sampling lattice [Eq. (4) in Ref. 8]. Research is still underway on the specification of MTF for sampled-image systems.

4.5 MTF MEASUREMENTS

In any situation where the measurement of MTF involves the detection of the image-plane flux, one component of the measurement-system MTF is caused by the finite aperture of the detector, which can be accounted for in the calibration of the instrument by dividing out the detector MTF seen in Eq. (12).

When OTF is measured with a point-source object, the image formed by the system under test is the impulse response. The two-dimensional impulse response can be Fourier transformed in two dimensions to yield OTF (ξ, η). If an illuminated pinhole is used, it should be as small as possible. However, flux-detection considerations dictate a finite size for any source. The object is small enough not to affect the measurement if its angular subtense is much smaller than the angular subtense of the impulse response, when both are viewed from the aperture stop of the system. For sources of larger extent, a Fourier analysis can be made of the object, and an OTF can be calculated using Eq. (2), over the range of spatial frequencies provided by the source.

If higher flux levels are needed to maintain signal-to-noise ratio, a line response can be measured. The system under test is presented with an illuminated line source, which acts as a delta function in one direction and a constant in the other: $\delta(x)1(y)$. The system forms an image, the

line response $\ell(x)$, which is a summation of vertically displaced impulse responses. In general $\ell(x) \neq h(x, 0)^2$. The line response only yields information about one profile of OTF (ξ, η) . The one-dimensional Fourier transform of the line response produces the corresponding profile of the two-dimensional OTF: $\mathcal{F}\{\ell(x)\} = \text{OTF}(\xi, 0)$. To obtain other profiles of the OTF, the line source is reoriented. Line response data are also available from the response of the system to a point source, using a receiver that integrates the impulse response along one direction: a detector that is long in one dimension and is scanned perpendicularly, or a long slit that is scanned in front of a large-area detector.

Another measurement of OTF uses the edge response $e(x)$, which is the response of the system to an illuminated knife edge. Each line in the open part of the aperture produces a displaced line response, so $e(x)$ is a cumulative distribution, related to the line response as follows: $d/dx\{e(x)\} = \ell(x)$, which Fourier transforms to the ξ profile of the OTF. The derivative operation increases the effect of noise. Any digital filter used for data smoothing has its own impulse response, and hence its own OTF contribution. The edge response can also be measured by using a scanning knife edge in front of a detector in the image plane, with a point-source or a line-source object.

An MTF calculated from a measured profile is the product of a diffraction MTF and a geometrical-aberration MTF. When combining the separately-measured MTFs of several optical subsystems, care should be taken to ensure that the diffraction MTF (determined by the aperture stop of the combined system) contributes only once to the calculation. The geometrical-aberration MTFs for each subsystem will cascade if each subsystem operates independently on an irradiance basis, with no partial coherence effects.⁹ The major exception to this condition occurs when two subsystems are designed to correct for each other's aberrations, and the MTF of the combined system is better than the individual MTFs would indicate.

MTF can also be obtained by the system's response to a sine-wave target, where the image modulation depth is measured as a function of spatial frequency. PTF can also be measured from the position of the waveform maxima as a function of frequency. Sine-wave targets are available as photographic prints or transparencies, which are suitable for testing visible-wavelength systems. Careful control in their manufacture is exercised¹⁰ to avoid harmonic distortions, including a limitation to relatively small modulation depths. Sine-wave targets are difficult to fabricate for testing infrared systems, and require the use of half-tone techniques.¹¹

A more convenient target to manufacture is the three- or four-bar target of equal line and space width, with a binary transmission or reflection characteristic. These are widely used for testing both visible-wavelength and infrared systems. The square-wave response is called the contrast transfer function (CTF) and is not equivalent to the sine-wave response for which MTF is defined. CTF is a function of the fundamental spatial frequency ξ_f (inverse of the bar-to-bar spacing) and is measured on the peak-to-valley variation of image irradiance. For any particular fundamental frequency, the measured response to bar targets will be higher than that measured for sinewaves of the same frequency, because additional harmonic components contribute to the modulation. For a square-wave pattern of infinite extent, an analytical relationship exists¹² between $\text{CTF}(\xi_f)$ and $\text{MTF}(\xi)$. Each Fourier component of the square wave has a known transfer factor given by $\text{MTF}(\xi)$, and the modulation depth as a function of ξ_f of the resultant waveform can be calculated by Eq. (5). This process yields the following series:

$$\text{CTF}(\xi_f) = \frac{4}{\pi} \left\{ \text{MTF}(\xi_f) - \frac{1}{3} \text{MTF}(3\xi_f) + \frac{1}{5} \text{MTF}(5\xi_f) - \frac{1}{7} \text{MTF}(7\xi_f) + \frac{1}{9} \text{MTF}(9\xi_f) - \dots \right\} \quad (13)$$

CTFs for the practical cases of three- and four-bar targets are slightly higher than the CTF curve for an infinite square wave. Figure 4¹³ compares the MTF for a diffraction-limited circular-aperture system with CTFs obtained for infinite, three- and four-bar targets. Because of the broad spectral features associated with bar patterns of limited extent, a finite level of modulation is present in the image, even when the fundamental frequency of the bar pattern equals the cutoff frequency of the system MTF.¹⁴ The inverse process of expressing the MTF in terms of CTFs is more difficult analytically,

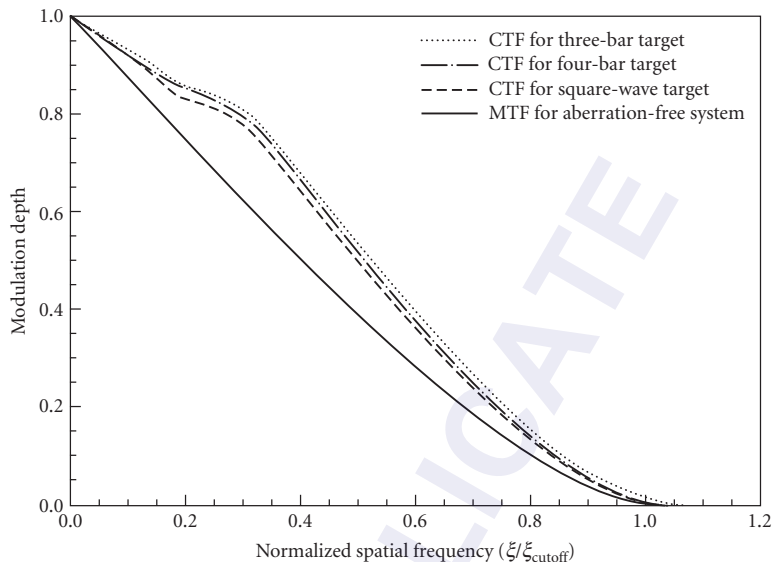


FIGURE 4 Comparison of MTF to CTFs obtained with infinite square wave, four-bar, and three-bar targets for a diffraction-limited system with circular pupil.

since square waves are not an orthogonal basis set for the expansion of sinusoids. A term-by-term series subtraction¹² yields the following:

$$\text{MTF}(\xi_f) = \frac{\pi}{4} \left\{ \text{CTF}(\xi_f) + \frac{1}{3} \text{CTF}(3\xi_f) - \frac{1}{5} \text{CTF}(5\xi_f) + \frac{1}{7} \text{CTF}(7\xi_f) + \frac{1}{11} \text{CTF}(11\xi_f) - \dots \right\} \quad (14)$$

Narrowband electronic filtering can be used to isolate the fundamental spatial-frequency component for systems where the image data are available as a time-domain waveform. These systems do not require the correction of Eq. (14), because the filter converts bar-target data to sinewave data.

The MTF can also be measured by the response of the system to a random object. Laser speckle provides a convenient means to generate a random object distribution of known spatial-frequency content. The MTF relates the input and output spatial-frequency power spectra of the irradiance waveforms:

$$S_{\text{output}}(\xi, \eta) = |\text{MTF}(\xi, \eta)|^2 \times S_{\text{input}}(\xi, \eta) \quad (15)$$

This method is useful in the measurement of an average MTF for sampled-image systems,¹⁵ since the speckle pattern has a random position with respect to the sampling sites.

A number of interferometric methods have been developed for measuring MTF.¹⁶ An interferogram of the wavefront exiting the system is reduced to find the phase map. The distribution of amplitude and phase across the exit pupil contains the information necessary for calculation of OTF by pupil autocorrelation.

4.6 REFERENCES

1. J. W. Goodman, *Introduction to Fourier Optics*, McGraw-Hill, New York, 1968, pp. 116–120.
2. J. D. Gaskill, *Linear Systems, Fourier Transforms, and Optics*, Wiley, New York, 1978, pp. 305–307.
3. V. N. Mahajan, *Optical Imaging and Aberrations*, SPIE Press, Bellingham, WA, 1998, Part II, chap. 2.
4. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1966, p. 323.

5. H. L. Snyder, "Image Quality and Observer Performance," in L. M. Biberman (ed.), *Perception of Displayed Information*, Plenum, New York, 1973.
6. G. D. Boreman and A. E. Plogstedt, "Spatial Filtering by a Line-Scanned Nonrectangular Detector-Application to SPRITE Readout MTF," *Appl. Opt.* **28**:1165–1168, 1989.
7. W. Wittenstein, J. C. Fontanella, A. R. Newbery, and J. Baars, "The Definition of the OTF and the Measurement of Aliasing for Sampled-Imaging Systems," *Optica Acta* **29**(1):41–50, 1982.
8. S. K. Park, R. Schowengerdt, and M. Kaczynski, "Modulation-Transfer-Function Analysis for Sampled Image Systems," *Appl. Opt.* **23**:2572, 1984.
9. J. B. DeVelis and G. B. Parrent, "Transfer Function for Cascaded Optical Systems," *J. Opt. Soc. Am.* **57**:1486–1490, 1967.
10. R. L. Lamberts, "The Production and Use of Variable-Transmittance Sinusoidal Test Objects," *Appl. Opt.* **2**:273–276, 1963.
11. A. Daniels, G. Boreman, A. Ducharme, and E. Sapiro, "Random Transparency Targets for Modulation Transfer Function Measurement in the Visible and Infrared," *Opt. Eng.* **34**:860–868, 1995.
12. J. W. Coltman, "The Specification of Imaging Properties by Response to a Sine Wave Input," *J. Opt. Soc. Am.* **44**:468, 1954.
13. G. Boreman and S. Yang, "Modulation Transfer Function Measurement Using Three- and Four-Bar Targets," *Appl. Opt.* **34**:8050–8052, 1995.
14. D. H. Kelly, "Spatial Frequency, Bandwidth, and Resolution," *Appl. Opt.* **4**:435, 1965.
15. G. D. Boreman, *Modulation Transfer Function in Optical and Electro-Optical Systems*, SPIE Press, Bellingham, WA, 2001, chap. 4.
16. D. Malacara, *Optical Shop Testing*, Wiley, New York, 1978, chap. 3.

This page intentionally left blank.

DO NOT DUPLICATE

COHERENCE THEORY

William H. Carter*

*Naval Research Laboratory
Washington, D.C.*

5.1 GLOSSARY

I	intensity (use irradiance or field intensity)
k	radian wave number
\mathbf{p}	unit propagation vector
t	time
U	field amplitude
u	Fourier transform of U
W_ω	cross-spectral density function
\mathbf{x}	spatial vector
$\Gamma_{12}(\tau)$	mutual coherence function
$\Delta\ell$	coherence length
$\Delta\tau$	coherence time
μ_ω	complex degree of spatial coherence
ϕ	phase
ω	radian frequency
Real ()	real part of ()

5.2 INTRODUCTION

Classical Coherence Theory

All light sources produce fields that vary in time with highly complicated and irregular waveforms. Because of diffraction, these waveforms are greatly modified as the fields propagate. All light detectors measure the intensity time averaged over the waveform. This measurement depends on the integration

*Deceased. The author was a visiting scientist at the Johns Hopkins University Applied Physics Laboratory when this chapter was written.

time of the detector and the waveform of the light at the detector. Generally this waveform is not precisely known. Classical coherence theory¹⁻⁸ is a mathematical model which is very successful in describing the effects of this unknown waveform on the observed measurement of time-averaged intensity. It is based on the electromagnetic wave theory of light as formulated from Maxwell's equations, and uses statistical techniques to analyze the effects due to fluctuations in the waveform of the field in both time and space.

Quantum Coherence Theory

Classical coherence theory can deal very well with almost all presently known optical coherence phenomena; however, a few laboratory experiments require a much more complicated mathematical model, quantum coherence theory⁹⁻¹¹ to explain them. This theory is not based on classical statistical theory, but is based on quantum electrodynamics.¹²⁻¹⁴ While the mathematical model underlying classical coherence theory uses simple calculus, quantum coherence theory uses the Hilbert space formulation of abstract linear algebra, which is very awkward to apply to most engineering problems. Fortunately, quantum coherence theory appears to be essential as a model only for certain unusual (even though scientifically very interesting) phenomena such as squeezed light states and photon antibunching. All observed naturally occurring phenomena outside of the laboratory appear to be modeled properly by classical coherence theory or by an approximate semiclassical quantum theory. This chapter will deal only with the simple classical model.

5.3 SOME ELEMENTARY CLASSICAL CONCEPTS

Analytical Signal Representation

Solutions of the time-dependent, macroscopic Maxwell's equations yield six scalar components of the electric and the magnetic fields which are functions of both time and position in space. As in conventional diffraction theory, it is much more convenient to treat monochromatic fields than it is to deal with fields that have complicated time dependencies. Therefore, each of these scalar components is usually represented at some typical point in space (given with respect to some arbitrary origin by the radius vector $\mathbf{x} = (x, y, z)$) by a superposition of monochromatic real scalar components. Thus the field amplitude for a typical monochromatic component of the field with radial frequency ω is given by

$$U_r(\mathbf{x}, \omega) = U_0(\mathbf{x}) \cos[\phi(\mathbf{x}) - \omega t] \quad (1)$$

where $U_0(\mathbf{x})$ is the field magnitude and $\phi(\mathbf{x})$ is the phase. Trigonometric functions like that in Eq. (1) are awkward to manipulate. This is very well known in electrical circuit theory. Thus, just as in circuit theory, it is conventional to represent this field amplitude by a "phasor" defined by

$$U(\mathbf{x}, \omega) = U_0(\mathbf{x}) e^{i\phi(\mathbf{x})} \quad (2)$$

The purpose for using this complex field amplitude, just as in circuit theory, is to eliminate the need for trigonometric identities when adding or multiplying field amplitudes. A time-dependent complex analytic signal (viz., Ref. 15, sec. 10.2) is usually defined as the Fourier transform of this phasor, i.e.,

$$u(\mathbf{x}, t) = \int_0^{\infty} U(\mathbf{x}, \omega) e^{-i\omega t} d\omega \quad (3)$$

The integration in Eq. (3) is only required from zero to infinity because the phasor is defined with hermitian symmetry about the origin, i.e., $U(-\mathbf{x}, \omega) = U^*(\mathbf{x}, \omega)$. Therefore, all of the information is contained within the domain from zero to infinity. To obtain the actual field component from the analytical signal just take the real part of it. The Fourier transform in Eq. (3) is well defined if the analytical signal represents a deterministic field. However, if the light is partially coherent, then the analytical signal is usually taken to be a stationary random process. In this case the Fourier inverse of

Eq. (3) does not exist. It is then possible to understand the spectral decomposition given by Eq. (3) only within the theory of generalized functions (viz., see Refs. 16, the appendix on generalized functions, and 17, pp. 25–30).

Scalar Field Amplitude

Each monochromatic component of an arbitrary deterministic light field propagating through a homogeneous, isotropic medium can always be represented using an angular spectrum of plane waves for each of the six scalar components of the vector field. The six angular spectra are coupled together by Maxwell's equations so that only two are independent.^{18–20} Any two of the six angular spectra can be used to define two scalar fields from which the complete vector field can be determined. A polarized light field can be represented in this way by only one scalar field.^{20,21} Thus it is often possible to represent one polarized component of a vector electromagnetic field by a single scalar field. It has also been found useful to represent completely unpolarized light by a single scalar field. In more complicated cases, where the polarization properties of the light are important, a vector theory is sometimes needed as discussed later under "Explicit Vector Representations."

Temporal Coherence and Coherence Time

Within a short enough period of time, the time dependence of any light field at a point in space can be very closely approximated by a sine wave (Ref. 15, sec. 7.5.8). The length of time for which this is a good approximation is usually called the *coherence time* $\Delta\tau$. The coherence time is simply related to the spectral bandwidth for any light wave by the uncertainty principle, i.e.,

$$\Delta\tau \Delta\omega \geq 1 \quad (4)$$

For a light wave which is also highly directional within some region of space (like a beam) so that it propagates generally in some fixed direction (given by the unit vector \mathbf{p}), the field amplitude is given by

$$u(\mathbf{x}, t) = f(\mathbf{p} \cdot \mathbf{x} - ct) \quad (5)$$

Such a traveling wave will be approximately sinusoidal (and hence coherent) over some coherence length $\Delta\ell$ in the direction of \mathbf{p} where from Eq. (4) we see that

$$\Delta\ell = c \Delta\tau \approx c/\Delta\omega \quad (6)$$

so that the coherence length varies inversely with bandwidth.

Spatial Coherence and Coherence Area

The time-dependent waveform for any light field is approximately the same at any point within a sufficiently small volume of space called the *coherence volume*. The projection of this volume onto a surface is termed a *coherence area*. If we have a field that, within some region, is roughly directional so that its field amplitude is given by Eq. (5), then the coherence length gives the dimension of the coherence volume in the direction of propagation \mathbf{p} , and the coherence area gives the dimensions of the coherence volume normal to this direction.

Measurements of Coherence

Coherence is usually measured by some form of interferometer that takes light from two test points in a light field, \mathbf{x}_1 and \mathbf{x}_2 , and then allows them to interfere after introducing a time advance τ in the light from \mathbf{x}_1 relative to that from \mathbf{x}_2 . If the field intensity of the interference pattern is measured as a function of τ , then in general it has the form (see Ref. 15, sec. 10.3.1)

$$I(\tau) = I(\mathbf{x}_1) + I(\mathbf{x}_2) + 2 \operatorname{Real}(\Gamma_{12}(\tau)) \quad (7)$$

where $I(\mathbf{x}_i)$ is the intensity at the i th test point, and $\Gamma_{12}(\tau)$ is the mutual coherence function which measures the τ advanced correlation between the waveforms at the two test points (as subsequently defined under “Mutual Coherence Function”). There are many interferometers which have been developed to measure $\Gamma_{12}(\tau)$ in this way. One of the earliest techniques was developed by Thompson and Wolf.²² They used a diffractometer to measure the coherence over a surface normal to the direction of propagation for a collimated beam from a partially coherent source. More recently, Carter²³ used an interferometer made from a grating and a microscope to similarly measure the coherence of a beam transverse to the direction of propagation.

5.4 DEFINITIONS OF COHERENCE FUNCTIONS

Mutual Coherence Function

In an older form of coherence theory¹ the principal coherence function was the mutual coherence function defined by²⁴

$$\Gamma_{12}(\tau) \approx \frac{1}{T \rightarrow \infty 2T} \int_{-T}^T u(\mathbf{x}_1, t + \tau) u^*(\mathbf{x}_2, t) dt \quad (8)$$

where $u(\mathbf{x}, t)$ represents the complex analytic time-dependent signal at some point \mathbf{x} and some time t as defined in Eq. (3). This definition was originally motivated by the fact that the intensity, as actually measured, is precisely this time averaged function with $\mathbf{x}_1 = \mathbf{x}_2$ and $\tau = 0$, and that this function is the most readily measured since it appears directly in Eq. (7). Thus it was clearly possible to measure $\Gamma_{12}(\tau)$ over some input plane, propagate it to an output plane,²⁵ and then find the intensity over the output plane from $\Gamma_{12}(\tau)$. It was assumed in this definition in Eq. (8) that $u(\mathbf{x}, t)$ is stationary in time so that $\Gamma_{12}(\tau)$ is only a function of τ and not of t . In most of the older literature, sharp brackets were usually used to represent this time average rather than an ensemble average (see Ref. 15, sec. 10.3.1). In the early 1960s it was found to be much more convenient to treat $u(\mathbf{x}, t)$ as an ergodic, stationary random process so that Eq. (8) could be replaced by

$$\Gamma_{12}(\tau) = \langle u(\mathbf{x}_1, t + \tau) u^*(\mathbf{x}_2, t) \rangle \quad (9)$$

where (everywhere in this chapter) the sharp brackets denote an ensemble average. After the change to ensemble averages the cross-spectral density function (to be defined shortly) became the most used correlation function in the coherence literature, because of the simpler and more general rules for its propagation (as discussed later under “Representations” on p. 5.16).

Complex Degree of Coherence

To obtain a function that depends only on the coherence properties of a light field it is often useful to normalize the mutual coherence function in the manner of

$$\gamma_{12}(\tau) = \frac{\langle u(\mathbf{x}_1, t + \tau) u^*(\mathbf{x}_2, t) \rangle}{\sqrt{\langle u(\mathbf{x}_1, t) u^*(\mathbf{x}_1, t) \rangle \langle u(\mathbf{x}_2, t) u^*(\mathbf{x}_2, t) \rangle}} \quad (10)$$

This is called the *complex degree of coherence*. It is a properly normalized correlation coefficient, so that $\gamma_{11}(0) = \gamma_{22}(0) = 1$. This indicates that the field at a point in space must always be perfectly coherent with itself. All other values of $\gamma_{12}(\tau)$ are generally complex with an amplitude less than one. This indicates that the fields at two different points, or at the same point after a time delay τ , are generally less than perfectly coherent with each other. The magnitude of the complete degree of spatial coherence (from zero to one) is a measure of the mutual coherence between the fields at the two test points and after a time delay τ .

Cross-Spectral Density Function

Just as in classical diffraction theory, it is much easier to propagate monochromatic light than light with a complicated time waveform. Thus the most frequently used coherence function is the cross-spectral density function, $W_\omega(\mathbf{x}_1, \mathbf{x}_2)$, which is the ensemble-averaged correlation function between a typical monochromatic component of the field at some point \mathbf{x}_1 with the complex conjugate of the same component of the field at some other point \mathbf{x}_2 . It may be defined by

$$\delta(\omega - \omega') W_\omega(\mathbf{x}_1, \mathbf{x}_2) = \langle U(\mathbf{x}_1, \omega) U^*(\mathbf{x}_2, \omega') \rangle \quad (11)$$

The amplitude $U(\mathbf{x}, \omega)$ for a field of arbitrary coherence is taken to be a random variable. Thus $U(\mathbf{x}, \omega)$ represents an ensemble of all of the possible fields, each of which is represented by a complex phasor amplitude like that defined in Eq. (2). The sharp brackets denote an ensemble average over all of these possible fields weighted by the probability for each of them to occur. The correlation functions defined by Eqs. (11) and (9) are related by the Fourier transform pairs

$$\Gamma_{12}(\tau) = \int_0^\infty W_\omega(\mathbf{x}_1, \mathbf{x}_2) e^{-i\omega\tau} d\omega \quad (12)$$

and

$$W_\omega(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{2\pi} \int_{-\infty}^\infty \Gamma_{12}(\tau) e^{i\omega\tau} d\tau \quad (13)$$

which is easily shown, formally, by substitution from Eq. (3) into (9) and then using (11). These relations represent a form of the generalized Wiener-Khinchine theorem (see Ref. 26, pp. 107–108).

Complex Degree of Spectral Coherence

Because the cross-spectral density function contains information about both the intensity (see “Intensity,” which follows shortly) and the coherence of the field, it is useful to define another coherence function which describes the coherence properties only. This is the complex degree of spectral coherence (not to be confused with the complex degree of *spatial* coherence, which is a totally different function), which is usually defined by²⁷

$$\mu_\omega(\mathbf{x}_1, \mathbf{x}_2) = \frac{W_\omega(\mathbf{x}_1, \mathbf{x}_2)}{\sqrt{W_\omega(\mathbf{x}_1, \mathbf{x}_1) W_\omega(\mathbf{x}_2, \mathbf{x}_2)}} \quad (14)$$

It is easy to show that this function is a properly normalized correlation coefficient which is always equal to unity if the field points are brought together, and is always less than or equal to unity as they are separated. If the magnitude of $\mu_\omega(\mathbf{x}_1, \mathbf{x}_2)$ is unity, it indicates that the monochromatic field component with radial frequency ω is perfectly coherent between the two points \mathbf{x}_1 and \mathbf{x}_2 . If the magnitude of this function is less than unity it indicates less-than-perfect coherence. If the magnitude is zero it indicates complete incoherence between the field amplitudes at the two test points. For most partially coherent fields the cross-spectral density function has significantly large values only for point separations which keep the two field points within the same coherence volume. This function depends only on the positions of the points and the single radial frequency that the field components at the two points share. Field components of different frequency are always uncorrelated (and therefore incoherent), even at the same point.

Spectrum and Normalized Spectrum

Recently, the changes in the spectrum of light due to propagation have been studied using coherence theory. It is therefore useful to define the spectrum of light as just the monochromatic intensity

(which is just the trace of the cross-spectral density function) as a function of ω , and the spectrum of a primary source as a very similar function, i.e.,

$$\begin{aligned} S_U(\mathbf{x}, \omega) &= \langle U_\omega(\mathbf{x}) U_\omega^*(\mathbf{x}) \rangle = W_U(\mathbf{x}, \mathbf{x}) \\ S_Q(\mathbf{x}, \omega) &= \langle \rho_\omega(\mathbf{x}) \rho_\omega^*(\mathbf{x}) \rangle = W_Q(\mathbf{x}, \mathbf{x}) \end{aligned} \quad (15)$$

where the subscript Q indicates that this is a primary source spectrum, and the subscript U indicates that this is a field spectrum. The spectrum for the primary source is a function of the phasor $\rho_\omega(\mathbf{x})$ which represents the currents and charges in this source as discussed under "Primary Sources" in the next section. It is also useful to normalize these spectra in the manner

$$s_A(\mathbf{x}, \omega) = \frac{S_A(\mathbf{x}, \omega)}{\int_0^\infty S_A(\mathbf{x}, \omega) d\omega} \quad (16)$$

where the subscript A can indicate either U or Q , and the normalized spectrum has the property

$$\int_0^\infty s_A(\mathbf{x}, \omega) d\omega = 1 \quad (17)$$

so that it is independent of the total intensity.

Angular Correlation Function

A new coherence function, introduced for use with the angular spectrum expansion of a monochromatic component of the field,²⁸ is the angular correlation function defined by

$$\mathcal{A}_\omega(\mathbf{p}_1, \mathbf{p}_2) = \langle A_\omega(\mathbf{p}_1) A_\omega^*(\mathbf{p}_2) \rangle \quad (18)$$

where $A_\omega(\mathbf{p}_i)$ is the angular spectrum which gives the complex amplitude of the plane wave component of the field which propagates in the direction given by the unit vector \mathbf{p}_i . This is related to the cross-spectral density function over the $z = 0$ plane by the equation

$$\mathcal{A}(\mathbf{p}_1, \mathbf{p}_2) = \frac{1}{\lambda^4} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W_\omega^{(0)}(\mathbf{x}'_1, \mathbf{x}'_2) e^{-ik(\mathbf{p}_1 \cdot \mathbf{x}'_1 - \mathbf{p}_2 \cdot \mathbf{x}'_2)} d^2 \mathbf{x}'_1 d^2 \mathbf{x}'_2 \quad (19)$$

This is the four-dimensional Fourier transform of the cross-spectral density function over the $z = 0$ plane. It represents a correlation function between the complex amplitudes of two plane wave components of the field propagating in the directions given by the unit vectors \mathbf{p}_1 and \mathbf{p}_2 , respectively. It can be used to calculate the cross-spectral density function (as described later under "Angular Spectrum Representation") between any pair of points away from the $z = 0$ plane, assuming that the field propagates in a source-free homogeneous medium. In this chapter we will use single-primed vectors, as in Eq. (19), to represent radius vectors from the origin to points within the $z = 0$ plane, i.e., $\mathbf{x}' = (x', y', 0)$, as shown in Fig. 1.

All other vectors, such as \mathbf{x} or \mathbf{x}'' , are to be taken as three-dimensional vectors. Generally \mathbf{s} and \mathbf{p} are three-dimensional unit vectors indicating directions from the origin, a superscript (0) on a function indicates that it is the boundary condition for that function over the $z = 0$ plane, and a superscript (∞) on a function indicates that it is the asymptotic value for that function on a sphere of constant radius R from the origin as $R \rightarrow \infty$.

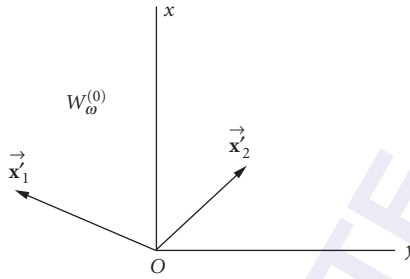


FIGURE 1 Illustrating the coordinate system and notation used for planar sources in the $z = 0$ plane.

Intensity

The intensity is usually considered to be the observable quantity in coherence theory. Originally it was defined to be the trace of the mutual coherence function as defined by Eq. (8), i.e.,

$$I(\mathbf{x}_1) \triangleq \Gamma_{11}(0) \approx \frac{1}{T \rightarrow \infty} \int_{-T}^T u(\mathbf{x}_1, t) u^*(\mathbf{x}_1, t) dt \quad (20)$$

which is always real. Thus it is the time-averaged square magnitude of the analytic signal. This represents the measurement obtained by the electromagnetic power detectors always used to detect light fields. Since the change to ensemble averages in coherence theory, it is almost always assumed that the analytic signal is an ergodic random process so that the intensity can be obtained from the equation

$$I(\mathbf{x}_1) \triangleq \Gamma_{11}(0) = \langle u(\mathbf{x}_1, t) u^*(\mathbf{x}_1, t) \rangle \quad (21)$$

where the sharp brackets indicate an ensemble average. Usually, in most recent coherence-theory papers, the intensity calculated is actually the spectrum, which is equivalent to the intensity of a single monochromatic component of the field which is defined as the trace of the cross-spectral density function as given by

$$I_\omega(\mathbf{x}_1) \triangleq W_\omega(\mathbf{x}_1, \mathbf{x}_1) = \langle U(\mathbf{x}_1, \omega) U^*(\mathbf{x}_1, \omega) \rangle \quad (22)$$

Since different monochromatic components of the field are mutually incoherent and cannot interfere, we can always calculate the intensity of the total field as the sum over the intensities of its monochromatic components in the manner

$$I(\mathbf{x}_1) \triangleq \int_0^\infty I_\omega(\mathbf{x}_1) d\omega \quad (23)$$

Since in most papers on coherence theory the subscript omega is usually dropped, the reader should be careful to observe whether or not the intensity calculated is for a monochromatic component of the field only. If it is, the total measurable intensity can be obtained simply by summing over all omega, as indicated in Eq. (23).

Radiant Emittance

In classical radiometry the *radiant emittance* is defined to be the power radiated into the far field by a planar source per unit area. A wave function with some of the properties of radiant emittance has been defined by Marchand and Wolf using coherence theory [see Refs. 29, eq. (32), and 30].

However, because of interference effects, the far-field energy cannot be subdivided into components that can be traced back to the area of the source that produced them. The result is that the radiant emittance defined by Marchand and Wolf is not nonnegative definite (as it is in classical radiometry) except for the special case of a completely incoherent source. Since, as discussed in the next section under “Perfectly Incoherent Source,” perfectly incoherent sources exist only as a limiting case, radiant emittance has not been found to be a very useful concept in coherence theory.

Radiant Intensity

In classical radiometry the *radiant intensity* is defined to be the power radiated from a planar source into a unit solid angle with respect to an origin at the center of the surface. This can be interpreted in coherence theory as the field intensity over a portion of a surface in the far field of the source, in some direction given by the unit vector \mathbf{s} , which subtends a unit solid angle from the source. Thus the radiant intensity for a monochromatic component of the field can be defined in coherence theory as

$$J_{\omega}(\mathbf{s}) \approx W_{\omega}^{(\infty)}(R\mathbf{s}, R\mathbf{s})R^2 \quad (24)$$

To obtain the total radiant intensity we need only sum this function over all omega.

Radiance

In classical coherence theory the radiance function is the power radiated from a unit area on a planar source into a unit solid angle with respect to an origin at the center of the source. In the geometrical optics model, from which this concept originally came, it is consistent to talk about particles of light leaving an area at a specified point on a surface to travel in some specified direction. However, in a wave theory, wave position and wave direction are Fourier conjugate variables. We can have a configuration space wave function (position) or a momentum space wave function (direction), but not a wave function that depends independently on both position and direction. Thus the behavior of a wave does not necessarily conform to a model which utilizes a radiance function.³¹ Most naturally occurring sources are quasi-homogeneous (discussed later). For such sources, a radiance function for a typical monochromatic component of the field can be defined as the Wigner distribution function^{32–35} of the cross-spectral density function over the $z = 0$ plane, which is given by

$$B_{\omega}(\mathbf{x}'_+, \mathbf{s}) = \frac{\cos\theta}{\lambda^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W_{\omega}^{(0)}(\mathbf{x}'_+ + \mathbf{x}'_-/2, \mathbf{x}'_+ - \mathbf{x}'_-/2) e^{-i\mathbf{k}\mathbf{s}\cdot\mathbf{x}'_-} d^2\mathbf{x}'_- \quad (25)$$

where θ is the angle that the unit vector \mathbf{s} makes with the $+z$ axis. For quasi-homogeneous fields this can be associated with the energy radiated from some point \mathbf{x}'_+ into the far field in the direction \mathbf{s} . Such a definition for radiance also works approximately for some other light fields, but for light which does not come from a quasi-homogeneous source, no such definition is either completely equivalent to a classical radiance function³¹ or unique as an approximation to it. Much progress has been made toward representing a more general class of fields using a radiance function.³⁶ In general, waves do not have radiance functions.

Higher-Order Coherence Functions

In general, the statistical properties of a random variable are uniquely defined by the probability density function which can be expanded into a series which contains correlation functions of all orders. Thus, in general, all orders of correlation functions are necessary to completely define the statistical properties of the field. In classical coherence theory we usually assume that the partially coherent fields arise from many independent sources so that, by the central limit theorem of statistics,

the probability distribution function for the real and imaginary components of the phasor field amplitude are zero-mean gaussian random variables (See Ref. 8, sec. 2.72e). Thus, from the gaussian moment theorem, the field correlation functions of any order can be calculated from the second-order correlation functions, for example,

$$\langle U(\mathbf{x}_1, \omega) U^*(\mathbf{x}_1, \omega) U(\mathbf{x}_2, \omega) U^*(\mathbf{x}_2, \omega) \rangle = I_\omega(\mathbf{x}_1) I_\omega(\mathbf{x}_2) + |W_\omega(\mathbf{x}_1, \mathbf{x}_2)|^2 \quad (26)$$

Thus, for gaussian fields, the second-order correlation functions used in coherence theory completely define the statistical properties of the field. Some experiments, such as those involving intensity interferometry, actually measure fourth- or higher-order correlation functions.³⁷

5.5 MODEL SOURCES

Primary Sources

In coherence theory it is useful to talk about primary and secondary sources. A primary source distribution is the usual source represented by the actual charge and current distribution which give rise to the field. For propagation from a primary source through a source-free media, the field amplitude is defined by the inhomogeneous Helmholtz equation [see Ref. 38, eq. (6.57)], i.e.,

$$\left(\nabla^2 + \frac{\omega^2}{c^2} \right) U_\omega(\mathbf{x}) = -4\pi\rho_\omega(\mathbf{x}) \quad (27)$$

where $\rho_\omega(\mathbf{x})$ represents the charge-current distribution in the usual manner. A solution to this wave equation gives

$$W_U(\mathbf{x}_1, \mathbf{x}_2) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W_Q(\mathbf{x}_1'', \mathbf{x}_2'') K(\mathbf{x}_1, \mathbf{x}_1'') K^*(\mathbf{x}_2, \mathbf{x}_2'') d^3\mathbf{x}_1'' d^3\mathbf{x}_2'' \quad (28)$$

for the cross-spectral density function, where $k = \omega/c = 2\pi/\lambda$

$$K(\mathbf{x}, \mathbf{x}'') = \frac{e^{ik|\mathbf{x}-\mathbf{x}''|}}{|\mathbf{x}-\mathbf{x}''|} \quad (29)$$

is the free-space propagator for a primary source, $W_U(\mathbf{x}_1, \mathbf{x}_2)$ is the cross-spectral density function for the fields (with suppressed ω -dependence), as defined by Eq. (11), and

$$W_Q(\mathbf{x}_1'', \mathbf{x}_2'') = \langle \rho_\omega(\mathbf{x}_1'') \rho_\omega^*(\mathbf{x}_2'') \rangle \quad (30)$$

is the cross-spectral density function for the source. This three-dimensional primary source can be easily reduced to a two-dimensional source over the $z = 0$ plane by simply defining the charge current distribution to be

$$\rho_\omega(\mathbf{x}) = \rho'_\omega(x, y) \delta(z) \quad (31)$$

where $\delta(z)$ is the Dirac delta function.

Secondary Sources

Often, however, it is more convenient to consider a field which arises from sources outside of the region in space in which the fields are of interest. Then it is sometime useful to work with boundary conditions for the field amplitude over some surface bounding the region of interest. In many coherence problems these boundary conditions are called a *planar secondary source* even though they are actually treated as boundary conditions. For example, most conventional diffraction equations

assume the boundary condition over the $z = 0$ plane is known and use it to calculate the fields in the $z > 0$ half-space. Then the field obeys the homogeneous Helmholtz equation [See Ref. 38, eq. (7.8)], i.e.,

$$\left(\nabla^2 + \frac{\omega^2}{c^2}\right)U_\omega(\mathbf{x}) = 0 \quad (32)$$

which has the solution

$$W_\omega(\mathbf{x}_1, \mathbf{x}_2) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W_\omega^{(0)}(\mathbf{x}'_1, \mathbf{x}'_2) h(\mathbf{x}_1, \mathbf{x}'_1) h^*(\mathbf{x}_2, \mathbf{x}'_2) d^2\mathbf{x}'_1 d^2\mathbf{x}'_2 \quad (33)$$

where $W_\omega^{(0)}(\mathbf{x}'_1, \mathbf{x}'_2)$ is the boundary condition for the cross-spectral density function of the fields over the $z = 0$ plane (as shown in Fig. 1), $W_\omega(\mathbf{x}_1, \mathbf{x}_2)$ is the cross-spectral density function anywhere in the $z > 0$ half-space, and

$$h(\mathbf{x}, \mathbf{x}') = \frac{-1}{2\pi} \frac{d}{dz} \frac{e^{ik|\mathbf{x}-\mathbf{x}'|}}{|\mathbf{x}-\mathbf{x}'|} \quad (34)$$

is the free-space propagator for the field amplitude. This is a common example of a secondary source over the $z = 0$ plane.

Perfectly Coherent Source

The definition of a perfectly coherent source, within the theory of partial coherence, is somewhat complicated. The light from an ideal monochromatic source is, of course, always perfectly coherent. Such a field produces high-contrast interference patterns when its intensity is detected by any suitable instrument. However, it is possible that light fields exist that are not monochromatic but can also produce similar interference patterns and therefore must be considered coherent. The ability of light fields to produce interference patterns at a detector is measured most directly by the complex degree of coherence $\gamma_{12}(\tau)$, defined by Eq. (10). If a field has a complex degree of coherence that has unit magnitude for every value of τ and for every point pair throughout some domain D , then light from all points in D will combine to produce high-contrast interference fringes.³⁹ Such a field is defined to be perfectly coherent within D . Mandel and Wolf²⁷ have shown that the mutual coherence function for such a field factors within D in the manner

$$\Gamma_{12}(\tau) = \psi(\mathbf{x}_1) \psi^*(\mathbf{x}_2) e^{-i\omega\tau} \quad (35)$$

We will take Eq. (35) to be the definition of a perfectly coherent field. Coherence is not as easily defined in the space-frequency domain because it depends on the spectrum of the light as well as on the complex degree of spectral coherence. For example, consider a field for which every monochromatic component is characterized by a complex degree of spectral coherence which has unit magnitude between all point pairs within some domain D . Mandel and Wolf⁴⁰ have shown that for such a field the cross-spectral density function within D factors in

$$W_\omega(\mathbf{x}_1'' \mathbf{x}_2'') = \bar{U}(\mathbf{x}_1'', \omega) \bar{U}^*(\mathbf{x}_2'', \omega) \quad (36)$$

However, even if Eq. (36) holds for this field within D , the field may not be perfectly coherent [as perfect coherence is defined by Eq. (35)] between all points within the domain. In fact it can be completely incoherent between some points within D .⁴¹ A secondary source covering the $z = 0$ plane with a cross-spectral density function over that plane which factors as given by Eq. (36) will produce a field in the $z > 0$ half-space (filled with free space or a homogeneous, isotropic dielectric) which has a cross-spectral density function that factors in the same manner everywhere in the half-space. This can be easily shown by substitution from Eq. (36) into Eq. (28), using Eq. (29). But, even if this is true for every monochromatic component of the field, Eq. (35) may not hold within the half-space for every point pair, so we cannot say that the field is perfectly coherent there. Perfectly coherent light

sources never actually occur, but sometimes the light from a laser can behave approximately in this way over some coherence volume which is usefully large. Radio waves often behave in this manner over very large coherence volumes.

Quasi-Monochromatic Source

In many problems it is more useful not to assume that a field is strictly monochromatic but instead to assume that it is only quasi-monochromatic so that the time-dependent field amplitude can be approximated by

$$u(\mathbf{x}, t) = u_0(\mathbf{x}, t)e^{-i\omega t} \quad (37)$$

where $u_0(\mathbf{x}, t)$ is a random process which varies much more slowly in time than $e^{-i\omega t}$. Then the mutual coherence function and the complex degree of spatial coherence can be usefully approximated by (see Ref. 15, sec. 10.4.1)

$$\begin{aligned} \Gamma_{12}(\boldsymbol{\tau}) &= \Gamma_{12}(0)e^{-i\omega\boldsymbol{\tau}} \\ \gamma_{12}(\boldsymbol{\tau}) &= \gamma_{12}(0)e^{-i\omega\boldsymbol{\tau}} \end{aligned} \quad (38)$$

within coherence times much less than the reciprocal bandwidth of the field, i.e., $\Delta\boldsymbol{\tau} \ll 1/\Delta\omega$. In the pre-1960 coherence literature, $\Gamma_{12}(0)$ was called the *mutual intensity*. Monochromatic diffraction theory was then used to define the propagation properties of this monochromatic function. It was used instead of the cross-spectral density function to formulate the theory for the propagation of a partially coherent quasi-monochromatic field. While this earlier form of the theory was limited by the quasi-monochromatic approximation and was therefore not appropriate for very wideband light, the newer formulation (in terms of the cross-spectral density function) makes no assumptions about the spectrum of the light and can be applied generally. The quasi-monochromatic approximation is still very useful when dealing with radiation of very high coherence.⁴²

Schell Model Source

Other, more general, source models have been developed. The most general of these is the Schell model source (see Refs. 43, sec. 7.5, and 44), for which we only assume that the complex degree of spectral coherence for either a primary or secondary source is stationary in space, so that from Eq. (14) we have

$$W_A(\mathbf{x}_1, \mathbf{x}_2) = \mu_A(\mathbf{x}_1 - \mathbf{x}_2) \sqrt{W_A(\mathbf{x}_1, \mathbf{x}_1)W_A(\mathbf{x}_2, \mathbf{x}_2)} \quad (39)$$

where the subscript A stands for U in the case of a Schell model secondary source and Q in the case of a Schell model primary source. The Schell model does not assume low coherence and, therefore, can be applied to spatially stationary light fields of any state of coherence. The Schell model of the form shown in Eq. (39) has been used to represent both three-dimensional primary sources^{45,46} and two-dimensional secondary sources.^{43,47,48}

Quasi-Homogeneous Source

If the intensity of a Schell model source is essentially constant over any coherence area, then Eq. (39) may be approximated by

$$W_A(\mathbf{x}_1, \mathbf{x}_2) = \mu_A(\mathbf{x}_1 - \mathbf{x}_2) I_A[(\mathbf{x}_1 + \mathbf{x}_2)/2] \quad (40)$$

where the subscript A can be either U , for the case of a quasi-homogeneous secondary source or Q , for the case of a quasi-homogeneous primary source. This equation is very useful in coherence theory because of the important exact mathematical identity.⁴⁹

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mu_{\omega}^{(0)}(\mathbf{x}'_1 - \mathbf{x}'_2) I_{\omega}^{(0)}[(\mathbf{x}'_1 + \mathbf{x}'_2)/2] e^{-ik(\mathbf{x}'_1 \cdot \mathbf{p}_1 - \mathbf{x}'_2 \cdot \mathbf{p}_2)} dx'_1 dy'_1 dx'_2 dy'_2 \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mu_{\omega}^{(0)}(\mathbf{x}'_+) e^{-ik(\mathbf{x}'_+ \cdot \mathbf{p}_+)} dx'_+ dy'_+ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_{\omega}^{(0)}(\mathbf{x}'_+) e^{-ik(\mathbf{x}'_+ \cdot \mathbf{p}_-)} dx'_+ dy'_+ \end{aligned} \quad (41)$$

where

$$\mathbf{x}'_+ = (\mathbf{x}'_1 + \mathbf{x}'_2)/2 \quad \mathbf{x}'_- = \mathbf{x}'_1 - \mathbf{x}'_2$$

and

$$\mathbf{p}_+ = (\mathbf{p}_1 + \mathbf{p}_2)/2 \quad \mathbf{p}_- = \mathbf{p}_1 - \mathbf{p}_2$$

which allows the four-dimensional Fourier transforms that occur in propagating the correlation functions for secondary sources [for example, Eqs. (49) or (54)] to be factored into a product of two-dimensional Fourier transforms. An equivalent identity also holds for the six-dimensional Fourier transform of the cross-spectral density function for a primary source, reducing it to the product of two three-dimensional Fourier transforms. This is equally useful in dealing with propagation from primary sources [for example, Eq. (53)]. This model is very good for representing two-dimensional secondary sources with sufficiently low coherence that the intensity does not vary over the coherence area on the input plane.^{49,50} It has also been applied to primary three-dimensional sources,^{45,46} to primary and secondary two-dimensional sources,^{51,52} and to three-dimensional scattering potentials.^{53,54}

Perfectly Incoherent Source

If the coherence volume of a field becomes much smaller than any other dimensions of interest in the problem, then the field is said to be *incoherent*. It is believed that no field can be incoherent over dimensions smaller than the order of a light wavelength. An incoherent field can be taken as a special case of a quasi-homogeneous field for which the complex degree of spectral coherence is approximated by a Dirac delta function, i.e.,

$$\begin{aligned} W_A(\mathbf{x}_1, \mathbf{x}_2) &= I(\mathbf{x}_1) \delta^2(\mathbf{x}_1 - \mathbf{x}_2) \\ W_Q(\mathbf{x}_1, \mathbf{x}_2) &= I(\mathbf{x}_1) \delta^3(\mathbf{x}_1 - \mathbf{x}_2) \end{aligned} \quad (42)$$

where the two-dimensional Dirac delta function is used for any two-dimensional source and a three-dimensional Dirac delta function is used only for a three-dimensional primary source. Even though this approximation is widely used, it is not a good representation for the thermal sources that predominate in nature. For example, the radiant intensity from a planar, incoherent source is not in agreement with Lambert's law. For thermal sources the following model is much better.

Thermal (Lambertian) Source

For a planar, quasi-homogeneous source to have a radiant intensity in agreement with Lambert's law it is necessary for the complex degree of spectral coherence to have the form

$$\mu_A(\mathbf{x}_1 - \mathbf{x}_2) = \frac{\sin(k|\mathbf{x}_1 - \mathbf{x}_2|)}{k|\mathbf{x}_1 - \mathbf{x}_2|} \quad (43)$$

to which arbitrary spatial frequency components with periods less than a wavelength can be added since they do not affect the far field.⁵⁵ It can also be shown that, under frequently obtained conditions, blackbody radiation has such a complex degree of spectral coherence.⁵⁶ It is believed that most naturally occurring light can be modeled as quasi-homogeneous with this correlation function.

5.6 PROPAGATION

Perfectly Coherent Light

Perfectly coherent light propagates according to conventional diffraction theory. By substitution from Eq. (36) into Eq. (33) using Eq. (34) we obtain

$$U_{\omega}(\mathbf{x}) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} U_{\omega}^{(0)}(\mathbf{x}') \frac{1}{2\pi} \frac{d}{dz'} \frac{e^{ik|\mathbf{x}-\mathbf{x}'|}}{|\mathbf{x}-\mathbf{x}'|} d^2\mathbf{x}' \quad (44)$$

which is just Rayleigh's diffraction integral of the first kind. Perfectly coherent light does not lose coherence as it propagates through any medium which is time-independent. For perfectly coherent light propagating in time-independent media, coherence theory is not needed.

Hopkin's Formula

In 1951 Hopkins⁵⁷ published a formula for the complex degree of spatial coherence for the field from a planar, secondary, incoherent, quasi-monochromatic source after propagating through a linear optical system with spread function $h(\mathbf{x}, \mathbf{x}')$, i.e.,

$$\gamma_{12}(0) = \frac{1}{\sqrt{I(\mathbf{x}_1)I(\mathbf{x}_2)}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_{\omega}^{(0)}(\mathbf{x}') h(\mathbf{x}_1 - \mathbf{x}') h^*(\mathbf{x}_2 - \mathbf{x}') d^2\mathbf{x}' \quad (45)$$

where $I_{\omega}^{(0)}(\mathbf{x}')$ is the intensity over the source plane. This formula can be greatly generalized to give the complex degree of spectral coherence for the field from any planar, quasi-homogeneous, secondary source⁵⁸ after transmission through this linear optical system, i.e.,

$$\mu_{\omega}(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{\sqrt{I(\mathbf{x}_1)I(\mathbf{x}_2)}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_{\omega}^{(0)}(\mathbf{x}') h(\mathbf{x}_1, \mathbf{x}') h^*(\mathbf{x}_2, \mathbf{x}') d^2\mathbf{x}' \quad (46)$$

provided that the spread function $h(\mathbf{x}, \mathbf{x}')$ can be assumed to be constant in its \mathbf{x}' dependence over any coherence area in the source plane.

van Cittert-Zernike Theorem

Hopkins' formula can be specialized for free-space propagation to calculate the far-field coherence properties of planar, secondary, quasi-homogeneous sources of low coherence. In 1934 van Cittert⁵⁹ and, later, Zernike⁶⁰ derived a formula equivalent to

$$\mu_{\omega}(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{\sqrt{I(\mathbf{x}_1)I(\mathbf{x}_2)}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_{\omega}^{(0)}(\mathbf{x}') \frac{e^{ik[|\mathbf{x}_1-\mathbf{x}'|+|\mathbf{x}_2-\mathbf{x}'|]}}{|\mathbf{x}_1-\mathbf{x}'||\mathbf{x}_2-\mathbf{x}'|} d^2\mathbf{x}' \quad (47)$$

for the complex degree of spectral coherence between any pair of points in the field radiated from an incoherent planar source, assuming that the points are not within a few wavelengths of the source. We can obtain Eq. (47) by substitution from Eq. (34) into Eq. (46) and then approximating

the propagator in a standard manner [see Ref. 61, eq. (7)]. Assume next, that the source area is contained within a circle of radius a about the origin in the source plane as shown in Fig. 4. Then, if the field points are both located on a sphere of radius R , which is outside of the Rayleigh range of the origin (i.e., $|\mathbf{x}_1| = |\mathbf{x}_2| = R \gg ka^2$), we can apply the Fraunhofer approximation to Eq. (47) to obtain

$$\mu_{12}(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{\sqrt{I(\mathbf{x}_1)I(\mathbf{x}_2)}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_{\omega}^{(0)}(\mathbf{x}') e^{ik\mathbf{x}' \cdot (\mathbf{x}_1 - \mathbf{x}_2)/R} d^2\mathbf{x}' \quad (48)$$

This formula is very important in coherence theory. It shows that the complex degree of spectral coherence from any planar, incoherent source with an intensity distribution $I_{\omega}^{(0)}(\mathbf{x}')$ has the same dependence on $(\mathbf{x}_1 - \mathbf{x}_2)$, over a sphere of radius R in the far field, as the diffraction pattern from a closely related perfectly coherent planar source with a real amplitude distribution proportional to $I_{\omega}^{(0)}(\mathbf{x}')$ (see Ref. 15, sec. 10.4.2a). Equation (48) can also be applied to a planar, quasi-homogeneous source⁴⁹ that is not necessarily incoherent, as will be shown later under “Reciprocity Theorem.”

Angular Spectrum Representation

Much more general equations for propagation of the cross-spectral density function can be obtained. Equation (33) is one such expression. Another can be found if we expand the fields in a source-free region of space into an angular spectrum of plane waves. Then we find that the cross-spectral density function over any plane can be calculated from the same function over a parallel plane using a linear systems approach. For an example, consider the two planes illustrated in Fig. 2.

We assume that the cross-spectral density function is known over the $z = 0$ plane in the figure, and we wish to calculate this function over the $z = d$ plane. To do this we first take the Fourier transform of the cross-spectral density function over the $z = 0$ plane according to Eq. (19)

$$\mathcal{A}_{\text{in}}(\mathbf{p}_1, \mathbf{p}_2) = \frac{1}{\lambda^4} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W_{\omega}^{(0)}(\mathbf{x}'_1, \mathbf{x}'_2) e^{-ik(\mathbf{p}_1 \cdot \mathbf{x}'_1 - \mathbf{p}_2 \cdot \mathbf{x}'_2)} d^2\mathbf{x}'_1 d^2\mathbf{x}'_2 \quad (49)$$

to obtain the angular correlation function in which all of the phase differences between the plane wave amplitudes are given relative to the point at the origin. Second, we shift the phase reference

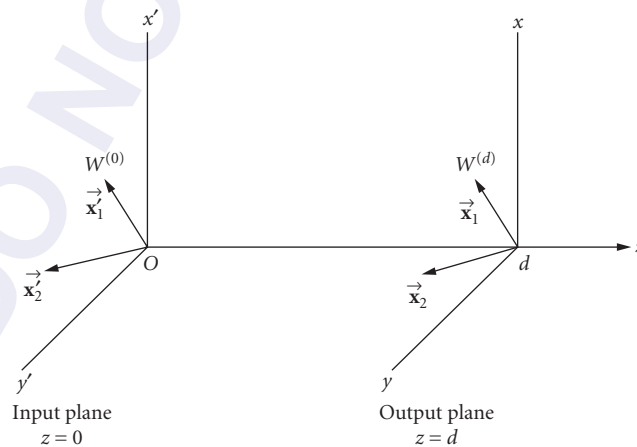


FIGURE 2 Illustrating the coordinate system for propagation of the cross-spectral density function using the angular spectrum of plane waves.

from the origin to the point on the z axis in the output plane by multiplying the angular correlation function by a transfer function, i.e.,

$$\mathcal{A}_{\text{out}}(\mathbf{p}_1, \mathbf{p}_2) = \mathcal{A}_{\text{in}}(\mathbf{p}_1, \mathbf{p}_2) \exp [ik(m_1 - m_2)d] \quad (50)$$

where d is the distance from the input to the output plane along the z axis (for back propagation d will be negative) and m_i is the third component of the unit vector $\mathbf{p}_i = (p_i, q_i, m_i)$, $i = 1$ or 2 , which is defined by

$$\begin{aligned} m_i &= \sqrt{1 - p_i^2 - q_i^2} & \text{if } p_i^2 + q_i^2 \leq 1 \\ &= i\sqrt{p_i^2 + q_i^2 - 1} & \text{if } p_i^2 + q_i^2 > 1 \end{aligned} \quad (51)$$

and is the cosine of the angle that \mathbf{p}_i makes with the $+z$ axis for real m_i . Finally, to obtain the cross-spectral density function over the output plane, we simply take the Fourier inverse of $\mathcal{A}_{\text{out}}(\mathbf{p}_1, \mathbf{p}_2)$, i.e.,

$$W_{\omega}^{(d)}(\mathbf{x}_1, \mathbf{x}_2) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mathcal{A}_{\text{out}}(\mathbf{p}_1, \mathbf{p}_2) e^{ik(\mathbf{p}_1 \cdot \mathbf{x}_1 - \mathbf{p}_2 \cdot \mathbf{x}_2)} d^2\mathbf{p}_1 d^2\mathbf{p}_2 \quad (52)$$

where, in this equation only, we use \mathbf{x}_i to represent a two-dimensional radius vector from the point $(0, 0, d)$ to a field point in the $z = d$ plane, as shown in Fig. 2. This propagation procedure is similar to the method usually used to propagate the field amplitude in Fourier optics. In coherence theory, it is the cross-spectral density function for a field of any state of coherence that is propagated between arbitrary parallel planes using the linear systems procedure. The only condition for the validity of this procedure is that the volume between the planes must either be empty space or a uniform dielectric medium.

Radiation Field

The cross-spectral density function far away from any source, which has finite size, can be calculated using a particularly simple equation. Consider a primary, three-dimensional source located within a sphere of radius a , as shown in Fig. 3. For field points, \mathbf{x}_1 and \mathbf{x}_2 which are much farther from

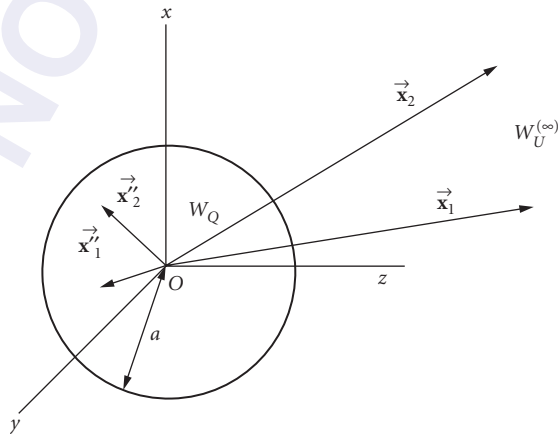


FIGURE 3 Illustrating the coordinate system used to calculate the cross-spectral density function in the far field of a primary, three-dimensional source.

the origin than the Rayleigh range ($|\mathbf{x}_i| \gg ka^2$), in any direction, a form of the familiar Fraunhofer approximation can be applied to Eq. (28) to obtain [see Ref. 62, eq. (2.5)]

$$W_U^{(\infty)}(\mathbf{x}_1, \mathbf{x}_2) = \frac{e^{ik(|\mathbf{x}_1| - |\mathbf{x}_2|)}}{|\mathbf{x}_1||\mathbf{x}_2|} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W_Q(\mathbf{x}'_1, \mathbf{x}'_2) \times e^{-ik[|\mathbf{x}_1, \mathbf{x}'_1|/|\mathbf{x}_1| - |\mathbf{x}_2, \mathbf{x}'_2|/|\mathbf{x}_2|]} d^3\mathbf{x}'_1 d^3\mathbf{x}'_2 \quad (53)$$

Thus the cross-spectral density function of the far field is proportional to the six-dimensional Fourier transform of the cross-spectral density function of its sources. A very similar expression can also be found for a two-dimensional, secondary source distribution over the $z = 0$ plane, as illustrated in Fig. 4.

If the sources are all restricted to the area within a circle about the origin of radius a , then the cross-spectral density function for all points which are outside of the Rayleigh range ($|\mathbf{x}_i| \gg ka^2$) from the origin and also in the $z > 0$ half-space can be found by using a different form of the Fraunhofer approximation in Eq. (33) [see Ref. 62, eq. (3.3)] to get

$$W_U^{(\infty)}(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{\lambda^2} \frac{z_1}{|\mathbf{x}_1|} \frac{z_2}{|\mathbf{x}_2|} \frac{e^{ik(|\mathbf{x}_1| - |\mathbf{x}_2|)}}{|\mathbf{x}_1||\mathbf{x}_2|} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W_U^{(0)}(\mathbf{x}'_1, \mathbf{x}'_2) \times e^{-ik[|\mathbf{x}_1, \mathbf{x}'_1|/|\mathbf{x}_1| - |\mathbf{x}_2, \mathbf{x}'_2|/|\mathbf{x}_2|]} d^2\mathbf{x}'_1 d^2\mathbf{x}'_2 \quad (54)$$

Because of both their relative simplicity (as Fourier transforms) and great utility, Eqs. (53) and (54) are very important in coherence theory for calculating both the radiant intensity and the far-field coherence properties of the radiation field from any source.

Representations

Several equations have been described here for propagating the cross-spectral density function. The two most general are Eq. (33), which uses an expansion of the field into spherical Huygens wavelets,

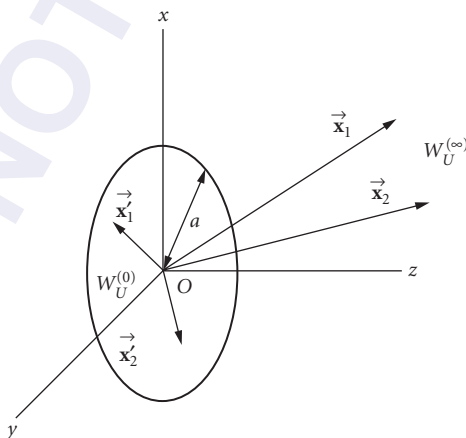


FIGURE 4 Illustrating the coordinate system for calculating the far field cross-spectral density function for a planar, secondary source distribution in the $z = 0$ plane. Single-primed coordinates indicate radius vectors from the origin to points within the $z = 0$ plane.

and Eqs. (49), (50), and (52), which use an expansion of the field into an angular spectrum of plane waves. These two formulations are completely equivalent. Neither uses any approximation not also used by the other method. The choice of which to use for a particular calculation can be made completely on a basis of convenience. The far-field approximations given by Eqs. (53) and (54), for example, can be derived from either representation. There are two bridges between the spherical and plane wave representations, one given by Weyl's integral, i.e.,

$$\frac{e^{ik|\mathbf{x}-\mathbf{x}''|}}{|\mathbf{x}-\mathbf{x}''|} = \frac{i}{\lambda} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{m} e^{ik[p(x-x'')+q(y-y'')+m|z-z''|]} dp dq \quad (55)$$

where

$$\begin{aligned} m &= \sqrt{1-p^2-q^2} & \text{if } p^2+q^2 \leq 1, \\ &= i\sqrt{p^2+q^2-1} & \text{if } p^2+q^2 > 1 \end{aligned} \quad (56)$$

and the other by a related integral

$$\frac{-1}{2\pi} \frac{d}{dz} \frac{e^{ik|\mathbf{x}-\mathbf{x}''|}}{|\mathbf{x}-\mathbf{x}''|} = \frac{\pm 1}{\lambda^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{ik[p(x-x'')+q(y-y'')+m|z-z''|]} dp dq \quad (57)$$

which can be easily derived from Eq. (55). In Eq. (57) the \pm sign holds according to whether $(z-z'') \geq 0$. With these two equations it is possible to transform back and forth between these two representations.

Reciprocity Theorem

The radiation pattern and the complex degree of spectral coherence obey a very useful reciprocity theorem for a quasi-homogeneous source. By substituting from Eq. (40) into Eq. (53), and using Eq. (24) and the six-dimensional form of Eq. (41), it has been shown that the radiant intensity in the direction of the unit vector \mathbf{s} from any bounded, three-dimensional, quasi-homogeneous primary source distribution is given by [see Ref. 45, eq. (3.11)]

$$J_{\omega}(\mathbf{s}) = J_0 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mu_Q(\mathbf{x}''_{-}) e^{-ik\mathbf{x}''_{-} \cdot \mathbf{s}} d^3\mathbf{x}''_{-} \quad (58)$$

where $\mu_Q(\mathbf{x}''_{-})$ is the (spatially stationary) complex degree of spectral coherence for the source distribution as defined in Eq. (40), and

$$J_0 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_Q(\mathbf{x}''_{+}) d^3\mathbf{x}''_{+} \quad (59)$$

where $I_Q(\mathbf{x}''_{+})$ is the intensity defined in Eq. (40). Note that the far-field radiation pattern depends, not on the source intensity distribution, but only on the source coherence. We also find from this calculation that the complete degree of spectral coherence between any two points in the far field of this source is given by [see Ref. 45, eq. (3.15)]

$$\mu^{(\infty)}u(R_1\mathbf{s}_1, R_2\mathbf{s}_2) = \frac{e^{ik(R_1-R_2)}}{J_0} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_Q(\mathbf{x}''_{+}) e^{ik\mathbf{x}''_{+} \cdot (\mathbf{s}_1 - \mathbf{s}_2)} d^3\mathbf{x}''_{+} \quad (60)$$

Note that the coherence of the far field depends, not on the source coherence, but rather on the source intensity distribution, $I_Q(\mathbf{x}''_{+})$. Equation (60) is a generalization of the van Cittert–Zernike theorem to three-dimensional, primary quasi-homogeneous sources, which are not necessarily incoherent. Equation (59) is a new theorem, reciprocal to the van Cittert–Zernike theorem, which was discovered by Carter and Wolf.^{45,49} Equations (58) and (60), taken together, give a reciprocity

relation. For quasi-homogeneous sources, far-field coherence is determined by source intensity alone and the far-field intensity pattern is determined by source coherence alone. Therefore, coherence and intensity are reciprocal when going from the source into the far field. Since most sources which occur naturally are believed to be quasi-homogeneous, this is a very useful theorem. This reciprocity theorem has been found to hold much more generally than just for three-dimensional primary sources. For a planar, secondary source, by substitution from Eq. (40) into Eq. (54) and then using Eq. (41), we obtain⁴⁹

$$J_{\omega}(\mathbf{s}) = J'_0 \cos^2 \theta \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mu_U^{(0)}(\mathbf{x}'_-) e^{-ik\mathbf{x}'_- \cdot \mathbf{s}} d^2\mathbf{x}'_- \quad (61)$$

where θ is the angle that \mathbf{s} makes with the $+z$ axis, and where

$$J'_0 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_U^{(0)}(\mathbf{x}'_+) d^2\mathbf{x}'_+ \quad (62)$$

and we also obtain

$$\mu_U^{(\infty)}(R_1\mathbf{s}_1, R_2\mathbf{s}_2) = \frac{1}{J'_0} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_U^{(0)}(\mathbf{x}'_+) e^{-ik\mathbf{x}'_+ \cdot (\mathbf{s}_1 - \mathbf{s}_2)} d^2\mathbf{x}'_+ \quad (63)$$

Very similar reciprocity relations also hold for scattering of initially coherent light from quasi-homogeneous scattering potentials,^{53,54} and for the scattering of laser beams from quasi-homogeneous sea waves.⁶³ Reciprocity relations which apply to fields that are not necessarily quasi-homogeneous have also been obtained.⁶⁴

Nonradiating Sources

One additional important comment must be made about Eq. (58). The integrals appearing in this equation form a three-dimensional Fourier transform of $\mu_Q(\mathbf{x}''_-)$. However, $J_{\omega}(\mathbf{s})$, the radiant intensity that this source radiates into the far field at radial frequency ω , is a function of only the direction of \mathbf{s} , which is a unit vector with a constant amplitude equal to one. It then follows that only the values of this transform over a spherical surface of unit radius from the origin affect the far field. Therefore, sources which have a complex degree of spectral coherence, $\mu_Q(\mathbf{x}''_-)$, which do not have spatial frequencies falling on this sphere, do not radiate at frequency ω . It appears possible from this fact to have sources in such a state of coherence that they do not radiate at all. Similar comments can also be made in regard to Eq. (53) which apply to all sources, even those that are not quasi-homogeneous. From Eq. (53) it is clear that only sources which have a cross-spectral density function with spatial frequencies (with respect to both \mathbf{x}''_1 and \mathbf{x}''_2) on a unit sphere can radiate. It is believed that this is closely related to the phase-matching condition in nonlinear optics.

Perfectly Incoherent Sources

For a completely incoherent primary source the far-field radiant intensity, by substitution from Eq. (42) into Eq. (58), is seen to be the same in all directions, independent of the shape of the source. By substitution from Eq. (40) into Eq. (28), and then using Eqs. (14), (29), and (42), we find that the complex degree of spectral coherence for this incoherent source is given by

$$\mu_U(\mathbf{x}_1, \mathbf{x}_2) = \frac{1}{\sqrt{I(\mathbf{x}_1)I(\mathbf{x}_2)}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_Q(\mathbf{x}'') \frac{e^{ik[|\mathbf{x}_1 - \mathbf{x}''| - |\mathbf{x}_2 - \mathbf{x}''|]}}{|\mathbf{x}_1 - \mathbf{x}''| |\mathbf{x}_2 - \mathbf{x}''|} d^3\mathbf{x}'' \quad (64)$$

Thus it is only the source intensity distribution that affects the field coherence. Comparison of this equation with Eq. (47) shows that this is a generalization of the van Cittert–Zernike theorem to primary

sources. It is clear from this equation that the complex degree of spectral coherence depends only on the shape of the source and not the fact that the source is completely incoherent. For a completely incoherent, planar, secondary source the radiant intensity is given by

$$J_{\omega}(\mathbf{s}) = J'_0 \cos^2 \theta \quad (65)$$

independent of the shape of the illuminated area on the source plane, where θ is the angle that \mathbf{s} makes with the normal to the source plane. This can be proven by substitution from Eq. (42) into Eq. (61). Note that such sources do not obey Lambert's law. The far-field coherence, again, depends on the source intensity as given by the van Cittert–Zernike theorem [see Eq. (47)] and not on the source coherence.

Spectrum

For a quasi-homogeneous, three-dimensional primary source the spectrum of the radiation $S_U^{(\infty)}(R\mathbf{s}, \omega)$ at a point $R\mathbf{s}$ in the direction \mathbf{s} (unit vector) and at a distance R from the origin, in the far field of the source, can be found, as a function of the source spectrum $S_Q^{(0)}(\omega)$, by substitution from Eqs. (24) and (15) into Eq. (58) to get

$$S_U^{(\infty)}(R\mathbf{s}, \omega) = \frac{c^3 S_Q^{(0)}(\omega)}{\omega^3 R^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mu_Q(\mathbf{x}''_-, \omega) e^{-ik\mathbf{x}''_- \cdot \mathbf{s}} d^3(k\mathbf{x}''_-) \quad (66)$$

where we explicitly indicate the dependence of the complex degree of spectral coherence on frequency. Notice that the spectrum of the field is not necessarily the same as the spectrum of the source and, furthermore, that it can vary from point to point in space. The field spectrum depends on the source coherence as well as on the source spectrum. A very similar propagation relation can be found for the far-field spectrum from a planar, secondary source in the $z = 0$ plane. By substitution from Eqs. (24) and (15) into Eq. (61) we get

$$S_U^{(\infty)}(R\mathbf{s}, \omega) = \frac{S_U^{(0)}(\omega)}{R^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mu_U^{(0)}(\mathbf{x}'_-, \omega) e^{-ik\mathbf{x}'_- \cdot \mathbf{s}} d^2(k\mathbf{x}'_-) \quad (67)$$

This shows that the spectrum for the field itself is changed upon propagation from the $z = 0$ plane into the far field and is different in different directions \mathbf{s} from the source.⁶⁵

5.7 SPECTRUM OF LIGHT

Limitations

The complex analytic signal for a field that is not perfectly coherent, as defined in Eq. (3), is usually assumed to be a time-stationary random process. Therefore, the integral

$$\int_{-\infty}^{\infty} u(\mathbf{x}, t) e^{i\omega t} dt \quad (68)$$

does not converge, so that the analytic signal does not have a Fourier transform. Therefore, it is only possible to move freely from the space-time domain to the space-frequency domain along the path shown in Fig. 5. Equation (3) does not apply to time-stationary fields within the framework of ordinary function theory.

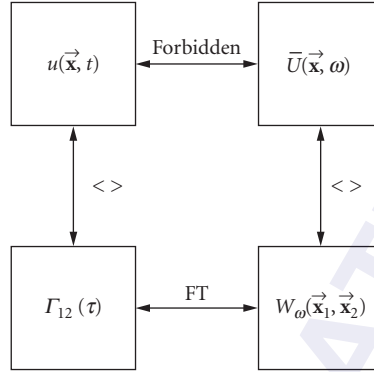


FIGURE 5 Illustrating the transformations which are possible between four functions that are used in coherence theory.

Coherent Mode Representation

Coherent mode representation (see Wolf^{66,67}) has shown that any partially coherent field can be represented as the sum over component fields that are each perfectly self-coherent, but mutually incoherent with each other. Thus the cross-spectral density function for any field can be represented in the form

$$W_A(\mathbf{x}_1, \mathbf{x}_2) = \sum_n \lambda_n \phi_{A,n}(\mathbf{x}_1) \phi_{A,n}^*(\mathbf{x}_2) \quad (69)$$

where $\phi_{A,n}(\mathbf{x}_i)$ is a phasor amplitude for its n th coherent component. This representation can be used either with primary ($A = Q$), or secondary ($A = U$) sources. The phasor amplitudes $\phi_{A,n}(\mathbf{x}_i)$ and the complex expansion coefficients λ_n in Eq. (69) are eigenfunctions and eigenvalues of $W_A(\mathbf{x}_1, \mathbf{x}_2)$, as given by the equation

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} W_A(\mathbf{x}_1, \mathbf{x}_2) \phi_{A,n}(\mathbf{x}_2) d^3 \mathbf{x}_2 = \lambda_n(\omega) \phi_{A,n}(\mathbf{x}_1) \quad (70)$$

Since $W_A(\mathbf{x}_1, \mathbf{x}_2)$ is hermitian, the eigenfunctions are complete and orthogonal, i.e.,

$$\sum_n \phi_{A,n}(\mathbf{x}_1) \phi_{A,n}^*(\mathbf{x}_2) = \delta^3(\mathbf{x}_1 - \mathbf{x}_2) \quad (71)$$

and

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \phi_{A,n}(\mathbf{x}_1) \phi_{A,m}^*(\mathbf{x}_1) d^3 \mathbf{x}_1 = \delta_{nm} \quad (72)$$

where δ_{nm} is the Kronecker delta function, and the eigenvalues are real and nonnegative definite, i.e.,

$$\begin{aligned} \text{Real} \{ \lambda_n(\omega) \} &\geq 0 \\ \text{Imag} \{ \lambda_n(\omega) \} &= 0 \end{aligned} \quad (73)$$

This is a very important concept in coherence theory. Wolf has used this representation to show that the frequency decomposition of the field can be defined in a different manner than was done in Eq. (3). A new phasor,

$$\bar{U}_A(\mathbf{x}, \omega) = \sum_n a_n(\omega) \phi_{A,n}(\mathbf{x}, \omega) \quad (74)$$

where $a_n(\omega)$ are random coefficients such that

$$\langle a_n(\omega)a_m^*(\omega) \rangle = \lambda_n(\omega)\delta_{nm} \quad (75)$$

can be introduced to represent the spectral components of the field. It then follows⁶⁶ that, if the $\phi_{A,n}(\mathbf{x}_i)$ are eigenfunctions of the cross-spectral density function, $W_A(\mathbf{x}_1, \mathbf{x}_2)$ then the cross-spectral density function can be represented as the correlation function between these phasors at the two spatial points \mathbf{x}_1 and \mathbf{x}_2 , i.e.,

$$W_A(\mathbf{x}_1, \mathbf{x}_2) = \langle \bar{U}_A(\mathbf{x}_1, \omega)\bar{U}_A^*(\mathbf{x}_2, \omega) \rangle \quad (76)$$

in a manner very similar to the representation given in Eq. (11) in respect to the older phasors. Notice, by comparison of Eqs. (11) and (76), that the phasors $\bar{U}_A(\mathbf{x}_i, \omega)$ and $U_A(\mathbf{x}_i, \omega)$ are not the same. One may formulate coherence theory either by defining $u(\mathbf{x}, t)$ in Fig. 5 and then moving in a counterclockwise direction in this figure to derive the correlation functions using the Wiener-Khinchene theory or, alternatively, defining $\bar{U}_A(\mathbf{x}_i, \omega)$ and moving in a clockwise direction to define the correlation functions using the coherent mode expansion.

Wolf Shift and Scaling law

From Eqs. (66) and (67) it is clear that the spectrum of a radiation field may change as the field propagates. It is not necessarily equal to the source spectrum as it is usually assumed. This brings up an important question. Why do the spectra for light fields appear so constant, experimentally, that a change with propagation was never suspected? Wolf⁶⁵ has provided at least part of the answer to this question. He discovered a scaling law that is obeyed by most natural fields and under which the normalized spectrums for most fields do remain invariant as they propagate. We can derive this scaling law by substitution from Eqs. (66) and (67) into Eq. (16). We then find that, if the complex degree of spectral coherence for the source is a function of $k\mathbf{x}_-$ only, i.e.,

$$\mu_A(\mathbf{x}_-, \omega) = f(k\mathbf{x}_-) \quad (77)$$

(so that this function is the same for each frequency component of the field, provided that the spatial separations of the two test points are always scaled by the wavelength), then the normalized spectrum in the far field is given by

$$s_U^{(\infty)}(R\mathbf{s}, \omega) = \frac{S_Q(\omega)\omega^3}{\int_0^\infty S_Q(\omega)\omega^3 d\omega} \neq f(R\mathbf{s}) \quad (78)$$

if $S_Q(\omega)$ is the complex degree of spectral coherence for a primary, quasi-homogeneous source [see Additional Reading, Ref. 3, eq. (65)], and

$$s_U^{(\infty)}(R\mathbf{s}, \omega) = \frac{S_U^{(0)}(\omega)}{\int_0^\infty S_U^{(0)}(\omega) d\omega} = s_U^{(0)}(\omega) \neq f(R\mathbf{s}) \quad (79)$$

if $S_U^{(0)}(\omega)$ is the normalized spectrum for a secondary, quasi-homogeneous source [see Additional Reading, Ref. 3, eq. (51)]. In each case the field spectrum does not change as the field propagates. Since the cross-spectral density function for a thermal source [see Eq. (43)] obeys this scaling law, it is not surprising that these changes in the spectrum of a propagating light field were never discovered experimentally. The fact that the spectrum can change was verified experimentally only after coherence theory pointed out the possibility.⁶⁸⁻⁷²

5.8 POLARIZATION EFFECTS

Explicit Vector Representations

As discussed earlier under “Scalar Field Amplitude,” the scalar amplitude is frequently all that is required to treat the vector electromagnetic field. If polarization effects are important, it might be necessary to use two scalar field amplitudes to represent the two independent polarization components. However, in some complicated problems it is necessary to consider all six components of the vector field explicitly. For such a theory, the correlation functions between vector components of the field become tensors,^{73,74} which propagate in a manner very similar to the scalar correlation functions.

5.9 APPLICATIONS

Speckle

If coherent light is scattered from a stationary, rough surface, the phase of the light field is randomized in space. The diffraction patterns observed with such light displays a complicated granular pattern usually called *speckle* (see Ref. 8, sec. 7.5). Even though the light phase can be treated as a random variable, the light is still perfectly coherent. Coherence theory deals with the effects of time fluctuations, not spatial variations in the field amplitude or phase. Despite this, the same statistical tools used in coherence theory have been usefully applied to studying speckle phenomena.^{75,76} To treat speckle, the ensemble is usually redefined, not to represent the time fluctuations of the field but, rather, to represent all of the possible speckle patterns that might be observed under the conditions of a particular experiment. An observed speckle pattern is usually due to a single member of this ensemble (unless time fluctuations are also present), whereas the intensity observed in coherence theory is always the result of a weighted average over all of the ensemble. To obtain the intensity distribution over some plane, as defined in coherence theory, it would be necessary to average over all of the possible speckle patterns explicitly. If this is done, for example, by moving the scatterer while the intensity of a diffraction pattern is time-averaged, then time fluctuations are introduced into the field during the measurement; the light becomes partially coherent; and coherence theory can be properly used to model the measured intensity. One must be very careful in applying the coherence theory model to treat speckle phenomena, because coherence theory was not originally formulated to deal with speckle.

Statistical Radiometry

Classical radiometry was originally based on a mechanical treatment of light as a flux of particles. It is not totally compatible with wave theories. Coherence theory has been used to incorporate classical radiometry into electromagnetic theory as much as has been found possible. It has been found that the usual definitions for the radiance function and the radiant emittance cause problems when applied to a wave theory. Other radiometric functions, such as the radiant intensity, have clear meaning in a wave theory.

Spectral Representation

It was discovered, using coherence theory, that the spectrum of a light field is not the same as that of its source and that it can change as the light field propagates away from its source into the radiation field. Some of this work was discussed earlier under “Perfectly Incoherent Sources” and “Coherent Mode Representation.” This work has been found very useful for explaining troublesome experimental discrepancies in precise spectroradiometry.⁷¹

Laser Modes

Coherence theory has been usefully applied to describing the coherence properties of laser modes.⁷⁷ This theory is based on the coherent mode representation discussed under “Spectrum of Light.”

Radio Astronomy

Intensity interferometry was used to apply techniques from radio astronomy to observation with light.³⁷ A lively debate ensued as to whether optical interference effects (which indicate partial coherence) could be observed from intensity correlations.⁷⁸ From relations like Eq. (26) and similar calculations using quantum coherence theory, it quickly became clear that they could.^{79,80} More recently, coherence theory has been used to model a radio telescope⁴² and to study how to focus an instrument to observe emitters that are not in the far field of the antenna array.⁸¹ It has been shown⁸² that a radio telescope and a conventional optical telescope are very similar, within a coherence theory model, even though their operation is completely different. This model makes the similarities between the two types of instruments very clear.

Noncosmological Red Shift

Cosmological theories for the structure and origin of the universe make great use of the observed red shift in the spectral lines of the light received from distant radiating objects, such as stars.⁸³ It is usually assumed that the spectrum of the light is constant upon propagation and that the observed red shift is the result of simple Doppler shift due to the motion of the distant objects away from the earth in all directions. If this is true, then clearly the observed universe is expanding and must have begun with some sort of explosion, called “the big bang.” The size of the observable universe is estimated based on the amount of this red shift. A new theory by Wolf^{84–87} shows that red shifts can occur naturally without Doppler shifts as the light propagates from the source to an observer if the source is not in thermal equilibrium, i.e., a thermal source as discussed earlier under “Thermal (Lambertian) Source.” The basis of Wolf’s theory was discussed in this paper.^{84,85}

5.10 REFERENCES

1. L. Mandel and E. Wolf, “Coherence Properties of Optical Fields,” *Rev. Mod. Phys.* **37**, 1965, pp. 231–287.
2. L. Mandel and E. Wolf, *Selected Papers on Coherence and Fluctuations of Light*, vol. I, Dover, New York, 1970. Reprinted in Milestone Series, SPIE press, Bellingham, Wash., 1990.
3. L. Mandel and E. Wolf, *Selected Papers on Coherence and Fluctuations of Light*, vol. II, Dover, New York, 1970. Reprinted in Milestone Series, SPIE press, Bellingham, Wash., 1990.
4. M. J. Beran and G. B. Parrent, *Theory of Partial Coherence*, Prentice-Hall, Engelwood Cliffs, N.J., 1964.
5. B. Crosignani and P. Di Porto, *Statistical Properties of Scattered Light*, Academic Press, New York, 1975.
6. A. S. Marathay, *Elements of Optical Coherence Theory*, Wiley, New York, 1982.
7. J. Perina, *Coherence of Light*, Reidel, Boston, 1985.
8. J. W. Goodman, *Statistical Optics*, Wiley, New York, 1985.
9. R. J. Glauber, “The Quantum Theory of Optical Coherence,” *Phys. Rev.* **130**, 1963, pp. 2529–2539.
10. R. J. Glauber, “Coherent and Incoherent States of the Radiation Field,” *Phys. Rev.* **131**, 1963, pp. 2766–2788.
11. C. L. Mehta and E. C. G. Sudarshan, “Relation between Quantum and Semiclassical Description of Optical Coherence,” *Phys. Rev.* **138**, B274–B280, 1965.
12. W. Louisel, *Radiation and Noise in Quantum Electronics*, McGraw-Hill, New York, 1964.
13. W. H. Louisell, *Quantum Statistical Properties of Radiation*, Wiley, New York, 1973.

14. M. Sargent, M. O. Scully, and W. E. Lamb, *Laser Physics*, Addison-Wesley, Reading, MA, 1974.
15. M. Born and E. Wolf, *Principles of Optics*, Pergamon Press, New York, 1980.
16. H. M. Nussenvieg, *Causality and Dispersion Relations*, Academic Press, New York, 1972.
17. A. M. Yaglom, *An Introduction to the Theory of Stationary Random Functions*, Prentice-Hall, Engelwood Cliffs, N.J., 1962.
18. G. Borgiotti, "Fourier Transforms Method in Aperture Antennas Problems," *Aha Frequenza* **32**, 1963, pp. 196–205.
19. D. R. Rhodes, "On the Stored Energy of Planar Apertures," *IEEE Trans.* **AP14**, 1966, pp. 676–683.
20. W. H. Carter, "The Electromagnetic Field of a Gaussian Beam with an Elliptical Cross-section," *J. Opt. Soc. Am.* **62**, 1972, pp. 1195–1201.
21. W. H. Carter, "Electromagnetic Beam Fields," *Optica Acta* **21**, 1974, pp. 871–892.
22. B. J. Thompson and E. Wolf, "Two-Beam Interference with Partially Coherent Light," *J. Opt. Soc. Am.* **47**, 1957, pp. 895–902.
23. W. H. Carter, "Measurement of Second-order Coherence in a Light Beam Using a Microscope and a Grating," *Appl. Opt.* **16**, 1977, pp. 558–563.
24. E. Wolf, "A Macroscopic Theory of Interference and Diffraction of Light from Finite Sources II: Fields with Spatial Range and Arbitrary Width," *Proc. Roy. Soc.* **A230**, 1955, pp. 246–265.
25. G. B. Parrent, "On the Propagation of Mutual Coherence," *J. Opt. Soc. Am.* **49**, 1959, pp. 787–793.
26. W. B. Davenport and W. L. Root, *An Introduction to the Theory of Random Signals and Noise*, McGraw-Hill, New York, 1960.
27. L. Mandel and E. Wolf, "Spectral Coherence and the Concept of Cross-spectral Purity," *J. Opt. Soc. Am.* **66**, 1976, pp. 529–535.
28. E. W. Marchand and E. Wolf, "Angular Correlation and the Far-zone Behavior of Partially Coherent Fields," *J. Opt. Soc. Am.* **62**, 1972, pp. 379–385.
29. E. W. Marchand and E. Wolf, "Radiometry with Sources of Any State of Coherence," *J. Opt. Soc. Am.* **64**, 1974, pp. 1219–1226.
30. E. Wolf, "The Radiant Intensity from Planar Sources of any State of Coherence," *J. Opt. Soc. Am.* **68**, 1978, pp. 1597–1605.
31. A. T. Friberg, "On the Existence of a Radiance Function for Finite Planar Sources of Arbitrary States of Coherence," *J. Opt. Soc. Am.* **69**, 1979, pp. 192–198.
32. E. Wigner, "Quantum Correction for Thermodynamic Equilibrium" *Phys. Rev.* **40**, 1932, pp. 749–760.
33. K. Imre, E. Ozizmir, M. Rosenbaum, and P. F. Zweifer, "Wigner Method in Quantum Statistical Mechanics," *J. Math. Phys.* **8**, 1967, pp. 1097–1108.
34. A. Papoulis, "Ambiguity Function in Fourier Optics," *J. Opt. Soc. Am.* **64**, 1974, pp. 779–788.
35. A. Walther, "Radiometry and Coherence," *J. Opt. Soc. Am.* **63**, 1973, pp. 1622–1623.
36. J. T. Foley and E. Wolf, "Radiometry as a Short Wavelength Limit of Statistical Wave Theory with Globally Incoherent Sources," *Opt. Commun.* **55**, 1985, pp. 236–241.
37. R. H. Brown and R. Q. Twiss, "Correlation between Photons in Two Coherent Beams of Light," *Nature* **177**, 1956, pp. 27–29.
38. J. D. Jackson, *Classical Electrodynamics*, Wiley, New York, 1962.
39. C. L. Mehta and A. B. Balachandran, "Some Theorems on the Unimodular Complex Degree of Optical Coherence," *J. Math. Phys.* **7**, 1966, pp. 133–138.
40. L. Mandel and E. Wolf, "Complete Coherence in the Space Frequency Domain," *Opt. Commun.* **36**, 1981, pp. 247–249.
41. W. H. Carter, "Difference in the Definitions of Coherence in the Space-Time Domain and in the Space-Frequency Domain," *J. Modern Opt.* **39**, 1992, pp. 1461–1470.
42. W. H. Carter and L. E. Somers, "Coherence Theory of a Radio Telescope," *IEEE Trans.* **AP24**, 1976, pp. 815–819.
43. A. C. Schell, *The Multiple Plate Antenna*, doctoral dissertation, M.I.T., 1961.
44. A. C. Schell, "A Technique for the Determination of the Radiation Pattern of a Partially Coherent Aperture," *IEEE Trans.* **AP-15**, 1967, p. 187.

45. W. H. Carter and E. Wolf, "Correlation Theory of Wavefields Generated by Fluctuating Three-dimensional, Primary, Scalar Sources: II. Radiation from Isotropic Model Sources," *Optica Acta* **28**, 1981, pp. 245–259.
46. W. H. Carter and E. Wolf, "Correlation Theory of Wavefields Generated by Fluctuating Three-dimensional, Primary, Scalar Sources: I. General Theory," *Optica Acta* **28**, 1981, pp. 227–244.
47. E. Collett and E. Wolf, "Is Complete Coherence Necessary for the Generation of Highly Directional Light Beams?" *Opt. Letters* **2**, 1978, pp. 27–29.
48. W. H. Carter and Bertolotti, "An Analysis of the Far Field Coherence and Radiant Intensity of Light Scattered from Liquid Crystals," *J. Opt. Soc. Am.* **68**, 1978, pp. 329–333.
49. W. H. Carter and E. Wolf, "Coherence and Radiometry with Quasi-homogeneous Planar Sources," *J. Opt. Soc. Am.* **67**, 1977, pp. 785–796.
50. W. H. Carter and E. Wolf, "Inverse Problem with Quasi-Homogeneous Sources," *J. Opt. Soc. Am.* **2A**, 1985, pp. 1994–2000.
51. E. Wolf and W. H. Carter, "Coherence and Radiant Intensity in Scalar Wave Fields Generated by Fluctuating Primary Planar Sources," *J. Opt. Soc. Am.* **68**, 1978, pp. 953–964.
52. W. H. Carter, "Coherence Properties of Some Isotropic Sources," *J. Opt. Soc. Am.* **A1**, 1984, pp. 716–722.
53. W. H. Carter and E. Wolf, "Scattering from Quasi-homogeneous Media," *Opt. Commun.* **67**, 1988, pp. 85–90.
54. W. H. Carter, "Scattering from Quasi-homogeneous Time Fluctuating Random Media," *Opt. Commun.* **77**, 1990, pp. 121–125.
55. W. H. Carter and E. Wolf, "Coherence Properties of Lambertian and Non-Lambertian Sources," *J. Opt. Soc. Am.* **65**, 1975, pp. 1067–1071.
56. J. T. Foley, W. H. Carter, and E. Wolf, "Field Correlations within a Completely Incoherent Primary Spherical Source," *J. Opt. Soc. Am.* **A3**, 1986, pp. 1090–1096.
57. H. H. Hopkins, "The Concept of Partial Coherence in Optics," *Proc. Roy. Soc.* **A208**, 1951, pp. 263–277.
58. W. H. Carter, "Generalization of Hopkins' Formula to a Class of Sources with Arbitrary Coherence," *J. Opt. Soc. Am.* **A2**, 1985, pp. 164–166.
59. P. H. van Cittert, "Die wahrscheinliche Schwingungsverteilung in einer von einer Lichtquelle direkt oder mittels einer Linse beleuchteten Ebene," *Physica* **1**, 1934, pp. 201–210.
60. F. Zernike, "The Concept of Degree of Coherence and Its Application to Optical Problems," *Physica* **5**, 1938, pp. 785–795.
61. W. H. Carter, "Three Different Kinds of Fraunhofer Approximations. I. Propagation of the Field Amplitude," *Radio Science* **23**, 1988, pp. 1085–1093.
62. W. H. Carter, "Three Different Kinds of Fraunhofer Approximations. II. Propagation of the Cross-spectral Density Function," *J. Mod. Opt.* **37**, 1990, pp. 109–120.
63. W. H. Carter, "Scattering of a Light Beam from an Air-Sea Interface," *Applied Opt.* **32**, 3286–3294 (1993).
64. A. T. Friberg and E. Wolf, "Reciprocity Relations with Partially Coherent Sources," *Opt. Acta* **36**, 1983, pp. 417–435.
65. E. Wolf, "Invariance of the Spectrum of Light on Propagation," *Phys. Rev. Letters* **56**, 1986, pp. 1370–1372.
66. E. Wolf, "New Theory of Partial Coherence in the Space-frequency Domain. Part I: Spectra and Cross Spectra of Steady-state Sources," *J. Opt. Soc. Am.* **72**, 1982, pp. 343–351.
67. E. Wolf, "New Theory of Partial Coherence in the Space-frequency Domain. Part II: Steady-state Fields and Higher-order Correlations," *J. Opt. Soc. Am.* **A3**, 1986, pp. 76–85.
68. G. M. Morris and D. Faklis, "Effects of Source Correlation on the Spectrum of Light," *Opt. Commun.* **62**, 1987, pp. 5–11.
69. M. F. Bocko, D. H. Douglass, and R. S. Knox, "Observation of Frequency Shifts of Spectral Lines Due to Source Correlations," *Phys. Rev. Letters* **58**, 1987, pp. 2649–2651.
70. F. Gori, G. Guattari, and C. Palma, "Observation of Optical Redshift and Blueshifts Produced by Source Correlations," *Opt. Commun.* **67**, 1988, pp. 1–4.
71. H. C. Kandpal, J. S. Vashiya, and K. C. Joshi, "Wolf Shift and Its Application in Spectro-radiometry," *Opt. Commun.* **73**, 1989, pp. 169–172.
72. G. Ingedouw, "Synthesis of Polychromatic Light Sources with Arbitrary Degree of Coherence: From Experiments," *J. Mod. Opt.* **36**, 1989, pp. 251–259.

73. W. H. Carter, "Properties of Electromagnetic Radiation from a Partially Correlated Current Distribution," *J. Opt. Soc. Am.* **70**, 1980, pp. 1067–1074.
74. W. H. Carter and E. Wolf, "Far-zone Behavior of Electromagnetic Fields Generated by Fluctuating Current Distributions," *Phys. Rev.* **36A**, 1987, pp. 1258–1269.
75. J. C. Dainty, "The Statistics of Speckle Patterns," in E. E. Wolf, *Progress in Optics*, North Holland, Amsterdam, 1976.
76. J. C. Dainty, *Laser Speckle and Related Phenomena*, Springer Verlag, Berlin, 1975.
77. E. Wolf and G. S. Agarwal, "Coherence Theory of Laser Resonator Modes," *J. Opt. Soc. Am.* **A1**, 1984, pp. 541–546.
78. E. M. Purcell, "The Question of Correlation Between Photons in Coherent Light Rays," *Nature* **178**, 1956, pp. 1449–1450.
79. R. H. Brown and R. Q. Twiss, "Interferometry of the Intensity Fluctuations in Light: I. Basic Theory: the Correlation between Photons in Coherent Beams of Radiation," *Proc. Roy. Soc.* **242**, 1957, pp. 300–324.
80. R. H. Brown and R. Q. Twiss, "Interferometry of the Intensity Fluctuations in Light: II. An Experimental Test of the Theory for Partially Coherent Light," *Proc. Roy. Soc.* **243**, 1957, pp. 291–319.
81. W. H. Carter, "Refocusing a Radio Telescope to Image Sources within the Near Field of the Antenna Array," NRL Report no. 9141, August 17, 1988.
82. W. H. Carter, "A New Theory for Image Formation with Quasi-homogeneous Sources and Its Application to the Detection and Location of ICBM Launches," NRL Report no. 9336 (September 17, 1991).
83. J. V. Narliker, "Noncosmological Redshifts," *Space Sci. Revs.* **50**, 1989, pp. 523–614.
84. E. Wolf, "Non-cosmological Redshifts of Spectral Lines," *Nature* **326**, 1987, pp. 363–365.
85. E. Wolf, "Red Shifts and Blue Shifts of Spectral Lines Emitted by Two Correlated Sources," *Phys. Rev. Letters* **58**, 1987, pp. 2646–2648.
86. A. Gamliel and E. Wolf, "Spectral Modulation by Control of Source Correlations," *Opt. Commun.* **65**, 1988, pp. 91–96.
87. E. Wolf, J. T. Foley, and F. Gori, "Frequency Shifts of Spectral Lines Produced by Scattering from Spatially Random Media," *J. Opt. Soc. Am.* **A6**, 1989, pp. 1142–1149.

5.11 ADDITIONAL READING

1. W. H. Carter, "Difference in the Definition of Coherence in the Space-Time Domain and in the Space-Frequency Domain," *J. Mod. Opt.* **39**, 1992, pp. 1461–1470.
2. W. H. Carter, "Coherence Properties of Some Broadband Highly Coherent Light Fields," *J. Opt. Soc. Am.* **A10**, 1993, pp. 1570–1578.
3. W. H. Carter, "A Study of the Propagation of Optical Spectra Using a Wavelength Independent Formulation of Diffraction Theory and Coherence Theory," *J. Mod. Opt.* **40**, 1993, pp. 2433–2449.

COHERENCE THEORY: TOOLS AND APPLICATIONS

Gisele Bennett

*Electro-Optical Systems Laboratory and School of Electrical and Computer Engineering
Georgia Tech Research Institute
Georgia Institute of Technology
Atlanta, Georgia*

William T. Rhodes

*School of Electrical and Computer Engineering
Georgia Institute of Technology, and
Department of Electrical Engineering and Imaging Technology Center
Florida Atlantic University
Boca Raton, Florida*

J. Christopher James

*Electro-Optical Systems Laboratory
Georgia Tech Research Institute
Georgia Institute of Technology
Atlanta, Georgia*

6.1 GLOSSARY

\mathcal{F}	Fourier transform operator
I	Intensity
k	wave number
u	complex field amplitude
Γ_{12}	mutual coherence function
J_{12}	mutual intensity function
μ_{12}	complex degree of coherence
\mathbf{x}	spatial vector

6.2 INTRODUCTION

The formalisms of coherence theory have been extensively developed over the past century, with roots extending back much further. An understanding of key concepts of the coherence properties of light waves and of their accompanying mathematical models can be exploited in specific applications to extract information about objects (e.g., stellar interferometry, microscopy^{1,2}), to encrypt signals transmitted by optical fiber (coherence modulation), to explain puzzling phenomena observed in nature (sunlight-produced speckle, enhanced backscatter), and to prevent significant errors in the measurement of optical quantities (measurement of the reflectivity of diffuse surfaces), among other examples. There are cases where using a coherence-theory-based analysis provides not only insight into understanding system behavior but simplifies the analysis through reduced computational steps.

In Chap. 5, Carter provides a compendium of basic definitions from coherence theory, introduces mathematical models for light sources of fundamental importance, discusses how the coherence properties of optical wave fields evolve under propagation, and shows how the spectral properties of light waves relate to their coherence properties.³ He also notes briefly several areas in which coherence theory has been exploited to advance significantly the understanding of optical phenomena. References 1-8 in Chap. 5 provide a suitable introduction to most aspects of classical (i.e., nonquantum) coherence theory. Since the publication of Vol. I of this *Handbook* series, Mandel and Wolf have written a comprehensive review of coherence theory and a brief discussion of specific applications.⁴ Other informative sources include a compilation of early papers on coherence edited by Mandel and Wolf.⁵

The formal structures of coherence theory can be intimidating to the nonexpert trying to ascertain, for example, the effect of a particular light source on a given optical system. One has only to review the myriad of defined quantities or look at a six- or seven-dimensional integral expression describing an image-forming system to understand why. Nevertheless, the basic concepts of coherence theory, and certain of their key formulas, can be helpful even to the nonexpert when they are applied to the understanding of a variety of optical systems and phenomena. That said, it should be noted that sometimes the application of coherence theory formalisms only serves to obfuscate the operation of a given system. One objective of this chapter is to provide guidance as to when coherence theory can help and when it might serve largely to confuse. We note that coherence theory is not generally needed if the optical wave fields of concern can be modeled as being monochromatic or quasi-monochromatic. The subject of interferometry is easily treated without the inclusion of coherence theory if the source of light is a laser operating in a single mode, and optical coherence tomography is easily understood without appeal to much more than an understanding of temporal coherence theory.⁶⁻¹⁰ It is with this latter point in mind that we emphasize on analyses involving the spatial coherence properties of the light.

6.3 KEY DEFINITIONS AND RELATIONSHIPS

We begin with the presentation of a set of tools that are critical to the analysis of certain applications. Although the cross-spectral density function^{3,4} is often preferred in the modeling of the spatio-temporal coherence properties of light, in a discussion of the applications of coherence theory where physical understanding is important we find it preferable to work instead with the mutual coherence function, defined by¹¹

$$\Gamma_{12}(\tau) = \langle u(\mathbf{x}_1, (t+\tau))u^*(\mathbf{x}_2, t) \rangle \quad (1)$$

where $u(\mathbf{x}, t)$ denotes the complex analytic signal associated with the scalar wave amplitude at position $\mathbf{x} = (x, y, z)$ and where the angular brackets $\langle \cdot \rangle$ denote a suitable time average.*

*The nature of the time average is discussed later in this chapter. Note that in Chap. 5 Carter uses angular brackets to denote ensemble averages rather than time averages.

This function has the advantage of relating directly to quantities that can be measured by an easily visualized Young's two-pinhole interferometer. In theoretical developments $u(\mathbf{x}, t)$ is usually treated as an ergodic random process, and the time average is replaced with an ensemble average. However, because such conditions are not satisfied in all applications of interest, we maintain specific reference to time averages. The mutual intensity function, unlike the complex amplitude itself, is a measurable quantity. The temporal frequency bandwidth of optical signals typically far exceeds the bandwidth of any optical detector, thus preventing the direct measurement of $u(\mathbf{x}, t)$. As discussed by Greivencamp, a Young's two-pinhole interferometer can be used, at least conceptually, in the measurement process.¹²

In many cases we can assume the light to be quasi-monochromatic.¹³ The narrowband condition $\Delta\lambda \ll \lambda$ is satisfied, where $\Delta\lambda$ denotes the spectral bandwidth of the light and where λ is the mean wavelength, and the coherence length of the light,³ $\Delta l \sim \lambda^2/\Delta\lambda$, is much greater than the maximum path length difference encountered in the passage of light from a source of interest to the measurement plane of concern. Under such conditions¹¹ the mutual coherence function is given by

$$\Gamma_{12}(\tau) = J_{12} e^{-i\omega\tau} \quad (2)$$

where ω is the mean angular frequency of the radiation and where J_{12} , the mutual intensity function, is given by

$$J_{12} = \Gamma_{12}(0) \quad (3)$$

Note that J_{12} depends on the spatial coordinates \mathbf{x}_1 and \mathbf{x}_2 but not on τ .

Effect of Transmissive Planar Object

In most of the applications considered in this chapter we are concerned with the behavior of J_{12} in a two-dimensional systems framework. Assume, therefore, that \mathbf{x}_1 and \mathbf{x}_2 represent the coordinates of a pair of points in a plane of constant z . The complex amplitude $u(\mathbf{x}, t)$ can then be written as $u(x, y, t)$, and J_{12} can be written in the form $J(x_1, y_1; x_2, y_2)$. If the wave field $u_{\text{inc}}(x, y, t)$ is incident upon a thin transmissive object of complex transmittance $t(x, y)$, the transmitted wave $u_{\text{trans}}(x, y, t)$ has amplitude $u_{\text{inc}}(x, y, t)t(x, y)$, and it is easily shown that the corresponding mutual intensity is given by

$$J_{\text{trans}}(x_1, y_1; x_2, y_2) = J_{\text{inc}}(x_1, y_1; x_2, y_2) t(x_1, y_1) t^*(x_2, y_2) \quad (4)$$

Of particular interest later is the case when the object is a thin spherical lens. The transmittance function of such a lens is given, in the paraxial approximation and ignoring an accompanying pupil function, by

$$t_{\text{lens}}(x, y) = \exp\left[-i\frac{k}{2f}(x^2 + y^2)\right] \quad (5)$$

where f denotes the focal length of the lens and where $k = 2\pi/\lambda$. Assuming a lens that is not the limiting aperture in a system and substitution in Eq. (4) yields

$$J_{\text{trans}}(x_1, y_1; x_2, y_2) = J_{\text{inc}}(x_1, y_1; x_2, y_2) \exp\left\{-i\frac{k}{2f}\left[(x_1^2 - x_2^2) + (y_1^2 - y_2^2)\right]\right\} \quad (6)$$

Effect of a General Linear System

If the wave amplitude $u(x, y, t)$ is input to a space-invariant two-dimensional linear optical system with spatial impulse response $h(x, y)$, the output wave amplitude has the form

$u_{\text{out}}(x, y; t) = u_{\text{in}}(x, y; t) \otimes \otimes h(x, y)$, where $\otimes \otimes$ denotes a two-dimensional convolution operation. The relationship for the corresponding mutual intensity functions is given by*

$$J_{\text{out}}(x_1, y_1; x_2, y_2) = J_{\text{in}}(x_1, y_1; x_2, y_2) \otimes \otimes h(x_1, y_1) ** h^*(x_2, y_2) \quad (7)$$

The cascade of two-dimensional convolutions in this equation can, if desired, be rewritten in the form of a four-dimensional convolution of $J_{\text{in}}(x_1, y_1; x_2, y_2)$ with a separable four-dimensional kernel:

$$J_{\text{out}}(x_1, y_1; x_2, y_2) = J_{\text{in}}(x_1, y_1; x_2, y_2) \otimes \otimes \otimes \otimes [h(x_1, y_1) h^*(x_2, y_2)] \quad (8)$$

Propagation of the Mutual Intensity

It is a straightforward matter to show how the mutual intensity of a wave field propagates.^{14–16} In most applications, concern is with propagation in the Fresnel (paraxial) regime, in which case $h(x, y)$ in Eq. (7), denoted now by $h_z(x, y)$, is given by

$$h_z(x, y) = \frac{\exp(ikz)}{i\lambda z} \exp\left[\frac{ikz(x^2 + y^2)}{2z}\right] \quad (9)$$

where z is the plane-to-plane propagation distance. The factor $(-i)^{-1} \exp(ikz)$ can often be ignored.

If we assume the source distribution to be planar, quasi-monochromatic, and spatially incoherent [the latter condition implying that $J_{12} = \gamma_o I(\mathbf{x}_1) \delta(\mathbf{x}_1 - \mathbf{x}_2)$, where I denotes the optical intensity of the wave field and where γ_o is a constant that depends on characteristics of the source at the scale of a wavelength] with optical intensity $I_o(x, y)$, it can be shown by combining Eqs. (8) and (9) that the mutual intensity a distance z away from the source is given by

$$\begin{aligned} J_z(x_1, y_1; x_2, y_2) &= \frac{\exp(-i\phi)}{(\lambda z)^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_o(\xi, \eta) \exp\left[\frac{i2\pi}{\lambda z}(\Delta x \xi + \Delta y \eta)\right] d\xi d\eta \\ &= \frac{\exp(-i\phi)}{(\lambda z)^2} \hat{I}_o\left(\frac{\Delta x}{\lambda z}, \frac{\Delta y}{\lambda z}\right) \end{aligned} \quad (10)$$

where $\hat{I}_o(x, y)$ denotes the two-dimensional Fourier transform of the source-plane intensity, $\Delta x = x_1 - x_2$, $\Delta y = y_1 - y_2$, and where

$$\phi = \frac{\pi}{\lambda z} [(x_2^2 + y_2^2) - (x_1^2 + y_1^2)] \quad (11)$$

Equation (10) is a compact statement of the van Cittert-Zernike theorem. If the wave field corresponding to J_z in Eq. (10) is incident on a thin spherical lens of focal length $f = z$, then, through Eqs. (10) and (4), the quadratic phase factors of Eq. (11) are removed and the resulting mutual intensity has the simple form $(1/\lambda z)^2 \hat{I}_o(\Delta x/\lambda z, \Delta y/\lambda z)$. This distribution, it is noted, is a function of Δx and Δy alone, that is, it is a function only of the vector separation between the two sample points \mathbf{x}_1 and \mathbf{x}_2 .

The Issue of Time Averages

The nature of the time average implied by the angular brackets in Eq. (1) will depend on the particular situation of concern. If, for example, $u(x, t)$ corresponds to the wave field produced by a white light source, adequate averaging is obtained over a time interval of roughly a picosecond (10^{-12} s). If, on the other hand, the light source is a stabilized laser with a temporal frequency bandwidth of several tens of kilohertz, the period over which the average is evaluated might be milliseconds in duration. Often the bandwidth of a detector determines the appropriate duration of the averaging process. If we are viewing white-light interference fringes by eye, the integration interval can be anywhere

* Although this expression does not appear in the standard textbooks on coherence theory, it is easily derived.

from a picosecond to perhaps a thirtieth of a second. If the same fringes are being observed with a megahertz-bandwidth measurement system, on the other hand, integration for times greater than a microsecond is inappropriate.

6.4 PROPAGATION, DIFFRACTION, AND SCATTERING: ENHANCED BACKSCATTER AND THE LAU EFFECT

Coherence theory can provide a useful tool for describing certain phenomena encountered in connection with wave field propagation, scattering, and diffraction by gratings. Two examples of cases where a coherence theory approach is especially enlightening are the phenomenon of enhanced backscatter (EBS) and the Lau effect, as discussed in the following subsections.

Enhanced Backscatter

Enhanced backscatter is observed when a laser beam, after passing through a fine-grained moving diffuser, is reflected back through the same moving diffuser. As illustrated in Fig. 1a, the far field contains, in addition to a broad background irradiance distribution, a tightly focused spot of light, corresponding to the diffraction-limited focusing of the partially recollimated incident beam. This result contrasts with that obtained with two *different* moving diffusers, as shown in Fig. 1b, in which case there is no focused light spot. Although a ray-optics model can provide some insight

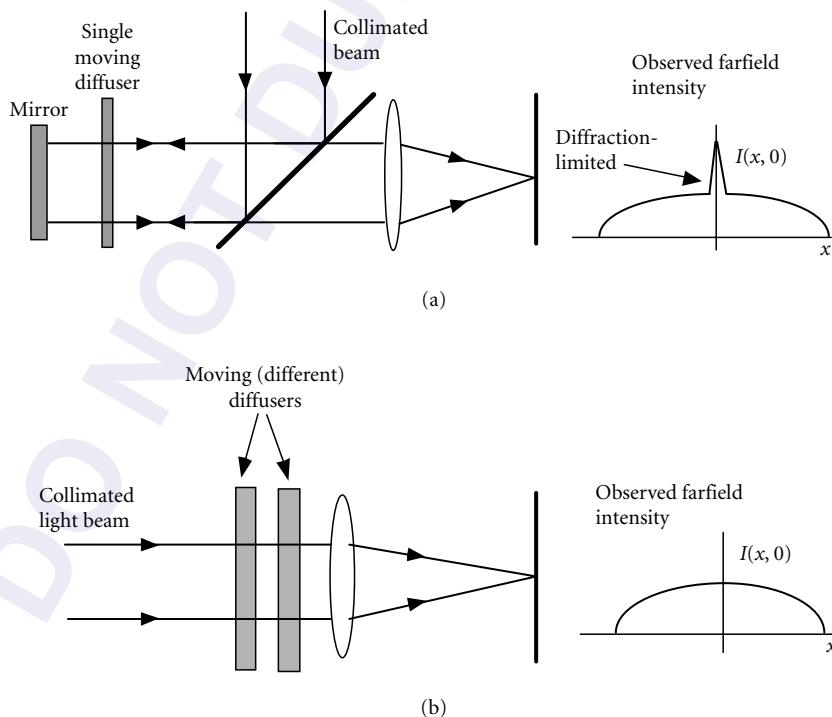


FIGURE 1 Scattering by moving diffusers: (a) enhanced backscatter (partial phase conjugation scattering) produced by double passage of light through the same diffuser and (b) conventional scattering produced by two different diffusers.

into this counter-intuitive phenomenon, coherence theory allows us to develop much greater understanding.

Consider the case where the incident wave field is quasi-monochromatic and of arbitrary spatial coherence properties. To simplify the math we work with a one-dimensional model. The moving diffuser is represented by a thin dynamic random phase screen with complex amplitude transmittance $\phi(x, t)$. Transmission of the incident wave $u(x, t)$ through the diffuser results in wave field $u(x, t)\phi(x, t)$, which propagates a distance z to be scattered a second time by an identical random phase screen. The doubly-scattered wave amplitude $u'(x, t)$ is given by¹⁷

$$u'(x, t) = \phi(x, t) \int_{-\infty}^{\infty} u(\xi, t) \phi(\xi, t) h_z(x - \xi) d\xi \quad (12)$$

where $h_z(x)$ represents the wave propagation kernel of Eq. (9). The mutual intensity of the doubly-scattered wave field is calculated as

$$\begin{aligned} J'(x_1, x_2) &= \langle u'(x_1, t) u'^*(x_2, t) \rangle \\ &= \left\langle \left[\phi(x_1, t) \int_{-\infty}^{\infty} u(\xi, t) \phi(\xi, t) h_z(x_1 - \xi) d\xi \right] \left[\phi^*(x_2, t) \int_{-\infty}^{\infty} u^*(\eta, t) \phi^*(\eta, t) h_z^*(x_2 - \eta) d\eta \right] \right\rangle \quad (13) \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \langle u(\xi, t) u^*(\eta, t) \rangle \langle \phi(x_1, t) \phi^*(x_2, t) \phi(\xi, t) \phi^*(\eta, t) \rangle h_z(x_1 - \xi) h_z^*(x_2 - \eta) d\xi d\eta \end{aligned}$$

where in expressing the effect of the time average it has been assumed that u and ϕ vary independently. To simplify this expression we note that $\langle u(\xi, t) u^*(\eta, t) \rangle$ is the mutual intensity $J(\xi, \eta)$ of the input wave field. Additionally, if $\phi(x, t)$ represents a spatially stationary, delta-correlated scatterer with suitably large excursions in phase, the second term in brackets in the integral of Eq. (13) can be modeled by

$$\langle \phi(x_1, t) \phi^*(x_2, t) \phi(\xi, t) \phi^*(\eta, t) \rangle = \gamma^2 [\delta(x_1 - x_2) \delta(\xi - \eta) + \delta(x_1 - \eta) \delta(\xi - x_2)] \quad (14)$$

where

$$\gamma = \int_{-\infty}^{\infty} \langle \phi(\xi, t) \phi^*(0, t) \rangle d\xi \quad (15)$$

Making the appropriate substitutions in Eq. (13) yields

$$J'(x_1, x_2) = a(x_1) \delta(x_1 - x_2) + b(x_1, x_2) J^*(x_1, x_2) \quad (16)$$

where

$$a(x_1) = \gamma^2 I(x_1) \otimes |h_z(x_1)|^2 \quad (17)$$

$$b(x_1, x_2) = \gamma^2 h_z(x_1 - x_2) h_z^*(x_1 - x_2) \quad (18)$$

and where the function $I(x)$ in Eq. (17) is the optical intensity of the incident wave. If the one-dimensional form of the Fresnel propagation kernel is substituted for $h_z(x)$,

$$h_z(x) = \frac{1}{\sqrt{i\lambda z}} \exp\left[\frac{i\pi}{\lambda z} x^2\right] \quad (19)$$

calculation of $a(x_1)$ and $b(x_1, x_2)$ yields

$$J'(x_1, x_2) = \delta(x_1 - x_2) \gamma^2 \int_{-\infty}^{\infty} I_{\text{inc}}(\xi) d\xi + \gamma^2 J_{\text{inc}}^*(x_1, x_2) \quad (20)$$

The doubly-scattered wave, $u'(x, t)$, can thus be thought of as consisting of two components. One component has the characteristics of a spatially incoherent wave and produces a general background glow in the far field. The second component, on the other hand, effectively replicates coherence properties of the incident wave, having mutual intensity $J_{\text{inc}}^*(x_1, x_2)$, the complex conjugate of that of the incident wave. A wave field with mutual intensity $J_{\text{inc}}^*(x_1, x_2)$ behaves like a time-reversed (or “backward-propagating”) version of a wave field with mutual intensity $J_{\text{inc}}(x_1, x_2)$. Thus, if the incident wave is diverging, the doubly-scattered wave will contain, in addition to an incoherent component, a coherence-replicated component that is converging at the same rate. For the case illustrated in Fig. 1a, since the incident beam is collimated, the coherence-replicated component of the doubly-scattered wave also behaves like a collimated light beam.

It should be noted that evaluation of the time average for the term $\langle u(\xi, t)u^*(\eta, t) \rangle$ in Eq. (13) may be satisfactorily complete in a fraction of a microsecond—that is, after an interval large compared to the reciprocal bandwidth of the light incident on the diffuser—whereas calculation of the second brackets term, $\langle \phi(x_1, t)\phi^*(x_2, t)\phi(\xi, t)\phi^*(\eta, t) \rangle$, may require milliseconds or even seconds, depending on how rapidly the random phase screen evolves with time. What is essential is that $\phi(x, t)$ evolves satisfactorily over the duration of the time average. In the terminology of random processes, we require that $\phi(x, t)$ goes through a sufficient number of realizations as to provide a good statistical average of the bracketed quantity.

The enhanced backscatter phenomenon can be exploited, at least in theory, in the imaging of diffuser-obscured objects. Let the mirror in the system of Fig. 1a be replaced by a planar object with amplitude reflectance $\mathbf{r}(x)$, assuming that the incident wave is monochromatic and planar. Through an analysis quite similar to that above, one can show that the mutual intensity of the doubly-scattered wave field $J_{\text{ds}}(x_1, x_2)$ again contains two components, a coherent one and a spatially incoherent one. The coherent component, which can be measured by interferometric means, is proportional to the modulus of the Fresnel transform of the object reflectance function¹⁸

$$J_{\text{dscoh}}(x_1, x_2) \propto \kappa^2 \left| \tilde{\mathbf{r}}_{z/2} \left(\frac{x_1 + x_2}{2} \right) \right|^2 \quad (21)$$

where $\tilde{\mathbf{r}}_{z/2}(x)$, defined by

$$\tilde{\mathbf{r}}_{z/2}(x) = \int_{-\infty}^{\infty} \mathbf{r}(\eta) \exp \left[\frac{-ik(x-\eta)^2}{z} \right] d\eta \quad (22)$$

is proportional to the wave field that would result from illuminating the object with a normally incident plane wave and propagating the reflected wave a distance $z/2$. The resulting distribution is the object function blurred by Fresnel diffraction. This Fresnel-blurred image is less distinct than a normal image of the object; however, it is significantly less blurred than would be a conventional image obtained through the diffuser.

Lau Effect

In the Lau effect, an incoherently illuminated amplitude grating with grating constant d is followed by a second identical grating at a distance z_0 .¹⁹ If the second grating is followed by a converging spherical lens, fringes appear in the back focal plane of the lens whenever the distance z_0 is an integral multiple of $d^2/2\lambda$. The experimental geometry for observing this effect is illustrated in Fig. 2. If white light illumination is used, colored fringes are observed. Although the effect can be explained in terms of geometrical optics and scalar diffraction theory,²⁰ a much more elegant explanation results from a straightforward application of coherence theory.^{21–25}

The essence of a coherence-based analysis using the mutual intensity function (expressed in one-dimensional coordinates for simplicity) is as follows: The incoherently illuminated

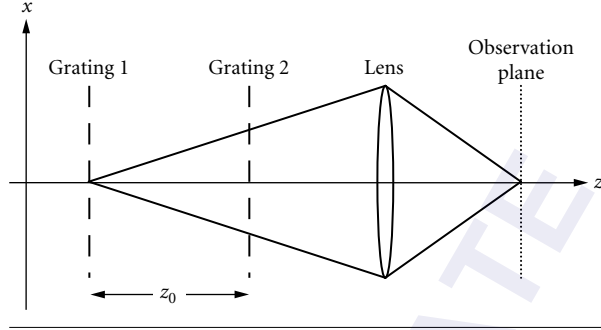


FIGURE 2 Geometry for examination of the Lau effect.

grating constitutes a periodic spatially incoherent source with period d with an intensity function described by

$$I_0(x) = \text{rect}\left(\frac{2x}{d}\right) \otimes \sum_{m=-\infty}^{\infty} \delta(x-md) \quad (23)$$

From the van Cittert-Zernike theorem [Eq. (10)] we know that the mutual intensity function of the propagated wave field arriving at the second grating can be expressed as

$$\begin{aligned} J_1(x_1, x_2) &= \frac{\exp(-i\varphi)}{(\lambda z_0)^2} \mathcal{F} \left\{ \text{rect}\left(\frac{2x}{d}\right) \otimes \sum_{m=-\infty}^{\infty} \delta(x-md) \right\} \Big|_{u=\frac{\Delta x}{\lambda z_0}} \\ &= \frac{\exp(-i\varphi)}{(\lambda z_0)^2} \left[\frac{d}{2} \text{sinc}\left(\frac{d\Delta x}{2\lambda z_0}\right) \sum_{m=-\infty}^{\infty} \delta\left(\Delta x - \frac{m\lambda z_0}{d}\right) \right] \end{aligned} \quad (24)$$

where \mathcal{F} denotes the Fourier transform operation. The presence of the Dirac delta functions in the above expression indicates that the propagated wave reaching the second grating is spatially coherent only for transverse separations that are integer multiples of the distance $\lambda z_0/d$. If the second grating is treated as a transmissive planar object with a transmission function analytically equivalent to the previously defined intensity function $I_0(x)$, Eq. (4) can be used to describe the mutual intensity of the light leaving the grating, obtaining

$$\begin{aligned} J_2(x_1, x_2) &= \left[\text{rect}\left(\frac{2x_1}{d}, \frac{2x_2}{d}\right) \otimes \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} \delta(x_1 - kd, x_2 - ld) \right] \\ &\quad \times \frac{\exp(-i\varphi)}{(\lambda z_0)^2} \left[\frac{d}{2} \text{sinc}\left(\frac{d\Delta x}{2\lambda z_0}\right) \sum_{m=-\infty}^{\infty} \delta\left(\Delta x - \frac{m\lambda z_0}{d}\right) \right] \end{aligned} \quad (25)$$

The first half of the above equation describes a two-dimensionally periodic spatial filter in the coherence domain with a period equal to d in both directions acting upon the mutual intensity of the light incident upon the grating. Maximum coherence (and hence maximum interference fringe-forming capability) is preserved when the period of the one-dimensional Dirac delta function in the second half of the equation (i.e., the nonzero coherence components of the incident light) is equal to the period of the two-dimensional Dirac delta function describing the coherence filter or $d = \lambda z_0/d$. A grating separation of $z_0 = d^2/\lambda$ meets this condition which is a special case of Lau's more general condition.

6.5 IMAGE FORMATION: LUKOSZ-TYPE SUPER-RESOLVING SYSTEM

It was in the analysis of imaging systems that coherence theory saw its first major application. Frits Zernike, in his classic 1938 paper on “The concept of degree of coherence and its application to optical problems,” showed how the mutual coherence function could be employed in the wave-optics analysis of systems of lenses and apertures in such a way that modification of the system through extension, for example, through the addition of a field lens, did not necessarily require a return to the starting point in the analysis, that is, to the source of illumination.¹⁶ In this section we illustrate the utility of Zernike’s approach through application to a particular super-resolving imaging system.

In the 1960s Lukosz published a pair of papers on super-resolution imaging systems that continue to interest and intrigue people.^{26,27} Lukosz-type systems can be used to increase the space-bandwidth product—in contemporary terms, the number of pixels—that can be transmitted from the object plane to the image plane with a system of a given numerical aperture. The intriguing nature and popularity of Lukosz-type imaging systems notwithstanding, they are surprisingly difficult to analyze in detail, and Lukosz’s original papers present challenging reading for one who wants to understand quickly just how to model the systems mathematically. We show in this section how coherence theory tools can be used to good effect in the analysis of a super-resolution imaging system of our invention that is closely related to those described by Lukosz.²⁸ Closely related systems are described in Refs. 29 and 30.

The system of interest, illustrated in Fig. 3, consists of a quasi-monochromatic object, a pair of identical time-varying diffusers, a pair of lenses, and a pinhole aperture. The complex wave amplitude following Mask 2 (one-dimensional notation is used for simplicity) for an arbitrary pupil (aperture) function is given by

$$\mathbf{u}'_3(x, t) = \mathbf{M}(x, t) \int_{-\infty}^{\infty} \mathbf{M}(-\xi, t) \mathbf{u}_{\text{inc}}(-\xi, t) h_{\text{sys}}(x - \xi) d\xi \quad (26)$$

where $h_{\text{sys}}(x)$ is the complex-amplitude point spread function for the imaging system and $\mathbf{M}(x, t)$ is the complex amplitude transmittance of the masks. The wave field incident on Mask 1 is related to the object complex amplitude by $\mathbf{u}_{\text{inc}}(x) = \mathbf{u}_{\text{obj}}(x) \otimes h_d(x)$, where $h_d(x)$ is the Fresnel kernel appropriate for propagation through a distance d . Because of the dynamic (moving) diffusers and the effect of the small aperture, calculation of the corresponding optical intensity tells us virtually nothing about the object distribution. The mutual intensity in the output plane, on the other hand, can tell us a great deal about the object.

The mutual intensity of the wave incident on Mask 1 contains information about the mutual intensity of the object distribution through the relationship:

$$J_{\text{inc}}(x_1, x_2) = J_{\text{obj}}(x_1, x_2) \otimes h_d(x_1) \otimes h_d^*(x_2) \quad (27)$$

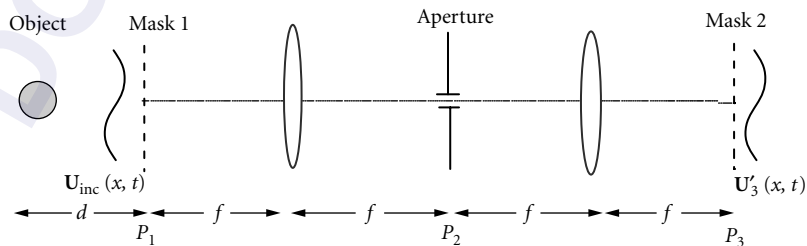


FIGURE 3 Geometry for Lukosz-type super-resolution optical system.

Because of the reversibility of the wave propagation phenomenon (manifested by the absence of nulls in the wave propagation transfer function), the mutual intensity of the object may be inferred from the mutual intensity of this incident wave field through

$$J_{\text{obj}}(x_1, x_2) = J_{\text{inc}}(x_1, x_2) \otimes h_d^*(x_1) \otimes h_d(x_2) \quad (28)$$

It is the nature of the system of Fig. 3 that this information can be transmitted to the output plane. As was done in the enhanced backscatter analysis, let the temporal behavior of the incident complex wave amplitude and the masks be statistically independent so that $\langle \mathbf{u}_{\text{inc}}(x, t) \mathbf{M}(x, t) \rangle = \langle \mathbf{u}_{\text{inc}}(x, t) \rangle \langle \mathbf{M}(x, t) \rangle$. In addition, assume that the masks are statistically homogeneous random phase screens that satisfy the fourth-moment theorem, as reflected by Eq. (14), so that the time average with respect to the diffusers within the integral simplifies to a product of second-order moments and the autocorrelation function of the mask is sufficiently narrow that it can be modeled by a delta function within an integral.³¹ Exploiting these conditions and exploiting the sifting property of the delta function, we can calculate the mutual intensity of the wave immediately following the second mask, that is, in plane 3:

$$J_3(x_1, x_2) = \kappa^2 J_{\text{inc}}^*(x_1, x_2) \left| \hat{\mathbf{P}} \left(\frac{x_1 + x_2}{\lambda f} \right) \right|^2 + \kappa \gamma(x_1 - x_2) \int_{-\infty}^{\infty} I_{\text{inc}}(-\xi) \hat{\mathbf{P}} \left(\frac{x_1 - \xi}{\lambda f} \right) \hat{\mathbf{P}}^* \left(\frac{x_2 - \xi}{\lambda f} \right) d\xi \quad (29)$$

where $\gamma(x)$ is the autocorrelation of the mask function and where I_{inc} is the optical intensity of the incident wave field.

The first term on the right in this equation is proportional to $|\hat{\mathbf{P}}|^2$, the modulus of the imaging system impulse response ($\hat{\mathbf{P}}$ denotes the Fourier transform of the pupil function associated with the aperture), times the mutual intensity of the wave field incident on Mask 1. The second term corresponds to a spatially incoherent wave field. Through interferometric measurements, it is possible to infer the mutual intensity of the incident wave field, that is, to determine J_{inc} and thus J_{obj} .

6.6 EFFICIENT SAMPLING OF COHERENCE FUNCTIONS

Several schemes for optical imaging rely on the measurement or transmission of the spatiotemporal coherence function produced at an aperture by the waves from a distant object.^{32,33} If the object is planar and spatially incoherent, this coherence function, described by the van Cittert-Zernike theorem,³⁴ is, to within a known quadratic phase factor, a function of the vector separation $(\Delta x, \Delta y)$ between pairs of points in the measurement plane, as illustrated in Eq. (10). As a consequence, the number of measurements required to characterize the function in a Nyquist sense is comparatively small. It is not necessary to sample the optical wave field at all possible pairs of points on a grid but, rather, only at a single pair of points for a given spacing $(\Delta x, \Delta y)$. For nonplanar objects, on the other hand, the coherence function has a more general form,³⁵ and the number of samples required of the coherence function necessarily increases. In order to keep the measurement time as small as possible, an efficient sampling strategy is desirable.

Nonuniform sampling grids have been shown to reduce the total number of samples required to unambiguously characterize the complex amplitude associated with a quasi-monochromatic optical wave field in certain cases.^{36,37} In the following analysis it is shown that a nonuniform sampling scheme can also be effectively applied to the problem of sampling the mutual intensity produced in an aperture by a quasi-monochromatic, spatially incoherent three-dimensional object distribution. The analysis presented is of brief out of necessity, missing details being presented elsewhere.³⁸

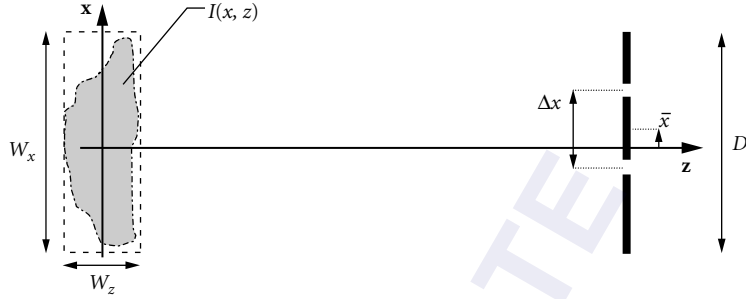


FIGURE 4 The optical coherence of the wave field generated by the source intensity distribution $I(x, z)$ is sampled in a plane some distance away from the center of the rectangular area bounding $I(x, z)$.

The source object, described by optical intensity $I(x, z)$, is assumed to be contained within a rectangular area of transverse width W_x and longitudinal depth W_z centered about the origin of the coordinate system, as shown in Fig. 4 (only one transverse dimension is considered for simplicity, extension to two being straightforward). If Fresnel-regime conditions are satisfied, the mutual intensity in the measurement plane can be shown to be given by the equation

$$J(x_1, x_2) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \frac{1}{\lambda(z-\zeta)} I(\xi, \zeta) \exp \left\{ j \frac{\pi}{\lambda(z-\zeta)} [(x_1 - \xi)^2 - (x_2 - \xi)^2] \right\} d\xi d\zeta \quad (30)$$

This function is conveniently expressed in terms of parameters $\bar{x} = (x_1 + x_2)/2$ and $\Delta x = (x_1 - x_2)$, with the result

$$J(\bar{x}, \Delta x) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \frac{1}{\lambda(z-\zeta)} I(\xi, \zeta) \exp \left[j \frac{2\pi\Delta x}{\lambda(z-\zeta)} (\bar{x} - \xi) \right] d\xi d\zeta \quad (31)$$

Criteria for both uniform and nonuniform sampling of this function are most easily determined if it is Fourier transformed, once with respect to the sample separation Δx and once with respect to sample-pair center coordinate \bar{x} :

$$\hat{J}_{\Delta x}(\bar{x}, \nu) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \frac{1}{\lambda(z-\zeta)} I(\xi, \zeta) \delta \left[\nu - \frac{(\bar{x} - \xi)}{\lambda(z-\zeta)} \right] d\xi d\zeta \quad (32)$$

$$\hat{J}_{\bar{x}}(\Delta\nu, \Delta x) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \frac{1}{\lambda(z-\zeta)} I(\xi, \zeta) \exp \left[-j \frac{2\pi\Delta x \xi}{\lambda(z-\zeta)} \right] \delta \left[\Delta\nu + \frac{\Delta x}{\lambda(z-\zeta)} \right] d\xi d\zeta \quad (33)$$

If the finite support of the object in x and z is taken into account, these functions can be shown to have regions of support bounded by lines given by the equations

$$\nu = \frac{1}{\lambda} \left(z \pm \frac{W_z}{2} \right)^{-1} \left(\bar{x} \pm \frac{W_x}{2} \right) \quad \Delta\nu = \frac{\Delta x}{\lambda} \left(z \pm \frac{W_z}{2} \right)^{-1} \quad (34)$$

For distances z satisfying the condition $z \gg W_x$, the corresponding information bandwidths are well approximated by the expressions

$$B_\nu = \frac{W_x}{\lambda z} \quad B_{\Delta\nu} = \frac{W_z \Delta x}{\lambda z^2} \quad (35)$$

A straightforward analysis based on the standard Nyquist criterion shows that uniform sample locations should be located according to the equations (letting D denote the size of the aperture in the measurement plane)

$$\bar{x}_m = \frac{m\lambda z}{2D} \quad \Delta x_n = \frac{n\lambda z}{D+W_x} \quad (36)$$

whereas a generalized Nyquist analysis, appropriate for nonuniform sampling, yields the results

$$\bar{x}_{m,n} = \frac{mW_x z}{nW_z} \quad \Delta x_n = \frac{n\lambda z}{W_x} \quad (37)$$

where the double subscripts m, n denote the sample order. The number of samples required is found by counting all valid indices (m, n) , given the extent D of the measurement area. A comparison of the two sampling regimes shows that nonuniform sampling reduces the number of samples required by a factor $[(1+D/W_x)(1+D^2/\lambda z)]^{-1}$ and that the minimum separation between measurement points is increased by a factor $(1+D/W_x)$. This latter difference is beneficial because a constraint on the measurement apparatus is relaxed. Numerical demonstration based on system specifications given in Ref. 35 yields excellent results.

6.7 AN EXAMPLE OF WHEN NOT TO USE COHERENCE THEORY

Coherence theory can provide an extremely valuable and sometimes indispensable tool in the analysis of certain phenomena. On the other hand, it can be applied unnecessarily to certain optical systems where it tends to obfuscate rather than clarify system operation. An example is given by the Koehler-illumination imaging system modeled in Fig. 5. In this system, a uniform, spatially incoherent source is imaged by condenser lens 1 into the pupil plane of the imaging optics, the latter formed by lenses 2 and 3. Beginning with the source distribution, one can calculate, through a succession of operations based on relationships presented in Sec. 6.3, the coherence functions appropriate for the object plane, the pupil plane, and the image plane. The effect on the irradiance distributions in these planes can through this means be investigated. But is this the best way to proceed? In fact, if only the image-plane irradiance is of ultimate interest, it is not, the method being unnecessarily complicated.

Frits Zernike made an important but sometimes forgotten observation in his 1938 paper on the concept of the degree of coherence.¹⁶ A coherence-function-based analysis of a system consisting of a cascade of subsystems, such as that just considered, can be advantageous if an objective is to calculate the irradiance of the wave field in more than a single plane. Thus, with regard to the system of Fig. 5, if calculations of the object illumination, the pupil-plane irradiance, and the image distribution are all

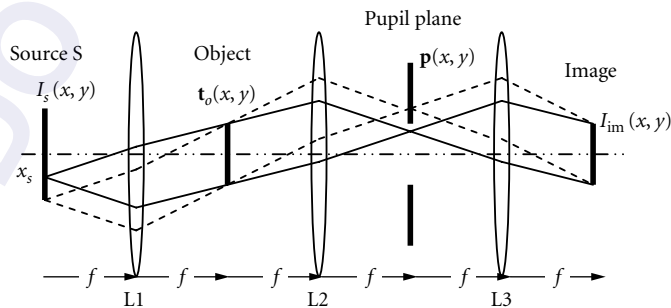


FIGURE 5 Koehler-illumination imaging model.

desired, then, assuming a quasi-monochromatic source, an analysis based on the plane-to-plane propagation of the mutual intensity of the light is the proper choice. In each plane of interest, the irradiance can be calculated from the mutual intensity without loss of the ability to propagate that mutual intensity function further in the system. If, on the other hand, the objective is simply the calculation of the image-plane irradiance distribution, given the source, object, and pupil functions, then, so long as the source is spatially incoherent, the pre-coherence theory approach to analysis is preferable. The image plane irradiance produced by light from a single source point is calculated, with an integration over the source distribution following. In many cases, this form of analysis provides greater insight into the operation of the system under investigation. For example, in the case of Koehler illumination, it shows that the source distribution should be uniform and that its image should overfill the aperture stop (pupil) by an amount proportional to the spatial-frequency bandwidth of the object wave amplitude transmittance function. Such insight is not readily available through a coherence-theory-based analysis.

6.8 CONCLUDING REMARKS

Some important points are summarized here:

1. Although it often appears in the analysis of optical systems, the complex amplitude associated with an optical wave field cannot be measured directly because of the very high frequencies of the wave field oscillations.
2. The optical intensity of the wave, which involves a time average, *can* be measured.
3. The mutual intensity of the wave can also be determined from measurements made using an interferometer.
4. Knowledge of the optical intensity of the wave in a given plane of an optical system does not, in general allow calculation of the mutual intensity and/or optical intensity of the wave in subsequent planes. (An exception is when the wave in that plane is spatially incoherent.)
5. By way of contrast, knowledge of the mutual intensity of the wave in a given plane *does* allow calculation of the mutual intensity and/or the optical intensity of the wave in subsequent planes.
6. If circular complex gaussian statistics for the wave field can be assumed, third- and higher-order moments of the wave field can be inferred from the second-order statistics—that is, from the mutual intensity—of the wave.^{11,31}

Taken together, these statements imply that the mutual intensity of the quasi-monochromatic wave field in an optical system conveys all of the information that is measurable—and, in a valid sense, meaningful—through the system. Restating point (5) from this perspective, if the mutual intensity of a wave in one plane of an optical system is known, the mutual intensity—and, hence, the information content of the wave—can be determined in subsequent planes of that system.

6.9 REFERENCES

1. H. H. Hopkins, "Applications of Coherence Theory in Microscopy and Interferometry," *J. Opt. Soc. Am.* **47** (6), June 1957, pp. 508–526.
2. C. Paul Montgomery, J. M. Lussert, P. Vabre, and D. Benhaddou, "Super-Resolution 3D Optical Imaging of Semiconductors Using Coherence Microscopy," *SPIE Proc. Vol. 2412 Three-Dimensional Microscopy: Image Acquisition and Processing II*, T. Wilson and C. J. Cogswell, eds., March 1995, pp. 88–94.
3. M. Bass ed., *Handbook of Optics*, Vol. 1, McGraw-Hill, Inc., New York, 1995, Chapter 4.
4. L. Mandel and E. Wolf, *Optical Coherence and Quantum Optics*, Cambridge University Press, Cambridge, 1995.
5. L. Mandel and E. Wolf, eds., "Selected Papers on Coherence and Fluctuations of Light, (1850–1966)," SPIE Press Book, **MS19**, 1990.

6. H. Park, M. Chodorow, and R. Kompfner, "High Resolution Optical Ranging System," *App. Opt.* **20**, 1981, pp. 2389–2394.
7. R. Youngquist, S. Carr, and D. E. N. Davies, "Optical Coherence-Domain Reflectometry: A New Optical Evaluation Technique," *Opt. Lett.* **12**, 1987, pp. 158–160.
8. H. H. Gilgen, R. P. Novak, R. P. Salathe, W. Hodel, and P. Beaud, "Submillimeter Optical Reflectometry," *J. of Lightwave Tech.* **7**, 1989, pp. 1225–1233.
9. B. E. Bouma and G. J. Tearney, *Handbook of Optical Coherence Tomography*, Informa Health Care, New York, 2001.
10. M. E. Brezinski, *Optical Coherence Tomography: Principles and Applications*, Academic Press, New York, 2006.
11. J. W. Goodman, *Statistical Optics*, Wiley, New York, 1985.
12. M. Bass, ed., *Handbook of Optics*, Vol. 1, McGraw-Hill, Inc., New York, 1995, Chapter 2.
13. A. S. Marathay, *Elements of Optical Coherence Theory*, Wiley, New York, 1982, pp. 71–73.
14. G. B. Parrent Jr., "On the Propagation of Mutual Coherence," *J. Opt. Soc. Am. A.* **49**, 1959, pp. 787–793. In: *Selected Papers on Coherence and Fluctuations of Light*, eds. L. Mandel and E. Wolf, Vol 1. Dover Publications, Inc. New York.
15. L. Mandel and E. Wolf, "Some Properties of Coherent Light," *J. Opt. Soc. Am. A.* **51**, 1961, pp. 815–819. In: *Selected Papers on Coherence and Fluctuations of Light*, eds. L. Mandel and E. Wolf, Vol 1. Dover Publications, Inc. New York.
16. F. Zernike, "The Concept of Degree of Coherence and Its Application to Optical Problems," *Physica* **5**, 1938, pp. 785–795. In: *Selected Papers on Coherence and Fluctuations of Light*, eds. L. Mandel and E. Wolf, Vol 1. Dover Publications, Inc. New York.
17. W. T. Rhodes and G. Welch, "Determination of a Coherent Wave Field after Double Passage through a Diffuser," *J. Opt. Soc. Am. A.* **9**(2), Feb 1992, pp. 341–343.
18. G. Welch and W. T. Rhodes, "Imaging Diffuser-Obscured Objects by Making Coherence Measurements on Enhanced Backscatter Radiation," *J. Opt. Soc. Am. A.* **9**(4), April 1992, pp. 539–542.
19. E. Von Lau, "Beugungerscheinungen an Doppelrastern," *Ann. Phys.* **6**, 1948, pp. 417.
20. J. Jahns and A. W. Lohmann, "The Lau Effect (a Diffraction Experiment with Incoherent Illumination)" *Opt. Comm.* **28**(3), 1979, pp. 263–267.
21. F. Gori, "Lau Effect and Coherence Theory," *Opt. Comm.* **31**, 1979, pp. 4–8.
22. R. Sudol and B. J. Thompson, "An Explanation of the Lau Effect Based on Coherence Theory," *Opt. Comm.* **31**, 1979, pp. 105–110.
23. R. Sudol and B. J. Thompson, "Lau Effect: Theory and Experiment," *Appl. Opt.* **20**, 1981, pp. 1107–1116.
24. S. Chitralakha, K. V. Avudainayagam, and S. V. Pappu, "Role of Spatial Coherence on the Rotation Sensitivity of Lau Fringes: An Experimental Study," *Appl. Opt.* **28**, 1989, pp. 345–349.
25. S. Cartwright, "Incoherent Optical Processing: A Coherence Theory Approach," *Appl. Opt.* **23**, 1984, pp. 318–323.
26. W. Lukosz, "Optical Systems with Resolving Powers Exceeding the Classical Limit," *J. Opt. Soc. Am. A.* **56** (11), November 1966, pp. 1463–1472.
27. W. Lukosz, "Optical Systems with Resolving Powers Exceeding the Classical Limit. II," *J. Opt. Soc. Am. A.* **57**, July 1967, pp. 932–941.
28. G. Welch, "Applications of Coherence to Enhanced Backscatter and Superresolving Imaging Systems," PhD dissertation. Georgia Institute of Technology, Atlanta, GA, 1995.
29. A. Cunha and E. N. Leith, "One-Way Phase Conjugation with Partially Coherent Light and Superresolution," *Opt. Lett.* **13**(12), 1988, pp. 1105–1107.
30. P. C. Sun and E. N. Leith, "Superresolution by Spatial-Temporal Encoding Methods," *Appl. Opt.* **31**(23), 1992, pp. 4857–4862.
31. I. S. Reed, "On a Moment Theorem for Complex Gaussian Processes," *IRE Trans. Inf. Theory* **IT-8**, 1962, pp. 194–195.
32. D. L. Marks, R. A. Stack, and D. J. Brady, "Three-Dimensional Coherence Imaging in the Fresnel Domain," *Appl. Opt.* **38**(8), 1999, pp.1332–1342.
33. G. Welch and W. T. Rhodes, "Image Reconstruction by Spatio-Temporal Coherence Transfer." In: *Free-Space Laser Communication and Laser Imaging*, D. G. Voelz and J. C. Ricklin, eds. (Proc. SPIE, Vol. 4489. 2001).

34. J. W. Goodman, *Statistical Optics*, Wiley, New York, 1985, Chapter 5.
35. J. Rosen and A. Yariv, "General Theorem of Spatial Coherence: Application to Three-Dimensional Imaging," *J. Opt. Soc. Am. A.* **13**(10), 1996, pp. 2091–2095.
36. A. Vanderlugt, "Optimum Sampling of Fresnel Transforms," *Appl. Opt.* **29**(23), 1990, pp. 3352–3361.
37. F. S. Roux, "Complex-Valued Fresnel-Transform Sampling," *Appl. Opt.* **34**(17), 1995, pp. 3128–3135.
38. J. C. James, G. Welch, and W. T. Rhodes, "Nonuniform Coherence Sampling for Fields Produced by Incoherent Three-Dimensional Sources," *J. Opt. Soc. Am. A.* **20**, April 2003, pp. 668–677.

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

SCATTERING BY PARTICLES

Craig F. Bohren

*Pennsylvania State University
University Park, Pennsylvania*

7.1 GLOSSARY

a	radius
a_n, b_n	scattering coefficients
C	cross section
D_n	logarithmic derivative, $d/d\rho[\ln\psi_n(\rho)]$
E	Electric field strength
\mathbf{e}_x	unit vector in the x direction
f	Nv
G	projected particle area
h	thickness
I	irradiance
I, Q, U, V	Stokes parameters
j	running index
k	imaginary part of the refractive index, $2\pi/\lambda$
m	relative complex refractive index
N	number
n	running index
P_n	associated Legendre functions of the first kind
p	phase function, normalized differential scattering cross section
Q	efficiencies or efficiency factors
r	distance
S	element of the amplitude-scattering matrix
v	volume
W	power
X	scattering amplitude

x	size parameter, ka
α	absorption coefficient
θ	angle
λ	wavelength
π_n	$P_n^1/\sin \theta$
τ_n	$dP_n^1/d\theta$
ψ, ζ	Riccati-Bessel functions
Ω	solid angle
ω	radian frequency
\parallel	Parallel
\perp	perpendicular
Re	real part of

Subscripts

abs	absorbed
ext	extinction
sca	scattered

7.2 INTRODUCTION

Light scattering by particles plays starring and supporting roles on a variety of stages: astronomy, cell biology, colloid chemistry, combustion engineering, heat transfer, meteorology, paint technology, solid-state physics—the list is almost endless. The best evidence of the catholicity of scattering by particles is the many journals that publish papers about it.

Scattering by single particles is the subject of monographs by van de Hulst,¹ Deirmendjian,² Kerker,³ Bayvel and Jones,⁴ Bohren and Huffman,⁵ Barber and Hill,⁶ and of a collection edited by Kerker.⁷ Two similar collections contain papers on scattering by atmospheric particles⁸ and by chiral particles⁹ (ones not superposable on their mirror images); scattering by chiral particles is also treated by Lakhtakia et al.¹⁰ Papers on scattering by particles are included in collections edited by Gouesbet and Gréhan¹¹ and by Barber and Chang.¹² Within this *Handbook* scattering by particles is touched upon in Chap. 9, “Volume Scattering in Random Media,” in this volume and Chap. 3, “Atmospheric Optics,” in Vol. V. A grand feast is available for those with the juices to digest it. What follows is a mere snack.

A particle is an aggregation of sufficiently many molecules that it can be described adequately in macroscopic terms (i.e., by constitutive parameters such as permittivity and permeability). It is a more or less well-defined entity unlike, say, a density fluctuation in a gas or a liquid. Single molecules are not particles, even though scattering by them is in some ways similar (for a clear but dated discussion of molecular scattering, see Martin¹³).

Scattering by single particles is discussed mostly in the wave language of light, although multiple scattering by incoherent arrays of many particles can be discussed intelligibly in the photon language. The distinction between single and multiple scattering is observed more readily on paper than in laboratories and in nature. Strict single scattering can exist only in a boundless void containing a lone scatterer illuminated by a remote source, although single scattering often is attained to a high degree of approximation. A distinction made less frequently is that between scattering by coherent and incoherent arrays. In treating scattering by coherent arrays, the wave nature of light cannot be ignored: phases *must* be taken into account. But in treating scattering by incoherent arrays, phases *may* be ignored.

Pure water is a coherent array of water molecules; a cloud is an incoherent array of water droplets. In neither of these arrays is multiple scattering negligible, although the theories used to describe them may not explicitly invoke it.

The distinction between incoherent and coherent arrays is not absolute. Although a cloud of water droplets is usually considered to be an incoherent array, it is not such an array for scattering in the forward direction. And although most of the light scattered by pure water is accounted for by the laws of specular reflection and refraction, it also scatters light—weakly yet measurably—in directions not accounted for by these laws.¹³

A single particle is itself a coherent array of many molecules, but can be part of an incoherent array of many particles, scattering collectively in such a way that the phases of the waves scattered by each one individually are washed out. Although this section is devoted to single scattering, it must be kept in mind that multiple scattering is not always negligible and is not just scaled-up single scattering. Multiple scattering gives rise to phenomena inexplicable by single-scattering arguments.¹⁴

7.3 SCATTERING: AN OVERVIEW

Why is light scattered? No single answer will be satisfactory to everyone, yet because scattering by particles has been amenable to treatment mostly by classical electromagnetic theory, our answer lies within this theory.

Although palpable matter may appear to be continuous and is often electrically neutral, it is composed of discrete electric charges. Light is an oscillating electromagnetic field, which can excite the charges in matter to oscillate. Oscillating charges radiate electromagnetic waves, a fundamental property of such charges with its origins in the finite speed of light. These radiated electromagnetic waves are scattered waves, waves excited or driven by a source external to the scatterer: an incident wave from the source excites secondary waves from the scatterer; the superposition of all these waves is what is observed. If the frequency of the secondary waves is (approximately) that of the source, these waves are said to be *elastically scattered* (the term *coherently scattered* is also used).

Scientific knowledge grows like the accumulation of bric-a-brac in a vast and disorderly closet in a house kept by a sloven. Few are the attempts at ridding the closet of rusty or obsolete gear, at throwing out redundant equipment, at putting things in order. For example, spurious distinctions are still made between reflection, refraction, scattering, interference, and diffraction despite centuries of accumulated knowledge about the nature of light and matter.

Countless students have been told that specular reflection is localized at smooth surfaces, and that photons somehow rebound from them. Yet this interpretation is shaky given that even the smoothest surface attainable is, on the scale of a photon, as wrinkled as the back of a cowboy's neck. Photons conceived of as tiny balls would be scattered in all directions by such a surface, for which it is difficult even to define what is meant by an angle of incidence.

Why do we think of reflection occurring at surfaces rather than because of them whereas we usually do not think of scattering by particles in this way? One reason is that we can see the surfaces of mirrors and ponds. Another is the dead hand of traditional approaches to the laws of specular reflection and refraction.

The empirical approach arrives at these laws as purely geometrical summaries of what is observed—and a discreet silence is maintained about underlying causes. The second approach is by way of continuum electromagnetic theory: reflected and refracted fields satisfy the Maxwell equations. Perhaps because this approach, which also yields the Fresnel formulas, entails the solution of a boundary-value problem, reflected and refracted fields are mistakenly thought to originate from boundaries rather than from all the illuminated matter they enclose. This second approach comes to grips with the nature of light but not of matter, which is treated as continuous. The third approach is to explicitly recognize that reflection and refraction are consequences of scattering by discrete matter. Although this scattering interpretation was developed by Paul Ewald and Carl Wilhelm Oseen early in this century, it has diffused with glacial slowness. According to this interpretation, when the optically smooth interface between optically homogeneous dissimilar media is illuminated, the reflected and refracted waves are superpositions of vast numbers of secondary waves excited by the incident wave.

Thus reflected and refracted light is, at heart, an interference pattern of scattered light. Doyle¹⁵ showed that although the Fresnel equations are obtained from macroscopic electromagnetic theory, they can be dissected to reveal their microscopic underpinnings.

No optics textbook would be complete without sections on interference and diffraction, a distinction without a difference: there is no diffraction without interference. Moreover, diffraction is encumbered with many meanings. Van de Hulst¹ lists several: small deviations from rectilinear propagation; wave motion in the presence of an obstacle; scattering by a flat particle such as a disk; an integral relation for a function satisfying the wave equation. To these may be added scattering near the forward direction and by a periodic array.

Van de Hulst stops short of pointing out that a term with so many meanings has no meaning. Even the etymology of diffraction is of little help: it comes from a Latin root meaning to break.

There is no fundamental distinction between diffraction and scattering. Born and Wolf¹⁶ refer to scattering by a sphere as diffraction by a sphere. I leave it as a penance for the reader to devise an experiment to determine whether a sphere scatters light or diffracts it.

The only meaningful distinction is that between approximate theories. Diffraction theories obtain answers at the expense of obscuring the physics of the interaction of light with matter. For example, an illuminated slit in an opaque screen may be the mathematical source but it is not the physical source of a diffraction pattern. Only the screen can give rise to secondary waves that yield the observed pattern. Yet generations of students have been taught that empty space is the source of the radiation diffracted by a slit. To befuddle them even more, they also have been taught that two slits give an interference pattern whereas one slit gives a diffraction pattern.

If we can construct a mathematical theory (diffraction theory) that enables us to avoid having to explicitly consider the nature of matter, all to the good. But this mathematical theory and its quantitative successes should not blind us to the fact that we are pretending. Sometimes this pretense cannot be maintained, and when this happens a finger is mistakenly pointed at “anomalies,” whereas what is truly anomalous is that a theory so devoid of physical content could ever give adequate results.

A distinction must be made between a physical process and the superficially different theories used to describe it. There is no fundamental difference between specular reflection and refraction by films, diffraction by slits, and scattering by particles. All are consequences of light interacting with matter. They differ only in their geometries and the approximate theories that are sufficient for their quantitative description. The different terms used to describe them are encrustations deposited during the slow evolution of our understanding of light and matter.

7.4 SCATTERING BY PARTICLES: BASIC CONCEPTS AND TERMINOLOGY

A single particle can be considered a collection of tiny dipolar antennas driven to radiate (scatter) by an incident oscillating electric field. Scattering by such a coherent array of antennas depends on its size and shape, the observation angle (scattering angle), the response of the individual antennas (composition), and the polarization state and frequency of the incident wave. Geometry, composition, and the properties of the illumination are the determinants of scattering by particles.

Perhaps the only real difference between optics and electrical engineering is that electrical engineers can measure amplitudes and phases of fields whereas the primary observable quantity in optics is the time-averaged Poynting vector (irradiance), an amplitude squared. Several secondary observables are inferred from measurements of this primary observable. Consider, for example, a single particle illuminated by a beam with irradiance I_i . The total power scattered by this particle is W_{sca} . Within the realm of linear optics, the scattered power is proportional to the incident irradiance. This proportionality can be transformed into an equality by means of a factor C_{sca} :

$$W_{\text{sca}} = C_{\text{sca}} I_i \quad (1)$$

For Eq. (1) to be dimensionally homogeneous C_{sca} must have the dimensions of area, hence C_{sca} has acquired the name *scattering cross section*.

Particles absorb as well as scatter electromagnetic radiation. The rate of absorption W_{abs} by an illuminated particle, like scattered power, is proportional to the incident irradiance:

$$W_{\text{abs}} = C_{\text{abs}} I_i \quad (2)$$

where C_{abs} is the *absorption cross section*. The sum of these cross sections is the *extinction cross section*:

$$C_{\text{ext}} = C_{\text{sca}} + C_{\text{abs}} \quad (3)$$

Implicit in these definitions of cross sections is the assumption that the irradiance of the incident light is constant over lateral dimensions large compared with the size of the illuminated particle. This condition is necessarily satisfied by a plane wave infinite in lateral extent, which, much more often than not, is the source of illumination in light-scattering theories.

The extinction cross section can be determined (in principle) by measuring transmission by a slab populated by N identical particles per unit volume. Provided that multiple scattering is negligible, the incident and transmitted irradiances I_i and I_t are related by

$$I_t = I_i e^{-NC_{\text{ext}}h} \quad (4)$$

where h is the thickness of the slab. Only the sum of scattering and absorption can be obtained from transmission measurements. To separate extinction into its components requires additional measurements.

Equation (4) requires that all particles in the slab be identical. They are different if they differ in size, shape, composition, or orientation (incident beams are different if they differ in wavelength or polarization state). Equation (4) is generalized to a distribution of particles by replacing NC_{ext} with

$$\sum_j N_j C_{\text{ext},j} \quad (5)$$

where j denotes all parameters distinguishing one particle from another.

Instead of cross sections, normalized cross sections called *efficiencies* or *efficiency factors*, Q_{sca} , Q_{abs} , and Q_{ext} , often are presented. The normalizing factor is the particle's area G projected onto a plane perpendicular to the incident beam. No significance should be attached to efficiency used as shorthand for normalized cross section. The normalization factor is arbitrary. It could just as well be the total area of the particle or, to honor Lord Rayleigh, the area of his thumbnail.

Proper efficiencies ought to be less than unity, whereas efficiencies for scattering, absorption, and extinction are not so constrained. Moreover, some particles—soot aggregates, for example—do not have well-defined cross-sectional areas. Such particles have cross sections for scattering and absorption but the corresponding efficiencies are nebulous.

If any quantity deserves the designation efficiency it is the cross section per particle volume ν . Equation (4) can be rewritten to display this:

$$I_t = I_i e^{-fh(C_{\text{ext}}/\nu)} \quad (6)$$

where $f = N\nu$ is the total volume of particles per unit slab volume. For a given particle loading, specified by fh (volume of particles per unit slab area), transmission is a minimum when C_{ext}/ν is a maximum.

Each way of displaying extinction (or scattering) versus particle size or wavelength of the incident beam tells a different story. This is illustrated in Fig. 1, which shows the scattering cross section, scattering efficiency, and scattering cross section per unit volume of a silicate sphere in air illuminated by visible light. These curves were obtained with van de Hulst's simple *anomalous diffraction* approximation¹ (all that is anomalous about it is that it gives such good results). Each curve yields a different answer to the question, what size particle is most efficient at scattering light? And comparison of Figs. 1c and 2 shows that scattering by a particle and specular reflection are similar.

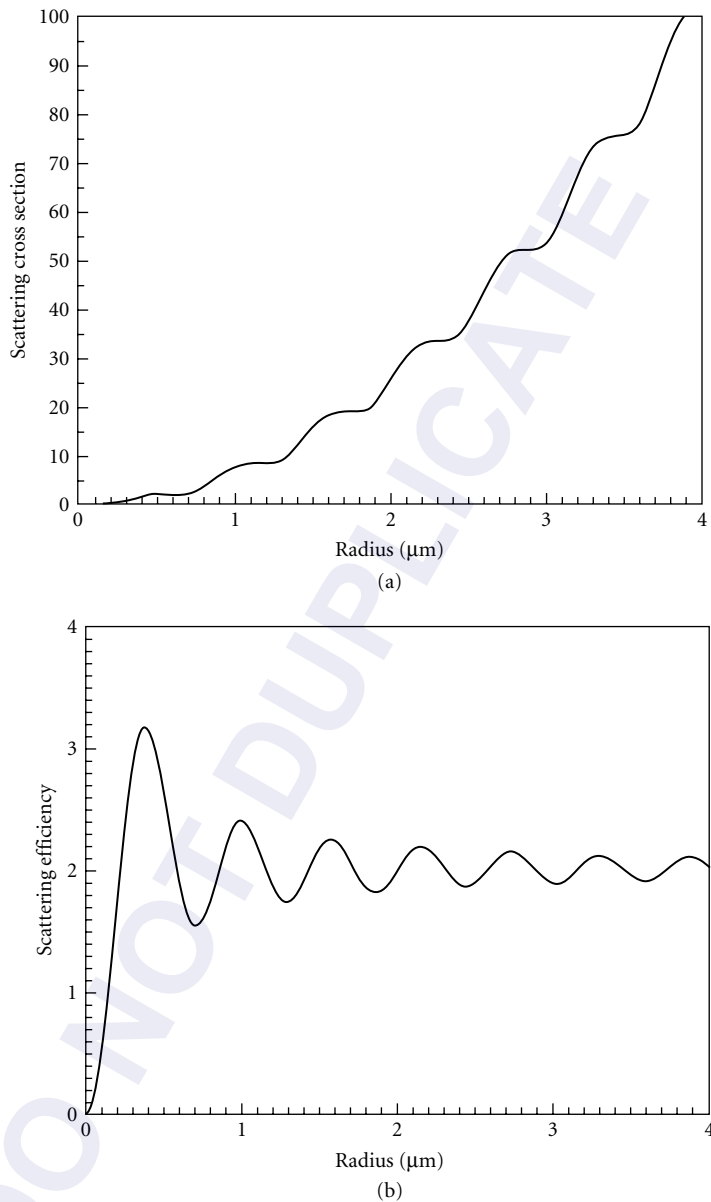
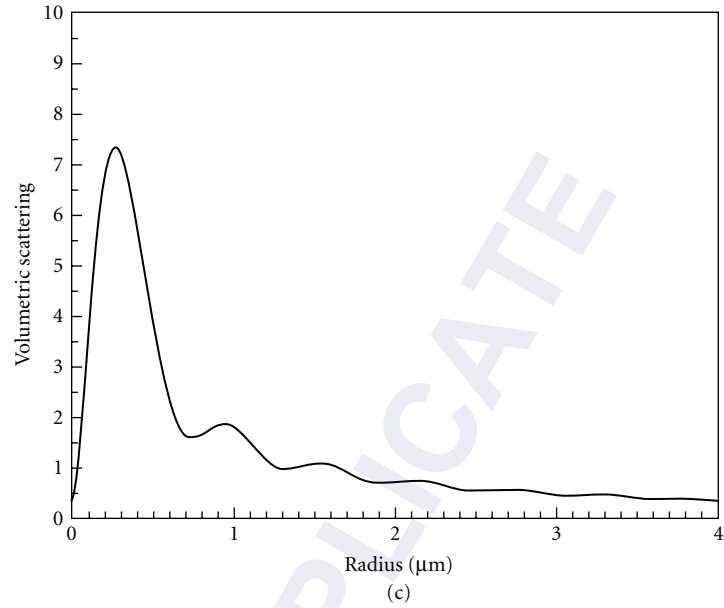
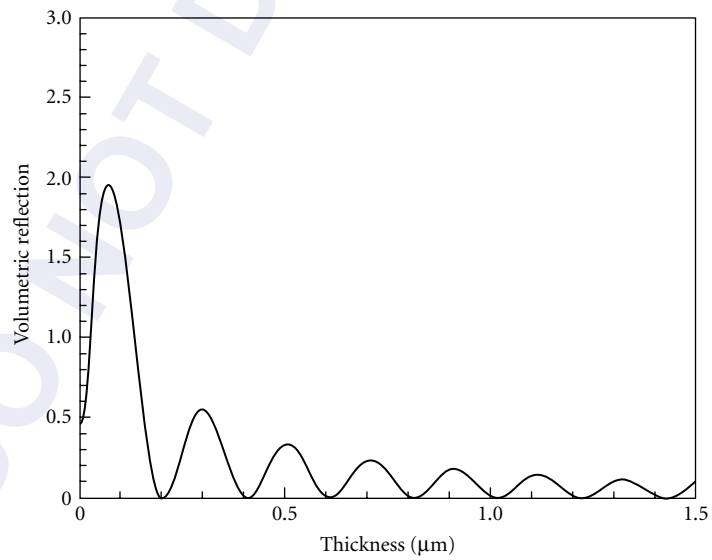


FIGURE 1 Scattering of visible light by a silicate sphere calculated using the anomalous diffraction approximation: (a) scattering cross section; (b) scattering efficiency (cross section normalized by projected area); and (c) volumetric scattering cross section (cross section per unit particle volume).

**FIGURE 1** (Continued)**FIGURE 2** Reflected power per unit incident irradiance and unit volume of a silicate slab normally illuminated by visible light (reflectance divided by slab thickness).

At sufficiently large distances r from a scatterer of bounded extent, the scattered field \mathbf{E}_s decreases inversely with distance and is transverse:

$$\mathbf{E}_s \sim \frac{e^{ik(r-z)}}{-ikr} \mathbf{X} E \quad (kr \gg 1) \quad (7)$$

where $k = 2\pi/\lambda$ is the wave number of the incident plane harmonic wave $\mathbf{E}_i = \mathbf{e}_x E$, $E = E_0 \exp(ikz)$ propagating along the z axis. The *vector-scattering amplitude* is written as \mathbf{X} as a reminder that the incident wave is linearly polarized along the x axis. Here and elsewhere the time-dependent factor $\exp(-i\omega t)$ is omitted.

The extinction cross section is related in a simple way to the scattering amplitude;

$$C_{\text{ext}} = \frac{4\pi}{k^2} \text{Re}\{(\mathbf{X} \cdot \mathbf{e}_x)_{\theta=0}\} \quad (8)$$

This remarkable result, often called the *optical theorem*, implies that plane-wave extinction depends only on scattering in the forward direction $\theta = 0$, which seems to contradict the interpretation of extinction as the sum of scattering in *all* directions and absorption. Yet extinction has two interpretations, the second manifest in the optical theorem: extinction is interference between incident and forward-scattered waves.

The scattering cross section is also obtained from the vector-scattering amplitude by an integration over all directions:

$$C_{\text{sca}} = \int_{4\pi} \frac{|\mathbf{X}|^2}{k^2} d\Omega \quad (9)$$

At wavelengths far from strong absorption bands, the scattering cross section of a particle small compared with the wavelength satisfies (approximately)

$$C_{\text{sca}} \propto \frac{v^2}{\lambda^4} \quad (ka \rightarrow 0) \quad (10)$$

where a is a characteristic linear dimension of the particle. This result was first obtained by Lord Rayleigh in 1871 by dimensional analysis (his paper is included in Ref. 8).

The extinction cross section of a particle large compared with the wavelength approaches the limit

$$C_{\text{ext}} \rightarrow 2G \quad (ka \rightarrow \infty) \quad (11)$$

The fact that C_{ext} approaches *twice* G instead of G is sometimes called the *extinction paradox*. This alleged paradox arises from the expectation that geometrical optics should become a better approximation as a particle becomes larger. But all particles have edges because of which extinction by them always has a component unaccounted for by geometrical optics. This additional component, however, may not be observed because it is associated with light scattered very near the forward direction and because all detectors have finite acceptance angles. Measured extinction is theoretical extinction reduced by the scattered light collected by the detector.

No particle scatters light equally in all directions; isotropic scatterers exist only in the dreams of inept theorists. The angular dependence of scattering can be specified by the *differential scattering cross section*, written symbolically as $dC_{\text{sca}}/d\Omega$ as a reminder that the *total* scattering cross section is obtained from it by integrating over all directions:

$$C_{\text{sca}} = \int_{4\pi} \frac{dC_{\text{sca}}}{d\Omega} d\Omega \quad (12)$$

The normalized differential scattering cross section p

$$p = \frac{1}{C_{\text{sca}}} \frac{dC_{\text{sca}}}{d\Omega} \quad (13)$$

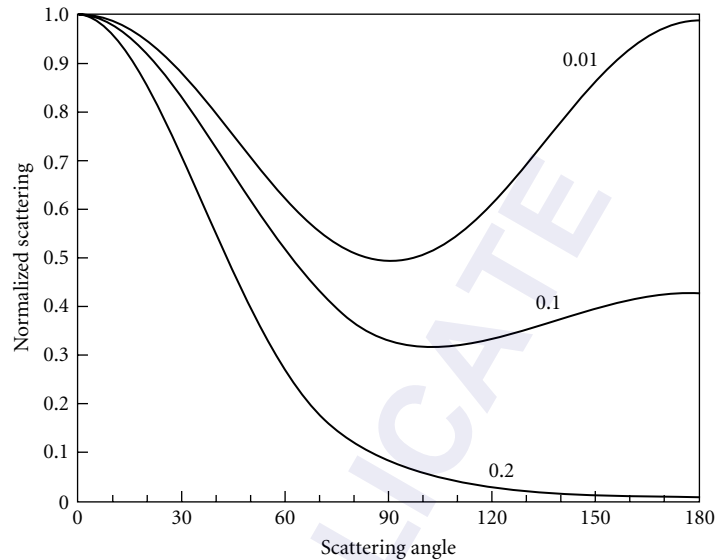


FIGURE 3 Scattering of unpolarized visible light by spheres of radii 0.01, 0.1, and 0.2 μm calculated according to the Rayleigh-Gans approximation.

is sometimes called the *phase function*. This coinage of astronomers (after the phases of astronomical bodies) confuses those who are perplexed by phase attached to a quantity from which phase in the usual optical sense is absent. To add to the confusion, the phase function is sometimes normalized to 4π instead of to unity.

A salient characteristic of scattering by particles is strong forward-backward asymmetry. Small metallic particles at far infrared wavelengths provide one of the few examples in which backscattering is larger than forward scattering. Except for very small particles, scattering is peaked in the forward direction; the larger the particle, the sharper the peak. Examples are given in Fig. 3, which shows differential scattering cross sections for unpolarized visible light illuminating spheres of various radii. These curves were obtained using the Rayleigh-Gans approximation,^{1,3,5} valid for particles optically similar to the surrounding medium. Forward scattering is much greater than backscattering even for a sphere as small as 0.2 μm .

A simple explanation of forward-backward asymmetry follows from the model of a scatterer as an array of N antennas. If we ignore mutual excitation (the antennas are excited solely by the external source), the total scattered field is the sum of N fields, the phases of which, in general, are different except in the forward direction. Scattering by noninteracting scatterers in this direction—and only in this direction—is in-phase regardless of their separation and the wavelength of the source. Thus as N increases, the scattered irradiance increases more rapidly in the forward direction than in any other direction.

Particles are miniature polarizers and retarders: they scatter differently the orthogonal components into which incident fields can be resolved. Similarly, an optically smooth surface can be both a polarizer and retarder. Just as polarization changes upon reflection are described by decomposing electric fields into components parallel and perpendicular to the plane of incidence, it is convenient to introduce a *scattering plane*, defined by the directions of the incident and scattered waves, for describing scattering by particles.

The incident plane wave is transverse, as is the scattered field at large distances. Thus these fields can be decomposed into two orthogonal components, one parallel, the other perpendicular to the scattering plane. The orthonormal basis vectors are denoted by \mathbf{e}_{\parallel} and \mathbf{e}_{\perp} and form a right-handed

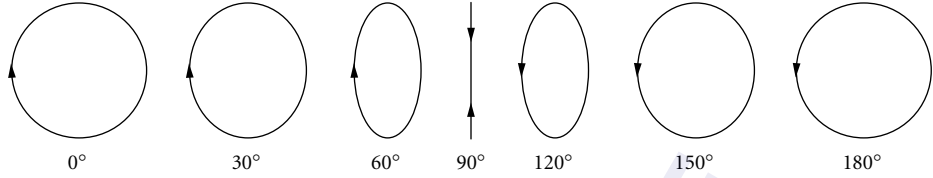


FIGURE 4 Vibration ellipses at various scattering angles for light scattered by a sphere much smaller than the wavelength of the incident right-circularly polarized light.

triad with the direction of propagation \mathbf{e}_p (of either the incident or scattered waves): $\mathbf{e}_\perp, \mathbf{x}, \mathbf{e}_\parallel = \mathbf{e}_p$. Incident and scattered fields are specified relative to different basis vectors. With this decomposition the relation between fields can be written^{1,5}

$$\begin{pmatrix} E_{\parallel s} \\ E_{\perp s} \end{pmatrix} = \frac{e^{ik(r-z)}}{-ikr} \begin{pmatrix} S_2 & S_3 \\ S_4 & S_1 \end{pmatrix} \begin{pmatrix} E_{\parallel i} \\ E_{\perp i} \end{pmatrix} \quad (14)$$

where i and s denote incident and scattered, respectively. The elements of this *amplitude scattering matrix* (or Jones matrix) are complex-valued functions of the scattering direction.

If a single particle is illuminated by completely polarized light, the scattered light is also completely polarized but possibly differently from the incident light, and differently in different directions. An example is given in Fig. 4, which shows vibration ellipses of light scattered by a small sphere. The polarization state of the scattered light varies from right-circular (the polarization state of the incident light) in the forward direction, to linear (perpendicular to the scattering plane) at 90°, to left-circular in the backward direction.

Just as unpolarized light can become partially polarized upon specular reflection, scattering of unpolarized light by particles can yield partially polarized light varying in degree and state of polarization in different directions. Unlike specular reflection, however, an ensemble of particles can transform completely polarized incident light into partially polarized scattered light if all the particles are not identical.

Transformations of polarized light upon scattering by particles are described most conveniently by the *scattering matrix* (or Mueller matrix) relating scattered to incident Stokes parameters:^{1,5}

$$\begin{pmatrix} I_s \\ Q_s \\ U_s \\ V_s \end{pmatrix} = \frac{1}{k^2 r^2} \begin{pmatrix} S_{11} & S_{12} & S_{13} & S_{14} \\ S_{21} & S_{22} & S_{23} & S_{24} \\ S_{31} & S_{32} & S_{33} & S_{34} \\ S_{41} & S_{42} & S_{43} & S_{44} \end{pmatrix} \begin{pmatrix} I_i \\ Q_i \\ U_i \\ V_i \end{pmatrix} \quad (15)$$

The scattering matrix elements S_{ij} for a single particle are functions of the amplitude-scattering matrix elements. Only seven of these elements are independent, corresponding to the four amplitudes and three phase differences of the S_j .

The scattering matrix for an ensemble of particles is the sum of matrices for each of them provided they are separated by sufficiently large random distances. Although all 16 matrix elements for an ensemble can be nonzero and different, symmetry reduces the number of matrix elements. For example, the scattering matrix for a rotationally and mirror symmetric ensemble has the form

$$\begin{pmatrix} S_{11} & S_{12} & 0 & 0 \\ S_{12} & S_{22} & 0 & 0 \\ 0 & 0 & S_{33} & S_{34} \\ 0 & 0 & -S_{34} & S_{44} \end{pmatrix} \quad (16)$$

7.5 SCATTERING BY AN ISOTROPIC, HOMOGENEOUS SPHERE: THE ARCHETYPE

An isotropic, homogeneous sphere is the simplest finite particle, the theory of scattering to which is attached the name of Gustav Mie.¹⁷ So firm is this attachment that in defiance of logic and history every particle under the sun has been dubbed a “Mie scatterer,” and Mie scattering has been promoted from a particular theory of limited applicability to the unearned rank of general scattering process.

Mie was not the first to solve the problem of scattering by an arbitrary sphere.¹⁸ It would be more correct to say that he was the last. He gave his solution in recognizably modern notation and also addressed a real problem: the colors of colloidal gold. For these reasons, his name is attached to the sphere-scattering problem even though he had illustrious predecessors, most notably Lorenz.¹⁹ This is an example in which eponymous recognition has gone to the last discoverer rather than to the first.

Mie scattering is not a physical process; Mie theory is one theory among many. It isn’t even exact because it is based on continuum electromagnetic theory and on illumination by a plane wave infinite in lateral extent.

Scattering by a sphere can be determined using various approximations and methods bearing little resemblance to Mie theory: Fraunhofer theory, geometrical optics, anomalous diffraction, coupled-dipole method, T -matrix method, etc. Thus, is a sphere a Mie scatterer or an anomalous diffraction scatterer or a coupled-dipole scatterer? The possibilities are endless. When a physical process can be described by several different theories, it is inadvisable to attach the name of one of them to it.

There is no distinct boundary between so-called Mie and Rayleigh scatterers. Mie theory includes Rayleigh theory, which is a limiting theory strictly applicable only as the size of the particle shrinks to zero. Even for spheres uncritically labeled “Rayleigh spheres,” there are always deviations between the Rayleigh and Mie theories. By hobbling one’s thinking with a supposed sharp boundary between Rayleigh and Mie scattering, one risks throwing some interesting physics out the window. Whether a particle is a Mie or Rayleigh scatterer is not absolute. A particle may be graduated from Rayleigh to Mie status merely by a change of wavelength of the illumination.

One often encounters statements about Mie scattering by cylinders, spheroids, and other non-spherical particles. Judged historically, these statements are nonsense: Mie never considered any particles other than homogeneous spheres.

Logic would seem to demand that if a particle is a Mie scatterer, then Mie theory can be applied to scattering by it. This fallacious notion has caused and will continue to cause mischief, and is probably the best reason for ceasing to refer to Mie particles or Mie scatterers. Using Mie theory for particles other than spheres is risky, especially for computing scattering toward the backward direction.

More often than not, a better term than Mie or Rayleigh scattering is available. If the scatterers are molecules, molecular scattering is better than Rayleigh scattering (itself an imprecise term):²⁰ the former term refers to an agent, the latter to a theory. Mie scatterer is just a needlessly aristocratic name for a humble sphere. Wherever Mie scatterer is replaced with sphere, the result is clearer. If qualifications are needed, one can add small or large compared with the wavelength or comparable to the wavelength.

Briefly, the solution to the problem of scattering by an arbitrary homogeneous sphere illuminated by a plane wave can be obtained by expanding the incident, scattered, and internal fields in a series of vector-spherical harmonics. The coefficients of these expansion functions are chosen so that the tangential components of the electric and magnetic fields are continuous across the surface of the sphere. Thus this scattering problem is formally identical to reflection and refraction because of interfaces, although the sphere problem is more complicated because the scattered and internal fields are not plane waves.

Observable quantities are expressed in terms of the coefficients a_n and b_n in the expansions of the scattered fields. For example, the cross sections are infinite series:

$$C_{\text{ext}} = \frac{2\pi}{k^2} \sum_{n=1}^{\infty} (2n+1) \text{Re}\{a_n + b_n\} \quad (17)$$

$$C_{\text{sca}} = \frac{2\pi}{k^2} \sum_{n=1}^{\infty} (2n+1) (|a_n|^2 + |b_n|^2) \quad (18)$$

If the permeability of the sphere and its surroundings are the same, the scattering coefficients can be written

$$a_n = \frac{[D_n(mx)/m+n/x]\psi_n(x) - \psi_{n-1}(x)}{[D_n(mx)/m+n/x]\xi_n(x) - \xi_{n-1}(x)} \quad (19)$$

$$b_n = \frac{[mD_n(mx)+n/x]\psi_n(x) - \psi_{n-1}(x)}{[mD_n(mx)+n/x]\xi_n(x) - \xi_{n-1}(x)} \quad (20)$$

ψ_n and ξ_n are Riccati-Bessel functions and the logarithmic derivative

$$D_n(\rho) = \frac{d}{d\rho} \ln \psi_n(\rho) \quad (21)$$

The *size parameter* x is ka , where a is the radius of the sphere and k is the wavenumber of the incident light in the surrounding medium, and m is the complex refractive index of the sphere relative to that of this (nonabsorbing) medium. Equations (19) and (20) are one of the many ways of writing the scattering coefficients, some of which are more suited to computations than others.

During the Great Depression mathematicians were put to work computing tables of trigonometric and other functions. The results of their labors now gather dust in libraries. Today, these tables could be generated more accurately in minutes on a pocket calculator. A similar fate has befallen Mie calculations. Before fast computers were inexpensive, tables of scattering functions for limited ranges of size parameter and refractive index were published. Today, these tables could be generated in minutes on a personal computer. The moral is to give algorithms rather than only tables of results, which are mostly useless except as checks for someone developing and testing algorithms.

These days it is not necessary to reinvent the sphere: documented Mie programs are readily available. The first widely circulated program was published as an IBM report by Dave in 1968, although it no longer seems to be available. A Mie program is given in Ref. 5. Reference 6 includes a diskette containing scattering programs for spheres (and other particles). Wiscombe^{21,22} suggested techniques for increasing the speed of computations, as did Lentz,²³ whose method makes use of continued fractions. Wang and van de Hulst²⁴ recently compared various scattering programs.

The primary tasks in Mie calculations are computing the functions in Eqs. (19) and (20) and summing series like Eqs. (17) and (18). Bessel functions are computed by recurrence. The logarithmic derivative, the argument of which can be complex, is usually computed by downward recurrence. $\psi_n(x)$ and $\xi_n(x)$ can be computed by upward recurrence if one does not generate more orders than are needed for convergence, approximately the size parameter x . When a program with no logical errors falls ill, it often can be cured by promoting variables from single to double precision.

Cross sections versus radius or wavelength convey physical information; efficiencies versus size parameter convey mathematical information. The size parameter is a variable with less physical content than its components, the whole being less than the sum of its parts. Moreover, cross section versus size parameter (or its inverse) is not equivalent to cross section versus wavelength. Except in the fantasy world of naive modelers, refractive indices vary with wavelength, and the Mie coefficients depend on x and m , wavelength being explicit in the first and implicit in the second.

The complex refractive index is written dozens of different ways, one of which is $n + ik$ (despite the risk of confusing the imaginary part with the wavenumber). The quantities n and k are called *optical constants*. But just as the Lord Privy Seal is neither a lord nor a privy nor a seal, optical constants are neither optical nor constant.

Few quantities in optics are more shrouded in myth and misconception than the complex refractive index. The real part for any medium is often defined as the ratio of the velocity of light c in free space to the phase velocity in the medium. This definition, together with notions that nothing can go faster than c , has engendered the widespread misconception that n must be greater than unity. But n can take on any value, even zero. The phase velocity is not the velocity of any palpable object or of any signal, hence is not subject to speed limits enforced by the special relativity police. The least physically relevant property of a refractive index is that it is a ratio of phase velocities. A refractive

index is a response function (or better, is simply related to response functions such as permittivity and permeability): it is a macroscopic manifestation of the microscopic response of matter to a periodic driving force.

When we turn to the imaginary part of the refractive index, we enter a ballroom in which common sense is checked at the door. It has been asserted countless times that an imaginary index of, say, 0.01 corresponds to a weakly absorbing medium (at visible and near-visible wavelengths). Such assertions are best exploded by expressing k in a more physically transparent way. The absorption coefficient α is

$$\alpha = \frac{4\pi k}{\lambda} \quad (22)$$

The inverse of α is the *e-folding distance* (or skin depth), the distance over which the irradiance of light propagating in an unbounded medium decreases by a factor of e . At visible wavelengths, the *e-folding distance* corresponding to $k = 0.01$ is about 5 μm . A thin sliver of such an allegedly weakly absorbing material would be opaque.

When can a particle (or any object) be said to be strongly absorbing? A necessary condition is that $\alpha d \gg 1$, where d is a characteristic linear dimension of the object. But this condition is not sufficient. As k increases, absorption increases—up to a point. As k approaches infinity, the absorption cross section of a particle or the absorptance of a film approaches zero.

One of the most vexing problems in scattering calculations is finding optical constants dispersed throughout dozens of journals. Palik²⁵ edited a compilation of optical constants for several solids. The optical constants of liquid water over a broad range were compiled by Hale and Querry;²⁶ Warren²⁷ published a similar compilation for ice. For other materials, you are on your own. Good hunting!

For small x and $|m|x$, the extinction and scattering efficiencies of a sphere are approximately

$$Q_{\text{ext}} = 4x \operatorname{Im} \left\{ \frac{m^2 - 1}{m^2 + 2} \right\} \quad (23)$$

$$Q_{\text{sca}} = \frac{8}{3} x^4 \left| \frac{m^2 - 1}{m^2 + 2} \right|^2 \quad (24)$$

These equations are the source of a nameless paradox, which is disinterred from time to time, a corpse never allowed eternal peace. If the sphere is nonabsorbing (m real), Eq. (23) yields a vanishing extinction cross section, whereas Eq. (24) yields a nonvanishing scattering cross section. Yet extinction never can be less than scattering. But note that Eq. (23) is only the first term in the expansion of Q_{ext} in powers of x . The first nonvanishing term in the expansion of Q_{sca} is of order x^4 . To be consistent, Q_{ext} and Q_{sca} must be expanded to the same order in x . When this is done, the paradox vanishes.

The amplitude-scattering matrix elements for a sphere are

$$S_1 = \sum_n \frac{2n+1}{n(n+1)} (a_n \pi_n + b_n \tau_n) \quad (25)$$

$$S_2 = \sum_n \frac{2n+1}{n(n+1)} (a_n \tau_n + b_n \pi_n) \quad (26)$$

where the angle-dependent functions are

$$\pi_n = \frac{P_n^1}{\sin \theta} \quad \tau_n = \frac{dP_n^1}{d\theta} \quad (27)$$

and P_n^1 are the associated Legendre functions of the first kind. The off-diagonal elements of the amplitude-scattering matrix vanish, because of which the scattering matrix is block-diagonal

and $S_{12} = S_{21}$, $S_{43} = -S_{34}$, $S_{44} = S_{33}$. Thus, when the incident light is polarized parallel (perpendicular) to the scattering plane, so is the scattered light, a consequence of the sphere's symmetry.

7.6 SCATTERING BY REGULAR PARTICLES

The field scattered by any spherically symmetric particle has the same form as that scattered by a homogeneous, isotropic sphere; only the scattering coefficients are different. One such particle is a uniformly coated sphere. Scattering by a sphere with a single layer was first treated by Aden and Kerker.²⁸ Extending their analysis to multilayered spheres is straightforward.²⁹

New computational problems arise in going from uncoated to coated spheres. The scattering coefficients for both contain spherical Bessel functions, which are bounded only if their arguments are real (no absorption). Thus, for strongly absorbing particles, the arguments of Bessel functions can be so large that their values exceed computational bounds. This does not occur for uncoated spheres because the only quantity in the scattering coefficients with complex argument is the logarithmic derivative, a ratio of Bessel functions computed as an entity instead of by combining numerator and denominator, each of which separately can exceed computational bounds. It is not obvious how to write the scattering coefficients for a coated sphere so that only ratios of possibly large quantities are computed explicitly. For this reason the applicability of the coated-sphere program in Ref. 5 is limited. Toon and Ackerman,³⁰ however, cast the coated-sphere coefficients in such a way that this limitation seems to have been surmounted.

Bessel functions of large complex argument are not the only trap for the unwary. A coated sphere is two spheres. The size parameter for the outer sphere determines the number of terms required for convergence of series. If the inner sphere is much smaller than the outer, the various Bessel functions appropriate to the inner sphere are computed for indices much greater than needed. More indices are not always better. Beyond a certain number, round-off error can accumulate to yield terms that should make ever smaller contributions to sums but may not.

Scattering by spheres and by infinitely long circular cylinders illuminated normally to their axes are in some ways similar. Spherical Bessel functions in the sphere-scattering coefficients correspond to cylindrical Bessel functions in the cylinder-scattering coefficients. Unlike a sphere, however, an infinitely long cylinder cannot be enclosed in a finite volume. As a consequence, the field scattered by such a cylinder decreases inversely as the square root of distance r instead of inversely as r (for sufficiently large r).

Infinite particles may be mathematically tractable but they are physically unrealizable. In particular, cross sections for infinite cylinders are infinite. But cross sections per unit length of infinite cylinders are finite. Such cross sections may be applied to a finite cylindrical particle by multiplying its length by the cross section per unit length of the corresponding infinite particle. If the aspect ratio (length/diameter) of the finite particle is sufficiently large, what are vaguely called "end effects" may be negligible. Because no exact theory for a finite cylinder exists, the aspect ratio at which differences between finite and infinite cylinders become negligible is not known with certainty, although the value 10 is bruited about. Nevertheless, there always will be differences between scattering by finite and infinite particles, which may or may not be of concern depending on the application.

A physical difference between scattering by spheres and by cylinders is that cross sections for cylinders depend on the polarization state of the incident plane wave. But normally incident light illuminating an infinite cylinder and polarized perpendicular (parallel) to the plane defined by the incident wave and the cylinder axis excites only scattered light polarized perpendicular (parallel) to the plane defined by the scattered wave and the cylinder axis. Obliquely incident linearly polarized light can, however, excite scattered light having both copolarized and cross-polarized components.

Obliquely illuminated uncoated cylinders pose no special computational problems. Coated cylinders, however, pose the same kinds of problems as coated spheres and are even more difficult to solve. Toon and Ackerman's³⁰ algorithm for coated spheres is based on the fact that spherical Bessel functions can be expressed in a finite number of terms. Because cylindrical Bessel functions cannot

be so expressed, this algorithm cannot be extended to coated cylinders, for which Bessel functions must be computed separately rather than as ratios and can have values beyond computational bounds. Even if such bounds are not exceeded, problems still can arise.

Although Barabás³¹ discussed in detail scattering by coated cylinders, Salzman and Bohren³² found that his computational scheme is unsuitable when absorption is large. They attempted, with only partial success, to write programs for *arbitrary*-coated cylinders. Care must be taken in computing Bessel functions. The often-used Miller algorithm can be inadequate for large, complex arguments.

The simplest nonspherical, finite particle is the spheroid, prolate, or oblate. Because the scalar wave equation is separable in spheroidal coordinates, scattering by spheroids can be solved in the same way as for spheres and cylinders. The expansion functions are based on spheroidal rather than spherical or cylindrical wave functions. Asano and Yamamoto³³ were the first to solve in this way the problem of scattering by an arbitrary spheroid. Although Asano³⁴ subsequently published an extensive set of computations based on this solution, it has not seen widespread use, possibly because of the intractability and exoticness of spheroidal functions.

Computational experience with spheroids and even simpler particles such as spheres and cylinders leads to the inescapable conclusion that hidden barriers lie between a mathematical solution to a scattering problem and an algorithm for reliably and quickly extracting numbers from it.

7.7 COMPUTATIONAL METHODS FOR NONSPHERICAL PARTICLES

The widespread notion that randomly oriented nonspherical particles are somehow equivalent to spheres is symptomatic of a failure to distinguish between the symmetry of an ensemble and that of its members. Considerable effort has been expended in seeking prescriptions for equivalent spheres. This search resembles that for the Holy Grail—and has been as fruitless.

From extensive studies of scattering by nonspherical particles, Mugnai and Wiscombe³⁵ concluded that “after examining hundreds of nonspherical results and observing that they all cluster relatively close together, relatively far from the equivolume spheres (except at forward angles), we have come to regard nonspherical particles as normal, and spheres as the most unrepresentative shape possible—almost a singularity.” This serves as a warning against using Mie theory for particles of all shapes and as a spur to finding methods more faithful to reality. We now turn to some of these methods. Keep in mind that no matter how different they may appear on the surface, they are all linked by the underlying Maxwell equations.

The *T-matrix method* is based on an integral formulation of scattering by an arbitrary particle. It was developed by Waterman, first for a perfect conductor,³⁶ then for a particle with less restricted properties.³⁷ It subsequently was applied to scattering problems under the name *extended boundary condition method* (EBCM).³⁸ Criticism of the *T-matrix method* was rebutted by Varadan et al.,³⁹ who cite dozens of papers on this method applied to electromagnetic scattering. Another source of papers and references is the collection edited by Varadan and Varadan.⁴⁰ Reference 6 is accompanied by a diskette containing *T-matrix* programs.

Linearity of the field equations and boundary conditions implies that the coefficients in the spherical harmonic expansion of the field scattered by any particle are linearly related to those of the incident field. The linear transformation connecting these two sets of coefficients is called the *T* (for transition) matrix.

The *T-matrix* elements are obtained by numerical integration. Computational difficulties arise for particles with high absorption or large aspect ratios. These limitations of the original *T-matrix* method have been surmounted somewhat by Iskander et al.,⁴¹ whose extension is dubbed the iterative extended boundary condition method.

Although the *T-matrix* method is not restricted to axisymmetric particles, it almost exclusively has been applied to spheroids and particles defined by Chebyshev polynomials.^{35,42,43}

Despite its virtues, the *T-matrix* method is not readily grasped in one sitting. Another method, variously called the Purcell-Pennypacker,⁴⁴ coupled-dipole,⁴⁵ digitized Green's function⁴⁶ method and

discrete dipole approximation,⁴⁷ is mathematically much simpler—the most complicated function entering into it is the exponential—and physically transparent. Although originally derived by heuristic arguments, the coupled-dipole method was put on firmer analytical foundations by Lakhtakia.⁴⁸

In this method, a particle is approximated by a lattice of N dipoles small compared with the wavelength but still large enough to contain many molecules. The dipoles often are, but need not be, identical and isotropic. Each dipole is excited by the incident field and by the fields of all the other dipoles. Thus the field components at each site satisfy a set of $3N$ linear equations. These components can be calculated by iteration^{44,49} or by inverting the $3N \times 3N$ coefficient matrix.⁴⁵ The coefficient matrix for only one particle orientation need be inverted. This inverse matrix then can be used to calculate scattering for other orientations.⁵⁰ A disadvantage of matrix inversion is that the number of dipoles is limited by computer storage.

Arrays of coupled dipoles were considered long before Purcell and Pennypacker entered the scene. More than half a century ago Kirkwood⁵¹ treated a dielectric as an array of molecules, the dipole moment of each of which is determined by the external field and by the fields of all the other molecules. What Purcell and Pennypacker did was to apply the coupled-dipole method to absorption and scattering by optically homogeneous particles. They bridged the gap between discrete arrays and continuous media with the Clausius-Mosotti theory. Because this theory, like every effective-medium theory, is not exact, critics of their method have pronounced it guilty by association. But the Clausius-Mosotti theory is merely the effective-medium theory that Purcell and Pennypacker happened to use. Whatever flaws their method may have, one of them is not that it is forever chained to the ghosts of Clausius and Mosotti. Alleged violations of the optical theorem are easily remedied by using the exact expression for the polarizability of a finite sphere,⁵² which in no way changes the structure of the method.

Draine⁴⁷ applied this method (under the label *discrete dipole approximation*) to extinction by interstellar grains, obtaining the field components with the conjugate gradient method. An outgrowth of his paper is that by Flatau et al.,⁵³ who considered scattering by rectangular particles. Goedecke and O'Brien⁴⁶ baptized their version of the digitized Green's function method and applied it to scattering of microwave radiation by snowflakes.⁵⁴ Varadan et al.⁵⁵ applied the method to scattering by particles with anisotropic optical constants. It also has been applied to scattering by helices,⁵⁶ by a cylinder on a reflecting surface,⁵⁷ and extended to intrinsically optically active particles.⁵⁸

Although Yung's analysis⁵⁹ of a large (15,600) array of dipoles representing a sphere suggests that there are no intrinsic limitations to the coupled-dipole method, it is plagued with practical limitations, most notably its inability to treat particles (especially compact ones or ones with large complex refractive indices) much larger than the wavelength of the illumination. Chiapetta,⁶⁰ then Singham and Bohren,⁶¹ reformulated the coupled-dipole method, expressing the total field at each dipole as the sum of the incident field and the fields scattered once, twice, and so on by all the other dipoles. Although this formulation is appealing because each term in the scattering-order series has a simple physical interpretation, the series can diverge. The greater the refractive index, the smaller the particle for which the series diverges. For a particle of given volume, fewer terms are needed for convergence the more the particle departs from sphericity. The greater the average separation between dipoles, the weaker the average interaction.

Except for improvements and refinements^{52,62,63} that increase accuracy and speed but do not remove barriers imposed by particle size and composition, the coupled-dipole method has not changed much since it first was used by Purcell and Pennypacker. It is not, of course, limited to optically homogeneous particles. It can be applied readily to aggregates of small particles. Indeed, it is best suited to aggregates with low fractal dimension. Berry and Percival⁶⁴ considered scattering by fractal aggregates using what they call the *mean-field approximation*, which is essentially the Rayleigh-Gans approximation, in turn a form of the scattering-order formulation of the coupled-dipole method in which the dipoles are excited only by the incident field.

The arbitrary border separating electrical engineering from optics is never more obvious than when it comes to methods for computing scattering. The engineers have theirs, the optical scientists have theirs, and rarely do the twain meet. In the hope of promoting smuggling, even illegal immigration, I must at least mention two methods that fall almost exclusively in the domain of electrical engineering: the method of moments and the finite-difference time-domain technique (FDTD).

Anyone interested in the method of moments must begin with Harrington's book,⁶⁵ a focal point from which paths fan out in all directions.

As its name implies, the FDTD technique is applied to what electrical engineers call the *time domain* (as opposed to the *frequency domain*, which most optical scientists inhabit even though they may not know it) and is explicitly labeled a finite-difference method (all methods for particles other than those of simple shape entail discretization in one form or another). Papers by Yee,⁶⁶ Holland et al.,⁶⁷ Mur,⁶⁸ Luebbers et al.,⁶⁹ and references cited in them will get you started on the FDTD technique.

When pondering the welter of species and subspecies of methods keep in mind that the differences among them and their ranges of validity are probably smaller than their adherents think or are willing to admit. There is no method that will happily compute scattering of arbitrary waves by particles of arbitrary size and composition in a finite amount of time. Moreover, each method, whatever its merits and demerits, often requires a tedious climb up a learning curve.

7.8 REFERENCES

1. H. C. van de Hulst, *Light Scattering by Small Particles*, Wiley, New York, 1957. Reprinted by Dover, New York, 1981.
2. D. Deirmendjian, *Electromagnetic Scattering on Polydispersions*, Elsevier, New York, 1969.
3. M. Kerker, *The Scattering of Light and Other Electromagnetic Radiation*, Academic, New York, 1969.
4. L. P. Bayvel and A. R. Jones, *Electromagnetic Scattering and Its Applications*, Elsevier, London, 1981.
5. C. F. Bohren and D. R. Huffman, *Absorption and Scattering of Light by Small Particles*, Wiley-Interscience, New York, 1983.
6. P. W. Barber and S. C. Hill, *Light Scattering by Particles: Computational Methods*, World Scientific, Singapore, 1990.
7. *Selected Papers on Light Scattering*, SPIE vol. 951, pts. 1 and 2, Milton Kerker (ed.), SPIE Optical Engineering Press, Bellingham, Wash., 1988.
8. *Selected Papers on Scattering in the Atmosphere*, SPIE vol. MS 7, Craig F. Bohren (ed.), SPIE Optical Engineering Press, Bellingham, Wash., 1989.
9. *Selected Papers on Natural Optical Activity*, SPIE vol. MS 15, Akhlesh Lakhtakia (ed.), SPIE Optical Engineering Press, Bellingham, Wash., 1990.
10. A. Lakhtakia, V. K. Varadan, and V. V. Varadan, *Time-Harmonic Electromagnetic Fields in Chiral Media*, Springer-Verlag, Berlin, 1989.
11. *Optical Particle Sizing: Theory and Practice*, Gérard Gouesbet and Gérard Gréhan (eds.), Plenum, New York, 1988.
12. *Optical Effects Associated with Small Particles*, P. W. Barber and R. K. Chang (eds.), World Scientific, Singapore, 1988.
13. W. H. Martin, "The Scattering of Light in One-Phase Systems," in *Colloid Chemistry: Theoretical and Applied*, vol. I, Jerome Alexander (ed.), Chemical Catalog Co., New York, 1926, pp. 340–352.
14. C. F. Bohren, "Multiple Scattering of Light and Some of Its Observable Consequences," *Am. J. Phys.* **55**:524 (1987).
15. W. T. Doyle, "Scattering Approach to Fresnel's Equations and Brewster's Law," *Am. J. Phys.* **53**:463 (1985).
16. Max Born and E. Wolf, *Principles of Optics*, 3d ed., Pergamon, Oxford, 1965.
17. P. Lilienfeld, "Gustav Mie: The Person," *Appl. Opt.* **30**:4696 (1991).
18. N. A. Logan, "Survey of Some Early Studies of the Scattering of Plane Waves by a Sphere," *Proc. IEEE* **53**:773 (1965).
19. Helge Kragh, "Ludvig Lorenz and Nineteenth Century Optical Theory: The Work of a Great Danish Scientist," *Appl. Opt.* **30**:4688 (1991).
20. A. T. Young, "Rayleigh Scattering," *Phys. Today* **35**:2 (1982).
21. W. J. Wiscombe, "Mie Scattering Calculations: Advances in Techniques and Fast, Vector-Speed Computer Codes," NCAR/TN-140 + STR, National Center for Atmospheric Research, Boulder, Colo., 1979.

22. W. J. Wiscombe, "Improved Mie Scattering Algorithms," *Appl. Opt.* **19**:1505 (1980).
23. W. J. Lentz, "Generating Bessel Functions in Mie Scattering Calculations Using Continued Fractions," *Appl. Opt.* **15**:668 (1976).
24. Ru T. Wang and H. C. van de Hulst, "Rainbows: Mie Computations and the Airy Approximation," *Appl. Opt.* **30**:106 (1991).
25. E. W. Palik (ed), *Handbook of Optical Constants of Solids*, Academic, New York, 1985.
26. G. M. Hale and M. Querry, "Optical Constants of Water in the 200-nm to 200- μm Wavelength Region," *Appl. Opt.* **12**:555 (1973).
27. S. G. Warren, "Optical Constants of Ice from the Ultraviolet to the Microwave," *Appl. Opt.* **23**:1206 (1984).
28. A. L. Aden and M. Kerker, "Scattering of Electromagnetic Waves from Two Concentric Spheres," *J. Appl. Phys.* **22**:601 (1951).
29. R. Bhandari, "Scattering Coefficients for a Multilayered Sphere: Analytic Expressions and Algorithms," *Appl. Opt.* **24**:1960 (1985).
30. O. B. Toon and T. P. Ackerman, "Algorithms for the Calculation of Scattering by Stratified Spheres," *Appl. Opt.* **20**:3657 (1981).
31. M. Barabás, "Scattering of a Plane Wave by a Radially Stratified Tilted Cylinder," *J. Opt. Soc. Am.* **A4**:2240 (1987).
32. G. C. Salzman and C. F. Bohren, "Scattering by Infinite Coated Cylinders Illuminated at Oblique Incidence," TCN 90203, U.S. Army Research Office, Research Triangle Park, N.C.
33. S. Asano and G. Yamamoto, "Light Scattering by a Spheroidal Particle," *Appl. Opt.* **14**:29 (1975).
34. S. Asano, "Light Scattering Properties of Spheroidal Particles," *Appl. Opt.* **18**:712 (1979).
35. A. Mugnai and W. Wiscombe, "Scattering from Nonspherical Chebyshev Particles, 3. Variability in Angular Scattering Patterns," *Appl. Opt.* **28**:3061 (1989).
36. P. C. Waterman, "Matrix Formulation of Electromagnetic Scattering," *Proc. IEEE* **53**:805 (1965).
37. P. C. Waterman, "Symmetry, Unitarity, and Geometry in Electromagnetic Scattering," *Phys. Rev.* **D3**:825 (1971).
38. P. Barber and C. Yeh, "Scattering of Electromagnetic Waves by Arbitrarily Shaped Dielectric Bodies," *Appl. Opt.* **14**:2864 (1975).
39. V. V. Varadan, A. Lakhtakia, and V. K. Varadan, "Comments on Recent Criticism of the T-Matrix Method," *J. Acoust. Soc. Am.* **84**:2280 (1988).
40. *Acoustic, Electromagnetic and Elastic Wave Scattering—Focus on the T-Matrix Approach*, V. K. Varadan and V. V. Varadan (eds.), Pergamon, New York, 1980.
41. M. F. Iskander, A. Lakhtakia, and C. H. Durney, "A New Procedure for Improving the Solution Stability and Extending the Frequency Range of the EBCM," *IEEE Trans. Antennas Propag.* **AP-31**:317 (1983).
42. A. Mugnai and W. Wiscombe, "Scattering from Nonspherical Chebyshev Particles, 1. Cross Sections, Single-Scattering Albedo, Asymmetry Factor, and Backscattered Fraction," *Appl. Opt.* **25**:1235 (1986).
43. W. J. Wiscombe and A. Mugnai, "Scattering from Nonspherical Chebyshev Particles, 2. Means of Angular Scattering Patterns," *Appl. Opt.* **27**:2405 (1988).
44. E. M. Purcell and C. R. Pennypacker, "Scattering and Absorption of Light by Nonspherical Dielectric Grains," *Astrophys. J.* **186**:705 (1973).
45. Shermila Brito Singham and G. C. Salzman, "Evaluation of the Scattering Matrix of an Arbitrary Particle Using the Coupled-Dipole Approximation," *J. Chem. Phys.* **84**:2658 (1986).
46. G. H. Goedecke and S. G. O'Brien, "Scattering by Irregular Inhomogeneous Particles via the Digitized Green's Function Algorithm," *Appl. Opt.* **27**:2431 (1988).
47. B. T. Draine, "The Discrete-Dipole Approximation and Its Application to Interstellar Graphite Grains," *Astrophys. J.* **333**:848 (1988).
48. Akhlesh Lakhtakia, "Macroscopic Theory of the Coupled Dipole Approximation Method," *Opt. Comm.* **79**:1 (1990).
49. S. D. Druger, M. Kerker, D.-S. Wang, and D. D. Cooke, "Light Scattering by Inhomogeneous Particles," *Appl. Opt.* **18**:3888 (1979).

50. M. K. Singham, S. B. Singham, and G. C. Salzman, "The Scattering Matrix for Randomly Oriented Particles," *J. Chem. Phys.* **85**:3807 (1986).
51. J. G. Kirkwood, "On the Theory of Dielectric Polarization," *J. Chem. Phys.* **4**:592 (1936).
52. C. E. Dungey and C. F. Bohren, "Light Scattering by Nonspherical Particles: A Refinement to the Coupled-Dipole Method," *J. Opt. Soc. Am.* **A8**:81 (1991).
53. P. J. Flatau, G. L. Stephens, and B. T. Draine, "Light Scattering by Rectangular Solids in the Discrete-Dipole approximation: A New Algorithm Exploiting the Block-Toeplitz Structure," *J. Opt. Soc. Am.* **A7**:593 (1990).
54. S. G. O'Brien and G. H. Goedecke, "Scattering of Millimeter Waves by Snow Crystals and Equivalent Homogeneous Symmetric Particles," *Appl. Opt.* **27**:2439 (1988).
55. V. V. Varadan, A. Lakhtakia, and V. K. Varadan, "Scattering by Three-Dimensional Anisotropic Scatterers," *IEEE Trans. Antennas Propag.* **37**:800 (1989).
56. S. B. Singham, C. W. Patterson, and G. C. Salzman, "Polarizabilities for Light Scattering from Chiral Particles," *J. Chem. Phys.* **85**:763 (1986).
57. M. A. Taubenblatt, "Light Scattering from Cylindrical Structures on Surfaces," *Opt. Lett.* **15**:255 (1990).
58. Shermila Brito Singham, "Intrinsic Optical Activity in Light Scattering from an Arbitrary Particle," *Chem. Phys. Lett.* **130**:139 (1986).
59. Y. L. Yung, "Variational Principle for Scattering of Light by Dielectric Particles," *Appl. Opt.* **17**:3707 (1978).
60. P. Chiapetta, "Multiple Scattering Approach to Light Scattering by Arbitrarily Shaped Particles," *J. Phys. A* **13**:2101 (1980).
61. Shermila Brito Singham and C. F. Bohren, "Light Scattering by an Arbitrary Particle: The Scattering Order Formulation of the Coupled-Dipole Method," *J. Opt. Soc. Am.* **A5**:1867 (1988).
62. Shermila Brito Singham and C. F. Bohren, "Hybrid Method in Light Scattering by an Arbitrary Particle," *Appl. Opt.* **28**:517 (1989).
63. J. J. Goodman, B. T. Draine, and P. J. Flatau, "Application of Fast-Fourier Transform Techniques to the Discrete-Dipole Approximation," *Opt. Lett.* **16**:1198 (1991).
64. M. V. Berry and I. C. Percival, "Optics of Fractal Clusters Such as Smoke," *Opt. Acta* **33**:577 (1986).
65. R. F. Harrington, *Field Computation by Moment Methods*, Robert E. Krieger, Malabar, Fla., 1987. Reprint of 1968 edition.
66. K. S. Yee, "Numerical Solution of Initial Boundary Value Problems Involving Maxwell's Equations in Isotropic Media," *IEEE Trans. Antennas Propag.* **AP14**:302 (1966).
67. R. Holland, L. Simpson, and K. S. Kunz, "Finite Difference Analysis of EMP Coupling to Lossy Dielectric Structures," *IEEE Trans. Electromag. Compat.* **EMC22**:203 (1980).
68. G. Mur, "Absorbing Boundary Conditions for the Finite Difference Approximation of the Time Domain Electromagnetic Field Equations," *IEEE Trans. Electromag. Compat.* **EMC23**:377 (1981).
69. R. Luebbers, F. P. Hunsberger, K. S. Kunz, R. B. Standler, and M. Schneider, "A Frequency-Dependent Finite-Difference Time-Domain Formulation for Dispersive Materials," *IEEE Trans. Electromag. Compat.* **EMC32**:222 (1990).

This page intentionally left blank.

DO NOT DUPLICATE

Eugene L. Church and Peter Z. Takacs

*Instrumentation Division
Brookhaven National Laboratory
Upton, New York*

8.1 GLOSSARY OF PRINCIPAL SYMBOLS

A_o	illuminated area
A, B, C	model parameters
BRDF	bidirectional reflectance distribution function
$C(\tau)$	autocovariance function or ACF
$D(\tau)$	structure function
f_{xy}, f_x	spatial frequencies
g, g'	Strehl parameter or Rayleigh index
J_n	ordinary Bessel function
K_n	modified Bessel function
L_o	illuminated length
L_x, L_y	lengths of illuminated area
N	fractal index
P	incident or scattered power
R_α, r_α	Fresnel reflection coefficients
RCS	radar cross-section
S	two-sided power spectral density (PSD)
S^s	one-sided profile PSD
T	topothesy
$Z(x, y)$	topographic surface roughness
$Z(x)$	surface profile
1D, 2D	dimensionalities
$\langle \dots \rangle$	ensemble average
α, β	linear-polarization parameters
ρ	correlation length
σ	root-mean-square roughness
τ	lag variable

8.2 INTRODUCTION

Imperfect surface finish degrades optical performance. The connection between the two is electromagnetic scattering theory, which is the subject of this chapter. Because of space limitations we limit our attention to the elemental case of simple highly reflective surfaces such as those used in x-ray imaging and other high-performance applications.

This chapter is divided into two parts. The first part discusses elementary physical optics and first-order perturbation theory of the scattering of electromagnetic waves by a randomly rough surface, ending with Eq. (35).

The second part discusses the interpretation of scattering measurements in terms of surface finish models, and the related question of surface-finish specification. Because many manufactured surfaces show an inverse-power-law behavior, we have concentrated on models that show this behavior. Two types of surfaces are considered throughout—1D or grating-like surfaces, and isotropic 2D surfaces.

The field of topographic surface scattering has been highly developed in radio physics, optics, and instrumentation over the past 50 years. There is a staggering volume of literature in each area, expressed in a bewildering variety of notations, which are summarized in Sec. 8.3. Important in-depth texts in these fields are the early reviews of Beckmann and Spizzichino¹ and the *Radar Cross-Section Handbook*², the more recent work of Ishimaru³, Fung⁴, the textbooks of Nieto-Vesperinas⁵ and Voronovich⁶, and most recently, the publications of Maradudin et al.⁷

Scattering measurements are discussed in the works of Stover⁸ and Germer⁹, and profile measurements using various optical techniques, including the Long Trace Profiler (LTP) are considered by Takacs et al.¹⁰

A number of important subjects have been omitted in the present review. These include instrumentation, statistical estimation, the effects of detrending, figure error, standards, and multilayer surfaces. Discussions of some of these and related subjects can be found in Chap. 11, “Analog Optical and Image Processing,” in this volume and Chap. 4, “Imaging through Atmospheric Turbulence,” Chap. 44, “Reflective Optics,” and Chap. 46, “X-Ray Mirror Metrology,” in Vol. V.

8.3 NOTATION

The scattering geometry is sketched in Fig. 1. Note that the angles of incidence, θ_i and θ_f , are always positive, and that the azimuthal angle φ_i is understood to be $\varphi_f - \varphi_i$ with $\varphi_i = 0$. Specular reflection occurs at $\theta_f = \theta_i$ and $\varphi_f = 0$, and backscatter at $\theta_f = \theta_i$ and $\varphi_f = \pi$.

The initial and final wavenumbers are

$$\mathbf{k}_i = \frac{2\pi}{\lambda} \begin{pmatrix} +\sin(\theta_i) \\ 0 \\ -\cos(\theta_i) \end{pmatrix} \quad \mathbf{k}_f = \frac{2\pi}{\lambda} \begin{pmatrix} +\sin(\theta_f) \cos(\varphi_f) \\ +\sin(\theta_f) \sin(\varphi_f) \\ +\cos(\theta_f) \end{pmatrix} \quad (1)$$

and the spatial frequency vectors are

$$\mathbf{f} = \frac{\mathbf{k}_f - \mathbf{k}_i}{2\pi} = \begin{pmatrix} f_x \\ f_y \\ f_z \end{pmatrix} = \frac{1}{\lambda} \begin{pmatrix} \sin(\theta_f) \cos(\varphi_f) - \sin(\theta_i) \\ \sin(\theta_f) \sin(\varphi_f) \\ \cos(\theta_f) + \cos(\theta_i) \end{pmatrix} \quad \mathbf{f}_{xy} = \begin{pmatrix} f_x \\ f_y \end{pmatrix} \quad (2)$$

These can be viewed as a generalization of the grating equation for first-order diffraction from a grating with the spatial wavelength $d = 1/f$.

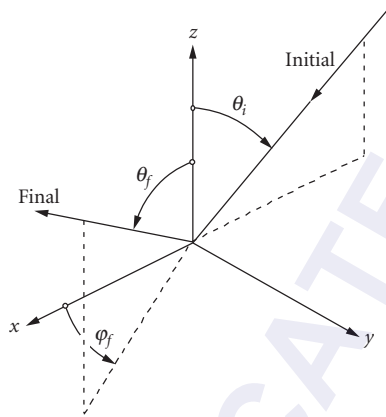


FIGURE 1 Coordinate system for 1D and 2D scattering. The x - y plane is the surface plane, with air/vacuum above and reflective material below. The planes of incidence and scattering are defined by the z axis and the initial and final directions of propagation.

The Jacobian relating configuration, wavenumber, and frequency space is

$$\cos(\theta_f) d\omega_f = \frac{1}{k^2} dk_x dk_y = \lambda^2 df_x df_y, \quad d\omega_f = \sin(\theta_f) d\theta_f d\phi_f \quad (3)$$

The two orthogonal states of linear polarization have a variety of notations in the literature:

$$\begin{aligned} s(\text{senkrecht}) &= h(\text{horizontal}) = \text{TE} = E \perp (\text{perpendicular}) \\ &\rightarrow \text{Dirichlet (soft) boundary condition} \\ p(\text{parallel}) &= v(\text{vertical}) = \text{TM} = H \parallel (\text{parallel}) \\ &\rightarrow \text{Neumann (hard) boundary condition} \end{aligned} \quad (4)$$

Designation of initial (α) to final (β) states of polarization:

$$A_{\beta\alpha} = A_{\alpha \rightarrow \beta}$$

We use the optical or “p, s” notation hereafter. Circular and elliptically polarized results are obtained by taking appropriate linear combinations of these linear forms. Unpolarized results are obtained by summing over the final states of polarization and averaging over the initial states.

The polarization vector is a unit vector in the direction of the electric vector. A natural sign convention is

$$\hat{\mathbf{s}} = \hat{\mathbf{z}} \times \hat{\mathbf{k}} \quad \hat{\mathbf{p}} = \hat{\mathbf{k}} \times \hat{\mathbf{s}} \quad \hat{\mathbf{k}} = \hat{\mathbf{s}} \times \hat{\mathbf{p}} \quad (5)$$

but this is not universal.

One can make a distinction between scattering from a surface that is rough in both the x and y directions, and a grating like surface that is rough only along the x axis. The first generates a 2D or bistatic scattering pattern in the upper hemisphere, while the second scatters only into an infinitesimal region about the plane of incidence. Randomly-rough 1D surfaces are of special interest for research purposes since they are easier to fabricate with prescribed statistical properties than are 2D surfaces.^{11–13} In the following discussions include 1D and 2D results side by side, wherever practical.

The distribution of the reflected and scattered power has a wide variety of notations in the optics and radar literature. For a 2D scatterer

$$\begin{aligned}
\frac{1}{P_i} \left(\frac{dP}{d\omega_f} \right)_{\alpha \rightarrow \beta}^{(2D)} &= \frac{1}{A_o \cos(\theta_i)} \cdot \lim_{r \rightarrow \infty} \left[r^2 \left| \frac{E_f}{E_i} \right|^2 \right] = \frac{|f_{\beta\alpha}^{(2D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i)|^2}{A_o \cos(\theta_i)} \\
&= \frac{\sigma_{\beta\alpha}^{(2D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i)}{A_o \cos(\theta_i)} = \frac{\text{RCS}_{\beta\alpha}^{(2D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i)}{4\pi A_o \cos(\theta_i)} = \frac{\gamma_{\beta\alpha}^{(2D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i)}{4\pi \cos(\theta_i)} \\
&= \text{DSC}_{\alpha \rightarrow \beta}^{(2D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i) = \cos(\theta_f) \cdot \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i) \\
&= \cos(\theta_i) \frac{A_o}{\lambda^2} D\{\rho\}^{(2D)} = I_{\alpha\beta}(Q)^{(2D)}
\end{aligned} \tag{6}$$

and for a 1D scatterer

$$\begin{aligned}
\frac{1}{P_i} \left(\frac{dP}{d\theta_f} \right)_{\alpha \rightarrow \beta}^{(1D)} &= \frac{1}{L_o \cos(\theta_i)} \cdot \lim_{r \rightarrow \infty} \left[r \left| \frac{E_f}{E_i} \right|^2 \right] = \frac{2}{\pi k} \frac{|f_{\beta\alpha}^{(1D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i)|^2}{L_o \cos(\theta_i)} \\
&= \frac{\sigma_{\beta\alpha}^{(1D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i)}{L_o \cos(\theta_i)} = \frac{\text{RCS}_{\beta\alpha}^{(1D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i)}{2\pi L_o \cos(\theta_i)} = \frac{\gamma_{\beta\alpha}^{(1D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i)}{2\pi \cos(\theta_i)} \\
&= \text{DSC}_{\alpha \rightarrow \beta}^{(1D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i) = \cos(\theta_f) \cdot \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)}(\hat{\mathbf{k}}_f, \hat{\mathbf{k}}_i) \\
&= \cos(\theta_i) \frac{L_o}{\lambda} D\{\rho\}^{(1D)} = I_{\alpha\beta}(Q)^{(1D)}
\end{aligned} \tag{7}$$

In Eqs. (6) and (7), P_i is the power incident on the surface; that is, $A_o \cos(\theta_i) |E_i|^2 / (2\eta_o)$, where A_o is the total illuminated area of the surface for a 2D scatterer, and $A_o = 1 \times L_o$ for a 1D or grating-like scatterer. The f 's here are the scattering amplitudes for spherical waves in 2D and cylindrical waves in 1D. σ is the bistatic cross-section, RCS is the bistatic radar cross-section, γ is the bistatic scattering coefficient, DSC is the differential scattering cross-section, and BRDF is the all-important radiometric quantity, the bidirectional reflectance distribution function. This is the quantity we will focus on in the detailed calculations described later in the chapter.

The cross-sections σ and RCS have the dimensions of area and length or “width.” In practice, these are frequently converted to dimensionless forms by normalizing them to the appropriate power of the radiation wavelength or some characteristic physical dimension of the scattering object. On the other hand, γ and BRDF are dimensionless. The former is more appropriate for isolated objects, and the latter for distributed surface roughness. The γ 's used here are the radar cross-section normalized to the illuminated surface area or length, although alternative definitions appear in the literature.

In Eq. (7), $D\{\rho\}$ and $I(Q)$ are quantities used in the books of Beckmann and Spizzichino¹ and Nieto-Vesperinas.⁵

At this point we do not distinguish deterministic and ensemble-average quantities, or coherent and incoherent scattering. These distinctions are discussed with respect to specific model calculations later in the chapter.

These full bistatic-scattering expressions may contain more information than is needed in practice. For example, scatterometry measurements^{8,9} are usually performed in the plane of incidence ($\varphi_f = 0$ or π), and surface isotropy is checked by rotating the sample to intermediate values of φ between measurements. Similarly, the radar literature is frequently concerned only with measurements in the backscatter (retroscatter) direction, ($\theta_f = \theta_i$, $\varphi_f = \pi$).

In principle, relating the BRDF to topographic surface features is a straightforward application of Maxwell's equations, but approximations are frequently necessary to get practically useful results. The best known of these approximations are the physical-optics (Fresnel-Kirchhoff) and the small-perturbation (Rayleigh-Rice) methods. These are discussed in Secs. 8.4 and 8.5.

8.4 THE FRESNEL-KIRCHHOFF APPROXIMATION

Introduction

The Fresnel-Kirchhoff or physical-optics approximation is also known as the tangent-plane or Huygens wave approximation. It is valid when the surface roughness has a large radius of curvature, that is, it is locally flat, but places no direct restrictions on the surface height or slope. It is an inherently paraxial approximation, and in its simplest form discussed below [Eq. (8)] omits multiple scattering and shadowing effects. As a result, it does not generally satisfy the conservation of energy.

In this approximation the ensemble-average value of the BRDF is in 2D:

$$\begin{aligned} \langle \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)} \rangle^{\text{FK}} &= \frac{1}{\lambda^2} \cdot \frac{1}{\cos(\theta_i) \cos(\theta_f)} \cdot Q_{\alpha \rightarrow \beta}^{\text{FK}} \cdot \langle \mathfrak{S}(\mathbf{f}_{xy})^{(2D)} \rangle \\ \langle \mathfrak{S}(\mathbf{f}_{xy})^{(2D)} \rangle &= \int_{-\infty}^{+\infty} d\boldsymbol{\tau}_{xy} \exp[i2\pi \mathbf{f}_{xy} \cdot \boldsymbol{\tau}_{xy}] \cdot \exp\left[-\frac{1}{2}(2\pi f_z)^2 \langle D(\boldsymbol{\tau}_{xy}) \rangle\right] \end{aligned} \quad (8a)$$

and in 1D:

$$\begin{aligned} \langle \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)} \rangle^{\text{FK}} &= \frac{1}{\lambda} \cdot \frac{1}{\cos(\theta_i) \cos(\theta_f)} \cdot Q_{\alpha \rightarrow \beta}^{\text{FK}} \cdot \langle \mathfrak{S}(f_x)^{(1D)} \rangle \\ \langle \mathfrak{S}(f_x)^{(1D)} \rangle &= \int_{-\infty}^{+\infty} d\tau_x \exp[i2\pi f_x \tau_x] \cdot \exp\left[-\frac{1}{2}(2\pi f_z)^2 \langle D(\tau_x) \rangle\right] \end{aligned} \quad (8b)$$

In these expressions the Q carries the information on the polarization of the radiation and the properties of the surface material, and the real quantity $\langle \mathfrak{S} \rangle$ carries the information about the statistical properties of the surface roughness. This occurs through the dependence on the ensemble- average surface structure function, $\langle D(\boldsymbol{\tau}) \rangle$

$$\langle D(\boldsymbol{\tau}_{xy}) \rangle = \langle D(|\boldsymbol{\tau}_{xy}|) \rangle = \langle [Z(\mathbf{r}_{xy} + \boldsymbol{\tau}_{xy}) - Z(\mathbf{r}_{xy})]^2 \rangle \quad (9)$$

where \mathbf{r}_{xy} is the position vector in the surface plane. The existence of $\langle D(\boldsymbol{\tau}) \rangle$ requires that the surface roughness have statistically stationary first differences. In the limit of a perfectly smooth surface, $D(\boldsymbol{\tau})$ vanishes and the \mathfrak{S} 's become a delta function in the specular direction.

The exponential dependence on the roughness structure function in Eq. (8) is a consequence of the usual assumption that the height fluctuations, $Z(x, y)$, have a gaussian bivariate distribution. Nongaussian effects, which become manifest only for rough surfaces, are discussed in the literature.^{11,13}

The polarization-materials factor can be written approximately as

$$\begin{aligned} Q_{\alpha \rightarrow \beta}^{\text{FK}} &= \left| A_{\alpha \rightarrow \beta}^{\text{FK}} \right|^2 \cdot R_{\alpha \rightarrow \beta}(\boldsymbol{\theta}) \\ A_{s \rightarrow s}^{\text{FK}} &= -A_{p \rightarrow p}^{\text{FK}} = \frac{(1 + \cos(\theta_i) \cos(\theta_f)) \cos(\varphi_f) - \sin(\theta_i) \sin(\theta_f)}{\cos(\theta_i) + \cos(\theta_f)} \\ A_{s \rightarrow p}^{\text{FK}} &= -A_{p \rightarrow s}^{\text{FK}} = \sin(\varphi_f) \end{aligned} \quad (10)$$

Although these results are correct for a perfect reflector, $R = 1$, the dependence on the surface reflectivity for an imperfect reflector is quite complicated.^{2,4} Here we follow custom by arbitrarily factoring out a reflection coefficient R in the first line, where the R values are the Fresnel intensity reflectivities,

$$R_{\alpha \rightarrow \alpha}(\theta) = |r_\alpha(\theta)|^2 \quad R_{\alpha \rightarrow \beta}(\theta) = \left| \frac{1}{2}(r_\alpha(\theta) - r_\beta(\theta)) \right|^2 \quad (11)$$

and the r 's are the amplitude reflection coefficients

$$r_s(\theta) = \frac{\mu \cos(\theta) - \sqrt{\mu \varepsilon - \sin^2(\theta)}}{\mu \cos(\theta) + \sqrt{\mu \varepsilon - \sin^2(\theta)}} \quad r_p(\theta) = \frac{\varepsilon \cos(\theta) - \sqrt{\mu \varepsilon - \sin^2(\theta)}}{\varepsilon \cos(\theta) + \sqrt{\mu \varepsilon - \sin^2(\theta)}} \quad (12)$$

For simplicity in presentation we cavalierly evaluate these reflectivity factors at the local angle of incidence. That is, $R(\theta) = R(\theta_{\text{loc}})$, where $\theta_{\text{loc}} = |\theta_f \cos(\varphi_f) + \theta_i|/2$ in the plane of incidence.²

The material parameters ε and μ are the electric permeability and magnetic permittivity of the surface relative to vacuum. In the case of a perfect electrical conductor (PEC), $\varepsilon \rightarrow \infty$ and the R 's are unity.

In the forward-scattering direction, $\varphi_f^2 \ll 1$, the Q 's simplify to

$$Q_{s \rightarrow s}^{\text{FK}} = Q_{p \rightarrow p}^{\text{FK}} = \left[\frac{\cos((\theta_f + \theta_i)/2)}{\cos((\theta_f - \theta_i)/2)} \right]^2 \cdot R_{\alpha \rightarrow \alpha}(\theta_{\text{loc}}) \quad Q_{s \rightarrow p}^{\text{FK}} = Q_{p \rightarrow s}^{\text{FK}} = 0 \quad (13)$$

and a related but different expression for $Q_{\alpha\alpha}$ in the retro-scattering direction, $\varphi_f = \pi$.

Statistically Stationary Surfaces

A random variable is statistically stationary if it has a finite mean-square value or variance, $\sigma^2 = \langle Z^2 \rangle$.¹⁴ In that case the structure function can be written as

$$\langle D(\tau_{xy}) \rangle = 2[\sigma^2 - \langle C(\tau_{xy}) \rangle] \quad \sigma^2 = \langle C(0) \rangle \quad (14)$$

$$\langle C(\tau_{xy}) \rangle = \langle C(|\tau_{xy}|) \rangle = \langle Z(\mathbf{r}_{xy} + \tau_{xy}) Z(\mathbf{r}_{xy}) \rangle$$

where $\langle C(\tau) \rangle$ is the autocovariance function. The scattering integrals in Eq. (8) can then be written as

$$\begin{aligned} \langle \mathfrak{S}(\mathbf{f}_{xy})^{(2D)} \rangle &= e^{-g'} \int_{-\infty}^{+\infty} d\tau_{xy} \exp[i2\pi \mathbf{f}_{xy} \cdot \tau_{xy}] \cdot \exp[(2\pi f_z)^2 \langle C(\tau_{xy}) \rangle] \\ \langle \mathfrak{S}(f_x)^{(1D)} \rangle &= e^{-g'} \int_{-\infty}^{+\infty} d\tau_x \exp[i2\pi f_x \tau_x] \cdot \exp[(2\pi f_z)^2 \langle C(\tau_x) \rangle] \\ g' = (2\pi f_z \sigma)^2 &= \left[4\pi \left(\frac{\cos(\theta_i) + \cos(\theta_f)}{2} \right) \frac{\sigma}{\lambda} \right]^2 \quad g = \left[4\pi \cos(\theta_i) \frac{\sigma}{\lambda} \right]^2 \end{aligned} \quad (15)$$

The quantity g , the Strehl or Rayleigh index, is an important dimensionless measure of the degree of surface roughness^{1,15} which will appear throughout the following discussions.

We now examine forms of the $\langle \text{BRDF} \rangle$ that follow from Eqs. (8) and (15) in four roughness regimes.

Perfectly Smooth Surfaces In perfectly smooth surfaces, where $g = 0$, Eq. (8) reduces to

$$\langle \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)} \rangle^{\text{FK}} = \frac{1}{\lambda^2} \cdot R_\alpha(\theta_i) \cdot \delta(\mathbf{f}_{xy}) \cdot \delta_{\alpha\beta} \quad (16)$$

$$\langle \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)} \rangle^{\text{FK}} = \frac{1}{\lambda} \cdot R_\alpha(\theta_i) \cdot \delta(f_x) \cdot \delta_{\alpha\beta}$$

that is, a sharp spike in the specular direction.

This spike is a delta function in spatial-frequency space and follows from the assumption of an infinite illumination area in Eq. (8). Later, in Sec. 8.6, we extend this idealized result to include the important effects of a finite-sized illumination area.

Slightly Rough Surfaces In slightly rough surfaces, $0 < g \ll 1$. In this case, expand the second exponent in Eq. (15) in a power series and keep the first two terms. These two terms separate the $\langle \text{BRDF} \rangle$ into “coherent” and “incoherent” parts:

$$\begin{aligned} \left\langle \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)} \right\rangle_{\text{coherent}}^{\text{FK}} &= (1-g) \cdot \frac{1}{\lambda^2} \cdot R_{\alpha}(\theta_i) \cdot \delta(\mathbf{f}_{xy}) \cdot \delta_{\alpha,\beta} \\ \left\langle \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)} \right\rangle_{\text{incoherent}}^{\text{FK}} &= \frac{16\pi^2}{\lambda^4} \frac{[\cos(\theta_i) + \cos(\theta_f)]^2}{4\cos(\theta_i)\cos(\theta_f)} \cdot \left| A_{\alpha \rightarrow \beta}^{\text{FK}} \right|^2 \cdot R_{\alpha,\beta}(\theta_{\text{loc}}) \cdot \langle S(\mathbf{f}_{xy})^{(2D)} \rangle \end{aligned} \quad (17a)$$

and

$$\begin{aligned} \left\langle \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)} \right\rangle_{\text{coherent}}^{\text{FK}} &= (1-g) \cdot \frac{1}{\lambda} \cdot R_{\alpha}(\theta_i) \cdot \delta(f_x, 0)_{f_x,0} \cdot \delta_{\alpha,\beta} \\ \left\langle \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)} \right\rangle_{\text{incoherent}}^{\text{FK}} &= \frac{16\pi^2}{\lambda^3} \frac{[\cos(\theta_i) + \cos(\theta_f)]^2}{4\cos(\theta_i)\cos(\theta_f)} \cdot \left| A_{\alpha \rightarrow \beta}^{\text{FK}} \right|^2 \cdot R_{\alpha}(\theta_{\text{loc}}) \cdot \delta_{\alpha,\beta} \cdot \langle S(f_x)^{(1D)} \rangle \end{aligned} \quad (17b)$$

The statistical properties of the surface appear here in the all-important quantities $\langle S^{(1D)} \rangle$ and $\langle S^{(2D)} \rangle$ —the ensemble-averages of the 1D and 2D power spectral densities of the surface roughness. These are the Fourier transforms of the corresponding covariance functions:

$$\begin{aligned} \langle S(f_x)^{(1D)} \rangle &= \int_{-\infty}^{+\infty} d\tau_x \exp[i2\pi f_x \tau_x] \cdot \langle C(\tau_x) \rangle \\ \langle S(\mathbf{f}_{xy})^{(2D)} \rangle &= \int_{-\infty}^{+\infty} d\tau_{xy} \exp[i2\pi \mathbf{f}_{xy} \cdot \tau_{xy}] \cdot \langle C(\tau_{xy}) \rangle \end{aligned} \quad (18)$$

These are discussed in detail in Sec. 8.7.

As a final note, the incoherent terms reduce to somewhat simpler forms when observations are made in the forward direction, $\varphi_f^2 \ll 1$:

$$\begin{aligned} \left\langle \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)} \right\rangle_{\text{incoherent}}^{\text{FK}} &= \frac{16\pi^2}{\lambda^4} \frac{\cos^4[(\theta_f + \theta_i)/2]}{\cos(\theta_i)\cos(\theta_f)} \cdot \delta_{\alpha,\beta} \cdot R(\theta_{\text{loc}}) \cdot \langle S(f_x, 0)^{(2D)} \rangle \\ \left\langle \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)} \right\rangle_{\text{incoherent}}^{\text{FK}} &= \frac{16\pi^2}{\lambda^3} \frac{\cos^4[(\theta_f + \theta_i)/2]}{\cos(\theta_i)\cos(\theta_f)} \cdot \delta_{\alpha,\beta} \cdot R(\theta_{\text{loc}}) \cdot \langle S(f_x)^{(1D)} \rangle \end{aligned} \quad (19)$$

Moderately Rough Surfaces In moderately rough surfaces, where $g \approx 1$, the power-series expansion of the second exponential in Eq. (15) cannot be cut off after the first two terms, and higher terms must be included. Beckmann and Spizzichino¹ illustrate this procedure for a gaussian $\langle C(\tau) \rangle$ and a gaussian bivariate distribution (see also Refs. 11 and 13).

A gaussian autocovariance function is rarely, if ever, observed in polished optical surfaces, although it is very convenient for analytic and experimental investigations. For example, O’Donnell and Mendez have made artificially rough grating-like surfaces of this type by superimposing speckle patterns, and have used them to make ingenious scattering studies of rough surfaces.¹¹

Very Rough Surfaces In very rough surfaces, where $g \gg 1$, the coherent term vanishes for a gaussian bivariate distribution, and we revert to Eq. (8) which then involves the limit of the scattering integral:

$$\lim_{\langle D(\tau_{xy}) \rangle \rightarrow \infty} \left\{ \int_{-\infty}^{+\infty} d\tau_{xy} \exp[i2\pi \mathbf{f}_{xy} \cdot \tau_{xy}] \cdot \exp \left[-\frac{1}{2} (2\pi f_z)^2 \langle D(\tau_{xy}) \rangle \right] \right\} \quad (20)$$

which is determined by the indicial behavior of $\langle D(\tau_{xy}) \rangle$. In the special case where the structure function is isotropic and quadratic

$$\lim_{\tau \rightarrow 0} \langle D(\tau_{xy}) \rangle = \sigma_M^2 \tau_{xy}^2 \quad (21)$$

where σ_M is the dimensionless root mean square (rms) value of the surface gradient of a 2D surface or the rms value of surface slope of a 1D surface. It is easy to see that this quadratic dependence leads to values of $\langle \mathfrak{S} \rangle$ that are proportional to λ^2 in 2D and λ^1 in 1D, in which case the $\langle \text{BRDF} \rangle$ is independent of the radiation wavelength. In other words, a quadratic structure function leads to geometric-optics results.

In fact, the resulting scattering pattern is a mapping of the slope distribution of the surface roughness, including the doubling of the deflection angle on reflection.^{1,2} The form of the scattering pattern depends on the form of the bivariate distribution involved. For example, a gaussian bivariate distribution leads to a gaussian pattern, while a gamma distribution leads to scattering patterns involving modified Bessel functions.¹³

If the indicial behavior of $\langle D \rangle$ is not quadratic, the form of the $\langle \text{BRDF} \rangle$'s depend on λ and the elegant geometrical-optics limits are not achieved. Mathematically, this occurs because such surfaces are not differentiable at $\tau = 0$ and so have no well-defined "slope."

Fractal Surfaces

Introduction Fractal surfaces have structure functions of the form

$$\langle D(\tau_{xy}) \rangle_{\text{fractal}} = T^2 \left| \frac{\tau_{xy}}{T} \right|^N \quad 0 < N < 2 \quad (22)$$

where T is a length parameter called the topothesy.^{16,17} Physically, T is the separation of surface points whose connecting chord has an rms slope of unity. Because Eq. (22) is proportional to T^{2-N} where $N < 2$, a perfectly smooth surface occurs when $T \rightarrow 0$. On the other hand, because $N = 2$ is excluded, very rough fractal surfaces do not lead to a geometrical-optics result.

The number N is the fractal index, which is related to the Hausdorff or Hausdorff-Besicovitch dimension D' ,^{16,17}

$$D' = (4 - N)/2 \quad N = 4 - 2D' \quad (23)$$

$$1 < D' < 2 \quad 0 < N < 2$$

D' can also be expressed in terms of the Hurst dimension or coefficient, but the connection depends on the dimensionality of the problem.¹⁸

Fractal Forms of the Scattering Integrals In the case of statistically isotropic surfaces, structure functions of the form of Eq. (22) lead to expressions for the incoherent scattering integrals which may be written as

$$\langle \mathfrak{S}(\mathbf{f}_{xy})^{(2D)} \rangle = 2\pi \int_0^{\infty} \tau_{xy} d\tau_{xy} J_0(2\pi f_{xy} \tau_{xy}) \cdot \sigma^2 \exp \left[- \left| \frac{\tau_{xy}}{\rho} \right|^N \right] \quad (24)$$

$$\langle \mathfrak{S}(f_x)^{(1D)} \rangle = \int_{-\infty}^{+\infty} d\tau_x \exp [i2\pi f_x \tau_x] \cdot \sigma^2 \exp \left[- \left| \frac{\tau_x}{\rho} \right|^N \right]$$

where

$$\sigma = 1 \text{ (dimensionless)} \quad \rho = \left[\frac{2}{(2\pi f_z T)^2} \right]^{1/N} \cdot T \quad (25)$$

are “pseudo” roughness and correlation length parameters. These do not play true physical roles in the case of fractals, but are artificial quantities introduced here to enable us to write Eq. (24) in form of a Fourier transform of a covariance function that will appear later as the power exponential (PEX) model.

The integrals in Eq. (24) are symmetric bell-shaped functions of f_{xy} and f_x that are flat at low spatial frequencies and fall off with inverse-power-law tails at high spatial frequencies. The 2D expression cannot be expressed in terms of known functions, while the 1D form can be written in terms of centered symmetric Lévy stable distributions of order α , that is, $L_\alpha(X)$, where:¹⁷

$$\langle \mathfrak{S}(f_x)^{(1D)} \rangle = 2\pi\sigma^2\rho \cdot I_N(Y) \quad Y = 2\pi f_x\rho \quad (26)$$

Nolan¹⁹ gives the computer program “STABLE” for this in terms of general stable distribution functions f_{Nolan} :

$$\langle \mathfrak{S}(f_x)^{(1D)} \rangle = 2\pi\sigma^2\rho \cdot f_{\text{Nolan}}(x|\alpha, \beta, \gamma, \delta, k) \quad (27)$$

$$x = |f_x|, \quad \alpha = N, \quad \beta = 0, \quad \gamma = 1/(2\pi\rho), \quad \delta = 0, \quad k = 0 \text{ or } 1$$

On the other hand, the low- and high-frequency limits of the 1D and 2D forms can be expressed in simple closed form. In the 1D case:

$$\begin{aligned} \langle \mathfrak{S}(|f_x| \ll 1/2\pi\rho)^{(1D)} \rangle &= \sigma^2\rho \cdot 2\Gamma((1+N)/N) \\ \langle \mathfrak{S}(|f_x| \gg 1/2\pi\rho)^{(1D)} \rangle &= \sigma^2\rho \cdot 2^N \sqrt{\pi} \frac{\Gamma((1+N)/2)}{\Gamma((2-N)/2)} \cdot \frac{N}{(2\pi|f_x|\rho)^{N+1}} \end{aligned} \quad (28a)$$

and in the isotropic 2D case:

$$\begin{aligned} \langle \mathfrak{S}(|f_{xy}| \ll 1/2\pi\rho)^{(2D)} \rangle &= \sigma^2\rho^2 \cdot \pi \Gamma((2+N)/N) \\ \langle \mathfrak{S}(|f_{xy}| \gg 1/2\pi\rho)^{(2D)} \rangle &= \sigma^2\rho^2 \cdot 2^{N+1} \pi \frac{\Gamma((2+N)/2)}{\Gamma((2-N)/2)} \cdot \frac{N}{(2\pi|f_{xy}|\rho)^{N+2}} \end{aligned} \quad (28b)$$

It follows from the Fourier transform nature of Eq. (8) that the areas under each of these curves in frequency space is simply σ^2 . This plus the asymptotic properties of the $\langle \mathfrak{S} \rangle$'s given above capture most of the physical properties of the $\langle \text{BRDF} \rangle$ of interest.

The case $N = 1$, $D' = 3/2$ is called the Brownian fractal, which falls between the Cantor set ($D' = 0.63093 \dots$) and the Sierpinski gasket ($D' = 1.5850 \dots$). The Brownian fractal has the virtue of leading to the simple analytic expressions for the diffraction integrals valid for all spatial frequencies:

$$\langle \mathfrak{S}(f_x)^{(1D)} \rangle_{N=1} = \frac{2\sigma^2\rho}{1+(2\pi f_x\rho)^2} \quad \langle \mathfrak{S}(f_{xy})^{(2D)} \rangle_{N=1} = \frac{2\pi\sigma^2\rho^2}{[1+(2\pi f_{xy}\rho)^2]^{3/2}} \quad (29)$$

where, for fractal surfaces, σ and ρ are given by Eq. (25).

8.5 THE RAYLEIGH-RICE (RR) OR SMALL-PERTURBATION APPROXIMATION

Results

The small-perturbation method is an alternative to the Fresnel-Kirchhoff method discussed above. Its first-order form was originally derived by Rice²⁰ using the Rayleigh hypothesis, and hence the name Rayleigh-Rice. Peake^{2,21} was the first to derive the expression for an arbitrary

surface material and the results have been rederived many times in the literature. These Rayleigh-Rice results have been extended to higher orders²²⁻²⁴, and in its more general form is called the small-perturbation method.

The lowest-order perturbation theory results are²

$$\begin{aligned}\left\langle \text{BRDF}(\mathbf{f}_{xy})_{\alpha \rightarrow \beta}^{(2D)} \right\rangle_{\text{incoherent}}^{\text{RR}} &= \frac{16\pi^2}{\lambda^4} \cos(\theta_i) \cos(\theta_f) \cdot Q_{\alpha \rightarrow \beta}^{\text{RR}} \cdot \langle S(\mathbf{f}_{xy})^{(2D)} \rangle \\ \left\langle \text{BRDF}(f_x)_{\alpha \rightarrow \beta}^{(1D)} \right\rangle_{\text{incoherent}}^{\text{RR}} &= \frac{16\pi^2}{\lambda^3} \cos(\theta_i) \cos(\theta_f) \cdot Q_{\alpha \rightarrow \beta}^{\text{RR}} \cdot \delta_{\alpha, \beta} \cdot \langle S(f_x)^{(1D)} \rangle\end{aligned}\quad (30)$$

The RR form of the coherent term is complicated and can be found in the literature.^{5,25}

The Q 's are the material-polarization factors similar to those appearing in the FK calculations. In contrast with the approximation made in the FK case, however, they do not separate into distinct angular and reflectivity factors. In particular,²

$$\begin{aligned}Q_{\alpha \rightarrow \beta}^{\text{RR}} &= \left| A_{\alpha \rightarrow \beta}^{\text{RR}} \right|^2 \\ A_{s \rightarrow s}^{\text{RR}} &= \frac{[\mu - 1][B(\theta_i)B(\theta_f)\cos(\varphi_f) - \mu \sin(\theta_i)\sin(\theta_f)] - \mu^2(\varepsilon - 1)\cos(\varphi_f)}{[B(\theta_i) + \mu \cos(\theta_i)][B(\theta_f) + \mu \cos(\theta_f)]} \\ A_{s \rightarrow p}^{\text{RR}} &= \frac{\varepsilon(\mu - 1)B(\theta_i) - \mu(\varepsilon - 1)B(\theta_f)}{[B(\theta_i) + \mu \cos(\theta_i)][B(\theta_f) + \varepsilon \cos(\theta_f)]} \sin(\varphi_f) \\ A_{p \rightarrow s}^{\text{RR}} &= \frac{\mu(\varepsilon - 1)B(\theta_i) - \varepsilon(\mu - 1)B(\theta_f)}{[B(\theta_i) + \varepsilon \cos(\theta_i)][B(\theta_f) + \mu \cos(\theta_f)]} \sin(\varphi_f) \\ A_{p \rightarrow p}^{\text{RR}} &= \frac{(\varepsilon - 1)[\varepsilon \sin(\theta_i)\sin(\theta_f) - B(\theta_i)B(\theta_f)\cos(\varphi_f)] + \varepsilon^2(\mu - 1)\cos(\varphi_f)}{[B(\theta_i) + \varepsilon \cos(\theta_i)][B(\theta_f) + \varepsilon \cos(\theta_f)]}\end{aligned}\quad (31)$$

where $B(\theta) = \sqrt{\varepsilon\mu - \sin^2(\theta)}$ and index of refraction $= \sqrt{\varepsilon\mu}$. The full ε - μ dependencies are useful for checking the duality of the results.

The Q^{RR} 's are closely related to the Fresnel reflection coefficients in Eqs. (11) and (12). For non-magnetic materials, $\mu = 1$ and

$$\begin{aligned}Q_{s \rightarrow s}^{\text{RR}} &= \sqrt{R_s(\theta_i)R_s(\theta_f)} \cdot \cos^2(\varphi_f) \\ Q_{\alpha \rightarrow \beta}^{\text{RR}}(\text{specular}) &= R_\alpha(\theta_i) \cdot \delta_{\alpha, \beta}\end{aligned}\quad (32)$$

On the other hand, for a perfectly reflecting (PEC) surface,

$$\begin{aligned}Q_{s \rightarrow s}^{\text{RR}} &= \cos^2(\varphi_f) \quad Q_{s \rightarrow p}^{\text{RR}} = \left[\frac{\sin(\varphi_f)}{\cos(\theta_f)} \right]^2 \\ Q_{p \rightarrow p}^{\text{RR}} &= \left[\frac{\cos(\varphi_f) - \sin(\theta_i)\sin(\theta_f)}{\cos(\theta_i)\cos(\theta_f)} \right]^2 \quad Q_{p \rightarrow s}^{\text{RR}} = \left[\frac{\sin(\varphi_f)}{\cos(\theta_i)} \right]^2\end{aligned}\quad (33)$$

In the special case of a perfectly reflecting surface measured in the plane of incidence we get the elegant results:

$$\begin{aligned}
 \langle \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)} \rangle_{\text{incoherent}}^{\text{RR}} &= \frac{16\pi^2}{\lambda^4} \cos(\theta_i) \cos(\theta_f) \cdot Q_{\alpha \rightarrow \beta}^{\text{RR}} \cdot \langle S(f_x, 0)^{(2D)} \rangle \\
 \langle \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)} \rangle_{\text{incoherent}}^{\text{RR}} &= \frac{16\pi^2}{\lambda^3} \cos(\theta_i) \cos(\theta_f) \cdot Q_{\alpha \rightarrow \beta}^{\text{RR}} \cdot \delta_{\alpha, \beta} \cdot \langle S(f_x)^{(1D)} \rangle \\
 Q_{s \rightarrow s}^{\text{RR}} = 1 \quad Q_{s \rightarrow p}^{\text{RR}} = Q_{p \rightarrow s}^{\text{RR}} = 0 \quad Q_{p \rightarrow p}^{\text{RR}} &= \left(\frac{1 - \sin(\theta_i) \sin(\theta_f)}{\cos(\theta_i) \cos(\theta_f)} \right)^2
 \end{aligned} \tag{34}$$

which may be compared with the FK result in Eq. (19) with $R(\theta) = 1$.

Again, the reader is alerted to the fundamental distinction between $\langle S(f_x, 0)^{(2D)} \rangle$ and $\langle S(f_x)^{(1D)} \rangle$, appearing here. They may be mathematically related for an isotropically rough surface, but they are never equal.

Comparison of RR and FK Results

The Rayleigh-Rice results are inherently a smooth-surface approximation for a statistically stationary random surface, so that the proper comparison is with the Fresnel-Kirchhoff results for slightly rough surfaces in Eq. (17) and the RR results in Eq. (30).

In the limit of paraxial scattering, $\theta_f \approx \theta_i$ and $\varphi_f \approx 0$, the two sets of results become identical and can be written in the common form:

$$\begin{aligned}
 \langle \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)} \rangle_{\text{incoherent}}^{\text{FK, RR}} &\rightarrow \frac{16\pi^2}{\lambda^4} \cos^2(\theta_i) \cdot R_\alpha(\theta_i) \cdot \left\{ \begin{array}{l} 1 \quad \alpha = \beta \\ \left(\frac{\sin(\varphi_f)}{\cos(\theta_i)} \right)^2 \quad \alpha \neq \beta \end{array} \right\} \cdot \langle S(f_x, f_y)^{(2D)} \rangle \\
 \langle \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)} \rangle_{\text{incoherent}}^{\text{FK, RR}} &\rightarrow \frac{16\pi^2}{\lambda^3} \cos^2(\theta_i) \cdot R_\alpha(\theta_i) \cdot \delta_{\alpha, \beta} \cdot \langle S(f_x)^{(1D)} \rangle
 \end{aligned} \tag{35a}$$

where $R_\alpha(\theta_i)$ is given by Eq. (11).

In the case of fractal surfaces the paraxial results are

$$\begin{aligned}
 \langle \text{BRDF}_{\alpha \rightarrow \beta}^{(2D)} \rangle_{\text{incoherent}}^{\text{FK, RR}} &\rightarrow \frac{1}{\lambda^2} \cdot R_\alpha(\theta_i) \cdot \left\{ \begin{array}{l} 1 \quad \alpha = \beta \\ \left(\frac{\sin(\varphi_f)}{\cos(\theta_i)} \right)^2 \quad \alpha \neq \beta \end{array} \right\} \cdot \langle \mathfrak{S}(f_x, f_y)^{(2D)} \rangle \\
 \langle \text{BRDF}_{\alpha \rightarrow \beta}^{(1D)} \rangle_{\text{incoherent}}^{\text{FK, RR}} &\rightarrow \frac{1}{\lambda} \cdot R_\alpha(\theta_i) \cdot \delta_{\alpha, \beta} \cdot \langle \mathfrak{S}(f_x)^{(1D)} \rangle
 \end{aligned} \tag{35b}$$

where the $\langle \mathfrak{S} \rangle$'s are given in Eqs. (24) et seq. Later, in "The J-K model" section we discuss the high-frequency forms of the $\langle \text{BRDF} \rangle$ that follow from these results.

A nice feature of the paraxial results in Eq. (35) is that they satisfy the conservation of energy. In the case of statistically stationary surfaces this occurs via

$$\begin{aligned}
 \int d\omega_s \frac{dP}{d\omega_s} &= \frac{\lambda^2}{R(\theta_i)} \int d\mathbf{f}_{xy} \langle \text{BRDF}^{(2D)} \rangle_{\text{total}}^{\text{FK, RR}} = (1-g) + g = 1 \\
 \int d\theta_s \frac{dP}{d\theta_s} &= \frac{\lambda}{R(\theta_i)} \int df_x \langle \text{BRDF}^{(1D)} \rangle_{\text{total}}^{\text{FK, RR}} = (1-g) + g = 1
 \end{aligned} \tag{36}$$

where the terms $(1-g)$ and g come from the coherent and incoherent components.

The differences between the FK and RR results that appear at larger deflection angles are attributed to the inherently paraxial approximation of the FK calculations. In particular, the FK results for rougher surfaces do not satisfy the conservation of energy due to their neglect of multiple scattering and shadowing effects.²⁶

On the other hand, the RR or lowest-order perturbation theory results are known to have intrinsic limitations at grazing angles, especially for p-polarized radiation.²⁷ The RR results also may not satisfy the conservation of energy because of the phenomenon of roughness-induced absorption.²⁸

8.6 EFFECTS OF FINITE ILLUMINATION AREA

The discussion above has assumed that the illuminated surface area is infinite, which is the meaning of the limits $\pm \infty$ in the integrals, $\langle \mathfrak{S} \rangle$, appearing in Eqs. (8) et seq.

In general, the effects of the finite rectangular illumination area, $L_x \times L_y$, can be taken into account by convolving the infinite-illumination forms of the $\langle \text{BRDF} \rangle$ with the system-response (SR) or point-spread function

$$\text{System response}(f_x, f_y) = \frac{L_x L_y}{\lambda^2} \left[\frac{\sin(\pi f_x L_x)}{\pi f_x L_x} \right]^2 \left[\frac{\sin(\pi f_y L_y)}{\pi f_y L_y} \right]^2 \quad (37)$$

which becomes $\delta(f_x)\delta(f_y)$ in the limit $L_x, L_y \gg \lambda$.²⁹

When the L 's are much larger than the correlation length of the surface roughness, the incoherent scattering pattern is a broad and "smooth" function of spatial frequency and is unaffected by the convolution with the relatively sharp system response function.

This argument breaks down when there are sharp features in the $\langle \text{BRDF} \rangle$. An example of this occurs in the smooth-surface limit of statistically stationary surfaces, which exhibit a delta-function coherent-scattering peak as in Eq. (17). That delta function is then smeared into the expected sinc^2 pattern by convolution with the system response.

Fractal scattering does not display a separable coherent component but does become increasingly bunched in the specular direction when the topography is sufficiently small. In that limit the observed scattering is again affected significantly by convolution with the system response.

The non-vanishing width of the system response function plays an important role in the discussion of surface-finish specification in Sec. 8.8.

8.7 SURFACE STATISTICS

Second-Order Statistical Functions

The ensemble-average scattered power can be written in terms of three second-order statistical functions: the structure function, $\langle D(\tau) \rangle$, the surface autocovariance function (ACF), $\langle C(\tau) \rangle$, and the power spectral density (PSD) $\langle S(f) \rangle$. Equation (18) gives the PSD in terms of the ACF, but it can also be written directly in terms of the surface profile according to

$$\begin{aligned} \langle S(f_x)^{(1D)} \rangle &= \lim_{L_o \rightarrow \infty} \left\langle \left| \frac{1}{L_o} \int_{-L_o/2}^{+L_o/2} dx \exp[i2\pi f_x x] \cdot Z(x) \right|^2 \right\rangle \\ \langle S(\mathbf{f}_{xy})^{(2D)} \rangle &= \lim_{A_o \rightarrow \infty} \left\langle \left| \frac{1}{A_o} \int_{A_o} d\mathbf{r}_{xy} \exp[i2\pi \mathbf{f}_{xy} \cdot \mathbf{r}_{xy}] \cdot Z(\mathbf{r}_{xy}) \right|^2 \right\rangle \end{aligned} \quad (38)$$

where $-\infty < f_x, f_{xy} < +\infty$. This is the basis for the periodogram estimate of the one-sided profile spectrum³⁰:

$$\hat{S}^*(f_n) = 2 \cdot \frac{D}{N} \left| \sum_{m=0}^{N-1} \exp[i2\pi mn/N] \cdot W(m) \cdot Z(mD) \right|^2 K(m) \quad (39)$$

$$f_n = \frac{n}{ND} \quad n = 1, \dots, \frac{N}{2}$$

where the “hat” on S means that it is an estimate, and the asterisk and the factor of two on the right mean that the negative spatial frequencies have been folded over and added into the positive frequencies. $W(n)$ is a real bell-shaped window function that eliminates the ringing that would otherwise appear in the estimate due to the sharp edges of the data window, and $K(m)$ is a book-keeping factor that equals unity everywhere except at the end points, where it equals 1/2.

In the above, D is the sample spacing of the measured data, and the range of spatial frequencies included in the measurement is $(1/ND) \leq f_n \leq 1/2D$, where $1/2D$ is the Nyquist frequency of the measurement. As written, N is even, although a similar expression holds for odd N .

Properties of Power Spectra

The power spectra show how the variance of the surface roughness is distributed over surface spatial frequencies. In particular,

$$\langle \sigma^2 \rangle^{(2D)} = \langle Z(x, y)^2 \rangle = \int_{-\infty}^{+\infty} df_{xy} \langle S(f_{xy})^{(2D)} \rangle = \langle C(0) \rangle \quad (40)$$

$$\langle \sigma^2 \rangle^{(1D)} = \langle Z(x)^2 \rangle = \int_{-\infty}^{+\infty} df_x \langle S(f_x)^{(1D)} \rangle = \langle C(0) \rangle$$

In other words, the roughness of a statistically stationary surface measured over a surface area is the same as that measured over any linear profile across it.

The 2D spectrum is what appears in surface-scattering measurements, while 1D spectrum appears in surface profile measurements. Both depend only on the magnitude of their spatial-frequency arguments. However, 1D and 2D spectra are distinctly different—they even have different dimensions: $\langle S^{(2D)} \rangle$ is $[L^4]$, while $\langle S^{(1D)} \rangle$ is $[L^3]$. What is not true is that the 1D form is a simple slice of the 2D form; that is, $\langle S(f_x)^{(1D)} \rangle$ does not equal $\langle S(f_x, 0)^{(2D)} \rangle$. Instead, the 1D form can be derived from the 2D form by integration,

$$\langle S(f_x)^{(1D)} \rangle = \int_{-\infty}^{+\infty} df_y \langle S(f_x, f_y)^{(2D)} \rangle \quad (41)$$

but the 2D form cannot be derived from the 1D form without providing further information about the 2D form. This is usually given in terms of its symmetry properties.

Incidentally, it follows from Eqs. (35) and (41) that the 1D scattering pattern equals the 2D pattern integrated over an long, narrow slit parallel to the y axis in Fig. 1.

If the surface is statistically isotropic, $\langle S^{(1D)} \rangle$ and $\langle S^{(2D)} \rangle$ are related by the integral transforms

$$\langle S(f_{xy})^{(2D)} \rangle = -\frac{1}{\pi} \int_{f_{xy}}^{\infty} \frac{df_x}{\sqrt{f_x^2 - f_{xy}^2}} \frac{d}{df_x} \langle S(f_x)^{(1D)} \rangle$$

$$\langle S(f_x)^{(1D)} \rangle = 2 \int_{f_x}^{\infty} \frac{f_{xy} df_{xy}}{\sqrt{f_{xy}^2 - f_x^2}} \langle S(f_{xy})^{(2D)} \rangle \quad (42)$$

The first is the inverse-Abel or “half derivative” transform and the second is the Abel transform or “half integral” transform.³¹

The transforms in Eq. (42) are useful since they allow profile measurements, which are inherently 1D, to be translated into the 2D spectra which appear in scatterometry and practical applications.³² They also permit the transformation of the high-frequency behavior of the spectra of one dimensionality to be transformed into the high-frequency behavior of the other without knowledge of their low-frequency behavior.³²

An important example of this is the case of inverse-power-law high-frequency tails,

$$\begin{aligned} \langle S(f_{xy})^{(2D)} \rangle &= \frac{J}{f_{xy}^m} & \langle S(f_x)^{(1D)} \rangle &= \frac{K}{|f_x|^{m-1}} \\ K &= \frac{\Gamma[1/2]\Gamma[(m-1)/2]}{\Gamma[m/2]} \cdot J & m > 2 \end{aligned} \quad (43)$$

which is the essence of the so-called J-K model discussed in the “The J-K Model” section.

Finish Models

General Remarks The magnitudes of the spatial frequencies $|f|$ appearing in the PSDs discussed above cover the range from 0 to $+\infty$. Real world measurements, however, include only a lesser range of spatial frequencies, $f_{\min} < |f| < f_{\max}$, where f_{\min} and f_{\max} are determined by the details of the measurement process.

In scatterometry the bandwidth limits are determined by the radiation wavelength and the maximum and minimum collection angles according to Eq. (2). In profilometry the minimum spatial frequency is the reciprocal of the trace length of the measurement, and the maximum is the reciprocal of twice the uniform sampling interval—that is, the Nyquist frequency of the measurement.

Surface finish models are parametric models that are fitted to experimental data. This condenses the measured data into a set of discrete finish parameters, smooths the measured data within the measurement bandpass, and—depending on one’s degree of trust in the model—can be used to extrapolate the data outside the measurement range.

In this chapter we consider four elementary models: The fractal model discussed above, and the ABC, PEX and J-K models considered below. These models have been chosen since each shows an inverse-power-law high-frequency tail displayed by many real surfaces.

The Fractal Model The fractal model is defined by the structure function $D(\tau) = T^2 |\tau/T|^N$ appearing in Eqs. (8) and (9), et seq. This two-parameter “model” follows from the geometrical scaling of fractal or self-affine surface roughness. Since there is no intrinsic limitation on its degree of roughness in this case, it is more properly described by the FK rather than the RR calculations, leading to the results given in the “Fractal Surfaces” section.

The ABC Model A very useful pair of spectra that satisfy the integral transforms in Eq. (42) is the “ABC” model:

$$\begin{aligned} \langle S(f_x)^{(1D)} \rangle &= \frac{A}{[1+(Bf_x)^2]^{C/2}} & \langle S(f_{xy})^{(2D)} \rangle &= \frac{A'}{[1+(Bf_{xy})^2]^{(C+1)/2}} \\ A' &= \frac{\Gamma[(C+1)/2]}{\Gamma[1/2]\Gamma[C/2]} \cdot AB & \sigma^2 &= \frac{\Gamma[1/2]\Gamma[(C-1)/2]}{\Gamma[C/2]} \cdot \frac{A}{B} \quad C > 1 \end{aligned} \quad (44)$$

Here σ is the rms roughness parameter that describes its vertical character, $B/(2\pi)$ is the roughness “correlation length” that characterizes its transverse character, and C is the spectral index that determines its high-frequency behavior.

The corresponding ACE, obtained by taking the Fourier transform of $\langle S(f_x)^{(1D)} \rangle$, is

$$\langle C(\tau) \rangle_{ABC} = \sqrt{2\pi} \frac{2A}{B} \frac{2^{-C/2}}{\Gamma[C/2]} \left(2\pi \frac{|\tau|}{B} \right)^{(C-1)/2} \cdot K_{(C-1)/2} \left(2\pi \frac{|\tau|}{B} \right) \quad C > 0 \quad (45)$$

where K_n is a modified Bessel function.³³ For this reason the “ABC” model is also called the K-correlation model.³⁰ When $C = 2, 4, 6, \dots$, it reduces to simple algebraic expressions. For $C = 2$, and 4, for example,

$$\begin{aligned} \langle C(\tau) \rangle_{ABC} |_{C=2} &= \pi \frac{A}{B} \exp \left[-2\pi \frac{|\tau|}{B} \right] \\ \langle C(\tau) \rangle_{ABC} |_{C=4} &= \frac{\pi}{2} \frac{A}{B} \exp \left[-2\pi \frac{|\tau|}{B} \right] \cdot \left[1 + 2\pi \frac{|\tau|}{B} \right] \end{aligned} \quad (46)$$

where the $C = 2$ form is the well-known two-sided exponential. When $C \rightarrow \infty$ the covariance function in Eq. (45) becomes gaussian.

The PEX Model Another three-parameter model with an inverse power-law power spectrum is the power exponential model:

$$\langle C(\tau) \rangle_{PEX} = \sigma^2 \exp \left[\frac{|\tau|^N}{\rho} \right] \quad 0 < N < 2 \quad (47)$$

This model is identical with the ABC model for $N = 1$ and $C = 2$, but otherwise they are different. Some of the mathematical properties of the 1D and 2D spectra of the PEX model have been given earlier in connection with fractal surfaces. In particular, Eqs. (28a) and (28b) give the low- and high-frequency forms, and Eq. (29) gives explicit results for $N = 1$.

Note that the ABC model only requires that $C > 1$ to be statistically stationary. This admits faster high-frequency falloffs than are permitted for the fractal and PEX models, which translate to $1 < C < 3$. This means that surfaces with $C \geq 3$ only fit the ABC model.

The J-K Model In the limit of long correlation lengths, the ABC, PEX, and fractal model each reduces to what we call the J-K model, defined in Eq. (43). This can be viewed as a new two-parameter finish model, where the magnitudes of J and m , or equivalently, K and n , are determined directly from experimental data without reference to any underlying model. On the other hand, the values of J and K can be written in terms of the parameters of the ABC, PEX, or fractal models leading to these high-frequency forms. For the fractal surfaces, for example,

$$K = \frac{m-2}{4 \cdot \pi^{m-1}} \cdot \frac{\Gamma[1/2] \Gamma[(m-1)/2]}{\Gamma[(4-m)/2]} \cdot T^{4-m} \quad S(f_x) = \frac{K}{|f_x|^{m-1}} \quad (48)$$

where $2 < m < 4$. This allows the topohesy T of a fractal scatterer to be determined from its high-frequency tail.

The observation of an inverse-power-law spectrum over a limited bandwidth does not guarantee that the surface is fractal. To do that one would have to observe low-frequency behavior conforming to Eqs. (24) and (25) as well, and confirm the unique wavelength and angular dependencies they imply.³⁴ Until then, power-law scattering over a finite bandwidth can at best be called “fractal-like” rather than fractal.

8.8 SURFACE FINISH SPECIFICATION

Performance Measures

It is desirable to express performance requirements on reflecting surfaces in terms of surface-finish parameters that can be measured by scatterometry or profilometry. One can concoct a wide variety of image-based performance measures.³⁵ For example, earlier we invoked an ambitious measure of image degradation based on the reduction of the on-axis image intensity, rather than the simpler measure based on the integrated image intensity suggested below.^{36,37}

The simplest practical measure appears to place an upper limit on the average total scattering outside the specular core,

$$\begin{aligned}\langle \mathcal{E}^{\text{FK, RR}} \rangle^{(2\text{D})} &= \frac{1}{R(\theta_i)} \int_{\text{core}}^{\infty} d\omega_f \cos(\theta_f) \langle \text{BRDF}_{\alpha \rightarrow \alpha}^{(2\text{D})} \rangle^{\text{FK, RR}} = \frac{\lambda^2}{R(\theta_i)} \int_{\text{core}}^{\infty} d\mathbf{f}_{xy} \langle \text{BRDF}_{\alpha \rightarrow \alpha}^{(2\text{D})} \rangle^{\text{FK, RR}} \\ \langle \mathcal{E}^{\text{FK, RR}} \rangle^{(1\text{D})} &= \frac{1}{R(\theta_i)} \int_{\text{core}}^{\infty} d\theta_f \cos(\theta_f) \langle \text{BRDF}_{\alpha \rightarrow \alpha}^{(1\text{D})} \rangle^{\text{FK, RR}} = \frac{\lambda}{R(\theta_i)} \int_{\text{core}}^{\infty} df_x \langle \text{BRDF}_{\alpha \rightarrow \alpha}^{(1\text{D})} \rangle^{\text{FK, RR}}\end{aligned}\quad (49)$$

For small correlation lengths the BRDF is flat at low frequencies, the omitted parts of the integrations in Eq. (49) can be neglected, and so

$$\langle \mathcal{E}^{\text{FK, RR}} \rangle^{(2\text{D})} = \langle \mathcal{E}^{\text{FK, RR}} \rangle^{(1\text{D})} = \left(4\pi \cos(\theta_i) \frac{\sigma}{\lambda} \right)^2 = g \quad (50)$$

where σ is the intrinsic rms roughness of the surface given by Eq. (40). In the earlier literature this result is called the total integrated scatter (TIS). This simple and beautiful expression for the image error permeates the wave-optics and microwave literature.¹⁵

In the limit of large correlation lengths, which is required to display inverse-power-law behavior, the BRDFs appearing in Eq. (49) diverge at low spatial frequencies. This leads to very different forms for the $\langle \mathcal{E} \rangle$'s, which now depend on the bandwidth limits included in the error calculation.

In the J-K language of Eq. (43),

$$\langle \mathcal{E}_{\rho > L_{xy}}^{(2\text{D})} \rangle = \left(4\pi \frac{\cos(\theta_i)}{\lambda} \right)^2 \cdot \frac{2\pi}{m-2} \cdot J L_{xy}^{m-2} \quad \langle \mathcal{E}_{\rho > L_x}^{(1\text{D})} \rangle = \left(4\pi \frac{\cos(\theta_i)}{\lambda} \right)^2 \cdot \frac{2}{m-2} \cdot K L_x^{m-2} \quad (51)$$

where L_{xy} and L_x are the *radii* of the excluded regions in Eq. (49).

Note that for statistically stationary surfaces the surface becomes perfectly smooth when $\sigma \rightarrow 0$, while for fractal surfaces this occurs when $T \rightarrow 0$.

Numerical Illustration

In an earlier paper we discussed the measurements of a silicon cylinder made with the LTP profiling instrument at Brookhaven National Laboratory.³⁷ These were fitted to the two-sided profile spectrum

$$S(f_x)^{(1\text{D})} = \frac{K}{|f_x|^{m-1}} = \frac{3.32 \times 10^{-9}}{|f_x|^{1.61}} \mu\text{m}^3 \quad 10^{-5} < |f_x| < 10^{-1} \mu\text{m}^{-1} \quad (52)$$

If fractal, this corresponds to $T = 8.36 \times 10^{-3}$ nm. If the surface is isotropic, the corresponding (2D) spectrum is

$$S(f_{xy})^{(2\text{D})} = \frac{J}{f_{xy}^m} = \frac{1.45 \times 10^{-9}}{f_{xy}^{2.61}} \mu\text{m}^4 \quad (53)$$

These allow us to evaluate the errors associated with the use of this mirror in 1D and 2D geometries.

If $\lambda = 10(-4) \mu\text{m}$ (0.1 nm), $\theta = (\pi/2) - 10^{-3}$ (i.e., 1 mrad glancing incidence), and $L_x = L_o = 0.1 \text{ m}$,

$$\langle \epsilon^{(2D)} \rangle = 26.5\% \quad \langle \epsilon^{(1D)} \rangle = 19.3\% \quad (54)$$

The corresponding errors at $\lambda = 0.6328 \mu\text{m}$ and at normal incidence are smaller than these by a factor of 0.025.

These results indicate that for the parameters used, this mirror would perform marginally well as a glancing-incidence x-ray mirror, but very well as a normal-incidence mirror at visible wavelengths, both in 1D and 2D applications.

Statistical Fluctuations

The discussion to this point has been concerned with ensemble average quantities $\langle U \rangle$. Individual, deterministic measurements of U will fluctuate about this average by an amount measured by the dimensionless quantity

$$\gamma_U^2 = \frac{\langle (U - \langle U \rangle)^2 \rangle}{\langle U \rangle^2} = \frac{\langle U^2 \rangle}{\langle U \rangle^2} - 1 \quad (55)$$

The limit $\gamma_U \rightarrow 0$ indicates a perfect measurement while $\gamma_U \rightarrow 1$ is particularly bad.

An important example of this is the inverse-power-law sum that appears in the periodogram and effective mean-square roughness in Eqs. (49) through (51) for the J-K model:

$$U(n) = \sum_{j=1}^{\infty} \frac{1}{j^n} \quad \gamma_U = \frac{\sqrt{\zeta(2n)}}{\zeta(n)} \quad (56)$$

where $\zeta(n)$ is the Riemann zeta function. For the brownian fractal, $n = 2$, for example, $\gamma_U = \sqrt{2/5} = 0.632 \dots$. The Tchebycheff inequality then indicates that the sum U may be rather broadly distributed about $\langle U \rangle$ except for n near unity.

γ_U vanishes when $n \rightarrow 1$ because the number of degrees of freedom included in the sum becomes infinite, and it approaches unity in the opposite limit of $n \rightarrow \infty$, since the sum then includes only a single term with two degrees of freedom, or a "single speckle."

If the fluctuations in $\langle \epsilon \rangle$ are unacceptably large, the lowest-order statistical description must be replaced by a deterministic version. Steps in this direction have been taken by Mori et al., Mimura et al., and Yamauchi et al.³⁸⁻⁴⁰, for very high-performance x-ray mirrors. They have reported very accurate profile measurements of highly polished mirrors which they have correlated with image quality measurements made using the 1-km beam line at SPring-8. Unfortunately, the discussion of this interesting and important work lies outside the scope of the present chapter.

8.9 RETROSPECT AND PROSPECT

This chapter outlines methods for understanding the performance of mirror surfaces in terms of statistical models of their surface topography. We illustrate this making use of simple models that exhibit an inverse-power-law fall off at high spatial frequencies. Obvious follow-on steps are to expand the database, to relate residual roughness to finishing methods, to explore different and composite models, and examine measurement errors and instrumental effects in metrology.⁴¹⁻⁴⁴

It is pointed out that the utility of certain statistical parameters may be limited by broad confidence limits. This suggests that a deterministic approach may be desirable for demanding applications, especially involving one-of-a-kind surfaces. Even so, the statistical approach provides deep insight into the physics involved, offers a pre-screening methodology, and will remain the lingua franca of this wide and diverse field in the future.

8.10 REFERENCES AND ENDNOTES

1. P. Beckmann and A. Spizzichino, *The Scattering of Electromagnetic Waves from Rough Surfaces*, Pergamon Press, New York, NY, 1963.
2. G. T. Ruck, D. E. Barrick, W. D. Stuart, and C. K. Krichbaum, *Radar Cross Section Handbook*, 2 vols., Pergamon Press, New York, NY, 1970. Especially D. E. Barrick, Chap. 9, "Rough Surfaces."
3. A. Ishimaru, *Electromagnetic Wave Propagation, Radiation and Scattering*, Prentice-Hall Inc., Englewood Cliffs, NJ, 1991; A. Ishimaru, *Wave Propagation and Scattering in Random Media*, 2 vols., Academic Press, New York, NY, 1978.
4. A. K. Fung, *Microwave Scattering and Emission Models and their Application*, Artech House, Boston, MA, 1994.
5. M. Nieto-Vesperinas, *Scattering and Diffraction in Physical Optics*, 2d ed., World Scientific, Hackensack, NJ, 2006.
6. A. G. Voronovich, *Wave Scattering from Rough Surfaces*, 2d ed, Springer Verlag, New York, NY, 1999. See also Chaps. 2 and 4 in Ref. 7 below.
7. A. A. Maradudin, ed., *Light Scattering and Nanoscale Surface Roughness*, Springer Verlag, New York, NY, 2007.
8. J. C. Stover, *Optical Scattering; Measurement and Analysis*, 2d ed., SPIE: Optical Engineering Press, Bellingham, WA, 1995.
9. T. A. Germer, "Measuring Interfacial Roughness by Polarized Optical Scattering," in Ref. 7 above, and references therein.
10. P. Z. Takacs, S. K. Feng, E. L. Church, S. Qian, and W. -M. Liu, "Long Trace Profile Measurements on Cylindrical Aspheres," *Proc. SPIE* **966**:354–364, 1988; see also P. Z. Takacs, K. Furenlid, R. DeBiasse, and E. L. Church, "Surface Topography Measurements over the 1 Meter to 10 Micrometer Spatial Period Bandwidth," *Proc. SPIE* **1164**:203–211, 1989.
11. K. A. O'Donnell and E. R. Mendez, "Experimental Study of Scattering from Characterized Random Surfaces," *J. Opt. Soc. Am. A* **4**:1194–1205, 1987; See also Chap. 9 in Ref. 7 above.
12. E. R. Mendez and D. Macías, "Inverse Problems in Optical Scattering," Chap. 16 in Ref. 7 above.
13. M. J. Kim, E. R. Mendez, and K. A. O'Donnell, "Scattering from Gamma-Distributed Surfaces," *J. Mod. Opt.* **34**:1107–1119, 1987.
14. H. Cramér and M. R. Leadbetter, *Stationary and Related Stochastic Processes*, Dover Pubs. Mineola, NY, 2004.
15. M. Born and E. Wolf, *Principles of Optics*, 4th ed., Pergamon, New York, NY, 1970. Chap. IX, Sec. 9.1 derives the deterministic "Strehl" criterion for focusing optics. These considerations are very well known to the visible-optics community. X-ray-opticians, on the other hand, often refer to rms slope errors, which is the square root of the bandwidth-limited value of the second moment of the power-spectral density.
16. M. V. Berry, "Diffractals," *J. Phys. A: Math. Gen.* **12**:781–797, 1979.
17. I. Simonsen, D. Vandembroucq, and S. Roux, "Wave Scattering from Self-Affine Surfaces," *Phys. Rev. E* **61**:5914–5917, 2000; see also *ibid* "Electromagnetic Wave Scattering from Conducting Self-Affine Surfaces: An Analytic and Numerical Study," *J. Opt. Soc. Am. A* **18**:1101–1111, 2001.
18. B. B. Mandelbrot, *The Fractal Geometry of Nature*, Freeman, New York, NY, 1975. The connection between the Hausdorff and Hurst dimensions is given as $D' + H' = 1 + \text{dimensionality}$. We prefer the descriptor D' because of its direct connection with the geometric scaling properties of the profile. See Ref. 16 for specifics.
19. J. P. Nolan, *An Introduction to Stable Distributions*, to be published. The program STABLE is available from J. P. Nolan's Web site: academic2.american.edu/~jpnolan.
20. S. O. Rice, "Reflection of Electromagnetic Waves from Slightly Rough Surfaces," *Commun. Pure Appl. Math.* **4**:361–378, 1951.
21. W. H. Peake, "Theory of Radar Return from Terrain," *IRE Conv. Rec.* **7**(1):34–41, 1959.
22. K. A. O'Donnell, "Small-Amplitude Perturbation Theory for One-Dimensionally Rough Surfaces," Chap. 5 in Ref. 7; see also K. A. O'Donnell, "High Order Perturbation Theory for Light Scattering from a Rough Metal Surface," *J. Opt. Soc. Am. A* **18**:1507–1518, 2001.
23. G. Berginc, "Small-Amplitude Perturbation Theory for Two-Dimensional Surfaces," Chap. 6 in Ref. 7 above.
24. J. T. Johnson, "Third Order Small Perturbation Method for Scattering from Dielectric Rough Surfaces," *J. Opt. Soc. Am.* **16**:2720–2736, 1999; see also J. T. Johnson, "Computer Simulations of Rough Surface Scattering," Chap. 7 in Ref. 7 above.

25. G. R. Valenzuela, short communication, *Trans. IEEE*, ca. 1972. Precise title and reference unavailable.
26. L. Tsang, J. A. Kong, and R. T. Shin, *Theory of Microwave Remote Sensing*, John Wiley and Sons, New York, NY, 1985.
27. D. Barrick and R. Fitzgerald, "The Failure of 'Classic' Perturbation Theory at a Rough Neumann Boundary near Grazing," *IEEE Trans. On Antennas and Propagation* **48**:1452–1460, 2000.
28. E. L. Church and J. C. Stover, "Roughness Effects on Emissivity," unpublished notes, 1997.
29. The literature frequently uses a gaussian system response function, which leads to algebraic simplifications with no essential loss of generality.
30. E. L. Church and P. Z. Takacs, "The Optimal Estimation of Finish Parameters," *Proc. SPIE*.**1530**:71–85, 1992.
31. A. Erdelyi, *Higher Transcendental Functions*, The Batemann Manuscript Project. Tables of Integral Transforms, 2 vols., McGraw-Hill, New York, NY, 1953. In particular, q.v. Weyl fractional calculus.
32. E. L. Church, T. A. Leonard, and P. Z. Takacs, "The Prediction of BRDFs from Surface Profile Measurements," *Proc. SPIE*.**1165**:136–150, 1991.
33. I. S. Gradshteyn, I. M. Ryzhik, and A. Jeffrey, eds., *Table of Integrals, Series and Products*, 5th ed. Academic Press, Inc., Boston, MA, 1980. In particular, entries 6.565.4, 6.576.7, 6.699.12, and 8.432.5.
34. I. A. Popov, L. A. Glushchenko, and J. Uozumi, "The Study of Fractal Structure of Ground Glass Surface by Means of Angle Resolved Scattering of Light," *Optics Comm.* **203**:191–196, 2002.
35. E. L. Church and P. Z. Takacs, "Measures of Surface Quality," unpublished notes, 2007.
36. E. L. Church and P. Z. Takacs, "Specification of Surface Figure and Finish in Terms of System Performance," *Appl. Opt.* **32**:3344–3353, 1993.
37. E. L. Church and P. Z. Takacs, "Specification of Glancing- and Normal-Incidence X-Ray Mirrors," *Opt. Eng.* **34**:353–359, 1995; erratum, *ibid* **34**:3348, 1995.
38. Y. Mori, K. Yamauchi, K. Yamamura, et al., "Sub-Micron Focusing of Hard X-Ray Beam by Elliptically Figured Mirrors for Scanning X-Ray Microscopy," *Proc. SPIE*. **4782**:58–64, 2002, and references therein.
39. H. Mimura, K. Yamauchi, K. Yamamura, et al., "Image Quality Improvement in Hard X-Ray Projection Microscope Using Total Reflection Mirror Optics," *J. Synchrotron Rad.* **11**:343–346, 2004; see also, H. Mimura, H. Yumoto, S. Matsuyama, et al., "Relative Angle Determinable Stitching Interferometry for Hard X-Ray Reflective Optics," *Rev. Sci. Instrum.* **76**:045102, 2005.
40. K. Yamauchi, K. Yamamura, H. Mimura, et al., "Wave-Optical Evaluation of Interference Fringes and Wavefront Phase in a Hard-X-Ray Beam Totally Reflected by Mirror Optics," *Appl. Opt.* **44**:6927–6932, 2005.
41. J. M. Bennett, "Characterization of Surface Roughness," Chap. 1 in Ref. 7 above.
42. E. L. Church, "Fractal Surface Finish," *Appl. Opt.* **27**:1518–1526, 1988.
43. V. V. Yashchuk, E. L. Church, M. Howells, W. R. McKinney, and P. Z. Takacs, "21st Century Metrology for Synchrotron Radiation Optics Is Understanding How to Specify and Characterize Optics Completely," *SRI-2006 III Workshop on Optical Metrology*, Daegu Korea, 28 May 2006.
44. V. V. Yashchuk, S. C. Irick, A. A. McDowell, W. R. McKinney, and P. Z. Takacs, "Air Convection Noise of Pencil-Beam Interferometer for Long Trace Profiler," Lawrence Berkeley National Laboratory Notes, 2008; Also, V. V. Yashchuk, "Positioning Errors of Pencil-Beam Interferometers for Long Trace Profilers," *ibid.*, 2008.

This page intentionally left blank.

DO NOT DUPLICATE

VOLUME SCATTERING IN RANDOM MEDIA

Aristide Dogariu and Jeremy Ellis

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

9.1 GLOSSARY

A	transversal area
C	correlation function
c	speed of light
D	diffusion constant
D	degree of polarization
E	field amplitude
e	polarization direction
F	scattering form amplitude
f	forward-scattering amplitude
G	pair correlation function
g	asymmetry parameter
I	field intensity
l	Stokes vector
J	diffuse flux
k	wave vector
K_B	Boltzman constant
K_{K-M}	effective absorption coefficient
l_a	absorption length
l_s	scattering length
l^*	transport mean free path
N	number of particles
n	refractive index
P	scattering phase function
q	scattering wavevector

\mathbf{r}, \mathbf{R}	vector position
S	static structure factor
S	scattering potential
\mathbf{s}	specific direction
S_{K-M}	effective scattering coefficient
T	absolute temperature
t	time
T_{ab}	transmission coefficient
U	diffuse energy density
V	volume
v	volume fraction
ϵ	dielectric constant
η	shear viscosity
θ	scattering angle
λ	wavelength of light
μ_a	absorption coefficient
μ_s	scattering coefficient
ρ	number density
σ	scattering cross-section
ω	frequency
Ω	solid angle

9.2 INTRODUCTION

Electromagnetic radiation impinging on matter induces oscillating charges that can be further regarded as secondary sources of radiation. The morphological details and the microscopical structure of the probed medium determine the frequency, intensity, and polarization properties of this re-emitted (scattered) radiation. This constitutes the basis of a long history of applications of light scattering as a characterization tool in biology, colloid chemistry, solid state physics, and so on.

A substantial body of applications deals with light scattering by particles. These smaller or larger ensembles of molecules have practical implications in many industries where they are being formed, transformed, or manipulated. Since the early works of Tyndall and Lord Rayleigh,¹ the study of light scattering by molecules and small particles has been consistently in the attention of many investigators. Classical reviews of the field are the books by van de Hulst,² Kerker,³ Bayvel and Jones,⁴ Bohren and Huffman;⁵ we also note the recent survey of techniques and theoretical treatments by Jones.⁶ Why and how the light is scattered by small particles has already been described in Chap. 7, “Scattering by Particles,” by Craig F. Bohren in this volume. Discussions on subjects related to light scattering can also be found in chapters like Chap. 5, “Coherence Theory,” by William H. Carter and Chap. 12, “Polarization,” by Jean M. Bennett.

For the topic of this chapter, light scattering by individual particles constitutes the building block of more complicated physical situations. When the three-dimensional extent of the medium that scatters the light is much larger than the typical size of a local inhomogeneity (scattering center), the physical process of wave interaction with matter can be classified as volume scattering. In this regime, the measured radiation originates from many different locations dispersed throughout the volume. Depending upon the structural characteristics of the medium, various scattering centers can act as secondary, independent sources of radiation (*incoherent scattering*) or they can partially add their contributions in a collective manner (*coherent scattering*). Another situation of interest happens when, for highly disordered systems, light is scattered successively at many locations throughout the volume (*multiple scattering*). All these three aspects of volume scattering will be discussed in this chapter.

In practice, particles very rarely exist singly and, depending on the illuminated volume or the volume seen by the detection system, scattering by a large number of particles needs to be considered. The simplest situation is that of *incoherent scattering*. When the fields scattered by different centers are completely independent, the measured intensity results from an intensity-based summation of all individual contributions. The ensemble of particles is described by temporal and spatial statistics which does not show up in scattering experiments; one can say that the volume scattering does not resolve the spatial arrangement of scattering centers.

When the scattering centers are sufficiently close, the phases of wavelets originating from individual scattering centers are not independent. This is the case of collective or *coherent scattering*. One faces the problem of expanding the scattering theories to ensembles of particles that can have certain degree of spatial or temporal correlations. A transition sets in from independent to coherent scattering regime. The situation is common for gels or composites which scatter light due to local inhomogeneities of the refractive index with length scales of the order of wavelength and where spatial correlations between the scattering centers are encountered. This is the basis of one of the most successful application of volume scattering: the observation of structural characteristics of inhomogeneous systems.

In the case of highly disordered systems, the light propagation can be subject of scattering at many different locations within the probed volume and a *multiple-scattering* regime sets in. For a long time, the intensity and phase fluctuations determined by multiple light scattering were regarded as optical “noise” that degrades the radiation by altering its coherence, broadening the beam, and decreasing its intensity. Experimentalists were trying to avoid it as much as possible and the development of comprehensive theories was not sufficiently motivated. Over the last two decades, however, remarkable advances in fundamental understanding and experimental methodologies proved that multiple scattering of waves is a source for unexplored physics leading to essentially new applications. New phenomena have been discovered and a series of experimental techniques have been implemented using particular coherent, polarization, temporal, and spectral properties of multiple scattered light. This revival of interest has been stimulated by the use of highly coherent sources in remote sensing and, especially, by considerable advances in solid-state physics. Many features of multiple scattered light are common to other classical waves like sound, heat, or microwaves but several analogies with electron transport phenomena have been at the core of this renewed interest in the propagation of optical waves in random systems.⁷

There is also another situation, which is often encountered in optics, when waves propagate through media with abrupt changes in their optical properties. Waves passing through inhomogeneous media with defined boundaries usually suffer surface scattering. In principle, scattering at rough surfaces can be considered as a limiting case of wave propagation and it is significant in various practical situations; this topic has been separately discussed in Chap. 8, “Surface Scattering,” by Eugene L. Church and Peter Z. Takacs in this volume.

9.3 GENERAL THEORY OF SCATTERING

The schematic of a typical scattering experiment is depicted in Fig. 1, where a plane wave \mathbf{E}^0 with the wavevector $k = \omega/c$ is incident on a spatially random medium occupying a finite volume V . Light is scattered by local inhomogeneities of the dielectric constant $\epsilon(\mathbf{r})$ and a basic theory of scattering aims at providing the link between the experimentally accessible intensity $I_s(\mathbf{R}) = |\mathbf{E}_s(\mathbf{R})|^2$ and the microscopic structure of the random medium.

The starting point of the theory is to describe the total electric field $\mathbf{E}(\mathbf{r})$ as a summation of the incoming and scattered fields and to consider that it satisfies the equation $(\nabla^2 + k^2)\mathbf{E}(\mathbf{r}) = -4\pi\mathbf{S}(\mathbf{r})\mathbf{E}(\mathbf{r})$, where $\mathbf{S}(\mathbf{r})$ represents a generic scattering potential. This equation can be converted into an integral one and, for \mathbf{R} sufficiently far from the scattering volume where the associated Green function simplifies, one eventually obtains the general result^{3,8}

$$\mathbf{E}_s(\mathbf{R}) = \frac{e^{ikR}}{R} \frac{k^2}{4\pi} \int_{(V)} d\mathbf{r} \{-\mathbf{k}_s \times [\mathbf{k}_s \times (\epsilon(\mathbf{r}) - 1) \cdot \mathbf{E}(\mathbf{r})]\} e^{-i\mathbf{k}_s \cdot \mathbf{r}} \quad (1)$$

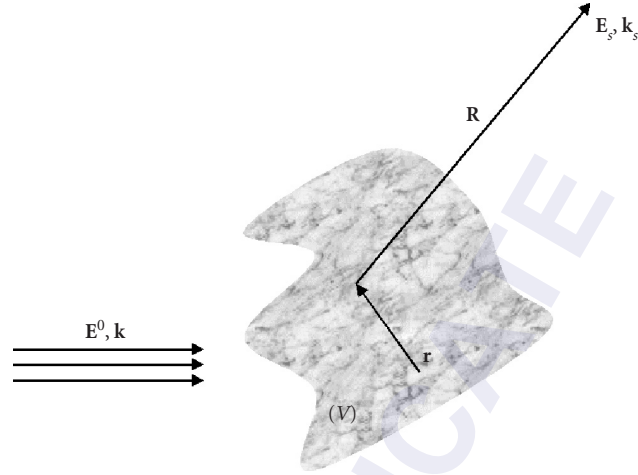


FIGURE 1 Schematic representation of the incident \mathbf{k} and scattered \mathbf{k}_s beams in a generic scattering experiment.

This expression represents the scattered field as an outgoing spherical wave that depends on the direction and magnitude of the total field inside the scattering volume V .

Approximate solutions can be obtained for the case of weak fluctuations of the dielectric constant. One can expand the field $\mathbf{E}(\mathbf{r}) = \mathbf{E}^0(\mathbf{r}) + \mathbf{E}_1(\mathbf{r}) + \mathbf{E}_2(\mathbf{r}) + \dots$ in terms of increasing orders of the scattering potential and use successive approximations of $\mathbf{E}(\mathbf{r})$ in Eq. (1) to obtain the so-called Born series. In the spirit of a first iteration, one replaces $\mathbf{E}(\mathbf{r})$ with $\mathbf{E}^0(\mathbf{r})$ and obtains the first Born approximation that describes the regime of single scattering.

Alternatively, one can write $\mathbf{E}(\mathbf{r}) = \exp[\Psi(\mathbf{r})]$ and develop the series solution for $\Psi(\mathbf{r})$ in terms of increasing orders of the scattering potential. This is the Rytov's series of exponential approximations, an alternative to the algebraic series representation of the Born method. The two approaches are almost equivalent, however, preference is sometimes given to Rytov's method because an exponential representation is believed to be more appropriate to describe waves in line-of-sight propagation problems.⁹⁻¹¹

It is worth pointing out here that Eq. (1) can be regarded as an integral equation for the total field and, because the total field is a superposition of the incident field and contributions originating from scattering from all volume V , this generic equation includes all possible multiple scattering effects.^{8,12}

9.4 SINGLE SCATTERING

Incoherent Scattering

When a sparse distribution of scattering centers is contained in the volume V , all the scattering centers are practically exposed only to the incident field, $\mathbf{E}(\mathbf{r}) = \mathbf{E}^0(\mathbf{r}) = E^0 \mathbf{e}^0 e^{i\mathbf{k}\cdot\mathbf{r}}$ where E^0 is the field magnitude, \mathbf{e}^0 is the polarization direction, and \mathbf{k} is the wave vector. A considerable simplification is introduced when the magnitude of the scattered field is much smaller than that of the incident field, such that the total field inside the medium can be everywhere approximated with the incident field. This is the condition of the first Born approximation for tenuous media that practically neglects multiple scattering effects inside the scattering volume V . From Eq. (1), it follows that

$$\mathbf{E}_s(\mathbf{R}) = E^0 \frac{e^{i\mathbf{k}_s \cdot \mathbf{R}}}{R} \frac{k^2}{4\pi} \int_{(V)} d\mathbf{r} [\mathbf{e}^s \cdot (\boldsymbol{\varepsilon}(\mathbf{r}) - 1) \cdot \mathbf{e}^0] e^{-i\mathbf{k}_s \cdot \mathbf{r}} \quad (2)$$

For instance, in the case of a system of N identical particles, Eq. (2) is evaluated to give

$$\mathbf{E}_s(\mathbf{R}) = E^0 \frac{e^{ikR}}{R} \frac{k^2}{4\pi} \sum_{j=1}^N e^{-ik_s \cdot \mathbf{r}_j} \int_{(V_j)} d\mathbf{r} [\mathbf{e}^s \cdot (\boldsymbol{\varepsilon}(\mathbf{r}) - 1) \cdot \mathbf{e}^0] e^{-ik_s \cdot \mathbf{r}} \quad (3)$$

where V_j is the volume of the j th particle located at \mathbf{r}_j . In terms of the scattering wavevector $\mathbf{q} = \mathbf{k}_s - \mathbf{k}$, the integral in Eq. (3) has the meaning of single-scattering amplitude of the j th particle, $F(\mathbf{q}) = |f| \cdot B(\mathbf{q})$, and depends on the forward scattering amplitude f and a scattering amplitude normalized such that $B(0) = 1$. As explained in Chap. 7 in this volume, the ratio between the refractive index of the particles n_p and of the suspending medium n_s determine the scattering strength f of an individual scatterer while its shape and size are accounted for in $B(\mathbf{q})$.

For a collection of discrete scattering centers the total scattered intensity of Eq. (3) factorizes like:¹³

$$I(\mathbf{q}) = \left(E_i \frac{e^{ikR}}{R} \frac{k^2}{4\pi} \right)^2 N F^2(\mathbf{q}) \sum_{i,j=1}^N e^{-iq \cdot (\mathbf{r}_i - \mathbf{r}_j)} : NP(\mathbf{q}) S(\mathbf{q}) \quad (4)$$

where we separated the single-scattering form factor $P(\mathbf{q})$ from the interference function or the static structure factor $S(\mathbf{q}) = \sum_{i,j=1}^N e^{-iq \cdot (\mathbf{r}_i - \mathbf{r}_j)}$. The structure factor quantifies the phase-dependent contributions due to different locations of scattering centers. When an ensemble average is taken over the volume V and after separating out the diagonal terms, the static structure factor can be written as^{13,14}

$$S(\mathbf{q}) = 1 + \langle e^{-iq \cdot (\mathbf{r}_i - \mathbf{r}_j)} \rangle = 1 + \rho \int G(\mathbf{r}) e^{-iq \cdot \mathbf{r}} d\mathbf{r} \quad (5)$$

in terms of the pair-correlation function $G(\mathbf{r})$ (where \mathbf{r} is the vectorial distance between two particles), describing the statistical properties of the spatial arrangement of scattering centers. It is through $S(\mathbf{q})$ that the link is made between the statistical mechanics description of the inhomogeneities and the measurable quantities in a scattering experiment.

As can be seen from Eq. (5), for the case of particles separated by distances much larger than the wavelength, $S(\mathbf{q})$ becomes unity; the situation corresponds to $G(\mathbf{r}) \equiv 1$, i.e., constant probability to find scattering centers anywhere in the scattering volume. This regime characterized by $S(\mathbf{q}) = 1$ is also called the incoherent case of volume scattering where $I(\mathbf{q})$ is a simple, intensity-based summation of individual contributions originating from different scattering centers.

Coherent Scattering

For higher volume fractions of particles, the pair-correlation function depends on both the particle size and their concentration. The estimation of the pair-correlation functions—and, therefore, the evaluation of an explicit form for the structure factor—is a subject of highest interest and is usually approached through various approximations.¹⁵

Typical structure factors are shown in Fig. 2 for increasing volume fractions of spherical particles with radius r_0 . These results are based on the Percus-Yevick approximation,^{16,17} which has the advantage of being available in a closed mathematical form but the trend shown in Fig. 2 is rather general. Note that the strength of the interparticle interactions is practically measured by the magnitude of the first peak in $S(q)$.^{13,14}

In general, wave-scattering experiments can be used to infer the pair-correlation function of the scattering centers in a random medium and, meanwhile, the characteristic of an individual scattering event. In an ideal experiment where one has access to all the values of $I(\mathbf{q})$ and the single-scattering phase function $P(\mathbf{q})$ is known, Eq. (3) can be inverted by standard methods.^{17,18} Capitalizing on conventional types of single-scattering form factors, various inversion schemes have been implemented to extract the structure factor from the values of angular-resolved intensity.¹⁹

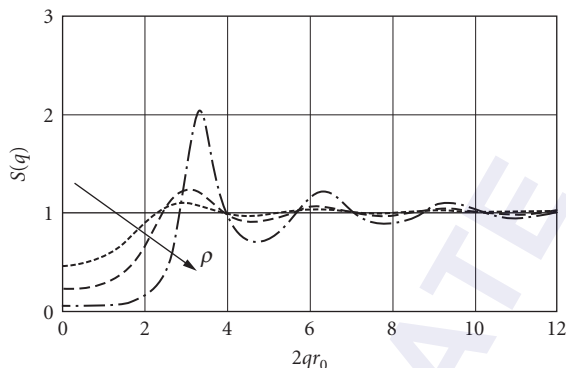


FIGURE 2 Typical structure factors corresponding to systems of identical spherical particles with increasing volume fractions as indicated. The horizontal line at $S(q)=1$ corresponds to the independent scattering approximation.

Convergence and stability in the presence of noise are the major requirements for a successful inversion procedure and a subsequent Fourier analysis to provide a description of the pair-correlation function. Of course, for systems of nonspherical particles or when the system exhibits structural anisotropies, a fully vectorial inverse problem needs to be approached.

The factorization approximation of Eq. (4) has found numerous applications in optical-scattering experiments for characterization of colloidal, polymeric, and complex micellar systems. Numerous static and dynamic techniques were designed to probe the systems at different length and time scales. For example, reaction kinetics or different phase transitions have been followed on the basis of angular and/or temporal dependence of the scattered intensity.¹⁹ Another significant body of applications deals with light-scattering studies of aggregation phenomena.^{20–22}

We should not conclude this section without mentioning here the similarities between volume light scattering and other scattering-based procedures such as x-ray and neutron scattering. Of course, the “scattering potentials” and the corresponding length scales are different in these cases but the collective scattering effects can be treated in a similar manner. From a practical viewpoint, however, light scattering has the appealing features of being noninvasive and, most of the time, easier to implement.

The average power scattered by a single particle is usually evaluated using the physical concept of total scattering cross-section $\sigma = k^{-4} \int_0^{2k} P(q) q dq$.^{2,3,18} For a system with a number density $\rho = N/V$ of scattering centers, the regime of scattering is characterized by a *scattering length* $l_s = 1/\rho\sigma$. When the extent over which the wave encounters scattering centers is less than this characteristic scattering length l_s , we deal with the classical single-scattering regime. Note that the system of particles can be in the single-scattering regime and exhibit both independent or collective scattering. From the general theory of scattering it follows that the details of the scattering-form factor depend on the size of scattering particle compared to the wavelength. The deviation of the scattering-form factor from an isotropic character is characterized by the asymmetry parameter $g = \langle \cos(\theta) \rangle = 1 - 2\langle q^2 \rangle / (2k)^2$, where $\langle q^2 \rangle = k^{-2} \int_0^{2k} P(q) q^3 dq$.^{2,18}

As we have seen in the preceding section, when the particles are closely packed the scattering centers are not independent, i.e., they are sufficiently close that the fields scattered by different centers are partially in phase, collective scattering is considered through the static structure factor $S(q)$. The effect of wave coupling to an individual particle can therefore be isolated from the statistical mechanics description of the particle locations in the volume. In this regime, one can consider the dispersion of correlated particles as a collection of pseudo-scattering centers, *equivalent particles*, that are characterized by a modified single-scattering form factor $P(q)S(q)$. There is no interaction

between these fictitious particles and a corresponding single-scattering phase function can also be defined to be:

$$\widetilde{P}(q) = \frac{k^2 |f|^2 \cdot |B(q)|^2 S(q)}{\int_0^{2k} |f|^2 \cdot |B(q)|^2 S(q) q dq} \quad (6)$$

The asymmetry parameter of these pseudo-particles can be written as

$$\widetilde{g} = 1 - \frac{\int_0^{2k} |f|^2 \cdot |B(q)|^2 S(q) q^3 dq}{2k^2 \int_0^{2k} |f|^2 \cdot |B(q)|^2 S(q) q dq} \quad (7)$$

The equivalent-particle concept is illustrated in Fig. 3, where both the phase function and asymmetry parameters are presented for the specific case of silica particles suspended in water. This simplifying representation of coherent scattering effects is useful to further interpret complex, multiple scattering phenomena.

Dynamic Scattering

So far, we limited the discussion to the case where the scattering potential varies across the volume V but, at a certain location, it remains constant in time. However, many physical systems are such that the volume distribution of scattering centers fluctuates in time and, therefore, gives rise to temporal fluctuations of the scattered radiation, i.e., to dynamic scattering. A complete analysis of light scattering involves the autocorrelation function of the dielectric constant fluctuations $\langle \varepsilon^*(\mathbf{r}, 0), \varepsilon(\mathbf{r}, t) \rangle$ that manifests itself in the statistics of the temporal fluctuations of the measured light intensity. Evaluation of such correlation functions requires knowledge of the transport properties of the volume medium and is based on statistical mechanics and many-body theory. In spite of the rather complex phenomenology, different photon correlation techniques have been successfully implemented in studies of reacting systems, molecular dynamics, or atmospheric scintillations.

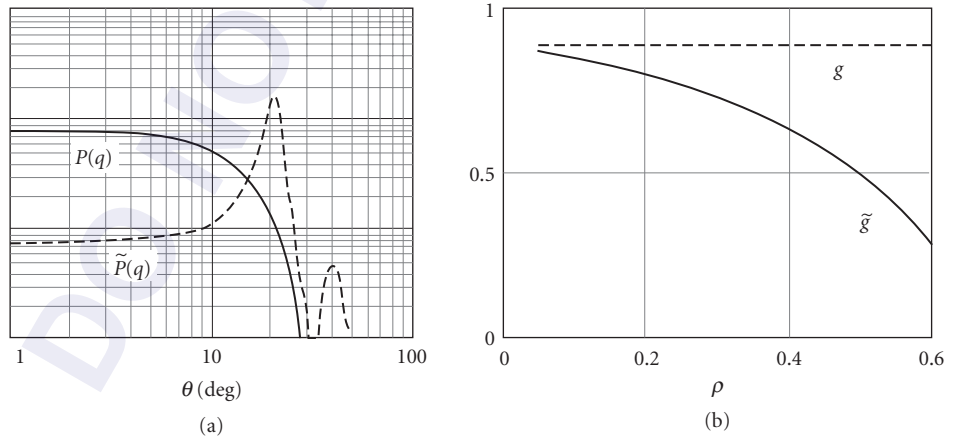


FIGURE 3 (a) Single scattering form factor $P(q)$ corresponding to silica particles of $0.476 \mu\text{m}$ placed in water and illuminated with $\lambda = 0.633 \mu\text{m}$ and the form factor $\widetilde{P}(q)$ of an equivalent particle corresponding to a collection of such particles at the volume fraction $\rho = 0.5$. (b) Values of the asymmetry parameter for one silica particle and for a collection of particles with an increasing volume fraction as indicated.

The benchmarks of the dynamic light scattering have been set in the seventies.^{23,24} An interesting and useful particularity stems from the fact that, based on dynamic scattering, one can actually measure mechanical properties in the scattering volume without knowledge of the refractive index. Random motion of scattering centers induces small Doppler frequency shifts which, in turn, produce an overall broadening of the incident spectrum of light. The detection methods for such spectral changes depend primarily on the time scales of interest and range from high-resolution spectroscopy for very fast phenomena to various mixing or beating techniques for processes slower than about 1 ms.

Based on comparing the scattering signal with itself at increasing time intervals, photon correlation spectroscopy (PCS) has emerged as a successful technique for the study of volume distributions of small particles suspended in fluids.²⁵ Practically, one deals with the time-dependent fluctuations of the speckle pattern and the goal is to determine the temporal autocorrelation function $\langle E(0)E^*(t) \rangle$, which probes the microscopic dynamics. For instance, in the simple case of a monodisperse and noninteracting system of brownian particles, $\langle E(0)E^*(t) \rangle$ is a single exponential that depends on the diffusion constant $D = k_B T / 6\pi\eta r_0$. Knowing the absolute temperature T and the shear viscosity η of the solvent, one can infer the particle radius r_0 . Similar to the case of static scattering discussed previously, refinements can be added to the analysis to account for possible dynamic structuring, polydispersivity as well as asphericity effects.^{19,23,26}

9.5 MULTIPLE SCATTERING

When optical waves propagate through media with random distributions of the dielectric constant or when they encounter extended regions containing discrete scatterers or random continuum, one needs to solve a wave equation in the presence of large number of scattering centers; this is a difficult task and, as we will discuss, a series of simplifying approaches have been proposed. A survey of multiple scattering applications is presented by van de Hulst.²⁷

Effective-Medium Representation

Some physical insight is given by a simple model, which describes the wave attenuation due to scattering and absorption in terms of an effective dielectric constant. Without explicitly involving multiple scattering, one considers the wave propagation through an homogeneous *effective-medium* which is defined in terms of averaged quantities. The effective permittivity ϵ_{eff} is calculated by simply considering the medium as a distribution of spheres with permittivity ϵ embedded in a continuum background of permittivity ϵ_0 . If a wave with the wavelength much larger than the characteristic length scales (size of inhomogeneity and mean separation distance) propagates through a volume random medium, the attenuation due to scattering can be neglected; therefore, a frequency-independent dielectric constant will appropriately describe the medium. Based on an induced dipoles model, the Maxwell-Garnett mixing formula relates the effective permittivity to the volume fraction ν of spheres $\epsilon_{\text{eff}} = \epsilon(1 + 2\nu a)/(1 - \nu a)$, where $a = (\epsilon - \epsilon_0)/(\epsilon + 2\epsilon_0)$.²⁸

In recent developments, the wave propagation through highly scattering media has been described by including multiple scattering interactions in a mean-field approach. This is simply done by considering that the energy density is uniform when averaged over the correlation length of the microstructure.²⁹ Through this effective medium, a “coherent” beam propagates with a propagation constant that includes the attenuation due to both absorption and scattering. In general, the propagation of the coherent beam is characterized by a complex index of refraction associated with the effective medium and a nontrivial dispersion law can be determined by resonant scattering.³⁰ It is worth mentioning that, in the long-wavelength limit, the attenuation due to scattering is negligible and the classical mixture formula for the effective permittivity applies.

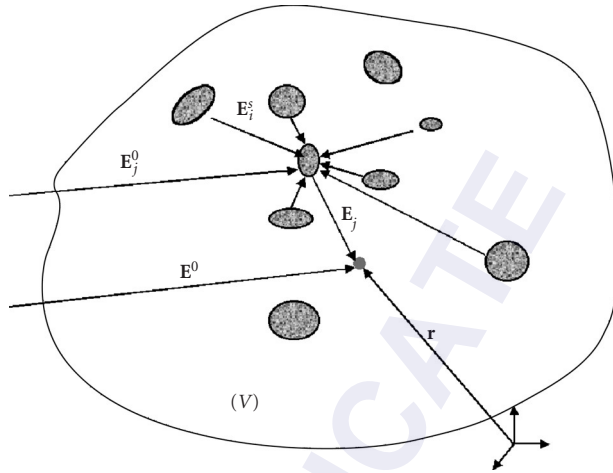


FIGURE 4 The field at \mathbf{r} is a summation of the incident field \mathbf{E}^0 and contributions \mathbf{E}_j from all other scattering centers; in turn, the effective field on each particle consists of an incident contribution \mathbf{E}^0 and contributions \mathbf{E}_i^s from all other scattering centers.

Analytical Theory of Multiple Scattering

A rigorous description of multiple light-scattering phenomena can be made if statistical considerations are introduced for quantities such as variances and correlation functions for $\epsilon(\mathbf{r})$ and general wave equations are subsequently produced. The advantage of an analytical theory is that the general formulation does not require a priori assumptions about the strength of individual scattering events nor about the packing fraction of scattering centers. The drawback, however, is that, in order to deal with the complexity of the problem, one needs to use quite involved approximations and, sometimes, rather formal representations.

Multiple Scattering Equations As shown in Fig. 4, when the field \mathbf{E}^0 is incident on a random distribution of N scattering centers located at $\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_N$ throughout the scattering volume V , the total field at one particular location inside V is the sum of the incident wave and the contributions from all the other particles

$$\mathbf{E} = \mathbf{E}^0 + \sum_{j=1}^N \mathbf{E}_j \quad (8)$$

the field scattered from the j th particle depends on the effective field incident on this particle and its characteristics (scattering potential) S_j

$$\mathbf{E}_j = S_j \left(\mathbf{E}_j^0 + \sum_{i=1, i \neq j}^N \mathbf{E}_i^s \right) \quad (9)$$

It follows from Eqs. (8) and (9) that the total field can be formally written as

$$\mathbf{E} = \mathbf{E}^0 + \sum_{j=1}^N S_j \left(\mathbf{E}_j^0 + \sum_{j=1}^N \sum_{i=1, i \neq j}^N S_i \mathbf{E}_i^s \right) \quad (10)$$

which is a series of contributions from the incident field, single scattering, and increasing orders of multiple scattering.¹⁸ In principle, knowing the scattering characteristics S_j of individual centers (this includes strength and location), one can develop the analytical approach by involving chains of successive scattering paths. From the summations in Eq. (10), by neglecting the scattering contributions that contain a scatterer more than once, Twersky has developed an expanded representation of multiple scattering which is practical only for cases of low-order scattering.³¹

Approximations of Multiple Scattering Rigorous derivation of multiple scattering equations using the Green's function associated with the multiple-scattering process and a system-transfer-operator approach (the T-matrix formalism) can be found in.³²⁻³⁴ However, in many cases of practical interest, a very large number of scattering centers needs to be involved and it is impossible to obtain accurate descriptions for either the T-matrix or the Green's function.

When large ensembles of scatterers are involved, a statistical description of the multiple scattering equations is appropriate. Probability density functions for finding scattering centers at different locations are determined by radial distribution function in a similar way as discussed in the context of correlated scattering. This expresses both the constraint on the particle locations and the fact that they are impenetrable. By using configuration-averaging procedures, self-consistent integral equations can be obtained for the average and fluctuating parts of the field produced through multiple scattering. When such a procedure is applied to Eq. (10),

$$\langle \mathbf{E} \rangle = \mathbf{E}^0 + \int S_j \mathbf{E}_j^0 G(\mathbf{r}_j) d\mathbf{r}_j + \iint S_i S_j \mathbf{E}_i^0 G(\mathbf{r}_i) G(\mathbf{r}_j) d\mathbf{r}_i d\mathbf{r}_j + \dots \quad (11)$$

a hierarchy of integral equations is generated and successive approximations are obtained by truncating the series at different stages. In fact, this expanded form has a more physically understandable representation in terms of the average field $\langle \mathbf{E}_j \rangle$ at the location of a generic particle j :

$$\langle \mathbf{E} \rangle = \mathbf{E}^0 + \int S_j \langle \mathbf{E}_j \rangle G(\mathbf{r}_j) d\mathbf{r}_j \quad (12)$$

Foldy was the first to introduce the concept of configurational averaging and used the joint probability distribution for the existence of a given configuration of scattering centers to average the resulting wave over all possible configurations.³⁵ However, Foldy's approximation is appropriate for wave propagation in sparse media with a small fractional volume of scatterers.

A comprehensive discussion on multiple scattering equations and various approximations can be found in Tsang's book³⁴ including the quasi-crystalline approximation³⁶ and coherent potential³⁷ approximations.

Radiative Transfer

The drawback of an analytical theory of multiple scattering is that it is too complicated; for systems with volume disorder often encountered in realistic situations, the scattering phenomena depends essentially on the ratio between the characteristic length scales of the system and the radiation wavelength. A statistical description in terms of such characteristic scattering lengths is usually sufficient. In general, the particular location, orientation, and size of a scattering center is irrelevant and the underlying wave character seems to be washed out. Because energy is transported through multiple scattering processes, what matters is only the energy balance. Of course, this approach cannot account for subtle interference and correlation effects but refinements can be developed on the basis of a microscopic interpretation of radiative transfer.³⁸

A comprehensive mathematical description of the nonstationary radiative transport is given by both Chandrasekhar³⁹ and Ishimaru.¹⁸ The net effect of monochromatic radiation flow through a medium with a density ρ of scattering centers is expressed in terms of a modified Stokes vector

$$\mathbf{l}(\mathbf{r}, \mathbf{s}, t) = [\langle E_1^* E_1 \rangle + \langle E_2^* E_2 \rangle, \langle E_1^* E_1 \rangle - \langle E_2^* E_2 \rangle, \text{Re}\{\langle E_1^* E_2 \rangle\}, \text{Im}\{\langle E_1^* E_2 \rangle\}]^T$$

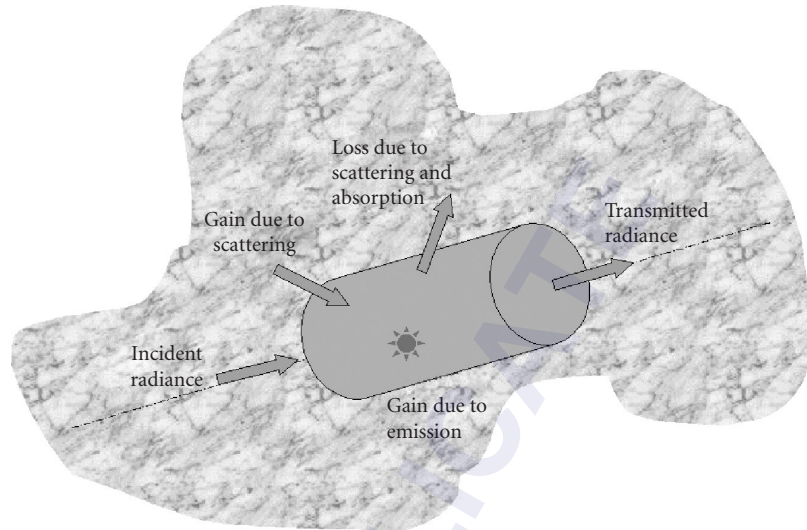


FIGURE 5 Loss-gain balance in a differential volume element; the incident radiance is attenuated by both absorption inside and scattering out of the volume element but it can be increased by ambient scattering and emission within the volume.

This quantity is a vectorial equivalent to the radiance where each element can be defined as the amount of energy in a given state that, at the position \mathbf{r} , flows per second and per unit area in the direction \mathbf{s} . When radiation propagates over the distance ds , there is a loss of specific intensity due to both scattering and absorption $dI = -\rho(\sigma_{sc} + \sigma_{abs})I ds$. In the mean time, there is a gain of specific intensity due to scattering from a generic direction \mathbf{s}' into the direction \mathbf{s} quantified by the tensorial scattering phase function (Mueller matrix) $P(\mathbf{s}', \mathbf{s})$. Also, there could be an increase, $\epsilon(\mathbf{r}, \mathbf{s}, t)$, of specific intensity due to emission within the volume of interest and the net loss-gain balance which is illustrated in Fig. 5, and represents the nonstationary radiative transfer equation:¹⁸

$$\left[\frac{1}{c} \frac{\partial}{\partial t} + \mathbf{s} \cdot \nabla + \rho(\sigma_{sc} + \sigma_{abs}) \right] I(\mathbf{r}, \mathbf{s}, t) = \rho \sigma_{sc} \int P(\mathbf{s}', \mathbf{s})(\mathbf{r}, \mathbf{s}, t) d\Omega + \epsilon(\mathbf{r}, \mathbf{s}, t) \quad (13)$$

No analytical solution exists for the transfer equation and, in order to solve specific problems, one needs to assume functional forms for both the phase function and the specific intensity. Various methods have been used to approach the transient scalar RTE.^{40–43} Cheung and Ishimaru⁴⁴ and Kim et al.⁴⁵ approached the steady state vector RTE using a Fourier analysis. Vaillon et al.⁴⁶ used a vector Monte Carlo method to analyze the radiative transfer in a particle-laden semitransparent medium. Jiang et al. presented a model for the atmospheric radiative transfer with polarization for remote-sensing applications.⁴⁷ Ma and Ishimaru used an eigenvalue-eigenfunction technique to solve numerically the vector radiative transfer equation.⁴⁸ To solve the problem of polarized pulse propagation in random media, Ishimaru et al.⁴⁹ used the discrete-ordinates method by expanding the Stokes vector in a Fourier series. Successive orders of approximation are obtained by spherical harmonic expansion of the specific intensity; for instance, the so-called P_1 approximation is obtained when the diffuse radiance is expressed as a linear combination of an isotropic radiance and a second term modulated by a cosine.¹⁸

Diffusion Approximation Perhaps one of the most widely used treatment for multiple light scattering is the diffusion approach. When (i) absorption is small compared to scattering, (ii) scattering is almost isotropic, and (iii) the radiance is not needed close to the source or boundaries

then the diffusion theory can be used as an approximation following from the general radiative transfer theory. To get insight into the physical meaning of this approximation it is convenient to define quantities that are directly measurable such as the diffuse energy density (average radiance) $U(\mathbf{r}, t) = \int_{4\pi} I_0(\mathbf{r}, \mathbf{s}, t) d\Omega$ and the diffuse flux $\mathbf{J}(\mathbf{r}, t) = \int_{4\pi} I(\mathbf{r}, \mathbf{s}, t) \mathbf{s} d\Omega$. In the diffusion approximation, the diffuse radiance is approximated by the first two terms of a Taylor's expansion:^{8,50}

$$I_0(\mathbf{r}, \mathbf{s}, t) \approx U(\mathbf{r}, t) + \frac{3}{4\pi} \mathbf{J}(\mathbf{r}, t) \cdot \mathbf{s} \quad (14)$$

and the following differential equation can be written for the average radiance

$$D\nabla^2 U(\mathbf{r}, t) - \mu_a U(\mathbf{r}, t) - \frac{\partial U(\mathbf{r}, t)}{\partial t} = S(\mathbf{r}, t) \quad (15)$$

The isotropic source density is denoted by $S(\mathbf{r}, t)$ and D is the diffusion coefficient which is defined in units of length as

$$D = \frac{1}{3[\mu_a + \mu_s(1-g)]} \quad (16)$$

in terms of the absorption μ_a and μ_s scattering coefficients. The diffusion equation is solved subject to boundary conditions and source specifics; most appealing is the fact that analytical solutions can be obtained for reflectance and transmittance calculations.

Because the phase function is characterized by a single anisotropy factor, the diffusion approximation provides mathematical convenience. Through renormalization, an asymmetry-corrected scattering cross-section that depends only on the average cosine of scattering angle defines the diffusion coefficient in Eq. (16) and, therefore, an essentially anisotropic propagation problem is mapped into an almost isotropic (diffusive) model.

The photon migration approach based on the diffusion approximation has been very successful in describing the interaction between light and complex fluids⁵¹ or biological tissues.^{52,53} It is instructive to note that three length scales characterize the light propagation in this regime: the absorption length $l_a = \mu_a^{-1}$ which is the distance traveled by a photon before it is absorbed, the scattering length $l_s = \mu_s^{-1}$ which is the average distance between successive scattering events, and the transport mean free path $l^* = l_s/(1-g)$ that defines the distance traveled before the direction of propagation is randomized. In experiments that are interpreted in the frame of the diffusion approximation, l^* is the only observable quantity and, therefore, the spatial and temporal resolution are limited by l^* and l^*/c , respectively.

Under appropriate boundary conditions, such as a mixed boundary condition in which the diffuse energy density vanishes linearly on a plane, the steady state diffusion equation can be solved and the photon flux is obtained from Fick's law.⁵⁴ Assuming an average energy transport velocity, the path length dependence of the energy flux can be evaluated yielding a path length probability distribution, $p(s)$, which can be regarded as the probability distribution of optical path lengths that correspond to waves that have traveled through the medium along closed loops and have also accumulated a total momentum transfer equal to $4\pi/\lambda$.

Low-Order Flux Models for Radiative Transfer In an effort to describe the optical properties of highly scattering materials while reducing the computational difficulties, simplifying *flux models* have been designed for the radiative energy transport. A volume scattering medium consisting of a collection of scattering centers is described as homogeneous material characterized by effective scattering and absorption properties that are determined by its internal structure.

In this approach, the fundamental equation of radiative transfer is based on the balance between the net flux change, the flux input, and flux continuing out in an infinitesimal volume. Assuming two diffusing components, a one-dimensional model based on plane symmetry for unit cross-section has been initially proposed by Schuster.⁵⁵ One of the most successful extensions of this model is the so-called

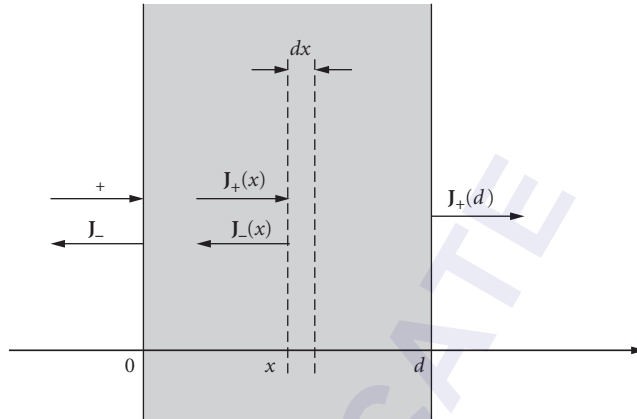


FIGURE 6 The two-flux approximation applied to radiative transfer through a diffusive layer.

Kubelka-Munk theory⁵⁶ which relates the phenomenological, effective scattering S_{K-M} and absorption K_{K-M} coefficients to measurable optical properties such as diffuse reflectance or transmittance.

The two-flux model, Kubelka-Munk theory, is schematically illustrated in Fig. 6. Diffuse radiation is assumed to be incident on the slab; the diffuse radiant flux in the positive x direction is J_+ while the one returning as a result of scattering is J_- . The net flux balance at a distance x across an infinitesimal layer of thickness dx is

$$\begin{aligned} dJ_+ &= -(K_{K-M} + S_{K-M})J_+ dx + S_- dx \\ dJ_- &= +(K_{K-M} + S_{K-M})J_- dx - S_+ dx \end{aligned} \quad (17)$$

where the coefficient K_{K-M} determines the flux attenuation due to absorption while S_{K-M} accounts for the net flux scattered between forward and backward directions. The solution of the simultaneous differential equations of the first order in one dimension is obtained by applying the boundary conditions for the layer shown in Fig. 6. The diffuse reflectance of the substrate at $x=d$ can also be accounted for and expressions for the total diffuse reflection and transmission are found for a specific application. In practical applications of the Kubelka-Munk theory, the effective scattering and absorption parameters are inferred by iteration from measurements of diffuse reflectance or transmission.

A considerable body of work was dedicated to relate the Kubelka-Munk parameters to microstructure and to incorporate both the single- and multiple-scattering effects. Refinements and higher-order flux models have also been developed. A more accurate model that accounts for the usual condition of collimated incident radiation was elaborated by J. Reichman.⁵⁷ A four-flux model has been developed that includes certain anisotropy of the scattered radiation.⁵⁸ A six-flux model was implemented to incorporate the effect of particle shape and interparticles correlation.⁵⁹ In spite of the fact that it is based on empirical determination of coefficients and that its range of applicability is rather unclear, the simple-to-implement Kubelka-Munk theory makes reasonably good description of experiments and has found applications in areas such as coatings, paper, paints, pigments, medical physics, and atmospheric physics.

Specific Effects in Multiple Light Scattering

The effects associated with light propagation through multiple scattering media depend on the scale of observation. An interesting analogy exists between optical and electron waves propagation in mesoscopic systems and, based on recent developments in solid state physics, useful insights have

been provided for a range of multiple scattering phenomena.^{7,60} It is clear now that multiple light scattering in volume disorder does not merely scramble a coherent incident wave. In fact, the apparent randomness of the scattered field conceals intriguing and sometimes counter intuitive effects such as the enhancement of the backscattered intensity and a range of correlations and statistics of vector waves.

Weak Localization Until recently, the coherent light propagating through random media has been considered to be somehow *degraded*, losing its coherence properties. However, recent experiments have brought evidence of the enhanced backscattering of light due to interference between the waves taking time-reversed paths,⁶¹ an effect which is associated with the more general phenomenon for the weak localization of waves in random media.^{62,63} First found in solid state physics more than 30 years ago, this phenomenon is applicable to all classical waves. Surveys of the state of the art in the optical counterpart of weak localization can be found in Refs. 64 and 65.

When coherent light is scattered by a medium with volume randomness, interference effects between the scattered waves which traveled through the medium along different paths occur. The result is a random pattern of interference called *laser speckle*. Because the correlations in this granular pattern extend over angular scales of typically 10^{-3} rad or less, when the individual scatterers are allowed to move over distances of the order of the wavelength or more, the distribution of intensities in the speckle pattern is rapidly averaged out and becomes essentially flat. So far, however, it is well understood and widely recognized that one kind of interference still survives in this average. This is the interference of the waves emerging from the medium in directions close to exact backscattering and which have traveled along the *same path* but in *opposite* directions.

The pairs of time-reversed light paths have some particularities which can be easily understood in the general context of the waves scattered by random media. In Fig. 7, the main contributions to scattering from a dense distribution of scatterers are presented together with their angular dependence. In the narrow angle of interest, the single scattering $I^{(s)}$ and the ladder term of multiple scattering $I^{(ml)}$ are practically constant. The third, cyclical term $I^{(mc)}$, however, corresponds to the paired (or coherent) scattering channels and, being an interference term, has a definite angular structure.

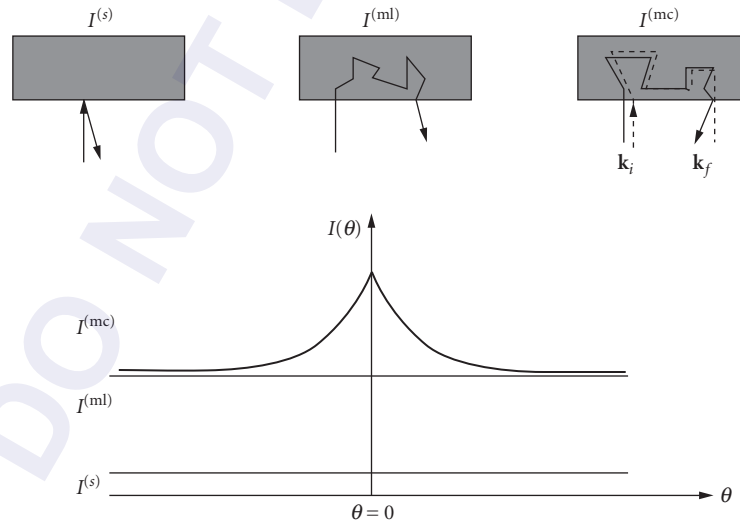


FIGURE 7 A schematic illustration of the origin of coherent backscattering. The classical contributions to backscattering are $I^{(s)}$ and $I^{(ml)}$; in addition, constructive interference occurs between reversed multiple scattering paths that have the same incident wave vector \mathbf{k}_i and exit wave vector \mathbf{k}_f .

The existence of such a cyclical term (an idea originated from Watson⁶⁶) is based on the fact that each scattering channel, involving a multitude of scattering centers, has its own coherent channel corresponding to the same sequence of scattering centers but time reversed. In the backward direction, the contribution of $I^{(\text{mc})}$ equals that of $I^{(\text{ml})}$ but its magnitude vanishes quickly away from this direction.

The angular profile of $I^{(\text{mc})}$ can be qualitatively described by taking into account the interference between two light paths as shown in Fig. 7. The two time reversed paths have the same initial and final wave vectors, \mathbf{k}_i and \mathbf{k}_f , and develop inside the medium through the same scattering centers. Under stationary conditions, the two outgoing waves corresponding to these paths are coherent and can interfere constructively for a special choice of \mathbf{k}_i and \mathbf{k}_f . If the positions of the first and the last scatterer in the sequence are \mathbf{r}_i and \mathbf{r}_f , respectively, the total phase shift between the two waves is $\mathbf{q}(\mathbf{r}_i - \mathbf{r}_f)$, where \mathbf{q} is the momentum transfer $\mathbf{k}_i + \mathbf{k}_f$. Close to the backward direction ($\mathbf{k}_f = -\mathbf{k}_i$), the two waves add coherently and the interference may be described by a weighting factor $\cos[(\mathbf{k}_i + \mathbf{k}_f)(\mathbf{r}_i - \mathbf{r}_f)]$, which is controlled by the interparticle distance $|\mathbf{r}_i - \mathbf{r}_f|$. One can say that the coherence between the time-reversed sequences is lost for angles $\theta > \lambda/|\mathbf{r}_i - \mathbf{r}_f|$ and this actually sets the angular width of the cyclical term $I^{(\text{mc})}$. The detailed angular profile of $I^{(\text{mc})}$ is determined by the probability distribution function for $|\mathbf{r}_i - \mathbf{r}_f|$ and, based on a diffusive model for light propagation in a semi-infinite medium, an approximate formula was given in Ref. 67:

$$I^{(\text{mc})} = [1 - e^{-3.4ql^*}] / 3.4ql^* \quad (18)$$

It should be noted that the intensities which contribute to classical backscattering, i.e., $I^{(s)}$ and $I^{(\text{ml})}$, correspond to incoherent channels of scattering, they add up on an intensity basis and, upon ensemble average, all angular dependences are washed out as seen in Fig. 7. As can be seen from Eq. (18), the angular shape of the coherent backscattering peak can be used to measure l^* for a specific multiple scattering medium.

Correlations in Speckle Patterns One of the significant discoveries in mesoscopic physics is the phenomenon of conductance fluctuations which arises from correlations between the transmission probabilities for different output and input modes. In multiple light scattering, the same phenomenon shows up in the form of correlations between the intensities of light transmitted in different directions such as the case schematically depicted in Fig. 8. A nice feature of the optical scattering is that in contrast with electronic conductance experiments, one has access to both angular

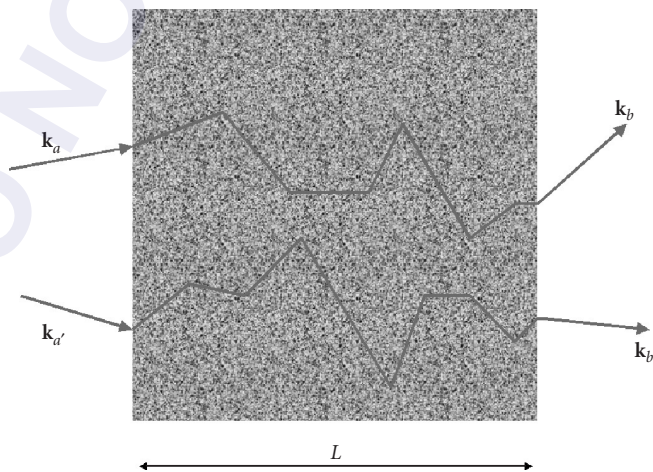


FIGURE 8 Scattering geometry containing two different wave trajectories.

dependence and angular integration of transmittance. A rich assortment of correlation functions of the different transmission quantities can be studied and their magnitudes, decay rates, etc., open up novel possibilities for multiple light scattering-based characterization and tomographic techniques.^{68,69} These measurements exploit the spatial, spectral, and temporal resolution accessible in optical experiments.⁷⁰

In a medium with volume disorder, optical waves can be thought to propagate through *channels* (propagating eigenmodes) defined angularly by one coherence area having the size of a speckle spot. The energy is coupled in and out of the medium through a number of $2\pi A/\lambda^2$ such channels where A is the transversal size of the medium. The transmission coefficient T_{ab} of one channel is defined as the ratio between the transmitted intensity in mode b and the incident intensity in mode a at a fixed optical frequency ω . An average transmission through the medium $\langle T_{ab} \rangle$ can be evaluated by ensemble averaging over different realizations of disorder and, for purely elastic scattering, it can be shown that $\langle T_{ab} \rangle \sim I^*/L$. The motion of scatterers, the frequency shift of the incident wave, and the variation of angle of incidence and/or detection, introduce phase shifts (momentum differences) in different propagating channels. Surprising and novel features are found when one studies carefully how various T_{ab} channels are correlated with one another. A general function $C_{aba'b'} = \langle \delta T_{ab} \delta T_{a'b'} \rangle$ can be designed to evaluate the correlation between changes in the transmission coefficients $\delta T_{ab} = T_{ab} - \langle T_{ab} \rangle$ and, to the lowest order in the disorder parameter $1/kl$, it can be shown to be a summation of three different terms.⁷¹ The first term C^1 , *short-range correlations*, represents the large local intensity fluctuations specific to speckle patterns and exhibits an angular *memory* effect: when the incident beam is tilted by a small angle, the transmitted speckle pattern will, in average, follow provided that the tilt angle is not too large.^{72,73} The second term C^2 , *long-range correlations*, arises from paths that cross at a certain scattering site; it is smaller than C^1 by a factor $1/\Sigma T_{ab}$ and decays very slowly with the momentum difference.⁷⁴ Finally, a uniform positive correlation which is independent of momentum differences is included in C^3 , *conductance correlations*, and it is determined by the less probable case of two crossings in propagation channels. This small correlation term (of the order of $1/\Sigma (T_{ab})^2$) causes just a shift in background, i.e., the spatially averaged intensity in each speckle pattern is always a little darker or brighter than the total intensity averaged over many disorder realizations in the sample. These fluctuations do not decrease when averaging is done over larger and larger spatial regions and they are the optical analogue of the universal conductance fluctuations in electronic systems.^{71,75}

Depolarization A common interpretation of multiple scattering effects assumes depolarization of the incident field. However, when coherent radiation interacts with scattering media, there is always a fixed phase and amplitude relationship between orthogonal electric field components at a given frequency at any point in time and space. Of course, these relationships may vary as a function of time, spatial coordinate, frequency, and material morphology. If, upon ensemble average, there is no correlation between any pair of orthogonal field components, the field is said to be unpolarized. A degree of polarization is defined as

$$D = \frac{(I_1^2 + I_2^2 + I_3^2)^{1/2}}{I_0} \quad (19)$$

where I_i are the ensemble averaged Stokes vector elements. This ensemble average can take the form of spatial, temporal, or frequency average or can also be the result of averaging over different material realizations.⁷⁶ For a static medium illuminated by polarized light, depolarization occurs when path length differences contributing to a particular point exceed the coherence length of the illuminating light. It is interesting to note that the measured degree of polarization will be a function of the detection geometry, decreasing with increasing the detector's size and integration time.

Based on symmetry considerations, van de Hulst² finds the Mueller matrix for single scattering on a collection of randomly oriented identical particles each of which has a plane of symmetry to be diagonal with $P_{22} = P_{33} \neq P_{44}$. For spheres in exact forward scattering $P_{22} = P_{44} = 1$. In multiple scattering

however, this relation is not true anymore. It is expected that when light of arbitrary incident polarization impinges on an optically thick, multiple scattering medium it emerges diffusely and totally depolarized. When increasing the optical density, the transfer matrix evolves toward that of a total depolarizer which has all elements equal to zero except for P_{11} . This depolarization process will depend on the size parameter of the scattering particles. Owing to a smaller scattering anisotropy for the particles with the size parameter close to 1, the total depolarization stage is reached at higher optical densities than for larger particles. Two different regimes can be identified in terms of optical density d/l^* : (i) a steep decay for low optical densities, which corresponds to the attenuation of ballistic photons, and (ii) a slower decay for large optical densities, corresponding to the diffusive regime. The effective coefficient of attenuation depends only on the volume fraction of the scattering medium.

For an arbitrary input state of polarization, the output state of polarization can be obtained from $\mathbf{I}_{\text{out}} = \mathbf{P} \mathbf{I}_{\text{in}}$. In the case of a diagonal transfer matrix, the renormalized output Stokes vector is

$$\mathbf{I}_{\text{out}} = \begin{bmatrix} 1 & I_1 P_{22} & I_2 P_{33} & I_3 P_{44} \end{bmatrix}^T \quad (20)$$

The degree of polarization of the scattered light can be obtained from Eq. (20) using Eq. (19) for any input state. The characteristics of the depolarization will depend both on the input state of polarization and the dominant scattering regime, exhibiting different behavior in the Mie and Rayleigh ($ka \ll 1$) regimes. The depolarizing behavior of multiple scattering, as a function of sample thickness is exemplified in Fig. 9a and b for linear and circular inputs, respectively. As can be seen, for samples 2 and 3 in Fig. 9a and b the slope for linear input is always steeper than for circular input, indicating that circularly polarized light is less depolarized than linearly polarized light for the same sample thickness.

It is expected that, as soon as the diffusive regime is reached, multiple scattering will completely depolarize the incident optical wave. Knowing the complete Mueller matrix, the state of polarization of scattered light can be estimated for any input state of polarization. A detailed analysis can also predict which type of illumination is better preserved while propagating through the scattering medium. This is particularly important in applications such as long-range target identification where one must take into account depolarization effects due to propagation.

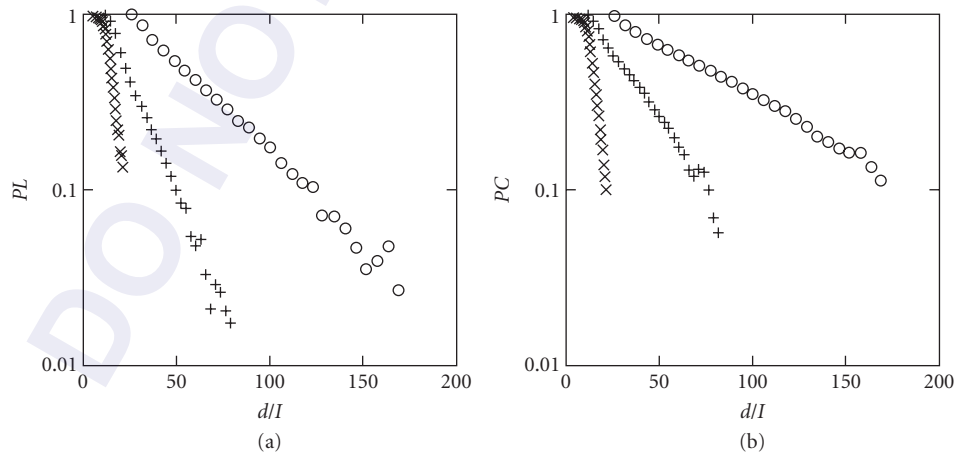


FIGURE 9 Degree of polarization of output light for (a) linear and (b) circular inputs. Symbols: x-Rayleigh scatterer, + and O-Mie particles.

9.6 REFERENCES

1. Lord Rayleigh, *Phil. Mag.* **41**:107, 274, 447, 1871; Tyndall, *Phil. Mag.* **37**:156, 1869.
2. H. C. van de Hulst, *Light Scattering by Small Particles*, Wiley, New York, 1957. Reprinted by Dover, New York, 1981.
3. M. Kerker, *The Scattering of Light*, Academic Press, New York, 1969.
4. L. P. Bayvel and A. R. Jones, *Electromagnetic Scattering and Applications*, Elsevier, London, 1981.
5. C. F. Bohren and D. R. Huffman, *Absorption and Scattering of Light by Small Particles*, Wiley-Interscience, New York, 1983.
6. A. R. Jones, "Light Scattering for Particle Characterization," *Prog. Energy Comb. Sci.* **25**:1–53, 1999.
7. W. van Haeringen and D. Lenstra, eds., *Analogies in Optics and Microelectronics*, Kluwer Academic Publishers, Netherlands, 1990.
8. A. Ishimaru, *Electromagnetic Wave Propagation, Radiation, and Scattering*, Prentice hall, New Jersey, 1991.
9. J. B. Keller, "Accuracy and Validity of the Born and Rytov Approximation," *J. Opt. Soc. Am.* **59**:1003, 1969.
10. B. Cairns and E. Wolf, "Comparison of the Born and Rytov approximations for scattering on quasi-homogeneous media," *Opt. Commun.* **74**:284, 1990.
11. T. M. Habashy, R. W. Groom, and B. R. Spies, "Beyond the Born and Rytov Approximations: A Nonlinear Approach to Electromagnetic Scattering," *J. Geophys. Res.* **98**:1759, 1993.
12. M. Born and E. Wolf, *Principles of Optics*, Cambridge University Press, Cambridge, 1999.
13. J. M. Ziman, *Models of Disorder*, Cambridge University Press, Cambridge, 1979.
14. N. E. Cusack, *The Physics of Structurally Disordered Matter*, Adam Hilger, Bristol, 1987.
15. Y. Waseda, *The Structure of Non-Crystalline Materials*, McGraw-Hill, New York, 1980.
16. J.-P. Hansen and I. McDonald, *Theory of Simple Liquids*, 2d ed., Academic Press, New York, 1986.
17. R. H. Bates, V. A. Smith, and R. D. Murch, "Manageable Multidimensional Inverse Scattering Theory," *Phys. Rep.* **201**:185–277, 1991.
18. A. Ishimaru, *Wave propagation and Scattering in Random Media*, Oxford University Press, Oxford, 1997.
19. W. Brown, ed., *Light Scattering Principles and Development*, Clarendon Press, Oxford, 1996.
20. V. Degiorgio, M. Corti, and M. Giglio, eds., *Light Scattering in Liquids and Macromolecular Solutions*, Plenum Press, New York, 1980.
21. J. Teixeira, "Experimental Methods for Studying Fractal Aggregates," p. 145, In: H. E. Stanley and N. Ostrowski, eds., *On Growth and Form*, Martinus Nijhoff, Boston, 1986.
22. J. G. Rarity, R. N. Seabrook, and R. J. G. Carr, "Light-Scattering Studies of Aggregation," *Proc. R. Soc. London A* **423**:89–102, 1989.
23. B. J. Berne and R. Pecora, *Dynamic Light Scattering*, John Wiley & Sons, New York, 1976.
24. H. Z. Cummins and E. R. Pike, eds., *Photon Correlation and Light Beating Spectroscopy*, Plenum Press, New York, 1974.
25. W. Brown, ed., *Dynamic Light Scattering*, Oxford Univ. Press., New York, 1993.
26. B. Chu, *Laser Light Scattering, Basic Principles and Practice*, Academic Press, San Diego, 1991.
27. H. C. van de Hulst, *Multiple Light Scattering—Tables, Formulas, and Applications*, Academic Press, New York, (1980).
28. J. C. Maxwell-Garnett, "Colours in Metal Glasses and in Metallic Films," *Phil. Trans. R. Soc. Lond.* **A203**:385, 1904.
29. E. N. Economou and C. M. Soukoulis, "Optical Localization: Computational Techniques and Results," In: P. Sheng ed., *Scattering and Localization of Classical Waves in Random Media*, World Scientific, Singapore, 1990.
30. A. Lagendijk and B. A. van Tiggelen, "Resonant Multiple Scattering of Light," *Phys. Reports* **270**:143–215, 1996.
31. V. Twerski, "On Propagation in Random Media of Discrete Scatterers," *Proc. Symp. Appl. Math* **16**:84–116, 1964.
32. P. C. Waterman, "Symmetry, Unitarity, and Geometry in Electromagnetic Scattering," *Phys. Rev. D* **3**:825, 1971.
33. V. K. Varadan and V. V. Varadan, eds., *Acoustic, Electromagnetic and Elastic Wave Scattering—Focus on the T-Matrix Approach*, Pergamon, New York, 1980.

34. L. Tsang, J. A. Kong, and R. T. Shin, *Theory of Microwave Remote Sensing*, Wiley, New York, 1985.
35. L. L. Foldy, "The Multiple Scattering of Waves. I. General Theory of Isotropic Scattering by Randomly Distributed Scatterers," *Phys. Rev.* **67**:107, 1945.
36. M. Lax, "Multiple Scattering of Waves. II. The Effective Field in Dense," *Phys. Rev.* **88**:621, 1952.
37. P. Soven, "Coherent-Potential Model of Substitutional Disordered Alloys," *Phys. Rev.* **156**:809–813, (1967).
38. M. C. W. van Rossum and Th. M. Nieuwenhuizen, "Multiple Scattering of Classical Waves: Microscopy, Mesoscopy, and Diffusion," *Rev. Mod. Phys.* **71**:313–371, 1999.
39. S. Chandrasekhar, *Radiative Transfer*, Dover Publ., New York, 1960.
40. R. Elaloufi, R. Carminati, and J. J. Greffet, "Time-Dependent Transport through Scattering Media: From Radiative Transfer to Diffusion," *J. Opt. A, Pure Appl. Opt.* **4**:S103–S108, 2002.
41. K. Mitra and S. Kumar, "Development and Comparison of Models for Light-Pulse Transport through Scattering-Absorbing Media," *Appl. Opt.* **38**:188–196, 1999.
42. Z. M. Tan and P. F. Hsu, "An Integral Formulation of Transient Radiative Transfer," *J. Heat Trans.* **123**:466–475, 2001.
43. M. Sakami, K. Mitra, and P. F. Hsu, "Analysis of Light-Pulse Transport through Two-Dimensional Scattering and Absorbing Media," *J. Quant. Spectrosc. Radiat. Transf.* **73**:169–179, 2002.
44. R. L. T. Cheung and A. Ishimaru, "Transmission, Backscattering, and Depolarization of Waves in Randomly Distributed Spherical Particles," *Appl. Opt.* **21**:3792–3798, 1982.
45. A. D. Kim, S. Jaruwatanadilok, A. Ishimaru, and Y. Kuga, "Polarized Light Propagation and Scattering in Random Media," *Proc. SPIE* **4257**:90–100, 2001.
46. R. Vaillon, B.T. Wong, and M. P. Menguc, "Polarized Radiative Transfer in a Particle-Laden Semi-Transparent Medium via a Vector Monte Carlo Method," *J. Quant. Spectrosc. Radiat. Transf.* **84**:383–394, 2004.
47. Y. Jiang, Y. L. Yung, S. P. Sander, and L. D. Travis, "Modeling of Atmospheric Radiative Transfer with Polarization and Its Application to the Remote Sensing of Tropospheric Ozone," *J. Quant. Spectrosc. Radiat. Transf.* **84**:169–179, 2004.
48. Q. Ma and A. Ishimaru, "Scattering and Depolarization of Waves Incident upon a Slab of Random Medium with Refractive Index Different from That of the Surrounding Media," *Radio Sci.* **25**:419–426, 1990.
49. A. Ishimaru, S. Jaruwatanadilok, and Y. Kuga, "Polarized Pulse Waves in Random Discrete Scatterers," *Appl. Opt.* **40**:5495–5502, 2001.
50. P. Sheng, *Introduction to Wave Scattering, Localization, and Mesoscopic Phenomena*, Academic Press, New York, 1995.
51. G. Maret, "Recent Experiments on Multiple Scattering and Localization of Light," In: E. Akkermanns, G. Montambaux, J.-L. Pichard, J. Zinn-Justin, eds., *Mesosopic Quantum Physics*, Elsevier Science, 1995.
52. For a collection of recent results and applications see *OSA Trends in Optics and Photonics vol. 2*, "Advances in Optical Imaging and Photon Migration," R. R. Alfano and J. M. Fujimoto eds., Optical Society of America, Washington, DC, 1996; *OSA Trends in Optics and Photonics vol. 21*, "Advances in Optical Imaging and Photon Migration," J. M. Fujimoto and M. S. Patterson eds., Optical Society of America, Washington, DC, 1998; *OSA Trends in Optics and Photonics vol. 22*, "Biomedical Optical Spectroscopy and Diagnostics," E. M. Sevick-Muraca and J. Izatt eds., Optical Society of America, Washington, DC, 1998.
53. S. A. Arridge, "Optical Tomography in Medical Imaging," *Inv. Probs.* **15**:41–93, 1999.
54. J. J. Duderstadt and L. J. Hamilton, *Nuclear Reactor Analysis*, Wiley, New York, 1976.
55. A. Schuster, "Radiation through a Foggy Atmosphere," *Astrophys. J.* **21**:1, 1905.
56. P. Kubelka and F. Munk, "Ein Beitrag zur Optik der Farbanstriche," *Z. Tech. Phys.* **12**:593–601, 1931.
57. J. Reichman, "Determination of Scattering and Absorption Coefficients for Nonhomogeneous Media," *Appl. Opt.* **12**:1811, 1973.
58. B. Maheu, J. N. Letoulouzan, and G. Gouesbet, "Four-Flux Models to Solve the Scattering Transfer Equations in Terms of Lorentz-Mie Parameters," *Appl. Opt.* **26**:3353–3361, 1984.
59. A. G. Emslie and J. R. Aronson, "Spectral Reflectance and Emittance of Particulate Materials," *Appl. Opt.* **12**:2563, 1973.
60. S. Datta, *Electronic Transport in Mesoscopic Systems*, Cambridge University Press, 1995.
61. K. Kuga and A. Ishimaru, "Retroreflectance from a Dense Distribution of Spherical Particles," *J. Opt. Soc. Am.* **A1**:831–835, 1984.

62. M. P. Albada and A. Lagendijk, "Observation of Weak Localization of Light in a Random Medium," *Phys. Rev. Lett.* **55**:2692–2695, 1985.
63. P. Wolf and G. Maret, "Weak Localization and Coherent Backscattering of Photons in Disordered Media," *Phys. Rev. Lett.* **55**:2696–2699, 1985.
64. M. P. Albada, M. B. van der Mark, and A. Lagendijk, "Experiments on Weak Localization and Their Interpretation," In: P. Sheng ed., *Scattering and Localization of Classical Waves in Random Media*, World Scientific, Singapore, 1990.
65. Y. Barabankov, Y. Kravtsov, V. D. Ozrin, and A. I. Saicev, "Enhanced Backscattering in Optics," In: E. Wolf, ed., *Progress in Optics XXIX*, North-Holland, Amsterdam, 1991.
66. K. Watson, "Multiple Scattering of Electromagnetic Waves in an Underdense Plasma," *J. Math. Phys.* **10**:688–702, 1969.
67. E. Akkermas, P. E. Wolf, R. Maynard, and G. Maret, "Theoretical Study of the Coherent Backscattering of Light in Disordered Media," *J. Phys. France* **49**:77–98, 1988.
68. R. Berkovits and S. Feng, "Correlations in Coherent Multiple Scattering," *Phys. Reports* **238**:135–172, 1994.
69. I. Freund, "'1001' Correlations in Random Wave Fields," *Waves in Random Media* **8**:119–158, 1998.
70. A. Z. Genack, "Fluctuations, Correlation and Average Transport of Electromagnetic Radiation in Random Media," In: P. Sheng ed., *Scattering and Localization of Classical Waves in Random Media*, World Scientific, Singapore, 1990.
71. S. Feng and P. A. Lee, "Mesoscopic Conductors and Correlations in Laser Speckle Patterns," *Science* **251**:633–639, 1991.
72. M. Rosenbluh, M. Hoshen, I. Freund, and M. Kaveh, "Time Evolution of Universal Optical Fluctuations," *Phys. Rev. Lett.* **58**:2754–2757, 1987.
73. J. H. Li and A. Z. Genack, "Correlation in Laser Speckle," *Phys. Rev. E* **49**:4530–4533, 1994.
74. M. P. van Albada, J. F. de Boer, and A. Lagendijk, "Observation of Long-Range Intensity Correlation in the Transport of Coherent Light through a Random Medium," *Phys. Rev. Lett.* **64**:2787–2790, 1990.
75. S. Feng, C. Kane, P. A. Lee, and A. D. Stone, "Correlations and Fluctuations of Coherent Wave Transmission through Disordered Media," *Phys. Rev. Lett.* **61**:834–837, 1988.
76. E. Collett, *Polarized Light in Fiber Optics*, PolaWave Group, Lincroft, New Jersey, 2003.

OPTICAL SPECTROSCOPY AND SPECTROSCOPIC LINESHAPES

Brian Henderson

*Department of Physics and Applied Physics
University of Strathclyde
Glasgow, United Kingdom*

10.1 GLOSSARY

A_{ba}	Einstein coefficient for spontaneous emission
a_0	Bohr radius
B_{if}	Einstein coefficient between initial state, $ i\rangle$, and final state, $ f\rangle$
E_{DC}	Dirac Coulomb term
E_{hf}	hyperfine energy
E_n	eigenvalue of quantum state n
$E(t)$	electric field at time t
$E(\omega)$	electric field at frequency ω
e	charge on the electron
ED	electric dipole term
EQ	electric quadrupole term
$\langle f V' i\rangle$	matrix element of perturbation V'
g_a	degeneracy of ground level
g_b	degeneracy of excited level
g_N	gyromagnetic ratio of nucleus
H_{so}	spin-orbit interaction Hamiltonian
\hbar	Planck's constant
I	nuclear spin
$I(t)$	emission intensity at time t
\mathbf{j}	total angular momentum vector given by $\mathbf{j} = \mathbf{l} \pm \frac{1}{2}$
l_i	orbital state
M_N	mass of nucleus N
MD	magnetic dipole term
m	mass of the electron

$n_{\omega}(T)$	equilibrium number of photons in a blackbody cavity radiator at angular frequency ω and temperature T
QED	quantum electrodynamics
R_{∞}	Rydberg constant for an infinitely heavy nucleus
$R_{nl}^{(r)}$	radial wavefunction
s	spin quantum number with the value $\frac{1}{2}$
s_i	electron spin
T	absolute temperature
W_{ab}	transition rate in absorption transition between states $ a\rangle$ and $ b\rangle$
W_{ba}	transition rate in emission transition from state $ b\rangle$ to state $ a\rangle$
Z	charge on the nucleus
$\alpha = e^2/4\pi\epsilon_0\hbar c$	fine structure constant
ϵ_0	permittivity of free space
μ_B	Bohr magneton
$\rho(\omega)$	energy density at frequency ω
$\zeta(r)$	spin-orbit parameter
τ_R	radiative lifetime
ω	angular frequency
$\Delta\omega$	natural linewidth of the transition
$\Delta\omega_D$	Doppler width of transition
ω_k	mode k with angular frequency ω

Spectroscopic measurements have played a key role in the development of quantum theory. This chapter presents a simple description of the quantum basis of spectroscopic phenomena, as a prelude to a discussion of the application of spectroscopic principles in atomic, molecular, and solid-state physics. A brief survey is presented of the multielectron energy-level structure in the three phases of matter and of the selection rules which determine the observation of optical spectra. Examples are given of the fine-structure, hyperfine-structure, and spin-orbit splittings in the spectra of atoms, molecules, and solids. Solid-state phenomena considered will include color center, transition metal, and rare earth ion spectra.

The intrinsic or homogeneous lineshapes of spectra are determined by lifetime effects. Other dephasing processes, including rotational and vibrational effects, lead to splitting and broadening of spectra. There are also sources of inhomogeneous broadening associated with Doppler effects in atomic and molecular spectra and crystal field disorder in solids. Methods of recovering the homogeneous lineshape include sub-Doppler laser spectroscopy of atoms, optical hole burning, and fluorescence line narrowing.

Finally, the relationship between linewidth and lifetime are discussed and the effects of time-decay processes outlined. The consequences of measurements in the picosecond and subpicosecond regime are described. Examples of vibrational relaxation in molecular and solid-state spectroscopy are reviewed.

10.2 INTRODUCTORY COMMENTS

Color has been used to enhance the human environment since the earliest civilizations. Cave artists produced spectacular colorations by mixing natural pigments. These same pigments, burned into the surfaces of clays to produce color variations in pottery, were also used to tint glass. The explanation of the coloration process in solids followed from Newton's observation that white light contains all the colors of the rainbow,¹ the observed color of a solid being complementary to that absorbed

from white light by the solid. Newton measured the wavelength variation of the refractive index of solids, which is responsible for dispersion, and his corpuscular theory of light explained the laws of reflection and refraction.¹ The detailed interpretation of polarization, diffraction, and interference followed from the recognition that light was composed of transverse waves, the directions of which were related to the direction of the electric field in Maxwell's electromagnetic theory:^{2,3} the electronic constituents of matter are set into transverse oscillation relative to the propagating light beam. Subsequently, Einstein introduced the photon in explaining the photoelectric effect.⁴ Thus the operating principles of optical components in spectrometers, such as light sources, mirrors, lenses, prisms, polarizers, gratings, and detectors, have been with us for a long time.

Many significant early developments in quantum physics led from optical spectroscopic studies of complex atoms. After Bohr's theory of hydrogen,⁵ the quantum basis of atomic processes developed apace. One of Schrödinger's first applications of wave mechanics was in the calculations of atomic energy levels and the strengths of spectroscopic transitions.⁶ Schrödinger also demonstrated the formal equivalence of wave mechanics and Heisenberg's matrix mechanics. Extensions of spectroscopy from atomic physics to molecular physics and solid-state physics more or less coincided with the early applications of quantum mechanics in these areas.

A survey of the whole of spectroscopy, encompassing atoms, molecules, and solids, is not the present intent. Rather it is hoped that by choice of a few critical examples the more general principles linking optical spectroscopy and the structure of matter can be demonstrated. The field is now vast: Originally the exclusive domain of physicists and chemists, optical spectroscopy is now practiced by a variety of biophysicists and biochemists, geophysicists, molecular biologists, and medical and pharmaceutical chemists with applications to proteins and membranes, gemstones, immunoassay, DNA sequencing, and environmental monitoring.

10.3 THEORETICAL PRELIMINARIES

The outstanding success of the Bohr theory was the derivation of the energy-level spectrum for hydrogenic atoms:

$$E_n = -\frac{mZ^2e^4}{2(4\pi\epsilon_0)^2n^2\hbar^2} = -\frac{Z^2hc}{n^2}R_\infty \quad (1)$$

Here the principal quantum number n is integral; $h = 2\pi\hbar$ is Planck's constant; Z is the charge on the nucleus, m and e are, respectively, the mass and charge on the electron; and ϵ_0 is the permittivity of free space. The Rydberg constant for an infinitely heavy nucleus, R_∞ , is regarded as a fundamental atomic constant with approximate value $10,973,731 \text{ m}^{-1}$. Equation (1) is exactly the relationship that follows from the boundary conditions required to obtain physically realistic solutions for Schrödinger's time-independent equation for one-electron atoms. However, the Schrödinger equation did not account for the *fine structure* in the spectra of atoms nor for the splittings of spectral lines in magnetic or electric fields.

In 1927 Dirac developed a relativistic wave equation,⁷ which introduced an additional angular momentum for the spinning electron of magnitude $s\hbar$, where $s^* = \sqrt{s(s+1)}$ and the spin quantum number s has the value $s = \frac{1}{2}$. The orbital and spin angular momenta are coupled together to form a total angular momentum vector \mathbf{j} , given by $\mathbf{j} = \mathbf{l} \pm \frac{1}{2}$. In the hydrogenic ground state $|n\mathbf{l}\rangle = |10\rangle$, this spin-orbit coupling yields a value of $\mathbf{j} = \frac{1}{2}$ only, giving the $1S_{1/2}$ level. In the first excited state, for which $n = 2$ the $l = 0$, $\mathbf{j} = \frac{1}{2}$ is represented by $2S_{1/2}$, while $l = 1$ leads to $\mathbf{j} = \frac{3}{2}$ ($2P_{3/2}$) and $\mathbf{j} = \frac{1}{2}$ ($2P_{1/2}$) levels, these two levels being separated by the fine structure interval.⁸ The Dirac form of the Coulomb energy, expressed as an expansion in powers of $Z \times$ the fine structure constant, $\alpha = (e^2/4\pi\epsilon_0\hbar c)$, is then

$$E_{\text{DC}} = -\frac{Z^2}{n^2}R_\infty hc \left[1 + \frac{(Z\alpha)^2}{n} \left(\frac{1}{\mathbf{j} + \frac{1}{2}} - \frac{3}{4n} \right) + O((Z\alpha)^4) \right] \quad (2)$$

The second term in the bracket in Eq. (2) is the spin-orbit correction to the energies, which scales as $(Z^4\alpha^2)/n^3$. In the case of hydrogenic atoms this relativistic coupling removes the $nP_{1/2} - nP_{3/2}$ and $nD_{3/2} - nD_{5/2}$ degeneracy, but does not split the $nS_{1/2}$ level away from the $nP_{1/2}$ level. A further relativistic correction to Eq. (1) involves replacing the electronic mass in R_∞ by the reduced mass of the electron $\mu = mM/(M + m)$, which introduces a further shift of order $(m/M)_N(1 - (Z\alpha/2n)^2)E_{DC}$. Here M_N is the mass of the nucleus.

There are two further energy-level shifts.⁹ The so-called quantum electrodynamic (QED) shifts include contributions due to finite nuclear size, relativistic recoil, and radiative corrections, collectively described as the Lamb shift, as well as terms due to electron self-energy and vacuum polarization. The Lamb shift raises the degeneracy of the $nS_{1/2} - nP_{1/2}$ levels. Overall, the QED shift scales as $\alpha(Z\alpha)^4/n^3$.¹⁰ The interaction of the electronic and nuclear magnetic moments gives rise to hyperfine structure in spectra. The hyperfine contribution to the electronic energies for a nucleus of mass M_N , nuclear spin I , and gyromagnetic ratio g_N is given by

$$E_{hf} = \alpha^2 \left(\frac{Z^3}{n^3} \right) \left(\frac{g_N m}{M_N} \right) hcR_\infty \frac{F(F-1) - I(I+1) - j(j+1)}{j(j+2)(2I+1)} \quad (3)$$

where $\mathbf{j} = \mathbf{l} + \mathbf{s}$ is the total electronic angular momentum and $\mathbf{F} = \mathbf{I} + \mathbf{j}$ is the total atomic angular momentum. E_{hf} scales as $Z^3\alpha^2/n^3$ and is larger for S-states than for higher-orbit angular momentum states. More generally, all the correction terms scale as some power of Z/n , demonstrating that the shifts are greatest for $n = 1$ and larger nuclear charge. Experiments on atomic hydrogen are particularly important, since they give direct tests of relativistic quantum mechanics and QED.

10.4 RATES OF SPECTROSCOPIC TRANSITION

The rates of transitions may be determined using time-dependent perturbation theory. Accordingly, it is necessary to consider perturbations which mix stationary states of the atom. The perturbations are real and oscillatory in time with angular frequency ω and have the form

$$H_1 = V \exp(-i\omega t) + V^* \exp(i\omega t) \quad (4)$$

where V is a function only of the spatial coordinates of the atom. In the presence of such a time-dependent perturbation, the Schrödinger equation

$$(H_0 + H_1)\Psi = i\hbar \frac{\delta\Psi}{\delta t} \quad (5)$$

has eigenstates

$$\Psi = \sum_n c_n(t) |nlm\rangle \exp(-iE_n t/\hbar) \quad (6)$$

which are linear combinations of the n stationary solutions of the time-independent Schrödinger equation, which have the eigenvalues E_n . The time-dependent coefficients, $c_m(t)$, indicate the extent of mixing between the stationary state wavefunctions $|nlm\rangle$. The value of $|c_j(t)|^2$, the probability that the electronic system, initially in state i , will be in a final state f after time t is given by

$$|c_f(t)|^2 = 4 \left| \frac{V_{fi}}{\hbar} \right|^2 \frac{\sin^2 \frac{1}{2}(\omega_{fi} - \omega)t}{(\omega_{fi} - \omega)^2} \quad (7)$$

for an absorption process in which the final state f is higher in energy than the initial state i . This expression defines the Bohr frequency condition,

$$\hbar\omega = E_f - E_i \quad (8)$$

and $\omega_{fi} = (E_f - E_i)/\hbar$. Obviously, $|c_f(t)|^2$ has a maximum value when $\omega_{fi} = \omega$, showing that the probability of an absorption transition is maximum when $E_f - E_i = \hbar\omega$. The emission process comes from the V^* exp (ωt) term in Eq. (4): the signs in the numerator and denominator of Eq. (7) are then positive rather than negative. For the probability to be significant then requires that $\omega_{fi} + \omega = 0$, so that the final state f is lower in energy than the initial state. If the radiation field has a density of oscillatory modes $u(\omega)$ per unit frequency range, then Eq. (7) must be integrated over the frequency distribution. The transition rate is then

$$W_{fi} = \frac{2\pi}{\hbar^2} |V_{fi}^\omega|^2 u(\omega_{fi}) \quad (9)$$

in which the V_{fi}^ω indicates that only a narrow band of modes close to $\omega = \omega_{fi}$ has been taken into account in the integration. This equation, which gives the probability of a transition from $|i\rangle \rightarrow |f\rangle$ per unit time, is known as Fermi's golden rule.

In Eq. (9) $V_{fi} = \langle f|V|i\rangle$ and $V_{fi}^* = \langle f|V^*|i\rangle$ determine the transition probabilities for absorption and emission between the initial i and final states f . In fact $|V_{fi}^\omega|$ and $|V_{if}^{\omega*}|$ are identical and the transition probabilities for absorption and emission are equal. For the k th mode with angular frequency ω_k the perturbation takes the form

$$V_k^\omega \equiv \sum_i \left(er_i \cdot E_k^0 + \frac{e}{2m} (l_i + 2s_i) \cdot B_k^0 + \frac{1}{2} er_i \cdot r_i \cdot k E_k^0 \right) \quad (10)$$

The first term in Eq. (10) is the electric dipole (ED) term. The second and third terms are the magnetic dipole (MD) and electric quadrupole terms (EQ), respectively. The relative strengths of these three terms are in the ratio $(ea_0)^2 : (\mu_b/c)^2 : ea_0^2/\lambda^2$ where a_0 and μ_b are the Bohr radius and Bohr magneton, respectively. These ratios are then approximately $1 : 10^{-5} : 10^{-7}$. Since the electromagnetic energy per unit volume contained in each mode, including both senses of polarization, is given by $2\epsilon_0\kappa|E_k^0|^2$, the energy density, $\rho(\omega)$, per unit volume per unit angular frequency is just $4\epsilon_0\kappa|E_k^0|^2 u_k(\omega)$. Hence, from Eq. (9) and using only the first term of Eq. (10) the electric dipole transition rate is determined as

$$W_{if} = \frac{\pi}{2\epsilon_0\kappa\hbar^2} \sum_k \left| \langle f | \sum_i er_i \cdot \hat{\epsilon}_k | i \rangle \right|^2 \rho(\omega) \quad (11)$$

where the summations are over the numbers of electrons, i , and polarization vectors, \mathbf{k} . For randomly polarized radiation Eq. (11) becomes

$$W_{if} = \frac{2\pi}{6\epsilon_0\kappa\hbar^2} \left| \langle f | \sum_i er_i | i \rangle \right|^2 \rho(\omega) \quad (12)$$

If the radiation has all the \mathbf{E}_k vectors pointing along the z -direction then only this mode is taken into account and

$$W_{if} = \frac{\pi}{2\epsilon_0\kappa\hbar^2} \left| \langle f | \sum_i ez_i | i \rangle \right|^2 \rho(\omega) \quad (13)$$

These relationships, Eqs. (12) and (13), are used subsequently in discussing experimental techniques for measuring optical absorption and luminescence spectra. They result in the selection rules that govern both polarized and unpolarized optical transitions.

For the most part the succeeding discussion is concerned with radiative transitions between the ground level a and an excited level b . These levels have degeneracies g_a and g_b with individual ground and excited states labeled by $|a_n\rangle$ and $|b_m\rangle$, respectively. The probability of exciting a transition from state $|a_n\rangle$ to state $|b_m\rangle$ is the same as that for a stimulated transition from $|b_m\rangle$ to $|a_n\rangle$. The transition rates in absorption, W_{ab} , and in emission, W_{ba} , are related through

$$g_a W_{ab} = g_b W_{ba} \quad (14)$$

assuming the same energy density for the radiation field in absorption and emission. Since the stimulated transition rate is defined by

$$W_{ab} = B_{ab}\rho(\omega) \quad (15)$$

the Einstein coefficient B_{ab} for stimulated absorption is directly related to the squared matrix element $|\langle b_m | \Sigma_i | e r_i | a_n \rangle|^2$. Furthermore, the full emission rate is given by

$$W_{ba} = A_{ba}[1 + n_\omega(T)] \quad (16)$$

where $n_\omega(T)$ is the equilibrium number of photons in a blackbody cavity radiator at angular frequency ω and temperature T . The first term in Eq. (16) (i.e., A_{ba}) is the purely spontaneous emission rate, related to the stimulated emission rate by

$$A_{ba} = 2B_{ba} u_k(\omega) \hbar \omega_k \quad (17)$$

Equation (17) shows that the spontaneous transition probability is numerically equal to the probability of a transition stimulated by one photon in each electromagnetic mode, \mathbf{k} . Similarly the stimulated absorption rate is given by

$$W_{ab} = B_{ab}\rho(\omega) = \frac{g_b}{g_a} A_{ba} n_\omega(T) \quad (18)$$

These quantum mechanical relationships show how the experimental transition rates, for both polarized and unpolarized radiation, are determined by the mixing of the states by the perturbing oscillatory electric field. Since the radiative lifetime τ_R is the reciprocal of the Einstein A coefficient for spontaneous emission [i.e., $\tau_R = (A_{ba})^{-1}$] we see the relationship between luminescence decaytime and the selection rules via the matrix element $|\langle b_m | \Sigma_i e r_i | a_n \rangle|$.

10.5 LINESHAPES OF SPECTRAL TRANSITIONS

Consider the excitation of optical transitions between two nondegenerate levels $|a\rangle$ and $|b\rangle$. The instantaneous population of the upper level at some time t after the atomic system has been excited with a very short pulse of radiation of energy $\hbar\omega_{ab}$ is given by

$$N_b(t) = N_b(0) \exp(-A_{ba}t) \quad (19)$$

where A_{ba} is the spontaneous emission rate of photons from level $|b\rangle$ to level $|a\rangle$. Since the energy radiated per second $I(t) = A_{ba} N_b(t) \hbar\omega_{ab}$, the emission intensity at time t , and frequency ω_{ba} is given by $I(t) = I(0) \exp(-t/\tau_R)$, where the radiative decaytime τ_R is defined as the reciprocal of the *spontaneous decay rate*, i.e., $\tau_R = (1/A_{ba})$. The expectation value of the time that the electron spends in the excited state, $\langle t \rangle$, is calculated from

$$\langle t \rangle = \frac{1}{N_b(0)} \int_{-s}^{+\infty} N_b(t) dt = (A_{ba})^{-1} = \tau_R \quad (20)$$

This is just the average time, or *lifetime*, of the electron in the excited state. In consequence, this simple argument identifies the radiative decaytime with the lifetime of the electron in the excited state. Typically, for an allowed electric dipole transition $\tau_R \sim 10^{-8}$ s.

The radiation from a collection of atoms emitting radiation at frequency ω_{ba} at time $t > 0$ has an associated electric field at some nearby point given by

$$E(t) = E_0 \exp(i\omega_{ba}t) \exp(-t/2\tau_R) \quad (21)$$

(i.e., the electric field oscillates at the central frequency of the transition on the atom). The distribution of frequencies in $E(t)$ is obtained by Fourier analyzing $E(t)$ into its frequency spectrum, from which

$$E(\omega) = \frac{E_0}{\sqrt{(\omega_{ab} - \omega)^2 + (2\tau_R)^{-2}}} \exp[i\phi(\omega)] \quad (22)$$

where $\phi(\omega)$ is a constant phase factor. Since $I(t) \approx E(t)^2$ we obtain the intensity distribution of frequencies given by

$$I(\omega) = \frac{I_0}{(\omega_{ab} - \omega)^2 + (2\tau_R)^{-2}} \quad (23)$$

This classical argument shows that the distribution of frequencies in the transition has a *lorentzian* shape with full width at half maximum (FWHM), $\Delta\omega$, given by

$$\Delta\omega = \frac{1}{\tau_R} = A_{ba} \quad (24)$$

An identical lineshape is derived from wave mechanics using time-dependent perturbation theory. This relationship between the *natural linewidth* of the transition, $\Delta\omega$, and the radiative decaytime, τ_R , is related to the uncertainty principle. The time available to measure the energy of the excited state is just $\langle t \rangle$; the width in energy of the transition is $\Delta E = \hbar\Delta\omega$. Hence $\Delta E\langle\tau\rangle = \hbar\Delta\omega\tau_R = \hbar$ follows from Eq. (24). For $\tau_R \approx 10^{-8}$ s the energy width $\Delta E/c \approx 5 \times 10^{-2} \text{ m}^{-1}$. Hence the natural linewidth of a transition in the visible spectrum is $\Delta\lambda \approx 2 \times 10^{-3} \text{ nm}$.

The broadening associated with the excited state lifetime is referred to as *natural* or *homogeneous* broadening. There are other processes which modulate the energy levels of the atom thereby contributing to the overall decay rate, τ^{-1} . It is this overall decay rate, $\tau^{-1} > \tau_R^{-1}$, which determines the width of the transition. Examples of such additional processes include lattice vibrations in crystals and the vibrations/rotations of molecules. In gas-phase spectroscopy, random motion of atoms or molecules leads to inhomogeneous broadening via the Doppler effect. This leads to a Gaussian spectral profile of FWHM given by

$$\Delta\omega_D = \frac{4\pi}{c} \left(\frac{2kT}{M} \ln 2 \right)^{1/2} \omega_{ba} \quad (25)$$

showing that the Doppler width varies as the square root of temperature and is smaller in heavier atoms. In solids, distortions of the crystal field by defects or growth faults lead to strain which is manifested as inhomogeneous broadening of spectra. The resulting lineshape is also gaussian. The great power of laser spectroscopy is that spectroscopists recover the true homogeneous width of a transition against the background of quite massive inhomogeneous broadening.

10.6 SPECTROSCOPY OF ONE-ELECTRON ATOMS

Figure 1a shows the energy level structure of atomic hydrogen for transitions between $n = 3$ and $n = 2$ states, i.e., the Balmer α -transition. Electric dipole transitions are indicated by vertical lines. The relative strengths of the various lines are indicated by the lengths of the vertical lines in Fig. 1b. Also shown in Fig. 1b is a conventional spectrum obtained using a discharge tube containing deuterium atoms cooled to $T = 50 \text{ K}$. This experimental arrangement reduces the Doppler width of the Balmer α -transition at 656 nm to 1.7 GHz. Nevertheless, only three transitions $2P_{3/2} \rightarrow 3D_{5/2}$, $2S_{1/2} \rightarrow 3P_{1/2}$, and $2P_{1/2} \rightarrow 3D_{3/2}$ are resolved.¹¹ However, sub-Doppler resolution is possible using laser saturation spectroscopy.¹² Figure 1c shows the Doppler-free Balmer α -spectrum to comprise

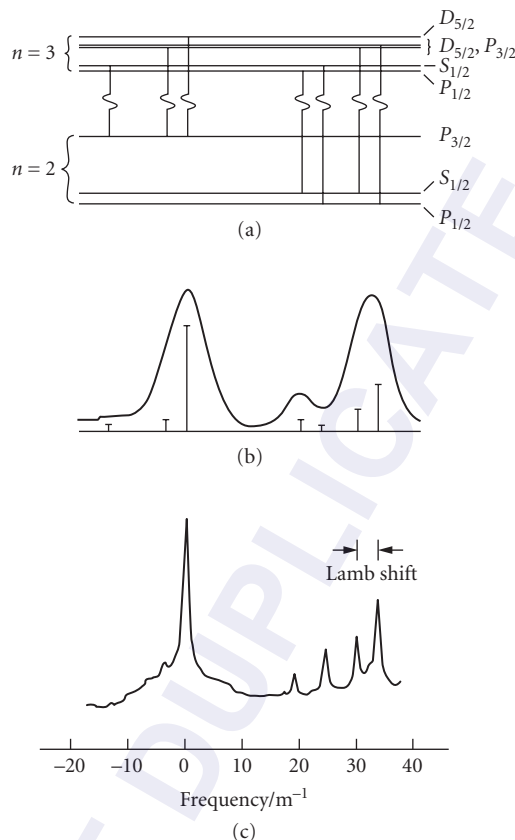


FIGURE 1 (a) The structure of the $n = 3$ and $n = 2$ levels of hydrogen showing the Balmer α -transitions. (b) A low resolution discharge spectrum of deuterium cooled to 50 K. (c) A Doppler-free spectrum of the Balmer α -spectrum of hydrogen. (After Hänsch et al.¹²)

comparatively strong lines due to $2P_{3/2} \rightarrow 3D_{5/2}$, $2S_{1/2} \rightarrow 3P_{1/2}$, $2S_{1/2} \rightarrow 3P_{3/2}$, and $2P_{1/2} \rightarrow 3D_{3/2}$ transitions, as well as a very weak $2P_{3/2} \rightarrow 3D_{3/2}$ transition. Also evident is a cross-over resonance between the two transitions involving a common lower level $2S_{1/2}$, which is midway between the $2S_{1/2} \rightarrow 3P_{1/2}$, $3P_{3/2}$ transitions. The splitting between $2P_{3/2}$, $2P_{1/2} \rightarrow 3D_{3/2}$ transitions measures the spin-orbit splitting in the $n = 2$ state, which from Eq. (2) is about 36.52 m^{-1} . The Lamb shift is measured from the splitting between $2S_{1/2} \rightarrow 3P_{3/2}$ and $2P_{1/2} \rightarrow 3D_{3/2}$ lines to be 3.53 m^{-1} , which compares well with the original microwave measurement (3.537 m^{-1}).¹³ Subsequently, Hänsch et al.¹⁴ made an interferometric comparison of a Balmer α -line with the 632.8-nm line from He-Ne locked to a component of $^{129}\text{I}_2$, thereby deriving a value of R_∞ of $10973731.43(10) \text{ m}^{-1}$, at that time an order of magnitude improvement in accuracy on previous values. Neither the $2S_{1/2}$ nor $2P_{1/2}$ hfs was resolved in this experiment, both splittings being less than the system resolution of $ca 0.05 \text{ m}^{-1}$. This probably followed from the use of pulsed dye lasers where the laser linewidth exceeds by factors of 10 the linewidth available from single-frequency continuous wave (CW) dye lasers. Subsequent measurements using CW dye lasers standardized against I_2 -stabilized He-Ne lasers gave further improvements in the value of R_∞ .^{15,16} Also using the Balmer α -transition Stacey et al.¹⁷ have studied the isotope shifts between spectra from

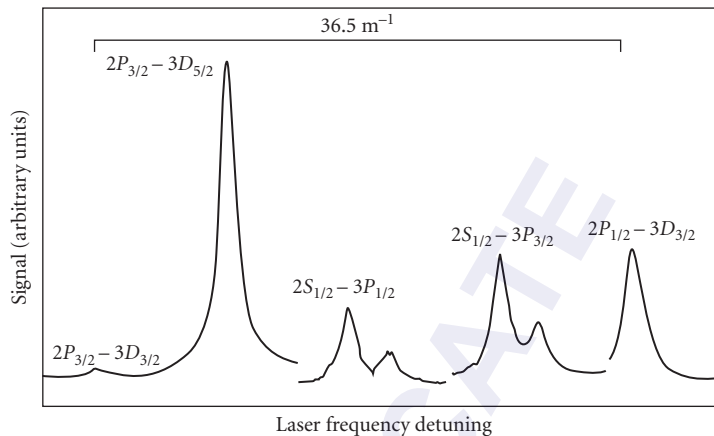


FIGURE 2 Showing hyperfine structure splittings of the $2S_{1/2} \rightarrow 3P_{3/2}, 3D_{3/2}$ transition in hydrogen. (After Stacey.¹⁷)

hydrogen, deuterium, and tritium. The spectrum shown in Fig. 2 reveals the $2S_{1/2}$ hyperfine splitting on the $2S_{1/2} \rightarrow P_{3/2}$ and $2S_{1/2} \rightarrow 3D_{3/2}$ transitions. These measurements yield isotope shifts of 124259.1(1.6) MHz and 41342.7(1.6) MHz for *H-D* and *D-T*, respectively, which accord well with theoretical values.

The literature on H-atom spectroscopy is vast, no doubt fueled by the unique relationship between experimental innovation and fundamental tests of relativistic quantum mechanics and QED. This chapter is not a comprehensive survey. However, it would be seriously remiss of the author to omit mention of three other categories of experimentation. The first experimental arrangement uses crossed atomic and laser beams: a well-collimated beam of atoms propagates perpendicular to a laser beam which, after traversing the atomic beam, is reflected to propagate through the atomic beam again in the opposite direction. The interaction between the counter-propagating beams and the atoms in the atomic beam is signaled by a change in the beam flux. The atomic beam replaces the discharge unit used in conventional atomic spectroscopy, thereby reducing errors due to the electric fields in discharges. This experiment is the optical analogue of the Lamb-Retherford radio-frequency experiment¹³ and has been much used by workers at Yale University.¹⁸ They reported a value of $R_{\infty} = 10973731.573(3) \text{ m}^{-1}$ in experiments on the Balmer β -transition ($n = 2$ to $n = 4$).

There have been several other studies of the Balmer β -transition, which has a narrower natural linewidth than the Balmer α -transition. However, because it is weaker than the Balmer α -transition, Wieman and Hänsch used a polarization scheme to enhance the sensitivity of the saturation absorption scheme.¹⁹ Finally, the metastable $2S_{1/2}$ level may decay spontaneously to the $1S_{1/2}$ ground state with the emission of two photons with energies that sum to the energy separation between $1S_{1/2}$ and $2S_{1/2}$. Such a process has a radiative decaytime of 0.14 s, giving a natural linewidth for the $2S_{1/2} \rightarrow 1S_{1/2}$ transition of order 1 Hz! The probability of a two-photon absorption transition is quite low. However, as with laser absorption saturation spectroscopy, two-photon absorption experiments are made feasible by Doppler-free resolution. Wieman and Hänsch²⁰ used an amplified CW laser beam at 243 nm to excite a two-photon absorption transition, which they detected by observing the Lyman α -emission at 121 nm. In addition, part of the laser beam was split off and used to measure simultaneously the Balmer β -spectrum. The coupled experiment permitted a direct measurement of the ground-state Lamb shift of 8161(29) MHz. Ferguson and his colleagues developed standard cells using $^{130}\text{Te}_2$ lines for studies of the Balmer β - and $1S_{1/2} \rightarrow 2S_{1/2}$ transitions.^{10,21}

10.7 MULTIELECTRON ATOMS

In order to calculate the energy level spectrum of a multielectron atom we require a suitable hamiltonian describing the interaction of all the electrons with the nucleus and with each other. A convenient starting point is the simplified Dirac equation for a one-electron atom, viz.,

$$H = H_0 + H_{so} \quad (26)$$

where H_0 is the simplified hamiltonian for the electron in the field of a nucleus of charge Ze at rest, i.e.,

$$H_0 = \frac{p^2}{2m} - \frac{Ze^2}{4\pi\epsilon_0 r} \quad (27)$$

and $H_{so} = -\zeta(r)l \cdot s$ is the spin-orbit hamiltonian. Wavefunctions which satisfy Eq. (27) are

$$\psi_{nlm_l}(r) = R_{nl}^{(r)} Y_l^m(\theta\phi) \quad (28)$$

where the labels n , l , and m are quantum numbers which characterize the eigenstates. The eigenvalues, given in Eq. (1), depend only on the principal quantum number n , which takes positive integral values. The quantum number l characterizing the orbital angular momentum also takes integral values, $l = 0, 1, 2, 3, \dots, (n-1)$, whereas m_l measures the z -component of the orbital angular momentum. There are $2l+1$ integral values of l given by $m = 1, (l-1), (l-2), \dots, -(l-1), -l$, and for a given value of n there are several different orbital states with identical energy. H_{so} , among other interactions, raises this degeneracy.

It is convenient to represent the orbital wavefunction ψ_{nlm} by the ket $|nlm\rangle$. Including spin angular momentum we represent a *spin orbital* by $|nlsmm_s\rangle$ or more simply by $|nlmm_s\rangle$. Recalling the brief discussion of the coupled representation, whereby $j = l + s$, an equally valid representation is $|nljm_j\rangle$. Indeed the new basis states $|nljm_j\rangle$ are just linear combinations of the $|nlmm_s\rangle$ basis states.⁸ Each wavefunction has a definite parity. The parity of wavefunctions is important in determining the selection rules of spectra. The inversion operator P_i , defined by $P_i f(r) = f(-r)$ for any function of r , gives the following result

$$P_i |nlm\rangle = (-1)^l |nlm\rangle \quad (29)$$

Hence, for even values of l the wavefunctions are said to have even parity since they do not change sign under inversion of coordinates. For l odd the wavefunctions have odd parity. The strength of an optical transition is determined by a matrix element $\langle\psi_b|\mu|\psi_a\rangle$, where the integration is taken over the volume of the atom. In the electric dipole approximation, $\mu = -er$ so that the matrix element is zero except that the wavefunction ψ_a and ψ_b have opposite parity. This defines the Laporte selection rule which states that the parity of a state must change from odd to even (or vice versa) in an electric dipole transition.

The hamiltonian for multielectron atoms is a sum over all N electrons of one-electron operators [see Eq. (1)] plus an electron-electron Coulomb repulsion between electrons i and j separated by a distance r_{ij} . Hence we may write this as

$$H = \sum_i \left(\frac{p_i^2}{2m} - \frac{Ze^2}{4\pi\epsilon_0 r_i} + \zeta(r_i) l_i \cdot s_i \right) + \sum_{i \neq j} \frac{e^2}{4\pi\epsilon_0 r_{ij}} \quad (30)$$

The computational complexity militates in favor of an approximate solution because the spin-orbit and electron-electron interactions are not spherically symmetric. In consequence, the first stage of the approximation to Eq. (30) is in the form

$$H = \sum_i \frac{p_i^2}{2m} + V_i'(r_i) + \zeta(r_i) l_i \cdot s_i \quad (31)$$

where $V'_i(r_i)$ is a spherically symmetric one-electron operator which represents the potential energy of the i th electron in the field of the nucleus and all the electrons. The first two terms in this sum constitute the orbital hamiltonian, H_0 , a sum of one-electron hydrogen-like hamiltonians [Eq. (27)], but with a more complicated radial potential energy function, $V'(r)$. The radial and angular parts of each one-electron hamiltonian are separable and we write orbital functions

$$R'_{nl}(r_i)Y_l^m(\theta, \phi) = |nlm\rangle \quad (32)$$

However, $R'_{nl}(r_i)$ is the solution of the radial equation involving the central potential, $V'(r_i)$, which is characterized by the quantum numbers n and l . In consequence, the energy of the one-electron state also depends on both n and l . The complete spin orbital is characterized by four quantum numbers including spin (i.e., $u = |nlmm_s\rangle$) and the many electron eigenstate of H_0 is a product of one-electron states

$$U = \prod_i |nlmm_s\rangle_i \quad (33)$$

The energy E_u of this product state is

$$E_u = \sum_i E_{n_i} \quad (34)$$

which depends on the set of $n_i l_i$ values. However, since E_u does not depend on m_i and m_s these eigenstates have a large degeneracy.

Experimentally, the complete wavefunctions of electrons are antisymmetric under the exchange of orbital and spin coordinates of any two electrons. The product wavefunction, Eq. (33), does not conform to the requirement of interchange symmetry. Slater solved this problem by organizing the spin orbitals into an antisymmetric N -electron wavefunction in determinantal form.²² The application of the Hartree-Fock variational approach to determine the central field potential consistent with the best Slater wavefunctions is described in detail by Tinkham.²³ There are many different sets of energy eigenfunctions that can be chosen; the net result is that the eigenstates of H_0 can be classified by a set of quantum numbers $LSM_L M_S$ for each $(n_i l_i)$ electron configuration, where $L = \sum l_i$ and $S = \sum s_i$. That the eigenfunctions must be antisymmetric restricts the number of possible L and S values for any given configuration. Since $J = L + S$ is also a solution of H_0 , we can represent the eigenstates of the configuration by the ket $|LSM_L M_S\rangle$ or alternatively by $|LSJM_J\rangle$ where the eigenstates of the latter are linear combinations of the former.

There is a particular significance to the requirement of antisymmetric wavefunctions in the Slater determinantal representation. A determinant in which any two rows or columns are identical has the value of zero. In the present context, if two one-electron states are identical, then two columns of the Slater determinant are identical, and the wavefunction is identically zero. This is a statement of the Pauli exclusion principle: no two electrons in an atom can occupy identical states (i.e., can have the same four quantum numbers). The Slater wavefunctions indicate those one-electron states which are occupied by electrons. To see how this works consider two equivalent electrons in p -states on an atom. For $n = n'$ for both electrons, the l -values may be combined vectorially to give $L = 2, 1$, and 0 . Similarly, the two electron spins, $s = \frac{1}{2}$, may be combined to give $S = 1$ or 0 . The antisymmetric requirement on the total wavefunction means that the symmetric orbitals $D(L = 2)$ and $S(L = 0)$ states can only be combined with the antisymmetric spin singlet, $S = 0$. The resulting spin orbitals are represented by 1D and 1S . In contrast the antisymmetric P state must be combined with the spin triplet, which is a symmetric function, yielding the antisymmetric spin orbital, 3P . However, for two inequivalent p -electrons in an excited state of the atom both spin singlet and spin triplet states are possible for the S, P , and D orbital states.

10.8 OPTICAL SPECTRA AND THE OUTER ELECTRONIC STRUCTURE

Optical spectroscopy probes those electronic transitions associated with a small number of electrons outside the closed shells of electrons. This gives further simplification to the computational problem since the multielectron hamiltonian

$$H = \sum_i \frac{p_i^2}{2m} + V'(r_i) + \zeta(r_i) l_i \cdot s_i + \sum_{i>j} \left(\frac{e^2}{4\pi\epsilon_0 r_{ij}} \right) \quad (35)$$

is summed only over the outer electrons where each of these electrons moves in the central field of the nucleus and the inner closed-shell electrons, $V'(r_i)$. Neglecting the smallest term due to the spin-orbit interaction, the hamiltonian in Eq. (35) takes the form $H_0 + H'$, where H_0 is a sum of one-electron hamiltonian with the radial potential functions $V'(r_i)$ and H' is the energy of the Coulomb interaction between the small number of outer electrons. The corrections to the one-electron energies are then the diagonal matrix elements $\langle n_i l_i m_i m_{s_i} | H' | n_i l_i m_i m_{s_i} \rangle$, expressed either in terms of Racah parameters, A, B, C , or Slater parameters, F_0, F_2, F_4, \dots .²⁴ Transition metal ion energy levels are normally described in terms of the Racah parameters and rare earth ion energy levels in terms of the Slater functions. The effect of H' is to split each configuration $(n_i l_i)$ into a number of LS terms for each of which there are $(2L + 1)(2S + 1)$ distinct energy eigenstates. We represent the energy eigenstates by the kets $|(n, l)LSJM_j\rangle$, which are defined as linear combinations of the $|(n_i l_i)LSM_L M_S\rangle$ states.

Returning briefly to the $(np)^2$ configuration (characteristic of the Group-4 elements, C, Si, Ge, etc. of the Periodic Table) it is noted that the diagonal matrix elements evaluated in terms of the Slater parameters are given by $E(^1D) = F_0 + F_2$, $E(^1S) = F_0 + 10F_2$, and $E(^3P) = F_0 - 5F_2$. For different atoms it is the relative values of F_0 and F_2 which change through the series $(2p)^2$, $(3p)^2$, $(4p)^2$, etc. Note that it is the term with maximum multiplicity, 3P , which is lowest in energy in conformity with Hund's rule. A similar situation arises for the $(np)^4$ configuration of, for example, atomic O, S, and Se which might equally and validly be considered as deriving from two holes in the $(np)^2$ configuration. The general conclusion from this type of analysis is that the energy level structures of atoms in the same period of the Periodic Table are identical, with the principal differences being the precise energies of the eigenstates. This is evident in Fig. 3, the term scheme for atomic Li, which has the outer electron configuration $(2S)^1$; this may be looked on as a pseudo-one-electron atom. There is a general spectroscopic similarity with atomic hydrogen, although the $(ns)^1 - (np)^1$ splittings are much larger than in hydrogen. The $3S_{1/2} \leftrightarrow 2S_{1/2}$ transition is forbidden in Li just as the $2S_{1/2} \leftrightarrow 1S_{1/2}$ transition is in hydrogen. However, it is unlikely to be observed as a two-photon process in emission because the $3S_{1/2} \rightarrow 2P_{3/2, 1/2}$ and $2P_{3/2, 1/2} \rightarrow 2S_{1/2}$ transitions provide a much more effective pathway to the ground state.

The comparison between Li and Na is much more complete as Figs. 3 and 4 show: assuming as a common zero energy the ground state, the binding energies are $E(2S_{1/2}, \text{Li}) = 43,300 \text{ cm}^{-1}$ and $E(3S_{1/2}, \text{Na}) = 41,900 \text{ cm}^{-1}$. The energies of the corresponding higher lying nS , nP , and nD levels on Li are quite similar to those of the $(n + 1)S$, $(n + 1)P$, and $(n + 1)D$ levels on Na. This similarity in the energy level structure is reflected in the general pattern of spectral lines, although the observed wavelengths are a little different. For example, the familiar D lines in the emission spectrum of Na occur at $\lambda \approx 589.9 \text{ nm}$ whereas the corresponding transition in Li occur at 670.8 nm .

Examination of the term schemes for a large number of elements reveal striking similarities between elements in the same group of the Periodic Table, due to this structure being determined by the number of electrons outside the closed shell structure. These term diagrams also reveal a very large number of energy levels giving rise to line spectra in the ultraviolet, visible, and infrared regions of the spectrum. The term diagrams are not very accurate. Tables of accurate energy levels determined from line spectra have been compiled by Moore²⁵ for most neutral atoms and for a number of ionization states of the elements. This comprehensive tabulation reports energy levels to a very large principal quantum number ($n \cong 10, 11$).

The term diagram of the neutral Tl atom, $(6p)^1$ configuration (see Fig. 5) shows two interesting features in comparison with those for alkali metals (see Figs. 3 and 4). In the $(6p)^1$ state the spin-orbit splitting into $6P_{1/2}$ and $6P_{3/2}$ levels amounts to almost 8000 cm^{-1} whereas the spin-orbit splitting

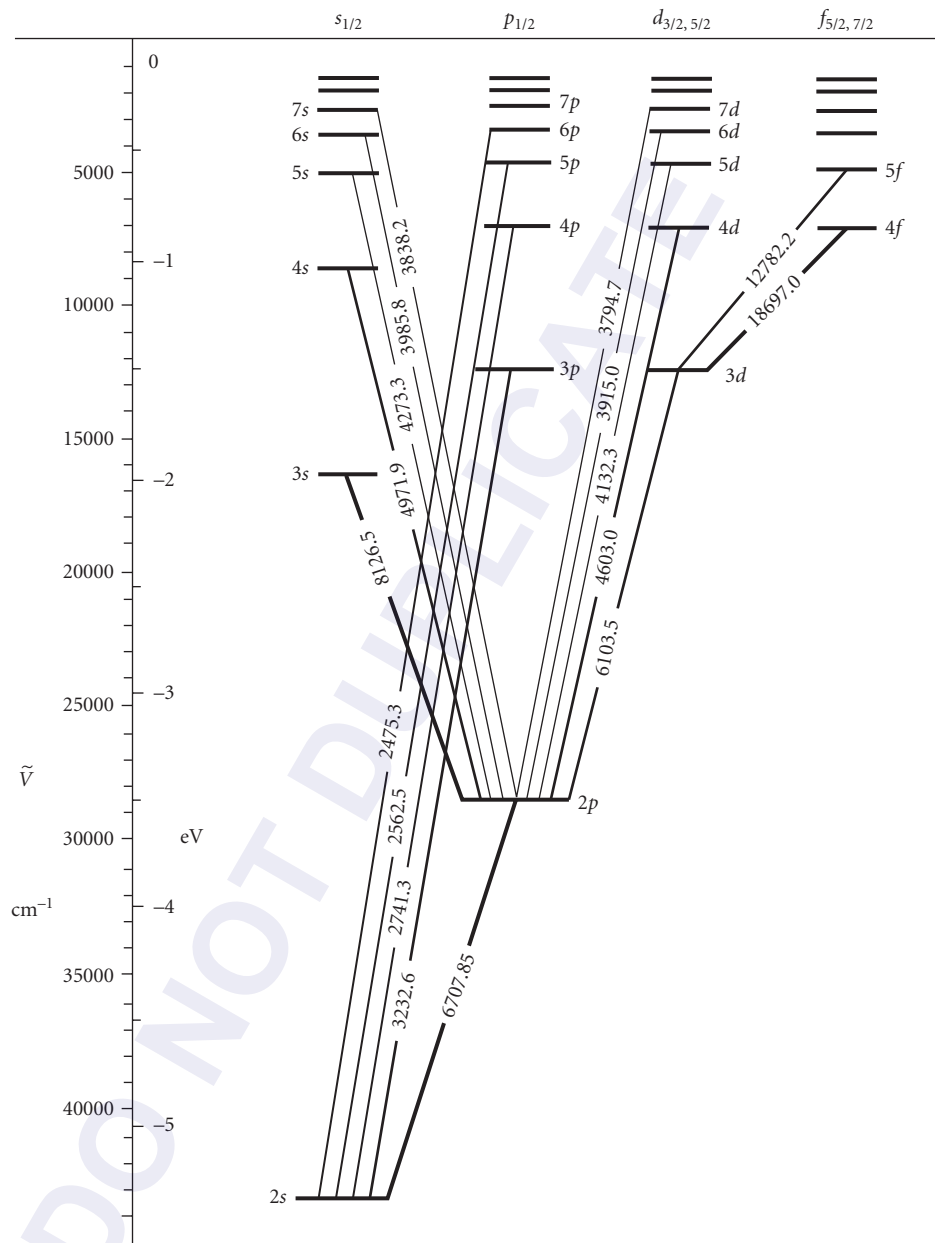


FIGURE 3 The term diagram of atomic Li, in which the slanted lines indicate the observed electric dipole transitions and the numbers on the lines are the wavelengths in Ångström units. (After Grotian.²⁶)

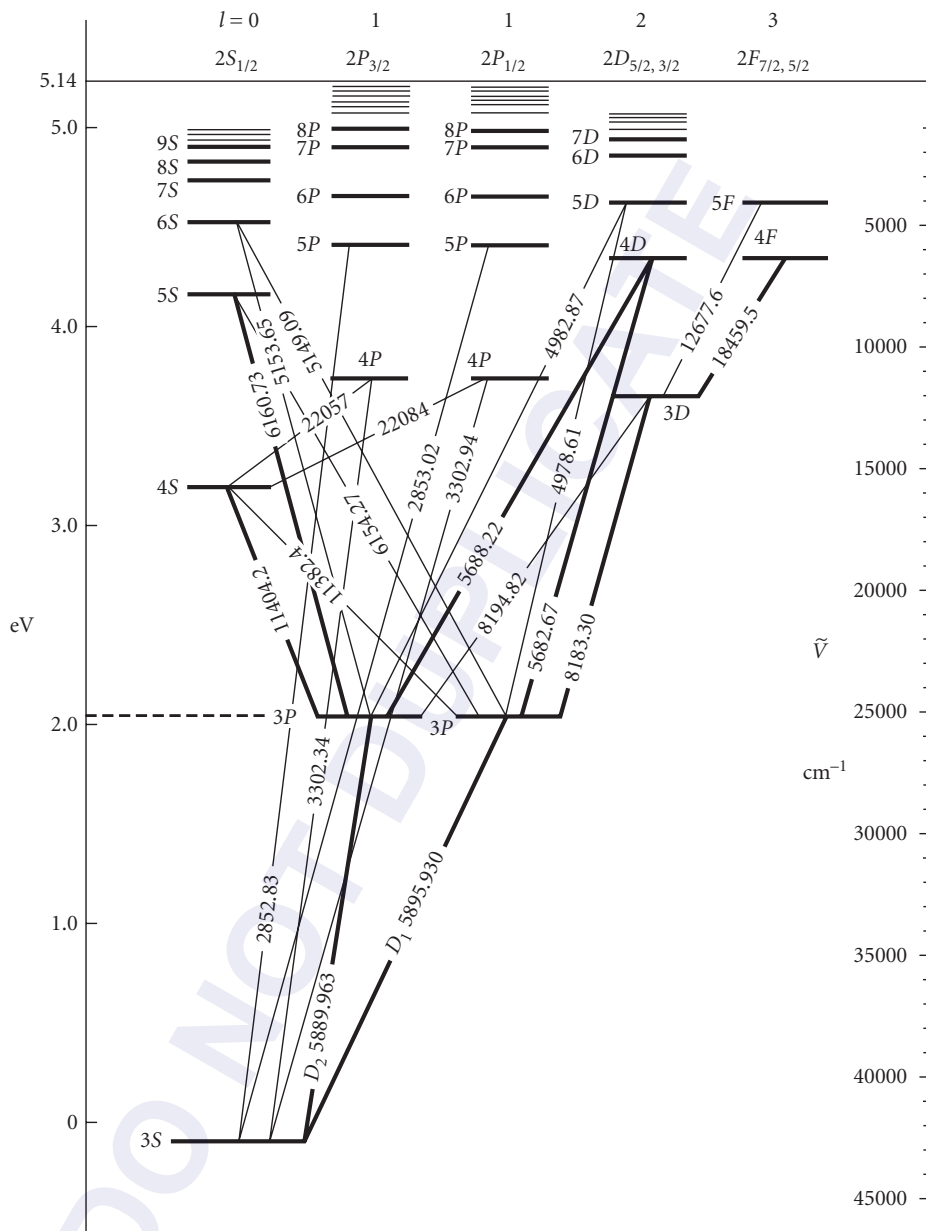


FIGURE 4 The term diagram of atomic Na. (After Grotian.²⁶)

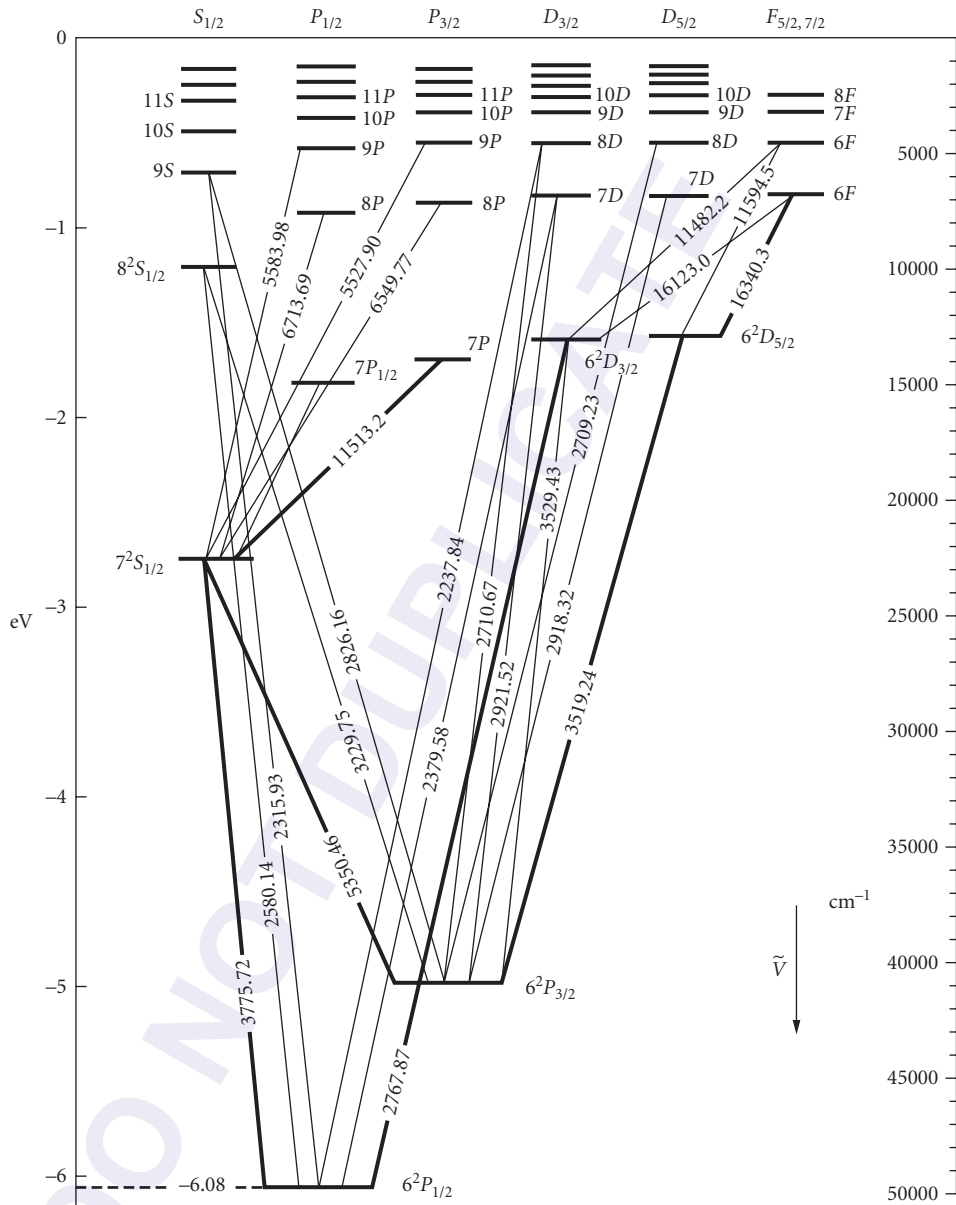


FIGURE 5 The term diagram of neutral Tl. (After Grotian.²⁶)

between the $3P_{1/2}$ and $3P_{3/2}$ levels of Na is only 17 cm^{-1} . This reflects the (Z^4/n^3) dependence of the spin-orbit coupling constant. Furthermore, when Tl is in the $7S_{1/2}$ state it can decay radiatively via transitions to either $6P_{3/2}$ or $6P_{1/2}$, each with a distinct transition probability. The relative probability of these transitions is known as the branching ratio for this mode of decay. Other examples of branching are apparent in Fig. 5: the branching ratios are intrinsic properties of the excited state.

10.9 SPECTRA OF TRI-POSITIVE RARE EARTH IONS

The rare earth elements follow lanthanum ($Z = 57$) in the Periodic Table from cerium ($Z = 58$), which has the outer electron configuration $4f^1 5d^1 6s^2$ to ytterbium ($Z = 70$) with electron configurations $4f^{13} 5d^1 6s^2$. In the triply charged state in ionic crystals all $5d$ and $6s$ electrons are used in ionic bonding and the many energy levels of these rare earth (RE) $^{3+}$ ions are due to the partially filled $4f$ shell. The number of electrons in the $4f$ shell for each trivalent ion and the ground-state configuration is indicated in Table 1. The energy levels of the unfilled $4f^n$ shells spread out over some $40,000 \text{ cm}^{-1}$, giving rise to numerous radiative transitions with energies in the visible region. A remarkable feature of the $4f^n$ electrons is that they are shielded by the outer $5s$ and $5d$ shells of electrons, with the result that $4f$ electrons are not strongly affected by interactions with neighboring ions in crystals. In consequence, the energy levels of the $4f$ electrons in crystals are essentially the free ion levels characterized by quantum numbers L , S , and J . As with the free ions the R.E. $^{3+}$ ions in crystals have very sharp energy levels which give rise to very sharp line spectra. The crystal field interaction does split the RE $^{3+}$ ion levels, but this splitting is very much smaller than the splittings between the free-ion levels. Hence for rare earth ions in different crystals the gross features of the optical spectra are unchanged.

As discussed in Sec. 10.6 the eigenstates of the $4f$ electrons are calculated using the central field approximation from which each $4f$ electron state is characterized by the ket $|n=4, l=3m_l m_s\rangle$. The effect of the Coulomb repulsion between electrons, $H' = \sum_i e^2 / 4\pi\epsilon_0 r_{ij}$, is to split the energy levels of the $4f^n$ configuration into different LS terms, with wavefunctions characterized by kets $|LSM_L M_S\rangle$. The magnitudes of the electrostatic interaction for each LS level are expressed as sums of Slater electron-electron integrals F_k , with $k = 0, 2, 4$, and 6 for $4f$ electrons. Since F_0 contributes equally to all LS states of the same $4f^n$ configuration, this term can be neglected. Generally, these Slater integrals are regarded as adjustable parameters with magnitudes determined by fitting to the measured line spectra. The values of the F_k integrals for $4f^n$ ions in many crystals vary by only about 2 percent from those obtained for free ions; they also vary slightly depending on the nature of the surrounding ions. The next largest term in the hamiltonian, after H' , is spin-orbit coupling. If the spin-orbit coupling

TABLE 1 The Number of Electrons (n) and Ground State of Tri-positive Rare Earth Ions

Ion	n (in ed^n)	Ground state
Ce $^{3+}$	1	$2F_{5/2}$
Pr $^{3+}$	2	$3H_4$
Nd $^{3+}$	3	$4I_{9/2}$
Pm $^{3+}$	4	$5I_4$
Sm $^{3+}$	5	$6H_{5/2}$
Eu $^{3+}$	6	$7F_0$
Gd $^{3+}$	7	$8S_0$
Tb $^{3+}$	8	$7F_6$
Dy $^{3+}$	9	$6H_{15/2}$
Ho $^{3+}$	10	$5I_8$
Er $^{3+}$	11	$4I_{15/2}$
Tm $^{3+}$	12	$3H_6$
Yb $^{3+}$	13	$2F_{7/2}$

energy is much smaller than the energy separation in the LS term then the spin-orbit interaction can be written as ζLS and the wavefunctions are characterized by $|LSJM_J\rangle$. The additional energy of the J -multiples are given by the Landé interval formula

$$E_J = \frac{\zeta}{2}[J(J+1) - L(L+1) - S(S+1)] \quad (36)$$

from which it is evident that the separation between adjacent levels is given by $J\zeta$ where J refers to the upper J value. However, deviations from this Landé interval rule do occur because of mixing of different LS -terms when spin-orbit coupling and electron-electron interaction are of similar magnitudes. Clear examples of this are obtained from the spectra of $\text{Pr}^{3+}(4f^2)$ and $\text{Tm}^{3+}(4f^{12})$.

As representative of the type of spectra observed from R.E.³⁺ ions in ionic crystals we consider just one example: $\text{Nd}^{3+}(4f^3)$ in $\text{Y}_3\text{Al}_5\text{O}_{12}$ (YAG). The Nd^{3+} -YAG material is important as the gain medium. Nd^{3+} has a multitude of levels many of which give rise to sharp line emission spectra. A partial energy level structure is shown in Fig. 6. The low-temperature emission is from the ${}^4F_{3/2}$ level

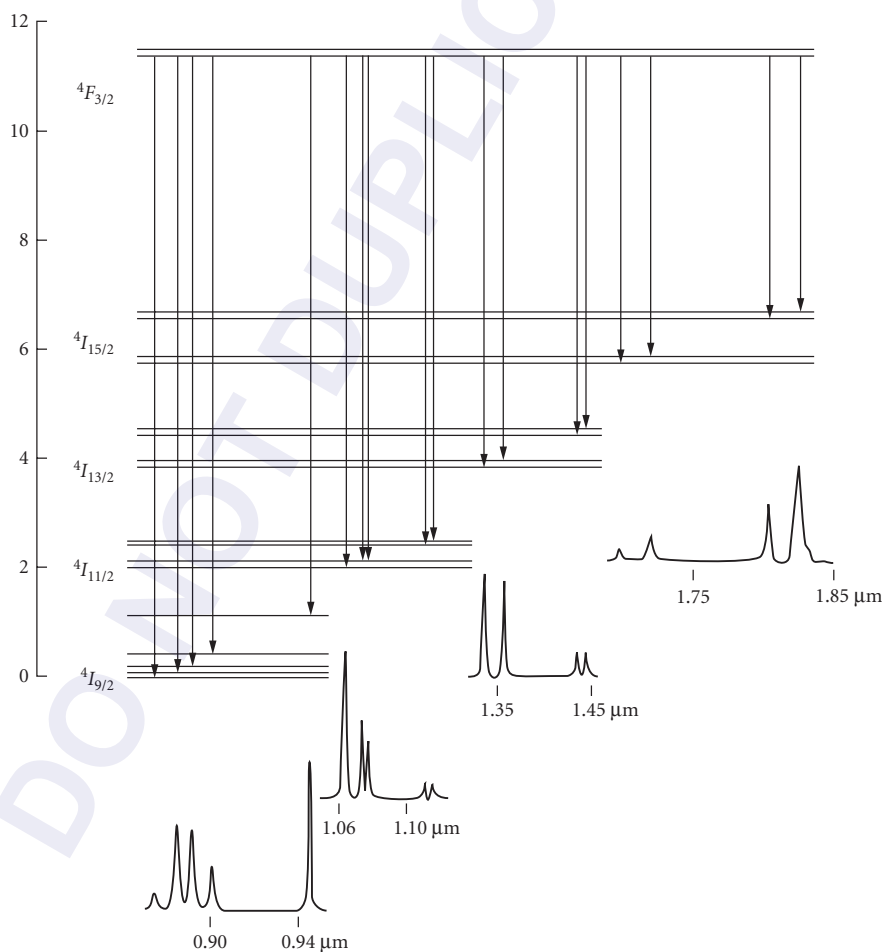


FIGURE 6 The low-temperature photoluminescence spectrum from the ${}^4F_{3/2}$ level of Nd^{3+} in $\text{Y}_3\text{Al}_5\text{O}_{12}$ and the corresponding energy level structure. (After Henderson and Imbusch.⁸)

to all the 4I_j levels of the Nd^{3+} ion. The spectra in Fig. 6 also show the splittings of the 4I_j levels by the crystal field, corresponding to the energy-level splitting patterns given in the upper portion of Fig. 6. Depending upon crystal quality these lines can be quite narrow with half-widths of order a few gigahertz. Nevertheless the low-temperature width in the crystal is determined by the distribution of internal strains and hence the lines are in homogeneously broadened. The natural linewidth of rare earth ion spectra is of the order of a few megahertz. By using optical hole burning (OHB), which is similar to the saturated absorption spectroscopy discussed earlier for atomic hydrogen, it is possible to eliminate the inhomogeneous broadening and recover the homogeneous lineshape of the spectrum. In principle, the natural width is determined by lifetime processes of which there are more numerous sources in crystals than in atomic vapors. Indeed the width may be determined by random modulation of the optical lineshape by photons and by both nuclear and electronic spins of neighboring ions. The two most general techniques for determining the homogeneous widths of optical transitions in solids are optical holeburning and fluorescence line narrow (FLN). Examples of these techniques are discussed in Chap. 2 of Vol. V.

10.10 VIBRATIONAL AND ROTATIONAL SPECTRA OF MOLECULES

A consultation of any one of the tables of data in Moore's compilation²⁵ shows that the energy level schemes of most atoms are complex. This is confirmed by the associated atomic spectra. Considerable interpretive simplification is afforded by the construction of term diagrams (e.g., Figs. 3 to 5), on which a very large number of lines may be associated with a much smaller number of terms, each term corresponding to an energy level of the atom. The observed spectral lines are due to transitions between pairs of terms (not all pairs) which occur subject to an appropriate selection rule. The spectra of even simple molecules measured with low-dispersion show characteristic band spectra which are even more complicated than the most complex atomic spectra. These band spectra, when studied at higher spectral resolution, are observed to consist of an enormous number of closely spaced lines. At first acquaintance, such band spectra appear to be so complex as to defy interpretation. Order can be brought to the riot of spectral components by constructing term schemes for molecules involving electronic, vibrational, and rotational energy terms, which enable the molecular spectroscopist to account for each and every line.

Molecular physics was a research topic to which quantum mechanics was applied from the very earliest times. Heitler and London developed the valence band theory of covalency in the H_2 -molecule in 1927.²⁷ The theory shows that with both electrons in 1s states there are two solutions to the hamiltonian

$$E_{\pm} = 2E(1\text{H}) + \frac{K \pm \mathfrak{S}}{1 \pm \mathcal{S}} \quad (37)$$

where $E(1\text{H})$ is the energy of an electron in the ground state of atomic hydrogen, K is the Coulomb interaction due to the mutual actions of charges distributed over each atom, \mathfrak{S} is the exchange energy, and \mathcal{S} is the overlap integral. The exchange energy is a purely quantum mechanical term, representing the frequency with which the deformation of the wavefunctions by their mutual interaction oscillates from one atom to another. The positive sign refers to the symmetric combination of orbital wavefunctions for the two hydrogen atoms. Since the overall wavefunction must be antisymmetric, the combined spin states must be antisymmetric (i.e., the spin singlet $S = 0$; this state is labeled $^1\Sigma_g$). The evaluation of the integrals in Eq. (37) as a function of internuclear separation leads to Fig. 7 for the spin singlet state. This $^1\Sigma_g$ ground state has a potential energy minimum of about 4.8 eV (experimentally) with the nuclei separated by 0.75 nm, relative to the total energy of the two hydrogen atoms at infinity. The theoretical value of the binding energy on the valence band model is only 3.5 eV.²⁷ The negative sign in Eq. (37) guarantees that in the state characterized by the antisymmetric combination of orbital states and $S = 1$, i.e., $^3\Sigma_u$, the energy is monotonically ascending,

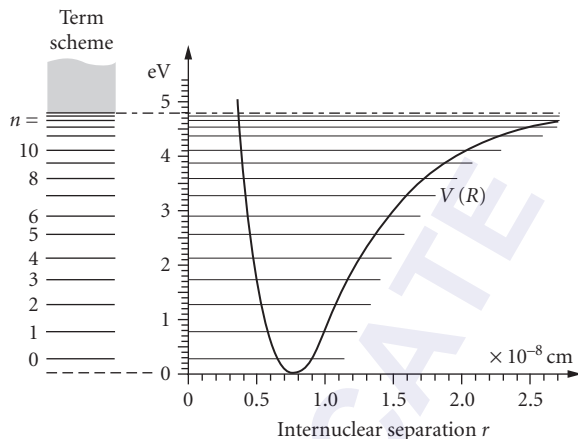


FIGURE 7 The internuclear potential, $V(R - R_0)$, in the ground state $^1\Sigma_g$ of the hydrogen molecule.

corresponding to repulsion between the two hydrogen atoms for all values of R . In such a state the molecule dissociates.

The energy versus internuclear separation curve in Fig. 7 represents the effective potential well, in which the protons oscillate about their mean positions. This is the potential used in the Born-Oppenheimer treatment of the molecular vibrations in the spectra of diatomic molecules.²⁸ The potential function may be written as a Taylor expansion

$$V(R - R_0) = V_0 + (R - R_0) \left(\frac{dV}{d(R - R_0)} \right)_0 + \frac{(R - R_0)^2}{2} \left(\frac{d^2V}{d(R - R_0)^2} \right)_0 + \frac{(R - R_0)^3}{6} \left(\frac{d^3V}{d(R - R_0)^3} \right)_0$$

where the subscript 0 refers to values of the differentials at $R = R_0$ and higher terms have been ignored. For a molecule in stable equilibrium the potential energy is a minimum and the force at $R = R_0$ must be zero. In consequence, the potential energy function may be written as

$$V(R - R_0) = \frac{1}{2} (R - R_0)^2 \left(\frac{d^2V}{d(R - R_0)^2} \right)_0 + \frac{1}{6} (R - R_0)^3 \left(\frac{d^3V}{d(R - R_0)^3} \right)_0 \quad (38)$$

after setting $V_0 = 0$.

The excited states of H_2 are constructed on the assumption that only one electron is excited in the transition: the appropriate configurations may be written as $(1s, 2s)$, $(1s, 2p)$, $(1s, 3s)$, etc. The electronic states of the molecule are then suitable combinations of the atomic states and categorized according to the total orbital angular momentum determined by vector addition of individual electronic orbital momenta. These orbital angular momentum states are designated as Σ , Π , Λ , etc. when the total angular momentum quantum number is 0, 1, 2, etc. Thus the electronic state in Fig. 7 is described as the $^1\Sigma_g$ (ground state) and the lowest lying excited state as $^3\Sigma_u$, where the subscripts g and u indicate even (gerade) and odd (ungerade) parity of the orbital states, respectively. The first excitation states associated with the $(1s, 2s)$ molecular configuration, designated as $^1\Sigma_g$ and $^3\Sigma_u$, have energies

$$E_{\pm} = E(1H) + E(2H) + \frac{K \pm \mathfrak{S}}{1 \pm \mathcal{G}} \quad (39)$$

Because the charge distributions are different in the atomic $1s$ and $2s$ states, the values of K and \mathfrak{S} are also different in the $(1s, 2s)$ configuration relative to the ground configuration. Obviously, as the orbital degeneracy of the molecular configuration increases so does the number of molecular levels resulting from the configuration. For example, there are six molecular levels associated with $(1s, np)$ configuration, ten associated with $(1s, nd)$, and so on. Roughly speaking, half of these levels will have a potential minimum in a graph of electronic potential energy versus internuclear separation (e.g., Fig. 7), and are associated with “bonding” orbitals of the molecule. In the other levels the force between the constituent atoms is repulsive for all values of the internuclear separation. In such orbitals the molecules will tend to dissociate.

The potential energy curve for a molecular bonding orbital (e.g., Fig. 7) is of the form given by Eq. (38). In the lowest order approximation this relation defines a potential well which is harmonic in the displacements from the equilibrium value. Of course, not all values of the potential energy, $V(R - R_0)$, are permitted: the vibrational energy of the molecule is quantized, the quantum number n having integral values from $n = 0$ to infinity. The energy levels are equally spaced with separations $h\nu_v$. The second term is an anharmonic term which distorts the parabolic shape relative to the harmonic oscillator. There are two important differences between the quantized harmonic and anharmonic oscillators. First, where there is an infinite number of levels in the former, there is only a finite number of vibrational states for an anharmonic molecule. Second, the levels are not equally spaced in the anharmonic oscillator, except close to the potential minimum. For this reason the anharmonic oscillator will behave like a harmonic oscillator for relatively small values of the vibrational quantum number. It is normal to assume harmonic vibrations in the Born-Oppenheimer approximation.²⁸

Molecular spectra fall into three categories, according as they occur in the far infrared (20 to 100 μm), near infrared (800 to 2000 nm), or visible/near ultraviolet (750 to 150 nm) region. Spectra excited by radiation in the infrared region are associated with changes in the rotational energy of the molecule. Spectra in the near-infrared region correspond to simultaneous changes in the rotational and vibrational energy of the molecule. Finally, the visible and ultraviolet spectra signal simultaneous changes in the rotational, vibrational, and electronic energies of the molecule. The latter category has the greatest potential complexity since in interpreting such visible/ultraviolet spectra we must, in principle, solve the molecular hamiltonian containing terms which represent electronic, vibrational, and rotational energies.

If we simplify the molecular problem somewhat we may represent the stationary states of a molecule by a linear sum of three energy terms: an electronic term determined by a quantum number n , a vibrational term determined by a quantum number v , and a rotational term determined by a quantum number r . The band spectra emitted by molecules are then characterized by three different kinds of excitations: electronic, vibrational, and rotational with frequencies in the ratio $\nu_e : \nu_v : \nu_r = 1 : \sqrt{m/M} : m/M$, where m is the electronic mass and M is the reduced molecular mass. In an optical transition the electronic configuration (nl) changes, and generally we expect the vibrational/rotational state to change also. In consequence, for a particular electronic transition, we expect splittings into a number of vibrational components, each of which is split into a number of rotational components. Using the rough ratio rule given previously, and assuming the Rydberg to be a characteristic electronic transition, we find spectral splittings of 500 cm^{-1} and 60 cm^{-1} to characterize the vibrational and rotational components, respectively. Since the quantized energy levels of the harmonic oscillator and rigid rotor are given by

$$E_v = \left(n + \frac{1}{2} \right) h\nu_v \quad (40)$$

$$E_r = r(r+1)h\nu_r \quad (41)$$

where $\nu_v = (1/2\pi)\sqrt{k/M}$ and $\nu_r = \hbar/4\pi I$, in which k is the “spring” constant and $I = Ml^2$ is the moment of inertia for a dumbbell-shaped molecule of length l , and applying to the usual selection rules for electronic ($\Delta l = \pm 1$, with $\Delta j = 0, \pm 1$), vibrational ($\Delta v = \pm 1, \pm 2, \dots$), and rotational ($\Delta r = 0, \pm 1$) transitions, we expect transitions at the following frequencies

$$\nu = \nu_0 + (n'' - n')/\nu_v + [r''(r''+1) - r'(r'+1)]\nu_r \quad (42)$$

where the superscript primes and double primes refer to final and initial states. For simplicity we have assumed that the vibrational spring constants and rotational moments of inertia are unchanged in the transitions. The number of potential transition frequencies is obviously very large and the spectrum consists of very many closely spaced lines. Of course, the vibrational/rotational structure may be studied directly in microwave spectroscopy.²⁹ An early account of the interpretation of such spectra was given by Ruark and Urey.³⁰ For many years such complex spectra were recorded photographically; as we discuss in later sections, Fourier transform spectroscopy records these spectra electronically and in their finest detail. Subsequent detailed accounts of the energy levels of molecular species and their understanding via spectroscopic measurements have been given by Slater, Hertzberg, and others.^{31–33}

An example of the complexities of band spectra, for a simple diatomic molecule, is given in Fig. 8, in this case for the photographically recorded β -bands ($2\Pi \rightarrow 3\Sigma$) of the nitric oxide molecule. Under low dispersion band spectra are observed (see Fig. 8a) which break down into line spectra under very high resolution (see Fig. 8b). This band system is emitted in transitions having common initial and final electronic states, with all the electronic states corresponding to a multiplet in atomic spectra. The electronic energy difference determines the general spectral range in which the bands are observed. However, the positions of the individual bands are determined by the changes in the vibrational quantum numbers. The spectra in Fig. 8 are emission spectra: the bands identified in Fig. 8a by (n', n'') are so-called *progressions* in which the transition starts on particular vibrational levels ($n' = 0$) of the upper electronic state and ends on a series of different vibrational levels (n'') of the lower electronic level. Such n' progressions measure the vibrational energy level differences in the lower electronic levels. Specifically identified in Fig. 8a are band progressions $n' = 0 \rightarrow n'' = 4, 5, 6, 7 \dots$ and $n' = 1 \rightarrow n'' = 10, 11, 12 \dots$. Also evident, but not identified in Fig. 8a, are *sequences* of lines in which the difference in vibrational quantum number $n' - n''$ is a constant. For example, either side of the (0, 4) band can be seen $\Delta n = -4$ and -3 sequences, respectively.

The (0, 7) band is shown under much higher spectral resolution in Fig. 8b, revealing a plethora of closely spaced lines associated with the rotational structure. For transitions in which the vibrational quantum number changes by some multiple of unity, i.e., $\Delta n = 1, 2, 3, \dots$, the rotational quantum number changes by $\Delta r = 0$ and ± 1 . The three branches of the rotational spectrum have the following positions (measured in cm^{-1}):

$$R = A + 2B(r+1) \quad \text{for } \Delta r = -1 \text{ and } r = 0, 1, 2 \dots \quad (43a)$$

$$Q = A + Cr \quad \text{for } \Delta r = 0 \text{ and } r = 0, 1, 2 \dots \quad (43b)$$

$$P = A - 2Br \quad \text{for } \Delta r = +1 \text{ and } r = 1, 2, 3 \dots \quad (43c)$$

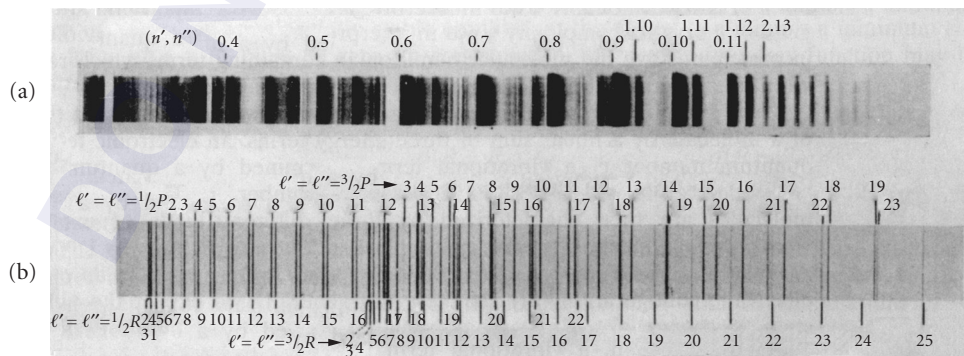


FIGURE 8 Characteristic band spectra of the diatomic NO molecule under conditions of (a) low resolution and (b) high resolution. (After Ruark and Urey.³⁰)

where in each case r refers to the rotational quantum number in the final electronic-vibrational state, $A = \bar{\nu}_e + \bar{\nu}_v$, and $B = \hbar/4\pi cl$. In consequence, if $C = 0$, the Q -branch consists of a single line at $Q(r) = \nu_e + \nu_v$, and the P - and R -branches shift linearly to lower and higher wavenumbers, respectively, relative to the Q -branch for increasing values of r . The spectra in Fig. 8 show that for the NO molecule $C \neq 0$ so that the lines in the P - and R -branches are shifted by $\pm 2B$ relative to the Q -branch lines. The particular Q -branches are defined by $l' = l'' = r + \frac{1}{2}$ with r being either 0 or 1, i.e., $l' = l'' = \frac{1}{2}$ or $\frac{3}{2}$. The $C \neq 0$ implies that the Eqs. (41) and (43) are inaccurate and additional rotational energy terms must be included even for vibrational states close to the bottom of the potential well (i.e., small n -values). Roughly the rotational term in Eq. (43) must be multiplied by $1 - (2\hbar^2(r+1)^2/4\pi^2 I^2)$, which has the further consequence that the separations of lines in the P - and R -branches corresponding to particular values of r increase with increasing r . A detailed analysis of Fig. 8b shows this to be the case.

10.11 LINESHAPES IN SOLID STATE SPECTROSCOPY

There are many modes of vibration of a crystal to which the optical center is sensitive. We will concentrate on one only, the breathing mode in which the ionic environment pulsates about the center. The variable for the lattice state, the so-called configurational coordinate, is labeled as Q . For the single mode of vibration the system oscillates about its equilibrium value Q_0^a in the ground state and Q_0^b in the excited state. The ground and excited state configurational coordinate curves are assumed to have identical harmonic shapes, and hence vibrational frequencies for the two states. This is illustrated in Fig. 9. In the Born-Oppenheimer approximation the optical center-plus-lattice system is represented in the electronic ground state by the product function²⁸

$$\Phi_a(r_i, Q) = \Psi_a(r_i, Q_0^a) \chi_a(Q) \quad (44)$$

and in the electronic excited state by

$$\Phi_b(r_i, Q) = \Psi_b(r_i, Q_0^b) \chi_b(Q) \quad (45)$$

The first term in the product is the electronic wavefunction, which varies with the electronic positional coordinate r_i , and hence is an eigenstate of $H_0 = \sum_i (p_i^2/2m) + V'(r_i)$, and is determined at the equilibrium separation Q_0^a . The second term in the product wavefunction is $\chi_a(Q)$, which is a function of the configurational coordinate Q . The entire ionic potential energy in state a is then given by

$$E^a(Q) = E_0^a + V_a(Q) \quad (46)$$

The $E^a(Q)$ in Fig. 9 is a harmonic potential function. This is an approximation to the potential used in Eq. (38) and illustrated in Fig. 7. A similar expression obtains for $E^b(Q)$. This representation of the Born-Oppenheimer approximation is often referred to as the *configurational coordinate model*.³⁹

Treating the electronic energy of the ground state E_0^a as the zero of energy, we write the ionic potential energy in the ground state as

$$E^a(Q) = \frac{1}{2} M \omega^2 (Q - Q_0^a)^2 \quad (47)$$

and in the excited state as

$$E_0^b(Q) = E_{ab} + \frac{1}{2} M \omega^2 (Q - Q_0^a)^2 - (2S)^{1/2} \hbar \omega \left(\frac{M \omega}{\hbar} \right)^{1/2} (Q - Q_0^a) \quad (48)$$

where E_{ab} is essentially the peak energy in an absorption transition between states $|a\rangle$ and $|b\rangle$ with the lattice at the coordinate, Q_0^a , where the Huang-Rhys parameter S is defined as

$$S = \frac{E_{dis}}{\hbar \omega} = \frac{1}{2} \frac{M \omega^2}{\hbar \omega} (Q_0^b - Q_0^a)^2 \quad (49)$$

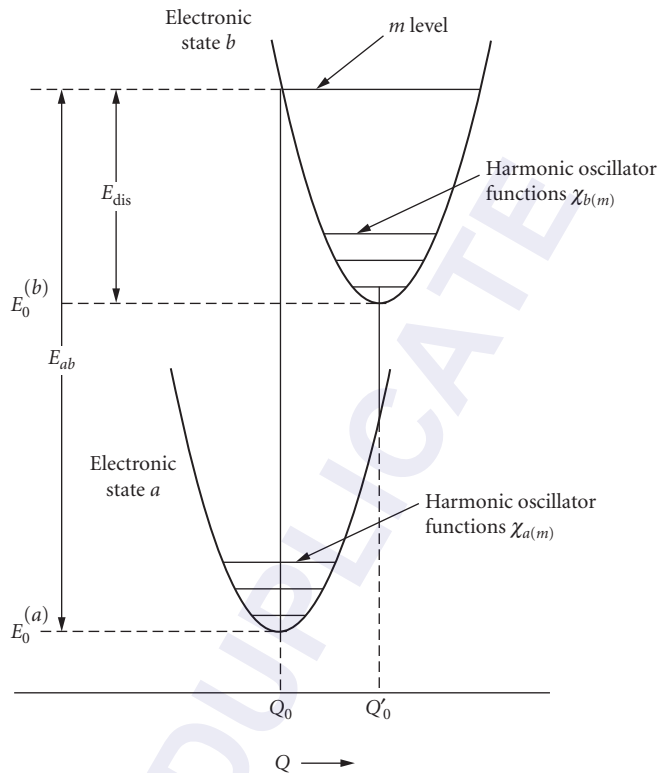


FIGURE 9 A configurational coordinate diagram showing the interionic potential in the harmonic approximation for electronic states $|a\rangle$ and $|b\rangle$. (After Henderson and Imbusch.⁸)

If the vertical line from $Q=Q_0^a$ in Fig. 9 intersects the upper configurational coordinate curve at the vibrational level n' then

$$E_{\text{dis}} = S\hbar\omega = \left(n' + \frac{1}{2}\right)\hbar\omega \quad (50)$$

The shapes of absorption and emission spectra are found to depend strongly on the difference in electron-lattice coupling between the states, essentially characterized by $(Q_0^b - Q_0^a)$ and by E_{dis} .

The radiative transition rate between the states $|a, n'\rangle \rightarrow |b, n''\rangle$, where n' and n'' are vibrational quantum numbers, is given by⁸

$$W(a, n' \rightarrow b, n'') = |\langle \psi_b(r_i, Q_0^b) \chi_b(n'') | \mu | \psi_a(r_i, Q_0^a) \chi_a(n') \rangle|^2 \quad (51)$$

It can be written as

$$W(a, n' \rightarrow b, n'') = W_{ab} |\langle \chi_b(n'') | \chi_a(n') \rangle|^2 \quad (52)$$

where W_{ab} is the purely electronic transition rate. The shape function of the transition is determined by the square of the vibrational overlap integrals, which are generally not zero. The

absorption bandshape at $T = 0$ K, where only the $n' = 0$ vibrational level is occupied, is then given by

$$I_{ab}(E) = I_0 \sum_{n''} \frac{S^{n''} \exp(-S)}{n''!} \delta(E_0 + n'' \hbar \omega - E) \quad (53)$$

where $E_0 = E_{b0} - E_{a0}$ is the energy of the transition between the zero vibrational levels of electronic states a and b . This is usually referred to as the *zero-phonon transition*, since $\sum_{n''} S^{n''} \exp(-S) / n''! = 1$. I_0 is the total intensity of the transition, which is independent of S . The intensity of the zero-phonon transition I_{00} is given by

$$I_{00} = I_0 \exp(-S) \quad (54)$$

so that if $S = 0$ all the intensity is contained in the zero-phonon transition. On the other hand, when S is large, the value of I_{00} tends to be zero and the intensity is concentrated in the *vibrational sidebands*. The single configurational coordinate model is relevant to the case of electron-vibrational structure in the spectra of molecules and the intensities so calculated fit observations rather well.

The net effect of this analysis is given in Fig. 10, which represents the absorption case when $S = 2$. At $T = 0$ we see a strong zero-phonon line, even stronger phonon-assisted transitions at $n'' = 1$ and 2, and then decreasing intensity in the phonon sidebands at $n'' = 3, 4, 5, \dots$. These individual transitions are represented by the vertical lines in this predicted spectrum. The envelope of these sidebands, given by the solid line, represents the effects of adding a finite width, $n'' \hbar \omega$, to each sideband feature.

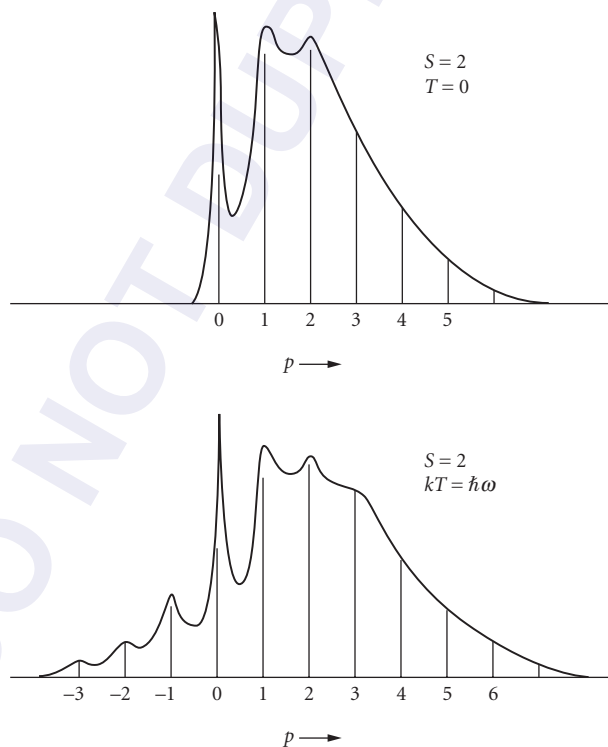


FIGURE 10 Showing the zero-phonon line and the Stokes shifted sideband on one side of the zero-phonon line at $T = 0$. When the temperature is raised the anti-Stokes sidebands appear on the other side of the zero-phonon line. (After Henderson and Imbusch.⁸)

The effect is to smear out the structure, especially at larger phonon numbers. Also in Fig. 10 we show the effect of temperature on the sideband shape in absorption, which is to reveal structure at lower energies than that of the zero-phonon line. Although the lineshape changes, the total intensity is independent of temperature, whereas the zero-phonon line intensity is given by

$$I_{00} = I_0 \exp \left[-S \coth \left(\frac{\hbar\omega}{2kT} \right) \right] \quad (55)$$

which decreases with increasing temperature. For values of $S \ll 1$, the phonon sideband intensity increases according to $I_0 S \coth(\hbar\omega/2kT)$. These effects are also shown in Fig. 10.

Three examples of the optical bandshapes in solid state spectra are considered. Figure 11 shows the luminescence spectrum of the molecular ion O_2^- in KBr measured at 77 K.³⁴ The O_2^- ion emits in the visible range from 400 to 750 nm, and the spectrum shown corresponds to the vibrational sidebands corresponding to $n'' = 5$ up to 13. The optical center is strongly coupled, $S = 10$, to the internal vibrational mode of the O_2^- ion with energy $\hbar\omega \approx 1000 \text{ cm}^{-1}$. However, as the detail of the $n'' = 8$ electronic vibrational transition is shown at $T = 4.2 \text{ K}$, there is also weak coupling $S \approx 1$ to the phonon spectrum of the KBr, where the maximum vibrational frequency is only about 200 cm^{-1} .

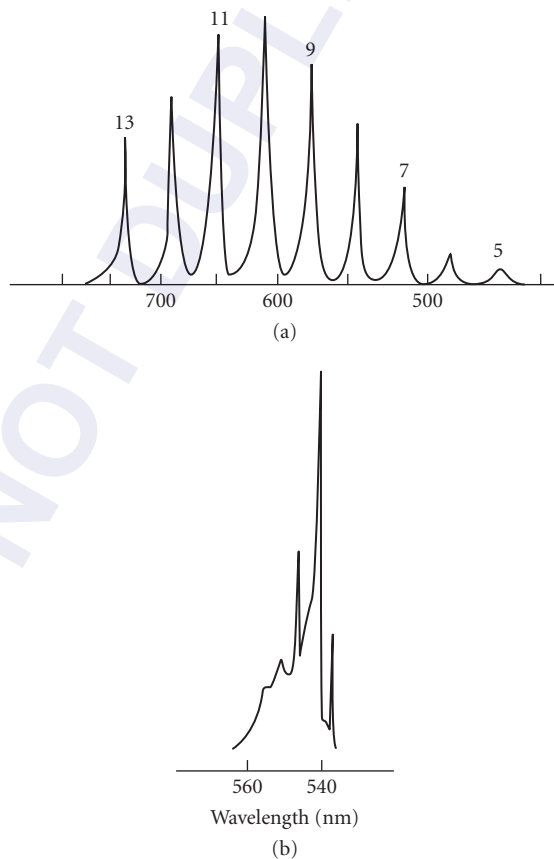


FIGURE 11 Photoluminescence spectrum of O_2^- in KBr at 77 K. (After Rebane and Rebane.³⁴)

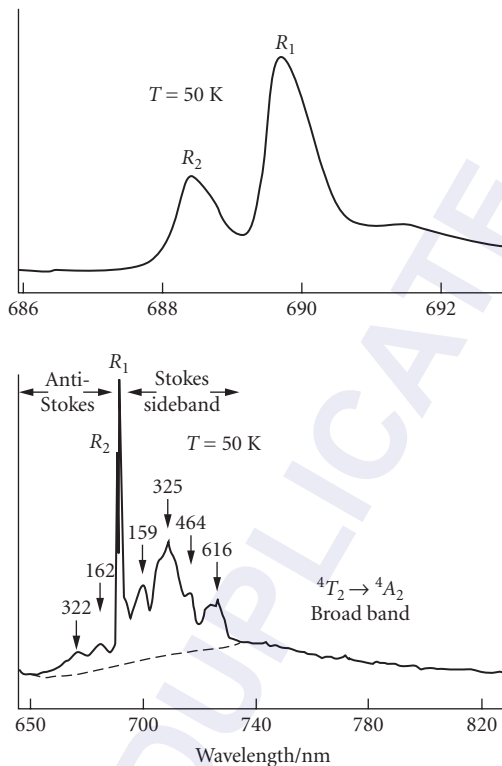


FIGURE 12 The photoluminescence spectrum of Cr^{3+} ions in $\text{Y}_3\text{Ga}_5\text{O}_{12}$.

The Cr^{3+} ion occupies a central position in the folklore of solid state spectroscopy. In yttrium gallium garnet, $\text{Y}_3\text{Ga}_5\text{O}_{12}$, the Cr^{3+} ion occupies a Ga^{3+} ion site at the center of six octahedrally-disposed O^{2-} ions. The crystal field at this site is intermediate in strength⁸ so that the 2E and 4T_2 states are mixed by the combined effects of spin-orbit coupling and zero-point vibrations. The emission then is a composite spectrum of ${}^2E \rightarrow {}^4A_2$ transition and ${}^4T_2 \rightarrow {}^4A_2$ transition, with a common radiative lifetime.³⁵ The composite spectrum, in Fig. 12, shows a mélange of the R -line and its vibronic sideband, ${}^2E \rightarrow {}^4A_2$ transition with $S \sim 0.3$, and the broadband ${}^4T_2 \rightarrow {}^4A_2$ transition for which $S \sim 6$.³⁶ Understanding this particular bandsape is complex.³⁵⁻³⁷

The final example, in Fig. 13, shows the absorption and emission spectra of F -centers in KBr .³⁸ This is a strongly coupled system with $S \approx 30$. Extension of the configurational coordinate model to the luminescence spectrum shows that the absorption and emission sidebands are mirror images of each other in the zero-phonon line. With S small (less than 6) there is structure in absorption and emission. However for S large, there is no structure, at least when a spectrum of vibrational modes interacts with the electronic states. The F -center represents strong coupling to a band of vibrational frequencies rather than to a single breathing mode of vibration. The effect of this is to broaden the spectrum to look like the envelope encompassing the spectrum of sharp sidebands shown in Fig. 10. In this case the zero-phonon line position is midway between the peaks in the absorption and emission bands of the F -center in KBr . Note also that as the temperature is raised the bands broaden and the peak shifts to longer wavelengths.

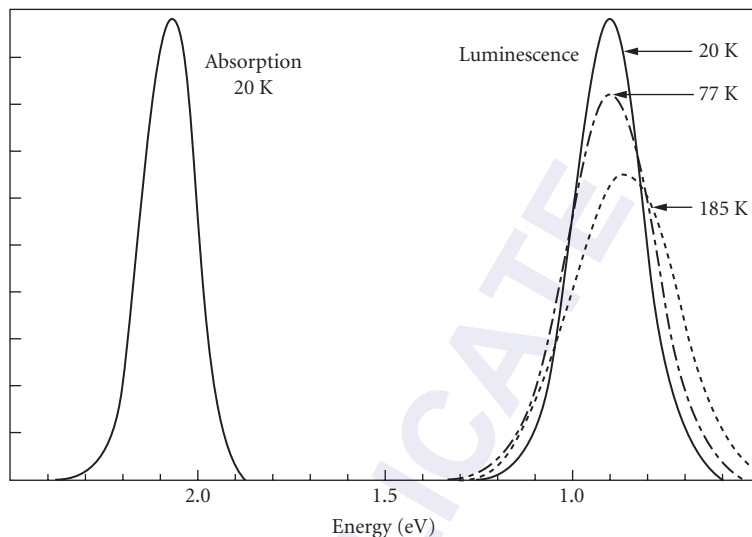


FIGURE 13 Optical absorption and photoluminescence of *F*-centers in KBr. (After Gebhardt and Kuhnert.³⁷)

10.12 REFERENCES

1. I. Newton, *Opticks* (1704).
2. T. Preston, *The Theory of Light*, 3d ed. (Macmillan, London, 1901). Gives an elegant account of the contributions of Huyghens, Young, and Fresnel to the development of the wave theory of light.
3. J. C. Maxwell, *Treatise on Electricity and Magnetism* (1873).
4. A. Einstein, *Ann. Phys.* **17**:132 (1905).
5. N. Bohr, *The Theory of Spectra and Atomic Constitution* (Cambridge University Press, 1922).
6. E. Schrödinger, *Ann. Phys.* **79**:361 (1926); **79**:489 (1926); **79**:734 (1926); **80**:437 (1926); **81**:109 (1926).
7. P. A. M. Dirac, *Quantum Mechanics* (Oxford University Press, 1927).
8. B. Henderson and G. F. Imbusch, *Optical Spectroscopy of Inorganic Solids* (Clarendon Press, Oxford, 1989).
9. W. R. Johnson and G. Soft, *At. Data Nucl. Data Tables* **33**:406 (1985).
10. A. I. Ferguson and J. M. Tolchard, *Contemp. Phys.* **28**:383 (1987).
11. B. P. Kibble, W. R. C. Rowley, R. E. Shawyer, and G. W. Series, *J. Phys. B* **6**:1079 (1973).
12. T. W. Hänsch, I. S. Shakin, and A. L. Schawlow, *Nature (Lond.)* **235**:63 (1972).
13. W. E. Lamb and R. C. Retherford, *Phys. Rev.* **72**:241 (1947); **79**:549 (1950); **81**:222 (1951); **86**:1014 (1952).
14. T. W. Hänsch, M. H. Nayfeh, S. A. Lee, S. M. Curry, and I. S. Shahin, *Phys. Ref. Lett.* **32**:1336 (1974).
15. J. E. M. Goldsmith, E. W. Weber, and T. W. Hänsch, *Phys. Lett.* **41**:1525 (1978).
16. B. W. Petley, K. Morris, and R. E. Shawyer, *J. Phys. B* **13**:3099 (1980).
17. D. N. Stacey, Private communication to A. I. Ferguson, cited in Ref. 10.
18. S. R. Amin, C. D. Caldwell, and W. Lichten, *Phys. Rev. Lett.* **47**:1234 (1981); P. Zhao, W. Lichten, H. P. Layer, and J. C. Bergquist, *ibid.* **58**:1293 (1987) and *Phys. Rev. A* **34**:5138 (1986).
19. C. E. Wieman and T. W. Hänsch, *Phys. Rev. Lett.* **36**:1170 (1976).
20. C. E. Wieman and T. W. Hänsch, *Phys. Rev. A* **22**:192 (1980).

21. J. R. M. Barr, J. M. Girkin, A. I. Ferguson, G. P. Barwood, P. Gill, W. R. C. Rowley, and R. C. Thompson, *Optics Commun.* **54**:217 (1985).
22. J. C. Slater, *Phys. Rev.* **35**:210 (1930).
23. M. Tinkham, *Group Theory and Quantum Mechanics* (McGraw-Hill, New York, 1964).
24. J. S. Griffith, *Theory of Transition Metal Ions* (Cambridge University Press, 1961).
25. C. E. Moore, *Atomic Energy Levels* [U.S. Nat. Bur. of Stand. Pub 467 (1950)].
26. W. Grotian *Graphische Darstellung der Spektren—on Atomen*, vol. 2 (Springer Verlag, Berlin, 1928).
27. W. Heitler and F. London, *Zeit. Phys.* **44**:455 (1927).
28. M. Born and J. P. Oppenheimer, *Ann. Phys.* **84**:457 (1927).
29. C. H. Townes and A. L. Schawlow, *Microwave Spectroscopy* (McGraw-Hill, New York, 1955). Gives an excellent account of the early days of what was then a novel sphere of study.
30. A. E. Ruark and H. C. Urey, *Atoms, Molecules and Quanta* (McGraw-Hill, New York, 1930).
31. J. C. Slater, *Quantum Theory of Molecules and Solids*, vol. 1 (McGraw-Hill, New York, 1963).
32. G. Herzberg, *Molecular Spectra and Molecular Structure* (vols I and II, Van Nostrand, New York 1953).
33. C. N. Banwell, *Fundamentals of Molecular Spectroscopy* (McGraw-Hill, U.K., 1983).
34. K. K. Rebane and L. A. Rebane in *Optical Properties of Ions in Solids*, B. Di Bartolo (ed.) (Plenum Press, New York, 1975). See also K. K. Rebane, *Impurity Spectra of Solids* (Plenum Press, New York, 1970).
35. M. Yamaga, B. Henderson, and K. P. O'Donnell, *J. Phys. (Cond. Matt.)* **1**:9175 (1989).
36. M. Yamaga, B. Henderson, K. P. O'Donnell, C. Trager-Cowan, and A. Marshall, *App. Phys.* **B50**:425 (1990).
37. M. Yamaga, B. Henderson, K. P. O'Donnell, F. Rasheed, Y. Gao, and B. Cockayne, *App. Phys. B.* **52**:225 (1991).
38. W. Gebhardt and H. Kuhnert, *Phys. Lett.* **11**:15 (1964).

ANALOG OPTICAL SIGNAL AND IMAGE PROCESSING

Joseph W. Goodman

*Department of Electrical Engineering
Stanford University
Stanford, California*

11.1 GLOSSARY

$C(\tau)$	cross-correlation function
d, z	distances
f	focal length
f_x, f_y	spatial frequencies
$H(f_x, f_y)$	transfer function
$I(x, y)$	intensity distribution
i	square root of negative one
M	matrix
$t_A(x, y)$	amplitude transmittance of a transparency
$U(x, y)$	phasor representation of a monochromatic field
u, v	vectors
V	velocity of propagation
x, y	spatial coordinates
θ_B	Bragg angle
λ	wavelength
Λ	period of grating
ν	optical frequency
τ	time delay

11.2 INTRODUCTION

The function of signal and image processing systems is the modification of signals and images to allow information extraction by a human observer or, alternatively, to allow fully automatic information extraction without human intervention. The origins of optical information processing are several, but certainly the invention of various techniques for visualizing the phase distribution of optical wavefronts qualifies (e.g., Ref. 1), as do the famous Abbe-Porter experiments.^{2,3} Starting in the 1950s, more general information processing tasks were undertaken with the help of optics.^{4,5} This chapter presents an overview of such methods.

Optical systems are of interest both for digital processing of information and for analog processing of information. Our attention here will be restricted only to analog processing operations, which are far more mature and well developed than digital optical methods.

Certain basic assumptions will be used throughout and are detailed here. First, monochromatic optical signals will be represented by complex phasor field distributions, with the direction of phasor rotation being assumed to be clockwise. Similarly, time-varying optical fields will be represented by complex analytic signals, again with the direction of rotation in the complex plane being clockwise. In both cases the underlying real signals are recoverable as the real part of the complex representations. In all cases, a small-angle assumption will be employed, allowing paraxial approximations to be used. Polarization effects will generally be ignored, it being assumed that a scalar theory of light propagation is sufficiently accurate for our purposes.⁶ The intensity of the optical waves, which is proportional to power density and is the observable quantity in an optical experiment, is defined as the squared magnitude of the complex fields.

It is very important to distinguish at the start between *coherent* and *incoherent* optical systems. For a review of optical coherence concepts, see Chap. 5 of this volume. For our purposes, we will regard an optical signal as coherent if the various optical contributions that produce an output add on an amplitude basis, with fixed and well-defined relative phases. Signals will be regarded as incoherent if the various contributions that add to produce an output at any point have time-varying phase relations, and therefore must add on an intensity or average-power basis.

11.3 FUNDAMENTAL ANALOG OPERATIONS

The fundamental components of any linear processing operation are addition and multiplication. We consider each of these operations separately.

Addition

Analog addition takes place in optical systems when light waves or wave components are superimposed. The exact nature of the addition depends on whether the optical components are mutually coherent or incoherent. In the coherent case, addition of complex phasor field components takes place. Thus if the $U_n(x, y)$ represent various optical field components that are superimposed at a given point (x, y) at the output, the resultant optical field $U(x, y)$ is given by

$$U(x, y) = \sum_n U_n(x, y) \quad (1)$$

Note that the result of such a superposition depends on the phases of the individual components.

On the other hand, if the various optical contributions at (x, y) are mutually incoherent, the addition takes place on an intensity basis. The resultant intensity $I(x, y)$ is given by

$$I(x, y) = \sum_n I_n(x, y) \quad (2)$$

where the $I_n(x, y)$ are the component intensity contributions. In this case the component intensity contributions are always positive and real, as is the resultant intensity.

In view of the above two equations, an important point can now be made. Coherent optical systems are *linear in complex amplitude*, while incoherent optical systems are *linear in intensity*. The design of an analog processing system thus depends fundamentally on whether the illumination used in the system is coherent or incoherent.

Multiplication

Analog multiplication takes place in optical systems as light passes through an absorbing or phase-shifting structure. If we define the complex amplitude transmittance $t_A(x, y)$ of a transmitting structure as the ratio of the transmitted complex field to the incident complex field, then analog multiplication in a coherent system is represented by

$$U_t(x, y) = t_A(x, y)U_i(x, y) \quad (3)$$

where $U_i(x, y)$ is the incident optical field and $U_t(x, y)$ is the transmitted optical field.

When the optical system is incoherent, then we define an intensity transmittance $t_I(x, y)$ as the ratio of the transmitted optical intensity to the incident optical intensity. Analog multiplication in such systems is represented by

$$I_t(x, y) = t_I(x, y)I_i(x, y) \quad (4)$$

Thus we have seen that the fundamental analog operations of addition and multiplication are quite naturally available in optical systems. It should be kept in mind that the operation of integration is just a generalization of addition, involving addition of an infinite number of infinitesimal components.

11.4 ANALOG OPTICAL FOURIER TRANSFORMS

Perhaps the most fundamental optical analog signal- and image-processing operation offered by optical systems is the Fourier transform. Such transforms occur quite simply and naturally with coherent optical systems. While Fourier sine and cosine transforms can be performed with incoherent light, the methods used are more cumbersome than in the coherent case and the numbers of resolvable spots involved in the images and transforms are more restricted. Therefore we focus here on Fourier transforms performed by coherent optical systems. The Fourier transform is normally two-dimensional in nature (image processing), although it can be restricted to a single dimension if desired (signal processing).

Focal-Plane-to-Focal-Plane Geometry

The optical system required to perform a two-dimensional Fourier transform is remarkably simple as shown in Fig. 1. We begin with a spatially coherent source of quasi-monochromatic light (a source that is both spatially and temporally coherent). The light from that point source is collimated by a positive lens and a transparency of the image to be Fourier transformed is introduced in the front

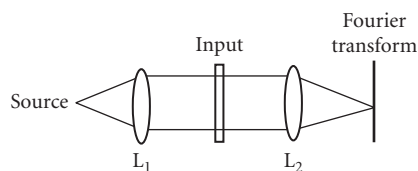


FIGURE 1 Simple optical system for performing a two-dimensional Fourier transform.

focal plane of a second positive lens (L_2). Under such conditions, the complex field appearing across the rear focal plane of that lens can be shown to be the two-dimensional Fourier transform of the complex field transmitted by the input transparency, as given by

$$U_f(x, y) = \frac{1}{i\lambda f} \iint_{-\infty}^{\infty} U_i(\xi, \eta) \exp\left[-i\frac{2\pi}{\lambda f}(x\xi + y\eta)\right] d\xi d\eta \quad (5)$$

Here λ is the optical wavelength, f is the focal length of the second lens, U_f is the field distribution across the back focal plane of lens L_2 , and U_i is the field transmitted by the transparency in the front focal plane.

An intuitive explanation for the occurrence of this elegant relationship between the fields in the two focal planes can be presented as follows. If we were to mathematically Fourier transform the fields transmitted by the input transparency, each such Fourier component could be recognized as a different plane wave component of the transmitted field. Each such Fourier component represents a plane wave traveling in a unique direction with respect to the optical axis. Such representations are the basis for the so-called angular spectrum of plane waves often used in the analysis of optical wavefields (see, for example, Ref. 6, p. 48). Now consider the effect of the positive lens on a single Fourier component, i.e., a plane wave traveling at a particular angle with respect to the optical axis. As that plane wave passes through the lens L_2 , it is changed into a spherical wave converging toward a focus in the rear focal plane, in a particular location determined by that plane wave's propagation direction. Thus the intensity of light at a given coordinate in the rear focal plane is proportional to the energy contained by the input wavefield at a particular Fourier spatial frequency. Hence the distribution of energy across the rear focal plane is a representation of the distribution of energy across the various spatial frequencies contained in input transparency.

Other Fourier Transform Geometries

A slightly more general configuration is one in which the input transparency is placed at an arbitrary distance d in front of the lens L_2 , while the field is again considered in the rear focal plane of that lens. The relation between the input and output fields remains of the general form of a two-dimensional Fourier transform, but with the complication that a multiplicative quadratic phase factor is introduced, yielding a relation between input and focal-plane fields given by

$$U_f(x, y) = \frac{\exp\left[i\frac{k}{2f}\left(1-\frac{d}{f}\right)(x^2 + y^2)\right]}{i\lambda f} \iint_{-\infty}^{\infty} U_i(\xi, \eta) \exp\left[-\frac{i2\pi}{\lambda f}(x\xi + y\eta)\right] d\xi d\eta \quad (6)$$

Three additional Fourier transform geometries should be mentioned for completeness. One is the case of an object transparency placed directly against the lens in Fig. 1, either in front or in back of the lens. This is a special case of Eq. (6), with d set equal to 0, yielding

$$U_f(x, y) = \frac{\exp\left[i\frac{k}{2f}(x^2 + y^2)\right]}{i\lambda f} \iint_{-\infty}^{\infty} U_i(\xi, \eta) \exp\left[-\frac{i2\pi}{\lambda f}(x\xi + y\eta)\right] d\xi d\eta \quad (7)$$

Another situation of interest occurs when the object transparency is located behind the lens L_2 , a distance d from the focal plane, as shown in Fig. 2. In this case the relationship between the fields transmitted by the object and incident on the focal plane becomes

$$U_f(x, y) = \frac{\exp\left[i\frac{k}{2d}(x^2 + y^2)\right]}{i\lambda d} \frac{f}{d} \iint_{-\infty}^{\infty} U_i(\xi, \eta) \exp\left[-\frac{i2\pi}{\lambda d}(x\xi + y\eta)\right] d\xi d\eta \quad (8)$$

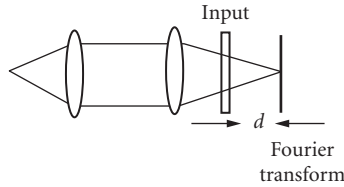


FIGURE 2 Fourier transform geometry with the object behind the lens.

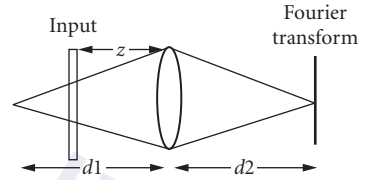


FIGURE 3 Fourier transform geometry using a single lens.

Note that now the scaling distance in the Fourier kernel is d , rather than the focal length f . Therefore, by moving the object toward or away from the focal plane, the transform can be made smaller or larger, respectively.

While the Fourier transform plane in all the above examples has been the rear focal plane of the lens L_2 , this is not always the case. The more general result states that the Fourier transform always appears in the plane where the original point-source of illumination is imaged by the optical system. In the previous examples, which all involved a collimating lens L_1 before the object transparency, the source was indeed imaged in the rear focal plane of L_2 , where we asserted the Fourier transform lies. However, in the more general case depicted in Fig. 3, the point source of light lies in plane P_1 and its image lies in plane P_2 , which for this geometry is the Fourier transform plane. A single lens L_1 performs both imaging of the source and Fourier transformation of the fields transmitted by the input transparency. If the input is placed to the right of the lens, at distance d from the image of the source, then the Fourier transform relation is identical to that presented in Eq. (8), for it does not matter what optical system illuminated the transparency with a converging spherical wave, only what distance exists between the input and the plane where the source is imaged.

If the input transparency is placed to the left of the single lens, as shown in Fig. 3, the resulting relationship between the fields transmitted by the object U_i and the fields across the plane where the source is imaged, U_f , becomes

$$U_f(x, y) = \frac{d_1}{i\lambda d_2(d_1 - z)} \exp \left\{ i \frac{k}{2} \left(\frac{1}{d_2} - \frac{z d_1}{d_2^2(d_1 - z)} \right) \right\} \times \iint_{-\infty}^{\infty} U_i(\xi, \eta) \exp \left\{ -i \frac{2\pi d_1}{\lambda d_2(d_1 - z)} [x\xi + y\eta] \right\} d\xi d\eta \quad (9)$$

where the meanings of z , d_1 , and d_2 are shown in Fig. 3, and $k = 2\pi/\lambda$. In this relation, d_1 and d_2 are connected through the lens law, $1/d_1 + 1/d_2 = 1/f$. It can be shown quite generally that the effective distance d associated with the Fourier transform kernel is $d = d_2/d_1(d_1 - z)$, while the quadratic phase factor is that associated with a diverging spherical wave in the transform plane that appears to have originated on the optical axis in the plane of the input transparency.

11.5 SPATIAL FILTERING

Given that Fourier transforms of optical fields occur so directly in coherent optical systems, it seems natural to consider the intentional manipulation of such spectra for the purposes of signal or image processing. Given that a signal or image has been introduced into the coherent optical system, either by means of photographic film or by means of an electronically or optically controlled spatial light modulator (see Chap. 6, "Acousto-Optic Devices," in Vol. V), the idea is to insert in the plane where the Fourier transform occurs a transparency (again either film or a spatial light modulator) which intentionally alters the fields transmitted through that plane. A second Fourier transforming lens

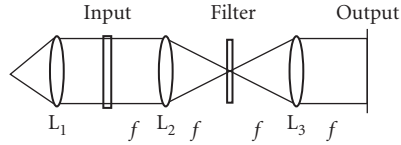


FIGURE 4 Spatial filtering system based on double Fourier transform.

then returns the observer to an image plane, where the filtered version of the input can be measured or extracted. The simplest geometry from the conceptual point of view is that shown in Fig. 4.

The lens L_1 is again a collimating lens, the lens L_2 is a Fourier transforming lens, and the lens L_3 is a second Fourier transforming lens. The fact that a sequence of two Fourier transforms takes place, rather than a Fourier transform followed by an inverse Fourier transform, results simply in an inversion of the image at the output of the system.

Systems of this type form the basis for coherent optical spatial filtering, although the detailed geometry of the layout may vary. We will discuss several such spatial filtering problems in later sections. For the moment it suffices to say that if a filtering system is desired to have a transfer function $H(f_x, f_y)$ then the amplitude transmittance of the transparency inserted in the Fourier plane should be

$$t_A(\xi, \eta) = H\left(\frac{\xi}{\lambda f}, \frac{\eta}{\lambda f}\right) \quad (10)$$

where λ has been defined, f is the focal length of the Fourier transforming lenses (assumed identical), and (ξ, η) represent the spatial coordinates in the filter plane.

11.6 COHERENT OPTICAL PROCESSING OF SYNTHETIC APERTURE RADAR DATA

The earliest serious application of coherent optics to signal processing was to the problem of processing data gathered by synthetic aperture radars. We explain the synthetic aperture principle, and then discuss optical signal-processing architectures that have been applied to this problem.

The Synthetic Aperture Radar Principle

The synthetic-aperture radar problem is illustrated in Fig. 5. An aircraft carrying a stable local oscillator and a side-looking antenna flies a straight-line path, illuminating the terrain with microwave energy and detecting the returned energy reflected and scattered from that terrain. In the simplest case, resolution in range (i.e., perpendicular to the aircraft flight path) is obtained by pulse echo timing, the usual radar range-measurement technique. Resolution in azimuth (the direction parallel to the flight path) is obtained by processing the Doppler-shifted returns, as will be explained. For the purpose of explaining the azimuth imaging, we neglect the pulsed nature of the radiation emitted by the aircraft, an approximation allowable because of the pulse-to-pulse coherence of the signals. The goal of the system is to obtain a two-dimensional image of the microwave reflectivity of the ground illuminated by the aircraft. The resolution of the system is not limited by the size of the antenna that is carried by the aircraft—in fact resolution increases as the size of the antenna is decreased. The system coherently combines the signals received along a portion of the flight path, thereby synthesizing the equivalent of a much longer antenna array.

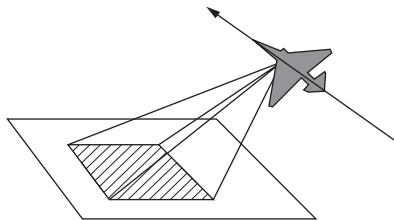


FIGURE 5 Aircraft flight path.

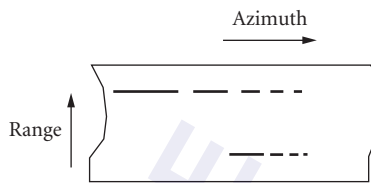


FIGURE 6 Recording format for signals received.

If we consider the signal received at the aircraft as a function of time, originating from a single point scatterer on the ground, that signal will suffer an upward frequency shift as the aircraft approaches the scatterer and a downward frequency shift as the aircraft flies away from the scatterer. This chirping signal is beat against the stable local oscillator in the aircraft, a bias is added, and the new signal is then recorded on a strip of film. Figure 6 shows the recording format. In the vertical direction, different scatterers are separated by the pulse echo timing, each being imaged on a separate horizontal line of the film. In the horizontal direction, the time histories of the chirping azimuth signals from different scatterers are recorded.

The signal recorded from a single scatterer is in fact an off-axis one-dimensional Fresnel zone plate, and as such is capable of imaging light in the horizontal direction to a focus. Such a focus constitutes the azimuthal image of the point scatterer that gave rise to this zone plate. However, the chirp rates, and therefore the focal lengths, of the zone plates produced by scatterers at different ranges are unfortunately not the same. The focal length is in fact proportional to the distance of the scatterer from the aircraft. Thus the focal points of scatterers at different ranges from the aircraft lie on a tilted plane with respect to the film plane, whereas the range images lie in the plane of the film. Thus the optical processing system must be designed to bring the two different images into coincidence.

Optical Processing Systems

The earliest system used for optical processing of synthetic aperture radar data is illustrated in Fig. 7.⁷ This processor uses a conical optical element, called an axicon, to change the focal lengths of all horizontal zone plates to infinity, thus moving the azimuth image to infinity. A cylindrical lens is placed one focal length from the film to likewise move the range image to infinity, and a spherical lens is placed one focal length from the final image plane to bring the infinitely distant azimuth and range planes back to a common focus.

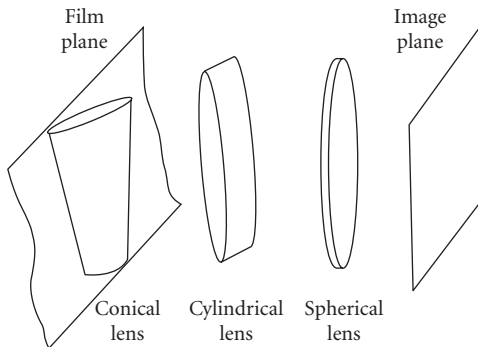


FIGURE 7 Processor using an axicon.

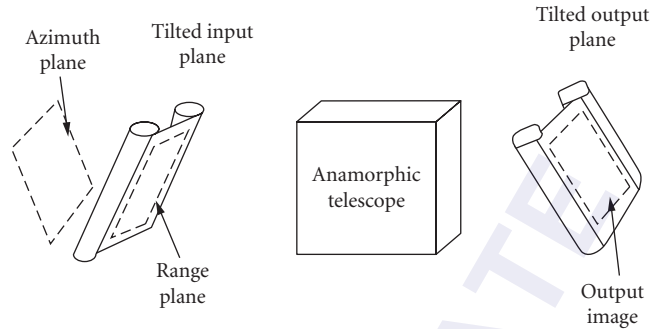


FIGURE 8 The tilted-plane processor.

The magnification achieved by such a system is a function of range, so the output is recorded through a vertical slit. As the input film is drawn through the system, an output film strip is likewise drawn past the slit in synchronism, with the result that an image with proper magnification is recorded at the output.

Following the use of such an optical system to produce images, a far more sophisticated processing system known as the “tilted-plane processor” was developed.⁸ The architecture of this system is illustrated in Fig. 8. In this case an anamorphic telescope is used to bring the range and azimuth planes into coincidence with a constant magnification, allowing a full two-dimensional image to be recorded at the output at one time. Again motion of the input film and the output film takes place in synchronism, but the throughput of the system is much higher due to the absence of the output slit.

From the very fundamental work on processing synthetic aperture radar signals at the University of Michigan during the late 1950s and early 1960s came a multitude of extraordinary inventions, including holograms with an off-axis reference wave and the holographic matched filter, or Vander Lugt filter, discussed in Sec. 11.8.

11.7 COHERENT OPTICAL PROCESSING OF TEMPORAL SIGNALS

An important subclass of information-processing operations is those that are applied to one-dimensional signals that are functions of time. Such signals can be processed by coherent optical-filtering systems once a suitable transducer is found to convert time-varying voltages representing the signals into space-varying wavefields. The best developed and most common of such transducers is the acousto-optic cell.^{9,10}

Acousto-Optic Cells for Inputting Signals

A time-varying electrical signal can be changed to an equivalent one-dimensional space-varying distribution of field strength by means of acousto-optic devices. In bulk form, such devices consist of a transducer, to which a time-varying voltage representing an RF signal is applied, and a transparent medium into which compressional waves are launched by the transducer. The RF signal is assumed to contain a carrier frequency, which generates a dynamic grating and, when illuminated by coherent light, produces a number of different diffraction orders, of which the +1 and -1 orders are of primary interest. Any modulation, in amplitude or phase, that may be carried by the RF signal is transferred to the spatial distributions of these diffraction orders.

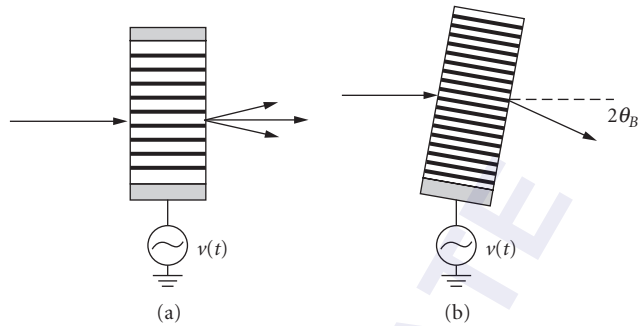


FIGURE 9 Acousto-optic diffraction in the (a) Raman-Nath and (b) Bragg regimes.

Acousto-optic diffraction is characterized as either *Raman-Nath* diffraction or *Bragg* diffraction, depending on the relations that exist between the cell thickness and the period of the acoustic wave generated by the RF carrier. For cells that are sufficiently thin, Raman-Nath diffraction takes place. The acousto-optic cell then acts like a thin phase grating, generating a multitude of diffraction orders. For cells that are sufficiently thick, Bragg diffraction takes place. In this case, high diffraction efficiency into a single grating order can be achieved if the acoustic grating is illuminated at the Bragg angle θ_B which satisfies

$$\sin \frac{\theta_B}{2} = \frac{\lambda}{2\Lambda} \quad (11)$$

In this case most of the optical power is transferred to the +1 diffraction order, and other orders, including the -1 and 0 orders can be neglected.

Figure 9 illustrates Raman-Nath and Bragg diffraction from an acousto-optic cell. $v(t)$ represents the voltage driving the cell transducer. For modern-day signal-processing applications, which involve very high microwave frequencies, the Bragg cell is invariably used, and the situation on the right-hand side of the figure is the one of interest.

The signal $v(t)$ is of the form (in complex notation)

$$v(t) = A(t) \exp[-i(2\pi\nu_0 t + \theta(t))] \quad (12)$$

where $A(t)$ is the amplitude modulation of the carrier, $\theta(t)$ is the phase modulation of the carrier, and ν_0 is the center frequency. If the speed of propagation of acoustic waves in the medium of the Bragg cell is V , then emerging from the right of that cell will be a spatial complex field distribution of the form

$$U(x, t) = A(x - Vt) \exp[-i\theta(x - Vt)] \quad (13)$$

where the dependence on y is suppressed due to uniformity of U in that dimension. Thus the temporal structure of the signal $v(t)$ has been changed to a spatial structure of the optical field $U(x, t)$.

The Bragg Cell Spectrum Analyzer

The most common use of coherent optics for signal processing is a method for finding and displaying the frequency (Fourier) spectrum of the electrical signal $v(t)$ applied to the cell. To construct such a spectrum analyzer, we follow the Bragg cell of Fig. 9 with a positive lens, which then Fourier transforms the wavefield emerging from the cell, as shown in Fig. 10. A detector array placed in the Fourier

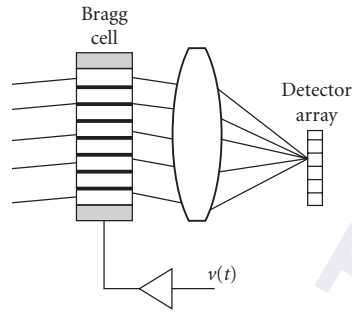


FIGURE 10 Bragg cell spectrum analyzer.

plane then measures the amount of signal power present in each frequency bin subtended by a detector element. Note the spectrum analysis is performed over a finite sliding window, namely, the window of time stored in the Bragg cell itself. Figure 10 shows a diagram illustrating the Bragg cell spectrum analyzer.

Assuming perfect optics, the resolution of such a spectrum analyzer is determined by the diffraction limit associated with the space window that is being transformed. The spatial dimension of a resolution element is given by $\Delta x = (\lambda f/L)$ where L is the length of the cell and f is again the focal length of the lens. Given the mapping from time to space that takes place in the cell, it follows that the temporal resolution of the spectrum analyzer (in hertz) is $\Delta \nu = (V/L)$.

Bragg cell spectrum analyzers have been built with center frequencies of more than 1 GHz, with bandwidths approaching 1 GHz, and time bandwidth products (equivalent to the number of resolvable spectral elements) of the order of 1000. While the vast majority of work on this type of spectrum analyzer has used bulk devices (e.g., bulk Bragg cells, discrete lenses, etc.), work has also been carried out on integrated versions. Such devices use planar waveguides rather than free-space propagation, surface acoustic waves rather than bulk acoustic waves, integrated optic lenses, etc. Such systems are more compact than those based on bulk approaches, but their performance is so far somewhat inferior to that of the more conventional bulk systems.

The chief difficulty encountered in realizing high-performance Bragg cell spectrum analyzers is the dynamic range that can be achieved. The dynamic range refers to the ratio of the largest spectral component that can be obtained within the limit of tolerable nonlinear distortion, to the smallest spectral component that can be detected above the noise floor.

Acousto-Optic Correlators

Many signal detection problems require the realization of correlators that produce cross-correlations between a local reference signal and an incoming unknown signal. A high cross-correlation between the reference and the unknown signal indicates a high degree of similarity between the two signals, while a low correlation indicates that the two signals are not very similar. Thus correlators play an important role in signal detection and recognition.

Given two complex-valued signals $v_1(t)$ and $v_2(t)$, the cross-correlation between those signals is defined as

$$C(\tau) = \int_{-\infty}^{\infty} v_1(t) v_2^*(t - \tau) dt \quad (14)$$

When v_1 and v_2 are identical, $C(\tau)$ achieves a peak value at the relative delay τ that causes the two signals to be identically aligned in time.

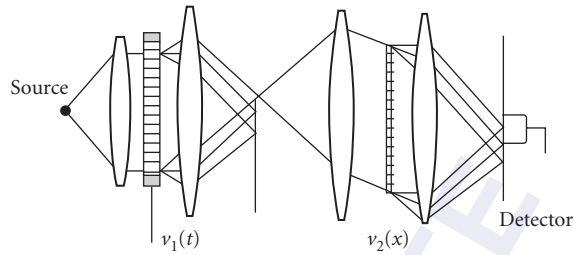


FIGURE 11 The time integrating correlator.

Two distinctly different architectures have been developed for using acousto-optic systems for cross-correlating wideband signals. We discuss each of these techniques separately.

The Space-Integrating Correlator The older of the two architectures is known as the space-integrating correlator. As the name indicates, the integration inherent in the correlation operation is carried out over space. The variable delay τ is achieved by allowing one signal to slip past the other in time.

Figure 11 shows the structure of a time-integrating correlator. One of the signals, $v_1(t)$, is introduced by means of an input Bragg cell. Spatial filtering is used to eliminate any residual of the zeroth and unwanted first diffraction orders, retaining only a single first order. The second signal, the reference $v_2(t)$, is stored on a transparency, complete with a spatial carrier frequency representing the center frequency and acting as a high-frequency amplitude- and phase-modulated grating. The integration over space is provided by the final output lens. The particular diffraction order used in the final transparency is chosen to yield the conjugate of $v_2(t)$. A point detector is used at the output, and different relative delays between the two signals are achieved simply by allowing $v_1(t)$ to slide through the Bragg cell.

The Time-Integrating Correlator A different approach to realizing the temporal cross-correlation operation is the so-called time-integrating correlator.^{11,12} The architecture of a time-integrating correlator is illustrated in Fig. 12. Spatial filtering selects one component that has undergone zeroth-order diffraction by the first cell and first-order diffraction by the second, and another component that has undergone first-order diffraction by the first cell and zeroth-order by the second. These two components interfere on a time-integrating detector array.

As the name implies, the correlation integration is in this case carried out by temporal integration on an array of time-integrating detectors. Note that the two electrical signals are introduced at opposite ends of the Bragg cells, with the result that at each spatial position on the Bragg cell pair the two signals have been delayed relative to one another by different amounts, thus introducing the

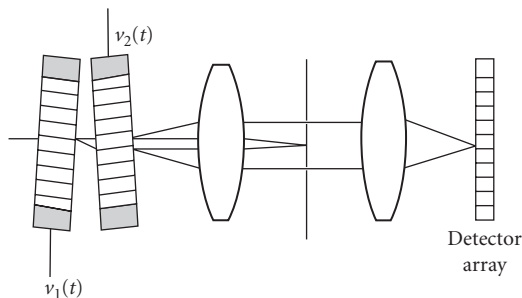


FIGURE 12 Time-integrating correlator.

relative delay required in the correlation integral. The lens on the right images the pair of Bragg cells onto the detector array. Thus different detector elements measure the interference of the two signals with different relative delays, one portion of which yields

$$\operatorname{Re}\{C(x)\} = \operatorname{Re} \left\{ \int_T v_1 \left[t - \left(\frac{x+L/2}{V} \right) \right] v_2^* \left[t + \left(\frac{x-L/2}{V} \right) \right] dt \right\} \quad (15)$$

which is the real part of the correlation integral of interest. Here L represents the length of the Bragg cells, V the velocity of propagation of acoustic waves, T the total integration time of the detector, and x is the position of a particular detector on the detector array at the output. Note that for the position x on the detector array the two signals have been delayed relative to each other by the amount

$$\tau = \frac{2x}{V} \quad (16)$$

Other variants of both space-integrating and time-integrating correlators are known but will not be presented here. Likewise, many architectures for other types of acousto-optic signal processing are known. The reader may wish to consult Ref. 13 for more details.

11.8 OPTICAL PROCESSING OF TWO-DIMENSIONAL IMAGES

Because optical systems are fundamentally two-dimensional in nature, they are well suited to processing two-dimensional data. The most important type of two-dimensional data is imagery. Thus we consider now the application of optical information processing systems to image processing. The applications of optical processing in this area can be divided into two categories: (1) pattern detection and recognition, and (2) image enhancement.

Optical Matched Filtering for Pattern Recognition

By far the most well-known approach to pattern recognition is by means of the matched filter.¹⁴ While this approach has many known defects in the pattern recognition application, it nonetheless forms the basis for many other more sophisticated approaches.

The Matched Filter A linear invariant filter is said to be “matched” to a certain spatial image $s(x, y)$ if the impulse response (point-spread function) $h(x, y)$ of that filter is of the form

$$h(x, y) = s^*(-x, -y) \quad (17)$$

When a general signal $v(x, y)$ is applied to the input of such a filter, the output (the convolution of the input and the impulse response) is given by

$$\begin{aligned} w(x, y) &= \iint_{-\infty}^{\infty} v(\xi, \eta) h(x-\xi, y-\eta) d\xi d\eta \\ &= \iint_{-\infty}^{\infty} v(\xi, \eta) s^*(\xi-x, \eta-y) d\xi d\eta \end{aligned} \quad (18)$$

which is the *cross-correlation* between the signals $v(x, y)$ and $s(x, y)$. Thus the output of a matched filter is the cross-correlation between the input signal and the signal for which the filter is matched.

In the frequency domain, the convolution relation becomes a simple product relation. The frequency domain equivalent of Eq. (18) is

$$W(f_x, f_y) = H(f_x, f_y)V(f_x, f_y) = S^*(f_x, f_y)V(f_x, f_y) \quad (19)$$

Thus the transfer function of the matched filter is the complex conjugate of the spectrum of the signal to which the filter is matched.

The coherent optical realization of the matched filter utilizes a system identical with that shown previously in Fig. 4, where the Fourier domain transparency is one with amplitude transmittance proportional to $S^*(f_x, f_y)$. The output of the filter, appearing at the plane on the far right in Fig. 4, consists of a bright spot at each location where the signal $s(x, y)$ is located within the input field.

Prior to 1964, a key difficulty in the realization of matched filtering systems was the construction of the Fourier domain filter with the proper amplitude transmittance. To control the amplitude and phase transmittance through the Fourier plane in a relatively complicated manner was often beyond the realm of possibility. However, in 1964 Vander Lugt published his classic paper on holographically recorded matched filters, and many new applications became possible.

The Vander Lugt Filter The method introduced by Vander Lugt¹⁵ for recording matched filters is shown in Fig. 13. It is assumed that a mask can be constructed with amplitude transmittance proportional to the desired impulse response $s^*(-x, -y)$, which in pattern recognition problems is often real and positive. That mask is illuminated by coherent light and Fourier transformed by a positive lens. In the Fourier domain, the spectrum $S^*(f_x, f_y)$ is allowed to interfere with an angularly inclined reference wave, often a plane wave. The result is an intensity pattern with a high spatial frequency carrier, which is amplitude modulated by the amplitude distribution associated with the incident spectrum, and phase modulated by the phase distribution of that spectrum. This recording is in fact a Fourier hologram of the desired point-spread function¹⁶ (see also Chap. 33, "Holography and Holographic Instruments," in this volume). The film used for recording in the Fourier domain responds to the incident optical intensity. With proper processing of the film, one of the grating orders of the resulting transparency yields a field component proportional to the desired field,

$$t_A(\xi, \eta) \approx S^*\left(\frac{\xi}{\lambda f}, \frac{\eta}{\lambda f}\right) \exp(-i2\pi\alpha\eta) \quad (20)$$

where (ξ, η) are the spatial coordinates in the filter plane, and α is the spatial frequency of the carrier. Thus the transmittance required for realizing the matched filter has been achieved, with the exception of the linear exponential term, which serves simply to shift the desired output off the optical axis.

The filter constructed as above is placed in the Fourier domain of the system in Fig. 4 and provided the correct region of the output plane in that figure is examined, the matched filter response is found. In a second region of the output plane, mirror symmetric with the matched filter region, the convolution of the signal $s(x, y)$ and the input $v(x, y)$ can be found.

Prior to Vander Lugt's invention the only matched filters that could be realized in practice were filters with very simple transfer functions $S^*(f_x, f_y)$. The significance of the Vander Lugt filter is that it extends the domain for which filters can be realized to those with reasonably simple impulse responses $s(-x, -y)$, a case more typically encountered in pattern recognition.

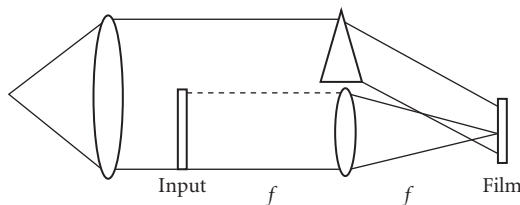


FIGURE 13 Recording a Vander Lugt filter.

Deficiencies of the Matched Filter Concept While the Vander Lugt filter provides an elegant solution to the problem of realizing coherent optical matched filters, nonetheless the use of coherent optics for pattern recognition has been very restricted in its applications. A major reason for this limited applicability can be traced to deficiencies of the matched filter concept itself, and is not due to the methods used for optical realization. The matched filter, in its original form, is found to be much too sensitive to parameters for which lack of sensitivity would be desired. This includes primarily rotation of the image and scale change of the image. Thus a matched filter that responds perfectly to the desired signal in its original rotation and scale size may not respond at all to that signal when it is rotated and magnified or demagnified.

Many attempts have been made to remove these undesired sensitivities of the matched filter (see for example, Refs. 17 and 18). These include the use of Mellin transforms and polar coordinate transformations to remove scale-size sensitivity and rotation sensitivity, and the use of circular harmonic decompositions to remove rotation sensitivity. These attempts have had varying degrees of success, but unfortunately they generally destroy one insensitivity that is present for a conventional matched filter, namely, insensitivity to translation of the signal. For a conventional matched filter, when an input $s(x, y)$ is translated, the resulting bright spot at the output translates in response. Realization of rotation invariance generally removes translation insensitivity, a serious loss.

Unfortunately, to date there have been no commercially successful applications of coherent optical matched filtering to pattern recognition, although several attempts to commercialize the technology have been made.

Coherent Optical Image Enhancement

Coherent optical spatial filtering systems can also be applied to the problem of image enhancement.⁴ Image enhancement problems come in a wide variety of types, ranging from simple operations, such as the suppression of periodic noise in a nonperiodic image, to more complex operations, such as restoring an image that has been degraded by a known blur. We focus here on image restoration, since it is the most challenging of these problems.

The Inverse Filter A common type of image restoration problem arises when an image produced by an incoherent imaging system has been blurred by a known, space-invariant, linear point-spread function. Let $i(x, y)$ represent the intensity of the blurred image, $o(x, y)$ represent the intensity of the object, and $b(x, y)$ represent the intensity point-spread function of the blur. These three quantities are related through a convolution equation,

$$i(x, y) = \iint_{-\infty}^{\infty} b(x - \xi, y - \eta) o(\xi, \eta) d\xi d\eta \quad (21)$$

The frequency domain equivalent is the relation

$$I(f_x, f_y) = B(f_x, f_y) O(f_x, f_y) \quad (22)$$

where I , O , and B are the Fourier transforms of the lowercase quantities. The quantity B is the transfer function of the blur, and is assumed to be perfectly known.

Examination of Eq. (22) suggests an obvious approach to image restoration. Convolve the blurred image $i(x, y)$ with a kernel that provides a deblurring transfer function that is the reciprocal of the blur transfer function, i.e., having a transfer function given by $H(f_x, f_y) = B^{-1}(f_x, f_y)$. For obvious reasons, such a filter is referred to as an *inverse filter*. The restored image then is given, in the frequency domain, by

$$\begin{aligned} R(f_x, f_y) &= [B(f_x, f_y) O(f_x, f_y)] H(f_x, f_y) \\ &= [B(f_x, f_y) O(f_x, f_y)] B^{-1}(f_x, f_y) \\ &= O(f_x, f_y) \end{aligned} \quad (23)$$

Returning to the space domain we see that result of the image restoration operation is perfect recovery of the original object $o(x, y)$.

The inverse filter is an elegant mathematical solution to the restoration problem, but it lacks practicality. Many problems exist, both with the concept and with its implementation. The conceptual flaw, which is the most serious drawback, arises because the problem formulation completely neglected the inevitable presence of noise in the image $i(x, y)$. The inverse filter boosts those spatial frequency components the most that were suppressed the most by the blur. In such regions of the frequency domain there is little or no image information to be boosted, but there is always noise, which then is amplified to the point that it dominates the restored image.

Other problems arise due to the very large dynamic range required of the deblurring filter transfer function in many cases. For the above reasons, the inverse filter is never used in practice, although it is an important concept to be aware of.

The Wiener Filter The Wiener filter overcomes many of the difficulties of the inverse filter by explicitly including noise in the basic imaging model. In this case the detected image intensity is represented by

$$i(x, y) = \iint_{-\infty}^{\infty} b(x-\xi, y-\eta) o(\xi, \eta) d\xi d\eta + n(x, y) \quad (24)$$

where $o(x, y)$ and $n(x, y)$ are regarded as statistically stationary random processes. The goal of the restoration process is now to choose a restoration filter that will minimize the mean-squared error between the restored image $r(x, y)$ and the original object $o(x, y)$. The solution to this problem can be shown to be a restoration filter having a transfer function of the form

$$H(f_x, f_y) = \frac{B^*(f_x, f_y)}{|B(f_x, f_y)|^2 + \frac{P_N(f_x, f_y)}{P_O(f_x, f_y)}} \quad (25)$$

where P_N and P_O represent the power spectral densities of the respective noise and object random processes.

The Wiener filter provides a balance between uncompensated blur and residual noise in just such a way as to minimize mean-squared error. Note that at spectral locations where the object power is much greater than the noise power, the Wiener filter approaches an inverse filter, while at spectral locations where the noise power dominates, the Wiener filter behaves as a matched filter with considerable attenuation.

Coherent Optical Realization of Inverse and Wiener Filters While the inverse filter is primarily of theoretical interest, nonetheless there is much to be learned from consideration of how one might realize an approximation to such a filter. In general, the transfer function $B(f_x, f_y)$ is complex-valued, or at least has sign reversals implying 180° phase shifts at some frequencies. This implies that the inverse filter must control both the magnitude and the phase of the transmitted fields. In most cases this implies a holographic filter and possibly a second absorbing filter.

The exact blur impulse response is assumed to be known. From a blurred image of a known point source, a photographic record of the blur impulse response can be obtained. If a filter is recorded in the geometry of Fig. 13, with the blur impulse response placed in the plane labeled "input," then an interferometrically generated transparency results, one component of amplitude transmittance being proportional to the conjugate of the blur transfer function

$$t_A(\xi, \eta) \approx B^*(f_x, f_y) \exp(-i2\pi\alpha\eta) \quad (26)$$

where α is again a carrier frequency introduced by the offset reference wave. Passage of the blurred image through a coherent optical filtering system with this transfer function will correct any

frequency-domain phase shifts associated with the blur, but will not restore the magnitude of the object spectrum correctly.

To correct the spectral magnitudes we require an additional transparency to be sandwiched with the above holographic filter. This filter can be generated in a number of ways, but the easiest to understand is a method that rests on properties of the photographic process. If a photographic emulsion is exposed to an optical intensity $I(\xi, \eta)$, then over a certain dynamic range the amplitude transmittance of the resulting negative transparency will be of the form

$$t_A(\xi, \eta) = K[I(\xi, \eta)]^{-\gamma/2} \quad (27)$$

where γ is the so-called gamma of the photographic process. If the intensity to which the emulsion is exposed is simply the intensity in the Fourier transform of the blur transfer function, as obtained by optically Fourier transforming the blur spread function (for example, as in the system of Fig. 13 but with the reference wave blocked), then if a gamma equal to 2 is achieved with the photographic processing, the second transparency will have amplitude transmittance

$$t_A(\xi, \eta) = K \left| B \left(\frac{\xi}{\lambda f}, \frac{\eta}{\lambda f} \right) \right|^{-2} \quad (28)$$

If the two transparencies discussed above are now placed in contact, the overall amplitude transmittance will be the product of the two individual transmittances, and the effective filter transfer function realized by the coherent optical processor will be

$$H(f_x, f_y) = \frac{B^*(f_x, f_y)}{|B(f_x, f_y)|^2} = \frac{1}{B(f_x, f_y)} \quad (29)$$

which is precisely the transfer function of the desired inverse filter. However, in practice there will be errors in this filter due to the limited dynamic range of the photographic media.

To realize an approximation to the Wiener filter, a different recording method can be used. In this case the full holographic recording system illustrated in Fig. 13 is used, including the reference beam. However, the intensity of the reference beam is made weak compared with the peak intensity of the $|B|^2$ component. Furthermore, the recording is arranged so that the exposure falls predominantly in a range where the amplitude transmittance of the developed transparency is proportional to the logarithm of the intensity incident during exposure. Now if amplitude transmittance is proportional to the logarithm of incident exposure, the changes of amplitude transmittance, which lead to diffraction of light by the transparency, will obey

$$\Delta t_A = \beta \Delta(\log E) \approx \beta \frac{\Delta E}{E} \quad (30)$$

where ΔE represents changes in exposure, \bar{E} represents the average exposure about which the fluctuations occur, and β is a proportionality constant. Restricting attention to the proper portion of the output plane, the following identifications can be made:

$$\begin{aligned} \Delta E &\propto \hat{B}^* \exp(-i2\pi\alpha\eta) \\ \bar{E} &\propto |\hat{B}|^2 + K \end{aligned} \quad (31)$$

where $|\hat{B}|^2$ is the squared magnitude of the blur transfer function, normalized to unity at the origin, while K is the ratio between the reference beam intensity and the maximum value of $|B|^2$. Neglecting the exponential term which leads to offset from the origin in the output plane, the amplitude transmittance of the deblurring filter becomes

$$\Delta t_A = \frac{\hat{B}^*}{|\hat{B}|^2 + K} \quad (32)$$

which is precisely the form of the Wiener filter for a constant ratio K of noise power spectral density to signal power spectral density.

Thus the Wiener filter has been achieved with a single holographic filter. If the signal-to-noise ratio in the blurred image is high, then the reference beam intensity should be much less than the object beam intensity ($K \ll 1$). Bleached filters of this kind can also be made.

11.9 INCOHERENT PROCESSING OF DISCRETE SIGNALS

The previous problems examined have all involved signals and images that are continuous functions of time or space. We turn attention now to signals that are sampled or discrete functions of time or space.

Background

A continuous signal $u(t)$ is sampled at times separated by Δt yielding a set of P samples $u(k \Delta t)$, which we represent by the column vector

$$\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_P \end{bmatrix} \quad (33)$$

For discrete signals, the continuous operations associated with convolution and correlation become matrix-vector operations. Thus any linear transformation of an input signal \mathbf{u} is represented by

$$\mathbf{v} = \mathbf{M}\mathbf{u} \quad (34)$$

where \mathbf{v} is a length Q output vector containing samples of the output signal, and \mathbf{M} is a $P \times Q$ matrix

$$\mathbf{M} = \begin{bmatrix} m_{11} & m_{12} & \cdots & m_{1P} \\ m_{21} & m_{22} & \cdots & m_{2P} \\ \vdots & \vdots & \vdots & \vdots \\ m_{Q1} & m_{Q2} & \cdots & m_{QP} \end{bmatrix} \quad (35)$$

In the sections that follow we examine some of the optical approaches that have been proposed and demonstrated for this kind of operation.

The Serial Incoherent Matrix-Vector Multiplier

An important starting point is provided by the serial incoherent matrix-vector multiplier invented by Bocker¹⁹ (see also Ref. 20), and illustrated in Fig. 14. The elements of the vector \mathbf{u} are applied as current pulses, with heights proportional to the u_i , to an LED. Light from the LED floods the matrix mask, which contains $Q \times P$ cells, each with an intensity transmittance proportional to a different m_{ij} . The light transmitted by the matrix mask then falls on a two-dimensional CCD detector array, used in an unusual mode of operation. Charges are transferred horizontally along the rows of the detector array. In the first clock cycle, when the first element of the input vector is generated by the LED, the charge deposited in the first column of the detector array can be represented by a vector with elements $c_{1j} = m_{1j}u_1$. This set of charge packets is now transferred one column to the right, and

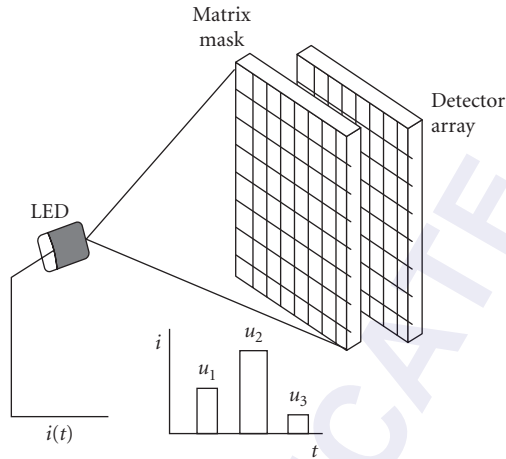


FIGURE 14 Serial matrix-vector multiplier.

the second pulse of light, proportional to u_2 is emitted. In the second column of the detector array a new charge is added to each of the existing charges, yielding a new set of charges $c_{2j} = m_{1j}u_1 + m_{2j}u_2$. After P clock cycles the column on the far right-hand side of the detector array contains a charge vector $\mathbf{c} = \mathbf{M}\mathbf{u}$, which within a proportionality constant is the desired output vector \mathbf{v} .

Thus the elements of the output vector are obtained in parallel from the right-hand column of the detector array. To compute the output vector, P cycles of the system are necessary, one for each element of the input vector. Multiplications are performed optically by passage of light through the matrix mask, while additions are performed electrically by charge addition.

The Parallel Matrix-Vector Multiplier

A fundamentally faster system for performing the matrix-vector product was discovered in 1978.²¹ The architecture of this system is shown in Fig. 15.

The elements of the vector \mathbf{u} are now entered in parallel as brightness values on an array of LEDs or laser diodes. The optics, not shown here, spread the light from each source in the vertical direction

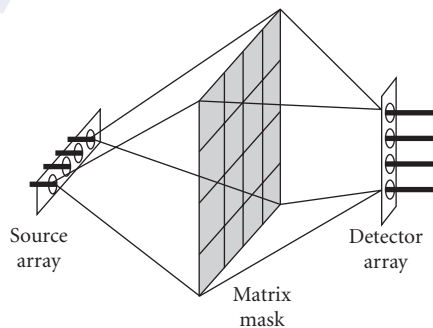


FIGURE 15 Parallel matrix-vector multiplier.

to cover the height of the matrix mask, while imaging each source onto an individual column in the horizontal direction. Passage of the light through the matrix mask multiplies the input vector, element by element, by the elements of the row vectors of the mask. The second set of optics, again not shown, adds the light transmitted by each row of the mask, placing an intensity on each element of the detector array that is proportional to the sum of the products produced by one row of the mask or, equivalently, the inner product of the input vector and a unique row vector of the matrix. In this case the detectors are of the nonintegrating type, and nearly instantaneously produce an output current proportional to an element of the output vector \mathbf{v} . In this way a series of input vectors can be flowed through the system at high speed.

In this case both the multiplications and the additions are performed optically. A different output vector can be obtained with each cycle of the system. The result is a fundamentally faster system.

Systems of this type have had a broad impact on optical signal processing, with applications ranging from photonic switching²² to neural networks.²³

The Outer Product Processor

Another fundamentally different architecture is represented by the outer product processor,²⁴ shown in Fig. 16.

The goal in this case is to calculate the outer product \mathbf{C} of two matrices \mathbf{A} and \mathbf{B} . Illustrating with a simple 3×3 example, the outer product is defined by the equation

$$\mathbf{C} = \begin{bmatrix} a_{11} \\ a_{21} \\ a_{31} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} & b_{13} \end{bmatrix} + \begin{bmatrix} a_{12} \\ a_{22} \\ a_{32} \end{bmatrix} \begin{bmatrix} b_{21} & b_{22} & b_{23} \end{bmatrix} + \begin{bmatrix} a_{13} \\ a_{23} \\ a_{33} \end{bmatrix} \begin{bmatrix} b_{31} & b_{32} & b_{33} \end{bmatrix} \quad (36)$$

The system of Fig. 16 accomplishes this operation by use of two Bragg cell arrays oriented in orthogonal directions, and a time-integrating two-dimensional detector array. A column of \mathbf{A} is entered in parallel on the first Bragg cell array, and a row of \mathbf{B} on the second. The first box labeled "optics" images one array onto the other (with appropriate spatial filtering as needed to convert phase to intensity). The second box labeled "optics" images that product onto the detector array. In one cycle of the system, one of the outer products in the summation of Eq. (36) is accumulated on the elements of the detector array. In this example, three cycles are necessary, with addition of charge at the detector, to accumulate the full outer product of the matrices. More generally, if \mathbf{A} is $P \times Q$ (i.e., P rows and Q columns) and \mathbf{B} is $Q \times P$, then the detector array should be $P \times P$ in size, and Q cycles are required to obtain the full outer product.

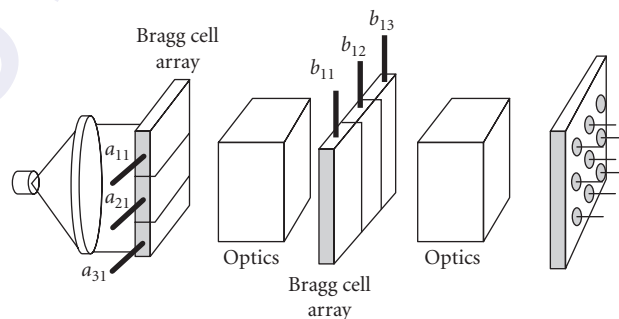


FIGURE 16 Outer product processor.

Other Discrete Processing Systems

A multitude of other discrete processing systems have been proposed throughout the 1980s and 1990s. Worth special mention here are the optical systolic matrix-vector processor of Caulfield et al.²⁵ and the SAOBIC processor of Guilfoyle.²⁶ We refer the reader to the original references for details.

11.10 CONCLUDING REMARKS

Analog optical signal and image processing were strong areas of research during three decades, the 1960s through the 1980s. Many ingenious systems were devised, each motivated by one or more applications. With some exceptions, these systems seldom survived for the particular application they were conceived for, but often they led to new applications not envisioned by their inventors. The majority of applications of this technology have been to defense-related problems. Research emphasis has shifted away from analog signal processing, as described above, towards the application of optics to providing interconnects between and within digital computers. However, the intellectual base formed by previous analog processing experience continues to strongly influence work in other, more modern disciplines, including integrated optics, modern microscopy, coherence tomography, ultrafast optical pulses, and digital image processing.

11.11 REFERENCES

1. F. Zernike, "Das Phasenkontrastverfahren bei der Mikroskopischen beobachtung," *A. Tech. Phys.* **16**:454 (1935).
2. E. Abbe, "Beitrage zur theorie des mikroskops und der mikroskopischen wahrnehmung," *Archiv. Mikroskopische Anat.* **9**:413–468 (1873).
3. A. B. Porter, "On the Diffraction Theory of Microscope Vision," *Phil. Mag.* **11**:154 (1906).
4. A. Marechal and P. Croce, "Un filtre de frequences spatiales pour l'amelioration du contraste des images optiques," *C.R. Acad. Sci.* **127**:607 (1953).
5. E. L. O'Neill, "Spatial Filtering in Optics," *IRE Trans. on Info. Theory* **IT-2**:56–65 (1956).
6. J. W. Goodman, *Introduction to Fourier Optics*, 3rd ed, Roberts & Company Publishers, Greenwood Village, CO, 2005.
7. L. J. Cutrona, et al., "On the Application of Coherent Optical Processing Techniques to Synthetic-Aperture Radar," *Proc. IEEE* **54**:1026–1032 (1966).
8. A. Kozma, E. N. Leith, and N. G. Massey, "Tilted Plane Optical Processor," *Appl. Opt.* **11**:1766–1777 (1972).
9. L. Slobodin, "Optical Correlation Techniques," *Proc. IEEE* **51**:1782 (1963).
10. M. Arm, L. Lambert, and I. Weissman, "Optical Correlation Technique for Radar Pulse Compression," *Proc. IEEE* **52**:842 (1964).
11. R. M. Montgomery, "Acousto-Optical Signal Processing System," U.S. Patent 3,634,749, 1972.
12. R. A. Sprague and C. L. Koliopoulos, "Time Integrating Acousto-Optic Correlator," *Appl. Opt.* **15**:89–92 (1975).
13. A. Vander Lugt, *Optical Signal Processing*, John Wiley & Sons, New York, 1992.
14. G. L. Turin, "An Introduction to Matched Filters," *IRE Trans. on Info. Theo.* **IT-6**:311 (1960).
15. A. B. Vander Lugt, "Signal Detection by Complex Spatial Filtering," *IEEE Trans. on Info.Theo.* **IT-10**:139–145 (1964).
16. R. J. Collier, C. B. Burkhart, and L. H. Lin, *Optical Holography*, Academic Press, New York, 1971.
17. D. Casasent and D. Psaltis, "Position, Rotation and Scale-Invariant Optical Correlation," *Appl. Opt.* **15**:1795–1799 (1976).
18. D. Casasent and D. Psaltis, "New Optical Transforms for Pattern Recognition," *Proc. IEEE* **65**:770–784 (1977).

19. R. P. Bocker, "Matrix Multiplication Using Incoherent Optical Techniques," *Appl. Opt.* **13**:1670–1676 (1974).
20. K. Bromley, "An Incoherent Optical Correlator," *Opt. Act.* **21**:35–41 (1974).
21. J. W. Goodman, A. R. Dias, and L. M. Woody, "Fully Parallel, High-Speed Incoherent Optical Method for Performing the Discrete Fourier Transform," *Opt. Lett.* **2**:1–3 (1978).
22. A. R. Dias, R. F. Kalman, J. W. Goodman, and A. A. Sawchuk, "Fiber-Optic Crossbar Switch with Broadcast Capability," *Proc. SPIE* **825**:170–177 (1987).
23. N. H. Farhat, D. Psaltris, A. Prata, and E. Paek, "Optical Implementation of the Hopfield Model," *Appl. Opt.* **24**:1469–1475 (1985).
24. R. A. Athale and W. C. Collins, "Optical Matrix-Matrix Multiplier Based on Outer Product Decomposition," *Appl. Opt.* **21**:2089–2090 (1982).
25. H. J. Caulfield, et al., "Optical Implementation of Systolic Array Processing," *Opt. Comm.* **40**:86–90 (1982).
26. P. S. Guilfoyle, "Systolic Acousto-Optic Binary Convolver," *Opt. Eng.* **23**:20–25 (1984).

This page intentionally left blank.

DO NOT DUPLICATE

PART

3

POLARIZED LIGHT

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

Jean M. Bennett[†]

*Research Department
 Michelson Laboratory
 Naval Air Warfare Center
 China Lake, California*

12.1 GLOSSARY

c	velocity of light
d	thickness
E	electric field
\mathbf{k}	wave vector ($k = 2\pi/\lambda$)
k	extinction coefficient
m	number of reflections
N	retardation per wavelength
n	real refractive index
\tilde{n}	complex refractive index
$\hat{\mathbf{n}}$	unit normal vector
P	degree of polarization
p	parallel polarization
R	intensity reflection coefficient
r	amplitude reflection coefficient
\mathbf{r}	position vector
s	senkrecht or perpendicular polarization
t	amplitude transmission coefficient
t	time
z	cartesian coordinate
$\alpha, \beta, a, b, c, d$	intermediate parameters
α	absorption coefficient

*This material was originally prepared under the auspices of the U.S. government and is not subject to copyright.

[†]Deceased.

γ	$2\pi nd \cos \theta / \lambda$
δ	phase angle
ϵ	dielectric constant
η	effective refractive index
θ_B	Brewster angle
θ	angle
κ	absorption index
λ	wavelength
ρ	extinction ratio
σ	conductivity
ω	radian or angular frequency
∇	laplacian operator
0	first medium
1	second medium

The material on polarization is abridged from the much more complete treatment by Bennett and Bennett.¹ Information on polarizers is found in Chap. 13, "Polarizers," in this volume.

12.2 BASIC CONCEPTS AND CONVENTIONS

Optical polarization was discovered by E. L. Malus in 1808. A major triumph of nineteenth- and early twentieth-century theoretical physics was the development of electromagnetic theory and the demonstration that optical polarization is completely described by it. This theory is phenomenological in that instead of trying to explain why materials have certain fundamental characteristics, it concentrates on the resulting properties which any material with those characteristics will display. In the optical case, the polarization and all other optical properties of a material are determined by two or more phenomenological parameters called *optical constants*. Electromagnetic theory has little or nothing to say about *why* a material should have these particular optical constants or *how* they are related to its atomic character. This problem has been extensively investigated in twentieth-century solid-state physics and is still only partially understood. It is clear, however, that the optical constants are a function not only of the atomic nature of the material, i.e., its position in the periodic table, but are also quite sensitive to how it is prepared. Perhaps *optical parameters* would be a better term than optical constants. Nevertheless, the concept of optical constants is an extremely useful one and makes it possible to predict quantitatively the optical behavior of a material and, under certain conditions, to relate this behavior to nonoptical parameters.

Since the optical constants are so fundamental, differences in their definition are particularly unfortunate. The most damaging of these differences arise from an ambiguity in the initial derivation. Maxwell's equations, which form the basis of electromagnetic theory, result in the wave equation, which in mks units is

$$\nabla^2 \mathbf{E} = \frac{\epsilon}{c^2} \frac{\partial^2 \mathbf{E}}{\partial t^2} + \frac{4\pi\sigma}{c^2} \frac{\partial \mathbf{E}}{\partial t} \quad (1)$$

where ∇^2 = laplacian operator

\mathbf{E} = electric field vector of traveling wave

t = time

c = velocity of light

σ = conductivity of material at frequency of wave motion

ϵ = dielectric constant of material at frequency of wave motion

A solution to this equation is

$$\mathbf{E} = \mathbf{E}_0 \exp [i(\omega t + \delta)] \exp (-i\mathbf{k} \cdot \mathbf{r}) \exp \left(-\frac{\alpha z}{2} \right) \quad (2)$$

where \mathbf{E}_0 = amplitude of wave
 ω = angular frequency of wave
 δ = phase vector
 \mathbf{k} = wave vector
 \mathbf{r} = position vector
 z = direction wave is traveling
 α = absorption coefficient

The wave vector \mathbf{k} is assumed to be real and equal to $(2\pi/\lambda_m)\hat{\mathbf{n}}$, where λ_m is the wavelength in the medium in which the wave is traveling and $\hat{\mathbf{n}}$ is a unit vector in the k direction.* Equation (2) can also be written in terms of \tilde{n} , the complex index of refraction, defined as

$$\tilde{n} = n - ik \quad (3)$$

where n is the index of refraction and k the extinction coefficient. In this form, Eq. (2) is

$$E = E_0 \exp \left[i\omega \left(t - \frac{\tilde{n}z}{c} \right) \right] \quad (4)$$

when $\delta = 0$. By setting the imaginary part of the exponent equal to zero one obtains

$$z = \frac{c}{n}t \quad (5)$$

To show that Eq. (4) represents a wave traveling in the positive z direction with phase velocity c/n , we note that the phase ϕ of the wave in Eq. (4) is $\omega t - (\omega\tilde{n}z)/c = \phi$. For a wave propagating with a constant phase, $d\phi = 0$, so that $\omega dt - (\omega\tilde{n}/c) dz = d\phi = 0$, and hence the phase velocity $v_p = dz/dt = c/n$.² The amplitude of the wave at z is, from Eq. (4),

$$|E| = E_0 e^{-2\pi k z / \lambda} \quad (6)$$

where λ is the wavelength in vacuum. The wave is thus exponentially damped, and the amplitude penetration depth, or distance below an interface at which the *amplitude* of the wave falls to $1/e$ times its initial value, is $z = \lambda/2\pi k$. The absorption coefficient α , or the reciprocal of the distance in which the *intensity* of the wave falls to $1/e$ times its initial value, is

$$\alpha = \frac{4\pi k}{\lambda} \quad (7)$$

This development follows that commonly given by those working at optical or radio frequencies. The confusion in the definition of the optical constants arises because an equally valid solution to Eq. (1) is

$$E' = E_0 \exp \left[-i\omega \left(t - \frac{\tilde{n}'z}{c} \right) \right] \quad (8)$$

*Frequently the wave vector is taken to be complex, that is, $\tilde{\mathbf{k}} = (2\pi/\lambda_m - i\alpha/2)\mathbf{n}$, and Eq. (2) is written $\mathbf{E} = \mathbf{E}_0 \exp [i(\omega t + \delta)] \exp (-i\tilde{\mathbf{k}} \cdot \mathbf{r})$.

which also represents an exponentially damped wave traveling in the $+z$ direction *provided that the complex index of refraction is defined to be*

$$\tilde{n}' = n + ik \quad (9)$$

where the primes indicate the alternative solution. When the wave equation arises in quantum mechanics, the solution chosen is generally the negative exponential, i.e., Eq. (8) rather than Eq. (4). Solid-state physicists working in optics thus often define the complex index of refraction as the form given in Eq. (9) rather than that in Eq. (3). Equally valid, self-consistent theories can be built up using either definition, and as long as only intensities are considered, the resulting expressions are identical. However, when phase differences are calculated, the two conventions usually lead to contradictory results. Even worse, an author who is not extremely careful may not consistently follow either convention, and the result may be pure nonsense. Some well-known books might be cited in which the authors are not even consistent from chapter to chapter.

There are several other cases in optics in which alternative conventions are possible and both are found in the literature. Among these, the most distressing are the use of a left-handed rather than a right-handed coordinate system, which makes the p and s components of polarized light have the same phase change at normal incidence (see Sec. 12.3), and defining the optical constants so that they depend on the angle of incidence, which makes the angle of refraction given by Snell's law real for an absorbing medium. There are many advantages to be gained by using a single set of conventions in electromagnetic theory. In any event, an author should *clearly* state the conventions being used and then *stay with them*.

Finally, the complex index of refraction is sometimes written

$$\tilde{n} = n(1 - i\kappa) \quad (10)$$

In this formulation the symbol κ is almost universally used instead of k , which is reserved for the imaginary part of the refractive index. Although k is more directly related to the absorption coefficient α than κ [see Eq. (7)] and usually makes the resulting expressions slightly simpler, in areas such as attenuated total reflection, the use of κ results in a simplification. To avoid confusion between k and κ , if Eq. (10) is used, κ could be called the *absorption index* to distinguish it from the extinction coefficient k , and the absorption coefficient α .

12.3 FRESNEL EQUATIONS

The Fresnel equations are expressions for the reflection and transmission coefficients of light at nonnormal incidence. In deriving these equations, the coordinate system assumed determines the signs in the equations and therefore the phase changes on reflection of the p and s components. In accordance with the Muller convention,³ we shall assume that the coordinate system is as shown in Fig. 1. In this system, the angle of incidence is θ_0 , and the angle of refraction is θ_1 . The s component of polarization is the plane of vibration of the E wave which is perpendicular to the plane of the paper, and the p component is the plane of vibration which is in the plane of the paper.* (The plane of incidence is in the plane of the paper.) The positive directions for the vibrations are indicated in Fig. 1 by the dots for E_s , E'_s , and E''_s and by the arrows for the corresponding p components. Note that the positive direction for E''_p is as shown in the figure because of the *mirror-image effect*.

*Unfortunately, when Malus discovered that light reflected at a certain angle from glass is, as he said, "polarized," he defined the plane of polarization" of the reflected light as the plane of incidence. Since the reflected light in this case has its E vector perpendicular to the plane of incidence, the "plane of polarization" is perpendicular to the plane in which the E vector vibrates. This nomenclature causes considerable confusion and has been partially resolved in modern terminology by discussing the *plane of vibration* of the E vector and avoiding, insofar as possible, the term *plane of polarization*. In this chapter, when specifying the direction in which light is polarized, we shall give the direction of vibration, *not* the direction of polarization.

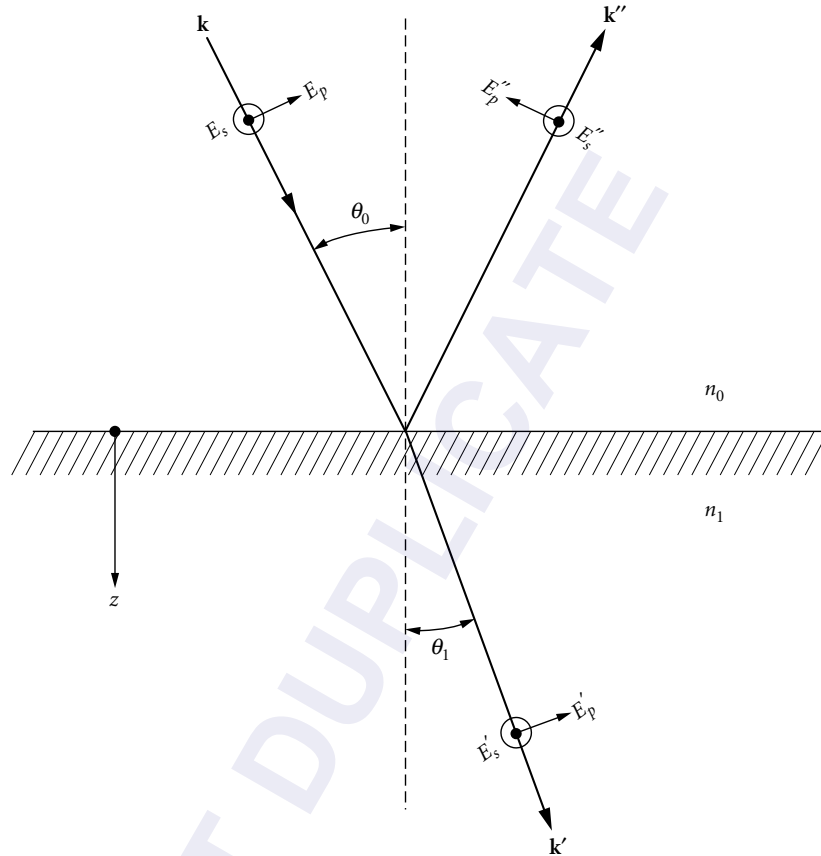


FIGURE 1 Coordinate system for measuring the E vectors of a plane wave reflected and refracted at a boundary between a medium of refractive index n_0 and a medium of refractive index n_1 (may be absorbing). The positive direction for the coordinates of the E_s , E_s' , and E_s'' components is out of the paper, and that for the coordinates of the E_p components is in the plane of the paper, as indicated by the arrows. The wave vector \mathbf{k} , the direction the wave is traveling z , and angles of incidence and refraction θ_0 and θ_1 are also shown. (Modified from Mutter.³)

By convention, one always looks *against the direction of propagation of the wave* so that the positive direction of E_p is to the right and the positive direction of E_p'' is also to the right. The positive directions of the reflected E vectors are not the same as the actual directions of the reflected E vectors. These latter directions will depend on the refractive index of the material and may be either positive or negative. For example, if $n_1 > n_0$, at normal incidence E_s'' will be in the negative direction and E_p'' will be in the positive direction. Thus we say that there is a phase change on reflection of 180° for the s wave and a phase change of 0° for the p wave.

With this coordinate system, the Fresnel amplitude reflection coefficients for a single interface, obtained from Eq. (4) by setting up and solving the boundary-value problem, can be written

$$\frac{E_s''}{E_s} \equiv r_s = \frac{n_0 \cos \theta_0 - n_1 \cos \theta_1}{n_0 \cos \theta_0 + n_1 \cos \theta_1} \quad (11)$$

and

$$\frac{E_p''}{E_p} \equiv r_p = \frac{n_1 \cos \theta_0 - n_0 \cos \theta_1}{n_1 \cos \theta_0 + n_0 \cos \theta_1} \quad (12)$$

The amplitude transmission coefficients are

$$\frac{E_s'}{E_s} \equiv t_s = \frac{2n_0 \cos \theta_0}{n_0 \cos \theta_0 + n_1 \cos \theta_1} \quad (13)$$

and

$$\frac{E_p'}{E_p} \equiv t_p = \frac{2n_0 \cos \theta_0}{n_1 \cos \theta_0 + n_0 \cos \theta_1} \quad (14)$$

Other forms of the Fresnel amplitude reflection and transmission coefficients containing only the angles of incidence and refraction are somewhat more convenient. These relations can be derived using Snell's law

$$\frac{\sin \theta_0}{\sin \theta_1} = \frac{n_1}{n_0} \quad (15)$$

to eliminate n_0 and n_1 from Eqs. (1) to (14):

$$r_s = \frac{-\sin(\theta_0 - \theta_1)}{\sin(\theta_0 + \theta_1)} \quad (16)$$

$$r_p = \frac{\tan(\theta_0 - \theta_1)}{\tan(\theta_0 + \theta_1)} \quad (17)$$

$$t_s = \frac{2 \sin \theta_1 \cos \theta_0}{\sin(\theta_0 + \theta_1)} \quad (18)$$

$$t_p = \frac{2 \sin \theta_1 \cos \theta_1}{\sin(\theta_0 + \theta_1) \cos(\theta_0 - \theta_1)} \quad (19)$$

For nonabsorbing materials the intensity reflection coefficients R_s and R_p are simply the squares of Eqs. (16) and (17):

$$R_s = r_s^2 = \frac{\sin^2(\theta_0 - \theta_1)}{\sin^2(\theta_0 + \theta_1)} \quad (20)$$

$$R_p = r_p^2 = \frac{\tan^2(\theta_0 - \theta_1)}{\tan^2(\theta_0 + \theta_1)} \quad (21)$$

and, at normal incidence,

$$R_s = R_p = \frac{(n_0 - n_1)^2}{(n_0 + n_1)^2} \quad (22)$$

from Eqs. (11) and (12). In the lower part of Fig. 2, R_s and R_p are given as a function of angle of incidence for various values of the refractive-index ratio n_1/n_0 with k_1 for the material equal to zero.

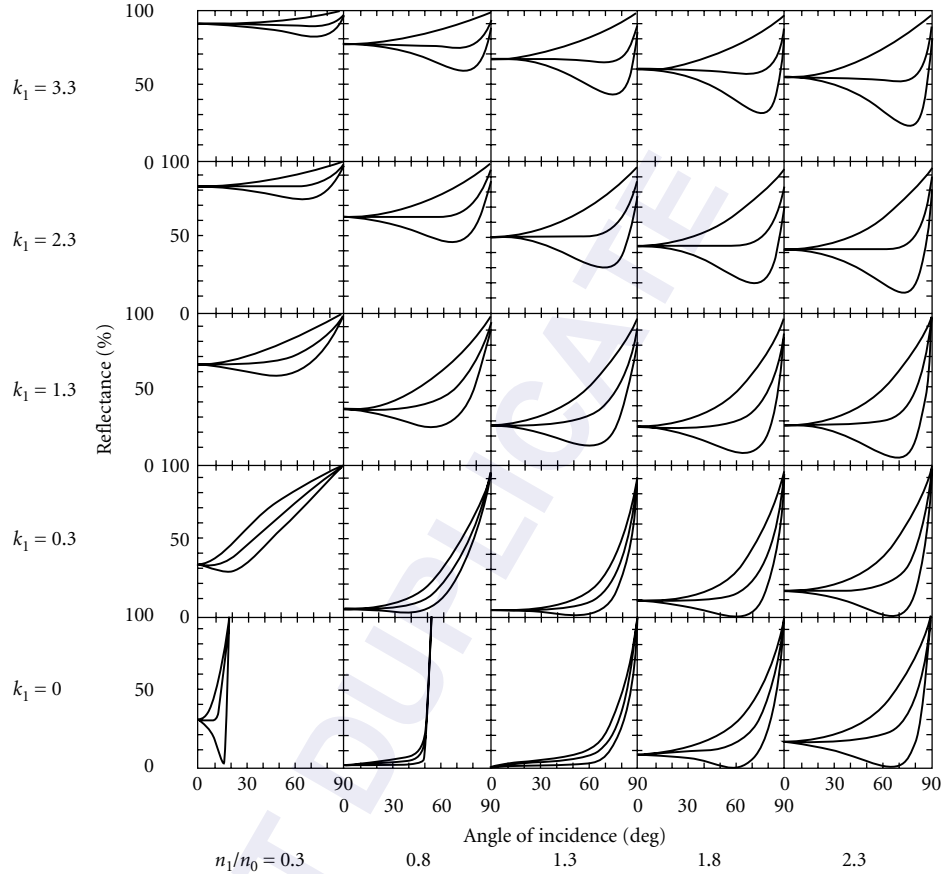


FIGURE 2 R_s (upper curves), R_p (lower curves), and $R_{av} = (R_s + R_p)/2$ as a function of angle of incidence for various values of the refractive-index ratio n_1/n_0 and k_1 . The incident medium, having refractive index n_0 , is assumed to be nonabsorbing. (Modified from Hunter.⁴)

The curves for $n_1/n_0 = 1.3, 1.8,$ and 2.3 show that the normal-incidence reflectance increases as n_1 increases. The curves for $n_1/n_0 = 0.3$ and 0.8 and $k_1 = 0$ have no physical significance as long as the incident medium is air. However, they are representative of *internal reflections* in materials of refractive index $n_0 = 3.33$ and 1.25 , respectively when the *other* medium is air ($n_1 = 1$).

The intensity transmission coefficients T_s and T_p are obtained from the Poynting vector and for nonabsorbing materials are

$$T_s = 1 - R_s = \frac{n_1 \cos \theta_1}{n_0 \cos \theta_0} t_s^2 = \frac{4n_0 n_1 \cos \theta_0 \cos \theta_1}{(n_0 \cos \theta_0 + n_1 \cos \theta_1)^2} = \frac{\sin 2\theta_0 \sin 2\theta_1}{\sin^2(\theta_0 + \theta_1)} \quad (23)$$

$$\begin{aligned} T_p = 1 - R_p &= \frac{n_1 \cos \theta_1}{n_0 \cos \theta_0} t_p^2 = \frac{4n_0 n_1 \cos \theta_0 \cos \theta_1}{(n_1 \cos \theta_0 + n_0 \cos \theta_1)^2} \\ &= \frac{\sin 2\theta_0 \sin 2\theta_1}{\sin^2(\theta_0 + \theta_1) \cos^2(\theta_0 - \theta_1)} \end{aligned} \quad (24)$$

These coefficients are for light passing through a single boundary and hence are of limited usefulness. In actual cases, the light is transmitted through a slab of material where there are two boundaries generally multiple reflections within the material, and sometimes interference effects when the boundaries are smooth and plane-parallel.

The intensity transmission coefficient T_{sample} for a slab of transparent material in air is given by the well-known Airy equation⁵ when the sample has smooth, plane-parallel sides and coherent multiple reflections occur within it:

$$T_{\text{sample}} = \frac{1}{1 + [4R_{s,p}/(1 - R_{s,p})^2] \sin^2 \gamma} \quad (25)$$

where

$$\gamma = \frac{2\pi n_1 d \cos \theta_1}{\lambda} \quad (26)$$

The values of R_s and R_p can be determined from Eqs. (20) to (22); d is the sample thickness, λ the wavelength, n_1 the refractive index of the material, and θ_1 the angle of refraction. Equation (25) holds for all angles of incidence including the Brewster angle, where $R_p = 0$ [see Eq. (48)]. The Airy equation predicts that at a given angle of incidence the transmission of the sample will vary from a maximum value of 1 to a minimum value of $(1 - R_{s,p})^2 / (1 + R_{s,p})^2$ as the wavelength or the thickness is changed. If the sample is very thick, the oscillations in the transmittance will occur at wavelengths very close together and hence will be unresolved. A complete oscillation occurs every time γ changes by π , so that the wavelength interval $\Delta\lambda$ between oscillations is

$$\Delta\lambda \approx \frac{\lambda^2}{2n_1 d \cos \theta_1} \quad (27)$$

As an example, a sample 1 mm thick with an index of 1.5 at 5000 Å will have transmission maxima separated by 0.83 Å when measured at normal incidence ($\cos \theta_1 = 1$). These maxima would not be resolved by most commercial spectrophotometers. In such a case, one would be measuring the average transmission $T_{\text{sample,av}}$:

$$T_{\text{sample,av}} = \frac{1 - R_{s,p}}{1 + R_{s,p}} \quad (28)$$

For nonabsorbing materials, this is the same value as that which would be obtained if the multiply reflected beams did not coherently interfere within the sample. If the sample is wedge-shaped, so that no multiply reflected beams contribute to the measured transmittance, T_{sample} is simply T_s^2 or T_p^2 and can be calculated from Eq. (23) or (24).

When the material is absorbing, i.e., has a complex refractive index, it is not so easy to calculate the reflectance and transmittance since the angle of refraction is complex. However, Snell's law [Eq. (15)] and Fresnel's equations (11) and (12) are sometimes used with complex values of n_1 and θ_1 . The resulting amplitude reflection coefficients are written

$$r_s = |r_s| e^{i\delta_s} \quad (29)$$

and

$$r_p = |r_p| e^{i\delta_p} \quad (30)$$

where $|r_s|$ and $|r_p|$ are the magnitudes of the reflectances and δ_s and δ_p are the phase changes on reflection. The intensity reflection coefficients are

$$R_{s,p} = r_{s,p} r_{s,p}^* \quad (31)$$

An alternative approach is to use the method of effective indexes to calculate R_s and R_p . In the medium of incidence, which is assumed to be nonabsorbing, the effective indexes η_{0s} and η_{0p} for the s and p components are

$$\eta_{0s} = n_0 \cos \theta_0 \quad (32)$$

$$\eta_{0p} = \frac{n_0}{\cos \theta_0} \quad (33)$$

where n_0 generally equals 1 for air. In the absorbing material both η 's are complex and can be written, according to the Bernings,^{6,7}

$$\tilde{\eta}_{1s} = \tilde{n}_1 \cos \theta_1 \quad (34)$$

$$\tilde{\eta}_{1p} = \frac{\tilde{n}_1}{\cos \theta_1} \quad (35)$$

where $\tilde{n}_1 = n_1 - ik_1$ is the complex refractive index of the material, and

$$\cos \theta_1 = \left[\frac{(\alpha_1^2 + \beta_1^2)^{1/2} + \alpha_1}{2} \right]^{1/2} - i \left[\frac{(\alpha_1^2 + \beta_1^2)^{1/2} - \alpha_1}{2} \right]^{1/2} \quad (36)$$

$$\alpha_1 = 1 + \left(\frac{n_0 \sin \theta_0}{n_1^2 + k_1^2} \right)^2 (k_1^2 - n_1^2) \quad (37)$$

and

$$\beta_1 = -2n_1 k_1 \left(\frac{n_0 \sin \theta_0}{n_1^2 + k_1^2} \right)^2 \quad (38)$$

Abelès' method⁸ also uses effective indexes for the absorbing material, but they are calculated differently:

$$\tilde{\eta}_{1s} = a - ib \quad (39)$$

$$\tilde{\eta}_{1p} = c - id \quad (40)$$

where

$$a^2 - b^2 = n_1^2 - k_1^2 - n_0^2 \sin^2 \theta_0 \quad (41)$$

$$ab = n_1 k_1 \quad (42)$$

$$c = a \left(1 + \frac{n_0^2 \sin^2 \theta_0}{a^2 + b^2} \right) \quad (43)$$

$$d = b \left(1 - \frac{n_0^2 \sin^2 \theta_0}{a^2 + b^2} \right) \quad (44)$$

In both methods employing effective indexes, the amplitude reflection coefficients are

$$r_s = \frac{\eta_{0s} - \eta_{1s}}{\eta_{0s} + \eta_{1s}} \quad (45)$$

$$r_p = \frac{\eta_{1p} - \eta_{0p}}{\eta_{1p} + \eta_{0p}} \quad (46)$$

which are equivalent to Eqs. (29) and (30) and reduce to Eqs. (11) and (12) when $k_1 = 0$. The intensity reflection coefficients are given by Eq. (31), as before. At normal incidence,

$$R_s = R_p = \frac{(n_0 - n_1)^2 + k_1^2}{(n_0 + n_1)^2 + k_1^2} \quad (47)$$

Values of R_s and R_p are plotted as a function of angle of incidence in Fig. 2 for various values of n_1 and k_1 . (The incident medium is assumed to be air with $n_0 = 1$ unless otherwise noted.) As n_1 increases with $k_1 > 0$ held constant, the magnitudes of R_s and R_p at normal incidence both decrease. As k_1 increases with n_1 held constant, the magnitudes of R_s and R_p at normal incidence both increase. Tables of R_s and R_p for various values of n_1 and k_1 are given for angles of incidence from 0 to 85° by Holl.⁹

The absolute phase changes on reflection δ_s and δ_p are also of interest in problems involving polarization. When the material is nonabsorbing, the phase changes can be determined from the amplitude reflection coefficients, Eqs. (11) and (12); when $\theta_0 = 0$ and $n_1 > n_0$, $\delta_s = 180^\circ$ and $\delta_p = 360^\circ$.^{*} This is an apparent contradiction since at normal incidence the s and p components should be indistinguishable. However, the problem is resolved by recalling that by convention we are always looking against the direction of propagation of the light (see Fig. 1). To avoid complications, the phase change on reflection at normal incidence (often defined as β) is identified with δ_s .

For a dielectric, if $n_0 < n_1$, δ_s remains 180° for all angles of incidence from 0 to 90°, as can be seen from the numerator of Eq. (11). However, there is an abrupt discontinuity in δ_p , as can be seen from Eq. (12). If $n_0 < n_1$, $\delta_p = 360^\circ$ at normal incidence and at larger angles for which the numerator of Eq. (12) is positive. Since $\cos \theta_0$ becomes increasingly less than $\cos \theta_1$ as θ_0 increases, and since $n_1 > n_0$, there will be an angle for which $n_1 \cos \theta_0 = n_0 \cos \theta_1$. At this angle δ_p undergoes an abrupt change from 360° to 180°, and it remains 180° for larger angles of incidence. At the transition value of θ_0 , which is called the *Brewster angle* θ_B since $R_p = 0$,

$$\tan \theta_B = \frac{n_1}{n_2} \quad (48)$$

(This angle is also called the *polarizing angle* since $\theta_0 + \theta_1 = 90^\circ$.)

The phase changes δ_s and δ_p are not simply 360° or 180° for an absorbing material. At normal incidence it follows from Eq. (45) that

$$\tan \delta_s = \frac{2n_0 k_1}{n_0^2 - n_1^2 - k_1^2} \quad (49)$$

so that $\delta_s = 180^\circ$ only if $k_1 = 0$. As before, $\delta_p = \delta_s + 180^\circ$, as seen by comparing Eqs. (45) and (46). At nonnormal incidence

$$\tan \delta_s = \frac{2\eta_{0s} b}{n_{0s}^2 - a^2 - b^2} \quad (50)$$

^{*}Since 360° and 0° are indistinguishable, many optics books state that $\delta_p = 0^\circ$ for dielectrics at normal incidence, but this makes the ellipsometric parameter $\Delta = \delta_p - \delta_s < 0$, which is incompatible with ellipsometric conventions—see Chap. 16, “Ellipsometry.”

and

$$\tan \delta_p = \frac{-2\eta_{0p}d}{c^2 + d^2 - \eta_{0p}^2} \quad (51)$$

where the relations for a , b , c , and d have been given in Eqs. (41) to (44). The following relations between these quantities may also prove helpful:

$$a^2 + b^2 = [(n_1^2 - k_1^2 - n_0^2 \sin^2 \theta_0)^2 + 4n_1^2 k_1^2]^{1/2} \quad (52)$$

$$c^2 + d^2 = \frac{(n_1^2 + k_1^2)^2}{a^2 + b^2} \quad (53)$$

$$b^2 = \frac{n_1^2 - k_1^2 - n_0^2 \sin^2 \theta_0}{2} + \frac{a^2 + b^2}{2} \quad (54)$$

Figure 3 shows how δ_s and δ_p change as a function of angle of incidence for an absorbing material. At normal incidence they are 180° apart because of the mirror-image effect, mentioned previously. As the angle of incidence increases, δ_p approaches δ_s , and at the *principal angle* $\bar{\theta}$ the two quantities differ by only 90° . At grazing incidence they coincide.

The reflectance R_p does not reach zero for an absorbing material as it does for a dielectric, but the angle for which it is a minimum is called the *pseudo-Brewster angle* θ'_B . Two other angles closely associated with the pseudo-Brewster are also of interest. The angle for which the ratio R_p/R_s is a minimum is sometimes called the *second Brewster angle*. It is generally only slightly larger than θ'_B . The *principal angle* $\bar{\theta}$, at which $\delta_p - \delta_s = 90^\circ$, is always larger than the second Brewster angle and θ'_B . For most metals θ'_B and $\bar{\theta}$ are only a fraction of a degree apart but it is possible for them to differ by as much as 45° .⁹ There is no polarizing angle as such for an absorbing material because the angle of refraction is complex.

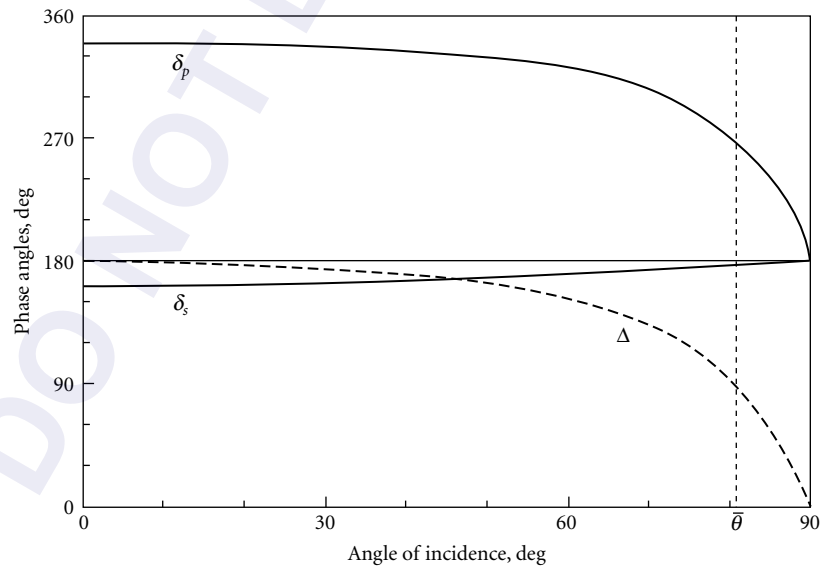


FIGURE 3 Phase changes on reflection δ_s and δ_p and phase difference $\Delta = \delta_p - \delta_s$ as a function of angle of incidence for an absorbing material. The principal angle, for which $\Delta = 90^\circ$, is also shown. (Bennett and Bennett.¹⁰)

12.4 BASIC RELATIONS FOR POLARIZERS

A linear* polarizer is anything which when placed in an incident unpolarized beam produces a beam of light whose electric vector is vibrating primarily in one plane, with only a small component vibrating in the plane perpendicular to it. If a polarizer is placed in a plane-polarized beam and is rotated about an axis parallel to the beam direction, the transmittance T will vary between a maximum value T_1 and a minimum value T_2 according to the law

$$T = (T_1 - T_2) \cos^2 \theta + T_2 \quad (55)$$

Although the quantities T_1 and T_2 are called the *principal transmittances*, in general $T_1 \gg T_2$; θ is the angle between the plane of the principal transmittance T_1 and the plane of vibration (of the electric vector) of the incident beam. If the polarizer is placed in a beam of unpolarized light, its transmittance is

$$T = \frac{1}{2}(T_1 + T_2) \quad (56)$$

so that a perfect polarizer would transmit only 50 percent of an incident unpolarized beam.†

When two identical polarizers are placed in an unpolarized beam, the resulting transmittance will be

$$T_{\parallel} = \frac{1}{2}(T_1^2 + T_2^2) \quad (57)$$

when their principal transmittance directions are parallel and will be

$$T_{\parallel} = T_1 T_2 \quad (58)$$

when they are perpendicular. In general, if the directions of principal transmittance are inclined at an angle θ to each other, the transmittance of the pair will be

$$T_{\theta} = \frac{1}{2}(T_1^2 + T_2^2) \cos^2 \theta + T_1 T_2 \sin^2 \theta \quad (59)$$

The polarizing properties of a polarizer are generally defined in terms of its *degree of polarization* $P^{\ddagger, \S}$

$$P = \frac{T_1 - T_2}{T_1 + T_2} \quad (60)$$

or its *extinction ratio* ρ_p

$$\rho_p = \frac{T_2}{T_1} \quad (61)$$

When one deals with nonnormal-incidence reflection polarizers, one generally writes P and ρ_p in terms of R_p and R_s , the reflectances of light polarized parallel and perpendicular to the plane of incidence, respectively. As will be shown in Sec. 12.5, R_s can be equated to T_1 and R_p to T_2 , so that

*Circular polarizers are discussed in Sec. 12.7.

†Jones¹¹ has pointed out that a perfect polarizer can transmit more than 50 percent of an incident unpolarized beam under certain conditions.

‡Bird and Shurcliff² distinguish between *degree of polarization*, which is a constant of the light beam, and *polarizance*, which is a constant of the polarizer. The polarizance is defined as being equal to the degree of polarization the polarizer produces in an incident monochromatic beam that is unpolarized. In practice, incident beams are often slightly polarized, so that the polarizance values differ slightly from the ideal degree of polarization. Other authors have not followed this distinction.

§Authors dealing with topics such as scattering from aerosols sometimes define *degree of polarization* (of the scattered light) in terms of the Stokes vectors (Sec. 12.8) as $P = (S_1^2 + S_2^2 + S_3^2)^{1/2} / S_0$.

Eqs. (60) and (61) become $P = (R_s - R_p)/(R_s + R_p)$ and $\rho_p = R_p/R_s$. If either ρ_p or P is known, the other can be deduced since

$$P = \frac{1 - \rho_p}{1 + \rho_p} \quad (62)$$

and

$$\rho_p = \frac{1 - P}{1 + P} \quad (63)$$

If one is determining the degree of polarization or the extinction ratio of a polarizer, the ratio of T_{\perp} to T_{\parallel} can be measured for two identical polarizers in unpolarized light. From Eqs. (57) and (58),

$$\frac{T_{\perp}}{T_{\parallel}} = \frac{T_1 T_2}{(T_1^2 + T_2^2)/2} \approx \frac{2T_2}{T_1} = 2\rho_p \quad (64)$$

if $T_2^2 \ll T_1^2$. If a perfect polarizer or a source of perfectly plane-polarized light is available, T_2/T_1 can be determined directly by measuring the ratio of the minimum to the maximum transmittance of the polarizer. Other relations for two identical partial polarizers are given by West and Jones,¹³ as well as the transmittance $T_{\theta ab}$ of two dissimilar partial polarizers a and b whose principal axes are inclined at an angle θ with respect to each other. This latter expression is

$$T_{\theta ab} = \frac{1}{2}(T_{1a}T_{1b} + T_{2a}T_{2b})\cos^2\theta + \frac{1}{2}(T_{1a}T_{2b} + T_{1b}T_{2a})\sin^2\theta \quad (65)$$

where the subscripts 1 and 2 refer to the principal transmittances, as before.

Spectrophotometric measurements can involve polarizers and dichroic samples. Dichroic (optically anisotropic) materials are those which absorb light polarized in one direction more strongly than light polarized at right angles to that direction. (Dichroic materials are to be distinguished from birefringent materials, which may have different refractive indices for the two electric vectors vibrating at right angles to each other but similar, usually negligible, absorption coefficients.) When making spectrophotometric measurements, one should know the degree of polarization of the polarizer and how to correct for instrumental polarization. This latter quantity may arise from nonnormal-incidence reflections from a grating, dispersing prism, or mirrors. Light sources are also sometimes polarized. Simon,¹⁴ Charney,¹⁵ Gonatas et al.,¹⁶ and Wizinowich¹⁷ suggest methods for dealing with imperfect polarizers, dichroic samples, and instrumental polarization. In addition, when a dichroic sample is placed between a polarizer and a spectrophotometer which itself acts like an imperfect polarizer, one has effectively three polarizers in series. This situation has been treated by Jones,¹⁸ who showed that anomalies can arise when the phase retardation of the polarizers is taken on certain values. Mielenz and Eckerle¹⁹ have discussed the accuracy of various types of polarization attenuators.

12.5 POLARIZATION BY NONNORMAL-INCIDENCE REFLECTION (PILE OF PLATES)

Pile-of-plates polarizers make use of reflection or transmission of light at nonnormal incidence, frequently near the Brewster or polarizing angle [Eq. (48) in Sec. 12.3]. The extinction ratio and "transmittance" of these polarizers can be calculated directly from the Fresnel equations. Some simplifications occur for nonabsorbing or slightly absorbing plates. Equations (20) and (21) give the values of the intensity reflection coefficients R_s and R_p for light vibrating perpendicular to the plane of incidence (s component) and parallel to the plane of incidence (p component). The angle of refraction θ_1 in those equations is related to the refractive index n of the material by

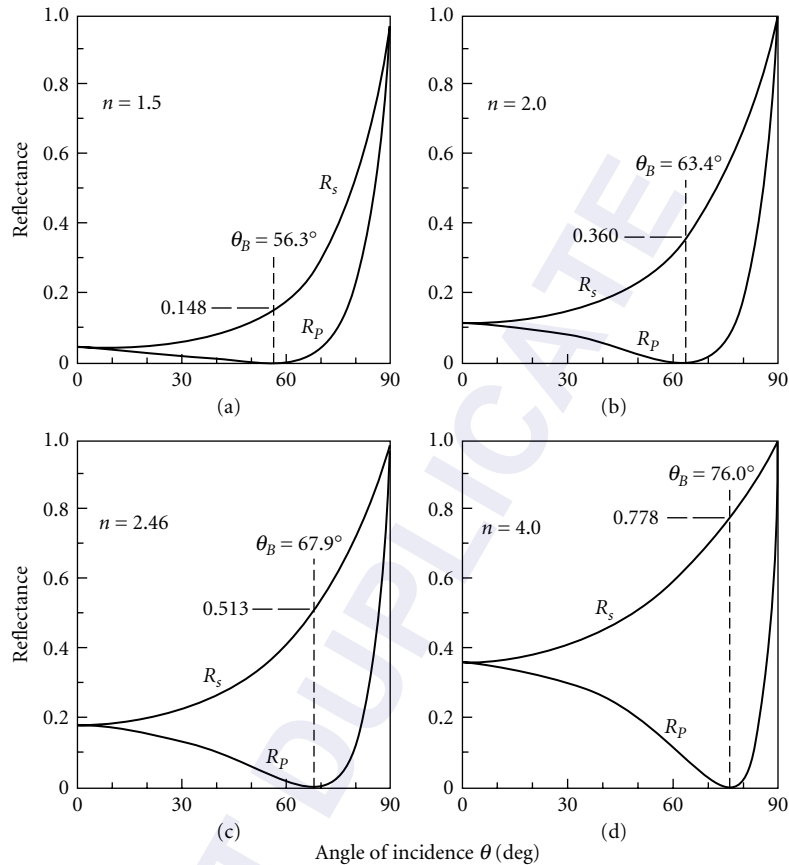


FIGURE 4 Reflectance of light polarized parallel R_p and perpendicular R_s to the plane of incidence from materials of different refractive index n as a function of angle of incidence: (a) $n = 1.5$ (alkali halides in ultraviolet and sheet plastics in infrared); (b) $n = 2.0$ (AgCl in infrared); (c) $n = 2.46$ (Se in infrared); and (d) $n = 4.0$ (Ge in infrared). The Brewster angle θ_B (at which R_p goes to 0) and the magnitude of R_s at θ_B are also indicated.

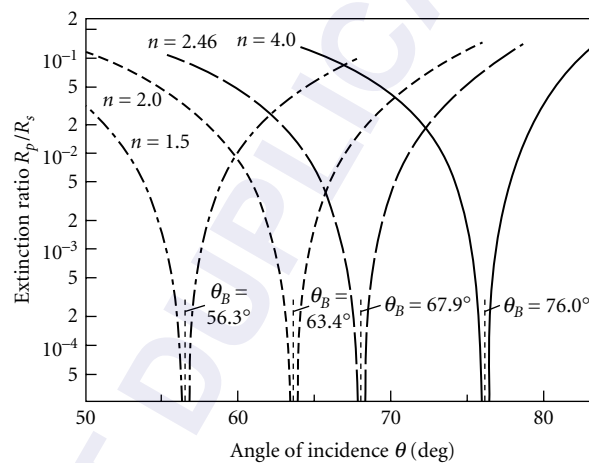
Snell's law [Eq. (15)*]. At the Brewster angle $R_p = 0$, so that the reflected light is, in principle, completely plane-polarized. This is the basis for all Brewster angle reflection polarizers.

Let us now see how the characteristics of a reflection polarizer depend on its refractive index. In Fig. 4 the reflectances R_s and R_p have been plotted for different values of the refractive index, roughly representing alkali halides in the ultraviolet and sheet-plastic materials, silver chloride, selenium, and germanium in the infrared. The Brewster angle, given by Eq. (48), is also indicated, as well as the magnitude of R_s at the Brewster angle. We note from these graphs that if light is polarized by a single reflection from a nonabsorbing material, the polarizer with the highest refractive index will have the largest throughput. In reflection polarizers, the quantity R_s is essentially the principal "transmittance" of the polarizer [T_r in Eqs. (55) to (65)] except that it must be multiplied by the reflectance of any other mirrors used to return the beam to its axial position.

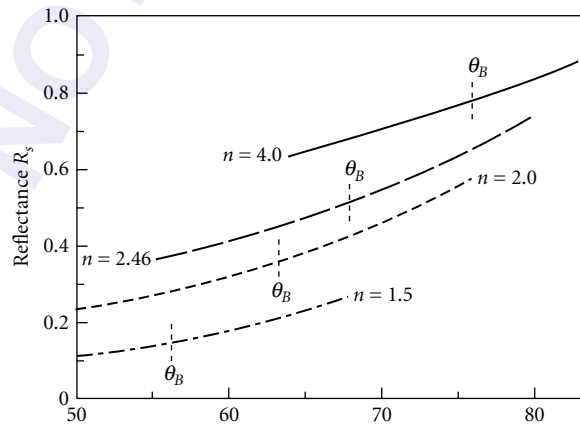
*Since we are assuming that the medium of incidence is air, $n_0 = 1$ and $n_1 = n$, the refractive index of the material.

The reflectance R_p can be equated to T_s , the minimum “transmittance” of the polarizer, so that the extinction ratio ρ_p of a reflection polarizer [Eq. (61)] is $\rho_p = R_p/R_s$. If R_p is really zero at the Brewster angle, the extinction ratio will be zero for all materials *independent of the value of n* . If a given extinction ratio is desired, for example, 10^{-3} [corresponding to 99.8 percent polarization; see Eq. (62)], then the convergence angle of the light beam must be small so that all the angles of incidence lie within about $\pm 1^\circ$ of the Brewster angle. The convergence angle depends only weakly on the refractive index for this case, varying from $\pm 1.2^\circ$ for $n = 1.5$ to $\pm 0.8^\circ$ for $n = 4.0$.

If a good extinction ratio is required for a beam of larger convergence angle, two polarizing reflections may be used. Then all the exponents in Fig. 5a are doubled, and the convergence angles for a given extinction ratio are greatly increased. To obtain a value of 10^{-3} with two reflections, the angle of incidence must be within about $\pm 6^\circ$ of the Brewster angle for values of n less than 3.5; for $n = 4$ it is reduced slightly and becomes more asymmetric ($+4.0$ and -5.2°). A disadvantage of having two reflections from the polarizing materials is that the throughput is reduced. All the values of



(a)



(b)

FIGURE 5 (a) Reflectance R_s and (b) extinction ratio R_p/R_s for materials of different refractive index at angles near the Brewster angle θ_B . A single surface of the material is assumed.

R_s in Fig. 5b are squared, so that for $n = 4$, $R_s = 0.78$ but $R_s^2 = 0.61$; for smaller refractive indexes the reduction in throughput is much greater.

The information shown graphically in Figs. 4 and 5 is given analytically in a paper by Azzam²⁰ who is concerned about the angular sensitivity of Brewster-angle reflection polarizers, particularly those made with silicon or germanium plates. Also, Murty and Shukla²¹ show analytically that the shadowy extinction patterns sometimes seen with a crossed Brewster angle reflection polarizer and analyzer are caused by light incident on the surfaces at angles different from the Brewster angle.

Although in many cases multiple reflections within a plate degrade its polarizing properties, this is not true for Brewster angle reflection polarizers. For multiple reflections within a plane-parallel plate of material

$$(R_{s,p})_{\text{plate}} = \frac{2R_{s,p}}{1 + R_{s,p}} \quad (66)$$

assuming no interference or absorption; R_s and R_p are given by Eqs. (20) and (21). Multiple reflections have a minor effect on the extinction ratio but the increase in R_s is appreciable. To fulfill the conditions of Eq. (66), the plate must have plane-parallel sides and be unbacked. We are also assuming that the plate is thick or nonuniform enough for interference effects within it to be neglected.

All the preceding discussion applies only to nonabsorbing materials. If a small amount of absorption is present, R_p will have a minimum that is very close to zero and the material will still make a good reflection polarizer. However, if the extinction coefficient k becomes appreciable, the minimum in R_p will increase and the polarizing efficiency will be degraded. By referring to Fig. 2 one can see roughly what the ratio of R_p to R_s will be for a given set of optical constants. Exact values of R_p and R_s can be calculated from n and k using Eqs. (45), (46), (31), and the other pertinent relations in Sec. 12.3. When choosing materials for possible use as metallic reflection polarizers, one wants the largest difference between R_s and R_p and the smallest magnitude of R_p at the minimum. Thus, ideally n should be much larger than k .

The Abelès condition²² applies to the amplitude reflectances r_s and r_p for either dielectrics or metals at 45° angle of incidence. At this angle

$$r_s^2 = r_p \quad (67)$$

and

$$2\delta_s = \delta_p \quad (68)$$

where the δ 's are the absolute phase changes on reflection for the p and s components (see Sec. 12.3). Relation in Eq. (67) is frequently applied to the intensity reflectances R_s and R_p , which are directly related to the amplitude reflectances [Eqs. (20), (21), and (31)].

12.6 POLARIZATION BY NONNORMAL-INCIDENCE TRANSMISSION (PILE OF PLATES)

The theory of Brewster angle transmission polarizers follows directly from that given for reflection polarizers. Table 1 lists the relations giving the s and p transmittances of the polarizers with various assumptions about multiple reflections, interference, absorption, etc.* All these relations contain R_s and R_p , the reflectances at a single interface, which are given at the bottom of the table.

At the Brewster angle, R_p at a single interface equals zero, and the transmittances of the plates can be expressed in terms of the refractive index of the material and the number of plates. The relations for the s and p transmittances at this angle are given in Table 2. Most references that contain

*Transmission polarizers in which the multiply internally reflected beams are coherent and produce interference effects are discussed in Chap. 13, "Polarizers."

TABLE 1 Transmittances and Degree of Polarization for a Single Plate and Multiple Plates at any Angle of Incidence in Terms of R_s and R_p for a Single Surface*

	m Plates ($2m$ Surfaces)	
	One Plate (Two Surfaces)	$(T_{s,p})_{\text{sample}}$
Single transmitted beam, no multiple reflections, no absorption	$(1 - R_{s,p})^2$	$\frac{T_p - T_s}{T_p + T_s} = \frac{1 - T_s/T_p}{1 + T_s/T_p}$
Multiple reflections within plate, no interference effects, no absorption	$\frac{1 - R_{s,p}^{\S}}{1 + R_{s,p}^{\S}}$	$\frac{1 - \cos^{4m}(\theta_0 - \theta_1)}{1 + \cos^{4m}(\theta_0 - \theta_1)}$
Multiple reflections within plate, no interference effects, small absorption	$\frac{1 - R_{s,p}^{\S} e^{-\alpha d}}{1 - R_{s,p}^2 e^{-2\alpha d}}$	$\frac{1 - \cos^{4m}(\theta_0 - \theta_1)}{1 + \cos^{4m}(\theta_0 - \theta_1)}$
Multiple reflections within plate, interference within plate, no absorption	$\frac{1}{1 + \frac{4R_{s,p}}{(1 - R_{s,p})^2} \sin^2 \gamma}$	$\frac{1 - \left[\frac{1 - R_{s,p}}{1 + R_{s,p}} \right]^{m \dagger, \S}}{1 + \left[\frac{1 - R_{s,p}}{1 + R_{s,p}} \right]^{m \dagger, \S}}$
Single surface	$R_s = \frac{\sin^2(\theta_0 - \theta_1)}{\sin^2(\theta_0 + \theta_1)}$ $T_s = 1 - R_s = \frac{\sin 2\theta_0 \sin 2\theta_1}{\sin^2(\theta_0 + \theta_1)}$	$\frac{1 - \cos^{4m}(\theta_0 - \theta_1)}{1 + \cos^{4m}(\theta_0 - \theta_1)}$ $P = \frac{1 - \cos^2(\theta_0 - \theta_1)}{1 + \cos^2(\theta_0 - \theta_1)}$
	$R_p = \frac{\tan^2(\theta_0 - \theta_1)}{\tan^2(\theta_0 + \theta_1)}$ $T_p = 1 - R_p = \frac{\sin 2\theta_0 \sin 2\theta_1}{\sin^2(\theta_0 + \theta_1) \cos^2(\theta_0 - \theta_1)}$	$\frac{1 - \cos^{4m}(\theta_0 - \theta_1)}{1 + \cos^{4m}(\theta_0 - \theta_1)}$ $P = \frac{1 - \cos^2(\theta_0 - \theta_1)}{1 + \cos^2(\theta_0 - \theta_1)}$

* $\alpha = 4\pi k' / (\lambda \cos \theta)$, $\gamma = 2\pi n d \cos \theta / \lambda$, θ_0 = angle of incidence, θ_1 = angle of refraction, n = refractive index = $(\sin \theta_0) / (\sin \theta_1)$, k = extinction coefficient, d = plate thickness, λ = wavelength.

[†]No multiple reflections between plates.

[‡]Multiple reflections between plates.

[§]Also holds for coherent multiple reflections averaged over one period of $\sin^2 \gamma$.

TABLE 2 Transmittances and Degree of Polarization for a Single Plate and Multiple Plates at the Brewster Angle θ_p , where $\tan \theta_p = n$ and $\theta_0 + \theta_1 = 90^\circ$ **

	One Plate (Two Surfaces)		<i>m</i> Plates (2 <i>m</i> Surfaces)	
	$(T_p)_{\text{sample}}$	$(T_s)_{\text{sample}}$	$(T_p)_{\text{sample}}$	$(T_s)_{\text{sample}}$
Single transmitted beams, no multiple reflections, no absorption	1	$\left(\frac{2n}{n^2+1}\right)^4$	1	$\left[\frac{2n}{n^2+1}\right]^{4mt}$
Multiple reflections within plate, no interference effects, no absorption	1	$\frac{2n^2}{n^4+1}$	1	$\frac{1-[2n/(n^2+1)]^{4mt}}{1+[2n/(n^2+1)]^{4m}}$
Multiple reflections within plate, no interference effects, small absorption	$e^{-\alpha d}$	$\frac{\left(\frac{2n}{n^2+1}\right)^4 e^{-\alpha d}}{1-\left(\frac{n^2-1}{n^2+1}\right)^4 e^{-2\alpha d}}$	1	$\frac{1-[2n^2/(n^4+1)]^{mt} e^{-\alpha d}}{1+[2n^2/(n^4+1)]^{mt}}$
Multiple reflections within plate, interference within plate, no absorption	1	$\frac{1}{1+\frac{(n^2-1)^2}{4n^4} \sin^2 \gamma}$	1	$\frac{1}{1+m(n^2-1)^2/2n^2}$
Single surface	$R_p = 0$	$R_s = \left(\frac{n^2-1}{n^2+1}\right)^2$	$R_p = 0$	$R_s = \left(\frac{n^2-1}{n^2+1}\right)^2$
	$T_p = 1 - R_p = 1$	$T_s = 1 - R_s = \left(\frac{2n}{n^2+1}\right)^2$		$T_p = \frac{1-[2n/(n^2+1)]^{2m}}{1+[2n/(n^2+1)]^{2m}}$

**Where $\alpha = 4\pi k(n^2+1)^{1/2}/\lambda n$, $\gamma = 2\pi n^2 d/\lambda(n^2+1)^{1/2}$, n = refractive index, k = extinction coefficient, d = plate thickness, λ = wavelength.

†No multiple reflections between plates.

‡Formula of Provostaye and Desains.²³

*Also holds for coherent multiple reflections averaged over one period of $\sin^2 \gamma$.

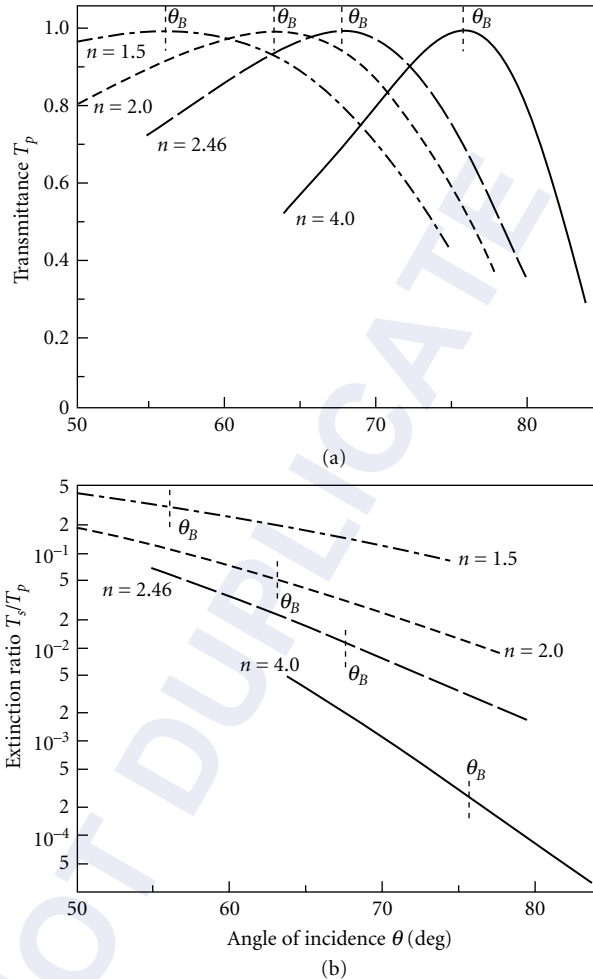


FIGURE 6 (a) Transmittance and (b) extinction ratio of four plane-parallel plates of refractive index n as a function of angle of incidence, for angles near the Brewster angle. Assumptions are multiple reflections but no interference within each plate and no reflections between plates.

an expression for the degree of polarization of a pile of plates give the formula of Provostaye and Desains,²³ which assumes an infinite series of multiple reflections between all surfaces, i.e., multiple reflections within and between plates. This assumption is not valid for most real transmission polarizers (see Chap. 13, “Polarizers,” specifically Brewster Angle Transmission Polarizers).

For most parallel-plate polarizers it is reasonable to assume incoherent multiple reflections within each plate and no reflections between plates. Figure 6 shows the principal transmittance (p component) and extinction ratio for several four-plate polarizers having the refractive indexes indicated.* The extinction ratio improves considerably with increasing refractive index. It is also

*The extinction ratio of a pile of m plates (no multiple reflections between plates) is simply the product of the extinction ratios of the individual plates.

improved by using the plates at an angle of incidence slightly above the Brewster angle. This procedure, which is most helpful for high refractive index plates, reduces the transmission per plate so that a trade-off is required between losses resulting from absorption or scattering when many plates are used and the reflectance loss per plate when only a few plates are used above the Brewster angle. In some cases significant improvements have been achieved by following the latter course.²⁴

When the number of plates of a given refractive index is increased, the transmittance is unaffected (in the absence of absorption) and the extinction ratio is greatly increased, as shown in the earlier polarization chapter.¹ In the absence of absorption, comparable transmittances and extinction ratios are obtained with a large number of low-refractive-index plates or a small number of high refractive index plates. Small amounts of absorption decrease the transmittance, but have little effect on the extinction ratio.¹ Tuckerman²⁵ has derived exact expressions for light reflected from or transmitted through a pile of absorbing plates. He has also noted mistakes that have been perpetuated in some of the formulas for light reflected from or transmitted through a pile of nonabsorbing plates.

A figure of merit giving the variation of the extinction ratio with angle of incidence can be defined as in Fig. 7, where the ordinate is the extinction ratio at a given angle of incidence divided by the extinction ratio at the Brewster angle. The angles of incidence are referred to the Brewster angle, and curves for different values of the refractive index are shown. These curves are calculated from the ratio

$$\frac{\left(\frac{T_s}{T_p}\right)_\theta}{\left(\frac{T_s}{T_p}\right)_{\theta_B}} = \frac{\left(\frac{1-R_s}{1+R_s} \frac{1+R_p}{1-R_p}\right)_\theta}{\left(\frac{1-R_s}{1+R_s}\right)_{\theta_B}} \quad (69)$$

and are for a single transparent film or plate having multiple incoherent internal reflections within the material. As an example of how to use the graphs, consider an optical system having a two-plate germanium polarizer with a refractive index of 4.0. If the angles of incidence vary from -1.4 to $+1.5^\circ$ around the Brewster angle, the ratio of the extinction ratios will vary between $1.10^2 = 1.21$ and

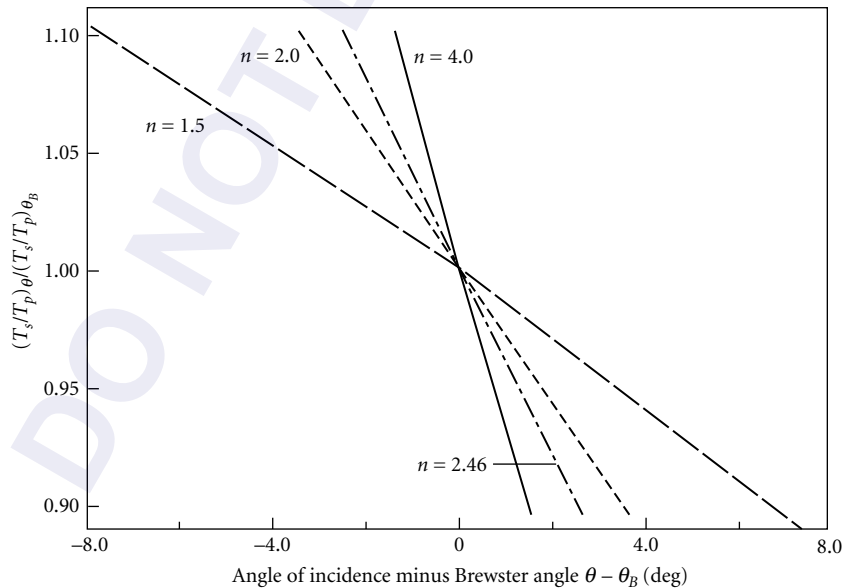


FIGURE 7 Variation of extinction ratio (per film) as a function of angle near the Brewster angle $\theta - \theta_B$. The ordinate is the extinction ratio at θ divided by the extinction ratio at θ_B .

$0.90^2 = 0.81$, respectively. (For m plates it would be 1.10^m and 0.90^m .) Thus, in order to restrict the percent variation of the extinction ratio to a given value, one must use a smaller acceptance angle when using more plates.

We have assumed that there are multiple incoherent reflections within each plate and no multiple reflections between plates. The difference in extinction ratios for a series of four plates with and without internal reflections is shown in Fig. 8. The principal transmittance is essentially the same as in Fig. 6 for values of T_p above 0.70 (and only about 0.025 lower when T_p drops to 0.30). However, the extinction ratio of high-refractive-index materials is much better without multiple internal reflections; for low-refractive-index materials the difference in extinction ratios is small.

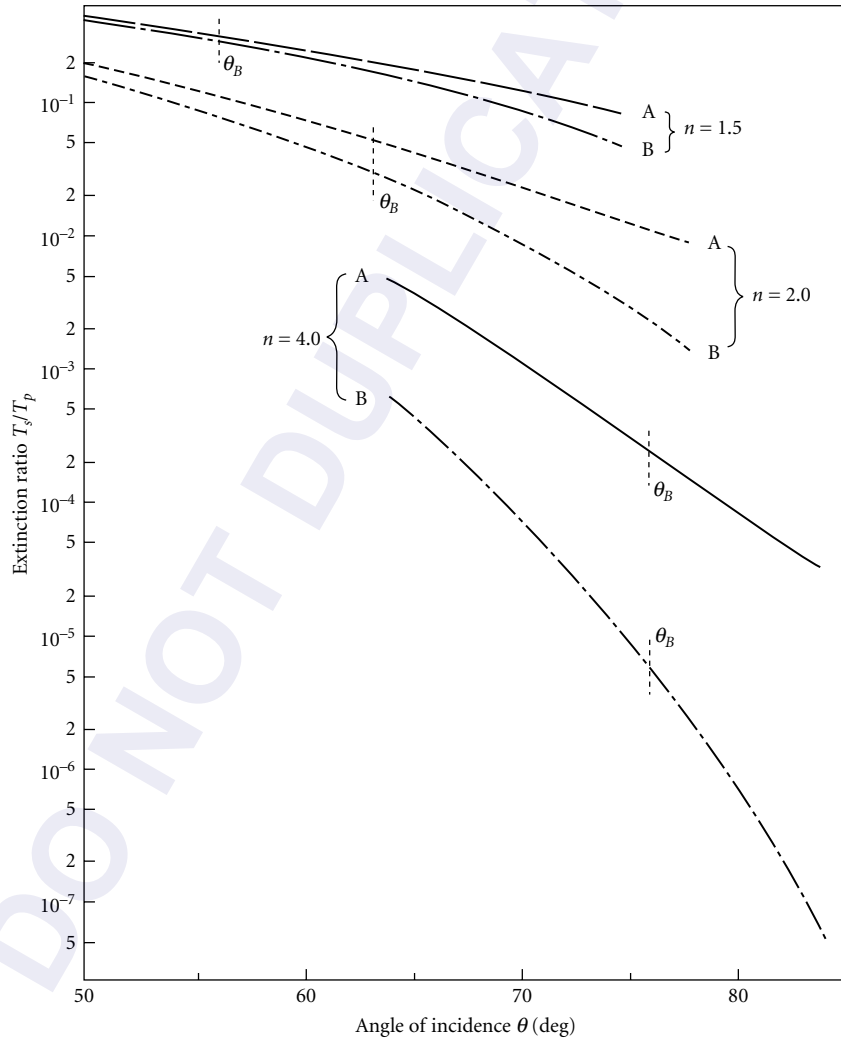


FIGURE 8 Extinction ratio of four plane-parallel plates of refractive index n as a function of angle of incidence for angles near the Brewster angle. Assumptions are A, multiple reflections but no interference within each plate and no reflections between plates; B, no multiple reflections within each plate or between plates. The transmittances for conditions A and B are essentially identical (see Fig. 6a).

The effect of multiple reflections on the extinction ratio can readily be seen from the three relations for the transmittances of the p and s components:

No multiple reflections:

$$(T_{s,p})_{\text{sample}} = (1 - R_{s,p})^{2m} = 1 - 2mR_{s,p} + 2m^2R_{s,p}^2 - mR_{s,p}^2 + \dots \quad (70)$$

Multiple reflections within plates:

$$(T_{s,p})_{\text{sample}} = \left(\frac{1 - R_{s,p}}{1 + R_{s,p}} \right)^m = 1 - 2mR_{s,p} + 2m^2R_{s,p}^2 + \dots \quad (71)$$

Multiple reflections within and between plates:

$$(T_{s,p})_{\text{sample}} = \frac{1 - R_{s,p}}{1 + (2m - 1)R_{s,p}} = 1 - 2mR_{s,p} + 4m^2R_{s,p}^2 - 2mR_{s,p}^2 + \dots \quad (72)$$

At the Brewster angle, $R_p = 0$, $T_p = 1$, and the extinction ratio will be smallest, i.e., highest degree of polarization, for the smallest values of the s transmittance. The first three terms in Eqs. (70) and (71) are identical, but Eq. (70) has an additional negative term in R_s^2 and so it will give a slightly smaller value of the s transmittance. Equation (72), from which the formula of Provostaye and Desains was derived, has twice as large a third term as the other two equations, and the negative fourth term is only $1/2m$ of the third term, so that it does not reduce the overall value of the expression appreciably. Thus, Eq. (72) gives an appreciably larger value of the s transmittance, but fortunately it is a limiting case and is rarely encountered experimentally.

12.7 QUARTER-WAVE PLATES AND OTHER PHASE RETARDATION PLATES

A retardation plate is a piece of birefringent, uniaxial (or uniaxial-appearing) material in which the ordinary and extraordinary rays travel at different velocities. Thus, one ray is retarded relative to the other, and the path $N\lambda$ between the two rays is given by

$$N\lambda = \pm d(n_e - n_o) \quad (73)$$

where n_o = refractive index of ordinary ray
 n_e = refractive index of extraordinary ray
 d = physical thickness of plate
 λ = wavelength

The positive sign is used when $n_e > n_o$, that is, a positive uniaxial crystal, and the negative sign is used for a negative uniaxial crystal, for which $n_e < n_o$. Since $N\lambda$ is the path difference between the two rays, N can be considered the retardation expressed in fractions of a wavelength. For example, $N = 1/4$ for a quarter-wave (or $\lambda/4$) plate, $1/2$ for a half-wave (or $\lambda/2$) plate, $3/4$ for a three-quarter-wave (or $3\lambda/4$) plate, etc.

The phase difference between two rays traveling through a birefringent material is $2\pi/\lambda$ times the path difference, so that the phase retardation δ is

$$\delta = 2\pi N = \pm \frac{2\pi d(n_e - n_o)}{\lambda} \quad (74)$$

Thus, phase differences of $\pi/2$, π , and $3\pi/2$ are introduced between the two beams in quarter-wave, half-wave, and three-quarter-wave plates, respectively.

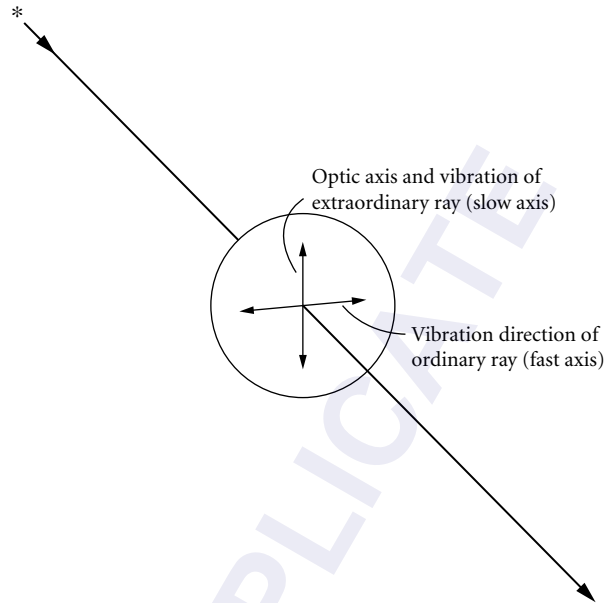


FIGURE 9 Light incident normally on the front surface of a retardation plate showing the vibration directions of the ordinary and extraordinary rays. In a positive uniaxial crystal, the fast and slow axes are as indicated in parentheses; in a negative uniaxial crystal, the two axes are interchanged.

A retardation plate can be made from a crystal which is cut so that the optic axis lies in a plane parallel to the face of the plate, as shown in Fig. 9. Consider a beam of unpolarized or plane-polarized light normally incident on the crystal. It can be resolved into two components traveling along the same path through the crystal but vibrating at right angles to each other. The ordinary ray vibrates in a direction perpendicular to the optic axis, while the extraordinary ray vibrates in a direction parallel to the optic axis. In a positive uniaxial crystal $n_e > n_o$, so that the extraordinary ray travels more slowly than the ordinary ray. The fast axis is defined as the direction in which the faster-moving ray vibrates; thus in a positive uniaxial crystal, the fast axis (ordinary ray) is perpendicular to the optic axis while the slow axis, (extraordinary ray) coincides with the optic axis. For a negative uniaxial crystal the fast axis coincides with the optic axis.

Figure 10 shows how the state of polarization of a light wave changes after passing through retardation plates of various thicknesses when the incident light is plane-polarized at an azimuth of 45° to the fast axis of the plate. If the plate has a retardation of $\lambda/8$, which means that the ordinary and extraordinary waves are out of phase by $\pi/4$ with each other, the transmitted light will be elliptically polarized with the major axis of the ellipse coinciding with the axis of the original plane-polarized beam. As the retardation gradually increases (plate gets thicker for a given wavelength or wavelength gets shorter for a given plate thickness), the ellipse gradually turns into a circle, but its major axis remains at 45° to the fast axis of the retardation plate. For a retardation of $\lambda/4$, the emerging light is right-circularly polarized. As the retardation continues to increase, the transmitted light becomes elliptically polarized with the major axis of the ellipse lying perpendicular to the plane of the incident polarized beam, and then the minor axis of the ellipse shrinks to zero and plane-polarized light is produced when the retardation becomes $\lambda/2$. As the retardation increases further, the patterns change in opposite order and the polarized light is left-circularly polarized when the retardation equals $3\lambda/4$. Finally, when the retardation is a full wave, the incident plane-polarized light is transmitted unchanged although the slow wave has now been retarded by a full wavelength relative to the fast wave.

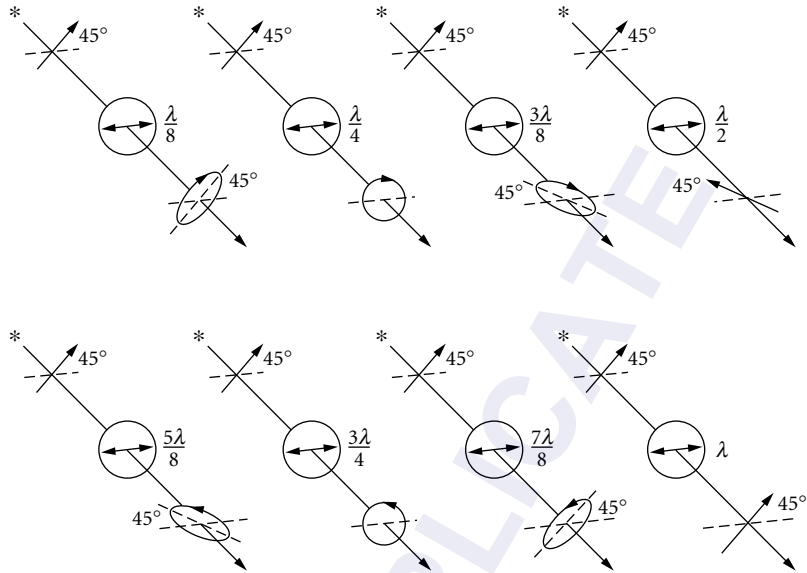


FIGURE 10 State of polarization of a light wave after passing through a crystal plate whose retardation is indicated in fractions of a wavelength (phase retardation $2\pi/\lambda$ times these values) and whose fast axis is indicated by the double arrow. In all cases the incident light is plane-polarized at an azimuth of 45° to the direction of the fast axis.

The most common type of retardation plate is the quarter-wave plate. Figure 11 shows how this plate affects the state of polarization of light passing through it when the fast axis is positioned in the horizontal plane and the azimuth of the incident plane-polarized light is changed from $\theta = 0^\circ$ to $\theta = 90^\circ$. When $\theta = 0^\circ$, only the ordinary ray (for a positive birefringent material) passes through the plate, so that the state of polarization of the beam is unchanged. When θ starts increasing, the transmitted beam is elliptically polarized with the major axis of the ellipse lying along the fast axis of the $\lambda/4$ plate; $\tan \theta = b/a$, the ratio of the minor to the major axis of the ellipse. In the next case, $\theta = 15^\circ$ and $\tan \theta = 0.268$, and so the ellipse is long and narrow. When the plane of vibration has rotated to an azimuth of 45° , the emerging beam is right-circularly polarized (the same situation as that shown in the second part of Fig. 10). For values of θ between 45° and 90° , the light is again elliptically polarized, this time with the major axis of the ellipse lying along the direction of the slow axis of the $\lambda/4$ plate. The angle shown in the figure is 60° , and $\tan 60^\circ = 1.732$, so that b/a (referred to the fast axis) is greater than unity. When θ increases to 90° , the plane of vibration coincides with the slow axis and the transmitted

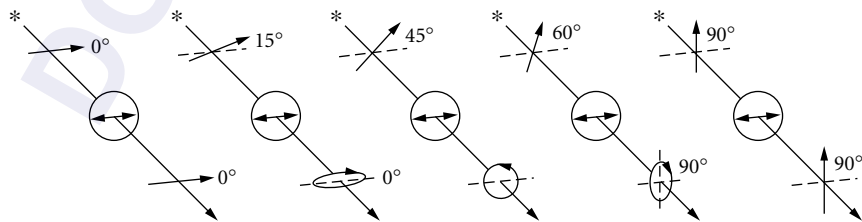


FIGURE 11 State of polarization of a light wave after passing through a $\lambda/4$ plate (whose fast axis is indicated by the double arrow) for different azimuths of the incident plane-polarized beam.

light is again plane-polarized. As θ continues to increase, the transmitted patterns repeat those already described and are symmetric about the slow axis, but the direction of rotation in the ellipse changes from right-handed, to left-handed, so that left-circularly polarized light is produced when $\theta = 135^\circ$.

The definition of right- and left-circularly polarized light should be clear from Figs. 10 and 11. When the rotation is *clockwise* with the observer looking *opposite to the direction of propagation*, the light is called *right-circularly polarized*; if the rotation is *counterclockwise*, the light is called *left-circularly polarized*.²⁶ When circularly polarized light is reflected from a mirror, the direction of propagation is reversed, so that the sense of the circular polarization changes; i.e., left-circularly polarized light changes on reflection into right-circularly polarized light and vice versa. Therefore, in experiments involving magnetic fields in which the sense of the circularly polarized light is important,^{27,28} it is important to know which kind one started with and how many mirror reflections occurred in the rest of the light path. Cyclotron resonance experiments can sometimes be used to determine the sense of the circular polarization.²⁸ Another method utilizing a polarizer and $\lambda/4$ plate has been described by Wood.²⁹

The behavior of a half-wave plate in a beam of plane-polarized light is completely different from that of a quarter-wave plate; the transmitted light is always plane-polarized. If the incident plane of vibration is at an azimuth θ with respect to the fast axis of the $\lambda/2$ plate, the transmitted beam will be rotated through an angle 2θ relative to the azimuth of the incident beam. The case showing $\theta = 45^\circ$ where the phase of vibration is rotated through 90° is illustrated in the fourth part of Fig. 10. In this situation the extraordinary beam is retarded by half a wavelength relative to the ordinary beam (for a positive birefringent material), hence the name, half-wave plate. If the polarizer is fixed and the $\lambda/2$ plate is rotated (or vice versa), the plane of vibration of the transmitted beam will rotate at twice the frequency of rotation of the $\lambda/2$ plate.

Quarter-wave plates are useful for analyzing all kinds of polarized light. In addition, they are widely employed in experiments using polarized light, e.g., measurements of the thickness and refractive index of thin films by ellipsometry or measurements of optical rotary dispersion, circular dichroism, or strain birefringence. Polarizing microscopes, interference microscopes, and petrographic microscopes are usually equipped with $\lambda/4$ plates. In some applications the $\lambda/4$ plate is needed only to produce circularly polarized light, e.g., for optical pumping in some laser experiments, or to convert a partially polarized light source into one which appears unpolarized, i.e., has equal amplitudes of vibration in all azimuths. For these and similar applications, one can sometimes use a circular polarizer which does not have all the other properties of a $\lambda/4$ plate (see Pars. 73 to 76 in Ref. 1).

The customary application for a $\lambda/2$ plate is to rotate the plane of polarization through an angle of 90° . In other applications the angle of rotation can be variable. Automatic-setting ellipsometers or polarimeters sometimes employ rotating $\lambda/2$ plates in which the azimuth of the transmitted beam rotates at twice the frequency of the $\lambda/2$ plate.

12.8 MATRIX METHODS FOR COMPUTING POLARIZATION

In dealing with problems involving polarized light, it is often necessary to determine the effect of various types of polarizers (linear, circular, elliptical, etc.), rotators, retardation plates, and other polarization-sensitive devices on the state of polarization of a light beam. The Poincaré sphere construction is helpful for giving a qualitative understanding of the problem; for quantitative calculations, one of several forms of matrix calculus can be used. The matrix methods are based on the fact that the effect of a polarizer or retarder is to perform a linear transformation (represented by a matrix) on the vector representation of a polarized light beam. The advantage of these methods over conventional techniques is that problems are reduced to simple matrix operations; thus since one does not have to think through the physics of every problem, the probability of making an error is greatly reduced. The most common forms of matrix calculus are the Mueller calculus and the Jones calculus, but the coherency-matrix formulation is also gaining popularity for dealing with problems involving partially polarized light. We give here a brief description of the Poincaré sphere and the

different matrix methods, indicating how they are used, the different types of problems for which they are helpful, and where complete descriptions of each may be found.

The *Poincaré sphere* is a useful device for visualizing the effects of polarizers and retarders on a beam of polarized light. The various states of polarization are represented on the sphere as follows. The equator represents various forms of linear polarization, the poles represent right- and left-circular polarization,* and other points on the sphere represent elliptically polarized light. Every point on the sphere corresponds to a different polarization form. The radius of the sphere indicates the intensity of the light beam (which is usually assumed to be unity). The effects of polarizers and retarders are determined by appropriate displacements on the sphere. Partially polarized light or absorption may be dealt with approximately by ignoring the intensity factor, since one is generally interested only in the state of polarization; however, the construction is most useful when dealing with non-absorbing materials. Good introductory descriptions of the Poincaré sphere, including references, can be found in *Polarized Light* by Shurcliff,^{30†} *Ellipsometry and Polarized Light* by Azzam and Bashara,³¹ and *Polarized Light in Optics and Spectroscopy* by Klinger Lewis and Randall,^{32‡} illustrative examples and problems are given in Sutton and Panati.³³ More comprehensive treatments are given by Ramachandran and Ramaseshan^{34‡} and Jerrard^{35‡} and include numerous examples of applications to various types of problems. The new book *Polarized Light, Fundamentals and Applications* by Collett³⁶ has a comprehensive 35-page chapter on the mathematical aspects of the Poincaré sphere; this material can be best understood after reading some of the introductory descriptions of the Poincaré sphere. The main advantage of the Poincaré sphere, like other graphical methods, is to reveal by essentially a physical argument which terms in exceedingly complex equations are negligible or can be made negligible by modifying the experiment. It is characteristic of problems in polarized light that the trigonometric equations are opaque to inspection and yield useful results only after exact calculation with the aid of a computer or after complex manipulation and rather clever trigonometric identities. The Poincaré sphere thus serves as a guide to the physical interpretation of otherwise obscure polarization phenomena. It can be used for solving problems involving retarders or combinations of retarders,^{30,32,36–39} compensators, half-shade devices, and depolarizers,³⁴ and it has also been applied to ellipsometric problems⁴⁰ and stress-optical measurements.⁴¹

The Poincaré sphere is based on the Stokes vectors, which are sometimes designated $S_0, S_1, S_2,$ and S_3 . The physical interpretation of the vectors is as follows. S_0 is the intensity of the light beam, corresponding to the radius of the Poincaré sphere. S_1 is the difference in intensities between the horizontal and vertical polarization components of the beam; when S_1 is positive, the preference is for horizontal polarization, and when it is negative, the preference is for vertical polarization.[§] S_2 indicates preference for $+45^\circ$ or -45° polarization, depending upon whether it is positive or negative, and S_3 gives the preference for right- or left-circular polarization. The Stokes vectors $S_1, S_2,$ and S_3 are simply the three cartesian coordinates of a point on the Poincaré sphere: S_1 and S_2 are perpendicular to each other in the equatorial plane, and S_3 points toward the north pole of the sphere.[¶] Thus, any state of polarization of a light beam can be specified by these three Stokes vectors. The intensity vector S_0 is related to the other three by the relation $S_0^2 = S_1^2 + S_2^2 + S_3^2$ when the beam is completely polarized. If the beam is partially polarized, $S_0^2 > S_1^2 + S_2^2 + S_3^2$. Good introductory material on Stokes vectors is given by Shurcliff,³⁰ Azzam and Bashara,³¹ Klinger et al.,³² Sutton and Panati,³³ and Walker.⁴² A comprehensive discussion of the Stokes vectors has been given by Collett.³⁶ Rigorous definitions of the simple vectors and those for partially coherent light can be found in Born and Wolf;⁴³ other authors are cited by Shurcliff³⁰ and Collett.³⁶ Stokes vectors are generally used in conjunction with the Mueller calculus, and some examples of applications will be given there. We note here that Budde⁴⁴ has demonstrated a method for experimentally determining the Stokes vectors and other polarization parameters from a Fourier analysis of measured quantities. Ioshpa and Obridko⁴⁵ have proposed a photoelectric

*Right-circularly polarized light is defined as a *clockwise* rotation of the electric vector when the observer is looking *against* the direction the wave is traveling.

†Schurcliff and Klinger, Lewis, and Randall have the S_3 axis pointing down, so that the upper pole represents left-circular polarization. The more logical convention, followed by most others, is for the upper pole to represent right-circular polarization.

‡The notation is similar to that used by Schurcliff,³⁰ with the upper pole representing left-circular polarization.

§Some authors dealing with light scattering from aerosols define S_1 as positive when the preference is for vertical polarization.

¶See Schurcliff and Klinger, Lewis and Randall footnote, p. 5.26.

method for simultaneously and independently measuring the four Stokes parameters. Collett⁴⁶ has developed a method for measuring the four Stokes vectors using a single circular polarizer. Azzam and coworkers^{47–51} have built, tested, analyzed, and calibrated a four-detector photopolarimeter for measuring normalized Stokes vectors of a large number of polarization states, and have given a physical meaning to the rows and columns in the instrument matrix. Other methods for measuring Stokes parameters are discussed by Collett.³⁶ Hauge⁵² has surveyed different types of methods for completely determining the state of polarization of a light beam using combinations of Stokes vectors.

The matrix methods for solving problems involving polarized light have certain properties in common. All use some type of representation for the original light beam (assumed to be a plane wave traveling in a given direction) that uniquely describes its state of polarization. Generally the beam is completely polarized, but for some of the matrix methods it can also be unpolarized or partially polarized or its phase may be specified. The beam encounters one or more devices which change its state of polarization. These are called *instruments* and are represented by appropriate matrices. After the instruments operate on the light beam, it emerges as an outgoing plane wave in an altered state of polarization. The basic problem for all the methods is to find a suitable representation for the incident plane wave (usually a two- or four-component column vector), and the correct matrices (2×2 or 4×4) to represent the instruments. Once the problem is set up, one can perform the appropriate matrix operations to obtain a representation for the outgoing plane wave. Its properties are interpreted in the same way as the properties of the incident plane wave.

An introduction to the Jones and Mueller calculus is given by Shurcliff,³⁰ Azzam and Bashara,³¹ and Klier et al.,³² and an excellent systematic and rigorous discussion of all the matrix methods has been given by O'Neill⁵³ and Collett.³⁶ All references contain tables of vectors for the various types of polarized beams and tables of instrument matrices. More complete tables are given by Sutton and Panati.³³ In the Mueller calculus the beam is represented by the four-component Stokes vector, written as a column vector. This vector has all real elements and gives information about *intensity* properties of the beam. Thus it is not able to handle problems involving phase changes or combinations of two beams that are coherent. The instrument matrix is a 4×4 matrix with all real elements. In the Jones calculus, the Jones vector is a two-component column vector that generally has complex elements. It contains information about the *amplitude* properties of the beam and hence is well suited for handling coherency problems. However, it cannot handle problems involving depolarization, as the Mueller calculus can. The Jones instrument matrix is a 2×2 matrix whose elements are generally complex.

Shurcliff³⁰ has noted some additional differences between Jones calculus and Mueller calculus. The Jones calculus is well suited to problems involving a large number of similar devices arranged in series in a regular manner and permits an investigator to arrive at an answer expressed explicitly in terms of the number of such devices. The Mueller calculus is not suited for this type of problem. The Jones instrument matrix of a train of transparent or absorbing nondepolarizing polarizers and retarders contains no redundant information. The matrix contains four elements each of which has two parts, so that there are a total of eight constants, none of which is a function of any other. The Mueller instrument matrix of such a train contains much redundancy; there are 16 constants but only 7 of them are independent.

In order to handle problems involving partially coherent polarized light, coherency-matrix formalism has been developed. In this system the beam is represented by a 4×4 matrix called a *coherency* or *density matrix*, which is the time average of the product of the Jones vector with its hermitian conjugate. The instrument matrices are the same as those used in the Jones calculus. O'Neill⁵³ and Born and Wolf⁴³ have good basic descriptions of coherency-matrix formalism; later extensions of the theory are given by Marathay.^{54,55}

There have been some modifications of the various matrix methods. Priebe⁵⁶ has introduced an operational notation for the Mueller matrices that facilitates the analysis by simplifying the functional description of a train of optical components. Collins and Steel⁵⁷ have suggested a modification of the Jones calculus in which the light vector is expressed as the sum of two circularly polarized (rather than linearly polarized) components. Schmieder⁵⁸ has given a unified treatment of Jones calculus and Mueller calculus including the coherency matrix and has shown that if the Stokes parameters are ordered in a different way from that customarily used, familiar relationships

are preserved and the rotation matrix looks like a rotation matrix rather than like a rearranged one. Tewarson⁵⁹ presents a generalized reciprocity equation expressing an algebraic relationship between the parameters of an optical system and its reciprocal system and has verified the equation for both plane-polarized and circularly polarized light beams. Since his equation follows from the reciprocity law in the Mueller calculus, that law is verified also. Cernosek⁶⁰ presents a simple geometric method based on the properties of quaternions to give a quick, quantitative analysis of the effect of any combination of linear retarders and rotators on the state of polarization of a system.

Among the applications of Mueller calculus and Jones calculus to problems involving polarized light, McCrackin⁶¹ has used both matrix methods to analyze instrumental errors in ellipsometry, and Hellerstein⁶² has used Mueller calculus to study the passage of linearly, circularly, and elliptically polarized light through a Sénarmont polariscope. Azzam and Bashara⁶³ have used Jones calculus to give a unified analysis of errors in ellipsometry, including effects of birefringence in cell windows, imperfect components, and incorrect azimuth angles. Azzam⁶⁴ also describes a simple photopolarimeter with rotating polarizer and analyzer for measuring Jones and Mueller matrices.

12.9 REFERENCES

1. H. E. Bennett and J. M. Bennett, "Polarization," in *Handbook of Optics*, W. G. Driscoll and W. Vaughan, eds. (McGraw-Hill, New York, 1978), pp. 10-1 to 10-164.
2. E. Collett, private communication, 1992.
3. R. H. Muller, Proc. Symp Recent Dev. Ellipsometry, *Surf. Sci.* **16**:14-33 (1969).
4. W. R. Hunter, *J. Opt. Soc. Am.* **55**:1197-1204 (1965).
5. W. E. Williams, *Applications of Interferometry*, 4th ed. (Wiley, New York, 1950), pp. 77-78.
6. J. A. Berning and P. H. Berning, *J. Opt. Soc. Am.* **50**:813-815 (1960).
7. P. H. Berning, "Theory and Calculations of Optical Thin Films," in *Physics of Thin Films*, vol. 1, G. Hass ed. (Academic Press, New York, 1963), pp. 78-81.
8. F. Abelès, *Prog. Opt.* **2**:251-288 (1963).
9. H. B. Holl, *The Reflection of Electromagnetic Radiation*, vols. 1 and 2 (U.S. Army Missile Command, Redstone Arsenal, Huntsville, Alabama, 1963), Report RF-63-4.
10. H. E. Bennett and J. M. Bennett, "Precision Measurements in Thin Film Optics," in *Physics of Thin Films*, vol. 4, G. Hass and R. E. Thun, eds. (Academic Press, New York, 1967), pp. 69-78.
11. R. C. Jones, *J. Opt. Soc. Am.* **52**:747-752 (1962).
12. G. R. Bird and W. A. Shurcliff, *J. Opt. Soc. Am.* **49**:235-237 (1959).
13. C. D. West and R. C. Jones, *J. Opt. Soc. Am.* **41**:976-982 (1951).
14. I. Simon, *J. Opt. Soc. Am.* **41**:336-345 (1951).
15. E. Charney, *J. Opt. Soc. Am.* **45**:980-983 (1955).
16. D. P. Gonatas, X. D. Wu, G. Novak, and R. H. Hildebrand, *Appl. Opt.* **28**:1000-1006 (1989).
17. P. L. Wizinowich, *Opt. Eng.* **28**:157-159 (1989).
18. R. C. Jones, *J. Opt. Soc. Am.* **46**:528-533 (1956).
19. K. D. Mielenz and K. L. Eckerle, *Appl. Opt.* **11**:594-603 (1972).
20. R. M. A. Azzam, *Appl. Opt.* **26**:2847-2850 (1987).
21. M. V. R. K. Murty and R. P. Shukla, *Appl. Opt.* **22**:1094-1098 (1983).
22. F. Abelès, *C. R. Acad. Sci.* **230**:1942-1943 (1950).
23. M. F. de la Provostaye and P. Desains, *Ann. Chim. Phys.*, ser. 3, **30**:158 (1850).
24. H. E. Bennett, J. M. Bennett, and M. R. Nagel, *J. Opt. Soc. Am.* **51**:237 (1961).
25. L. B. Tuckerman, *J. Opt. Soc. Am.* **37**:818-825 (1947).
26. F. A. Jenkins and H. E. White, *Fundamentals of Optics*, 3d ed. (McGraw-Hill, New York, 1957), p. 229.

27. E. D. Palik, *Appl. Opt.* **2**:527–539 (1963).
28. P. L. Richards and G. E. Smith, *Rev. Sci. Instrum.* **35**:1535–1537 (1964).
29. R. W. Wood, *Physical Optics*, 3d ed. (Macmillan, New York, 1934), p. 360.
30. W. A. Shurcliff, *Polarized Light* (Harvard University Press, Cambridge, Mass., 1962), pp. 15–29, 95–98, 109–123.
31. R. M. A. Azzam and N. M. Bashara, *Ellipsometry and Polarized Light* (North-Holland Publishing Company, Amsterdam, 1977), Chapters 1 and 2.
32. D. S. Kliger, J. W. Lewis, and C. E. Randall, *Polarized Light in Optics and Spectroscopy* (Academic Press, Inc., San Diego, 1990), Chapters 1 to 5.
33. A. M. Sutton and C. F. Panati, Lecture notes, Modern Optics, RCA Institute Clark, N.J. (1969).
34. G. N. Ramachandran and S. Ramaseshan, “Crystal Optics,” in *Handbuch der Physik*, S. Flügge, ed. vol. 25/1 (Springer, Berlin, 1961), pp. 1–54.
35. H. G. Jerrard, *J. Opt. Soc. Am.* **44**:634–640 (1954).
36. E. Collett, *Polarized Light: Fundamentals and Applications* (Marcel Dekker, Inc., New York, 1992).
37. C. J. Koester, *J. Opt. Soc. Am.* **49**:405–409 (1959).
38. S. Pancharatnam, *Proc. Indian Acad. Sci.* **A41**:130–136 (1955).
39. S. Pancharatnam, *Proc. Indian Acad. Sci.* **A41**:137–144 (1955).
40. F. L. McCrackin, E. Passaglia, R. R. Stromberg, and H. L. Steinberg, *J. Res. Natl. Bur. Stand (U.S.)* **67A**:363–377 (1963).
41. A. J. Michael, *J. Opt. Soc. Am.* **58**:889–894 (1968).
42. M. J. Walker, *Am. J. Phys.* **22**:170–174 (1954).
43. M. Born and E. Wolf, *Principles of Optics*, 6th ed. (Pergamon Press, New York, 1980), pp. 30–32, 544–555.
44. W. Budde, *Appl. Opt.* **1**:201–205 (1962).
45. B. A. Iospha and V. N. Obridko, *Opt. Spectrosc. (USSR)* **15**:60–62 (1963).
46. E. Collett, *Opt. Commun.* **52**:77–80 (1984).
47. R. M. A. Azzam, *Opt. Lett.* **10**:110–112 (1985).
48. R. M. A. Azzam, I. M. Elminyawi, and A. M. El-Saba, *J. Opt. Soc. Am.* **A5**:681–689 (1988).
49. R. M. A. Azzam, E. Masetti, I. M. Elminyawi, and F. G. Grosz, *Rev. Sci. Instr.* **59**:84–88 (1988).
50. R. M. A. Azzam and A. G. Lopez, *J. Opt. Soc. Am.* **A6**:1513–1521 (1989).
51. R. M. A. Azzam, *J. Opt. Soc. Am.* **A7**:87–91 (1990).
52. P. S. Hauge, *Proc. Soc. Photo-Opt. Instrum. Eng.* **88**:3–10 (1976).
53. E. L. O’Neill, *Introduction to Statistical Optics* (Addison-Wesley, Reading, Mass, 1963), pp. 133–156.
54. A. S. Marathay, *J. Opt. Soc. Am.* **55**:969–980 (1965).
55. A. S. Marathay, *J. Opt. Soc. Am.* **56**:619–623 (1966).
56. J. R. Priebe, *J. Opt. Soc. Am.* **59**:176–180 (1969).
57. J. G. Collins and W. H. Steele, *J. Opt. Soc. Am.* **52**:339 (1962).
58. R. W. Schmieder, *J. Opt. Soc. Am.* **59**:297–302 (1969).
59. S. P. Tewarson, *Indian J. Phys.* **40**:281–293, 562–566 (1966).
60. J. Cernosek, *J. Opt. Soc. Am.* **61**:324–327 (1971).
61. F. L. McCrackin, *J. Opt. Soc. Am.* **60**:57–63 (1970).
62. D. Hellerstein, *Appl. Opt.* **2**:801–805 (1963).
63. R. M. A. Azzam and N. M. Bashara, *J. Opt. Soc. Am.* **61**:600–607, 773–776, 1380–1391 (1971).
64. R. M. A. Azzam, *Opt. Commun.* **25**:137–140 (1978).

This page intentionally left blank.

DO NOT DUPLICATE

POLARIZERS

Jean M. Bennett*

*Research Department
Michelson Laboratory
Naval Air Warfare Center
China Lake, California*

13.1 GLOSSARY

D	optical density
d	grid spacing
e	extraordinary
i	angle of incidence
i	semicone angle
M	positive integer
m	number of plates
N	1/4, 1/2
n	refractive index
o	ordinary
S	cut angle
T	intensity transmittance
α	faces angle or angle between normal and optical axis
α_i	absorption coefficient for i th component
β	angle between normal and optical axis
γ	maximum variation plane of vibration
δ	deviation angle
Δn	change in retardation
Δn	$n_e - n_o$
λ	wavelength
ν	frequency (wave number)
ϕ	angle ϕ
ϕ	angle to wave normal

*Deceased.

13.2 PRISM POLARIZERS

The material on prism polarizers is abridged from the much more complete treatment by Bennett and Bennett.¹ Basic relations for polarizers are given in Sec. 12.4 of Chap. 12, "Polarization."

Double Refraction in Calcite

Although many minerals, specifically those which do not have a cubic crystal structure, are doubly refracting, nearly all polarizing prisms used in the visible, near-ultraviolet, and near-infrared regions of the spectrum are made from optical calcite, which exhibits strong birefringence over a wide wavelength range. Polarizing prisms made from other birefringent crystals are used primarily in the ultraviolet and infrared at wavelengths for which calcite is opaque (see Sec. 13.7).

Next to quartz, calcite is the most widely distributed of all minerals and usually occurs in an impure polycrystalline form as marble, limestone, or chalk. Optical calcite, or Iceland spar, which is quite rare, originally came from a large deposit on the east coast of Iceland. This source is now exhausted, and optical calcite now comes principally from Mexico, Africa, and Siberia. It has been grown artificially by a hybrid gel-solution method,² but maximum edge lengths are only 3 to 4 mm.

Although calcite is much softer than glass, with care it can be worked to an excellent polish. Surfaces flat to one-fifth fringe, or even, with care, one-tenth fringe, which are free from surface defects or perceptible turned edges can be produced using more or less conventional pitch-polishing techniques.³ Such techniques fail only for surfaces normal to the optic axis, in which case pitch polishing tends to cleave out small tetrahedra. Such surfaces can be polished to a lower surface quality using cloth polishers.

Crystals of calcite are negative uniaxial and display a prominent double refraction. The material can easily be cleaved along three distinct planes, making it possible to produce rhombs of the form shown in Fig. 1. At points *B* and *H*, a given face makes an angle of $101^{\circ}55'$ with each of the other two. At all the other points, two of the angles are $78^{\circ}55'$ and one is $101^{\circ}55'$. The *optic axis* *HI*, the direction in the crystal along which the two sets of refracted waves travel at the same velocity, makes equal

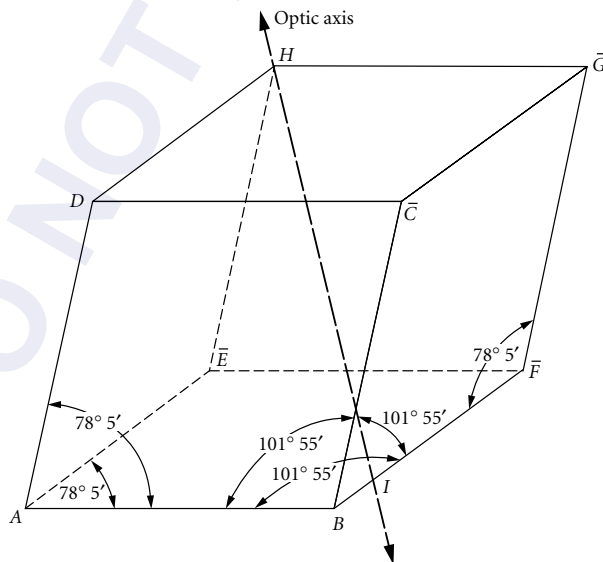


FIGURE 1 Schematic representation of a rhombohedral calcite crystal showing the angles between faces. The optic axis passes through corner *H* and point *I* on side *BF*.

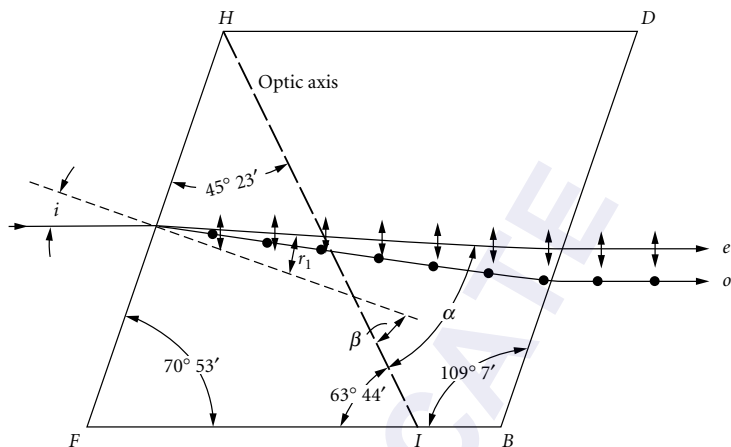


FIGURE 2 Side view of a principal section for the calcite rhomb in Fig. 1. The direction of the optic axis and the angles of the principal section are indicated. The angle of incidence is i , angle of refraction is r , angle between the e ray and the optic axis is α , and angle between the normal to the surface and the optic axis is β . The directions of vibration of the e and o rays are in the plane of the paper and perpendicular to it, respectively.

angles with all three faces at point H .^{*} Any plane, such as $DBFH$, which contains the optic axis and is perpendicular to the two opposite faces of the rhomb $ABCD$ and $EFGH$ is called a *principal section*. A side view of the principal section $DBFH$ is shown in Fig. 2. If light is incident on the rhomb so that the plane of incidence coincides with a principal section, the light is broken up into two components polarized at right angles to each other. One of these, the ordinary ray o , obeys Snell's law and has its plane of vibration (of the electric vector) perpendicular to the principal section. The second, the extraordinary ray e , has its plane of vibration parallel to the principal section. The refraction of the extraordinary ray in some cases violates Snell's law, at least in its simple form. The anomalous deflection of the ray is caused by the wavefront becoming ellipsoidal, so that the direction of propagation of the light is not along the wave normal. This ellipticity causes the velocity of the light in the crystal, and hence its refractive index, to be a function of angle. If light is incident on rhomb face $EFGH$ parallel to edge BF of the rhomb, the o and e rays, both of which lie in a principal section, are as shown in Fig. 2. As the angle of incidence is changed in Fig. 2 so that the direction taken by the o ray approaches that of the optic axis HI , the separation between the e and o rays decreases. If the rhomb is rotated about an axis parallel to HD , the e ray will precess about the o ray. However, unlike the o ray, it will not remain in the plane of incidence unless this plane coincides with the principal section.

The plane containing the o ray and the optic axis is defined as the *principal plane of the o ray*, and that containing the e ray and the optic axis as the *principal plane of the e ray*. In the case discussed earlier, the two principal planes and the principal section coincide. In the general case, they may all be different. However, in all cases, the o ray is polarized with its plane of vibration perpendicular to its principal plane and the e ray with its plane of vibration in its principal plane (see Fig. 2). In all cases, the vibration direction of the e ray remains perpendicular to that of the o ray.

The value of the index of refraction of the e ray which differs most from that of the o ray, i.e., the index when the e ray vibrations are parallel to the optic axis, is called the *principal index for the extraordinary ray* n_e . Snell's law can be used to calculate the path of the e ray through a prism for this case. Snell's law can always be used to calculate the direction of propagation of the ordinary ray.

Table 1 lists values of n_o and n_e for calcite, along with the two absorption coefficients a_o and a_e , all as a function of wavelength. Since $n_e < n_o$ in the ultraviolet, visible and infrared regions, calcite is a

^{*}The direction of the optic axis in a uniaxial crystal such as calcite or crystalline quartz can be determined by observing the crystal between crossed polarizers. If the alignment is correct, so that the optic axis is parallel to the line of sight, there will be concentric colored circles with a black cross superimposed.⁴

TABLE 1 Refractive Indices^a and Absorption Coefficients^a for Calcite

λ , μm	n_o	α_o	n_e	α_e	λ , μm	n_o	α_o	n_e	α_e
0.1318	1.56 ^b	534,000 ^b	1.80 ^b	477,000 ^b	0.3195	—	0.059	—	—
0.1355	1.48	473,000	1.84	380,000	0.327	—	0.028	—	—
0.1411	1.40	561,000	1.82	196,000	0.330	1.70515	—	1.50746	—
0.1447	1.48	669,000	1.80	87,000	0.3355	—	0.028	—	—
0.1467	1.51	711,000	1.75	20,500	0.340	1.70078	—	1.50562	—
0.1478 ₅	1.54	722,000	1.75	17,000	0.3450	—	0.0170	—	—
0.1487	1.58	735,000	1.75	14,400	0.346	1.69833	—	1.50450	—
0.1495 ₅	1.62	714,000	1.75	12,600	0.3565	—	0.0112	—	—
0.1513	1.68	756,000	1.75	8,300	0.361	1.69316	—	1.50224	—
0.1518 ₅	1.72	753,000	1.74	10,700	0.3685	—	0.0056	—	—
0.1536	1.80	761,000	1.74	9,000	0.3820	—	0.0056	—	—
0.1544 ₅	1.87	748,000	1.74	6,500	0.394	1.68374	—	1.49810	—
0.1558 ₅	1.92	766,000	1.74	8,100	0.397	—	0.000	1.49640 ^c	—
0.1581 ₅	2.02	715,000	1.73	11,100	0.410	1.68014 ^c	—	1.49430	—
0.1596	2.14	669,000	1.72	12,600	0.434	1.67552	—	1.49373	—
0.1608	2.20	594,000	1.70	13,300	0.441	1.67423	—	1.48956	—
0.1620	2.10	566,000	1.65	14,000	0.508	1.66527	—	1.48841	—
0.1633	2.00	608,000	1.65	10,800	0.533	1.66277	—	1.48736	—
0.1662	2.00	559,000	1.64	7,500	0.560	1.66046	—	1.48640	—
0.1700	1.94	414,000	1.63	≤4,400	0.589	1.65835	—	1.48490	—
0.1800	1.70	391,000	1.61	≤1,400	0.643	1.65504	—	1.48459	—
0.1900	1.72	278,000	1.59	≤321 ^d	0.656	1.65437	—	1.48426	—
0.198	—	—	1.57796 ^c	—	0.670	1.65367	—	1.48353	—
0.200	1.90284 ^c	257,000	1.57649	133	0.706	1.65207	—	1.48259	—
0.204	1.88242	—	1.57081	—	0.768	1.64974	—	1.48215	—
0.208	1.86733	149,000	1.56640	—	0.795	1.64886	—	1.48216	—
0.211	1.85692	—	1.56327	—	0.801	1.64869	—	1.48176	—
0.214	1.84558	—	1.55976	~0.1	0.833	1.64772	—	1.48137	—
0.219	1.83075	—	1.55496	—	0.867	1.64676	—	1.48098	—
0.226	1.81309	—	1.54921	—	0.905	1.64578	—	1.48060	—
0.231	1.80233	—	1.54541	—	0.946	1.64480	—	1.48022	—
0.242	1.78111	—	1.53782	—	0.991	1.64380	—	1.47985	—
0.2475	—	0.159 ^c	—	—	1.042	1.64276	—	1.47948	—
0.2520	—	0.125	—	—	1.097	1.64167	—	1.47910	—
0.256	—	0.109	—	—	1.159	1.64051	—	1.47870	—
0.257	1.76038	—	1.53005	—	1.229	1.63926	—	—	—
0.2605	—	0.102	—	—	1.273	1.63849	—	1.47831	—
0.263	1.75343	—	1.52736	—	1.307	1.63789	—	—	—
0.265	—	0.096	—	—	1.320	1.63767	—	—	—
0.267	1.74864	—	1.52547	—	1.369	1.63681	—	—	—
0.270	—	0.096	—	—	1.396	1.63637	—	1.47789	—
0.274	1.74139	—	1.52261	—	1.422	1.63590	—	—	—
0.275	—	0.102	—	—	1.479	1.63490	—	—	—
0.2805	—	0.096	—	—	1.497	1.63457	—	1.47744	—
0.286	—	0.102	—	—	1.541	1.63381	—	—	—
0.291	1.72774	—	1.51705	—	1.6	—	0.05 ^f	—	—
0.2918	—	0.109	—	—	1.609	1.63261	—	—	—
0.2980	—	0.118	—	—	1.615	—	—	1.47695	—
0.303	1.71959	—	1.51365	—	1.682	1.63127	—	—	—
0.305	—	0.118	—	—	1.7	—	0.09	—	—
0.312	1.71425	0.096	1.51140	—	1.749	—	—	1.47638	—

TABLE 1 Refractive Indices^a and Absorption Coefficients^a for Calcite (*Continued*)

λ , μm	n_o	α_o	n_e	α_e	λ , μm	n_o	α_o	n_e	α_e
1.761	1.62974				2.4	—	2.3	—	0.09
1.8	—	0.16			2.5	—	2.7	—	0.14
1.849	1.62800				2.6	—	2.5	—	0.07
1.9	—	0.23			2.7	—	2.3	—	0.07
1.909	—	—	1.47573		2.8	—	2.3	—	0.09
1.946	1.62602				2.9	—	2.8	—	0.18
2.0	—	0.37			3.0	—	4.0	—	0.28
2.053	1.62372				3.1	—	6.7	—	0.46
2.100	—	0.62	1.47492	0.02 ^f	3.2	—	10.6	—	0.69
2.172	1.62099				3.3	—	15.0	—	0.92
2.2	—	1.1	—	0.05	3.324	—	—	1.47392	
2.3	—	1.7	—	0.07	3.4	—	19.0	—	1.2

^aRefractive indexes n_o and n_e are the ordinary and extraordinary rays, respectively, and the corresponding absorption coefficients are $\alpha_o = 4\pi k_o / \lambda \text{ cm}^{-1}$ and $\alpha_e = 4\pi k_e / \lambda \text{ cm}^{-1}$, where the wavelength λ is in centimeters. In the table, the wavelength is in micrometers.

^bUzan et al., Ref. 5; α_o and α_e were calculated from the reported values of k_o and k_e .

^cBallard et al., Ref. 6.

^dSchellman et al., Ref. 7; α_e was calculated from the optical density for the extraordinary ray.

^eBouriau and Lenoble, Ref. 8; reported absorption coefficient in this paper was for both o and e rays. α_o was calculated by assuming $\alpha_e = 0$.

^fBallard et al., Ref. 9.

negative uniaxial crystal. However, at wavelengths shorter than 1520 Å in the vacuum ultraviolet, the birefringence $n_e - n_o$ becomes positive, in agreement with theoretical predictions.^{5,10} For additional data in the 0.17- to 0.19- μm region, see Uzan et al.¹¹ The range of transparency of calcite is approximately from 0.214 to 3.3 μm for the extraordinary ray but only from about 0.23 to 2.2 μm for the ordinary ray.

If the principal plane of the e ray and the principal section coincide (Fig. 2), the wave normal (*but not the e ray*) obeys Snell's law, except that the index of refraction n_ϕ of this wave is given by^{12,13}

$$\frac{1}{n_\phi^2} = \frac{\sin^2 \phi}{n_e^2} + \frac{\cos^2 \phi}{n_o^2} \quad (1)$$

where ϕ is the angle between the direction of the *wave normal* and the optic axis ($\phi \leq 90^\circ$). When $\phi = 0^\circ$, $n_\phi = n_o$, and when $\phi = 90^\circ$, $n_\phi = n_e$. The angle of refraction for the wave normal is $\phi - \beta$, where β is the angle the normal to the surface makes with the optic axis. Snell's law for the extraordinary-ray *wave normal* then becomes

$$n \sin i = \frac{n_e n_o \sin(\phi - \beta)}{(n_o^2 \sin^2 \phi + n_e^2 \cos^2 \phi)^{1/2}} \quad (2)$$

where i is the angle of incidence of light in a medium of refractive index n . Since all other quantities in this equation are known, ϕ is uniquely determined but often must be solved for by iteration. Once ϕ is known, the angle of refraction r for the extraordinary ray can be determined as follows. If α is the angle the ray makes with the optic axis ($\alpha \leq 90^\circ$), then $r = \alpha - \beta$ and¹³

$$\tan \alpha = \frac{n_o^2}{n_e^2} \tan \phi \quad (3)$$

Although the angle of refraction of the extraordinary ray determines the path of the light beam through the prism, one must use the angle of refraction of the *wave normal*, $\phi - \beta$, in Fresnel's equation [Eq. (21) in Chap. 12, "Polarization"] when calculating the reflection loss of the e ray at the surface of the prism.

For the special case in which the optic axis is parallel to the surface as well as in the plane of incidence, α and ϕ are the complements of the angles of refraction of the ray and wave normal, respectively. If the light is normally incident on the surface, ϕ and α are both 90° and the extraordinary ray is undeviated and has its minimum refractive index n_o . In other cases for which the optic axis is not parallel to the surface, the extraordinary ray is refracted even for normal incidence.

If the plane of incidence is neither in a principal section nor perpendicular to the optic axis, it is more difficult to determine the angle of refraction of the extraordinary ray. In such cases, Huygens' construction is helpful.¹³⁻¹⁵

Types of Polarizing Prisms and Definitions

In order to make a polarizing prism out of calcite, some way must be found to separate the two polarized beams. In wavelength regions where calcite is absorbing (and hence only a minimum thickness of calcite can be used), this separation has been made simply by using a very thin calcite wedge cut so that the optic axis is parallel to the faces of the wedge to enable the e and o rays to be separated by a maximum amount. The incident light beam is restricted to a narrow pencil. Calcite polarizers of this type can be used at wavelengths as short as 1900 \AA .¹⁶ In more favorable wavelength regions, where the amount of calcite through which the light passes is not so critical, more sophisticated designs are usually employed. Such prisms can be divided into two main categories, *conventional polarizing prisms* (Secs. 13.3 and 13.4) and *polarizing beam-splitter prisms* (Sec. 13.5), and a third category, *Feussner prisms* (Sec. 13.6).

In conventional polarizing prisms, only light polarized in one direction is transmitted. This is accomplished by cutting and cementing the two halves of the prism together in such a way that the other beam suffers total internal reflection at the cut. It is usually deflected to the side, where it is absorbed by a coating containing a material such as lampblack. Since the ordinary ray, which has the higher index, is the one usually deflected, the lampblack is often mixed in a matching high-index binder such as resin of aloes ($n_D = 1.634$) or balsam of Tolu ($n_D = 1.628$) to minimize reflections.¹⁷ When high-powered lasers are used, the coating is omitted to avoid overheating the prism, and the light is absorbed externally.

Conventional polarizing prisms fall into two general categories: *Glan types* (Sec. 13.3) and *Nicol types* (Sec. 13.4), which are illustrated in Fig. 3. Glan types have the optic axis in the plane of the entrance face. If the principal section is parallel to the plane of the cut, the prism is a Glan-Thompson design (sometimes called a Glazebrook design); if perpendicular, a Lippich design; and if 45° , a Frank-Ritter design. In Nicol-type prisms, which include the various Nicol designs and the Hartnack-Prazmowsky, the principal section is perpendicular to the entrance face, but the optic axis is neither parallel nor perpendicular to the face.

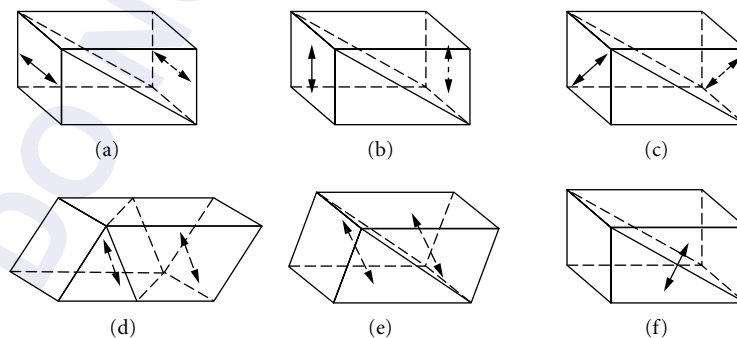


FIGURE 3 Types of conventional polarizing prisms. Glan types: (a) Glan-Thompson, (b) Lippich, and (c) Frank-Ritter; Nicol types: (d) conventional Nicol, (e) Nicol, Halle form, and (f) Hartnack-Prazmowsky. The optic axes are indicated by the double-pointed arrows.

Air-spaced prisms can be used at shorter wavelengths than cemented prisms, and special names have been given to some of them. An air-spaced Glan-Thompson prism is called a Glan-Foucault, and an air-spaced Lippich prism, a Glan-Taylor. In common practice, either of these may be called a Glan prism. An air-spaced Nicol prism is called a Foucault prism. Double prisms can also be made, thus increasing the prism aperture without a corresponding increase in length. Most double prisms are referred to as double Frank-Ritter, etc., but a double Glan-Thompson is called an Ahrens prism.

In polarizing beam-splitter prisms, two beams, which are polarized at right angles to each other, emerge but are separated spatially. The prisms have usually been used in applications for which both beams are needed, e.g., in interference experiments, but they can also be used when only one beam is desired. These prisms are also of two general types, illustrated in Fig. 10; those having the optic axis in the two sections of the prism perpendicular and those having them parallel. Prisms of the first type include the Rochon, Sénarmont, Wollaston, double Rochon, and double Sénarmont. Prisms of the second type are similar to the conventional polarizing prisms but usually have their shape modified so that the two beams emerge in special directions. Examples are the Foster, the beam-splitting Glan-Thompson, and the beam-splitting Ahrens.

The Feussner-type prisms, shown in Fig. 12, are made of isotropic material, and the film separating them is birefringent. For negative uniaxial materials the ordinary ray rather than the extraordinary ray is transmitted. These prisms have the advantage that much less birefringent material is required than for the other types of polarizing prisms, but they have a more limited wavelength range when calcite or sodium nitrate is used because, for these materials, the extraordinary ray is transmitted over a wider wavelength range than the ordinary ray.

The amount of flux which can be transmitted through a prism or other optical element depends on both its angular aperture and its cross-sectional area. The greater the amount of flux which can be transmitted, the better the *throughput* or *light-gathering power* (sometimes called *étendue* or *luminosity*) of the system.^{18,19} If a pupil or object is magnified, the convergence angle of the light beam is reduced in direct ratio to the increase in size of the image. The maximum throughput of a prism is thus proportional to the product of the prism's solid angle of acceptance and its cross-sectional area perpendicular to the prism axis. Hence, a large Glan-Taylor prism having an 8° field angle may, if suitable magnification is used, have a throughput comparable to a small Glan-Thompson prism with a 26° field angle. In general, to maximize prism throughput in an optical system, both the angular aperture and clear aperture (diameter of the largest circle perpendicular to the prism axis which can be included by the prism) should be as large as possible.

The quantities normally specified for a prism are its clear aperture, field angle, and length-to-aperture (L/A) ratio. The *semifield angle* is defined as the maximum angle to the prism axis* at which a ray can strike the prism and still be completely polarized *when the prism is rotated about its axis*. The field angle is properly twice the semifield angle.[†] (Some manufacturers quote a "field angle" for their polarizing prisms which is not symmetric about the prism axis and is thus in most cases unusable.) The *length-to-aperture (L/A) ratio* is the ratio of the length of the prism base (parallel to the prism axis) to the minimum dimension of the prism measured perpendicular to the prism base. For a square-ended prism, the L/A ratio is thus the ratio of prism length to width.

In determining the maximum angular spread a light beam can have and still be passed by the prism, both the field angle and the L/A ratio must be considered, as illustrated in Fig. 4. If the image of a point source were focused at the center of the prism, as in Fig. 4a, the limiting angular divergence of the beam would be determined by the field angle $2i$ of the prism.[‡] However, if an extended

*The prism axis, which is parallel to its base, is not to be confused with the optic axis of the calcite.

†In many prism designs, there is asymmetry about the prism axis, so that although light transmitted at a given angle to the prism axis may be completely polarized for one prism orientation, it will not be completely polarized when the prism is rotated about its axis. Thus, the semifield angle is not necessarily the largest angle at which completely polarized light can be transmitted by the prism in any orientation.

‡We are assuming that the prism is wide enough to ensure that the sides of the prism do not limit the angular width of the beam.

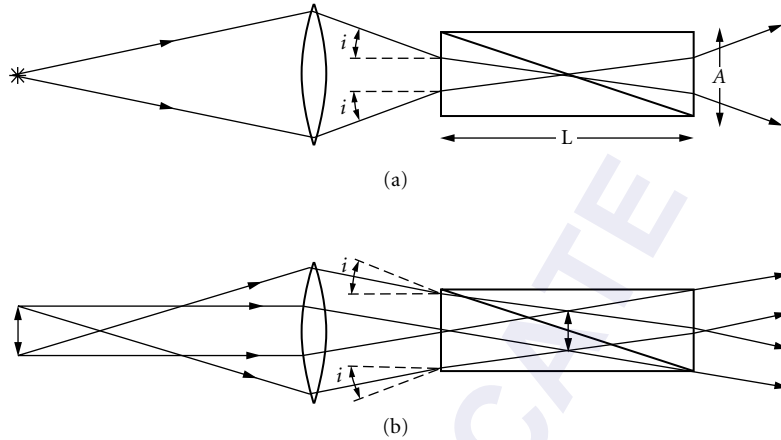


FIGURE 4 The effect of field angle and length-to-aperture ratio of a prism polarizer on the maximum angular beam spread for (a) a point source and (b) an extended source. The field angle is $2i$, and $L/A = 3$. The field angle is exaggerated for clarity.

source were focused there (Fig. 4b), the limiting angular divergence would be determined by the L/A ratio, not the field angle.

The field angle of a polarizing prism is strongly wavelength-dependent. For example, a Glan prism having an 8° field angle at $0.4 \mu\text{m}$ has only a 2° field angle at $2 \mu\text{m}$. In designing optical systems in which polarizing prisms are to be used, the designer must allow for this variation in field angle. If he does not, serious systematic errors may occur in measurements made with the system.

13.3 GLAN-TYPE PRISMS

Most prisms used at the present time are of the Glan type. Although they require considerably more calcite than Nicol types of comparable size, they are optically superior in several ways: (1) Since the optic axis is perpendicular to the prism axis, the index of the extraordinary ray differs by a maximum amount from that of the ordinary ray. Thus, a wider field angle or a smaller L/A ratio is possible than with Nicol types. (2) The light is nearly uniformly polarized over the field; it is not for Nicol types. (3) There is effectively no lateral displacement in the apparent position of an axial object viewed through a (perfectly constructed) Glan-type prism. Nicol types give a lateral displacement. (4) Since off-axis wander results in images which have astigmatism when the prism is placed in a converging beam, Glan types have slightly better imaging qualities than Nicol types.

Two other often-stated advantages of Glan-type prisms over Nicol types appear to be fallacious. One is that the slanting end faces of Nicol-type prisms have higher reflection losses than the square-ended faces of Glan types. Since the extraordinary ray vibrates in the plane of incidence and hence is in the p direction, increasing the angle of incidence toward the polarizing angle should decrease the reflection loss. However, the index of refraction for the extraordinary ray is higher in Nicol-type prisms (Glan types have the minimum value of the extraordinary index), so the reflection losses are actually almost identical in the two types of prisms. The second "advantage" of Glan-type prisms is that the slanting end faces of the Nicol type supposedly induce elliptical polarization. This widely stated belief probably arises because in converging light the field in Nicol-type polarizers is not uniformly polarized, an effect which could be misinterpreted as ellipticity (see "Landolt Fringe" in Sec. 13.4). It is possible that strain birefringence could be introduced in the surface layer of a calcite prism by some optical polishing techniques resulting in

ellipticity in the transmitted light, but there is no reason why Nicol-type prisms should be more affected than Glan types.

Glan-Thompson-Type Prisms

Glan-Thompson-type prisms may be either cemented or air-spaced. Since, as was mentioned previously, an air-spaced Glan-Thompson-type prism is called a Glan-Foucault or simply a Glan prism,* the name Glan-Thompson prism implies that the prism is cemented. Both cemented and air-spaced prisms, however, have the same basic design. The cemented prisms are optically the better design for most applications and are the most common type of prisms in use today. The Glan-Thompson prism is named for P. Glan,²⁰ who described an air-spaced Glan-Thompson-type prism in 1880, and for S. P. Thompson,²¹ who constructed a cemented version in 1881 and modified it to its present square-ended design in 1882.²² These prisms are also sometimes called Glazebrook prisms because R. T. Glazebrook²³ demonstrated analytically in 1883 that when rotated about its axis, this prism gives the most uniform rotation of the plane of polarization for a conical beam of incident light. The cut in a Glan-Thompson-type prism is made parallel to the optic axis, which may either be parallel to two sides, as in Fig. 3a, or along a diagonal. The end faces are always perpendicular to the axis of the prism and contain the optic axis.

The extinction ratio[†] obtainable with a good Glan-Thompson-type prism equals or exceeds that of any other polarizer. Ratios of 5 parts in 100,000 to 1 part in 1 million can be expected although values as high as 1 part in 3×10^7 have been reported for small selected apertures of the prism.²⁴ The small residuals result mainly from imperfections in the calcite or from depolarization by scattering from the prism faces,²⁴ although if the optic axis is not strictly in the plane of the end face, or if the optic axes in the two halves of the prism are not accurately parallel, the extinction ratio will be reduced. Also, the extinction ratio may depend strongly upon which end of the prism the light is incident. When prisms are turned end for end, changes in the extinction ratio of as much as a factor of 6 have been reported.²⁴

When measuring the extinction ratio, it is essential that none of the unwanted ordinary ray, which is internally reflected at the interface and absorbed or scattered at the blackened side of the prism, reach the detector. King and Talim²⁵ found that they had to use two 4-mm-diameter apertures and a distance of 80 mm between the photomultiplier detector and prism to eliminate the o-ray scattered light. With no limiting apertures and a 20-mm distance, their measured extinction ratio was in error by a factor of 80.

The field angle of the prism depends both on the cement used between the two halves and on the angle of the cut, which is determined by the L/A ratio. Calculation of the field angle is discussed in "Field Angle" section on p. 13.12 and by Bennett and Bennett.¹ Very large field angles can be obtained with Glan-Thompson prisms. For example, if the L/A ratio is 4, the field angle can be nearly 42°. Normally, however, smaller L/A ratios are used. The most common types of cemented prisms are the long form, having an L/A ratio of 3 and a field angle of 26°, and the short form, having an L/A ratio of 2.5 and a field angle of 15°.

Transmission In Fig. 5 the transmission of a typical Glan-Thompson prism is compared with curves for a Glan-Taylor prism and a Nicol prism. The Glan-Thompson is superior over most of the range, but its transmission decreases in the near ultraviolet, primarily because the cement begins to absorb. Its usable transmission range can be extended to about 2500 Å by using an ultraviolet-transmitting cement. Highly purified glycerin, mineral oil, castor oil, and Dow Corning DC-200 silicone oil, which because of its high viscosity is not as subject to seepage as lighter oils, have been used as cements in the ultraviolet, as have dextrose, glucose, and *gédamine* (a urea formaldehyde resin in butyl alcohol).

*An air-spaced Lippich prism, the Glan-Taylor (see "Glan-Taylor Prism" section on p. 13.12), has similar optical properties to the Glan-Foucault prism but better transmission. It is also called a Glan prism.

[†]The extinction ratio is the ratio of the maximum to the minimum transmittance when a polarizer is placed in a plane polarized beam and is rotated about an axis parallel to the beam direction.

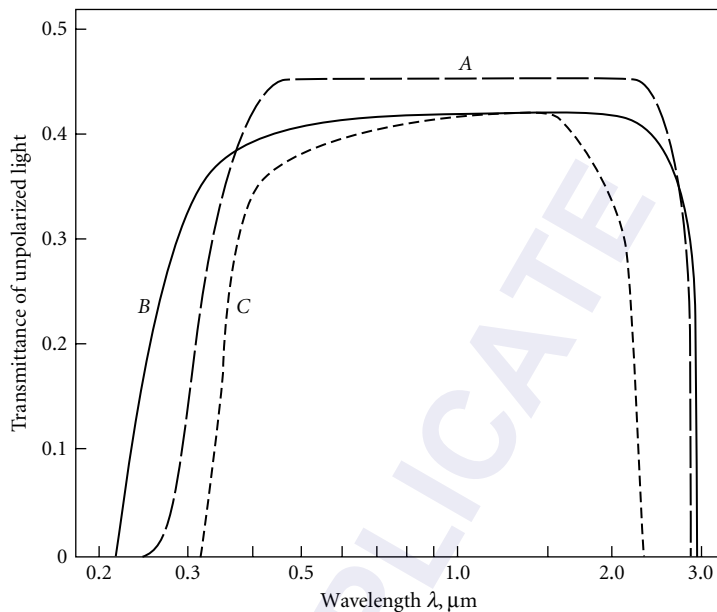


FIGURE 5 Transmittance curves for typical polarizing prisms: A, Glan-Thompson, B, Glan-Taylor, and C, Nicol prism. (Measured by D. L. Decker, *Michelson Laboratory*.) In the visible and near-infrared regions the Glan-Thompson has the best energy throughput. In the near ultraviolet the Glan-Thompson may still be superior because the Glan-Taylor has such an extremely small field angle that it may cut out most of the incident beam.

Transmission curves for 1-mm thicknesses of several of these materials are shown in Fig. 6, along with the curve for Canada balsam, a cement formerly widely used for polarizing prisms in the visible region.⁸ *Gédamine*, one of the best of the ultraviolet-transmitting cements, has an index of refraction $n_D = 1.465$, and can be fitted to the dispersion relation⁸

$$n = 1.464 + \frac{0.0048}{\lambda^2} \quad (4)$$

where the wavelength λ is in micrometers.

Figure 7 shows ultraviolet transmission curves for Glan-Thompson prisms with L/A ratios of 2.5 and 3 which are probably cemented with *n*-butyl methacrylate, a low-index polymer that has largely replaced Canada balsam. Better ultraviolet transmission is obtained with a Glan-Thompson prism cemented with DC-200 silicone oil. Air-spaced prisms can be used to nearly 2140 Å in the ultraviolet, where calcite begins to absorb strongly. Transmission curves for two such prisms are shown in Fig. 7. The Glan-Taylor, which is an air-spaced prism of the Lippich design, has a higher ultraviolet transmission than the Glan-Foucault, an air-spaced Glan-Thompson prism. The reason for this difference is that multiple reflections occur between the two halves of the Glan-Foucault prism, resulting in a lowered transmission, but are largely absent in the Glan-Taylor design (see “Glan-Taylor Prism” section on p. 13.12).

The infrared transmission limit of typical Glan-Thompson prisms is about 2.7 μm although they have been used to 3 μm.²⁶ The same authors report using a 2.5-cm-long Glan-Thompson prism in the 4.4- to 4.9-μm region.

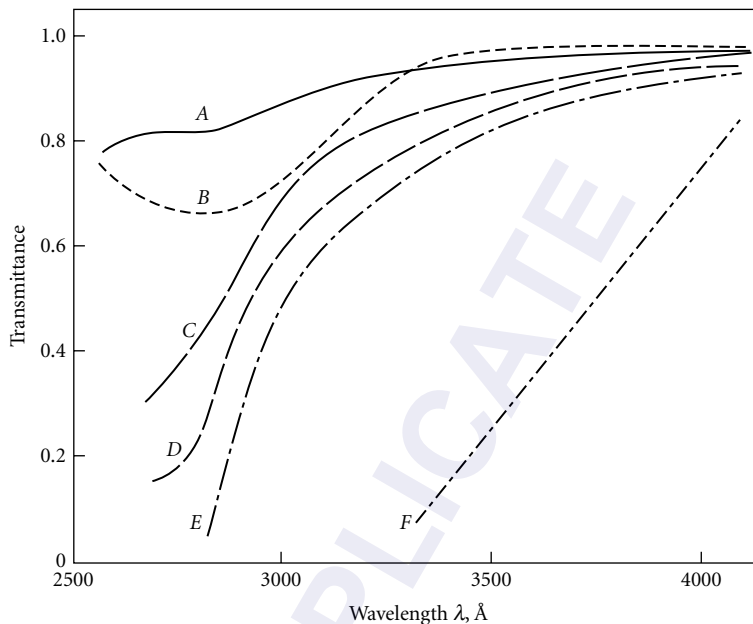


FIGURE 6 Transmittance curves for 1-mm thicknesses of various cements: A, crystalline glucose, B, glycerine, C, gédamine (urea formaldehyde resin in butyl alcohol), D, Rhodopas N60A (polymerized vinyl acetate in alcohol), E, urea formaldehyde, and F, Canada balsam. (Modified from Bouriau and Lenoble.⁸) The transmittance of these materials is adequate at longer wavelengths.

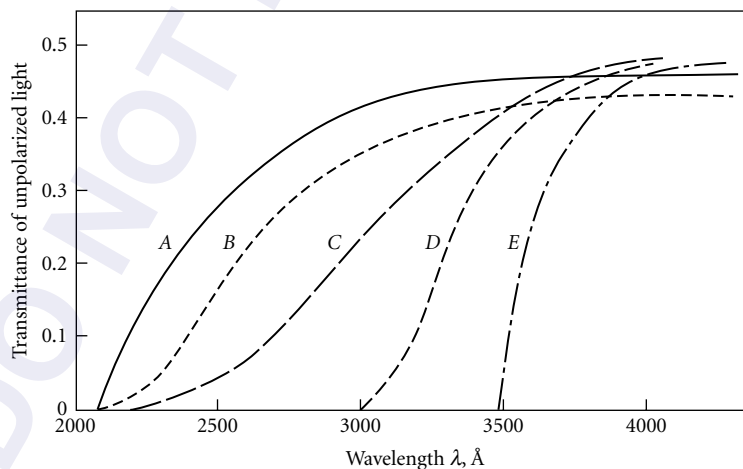


FIGURE 7 Ultraviolet transmittance curves for various Glan-Thompson and air-spaced prisms: A, Glan-Taylor (air-spaced Lippich-type prism), B, Glan-Foucault (air-spaced Glan-Thompson prism), C, Glan-Thompson prism with L/A ratio of 2 cemented with DC-200 silicone oil, D, Glan-Thompson prism with L/A ratio of 2.5 probably cemented with *n*-butyl methacrylate, and E, Glan-Thompson prism similar to D except with $L/A = 3$. (Modified from curves supplied by Karl Lambrecht Corporation, Chicago.)

Field Angle Since many prism polarizers are used with lasers that have parallel beams of small diameter, field-angle effects are not as important as previously when extended area sources were used. Extensive calculations of the field angles for a Glan-Thompson prism are included in the earlier polarization chapter.¹

Other Glan-Thompson-Type Prisms Other types of Glan-Thompson-type prisms include the Ahrens prism (two Glan-Thompson prisms placed side-by-side), Glan-Foucault prism (an air-spaced Glan-Thompson prism), Grosse prism (an air-spaced Ahrens prism), and those constructed of glass and calcite. Information about these prisms can be found in the earlier polarization chapter.¹

Lippich-Type Prisms

Lippich²⁷ (1885) suggested a polarizing-prism design similar to the Glan-Thompson but with the optical axis in the entrance face and at right angles to the intersection of the cut with the entrance face (Fig. 3*b*).^{*} For this case, the index of refraction of the extraordinary ray is a function of angle of incidence and can be calculated from Eq. (1) after ϕ , the complement of the angle of refraction of the wave normal is determined from Eq. (2). In the latter equation, β , the angle normal to the surface makes with the optic axis, is 90° since the optic axis is parallel to the entrance face. Since the directions of the ray and the wave normal no longer coincide, the ray direction must be calculated from Eq. (3). Lippich prisms are now little-used because they have small field angles, except for two; the air-spaced Lippich, often called a Glan-Taylor prism, and the Marple-Hess prism (two Glan-Taylor prisms back-to-back) that is described in “Marple-Hess Prism” section on p. 13.13. Further information about all Lippich-type prisms is given in the earlier polarization chapter.¹

Glan-Taylor Prism The Glan-Taylor prism, first described in 1948 by Archard and Taylor,²⁹ has substantial advantages over its Glan-Thompson design counterpart, the Glan-Foucault prism (see “Other Glan-Thompson-Type Prisms” section earlier). Since air-spaced prisms have a very small field angle, the light must be nearly normally incident on the prism face, so that the difference in field angles between the Glan-Taylor and Glan-Foucault prisms (caused by the difference in the refractive index of the extraordinary ray) is negligible.

The major advantages of the Glan-Taylor prism are that its calculated transmission is between 60 and 100 percent higher than that of the Glan-Foucault prism and the intensity of multiple reflections between the two sides of the cut always a principal drawback with air-spaced prisms, is reduced to less than 10 percent of the value for the Glan-Foucault prism.

The calculated and measured transmittances of a Glan-Taylor prism are in reasonable agreement, but the measured transmittance of a Glan-Foucault prism (Fig. 7) may be considerably higher than its theoretical value.²⁹ Even so, the transmission of the Glan-Taylor prism is definitely superior to that of the Glan-Foucault prism, as can be seen in Fig. 7. Extinction ratios of better than 1 part in 10^3 are obtainable for the Glan-Taylor prism.³⁰

A final advantage of the Glan-Taylor prism is that it can be cut in such a way as to conserve calcite. Archard and Taylor²⁹ used the Ahrens method of spar cutting described by Thompson²² and found that 35 percent of the original calcite rhomb could be used in the finished prism.

In a modified version of the Glan-Taylor prism becoming popular for laser applications, the cut angle[†] is increased, the front and back faces are coated with antireflection coatings, and portions of the sides are either covered with absorbing black glass plates or highly polished to let the unwanted beams escape.³⁰ The effect of increasing the cut angle is twofold: a beam normally incident on the prism face will have a smaller angle of incidence on the cut and hence a smaller reflection loss at

^{*}The Lippich prism should not be confused with the Lippich half-shade prism, which is a device to determine a photometric match point. The half-shade prism consists of a Glan-Thompson or Nicol prism placed between the polarizer and analyzer such that it intercepts half the beam and is tipped slightly in the beam. The prism edge at the center of the field is highly polished to give a sharp dividing line. The eye is focused on this edge; the disappearance of the edge gives the photometric match point.²⁸

[†]The cut angle is the acute angle the cut makes with the prism base.

the cut than a standard Glan-Taylor prism, but, at the same time, the semifield angle will be reduced throughout most of the visible and near-infrared regions.

A new type of air-spaced prism³¹ has a very high transmittance for the extraordinary ray. It resembles the Glan-Taylor prism in that the optic axis is parallel to the entrance face and at right angles to the intersection of the cut with the entrance face. However, instead of striking the prism face at normal incidence, the light is incident at the Brewster angle for the extraordinary ray (54.02° for the $6328\text{-}\text{\AA}$ helium-neon laser wavelength), so that there is no reflection loss for the e ray at this surface. Since the ordinary ray is deviated about 3° more than the extraordinary ray and its critical angle is over 4° less, it can be totally reflected at the cut with tolerance to spare while the extraordinary ray can be incident on the cut at only a few degrees beyond its Brewster angle. Thus this prism design has the possibility of an extremely low light loss caused by reflections at various surfaces. A prototype had a measured transmission of 0.985 for the extraordinary ray at 6328 \AA .³¹ If the prism is to be used with light sources other than lasers, its semifield angle can be calculated.¹

A major drawback to the Brewster angle prism is that since the light beam passes through a plane-parallel slab of calcite at nonnormal incidence, it is displaced by an amount that is proportional to the total thickness of the calcite. Some of the prisms are made with glass in place of calcite for the second element. In this case, the beam will usually be deviated in addition to being displaced. Measurements on a calcite-glass prototype at 6328 \AA showed that the output beam was laterally displaced by several millimeters with an angular deviation estimated to be less than 0.5° .³¹

Marple-Hess Prism If a larger field angle is required than can be obtained with a Glan-Taylor prism, a Marple-Hess prism may be used. This prism, which was first proposed in 1960 as a double Glan-Foucault by D. T. F. Marple of the General Electric Research Laboratories and modified to the Taylor design by Howard Hess of the Karl Lambrecht Corporation,³² is effectively two Glan-Taylor prisms back-to-back. The analysis for this prism is made in the same way as for the Glan-Taylor prism (see "Glan-Taylor Prism" section earlier) and Lippich-type prisms in general, keeping in mind that the refractive index of the "cement" is 1 since the components are air-spaced.

Since the ordinary ray is totally reflected for all angles of incidence by one or the other of the two cuts, the field angle is symmetric about the longitudinal axis of the prism and is determined entirely by the angle at which the extraordinary ray is totally reflected at one of the two cuts. This angle can be readily calculated.¹ The field angle is considerably larger than for the Glan-Foucault or Glan-Taylor prism and does not decrease as the wavelength increases.

Unlike the Glan-Foucault or Glan-Taylor prisms, which stop being efficient polarizers when the angle of incidence on the prism face becomes too large, the Marple-Hess prism continues to be an efficient polarizer as long as the axial ordinary ray is not transmitted. If the prism is used at a longer wavelength than the longest one for which it was designed (smallest value of n_o used to determine the cut angle), the value of n_o will be still smaller and the critical angle for the axial ordinary ray will not be exceeded. Thus the axial o ray will start to be transmitted before off-axis rays get through. When this situation occurs, it only makes matters worse to decrease the convergence angle. Thus, there is a limiting long wavelength, depending on the cut angle, beyond which the Marple-Hess prism is not a good polarizer. At wavelengths shorter than the limiting wavelength, the Marple-Hess prism has significant advantages over other air-spaced prism designs.

It is not easy to make a Marple-Hess prism, and the extinction ratio in the commercial model is given as between 1×10^{-4} and 5×10^{-5} , somewhat lower than for a Glan-Taylor prism.³⁰ On the other hand, even though the Marple-Hess prism has an increased L/A ratio, 1.8 as compared to 0.85 for a Glan-Taylor prism, its ultraviolet transmission is still superior to commercially available ultraviolet-transmitting Glan-Thompson prisms of comparable aperture.

Frank-Ritter-Type Prisms

The third general category of Glan-type polarizing prisms is the Frank-Ritter design. Prisms of this type are characterized by having the optic axis in the plane of the entrance face, as in other Glan-type prisms, but having the cut made at 45° to the optic axis (Fig. 3c) rather than at 0° , as in Glan-Thompson

prisms, or at 90° , as in Lippich prisms. Frank-Ritter prisms are particularly popular in the Soviet Union, and over 80 percent of the polarizing prisms made there have been of this design.³³ Usually double prisms comparable to the Ahrens modification of the Glan-Thompson are used,¹ primarily because from a rhombohedron of Iceland spar two Frank-Ritter double prisms can be obtained but only one Ahrens of the same cross-section or one Glan-Thompson of smaller cross-section.³³ However, this apparent advantage can be illusory since Iceland spar crystals often are not obtained as rhombs. For example, if the natural crystal is in the form of a plate, it may be less wasteful of material to make a Glan-Thompson or Ahrens prism than a Frank-Ritter prism.³³

Optically Frank-Ritter prisms should be similar to Glan-Thompson and Ahrens types, although the acceptance angle for a given L/A ratio is somewhat smaller since the refractive index of the extraordinary ray is larger than n_e in the prism section containing the longitudinal axis and perpendicular to the cut. In practice, the degree of polarization for a Frank-Ritter prism seems to be quite inferior to that of a good Glan-Thompson or even an Ahrens prism.³³

Use of Glan-Type Prisms in Optical Systems

Several precautions should be taken when using Glan-type prisms in optical systems: (1) the field angle of the prism should not be exceeded, (2) there should be an adequate entrance aperture so that the prism does not become the limiting aperture of the optical system, and (3) baffles should be placed preceding and following the prism to avoid incorrect collection of polarized light or extraneous stray light. The reason why these precautions are important are discussed in the earlier polarization chapter.¹

Common Defects and Testing of Glan-Type Prisms

Several common defects are found in the construction of Glan-type prisms and limit their performance:

1. The axial beam is displaced as the prism is rotated. This defect called *squirm*, results when the optic axes in the two halves of the prism are not strictly parallel. A line object viewed through the completed prism will oscillate as the prism is turned around the line of sight.³⁴

2. The axial ray is deviated as the prism is rotated. This defect is caused by the two prism faces not being parallel. A residual deviation of 3 minutes of arc is a normal tolerance for a good Glan-Thompson prism; deviations of 1 minute or less can be obtained on special order.

3. The optic axis does not lie in the end face. This error is often the most serious, since if the optic axis is not in the end face and the prism is illuminated with convergent light, the planes of vibration of the transmitted light are no longer parallel across the face of the prism. This effect, which in Nicol-type prisms gives rise to the Landolt fringe, is illustrated in the following practical case.³⁵ For a convergent beam of light of semicone angle i , the maximum variation of the plane of vibration of the emergent beam is $\pm\gamma$, where, approximately,

$$\tan \gamma = n_e \sin i \tan \phi \quad (5)$$

and ϕ is the angle of inclination of the optic axis to the end face, caused by a polishing error. For $i = 3^\circ$ and $p = 5^\circ$, the plane of vibration of the emergent beam varies across the prism face by ± 23 minutes of arc. Thus, good extinction cannot be achieved over the entire aperture of this prism even if nearly parallel light is incident on it. The field angle is also affected if the optic axis is not in the end face or is not properly oriented in the end face, but these effects are small.

4. The cut angle is incorrect or is different in the two halves of the prism. If the cut angle is slightly incorrect, the field angle may be decreased. This error is particularly important in Glan-Foucault or Glan-Taylor prisms, for which the angular tolerances are quite severe, and a small change in cut angle for these prisms may greatly alter the field angle, as discussed in "Glan-Taylor Prism" section on p. 13.12 and Ref. 1. If the cut angles are different in the two halves of the prism, the field angle will change when the prism is turned end-for-end. The field angle is determined by the cut angle in the half of the prism toward the incident beam. Differences in the two cut angles may also cause a beam deviation. If the angles in the two halves differ by a small angle α that makes the end

faces nonparallel, the beam will be deviated by an angle $\delta = \alpha(n_c - 1)$.³⁵ If instead, the end faces are parallel and the difference in cut angle is taken up by the cement layer which has a refractive index of approximately n_c , there will be no deviation. However, if the prism is air-spaced, the deviation δ' caused by a nonparallel air film is approximately $\delta' = \alpha n_c$, illustrating one reason why air-spaced prisms are harder to make than conventional Glan-Thompson prisms.³⁵

5. The transmittance is different when the prism is rotated through 180° . A potentially more serious problem when one is making photometric measurements is that the transmission of the prism may not be the same in two orientations exactly 180° apart.³⁶ This effect may be caused by the presence of additional light outside the entrance or exit field angle, possibly because of strain birefringence in the calcite.

Two factors which limit other aspects of polarizer performance in addition to the extinction ratio are *axis wander*, i.e., variation of the azimuth of the transmitted beam over the polarizer aperture, and the ellipticity of the emergent polarized beams²⁵ caused by material defects in the second half of the prism. Further details are discussed in the earlier polarization chapter.¹

In order to determine the cut angle, field angle, parallelism of the prism surfaces, thickness and parallelism of the air film or cement layer, and other prism parameters, one can use the testing procedures outlined by Decker et al.,³⁷ which require a spectrometer with a Gauss eyepiece, laser source, and moderately good polarizer. (Other testing procedures have been suggested by Archard.³⁵) Rowell et al.³⁸ have given a procedure for determining the absolute alignment of a prism polarizer. However, they failed to consider some polarizer defects, as pointed out by Aspnes³⁹ who gives a more general alignment procedure that compensates for the prism defects. (There is also a response from Rowell.⁴⁰) Further information about testing Glan-type prisms and reasons why prism errors are important can be found in the earlier polarization chapter.¹

13.4 NICOL-TYPE PRISMS

Nicol-type prisms are not generally used at the present time, as Glan types are optically preferable. However, they were the first kind made and were once so common that Nicol became a synonym for polarizer. There is much more calcite wastage in making Glan-type prisms than in making the simpler Nicol types so that, even though Glan polarizers were developed in the nineteenth century, it was only following the recent discoveries of new calcite deposits that they became popular. Many of the older instruments are still equipped with Nicol prisms so they will be briefly described here.

Conventional Nicol Prism

The first polarizing prism was made in 1828 by William Nicol⁴¹ a teacher of physics in Edinburgh. By cutting a calcite rhomb diagonally and symmetrically through its blunt corners and then cementing the pieces together with Canada balsam, he could produce a better polarizer than any known up to that time. A three-dimensional view of Nicol's prism is shown in Fig. 3*d*. The cut is made perpendicular to the principal section (defined in "Double Refraction in Calcite" in Sec. 13.2), and the angle is such that the ordinary ray is totally reflected and only the extraordinary ray emerges. When the rhomb is intact, the direction of polarization can be determined by inspection. However, the corners are sometimes cut off, making the rhomb difficult to recognize.

The principal section of Nicol's original prism is similar to that shown in Fig. 2 except that the ordinary ray is internally reflected at the cut along diagonal *BH*. The cut makes an angle of $19^\circ 8'$ with edge *BF* in Fig. 2 and an angle of about 90° with the end face of the rhomb. Since the obtuse angle is $109^\circ 7'$ (Fig. 3*d*), the angle between the cut and the optic axis is $44^\circ 36'$. The field of the prism is limited on one side by the angle at which the ordinary ray is no longer totally reflected from the balsam film, about 18.8° from the axis of rotation of the prism, and on the other by the angle at which the extraordinary ray is totally reflected by the film, about 9.7° from the axis. Thus the total angle is about 28.5° but is not by any means symmetric about the axis of rotation; the field angle (see "Types of Polarizing Prisms and Definitions" in Sec. 13.2) is only $2 \times 9.7^\circ = 19.4^\circ$.

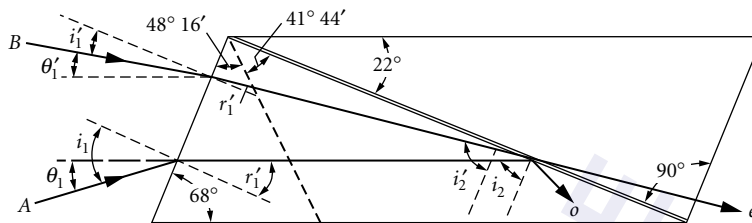


FIGURE 8 Principal section of a conventional Nicol prism with slightly trimmed end faces. Ray A gives the limiting angle θ_1 beyond which the ordinary ray is no longer totally internally reflected at the cut; ray B gives the limiting angle θ'_1 for which the extraordinary ray starts to be totally internally reflected at the cut.

In order to produce a somewhat more symmetric field and increase the field angle, the end faces of Nicol prisms are usually trimmed to an angle of 68° . This practice was apparently started by Nicol himself.²² If the cut is made at 90° to the new face, as shown in Fig. 8, the new field angle is twice the smaller of θ_1 and θ'_1 . The field angles are computed as described in the earlier polarization chapter.¹

Trimmed Nicol-Type Prisms

The angle at which the cut is made in a Nicol-type prism is not critical. The field angle is affected, but a useful prism will probably result even if the cut is made at an angle considerably different from 90° . The conventional trimmed Nicol, discussed in “Conventional Nicol Prism” section earlier, is shown again in Fig. 9a. In this and the other five parts of the figure, principal sections of various prisms are

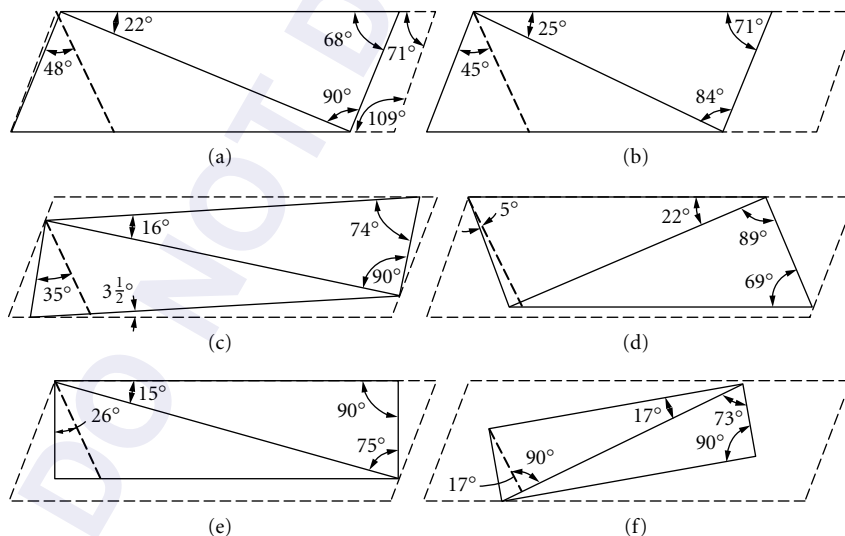


FIGURE 9 Principal sections of various types of trimmed cemented Nicol prisms shown superimposed on the principal section of a cleaved calcite rhomb (see Fig. 2): (a) conventional trimmed Nicol; (b) Steeg and Reuter shortened Nicol (*Thompson*²²); (c) Ahrens Nicol (*Thompson*²²); (d) Thompson reversed Nicol (*Thompson*⁴²); (e) square-ended Nicol; and (f) Hartnack-Prazmowski reversed Nicol. In all cases, the angle between the prism face and the optic axis (*heavy dashed line*), the angle of the cut, and the acute angle of the rhomb are indicated.

shown superimposed on the principal section of the basic calcite rhomb (Fig. 2). Thus, it is clear how much of the original rhomb is lost in making the different types of trimmed Nicols.

In the Steeg and Reuter Nicol shown in Fig. 9b, the rhomb faces are not trimmed, and the cut is made at 84° to the faces instead of 90° , giving a smaller L/A ratio. The asymmetry of the field which results is reduced by using a cement having a slightly higher index than Canada balsam.

Alternately, in the Ahrens Nicol shown in Fig. 9c, the ends are trimmed in the opposite direction, increasing their angles with the long edges of the rhomb from $70^\circ 53'$ to $74^\circ 30'$ or more. By also trimming the long edges by $3^\circ 30'$, the limiting angles are made more symmetric about the prism axis.

Thompson Reversed Nicol In the Thompson reversed Nicol shown in Fig. 9d, the ends are heavily trimmed so that the optic axis lies nearly in the end face. As a result, the blue fringe is thrown farther back than in a conventional Nicol, and although the resulting prism is shorter, its field angle is actually increased.

Nicol Curtate, or Halle, Prism The sides of the calcite rhomb may also be trimmed so that they are parallel or perpendicular to the principal section. Thus, the prism is square (or sometimes octagonal). This prism is of the Halle type^{43,44} and was shown in Fig. 3e. Halle, in addition, used thickened linseed oil instead of Canada balsam and altered the angle of the cut. In this way he reduced the length-to-aperture ratio from about 2.7 to 1.8 and the total acceptance angle from 25° to about 17° . Such shortened prisms cemented with low-index cements are often called Nicol curtate prisms (curtate means shortened).

Square-Ended Nicol The slanting end faces on conventional Nicol prisms introduce some difficulties, primarily because the image is slightly displaced as the prism is rotated. To help correct this defect, the slanting ends of the calcite rhomb can be squared off, as in Fig. 9e, producing the so-called square-ended Nicol prism. The angle at which the cut is made must then be altered since the limiting angle θ_1 for an ordinary ray depends on the angle of refraction at the end face in a conventional prism, in which the limiting ray travels nearly parallel to the prism axis inside the prism (ray A in Fig. 8). If the cut remained the same, the limiting value of θ_1 would thus be zero. However, if the cut is modified to be 15° to the sides of the prism, the total acceptance angle is in the 24° to 27° range, depending on the type of cement used.²²

Some image displacement will occur even in square-ended Nicol prisms since the optic axis is not in the plane of the entrance face. Therefore, the extraordinary ray will be bent even if light strikes the entrance face of the prism at normal incidence. There is considerable confusion on this point in the literature.^{22,45}

Hartnack-Prazmowski Prism A reversed Nicol which has the cut at 90° to the optic axis⁴⁶ is shown in Figs. 3f and 9f. If it is cemented with linseed oil, the optimum cut angle calculated by Hartnack is 17° to the long axis of the prism, giving a total acceptance angle of 35° and an L/A ratio of 3.4.²² If Canada balsam is used, the cut should be 11° , in which case the total acceptance angle is 33° and the L/A ratio is 5.2.

Foucault Prism A modified Nicol prism in which an air space is used between the two prism halves instead of a cement layer⁴⁷ consists of a natural-cleavage rhombohedron of calcite which has been cut at an angle of 51° to the face. The cut nearly parallels the optic axis. Square-ended Foucault-type prisms, such as the Hofmann prism, have also been reported.²² The angle at which the cut is made can be varied slightly in both the normal Foucault prism and the Hofmann variation of it. In all designs the L/A ratio is 1.5 or less, and the total acceptance angle about 8° or less. The prisms suffer somewhat from multiple reflections, but the principal trouble, as with all Nicol prisms, is that the optic axis is not in the plane of the entrance face. This defect causes various difficulties, including nonuniform polarization across the field and the occurrence of a Landolt fringe (discussed next and Ref. 1) when two Nicol-type prisms are crossed.

Landolt Fringe If an intense extended light source is viewed through crossed polarizing prisms, careful observation will reveal that the field is not uniformly dark. In Nicol-type prisms the darkened field is crossed by a darker line whose position is an extremely sensitive function of the angle between the polarizer and analyzer. Other types of polarizing prisms also exhibit this anomaly but to a lesser extent. The origin of the Landolt fringe is given in the earlier polarization chapter¹ and the references cited therein.

13.5 POLARIZING BEAM-SPLITTER PRISMS

The three classic polarizing beam-splitter prisms are the Rochon, Sénarmont, and Wollaston, shown in perspective in Fig. 10a to 10c and in side view in Fig. 11a to 11c. In addition, any polarizing prism can be used as a polarizing beam splitter by changing the shape of one side and removing the absorbing coating from its surface. Two examples of such prisms are the Foster prism, in which the ordinary and extraordinary rays emerge at right angles to each other, and the beam-splitting Glan-Thompson prism, in which the ordinary ray emerges normal to one side (Figs. 10d and e and 11d and e). Another prism of this type, the beam-splitting Ahrens prism, is a double beam-splitting Glan-Thompson prism (see “Other Glan-Thompson-Type Prisms” in Sec. 13.3).

In polarizing prisms, the optic axes are always parallel to each other in the two halves of the prism. By contrast, the optic axes in the two halves of the Rochon, Sénarmont, and Wollaston polarizing beam-splitter prisms are at right angles to each other. Crystal quartz is often used to make these beam splitters, and such prisms can be used down to the vacuum ultraviolet. In applications not requiring such short wavelengths, calcite is preferable because it gives a greater angular separation of the beams (typically 10° as compared to 0.5° for quartz) and does not produce optical rotation.

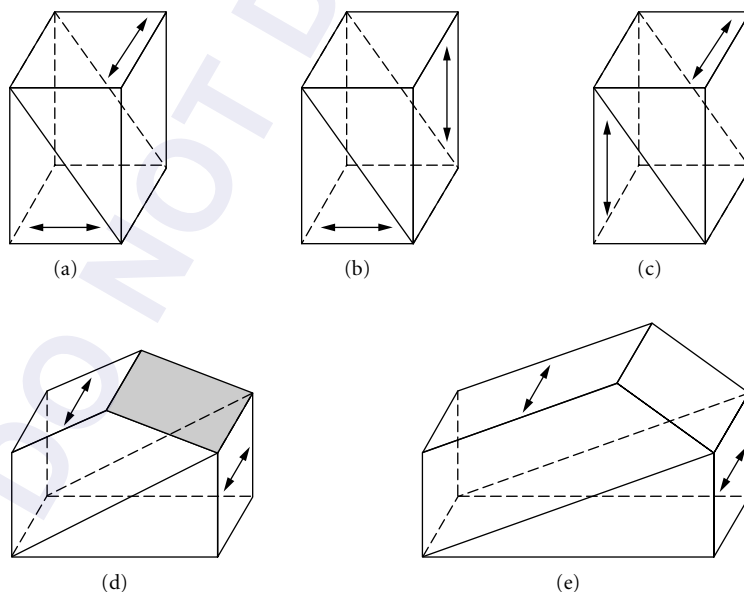


FIGURE 10 Three-dimensional views of various types of polarizing beam-splitter prisms: (a) Rochon; (b) Sénarmont; (c) Wollaston; (d) Foster (shaded face is silvered); and (e) beam-splitting Glan-Thompson.

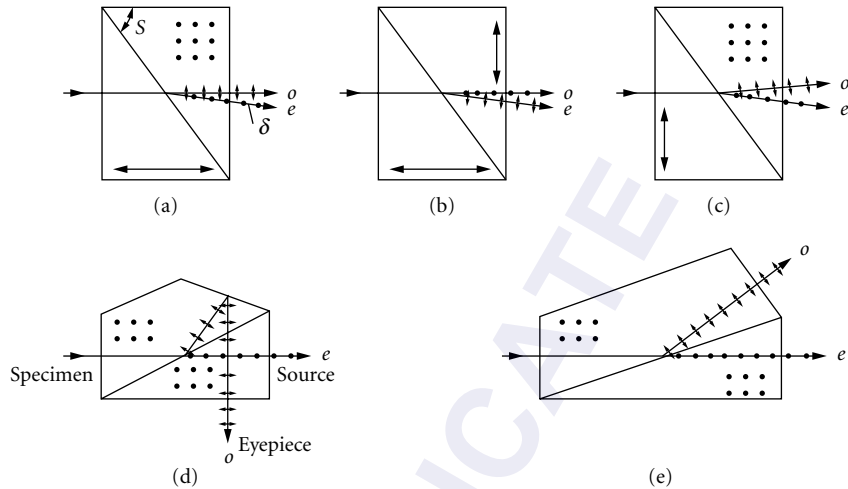


FIGURE 11 Side views of the polarizing beam-splitter prisms in Fig. 10. The directions of the optic axes are indicated by the dots and the heavy double-pointed arrows. The angle of the cut for the Rochon prism is S . When the Foster prism is used as a microscope illuminator, the source, specimen, and eyepiece are in the positions indicated.

Rochon Prism

The Rochon prism, invented in 1783,⁴⁸ is the most common type of polarizing beam splitter. It is often used in photometric applications in which both beams are utilized. It is also used as a polarizing prism in the ultraviolet, in which case one of the beams must be eliminated, e.g., by imaging the source beyond the prism and blocking off the deviated image.

The paths of the two beams through the prism are shown in Fig. 11a. A ray normally incident on the entrance face travels along the optic axis in the first half of the prism, so that both ordinary and extraordinary rays are undeviated and have the same refractive index n_o . The second half of the prism has its optic axis at right angles to that in the first half, but the ordinary ray is undeviated since its refractive index is the same in both halves. The extraordinary ray, however, has its minimum index in the second half, so that it is refracted at the cut according to Snell's law (see "Double Refraction in Calcite" in Sec. 13.2). Since the deviation angle depends on the ratio n_e/n_o , it is a function of wavelength. If the angle of the cut is S , to a good approximation the beam deviation δ of the extraordinary ray depends on the cut angle in the following manner, according to Steinmetz et al.,⁴⁹

$$\tan S = \frac{n_o - n_e}{\sin \delta} + \frac{\sin \delta}{2n_e} \quad (6)$$

This relation holds for light normally incident on the prism face. The semifield angle i_{\max} is given by⁴⁹

$$\tan i_{\max} = \frac{1}{2}(n_e - n_o) \cot S \quad (7)$$

If the prism is to be used as a polarizer, the light should be incident as shown. Rochon prisms also act as polarizing beam splitters when used backward, but the deviation of the two beams is then slightly less.

When a Rochon prism is used backward, both the dispersion and the optical activity (for quartz) will adversely affect the polarization. Thus, one generally uses a Rochon in the normal manner. However, an exception occurs when a quartz Rochon is to be used as an analyzer. In this case it is

best to reverse the prism and use a detector that is insensitive to polarization to monitor the relative intensities of the two transmitted beams.

A Rochon prism is achromatic for the ordinary ray but chromatic for the extraordinary ray. Since total internal reflection does not occur for either beam, the type of cement used between the two halves of the prism is less critical than that used for conventional polarizing prisms. Canada balsam is generally used, although the two halves are sometimes optically contacted for high-power laser applications or for use in the ultraviolet at wavelengths shorter than 3500 Å. Optically contacted crystalline-quartz Rochon prisms can be used to wavelengths as short as 1700 Å, and a double Rochon of MgF_2 has been used to 1300 Å in the vacuum ultraviolet.⁴⁹ Optically contacted single Rochon prisms of MgF_2 have also been constructed, and the transmission of one has been measured from 1400 Å to 7 μm .⁵⁰ Ultraviolet-transmitting cements such as *gédamine* can be used to extend the short-wavelength limit of calcite prisms to about 2500 Å (see “Transmission” in Sec. 13.3).

Defects Quartz and calcite Rochon prisms suffer from several defects. Quartz exhibits optical activity when light is transmitted through it parallel to the optic axis, and although two mutually perpendicular, polarized beams will emerge from a quartz Rochon prism used in the conventional direction, their spectral composition will not faithfully reproduce the spectral compositions of the horizontal and vertical components of the input. If such a prism is used backward, different wavelengths emerge from the prism vibrating in different planes. Hence the output consists of many different polarizations instead of the desired two.⁵¹

Calcite Rochon prisms do not exhibit optical activity but are difficult to make, since when calcite surfaces are cut normal to the optic axis, small tetrahedra tend to cleave out from the surface during pitch polishing. These tetrahedra may also cleave out during attempts to clean the prisms, and occasionally glass plates are cemented to such surfaces to prevent damage. Some image distortion will occur in calcite prisms; if nonnormally incident rays pass through the prism, both beams will be distorted along their directions of vibration; i.e., the undeviated beam (*o* ray), which vibrates in a vertical plane, will be distorted vertically, and the deviated beam (*e* ray), which vibrates in a horizontal plane, will be distorted horizontally.⁵¹

Glass-Calcite Rochons Some of the difficulties mentioned in the preceding section can be minimized or eliminated by making the entrance half of the Rochon prism out of glass of matching index instead of quartz or calcite. Both *o* and *e* rays travel along the same path and have the same reflective index in this half of the prism, so that the birefringent qualities of the quartz or calcite are not being used and an isotropic medium would serve just as well. By properly choosing the index of the glass, either the ordinary or the extraordinary ray can be deviated, and glasses are available for matching either index of calcite reasonably well over much of the visible region.⁵¹ The extraordinary ray always suffers some distortion in its direction of vibration, but the distortion of the ordinary ray can be eliminated in the glass-calcite construction. By properly choosing the refractive index of the glass we can determine whether the *e* ray will be the deviated or the undeviated beam. (Some distortion also arises for deviated beams in the direction of the deviation because of Snell's law and cannot be corrected in this way.) Another method of obtaining an undeviated beam was used by Hardy;⁵² unable to find a glass with refractive index and dispersion matching those of calcite, he selected a glass with the correct dispersive power and then compensated for the difference in refractive index by putting a slight wedge angle on the calcite surface. Now a wider selection of glasses is available, but glass-calcite prisms cannot be made strictly achromatic over an extended wavelength range, and thermally induced strains caused by the difference in expansion coefficients in the two parts of the prism may be expected unless the cement yields readily.

Total Internal Reflection in Rochons When normal Rochon prisms are used as polarizers, one of the beams must be screened off and eliminated. This restriction might be removed by making the cut between halves of the prism at a sufficiently small angle for the extraordinary ray to be totally reflected. Calculations indicate that this approach should be feasible,⁵³ but it has apparently not been followed.

Sénarmont Prism

The Sénarmont polarizing beam splitter, shown in Figs. 10*b* and 11*b*, is similar to the Rochon prism except that the optic axis in the exit half of the prism is coplanar with the optic axis in the entrance half, i.e., at right angles to the Rochon configuration. As a result, light whose plane of vibration is initially vertical is deviated in the Sénarmont prism, while in the Rochon prism the deviated beam has its plane of vibration horizontal (assuming no optical activity in either case) (compare Fig. 11*a* and *b*). The amount of the deviation in the Sénarmont prism is slightly less than in the Rochon because the extraordinary ray does not have its minimum refractive index [Eq. (1)].

An alternate form of Sénarmont prism, the right-angle Sénarmont or Cotton polarizer,⁵⁴ consists of only the first half of the Sénarmont prism. Unpolarized light normally incident on the prism face is totally internally reflected at the hypotenuse and is then resolved into two planes of vibration, one parallel to the optic axis and the other perpendicular to it. Double refraction will then occur just as in a normal Sénarmont prism. Such a prism has a transmission equivalent to that of an optically contacted Sénarmont or Rochon but is much less expensive.

Wollaston Prism

The Wollaston prism (Figs. 10*c* and 11*c*) is a polarizing beam splitter, also used as a polarizing prism in the vacuum ultraviolet,⁵⁵ that deviates both transmitted beams. The deviations, indicated in Fig. 11*c*, are nearly symmetrical about the incident direction, so that the Wollaston has about twice the angular separation of a Rochon or Sénarmont prism. A normally incident beam is undeviated upon entering the prism, but the *o* ray, vibrating perpendicular to the optic axis, has a refractive index n_o , while the *e* ray, vibration parallel to the optic axis has its minimum (or principal) index n_e . At the interface the *e* ray becomes the *o* ray and vice versa because the direction of the optic axis in the second half is perpendicular to its direction in the first half. Thus the original *o* ray enters a medium of lower refractive index and is refracted away from the normal at the cut, while the original *e* ray passes into a medium of higher refractive index and is refracted toward the normal. On leaving the second half of the prism, both rays are refracted away from the normal, so that their divergence increases.

The deviation of each beam is chromatic in Wollaston prisms, which are most commonly used to determine the relative intensities of two plane-polarized components. Since the light never travels along the optic axis, optical activity does not occur and the relative intensities of the two beams are always proportional to the intensities of the horizontal and vertical polarization components in the incident beam. For an L/A ratio of 1.0, the angular separation between beams is about 1° for a crystalline-quartz Wollaston prism; it can be as high as $3^\circ 30'$ for an L/A ratio of 4.0. With a calcite prism, the beams would have an angular separation of about 19° for an L/A ratio of 1.0, but severe image distortion and lateral chromatism results when such large angular separations are used. These effects can be minimized or the angular separation can be increased for a given L/A ratio by using a three-element Wollaston prism, a modification, apparently suggested by Karl Lambrecht.³⁰ Divergences as large as 30° can be obtained.¹

The ellipticity in the emergent polarized beams has been measured by King and Talim.²⁵ For calcite Wollaston prisms, the ellipticities were in the 0.004° to 0.025° range, comparable to those of Glan-Thompson prisms ("Common Defects and Testing of Glan-Type Prism" in Sec. 13.3). Larger values, between 0.12° and 0.16° , were measured for crystalline-quartz Wollaston prisms. The major contribution, which was from the combined optical activity and birefringence in the quartz rather than from defects within the crystal, cannot be avoided in quartz polarizers.

Foster Prism

This prism, shown in a three-dimensional view in Fig. 10*d* and in cross-section in Fig. 11*d*, can be used to form two plane-polarized beams separated by 90° from each other.⁵⁶ Its construction is similar to that of a Glan-Thompson prism except that one side is cut at an angle and silvered to reflect the ordinary ray out the other side.

The Foster prism is often used backward as a polarizing microscope illuminator for observing reflecting specimens. For this application, the light source is at e in Fig. 11*d*, and unpolarized light enters the right-hand face of the prism. The ordinary ray (not shown) is reflected at the cut and absorbed in the blackened side of the prism, while the extraordinary ray is transmitted undeviated out the left face of the prism. It then passes through the microscope objective and is reflected by the specimen, returning on its same path to the prism. Light that is unchanged in polarization will be transmitted undeviated by the prism along the path to the light source. If, however, the plane of vibration has been rotated so that it is at right angles to the optic axis (in the plane of the figure), the light will be reflected into the eyepiece. The prism thus acts like a crossed polarizer-analyzer combination.

If a correctly oriented quarter-wave plate is inserted in the beam between the prism and the microscope objective, the light striking the sample will be circularly polarized, and, after being reflected back through the quarter-wave plate, it will be linearly polarized again but with the plane of vibration rotated by 90° . This light is vibrating perpendicular to the optic axis and will be reflected into the eyepiece, giving bright-field illumination. Foster prisms used in this manner introduce no astigmatism since the light forming the image enters and leaves the prism normal to the prism faces and is reflected only by plane surfaces.

Beam-Splitting Glan-Thompson Prism

If a prism design similar to the Foster is used but the side of the prism is cut at an angle so that the ordinary ray, which is deflected, passes out normal to the surface of the prism rather than being reflected, the prism is called a beam-splitting Glan-Thompson prism (Figs. 10*e* and 11*e*). Since no refraction occurs for either beam, the prism is achromatic and nearly free from distortion. The angle between the two emerging beams is determined by the angle of the cut between the two halves of the prism and hence depends on the L/A ratio of the prism. For an L/A ratio of 2.414, the angle is 45° . The field angle around each beam is calculated for different L/A ratios just as for a conventional Glan-Thompson prism. By making the prism double, i.e., a beam-splitting Ahrens prism, the incident beam can be divided into three parts, one deflected to the left, one to the right, and one undeviated.

13.6 FEUSSNER PRISMS

The polarizing prisms discussed so far require large pieces of birefringent material, and the extraordinary ray is the one usually transmitted. Feussner⁵⁷ suggested an alternate prism design in which only thin plates of birefringent material are required and the ordinary ray rather than the extraordinary ray is transmitted for negative uniaxial materials. A similar suggestion was apparently made by Sang in 1837, although he did not publish it until 1891.⁵⁸ In essence, Feussner's idea was to make the prisms isotropic and the film separating them birefringent, as shown in Fig. 12. The isotropic

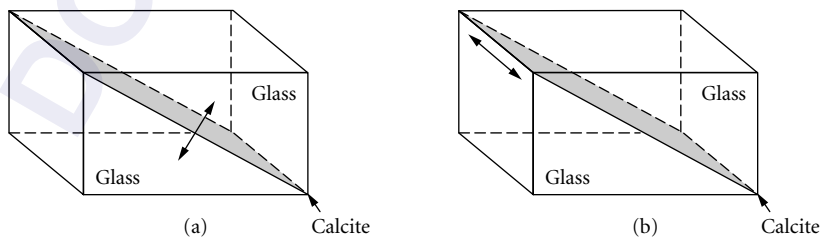


FIGURE 12 Types of Feussner prisms: (a) original Feussner prism and (b) Bertrand type. The arrows indicate the orientation of the optic axis in the calcite (or other birefringent material).

prisms should have the same refractive index as the higher index of the birefringent material so that for negative uniaxial materials e.g., calcite or sodium nitrate, the ordinary ray is transmitted and the extraordinary ray totally internally reflected. Advantages of this design are (1) since the ordinary ray is transmitted, the refractive index does not vary with angle of incidence and hence the image is anastigmatic, (2) large field angles or prisms of compact size can be obtained, and (3) the birefringent material is used economically. Furthermore, because the path length of the ray through the birefringent material is short, a lower-quality material can be used.

Disadvantages are (1) for both calcite and sodium nitrate, the extraordinary ray is transmitted over a larger wavelength range than the ordinary ray so that Feussner prisms do not transmit over as large a wavelength range as conventional prisms, and (2) the thermal-expansion coefficients of the isotropic and birefringent materials are different, making thermally induced strains likely. Solutions to the second problem are to use a thixotropic cement, which flows more readily with increasing stress, or to enclose the system in a metal sleeve and use oil instead of cement. If the ordinary index is matched by the oil, the birefringent material does not even need to be polished very well. Even a cleavage section of calcite can be used, with only a tolerable loss in angular field.⁵⁹

Feussner suggested orienting the optic axis of the birefringent slab perpendicular to the cut, as indicated in Fig. 12a. Since the thermal expansion of the slab is the same in all directions perpendicular to the optic axis, thermally induced strains are minimized in this way. Field angles for Feussner prisms employing calcite and sodium nitrate slabs are given in the earlier polarization chapter.¹

Shortly after Feussner's article was published, Bertrand⁶⁰ pointed out that the optic axis of the birefringent slab should be parallel to the entrance face of the prism to give the maximum difference between the refractive indices of the ordinary and extraordinary rays. A prism made in this way, sometimes called a Bertrand-type Feussner prism, is shown in Fig. 12b.

Since sodium nitrate is easily obtainable and has a birefringence even larger than that of calcite, attempts have been made to produce polarizing prisms of this material by Wulff,⁶¹ Stöber,⁶²⁻⁶⁴ Tzekhovitzer,⁶⁵ West,⁶⁶ Huot de Longchamp,⁶⁷ and Yamaguti.^{68,69} However, it is not only deliquescent but also very soft, so that although large single crystals can be obtained, they are difficult to work. They can be crystallized in the desired orientation from a melt using a technique discovered by West.⁶⁶ When sodium nitrate crystallizes from a melt on a mica cleavage surface, one of its basal planes is oriented parallel to the mica cleavage and hence its optic axis is perpendicular to the mica surface. West reports growing single crystals as large as $38 \times 19 \times 2$ cm using this technique. Yamaguti^{68,69} has produced polarizing prisms of sodium nitrate by placing thin, closely spaced glass plates on edge on a mica sheet and then immersing the assembly in a melt of sodium nitrate. The thin single crystal thus formed was annealed and cemented between glass prisms to form a Bertrand-type Feussner prism. Conceivably, the sodium nitrate could have been grown directly between the glass prisms themselves, but when such thick pieces of glass are used, it is difficult to avoid setting up strains in the crystal and consequently reducing the polarization ratio. Yamaguti used SK5 glass prisms ($n_D = 1.5889$) cut at an angle of 23° to form his polarizing prism and reports a field of view of 31° , symmetric about the normal to the entrance face.

Another possible birefringent material suitable for a Feussner prism is muscovite mica, and such prisms have actually been constructed and tested.^{70,71} A 6° field angle can be obtained,⁵⁹ which is adequate for many optical systems illuminated by lasers.

13.7 NONCALCITE POLARIZING PRISMS

Polarizing prisms made of materials other than calcite have been used primarily in the ultraviolet region at wavelengths for which calcite is opaque. Prism materials used successfully in this region include crystalline quartz, magnesium fluoride, sodium nitrate, and ammonium dihydrogen phosphate. Rutile polarizing prisms have been used beyond the calcite cutoff in the infrared. A new prism material, yttrium orthovanadate, has been used to make high-transmission polarizers for the

visible and near-infrared spectral regions.⁷² Properties of this material were described in the earlier polarization chapter.¹

Rochon or Wollaston prisms (see “Rochon Prism” and “Wollaston Prism” in Sec. 13.5) are sometimes made of crystalline quartz for use in the far ultraviolet. The short-wavelength cutoff of the quartz is variable, depending on the impurities present, but can be as low as 1600 Å.

By utilizing magnesium fluoride instead of quartz for the polarizing prisms, the short-wavelength limit can be extended to 1300 Å. Magnesium fluoride transmits to about 1125 Å, but below 1300 Å its birefringence decreases rapidly and changes sign at 1194 Å.^{55,73} Although it is the most birefringent material available in this region, MgF₂ has a much smaller birefringence than that of calcite; hence, a small cut angle and large L/A ratio for the prism are unavoidable. Since absorption does occur, it is desirable to minimize the length of the prism. Johnson⁵⁵ solved this problem by constructing a MgF₂ Wollaston prism which requires only half the path length necessary for a Rochon prism. However, both beams are deviated, creating instrumental difficulties.

Steinmetz et al.⁴⁹ constructed a double Rochon prism of MgF₂ which has the same L/A ratio as the Wollaston prism but does not deviate the desired beam. Problems with the prism included fluorescence, scattered light, and nonparallelism of the optic axes.¹ In principle, however, a MgF₂ double Rochon polarizing prism should be an efficient, high-extinction-ratio, on-axis polarizer for the 1300- to 3000-Å wavelength range and should also be useful at longer wavelengths. Morris and Abramson⁵⁰ reported on the characteristics of optically contacted MgF₂ single Rochon prisms.

A different type of polarizer suggested by Chandrasekharan and Damany⁷⁴ to take the place of a Rochon or Wollaston prism in the vacuum ultraviolet consisted of a combination of two MgF₂ lenses, one planoconcave and the other planoconvex of the same radius of curvature, combined so that their optic axes were crossed. The combination acted as a convergent lens for one polarization and as a divergent lens for the other. It had the advantage that the polarized beam remained on axis and was focused. A measured degree of polarization of 98.5 percent was obtained at 1608 Å, in good agreement with the calculated value.

Prism polarizers can also be constructed for use in the infrared at wavelengths longer than those transmitted by calcite. Rutile, TiO₂, a positive uniaxial mineral with a large birefringence and good transmittance to 5 μm in the infrared, has been used by Landais⁷⁵ to make a Glan-Foucault-type crystal polarizer. Since rutile has a positive birefringence (in contrast to the negative birefringence of calcite), the ordinary ray is transmitted undeviated and the extraordinary ray is reflected out one side. Other characteristics are given in the earlier Polarization chapter.¹

13.8 DICHROIC AND DIFFRACTION-TYPE POLARIZERS

Some of the most useful polarizers available employ either dichroism or diffraction effects. These polarizers come in sheet form, sometimes in large sizes, are easily rotated, and produce negligible beam deviation. Also, they are thin, lightweight, and rugged, and most can be made in any desired shape. The cost is generally much less than that of a prism-type polarizer. Furthermore, both types are insensitive to the degree of collimation of the beam, so that dichroic or diffraction-type polarizers can be used in strongly convergent or divergent light.

A dichroic* material is one which absorbs light polarized in one direction more strongly than light polarized at right angles to that direction. Dichroic materials are to be distinguished from birefringent materials, which may have different refractive indexes for the two electric vectors vibrating at right angles to each other but similar (usually negligible) absorption coefficients.

*The term *dichroic* is also used in three other ways: (1) to denote the change in color of a dye solution with change in concentration, (2) to denote a color filter that has two transmission bands in very different portions of the visible region and hence changes color when the spectral distribution of the illuminating source is changed, and (3) to denote an interference filter that appears to be of a different color when viewed in reflected or transmitted light.

Various materials are dichroic, either in their natural state or in a stretched condition. The most common materials used as dichroic polarizers are stretched polyvinyl alcohol sheets treated with absorbing dyes or polymeric iodine, commonly marketed under the trade name Polaroid. These and similar materials are discussed in the following section. Another type of dichroic polarizer is prepared by rubbing a glass or plastic surface in a single direction and then treating it with an appropriate dye. Polarizers of this type are sold under the trade name Polacoat and will be described in "Dichroic Polarizing Coating" section on p. 13.28. In certain portions of the infrared spectral region, calcite is strongly dichroic and makes an excellent high-extinction polarizer.⁷⁶ Pyrolytic graphite is electrically and optically anisotropic and has been successfully used as an infrared polarizer; it is described in "Pyrolytic-Graphite Polarizers" section on p. 13.28. Other materials which exhibit dichroism in the infrared include single-crystal tellurium,⁷⁷ ammonium nitrate,⁷⁸ mica, rubber under tension, polyvinyl alcohol, and polyethylene.⁷⁹ In the visible region, gold, silver, and mercury in the form of microcrystals,⁸⁰ needles of tellurium,⁸¹ graphite particles,⁸² and glasses containing small elongated silver particles⁸³ are all dichroic.

A sodium nitrate polarizer described by Yamaguti⁸⁴ is not dichroic in the strict sense of the word but acts like a dichroic polarizer. Roughened plates of SK5 glass are bonded together by a single crystal of sodium nitrate, which has a refractive index for the ordinary ray nearly equal to that of the glass. The extraordinary ray has a much lower index, so that it is scattered out of the beam by the rough surfaces, leaving the ordinary ray to be transmitted nearly undiminished. (Yamaguti has also made Feussner prisms out of single-crystal sodium nitrate described in Sec. 13.6.)

Diffraction-type polarizers include diffraction gratings, echelettes, and wire grids. These are all planar structures that have properties similar to those of dichroic polarizers except that they transmit one component of polarization and reflect the other when the wavelength of the radiation is much longer than the grating or grid spacing. Wire grid and grating polarizers are covered in "Wire-Grid and Grating Polarizers" section on p. 13.30.

None of these polarizers has as high a degree of polarization as the prism polarizers of Secs. 13.1 to 13.5. Thus it is frequently necessary to measure the polarizing properties of the particular polarizer used. A source of plane-polarized light is desirable for such a measurement. Lacking that, one of the procedures described in "Measuring Polarization of Imperfect Polarizers" section on p. 13.33 can be followed if there are two identical imperfect polarizers. Alternate methods are also described which are applicable to two nonidentical imperfect polarizers.

Sheet Polarizers

Various types of sheet polarizers have been developed by Edwin H. Land and coworkers at the Polaroid Corporation, Cambridge, Mass. Sheet polarizers are also available from several European companies. The J sheet polarizer, the first type available in America (around 1930), consisted of submicroscopic needles of herapathite oriented parallel to one another in a sheet of cellulose acetate. Since this type of polarizer, being microcrystalline, had some tendency to scatter light, it was superseded by H and K sheet molecular polarizers, which exhibit virtually no scattering. The most widely used sheet polarizer is the H type, which consists of a sheet of polyvinyl alcohol that has been unidirectionally stretched and stained with iodine in a polymeric form. The K type is made by heating a sheet of polyvinyl alcohol in the presence of a catalyst to remove some of the water molecules and produce the dichromophore polyvinylene. It was developed primarily for applications where resistance to high temperature and high humidity are necessary. Another type of polarizing sheet, made from a combination of the H and K types, has an absorption maximum at about 1.5 μm in the infrared and is designated as HR Polaroid.

The history of the development of the various kinds of sheet polarizers has been given by Land,⁸¹ their chemical composition by Land and West,⁸⁰ and their optical performance by Shurcliff,⁸² Baumeister and Evans,⁸⁵ Land and West,⁸⁰ and Land.⁸¹ In addition, Blake et al.⁸⁶ mention the HR infrared polarizer, and Makas⁸⁷ describes the modified H-film polarizer for use in the near ultraviolet. Baxter et al.⁸⁸ describe a technique for measuring the optical density of high-extinction polarizers in the presence of instrumental polarization.

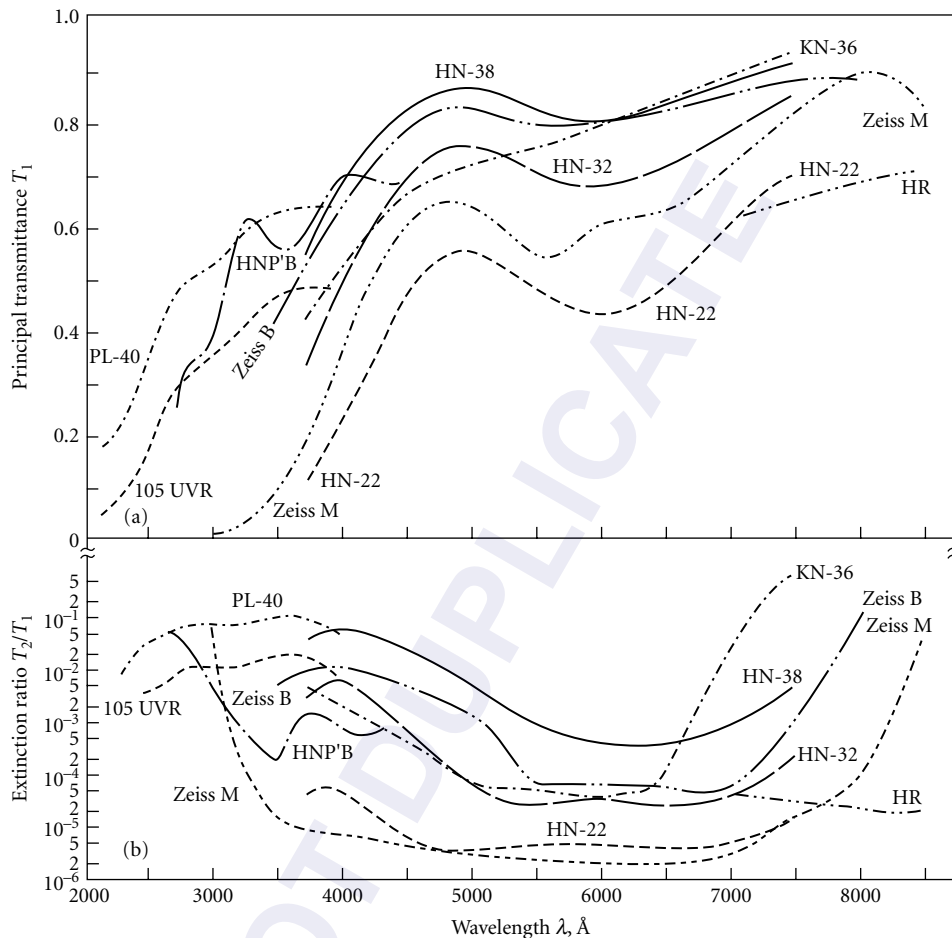


FIGURE 13 (a) Principal transmittance and (b) extinction ratio for various types of dichroic polarizers: Polaroid sheet polarizers HN-22, HN-32, HN-38, and KN-36; Zeiss (Oberkochen) Bernotar and Micro Polarization filters; and Polacoat PL-40 and 105 UVR polarizing filters. The last is stated to have a transmittance (for unpolarized light) of 32 percent at 5460 Å. (Modified from curves of Shurcliff,⁸² Baumeister and Evans,⁸⁵ Jones,⁸⁹ Haase,⁹⁰ and McDermott and Novick.⁹¹)

Figure 13 shows the principal transmittance T_1 and extinction ratio T_2/T_1 of various types of H and K sheet polarizers used in the visible and near ultraviolet.^{82,85,89} In addition, curves for two sheet polarizers manufactured by Zeiss and two types of polarizing filters from Polacoat (see “Dichroic Polarizing Coatings” section on p. 13.28) are shown. The letter N in the designation of the Polaroid sheets stands for neutral (to distinguish them from sheet polarizers prepared from colored dyes), and the number 22, 32, etc., indicates the approximate transmittance of unpolarized visible light. Figure 14 gives the principal transmittance and extinction ratio of a typical plastic laminated HR infrared polarizer.^{82,89} Sometimes the optical density D of a polarizer is plotted instead of its transmittance. The relation between these two quantities is

$$D = \log \frac{1}{T} \quad (8)$$

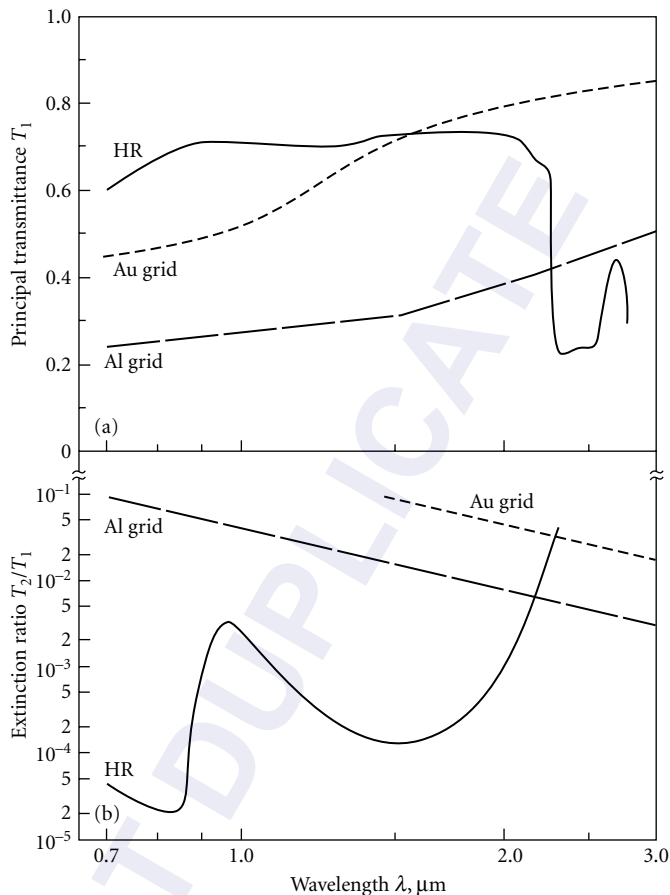


FIGURE 14 (a) Principal transmittance and (b) extinction ratio for plastic laminated HR infrared polarizer (modified from curves of Shurcliff,⁸² and Jones,⁸⁹) and two wire grid polarizers with 0.463- μm grating spacings (Bird and Parrish,⁹²).

The extinction ratio of the HN-22 Polaroid compares favorably with that of Glan-Thompson prisms throughout the visible region, but the transmission of the Glan-Thompson is superior. In the ultraviolet, the new HNP'B material has a reasonably good extinction ratio (about 10^{-3} or better) for wavelengths longer than 3200 Å. It is a specially purified form of HN-32, and its properties match those of the standard HNT-32 Polaroid at wavelengths longer than 4500 Å. Optical properties of various types of Polaroid dichroic polarizers have been described by Trapani.⁹³ According to West and Jones,⁴⁸ the extinction ratio for a dichroic polarizer of the Polaroid type has a practical limit of about 10^{-5} because, as the concentration of dichromophore is increased beyond a certain value, the optical density no longer increases proportionately. Gunning and Foschaar⁹⁴ have described a method for the controlled bleaching of the iodine dichromophore in iodine-polyvinyl alcohol polarizers to achieve an increased internal transmission of up to 95 percent for the principal transmittance of linearly polarized light in the 5000- to 6000-Å wavelength region. This is achieved at the expense of degrading the extinction ratio and drastically affecting the short wavelength performance of the polarizer. Baum⁹⁵ describes the application of sheet polarizers to liquid crystal displays and problems encountered in this application.

If Polaroids are used in applications where beam deviation is important, they should be checked for possible deviation. Most Polaroids, which are laminated in plastic sheets, do produce a slight beam deviation that can be observed through a telescope as a shift in the image position when the Polaroid is rotated. The amount of the deviation varies from point to point on the Polaroid and can be much worse if the material is mounted between glass plates. It is possible to order specially selected sheet Polaroid laminated between polished glass plates that deviates the beam by only about 5 seconds of arc.

Sheet polarizers made of stretched polyvinyl alcohol that has been stained with iodine or various dyes are also made in countries outside the United States, as described in the earlier polarization chapter.¹

King and Talim²⁵ have measured the axis wander and ellipticity of beams transmitted by various types of sheet polarizers in the same way as for Glan-Thompson prisms ("Common Defects and Testing of Glan-Type Prisms" in Sec. 13.3). They found considerable variations from one type of sheet polarizer to another and also over a single sheet. Details are given in the earlier chapter on polarization.¹

Dichroic Polarizing Coatings

Beilby-layer polarizers⁸² are dichroic coatings that can be applied to the surface of glass or plastic. The process was developed by Dreyer,⁹⁶ who founded the company which manufactures Polacoat polarizing filters. There are three main steps in the production of these polarizers. First, the substrate (quartz, glass, plastic, etc.) is rubbed along parallel lines with filter paper, cotton, or rouge to produce a preferred surface orientation. (The affected region of minute scratches extends to a depth of less than 1 μm .) Then the sheet is rinsed and treated with a solution of dichroic molecules e.g., a 0.5 percent solution of methylene blue in ethanol or one or more azo dyes, and then dried in a controlled fashion. Presumably the molecules line up preferentially along the rubbing direction, resulting in a greater absorption for light, polarized in that direction. As a final step, the surface is treated with an acidic solution, often that of a metallic salt such as stannous chloride, which can increase the dichroism and produce a more neutral color. A protective coating over the polarized surface provides mechanical protection for the fragile layer with no loss in transmission. McDermott and Novick⁹¹ give a somewhat more complete description of the Polacoat process, and Anderson⁹⁷ has investigated the absorption of methylene blue molecules on a unidirectionally polished surface. References to patents and related work are given by Shurcliff.⁸²

The principal transmittance and extinction ratio of two standard Polacoat coatings, PL-40 and 105 UVR (32 percent transmission of unpolarized light at 5460 \AA), are shown in Fig. 13. These curves are taken from the data of McDermott and Novick.⁹¹ Polacoat 105 UVR coating comes in various densities; the data shown are for the highest-density material with the best extinction ratio.* A major advantage of Polacoat over sheet Polaroid is that it does not bleach upon exposure to intense ultraviolet radiation.

Kyser⁹⁹ tested a stock PL40 polarizing filter on fused quartz and found that it produced a large quantity of scattered light of the unwanted component. This light was dispersed spectrally and was scattered at angles up to about 20° as though the scratches on the rubbed surface were acting like rulings on a diffraction grating. There was relatively little of the unwanted component on axis; most of it was scattered at larger angles. Despite these difficulties, Polacoat PL40 polarizers appear to be the best large-aperture transmission-type polarizers available for work in the 2000- to 3000- \AA wavelength range in the ultraviolet.

Pyrolytic-Graphite Polarizers

Pyrolytic graphite has a large anisotropy in both the electric conductivity and in the optical properties. If the E vector of an electromagnetic wave is pointing in the direction of the c -axis of the graphite, the absorption coefficient is a minimum, the reflectance is also a minimum, and hence

*The company literature⁹⁸ is somewhat misleading in that the transmittance of this material is stated to be 35 percent, but the transmission curve (for unpolarized light) given in the bulletin does not rise above 30 percent until the wavelength becomes longer than 6500 \AA

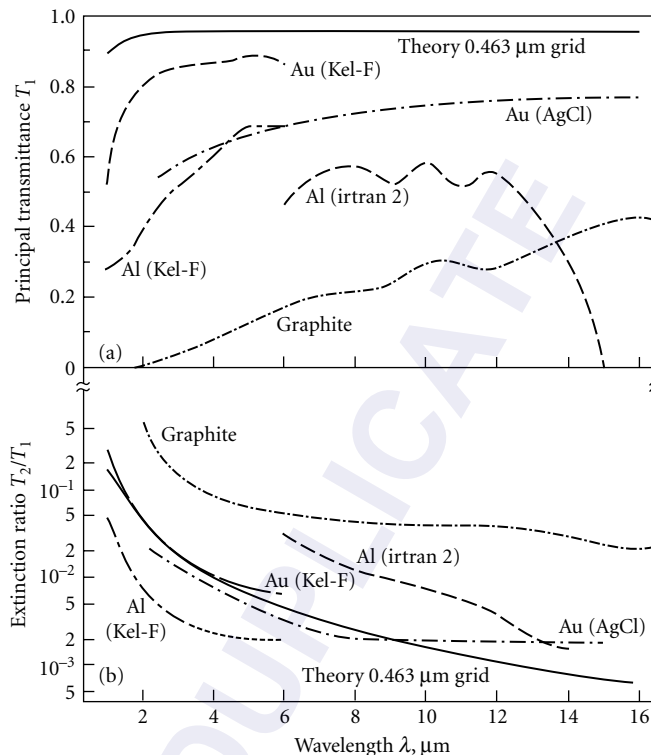


FIGURE 15 (a) Principal transmittance and (b) extinction ratio for a pyrolytic-graphite polarizer (Rupprecht et al.¹⁰⁰) and various wire-grid polarizers. (Bird and Parrish,⁹² Perkin-Elmer,¹⁰¹ and Young et al.¹⁰²) The substrate materials and metals used for the grids are indicated. Theoretical curves (solid lines) calculated from relations given in Ref. 1 with $n = 1.5$ and $d = 0.463$ are also shown for comparison.

the transmittance is a maximum. If the E vector lies in the plane perpendicular to the c direction, the absorption is a maximum, reflectance is a maximum, and transmittance is a minimum. Thus, pyrolytic graphite should be a good material from which to make a dichroic polarizer if a thin foil is cut and polished to contain the c -axis. Several such polarizers have been made by Rupprecht et al.¹⁰⁰; two had thicknesses of 9.2 μm , and a third was 4.2 μm thick. The transmittances T_1 of the thinner one and T_1 and T_2 of the two thicker ones were determined using one of the methods described in Par. 49 of the earlier Polarization chapter.¹ The principal transmittance and extinction ratio for one of the 9.2- μm -thick ones are shown in Fig. 15 for infrared wavelengths from 2 to 16 μm , along with curves for various wire-grid polarizers (see “Wire Grid and Grating Polarizers” section next). In the far infrared out to 600 μm , T_1 gradually increases to 0.50 and T_2/T_1 drops down to the 10^{-3} range.¹⁰⁰ The transmittance of the thinner pyrographite polarizer was larger than the curve shown, but its extinction ratio, although not given, was probably poorer. Pyrolytic-graphite polarizers have the advantages of being planar and thus easily rotatable, having large acceptance angles, and having reasonably high transmittances and good extinction ratios in the far infrared. However, in the shorter-wavelength region shown in Fig. 15, they are inferior to all the wire-grid polarizers. In addition, they are fragile, and the largest clear aperture obtained by Rupprecht et al.¹⁰⁰ was about 12 mm diameter.

Wire-Grid and Grating Polarizers

Wire grids* have a long history of use as optical elements to disperse radiation and detect polarization in far-infrared radiation and radio waves.⁹² They transmit radiation whose E vector is vibrating perpendicular to the grid wires and reflect radiation with the E vector vibrating parallel to the wires when the wavelength λ is much longer than the grid spacing d . When λ is comparable to d , both components are transmitted. For grids made of good conductors, absorption is negligible. Various aspects of the theory of reflection and transmission of radiation by wire grids are summarized in the earlier polarization chapter.¹ In addition to that theoretical treatment, Casey and Lewis^{104,105} considered the effect of the finite conductivity of the wires on the transmission and reflection of wire-grid polarizers when the light was polarized parallel to the wires. Mohebi, Liang, and Soileau¹⁰⁶ extended the treatment to the case for which light was polarized both parallel and perpendicular to the wires; they also calculated the absorption of the wire grids as a function of d/λ . In addition, they measured the absorption and surface damage of wire-grid polarizers consisting of aluminum strips (0.84- μm period) deposited on ZnSe substrates at 10.6 μm , 1.06 μm , and 0.533 μm . Stobie and Dignam¹⁰⁷ calculated the amplitude transmission coefficients for parallel and perpendicular components and relative phase retardation between them, both as a function of λ/d . Burton¹⁰⁸ proposed using wire-grid polarizers in the form of cylinders and paraboloids instead of planar structures in infrared interferometers, but did not show any experimental measurements.

Figure 16 shows values of the calculated principal transmittance and extinction ratio for various values of the refractive index n as a function of λ/d . These curves were calculated from relations given in the earlier polarization chapter.¹ It is clear that the shortest wavelength for which a given grid will act as a useful polarizer is $\lambda \approx 2d$. Also, the best performance is obtained with the lowest refractive index substrate. Since absorption in the substrate material has been neglected, principal transmittances measured for real materials will be lower than the calculated values, but the extinction ratios should be unaffected. If one must use a high refractive index substrate such as silicon or germanium, the performance of the grid can be considerably improved by applying an antireflection coating to the substrate *before* depositing the conducting strips, since a perfectly antireflected substrate acts like an unsupported grid.¹⁰⁹ However, if the antireflecting layer is laid down *over* the grid strips, the performance of the wire-grid polarizer is degraded.¹

Many people have built and tested wire-grid polarizers including Bird and Parrish,⁹² Young et al.,¹⁰² Hass and O'Hara,¹¹⁰ Hilton and Jones,¹¹¹ Auton,¹⁰⁹ Vickers et al.,¹¹² Cheo and Bass,¹¹³ Auton and Hutley,¹¹⁴ Costley et al.,¹¹⁵ Beunen et al.,¹¹⁶ Leonard,¹¹⁷ Sonek et al.,¹¹⁸ Eichhorn and Magner,¹¹⁹ and Novak et al.¹²⁰ In addition, two types of wire grids are manufactured commercially by Buckbee Mears (see Ref. 110) and Perkin-Elmer,¹⁰¹ and a third type composed of 152 μm -diameter tungsten wires spaced 800 to the inch has been mentioned, but no performance characteristics have been given.¹²¹ Hwang and Park¹²² measured the polarization characteristics of two-dimensional wire mesh (64 μm and 51 μm spacings) at a laser wavelength of 118.8 μm . The different wire-grid polarizers are listed in Table 2, and the principal transmittances and extinction ratios of several are shown in Figs. 14 and 15.

The polarizers with grid spacings of 1.69 μm and less were all made by evaporating the grid material at a very oblique angle onto a grating surface which had been prepared either by replicating a diffraction grating with the appropriate substrate material (silver bromide, Kel-F, polymethyl methacrylate, etc.) or by ruling a series of lines directly onto the substrate (Irtran 2 and Irtran 4). The oblique evaporation (8° to 12° from the surface) produced metallic lines on the groove tips which acted like the conducting strips of the theory, while the rest of the surface was uncoated and became the transparent region between strips. Larger grid spaces (4 to 25.4 μm) were produced by a photoetching process, and one 25.4- μm grid was made by an electroforming process. Still larger grid spacings were achieved by wrapping wires around suitable mandrels.

If a wire-grid polarizer is to be used in the near infrared, it is desirable to have the grid spacing as small as possible. Bird and Parrish⁹² succeeded in obtaining a very good extinction ratio in the 2- to 6- μm wavelength region with an aluminum-coated Kel-F substrate (Figs. 14 and 15). Unfortunately,

*Wire grid is being used here, as is customary, to denote a planar structure composed of a series of parallel wires or strips. Renk and Genzel¹⁰³ and a few others use the term to designate a two-dimensional array with two series of elements arranged at right angles to each other. They call a one-dimensional array a wire or strip grating.

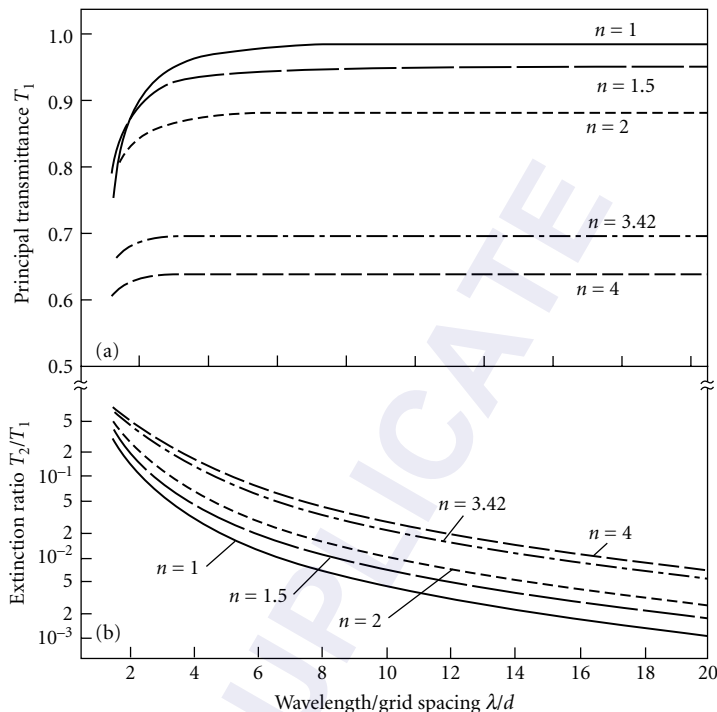


FIGURE 16 (a) Principal transmittance and (b) extinction ratio as a function of λ/d calculated from relations given in Ref. 1 for various values of n for the substrate. Substrate indexes correspond approximately to an antireflected substrate of air, organic plastic, silver chloride, silicon, and germanium.

Kel-F (CF_2CFCl)_n, has absorption bands at 7.7 to 9.2 μm and 10.0 to 11.0 μm , making the polarizer useless in these regions, but it can be used at longer wavelengths out to 25 μm .⁹² Polyethylene would be an excellent substrate material since it has fewer absorption bands than Kel-F, but its insolubility in common solvents makes it much more difficult to use for replicating gratings.¹¹⁰ It does, however, make an excellent substrate material for photoetched grids.¹⁰⁹

For infrared wavelengths longer than about 24 μm , a photoetched grid with 1- μm -wide lines (close to the present limit for the photoetching process) and a 2- μm spacing should have an extinction ratio of 5×10^{-3} or better if the refractive index of the substrate is about 1.5—for example, polyethylene. The extinction ratio would continue to decrease; i.e., the polarization properties would improve as the wavelength is increased. At very long wavelengths, grids with a larger spacing would have a high degree of polarization. The important factor is the ratio of wavelength to grid spacing, which should be kept as large as possible (Fig. 16b).

One definite advantage of the wire-grid polarizer is that it can be used in sharply converging beams, i.e., systems with high numerical apertures. Young et al.¹⁰² found no decrease in percent of polarization for an Irtran 2 polarizer at 12 μm used at angles of incidence from 0° to 45°. They did find, however, that the transmittance decreased from 0.55 at normal incidence to less than 0.40 at 45° incidence.

If a grid were to be used at a single wavelength, one might possibly make use of interference effects in the substrate to increase the transmission.¹⁰⁹ If the substrate has perfectly plane-parallel surfaces, it will act like a Fabry-Perot interferometer and transmit a maximum amount of light when twice the product of the thickness and refractive index is equal to an integral number of wavelengths. The 0.25-mm-thick pressed polyethylene substrates used by Auton¹⁰⁹ were not uniform enough to show interference effects, but the Mylar film backing on the Buckbee Mears electroformed grid did show interference effects.¹¹⁰

TABLE 2 Types of Wire-Grid Polarizers

Grid spacing, μm	Grid material	Substrate	Wavelength range, μm	Reference
0.115	Evaporated Al	Quartz	0.2–0.8	Sonek et al. ¹¹⁸
0.347	Evaporated Au	Silverchloride	2.5–30	Perkin-Elmer ¹⁰¹
0.22–0.71	Evaporated Al	KRS-5	3–15, 3.39, 10.6	Auton and Hutley ¹¹⁴
0.22–0.45	Evaporated Al	CaF ₂	3–10, 3.39	Auton and Hutley ¹¹⁴
0.42	Evaporated Al	Glass	3–5	Auton and Hutley ¹¹⁴
0.463	Evaporated Au	Kel-F	1.5–10*	Bird and Parrish ⁹²
0.463	Evaporated Al	Kel-F	0.7–15*	Bird and Parrish ⁹²
0.463	Evaporated Al	Polymethyl methacrylate	1–4000 [†]	Hass and O'Hara ¹¹⁰
1.67	Evaporated Al	Irtran 2	6–14	Young et al. ¹⁰²
1.67	Evaporated Al	Irtran 4	8–19	Young et al. ¹⁰²
1.69	Evaporated Al	Polyethylene	2.9–200 [‡]	Hass and O'Hara ¹¹⁰
2	Evaporated Cr	Silicon	10.6	Cheo and Bass ¹¹³
?	?	BaF ₂	2–12	Leonard ¹¹⁷
?	?	ZnSe	3–17	Leonard ¹¹⁷
4	Photoetched Al	Polyethylene	>16	Auton ¹⁰⁹
5.1	Photoetched Al	Silicon	54.6	Hilton and Jones ¹¹¹
10	Photoetched Al	Polyethylene	>16	Auton ¹⁰⁹
25.4	Photoetched Al	Silicon	54.6	Hilton and Jones ¹¹¹
25.4	Evaporated Au	Mylar	>60	Hass and O'Hara ¹¹⁰
25	Stainless steel wire 8 μm diam.	Air	80–135	Novak et al. ¹²⁰
32.4	Gold-coated W wire 21 μm diam.	Air	100–10,000	Eichhorn and Magner ¹¹⁹
64, 51	Wire mesh (2D)	Air	118.8	Hwang and Park ¹²²
?	Stainless steel wire 50 μm diam.	Air	200–1000	Vickers et al. ¹¹²
317	W wire 152 μm diam.	Air	40–300	Roberts and Coon ¹²¹
25–1800	W wire 10 μm diam.	Air	>50	Costley et al. ¹¹⁵
30–65	W wire 10 μm diam.	Air	22–500, 337	Beunen et al. ¹¹⁶

*Strong absorption bands near 8.3 and 10.5 μm .

[†]Strong absorption bands between 5.7 and 12.5 μm .

[‡]Absorption bands between 6 and 15.5 μm .

Lamellar eutectics of two phases consist of thin needles of a conducting material embedded in a transparent matrix. The material is made by a controlled cooling process in which there is a unidirectional temperature gradient. This method of cooling orients conducting needles parallel to the temperature gradient, and hence the material can act like a wire-grid polarizer. Weiss and coworkers^{123–125} have grown eutectic alloys of InSb and NiSb in which the conducting needles of NiSb are approximately 1 μm in diameter and approximately 50 μm long. A degree of polarization of more than 99 percent has been reported. Other eutectic alloys of InAs, GaSb, and InSb containing conducting needlelike crystals of Ni, Fe, Mn, Cr, and Co (or their compounds) have also been investigated. An advantage of this type of polarizer is that its performance can be optimized at a specific wavelength, e.g., that of a CO₂ laser line, by choosing the thickness of the crystalline film so that there will be an interference maximum at the desired wavelength.¹²⁶ Recently, Saito and Miyagi¹²⁷ have proposed using a thin film of anodized aluminum with implanted metallic columns to make a high-performance polarizer. Their theoretical calculations suggest that this type of polarizer should have a large extinction ratio and low loss in the infrared.

In summary, wire grids are very useful infrared polarizers, particularly for wavelengths much greater than the grid spacing. They are compact and easily rotatable and can be used with sharply converging beams. A major advantage is the extreme breadth of the wavelength band over which they have good polarizing properties. The long-wavelength limit is set by the transmission of the substrate material rather than by the loss of polarization of the grid. The short-wavelength limit is

determined by the grid spacing; if gratings with smaller spacings could be successfully replicated and coated, the short-wavelength limit could be pushed closer to the visible region.

Another possible method of producing plane-polarized light is by using diffraction gratings or echelette gratings. Light reflected from diffraction gratings has long been known to be polarized, but the effect is generally small and extremely wavelength-dependent.^{128,129} However, Roumiguieres¹³⁰ predicted that under certain conditions (rectangular groove grating with equal groove and land spacings and small groove depth), a high polarizing efficiency could be obtained. For wavelengths in the range $1.1 < \lambda/d < 1.7$, over 80 percent of the light polarized parallel to the grooves should be reflected in the zero order at a 50° angle of incidence and less than 5 percent of the other polarization. His predictions were verified by Knop¹³¹ who fabricated gold-coated photoresist gratings as well as an electroplated nickel master grating. Knop's measured reflectances of the two polarized components were within ± 3 percent of the predicted values. In general, one tries to avoid polarization in the diffracted light to obtain high efficiencies in a blazed grating since polarization effects are frequently associated with grating anomalies.^{132,133}

In contrast to diffraction gratings, echelette gratings, have been found to produce an appreciable amount of plane-polarized light. Experimental studies have been made by Peters et al.,¹³⁴ Hadni et al.,^{135,136} and Mitsubishi et al.,¹³⁷ as discussed in the earlier Polarization chapter.¹ The theory of the polarization of light reflected by echelette gratings in the far-infrared and microwave regions has been given by Janot and Hadni¹³⁸ and Rohrbaugh et al.¹³⁹ A general numerical technique published by Kalhor and Neureuther¹⁴⁰ should be useful for calculating the polarization effects of echelette gratings of arbitrary groove shape used in the visible region.

Measuring Polarization of Imperfect Polarizers

In determining the principal transmittance, extinction ratio, and other properties of an imperfect polarizer, the effects of source polarization, instrumental polarization, and sensitivity of the detector to the plane of polarization must either be measured or eliminated from the calculations. This is easy if an auxiliary polarizer is available that has a much higher degree of polarization than the one to be measured. In such a case, the "perfect" polarizer can be placed in the beam, and the transmittances T_1 and T_2 for the unknown polarizer can be measured directly.* Source polarization, instrumental polarization, and variation of detector response with plane of polarization can all be lumped together as a product. If this product is different in the horizontal and vertical planes, the ratio of the signals obtained when the "perfect" polarizer is oriented horizontally and vertically will not equal unity. One should always take more than the minimum number of measurements, i.e., introduce redundancy, to make sure that no systematic errors are present.

If a high-quality polarizer is not available, two polarizers having unknown properties may be used instead. Several procedures have been described in detail in the earlier polarization chapter.¹ The method of Hamm et al.¹⁴¹ which yields the extinction ratio of each polarizer and the instrumental polarization was described in detail and a brief summary of the method of Kudo et al.¹⁴² was given. The methods of Hamm et al.,¹⁴¹ Horton et al.,¹⁴³ and Schledermann and Skibowski¹⁴⁴ were specifically developed for non-normal incidence reflection polarizers (see "Brewster Angle Reflection Polarizers" section in Sec. 13.9).

13.9 NON-NORMAL-INCIDENCE REFLECTION AND TRANSMISSION POLARIZERS

By far the largest class of polarizers used in the infrared and ultraviolet spectral regions (where dichroic sheet polarizers and calcite polarizing prisms cannot be used) is the so-called *pile-of-plates polarizers* from which light is reflected (or transmitted) at non-normal incidence. Since most of these polarizers operate at angles near the Brewster or polarizing angle [see Eq. (48) in Chap. 12, "Polarization"], they are frequently called Brewster angle polarizers. The plane-parallel plates which are used for Brewster angle transmission polarizers (see "Brewster Angle Transmission Polarizers" section on p. 13.38) are generally thick enough to ensure that although multiple reflections occur within each plate, the

*When using an air-spaced polarizing prism, extreme care should be taken not to exceed the acceptance angle of the prism.

coherence of the light beam is lost and there are no interference effects. However, another class of non-normal-incidence transmission polarizers makes use of interference effects to enhance their polarizing properties. These include interference polarizers (see “Interference Polarizers” section on p. 13.39) and polarizing beam splitters (see “Polarizing Beam Splitters” section on p. 13.41). These thin-film devices are discussed in much more detail in Chap. 7, “Optical Properties of Films and Coatings,” in Vol. IV by Jerzy A. Dobrowolski. A relation which is frequently used in connection with non-normal-incidence reflectance measurements is the Abelès condition, discussed in the following section.

Brewster Angle Reflection Polarizers

Most reflection-type polarizers are made of plates which are either nonabsorbing or only slightly absorbing. The angle of incidence most often used is the Brewster angle at which the reflection of the p component, light polarized parallel to the plane of incidence, goes to 0. Thus the reflected light is completely plane polarized with the electric vector vibrating perpendicular to the plane of incidence (s component). Curves showing the reflectance and extinction ratio for various materials and angles near the Brewster angle are given in Fig. 5 of Chap. 12, “Polarization.” The polarizing efficiency of reflection-type polarizers can be experimentally determined using any of the methods given in Par. 49 of the earlier polarization chapter;¹ the methods of Hamm et al.,¹⁴¹ Horton et al.,¹⁴³ and Schleder and Skibowski¹⁴⁴ were specifically developed for polarizers of this type.

Brewster angle reflection polarizers for the infrared are made from the semiconductors silicon, germanium, and selenium which are transparent beyond their absorption edges and have high refractive indexes. Table 3 lists various infrared polarizers which have been described in the literature. All involve external reflections except the Ge-Hg polarizer described by Harrick,¹⁴⁵ in which light undergoes two or four reflections within a bar of germanium. While Harrick’s polarizer has attractive features, it depends on maintaining polarization in the germanium, so that great care must be taken to obtain material with a minimum of strain birefringence.

In the ultraviolet, materials such as LiF, MgF₂, CaF₂, and Al₂O₃, can be used as polarizers. Biotite, a form of mica, has also been found to perform very well in the 1000- to 6000-Å region. In the extreme ultraviolet, metallic films, particularly Au, Ag, and Al, have been used as polarizers. Table 4 lists various non-normal-incidence ultraviolet reflection polarizers as well as authors who have made calculations and measurements on various materials for ultraviolet polarizers.

TABLE 3 Infrared Brewster Angle Reflection Polarizers

Material	Description	Reference
Ge-Hg	Multiple internal reflections in Ge immersed in Hg	Harrick ¹⁴⁵
Ge	Single external reflection from 1-cm-thick polished Ge single crystal	Edwards and Bruemmer ¹⁴⁶
Ge	Proposed parallel and antiparallel arrangements of two Ge plates	Krizek ¹⁴⁷
Ge	Double-beam system: beam 1, single reflection; beam 2, one transmission, one reflection	Craig et al. ¹⁴⁸
Ge	Axial arrangement with reflections from two Ge wedges and two Al mirrors	Bor and Brooks ¹⁴⁹
Se	Reflections from two cast-Se films on roughened glass plates	Pfund ¹⁵⁰
Se	Axial arrangement with reflections from two Se films evaporated on NaCl and one Ag mirror	Barchewitz and Henry ¹⁵¹
Se	Large-aperture, axial, venetian-blind arrangement with one or two reflections from evaporated Se films on roughened glass plates (additional reflections from Al mirrors)	Takahashi ¹⁵²
Si	Single reflection from polished single crystal Si	Walton and Moss ¹⁵³
Si	Axial arrangement with reflection from two Al mirrors and polished Si plate with roughened back	Baumel and Schnatterly ¹⁵⁴
PbS	Axial arrangement with reflections from two chemically deposited PbS films and one Al film	Grechushnikov and Petrov ¹⁵⁵
CdTe	Single plate	Leonard ¹¹⁷
Al + Al ₂ O ₃	Multiple reflections from Al ₂ O ₃ coated with metal at 10.6 μm (calculations only)	Cox and Hass ¹⁵⁶
Ti + SiO ₂	Multiple reflections from dielectric coated Ti at 2.8 μm (calculations only)	Thonn and Azzam ¹⁵⁷

TABLE 4 Ultraviolet Reflection Polarizers and Polarization Measurements

Material	Description	Wavelength range, Å	Reference
Al ₂ O ₃ , Al, Au, ZnS, glass, and others	Calculated values of R_s and $(R_s/R_p)_{\max}$ vs. wavelength for a single reflection	500–2000	Hunter ¹⁵⁸
Al ₂ O ₃ , Al, glass, and others	Calculated values of R_s and $(R_s - R_p)/(R_s + R_p)$ vs. angle of incidence for a single reflection; also principal angle and related angles	584	Damany ¹⁵⁹
Al ₂ O ₃ , CaF ₂ , LiF, and Pyrex	Measured optical constants; calculated R_s and $(R_s - R_p)/(R_s + R_p)$ vs. angle of incidence and wavelength for a single reflection	200–2000	Stephan et al. ¹⁶⁰
Al ₂ O ₃ and CaF ₂	Measured $(R_s - R_p)/(R_s + R_p)$ vs. angle of incidence at selected wavelengths for a single reflection; used both materials as single-reflection polarizers	1026–1600	de Chelle and Merdy ¹⁶¹
LiF, Al ₂ O ₃ , MgF ₂ , SiO ₂ , ZnS	Used single-reflection LiF polarizer at Brewster angle to measure R_s and R_p/R_s for various materials; best polarizers were Al ₂ O ₃ and MgF ₂	1216	McIlraith ¹⁶²
Al, Ag, Au, MgF ₂ , SiO ₂ , ZnS	Measured polarization of uncoated aluminum grating and optical constants of all materials listed	304–1216	Cole and Oppenheimer ¹⁶³
Al, Au	Determined polarization of Au- and Al-coated gratings by measuring reflectance of Au and fused-silica mirrors at 45°	600–2000	Uzan et al. ¹⁶⁴
Al, Au, glass	Measured average reflectance and degree of polarization of Al, Au, and glass as a function of angle of incidence; measured polarization of a glass grating and an Al-coated grating	584	Rabinovitch et al. ¹⁶⁵
MgF ₂	Measured R_p/R_s at 60° for a single reflection and compared it with calculated values	916, 1085, 1216	Sasaki and Fukutani ¹⁶⁶
MgF ₂ + Al	Calculated performance, constructed axial triple-reflection polarizer and analyzer of MgF ₂ -coated Al, and measured transmission	1216	Winter et al. ¹⁶⁷
MgF ₂ , MgF ₂ + Al	Calculated performance, constructed triple-reflection polarizer of a MgF ₂ plate and two MgF ₂ -coated Al mirrors, and measured transmission	300–2000	Hass and Hunter ¹⁶⁸
MgF ₂ , MgF ₂ + Al	Constructed four-reflection polarizers of a MgF ₂ plate and three MgF ₂ -coated Al mirrors, no performance properties measured; polarizer part of the UV spectrometer and polarimeter for the NASA Solar Maximum Mission	1150-visible	Spencer et al. ¹⁶⁹

TABLE 4 Ultraviolet Reflection Polarizers and Polarization Measurements (*Continued*)

Material	Description	Wavelength range, Å	Reference
MgF ₂ , Au, and other metals	Calculated values of R_s and $(R_s/R_p)_{\max}$ vs. angle of incidence and wavelength for one and more reflections for a variety of materials	300–2000	Hunter ⁷⁰
Au, Ag	Measured R_p/R_s at 45° for a single reflection and compared it with calculated values; measured R_p/R_s for a platinumized grating	500–1300	Hamm et al. ¹⁴¹
Au	Used single-reflection Au mirror at 60° as polarizer (Brewster angle about 55°)	600–1200	Ejiri ¹⁷¹
Au	Used axial arrangement of eight Au mirrors at 60° as polarizer and analyzer to measure polarization of synchrotron radiation; determined polarizing properties of each polarizer	500–1000	Rosenbaum et al. ¹⁷²
Au	Constructed axial triple-reflection Au polarizer; measured extinction ratio and transmission for different angles of incidence on Au plates	500–5000	Horton et al. ¹⁴³
Au	Reflection from two cylindrical gold mirrors in a Seya-Namioka monochromator; measured polarization ratio	1200–3000	Rehfeld et al. ¹⁷³
Au	Calculated performance of a polarizer made of 2 concave Au-coated spherical mirrors used off axis, constructed polarizer, no measurements made of polarization or transmission	584	Van Hoof ¹⁷⁴
Au	Constructed a 4-reflection Au-coated polarizer of Van Hoof's design (2 plane, 2 spherical mirrors), measured transmission and degree of polarization	400–1300	Hibst and Bukow ¹⁷⁵
Au	Constructed a single-reflection Au-coated polarizer and measured the polarizing efficiency	584	Khakoo et al. ¹⁷⁶
Biotite	Constructed axial polarizer and analyzer each with 61° Brewster angle reflection from biotite and two reflections from MgF ₂ -coated Al mirrors; measured transmission and extinction ratio	1100–6000	Robin et al. ¹⁷⁷
Biotite	Constructed two polarizers: (1) axial polarizer with two 60° reflections from biotite and 30° reflection from MgF ₂ -coated Al mirror; (2) displaced-beam polarizer with 60° reflections from two biotite plates; measured degree of polarization of various gratings	1000–2000	Matsui and Walker ¹⁷⁸

The most versatile non-normal-incidence reflection polarizer would be one which does not deviate or displace the beam from its axial position. One convenient arrangement would be a symmetric three-reflection system in which the light is incident on one side of a triangle, reflected to a plane mirror opposite the apex, and back to the other side of the triangle, as was done by Horton et al.,¹⁴³ and Barchewitz and Henry.¹⁵¹ If the polarizer must have a good extinction ratio and the light beam is highly convergent, two of the reflections could be from the polarizing material and the third from a silvered or aluminized mirror. If the beam is highly collimated or more throughput is required, only one reflection may be from the polarizing material. The throughput can also be increased by using a plane-parallel plate for the polarizing reflection. The major drawback to a reflection polarizer is the extreme length of the device required to accommodate a beam of large cross-sectional area. For example, if a germanium polarizer were used at the Brewster angle (76°) and the beam width were about 25 mm, each Ge plate would have to be about $25 \text{ mm} \times 100 \text{ mm}$ and the overall length of the polarizer would be greater than 200 mm if a three-reflection axial arrangement such as that described above were used.

The Abelès condition,¹⁷⁹ which applies to the amplitude reflectance at 45° angle of incidence (see Sec. 12.5 in Chap. 12, "Polarization") is useful for testing the quality of reflection polarizers. Schulz and Tangherlini¹⁸⁰ apparently rediscovered the Abelès condition and used the ratio $R_s^2/R_p = 1$ as a test to evaluate their reflecting surfaces. They found that surface roughness made the ratio too small but annealing the metal films at temperatures higher than 150°C made the ratio larger than unity. Rabinovitch et al.¹⁶⁵ made use of the Abelès condition to determine the polarization of their Seya-Namioka vacuum-ultraviolet monochromator. They measured the reflectance at 45° of a sample whose plane of incidence was perpendicular or parallel to the exit slit. From these measurements they deduced the instrumental polarization by assuming the Abelès condition. Values of instrumental polarization obtained using carefully prepared gold and fused-silica samples were in excellent agreement, showing that neither of these materials had surface films which invalidated the Abelès condition. Surface films usually have relatively little effect on the Abelès condition in the visible region¹⁸¹ but become important in the vacuum ultraviolet. Hamm et al.¹⁴¹ eliminated the effect of instrumental polarization from their measurements of the reflectance of a sample in unpolarized light at 45° angle of incidence by making use of the Abelès condition. Although McIlrath¹⁶² did not refer to the Abelès condition as such, he used it to determine the instrumental polarization of his vacuum-ultraviolet apparatus so he could measure the absolute reflectance of a sample at 45° angle of incidence. Thonn and Azzam¹⁵⁷ have calculated the polarizing properties of dielectric-coated metal mirrors at $2.8 \mu\text{m}$ in the infrared. Reflections from 2, 3, or 4 such mirrors at the Brewster angle should give excellent performance, although the polarizer would be quite long.

Brewster Angle Transmission Polarizers

To help overcome the beam-deviation problem and the extreme length of reflection-type polarizers, Brewster angle polarizers are often used in transmission, particularly in the infrared, where transparent materials are available. At the Brewster angle, all of the p component and an appreciable fraction of the s component are transmitted. Thus, several plates must be used to achieve a reasonable degree of polarization. The higher the refractive index of the plates, the fewer are required.

Tables 1 and 2 in Chap. 12, "Polarization" give equations for the transmittances and degree of polarization for a single plate and multiple plates at any angle of incidence in terms of R_s and R_p for a single surface, as well as these same quantities at the Brewster angle. Conn and Eaton¹⁸² have shown that the formulas which assume incoherent multiple reflections within each plate and none between plates give the correct degree of polarization for a series of Zapon lacquer films ($n = 1.54$) and also for a series of eight selenium films, whereas the formula of Provostaye and Desains¹⁸³ predicted values which were much too low. These authors also point out that the number of multiply reflected beams between plates that enter the optical system depends on the spacing between plates and the diaphragm used to limit the number of beams. One can use a fanned arrangement, as suggested by Bird and Shurcliff,¹⁸⁴ to eliminate these multiply reflected beams. Internal reflections within each plate can be removed by wedging the plates.¹⁸⁴

Most of the infrared Brewster angle transmission polarizers described in the literature have been made of selenium, silver chloride, or polyethylene sheet; they are listed in Table 5. For wavelengths

TABLE 5 Infrared Brewster Angle Transmission Polarizers

Material	Description	Wavelength range, μm	Reference
Se	5 or 6 unbacked films (4 μm thick) at 65° angle of incidence (Brewster angle 68.5°)	2–14	Elliott and Ambrose ¹⁸⁵ Elliott et al. ¹⁸⁶
Se	5 unbacked films at 65° incidence (different method of preparation from above)		Ames and Sampson ¹⁸⁷
Se	8 unbacked films at the Brewster angle		Conn and Eaton ¹⁸²
Se	Se films (3–8 μm thick) evaporated on one side of collodion films; 68.5° angle of incidence	1–15	Barhewitz and Henry ¹⁵¹
Se	1–6 unbacked films (1.44–8 μm thick) at 68° angle of incidence	Visible–20	Duverney ¹⁸⁸
Se	3 unbacked films (0.95 μm thick) at 71° angle of incidence	6–17	Hertz ¹⁸⁹
Se	5 Formvar films coated on both sides with Se (various thicknesses) at 65° angle of incidence	1.8–3.2	Buijs ¹⁹⁰
Se	Unbacked films (different method of preparation from Elliott et al. ¹⁸⁶)		Bradbury and Elliott ¹⁹¹
Se	4 to 6 plated unsupported films (4–8 μm thick) at the Brewster angle	2.5–25	Greenler et al. ¹⁹²
AgCl	3 plates (1 mm thick) at 63.5° angle of incidence	Visible–15	Wright ¹⁹³
AgCl	6–12 plates (0.05 mm thick) at 60–75° angle of incidence	2–20	Newman and Halford ⁷⁸
AgCl	6 plates (0.5 mm thick) at 63.5° stacked in alternate directions		Makas and Shurcliff ⁹⁴
AgCl	Suggest 6 wedge-shaped plates at 68° stacked in alternate directions in a fanned arrangement		Bird and Shurcliff ¹⁸⁴
AgCl	2 V-shaped plates (3.2 mm thick) at 77° angle of incidence; large aperture		Bennett et al. ¹⁹⁵
KRS-5	1–3 thallium bromide-iodide plates (1 and 4 mm thick) at polarizing angle	1–15	Lagemann and Miller ¹⁹⁶
ZnS	4 glass plates (0.1 mm thick) coated on both sides with uniform ZnS films of same thickness (several sets of plates to cover extended wavelength range)	Visible–6	Huld and Staffin ¹⁹⁷
ZnSe	6 plates, extinction ratio of 800 at 4 μm	4	Leonard et al. ¹⁹⁸
Ge	1 single-crystal Ge plate (0.8 mm thick) at 76° angle of incidence		Meier and Günthard ¹⁹⁹
Ge	2 plates (1 mm thick) in an X-shaped arrangement at 76° angle of incidence		Harrick ²⁰⁰
Ge	3 plates (2 wedged) at the Brewster angle	2–6	Murarka and Wilner ²⁰¹
Polyethylene	12 sheets (8 μm thick) at the Brewster angle	6–20 (absorption bands 6–14)	Smith et al. ²⁰²
Polyethylene	9–15 sheets (20–50 μm thick) at the Brewster angle (55°)	30–200	Mitsuishi et al. ¹³⁷
Polyethylene	4 sheets at the Brewster angle	200–350	Hadni et al. ¹³⁶
Polyethylene	12 sheets (5 μm thick) at the Brewster angle	1.5–13 (selected wavelengths)	Walton and Moss ¹⁵³
Polyethylene	1–15 stretched sheets (12.7 μm thick) at the Brewster angle	10.6	Rampton and Grow ²⁰³
Polyethylene	20 sheets (30 μm thick) at the Brewster angle	45–200	Munier et al. ²⁰⁴
Polyethylene	25–30 sheets at the Brewster angle	54.6	Hilton and Jones ¹¹¹
Melinex	11–13 polyethylene terephthalate (Melinex) sheets (4.25–9 μm thick) at the Brewster angle	1–5	Walton et al. ²⁰⁵

TABLE 6 Ultraviolet Brewster Angle Transmission Polarizers

Material	Description	Wavelength range, Å	Reference
LiF	4–8 plates (0.3–0.8 mm thick) at 60° angle of incidence (Brewster angle 55.7–58.7°) stacked in alternate directions	1200–2000	Walker ²⁰⁶
LiF	8 plates	1100–3000	Hinson ²⁰⁷
LiF	8 plates (0.25–0.38 mm thick) at 60° angle of incidence stacked in groups of 4 in alternate directions	1200–2000	Heath ²⁰⁸
CaF ₂	4 to 8 wedged plates stacked in alternate directions in fanned arrangement at 65° angle of incidence (Brewster angle 56.7)	1500–2500	Schellman et al. ²⁰⁹
Al	Calculations of polarizing efficiency for 1000-Å-thick unbacked Al film, 1000- and 500-Å Al films each covered with 30-Å Al ₂ O ₃ and 100-Å Au films	300–800	Hunter ¹⁵⁸

longer than 3 μm , where calcite polarizing prisms become highly absorbing, to about 10 μm , beyond which wire-grid polarizers have good extinction ratios, Brewster angle transmission polarizers are the most useful, since the better-extinction, on-axis reflection-type polarizers (see “Brewster Angle Reflection Polarizers” section on p. 13.34) are impossibly long. Some of the interference polarizers described in the following sections “Interference Polarizers” and “Polarizing Beam Splitters” are superior if the beam-convergence angle is small. Ultraviolet Brewster angle transmission polarizers are not nearly as common; LiF and CaF₂ have mainly been used from about 1500 to 2500 Å (see Table 6). In the wavelength region where calcite polarizing prisms are usable (>2140 Å), Brewster angle polarizers have the advantage of a larger linear aperture and less absorption.

Low-absorption glass pile-of-plates polarizers have been used in the visible spectral region by Weiser,²¹⁰ in preference to more absorbing Glan-Thompson prism polarizers, to increase the power output of giant-pulse ruby lasers. Weinberg²¹¹ calculated the degree of polarization of glass and silver chloride plates, but he did not calculate the transmittance of his polarizers.

Interference Polarizers

When the sheets or films constituting a non-normal-incidence transmission polarizer are thin and have very smooth surfaces, the internally reflected beams can interfere constructively or destructively. In this case, the transmittance of the p component remains unity at the Brewster angle (where $R_p = 0$) and only oscillates slightly (with respect to wavelength) for angles close to the Brewster angle. However, the s transmittance varies from a maximum of unity to a minimum of $(1 - R_s)^2 / (1 + R_s)^2$ whenever λ changes by an amount that will make the quantity $(nd \cos \theta_1) / \lambda$ in Eq. (26) in Chap. 12, “Polarization” change by $1/2$.^{*} These transmittance oscillations are only ± 0.225 for a single film of refractive index 1.5 but can become as large as ± 0.492 when $n = 4.0$. Since the p transmittance remains essentially constant, the extinction ratio will vary cyclically with the s transmittance, as can be seen in the upper curve of Fig. 17 for a 2.016- μm -thick selenium film.

If a transmission polarizer with a good extinction ratio is needed for use over a limited wavelength range, it can be made of several uniform films of a thickness that yields a minimum extinction ratio in the given wavelength region. The extinction ratio for a series of m films is $(T_s/T_p)^m$ when there are no multiple reflections between them. In this way only *half* as many films would be needed to achieve a given extinction ratio as would be necessary if interference effects were not present. This rather surprising result can be seen from the expressions for $(T_s)_{\text{sample}}$ for m plates with and without interference effects in Table 2 in Chap. 12, “Polarization.” Assuming no multiple

^{*}The approximate expression for this wavelength interval $\Delta\lambda$ (assuming that the oscillations are sufficiently close together for $\lambda_1, \lambda_2 \approx \lambda_2$) is given in Eq. (27) in Chap. 12, “Polarization.”

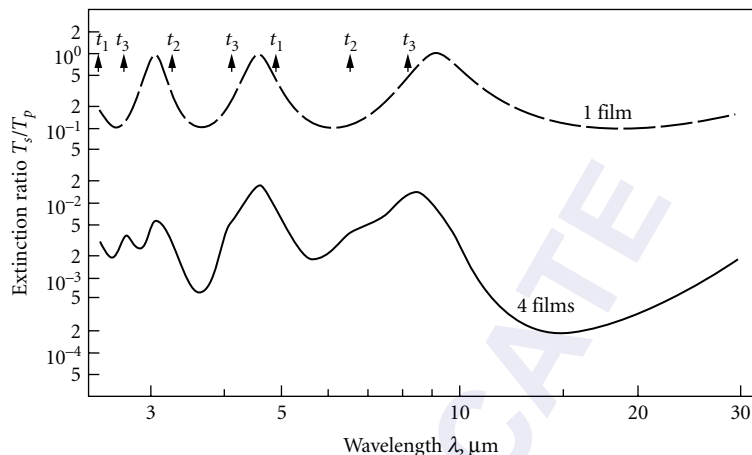


FIGURE 17 Calculated extinction ratios for a series of selenium films ($n = 2.46$) as a function of wavelength from 2.5 to 30 μm . Light is incident at the Brewster angle, 67.9° , and multiply reflected beams interfere within the film. The upper curve is for a single film 2.016 μm thick; arrows indicate positions of maxima for three thinner films: $t_1 = 1.080 \mu\text{m}$, $t_2 = 1.440 \mu\text{m}$, and $t_3 = 1.800 \mu\text{m}$. The lower curve is the extinction ratio for the four films in series assuming no reflections between films. The calculated p transmittance for each film (and for four films in series) is unity at the Brewster angle.

reflections between plates, the expressions are $[2n^2/(n^4 + 1)]^{2m}$ and $[2n^2/(n^4 + 1)]^m$, respectively. Hertz¹⁸⁹ achieved a degree of polarization of 99.5 percent in the 6- to 17- μm region using three unbacked selenium films 0.95 μm thick. Conn and Eaton¹⁸² obtained only a slightly better performance with eight thicker nonuniform selenium films.

As can be seen in Fig. 17, the calculated extinction ratio for the 2.016- μm -thick film goes to unity at 3.0, 4.6, and 9.2 μm , indicating that the s as well as the p transmittance at these wavelengths is unity. This ratio will remain unity at the above wavelengths if there are several nonabsorbing films of the same thickness. Even if the films have slightly different thicknesses, or if their surfaces are somewhat rough, interference effects may still persist, adversely affecting polarizer performance. Such effects have been observed by Elliott et al.,¹⁸⁶ Barchewitz and Henry,¹⁵¹ Duverney,¹⁸⁸ Mitsuishi et al.,¹³⁷ and Walton et al.²⁰⁵

By choosing films of appropriate thicknesses, interference effects can be used to advantage. The lower curve in Fig. 17 shows the extinction ratio obtained if four selenium films of thicknesses 1.08, 1.44, 1.80, and 2.02 μm are used at the Brewster angle as a transmission polarizer. (The wavelengths at which maxima occur for the three thinner films are indicated by arrows in the upper portion of the figure.) In this example the extinction ratio for the four films in series is better than 2×10^{-2} from 2.5 to 30 μm and at most wavelengths is better than 10^{-2} (corresponding to a degree of polarization in excess of 98 percent). In the 11- to 27- μm wavelength region the extinction ratio is better than 10^{-3} . Four thick or nonuniform selenium films without interference effects have a calculated extinction ratio of about 10^{-2} , and six films are required to change this ratio to 10^{-3} . Thus, in the 11- to 27- μm wavelength region, *four* selenium films of appropriate thicknesses *with interference* have a superior extinction ratio to *six* selenium films *without interference*. If one wishes to optimize the extinction ratio over a more limited wavelength range, the film thicknesses can be adjusted accordingly and the extinction ratio improved. Unfortunately, the gain in extinction ratio is offset by a more sensitive angular function than that shown in Fig. 7 in Chap. 12., "Polarization," so that the incident beam must be very well collimated.

Interference effects can also be used to advantage in other types of non-normal-incidence polarizers. Bennett et al.²¹² made a transmission polarizer from a series of four germanium films

(ranging in thickness from 0.164 to 0.593 μm), evaporated onto strain-free plates of sodium chloride. The plates were inclined at the Brewster angle for germanium and arranged in the form of an X so that the polarizer would have a large square aperture and would not deviate the beam. An extinction ratio better than 3×10^{-3} was measured at 2.5 μm , and the plates transmitted from 2 to 13 μm . (Calculated extinction ratios in this wavelength range vary from 1×10^{-3} to 2×10^{-4} for radiation incident at the Brewster angle.)

Polarizers consisting of a high refractive index transparent film on a lower refractive index transparent substrate have been suggested for use in the visible wavelength region by Schröder²¹³ and Abelès.²¹⁴ These still have a Brewster angle where $R_p = 0$, and furthermore R_s at this angle is greatly increased over its value for an uncoated low refractive index substrate. Thus, a large-aperture, high-efficiency polarizer with no absorption losses is possible, which should find numerous applications in laser systems. One polarizer of this type, suggested independently by Schröder and by Abelès, would consist of high refractive index titanium dioxide films ($n \approx 2.5$) evaporated onto both sides of a glass substrate ($n = 1.51$). At the Brewster angle, 74.4° , $R_s \approx 0.8$, making this polarizer equivalent to one made from a material of refractive index 4 ($\theta_B = 76.0^\circ$ as shown in Fig. 4 in Chap. 12, "Polarization").* Two glass plates coated on both sides with TiO_2 films should have an extinction ratio of about 1.6×10^{-3} at 5500 Å and about twice that value at the extreme ends of the visible region, according to Abelès.²¹⁴ Schröder²¹³ measured the degree of polarization as a function of angle of incidence for one such TiO_2 -coated glass plate and found values comparable to the calculated ones. Kubo²¹⁵ calculated the degree of polarization, reflectance, and transmittance (as a function of angle of incidence and wavelength) of a glass plate ($n = 1.50$) covered with a thin transparent film of index, 2.20. His results are similar to those of Abelès and Schröder.

Schopper,²¹⁶ Ruiz-Urbieto and Sparrow,²¹⁷⁻²¹⁹ and Abelès²²⁰ have also investigated making non-normal-incidence reflection polarizers from a thin transparent or absorbing film deposited onto an absorbing substrate. Zaghoul and Azzam²²¹ proposed using silicon films on fused silica substrates as reflection polarizers for different mercury spectral lines in the visible and ultraviolet regions. Abelès designed some specialized reflection polarizers for use in the vacuum ultraviolet. Unfortunately the wavelength range covered by such a polarizer is very narrow; for one polarizer it was 25 Å at a wavelength of 1500 Å. However, the spectral range could possibly be increased by using several thin films instead of one.

Multilayer film stacks have also been used to produce non-normal-incidence reflection or transmission polarizers by Buchman et al.²²² Buchman²²³ later improved the design performance of his polarizers by adding antireflection layers between the repeating groups of layers. Although this type of polarizer has a relatively narrow operating bandwidth, a small angular acceptance, tight wavelength centering, and layer thickness uniformity requirements, it can be used successfully in high power laser systems as shown by Refermat and Eastman.²²⁴ Songer²²⁵ described how to design and fabricate a Brewster angle multilayer interference polarizer out of a titanium dioxide, silicon dioxide multilayer on BK 7 glass for use in a 1.06- μm laser beam. Blanc, Lissberger, and Roy²²⁶ designed, built, and tested multilayer zinc sulfide-cryolite-coated glass and quartz polarizers for use with a pulsed 1.06- μm laser. Recently, Maehara et al.²²⁷ have reported excellent performance for a pair of polarizers coated with 21 ruthenium and silicon films on a silicon wafer over a wide wavelength range in the soft x-ray region. In several designs of multilayer film stacks, both the reflected and transmitted beams are used; they are discussed in the following section.

Polarizing Beam Splitters

Polarizing beam splitters are a special form of non-normal-incidence interference polarizers in which the beam is incident on a multilayer dielectric stack at 45° . The transmitted beam is almost entirely plane-polarized in the p direction, while the reflected beam is nearly all plane-polarized

*We are assuming no multiply reflected beams within the substrate in either case.

in the s direction. Generally the alternating high and low refractive index dielectric layers are deposited onto the hypotenuses of two right-angle prisms, which are then cemented together to form a cube. The beam enters a cube face normally and strikes the multilayers on the hypotenuse (the high refractive index layer is next to the glass), and the reflected and transmitted beams emerge normal to cube faces, being separated by 90° . Clapham et al.²²⁸ have a good discussion of polarizing beam splitters, which were invented by S. M. MacNeille²²⁹ and developed by Banning.²³⁰ Banning's beam splitter was made with three zinc sulfide and two cryolite layers on each prism; the polarization for white light was greater than 98 percent over a 5° -angle on each side of the normal to the cube face for both the reflected and transmitted beams. Variations on this design have since been proposed by Dobrowolski and Waldorf,²³¹ Monga et al.,²³² and Mouchart et al.,²³³ primarily to improve the laser damage resistance of the device and increase the angular field of view. Dobrowolski and Waldorf²³¹ designed and built a polarizing beam splitter consisting of a multilayer coating of HfO_2 and SiO_2 deposited onto fused silica and immersed in a water cell that acted like the MacNeille cube. Tests with a $0.308\ \mu\text{m}$ excimer laser showed a high laser damage threshold. The multi-wavelength polarizing beam splitters designed by Monga et al.²³² could be made in large sizes and could withstand high laser power levels. The modified MacNeille cube polarizers designed by Mouchart et al.²³³ had angular fields of view that could be increased to about $\pm 10^\circ$ when the polarizers were used with monochromatic light sources.

Lees and Baumeister²³⁴ designed a frustrated total internal reflection beam splitter that had a multilayer dielectric stack deposited onto the hypotenuse of a prism. Their designs, for use in the infrared spectral region, consisted of multilayer stacks of PbF, and Ge deposited onto a germanium prism and covered by a second germanium prism. Azzam²³⁵ designed polarization-independent beam splitters for $0.6328\ \mu\text{m}$ and $10.6\ \mu\text{m}$ using single-layer coated zinc sulfide and germanium prisms. The devices were found to be reasonably achromatic and their beam-splitting ratio could be varied over a wide range with little degradation in polarization properties. Azzam²³⁶ also proposed coating a low-refractive-index dielectric slab on both sides with high-refractive-index dielectric films to make an infrared polarizing beam splitter.

Various high- and low-refractive-index materials have been successfully used in the multilayer stacks. In addition to zinc sulfide and cryolite on glass by Banning²³⁰ and Schröder and Schläfer,²³⁷ layers of a controlled mixture of silicon dioxide and titanium dioxide have been alternated with pure titanium dioxide on fused-silica prisms by Pridatko and Krylova,²³⁸ thorium dioxide and silicon dioxide have been used on fused-silica prisms by Sokolova and Krylova,²³⁹ chiolite (a mixture of sodium and aluminum fluorides) and lead fluoride have been used on fused-silica prisms by Turner and Baumeister,²⁴⁰ bismuth oxide and magnesium fluoride have been used on EDF glass prisms by Clapham et al.,²²⁸ and zirconium oxide and magnesium fluoride have been used on dense flint-glass prisms by Clapham et al.²²⁸ The calculations involved in optimizing these beam splitters for good polarizing characteristics, achromaticity, and relative insensitivity to angle of incidence are quite involved. Clapham et al.²²⁸ and Turner and Baumeister²⁴⁰ discuss various calculational techniques frequently used. Clapham²⁴¹ also gives the measured characteristics of a high-performance achromatic polarizing beam splitter made with zirconium oxide and magnesium fluoride multilayers.

Although polarizing beam splitters are generally designed so that the s and p polarized beams emerge at right angles to each other, Schröder and Schläfer²³⁷ have an ingenious arrangement in which a half-wave plate and mirror are introduced into the path of the reflected beam to make it parallel to the transmitted beam and of the same polarization. Other optical schemes to accomplish the same purpose have been described in a later paper.²⁴²

For some purposes it is desirable to have a beam splitter that is insensitive to the polarization of the incident beam. Baumeister²⁴³ has discussed the design of such beam splitters made from multilayer dielectric stacks of alternating low- and high-refractive-index materials. One of his designs is composed of six dielectric layers for which the extinction ratio T_s/T_p varies from 0.93 to 0.99 in a bandwidth of about $800\ \text{Å}$, with a $\pm 1^\circ$ variation in the angle of incidence. In principle, any multilayer filter which is nonreflecting at normal incidence will be nonpolarizing at all angles of incidence, according to Baumeister.²⁴⁴ Costich²⁴⁵ has described filter designs for use in the near infrared which are relatively independent of polarization at 45° angle of incidence.

13.10 RETARDATION PLATES

Introduction

The theory of retardation plates and especially quarter-wave retarders is given in Chap. 12, "Polarization." The basic relation for retardation plates, Eq. (73) in that section, is

$$N\lambda = d(n_e - n_o) \quad (9)$$

where n_o = refractive index of the ordinary ray, n_e = refractive index of the extraordinary ray, d = physical thickness of the plate, and λ = wavelength.

Retardation plates are generally made of mica, stretched polyvinyl alcohol, and quartz, although other stretched plastics such as cellophane, Mylar, cellulose acetate, cellulose nitrate, sapphire, magnesium fluoride, and other materials can also be used (see West and Makas²⁴⁶). Polyvinyl alcohol in sheet form transmits well into the ultraviolet beyond the cutoff for natural mica and is thus particularly useful for ultraviolet retardation plates, according to McDermott and Novick.⁹¹ As suggested by Jacobs et al.,²⁴⁷ permanent birefringence can be thermomechanically induced in the borosilicate optical glass ARG-2, making it an attractive alternate to natural crystalline quartz and mica for large aperture wave plates for laser systems. Refractive indexes and birefringences of some materials are listed in Tables 7 and 8. The birefringences reported for mica and apophyllite should be considered as approximate, since they are measurements made on single samples. There is good reason to believe that the birefringence of apophyllite may be different for other samples, (see "Achromatic Retardation Plates" section on p. 13.48). Although calcite would seem at first to be a good material for retardation plates, its birefringence is so high that an extremely thin piece, less than 1 μm , would be required for a single $\lambda/4$ retardation plate. If a "first-order" or multiple-order plate were constructed (see sections "First-Order Plates" on p. 13.46 and "Multiple-Order Plates" on p. 13.47), or if calcite were used as one component of an achromatic retardation plate (see "Achromatic Retardation Plates" section), the tolerance on the thickness would be very stringent.

Retardation plates are generally made of a single piece of material, although when the thickness required for a plate is too small, two thicker pieces may be used with the fast axis of one aligned parallel to the slow axis of the other to cancel out all but the desired retardation. Plates which are a little too thin or a little too thick may be rotated about an axis parallel or perpendicular to the optic axis

TABLE 7 Refractive Indices of Selected Materials at 5893 Å²⁴⁸

Material	n_o	n_e
Positive Uniaxial Crystals		
Ice, H ₂ O	1.309	1.313
Sellaite, MgF ₂	1.378	1.390
Apophyllite, 2[KCa ₄ Si ₈ O ₂₀ (F, OH)·8H ₂ O]	1.535±	1.537±
Crystalline quartz, SiO ₂	1.544	1.553
Diopase, CuSiO ₃ ·H ₂ O	1.654	1.707
Zircon, ZrSiO ₄	1.923±	1.968±
Rutile, TiO ₂	2.616	2.903
Negative Uniaxial Crystals		
Beryl (emerald), Be ₃ Al ₂ (SiO ₃) ₆	1.581±	1.575±
Sodium nitrate, NaNO ₃	1.584	1.336
Muscovite mica (complex silicate)	1.5977±	1.5936±
Apatite, Ca ₁₀ (F, Cl) ₂ (PO ₄) ₆	1.634	1.631
Calcite, CaCO ₃	1.658	1.486
Tourmaline (complex silicate)	1.669±	1.638±
Sapphire, Al ₂ O ₃	1.768	1.760

TABLE 8 Birefringence $n_e - n_o$ of Various Optical Materials^a

Wave-length, μm	Rutile ^b TiO ₂	CdSe ^b	Cryst- alline quartz ^{c-g}	MgF ₂ ^{cd,hi}	CdS ^{f-l}	Apo- phyllite ^m	ZnS/ (Wurt- zite)	Calcite ^e	LiNbO ₃ ⁿ	BaTiO ₃ ^o	AdP ^p	KDP ^p	Sap- phire ^{h,q,r} (Al ₂ O ₃)	Mica ^f
0.15	—	—	0.0214	0.0143	—	—	—	-0.326	—	—	-0.0613	-0.0587	-0.0111	—
0.20	—	—	0.0130	0.0134	—	—	—	-0.234	—	—	-0.0543	-0.0508	-0.0097	—
0.25	—	—	0.0111	0.0128	—	—	—	-0.206	—	—	-0.0511	-0.0474	-0.0091	—
0.30	—	—	0.0103	0.0125	—	—	—	-0.193	—	—	-0.0492	-0.0456	-0.0087	—
0.35	—	—	0.0098	0.0122	—	—	—	-0.184	—	—	-0.0482	-0.0442	-0.0085	—
0.40	—	—	0.0096	0.0121	—	—	0.004	-0.179	—	—	-0.0473	-0.0432	-0.0083	-0.00457
0.45	0.338	—	0.00937	0.0120	—	0.0019	0.004	-0.176	-0.1049	-0.097	-0.0465	-0.0424	-0.0082	-0.00468
0.50	0.313	—	0.00925	0.0119	—	0.0022	0.004	-0.173	-0.0998	-0.079	-0.0458	-0.0417	-0.0081	-0.00476
0.55	0.297	—	0.00917	0.0118	0.014	0.0024	0.004	-0.172	-0.0947	-0.070	-0.0451	-0.0410	-0.0081	-0.00480
0.60	0.287	—	0.00909	0.0118	0.018	0.0026	0.004	-0.170	-0.0919	-0.064	-0.0444	-0.0403	-0.0080	-0.00482
0.65	0.279	—	0.00903	0.0117	0.018	0.0028	0.004	-0.169	-0.0898	—	-0.0438	-0.0396	-0.0080	-0.00483
0.70	0.274	—	0.00898	0.0117	0.018	—	0.004	-0.167	-0.0882	—	-0.0425	-0.0382	-0.0079	—
0.80	0.265	—	0.0089	0.0116	0.018	—	0.004	-0.165	-0.0857	—	-0.0411	-0.0367	—	—
0.90	0.262	—	0.0088	0.0115	0.018	—	0.004	-0.164	-0.0840	—	-0.0396	-0.0350	—	—
1.00	0.259	0.0195	0.0088	0.0114	0.018	—	0.004	-0.162	-0.0827	—	-0.0379	-0.0332	—	—
1.10	0.256	0.0195	0.0087	0.0114	0.018	—	0.004	-0.161	-0.0818	—	-0.0361	-0.0313	—	—
1.20	0.254	0.0195	0.0087	0.0114	0.017	—	0.004	-0.160	-0.0810	—	-0.342	-0.0292	—	—
1.30	0.252	0.0195	0.0086	0.0113	0.017	—	0.004	-0.158	-0.0804	—	-0.0321	-0.0269	—	—
1.40	0.251	0.0195	0.0085	0.0113	0.017	—	0.004	-0.157	-0.0798	—	-0.0298	-0.0245	—	—
1.50	0.250	0.0195	0.0085	0.0113	—	—	—	-0.156	-0.0788	—	-0.0274	-0.0219	—	—
1.60	0.249	0.0195	0.0084	0.0112	—	—	—	-0.154	-0.0782	—	-0.0248	-0.0191	—	—
1.70	0.248	0.0195	0.0084	0.0112	—	—	—	-0.153	-0.0777	—	-0.0221	-0.0162	—	—
1.80	0.247	0.0195	0.0083	0.0112	—	—	—	-0.151	-0.0774	—	-0.0192	-0.0130	—	—
1.90	0.246	0.0195	0.0082	0.0112	—	—	—	-0.150	-0.0771	—	-0.0161	-0.0097	—	—
2.00	0.246	0.0195	0.0081	0.0111	—	—	—	-0.148	-0.0766	—	—	—	—	—
2.10	0.245	0.0195	0.0081	0.0111	—	—	—	—	-0.0761	—	—	—	—	—
2.20	0.244	0.0195	0.0080	0.0111	—	—	—	—	-0.0752	—	—	—	—	—
2.30	0.243	0.0195	0.0079	0.0110	—	—	—	—	-0.0744	—	—	—	—	—
2.40	0.241	0.0195	0.0078	0.0110	—	—	—	—	-0.0739	—	—	—	—	—
2.50	—	0.0195	0.0077	0.0110	—	—	—	—	-0.0734	—	—	—	—	—
2.60	—	0.0195	0.0076	0.0110	—	—	—	—	—	—	—	—	—	—

^aCalculated values at 24.8°C obtained from analytical expressions are given for crystalline quartz, MgF₂, calcite, ADP, and KDP by Beckers³⁹
^bBond²⁵⁰

^cBallard et al.²⁵⁵

^dChandrasekharan and Damany²⁶¹

^eEnnos and Opperman²⁵⁶

^fPalik²⁵⁷

^gGobrecht and Bartschat²⁵⁸

^hShumate²⁵⁹

ⁱLoewenstein²⁶⁰

^jShields and Ellis²⁵¹

^kChandrasekharan and Damany²⁵³

^lPalik and Henvis²⁵²

^mBoyd et al.²⁵³

ⁿJeppesen²⁵⁴

^oMaillard²⁶²

^pBieniewski and Czyzak²⁶³

^qFrançon et al.²⁶⁴

^rZernike²⁶⁵

^s422. J. Zernike²⁶⁵

^tEinspoorn²⁶⁶

to change the retardation to the desired amount, as suggested by Gieselmann et al.,²⁶⁷ and Daniels.²⁶⁸ There are also some novel circular polarizers and polarization rotators for use in the far ultraviolet (see the papers by McIlrath,¹⁶² Saito et al.,²⁶⁹ and Westerveld et al.²⁷⁰), far infrared (Richards and Smith,²⁷¹ Johnston,²⁷² and Gonates et al.²⁷³), and visible region (Lostis,²⁷⁴ and Greninger²⁷⁵).

Achromatic retardation plates which have the same retardation over a range of wavelengths can be made from two or more different materials or from two or more plates of the same material whose axes are oriented at appropriate angles with respect to each other. These latter devices are known as composite plates (see “Composite Retardation Plates” section on p. 13.52 and the earlier polarization chapter¹), and although they can change plane-polarized light into circularly polarized light, they do not have all the other properties of true retardation plates. By far the most achromatic $\lambda/4$ retarders are devices, such as the Fresnel rhomb, which obtain their retardation from internal reflections at angles greater than the critical angle.

Mica retardation plates are mentioned in “Mica Retardation Plates” section and are discussed in detail in the earlier polarization chapter,¹ which includes the theory of multiple reflections; “Crystalline Quartz Retardation Plates” section on p. 13.46 is devoted to various types of crystalline-quartz retardation plates, and “Achromatic Retardation Plates” section on p. 13.48 covers all achromatic retardation plates, except those of the rhomb-type; the latter are mentioned in “Rhombs as Achromatic $\lambda/4$ Retarders” section on p. 13.52 and in detail by Bennett²⁷⁶ and also in the earlier polarization chapter.¹ Various types of composite plates and unusual retardation plates are also described in detail in Ref. 1.

Methods for making and testing quarter-wave plates including ways of splitting mica, how to distinguish between fast and slow axes, methods for measuring retardations close to $\lambda/4$, and the tolerance on plate thickness have all been described in detail in the earlier polarization chapter.¹ An additional paper by Nakadate²⁷⁷ shows how Young’s fringes can be used for a highly precise measurement of phase retardation.

Waveplates are all sensitive to some degree to temperature changes, variations in the angle of incidence, coherence effects in the light beam, and wavelength variations. Multiple-order plates are much more sensitive than “first-order” or single-order plates. Hale and Day²⁷⁸ discuss these effects for various types of waveplates and suggest designs that are less sensitive to various parameters.

Most retardation plates are designed to be used in transmission, generally at normal incidence. However, there are also reflection devices that act as quarter-wave and half-wave retarders and polarization rotators. In the vacuum ultraviolet, Westerveld et al.²⁷⁰ produced circularly polarized light by using Au-coated reflection optics. Saito et al.²⁶⁹ used an evaporated Al mirror as a retardation plate at 1216 Å, Lyman α radiation, following earlier work by McIlrath.¹⁶² Greninger²⁷⁵ showed that a three-mirror device could be used in place of a half-wave plate to rotate the plane of polarization of a plane-polarized beam and preserve the collinearity of input and output beams. Johnston²⁷² used a different three-mirror arrangement for the same application in the far-infrared. Thonn and Azzam¹⁵⁷ designed three-reflection half-wave and quarter-wave retarders from single-layer dielectric coatings on metallic film substrates. They showed calculations for ZnS-Ag film-substrate retarders used at 10.6 μm . Previously Zaghoul, Azzam, and Bashara^{279,280} had proposed using a SiO_2 film on Si as an angle-of-incidence tunable reflection retarder for the 2537-Å mercury line in the ultraviolet spectral region. Kawabata and Suzuki²⁸¹ showed that a film of MgF_2 on Ag was superior to Zaghoul et al.’s design at 6328 Å. They also performed calculations using Al, Cu, and Au as the metals and concluded that Ag worked best.

Mica Retardation Plates

Mica quarter-wave plates can be made by splitting thick sheets of mica down to the appropriate thickness, as described by Chu et al.,²⁸² and in the earlier polarization chapter.¹ Since the difference between the velocities of the ordinary and extraordinary rays is very small, the mica sheets need not be split too thin; typical thicknesses lie in the range 0.032–0.036 mm for yellow light. The fast and slow axes of a mica quarter-wave plate can be distinguished using Tutton’s test, as mentioned in Strong’s book,²⁸³ and the retardation can be measured using one of several rather simple tests.¹

If the mica sheets are used without glass cover plates, multiply reflected beams in the mica can cause the retardation to oscillate around the value calculated from the simple theory, as described in the earlier polarization chapter.¹ Fortunately this effect can be eliminated in one of several ways.¹ Mica does have one serious drawback. There are zones in the cleaved mica sheets which lie at angles to each other and which do not extinguish at the same angle, as noted by Smith.²⁸⁴ Thus, extinction cannot be obtained over the whole sheet simultaneously. In very critical applications such as ellipsometry, much better extinction can be obtained using quarter-wave plates made of crystalline quartz (“Crystalline Quartz Retardation Plates” section next), which do not exhibit this effect. Properties of mica quarter-wave plates and methods for making and testing all $\lambda/4$ plates are discussed in detail in the earlier polarization chapter.¹

Crystalline-Quartz Retardation Plates

Crystalline quartz is also frequently used for retardation plates, particularly those of the highest quality. It escapes the problem of zones with different orientations like those found in mica. The thickness of quartz required for a single quarter-wave retardation at the 6328-Å helium-neon laser line is about 0.017 mm, much too thin for convenient polishing. If the plate is to be used in the infrared, single-order quarter-wave plates are feasible (see “Single-Order Plates in the Infrared” section). Two types of quartz retardation plates are generally employed in the visible and ultraviolet regions: so-called “first-order” plates made of two pieces of material “First-Order Plates section” which are the best for critical applications, and multiple-order plates made of one thick piece of crystalline quartz (see “Multiple-Order Plates,” “Sensitivity to Temperature Changes,” and “Sensitivity to Angle of Incidence” sections). The multiple-order plates are generally not used for work of the highest accuracy since they are extremely sensitive to small temperature changes (see “Sensitivity to Temperature Changes” section) and to angle of incidence. Also, they have $\lambda/4$ retardation only at certain wavelengths; at other wavelengths the retardation may not even be close to $\lambda/4$.

When using any of the different types of retardation plates at a single wavelength, the methods for measuring the retardation and for distinguishing between fast and slow axes given in the earlier polarization chapter¹ can be used.

“First-Order” Plates A so-called “first-order” plate is made by cementing together two nearly equal thicknesses of quartz such that the fast axis of one is aligned parallel to the slow axis of the other (both axes lie in planes parallel to the polished faces). The plate is then polished until the difference in thickness between the two pieces equals the thickness of a single $\lambda/4$ plate. The retardation of this plate can be calculated from Eq. (9) by setting d equal to the *difference in thickness* between the two pieces. The “first-order” plate acts strictly like a single-order quarter-wave plate with respect to the variation of retardation with wavelength, temperature coefficient of retardation, and angle of incidence.

The change in phase retardation with temperature at 6328 Å, as calculated from equations given in the earlier polarization chapter,¹ is 0.0091°/°C, less than one-hundredth that of the 1.973-mm multiple-order plate discussed in “Sensitivity to Temperature Changes” section on p. 13.48. The change in retardation with angle of incidence* at this wavelength is also small: $(\Delta N)_{10^\circ} = 0.0016$, as compared with 0.18 for the thick plate (see “Sensitivity to Angle of Incidence” section on p. 13.48).

A “first-order” quartz $\lambda/4$ plate has several advantages over a mica $\lambda/4$ plate: (1) Crystalline quartz has a uniform structure, so that extinction can be obtained over the entire area of the plate at a given angular setting. (2) Since the total plate thickness is generally large, of the order of 1 mm or so, the coherence of the multiple, internally reflected beams is lost and there are no oscillations in the transmitted light or in the phase retardation. (3) Crystalline quartz is not pleochroic, except in the infrared, so that the intensity transmitted along the two axes is the same. (4) Crystalline quartz transmits farther into the ultraviolet than mica, so that “first-order” plates can be used from about 0.185–2.0 μm (see Table 8).

*Grechushnikov²⁸⁵ has an incorrect relation for the change in phase retardation with angle of incidence [his Eq. (2)]. He assumed that the retardations in the two halves of the plate add rather than subtract, yielding a retardation comparable to that of a thick quartz plate.

Single-Order Plates in the Infrared Although a crystalline-quartz retardation plate which is $\lambda/4$ in the visible is too thin to make from a single piece of material, the thickness required for such a plate is larger in the infrared. Jacobs and coworkers²⁶⁷ describe such a $\lambda/4$ plate for use at the 3.39- μm helium-neon laser line. They measured the birefringence of quartz at this wavelength and found it to be 0.0065 ± 0.0001 , so that the thickness required for the plate was 0.1304 mm. The actual plate was slightly thinner (0.1278 mm), so that it was tipped at an angle of 10° (rotating it about an axis parallel to the optic axis) to give it exactly $\lambda/4$ retardation (see “Sensitivity to Angle of Incidence” section on p. 13.48). Maillard²⁶² has also measured the birefringence of quartz at 3.39 μm and 3.51 μm and obtained values of 0.00659 and 0.00642, respectively (both ± 0.00002), in agreement with Jacobs’ value. These data lie on a smooth curve extrapolated from the values of Shields and Ellis.²⁵¹

A problem encountered when using crystalline quartz in the infrared is that, in general, the ordinary and extraordinary rays have different absorption coefficients; thus it may be impossible to construct a perfect wave plate regardless of the relative retardation between the rays. For an absorbing wave plate to have a retardation of exactly $\lambda/4$, the requirement

$$\left(\frac{n_o + 1}{n_e + 1}\right)^2 \exp\left[-\frac{(\alpha_e - \alpha_o)\lambda}{8(n_e - n_o)}\right] = 1 \quad (10)$$

must be met;²⁶⁷ α_e and α_o are the absorption coefficients for the extraordinary and ordinary rays, respectively. At wavelengths shorter than 3.39 μm , the birefringence is small enough for it to be possible to approximate the condition in Eq. (10) closely whenever $\alpha_e \approx \alpha_o$. Values of these quantities are given by Drummond.²⁸⁶ Gonatas et al.²⁷³ concluded that, in the far infrared and submillimeter wavelength region, the effect of different absorption coefficients in the crystalline quartz was small and could be corrected.

Another problem which occurs for crystalline quartz and also for sapphire²⁷³ in the infrared is that the Fresnel reflection coefficients are slightly different for the ordinary and extraordinary rays since the refractive indexes and absorption coefficients are in general different. One possible solution is to deposit isotropic thin films on the crystal surfaces.²⁷³ The refractive index of these films is chosen to balance the anisotropic absorption effect by making the Fresnel reflection coefficients appropriately anisotropic. On the other hand, if anisotropic Fresnel reflection proves to be undesirable, it can be greatly diminished by using an antireflection coating, as suggested by Gieselmann et al.²⁶⁷

If a single-order crystalline-quartz plate is to be used for a continuous range of wavelengths, both the phase retardation and the transmittance of the ordinary and extraordinary rays will oscillate as a function of wavelength because of multiple coherent reflections in the quartz. The separation between adjacent maxima in the phase retardation can be calculated from Eq. (144) in the earlier polarization chapter.¹ Using $\lambda = 3.3913 \mu\text{m}$, $n \approx 1.4881$, and $d = 127.8 \mu\text{m}$, $\Delta\lambda = 0.03024 \mu\text{m}$, an amount which should be well-resolved with most infrared instruments. Thus, if a wave plate is to be used over a range of wavelengths, it would be well to antireflect the surfaces to eliminate the phase oscillations.

Multiple-Order Plates Thick plates made from crystalline quartz are sometimes used to produce circularly polarized light at a single wavelength or a discrete series of wavelengths. The plate thickness is generally of the order of one or more millimeters so that the retardation is an integral number of wavelengths plus $\lambda/4$, hence the name multiple-order wave plate. This plate acts like a single $\lambda/4$ plate provided it is used only at certain specific wavelengths; at other wavelengths it may not even approximate the desired retardation. For example, a 1.973-mm-thick quartz plate was purchased which had an order of interference $N = 28.25$ at 6328 \AA . From Eq. (9) and Table 8, this plate would have $N = 30.52$ at 5890 \AA , and would thus be an almost perfect half-wave plate at this latter wavelength.

If a multiple-order plate is used to produce circularly polarized light at unspecified discrete wavelengths e.g., to measure circular or linear dichroism, it can be placed following a polarizer and oriented at 45° to the plane of vibration of the polarized beam. When the wavelengths are such that N calculated from Eq. (9) equals $1/4$, $3/4$, or in general $(2M - 1)/4$ (where M is a positive integer),

the emerging beam will be alternately right and left circularly polarized. The frequency interval $\Delta\nu$ between wavelengths at which circular polarization occurs is

$$\Delta\nu = \frac{1}{2d(n_e - n_o)} \quad (11)$$

where $\nu = 1/\lambda$. If the birefringence is independent of wavelength, the retardation plate will thus produce circularly polarized light at equal intervals on a frequency scale and can conveniently be used to measure circular dichroism, as described by Holzwarth.²⁸⁷

In order to approximately calibrate a multiple-order retardation plate at a series of wavelengths, it can be inserted between crossed polarizers and oriented at 45° to the polarizer axis. Transmission maxima will occur when the plate retardation is $\lambda/2$ or an odd multiple thereof; minima will occur when the retardation is a full wave or multiple thereof. If the axes of the two polarizers are parallel, maxima in the transmitted beam will occur when the plate retardation is a multiple of a full wavelength. The birefringence of the retardation plate can be determined by measuring the wavelengths at which maxima or minima occur if the plate thickness is known. Otherwise d can be measured with a micrometer, and an approximate value of $n_e - n_o$ can be obtained.

Palik²⁸⁸ made and tested a 2.070-mm-thick CdS plate for the 2- to 15- μm infrared region and also made thick retardation plates of SnSe, sapphire, and crystalline quartz to be used in various parts of the infrared. Holzwarth²⁸⁷ used a cultured-quartz retardation plate 0.8 mm thick to measure circular dichroism in the 1850- to 2500- \AA region of the ultraviolet; Jaffe et al.²⁸⁹ measured linear dichroism in the ultraviolet using a thick quartz plate and linear polarizer.

Sensitivity to Temperature Changes Small temperature changes can have a large effect on the retardation of a multiple-order plate. The method for calculating this effect was given earlier in the polarization chapter.¹ For the 1.973-mm-thick quartz plate mentioned in Multiple-Order Plates' section ($N = 28.25$ at 6328 \AA), the phase retardation will decrease 1.03° for each Celsius degree increase in temperature. If the temperature of the wave plate is not controlled extremely accurately, the large temperature coefficient of retardation can introduce sizable errors in precise ellipsometric measurements in which polarizer and analyzer settings can be made to $\pm 0.01^\circ$.

Sensitivity to Angle of Incidence The effect of angle of incidence (and hence field angle) on the retardation was calculated in the earlier polarization chapter.¹ It was shown there that the change in phase retardation with angle of incidence, $2\pi(\Delta N)_\theta$, is proportional to the total thickness of the plate (which is incorporated into N) and the square of the angle of incidence when the rotation is about an axis parallel to the optic axis. If the 1.973-mm-thick plate mentioned previously is rotated parallel to the optic axis through an angle of 10° at a wavelength of 6328 \AA , the total retardation changes from 28.25 to 28.43, so that the $\lambda/4$ plate is now nearly a $\lambda/2$ plate.

If the plate had been rotated about an axis *perpendicular* to the direction of the optic axis, in the limit when the angle of incidence is 90° , the beam would have been traveling along the optic axis; in this case the ordinary and extraordinary rays would be traveling with the same velocities, and there would have been *no retardation* of one relative to the other. For any intermediate angle of incidence the retardation would have been *less than* the value at normal incidence. The relation for the retardation as a function of angle of incidence is not simple, but the retardation will be approximately as angle-sensitive as it was in the other case. An advantage of rotation about either axis is that, with care, one can adjust the retardation of an inexact wave plate to a desired value. Rotation about an axis *parallel* to the optic axis will *increase* the retardation, while rotation about an axis *perpendicular* to the optic axis will *decrease* the retardation.

Achromatic Retardation Plates

Achromatic retardation plates are those for which the phase retardation is independent of wavelength. The name arose because when a plate of this type is placed between polarizers, it does not appear colored and hence is achromatic, as shown by Gaudefroy.²⁹⁰ In many applications, a truly

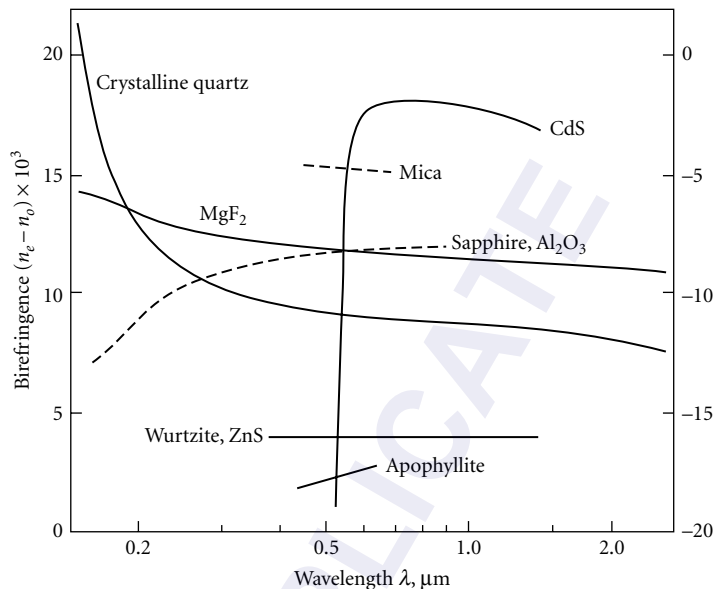


FIGURE 18 Birefringence of various optical materials as a function of wavelength. The scale at the left is for materials having a positive birefringence (*solid curves*), and the scale at the right is for materials with a negative birefringence (*dashed curves*).

achromatic retardation plate is not required. Since the wavelength of light changes by less than a factor of 2 across the visible region, a quarter- or half-wave mica plate often introduces only tolerable errors even in white light. The errors that do occur cancel out in many kinds of experiments.

Achromatic retardation plates can be made in various ways. The most achromatic are based on the principle of the Fresnel rhomb, in which the phase retardation occurs when light undergoes two or more total internal reflections (see next section “Rhombs as Achromatic $\lambda/4$ Retarders” and Ref. 1). A material with the appropriate variation of birefringence with wavelength can also be used. Such materials are uncommon, but plates of two or more different birefringent materials can be combined to produce a reasonably achromatic combination. Composite plates, consisting of two or more plates of the same material whose axes are oriented at the appropriate angles, can be used as achromatic circular polarizers or achromatic polarization rotators,¹ although they do not have all the properties of true $\lambda/4$ or $\lambda/2$ plates. One unusual achromatic half-wave plate is described in the earlier polarization chapter.¹

The simplest type of achromatic retardation plate could be made from a single material if its birefringence satisfied the requirement that $(n_e - n_o)/\lambda$ be independent of wavelength, i.e., that $n_e - n_o$ be directly proportional to λ . This result follows from Eq. (9) since $d(n_e - n_o)/\lambda$ must be independent of λ to make N independent of wavelength. (The plate thickness d is constant.) The birefringences of various materials are listed in Table 8 and plotted in Figs. 18 and 19. Only one material, the mineral apophyllite, has a birefringence which increases in the correct manner with increasing wavelength.^{264*} A curve of the phase retardation vs. wavelength for a quarter-wave apophyllite plate is shown as curve *D* in Fig. 20. Also included are curves for other so-called achromatic $\lambda/4$ plates as well as for simple $\lambda/4$ plates of quartz and mica. The phase retardation of apophyllite is not as constant with λ as that of the rhomb-type retarders, but it is considerably more constant than that of the other “achromatic” $\lambda/4$ plates. Since the birefringence of apophyllite is small, a $\lambda/4$ plate needs

*For materials having a negative birefringence the requirement is that $-(n_e - n_o)$ be proportional to λ .

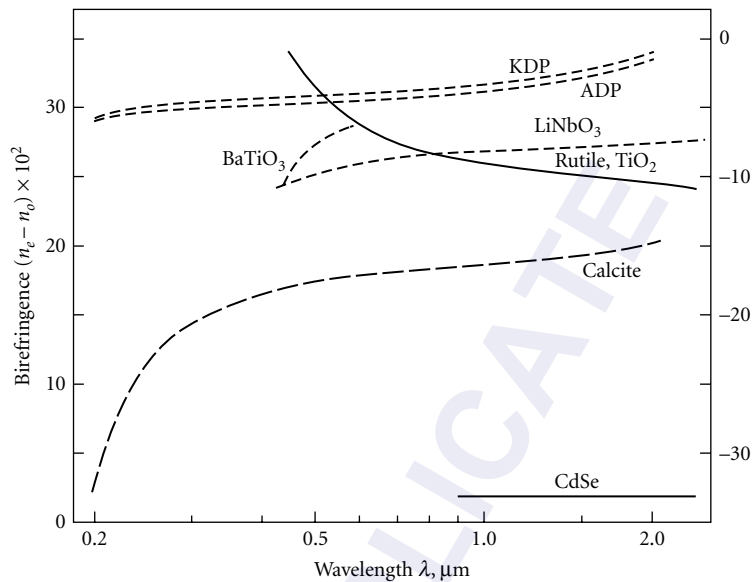


FIGURE 19 Birefringence of various optical materials which have larger birefringences than those shown in Fig. 18. The scale at the left is for materials having a positive birefringence (*solid curves*), and the scale at the right is for materials with a negative birefringence (*dashed curves*).

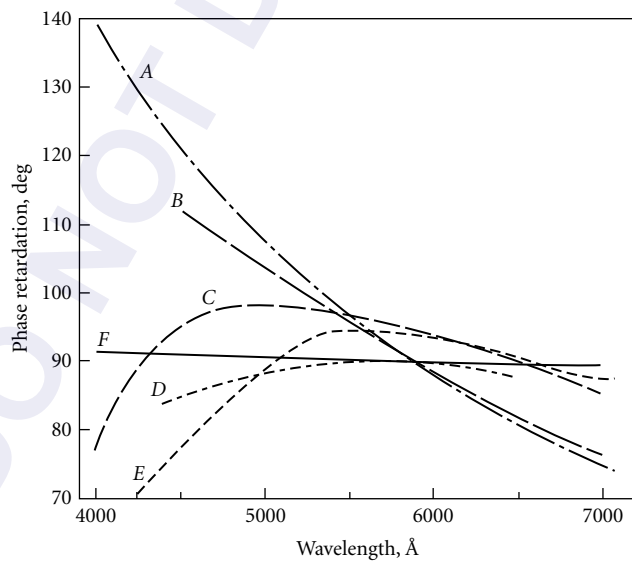


FIGURE 20 Curves of the phase retardation vs. wavelength for $\lambda/4$ plates: A, quartz; B, mica; C, stretched plastic film; D, apophyllite; and E, quartz-calcite achromatic combination. Curve F is for a Fresnel rhomb but is representative of all the rhomb-type devices. (After Bennett.²⁷⁶)

a thickness of about 56.8 μm , which is enough for it to be made as a single piece rather than as a “first-order” plate. Unfortunately optical-grade apophyllite is rare, the sample for which data are reported here having come from Sweden. There is some indication that the optical properties of other apophyllite samples may be different. Isotropic, positive, and negative-birefringent specimens have been reported by Deer et al.²⁹¹ According to them, the optical properties of apophyllite are often anomalous, some specimens being isotropic, uniaxial negative, or even biaxial with crossed dispersion of optic axial planes. Whether many samples have the favorable birefringence of the Swedish sample is uncertain.

Certain types of plastic film stretched during the manufacturing process have birefringences which are nearly proportional to wavelength and can serve as achromatic retardation plates if they have the proper thickness, as pointed out by West and Makas.²⁴⁶ Curve *C* in Fig. 20 is the retardation of a stretched cellulose nitrate film as measured by West and Makas.²⁴⁶ A combination of stretched cellulose acetate and cellulose nitrate sheets with their axes parallel will also make a reasonably achromatic $\lambda/4$ plate over the visible region. The advantages of using stretched plastic films for retardation plates are that they are cheap, readily available, have a retardation which is uniform over large areas, and can be used in strongly convergent light. However, each sheet must be individually selected since the birefringence is a strong function of the treatment during the manufacturing process and the sheets come in various thicknesses, with the result that their retardations are not necessarily $\lambda/4$ or $\lambda/2$. Also, Ennos²⁹² found that while the magnitude of the retardation was uniform over large areas of the sheets, the direction of the effective crystal axis varied from point to point by as much as 1.5° on the samples he was testing. Thus, film retarders appear to be excellent for many applications but are probably not suitable for measurements of the highest precision.

A reasonably achromatic retardation plate can be constructed from pairs of readily available birefringent materials such as crystalline quartz, sapphire, magnesium fluoride, calcite, or others whose birefringences are listed in Table 8. Assume that the plate is to be made of materials *a* and *b* having thicknesses d_a and d_b , respectively (to be calculated), and that it is to be achromatized at wavelengths λ_1 and λ_2 . From Eq. (9) we can obtain the relations

$$\begin{aligned} N\lambda_1 &= d_a \Delta n_{1a} + d_b \Delta n_{1b} \\ N\lambda_2 &= d_a \Delta n_{2a} + d_b \Delta n_{2b} \end{aligned} \quad (12)$$

where $N = 1/4$ for a $\lambda/4$ plate, $1/2$ for a $\lambda/2$ plate, etc., and the Δn 's are values of $n_e - n_o$ for the particular materials at the wavelengths specified; Δn will be positive for a positive uniaxial crystal and negative for a negative uniaxial crystal. (A positive uniaxial material can be used with its fast axis crossed with that of another positive uniaxial material; in this case the first material will have a negative Δn .) Equations (12) can be solved for d_a and d_b :

$$d_a = \frac{N(\lambda_1 \Delta n_{2b} - \lambda_2 \Delta n_{1b})}{\Delta n_{1a} \Delta n_{2b} - \Delta n_{1b} \Delta n_{2a}} \quad d_b = \frac{N(\lambda_2 \Delta n_{1a} - \lambda_1 \Delta n_{2a})}{\Delta n_{1a} \Delta n_{2b} - \Delta n_{1b} \Delta n_{2a}} \quad (13)$$

As an example of a compound plate, let us design a $\lambda/4$ plate of crystalline quartz and calcite and achromatize it at wavelengths $\lambda_1 = 0.508 \mu\text{m}$ and $\lambda_2 = 0.656 \mu\text{m}$. Quartz has a positive birefringence and calcite a negative birefringence (Table 8) so that Δn_{1a} and Δn_{2a} (for quartz) are positive Δn_{1b} and Δn_{2b} (for calcite) are negative. Equations (13) are satisfied for $d_{\text{qtz}} = 426.2 \mu\text{m}$ and $d_{\text{calc}} = 21.69 \mu\text{m}$; thus the phase retardation is exactly 90° at these two wavelengths. An equation of the form of those in Eqs. (12) is now used to calculate N for all wavelengths in the visible region using birefringence values listed in Table 8, and the results are plotted as curve *E* in Fig. 20. Although the achromatization for this quartz-calcite combination is not as good as can be obtained with a rhomb-type device or apophyllite, the phase retardation is within $\pm 5^\circ$ of 90° in the wavelength region 4900–7000 \AA and is thus much more constant than the retardation of a single mica or quartz $\lambda/4$ plate. Better two-plate combinations have been calculated by Beckers,²⁴⁹ the best being MgF_2 -ADP and MgF_2 -KDP, which have maximum deviations of ± 0.5 and ± 0.4 percent, respectively, compared with ± 7.2 percent

for a quartz-calcite combination over the same 4000- to 7000-Å wavelength region. The thicknesses of the materials which are required to produce $\lambda/4$ retardation* are $d_{\text{MgF}_2} = 113.79 \mu\text{m}$, $d_{\text{ADP}} = 26.38 \mu\text{m}$, and $d_{\text{MgF}_2} = 94.47 \mu\text{m}$, $d_{\text{KDP}} = 23.49 \mu\text{m}$. Since the ADP and KDP must be so thin, these components could be made in two pieces as “first-order” plates.

Other two-component compound plates have been proposed by Chandrasekharan and Damany,²⁶¹ Gaudefroy,²⁹⁰ Ioffe and Smirnova,²⁹³ and Mitchell.²⁹⁴ The paper by Ioffe and Smirnova describes a quartz-calcite combination similar to the one illustrated earlier, but it contains various numerical errors which partially invalidate the results.

If better achromatization is desired and one does not wish to use a rhomb-type $\lambda/4$ device, three materials can be used which satisfy the relations

$$\begin{aligned} N\lambda_1 &= d_a \Delta n_{1a} + d_b \Delta n_{1b} + d_c \Delta n_{1c} \\ N\lambda_2 &= d_a \Delta n_{2a} + d_b \Delta n_{2b} + d_c \Delta n_{2c} \\ N\lambda_3 &= d_a \Delta n_{3a} + d_b \Delta n_{3b} + d_c \Delta n_{3c} \end{aligned} \quad (14)$$

where the Δn 's are birefringences of the various materials at wavelengths λ_1 , λ_2 , and λ_3 .

Instead of using only three wavelengths, Beckers²⁴⁹ suggested that the thicknesses can be optimized such that the maximum deviations from achromatization are minimized over the entire wavelength interval desired. In this way, he obtained a three-component combination of quartz, calcite, and MgF_2 which has a retardation of a full wavelength and a maximum deviation of only ± 0.2 percent over the 4000- to 7000-Å wavelength region. The maximum deviation of slightly different thicknesses of these same three materials rises to ± 2.6 percent if the wavelength interval is extended to 3000–11,000 Å. Chandrasekharan and Damany²⁶¹ have designed a three-component $\lambda/4$ plate from quartz, MgF_2 , and sapphire for use in the vacuum ultraviolet. Title²⁹⁵ has designed achromatic combinations of three-element, four-element, nine-element, and ten-element waveplates using Jones matrix techniques. The nine-element combination is achromatic to within 1° from 3500 to 10,000 Å. He constructed and tested several waveplate combinations, and they performed as designed.

Rhombs as Achromatic $\lambda/4$ Retarders

The simplest stable, highly achromatic $\lambda/4$ retarder with a reasonable acceptance angle and convenient size appears to be a rhomb-type retarder. Several types are available; the choice of which one to use for a specific application depends on (1) the geometry of the optical system (can a deviated or displaced beam be tolerated?), (2) wavelength range, (3) degree of collimation of the beam, (4) beam diameter (determining the aperture of the retarder), (5) space available, and (6) accuracy required. Table 9 summarizes the properties of the various achromatic rhombs. This subject has been covered in detail by Bennett²⁷⁶ and is condensed from that reference in the earlier polarization chapter.¹ Anderson²⁹⁶ has compared the retardation of a CdS $\lambda/4$ plate and a Fresnel rhomb in the 10- μm CO_2 laser emission region. Wizinowich²⁹⁷ used a Fresnel rhomb along with some additional optics to change an unpolarized light beam from a faint star object into linearly polarized light to improve the throughput of a grating spectrograph and make it independent of the input polarization.

Composite Retardation Plates

A composite retardation plate is made up of two or more elements of the same material combined so that their optic axes are at appropriate angles to each other. Some of the composite plates have nearly all the properties of a true retardation plate, whereas others do not. In the earlier polarization chapter,¹

*Beckers' tables II to V give the thickness of materials required to produce one full-wave retardation. To obtain values of thicknesses for $\lambda/4$ retardation, for example, multiply all d values in the table by 0.25. The percent deviations should remain unchanged.

TABLE 9 Properties of Achromatic Rhombs²⁷⁵

Name	Light Path	Internal Angle of Incidence, deg	Material	Refractive Index		Variation of Phase Retardation			
				<i>n</i>	Wave-Length, Å	With Wavelength		With Angle of Incidence	
						Var., deg	Wave-Length, Å	Var., deg	Angle, deg
Fresnel rhomb	Translated	54.7	Crown glass	1.511	5893	2.5	3650–7682	9.1	-7 to +7
Coated Fr. rhomb	Translated	51.5	Crown glass	1.5217	5461	0.4	3341–5461	2.5	-4 to +6
	Translated	54.0	Fused quartz	1.4880	3000	0.7	2148–3341	<0.5	-1.5 to +1.5
Mooney rhomb	Deviated	60.0	Flint glass	1.650	5893	1.9	4047–6708	0.7	-7 to +7
AD-1	Undeviated	74.3	Fused quartz	1.4702	4000	2.0	3000–8000	0.7	-7 to +7
AD-2	Undeviated	73.2, 56.4	Fused quartz	1.4702	4000	2.9	3000–8000	13.2	-3 to +3
Coated AD-2	Undeviated	72.2	Fused quartz	1.4601	5461	0.3	2140–5461	6.0	-1.5 to +1.5
AD	Undeviated	53.5	Crown glass	1.511	5893	1.6	3650–7682	9.4	-7 to +7

composite plates were described which produced circularly polarized light at a given wavelength, those which acted as achromatic circular polarizers, and those which acted as achromatic polarization rotators or pseudo $\lambda/2$ plates. The effect of combining several birefringent plates with their axes at arbitrary angles to each other can be easily understood using the Poincaré sphere. A general treatment of this subject has been given by Ramachandran and Ramaseshan.²⁹⁸

13.11 VARIABLE RETARDATION PLATES AND COMPENSATORS

Variable retardation plates can be used to modulate or vary the phase of a beam of plane-polarized light, to measure birefringence in mineral specimens, flow birefringence, or stress in transparent materials, or to analyze a beam of elliptically polarized light such as might be produced by transmission through a birefringent material or by reflection from a metal or film-covered surface. The term compensator is frequently applied to a variable retardation plate since it can be used to compensate for the phase retardation produced by a specimen. Common types of variable compensators include the Babinet and Soleil compensators, in which the total thickness of birefringent material in the light path is changed, the Sénarmont compensator,¹ which consists of a fixed quarter-wave plate and rotatable analyzer to compensate for varying amounts of ellipticity in a light beam, and tilting-plate compensators,¹ with which the total thickness of birefringent material in the light beam is changed by changing the angle of incidence. Electro-optic and piezo-optic modulators can also be used as variable retardation plates since their birefringence can be changed by varying the electric field or pressure. However, they are generally used for modulating the amplitude, phase, frequency, or direction of a light beam, in particular a laser beam, at frequencies too high for mechanical shutters or moving mirrors to follow. Information on electro-optic materials and devices is contained in the chapter on electro-optic modulators by Georgeanne M. Purvinis and Theresa A. Maldonado (Chap. 7 in Vol. V) and in the earlier polarization chapter.¹

Babinet Compensator

There are many devices which compensate for differences in phase retardation by having a variable thickness of a birefringent material (such as crystalline quartz) in the light beam, as discussed by Johansen,²⁹⁹ and Jerrard.³⁰⁰ One such device, described by Hunt,³⁰¹ can compensate for a residual

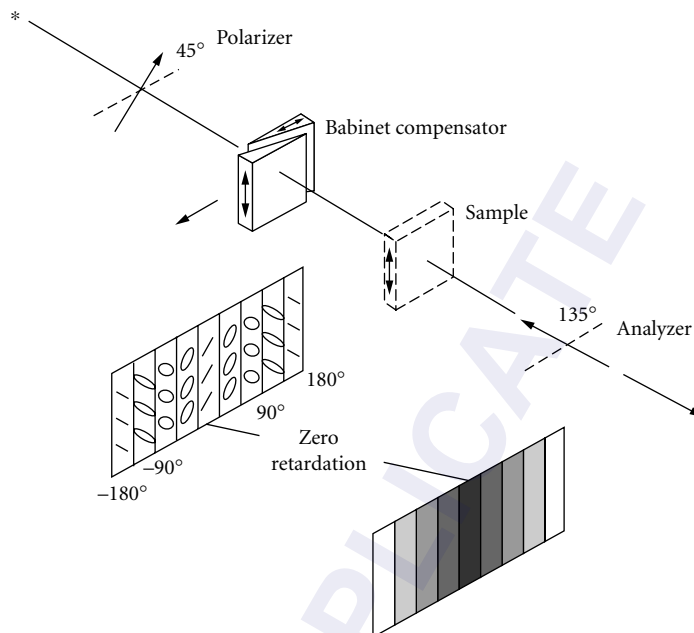


FIGURE 21 Arrangement of a Babinet compensator, polarizer, and analyzer for measuring the retardation of a sample. The appearance of the field after the light has passed through the compensator is shown to the left of the sample position. Retardations are indicated for alternate regions. After the beam passes through the analyzer, the field is crossed by a series of dark bands, one of which is shown to the left of the analyzer.

wedge angle between the entrance and exit faces of birefringent optical components such as optical modulators and waveplates.

The most common variable retardation plates are the Babinet compensator and the Soleil compensator. The Babinet compensator was proposed by Babinet in 1837 and later modified by Jamin; references to the voluminous early literature are given by Partington.³⁰² Ellerbroek and Groosmuller³⁰³ have a good description of the theory of operation (in German), and Jerrard³⁰⁴⁻³⁰⁶ and Archard³⁰⁷ describe various optical and mechanical defects of Babinet compensators.

The Babinet compensator, shown schematically in Fig. 21, consists of two crystalline-quartz wedges, each with its optic axis in the plane of the face but with the two optic axes exactly 90° apart. One wedge is stationary, and the other is movable by means of a micrometer screw in the direction indicated by the arrow, so that the total amount of quartz through which the light passes can be varied uniformly. In the first wedge, the extraordinary ray vibrates in a horizontal plane and is retarded relative to the ordinary ray (crystalline quartz has a positive birefringence; see Table 8). When the rays enter the second wedge, the ray vibrating in the horizontal plane becomes the ordinary ray and is advanced relative to the ray vibrating in the vertical plane. Thus, the total retardation is proportional to the difference in thickness between the two wedges:

$$N\lambda = (d_1 - d_2)(n_e - n_o) \quad (15)$$

where N = retardation in integral and fractional parts of a wavelength
 d_1, d_2 = thickness of the first and second wedges where light passes through
 n_o, n_e = ordinary and extraordinary refractive indexes for crystalline quartz

If light polarized at an angle of 45° to one of the axes of the compensator passes through it, the field will appear as shown in Fig. 21; the wedges have been set so there is zero retardation at the center of the field. (If the angle α of the incident plane-polarized beam were different from 45° , the beam retarded or advanced by 180° in phase angle would make an angle of 2α instead of 90° with the original beam.) When an analyzer whose axis is crossed with that of the polarizer is used to observe the beam passing through the compensator, a series of light and dark bands is observed in monochromatic light. In white light only one band, that for which the retardation is zero, remains black. All the other bands are colored. These are the bands for which the retardation is multiples of 2π (or, expressed in terms of path differences, integral numbers of wavelengths). On one side of the central black band one ray is advanced in phase relative to the other ray; on the other side it is retarded. If one wedge is moved, the whole fringe system translates across the field of view. The reference line is scribed on the stationary wedge so that it remains in the center of the field. Information on calibrating and using a Babinet compensator is given in the earlier polarization chapter.¹

Soleil Compensator

The Soleil compensator (see Wood³⁰⁸ and Ditchburn³⁰⁹), sometimes called a Babinet-Soleil compensator, is shown in Fig. 22. It is similar to the Babinet compensator in the way it is used, but instead of having a field crossed with alternating light and dark bands in monochromatic light, the

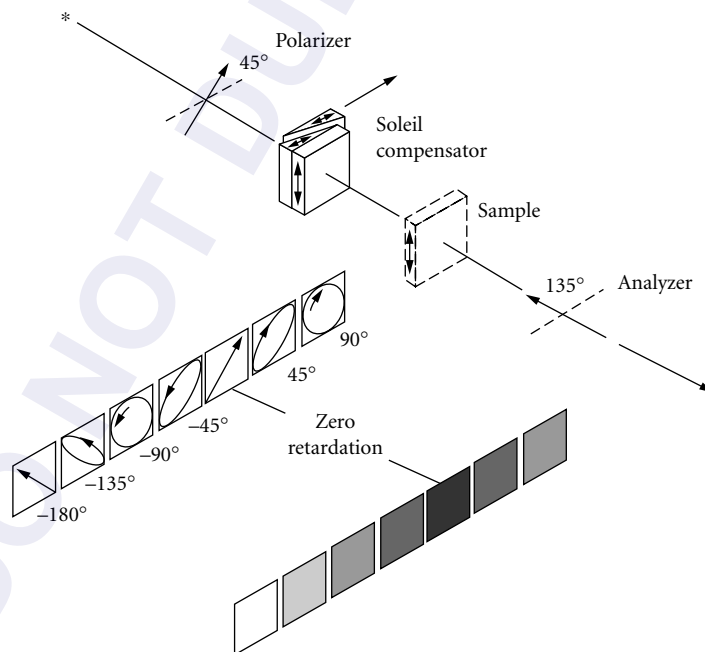


FIGURE 22 Arrangement of a Soleil compensator, polarizer, and analyzer for measuring the retardation of a sample. The appearance of the field after the light has passed through the compensator is shown to the left of the sample position. After the beam passes through the analyzer, the field appears as one of the shades of gray shown to the left of the analyzer.

field has a uniform tint if the compensator is constructed correctly. This is because the ratio of the thicknesses of the two quartz blocks (one composed of a fixed and a movable wedge) is the same over the entire field. The Soleil compensator will produce light of varying ellipticity depending on the position of the movable wedge. Calibration of the Soleil compensator is similar to that of the Babinet compensator.¹ The zero-retardation position is found in the same manner except that now the entire field is dark. The compensator is used in the same way as a Babinet compensator with the uniformly dark field (in white light) of the Soleil corresponding to the black zero-retardation band in the Babinet.

The major advantage of the Soleil compensator is that a photoelectric detector can be used to make the settings. The compensator is offset a small amount on each side of the null position so that equal-intensity readings are obtained. The average of the two drum positions gives the null position. Photoelectric setting can be much more precise than visual setting, but this will not necessarily imply increased accuracy unless the compensator is properly constructed. Since Soleil compensators are composed of three pieces of crystalline quartz, all of which must be very accurately made, they are subject to more optical and mechanical defects than Babinet compensators. Jerrard³¹⁰⁻³¹² has described many of these defects in detail. Ives and Briggs³¹³ found random departures of about $\pm 1.5^\circ$ from their straight-line calibration curve of micrometer reading for extinction vs. wedge position. This variation was considerably larger than the setting error with a half-shade plate and was attributed to variations in thickness of the order of $\pm \lambda/4$ along the quartz wedges.

Soleil compensators have been used for measurements of retardation in the infrared. They have been made of crystalline quartz, cadmium sulfide, and magnesium fluoride (see the work of Palik^{257,314} and Palik and Hennis²⁵²). A by-product of this work was the measurement of the birefringence of these materials in the infrared.

Two other uniform-field compensators have been proposed. Jerrard,³¹¹ following a suggestion by Soleil, has taken the Babinet wedges and reversed one of them so the light passes through the thicker portions of each wedge. This reversed Babinet compensator is less subject to mechanical imperfections than the Soleil compensator but does produce a small deviation of the main beam. Hariharan and Sen³¹⁵ suggest double-passing a Babinet compensator (with a reflection between the two passes) to obtain a uniform field.

13.12 HALF-SHADE DEVICES

It is sometimes necessary to measure accurately the azimuth of a beam of plane-polarized light, i.e., the angle the plane of vibration makes with a reference coordinate system. This can be done most easily by using a polarizer as an analyzer and rotating it to the position where the field appears the darkest. The analyzer azimuth is then exactly 90° from the azimuth of the plane-polarized beam. A more sensitive method is to use a photoelectric detector and offset on either side of the extinction position at angles where the intensities are equal. The average of these two angles is generally more accurate than the value measured directly, but care must be taken to keep the angles small so that asymmetries will not become important.

Before the advent of sensitive photoelectric detectors, the most accurate method of setting on a minimum was to use a half-shade device as the analyzer or in conjunction with the analyzer. The device generally consisted of two polarizers having their axes inclined at an angle α to each other (angle fixed in some types and variable in others). As the device was rotated, one part of the field became darker while the other part became lighter. At the match position, both parts of the field appeared equally bright. The Jellett-Cornu prism, Lippich, and Laurent half shades, Nakamura biplate, and Savart plate are examples of half-shade devices.¹

Ellipticity half-shade devices are useful for detecting very small amounts of ellipticity in a nominally plane-polarized beam and hence can indicate when a compensator has completely converted elliptically polarized light into plane-polarized light. Two of these devices are the Bravais biplate

and the Brace half-shade plate. Half-shade devices for both plane and elliptically polarized light are described in detail in the earlier polarization chapter.¹

13.13 MINIATURE POLARIZATION DEVICES

Polarization Devices for Optical Fibers

Single-mode optical fiber-type polarizers are important devices for optical fiber communication and fiber sensor systems. These polarizers have been made by a variety of techniques. Polarizers have been made by bending³¹⁶ or by tapering³¹⁷ a birefringent fiber to induce differential attenuation in the orthogonal modes. In most cases a fiber was polished laterally and some device was placed in contact with the exposed guiding region of the fiber to couple out the unwanted polarization. Bergh et al.³¹⁸ used a birefringent crystal as the outcoupling device and obtained a high extinction ratio polarizer. Optical fiber polarizers made with a metal film coated onto the polished area to eliminate the unwanted polarization state seem to be preferred because they are stable and rugged. The original version by Eickhoff³¹⁹ used the thin cladding remaining after polishing as the buffer layer, but it had an insufficient extinction ratio. Other designs using metal coatings were suggested by Gruchmann et al.,³²⁰ and Hosaka et al.³²¹ Feth and Chang³²² used a fiber polished into its core to which a superstrate coated with a very thin metal layer was attached by an index-matching oil. Yu and Wu³²³ gave a theoretical analysis of metal-clad single-mode fiber-type polarizers. Dyott et al.³²⁴ made a metal-fiber polarizer from an etched D-shaped fiber coated with indium.

In the above approaches, either expensive components are used or the structure of the polarizer is complicated and fragile. Lee and Chen³²⁵ suggested a new way of fabricating high-quality metal-clad polarizers by polishing a fiber $\sim 0.4 \mu\text{m}$ into its core and then overcoating it with a 265-nm MgF_2 film as the buffer layer followed by a 100-nm Al film. Polarizers fabricated in this way had an average extinction ratio of 28 dB with a 2-dB insertion loss at a 0.63- μm wavelength or a 34-dB extinction ratio with a 3-dB insertion loss at 0.82 μm .³²⁵

Other devices for optical fibers have also been designed. Ulrich and Johnson³²⁶ made a single-mode fiber-optical polarization rotator by mechanically twisting successive half-wave fiber sections in alternating directions; Hosaka et al.'s fiber circular polarizer³²⁷ was composed of a metal-coated fiber polarizer and a $\lambda/4$ platelet fabricated on a birefringent fiber; polished-type couplers acting as polarizing beam splitters were made by Snyder and Stevenson.³²⁸ The patent literature contains references to other polarization devices for optical fibers.

Polarization Devices for Integrated Circuits

Small and highly efficient polarization devices are also needed for integrated circuits. Some such devices have been proposed and fabricated. Uehara et al.³²⁹ made an optical waveguiding polarizer for optical fiber transmission out of a plate of calcite attached to borosilicate glass into which a three-dimensional high-index region had been formed by ion migration to act as the waveguide. Mahlein³³⁰ deposited a multilayer dielectric film onto a glass superstrate which was then contacted to a planar waveguide to couple out the TM polarization. This paper contains a good description of the polarizer design as well as extensive references. Suchoski et al.³³¹ fabricated low-loss, high-extinction polarizers in LiNbO_3 by proton exchange. Noé et al.³³² achieved automatic endless polarization control with integrated optical Ti:LiNbO_3 polarization transformers. This was a better method of matching polarization states between two superposed waves than techniques that had been used previously. Finally Baba et al.³³³ proposed making a polarizer for integrated circuits out of periodic metal-dielectric laminated layers (Lamipol structures). Their experiments with Al-SiO₂ structures were encouraging. Patents have been filed for other polarization devices for integrated circuits.

13.14 REFERENCES*

1. H. E. Bennett and J. M. Bennett, "Polarization," in W. G. Driscoll and W. Vaughan, (eds.), *Handbook of Optics*, 1st ed., McGraw-Hill, New York, 1978, pp. 10-1–10-164.
2. H. J. Nickl and H. K. Henisch, *J. Electrochem. Soc.* **116**:1258–1260, 1969.
3. R. N. Smartt, *J. Sci. Instrum.* **38**:165, 1961.
4. E. E. Wahlstrom, *Optical Crystallography*, 4th ed., Wiley, New York, 1969, pp. 236–267.
5. E. Uzan, H. Damany, and V. Chandrasekharan, *Opt. Commun.* **1**:221–222, 1969.
6. S. S. Ballard, J. S. Browder, and J. F. Ebersole, in D. E. Gray (ed.), *American Institute of Physics Handbook*, 3d ed., McGraw-Hill, New York, 1972, pp. 6–20.
7. J. Schnellman, V. Chandrasekharan, H. Damany, and J. Romand, *C. R. Acad. Sci.* **260**:117–120, 1965.
8. Y. Bouriau and J. Lenoble, *Rev. Opt.* **36**:531–543, 1957.
9. S. S. Ballard, J. S. Browder, and J. F. Ebersole, in D. E. Gray, (ed.), *American Institute of Physics Handbook*, 3d ed., McGraw-Hill, New York, 1972, p. 6–65.
10. H. Damany, Laboratoire des Hautes Pressions, Bellevue, France, private communication, 1970.
11. E. Uzan, H. Damany, and V. Chandrasekharan, *Opt. Commun.* **2**:273–275, 1970.
12. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon Press, New York, 1980, p. 680.
13. L. C. Martin, *Technical Optics*, vol. 1, Pitman, London, 1948, pp. 196–198.
14. R. W. Ditchburn, *Light*, 2d ed., Interscience, New York, 1963, pp. 595–616.
15. A. Schuster, *Theory of Optics*, 2nd ed., Arnold, London, 1920, pp. 168–187.
16. R. Müller, *Optik* **20**:510–511, 1963.
17. C. Dévé, *Optical Workshop Principles*, T. L. Tippell (trans.), Hilger & Watts, London, 1954, p. 295.
18. P. Jacquinot, *J. Opt. Soc. Am.* **44**:761–765, 1954.
19. L. Mertz, *Transformations in Optics*, Wiley, New York, 1965, pp. 15–16.
20. P. Glan, *Carl's Repert.* **16**:570, 1880.
21. S. P. Thompson, *Phil. Mag.* **12**(5):349, 1881.
22. S. P. Thompson, *Proc. Opt. Conv.* 216–235, 1905.
23. R. T. Glazebrook, *Phil. Mag.* **15**(5):352, 1883.
24. C. E. Moeller and D. R. Grieser, *Appl. Opt.* **8**:206–207, 1969.
25. R. J. King and S. P. Talim, *J. Phys.* (GB) **4**(E):93–96, 1971.
26. J. W. Ellis and J. Bath, *J. Chem. Phys.* **6**:221–222, 1938.
27. F. Lippich, *Wien Akad. Sitzungsber.* **91**(III):1059, 1885.
28. A. C. Hardy and F. H. Perrin, *The Principles of Optics*, McGraw-Hill, New York, 1932, p. 611.
29. J. F. Archard and A. M. Taylor, *J. Sci. Instrum.* **25**:407–409, 1948.
30. Karl Lambrecht Corp., Bull. P-73, Chicago, 1973.
31. J. Swartz, D. K. Wilson, and R. J. Kapash, *High Efficiency Laser Polarizers*, Electro-Opt. 1971 West Conf., Anaheim Calif., May, 1971.
32. A. Lambrecht, Karl Lambrecht Corp., Chicago, Ill., private communication, 1969.
33. A. V. Shustov, *Sov. J. Opt. Technol.* **34**:177–181, 1967.
34. F. Twyman, *Prism and Lens Making*, 2d ed., Hilger & Watts, London, 1952, pp. 244, 599.
35. J. F. Archard, *J. Sci. Instrum.* **26**:188–192, 1949.
36. H. E. Bennett and J. M. Bennett, "Precision Measurements in Thin Film Optics," in G. Hass and R. E. Thun (eds.), *Physics of Thin Films*, vol. 4, Academic Press, New York, 1967, pp. 69–78.
37. D. L. Decker, J. L. Stanford, and H. E. Bennett, *J. Opt. Soc. Am.* **60**:1557A, 1970.

*In all references to the Russian literature, volume and pages cited are for the English translation.

38. R. L. Rowell, A. B. Levit, and G. M. Aval, *Appl. Opt.* **8**:1734, 1969.
39. D. E. Aspnes, *Appl. Opt.* **9**:1708–1709, 1970.
40. R. L. Rowell, *Appl. Opt.* **9**:1709, 1970.
41. W. Nicol, *Edinb. New Phil. J.* **6**:83, 1828–1829, as quoted in A. Johannsen, *Manual of Petrographic Methods*, 2d ed., Hafner, New York, 1968, p. 158; (originally published in 1918).
42. S. P. Thompson, *Phil. Mag.* **21**(5):476, 1886.
43. B. Halle, *Dtsch. Mech. Z.* **1**:6–7, Jan. 1, 1908.
44. B. Halle, *Dtsch. Mech. Z.* **2**:16–19, Jan. 15, 1908.
45. A. B. Dale, in R. Glazebrook (ed.), *A Dictionary of Applied Physics*, vol. 4, Macmillan, London, 1923, pp. 496–497.
46. Hartnack and Prazmowski, *Ann. Chim. Phys.* **7**(4):181, 1866.
47. L. Foucault, *C. R. Acad. Sci.* **45**:238, 1857.
48. C. D. West and R. C. Jones, *J. Opt. Soc. Am.* **41**:976–982, 1951.
49. D. L. Steinmetz, W. G. Phillips, M. Wirick, and F. F. Forbes, *Appl. Opt.* **6**:1001–1004, 1967.
50. G. C. Morris and A. S. Abramson, *Appl. Opt.* **8**:1249–1250, 1969.
51. E. O. Ammann and G. A. Massey, *J. Opt. Soc. Am.* **58**:1427–1433, 1968.
52. A. C. Hardy, *J. Opt. Soc. Am.* **25**:305–311, 1935.
53. C. Bouhet and R. LaFont, *Rev. Opt.* **28**:490–493, 1949.
54. A. Cotton, *C. R. Acad. Sci.* **193**:268–271, 1931.
55. W. C. Johnson, Jr., *Rev. Sci. Instrum.* **35**:1375–1376, 1964.
56. L. V. Foster, *J. Opt. Soc. Am.* **28**:124–126, 127–129, 1938.
57. K. Feussner, *Z. Instrumentenk.* **4**:41, 1884.
58. A. Johannsen, *Manual of Petrographic Methods*, 2d ed., Hafner, New York, 1968, pp. 169, 283–285; (originally published in 1918).
59. W. L. Hyde, New York Univ., Bronx N.Y., private communication 1970.
60. E. Bertrand, *C. R. Acad. Sci.* **49**:538, 1884.
61. L. Wulff, *Sitz. Preuss. Akad. Wiss.* **135**:879, 1896.
62. P. Stöber, *Z. Krist.* **61**:299, 1924.
63. P. Stöber, *Neues Jahrb. Mineral.* **A57**:139, 1928.
64. P. Stöber, *Chem. Erde* **6**:357, 453, 1930.
65. E. Tzschknovitzer, *J. Phys. Chem. (USSR)* **5**:1452, 1934.
66. C. D. West, *J. Opt. Soc. Am.*, **35**:26–31, 1945.
67. M. Huot de Longchamp, *Rev. Opt.* **26**:94–98, 1947.
68. T. Yamaguti, *J. Phys. Soc. Japan* **10**:219–221, 1955.
69. T. Yamaguti, I. Makino, S. Shinoda, and I. Kuroha, *J. Phys. Soc. Japan* **14**:199–201, 1959.
70. F. J. Dumont and R. N. Smartt, *J. Opt. Soc. Am.* **59**:1541A, 1969.
71. F. J. Dumont, *J. Opt. Soc. Am.* **60**:719A, 1970.
72. L. G. DeShazer, Dept. of Electrical Engineering Univ. Southern California, Los Angeles, private communication, 1971.
73. V. Chandrasekharan and H. Damany, *Appl. Opt.* **8**:675–675, 1969.
74. V. Chandrasekharan and H. Damany, *Appl. Opt.* **10**:681–682, 1971.
75. E. Landais, *Bull. Soc. Fr. Mineral. Cristallogr.* **91**:350–354, 1968.
76. T. J. Bridges and J. W. Kluver, *Appl. Opt.* **4**:1121–1125, 1965.
77. J. J. Loferski, *Phys. Rev.* **87**:905–906, 1952.
78. R. Newman and R. S. Halford, *Rev. Sci. Instrum.* **19**:270–271, 1948.
79. W. L. Hyde, *J. Opt. Soc. Am.* **38**:663A, 1948.

80. E. H. Land and C. D. West, in J. Alexander (ed.), *Colloid Chemistry*, vol. 6, Reinhold, New York, 1946, pp. 160–190.
81. E. H. Land, *J. Opt. Soc. Am.* **41**:957–963, 1951.
82. W. A. Shurcliff, *Polarized Light*, Harvard University Press, Cambridge, Mass., 1962, pp. 43–64.
83. S. D. Stookey and R. J. Araujo, *Appl. Opt.* **7**:777–779, 1968.
84. T. Yamaguti, *J. Opt. Soc. Am.* **45**:891–892, 1955.
85. P. Baumeister and J. Evans, in D. E. Gray (ed.), *American Institute of Physics Handbook*, 3rd ed., McGraw-Hill, New York, 1972, pp. 6-171–6-172.
86. R. P. Blake, A. S. Makas, and C. D. West, *J. Opt. Soc. Am.* **39**:1054A, 1949.
87. A. S. Makas, *J. Opt. Soc. Am.* **52**:43–44, 1962.
88. L. Baxter, A. S. Makas, and W. A. Shurcliff, *J. Opt. Soc. Am.* **46**:229, 1956.
89. R. C. Jones, Polaroid Corporation, Cambridge, Mass., private communication, 1970.
90. M. Haase, *Zeiss-Mitt.* **2**:173, 1961.
91. M. N. McDermott and R. Novick, *J. Opt. Soc. Am.* **51**:1008–1010, 1961.
92. G. R. Bird and M. Parrish, Jr., *J. Opt. Soc. Am.* **50**:886–891, 1960.
93. G. B. Trapani, *Proc. Soc. Photo-Opt. Instrum. Eng.* **88**:105–113, 1976.
94. W. J. Gunning and J. Foschaar, *Appl. Opt.* **22**:3229–3231, 1983.
95. S. J. Baum, *Proc. Soc. Photo-Opt. Instrum. Eng.* **88**:50–56, 1976.
96. J. F. Dreyer, *J. Opt. Soc. Am.* **37**:983A, 1947.
97. S. Anderson, *J. Opt. Soc. Am.* **39**:49–56, 1949.
98. Polacoat Bull. P-108 and P-112, Polacoat, Inc., Cincinnati, Ohio, 1967.
99. D. S. Kyser, Michelson Laboratory, Naval Weapons Center, China Lake, Calif., private communication, 1970.
100. G. Rupprecht, D. M. Ginsberg, and J. D. Leslie, *J. Opt. Soc. Am.* **52**:665–669, 1962.
101. Perkin-Elmer Corp., Instrument Division, Norwalk, Conn., “Wire Grid Polarizer Accessory,” Sheet no. D-454, 1966.
102. J. B. Young, H. A. Graham, and E. W. Peterson, *Appl. Opt.* **4**:1023–1026, 1965.
103. K. F. Renk and L. Genzel, *Appl. Opt.* **1**:643–648, 1962.
104. J. P. Casey and E. A. Lewis, *J. Opt. Soc. Am.* **42**:971–977, 1952.
105. E. A. Lewis and J. P. Casey, *J. Appl. Phys.* **23**:605–608, 1952.
106. M. Mohebi, J. Q. Liang, and M. J. Soileau, *Appl. Opt.* **28**:3681–3683, 1989.
107. R. W. Stobie and J. J. Dignam, *Appl. Opt.* **12**:1390–1391, 1973.
108. C. H. Burton, *Appl. Opt.* **18**:420–422, 1979.
109. J. P. Auton, *Appl. Opt.* **6**:1023–1027, 1967.
110. M. Hass and M. O’Hara, *Appl. Opt.* **4**:1027–1031, 1965.
111. A. R. Hilton and C. E. Jones, *J. Electrochem. Soc.* **113**:472–478, 1966.
112. D. G. Vickers, E. I. Robson, and J. E. Beckman, *Appl. Opt.* **10**:682–684, 1971.
113. P. K. Cheo and C. D. Bass, *Appl. Phys. Lett.* **18**:565–567, 1971.
114. J. P. Auton and M. C. Hutley, *Infrared Phys.*, **12**:95–100, 1972.
115. A. E. Costley, K. H. Hursey, G. F. Neill, and J. M. Ward, *J. Opt. Soc. Am.* **67**:979–981, 1977.
116. J. A. Beunen, A. E. Costley, G. F. Neill, C. L. Mok, T. J. Parker, and G. Tait, *J. Opt. Soc. Am.* **71**:184–188, 1981.
117. T. A. Leonard, *Soc. Photo-Opt. Instrum. Eng.* **288**:129–135, 1981.
118. G. J. Sonek, D. K. Wanger, and J. M. Ballantyne, *Appl. Opt.* **22**:1270–1272, 1983.
119. W. L. Eichhorn and T. J. Magner, *Opt. Eng.* **25**:541–544, 1986.
120. G. Novak, R. J. Pernic, and J. L. Sundwall, *Appl. Opt.* **28**:3425–3427, 1989.
121. S. Roberts and D. D. Coon, *J. Opt. Soc. Am.* **52**:1023–1029, 1962.
122. Y. S. Hwang and H. K. Park, *Appl. Opt.* **28**:4999–5001, 1989.

123. H. Weiss and M. Wilhelm, *Z. Phys.* **176**:399–408, 1963.
124. B. Paul, H. Weiss, and M. Wilhelm, *Solid State Electron.* (GB) **7**:835–842, 1964.
125. A. Mueller and M. Wilhelm, *J. Phys. Chem. Solids* **26**:2029, 1965.
126. N. M. Davis, A. R. Clawson and H. H. Wieder, *Appl. Phys. Lett.* **15**:213–215, 1969.
127. M. Saito and M. Miyagi, *Appl. Opt.* **28**:3529–3533, 1989.
128. A. Hidalgo, J. Pastor and J. M. Serratos, *J. Opt. Soc. Am.* **52**:1081–1082, 1962.
129. T. G. R. Rawlins, *J. Opt. Soc. Am.* **54**:423–424, 1964.
130. J.-L. Roumiguieres, *Opt. Commun.* **19**:76–78, 1976.
131. K. Knop, *Opt. Commun.* **26**:281–283, 1978.
132. G. W. Stroke, *Phys. Lett.* (Neth.) **5**:45–48, 1963.
133. G. W. Stroke, *J. Opt. Soc. Am.* **54**:846, 1964.
134. C. W. Peters, T. F. Zipf, and P. V. Deibel, *J. Opt. Soc. Am.* **43**:816A, 1953.
135. A. Hadni, E. Décamps, and P. Delorme, *J. Phys. Radium* **19**(8):793–794, 1958.
136. A. Hadni, E. Décamps, D. Grandjean, and C. Janot, *C. R. Acad. Sci.* **250**:2007–2009, 1960.
137. A. Mitsuishi, Y. Yamada, S. Fujita, and H. Yoshinaga, *J. Opt. Soc. Am.* **50**:433–436, 1960.
138. C. Janot and A. Hadni, *J. Phys. Radium* **24**(8):1073–1077, 1963.
139. J. H. Rohrbaugh, C. Pine, W. G. Zoellner, and R. D. Hatcher, *J. Opt. Soc. Am.* **48**:710–711, 1958; [see also R. D. Hatcher and J. H. Rohrbaugh, *J. Opt. Soc. Am.* **46**:104–110, 1956 and J. H. Rohrbaugh and R. D. Hatcher, *J. Opt. Soc. Am.* **48**:704–709, 1958].
140. H. A. Kalthor and A. R. Neureuther, *J. Opt. Soc. Am.* **61**:43–48, 1971.
141. R. N. Hamm, R. A. MacRae, and E. T. Arakawa, *J. Opt. Soc. Am.* **55**:1460–1463, 1965.
142. K. Kudo, T. Arai, and T. Ogawa, *J. Opt. Soc. Am.* **60**:1046–1050, 1970.
143. V. G. Horton, E. T. Arakawa, R. N. Hamm, and M. W. Williams, *Appl. Opt.* **8**:667–670, 1969.
144. M. Schledermann and M. Skibowski, *Appl. Opt.* **10**:321–326, 1971.
145. N. J. Harrick, *J. Opt. Soc. Am.* **49**:376–379, 379–380, 1959.
146. D. F. Edwards and M. J. Bruemmer, *J. Opt. Soc. Am.* **49**:860–861, 1959.
147. M. Krizek, *Czech. J. Phys.* **13B**:599–610, 683–691, 1963.
148. J. P. Craig, R. F. Gribble, and A. A. Dougal, *Rev. Sci. Instrum.* **35**:1501–1503, 1964.
149. J. Bor, and L. A. Brooks, *J. Sci. Instrum.* **43**:944, 1966.
150. A. H. Pfund, *J. Opt. Soc. Am.* **37**:558–559, 1947.
151. P. Barchewitz and L. Henry, *J. Phys. Radium* ser. 8, **15**:639–640, 1954.
152. S. Takahashi, *J. Opt. Soc. Am.* **51**:441–444, 1961.
153. A. K. Walton and T. S. Moss, *Proc. Phys. Soc.* (GB) **78**:1393–1407, 1961
154. R. T. Baumel and S. E. Schnatterly, *J. Opt. Soc. Am.* **61**:832–833, 1971.
155. B. N. Grechushnikov and I. P. Petrov, *Opt. Spectrosc.* (USSR) **14**:160–161, 1963.
156. J. T. Cox and G. Hass, *Appl. Opt.* **17**:1657–1658, 1978.
157. T. F. Thonn and R. M. A. Azzam, *Opt. Eng.* **24**:202–206, 1985.
158. W. R. Hunter, *Japan J. Appl. Phys.* **4**(1):520, 1965; (*Proc. Conf. Photogr. Spectrosc. Opt.*, 1964).
159. H. Damany, *Opt. Acta* **12**:95–107, 1965.
160. G. Stephan, J.-C. Lemonnier, Y. LeCalvez, and S. Robin, *C. R. Acad. Sci.* **262B**:1272–1275, 1966.
161. F. de Chelle and H. Merdy, *C. R. Acad. Sci.* **265B**:968–971, 1967.
162. T. J. McIlrath, *J. Opt. Soc. Am.* **58**:506–510, 1968.
163. T. T. Cole and F. Oppenheimer, *Appl. Opt.* **1**:709–710, 1962.
164. E. Uzan, H. Damany, and J. Romand, *C. R. Acad. Sci.* **260**:5735–5737, 1965.
165. K. Rabinovitch, L. R. Canfield, and R. P. Madden, *Appl. Opt.* **4**:1005–1010, 1965.
166. T. Sasaki and H. Fukutani, *Japan J. Appl. Phys.* **3**:125–126, 1964.

167. H. Winter, H. H. Bukow, and P. H. Heckmann, *Opt. Commun.* **11**:299–300, 1974.
168. G. Hass and W. R. Hunter, *Appl. Opt.* **17**:76–82, 1978.
169. R. S. Spencer, G. J. Bergen, C. M. Fleetwood, H. Herzig, L. Miner, S. H. Rice, E. Smigocki, B. E. Woodgate, and J. J. Zaniewski, *Opt. Eng.* **24**:548–554, 1985.
170. W. R. Hunter, *Appl. Opt.* **17**:1259–1270, 1978.
171. A. Ejiri, *J. Phys. Soc. Japan* **23**:901, 1967.
172. G. Rosenbaum, B. Feuerbacher, R. P. Godwin, and M. Skibowski, *Appl. Opt.* **7**:1917–1920, 1968.
173. N. Rehfeld, U. Gerhardt, and E. Dietz, *Appl. Phys.* **1**:229–232, 1973.
174. H. A. Van Hoof, *Appl. Opt.* **19**:189–190, 1980.
175. R. Hibst and H. H. Bukow, *Appl. Opt.* **28**:1806–1812, 1989.
176. M. A. Khakoo, P. Hammond, and J. W. McConkey, *Appl. Opt.* **26**:3492–3494, 1987.
177. M. B. Robin, N. A. Kuebler, and Y.-H. Pao, *Rev. Sci. Instrum.* **37**:922–924, 1966.
178. A. Matsui and W. C. Walker, *J. Opt. Soc. Am.* **60**:64–65, 1970.
179. F. Abelès, *C. R. Acad. Sci.* **230**:1942–1943, 1950.
180. L. G. Schulz and F. R. Tangherlini, *J. Opt. Soc. Am.* **44**:362–368, 1954.
181. D. K. Burge and H. E. Bennett, *J. Opt. Soc. Am.* **54**:1428–1433, 1964.
182. G. K. T. Conn and G. K. Eaton, *J. Opt. Soc. Am.* **44**:553–557, 1954.
183. M. F. de la Provostaye and P. Desains, *Ann. Chim. Phys.* **30**(3):158, 1850.
184. G. R. Bird and W. A. Shurcliff, *J. Opt. Soc. Am.* **49**:235–237, 1959.
185. A. Elliott and E. J. Ambrose, *Nature* **159**:641–642, 1947.
186. A. Elliott, E. J. Ambrose, and R. Temple, *J. Opt. Soc. Am.* **38**:212–216, 1948.
187. J. Ames and A. M. D. Sampson, *J. Sci. Instrum.* **26**:132, 1949.
188. R. Duverney, *J. Phys. Radium* **20**(8)(suppl. 7):66A, 1959.
189. J. H. Hertz, *Exper. Tech. der Phys.* **7**:277–280, 1959.
190. K. Buijs, *Appl. Spectrosc.* **14**:81–82, 1960.
191. E. M. Bradbury and A. Elliott, *J. Sci. Instrum.* **39**:390, 1962.
192. R. G. Greenler, K. W. Adolph, and G. M. Emmons, *Appl. Opt.* **5**:1468–1469, 1966.
193. N. Wright, *J. Opt. Soc. Am.* **38**:69–70, 1948.
194. A. S. Makas and W. A. Shurcliff, *J. Opt. Soc. Am.* **45**:998–999, 1955.
195. H. E. Bennett, J. M. Bennett, and M. R. Nagel, *J. Opt. Soc. Am.* **51**:237, 1961.
196. R. T. Lagemann and T. G. Miller, *J. Opt. Soc. Am.* **41**:1063–1064, 1951.
197. L. Huldtt and T. Staflin, *Opt. Acta* **6**:27–36, 1959.
198. T. A. Leonard, J. Loomis, K. G. Harding, and M. Scott, *Opt. Eng.* **21**:971–975, 1982.
199. R. Meier and H. H. Günthard, *J. Opt. Soc. Am.* **49**:1122–1123, 1959.
200. N. J. Harrick, *J. Opt. Soc. Am.* **54**:1281–1282, 1964.
201. N. P. Murarka and K. Wilner, *Appl. Opt.* **20**:3275–3276, 1981.
202. S. D. Smith, T. S. Moss, and K. W. Taylor, *J. Phys. Chem. Solids* **11**:131–139, 1959.
203. D. T. Rampton and R. W. Grow, *Appl. Opt.* **15**:1034–1036, 1976.
204. J.-M. Munier, J. Claudel, E. Décamps, and A. Hadni, *Rev. Opt.* **41**:245–253, 1962.
205. A. K. Walton, T. S. Moss, and B. Ellis, *J. Sci. Instrum.* **41**:687–688, 1964.
206. W. C. Walker, *Appl. Opt.* **3**:1457–1459, 1964.
207. D. C. Hinson, *J. Opt. Soc. Am.* **56**:408, 1966.
208. D. F. Heath, *Appl. Opt.* **7**:455–459, 1968.
209. J. Schellman, V. Chandrasekharan, and H. Damany, *C. R. Acad. Sci.* **259**:4560–4563, 1964.
210. G. Weiser, *Proc. IEEE* **52**:966, 1964.

211. J. L. Weinberg, *Appl. Opt.* **3**:1057–1061, 1964.
212. J. M. Bennett, D. L. Decker, and E. J. Ashley, *J. Opt. Soc. Am.* **60**:1577A, 1970.
213. H. Schröder, *Optik* **3**:499–503, 1948.
214. F. Abelès, *J. Phys. Radium* **11**(8):403–406, 1950.
215. K. Kubo, *Tokyo Inst. Phys. Chem. Res. J. Sci. Res. Instrum.* **47**:1–6, 1953.
216. H. Schopper, *Optik* **10**:426–438, 1953.
217. M. Ruiz-Urbietta and E. M. Sparrow, *J. Opt. Soc. Am.* **62**:1188–1194, 1972.
218. M. Ruiz-Urbietta and E. M. Sparrow, *J. Opt. Soc. Am.* **63**:194–200, 1973.
219. M. Ruiz-Urbietta, E. M. Sparrow, and G. W. Goldman, *Appl. Opt.* **12**:590–596, 1973.
220. F. Abelès, *Proc. Conf. Photogr. Spectrosc. Opt.*, 1964, *Japan J. Appl. Phys.* **4**(suppl. 1):517, 1965.
221. A.-R. M. Zaghoul and R. M. A. Azzam, *Appl. Opt.* **16**:1488–1489, 1977.
222. W. W. Buchman, S. J. Holmes, and F. J. Woodberry, *J. Opt. Soc. Am.* **61**:1604–1606, 1971.
223. W. W. Buchman, *Appl. Opt.* **14**:1220–1224, 1975.
224. S. Refermat and J. Eastman, *Proc. Soc. Photo-Opt. Instrum. Eng.* **88**:28–33, 1976.
225. L. Songer, *Optical Spectra* **12**(10):49–50, October 1978.
226. D. Blanc, P. H. Lissberger, and A. Roy, *Thin Solid Films* **57**:191–198, 1979.
227. T. Maehara, H. Kimura, H. Nomura, M. Yanagihara, and T. Namioka, *Appl. Opt.* **30**:5018–5020, 1991.
228. P. B. Clapham, M. J. Downs, and R. J. King, *Appl. Opt.* **8**:1965–1974, 1969. [See also P. B. Clapham, *Thin Solid Films* **4**:291–305, 1969].
229. S. M. MacNeille, U.S. Patent 2,403,731, July 9, 1946.
230. M. Banning, *J. Opt. Soc. Am.* **37**:792–797, 1947.
231. J. A. Dobrowolski and A. Waldorf, *Appl. Opt.* **20**:111–116, 1981.
232. J. C. Monga, P. D. Gupta, and D. D. Bhawalkar, *Appl. Opt.* **23**:3538–3540, 1984.
233. J. Mouchart, J. Begel, and E. Duda, *Appl. Opt.* **28**:2847–2853, 1989.
234. D. Lees and P. Baumeister, *Opt. Lett.* **4**:66–67, 1979.
235. R. M. A. Azzam, *Opt. Lett.* **10**:110–112, 1985.
236. R. M. A. Azzam, *Appl. Opt.* **25**:4225–4227, 1986.
237. H. Schröder and R. Schläfer, *Z. Naturforsch.* **4a**:576–577, 1949.
238. G. Pridatko and T. Krylova, *Opt.-Mekh. Prom.* **3**:23, 1958.
239. R. S. Sokolova and T. N. Krylova, *Opt. Spectrosc. (USSR)* **14**:213–215, 1963.
240. A. F. Turner and P. W. Baumeister, *Appl. Opt.* **5**:69–76, 1966.
241. P. B. Clapham, *Opt. Acta* **18**:563–575, 1971.
242. H. Schröder, *Optik* **13**:158–168, 169–174, 1956.
243. P. Baumeister, *Opt. Acta* **8**:105–119, 1961.
244. P. Baumeister, Institute of Optics, Univ. Rochester, Rochester, N.Y., private communication, 1971.
245. V. R. Costich, *Appl. Opt.* **9**:866–870, 1970.
246. C. D. West and A. S. Makas, *J. Opt. Soc. Am.* **39**:791–794, 1949.
247. S. D. Jacobs, Y. Asahara, and T. Izumitani, *Appl. Opt.* **21**:4526–4532, 1982.
248. B. H. Billings, in D. E. Gray, (ed.), *American Institute of Physics Handbook*, 3rd ed., McGraw-Hill, New York, 1972, pp. 6-37, 6-40, 6-46, 6-112, and 6-113.
249. J. M. Beckers, *Appl. Opt.* **10**:973–975, 1971.
250. W. L. Bond, *J. Appl. Phys.* **36**:1674–1677, 1965.
251. J. H. Shields and J. W. Ellis, *J. Opt. Soc. Am.* **46**:263–265, 1956.
252. E. D. Palik and B. W. Hennis, *Appl. Opt.* **6**:2198–2199, 1967.
253. G. D. Boyd, W. L. Bond, and H. L. Carter, *J. Appl. Phys.* **38**:1941–1943, 1967.
254. M. A. Jeppesen, *J. Opt. Soc. Am.* **48**:629–632, 1958.

255. S. S. Ballard, J. S. Browder, and J. F. Ebersole, in D. E. Gray (ed.), *American Institute of Physics Handbook*, 3d ed., McGraw-Hill, New York, 1972, pp. 6–20, 6–27, and 6–35.
256. A. E. Ennos and K. W. Opperman, *Appl. Opt.* **5**:170, 1966.
257. E. D. Palik, *Appl. Opt.* **7**:978–979, 1968.
258. H. Gobrecht and A. Bartschat, *Z. Phys.* **156**:131–143, 1959.
259. M. S. Shumate, *Appl. Opt.* **5**:327–332, 1966.
260. E. V. Loewenstein, *J. Opt. Soc. Am.* **51**:108–112, 1961.
261. V. Chandrasekharan and H. Damany, *Appl. Opt.* **7**:939–941, 1968.
262. J.-P. Maillard, *Opt. Commun.* **4**:175–177, 1971.
263. T. M. Bieniewski and S. J. Czyzak, *J. Opt. Soc. Am.* **53**:496–497, 1963.
264. M. Françon, S. Mallick and J. Vulmière, *J. Opt. Soc. Am.* **55**:1553, 1965.
265. F. Zernicke, Jr., *J. Opt. Soc. Am.* **54**:1215–1220, 1964, [erratum in *J. Opt. Soc. Am.* **55**:210–211, 1965].
266. E. Einsporn, *Phys. Z.* **37**:83–88, 1936.
267. E. L. Gieszelmann, S. F. Jacobs, and H. E. Morrow, *J. Opt. Soc. Am.* **59**:1381–1383, 1969 [erratum in *J. Opt. Soc. Am.* **60**:705, 1970].
268. J. M. Daniels, *Rev. Sci. Instrum.* **38**:284–285, 1967.
269. T. Saito, A. Ejiri, and H. Onuki, *Appl. Opt.* **29**:4538–4540, 1990.
270. W. B. Westerveld, K. Becker, P. W. Zetner, J. J. Corr, and J. W. McConkey, *Appl. Opt.* **14**:2256–2262, 1985.
271. P. L. Richards and G. E. Smith, *Rev. Sci. Instrum.* **35**:1535–1537, 1964.
272. L. H. Johnston, *Appl. Opt.* **16**:1082–1084, 1977.
273. D. P. Gonatas, X. D. Wu, G. Novak, and R. H. Hildebrand, *Appl. Opt.* **28**:1000–1006, 1989.
274. P. Lostis, *J. Phys. Radium*, **18**(8):51S, 1957.
275. C. E. Greninger, *Appl. Opt.* **27**:774–776, 1988.
276. J. M. Bennett, *Appl. Opt.* **9**:2123–2129, 1970.
277. S. Nakadate, *Appl. Opt.* **29**:242–246, 1990.
278. P. D. Hale and G. W. Day, *Appl. Opt.* **27**:5146–5153, 1988.
279. A.-R. M. Zaghoul, R. M. A. Azzam, and N. M. Bashara, *Opt. Commun.* **14**:260–262, 1975.
280. A.-R. M. Zaghoul, R. M. A. Azzam, and N. M. Bashara, *J. Opt. Soc. Am.* **65**:1043–1049, 1975.
281. S. Kawabata and M. Suzuki, *Appl. Opt.* **19**:484–485, 1980.
282. S. Chu, R. Conti, P. Bucksbaum, and E. Commins, *Appl. Opt.* **18**:1138–1139, 1979.
283. J. Strong, *Procedures in Experimental Physics*, Prentice-Hall, Englewood Cliffs, N.J., 1938, pp. 388–389.
284. P. H. Smith, *Proc. Symp. Recent Dev. Ellipsometry, Surf. Sci.* **16**:34–66, 1969.
285. B. N. Grechushnikov, *Opt. Spectrosc. (USSR)* **12**:69, 1962.
286. D. G. Drummond, *Proc. Roy. Soc. (Lond.)* **153A**:318–339, 1936.
287. G. Holzwarth, *Rev. Sci. Instrum.* **36**:59–63, 1965.
288. E. D. Palik, *Appl. Opt.* **2**:527–539, 1963.
289. J. H. Jaffe, H. Jaffe, and K. Rosenbeck, *Rev. Sci. Instrum.* **38**:935–938, 1967.
290. C. Gaudefroy, *C. R. Acad. Sci.* **189**:1289–1291, 1929.
291. W. A. Deer, R. A. Howie, and J. Zussman, *Rock-forming Minerals 3; Sheet Silicates*, Wiley, New York, 1962, pp. 258–262.
292. A. E. Ennos, *J. Sci. Instrum.* **40**:316–317, 1963.
293. S. B. Ioffe and T. A. Smirnova, *Opt. Spectrosc. (USSR)* **16**:484–485, 1964.
294. S. Mitchell, *Nature* **212**:65–66, 1966.
295. A. M. Title, *Appl. Opt.* **14**:229–237, 1975.
296. R. Anderson, *Appl. Opt.* **27**:2746–2747, 1988.
297. P. L. Wizinowich, *Opt. Eng.* **28**:157–159, 1989.

298. G. N. Ramachandran and S. Ramaseshan, "Crystal Optics," in S. Flügge (ed.), *Handbuch der Physik* **25/1**, Springer, Berlin, 1961, pp. 156–158.
299. A. Johannsen, *Manual of Petrographic Methods*, 2d ed., Hafner, New York, 1968, pp. 369–385 (originally published in 1918).
300. H. G. Jerrard, *J. Opt. Soc. Am.* **38**:35–59, 1948.
301. R. P. Hunt, *Appl. Opt.* **9**:1220–1221, 1970.
302. J. R. Partington, *An Advanced Treatise on Physical Chemistry*, **4**, Wiley, New York, 1953, pp. 173–177.
303. J. Ellerbroek and J. T. Groosmuller, *Phys. Z.* **27**:468–471, 1926.
304. H. G. Jerrard, *J. Opt. Soc. Am.* **39**:1031–1035, 1949.
305. H. G. Jerrard, *J. Sci. Instrum.* **26**:353–357, 1949.
306. H. G. Jerrard, *J. Sci. Instrum.* **27**:62–66, 1950.
307. J. F. Archard, *J. Sci. Instrum.* **27**:238–241, 1950.
308. R. W. Wood, *Physical Optics*, 3d ed., Macmillan, New York, 1934, pp. 356–361.
309. R. W. Ditchburn, *Light*, 2d ed., Interscience, New York, 1963*b*, pp. 483–485.
310. H. G. Jerrard, *J. Sci. Instrum.* **27**:164–167, 1950.
311. H. G. Jerrard, *J. Sci. Instrum.* **28**:10–14, 1951.
312. H. G. Jerrard, *J. Sci. Instrum.* **30**:65–70, 1953.
313. H. E. Ives and H. B. Briggs, *J. Opt. Soc. Am.* **26**:238–246, 1936.
314. E. D. Palik, *Appl. Opt.* **4**:1017–1021, 1965.
315. P. Hariharan and D. Sen, *J. Sci. Instrum.* **37**:278–281, 1960.
316. M. P. Varnham, D. N. Payne, A. J. Barlow, and E. J. Tarbox, *Opt. Lett.* **9**:306–308, 1984.
317. C. A. Villarruel, M. Abebe, W. K. Burns, and R. P. Moeller, in *Digest of the Seventh Topical Conference on Optical Fiber Communication*, **84.1**, Optical Society of America, Washington, D.C., 1984.
318. R. A. Bergh, H. C. Lefevre, and H. J. Shaw, *Opt. Lett.* **5**:479–481, 1980.
319. W. Eickhoff, *Electron. Lett.* **16**:762–763, 1980.
320. D. Gruchmann, K. Petermann, L. Satandigel, and E. Weidel, in *Proceedings of the European Conference on Optical Communication*, North Holland, Amsterdam, 1983, pp. 305–308.
321. T. Hosaka, K. Okamoto, and T. Edahiro, *Opt. Lett.* **8**:124–126, 1983.
322. J. R. Feth and C. L. Chang, *Opt. Lett.* **11**:386–388, 1986.
323. T. Yu and Y. Wu, *Opt. Lett.* **13**:832–834, 1988.
324. R. B. Dyott, J. Bello, and V. A. Handerek, *Opt. Lett.* **12**:287–289, 1987.
325. S. C. Lee and J.-I. Chen, *Appl. Opt.* **29**:2667–2668, 1990.
326. R. Ulrich and M. Johnson, *Appl. Opt.* **18**:1857–1861, 1979.
327. T. Hosaka, K. Okamoto, and T. Edahiro, *Appl. Opt.* **22**:3850–3858, 1983.
328. A. W. Snyder and A. J. Stevenson, *Opt. Lett.* **11**:254–256, 1986.
329. S. Uehara, T. Izawa, and H. Nakagome, *Appl. Opt.* **13**:1753–1754, 1974.
330. H. F. Mahlein, *Opt. Commun.* **16**:420–424, 1976.
331. P. G. Suchoski, T. K. Findakly, and F. J. Leonberger, *Opt. Lett.* **13**:172–174, 1988.
332. R. Noé, H. Heidrich and D. Hoffmann, *Opt. Lett.* **13**:527–529, 1988.
333. K. Baba, K. Shiraishi, K. Obi, T. Kataoka, and S. Kawakami, *Appl. Opt.* **27**:2554–2560, 1988.

This page intentionally left blank.

DO NOT DUPLICATE

MUELLER MATRICES

Russell A. Chipman

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

14.1 GLOSSARY

A	analyzer vector
d	diattenuation parameters, set of three
d_i	magnitude of polarization parameters via matrix roots
D	diattenuation
D	diagonal matrix
D_i	differential polarization parameters via matrix roots
DD	diagonal depolarizer Mueller matrix
Dep	depolarization
DI	depolarization index
DoCP	degree of circular polarization
DoLP	degree of linear polarization
DoP	degree of polarization
e	ellipticity
E	extinction ratio
E	Jones vector
E_x, E_y	electric field components
ED	elliptical diattenuator Mueller matrix
EP	elliptical polarizer Mueller matrix
ER	elliptical retarder Mueller matrix
G	Mueller matrix generators for polarization properties
H	hermitian coherency matrix
\hat{H}	horizontal polarized Stokes vector, 0° , normalized
H_p	hermitian coherency matrix of a physical Mueller matrix
HLP	horizontal linear polarizer Mueller matrix
HQWLR	quarter-wave linear retarder Mueller matrix, horizontal fast axis

HWLR	half-wave linear retarder Mueller matrix
i	$\sqrt{-1}$
I	inhomogeneity of a Mueller matrix
I	identity matrix
ID	ideal depolarizer Mueller matrix
J	Jones matrix
$j_{xx}, j_{xy}, j_{yx}, j_{yy}$	Jones matrix elements
\mathbf{k}	propagation vector
LCP	left circular polarizer Mueller matrix
LD	linear diattenuator operator
LD	linear diattenuator Mueller matrix
LDR	homogeneous linear diattenuator and retarder Mueller matrix
LP(θ)	linear polarizer Mueller matrix transmitting along axis θ
M	Mueller matrix
$\hat{\mathbf{M}}$	normalized Mueller matrix
\mathbf{M}_D	diattenuator Mueller matrix
\mathbf{M}_N	nondepolarizing Mueller matrix, Mueller-Jones matrix
\mathbf{M}_P	physical Mueller matrix
\mathbf{M}_R	retarder Mueller matrix
$\mathbf{M}_{\text{refl}}, \mathbf{M}_{\text{refr}}$	Mueller matrices for reflection and refraction
$m_{00}, m_{01}, \dots, m_{33}$	Mueller matrix elements
n_1, n_2	refractive indices of modes
O	orthogonal matrix
P	polarizance
P	high-order matrix root of Mueller matrix
PD	partial depolarizer Mueller matrix
PDL	polarization-dependent loss
q	index for a sequence of polarization elements
q	index for mode order
Q	index limit
QWLR	quarter-wave linear retarder Mueller matrix
QWRCR, QWLRCR	quarter-wave circular retarder Mueller matrix, right and left fast mode
\mathbf{R}_M	rotational change of basis matrix for Stokes vectors
RCP	right circular polarizer Mueller matrix
Re	real part
\mathbf{s}	Stokes three-vector
S	Stokes vector
$\hat{\mathbf{S}}$	normalized polarized Stokes vector
\mathbf{S}'	exiting Stokes vector
$\mathbf{S}_{\text{max}}, \mathbf{S}_{\text{min}}$	incident Stokes vectors of maximum and minimum intensity transmittance
S_0, S_1, S_2, S_3	Stokes vector elements
t	thickness
t	time
T	transpose, superscript
T	intensity transmittance

T_{avg}	intensity transmission averaged over all incident polarization states
T_{max}	maximum intensity transmittance
T_{min}	minimum intensity transmittance
Tr	trace of a matrix
U	unitary matrix
U	Jones / Mueller transformation matrix
$\hat{\mathbf{U}}$	unpolarized Stokes vector, normalized
v	eigenvectors
$\hat{\mathbf{V}}$	vertical polarized Stokes vector, 90°, normalized
VD	variable partial depolarizer Mueller matrix
VQWLR	quarter-wave linear retarder Mueller matrix, vertical fast axis
$\{\hat{\mathbf{x}}, \hat{\mathbf{y}}\}$	Stokes vector coordinate basis
z	spatial coordinate
δ	retardance, radians
$\delta_H, \delta_{45}, \delta_R$	retardance components: horizontal, 45°, right
ε	latitude on Poincaré sphere
θ	orientation angle: (1) polarizer axis, (2) retarder fast axis, (3) major axis of polarization ellipse
λ	wavelength
λ	eigenvalue
ρ	amplitudes of a complex number
σ_i	identity matrix, $i = 0$, and normalized Pauli spin matrices, $i = 1, 2, 3$
ϕ	phase of a complex number
χ	angle between the eigenpolarizations on the Poincaré sphere
ω	angular frequency
\otimes	tensor product
\dagger	hermitian adjoint
\cdot	dot product, matrix multiplication
$\ \ $	norm of a vector

14.2 CONVENTIONS

All angles are in radians unless the degree sign ($^\circ$) is used. Retardance is specified in radians throughout. The last Stokes vector element, S_3 , is positive for a right circularly polarized component and negative for a left circularly polarized component. All vectors and matrices are represented by bold characters.

14.3 OBJECTIVES

This chapter surveys the Mueller matrix and its properties. The Mueller matrix has become the principal quantity used in polarimetric measurements of optical and polarization elements. For optical design and theoretical analyses, particularly of interferometers, Jones matrices and coherence matrices are often preferred. Mueller matrices are straightforward to measure and dominate experimental studies.

Despite the Mueller matrix's straightforward definition, the relation between the polarization properties and the matrix elements is complex, particularly when depolarization is involved.

Several issues are addressed:

1. Determining the Mueller matrix from the specification of a polarization element.
2. Given the Jones matrix for a polarization element, determine the corresponding Mueller matrix.
3. Given a Mueller matrix, determine if there is a corresponding Jones matrix and calculate this Jones matrix.
4. Given a Mueller matrix, determine the corresponding polarization properties of diattenuation, retardance, and depolarization.
5. Given a 4×4 matrix which violates the constraints on Mueller matrices, find the closest Mueller matrices.

This chapter supports Chap. 15, "Polarimetry" which assumes much of the material here. Chapter 15 is principally concerned with measuring Stokes parameters and Mueller matrices; this chapter treats Mueller matrix calculations and data reduction.

14.4 STOKES PARAMETERS AND MUELLER MATRICES

Several calculi have been developed for analyzing polarization, including those based on the Jones matrix, coherency matrix, Mueller matrix, and other matrices.¹⁻¹⁰ Of these methods, the Mueller calculus is most generally suited for describing irradiance-measuring instruments, including most polarimeters, radiometers, and spectrometers.

The set of four linear equations relating incident and exiting Stokes parameters was first introduced by Soleillet in 1929.¹¹ Hans Müller of MIT formulated these equations as a 4×4 matrix times a Stokes vector in his class notes, an Optical Society of America meeting abstract, and in a technical report but he never published a journal article.¹² His graduate student Parke developed the matrix properties in great detail.^{13,14} R. Clark Jones became aware of the work prior to publication and was the first to use the term and spelling Mueller matrix in a journal article.¹⁵ In keeping with current practice, we will refer to this matrix as the Mueller matrix, although it is sometimes referred to as the Stokes matrix.¹⁶

In the Mueller calculus, the Stokes vector \mathbf{S} describes the polarization state of a light beam, and the Mueller matrix \mathbf{M} describes the polarization-altering characteristics of a sample. This sample may be a surface, a polarization element, an optical system, or some other light/matter interaction which produces a reflected, refracted, diffracted, or scattered light beam. Chapter 15 "Polarimetry" contains a detailed description of Stokes vector properties.

14.5 THE STOKES PARAMETERS AND THE POINCARÉ SPHERE

The Stokes parameters (Stokes vector) can be normalized by its flux S_0 , and used to define the Stokes three-vector \mathbf{s} ,

$$\hat{\mathbf{S}} = \frac{\mathbf{S}}{S_0} = \begin{pmatrix} 1 \\ s_1 \\ s_2 \\ s_3 \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{s} \end{pmatrix} \quad \mathbf{s} = \{s_1, s_2, s_3\} \quad (1)$$

For unpolarized light $\mathbf{s} = \{0, 0, 0\}$ and the magnitude of \mathbf{s} ,

$$|\mathbf{s}| = \sqrt{s_1^2 + s_2^2 + s_3^2} = 0 \quad (2)$$

For a completely polarized polarization state, $|\mathbf{s}| = 1$. The vector \mathbf{s} is the coordinates of $\hat{\mathbf{S}}$ on the unit Poincaré sphere. The introduction of \mathbf{s} allows the polarization state to be specified irrespective of the flux. Completely polarized states \mathbf{s}_a and \mathbf{s}_b are orthogonal when

$$\mathbf{s}_a = -\mathbf{s}_b \quad (3)$$

Orthogonal states have opposite helicity (the electric fields rotate clockwise and counterclockwise) and the orientations of the polarization ellipse major axes are 90° apart.

The Poincaré sphere is a geometrical construction for the representation of Stokes vectors and polarization ellipses where the Stokes three-vector is plotted in a three-dimensional space with axes $\{S_1, S_2, S_3\}$ as shown in Fig 1.

Each point on the surface of the Poincaré sphere can be parameterized by angles $\{\theta, \varepsilon\}$ where θ is the orientation of the major axis of the polarization ellipse and ε is the latitude; $\sin \varepsilon$ is the degree of circular polarization,

$$\mathbf{S}(\theta, \varepsilon) = \begin{pmatrix} 1 \\ \cos(2\theta)\cos\varepsilon \\ \sin(2\theta)\cos\varepsilon \\ \sin\varepsilon \end{pmatrix} \quad (4)$$

Linearly polarized Stokes parameters are located around the Poincaré sphere equator, $\{\cos(2\theta), \sin(2\theta), 0\}$, where θ is the orientation of linear polarization. The north pole $\{0, 0, 1\}$ represents right circularly polarized light and the south pole $\{0, 0, -1\}$ left circularly polarized light. The sphere's center $\{0, 0, 0\}$ represents unpolarized light. The set of spheres centered on the origin each contain partially polarized states with a degree of polarization equal to the radius r . The surface at radius one,

$$r = 1 = \sqrt{s_1^2 + s_2^2 + s_3^2} \quad (5)$$

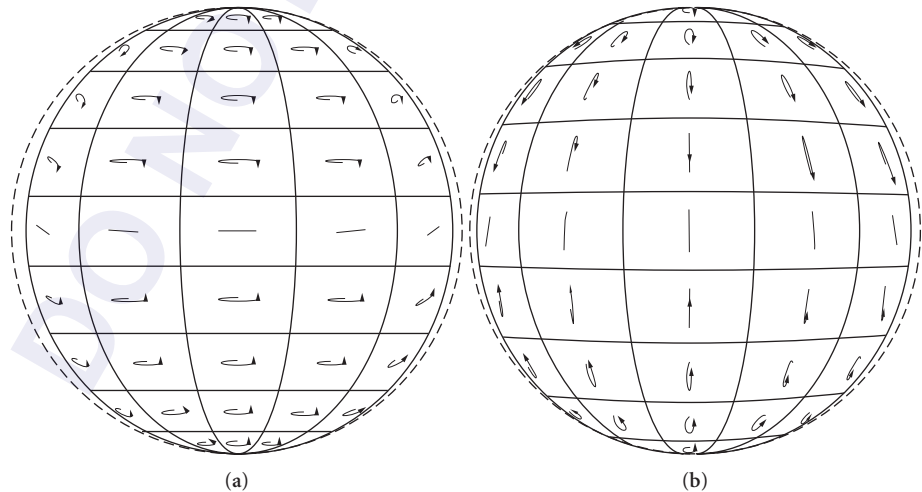


FIGURE 1 View of the Poincaré sphere along the $+S_1$ axis (a) and along the $-S_1$ axis (b) with the polarization ellipses associated with different locations indicated.

represents all possible states of (completely) polarized light. θ is one half the latitude on a traditional globe because, for incoherent light, a rotation of 180° returns the polarization state to its initial state; rotate a linear polarizer through 180° and the polarization states repeat. Notice that on the Poincaré sphere, the locations for 0° linearly polarized light (nominally horizontal $\{1, 0, 0\}$) and 45° linearly polarized light $\{0, 1, 0\}$ are 90° apart while horizontal and vertical $\{-1, 0, 0\}$ linearly polarized light are diametrically opposite, 180° apart. Orthogonal polarization states are at opposite points on the sphere surface.

The Stokes parameters have an unusual coordinate system because the S_1 , S_2 , and S_3 axes do not represent polarization states 90° apart, the traditional definition of orthogonal vectors. The Stokes parameter coordinate system is a clever and effective representation of incoherent light because equal amounts of orthogonal polarized fluxes, when combined, yield unpolarized light. Due to the properties of this coordinate system, the four Stokes parameters do not transform as a vector, and cannot be considered as a true *vector*. The Stokes parameters do add as vectors and are operated on by Mueller matrices like vectors; thus the widespread use of the term Stokes vector.

In current practice, the Poincaré sphere is used three different ways:

1. To represent a polarization state
2. To represent diattenuation by indicating the Stokes parameters of maximum transmission, and the diattenuation magnitude as distance from the origin
3. To represent retardance by indicating the axis through the origin of the fast and slow axes (eigenpolarizations), or by representing retardance within a three dimensional *retardance space* with components $\{\delta_x, \delta_{45}, \delta_R\}$

14.6 MUELLER MATRICES

The Mueller matrix is a 4×4 matrix with real-valued elements.^{10,14} The Mueller matrix \mathbf{M} for a polarization-altering device is defined as the matrix which transforms an incident Stokes vector \mathbf{S} into the exiting (reflected, transmitted, or scattered) Stokes vector \mathbf{S}' ,

$$\mathbf{M} \cdot \mathbf{S} = \mathbf{S}' = \begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \cdot \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} = \begin{pmatrix} S'_0 \\ S'_1 \\ S'_2 \\ S'_3 \end{pmatrix} \quad (6)$$

$$= \begin{pmatrix} S_0 m_{0,0} + S_1 m_{0,1} + S_2 m_{0,2} + S_3 m_{0,3} \\ S_0 m_{1,0} + S_1 m_{1,1} + S_2 m_{1,2} + S_3 m_{1,3} \\ S_0 m_{2,0} + S_1 m_{2,1} + S_2 m_{2,2} + S_3 m_{2,3} \\ S_0 m_{3,0} + S_1 m_{3,1} + S_2 m_{3,2} + S_3 m_{3,3} \end{pmatrix}$$

Each element of the incident \mathbf{S} is related to the four elements of \mathbf{S}' by the elements of \mathbf{M} . Since the elements of \mathbf{S} and \mathbf{S}' are irradiances, the elements of \mathbf{M} are dimensionless ratios of irradiances. Since irradiances are real, the elements of \mathbf{M} are real valued, not complex numbers. When the Mueller matrix is known, then the exiting polarization state is known for an arbitrary incident polarization state. Our convention numbers the subscripts from 0 to 3 to match the corresponding Stokes vector subscripts.

The Mueller matrix $\mathbf{M}(\mathbf{k}, \lambda)$ for a device is always a function of the direction of propagation \mathbf{k} and wavelength λ .

14.7 SEQUENCES OF POLARIZATION ELEMENTS

The Mueller matrix \mathbf{M} associated with a beam path through a sequence (cascade) of polarization elements $q = 1, 2, \dots, Q$ is the right-to-left product of the individual matrices \mathbf{M}_q ,

$$\mathbf{M} = \mathbf{M}_Q \cdot \mathbf{M}_{Q-1} \cdot \dots \cdot \mathbf{M}_q \cdot \dots \cdot \mathbf{M}_2 \cdot \mathbf{M}_1 = \prod_{q=Q, -1}^1 \mathbf{M}_q \quad (7)$$

In evaluating cascades of Mueller matrices, the associative rule for matrix multiplication can be applied,

$$(\mathbf{M}_3 \cdot \mathbf{M}_2) \cdot \mathbf{M}_1 = \mathbf{M}_3 \cdot (\mathbf{M}_2 \cdot \mathbf{M}_1) \quad (8)$$

and adjacent matrices grouped in any order for multiplication.

14.8 POLARIZATION ELEMENTS' PROPERTIES IN THE MUELLER CALCULUS

For ideal polarization elements, the polarization properties are readily defined. For real polarization elements, the precise description of the polarization properties is more complex. Polarization elements such as polarizers, retarders, and depolarizers have three general polarization properties: diattenuation, retardance, and depolarization, and a typical element displays some amount of all three. Diattenuation arises when the intensity transmittance of an element is a function of the incident polarization state.¹⁷ The diattenuation D of a device is defined in terms of the maximum T_{\max} and minimum T_{\min} intensity transmittances,

$$D = \frac{T_{\max} - T_{\min}}{T_{\max} + T_{\min}} \quad (9)$$

For an ideal polarizer, $D = 1$. When $D = 0$, all incident polarization states are transmitted with equal attenuation. For an ideal retarder the polarization states change upon transmission but T_{\max} and T_{\min} are equal and $D = 0$. The quality of a polarizer is often expressed in terms of the related quantity, the extinction ratio E ,

$$E = \frac{T_{\max}}{T_{\min}} = \frac{1+D}{1-D} \quad (10)$$

where the ideal polarizer has $E = \infty$.

Retardance is the phase change a device introduces between its eigenpolarizations (eigenstates). For a birefringent retarder with refractive indices n_1 and n_2 , and thickness t , the retardance δ expressed in radians is

$$\delta = \frac{2\pi(n_1 - n_2)t}{\lambda} \quad (11)$$

Depolarization describes the coupling by a device of incident polarized light into depolarized light in the exiting beam. For example, depolarization occurs when light transmits through milk or scatters from clouds. Multimode optical fibers generally depolarize the light. Depolarization is intrinsically associated with scattering and a loss of coherence in the polarization state. A small amount of depolarization is associated with the scattered light from all optical components.

14.9 ROTATION OF AN ELEMENT ABOUT THE OPTICAL AXIS

When a polarization element with Mueller matrix \mathbf{M} is rotated about the incident beam of light by an angle θ such that the angle of incidence is unchanged (for example, for a normal-incidence beam, rotating the element about the normal), the resulting Mueller matrix $\mathbf{M}(\theta)$ is

$$\mathbf{M}(\theta) = \mathbf{R}_M(\theta) \cdot \mathbf{M} \cdot \mathbf{R}_M(-\theta) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos 2\theta & -\sin 2\theta & 0 \\ 0 & \sin 2\theta & \cos 2\theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} m_{00} & m_{01} & m_{02} & m_{03} \\ m_{10} & m_{11} & m_{12} & m_{13} \\ m_{20} & m_{21} & m_{22} & m_{23} \\ m_{30} & m_{31} & m_{32} & m_{33} \end{pmatrix} \quad (12)$$

$$\cdot \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos 2\theta & \sin 2\theta & 0 \\ 0 & -\sin 2\theta & \cos 2\theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

where \mathbf{R}_M is the rotational change of basis matrix for Stokes vectors and Mueller matrices.

14.10 NONPOLARIZING MUELLER MATRICES

A nonpolarizing matrix does not change the polarization state of any incident polarization vector; only the amplitude and/or phase change. The Mueller Matrix for a nonabsorbing, nonpolarizing sample is the identity matrix \mathbf{I} ,

$$\mathbf{I} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (13)$$

\mathbf{I} is the Mueller matrix for vacuum and the approximate Mueller matrix for air. For a neutral density filter or polarization-independent absorption or loss, the Mueller matrix has $T_{\max} = T_{\min} = T$, and the resulting Mueller matrix is proportional to the identity matrix and can be expressed in terms of our notation for linear diattenuators, $\mathbf{LD}(T_{\max}, T_{\min}, \theta)$, as

$$\mathbf{LD}(T, T, 0) = T \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (14)$$

14.11 MUELLER MATRICES OF IDEAL POLARIZERS

First, the properties of an example ideal polarizer are examined. Tables of ideal polarizer Mueller matrices are presented followed by equations for linear and elliptical ideal polarizers.

An ideal polarizer has a transmittance of one for its principal state and a transmittance of zero for the orthogonal "blocked" state.

Consider the Mueller matrix for a horizontal linear polarizer (**HLP**), which we also express as **LD** (1, 0, 0) for linear diattenuator, $T_{\max} = 1$, $T_{\min} = 0$, orientation of transmission axis 0,

$$\mathbf{HLP} = \mathbf{LD} (1, 0, 0) = \frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (15)$$

When operating on the Stokes vector for horizontal linearly polarized light,

$$\mathbf{HLP} \cdot \hat{\mathbf{H}} = \frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} \quad (16)$$

horizontally polarized light exits without loss. Vertically polarized incident light is completely blocked,

$$\mathbf{HLP} \cdot \hat{\mathbf{V}} = \frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ -1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad (17)$$

In fact, **H** and **V** are the two eigenpolarizations of **HLP**. Eigenpolarizations are the eigenvectors which correspond to physically realizable Stokes vectors. For **HLP** the remaining two eigenvectors are nonphysical as Stokes vectors since $S_0 = 0$,

$$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix} \quad (18)$$

For an arbitrary incident Stokes vector,

$$\mathbf{HLP} \cdot \mathbf{S} = \mathbf{S}' = \frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} S_0 + S_1 \\ S_0 + S_1 \\ 0 \\ 0 \end{pmatrix} = \frac{S_0 + S_1}{2} \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} \quad (19)$$

Since the first two rows of **M** are equal, the first two elements of **S'** are equal, the S_2 and S_3 characteristics of the incident light are lost, and the exiting light is always horizontally linearly polarized. Table 1 lists the Mueller matrices for the six basis polarization states and the general linear polarizer.

TABLE 1 Mueller Matrices for Ideal Polarizers for the Basis Polarization States

Type of Polarizer	Symbol	Mueller Matrix
Horizontal linear polarizer	HLP	$\frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$
Vertical linear polarizer	VLP	$\frac{1}{2} \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$
45° linear polarizer	LP(45°)	$\frac{1}{2} \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$
135° linear polarizer	LP(135°)	$\frac{1}{2} \begin{pmatrix} 1 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$
Right circular polarizer	RCP	$\frac{1}{2} \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \end{pmatrix}$
Left circular polarizer	LCP	$\frac{1}{2} \begin{pmatrix} 1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 \end{pmatrix}$
Linear polarizer with transmission axis oriented at angle θ , measured positive when rotating counterclockwise from the x -axis looking into the beam	LP(θ)	$\frac{1}{2} \begin{pmatrix} 1 & \cos 2\theta & \sin 2\theta & 0 \\ \cos 2\theta & \cos^2 2\theta & \sin 2\theta \cos 2\theta & 0 \\ \sin 2\theta & \sin 2\theta \cos 2\theta & \sin^2 2\theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$

For an elliptical polarizer which transmits a polarization state with the major axis of the ellipse oriented at θ located at latitude ϕ on the Poincaré sphere, $-\pi/2 \leq \phi \leq \pi/2$, the Mueller matrix $\mathbf{EP}(\theta, \phi)$ is

$$\mathbf{EP}(\theta, \phi) = \frac{1}{2} \begin{pmatrix} 1 & \cos 2\theta \cos \phi & 2\cos \theta \sin \theta \cos \phi & \sin \phi \\ \cos 2\theta \cos \phi & \cos^2 2\theta \cos^2 \phi & \frac{1}{2} \sin 4\theta \cos^2 \phi & \frac{1}{2} \cos 2\theta \cos \phi \sin \phi \\ 2\cos \theta \sin \theta \cos \phi & \frac{1}{2} \sin 4\theta \cos^2 \phi & \sin^2 2\theta \cos^2 \phi & 2\cos \theta \sin \theta \cos \phi \sin \phi \\ \sin \phi & \frac{1}{2} \cos 2\theta \sin 2\phi & 2\cos \theta \sin \theta \cos \phi \sin \phi & \sin^2 \phi \end{pmatrix} \quad (20)$$

14.12 RETARDER MUELLER MATRICES

Retarders have two polarization states which are transmitted in the incident polarization state (eigenpolarizations) but with different optical path lengths (phases). Birefringent retarders divide incident light into two modes with orthogonal polarizations and delay one mode with respect to the other due to birefringence, the refractive index difference between the modes. Other retarding interactions include the following: reflections from metals, reflection and transmission through multilayer thin films, stress birefringence, and interactions with diffraction gratings. These interactions also are often diattenuating.

Retarders are specified by the optical path difference between the eigenpolarizations (the retardance δ) and the eigenpolarization states, either the state with the smaller optical path length (the fast axis) or the larger optical path (the slow axis). Retardance is specified in this chapter in radians, so $\delta = 2\pi$ indicates one wavelength of optical path difference. Note that *axis* implies a linear polarization state, but the fast eigenpolarization may be elliptical or circular and the term “axis” is still applied. The most common retarders in practice are quarter-wave linear retarders and half-wave linear retarders. Quarter-wave linear retarders are most commonly used to convert between linear and circularly polarized light. Half-wave linear retarders are most commonly used to rotate the plane of linear polarization.

In the Mueller calculus, retarders are represented by real unitary matrices of the form

$$\mathbf{M}_{\text{retarder}} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \text{3} \times \text{3 rotation} & & \\ 0 & \text{matrix} & & \\ 0 & & & \end{pmatrix} \quad (21)$$

where, except for the $M_{0,0}$ element, the first row and column are zero. Real unitary matrices are called *orthogonal matrices*. The definition of a unitary matrix \mathbf{U} is a matrix whose hermitian adjoint (complex conjugate of the matrix transpose) equals its matrix inverse,

$$\mathbf{U}^\dagger = (\mathbf{U}^T)^* = \mathbf{U}^{-1} \quad (22)$$

For a real matrix, the complex conjugate of a matrix equals the matrix, so the transpose of an orthogonal matrix \mathbf{O} equals its inverse,

$$\mathbf{O}^T = \mathbf{O}^{-1} \quad (23)$$

This equation *tests if a Mueller matrix is a pure retarder*. Orthogonal matrices such as retarder Mueller matrices are rotation matrices. The lower right 3×3 elements form a rotation matrix in $\{S_1, S_2, S_3\}$ space showing how retarders operate on Stokes vectors as a rotation of the Poincaré sphere. The retardance δ of a pure retarder Mueller matrix in radians is

$$\delta = \arccos\left(\frac{m_{0,0} + m_{1,1} + m_{2,2} + m_{3,3}}{2} - 1\right) = \arccos\left(\frac{\text{Tr}(\mathbf{M})}{2} - 1\right) \quad (24)$$

The Mueller matrices for quarter wave retarders with fast axes corresponding to the six basis polarization states are given in Table 2. Table 3 lists the half wave retarder Mueller matrices corresponding to the basis polarization states.

The Mueller matrix for a quarter-wave linear retarder with fast axis at angle θ , $\mathbf{QWLR}(\theta)$ is

$$\mathbf{QWLR}(\theta) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos^2 2\theta & \cos 2\theta \sin 2\theta & -\sin 2\theta \\ 0 & \cos 2\theta \sin 2\theta & \sin^2 2\theta & \cos 2\theta \\ 0 & \sin 2\theta & -\cos 2\theta & 0 \end{pmatrix} \quad (25)$$

TABLE 2 Quarter Wave Retarder Mueller Matrices for the Basis Polarization States

Type of Retarder	Symbol	Mueller Matrix
Horizontal quarter-wave Linear retarder	HQWLR	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & -1 & 0 \end{pmatrix}$
Vertical quarter-wave Linear retarder	VQWLR	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 0 \end{pmatrix}$
45° quarter-wave Linear retarder	QWLR(45°)	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$
135° quarter-wave Linear Retarder	QWLR(135°)	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & -1 & 0 & 0 \end{pmatrix}$
Quarter-wave right Circular retarder	QWRCR	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$
Quarter-wave left circular retarder	QWLCR	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$

TABLE 3 Half-Wave Retarder Mueller Matrices for the Basis Polarization States

Type of Retarder	Symbol	Mueller Matrix
Horizontal or vertical half-wave linear retarder (same matrix)	HHWLR	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix}$
45° or 135° half-wave linear retarder	HWLR(45°)	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix}$
Right or left half-wave circular retarder	RHWCR	$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$

Similarly a half wave linear retarder with fast axis at angle θ , **HWLR**(θ), has the matrix

$$\mathbf{HWLR}(\theta) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos 4\theta & \sin 4\theta & 0 \\ 0 & \sin 4\theta & -\cos 4\theta & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix} = \mathbf{HWLR}\left(\theta + \frac{\pi}{2}\right) \quad (26)$$

The Mueller matrix is the same for a horizontal half-wave linear retarder and a vertical half-wave linear retarder or for other pairs 90° apart because both half-wave retarders perform the same transformation on Stokes vectors.

The Mueller matrix **LR**(δ, θ) for a linear retarder with retardance δ and fast axis oriented at an angle θ is

$$\mathbf{LR}(\delta, \theta) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos^2(2\theta) + \cos(\delta) \sin^2(2\theta) & (1 - \cos(\delta)) \cos(2\theta) \sin(2\theta) & -\sin(\delta) \sin(2\theta) \\ 0 & (1 - \cos(\delta)) \cos(2\theta) \sin(2\theta) & \cos(\delta) \cos^2(2\theta) + \sin^2(2\theta) & \cos(2\theta) \sin(\delta) \\ 0 & \sin(\delta) \sin(2\theta) & -\cos(2\theta) \sin(\delta) & \cos(\delta) \end{pmatrix} \quad (27)$$

Elliptical retarders (**ER**) are commonly specified in two ways: (1) by specifying horizontal, 45° , and circular retardance components: $\delta_H, \delta_{45}, \delta_R$, or (2) by retardance, δ , orientation, θ , and latitude, ε , of the fast eigenstates using Poincaré sphere coordinates. In terms of retardance components, the magnitude of the retardance is

$$\delta = \sqrt{\delta_H^2 + \delta_{45}^2 + \delta_R^2} \quad \mathbf{S}_{\text{fast}} = \begin{pmatrix} 1 \\ \delta_H/\delta \\ \delta_{45}/\delta \\ \delta_R/\delta \end{pmatrix} \quad (28)$$

The ideal retarder Mueller matrix expressed in terms of retardance components is

$$\mathbf{ER}(\delta_H, \delta_{45}, \delta_R) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \frac{\delta_H^2 + (\delta_{45}^2 + \delta_R^2)C}{\delta^2} & \frac{\delta_{45} \delta_H T - \delta_R S}{\delta^2} & \frac{\delta_H \delta_R T + \delta_{45} S}{\delta^2} \\ 0 & \frac{\delta_{45} \delta_H T + \delta_R S}{\delta^2} & \frac{\delta_{45}^2 + (\delta_R^2 + \delta_H^2)C}{\delta^2} & \frac{\delta_R \delta_{45} T - \delta_H S}{\delta^2} \\ 0 & \frac{\delta_H \delta_R T - \delta_{45} S}{\delta^2} & \frac{\delta_R \delta_{45} T + \delta_H S}{\delta^2} & \frac{\delta_R^2 + (\delta_{45}^2 + \delta_H^2)C}{\delta^2} \end{pmatrix} \quad (29)$$

$$C = \cos \delta, S = \sin \delta, T = 1 - \cos \delta$$

In terms of the Poincaré sphere parameters: retardance, δ , orientation, θ , and latitude, ε , of the fast eigenstates, the elliptical retarder Mueller matrix has a long equation given by the following matrix product:

$$\mathbf{ER}(\delta, \theta, \varepsilon) = \mathbf{LR}(\varepsilon, \theta + 45^\circ) \cdot \mathbf{LR}(\delta, \theta) \cdot \mathbf{LR}(\varepsilon, \theta - 45^\circ) \quad (30)$$

Here is a list of the matrix elements,

$$\begin{pmatrix} m_{0,0} \\ m_{0,1} \\ m_{0,2} \\ m_{0,3} \\ m_{1,0} \\ m_{1,1} \\ m_{1,2} \\ m_{1,3} \\ m_{2,0} \\ m_{2,1} \\ m_{2,2} \\ m_{2,3} \\ m_{3,0} \\ m_{3,1} \\ m_{3,2} \\ m_{3,3} \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \frac{1}{8} \left(2 + 6\cos\delta - \cos(\delta - 2\varepsilon) + 2\cos 2\varepsilon - \cos(\delta + 2\varepsilon) + 8\cos^2(\varepsilon)\cos 4\theta \sin^2\left(\frac{\delta}{8}\right) \right) \\ \sin\delta \sin\varepsilon + \cos^2\varepsilon \sin^2\frac{\delta}{2} \sin 4\theta \\ \cos(\varepsilon) \left(2\cos 2\theta \sin^2\frac{\delta}{2} \sin\varepsilon - \sin\delta \sin 2\theta \right) \\ 0 \\ -\sin\delta \sin\varepsilon + \cos^2\varepsilon \sin^2\frac{\delta}{2} \sin 4\theta \\ \frac{1}{8} \left(2 + 6\cos\delta - \cos(\delta - 2\varepsilon) + 2\cos 2\varepsilon - \cos(\delta + 2\varepsilon) - 8\cos^2(\varepsilon)\cos 4\theta \sin^2\left(\frac{\delta}{8}\right) \right) \\ \cos\varepsilon \left(\cos 2\theta \sin\delta + 2\sin^2\frac{\delta}{2} \sin\varepsilon \sin 2\theta \right) \\ 0 \\ \cos\varepsilon \left(2\cos 2\theta \sin^2\frac{\delta}{2} \sin\varepsilon + \sin\delta \sin 2\theta \right) \\ \cos\varepsilon \left(-\cos 2\theta \sin\delta + 2\sin^2\frac{\delta}{2} \sin\varepsilon \sin 2\theta \right) \\ \cos^2\frac{\delta}{2} - \cos 2\varepsilon \sin^2\frac{\delta}{2} \end{pmatrix} \quad (31)$$

The Mueller matrix for half-wave elliptical and linear retarders **HWR** simplifies to the following form:

$$\mathbf{HWR}(\delta_H, \delta_{45}, \delta_R) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 + 2d_1^2 & 2d_2d_1 & 2d_1d_3 \\ 0 & 2d_2d_1 & -1 + 2d_2^2 & 2d_2d_3 \\ 0 & 2d_3d_1 & 2d_2d_3 & -1 + 2d_3^2 \end{pmatrix} \quad (32)$$

$$\sqrt{d_1^2 + d_2^2 + d_3^2} = 1 \quad d_1 = \frac{\delta_H}{\pi} \quad d_2 = \frac{\delta_{45}}{\pi} \quad d_3 = \frac{\delta_R}{\pi}$$

14.13 RETARDER MUELLER MATRICES AMBIGUITIES AND RETARDER SPACE

Retarders can be represented as points in a three-dimensional retarder space $\{\delta_H, \delta_{45}, \delta_R\}$ as in Fig. 2. In this space, all quarter-wave elliptical retarders lie on a sphere of radius $\pi/2$, all half-wave retarders on a sphere of radius π , and so on. All Mueller matrices on spheres of radius $2\pi n$, where n is an integer, the retarder order, have the identity matrix as their Mueller matrix, as does the point at the origin. The retarder space is similar to the Poincaré sphere except the retardance components are plotted instead of the Stokes vector, so the size of the space is not limited to a radius of one.

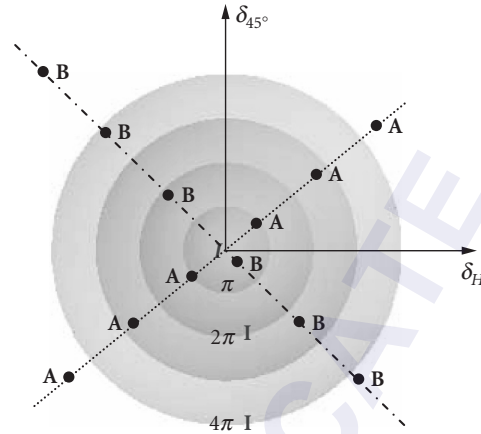


FIGURE 2 Retarder space is a Poincaré sphere like space which represents retarder fast axes $\{\delta_H, \delta_{45}, \delta_R\}$ as points in a three-dimensional space. All half-wave retarders at the points indicated by A's have the same Mueller matrix; all retarders at points B share the same Mueller matrix.

Retarder Mueller matrices have an ambiguity with regard to retarder order n , the integer part or half integer part of the number of waves of retardance. The Mueller matrix relates incident Stokes vectors to transmitted Stokes vectors. There is a family of retarders which will perform the same transformation on all polarization states. For example, a retarder with retardance $\delta = 0$ leaves all polarization states unchanged. Similarly when all polarization states exit a retarder with $\delta = 2\pi$, or $n2\pi$, all polarization states are returned to the incident polarization state. As another example, a quarter-wave retarder rotates the Poincaré sphere $\pi/2$ radians clockwise about an axis. A three-quarter-wave retarder with the orthogonal axis rotates the Poincaré sphere $3\pi/2$ radians counterclockwise, has the same Mueller matrix and transmits the same Stokes vectors. As a third example, half-wave retarders with orthogonal axes rotate the Poincaré sphere by half a rotation in opposite directions and have the same Mueller matrix. So in general all Mueller matrices with retardance $2\pi n + \delta$ and a particular normalized fast axis $\{\delta_H, \delta_{45}, \delta_R\}$ and all Mueller matrices with retardance $2\pi m - \delta$ and the orthogonal normalized fast axis $\{-\delta_H, -\delta_{45}, -\delta_R\}$ have the same Mueller matrix (m and n integers). This is shown in Fig. 2 by the set of half-wave retarder locations with A's. Another set of elliptical retarders with the same Mueller matrix are indicated by the set of B's.

The surfaces of various sets of retarders in the 16-dimensional space of Mueller matrices have an interesting topology due to the retarder ambiguity. In a narrow Möbius strip, the edge is nearly a circle which goes around twice before returning to its starting point. Similarly going around the equatorial plane of the half-wave retarders, the Mueller matrices for the corresponding half-wave linear retarders repeat twice. When plotted in the 16-dimensions of the Mueller matrix space, the half-wave retarders circle twice in a perfect circle for one 180° rotation of the fast axis. The linear retarder Mueller matrices for a retardance slightly different from π circle twice slightly offset and outline the edges of a Möbius strip. This doubling is a consequence of the fact that when all the linear retarders are plotted as points in the Mueller matrix space, they form a two-dimensional surface which is topologically equivalent to a Klein bottle, a single-sided surface with Möbius strip cross-sections. Similarly when all elliptical retarders are plotted in the Mueller matrix space, they form a higher-dimensional Klein bottle with a three-dimensional surface. To summarize, any plane through the origin of the retarder space in Fig. 2 maps to a Klein bottle in the Mueller matrix space, and the entire space maps to the higher-dimensional Klein bottle.

14.14 TRANSMITTANCE AND DIATTENUATION

Polarizers and partial polarizers are characterized by the property diattenuation, which describes the magnitude of the variation of the transmitted irradiance as a function of the incident polarization state. The diattenuation magnitude D , usually referred to as the *diattenuation*, is a function of the maximum, T_{\max} , and minimum, T_{\min} , transmittances of a polarization element,

$$D = \frac{T_{\max} - T_{\min}}{T_{\max} + T_{\min}} = \frac{\sqrt{m_{1,0}^2 + m_{2,0}^2 + m_{3,0}^2}}{m_{0,0}} \quad 0 \leq D \leq 1 \quad (33)$$

The diattenuation has the useful property that D varies from 1 for a polarizer to 0 for an element which transmits all polarization states equally, such as a retarder or a nonpolarizing interaction.

The transmitted irradiance of a Mueller matrix and its diattenuation depends only on the first row, $m_0 = \{m_{0,0}, m_{0,1}, m_{0,2}, m_{0,3}\}$, because these are the only elements which affect S'_0 . The diattenuation is not linear in T_{\min}/T_{\max} as shown in Fig. 3.

To find T_{\max} and T_{\min} , first the incident Stokes vector is normalized so $S_0 = 1$. The transmittance $T(\mathbf{S})$ of a device with Mueller matrix \mathbf{M} depends on the incident polarization state and is the ratio of the exiting flux to the incident flux,

$$T(\mathbf{M}, \hat{\mathbf{S}}) = \frac{S'_0}{S_0} = \frac{(\mathbf{M} \cdot \hat{\mathbf{S}})_0}{S_0} = \frac{m_{0,0}s_0 + m_{0,1}s_1 + m_{0,2}s_2 + m_{0,3}s_3}{s_0} \quad (34)$$

which depends on the dot product of the first row of the Mueller matrix with the incident Stokes vector. The dependence of the transmission on incident polarization state is characterized by a set of three diattenuation parameters, \mathbf{d} , defined as

$$\mathbf{d} = \frac{\{m_{0,1}, m_{0,2}, m_{0,3}\}}{m_{0,0}} = \{d_H, d_{45}, d_R\} \quad (35)$$

The diattenuation parameters have three components corresponding to the three components of the Stokes vector, x/y , $45^\circ/135^\circ$, right/left, each of which characterizes how the transmission varies with each of the Stokes vector component. The diattenuation parameter set \mathbf{d} is often called the

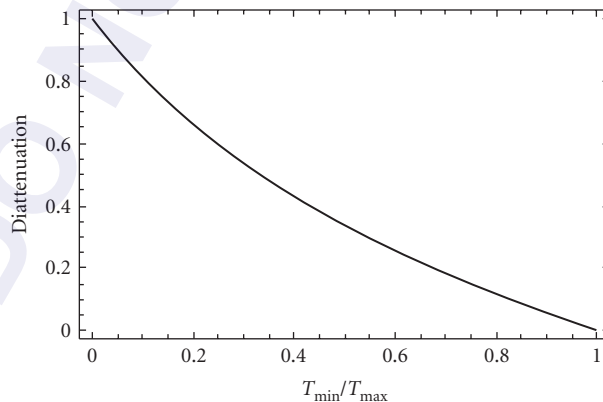


FIGURE 3 Relationship between the diattenuation and the extinction ratio, T_{\min}/T_{\max} .

diattenuation vector, but like the Stokes vector, \mathbf{d} is not a true vector. Diattenuation parameters do not add. For a Stokes three-vector \mathbf{s} , the transmission function T is

$$T(\mathbf{M}, \mathbf{s}) = \frac{S'_0}{S_0} = m_{0,0} + m_{0,1} s_1 + m_{0,2} s_2 + m_{0,3} s_3 = m_{0,0}(1 + \mathbf{d} \cdot \mathbf{s}) \quad (36)$$

The average transmission, formed by averaged over all polarized Stokes vectors, is $m_{0,0}$. The average transmission is also the transmission for unpolarized incident light, $\mathbf{s}_U = \{0,0,0\}$. The polarization-dependent variation of the transmission is contained in the dot product term between the incident Stokes three-vector and the diattenuation vector, $\mathbf{s} \cdot \mathbf{d}$. The maximum transmission, T_{\max} , occurs when the dot product is maximized, which occurs when \mathbf{s} and \mathbf{d} are parallel, and the magnitude of S'_0 is as large as possible. The incident normalized Stokes vectors with maximum transmittance, S_{\max} , and minimum transmittance, S_{\min} , are

$$\mathbf{s}_{\max} = \frac{\mathbf{d}}{|\mathbf{d}|} \quad \hat{\mathbf{S}}_{\max} = \frac{1}{D} \begin{pmatrix} D \\ d_x \\ d_{45} \\ d_R \end{pmatrix} \quad \mathbf{s}_{\min} = \frac{-\mathbf{d}}{|\mathbf{d}|} \quad \hat{\mathbf{S}}_{\min} = \frac{1}{D} \begin{pmatrix} D \\ -d_x \\ -d_{45} \\ -d_R \end{pmatrix} \quad (37)$$

yielding

$$T_{\max} = m_{0,0}(1+D) \quad T_{\min} = m_{0,0}(1-D) \quad (38)$$

Therefore the diattenuation of any Mueller matrix is

$$D(\mathbf{M}) = \frac{T_{\max} - T_{\min}}{T_{\max} + T_{\min}} = \frac{\sqrt{m_{0,1}^2 + m_{0,2}^2 + m_{0,3}^2}}{m_{0,0}} \quad (39)$$

For an ideal polarizer the minimum transmission is zero, $D=1$, $T_{\min} = m_{0,0}(1-D) = 0$.

Linear polarization sensitivity or *linear diattenuation* $LD(\mathbf{M})$ characterizes the variation of intensity transmittance with incident linear polarization states:

$$LD(\mathbf{M}) = \frac{\sqrt{m_{0,1}^2 + m_{0,2}^2}}{m_{0,0}} \quad (40)$$

Linear polarization sensitivity is frequently specified as a performance parameter in remote sensing systems designed to measure incident power independently of any linearly polarized component present in scattered earth-light.¹⁸ Note that $LD(\mathbf{M}) = 1$ identifies \mathbf{M} as a linear analyzer; \mathbf{M} is not necessarily a linear polarizer, but may represent a linear polarizer followed by some other polarization element.

Diattenuation in fiber optic components and systems is often characterized by the *polarization-dependent loss* (*PDL*) specified in decibels:

$$PDL(\mathbf{M}) = 10 \log_{10} \frac{T_{\max}}{T_{\min}} \quad (41)$$

14.15 POLARIZANCE

The polarizance ($P(\mathbf{M})$) is the degree of polarization (DoP) of the transmitted light when unpolarized light $\hat{\mathbf{U}}$ is incident,³

$$P(\mathbf{M}) = \text{DoP}(\mathbf{M} \cdot \hat{\mathbf{U}}) = \frac{\sqrt{m_{1,0}^2 + m_{2,0}^2 + m_{3,0}^2}}{m_{0,0}} \quad (42)$$

The exiting polarization state, $\mathbf{S}_p(\mathbf{M})$, is the first column of \mathbf{M} ,

$$\mathbf{S}_p(\mathbf{M}) = \mathbf{M} \cdot \mathbf{U} = \begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} m_{0,0} \\ m_{1,0} \\ m_{2,0} \\ m_{3,0} \end{pmatrix} \quad (43)$$

The polarizance does not necessarily equal the diattenuation. Nor does \mathbf{S}_p necessarily equal, $\bar{\mathbf{S}}_{\max}$, the incident state of maximum transmittance.

14.16 MUELLER MATRICES OF DIATTENUATORS

The Mueller matrix for a partial polarizer (homogeneous diattenuator) with intensity transmittances T_x and T_y and eigenpolarizations along the x and y axes, $\mathbf{LD}(T_x, T_y, 0)$ is

$$\mathbf{LD}(T_x, T_y, 0) = \frac{1}{2} \begin{pmatrix} T_x + T_y & T_x - T_y & 0 & 0 \\ T_x - T_y & T_x + T_y & 0 & 0 \\ 0 & 0 & 2\sqrt{T_x T_y} & 0 \\ 0 & 0 & 0 & 2\sqrt{T_x T_y} \end{pmatrix} \quad (44)$$

Ideal diattenuators have two different intensity transmittances T_{\max} and T_{\min} , for two orthogonal linear eigenpolarizations; thus the name “di” “attenuator”. A linear diattenuator oriented at angle θ has the Mueller matrix

$$\mathbf{LD}(T_{\max}, T_{\min}, \theta) = \frac{1}{2} \begin{pmatrix} A & B \cos 2\theta & B \sin 2\theta & 0 \\ B \cos 2\theta & A \cos^2 2\theta + C \sin^2 2\theta & (A - C) \cos 2\theta \sin 2\theta & 0 \\ B \sin 2\theta & (A - C) \cos 2\theta \sin 2\theta & C \cos^2 2\theta + A \sin^2 2\theta & 0 \\ 0 & 0 & 0 & C \end{pmatrix} \quad (45)$$

where

$$A = T_{\max} + T_{\min} \quad B = T_{\max} - T_{\min} \quad C = 2\sqrt{T_{\max} T_{\min}} \quad (46)$$

Ideal diattenuators have no retardance, although in practice most diattenuators have some retardance. An example of a pure linear diattenuator without retardance is transmission into a transparent dielectric; T_{\max} and T_{\min} are then given by intensity Fresnel coefficients. Reflection at metal surfaces acts as a diattenuator with retardance.

Ideal diattenuator Mueller matrices are hermitian matrices; they have real eigenvalues. A hermitian matrix equals the complex conjugate of its transpose, its hermitian adjoint, $\mathbf{H} = \mathbf{H}^\dagger = (\mathbf{H}^T)^*$. But since Mueller matrices are real, $\mathbf{H}^* = \mathbf{H}$, ideal diattenuator Mueller matrices equal their transpose,

$$\mathbf{H} = \mathbf{H}^T \quad (47)$$

The general equation for a diattenuator, either linear, elliptical, or circular, expressed in terms of the first row of the Mueller matrix is

$$\begin{aligned} \text{Diattenuator}(d_H, d_{45}, d_R, T_{\text{avg}}) = \\ T_{\text{avg}} \begin{pmatrix} 1 & d_H & d_{45} & d_R \\ d_H & A & 0 & 0 \\ d_{45} & 0 & A & 0 \\ d_R & 0 & 0 & A \end{pmatrix} + \frac{T_{\text{avg}}(1-A)}{D^2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & d_H^2 & d_{45}d_H & d_Hd_R \\ 0 & d_{45}d_H & d_{45}^2 & d_{45}d_R \\ 0 & d_Hd_R & d_{45}d_R & d_R^2 \end{pmatrix} \quad (48) \\ D = \sqrt{d_H^2 + d_{45}^2 + d_R^2} \quad A = \sqrt{1 - d_H^2 - d_{45}^2 - d_R^2} \quad T_{\text{avg}} = \frac{T_{\text{max}} + T_{\text{min}}}{2} \end{aligned}$$

14.17 NORMALIZING A MUELLER MATRIX

Mueller matrices are normalized by dividing \mathbf{M} by $m_{0,0}$, the transmission for unpolarized light, and also the transmission when the input state is averaged over the entire Poincaré sphere. The normalized Mueller matrix $\bar{\mathbf{M}}$ has an average transmission of one,

$$\bar{\mathbf{M}} = \frac{\mathbf{M}}{m_{0,0}} = \begin{pmatrix} 1 & m_{0,1}/m_{0,0} & m_{0,2}/m_{0,0} & m_{0,3}/m_{0,0} \\ m_{1,0}/m_{0,0} & m_{1,1}/m_{0,0} & m_{1,2}/m_{0,0} & m_{1,3}/m_{0,0} \\ m_{2,0}/m_{0,0} & m_{2,1}/m_{0,0} & m_{2,2}/m_{0,0} & m_{2,3}/m_{0,0} \\ m_{3,0}/m_{0,0} & m_{3,1}/m_{0,0} & m_{3,2}/m_{0,0} & m_{3,3}/m_{0,0} \end{pmatrix} \quad (49)$$

Normalization limits all element values to the range $-1 \leq m_{i,j} \leq 1$. Measured Mueller matrix data is frequently normalized to facilitate comparison of Mueller matrix polarization properties with the average transmission removed. It facilitates comparison to Mueller matrices tabulated in the literature. Normalizing Mueller matrix images or spectra simplifies *eyeballing* the data; with the flux variations removed, diattenuation, retardance, and depolarization variations are easier to see.

Normalizing by a different value, $k/m_{0,0}$, sets the average transmission to k , such as $k = 1/2$ for an ideal polarizer. A Mueller matrix $\hat{\mathbf{M}}$ normalized, so the maximum transmission is one is

$$\hat{\mathbf{M}} = \frac{\mathbf{M}}{m_{0,0} + \sqrt{m_{0,1}^2 + m_{0,2}^2 + m_{0,3}^2}} = \frac{\mathbf{M}}{T_{\text{max}}} \quad (50)$$

14.18 COORDINATE SYSTEM FOR THE MUELLER MATRIX

Consider a Mueller polarimeter consisting of a polarization generator which illuminates a sample, and a polarization analyzer which collects the light exiting the sample in a particular direction. We wish to characterize the polarization modification properties of the sample for a particular incident and exiting beam through the Mueller matrix. The incident polarization states are specified by Stokes vectors defined relative to an $\{\hat{\mathbf{x}}, \hat{\mathbf{y}}\}$ coordinate system orthogonal to the propagation

direction of the incident light. Similarly, the exiting light's Stokes vector is defined relative to an $\{\hat{x}', \hat{y}'\}$ coordinate system orthogonal to its propagation direction. For transmission measurements where the beam exits undeviated, the orientations of $\{\hat{x}, \hat{y}\}$ and $\{\hat{x}', \hat{y}'\}$ will naturally be chosen to be aligned, ($\hat{x}=\hat{x}', \hat{y}=\hat{y}'$). The global orientation of $\{\hat{x}, \hat{y}\}$ is arbitrary, and the measured Mueller matrix varies systematically if $\{\hat{x}, \hat{y}\}$ and $\{\hat{x}', \hat{y}'\}$ are rotated together.

When the exiting beam emerges in a different direction from the incident beam, orientations must be specified for both sets of coordinates. For measurements of reflection from a surface, a logical choice sets $\{\hat{x}, \hat{y}\}$ and $\{\hat{x}', \hat{y}'\}$ to the $\{\hat{s}, \hat{p}\}$ orientations for the two beams. Other Mueller matrix measurement configurations may have other obvious arrangements for the coordinates. All choices, however, are arbitrary, and lead to different Mueller matrices. Let a Mueller matrix \mathbf{M} be defined relative to a particular $\{\hat{x}, \hat{y}\}$ and $\{\hat{x}', \hat{y}'\}$. Let another Mueller matrix $\mathbf{M}(\theta_1, \theta_2)$ for the same measurement conditions have its \hat{x} axis rotated by θ_1 and x' axis rotated by θ_2 , where $\theta > 0$ indicates a counterclockwise rotation looking into the beam (\hat{x} into \hat{y}). These Mueller matrices are related by the equation

$$\mathbf{M}(\theta_1, \theta_2) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos 2\theta_2 & -\sin 2\theta_2 & 0 \\ 0 & \sin 2\theta_2 & \cos 2\theta_2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos 2\theta_1 & \sin 2\theta_1 & 0 \\ 0 & -\sin 2\theta_1 & \cos 2\theta_1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (51)$$

When $\theta_1 = \theta_2$, the coordinates rotate together, the eigenvalues are preserved, the circular polarization properties are preserved, and the linear properties are shifted in orientation. When $\theta_1 \neq \theta_2$, the matrix properties are qualitatively different; the eigenvalues of the matrix change. If the eigenpolarizations of \mathbf{M} were orthogonal, they may not remain orthogonal. After we perform data reduction on the matrix, the basic polarization properties couple in a complex fashion. For example, linear diattenuation in \mathbf{M} yields a circular retardance component in $\mathbf{M}(\theta_1, \theta_2)$. The selection of the coordinate systems for the incident and exiting beams is not important for describing exiting polarization states, but is crucial for properly identifying polarization characteristics of the sample.

14.19 MUELLER MATRICES FOR REFRACTION

Reflections and refractions at homogenous and isotropic interfaces, typical glass or metal interfaces, have s and p eigenpolarizations. The polarization is a combination of diattenuation and retardance, with the eigenpolarizations aligned with the s and p planes. Let s be aligned with x and p with y . T_s is the s -intensity reflectance or transmittance and T_p is the p -intensity reflectance or transmittance. The retardance between the s and p states is δ . T_s , T_p , δ are determined from Fresnel equations or from a thin-film coating calculation. The Mueller matrix is the product of the diattenuator and retarder Mueller matrices,

$$\begin{aligned} \text{LDR}(D, \delta, 0) &= \frac{T_{\max} + T_{\min}}{2} \begin{pmatrix} 1 & D & 0 & 0 \\ D & 1 & 0 & 0 \\ 0 & 0 & \sqrt{1-D^2} \cos \delta & \sqrt{1-D^2} \sin \delta \\ 0 & 0 & -\sqrt{1-D^2} \sin \delta & \sqrt{1-D^2} \cos \delta \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} T_s + T_p & T_s - T_p & 0 & 0 \\ T_s - T_p & T_s + T_p & 0 & 0 \\ 0 & 0 & 2\sqrt{T_s T_p} \cos \delta & 2\sqrt{T_s T_p} \sin \delta \\ 0 & 0 & -2\sqrt{T_s T_p} \sin \delta & 2\sqrt{T_s T_p} \cos \delta \end{pmatrix} \end{aligned} \quad (52)$$

For transmission at an uncoated interface, $\delta = 0$, but for thin-film-coated interfaces, such as anti-reflection coatings or beam-splitter coatings, the retardance is nonzero. If the plane of incidence is not vertical, the Mueller rotation operator is applied. Refraction Mueller matrices are homogeneous; the eigenpolarizations, the s and p polarizations, are orthogonal.

14.20 MUELLER MATRICES FOR REFLECTION

In most optics notation, including this chapter, a sign change occurs in the coordinate system after reflection to maintain a right-handed coordinate system after the propagation vector has changed direction. After reflection, the S_2 component of Stokes vectors (linearly polarized light at $45^\circ/135^\circ$) and the S_3 component (circularly polarized light) change sign. The S_2 component changes sign during reflection (diffuse or specular) because the z -component of the light propagation vector (the component parallel to the sample surface normal) changes sign. To maintain a right-handed coordinate system, one of the transverse coordinates must change sign as well. Choosing x , the spatial coordinates (x, y, z) switch to $(-x, y, -z)$ after reflection or scatter from a sample; z is the direction of propagation before reflection which changes to $-z$ after reflection. The change of coordinates dictates that a beam polarized at an angle of 45° which reflects polarized in the same global plane is described as having a 135° orientation in the coordinates following reflection.

Picture a Stokes polarimeter measuring in transmission; now rotate that polarimeter around z axis and move it to measure in reflection and you should see how the 45° component has changed sign. In addition, the helicity (i.e., handedness) of all circular and elliptical states changes sign upon reflection. Right circular polarization reflects as left circular polarization, and vice versa.

First let s be aligned with x and p with y . R_s is the s -intensity reflectance and R_p is the p -intensity reflectance. The retardance between the s and p states is δ . R_s, R_p, δ are determined from Fresnel equations or from a thin-film coating calculation. The reflection Mueller matrix is

$$\frac{1}{2} \begin{pmatrix} R_s + R_p & R_s - R_p & 0 & 0 \\ R_s - R_p & R_s + R_p & 0 & 0 \\ 0 & 0 & -2\sqrt{R_s R_p} \cos \delta & -2\sqrt{R_s R_p} \sin \delta \\ 0 & 0 & 2\sqrt{R_s R_p} \sin \delta & -2\sqrt{R_s R_p} \cos \delta \end{pmatrix} \quad (53)$$

With this convention for reflection, the equation for rotating the Mueller matrix, \mathbf{M} , of a sample measured by a polarimeter in a reflection configuration about its normal is

$$\mathbf{M}_R(\theta) = \mathbf{R}(\theta) \cdot \mathbf{M} \cdot \mathbf{R}(\theta)$$

$$= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos 2\theta & -\sin 2\theta & 0 \\ 0 & \sin 2\theta & \cos 2\theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos 2\theta & -\sin 2\theta & 0 \\ 0 & \sin 2\theta & \cos 2\theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (54)$$

compared to

$$\mathbf{M}_T(\theta) = \mathbf{R}(\theta) \cdot \mathbf{M} \cdot \mathbf{R}(-\theta) \quad (55)$$

for Mueller matrices in transmission. For example, the Mueller matrix of a transmission polarizer with its transmission axis oriented at 20° and the Mueller matrix of a reflection polarizer oriented at 20° are different since polarized light exits the reflection polarizer oriented at -20° in the reflection coordinates (20° in the incident coordinates). In essence the reflection polarizer Mueller matrix is

analyzing at 20° but polarizing at -20° . For the special cases of linear polarizer matrices oriented at 0° or 90° and linear retarders oriented at 0° or 90° , this transformation results in Mueller matrices which are the same for transmission and reflection.

The normalized reflection Mueller matrices for weakly polarizing diffuse reflecting samples, those with diattenuation, retardance, and depolarization close to zero, are close to the Mueller matrix for an ideal reflector,

$$\mathbf{M}_{\text{refl}} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix} \quad (56)$$

14.21 CONVERSION BETWEEN MUELLER MATRICES AND JONES MATRICES

Jones matrices form an alternative and very useful representation of sample polarization, particularly because Jones matrices have simpler properties and are more easily manipulated and interpreted. The complication in mapping Mueller matrices onto Jones matrices and vice versa is that Mueller matrices cannot represent absolute phase and Jones matrices cannot represent depolarization. Thus, only nondepolarizing Mueller matrices have corresponding Jones matrices. All Jones matrices have corresponding Mueller matrices, but because the absolute phase is not represented, the mapping is many Jones matrices to one Mueller matrix.

Both Jones matrices and Mueller matrices can calculate the polarization properties of sequences of nondepolarizing interactions, the effect of cascading a series of diattenuators and retarders. When this same polarization element sequence is calculated by Jones matrices and alternatively by Mueller matrices, the answer is the same diattenuating and retarding properties. Either method is suitable.

Jones vectors and Jones matrices are commonly represented with two different sign conventions for the phase. Electromagnetic waves are commonly written with two different conventions, either the phase decreases with time ($kz - \omega t - \phi$), the convention adopted here, or the phase increases with time ($\omega t - kz + \phi$). Depending on the choice, various plus and minus signs must be adjusted in the Jones vectors for circularly and elliptically polarized light and in the various Jones matrices. Both conventions are in widespread use so care is necessary when taking Jones matrices from different sources. In this chapter, the phase decreases with time, so a monochromatic plane wave propagating in the z direction has the form

$$\mathbf{E}(z, t) = \text{Re} \left\{ \begin{pmatrix} E_x \\ E_y \end{pmatrix} e^{i(kz - \omega t - \phi)} \right\} = \begin{pmatrix} E_x \\ E_y \end{pmatrix} \cos(kz - \omega t - \phi) \quad (57)$$

A wave is *advanced* by *subtracting* from the phase. A wave is delayed or retarded by adding to the phase.

A Jones matrix \mathbf{J} is transformed into a Mueller matrix by the relation

$$\mathbf{M} = \mathbf{U}(\mathbf{J} \otimes \mathbf{J}^*) \mathbf{U}^{-1} \quad (58)$$

in which \otimes represents the tensor product and \mathbf{U} is the Jones/Mueller transformation matrix:^{9,10,19,20}

$$\mathbf{U} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & -1 \\ 0 & 1 & 1 & 0 \\ 0 & i & -i & 0 \end{pmatrix} = (\mathbf{U}^{-1})^\dagger \quad (59)$$

where the hermitian adjoint is represented by †. All Jones matrices of the form $\mathbf{J}' = e^{i\phi}\mathbf{J}$ transform to the same Mueller matrix. Consider a Jones matrix with complex elements expressed in polar coordinate form

$$\mathbf{J} = \begin{pmatrix} \phi_{x,x} & \phi_{x,y} \\ \phi_{y,x} & \phi_{y,y} \end{pmatrix} = \begin{pmatrix} \rho_{x,x}e^{i\phi_{x,x}} & \rho_{x,y}e^{i\phi_{x,y}} \\ \rho_{y,x}e^{i\phi_{y,x}} & \rho_{y,y}e^{i\phi_{y,y}} \end{pmatrix} \quad (60)$$

The tensor product $(\mathbf{J} \otimes \mathbf{J}^*)$ gives a fourth rank second-order tensor $\{2 \times 2 \times 2 \times 2\}$,

$$(\mathbf{J} \otimes \mathbf{J}^*) = \begin{pmatrix} \rho_{x,x}e^{i\phi_{x,x}} \begin{pmatrix} \rho_{x,x}e^{-i\phi_{x,x}} & \rho_{x,y}e^{-i\phi_{x,y}} \\ \rho_{y,x}e^{-i\phi_{y,x}} & \rho_{y,y}e^{-i\phi_{y,y}} \end{pmatrix} & \rho_{x,y}e^{i\phi_{x,y}} \begin{pmatrix} \rho_{x,x}e^{-i\phi_{x,x}} & \rho_{x,y}e^{-i\phi_{x,y}} \\ \rho_{y,x}e^{-i\phi_{y,x}} & \rho_{y,y}e^{-i\phi_{y,y}} \end{pmatrix} \\ \rho_{y,x}e^{i\phi_{y,x}} \begin{pmatrix} \rho_{x,x}e^{-i\phi_{x,x}} & \rho_{x,y}e^{-i\phi_{x,y}} \\ \rho_{y,x}e^{-i\phi_{y,x}} & \rho_{y,y}e^{-i\phi_{y,y}} \end{pmatrix} & \rho_{y,y}e^{i\phi_{y,y}} \begin{pmatrix} \rho_{x,x}e^{-i\phi_{x,x}} & \rho_{x,y}e^{-i\phi_{x,y}} \\ \rho_{y,x}e^{-i\phi_{y,x}} & \rho_{y,y}e^{-i\phi_{y,y}} \end{pmatrix} \end{pmatrix} \quad (61)$$

This tensor is contracted to a second rank fourth-order (4×4) tensor.

$$(\mathbf{J} \otimes \mathbf{J}^*) = \begin{pmatrix} \rho_{x,x}^2 & \rho_{x,x}\rho_{x,y}e^{i(\phi_{x,x}-\phi_{x,y})} & \rho_{x,y}\rho_{x,x}e^{i(\phi_{x,y}-\phi_{x,x})} & \rho_{x,y}^2 \\ \rho_{x,x}\rho_{y,x}e^{i(\phi_{x,x}-\phi_{y,x})} & \rho_{x,x}\rho_{y,y}e^{i(\phi_{x,x}-\phi_{y,y})} & \rho_{x,y}\rho_{y,x}e^{i(\phi_{x,y}-\phi_{y,x})} & \rho_{x,y}\rho_{y,y}e^{i(\phi_{x,y}-\phi_{y,y})} \\ \rho_{y,x}\rho_{x,x}e^{i(\phi_{y,x}-\phi_{x,x})} & \rho_{y,x}\rho_{x,y}e^{i(\phi_{y,x}-\phi_{x,y})} & \rho_{y,y}\rho_{x,x}e^{i(\phi_{y,y}-\phi_{x,x})} & \rho_{y,y}\rho_{x,y}e^{i(\phi_{y,y}-\phi_{x,y})} \\ \rho_{y,x}^2 & \rho_{y,x}\rho_{y,y}e^{i(\phi_{y,x}-\phi_{y,y})} & \rho_{y,y}\rho_{y,x}e^{i(\phi_{y,y}-\phi_{y,x})} & \rho_{y,y}^2 \end{pmatrix} \quad (62)$$

which when multiplied with \mathbf{U} gives the Mueller matrix elements as follows:

$$\mathbf{U} \cdot (\mathbf{J} \otimes \mathbf{J}^*) \cdot \mathbf{U}^{-1} = \begin{pmatrix} m_{0,0} \\ m_{0,1} \\ m_{0,2} \\ m_{0,3} \\ m_{1,0} \\ m_{1,1} \\ m_{1,2} \\ m_{1,3} \\ m_{2,0} \\ m_{2,1} \\ m_{2,2} \\ m_{2,3} \\ m_{3,0} \\ m_{3,1} \\ m_{3,2} \\ m_{3,3} \end{pmatrix} = \begin{pmatrix} \frac{1}{2}(\rho_{x,x}^2 + \rho_{x,y}^2 + \rho_{y,x}^2 + \rho_{y,y}^2) \\ \frac{1}{2}(\rho_{x,x}^2 - \rho_{x,y}^2 + \rho_{y,x}^2 - \rho_{y,y}^2) \\ \rho_{x,x}\rho_{x,y} \cos(\phi_{x,x} - \phi_{x,y}) + \rho_{y,x}\rho_{y,y} \cos(\phi_{y,x} - \phi_{y,y}) \\ \rho_{x,x}\rho_{x,y} \sin(\phi_{x,x} - \phi_{x,y}) + \rho_{y,x}\rho_{y,y} \sin(\phi_{y,x} - \phi_{y,y}) \\ \frac{1}{2}(\rho_{x,x}^2 + \rho_{x,y}^2 - \rho_{y,x}^2 - \rho_{y,y}^2) \\ \frac{1}{2}(\rho_{x,x}^2 - \rho_{x,y}^2 - \rho_{y,x}^2 + \rho_{y,y}^2) \\ \rho_{x,x}\rho_{x,y} \cos(\phi_{x,x} - \phi_{x,y}) - \rho_{y,x}\rho_{y,y} \cos(\phi_{y,x} - \phi_{y,y}) \\ \rho_{x,x}\rho_{x,y} \sin(\phi_{x,x} - \phi_{x,y}) - \rho_{y,x}\rho_{y,y} \sin(\phi_{y,x} - \phi_{y,y}) \\ \rho_{x,x}\rho_{y,x} \cos(\phi_{x,x} - \phi_{y,x}) + \rho_{x,y}\rho_{y,y} \cos(\phi_{x,y} - \phi_{y,y}) \\ \rho_{x,x}\rho_{y,x} \cos(\phi_{x,x} - \phi_{y,x}) - \rho_{x,y}\rho_{y,y} \cos(\phi_{x,y} - \phi_{y,y}) \\ \rho_{x,y}\rho_{y,x} \cos(\phi_{x,y} - \phi_{y,x}) + \rho_{x,x}\rho_{y,y} \cos(\phi_{x,x} - \phi_{y,y}) \\ -\rho_{x,y}\rho_{y,x} \sin(\phi_{x,y} - \phi_{y,x}) + \rho_{x,x}\rho_{y,y} \sin(\phi_{x,x} - \phi_{y,y}) \\ -\rho_{x,x}\rho_{y,x} \sin(\phi_{x,x} - \phi_{y,x}) - \rho_{x,y}\rho_{y,y} \sin(\phi_{x,y} - \phi_{y,y}) \\ -\rho_{x,x}\rho_{y,x} \sin(\phi_{x,x} - \phi_{y,x}) + \rho_{x,y}\rho_{y,y} \sin(\phi_{x,y} - \phi_{y,y}) \\ -\rho_{x,y}\rho_{y,x} \sin(\phi_{x,y} - \phi_{y,x}) - \rho_{x,x}\rho_{y,y} \sin(\phi_{x,x} - \phi_{y,y}) \\ -\rho_{x,y}\rho_{y,x} \cos(\phi_{x,y} - \phi_{y,x}) + \rho_{x,x}\rho_{y,y} \cos(\phi_{x,x} - \phi_{y,y}) \end{pmatrix} \quad (63)$$

An equivalent method to convert Jones matrices to Mueller matrices utilizes dot products with two Pauli spin matrices to determine each Mueller matrix element, $m_{i,p}$

$$m_{i,j} = \frac{1}{2} \text{Tr}(\mathbf{J} \cdot \boldsymbol{\sigma}_i \cdot \mathbf{J}^\dagger \cdot \boldsymbol{\sigma}_j) \quad (64)$$

where Tr is the trace of the matrix and

$$\boldsymbol{\sigma}_i = \left\{ \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} \right\} \quad i=0, 1, 2, 3 \quad (65)$$

are the identity matrix and normalized Pauli spin matrices.

Nondepolarizing Mueller matrices are transformed into the equivalent Jones matrices using the following relations:

$$\mathbf{J} = \begin{pmatrix} j_{x,x} & j_{x,y} \\ j_{y,x} & j_{y,y} \end{pmatrix} = \begin{pmatrix} \rho_{x,x} e^{i\phi_{x,x}} & \rho_{x,y} e^{i\phi_{x,y}} \\ \rho_{y,x} e^{i\phi_{y,x}} & \rho_{y,y} e^{i\phi_{y,y}} \end{pmatrix} \quad (66)$$

where the amplitudes are

$$\begin{aligned} \rho_{x,x} &= \frac{1}{\sqrt{2}} \sqrt{m_{0,0} + m_{0,1} + m_{1,0} + m_{1,1}} & \rho_{x,y} &= \frac{1}{\sqrt{2}} \sqrt{m_{0,0} - m_{0,1} + m_{1,0} - m_{1,1}} \\ \rho_{y,x} &= \frac{1}{\sqrt{2}} \sqrt{m_{0,0} + m_{0,1} - m_{1,0} - m_{1,1}} & \rho_{y,y} &= \frac{1}{\sqrt{2}} \sqrt{m_{0,0} - m_{0,1} - m_{1,0} + m_{1,1}} \end{aligned} \quad (67)$$

and the relative phases are

$$\begin{aligned} \phi_{x,x} - \phi_{x,y} &= \arctan \left(\frac{m_{0,3} + m_{1,3}}{m_{0,2} + m_{1,2}} \right) & \phi_{y,x} - \phi_{x,x} &= \arctan \left(\frac{m_{3,0} + m_{3,1}}{m_{2,0} + m_{2,1}} \right) \\ \phi_{y,y} - \phi_{x,x} &= \arctan \left(\frac{m_{3,2} - m_{2,3}}{m_{2,2} + m_{3,3}} \right) \end{aligned} \quad (68)$$

The phase $\phi_{x,x}$ is not determined but is the “reference phase.” These equations are not unique and other equivalent forms can be derived.

For example, a special case occurs if $j_{x,x} = 0$; then both the numerator and denominator of the arctan are zero and the phase equations fail. The transformation equations can then be recast in closely related forms to use the phase of another Jones matrix element as the “reference phase.”

14.22 NONDEPOLARIZING MUELLER MATRICES AND MUELLER-JONES MATRICES

Nondepolarizing Mueller matrices are the set of Mueller matrices for which completely polarized incident light [DoP(S)=1] is transmitted as completely polarized for all incident polarization states and have a depolarization index of one [Eq. (96)]. Nondepolarizing Mueller matrices are a subset of the Mueller matrices. Jones matrices can only represent nondepolarizing interactions. The nondepolarizing Mueller matrices are Mueller matrices with corresponding Jones matrices; thus nondepolarizing Mueller matrices are also called Mueller-Jones matrices.

An ideal polarizer is nondepolarizing when if the incident beam is polarized, the exiting beam is also polarized. Similarly an ideal retarder is nondepolarizing. The nondepolarizing Mueller matrices comprise the Mueller matrices for the matrix product of all arbitrary sequences of diattenuation and retardance. A Mueller-Jones matrix must satisfy the following condition for all θ and ϕ ,

$$\text{DoP}(\mathbf{M} \cdot \mathbf{S}) = \text{DoP} \left(\begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \cdot \begin{pmatrix} 1 \\ \cos 2\theta \cos \phi \\ \sin 2\theta \cos \phi \\ \sin \phi \end{pmatrix} \right) = 1 \quad (69)$$

One necessary, but not sufficient, condition for nondepolarizing Mueller matrices is²¹

$$\begin{aligned} \text{Tr}(\mathbf{M} \cdot \mathbf{M}^T) = 4 m_{0,0}^2 = & m_{0,0}^2 + m_{0,1}^2 + m_{0,2}^2 + m_{0,3}^2 + m_{1,0}^2 + m_{1,1}^2 + m_{1,2}^2 + m_{1,3}^2 \\ & + m_{2,0}^2 + m_{2,1}^2 + m_{2,2}^2 + m_{2,3}^2 + m_{3,0}^2 + m_{3,1}^2 + m_{3,2}^2 + m_{3,3}^2 \end{aligned} \quad (70)$$

In a typical imaging optical system, depolarization is an undesirable characteristic for lens and mirror surfaces, filters, and polarization elements. Depolarization is associated with scattering, and optical surfaces are carefully fabricated and coated to minimize scattering. Depolarization is generally very small from high-quality optical surfaces. Thus the majority of optical surfaces are well described by nondepolarizing Mueller matrices.

14.23 HOMOGENEOUS AND INHOMOGENEOUS POLARIZATION ELEMENTS

A nondepolarizing Mueller matrix is *homogeneous* if the two Stokes vector eigenpolarizations are orthogonal, and *inhomogeneous* otherwise. A nondepolarizing Mueller matrix \mathbf{M}_N can be factored into a cascade of a diattenuator Mueller matrix \mathbf{M}_D followed by a retarder Mueller matrix \mathbf{M}_R or into a cascade of the same retarder followed by a different diattenuator \mathbf{M}'_D .²²

$$\mathbf{M}_N = \mathbf{M}_R \mathbf{M}_D = \mathbf{M}'_D \mathbf{M}_R \quad (71)$$

The magnitude of the diattenuation of \mathbf{M}_D and \mathbf{M}'_D are equal. We define the diattenuation of \mathbf{M} as the diattenuation of \mathbf{M}_D , and the retardance of \mathbf{M} as the retardance of \mathbf{M}_R . For a homogeneous device, $\mathbf{M}_D = \mathbf{M}'_D$ and the eigenvectors of \mathbf{M}_R and \mathbf{M}_D are equal. Thus the retardance and diattenuation of a homogeneous Mueller matrix are “aligned,” giving it substantially simpler properties than the inhomogeneous Mueller matrices. A necessary condition for a homogeneous Mueller matrix is

$$m_{0,1} = m_{1,0}, m_{0,2} = m_{2,0}, m_{0,3} = m_{3,0} \quad (72)$$

then, $P(\mathbf{M}) = D(\mathbf{M})$.

The inhomogeneity of a Mueller matrix is characterized by the inhomogeneity index $I(\mathbf{M})$,

$$I(\mathbf{M}) = \frac{\sqrt{\hat{\mathbf{S}}_1 \cdot \hat{\mathbf{S}}_2}}{2} = \cos(\chi/2) \quad (73)$$

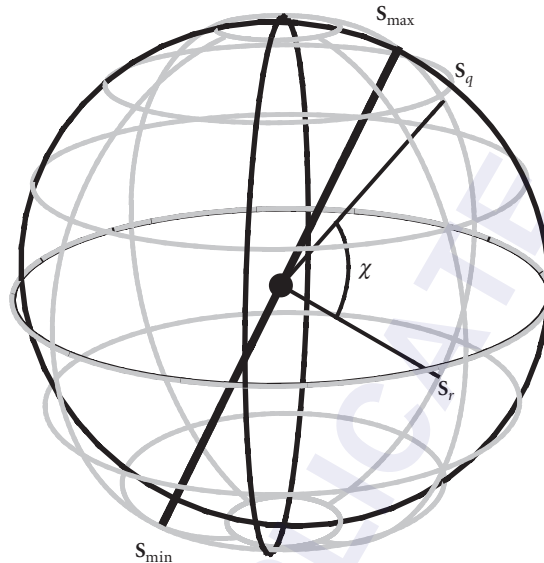


FIGURE 4 The principal Stokes vectors associated with an inhomogeneous polarization element mapped on the Poincaré sphere. The incident Stokes vectors of maximum S_{\max} and minimum S_{\min} intensity transmittance are diametrically opposite on the Poincaré sphere (indicating orthogonal polarization states) while the eigenpolarizations S_q and S_r are separated by the angle χ .

where \hat{S}_q and \hat{S}_r are normalized polarized Stokes vector eigenpolarizations of a Mueller matrix and χ is the angle between the eigenpolarizations on the Poincaré sphere measured from the center of the sphere as illustrated in Fig. 4. $I(\mathbf{M})$ varies from zero for orthogonal eigenpolarizations to one for degenerate (equal) eigenpolarizations.

The product of an arbitrary sequence of nondepolarizing Mueller matrices is another nondepolarizing Mueller matrix.

14.24 MUELLER MATRICES NEAR THE IDENTITY MATRIX, WEAK POLARIZATION ELEMENTS

The Mueller matrices of weak polarization elements are close to the 4×4 identity matrix. The properties of weak Mueller matrices are much simpler than general Mueller matrices because the retardance, diattenuation, and depolarization are close to zero. These simpler properties will be utilized later in the analysis of Mueller matrix properties by their matrix roots.

Some important examples of such weakly polarizing elements would be the lens surfaces and mirror surfaces in lenses, microscopes, and telescopes, where the polarization properties are not zero due to Fresnel equations, antireflection coatings, or mirrored surfaces, but the effects are often well below 5 percent.

The structure of the Mueller calculus and the properties of these weak elements can be explored by performing Taylor series on the Mueller matrix expressions with respect to diattenuation or retardance. Weak retarders have a retardance near zero. Performing a Taylor series expansion on

the general equation for an elliptical retarder and keeping the first-order terms yields the following simple expression for weak retarders:

$$\lim_{\delta_H, \delta_{45}, \delta_R \rightarrow 0} \mathbf{LR}(\delta_H, \delta_{45}, \delta_R) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & \delta_R & -\delta_{45} \\ 0 & -\delta_R & 1 & \delta_H \\ 0 & \delta_{45} & -\delta_H & 1 \end{pmatrix} \quad (74)$$

Similarly, a first-order Taylor series expansion on the general diattenuator expression yields

$$\lim_{d_H, d_{45}, d_R \rightarrow 0} \mathbf{D}(d_H, d_{45}, d_R, T_{\text{avg}}) = T_{\text{avg}} \begin{pmatrix} 1 & d_H & d_{45} & d_R \\ d_H & 1 & 0 & 0 \\ d_{45} & 0 & 1 & 0 \\ d_R & 0 & 0 & 1 \end{pmatrix} \quad (75)$$

Combining these yields the weak diattenuators and retarder Mueller matrix form

$$\mathbf{WDR}(d_H, d_{45}, d_R, \delta_H, \delta_{45}, \delta_R, T_{\text{avg}}) = T_{\text{avg}} \begin{pmatrix} 1 & d_H & d_{45} & d_R \\ d_H & 1 & \delta_R & -\delta_{45} \\ d_{45} & -\delta_R & 1 & \delta_H \\ d_R & \delta_{45} & -\delta_H & 1 \end{pmatrix} \quad (76)$$

These three equations are only correct to first order. Higher-order terms, which are present when these parameters are not infinitesimal, are calculated from the exact equations presented earlier.

So weak diattenuators are symmetric in the top row and first column. Weak retarders are antisymmetric in the off-diagonal lower right 3×3 elements. The presence of antisymmetric components in the top row and column and symmetric components in the lower right 3×3 elements of weak polarization element Mueller matrices indicates the presence of depolarization.

14.25 MATRIX ROOTS OF NONDEPOLARIZING MUELLER MATRICES

In this section an order-independent representation of Mueller-Jones matrices is developed using *matrix generators* to provide additional insights into the polarization properties of Mueller matrices. This matrix-generator approach is extended in the Sec. 14.31.

The matrix decompositions in Eq. (71) are order-dependent; the retardance occurs before the diattenuation or vice versa. Because the retardance and diattenuation components in general do not commute, the result is order-dependent.

An order-independent representation was developed for the Jones calculus by Jones with his N matrices, which represented Jones matrices with differential amounts of the three diattenuation and three retardance degrees of freedom. Jones approached the problem by considering propagation through a dichroic (diattenuating) and birefringent (retarding) anisotropic crystal and analyzed the Jones matrices as the path length t is cut into many (N) short lengths of t/N as N approaches infinity.

The same method is applicable to nondepolarizing Mueller-Jones matrices \mathbf{M}_N . Dividing the path length in half ($N = 2$) corresponds to taking the square root of \mathbf{M}_N ; this matrix square root has half the magnitude of the polarization properties of \mathbf{M}_N . For example, the square root of a quarter-wave retarder is an eighth-wave retarder with the same eigenpolarizations. The square root of a diattenuator with transmissions 1 and 0.64 has transmissions 1 and 0.8. For a square matrix, multiple

matrix square roots exist; a 4×4 matrix has 2^4 square roots. A real matrix square root always exists which is closer than \mathbf{M}_N to the identity matrix. Taking the series of these matrix roots

$$\sqrt[n]{\mathbf{M}_N} \quad n=2, 3, 4, \dots \quad (77)$$

yields a matrix sequence which approaches the identity matrix. The direction in which this sequence approaches the identity matrix depends on the ratio of polarization parameters and identifies the \mathbf{M}_N as linearly or elliptically diattenuating or retarding. The matrix properties in the vicinity of the identity matrix, as shown in Sec. 14.24 are particularly simple.

The differential Mueller matrices for diattenuation and retardance corresponding to the N matrices are presented below. They are obtained by taking high-order roots of Mueller matrices or by using differential values for the diattenuation and retardance in the equations for the basis retarder and diattenuator Mueller matrices. Such differential matrices are also known as *generators* and differ from the identity matrix by infinitesimal amounts which *point* in the direction of a particular polarization property.

The Mueller matrix generators describe infinitesimal amounts of each polarization property. There are three generators for diattenuation, $\mathbf{G}_1(d_1)$, $\mathbf{G}_2(d_2)$, and $\mathbf{G}_3(d_3)$, and three generators for retardance $\mathbf{G}_4(d_4)$, $\mathbf{G}_5(d_5)$, and $\mathbf{G}_6(d_6)$. The three diattenuation generators along with their first-order expansions and second-order terms are the following:

$$\mathbf{G}_1(d_1) = \begin{pmatrix} 1 & d_1 & 0 & 0 \\ d_1 & 1 & 0 & 0 \\ 0 & 0 & \sqrt{1-d_1} & 0 \\ 0 & 0 & 0 & \sqrt{1-d_1} \end{pmatrix} \approx \begin{pmatrix} 1 & d_1 & 0 & 0 \\ d_1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & -d_1^2 & 0 \\ 0 & 0 & 0 & -d_1^2 \end{pmatrix} \quad (78)$$

$$\mathbf{G}_2(d_2) = \begin{pmatrix} 1 & 0 & d_2 & 0 \\ 0 & \sqrt{1-d_2} & 0 & 0 \\ d_2 & 0 & 1 & 0 \\ 0 & 0 & 0 & \sqrt{1-d_2} \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & d_2 & 0 \\ 0 & 1 & 0 & 0 \\ d_2 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -d_2^2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -d_2^2 \end{pmatrix} \quad (79)$$

$$\mathbf{G}_3(d_3) = \begin{pmatrix} 1 & 0 & 0 & d_3 \\ 0 & \sqrt{1-d_3} & 0 & 0 \\ 0 & 0 & \sqrt{1-d_3} & 0 \\ d_3 & 0 & 0 & 1 \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 & d_3 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ d_3 & 0 & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -d_3^2 & 0 & 0 \\ 0 & 0 & -d_3^2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (80)$$

The equation for the general diattenuator is

$$\begin{aligned} \mathbf{M}_D(d_0, d_1, d_2, d_3) &= d_0 \lim_{N \rightarrow \infty} \left\{ \left[\mathbf{G}_1\left(\frac{d_1}{N}\right) \mathbf{G}_2\left(\frac{d_2}{N}\right) \mathbf{G}_3\left(\frac{d_3}{N}\right) \right]^N \right\} \\ &= d_0 \lim_{N \rightarrow \infty} \left\{ \left(\prod_{i=1}^3 \mathbf{G}_i\left(\frac{d_i}{N}\right) \right)^N \right\} \end{aligned} \quad (81)$$

The d_i are the diattenuation parameters and are restricted to the range

$$-1 \leq d_1, d_2, d_3 \leq 1 \quad (82)$$

and d_0 is an overall constant. For ideal polarizers,

$$\sqrt{d_1^2 + d_2^2 + d_3^2} = 1 \quad d_0 = \frac{1}{2} \quad (83)$$

The generators for retarders are as follows:

$$\mathbf{G}_4(d_4) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & \cos d_4 & \sin d_4 \\ 0 & 0 & -\sin d_4 & \cos d_4 \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & d_4 \\ 0 & 0 & -d_4 & 1 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & d_4^2 & 0 \\ 0 & 0 & 0 & d_4^2 \end{pmatrix} \quad (84)$$

$$\mathbf{G}_5(d_5) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos d_5 & 0 & -\sin d_5 \\ 0 & 0 & 1 & 0 \\ 0 & \sin d_5 & 0 & \cos d_5 \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & -d_5 \\ 0 & 0 & 1 & 0 \\ 0 & d_5 & 0 & 1 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & d_5^2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & d_5^2 \end{pmatrix} \quad (85)$$

$$\mathbf{G}_6(d_6) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos d_6 & \sin d_6 & 0 \\ 0 & -\sin d_6 & \cos d_6 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & d_6 & 0 \\ 0 & -d_6 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & d_6^2 & 0 & 0 \\ 0 & 0 & d_6^2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (86)$$

The parameters d_1 , d_2 , and d_3 are the retardance components in radians and are not limited in range. The general elliptical retarder has the form

$$\mathbf{M}_R(d_4, d_5, d_6) = \lim_{N \rightarrow \infty} \left\{ \left[\mathbf{G}_4\left(\frac{d_4}{N}\right) \mathbf{G}_5\left(\frac{d_5}{N}\right) \mathbf{G}_6\left(\frac{d_6}{N}\right) \right]^N \right\} = \lim_{N \rightarrow \infty} \left\{ \left(\prod_{i=4}^6 \mathbf{G}_i\left(\frac{d_i}{N}\right) \right)^N \right\} \quad (87)$$

The retardance generators are periodic in the retardances d_i ,

$$\mathbf{G}_i(d_i + \pi) = \mathbf{G}_i(d_i) \quad i = 4, 5, 6 \quad (88)$$

Combining these, the general nondepolarizing Mueller matrix \mathbf{M}_N (Mueller-Jones matrix), an inhomogeneous diattenuator and retarder, has the form

$$\mathbf{M}_N(d_0, \dots, d_6) = d_0 \lim_{N \rightarrow \infty} \left(\prod_{i=1}^6 \mathbf{G}_i\left(\frac{d_i}{N}\right) \right)^N \quad (89)$$

This representation is order-independent; the product of the six $\mathbf{G}_i(d_i)$ with differential d_i/N can be taken in any order of the six i , prior to raising to the N th power.

The d_i for a Mueller-Jones matrix \mathbf{M}_N are determined by calculating a high-order matrix root. For the majority of Mueller-Jones matrices, as the root order $N \rightarrow \infty$ the matrix root approaches a constant times the identity matrix with small first-order deviations, D_i , corresponding to the first-order terms above, as

$$\sqrt[N]{\mathbf{M}_N} = d_0 \begin{pmatrix} 1 & D_1 & D_2 & D_3 \\ D_1 & 1 & D_6 & -D_5 \\ D_2 & -D_6 & 1 & D_4 \\ D_3 & D_5 & -D_4 & 1 \end{pmatrix} \quad (90)$$

Due to the inexact nature of computer arithmetic, in practice N is set large enough that the matrix root is very close to the identity matrix without losing accuracy due to round off errors in calculations, which occurs when N is very large. Setting $10^4 < q < 10^9$ usually works well. The D_i are scaled by the matrix root order N ,

$$d_i = ND_i \quad (91)$$

yielding the magnitudes of the diattenuation and retardance parameters.

For any Mueller-Jones matrix, $\mathbf{M}_N(d_0, \dots, d_6)$, varying one of the d_i generates a family of \mathbf{M}_N where only a single polarization property varies. Varying each of the d_i one at a time generates a family of orthogonal trajectories through \mathbf{M}_N . Thus the set $\{d_0, d_1, d_2, d_3, d_4, d_5, d_6\}$ comprises a coordinate system for the Mueller-Jones matrices. Within the 15-dimensional space of normalized Mueller matrices, the Mueller-Jones matrices, parameterized by $\{d_1, d_2, d_3, d_4, d_5, d_6\}$ form an open six-dimensional surface embedded on a 15-dimensional hypersphere of radius $\sqrt{3}$, the surface where the depolarization index equals one.

A homogeneous nondepolarizing Mueller matrix has the same eigenvectors for its diattenuation (hermitian) and retarding (unitary) parts, and the corresponding condition on the Mueller roots is

$$\{d_1, d_2, d_3\} = k\{d_4, d_5, d_6\} \quad (92)$$

where k is a real constant, the ratio of the diattenuation to the retardance.

14.26 DEPOLARIZATION AND THE DEPOLARIZATION INDEX

Depolarization is the reduction of the degree of polarization of light. In the Mueller calculus depolarization can be pictured as a coupling of polarized into unpolarized light, where polarized light is incident and the exiting Stokes vector can be mathematically separated into a fully polarized and an unpolarized Stokes vector. Lenses, mirrors, filters, and other typical optical elements exhibit very small amounts of depolarization, typically less than a few tenths of a percent. In contrast, the depolarization of most diffusely reflecting objects such as paints, metal and wood surfaces, natural materials, and the like is significant, varying from a few to 100 percent (i.e., complete depolarization).

Two single-valued depolarization matrices, the depolarization index and the average degree of polarization, have been introduced to describe the degree to which a Mueller matrix depolarizes incident states.^{23–25} However, such single-number matrices cannot describe the complexity of depolarization associated with a Mueller matrix.

Consider three Mueller matrices of the following forms:

$$\mathbf{ID} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad \mathbf{PD} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & a & 0 & 0 \\ 0 & 0 & a & 0 \\ 0 & 0 & 0 & a \end{pmatrix} \quad \mathbf{DD} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & a & 0 & 0 \\ 0 & 0 & b & 0 \\ 0 & 0 & 0 & c \end{pmatrix} \quad (93)$$

Matrix **ID** is the ideal depolarizer; only unpolarized light exits the depolarizer. Matrix **PD** is the partial depolarizer; all fully polarized incident states exit with their incident polarization ellipse, but with a degree of polarization $\text{DoP}(\mathbf{PD} \cdot \mathbf{S}) = a$. The diagonal depolarizer matrix **DD** represents a variable partial depolarizer; the degree of polarization of the exiting light is a function of the incident state. Physically, depolarization is closely related to scattering and usually has its origin in retardance or diattenuation which is rapidly varying in time, space, or wavelength.

Of the 16 degrees of freedom in the Mueller matrix, 1 corresponds to loss, 3 to diattenuation, and 3 to retardance. The remaining 9 degrees of freedom describe depolarization.

14.27 DEGREE OF POLARIZATION SURFACES AND MAPS

The degree of polarization is a measure of the randomness of polarization in a light beam, a property characterized by how much of this beam may be blocked by a polarizer. Degree of polarization maps and surfaces represent this dependence of depolarization on incident polarization state.²⁶ For the typical depolarizer, different incident polarization states are depolarized by different amounts.

$$\text{DoP}(\mathbf{M} \cdot \mathbf{S}(\theta, \varepsilon)) = \text{DoP} \left(\begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \cdot \begin{pmatrix} 1 \\ \cos 2\theta \cos \varepsilon \\ \sin 2\theta \cos \varepsilon \\ \sin \varepsilon \end{pmatrix} \right) \quad (94)$$

Degree of polarization surfaces are formed from a nonuniform contraction of the Poincaré sphere corresponding to the depolarization properties of a Mueller matrix. The DoP surface for a Mueller matrix \mathbf{M} is formed by moving normalized Stokes vectors \mathbf{S} on the surface of the Poincaré sphere radially inward to a distance $\text{DoP}(\mathbf{S}' = \mathbf{M} \cdot \mathbf{S})$ from the origin, plotted for all incident \mathbf{S} on the surface of the Poincaré sphere. The DoP surface results from the product of a scalar, the DoP, and a vector, (s_1, s_2, s_3) , formed from the Stokes three-vector,

$$\text{DoP surface}(\mathbf{M}, \mathbf{S}) = \frac{\sqrt{S'_1(\mathbf{M}, \mathbf{S})^2 + S'_2(\mathbf{M}, \mathbf{S})^2 + S'_3(\mathbf{M}, \mathbf{S})^2}}{S'_0(\mathbf{M}, \mathbf{S})} (s_1, s_2, s_3) \quad (95)$$

for all $(s_1^2 + s_2^2 + s_3^2)^{1/2} = 1$.

The DoP map for a Mueller matrix is a two-dimensional plot, such as a contour plot or density plot, of the DoP of exiting light as a function of the incident polarized state and represents a “flattened” DoP surface. In this paper, the DoP map is plotted with axes θ (polarization ellipse major axis orientation) and DoCP, but there is some flexibility in the parameterization of the polarized Stokes vectors. In general the DoP map provides easier visualization of maxima, minima, saddles, and other features of the depolarization variation than the DoP surface.

Fig. 5 shows depolarization maps for a liquid crystal cell’s Mueller matrices measured at 2, 3, and 4 V at 550 nm in laboratory.²⁷ The cell is an untwisted nematic cell (Fredericksz cell) with fast and slow axes at $45^\circ/135^\circ$ used as a variable linear retarder for polarization control. The depolarization characteristics change significantly with applied voltage and incident polarization state. At 2 V, the maximum DoP (white area) occurs for a slightly elliptical state located just above the 135° linear state on the Poincaré sphere; a second DoP maximum occurs near 45° . At 0 V the liquid crystal (LC) directors will be parallel to one of these axes throughout the cell. So the least depolarization is occurring along the direction the LC molecules are aligned at the two surfaces. As the voltage increases, the molecules near the center of the cell gap rotate toward the cell normal while simultaneously the retardance decreases. Both maxima drift away from the fast/slow axes as the voltage increases while the depolarization, a measure of order, decreases the most around 90° . Unfortunately, polarization controllers are usually operated midway between the fast and slow axes and miss the DoP maxima.

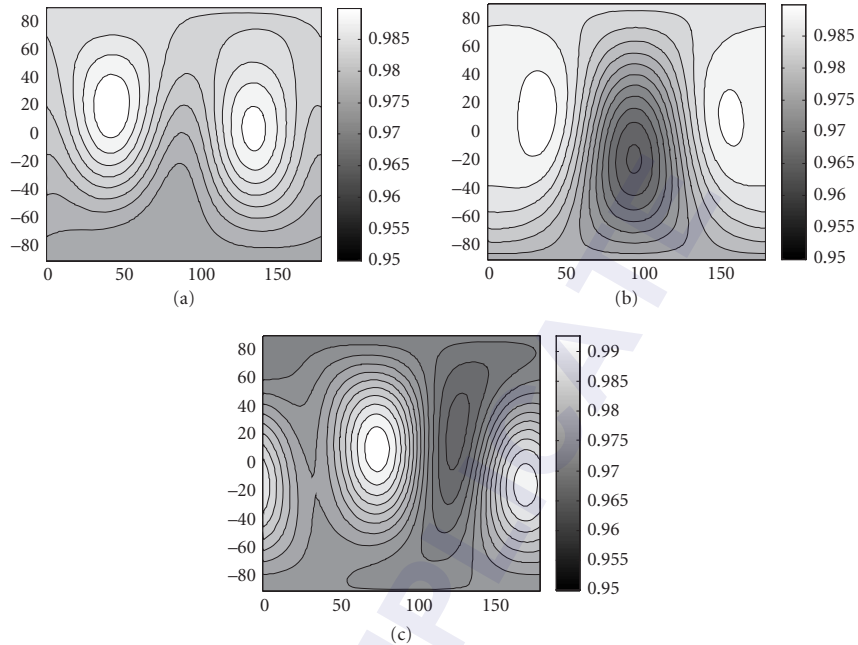


FIGURE 5 DoP maps for an untwisted nematic liquid crystal cell at (a) 2 V, (b) 3 V, and (c) 4 V. The DoP varies between 0.965 and 0.991 as the voltage and incident polarization state varies. At low voltages, the low depolarization states are close to the fast and slow axes. As the voltage increases, these DoP maxima (lighter gray) drift away.

14.28 THE DEPOLARIZATION INDEX

One figure of merit of the depolarization characteristics of a Mueller matrix is the depolarization index $DI(\mathbf{M})$ introduced by Gil and Bernabeu.^{23,24} $DI(\mathbf{M})$ is the euclidian distance of the normalized Mueller matrix $\mathbf{M}/m_{0,0}$ from the ideal depolarizer:

$$DI(\mathbf{M}) = \left\| \frac{\mathbf{M}}{m_{0,0}} - \mathbf{ID} \right\| = \frac{\sqrt{\left(\sum_{i,j} m_{i,j}^2 \right) - m_{0,0}^2}}{\sqrt{3} m_{0,0}} \quad (96)$$

$DI(\mathbf{M})$ varies from zero for the ideal depolarizer to 1 for all nondepolarizing Mueller matrices, including all pure diattenuators, pure retarders, and any sequences composed from them.

14.29 THE AVERAGE DEGREE OF POLARIZATION

The average degree of polarization, averageDoP, is the arithmetic mean of the degree of polarization of the exiting light when averaged over the Poincaré sphere,²⁵

$$\text{averageDoP}(\mathbf{M}) = \frac{\int_0^\pi \int_{-\pi/2}^{\pi/2} \text{DoP}(\mathbf{M} \cdot \mathbf{S}(\theta, \epsilon)) \cos(\epsilon) d\theta d\epsilon}{4\pi} \quad (97)$$

$S(\theta, \varepsilon)$ the Stokes vector parameterized by orientation of polarization θ and latitude ε , in radians on the Poincaré sphere, is

$$S(\theta, \varepsilon) = \begin{pmatrix} 1 \\ \cos(2\theta)\cos(\varepsilon) \\ \sin(2\theta)\cos(\varepsilon) \\ \sin(\varepsilon) \end{pmatrix} \quad (98)$$

The averageDoP varies from zero to one and provides a summary of the depolarizing property in a single number. When averageDoP is equal to one the exiting light is always completely polarized indicating a nondepolarizing Mueller matrix. Values near one indicates little depolarization. When averageDoP equals zero the exiting light is completely depolarized; only unpolarized light exits the interaction.

The DI and the averageDoP are often equal and usually quite close. The averageDoP is the easier metric to understand; it provides the mean DoP of the exiting light averaged over the Poincaré sphere, the expected value. The DI has a clear geometric meaning in the Mueller matrix configuration space, being the fractional distance of a Mueller matrix along a line segment from the ideal depolarizer to the hypersphere of nondepolarizing Mueller matrices, so it remains a useful and meaningful parameter, but more useful for theoretical studies of the Mueller calculus than for representing the depolarization of an optical element.

14.30 DETERMINING MUELLER MATRIX PROPERTIES

Given a Mueller matrix, measured or calculated, a natural question is “what are the polarization properties of this matrix?” Algorithms for the properties of Jones matrices were developed decades ago.^{1,15,28–30} For Mueller matrices, the algorithm development was challenging and work continues in this area.^{13,14,16,24,31–37}

A matrix decomposition expresses a matrix as a function of other matrices which indicate useful properties. For example, the polar decomposition of a Jones matrix expresses J as the matrix product of a hermitian matrix and a unitary matrix, corresponding to a diattenuator (partial polarizer) and retarder.^{21,22} Jones’ N-matrix decomposition of the Jones matrix divided the Jones matrix into a large number of identical matrix components infinitesimally close to the identity matrix and provided a simple description of Jones matrix properties.¹⁵

The Mueller matrix has 16 degrees of freedom and thus can be described by 16 unique properties. The Mueller roots decomposition presented in the following section is a generalization of Jones’ method to the Mueller matrix and provides an order-independent decomposition. The decomposition of Lu and Chipman, presented in the following section, expresses the Mueller matrix as the product of pure diattenuator, pure retarder, and depolarizer Mueller matrices.³⁰ The three components could occur in any specified order but the values of the components changed based on the order.

14.31 GENERATORS FOR DEPOLARIZATION

The generator method for describing polarization properties of Mueller-Jones matrices extends to depolarizing Mueller matrices by adding nine additional generators to span the remaining nine degrees of freedom.³⁸ These nine generators have been divided into three families of three generators corresponding to the three diattenuation and retardance degrees

of freedom. The nine generators for depolarization in exact, first- and second-order representations are as follows:

$$\mathbf{G}_7(d_7) = \begin{pmatrix} 1 & d_7 & 0 & 0 \\ -d_7 & 1-d_7 & 0 & 0 \\ 0 & 0 & \sqrt{1-d_7^2} & 0 \\ 0 & 0 & 0 & \sqrt{1-d_7^2} \end{pmatrix} \quad (99)$$

$$\approx \begin{pmatrix} 1 & d_7 & 0 & 0 \\ -d_7 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & -d_7^2 & 0 \\ 0 & 0 & 0 & -d_7^2 \end{pmatrix}$$

$$\mathbf{G}_8(d_8) = \begin{pmatrix} 1 & 0 & d_8 & 0 \\ 0 & \sqrt{1-d_8^2} & 0 & 0 \\ -d_8 & 0 & 1 & 0 \\ 0 & 0 & 0 & \sqrt{1-d_8^2} \end{pmatrix} \quad (100)$$

$$\approx \begin{pmatrix} 1 & 0 & d_8 & 0 \\ 0 & 1 & 0 & 0 \\ -d_8 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -d_8^2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -d_8^2 \end{pmatrix}$$

$$\mathbf{G}_9(d_9) = \begin{pmatrix} 1 & 0 & 0 & d_9 \\ 0 & \sqrt{1-d_9^2} & 0 & 0 \\ 0 & 0 & \sqrt{1-d_9^2} & 0 \\ -d_9 & 0 & 0 & 1 \end{pmatrix} \quad (101)$$

$$\approx \begin{pmatrix} 1 & 0 & 0 & d_9 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ -d_9 & 0 & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -d_9^2 & 0 & 0 \\ 0 & 0 & -d_9^2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\mathbf{G}_{10}(d_{10}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & \sqrt{1-d_{10}^2} & d_{10} \\ 0 & 0 & d_{10} & \sqrt{1-d_{10}^2} \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & d_{10} \\ 0 & 0 & d_{10} & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & -d_{10}^2 & 0 \\ 0 & 0 & 0 & -d_{10}^2 \end{pmatrix} \quad (102)$$

$$\mathbf{G}_{11}(d_{11}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \sqrt{1-d_{11}^2} & 0 & d_{11} \\ 0 & 0 & 1 & 0 \\ 0 & d_{11} & 0 & \sqrt{1-d_{11}^2} \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & d_{11} \\ 0 & 0 & 1 & 0 \\ 0 & d_{11} & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -d_{11}^2 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -d_{11}^2 \end{pmatrix} \quad (103)$$

$$\mathbf{G}_{12}(d_{12}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \sqrt{1-d_{12}^2} & d_{12} & 0 \\ 0 & d_{12} & \sqrt{1-d_{12}^2} & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \approx \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & d_{12} & 0 \\ 0 & d_{12} & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -d_{12}^2 & 0 & 0 \\ 0 & 0 & -d_{12}^2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (104)$$

The final three degrees of freedom describe the depolarization along the matrix diagonal elements: $m_{1,1}, m_{2,2}, m_{3,3}$:

$$\mathbf{G}(d_{13}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -\sqrt{\frac{3}{2}}d_{13} + \sqrt{1-d_{13}^2} & 0 & 0 \\ 0 & 0 & \sqrt{\frac{3}{2}}d_{13} + \sqrt{1-d_{13}^2} & 0 \\ 0 & 0 & 0 & \sqrt{1-d_{13}^2} \end{pmatrix} \quad (105)$$

$$= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 - \sqrt{\frac{3}{2}}d_{13} & 0 & 0 \\ 0 & 0 & 1 + \sqrt{\frac{3}{2}}d_{13} & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} - \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & d_{13}^2 & 0 & 0 \\ 0 & 0 & d_{13}^2 & 0 \\ 0 & 0 & 0 & -d_{13}^2 \end{pmatrix}$$

$$\mathbf{G}(d_{14}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \frac{d_{14}}{\sqrt{2}} + \sqrt{1-d_{14}^2} & 0 & 0 \\ 0 & 0 & \frac{d_{14}}{\sqrt{2}} + \sqrt{1-d_{14}^2} & 0 \\ 0 & 0 & 0 & -\sqrt{2}d_{14} + \sqrt{1-d_{14}^2} \end{pmatrix} \quad (106)$$

$$= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 + \frac{d_{14}}{\sqrt{2}} & 0 & 0 \\ 0 & 0 & 1 + \frac{d_{14}}{\sqrt{2}} & 0 \\ 0 & 0 & 0 & 1 - \sqrt{2}d_{14} \end{pmatrix} - \frac{1}{2} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & d_{14}^2 & 0 & 0 \\ 0 & 0 & d_{14}^2 & 0 \\ 0 & 0 & 0 & d_{14}^2 \end{pmatrix}$$

$$\mathbf{G}(d_{15}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1-d_{15} & 0 & 0 \\ 0 & 0 & 1-d_{15} & 0 \\ 0 & 0 & 0 & 1-d_{15} \end{pmatrix} \quad (107)$$

TABLE 4 Polarizing and Depolarizing Degrees of Freedom

	Horizontal and Vertical	45° and 135°	Right and Left
Diattenuation	d_1	d_2	d_3
Retardance	d_4	d_5	d_6
Amplitude depolarization	d_7	d_8	d_9
Phase depolarization	d_{10}	d_{11}	d_{12}
Diagonal depolarization	Orthogonal 1 d_{13}	Orthogonal 2 d_{14}	Radial d_{15}

Combining the 9 depolarizing generators with the 6 nondepolarizing generators yields an order-independent equation for Mueller matrices in terms of 15 polarization parameters,

$$\begin{aligned}
 \mathbf{M} &= d_0 \lim_{N \rightarrow \infty} \left[\mathbf{G}_1 \left(\frac{d_1}{N} \right) \mathbf{G}_2 \left(\frac{d_2}{N} \right) \mathbf{G}_3 \left(\frac{d_3}{N} \right) \mathbf{G}_4 \left(\frac{d_4}{N} \right) \mathbf{G}_5 \left(\frac{d_5}{N} \right) \mathbf{G}_6 \left(\frac{d_6}{N} \right) \mathbf{G}_7 \left(\frac{d_7}{N} \right) \mathbf{G}_8 \left(\frac{d_8}{N} \right) \right. \\
 &\quad \left. \mathbf{G}_9 \left(\frac{d_9}{N} \right) \mathbf{G}_{10} \left(\frac{d_{10}}{N} \right) \mathbf{G}_{11} \left(\frac{d_{11}}{N} \right) \mathbf{G}_{12} \left(\frac{d_{12}}{N} \right) \mathbf{G}_{13} \left(\frac{d_{13}}{N} \right) \mathbf{G}_{14} \left(\frac{d_{14}}{N} \right) \mathbf{G}_{15} \left(\frac{d_{15}}{N} \right) \right]^N \\
 &= d_0 \lim_{N \rightarrow \infty} \left(\prod_{i=1}^{15} \mathbf{G}_i \left(\frac{d_i}{N} \right) \right)^N
 \end{aligned} \tag{108}$$

As any single parameter is scanned in value, it generates a trajectory through the space of Mueller matrices where just one polarization property is changing. Thus Eq. (108) provides a coordinate system for depolarizing and nondepolarizing Mueller matrices in terms of the individual polarization parameters. The names of the polarization properties are given in Table 4. \mathbf{G}_7 , \mathbf{G}_8 , and \mathbf{G}_9 share the same first-order matrix elements as diattenuation and affect the flux of the light; thus the name amplitude diattenuation. \mathbf{G}_{10} , \mathbf{G}_{11} , and \mathbf{G}_{12} share the same first-order matrix elements as retardance and do not affect the flux of the light; thus the name phase diattenuation.

The final three depolarization degrees of freedom lie along the matrix diagonal. In most depolarizing samples these are the largest components of the depolarization. Several different bases can be considered for these degrees of freedom. The uniform depolarizer is a depolarizer which depolarized all polarization states equally. Combining the uniform depolarizer with any nondepolarizing Mueller matrix should move that matrix straight toward the ideal depolarizer in a radial direction in the 15-dimensional normalized Mueller matrix space. This uniform depolarizer generator is chosen as the final generator, \mathbf{G}_{15} , \mathbf{G}_{13} , \mathbf{G}_{14} , and \mathbf{G}_{15} are an orthogonal basis for the diagonal matrix elements.

The degree of polarization maps associated with the individual depolarization generators are shown in Fig. 6. Nonphysical regions with the degree of polarization greater than one are hatched; physical regions are solid. All of the generators except \mathbf{G}_{15} have DoP maps where half the area has a value below one and half is above one (hatched). Thus \mathbf{G}_7 through \mathbf{G}_{14} are nonphysical Mueller matrices by themselves. Only when they are combined with an equal or greater amount of \mathbf{G}_{15} , do they form physical Mueller matrices with DoP values always equal to or less than one. \mathbf{G}_{15} is an essential component of any depolarizing Mueller matrix.

The algorithm for the decomposition of Mueller matrix \mathbf{M} into matrix root parameters d_i has four steps. First, \mathbf{M} is normalized so its average transmission is one,

$$\hat{\mathbf{M}} = \frac{\mathbf{M}}{m_{0,0}} = \begin{pmatrix} 1 & m_{0,1}/m_{0,0} & m_{0,2}/m_{0,0} & m_{0,3}/m_{0,0} \\ m_{1,0}/m_{0,0} & m_{1,1}/m_{0,0} & m_{1,2}/m_{0,0} & m_{1,3}/m_{0,0} \\ m_{2,0}/m_{0,0} & m_{2,1}/m_{0,0} & m_{2,2}/m_{0,0} & m_{2,3}/m_{0,0} \\ m_{3,0}/m_{0,0} & m_{3,1}/m_{0,0} & m_{3,2}/m_{0,0} & m_{3,3}/m_{0,0} \end{pmatrix} \tag{109}$$

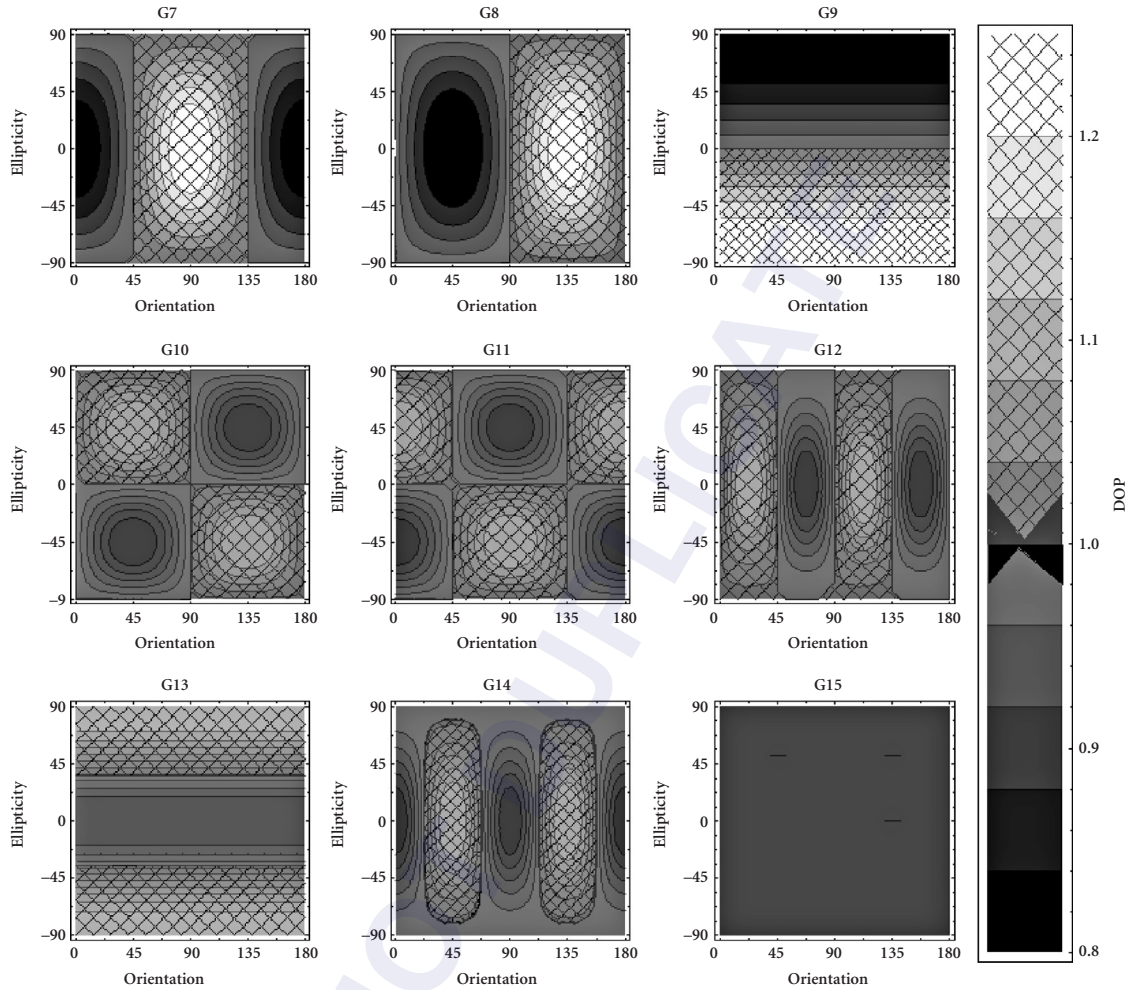


FIGURE 6 Degree of polarization maps associated with the depolarization generators G_7 – G_{15} for a value of $d_i = 0.2$. Hatched areas indicate incident polarization states where the exiting Stokes vector has a DoP > 1. G_{15} , radial depolarization, has a constant map.

Normalizing \mathbf{M} separates the property of the average transmission (nonpolarizing) from the remaining 15 polarizing properties. Second, a high-order matrix root of \mathbf{M} is calculated,

$$\mathbf{P} = \sqrt[q]{\hat{\mathbf{M}}} \quad (110)$$

For the majority of Mueller matrices, \mathbf{P} approaches a constant d_0 times the identity matrix as $q \rightarrow \infty$

$$\lim_{q \rightarrow \infty} \sqrt[q]{\hat{\mathbf{M}}} = d_0 \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (111)$$

Due to the inexact nature of computer arithmetic, in practice q is selected as a large enough number such that \mathbf{P} is very close to the identity matrix without losing accuracy due to round off errors in calculations. Setting $10^4 < q < 10^9$ usually works well.

Third, the difference between \mathbf{P} and the identity matrix is decomposed into 15 terms, corresponding to 15 polarization properties,

$$\mathbf{P} = \begin{pmatrix} P_{0,0} & P_{0,1} & P_{0,2} & P_{0,3} \\ P_{1,0} & P_{1,1} & P_{1,2} & P_{1,3} \\ P_{2,0} & P_{2,1} & P_{2,2} & P_{2,3} \\ P_{3,0} & P_{3,1} & P_{3,2} & P_{3,3} \end{pmatrix} = d_0 \begin{pmatrix} 1 & D_1 + D_7 & D_2 + D_8 & D_3 + D_9 \\ D_1 - D_7 & 1 - f_{13} & D_6 + D_{12} & -D_5 + D_{11} \\ D_2 - D_8 & -D_6 + D_{12} & 1 - f_{14} & D_4 + D_{10} \\ D_3 - D_9 & D_5 + D_{11} & -D_4 + D_{10} & 1 - f_{15} \end{pmatrix} \quad (112)$$

where the D_i are very small numbers which decrease linearly with large N . D_{13} , D_{14} , and D_{15} are more complex due to our selection of diagonal depolarization generators,

$$d_{13} = \frac{f_{13} - f_{14}}{\sqrt{6}} \quad d_{14} = \frac{-\sqrt{2}f_{13} + \sqrt{2}f_{14} + 2\sqrt{2}f_{15}}{6} \quad d_{15} = \frac{f_{13} + f_{14} + f_{15}}{3} \quad (113)$$

Finally, the D_i are scaled by the matrix root order q ,

$$d_i = qD_i \quad (114)$$

yielding the magnitudes of the polarization parameters.

One method to ensure the proper matrix root is calculated is to calculate the appropriate matrix square or cube root, then repeatedly apply the square or cube root function to generate high-order roots. Higham has addressed finding the real square roots of real matrices.^{39,40}

Several classes of Mueller matrices and their matrix roots require special consideration. For ideal polarizers ($T_{\max} = 1$, $T_{\min} = 0$), the Mueller matrix squared equals the Mueller matrix, then the high-order matrix roots do not approach the identity matrix. But when T_{\min} is infinitesimally increased above zero, the high-order matrix roots do approach the identity matrix.

Some highly depolarizing Mueller matrices have a negative determinant. The Mueller matrices generated by Eq. (108) span most of the space of physically realizable Mueller matrices. The determinant of normalized Mueller matrices vary over the range

$$\frac{-1}{27} \leq \det(\mathbf{M}) \leq 1 \quad (115)$$

Since \mathbf{P} has a positive determinant, its high-order powers can never equal a matrix with negative determinant. Unfortunately, the Mueller matrices spanned by Eq. (108) do not cover the small number of negative determinant Mueller matrices, an issue beyond the scope of this chapter.

The high-order roots of the majority of Mueller matrices are well behaved in approaching the identity matrix. The boundary between these classes is an area of ongoing investigation. An alternative form for the depolarization generators for the matrix diagonal is the following:

$$\mathbf{G}'(d_{13}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 - d_{13} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (116)$$

$$\mathbf{G}'(d_{14}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 - d_{14} & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (117)$$

$$\mathbf{G}'(d_{15}) = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1-d_{15} \end{pmatrix} \quad (118)$$

14.32 INTERPRETATION OF ARBITRARY MUELLER MATRICES, THE POLAR DECOMPOSITION OF MUELLER MATRICES

When simulating polarization elements, polarimeters, and optical systems, typically the Mueller matrices are nondepolarizing. Measured Mueller matrices have depolarization at some level; usually low for optical systems or much higher for light scattered from rough surfaces. To understand these matrices, it is desirable to interpret the Mueller matrices in terms of standard polarization elements. Given a Mueller matrix and access to a well-stocked inventory of polarization elements, diattenuator, retarder, and depolarizer could be combined to perform equivalent polarization transformations.

One common algorithm for Mueller matrix decomposition is the Lu-Chipman polar decomposition which represents an arbitrary depolarizing Muller matrix as the product of a pure diattenuator, a pure retarder, and a depolarizer Mueller matrices.²² The algorithm is complex; only a few steps will be outlined here; the reader is referred to other treatments of this algorithm and comments on its properties.^{10,37,41}

In this generalized polar decomposition, depolarizing Mueller matrices are expressed as the product of the three matrix factors: diattenuation, retardance, and depolarization,

$$\mathbf{M} = \mathbf{M}_\Delta \cdot \mathbf{M}_R \cdot \mathbf{M}_D \quad (119)$$

\mathbf{M}_D is the diattenuator Mueller matrix and \mathbf{M}_R is the retarder Mueller matrix. For the purpose of the Lu-Chipman decomposition algorithm, the depolarization Mueller matrix \mathbf{M}_Δ has the form

$$\mathbf{M}_\Delta = \begin{pmatrix} 1 & 0 & 0 & 0 \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{1,2} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{1,3} & m_{2,3} & m_{3,3} \end{pmatrix} = \begin{pmatrix} 1 & \vec{0}^T \\ \vec{\mathbf{P}}_\Delta & \mathbf{m}_\Delta \end{pmatrix} \quad (120)$$

The submatrix \mathbf{m}_Δ is symmetric, so it does not contain any retardance. The first column contains a polarizance term. The first step is to form \mathbf{M}_D from the top row of \mathbf{M} and then remove the diattenuation via

$$\mathbf{M}' = \mathbf{M} \cdot \mathbf{M}_D^{-1} = \mathbf{M}_\Delta \cdot \mathbf{M}_R \quad (121)$$

Extensive manipulations are then required to form \mathbf{M}_R and \mathbf{M}_Δ , and these are used to define the retardance and depolarization.

The Lu-Chipman decomposition algorithm has several disadvantages: (1) It is an order-dependent representation, so the polarization properties depend on which of the six polarization element sequence permutations is chosen. (2) Some highly depolarizing Mueller matrices have negative determinants; these negative determinant Mueller matrices form a very small subset of Mueller matrices, but they cannot be decomposed with the generalized polar decomposition.⁴¹ (3) The form for the depolarizer \mathbf{M}_Δ is peculiar, combining depolarization and polarizance, and sometimes \mathbf{M}_Δ which are by themselves nonphysical Mueller matrices are generated.

The advantage of this algorithm is that as a well-defined procedure, the values of diattenuation and retardance returned are carefully specified and reproducible. The results can be readily communicated between different groups.

Regarding the order-dependence of the algorithm, it makes sense to put the depolarizing part of the decomposition last (on the left). When fully polarized light is incident, the degree of polarization is not changed by \mathbf{M}_D or \mathbf{M}_R , so \mathbf{M}_Δ performs all the depolarizing. Two of the other five permutations are obtained by straightforward manipulations:

$$\mathbf{M} = \mathbf{M}_{\Delta 2} \cdot \mathbf{M}_{D 2} \cdot \mathbf{M}_{R 2} = \mathbf{M}_\Delta \cdot (\mathbf{M}_R \cdot \mathbf{M}_D \cdot \mathbf{M}_R^T) \cdot \mathbf{M}_R \quad (122)$$

$$\mathbf{M} = \mathbf{M}_{R 3} \cdot \mathbf{M}_{\Delta 3} \cdot \mathbf{M}_{D 3} = \mathbf{M}_R \cdot (\mathbf{M}_R^T \cdot \mathbf{M}_\Delta \cdot \mathbf{M}_R) \cdot \mathbf{M}_D \quad (123)$$

Straightforward relationships for the other three permutations are not easily obtained:

$$\mathbf{M} = \mathbf{M}_{D 4} \cdot \mathbf{M}_{R 4} \cdot \mathbf{M}_{\Delta 4} \quad (124)$$

$$\mathbf{M} = \mathbf{M}_{R 5} \cdot \mathbf{M}_{\Delta 5} \cdot \mathbf{M}_{D 5} \quad (125)$$

$$\mathbf{M} = \mathbf{M}_{D 6} \cdot \mathbf{M}_{\Delta 6} \cdot \mathbf{M}_{R 6} \quad (126)$$

14.33 PHYSICALLY REALIZABLE MUELLER MATRICES

Mueller matrices form a subset of the 4×4 real matrices. A 4×4 real matrix is not a physically realizable Mueller matrix if it can operate on an incident Stokes vector to produce a vector with degree of polarization greater than one ($S_0^2 < S_1^2 + S_2^2 + S_3^2$), which represents a physically unrealizable polarization state. Similarly, a Mueller matrix cannot output a state with negative flux. Conditions for physical realizability have been studied extensively in the literature, and many necessary conditions have been published.^{33,34,42–47} The following four necessary conditions for physical realizability are among the more general of those published:^{31,44}

1. $\text{Tr}(\mathbf{M}\mathbf{M}^T) \leq 4m_{0,0}^2$
2. $m_{0,0} \geq |m_{i,j}|$
3. $m_{0,0}^2 \geq b^2$
4. $(m_{0,0} - b)^2 \geq \sum_{j=1}^3 (m_{0,j} - \sum_{k=1}^3 m_{j,k} a_k)^2$

where $b = \sqrt{m_{0,1}^2 + m_{0,2}^2 + m_{0,3}^2}$, $a_j = m_{0,j}/b$, and Tr indicates the trace of a matrix.

Another condition for physical realizability is that the matrix can be expressed as a sum of nondepolarizing Mueller matrices. The Mueller matrix for a passive device $T_{\max} \leq 1$, a device without gain, must satisfy the relation $T_{\max} m_{0,0} = \sqrt{m_{0,1}^2 + m_{0,2}^2 + m_{0,3}^2} \leq 1$.

In the 16-dimensional space of Mueller matrices, the matrices for ideal polarizers, ideal retarders, and other nondepolarizing elements lie on the boundary between the physically realizable Mueller matrices and the unrealizable nonphysical matrices. Thus, a small amount of noise in the measurement of a Mueller matrix for a polarizer or retarder may yield a marginally unrealizable matrix.

When calculating a Mueller matrix \mathbf{M} from a set of flux measurements, there is error present due to nonidealities in the system. When error is present, it is often the case that the reconstructed \mathbf{M} is nonphysical, i.e., it is not possible to generate this \mathbf{M} using real components such as polarizers, retarders, and depolarizers. In this case degree of polarization is outside the range of 0 to 1, and/or the intensity is negative. The “nearest” physical matrix to \mathbf{M} may be found, both to reduce the error when extracting polarization parameters, and to give a quantifiable metric for the error in the measurement of \mathbf{M} .

Depolarization in Mueller matrices results from the addition of nondepolarizing Mueller matrices. The set of all normalized physically realizable Mueller matrices is thus formed from the convex hull of all the nondepolarizing Mueller matrices.

A requirement for a physically realizable Mueller matrix is that the complex hermitian matrix \mathbf{H} known as the coherency matrix derived from \mathbf{M} has non-negative eigenvalues.^{39,48}

$$\mathbf{H} = \frac{1}{2} \sum_{j=0}^3 \sum_{i=0}^3 m_{ij} (\boldsymbol{\sigma}_i \otimes \boldsymbol{\sigma}_j^*) \quad (127)$$

where the $\boldsymbol{\sigma}$'s are the normalized Pauli matrices:

$$\boldsymbol{\sigma}_0 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad \boldsymbol{\sigma}_1 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \quad \boldsymbol{\sigma}_2 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad \boldsymbol{\sigma}_3 = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} \quad (128)$$

and \otimes indicates the outer product function flattened into a matrix, i.e.,

$$\boldsymbol{\sigma}_i \otimes \boldsymbol{\sigma}_j^* = \begin{pmatrix} \boldsymbol{\sigma}_{i,0,0} \boldsymbol{\sigma}_{j,0,0}^* & \boldsymbol{\sigma}_{i,0,0} \boldsymbol{\sigma}_{j,0,1}^* & \boldsymbol{\sigma}_{i,0,1} \boldsymbol{\sigma}_{j,0,0}^* & \boldsymbol{\sigma}_{i,0,1} \boldsymbol{\sigma}_{j,0,1}^* \\ \vdots & \ddots & \vdots & \vdots \\ \boldsymbol{\sigma}_{i,1,0} \boldsymbol{\sigma}_{j,1,0}^* & \cdots & \cdots & \boldsymbol{\sigma}_{i,1,1} \boldsymbol{\sigma}_{j,1,1}^* \end{pmatrix} \quad (129)$$

The Mueller matrix is reconstructed from \mathbf{H} as

$$m_{i,j} = \text{Tr}\{\mathbf{H} \cdot (\boldsymbol{\sigma}_i \otimes \boldsymbol{\sigma}_j^*)\} \quad (130)$$

If any of the four eigenvalues of \mathbf{H} are negative, then \mathbf{M} is not physical. *This is the fundamental test for physical \mathbf{M} .* Often when a Mueller matrix is measured in the presence of noise, the Mueller Matrix is close to physical. Usually three of the four eigenvalues are positive, and one is small and negative.

Two methods are presented to construct a physical Mueller matrix \mathbf{M}_p from a physical Mueller matrix \mathbf{M} : (1) One method first calculates the hermitian matrix \mathbf{H}_p corresponding to the closest \mathbf{M}_p via optimization in the coherency matrix domain.⁴⁸ Various metrics can define the meaning of *closest*. (2) A faster method finds \mathbf{H}_p which has all nonnegative eigenvalues (i.e., is positive semidefinite).

Method 1

To calculate the closest \mathbf{M}_p , the first step is to construct a positive semidefinite matrix \mathbf{H}_{opt} from a Cholesky decomposition of the Mueller matrix \mathbf{M} ,

$$\mathbf{M} = \mathbf{C}^\dagger \cdot \mathbf{C} \quad (131)$$

where \dagger indicates conjugate transpose, and \mathbf{C} is an upper triangular matrix. If the Cholesky decomposition of \mathbf{M} exists, then \mathbf{M} must be positive semidefinite. To preserve the magnitude of \mathbf{H}_{opt} , its Cholesky decomposition is normalized by a constant,

$$\mathbf{H}_{\text{opt}} = \frac{2\mathbf{C}^\dagger \cdot \mathbf{C}}{\text{Tr}(\mathbf{C}^\dagger \cdot \mathbf{C})} \quad (132)$$

where \mathbf{C} is given by

$$\mathbf{C} = \begin{pmatrix} h_1 & h_5 + ih_6 & h_{1,1} + ih_{1,2} & h_{1,5} + ih_{1,6} \\ 0 & h_2 & h_7 + ih_8 & h_{1,3} + ih_{1,4} \\ 0 & 0 & h_3 & h_9 + ih_{1,0} \\ 0 & 0 & 0 & h_4 \end{pmatrix} \quad (133)$$

To perform the optimization, \mathbf{H}_{opt} (having 16 variable parameters $h_1 - h_{16}$) is formed via Eqs. (132) and (133), and the 16 h parameters are varied until the difference metric between \mathbf{M}_p and \mathbf{M} is minimized. A starting point for the 16 parameters can be selected as \mathbf{H}_{init} ,

$$\mathbf{H}_{\text{init}} = \frac{2\sqrt{\mathbf{H}}(\sqrt{\mathbf{H}})^\dagger}{\text{Tr}(\sqrt{\mathbf{H}}(\sqrt{\mathbf{H}})^\dagger)} \quad (134)$$

The Cholesky decomposition of \mathbf{H}_{init} yields the upper triangular matrix \mathbf{C}_{init} , and the 16 starting values. By minimizing the Frobenius distance F_d between \mathbf{M}_p and \mathbf{M} ,

$$F_d = \frac{\|\mathbf{M} - \mathbf{M}_p\|}{\|\mathbf{M} + \mathbf{M}_p\|} \quad (135)$$

where $\|\cdot\|$ is any appropriate matrix norm. The euclidean norm for Mueller matrix \mathbf{M} is

$$\|\mathbf{M}\| = \sqrt{\sum_{i=0}^3 \sum_{j=0}^3 m_{i,j}^2} \quad (136)$$

Method 2

The Cholesky decomposition method is computationally intensive. A more expedient method for generating a physical \mathbf{M}_p is to find the eigenvalues $\{\lambda_1, \lambda_2, \lambda_3, \lambda_4\}$ and eigenvectors $\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4\}$ of \mathbf{H} . Any negative eigenvalues are set to zero. The set of positive definite eigenvalues are formed into a diagonal matrix \mathbf{D} :

$$\mathbf{D} = \begin{pmatrix} \lambda'_1 & 0 & 0 & 0 \\ 0 & \lambda'_2 & 0 & 0 \\ 0 & 0 & \lambda'_3 & 0 \\ 0 & 0 & 0 & \lambda'_4 \end{pmatrix} \quad (137)$$

$$\text{where } \lambda'_i = \begin{cases} \lambda_i & \text{if } \lambda_i \geq 0 \\ 0 & \text{if } \lambda_i < 0 \end{cases}$$

and the eigenvectors are formed into a square unitary matrix \mathbf{U} ,

$$\mathbf{U} = \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{pmatrix} \quad (138)$$

A new physically realizable \mathbf{H}_p is formed as the product

$$\mathbf{H}_p = \mathbf{U} \cdot \mathbf{D} \cdot \mathbf{U}^{-1} \quad (139)$$

per the Eigen decomposition theorem.^{49,50} Finally \mathbf{M}_p is derived from \mathbf{H}_p via Eq. (130) yielding a dramatically faster algorithm. This \mathbf{M}_p is not necessarily the closest one to the original nonphysical \mathbf{M} .

14.34 ACKNOWLEDGMENTS

In preparing this chapter extensive use has been made of the works of some of the author's students and other collaborators and they deserve particular credit: Shih-Yau Lu, Karlton Crabtree, Brian deBoo, Karen Tweetmeyer, Garam Yun, and Neil Beaudry.

14.35 REFERENCES

1. R. C. Jones, "A New Calculus for the Treatment of Optical Systems: I. Description and Discussion of the Calculus," *J. Opt. Soc. Am.* **31**:488–493 (1941).
2. E. Wolf, "Coherence Properties of Partially Polarized Electromagnetic Radiation," *Nuovo Cim.* **13**:1165–1181 (1959).
3. W. A. Shurcliff, *Polarized Light—Production and Use*, Harvard University Press, Cambridge, Mass (1962).
4. A. Gerrard and J. M. Burch, *Introduction to Matrix Methods in Optics*: Wiley, London (1975).
5. P. S. Theocaris, and E. E. Gdoutos, *Matrix Theory of Photoelasticity*: Springer-Verlag, Heidelberg, Germany (1979).
6. R. M. A. Azzam and N. M. Bashara, *Ellipsometry and Polarized Light*: Elsevier, Amsterdam (1987).
7. K. L. Coulson, *Polarization and Intensity of Light in the Atmosphere*: A. Deepak, Hampton (1988).
8. W. Egan, ed., "Polarization in Remote Sensing," *Proc. Soc. Photo—Opt. Instrum. Eng.*, **1747**:2–48 (1992).
9. C. Brosseau, "Fundamentals of Polarized Light: A Statistical Optics Approach," Wiley-Interscience, Malden, Mass (1998).
10. D. Goldstein, *Polarized Light*: Marcel Dekker, New York (2003).
11. P. Soleillet, "Parameters Characterising Partially Polarized Light in Fluorescence Phenomena," *Ann. Phys.* **12**(10):23–97 (1929).
12. H. Mueller, "The Foundation of Optics," *J. Opt. Soc. Am.* **38**:661 (1948).
13. N. G. Parke, III, *Matrix Optics*, MIT, Ph.D (1948).
14. N. G. Parke, III, "Optical Algebra," *J. Math. Phys.* **28**(2):131–139 (1949).
15. R. C. Jones, "A New Calculus for the Treatment of Optical Systems: VII. Properties of the N-Matrices," *J. Opt. Soc. Am.* **38**:671–685 (1948).
16. E. S. Fry and G. W. Kattawar, "Relationships Between Elements of the Stokes Matrix," *Appl. Opt.* **20**:2811–2814 (1981).
17. R. A. Chipman, "Polarization Analysis of Optical Systems," *Opt. Eng.* **28**(2):90–99 (1989).
18. P. W. Maymon and R. A. Chipman, "Linear Polarization Sensitivity Specifications for Spaceborne Instruments. Polarization Analysis and Measurement," *Proc. SPIE* **1746**:148–156 (1992).
19. R. Simon, "The Connection Between Mueller and Jones Matrices of Polarization Optics," *Opt. Commun.* **42**:293–297 (1982).
20. K. Kim, L. Mandel, et al., "Relationship Between Jones and Mueller Matrices for Random Media," *J. Opt. Soc. Am. A.* **4**:433–437 (1987).
21. J. J. Gil and E. Bernabeu, "Obtainment of the Polarizing and Retardation Parameters of a Non-Depolarizing Optical System from the Polar Decomposition of Its Mueller Matrix," *Optik* **76**:67–71 (1987).
22. S. Y. Lu and R. A. Chipman, "Homogeneous and Inhomogeneous Jones Matrices," *J. Opt. Soc. Am. A.* **11**(2):766–773 (1994).
23. J. J. Gil and E. Bernabeu, "A Depolarization Criterion in Mueller Matrices," *Opt. Acta.* **32**:259–261 (1985).
24. J. J. Gil and E. Bernabeu, "Depolarization and Polarization Indices of an Optical System," *Opt. Acta.* **33**:185–189 (1986).
25. R. A. Chipman, "Depolarization Index and the Average Degree of Polarization," *Appl. Opt.* **44**(13):2490–2495 (2004).
26. De Boo, B. J. Sasian, et al., "Degree of Polarization Surfaces and Maps for Analysis of Depolarization," *Opt. Exp.* **12**(20):4941–4958 (2004).
27. J. Wolfe and R. A. Chipman, "Polarimetric Characterization of Liquid Crystal on Silicon Panels," *Appl. Opt.* **44**(36):7853–7857 (2005).
28. R. C. Jones, "A New Calculus for the Treatment of Optical Systems: III. The Sohncke Theory of Optical Activity," *J. Opt. Soc. Am.* **31**:500–503 (1941).
29. R. C. Jones, "A New Calculus for the Treatment of Optical Systems: IV," *J. Opt. Soc. Am.* **32**:486–493 (1942).
30. S.-Y. Lu and R. A. Chipman, "An Interpretation of Mueller Matrices Based upon the Polar Decomposition," *J. Opt. Soc. Am. A.* **13**(5):1106–1113 (1996).

31. R. Barakat, "Bilinear Constraints Between Elements of the 4×4 Mueller-Jones Transfer Matrix of Polarization Theory," *Opt. Commun.* **38**:159–161 (1981).
32. S. R. Cloude, "Group Theory and Polarization Algebra," *Optik* **75**:26–36 (1986).
33. R. Barakat, "Conditions for the Physical Realizability of Polarization Matrices Characterizing Passive Systems," *J. Mod. Opt.* **34**:1535–1544 (1987).
34. S. R. Cloude, "Conditions for the Physical Realizability of Matrix Operators in Polarimetry," *Proc. Soc. Photo—Opt. Instrum. Eng.* **1166**:177–185 (1989).
35. R. A. Chipman, "The Structure of the Mueller Calculus," *Polarization Measurement and Analysis, Proc. SPIE* **4133**:1–9, San Diego, Calif (2000).
36. J. J. Gil, "Characteristic Properties of Mueller Matrices," *J. Opt. Soc. Am. A.* **17**(2):328–334 (2000).
37. J. J. Gil, "Polarimetric Characterization of Light and Media," *Eur. Phys. J. Appl. Phys* **40**:1–47 (2007).
38. R. A. Chipman, "Depolarization in Mueller Calculus," *Proc. SPIE* **5158**:184–192 (2003).
39. N. J. Higham, "Computing Real Square Roots of a Real Matrix," *Linear Algebra and Its Applications* **88–89**:405–430 (1987).
40. N. J. Higham, "Stable Iterations for the Matrix Square Root," *Num. Algo.* **15**:227–242 (1987).
41. R. Ossikovski, A. De Martino, S. Guyot, "Forward and Reverse Product Decompositions of Depolarizing Mueller Matrices," *Opt. Lett.* **32**(6): 689–691 (2007).
42. J. W. Hovenier, H. C. Van de Hulst, et al., "Conditions for the Elements of the Scattering Matrix," *Astron. Astrophys.* **157**:301–310 (1986).
43. S. S. Girgel, "Structure of the Mueller Matrices of Depolarized Optical Systems," *Sov. Phys. Crystallogr.* **36**:890–891 (1991).
44. A. B. Kostinski, C. R. Clark, et al., "Some Necessary Conditions on Mueller Matrices," *Proc. Soc. Photo-Opt. Instrum. Eng.* **1746**:213–220 (1992).
45. M. S. Kumar and R. Simon, "Characterizing of Mueller Matrices in Polarization Optics," *Opt. Commun.* **88**:464–470 (1992).
46. C. V. M. van der Mee and J. W. Hovenier, "Structure of Matrices Transforming Stokes Parameters," *J. Math. Phys.* **33**:3574–3584 (1992).
47. Z.-F. Xing, "On the Deterministic and Non-Deterministic Mueller Matrix," *J. Mod. Opt.* **39**:461–484 (1992).
48. A. Aiello, G. Puentes, D. Voigt, J. P. Woerdman, "Maximum-Likelihood Estimation of Mueller Matrices," *Opt. Lett.* **31**:817–819 (2006).
49. E. W. Weisstein, "Eigen Decomposition Theorem" from <http://mathworld.wolfram.com/EigenDecompositionTheorem.html> (2008).
50. W. H. Press, S. A. Teukolsky, et al. *Numerical Recipes in C*: Cambridge University Press. Cambridge, UK (1988).

Russell A. Chipman

College of Optical Sciences
University of Arizona
Tucson, Arizona

15.1 GLOSSARY

A	analyzer vector
<i>a</i>	analyzer vector element
<i>a</i>	semimajor axis of ellipse
A, B, C	principal axes for a dielectric tensor
BRDF	bidirectional reflectance distribution function
<i>b</i>	semiminor axis of ellipse
C_M	covariance matrix
<i>d</i>	liquid crystal cell gap
D	diagonal matrix
<i>D</i>	diattenuation
DoCP	degree of circular polarization
DoLP	degree of linear polarization
DoP	degree of polarization
<i>E</i>	extinction ratio
<i>e</i>	ellipticity
EM	error metric
<i>I</i>	inhomogeneity of a Mueller matrix
<i>I</i>	first Stokes element, S_0 , flux
<i>i, j, k</i>	summation indices
k	propagation vector
<i>L, L₂</i>	condition number of a matrix
LP	linear polarizer
M	Mueller matrix
\vec{M}	Mueller vector
MMBRDF	Mueller matrix bidirectional reflectance distribution function
M_R	polarimeter's estimate of the Mueller matrix
M_R, M_T	beamsplitter Mueller matrix in reflection and transmission

\mathbf{M}_s	sample Mueller matrix
$m_{00}, m_{01}, \dots, m_{33}$	Mueller matrix elements
Δn	birefringence
n_1, n_2	refractive indices of a birefringent medium
\mathbf{P}	flux measurement vector
P	polarizance
P	flux measurement
PAF	polarization aberration function
PSM	point spread matrix
QWLR	quarter-wave linear retarder
Q, R	index limit
Q	second Stokes vector element, S_1
q	index for a sequence of polarization elements
r	index for polarimeter variables
\mathbf{S}	Stokes vector
\mathbf{S}'	exiting Stokes vector
\mathbf{S}_m	measured Stokes vector
S_0, S_1, S_2, S_3	Stokes vector elements
$\mathbf{S}_p, \mathbf{S}_u$	polarized and unpolarized part of Stokes vector
SD	standard deviation
T	transpose, superscript
T_{\max}	maximum intensity transmittance
T_{\min}	minimum intensity transmittance
t	time
t	thickness
\mathbf{U}	Jones/Mueller transformation matrix
U	third Stokes vector element, S_2
\mathbf{U}, \mathbf{V}	unitary matrices
V	fourth Stokes vector element, S_3
\mathbf{W}	polarimetric measurement matrix
\mathbf{W}^{-1}	polarimetric data-reduction matrix
\mathbf{W}_p^{-1}	pseudoinverse of \mathbf{W}
α	liquid crystal cell rubbing direction
β	bistatic angle
γ, δ	angles of scatter
δM_k	error in Mueller matrix element
$\langle \delta \mathbf{M} \rangle$	mean
$\langle \delta \mathbf{W} \rangle$	corrections to polarimetric measurement matrix
δ	retardance
ε	eccentricity
$\tilde{\mathbf{\epsilon}}$	dielectric tensor
η	azimuth of ellipse
θ	Euler angle
θ	orientation angle
θ_p, ϕ_i	angles of incidence

θ_s, ϕ_s	angles of scatter
Θ_1, Θ_2	pretilt at liquid crystal cell input and output
κ_p	condition number
λ	wavelength
μ_k	singular values
φ	Euler angle
Φ	liquid crystal cell twist angle
ϕ	phase of a complex number
ψ	Euler angle
Ω	solid angle, steradians
•	dot product, matrix multiplication

15.2 OBJECTIVES

The principles of polarization measurements are surveyed in this chapter. One of the primary difficulties in performing accurate polarization measurements is the systematic errors due to non-ideal polarization elements. Therefore, the polarimetric measurement and data-reduction process is formulated to incorporate arbitrary polarization elements calibrated by measurement of their transmitted and analyzed Stokes vectors. Polarimeter optimization is addressed through the minimization of the condition number. First derivatives of the polarimetric measurement matrix provide an error analysis method. Methods for polarization modulation are compared. The chapter concludes with a survey of polarimeter applications including the following sections: “Ellipsometry and Generalized Ellipsometry,” “Liquid Crystal Cell and System Testing,” “Polarization Aberrations,” “Remote Sensing,” “Polarization Light Scattering,” “Ophthalmic Polarimetry.”

Throughout this chapter, quantities are formulated in terms of the Stokes vector and Mueller matrix, as these usually comprise the most appropriate representation of polarization for radiometric measurements.

15.3 POLARIMETERS

Polarimeters are optical instruments for measuring the polarization properties of light beams and samples. Polarimetry, the science of polarization measurement, is most simply characterized as *radiometry with polarization elements*. Accurate polarimetry requires careful attention to all the issues necessary for accurate radiometry, together with many additional polarization issues which must be mastered to accurately determine polarization properties from polarimetric measurements.

Typical applications of polarimeters include the following: remote sensing of the earth and astronomical bodies, calibration of polarization elements, measurements of the thickness and refractive indices of thin films (ellipsometry), spectropolarimetric studies of materials, and alignment of polarization-critical optical systems such as liquid crystal displays and projectors. Polarimeters are divided into several categories.

15.4 LIGHT-MEASURING POLARIMETERS

Light-measuring polarimeters measure the polarization state of a beam of light and its polarization characteristics: the Stokes vector, the direction of oscillation of the electric field vector for a linearly polarized beam, the helicity of a circularly polarized beam, the elliptical parameters of an elliptically polarized beam, and the degree of polarization.

A light-measuring polarimeter utilizes a set of polarization elements placed in a beam of light in front of a radiometer. The light beam is analyzed by this set of polarization state analyzers, and a set of flux measurements is acquired. The polarization characteristics of the light beam are determined from these measurements by data-reduction algorithms (see “Data Reduction for Light-measuring Polarimeters”).

15.5 SAMPLE-MEASURING POLARIMETERS

Sample-measuring polarimeters determine the relationship between the polarization states of incident and exiting beams for a sample. The term *exiting beam* is general and in different measurements might describe beams which are transmitted, reflected, diffracted, scattered, or otherwise modified. The term sample is also an inclusive term used in a broad sense to describe a general light-matter interaction or sequence of such interactions and applies to practically anything.

Measurements are acquired using a set of polarization elements located between a source and sample and the exiting beams are analyzed with a separate set of polarization elements between the sample and radiometer. Samples of great interest include surfaces, thin films on surfaces, polarization elements, optical elements, optical systems, natural scenes, biological samples, and industrial samples.

Accurate polarimetric measurements can be made only if the polarization generator and polarization analyzer are well calibrated. To perform accurate polarimetry, the polarization elements need not be ideal or of the highest quality. If the Mueller matrices of the polarization components are known from careful calibration, the systematic errors due to nonideal polarization elements are removed during the data reduction (see “Polarimetric Measurement Equation and Polarimetric Data-Reduction Equation”).

15.6 COMPLETE AND INCOMPLETE POLARIMETERS

A light-measuring polarimeter is *complete* if a Stokes vector can be determined from its measurements. An *incomplete* light-measuring polarimeter cannot be used to determine a Stokes vector. For example, a polarimeter which employs a rotating polarizer in front of a detector does not determine the circular polarization content of a beam, and is incomplete. Similarly, a sample-measuring polarimeter is complete if it is capable of measuring the full Mueller matrix, and incomplete otherwise. Complete polarimeters are referred to as Stokes polarimeters or Mueller polarimeters.

15.7 POLARIZATION GENERATORS AND ANALYZERS

A *polarization generator* consists of a light source, optical elements, and polarization elements to produce a beam of known polarization state. A polarization generator is specified by the Stokes vector S of the exiting beam. A *polarization analyzer* is a configuration of polarization elements, optical elements, and a detector which performs a flux measurement of a particular polarization component in an incident beam. A polarization analyzer is characterized by a Stokes-like *analyzer vector* A which specifies the incident polarization state which is analyzed, the state which produces the maximal response at the detector. Sample-measuring polarimeters require polarization generators and polarization analyzers, while light-measuring polarimeters only require polarization analyzers. Frequently the terms “polarization generator” and “polarization analyzer” refer just to the polarization elements in the generator and analyzer. It is important to distinguish between elliptical (and circular) generators and elliptical analyzers for a given state because they generally have different

polarization characteristics and different Mueller matrices (see “Elliptical and Circular Polarizers and Analyzers”).

15.8 CLASSES OF POLARIMETERS

Polarimeters operate by acquiring measurements with a set of polarization analyzers. The following sections classify polarimeters by the four broad methods by which these multiple measurements are most often acquired. A complete Stokes polarimeter requires a minimum of four flux measurements with a set of linearly independent polarization generators in order to set up four equations in four unknowns, the four Stokes parameters. Many Stokes polarimeters use more than four flux measurements to improve signal to noise and/or reduce systematic errors.

15.9 TIME-SEQUENTIAL MEASUREMENTS

In a time-sequential polarimeter, the series of flux measurements are taken sequentially in time. Between measurements, the polarization generator and analyzer are changed. Time-sequential polarimeters frequently employ rotating polarization elements or filter wheels containing a set of analyzers. A time-sequential polarimeter generally employs a single source and single detector or focal plane.

15.10 POLARIZATION MODULATION

Polarization modulation polarimeters contain a polarization modulator, a rapidly changing polarization element. The output of the analyzer is a rapidly fluctuating irradiance on which polarization information is encoded. Polarization parameters are determined by lock-in amplifiers or by frequency-domain digital signal processing techniques. For example, a rapidly spinning polarizer produces a modulated output which allows the flux and the degree of linear polarization to be read with a DC voltmeter and an AC voltmeter. The most common high-speed polarization modulator in general use is the photoelastic modulator.

15.11 DIVISION OF APERTURE

Division of aperture polarimeters use multiple polarization analyzers operating side by side. The aperture of the polarimeter beam is subdivided. Each beam propagates through a separate polarization analyzer to a separate detector. The detectors are usually synchronized to acquire measurements simultaneously. This is similar in principle to the polarizing glasses used in three-dimensional movie systems, where different analyzers are placed over each eye, sometimes a 45° and a 135° polarizer, sometimes a right and left circular analyzer, presenting two different perspective views simultaneously to each eye.

15.12 DIVISION OF AMPLITUDE

Division-of-amplitude polarimeters utilize beam splitters to divide the measured beam and direct the component beams to multiple analyzers and detectors. A division-of-amplitude polarimeter can acquire its measurements simultaneously, providing advantages for rapidly changing scenes

or measurements from moving platforms. Many division-of-amplitude polarimeters use polarizing beam splitters to simultaneously divide and analyze the beam.

15.13 SPECTROPOLARIMETERS

Polarimeters can be combined with monochromators or spectrometers to measure Stokes vector spectra or Mueller matrix spectra. Because grating monochromators have large diattenuation which varies rapidly with wavelength, the monochromator should be configured to use only a single polarization.

15.14 IMAGING POLARIMETERS

When the polarimeter's detector is a focal plane array, a series of images acquired with different analyzers (the raw images) can be reduced to measure a Stokes vector image or a Mueller matrix image.

Imaging polarimeters are particularly susceptible to misalignment of the raw images since polarization properties are determined from the difference between flux measurements. Such misalignment causes polarization artifacts in the image on account of spurious polarization mixed with the actual polarization.

Raw image misalignments occur due to source motion, polarimeter motion, vibration, and beam wander from slight wedge in rotating components. Polarization artifacts are largest in areas where the image intensity is changing the fastest, around object edges and near point sources. The edges of objects are usually where the angles of incidence and angles of scatter are larger. The largest polarization is typically expected around these areas, but due to vibration, image motion, and image misalignment, these are also the areas where the data is most suspect. Other errors result from imperfect polarization elements, and detector noise.

When the source flux fluctuates between raw images, a uniform polarization error occurs across the entire image. Source fluctuations are a serious problem in outdoor Stokes imagery because sunlight fluctuates due to cloud motion.

Many of the polarization images and spectra presented in conferences and publications are inaccurate. In our polarization laboratory where rigorous polarimeter-operating procedures are in place, still about a quarter of our data is discarded as dubious and remeasured. It is recommended that all polarization measurements be approached with a degree of skepticism until the measurement system and measurement circumstances are clearly understood and appropriate tests and calibrations are provided.

15.15 DEFINITIONS

Analyzer an element whose intensity transmission is proportional to the content of a specific polarization state in the incident beam. Analyzers are placed before the detector in polarimeters. The polarization state emerging from an analyzer is not necessarily the same as the incident state being analyzed.

Birefringence the material property of having two refractive indices associated with one propagation direction. For each propagation direction within a birefringent medium, there are two modes of propagation with orthogonal polarization states and with different refractive indices n_1 and n_2 . The birefringence Δn is $|n_1 - n_2|$.

Depolarization any process which couples polarized light into partially polarized light. Depolarization is intrinsically associated with scattering and with diattenuation and retardance which vary in space, time, and/or wavelength.

Diattenuation the property of an optical element or system whereby the flux of the exiting beam depends on the polarization state of the incident beam. The transmitted intensity is a maximum P_{\max} for one incident state, and a minimum P_{\min} for the orthogonal state. The diattenuation is defined as $(P_{\max} - P_{\min}) / (P_{\max} + P_{\min})$. Polarizers have a diattenuation of one, while ideal retarders have a diattenuation of zero. Diattenuation is an essential property of analyzers.

Diattenuator any polarization element which displays diattenuation. Polarizers have a diattenuation very close to one, but nearly all optical interfaces are weak diattenuators. Examples of diattenuators include the following: polarizers and dichroic materials, as well as metal and dielectric interfaces with reflection and transmission differences described by Fresnel equations; thin films (homogeneous and isotropic); and diffraction gratings.

Eigenpolarization a polarization state transmitted unaltered by a polarization element except for a change of amplitude and phase. Every nondepolarizing polarization element has two eigenpolarizations. Any incident light not in an eigenpolarization state is transmitted in a polarization state different from the incident state. Eigenpolarizations are eigenvectors of the corresponding Mueller or Jones matrix which correspond to physical polarization states.

Ellipsometry a polarimetric technique which uses the change in the polarization state upon reflection or transmission to characterize the complex refractive index of surfaces and interfaces, and refractive indices and thicknesses of thin films.¹

Fast axis the eigenpolarization of a retarder which exits the device first. For a linear retarder, the axis is a line at a particular angle, such as 0° and 180° . For an elliptical or circular retarder, it is the corresponding elliptical polarization.

Homogeneous polarization element an element whose eigenpolarizations are orthogonal. Its eigenpolarizations are the states of maximum and minimum transmittance and also of maximum and minimum optical path length. A homogeneous element is classified as linear, circular, or elliptical depending on the form of the eigenpolarizations.

Inhomogeneous polarization element an element whose eigenpolarizations are not orthogonal. The diattenuation axis and retardance axis are not aligned. Such an element will also display different polarization characteristics for forward and backward propagating beams. The eigenpolarizations are generally not the states of maximum and minimum transmittance. Often inhomogeneous elements cannot be simply classified as linear, circular, or elliptical.

Ideal polarizer a polarizer with an intensity transmittance of one for its principal state and zero for its orthogonal state.

Linear polarizer a device which, when placed in an incident unpolarized beam, produces a beam of light whose electric field vector is oscillating primarily in one plane, with only a small component in the perpendicular plane.²

Nonpolarizing element an element which does not change the polarization state of light beams. The polarization state of the output light is equal to the polarization state of the incident light. The Jones matrix or Mueller matrix of a nonpolarizing element is proportional to the identity matrix.

Partially polarized light light containing an unpolarized component; cannot be extinguished by an ideal polarizer.

Polarimeter an optical instrument for the determination of the polarization state of a light beam, or the polarization-altering properties of a sample.

Polarimetry the science of measuring the polarization state of a light beam and the diattenuating, retarding, and depolarizing properties of materials.

Polarization (1) the polarization state of a light beam; (2) any process which alters the polarization state of a beam of light, including diattenuation, retardance, depolarization.

Polarization coupling any conversion of light from one polarization state into another state.

Polarized light light in a fixed, elliptically (including linearly or circularly) polarized state. Polarized light can be extinguished by an ideal polarizer. For polychromatic polarized light, the polarization ellipses associated with each spectral component have identical ellipticity, orientation, and helicity.

Polarizer a strongly diattenuating optical element designed to transmit light in a specified polarization state independent of the incident polarization state. The transmission of the extinguished eigenpolarization is near zero.

Polarization element any optical element used to control the polarization state of light. This includes polarizers, retarders, and depolarizers.

Pure diattenuator a diattenuator with zero retardance and no depolarization.

Pure retarder a retarder with zero diattenuation and no depolarization.

Retardance a polarization-dependent phase change associated with a polarization element or system. The phase (optical path length) of the output beam depends on the polarization state of the input beam. The optical path length (phase) of the transmitted beam is a maximum for one eigenpolarization, and a minimum for the other eigenpolarization. Other states show polarization coupling. Strictly speaking, retardance is measured in radians, but it may also be expressed equivalently as an optical path difference (length) or in fractions of a wavelength (unitless).

Retardation plate a retarder constructed from a plane parallel plate or plates of linearly birefringent material.

Retarder a polarization element designed to produce a specified phase difference between the exiting beams for two orthogonal incident polarization states (the eigenpolarizations of the element). For example, a quarter-wave linear retarder has as its eigenpolarizations two orthogonal linearly polarized states which are transmitted in their incident polarization states but with a 90° (quarter wavelength) relative phase difference (optical path length difference) introduced.

Slow axis the eigenpolarization of a retarder orthogonal to the fast axis.

Spectropolarimetry the spectroscopic study of the polarization properties of materials. Conventional spectroscopy measures the reflectance or transmission of a sample as a function of wavelength. Spectropolarimetry also measures the diattenuating, retarding, and depolarizing properties as a function of wavelength. Complete characterization of these properties is obtained by measuring the Mueller matrix of the sample as a function of wavelength.

Waveplate a retardation plate.

15.16 STOKES VECTORS AND MUELLER MATRICES

Several systematic methods of calculation have been developed for analyzing polarization, including those based on the Jones matrix, coherency matrix, and Mueller matrix.³⁻⁸ Of these methods, the Mueller calculus is most suited for describing irradiance-measuring instruments, including most

polarimeters, radiometers, and spectrometers. The Mueller calculus is primarily used in this chapter. The properties of the Mueller matrix are described in Chap. 14 “Mueller Matrices.”

In the Mueller calculus, a Stokes vector \mathbf{S} describes the polarization state of a light beam, and a Mueller matrix \mathbf{M} describes the polarization-altering characteristics of a sample. This sample may be a surface, a polarization element, an optical system, or some other light/matter interaction which produces a reflected, refracted, diffracted, or scattered light beam. Vectors and matrices are represented with bold characters.

15.17 PHENOMENOLOGICAL DEFINITION OF THE STOKES VECTOR

The Stokes vector \mathbf{S} describes the polarization state of a light beam. \mathbf{S} is defined relative to the following six flux measurements P performed with ideal polarizers in front of a radiometer.³

P_H	horizontal linear polarizer (0°)
P_V	vertical linear polarizer (90°)
P_{45}	45° linear polarizer
P_{135}	135° linear polarizer
P_R	right circular polarizer
P_L	left circular polarizer

The Stokes vector is defined as

$$\mathbf{S} = \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} = \begin{pmatrix} P_H + P_V \\ P_H - P_V \\ P_{45} - P_{135} \\ P_R - P_L \end{pmatrix} = \begin{pmatrix} I \\ Q \\ U \\ V \end{pmatrix} \quad (1)$$

where S_0 , S_1 , S_2 , and S_3 (or alternatively I , Q , U , and V) are the four Stokes vector elements or Stokes parameters. The Stokes vector does not need to be measured by these six ideal measurements; the method must reproduce the Stokes vector defined by these measurements. Ideal polarizers are not required. Further, the Stokes vector is a function of wavelength, position on the object, and the light's direction of emission or scatter. Thus, a Stokes vector measurement is an average over area, solid angle, and wavelength, as is any radiometric measurement. The Stokes vector is defined relative to a local x - y coordinate system in the plane perpendicular to the light's propagation vector, established by the polarimeter. The coordinate system is right-handed; the cross-product $\hat{x} \times \hat{y}$ of the basis vectors points in the direction of propagation of the beam.

15.18 POLARIZATION PROPERTIES OF LIGHT BEAMS

From the Stokes vector, the following polarization parameters are defined.^{6,9-11}

$$\text{Flux} \quad P = S_0 \quad (2)$$

$$\text{Degree of polarization} \quad \text{DoP} = \frac{\sqrt{S_1^2 + S_2^2 + S_3^2}}{S_0} \quad (3)$$

$$\text{Degree of linear polarization} \quad \text{DoLP} = \frac{\sqrt{S_1^2 + S_2^2}}{S_0} \quad (4)$$

$$\text{Degree of circular polarization} \quad \text{DoCP} = \frac{|S_3|}{S_0} \quad (5)$$

The Stokes vector for a partially polarized beam ($\text{DoP} < 1$) can be considered as a superposition of a completely polarized Stokes vector \mathbf{S}_p and an unpolarized Stokes vector \mathbf{S}_U .¹¹

$$\mathbf{S} = \mathbf{S}_p + \mathbf{S}_U = \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} = S_0 \text{DoP} \begin{pmatrix} 1 \\ S_1/(S_0 \text{DoP}) \\ S_2/(S_0 \text{DoP}) \\ S_3/(S_0 \text{DoP}) \end{pmatrix} + (1 - \text{DoP}) S_0 \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad (6)$$

There is no polarization element which performs this separation into polarized and unpolarized light components; a polarizer will always transmit half the unpolarized component. The polarized portion of the beam represents a net polarization ellipse traced by the electric field vector as a function of time as shown in Fig. 1. The polarization ellipse is uniquely described by four parameters. One is the phase at $t=0$, and the other three are selected from the following list: the magnitude of the semimajor axis a , semiminor axis b , orientation of the major axis ψ (azimuth of the ellipse) measured counterclockwise from the x axis, and eccentricity, and ellipticity.

$$\text{Ellipticity} \quad e = \frac{b}{a} = \frac{S_3}{S_0 + \sqrt{S_1^2 + S_2^2}} \quad (7)$$

$$\text{Orientation of major axis, azimuth} \quad \psi = \frac{1}{2} \arctan\left(\frac{S_2}{S_1}\right) \quad (8)$$

$$\text{Eccentricity} \quad \varepsilon = \sqrt{1 - e^2} \quad (9)$$

The ellipticity is the ratio of the minor to the major axis of the corresponding electric field polarization ellipse, and varies from 0 for linearly polarized light to 1 for circularly polarized light. The polarization ellipse is alternatively described by its eccentricity, which is zero for circularly polarized light, increases as the polarization ellipse becomes thinner, and is one for linearly polarized light. The polarization ellipse strictly refers to the light's electric field.

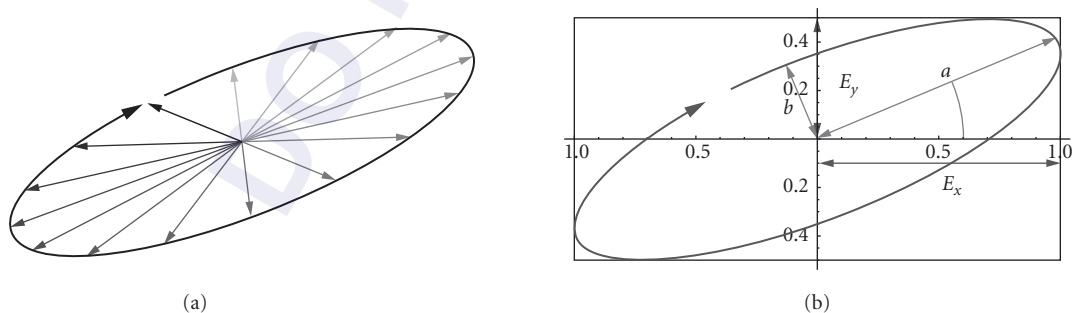


FIGURE 1 The tip of the electric field vector rotating as a function of time traces the polarization ellipse (a). The polarization ellipse parameters (b).

15.19 MUELLER MATRICES

The Mueller matrix \mathbf{M} for a polarization-altering device is defined as the matrix which transforms an incident Stokes vector \mathbf{S} into the exiting Stokes vector \mathbf{S}' ,

$$\mathbf{S}' = \begin{pmatrix} S'_0 \\ S'_1 \\ S'_2 \\ S'_3 \end{pmatrix} = \mathbf{M} \cdot \mathbf{S} = \begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \cdot \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} \quad (10)$$

The Mueller matrix is a 4×4 matrix with real-valued elements. The Mueller matrix $\mathbf{M}(\mathbf{k}, \lambda)$ for a device is always a function of the direction of propagation \mathbf{k} and wavelength λ . The Mueller matrix is an appropriate formalism for characterizing polarization measurements because it contains within its elements all of the polarization properties (diattenuation, retardance, depolarization) and their form (linear, circular, elliptical). When the Mueller matrix is known, then the exiting polarization state is known for an arbitrary incident polarization state. Chapter 14, "Mueller Matrices" contains tables of Mueller matrices for common polarization elements. Other Mueller matrix tables are found in many references including the following: Shurcliff,³ Gerrard and Burch,⁴ Azzam and Bashara,⁹ Theocaris and Gdoutos,⁵ and Goldstein.¹² The Mueller Matrix chapter contains a detailed discussion of the polarization properties and methods for calculating these properties from the Mueller matrix.

The Mueller matrix \mathbf{M} associated with a beam path through a sequence (cascade) of polarization elements $q = 1, 2, \dots, Q$ is the right-to-left product of the individual matrices \mathbf{M}_q ,

$$\mathbf{M} = \mathbf{M}_Q \cdot \mathbf{M}_{Q-1} \cdot \dots \cdot \mathbf{M}_q \cdot \dots \cdot \mathbf{M}_2 \cdot \mathbf{M}_1 = \prod_{q=Q,-1}^1 \mathbf{M}_q \quad (11)$$

15.20 DATA REDUCTION FOR LIGHT-MEASURING POLARIMETERS

This section presents a general formulation of the measurement and data-reduction procedure for a light-measuring polarimeter. The objective of Stokes polarimetry is to determine the Stokes parameters from a series of radiometric measurements. The data reduction is a linear estimation process, and lends itself to efficient solution using linear algebra, usually with a least-squares estimator to find the best match to the data. Similar developments are found in Thie,¹³ Azzam,¹⁴ and Stenflo.¹⁵

Stokes vectors and related polarization parameters for a beam are determined by measuring the flux transmitted through a set of polarization analyzers. Each analyzer determines the flux of one polarization component in the incident beam. Since a polarization analyzer does not contain ideal polarization elements, the analyzer must be calibrated, and the calibration data used in the data reduction. The polarizer in an analyzer does not need T_{\min} to equal zero; it never does, and this leakage will be corrected in the data reduction. The measured values are related to the incident Stokes vector and the analyzers by the polarimetric measurement equation. A set of linear equations, the data-reduction equations, is then solved to determine the Stokes parameters for the beam.

The *polarization analyzer* consists of the polarization elements used for analyzing the polarization state, any other optical elements (lenses, mirrors, etc.) following the analyzer, and the polarimeter's detector. The polarization effects from all elements are included in the measurement and data-reduction procedures. A polarization analyzer is characterized by an *analyzer vector* containing four elements, defined analogously to a Stokes vector. Let P_H be the flux measurement measured by the detector (the current or voltage generated) when one unit of horizontally polarized light is incident.

Similarly P_V , P_{45° , P_{135° , P_R and P_L are the detector's flux measurements for the corresponding incident polarized beams with unit flux. Then the analyzer vector \mathbf{A} is

$$\mathbf{A} = \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{pmatrix} = \begin{pmatrix} P_H + P_V \\ P_H - P_V \\ P_{45^\circ} - P_{135^\circ} \\ P_R - P_L \end{pmatrix} \quad (12)$$

Note that in the absence of noise, $P_H + P_V = P_{45^\circ} + P_{135^\circ} = P_R + P_L$. The response P of the polarization analyzer to an arbitrary polarization state \mathbf{S} is a dot product

$$P = \mathbf{A} \cdot \mathbf{S} = a_0 S_0 + a_1 S_1 + a_2 S_2 + a_3 S_3 \quad (13)$$

A Stokes vector measurement is a set of measurements acquired with a set of polarization analyzers placed into the beam. Let the total number of analyzers be Q , with each analyzer \mathbf{A}_q specified by index $q = 0, 1, \dots, Q - 1$. We assume the incident Stokes vector is the same for all polarization analyzers. The q th measurement generates an output $P_q = \mathbf{A}_q \cdot \mathbf{S}$. A polarimetric measurement matrix \mathbf{W} is defined as a $Q \times 4$ matrix with the q th row containing the analyzer vector \mathbf{A}_q ,

$$\mathbf{W} = \begin{pmatrix} a_{0,0} & a_{0,1} & a_{0,2} & a_{0,3} \\ a_{1,0} & a_{1,1} & a_{1,2} & a_{1,3} \\ \vdots & \vdots & \vdots & \vdots \\ a_{Q-1,0} & a_{Q-1,1} & a_{Q-1,2} & a_{Q-1,3} \end{pmatrix} \quad (14)$$

The Q measured fluxes are arranged in a measurement vector $\mathbf{P} = \{P_0, P_1, \dots, P_{Q-1}\}^T$. \mathbf{P} is related to \mathbf{S} by the polarimetric measurement equation

$$\mathbf{P} = \begin{pmatrix} P_0 \\ P_1 \\ \vdots \\ P_{Q-1} \end{pmatrix} = \mathbf{W} \cdot \mathbf{S} = \begin{pmatrix} a_{0,0} & a_{0,1} & a_{0,2} & a_{0,3} \\ a_{1,0} & a_{1,1} & a_{1,2} & a_{1,3} \\ \vdots & \vdots & \vdots & \vdots \\ a_{Q-1,0} & a_{Q-1,1} & a_{Q-1,2} & a_{Q-1,3} \end{pmatrix} \cdot \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} \quad (15)$$

During calibration of the polarimeter, the objective is the accurate determination of \mathbf{W} .

To calculate the Stokes vector from the data, the inverse of \mathbf{W} is determined and applied to the measured data. The measured value for the incident Stokes vector \mathbf{S}_m is related to the data by the polarimetric data-reduction matrix \mathbf{W}^{-1} ,

$$\mathbf{W}^{-1} \cdot \mathbf{P} = \mathbf{S}_m \quad (16)$$

This is the *polarimetric measurement equation*. Three considerations in the solution of this equation are the existence, rank, and uniqueness of the matrix inverse \mathbf{W}^{-1} .

The simplest case is when four measurements are performed. If $Q = 4$ and if four linearly independent analyzer vectors are used, then \mathbf{W} is of rank four, and \mathbf{W}^{-1} exists, is unique and nonsingular. Data reduction is performed by Eq. (16); the polarimeter measures all four elements of the incident Stokes vector.

When $Q > 4$, \mathbf{W} is not square and \mathbf{W}^{-1} is not unique; multiple \mathbf{W}^{-1} exist. \mathbf{S}_m is overdetermined; there are more equations than unknowns. In the absence of noise, the different \mathbf{W}^{-1} all yield the same \mathbf{S}_m . Because noise is always present, the optimum \mathbf{W}^{-1} is desired. The least squares estimate for \mathbf{S}_m utilizes a particular matrix inverse, the pseudoinverse \mathbf{W}_p^{-1} of \mathbf{W} ,

$$\mathbf{W}_p^{-1} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \quad (17)$$

The optimal estimate of \mathbf{S} is

$$\mathbf{S} = \mathbf{W}_p^{-1} \cdot \mathbf{P} = (\mathbf{W}^T \cdot \mathbf{W})^{-1} \cdot \mathbf{W}^T \cdot \mathbf{P} \quad (18)$$

When \mathbf{W} is of rank three or less, the polarimeter is *incomplete*. The optimal matrix inverse is the pseudoinverse, but only three or fewer *properties* of the Stokes vector elements are determined; the projection of the Stokes vector onto three or fewer directions is measured. If these directions align with the Stokes basis vectors, then these Stokes vector elements are measured, but in general, linear combinations of elements are measured.

15.21 SAMPLE-MEASURING POLARIMETERS FOR MEASURING MUELLER MATRIX ELEMENTS

This section contains a general formulation of Mueller matrix measurement which is an extension of the Stokes vector method of the preceding section. The Mueller matrix is always a function of wavelength, angle of incidence, and location on the sample. These are assumed fixed here for simplicity; this method can be generalized to these more general cases. Figure 2 is a block diagram of a sample-measuring polarimeter. The polarization state generator (PSG) prepares polarization states incident on a sample. The light beam exiting the sample is analyzed by the polarization state analyzer (PSA), and the flux at the detector measured.

The objective is to determine several or all of the sample's Mueller matrix \mathbf{M} elements through a sequence $q = 0, 1, \dots, Q - 1$ of polarimetric measurements. The polarization generator prepares a set of polarization states with a sequence of Stokes vectors \mathbf{S}_q . The Stokes vectors exiting the sample are $\mathbf{M} \cdot \mathbf{S}_q$. These exiting states are analyzed by the q th polarization state analyzer \mathbf{A}_q , yielding the q th measured flux $P_q = \mathbf{A}_q^T \cdot \mathbf{M} \cdot \mathbf{S}_q$. Each measured flux is assumed to be a linear function of the sample's Mueller matrix elements (nonlinear optical interactions such as frequency doubling are not treated by the present formulation). A set of linear equations is developed from the set of polarimetric measurements to solve for the Mueller matrix elements.

For example, consider a measurement taken with horizontal linear polarizers as the generator and the analyzer. As Eq. (11) shows, the measured flux only depends on the Mueller matrix elements $m_{0,0}$, $m_{0,1}$, $m_{1,0}$, and $m_{1,1}$,

$$P = \mathbf{A}^T \cdot \mathbf{M} \cdot \mathbf{S} = \frac{1}{2} (1, 1, 0, 0) \begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix} \quad (19)$$

$$= \frac{m_{0,0} + m_{0,1} + m_{1,0} + m_{1,1}}{2}$$

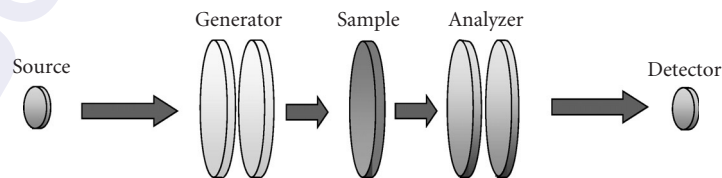


FIGURE 2 A sample-measuring polarimeter consists of a source, polarization state generator (PSG), the sample, a polarization state analyzer (PSA), and the detector. (See also color insert.)

As another example, the four Mueller matrix elements $m_{0,0}$, $m_{0,1}$, $m_{1,0}$, and $m_{1,1}$ can be measured using four measurements with ideal horizontal (**H**) and vertical (**V**) linear polarizers. Four measurements P_0 , P_1 , P_2 , and P_3 are taken with (generator/analyzer) settings of (**H/H**), (**V/H**), (**H/V**), and (**V/V**), determining the following combinations of Mueller matrix elements,

$$\begin{aligned} P_0 &= (m_{0,0} + m_{0,1} + m_{1,0} + m_{1,1})/4 & P_1 &= (m_{0,0} - m_{0,1} + m_{1,0} - m_{1,1})/4 \\ P_2 &= (m_{0,0} + m_{0,1} - m_{1,0} - m_{1,1})/4 & P_3 &= (m_{0,0} - m_{0,1} - m_{1,0} + m_{1,1})/4 \end{aligned} \quad (20)$$

These four equations are solved for the Mueller matrix elements as

$$\begin{pmatrix} m_{0,0} \\ m_{0,1} \\ m_{1,0} \\ m_{1,1} \end{pmatrix} = \begin{pmatrix} P_0 + P_1 + P_2 + P_3 \\ P_0 - P_1 + P_2 - P_3 \\ P_0 + P_1 - P_2 - P_3 \\ P_0 - P_1 - P_2 + P_3 \end{pmatrix} \quad (21)$$

Other Mueller matrix elements are determined using different combinations of generator and analyzer states. The four matrix elements at the corners of a rectangle in the Mueller matrix $\{m_{0,0}, m_{0,i}, m_{j,0}, m_{j,i}\}$ can be determined from four measurements using a $\pm i$ -generator and $\pm j$ -analyzer. For example, a pair of right and left circularly polarizing generators and a pair of 45° and 135° oriented analyzers determine elements $m_{0,0}, m_{0,3}, m_{2,0}, m_{2,3}$.

In practice, the data-reduction equations are far more complex than these examples because many more measurements are involved with nonideal polarization elements. The following section contains a systematic method for calculation of data-reduction equations based on calibration data for the generator and analyzer.

15.22 POLARIMETRIC MEASUREMENT EQUATION AND POLARIMETRIC DATA-REDUCTION EQUATION

This section develops data-reduction equations to calculate Mueller matrices from arbitrary sequences of measurements. The algorithm uses either ideal or calibrated values for the polarization generator and analyzer vectors. The data-reduction equations are straightforward matrix-vector multiplication on a data vector. This method is an extension of the data-reduction methods presented in “Data Reduction for Light-Measuring Polarimeters”.

A Mueller matrix polarimeter takes Q measurements identified by index $q = 0, 1, \dots, Q-1$. For the q th measurement, the generator produces a beam with Stokes vector S_q and the beam exiting the sample is analyzed by analyzer vector A_q . The measured flux P_q is related to the sample Mueller matrix by

$$P_q = \mathbf{A}_q^T \mathbf{M} \mathbf{S}_q = (a_{q,0} \quad a_{q,1} \quad a_{q,2} \quad a_{q,3}) \begin{pmatrix} m_{0,0} & m_{0,1} & m_{0,2} & m_{0,3} \\ m_{1,0} & m_{1,1} & m_{1,2} & m_{1,3} \\ m_{2,0} & m_{2,1} & m_{2,2} & m_{2,3} \\ m_{3,0} & m_{3,1} & m_{3,2} & m_{3,3} \end{pmatrix} \begin{pmatrix} S_{q,0} \\ S_{q,1} \\ S_{q,2} \\ S_{q,3} \end{pmatrix} \quad (22)$$

$$= \sum_{j=0}^3 \sum_{k=0}^3 a_{q,j} m_{j,k} S_{q,k}$$

This equation is rewritten as a vector-vector dot product.^{16,17} First, the Mueller matrix is flattened into a 16×1 *Mueller vector* $\vec{\mathbf{M}} = (m_{0,0}, m_{0,1}, m_{0,2}, m_{0,3}, m_{1,0}, \dots, m_{3,3})^T$. A 16×1 polarimetric measurement vector \mathbf{W}_q for the q th measurement is defined as follows:

$$\begin{aligned} \mathbf{W}_q &= (w_{q,0,0}, w_{q,0,1}, w_{q,0,2}, w_{q,0,3}, w_{q,1,0}, \dots, w_{q,3,3})^T \\ &= (a_{q,0} S_{q,0}, a_{q,0} S_{q,1}, a_{q,0} S_{q,2}, a_{q,0} S_{q,3}, a_{q,1} S_{q,0}, \dots, a_{q,3} S_{q,3})^T \end{aligned} \quad (23)$$

where $w_{q,j,k} = a_{q,j} s_{q,k}$. The q th measured flux is the dot product

$$P_q = \mathbf{W}_q \cdot \vec{\mathbf{M}} = \begin{pmatrix} a_{q,0} S_{q,0} \\ a_{q,0} S_{q,1} \\ a_{q,0} S_{q,2} \\ a_{q,0} S_{q,3} \\ a_{q,1} S_{q,0} \\ a_{q,1} S_{q,1} \\ \vdots \\ a_{q,3} S_{q,3} \end{pmatrix} \cdot \begin{pmatrix} m_{0,0} \\ m_{0,1} \\ m_{0,2} \\ m_{0,3} \\ m_{1,0} \\ m_{1,1} \\ \vdots \\ m_{3,3} \end{pmatrix} = \begin{pmatrix} w_{q,0,0} \\ w_{q,0,1} \\ w_{q,0,2} \\ w_{q,0,3} \\ w_{q,1,0} \\ w_{q,1,1} \\ \vdots \\ w_{q,3,3} \end{pmatrix} \cdot \begin{pmatrix} m_{0,0} \\ m_{0,1} \\ m_{0,2} \\ m_{0,3} \\ m_{1,0} \\ m_{1,1} \\ \vdots \\ m_{3,3} \end{pmatrix} \quad (24)$$

The full sequence of Q measurements is described by the $Q \times 16$ polarimetric measurement matrix \mathbf{W} , where the q th row is \mathbf{W}_q . The polarimetric measurement equation relates the measurement vector \mathbf{P} to the sample Mueller vector as

$$\mathbf{P} = \mathbf{W} \cdot \vec{\mathbf{M}} = \begin{pmatrix} P_0 \\ P_1 \\ \vdots \\ P_{Q-1} \end{pmatrix} = \begin{pmatrix} w_{0,0,0} & w_{0,0,1} & \dots & w_{0,3,3} \\ w_{1,0,0} & w_{1,0,1} & \dots & w_{1,3,3} \\ \vdots & \vdots & \vdots & \vdots \\ w_{Q-1,0,0} & w_{Q-1,0,1} & \dots & w_{Q-1,3,3} \end{pmatrix} \cdot \begin{pmatrix} m_{0,0} \\ m_{0,1} \\ \vdots \\ m_{3,3} \end{pmatrix} \quad (25)$$

If \mathbf{W} contains 16 linearly independent rows, all 16 elements of the Mueller matrix can be determined. When $Q = 16$, the matrix inverse is unique and the Mueller matrix elements are determined from the polarimetric data-reduction equation

$$\vec{\mathbf{M}} = \mathbf{W}^{-1} \cdot \mathbf{P} \quad (26)$$

Often, $Q > 16$, and $\vec{\mathbf{M}}$ is overdetermined. The optimal (least-squares) polarimetric data-reduction equation for $\vec{\mathbf{M}}$ uses the pseudoinverse \mathbf{W}_p^{-1} of \mathbf{W} ,

$$\vec{\mathbf{M}} = (\mathbf{W}^T \cdot \mathbf{W})^{-1} \cdot \mathbf{W}^T \cdot \mathbf{P} = \mathbf{W}_p^{-1} \cdot \mathbf{P} \quad (27)$$

The advantages of this procedure are as follows: First, this procedure does not assume that the set of states of polarization state generator and analyzer have any particular form. For example, the polarization elements in the generator and analyzer do not need to be rotated in uniform angular increments, but can comprise an arbitrary sequence. Second, the polarization elements are not assumed to be ideal polarization elements or have any particular imperfections. If the polarization generator and analyzer vectors are determined through a calibration procedure, the effects of non-ideal polarization elements are corrected in the data reduction. Third, the procedure readily treats overdetermined measurement sequences (more than 16 measurements for the full Mueller matrix), providing a least-squares solution. Finally, a matrix-vector form of data reduction is readily implemented and easily understood.

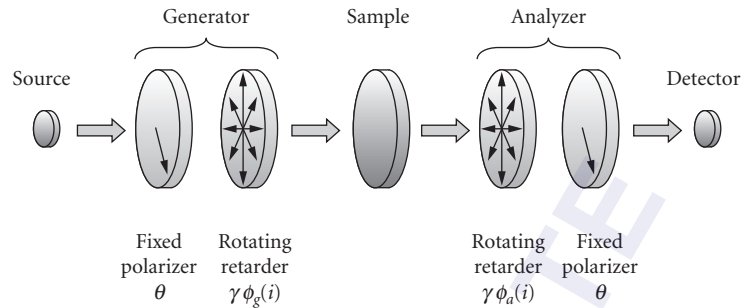


FIGURE 3 The dual rotating retarder polarimeter consists of a source, a fixed linear polarizer, a retarder which rotates in steps, the sample, a second retarder which rotates in steps, a fixed linear polarizer, and the detector. (See also color insert.) (After Ref. 18.)

15.23 DUAL ROTATING RETARDER POLARIMETER

The dual rotating retarder Mueller matrix polarimeter shown in Fig. 3¹⁸ is one of the most common Mueller polarimeters. Light from the source passes first through a fixed linear polarizer, then through a rotating linear retarder, the sample, a rotating linear retarder, and finally through a fixed linear polarizer. In the most common configuration, the polarizers are parallel, and the retarders are rotated in angular increments of five-to-one. This five-to-one ratio encodes all 16 Mueller matrix elements onto the amplitudes and phases of 12 distinct frequencies in the detected signal. This configuration was first described by Azzam¹⁶ who provides an explanation of how the ratios one-to-one, two-to-one, three-to-one, and four-to-one all yield incomplete polarimeters. Thus five-to-one is the first integer ratio yielding a complete Mueller matrix polarimeter. The data reduction can be performed using the polarimetric data-reduction matrix method of the preceding section, or alternatively the detected signal can be Fourier analyzed, and the Mueller matrix elements calculated from the Fourier coefficients.¹⁹

This polarimeter configuration has several design advantages. Since the polarizers do not move, the polarizer in the generator accepts only one polarization state from the source optics, making the measurement immune to source polarization and polarization aberrations from optics prior to the polarizer. If the polarizer did rotate, and if the beam incident on it were elliptically polarized, a systematic modulation of intensity would be introduced which would require compensation in the data reduction. Similarly, the polarizer in the analyzer does not rotate; only one polarization state is transmitted through the analyzing optics and onto the detector. Any diattenuation in the analyzing optics and any polarization sensitivity in the detector will not affect the measurements.

Optimal values for the retardances are near $2\pi/3$ rad ($\lambda/3$ waveplates).¹⁹ If $\delta_1 = \delta_2 = \pi$ rad (half-wave linear retarders), only linear states are generated and analyzed, and the last row and column of the sample Mueller matrix are not measured.

Hauge and Broch^{20,21} developed an algorithm to compensate for the linear diattenuation and linear retardance of the retarders. Goldstein and Chipman²² treat five errors, the retardances of the two retarders, and orientation errors of the two retarders and one of the polarizers, in a small angle approximation good for small errors. Chenault, Pezzaniti, and Chipman²³ extended this method to larger errors.

15.24 INCOMPLETE SAMPLE-MEASURING POLARIMETERS

Incomplete sample-measuring polarimeters do not measure the full Mueller matrix of a sample and thus provide incomplete information regarding the polarization properties of a sample. Often the full Mueller matrix is not needed. For example, many birefringent samples have considerable linear retardance and insignificant amounts of other polarization forms. The magnitude of the retardance

can be measured, assuming all the other polarization effects are small, using much simpler configurations than a Mueller matrix polarimeter, such as the circular polariscope.⁵ Similarly, homogeneous and isotropic interfaces, such as dielectrics, metals, and thin films, should only display linear diattenuation and linear retardance aligned with the s - p planes. These interfaces do not need characterization of their circular diattenuation and circular retardance. So most ellipsometers characterize the diattenuation and retardance associated with s and p without providing the full Mueller matrix.^{1,6,9}

15.25 NONIDEAL POLARIZATION ELEMENTS

Polarization elements used in polarimetry require a level of characterization beyond what is normally provided by vendors at the time of this writing. For retarders, vendors usually only specify the linear retardance. For polarizers, usually only the two principal transmittances or the extinction ratio is given. For polarization critical applications, this is inadequate. In the following sections, common defects of polarization elements are described. The Mueller calculus is recommended as an appropriate means of describing complex behaviors and shortcomings.

For ideal polarization elements, the polarization properties are readily defined. For real polarization elements, the precise description is more complex. Polarization elements such as polarizers, retarders, and depolarizers have three general polarization properties: diattenuation, retardance, and depolarization; a typical element displays some amount of all three. Diattenuation occurs when the intensity transmittance is a function of the incident polarization state.²⁴ The diattenuation D is defined in terms of the maximum T_{\max} and minimum T_{\min} intensity transmittances, as

$$D = \frac{T_{\max} - T_{\min}}{T_{\max} + T_{\min}} \quad (28)$$

For an ideal polarizer, $D = 1$. When $D = 0$, all incident polarization states are transmitted with equal attenuation. The quality of a polarizer is often expressed in terms of the related quantity, the extinction ratio E ,

$$E = \frac{T_{\max}}{T_{\min}} = \frac{1 + D}{1 - D} \quad (29)$$

Retardance is the phase difference a device introduces between its eigenpolarizations (eigenstates). For a birefringent retarder with refractive indices n_1 and n_2 and thickness t , the retardance δ expressed in radians is

$$\delta = \frac{2\pi |n_1 - n_2| t}{\lambda} \quad (30)$$

Depolarization describes the coupling by a device of incident polarized light into depolarized light in the exiting beam. Depolarization occurs when light transmits through milk or scatters from clouds. Multimode optical fibers generally depolarize the light. Depolarization is intrinsically associated with scattering and a loss of coherence of the polarization state. A small amount of depolarization is associated with the scattered light from all optical components. The depolarization varies as a function of the incident polarization state.²⁵

15.26 ELLIPTICAL AND CIRCULAR POLARIZERS AND ANALYZERS

There are few good and convenient circularly or elliptically polarizing mechanisms, whereas linear polarizers are simple, inexpensive, and of high quality. Therefore, most circular and elliptical polarizers incorporate linear polarizers to perform the polarizing, and retarders to convert between

polarization states. For such compound devices, the distinction between a polarizer and an analyzer is significant. This is illustrated by three examples: (1) a left circular polarizer (and horizontal linear analyzer) constructed from a horizontal linear polarizer $\text{LP}(0^\circ)$ followed by a quarter-wave linear retarder with the fast axis oriented at 135° , $\text{QWLR}(135^\circ)$ Eq. (31), (2) a left circular analyzer (and horizontal linear polarizer) constructed from a $\text{QWLR}(45^\circ)$ followed by a horizontal linear polarizer $\text{LP}(0^\circ)$ Eq. (32), and (3) a homogeneous left circular polarizer (and left circular analyzer) constructed from a $\text{QWLR}(135^\circ)$, then an $\text{LP}(0^\circ)$, followed by a $\text{QWLR}(45^\circ)$ Eq. (33). The three Mueller matrix equations and the exiting polarization states for arbitrary incident states are as follows:

$$\begin{aligned} \text{QWLR}(135^\circ) \cdot \text{LP}(0^\circ) \cdot \mathbf{S} &= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & -1 & 0 & 0 \end{pmatrix} \cdot \frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & -1 & 0 & 0 \end{pmatrix} \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} S_0 + S_1 \\ 0 \\ 0 \\ -S_0 - S_1 \end{pmatrix} \end{aligned} \quad (31)$$

$$\begin{aligned} \text{LP}(0^\circ) \cdot \text{QWLR}(45^\circ) \cdot \mathbf{S} &= \frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} 1 & 0 & 0 & -1 \\ 1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} S_0 - S_3 \\ S_0 - S_3 \\ 0 \\ 0 \end{pmatrix} \end{aligned} \quad (32)$$

$$\text{QWLR}(135^\circ) \cdot \text{LP}(0^\circ) \cdot \text{QWLR}(45^\circ) \cdot \mathbf{S}$$

$$\begin{aligned} &= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & -1 & 0 & 0 \end{pmatrix} \frac{1}{2} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \frac{1}{2} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} 1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} S_0 \\ S_1 \\ S_2 \\ S_3 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} S_0 - S_3 \\ 0 \\ 0 \\ -S_0 + S_3 \end{pmatrix} \end{aligned} \quad (33)$$

The device in Eq. (31) transmits only left circularly polarized light, because the S_0 and S_3 have equal magnitude and opposite sign; thus it is a left circular polarizer. However, the transmitted flux $(S_0 + S_1)/2$ is the flux of horizontal linearly polarized light in the incident beam, making it a horizontal linear analyzer.

Similarly, the transmitted flux in Eq. (32), $(S_0 - S_3)/2$, is the flux of incident left circularly polarized light, making this combination a left circular analyzer. The final polarizer makes the Eq. (32) device a horizontal linear polarizer, although this is not the standard Mueller matrix for horizontal linear polarizers found in tables. Thus an analyzer for a state does not necessarily transmit that state;

its transmitted flux is proportional to the amount of the analyzed state in the incident beam. The examples in Eqs. (31) and (32) are inhomogeneous polarization elements because the eigenpolarizations are not orthogonal. Equation (31) has left circular and vertically polarized eigenpolarizations. Equation (32) has horizontal and right circularly polarized eigenpolarizations. The characteristics of both devices are different for propagation in opposite directions.

The Eq. (33) device is both a left circular polarizer and a left circular analyzer; it has the same characteristics for propagation in opposite directions. The eigenpolarizations are orthogonal, left and right circularly polarized, so this device is a homogeneous left circular polarizer. This is the left circular polarizer Mueller matrix commonly found in tables; however, it is not the most common implementation.

15.27 COMMON DEFECTS OF POLARIZATION ELEMENTS

Here we list some common defects found in real polarization elements.

1. Polarizers have nonideal diattenuation since $T_{\min} > 0$ and also nonideal transmission since $T_{\max} < 1$.^{2,25,26}
2. Retarders have the incorrect retardance. Thus, there will be some deviation from a quarter-wave or a half-wave of retardance, for example, because of fabrication errors or a change in wavelength or temperature.
3. Retarders usually have some diattenuation. This may occur due to differences in the absorption coefficients (dichroism). Birefringent retarders have diattenuation due to the difference of the Fresnel coefficients at normal incidence for the two eigenpolarizations since $n_1 \neq n_2$. This can be reduced by antireflection coatings.
4. Polarizers usually have some retardance; there is an optical path length difference between the transmitted (principal) eigenpolarization and the small amount leaked of the secondary eigenpolarization. For example, sheet polarizers and wire-grid polarizers show substantial retardance when the secondary state is not completely extinguished.
5. The polarization properties vary with angle of incidence; for example, Glan-Thompson polarizers polarize over only a 4° field of view.² Birefringent retarders commonly show a quadratic variation of retardance with angle of incidence; the retardance increases along one axis and decreases along the orthogonal axis.^{27,28} For polarizing beam-splitter cubes, the axis of the transmitted linear polarization rotates when the light is incident from outside of the normal plane (the plane of the face normals and the beam-splitting interface normal).
6. The polarization properties vary with wavelength; for simple retarders made from a single birefringent plate, the retardance varies approximately as 1/wavelength. Other components have more complex dependence.
7. For polarizers, the *analyzed state* and the *transmitted state* can be different. Consider a polarizing device formed from a linear polarizer oriented at 0° followed by a linear polarizer oriented at 2°. Incident light linearly polarized at 0° has the highest transmittance for all possible polarization states and is the analyzed state. The exiting beam is linearly polarized at 2°, the only state exiting the device. The transmitted state is an eigenpolarization; the analyzed state isn't. This *rotation* between the analyzed and transmitted states of a polarizer frequently occurs, for example, when the crystal axes of a birefringent polarizing prism assembly, such as a Glan-Thompson polarizer, are misaligned.
8. A nominally "linear" element may be slightly elliptical (have elliptical eigenpolarizations). For example, a crystal quartz linear retarder waveplate with the optical crystal axis not quite parallel with the surface becomes an elliptical retarder due to quartz's optical activity. A compound waveplate with two or more birefringent crystals whose fast axes are not exactly 0° or 90°

apart is a retarder with slightly elliptical eigenpolarizations. Similarly a circular element may be slightly elliptical. For example, an inhomogeneous circular polarizer formed from a linear polarizer followed by a quarter-wave linear retarder at 45° [see Eq. (31)] becomes an elliptical polarizer as the retarder's fast axis is rotated or as the retardance changes with wavelength.

9. The eigenpolarizations of a polarization element may not be orthogonal; that is, a polarizer may transmit linearly polarized light at 0° without change of polarization while extinguishing linearly polarized light oriented at 88° . Such a polarization element is referred to as *inhomogeneous*.^{3,29} Sequences of polarization elements, such as optical isolator assemblies, often are inhomogeneous. The circular polarizer in Eq. (31) is inhomogeneous.
10. A polarization element may depolarize, coupling polarized light into an unpolarized component. A polarizer or retarder with a small amount of depolarization, when illuminated by a completely polarized beam, will have a small amount of unpolarized light in the transmitted beam. Such a transmitted beam can no longer be extinguished by an ideal polarizer. Depolarization results from fabrication errors such as surface roughness, bulk scattering, random strains and dislocations, and thin-film microstructure. Pinholes in a polarizer allow unpolarized light into the transmitted beam.
11. Multiply reflected beams and other "secondary" beams may be present with undesired polarization properties. For example, the multiply reflected beams inside a birefringent waveplate have various retardances. Antireflection coatings will reduce this effect in one waveband, but may increase these problems with multiple reflections in other wavebands.

The preceding list is by no means comprehensive but should serve as a warning to those with demanding polarization element applications. In particular, the performance of polarizing beam-splitting cubes and of liquid crystal cells have been found to be far from ideal.³⁰

15.28 POLARIZATION MODULATORS, RETARDANCE MODULATORS

Time-sequential polarimeters require rapid variation of the polarization state in a controlled manner. This section reviews the principal polarization modulation technologies. Several varieties of retardance modulators are in widespread use. The only common diattenuation modulator is the spinning polarizer.

Variable retarders generally have either a fixed retardance with variable axis (i.e., motor driven rotating waveplate), or a variable retardance with a fixed axis (i.e., liquid crystal retarder, electro-optical modulator, or photoelastic modulator).

In a retardance modulator, at least one of the two (often degenerate or equal) modes' refractive indices change. If polarized light is launched into the mode with varying refractive index, a phase modulator results. For polarization modulation, the incident state must be in a combination of the modes, usually equally distributed between the two modes. An amplitude (intensity) modulator is produced by placing an additional polarizer after a polarization modulator, oriented between the two modes.

15.29 ROTATING RETARDERS

Retarders with fast axes rotated by rotary stages are the gold standard for accurate polarimetry. Crystal and polymer retarders are fabricated to high accuracy and uniformity. Rotary stages can locate the fast axis angle to one arc second or better. Alternatively retarders can be continuously rotated, usually using DC brushless motors, with high uniformity and repeatability. Smith¹⁹ discusses the optimization of rotating retarder Mueller matrix polarimeters.

The disadvantages of rotating retarders are the size of the motors, their cost, and the relatively low rotation speeds (less than 1000 revolutions per second).

15.30 PHOTO-ELASTIC MODULATORS

Photo-elastic modulators (PEMs) use oscillating stress birefringence in a resonant crystal driven by a sound wave. An isotropic optical material such as glass becomes birefringent when compressed along one axis.^{31,32} This is stress-induced birefringence, or the photoelastic effect. A variable retarder can be constructed by compressing glass, but a large amount of power is needed to slowly modulate stress-induced birefringence. PEMs use a mechanically resonant bar with a high mechanical quality factor Q of 10^3 to 10^4 . A piezoelectric transducer (PZT) is coupled to the glass or fused silica bar, and a standing sound wave that oscillates at the bar's fundamental frequency is induced, causing a rapid sinusoidal modulation of the birefringence. This reduces power requirement for a quarter wave or half wave of retardance to less than 0.5 W.³³ The positive and negative parts of the sine correspond to retardance fast axes 90° apart.

PEMs have been in use for over 25 years as a method of polarization modulation in a variety of research and industrial applications. The principal supplier of PEMs is Hinds Instruments (Hillsboro, Oregon). The benefits of PEMs include low operating voltages, large apertures, and wide angular acceptance.^{34,35} Because PEMs are constructed from glass, fused silica, and other transparent materials, transmittance over a wide spectral range is straightforward. Polarimetric sensitivities (i.e., precision) of about 3 parts in 10^6 have been obtained for solar astronomy applications.^{36,37,38}

A single PEM oscillating between 0° and 90° fast axes in front of a polarizer can measure S_0 , S_2 , and S_3 , but not S_1 . Complete Stokes vector measurement requires two PEMs, optimally 45° apart. Complete Mueller matrix measurement requires two PEMs in the generator with axes nominally 45° apart and two PEMs in the analyzer nominally 45° apart. Another common configuration measures S_0 , S_1 , and S_2 but not S_3 by placing a PEM between two quarter-wave linear retarders whose axes are at $\pm 45^\circ$ to the PEM axis. When linearly polarized light is incident, linearly polarized light exits with a rapidly modulated orientation. This retarder/PEM/retarder assembly operates as a circular retardance modulator.

Typical PEM frequencies are in the tens of kHz for glass elements several centimeters in size. Smaller elements have higher frequencies and larger elements lower frequencies. The instantaneous retardance is spatially nonuniform, varying as a half period cosine across the aperture, thus varying quadratically about the center of the aperture.

PEMs which modulate at frequencies suitable for interfacing with cameras are impractically large, so PEMs are used almost exclusively with single channel detectors acquiring hundreds of thousands of measurements per second.

Because of the PEM's high Q and extremely stable frequency operation, they excel at the measurement of low birefringence in glass, such as for glass for liquid crystal (LC) cells.^{39,40,41}

15.31 LIQUID CRYSTAL RETARDERS

Two types of liquid crystal cells for polarization modulation are in widespread use: untwisted nematic cells and ferroelectric cells.

Untwisted nematic liquid crystal cells, Fredericksz cells, are available as polarization modulators. These liquid crystal variable retarders (LCVRs) are electrically tunable waveplates with retardance in the range of zero to several waves. The Fredericksz cell configuration is different from the twisted nematic configuration typically used in liquid crystal displays. There are four components (Fig. 4): two glass plates which form a cavity, indium tin oxide transparent electrodes coating on the outside of the plates, a polyimide layer on the inside of each plate which acts to align the liquid crystal molecules parallel to the plates, and a high birefringence liquid crystal material sandwiched between the plates. When no voltage is applied, the liquid crystal molecules' directions are aligned in one direction parallel to the plates and the retardance is at a maximum. When a voltage differential is applied between the plates, an electric field is induced which supplies a torque to the liquid crystal molecules, increasing the angle of the molecules with respect to the plates; the retardance is decreased, as shown in Fig. 5. When the majority of the molecules are nearly perpendicular to the

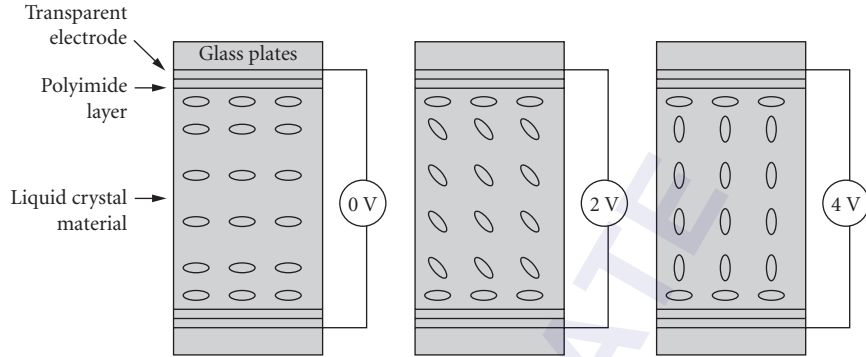


FIGURE 4 Liquid crystal cell composition and liquid crystal orientation as a function of applied voltage, varying from maximum retardance (left) to minimum (right).

plates, retardance is at a minimum. The molecules near the plates are unable to fully rotate, and so the retardance doesn't quite get to zero.

If a DC electric field is applied continuously, impurity ions in the liquid crystal material migrate toward the plates and may damage the liquid crystal structure. Once at the plates, the ions create a permanent electric field which reduces the dynamic range of the device. To avoid this problem an alternating square wave voltage at approximately 1 kHz is applied.

Retardance is a nonlinear function of applied voltage, and the relationship is highly variable from device to device, so that individual calibration of each device is required. Although optical quality and surface figure are generally very good, spatial uniformity is often poor due to the difficulty of obtaining an even distribution of the liquid crystal material between the plates. Transmission of LC

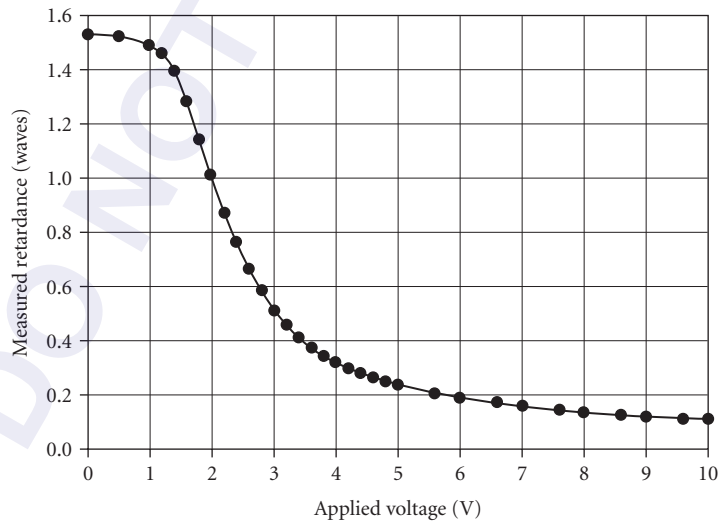


FIGURE 5 Typical retardance vs. voltage curve for a liquid crystal variable retarder.

cells without polarizers is fairly high (typically 80 to 90 percent), but can vary with applied voltage. Temperature dependence is typically 0.5 percent per °C, a significant issue. The variation of retardance with angle of incidence is large, typically 2° to 4° of retardance per degree angle of incidence! Also the liquid crystal material scatters, causing some depolarization.

Switching time in LCVRs is highly variable, depending on cell thickness, liquid crystal viscosity, temperature, and applied voltage. Switching time is asymmetric: when voltage and thus electric field is increased, the torque applied to the molecules determines the response time, but when voltage is decreased a slower mechanical restoring force bringing the molecules to equilibrium determines response time. LCVR switching time is measured by placing the LC cell between two polarizers with its axis at 45°, sending light through the assembly onto a photodiode, then modulating the cell from low to high voltage amplitude with a sinusoidal or square wave. The detected signal is observed with an oscilloscope while sweeping the signal frequency. The onset of hysteresis indicates the switching time. LCVR switching times are typically on the order of 5 to 100 ms. Vendor specifications are often misleading and should be independently tested.

LCVRs exhibit significant nonuniformity in polarization parameters across the clear aperture, with 5 percent typical. Retardance variation has many causes: cell thickness variations, temperature nonuniformity, variation of surface charge on the electrodes, and nonuniform squeezing when charge is applied. The depolarization index is usually significant, typically between 1 and 10 percent. Depolarization has several causes: bulk scattering, the glass spacer balls, and a small high-frequency oscillation of the LC molecules in response to the kHz square wave drive voltage. The large variation in uniformity and depolarization index observed among individual lots of LCVRs results from the handmade quality of most cells (Refs. 42 and 43, Chap. 7).

Fredericksz cells are small and inexpensive relative to the other retardance modulators: rotating retarders, electro-optical modulators, PEMs, and magneto-optical modulators. Thus LCVRs appear to be the ideal modulators for most applications. But they have many difficult and nonideal characteristics. Thus they are relegated to mainly qualitative applications, such as intensity modulation or low-accuracy polarization state control. It is likely that more polarimeter development projects based on LC cells have failed, than those using any other polarization modulation technology. Significant time and resources are necessary to develop accurate LC-based polarimetry.

15.32 ELECTRO-OPTICAL MODULATORS

Electro-optic modulators use the electric field across the modulating material to induce retardance. The two principal mechanisms are the Pockels effect and the Kerr effect. The electric field is generally produced by placing the crystal within a capacitor. The electro-optic effects are relatively weak so the modulator crystal aperture is generally small, the path length long, and the associated voltages large, hundreds or thousands of volts. Modulation speeds can be very high, in the hundreds of megahertz, or when operating in waveguides, modulation can be produced in the tens of gigahertz. Lithium niobate, potassium dihydrogen phosphate (KDP), and ammonium dihydrogen phosphate (ADP) are common electro-optic modulator materials.

15.33 MAGNETO-OPTICAL MODULATORS

Circular retardance modulators using the magneto-optical effect are produced by placing a high Verdet coefficient material, such as yttrium iron garnet, in a solenoid and varying the magnetic field. Large apertures are easily achieved. High currents are required and switching times are fairly slow.

Magneto-optic materials are primarily used in Faraday isolators, which allow light to pass in one direction and block the counter-propagating light. The corresponding magneto-optical modulators have not been widely commercialized.

15.34 FIBER SQUEEZERS

For fiber optic polarimetry, fiber squeezers are a fast, economical, widely deployed retardance modulator.^{44,45} When a fiber is squeezed mechanically, retardance is introduced both because the core becomes elliptical and due to the stress-optic effect. Retardance is linearly proportional to the force applied. Piezoelectric transducers can modulate fiber squeezer polarization at rates up to 30 KHz with low insertion loss.

The polarization state through long fibers, such as fiber communication links between cities, tends to drift as a function of time, quickly if the fibers are moved. One important application of fiber squeezers is to maintain the exiting polarization in a fixed state.

General Photonics (Chino, California) is a leading supplier of fiber squeezers and associated polarimeters, polarization mode dispersion controllers, depolarizers, and other fiber squeezer-based devices.

15.35 POLARIMETER DESIGN METRICS

Several methods have been developed for evaluating the suitability of a polarimeter configuration for Stokes or Mueller matrix measurement. Such methods are needed to select the sets of generators and analyzer states, determine optimum values for retarders and rotation angles, and obtain a deeper understanding of how the polarization parameters will be measured by a particular polarimeter. The following development closely follows Twietmeyer⁴³ and Twietmeyer and Chipman.⁴⁶

The rank and null space of the polarimetric measurement matrix \mathbf{W} identifies a polarimeter as complete or incomplete. The rank of \mathbf{W} should be four for a complete Stokes polarimeter and 16 for a complete Mueller matrix polarimeter. Any polarization state which lies partially or wholly in the null space of \mathbf{W} cannot be measured. A complete polarimeter has no null space. When \mathbf{M} has components in the null space, the data reduction returns a nearby reconstruction in the range of \mathbf{W} .

Each row of \mathbf{W} forms one basis vector in the reconstruction of \mathbf{M} , that is, the measured intensity at each polarimeter state is the projection of \mathbf{M} onto the corresponding basis vector. For an effective reconstruction, there should be minimum correlation between basis vectors; they should be linearly independent, widely distributed, and well balanced in magnitude. For an overspecified system with $Q > 16$, the basis vectors provide redundant coverage of the polarization space, improving performance in the presence of noise. Basis states may be chosen to lie more densely in directions where most information about \mathbf{M} is desired. For example, polarimeters to measure stress birefringence are most interested in linear retardance, so the basis states can be selected to improve the signal to noise on those parameters at the expense of diattenuation and depolarization accuracy.

For a general purpose polarimeter which measures a wide variety of arbitrary \mathbf{M} , the polarimetric measurement matrix should be as far from singular as possible; it should be *well conditioned*. Various linear algebra metrics quantify this distance from singular. The most widely used is κ_p , the condition number based on the L_p norm of the matrix \mathbf{W} , defined as⁴⁷

$$\kappa_p(\mathbf{W}) = \|\mathbf{W}\|_p \|\mathbf{W}^{-1}\|_p \quad (34)$$

where the bars signify the p -norm

$$\|\mathbf{W}\|_p = \sup_{\mathbf{x} \in D(\mathbf{W})} \frac{\|\mathbf{W} \cdot \mathbf{x}\|_p}{\|\mathbf{x}\|_p} \quad (\text{matrix } p\text{-norm}) \quad (\|\mathbf{x}\|_p)^p = \sum_i x_i^p \quad (\text{vector } p\text{-norm}) \quad (35)$$

and where \mathbf{x} is a vector, $D(\mathbf{W})$ is the domain of \mathbf{W} , and sup is the supremum (limiting maximum value). Minimization of the condition number of \mathbf{W} is a standard optimization method for polarimeters.

Four different condition number definitions are in general use: the L_1 condition number ($p = 1$) based on the maximum absolute column sum; the L_∞ condition number ($p = \infty$) based on the maximum absolute row sum; the L_2 condition number ($p = 2$) based on the euclidean length of the rows of \mathbf{W} ; and the frobenius norm (applicable where \mathbf{W} is square and invertible) based on the determinant of \mathbf{W} . Though the various condition numbers differ for a given matrix, they are similarly bounded,⁴⁷ and so provide equivalent utility. In polarimetry the L_2 condition number is preferred. The range of the L_2 condition number varies from 1 (perfect conditioning, the identity matrix) to infinity (singular matrices).

Further insight into the conditioning of \mathbf{W} is obtained from its singular value decomposition (SVD) which was introduced to polarimeter design by Tyo⁴⁸ and Sabatke et al.⁴⁹ The SVD factors any $N \times K$ matrix \mathbf{W} as

$$\mathbf{W} = \mathbf{U} \cdot \mathbf{D} \cdot \mathbf{V}^T = \mathbf{U} \cdot \begin{pmatrix} \mu_1 & & & & \\ & \mu_2 & & & \\ & & \ddots & & \\ & & & \mu_{K-1} & \\ & & & & \mu_K \\ 0 & 0 & \cdots & 0 & 0 \\ & & & \vdots & \end{pmatrix} \cdot \mathbf{V}^T \quad (36)$$

where \mathbf{U} and \mathbf{V} are $N \times N$ and $K \times K$ unitary matrices, and \mathbf{D} is an $N \times K$ diagonal matrix. The diagonal elements μ_k are the singular values. The rank of \mathbf{W} is the number of nonzero singular values. Those columns of \mathbf{U} associated with nonzero singular values form an orthonormal basis for the range of \mathbf{W} ; those columns of \mathbf{V} associated with zero-valued singular values form an orthonormal basis for the null space of \mathbf{W} . The columns of \mathbf{V} associated with nonzero singular values form an orthonormal basis which spans the full vector space of \mathbf{W} and thus reconstructs \mathbf{M} . Each singular value gives the relative strength of the corresponding vector in this basis set, and the columns of \mathbf{U} form a mapping from the \mathbf{V} basis set back to the original basis set of \mathbf{W} . Further, since

$$\mathbf{P} = \mathbf{W} \cdot \mathbf{M} = \mathbf{U} \cdot \begin{pmatrix} u_1 & 0 & 0 & 0 \\ 0 & u_2 & 0 & 0 \\ & & \ddots & \\ 0 & 0 & 0 & u_{16} \\ & & & \vdots \end{pmatrix} \cdot \mathbf{V}^T \cdot \vec{\mathbf{M}} \quad (37)$$

the rows of \mathbf{U} corresponding to zero-valued singular values describe sets of flux measurements which are not generated by any Mueller matrix, so their presence in a polarimetric measurement can only be due to noise. Based on this interpretation, any basis vector in \mathbf{V} which is associated with a relatively small singular value is near the null space and likely has little information content; such small singular values predominantly amplify noise into the reconstruction of \mathbf{M} . Error sources which produce projections (flux vectors) which are similar to the flux vectors generated by the basis vectors in \mathbf{V} (particularly those which correspond to large singular values) will couple strongly into the reconstruction of \mathbf{M} . The L_2 condition number is equal to the ratio of the largest to smallest singular values,⁵⁰ and thus minimizing the condition number is equivalent to equalizing, to the extent possible, the range of singular values so that the basis vectors have wide distribution and similar weight.

For a four measurement Stokes polarimeter, the Stokes vectors representing each of the four analyzer states, when plotted on the Poincaré sphere, define a tetrahedron which is generally irregular. The volume of the tetrahedron is proportional to the determinant of \mathbf{W} , and is maximized when the vertices form a regular tetrahedron. In this case the maximum distance from a vertex to any point on the sphere is minimized, and the condition number is also at a minimum.

15.36 SINGULAR VALUE DECOMPOSITION EXAMPLES

Two examples of the application of the condition number to Mueller matrix polarimeters follow: the first is an example of an optimum polarimeter, the second is nearly singular. Consider a Mueller matrix polarimeter which uses four generator states $\mathbf{V}_1, \mathbf{V}_2, \mathbf{V}_3, \mathbf{V}_4$, located at the vertices of a regular tetrahedron on the Poincaré sphere, with associated Stokes vectors

$$\mathbf{V}_1 = \{1, 1, 0, 0\} \quad \mathbf{V}_2 = \left\{1, -\frac{1}{3}, \frac{2\sqrt{2}}{3}, 0\right\} \quad \mathbf{V}_3 = \left\{1, -\frac{1}{3}, \frac{-\sqrt{2}}{3}, \sqrt{\frac{2}{3}}\right\} \quad \mathbf{V}_4 = \left\{1, -\frac{1}{3}, -\frac{\sqrt{2}}{3}, -\sqrt{\frac{2}{3}}\right\} \quad (38)$$

The analyzer states are also chosen as $\mathbf{V}_1, \mathbf{V}_2, \mathbf{V}_3, \mathbf{V}_4$. Sixteen measurements are acquired at each of the combinations of generator and analyzer. This polarimeter is one member of the set of 16-measurement Mueller matrix polarimeters with minimum condition number, so it can be considered an optimum configuration. The corresponding 16 singular values are

$$\left\{4, \frac{4\sqrt{3}}{3}, \frac{4\sqrt{3}}{3}, \frac{4\sqrt{3}}{3}, \frac{4\sqrt{3}}{3}, \frac{4\sqrt{3}}{3}, \frac{4\sqrt{3}}{3}, \frac{4}{3}, \frac{4}{3}, \frac{4}{3}, \frac{4}{3}, \frac{4}{3}, \frac{4}{3}, \frac{4}{3}, \frac{4}{3}\right\} \quad (39)$$

and the condition number, equal to the quotient of the first and last singular values, is 3. Each of the 16 columns of \mathbf{U} represents a different orthogonal component used to reconstruct a measured Mueller matrix. In the presence of white noise, the Mueller matrix component corresponding to the first column will be measured with the highest signal to noise, about $\sqrt{3}$ times better than the next six components (columns) from \mathbf{U} , and about three times better than the last nine Mueller matrix components.

As an example of a polarimeter with a nearly singular polarimetric measurement matrix, the second row of \mathbf{W} , (generate \mathbf{V}_1 , analyze \mathbf{V}_2),

$$\{1.0, -0.333, 0.943, 0, 1.0, -0.333, 0.943, 0, 0, 0, 0, 0, 0, 0, 0, 0\} \quad (40)$$

will be replaced with a vector

$$\{1, 1.0, 0.0005, 0, 1.0, 1.0, 0.0005, 0, 0, 0, 0, 0, 0, 0, 0, 0\} \quad (41)$$

nearly equal to the first row of \mathbf{W} ,

$$\{1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0\} \quad (42)$$

so that these two rows are nearly linearly dependent. Examining the resulting singular values,

$$\{4.0557, 2.7273, 2.3094, 2.3094, 2.3094, 2.3094, 2.0221, 1.4988, 1.3333, 1.3333, 1.3333, 1.3333, 1.3333, 1.3333, 1.3333, 0.0004\} \quad (43)$$

the last singular value is close to zero and the condition number is about 10,000. Whenever the measured flux contains the pattern corresponding to the last row of \mathbf{V}^T , this component will be amplified by about 10,000 during the data reduction relative to the other 15 components of the Mueller

matrix and will usually dominate the measurement. In the presence of random noise, the measured Mueller matrix will be close to the 16th column of \mathbf{U} (partitioned into a 4×4 “Mueller matrix”), and so the measurement will be inaccurate.

In summary, components corresponding to very small singular values are greatly amplified in the matrix inverse and can overwhelm the remainder of the Mueller matrix in the polarimetric data reduction.

15.37 POLARIMETER ERROR ANALYSIS

When operating a polarimeter, \mathbf{W} is not known exactly and may have changed since calibration. For example, a rotating retarder may have inconsistent orientation, rays may take different paths through the polarimeter for different samples, or the spectral distribution of the measured light may vary. Measurement error is present due to detector noise and source fluctuation. Eqs. (18) and (27) may be modified to include these effects as follows:

$$(\mathbf{W} + \delta\mathbf{W}) \cdot \mathbf{M} + \delta\mathbf{P} = \mathbf{P}_M \quad (44)$$

where $\delta\mathbf{W}$ is an $N \times 16$ matrix representing the difference between the actual and calibrated \mathbf{W} , $\delta\mathbf{P}$ is an $N \times 1$ vector representing intensity measurement error, and \mathbf{P}_M is the $N \times 1$ vector of fluxes measured in the presence of error. \mathbf{M}_R , the polarimeter’s estimate of \mathbf{M} , is then calculated using the calibration data as

$$\begin{aligned} \mathbf{M}_R &= \mathbf{W}_p^{-1} \cdot \mathbf{P}_M \\ &= \mathbf{W}_p^{-1} \cdot ((\mathbf{W} + \delta\mathbf{W}) \cdot \mathbf{M} + \delta\mathbf{P}) \\ &= \mathbf{M} + \delta\mathbf{M} = \mathbf{M} + \mathbf{W}_p^{-1} \cdot (\delta\mathbf{W} \cdot \mathbf{M} + \delta\mathbf{P}) \end{aligned} \quad (45)$$

where $\delta\mathbf{M}$ is the difference between the measured \mathbf{M}_R and the actual \mathbf{M} . There are two error terms in $\delta\mathbf{M}$, one dependent on $\delta\mathbf{W}$ and \mathbf{M} , and the other on $\delta\mathbf{P}$.

Small errors may be described by a first-order Taylor expansion. The error for the j th component of the i th polarimeter state, having R variables x_r which may be subject to error (such as a retardance magnitude), each with nominal value ϕ_r , and error magnitude δ_r , is

$$\delta\mathbf{W}_{i,j} \approx \sum_{r=1}^R \delta_r \frac{\partial \mathbf{W}_{i,j}}{\partial x_r} \Big|_{x_r=\phi_r} \quad (46)$$

The error in the fluxes \mathbf{P} is assumed independent of the polarimeter elements, and is given by

$$\delta\mathbf{P}_i = \varepsilon_i \quad (47)$$

where ε_i is the error in the i th intensity measurement. The error in reconstructing each of the $k = 1, \dots, 16$ elements of \mathbf{M} in terms of the errors in the instrument and detection process is then

$$\delta\mathbf{M}_k = \sum_{i=1}^N \mathbf{W}_{k,i}^{-1} \cdot \left[\sum_{j=1}^{16} \sum_{r=1}^R \delta_r \frac{\partial \mathbf{W}_{i,j}}{\partial x_r} \Big|_{x_r=\phi_r} \cdot \mathbf{M}_j \right] + \sum_{i=1}^N \mathbf{W}_{k,i}^{-1} \varepsilon_i \quad k=1, \dots, 16 \quad (48)$$

The mean (expectation) and standard deviation (SD) of the error [$\langle \delta \mathbf{M} \rangle$ and $\text{SD}(\delta \mathbf{M})$] may be estimated when the polarimeter state variables and the statistics of the error sources are approximately known.

The error due to a known systematic (nonzero mean error) source may be compensated by estimating $\delta \mathbf{W}$ (e.g., using ideal Mueller matrices to model the polarimeter) and then forming a new polarimetric measurement matrix $\mathbf{W}_q = \mathbf{W} + \langle \delta \mathbf{W} \rangle$. For example, when using a liquid crystal retarder with a known profile of retardance magnitude as a function of temperature, a new \mathbf{W}_q may be recalculated at every use given the ambient temperature.

A covariance matrix can optimize a polarimeter in the presence of known error. This method has been applied to random measurement noise in Stokes polarimetry by Sabatke et al.⁴⁹ and Twietmeyer,⁴⁶ and to random instrument noise in Stokes polarimetry by Tyo.⁴⁸ The covariance matrix is a symmetric matrix which describes the correlation between random variables which have been centered about their means. For Mueller matrix polarimetry, the elements of the 16×16 covariance matrix, \mathbf{C}_M , are

$$C_{M,j,k} = \langle \delta M_j \delta M_k \rangle - \langle \delta M_j \rangle \langle \delta M_k \rangle \quad j, k=1, \dots, 16 \quad (49)$$

One useful error metric (EM) is the sum of the diagonal elements,

$$\text{EM} = \sum_i C_{M,ii} \quad (50)$$

EM is a function of the polarimeter configuration, the number of states, the Mueller matrix of the sample, and the statistical properties of the error sources. Minimization of EM with respect to a polarimeter variable may be used to compute the variable's optimal value in the presence of known error.

15.38 THE MUELLER MATRIX FOR POLARIZATION COMPONENT CHARACTERIZATION

The Mueller matrix provides detailed characterization of a polarization element.^{3,9} Using Mueller matrix functions, all of the previous performance defects and more can be specified. Thus, when using polarization elements in critical applications such as polarimetry, knowledge of its Mueller matrix is desirable. This is analogous to having the interferogram of a lens to ensure that it is of suitable quality for incorporation into a critical imaging system.

15.39 RETRO-REFLECTION TESTING AND CORRECTION FOR SUPPLEMENTAL OPTICS

Some reflective optical components are tested near normal incidence, such as corner cube retro reflectors, liquid crystal on silicon panels (LCOS), and other reflective spatial light modulators. Retro-reflection testing requires the insertion of a low polarization, ideally nonpolarizing, beam splitter in front of the sample, as shown in Fig. 6. The polarimeter measures the Mueller matrix of everything between the generator and the analyzer. This is the *polarization critical region*, where any significant polarization from beam splitters, mirrors, lenses, and the like, needs to be characterized and accounted for in data reduction.

In Fig. 6, a portion of the beam from the polarization state generator reflects from a nonpolarizing beam splitter and is normally incident on the sample; the remainder is removed in a beam dump. The light reflected from the sample divides at the beam splitter and the transmitted portion continues through the polarization analyzer to the focal plane. The focal plane acquires a series of raw images of the sample, and from the set of raw images the Mueller matrix image of all the optics in the polarization critical region is calculated pixel by pixel.

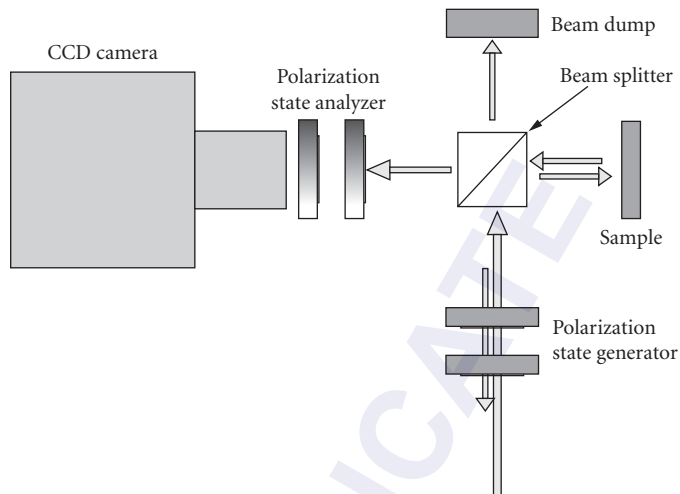


FIGURE 6 Imaging polarimeter configured for retro reflection testing using a nonpolarizing beam splitter and beam dump. (See also color insert.)

To obtain the Mueller matrix image of the sample, contributions from the reflection off the nonpolarizing beam splitter and transmission through the nonpolarizing beam splitter must be calibrated and removed. The ideal nonpolarizing beam splitter should have no polarization, its retardance and diattenuation should be zero; the Mueller matrix would be the identity matrix for both reflection and transmission. In practice, commercially available nonpolarizing beam splitters always have some diattenuation and retardance.

The sample Mueller matrix \mathbf{M}_S is determined from the measured Mueller matrix, $\mathbf{M}_{\text{measured}}$, where \mathbf{M}_T is the beam splitter in transmission and \mathbf{M}_R is the beam splitter in reflection,

$$\mathbf{M}_{\text{measured}} = \mathbf{M}_T \cdot \mathbf{M}_S \cdot \mathbf{M}_R \quad (51)$$

\mathbf{M}_T and \mathbf{M}_R are measured during sample compartment calibration at each wavelength. \mathbf{M}_S is determined as

$$\mathbf{M}_S = (\mathbf{M}_T)^{-1} \cdot \mathbf{M}_{\text{measured}} \cdot (\mathbf{M}_R)^{-1} \quad (52)$$

The compensation must be cautiously applied, with all instrumental variables such as collimation, vignetting, stray light, and angle of incidence carefully considered.

The same method is applicable to lenses, mirrors, and other supplemental optics used to manipulate the beams through the sample compartment. Once the Mueller matrices for the optics before \mathbf{M}_1 and after \mathbf{M}_2 are calibrated, their matrix inverses can be applied during data reduction,

$$\mathbf{M}_S = (\mathbf{M}_2)^{-1} \cdot \mathbf{M}_{\text{measured}} \cdot (\mathbf{M}_1)^{-1} \quad (53)$$

15.40 APPLICATIONS OF POLARIMETRY

Polarimetry and ellipsometry have found application in nearly all areas of science and technology with several tens of thousands of papers detailing various applications. The following summarizes a few of the principal applications and introduces some of the books, reference works, and review papers which provide gateways into the subject.

15.41 ELLIPSOMETRY AND GENERALIZED ELLIPSOMETRY

Ellipsometry is the application of polarimetry to the determination of the optical properties of surfaces and interfaces. Example applications are refractive index and thin-film thickness measurement, and investigations of processes at surfaces such as contamination and corrosion. Chapter 16, “Ellipsometry,” by Rasheed M. A. Azzam treats ellipsometry fundamentals. A more extensive treatment is found in the textbook by Azzam and Bashara.^{6,9} SPIE Milestone Series by Azzam⁵¹ is a collection of historical papers. Calculation of the polarization properties of thin films is presented in the chapter by Dobrowolski,⁵² and also in the text by Macleod.⁵³

Ellipsometry is a well-established technique for determining optical properties such as refractive indices, absorption coefficients, and film thicknesses of material samples by measuring polarization changes that occur on reflection and refraction.⁶ In ellipsometer systems, the measurement configuration is varied, and the polarization change measured. Configuration changes include illumination angle, wavelength, and sample orientation. A forward calculation based on a model, such as the thin-film reflectance or transmission equations, has its free parameters optimized to provide the best fit to the data.

The recent development of generalized ellipsometry or biaxial ellipsometry uses measurements of the complete Jones matrix or Mueller matrix to determine the optical properties of more general anisotropic structures such as birefringent crystals and polarizing films.^{54–68} Generalized ellipsometry measures the optical constants of materials such as anisotropic films, and multilayer stacks of anisotropic films, birefringent crystals, and polarizing materials. With the widespread adoption of biaxial multilayer films in liquid crystal projectors for control of retardance as a function of field of view and wavelength, accurate characterization of anisotropic materials has become more important. The need for such types of ellipsometric instruments has increased with the rapid evolution of liquid crystal displays, new materials and fabrication techniques, and nanostructured materials.

The optical constants of anisotropic materials are conveniently expressed in the form of the dielectric tensor $\tilde{\epsilon}$. For nonoptically active, non-magneto-optic materials with aligned retardance and diattenuation, the dielectric tensor is symmetric and can be expressed as a rotated diagonal matrix of the form

$$\tilde{\epsilon} = \mathbf{R}^{-1}(\varphi, \theta, \psi) \cdot \begin{bmatrix} (n_x + i\kappa_x)^2 & 0 & 0 \\ 0 & (n_y + i\kappa_y)^2 & 0 \\ 0 & 0 & (n_z + i\kappa_z)^2 \end{bmatrix} \cdot \mathbf{R}(\varphi, \theta, \psi) \quad (54)$$

where $n_x + i\kappa_x$, $n_y + i\kappa_y$, and $n_z + i\kappa_z$ are the complex refractive indices along the principal axes, and \mathbf{R} is a rotation matrix through Euler angles φ , θ , and ψ with respect to the laboratory coordinates as represented in Fig. 7.⁶⁹

Characterizing such an anisotropic multilayer thin film requires measuring up to 10 parameters for each layer, 9 that specify the dielectric tensor and 1 for thickness. Multiple measurements must be acquired which span a suitably large range of incident and azimuthal (about the surface normal) angles so that each optical constant to be determined has a distinct effect on the measurements. Changes to the optical constants need to cause distinct changes to the ellipsometric dataset. Each dielectric tensor and thickness parameter must significantly change the polarization within the range of illumination angles, so that ellipsometric data points are not a linear combination of previous measurements. By simultaneously measuring a large range of both incident and azimuthal angles, the components of the dielectric tensor can be determined from a single Mueller matrix image.⁶⁸ Figures 8 and 9 show Mueller matrix imaging polarimeters with converging beams operating in reflection and transmission for generalized ellipsometry. Figure 10 is an example Mueller matrix image of an LC projector biaxial field correcting film.

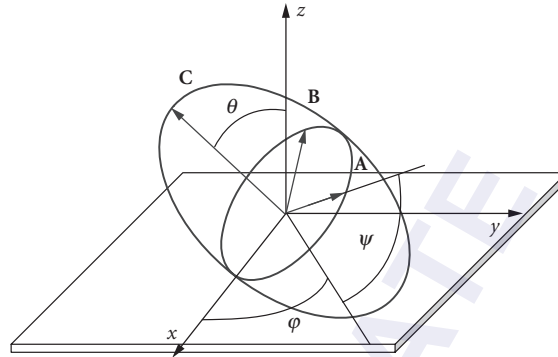


FIGURE 7 Example of a biaxial index ellipsoid with principal axes oriented at an arbitrary orientation along orthogonal vectors **A**, **B**, and **C**.

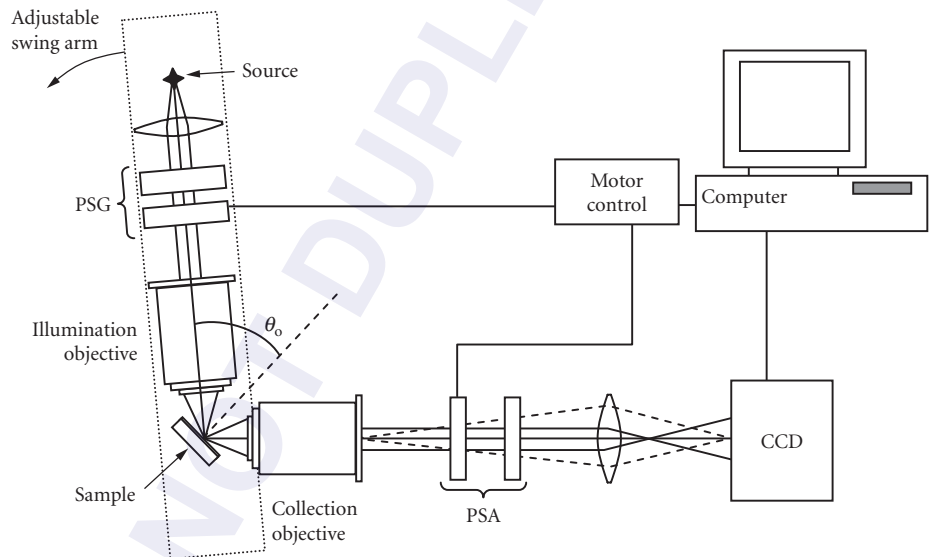


FIGURE 8 Mueller matrix imaging polarimeter configured for reflection generalized ellipsometry with the inclusion of two microscope objectives in the sample compartment. The microscope objective's exit pupil is imaged onto the CCD so that each pixel receives light which reflected at a different angle of incidence and azimuth.

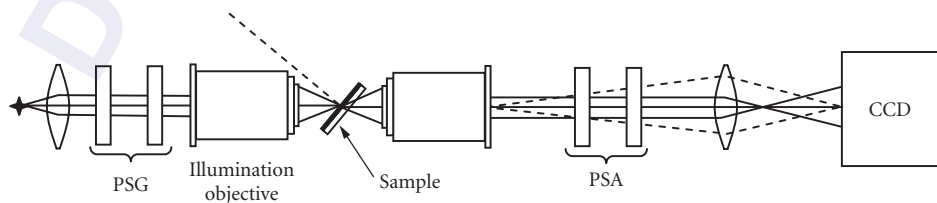


FIGURE 9 This Mueller matrix imaging polarimeter configured for transmission-generalized ellipsometry uses two microscope objectives to obtain polarization change as a function of angle of incidence.

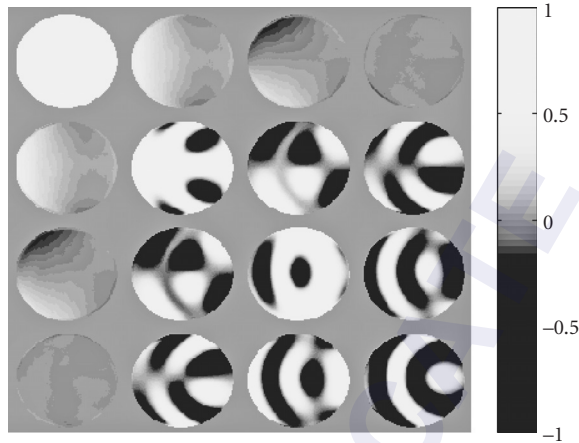


FIGURE 10 Angle of incidence Mueller matrix image of a field widening film for liquid crystal projector systems. The retardance is engineered to vary opposite to the LC, greatly increasing the field of view before artifacts, such as color shifts, occur in the display. The film's optical properties, such as thickness and refractive indices, can be reverse engineered from such measurements by generalized ellipsometry.

Biaxial materials have tensor components that vary with direction and can only be fully characterized if measurements are performed while both the incident and the azimuthal (about the normal) angle of the illuminating light varies with respect to the sample. Dielectric tensor information can also be obtained by measuring at multiple wavelengths and fitting a parameterized dispersion relationship based on a physical model of the dielectric tensor.^{66,70}

Several methods for calculating the reflection and transmission properties for arbitrary anisotropic multilayer structures have been developed, including the Berreman calculus and related methods derived by Yeh, Mansuripur, and Schubert.⁷¹⁻⁷⁴

15.42 LIQUID CRYSTAL CELL AND SYSTEM TESTING

Liquid crystal (LC) displays of all types are polarization critical optical systems, where the systems are readily put out of specifications by misalignment of polarization elements, poor polarization element quality, stress birefringence, depolarization and scattering, LC cell defects, and a myriad of other issues. Such displays include laptop displays, computer monitors, conference room projectors, direct view and projection televisions, and the myriad of small displays in watches, calculators, cell phones, and the like. Mueller matrix polarimetry and imaging polarimetry are important methods to provide detailed diagnostics of LCs and the associated optical systems.

Polarization measurements made on liquid crystal (LC) cells are particularly useful for determining the key physical parameters of the cell, namely the cell gap, rubbing direction, twist angle, and pretilt angles of the LC as shown in Fig. 11. The cell gap is the thickness of the LC layer between the two glass plates. The rubbing direction describes the azimuth angle (orientation angle) of the LC director at the top glass surface, and the twist angle describes the change in orientation angle of the LC director through the thickness of the cell, such that the orientation angle at the bottom glass is the top-glass rubbing direction plus the twist angle. The pretilt

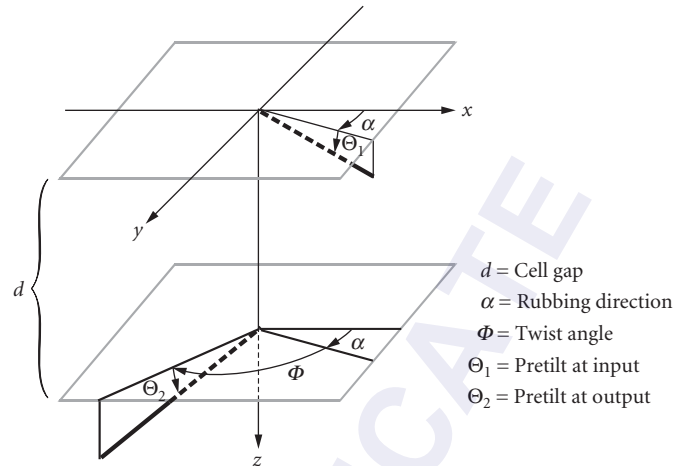


FIGURE 11 Parameters describing a twisted nematic liquid crystal cell can be measured by multiangle Mueller matrix polarimetry.

angles describe the polar angle (tilt angle) of the LC directors at the glass plate interfaces. By adjusting these parameters, LC panel manufacturers can tune their design for the desired response time, color properties, and useable field-of-view. By making polarization measurements of the cell at a variety of incident angles and/or wavelength, and applying curve-fitting techniques, these key parameters can be measured, such as by the AxoScan Mueller matrix polarimeter from Axometrics, Inc. (Huntsville, Alabama).

LCs operate as electrically addressable variable retarders. With the inclusion of polarizers and color filters, pixilated LC arrays (LC panels) serve as color image generators. Several performance specifications of an LC system are critical. The contrast ratio is the ratio of the on-screen illuminance in the white state to the intensity in the dark state. A high contrast ratio depends on the dark state intensity being nearly zero, which requires that all polarization properties are nearly ideal so the beam is well extinguished at the final polarizer. The efficiency of the system is the ratio of the output intensity to the source intensity in the white state. Spatial uniformity characterizes the variations in brightness of the dark and bright states across the aperture.

The contrast ratio and efficiency of the LC projector are a function of the polarization properties of the components in the light valve: the polarizer, beam splitter, LC, and analyzer, as well as the étendue, color balancing, and many other factors. Determination of the polarization properties of each component allows for modeling of system performance, and, in the case of poor system performance, diagnosing which elements are the source of the problem.

The Mueller matrix is especially useful for characterization of LC display system components because sequences of elements behave as the product of their Mueller matrices. Thus the Mueller matrix of an LC panel can be combined with the matrices for beam splitters, dichroic filters, and other components to understand their polarization interactions and tolerance of combinations of optical elements.

Polarization testing of LCs provides important information augmenting radiometric testing. Polarization testing of LC panels requires illuminating the panels with a variety of incident polarization states and measuring the corresponding output polarization states. These additional states are not normally used during LC system operation so their usage for testing appears extraneous until the diagnostic value of retardance and depolarization maps is realized; these parameters directly determine performance properties of the LCs. Radiometric testing measures the LC performance, treating the LC as a black box. Imaging polarimetry quantifies desired and undesired polarization properties enabling better diagnosis of LC problems, problems more difficult to isolate using

TABLE 1 Performance Specification Issues Related to Polarization Causes and LC Defects

Defect	Possible Polarization Causes
Low contrast	Depolarization Incorrect trim retarder Misalignment of liquid crystal, trim retarder, or PBS
Low brightness	Incorrect retardance Oxidized reflector High levels of depolarization
Poor uniformity	Spatial variation in retardance magnitude Spatial variation of retardance orientation Temperature-induced retardance variations

radiometric testing alone. Table 1 and Table 2 summarize the relationship of these polarization properties with several LC performance defects. In Table 1, LC system defects are paired with the related polarization properties which can cause the defect. Many display problems are the direct result of nonideal polarization properties. Table 2 lists some nonideal LC polarization properties and the associated effects.

Depolarization adversely affects LC system performance in different ways than incorrect retardance or retardance nonuniformity. With depolarization, a fraction of the exiting light can be treated as unpolarized light [Eq. (6)]; this is the depolarized component. Fifty percent of the depolarized light will pass through the analyzer and 50 percent will be blocked, so the fraction of leaked light is, at minimum,

$$\text{Leakage} = \frac{1 - \text{DoP}}{2} \quad (55)$$

In the dark state, the leaked depolarized light increases the dark state intensity, and if significant, has a severe effect on the contrast ratio. In the white state half of the depolarized light is blocked by the analyzer decreasing the white state brightness, a less critical problem than dark state leakage. For high contrast, the LC panel must have very low levels of depolarization.

Scattering is a common cause of depolarization in liquid crystals. Liquid crystal depolarization also arises from spatial averaging; micron-scale retardance variations cause adjacent parts of the beam to emerge with different polarization states which average at the polarimeter resulting in a depolarized component in the measurement. An imaging polarimeter measures the average retardance within each of its pixels and any subpixel retardance variations are measured as depolarization. Temperature variations, electric field variations, edge effects in pixels, and disclinations in the LC all cause depolarization.

In any polarimeter measurement, small values of depolarization need to be critically evaluated to ensure the depolarization is due to the device under test and is not due to noise or calibration error within the polarimeter. All polarimeter measurements have some depolarization noise or bias.

TABLE 2 Polarization Defects and Resulting Onscreen Effects

Polarization Property	On-Screen Effect
Retardance spatial nonuniformity	Spatial variation of brightness, color, or contrast
Nonzero dark state retardance	Reduced contrast and color saturation
Incorrect retardance orientation or magnitude	Reduced brightness and contrast
Spectral variation of retardance	Wavelength-dependant contrast and brightness
Depolarization	Reduced contrast and brightness

Typical twisted nematic LC cells cannot be driven to zero retardance, suffering some residual retardance due to thin boundary layers of liquid crystal along the alignment layers as shown in Fig. 3. This residual retardance is usually compensated by placing an additional “trim retarder” over the LC panel.⁷⁵ The trim retarder introduces retardance equal in magnitude to the LC’s single pass retardance with the retardance axis rotated 90°. ⁷⁵ Such a retarder combines with the LC retardance yielding a retardance of zero.

While a trim retarder can be used to reduce the effect of dark state retardance, depolarization cannot be compensated; it must be reduced to acceptable levels during LC device development and fabrication.

15.43 POLARIZATION ABERRATIONS

Polarimetry is useful in optical metrology for measuring the polarization aberrations of optical systems and for characterizing optical and polarization components. Optical systems modify the polarization state of light due to reflections, refractions, and other interactions. Lenses and mirrors have polarization properties described by the Fresnel equations and associated multilayer thin-film equations. For many optical systems, such as camera lenses and Cassegrain telescopes, these polarization aberrations are small, but not necessarily negligible. Other optical systems, with large angles of incidence, diffraction gratings, beam splitters, or other significantly polarizing components, have significant and often troublesome polarization aberrations.

Each ray path through the optical system can be characterized by its polarization matrix. Polarization ray-tracing is the technique of calculating the polarization matrices for ray paths through optical systems.^{24,76–79} Diffraction image formation of polarization-aberrated beams is then handled by vector extensions to diffraction theory.^{80–85} Polarimeters, particularly imaging polarimeters, can measure the Mueller matrices of ray paths through optical systems determining the polarization aberrations. These polarization aberrations frequently have similar functional forms to the geometrical aberrations, since they arise from similar geometrical considerations of surface shape and angle of incidence variation.^{85–92}

Optical system polarization aberrations can be measured by placing the system in the sample compartment of a Mueller matrix imaging polarimeter. Usually the exit pupil is imaged, giving the polarization aberration function (PAF) a Mueller matrix as a function of pupil coordinates. Then maps are generated of linear diattenuation, linear retardance, and other metrics. Figure 12 shows the diattenuation and retardance polarization aberrations measured through a pair of 0.55 numerical aperture microscope objectives; collimated light enters the pupil of the first objective, focuses at the focal point of the second objective, and is recollimated, like Fig. 9 without the sample. The lengths of the lines in the images correspond to the magnitude of the diattenuation and retardance.

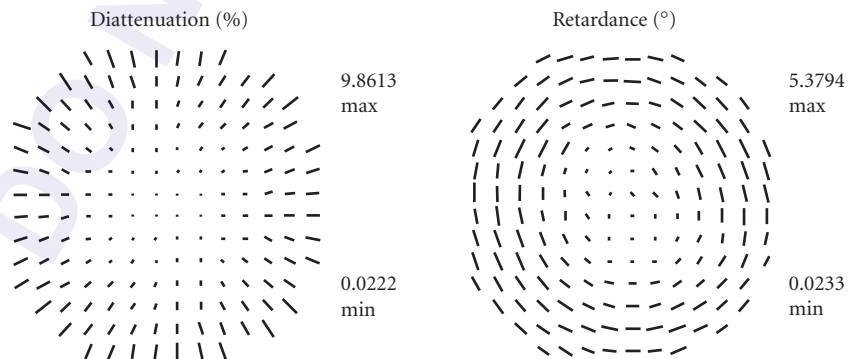


FIGURE 12 The polarization aberrations transmitting through a pair of microscope objectives is represented by these linear diattenuation and linear retardance pupil maps.

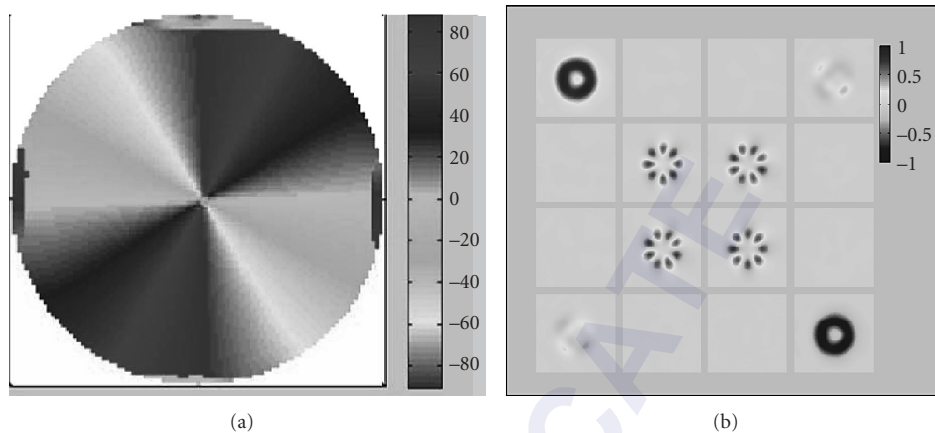


FIGURE 13 The orientation of the fast axis of the half wave vortex retarder rotates by 360° around the pupil (a). The Point spread matrix describes the polarization dependence of the point spread function as a Mueller matrix image (b).

The microscope objective pair has up to 5.4° of spatially varying retardance and 0.1 of spatially varying diattenuation. When placed between crossed linear polarizers, this pair of objectives will leak about 0.15 percent of the incident flux, averaged over the pupil.

When significant polarization aberrations are present, an optical system illuminated with a uniform polarization state will have polarization variations within the point spread function. To characterize these variations and the dependence of the point spread function on the incident polarization state, a Mueller matrix imaging polarimeter focusses on the image of a point object and measures the point spread matrix (**PSM**) as a Mueller matrix image. Measured **PSM** with large polarization aberration is shown in Fig. 13. A vortex retarder was placed in the pupil of an imaging system with a large $f/\#$ image on a camera focal plane, and a Mueller matrix image acquired. This vortex retarder is half-wave linear retarder whose fast axis varies as a function of pupil angle.⁹³ The pupil image on the left side shows the retardance orientation varying by 360° around the pupil. The right side contains the **PSM**. When the Stokes vector of the incident light is multiplied by the **PSM**, the resulting Stokes vector function describes the flux (point spread function) and polarization state variations within the image as a Stokes vector image. Figure 14 shows the point spread function for a fixed incident polarization state and several analyzers, demonstrating the polarization variations within the point spread function.

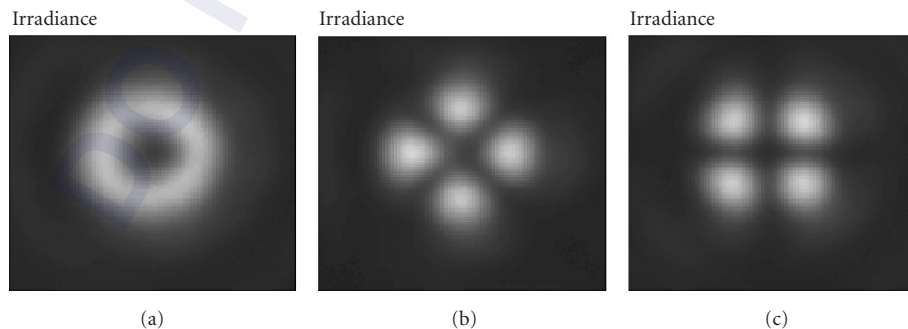


FIGURE 14 The measured point spread function of the vortex retarder completely changes with the analyzed polarization state: (a) no analyzer; (b) horizontal linear analyzer; and (c) vertical linear analyzer. Horizontal linearly polarized light is input.

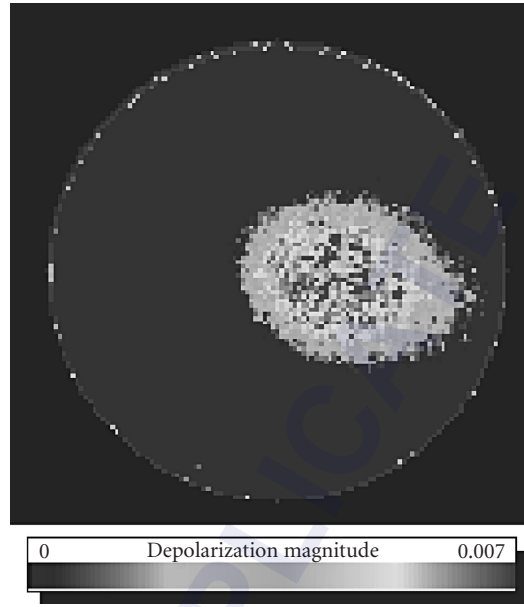


FIGURE 15 Depolarization index of lens with coating damage on the right center causing about 0.005 depolarization.

Figure 15 shows another polarization aberration measurement example.⁹⁴ The lens coatings became damaged by heat and began flaking off. The resulting Mueller matrix pupil image shows a few tenths of a percent depolarization in the damaged area. The undamaged area has a depolarization of only a few hundredths of a percent, more typical of coated lenses.

15.44 REMOTE SENSING

Polarimetry is an important remote sensing technique which complements spectroscopic or hyperspectral imaging. The sunlight which illuminates the earth is essentially unpolarized, but the scattered light has a surprisingly large degree of polarization, which is mostly linear polarization.^{7,94-97} Visible light scattered from forest canopy, cropland, meadows, and similar features frequently has a degree of polarization of 20 percent or greater in the visible range.^{98,99} Light reflecting from mudflats and water can have a degree of polarization greater than 50 percent, particularly for light incident near Brewster's angle. The magnitude of the degree of linear polarization depends on many variables, including the angle of incidence, the angle of scatter, the wavelength, and the weather. The polarization from a site varies from day to day even if the angles of incidence and scatter remain the same; these variations are caused by changes in the earth's vegetation, cloud cover, humidity, rain, and standing water. Polarization is complex to interpret but it conveys useful information.

Light scattered from dense white clouds is nearly unpolarized due to multiple scattering.^{7,96} Scattering from thin aerosols is partially polarized. Hyperspectral imaging combined with Mie scattering theory can determine the mean particle size and the imaginary part of the refractive index of an aerosol. Adding multiangle polarimetric data at visible and shortwave infrared wavelengths provides additional information on the real part of the aerosol refractive index, n_r , and particle size variances, with greater sensitivity than intensity measurements alone. This has been demonstrated with the airborne research scanning polarimeter (RSP),^{100,101} through theoretical sensitivity studies,¹⁰² and with the space-borne Polarization and Directionality of Earth's Reflectances (POLDER) instrument.¹⁰³ POLDER spatial resolution is 6 to 7 km, with degree of linear polarization (DoLP) uncertainty of ~ 2 percent.¹⁰⁴

The Aerosol Polarimeter Sensor (APS) instrument for NASA's Glory mission, using similar design concepts as the airborne RSP, will provide very accurate multi-angle polarimetric measurements (linear polarization uncertainty ~ 0.2 percent), but in a coarse resolution (6 to 20 km) due to nonimaging operation.¹⁰⁵

Many factors affect the accuracy of imaging polarimeters.¹⁰⁶ Polarization aberrations of the optics (instrumental polarization) is addressed through accurate calibration and removal of systematic errors. Many remote sensing polarimeters use different analyzers over different detectors whose signals are then subtracted to measure polarization, and are thus susceptible to gain variations and pixel sensitivity drift. Ongoing detector cross-calibration is desirable. Spatial displacements on the ground between the locations where different polarization orientations are measured gives rise to *polarization artifacts*, also known as *false polarization*. Spatial misregistration between the measurements comprising a polarization measurement is particularly problematic in the presence of scene gradients.

15.45 POLARIZATION LIGHT SCATTERING

Polarization light scattering is the application of polarimetry to scattered light.^{107,108} The scattering characteristics of a sample are generally described by its bidirectional reflectance distribution function, BRDF $(\theta_p, \phi_p, \theta_s, \phi_s, \lambda)$, depicted in Fig. 16, which is the ratio of the scattered flux in a particular direction (θ_s, ϕ_s) to the flux of an incident beam from direction (θ_p, ϕ_p) ,¹⁰⁹

$$\text{BRDF}(\theta_i, \phi_i, \theta_s, \phi_s) = \frac{dL_s(\theta_s, \phi_s)}{dE_i(\theta_i, \phi_i)} \quad (56)$$

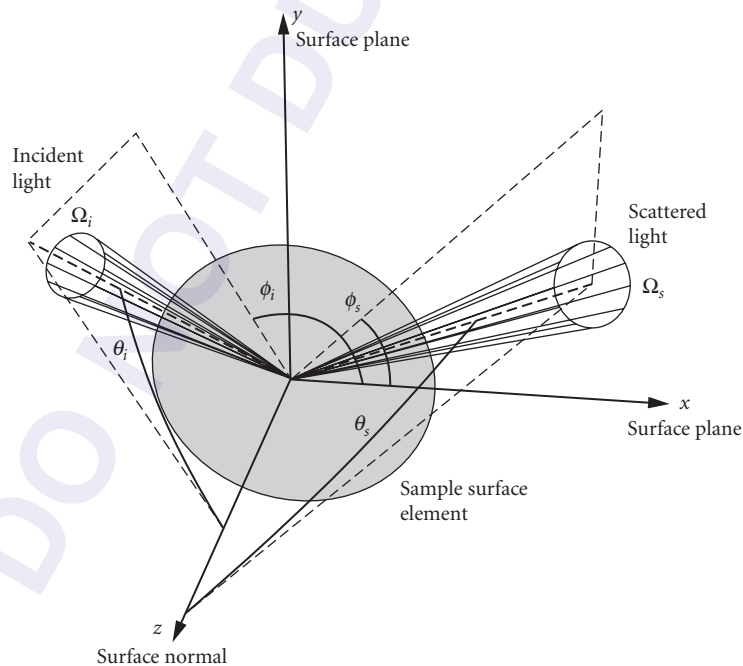


FIGURE 16 The BRDF angle nomenclature: The incident light has an angle of incidence θ_i and azimuth angle ϕ_i and subtends solid angle Ω_i . The scattered light has an angle of scatter θ_s and azimuth angle ϕ_s and subtends solid angle Ω_s .

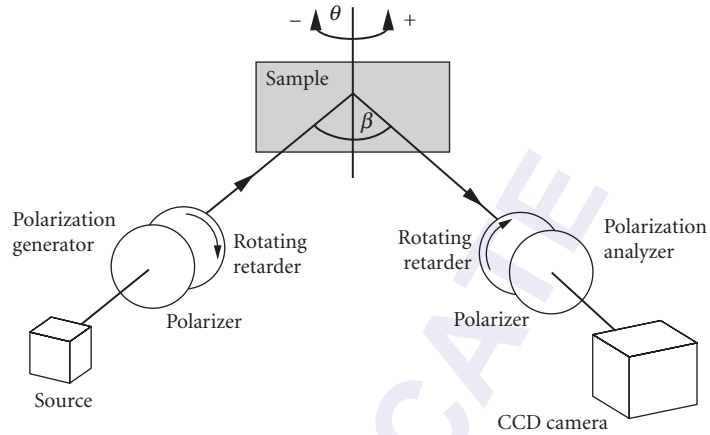


FIGURE 17 Mueller matrix imaging polarimeter for in-plane **MMBRDF** measurements, with bistatic angle β between the polarization generator arm and analyzer arm.

This standard BRDF definition makes no reference to the incident or scattered polarization state so the BRDF function contains no polarization information. The BRDF can be generalized to a Mueller matrix bidirectional reflectance distribution function, or **MMBRDF**($\theta_i, \varphi_i, \theta_s, \varphi_s, \lambda$), the Mueller matrix relating incident and scattered beams in arbitrary directions,²⁵

$$\mathbf{MMBRDF}(\theta_i, \varphi_i, \theta_s, \varphi_s) = \begin{pmatrix} m_{0,0}(\theta_i, \varphi_i, \theta_s, \varphi_s) & \cdots & m_{0,3}(\theta_i, \varphi_i, \theta_s, \varphi_s) \\ \vdots & \ddots & \vdots \\ m_{3,0}(\theta_i, \varphi_i, \theta_s, \varphi_s) & \cdots & m_{3,3}(\theta_i, \varphi_i, \theta_s, \varphi_s) \end{pmatrix} \quad (57)$$

Then the BRDF function is the $m_{0,0}$ element of the **MMBRDF**($\theta_i, \varphi_i, \theta_s, \varphi_s, \lambda$).

Scattered light is a sensitive indicator of surface conditions; a small amount of surface roughness may reduce the specular power by less than a percent while increasing the scattered power by orders of magnitude. The retardance, diattenuation, and depolarization of the scattered light similarly provide sensitive indicators of light scattering conditions, such as uniformity of refractive index, orientation of surface defects, texture, strain and birefringence at an interface, subsurface damage, coating microstructure, and the degree of multiple scattering. Figure 17 depicts a polarimeter configured for polarization scattered light measurements.

Frequently the last 15 **MMBRDF** elements are normalized (divided) by the $m_{0,0}$ element which simplifies the interpretation of polarization properties associated with scattering by adjusting these elements to a -1 to 1 scale.

Figure 18 shows two examples **MMBRDF** from DeBoo.²⁵ Concrete is nearly lambertian and the $m_{1,1}$ element (labeled $m_{0,0}$ elsewhere in this chapter) varies little with angle. The gold-coated diffuser has a more distinct specular peak at $\theta = 0$, and becomes more diattenuating and depolarizing as θ varies away from zero.

15.46 OPHTHALMIC POLARIMETRY

The human visual system is polarization insensitive; an observer cannot discern between unpolarized light and polarized light of various states. The structures of the eye are, however, diattenuating, retarding, and depolarizing.

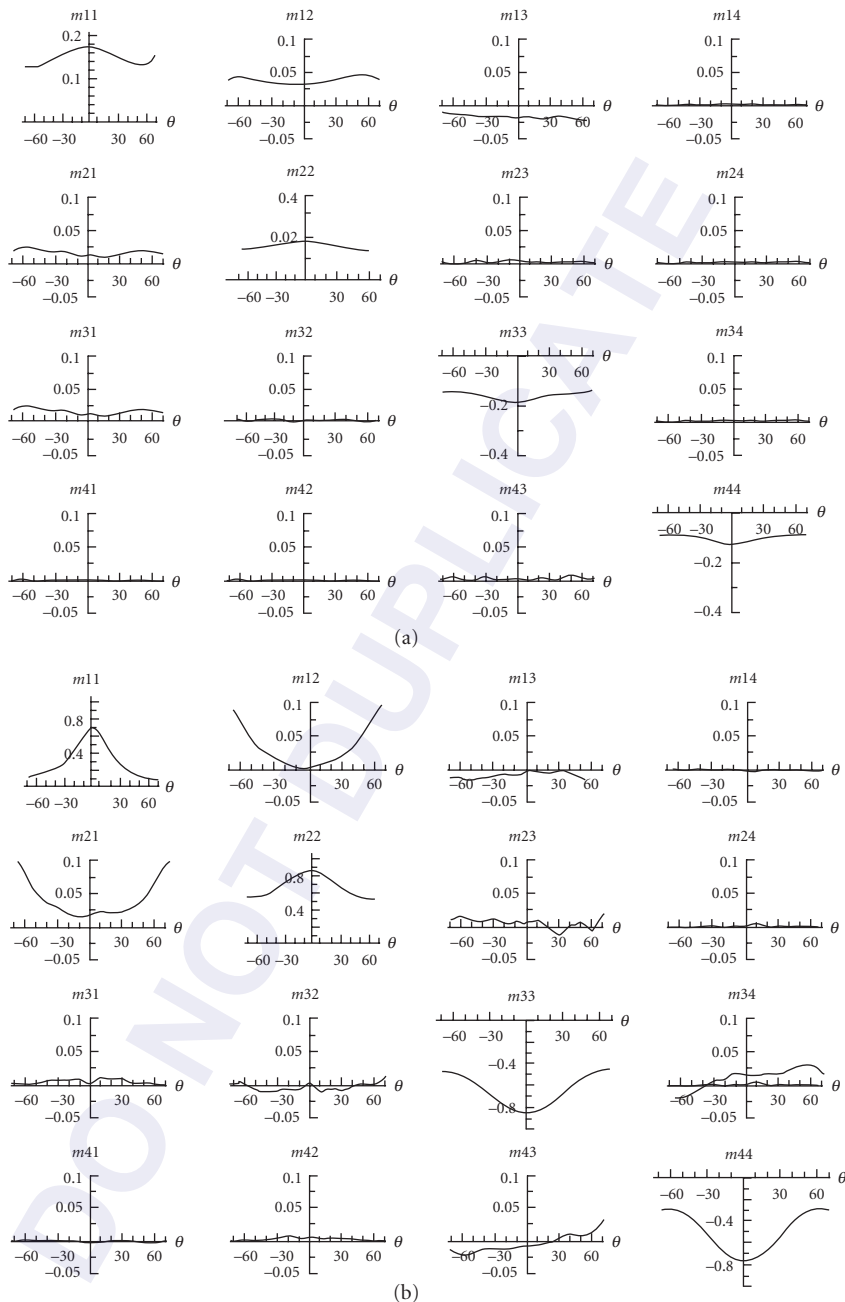


FIGURE 18 (a) In-plane Mueller matrix BRDF for scattering from concrete measured at 808 nm as a normalized Mueller matrix spectrum. In this figure, the Mueller matrix index runs from 1 to 4. The $m_{2,4}$, $m_{4,2}$, $m_{3,4}$, and $m_{4,3}$ elements are nearly zero indicating the absence of linear retardance. The positive $m_{1,2}$ and $m_{2,1}$ elements indicate diattenuation in the s - p orientations such that s has the larger diffuse reflectance. (b) The same for a gold-coated diffuser. The $m_{3,4}$ and $m_{4,3}$ elements are nonzero indicating linear retardance between the s - p components.

The eye's strongest polarization effects are found in the cornea, retinal nerve fiber layer, and Henle's layer. These anisotropic structures contain long thin parallel-oriented cylinders (such as collagen fibrils or microtubules), uniformly distributed within the surrounding medium, with dimensions smaller than the wavelength of visible light. Wiener¹¹⁰ demonstrated in his theory of mixed dielectrics that due to the small difference in refractive index between cylinders and medium, this type of structure has different refractive indices for light polarized parallel and perpendicular to the cylinder axes, an effect termed *form birefringence*. The retardance increases linearly with propagation distance.¹¹¹ Hemenger deduced that these structures also have greater absorption for light oscillating parallel to the cylinders (similar to a wire grid polarizer), so that diattenuation increases with propagation distance, an effect termed *form dichroism*.¹¹² Other models have described diattenuation and retardance with structured ocular tissues.^{113,114}

The interaction of polarized light with retinal tissue has been actively explored to detect subtle changes in the tissue microstructure. A healthy retina has an ordered microstructure.¹¹⁵ The more ordered a structure, the larger the diattenuation and retardance should be. As these cellular structures become disordered in certain disease states, the diattenuation and retardance are expected to decrease and the depolarization to increase.

Direct measurements of retinal polarization have been performed using a variety of techniques. Van Blokland¹¹⁶ was the first to obtain a complete Mueller matrix with a single pixel and demonstrated significant retinal depolarization. Imaging methods include camera-based retinal polarimeters,^{117–119} scanning laser polarimetry methods, and polarization-sensitive optical coherence tomography. Retinal images are assembled through appropriate reconstruction of the detector signal.^{120–122}

A retinal polarimeter, the GDx Nerve Fiber Analyzer (Carl Zeiss Meditec, Dublin California) has been commercially available since the late 1990s and is FDA approved for the measurement of retinal nerve fiber layer thickness and its thinning for the purpose of diagnosing the progression of glaucoma. The GDx is a scanning laser ophthalmoscope measuring linear retardance only, an incomplete polarimeter. The linear retardance is used to estimate the thickness of the retinal nerve fiber layer, which aids in diagnosis and monitoring of glaucoma.^{123–125} As an incomplete polarimeter, the GDx is increasingly inaccurate as depolarization and diattenuation increase.¹¹⁸

Elsner, Burns, and their coworkers have demonstrated the utility of depolarization images to provide higher contrast for deep tissue scattering abnormalities that occur in age-related macular degeneration, central serous chorioretinopathy, and other maculopathies. These abnormalities include drusen, pools of fluid, pigmentation changes, and abnormal vasculature. Lara and Dainty¹²⁶ have reported a complete retinal polarimeter incorporating multiple polarizing beam splitters and detectors.

Polarization-sensitive optical coherence tomography (OCT) generates three-dimensional retinal polarization images.^{127–129} The OCT repeatedly scans with different polarization states illuminating the sample. De Boer has measured the birefringence distribution through the retina and demonstrated that the birefringence of the nerve fiber layer is not uniform. The polarization of the reference beam must be closely matched at the detector for useful fringe visibility. Thus polarization OCT is limited in the measurement of depolarization since multiply scattered incoherent light is rejected in OCT.

15.47 ACKNOWLEDGMENTS

In preparing this chapter, extensive use was made of the works of some of the author's students and other collaborators and they deserve particular credit: Brian de Boo, Karen Twietmeyer, Justin Wolfe, Neil Beaudry, Ann Elsner (Univ. Indiana), Matt Smith (Axometrics), and Dave Diner (NASA/JPL).

15.48 REFERENCES

1. R. M. A. Azzam, "Ellipsometry," In: M. Bass. *Handbook of Optics*, vol 2: McGraw-Hill, New York (1994).
2. J. M. Bennett, "Polarization," In: M. Bass. *Handbook of Optics*, vol 2: McGraw-Hill, New York (1994).
3. W. A. Shurcliff, *Polarized Light—Production and Use*: Harvard University Press, Cambridge, MA (1962).
4. A. Gerrard and J. M. Burch, *Introduction to Matrix Methods in Optics*: Wiley, London (1975).

5. P. S. Theocaris and E. E. Gdoutos, *Matrix Theory of Photoelasticity*: Springer-Verlag, New York (1979).
6. R. M. A. Azzam, (ed.), "Selected Papers on Ellipsometry," *SPIE Milestone Series MS(27)* (1991).
7. K. L. Coulson, *Polarization and Intensity of Light in the Atmosphere*: A. Deepak, Hampton (1988).
8. W. Egan, (ed.), "Polarization in Remote Sensing," *Proc. SPIE 1747* (1992).
9. R. M. A. Azzam and N. M. Bashara, *Ellipsometry and Polarized Light*: Elsevier, Amsterdam, North-Holland (1977).
10. D. S. Kliger, J. W. Lewis, et al., *Polarized Light in Optics and Spectroscopy*: Academic Press, Boston (1990).
11. E. Collett, *Polarized Light*: Marcel Dekker, Inc, New York (1992).
12. D. Goldstein, *Polarized Light*: Marcel Dekker, New York (2003).
13. M. A. F. Thiel, "Error Calculation of Polarization Measurements," *J. Opt. Soc. Am.* **66**(1):65–67 (1976).
14. R. M. A. Azzam, "Instrument Matrix of the Four-Detector Photopolarimeter: Physical Meaning of its Rows and Columns and Constraints on its Elements," *J. Opt. Soc. Am. A.* **7**:87–91 (1990).
15. J. O. Stenflo, "Optimization of the LEST Polarization Modulation System," *LEST Found., Tech. Rep.* **44** (1991).
16. R. M. A. Azzam, "Photopolarimetric Measurement of the Mueller Matrix by Fourier Analysis of a Single Detected Signal," *Opt. Lett.* **2**(6):148–150 (1978).
17. D. H. Goldstein, "Mueller Matrix Dual-Rotating Retarder Polarimeter," *Appl. Opt.* **31**(31):6676–6683 (1992).
18. D. B. Chenault, *Infrared Spectropolarimetry*: University of Alabama, Huntsville, Ala (1992).
19. M. H. Smith, "Optimization of a Dual-Rotating-Retarder Mueller Matrix Polarimeter," *Appl. Opt.* **41**(13):2488–2493 (2002).
20. P. S. Hauge, "Mueller Matrix Ellipsometry with Imperfect Compensators," *J. Opt. Soc. Am.* **68**(11):1519–1528 (1978).
21. L. Broch and L. Johann, "Optimizing Precision of Rotating Compensator Ellipsometry," *Phys. Stat. Sol. C* **5**(5):1036–1040 (2008).
22. D. H. Goldstein and R. A. Chipman, "Error analysis of a Mueller Matrix Polarimeter," *J. Opt. Soc. Am.* **7**(4):693–700 (1990).
23. D. B. Chenault, J. L. Pezzaniti, et al., "Mueller Matrix Algorithms," *Proc. SPIE 1746*:231–246 (1992).
24. R. A. Chipman, "Polarization Analysis of Optical Systems," *Opt. Eng.* **28**(2):90–99 (1989).
25. B. De Boo, J. Sasian, et al., "Degree of Polarization Surfaces and Maps for Analysis of Depolarization," *Opt. Exp.* **12**(20):4941–4958 (2004).
26. R. J. King and S. P. Talim, "Some Aspects of Polarizer Performance," *J. Physics, Ser. E* **4**:93–96 (1971).
27. A. M. Title and W. J. Rosenberg, "Improvements in Birefringent Filters. 5: Field of View Effects," *Appl. Opt.* **18**(20):3443–3456 (1979).
28. P. D. Hale and G. W. Day, "Stability of Birefringent Linear Retarders (Waveplates)," *Appl. Opt.* **27**(24):5146–5153 (1988).
29. S. Y. Lu and R. A. Chipman, "Generalized Diattenuation and Retardance for Inhomogeneous Polarization Elements," *Proc. SPIE 1746*:197–200 (1992).
30. J. L. Pezzaniti and R. A. Chipman, "Off Axis Polarizing Properties of Polarizing Beam Splitter Cubes." "Polarization Analysis and Measurement," *Proc. SPIE 1746*:343–347 (1992).
31. T. Oakberg, "Relative Variation of Stress-Optic Coefficient with Wavelength in Fused Silica and Calcium Fluoride," *Proc. SPIE 3754*:226–234 (1999).
32. C. U. Keller, "Instrumentation for Astrophysical Spectropolarimetry," In: J. Trujillo-Bueno, F. Moreno-Insertis and Sánchez. *Astrophysical Spectropolarimetry*: Cambridge University Press, Cambridge: 303–354 (2002).
33. T. Oakberg and A. Bryan, "Detectors with Photoelastic Modulators," *Proc. SPIE 4819*:98–106 (2002).
34. J. O. Stenflo and H. Povel, "Astronomical Polarimeter with 2-D Detector Arrays," *Appl. Opt.* **24**(22):3893–3898 (1985).
35. B. Wang, J. List, et al., "Stokes Polarimeter Using Two Photoelastic Modulators. Polarization Measurement, Analysis, and Applications V," *Proc. SPIE 4819*:1–8 (2002).
36. H. Povel, H. Aebersold, et al., "Charge-Coupled Device Image Sensor as a Demodulator in a 2-D Polarimeter with a Piezoelectric Modulator," *Appl. Opt.* **29**(8):1186–1190 (1990).

37. H. P. Povel, C. U. Keller, et al., "Two-Dimensional Polarimeter with a Chargecoupled-Device Image Sensor and a Piezoelectric Modulator," *Appl. Opt.* **33**(19):4254–4260 (1994).
38. A. M. Gandorfer and H. P. Povel, "First Observations with a New Imaging Polarimeter," *Astron. Astrophys.* **328**:381–389 (1997).
39. T. C. Oakberg, "Measurement of Low-Level Strain Birefringence in Optical Elements Using a Photoelastic Modulator," *Proc. SPIE* **2873**:17–20 (1996).
40. T. Oakberg, "Measurement of Low-Level Strain Retardation in Optical Materials," *Proc. SPIE* **3121**:19–22 (1997).
41. B. Wang, "An Improved Method for Measuring Low-Level Linear Birefringence in Optical Materials," *Proc. SPIE* **3424**:120–124 (1998).
42. J. Bueno, "Polarimetry Using Liquid-Crystals Variable Retarders: Theory and Calibration," *J. Opt.: Pure Appl. Opt.* **2**:216–222 (2000).
43. K. Twietmeyer, "GDx-MM: An Imaging Mueller Matrix Retinal Polarimeter," College of Optical Sciences, Tucson, The University of Arizona, Ph.D.: 347 (2007).
44. X. S. Yao, L. Yan, et al., "Highly Repeatable All-Solid-State Polarization-State Generator," *Opt. Lett.* **30**(11):1324–1326 (2005).
45. L.-S. Yan, X. S. Yao, et al., "High-Speed and Highly Repeatable Polarization-State Analyzer for 40-gb/s System Performance Monitoring," *IEEE Photonics Technol. Lett.* **18**(4):643–645 (2006).
46. K. Twietmeyer and R. A. Chipman, "Optimization of Mueller Matrix Polarimeters in the Presence of Error Sources," *Opt. Exp.* **16**(15):11589–11603 (2008).
47. G. H. Golub and C. F. Van Loan, *Matrix Computations*: John Hopkins University Press, Baltimore (1983).
48. J. S. Tyo, "Design of Optimal Polarimeters: Maximization of Signal-to-Noise Ratio and Minimization of Systematic Error," *Appl. Opt.* **41**(4):619–630 (2002).
49. D. S. Sabatke, A. M. Locke, et al., "Figures of Merit for Complete Stokes Polarimeter Optimization. Polarization Analysis, Measurement, and Remote Sensing III," *Proc. SPIE* **4133**:75–81 (2000).
50. R. A. Horn and C. R. Johnson, *Matrix Analysis*: Cambridge University Press, Cambridge (1985).
51. R. M. A. Azzam and N. M. Bashara, *Ellipsometry and Polarized Light*: Elsevier, Amsterdam, North-Holland (1987).
52. G. Dobrowolski, "Optical Properties of Films and Coatings," In: M. Bass. *Handbook of Optics*, vol 1: McGraw-Hill, New York (1994).
53. H. A. Macleod, *Thin Film Optical Filters*: Macmillan, New York (1986).
54. J. Lekner, "Ellipsometry of Anisotropic Media," *J. Opt. Soc. Am. A.* **10**(7):1579–1581 (1993).
55. M. Schubert, B. Rheinlander, et al., "Generalized Transmission Ellipsometry for Twisted Biaxial Dielectric Media: Application to Chiral Liquid Crystals," *J. Opt. Soc. Am. A.* **13**(9):1930–1940 (1996).
56. M. Schubert, B. Rheinlander, et al., "Extension of Rotating-Analyzer Ellipsometry to Generalized Ellipsometry: Determination of the Dielectric Function Tensor from Uniaxial TiO_2 ," *J. Opt. Soc. Am. A.* **13**(4):875–883 (1996).
57. G. E. Jellison Jr. and F. A. Modine, "Two-Modulator Generalized Ellipsometry: Experiment and Calibration," *Appl. Opt.* **36**(31):8184–8189 (1997).
58. G. E. Jellison Jr. and F. A. Modine, "Two-Modulator Generalized Ellipsometry: Theory," *Appl. Opt.* **36**(31):8190–8198 (1997).
59. G. E. Jellison Jr., F. A. Modine, et al., "Measurement of the Optical Functions of Uniaxial Materials by Two-Modulator Generalized Ellipsometry: Rutile (TiO_2)," *Opt. Lett.* **22**(23):1808–1810 (1997).
60. M. Schubert, B. Rheinlander, et al., "Anisotropy of Boron Nitride Thin-Film Reflectivity Spectra by Generalized Ellipsometry," *App. Phys. Lett.* **70**(14):1819–1821 (1997).
61. J. F. Elman, J. G. U., et al., "Characterization of Biaxially-Stretched Plastic Films by Generalized Ellipsometry," *Thin Solid Films* **313–314**:814–818 (1998).
62. A. En-Naciri, L. Johann, et al., "Spectroscopic Ellipsometry of Anisotropic Materials: Application to the Optical Constants of HgI_2 ," *Appl. Opt.* **38**(4):647–654 (1998).
63. J.-D. Hecht, A. Eifler, et al., "Birefringence and Reflectivity of Single-Crystal CdAl_2Se_4 by Generalized Ellipsometry," *Phys. Rev. B* **57**(12):7037–7042 (1998).

64. A. En-Naciri, L. Johann, et al., "Spectroscopic Generalized Ellipsometry Based on Fourier Analysis," *Appl. Opt.* **38**(22):4802–4811 (1999).
65. M. Schubert, T. E. Tiwald, et al., "Explicit Solutions for the Optical Properties of Arbitrary Magneto-Optic Materials in Generalized Ellipsometry," *Appl. Opt.* **38**(1):177–187 (1999).
66. M. Schubert, A. Kasic, et al., "Generalized Ellipsometry of Complex Mediums in Layered Systems," *Proc. SPIE* **4806**:264–276 (2002).
67. R. A. Chipman, "Biaxial Ellipsometry," *Proc. SPIE* **5875**:587506 (2005). DOI:10.1117/12.623405.
68. N. Beaudry, Y. Zhao, and R. Chipman, "Dielectric Tensor Measurement from a Single Mueller Matrix Image," *J. Opt. Soc. Am. A.* **24**(3):814–824 (2007).
69. M. Born and E. Wolf, *Principles of Optics*: Cambridge University Press, Cambridge, UK (1980).
70. G. E. Jellison, "Spectroscopic Ellipsometry Data Analysis: Measured versus Calculated Quantities," *Thin Solid Films* **313–314**:33–39 (1998).
71. D. W. Berreman, "Optics in Stratified and Anisotropic Media: 4×4 -Matrix Formulation," *J. Opt. Soc. Am.* **62**(4):502–510 (1972).
72. P. Yeh, "Optics of Anisotropic Layered Media: A New 4×4 Matrix Algebra," *Surf. Sci.* **96**:41–53 (1980).
73. M. Mansuripur, "Analysis of Multilayer Thin-Film Structures Containing Magneto-optic and Anisotropic Media at Oblique Incidence Using 2×2 Matrices," *J. App. Phys.* **67**(10):6466–6475 (1990).
74. M. Schubert, "Polarization-Dependent Optical Parameters of Arbitrarily Anisotropic Homogeneous Layered Systems," *Phys. Rev. B* **53**(8):4265–4274 (1996).
75. D. J. McKnight, K. M. Johnson, et al., "256 by 256 Liquid-Crystal-on-Silicon Spatial Light Modulator," *Appl. Opt.* **33**:2775–2784 (1994).
76. T. J. Bruegge, "Analysis of Polarization in Optical Systems," *Proc. SPIE* **1166**:165–176 (1989).
77. E. Waluschka, "A Polarization Ray Trace," *Opt. Eng.* **28**:86–89 (1989).
78. R. A. Chipman, "Polarization Analysis of Optical Systems II," *Proc. SPIE* **1166**:79–99 (1989).
79. L. B. Wolff, and D. J. Kurlander, "Ray Tracing with Polarization Parameters," *IEEE Computer Graphics and Appl.*:44–55 (1990).
80. H. Kuboda, and S. Inouè, "Diffraction Images in the Polarizing Microscope," *J. Opt. Soc. Am.* **49**(2):191–198 (1959).
81. W. Urbanczyk, "Optical Imaging Systems Changing the State of Light Polarization," *Optik* **66**:301–309 (1984).
82. W. Urbanczyk, "Optical Transfer Functions for Imaging Systems Which Change the State of Light Polarization," *Opt. Acta.* **33**:53–62 (1986).
83. J. P. McGuire Jr. and R. A. Chipman, "Diffraction Image Formation in Optical Systems with Polarization Aberrations I: Formulation and Example," *J. Opt. Soc. Am. A.* **7**(9):1614–1626 (1990).
84. M. Mansuriper, "Effects of High-Numerical-Aperture Focusing on the State of Polarization in Optical and Magneto-Optical Data Storage Systems," *Appl. Opt.* **30**(22):3154–3162 (1991).
85. J. P. McGuire Jr. and R. A. Chipman, "Diffraction Image Formation in Optical Systems with Polarization Aberrations II: Amplitude Response Matrices for Rotationally Symmetric Systems," *J. Opt. Soc. Am. A.* **8**:833–840 (1991).
86. R. A. Chipman, "Polarization Aberrations," *Optical Sciences*: Univ. Of Arizona, Tucson (1987).
87. J. P. McGuire Jr. and R. A. Chipman, "Polarization Aberrations in Optical Systems" In: R. Fischer and W. Smith (eds.) "Current Developments in Optical Engineering II," *Proc. SPIE* **818**:240–245 (1987).
88. E. W. Hansen, "Overcoming Polarization Aberrations in Microscopy," In: R. Chipman (ed.) *Polarization Considerations for Optical Systems. Proc. SPIE* **891**:190–197 (1988).
89. R. A. Chipman, and L. J. Chipman, "Polarization Aberration Diagrams," *Opt. Eng.* **28**(2):100–106 (1989).
90. J. P. McGuire Jr. and R. A. Chipman, "Polarization Aberrations in the Solar Activity Measurements Experiments (SAMEX) Solar Vector Magnetograph," *Opt. Eng.* **28**(2):141–147 (1989).
91. J. P. McGuire Jr. and R. A. Chipman, "Diffraction Image Formation in Optical Systems with Polarization Aberrations I: Formulation and Example," *J. Opt. Soc. Am. A.* **7**(9):1614–1626 (1990).
92. J. P. McGuire Jr. and R. A. Chipman, "Analysis of Spatial Pseudodepolarizers in Imaging Systems," *Opt. Eng.* **29**(12):1478–1484 (1990).

93. S. C. McEldowney, D. M. Shemo, et al., "Vortex Retarders Produced from Photoaligned Liquid Crystal Polymers," *Opt. Exp.* **16**(10):7295–7308 (2008).
94. J. Wolfe and R. A. Chipman, "Reducing Symmetric Polarization Aberrations in a Lens by Annealing," *Opt. Exp.* **12**(15):3443–3451 (2004).
95. W. G. Egan, *Photometry and Polarization in Remote Sensing*: Elsevier, New York (1985).
96. G. P. Konnen, *Polarized Light in Nature*: Cambridge University Press, Cambridge (1985).
97. K. L. Coulson, "Polarization of Light in the Natural Environment," *Proc. SPIE* **1166**:2–10 (1989).
98. P. J. Curran, "Polarized Visible Light as an Aid to Vegetation Classification," *Remote Sensing of the Environment* **12**:491–499 (1982).
99. M. J. Duggin, S. A. Israel, et al., "Use of Polarization Methods in Earth Resource Investigations," *Proc. SPIE* **1166**:42–51 (1989).
100. L. A. Remer, D. Tanré, et al., "Validation of MODIS Aerosol Retrieval over Ocean," *Geophys. Res. Lett.* **29**:12 (2002). DOI: 10.1029/2001GL013204.
101. H. Yu, Y. J. Kaufman, et al., "A Review of Measurement-Based Assessments of the Aerosol Direct Radiative Effect and Forcing," *Atmos. Chem. Phys.* **6**:613–666 (2006).
102. D. A. Chu, Y. J. Kaufman, et al., "Validation of MODIS Aerosol Optical Depth Retrieval Over Land," *Geophys. Res. Lett.* **29**(12):8007 (2002). DOI: 10.1029/2002GL013205.
103. O. Torres, P. K. Bhartia, et al., "A Long Term Record of Aerosol Optical Depth from TOMS Observations and Comparison to AERONET Measurements," *J. Atm. Sci.* **59**:398–413 (2002).
104. O. Torres, R. Decae, et al., "OMI Aerosol Retrieval Algorithm," In: P. Stammes and R. Noordhoek, *OMI Algorithm Theoretical Basis Document*, vol. III (2002).
105. P. Veeffkind, G. de Leeuw, et al., "Retrieval of Aerosol Optical Depth over Land Using Two-Angle View Satellite Radiometry during TARFOX," *Geophys. Res. Lett.* **25**:3135–3138 (1998).
106. J. V. Martonchik, D. J. Diner, et al., "Regional Aerosol Retrieval Results from MISR," *IEEE Trans. Geosci. Remote Sens.* **40**:1520–1531 (2002).
107. H. C. van de Hulst, *Light Scattering by Small Particles*: John Wiley and Sons, New York (1957).
108. J. C. Stover, *Optical Scattering, Measurement and Analysis*: McGraw-Hill Inc, New York (1990).
109. J. C. Stover, *Optical Scatter Measurements and Analysis*, 2d ed.: SPIE Optical Engineering Press, Bellingham (1995).
110. O. Wiener, "Die Theorie des Mischkorpers für das Feld der stationären Stromung," *Abh. Math.-Phys. Klasse Königlich Sachsische Des. Wiss.* **32**:509–604 (1912).
111. F. Horowitz, Structure-Induced Optical Anisotropy in Thin Films. Tucson, University of Arizona. Ph.D (1983).
112. R. P. Hemenger, "Dichroism of the Macular Pigment and Haidinger's Brushes," *J. Opt. Soc. Am.* **72**(6):734–737 (1982).
113. X.-R. Huang and R. W. Knighton, "Theoretical Model of the Polarization Properties of the Retinal Nerve Fiber Layer in Reflection," *Appl. Opt.* **42**(28):5726–5736 (2003).
114. N. J. Kemp, H. N. Zaatari, et al., "Form-Biattenuance in Fibrous Tissues Measured with Polarization-Sensitive Optical Coherence Tomography (PS-OCT)," *Opt. Exp.* **13**(12):4611–4628 (2005).
115. M. Miura, A. E. Elsner, et al., "Imaging Polarimetry in Central Serous Chorioretinopathy," *Am. J. Ophthalmol.* **140**(6):1014–1019 (2005).
116. G. J. van Blokland, "Ellipsometry of the Human Retinal in Vivo: Preservation of Polarization," *J. Opt. Soc. Am. A.* **2**(3):72–75 (1985).
117. X.-R. Huang and R. W. Knighton, "Linear Birefringence of the Retinal Nerve Fiber Layer Measured In Vitro with a Multispectral Imaging Micropolarimeter," *J. Biomed. Opt.* **7**(2):199–204 (2002).
118. J. M. Bueno, "The Influence of Depolarization and Corneal Birefringence on Ocular Polarization," *J. Opt. A. Pure Appl. Opt.* **6**:S91–S99 (2004).
119. J. M. Bueno, E. Berrio, et al., "Degree of Polarization as an Objective Method of Estimating Scattering," *J. Opt. Soc. Am. A.* **21**(7):1316–1321 (2004).
120. R. N. Weinreb, C. Bowd, et al., "Measurement of the Magnitude and Axis of Corneal Polarization with Scanning Laser Polarimetry," *Arch. Ophthalmol.* **120**:901–906 (2002).

121. X.-R. Huang, H. Bagga, et al., "Variation of Peripapillary Retinal Nerve Fiber Layer Birefringence in Normal Human Subjects," *IOVS* **45**(9):3073–3080 (2004).
122. J. M. Bueno and B. Vohnsen, "Polarimetric High-Resolution Confocal Scanning Laser Ophthalmoscope," *Vis. Res.* **45**:3526–3534 (2005).
123. R. N. Weinreb, L. Zangwill, et al., "Detection of Glaucoma with Scanning Laser Polarimetry," *Arch. Ophthalmol* **116**:1583–1589 (1998).
124. A. Weber, A. E. Elsner, et al., "Relationship between Foveal Birefringence and Visual Acuity in Neovascular Age-Related Macular Degeneration," *Eye* **21**(3):353–361 (2006).
125. A. E. Elsner, A. Weber, et al., "Imaging Polarimetry in Patients with Neovascular Age-Related Macular Degeneration," *J. Opt. Soc. Am. A.* **24**(5):1468–1480 (2007).
126. D. Lara and C. Dainty, "Axially Resolved Complete Mueller Matrix Confocal Microscopy," *Appl. Opt.* **45**(9):1917–1930 (2006).
127. J. F. de Boer, T. E. Milner, et al., "Two Dimensional Birefringence Imaging in Biological Tissue by Polarization-Sensitive Optical Coherence Tomography," *Opt. Lett.* **22**:934–936 (1997).
128. J. F. de Boer, T. E. Milner, et al., "Determination of the Depth-Resolved Stokes Parameters of Light Backscattered from Turbid Media by Use of Polarization-Sensitive Optical Coherence Tomography," *Opt. Lett.* **24**(5):300–302 (1998).
129. J. E. Roth, J. A. Kozak, et al., "Simplified Method for Polarization-Sensitive Optical Coherence Tomography," *Opt. Lett.* **26**:1069–1071 (2001).

ELLIPSOMETRY

Rasheed M. A. Azzam

*Department of Electrical Engineering
University of New Orleans
New Orleans, Louisiana*

16.1 GLOSSARY

A	instrument matrix
D_ϕ	film thickness period
d	film thickness
E	electrical field
\mathbf{E}_0	constant complex vector
$f()$	function
I	interface scattering matrix
k	extinction coefficient
L	layer scattering matrix
N	complex refractive index = $n - jk$
n	real part of the complex refractive index
R	reflection coefficient
r	reflection coefficient
S_{ij}	scattering matrix elements
s, p	subscripts for polarization components
X	$\exp(-j2\pi d/D_\phi)$
Δ	ellipsometric angle
ϵ	dielectric function
$\langle \epsilon \rangle$	psuedo dielectric function
ρ	$R_p/R_s = \tan \psi \exp(j\Delta) = \chi_i/\chi_r$
ϕ	angle of incidence
χ_i	E_{is}/E_{ip}
χ_r	E_{rs}/E_{rp}
ψ	ellipsometric angle

16.2 INTRODUCTION

Ellipsometry is a nonperturbing optical technique that uses the change in the state of polarization of light upon reflection for the in-situ and real-time characterization of surfaces, interfaces, and thin films. In this chapter we provide a brief account of this subject with an emphasis on modeling and instrumentation. For extensive coverage, including applications, the reader is referred to several monographs,¹⁻⁴ handbook,⁵ collected reprints,⁶ conference proceedings,⁷⁻¹⁵ and general and topical reviews.¹⁶⁻³²

In ellipsometry, a collimated beam of monochromatic or quasi-monochromatic light, which is polarized in a known state, is incident on a sample surface under examination, and the state of polarization of the reflected light is analyzed. From the incident and reflected states of polarization, ratios of complex reflection coefficients of the surface for the incident orthogonal linear polarizations parallel and perpendicular to the plane of incidence are determined. These ratios are subsequently related to the structural and optical properties of the ambient-sample interface region by invoking an appropriate model and the electromagnetic theory of reflection. Finally, model parameters of interest are determined by solving the resulting inverse problem.

In ellipsometry, one of the two copropagating orthogonally polarized waves can be considered to act as a reference for the other. Inasmuch as the state of polarization of light is determined by the superposition of the orthogonal components of the electric field vector, an ellipsometer may be thought of as a common-path polarization interferometer. And because ellipsometry involves only relative amplitude and relative phase measurements, it is highly accurate. Furthermore, its sensitivity to minute changes in the interface region, such as the formation of a submonolayer of atoms or molecules, has qualified ellipsometry for many applications in surface science and thin-film technologies.

In a typical scheme, Fig. 1, the incident light is linearly polarized at a known but arbitrary azimuth and the reflected light is elliptically polarized. Measurement of the ellipse of polarization of the reflected light accounts for the name ellipsometry, which was first coined by Rothen.³³ (For a discussion of light polarization, the reader is referred to Chap. 12 in this volume. For a historical background on ellipsometry, see Rothen³⁴ and Hall.³⁵)

For optically isotropic structures, ellipsometry is carried out only at oblique incidence. In this case, if the incident light is linearly polarized with the electric vector vibrating parallel p or perpendicular s to the plane of incidence, the reflected light is likewise p - and s -polarized, respectively. In other words, the p and s linear polarizations are the eigenpolarizations of reflection.³⁶ The associated eigenvalues are the complex amplitude reflection coefficients R_p and R_s . For an arbitrary input state with phasor electric-field components E_{ip} and E_{is} , the corresponding field components of the reflected light are given by

$$E_{rp} = R_p E_{ip} \quad E_{rs} = R_s E_{is} \quad (1)$$

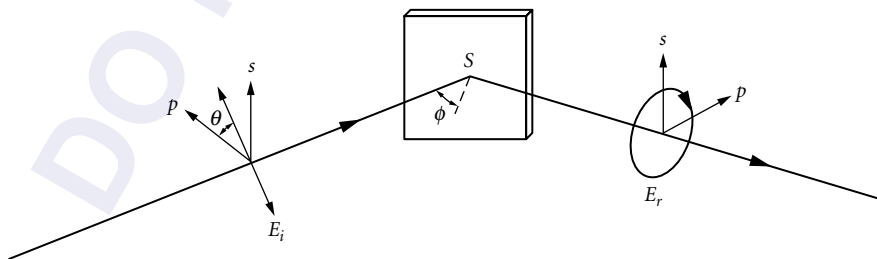


FIGURE 1 Incident linearly polarized light of arbitrary azimuth θ is reflected from the surface S as elliptically polarized. p and s identify the linear polarization directions parallel and perpendicular to the plane of incidence and form a right-handed system with the direction of propagation. ϕ is the angle of incidence.

By taking the ratio of the respective sides of these two equations, one gets

$$\rho = \chi_i / \chi_r \quad (2)$$

where

$$\rho = R_p / R_s \quad (3)$$

$$\chi_i = E_{is} / E_{ip} \quad \chi_r = E_{rs} / E_{rp} \quad (4)$$

χ_i and χ_r of Eqs. (4) are complex numbers that succinctly describe the incident and reflected polarization states of light;³⁷ their ratio, according to Eqs. (2) and (3), determines the ratio of the complex reflection coefficients for the p and s polarizations. Therefore, ellipsometry involves pure polarization measurements (without account for absolute light intensity or absolute phase) to determine ρ . It has become customary in ellipsometry to express ρ in polar form in terms of two *ellipsometric angles* ψ and Δ ($0 \leq \psi \leq 90^\circ$, $0 \leq \Delta < 360^\circ$) as follows

$$\rho = \tan \psi \exp(j\Delta) \quad (5)$$

$\tan \psi = |R_p| / |R_s|$ represents the relative amplitude attenuation and $\Delta = \arg(R_p) - \arg(R_s)$ is the differential phase shift of the p and s linearly polarized components upon reflection.

Regardless of the nature of the sample, ρ is a function,

$$\rho = f(\phi, \lambda) \quad (6)$$

of the angle of incidence ϕ and the wavelength of light λ . Multiple-angle-of-incidence ellipsometry³⁸⁻⁴³ (MAIE) involves measurement of ρ as a function of ϕ , and spectroscopic ellipsometry^{3,22,27-31} (SE) refers to the measurement of ρ as a function of λ . In variable-angle spectroscopic ellipsometry⁴³ (VASE) the ellipsometric function ρ of the two real variables ϕ and λ is recorded.

16.3 CONVENTIONS

The widely accepted conventions in ellipsometry are those adopted at the 1968 Symposium on Recent Developments in Ellipsometry following discussions of a paper by Muller.⁴⁴ Briefly, the electric field of a monochromatic plane wave traveling in the direction of the z axis is taken as

$$\mathbf{E} = \mathbf{E}_0 \exp(-j2\pi Nz/\lambda) \exp(j\omega t) \quad (7)$$

where \mathbf{E}_0 is a constant complex vector that represents the transverse electric field in the $z = 0$ plane, N is the complex refractive index of the optically isotropic medium of propagation, ω is the angular frequency, and t is the time. N is written in terms of its real and imaginary parts as

$$N = n - jk \quad (8)$$

where $n > 0$ is the refractive index and $k \geq 0$ is the extinction coefficient. The positive directions of p and s before and after reflection form a right-handed coordinate system with the directions of propagation of the incident and reflected waves, Fig. 1. At normal incidence ($\phi = 0$), the p directions in the incident and reflected waves are antiparallel, whereas the s directions are parallel. Some of the consequences of these conventions are as follows:

1. At normal incidence, $R_p = -R_s$, $\rho = -1$, and $\Delta = \pi$.
2. At grazing incidence, $R_p = R_s$, $\rho = 1$, and $\Delta = 0$.
3. For an abrupt interface between two homogeneous and isotropic semi-infinite media, Δ is in the range $0 \leq \Delta \leq \pi$, and $0 \leq \psi \leq 45^\circ$.

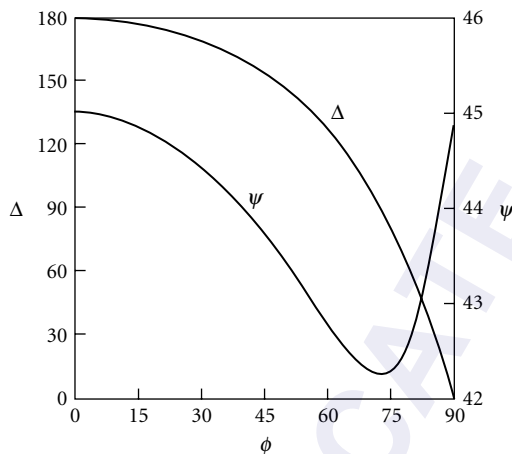


FIGURE 2 Ellipsometric parameters ψ and Δ of an air/Au interface as functions of the angle of incidence ϕ . The complex refractive index of Au is assumed to be $0.306 - j2.880$ at 564-nm wavelength. ψ , Δ , and ϕ are in degrees.

As an example, Fig. 2 shows ψ and Δ vs. ϕ for light reflection at the air/Au interface, assuming $N = 0.306 - j2.880$ for Au⁴⁵ at $\lambda = 564$ nm.

16.4 MODELING AND INVERSION

The following simplifying assumptions are usually made or implied in conventional ellipsometry: (1) the incident beam is approximated by a monochromatic plane wave; (2) the ambient or incidence medium is transparent and optically isotropic; (3) the sample surface is a plane boundary; (4) the sample (and ambient) optical properties are uniform laterally but may change in the direction of the normal to the ambient-sample interface; (5) the coherence length of the incident light is much greater than its penetration depth into the sample; and (6) the light-sample interaction is linear (elastic), hence frequency-conserving.

Determination of the ratio of complex reflection coefficients is rarely an end in itself. Usually, one is interested in more fundamental information about the sample than is conveyed by ρ . In particular, ellipsometry is used to characterize the optical and structural properties of the interfacial region. This requires that a stratified-medium model (SMM) for the sample under measurement be postulated that contains the sample physical parameters of interest. For example, for visible light, a polished Si surface in air may be modeled as an optically opaque (semi-infinite) Si substrate which is covered with a SiO₂ film, with the Si and SiO₂ phases assumed uniform, and the air/SiO₂ and SiO₂/Si interfaces considered as parallel planes. This is often referred to as the three-phase model. More complexity (and more layers) can be built into this basic SMM to represent such finer details as the interfacial roughness and phase mixing, a damage surface layer on Si caused by polishing, or the possible presence of an outermost contamination film. Effective medium theories^{46–54} (EMTs) are used to calculate the dielectric functions of mixed phases based on their microstructure and component volume fractions; and the established theory of light reflection by stratified structures^{55–60} is employed to calculate the ellipsometric function for an assumed set of model parameters. Finally, values of the model parameters are sought that best match the measured and computed values of ρ . Extensive data (obtained, e.g., using VASE) is required to determine the parameters of more complicated samples. The latter task, called the

inverse problem, usually employs linear regression analysis,^{61–63} which yields information on parameter correlations and confidence limits. Therefore, the full practice of ellipsometry involves, in general, the execution and integration of three tasks: (1) polarization measurements that yield ratios of complex reflection coefficients, (2) sample modeling and the application of electromagnetic theory to calculate the ellipsometric function, and (3) solving the inverse problem to determine model parameters that best match the experimental and theoretically calculated values of the ellipsometric function.

Confidence in the model is established by showing that complete spectra can be described in terms of a few wavelength-independent parameters, or by checking the predictive power of the model in determining the optical properties of the sample under new experimental conditions.²⁷

The Two-Phase Model

For a single interface between two homogeneous and isotropic media, 0 and 1, the reflection coefficients are given by the Fresnel formulas¹

$$r_{01p} = (\epsilon_1 S_0 - \epsilon_0 S_1) / (\epsilon_1 S_0 + \epsilon_0 S_1) \quad (9)$$

$$r_{01s} = (S_0 - S_1) / (S_0 + S_1) \quad (10)$$

in which

$$\epsilon_i = N_i^2 \quad i = 0, 1 \quad (11)$$

is the dielectric function (or dielectric constant at a given wavelength) of the i th medium,

$$S_i = (\epsilon_i - \epsilon_0 \sin^2 \phi)^{1/2} \quad (12)$$

and ϕ is the angle of incidence in medium 0 (measured from the interface normal). The ratio of complex reflection coefficients which is measured by ellipsometry is

$$\rho = [\sin \phi \tan \phi - (\epsilon - \sin^2 \phi)^{1/2}] / [\sin \phi \tan \phi + (\epsilon - \sin^2 \phi)^{1/2}] \quad (13)$$

where $\epsilon = \epsilon_1 / \epsilon_0$. Solving Eq. (13) for ϵ gives

$$\epsilon_1 = \epsilon_0 \{ \sin^2 \phi + \sin^2 \phi \tan^2 \phi [(1 - \rho) / (1 + \rho)]^2 \} \quad (14)$$

For light incident from a medium (e.g., vacuum, air, or an inert ambient) of known ϵ_0 , Eq. (14) determines, concisely and directly, the complex dielectric function ϵ_1 of the reflecting second medium in terms of the measured ρ and the angle of incidence ϕ . This accounts for an important application of ellipsometry as a means of determining the optical properties (or *optical constants*) of bulk absorbing materials and opaque films. This approach assumes the absence of a transition layer or a surface film at the two-media interface. If such a film exists, ultrathin as it may be, ϵ_1 as determined by Eq. (14) is called the pseudo dielectric function and is usually written as $\langle \epsilon_1 \rangle$. Figure 3 shows lines of constant ψ and lines of constant Δ in the complex ϵ plane at $\phi = 75^\circ$.

The Three-Phase Model

This often-used model, Fig. 4, consists of a single layer, medium 1, of parallel-plane boundaries which is surrounded by two similar or dissimilar semi-infinite media 0 and 2. The complex amplitude reflection coefficients are given by the Airy-Drude formula^{64,65}

$$R = (r_{01v} + r_{12v} X) / (1 + r_{01v} r_{12v} X) \quad v = p, s \quad (15)$$

$$X = \exp[-j2\pi(d/D_\phi)] \quad (16)$$

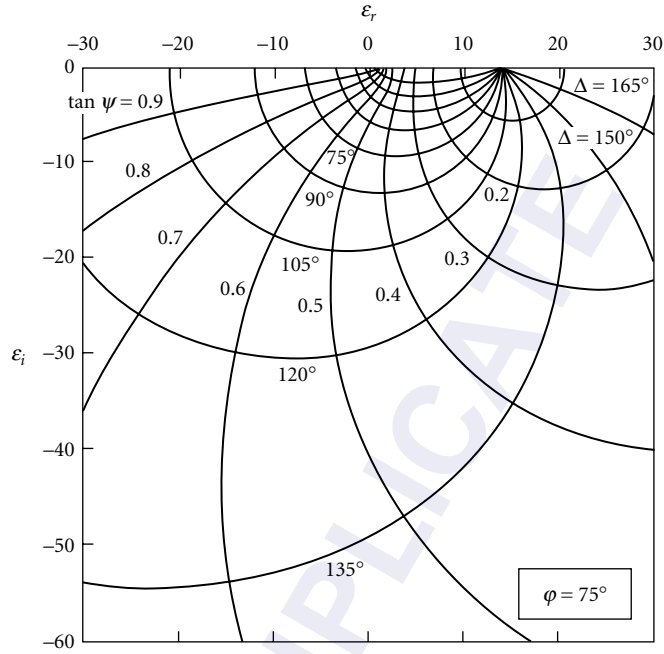


FIGURE 3 Contours of constant $\tan \psi$ and constant Δ in the complex plane of the relative dielectric function ϵ of a transparent medium/absorbing medium interface.

r_{ij}^v is the Fresnel reflection coefficient of the ij interface ($ij = 01$ and 12) for the v polarization, d is the layer thickness, and

$$D_\phi = (\lambda/2)(1/S_1) \tag{17}$$

where λ is the vacuum wavelength of light and S_1 is given by Eq. (12). The ellipsometric function of this system is

$$\rho = (A + BX + CX^2)/(D + EX + FX^2) \tag{18}$$

$$\begin{aligned} A &= r_{01p} & B &= r_{12p} + r_{01p}r_{01s}r_{12s} & C &= r_{12p}r_{01s}r_{12s} \\ D &= r_{01s} & E &= r_{12s} + r_{01p}r_{01s}r_{12p} & F &= r_{12s}r_{01p}r_{12p} \end{aligned} \tag{19}$$

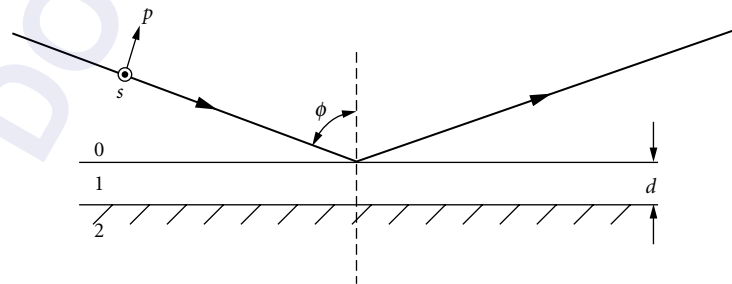


FIGURE 4 Three-phase, ambient-film-substrate system.

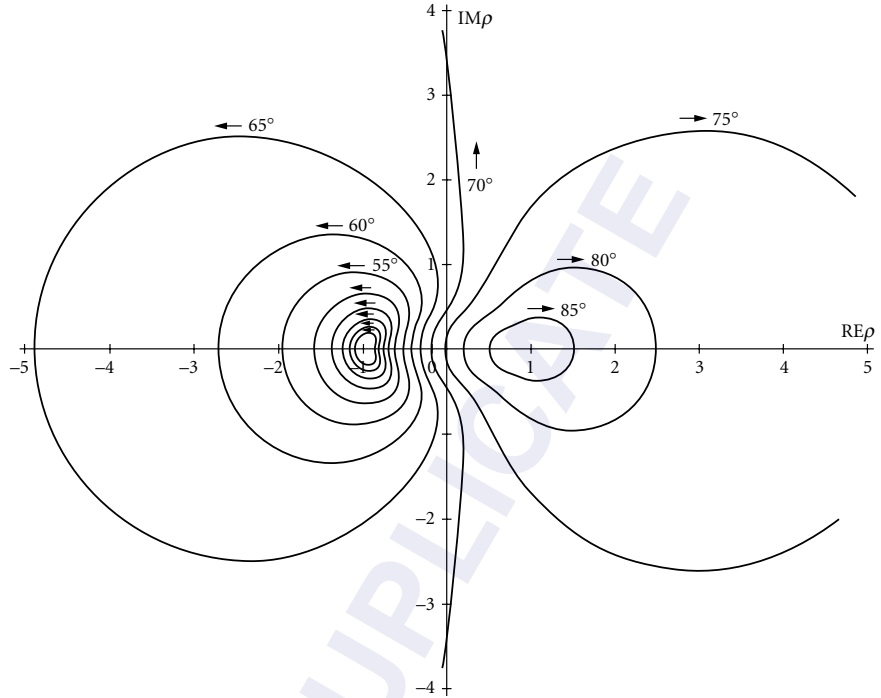


FIGURE 5 Family of constant-angle-of-incidence contours of the ellipsometric function ρ in the complex plane for light reflection in air by the SiO_2/Si film-substrate system at 633-nm wavelength. The contours are for angles of incidence from 30° to 85° in steps of 5° . The arrows indicate the direction of increasing film thickness.⁶⁶

For a transparent film, and with light incident at an angle ϕ such that $\epsilon_1 > \epsilon_0 \sin^2 \phi$ so that total reflection does not occur at the 01 interface, D_ϕ is real, and X , R_p , R_s , and ρ become periodic functions of the film thickness d with period D_ϕ . The locus of X is the unit circle in the complex plane and its multiple images through the conformal mapping of Eq. (18) at different values of ϕ give the constant-angle-of-incidence contours of ρ . Figure 5 shows a family of such contours⁶⁶ for light reflection in air by the SiO_2 -Si system at 633-nm wavelength at angles from 30° to 85° in steps of 5° . Each and every value of ρ , corresponding to all points in the complex plane, can be realized by selecting the appropriate angle of incidence and the SiO_2 film thickness (within a period).

If the dielectric functions of the surrounding media are known, the dielectric function ϵ_1 and thickness d of the film are obtained readily by solving Eq. (18) for X ,

$$X = \{-(B - \rho E) \pm [(B - \rho E)^2 - 4(C - \rho F)(A - \rho D)]^{1/2}\} / 2(C - \rho F) \quad (20)$$

and requiring that^{66,67}

$$|X| = 1 \quad (21)$$

Equation (21) is solved for ϵ_1 as its only unknown by numerical iteration. Subsequently, d is given by

$$d = [-\arg(X) / 2\pi] D_\phi + m D_\phi \quad (22)$$

where m is an integer. The uncertainty of an integral multiple of the film thickness period is often resolved by performing measurements at more than one wavelength or angle of incidence and requiring that d be independent of λ or ϕ .

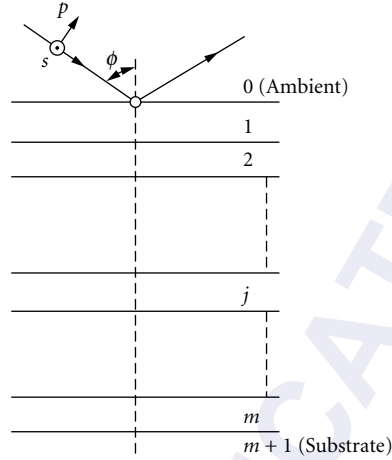


FIGURE 6 Light reflection by a multilayer structure.¹

When the film is absorbing (semitransparent), or the optical properties of one of the surrounding media are unknown, more general inversion methods⁶⁸⁻⁷² are required which are directed toward minimizing an error function of the form

$$f = \sum_{i=1}^N [(\psi_{im} - \psi_{ic})^2 + (\Delta_{im} - \Delta_{ic})^2] \quad (23)$$

where ψ_{im} , ψ_{ic} and Δ_{im} , Δ_{ic} denote the i th measured and calculated values of the ellipsometric angles, and N is the total number of independent measurements.

Multilayer and Graded-Index Films

For an outline of the matrix theory of multilayer systems refer to Chap. 7, "Optical Properties of Films and Coatings," in Vol. IV. For our purposes, we consider a multilayer structure, Fig. 6, that consists of m plane-parallel layers sandwiched between semi-infinite ambient and substrate media (0 and $m + 1$, respectively). The relationships between the field amplitudes of the incident (i), reflected (r), and transmitted (t) plane waves for the p or s polarizations are determined by the scattering matrix equation⁷³

$$\begin{bmatrix} E_i \\ E_r \end{bmatrix} = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix} \begin{bmatrix} E_t \\ 0 \end{bmatrix} \quad (24)$$

The complex-amplitude reflection and transmission coefficients of the entire structure are given by

$$\begin{aligned} R &= E_r/E_i = S_{21}/S_{11} \\ T &= E_t/E_i = 1/S_{11} \end{aligned} \quad (25)$$

The scattering matrix $\mathbf{S} = (S_{ij})$ is obtained as an ordered product of all the interface \mathbf{I} and layer \mathbf{L} matrices of the stratified structure,

$$\mathbf{S} = \mathbf{I}_{01} \mathbf{L}_1 \mathbf{I}_{12} \mathbf{L}_2 \cdots \mathbf{I}_{(j-1)j} \mathbf{L}_j \cdots \mathbf{L}_m \mathbf{I}_{m(m+1)} \quad (26)$$

and the numbering starts from layer 1 (in contact with the ambient) to layer m (adjacent to the substrate) as shown in Fig. 6. The interface scattering matrix is of the form

$$\mathbf{I}_{ab} = (1/t_{ab}) \begin{bmatrix} 1 & r_{ab} \\ r_{ab} & 1 \end{bmatrix} \quad (27)$$

where r_{ab} is the local Fresnel reflection coefficient of the $ab[j(j+1)]$ interface evaluated [using Eqs. (9) and (10) with the appropriate change of subscripts] at an incidence angle in medium a which is related to the external incidence angle ϕ in medium 0 by Snell's law. The associated interface transmission coefficients for the p and s polarizations are

$$t_{abp} = 2(\epsilon_a \epsilon_b)^{1/2} S_a / (\epsilon_b S_a + \epsilon_a S_b) \quad (28)$$

$$t_{abs} = 2S_a / (S_a + S_b)$$

where S_j is defined in Eq. (12). The scattering matrix of the j th layer is

$$\mathbf{L}_j = \begin{bmatrix} Y_j^{-1} & 0 \\ 0 & Y_j \end{bmatrix} \quad (29)$$

$$Y_j = X_j^{1/2} \quad (30)$$

and X_j is given by Eqs. (16) and (17) with the substitution $d = d_j$ for the thickness, and $\epsilon_1 = \epsilon_j$ for the dielectric function of the j th layer.

Except in Eqs. (28), a polarization subscript $v = p$ or s has been dropped for simplicity. In reflection and transmission ellipsometry, the ratios $\rho_r = R_p/R_s$ and $\rho_t = T_p/T_s$ are measured. Inversion for the dielectric functions and thicknesses of some or all of the layers requires extensive data, as may be obtained by VASE, and linear regression analysis to minimize the error function of Eq. (23).

Light reflection and transmission by a graded-index (GRIN) film is handled using the scattering matrix approach described here by dividing the inhomogeneous layer into an adequately large number of sublayers, each of which is approximately homogeneous. In fact, this is the most general approach for a problem of this kind because analytical closed-form solutions are only possible for a few simple refractive-index profiles.⁷⁴⁻⁷⁶

Dielectric Function of a Mixed Phase

For a microscopically inhomogeneous thin film that is a mixture of two materials, as may be produced by coevaporation or cosputtering, or a thin film of one material that may be porous with a significant void fraction (of air), the dielectric function is determined using EMTs.⁴⁶⁻⁵⁴ When the scale of the inhomogeneity is small relative to the wavelength of light, and the domains (or grains) of different dielectrics are of nearly spherical shape, the dielectric function of the mixed phase ϵ is given by

$$\frac{\epsilon - \epsilon_h}{\epsilon + 2\epsilon_h} = v_a \frac{\epsilon_a - \epsilon_h}{\epsilon_a + 2\epsilon_h} + v_b \frac{\epsilon_b - \epsilon_h}{\epsilon_b + 2\epsilon_h} \quad (31)$$

where ϵ_a and ϵ_b are the dielectric functions of the two component phases a and b with volume fractions v_a and v_b and ϵ_h is the host dielectric function. Different EMTs assign different values to ϵ_h . In the Maxwell Garnett EMT,^{47,48} one of the phases, say b , is dominant ($v_b \gg v_a$) and $\epsilon_h = \epsilon_b$. This reduces the second term on the right-hand side of Eq. (31) to zero. In the Bruggeman EMT,⁴⁹ v_a and v_b are comparable, and $\epsilon_h = \epsilon$, which reduces the left-hand side of Eq. (31) to zero.

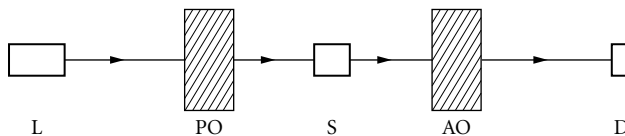


FIGURE 7 Generic ellipsometer with polarizing optics PO and analyzing optics AO. L and D are the light source and photodetector, respectively.

16.5 TRANSMISSION ELLIPSOMETRY

Although ellipsometry is typically carried out on the reflected wave, it is possible to also monitor the state of polarization of the transmitted wave, when such a wave is available for measurement.^{77–81} For example, by combining reflection and transmission ellipsometry, the thickness and complex dielectric function of an absorbing film between transparent media of the same refractive index (e.g., a solid substrate on one side and an index-matching liquid on the other) can be obtained analytically.^{79,80} Polarized light transmission by a multilayer was discussed previously under “Multilayer and Graded-Index Films.” Transmission ellipsometry can be carried out at normal incidence on optically anisotropic samples to determine such properties as the natural or induced linear, circular, or elliptical birefringence and dichroism. However, this falls outside the scope of this chapter.

16.6 INSTRUMENTATION

Figure 7 is a schematic diagram of a generic ellipsometer. It consists of a source of collimated and monochromatic light L, polarizing optics PO on one side of the sample S, and polarization analyzing optics AO and a (linear) photodetector D on the other side. An apt terminology²⁵ refers to the PO as a polarization state generator (PSG) and the AO plus D as a polarization state detector (PSD).

Figure 8 shows the commonly used polarizer-compensator-sample-analyzer (PCSA) ellipsometer arrangement. The PSG consists of a linear polarizer with transmission-axis azimuth P and a linear retarder, or compensator, with fast-axis azimuth C . The PSD consists of a single linear polarizer, that functions as an analyzer, with transmission-axis azimuth A followed by a photodetector D.

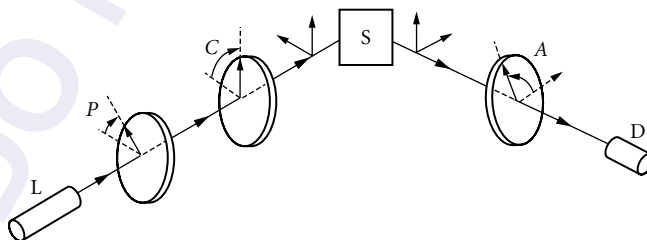


FIGURE 8 Polarizer-compensator-sample-analyzer (PCSA) ellipsometer. The azimuth angles P of the polarizer, C of the compensator (or quarter-wave retarder), and A of the analyzer are measured from the plane of incidence, positive in a counterclockwise sense when looking toward the source.¹

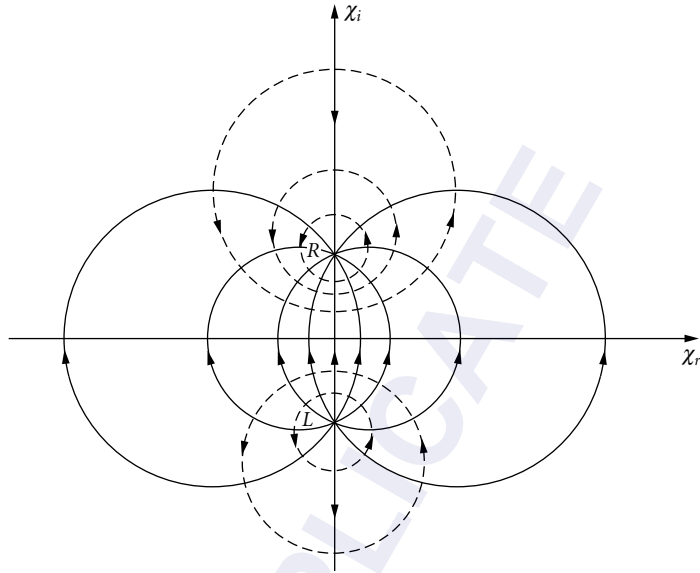


FIGURE 9 Constant C , variable P contours (continuous lines), and constant $P - C$, variable C contours (dashed lines) in the complex plane of polarization for light transmitted by a polarizer-compensator (PC) polarization state generator.¹

All azimuths P , C , and A , are measured from the plane of incidence, positive in a counterclockwise sense when looking toward the source. The state of polarization of the light transmitted by the PSG and incident on S is given by

$$\chi_i = [\tan C + \rho_c \tan(P - C)] / [1 - \rho_c \tan C \tan(P - C)] \quad (32)$$

where $\rho_c = T_{cs}/T_{cf}$ is the ratio of complex amplitude transmittances of the compensator for incident linear polarizations along the slow s and fast f axes. Ideally, the compensator functions as a quarter-wave retarder (QWR) and $\rho_c = -j$. In this case, Eq. (32) describes an elliptical polarization state with major-axis azimuth C and ellipticity angle $-(P - C)$. (The tangent of the ellipticity angle equals the minor-axis-to-major-axis ratio and its sign gives the handedness of the polarization state, positive for right-handed states.) All possible states of total polarization χ_i can be generated by controlling P and C . Figure 9 shows a family of constant C , variable P contours (continuous lines) and constant $P - C$, variable C contours (dashed lines) as orthogonal families of circles in the complex plane of polarization. Figure 10 shows the corresponding contours of constant P and variable C . The points R and L on the imaginary axis at $(0, +1)$ and $(0, -1)$ represent the right- and left-handed circular polarization states, respectively.

Null Ellipsometry

The PCSA ellipsometer of Fig. 8 can be operated in two different modes. In the null mode, the output signal of the photodetector D is reduced to zero (a minimum) by adjusting the azimuth angles P of the polarizer and A of the analyzer with the compensator set at a fixed azimuth C . The choice $C = \pm 45^\circ$ results in rapid convergence to the null. Two independent nulls are reached for each compensator setting. The two nulls obtained with $C = +45^\circ$ are usually referred to as the nulls in zones 2 and 4; those for $C = -45^\circ$ define zones 1 and 3. At null, the reflected polarization is linear and is

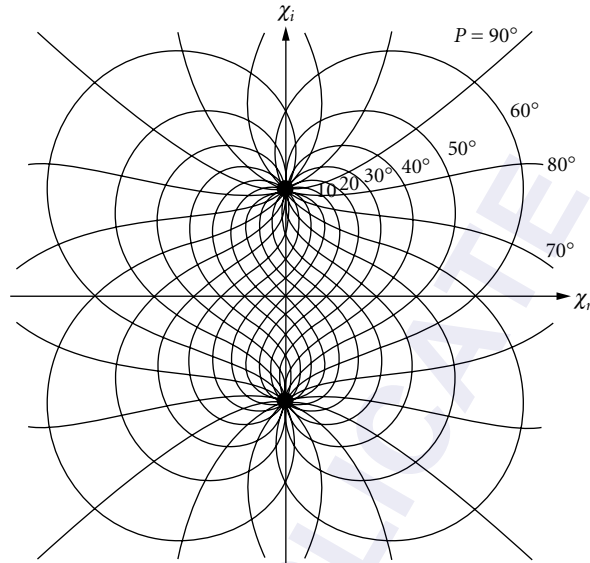


FIGURE 10 Constant P , variable C contours in the complex plane of polarization for light transmitted by a polarizer-compensator (PC) polarization state generator.¹

crossed with the transmission axis of the analyzer; therefore, the reflected state of polarization is given by

$$\chi_r = -\cot A \quad (33)$$

where A is the analyzer azimuth at null. With the incident and reflected polarizations determined by Eqs. (32) and (33), the ratio of complex reflection coefficients of the sample for the p and s linear polarizations ρ is obtained by Eq. (2). Whereas a single null is sufficient to determine ρ in an ideal ellipsometer, results from multiple nulls (in two or four zones) are usually averaged to eliminate the effect of small component imperfections and azimuth-angle errors. Two-zone measurements are also used to determine ρ of the sample and ρ_c of the compensator simultaneously.⁸²⁻⁸⁴ The effects of component imperfections have been considered extensively.⁸⁵

The null ellipsometer can be automated by using stepping or servo motors^{86,87} to rotate the polarizer and analyzer under closed-loop feedback control; the procedure is akin to that of nulling an ac bridge circuit. Alternatively, Faraday cells can be inserted after the polarizer and before the analyzer to produce magneto-optical rotations in lieu of the mechanical rotation of the elements.⁸⁸⁻⁹⁰ This reduces the measurement time of a null ellipsometer from minutes to milliseconds. Large ($\pm 90^\circ$) Faraday rotations would be required for limitless compensation. Small ac modulation is often added for the precise localization of the null.

Photometric Ellipsometry

The polarization state of the reflected light can also be detected photometrically by rotating the analyzer⁹¹⁻⁹⁵ of the PCSA ellipsometer and performing a Fourier analysis of the output signal I of the linear photodetector D . The detected signal waveform is simply given by

$$I = I_0(1 + \alpha \cos 2A + \beta \sin 2A) \quad (34)$$

and the reflected state of polarization is determined from the normalized Fourier coefficients α and β by

$$\chi_r = [\beta \pm (1 - \alpha^2 - \beta^2)^{1/2}] / (1 + \alpha) \quad (35)$$

The sign ambiguity in Eq. (35) indicates that the rotating-analyzer ellipsometer (RAE) cannot determine the handedness of the reflected polarization state. In the RAE, the compensator is not essential and can be removed from the input PO (i.e., the PSA instead of the PCSA optical train is used). Without the compensator, the incident linear polarization is described by

$$\chi_i = \tan P \quad (36)$$

Again, the ratio of complex reflection coefficients of the sample ρ is determined by substituting Eqs. (35) and (36) in Eq. (2). The absence of the wavelength-dependent compensator makes the RAE particularly qualified for SE. The dual of the RAE is the rotating-polarizer ellipsometer which is suited for real-time SE using a spectrograph and a photodiode array that are placed after the fixed analyzer.³¹

A photometric ellipsometer with no moving parts, for fast measurements on the microsecond time scale, employs a photoelastic modulator⁹⁶⁻¹⁰⁰ (PEM) in place of the compensator of Fig. 8. The PEM functions as an oscillating-phase linear retarder in which the relative phase retardation is modulated sinusoidally at a high frequency (typically 50 to 100 kHz) by establishing an elastic ultrasonic standing wave in a transparent solid. The output signal of the photodetector is represented by an infinite Fourier series with coefficients determined by Bessel functions of the first kind and argument equal to the retardation amplitude. However, only the dc, first, and second harmonics of the modulation frequency are usually detected (using lock-in amplifiers) and provide sufficient information to retrieve the ellipsometric parameters of the sample.

Numerous other ellipsometers have been introduced²⁵ that employ more elaborate PSDs. For example Fig. 11 shows a family of rotating-element photopolarimeters²⁵ (REP) that includes the RAE. The column on the right represents the Stokes vector and the fat dots identify the Stokes

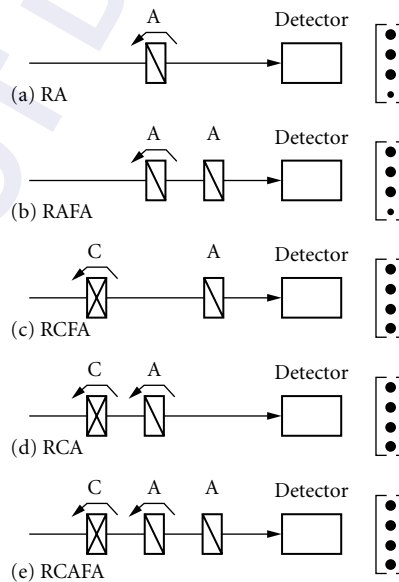


FIGURE 11 Family of rotating-element photopolarimeters (REP) and the Stokes parameters that they can determine.²⁵

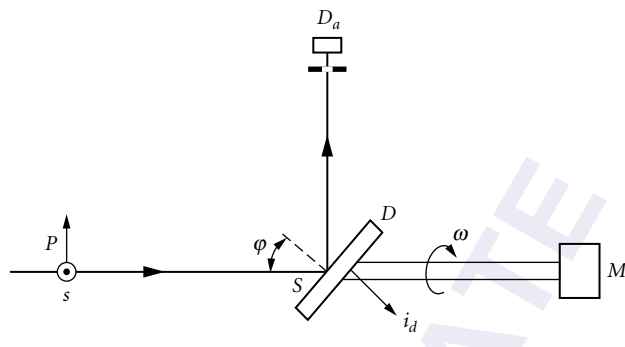


FIGURE 12 Rotating-detector ellipsometer (RODE).¹⁰²

parameters that are measured. (For a discussion of the Stokes parameters, see Chap. 12 in this volume of the *Handbook*.) The simplest complete REP, that can determine all four Stokes parameters of light, is the rotating-compensator fixed-analyzer (RCFA) photopolarimeter originally invented to measure skylight polarization.¹⁰¹ The simplest handedness-blind REP for totally polarized light is the rotating-detector ellipsometer^{102,103} (RODE), Fig. 12, in which the tilted and partially reflective front surface of a solid-state (e.g., Si) detector performs as polarization analyzer.

Ellipsometry Using Four-Detector Photopolarimeters

A new class of fast PSDs that measure the general state of partial or total polarization of a quasi-monochromatic light beam is based on the use of four photodetectors. Such PSDs employ the division of wavefront, the division of amplitude, or a hybrid of the two, and do not require any moving parts or modulators. Figure 13 shows a division-of-wavefront photopolarimeter (DOWP)¹⁰⁴ for

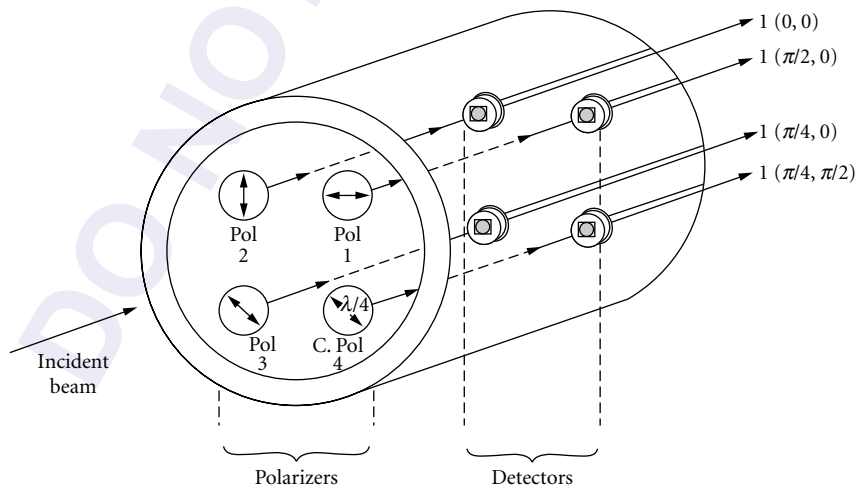


FIGURE 13 Division-of-wavefront photopolarimeter for the simultaneous measurement of all four Stokes parameters of light.¹⁰⁴

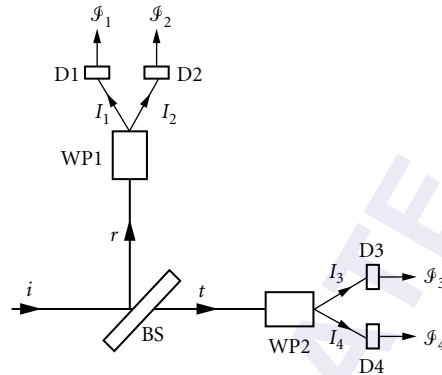


FIGURE 14 Division-of-amplitude photopolarimeter (DOAP) for the simultaneous measurement of all four Stokes parameters of light.¹⁰⁷

performing ellipsometry with nanosecond laser pulses. The DOWP has been adopted recently in commercial automatic polarimeters for the fiber-optics market.^{105,106}

Figure 14 shows a division-of-amplitude photopolarimeter^{107,108} (DOAP) with a coated beam splitter BS and two Wollaston prisms WP1 and WP2, and Fig. 15 represents a recent implementation¹⁰⁹ of that technique. The multiple-beam-splitting and polarization-altering properties of grating diffraction are also well-suited for the DOAP.^{110,111}

The simplest DOAP consists of a spatial arrangement of four solid-state photodetectors Fig. 16, and no other optical elements. The first three detectors (D_0 , D_1 , and D_2) are partially specularly reflecting and the fourth (D_3) is antireflection-coated. The incident light beam is steered in such a way that the plane of incidence is rotated between successive oblique-incidence reflections, hence

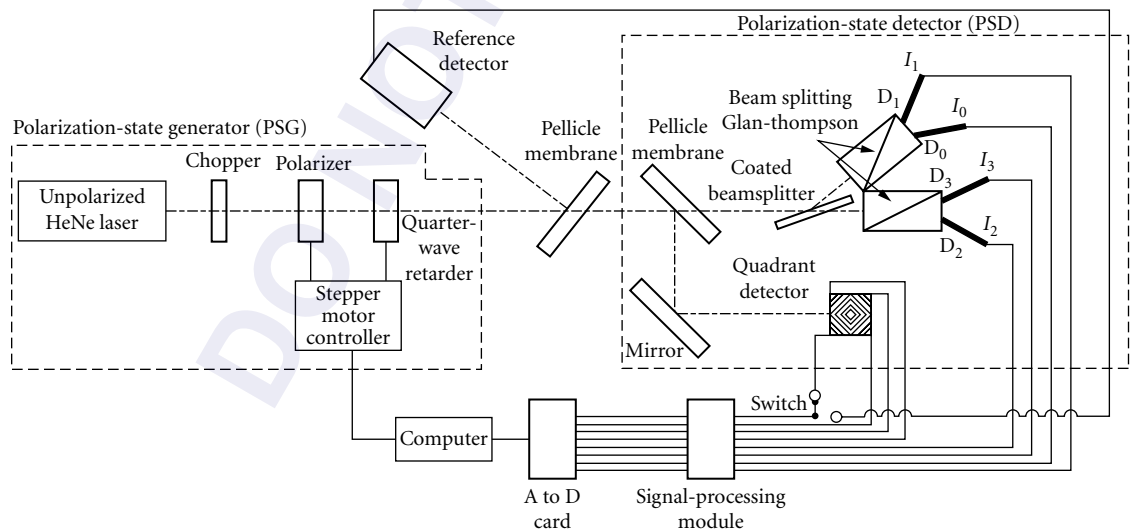


FIGURE 15 Recent implementation of DOAP.¹⁰⁹

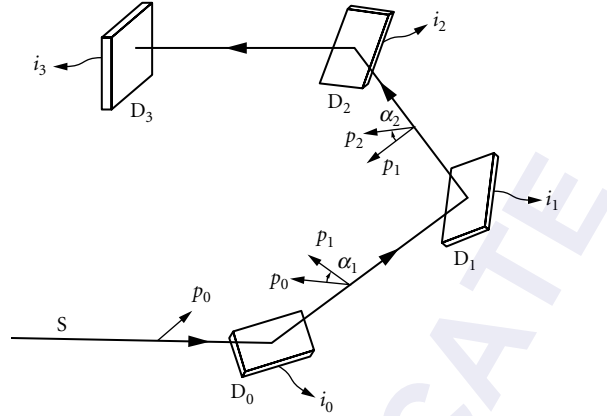


FIGURE 16 Four-detector photopolarimeter for the simultaneous measurement of all four Stokes parameters of light.¹¹²

the light path is nonplanar. In this four-detector photopolarimeter^{112–117} (FDP), and in other DOAPs, the four output signals of the four linear photodetectors define a current vector $\mathbf{I} = [I_0 \ I_1 \ I_2 \ I_3]^t$ which is linearly related,

$$\mathbf{I} = \mathbf{A}\mathbf{S} \quad (37)$$

to the Stokes vector $\mathbf{S} = [S_0 \ S_1 \ S_2 \ S_3]^t$ of the incident light, where t indicates the matrix transpose. The 4×4 instrument matrix \mathbf{A} is determined by calibration¹¹⁵ (using a PSG that consists of a linear polarizer and a quarter-wave retarder). Once \mathbf{A} is determined, \mathbf{S} is obtained from the output signal vector by

$$\mathbf{S} = \mathbf{A}^{-1}\mathbf{I} \quad (38)$$

where \mathbf{A}^{-1} is the inverse of \mathbf{A} . When the light under measurement is totally polarized (i.e., $S_0^2 = S_1^2 + S_2^2 + S_3^2$), the associated complex polarization number is determined in terms of the Stokes parameters as¹¹⁸

$$\chi = (S_2 + jS_3)/(S_0 + S_1) = (S_0 - S_1)/(S_2 - jS_3) \quad (39)$$

For further information on polarimetry, see Chap. 15 in this volume.

Ellipsometry Based on Azimuth Measurements Alone

Measurements of the azimuths of the elliptic vibrations of the light reflected from an optically isotropic surface, for two known vibration directions of incident linearly polarized light, enable the ellipsometric parameters of the surface to be determined at any angle of incidence. If θ_i and θ_r represent the azimuths of the incident linear and reflected elliptical polarizations, respectively, then^{119–121}

$$\tan 2\theta_r = (2 \tan \theta_i \tan \psi \cos \Delta) / (\tan^2 \psi - \tan^2 \theta_i) \quad (40)$$

A pair of measurements $(\theta_{i1}, \theta_{r1})$ and $(\theta_{i2}, \theta_{r2})$ determines ψ and Δ via Eq. (40). The azimuth of the reflected polarization is measured precisely by an ac-null method using an ac-excited Faraday

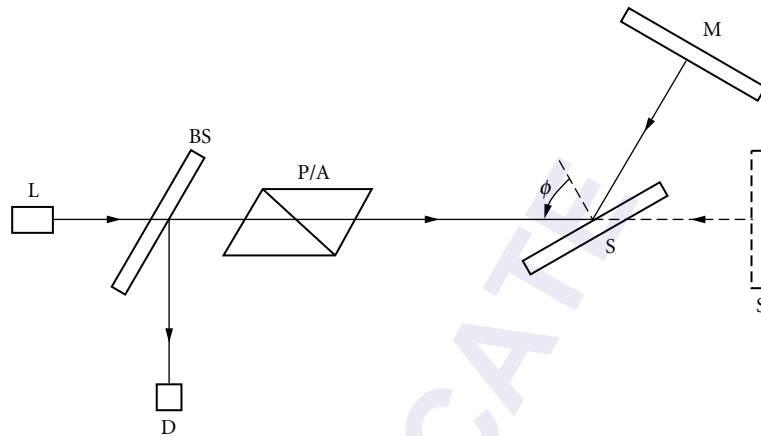


FIGURE 17 Return-path ellipsometer. The dashed lines indicate the configuration for perpendicular-incidence ellipsometry on optically anisotropic samples.¹²⁶

cell followed by a linear analyzer.¹¹⁹ The analyzer is rotationally adjusted to zero the fundamental-frequency component of the detected signal; this aligns the analyzer transmission axis with the minor or major axis of the reflected polarization ellipse.

Return-Path Ellipsometry

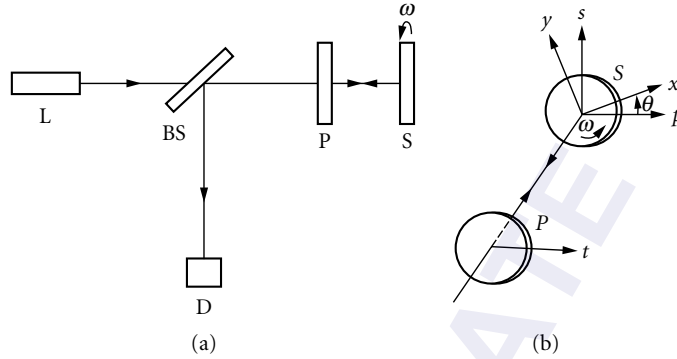
In a return-path ellipsometer (RPE), Fig. 17, an optically isotropic mirror *M* is placed in, and perpendicular to, the reflected beam. This reverses the direction of the beam, so that it retraces its path toward the source with a second reflection at the test surface *S* and second passage through the polarizing/analyzing optics *P/A*. A beam splitter *BS* sends a sample of the returned beam to the photodetector *D*. The RPE can be operated in the null or photometric mode.

In the simplest RPE,^{122,123} the *P/A* optics consists of a single linear polarizer whose azimuth and the angle of incidence are adjusted for a zero detected signal. At null, the angle of incidence is the principal angle, hence $\Delta = \pm 90^\circ$, and the polarizer azimuth equals the principal azimuth, so that the incident linearly polarized light is reflected circularly polarized. Null can also be obtained at a general and fixed angle of incidence by adding a compensator to the *P/A* optics. Adjustment of the polarizer azimuth and the compensator azimuth or retardance produces the null.^{124,125} In the photometric mode,¹²⁶ an element of the *P/A* is modulated periodically and the detected signal is Fourier-analyzed to extract ψ and Δ .

RPEs have the following advantages: (1) the same optical elements are used as polarizing and analyzing optics; (2) only one optical port or window is used for light entry into and exit from the chamber in which the sample may be mounted; and (3) the sensitivity to surface changes is increased because of the double reflection at the sample surface.

Perpendicular-Incidence Ellipsometry

Normal-incidence reflection from an optically isotropic surface is accompanied by a trivial change of polarization due to the reversal of the direction of propagation of the beam (e.g., right-handed circularly polarized light is reflected as left-handed circularly polarized). Because this change of polarization is not specific to the surface, it cannot be used to determine the properties of the


FIGURE 18 Normal-incidence rotating-sample ellipsometer (NIRSE).¹²⁸

reflecting structure. This is why ellipsometry of isotropic surfaces is performed at oblique incidence. However, if the surface is optically anisotropic, perpendicular-incidence ellipsometry (PIE) is possible and offers two significant advantages: (1) simpler single-axis instrumentation of the return-path type with common polarizing/analyzing optics, and (2) simpler inversion for the sample optical properties, because the equations that govern the reflection of light at normal incidence are much simpler than those at oblique incidence.^{127,128}

Like RPE, PIE can be performed using null or photometric techniques.^{126–132} For example, Fig. 18 shows a simple normal-incidence rotating-sample ellipsometer¹²⁸ (NIRSE) that is used to measure the ratio of the complex principal reflection coefficients of an optically anisotropic surface S with principal axes x and y . (The incident linear polarizations along these axes are the eigenpolarizations of reflection.) If we define

$$\eta = (R_{xx} - R_{yy}) / (R_{xx} + R_{yy}) \quad (41)$$

then

$$\eta = \{a_2 \pm j[8a_4(1-a_4) - a_2^2]^{1/2}\} / 2(1-a_4) \quad (42)$$

R_{xx} and R_{yy} are the complex-amplitude principal reflection coefficients of the surface, and a_2 and a_4 are the amplitudes of the second and fourth harmonic components of the detected signal normalized with respect to the dc component. From Eq. (41), we obtain

$$\rho = R_{yy} / R_{xx} = (1 - \eta) / (1 + \eta) \quad (43)$$

PIE can be used to determine the optical properties of bare and coated uniaxial and biaxial crystal surfaces.^{127–130,133}

Interferometric Ellipsometry

Ellipsometry using interferometer arrangements with polarizing optical elements has been suggested and demonstrated.^{134–136} Compensators are not required because the relative phase shift is obtained by the unbalance between the two interferometer arms; this offers a distinct advantage for SE. Direct display of the polarization ellipse is possible.^{134–136}

16.7 JONES-MATRIX GENERALIZED ELLIPSOMETRY

For light reflection at an anisotropic surface, the p and s linear polarizations are not, in general, the eigenpolarizations of reflection. Consequently, the reflection of light is no longer described by Eqs. (1). Instead, the Jones (electric) vectors of the reflected and incident waves are related by

$$\begin{bmatrix} E_{rp} \\ E_{rs} \end{bmatrix} = \begin{bmatrix} R_{pp} & R_{ps} \\ R_{sp} & R_{ss} \end{bmatrix} \begin{bmatrix} E_{ip} \\ E_{is} \end{bmatrix} \quad (44)$$

or, more compactly,

$$\mathbf{E}_r = \mathbf{R}\mathbf{E}_i \quad (45)$$

where \mathbf{R} is the nondiagonal reflection Jones matrix. The states of polarization of the incident and reflected waves, described by the complex variables χ_i and χ_r of Eqs. (4), are interrelated by the bilinear transformation^{85,137}

$$\chi_r = (R_{ss}\chi_i + R_{sp}) / (R_{ps}\chi_i + R_{pp}) \quad (46)$$

In generalized ellipsometry (GE), the incident wave is polarized in at least three different states (χ_{i1} , χ_{i2} , χ_{i3}) and the corresponding states of polarization of the reflected light (χ_{r1} , χ_{r2} , χ_{r3}) are measured. Equation (46) then yields three equations that are solved for the normalized Jones matrix elements, or reflection coefficients ratios,¹³⁸

$$\begin{aligned} R_{pp}/R_{ss} &= (\chi_{i2} - \chi_{i1}H) / (-\chi_{r1} + \chi_{r2}H) \\ R_{ps}/R_{ss} &= (H-1) / (-\chi_{r1} + \chi_{r2}H) \\ R_{sp}/R_{ss} &= (\chi_{i2}\chi_{r1} - \chi_{i1}\chi_{r2}H) / (-\chi_{r1} + \chi_{r2}H) \\ H &= (\chi_{r3} - \chi_{r1})(\chi_{i3} - \chi_{i2}) / (\chi_{i3} - \chi_{i1})(\chi_{r3} - \chi_{r2}) \end{aligned} \quad (47)$$

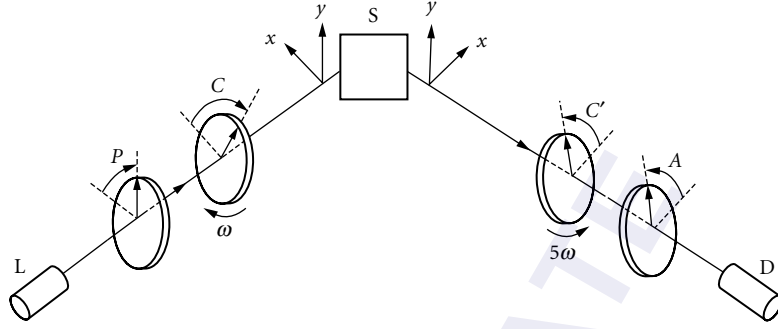
Therefore, the nondiagonal Jones matrix of any optically anisotropic surface is determined, up to a complex constant multiplier, from the mapping of three incident polarizations into the corresponding three reflected polarizations. A PCSA null ellipsometer can be used. The incident polarization χ_i is given by Eq. (32) and the reflected polarization χ_r is given by Eq. (33). Alternatively, the Stokes parameters of the reflected light can be measured using the RCFA photopolarimeter, the DOAP, or the FDP, and χ_r is obtained from Eq. (39). More than three measurements can be taken to overdetermine the normalized Jones matrix elements and reduce the effect of component imperfections and measurement errors. GE can be performed based on azimuth measurements alone.¹³⁹ The main application of GE has been the determination of the optical properties of crystalline materials.^{138–143}

16.8 MUELLER-MATRIX GENERALIZED ELLIPSOMETRY

The most general representation of the transformation of the state of polarization of light upon reflection or scattering by an object or sample is described by¹

$$\mathbf{S}' = \mathbf{M}\mathbf{S} \quad (48)$$

where \mathbf{S} and \mathbf{S}' are the Stokes vectors of the incident and scattered radiation, respectively, and \mathbf{M} is the real 4×4 Mueller matrix that succinctly characterizes the linear (or elastic) light-sample interaction.


FIGURE 19 Dual-rotating-retarder Mueller-matrix photopolarimeter.¹⁴⁵

For light reflection at an optically isotropic and specular (smooth) surface, the Mueller matrix assumes the simple form¹⁴⁴

$$\mathbf{M} = r \begin{bmatrix} 1 & a & 0 & 0 \\ a & 1 & 0 & 0 \\ 0 & 0 & b & c \\ 0 & 0 & -c & b \end{bmatrix} \quad (49)$$

In Eq. (49), r is the surface power reflectance for incident unpolarized or circularly polarized light, and a, b, c are determined by the ellipsometric parameters ψ and Δ as:

$$a = -\cos 2\psi \quad b = \sin 2\psi \cos \Delta \quad \text{and} \quad c = \sin 2\psi \sin \Delta \quad (50)$$

and satisfy the identity $a^2 + b^2 + c^2 = 1$.

In general (i.e., for an optically anisotropic and rough surface), all 16 elements of \mathbf{M} are nonzero and independent.

Several methods for Mueller matrix measurements have been developed.^{25,145-149} An efficient scheme¹⁴⁵⁻¹⁴⁷ uses the PCSCA ellipsometer with symmetrical polarizing (PC) and analyzing ($C'A$) optics, Fig. 19. All 16 elements of the Mueller matrix are encoded onto a single periodic detected signal by rotating the quarter-wave retarders (or compensators) C and C' at angular speeds in the ratio 1:5. The output signal waveform is described by the Fourier series

$$I = a_0 + \sum_{n=1}^{12} (a_n \cos nC + b_n \sin nC) \quad (51)$$

where C is the fast-axis azimuth of the slower of the two retarders, measured from the plane of incidence. Table 1 gives the relations between the signal Fourier amplitudes and the elements of the

TABLE 1 Relations Between Signal Fourier Amplitudes and Elements of the Scaled Mueller Matrix \mathbf{M}'

n	0	1	2	3	4	5	6
a_n	$m'_{11} + \frac{1}{2}m'_{12}$ $+\frac{1}{2}m'_{21} + \frac{1}{4}m'_{22}$	0	$\frac{1}{2}m'_{12} + \frac{1}{4}m'_{22}$	$-\frac{1}{4}m'_{43}$	$-\frac{1}{2}m'_{44}$	0	$\frac{1}{2}m'_{44}$
b_n		$m'_{14} + \frac{1}{2}m'_{24}$	$\frac{1}{2}m'_{13} + \frac{1}{4}m'_{23}$	$-\frac{1}{4}m'_{42}$	0	$-m'_{41} - \frac{1}{2}m'_{42}$	0
n	7	8	9	10	11	12	
a_n	$\frac{1}{4}m'_{43}$	$\frac{1}{8}m'_{22} + \frac{1}{8}m'_{33}$	$\frac{1}{4}m'_{34}$	$\frac{1}{2}m'_{21} + \frac{1}{4}m'_{22}$	$-\frac{1}{4}m'_{34}$	$\frac{1}{8}m'_{22} - \frac{1}{8}m'_{33}$	
b_n	$-\frac{1}{4}m'_{42}$	$-\frac{1}{8}m'_{23} + \frac{1}{8}m'_{32}$	$-\frac{1}{4}m'_{24}$	$\frac{1}{2}m'_{31} + \frac{1}{4}m'_{32}$	$\frac{1}{4}m'_{24}$	$\frac{1}{8}m'_{23} + \frac{1}{8}m'_{32}$	

The transmission axes of the polarizer and analyzer are assumed to be set at 0 azimuth, parallel to the scattering plane or the plane of incidence.

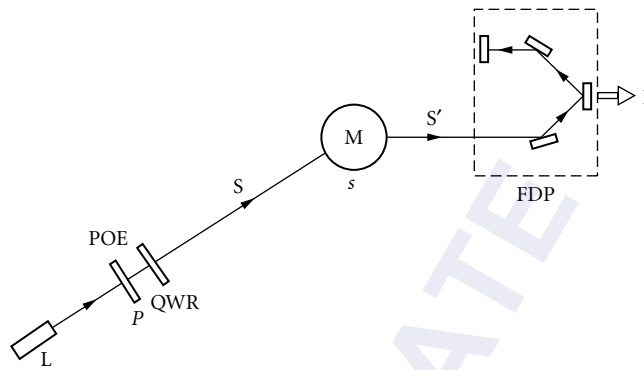


FIGURE 20 Scheme for Mueller-matrix measurement using the four-detector photopolarimeter.¹⁵²

Mueller matrix \mathbf{M}' which differs from \mathbf{M} only by a scale factor. Inasmuch as only the normalized Mueller matrix, with unity first element, is of interest, the unknown scale factor is immaterial. This dual-rotating-retarder Mueller-matrix photopolarimeter has been used to characterize rough surfaces¹⁵⁰ and the retinal nerve-fiber layer.¹⁵¹

Another attractive scheme for Mueller-matrix measurement is shown in Fig. 20. The FDP (or equivalently, any other DOAP) is used as the PSD. Fourier analysis of the output current vector of the FDP, $\mathbf{I}(C)$, as a function of the fast-axis azimuth C of the QWR of the input PO readily determines the Mueller matrix \mathbf{M} , column by column.^{152,153}

16.9 APPLICATIONS

The applications of ellipsometry are too numerous to try to cover in this chapter. The reader is referred to the books and review articles listed in the bibliography. Suffice it to mention the general areas of application. These include: (1) measurement of the optical properties of materials in the visible, IR, and near-UV spectral ranges. The materials may be in bulk or thin-film form and may be optically isotropic or anisotropic.^{3,22,27-31} (2) Thin-film thickness measurements, especially in the semiconductor industry.^{2,5,24} (3) Controlling the growth of optical multilayer coatings¹⁵⁴ and quantum wells.^{155,156} (4) Characterization of physical and chemical adsorption processes at the vacuum/solid, gas/solid, gas/liquid, liquid/liquid, and liquid/solid interfaces.^{26,157} (5) Study of the oxidation kinetics of semiconductor and metal surfaces in various gaseous or liquid ambients.¹⁵⁸ (6) Electrochemical investigations of the electrode/electrolyte interface.^{18,19,32} (7) Diffusion and ion implantation in solids.¹⁵⁹ (8) Biological and biomedical applications.^{16,20,151,160}

16.10 REFERENCES

1. R. M. A. Azzam and N. M. Bashara, *Ellipsometry and Polarized Light*, North-Holland, Amsterdam, 1987.
2. K. Riedling, *Ellipsometry for Industrial Applications*, Springer-Verlag, New York, 1988.
3. R. Röseler, *Infrared Spectroscopic Ellipsometry*, Akademie-Verlag, Berlin, 1990.
4. H. G. Tompkins, and W. A. McGahan, *Spectroscopic Ellipsometry and Reflectometry: A User's Guide*, Wiley, New York, 1999.
5. H. G. Tompkins and E. A. Irene (eds.), *Handbook of Ellipsometry*, William Andrew, Norwich, New York, 2005.
6. R. M. A. Azzam (ed.), *Selected Papers on Ellipsometry*, vol. MS 27 of the Milestone Series, SPIE, Bellingham, Wash., 1991.

7. E. Passaglia, R. R. Stromberg, and J. Kruger (eds.), *Ellipsometry in the Measurement of Surfaces and Thin Films*, NBS Misc. Publ. 256, USGPO, Washington, D.C., 1964.
8. N. M. Bashara, A. B. Buckman, and A. C. Hall (eds.), *Recent Developments in Ellipsometry*, Surf. Sci. vol. 16, North-Holland, Amsterdam, 1969.
9. N. M. Bashara and R. M. A. Azzam (eds.), *Proceedings of the Third International Conference on Ellipsometry*, Surf. Sci. vol. 56, North-Holland, Amsterdam, 1976.
10. R. H. Muller, R. M. A. Azzam, and d. E. Aspnes (eds.), *Proceedings of the Fourth International Conference on Ellipsometry*, Surf. Sci. vol. 96, North-Holland, Amsterdam, 1980.
11. *Proceedings of the International Conference on Ellipsometry and Other Optical Methods for Surface and Thin Film Analysis*, J. de Physique, vol. 44, Colloq. C10, Les Editions de Physique, Paris, 1984.
12. A. C. Boccara, C. Pickering, and J. Rivory (eds.), *Proceedings of the First International Conference on Spectroscopic Ellipsometry*, Thin Solid Films, vols. 233 and 234, Elsevier, Amsterdam, 1993.
13. R. W. Collins, D. E. Aspnes, and E. A. Irene (eds.), *Proceedings of the 2nd International Conference on Spectroscopic Ellipsometry*, Thin Solid Films, vols. 313 and 314, Elsevier, Amsterdam, 1998.
14. M. Fried, K. Hingerl, and J. Humlicek (eds.), *Proceedings of the 3rd International Conference on Spectroscopic Ellipsometry*, Thin Solid Films, vols. 455 and 456, Elsevier, Amsterdam, 2004.
15. H. Arwin, U. Beck, and M. Schubert (eds.), *Proceedings of the 4th International Conference on Spectroscopic Ellipsometry*, Wiley-VCH, Weinheim, 2008.
16. G. Poste and C. Moss, "The Study of Surface Reactions in Biological Systems by Ellipsometry," in S. G. Davison (ed.), *Progress in Surface Science*, vol. 2, pt. 3, Pergamon, New York, 1972, pp. 139–232.
17. R. H. Muller, "Principles of Ellipsometry," in R. H. Mueller (ed.), *Advances in Electrochemistry and Electrochemical Engineering*, vol. 9, Wiley, New York, 1973, pp. 167–226.
18. J. Kruger, "Application of Ellipsometry in Electrochemistry," in R. H. Muller (ed.), *Advances in Electrochemistry and Electrochemical Engineering*, vol. 9, Wiley, New York, 1973, pp. 227–280.
19. W.-K. Paik, "Ellipsometric Optics with Special Reference to Electrochemical Systems," in J. O'M. Bockris (ed.), *MTP International Review of Science, Physical Chemistry*, series 1, vol. 6, Butterworths, Univ. Park, Baltimore, 1973, pp. 239–285.
20. A. Rothen, "Ellipsometric Studies of Thin Films," in D. A. Cadenhead, J. F. Danielli, and M. D. Rosenberg (eds.), *Progress in Surface and Membrane Science*, vol. 8, Academic, New York, 1974, pp. 81–118.
21. R. H. Muller, "Present Status of Automatic Ellipsometers," *Surf. Sci.* **56**:19–36 (1976).
22. D. E. Aspnes, "Spectroscopic Ellipsometry of Solids," in B. O. Seraphin (ed.), *Optical Properties of Solids: New Developments*, North-Holland, Amsterdam, 1976, pp. 799–846.
23. W. E. J. Neal, "Application of Ellipsometry to Surface Films and Film Growth," *Surf. Technol.* **6**:81–110 (1977).
24. A. V. Rzhanov and K. K. Svitashv, "Ellipsometric Techniques to Study Surfaces and Thin Films," in L. Marton and C. Marton (eds.), *Advances in Electronics and Electron Physics*, vol. 49, Academic, New York, 1979, pp. 1–84.
25. P. S. Hauge, "Recent Developments in Instrumentation in Ellipsometry," *Surf. Sci.* **96**:108–140 (1980).
26. F. H. P. M. Habraken, O. L. J. Gijzeman, and G. A. Bootsma, "Ellipsometry of Clean Surfaces, Submonolayer and Monolayer Films," *Surf. Sci.* **96**:482–507 (1980).
27. D. E. Aspnes "Characterization of Materials and Interfaces by Visible-Near UV Spectrophotometry and Ellipsometry," *J. Mat. Educ.* **7**:849–901 (1985).
28. P. J. McMarr, K. Vedam, and J. Narayan, "Spectroscopic Ellipsometry: A New Tool for Nondestructive Depth Profiling and Characterization of Interfaces," *J. Appl. Phys.* **59**:694–701 (1986).
29. D. E. Aspnes, "Analysis of Semiconductor Materials and Structures by Spectroellipsometry," *SPIE Proc.* **946**:84–97 (1988).
30. R. Drevillon, "Spectroscopic Ellipsometry of Ultrathin Films: From UV to IR," *Thin Solid Films* **163**:157–166 (1988).
31. R. W. Collins and Y.-T. Kim, "Ellipsometry for Thin-Film and Surface Analysis," *Anal. Chem.* **62**:887A–900A (1990).
32. R. H. Muller, "Ellipsometry as an In Situ Probe for the Study of Electrode Processes," in R. Varma and J. R. Selman (eds.), *Techniques for Characterization of Electrode Processes*, Wiley, New York, 1991.

33. A. Rothen, *Rev. Sci. Instrum.* **16**:26–30 (1945).
34. A. Rothen, in Ref. 7, pp. 7–21.
35. A. C. Hall, *Surf. Sci.* **16**:1–13 (1969).
36. R. M. A. Azzam and N. M. Bashara, Ref. 1, sec. 2.6.1.
37. R. M. A. Azzam and N. M. Bashara, Ref. 1, sec. 1.7.
38. M. M. Ibrahim and N. M. Bashara, *J. Opt. Soc. Am.* **61**:1622–1629 (1971).
39. O. Hunderi, *Surface Sci.* **61**:515–520 (1976).
40. J. Humlíček, *J. Opt. Soc. Am. A* **2**:713–722 (1985).
41. Y. Gaillyová, E. Schmidt, and J. Humlíček, *J. Opt. Soc. Am. A* **2**:723–726 (1985).
42. W. H. Weedon, S. W. McKnight, and A. J. Devaney, *J. Opt. Soc. Am. A* **8**:1881–1891 (1991).
43. J. A. Woollam Co., Lincoln, NE 68508.
44. R. H. Muller, *Surf. Sci.* **16**:14–33 (1969).
45. E. D. Palik, *Handbook of Optical Constants of Solids*, Academic, New York, 1985, p. 294.
46. L. Lorenz, *Ann. Phys. Chem.* (Leipzig) **11**:70–103 (1880).
47. J. C. Maxwell Garnett, *Philos. Trans. R. Soc. London* **203**:385–420 (1904).
48. J. C. Maxwell Garnett, *Philos. Trans. R. Soc. London* **205**:237–288 (1906).
49. D. A. G. Bruggeman, *Ann. Phys.* (Leipzig) **24**:636–679 (1935).
50. C. G. Granqvist and O. Hunderi, *Phys. Rev.* **B18**:1554–1561 (1978).
51. C. Grosse and J.-L. Greffe, *J. Chim. Phys.* **76**:305–327 (1979).
52. D. E. Aspnes, J. B. Theeten and F. Hottier, *Phys. Rev.* **B20**:3292–3302 (1979).
53. D. E. Aspnes, *Am. J. Phys.* **50**:704–709 (1982).
54. D. E. Aspnes, *Physica* **117B/118B**:359–361 (1983).
55. F. Abelès *Ann. de Physique* **5**:596–640 (1950).
56. P. Drude, *Theory of Optics*, Dover, New York, 1959.
57. J. R. Wait, *Electromagnetic Waves in Stratified Media*, Pergamon, New York, 1962.
58. O. S. Heavens, *Optical Properties of Thin Solid Films*, Dover, New York, 1965.
59. Z. Knittl, *Optics of Thin Films*, Wiley, New York, 1976.
60. J. Lekner, *Theory of Reflection*, Marinus Nijhoff, Dordrecht, 1987.
61. E. S. Keeping, *Introduction to Statistical Inference*, Van Nostrand, Princeton, 1962, chap. 12.
62. D. W. Marquardt, *SIAM J. Appl. Math.* **11**:431–441 (1963).
63. J. R. Rice, in *Numerical Solutions of Nonlinear Problems*, Computer Sci. Center, University of Maryland, College Park, 1970.
64. G. B. Airy, *Phil. Mag.* **2**:20 (1833).
65. P. Drude, *Annal. Phys. Chem.* **36**:865–897 (1889).
66. R. M. A. Azzam, A.-R. M. Zaghoul, and N. M. Bashara, *J. Opt. Soc. Am.* **65**:252–260 (1975).
67. A. R. Reinberg, *Appl. Opt.* **11**:1273–1274 (1972).
68. F. L. McCrackin and J. P. Colson, in Ref. 7, pp. 61–82.
69. M. Malin and K. Vedam, *Surf. Sci.* **56**:49–63 (1976).
70. D. I. Bilenko, B. A. Dvorkin, T. Y. Druzhinina, S. N. Krasnobaev, and V. P. Polyanskaya, *Opt. Spedrosc.* **55**:533–536 (1983).
71. G. H. Bu-Abbud, N. M. Bashara and J. A. Woollam, *Thin Solid Films* **138**:27–41 (1986).
72. G. E. Jellison Jr., *Appl. Opt.* **30**:3354–3360 (1991).
73. R. M. A. Azzam and N. M. Bashara, Ref. 1 sec. 4.6, and references cited therein.
74. F. Abelès, Ref. 7, pp. 41–58.
75. J. C. Charmet and P. G. de Gennes, *J. Opt. Soc. Am.* **73**:1777–1784 (1983).
76. J. Lekner, Ref. 60, chap. 2.

77. R. M. A. Azzam, M. Elshazly-Zaghloul, and N. M. Bashara, *Appl. Opt.* **14**:1652–1663 (1975).
78. I. Ohlídal and F. Lukeš, *Thin Solid Films* **85**:181–190 (1981).
79. R. M. A. Azzam, *J. Opt. Soc. Am.* **72**:1439–1440 (1982).
80. R. M. A. Azzam, Ref. 11, pp. 67–70.
81. I. Ohlídal and F. Lukeš, *Thin Solid Films* **115**:269–282 (1984).
82. F. L. McCrackin, *J. Opt. Soc. Am.* **60**:57–63 (1970).
83. J. A. Johnson and N. M. Bashara, *J. Opt. Soc. Am.* **60**:221–224 (1970).
84. R. M. A. Azzam and N. M. Bashara, *J. Opt. Soc. Am.* **62**:222–229 (1972).
85. R. M. A. Azzam and N. M. Bashara, Ref. 1, sees. 3.7 and 3.8 and references cited therein.
86. H. Takasaki, *Appl. Opt.* **5**:759–764 (1966).
87. J. L. Ord, *Appl. Opt.* **6**:1673–1679 (1967).
88. A. B. Winterbottom, Ref. 7, pp. 97–112.
89. H. J. Mathiu, D. E. McClure, and R. H. Muller, *Rev. Sci. Instrum.* **45**:798–802 (1974).
90. R. H. Muller and J. C. Farmer, *Rev. Sci. Instrum.* **55**:371–374 (1984).
91. W. Budde, *Appl. Opt.* **1**:201–205 (1962).
92. B. D. Cahan and R. F. Spanier, *Surf. Sci.* **16**:166–176 (1969).
93. R. Greef, *Rev. Sci. Instrum.* **41**:532–538 (1970).
94. P. S. Hauge and F. H. Dill, *IBM J. Res. Develop.* **17**:472–489 (1973).
95. D. E. Aspnes and A. A. Studna, *Appl. Opt.* **14**:220–228 (1975).
96. M. Billardon and J. Badoz, *C. R. Acad. Sci. (Paris)* **262**:1672 (1966).
97. J. C. Kemp, *J. Opt. Soc. Am.* **59**:950–954 (1969).
98. S. N. Jasperson and S. E. Schnatterly, *Rev. Sci. Instrum.* **40**:761–767 (1969).
99. J. Badoz, M. Billardon, J. C. Canit, and M. F. Russel, *J. Opt. (Paris)* **8**:373–384 (1977).
100. V. M. Bermudez and V. H. Ritz, *Appl. Opt.* **17**:542–552 (1978).
101. Z. Sekera, *J. Opt. Soc. Am.* **47**:484–490 (1957).
102. R. M. A. Azzam, *Opt. Lett.* **10**:427–429 (1985).
103. D. C. Nick and R. M. A. Azzam, *Rev. Sci. Instrum.* **60**:3625–3632 (1989).
104. E. Collett, *Surf. Sci.* **96**:156–167 (1980).
105. Lightwave Polarization Analyzer Systems, Hewlett Packard Co., Palo Alto, California 94303.
106. Polarscope, Electro Optic Developments Ltd, Basildon, Essex SS14 3BE, England.
107. R. M. A. Azzam, *Opt. Acta* **29**:685–689 (1982).
108. R. M. A. Azzam, *Opt. Acta* **32**:1407–1412 (1985).
109. S. Krishnan, *J. Opt. Soc. Am. A* **9**:1615–1622 (1992).
110. R. M. A. Azzam, *Appl. Opt.* **31**:3574–3576 (1992).
111. R. M. A. Azzam and K. A. Giardina, *J. Opt. Soc. Am. A* **10**:1190–1196 (1993).
112. R. M. A. Azzam, *Opt. Lett.* **10**:309–311 (1985); U.S. Patent 4,681,450, July 21, 1987.
113. R. M. A. Azzam, E. Masetti, I. M. Elminyawi, and F. G. Grosz, *Rev. Sci. Instrum.* **59**:84–88 (1988).
114. R. M. A. Azzam, I. M. Elminyawi, and A. M. El-Saba, *J. Opt. Soc. Am. A* **5**:681–689 (1988).
115. R. M. A. Azzam and A. G. Lopez, *J. Opt. Soc. Am. A* **6**:1513–1521 (1989).
116. R. M. A. Azzam, *J. Opt. Soc. Am. A* **7**:87–91 (1990).
117. The Stokesmeter, Gaertner Scientific Co., Chicago, Illinois 60614.
118. P. S. Hauge, R. H. Muller, and C. G. Smith, *Surf. Sci.* **96**:81–107 (1980).
119. J. Monin and G.-A. Boutry, *Nouv. Rev. Opt.* **4**:159–169 (1973).
120. C. Som and C. Chowdhury, *J. Opt. Soc. Am.* **62**:10–15 (1972).
121. S. I. Idnurm, *Opt. Spectrosc.* **42**:210–212 (1977).

122. H. M. O'Bryan, *J. Opt. Soc. Am.* **26**:122–127 (1936).
123. M. Yamamoto, *Opt. Commun.* **10**:200–202 (1974).
124. T. Yamaguchi and H. Takahashi, *Appl. Opt.* **15**:677–680 (1976).
125. R. M. A. Azzam, *Opt. Acta* **24**:1039–1049 (1977).
126. R. M. A. Azzam, *J. Opt. (Paris)* **9**:131–134 (1978).
127. R. M. A. Azzam, *Opt. Commun.* **19**:122–124 (1976); *Opt. Commun.* **20**:405–408 (1977)
128. Y. Cui and R. M. A. Azzam, *Appl. Opt.* **35**:2235–2238 (1996).
129. R. H. Young and E. I. P. Walker, *Phys. Rev. B* **15**:631–637 (1977).
130. D. W. Stevens, *Surf. Sci.* **96**:174–201 (1980).
131. R. M. A. Azzam, *J. Opt. (Paris)* **12**:317–321 (1981).
132. R. M. A. Azzam, *Opt. Eng.* **20**:58–61 (1981).
133. R. M. A. Azzam, *Appl. Opt.* **19**:3092–3095 (1980).
134. A. L. Dmitriev, *Opt. Spectrosc.* **32**:96–99 (1972).
135. H. F. Hazebroek and A. A. Holscher, *J. Phys. E: Sci. Instrum.* **6**:822–826 (1973).
136. R. Calvani, R. Caponi, and F. Cisternino, *J. Light. Technol.* **LT4**:877–883 (1986).
137. R. M. A. Azzam and N. M. Bashara, *J. Opt. Soc. Am.* **62**:336–340 (1972).
138. R. M. A. Azzam and N. M. Bashara, *J. Opt. Soc. Am.* **64**:128–133 (1974).
139. R. M. A. Azzam, *J. Opt. Soc. Am.* **68**:514–518 (1978).
140. D. J. De Smet, *J. Opt. Soc. Am.* **64**:631–638 (1974); *J. Opt. Soc. Am.* **65**:542–547 (1975).
141. P. S. Hauge, *Surf. Sci.* **56**:148–160 (1976).
142. M. Elshazly-Zaghloul, R. M. A. Azzam, and N. M. Bashara, *Surf. Sci.* **56**:281–292 (1976).
143. M. Schubert, *Theory and Applications of Generalized Ellipsometry*, in Ref. 5, pp. 637–717.
144. R. M. A. Azzam and N. M. Bashara, Ref. 1, p. 491.
145. R. M. A. Azzam, *Opt. Lett.* **2**:148–150 (1978); U. S. Patent 4, 306, 809, December 22, 1981.
146. P. S. Hauge, *J. Opt. Soc. Am.* **68**:1519–1528 (1978).
147. D. H. Goldstein, *Appl. Opt.* **31**:6676–6683 (1992).
148. A. M. Hunt and D. R. Huffman, *Appl. Opt.* **17**:2700–2710 (1978).
149. R. C. Thompson, J. R. Bottiger, and E. S. Fry, *Appl. Opt.* **19**:1323–1332 (1980).
150. D. A. Ramsey, *Thin Film Measurements on Rough Substrates using Mueller-Matrix Ellipsometry*, Ph.D. thesis, The University of Michigan, Ann Arbor, 1985.
151. A. W. Dreher, K. Reiter, and R. N. Weinreb, *Appl. Opt.* **31**:3730–3735 (1992).
152. R. M. A. Azzam, *Opt. Lett.* **11**:270–272 (1986).
153. R. M. A. Azzam, K. A. Giardina, and A. G. Lopez, *Opt. Eng.* **30**:1583–1589 (1991).
154. Ph. Houdy, *Rev. Phys. Appl.* **23**:1653–1659 (1988).
155. J. B. Theeten, F. Hottier, and J. Hallais, *J. Crystal Growth* **46**:245–252 (1979).
156. D. E. Aspnes, W. E. Quinn, M. C. Tamargo, M. A. A. Pudensi, S. A. Schwarz, M. J. S. Brasil, R. E. Nahory, and S. Gregory, *Appl. Phys. Lett.* **60**:1244–1247 (1992).
157. R. M. A. Azzam and N. M. Bashara, Ref. 1, sec. 6.3.
158. R. M. A. Azzam and N. M. Bashara, Ref. 1, sec. 6.4.
159. D. E. Aspnes and A. A. Studna, *Surf. Sci.* **96**:294–306 (1980).
160. H. Arwin, *Ellipsometry in Life Sciences*, in Ref. 5, pp. 799–855.

This page intentionally left blank.

DO NOT DUPLICATE

PART

4

COMPONENTS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

LENSES

R. Barry Johnson

*Consultant
Huntsville, Alabama*

17.1 GLOSSARY

a, b	first and second lenses
AChr	axial chromatic aberration
AST	astigmatism
b	factor
bfl	back focal length
C_o	scaling factor
c	curvature
C_1	scaling factor
C_2	scaling factor
CC	conic constant
CMA_s	sagittal coma
CMA_t	tangential coma
D_{ep}	diameter of entrance pupil
d_o	distance from object to loupe
d_e	distance from loupe to the eye
E	irradiance
efl	effective focal length
ep	eyepiece
FN	F-number
f	focal length
h	height above axis
H_i	height of ray intercept in image plane
K	shape factor
i	image
$J_1()$	Bessel function of the first kind
k	$2\pi/\lambda$

L	length
MP	magnifying power [cf. linear lateral longitudinal magnification]
m	linear, lateral magnification
m_{oc}	nodal-point to optical-center magnification
\bar{m}	linear, longitudinal magnification
MTF	modulation transfer function
M	factor
n	refractive index
NA	numerical aperture
o	object
obj	objective
P	partial dispersion
P_i	principal points
p	$=s_d/f_a$
$\tilde{\mathfrak{K}}$	peak normalized spectral weighting function
s	object to image distance
SA3	third-order spherical aberration
SAC	secondary angular spectrum
s_i	image distance
s_{ot}	optical tube length
s_o	object distance
TPAC	transverse primary chromatic aberration
t	thickness
u	slope
V	Abbe number or reciprocal dispersion
v	φ -normalized reciprocal object distance $1/s_o\varphi$
x, y, z	cartesian coordinates
β	angular blur diameter
δ	depth of focus
ζ	sag
$\Delta\theta$	angular blur tolerance
θ	field of view
λ	wavelength
ν	spatial frequency
ϕ	lens power
r	radius
σ	standard deviation of the irradiance distribution
τ	transmission
Ω	normalized spatial frequency

17.2 INTRODUCTION

This chapter provides a basic understanding of using lenses for image formation and manipulation. The principles of image formation are reviewed first. The effects of lens shape, index of refraction, magnification, and F-number on the image quality of a singlet lens are discussed in some detail. Achromatic doublets and more complex lens systems are covered next. A representative variety of lenses is analyzed

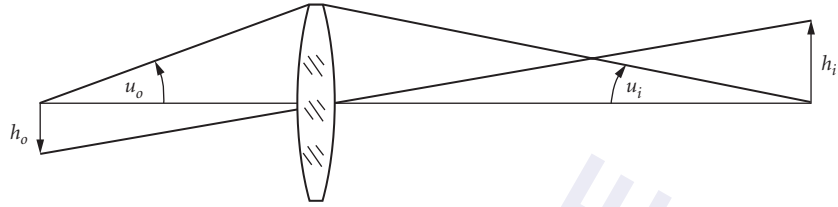


FIGURE 1 Imaging by a simple lens.

and discussed. Performance that may be expected of each class of lens is presented. The section concludes with several techniques for rapid estimation of the performance of lenses. Refer to Chap. 1 “Geometrical Optics,” in this volume for further discussion of geometrical optics and aberrations.

17.3 BASICS

Figure 1 illustrates an image being formed by a simple lens. The object height is h_o and the image height is h_i , with u_o and u_i being the corresponding slope angles. It follows from the Lagrange invariant that the *lateral magnification* is defined to be

$$m \equiv \frac{h_i}{h_o} = \frac{(nu)_o}{(nu)_i} \quad (1)$$

where n_o and n_i are the refractive indices of the medium in which the object and image lie, respectively. By convention, a height is positive if above the optical axis and a ray angle is positive if its slope angle is positive. Distances are positive if the ray propagates left to right. Since the Lagrange invariant is applicable for paraxial rays, the angle nu should be understood to mean $n \tan u$. This interpretation applies to all paraxial computations. For an *aplanatic* lens, which is free of spherical aberration and linear coma, the magnification can by the *optical sine theorem* be given by

$$m \equiv \frac{h_i}{h_o} = \frac{n_o \sin u_o}{n_i \sin u_i} \quad (2)$$

If the object is moved a small distance ∂s_o longitudinally, the corresponding displacement of the image ∂s_i can be found by the differential form of the basic imaging equation and leads to an equation analogous to the Lagrange invariant. The *longitudinal magnification* is then defined as

$$\bar{m} \equiv \frac{\partial s_i}{\partial s_o} = \frac{(nu^2)_o}{(nu^2)_i} = m^2 \left[\frac{n_i}{n_o} \right] \quad (3)$$

The following example will illustrate one application of m and \bar{m} . Consider that a spherical object of radius r_o is to be imaged as shown in Fig. 2. The equation of the object is $r_o^2 = y_o^2 + z^2$, where z is measured along the optical axis and is zero at the object's center of curvature. Letting the surface sag

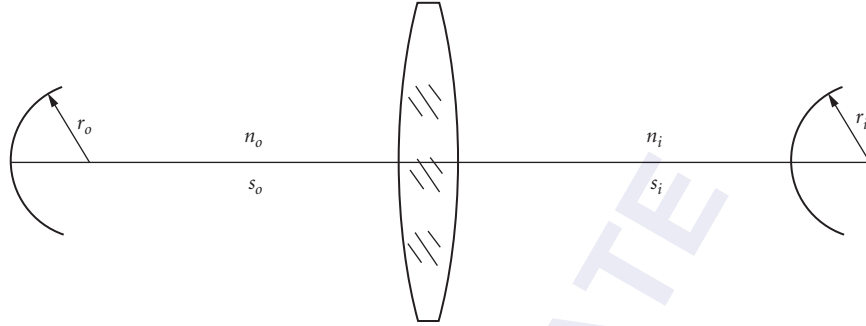


FIGURE 2 Imaging of a spherical object by a lens.

as measured from the vertex plane of the object be denoted as ζ_o , the equation of the object becomes $r_o^2 = (r_o - \zeta_o)^2 + y_o^2$ since $z = r_o - \zeta_o$. In the region near the optical axis, $\zeta_o^2 \ll r_o^2$, which implies that $r_o \approx y_o^2 / 2\zeta_o$. The image of the object is expressed in the transverse or lateral direction by $y_i = m y_o$ and in the longitudinal or axial direction by $\zeta_i = \bar{m} \zeta_o = \zeta_o m^2 (n_i / n_o)$. In a like manner, the image of the spherical object is expressed as $r_i \approx (y_i)^2 / 2\zeta_i$. By substitution, the sag of the image is expressed by

$$r_i \equiv \frac{n_o y_o^2}{2n_i \zeta_o} = r_o \left[\frac{n_o}{n_i} \right] \tag{4}$$

Hence, in the paraxial region about the optical axis, the radius of the image of a spherical object is independent of the magnification and depends only on the ratio of the refractive indices of the object and image spaces.

When an optical system as shown in Fig. 3 images a tilted object, the image will also be tilted. By employing the concept of lateral and longitudinal magnification, it can be easily shown that the

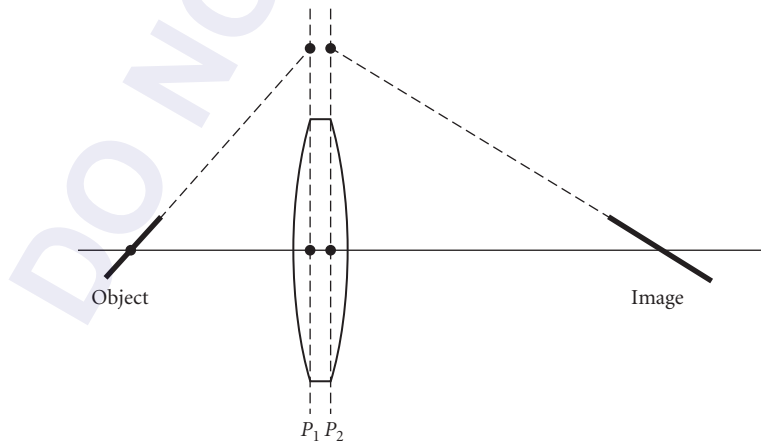


FIGURE 3 Imaging of a tilted object illustrating the Scheimpflug condition.

intersection height of the object plane with the first principal plane P_1 of the lens must be the same as the intersection height of the image plane with the second principal plane P_2 of the lens. This principle is known as the *Scheimpflug condition*.

The object-image relationship of a lens system is often described with respect to its *cardinal points*, which are the *principal and focal points* (conceived by Carl Gauss in 1841) and *nodal points* (conceived by Johann Listing in 1845). Gauss demonstrated that, so far as paraxial rays are concerned, a lens of any degree of complexity can be replaced by its cardinal points, viz., two principal points and two focal points, where the distances from the principal points to their respective focal points being the focal lengths of the lens. This was the first formal definition of the focal length of a lens system. The properties of the cardinal points and related points are as follows:

- *Principal points*: the axial intersection points of conjugate planes P_1 and P_2 where these *principal planes* are related by unit lateral magnification. A ray incident at height h on P_1 will exit P_2 at the same height. When $n_i \neq n_o$, a ray incident at the first principal point with angle u will exit the second principal point with angle $u(n_i/n_o)$.
- *Nodal points*: conjugate points related by unit angular magnification ($m = u_i/u_o$). Nodal points are the axial intersection points of conjugate planes N_1 and N_2 where these *nodal planes* are related, by application of the Lagrange invariant, lateral magnification ($m = n_i/n_o$). A ray incident at height h on N_1 will exit N_2 at $h(n_i/n_o)$.
- *Focal points*: anterior or front (f_1) and posterior or rear (f_2) focal points are the axial intersections of the respective focal planes, which are not conjugate. Any ray parallel to the optical axis, and entering the lens from the left, will intersect the axis at f_2 . Any ray parallel to the optical axis, and entering the lens from the right, will intersect the axis at f_1 . When $n_i \neq n_o$, the distance $f_1 N_1$ equals the posterior focal length $P_2 f_2$ and the distance $N_2 f_2$ equals the anterior focal length $P_1 f_1$.
- *Antiprincipal points*: the axial intersection points of conjugate planes where these *antiprincipal planes* are related by negative unit lateral magnification. An example is a lens used at $m = -1$ where the object and image planes are located at twice the focal length from the principal points.
- *Antinodal points*: conjugate points related by negative unit angular magnification ($m = -u_i/u_o$).

Although imaging can be done solely using n_i/n_o and either the principal point and planes, or the nodal point and planes, it is customary and easier to use a mixed set, that is, principal planes and nodal points. In this manner, the lateral and angular magnifications are both unity. This is particularly useful when performing graphical ray tracing. In addition, the image location and magnification can be determined in the following manner:

- Trace a horizontal ray from the object tip to the first principal plane.
- Transfer to the second principal plane.
- Trace the ray through the second focal point.
- Trace a second ray from the object tip through the first focal point and to the intersection with the first principal plane.
- Transfer to the second principal plane.
- Project this ray horizontally until it intersects the first ray.

The intersection is the image height and locates the image plane. Magnification is the ratio of image height to object height. There are several alternative graphical ray-tracing method that can determine image height and location. One alternative is to trace a ray from the object tip to the first nodal point, and then project the ray exiting the second nodal point to intersect with the second ray mentioned above.

The focal length of a lens is related to the power of the lens by

$$\phi = \frac{n_o}{f_o} = \frac{n_i}{f_i} \quad (5)$$

This relationship is important in such optical systems as underwater cameras, cameras in space, and the like. For example, it is evident that the field of view is decreased for a camera in water.

The *lens law* can be expressed in several forms. If s_o and s_i are the distance from the object to the first principal point and the distance from the second principal point to the image, then the relationship between the object and the image is given by

$$\phi = \frac{n_i}{s_i} + \frac{n_o}{s_o} \quad (6)$$

Should the distance be measured with respect to the nodal points, the imaging equation becomes

$$\phi = \frac{n_o}{s_i} + \frac{n_i}{s_o} \quad (7)$$

When the distances are measured from the focal points, the image relationship, known as the *Newtonian imaging equation*, is given by

$$f_1 f_2 = s_o s_i \quad (8)$$

The power of a spherical refracting surface, with curvature c and n being the refractive index following the surface, is given by

$$\phi = c(n - n_o) \quad (9)$$

It can be shown that the power of a single thick lens in air is

$$\phi_{\text{thick}} = \phi_1 + \phi_2 - \phi_1 \phi_2 \frac{t}{n} \quad (10)$$

where t is the thickness of the lens. The distance from the first principal plane to the first surface is $-(t/n)\phi_2 f_1$ and the distance from the second principal point to the rear surface is $-(t/n)\phi_1 f_2$. The power of a thin lens ($t \rightarrow 0$) in air is given by

$$\phi_{\text{thin}} = (n-1)(c_1 - c_2) \quad (11)$$

17.4 STOPS AND PUPILS

The *aperture stop* or *stop* of a lens is the limiting aperture associated with the lens that determines how large an axial beam may pass through the lens. The stop is also called an *iris*. The *marginal ray* is the extreme ray from the axial point of the object through the edge of the stop. The *entrance pupil* is the image of the stop formed by all lenses preceding it when viewed from object space. The *exit pupil* is the image of the stop formed by all lenses following it when viewed from image space. These pupils and the stop are all images of one another. The *principal ray* is defined as the ray emanating from an off-axis object point that passes through the center of the stop. In the absence of pupil aberrations, the principal ray also passes through the center of the entrance and exit pupils.

As the obliquity angle of the principal ray increases, the defining apertures of the components comprising the lens may limit the passage of some of the rays in the entering beam thereby causing the stop not to be filled with rays. The failure of an off-axis beam to fill the aperture stop is called *vignetting*. The ray centered between the upper and lower rays defining the oblique beam is called the *chief ray*. When the object moves to large off-axis locations, the entrance pupil often has a highly distorted shape, may be tilted, and/or displaced longitudinally and transversely. Due to the vignetting and pupil aberrations, the chief and principal rays may become displaced from one another. In some cases, the principal ray is vignetted.

The *field stop* is an aperture that limits the passage of principal rays beyond a certain field angle. The image of the field stop when viewed from object space is called the *entrance window* and is called the *exit window* when viewed from image space. The field stop effectively controls the field of view of the lens system. Should the field stop be coincident with an image formed within or by the lens system, the entrance and exit windows will be located at the object and/or image(s).

A *telecentric stop* is an aperture located such that the entrance and/or exit pupils are located at infinity. This is accomplished by placing the aperture in the focal plane. Consider a stop placed at the front focal plane of a lens. The stop image or exit pupil is located at infinity and the principal ray exits the lens parallel to the optical axis. This feature is often used in metrology since the measurement error is reduced when compared to conventional lens systems because the centroid of the blur remains at the same height from the optical axis even as the focus is varied.

17.5 F-NUMBER AND NUMERICAL APERTURE

The focal ratio or F-number (FN) of a lens is defined as the effective focal length divided by the entrance pupil diameter D_{ep} . When the object is not located at infinity, the effective FN is given by

$$FN_{\text{eff}} = FN_{\infty}(1-m) \quad (12)$$

where m is the magnification. For example, for a simple positive lens being used at unity magnification ($m = -1$), the $FN_{\text{eff}} = 2FN_{\infty}$. The *numerical aperture* of a lens is defined as

$$NA = n_i \sin U_i \quad (13)$$

where n_i is the refractive index in which the image lies and U_i is the slope angle of the marginal ray exiting the lens. If the lens is aplanatic, then

$$FN_{\text{eff}} = \frac{1}{2NA} \quad (14)$$

The T-number of a lens is the effective FN divided by the square root of the transmittance of the lens and is used for radiometric computations; however, the FN should be used when computing depth of focus and depth of field discussed in Sec. 17.21.

17.6 MAGNIFIER OR EYE LOUPE

The typical magnifying glass, or *loupe*, comprises a singlet lens and is used to produce an erect but virtual magnified image of an object. The magnifying power of the loupe is stated to be the ratio of the angular size of the image when viewed through the magnifier to the angular size without the magnifier. By using the thin-lens model of the human eye, the magnifying power (MP) can be shown to be given by

$$MP = \frac{25 \text{ cm}}{d_e + d_o - \phi d_e d_o} \quad (15)$$

where d_o is the distance from the object to the loupe, d_e is the separation of the loupe from the eye, and $\phi = 1/f$ is the power of the magnifier. When d_o is set to the focal length of the lens, the virtual image is placed at infinity and the magnifying power reduces to

$$MP = \frac{25 \text{ cm}}{f} \quad (16)$$

Should the virtual image be located at the near viewing distance of the eye (about 25 cm), then

$$MP = \frac{25 \text{ cm}}{f} + 1 \quad (17)$$

Typically simple magnifiers are difficult to make with magnifying powers greater than about 10 \times .

17.7 COMPOUND MICROSCOPES

For magnifying power greater than that of a simple magnifier, a compound microscope, which comprises an objective lens and an eyepiece, may be used. The objective forms an aerial image of the object at a distance s_{ot} from the rear focal point of the objective. The distance s_{ot} is called the *optical tube length* and is typically 160 mm. The objective magnification is

$$MP_{\text{obj}} = \frac{s_{\text{ot}}}{f_{\text{obj}}} \quad (18)$$

The image formed is further magnified by the eyepiece which has a $MP_{\text{ep}} = 250 \text{ mm}/f_{\text{ep}}$. The total magnifying power of the compound microscope is given by

$$\begin{aligned} MP &= MP_{\text{obj}} \cdot MP_{\text{ep}} \\ &= \frac{160}{f_{\text{obj}}} \cdot \frac{250}{f_{\text{ep}}} \end{aligned} \quad (19)$$

Typically, $f_{\text{ep}} = 25 \text{ mm}$, so its MP 10. Should the objective have a focal length of 10 mm, the total magnifying power of the microscope is 16 \times times 10 \times , or 160 \times .

17.8 FIELD AND RELAY LENSES

Field lenses are placed at (or near) an image location for the purpose of optically relocating the pupil or to increase the field-of-view of the optical system. For example, a field lens may be used at the image plane of an astronomical telescope such that the field lens images the objective lens onto the eyepiece. In general, the field lens does not contribute to the aberrations of the system except for distortion and field curvature. Since the field lens must be positive, it adds inward curving Petzval. For systems having a small detector requiring an apparent increase in size, the field lens is a possible solution. The detector is located beyond the image plane such that it subtends the same angle as the objective lens when viewed from the image point. The field lens images the objective lens onto the detector.

Relay lenses are used to transfer an image from one location to another such as in a submarine periscope or borescope. It is also used as a means to erect an image in many types of telescopes and other such instruments. Often relay lenses are made using two lens groups spaced about a stop, or an image of the system stop, in order to take advantage of the principle of symmetry, thereby minimizing the comatic aberrations and lateral color. The relayed image is frequently magnified.

17.9 APLANATIC SURFACES AND IMMERSION LENSES

Abbe called a lens an aplanat that has an equivalent refractive surface which is a portion of a sphere with a radius r centered about the focal point. Such a lens satisfies the Abbe sine condition and implies that the lens is free of spherical and coma near the optical axis. Consequently, the maximum

possible numerical aperture (NA) of an aplanat is unity, or an FN = 0.5. In practice, an FN less than 0.6 is difficult to achieve. For an aplanat,

$$\text{FN} = \frac{1}{2 \cdot \text{NA}} \quad (20)$$

It can be shown that three cases exist where the spherical aberration is zero for a spherical surface. These are: (1) the trivial case where the object and image are located at the surface, (2) the object and image are located at the center of curvature of the surface, and (3) the object is located at the aplanatic point. The third case is of primary interest. If the refractive index preceding the surface is n_o and following the surface is n_i , then the object is located a distance s_o from the surface as expressed by

$$s_o = \frac{r(n_o + n_i)}{n_o} \quad (21)$$

and the image is located at

$$s_i = \frac{r(n_o + n_i)}{n_i} \quad (22)$$

An immersion lens or contact lens can be formed from an aplanatic surface and a plano surface. Figure 4 illustrates a hemispherical magnifier that employs the second aplanatic case. The resultant magnification is n_i if in air or n_i/n_o otherwise. A similar magnifier can be constructed by using a hyperhemispherical surface and a plano surface as depicted in Fig. 5. The lateral magnification is n_i^2 . This lens, called an *Amici lens*, is based upon the third aplanatic case. The image is free of all orders of spherical aberration, third-order coma, and third-order astigmatism. Axial color is also absent from the hemispherical magnifier. These magnifiers are often used as a means to make a detector appear larger and as the first component in microscope objectives.

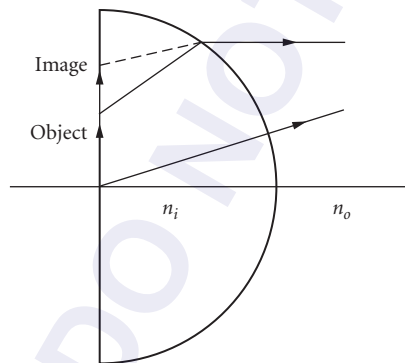


FIGURE 4 Aplanatic hemispherical magnifier with the object and image located at the center of curvature of the spherical surface. This type of magnifier has a magnification of n_i/n_o which can be used as a contact magnifier or as an immersion lens.

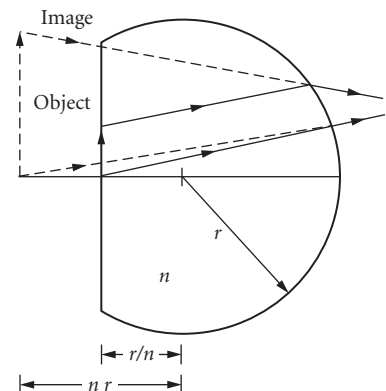


FIGURE 5 Aplanatic hyperhemispherical magnifier or Amici lens has the object located at the aplanatic point. The lateral magnification is $(n_i/n_o)^2$.

17.10 SINGLE ELEMENT LENS

It is well known that the spherical aberration of a lens is a function of its shape factor or bending. Although several definitions for the shape factor have been suggested, a useful formulation is

$$\mathcal{H} = \frac{c_1}{c_1 - c_2} \quad (23)$$

where c_1 and c_2 are the curvatures of the lens with the first surface facing the object. By adjusting the lens bending, the spherical aberration can be seen to have a minimum value.

The power of a thin lens or the reciprocal of its focal length is given by

$$\phi = \frac{(n-1)c_1}{\mathcal{H}} \quad (24)$$

When the object is located at infinity, the shape factor for minimum spherical aberration can be represented by

$$\mathcal{H} = \frac{n(2n+1)}{2(n+2)} \quad (25)$$

The resultant third-order spherical aberration of the marginal ray in angular units is

$$\text{SA3} = \frac{n^2 - (2n+1)\mathcal{H} + (1+2/n)\mathcal{H}^2}{16(n-1)^2(\text{FN})^3} \quad (26)$$

or after some algebraic manipulations,

$$\text{SA3} = \frac{n(4n-1)}{64(n+2)(n-1)^2(\text{FN})^3} \quad (27)$$

where, for a thin lens, the FN is the focal length f divided by the lens diameter, which in this case is the same as entrance pupil diameter D_{ep} . Inspection of this equation illustrates that smaller values of spherical aberration are obtained as the refractive index increases.

When the object is located at a finite distance s_o , the equations for the shape factor and residual spherical aberration are more complex. Recalling that the magnification m is the ratio of the object distance to the image distance and that the object distance is negative if the object lies to the left of the lens, the relationship between the object distance and the magnification is

$$\frac{1}{s_o \phi} = \frac{m}{1-m} \quad (28)$$

where m is negative if the object distance and the lens power have opposite signs. The term $1/s_o \phi$ represents the reduced or ϕ -normalized reciprocal object distance ν , that is, s_o is measured in units of focal length ϕ^{-1} . The shape factor for minimum spherical aberration is given by

$$\mathcal{H} = \frac{n(2n+1)}{2(n+2)} + \frac{2(n^2-1)}{n+2} \left(\frac{m}{1-m} \right) \quad (29)$$

and the resultant third-order spherical aberration of the marginal ray in angular units is

$$\begin{aligned} \text{SA3} = & \frac{1}{16(n-1)^2(\text{FN})^3} \left[n^2 - (2n+1)\mathcal{H} + \frac{n+2}{n}\mathcal{H}^2 + (3n+1)(n-1) \left(\frac{m}{1-m} \right) \right. \\ & \left. - \frac{4(n^2-1)}{n} \left(\frac{m}{1-m} \right) \mathcal{H} + \frac{(3n+2)(n-1)^2}{n} \left(\frac{m}{1-m} \right)^2 \right] \quad (30) \end{aligned}$$

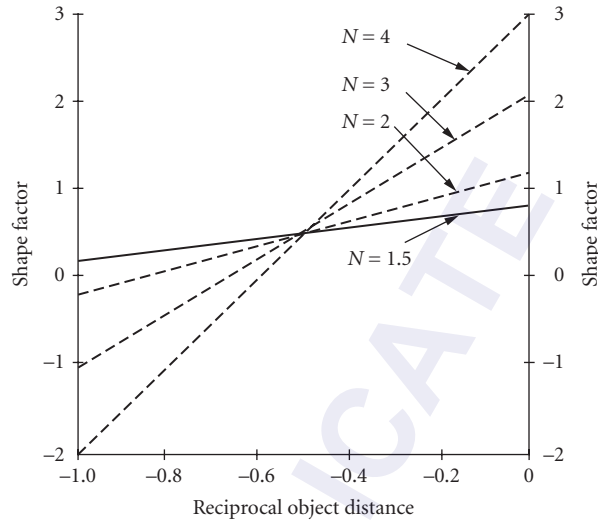


FIGURE 6 The shape factor for a single lens is shown for several refractive indexes as a function of reciprocal object distance v where the distance is measured in units of focal length.

where FN is the effective focal length of the lens f divided by its entrance pupil diameter. When the object is located at infinity, the magnification becomes zero and the above two equations reduce to those previously given.

Figure 6 illustrates the variation in shape factor as a function of v for refractive indices of 1.5 to 4 for an $FN = 1$. As can be seen from the figure, lenses have a shape factor of 0.5 regardless of the refractive index when the magnification is -1 or $v = -0.5$. For this shape factor, all lenses have biconvex surfaces with equal radii. When the object is at infinity and the refractive index is 4, lenses have a meniscus shape toward the image. For a lens with a refractive index of 1.5, the shape is somewhat biconvex, with the second surface having a radius about 6 times greater than the first surface radius.

Since the minimum-spherical lens shape is selected for a specific magnification, the spherical aberration will vary as the object-image conjugates are adjusted. For example, a lens having a refractive index of 1.5 and configured for $m = 0$ exhibits a substantial increase in spherical aberration when the lens is used at a magnification of -1 . Figure 7 illustrates the variation in the angular spherical aberration as both a function of refractive index and reciprocal object distance when the lens bending is for minimum spherical aberration with the object located at infinity. As can be observed from Fig. 7, the ratio of the spherical aberration, when $m = -0.5$ and $m = 0$, increases as n increases. Figure 8 shows the variation in angular spherical aberration when the lens bending is for minimum spherical aberration at a magnification of -1 . In a like manner, Fig. 9 presents the variation in angular spherical aberration for a convex-plano lens with the plano side facing the image. The figure can also be used when the lens is reversed by simply replacing the object distance with the image distance.

Figures 7 to 9 may provide useful guidance in setting up experiments when the three forms of lenses are available. The so-called “off-the-shelf” lenses that are readily available from a number of vendors often have the convex-plano, equal-radii biconvex, and minimum spherical shapes.

Figure 10 shows the relationship between the third-order spherical aberration and coma, and the shape factor for a thin lens with a refractive index of 1.5, stop in contact, and the object at infinity. The coma is near zero at the minimum spherical aberration shape. The shape of the lens as a function of shape factor is shown at the top of the figure.

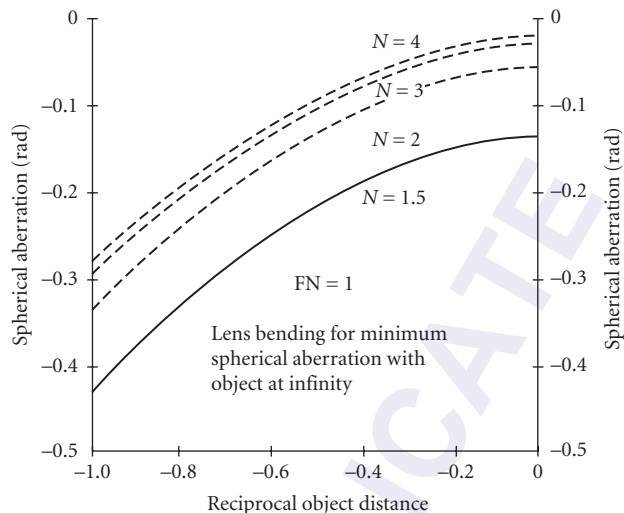


FIGURE 7 Variation of angular spherical aberration as a function of reciprocal object distance ν for various refractive indices when the lens is shaped for minimum spherical aberration with the object at infinity. Spherical aberration for a specific FN is determined by dividing the aberration value shown by $(FN)^3$.

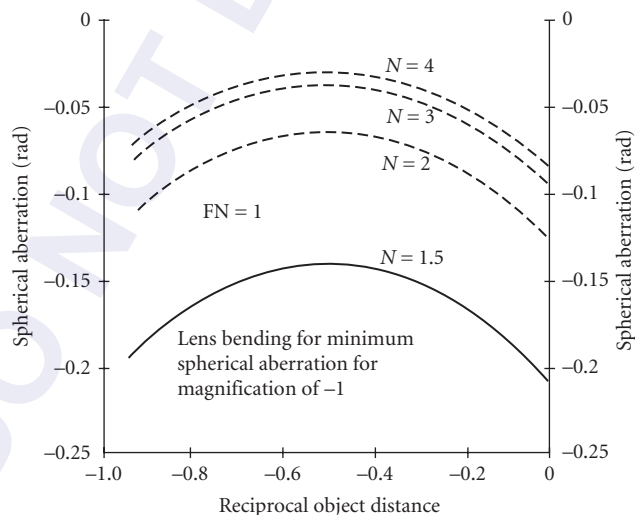


FIGURE 8 Variation of angular spherical aberration as a function of reciprocal object distance ν for various refractive indices when the lens is shaped for minimum spherical aberration for a magnification of -1 . Spherical aberration for a specific FN is determined by dividing the aberration value shown by $(FN)^3$.

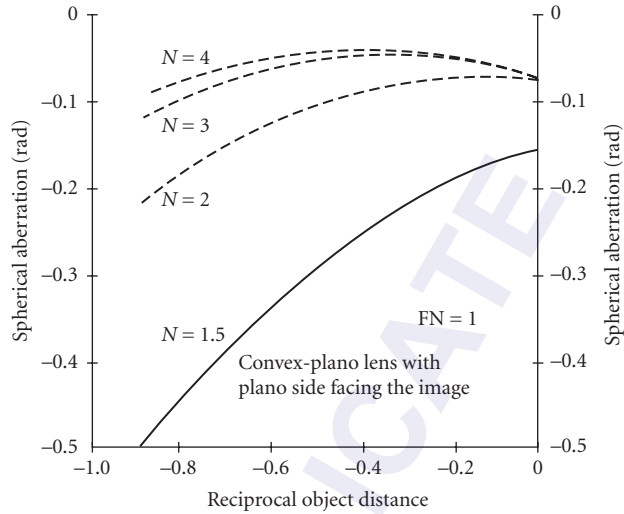


FIGURE 9 Variation of angular spherical aberration as a function of reciprocal object distance v for various refractive indices when the lens has a convex-plano shape with the plano side facing the object. Spherical aberration for a specific FN is determined by dividing the aberration value shown by $(FN)^3$.

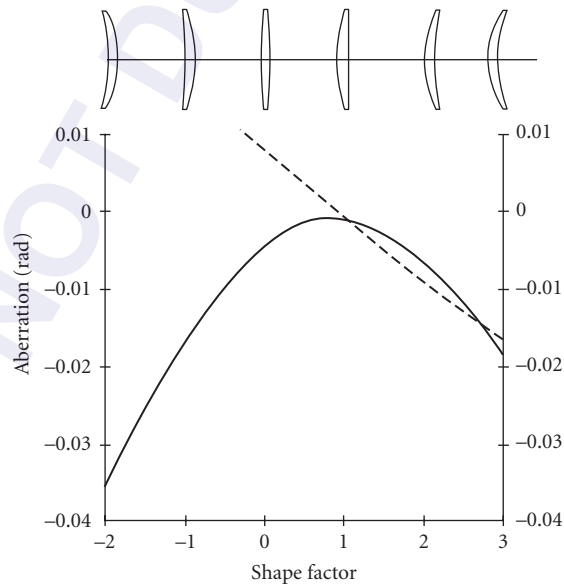


FIGURE 10 Variation of spherical aberration (solid curve) and coma (dashed line) as a function of shape factor for a thin lens with a refractive index of 1.5, stop in contact with the lens, and the object at infinity. The shape of the lens as the shape factor changes is shown at the top of the figure.

For certain cases, it is desirable to have a single lens with no spherical aberration. A useful form is the plano-convex, with the plano side facing the object, if the convex side is figured as a conic surface with a conic constant of $-n^2$. Caution should be exercised when using this lens form at other than infinite object distances; however, imaging at finite conjugates can be accomplished by using two lenses with their plano surfaces facing one another and the magnification being determined by the ratio of the focal lengths. It should be noted that for this lens form, the actual thickness of the lenses is not important and that the inclusion of the conic surface does not alter the focal length.

The off-axis performance of a lens shaped for minimum spherical aberration with the object at infinity can be estimated by using the following equations. Assuming that the stop is in contact with the lens, the third-order angular sagittal coma is given by

$$\text{CMA}_s = \frac{\theta}{16(n+2)(\text{FN})^2} \quad (31)$$

where the field angle θ is expressed in radians. The tangential coma is three times the sagittal coma or $\text{CMA}_t = 3 \cdot \text{CMA}_s$. The diameter of the angular astigmatic blur formed at best focus is expressed by

$$\text{AST} = \frac{\theta^2}{\text{FN}} \quad (32)$$

The best focus location lies midway between the sagittal and tangential foci. An estimate of the axial angular chromatic aberration is given by

$$\text{AChr} = \frac{1}{2V(\text{FN})} \quad (33)$$

where V is the Abbe number of the glass and $V = (n_2 - 1)/(n_3 - n_1)$ with $n_1 < n_2 < n_3$.

If a singlet is made with a conic or fourth-order surface, the spherical aberration is corrected by the aspheric surface, and the bending can be used to remove the coma. With the stop in contact with the lens, the residual astigmatism and chromatic errors remain as expressed by the preceding equations. Figure 11 depicts the shapes of such singlets for refractive indices of 1.5, 2, 3, and 4. Each lens has a unity focal length and an FN of 10. Table 1 presents the prescription of each lens where CC_2 is the conic constant of the second surface.

The *optical center* of a thick lens is located where a nodal ray crosses the optical axis of the lens. A *nodal ray* is aimed at the first nodal point, passes through the lens undeviated (although translated), and appears to emerge from the lens from the second nodal point. It can be shown that the distance from the first surface vertex to the optical center is $t/[1 - (c_1/c_2)]$ where t is the thickness of the lens.

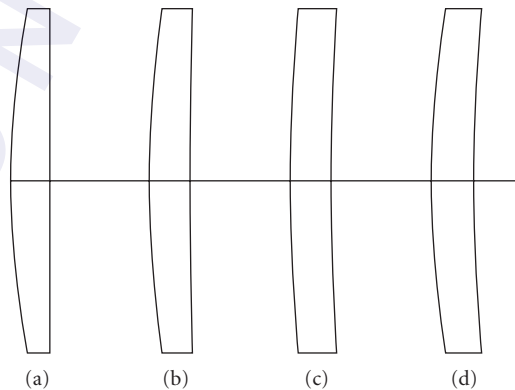


FIGURE 11 Variation of shape of singlets when the spherical aberration is corrected by the conic constant and the coma by the bending .

TABLE 1 Prescription of Singlets Corrected for Both Spherical Aberration and Coma

Lens	R_1	Thickness	R_2	Index	CC_2
<i>a</i>	0.55143	0.025	-5.27966	1.5	-673.543
<i>b</i>	0.74715	0.025	2.90553	2.0	23.2435
<i>c</i>	0.88729	0.025	1.56487	3.0	0.86904
<i>d</i>	0.93648	0.025	1.33421	4.0	0.24340

A remarkable property of the optical center is its wavelength independence (n does not appear in the preceding equation). This means that the spatial position of the optical center is fixed, where in contrast, the spatial positions of the six cardinal points are a function of wavelength because of their dependence upon n .

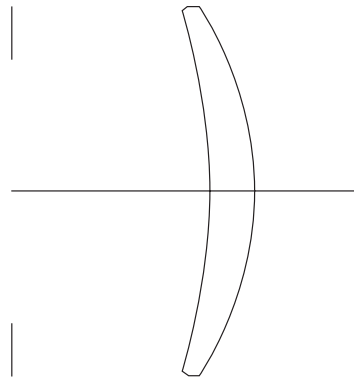
The optical center point (plane) is conjugate with the nodal points (planes); however, while the nodal points are related by unit angular magnification, the nodal-point to optical-center magnification (m_{OC}) is not necessarily unity. In general, m_{OC} is the ratio of the nodal ray slope angles at the first nodal point and the optical center. For a single thick lens, the magnification m_{OC} can be readily shown to be given by

$$m_{OC} = (r_1 - r_2) / [N(r_1 - r_2) - t(N - 1)]$$

All rotationally symmetric lenses have an optical center just as they possess the six cardinal points. Since the optical center is conjugate with N_1 and N_2 , the optical center can justifiably be considered also as a cardinal point. Should the aperture stop be located at the optical center, then the entrance pupil will be located at the first nodal point and the exit pupil will be located at the second nodal point with a unity pupil magnification. This statement is true whether the lens is of symmetrical or unsymmetrical design. When $n_o \neq n_i$, the exit pupil magnification will be n_o/n_i rather than unity.

17.11 LANDSCAPE LENSES AND THE INFLUENCE OF STOP POSITION

The first lens used for photography was designed in 1812 by the English scientist W. H. Wollaston about a quarter of a century before the invention of photography. He discovered that a meniscus lens with its concave surface toward the object could produce a much flatter image field than the simple biconvex lens commonly used at that time in the camera obscuras. This lens became known as the landscape lens and is illustrated in Fig. 12. Wollaston

**FIGURE 12** Landscape lens with the aperture stop located to the left of the lens.

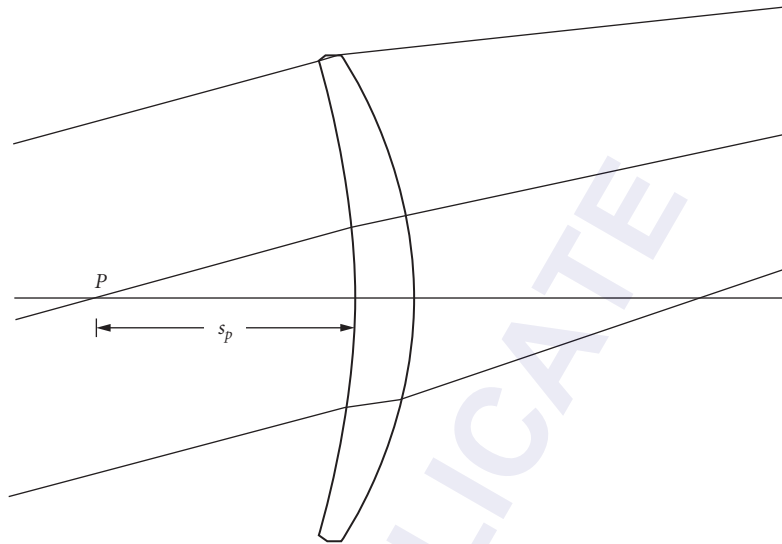


FIGURE 13 Rays traced at a given obliquity where the intersection of a given ray with the optical axis is P , located a distance s_p from the front surface of the lens.

realized that if the stop was placed an appropriate amount in front of the lens and the F-number was made to be modest, the image quality would be improved significantly over the biconvex lens.

The rationale for this can be readily seen by considering the influence on the residual aberrations of the lens by movement of the stop. Functionally, the stop allows certain rays in the oblique beam to pass through it while rejecting the rest. By simple inspection, it is clear that the movement of the stop (assuming a constant FN is maintained) will not affect the axial aberrations, while the oblique aberrations will be changed. In order to understand the influence of stop movement on the image quality, a graphical method was devised by R. Kingslake in which he traced a number of rays in the meridional plane at a given obliquity angle as illustrated in Fig. 13. A plot is generated that relates the intercept height of each real ray at the image plane H_i to the distance s_p from the intersection of the ray with optical axis P to the front surface of the lens. Each ray can be viewed as the principal ray when the stop is located at the intersection point P . This $H_i - s_p$ plot provides significant insight into the effect upon image quality incurred by placement of the stop. The shape of the curve provides information about the spherical aberration, coma, tangential field curvature, and distortion. Spherical aberration is indicated by an S-shaped curve, while the curvature at the principal ray point is a gauge of the coma. The coma is zero at inflection points. When the curve is a straight line, both coma and spherical aberration are essentially absent. The slope of the curve at the principal ray point is a measure of the tangential field curvature or the sag of the tangential field, that is, astigmatism. The difference in height of the real and gaussian principal rays in the image plane is distortion. For situations where the curve does not exhibit spherical aberration, it is impossible to correct the coma by shifting the stop.

Since a simple meniscus lens has stop position and lens bending as degrees of freedom, only two aberrations can be corrected. Typically, coma and tangential field curvature are chosen to be corrected, while axial aberrations are controlled by adjusting the FN of the lens. The $H_i - s_p$ plot for the lens shown in Fig. 13 is presented in Fig. 14, where the field angle is 10° and the image height is expressed as a percent of the gaussian image height. The lens has a unity focal length,

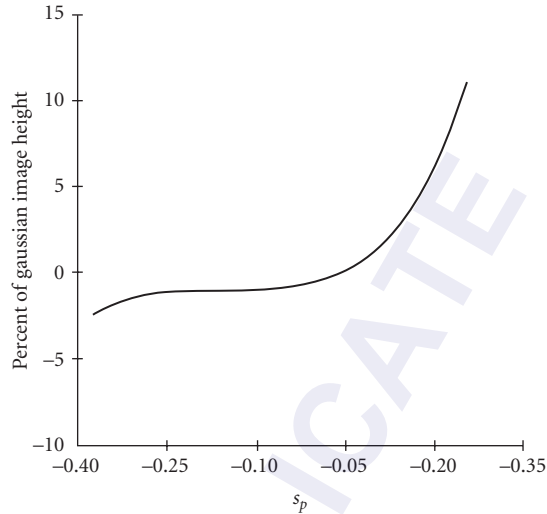


FIGURE 14 The image height H_i of each ray traced in Fig. 13 is plotted against the intersection length s_p to form the $H_i - s_p$ plot. H_i is expressed as a percent of the gaussian image height as a direct measure of distortion.

and the lens diameter is 0.275. Table 2 contains the prescription of the lens. Examination of this graph indicates that the best selection for stop location is when the stop is located at $s_p = -0.1505$ (left of the lens). For this selection, the coma and tangential astigmatism will be zero since the slope of the curve is zero and an inflection point is located at this stop position. Figure 15 shows the astigmatic field curves which clearly demonstrate the flat tangential image field for all field angles. Other aberrations cannot be controlled and must consequently be tolerated. When this lens is used at $F/11$, the angular blur diameter is less than $300 \mu\text{rad}$. It should be noted that this condition is generally valid for only the evaluated field-angle obliquity and will likely be different at other field angles. Nevertheless, the performance of this lens is often acceptable for many applications.

An alternate configuration can be used where the lens is in front of the stop. Such configuration is used to conserve space since the stop would be located between the lens and the image. The optical performance is typically poorer due to greater residual spherical aberration.

The principle demonstrated by the $H_i - s_p$ plot can be applied to lenses of any complexity as a means to locate the proper stop position. It should be noted that movement of the stop will not affect the coma if spherical aberration is absent nor will astigmatism be affected if both spherical aberration and coma have been eliminated.

TABLE 2 Prescription of Landscape Lens Shown in Fig. 13

Surface no.	Radius	Thickness	Index	Comment
1	Infinite	0.15050	1.0	Stop
2	-0.45759	0.03419	1.51680	BK7
3	-0.24887	0.99843	1.0	
4	Infinite			Image

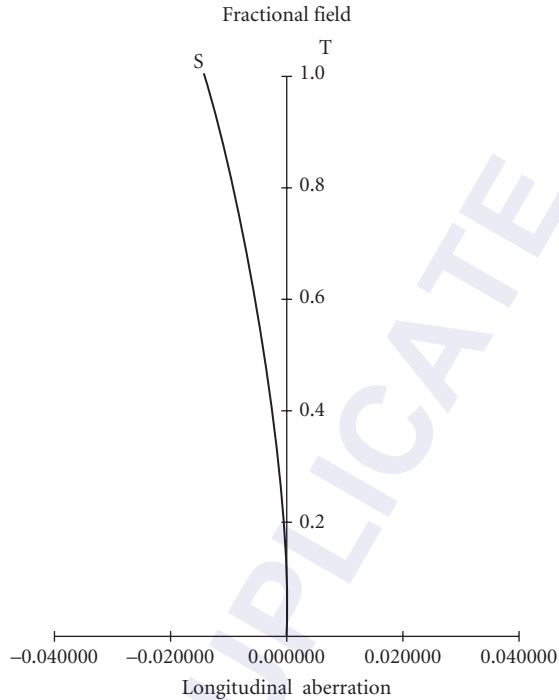


FIGURE 15 Astigmatic field curves for the landscape lens having the stop located at the zero slope location on the $H_i - s_p$ plot in Fig. 14, which is the flat tangential field position. S represents the sagittal astigmatic focus while T indicates the tangential astigmatic focus.

17.12 TWO-LENS SYSTEMS

Figure 16 illustrates the general imaging problem where an image is formed of an object by two lenses at a specified magnification and object-to-image distance. Most imaging problems can be solved by using two *equivalent* lens elements. An equivalent lens can comprise one lens or multiple lenses and may be represented by the principal planes and power of a single thick lens. All distances are measured from the principal points of each equivalent lens element. For simplicity, the lenses shown in Fig. 16 are thin lenses. If the magnification m , object-image distance s , and lens powers ϕ_a and ϕ_b are known, then the equations for s_1 , s_2 , and s_3 are given by

$$s_1 = \frac{\phi_b(s-s_2)-1+m}{m\phi_a + \phi_b}$$

$$s_2 = \frac{s}{2} \left[1 \pm \sqrt{1 - \frac{4[sm(\phi_a + \phi_b) + (m-1)^2]}{s^2 m \phi_a \phi_b}} \right] \quad (34)$$

$$s_3 = s - s_1 - s_2$$

The equation for s_2 indicates that zero, one, or two solutions may exist.

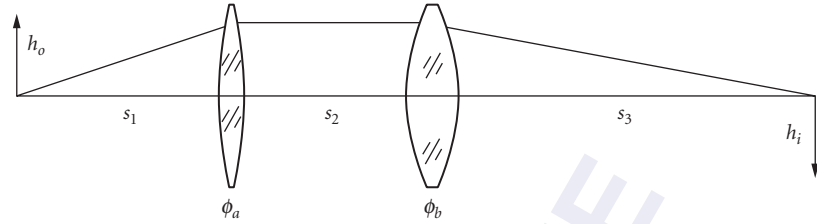


FIGURE 16 General imaging problem where the image is formed by two separated lenses.

If the magnification and the distances are known, then the lens powers can be determined by

$$\phi_a = \frac{s + (s_1 + s_2)(m - 1)}{ms_1s_2}$$

and

$$\phi_b = \frac{s + s_1(m - 1)}{s_2(s - s_1 - s_2)} \tag{35}$$

It can be shown that only certain pairs of lens powers can satisfy the magnification and separation requirements. Commonly, only the magnification and object-image distance are specified with the selection of the lens powers and locations to be determined. By utilizing the preceding equations, a plot of regions of all possible lens power pairs can be generated. Such a plot is shown as the shaded region in Fig. 17 where $s = 1$ and $m = -0.2$.

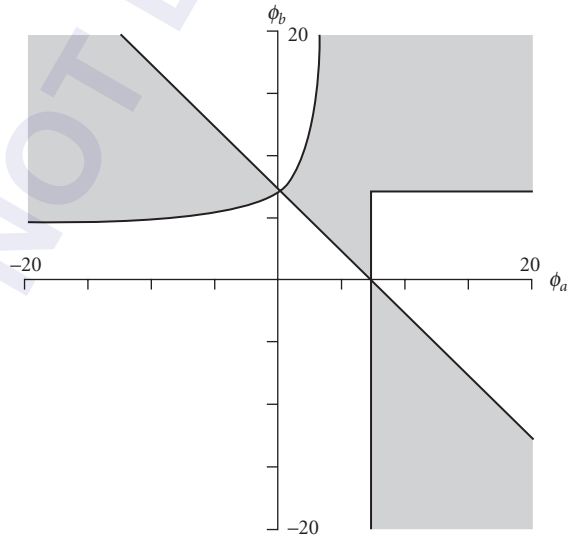


FIGURE 17 Shaded regions indicate all possible power pairs for the two lenses used for imaging. The solution space may be limited by physical considerations such as maximum aperture.

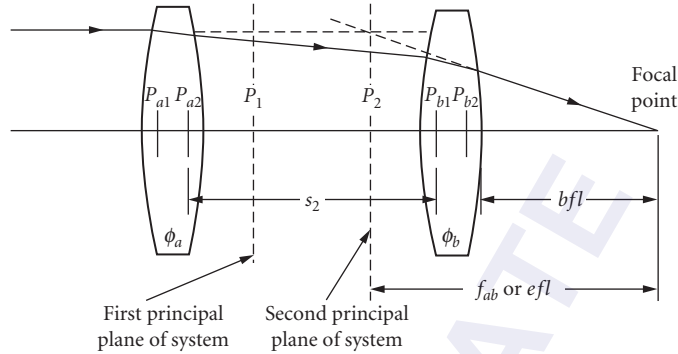


FIGURE 18 Combination of two thick lenses illustrating the principal points of each lens and the system, the f_{ab} or efl, and the bfl. Distances are measured from the principal points with the exception of the bfl.

Examination of this plot can assist in the selection of lenses that may likely produce better performance by, for example, selecting the minimum power lenses. The potential solution space may be limited by placing various physical constraints on the lens system. For example, the allowable lens diameters can dictate the maximum powers that are reasonable. Lines of maximum power can then be plotted to show the solution space.

When s_1 becomes very large compared to the effective focal length efl of the lens combination, the optical power of the combination of these lenses is expressed by

$$\phi_{ab} = \phi_a + \phi_b - s_2 \phi_a \phi_b \quad (36)$$

The effective focal length is ϕ_{ab}^{-1} or

$$f_{ab} = \frac{f_a f_b}{f_a + f_b - s_2} \quad (37)$$

and the back focal length is given by

$$\text{bfl} = f_{ab} \left(\frac{f_a - s_2}{f_a} \right) \quad (38)$$

The separation between lenses is expressed by

$$s_2 = f_a + f_b - \frac{f_a f_b}{f_{ab}} \quad (39)$$

Figure 18 illustrates the two-lens configuration when thick lenses are used. The principal points for the lens combination are denoted by P_1 and P_2 , P_{a1} and P_{a2} for lens a , and P_{b1} and P_{b2} for lens b . The distance between the principal points of a lens is called a *hiatus*. With the exception of the back focal length, all distances are measured from the principal points of each lens element or the combined lens system, as shown in the figure. For example, s_2 is the distance from P_{a2} to P_{b1} . The bfl is measured from the final surface vertex of the lens system to the focal point.

17.13 ACHROMATIC DOUBLET

The singlet lens suffers from axial chromatic aberration, which is determined by the Abbe number V of the lens material and its FN. A widely used lens form that corrects this aberration is the achromatic doublet as illustrated in Fig. 19. An achromatic lens has equal focal lengths in c and f light.

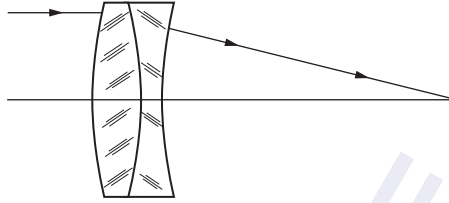


FIGURE 19 Typical achromatic doublet lens.

This lens comprises two lens elements where one element with a high V -number (crown glass) has the same power sign as the doublet and the other element has a low V -number (flint glass) with opposite power sign. Three basic configurations are used. These are the cemented doublet, broken contact doublet, and the widely airspaced doublet (dialyte). The degrees of freedom are two lens powers, glasses, and shape of each lens.

The resultant power of two thin lenses in close proximity, $s_2 \rightarrow 0$, is $\phi_{ab} = \phi_a + \phi_b$ and the transverse primary chromatic aberration TPAC is

$$\text{TPAC} = -yf_{ab} \left[\frac{\phi_a}{V_a} + \frac{\phi_b}{V_b} \right] \quad (40)$$

where y is the marginal ray height. Setting $\text{TPAC} = 0$ and solving for the powers of the lenses yields

$$\phi_a = \frac{V_a}{f_{ab}(V_a - V_b)} \quad (41)$$

and

$$\phi_b = \frac{-V_b \phi_a}{V_a} \quad (42)$$

The bending or shape of a lens is expressed by $c = c_1 - c_2$ and affects the aberrations of the lens. The bending of each lens is related to its power by $c_a = \phi_a / (n_a - 1)$ and $c_b = \phi_b / (n_b - 1)$. Since the two bendings can be used to correct the third-order spherical and coma, the equations for these aberrations can be combined to form a quadratic equation in terms of the curvature of the first surface c_1 . Solving for c_1 will yield zero, one, or two solutions for the first lens. A linear equation relates c_1 to c_2 of the second lens.

While maintaining the achromatic correction of a doublet, the spherical aberration as a function of its shape (c_1) is described by a parabolic curve. Depending upon the choices of glasses, the peak of the curve may be above, below, or at the zero spherical aberration value. When the peak lies in the positive spherical aberration region, two solutions with zero spherical aberration exist in which the solution with the smaller value of c_1 is called the left-hand solution (Fraunhofer or Steinheil forms) and the other is called the right-hand solution (Gaussian form). Two additional solutions are possible by reversal of the glasses. These two classes of designs are denoted as crown-in-front and flint-in-front designs. Depending upon the particular design requirements, one should examine all four configurations to select the most appropriate. The spherical aberration curve can be raised or lowered by the selection of the V difference or the n difference. Specifically, the curve will be lowered as the V difference is increased or if the n difference is reduced. As for the thin singlet lens, the coma will be zero for the configuration corresponding to the peak of the spherical aberration curve.

Although the primary chromatic aberration may be corrected, a residual chromatic error often remains and is called the secondary spectrum, which is the difference between the ray intercepts in d and c spectral lines. Figure 20a illustrates an F/5 airspaced doublet that exhibits well-corrected spherical light and primary chromatic aberrations and has notable secondary color. The angular secondary spectrum for an achromatic thin-lens doublet is given by

$$\text{SAC} = \frac{-(P_a - P_b)}{2(\text{FN})(V_a - V_b)} \quad (43)$$

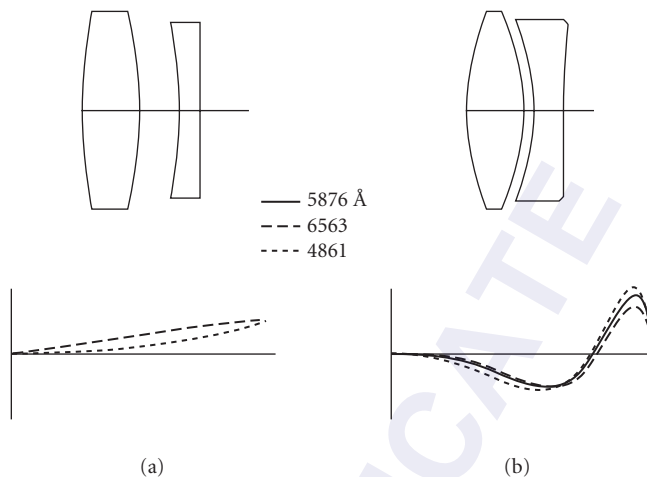


FIGURE 20 An F/5 airspaced doublet using conventional glasses is shown in (a) and exhibits residual secondary chromatic aberration. A similar lens is shown in (b) that uses a new glass to effectively eliminate the secondary color.

where $P = (n_d - n_c)(n_f - n_c)$ is the partial dispersion of a lens material. In general, the ratio $(P_a - P_b)/(V_a - V_b)$ is nearly a constant which means little can be done to correct the SAC. A few glasses exist that allow $P_a - P_b \approx 0$, but the $V_a - V_b$ is often small, which results in lens element powers of rather excessive strength in order to achieve achromatism. Figure 20b shows an F/5 airspaced doublet using a relatively new pair of glasses that have a small $P_a - P_b$ and a more typical $V_a - V_b$. Both the primary and secondary chromatic aberration are well corrected. Due to the relatively low refractive index of the crown glass, the higher power of the elements results in spherical aberration through the seventh order. Almost no spherochromatism (variation of spherical aberration with wavelength) is observed. The 80 percent blur diameter is almost the same for both lenses and is 0.007. Table 3 contains the prescriptions for these lenses.

TABLE 3 Prescriptions for Achromatic Doublets Shown in Fig. 20

Achromatic Doublet—1			
Surface No.	Radius	Thickness	Glass
1	49.331	6.000	BK7 517:642
2	-52.351	4.044	Air
3	-43.888	2.000	SF1 717:295
4	-141.706		Air
Achromatic Doublet—2			
Surface No.	Radius	Thickness	Glass
1	23.457	6.000	FK03 439:950
2	-24.822	1.059	Air
3	-22.516	3.000	BK7 517:642
4	94.310		Air

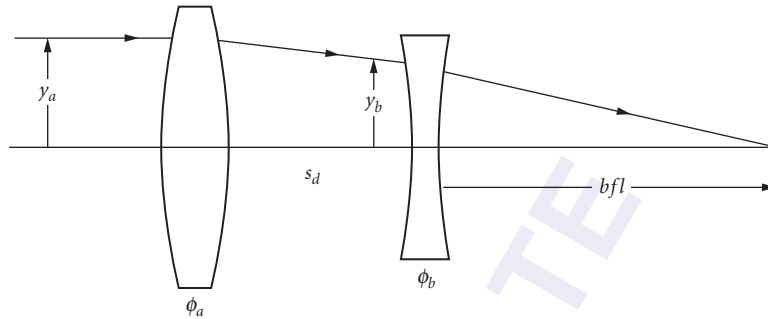


FIGURE 21 Widely separated achromatic doublet known as the dialyte lens.

When the separation between the lens elements is made a finite value, the resultant lens is known as a *dialyte* and is illustrated in Fig. 21. As the lenses are separated by a distance s_d , the power of the flint or negative lens increases rapidly. The distance s_d may be expressed as a fraction of the crown-lens focal length by $p = s_d/f_a$. Requiring the chromatic aberration to be zero implies that

$$\frac{y_a^2}{f_a V_a} + \frac{y_b^2}{f_b V_b} = 0 \quad (44)$$

By inspection of the figure and the definition of p , it is evident that $y_b = y_a(1 - p)$ from which it follows that

$$f_b V_b = -f_a V_a (1 - p)^2 \quad (45)$$

The total power of the dialyte is

$$\phi = \phi_a + \phi_b (1 - p) \quad (46)$$

Solving for the focal lengths of the lenses yields

$$f_a = f_{ab} \left[1 - \frac{V_b}{V_a (1 - p)} \right] \quad (47)$$

$$f_b = f_{ab} (1 - p) \left[1 - \frac{V_b (1 - p)}{V_b} \right] \quad (48)$$

The power of both lenses increases as p increases.

The typical dialyte lens suffers from residual secondary spectrum; however, it is possible to design an airspaced achromatic doublet with only one glass type that has significantly reduced secondary spectrum. Letting $V_a = V_b$ results in the former equations becoming

$$f_a = \frac{p f_{ab}}{p - 1} \quad f_b = -p f_{ab} (p - 1) \quad s_d = p f_a \quad bfl = -f_{ab} (p - 1) \quad (49)$$

When $f_{ab} > 0$, then p must be greater than unity, which means that the lens is quite long. The focal point lies between the two lenses, which reduces its general usefulness. This type of lens is known as the Schupmann lens, based upon his research in the late 1890s. Several significant telescopes, as well as eyepieces, have employed this configuration. For $f_{ab} < 0$, the lens can be made rather compact and is sometimes used as the rear component of some telephoto lenses.

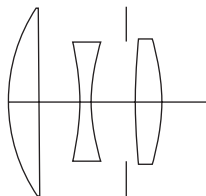


FIGURE 22 Typical triplet lens.

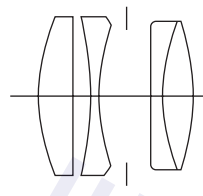


FIGURE 23 Typical Tessar lens.

17.14 TRIPLET LENSES

In 1893, a new type of triplet lens for photographic applications was invented by the English designer H. Dennis Taylor. He realized that the power of two lenses in contact of equal, but opposite, power is zero, as is its Petzval sum. As the lenses are separated, the system power becomes positive since the negative lens contributes less power. The Petzval sum remains zero, since it does not depend upon the marginal ray height. In order to overcome the large aberrations of such a configuration, Taylor split the positive lens into two positive lenses and placed one on each side of the negative lens. A stop is often located between the negative and rear-positive lenses. Figure 22 illustrates a typical triplet lens. The triplet can be used at reasonably large apertures ($>F/4$) and moderately large fields of view ($>\pm 25^\circ$).

The triplet has eight degrees of freedom which are the three powers, two airspaces, and three lens bendings. The lens powers and airspaces are used to control the axial and lateral chromatic aberrations, the Petzval sum, the focal length, and the ratio of the airspaces. Spherical aberration, coma, and astigmatism are corrected by the lens bendings. Distortion is usually controlled by the airspace ratio or the choice of glasses. Consequently, the triplet has exactly the number of degrees of freedom to allow correction of the basic aberrations and maintain the focal length.

The design of a triplet is somewhat difficult since a change of any surface affects every aberration. The choice of glass is important and impacts the relative aperture, field of view, and overall length. For example, a large ΔV produces a long system. It should be noted that a triplet corrected for third-order aberrations by using the degrees of freedom almost always leads to a lens with poor performance. A designer normally leaves a certain amount of residual third-order aberrations to balance the higher-order terms. The process for thin-lens predesign is beyond the scope of this *Handbook*; however, it may be found in various references comprising the bibliography.

A few years later, Paul Rudolph of Zeiss developed the Tessar, which resembles the triplet, with the rear lens replaced by an achromatic doublet. The Tessar shown in Fig. 23 was an evolution of Rudolph's anastigmats which were achromatic lenses located about a central stop. The advantage of the achromatic rear component is that it allows reduction of the zonal spherical aberration and the oblique spherical aberration, and reduces the separation of the astigmatic foci at other than the design maximum field angle. Performance of the Tessar is quite good and has generally larger relative apertures at equivalent field angles than the triplet. A variety of lenses were derived from the triplet and the Tessar in which the component lenses were made into doublets or cemented triplets.

17.15 SYMMETRICAL LENSES

In the early 1840s, it was recognized that lenses that exhibit symmetry afford various benefits to the lens designer. The first aberration acknowledged to be corrected by the symmetry principle was distortion. It can also be shown that coma and lateral color are necessarily corrected by a symmetrical lens construction. Although the principle of symmetry implies that the lens be

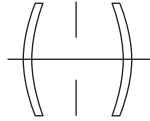


FIGURE 24 The periscopic lens illustrates the earliest form of symmetrical lenses. It is formed by placing two landscape lenses about a central stop. Symmetry removes the aberrations of coma, distortion, and lateral color.

operated at a magnification of -1 , the degree to which the aberrations are upset by utilizing the lens at other conjugates is remarkably small. This principle forms the basis of most wide-field-of-view lenses.

One of the earliest symmetrical lenses was the Periscopic (Periskop) lens invented by C. A. Steinheil in 1865. Figure 24 shows an F/11 Periscopic lens constructed from the landscape lens discussed previously. Symmetry corrects for coma and distortion, while the spacing of the lenses and their shapes are selected to produce a flat tangential astigmatic field. Since the stop position for the landscape lens was chosen to yield a flat tangential astigmatic field, essentially no change in the lens separation is necessary even though the Periscopic lens is being used at infinite conjugates. No correction for spherical aberration can be made. When used at other than unit magnification, some optical improvement can be achieved by making the stop slightly asymmetrical and/or having a different shape for the front or rear lens. This lens has continued to find application throughout this century.

By 1866, Dallmeyer in England and Steinheil and von Seidel in Germany both invented the Rapid Rectilinear lens that could be used at apertures of up to F/6. The lens has two cemented achromats about a central stop. Use of the doublet allows correction of the axial chromatic and spherical aberrations. Glass selection is of importance in the design. Typically, the Δn between the glasses should be large while the ΔV should be relatively small. The positive lens is located nearest the stop and has the lower refractive index. A notable characteristic of the lens is that the aberrations are reasonably stable over a broad range of object distances.

It should be noted that vignetting is often used in these and other lens types to control the higher-order aberrations that are often observed at large field angles. Although a loss in illumination occurs, the gain in resolution is often worthwhile.

The airspace dialyte lens comprises four lenses symmetrically arranged about a central stop. The rear portion of the lens is an achromatic doublet that has five degrees of freedom (an air space, two powers, and two bendings) which may be used to control the focal length, spherical aberration, axial chromatic aberration, astigmatism, and the Petzval sum. With a like pair of lenses mounted in front of the stop, the symmetry corrects the coma, distortion, and lateral color. When used at infinite conjugates, the resultant residuals of the aberrations can be controlled by deviating somewhat from perfect symmetry of the air spaces about the stop. Lenses of this type can provide useful performance with apertures approaching F/4 and fields of view of about $\pm 20^\circ$ or so.

17.16 DOUBLE-GAUSS LENSES

In the early 1800s, Gauss described a telescope objective comprising a pair of meniscus lenses with one having positive power and the other negative power. An interesting aspect of his lens is that the spherochromatism is essentially constant. Although this lens found little acceptance, in 1888, Alvan Clark of Massachusetts placed a pair of the Gauss lenses around a central stop to create a high-aperture, wide-field-of-view lens. This lens form is known as the Double-Gauss lens and is the basis of almost every

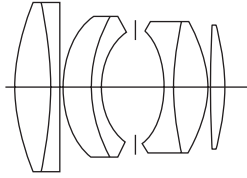


FIGURE 25 Unsymmetrical Double-Gauss or Biotar lens introduced as the Leica Summitar in 1939.

high-aperture lens developed to date. An example of this lens was patented by Richter in 1933 and can cover a field of view of $\pm 45^\circ$ at F/6.

In 1896, Paul Rudolph of Zeiss developed the Planar which reduces the often serious oblique spherical aberration and the separation of the astigmatic foci at intermediate field angles. Rudolph placed a buried surface into the thick negative elements to control the chromatic aberration. A buried surface is defined as the interface between two glasses that have the same refractive index n_d at the central wavelength, but have significantly different Abbe numbers. Such a surface has no effect upon the monochromatic aberrations or the lens system power, but does allow the inclusion of a wide range of chromatic aberration to compensate for that caused by the rest of the lens.

Many Double-Gauss lenses are symmetrical; however, it was discovered that if the lens was made unsymmetrical, then an improvement in performance could be realized. This lens form is often called the Biotar. A large portion of 35-mm camera lenses are based upon this design form or some modification thereof. Figure 25 shows the configuration of the Leica Summitar introduced in 1939.

It is the general nature of meniscus lens systems of this type to exhibit little coma, distortion, or lateral color; however, oblique spherical aberration is often observed to increase to significant levels as the field angle increases. Oblique spherical aberration can be recognized in transverse ray plots as the S shape of spherical aberration, but with the S becoming increasingly stronger as the field angle increases. As the aperture is increased beyond about F/8, the outer negative elements must be thickened dramatically and achromatic surfaces must necessarily be included.

17.17 PETZVAL LENSES

In 1839, Petzval designed a new type of lens that comprises a front objective with an achromatic, airspaced doublet as the rear elements. The Petzval lens has found great application in projectors and as a portrait lens. Both spherical aberration and coma can be well-corrected, but the lens configuration causes the Petzval sum to be undercorrected, which results in the field of view being limited by the astigmatism. The Petzval-field curves inward and may be corrected by including a *field flattener lens* in close proximity to the image plane. A typical example of a Petzval lens is shown in Fig. 26.

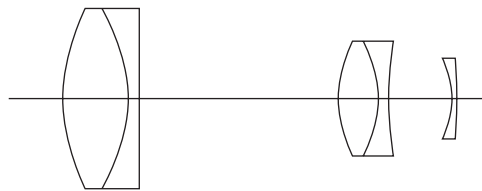


FIGURE 26 Typical Petzval lens.

17.18 TELEPHOTO LENSES

A telephoto lens provides an effective focal length efl that is longer than its overall length s_{ol} as measured from the front of the lens to the image plane. The telephoto ratio is defined as s_{ol}/efl , thus a lens with a ratio less than one is a telephoto lens. The basic concept of a telephoto lens is illustrated by the dialyte lens configuration in which a negative lens is inserted between the objective lens and the image plane. This concept goes back to Kepler, but Peter Barlow developed the idea in the early 1800s by including a negative achromat in telescopes to increase their magnification. Barlow type lenses are widely used today. As the telephoto ratio is made smaller, the design of the lens becomes more difficult, primarily due to the Petzval sum increasing.

When most telephoto lenses are used to view objects that are relatively close, the image quality degrades rapidly due to the typical unsymmetrical lens configuration. Some modern telephoto lenses include one or more elements that move as the lens is focused for the purpose of aberration correction.

17.19 INVERTED OR REVERSE TELEPHOTO LENSES

A reverse telephoto lens has a telephoto ratio greater than unity and exhibits a shorter focal length than its overall length, a larger bfl than is provided by normal lenses of the same efl , lenses with generally large apertures and wide fields of view, and lens elements of physically larger size that allow easier manufacture and handling. The basic configuration has a large negative lens located in front of a positive objective lens. Since the negative lens makes the object appear closer to the objective lens, the resultant image moves beyond the focal point, thereby making the bfl greater than the efl .

An extreme form of the reverse telephoto lens is the fish-eye or sky lens. Such lenses have a total field of view of 180° or more. The image formed by these lenses has very large barrel distortion. Recalling that the image height for a distortionless lens on a flat image surface is $f \tan \theta$, the reverse telephoto lens has mapping relationships such as $f\theta$ and $f \sin \theta$. When the barrel distortion is given by $f \sin \theta$, the illumination across the image will be constant if such effects as vignetting and stop/pupil distortion are absent. Barrel distortion has the effect of compressing the outer portions of the image toward the central portion, thereby increasing the flux density appropriately.

After World War II, the Russian designer M. M. Roosinov patented a double-ended reverse-telephoto lens that was nearly symmetrical with large negative lenses surrounding a pair of positive lenses with a central stop. Although the back focal length is quite short, it provides relatively large aperture with a wide field of view and essentially no distortion. Lenses of this type have found significant use in aerial photography and photogrammetry.

17.20 PERFORMANCE OF REPRESENTATIVE LENSES

Figures 27 to 38 present the performance of lenses, selected generally from the patent literature, representing a variety of lens types. The measures of performance provided in each figure have been selected for utilization purposes. Diffraction effects have not been included.

Each figure is divided into four sections *a* to *d*. Section *a* is a drawing of the lens showing the aperture stop. Section *b* contains two sets of plots. The *solid* line is for the distortion versus field of view (θ) in degrees while the *dashed* lines show the transmission of the lens versus field of view for three F-numbers. Transmission in this case is *one minus the fractional vignetting*. No loss for coatings, surface reflection, absorption, and the like is included. The rms diameter of the geometric point source image versus field-of-view for three F-numbers is presented in section *c*. The spot sizes are in angular units and were calculated for the central wavelength only, that is, monochromatic values. Note that the ordinate is logarithmic. The final section, *d*, contains angular transverse ray plots in all three colors for both the on-axis and near-extreme field angles with y_{ep} being measured in the

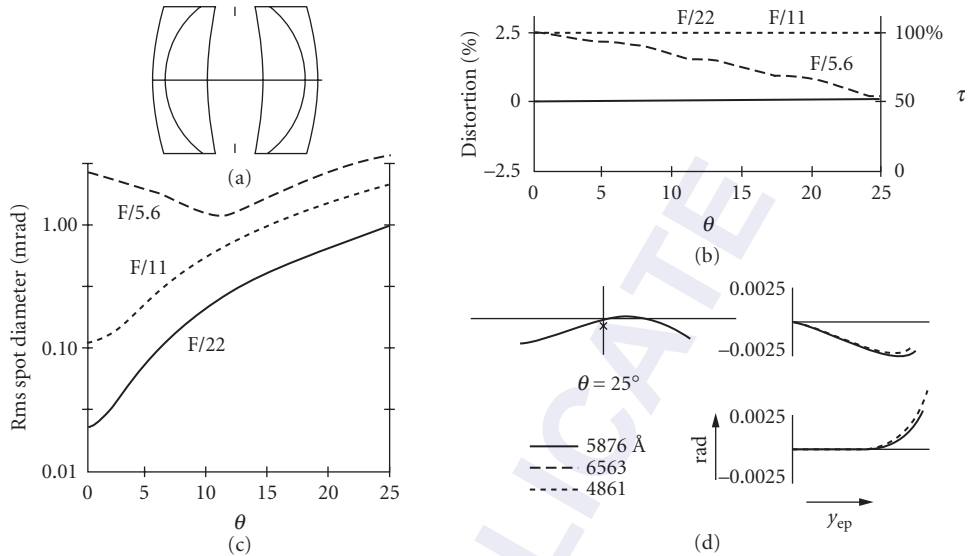


FIGURE 27 Rapid Rectilinear: This lens is an aplanat which is symmetrical with the rear half corrected for spherical aberration and flat tangential field. A compact configuration is realized by having a large amount of coma in each half. Symmetry removes the lens system coma, distortion, and lateral color. This type of lens is one of the most popular camera lenses ever made.

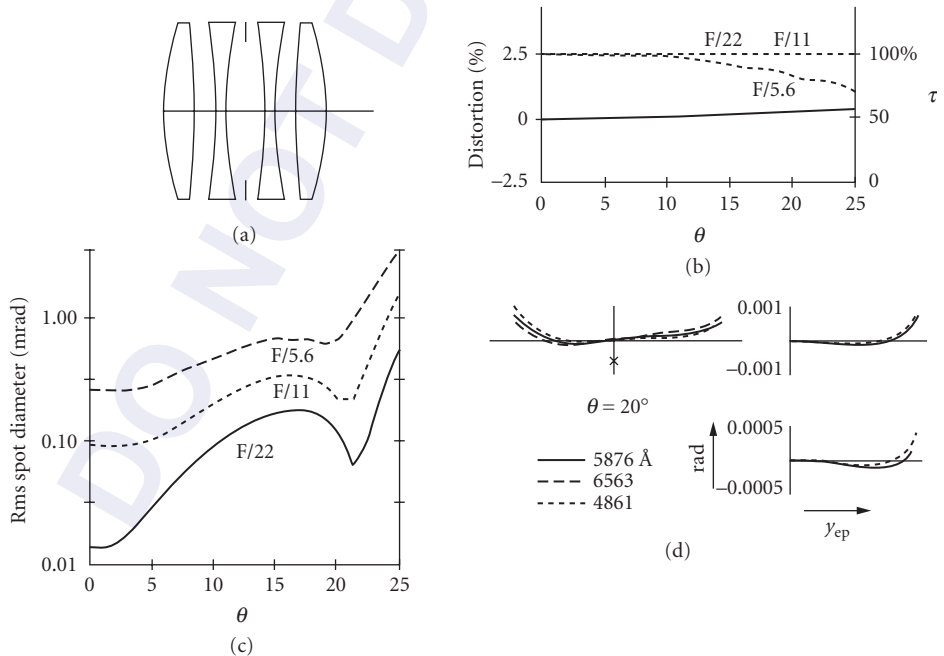


FIGURE 28 Celor : F/5.6 with 50° total field of view. Also known as an airspaced dialyte lens. (After R. Kingslake, *Lens Design Fundamentals*, Academic Press, New York, 1978, p. 243.)

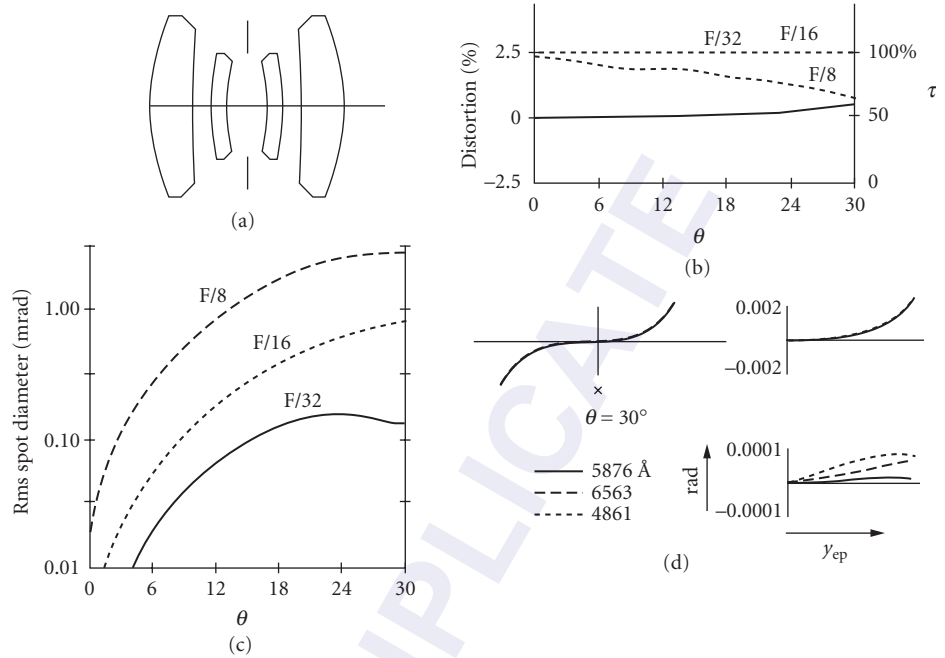


FIGURE 29 Symmetrical double anastigmat or gauss homocentric objective: basic form of Double-Gauss lens using a pair of gauss telescope objectives. First patented by Alvan Clark in 1888, USP 399,499. (After R. Kingslake, *Lens Design Fundamentals*, Academic Press New York, 1978, pp. 224–250.)

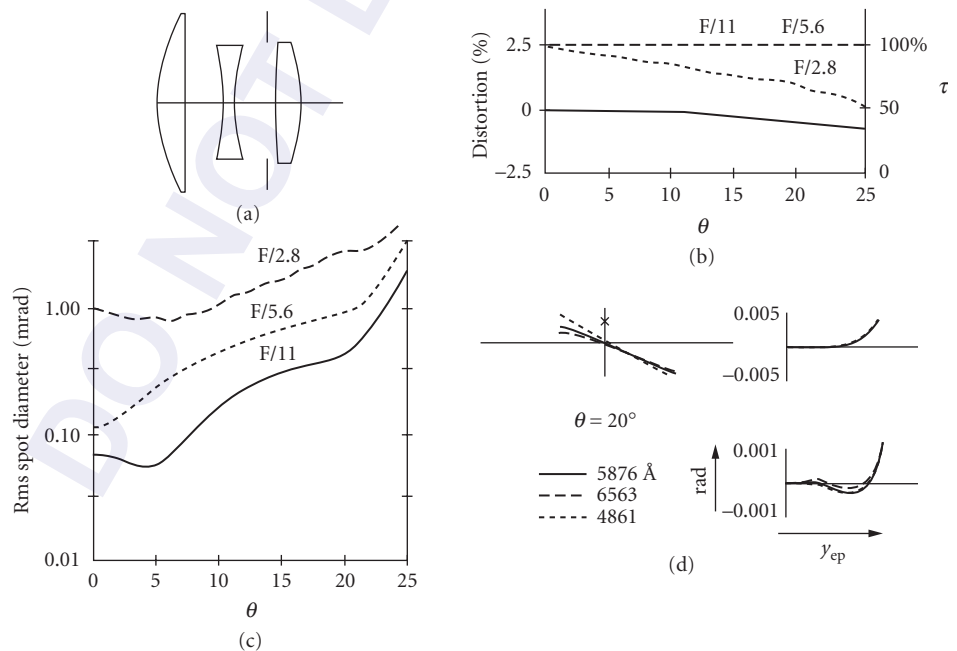


FIGURE 30 Triplet: F/2.8 with 50° total field of view. (Tronnier, USP 3,176,582.)

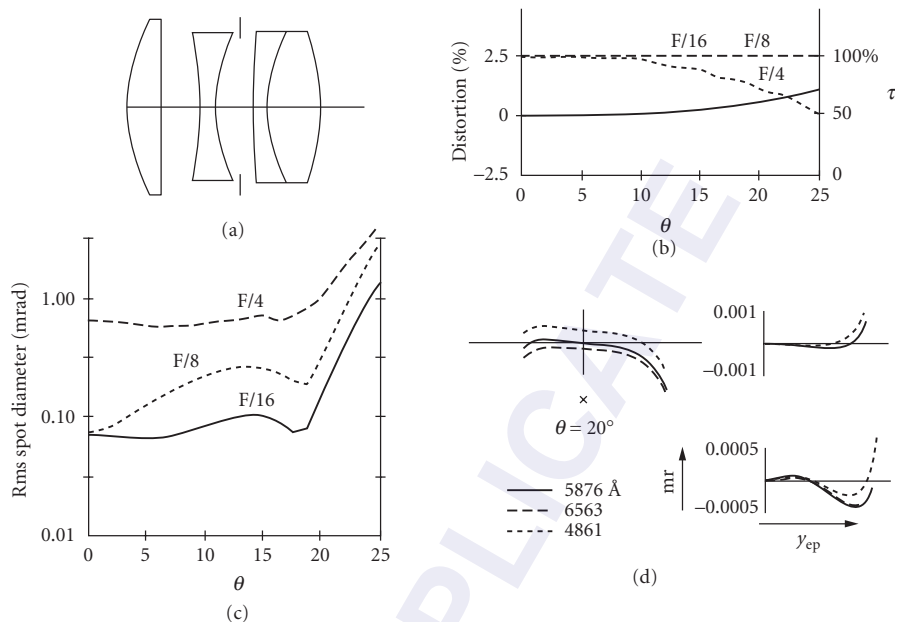


FIGURE 31 Tessar: F/4 with 50° total field of view. (Tronnier, USP 2,084,714, 1937.)

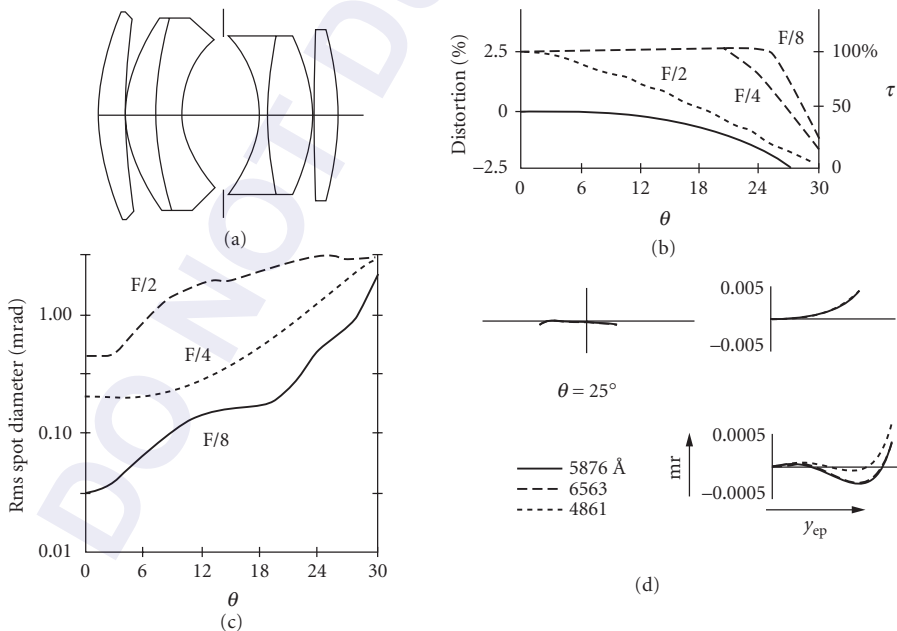


FIGURE 32 Unsymmetrical Double-Gauss: This lens was designed in 1933 for Leitz and was called the Summar. F/2 with 60° total field of view. This lens was replaced by the Leitz Summar in 1939, due to rapidly degrading off-axis resolution and vignetting. Compare this lens with the lens shown in Fig. 33. (Tronnier, USP 2,673,491.)

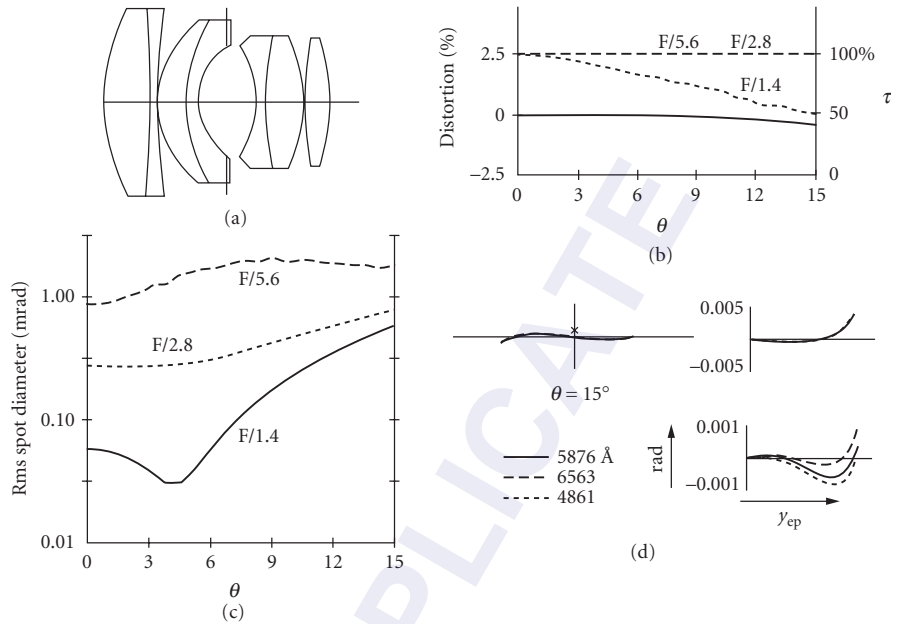


FIGURE 33 Unsymmetrical Double-Gauss: This lens type was designed in 1939 for Leitz and was called the F/2 Summarit. Kodak had a similar lens called the F/1.9 Ektar. A later example of this design form is shown and operates at F/1.4 with 30° total field of view. (*Klomp, USP 3,005,379.*)

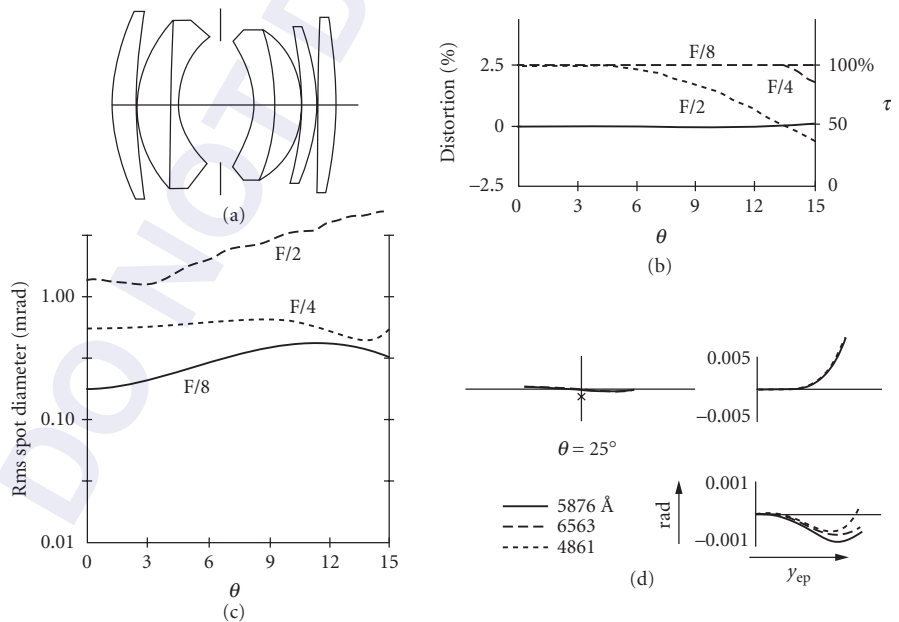


FIGURE 34 Unsymmetrical Double-Gauss: F/1.75 with 50° total field of view. Similar to the 1949 Leitz F/1.5 Summarit. This lens has a split rear element which produces improved resolution of the field of view and less vignetting than the earlier Summar type lens. (*Cook, USP 2,959,102.*)

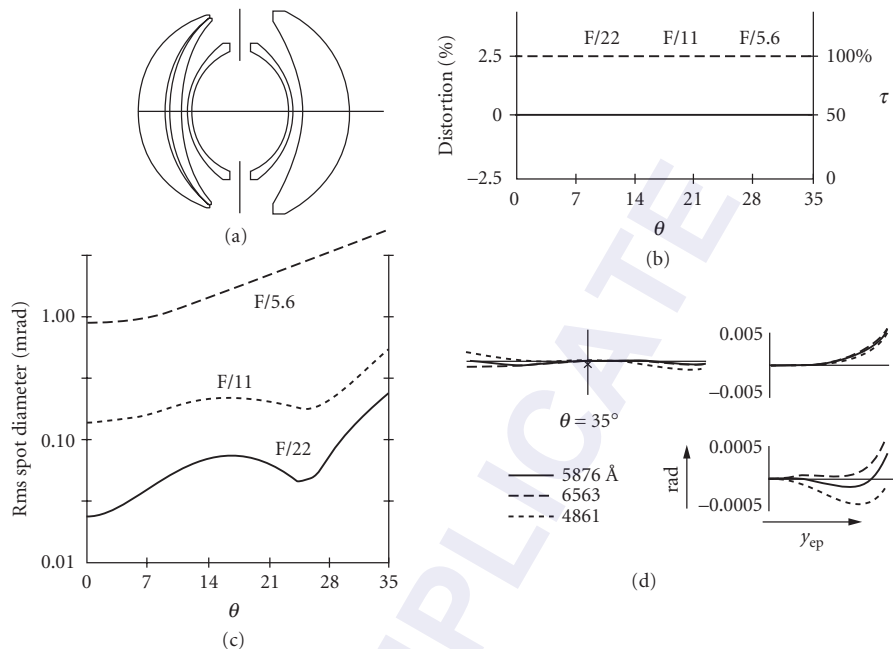


FIGURE 35 Unsymmetrical Double-Gauss: F/5.6 with 70° field of view. This lens is a variant of the 1933 Zeiss F/6.3 Topogon (USP 2,031,792) and is the Bausch & Lomb Metrogon. The principal difference is the splitting of the front element. (Rayton, USP 2,325,275.)

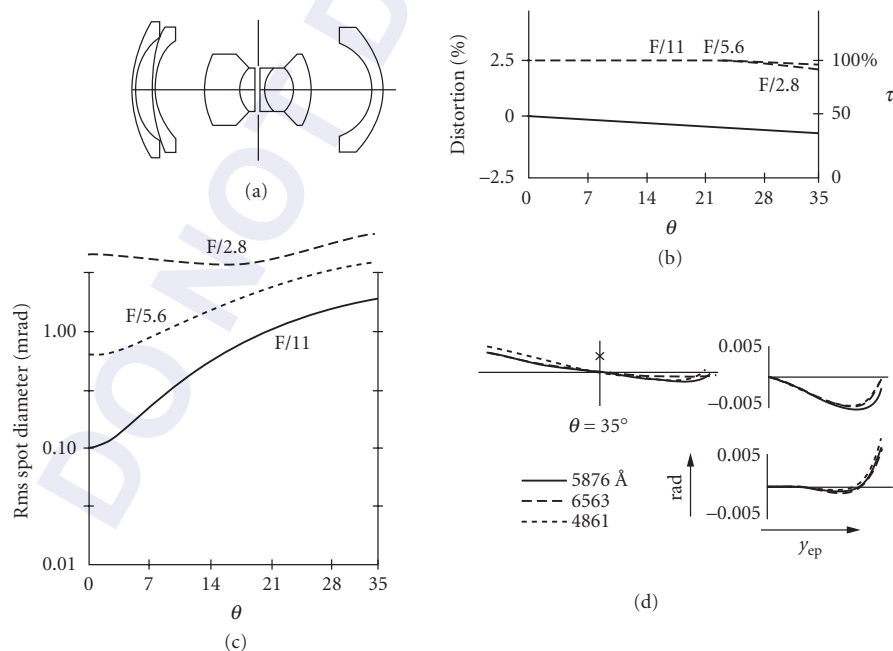


FIGURE 36 Reverse Telephoto: This lens was developed by Zeiss in 1951 and is known as the Biogon. It operates at F/2.8 with 70° field of view. This lens comprises two reverse-telephoto objectives about a central stop. (Bertele, USP 2,721,499.)

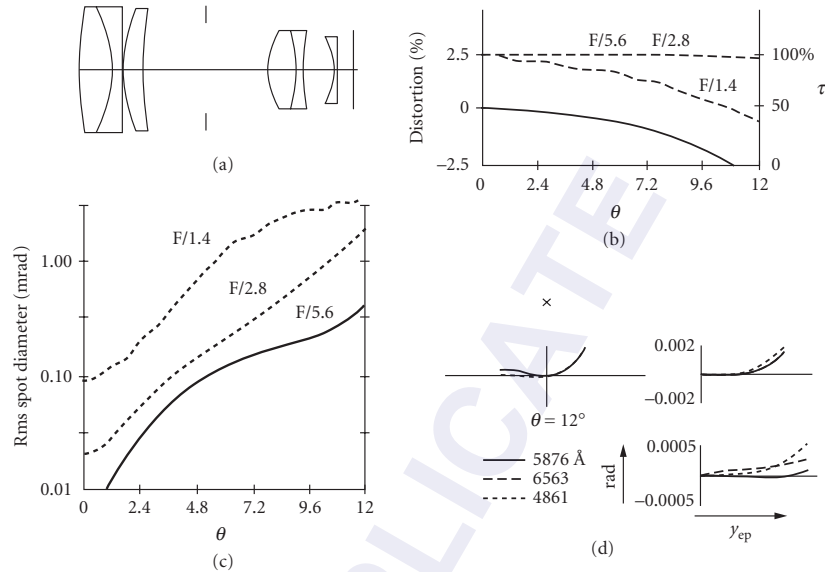


FIGURE 37 Petzval: Example of Kodak projector lens operating at F/1.4 with 24° total field of view. The front lens group has its power shared between a cemented doublet and a singlet for aberration correction. Note that the aperture stop is located between the front and rear groups rather than the more common location at the front group. Resolution in the region near the optical axis is very good although it falls off roughly exponentially. The limiting aberrations are oblique spherical and cubic coma. (Schade, USP 2,541,484.)

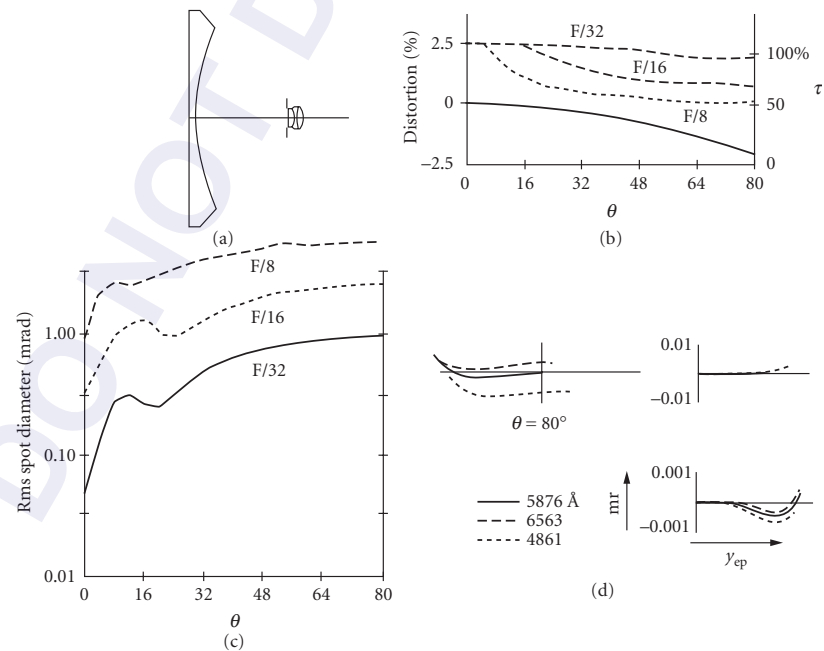


FIGURE 38 Fish-eye: The Hill Sky lens was manufactured by Beck of London in 1924. The lens has moderate resolution and enormous distortion characteristic of this type of lens. (Merte, USP 2,126,126.)

entrance pupil. The lower right plot shows the axial aberrations while the upper left plot represents the tangential/meridional aberrations and the upper right plot presents the sagittal aberrations. The X included on some of the tangential plots represents the location of the paraxial principal ray which also provides a measure of the distortion. The legend indicating the relationship between line type and wavelength is included.

The linear spot size is computed by multiplying the efl by the angular spot size. This value can be compared against the diffraction-limited spot size given by $2.44(\lambda/D_{ep})$. If the geometric spot is several times *smaller* than the diffraction-limited spot, then the lens may be considered to be diffraction-limited for most purposes. If the geometric spot is several times *larger*, then the lens performance is controlled by the geometric spot size for most applications.

17.21 RAPID ESTIMATION OF LENS PERFORMANCE

Singlet

Figure 39 is a nomogram that allows quick estimation of the performance of a single refracting lens, with the stop at the lens, as a function of refractive index N , dispersion V , F-number, and field of view θ . Chart A estimates the angular blur diameter β resulting from a singlet with bending for minimum spherical aberration. The angular chromatic blur diameter is given by Chart B. The three rows of FN values below the chart represent the angular blur diameter that contains the indicated

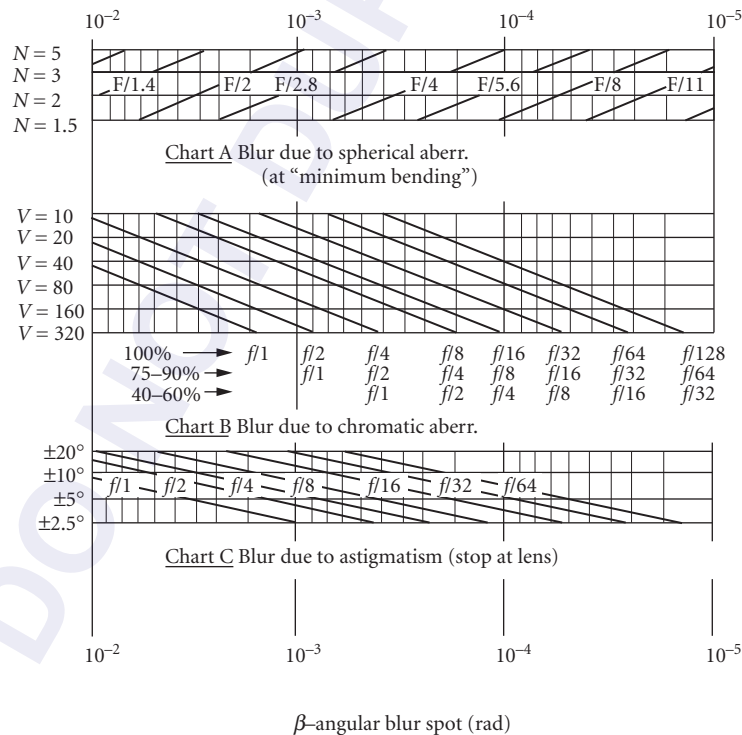


FIGURE 39 Estimation of single lens spot size as a function of refractive index, dispersion, F-number, and field of view. (Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1990, p. 458.)

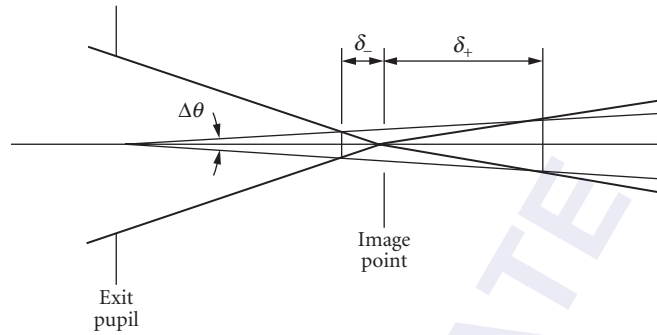


FIGURE 40 Geometric relationships for determining the geometric depth of focus of a lens .

percentage of the total energy. Chart C shows the blur diameter due to astigmatism. Coma for a singlet bent for minimum spherical aberration with the stop at the lens is approximately

$$\frac{\theta}{16 \cdot (N+2) \cdot (\text{FN})^2} \quad (50)$$

Depth of Focus

The *depth of focus* of an optical system is expressed as the axial displacement that the image may experience before the resultant image blur becomes excessive. Figure 40 shows the geometric relationship of the angular blur tolerance $\Delta\theta$ to the depth of focus δ_{\pm} . If the entrance pupil diameter is D_{ep} and the image distance is s_i , then the depth of focus is

$$\delta_{\pm} = \frac{s_i^2 \Delta\theta}{D_{\text{ep}} \pm s_i \Delta\theta} \quad (51)$$

or when $\delta \ll s_i$, the depth of focus becomes

$$\delta = \frac{s_i^2 \Delta\theta}{D_{\text{ep}}} \quad (52)$$

When $s_i = f$, then

$$\delta = f \Delta\theta \text{FN} \quad (53)$$

The *depth of field* is distance that the object may be moved without causing excessive image blur with a fixed image location. The distance at which a lens may be focused such that the depth of field extends to infinity is $s_o = D_{\text{ep}} / \Delta\theta$ and is called the hyperfocal distance.

If the lens system is diffraction-limited, then the depth of focus according to the Rayleigh criterion is given by

$$\delta = \pm \frac{\lambda}{2n_i \sin^2 u_i} \quad (54)$$

Diffraction-Limited Lenses

It is well known that the shape of the image irradiance of an incoherent, monochromatic point source formed by an aberration-free, circularly-symmetric lens system is described by the Airy pattern

$$E(r) = C_0 \left[\frac{2J_1(kD_{\text{ep}}r/2)}{kD_{\text{ep}}r} \right]^2 \quad (55)$$

where J_1 is the first-order Bessel function of the first kind, D_{ep} is the diameter of the entrance pupil, k is $2\pi/\lambda$, r is the radial distance from the center of the image to the observation point, and C_0 is a scaling factor. The angular radius β_{DL} of the first dark ring of the image is $1.22(\lambda/D_{ep})$. A common measure for the resolution is Lord Rayleigh's criterion that asserts that two point sources are just resolvable when the maximum of one Airy pattern coincides with the first dark ring of the second Airy pattern, that is, an angular separation of β_{DL} . Figure 41 presents a nomogram that can be used to make a rapid estimate of the diameter of angular or linear blur for a diffraction-limited system.

The modulation transfer function (MTF) at a specific wavelength λ for a circular entrance pupil can be computed by

$$MTF_{\lambda}(\Omega) = \frac{2}{\pi} \left[\arccos \Omega - \Omega \sqrt{1 - \Omega^2} \right] \quad \text{for } 0 \leq \Omega \leq 1 \quad (56)$$

where Ω is the normalized spatial frequency (ν/ν_{co}) with the maximum or cut-off frequency ν_{co} being $1/\lambda_o$ FN.

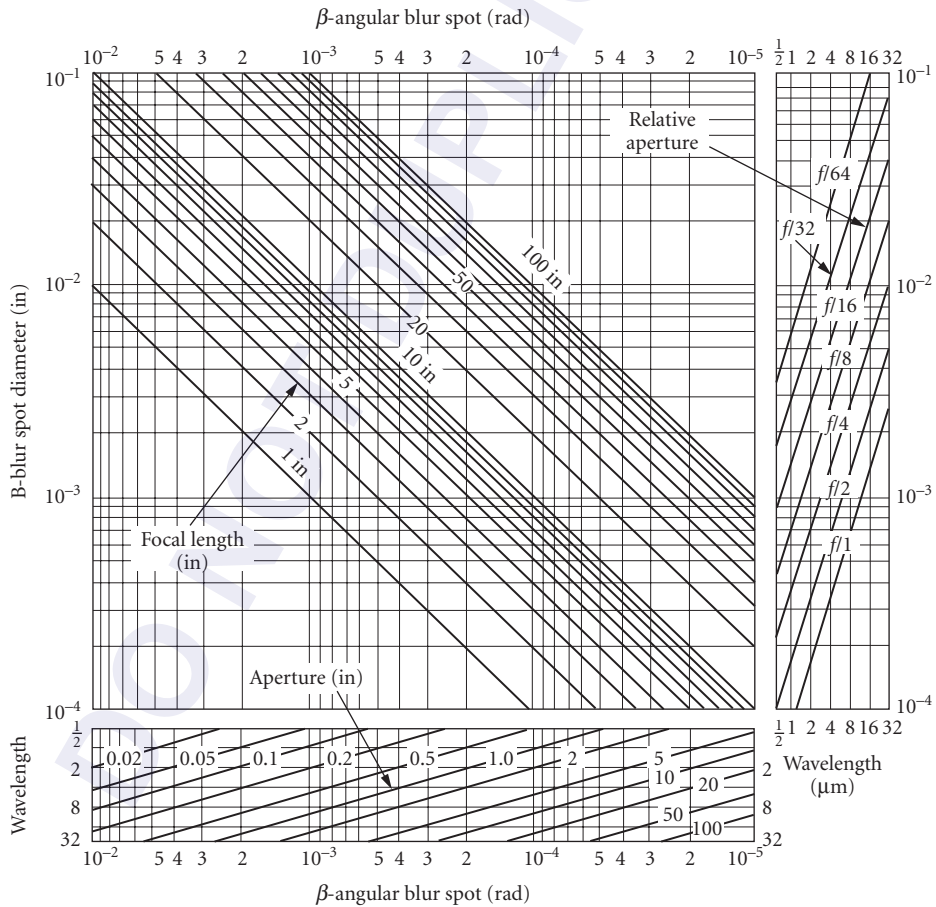


FIGURE 41 Estimation of the spot diameter for a diffraction-limited lens system. The diameter is that of the first dark ring of the Airy disk. (Smith, *Modern Optical Engineering*, McGraw-Hill, New York, 1990, p. 458.)

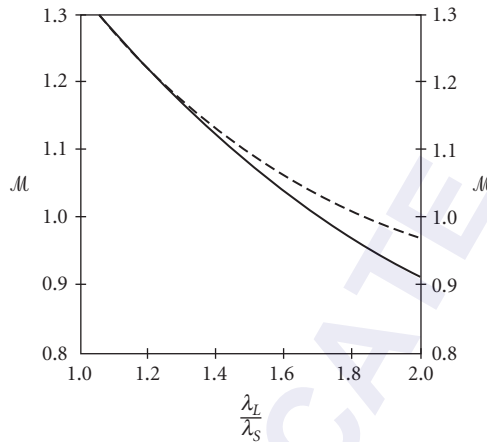


FIGURE 42 Variation of \mathcal{M} with λ_L/λ_S or $b + 1$ for $\tilde{\mathfrak{H}}(\lambda) = 1$ as the solid curve and $\tilde{\mathfrak{H}}(\lambda) = \lambda/\lambda_S$ as the dashed curve.

Should the source be polychromatic and the lens system be aberration-free, then the perfect-illumination irradiance distribution of a point source can be written as

$$E(r) = C_1 \int_0^\infty \tilde{\mathfrak{H}}(\lambda) \left[\frac{2J_1(kD_{\text{ep}} r/2)}{kD_{\text{ep}} r} \right]^2 d\lambda \quad (57)$$

where $\tilde{\mathfrak{H}}(\lambda)$ is the peak normalized spectral weighting factor and C_1 is a scaling factor.

A quick estimation of this ideal irradiance distribution can be made by invoking the central limit theorem to approximate this distribution by a Gaussian function, that is,

$$E(r) \approx C_2 e^{-(r^2/2\sigma^2)} \quad (58)$$

where C_2 is a scaling constant and σ^2 is the estimated variance of the irradiance distribution. When $\tilde{\mathfrak{H}}(\lambda) = 1$ in the spectral interval λ_S to λ_L and zero otherwise with $\lambda_S < \lambda_L$, an estimate of σ can be written as

$$\sigma = \frac{\mathcal{M}\lambda_L}{\pi D_{\text{ep}}} \quad (59)$$

where $\mathcal{M} = 1.335 - 0.625b + 0.025b^2 - 0.0465b^3$ with $b = (\lambda_L/\lambda_S) - 1$. Should $\tilde{\mathfrak{H}}(\lambda) = \lambda/\lambda_L$ in the spectral interval λ_S to λ_L and zero otherwise, which approximates the behavior of a quantum detector, $\mathcal{M} = 1.335 - 0.65b + 0.385b^2 - 0.099b^3$. The Gaussian estimate residual error is less than a few percent for $b = 0.5$ and remains useful even as $b \rightarrow 0$. Figure 42 contains plots of \mathcal{M} for both cases of $\tilde{\mathfrak{H}}(\lambda)$, where the abscissa is λ_L/λ_S . A useful estimation of the modulation transfer function for this polychromatic lens system is given by

$$\text{MTF}(v) \approx e^{-2(\pi\sigma v)^2} \quad (60)$$

where v is the spatial frequency. This approximation overestimates the MTF somewhat at lower spatial frequencies, while being rather a close fit at medium and higher spatial frequencies. The reason for this is that the central portion of the irradiance distribution is closely matched by the gaussian approximation, while the irradiance estimation beyond several Airy radii begins to degrade, therefore impacting the lower spatial frequencies. Nevertheless, this approximation can provide useful insight into expected performance limits.

17.22 BIBLIOGRAPHY

- Erfle, H., "Die optische abbildung durch kugelflaechen," Chap. 3 in S. Czapski und O. Eppenstein, *Grundzuege der Theorie der Optischen Instrumente nach Abbe*, 3d ed., H. Erfle and H. Boegehold, eds., Leipzig: Barth, 72–134 (1924).
- Goodman, D. S., "Basic Optical Instruments," Chap. 4 in *Geometrical and Instrumental Optics*, Daniel Malacara ed., *Methods of Experimental Physics* **25**, Academic Press, San Diego (1988).
- Hopkins, R. E., "The Components in the Basic Optical Systems," Chap. 3 in *Geometrical and Instrumental Optics*, Daniel Malacara ed., *Methods of Experimental Physics* **25**, Academic Press, San Diego (1988).
- Hopkins, R. E., "Geometrical Optics," Chap. 2 in *Geometrical and Instrumental Optics*, Daniel Malacara ed., *Methods of Experimental Physics* **25**, Academic Press, San Diego (1988).
- Johnson, R. B., "Basic Lenses," Chap. 13 in *Optical Engineer's Desk Reference*, W. L. Wolfe, ed., SPIE Press, Bellingham (2003).
- Johnson, R. B. and C. Feng, "A History of IR Lens Designs," *SPIE Critical Reviews* **CR37:3–18** (1991).
- Kingslake, R., "Basic Geometrical Optics," Chap. 6 in *Applied Optics and Optical Engineering*, **1**, Academic Press, New York (1965).
- Kingslake, R., *Lens Design Fundamentals*, Academic Press, New York (1978).
- Kingslake, R., *Optical System Design*, Academic Press, New York (1983).
- Kingslake, R., *A History of the Photographic Lens*, Academic Press, San Diego (1989).
- Kingslake, R., *Optics in Photography*, SPIE Press, Bellingham (1992).
- Kingslake, R., and R. Barry Johnson, *Lens Design Fundamentals*, 2d ed., Academic Press, New York (2009).
- Laikin, M., *Lens Design*, 4th ed., Marcel Dekker, New York (2006).
- Mahajan, V. N., *Aberration Theory Made Simple*, SPIE Press, Bellingham (1991).
- MIL-HDBK-141, *Optical Design*, Defense Supply Agency, Washington (1962).
- Mouroullis, R. and J. Macdonald, *Geometrical Optics and Optical Design*, Oxford University Press, New York (1997).
- Ray, S. F., *Applied Photographic Optics*, 2d ed., Focal Press, Oxford (1997).
- Schroeder, H., "Notiz betreffend die gaussischen hauptpunkte," *Astron. Nachrichten* **111:187–188** (1885).
- Smith, W. J., *Practical Optical System Layout*, McGraw-Hill, New York (1997).
- Smith, W. J., *Modern Optical Engineering*, 3d ed., McGraw-Hill, New York (2000).
- Smith, W. J., "Camera Lenses," Chap. 14 in *Optical Engineer's Desk Reference*, W. L. Wolfe, ed., SPIE Press, Bellingham (2003).
- Smith, W. J., *Modern Lens Design*, 2d ed., McGraw-Hill, New York (2005).

AFOCAL SYSTEMS

William B. Wetherell*

*Optical Research Associates
Framingham, Massachusetts*

18.1 GLOSSARY

BFL	back focal length
D	pupil diameter
ER_{cp}	eye relief common pupil position
ER_k	eye relief keplerian
e	exit pupil; eye space
F, F'	focal points
FFL	front focal length
h, h'	object and image heights
l, l'	object and image distances
M	angular magnification
m	linear, lateral magnification
n	refractive index
OR	object relief
o	entrance pupil; object space
P, P'	principal points
R	radius
TTL	total length
$\tan \alpha$	slope
x, y, z	cartesian coordinates
Δz	axial separation

*Retired.

18.2 INTRODUCTION

If collimated (parallel) light rays from an infinitely distant point source fall incident on the input end of a lens system, rays exiting from the output end will show one of three characteristics: (1) they will converge to a real point focus outside the lens system, (2) they will appear to diverge from a virtual point focus within the lens system, or (3) they will emerge as collimated rays that may differ in some characteristics from the incident collimated rays. In cases 1 and 2, the paraxial imaging properties of the lens system can be modeled accurately by a characteristic focal length and a set of fixed principal surfaces. Such lens systems might be called *focusing* or *focal* lenses, but are usually referred to simply as *lenses*. In case 3, a single finite focal length cannot model the paraxial characteristics of the lens system; in effect, the focal length is infinite, with the output focal point an infinite distance behind the lens, and the associated principal surface an infinite distance in front of the lens. Such lens systems are referred to as *afocal*, or without focal length. They will be called *afocal lenses* here, following the common practice of using “lens” to refer to both single element and multielement lens systems. They are the topic of this chapter.

The first afocal lens was the galilean telescope (to be described later), a visual telescope made famous by Galileo’s astronomical observations. It is now believed to have been invented by Hans Lipperhey in 1608.¹ Afocal lenses are usually thought of in the context of viewing instruments or attachments to change the effective focal length of focusing lenses, whose outputs are always collimated. In fact, afocal lenses can form real images of real objects. A more useful distinction between focusing and afocal lenses concerns which optical parameters are fixed, and which can vary in use. Focusing lenses have a fixed, finite focal length, can produce real images for a wide range of object distances, and have a linear magnification which varies with object distance. Afocal lenses have a fixed magnification which is independent of object distance, and the range of object distances yielding real images is severely restricted.

This chapter is divided into six sections, including this introduction. Section 18.3 reviews the gaussian (paraxial) image-forming characteristics of afocal lenses and compares them to the properties of focusing lenses. The importance of the optical invariant in designing afocal lenses is discussed. Section 18.4 reviews the keplerian telescope and its descendants, including both infinite conjugate and finite conjugate variants. Section 18.5 discusses the galilean telescope and its descendants. Thin-lens models are used in the Secs. 18.4 and 18.5 to define imaging characteristics and design principles for afocal lenses. Section 18.6 reviews relay trains and periscopes. The final section reviews reflecting and catadioptric afocal lenses.

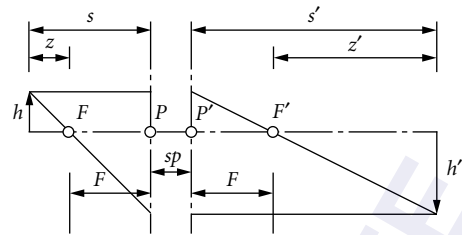
This chapter is based on an earlier article by Wetherell.² That article contains derivations of many of the equations appearing here, as well as a more extensive list of patents illustrating different types of afocal lens systems.

18.3 GAUSSIAN ANALYSIS OF AFOCAL LENSES

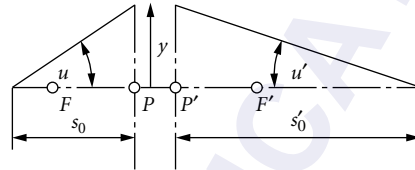
Afocal lenses differ from focusing lenses in ways that are not always obvious. It is useful to review the basic image-forming characteristics of focusing lenses before defining the characteristics unique to afocal lenses.

Focusing Lenses

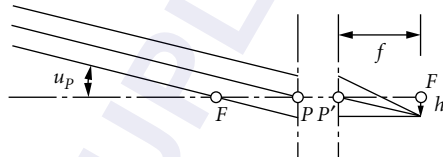
In this chapter, all lens elements are assumed to be immersed in air, so that object space and image space have the same index of refraction. Points in object space and image space are represented by two rectangular coordinate systems (x, y, z) and (x', y', z') , with the prime indicating image space. The z - and z' axes form a common line in space, the *optical axis* of the system. It is assumed, unless noted otherwise, that all lens elements are rotationally symmetric with respect to the optical axis. Under these conditions, the imaging geometry of a focusing lens can be defined in terms of two principal



(a) Finite conjugate model of focusing lens



(b) Computing image tilt with scheidpflug rule



(c) Infinite conjugate model of focusing lens

FIGURE 1 Imaging geometry of focusing lenses.

points P and P' , two focal points F and F' , and a single characteristic focal length f , as shown in Fig. 1. P , P' , F , and F' all lie on the optical axis.

The focal points F and F' , will be the origins for the coordinate systems (x, y, z) and (x', y', z') . If the origins are at P and P' , the coordinates will be given as (x, y, s) and (x', y', s') , where $s = z - f$ and $s' = z' + f$. Normal right-hand sign conventions are used for each set of coordinates, and light travels along the z axis from negative z toward positive z' , unless the optical system has internal mirrors. Figure 1a illustrates the terminology for finite conjugate objects.

Object points and image points are assumed to lie in planes normal to the optical axis, for paraxial computations. *Object distance* is specified by the axial distance to the object surface, z or s , and *image distance* by z' or s' . The two most commonly used equations relating image distance to object distance are

$$\frac{1}{s'} - \frac{1}{s} = \frac{1}{f} \quad (1)$$

and

$$zs' = -f^2 \quad (2)$$

For infinitely distant object points, $z' = 0$ and $s' = f$, and the corresponding image points will lie in the focal plane at F' .

To determine the actual distance from object plane to image plane, it is necessary to know the distance sp between P and P' . The value of sp is a constant specific to each real lens system, and may be either positive [moving object and image further apart than predicted by Eq. (1) or (2)] or negative (moving them closer together).

For rotationally symmetric systems, off-axis object and image coordinates can be expressed by the *object height* h and *image height* h' , where $h^2 = x^2 + y^2$ and $h'^2 = x'^2 + y'^2$. Object height and image height are related by the *linear magnification* m , where

$$m = \frac{h'}{h} = \frac{s'}{s} = \frac{z' + f}{z - f} \quad (3)$$

Since the product zz' is a constant, Eq. (3) implies that magnification varies with object distance.

The *principal surfaces* of a focusing lens intersect the optical axis at the principal points P and P' . In paraxial analysis, the principal surfaces are planes normal to the optical axis; for real lenses, they may be curved. The principal surfaces are conjugate image surfaces for which $m = +1.0$. This property makes the raytrace construction shown in Fig. 1a possible, since a ray traveling parallel to the optical axis in either object or image space must intersect the focal point in the conjugate space, and must also intersect both principal surfaces at the same height.

In real lenses, the object and image surfaces may be tilted or curved. Planes normal to the optical axis are still used to define object and image positions for off-axis object points, and to compute magnification. For tilted object surfaces, the properties of the principal surfaces can be used to relate object surface and image surface tilt angles, as shown in Fig. 1b. Tilted object and image planes intersect the optical axis and the two principal planes. The tilt angles with respect to the optical axis, u and u' , are defined by meridional rays lying in the two surfaces. The points at which conjugate tilted planes intersect the optical axis are defined by s_a and s'_a , given by Eq. (1). Both object and image planes must intersect their respective principal surfaces at the same height y , where $y = s_a \tan u = s'_a \tan u'$. It follows that

$$\frac{\tan u'}{\tan u} = \frac{s_a}{s'_a} = \frac{1}{m_a} \quad (4)$$

The geometry of Fig. 1b is known as the *Scheimpflug condition*, and Eq. (4) is the *Scheimpflug rule*, relating image to object tilt. The magnification m_a applies only to the axial image.

The height off axis of an infinitely distant object is defined by the principal ray angle u_p measured from F or P , as shown in Fig. 1c. In this case, the image height is

$$h' = f \tan u_p \quad (5)$$

A focusing lens which obeys Eq. (5) for all values of u_p within a specified range is said to be *distortion-free*: if the object is a set of equally spaced parallel lines lying in an object plane perpendicular to the optical axis, it will be imaged as a set of equally spaced parallel lines in an image plane perpendicular to the optical axis, with line spacing proportional to m .

Equations (1) through (5) are the basic gaussian imaging equations defining a perfect focusing lens. Equation (2) is sometimes called the *newtonian* form of Eq. (1), and is the more useful form for application to afocal lens systems.

Afocal Lenses

With afocal lenses, somewhat different coordinate system origins and nomenclature are used, as shown in Fig. 2. The object and image space reference points RO and RE are at conjugate image points. Since the earliest and most common use for afocal lenses is as an aid to the eye for viewing distant objects, image space is referred to as *eye space*. Object position is defined by a right-hand coordinate system (x_o, y_o, z_o) centered on reference point RO. Image position in eye space is defined by coordinates (x_e, y_e, z_e) centered on RE.

Because afocal lenses are most commonly used for viewing distant objects, their imaging characteristics are usually specified in terms of *angular magnification* M , *entrance pupil diameter* D_o , and total field of view. Figure 2a models an afocal lens used at infinite conjugates. Object height off axis is

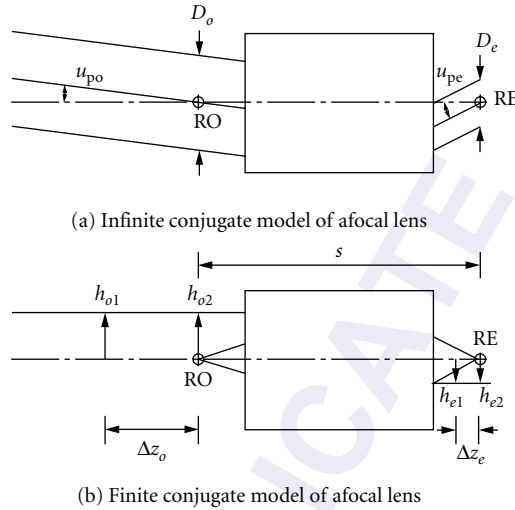


FIGURE 2 Imaging geometry of focusing lenses.

defined by the principal ray angle u_{po} , and the corresponding image height is defined by u_{pe} . Objects viewed through the afocal lens will appear to be magnified by a factor M , where

$$\tan u_{pe} = M \tan u_{po} \quad (6)$$

If M is negative, as in Fig. 2a, the image will appear to be inverted. [Strictly speaking, since RO and RE are separated by a distance S , the apparent magnification seen by an eye at RE with and without the afocal lens will differ slightly from that indicated by Eq. (6) for nearby objects.]²

The imaging geometry of an afocal lens for finite conjugates is illustrated in Fig. 2b. Since a ray entering the afocal lens parallel to the optical axis will exit the afocal lens parallel to the optical axis, it follows that the linear magnification m relating object height h_o and image height h_e must be invariant with object distance. The linear magnification m is the inverse of the angular magnification M :

$$m = \frac{h_e}{h_o} = \frac{1}{M} \quad (7)$$

The axial separation Δz_e of any two images h_{e1} and h_{e2} is related to the separation Δz_o of the corresponding objects h_{o1} and h_{o2} by

$$\Delta z_e = m^2 \Delta z_o = \frac{\Delta z_o}{M^2} \quad (8)$$

It follows that any convenient pair of conjugate image points can be chosen as reference points RO and RE. Given the location of RO, the reference point separation S , and the magnifications $m = 1/M$, the imaging geometry of a rotationally symmetric distortion-free afocal lens can be given as

$$x_e = mx_o = \frac{x_o}{M} \quad y_e = my_o = \frac{y_o}{M} \quad z_e = m^2 z_o = \frac{z_o}{M^2} \quad (9)$$

Equation (9) is a statement that coordinate transformation between object space and eye space is rectilinear for afocal lenses, and is solely dependent on the *afocal magnification* M and the location

of two conjugate reference points RO and RE. The equations apply (paraxially) to all object and image points independent of their distances from the afocal lens. Any straight line of equally spaced object points will be imaged as a straight line of equally spaced image points, even if the line does not lie in a plane normal to the optical axis. Either RO or RE may be chosen arbitrarily, and need not lie on the axis of symmetry of the lens system, so long as the z_o - and z_e axes are set parallel to the axis of symmetry.

A corollary of invariance in lateral and axial linear magnification is invariance in angular magnification. Equation (6) thus applies to any ray traced through the afocal system, and to tilted object and image surfaces. In the latter context, Eq. (6) can be seen as an extension of Eq. (4) to afocal lenses.

The *eye space pupil diameter* D_e is of special importance to the design of visual instruments and afocal attachments: D_e must usually be large enough to fill the pupil of the associated instrument or eye. The *object space pupil diameter* D_o is related to D_e by

$$D_e = \frac{D_o}{M} = mD_o \quad (10)$$

(The more common terminology *exit pupil* and *entrance pupil* will be used later in this chapter.)

Subjective Aspects of Afocal Imagery

The angular magnification M is usually thought of in terms of Eq. (6), which is often taken to indicate that an afocal lens projects an image which is M -times as large as the object. (See, for example, Fig. 5.88 in Hecht and Zajac.)³ Equation (9) shows that the image height is actually $1/M$ -times the object height (i.e., smaller than the object when $|M| > 1$). Equation (9) also shows, however, that the image distance is reduced to $1/M^2$ -times the object distance, and it is this combination of linear height reduction and quadratic distance reduction which produces the subjective appearance of magnification. Equation (6) can be derived directly from Eq. (9).

$$\tan u_{pe} = \frac{y_e}{z_e} = \frac{y_o/M}{z_o/M^2} = M \tan u_{po}$$

Equation (9) is therefore a more complete model than Eq. (6) for rotationally symmetric, distortion-free afocal lenses.

Figure 3 illustrates two subjective effects which arise when viewing objects through afocal lenses. In Fig. 3a for which $M = +3\times$, Eq. (9) predicts that image dimensions normal to the optical axis will be reduced by $1/3$, while image dimensions along the optical axis will be reduced by $1/9$. The image of the cube in Fig. 3a looks three times as tall and wide because it is nine times closer, but it appears compressed by a factor of 3 in the axial direction, making it look like a cardboard cutout. This subjective compression, most apparent when using binoculars, is intrinsic to the principle producing angular magnification, and is independent of pupil spacing in the binoculars.

Figure 3a assumes the optical axis is horizontal within the observer's reference framework. If the axis of the afocal lens is not horizontal, the afocal lens may create the illusion that horizontal surfaces are tilted. Figure 3b represents an $M = +7\times$ afocal lens whose axis is tilted 10° to a horizontal surface. Equation (6) can be used to show that the image of this surface is tilted approximately 51° to the axis of the afocal lens, creating the illusion that the surface is tilted 41° to the observer's horizon. This illusion is most noticeable when looking downward at a surface known to be horizontal, such as a body of water, through a pair of binoculars.

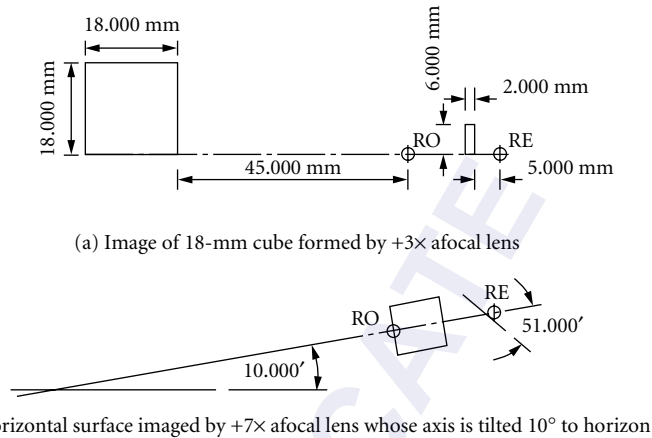


FIGURE 3 Subjective aspects of afocal imagery.

Afocal Lenses and the Optical Invariant

Equations (6) and (7) can be combined to yield

$$h_e \tan u_{pe} = h_o \tan u_{po} \quad (11)$$

which is a statement of the optical invariant as applied to distortion-free afocal lenses. Neither u_{po} nor u_{pe} is typically larger than 35° – 40° in distortion-free afocal lenses, although there are examples with distortion where $u_{po} \rightarrow 90^\circ$. Given a limit on one angle, Eq. (11) implies a limit on the other angle related to the ratio $h_o/h_e = D_o/D_e$. Put in words, *the ratio D_o/D_e cannot be made arbitrarily large without a corresponding reduction in the maximum allowable field of view*. All designers of afocal lens systems *must* take this fundamental principle into consideration.

18.4 KEPLERIAN AFOCAL LENSES

A simple afocal lens can be made up of two focusing lenses, an *objective* and an *eyepiece*, set up so that the rear focal point of the objective coincides with the front focal point of the eyepiece. There are two general classes of simple afocal lenses, one in which both focusing lenses are positive, and the other in which one of the two is negative. Afocal lenses containing two positive lenses were first described by Johannes Kepler in *Dioptrice*, in 1611,⁴ and are called *keplerian*. Lenses containing a negative eyepiece are called *galilean*, and will be discussed separately. Generally, afocal lenses contain at least two powered surfaces. The simplest model for an afocal lens consists of two thin lenses.

Thin-Lens Model of a Keplerian Afocal Lens

Figure 4 shows a thin-lens model of a keplerian telescope. The focal length of its objective is f_o and the focal length of its eyepiece is f_e . Its properties can be understood by tracing two rays, ray 1 entering the objective parallel to the optical axis, and ray 2 passing through F_o , the front focal

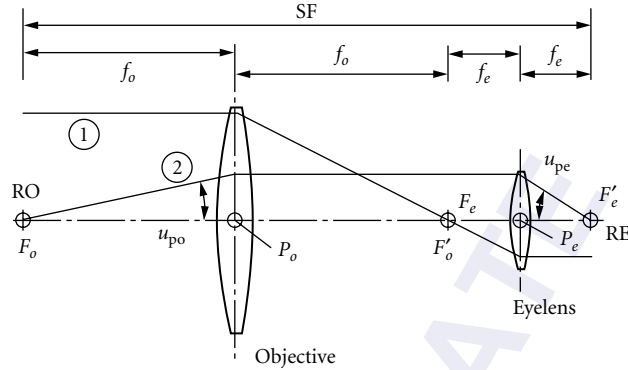


FIGURE 4 Thin-lens model of keplerian afocal lens.

point of the objective. Ray 1 leads directly to the linear magnification m , and ray 2 to the angular magnification M :

$$m = -\frac{f_e}{f_o} \quad M = -\frac{f_o}{f_e} = \frac{\tan u_{pe}}{\tan u_{po}} \quad (12)$$

Equation (12) makes the relationship of afocal magnification to the Scheimpflug rule of Eq. (4) more explicit, with focal lengths f_o and f_e substituting for s_a and s'_a .

The second ray shows that placing the reference point RO at F_o will result in the reference point RE falling on F'_e , the rear focal point of the eyepiece. The reference point separation for RO in this location is

$$SF = 2f_e + 2f_o = 2(1-M)f_e = 2(1-m)f_o \quad (13)$$

Equation (13) can be used as a starting point for calculating any other locations for RO and RE, in combination with Eq. (9).

One additional generalization can be drawn from Fig. 4: the ray passing through F_o will emerge from the objective parallel to the optical axis. It will therefore also pass through F'_e even if the spacing between objective and eyepiece is increased to focus on nearby objects. Thus the angular magnification remains invariant, if u_{po} is measured from F_o and u_{pe} is measured from F'_e , even when adjusting the eyepiece to focus on nearby objects makes the lens system depart from being strictly afocal.

The simple thin-lens model of the keplerian telescope can be extended to systems composed of two real focusing lenses if we know their focal lengths and the location of each lens' front and rear focal points. Equation (12) can be used to derive M , and SF can be measured. Equation (9) can then be used to compute both finite and infinite conjugate image geometry.

Eye Relief Manipulation

The earliest application of keplerian afocal lenses was to obtain magnified views of distant objects. To view distant objects, the eye is placed at RE. An important design consideration in such instruments is to move RE far enough away from the last surface of the eyepiece for comfortable viewing. The distance from the last optical surface to the exit pupil at RE is called the *eye relief* ER. One way to increase eye relief ER is to move the entrance pupil at RO toward the objective. Most telescopes and binoculars have the system stop at the first surface of the objective, coincident with the entrance pupil, as shown in Fig. 5a.

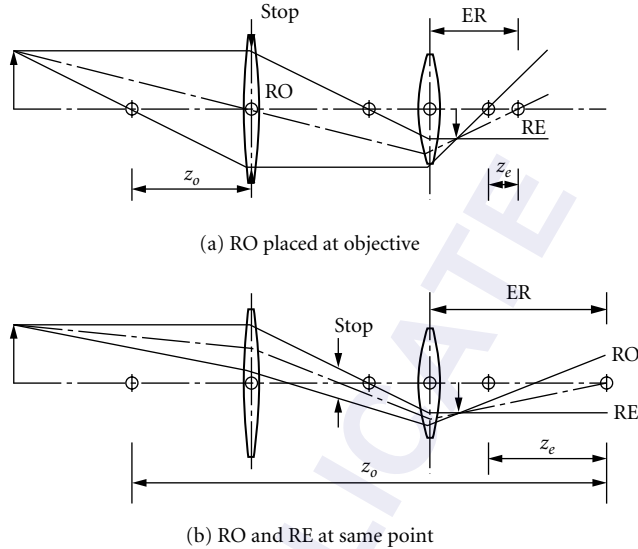


FIGURE 5 Increasing eye relief ER by moving stop.

In the thin-lens model of Fig. 5a, RO is moved a distance $z_o = f_o$ to place it at the objective. Thus RE must move a distance $z_e = f_o/M^2 = -f_e/M$, keeping in mind that M is negative in this example. Thus for a thin-lens keplerian telescope with its stop at the objective, the eye relief ER_k is

$$ER_k = \frac{(M-1)}{M} f_e \quad (14)$$

It is possible to increase the eye relief further by placing the stop inside the telescope, moving the location of RO into virtual object space. Figure 5b shows an extreme example of this, where the virtual location of RO has been matched to the real location of RE. For this common-pupil-position case, the eye relief ER_{cp} is

$$ER_{cp} = \frac{(M-1)}{(M+1)} f_e \quad (15)$$

A price must be paid for locating the stop inside the afocal lens, in that the elements ahead of the stop must be increased in diameter if the same field of view is to be covered without vignetting.

The larger the magnitude of M , the smaller the gain in ER yielded by using an internal stop. To increase the eye relief further, it is necessary to make the objective and/or the eyepiece more complex, increasing the distance between F_o and the first surface of the objective, and between the last surface of the eyepiece and F_e' . If this is done, placing RO at the first surface of the objective will further increase ER.

Figure 6 shows a thin-lens model of a telephoto focusing lens of focal length f_t . For convenience, a zero Petzval sum design is used, for which $f_1 = f$ and $f_2 = -f$. Given the telephoto's focal length f_t and the lens separation d , the rest of the parameters shown in Fig. 6 can be defined in terms of the constant $C = d/f_t$. The component focal length f , back focal length bfl, and front focal length ffl, are given by

$$f = f_t C^{1/2} \quad \text{bfl} = f_t(1 - C^{1/2}) \quad \text{ffl} = f_t(1 + C^{1/2}) \quad (16)$$

and the total physical length ttl and focal point separation sf are given by

$$\text{ttl} = f_t(1 + C - C^{1/2}) \quad \text{sf} = f_t(2 + C) \quad (17)$$

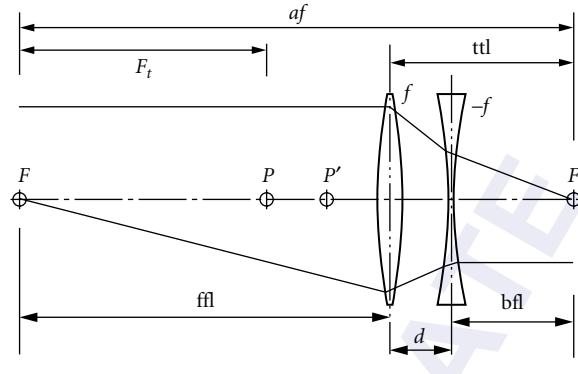


FIGURE 6 Zero Petzval sum telephoto lens.

The maximum gain in eye relief will be obtained by using telephoto designs for both objective and eyepiece, with the negative elements of each facing each other. Two cases are of special interest. First, t_{tl} can be minimized by setting $C = 0.25$ for both objective and eyepiece. In this case, the eye relief $ER_{t_{tl}}$ is

$$ER_{t_{tl}} = 1.5 \frac{(M-1)}{M} f_e = 1.5 ER_k \quad (18)$$

Second, s_f can be maximized by setting $C = 1.0$ for both objective and eyepiece. This places the negative element at the focal plane, merging the objective and eyepiece negative elements into a single negative field lens. The eye relief in this case, ER_{s_f} is

$$ER_{s_f} = 2.0 \frac{(M-1)}{M} f_e = 2.0 ER_k \quad (19)$$

Placing a field lens at the focus between objective and eyepiece can be problematical, when viewing distant objects, since dust or scratches on the field lens will be visible. If a reticle is required, however, it can be incorporated into the field lens. Equations (14), (18), and (19) show that significant gains in eye relief can be made by power redistribution. In the example of Eq. (18), the gain in ER is accompanied by a reduction in the physical length of the optics, which is frequently beneficial.

Terrestrial Telescopes

Keplerian telescopes form an inverted image, which is considered undesirable when viewing earth-bound objects. One way to make the image erect, commonly used in binoculars, is to incorporate erecting prisms. A second is to insert a relay stage between objective and eyepiece, as shown in Fig. 7. The added relay is called an *image erector*, and telescopes of this form are called *terrestrial telescopes*. (The keplerian telescope is often referred to as an *astronomical telescope*, to distinguish it from terrestrial telescopes, since astronomers do not usually object to inverted images. *Astronomical* has become ambiguous in this context, since it now more commonly refers to the very large aperture reflecting objectives found in astronomical observatories. *Keplerian* is the preferred terminology.) The terrestrial telescope can be thought of as containing an objective, eyepiece, and image erector, or as containing two afocal relay stages.

There are many variants of terrestrial telescopes made today, in the form of binoculars, theodolites, range finders, spotting scopes, rifle scopes, and other military optical instrumentation. All are offshoots of the keplerian telescope, containing a positive objective and a positive eyepiece, with intermediate

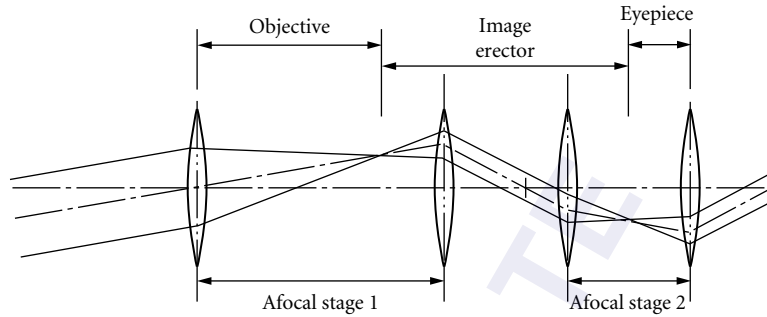


FIGURE 7 Terrestrial telescope.

relay stages to perform special functions. Patrick⁵ and Jacobs⁶ are good starting points for obtaining more information.

Field-of-View Limitations in Keplerian and Terrestrial Telescopes

The maximum allowable eye space angle u_{pe} and magnification M set an upper limit on achievable fields of view, in accordance with Eq. (11). MIL-HDBK-141⁷ lists one eyepiece design for which the maximum $u_{pe} = 36^\circ$. If $M = 7\times$, using that eyepiece allows a 5.9° maximum value for u_{po} . It is a common commercial practice to specify the total field of view FOV as the width in feet which subtends an angle $2u_{po}$ from 1000 yd away, even when the pupil diameter is given in millimeters. FOV is thus given by

$$\text{FOV} = 6000 \tan u_{po} = \frac{6000}{M} \tan u_{pe} \quad (20)$$

For our $7\times$ example, with $u_{pe} = 36^\circ$, $\text{FOV} = 620$ ft at 1000 yd. For commercial 7×50 binoculars ($M = 7\times$ and $D_o = 50$ mm), $\text{FOV} = 376$ ft at 1000 yd is more typical.

Finite Conjugate Afocal Relays

If an object is placed in contact with the front surface of the keplerian telescope of Fig. 5, its image will appear a distance ER_k behind the last surface of the eyepiece, in accordance with Eq. (14). There is a corresponding *object relief* distance $OR_k = M^2 ER_k$ defining the position of an object that will be imaged at the output surface of the eyepiece, as shown in Fig. 8. OR_k and ER_k define the portions of object space and eye space within which real images can be formed of real objects with a simple keplerian afocal lens.

$$OR_k = M(M-1)f_e \quad (21)$$

Object relief is enlarged by the power redistribution technique used to extend eye relief. Thus there is a minimum total length design corresponding to Eq. (18), for which the object relief OR_{tl} is

$$OR_{tl} = 1.5M(M-1)f_e \quad (22)$$

and a maximum eye relief design corresponding to Eq. (19), for which OR_{sf}

$$OR_{sf} = 2.0M(M-1)f_e \quad (23)$$

is also maximized.

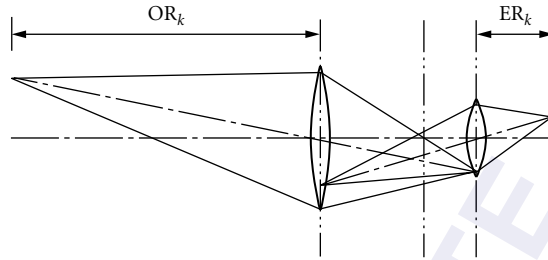


FIGURE 8 Finite conjugate keplerian afocal lens showing limits on usable object space and image space.

Figure 9 shows an example of a zero Petzval sum finite conjugate afocal relay designed to maximize OR and ER by placing a negative field lens at the central infinite conjugate image. Placing the stop at the field lens means that the lens is *telecentric* (principal rays parallel to the optical axis) in both object and eye space. As a result, magnification, principal ray angle of incidence on object and image surface, and cone angle are all invariant over the entire range of OR and ER for which there is no vignetting. Magnification and cone angle invariance means that object and image surfaces can be tilted with respect to the optical axis without introducing keystone or variation in image irradiance over the field of view. Having the principal rays telecentric means that object and image position can be adjusted for focus without altering magnification. It also means that the lens can be defocused without altering magnification, a property very useful for unsharp masking techniques used in the movie industry.

One potential disadvantage of telecentric finite conjugate afocal relays is evident from Fig. 9: to avoid vignetting, the apertures of both objective and eyepiece must be larger than the size of the associated object and image. While it is possible to reduce the diameter of either the objective or the eyepiece by shifting the stop to make the design nontelecentric, the diameter of the other lens group becomes larger. Afocal relays are thus likely to be more expensive to manufacture than focusing lens relays, unless object and image are small.

Finite conjugate afocal lenses have been used for alignment telescopes,⁸ for laser velocimeters,⁹ and for automatic inspection systems for printed circuit boards.¹⁰ In the last case, invariance of magnification, cone angle, and angle of incidence on a tilted object surface make it possible to measure the volume of solder beads automatically with a computerized video system. Finite conjugate afocal lenses are also used as Fourier transform lenses.¹¹ Brief descriptions of these applications are given in Wetherell.²

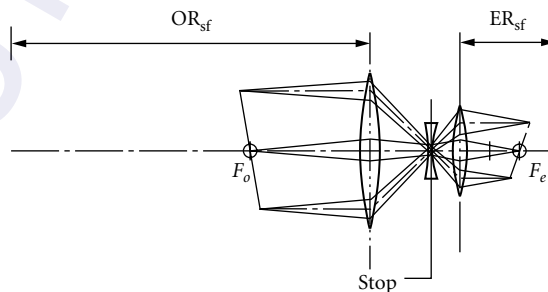


FIGURE 9 Finite conjugate afocal relay configured to maximize eye relief ER and object relief OR. Stop at common focus collimates principal rays in both object space and eye space.

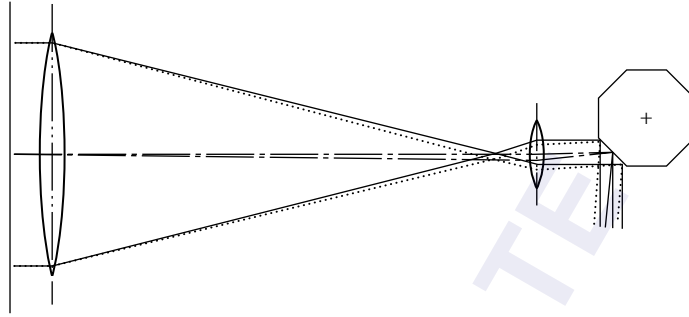


FIGURE 10 Afocal lens scanner geometry.

Afocal Lenses for Scanners

Many optical systems require scanners, and if the apertures of the systems are large enough, it is preferable to place the scanner inside the system. Although scanners have been designed for use in convergent light, they are more commonly placed in collimated light (see Chap. 30, “Scanners,” in this volume, Marshall,¹² and Chap. 7 in Lloyd,¹³ for descriptions of scanning techniques). A large aperture objective can be converted into a high magnification keplerian afocal lens with the aid of a short focal length eyepiece collimator, as shown in Fig. 10, providing a pupil in a collimated beam in which to insert a scanner. For the polygonal scanner shown, given the desired scan angle and telescope aperture diameter, Eq. (11) will define the combination of scanner facet size and number of facets needed to achieve the desired scanning efficiency. Scanning efficiency is the time it takes to complete one scan divided by the time between the start of two sequential scans. It is tied to the ratio of facet length to beam diameter, the amount of vignetting allowed within a scan, the number of facets, and the angle to be scanned.

Two limitations need to be kept in mind. First, the optical invariant will place an upper limit on M for the given combination of D_o and u_{po} , since there will be a practical upper limit on the achievable value of u_{pe} . Second, it may be desirable in some cases for the keplerian afocal relay to have enough barrel distortion so that Eq. (6) becomes

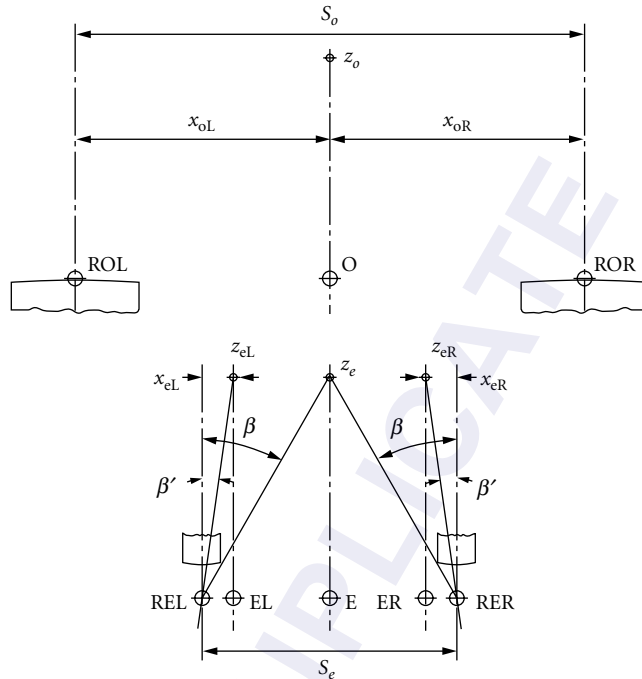
$$u_{pe} = Mu_{po} \quad (24)$$

An afocal lens obeying Eq. (24) will convert a constant rotation rate of the scan mirror into a constant angular scan rate for the external beam. The same property in “f-theta” focusing lenses is used to convert a constant angular velocity scanner rotation rate into a constant linear velocity rate for the recording spot of light.

The above discussion applies to scanning with a point detector. When the detector is a linear diode array, or when a rectangular image is being projected onto moving film, the required distortion characteristics for the optical system may be more complex.

Imaging in Binoculars

Most commercial binoculars consist of two keplerian afocal lenses with internal prismatic image erectors. Object and image space coordinates for binoculars of this type are shown schematically in Fig. 11. Equation (9) can be applied to Fig. 11 to analyze their imaging properties. In most binoculars, the spacing S_o between objectives differs from the spacing S_e between eyepieces, and S_o may be either larger or smaller than S_e . Each telescope has its own set of reference points, ROL and REL for the left telescope, and ROR and RER for the right. Object space is a single domain with a single origin O . The object point at z_o , midway between the objective axes, will be x_{oL} units to the right


FIGURE 11 Imaging geometry of binoculars.

of the left objective axis, and x_{oR} units to the left of the right objective axis. In an ideal binocular system, the images of the object formed by the two telescopes would merge at one point, z_o units in front of eye space origin E . This will happen if $S_o = MS_e$, so that $x_{eL} = x_{oL}/M$ and $x_{eR} = x_{oR}/M$. In most modern binoculars, however, $S_o \ll MS_e$, and separate eye space reference points EL and ER will be formed for the left and right eye. As a result, each eye sees its own eye space, and while they overlap, they are not coincident. This property of binoculars can affect stereo acuity² and eye accommodation for the user.

It is normal for the angle at which a person's left-eye and right-eye lines of sight converge to be linked to the distance at which the eyes focus. (In my case, this linkage was quite strong before I began wearing glasses.) Eyes focused for a distance z_e normally would converge with an angle β , as shown in Fig. 11. When $S_o \ll MS_e$, as is commonly the case, the actual convergence angle β' is much smaller. A viewer for whom focus distance is strongly linked to convergence angle may find such binoculars uncomfortable to use for extended periods, and may be in need of frequent focus adjustment for different object distances.

A related but more critical problem arises if the axes of the left and right telescopes are not accurately parallel to each other. Misalignment of the axes requires the eyes to twist in unaccustomed directions to fuse the two images, and refocusing the eyepiece is seldom able to ease the burden. Jacobs⁶ is one of the few authors to discuss this problem. Jacobs divides the axes misalignment into three categories: (1) misalignments requiring a divergence D of the eye axes to fuse the images, (2) misalignments requiring a convergence C , and (3) misalignments requiring a vertical displacement V . The tolerance on allowable misalignment in minutes of arc is given by Jacobs as

$$D = 7.5/(M-1) \quad C = 22.5/(M-1) \quad V = 8.0/(M-1) \quad (25)$$

Note that the tolerance on C , which corresponds to convergence to focus on nearby objects, is substantially larger than the tolerances on D and V .

18.5 GALILEAN AND INVERSE GALILEAN AFOCAL LENSES

The combination of a positive objective and a negative eyepiece forms a *galilean* telescope. If the objective is negative and the eyepiece positive, it is referred to as an *inverse galilean* telescope. The galilean telescope has the advantage that it forms an erect image. It is the oldest form of visual telescope, but it has been largely replaced by terrestrial telescopes for magnified viewing of distant objects, because of field-of-view limitations. In terms of number of viewing devices manufactured, there are far more inverse galilean than galilean telescopes. Both are used frequently as power-changing attachments to change the effective focal length of focusing lenses.

Thin-Lens Model of a Galilean Afocal Lens

Figure 12 shows a thin-lens model of a galilean afocal lens. The properties of galilean lenses can be derived from Eqs. (9), (12), and (13). Given that f_e is negative and f_o is positive, M is positive, indicating an erect image. If RO is placed at the front focal point of the objective, RE is a virtual pupil buried inside the lens system. In fact, galilean lenses cannot form real images of real objects under any conditions, and at least one pupil will always be virtual.

Field of View in Galilean Telescopes

The fact that only one pupil can be real places strong limitations on the use of galilean telescopes as visual instruments when $M \gg 1x$. Given the relationship $\Delta z_o = M^2 \Delta z_e$, moving RE far enough outside the negative eyepiece to provide adequate eye relief moves RO far enough into virtual object space to cause drastic vignetting at even small field angles. Placing RE a distance ER behind the negative lens moves RO to the position shown in Fig. 13, SF' units behind RE, where

$$SF' = (M^2 - 1)ER - (M - 1)^2 f_e \tag{26}$$

In effect, the objective is both field stop and limiting aperture, and vignetting defines the maximum usable field of view. The maximum acceptable object space angle u_{po} is taken to be that for the principal ray which passes just inside D_o , the entrance pupil at the objective. If the F-number of the objective is $FN_{ob} = f_o/D_o$, then

$$\tan u_{po} = \frac{-f_e}{2FN_{ob}(MER + f_e - Mf_e)} \tag{27}$$

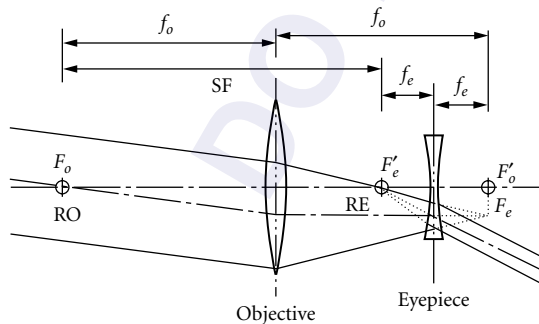


FIGURE 12 Thin-lens model of galilean afocal lens.

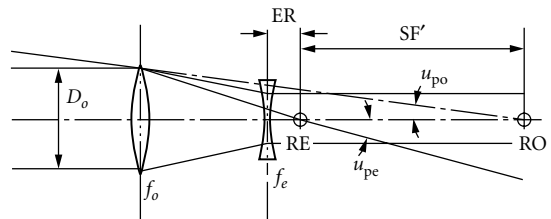


FIGURE 13 Galilean field-of-view limitations.

For convenience, assume $ER = -f_e$. In this case, Eq. (27) reduces to

$$\tan u_{po} = \frac{1}{2FN_{ob}(2M-1)} \quad (28)$$

For normal achromatic doublets, $FN_{ob} \geq 4.0$. For $M = 3x$, in this case, Eq. (28) indicates that $u_{po} \leq 1.43^\circ$ (FOV ≤ 150 ft at 1000 yd). For $M = 7x$, $u_{po} \leq 0.55^\circ$ (FOV ≤ 57.7 ft at 1000 yd). The effective field of view can be increased by making the objective faster and more complex, as can be seen in early patents by von Rohr¹⁴ and Erfle.¹⁵ In current practice, galilean telescopes for direct viewing are seldom made with M larger than $1.5x - 3.0x$. They are more typically used as power changers in viewing instruments, or to increase the effective focal length of camera lenses.¹⁶

Field of View in Inverse Galilean Telescopes

For inverse galilean telescopes, where $M \ll 1x$, adequate eye relief can be achieved without moving RO far inside the first surface of the objective. Inverse galilean telescopes for which $u_{po} \rightarrow 90^\circ$ are very common in the form of security viewers¹⁷ of the sort shown in Fig. 14, which are built into doors in hotel rooms, apartments, and many houses. These may be the most common of all optical systems more complex than eyeglasses. The negative objective lens is designed with enough distortion to allow viewing of all or most of the forward hemisphere, as shown by the principal ray in Fig. 14.

Inverse galilean telescopes are often used in camera view finders.¹⁸ These present reduced scale images of the scene to be photographed, and often have built in arrangements to project a frame of lines representing the field of view into the image. Inverse galilean power changers are also used to increase the field of view of submarine periscopes and other complex viewing instruments, and to reduce the effective focal length of camera lenses.¹⁹

Anamorphic Afocal Attachments

Afocal attachments can compress or expand the scale of an image in one axis. Such devices are called *anamorphosers*, or *anamorphic afocal attachments*. One class of anamorphoser is the cylindrical galilean telescope, shown schematically in Fig. 15a. Cox²⁰ and Harris²¹ have patented representative examples. The keplerian form is seldom if ever used, since a cylindrical keplerian telescope would introduce image inversion in one direction. Anamorphic compression can also be obtained using two prisms, as shown in Fig. 15b. The adjustable magnification anamorphoser patented by Luboshez²² is a good example of prismatic anamorphosers. Many anamorphic attachments were developed in the 1950s for the movie industry for use in wide-field cameras and projectors. An extensive list of both types will be found in Wetherell.²

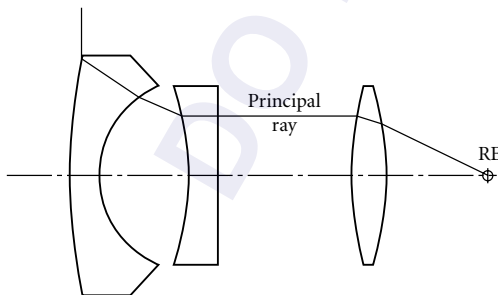


FIGURE 14 Inverse galilean security viewer with hemispheric field of view.

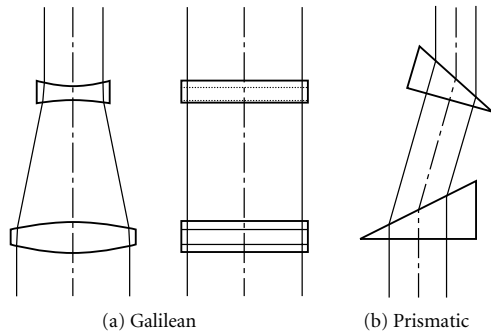


FIGURE 15 Anamorphic afocal attachments.

Equation (9) can be modified to describe anamorphic afocal lenses by specifying separate afocal magnifications M_x and M_y for the two axes. One important qualification is that separate equations are needed for object and image distances for the x and y planes. In general, anamorphic galilean attachments work best when used for distant objects, where any difference in x -axis and y -axis focus falls within the depth of focus of the associated camera lens. If it is necessary to use a galilean anamorphoser over a wide range of object distances, it may be necessary to add focus adjustment capabilities within the anamorphoser.

18.6 RELAY TRAINS AND PERISCOPEs

There are many applications where it is necessary to perform remote viewing because the object to be viewed is in an environment hostile to the viewer, or because the object is inaccessible to the viewer without unacceptable damage to its environment. Military applications³ fall into the former category, and medical applications²³ fall into the latter. For these applications, instrumentation is needed to *collect* light from the object, *transport* the light to a location more favorable for viewing, and *dispense* the light to the viewing instruments or personnel. Collecting and dispensing optical images is done with focusing lenses, typically. There are three image transportation techniques in common use today: (1) sense the image with a camera and transport the data electronically, (2) transport the light pattern with a coherent fiber optics bundle, and (3) transport the light pattern with a relay lens or train of relay lenses. The first two techniques are outside the scope of this chapter. *Relay trains*, however, are commonly made up of a series of unit power afocal lenses, and are one of the most important applications of finite conjugate afocal lenses.

Unit Power Afocal Relay Trains

Several factors are important in designing relay trains. First, it is desirable to minimize the number of relay stages in the relay train, both to maximize transmittance and to minimize the field curvature caused by the large number of positive lenses. Second, the outside diameter of the relay train is typically restricted (or a single relay lens could be used), so the choice of image and pupil diameter within the relay is important. Third, economic considerations make it desirable to use as many common elements as possible, while minimizing the total number of elements. Fourth, it is desirable to keep internal images well clear of optical surfaces where dust and scratches can obscure portions of the image. Fifth, the number of relay stages must be either odd or even to ensure the desired output image orientation.

Figure 16 shows thin-lens models of the two basic afocal lens designs which can be applied to relay train designs. Central to both designs is the use of symmetry fore and aft of the central stop to control coma, distortion, and lateral color, and matching the image diameter D_i and stop diameter D_s to

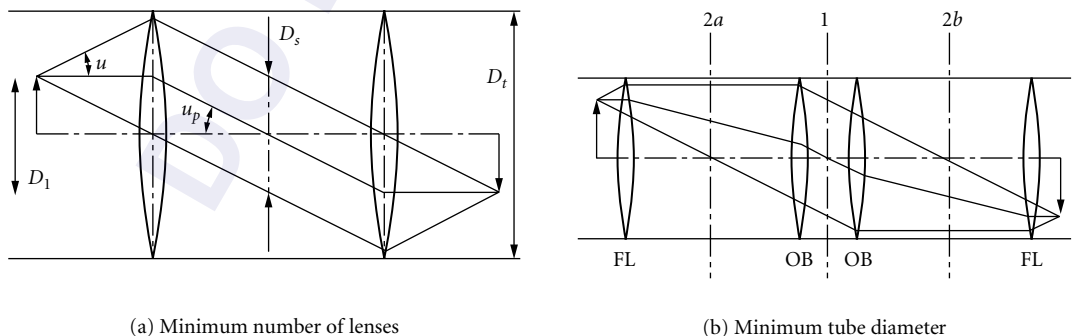


FIGURE 16 Basic unit power afocal relay designs.

maximize the stage length to diameter ratio. In paraxial terms, if $D_i = D_s$, then the marginal ray angle u matches the principal ray angle u_p , in accordance with the optical invariant. If the relay lens is both aplanatic and distortion free, a better model of the optical invariant is

$$D_i \sin u = D_s \tan u_p \quad (29)$$

and either the field of view $2u_p$ or the numerical aperture $NA = n \sin u$ must be adjusted to match pupil and image diameters. For some applications, maximizing the optical invariant which can pass through a given tube diameter D_i in a minimum number of stages is also critical.

If maximizing the ratio $D_i \sin u/D_i$ is not critical, Fig. 16a shows how the number of elements can be minimized by using a keplerian afocal lens with the stop at the common focus, eliminating the need for field lenses between stages. The required tube diameter in this example is at least twice the image diameter. If maximizing $D_i \sin u/D_i$ is critical, field lenses FL must be added to the objectives OB as shown in Fig. 16b, and the field lenses should be located as close to the image as possible within limits set by obstructions due to dirt and scratches on the field lens surfaces. Symmetry fore and aft of the central stop at 1 is still necessary for aberration balancing. If possible within performance constraints, symmetry of OB and FL with respect to the planes 2a and 2b is economically desirable, making OB and FL identical.

For medical instruments, where minimizing tube diameter is critical, variants of the second approach are common. The rod lens design²⁴ developed by H. H. Hopkins²⁵ can be considered an extreme example of either approach, making a single lens so thick that it combines the functions of OB and FL. Figure 17a shows an example from the first of two patents by McKinley.^{26,27} The central element in each symmetrical cemented triplet is a sphere. Using rod lenses does maximize the optical invariant which can be passed through a given tube diameter, but it does not eliminate field curvature. It also maximizes weight, since the relay train is almost solid glass, so it is most applicable to small medical instruments.

If larger diameter relays are permissible, it is possible to correct field curvature in each relay stage, making it possible to increase the number of stages without adding field curvature. Baker²⁸ has patented the lens design shown in Fig. 17b for such an application. In this case, field lens and objective are identical, so that an entire relay train can be built using only three different element forms. Pupil and image diameters are the same, and pupil and image are interchangeable.

For purposes of comparison, the two designs shown in Fig. 17 have been scaled to have the same image diameter (2.8 mm) and numerical aperture (0.10), with component focal lengths chosen so that $D_s = D_i$. Minimum tube diameter is 4.0 mm for the rod lens and 5.6 mm for the Baker relay. The image radius of curvature is about 20 mm for the rod relay and about -368 mm for the Baker relay (i.e., field curvature is overcorrected). Image quality for the rod relay is 0.011 waves rms on axis and 0.116 waves rms at full field, both values for best focus, referenced to 587-nm wavelength. For the Baker relay, the corresponding values are 0.025 and 0.056 waves rms, respectively. The Baker design used for this comparison was adapted from the cited patent, with modern glasses substituted for types no longer available. No changes were made to the design other than refocusing it and scaling it to match the first-order parameters of the McKinley design. Neither design necessarily represents the best performance which can be obtained from its design type, and both should be evaluated



(a) McKinley type 1 rod relay



(b) Baker flat field relay

FIGURE 17 Improved unit power afocal relays.

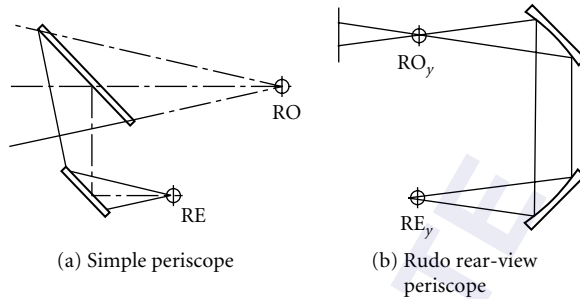


FIGURE 18 Basic reflecting periscopes.

in the context of a complete system design where, for example, the field curvature of the McKinley design may be compensated for by that of the collecting and dispensing objectives. Comparing the individual relay designs does, however, show the price which must be paid for either maximizing the optical invariant within a given tube diameter or minimizing field curvature.

Periscopes

Periscopes are relay trains designed to displace the object space reference point RO a substantial distance away from the eye space reference point RE. This allows the observer to look over an intervening obstacle, or to view objects in a dangerous environment while the observer is in a safer environment. The submarine periscope is the archetypical example. Many other examples can be found in the military⁵ and patent² literature.

The simplest form of periscope is the pair of fold mirrors shown in Fig. 18a, used to allow the viewer to see over nearby obstacles. Figure 18b shows the next higher level of complexity, in the form of a rear-view vehicle periscope patented²⁹ by Rudd.³⁰ This consists of a pair of cylindrical mirrors in a roof arrangement. The cylinders image one axis of object space with the principal purpose of compensating for the image inversion caused by the roof mirror arrangement. This could be considered to be a keplerian anamorphoser, except that it is usually a unit power magnifier, producing no anamorphic compression. Beyond these examples, the complexity of periscopes varies widely.

The optics of complex periscopes such as the submarine periscope can be broken down into a series of component relays. The core of a submarine periscope is a pair of fold prisms arranged like the mirrors in Fig. 18a. The upper prism can be rotated to scan in elevation, while the entire periscope is rotated to scan in azimuth, typically. The main optics is composed of keplerian afocal relays of different magnification, designed to transfer an erect image to the observer inside the submarine, usually at unit net magnification. Galilean and inverse galilean power changers can be inserted between the upper prism and main relay optics to change the field of view. Variants of this arrangement will be found in other military periscopes, along with accessories such as reticles or image intensifiers located at internal foci. Optical design procedures follow those for other keplerian afocal lenses.

18.7 REFLECTING AND CATADIOPTRIC AFOCAL LENSES

Afocal lenses can be designed with powered mirrors or combinations of mirrors and refractors. Several such designs have been developed in recent years for use in the photolithography of microcircuits. All-reflecting afocal lenses are classified here according to the number of powered mirrors they contain. They will be reviewed in order of increasing complexity, followed by a discussion of catadioptric afocal systems.

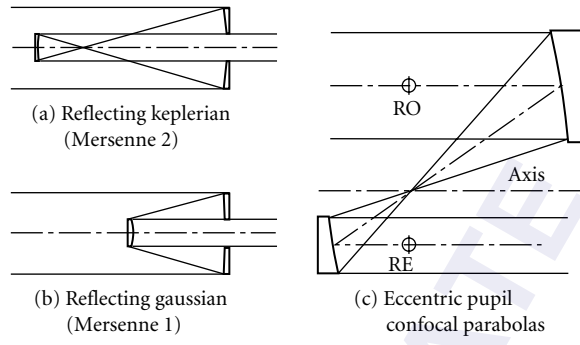


FIGURE 19 Reflecting Galilean.

Two-Powered-Mirror Afocal Lenses

The simplest reflecting afocal lenses are the variants of the galilean and keplerian telescopes shown in Fig. 19a and 19b. They may also be thought of as afocal cassegrainian and gregorian telescopes. The galilean/cassegrainian version is often called a *Mersenne* telescope. In fact, both galilean and keplerian versions were proposed by Mersenne in 1636,³¹ so his name should not be associated solely with the galilean variant.

Making both mirrors parabolic corrects all third-order aberrations except field curvature. This property of *confocal parabolas* has led to their periodic rediscovery,^{32,33} and to subsequent discussions of their merits and shortcomings.³⁴⁻³⁶ The problem with both designs, in the forms shown in Fig. 19a and 19b, is that their eyepieces are buried so deeply inside the design that their usable field of view is negligible. The galilean form is used as a laser beam expander,³⁷ where field of view and pupil location is not a factor, and where elimination of internal foci may be vital.

Eccentric pupil versions of the keplerian form of confocal parabolas, as shown in Fig. 19c, have proven useful as lens attachments.³⁸ RO, RE, and the internal image are all accessible when RO is set one focal length ahead of the primary, as shown. It is then possible to place a field stop at the image and pupil stops at RO and RE, which very effectively blocks stray light from entering the following optics. Being all-reflecting, confocal parabolas can be used at any wavelength, and such attachments have seen use in infrared designs.

Three-Powered-Mirror Afocal Lenses

The principle which results in third-order aberration correction for confocal parabolas also applies when one of the parabolas is replaced by a classical cassegrainian telescope (parabolic primary and hyperbolic secondary), as shown in Fig. 20, with two important added advantages. First, with one negative and two positive mirrors, it is possible to reduce the Petzval sum to zero, or to leave a small

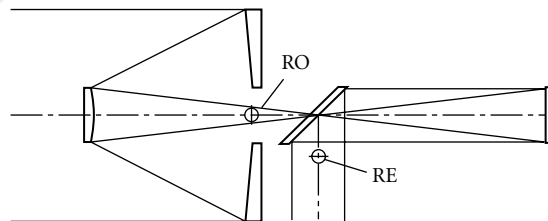


FIGURE 20 Three-powered-mirror afocal lenses.

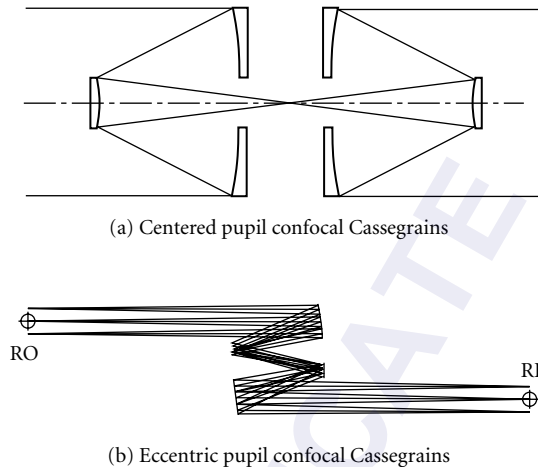


FIGURE 21 Four-powered-mirror afocal lenses.

residual of field curvature to balance higher-order astigmatism. Second, because the cassegrainian is a telephoto lens with a remote front focal point, placing the stop at the cassegrainian primary puts the exit pupil in a more accessible location. This design configuration has been patented by Offner,³⁹ and is more usefully set up as an eccentric pupil design, eliminating the central obstruction and increasing exit pupil accessibility.

Four-Powered-Mirror Afocal Lenses

The confocal parabola principle can be extrapolated one step further by replacing both parabolas with classical cassegrainian telescopes, as shown in Fig. 21a. Each cassegrainian is corrected for field curvature independently, and the image quality of such *confocal cassegrainians* can be quite good. The most useful versions are eccentric pupil. Figure 21b shows an example from Wetherell.⁴⁰ Since both objective and eyepiece are telephoto designs, the separation between entrance pupil RO and exit pupil RE can be quite large. An afocal relay train made up of eccentric pupil confocal cassegrainians will have very long collimated paths. If the vertex curvatures of the primary and secondary mirrors within each cassegrainian are matched, the relay will have zero field curvature, as well. In general, such designs work best at or near unit magnification.

Unit Power Finite Conjugate Afocal Lenses

The simplest catadioptric afocal lens is the cat's-eye retroreflector shown in Fig. 22a, made up of a lens with a mirror in its focal plane. Any ray entering the lens will exit parallel to the incoming ray but traveling in the opposite direction. If made with glass of index of refraction $n = 2.00$, a sphere with one hemisphere reflectorized (Fig. 22b) will act as a perfect retroreflector for collimated light entering the transparent hemisphere. Both designs are, in effect, unit power ($M = -1.00$) afocal lenses. Variations on this technique are used for many retroreflective devices.

Unit power relays are of much interest in photolithography, particularly for microcircuit manufacturing, which requires very high resolution, low focal ratio unit power lenses. In the *Dyson lens*,⁴¹ shown in Fig. 23a, the powered surfaces of the refractor and the reflector are concentric, with radii R and r given by

$$\frac{R}{r} = \frac{n}{(n-1)} \quad (30)$$

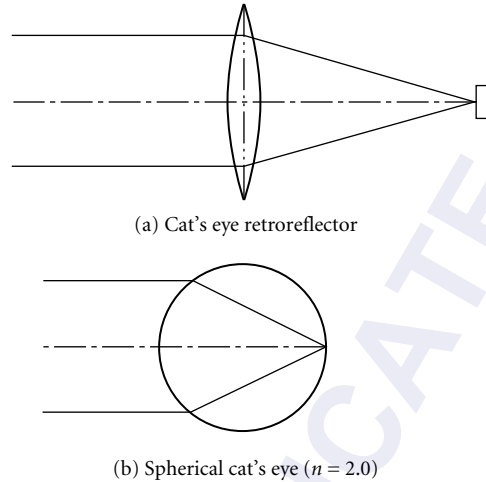


FIGURE 22 Afocal retroreflector designs.

where n is the index of refraction of the glass. At the center point, spherical aberration and coma are completely corrected. In the nominal design, object and image are on the surface intersecting the center of curvature, displaced laterally to separate object from image sensor (this arrangement is termed *eccentric field*, and is common to many multimirror lens systems). In practice, performance of the system is limited by off-axis aberrations, and it is desirable to depart from the nominal design to balance aberrations across the field of view.⁴²

The unit power all-reflecting concentric design shown in Fig. 23b is patented⁴³ by Offner.⁴⁴ It was developed for use in manufacturing microcircuits, and is one of the most successful finite conjugate afocal lens designs in recent years. The spheres are concentric and the plane containing object and image surfaces passes through the common center of curvature. It is an all-reflecting, unit power equivalent of the refracting design shown in Fig. 9. Object and image points are eccentric field, and this is an example of the *ring field* design concept, where axial symmetry ensures good correction throughout a narrow annular region centered on the optical axis. As with the Dyson lens, having an eccentric field means performance is limited by off-axis aberrations. Correction of the design can be improved at the off-axis point by departing from the ideal design to balance on-axis and off-axis aberrations.⁴⁵

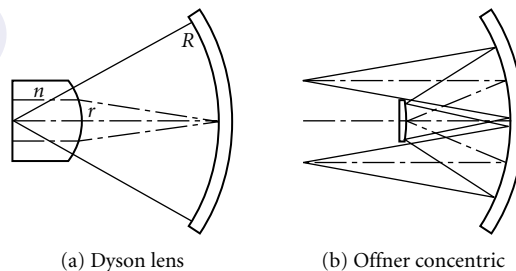


FIGURE 23 Concentric spheres unit power afocal lenses.

18.8 REFERENCES

1. A. van Helden, "The Invention of the Telescope," *Trans. Am. Philos. Soc.* **67**, part 4, 1977.
2. W. B. Wetherell, In "Afocal Lenses," R. R. Shannon and J. C. Wyant (eds.), *Applied Optics and Optical Engineering*, vol. X, Academic Press, New York, 1987, pp. 109–192.
3. E. Hecht and A. Zajac, *Optics*, Addison-Wesley, Reading, Mass., 1974, p. 152.
4. H. C. King, *The History of the Telescope*, Dover, New York, 1979, pp. 44–45.
5. F. B. Patrick, "Military Optical Instruments," In R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. V, Academic Press, New York, 1969, pp. 183–230.
6. D. H. Jacobs, *Fundamentals of Optical Engineering*, McGraw-Hill, New York, 1943.
7. Defense Supply Agency, *Military Standardization Handbook: Optical Design*, MIL-HDBK-141, Defense Supply Agency, Washington, D.C., 1962, section 14, p. 18.
8. A. König, Telescope, U.S. Patent 1,952,795, March 27, 1934.
9. D. B. Rhodes, Scanning Afocal Laser Velocimeter Projection Lens System, U.S. Patent 4,346,990, August 31, 1982.
10. J. C. A. Chastang and R. F. Koerner, Optical System for Oblique Viewing, U.S. Patent 4,428,676, January 31, 1984.
11. A. R. Shulman, *Optical Data Processing*, Wiley, New York, 1970, p. 325.
12. G. F. Marshall, (ed.), *Laser Beam Scanners*, Marcel Dekker, Inc., New York, 1970.
13. J. M. Lloyd, *Thermal Imaging Systems*, Plenum, New York, 1975.
14. M. von Ruhr, Galilean Telescope System, U.S. Patent 962,920, June 28, 1910.
15. H. Erfle, Lens System for Galilean Telescope, U.S. Patent 1,507,111, September 2, 1924.
16. H. Köhler, R. Richter, and H. Kaselitz, Afocal Lens System Attachment for Photographic Objectives, U.S. Patent 2,803,167, August 20, 1957.
17. J. C. J. Blossie, Optical Viewer, U.S. Patent 2,538,077, January 16, 1951.
18. D. L. Wood, View Finder for Cameras, U.S. Patent 2,234,716, March 11, 1941.
19. H. F. Bennett, Lens Attachment, U.S. Patent 2,324,057, July 13, 1943.
20. A. Cox, Anamorphosing Optical System, U.S. Patent 2,720,813, October 18, 1955.
21. T. J. Harris, W. J. Johnson, and I. C. Sandbeck, Wide Angle Lens Attachment, U.S. Patent 2,956,475, October 18, 1960.
22. B. J. Luboshez, Prism Magnification System Having Correction Means for Unilateral Color, U.S. Patent 2,780,141, February 5, 1957.
23. J. H. Hett, "Military Optical Instruments," In R. Kingslake, (ed.), *Applied Optics and Optical Engineering*, vol. V, Academic Press, New York, 1969, pp. 251–280.
24. S. J. Dobson and H. H. Hopkins, "A New Rod Lens Relay System Offering Improved Image Quality," *J. Phys. E: Sci. Instrum.* **22**, 1989, p. 481.
25. H. H. Harris, Optical System Having Cylindrical Rod-like Lenses, U.S. Patent 3,257,902, June 28, 1966.
26. H. R. McKinley, Endoscope Relay Lens, U.S. Patent 5,069,009, October 22, 1991.
27. H. R. McKinley, Endoscope Relay Lens Configuration, U.S. Patent 5,097,359, March 17, 1992.
28. J. G. Baker, Sighting Instrument Optical System, U.S. Patent 2,899,862, August 18, 1959.
29. M. O. Rudd, Two-Mirror System for Periscopic Rearward Viewing, U.S. Patent 4,033,678, July 5, 1977.
30. M. O. Rudd, "Rearview Periscope Comprised of Two Concave Cylindrical Mirrors," *Appl. Opt.* **17**, 1978, pp. 1687–1696.
31. H. C. King, *The History of the Telescope*, Dover, New York, 1979, pp. 48–49.
32. S. C. B. Gascoigne, "Recent Advances in Astronomical Optics," *Appl. Opt.* **12**, 1973, p. 1419.
33. S. Rosin, M. Amon, and P. Laakman, "Afocal Parabolic Reflectors," *Appl. Opt.* **13**, 1974, p. 741.
34. R. B. Annable, "Still More on Afocal Parabolic Reflectors," *Appl. Opt.* **13**, 1974, p. 2191.
35. S. Rosin, "Still More on Afocal Parabolic Reflectors," *Appl. Opt.* **13**, 1974, p. 2192.
36. W. B. Wetherell and M. P. Rimmer, "Confocal Parabolas: Some Comments," *Appl. Opt.* **13**, 1974, p. 2192.

37. W. K. Pratt, *Laser Communication Systems*, Wiley, New York, 1969, p. 53.
38. I. R. Abel, Wide Field Reflective Optical Apparatus, U.S. Patent 3,811,749, May 21, 1974.
39. A. Offner, Catoptric Anastigmatic Afocal Optical System, U.S. Patent 3,674,334, July 4, 1972.
40. W. B. Wetherell, "All-Reflecting Afocal Telescopes," In "Reflective Optics," *SPIE*, vol. 751, 1987, pp. 126–134.
41. J. Dyson, "Unit Magnification Optical System Without Seidel Aberrations," *J. Op. Soc. Am.* **49**, 1959, pp. 713–716.
42. R. M. H. New, G. Owen, and R. F. W. Pease, "Analytical Optimization of Dyson Optics," *Appl. Opt.* **31**, 1992, pp. 1444–1449.
43. A. Offner, Unit Power Imaging Cataptric Anastigmat, U.S. Patent 3,748,015, July 24, 1973.
44. A. Offner, "New Concepts in Projection Optics," *Opt. Eng.* **14**, 1975, p. 130.
45. C. S. Ih and K. Yen, "Imaging and Fourier Transform Properties of an Allspherical Mirror System," *Appl. Opt.* **19**, 1980, p. 4196.

William L. Wolfe

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

19.1 GLOSSARY

δ	angular deviation
ϕ	phase
ω	radian frequency of rotation
A, B, C, D, d	prism dimensions
t	time
x, y	rectangular components
α	angle
1, 2	prism number

19.2 INTRODUCTION

Prisms of various shapes and sizes are used for folding, inverting, reverting, displacing, and deviating a beam of light, whether it be collimated, converging, or diverging.

Prisms, rather than mirrors, are often used for the applications discussed here, since they make use of reflecting coatings at what amounts to an interior surface. The coatings can be protected on their backs by other means, and do not tarnish with age and exposure. Even better, some prisms do not need such coatings if the (internal) angle of incidence exceeds the critical angle.

In these applications, chromatism is to be avoided. Thus, the arrangements either make use of perpendicular incidence or compensating angles of incidence.

Almost all of these prisms are meant to be used with collimated beams. Most of the operations are somewhat equivalent to the use of plane parallel plates, which displace but do not deviate a

collimated beam. However, such plates have both chromatism and spherical aberration in a convergent beam.

Dispersing prisms are discussed in Chap. 20 by George J. Zissis and polarizing prisms in Chap. 13 by Jean M. Bennett, both in this volume of the *Handbook*.

19.3 INVERSION, REVERSION

A reverted image shifts the image left for right. An inverted image is upside down. A reinverted image or inverted-reverted image does both.

The best two references on this subject are the Frankford Arsenal book called *Design of Fire Control Optics*,¹ and Jacobs' book called *Optical Engineering*.² Many of the diagrams shown here have been taken from the former since they provide direct design information as well as descriptions of the performance of the prism.

19.4 DEVIATION, DISPLACEMENT

The beam, in addition to being inverted and/or reverted, can also be displaced and/or deviated. Displacement means that the beam has been translated in x or y , but it has not changed the direction in which it was traveling. Deviation indicates that the beam has been caused to change its direction. Deviation is measured in angular measure; displacement in linear measure. If a beam has been deviated, displacement is not important. If a beam has been displaced, it usually has not been deviated. Although the two can occur simultaneously, it seldom happens in optical prisms.

19.5 SUMMARY OF PRISM PROPERTIES

Table 1 is a listing of the prisms that are described in this section. The first column provides the name of the prism. The second column indicates whether the image is reverted. The third column indicates whether the image has been inverted. The fourth column indicates how much the image has been displaced as a ratio of the critical prism dimension A . The next column gives the amount of deviation, and the last column provides comments.

19.6 PRISM DESCRIPTIONS

Each diagram shows at least one view of the prism and a set of dimensions. The A dimension is a reference dimension. It is always 1.00, and the rest of the dimensions are related to it. The refractive index is almost always taken as 1.5170, a representative value for glass in the visible. Prism dimensions can change somewhat if the refractive index is modestly different from the chosen value. However, if a material like germanium is used, the prism might be drastically different.

TABLE 1 Summary of Prism Properties

Prism	Reverts	Inverts	Displaces	Deviates	Comments
Right angle	Yes	No	—	90	Simplest
	No	Yes			
Porro A	Yes	Yes	1.1A	0	Binocs
Porro B	Yes	Yes	1.1A	0	Binocs
Abbe A	Yes	Yes	0	0	
Abbe B	Yes	Yes	0	0	
Dove, double dove	Yes	No	0	0	Parallel light
Pechan	Yes	No	0	0	Nonparallel
Amici	Yes	No	—	90	Roof
Schmidt	Yes	Yes	—	45	
Leman	Yes	Yes	3A	0	
Penta	No	No	—	90	exactly
Reversion	Yes	Yes	0	0	Nonparallel
Wollaston	No	No	—	90	tracing
Zeiss	Yes	Yes	Design		
Goerz	Yes	Yes	Design		
Frank 1	Yes	Yes	—	115	
Frank 2	Yes	Yes	—	60	
Frank 3	Yes	Yes	—	45v 90h	
Frank 4	Yes	No	—	45v 90h	
Frank 5	No	Yes	—	60v 90h	
Frank 6	Yes	Yes	—	60v 90h	
Frank 7	No	No	—	45v 90h	
Hastings	Yes	Yes	0	0	
Rhomboid	No	No	A	0	
Risleys	No	No	No	0–180	
Retro	Yes	Yes	No	180	
D40	No	No	—	40–50	
D60	No	No	—	50–60	
D90	No	No	—	80–100	
D120	No	No	—	110–130	

Right-Angle Prism

Perhaps the simplest of the deviating prisms is the right-angle prism. Light enters one of the two perpendicular faces, as shown in Fig. 1a, reflects off the diagonal face, and emerges at 90° from the other perpendicular face. The beam has been rotated by 90°, and the image has been inverted. If the prism is used in the other orientation, shown in Fig. 1b, then the image is reverted. The internal angle of incidence is 45°, which is sufficient for total internal reflection (as long as the refractive index is greater than 1.42).

Porro Prism

A Porro prism, shown in Fig. 2, has a double reflection and may be considered to be two right-angle prisms together. They are often identical. Often, two Porro prisms are used together to invert and revert the image. The incidence angles are the same as with the right-angle prism, so that total internal reflection takes place with refractive indices larger than 1.42. It is a direct-vision prism.

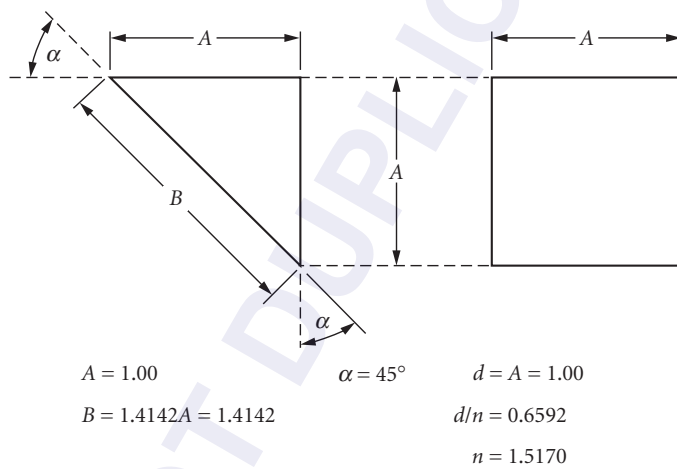
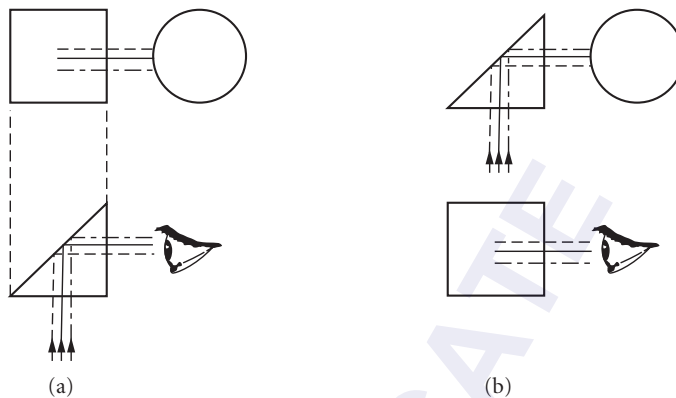
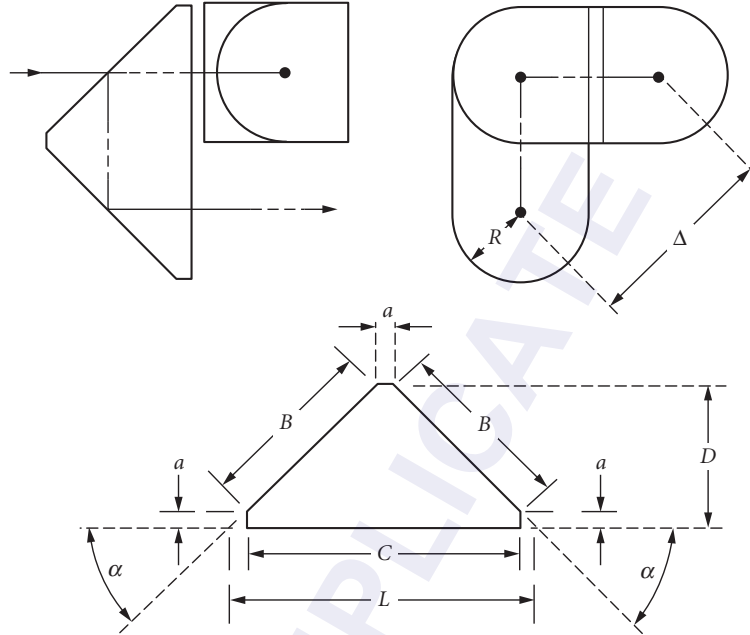


FIGURE 1 Right-angle prism.



$A = 1.00$	$\alpha = 45^\circ$ (These values are given)	$a = 0.10$ (chosen arbitrarily)
$B = 1.4142A = 1.4142$	$\Delta = 1.4142(A + a) = 1.5556$	$d = 2(2A + 3a) = 4.60$
$C = 2A + a = 2.1$		$d/n = 3.0324$
$D = A + a = 1.1$		$n = 1.5170$
$L = 2A + 3a = 2.30$		
$R = A/2 = 0.50$		

FIGURE 2 Porro prism.

Abbe's version of the Porro prism is shown in Fig. 3. The resultant beam is inverted and reverted and is directed parallel and in the same direction as incident beam.

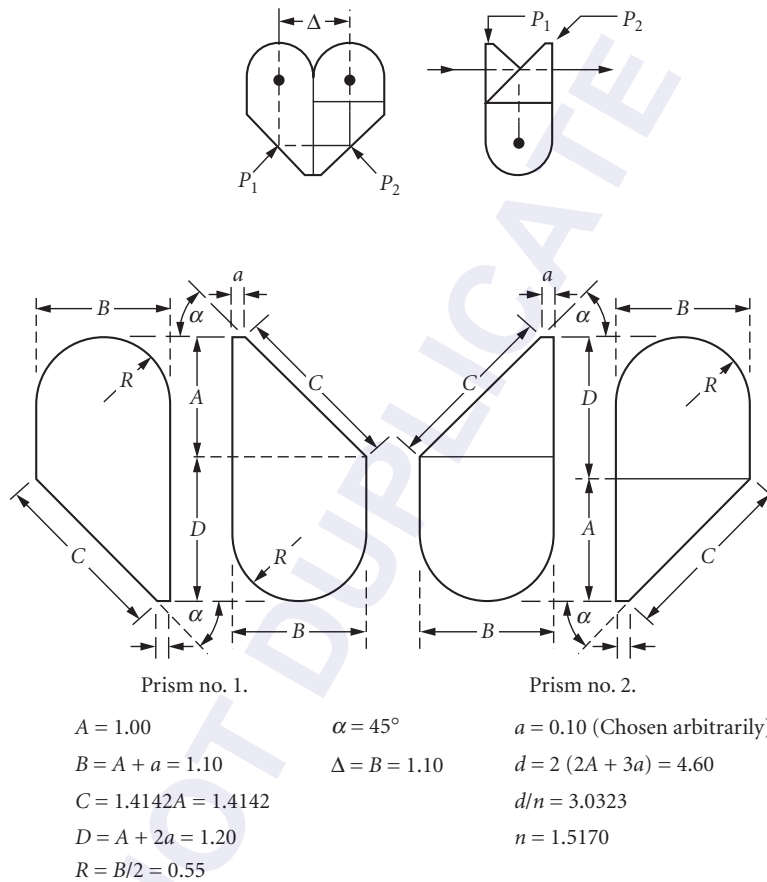


FIGURE 3 Abbe modification of Porro prisms for binoculars.

Abbe's Prisms

Two versions of prisms invented by Abbe are shown. They both are direct-vision prisms that revert and invert the image. One version is symmetrical; the other uses three different prism segments. They are shown in Figs. 4 and 5.

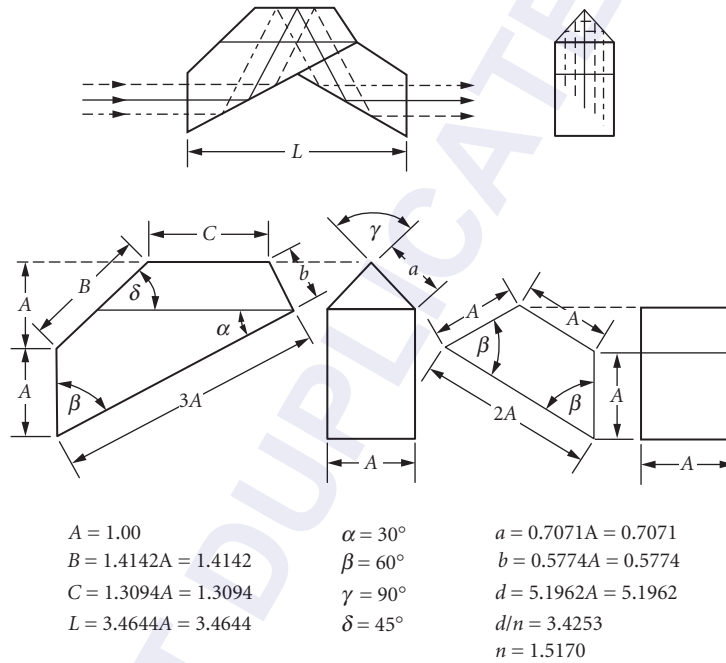


FIGURE 4 Abbe direct-vision prism—A.

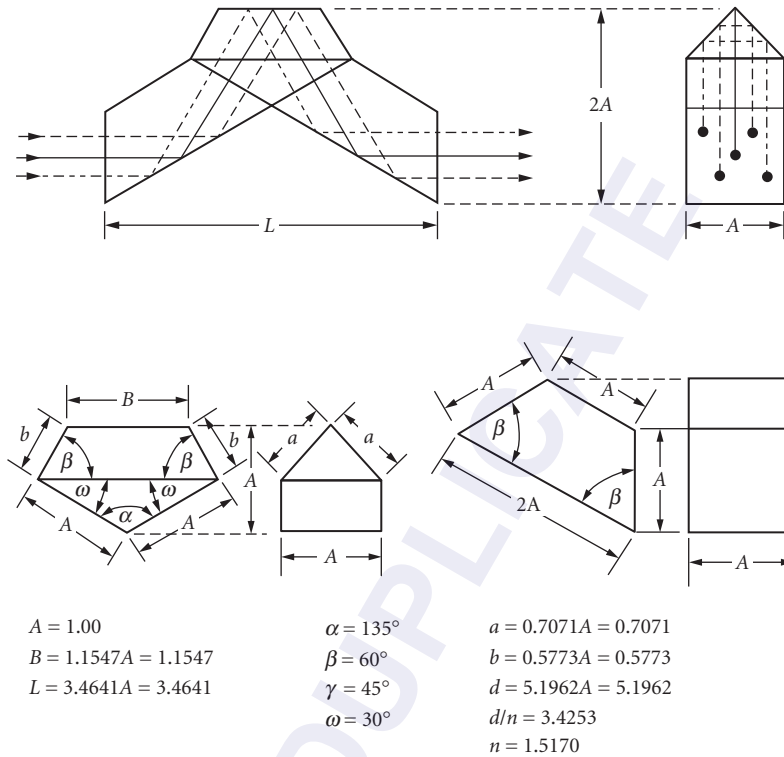
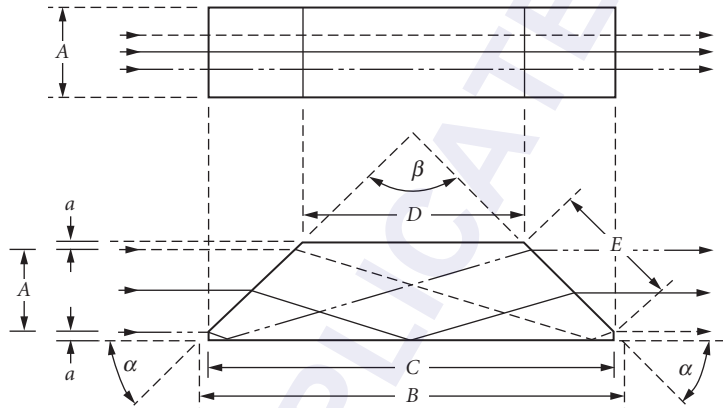


FIGURE 5 Abbe direct-vision prism—B.

Dove Prism

A Dove prism (also known as a Harting-Dove prism) does not deviate or displace an image but it can be used to either invert or revert an image. It must be placed in parallel light. Such a prism is shown in Fig. 6.

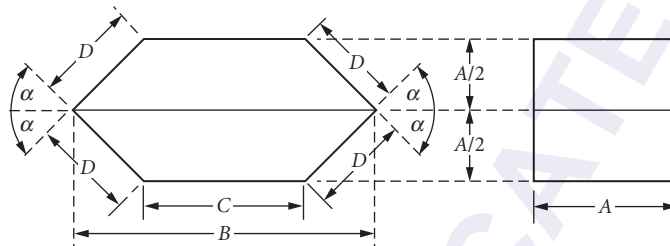


$$\begin{aligned}
 A &= 1.00 & \alpha &= 45^\circ & a &= 0.05 \\
 B &= (A + 2a) \left\{ \frac{\sqrt{n^2 - \sin^2 \alpha} + \sin \alpha}{\sqrt{n^2 - \sin^2 \alpha} - \sin \alpha} + 1 \right\} & \beta &= 90^\circ & d &= \frac{n(A + 2a)}{\sin \alpha \left\{ \sqrt{n^2 - \sin^2 \alpha} - \sin \alpha \right\}} \\
 &= 4.2271(A + 2a) = 4.6498 & & & &= 3.3787(A + 2a) = 3.7165 \\
 C &= B - 2a = 4.5498 & & & d/n &= 2.4499 \\
 D &= B - 2(A + 2a) = 2.4498 & & & n &= 1.5170 \\
 E &= \frac{a + A}{\cos \alpha} = 1.4142(A + a) = 1.4849
 \end{aligned}$$

FIGURE 6 Harting-Dove prism.

Double Dove

Two Dove prisms are glued together. The length is halved, but the height is doubled. It performs the same functions as a single Dove in almost the same way. It is shown in Fig. 7.



$$A = 1.00$$

$$\alpha = 45^\circ$$

$$d = \frac{nA}{2 \sin \alpha \{ \sqrt{n^2 - \sin^2 \alpha} - \sin \alpha \}}$$

$$B = \frac{A}{2} \left\{ \frac{\sqrt{n^2 - \sin^2 \alpha} + \sin \alpha}{\sqrt{n^2 - \sin^2 \alpha} - \sin \alpha} + 1 \right\}$$

$$= nAC = 1.6893$$

$$= 2.1136A = 2.1136$$

$$d/n = 1.1135$$

$$n = 1.5170$$

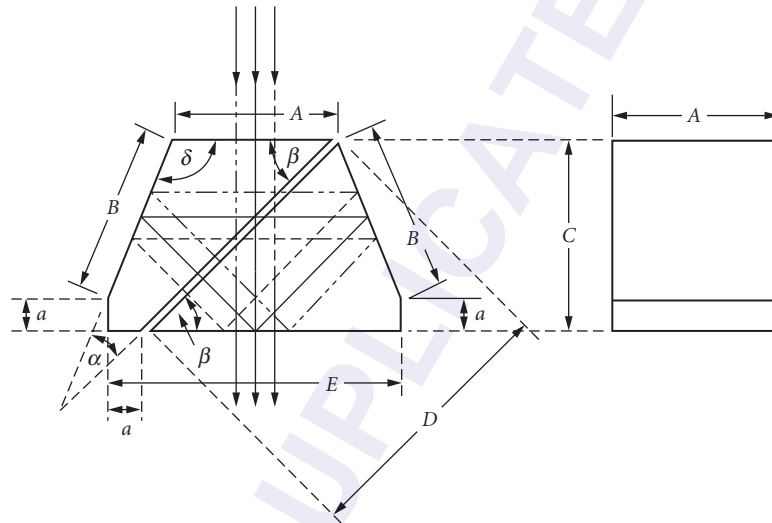
$$C = B - A = 1.1136$$

$$D = \frac{A}{2 \cos \alpha} = 0.7071A = 0.7071$$

FIGURE 7 Double Dove prism.

Pechan Prism

The Pechan prism (shown in Fig. 8) performs the same function as the Dove, but it can do it in converging or diverging beams. The surfaces marked *B* are silvered and protected. The surfaces bordering the air space are unsilvered.



$A = 1.00$	$\alpha = 22^\circ 30'$	$a = 0.2071A = 0.2071$
$B = 1.0824A = 1.0824$	$\beta = 45^\circ$	$d = 4.6213A = 4.6213$
$C = 1.2071A = 1.2071$	$\gamma = 67^\circ 30'$	$d/n = 3.0464$
$D = 1.7071A = 1.7071$	$\delta = 112^\circ 30'$	$n = 1.5170$
$E = 1.8284A = 1.8284$		

FIGURE 8 Pechan prism.

Amici (Roof) Prism

This more complex arrangement of surfaces inverts the image, reverts it, and deviates it 90° . It is shown in Fig. 9. Since this prism makes use of the roof effect, it has the same angles as both the right-angle and Porro prisms, and exhibits total internal reflection for refractive indices larger than 1.42.

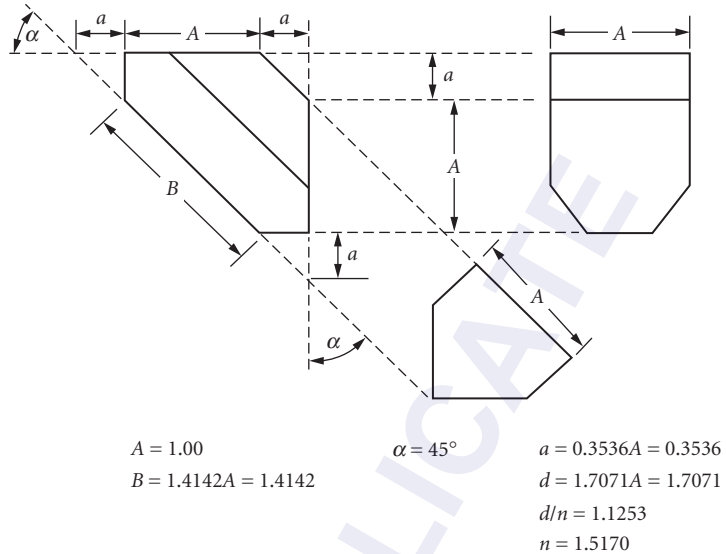


FIGURE 9 Amici (roof) prism.

Schmidt Prism

The prism will invert and revert the image, and it will deviate it through 45° . It is shown in Fig. 10.

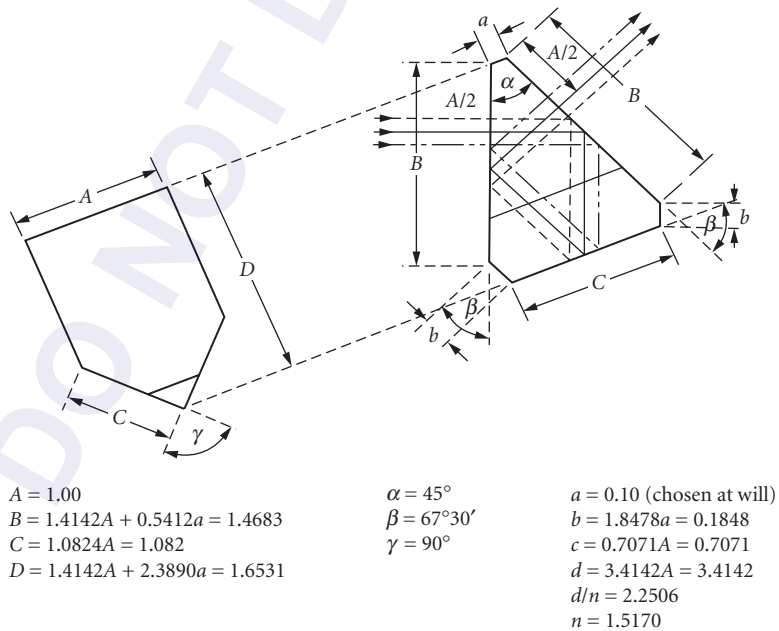


FIGURE 10 Schmidt prism.

Leman Prism

This rather strange looking device, shown in Fig. 11, reverts, inverts, and displaces by $3A$, an image.

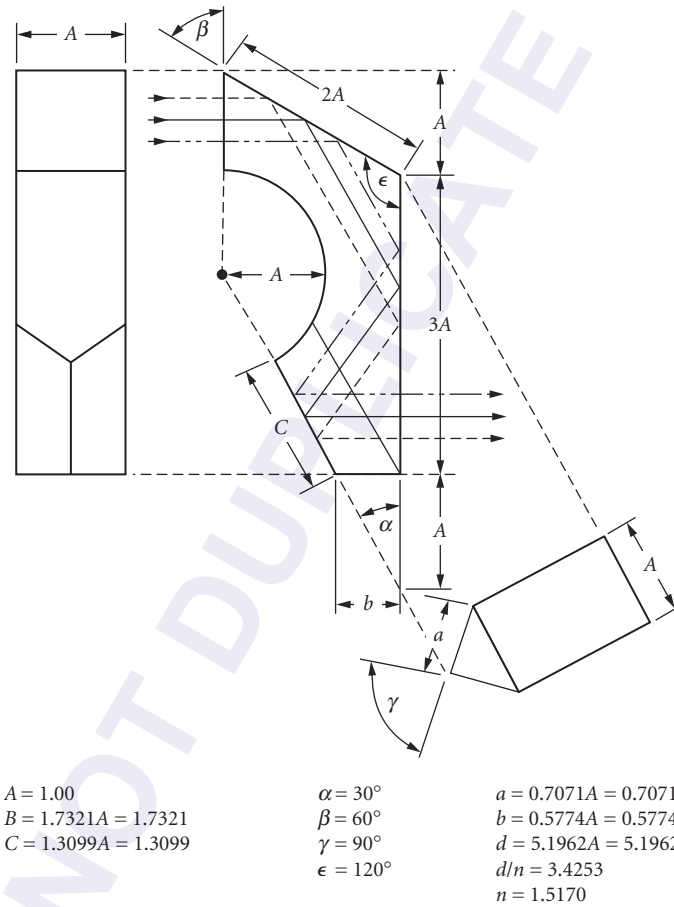
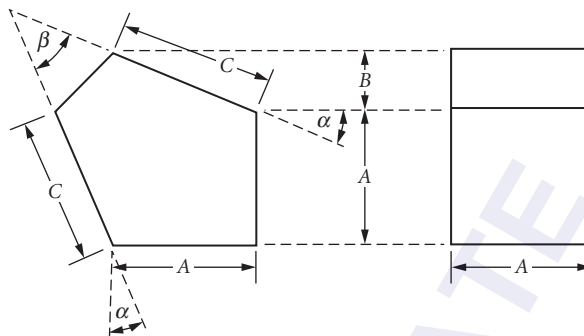


FIGURE 11 Leman prism.

Penta Prism

A penta prism has the remarkable property that it always deviates a beam by exactly 90° in the principal plane. This is akin to the operation of a cube corner. The two reflecting surfaces of the penta prism, shown in Fig. 12, must be reflectorized, as the angles are 22.5° and therefore require a refractive index of 2.62 or greater for total internal reflection. Some penta-prisms are equipped with a roof to revert the image.

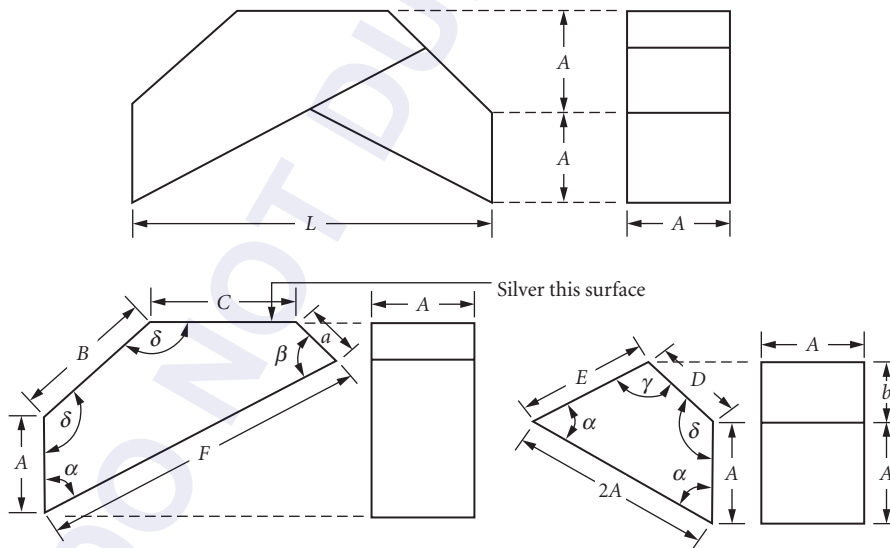


$A = 1.00$	$\alpha = 22^\circ 30'$	$a = 3.4142A = 3.4142$
$B = 0.4142A = 0.4142$	$\beta = 45^\circ$	$d/n = 2.2506$
$C = 1.0824A = 1.0824$		$n = 1.5170$

FIGURE 12 Penta prism.

Reversion Prism

This prism operates like an Abbe prism, type A, but does not require parallel light. It is shown in Fig. 13.

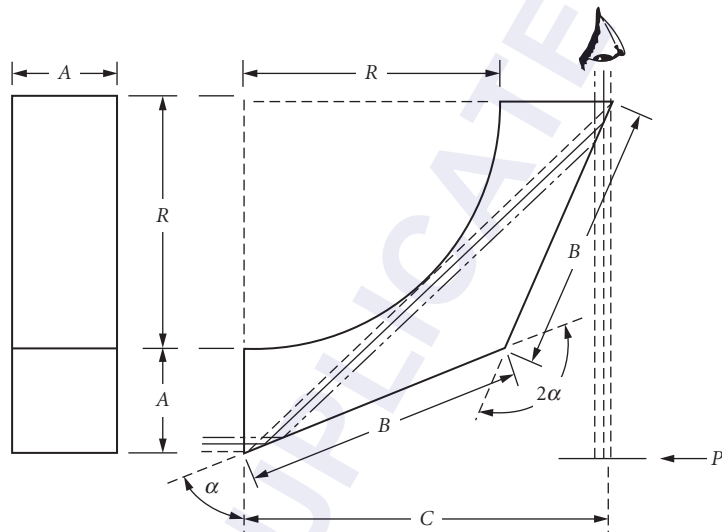


$A = 1.00$	$\alpha = 60^\circ$	$a = 0.5176A$
$B = 1.4142A$	$\beta = 75^\circ$	$b = 0.6340A$
$C = 1.4641A$	$\gamma = 105^\circ$	$d = 5.1962A$
$D = 0.8966A$	$\delta = 135^\circ$	$d/n = 3.4253A$
$E = 1.2679A$		$n = 1.5170$
$F = 3.2679A$		
$L = 3.4641A$		

FIGURE 13 Reversion prism.

Wollaston Prism

This prism does not invert, revert, or displace. It does deviate a beam by 90°, allowing a tracing to be made. It is shown in Fig. 14.



$A = 1.00$	$\alpha = 67^{\circ}30'$	$d = 2R = 4.8284A = 4.8284$
$B = 2.6131A = 2.6131$		$d/n = 3.1829$
$C = 3.4142A = 3.4142$		$n = 1.5170$
$R = 2.4142A = 2.4142$		

FIGURE 14 Wollaston prism.

Carl Zeiss Prism System

This arrangement of three prisms allows the image to be reverted, inverted, and displaced, but not deviated. The amount of displacement is adjustable. The system is shown in Fig. 15.

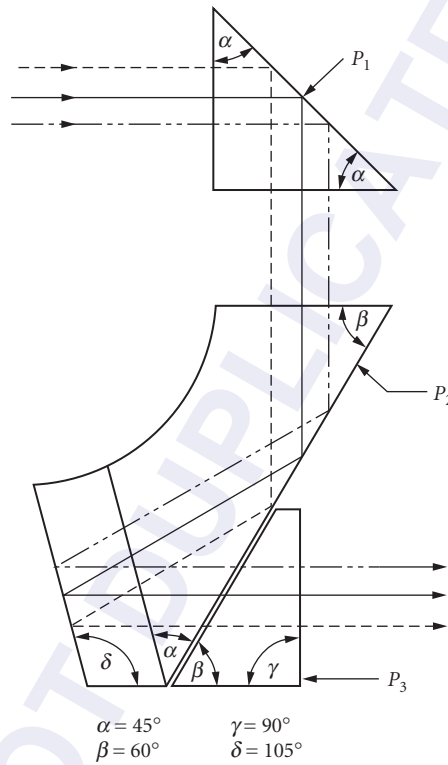


FIGURE 15 Carl Zeiss prism system.

Goerz Prism System

This is an alternate to the Zeiss system. It does the same things. It is shown in Fig. 16.

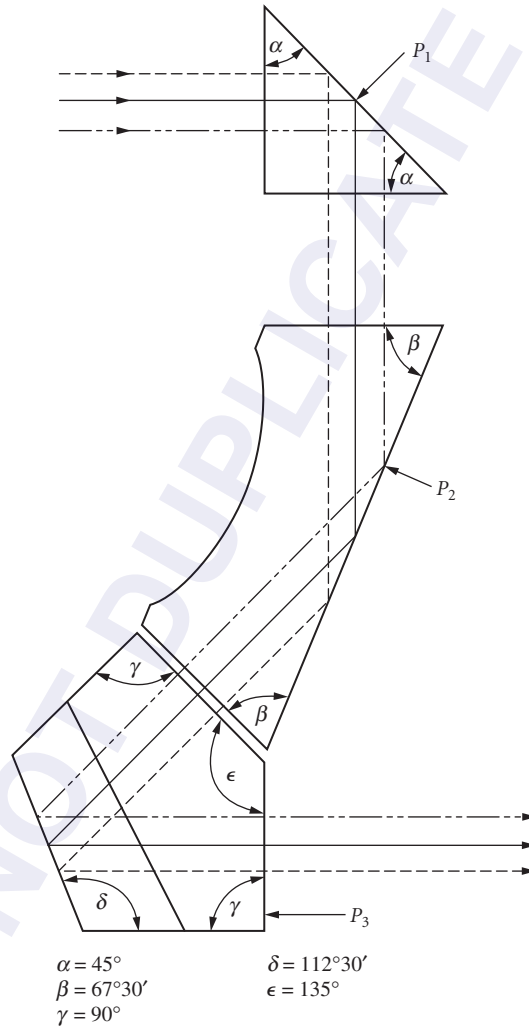
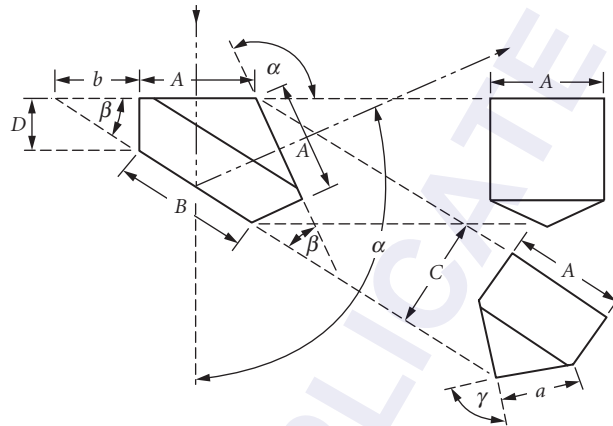


FIGURE 16 C. P. Goerz prism system.

Frankford Arsenal 1

This prism, shown in Fig. 17, reverts, inverts, and deviates through 115° .

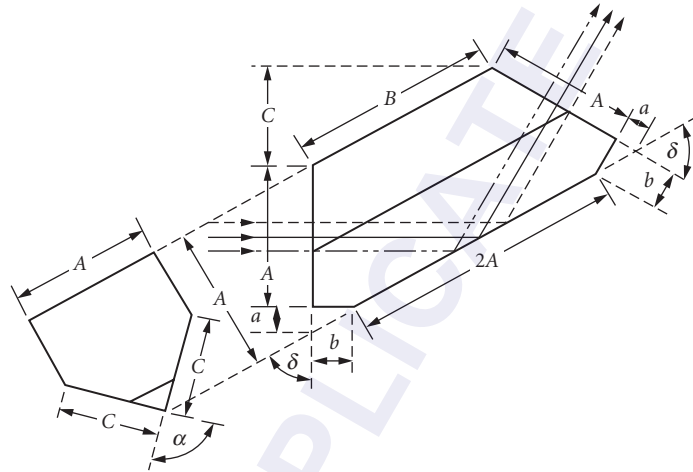


$A = 1.00$	$\alpha = 115^\circ$	$a = 0.7071A = 0.7071$
$B = 1.1857A = 1.1857$	$\beta = 32^\circ 30'$	$b = 0.7320A = 0.7320$
$C = 0.9306A = 0.9306$	$\gamma = 90^\circ$	$d = 1.5697A = 1.5697$
$D = 0.4613A = 0.4613$		$d/m = 1.0347$
		$n = 1.5170$

FIGURE 17 Frankford Arsenal prism 1.

Frankford Arsenal 2

This prism reverts, inverts, and deviates through 60° . It is shown in Fig. 18.

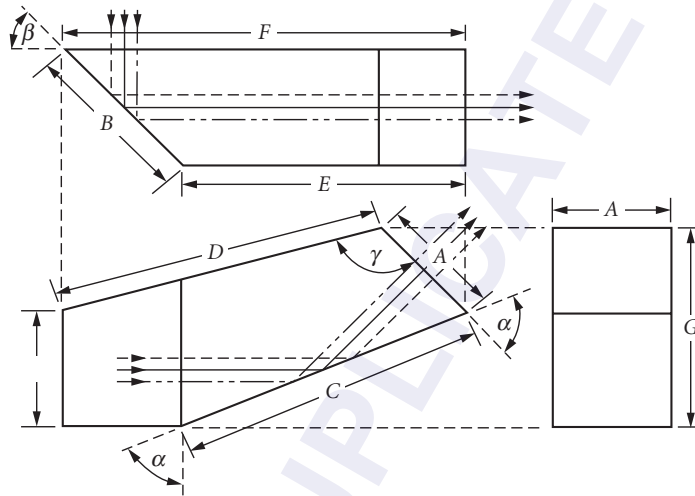


$A = 1.00$	$\alpha = 90^\circ$	$a = 0.1547A = 0.1547$
$B = 1.4641A = 1.4641$	$\delta = 60^\circ$	$b = 0.2680A = 0.2680$
$C = 0.7321A = 0.7321$		$d = 2.2690A = 2.2680$
		$d/n = 1.4951$
		$n = 1.5170$

FIGURE 18 Frankford Arsenal prism 2.

Frankford Arsenal 3

This prism reverts, inverts, and deviates through an angle of 45° upward and 90° horizontally. It is shown in Fig. 19.

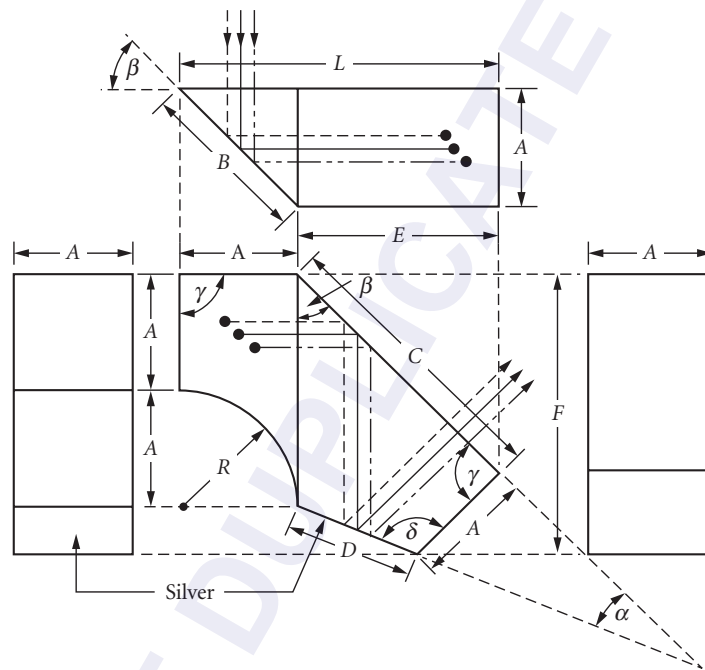


$A = 1.00$	$\alpha = 67^\circ 30'$	$d/n = 2.2506$
$B = 1.4142A = 1.4142$	$\beta = 45^\circ$	$d = 3.4142A = 3.4142$
$C = 2.6131A = 2.6131$	$\gamma = 120^\circ 21' 40''$	$n = 1.5170$
$D = 2.7979A = 2.7979$		
$E = 2.4142A = 2.4142$		
$F = 3.4142A = 3.4142$		
$G = 1.7071A = 1.7071$		

FIGURE 19 Frankford Arsenal prism 3.

Frankford Arsenal 4

This prism reverts the image and deviates it 45° upward and 90° sideways, like Frankford Arsenal 3. It is shown in Fig. 20.

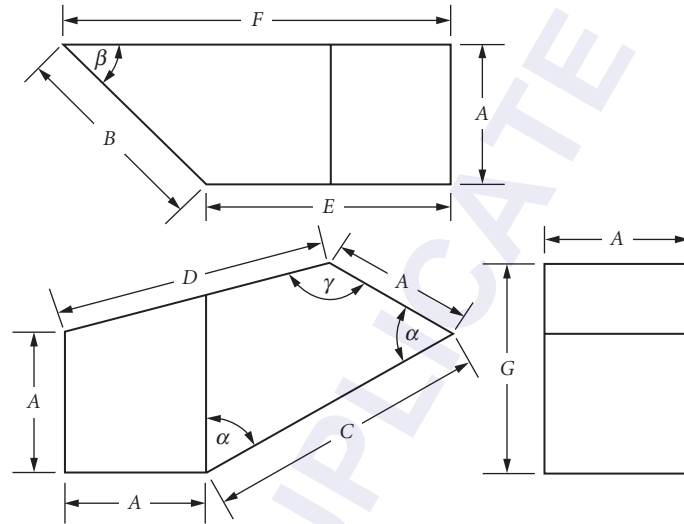


$A = 1.00$	$\alpha = 22^{\circ}30'$	$d = 4.4142A = 4.4142$
$B = 1.4142A = 1.4142$	$\beta = 45^{\circ}$	$d/n = 2.9098$
$C = 2.4142A = 2.4142$	$\gamma = 90^{\circ}$	$n = 1.5170$
$D = 1.0824A = 1.0824$	$\delta = 112^{\circ}30'$	
$E = 1.7071A = 1.7071$		
$F = 2.4142A = 2.4142$		
$L = 2.7071A = 2.7071$		
$R = A = 1.00$		

FIGURE 20 Frankford Arsenal prism 4.

Frankford Arsenal 5

This prism inverts the image while deviating it 90° sideways and 60° upward. It is shown in Fig. 21.



$A = 1.00$

$B = 1.4142A = 1.4142$

$C = 2.000A = 2.000$

$D = 1.9318A = 1.9318$

$E = 1.7321A = 1.7321$

$F = 2.7321A = 2.7321$

$G = 1.500A = 1.500$

$\alpha = 60^\circ$

$\beta = 45^\circ$

$\gamma = 135^\circ$

$d = 2.7437A = 2.7431$

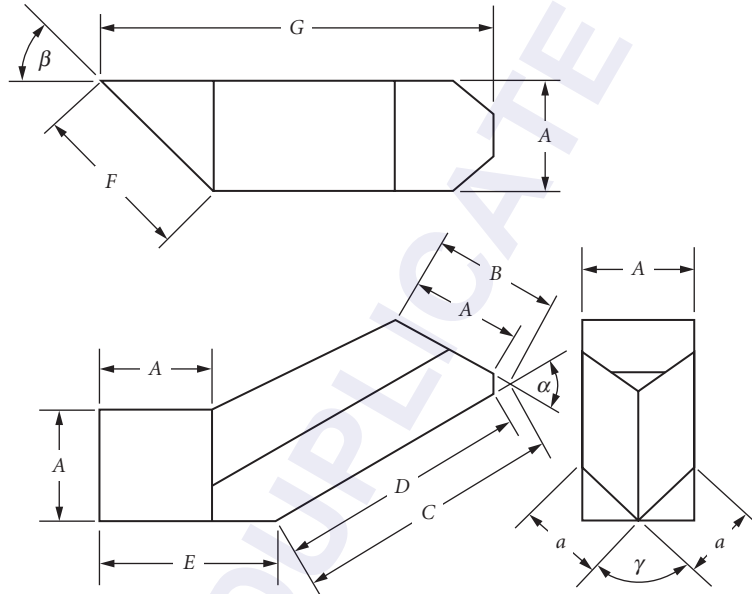
$d/n = 1.8086$

$n = 1.5170$

FIGURE 21 Frankford Arsenal prism 5.

Frankford Arsenal 6

This prism inverts, reverts, and deviates 90° horizontally and 60° vertically. It is shown in Fig. 22.

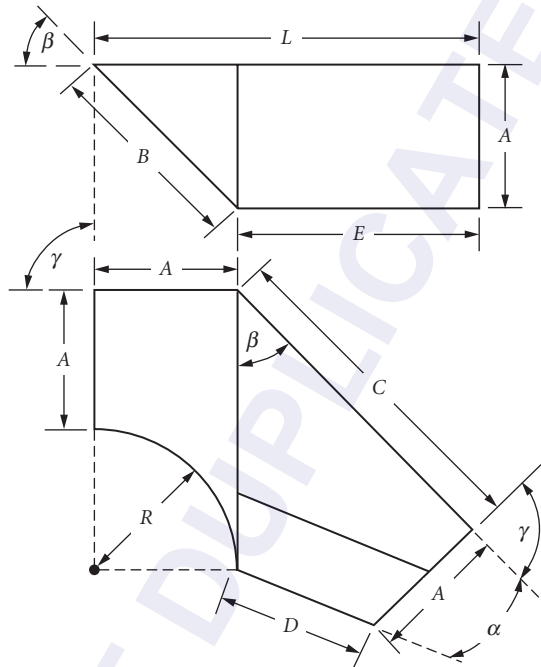


- | | | |
|------------------------|---------------------|------------------------|
| $A = 1.00$ | $\alpha = 60^\circ$ | $a = 0.7071A = 0.7071$ |
| $B = 1.2071A = 1.2071$ | $\beta = 45^\circ$ | $d = 3.6681A = 3.6681$ |
| $C = 2.4142A = 2.4142$ | $\gamma = 90^\circ$ | $d/n = 2.4180$ |
| $D = 2.2071A = 2.2071$ | | $n = 1.5170$ |
| $E = 1.5774A = 1.5774$ | | |
| $F = 1.4142A = 1.4142$ | | |
| $G = 3.4888A = 3.4888$ | | |
| $H = 1.8107A = 1.8187$ | | |

FIGURE 22 Frankford Arsenal prism 6.

Frankford Arsenal 7

This prism neither reverts nor inverts, but deviates 90° horizontally and 45° vertically. It is shown in Fig. 23.



$A = 1.00$	$\alpha = 22^{\circ}30'$	$d = 4.4142A = 4.4142$
$B = 1.4142A = 1.4142$	$\beta = 45^{\circ}$	$d/n = 2.9098$
$C = 2.4142A = 2.4142$	$\gamma = 90^{\circ}$	$n = 1.5170$
$D = 1.0824A = 1.0824$		
$E = 1.7071A = 1.7071$		
$L = 2.7071A = 2.7071$		
$R = A = 1.00$		

FIGURE 23 Frankford Arsenal prism 7.

Brashear-Hastings Prism

This device, shown in Fig. 24, inverts an image without changing the direction of the beam. Since this is a relatively complicated optical element, it does not see much use.

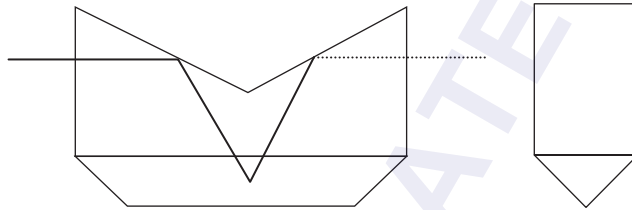


FIGURE 24 Brashear-Hastings prism.

Rhomboidal Prism

A rhomboidal prism, as shown in Fig. 25, displaces the beam without inverting, reverting, deviating, or otherwise changing things. The reflecting analog is a pair of mirrors at 45° .

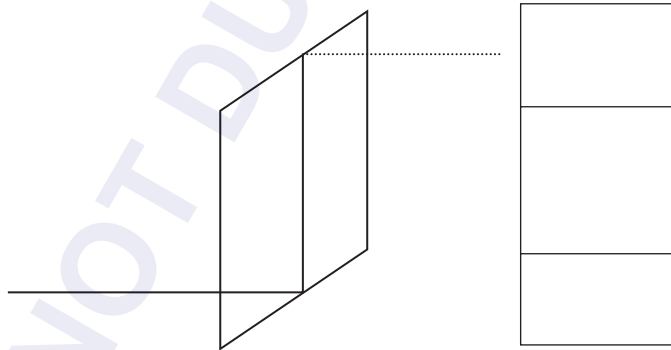


FIGURE 25 Rhomboidal prism.

Risley Prisms

Risley prisms are used in two ways. If they are slightly absorbing, they can be used as variable attenuators by translating one with respect to the other perpendicular to their apices.³ They can also be rotated to generate a variety of angular deviations.⁴ A single prism deviates the beam according to its wedge angle and refractive index. If rotated in a circle about an axis perpendicular to its face, it will rotate the beam in a similar circle. A second, identical prism in series with it, as shown in Fig. 26, can double the angle of the beam rotation and generate a circle of twice the radius. If they rotate in opposite directions, one motion is canceled and a line is generated. In fact, all sorts of Lissajous-type figures can be obtained; some are shown in Fig. 27. The equations that govern the patterns are

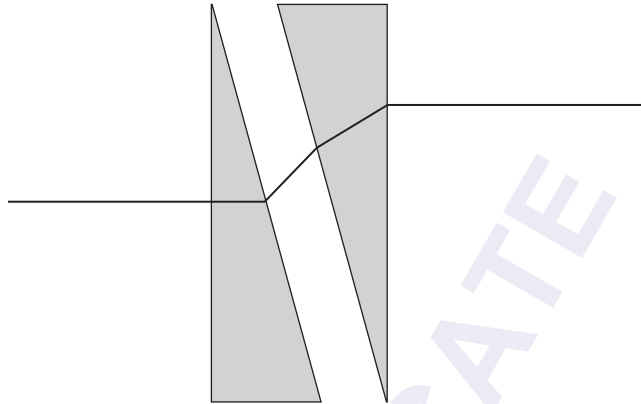


FIGURE 26 Risley prisms.

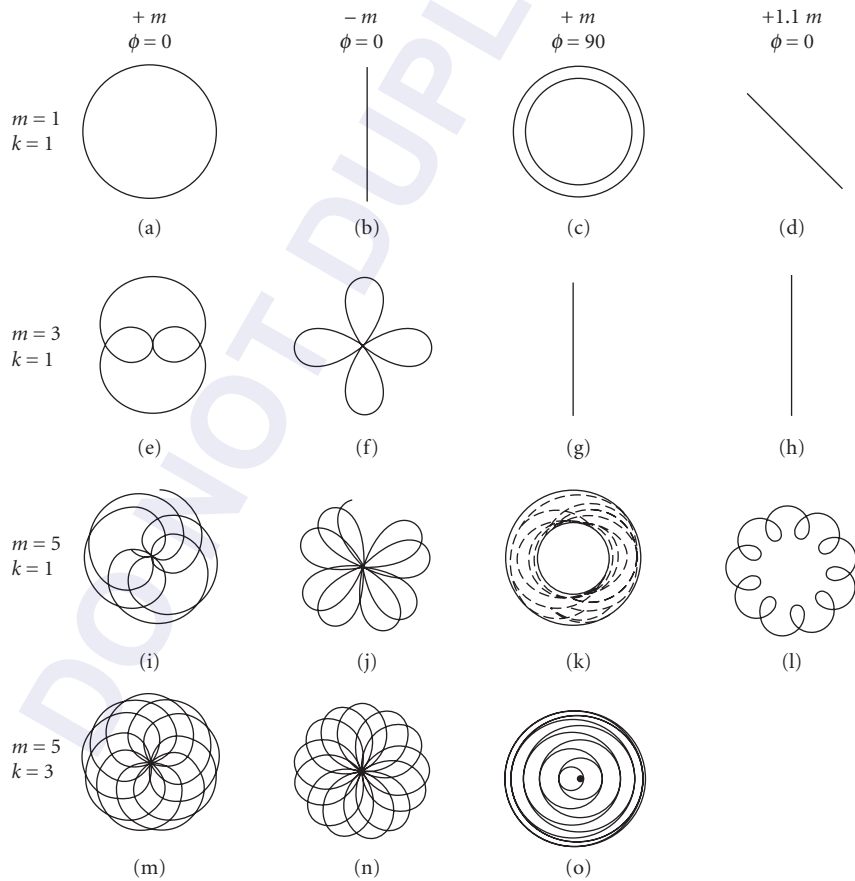


FIGURE 27 Risley prism patterns.

$$\delta_x = \delta_1 \cos \omega_1 t + \delta_2 \cos(\omega_2 t + \phi) \tag{1}$$

$$\delta_y = \delta_1 \sin \omega_1 t + \delta_2 \sin(\omega_2 t + \phi) \tag{2}$$

where δ_x and δ_y are the beam deviations, δ_1 and δ_2 are the individual prism deviations, ω is the rotation rate, t is time, and ϕ is the phase of the prism position. For relatively monochromatic applications, the prisms can be “fresnelled,” as shown in Fig. 28, and the mirror analogs, shown in Fig. 29, can also be used.



FIGURE 28 Fresnel Risleys.

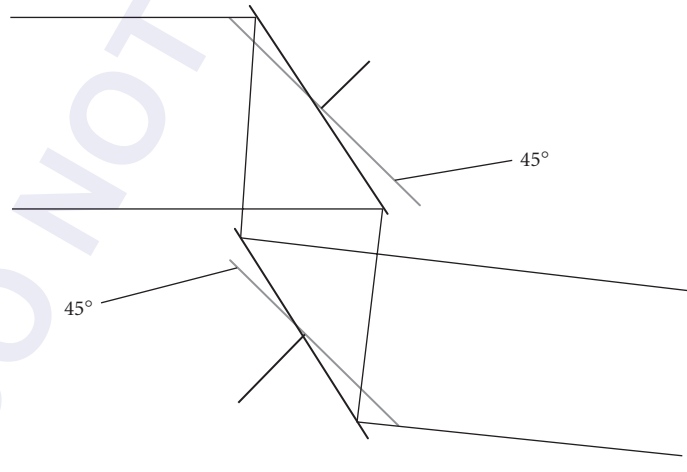


FIGURE 29 Risley mirrors.

Retroreflectors

The familiar reflective cube corner (not corner cube), that sends a ray back in the direction from which it came, has its refractive analog, as shown in Fig. 30. The angles are adjusted so that total internal reflection occurs. The angular acceptance range can be large.

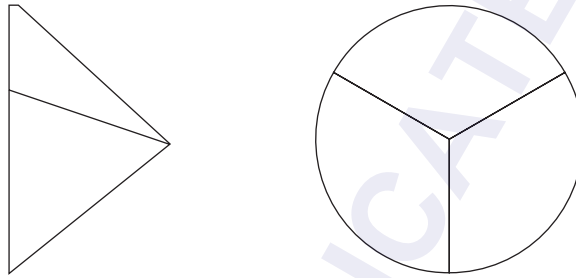


FIGURE 30 Retroreflectors.

General Deviation Prisms

Figure 31 shows a 40° -deviation prism. Other angles are obtainable with appropriate changes in the prism manufacture, as shown, for example, in Figs. 32 and 33.

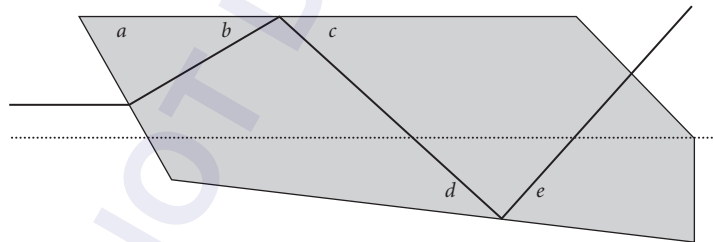


FIGURE 31 40° -deviation prism—D40.

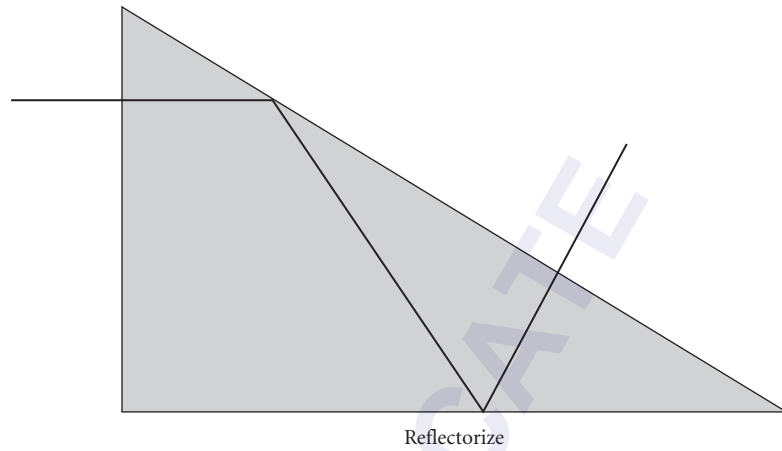


FIGURE 32 60°-deviation prism—D60.

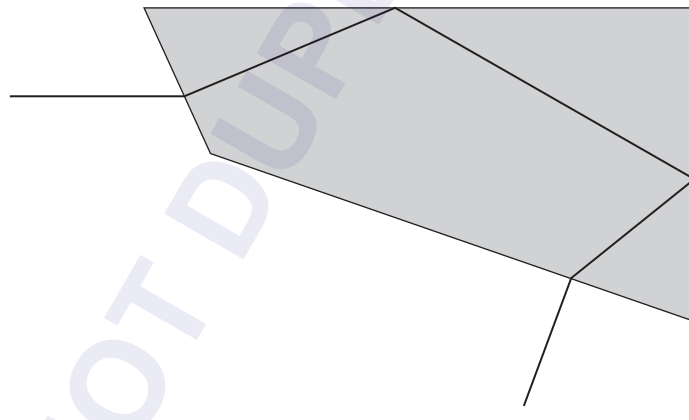


FIGURE 33 120°-deviation prism—D120.

19.7 REFERENCES

1. F. Arsenal, *Design of Fire Control Optics*, U.S. Government Printing Office, 1952.
2. D. H. Jacobs, *Fundamentals of Optical Engineering*, McGraw-Hill, 1943.
3. F. A. Jenkins and H. E. White, *Fundamental Optics*, 3d ed., McGraw-Hill, 1957.
4. W. L. Wolfe and G. J. Zissis, *The Infrared Handbook*, U.S. Government Printing Office, 1978.

This page intentionally left blank.

DO NOT DUPLICATE

DISPERSIVE PRISMS AND GRATINGS

George J. Zissis

*Environmental Research Institute of Michigan
Ann Arbor, Michigan*

20.1 GLOSSARY

A_p	prism angle
B	prism base
D_p	angle of minimum deviation
d	grating constant
E	irradiance
N	number of slits
n	refractive index
p	order number
RP	resolving power
r	angles
W	prism width
β	angle
γ	angle

20.2 INTRODUCTION

Spectroradiometers (Fig. 1) are radiometers designed specifically to allow determination of the wavelength distribution of radiation. This category of measurement systems usually consists of those in which separation of the radiation into its spectral components, or *dispersion*, is accomplished by the use of an optical element possessing a known functional dependence on wavelength—specifically, prisms and diffraction gratings. (Interferometers can also provide spectral dispersion as is discussed in Chap. 32, “Interferometers,” by Parameswaran Hariharan.)

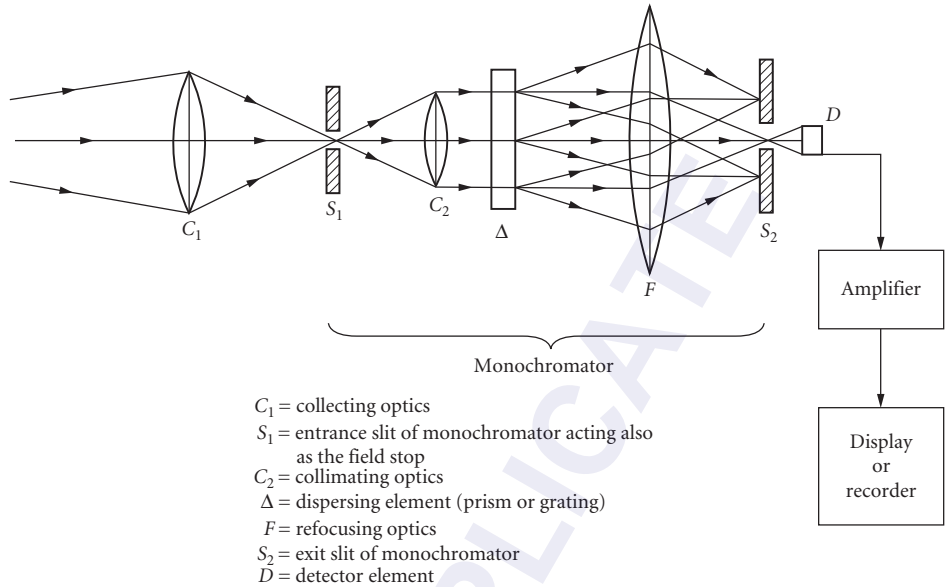


FIGURE 1 Basic spectroradiometer.

20.3 PRISMS^{1,2,3}

The wavelength dependence of the index of refraction is used in prism spectrometers. Such an optical element disperses parallel rays or collimated radiation into different angles from the prism according to wavelength. Distortion of the image of the entrance slit is minimized by the use of plane wave illumination. Even with plane wave illumination, the image of the slit is curved because not all of the rays from the entrance slit can traverse the prism in its principal plane. A prism is shown in the position of minimum angular deviation of the incoming rays in Fig. 2. At minimum angular deviation, maximum power can pass through the prism. For a prism adjusted to the position of minimum deviation,

$$r_1 = r_2 = A_p / 2 \tag{1}$$

and

$$i_1 = i_2 = (D_p + A_p) / 2 \tag{2}$$

where D_p = angle of minimum deviation for the prism
 A_p = angle of the prism
 r_1 and r_2 = internal angles of refraction
 i_1 and i_2 = angles of entry and exit

The angle of minimum deviation D_p varies with wavelength. The angular dispersion is defined as $dD_p/d\lambda$, while the linear dispersion is

$$dx/d\lambda = F dD_p/d\lambda \tag{3}$$

where F is the focal length of the camera or imaging lens and x is the distance across the image plane. It can be shown¹ that

$$dD_p/d\lambda = (B/W)(dn/d\lambda) \tag{4}$$

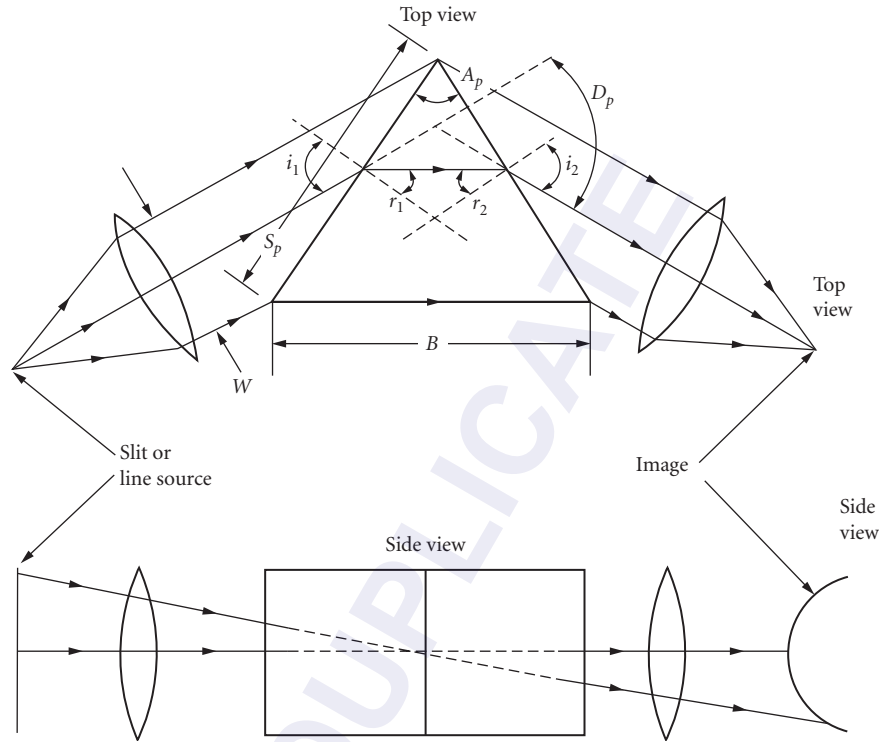


FIGURE 2 Elementary prism spectrometer schematic. W is the width of the entrance beam; S_p is the length of the prism face; and B is the prism base length.

where B = base length of the prism
 W = width of the illumination beam
 n = index of refraction
 $dx/d\lambda = F(B/W)(dn/d\lambda)$

The resolving power RP of an instrument may be defined as the smallest resolvable wavelength difference, according to the Rayleigh criterion, divided by the average wavelength in that spectral region. The limiting resolution is set by diffraction due to the finite beam width, or effective aperture of the prism, which is rectangular. Thus,

$$RP_p = B(dn/d\lambda) \quad (5)$$

If the entire prism face is not illuminated, then only the illuminated base length must be used for B .

20.4 GRATINGS

A grating is an n -slit system used in Fraunhofer diffraction with interference arising from division of the incident plane wave front. Thus it is a multiple beam interferometer described by

$$p\lambda = d(\sin \theta + \sin \phi) \quad (6)$$

where p = order number ($= 0, 1, 2, \dots$) of the principal maxima
 d = the grating constant or spacing (the distance between adjacent slits)
 ϕ = angle of incidence
 θ = angle of diffraction

The most common case is that of normal incidence, that is, $\phi = 0$, so that

$$p\lambda = d \sin \theta \quad (7)$$

and the irradiance distribution is

$$E = E_o \left\{ \frac{\sin [(\pi w \sin \theta)/\lambda]}{[(\pi w \sin \theta)/\lambda]} \right\}^2 \quad (8)$$

$$\times \left\{ \frac{\sin [(N\pi d \sin \theta)/\lambda]}{\sin [(\pi d \sin \theta)/\lambda]} \right\}^2$$

where w is the slit width and N is the number of slits or grooves. This equation is often written as

$$E = E_o [(\sin \beta)/\beta]^2 [(\sin N\gamma)/\sin \gamma]^2 \quad (9)$$

which can be considered to be

$$E = \text{constant} \times \text{single-slit diffraction function} \quad (10)$$

$$\times N\text{-slit interference function}$$

These considerations are for unblazed gratings. For a diffraction grating, the angular dispersion is given (for constant angle ϕ) by

$$dD_g/d\lambda \quad \text{or} \quad d\theta/d\lambda = p/(d \cos \theta) \quad (11)$$

The resolving power is given by

$$RP_g = pN \quad (12)$$

20.5 PRISM AND GRATING CONFIGURATIONS AND INSTRUMENTS

Classical

There are several basic prism and grating configurations and spectrometer designs which continue to be useful. One of the oldest spectrometer configurations is shown in Fig. 3.¹ Reflective interactions and prism combinations are used in Figs. 4, 5, and 6. Dispersion without deviation is realized in Figs. 7 and 8, while half-prisms are used in Fig. 9 in an arrangement which uses smaller prisms but still attains the same beam width. A few classical prism instrumental configurations are shown in Figs. 10, 11, and 12. Multiple-pass prism configurations are illustrated in Figs. 13 and 14.^{4,5}

A well-known example of a single beam double-pass prism infrared spectrometer was the Perkin-Elmer Model 112 instrument shown in Fig. 15. Infrared radiation from a source is focused by mirrors M_1 and M_2 on the entrance slit S_1 of the monochromator. The radiation beam from S_1 , path 1, is collimated by the off-axis paraboloid M_3 and a parallel beam traverses the prism for a first refraction. The beam is reflected by the Littrow mirror M_4 , through the prism for a second refraction, and focused by the paraboloid, path 2, at the corner mirror M_6 . The radiation returns along path 3, traverses the prism again, and is returned along path 4 for reflection by mirror M_7 to the exit slit S_2 . By this double dispersion, the radiation is spread out along the plane of S_2 . The radiation of the frequency interval which passes through S_2 is focused by mirrors M_8 and

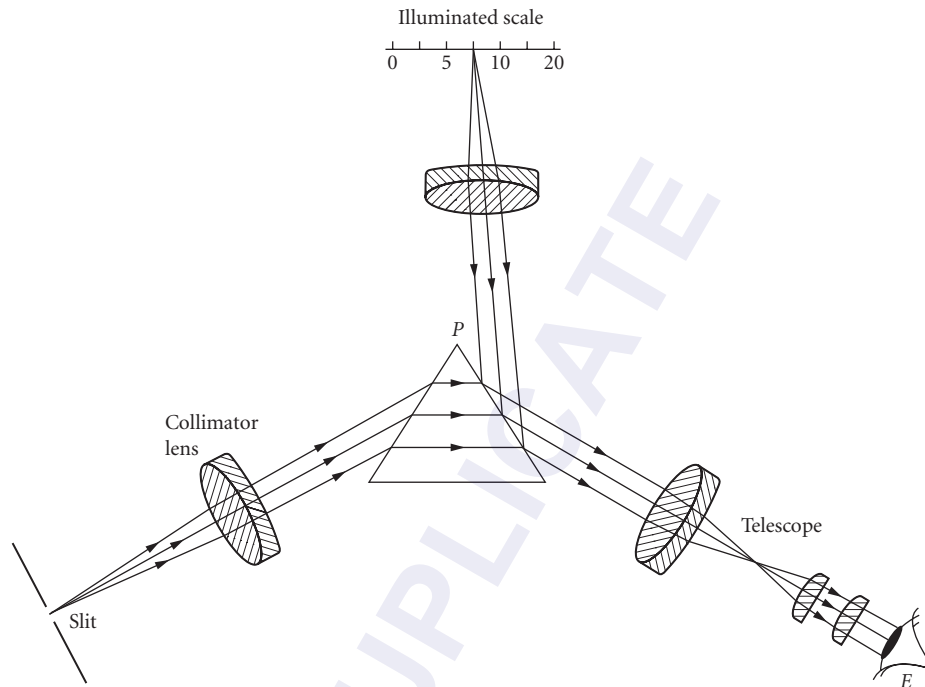


FIGURE 3 Bunsen-Kirchhoff spectrometer. An illuminated scale is reflected from the prism face into the telescope.

M_9 on the thermocouple TC. The beam is chopped by CH, near M_6 , to produce a voltage (at the thermocouple) which is proportional to the radiant power or intensity of the beam. This voltage is amplified and recorded by an electronic potentiometer. Motor-driven rotation of Littrow mirror M_4 causes the infrared spectrum to pass across exit slit S_2 permitting measurement of the radiant intensity of successive frequencies.

Gratings can be used either in transmission or reflection.⁶ Another interesting variation comes from their use in plane or concave reflection form. The last was treated most completely by Rowland, who achieved a useful combination of focusing and grating action. He showed that the radius of curvature of the grating surface is the diameter of a circle (called the Rowland circle). Any source placed on the

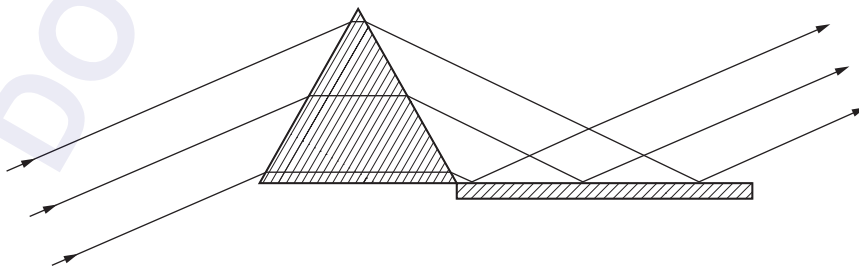


FIGURE 4 Wadsworth constant-deviation, prism-mirror arrangement. The beam enters the prism at minimum deviation and emerges displaced but not deviated from its original direction.

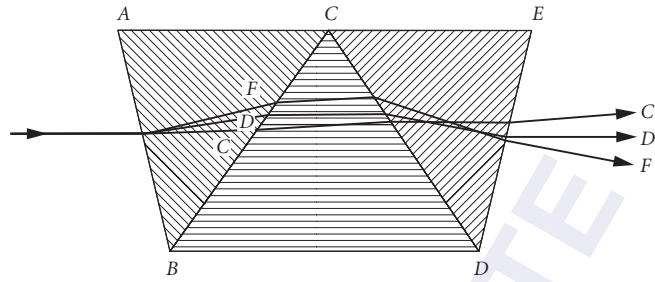


FIGURE 5 Amici prism. The central ray D enters and leaves parallel to the base. The C and F rays are deviated and dispersed.

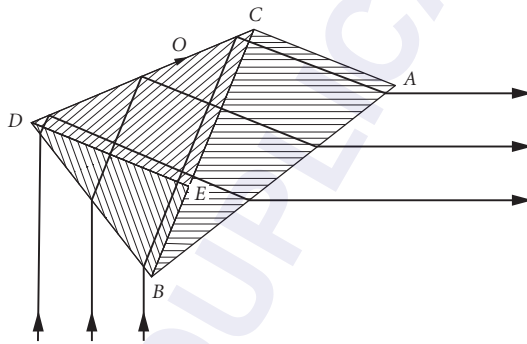


FIGURE 6 Pellin-Broca prism. The prism is equivalent to two 30° prisms, ABC and BED , and one 45° prism, DEC , but is made in one piece. The beam shown, entering at minimum deviation, emerges at 90° deviation to its entrance direction.

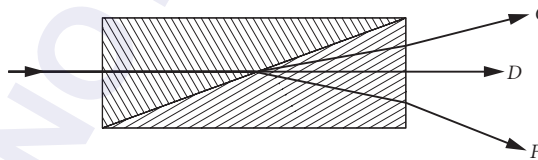


FIGURE 7 Zenger prism. The central ray D is undeviated. The C and F rays are deviated and dispersed.

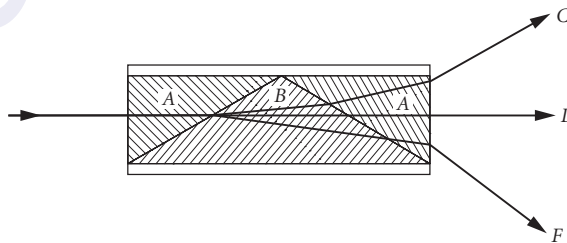


FIGURE 8 Wernicke prism. This arrangement is essentially two Zenger prisms, back-to-back.

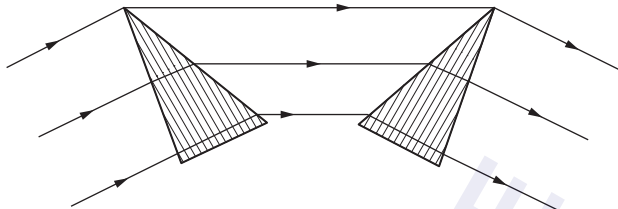


FIGURE 9 Young-Thollon half prisms. The passage of a beam at minimum deviation is shown.

circle will be imaged on the circle, with dispersion, if the rulings are made so that d is constant on the secant to the grating-blank (spherical) surface. The astigmatism acts so that a point source on a Rowland circle is imaged as a vertical line perpendicular to the plane of the circle. Rowland invented and constructed the first concave grating mounting, illustrated in Fig. 16.¹

If dispersion is sufficiently large, one may find overlapping of the lines from one order with members of the spectra belonging to a neighboring order. Errors and imperfections in the ruling of gratings can produce spurious images which are called “ghosts.” Also, the grooves in a grating can be shaped so as to send more radiation along a preferred direction corresponding to an order other than the zero order. Such gratings are said to be blazed in that order. These issues and many more involved in the production of gratings by ruling engines were thoroughly discussed by Harrison in his 1973 paper “The Diffraction Grating—An Opinionated Appraisal.”⁷

Six more grating configurations¹ which are considered to be “classics” are

1. *Paschen-Runge*, illustrated in Fig. 17. In this arrangement, one or more fixed slits are placed to give an angle of incidence suitable for the uses of the instrument. The spectra are focused along the Rowland circle PP' , and photographic plates, or other detectors, are placed along a large portion of this circle.

2. *Eagle*, shown in Fig. 18. This is similar to the Littrow prism spectrograph. The slit and plate holder are mounted close together on one end of a rigid bar with the concave grating mounted on the other end.

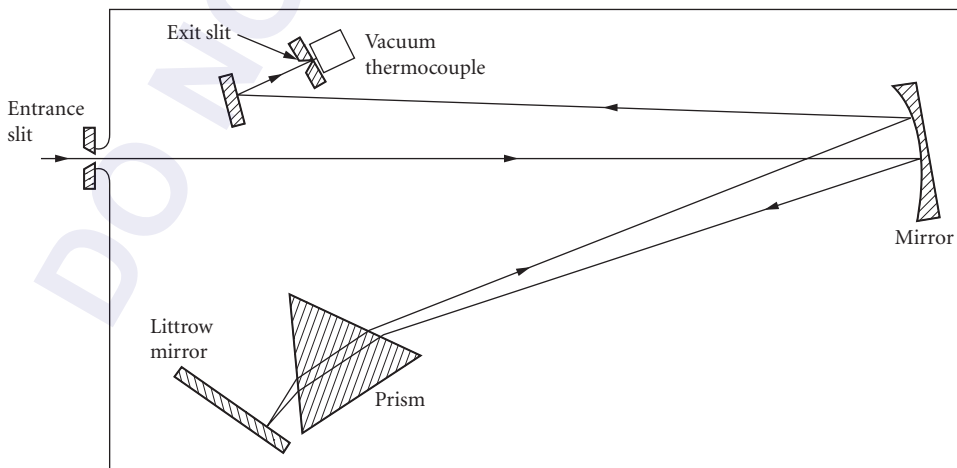


FIGURE 10 Infrared spectrograph of the Littrow-type mount with a rock salt prism.

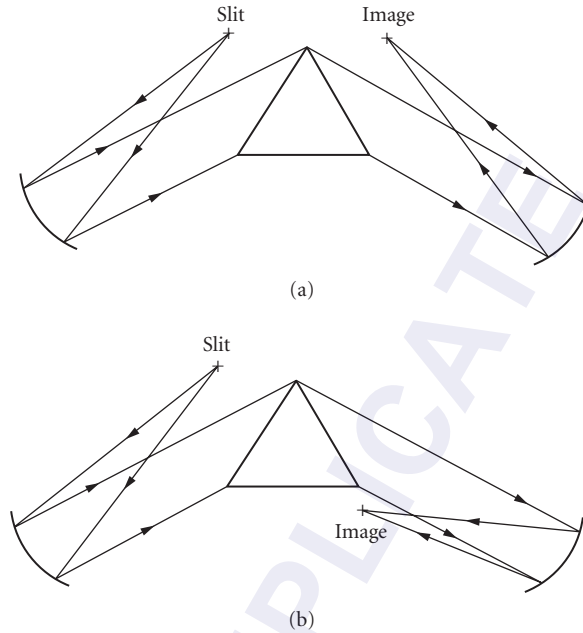


FIGURE 11 Mirror spectrometer with two choices of the location of the image. Arrangement (b) leads to smaller aberrations than arrangement (a) and is used in the Czerny-Turner mount.

3. *Wadsworth*, shown in Fig. 19. The Rowland circle is not used in this mounting in which the grating receives parallel light.

4. *Ebert-Fastie*, shown in Fig. 20. The Ebert-Fastie features a single, spherical, collimating mirror and a grating placed symmetrically between the two slits. The major advantage of the Ebert system is the fact that it is self-correcting for spherical aberration. With the use of curved slits, astigmatism is almost completely overcome.

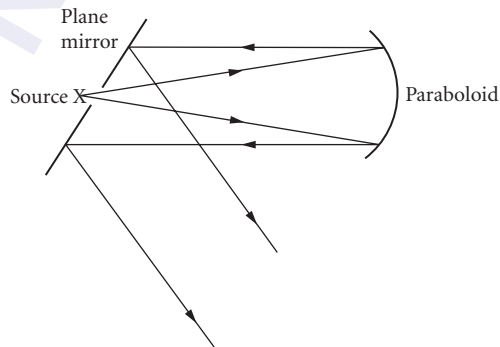


FIGURE 12 Pfund mirror. The use of a plane mirror to avoid astigmatism in the use of a paraboloidal mirror.

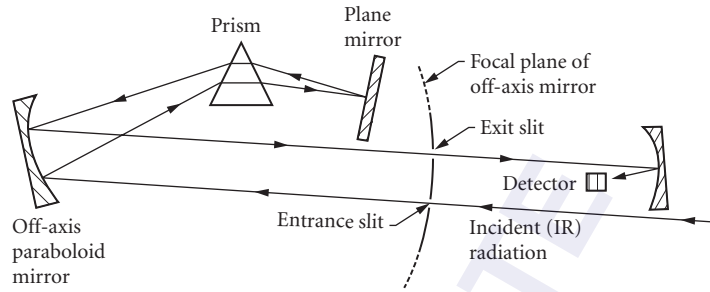


FIGURE 13 Double-pass monochromator.

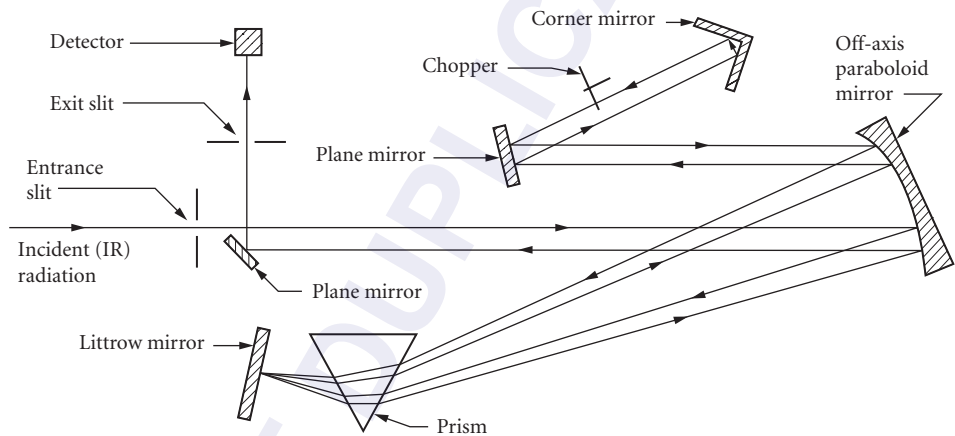


FIGURE 14 Perkin-Elmer Model 99 double-pass monochromator.

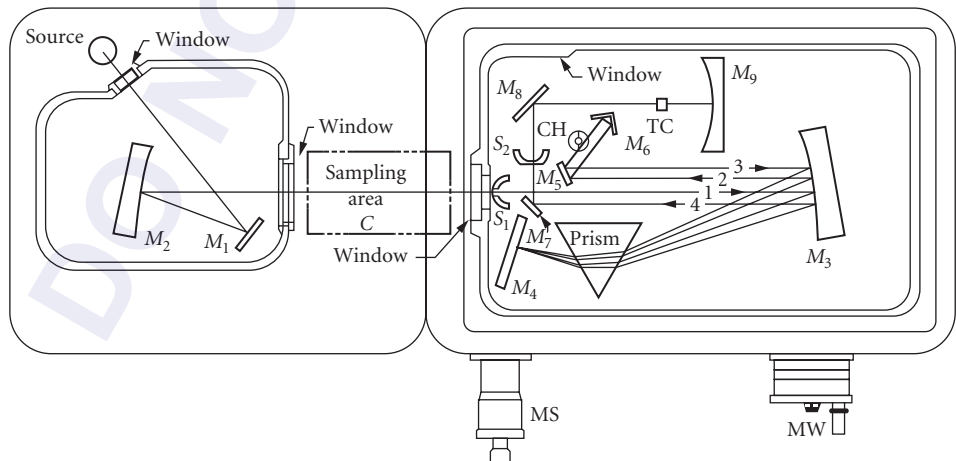


FIGURE 15 Perkin-Elmer Model 112 single-beam double-pass infrared spectrometer.

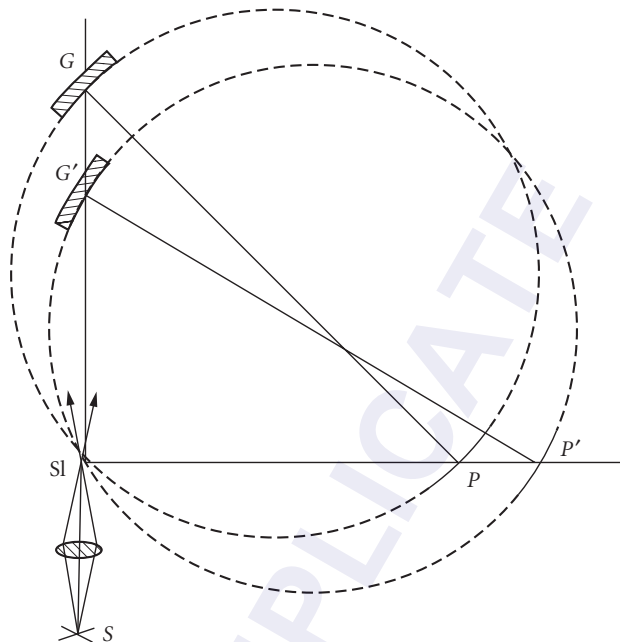


FIGURE 16 Rowland mounting of the concave grating. The grating plate-holder bar, which slides on the two perpendicular ways, is shown in two positions, GP and $G'P'$. The slit SI and source S remain fixed.

5. *Littrow*, shown in Fig. 10. The Littrow system has slits on the same side of the grating to minimize astigmatism. An advantage of the Littrow mount, therefore, is that straight slits can be used. In fact, such slits may be used even for a spherical collimating mirror if the aperture is not too large. Its greatest disadvantage is that it does not correct for spherical aberration—not too serious a defect for long focal-length/small-aperture instruments. If an off-axis parabola is used to collimate the light, aberrations are greatly reduced.

6. *Pfund*, shown in Figs. 12 and 21. This is an on-axis, Pfund-type grating instrument.⁵ Incident infrared radiation, focused by a collimating lens on the entrance slit and modulated by a chopper, passes through the central aperture of plane mirror M_1 . Reflected by the paraboloidal mirror P_1 , it emerges as a parallel beam of radiation, which is reflected by mirror M_1 to the grating. The grating is accurately located on a turntable, which may be rotated to scan the spectrum. From the grating, the diffracted beam, reflected by mirror M_2 , is focused by a second paraboloid P_2 through the central aperture of mirror M_2 to the exit slit. The emerging beam is then focused by the ellipsoidal mirror M_3 on the detector.

An off-axis, double-pass grating instrument is illustrated in Fig. 22.⁶

Combinations of prisms and gratings are not uncommon. An illustrative and complex prism-grating, double-monochromator spectrometer designed by Unicam Instruments, Ltd. is shown in Fig. 23.⁵ The prism monochromator has four interchangeable prisms, and the grating monochromator has two interchangeable gratings. The two monochromators, ganged by cams which are linear in wave number, are driven by a common shaft. The instrument can be used either as a prism-grating double monochromator, or as a prism spectrometer by blanking the grating monochromator. Gratings, prisms, and cams can be automatically interchanged by means of push buttons.

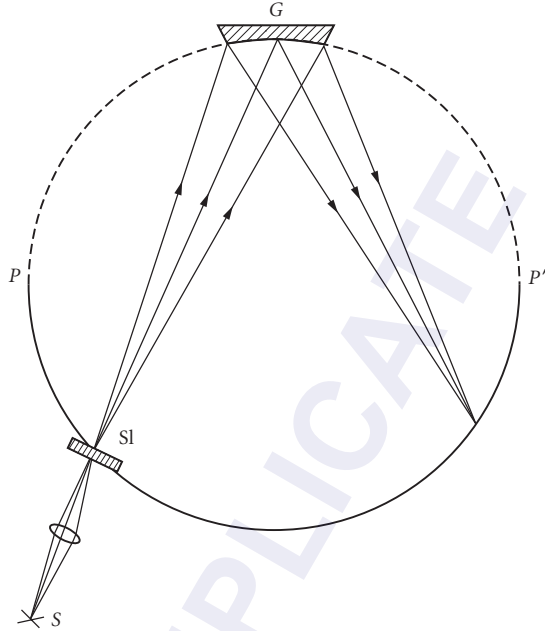


FIGURE 17 Paschen-Runge mounting of the concave grating. Sl is the slit, G is the grating, and S is the light source.

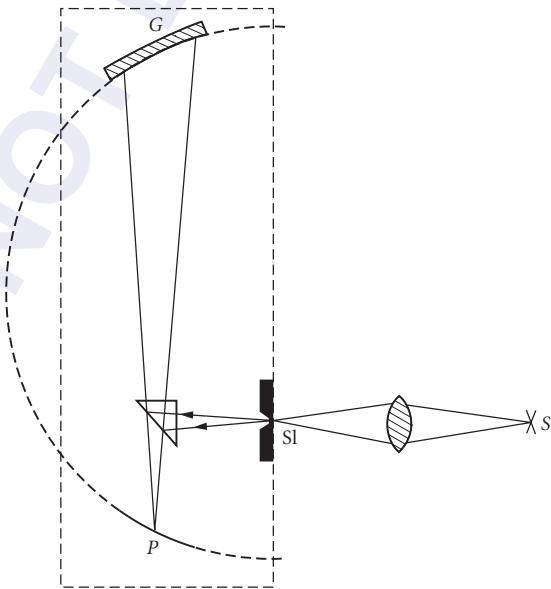


FIGURE 18 Eagle mounting on the concave grating. Sl is the slit, G is the grating, S is the light source, and P is the plate holder.

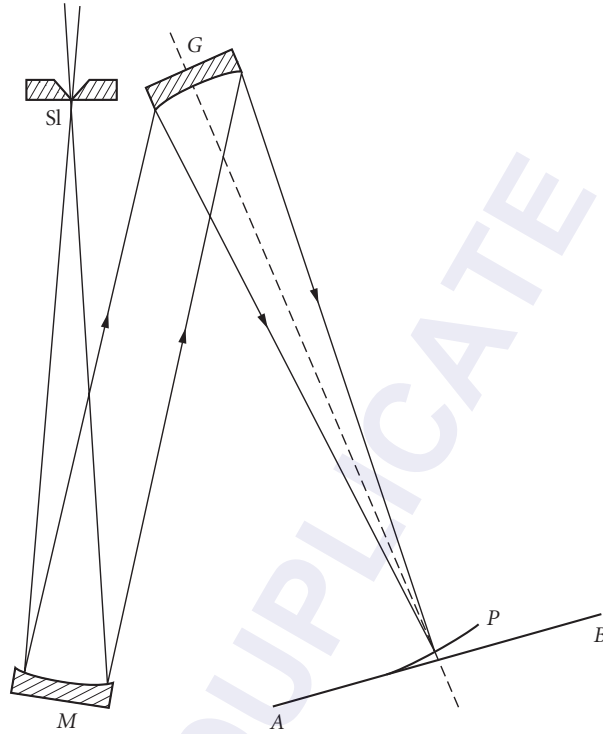


FIGURE 19 Wadsworth mounting of the concave grating. Sl is the entrance slit, G is the concave grating, M is the concave mirror, P is the plate holder, and AB is the rail for the plate holder. To minimize aberrations, one must locate the slit close to the grating.

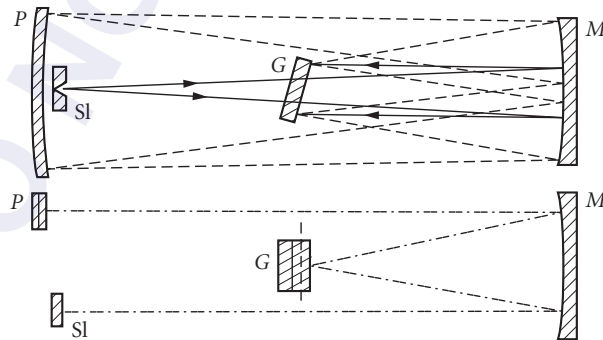


FIGURE 20 Ebert mounting of the plane grating designed by Fastie. Sl is the entrance slit, G is the grating, M is the concave mirror, and P is the photographic plate. The horizontal section is at the top and the vertical section is at the bottom.

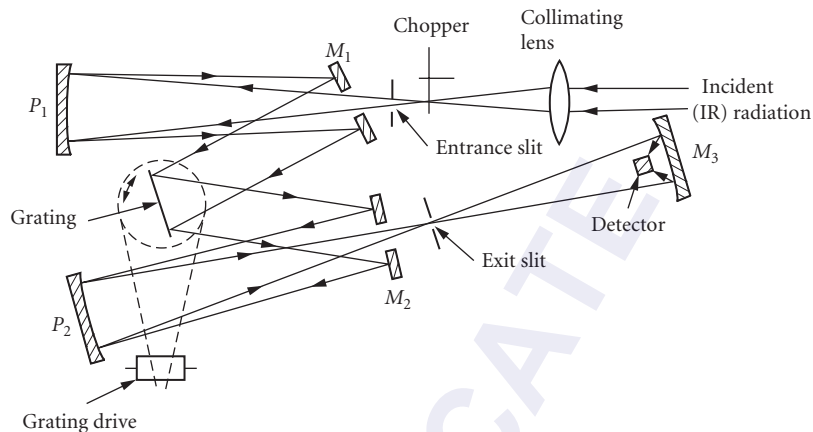


FIGURE 21 On-axis Pfund grating spectrograph.

Magnetically operated slits, programmed by a taped potentiometer, provide a constant energy background. A star-wheel, time-sharing, beam attenuator is used in the double-beam photometer.

Contemporary

In recent years there has been more attention paid to total system design and integration for specific purposes and applications, as for analytical atomic and molecular spectroscopy in analytical chemistry. Thus the conventional dispersive elements are often used in the classical configurations with variations. Innovations have come especially in designs tailored for complete computer control; introduction of one- and two-dimensional detector arrays as well as new detector types (especially for signal matching); the use of holographic optical elements either alone or combined with holographic gratings; and special data-processing software packages, displays, and data storage systems. This is the case also for interferometric systems as discussed in Chap. 32, "Interferometers," by Parameswaran Hariharan.

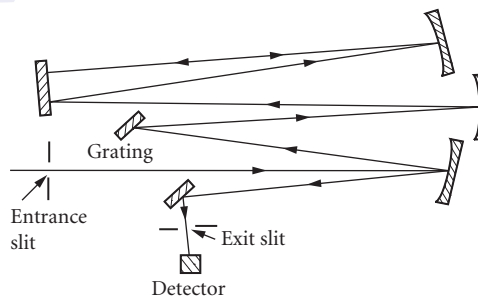


FIGURE 22 Off-axis, double-pass grating spectrograph.

TABLE 1 Examples of Prism/Grating Spectroradiometers

Manufacturer	Comments
ARC (Acton Research Corp.), Acton, Mass.	Czerny-Turner or Rowland systems with triple indexable Vac UV/IR gratings
ARIES (Acton Research Instrument & Equipment Services Inc.), QEI (Quantum Electronics Instruments Inc.), Concord, Mass.	Czerny-Turner variation with double or triple selectable gratings for 165-nm to 40- μ m regions
Beckman Instruments Inc., Fullerton, Calif.	DU Series 60 and 70 modular construction, computer-controlled spectrophotometers for analytical applications
CVI Laser Corp., Albuquerque, N. Mex.	Digikrom Monochrometers, 1/8-, 1/4-, and 1/2-m Czerny-Turner grating systems, 186 nm–20 μ m
Cary/Varian Instrument Group, San Fernando, Calif.	Cary 1, 3, 4, and 5 spectrophotometers for UV-Vis-IR; double beam, dual chopper/grating Littrow systems; attachments (e.g., reflectance) and applications software
CI Systems Ltd., New York City, N.Y. and Israel Infrared Systems, Inc., Orlando, Fla.	CVF spectroradiometers for 0.4 to 20- μ m scan CVF spectroradiometer
Instruments SA, Inc., J-Y Optical Systems, Edison, N.J.	Monochrometers, spectrometers for UV-Vis-IR, holographic gratings in Czerny-Turner or concave aberration-corrected holographic gratings and Rowland mounts; single and double pass; imaging spectrographs
LECO Corp., St. Joseph, Mich.	ICP (Inductively Coupled Plasma) spectrometer system with Pachen-Runge mount concave grating followed by an echelle and a linear detector array
Leeman Labs, Inc., Lowell, Mass.	ICP system with a fixed echelle grating followed by a prism with crossed order dispersion and scanned photomultipliers or detector arrays
McPherson, Division of SI Corp., Acton, Mass.	Double/triple monochrometers, spectroradiometers using gratings and/or prisms in Seya-Namioka, Czerny-Turner (C-T), crossed C-T, or Rowland configurations
Minirad Systems, Inc., Fairfield, Conn.	CVF and discrete filters in spectroradiometers for field measurements, 0.2 to 30 μ m
Optometrics Corp., Ayer, Mass.	Monochrometers, prism or grating, Ebert-Fastie systems for UV-Vis-NIR
Optronix Laboratories, Inc., A Subsidiary of Kollmorgen Corp., Orlando, Fla.	Spectroradiometers, UV-Vis-IR for precision measurements; filter wheels, gratings, and prisms in single/double monochrometer configurations
Oriel Corporation, Stratford, Conn.	Scanning monochrometers, rotation filter wheels, and detector array instruments
Perkin-Elmer Corporation, Norwalk, Conn.	Complete sets of UV-Vis-IR spectroscopic systems using gratings and prisms, or FT-IR, with software and hardware for computer control, and accessories for microscopy, reflectance measurement, etc.
Shimadzu Scientific Instruments, Inc., Columbia, Md.	UV-Vis-NIR spectroscopic systems using holographic gratings in Czerny-Turner mounts in single- and double-beam configurations, computer-controlled, with accessories for analyses
SPEX Industries, Inc., Edison, N.J.	UV through IR grating spectrometers, 1/2 and 1/4 m, with CCD or PDA multichannel detectors
Thermo Jarrell Ash Corp., A Subsidiary of Thermo Instrument Systems, Inc., Franklin, Mass.	Monochromators and spectroscopic systems for analyses, UV-Vis-IR with gratings (in 1942 in Wadsworth, then in 1953, Ebert, and now Paschen-Runge and crossed Czerny-Turner mounts); complete systems

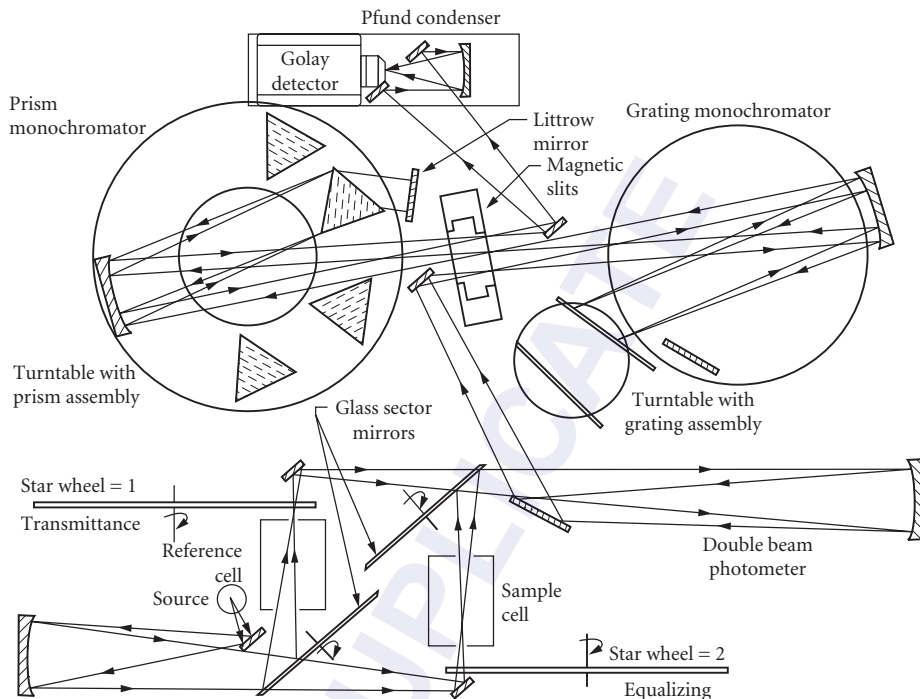


FIGURE 23 Unicam prism-grating double monochromator spectrometer.

Some examples found by a brief look through manufacturers' literature and journals such as *Spectroscopy*, *Physics Today*, *Laser Focus*, *Photonics Spectra*, and *Lasers & Optronics*,⁸ are presented in Table 1. Most of these systems are designed for analytical spectroscopy with techniques described in many texts such as Robinson's *Atomic Spectroscopy*.⁹

20.6 REFERENCES

1. R. A. Sawyer, *Experimental Spectroscopy*, 3d ed. esp. Chapters 4, 6, 7, and 11, Dover Press, New York, 1963.
2. F. A. Jenkins and H. E. White, *Fundamentals of Optics*, 4th ed., McGraw-Hill, New York, 1976.
3. E. Hecht, *Optics*, 2d ed., Addison-Wesley, Reading, MA, reprinted April 1988.
4. A. Walsh, "Multiple Monochromators II. Application of a Double Monochromator to Infrared Spectroscopy," *Journal of the Optical Society of America*, Optical Society of America, Washington, DC, vol. 42, 1952, p. 95.
5. H. L. Hackforth, *Infrared Radiation*, McGraw-Hill, New York, 1960, pp. 209, 211, 214.
6. A. H. Nielsen, "Recent Advances in IR Spectroscopy," Tech. Memo 53-2, Office of Ordnance Research, Durham, NC, December 1953.
7. G. R. Harrison, "The Diffraction Grating—An Opinionated Appraisal," *Applied Optics*, vol. 12, no. 9, 1973, p. 2039.
8. *Spectroscopy*, especially issues of May and June 1990, Aster Publishing Corp., Eugene, OR; *Physics Today*, Annual Buyers' Guide, 7 August 1990, American Institute of Physics, 335 East 45th St., New York; *Laser Focus World* and LFW's *The Buyers' Guide*, 25th ed., 1990, PennWell Publishing Co., Westford, MA; *Photonics Spectra* and *The Photonics Directory*, 4 vols., 36th ed., 1990, Laurin Publishing Company Inc., Pittsfield, MA; *Lasers & Optronics* and *L & O's 1990 Buying Guide*, Gordon Publications, Inc., Morris Plains, NJ.
9. J. W. Robinson, *Atomic Spectroscopy*, Marcel Dekker, Inc., New York, 1990.

This page intentionally left blank.

DO NOT DUPLICATE

INTEGRATED OPTICS

Thomas L. Koch

*Lehigh University
Bethlehem, Pennsylvania*

Frederick J. Leonberger

*MIT Center for Integrated Photonic Systems
Cambridge, Massachusetts*

Paul G. Suchoski

*Audigence Inc.
Melbourne, Florida*

21.1 GLOSSARY

APE	annealed proton exchange
CATV	cable television
CVD	chemical vapor deposition
CMOS	complementary metal oxide semiconductor
DBR	distributed Bragg reflector
DFB	distributed feedback
\vec{E}	electric field of propagating light
FOG	fiber optic gyroscope
Gb/s	gigabits per second
\vec{H}	magnetic field of propagating light
IO	integrated optics
IOC	integrated optic circuit
L_c	coupling length of directional coupler
LED	light-emitting diode
Mb/s	megabits per second
MMIC	monolithic millimeter-wave integrated circuit
MZ	Mach-Zehnder
n	index of refraction
OEIC	optoelectronic integrated circuit
r_{ij}	electro-optic tensor element
PIC	photonic integrated circuit
RF	radio frequency
SOI	silicon-on-insulator

t_{cutoff}	waveguide thickness for cutoff of first odd mode
TE	transverse electric mode
TM	transverse magnetic mode
V_{π}	voltage for π radian phase shift in electro-optic modulator
VLSI	very large scale integration
WDM	wavelength division multiplexing
β	propagation constant of waveguide mode
ϵ_m	field amplitude of mode m
η	coupling efficiency between modes
θ_{crit}	critical angle for total internal reflection
Λ	spatial period of periodic feature along waveguide
λ	vacuum wavelength of propagating light
λ_{PL}	photoluminescence peak wavelength of semiconductor

21.2 INTRODUCTION

The field of integrated optics is concerned with the theory, fabrication, and applications of guided wave optical devices and circuits. These structures guide light along the surface of a wafer typically using dielectric waveguides that confine light to lateral dimensions on the scale of the optical wavelength. Guided wave devices that perform passive operations analogous to classic optics, such as reflecting, beam splitting, attenuating and spectral filtering, can be formed using microelectronic-based fabrication techniques. By fabricating devices in active materials such as ferroelectrics, modulators, and switches based on the classic electro-optic effect can be formed. Compound semiconductors such as GaAs or InP additionally allow for the detection of light, and generation and amplification of light with light-emitting diodes (LEDs), lasers, and optical amplifiers. Extremely compact passive and active optical devices have also recently been demonstrated in silicon, taking explicit advantage of the highly advanced very large scale integration (VLSI) process technology developed for electronics. The monolithic integrate of optically interconnected passive and active devices in a multicomponent circuit is referred to as an *integrated optic circuit* (IOC) or a *photonic integrated circuit* (PIC), with the latter term usually applied to active semiconductor-based circuits. In semiconductor materials, purely electronic devices such as transistors can be integrated as well to form what is often referred to as an optoelectronic integrated circuit (OEIC).

Progress in the field of integrated optics has been rapid since its inception¹ in 1969. Much of this progress is due the availability of increasingly high-quality materials, microelectronic-processing equipment and techniques, and the overall rapid advancement and deployment of fiber optic systems. Interest in integrated optics stems from its numerous advantages over other optical technologies. Integrated optics devices interface efficiently with optical fibers, and can reduce cost in complex circuits by eliminating the need for separate, individual packaging of each circuit element. They also offer smaller size and weight, lower power consumption, improved reliability, and often larger electrical modulation bandwidths compared to their bulk-optic counterparts.

The applications for integrated optics are widespread. Generally these applications involve interfacing with single-mode fiber optic systems such as digital and analog communications, but also include RF signal processing using optical techniques, laser beam control, and navigational and biomedical sensors. Integrated optics is viewed in the marketplace as a key enabling technology for high-speed digital optical fiber telecommunications, cable television (CATV) signal distribution, and fiber optic gyroscopes (FOG), and it will certainly continue to have a major impact on progress in broadband information distribution and connectivity.

This chapter reviews the integrated optics (IO) field, beginning with a brief review of IO device physics and fabrication techniques. A phenomenological description of IO circuit elements, both passive and active, is given, followed by a discussion of IO applications and system demonstrations.

The chapter concludes with a short look at future trends. Due to the brevity of this chapter relative to the work in the field, much of the coverage is necessarily limited. The reader is referred to Refs. 2–17 for more detailed information at a variety of levels.

21.3 DEVICE PHYSICS

Optical Waveguides

Central to integrated optics is the concept of guiding light in dielectric waveguide structures with dimensions comparable to the wavelength of the guided light. In this section we present only a brief survey of the relevant physics and analysis techniques used to study their properties. The reader is referred to a number of excellent texts dealing with this topic for a more comprehensive treatment.^{2–6}

A dielectric waveguide confines light to the core of the waveguide by somehow reflecting power back toward the waveguide core that would otherwise diffract or propagate away. While any means of reflection can accomplish this end (e.g., glancing-incidence partial reflections from interfaces between different media can serve as the basis for *leaky* waveguides), the most common technique employs a 100 percent *total internal reflection* from the boundary of a high-index core and a lower-index cladding material. As light propagates down the axis of such a structure, the waveguide cross section can also be viewed as a lens-like phase plate that provides a larger retardation in the core region. Propagation down the guide then resembles a continuous refocusing of light that would otherwise diffract away.

The pedagogical structure used to illustrate this phenomenon is the symmetric slab waveguide, composed of three layers of homogeneous dielectrics as shown in Fig. 1. It is well known that propagation in slab structures can be analyzed using either a ray-optics approach, or through the use of interface boundary conditions applied to the simple solutions of Maxwell's equations in each homogeneous layer of the structure.^{2–6} In the ray-optics description, the rays represent the phase fronts of two intersecting plane waves propagating in the waveguide core region. Since the steady-state field has a well-defined phase at each point, a finite set of discrete modes arises from the self-consistency condition that, after propagation and two reflections from the core-cladding boundaries, any phase front must rejoin itself with an integral multiple of a 2π phase shift. For a given core thickness, there will be a limited discrete number of propagation angles in the core that satisfy this criterion, with the lower bound on the angle given by the critical angle for total internal reflection, $\theta_{\text{crit}} = \sin^{-1}(n_0/n_1)$. In general, a thicker and higher-index waveguide core will admit a larger number of confined solutions or *bound* modes. Figure 1 shows both the fundamental even mode and the first higher order,

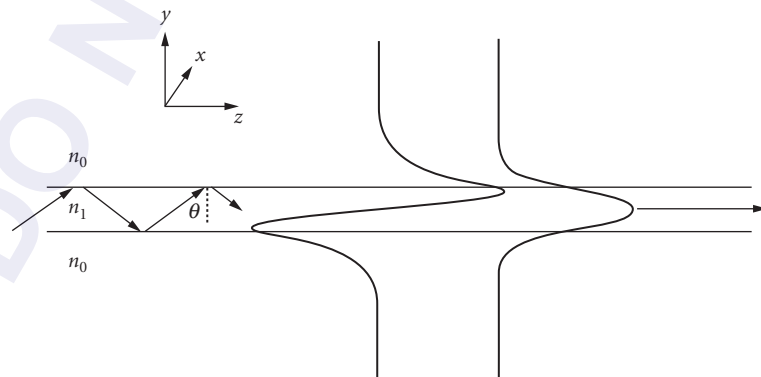


FIGURE 1 A symmetric three-layer slab waveguide. The fundamental even and first odd mode are shown.

odd mode. If the dimensions are small enough, only one bound mode for each polarization state will exist and the guide is termed a *single-mode* waveguide. Care must be exercised to include the angle-dependent phase shift upon total internal reflection, referred to as the Goos-Hanchen shift, that can be viewed as a displaced effective reflection plane.² The quantity $\beta = (2\pi n_1 / \lambda) \cdot \sin\theta$, referred to as the propagation constant, is the z projection of the wave vector and thus governs the phase evolution of the field along the length of the guide. In addition to the discrete set of bound modes, plane waves can also enter from one side and pass at an angle *through* such a structure, and form a continuous set of *radiation* modes. In an asymmetrical structure, some of the radiation modes may be propagating on one side of the guide, but evanescent on the other. From a mathematical point of view, the set of all bound and radiation modes in a nondissipative structure form a complete set for expansion of any electromagnetic field in the structure. Analysis techniques will be discussed in more detail below.

The slab waveguide in Fig. 1 employed total internal reflection from an abrupt index discontinuity for confinement. Some fabrication techniques for waveguides, particularly in glasses or electro-optic materials such as LiNbO₃,³ achieve the high-index core by impurity diffusion or implantation, leading to a *graded* index profile. Here a field solution will usually be required to properly describe the modes and the ray paths become curved, but total internal reflection is still responsible for confinement.

Most useful integrated optics devices require waveguide confinement not just in a two-dimensional slab, but in a stripe or channel geometry. While recognizing that the vector nature of the electromagnetic field makes the rigorous analysis of a particular structure quite cumbersome, the reader can appreciate that the same phenomenon of confinement by reflection will be operative in two dimensions as well. Figure 2 shows the cross sections of the most common stripe or channel waveguide types used in integrated optics. Common to the cross section for all these structures is a region on the waveguide axis containing higher index material than the surrounding cladding areas. The diffused waveguide may require a full two-dimensional analysis, but a common technique for the approximate analysis of high-aspect-ratio channel guides such as in Fig. 2a, b, c, and d, is the *effective index method*.^{2,3,18} In this technique a slab waveguide analysis is applied sequentially to the two dimensions. First three separate vertical problems are solved to obtain the modal phase index $n_{\text{mode}} \equiv \beta \cdot \lambda / 2\pi$ for a given polarization mode in each lateral region as if it were an infinite

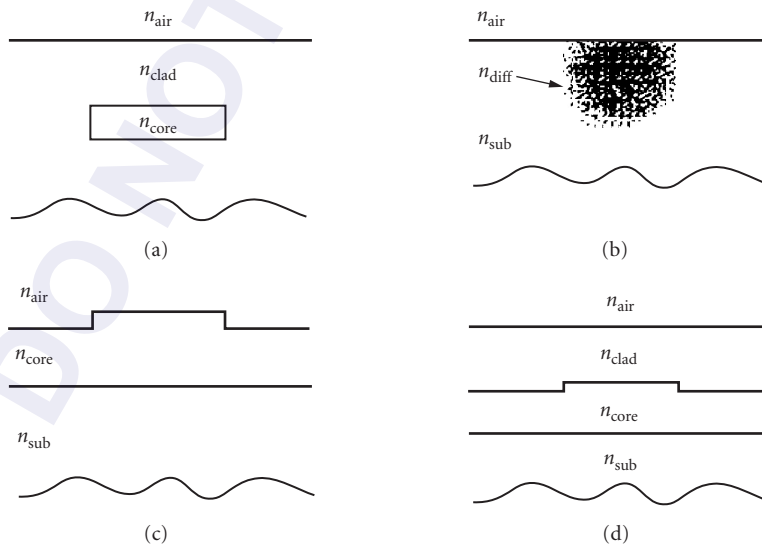


FIGURE 2 Various types of channel or stripe waveguides .

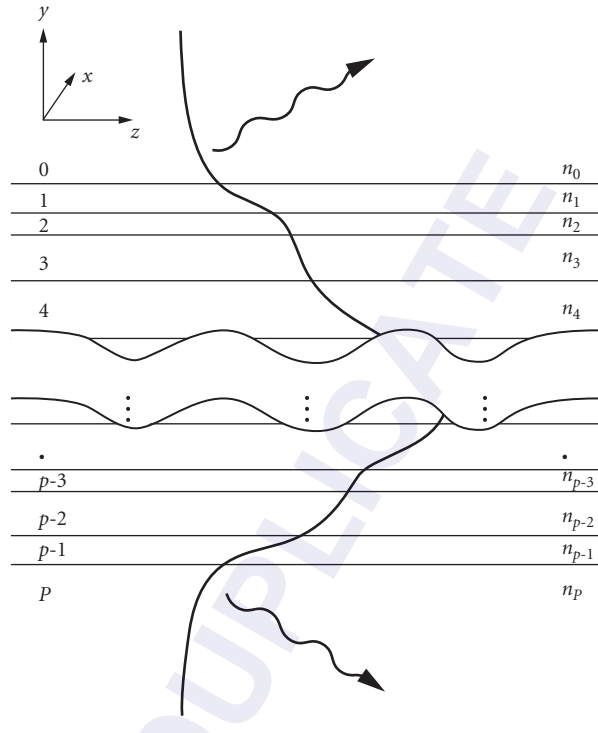


FIGURE 3 A general multilayer slab waveguide structure.

slab. These indices are then used as input to a final “effective” slab waveguide problem in the lateral dimension using the opposite polarization boundary conditions. Since the properties of multilayer slab waveguides play an important role in waveguide analysis, a more comprehensive general formulation is outlined below. This task is more tractable using the field solutions of Maxwell’s equations than the ray-optics approach.

A general multilayer slab is shown in Fig. 3. Since the wave equation in this case is separable, we need only consider a two-dimensional problem in the y direction perpendicular to the layers, and a propagation direction z . Here we will consider the concept of a mode in such a structure to be quantified in physical terms as a solution to Maxwell’s equations whose sole dependence on the coordinate in the propagation direction z is given by $e^{i\beta z}$. This translates to a requirement that the *shape* of the field distribution in the y direction, perpendicular to layers, remain unchanged with propagation. In this manner we can easily generalize to leaky structures or materials exhibiting loss or gain, where β may be complex to allow for the scaling of the mode amplitude with propagation, but the relative mode profile in the perpendicular y direction still remains constant. These latter solutions are not normalizable or “proper” in the sense of mathematical completeness, but are very useful in understanding propagation behavior in such structures.

Since the field in each homogeneous layer m is well known to be $e^{\pm i\vec{k}_m \cdot \vec{r}}$, with $|\vec{k}_m| = 2\pi n_m / \lambda$ for the (generally complex) index of refraction n_m , the general solution to the field amplitude in each layer m is

$$\epsilon_m = \left[a_m e^{i q_m y} + b_m e^{-i q_m y} \right] e^{i \beta z} \quad (1)$$

where $q_m \equiv [2\pi n_m / \lambda]^2 - \beta^2]^{1/2}$. Inspection of the vector Maxwell’s equations reveals that the general vector solution in the multilayer slab can be broken down into the superposition of a TE

(transverse electric) and a TM (transverse magnetic) solution.^{2,3} The TE (TM) solution is characterized by having only one component of the electric (magnetic) field that points in the x direction, parallel to the layers and perpendicular to the propagation direction z . The mode field amplitude ε_m in Eq. (1) refers to the E_x or the H_x field for the TE and TM case, respectively.

In a very simple exercise, for each of these cases one can successively match boundary conditions for continuous tangential \vec{E} and \vec{H} across the interfaces to provide the coefficients a_{m+1} and b_{m+1} in each layer $m+1$ based upon the value of the coefficients in the preceding layer m ,

$$\begin{bmatrix} a_{m+1} \\ b_{m+1} \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \left(1 + \frac{q_m \gamma_m}{q_{m+1} \gamma_{m+1}} \right) e^{-i(q_{m+1} - q_m) \gamma_m} & \frac{1}{2} \left(1 - \frac{q_m \gamma_m}{q_{m+1} \gamma_{m+1}} \right) e^{-i(q_{m+1} + q_m) \gamma_m} \\ \frac{1}{2} \left(1 - \frac{q_m \gamma_m}{q_{m+1} \gamma_{m+1}} \right) e^{i(q_{m+1} + q_m) \gamma_m} & \frac{1}{2} \left(1 + \frac{q_m \gamma_m}{q_{m+1} \gamma_{m+1}} \right) e^{i(q_{m+1} - q_m) \gamma_m} \end{bmatrix} \begin{bmatrix} a_m \\ b_m \end{bmatrix} \quad (2)$$

where γ_m are the coordinates of the interfaces between layers m and $m+1$, and $\gamma_m \equiv 1$ for TE modes and $\gamma_m \equiv n_m^{-2}$ for TM modes. The wave is assumed evanescently decaying or outward leaking on one initial side of the arbitrary stack of complex-index layers, that is $b_0 = 0$ on the uppermost layer. When the lowermost “cladding” layer $m = p$ is reached, one again demands that only the coefficient b_p of the evanescently decaying, or possibly the outward leaking, component be nonzero, which recursively provides the eigenvalue equation $a_p(\beta) = 0$ for the eigenvalues β_j . Arbitrarily letting $a_0 = 1$, this can be written explicitly as

$$a_p(\beta) = [1 \quad 0] \cdot \left[\prod_{m=p-1}^{m=0} \mathbf{M}_m(\beta) \right] \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 0 \quad (3)$$

where $\mathbf{M}_m(\beta)$ is the matrix appearing in Eq. (2). This method is essentially the one described for real propagation constants by Kogelnik,² who further provides important variational and perturbative expressions for determining the changes in propagation constants due to local changes in materials properties as might be used for modulation. In practice Eq. (3) is solved numerically in the form of two equations (for the real and imaginary parts of a_p) in two unknowns (the real and imaginary parts of β). Once the complex solutions β_j are obtained using standard root-finding routines, the spatial profiles are easily calculated for each mode j by actually evaluating the coefficients for the solutions using the relations above with $a_0 = 1$, for example.

Application of Eq. (3) to the simple symmetric slab of Fig. 1 with thickness t , and real core index n_1 and cladding index n_0 can be reduced with some trigonometric half-angle identities to a simple set of equations with intuitive solutions by graphical construction.¹⁹ Defining new independent variables $r \equiv t/2[(2\pi n_1/\lambda)^2 - \beta^2]^{1/2}$ and $s \equiv t/2[\beta^2 - (2\pi n_0/\lambda)^2]^{1/2}$, one must simultaneously solve for positive r and s the equation

$$r^2 + s^2 = (\pi t/\lambda)^2 (n_1^2 - n_0^2) \quad (4)$$

and either one of the following equations:

$$s = \frac{\gamma_1}{\gamma_0} \cdot r \cdot \begin{cases} \tan(r) & \text{(even modes)} \\ -\cot(r) & \text{(odd modes)} \end{cases} \quad (5)$$

where again $\gamma_m \equiv 1$ for TE modes and $\gamma_m \equiv n_m^{-2}$ for TM modes.

By plotting the circles described by Eq. (4) and the functions in Eq. (5) in the (r, s) plane, intersections provide the solutions (r_j, s_j) for mode j , yielding β_j from the definition of either r or s . This construction is shown in Fig. 4 for TE modes, where Eq. (4) has been parametrized with $u \equiv (\pi t/\lambda) (n_1^2 - n_0^2)^{1/2}$. Due to the presence of γ_m in Eq. (5), the TE and TM modes will have different propagation constants, leading to waveguide birefringence.

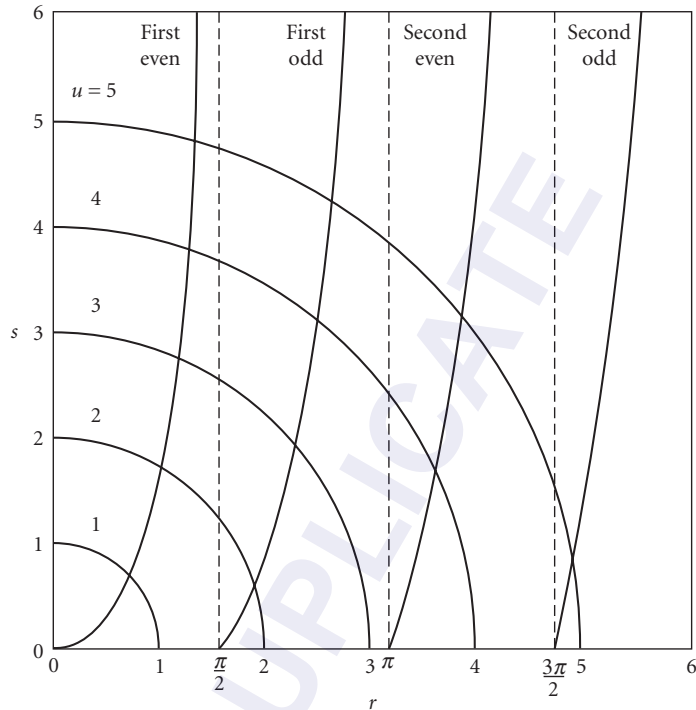


FIGURE 4 A graphical solution for the symmetric three-layer slab waveguide. For an arbitrary value of the parameter u , solutions are found at the intersections of the circular arcs and the transcendental functions as shown.

It is easy to see from the zero-crossing of Eq. (5) at $r = \pi/2$ that the cutoff of the first odd mode occurs when the thickness reaches a value

$$t_{\text{cutoff}} = (\lambda/2) \cdot (n_1^2 - n_0^2)^{-1/2} \quad (6)$$

Another important feature of the symmetric slab is the fact that neither the TE nor TM fundamental (even) mode is ever cut off. This is not true for asymmetric guides. More complicated structures are easily handled using Eq. (2), and one can also approximate graded-index profiles using the multi layer slab. However, analytical solutions also exist¹ for a number of interesting graded-index profiles, including parabolic, exponential, and “cosh⁻²” index profiles.

Once the modes of a waveguide are known, there are many physical phenomena of interest that can be easily calculated. Quite often the propagation constant is required to evaluate the phase evolution of guided-wave components. In other cases the frequency dependence of the propagation constant is required to get the group velocity, $v_g = [\partial\beta/\partial\omega]^{-1}$, which determines quantities such as the flow of energy or the mode spacing in resonators. In other cases, the designer may need to know the electric field amplitude, or perhaps the fraction of the propagating energy, that lies within a certain layer of the waveguide system.

One critical application lies in evaluating how light couples between guided-wave components with different modal structure, or from free space into a guided-wave component. For example, it is easy to show from the mathematical completeness of the waveguide modes²⁻⁴ that the energy

efficiency η of the coupling from a field ϵ_{inj} injected at the input facet of the waveguide into a particular TE mode ϵ_m of a waveguide is given by

$$\eta = \frac{|\int \epsilon_m^*(y) \epsilon_{\text{inj}}(y) dy|^2}{\int |\epsilon_m(y)|^2 dy \cdot \int |\epsilon_{\text{inj}}(y)|^2 dy} \quad (7)$$

Another mathematical formalism, called coupled-mode theory,^{2,3,9} is one of the most important design tools for the guided-wave device designer and allows for the calculation of the coupling between parallel waveguides as in a directional coupler. It also allows for the evaluation of coupling between distinct modes of a waveguide when a longitudinal perturbation along the propagation direction destroys the exact mode orthogonality. An important example of the latter is when the perturbation is periodic as in a corrugated-waveguide grating where the grating allows for phase matching between modes with different propagation constants. Waveguide gratings are used to form wavelength selective coupling as in Bragg reflectors, distributed feedback lasers, and other grating-coupled devices. The comprehensive treatment of these phenomena is beyond the scope of this chapter, and the reader is referred to the references cited above. In some instances, evaluating the performance of devices where radiation plays a significant role may be tedious using a modal analysis, and numerical techniques such as the *beam propagation method* (BPM) are used to actually launch waves through the structure to evaluate radiation losses in waveguide bends, branches, or complicated splitters.²⁰

Index of Refraction and Active Index-Changing Mechanisms

Index of Refraction Waveguide analysis and design requires precise knowledge of the material index of refraction. One of the most common electro-optic materials is LiNbO₃, a uniaxial birefringent crystal whose index can be characterized by providing the wavelength-dependent ordinary and extraordinary indices n_o and n_e . They are given by²¹

$$n_{o,e}^2 = A_{o,e} + \frac{B_{o,e}}{D_{o,e} - \lambda^2} + C_{o,e} \lambda^2 \quad (8)$$

where $A_o = 4.9048$
 $B_o = -0.11768$
 $C_o = -0.027169$
 $D_o = 0.04750$
 $A_e = 4.5820$
 $B_e = -0.099169$
 $C_e = -0.021950$
 $D_e = 0.044432$

Glasses are also common substrate materials, but compositions and indices are too varied to list here; indices usually lie in the range of 1.44 to 1.65, and are mildly dispersive. Commonly deposited dielectrics are SiO₂ and Si₃N₄, with indices of ~ 1.44 and ~ 2.0 at 1.55 μm . The reader is referred to various tables in the literature for more detail.²¹

InP and GaAs are by far the most common substrates for IO devices in semiconductors. The usual epitaxial material on GaAs substrates is Al_xGa_{1-x}As, which is nearly lattice-matched for all values of x , with GaAs at $x = 0$ providing the narrowest bandgap. For photon energies below the absorption edge, the index of refraction of this material system is given by²²

$$n_{\text{AlGaAs}}(E, x) = \left[1 + \gamma(E_f^4 - E_r^4) + 2\gamma(E_f^2 - E_r^2)E^2 + 2\gamma E^4 \cdot \ln \left(\frac{E_f^2 - E^2}{E_r^2 - E^2} \right) \right]^{1/2} \quad (9)$$

where $E = 1.2398/\lambda$ is the incident photon energy,

$$\gamma = \frac{E_d}{4E_0^3(E_0^2 - E_\Gamma^2)} \quad \text{and} \quad E_f = (2E_0^2 - E_\Gamma^2)^{1/2} \quad (10)$$

where $E_0(x) = 3.65 + 0.871x + 0.179x^2$

$$E_d(x) = 36.1 - 2.45x$$

$$E_\Gamma(x) = 1.424 + 1.266x + 0.26x^2$$

For devices fabricated on InP substrates, common in telecommunications applications for devices in the 1.3- μm and 1.55- μm bands, the most common epitaxial material is a quaternary alloy composition $\text{In}_x\text{Ga}_{1-x}\text{As}_y\text{P}_{1-y}$. In this case the material is only lattice matched for the specific combination $y = 2.917x$, and this lattice-matched alloy can be characterized by its photoluminescence wavelength λ_{PL} under low intensity optical excitation. The index of this quaternary alloy is given by²³

$$n_Q(E, E_{\text{PL}}) = \left(1 + \frac{A_1}{1 - \left(\frac{E}{E_{\text{PL}} + E_1} \right)^2} + \frac{A_2}{1 - \left(\frac{E}{E_{\text{PL}} + E_2} \right)^2} \right)^{1/2} \quad (11)$$

where $E = 1.2398/\lambda$ and $E_{\text{PL}} = 1.2398/\lambda_{\text{PL}}$ are, respectively, the incident photon energy and photoluminescence peak photon energy for λ in micrometer and $A_1(E_{\text{PL}})$, $A_2(E_{\text{PL}})$, E_1 , E_2 are fitted parameters given by

$$A_1 = 13.3510 - 5.4554 \cdot E_{\text{PL}} + 1.2332 \cdot E_{\text{PL}}^2$$

$$A_2 = 0.7140 - 0.3606 \cdot E_{\text{PL}}$$

$$E_1 = 2.5048 \text{ eV} \quad \text{and} \quad E_2 = 0.1638 \text{ eV}$$

For application to the binary InP, the value of the photoluminescence peak should be taken as $\lambda_{\text{PL}} = 0.939 \mu\text{m}$.

Many integrated optics devices rely on active phenomena such as the electro-optic effect to alter the real or imaginary index of refraction. This index change is used to achieve a different device state, such as the tuning of a filter, the switching action of a waveguide switch, or the induced absorption of an electroabsorption modulator. Below, a brief survey is provided of the most commonly exploited index-changing phenomena.

Linear Electro-Optic Effect The linear electro-optic or Pockels effect refers to the change in the optical dielectric permittivity experienced in noncentrosymmetric ordered materials that is *linear* with applied quasi-static electric field. This relation is commonly expressed using the dielectric impermeability $(1/n^2)_{ij} \equiv \epsilon_0 \partial E_i / \partial D_j$ appearing in the *index ellipsoid* equation for propagation in anisotropic crystals.²⁴ Convention employs symmetry arguments to contract all the $(1/n^2)_{ij}$ to only six independent values which are then denoted by a single subscript $(1/n^2)_i$ for $i = 1, \dots, 6$. In the *principal axes* coordinate system, the impermeability is diagonalized and $(1/n^2)_i = 0$ for $i = 4, 5, 6$ in the absence of an applied electric field, with the value of $(1/n^2)_{ij}$ providing the inverse square of the index for optical fields polarized along each axis $i = 1, 2, 3$. For an electric field expressed in the

principal axes coordinate system, the changes in the coefficients are then evaluated using the 6×3 electro-optic tensor \mathbf{r}

$$\Delta\left(\frac{1}{n^2}\right)_i = \sum_{j=1}^3 r_{ij} E_j \quad i=1, \dots, 6 \quad (12)$$

With an applied field, the equation for the index ellipsoid in general must be *redagonalized* to again yield $(1/n^2)_i = 0$ for $i = 4, 5, 6$. This provides a new set of principal axes and the coefficients in the new index ellipsoid equation provide the altered value of the refractive index along each new principal axis.

For a particular crystal, symmetry also limits the number of nonzero r_{ij} that are possible. In the cubic zinc-blend III-V compounds there are only three equal nonzero components $r_{63} = r_{52} = r_{41}$ and the analysis is relatively easy. As a specific example, consider a static field \mathbf{E} applied along the (001) direction, surface normal to the wafers commonly used for epitaxial growth. The re-diagonalized principal axes in the presence of the field become the (001) direction (z axis), the (011) direction (x axis), and the (01 $\bar{1}$) direction (y axis); the latter two directions are the cleavage planes and are thus common directions for propagation. The respective index values become

$$\begin{aligned} n_x &= n_0 - \frac{1}{2} n_0^3 r_{41} E \\ n_y &= n_0 + \frac{1}{2} n_0^3 r_{41} E \\ n_z &= n_0 \end{aligned} \quad (13)$$

If we thus consider a slab guide on a (001) substrate and propagation in the (011) direction, the applied field would produce a phase retardation for TE-polarized light of $\Delta\phi = \pi/\lambda n_0^3 r_{41} \mathbf{E} \cdot \mathbf{L}$ after a propagation length L . With values of $r_{41} \sim 1.4 \times 10^{-10}$ cm/V, micron-scale waveguide structures in GaAs or InP lead to retardations in the range of $10^\circ/\text{V}\cdot\text{mm}$. This TE retardation could be used as a phase modulator, or in a Mach-Zehnder (MZ) interferometer to provide intensity modulation. For fields applied in other directions such as the (011), the principal axes are rotated away from the (011) and (01 $\bar{1}$) directions. Propagation along a cleavage direction can then alter the polarization state, a phenomenon that also has device implications as will be discussed in more detail later.

In the case of LiNbO_3 and LiTaO_3 , two of the most common integrated optic materials for electro-optic devices, the dielectric tensor is more complex and the materials are also birefringent in the absence of an applied field. There are eight nonzero components to the electro-optic tensor, $r_{22} = -r_{12} = -r_{61}$, $r_{51} = r_{42}$, $r_{13} = r_{23}$, and r_{33} . For LiNbO_3 , the largest coefficient is $r_{33} \sim 30.8 \times 10^{-10}$ cm/V. Both retardation and polarization changes are readily achieved, and the reader is referred to Chap. 7, "Electro-Optic Modulators," Vol. V of this *Handbook* or the literature for a comprehensive treatment of the properties of these and other materials.^{24,25}

The electro-optic effect in these materials is associated with field-induced changes in the positions of the constituent atoms in the crystal, and the resulting change in the crystal polarizability. The absorption induced by the conventional electro-optic effect is thus negligible. Below, both free-carrier and field-induced absorption effects that are readily observed in semiconductors will be described.

Carrier Effects In semiconductors, other powerful index-changing mechanisms are available related to the interaction of the optical field with the free electrons or holes. The simplest of these is the plasma contribution resulting from the polarizability of the mobile carriers. This mechanism is important in active integrated optical devices in both compound semiconductors and in silicon, where in the latter case it provides one of the only viable technique for modulation. According to the simple Drude model,²⁶ this is given for each carrier species by $\Delta n \approx -N \cdot e^2 \lambda^2 / (8\pi^2 \epsilon_0 n c^2 m^*)$ in MKS units, where N and m^* are the carrier concentration and effective mass, e is the electronic charge, and ϵ_0 is the free-space permittivity. This can produce index changes approaching $\Delta n \sim -0.01$ at $10^{18}/\text{cm}^3$ electron-/hole-injected carrier density at 1.5- μm wavelengths. Static index shifts can also be achieved by impurity doping of the semiconductor, and can be exploited for waveguide design. Near the bandgap of the semiconductor, there are additional strong index changes with variations in the

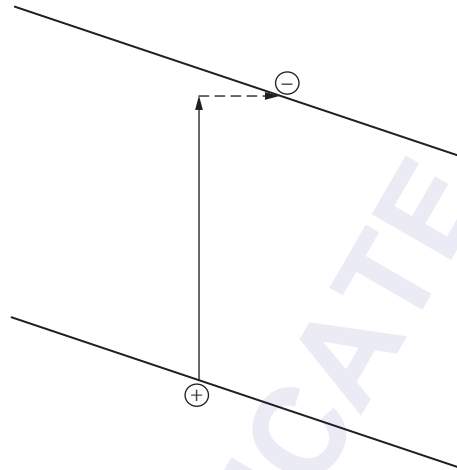


FIGURE 5 Franz-Keldysh effect. Electron can complete transition to the tilted conduction band by tunneling.

carrier density that arise from the associated dramatic changes in the optical loss or gain. Since these correspond to changes in the imaginary index, the Kramers-Kronig relation dictates that changes also occur in the real index. In the spectral region for gain in the semiconductor, these changes are comparable in magnitude and of the same sign as the free-carrier index contribution. These changes are responsible for chirping in the output of modulated semiconductor lasers, and also strongly influence their sensitivity to feedback. They can dramatically impact the static and dynamic performance of semiconductor devices and PICs that employ gain media.

In addition to the effects described above arising from changing carrier populations, the electronic transitions that generate the free carriers can be modified by an applied electric field. For optical frequencies close to these transitions, this can give rise both to electroabsorption and to an enhanced electro-optic effect, which shall be termed electrorefraction, due to the Stark effect on the carrier-generating transitions. In bulk material, the induced absorption is termed the Franz-Keldysh effect,²⁷ and can be viewed as a tunneling effect. For an electron in the valence band with insufficient energy to complete a transition to the conduction band, the field can be viewed as adding a potential to the bands that effectively tilts them in space as shown in Fig. 5. If the excited carrier also traversed a short distance down-field from its initial location, it would have sufficient energy to complete the transitions. This distance depends on the tilt, and thus the strength of the applied field. Since carriers can not be precisely localized according to the Heisenberg uncertainty principle, there is a finite amplitude for completing the transition that is an increasing function of electric field. For fields on the order of 10^5 V/cm, absorption values of ~ 100 cm^{-1} can be readily achieved in material that is quite transparent at zero field. According to the Kramers-Kronig relations, in addition to the absorption described above, this will also induce a change in the real index that will be positive below the band edge and can be used for phase modulation.

In the case of quantum wells, carrier-induced effects can be enhanced due to the forced carrier proximity arising from confinement in the wells. *Excitonic* effects, resulting from the coulombic attraction between electrons and holes, produce sharp features in the absorption spectrum near the band gap that can be dramatically altered by an applied electric field. This quantum-confined Stark effect (QCSE) can enhance both electroabsorptive and electrorefractive effects.^{28,29} This suggests that more compact, lower-drive voltage devices are possible when quantum wells are employed, a fact that has been confirmed experimentally. However, in both the bulk and especially the quantum well case, care must be taken to operate at an optical frequency where strong electroabsorptive or electrorefractive effects are operative but the zero-field background absorption is

not prohibitively high. Another issue that impacts the design of devices based on electroabsorption is the requirement for removal of the photogenerated carriers to prevent screening of the applied field and band-filling which result in saturation. While the vast majority of QCSE devices have been demonstrated in compound semiconductors due to their direct bandgap, recent work has demonstrated that the effect is also operative in quantum wells at the direct bandgap of indirect gap materials such as Ge, providing impetus for further studies of QCSE device applications in the silicon-germanium materials system.³⁰

Thermal Effects In virtually all materials, the optical path length of a given section of waveguide will increase with temperature. This is the combination of both the physical expansion of the material and the change of the index of refraction with temperature. While both are significant, in most integrated-optic materials the latter effect is dominant. In SiO₂ on Si, for example, this mechanism provides a useful means of index change, and numbers on the order of $\Delta n \sim 10^{-5}/^\circ\text{C}$ are observed. This effect has been used to form a variety of thermo-optic switches and filters, but a significant disadvantage for many applications is the inherent slow speed and high power dissipation. In semiconductors, this index change is more than an order of magnitude larger, and leads to frequency shifts in filter devices and single-longitudinal-mode laser of $\Delta f \sim 10 \text{ GHz}/^\circ\text{C}$.

Nonlinear Effects Another class of index changes results from the nonlinear optical effects caused by the incident optical fields themselves. This is treated in depth in other portions of this *Handbook*,³¹ but two phenomena will be mentioned here. The first is closely related to the electro-optic effect discussed earlier, where now the applied field giving rise to the index change is no longer “quasi-static” but is in fact the optical field itself. The response of the medium is in general not the same at optical frequencies, but the same symmetry arguments and contracted tensor approach is employed.

An exemplary phenomenon related to this kind of nonlinearity is second harmonic generation. The polarization resulting from the incident field at ω multiplied by the index oscillating at ω generates second harmonic components at 2ω . This frequency doubling can occur in waveguides, but great care must be taken to achieve *phase matching* where each locally generated frequency-doubled field propagates in such a way as to add coherently to frequency-doubled fields generated farther along the guide. This requires either that the dispersive properties of the materials and guide geometry allow $n(\omega) = n(2\omega)$, or that the frequency-doubled light radiates away from the guide at an angle to allow phase matching of the z component of the wavevector, or that a periodic domain reversal be introduced into the crystal to allow phase matching. This latter approach, often referred to as quasi-phase matching, has generated considerable interest recently. In this approach, the optic axis of the crystal is periodically reversed with a period equal to $\lambda/2$ the difference in the index of refraction of the fundamental and second harmonic. To date, the most promising results are in LiTaO₃, LiNbO₃, and KTP. In LiNbO₃, periodic domain reversal has been obtained by the application of 100 μs pulsed electric fields of 24 kV/mm using a 2.8 μm -period segmented electrode that is subsequently removed.³² The domain reversal in LiTaO₃ can be obtained on a few-micron scale by proton exchange or electron bombardment. KTP has a higher nonlinear coefficient, but the material is not as well developed. Lower efficiencies have been obtained.

This nonlinear mechanism can be used not just for second harmonic generation, but also for sum and difference frequency generation using input signals of different frequencies.

A second application of nonlinear optics to integrated structures involves a higher order of nonlinearity referred to four-wave mixing, or in some situations as self-phase modulation. The change in index in these cases arises from the *product* of two optical fields. If all fields are the same frequency, this is termed *degenerate*, and if only one field is present, it becomes self-phase modulation with the index change driven by the *intensity* of the incident wave. This nonlinearity has been of interest in research aimed at creating all-optical logic devices. Here the intensity of either the input signal or a separate gating signal can determine the output port of a Mach-Zehnder or directional coupler switch, for example.³³ Recent work has also shown that parametric amplification can be accomplished in guided-wave structures when the waveguide is properly constructed to allow phase matching of the resulting fields,³⁴ suggesting the potential for compact, low-noise amplification in new regions of the spectrum.

21.4 INTEGRATED OPTICS MATERIALS AND FABRICATION TECHNOLOGY

Ion-Exchanged Glass Waveguides

Passive integrated optic devices can be fabricated in certain glass substrates using the ion-exchange technique.^{35,36} In this fabrication process, a sodium-rich glass substrate is placed in a mixture of molten nitrate salts containing alkali cations, such as Cs^+ , Rb^+ , Li^+ , K^+ , Ag^+ , and Tl^+ . During this process, sodium ions at the surface of the host glass are replaced with the alkali cations, resulting in a local increase in the refractive index. Channel waveguides are realized by selectively masking the glass surface. The index change and the diffusion depth are a function of host material composition, the specific alkali cation being used, the temperature, and the diffusion time. The exchange process can be substantially enhanced by applying an electric field across the wafer while the substrate is immersed in the salt bath.

Multimode devices are typically fabricated using thallium ion exchange in borosilicate glasses.³⁷ The high polarizability of the thallium ions results in a large index change (>0.1) while providing low propagation losses (0.1 dB/cm). However, thallium-sodium ion exchange has two significant drawbacks. Thallium is extremely toxic, and it also has a large ionic radius compared to sodium (1.49 Å compared to 0.95 Å), resulting in low diffusion coefficients and low mobility. It is therefore necessary to process the waveguides at high bath temperatures approaching the glass transition temperature of the host material (500°C) for long diffusion times (10 hours) with high applied electric fields (>100 V/cm) to achieve deep multimode waveguides that efficiently couple to commercially available multimode fiber (50 to 100 μm core size). Finding suitable masking materials is a challenge.

Single-mode devices are typically realized using either Ag^+-Na^+ , K^+-Na^+ , or Cs^+-K^+ exchange.^{29,30} The first two processes have been extensively studied and are well understood; however, they each appear to have drawbacks. Ag^+-Na^+ exchanged waveguides are somewhat lossy (0.5 dB/cm) due to a tendency for silver reduction in the host material. K^+-Na^+ exchanged waveguides are typically highly stressed and prone to surface scattering that increases the propagation loss. Although not as extensively studied, Cs^+-K^+ exchanged waveguides show great promise. These waveguides are nearly stress free, low loss (<0.1 dB/cm), reproducible, and can be buried using a two-step exchange process. The two-step process further reduces the propagation loss and results in efficient fiber-waveguide coupling (<0.1 dB loss per interface).

Thin Film Oxides

In recent years there has been substantial interest in IO devices fabricated in thin-film dielectrics on silicon substrates. This is due in part to the excellent surface quality, large-area wafers and mechanical integrity of silicon itself. However, this interest also stems in some cases from the availability of mature silicon-processing technology developed by the electronic integrated circuit industry. IO technology on silicon substrates is usually carried out in SiO_2 , and there are two generic approaches to the Si/SiO_2 fabrication that have proven capable of producing very high performance IO devices. IO devices using both approaches are characterized by waveguides that are extremely low loss and are easily matched in mode characteristics to optical fibers used for transmission, thereby providing very efficient coupling.

The first approach borrows more from the technology of optical fiber manufacture than it does from the Si electronics industry.³⁸ Using a technique known as flame hydrolysis (FHD), a “soot” of SiO_2 is deposited on a Si wafer to a depth of 50 to 60 μm , followed by a thinner layer of a $\text{SiO}_2/\text{GeO}_2$ mix to form what will become the high-index waveguide core. This material is consolidated at $\sim 1300^\circ\text{C}$ for several hours down to roughly half its original thickness, and then the waveguide core layer is patterned using reactive ion etching to form square cross-section waveguide cores. Then FHD is again used, followed by more consolidation, to form the upper cladding layers. Typical index differences for the core material are in the range of 0.25 to 0.75 percent, with core dimensions of 6 to 8 μm square.

A measure of the material quality that was available using this approach is given by some of the extraordinary early devices results obtained. IO power splitters were fabricated to sizes of 1×128 using seven stages and a total of 127 Y-branch 1×2 splitters and a total device length of 5 cm. Total *fiber-to-fiber* excess loss for this device was 3.2 dB with a standard deviation of 0.63 dB.³⁹ A large variety of devices has been made using this technique.

Another technique for Si/SiO₂ fabrication employs film-deposition technology borrowed from silicon electronics processing.⁴⁰ First a base SiO₂ layer is deposited using high-pressure steam to a thickness of $\sim 15 \mu\text{m}$ to prevent leakage to the high-index Si substrate. The waveguide and cladding layers are deposited using low-pressure chemical vapor deposition, either from silane and oxygen, or from tetraethylorthosilane and ammonia. Phosphine is added to increase the index, with guide cores typically containing 6.5 to 8 percent phosphine. The wafer is usually annealed at 1000°C to relieve strain and to densify the films. Waveguide losses below 0.05 dB/cm being reported⁴¹ using this technique, and a large variety of devices have been demonstrated using this approach, including splitters, couplers, and WDM devices. One of the interesting features of this approach to fabrication is that it readily lends itself to the inclusion of other thin films common in Si processing. One such film is Si₃N₄ and this has been used as a high index core for waveguides with much larger core-cladding index step. Such waveguides can generate tightly confined modes that are a much closer match to the modes commonly found in active semiconductor components such as lasers. This feature has been used in a novel mode converter device⁴² that adiabatically transforms from the smaller mode into the larger, fiber-matched mode commonly employed in Si/SiO₂ IOCs.

In some instances, slow response *active* devices have been fabricated in Si/SiO₂ technology using thermal effects to achieve local index changes in one arm of a Mach-Zehnder interferometer. This can either be used as a thermo-optic switch or as a tuning element in WDM components. The heating elements in these devices comprises a simple metal-film-resistive heater deposited directly on the upper surface of the wafer.

Another characteristic feature of IOCs in Si/SiO₂ is a degree of birefringence that results from the compressive stress induced in the film by the Si substrate after cooling down from the high-temperature film deposition or consolidation. Typical amounts of birefringence are $n_{\text{TE}} - n_{\text{TM}} = 3 \times 10^{-4}$. This birefringence can cause wavelength shifts with input polarization in WDM components, and techniques to counteract it include stress removal by adding strain-relief grooves in the film, or stress compensation by adding a counteracting stress-inducing film on the surface of the guide, or the explicit introduction of birefringence-compensating waveguide elements.

The Si/SiO₂ technology has matured to a point where it is commonly used in telecommunications for passive devices such as wavelength demultiplexing circuits that will be discussed later in this chapter, and active thermo-optic device circuits continue to be explored.

Silicon Photonics Fabrication and Materials

Recent years have seen a strong increase in research and development of *silicon photonics* where light is guided in silicon layers that are transparent in the important telecommunications window from 1.2 to 1.6 μm .^{43–46} The most common approach to waveguide design directly uses silicon-on-insulator (SOI) technology, which provides for a layer of single-crystal silicon with thicknesses in the 0.1- to 5- μm range, separated from the underlying silicon substrate by a buried oxide (BOX) layer of SiO₂ with a thickness typically in the 1- to 2- μm range. Typical waveguide structures fabricated in the SOI system are shown in Fig. 6.

SOI has become a mainstream VLSI electronics technology offering reduced electrical parasitics for high-performance microprocessor applications, and the most common SOI wafer-fabrication method is the SMARTCUT technique, developed at CEA/LETI in France, and commercialized by Soitec.⁴⁷ This technique works by implanting hydrogen in a first substrate and growing oxide on a second substrate. After the two substrates are wafer bonded, a rapid annealing step expands the hydrogen, creating an accurate cleave plane and leaving a uniform silicon film on top of the oxide. High-quality SOI wafers are now available to the silicon photonics community in sizes up to 12 in (300 mm), and this materials system immediately brings to bear the unprecedented process

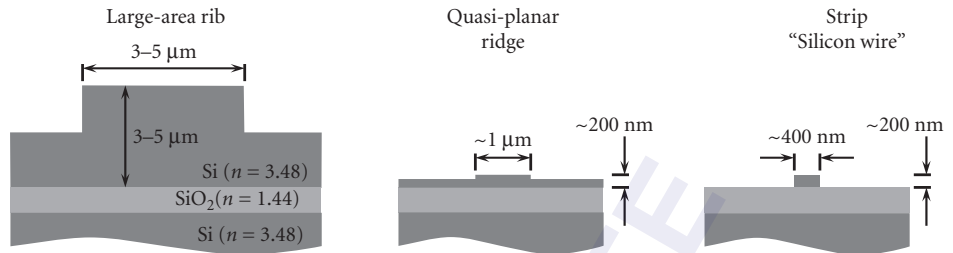


FIGURE 6 Typical silicon-on-insulator (SOI) waveguide structures.

precision associated with VLSI complementary metal-oxide-semiconductor (CMOS) technology. This precision is not just in lithography, etching, and feature size, now at the 65-nm node and still shrinking, but also in the remarkable versatility in materials sequencing that allows for exotic dielectrics and multilayer metallizations.

One of the distinctive features of the SOI system is the extremely high index contrast available for waveguides, for example, with $n_{\text{Si}} = 3.475$ and $n_{\text{SiO}_2} = 1.444$ at $1.55 \mu\text{m}$. As illustrated in Fig. 7, this allows for extremely small waveguides with cross-sectional areas below $0.1 \mu\text{m}^2$.⁴⁶ This feature enables extremely compact passive devices, and also enables high-performance active devices by concentrating the mode on the index-changing region. However, the same high index contrast also leads to significant scattering losses from fabrication-induced waveguide sidewall roughness, which becomes increasingly challenging in very small cross-section waveguides. Larger mode devices in SOI have shown losses in the $\sim 0.1 \text{ dB/cm}$ range,⁴⁸ while maximally vertically confined shallow ridge guides have losses of $\sim 0.4 \text{ dB/cm}$,⁴⁹ and several decibels per centimeter is more typical for most tightly confined “wire” waveguides.

The principle active index-changing mechanism used in silicon photonics is the plasma index contribution from free electrons and holes in the silicon. As will be illustrated later in discussions of devices, these free-carrier populations and locations can be controlled in a variety of diode and field-effect structures. The ability to harness the precision of VLSI CMOS technology to predictably

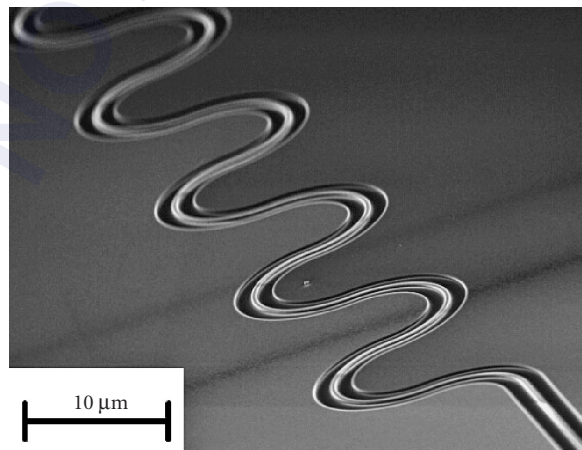


FIGURE 7 Example of losses for high-index-contrast “wire” waveguides: For $6.5\text{-}\mu\text{m}$ radius bends, losses are 0.0043 dB per 180° turn as demonstrated by Vlasov et al., IBM.⁸⁶

engineer very high confinement waveguides together with precision dynamic control of free carriers in transistor-like device geometries has recently resulted in a number of high-performance modulator designs. While historically silicon has not been viewed as a high-performance active optical material, these and other advances have caused a serious reexamination of the potential of silicon as a powerful medium for PICs.

LiNbO₃ and LiTaO₃

The majority of the integrated optics R&D from 1975 to 1985 and the majority of the currently commercial integrated optics product offerings utilize LiNbO₃ as the substrate material. A number of excellent review papers detailing R&D efforts in LiNbO₃ are available.⁵⁰⁻⁵⁴ LiNbO₃ is an excellent electro-optic material with high optical transmission in the visible and near infrared, a relatively large refractive index ($n = 2.15 - 2.2$), and a large electro-optic coefficient ($r_{33} = 30.8 \times 10^{-10}$ cm/V). Probably most important, but frequently overlooked, is the widespread availability of high-quality LiNbO₃ wafers. Hundreds of tons of LiNbO₃ is produced annually for the fabrication of surface acoustic wave (SAW) devices. This large volume has resulted in well-developed crystal growth and wafer-processing techniques. In addition, LiNbO₃ wafers are at least an order of magnitude less expensive than they would be if integrated optics was the only user of this material. High-quality 3- and 4-in optical-grade LiNbO₃ wafers are now available from multiple vendors.

LiNbO₃ is a uniaxial crystal which is capable of supporting an extraordinary polarization mode for light polarized along the optic axis (z axis) and an ordinary polarization mode for light polarized in the x - y plane. LiNbO₃ is slightly birefringent with $n_e = 2.15$ and $n_o = 2.20$. LiNbO₃ devices can be fabricated on x -, y -, and z -cut wafers. Phase modulators, fiber gyro circuits, and Mach-Zehnder interferometers are typically fabricated on x -cut, y -propagating wafers, and operate with the TE (extraordinary) mode. Push-pull devices, such as delta-beta directional coupler switches, are typically fabricated on z -cut, y -propagating wafers and operate with the TM (extraordinary) mode. Both configurations utilize the strong r_{33} electro-optic coefficient. Devices that require two phase-matched modes for operation are typically fabricated on x -cut, z -propagating wafers.

The majority of LiNbO₃ integrated optic devices demonstrated to date have been fabricated using the titanium in-diffusion process.⁵⁵ Titanium strips of width 3 to 10 μm and thickness 500 to 1200 \AA are diffused into the LiNbO₃ at 950 to 1050°C for diffusion times of 5 to 10 hours.^{56,57} The titanium diffusion results in a local increase in both the ordinary and extraordinary refractive indices so that both TE and TM modes can be supported for any crystal orientation. Titanium thickness and strip width typically need to be controlled to ± 1 percent and $\pm 0.1 \mu\text{m}$, respectively, for reproducible device performance. Due to the high processing temperatures that approach the Curie temperature of LiNbO₃, extreme care must be taken to prevent Li₂O out-diffusion^{58,59} and ferroelectric domain inversion, both of which significantly degrade device performance. Photorefractive optical damage⁶⁰ also needs to be considered when utilizing Ti-diffused devices for optical wavelengths shorter than 1 μm . Optical damage typically prevents the use of Ti-diffused devices for optical power greater than a few hundred micro Watts at 800-nm wavelength, although the problem can be reduced by utilizing MgO-doped LiNbO₃ wafers. Optical damage is typically not a problem at 1300 and 1550 nm for optical powers up to 100 mW.

An alternative process for fabricating high-quality waveguides in LiNbO₃ is the annealed proton exchange (APE) process.^{61,62} In the APE process, a masked LiNbO₃ wafer is immersed in a proton-rich source (benzoic acid is typically used) at temperatures between 150 and 245°C and times ranging from 10 to 120 minutes. The wafer is then annealed at temperatures between 350 and 400°C for 1 to 5 hours. During the initial acid immersion, lithium ions from the wafer are exchanged with hydrogen ions from the bath in the unmasked region, resulting in a stress-induced waveguide that supports only the extraordinary polarization mode. Proton-exchanged waveguides that are not subjected to further processing are practically useless due to temporal instabilities in the modal propagation constants, high propagation loss, DC drift, and a much-reduced electro-optic coefficient. However, it has been demonstrated⁶² that proper postannealing results in extremely high quality waveguides that are stable, low-loss, and electro-optically efficient.

The APE process has recently become the fabrication process of choice for the majority of applications currently in production. Since the APE waveguides only support the extraordinary polarization mode, they function as high-quality polarizers with polarization extinction in excess of 60 dB.⁶³ As described later in this chapter, high-quality polarizers are essential for reducing the drift in fiber optic gyroscopes and minimizing nonlinear distortion products in analog links. APE waveguides exhibit low propagation losses of 0.15 dB/cm for wavelengths ranging from 800 to 1550 nm. APE LiNbO₃ devices exhibit stable performance for optical powers of 10 mW at 800 nm and 200 mW at 1300 and 1550 nm. The APE process can also be used to fabricate devices in LiTaO₃ for applications requiring higher optical powers (up to 200 mW) at 800 nm.⁶⁴ In addition to offering performance advantages, the APE process also appears to be the more manufacturable process. It is relatively easy to scale the APE process so that it can handle 25-wafer lots with no degradation in device uniformity. The fiber pigtailling requirements are also substantially reduced when packaging APE devices since these devices only support a single polarization mode.

After the waveguides have been fabricated in the LiNbO₃ wafer, electrodes need to be deposited on the surface. One- μm -thick gold is typically used for lumped-electrode devices while 5- μm -thick gold is typically used for traveling-wave devices to reduce RF resistive losses. The lift-off process and electron-beam deposition is typically used for lumped-electrode devices while up-plating is typically used for realizing the thicker gold electrodes. Better than 0.5- μm layer-to-layer registration is required for optimum device performance. As shown in Fig. 10, electrodes on x-cut LiNbO₃ are usually placed along side the waveguide so that the horizontal component of the electric field interacts with the propagating TE mode. Electrodes on z-cut LiNbO₃ are placed on top of the waveguide so that the vertical component of the electric field interacts with the propagating TM mode. An SiO₂ buffer layer (0.1–1 μm thick) is required between the optical waveguide and the electrode on all z-cut devices to reduce metal-loading loss. A thick (1 μm) SiO₂ buffer layer is also utilized on some x- and z-cut devices to reduce the velocity mismatch between the microwave and optical waves in high-speed traveling-wave modulators. A thin layer of amorphous silicon is also utilized on some z-cut devices to improve device stability over temperature.

III-V Materials and Fabrication Technology

In this section we will briefly review some of the epitaxial growth and fabrication techniques that are used to make PICs in III-V materials, with a primary focus on InP-based devices.

III-V Epitaxial Crystal Growth The epitaxial growth of III-V optoelectronic materials has evolved during the last several decades from nearly exclusive use of manually controlled liquid-phase epitaxial (LPE) growth to a variety of highly versatile computer-automated vapor and beam-growth techniques. These include atmospheric-pressure and low-pressure metal-organic vapor-phase epitaxy (MOVPE), hydride and chloride vapor-phase epitaxy (VPE), molecular beam epitaxy (MBE), chemical beam epitaxy (CBE), and metal-organic molecular beam epitaxy (MOMBE). Detailed descriptions of reactor design and growth chemistry are beyond the scope of this section, and the interested reader is referred to recent texts and conference proceedings for the most current information.⁶⁵

One of the critical criteria for evaluating crystal growth is the uniformity, both in thickness and in epitaxial composition. Layer thickness changes of several percent can lead to nanometer-scale wavelength changes in grating-based lasers and filter devices. Similarly, compositional changes leading to a 10-nm shift in the photoluminescence peak wavelength of the guide layers, which is not at all uncommon, can also result in nanometer-scale wavelength shifts in distributed feedback (DFB) laser emission wavelengths, in addition to potential undesirable gain-peak mismatches that may result from the λ_{pl} shift itself.

Proper reactor geometry, sometimes with substrate rotation, have been shown capable of percent-level uniformity both in MOVPE and in the beam techniques. One difficulty associated with latter lies in the ballistic “line-of-sight” growth which prevents regrowth over reentrant mesa geometries or overhanging mask surfaces often encountered in PIC and laser fabrication, while

MOVPE and especially VPE offer outstanding coverage over a wide range of morphologies. Other criteria to be considered are the doping capabilities. The lower growth temperatures associated with MBE, CBE and MOMBE enable very abrupt changes in doping level, and highly concentrated doping sheets that are desirable for high-speed transistors in OEICs, for example. Both the vapor and beam techniques have successfully grown semi-insulating Fe-doped InP, a material that is playing an increasingly pivotal role in photonic devices.

The typical PIC processing involves the growth of a base structure that is followed by processing and regrowths. During both the base wafer and regrowths, selective area growth is often employed where a patterned dielectric film is used to prevent crystal growth over protected areas. This film is typically SiO_2 or Si_3N_4 deposited by CVD or plasma-assisted CVD. This technique is readily used with MOVPE, but care must be taken to keep a substantial portion of the field open for growth to avoid the formation of polycrystalline deposits on the dielectric mask. Caution must be exercised during regrowths over mesas or other nonplanar geometries, as well as in the vicinity of masked surfaces. Gross deviations from planarity can occur due to overshoots of crystal growth resulting from crystal-orientation-dependent growth rates on the various exposed surfaces.

III-V Etching Technology A fundamental step in III-V PIC processing is mesa etching for definition of the optical waveguides. This is usually accomplished by patterning a stripe etch mask on a base wafer that has a number of epitaxial layers already grown, and removing some of the layers in the exposed regions to leave a mesa comprised of several of the epitaxial layers. The etching process can either be a “wet” chemical etchant, or a “dry” plasma-type etch.

Wet etching refers to the use of an acid bath to attack the unprotected regions of a surface. The acids that are commonly used to etch III-V materials⁶⁶ also tend to significantly undercut a photoresist pattern, and hence photoresist is usually used only in broad-area features or in shallow etches where undercutting is not a concern. For precise geometries such as waveguide stripes, another masking material such as SiO_2 or Si_3N_4 is first deposited and patterned with photoresist and plasma etching, or HF etching (for SiO_2).

In some instances it is required that the etchants be nonselective, uniformly removing layers regardless of composition. This is usually the case when etching a mesa through a multilayer active region to form a buried heterostructure laser. Br-based etchants, such as bromine in percent-level concentration in methanol, tend to be very good in this regard. This etchant, along with many of the nonselective etchants, will form a reentrant 54.7° (111A) face mesa for stripes along the (011) direction (with a nonundercutting mask) and will form an outward-sloping 54.7° walled mesa for stripes along the (01 $\bar{1}$) direction. Other etchants, with varying degrees of nonselectivity and crystallographic behavior, include mixtures of HBr, CH_3COOH , or HCl, CH_3COOH , and H_2O_2 .⁶⁷

In fabricating precise geometries in III-V integrated optic or PIC devices, it is often desirable to remove specific layers while leaving others, or control mesa heights to a very high degree of precision. The combination of material-selective etchants and the inclusion of special etch-stop layers offers a convenient and precise means of achieving this. Hundred-Å-thick layers of InGaAsP can easily halt InP etches even after many microns of etching. Extensive compilations have been made of etches for the InP-based compounds,⁵² and the most common selective InP etches are HCl-based. Typical mixtures are HCl and H_3PO_4 in ratios ranging from 3:1 to 1:3, with the lower HCl content leading to less undercutting and slower etch rates. The HCl-based etchants are highly crystallographic in nature,⁶⁸ and can produce mesas with nearly vertical walls or outward sloping walls, depending on the mesa stripe orientation.

A common selective etch for removing InGaAsP or InGaAs while only weakly attacking InP are mixtures of H_2SO_4 , H_2O_2 , and H_2O in a ratio of X:1:1 with X typically ranging from 3 to 30. Selectivities in the range of 10:1 and typically much higher are readily achieved. Other selective etchants for InGaAsP are based on HNO_3 or mixtures of KOH, $\text{K}_2\text{Fe}(\text{CN})_6$, and H_2O .

Dry etching techniques, such as reactive ion etching (RIE) or other variants such as chemically assisted reactive ion beam etching (CAIBE), also play a key role in III-V PIC processing. These have often been carried out using Cl_2 -based mixtures with O_2 and Ar,⁶⁹ while in other cases the reactive chlorine is derived from compounds such as CCl_2F_2 . Excellent results have also been obtained with methane/hydrogen mixtures or ethane/hydrogen.⁷⁰ In these latter cases Ar is also often used as a

sputtering gas to remove interfering redeposited compounds. Reactive ion etching has been used both to form mesa and facet structures as well as in transferring grating patterns into semiconductors through an etch mask.

The appeal of reactive ion etching is the lack of mask undercutting that can usually be achieved, allowing very high lateral precision with the promise of reproducible submicron mesa features. In addition, the ability to create vertical wall-etched facets through a variety of different composition epitaxial layers suggests the possibility of integrated resonator or reflecting and coupling structures without the use of gratings. This approach has been used to form corner reflectors,⁷¹ square-geometry ring-resonators,⁷² and a variety of complex waveguide patterns using beam splitters.⁷³ Another recent application has been the use of etched-facet technology to create gratings, not as an interfacial corrugation *along* the waveguide, but as a grating in the other dimension *at the end surface* of a waveguide for two-dimensional “free-space” grating spectrometers.^{74,75}

Grating Fabrication Many of the PICs employ corrugated-waveguide grating-based resonators or filters, and the most common technique for fabricating these gratings involves a “holographic” or interferometric exposure using a short wavelength laser source. Here a thin (typically 500 to 1000Å thick) layer of photoresist is spun on a wafer surface and exposed with two collimated, expanded beams from a blue or UV laser at an appropriate angle to form high-contrast fringes at the desired pitch. Since the illuminating wavelength is precisely known, and angles are easily measured in the milliradian range, the typical corrugation in the 2000Å-period range can be fabricated to armstrong-level precision in period. The resist is developed and then functions as an etch mask for the underlying layers. This etching can be either a wet etch (commonly using *HBr*-based etchants), or a dry reactive ion etch. Commonly used lasers are HeCd at 325 nm or one of the UV lines of an Argon ion laser at 364 nm. Electron-beam lithography has also been successfully applied to the generation of gratings for III-V integrated optic devices.

Active-Passive Transitions Compound semiconductors are appealing for PICs in large part due to their ability to emit, amplify, and detect light. However, waveguide elements that perform these functions are not low loss without excitation, and are generally not suitable for providing passive interconnections between circuit elements. One of the most fundamental problems to overcome is the proper engineering and fabrication of the coupling between active waveguides, containing lower bandgap material, and passive waveguides composed of higher bandgap material.

Most PICs demonstrated to date have employed some form of butt-coupling, where an active waveguide of one vertical and/or lateral structure mates end-on with a passive waveguide of a different vertical and/or lateral structure. Butt-coupling offers design simplicity, flexibility, and favorable fabrication tolerances. The most straightforward approach for butt-coupling involves the selective removal of the entire active waveguide core stack using selective wet chemical etching, followed by a regrowth of a mated, aligned passive waveguide structure. The principal advantage of such an approach is the independent selection of compositional and dimensional design parameters for the two guides.

Another approach to butt-coupling, often called “offset quantum wells” or “active layer removal,” employs a largely continuous passive waveguide structure with a thin active layer (usually of quantum wells) residing on top, which is selectively removed on the portions of the structure which are to be passive. Using material-selective wet chemical etches, the thin MQW stack can be removed with very high reproducibility and precision, and the dimensional control is thus placed in the original computer-automated MOVPE growth of the base wafer. The removal of the thin active layer constitutes only a small perturbation of the continuous guide core constituted by the lower, thicker layer, and efficient coupling can be achieved.⁷⁶

Two powerful alternatives to the butt-coupling variants discussed above are *selective area epitaxy* and *quantum-well intermixing*. These techniques provide for longitudinal modification in the effective bandgap of quantum-well-containing waveguides by altering the dimensions or profile of the quantum well along the length of a waveguide. Selective area epitaxy accomplishes this by the inclusion of growth-inhibiting masks, such as SiO₂, laterally adjacent to the waveguide during MOVPE growth. The resulting local increase in vapor-phase reactants, combined with surface diffusion along

the growth-inhibiting mask, leads to an increased growth rate in regions adjacent to the lateral masks. This leads to thicker quantum wells with correspondingly lower bandgaps from quantum confinement, making these regions suitable for gain media or absorption along a waveguide which is transparent in regions with no masks. By controlling the degree of enhancement, regions of transparency, regions suitable for electroabsorption or QCSE, and regions suitable for gain or absorption can all be formed along a guide with remarkably smooth continuity between regions. This technique has been highly successful in PICs comprising an integrated DFB laser and electroabsorption modulator.⁷⁷

The quantum-well intermixing accomplishes a similar end by implanting regions and diffusing impurities or vacancies into the quantum-well region of a waveguide to intermix the boundaries of the wells and barriers. This intermixing smoothes the rectangular well into a rounded profile and also can change the average composition of the well, leading to a higher bandgap which can be used for longitudinally differentiating gain regions, modulation regions, and transparent regions.⁷⁸ Although not in as wide commercial use, this method has the advantage that the various regions are large-area materials modifications allowing for standard materials characterization techniques, in contrast to the selective area epitaxy approaches which require difficult characterization of narrow stripe regions.

Yet another approach to coupling between two different waveguides employs directional coupling in the vertical plane between epitaxial layers serving as the cores of the two distinct waveguides. This type of vertical coupling can either be accomplished using the principle of intersecting dispersion curves, or through the use of a corrugated-waveguide grating to achieve phase matching. Vertical coupler structures may be useful for wide-tuning applications, since a small change of effective index for one mode can lead to a large change in coupling wavelength.⁷⁹

Organic Polymers

Polymer films are a relatively newer class of materials for integrated optics.⁸⁰ Polymers offer much versatility, in that molecular engineering permits many different materials to be fabricated; they can be applied by coating techniques to many types of substrates, and their optical and electro-optical properties can be modified in a variety of ways. Applications range from optical interconnects, in which passive guides are used in an optical PC board arrangement, to equivalents of IOCs and OEICs. Polymer devices are also being explored for third-order nonlinear applications.

Numerous methods for fabricating polymer waveguide electro-optic devices have been reported. One attractive technique consists of spin-coating a three-layer polymer sandwich over a metal film, often on a semiconductor (Si) substrate. The three polymer layers form a symmetric planar waveguide; the middle layer is electro-optic, due to the presence of a guest molecule that imparts the electro-optic property, or the use of a side-chain polymer. The sample is overcoated with metal and the entire structure is heated near the glass transition temperature and poled at an electric field of typically 150 V/ μm . The poling aligns the nonlinear molecules in the middle polymer layer, thereby inducing the Pockels effect and a birefringence. Typical values of index and birefringence are 1.6 and 0.05, respectively. Electro-optic coefficients are in the 16 to 38 pm/V range. Channel waveguides are subsequently defined by a variety of methods. An attractive technique is photobleaching, in which the waveguide region is masked with a metal film and the surrounding area exposed to UV light. This exposure alters the molecules/linking in the active layer, thereby reducing the refractive index and providing lateral confinement. Losses in such guides are typically in the 1dB/cm range.

The basic IO modulators have been demonstrated in a variety of polymers. Of particular note are traveling wave modulators with a 3-dB bandwidth of 40 GHz and a low-frequency V pi of 6 V,⁸¹ and work illustrating the high-frequency potential of polymers up to modulation frequencies of 110 GHz.⁸² Relative to LiNbO₃, polymer modulators can have higher overlap factors because the lower metal layer provides vertical, well-confined signal fields. However, the relatively low index of polymers and their comparable electro-optic coefficient to LiNbO₃ implies a lower electro-optic efficiency. Polymers do provide a better velocity match of optical and electromagnetic velocities which can result in very high frequency performance as described above.

For polymers to fulfill their potential, a number of material and packaging issues must be addressed. First, it is highly desirable to develop polymers that overcome the long-term relaxation of the electro-optic effect typical of many of the early reports. Development of polymers with transition temperatures in the 300°C range (so they can withstand the temperatures typical of device processing and packaging) is also highly desirable. Work on polyimide is particularly promising in this area. Recent work has demonstrated encouraging results with devices processed at lower temperatures that exhibit long-term stability at elevated operating temperatures of 85°C.⁸³

21.5 CIRCUIT ELEMENTS

Passive Devices

This section provides a phenomenological description of the most common passive and active IO devices. Detailed descriptions of the device theoretical properties as well as typical characteristics can be found in Refs. 2–6.

Passive devices form many of the fundamental building blocks for IOCs and PICs, and passive waveguides also comprise the interconnections between devices. Passive devices are defined as those dielectric waveguide structures which involve neither application of electrical signals, generation or detection of light, nor nonlinear optical effects. This section will focus on the most important structures: waveguide bends, polarizers, fiber-to-chip coupling, power splitters, and filters.

Waveguide bends, such as those illustrated in Figs. 7, 10, 11, 12, 14, 15, 19, 22, 25, and 26, are needed to laterally offset modulators and other devices, and also to increase device-packing density. Waveguide bends lead to radiation or leakage from the guide, and analytical treatments reveal that the induced losses decrease exponentially with radius of curvature in a manner that is very sensitive to the index difference between the core and cladding of the guide.^{84,85} The most widely used bend is based on an S-bend geometry described by a raised cosine function.² This structure minimizes the tendency of light in a dielectric waveguide to “leak” as the guide’s direction is altered by starting with a small bend (large effective bend radius and then increasing the bend rate until the midpoint of the offset, then following the pattern in reverse through the bend completion.

Since the index difference between the guide and surrounding dielectric material is often small (10^{-3} to 10^{-4}) bends must be gradual (effectively a few degrees) to keep losses acceptably (< 0.5 dB) small. In LiNbO₃, offsets of 100 μm require linear distances of typically 3 mm. In semiconductor research device work, designs with high index steps are sometimes used to form small-radius bends and ion beam etching has been utilized to form reflective micromirrors⁷³ at 45° to the guide to create a right-angle bend. While compact, to date these have generally been relatively lossy compared to their bending waveguide counterparts.

As noted earlier, silicon photonics carries a very high index contrast of $\Delta n \sim 2$, which also allows for extremely tight waveguide radii, as illustrated in Fig. 7, without introducing the bending losses that often limit the size of weaker index contrast materials systems.⁸⁶ This can provide a dramatic increase in the density of possible PICs developed in the SOI system.

Polarizers are necessary for polarization-sensitive devices such as many electro-optic modulators and for polarization-sensitive applications such as fiber gyroscopes. Polarizers can be formed on dielectric waveguides that support both TE and TM propagation by forming overlays that selectively couple one polarization out of the guide. For example, a plasmon polarizer formed on LiNbO₃ by overcoating the guide with a Si₃N₄/Au/Ag thin-film sandwich selectively attenuates the TM mode.⁸⁷ In some materials it is possible to form waveguides that only support one polarization (the other polarization is not guided and radiates into the substrate). By inserting short (millimeter) lengths of such guides in circuits or alternatively forming entire circuits from these polarizing guides, high extinction can be obtained. For example, annealed proton exchange waveguides (APE) in LiNbO₃ exhibit polarization extinction ratios of at least 60 dB.⁸⁸

Guided wave devices for splitting light beams are essential for most IOCs. Figure 8 illustrates the two common splitters: a directional coupler and a Y junction. The figure illustrates 3-dB coupling (1X2),

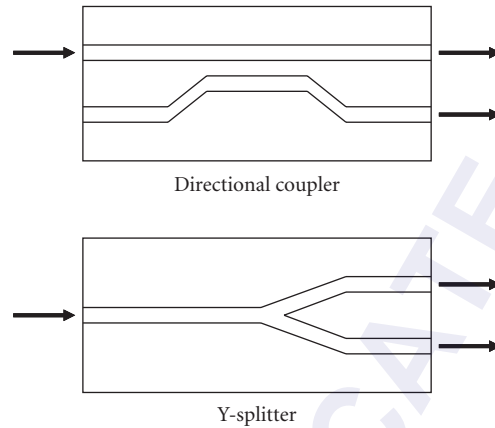


FIGURE 8 Passive directional coupler and Y-branch IO splitter devices.

and by cascading such devices and using variations on the basic designs it is possible to fabricate $N \times N$ structures. IO splitters of complexity 8×8 are commercially available in glass.

The operation of the directional coupler is analogous to the microwave coupler and is described by the coupled mode equations. The coupling strength is exponentially dependent of the ratio of the guide spacing and the effective tail length of the guided mode. Thus when guides are far apart (typically greater than $10 \mu\text{m}$ in weakly guided glass and ferroelectric devices) as in the left-most portion of the structure in Fig. 8, there is negligible coupling. When the guides are close together (typically a few microns), power will couple to the adjacent guide. The fraction of power coupled is sinusoidally dependent on the ratio of the interaction length to the coupling length L_c . L_c is typically 0.5–10 mm and is defined as the length for full power transfer from an incident guide to a coupled guide. The 3-dB coupling illustrated requires an interaction length of half L_c .²⁻⁶ Operation of this device is symmetric; light incident in any one of the four inputs will result in 3-dB splitting of the output light. However, if coherent light is incident on both input guides simultaneously, the relative power out of the two output guides will depend on the phase and power relationship of the incident signals.

The Y splitter illustrated in Fig. 8 operates on a modal evolution principle. Light incident on the junction from the left will divide symmetrically so that the fundamental mode of each output branch is excited. Branching circuit design follows closely from the design of waveguide bends. In low index contrast systems, the Y-junction angle is typically a few degrees and the interaction length is a few millimeter. Operation of this device is not symmetric with respect to loss. If coherent light is incident on both guides from the right, the amount of light exiting the single guide will depend on the power and phase relationship of the optical signals as they enter the junction area. If coherent light is only incident in one arm of the junction from the right it will experience a fundamental 3-dB loss in propagation to the left to the single guide. This is due to the asymmetric modal excitation of the junction. (see the next section, “Active Devices”).

An extremely important issue in integrated optics is the matching of the waveguide mode to the mode of the fiber coupled to the guide. Significant mode mismatch causes high insertion loss, whereas a properly designed waveguide mode can have coupling loss well under 1 dB. As illustrated in Eq. (7), the proper design requires optimizing the overlap integral of the optical fields in the two media. Often some sort of index matching between the two materials is also employed to minimize additional reflective losses. Figure 9 illustrates the issue with mode profiles in the two transverse directions for a Ti indiffused guide. In general the IO mode is asymmetric relative to the fiber mode. It should be noted that the loss obtained on a pigtailed fiber-chip joint is also highly determined by the precision and stability of the mechanical method of attaching the fiber to the chip. Most techniques use some sort of carrier block for the fiber (e.g., a Si V-groove) and attach the block to the IO chip.

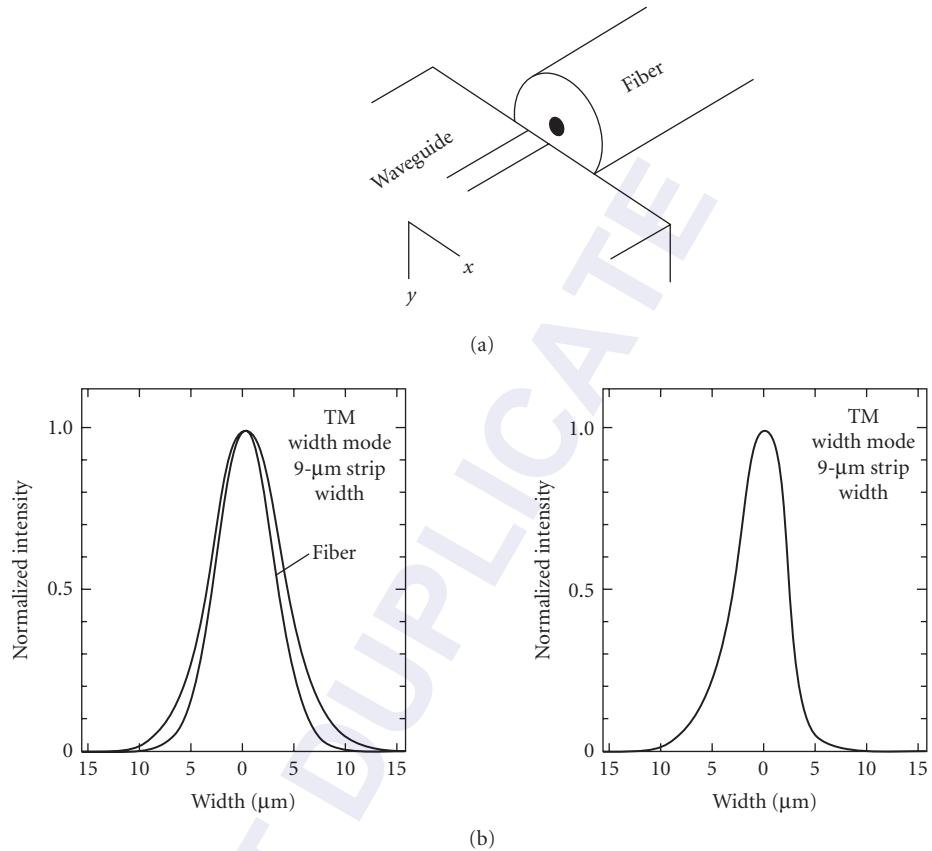


FIGURE 9 Mode-matching illustrated by coupling a Ti-indiffused guide to an optical fiber. Mode profiles are shown both in width and depth for the waveguide.

Performance on commercially available devices is typically <1-dB coupling loss per interface with robust performance over at least the -30 to 60°C range.

Filters are another critical passive IO component used to selectively pass, block, or route predetermined wavelengths or bands of the input optical spectrum. Filters most often use interferometric concepts to achieve the desired effect. For example, a simple unequal-arm Mach-Zehnder filter can be constructed by combining two Y-splitters where the two arms in the interconnecting region have different path lengths as illustrated in Fig. 10. Combining the two paths reveals trivially that the resulting transmission of the device will have an amplitude response that is sinusoidal in frequency, or an intensity transmission given by

$$I(\omega) = I_0 \cos^2 \left(\frac{n\omega(L_2 - L_1)}{2c} \right) \quad (14)$$

By cascading such filters, a large variety of passband characteristics can be realized.

IO filters are important for wavelength division multiplexed (WDM) optical communications, where they are commonly used to combine different frequency channels together in multiplexers and to separate the different frequency channels in demultiplexers. Sophisticated multipath variants of the unequal arm interferometer are commonly used for this purpose, and are referred to as

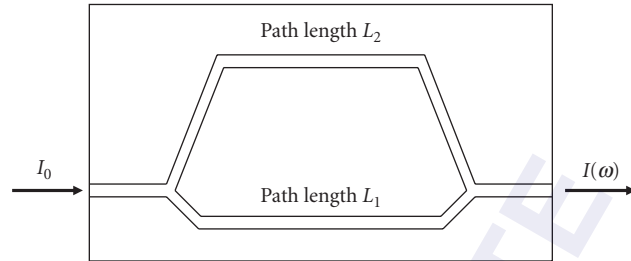


FIGURE 10 Unequal arm Mach-Zehnder interference filter.

arrayed-waveguide gratings (AWGs), waveguide grating routers (WGRs), or phasars.^{89–91} A multiport WGR is shown in Fig. 11 where each input to a primary star coupler expands in the laterally unguided portion of the star to uniformly illuminate each output waveguide of the primary star. The path lengths of each guide in the array between the primary and secondary star are incremented in length by an integral multiple of some base wavelength. At this wavelength, upon arrival at the input to the secondary star, each wave has the same phase relation to its neighbors that it had at the output to the first star coupler and reciprocity demands that this focus the beam, as a phased array, back to the corresponding single waveguide at the secondary star output. However, a slight change in wavelength will produce a phase tilt across this phased array, and the focus will go to a different output waveguide. In this manner, the device operates as a demultiplexer. The incremental length difference between adjacent guides in the region between the two stars functions just as a grating would in a bulk-optic equivalent of this device. The design illustrated here also has useful cyclical mapping between input and output ports that can be useful in wavelength-based switching schemes.

AWG designs have been successfully executed both in the Si/SiO₂ and InP-based technologies with extraordinary performance. Commercially available devices provide out of band rejection in excess of 30 dB, polarization-independent operation, typical channel spacings of 50 GHz or 100 GHz, and fiber-to-fiber insertion losses below 2.5 dB. Research devices have demonstrated as many as 4200 channels with spacings as narrow as 5 GHz.⁹² By properly designing the input ports in the coupler regions, nearly ideal flat bandpass characteristics can be obtained at the expense of a few additional dB of insertion loss. The precision waveguide spacing on the inputs and outputs also enables fiber

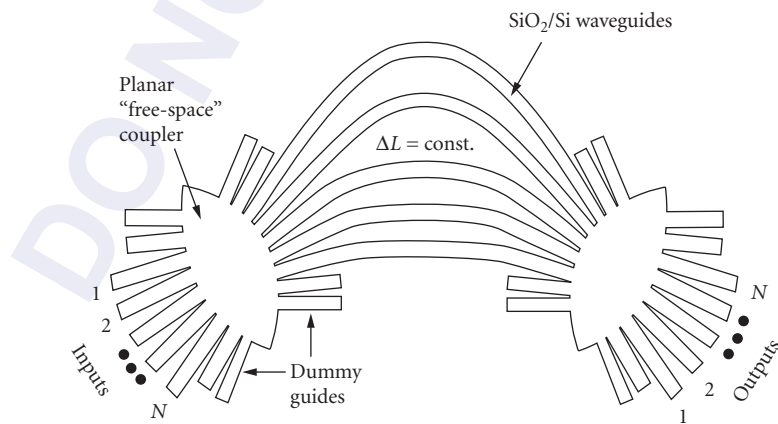


FIGURE 11 Multiport waveguide grating router filter architecture.

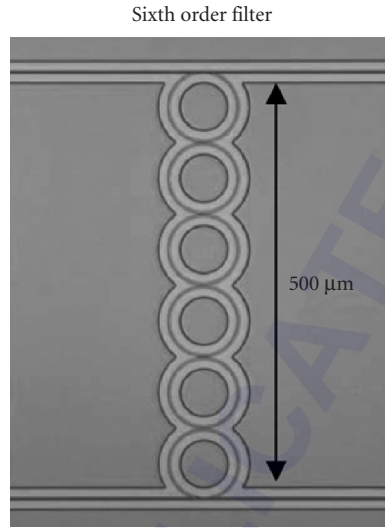


FIGURE 12 Experimental realization of a cascaded ring channel drop filter.

ribbon array connections. IO devices that perform this function have become instrumental for cost-effective deployment of high-channel-count WDM optical communications systems.

Another filter design that is becoming increasingly explored is based upon ring resonators, as illustrated in Fig. 12. Here weak directional couplers are used to couple the ring to the two straight waveguides, allowing coupling between the two straight waveguides only at the sharp resonant frequency of the rings, with other nonresonant frequencies bypassing the rings and remaining in their original input waveguide. As with the unequal-arm Mach-Zehnder filter, rings can be cascaded between two waveguides to provide flat bandpass characteristics with very sharp edges and high out-of-band rejection.⁹³

Another key technology used extensively for filtering in IO is gratings. By periodically modulating the index profile of a waveguide, typically done by corrugating the boundaries of the guide through etching technology, Bragg reflection can be achieved. Such reflectors typically exhibit a relatively flat passband with a width dependent upon the strength and length of the periodic perturbation. Using such a filter inside a laser resonator provides a highly frequency-selective resonator and forms the basis for distributed feedback (DFB) and distributed bragg reflector (DBR) lasers. By cascading sections of grating, reflectors with a periodic sequence of peaks can be achieved. Such designs have been used as the basis for wavelength selective reflectors in widely tunable lasers. Gratings are critical to PICs containing semiconductor lasers because they allow for on-chip resonators without the need for reflecting facets on the chip edge which would constrain the PIC chip size and component layout.

Active Devices

Active IO devices are those capable of having their state altered by an external applied voltage, current, or other stimulus. This may include electro-optic devices, or devices that generate, amplify, or detect light. Active IO devices in non-semiconducting dielectrics generally depend on the linear electro-optic effect, or Pockels effect,²⁴ which produces a change of the index of refraction of a material upon the application of an electric field as discussed earlier. Typical values for a variety of materials is about 10^{-4} for a field of 10^4 V/cm. This results in a phase change for light propagating in the field region and is the basis for a large family of modulators.

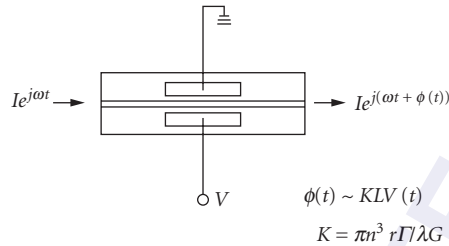


FIGURE 13 Top-down view of a typical LiNbO_3 phase modulator. Field is applied laterally across the guide by surface electrodes on each side.

The fundamental guided-wave modulator is a phase modulator, as illustrated in Fig. 13. In this device, electrodes are placed along side the waveguide and the lateral electric field determines the modulation. In other modulator designs, the vertical field component is used. For the geometry shown, the phase shift is KLV , where K is a constant, L is the electrode length, and V is the applied voltage. For LiNbO_3 , $K = \pi n^3 r_{63} \Gamma / g \lambda$ for the preferred orientation of field along the z (optic) axis and propagation along the y axis. Here, n is the index r_{63} is the electro-optic coefficient, λ is the wavelength, g is the electrode gap, and Γ is the overlap of the electrical and optical fields. In general the value of K is anisotropic and is determined by the electro-optic tensor.²⁴ It should be noted that modulators are often characterized by their V_π value. This is the voltage required for a pi-radian phase shift; in this nomenclature phase shift is written as $\Phi = \pi \cdot V / V_\pi$ where $V_\pi \equiv \pi / KL$. Due to the requirement that the optical field be aligned with a particular crystal axis (e.g., in LiNbO_3 and III-V semiconductors), the input fiber on modulators is generally polarization maintaining.

Modulators in LiNbO_3 typically have efficiencies at $1.3 \mu\text{m}$ of 50%/V, cm, a V_π of 5 V for a 1-GHz 3-dB bandwidth, and a fiber-to-fiber insertion loss of 2 to 3 dB. In semiconductors, modulation efficiencies can be significantly higher if one designs a tightly guided mode (i.e., one well suited for on-chip laser coupling, but having a relatively high fiber-to-chip mismatch coupling loss).

The modulation bandwidth of phase and intensity modulators is determined by the dielectric properties of the electro-optic material and the electrode geometry. For structures in which the electrode length is comparable to or shorter than a quarter RF wavelength, it is reasonable to consider the electrodes as lumped and to model the modulator as a capacitor with a parasitic resistance and inductance. In this case, the bandwidth is proportional to $1/L$. For most IO materials, lumped-element modulators have bandwidths less than 5 GHz to maintain reasonable drive voltages. For larger bandwidths, the electrodes are designed as transmission lines and the RF signal copropagates with the optical wave. This is referred to as a *traveling wave modulator*. The microwave performance of this type of structure is determined by the degree of velocity match of the optical and RF waves, the electrode microwave loss, the characteristic impedance of the electrodes and a variety of microwave packaging considerations. In general, semiconductor and polymer modulators fundamentally have better velocity match than LiNbO_3 and thus are attractive for highest frequency operation. Techniques and structures have been developed to substantially improve the velocity match in LiNbO_3 , however, intensity modulators with 50 GHz have been reported.⁹⁴

To achieve intensity modulation, it is generally necessary to incorporate a phase modulator into a somewhat more complex guided-wave structure. The two most common devices are the Mach-Zehnder and the directional coupler modulator. Figure 14 illustrates the MZ modulator. This device is the guided-wave analog of the classical MZ interferometer. The input and output Y junctions serve as 3-dB splitters and modulation is achieved in a push-pull manner by phase modulating both arms of the interferometer. The interferometer arms are spaced sufficiently that there is no coupling between them. When the applied voltage results in a pi-radian phase shift in light propagating in the two arms when they recombine at the output junction, the resultant odd field distribution corresponds to a second-order mode that cannot be guided and light radiates into the substrate. The output intensity I

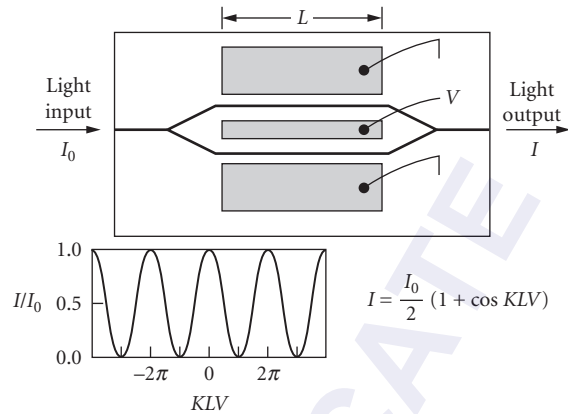


FIGURE 14 Geometry of lumped-element Mach-Zehnder modulator and transfer characteristic.

of this device is given by $I = I_0/2[1 + \cos(KLV)]$. The sinusoidal transfer characteristic is unique in IO modulators and provides the unique capability to “count fringes” by applying drive signals that are multiples of V_π . This feature has been exploited in a novel analog-to-digital converter.⁹⁵ The device can be operated about its linear bias point $\pi/2$ for analog applications and can also be used as a digital switch. A variation on this device is a balanced bridge modulator. In this structure the two Y junctions are replaced by 3-dB directional couplers. This structure retains a sinusoidal transfer characteristic, but can function as a 2×2 switch.

A directional coupler switch is shown in Fig. 15. In the embodiment illustrated, a set of electrodes is positioned over the entire coupler region. The coupler is chosen to be L_c , a coupling length long, so that in the absence of an applied voltage, all light incident in one guide will

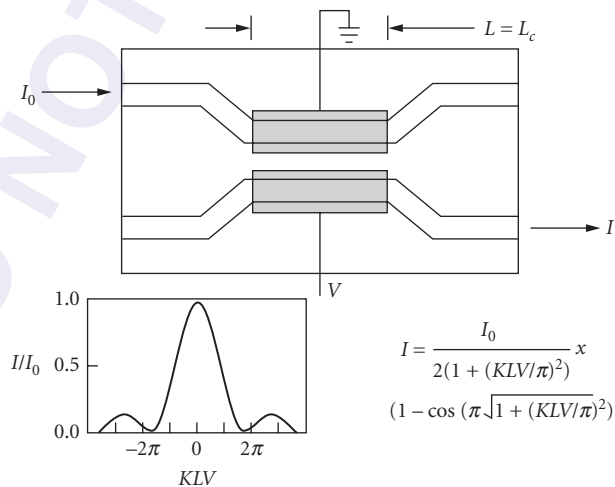


FIGURE 15 Geometry of lumped-element directional coupler switch and transfer characteristic.

cross over and exit the coupled guide. The performance of the directional coupler switch can be modeled by coupled mode theory. The application of an electric field spoils the synchronism of the guides, resulting in reduced coupling, and a shorter effective coupling length. For application of a voltage such that $KLV = \pi\sqrt{3}$, all light will exit the input guide. In general, the transfer characteristic is given by

$$I = \frac{I_0}{2(1+(KLV/\pi)^2)} \left[1 - \cos(\pi\sqrt{1+(KLV/\pi)^2}) \right] \quad (15)$$

Directional coupler switches can also be used for analog or digital modulation. They have also been fabricated in matrix arrays for applications in $N \times N$ switch arrays (see Sec. 21.6). To increase the fabrication tolerance of directional coupler switches, designs based on reversing the sign of index change (Δn) periodically along the coupler have been developed. The most common device consists of a device one to two coupling lengths long and a single reversal of the voltage formed by a two-section electrode.

Both Mach-Zehnder and directional coupler devices have been developed in semiconductors, LiNbO_3 and polymers. Devices are commercially available in LiNbO_3 . Drive voltages and bandwidths achieved are similar to the values quoted above for phase modulators. Additional effort has been focused in LiNbO_3 to make devices that are polarization insensitive so that they are compatible with conventional single-mode fiber.⁹⁶ Mach-Zehnder modulators have also been formed in glass waveguides. Here a resistive pad is heated to vary the index of the waveguide via the thermo-optic effect.

Another important IO component is the TE-to-TM mode converter. This device, illustrated in Fig. 16, depends on an off-diagonal component r_{51} of the electro-optic tensor in LiNbO_3 to convert incident TE (TM) light to TM (TE) polarization. In the converter a periodic electrode structure is used to create a periodic index change along the waveguide to provide phase matching, and thus coupling, between the TE and TM wave. The period Λ of this index change is given by $\Lambda = \lambda / (n_{\text{TE}} - n_{\text{TM}})$. The coupling efficiency at the phase-matched wavelength is given by $\sin^2(\kappa L)$ where $\kappa = \pi n^3 r_{51} E / \lambda$ and E is the applied field. This type of device can be used for polarization control. In a polarization controller, a phase modulator is placed before and after the converter so that signals of arbitrary input polarization can be converted into any desired output polarization. The converter also serves as the basis for a tunable wavelength filter.⁹⁷

There are numerous other types of IO intensity modulators that have been reported. These include a crossing channel switch, a digital optical switch, an acousto-optic tunable wavelength switch, and a cutoff modulator. The first two devices depend on modal interference and evolution effects. The acousto-optic switch utilizes a combination of acoustically induced TE-to-TM mode

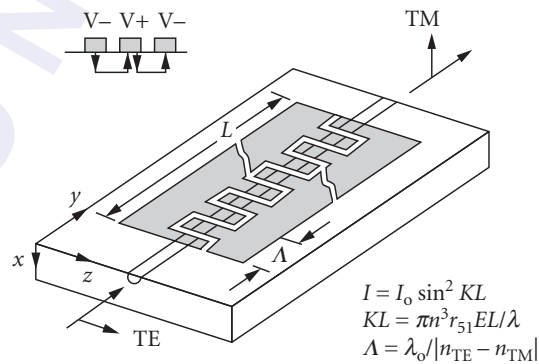


FIGURE 16 TE-TM mode converter using periodic electrodes to achieve phase matching.

conversion and TE-TM splitting couplers to switch narrow-optical-band signals. The cutoff modulator is simply a phase modulator designed near the cutoff of the fundamental mode such that an applied field effectively eliminates the guiding index change between the guide and the substrate. This results in light radiating into the substrate.

In addition to the electro-optic devices described above, another common modulation technique employed in III-V materials employs the electroabsorption or electrorefraction effects discussed in Sec. 21.3. Here the bandgap energy of a bulk medium or an appropriately engineered quantum-well medium is chosen to be somewhat higher than the energy of the propagating photons. An applied field directly induces absorption, or a large index shift associated with the change in absorption at higher energy. The latter effect is used interferometrically in directional couplers, Mach-Zehnder modulators, or other designs as an enhanced substitute for the conventional electro-optic effect. The former is used as a single-pass absorptive waveguide modulator.

To achieve low operating voltages, such modulators are usually designed with small waveguide dimensions for tight confinement of the optical mode. This usually leads to a significant insertion loss of approximately 2 to 3 dB/cm when coupling to optical fibers. However, the tight waveguide mode is very similar to the waveguides employed in semiconductor lasers, and hence one of the primary appeals of waveguide electroabsorption modulators lies in their potential for integration with semiconductor lasers on a single PIC chip.⁹⁸

A particular implementation used by Soda et al.,⁹⁹ is shown schematically in Fig. 17. A 1.55- μm DFB laser structure is mated to an electroabsorption modulator with an InGaAsP core layer having a photoluminescence wavelength of λ_{PL} approximately 1.40 μm . The entire structure uses a buried heterostructure waveguide with semi-insulating InP lateral cladding to provide good current blocking with low capacitance for high modulator bandwidth. Optimization of the modulator core λ_{PL} is very important in this device. With proper design, devices have yielded a good compromise between high output power and high modulator extinction ratios with low voltage drive. Typical device designs exhibit milli watt-level fiber-coupled output power with a -10-dB extinction ratio at drive levels of 2 to 4 V.

Semiconductor lasers comprise one of the most highly developed implementations of integrated optics and are the subject of Chap. 19 in Vol. II, "Semiconductor Lasers." Semiconductor lasers are routinely combined in PICs with other IO elements such as filters, modulators, amplifiers, and detectors. Common to most laser, detector or waveguide PICs is the inclusion of a guide containing an amplifying or gain medium, or an absorptive medium for detection, and the requirements of

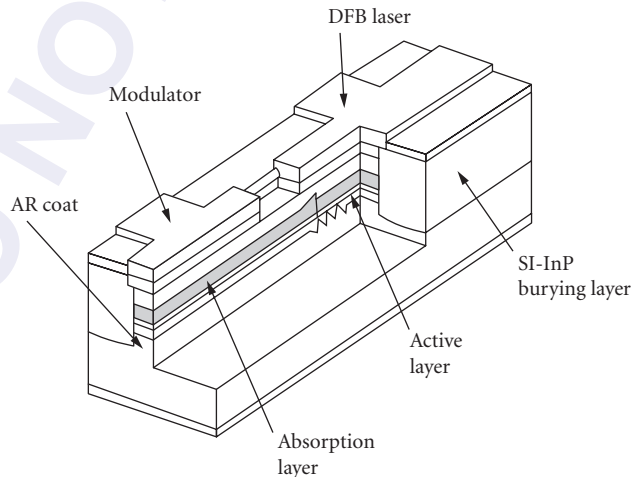


FIGURE 17 Integrated semiconductor laser/electroabsorption modulator PIC.

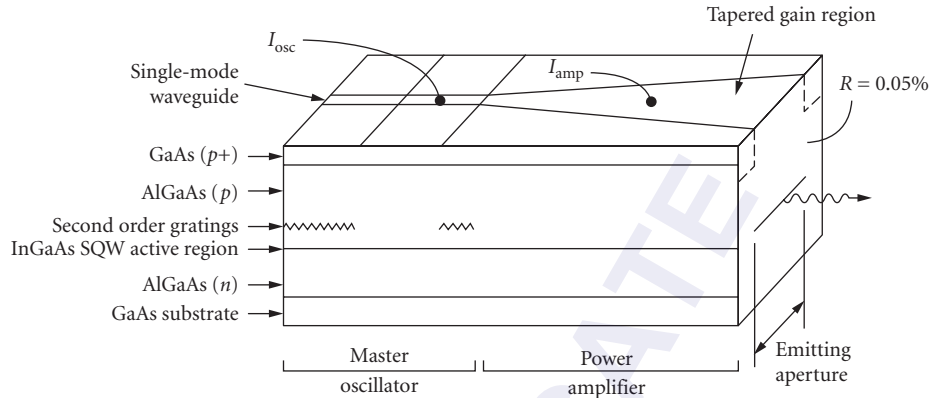


FIGURE 18 Integrated semiconductor master oscillator/power amplifier (MOPA) PIC.

current drive or extraction. The design and processing associated with these guided-wave components relies heavily on a relatively mature technology associated with semiconductor lasers.¹⁰⁰

The gain section of a semiconductor laser is usually fabricated in a buried heterostructure guide, as shown in Fig. 2a, and is driven through a forward biased $p-n$ junction where the layers are usually doped during the crystal growth. With zero or reverse bias, this same structure can function as a waveguide photodetector. In a DFB laser or a distributed Bragg reflector (DBR) laser, this feature can be used to provide an integrated detector, on the back end of the device external to the cavity, for monitoring laser power. Alternatively, a separate gain medium external to the laser cavity can be located external to the cavity on the *output* side to function as an integrated power amplifier for the laser. Such a structure is shown in Fig. 18, where an array of DBR laser is followed by a fan-shaped amplifier to keep the amplifier medium at a relatively constant state of saturation. These PICs are termed master-oscillator/power-amplifiers (MOPAs), and can provide watt-level single-frequency, diffraction-limited output beams from a single chip.¹⁰¹

The challenge of laser integration is to fabricate other guided-wave components without compromising the intricate processing sequence required to make high-performance lasers. Figure 19 shows an early implementation of a sophisticated PIC providing a balanced heterodyne receiver that might be used for coherent optical communications links.¹⁰² Here a tunable local oscillator is tuned to an optical frequency offset by a predetermined amount from one of potentially many incoming signals. The beat signal are generated in the integrated photodetectors, and whose signals can be subtracted for noise reduction, and then electrically amplified, filtered, and fed to a decision circuit. This PIC combines five different types, and a total of seven, guided-wave optical devices: two tunable Bragg reflection filters, an MQW optical gain section, an electrically adjustable phase shifter, a zero-gap directional coupler switch, and two MQW waveguide photodetectors. It also demonstrates self-aligned connections between the buried heterostructure guides, which offer current access and semi-insulating InP lateral current blocking, and the low-loss semi-insulating InP-clad rib guides used in the S-bends and input port. The processing sequence for PICs of this complexity has been described in some detail in the literature, and can be carried out following most of the same steps used in commercial semiconductor laser fabrication.¹⁰³

Tuning of the DBR lasers, as used in the PIC above, is accomplished by injecting current into the (transparent) Bragg reflectors, shifting their index via the plasma and anomalous dispersion effects discussed under Sec. 21.3. This shifts the wavelength of peak Bragg reflectivity, thereby selecting different longitudinal modes for the laser. The laser can also be *continuously* tuned by shifting the frequency of any particular longitudinal mode by injecting current to provide an index shift in the (transparent) phase section of the laser. Detectors in PICs of this type often employ for absorption the same layers used for gain in the laser, and typically have a capacitance of several picofarads dominated by the contact pads rather than the depletion capacitance of the $p-n$ junction.

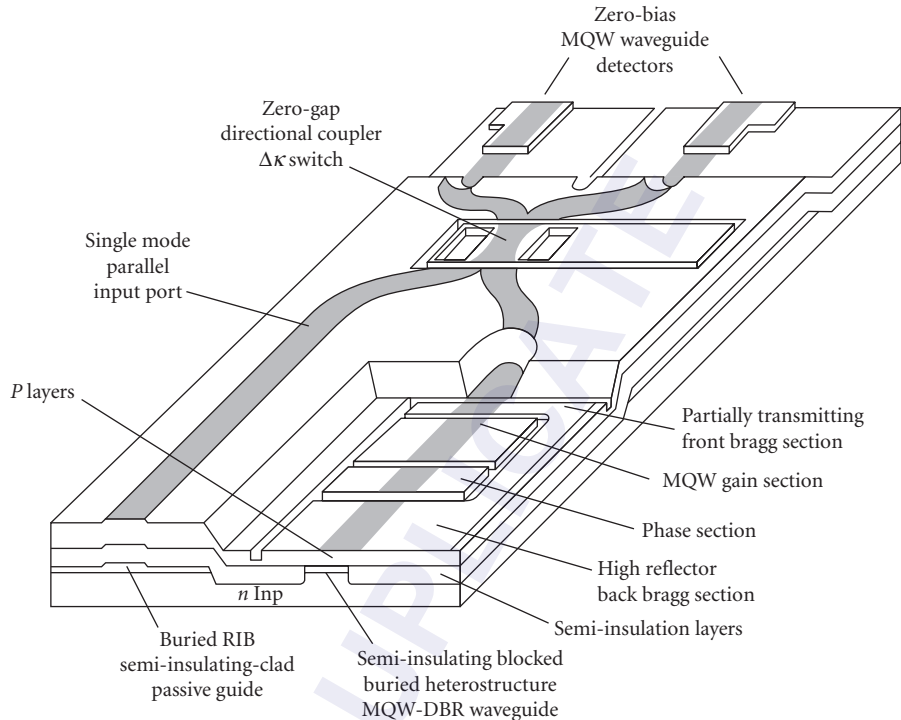


FIGURE 19 Integrated balanced heterodyne receiver PIC.

Early experimental prototypes of PICs of this type have demonstrated total on-chip losses including propagation losses, bending losses, radiation losses at the coupler and at the active/passive detector transitions, and any departures from 100 percent quantum efficiency in the detectors, of approximately 4 dB, providing encouragement for even more complex integrations. This PIC demonstrated error-free heterodyne reception of frequency-shift-keyed digital signals with sensitivities of -40 dBm at 200 Mb/s measured in free-space outside the chip. PICs such as this may in the future offer a cost-effective pathway for the efficient use of the optical spectrum via electrical filtering and encoding techniques, much in the same way cell phone systems operate today.

21.6 APPLICATIONS OF INTEGRATED OPTICS

Digital Transmission

The performance metric in digital optical fiber transmission is the ability of a transmitter to deliver a signal to the receiver at the end of the link in a manner such that the receiver can clearly distinguish between the “0” and “1” state in each time period or bit slot. Binary amplitude-shift-keyed transmission (ASK) is by far the most common format in commercial systems, but high performance systems today also employ differential phase-shift-keyed (DPSK) formats. A decision circuit at the receiver must distinguish between “0” and “1,” and this circuit will be more susceptible to noise when the “0” and “1” level difference is reduced, or when the time over which this difference is maintained is reduced below the full bit period.

The performance of a transmitter is thus governed by its rise and fall times, and its modulation bandwidth or flatness of response to avoid pattern-effects and its output power. Furthermore, the spectral characteristics of its optical output can impair transmission. Examples of the latter include unnecessary optical bandwidth, as might be present in an LED or a multilongitudinal-mode laser, that can produce pulse spreading of the digital pulses due to dispersion in the transmission fiber. While transmission sources include current-modulated LEDs, for speeds higher than approximately 100 Mb/s semiconductor lasers are used, and IO technology has played a fundamental role in the evolution of semiconductor laser technology. In very high-speed systems (typically > 1 Gb/s), dispersive pulse distortion can cause severe degradation with directly modulated lasers unless devices which emit only one longitudinal mode are employed. The incorporation of gratings in DFB and DBR lasers has produced sources with exceptional spectral purity and allow multi-gigabit per second transmission over intermediate distances (< 100 km) in conventional fiber.

The advent of fiber amplifiers has enabled longer transmission spans without digital regeneration, and here the unavoidable frequency excursions that result from directly modulating even a single-longitudinal-mode laser again lead to pulse spreading and distortion. In these instances, a continuous wave (CW) laser followed by an external modulator is a preferred source. The integrated DFB/electroabsorption modulator, as discussed in Sec. 23.5, provides such a source. These PICs have demonstrated error-free transmission in excess of 500 km in dispersive fiber at 2.5 Gb/s.¹⁰⁴ However, even these devices impose a small residual dynamic phase shift on the signal due to electrorefractive effects accompanying the induced absorption in the modulator. This can be especially problematic with typical percent-level antireflection coatings on the output facet, since this will provide phase-varying optical feedback into the laser and further disrupt its frequency stability.

The highest performance digital transmission has been achieved using external LiNbO₃ Mach-Zehnder modulators to encode a CW semiconductor laser. Modulators similar to that in Fig. 14 have been constructed to allow separate traveling-wave voltages to be applied to each arm of the modulator in a push-pull configuration. This device design can produce a digitally encoded signal with *zero* residual phase shift or chirp.¹⁰⁵ Such a source has only its information-limited bandwidth and generally provides nearly optimized performance in a dispersive environment. The Mach-Zehnder can also be driven to intentionally provide positive or negative chirping to optimize transmission characteristics in dispersive links.¹⁰⁶ Semiconductor lasers have also been monolithically integrated with Mach-Zehnder modulators using electrorefraction from the QCSE.¹⁰⁷

Recently there has been significant focus on using phase modulation formats such as DPSK, and higher level formats such as differential quadrature phase-shift keyed (DQPSK) modulation where each bit period can contain 2 bits of information. Such formats not only allow higher spectral efficiency, or information per unit bandwidth in the fiber, but also offer advantages in robustness against impairments from fiber nonlinearity which can seriously degrade the performance of long-distance amplified WDM systems. Just as in the case of the simpler integrated laser-modulator, PIC technology provides a powerful vehicle for the precision control of the pathlengths or phase delays required to encode and decode such formats in transmitters and receivers.^{108,109}

Analog Transmission

A second application area that is expected to use a large number of IOCs is analog fiber optic links. Analog fiber optic links are currently being used to transmit cable television (CATV) signals at the trunk and supertrunk level. They are also being used for both commercial and military antenna remoting. Analog fiber optic links are being fielded for these applications because of their low distortion, low loss, and low life-cycle cost when compared to more conventional coaxial cable-based transmission systems.

An analog fiber optic link using IOCs is typically configured as shown in Fig. 20. A high-power CW solid-state laser, such as a 150-mW diode-pumped YAG laser operating at 1319 nm, is typically used as the source in order to maximize dynamic range, carrier-to-noise ratio, and link gain. An interferometric modulator, such as a Mach-Zehnder interferometer or a Y-fed balanced bridge modulator, is typically used to modulate the light with the applied RF or microwave signal

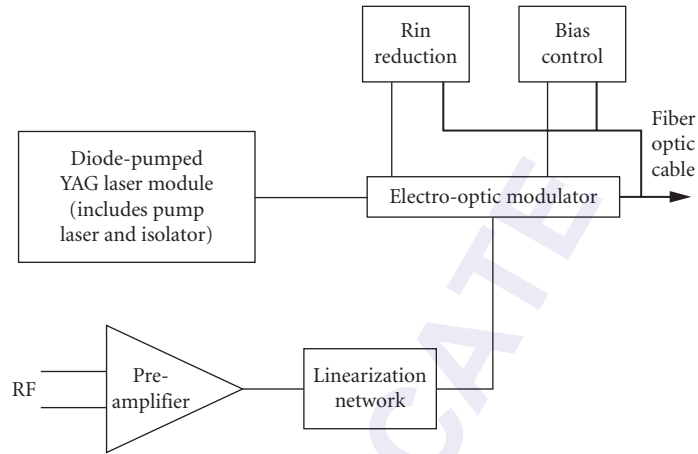


FIGURE 20 Standard configuration for externally modulated analog fiber optic link.

via the linear electro-optic effect. Current analog links for CATV signal distribution utilize a 1-GHz Y-fed balanced bridge modulator biased at the quadrature point (linear 3-dB point).^{110,111} A predistortion circuit is required to minimize third-order distortion associated with the interferometric modulator response. The CATV-modulated signal can be transmitted on both output fibers of the device. Analog links for antenna remoting typically fit into one of two categories. Certain applications require relatively narrow passbands in the UHF region while other microwave-carrier applications require very broadband (several GHz) performance. Y-fed balanced bridge modulators biased at the quadrature point are again used to perform the electrical-to-optical conversion, with the electrode structure tailored to the application. In both narrowband and broadband applications, 20- to 30-dB preamplifiers are typically utilized to minimize the noise figure and maximize the RF gain of the link.

Two important modulator parameters are the insertion loss and the half-wave drive voltage which both impact the link gain and dynamic range. Fully packaged Y-fed balanced bridge modulators with 2.5 to 4.0 dB insertion loss are now readily achieved in production for both UHF and microwave bandwidths. A trade-off typically needs to be made between half-wave voltage and bandwidth for both lumped-element and traveling-wave electrode structures. Commercially available lumped-element LiNbO_3 interferometric modulators typically have half-wave voltages of approximately 5 V for 600-MHz, 1-dB bandwidths. Commercially available traveling-wave LiNbO_3 interferometric modulators typically have half-wave voltages of approximately 8 V for 12-GHz, 3-dB bandwidths. The half-wave voltages of LiNbO_3 traveling-wave modulators can be reduced by more than a factor of two using a velocity-matched electrode structure as described in Ref. 93.

In order to be used in practical systems, it is critical that the integrated optical modulators have well-behaved, flat frequency responses in the band of interest. Modulators for CATV signal transmission and UHF antenna remoting typically required that the amplitude response and the phase response be flat to ± 0.25 dB and $\pm 2^\circ$, respectively. The frequency response of an integrated optical modulator is a function of both the device design and packaging parasitics. Care must be exercised in designing modulators since LiNbO_3 is both a piezoelectric and an acousto-optic material. Early LiNbO_3 modulators typically had 1 to 3 dB of ripple in the passband due to acoustic mode excitation. When packaging lumped-electrode devices, it is also critical to keep terminating resistors close and wire bonds short to minimize stray inductance and capacitance. When packaging traveling-wave modulators, it is critical to control the impedance of the launch, the transitions, the device, and the termination. Through proper device design and packaging, it is possible to achieve well-behaved frequency responses in both lumped-electrode and traveling-wave devices as shown in Fig. 21.

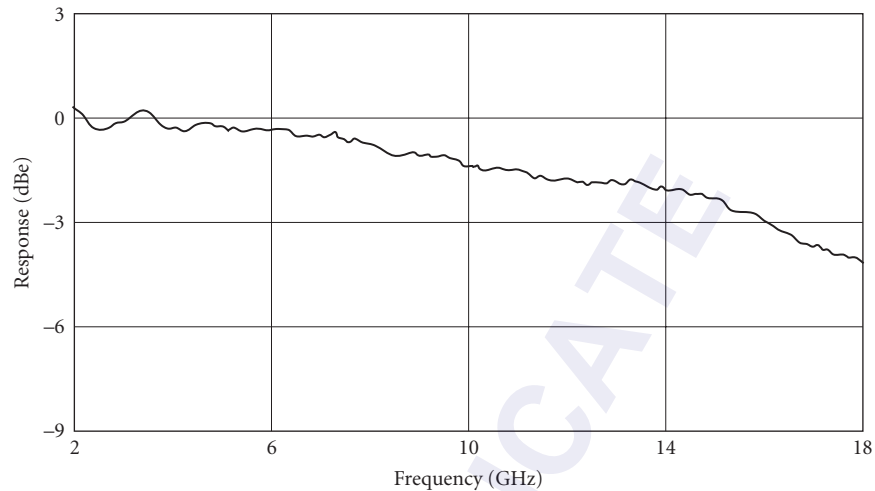


FIGURE 21 Frequency response of APE LiNbO₃ microwave interferometric modulator for broadband antenna remoting.

An additional issue that impacts IOC modulator design for analog links is harmonic and intermodulation distortion. Most modulators used in analog links are interferometric in nature with a sinusoidal transfer function. By operating the device precisely at the quadrature point, all even harmonics can be suppressed. Second-harmonic distortion less than -75 dBc is easily achieved using an electronic feedback loop for modulator bias. Alternative techniques are being investigated to laser trim one leg of a Mach-Zehnder to bring the modulator to quadrature. Third-order distortion due to the sinusoidal transfer function of the interferometric modulator also poses a problem, but the transfer functions are very well-behaved and predictable, and this distortion can be suppressed to acceptable levels using electronic predistortion or feed-forward techniques.

Forty- and eighty-channel CATV transmitters operating at 1300-nm wavelengths with APE LiNbO₃ modulators are currently being fielded. Compared to coaxial transmission which requires transmitters every 500 m, the fiber optic systems can transmit over distances up to 50 km without repeaters. Similarly, externally modulated analog fiber optic links are currently being fielded for military and commercial applications. UHF links with 115 dB/Hz^{2/3} dynamic range, 4-dB noise figure, and unity gain have been demonstrated using commercially available hardware. These systems will maintain this quality of transmission over temperature ranges of -25 to $+50^{\circ}\text{C}$. Microwave links with 2 to 18 GHz frequency response, 114-dB/Hz^{2/3} spurious-free dynamic range, and input noise figure of 22 dB can also be achieved using commercially available hardware.

Switching

Arrays of IO switches have been proposed and demonstrated for a variety of space switching and time-multiplexed switching (TMS) applications. In space switching, it is generally necessary to formulate the switches in a nonblocking geometry and the reconfiguration time can be relatively slow (seconds or longer). This requirement led to the development of cross-bar switches in which an $N \times N$ switch contains N^2 IO switches and $2N - 1$ stages and from 1 to $2N - 1$ cross points. Typically $N = 4$ in LiNbO₃ and in InP. More recently, much attention in IO switch arrays has shifted to the dilated Benes architecture which is only rearrangeably nonblocking but reconfigurable in short (ns) times suitable for TMS, and has the advantage of requiring substantially fewer switches and a constant number $2 \log_2 N$ of cross points.

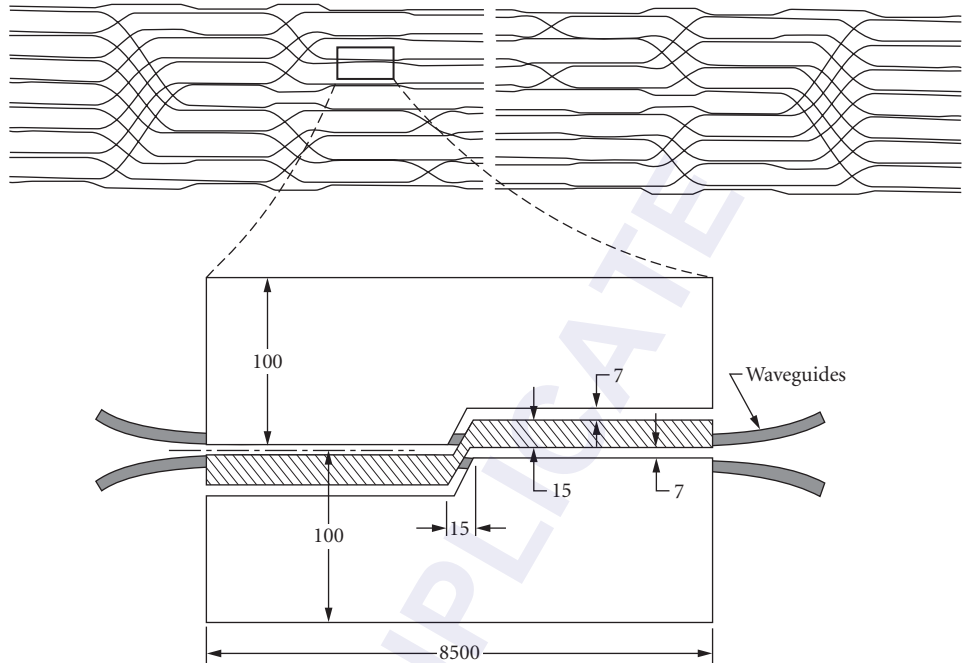


FIGURE 22 Architecture of 8×8 diluted Benes directional coupler switch array.

A schematic of a two-chip 8×8 diluted Benes switch making extensive use of crossovers is shown in Fig. 22.¹¹² The performance of switch arrays of the diluted Benes architecture are much more forgiving than crossbar switches to the degradation of individual switches. The device shown contains 48 delta beta switches driven by a single-voltage arrangement. The switching voltage at $1.3 \mu\text{m}$ was $9.4 \pm 0.2 \text{ V}$, the insertion loss varied from -8 to -11 dB (93 percent of the 256 paths through the switch were within $\pm 1 \text{ dB}$), the cross-talk levels in individual switches ranged from -22 to -45 dB , and the reconfiguration time was 2.5 ns . Larger 16×16 switches have also been demonstrated.¹¹³

An advantage of these types of IO switch arrays is that they are data rate transparent. That is, once the switch is reconfigured, the data stream through the device is simply the passage of light and can easily be multi-gigabit. Crossbar switches are now commercially available and other types of arrays continue to be explored.

Fiber Optic Gyroscopes

Another application that may require large quantities of integrated optical circuits is the fiber optic gyroscope (FOG)^{114–119}. A FOG is one form of a Sagnac interferometer, in which a rotation rate results in a phase shift between clockwise- and counterclockwise-propagating optical fields. The most frequently utilized FOG configuration, which was first proposed by workers at Thomson CSF in the mid 1980s,¹²⁰ is presented in Fig. 23.

FOG IOCs are typically fabricated in LiNbO_3 using the annealed proton exchange (APE) process⁶⁶ although titanium-diffused IOCs with surface plasmon polarizers have also been utilized. The IOC performs four primary functions in the fiber gyroscope. First, the Y-junction serves as the loop coupler splitting and recombining the clockwise- and counterclockwise-propagating optical fields. Second, the IOC functions as a high-quality polarizer. Third, a 90° phase dither (at the eigen frequency of the fiber coil) is typically applied to one of the integrated optical phase modulators.

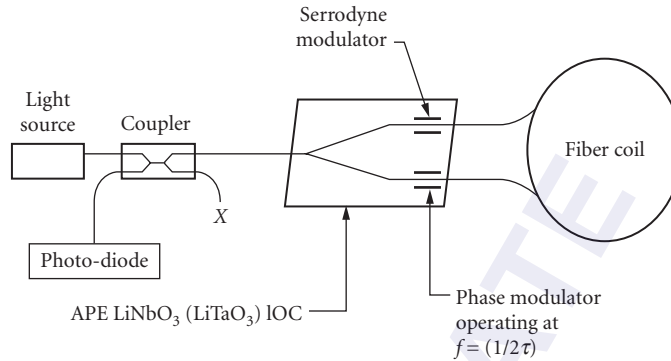


FIGURE 23 Standard configuration for fiber optic gyroscope incorporating a three-port integrated optical circuit.

This approach keeps the Sagnac interferometer biased at the 3-dB point where it is linear and most sensitive to rotation. Finally, in a closed-loop FOG configuration, one of the phase modulators functions as a frequency shifter. A serrodyne signal (saw-tooth wave) is applied to the phase modulator to effectively cancel the shift due to the rotation.

The output signal from a fiber gyro at rest is the sum of white receiver noise, primarily dependent on the amount of optical power arriving at the detector, and an additional long-term drift of the mean value. The long-term drift in a FOG associated with a residual lack of reciprocity typically limits the sensitivity of the FOG to measure low rotation rates. Another important characteristic of a gyro is the scale factor, which is a measure of the linearity between the actual rotation rate and the gyro response. The critical performance parameters for a FOG IOC are presented in Table 1. The performance of 800- and 1300-nm APE LiNbO₃ FOG IOCs that are currently in production is also presented in this table.

One application of the FOG is inertial guidance, requiring a FOG with a bias drift < 0.01°/h and a scale factor accuracy < 5 ppm. A 1300-nm LED or an erbium-doped fiber is typically used as the light source. A large coil of polarization-maintaining fiber (typically 1 km of fiber wound in a 15 to 20 cm diameter coil) and precise source spectral stability are required to achieve the desired sensitivity. The fiber is typically wound in a quadrupole configuration to minimize sensitivity to temperature gradients. With recent improvements in optical sources, integrated optics, and fiber coil-winding technology, it is now possible to achieve inertial grade FOG performance over a temperature range of -55 to +95°C.

TABLE 1 Critical Performance Parameters for 800- and 1300-nm APE LiNbO₃ FOG IOCs. Listed Values Are Maintained over a Temperature Range of -55 to +95°C and during Vibration up to 15 Grms.

Performance Parameter	1300-nm IOCs	800-nm IOCs
Insertion loss (pigtailed)	3 dB	4 dB
Polarization extinction	70 dB	60 dB
Y-junction split ratio (pigtailed)	48/52 to 52/48	45/55 to 55/45
Polarization crosstalk at fiber-waveguide interfaces	< -30 dB	< -25 dB
Optical back reflection	< -65 dB	< -65 dB
Half-wave voltage	4.5 V	3.5 V
Residual intensity modulation	0.02%	0.05%

A second tactical-grade FOG design is more typical to aerospace applications, with bias drift and scale factor accuracy requirements ranging from 0.1 to $10^\circ/\text{h}$ and 10 to 1000 ppm, respectively. These systems are typically designed for operation at 810 to 830 nm to make use of low-cost multimode 830-nm AlGaAs laser diodes as used in consumer electronic products. These systems typically utilize 2- to 5-cm-diameter fiber coils with 100 to 200 m of either polarization-maintaining or single-mode fiber. A third very low-cost, low-grade FOG design for automotive navigation is also nearing production. The required bias drift is only $1000^\circ/\text{h}$, and a closed-loop configuration is unnecessary since the scale factor accuracy is only 0.1 percent. Current designs to achieve low cost include low-performance IOCs, laser, and fiber couplers, fully automated FOG assembly and test procedures, and only approximately 50 m of single-mode fiber. More advanced IOCs, including four-port designs that integrate the source/detector coupler into the IOC, are also being considered to reduce component count.

WDM Systems

Wavelength division multiplexing (WDM), by encoding parallel data streams at different wavelengths on the same fiber, offers a technique to increase transmission capacity, or increase networking or switching flexibility, without requiring higher speed electronics to process each channel. As noted earlier, IO is an enabling technology in WDM demultiplexer design. However, due to the large number of component interconnections at the transmit and receive end of a WDM link, PIC technology offers great promise in cost reduction for both WDM transmitters and receivers by additionally eliminating separate packages for lasers, amplifiers, detectors, and the like.

One key application of PICs for WDM systems is the stable and efficient longitudinal connection of optical elements required to form a tunable laser. The tunable Bragg laser was briefly discussed earlier in the description of the balanced heterodyne receiver PIC. Figure 24 below shows a more complex sampled-grating distributed Bragg reflector PIC. Here the reflectivity of a grating with

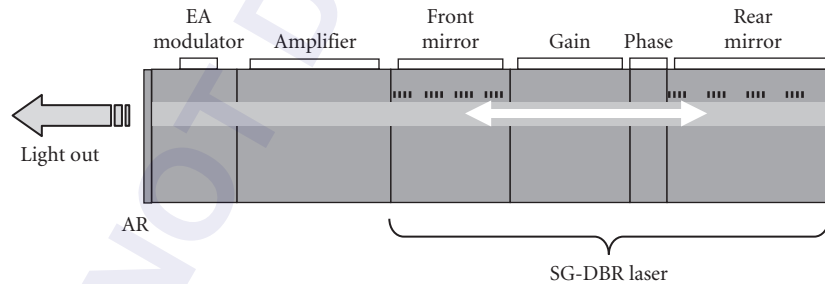


FIGURE 24 Tunable sampled-grating distributed Bragg reflector (SG-DBR) laser with integrated semiconductor optical amplifier and electroabsorption modulator.

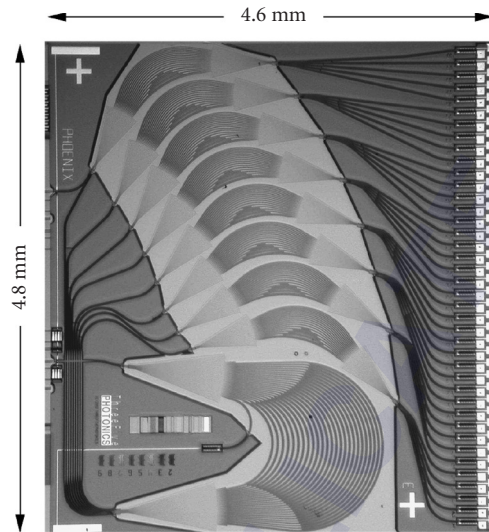


FIGURE 25 Integrated 40-channel optical channel monitor, comprising 9 different AWG demultiplexers and 40 photodetectors. (By III-V Photonics, Inc.)

sections periodically omitted provides a comb of reflection bands, and if the front and back combs have different spacings in frequency, the laser can tune over a wide range with a small index change because the laser will operate only at a frequency where the front and back reflection bands coincide in frequency.¹²¹ Tuning is accomplished by changing the index of the front or back section by current injection, for example.

Because the InP materials system readily allows detectors as well as high-quality passives, the AWG demultiplexer discussed earlier can be combined with an integrated detector for each wavelength output port to make a single-chip WDM receiver. Figure 25 illustrates a 40-channel WDM receiver PIC, comprising nine different AWG demultiplexers and 40 photodetectors within a 4.6×4.8 mm chip.¹²² This PIC had only 4-dB on-chip loss and provided less than 35-dB cross talk between different wavelength channels.

Figure 26 shows the layout of a WDM transmission PIC, which is perhaps the most sophisticated example of commercially deployed active IO demonstrated to date.¹²³ Following the concept initially introduced by Aiki et al.¹²⁴ in combining multiple DFB lasers, this PIC includes 10 frequency tunable DFB lasers with wavelength locking function, 10 electroabsorption modulators, 10 monitor detectors, 10 adjustable attenuators, an AWG for low-loss, frequency-selective combining into a single waveguide output port. This PIC provides a 100-Gb/s transmission capacity from a single chip, and together with a matching 10-channel WDM receiver PIC has been commercially deployed in telecommunication networks worldwide.¹²³ These PICs have demonstrated remarkable performance, and offer significant improvements in cost, size, power, and reliability compared to discrete component solutions.

Silicon Photonics Transceivers

As noted in the earlier discussions of IO materials, silicon has gained significant attention in recent years due to its ability to harness highly advanced CMOS fabrication technologies, combined with the very high index contrast available and effective modulation using the plasma index change from free carriers. This modulation is most effectively done at high speed using either a reverse-biased P-N junction, or using a structure similar to a field-effect-transistor (FET) where charge can accumulate on a gate oxide within the optical mode. Both of these designs have been demonstrated

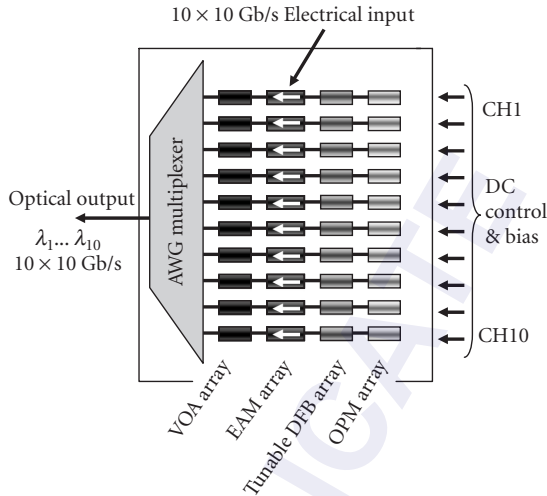


FIGURE 26 Chip architecture for commercially implemented PIC with nearly 50 elements, including 10 frequency tunable DFB lasers with wavelength-locking function, 10 electroabsorption modulators, 10 monitor detectors, 10 adjustable attenuators, and an integrated AWG combiner.

by several organizations^{125,126} and commercialized with fabrication conducted within a CMOS electronics foundry using nearly standard IC fabrication steps. PICs have been demonstrated that include 10 Gb/s transmitters and receivers, and the ability to include CMOS electronics has also allowed for on-chip modulator drivers, preamplifiers, clock and data-recovery circuits, and even pseudo-random bit-stream generators for on-chip testing.⁴⁶ Commercial silicon photonics transceivers have demonstrated low-power operation, and also 40-Gb/s data rates using integration to achieve parallelism with only one laser source.

21.7 FUTURE TRENDS

Shift from R&D to Manufacturing

Integrated optical circuit technology has now advanced from R&D into manufacturing as noted in the examples given. LiNbO₃ modulators and integrated DFB/electroabsorption modulators have been critical, high-volume products for more than a decade. There are now several companies producing more complex IO and PIC products in LiNbO₃, Si/SiO₂, and InP in moderate volumes (several thousand devices each per year) with considerable success in reducing manufacturing costs and achieving high reliability.

The majority of the LiNbO₃ devices described in this chapter can be fabricated using either existing or slightly modified semiconductor processing equipment. Clean room requirements are not nearly as tight as what is required for VLSI and MMIC wafer fabrication. Production yields of LiNbO₃ integrated optical circuits have been obtained well in excess of 95 percent. The majority of the defects were initially mechanical in nature (probes scratching electrodes or fibers damaging polished end faces), but these have been minimized as the processes have become more automated. At this time, it appears that all wafer-processing operations should be easily scalable to wafer batch processing, including the end-face optical polishing. The majority of the cost of an integrated optical circuit is currently associated with fiber attachment, packaging, and final testing. While these operations are

not fundamentally expensive, the limited production volumes do not always justify the capital expenditure required to fully automate processes with robotics and machine vision.

The second area that has seen improvement is device reliability. First-generation commercial integrated optic products were plagued by premature failures and poor performance over temperature. The majority of these problems resulted from poor fiber attachment and packaging techniques. These problems have been remedied by carefully selecting compatible material systems and incorporating proven hybrid electronic packaging techniques. For example, commercially available LiNbO_3 integrated optical circuits can be thermally cycled hundreds of times from -65 to $+125$ °C with less than 1 dB variation in insertion loss. The devices can also withstand 30 Grms random vibration testing. Long-term aging effects have been addressed by identifying potential failure mechanisms and then by developing physical models of these failures. LiNbO_3 integrated optical circuit reliability has now achieved a level where devices can be certified for operational lifetimes in excess of 25 years for properly designed, fabricated, and assembled devices. Obviously, this level of performance can only be guaranteed by fabricating the devices in well-controlled, well-documented production environments.

In the semiconductor area, reliability certification technology for PICs borrows heavily from the more mature level models and aging methodologies developed for semiconductor lasers. This includes sequences of purges at high currents or voltages combined with high temperatures, together with extended accelerated aging to identify activation energies and aging rates for various degradation mechanisms. Integrated laser modulators have been deployed in high volume with outstanding reliability, also easily in excess of 25 years including wavelength stability for WDM lasers. More complex PICs have also shown extremely high reliability, and as with the LiNbO_3 IO devices, any failures that do occur are most often associated with poorly designed packaging of the device rather than the PIC chip itself. This again points to one of the fundamental advantages of PICs, where many individual component packages are eliminated by lithographic connections and only one package is required for the final PIC or IO device.

Advanced Integration and New PIC Applications

The most active areas of IO and PIC research today are in telecommunications, datacom, and sensor applications. In InP materials, for example, there has been extensive work on all-optical networking technologies where wavelength-converter PICs have been demonstrated that encode incoming data onto a new wavelength determined by an on-chip tunable laser in conjunction with nonlinear Mach-Zehnder interferometers with SOA gain elements.^{127,128} The complexity, speed, and wavelength count of high speed WDM transmit and receive PICs continues to grow, with research demonstrations of 1.6 Tb/s from a single chip combining 40 tunable DFB lasers with 40 electroabsorption modulators running at 40 Gb/s each.¹²⁹ As transmission technology continues to evolve in sophistication and spectral efficiency, work on advanced modulation formats moves forward aggressively. A resurgence of interest in heterodyne-detection technology has shown that many system impairments can be effectively removed in the electronic domain when the information is received linearly in the optical electric field. It is quite likely that information coding in the optical domain will progress to a level that is now commonplace in the world of RF and cell phone technology.

Silicon photonics has made remarkable progress in a short period, and whereas true OEIC integration with electronics has only had very limited penetration in III-V materials, the potential for including sophisticated electronics in silicon has provided a strong incentive to explore and commercialize this technology. In datacom, for example, the proximity of driver circuits and modulator devices allows for low power by obviating the usual electrical transmission line and terminating resistor approaches. For sensors, the ability to combine preamplification, analog-to-digital converters, and digital signal processors is extremely compelling, and biochem sensors that can detect frequency shifts in suitably sensitized resonators are now emerging. It is also likely that the mapping of photonic functions into commercial CMOS electronics foundries will become more common, and if suitable design tools become more widely available, the reduced barriers to using silicon photonics solutions may lead to a growth of new market applications.

Analogies with massively repetitive logic gates in electronics integration are not particularly good, but there is indeed now a proliferation of applications where stable designs have emerged that require significant numbers of interconnected optical components. In these areas, the improvements in size, cost, power, and reliability are extremely compelling, and today it is safe to say that IO and PIC technologies have earned a secure footing in the world of optics.

REFERENCES

1. S. E. Miller, "Integrated Optics: An Introduction," *Bell Labs Tech. J.* **48**(7):2059–2070 (1969).
2. T. Tamir, ed., *Guided-Wave Optoelectronics*, 2d ed., New York: Springer-Verlag (1990).
3. R. G. Hunsperger, *Integrated Optics: Theory and Technology*, Berlin: Springer-Verlag (1985).
4. T. Tamir, ed., *Integrated Optics*, 2d ed., New York: Springer-Verlag (1979).
5. L. D. Hutchenson, ed., *Integrated Optical Circuits and Components*, New York: Marcel Dekker (1987).
6. K. Okamoto, *Fundamentals of Optical Waveguides*, San Diego, Calif.: Academic Press (2000).
7. H. Nishihara, M. Haruna, and T. Suhara, *Optical Integrated Circuits*, New York: McGraw-Hill (1985).
8. D. Marcuse, *Theory of Dielectric Optical Waveguide*, 2d ed., San Diego, Calif.: Academic Press Inc., pp. 307–318 (1991).
9. A. Yariv and P. Yeh, "Optical Waves in Crystals," New York, John Wiley & Sons Inc., pp. 177–201, 425–459 (1984).
10. I. P. Kaminow and S. E. Miller, ed., *Optical Fiber Telecommunications II*, San Diego, Calif.: Academic Press (1985).
11. I. P. Kaminow and T. L. Koch, ed., *Optical Fiber Telecommunications III*, San Diego, Calif.: Academic Press (1997).
12. I. P. Kaminow and T. Li, ed., *Optical Fiber Telecommunications IV*, San Diego, Calif.: Academic Press (2002).
13. Special Issue on Photonic Devices and Circuits, *IEEE J. Quantum Electron.* **QE-27**(3) (1991).
14. Special Issue on Integrated Optics, *J. Lightwave Technol.* **6**(6) (1988).
15. Special Issue on Integrated Optics, *IEEE J. Quantum Electron.* **QE-22**(6) (1986).
16. *Digests of IEEE/OSA Conferences on Integrated Photonics Research*, 1990–Present, OSA, Washington, D.C.
17. *Digests of IEEE Conference on Integrated and Guided Wave Optics*, 1972–1989, OSA, Washington, D.C.
18. G. B. Hocker and W. K. Burns, "Mode Dispersion in Diffused Channel Waveguides by the Effective Index Method," *Appl. Opt.* **16**:113–118 (1977).
19. J. F. Lotspeich, "Explicit General Eigenvalue Solutions for Dielectric Slab Waveguides," *Appl. Opt.* **14**:327–335 (1975).
20. M. D. Feit and J. A. Fleck, Jr., "Light Propagation in Graded-Index Optical Fibers," *Appl. Opt.* **17**:3990–3998 (1978).
21. M. J. Weber, ed., *CRC Handbook of Laser Science and Technology, Vol. IV, Optical Materials, Part 2*, Boca Raton: CRC Press Inc. (1986).
22. M. A. Fromowitz, "Refractive Index of $\text{Ga}_{1-x}\text{Al}_x\text{As}$," *Solid State Commun.* **15**:59–63 (1974).
23. C. H. Henry, L. F. Johnson, R. A. Logan, and D. P. Clarke, "Determination of the Refractive Index of InGaAsP Epitaxial Layers by Mode Line Luminescence Spectroscopy," *IEEE J. Quantum Electron.* **QE-21**:1887–1892 (1985).
24. T. Maldonado, "Electro-Optics Modulators," in vol. V, chap. 7 of this *Handbook*.
25. I. P. Kaminow, *An Introduction to Electro-Optic Devices*, Orlando: Academic Press (1974).
26. T. S. Moss, G. S. Burrell, and B. Ellis, *Semiconductor Opto-Electronics*, New York: John Wiley & Sons Inc. (1973).
27. K. Tharmalingam, "Optical Absorption in the Presence of a Uniform Field," *Phys. Rev.* **130**:2204–2206 (1963).

28. D. A. B. Miller, J. S. Weiner, and D. S. Chemla, "Electric-Field Dependence of Linear Optical Properties Quantum Well Structures: Waveguide Electroabsorption and Sum Rules," *IEEE J. Quantum Electron.* **QE-22**:1816–1830 (1986).
29. J. E. Zucker, I. Bar-Joseph, G. Sucha, U. Koren, B. I. Miller, and D. S. Chemla, "Electrorefraction in GaInAs/InP Multiple Quantum Well Heterostructures," *Electron. Lett.* **24**:458–460 (1988).
30. Y. -H. Kuo, Y. -K. Lee, Y. Ge, S. Ren, J. E. Roth, T. I. Kamins, D. A. B. Miller, and J. S. Harris, "Strong Quantum-Confined Stark Effect in Germanium Quantum-Well Structures on Silicon," *Nature* **437**:1334–1336 (2005).
31. C. L. Tang, "Nonlinear Optics," in vol. IV, chap. 10 of this *Handbook*.
32. M. Yamada, N. Nada, M. Saitoh, and K. Watanabe, "First-Order Quasi-Phase Matched LiNbO₃ Waveguide Periodically Poled by Applying an External Field for Efficient Blue Second-Harmonic Generation," *Appl. Phys. Lett.* **62**:435–436 (1993).
33. H. A. Haus, E. P. Ippen, and F. J. Leonberger, "Nonlinear Optical Waveguide Devices," in *Optical Signal Processing*, J. L. Horner (ed.), Orlando, FL: Academic Press, pp. 245–267 (1987).
34. M. A. Foster, A. C. Turner, J. E. Sharping, B. S. Schmidt, M. Lipson, and A. L. Gaeta, "Broad-Band Optical Parametric Gain on a Silicon Photonic Chip," *Nature* **441**:960–963 (2006).
35. T. K. Findakly, "Glass Waveguides by Ion-Exchange: A Review," *Opt. Eng.* **24**:244 (1985).
36. R. V. Ramaswamy and R. Srivastava, "Ion-Exchanged Glass Waveguides: A Review," *J. Lightwave Technol.* **LT-6**:984–1002 (1988).
37. T. Izawa and H. Nakogome, "Optical Waveguide formed by Electrically Induced Migration of Ions in Glass Plates," *Appl. Phys. Lett.* **21**:584 (1972).
38. M. Kawachi, "Silica Waveguides on Silicon and their Applications to Integrated-Optic Components," *Opt. Quant. Electron.* **22**:391–416 (1990).
39. H. Takahashi, Y. Ohnori, and M. Kawachi, "Design and Fabrication of Silica-Based Integrated-Optic 1 × 128 Power Splitter," *Electron. Lett.* **27**:2131–2133 (1991).
40. C. H. Henry, G. E. Blonder, and R. F. Kazarinov, "Glass Waveguides on Silicon for Hybrid Optical Packaging," *J. Lightwave Technol.* **7**:1530–1539 (1989).
41. R. Adar, Y. Shani, C. H. Henry, R. C. Kistler, G. E. Blonder, and N. A. Olsson, "Measurement of Very Low-Loss Silica on Silicon Waveguides with a Ring Resonator," *Appl. Phys. Lett.* **58**(5):444–445 (1991).
42. Y. Shani, C. H. Henry, R. C. Kistler, K. J. Orlovski, and D. A. Ackerman, "Efficient Coupling of a Semiconductor Laser to an Optical Fiber by Means of a Tapered Waveguide on Silicon," *Appl. Phys. Lett.* **55**:2389–2391 (1989).
43. R. Soref, "The Past, Present, and Future of Silicon Photonics," *IEEE J. Quantum Electron.* **QE-12**:1678–1687 (2006).
44. G. T. Reed and A. P. Knights, *Silicon Photonics: An Introduction*, West Sussex, England: John Wiley & Sons Ltd. (2004).
45. L. Pavesi and D. L. Lockwood, eds., *Silicon Photonics*, Berlin Heidelberg: Springer-Verlag (2004).
46. C. Gunn and T. L. Koch, "Silicon Photonics," in *Optical Fiber Telecommunications V*, I. P. Kaminow, T. Li, and A. E. Willner, (eds.), San Diego, Calif.: Academic Press (2008).
47. M. Bruel, "Silicon on Insulator Material Technology," *Electron. Lett.* **31**:1201–1202 (1995).
48. U. Fischer, T. Zinke, J.-R. Kropp, F. Amtdt, and K. Petermann, "0.1 dB/cm Waveguide Losses in Single-Mode SOI Rib Waveguides," *IEEE Phot. Tech. Lett.* **8**(5):647–648 (1996).
49. M. A. Webster, R. M. Pafchek, G. Sukumaran, and T. L. Koch, "Low-Loss Thin SOI Waveguides and High-Q Ring Resonators," Paper FTuV3 in Proceedings of the Optical Society of America Annual Meeting, Tucson, (Oct. 18, 2005).
50. R. C. Alferness, "Guided-Wave Devices for Optical Communications," *IEEE J. Quantum Electron.* **QE-17**(6):946–959 (1981).
51. H. F. Taylor, "Applications of Guided-Wave Optics in Signal Processing and Sensing," *Proc. IEEE* **75**(11):1524–1535 (1987).
52. E. Voges and A. Neyer, "Integrated-Optic Devices on LiNbO₃ for Optical Communication," *J. Lightwave Technol.* **LT-5**(9):1229–1238 (1987).
53. L. Thylen, "Integrated Optics in LiNbO₃: Recent Developments in Devices for Telecommunications," *J. Lightwave Technol.* **LT-6**(6):847–861 (1988).

54. R. C. Alferness, "Waveguide Electrooptic Switch Arrays," *IEEE J. Selected Areas in Communications* **6**(7):1117–1130 (1988).
55. R. V. Schmidt and I. P. Kaminow, "Metal Diffused Optical Waveguides in LiNbO₃," *Appl. Phys. Lett.* **15**:458–460 (1974).
56. M. Fukuma and J. Noda, "Optical Properties of Ti-Diffused LiNbO₃ Strip Waveguides and their Coupling-to-a-Fiber Characteristics," *Appl. Opt.* **19**:591–597 (1980).
57. S. K. Korotky and R. C. Alferness, "Ti:LiNbO₃ Integrated Optic Technology," *Integrated Optical Circuits and Components*, L. D. Hutcheson, (ed.), New York: Marcel Dekker (1987).
58. R. J. Esdaile, "Closed-Tube Control of Out-Diffusion during Fabrication of Optical Waveguides in LiNbO₃," *Appl. Phys. Lett.* **33**:733–734 (1978).
59. J. L. Jackel, V. Ramaswamy, and S. P. Lyman, "Elimination of Out-Diffused Surface Guiding in Titanium-Diffused LiNbO₃," *Appl. Phys. Lett.* **38**:509–511 (1981).
60. A. M. Glass, "The Photorefractive Effect," *Opt. Eng.* **17**:470–479 (1978).
61. J. L. Jackel, C. E. Rice, and J. J. Veselka, "Proton Exchange for High-Index Waveguides in LiNbO₃," *Appl. Phys. Lett.* **47**:607–608 (1982).
62. P. G. Suchoski, T. K. Findakly, and F. J. Leonberger, "Stable Low-Loss Proton-Exchanged LiNbO₃ Waveguide Devices with No Electrooptic Degradation," *Opt. Lett.* **13**:1050–1052 (1988).
63. P. G. Suchoski, T. K. Findakly, and F. J. Leonberger, "Low-Loss High Extinction Polarizers Fabricated in LiNbO₃ by Proton Exchange," *Opt. Lett.* **13**:172–174 (1988).
64. T. K. Findakly, P. G. Suchoski, and F. J. Leonberger, "High-Quality LiTaO₃ Integrated Optical Waveguides and Devices Fabricated by the Annealed-Proton-Exchange Technique," *Opt. Lett.* **13**:797–799 (1988).
65. See, for example, *Conference Proceedings of 20th International Conference on Indium Phosphide and Related Materials*, IEEE, Versailles (2008).
66. R. Clawson, "Reference Guide to Chemical Etching of InGaAsP and In_{0.53}Ga_{0.47}As Semiconductors," NOSC Tech. Note 1206, San Diego, Calif. (1982).
67. S. Adachi and H. Kawaguchi, "Chemical Etching Characteristics of (001) InP," *J. Electrochem. Soc.: Sol. St. Sci. and Tech.* **128**:1342–1349 (1981).
68. L. A. Coldren, K. Furuya, B. I. Miller, and J. A. Rentschler, "Etched Mirror and Groove-Coupled GaInAsP/InP Laser Devices for Integrated Optics," *IEEE J. Quantum Electron.* **QE-18**:1679–1688 (1982).
69. L. A. Coldren and J. A. Rentschler, "Directional Reactive-Ion-Etching of InP with Cl₂ Containing Gases," *J. Vac. Sci. Technol.* **19**:225–230 (1981).
70. J. W. McNabb, H. G. Craighead, and H. Temkin, "Anisotropic Reactive Ion Etching of InP in Methane/Hydrogen Based Plasma," Paper TuD15 in *Tech. Digest of Integrated Photonics Research 1991*, Monterey, pp. 26–27.
71. P. Buchmann and H. Kaufmann, "GaAs Single-Mode Rib Waveguides with Reactive Ion-Etched Totally Reflecting Corner Mirrors," *IEEE J. Lightwave Technol.* **LT-3**:785–788 (1985).
72. S. Oku, M. Okayasu, and M. Ikeda, "Low-Threshold CW Operation of Square-Shaped Semiconductor Ring Lasers (Orbiter Lasers)," *IEEE Phot. Tech. Lett.* **3**:588–590 (1991).
73. W. J. Grande, J. E. Johnson, and C. L. Tang, "AlGaAs Photonic Integrated Circuits Fabricated Using Chemically Assisted Ion Beam Etching," Paper. OE10.4/ThUU4 in *Conf. Digest of IEEE LEOS Annual Meeting*, Boston, p. 169 (1990).
74. J. B. D. Soole, A. Scherer, H. P. LeBlanc, N. C. Andreadakis, R. Bhat, and M. A. Koza, "Monolithic InP-Based Grating Spectrometer for Wavelength-Division Multiplexed Systems at 1.5 μm," *Electron. Lett.* **27**:132–134 (1991).
75. C. Cremer, G. Ebbinghaus, G. Heise, R. Muller-Nawrath, M. Shienle, and L. Stoll, "Grating Spectrograph in InGaAsP/InP for Dense Wavelength Division Multiplexing," *Appl. Phys. Lett.* **59**:627–629 (1991).
76. T. L. Koch and U. Koren, "Semiconductor Photonic Integrated Circuits," *IEEE J. Quantum Electron.* **QE-27**:641–653 (1991).
77. M. Aoki, M. Suzuki, T. Taniwatari, H. Sano, and T. Kawano, "New Photonic Device Integration by Selective-Area MOVPE and its Application to Optical Modulator/Laser Integration," *Microwave and Opt. Technol. Lett.* **7**(3):132–139 (1994).
78. J. H. Marsh, "Quantum Well Intermixing," *Semicon. Sci. Technol.* **8**:1136–1155 (1993).

79. R. C. Alferness, U. Koren, L. L. Buhl, B. I. Miller, M. G. Young, T. L. Koch, G. Raybon, and C. A. Burrus, "Broadly Tunable InGaAsP/InP Laser Based on a Vertical Coupler Filter with 57-nm Tuning Range," *Appl. Phys. Lett.* **60**:3209–3211 (1992).
80. E. Van Tomme, P. P. Van Daele, R.G. Baets, and P.E. Lagasse, "Integrated Optic Devices Based on Nonlinear Optical Polymers," *IEEE J. Quantum Electron.* **QE-27**:778 (1991).
81. C. C. Teng, "Traveling-Wave Polymeric Optical Intensity Modulator with More than 40 GHz of 3-dB Electrical Bandwidth," *Appl. Phys. Lett.* **60**:1538–1540 (1992).
82. D. Chen, H. R. Fetterman, A. Chen, W. H. Steier, L. R. Dalton, W. Wang, and Y. Shi, "Demonstration of 110 GHz Electro-Optic Polymer Modulators," *Appl. Phys. Lett.* **70**(25):3335–3337 (1997).
83. Y. -H. Kuo, J. Luo, W. H. Steier, and A. K. -Y. Jen, "Enhanced Thermal Stability of Electrooptic Polymer Modulators Using the Diels–Alder Crosslinkable Polymer," *IEEE Phot. Tech. Lett.* **18**(1):175–177 (2006).
84. D. Marcuse, "Bending Losses of the Symmetric Slab Waveguide," *Bell Syst. Tech. J.* **50**:2551–2563 (1971).
85. M. Heiblum and J. H. Harris, "Analysis of Curved Waveguides by Conformal Transformation," *IEEE J. Quantum Electron.* **QE-11**:75–83 (1975).
86. Yu. M. Vlasov, IBM Corporation. Private communication.
87. J. N. Polkay and G. L. Mitchell, "Metal-Clad Planar Dielectric Waveguide for Integrated Optics," *J. Opt. Soc. Am.* **64**:274–279 (1974).
88. P. G. Suchoski, T. K. Findakly, and F. J. Leonberger, "Stable Low-Loss Proton-Exchanged LiNbO₃ Waveguide Devices with No Electro-Optic Degradation," *Opt. Lett.* **13**:1050 (1988).
89. M. K. Smit, "New Focusing and Dispersive Planar Component Based on an Optical Phased Array," *Electron. Lett.* **7**(24):385–386 (1988).
90. H. Takahashi, S. Suzuki, K. Kato, and I. Nishi, "Arrayed-Waveguide Grating for Wavelength Division Multi/Demultiplexer with Nanometre Resolution," *Electron. Lett.* **26**(2):87–88 (1990).
91. C. Dragone, "An N X N Optical Multiplexer Using a Planar Arrangement of Two Star Couplers," *IEEE Photon. Tech. Lett.* **3**:812–815 (1991).
92. K. Takada, M. Abe, T. Shibata, and K. Okamoto, "5 GHz-Spaced 4200-Channel Two-Stage Tandem Demultiplexer for Ultra-Multi-Wavelength Light Source Using Supercontinuum Generation," *Electron. Lett.* **38**(12):575–576 (2002).
93. B. E. Little, S. T. Chu, P. P. Absil, J. V. Hryniewicz, F. G. Johnson, F. Seiferth, D. Gill, V. Van, O. King, and M. Trakalo, "Very High-Order Microring Resonator Filters for WDM Applications," *IEEE Phot. Tech. Lett.* **16**(10):2263–2265 (2004).
94. D.W. Dolfi and T.R. Ranganath, "50 GHz Velocity-Matched Broad Wavelength Ti:LiNbO₃ Mach-Zehnder Modulator with Multimode Active Section," *Electron. Lett.* **28**:1197–1198 (1992).
95. R. A. Becker, C. E. Woodward, F. J. Leonberger, and R. C. Williamson, "Wideband Electrooptic Guided-Wave Analog-to-Digital Converters," *IEEE Proc.* **72**:802 (1984).
96. R. C. Alferness, "Polarization-Independent Optical Directional Coupler Switch Using Weighted Coupling," *Appl. Phys. Lett.* **35**:748–750 (1979).
97. F. Heismann, L. L. Buhl, and R. C. Alferness, "Electro-Optically Tunable Narrowband Ti:LiNbO₃ Wavelength Filters," *Electron. Lett.* **23**:572–573 (1987).
98. M. Suzuki, Y. Noda, H. Tanaka, S. Akiba, Y. Kushiro, and H. Isshiki, "Monolithic Integration of InGaAsP/InP Distributed Feedback Laser and Electroabsorption Modulator by Vapor Phase Epitaxy," *IEEE J. Lightwave Technol.* **LT-5**:1277–1285 (1987).
99. H. Soda, M. Furutsu, K. Sato, N. Okazaki, Y. Yamazaki, H. Nishimoto, and H. Ishikawa, "High-Power and High-Speed Semi-Insulating BH Structure Monolithic Electroabsorption Modulator/DFB Laser Light Source," *Electron Lett.* **26**:9–10 (1990).
100. G. P. Agrawal and N. K. Dutta, *Long-Wavelength Semiconductor Lasers*, New York: Van Nostrand Reinhold Company Inc. (1986).
101. R. Parke, D. F. Welch, A. Hardy, R. Lang, D. Muhuys, S. O'Brien, K. Dzurko, and D. Scifres, "2.0 W CW, Diffraction-Limited Operation of a Monolithically Integrated Master Oscillator Power Amplifier," *IEEE Phot. Tech. Lett.* **5**:297–300 (1993).
102. T. L. Koch, F. S. Choa, U. Koren, R. P. Gnall, F. Hernandez-Gil, C. A. Burrus, M. G. Young, M. Oron, and B. I. Miller, "Balanced Operation of an InGaAs/InGaAsP Multiple-Quantum-Well Integrated Heterodyne Receiver," *IEEE Phot. Tech. Lett.* **2**:577–580 (1990).

103. T. L. Koch and U. Koren, "Semiconductor Photonic Integrated Circuits," *IEEE J. Quantum Electron.* **QE-27**:641–653 (1991).
104. P. D. Magill, K. C. Reichman, R. Jopson, R. Derosier, U. Koren, B. I. Miller, M. Young, and B. Tell, "1.3 Tbit·km/s Transmission Through Non-Dispersion Shifted Fiber by Direct Modulation of a Monolithic Modulator/Laser," Paper PD9 in *Technical Digest of OFC '92*, San Jose, CA: OSA, pp. 347–350 (1992).
105. S. K. Korotky, J. J. Vaselka, C. T. Kemmerer, W. J. Minford, D. T. Moser, J. E. Watson, C. A. Mattoe, and P. L. Stoddard, "High-Speed, Low Power Optical Modulator with Adjustable Chirp Parameter," Paper TuG2 in *Tech. Digest of 1991 Integrated Photonics Research Topical Meeting*, OSA, p. 53 (1991).
106. H. Gnauck, S. K. Korotky, J. J. Vaselka, J. Nagel, C. T. Kemmerer, W. J. Minford, and D. T. Moser, "Dispersion Penalty Reduction Using an Optical Modulator with Adjustable Chirp," *IEEE Phot. Tech. Lett.* **3**:916–918 (1991).
107. J. E. Zucker, K. L. Jones, M. A. Newkirk, R. P. Gnall, B. I. Miller, M. G. Young, U. Koren, C. A. Burrus, and B. Tell, "Quantum Well Interferometric Modulator Monolithically Integrated with 1.55 μm Tunable Distributed Bragg Reflector Laser," *Electron. Lett.* **28**:1888–1889 (1992).
108. C. R. Doerr, L. Zhang, S. Chandrasekhar, and L. L. Buhl, "Monolithic DQPSK Receiver in InP With Low Polarization Sensitivity," *Phot. Tech. Lett.* **19**(21):1765–1767 (2007).
109. S. Corzine, P. Evans, M. Kato, G. He, M. Fisher, M. Raburn, A. Dentai, et al., "10-Channel x 40Gb/s per Channel DQPSK Monolithically Integrated InP-Based Transmitter PIC," Paper PDP18 in the *Tech. Digest of the 2008 Optical Fiber Communication Conference (OFC 2008)*, San Diego, Calif. (2008).
110. G. S. Maurer, G. T. Tremblay, and S. R. McKenzie, "Transmitter Design Alternatives for CATV Distribution in FITL Deployment," *Proc. SPIE OE/Fibers* (1992).
111. R. B. Childs and V. A. O'Byrne, "Multichannel AM Video Transmission Using a High Power Nd:YAG Laser and a Linearized External Modulator," *IEEE J. on Selected Areas of Comm.* **8**:1369 (1990).
112. E. Watson, M. A. Milbrodt, K. Bahadori, M. F. Dautartas, C. T. Kemmerer, D. T. Moser, A. W. Schelling, T. O. Murphy, J. J. Vesselka, and D. A. Herr, "A Low-Voltage 8×8 TiLiNbO₃ Switch with a Dilated-Benes Architecture," *IEEE/OSA J. Lightwave Technol.* **8**:794–801 (1990).
113. S. S. Bergstein, A. F. Ambrose, B. H. Lee, M. T. Fatehi, E. J. Murphy, T. O. Murphy, G. W. Richards, et al., "A Fully Implemented Strickly Non-Blocking 16×16 Photonic Switching System," Paper PD30 in *Tech. Digest of OFC/IOOC '93*, San Jose, Calif, pp. 123–126 (1993).
114. H. Lefevre, *The Fiber-Optic Gyroscope*, Norwood, MA: Artech House Inc. (1993).
115. R. B. Smith, "Fiber-Optic Gyroscopes 1991: A Bibliography of Published Literature," *SPIE Proc.* **1585**:464–503 (1991).
116. S. Ezekiel and H.J. Arditty, eds., "Fiber-Optic Rotation Sensors and Related Technologies," Proc. First International Conference, *Springer Series in Optical Sciences* **32** (1981).
117. E. Udd, ed., "Fiber Optic Gyros: 10th Anniversary Conference," *Proc. SPIE* **719** (1986).
118. S. Ezekiel and E. Udd, eds., "Fiber Optic Gyros: 15th Anniversary Conference," *Proc. SPIE* **1585** (1991).
119. R. A. Bergh, H. C. Lefevre, and H. J. Shaw, "An Overview of Fiber-Optic Gyroscopes," *J. Lightwave Technol.* **LT-2**:91–107 (1984).
120. H. C. Lefevre, S. Vatoux, M. Papuchon, and C. Puech, "Integrated Optics: A Practical Solution for the Fiber-Optic Gyroscope," *Proc. SPIE* **719**:101–112 (1986).
121. Y. A. Akulova, G. A. Fish, P. -C. Koh, C. L. Schow, P. Kozodoy, A. P. Dahl, S. Nakagawa, et al., This Paper appears in: "Selected Topics in Quantum Electronics," *IEEE J. Sel. Topics In Quant. Electron.* **JSTQE** **8**(6):1349–1357 (2002).
122. See, for example, M. K. Smit, "Progress in AWG Design and Technology," Proc. WPOFC 05 Palermo, Italy, pp. 26–31 (2005).
123. R. Nagarajan, C. Joyner, R. Schneider, Jr. J. Bostak, T. Butrie, A. Dentai, V. Dominic et al., "Large-Scale Photonic Integrated Circuits," *IEEE J. Sel. Topics Quantum Electron.* **11**:50–65 (2005).
124. Aiki, M. Nakamura, and J. Umeda, "A Frequency-Multiplexing Light Source with Monolithically Integrated Distributed-Feedback Diode Lasers," *IEEE J. Quantum Electron.* **QE-13**:220–223 (1977).
125. A. Huang, C. Gunn, G.-L. Li, Y. Liang, S. Mirsaidi, A. Narasimha, and T. Pinguet, "A 10Gb/s Photonic Modulator and WDM MUX/DEMUX Integrated with Electronics in 0.13 μm SOI CMOS," in Proc. 2006 IEEE International Solid State Circuits Conference, ISSCC 2006, pp. 922–929 (2006).
126. L. Liao, D. Samara-Rubio, M. Morse, A. Liu, D. Hodge, D. Rubin, U. D. Keil, and T. Franck, "High Speed Silicon Mach-Zehnder Modulator," *Opt. Exp.* **13**:3129–3135 (2005).

127. P. Bernasconi, W. Yang, L. Zhang, N. Sauer, L. Buhl, I. Kang, S. Chandrasekhar, and D. Neilson, "Monolithically Integrated 40Gb/s Wavelength Converter with Multi-Frequency Laser," Paper PDP16 in Proc. of OFC/NFOEC 2005, OSA (2005).
128. M. L. Masanovic, V. Lal, J. A. Summers, J. S. Barton, E. J. Skogen, L. G. Rau, L. A. Coldren, and D. J. Blumenthal, "Widely Tunable Monolithically Integrated All-Optical Wavelength Converters in InP," *IEEE J. Lightwave Technol.* **23**(3):1350–1362 (2005).
129. R. Nagarajan, M. Kato, J. Pleumeekers, P. Evans, D. Lambert, A. Chen, V. Dominic, "Single Chip, 40-Channel x 40Gbit/s per Channel, InP Transmitter Photonic Integrated Circuit," Paper CPDB3 in Proc. of 2006 Conference on Lasers and Electro-Optics/Quantum Electronics and Laser Science Conference and Photonic Applications Systems Technologies, CLEO/QELS 2006, OSA (2006).

Tom D. Milster

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

Tomasz S. Tkaczyk

*Department of Bioengineering
Rice University
Houston, Texas*

22.1 GLOSSARY

A, B, C, D	constants
$A(r, z)$	converging spherical wavefront
c	curvature
D	diffusion constant
d	diffusion depth
EFL	effective focal length
f	focal length
g	gradient constant
h	radial distance from vertex
i	imaginary
k	conic constants
k	wave number
LA	longitudinal aberration
l_0	paraxial focal length
M	total number of zones
NA	numerical aperture
n	refractive index
r	radial distance from optical axis
r_{mask}	mask radius
r_m	radius of the m th zone
t	fabrication time
\bar{u}	slope
W_{ijk}	wavefront function
X	shape factor

x, y	Cartesian coordinates
y	height
Z	sag
z	optical axis
Δ	relative refractive difference
ρ	propagation distance
λ	wavelength
$\bar{\sigma}$	$\sigma_{\text{rms}}/2y$
σ_{rms}	rms wavefront error
Φ	phase
ψ	special function
FOV	field of view

22.2 INTRODUCTION

Optical components come in many sizes and shapes. A class of optical components that has become very useful in many applications is called micro-optics. We define micro-optics very broadly as optical components ranging in size from several millimeters to several hundred microns. In many cases, micro-optic components are designed to be manufactured in volume, thereby reducing cost to the customer. The following paragraphs describe micro-optic components that are potentially useful for large-volume applications. The discussion includes several uses of micro-optics, design considerations for micro-optic components, molded glass and plastic lenses, distributed-index planar lenses, micro-Fresnel lenses, laser printing, grayscale lithography, diamond turning and micromilling, and liquid tunable lenses. For further information the reader is directed to Refs. 1–3.

22.3 USES OF MICRO-OPTICS

Micro-optics are becoming an important part of many optical systems. This is especially true in systems that demand compact design and form factor. Some optical fiber-based applications include fiber-to-fiber coupling, laser diode-to-fiber connections, LED-to-fiber coupling, and fiber-to-detector coupling. Microlens arrays are useful for improving radiometric efficiency in focal-plane arrays, where relatively high numerical aperture (NA) microlenslets focus light onto individual detector elements. Microlens arrays can also be used for wavefront sensors, where relatively low-NA lenslets are required. Each lenslet is designed to sample the input wavefront and provide a deviation on the detector plane that is proportional to the slope of the wavefront over the lenslet area. Micro-optics are also used for coupling laser diodes to waveguides and collimating arrays of laser diodes. An example of a large-volume application of micro-optics is data storage, where the objective and collimating lenses are only a few millimeters in diameter.⁴ Recently micro-optics is also widely used in medical applications like endo-microscopy including confocal, multiphoton microscopy and optical coherence tomography (OCT). It also has a strong position in consumer market in applications like cell phone cameras, Blu-Ray readers, etc.

22.4 MICRO-OPTICS DESIGN CONSIDERATIONS

Conventional lenses made with bulk elements can exploit numerous design parameters, such as the number of surfaces, element spacings, and index/dispersion combinations, to achieve performance requirements for NA, operating wavelength, and field of view. However, fabricators of micro-optic

lenses seek to explore molded or planar technologies, and thus the design parameters tend to be more constrained. For example, refractive microlenses made by molding, ion exchange, mass transport process resemble single-element optics. Performance of these lenses is optimized by manipulating one or possibly two radii, the thickness, and the index or index distribution. Index choices are limited by the available materials. Distributed-index and graded-index lenses have a limited range of index profiles that can be achieved. Additional performance correction is possible by aspherizing one or both surfaces of the element. This is most efficiently done with the molding process, but molded optics are difficult to produce when the diameter of the lens is less than 1.0 mm. In general, one or two aberrations may be corrected with one or two aspheres, respectively.

Due to the single-element nature of microlenses, insight into their performance may be gained by studying the well-known third-order aberrations of a thin lens in various configurations. Lens bending and stop shift are the two parameters used to control aberrations for a lens of a given power and index. Bending refers to distribution of power between the two surfaces, i.e., the shape of the lens, as described in R. Barry Johnson's Chap. 17, "Lenses." The shape is described by the shape factor X that is

$$X = \frac{C_1 + C_2}{C_1 - C_2} \quad (1)$$

where C_1 and C_2 are the curvatures of the surfaces. The third-order aberrations as a function of X are shown in Fig. 1. These curves are for a lens with a focal length of 10.0 mm, an entrance pupil diameter of 1.0 mm, field angle $\bar{u} = 20^\circ$, an optical index of refraction of 1.5, $\lambda = 0.6328 \mu\text{m}$, and the object at infinity. For any given bending of the lens, there is a corresponding stop position that eliminates coma,⁵ and this is the stop position plotted in the figure. The stop position for which coma is zero is referred to as the *natural* stop shift, and it also produces the least curved tangential field for the given bending. Because the coma is zero, these configurations of the thin lens necessarily satisfy the Abbe sine condition. When the stop is at the lens (zero stop shift), the optimum shape to eliminate coma is approximately convex-plano ($X = +1$) with the convex side toward the object. The optimum shape is a function of the index, and the higher the index, the more the lens must be bent into a meniscus. Spherical aberration is minimized with the stop at the lens, but astigmatism is near its maximum. It is interesting to note that biaspheric objectives for data storage tend toward the convex-plano shape.

Astigmatism can be eliminated for two different lens-shape/stop-shift combinations, as shown in Fig. 1. The penalty is an increase in spherical aberration. Note that there is no lens shape for which

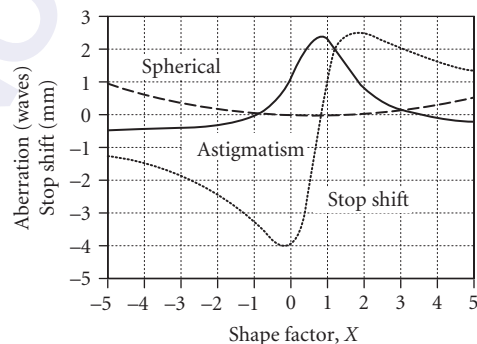


FIGURE 1 Third-order aberrations as a function of the shape factor, or bending, of a simple thin lens with focal length 10.0 mm, entrance pupil diameter of 1.0 mm, field angle 20° , $n = 1.5$, and object at infinity. The stop position shown is the *natural* stop shift, that is, the position that produces zero coma.

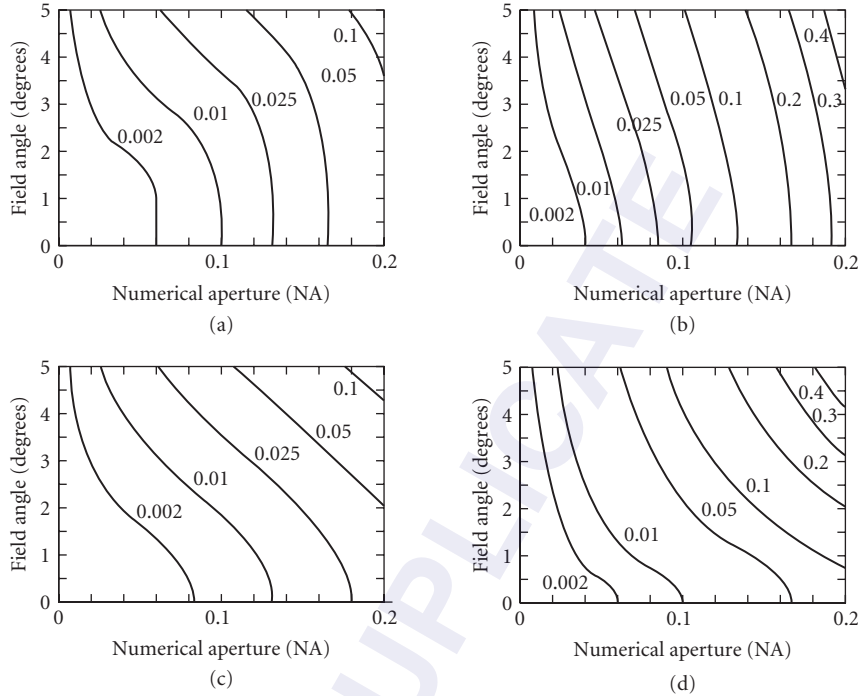


FIGURE 2 Contours of normalized rms wavefront deviation, $\bar{\sigma} = 1000\lambda\sigma_{\text{rms}}/2y$, versus field angle and NA, where $2y$ is the diameter of the stop. The stop is located at the lens. The focus is adjusted to give minimum rms deviation of the wavefront, so effects of Petzval curvature are not included. (a) $X = 1, n = 1.5$; (b) $X = -1, n = 1.5$; (c) $X = 1, n = 3.0$; (d) $X = -1, n = 3.0$.

spherical, coma, and astigmatism are simultaneously zero in Fig. 1, that is, there is no aplanatic solution when the object is at infinity. The aplanatic condition for a thin lens is only satisfied at finite conjugates.

The plano-convex shape ($X = -1$) that eliminates astigmatism is particularly interesting because the stop location is in front of the lens at the optical center of curvature of the second surface. All chief rays are normally incident at the second surface. Thus, the design is monocentric.⁶ (obviously, the first surface is not monocentric with respect to the center of the stop, but it has zero power and only contributes distortion.)

Two very common configurations of micro-optic lenses are $X = +1$ and $X = -1$ with the stop at the lens. Typically, the object is at infinity. In Fig. 2, we display contours of normalized rms wavefront deviation, $\bar{\sigma} = \sigma_{\text{rms}}/2y$, versus \bar{u} and NA, where $2y =$ diameter of the stop. Aberration components in σ_{rms} include third-order spherical, astigmatism, and coma. The focus is adjusted to give minimum rms deviation of the wavefront, so effects of Petzval curvature are not included. Tilt is also subtracted. As NA or field angle is increased, rms wavefront aberration increases substantially. The usable field of view of the optical system is commonly defined in terms of Maréchal's criterion⁷ as field angles less than those that produce $2y\bar{\sigma}/1000\lambda \leq 0.07$ wave. For example, if the optical system operates at $2y = 1.0$ mm, $\lambda = 0.6328$ μm , $\text{NA} = 0.1$, $X = +1$, $n = 1.5$, and $\bar{u} = 2^\circ$, the wavefront aberration due to third-order contributions is

$$\sigma_{\text{rms}} = \frac{2y\bar{\sigma}}{1000\lambda} \approx \frac{(1.0 \times 10^{-3} \text{ m})(0.015)}{(10^3)(0.6328 \times 10^{-6} \text{ m/wave})} = 0.024 \text{ wave} \quad (2)$$

which is acceptable for most situations. Note that the configuration for $X = -1$ yields $\sigma_{\text{rms}} \approx 0.079$ wave, which is beyond the acceptable limit. When large values of σ_{rms} are derived from Fig. 2, care must be taken in interpretation of the result because higher-order aberrations are not included in the calculation. Also, if field curvature is included in the calculation, the usable field of view is significantly reduced.

Coma and astigmatism are only significant if the image field contains off-axis locations. In many laser applications, like laser diode collimators, the micro-optic lens is designed to operate on axis with only a very small field of view. In this case, spherical aberration is very significant. A common technique that is used to minimize spherical aberration is to aspherize a surface of the lens. Third-, fifth-, and higher orders of spherical aberration may be corrected by choosing the proper surface shape. In some lens design codes, the shape is specified by

$$Z = \frac{ch^2}{1 + \sqrt{1 - (1+k)c^2h^2}} + Ah^4 + Bh^6 + Ch^8 + Dh^{10} \quad (3)$$

where Z is the sag of the surface, c is the base curvature of the surface, k is the conic constant ($k = 0$ is a sphere, $k = -1$ is a paraboloid, etc.), and $h = \sqrt{x^2 + y^2}$ is the radial distance from the vertex. The A , B , C , and D coefficients specify the amount of aspheric departure in terms of a polynomial expansion in h .

When a plane-parallel plate is inserted in a diverging or converging beam, such as the window glass of a laser diode or an optical disk, spherical aberration is introduced. The amount of aberration depends on the thickness of the plate, the NA of the beam, and to a lesser extent the refractive index of the plate,⁸ as shown in Fig. 3. The magnitude of all orders of spherical aberration is linearly proportional to the thickness of the plate. The sign is opposite that of the spherical aberration introduced by an $X = +1$ singlet that could be used to focus the beam through the plate. Therefore, the aspheric correction on the singlet compensates for the difference of the spherical aberration of the singlet and the plate. This observation follows the fact that minimum spherical aberration without aspheric correction is achieved with the smallest possible air gap between the lens and the plate. For high-NA singlet objectives, one or two aspheric surfaces are added to correct the residual spherical aberration.

Considerations on High-Performance Miniature Systems

In recent years, several designs of high-performance miniature optics were implemented.^{9–13} Since one of the major applications is miniature microscopy, many systems have similar characteristics. The diameter of miniature objectives is in the range of 1 to 8 mm, consists of multiple lenses (upto 10), NA is 0.4–1.0 and field of view (FOV) 250 to 500 μm . All these parameters place high requirements on system assembly and fabrication technologies. Examples of manufacturing techniques

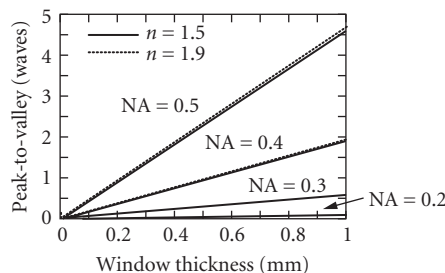


FIGURE 3 Effect of a window on wavefront distortion at $\lambda = 830$ nm.

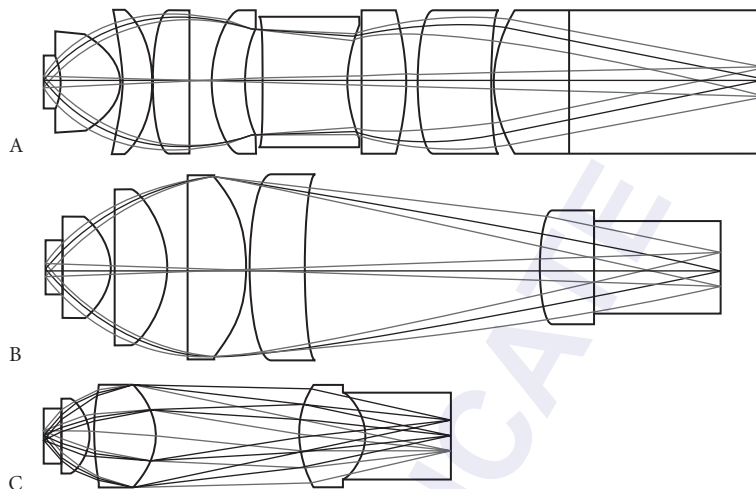


FIGURE 4 Comparison of designs of miniature NA = 1.0, 205 μm FOV microscope objectives. Design A¹⁰ is an all glass, spherical lenses system. Design B¹¹ uses a plastic injection molded objective lens made in Zeonex. Design C¹³ consists of a leading spherical glass lens followed by two aspherical Zeonex plastic lenses.

include (1) grinding and polishing methods, (2) injection glass and plastic molding, (3) grayscale lithography, and (4) direct diamond turning/micromilling. Grinding and polishing in glass is the most costly technique and requires a large number of components to provide good system correction (lenses are usually spherical). All other mentioned technologies allow easy application of aspherical lenses that significantly reduces the number of parts. On the other hand, these methods often rely on application of photosensitive polymers, solgels, or plastics, which limits selection of materials with different Abbe numbers. This limitation makes correction for chromatic aberrations more difficult.

An example of an effective approach is to build a high NA system as a combination of glass spherical lenses and aspherical plastic/polymer lenses. A comparison of three optical designs with single wavelength, NA = 1.0, FOV = 250 μm is shown in Fig. 4. Design A¹⁰ is an all glass, spherical lenses system. Design B¹¹ uses aspheric plastic injection molded lenses made in Zeonex. Design C¹³ consists of a leading spherical glass lens followed by two aspherical Zeonex plastic lenses, which can be injection molded or diamond turned. All objectives in Fig. 4 are presented in scale and the clear aperture of the smallest design (C) is 2.75 mm. These designs show progressive decrease of complexity based on material and surface shape choices.

A common method used for tolerance analysis of multicomponent systems is Monte Carlo (MC) simulation that statistically adjusts design parameters like radius error, thickness error, surface form error, surface decenter, surface tilt, and refractive index tolerances.^{11,13} MC can apply a normal distribution of errors and then evaluate the performance of the system based on the rms spot size or wavefront error. Possible compensation parameters can include the object position, image position, and object shape which is for three-dimensional volume imaging in biological applications. The acceptance criteria can be defined based on the percentage of results that have a rms spot size equal to or less than the diffraction limited Airy disk spot size.

As mentioned above, one of the most significant issues in building high-performance miniature systems is precision assembly. Several alignment-free approaches were recently proposed to simplify the manufacturing process. They include techniques like: micro-optical table (MOT),¹⁴ fabrication of kinematic mechanical features embedded in optical components,^{11,15,16} and self centering mounts.¹³ All these techniques rely on a zero-alignment concept. In practical terms, the

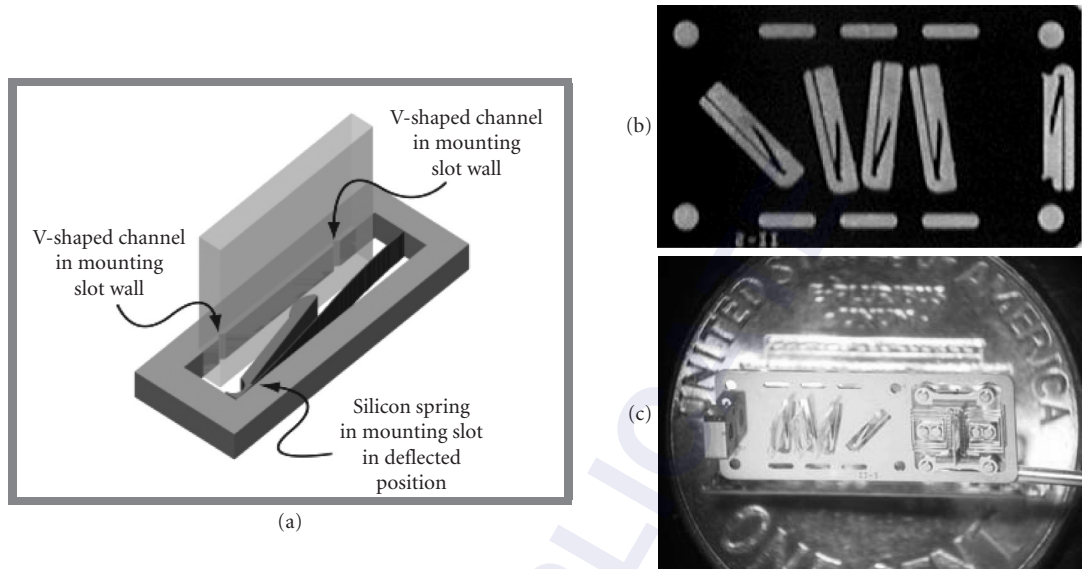


FIGURE 5 (a) Concept of MOT assembly;¹⁴ (b) MOT platform made with DXRL process; and (c) NA = 0.4 miniature microscope assembled on MOT platform.¹⁹

zero-alignment concept translates into assembly errors that are smaller than the tolerances on the performance of the optical system. Very low assembly errors are achieved through positioning features on each optomechanical and optical component. One technology to fabricate mechanical components is deep reactive ion etching (DRIE)¹⁷ in silicon. A less expensive and more robust alternative is deep x-ray lithography (DXRL).¹⁸ An example of MOT platform fabricated using DXRL technology in nickel-steel alloy is shown in Fig. 5*b*. The DXRL technology delivers sub-micron assembly precision and helps with packaging. Spring and groove features in MOT shown in the figure allow inserting optical components fabricated with grayscale lithography on thin glass substrates. Positioning precision is provided by assembly features on both MOT and optics. (Grayscale lithography allows fabrication of lens and mechanical features in one step.) Figure 5 shows (a) the MOT assembly concept, (b) MOT platform made with DXRL process, and (c) assembled NA = 0.4 miniature microscope.¹⁹ Table 1 provides assembly precision for MOT technology with DRIE and lithographically fabricated opto-mechanics.¹⁴ These parameters must be accommodated in optical design.

The principle of a kinematic mount fabricated with injection molding or a lithographic process is shown in Fig. 6. Note that V- or U-shaped grooves can be made in the same process as a lens. Precision spheres are then used for stacking consecutive lenses or lens layers, stops, and spacer

TABLE 1 Assembly Precision for MOT Technology with DRIE and Lithographically Fabricated Optomechanics¹⁴

Position	Measured Position
Translation along slot	$3 \mu\text{m} \pm 1 \mu\text{m}$
Yaw rotation (left)	$-5 \text{ arc min} \pm 2 \text{ arc min}$
Yaw rotation (right)	$-5 \text{ arc min} \pm 2 \text{ arc min}$
Pitch rotation	$0.44^\circ \pm 0.02^\circ$

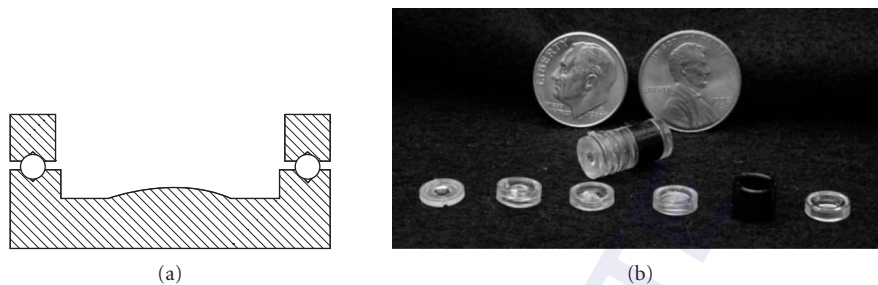


FIGURE 6 (a) The concept of kinematic lens mount and (b) NA = 1.0 microscope objective assembled using kinematic mount embedded in plastic injection molded lenses.¹¹

components.¹¹ Figure 6a shows the concept of kinematic mount and Fig. 6b shows NA = 1.0 microscope objective assembled using kinematic mounts embedded in plastic injection molded lenses.¹¹ (The design of the system is presented in Fig. 4b.)

A kinematic approach was also used to build arrays of complex optical systems. Examples of applications include microchemical chips using microlens arrays¹⁶ (used for illumination and detection) and array of high NA microscopes for high throughput digital telepathology.¹⁵

Another method to position lens components in the objective is to use a self-centering ring, which engages the lens surface, flexing away as it makes contact, while centering the lens with respect to its optical axis.¹⁴ This self-centering ring eliminates any decentration associated with manufacturing error and allows looser lens tolerance. Figure 7 shows a model self-centering lens spring. This approach was used for assembly of NA = 1.0 plastic-glass objective¹³ (see also Fig. 4c).

22.5 MOLDED MICROLENSES

Molded micro-optic components have found applications in several commercial products, which include compact disk players, bar-code scanners, and diode-to-fiber couplers. Molded lenses become especially attractive when one is designing an application that requires aspheric surfaces.

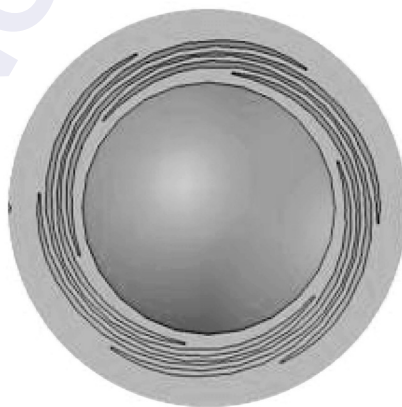


FIGURE 7 Model self-centering lens spring for lens alignment.

Conventional techniques for polishing and grinding lenses tend to be time-expensive and do not yield good piece-to-piece uniformity. Direct molding, on the other hand, eliminates the need for significant grinding or polishing. Another advantage of direct molding is that useful reference surfaces can be designed directly into the mold. The reference surfaces can take the form of flats.²⁰ The reference flats are used to aid in aligning the lens element during assembly into the optical device. Therefore, in volume applications that require aspheric surfaces, molding becomes a cost-effective and practical solution. The molding process utilizes a master mold, which is commonly made by single-point diamond turning and postpolishing to remove tooling marks and thus minimize scatter from the surface. The master can be tested with conventional null techniques, computer-generated null holograms,²¹ or null Ronchi screens.²² Two types of molding technology are described in the following paragraphs. The first is molded glass technology. The second is molded plastic technology.

Molded Glass

One of the reasons glass is specified as the material of choice is thermal stability. Other factors include low birefringence, high transmission over a broad wavelength band, and resistance to harsh environments.

Several considerations must be made when molding glass optics. Special attention must be made to the glass-softening point and refractive index.²³ The softening point of the glass used in molded optics is lower than that of conventional components. This enables the lenses to be formed at lower temperatures, thereby increasing options for cost-effective tooling and molding. The refractive index of the glass material can influence the design of the surface. For example, a higher refractive index will reduce the surface curvature. Smaller curvatures are generally easier to fabricate and are thus desirable.

An illustration is Coming's glass molding process.²³ The molds that are used for aspheric glass surfaces are constructed with a single-point diamond turning machine under strict temperature and humidity control. The finished molds are assembled into a precision-bored alignment sleeve to control centration and tilt of the molds. A ring member forms the outside diameter of the lens, as shown in Fig. 8. The glass material, which is called a preform, is inserted between the molds. Two keys to accurate replication of the aspheric surfaces are forming the material at high glass viscosity and maintaining an isothermal environment. After the mold and preform are heated to the molding temperature, a load is applied to one of the molds to press the preform into shape. After molding, the assembly is cooled to below the glass transformation point before the lens is removed. Optical performance characteristics of the finished lens are determined by the quality of the mold surfaces, the glass material, and the preform volume, which also determines the thickness of the lens when pressed.

An alternative process is used at Kodak, Inc., where molded optics are injection molded and mounted into precision lens cells.²⁴ In this process, a tuned production mold can reproduce intricate

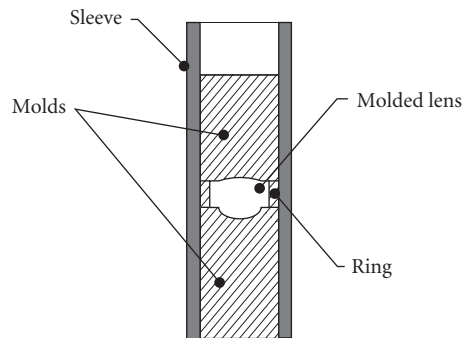


FIGURE 8 Mold for glass optics.

TABLE 2 Preferred and Possible Tolerances for Molded Glass Components²⁴

	Preferred	Possible
Center thickness (mm)	10.00 max 0.40 min ± 0.030 tol	25.00 0.35 ± 0.015
Diameter (mm)	25.00 max 4.00 min ± 0.10 tol	50.00 2.00 ± 0.01
Diameter of lens beyond clear aperture (mm)	2.00	0.50
Surface quality	80–50	40–20
Axis alignment	3×10^{-3} rad	2×10^{-3} rad
Radius (mm)— best fit sphere	5 to ∞	2 to ∞
Slope (λ /mm)	50 max	100 max
Wavelengths (λ) departure from BFS	≤ 250	≤ 500

mounting datum features and extremely well-aligned optics. It can also form a stop, baffle, or a film-plane reference in the system. Table 2 lists preferred and possible tolerances for molded glass components. The Kodak process has been tested with over 50 optical glasses, which include both crowns and flints. This provides a wide index-of-refraction range, $1.51 < n < 1.85$, to choose from.

Most of the molded glass microlenses manufactured to date have been designed to operate with infrared laser diodes at $\lambda = 780$ to 830 nm. The glass used to make the lenses is transparent over a much broader range, so the operating wavelength is not a significant factor if designing in the visible or near infrared. Figure 9 displays a chart of the external transmission of several optical materials versus wavelength. LaK09 (curve B) is representative of the type of glass used in molded optics. The external transmission from 300 to over 2200 nm is limited primarily by Fresnel losses due to the relatively high index of refraction ($n = 1.73$). The transmission can be improved dramatically with antireflection coatings. Figure 10 displays the on-axis operating characteristics of a Corning 350110 lens, which is used for collimating laser diodes. The rms wavefront variation and effective focal length (EFL) are shown versus wavelength. The highest aberration is observed at shorter wavelengths. As the wavelength increases, the EFL increases, which decreases the NA slightly. Table 3 lists several optical properties of molded optical materials. The trend in molded glass lenses is to make smaller, lighter, and higher NA components.²⁵ Reduction in mass and size allows for shorter access times in optical data storage devices, and higher NA improves storage density in such devices.

Molded Plastic

Molded plastic lenses are an inexpensive alternative to molded glass. In addition, plastic components are lighter than glass components. However, plastic lenses are more sensitive to temperatures and environmental factors. The most common use of molded plastic lenses is in compact disk (CD) players.

Precision plastic microlenses are commonly manufactured with injection molding equipment in high-volume applications. However, the classical injection molding process typically leaves some inhomogeneities in the material due to shear and cooling stresses.²⁶ Improved molding techniques can significantly reduce variations, as can compression molding and casting. The current state of the art in optical molding permits master surfaces to be replicated to an accuracy of roughly one fringe

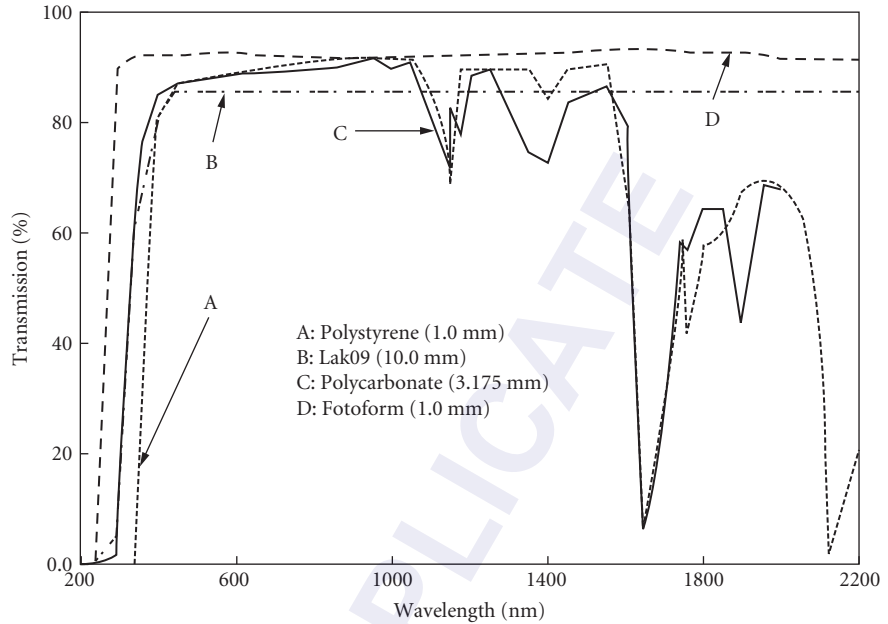


FIGURE 9 External transmission of several optical materials versus wavelength. (A) Polystyrene 1.0 mm thick, which is used for molded plastic lenses;²⁶ (B) LaK09 10.0 mm thick, which is used for molded glass lenses;¹⁰⁵ (C) Polycarbonate 3.175 mm thick, which is used for molded plastic lenses;²⁶ and (D) Fotoform glass 1.0 mm thick, which is used in the production of SMILE lenses.¹⁰⁹

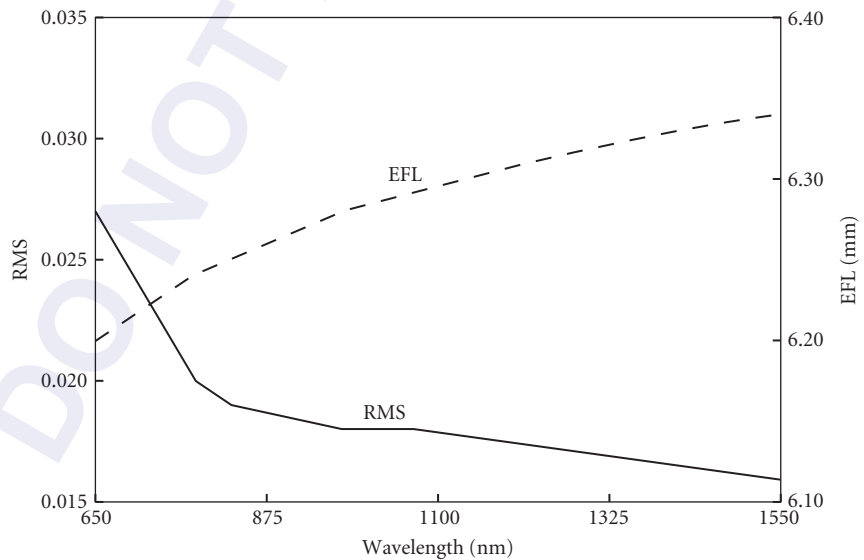


FIGURE 10 On-axis operating characteristics versus wavelength of a Corning 350110 lens, which is a molded glass aspheric used for collimating laser diodes.¹⁰⁶

TABLE 3 Properties of Materials Used for Molding Micro-Optics

Property	PMMA (Acrylic)	PMMA (Imide)	SSMA	Polycarbonate	Polystyrene	Zeonex	LaK09	BK7
Index (n_d)	1.491	1.528	1.564	1.586	1.589	1.525	1.734	1.517
Abbe # (V_d)	57.4	48	35	30	31	56.2	51.5	64.2
Density (g/mm ³)	1.19	1.21	1.09	1.20	1.06	1.01	4.04	2.51
Max service temp (°C)	72	142	87	121	75	80.0	500	500
Thermal expansion coefficient (1E-6 mm/mm °C)	67.9	-	56.0	65.5	50.0	67	5.5	7.1
Thermal index coefficient (1E-6/°C)	-105	-	-	-107	-	-	6.5	3
Young's modulus (10E4 kg/cm ²)	3.02	-	3.30	2.43	3.16	2.2	11.37	83.1
Impact strength	2	-	3	3	4	-	-	1
Abrasion resistance	4	-	3	3	2	-	-	5
Cost/lb	3	-	2	4	2	-	-	5
Birefringence	2	-	4	3	5	-	-	1

1 = lowest / 5 = highest.

per 25 mm diameter, or perhaps a bit better.²⁷ Detail as small as 5 nm may be transferred if the material properties and processing are optimum and the shapes are modest. Table 4 lists tolerances of injection-molded lenses.²⁸ The tooling costs associated with molded plastics are typically less than those associated with molded glass because of the lower transition temperature of the plastics. Also, the material cost is lower for polymers than for glass. Consequently, the costs associated with manufacture of molded plastic microlenses are much less than those for molded glass microlenses. The index of refraction for the plastics is less than that for the glass lenses, so the curvature of the surfaces must be greater, and therefore harder to manufacture, for comparable NA.

The glass map for molded plastic materials is shown in Fig. 11. The few polymers that have been characterized lie mainly outside the region containing the optical glasses and particularly far from the flint materials.²⁹ Data on index of refraction and Abbe number are particularly difficult to obtain for molded plastic. The material is supplied in pelletized form, so it must first be molded into a form suitable for measurement. The molding process subjects the material to a heating and annealing cycle that potentially affects the optical properties. Typically, the effect of the additional thermal history is to shift the dispersion curve upward or downward, leaving the shape unchanged. A more complete listing of optical plastics and their properties is given in Ref. 26. Additional information

TABLE 4 Injection Molding Tolerances for Plastic Lenses

Focal length	±0.5%
Radius of curvature	±0.5%
Spherical power	2 to 5 f*
Surface quality	60/40 (40/20 possible)
Vertex thickness (in)	±0.0005
Diameter (in per in DIA.)	±0.002–0.0005
Repeatability lens-to-lens	0.1–0.3%

*Tolerances given in optical fringes abbreviated by "f".

Vertex-to-edge thickness ratio

4:1 Difficult to mold

3:1 Moderately easy to mold

2:1 Easy to mold

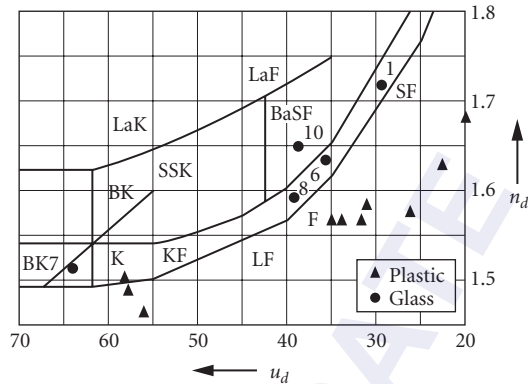


FIGURE 11 Glass map for molded plastic materials, which are shown as triangles in the figure. The few polymers that have been characterized lie mainly outside the region containing the optical glasses and particularly far from the flint materials.²⁷

can be obtained from the *Modern Plastics Encyclopedia*,³⁰ the *Plastics Technology, Manufacturing Handbook and Buyer's Guide*,³¹ and in John D. Lytle's Chap. 3, "Polymeric Optics," in Vol. IV.

Changes in dimension or refractive index due to thermal variations occur in both molded glass and molded plastic lenses. However, the effect is more pronounced in polymer optical systems because the thermal coefficients of refractive index and expansion are ten times greater than for optical glasses, as shown in Table 3. When these changes are modeled in a computer, a majority of the optical systems exhibit a simple defocus and a change of effective focal length and corresponding first-order parameters. An experimental study³² was made on an acrylic lens designed for a focal length of 6.171 mm at $\lambda = 780$ nm and 20°C. At 16°C, the focal length changed to 6.133 mm. At 60°C, the focal length changed to 6.221 mm. Thermal gradients, which can introduce complex aberrations, are a more serious problem. Therefore, more care must be exercised in the design of athermalized mounts for polymer optical systems.

The transmission of two common optical plastics, polystyrene and polycarbonate, are shown in Fig. 9. The useful transmittance range is from 380 to 1000 nm. The transmission curve is severely degraded above 1000 nm due to C-H vibrational overtone and recombination bands, except for windows around 1300 nm and 1500 nm. Sometimes, a blue dye is added to the resins to make the manufactured part appear "water clear," instead of slightly yellowish in color. It is recommended that resins be specified with *no* blue toner for the best and most predictable optical results.²⁶

The shape of the lens element influences how easily it can be manufactured. Reasonable edge thickness is preferred in order to allow easier filling. Weak surfaces are to be avoided because surface-tension forces on weak surfaces will tend to be very indeterminate. Consequently, more strongly curved surfaces tend to have better shape retention due to surface-tension forces. However, strongly curved surfaces are a problem because it is difficult to produce the mold. Avoid clear apertures that are too large of a percentage of the physical surface diameter. Avoid sharp angles on flange surfaces. Use a center/edge thickness ratio less than 3 for positive lenses (or 1/3 for negative lenses). Avoid cemented interfaces. Figure 12 displays a few lens forms. The examples that mold well are C, E, F, and H. Form A should be avoided due to a small edge thickness. Forms A and B should be avoided due to weak rear surfaces. Form D will mold poorly due to bad edge/center thickness ratio. Form G uses a cemented interface, which could develop considerable stress due to the fact that thermal differences may deform the pair, or possibly even destroy the bond.

The plastic injection-molding process may cause flow-induced birefringence in molded optics and significantly reduce system performance, as it creates highly localized regions of refractive

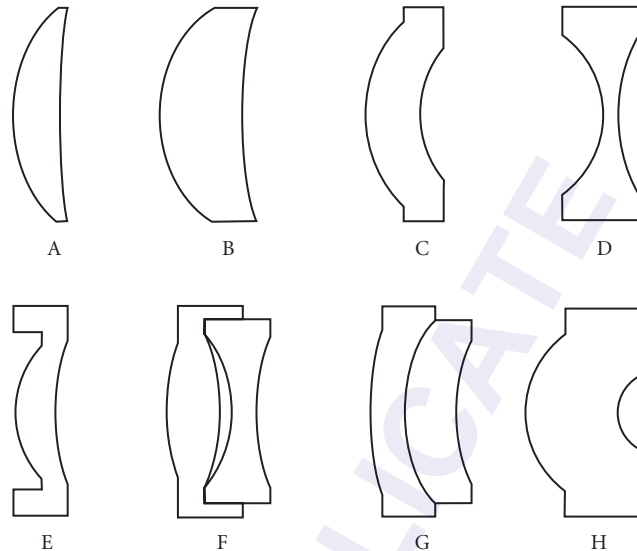


FIGURE 12 Example lens forms for molded plastic lenses. Forms C, E, F, and H mold well. Form A should be avoided due to small edge thickness. Forms A and B should be avoided due to weak rear surfaces. Form D will mold poorly due to bad edge/center ratio. Form G uses a cemented interface, which could develop stress.²⁹

index changes. The effect is more likely for high-power components with small edge thickness. Plane polariscope images of injection-molded lenses are presented in Fig. 13. The very left image is for a lens fabricated in a standard injection-molding process, while the center and right images are for the same lens type but after lowering fill speed and first-stage packing pressure. The center and right image had the same molding parameters, but the position of the optical molding insert pin was corrected. Figure 13 shows how significant birefringence effects may occur in the fabrication process. It also shows that it is possible to reduce birefringence effects through parameter adjustment. These effects should be taken into considerations during lens form design. Note that CA in the picture marks clear aperture of the fabricated lenses.



FIGURE 13 Plane polariscope images of injection-molded lenses. The lens edge thickness is 0.5 mm. The very left lens was fabricated in standard injection-molding process. Center and right lenses were fabricated with lower fill speed and lower first-stage packing pressure. CA is clear aperture of the lens.³³

Since polymers are generally softer than glass, there is concern about damage from ordinary cleaning procedures. Surface treatments, such as diamond films,³⁴ can be applied that greatly reduce the damage susceptibility of polymer optical surfaces.

A final consideration is the centration tolerance associated with aspheric surfaces. With spherical optics, the lens manufacturer is usually free to trade off tilt and decentration tolerances. With aspheric surfaces, this tradeoff is no longer possible. The centration tolerance for molded aspherics is determined by the alignment of the mold halves. A common specification is 4 to 6 μm , although 3 to 4 μm is possible.

22.6 DIAMOND TURNING

Ultra precision diamond tool machining is an important technology for fabrication of miniature optical systems. It can be used in two ways: (1) producing molds for glass and plastic injection molding, and (2) for direct part manufacturing. Diamond turning was originally developed for large-scale aspherical optics (reflective and IR optics) for astronomy and military applications.³⁵ Recently, it took an important place in prototyping and fabrication of small-scale systems. Currently diamond turning is capable of fabricating parts with diameters from a fraction of millimeter to over 500 mm.

Diamond tool machining is a technology using monocrystal diamond-cutting tools assembled on an ultraprecision numerically controlled lathe. Diamond tools have nanometer precision edges, which allow cutting a variety of materials and produce molds for injection molding (in metals) or components like mirrors (metals) and lenses (infrared materials and plastics). The depth of a single cut can be below 1 μm . Depending on the material, rough passes have cuts of 10 μm (metals) to approximately 100 μm (plastics). The finishing cut is usually in the 1 to 5 μm range. Modern diamond-turning instruments use air bearings and high-pressure oil bearings to eliminate direct contact of moving parts. After application of a counter-balancing mass, it is possible to achieve exceptionally high precision machining. Optical interpolators assembled on current instruments allow for 1-nm position monitoring. Effective positioning precision is in the range 10 to 50 nm. Therefore, the process is capable of producing submicrometer form precision and nanometer range roughness.

An important performance parameters is form error which is usually in the range 0.05 to 5 μm , while roughness (Ra) is between 1 and 10 nm. The lowest roughness can be reached for metals, where for high-purity aluminum and nickel-plated aluminum Ra is 1 to 5 nm. Slightly worse Ra of 5 to 10 nm can be obtained for plastics (PMMA, Zeonex), while Ra for infrared crystals is few tens of nanometers. Specifications for diamond-turning machines (for example, UPL 250 from Nanotech) guarantee form error of 0.125 μm over a 75-mm diameter part, and surface roughness of 3 nm (Ra) is achieved for spherical component made with high-purity aluminum. While these parameters are manufacturer specific, it is possible to reach values better than 0.05 μm form error and 1.2 nm roughness (Ra). No polishing is required after diamond cutting, and both form error and roughness depends greatly on the surface size, shape, and material.

While diamond-machining technology has many advantages, it also has two major limitations. One drawback is a periodic character of surface roughness of the fabricated part. Improper surface fabrication may result in diffractive effects. Small cutting depth and low feed rate help to mitigate this issue. A second disadvantage is due to material limitations. Diamond turning can machine nonferrous metals, infrared crystals, and selected polymers. It is not possible to machine glass, due to material microcracking. The list of diamond-turnable materials is shown in Table 5.

Current diamond-turning machines have up to 5° of freedom to position the part/tool. They are all called (in machining jargon axes) *X*, *Y*, *Z*, *C*, and *B*. *X*, *Y*, and *Z* correspond to linear translation of the tool in reference to the part. *C* axis relates to the angular coordinate of a spindle and in consequence angular position of a machined part. The *B* axis controls the angular position of the diamond tool. The variety of axes allows producing aspherical and nonrotational arbitrary

TABLE 5 Summary of Diamond Turnable Materials¹

Material Category	Material Name	Refractive Index ^{36,37}	Abbe Number ^{36,37}	Transmission ² (μm) ^{36,37}	Comments
Metals	Aluminum alloys				
	Brass				
	Copper				
	Gold				
	Nickel				Electroplated
	Silver				
Polymers	Acetyl				
	Acrylic (PMMA)	1.491	57.4	0.35–2.1	
	Fluoroplastic	1.320	92		
	Polycarbonate	1.586	30.1	0.30–2.1	Damages tool quickly
	Polypropylene				
	Polystyrene	1.571	41.2	0.37–2.1	Possible clouding
	Polysulfone	1.634	23.5		
Zeonex (polyolefin)	1.528	55.7			
IR Crystals	Cadmium sulfide*	2.537 (0.023)		0.51–14.8	*Denotes uniaxial crystal materials Refractive index given for birefringent materials is for ordinary wave. Value in parenthesis is material birefringence $\Delta n = n_e - n_o$
	Cadmium telluride	2.817		0.90–30	
	Calcium fluoride	1.562		0.13–12	
	Cesium iodide	1.781		0.25–62	
	Gallium arsenide	4.020		0.90–17.3	
	Germanium	4.052		1.80–15	
	Indium antimonite	5.130		6–25	
	Iridium				
	Lithium niobate*	2.183			
	Magnesium Fluoride*	(–0.076)		0.13–7.7	
	Potassium bromide	1.377		0.20–306	
	Potassium Phosphate	(0.0188)		1.1–6.5	
	Silicon	1.557		0.17–18	
	Sodium chloride	3.478		0.42–4	
	Tellurium dioxide*	1.531		0.5–20	
	Zinc selenide	2.584 (0.288)			
Zinc sulfide	2.591				
	2.536 (0.022)				

¹Values in the table are approximate and a general guidance.

²For IR Crystals transmission is given as a wavelength range of 1-mm-thick sample at 300 K, assuming that the value is larger than 10 percent.³⁶ For polymers transmission is given for 3.2-mm thick sample and at least 10 percent transmission. Note that polymers with listed transmission range have at least 80 percent (3.2-mm sample) transmission in 400 to 1100 nm wavelengths.

optical components, like Alvarez plates³⁸ or arrays of lenses. A diagram of a five-axis diamond machine geometry is presented in Fig. 14.

Only X and Z axes are required for machining aspherical and axially symmetrical parts. For this two-axis mode, the part is mounted on the vacuum chuck of the rotation spindle, while the diamond tool is mounted on the Z -axis stage. While the spindle rotates with constant speed, X and Z move to provide proper surface height (Z) for the related radius coordinate (X). An example of Zeonex aspherical lenses made for a $\text{NA} = 1.0$ miniature microscope cut in this two-axis mode is presented in Fig. 15. The diameter of the lenses is 2.75 mm.

Arbitrary (free) form machining can be obtained in two cutting modes: diamond turning with C axis (X, Z, C configuration) and micromilling (X, Y, Z). The B axis is used for changing tool angle and, in effect, allows increasing optical power of the fabricated components. Diamond turning using the C axis is based on moving the diamond-turning tool in the Z direction synchronously with the rotation

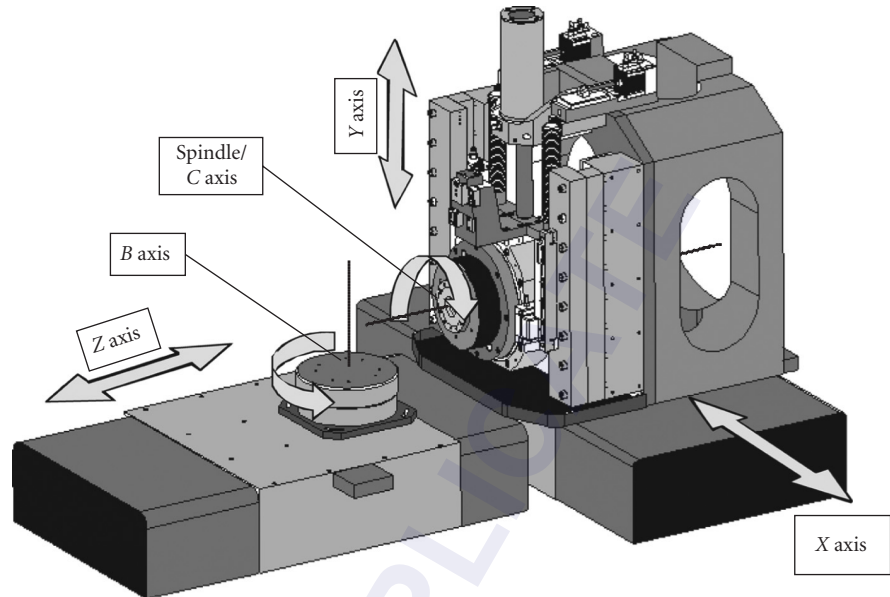


FIGURE 14 A geometry of five-axis diamond-turning machine. (Courtesy of Moore Nanotechnology System, LLC.)

angle of the part controlled with spindle controller (*C* axis). The *C* axis is used to change the spindle rate during the cutting process. *C*-axis machining can be used in slow-slide servo and fast servo modes. The slow-slide servo mode uses the *Z* stage of the machine to move the diamond tool. Fast-tool servo requires the diamond tool to be mounted on a piezo or voice-coil controller (the controller is placed on the *Z* stage) and move it with high frequency and small *Z* amplitude synchronously with angular position of spindle (*C* axis). The fast-tool servo, depending on the particular design, works at 100 to 700 Hz frequency and within the *Z* range of approximately 50–70 μm to 500 μm (Precitech, Moore Nanotechnology Systems). Large scanning range of $\pm 3000 \mu\text{m}$ is possible at low 20-Hz rate (Moore Nanotechnology Systems). The fast-tool servo allows higher fabrication speeds compared with slow-slide servo mode, which in some cases slows the spindle down (during process) below 1 rpm. Depending on the particular machine, spindle-rotation rates can reach upto 10,000 rpm. Common diamond-turning rates are 1000 to 3000 rpm. An example of a lens array fabricated with slow-slide servo



FIGURE 15 Zeonex aspherical lenses made for a $\text{NA} = 1.0$ miniature microscope cut in this two-axis mode.

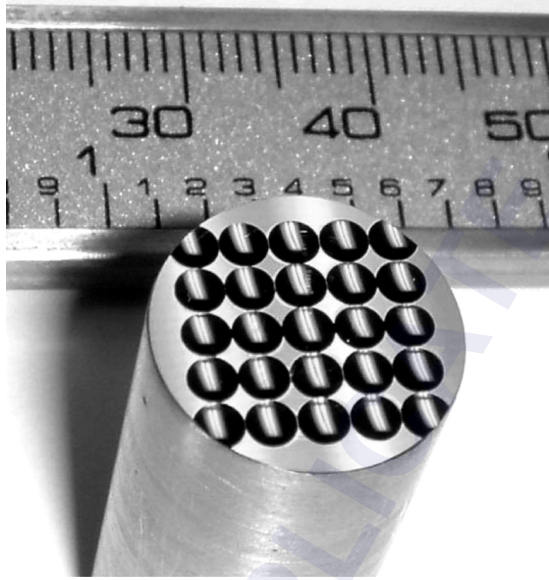


FIGURE 16 An example of a lens array fabricated with slow-slide servo *C*-axis machining.³⁹

C-axis machining is presented in Fig. 16.³⁹ There is a limitation on the slope angle for diamond turning of array components in *C*-axis mode. It depends on the primary clearance angle of the diamond tool. For a standard diamond-cutting tool the slope angle limit is approximately 20° . Some manufactures can supply tools with angles up to 32° , but these tools are very fragile. Note that the slope angle limitation can be mitigated with *B* axis or in micromilling mode.

Micromilling technology requires the diamond tool to be mounted on the spindle and the part now on the *Z*-axis stage. While the spindle (tool) rotates, it can be moved in *X*, *Y*, *Z* directions in reference to the fabricated part. The cutting time of free-form surfaces with diamond micromilling is significantly longer than for diamond turning in *C*-axis mode. For that reason, it is more appropriate for small optic applications. On the other hand, micromilling is very effective for fabrication of parts that conform to the shape of the tool. For example, making arrays of spherical lenses can be very fast ($1\text{-}\mu\text{m}$ cutting depth/revolution) and deliver high surface quality. Micromilling spindle rates are usually faster than for normal diamond-turning applications and typically are 6000 to 8000 rpm. An example of a lens array cut with micromilling is shown in Fig. 17. Micromilling can produce high-power components, like hemispherical lenses or arrays of such lenses.

Research on fabrication of aspherically shaped tools is currently being pursued.⁴⁰ If successful, it will provide tremendous benefits for micromilling fabrication and will allow fast production of arrays of high-power aspherical lenses.

22.7 LITHOGRAPHY FOR MAKING REFRACTIVE COMPONENTS

In recent years lithographic techniques became important for micro-optics fabrication. Since lithography is designed for making small, planar components it is more appropriate for fabrication of miniature and micro components rather than for large optics. While entire micro-optics systems can be made with lithographic processes, fabrication of large-scale devices is limited to phase

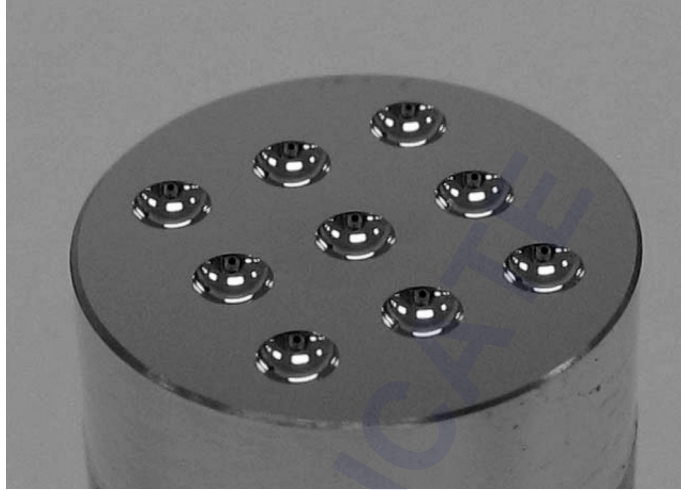


FIGURE 17 A lens array cut with diamond micromilling. (Courtesy of Moore Nanotechnology Systems, LLC.)

correctors/plates. As lithography was derived from electronics, it allows producing components in large quantities (parts can be diced out of the wafer) or arrays of systems. Additional benefits are obtained with use of grayscale lithography, which allows making components of arbitrary shape. In this process, fabrication of aspherical surfaces is as easy as making spherical ones (no rotational symmetry is necessary) and it gives the designer enormous work freedom. One example application of grayscale lithography for making arbitrary surfaces was proposed by J. Rogers in the design of miniature microscope.⁴¹ The goal of the design was to remove ghost reflections through lens tilts and to at the same time maintain high performance of $NA = 0.4$ optics. Shapes of lens surfaces made with grayscale lithography corrected for aberrations originating from tilts. Another benefit of using grayscale lithography is making optomechanical (assembly) features in the same process as optics. This combination can significantly improve assembly precision (see also Sec. 22.4).

While lithography has many advantages, it also has lens sag limitations (in effect limits lens power) currently between 100 and 150 μm .⁴² To improve lens parameters research on higher refractive index materials (for example, solgel) is currently being pursued.

Grayscale lithography can be divided into two major categories: (1) direct printing and (2) lithography using grayscale masks. Direct printing can be performed using laser, x-ray or e-beam writing in photosensitive material. In this chapter we will concentrate on fabrication using laser direct writing.

Laser Direct Write Fabrication of Micro-optical Elements

The maskless optical lithography technique referred to in this section is laser direct write (LDW) technology. Most LDW phototools have been developed by modulating a single laser spot focused by a microscope objective. The exposure media is moved under the single focused spot using a translation stage, as shown in Fig. 18.⁴³ Both rectilinear raster scan and rotational scan techniques have been studied. The rotational scan spins the media on a rotary table, and the objective translates across a radius, as shown in Fig. 19.⁴⁴ Rectilinear raster scan by means of a translating X-Y stage allows patterning of nonrotational profiles with the advantage of patterning more than one element.⁴³ This method adds a nonrotational pattern structure to the final image.

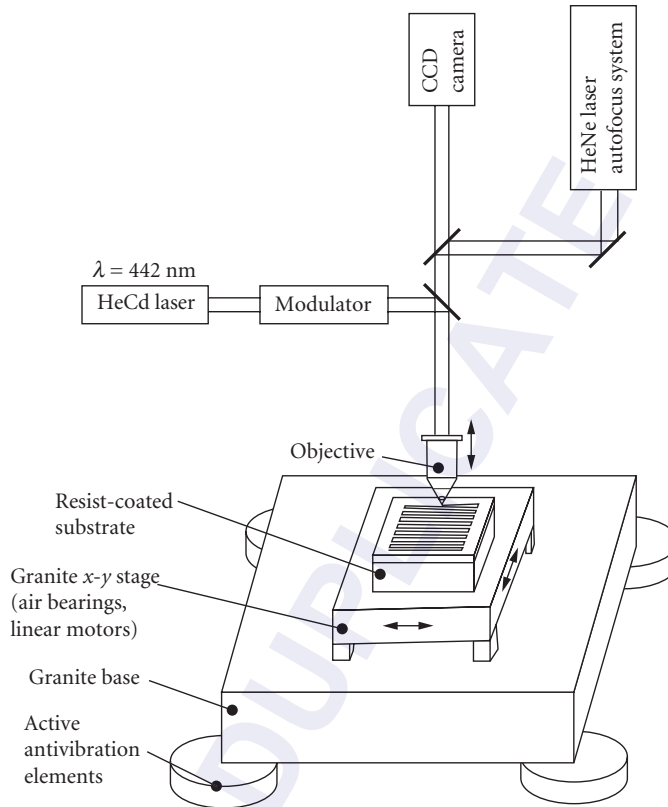


FIGURE 18 Typical laser direct write (LDW) system employing a stable base and x-y scan stage with a modulated laser to expose the photosensitive surface at each scan point.⁴³

The process steps used in a typical application using LDW technology are

1. Calculate the desired two-dimensional surface sag matrix in the lens material, based on the optical wavefront desired in transmission or reflection. In transmission, properties of the optical material used for the micro-optical element must be considered at the design wavelength. Exposure into simple photoresist is often sufficient for many applications, but the pattern may also be milled or etched into an underlying substrate, like fused silica or GaP.
2. If all other process steps are linear, the surface sag matrix is directly proportional to the depth of material removed with a positive-tone photoresist, which is the most common type of photoresist.
3. However, these process steps are rarely linear. Typically, the designer must compensate the surface sag by the nonlinear properties of the photoresist and the etching or milling processes. Nonlinear properties of the photoresist are best calibrated directly on the LDW phototool immediately before exposure. A calibration ramp is shown in Fig. 20, where responses of the exposure system and photoresist to a linear input ramp show a strongly nonlinear behavior. Depending on the milling or etching technique used, the additional correction for differential milling rate (etch rate of the optical material divided by the etch rate of the resist) must be accommodated. In dry-etch system, the differential etch rate is typically 6:1, but may be as high as 100:1 or as low as 1:1.

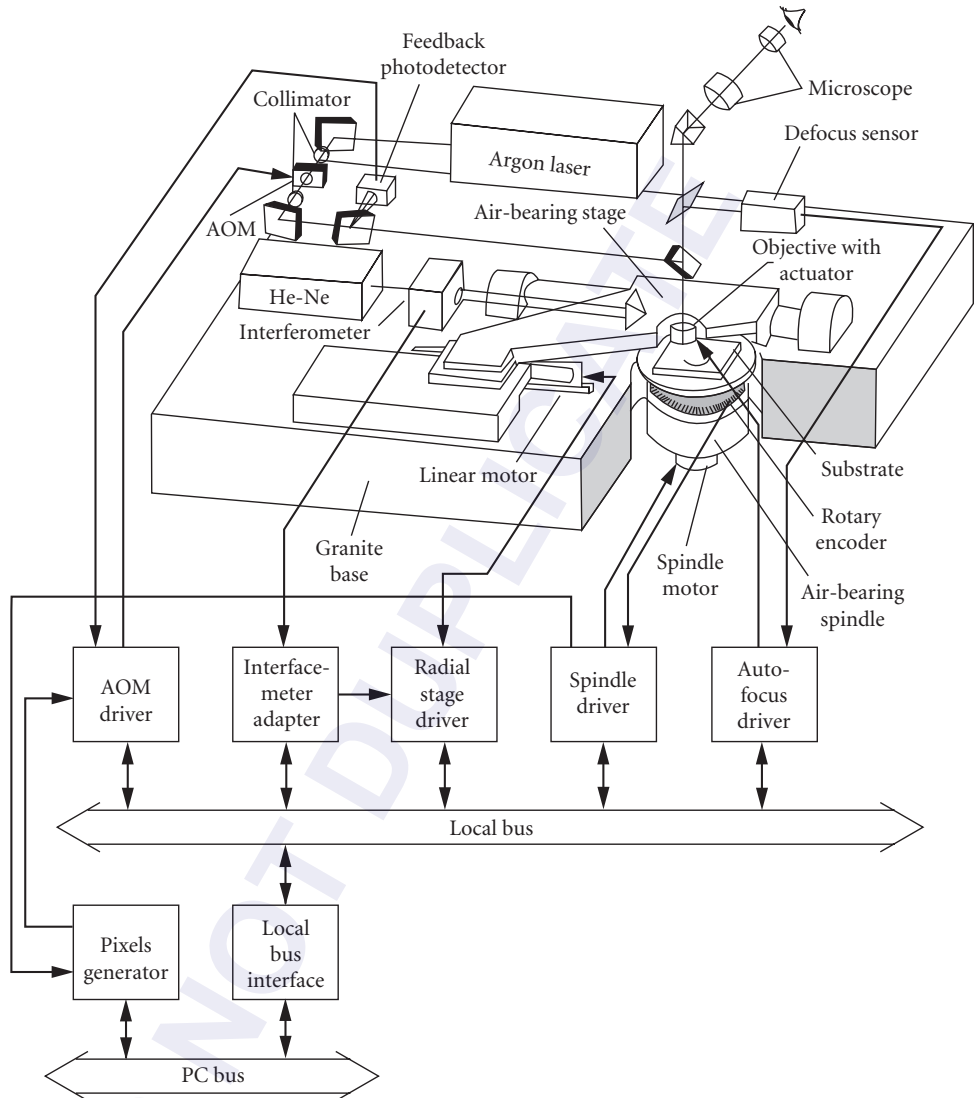


FIGURE 19 Polar coordinate writer for fabrication of micro-optical components.⁴⁴ The substrate rotating on the spindle is exposed with a focused and modulated laser beam along a diameter. A translation stage moves the focused beam along the spindle radius to expose different diameters. Good synchronization between the spindle and laser modulation can produce arbitrary gray-scale profiles.

4. Resist layer is deposited on the substrate, exposed and developed according to the compensated exposure profile. Certain photoresist process steps are required, such as prebake to remove residual solvents and postbake or postdevelop bake, which depend on the chemistry being used. If resist is used as the optical material, there are no more process steps required.
5. If the pattern is transferred into the substrate material, the substrate must be mounted into the milling chamber or etch bath for further processing.

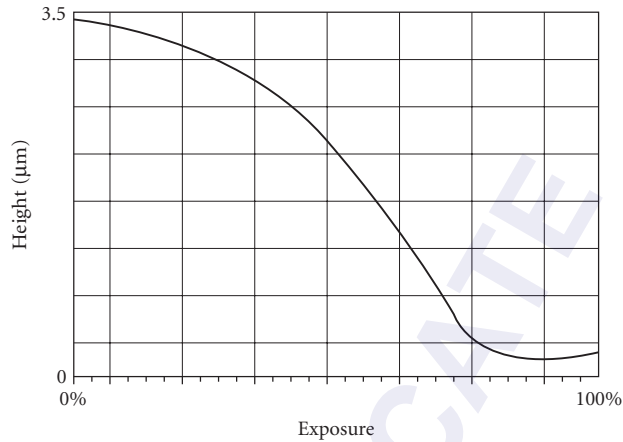


FIGURE 20 Response curve of positive-tone photoresist to a linearly increasing exposure signal (0 to 100 percent laser power) that takes into account nonlinearities of the exposure system and the photoresist.

Although exposure in photoresist is the most common method, other materials may be used for direct laser writing. For example, Poleshchuk et al. use a polar coordinate writer to change the absorption of amorphous silicon as a function of position for fabrication of gray-scale masks, as shown in Fig. 21.⁴⁴ The reduction in optical absorption is caused by an induced crystallization of the amorphous silicon.

Other types of LDW phototools have also been proposed. A system using x - y galvanometer mirrors has been constructed to fabricate micro-optical components by photodeposition of amorphous

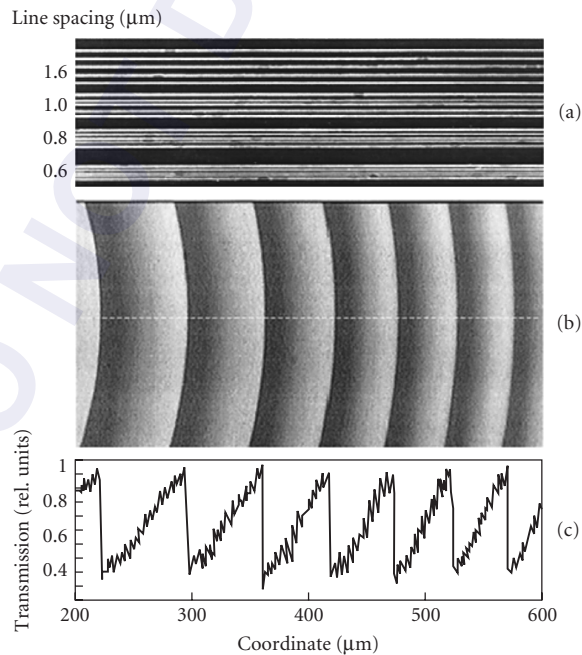


FIGURE 21 Directing laser writing of gray-scale patterns in amorphous silicon films using a polar coordinate writer.⁴⁴

selenium (a-Se).⁴⁵ The galvanometer mirrors are used instead of x - y stages to move the focused laser spot across the sample. Yet another system uses digital multi-micromirror device (DMD) array with a 10:1 lithographic reduction camera to expose micro-optical patterns.⁴⁶ In this system, the DMD is illuminated with a He g-line (436 nm) illumination system, and pixels of the DMD are turned “on” or “off” to direct light into or away from, respectively, the projection camera. A full image of the DMD array is exposed at one time with this system, thus increasing throughput. Pixels in the “on” state expose the resist. Still another system under development is the high-speed maskless lithography tool (MLT) at the University of Arizona.⁴⁷ In this system, an 8-bit modulated 370-nm laser beam is reflected from a rotating polygon and focused to a scan line (x dimension) in the photoresist. As the polygon rotates, the photoresist-coated sample is slowly translated with a linear stage perpendicular (y dimension) to the scan line. By the time the scan line is complete, the stage has moved a y distance corresponding to the width of 1 pixel in that dimension and is ready for the next scan line. A twelve-sided polygon rotating at 3000 rpm provides 12,000 pixels across one scan line. Pixel dimensions are 2.1 μm square, and the spot size is approximately 2.5 μm . The length of the written pattern in the y direction depends on the range of the stage and the size of the computer buffer memory. One 25 mm \times 25 mm substrate can be exposed in approximately 12 seconds.

Lithography with Gray-Scale Masks

Lithography with use of gray-scale masks is a three step process:

1. Substrate is coated with photoresist, polymer, or solgel glass.
2. Sample is exposed with UV light through gray-scale mask.
3. Sample is developed to create refractive lens component.

For transferring lens shape into glass, an additional step of deep reactive ion etching (DRIE) is required. A conceptual graph presenting the gray-scale lithography process is shown in Fig. 22. The process requires two critical components, which are a gray-scale mask and a photosensitive material with linear relationship between optical density and level of polymerization so the lens height will directly relate to energy delivered to the sample. Note that a nonlinear relation is also possible to use,

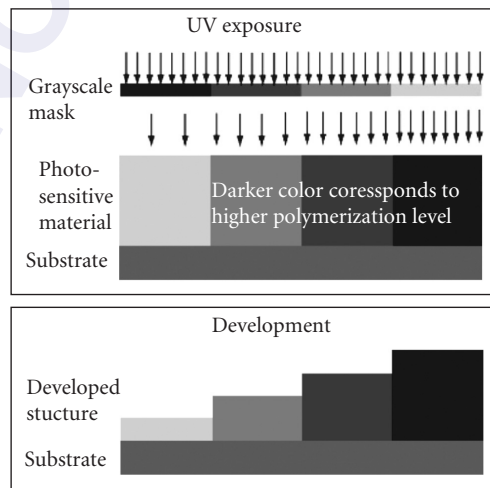


FIGURE 22 Principle of gray-scale lithography process.

TABLE 6 Comparison of Achievable Resolution for Different Type of Gray-Scale Mask³

Mask Type	Minimum Feature Size
Photographic masks	<5 μm
Halftoning	\cong 8 μm
HEBS masks	<1 μm

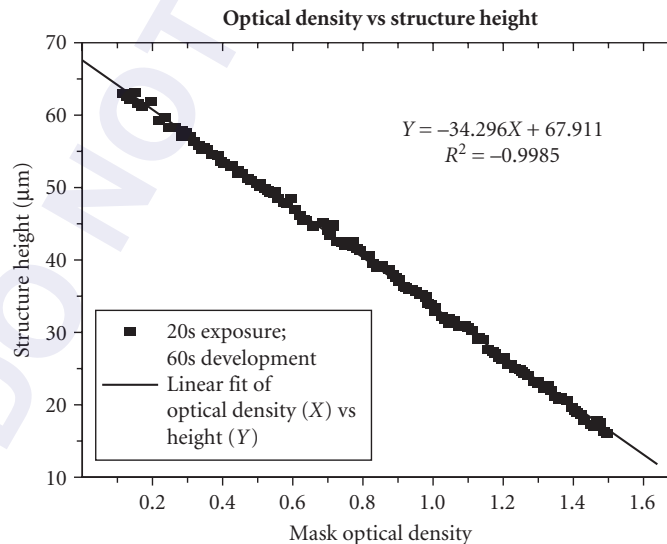
but it requires a careful calibration process for the mask fabrication, exposure, and development process. Examples of technologies to make gray-scale masks include³

- Photographic masks—can be made with photoreduction on photosensitive emulsion or with variable laser intensity writing.
- Halftoning—gray levels are obtained by different density and size of binary structure.
- HEBS masks—high energy beam sensitive glass is used to fabricate masks with e-beam writing.

Comparison of parameters for these mask techniques is presented in Table 6.

During the lithography process negative-tone photosensitive material (solgel, photoresist) is polymerized with UV light and immersed in developer. Regions less exposed to UV radiations dissolve more quickly and create variable height components. An example of structure height as a function of optical density is shown in Fig. 23.⁴⁸

Especially interesting is the gray-scale lithography process with thick solgel and photoresist materials, as it allows higher optical power and larger components. It is especially important for miniature imaging systems (miniature microscopes, endoscopes, cameras). These devices need to obtain large FOV and therefore their size cannot be too small. High-film thicknesses is a prerequisite to optical elements of greater optical power, that is, shorter focal length. Small root-mean-square (rms) surface roughness is demanded to minimize undesired scattering from lithographically

**FIGURE 23** Optical component surface height as a function of optical density for hybrid solgel glass.⁴⁸

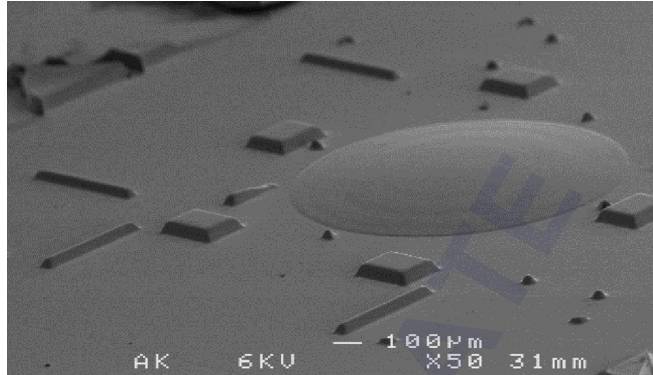


FIGURE 24 SEM image of the lens fabricated using solgel and gray-scale HEBS photomasks.⁴⁸

fabricated optical elements. Values of large thickness (up to 135 μm) at a low *rms* roughness of 10 nm were reported.⁴² An example of a solgel lens fabricated with gray-scale process is shown in Fig. 24.⁴⁸ Features around the lens were designed for optomechanical assembly and work in the MOT setting. (For more details refer to Sec. 22.4.)

22.8 MONOLITHIC LENSLET MODULES

Monolithic lenslet modules (MLMs) are micro-optic lenslets configured into close-packed arrays. Lenslets can be circular, square, rectangular, or hexagonal. Aperture sizes range from as small as 25 μm to 1.0 mm. Overall array sizes can be fabricated up to 68 \times 68 mm. These elements, like those described in the previous section, are fabricated from molds. Unlike molded glass and plastic lenses, MLMs are typically fabricated on only one surface of a substrate, as shown in the wavefront sensing arrangement of Fig. 25. An advantage of MLMs over other microlens array techniques is that the fill factor, which is the fraction of usable area in the array, can be as high as 95 to 99 percent.

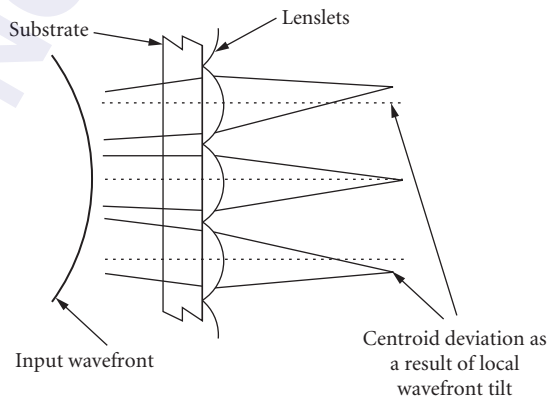


FIGURE 25 Monolithic lenslet modules (MLMs) configured for wavefront sensing.⁵⁷

Applications for MLMs include Hartman testing,⁴⁹ spatial light modulators, optical computing, video projection systems, detector fill-factor improvement,⁵⁰ and image processing.

There are three processes that are used to construct MLMs.⁵¹ All three techniques depend on using a master made of high-purity annealed and polished material. After the master is formed, a small amount of release agent is applied to the surface. In the most common fabrication process, a small amount of epoxy is placed on the surface of the master. A thin glass substrate is placed on top. The lenslet material is a single-part polymer epoxy. A slow-curing epoxy can be used if alignment is necessary during the curing process.⁵² The second process is injection molding of plastics for high-volume applications. The third process for fabrication of MLMs is to grow infrared materials, like zinc selenide, on the master by chemical vapor deposition. Also, transparent elastomers can be used to produce flexible arrays.

MLMs are advertised⁵³ to be diffraction-limited for lenslets with $NA < 0.10$. Since the lens material is only a very thin layer on top of the glass substrate, MLMs do not have the same concerns that molded plastic lenses have with respect to birefringence and transmission of the substrate. For most low-NA applications, individual lenslets can be analyzed as plano-convex lenses. Aspheres can be fabricated to improve imaging performance for higher NAs. Aspheres as fast as $NA = 0.5$ have been fabricated with spot sizes about twice what would be expected from a diffraction-limited system. The residual error is probably due to fabrication imperfections observed near the edges and corners of the lenslets.⁵⁴

22.9 DISTRIBUTED-INDEX PLANAR MICROLENSSES

A distributed-index planar microlens, which is also called a Luneberg lens,⁵⁵ is formed with a radially symmetric index distribution. The index begins at a high value at the center and decreases to the index value of the substrate at the edge of the lens. The function that describes axial and radial variation of the index is given by⁵⁶

$$n(r, z) \approx n(0, 0) \sqrt{1 - g^2 r^2 - \frac{2g \Delta n^2(0, 0)}{d} z^2} \quad (4)$$

where r is the radial distance from the optical axis, z is the axial distance, $n(0, 0)$ is the maximum index at the surface of the lens, g is a constant that expresses the index gradient, d is the diffusion depth, and $\Delta = (n(0, 0) - n_2)/n(0, 0)$, where n_2 is the substrate index. Typical values are $\Delta = 0.05$, $d = 0.4$ mm, $r_{\max} = 0.5$ mm, $n_2 = 1.5$, and $g = \sqrt{2\Delta}/r_{\max} = 0.63$ mm⁻¹. These lenses are typically fabricated on flat substrates and yield hemispherical index profiles, as shown in Fig. 26. Two substrates placed together will produce a spherical lens. Several applications of light coupling with distributed-index microlenses have recently been demonstrated.⁵⁷ These include coupling laser diodes to fibers, LEDs to fibers, fibers to fibers, and fibers to detectors. In the future, arrays of lenslets might aid in parallel communication systems.

One way to introduce the index gradient is through ion exchange.⁵⁸ As shown in Fig. 26, a glass substrate is first coated with a metallic film. The film is then patterned with a mask that allows ions to diffuse from a molten salt bath through open areas of the mask. Ions in the glass substrate are exchanged for other ions in the molten salt at high temperatures. The diffused ions change the refractive index of the substrate by an amount that is proportional to their electric polarizability and concentration. To increase the index, diffusing ions from the salt bath must have a larger electronic polarizability than that of the ions involved in the glass substrate. Since ions that have larger electron polarizability also have larger ionic radius, the selective ion exchange changes the index distribution and creates local swelling where the diffusing ion concentration is high. The swelling can be removed with polishing for a smooth surface. Alternatively, the swelling can be left to aid in the lensing action of the device. To obtain the proper index distribution, the mask radius and diffusion time must be chosen carefully.⁵⁹ If the mask radius, r_{mask} , is small compared to the diffusion depth, the derivative of the index distribution with respect to radial distance r monotonically

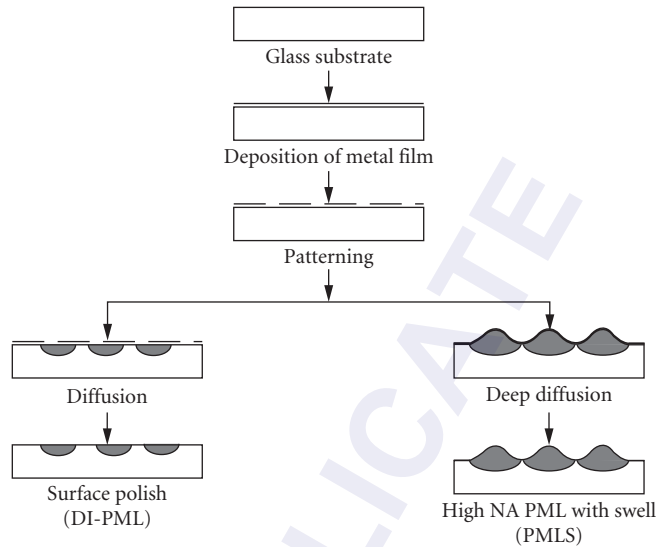


FIGURE 26 Planar distributed-index microlens array and fabrication process.⁵⁷

decreases. Since the curvature of a light ray passing through the medium is proportional to the gradient of the logarithm of the refractive index, the rays tend not to focus. A suitable combination of diffusion time t and mask radius is given by $Dt/r_{\text{mask}}^2 \approx 0.4$, where D is the diffusion constant of the dopant in the substrate. Table 7 displays the diffusion time necessary for making a planar microlens with a radius of 0.5 mm. Typically, the paraxial focal length in the substrate is $l_0 \approx 20 r_{\text{mask}}$, and the numerical aperture is $\text{NA} \approx n_2/20$.

Other fabrication techniques can also be used. Planar lens arrays in plastics are fabricated with monomer-exchange diffusion.^{60,61} Plastics are suitable for making larger-diameter lenses because they have large diffusion constants at relatively low temperatures (100°C). The electromigration technique⁶² is more effective for creating devices with short focal length. For example, by applying an electric field of 7 V/mm for 8 hours, it is possible to obtain a planar microlens with radius of 0.6 mm and focal length of 6.8 mm.⁵⁹ A distributed-index microlens array using a plasma chemical vapor deposition (CVD) method has also been reported.⁶³ In this process, hemispherical holes are etched into a planar glass substrate. The holes are filled with thin layers of a combination of SiO_2 and Si_3N_4 . These materials have different indices of refraction, and the composition is varied from the hemispherical outside shell to the center to provide a Luneburg index distribution.

TABLE 7 Summary of Diffusion Times for Planar DI Lenses⁶¹

Materials	W_n/n	D (m ² /s)	t (s)*
Plastics (DIA-MMA)	0.05	3×10^{-10}	3×10^2
Glass (TI) ion-exchange	0.05	4×10^{-13}	9×10^4
Glass (TI) electromigration	0.05	—	$3 \times 10^{4†}$

* $t = (r_m^2/D) \times 0.4$.

†Experimental data with radius of 0.6 mm.

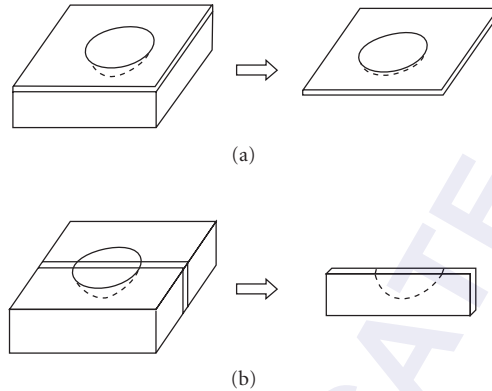


FIGURE 27 Slicing a lens to obtain a thin sample for interferometric characterization.⁵⁹ (a) Lateral slice and (b) longitudinal slice.

Shearing interferometry can be used to measure the index distribution from thinly sliced samples of lenslets. Samplers are acquired laterally or longitudinally, as shown in Fig. 27. Results of the measurement on a lateral section are shown in Fig. 28 for the ion-exchange technique. The solid line is the theoretical prediction, and the dotted line corresponds to measured data. The large discrepancy between measured and theoretical results is probably due to concentration-dependent diffusion or the interaction of the dopants. Figure 29 shows the two-dimensional index profile resulting from a deep electromigration technique.⁶⁴ These data correspond much more closely to the theoretical values in Fig. 28.

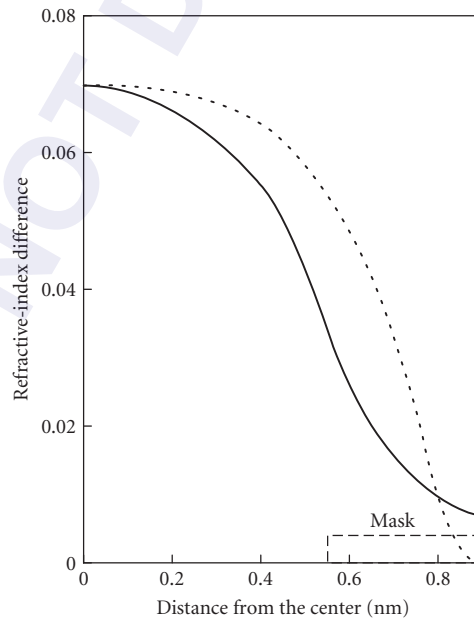


FIGURE 28 Surface index distribution of a planar microlens. Theoretical (—) and experimental (···).¹⁰⁷

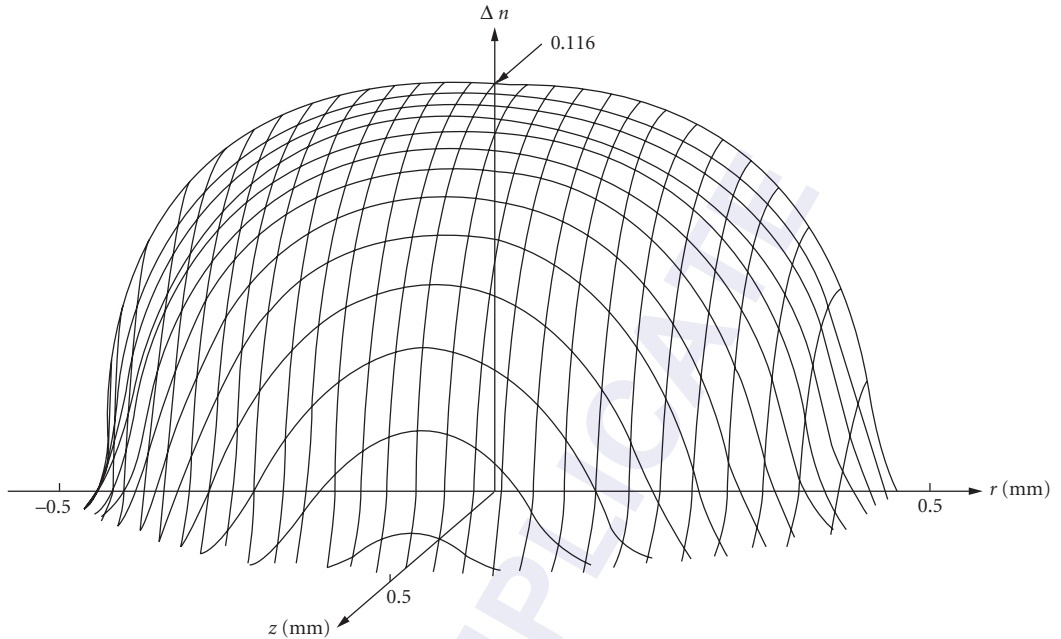


FIGURE 29 Two-dimensional index distribution of a distributed-index planar microlens prepared with the deep electromigration technique.⁶⁴

The ray aberration of a distributed-index lens is commonly determined by observing the longitudinal aberration at infinite conjugates, as shown in Fig. 30. The paraxial focusing length, l_0 , is given by

$$l_0 = d + \frac{\sqrt{1-2\Delta}}{g} \cot \left[\frac{gd \sin^{-1} \sqrt{2\Delta}}{\sqrt{2\Delta}} \right] \quad (5)$$

The amount of longitudinal aberration is defined by $LA = (l - l_0)/l_0$, where l is the distance at which a ray crosses the optical axis. LA increases with the radius r of the ray. In order to display the effects of different Δ and n_2 parameters, we define a normalized numerical aperture that is given by

$$\overline{NA} = \frac{NA}{n_2 \sqrt{2\Delta}} \quad (6)$$

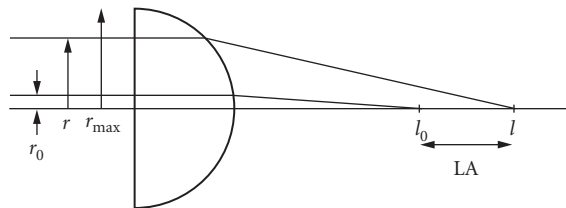


FIGURE 30 Longitudinal ray aberration, LA , of a distributed-index planar microlens. The object is at infinity. l_0 is the paraxial focal distance. LA increases with r .

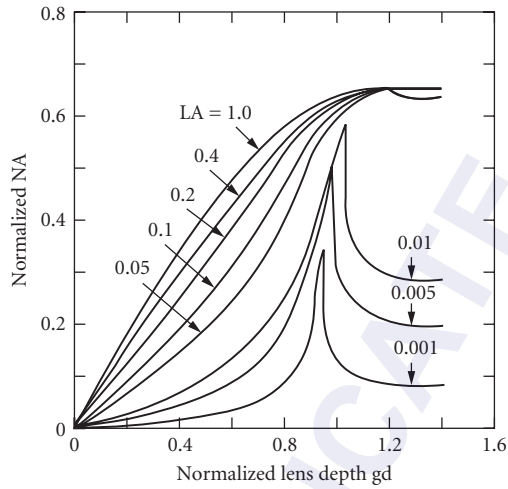


FIGURE 31 Normalized NA versus normalized depth of the distributed index region at several values of LA.⁵⁶

and is plotted in Fig. 31 versus diffusion depth for several values of LA. Notice that, for small values of LA, the maximum NA occurs at a diffusion depth of $d \approx 0.9/g$.

Wave aberration of a planar distributed index microlens is shown in Fig. 32. The large departure at the maximum radius indicates severe aberration if used at full aperture. Swelled-structure lenses can exhibit much improved performance.⁶⁵ It has been determined that the index distribution contributes

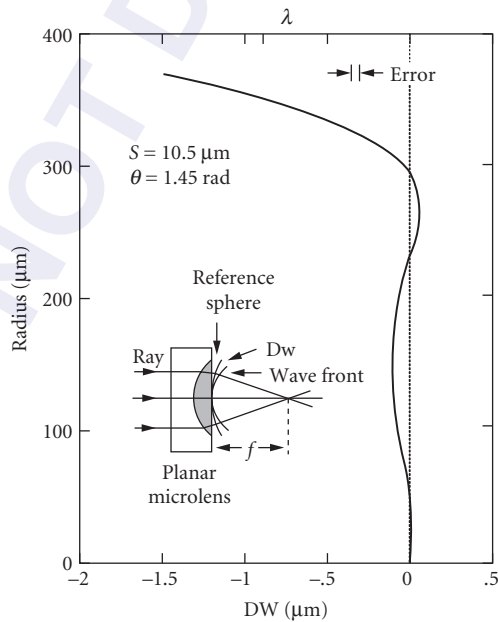


FIGURE 32 Wave aberration of a distributed-index planar microlens.¹⁰⁸

TABLE 8 Fundamental Characteristics of the Planar DI Microlens²⁸

	Diameter (μm)	NA	Focal Length (μm)
Planar	10–1000	0.02–0.25	20–4000
Swelled	50–400	0.4–0.6	55–500

very little to the power of the swelled-surface element. Most of the focusing power comes from the swelled surface-air interface. A few characteristics of ion-exchanged distributed-index microlenses are shown in Table 8.

22.10 MICRO-FRESNEL LENSES

The curvature of an optical beam's wavefront determines whether the beam is converging, diverging, or collimated. A bulk lens changes the wavefront curvature in order to perform a desired function, like focusing on a detector plane. The micro-Fresnel lens (MFL) performs the same function as a bulk lens, that is, it changes the curvature of the wavefront. In a simple example, the MFL converts a plane wavefront into a converging spherical wavefront, $A(x, y, z)$, as shown in Fig. 33. The difference between an MFL and a bulk lens is that the MFL must change the wavefront over a very thin surface.

A Fresnel lens is constructed of many divided annular zones, as shown in Fig. 34. Fresnel lenses are closely related to Fresnel zone plates.^{66,67} Both zone patterns are the same. However, unlike a Fresnel zone plate, the Fresnel lens has smooth contours in each zone, which delay the phase of the optical beam by 2π radians at the thickest point. In the central zone, the contour is usually smooth enough that it acts as a refractive element. Toward the edges, zone spacing can become close to the wavelength of light, so the Fresnel lens exhibits diffractive properties. Also, due to the quasi-periodical nature of the zones and the diffractive properties, Fresnel lenses have strong wavelength dependencies.

Advantages of the Fresnel lens are that they can be made small and light compared to bulk optical components. Note that binary optics, which are described in Michael W. Farn and Wilfrid B. Veldkamp's Chap. 23, "Binary Optics," are stepped approximations to the MFL smooth-zone contour.

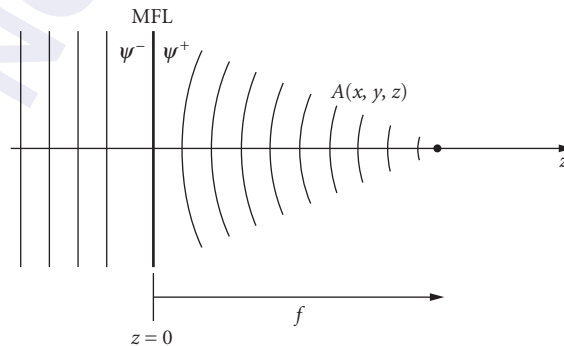


FIGURE 33 A micro-Fresnel lens (MFL) is often used to convert a planar wavefront into a converging spherical wave, $A(x, y, z)$, which focuses a distance f away from the MFL. The phase of the light in a plane on either side of the MFL is described by ψ^- and ψ^+ .

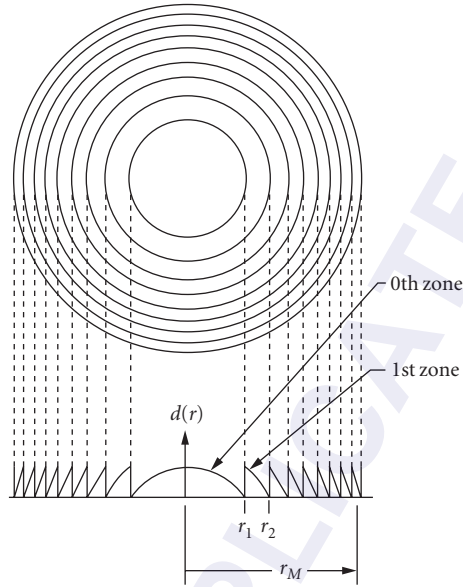


FIGURE 34 Fresnel lens construction. M -divided annular zones occur at radii r_i in the same manner as a Fresnel zone plate. The profiles of each zone are given by $d(r)$, and they are optimized to yield the maximum efficiency in the focused beam.

To understand the zonal profiles of the MFL, we return to our example problem illustrated in Fig. 33. Our development is similar to that described by Nishihara and Suhara.⁶⁸ The converging spherical wavefront is given by

$$A(r, z) = \frac{A_0}{\rho} \exp[-i(k\rho + \omega t)] = \frac{A_0}{\rho} \exp[i\phi(r, z)] \quad (7)$$

where A_0 is the amplitude of the wave, $\rho^2 = (z - f)^2 + r^2$, $r = \sqrt{x^2 + y^2}$, f is the focal length, and $k = 2\pi/\lambda$. The phase of $A(x, y, z)$ at $t = 0$ and in a plane just behind the MFL is given by

$$\phi(x, y, 0^+) = -k\sqrt{f^2 + r^2} \quad (8)$$

We could add a constant to Eq. (8) and not change any optical properties other than a dc phase shift. Let

$$\psi^+(r) = \phi(x, y, 0^+) + kf + 2\pi = 2\pi + k(f - \sqrt{f^2 + r^2}) \quad (9)$$

Zone radii are found by solving

$$k(f - \sqrt{f^2 + r_m^2}) = -2\pi m \quad (10)$$

where $m = 1, 2, 3, \dots$ is the zone number. The result is

$$r_m = \sqrt{2\lambda fm + (\lambda m)^2} \quad (11)$$

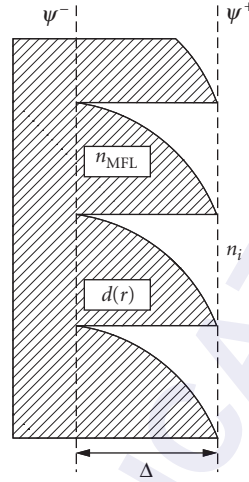


FIGURE 35 Portion of a Fresnel lens profile showing the thickness variation of the pattern. The thickness at any radius is given by $d(r)$, where r is the radial distance from the center of the lens. The phase shift that is added to wavefront ψ^- is determined by $d(r)$, the index of refraction of the substrate, n_{MFL} , and the index of refraction of the image space, n_i . The maximum thickness of the pattern is given by Δ . The resulting phase in a plane just after the MFL is given by ψ^+ .

Equation (9) becomes

$$\psi_m^+(r) = 2\pi(m+1) + k(f - \sqrt{f^2 + r^2}) \quad (12)$$

The job of the MFL is to provide a phase change so that the incident wavefront phase, $\psi^-(r)$, is changed into $\psi^+(r)$. The phase introduced by the MFL, $\psi_{\text{MFL}}(r)$, must be

$$\psi_{\text{MFL}}(r) = \psi^+(r) - \psi^-(r) \quad (13)$$

A phase change occurs when a wave is passed through a plate of varying thickness, as shown in Fig. 35. $\psi^+(r)$ is given by

$$\begin{aligned} \psi^+(r) &= \psi^-(r) + kn_{\text{MFL}} d(r) + kn_i [\Delta - d(r)] \\ &= \psi^-(r) + k(n_{\text{MFL}} - n_i) d(r) + kn_i \Delta \end{aligned} \quad (14)$$

where $d(r)$ is the thickness profile, n_i is the refractive index of the image space, n_{MFL} is the refractive index of the substrate, and Δ is the maximum thickness of the MFL pattern. $d(r)$ is found by substituting Eq. (14) into Eq. (12). Note that the factor Δ is a constant and only adds a constant phase shift to Eq. (14). Therefore, we will ignore Δ in the remainder of our development. If $n_i = 1$, the result is

$$d_m(r) = \frac{\lambda(m+1)}{n_{\text{MFL}} - 1} - \frac{\sqrt{f^2 + r^2} - f}{n_{\text{MFL}} - 1} \quad (15)$$

where we have arbitrarily set $\psi^-(r) = 0$. The total number of zones M for a lens of radius r_M is

$$M = \frac{r_M(1 - \sqrt{1 - \text{NA}^2})}{\lambda \text{NA}} \quad (16)$$

The minimum zone period, Λ_{\min} , occurs at the outermost part of the lens and is given by

$$\Lambda_{\min} = r_M - r_{M-1} = r_M \left(1 - \sqrt{1 - \frac{2\lambda f + (2M-1)\lambda^2}{2M\lambda f + (M\lambda)^2}} \right) \quad (17)$$

The following approximations may be used without significant error if $\text{NA} < 0.2$ and $M \gg 1$:

$$d_m(r) \approx \frac{m\lambda f - 0.5r^2}{f(n_{\text{MFL}} - 1)} \quad (18)$$

$$r_m \approx \sqrt{2m\lambda f} \quad (19)$$

$$M \approx \frac{r_M}{2\lambda} \text{NA} \quad (20)$$

and

$$\Lambda_{\min} \approx \frac{\lambda}{\text{NA}} \quad (21)$$

The consequence of using Eqs. (18) and (19) for $\text{NA} > 0.2$ is that a small amount of spherical aberration is introduced into the system.

The aberration characteristics of the MFL and the Fresnel zone plate are very similar. Aberrations of Fresnel zone plates have been discussed by Young.⁶⁹ For convenience, we describe a zone plate with the stop at the lens that is illuminated with a plane wave at angle α . For an MFL made according to Eq. (15) and used at the proper conjugates, there will be no spherical aberration or distortion. Coma, astigmatism, and field curvature are given by $W_{131} = \alpha r_M^3 / 2\lambda f^2$, $W_{222} = \alpha^2 r_M^2 / 2\lambda f$, and $W_{220} = \alpha^2 r_M^2 / 4\lambda f$, respectively. When $M \gg 1$ and α is small, the dominant aberration is coma, W_{131} . If the substrate of the zone plate is curved with a radius of curvature equal to the focal length, coma can be eliminated.⁷⁰ Chromatic variations in the focal length of the MFL are also similar to a Fresnel zone plate. For $\text{NA} < 0.2$,

$$\lambda f \approx \frac{r_M^2}{2M} \quad (22)$$

A focal-length-shift versus wavelength comparison of a Fresnel (hologram) lens and some single-element bulk-optic lenses are shown in Fig. 36. Note that the dispersion of the MFL is much greater than the bulk lens, and the dispersion of the MFL is opposite in sign to that of the bulk lenses. These facts have been used to design, hybrid achromats by combining bulk lenses and diffractive lenses into the same system.⁷¹ The thermal variations in MFLs primarily result in a change of focal length given by⁷²

$$\Delta f = 2f\alpha_g \Delta T \quad (23)$$

where f is the nominal focal length, α_g is the coefficient of thermal expansion for the substrate material, and ΔT is a uniform temperature change of the element. For most optical glasses, α_g ranges from $5 \times 10^{-4} \text{ } ^\circ\text{C}^{-1}$ to $10 \times 10^{-4} \text{ } ^\circ\text{C}^{-1}$.

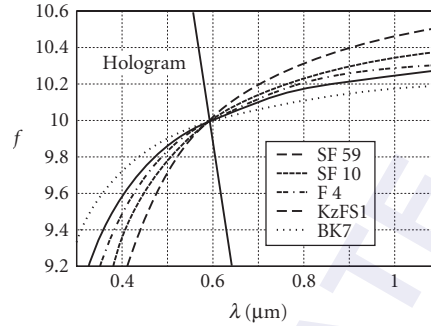


FIGURE 36 Single-element dispersions for a Fresnel (hologram) lens and refractive singlets. The focal lengths (arbitrary units) of thin lenses are plotted versus wavelength for refractive lenses of various optical glasses. Each lens was constructed to have a focal length of 10 at $\lambda_s = 0.5876 \mu\text{m}$.⁷¹

The diffraction efficiency, η , of an MFL is defined as the ratio of the power in the focused spot to the power in the unfocused beam transmitted through the lens. At best, Fresnel zone plates exhibit $\eta = 40.5$ percent.⁷³ Blazing the grating profile can significantly increase the efficiency of the lens. Theoretically, η of an MFL can be 100 percent with the proper profile. However, there are several process parameters that limit η , as shown in Fig. 37, where a perfect zone profile has width T and height $d_{\text{MAX}} = \lambda/(n_{\text{MFL}} - 1)$. Variation of film thickness, overetching, swell of the resist, and swell all exhibit sinc-squared

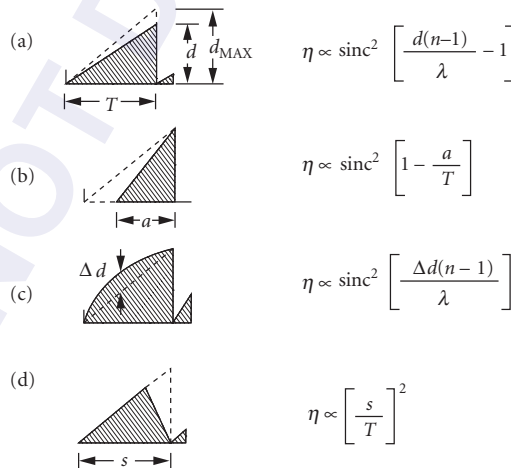


FIGURE 37 Four parameters that influence the diffraction efficiency of MFLs are: (a) film thickness variation; (b) overetching; (c) swell of the resist; and (d) imperfection of the shoulders. A profile of one zone is illustrated for each parameter. The ideal profile is shown as a dotted line, where d_{MAX} is the ideal height and T is the ideal period. The diffraction efficiency η of each profile is determined from extrapolating the result obtained from an infinite blazed grating.⁷⁵

dependency on the errors. Shoulder imperfection is the most critical parameter, with η proportional to $(s/T)^2$. For $\eta > 90$ percent, $s/T \geq 0.95$, which implies that the falling edge of the zone profile must take no more than 5 percent of the grating period. This is possible with low NA systems, where the zone spacing is large compared to the resolution of the exposure system, but it becomes difficult in high NA systems, where the zone spacing is on the order of several microns. Analysis of the three remaining parameters indicates fairly loose tolerances are acceptable. For $\eta > 98$ percent, tolerance on individual parameters are: $|d(n-1)/\lambda - 1| < 0.25$, $a/T > 0.50$, and $\Delta d(n-1)/\lambda < 0.50$. Due to the increasing difficulty in fabricating correct zone profiles with decreasing zone width, most MFLs exhibit a variation in diffraction efficiency versus radius. In the center of the zone pattern, where the zone spacing is large, the measured diffraction efficiency can be in excess of 90 percent. At the edge of the zone pattern, where the zone spacing can be on the order of a few wavelengths, the measured diffraction efficiency is much lower. One possible solution to this problem is to use “superzone” construction,⁷⁴ in which the zone radii are found from a modified form of Eq. (10), that is

$$k(f - \sqrt{f^2 + r_M^2}) = 2\pi Nm \quad (24)$$

where N is the superzone number. This results in a maximum thickness of $d_{\text{MAX}} = N\lambda/(n_{\text{MFL}} - 1)$. Note that $N = 1$ corresponds to the standard MFL. $N = 2$ implies that zones are spaced at every 4π phase transition boundary instead of at every 2π phase transition boundary. Although this makes the zones wider apart, the surface relief pattern must be twice as thick.

Molding provides a potentially valuable process for fabricating large quantities of MFLs economically. MFLs can be produced with conventional injection molding, but due to the large thermal expansion coefficient of polymers, the lenses are sensitive to thermal variations. An alternative MFL molding process is shown in Fig. 38, where a glass substrate is used to avoid large thermal effects.

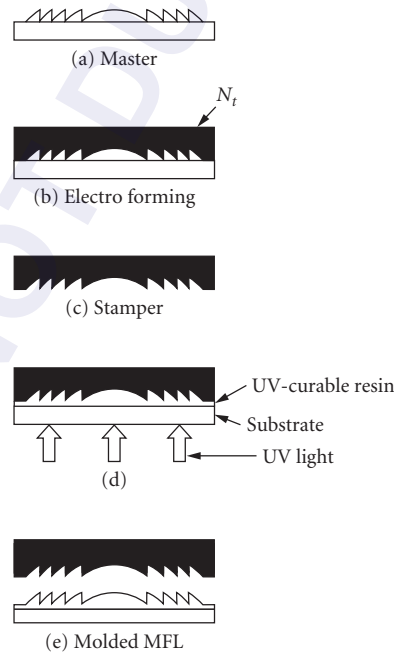


FIGURE 38 Molding process for a MFL on a glass substrate. First, a master is made by electron-beam lithography, then a stamper is electroformed from the master. MFLs are molded by potting a UV-curable resin between the stamper and the substrate and then exposing through the substrate.⁵³

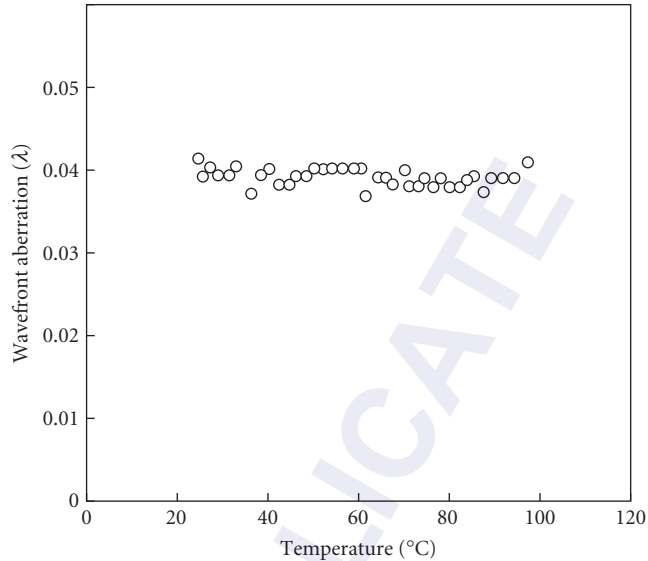


FIGURE 39 Wavefront aberration versus substrate temperature for a 0.25 NA molded MFL on a glass substrate designed to operate at $\lambda = 780 \text{ nm}$.⁷⁸

First, a master lens is formed with electron-beam or laser writing. A stamper is prepared using conventional nickel electro-forming methods.⁷⁶ After potting a UV-curable resin between the stamper and the glass substrate, the replica lenses are molded by the photopolymerization (2P) process.⁷⁷ The wavefront aberration versus temperature for a $\lambda = 780 \text{ nm}$, $\text{NA} = 0.25$, diameter = 0.5 mm lens formed with this technique is shown in Fig. 39.⁷⁸ A variation on this technique is to use the stamper as a substrate in an electron-beam evaporation device.⁷⁹ Inorganic materials of various refractive indices can be deposited on the stamper, resulting in a thin lens of high refractive index. The high refractive index of a material like ZnS ($n = 2.35$) can be used to lower the maximum thickness requirement of the lens, which makes fabrication of the master with electron-beam writing easier.

22.11 LIQUID LENSES

Integrated micro and miniature systems are often limited by fixed geometry, difficult realignment, and tuning capabilities. Therefore, extensive research was performed to develop miniature optical components with adjustable optical power. Even though zoom solutions, like applications of Alvarez-Humphrey plates mounted on microelectro mechanical system (MEMS) actuators, were recently suggested,⁸⁰ the most successful tunable miniature optic components are based on principles of liquid lenses or liquid crystals. In general, the concept of liquid lenses arises from the fact that it is possible to change the shape of the liquid volume and provide an optical power change. This power change can be done in several ways, including electrowetting,^{81,82} pressure,^{83–85} or temperature change.⁸⁶ Liquid crystal lenses⁸⁷ are based on changing an electric field to create different crystal orientations and, in consequence, refractive index distribution. One of the most successful approaches for tunable miniature lenses is the concept of an electrowetting lens.⁸¹ It was developed primarily for consumer market applications, like cell phone or credit card cameras. The limited amount of space and relatively high cost does not permit using traditional motor-driven systems in these products. Another application of electrowetting lenses is Blu-Ray Disk (BD) recording systems for dual layer disks.

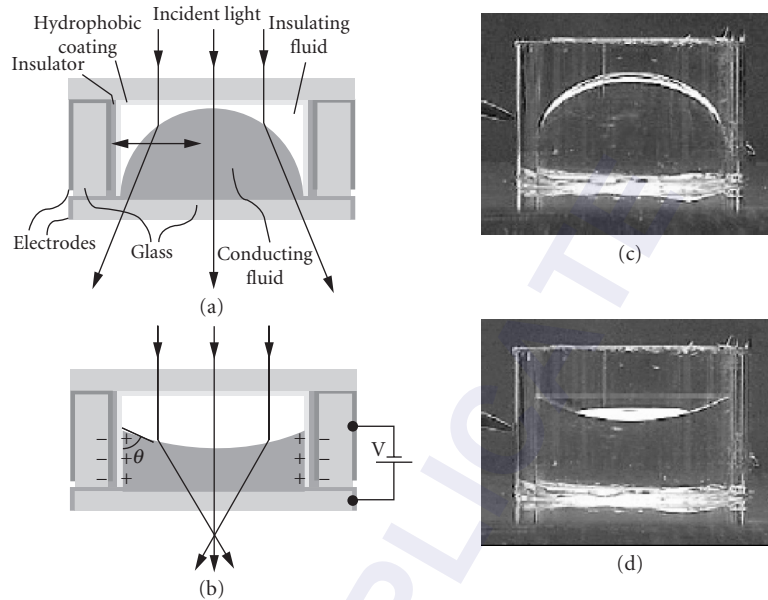


FIGURE 40 A-B schematic cross section of electrowetting lens for convex and concave lens, respectively. C-D photographs of electrowetting lenses.⁸¹

The principle of an electrowetting lens is presented in Fig. 40.⁸¹ It relies on using two non-mixable liquids, where one is conductive and one nonconductive. An example of conductive liquid is salted water, while a nonconductive example is nonpolar oil. Both liquids must have significantly different refractive index (typically $\Delta n = 0.15\text{--}0.20$) and similar density. Refractive index difference is required to provide optical power, while similar density makes the lens insensitive to vibrations and shocks. Liquid is placed in a cylinder where the cylinder wall and bottom are metal coated and act as two electrodes. Note that the cylinder wall is separated from the conductive liquid by an insulating layer. Changing potential between these two electrodes influences the shape of the conductive liquid and creates either a convex or a concave refractive lens.

An example of electrowetting camera lens from Philips is presented in Fig. 41.⁸¹ The outer diameter of the lens (Fig. 40) is 4 mm, inner diameter is 3 mm, and height 2.2 mm. Typical switching voltage is in range of 50 to 120 V. Note that a concave lens is a natural state when no voltage is applied (voltage application allows to decrease negative power or obtaining convex lenses). The switching time is about 10 ms. The Philips lens can achieve maximum negative power of -100 dioptres, and positive of $+50$ dioptres. This corresponds to focus lengths of -10 mm and 20 mm, respectively. Optical power D can be described as

$$D = \Delta n C \quad (25)$$

where C is lens curvature (reciprocal of lens radius).

Electrowetting lenses can provide sufficiently good optical performance. The main imperfections arise from asymmetric cylinder shape and nonuniform coating. Major wave aberrations in the lens are coma and astigmatism. It has been shown that they can be within diffraction limit of 0.07λ wavefront error⁸¹ within their operating range. Note that work on aspherical, asymmetric electrowetting optical components⁸² and reconfigurable lens arrays^{88,89} is currently an active research area. For the purpose of this chapter, however, we limit our discussion of electrowetting to a principle level analysis only.

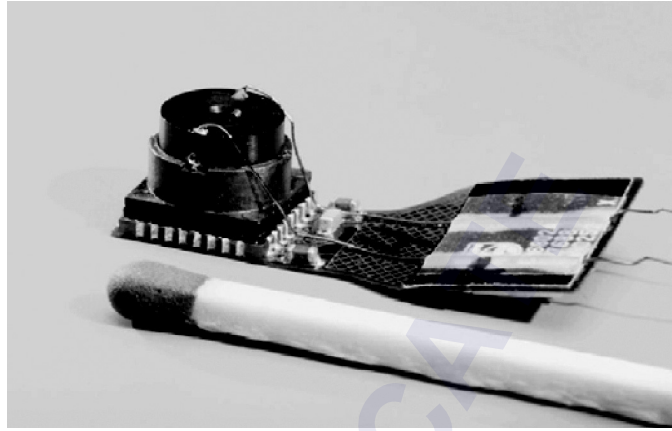


FIGURE 41 A prototype of camera module with electrowetting liquid lens.⁸¹

Another tunable lens concept is based on using liquids/gels and a pressure change to adjust optical power.^{83,84,85} In this case diameter of the cylinder containing the gel or liquid is covered with an elastic membrane (most often polydimethylsiloxane [PDMS]). The membrane is pushed out or pulled in by an increase/decrease of liquid pressure, respectively, and can create convex or concave lenses. The principle of the pressure-based liquid lens is presented in Fig. 42.

Depending on the liquid and lens diameter, pressure change lenses can obtain approximately 1.0 to 10 mm and -1.0 to -10 mm focal lengths. If the membrane is sufficiently thin, the same equation as used for electrowetted lenses can be used to estimate lens power. The pressure change technique allows easy population of lenses in micro fluidic systems and fabrication of tunable lens arrays.^{83,84,85} An array of 200- μm (diameter) lenses made on a substrate glass in PDMS is shown in Fig. 43.⁸³

Note that membrane lenses are subject to spherical aberrations, since the membrane shape produces an aspherical form due to lower membrane stiffness in the middle area of a lens that creates a nonuniform curvature. Liquid material used for an array prototype (Fig. 42) was a microscope immersion oil ($n = 1.51$) or UV-curable polymers (Norland 63, $n = 1.56$), but other materials are also possible.

While the majority of applications for pressure-based liquid lenses are “lab on a chip” use, they were also recently prototyped for use in optical coherence tomography (OCT). An example of liquid lens and complete head of an OCT probe is shown in Fig. 44.⁹⁰

Tunable lenses can also implement thermal change of the liquid volume.⁸⁵ This method is suitable for small millimeter or submillimeter components. It is due to the fact the volume change $\Delta V = \beta V \Delta T$.

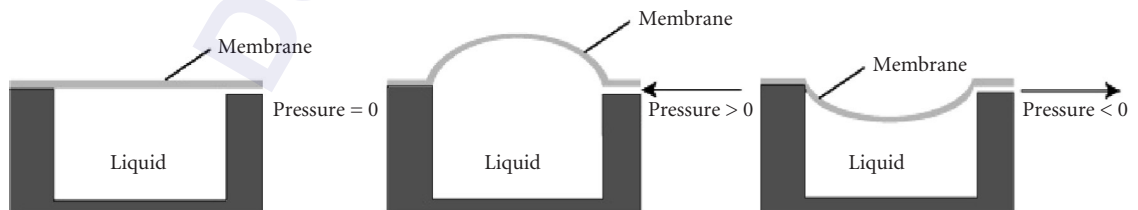


FIGURE 42 The principle of the pressure-based liquid lens.

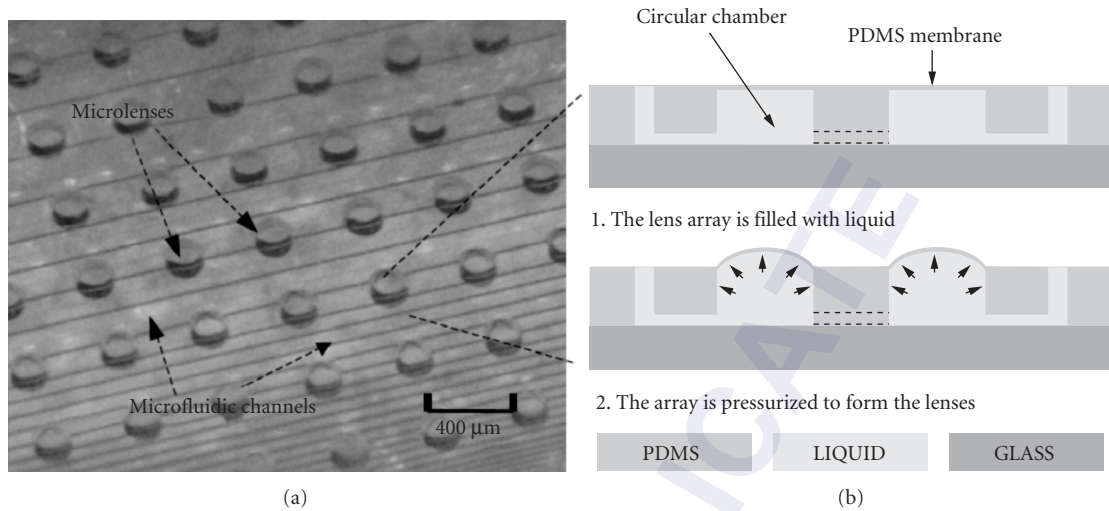


FIGURE 43 (a) Picture of the pressure base microfluidic lens array. (b) Schematic of the microlens array.⁸³

Larger masses require in linear proportion more heat energy. Therefore, this method is more effective for small liquid volumes. Note that β denotes volume expansion coefficient, V and ΔV material volume and volume change, respectively, and ΔT is the introduced temperature change. An example of a thermally controlled liquid lens is presented in Fig. 45.⁸⁶ The liquid chamber is covered with PDMS lens (gives system an initial optical power) made in a two-step replication process. PDMS was applied as a flexible and easily stretchable material. This lens has 1.9-mm clear aperture. The entire package is 8.5 mm \times 6.5 mm \times 1.5 mm. The lens requires voltage change (for thermal actuation) of 0 to 14 V, which corresponds to temperature change from about 20 to 50°C. In effect, focal length changes from 14.7 to 2.8 mm.

Entirely different technology for making tunable lenses is based on the liquid crystal principle⁸⁷ in which it is possible to create a continuously changing electric field that causes axially symmetric

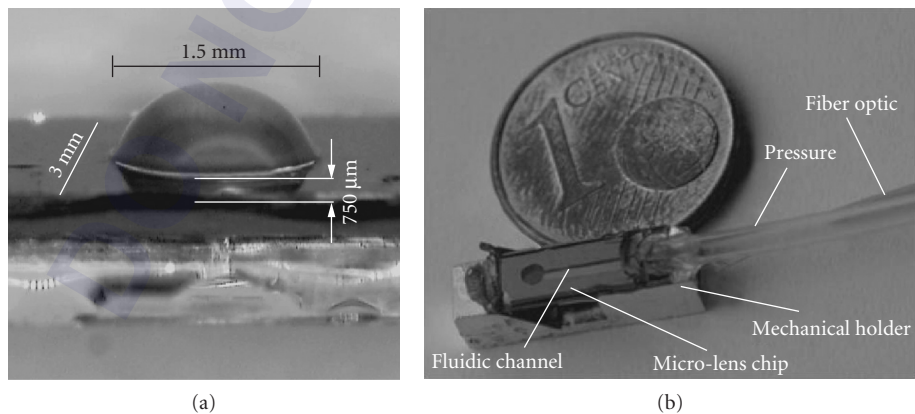


FIGURE 44 (a) Microfluidic pressure lens and (b) picture of entire OCT head.⁹⁰

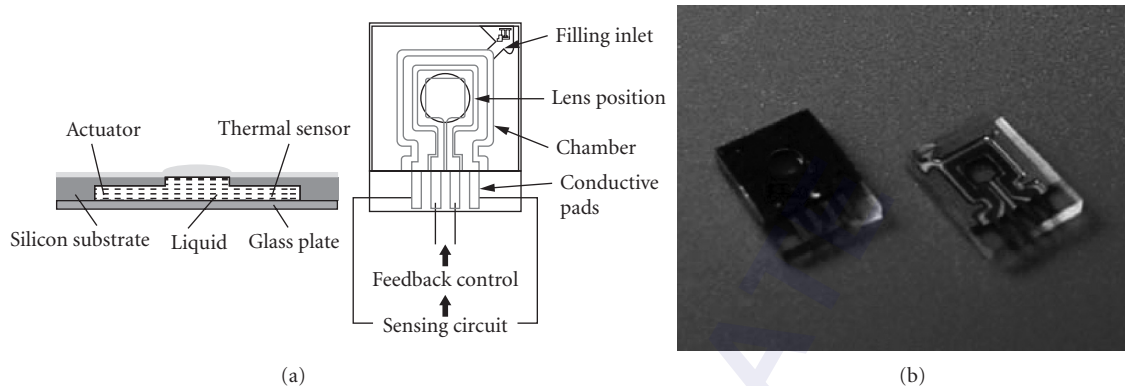


FIGURE 45 (a) Schematic of an integrated thermal volume change lens. (b) Photo of a complete packaged lens.⁸⁶

orientation changes of liquid crystals (LC). The crystal orientations results in a distribution of refractive index similar to that obtained in gradient index (GRIN) lenses. LC lenses are being used in machine vision, photonics, and eyeglasses.

The LC lens principle can be described following Naumov et al.⁸⁷ and his design of a liquid crystal lens. The LC is placed between two transparent electrodes deposited on glass substrates. One is the control electrode with distributed resistance much higher than the distributed resistance of the ground electrode. The LC layer placed between these two electrodes acts as a capacitor. In case of application of AC voltage to the single circular contact at the control electrode, the active impedance of the control electrode and the reactive impedance of the capacitor create a distributed voltage divider. This divider results in a distributed AC voltage over the layer of LC relatively close to a parabolic distribution. In consequence the axially symmetric distribution of liquid crystals provides a change of optical power. It is also possible to build cylindrical lenses by using two line contacts at the control electrode. LC lenses are very convenient, due to the fact that they are electrically controlled integrated components. Switching voltage of LC lenses, depending on design, is usually in 0–100 V range. Their major drawback arises from relatively low optical power and polarization effects due to the crystal structure.

Liquid crystal lens research is a very broad area. For the purpose of this discussion we concentrate on selected examples. One group of miniature liquid lenses includes relatively large components of 5 to 10 mm in diameter and quite low optical power with focal length between few hundred millimeters and infinity.⁸⁷ Multilayer liquid crystal lenses⁹¹ can increase optical power and obtain a range of 93 to 1230 mm. Examples of negative lenses with a glass lens embedded between liquid crystal layers were also demonstrated.⁹² Most recently, research on aspheric liquid crystal components (like axicons) is also being pursued.⁹³ A second group of liquid crystal lenses are microlenses⁹⁴ and arrays of microlenses.^{95,96} The diameter of a single lens is usually few hundred microns (200 to 600 μm) and these lenses can achieve short focal lengths in range ± 2 to ± 10 mm.

To compare various techniques for making tunable miniature or microlenses, a summary of parameters of various lens types are presented in Table 9. Values from literature are rounded to show achievable range rather than detailed numbers. Note that it is possible to find examples of systems going outside values summarized in Table 9.

Other tunable liquid lens approaches can be found in literature including adaptive liquid microlenses activated with different stimuli.⁹⁷ The optical power of these lenses may potentially be designed in a way so it will adapt to the environment, which means that optical power will change depending on conditions (PH, temperature, etc.). Detailed discussions of these quite new approaches are omitted, as they are in initial research stages.

TABLE 9 Comparison of Parameters of Different Tunable Lenses

Technique	Common Size of Clear Aperture	Focal Range	Possible Shapes	Range of Tuning Parameter	Switching Time
Electrowetted	0.2–5.0 mm	20 mm–∞ –10 mm–∞	Convex Concave	Voltage 50–120 V	10 ms
Pressure change	0.2–2 mm	1.0–10.0 mm –1.0––10 mm	Convex Concave	Pressure 0–50 KPa	500 ms– few seconds
Volume change	~2.0 mm, possible Decrease to 300 μm	2.0–15.0 mm	Convex	Temperature 20–50°C	50 seconds
Liquid crystal	5–10 mm 0.2–0.6 mm	100 mm–∞ 2–10 mm –2––10 mm	Convex Convex Concave	Voltage 0–100 V	500 ms– few seconds

22.12 OTHER TECHNOLOGIES

There are several other technologies that are potentially valuable for micro-optic components. Four particularly interesting technologies are melted-resin arrays,⁹⁸ laser-assisted chemical etching,⁹⁹ mass transport,¹⁰⁰ and drawn preform cylindrical lenses.⁹⁸

Melted-resin arrays of hemispherical ball lenses are formed with the process shown in Fig. 46. First, an Al film is deposited on a quartz substrate and patterned with holes that serve as aperture

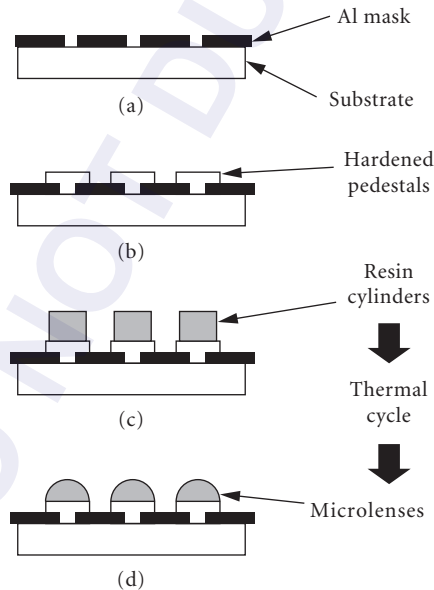


FIGURE 46 Process used to form melted-resin micro-lenses: (a) an Al film is deposited on the substrate and patterned with holes to serve as aperture stops for the array; (b) circular pedestals are formed on top of the aperture holes and hardened; (c) cylinders of resin are developed on top of the pedestals; and (d) pedestals are melted to form spherical surfaces.⁹⁸

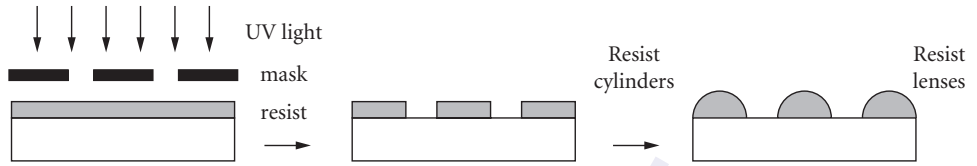


FIGURE 47 Principle of fabrication of microlens arrays using photoresist reflow technique.¹⁶

stops for the array. Next, circular pedestals are formed on top of the aperture holes. The pedestals are hardened so that they are insoluble and stable for temperatures in excess of 180°C. Cylinders of resin are then developed on top of the pedestals. The device is heated to 140°C to melt the resin. The pedestals serve to confine the melting resin. The lenses form into hemispherical shapes due to surface tension forces. Lens diameters have been demonstrated at 30 μm with good wave-front performance and uniformity.⁹⁸

Similar to the process for making melted resin arrays is fabrication of reflow lenses.^{3,16} The graphical representation of reflow process is presented in Fig. 47.¹⁶ Figure 48 shows example of fabricated lenslet array.¹⁶

Fabrication of reflow lenses consists of two major steps: (1) A layer of photoresist or solgel is exposed through lithographic binary mask and then developed; (2) photoresist cylinders are melted to create surfaces close to spherical shape. The third step of reactive ion etching^{3,16} or plasma etching¹⁰¹ can be added to transfer lens shapes into a substrate. For the melting step, the sample is heated to the glass temperature of the resist, which is in range of 150 to 200°C. After Sinzinger and Jahns,³ the focal length equation in Eq. (26) of the obtained lens is

$$f = \frac{r_c}{n-1} \quad (26)$$

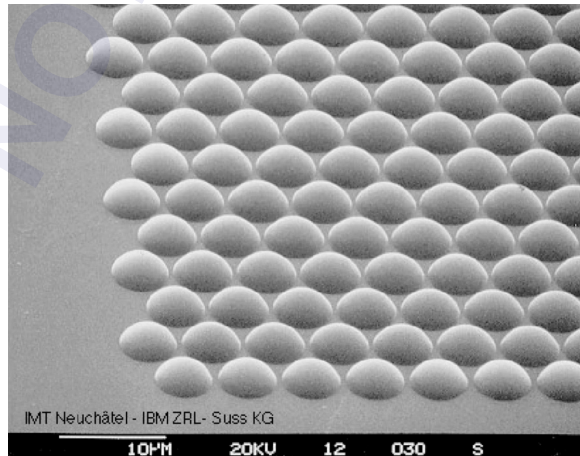


FIGURE 48 SEM picture of 5-μm diameter microlenses packed hexagonally.¹⁶ Lenses were made using reflow process in photoresist.

where r_c is radius of curvature derived from properties of binary cylinders:

$$r_c = \frac{h^2 + \frac{D^2}{4}}{2h} \quad (27)$$

Note that h is the height of the droplet and corresponds to lens sag and D is diameter of the cylinder.

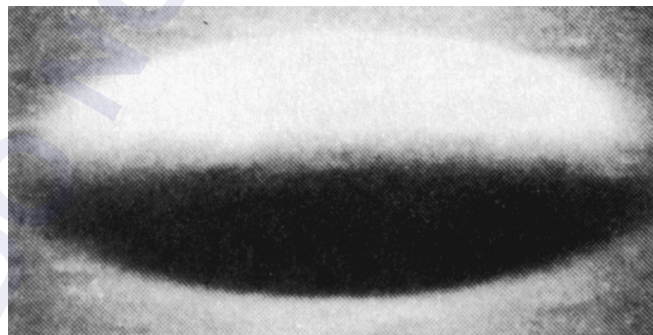
Knowledge of cylinder volume allows also to find a relation between cylinder thickness and the droplet height (sag),³

$$t = \frac{h}{6} \left(3 + 4 \frac{h^2}{D^2} \right) \quad (28)$$

Reflow lenses require very careful selection of photosensitive material (high viscosity, e.g., Hoechst AZ 4562), cylinder diameter, and its height. If the cylinder structure is too shallow, the melting process may create lenses deviating significantly from a sphere. The ratio between cylinder height and diameter is usually in range of 0.04 to 0.5. To improve quality of fabricated lenses, preshaping of



(a)



(b)

20 μm

FIGURE 49 SEM photographs showing perspective views of (a) etched multilevel mesa structure and (b) the microlens formed after mass transport.¹⁰⁰

photoresist can be applied. It is also possible to fabricate asymmetric lenses using either preshaping or noncircular masks.

The reflow technique can be further extended and combined with a molding process.^{102,103} After exposure, developing, and melting steps the lens array can be coated with PDMS and cured. In result, one obtains the negative impression of the photoresist array. Due to the fact that PDMS is hydrophobic and elastic such a created mold easily comes of the substrate and can be used in additional steps to create solgel (or other polymer) lens arrays.¹⁰³

The quality of reflow fabricated lenses is high, resulting in surface roughness in the range of $R_a = 1$ nm and deviation from perfect sphere of $\lambda/10$ (RMS) can be achieved. By application of aforementioned method¹⁰³ large sag and NA values can be obtained.

Diameters of reflow lenses usually vary between 5 and $2000 \mu\text{m}^{3,16}$ and their sag is within 2.5 and 60 to $70 \mu\text{m}$. Reported focal lengths were few and $9000 \mu\text{m}^3$ where longer focal lengths were obtained for larger lens diameters.

Laser-assisted chemical etching (LACE) can be used to make arrays of $F/0.7$ – $F/10$ lenslets with spacings of 50 to $300 \mu\text{m}$. Microlenses have been fabricated in glass, silicon, CdTe, and sapphire with 95 percent fill factors and figure quality better than $1/10$ th wave.⁹⁹ In the LACE process, a focused laser beam is scanned over a thick layer of photoresist. The irradiance of the laser beam is modulated in order to vary the exposure and thus the thickness of the developed resist. With the proper irradiance mapping, accurate lens profiles can be produced in the developed resist. If lenslet material other than photoresist is required, a pattern can be exposed in the photoresist and transferred into the new material by ion milling.

Our discussion of the mass-transport process follows discussion presented in Ref. 100. In the mass-transport process, a multilevel mesa structure is first etched into a semiconductor, as shown in Fig. 49a. The semiconductor must be a binary compound in which the evaporation rate of one element is negligible compared to that of the other. For example, InP has been used successfully. The mesa structure is placed in a furnace at an elevated temperature. Since some surface decomposition occurs with InP, a minimum phosphorus vapor pressure must be maintained in the gas ambient to prevent the sample from being transformed into metallic In. The decomposition produces free In atoms located at the crystal surface, which are in equilibrium with phosphorus in the vapor and InP in the crystal. The concentration of In in the vapor is negligible. The equilibrium concentration of free In atoms increases with increasing positive surface curvature, since the higher surface energy of the high-curvature regions translates into a lower bonding energy in the decomposition process. Consequently, a variation in curvature across the surface will result in diffusion of free In atoms from regions of high positive curvature, where the concentrations are high, to low-curvature regions, where the In-diffused atoms exceed the equilibrium concentration and have to be reincorporated into the crystal by reaction with P to form InP. (The diffusion of P in the vapor phase is presumably much faster than the diffusion of free In atoms on the surface. The latter is therefore assumed to be the rate-limiting process.) The mass transport of InP, resulting from decomposition in high-curvature regions, will continue until the difference in curvature is completely eliminated. After mass transport, a smooth profile is obtained, as shown in Fig. 49. The design of the mesa structure can result in very accurate control of the lens profile, as shown in Fig. 50. Mass-transport lens arrays have been used to collimate arrays of laser diodes, with diffraction-limited performance at $NA \sim 0.5$.

Very accurate cylindrical lenses can be drawn from preforms.¹⁰⁴ An SEM photo of an elliptical cylindrical lens is shown in Fig. 51. The first step in the process is to make a preform of a suitable glass material. Since the cross-sectional dimensions of the preform are uniformly reduced (typically ~ 50 to $100\times$) in the fiber drawing process, small manufacturing errors become optically insignificant. Therefore, standard numerically controlled grinding techniques can be utilized to generate a preform of any desired shape. Besides maintaining the preform shape, the drawing process also polishes the fiber. Results are presented⁹⁸ that demonstrate a $200\text{-}\mu\text{m}$ -wide elliptical cylindrical lens. The SFL6 preform was about 0.75 cm wide. The lens has a nominal focal length of $220 \mu\text{m}$ at $\lambda = 800 \mu\text{m}$. The lens is diffraction-limited over about a $150\text{-}\mu\text{m}$ clear aperture, or $NA \sim 0.6$. The application is to collimate the fast axis of laser diodes.

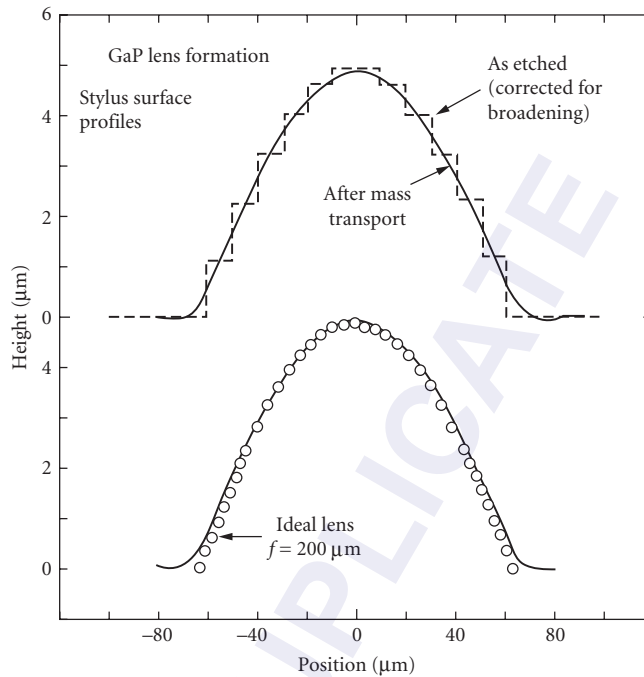


FIGURE 50 Stylus surface profiles of the multilevel mesa structure and the microlens formed after mass transport (upper half) and the comparison of the measured lens profile with an ideal one (lower half).¹⁰⁰

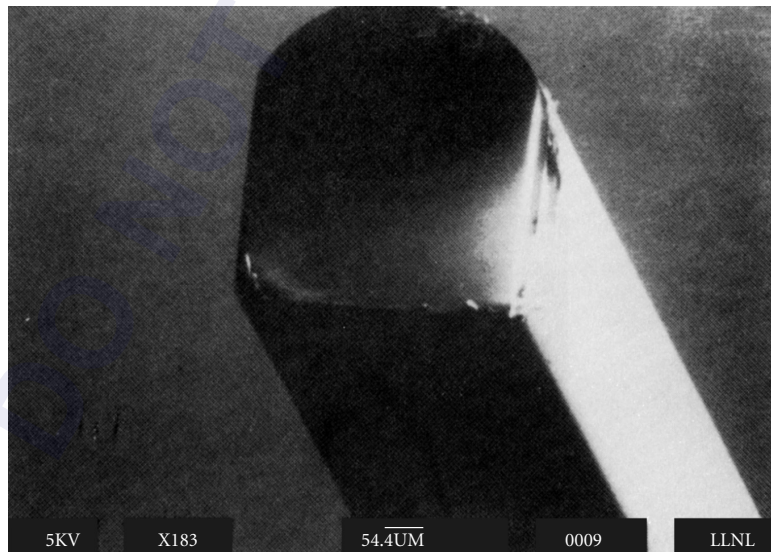


FIGURE 51 Scanning electron microscope photo of an elliptical cylindrical microlens. The lens width is $200\ \mu\text{m}$.¹⁰⁴

22.13 REFERENCES

1. M. and S. Kufner, *Micro-Optics and Lithography*, VUB Press, Brussels, 1997.
2. H. N. Herzig, ed., *Micro-Optics: Elements, Systems and Applications*, Taylor and Francis Ltd., 1997.
3. S. Sinzinger and J. Jahns, *Microoptics*, Wiley-VCH GmbH & Co. KGaA, Weinheim, Germany, 2003.
4. G. T. Sincerbox, "Miniature Optics for Optical Recording," *Proc. SPIE* **935**:63–76, 1988.
5. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York, p. 61, 1966.
6. D. K. Towner, "Scanning Techniques for Optical Data Storage," *Proc. SPIE* **695**:172–180, 1986.
7. M. Born and E. Wolf, *Principles of Optics*, Pergamon Press Ltd., Oxford, pp. 468–469, 1980.
8. J. F. Forkner and D. W. Kurtz, "Characteristics of Efficient Laser Diode Collimators," *Proc. SPIE* **740**:27–35, 1987.
9. A. R. Rouse, A. Kano, J. A. Udovich, S. M. Krotto, and A. F. Gmitro, "Design and Demonstration of a Miniature Catheter for a Confocal Microendoscope," *Appl. Opt.* **43**(31):5783, 2004.
10. C. Liang, K. B. Sung, R. Richards-Kortum, and M. R. Descour, "Fiber Confocal Reflectance Microscope (FCRM) for In-Vivo Imaging," *Opt. Exp.* **9**:821–830, 2001.
11. M. D. Chidley, K. Carlson, M. R. Descour, and R. Richards-Kortum, "Design, Assembly, and Optical Bench Testing of a High Numerical Aperture Miniature Injection-Molded Objective for Fiber-Optic Confocal Reflectance Microscopy," *Appl. Opt.* **45**:2545–2554, 2006.
12. J. D. Rogers, S. Landau, T. Tkaczyk, M. R. Descour, M. S. Rahman, R. Richards-Kortum, A. Kärkkäinen, and T. Christenson, "Performance and Imaging Results for Miniature Integrated Microendoscope," *J Biomed. Opt.* **13**(5): 054020-1-6, September/October 2008.
13. R. T. Kester, T. Tkaczyk, M. R. Descour, T. Christenson, and R. Richards-Kortum, "High Numerical Aperture Microendoscope Objective for a Fiber Confocal Reflectance Microscope," *Opt. Exp.* **15**:2409–2420, 2007. <http://www.opticsinfobase.org/abstract.cfm?URI=oe-15-5-2409>.
14. M. R. Descour, A. H. O. Kärkkäinen, J. D. Rogers, C. Liang, R. S. Weinstein, J. T. Rantala, B. Kilic, et al., "Toward the Development of Miniaturized Imaging Systems for Detection of Pre-Cancer," *IEEE J. Quantum Electron.* **38**(2), Feb. 2002.
15. A. Olszak and M. Descour, "Microscopy in Multiples," *OE Magazine*, May 2005.
16. Ph Nussbaum, R. Völkel, H. P. Herzig, M. Eisner, and S. Haselbeck, "Design, Fabrication and Testing of Microlens Arrays for Sensors and Microsystems," *Pure Appl. Opt.* **6**:617–636, 1997.
17. M. J. Madou, *Fundamentals of Microfabrication: The Science of Miniaturization*, 2d ed., CRC, Boca Raton, FL, March 13, 2002.
18. M. Gad-el-Hak, ed., *The MEMS Handbook*, 2d ed., vol. 2, "MEMS Design and Fabrication," Chap. 5, "X-Ray Based Fabrication," CRC Taylor & Francis, Boca Raton, FL, 2006.
19. T. Tkaczyk, J. D. Rogers, M. Rahman, T. Christenson, S. Gaalema, E. Dereniak, R. R. Kortum, and M. R. Descour, "Multi-Modal Miniature Microscope: 4M Device for Bio-Imaging Applications—An Overview of the System," *Proc. SPIE* **5959**:138–146, Medical Imaging, Warsaw, September 2005.
20. R. O. Maschmeyer, U.S. Patent 4,537,473, 1985.
21. D. Malacara, *Optical Shop Testing*, John Wiley and Sons, New York, 1978, Chap. 12, pp. 381–397, 1978.
22. G. W. Hopkins and R. N. Shagam, "Null Ronchi Gratings from Spot Diagrams," *Appl. Opt.* **16**:2602, 1977.
23. M. A. Fitch, "Molded Optics: Mating Precision and Mass Production," *Photonics Spectra*, 84–87, October, 1991.
24. T. Aquilina, D. Richards, and H. Pollicove, "Finished Lens Molding Saves Time and Money," *Photonics Spectra*, 73–80, September 1986.
25. M. A. Fitch and Y. K. Konishi, "Technical Directions in Molded Optics," *Proc. SPIE* **1139**:187–190, 1989.
26. D. L. Keyes, "Fundamental Properties of Optical Plastics," CRC Handbook of Laser Science and Technology, CRC Press Inc., Boca Raton, FL, in press.
27. J. D. Lytle, "Aspheric Surfaces in Polymer Optics," *Proc. SPIE* **381**:64–68, 1983.
28. D. Keyes, Private Communication.
29. R. M. Altman and J. D. Lytle, "Optical-Design Techniques for Polymer Optics," *Proc. SPIE* **237**:380–385, 1980.

30. *Modern Plastics Encyclopedia*, vol 68, McGraw-Hill, New York, 1991/1992.
31. *Plastics Technology, Manufacturing Handbook and Buyer's Guide*, vol. 37(8), Bill Communications Inc., New York, 1991/1992.
32. Diverse Optics Inc., Commercial Literature, 1992.
33. M. D. Chidley, T. Tkaczyk, R. Kestera, M. R. Descoura, "Flow-Induced Birefringence: The Hidden PSF Killer in High Performance Injection-Molded Plastic Optics," *Proc. SPIE* **6082**:60820E, 2006.
34. A. H. Deutchman, R. J. Partyka, and J. C. Lewis, "Dual Ion Beam Deposition of Diamond Films on Optical Elements," *Proc. SPIE* **1146**:124–134, 1989.
35. C. F. Cheung and W. B. Lee, *Surface Generation in Ultra-Precision: Modeling and Practices*, Professional Engineering Publishing, UK, 2003.
36. M. J. Weber, *Handbook of Optical Materials*, CRC Press, Boca Raton, FL, 2003.
37. S. Bäumer, ed., *Handbook of Plastic Optics*, Wiley-VCH, GmbH & Co. KGaA, Weinheim, Germany 2005.
38. L. W. Alvarez and W. E. Humphrey, "Variable Power Lens and System," Patent # 3,507,565 United States Patent Office, 1970.
39. A. Y. Yi and L. Li, "Design and Fabrication of a Microlens Array by Use of a Slow Tool Servo," *Opt. Lett.* **30**(13):1707, 2005.
40. W. C. Sweatt, D. D. Gill, D. P. Adams M. J. Vasile, and A. A. Claudet, "Diamond Milling of Micro-Optics," Aerospace Conference, *IEEE* 2006.
41. J. Rogers, T. Tkaczyk, M. Descour, A. Kärkkäinen, and R. Richards-Kortum, "Removal of Ghost Images by Using Tilted Element Optical Systems with Polynomial Surfaces for Aberration Compensation," *Opt. Lett.* **31**:504–506, 2006.
42. T. Tkaczyk, J. D. Rogers, M. Rahman, T. C. Christenson, S. Gaalema, E. L. Dereniak, R. Richards-Kortum, and M. R. Descour, "Multi-Modal Miniaturized Microscope: Successful Merger of Optical, MEMS, and Electronic Technologies," *Proc. SPIE* **6050**:605016-1–605016-8, 2005.
43. M. T. Gale, "Fabrication of Continuous-Relief Micro-Optical Elements by Direct Laser Writing in Photoresists," *Opt. Eng.* **33**(11):3556–3566, 1994.
44. A. G. Poleshchuk, E. G. Churin, V. P. Koronkevich, V. P. Korolkov, A. A. Kharissov, V. V. Cherkashin, V. P. Kiryanov, A. V. Kiryanov, S. A. Kokarev, and A. G. Verhoglyad, "Polar Coordinate Laser Pattern Generator for Fabrication of Diffractive Optical Elements with Arbitrary Structure," *Appl. Opt.* **38**(8):1295–1301, 1999.
45. I. Baal-Zedaka, S. Hava, N. Mirchin, R. Margolin, M. Zagon, I. Lapsker, J. Azoulay, and A. Peled, "Diffractive Optical Elements Written by Photodeposition," *Appl. Surface Science* **208–209**(1):226–232, 2003.
46. L. Erdmann, A. Deparnay, F. Wirth and R. Brunner, "MEMS Based Lithography for the Fabrication of Microoptical Components," in *Micromachining Technology for Micro-Optics and Nano-Optics II*, E. G. Johnson and G. P. Nordin eds., *Proc. SPIE* **5347**, SPIE, Bellingham, WA, 2004, pp. 79–84.
47. J. M. Tamkin, B. Bagwell, B. Kimbrough, G. Jabbour, and M. Descour, "High-Speed Gray Scale Laser Direct Write Technology for Micro-Optic Fabrication," in *Micromachining Technology for Micro-Optics and Nano-Optics*, E. G. Johnson ed., *Proc. SPIE* **4984**, SPIE, Bellingham, WA, 2003, pp. 210–218.
48. J. D. Rogers, A. H. O. Kärkkäinen, T. Tkaczyk, J. T. Rantala, and M. Descour, "Realization of Refractive Microoptics through Grayscale Lithographic Patterning of Photosensitive Hybrid Glass," *Opt. Exp.* **12**:1294–1303, 2004.
49. L. E. Schmutz, "Hartman Sensing at Adaptive Optics Associates," *Proc. SPIE* **779**:13–17, 1987.
50. D. D'Amato and R. Centamore, "Two Applications for Microlens Arrays: Detector Fill Factor Improvement and Laser Diode Collimation," *Proc. SPIE* **1544**:166–177, 1991.
51. Adaptive Optics Associates, Private communication.
52. D. D'Amato, S. Barletta, P. Cone, J. Hizny, R. Martinsen, and L. Schmutz, "Fabricating and Testing of Monolithic Lenslet Module (MLM)," in *Optical Fabrication and Testing*, pp. 74–77, 1990; *Technical Digest Series*, vol. 11, Conference ed., June 12–14, 1990, Monterey, California.
53. Adaptive Optics Associates, Commercial Literature.
54. T. D. Milster and J. N. Wong, "Modeling and Measurement of a Micro-Optic Beam Deflector," to be published in *Proc. SPIE* **1625**.
55. R. K. Luneburg, *A Mathematical Theory of Optics*, University of California Press, Berkeley, 1964, pp. 164–195.

56. K. Iga, M. Oikawa, and J. Banno, "Ray Traces in a Distributed-Index Planar Microlens," *Appl. Opt.* **21**:3451–3455, 1982.
57. M. Oikawa, H. Nemoto, K. Hamanaka, H. Imanishi, and T. Kishimoto, "Light Coupling Characteristics of Planar Microlens," *Proc. SPIE* **1544**:226–237, 1991.
58. K. Iga and S. Misawa, "Distributed-Index Planar Microlens and Stacked Planar Optics: A Review of Progress," *Appl. Opt.* **25**:3388, 1986.
59. K. Iga, Y. Kokobun, and M. Okawa, *Fundamentals of Microoptics*, Academic Press Inc., Tokyo, 1984.
60. Y. Ohtsuka, "Light-focusing Plastic Rod Prepared from Diallyl Isophthalate-Methyl Methacrylate Copolymerization," *Appl. Phys. Lett.* **23**:247, 1973.
61. K. Iga and N. Yamanamoto, "Plastic Focusing Fiber for Imaging Applications," *Appl. Opt.* **16**:1305, 1977.
62. T. Izawa and H. Nakagome, "Optical Waveguide Formed by Electrically Induced Migration of Ions in Glass Plates," *Appl. Phys. Lett.* **21**:584, 1972.
63. G. D. Khoe, H. G. Kock, J. A. Luijendijk, C. H. J. van den Brekel, and D. Küppers, "Plasma CVD Prepared $\text{SiO}_2/\text{Si}_3\text{N}_4$ Graded Index Lenses Integrated in Windows of Laser Diode Packages," in *Technical Digest, Seventh European Conference on Optical Communication*, Copenhagen, pp. 7.6-1–7.6-4, 1981.
64. K. Iga, M. Oikawa, and T. Sanada, *Electron. Lett.* **17**:452, 1981.
65. M. Oikawa, H. Nemoto, K. Hamanaka, and E. Okuda, "High Numerical Aperture Planar Microlenses with Swelled Structure," *Appl. Opt.* **29**:4077, 1990.
66. O. E. Meyres, "Studies of Transmission Zone Plates," *Am. J. Phys.* **19**:359, 1951.
67. M. Sussman, "Elementary Diffraction Theory of Zone Plates," *Am. J. Phys.* **28**:394, 1960.
68. H. Nishihara and T. Suhara, *Progress in Optics*, North-Holland, New York, vol. 24, pp. 1–37, 1987.
69. M. Young, "Zone Plates and Their Aberrations," *J. Opt. Soc. Am.* **62**:972–976, 1972.
70. T. D. Milster, R. M. Trusty, M. S. Wang, F. F. Froehlich, and J. K. Erwin, "Micro-Optic Lens for Data Storage," *Proc. SPIE* **1499**:286–292, 1991.
71. T. Stone and N. George, "Hybrid Diffractive-Refractive Lenses and Achromats," *Appl. Opt.* **27**:2960–2971, 1988.
72. G. P. Behrmann and J. P. Bowen, "Thermal Effects in Diffractive Lenses," in *Diffractive Optics: Design, Fabrication, and Applications, Technical Digest*, Optical Society of America, Washington, D.C., vol. 9, pp. 8–10, 1992.
73. R. Magnusson and T. K. Gaylord, "Diffraction Efficiencies of Thin Phase Gratings with Arbitrary Grating Shape," *J. Opt. Soc. Am.* **68**:806–809, 1978.
74. J. Futhey and M. Fleming, "Superzone Diffractive Lenses," in *Diffractive Optics: Design, Fabrication, and Applications Technical Digest*, Optical Society of America, Washington, D.C., vol. 9, pp. 4–6, 1992.
75. T. Fujita, H. Nishihara, and J. Koyama, "Blazed Gratings and Fresnel Lenses Fabricated by Electron-Beam Lithography," *Opt. Lett.* **7**:578–580, 1982.
76. R. W. Schneck, "Process Factors for Electroforming Video Disks," *Plating and Surface Finishing*, **71**(1):38, 1984.
77. K. Goto, K. Mori, G. Hatakoshi, and S. Takahashi, "Spherical Grating Objective Lenses for Optical Disk Pickups," *Proc. International Symposium on Optical Memory*, 1987, *Jpn. J. Appl. Phys.* **26**(suppl. 26–24):135, 1987.
78. M. Tanigami, S. Aoyama, T. Yamashita, and K. Imanaka, "Low Wavefront Aberration and High Temperature Stability Molded Micro-Fresnel Lens," *IEEE Photonics Technology Letters* **1**(11):384–385, 1989.
79. H. Hosokawa and T. Yamashita, "ZnS Micro-Fresnel Lens and Its Uses," *Appl. Opt.* **29**:5706–5710, 1990.
80. S. Rege, T. Tkaczyk, and M. Descour, "Application of the Alvarez-Humphrey Concept to the Design of a Miniaturized Scanning Microscope," *Opt. Exp.* **12**:2574–2588, 2004.
81. B. H. W. Hendriks, S. Kuiper, M. A. J. Van Aa, C. A. Renders, and T. W. Tukker, "Electrowetting-Based Variable-Focus Lens for Miniature Systems," *Opt. Rev.* **12**(3):255–259, 2005.
82. S. Kuiper, B. H. W. Hendriks, R. A. Hayes, B. J. Feenstra, and J. M. E. Baken, "Electrowetting-Based Optics," *Proc. SPIE* **5908**, SPIE, Bellingham, WA, 2005.
83. N. Chronis, G. L. Liu, K.-H. Jeong, and L. P. Lee, "Tunable Liquid-Filled Microlens Array Integrated with Microfluidic Network," *Opt. Exp.* **11**(19):2370, 2003.
84. K.-H. Jeong, G. L. Liu, N. Chronis, and L. P. Lee, "Tunable Microdoublet Lens Array," *Opt. Exp.* **12**(11):2494, 2004.

85. K.-S. Hong, J. Wang, A. Sharonov, D. Chandra, J. Aizenberg, and S. Yang, "Tunable Microfluidic Optical Devices with an Integrated Microlens Array," *J. Micromech. Microeng.* **16**:1660–1666, 2006.
86. W. Wang and J. Fang, "Design, Fabrication and Testing of a Micromachined Integrated Tunable Microlens," *J. Micromech. Microeng.* **16**:1221–1226, 2006.
87. A. F. Naumov, M. Yu. Loktev, I. R. Guralnik, and G. V. Vdovin. "Liquid Crystal Adaptive Lenses with Modal Control," *Opt. Lett.* **23**:992–994, 1998.
88. F. Krogmann, R. Shaik, W. Mönch, and H. Zappe, "Repositionable Liquid Micro-Lenses with Variable Focal Length," *IEEE 707*, MEMS 2007, Kobe, Japan, 21–25 January 2007.
89. F. Krogmann, H. Qu, W. Moench, and H. Zappe, "Light-Actuated Push/Pull Manipulation of Liquid Droplets," *IEEE*, 2006.
90. K. Aljaseem, A. Werber, A. Seifert, and H. Zappe, "Fiber Optic Tunable Probe for Endoscopic Optical Coherence Tomography," *J. Opt. A: Pure Appl. Opt.* **10**(044012):8, 2008.
91. M. Ye, B. Wang, and S. Sato, "Liquid-Crystal Lens with a Focal Length that is Variable in a Wide Range," *Appl. Opt.* **43**(35):6407, 2004.
92. B. Wang, M. Ye, and S. Sato, "Lens of Electrically Controllable Focal Length made by a Glass Lens and Liquid-Crystal Layers," *Appl. Opt.* **43**(17):3420, 2004.
93. A. K. Kirby, P. J. W. Hands, and G. D. Love, "Liquid Crystal Multi-Mode Lenses and Axicons Based on Electronic Phase Shift Control," *Opt. Exp.* **15**(21):13496, 2007.
94. O. Pishnyak, S. Sato, and O. D. Lavrentovich, "Electrically Tunable Lens Based on a Dual-Frequency Nematic Liquid Crystal," *Appl. Opt.* **45**(19):4576, 2006.
95. H. Ren, Y.-H. Fan, and S.-T. Wu, "Liquid-Crystal Microlens Arrays Using Patterned Polymer Networks," *Opt. Lett.* **29**(14), July 15, 2004.
96. C.-C. Cheng, C. A. Chang, C.-H. Liu, and J. A. Yeh, "A Tunable Liquid-Crystal Microlens with Hybrid Alignment," *J. Opt. A: Pure Appl. Opt.* **8**:S365–S369, 2006.
97. L. Dong, A. K. Agarwal, D. J. Beebe, and H. Jiang, "Adaptive Liquid Microlenses Activated by Stimuli-responsive Hydrogels," *Nature* **442**, 2006.
98. Z. D. Popovic, R. A. Sprague, and G. A. Neville Connell, "Technique for Monolithic Fabrication of Microlens Arrays," *Appl. Opt.* **27**:1281–1284, 1988.
99. E. J. Gratrix and C. B. Zarowin, "Fabrication of Microlenses by Laser Assisted Chemical Etching (LACE)," *Proc. SPIE* **1544**:238–243, 1991.
100. V. Diadiuk, Z. L. Liao, and J. N. Walpole, "Fabrication and Characterization of Semiconductor Microlens Arrays," *Proc. SPIE* **1354**:496–500, 1990.
101. H. W. Choi, E. Gu, C. Liu, C. Griffin, J. M. Girkin, I. M. Watson, and M. D. Dawson, "Fabrication of Natural Diamond Microlenses by Plasma Etching," *J. Vac. Sci. Technol. B* **23**(1):130–132, Jan/Feb 2005.
102. M. V. Kunnavaikkam, F. M. Houlihan, M. Schlax, J. A. Liddle, P. Kolodner, O. Nalamasu, and J. A. Rogers, "Low-Cost, Low-Loss Microlens Arrays Fabricated by Soft-Lithography Replication Process," *Appl. Phys. Lett.* **82**(8):1152–1154, 2003.
103. X.-C. Yuan, W. X. Yu, M. He, J. Bu, W. C. Cheong, H. B. Niu, and X. Peng, "Soft-Lithography-Enabled Fabrication of Large Numerical Aperture Refractive Microlens Array in Hybrid SiO₂-TiO₂ Sol-Gel Glass," *Appl. Phys. Lett.* **86**:114102, 2005.
104. J. J. Snyder, P. Reichert, and T. M. Baer, "Fast Diffraction-Limited Cylindrical Microlenses," *Appl. Opt.* **30**:2743–2747, 1991.
105. Ohara Optical Glass Catalog, Ohara Corporation, Somerville, New Jersey, 1990.
106. Corning, Incorporated, Precision Molded Optics Department, Corning, New York, Commercial Literature, 1990.
107. Y. Kokuban and K. Iga, *Appl. Opt.* **21**:1030, 1982.
108. Y. Kokuban, T. Usui, M. Oikawa, and K. Iga, "Wave Aberration Testing System for Microlenses by Shearing Interference Method," *Jap. J. of Appl. Phys.* **23**(1):101–104, 1984.
109. R. H. Bellman, N. F. Borrelli, L. G. Mann, and J. M. Quintal, "Fabrication and Performance of a One-to-One Erect Imaging Microlens Array for Fax," *Proc. SPIE* **1544**:209–217, 1991.

BINARY OPTICS

Michael W. Farn and Wilfrid B. Veldkamp

*MIT/Lincoln Laboratory
Lexington, Massachusetts*

23.1 GLOSSARY

A	aspheric
C	describes spherical aberration
C_m	Fourier coefficients
c	curvature
$c(x, y)$	complex transmittance
D	local period
f	focal length
k, l	running indices
l_i	paraxial image position
L, M	direction cosines
m	diffraction order
P	partial dispersion
s	spheric
t	thickness
V_d	Abbe number
x, y, z	Cartesian coordinates
λ	wavelength
η	diffraction efficiency
ξ_i	paraxial image height
$\phi(x, y)$	phase

23.2 INTRODUCTION

Binary optics is a surface-relief optics technology based on VLSI fabrication techniques (primarily photolithography and etching), with the “binary” in the name referring to the binary coding scheme used in creating the photolithographic masks. The technology allows the creation of new, unconventional optical elements and provides greater design freedom and new materials choices for conventional elements. This capability allows designers to create innovative components that can solve problems in optical sensors, optical communications, and optical processors. The technology has advanced sufficiently to allow the production of diffractive elements, hybrid refractive-diffractive elements, and refractive micro-optics which are satisfactory for use in cameras, military systems, medical applications, and other demanding areas.

The boundaries of the binary optics field are not clearly defined, so in this chapter, the concentration will be on the core of the technology: passive optical elements which are fabricated using VLSI technology. As so defined, binary optics technology can be broadly divided into the areas of optical design and VLSI-based fabrication. Optical design can be further categorized according to the optical theory used to model the element: geometrical optics, scalar diffraction theory, or vector diffraction theory; while fabrication is composed of two parts: translation of the optical design into the mask layout and the actual micromachining of the element. The following sections discuss each of these topics in some detail, with emphasis on optical design. For a more general overview, the reader is referred to Refs. 1 for many of the original papers, 2 and 3 for a sampling of applications and research, and 4 to 6 for a lay overview.

Directly related areas which are discussed in other chapters but not in this one include micro-optics and diffractive optics fabricated by other means (e.g., diamond turning, conventional manufacturing, or optical production), display holography (especially computer-generated holography), mass replication technologies (e.g., embossing, injection molding, or epoxy casting), integrated optics, and other micromachining technologies.

23.3 DESIGN—GEOMETRICAL OPTICS

In many applications, binary optics elements are designed by ray tracing and “classical” lens design principles. These designs can be divided into two classes: broadband and monochromatic. In broadband applications, the binary optics structure has little optical power in order to reduce the chromatic aberrations and its primary purpose is aberration correction. The device can be viewed as an aspheric aberration corrector, similar to a Schmidt corrector, when used to correct the monochromatic aberrations and it can be viewed as a material with dispersion an order of magnitude greater than and opposite in sign to conventional materials when used to correct chromatic aberrations. In monochromatic applications, binary optics components can have significant optical power and can be viewed as replacements for refractive optics.

In both classes of designs, binary optics typically offers the following key advantages:

- Reduction in system size, weight, and/or number of elements
- Elimination of exotic materials
- Increased design freedom in correcting aberrations, resulting in better system performance
- Generation of arbitrary lens shapes (including micro-optics) and phase profiles

Analytical Models

Representation of a Binary Optics Element As with any diffractive element, a binary optics structure is defined by its phase profile $\phi(x, y)$ (z is taken as the optical axis), design wavelength λ_0 , and the surface on which the element lies. For simplicity, this surface is assumed to be planar for the remainder of this chapter, although this is commonly not the case. For example, in many refractive/diffractive systems,

the binary optics structure is placed on a refractive lens which may be curved. The phase function is commonly represented by either explicit analytical expression or decomposition into polynomials in x and y (e.g., the HOE option in CODE V).

Explicit analytic expressions are used in simple designs, the two most common being lenses and gratings. A lens used to image point (x_o, y_o, z_o) to point (x_i, y_i, z_i) at wavelength λ_0 has a phase profile

$$\phi(x, y) = \frac{2\pi}{\lambda_0} \left[z_o \left(\sqrt{(x-x_o)^2/z_o^2 + (y-y_o)^2/z_o^2 + 1} - 1 \right) - z_i \left(\sqrt{(x-x_i)^2/z_i^2 + (y-y_i)^2/z_i^2 + 1} - 1 \right) \right] \quad (1)$$

where z_o and z_i are both taken as positive to the right of the lens. The focal length is given by the gaussian lens formula:

$$1/f_0 = 1/z_i - 1/z_o \quad (2)$$

with the subscript indicating that f_0 is the focal length at λ_0 . A grating which deflects a normally incident ray of wavelength λ_0 to the direction with direction cosines (L, M) is described by

$$\phi(x, y) = \frac{2\pi}{\lambda_0} (xL + yM) \quad (3)$$

Axicons are circular gratings and are described by

$$\phi(x, y) = \frac{2\pi}{\lambda_0} (\sqrt{x^2 + y^2} L) \quad (4)$$

where L now describes the radial deflection.

For historical reasons, the polynomial decomposition of the phase profile of the element commonly consists of a spheric term and an aspheric term:

$$\phi(x, y) = \phi_s(x, y) + \phi_A(x, y) \quad (5)$$

where

$$\phi_A(x, y) = \frac{2\pi}{\lambda_0} \sum_k \sum_l a_{kl} x^k y^l$$

and the spheric term $\phi_s(x, y)$ takes the form of Eq. (1). Since the phase profiles produced by binary optics technology are not constrained to be spheric, $\phi_s(x, y)$ is often set to zero by using the same object and image locations and the aspheric term alone is used to describe the profile. The binary optics element is then optimized by optimizing the polynomial coefficients a_{kl} . If necessary, the aspheric term can be forced to be radially symmetric by constraining the appropriate polynomial coefficients.

It is possible to describe the phase profile of a binary optics element in other ways. For example, $\phi(x, y)$ could be described by Zernike polynomials or could be interpolated from a two-dimensional look-up table.

Ray Tracing by the Grating Equation A binary optics element with phase $\phi(x, y)$ can be ray traced using the grating equation by modeling the element as a grating, the period of which varies with position. This yields

$$L' = L + \frac{m\lambda}{2\pi} \frac{\partial \phi}{\partial x} \quad (6)$$

$$M' = M + \frac{m\lambda}{2\pi} \frac{\partial \phi}{\partial y} \quad (7)$$

where m is the diffracted order, L, M are the direction cosines of the incident ray, and L', M' are the direction cosines of the diffracted ray.⁷ In geometrical designs, the element is usually blazed for the first order ($m = 1$). Note that it is the phase gradient $\nabla\phi(x, y)$ (a vector quantity proportional to the local spatial frequency) and not the phase $\phi(x, y)$ which appears in the grating equation. The magnitude of the local period is inversely proportional to the local spatial frequency and given by

$$D(x, y) = 2\pi / |\nabla\phi| \quad (8)$$

where $||$ denotes the vector magnitude. The minimum local period determines the minimum feature size of the binary optics structure, a concern in device fabrication (see “Fabrication” later in this chapter).

Ray Tracing by the Sweatt Model The Sweatt model,⁸ which is an approximation to the grating equation, is another method for ray tracing. The Sweatt approach models a binary optics element as an equivalent refractive element and is important since it allows results derived for refractive optics to be applied to binary optics. In the Sweatt model, a binary optics element with phase $\phi(x, y)$ at wavelength λ_0 is replaced by a refractive equivalent with thickness and refractive index given by

$$t(x, y) = \frac{\lambda_0}{n_0 - 1} \frac{\phi(x, y)}{2\pi} + t_0 \quad (9)$$

$$n(\lambda) - 1 = \frac{\lambda}{\lambda_0} (n_0 - 1) \quad (10)$$

Here, t_0 is a constant chosen to make $t(x, y)$ always positive and n_0 is the index of the material at wavelength λ_0 . The index n_0 is chosen by the designer and as $n_0 \rightarrow \infty$, the Sweatt model approaches the grating equation. In practice, a value of $n_0 = 10,000$ is sufficiently high for accurate results.⁹

In the special case of a binary optics lens described by Eq. (1), the more accurate Sweatt lens¹⁰ can be used. In this case, the element is modeled by two surfaces of curvature

$$c_o = 1 / [(1 - n_0)z_o] \quad (11)$$

$$c_i = 1 / [(1 - n_0)z_i] \quad (12)$$

and conic constant $-n_0^2$, with the axis of each surface passing through the respective point source. The refractive index is still modeled by Eq. (10).

Aberration Correction

Aberrations of a Binary Optics Singlet As a simple example of a monochromatic imaging system, consider a binary optics singlet which is designed to image the point $(0, 0, z_o)$ to the point $(0, 0, z_i)$

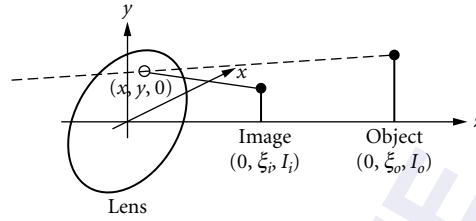


FIGURE 1 Primary aberrations of a binary optics lens.⁷

at wavelength λ_0 . The phase profile of this lens can be derived from Eq. (1) and the focal length f_0 from Eq. (2). Now consider an object point of wavelength λ located at $(0, \xi_o, l_o)$. The lens will form an image at $(0, \xi_i, l_i)$ (see Fig. 1), with the paraxial image position l_i and height ξ_i given by⁷

$$\frac{1}{l_i} = \frac{\lambda}{f_0 \lambda_0} + \frac{1}{l_o} \quad (13)$$

$$\xi_i / l_i = \xi_o / l_o \quad (14)$$

Note that the first equation is just the gaussian lens law but using a wavelength-dependent focal length of

$$f(\lambda) = f_0 \frac{\lambda_0}{\lambda} \quad (15)$$

The focal length being inversely proportional to the wavelength is a fundamental property of diffractive lenses. In addition, due to the wavelength shift and position change of the object point, the lens will form a wavefront with a primary aberration of⁷

$$\begin{aligned} W(x, y) = & \frac{1}{8} \left[\left(\frac{1}{l_i^3} - \frac{1}{l_o^3} \right) - \frac{\lambda}{\lambda_0} \left(\frac{1}{z_i^3} - \frac{1}{z_o^3} \right) \right] (x^2 + y^2)^2 \\ & - \frac{1}{2l_i} \left(\frac{1}{l_i^2} - \frac{1}{l_o^2} \right) \xi_i y (x^2 + y^2) \\ & + \frac{3}{4l_i^2} \left(\frac{1}{l_i} - \frac{1}{l_o} \right) \xi_i^2 y^2 + \frac{1}{4l_i^2} \left(\frac{1}{l_i} - \frac{1}{l_o} \right) \xi_i^2 x^2 \end{aligned} \quad (16)$$

where the ray strikes the lens at (x, y) . The first term is spherical aberration, the second is coma, and the last two are tangential and sagittal field curvature. As noted by Welford, all the off-axis aberrations can be eliminated if and only if $l_i = l_o$, a useless configuration. In most systems of interest, the limiting aberration is coma.

The performance of the binary optics singlet can be improved by introducing more degrees of freedom: varying the stop position, allowing the binary optics lens to be placed on a curved surface, using additional elements, etc. For a more detailed discussion, see Refs. 1, 7, and 11.

Chromatic Aberration Correction Binary optics lenses inherently suffer from large chromatic aberrations, the wavelength-dependent focal length [Eq. (15)] being a prime example. By themselves,

they are unsuitable for broadband imaging and it has been shown that an achromatic system consisting only of diffractive lenses cannot produce a real image.¹²

However, these lenses can be combined successfully with refractive lenses to achieve chromatic correction (for a more detailed discussion than what follows, see Refs. 4, 13, and 14). The chromatic behavior can be understood by using the Sweatt model, which states that a binary optics lens behaves like an ultrahigh index refractive lens with an index which varies linearly with wavelength [let $n_0 \rightarrow \infty$ in Eq. (10)]. Accordingly, they can be used to correct the primary chromatic aberration of conventional refractive lenses but cannot correct the secondary spectrum. For the design of achromats and apochromats, an effective Abbe number and partial dispersion can also be calculated. For example, using the C, d, and F lines, the Abbe number is defined as $V_d = [n(\lambda_d) - 1] / [n(\lambda_F) - n(\lambda_C)]$. Substituting Eq. (10) and letting $n_0 \rightarrow \infty$ yields

$$V_d = \lambda_d / (\lambda_F - \lambda_C) = -3.45 \quad (17)$$

In a similar fashion, the effective partial dispersion using the g and F lines is

$$P_{gF} = (\lambda_g - \lambda_F) / (\lambda_F - \lambda_C) = 0.296 \quad (18)$$

By using these effective values, the conventional procedure for designing achromats and apochromats¹⁵ can be extended to designs in which one element is a binary optics lens.

Figure 2 plots the partial dispersion P_{gF} versus Abbe number V_d for various glasses. Unlike all other materials, a binary optics lens has a negative Abbe number. Thus, an achromatic doublet can be formed by combining a refractive lens and a binary optics lens, both with positive power. This significantly reduces the lens curvatures required, allowing for larger apertures. In addition, the binary optics lens has a position in Fig. 2 which is not collinear with the other glasses, thus also allowing the design of apochromats with reduced lens curvatures and larger apertures.

Monochromatic Aberration Correction For a detailed discussion, the reader is referred to Refs. 1 and 11. As a simple example,⁴ consider a refractive system which suffers from third-order spherical aberration and has a residual phase given by

$$\phi_r(x, y) = \frac{2\pi}{\lambda} C(x^2 + y^2)^2 \quad (19)$$

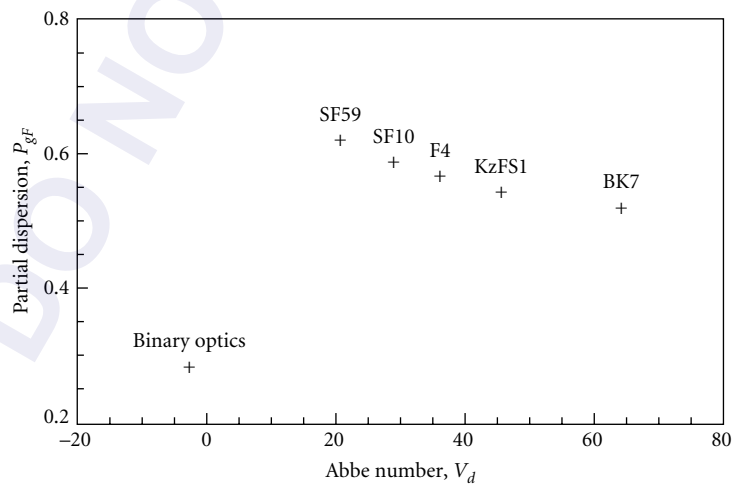


FIGURE 2 Partial dispersion vs. Abbe number.¹⁴

where C describes the spherical aberration. Then, a binary optics corrector with phase

$$\phi_b(x, y) = -\frac{2\pi}{\lambda_0} C(x^2 + y^2)^2 \quad (20)$$

will completely correct the aberration at wavelength λ_0 and will reduce the aberration at other wavelengths to

$$\phi_r + \phi_b = \frac{2\pi}{\lambda} C(1 - \lambda / \lambda_0)(x^2 + y^2)^2 \quad (21)$$

The residual aberration is spherochromatism.

Micro-Optics

Binary optics technology is especially suited for the fabrication of micro-optics and micro-optics arrays, as shown in Fig. 3. The advantages of binary optics technology include the following:

- *Uniformity and coherence.* If desired, all micro-optics in an array can be made identical to optical tolerances. This results in coherence over the entire array (see Fig. 4).
- *Refractive optics.* Binary optics is usually associated with diffractive optics. This is not a fundamental limit but results primarily from fabrication constraints on the maximum achievable depth (typically, $3 \mu\text{m}$ with ease and up to $20 \mu\text{m}$ with effort). However, for many micro-optics, this is sufficient to allow the etching of refractive elements. For example, a lens of radius R_0

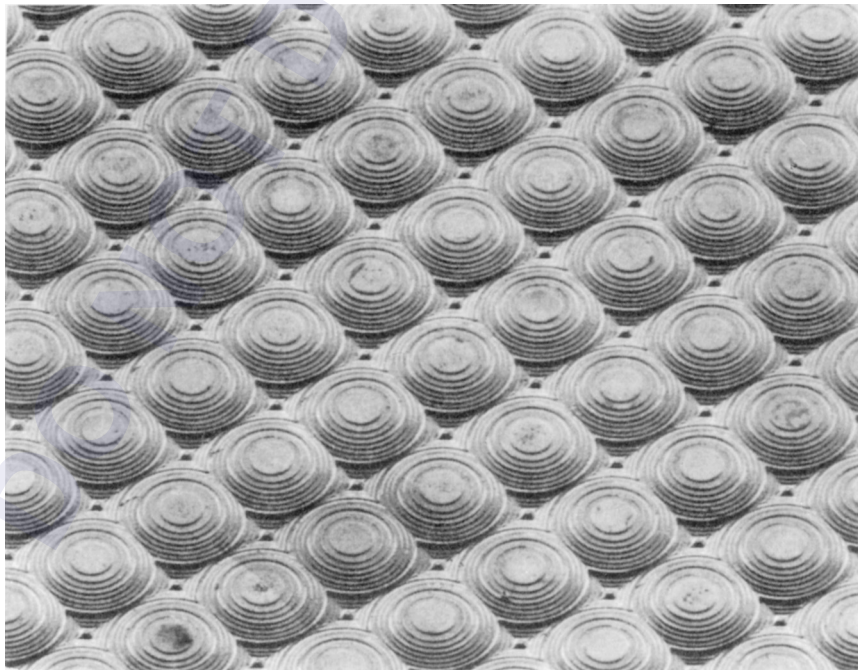


FIGURE 3 96×64 array of $51 \times 61 \mu\text{m}$ CdTe microlenses.

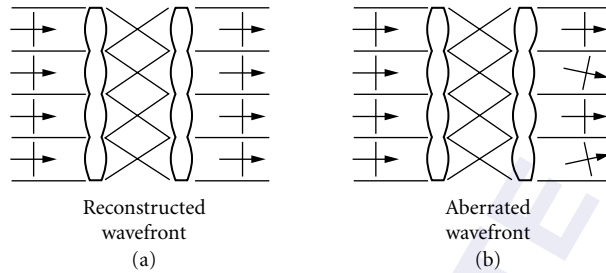


FIGURE 4 Micro-optic telescope using (a) coherent arrays and (b) incoherent arrays.

which is corrected for spherical aberration¹⁵ and focuses collimated light at a distance z_0 (see Fig. 5) has a thickness of

$$t_{\max} = n \left[\sqrt{R_0^2 + z_0^2} - z_0 \right] / (n-1) \quad (22)$$

where n is the index of the material.

- *Arbitrary phase profiles.* Binary optics can produce arbitrary phase profiles in micro-optics just as easily as in macro-optics. Fabricating arrays of anamorphic lenses to correct the astigmatism of semiconductor lasers, for example, is no more difficult than fabricating arrays of conventional spherical lenses.
- *100 percent fill factor.* While many technologies are limited in fill factor (e.g., round lenses on a square grid yield a 79 percent fill factor), binary optics can achieve 100 percent fill factor on any shape grid.
- *Spatial multiplexing.* Each micro-optic in an array can be different from its neighbors and the array itself can compose an arbitrary mosaic rather than a regular grid. For example, a binary optics array of individually designed micro-optics can be used to optimally mode-match one-dimensional laser arrays to laser cavities or optical fibers.¹⁶

Optical Performance

Wavefront Quality The wavefront quality of binary optics components is determined by the accuracy with which the lateral features of the element are reproduced. Since the local period (typically several μm) is usually much larger than the resolution with which it can be reproduced (of order $0.1 \mu\text{m}$), wavefront quality is excellent. In fact, wavefront errors are typically limited by the optical quality of the substrate rather than the quality of the fabrication.

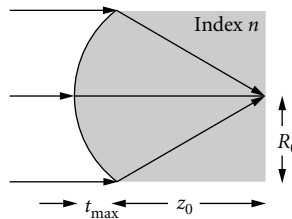


FIGURE 5 Thickness of a refractive lens.¹⁵

Diffraction Efficiency The diffraction efficiency of a device is determined by how closely the binary optics stepped-phase profile approximates a true blaze. The theoretical efficiency at wavelength λ of an element with I steps designed for use at λ_0 is⁴

$$\eta(\lambda, I) = \left| \text{sinc}(1/I) \frac{\sin(I\pi\alpha)}{I \sin \pi\alpha} \right|^2 \quad (23)$$

where $\text{sinc}(x) = \sin(\pi x)/(\pi x)$ and $\alpha = (\lambda_0/\lambda - 1)/I$.

This result is based on scalar theory, assumes perfect fabrication, and neglects any material dispersion. Figure 6 plots the efficiency $\eta(\lambda, I)$ for different numbers of steps I ; while Table 1 gives the average efficiency over the bandwidth $\Delta\lambda$ for a perfectly blazed element ($I \rightarrow \infty$).⁴ The efficiency equation is asymmetric in λ but symmetric in $1/\lambda$.

The use of scalar theory in the previous equation assumes that the local period $D(x, y)$ [see Eq. (8)] is large compared to the wavelength. As a rule of thumb, this assumption begins to lose validity when the period dips below 10 wavelengths (e.g., a grating with period less than $10\lambda_0$ or a lens faster than $F/5$) and lower efficiencies can be expected in these cases. For a more detailed discussion, see Ref. 17.

The efficiency discussed here is the diffraction efficiency of an element. Light lost in this context is primarily diffracted into other diffraction orders, which can also be traced through a system to determine their effect. As with conventional elements, binary optics elements will also suffer reflection losses which can be minimized in the usual manner.

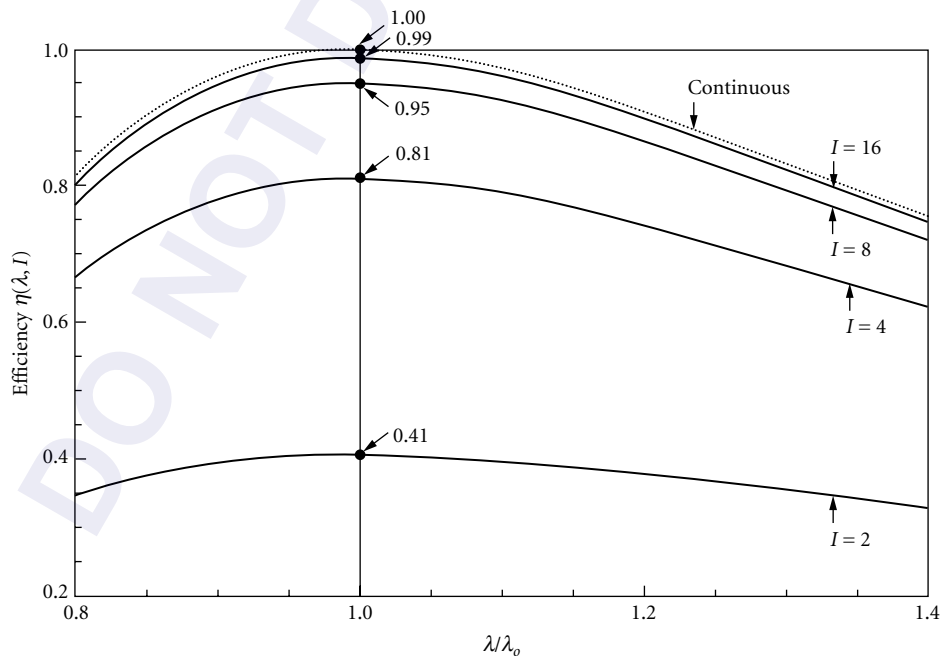


FIGURE 6 Diffraction efficiency of binary optics.⁴

TABLE 1 Average Diffraction Efficiency for Various Bandwidths⁴

$\Delta\lambda/\lambda_0$	$\bar{\eta}$
0.00	1.00
0.10	1.00
0.20	0.99
0.30	0.98
0.40	0.96
0.50	0.93
0.60	0.90

23.4 DESIGN—SCALAR DIFFRACTION THEORY

Designs via scalar diffraction theory are based on the direct manipulation of the phase of a wavefront. The incident wavefront is generally from a coherent source and the binary optics element manipulates the phase of each point of the wavefront such that the points interfere constructively or destructively, as desired, at points downstream of the element. In this regime, binary optics can perform some unique functions, two major applications being wavefront multiplexing and beam shaping.

Analytical Models

In the scalar regime, the binary optics component with phase profile $\phi(x, y)$ is modeled as a thin-phase screen with a complex transmittance of

$$c(x, y) = \exp[j\phi(x, y)] \quad (24)$$

The phase screen retards the incident wavefront and propagation of the new wavefront is modeled by the appropriate scalar formulation (e.g., angular spectrum, Fresnel diffraction, Fraunhofer diffraction) for nonperiodic cases, or by Fourier series decomposition for periodic cases.

The design of linear gratings is an important problem in the scalar regime since other problems can be solved by analogy. A grating with complex transmittance $c(x)$ and period D can be decomposed into its Fourier coefficients C_m , where

$$C_m = \frac{1}{D} \int_0^D c(x) \exp(-j2\pi mx/D) dx \quad (25)$$

$$c(x) = \sum_{m=-\infty}^{\infty} C_m \exp(j2\pi mx/D) \quad (26)$$

The relative intensity or efficiency of the m th diffracted order of the grating is

$$\eta_m = |C_m|^2 \quad (27)$$

Due to the fabrication process, binary optics gratings are piecewise flat. The grating transmission in this special case can be expressed as $c(x) = c_i$ for $x_i < x < x_{i+1}$, where c_i is the complex transmission of step i of I total steps, $x_0 = 0$, and $x_I = D$. The Fourier coefficients then take the form

$$C_m = \sum_{i=0}^{I-1} c_i \delta_i \exp(-j2\pi m \Delta_i) \text{sinc}(m \delta_i) \quad (28)$$

where $\delta_i = (x_{i+1} - x_i)/D$ and $\Delta_i = (x_{i+1} + x_i)/(2D)$.

The sinc term is due to the piecewise flat nature of the grating. If, in addition to the above, the grating transition points are equally spaced, then $x_i = iD/I$ and Eq. (28) reduces to

$$C_m = \exp(-j\pi m/I) \operatorname{sinc}(m/I) \left[\frac{1}{I} \sum_{i=0}^{I-1} c_i \exp(-j2\pi mi/I) \right] \quad (29)$$

The bracketed term is the FFT of c_i , which makes this case attractive for numerical optimizations. If the complex transmittance is also stepped in phase by increments of ϕ_0 , then $c_i = \exp(ji\phi_0)$ and Eq. (29) further reduces to¹⁸

$$C_m = \exp\{j\pi[(I-1)\alpha - m/I]\} \operatorname{sinc}(m/I) \frac{\sin(I\pi\alpha)}{I \sin \pi\alpha} \quad (30)$$

where $\alpha = \phi_0/(2\pi) - m/I$. This important case occurs whenever a true blaze is approximated by a stepped-phase profile. The efficiency equation [Eq. (23)] is a further specialization of this case.

Wavefront Multiplexers

Grating Designs Grating multiplexers (also known as beam-splitter gratings) split one beam into many diffracted beams which may be of equal intensity or weighted in intensity.¹⁹ Table 2 shows some common designs. In general, the designs can be divided into two categories: continuous phase and binary. Continuous phase multiplexers generally have better performance, as measured by the total efficiency and intensity uniformity of the diffracted beams, while binary multiplexers are easier to fabricate (with the exception of several naturally occurring continuous phase profiles). Upper bounds for the efficiency of both continuous and binary types are derived in Ref. 20.

If the phase is allowed to be continuous or nearly continuous (8 or 16 phase levels), then the grating design problem is analogous to the phase retrieval problem and iterative techniques are commonly used.²¹ A generic problem is the design of a multiplexer to split one beam into K equal intensity beams. Fanouts up to 1:50 with perfect uniformity and efficiencies of 90–100 percent are typical.

The complex transmittance of a binary grating has only two possible values [typically +1 and -1, or $\exp(j\phi_0)$ and $\exp(-j\phi_0)$], with the value changing at the transition points of the grating. By nature, the response of these gratings has the following properties:

- The intensity response is symmetric; that is, $\eta_m = \eta_{-m}$.
- The relative intensities of the nonzero orders are determined strictly by the transition points. That is, if the transition points are held constant, then the ratios η_m/η_n for all $m, n \neq 0$ will be constant, regardless of the actual complex transmittance values.
- The complex transmittance values only affect the balance of energy between the zero and non-zero orders.

TABLE 2 Grating Multiplexers of Period D , $0 < x < D$

Phase Profile	η_{-1}	η_0	η_1	Remarks
$\phi(x, y) = \begin{cases} 0 & x < D/2 \\ \pi & D/2 < x \end{cases}$	0.41	0	0.41	Binary 1:2 splitter
$\phi(x, y) = \begin{cases} 0 & x < D/2 \\ 2.01 & D/2 < x \end{cases}$	0.29	0.29	0.29	Binary 1:3 splitter
$\phi(x, y) = \pi x/D$	×	0.41	0.41	Continuous 1:2 splitter
$\phi(x, y) = \arctan[2.657 \cos(2\pi x/D)]$	0.31	0.31	0.31	Continuous 1:3 splitter

Binary gratings are usually designed via the Dammann approach or search methods and tables of binary designs have been compiled.^{22,23} Efficiencies of 60–90 percent are typical for the 1:K beam-splitter problem.

Multifocal Lenses The concepts used to design gratings with multiple orders can be directly extended to lenses and axicons to design elements with multiple focal lengths by taking advantage of the fact that while gratings are periodic in x , paraxial lenses are periodic in $(x^2 + y^2)$, nonparaxial lenses in $\sqrt{x^2 + y^2 + f_0^2}$, and axicons in $\sqrt{x^2 + y^2}$. For gratings, different diffraction orders correspond to plane waves traveling in different directions, but for a lens of focal length f_0 , the m th diffraction order corresponds to a lens of focal length f_0/m . By splitting the light into different diffraction orders, a lens with multiple focal lengths (even of opposite sign if desired) can be designed.

As an example, consider the paraxial design of a bifocal lens, as used in intraocular implants. Half the light should see a lens of focal length f_0 , while the other half should see no lens. This is a lens of focal length f_0 , but with the light split evenly between the 0 and +1 orders. The phase profile of a single focus lens is given by $\phi(r) = -2\pi r^2 / (2\lambda_0 f_0)$, where $r^2 = x^2 + y^2$. This phase, with the 2π ambiguity removed, is plotted in Fig. 7a as a function of r and in Fig. 7b as a function of r^2 , where the periodicity in r^2 is evident. To split the light between the 0 and +1 orders, the blaze of Fig. 7b is replaced by the 1:2 continuous splitter of Table 2, resulting in Fig. 7c. This is the final design and the phase profile is displayed in Fig. 7d as a function of r .

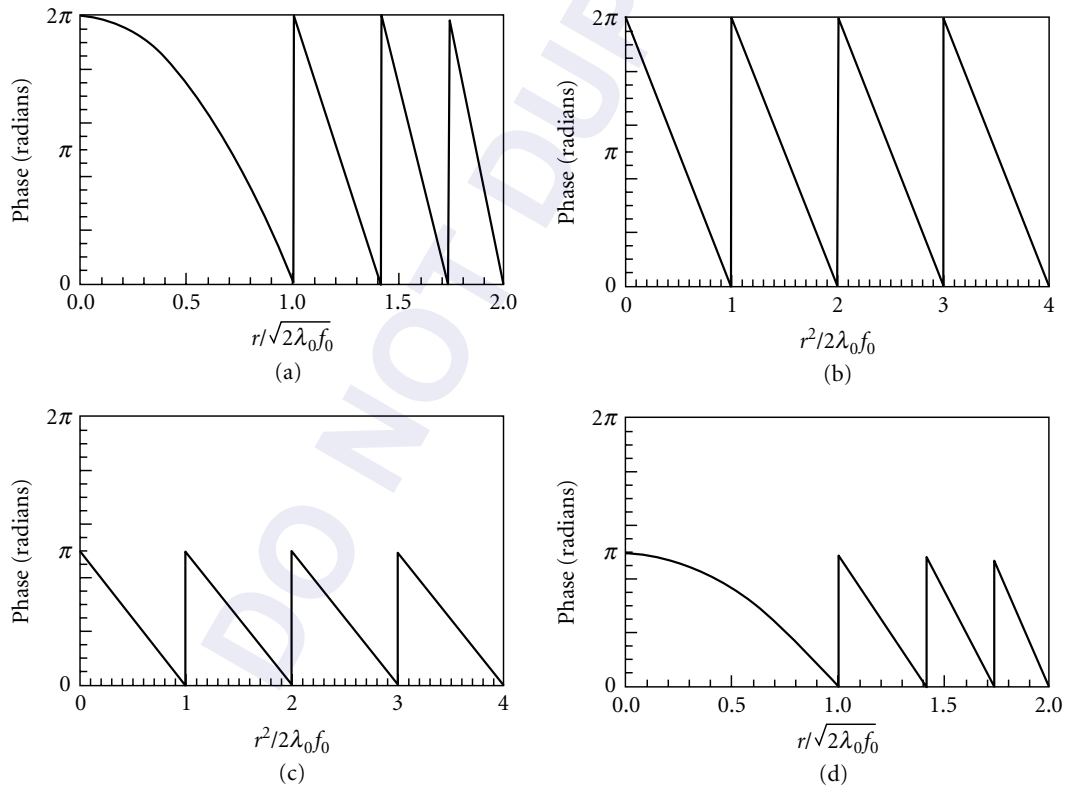


FIGURE 7 Designing a bifocal lens: (a) lens with a single focus; (b) same as (a), but showing periodicity in r^2 ; (c) substitution of a beam-splitting design; and (d) same as (c), but as a function of r .

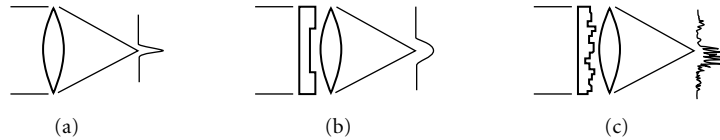


FIGURE 8 Reshaping a focused beam: (a) gaussian focus; (b) deterministic beam-shaper; and (c) statistical diffuser.

Beam Shapers and Diffusers

In many cases, the reshaping of a laser beam can be achieved by introducing the appropriate phase shifts via a binary optics element and then letting diffraction reshape the beam as it propagates. If the incoming beam is well characterized, then it is possible to deterministically design the binary optics element.²⁴ For example, Fig. 8a shows the focal spot of a gaussian beam without any beam-forming optics. In Fig. 8b, a binary optics element flattens and widens the focal spot. In this case, the element could be designed using phase-retrieval techniques, the simplest design being a clear aperture with a π phase shift over a central region. If the beam is not well-behaved, then a statistical design may be more appropriate.²⁵ For example, in Fig. 8c, the aperture is subdivided into randomly phased subapertures. The envelope of the resulting intensity profile is determined by the subaperture but is modulated by the speckle pattern from the random phasing. If there is some randomness in the system (e.g., changing laser wavefront), then the speckle pattern will average out and the result will be a design which reshapes the beam and is robust to variations in beam shape.

Other Devices

Other Fourier optics-based applications which benefit from binary optics include the coupling of laser arrays via filtering in the Fourier plane or other means²⁶ the fabrication of phase-only components for optical correlators,²⁷ and the implementation of coordinate transformations.^{16,28} In all these applications, binary optics is used to directly manipulate the phase of a wavefront.

23.5 DESIGN—VECTOR DIFFRACTION THEORY

Binary optics designs based on vector diffraction theory fall into two categories: grating-based designs and artificial index designs.

Grating-based designs rely on solving Maxwell's equations for diffraction by the element. This is practical for periodic structures. Two major methods for this analysis are the expansion in terms of space harmonics (coupled wave theory) and the expansion in terms of modes (modal theory).²⁹ In this category, optical design is difficult since it can be both nonintuitive and computationally intensive.

Artificial index designs are based on the following premise. When features on the component are small compared to the wavelength, then the binary optics element will behave as a material of some average index. Two common applications are shown in Fig. 9. In Fig. 9a, the device behaves as an antireflection coating (analogous to anechoic chambers) since, at different depths, the structure has a different average index, continuously increasing from n_1 to n_2 . In Fig. 9b, the regular, subwavelength structure exhibits form birefringence.³⁰ For light polarized with the electric vector perpendicular to the grooves, the effective index is

$$\frac{1}{n_{\text{eff}}^2} = p \frac{1}{n_1^2} + (1-p) \frac{1}{n_2^2} \quad (31)$$

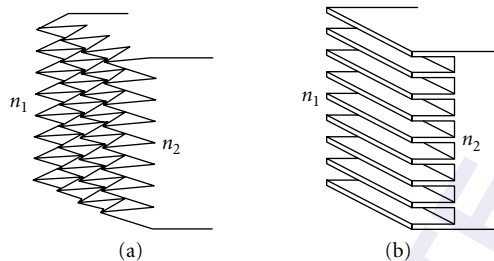


FIGURE 9 Artificial index designs: (a) antireflection layer and (b) form birefringence.

where p is the fraction of total volume filled by material 1. However, for light polarized with the electric vector parallel to the grooves,

$$n_{\text{eff}}^2 = pn_1^2 + (1-p)n_2^2 \quad (32)$$

In both these cases, the period of the structure must be much less than the wavelength in either medium so that only the zero order is propagating.

23.6 FABRICATION

Mask Layout

At the end of the optical design stage, the binary optics element is described by a phase profile $\phi(x, y)$. In the mask layout process, this profile is transformed into a geometrical layout and then converted to a set of data files in a format suitable for electron-beam pattern generation. From these files, a mask maker generates the set of photomasks which are used to fabricate the element.

The first step is to convert the phase profile $\phi(x, y)$ into a thickness profile (see Fig. 10a and b) by the relation

$$t(x, y) = \frac{\lambda_0}{2\pi(n_0 - 1)}(\phi \bmod 2\pi) \quad (33)$$

where λ_0 is the design wavelength and n_0 is the index of the substrate at λ_0 . The thickness profile is the surface relief required to introduce a phase shift of $\phi(x, y)$. The thickness varies continuously from 0 to t_0 , where

$$t_0 = \lambda_0 / (n_0 - 1) \quad (34)$$

is the thickness required to introduce one wave of optical path length difference.

To facilitate fabrication, $t(x, y)$ is approximated by a multilevel profile $t'(x, y)$ (Fig. 10c), which normally would require one processing cycle (photolithography plus etching) to produce each thickness level. However, in binary optics, a binary coding scheme is used so that only N processing cycles are required to produce

$$I = 2^N \quad (35)$$

thickness levels (hence the name binary optics).

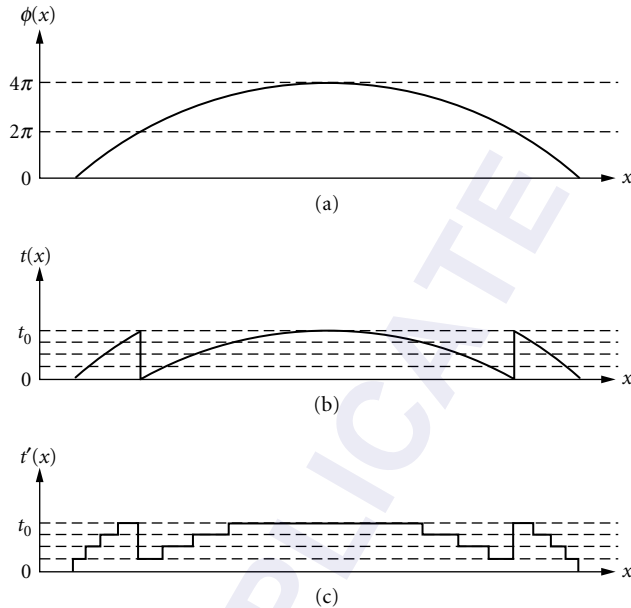


FIGURE 10 Translation from $\phi(x, y)$ to micromachined surface: (a) phase $\phi(x, y)$; (b) thickness $t(x, y)$; and (c) binary optics profile $t'(x, y)$.

The photomasks and etch depths required for each processing cycle are determined from contours of the thickness $t(x, y)$ or equivalently the phase $\phi(x, y)$, as shown in Table 3. The contours can be generated in several ways. For simple phase profiles, the contours are determined analytically. Otherwise, the contours are determined either by calculating the thickness at every point on a grid and then interpolating between points³¹ or by using a numerical contouring method,³² analogous to tracing fringes on an interferogram.

To generate the photomasks, the geometrical areas bounded by the contours must be described in a graphics format compatible with the mask vendor (see Fig. 11a and b). Common formats are GDSII and CIF,³³ both of which are high-level graphics descriptions which use the multisided polygon (often limited to 200 sides) as the basic building block. Hierarchical constructions (defining structures in terms of previously defined structures) and arraying of structures are also allowed.

The photomasks are usually written by electron-beam generators using the MEBES (moving electron beam exposure system) format as input. Most common high-level graphics descriptions can be translated or “fractured” to MEBES with negligible loss in fidelity via existing translation routines. Currently, commercial mask makers can achieve a minimum feature size or “critical dimension” (CD) of $0.8 \mu\text{m}$ with ease, $0.5 \mu\text{m}$ with effort, and $0.3 \mu\text{m}$ in special cases. The CD of a

TABLE 3 Processing Steps for Binary Optics

Layer	Etch Region, Defined by $t(x, y)$	Etch Region, Defined by $\phi(x, y)$	Etch Depth
1	$0 < t \bmod (t_0) < t_0/2$	$0 < \phi \bmod 2\pi < \pi$	$t_0/2$
2	$0 < t \bmod (t_0/2) < t_0/4$	$0 < \phi \bmod \pi < \pi/2$	$t_0/4$
3	$0 < t \bmod (t_0/4) < t_0/8$	$0 < \phi \bmod \pi/2 < \pi/4$	$t_0/8$
4	$0 < t \bmod (t_0/8) < t_0/16$	$0 < \phi \bmod \pi/4 < \pi/8$	$t_0/16$

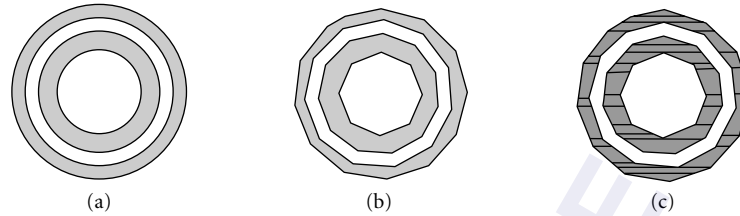


FIGURE 11 Mask layout descriptions: (a) mathematical description based on thickness contours; (b) high-level graphics description; and (c) MEBES.

binary optics element is determined by the minimum local period [see Eq. (8)] divided by the number of steps, D_{\min}/I . For lenses,

$$D_{\min} \doteq 2\lambda_0 F \quad (36)$$

where F is the F -number of the lens; while, for gratings, D_{\min} is the period of the grating.

In MEBES, all geometrical shapes are subdivided into trapezoids whose vertices lie on a fixed rectangular grid determined by the resolution of the electron-beam machine (see Fig. 11c). The resolution (typically $0.05 \mu\text{m}$) should not be confused with the CD achievable by the mask maker.

In summary, the description of the photomask begins as a mathematical description based on contours of the thickness profile and ends as a set of trapezoids whose vertices fall on a regular grid (see Fig. 11). This series of translations results in the following artifacts. First, curves are approximated by straight lines. The error introduced by this approximation (see Fig. 12) is

$$\delta = R(1 - \cos \theta/2) \doteq R\theta^2/8 \quad (37)$$

Normally, the maximum allowable error is matched to the electron-beam resolution. Second, all coordinates are digitized to a regular grid. This results in pixelization artifacts (which are usually negligible), analogous to the ziggurat pattern produced on video monitors when plotting gently sloped lines. Finally, the MEBES writing process itself has a preferred direction since it uses electrostatic beam deflection in one direction and mechanical translation in the other.

In addition to the digitized thickness profile, photomasks normally include the following features which aid in the fabrication process. Alignment marks³⁴ are used to align successive photomasks, control features such as witness boxes allow the measurement of etch depths and feature sizes without probing the actual device, and labels allow the fabricator to easily determine the mask name, orientation, layer, etc.

Micromachining Techniques

Binary optics uses the same fabrication technologies as integrated circuit manufacturing.^{34,35} Specifically, the micromachining of binary optics consists of two steps: replication of the photomasks

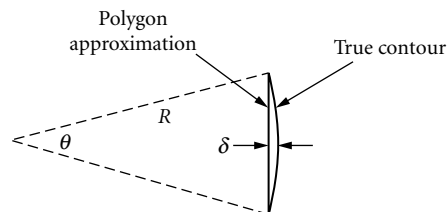


FIGURE 12 Quantization angle.

pattern into photoresist (photolithography) and the subsequent transfer of the pattern into the substrate material to a precise depth (etching or deposition).

The replication of the photomasks onto a photoresist-covered substrate is achieved primarily via contact, proximity, or projection optical lithography. Contact and proximity printing offer lower equipment costs and more flexibility in handling different substrate sizes and substrate materials. In contact printing, the photomask is in direct contact with the photoresist during exposure. Vacuum-contact photolithography, which pulls a vacuum between the mask and photoresist, results in the highest resolution (submicron features) and linewidth fidelity. Proximity printing, which separates the mask and photoresist by 5 to 50 μm , results in lower resolution due to diffraction. Both contact and proximity printing require 1:1 masks. In projection printing, the mask is imaged onto the photoresist with a demagnification from 1 \times to 20 \times . Projection printers are suitable for volume manufacturing and can take advantage of magnified masks. However, they also require expensive optics, strict environmental controls, and can only expose limited areas (typically 2 cm \times 2 cm).

Following exposure, either the exposed photoresist is removed (positive resist) or the unexposed photoresist is removed (negative resist) in a developer solution. The remaining resist serves as a protective mask during the subsequent etching step.

The most pertinent etching methods are reactive ion etching (RIE) and ion milling. In RIE, a plasma containing reactive neutral species, ions, and electrons is formed at the substrate surface. Etching of the surface is achieved through both chemical reaction and mechanical bombardment by particles. The resulting etch is primarily in the vertical direction with little lateral etching (an anisotropic etch) and the chemistry makes the etch attack some materials much more vigorously than others (a selective etch). Because of the chemistry, RIE is material-dependent. For example, RIE can be used to smoothly etch quartz and silicon, but RIE of borosilicate glasses results in micropatterned surfaces due to the impurities in the glass. In ion milling, a stream of inert gas ions (usually Ar) is directed at the substrate surface and removes material by physical sputtering. While ion milling is applicable to any material, it is usually slower than RIE.

For binary optics designed to be blazed for a single order (i.e., designs based on geometrical optics), the major effect of fabrication errors is to decrease the efficiency of the blaze. There is little or no degradation in the wavefront quality. Fabrication errors can be classified as lithographic errors, which include alignment errors and over/underexposure of photoresist, and etching errors, which include depth errors and nonuniform etching of the substrate. As a rule of thumb, lithographic errors should be held to less than 5 percent of the minimum feature size ($<0.05 D_{\min}/l$), which can be quite challenging; while etching errors should be held to less than 5 percent of t_0 , which is usually not too difficult. For binary optics designed via scalar or vector diffraction theory, manufacturing tolerances are estimated on a case-by-case basis through computer simulations.

23.7 REFERENCES

1. "Holographic and Diffractive Lenses and Mirrors," *Proc. Soc. Photo-Opt. Instrum. Eng.* Milestone Series **34**, 1991.
2. "Computer and Optically Generated Holographic Optics," *Proc. Soc. Photo-Opt. Instrum. Eng.* series, **1052**, 1989; **1211**, 1990; **1555**, 1991.
3. "Miniature and Microoptics," *Proc. Soc. Photo-Opt. Instrum. Eng.* series, **1544**, 1991 and **1751**, 1992.
4. G. J. Swanson, "Binary Optics Technology: The Theory and Design of Multi-Level Diffractive Optical Elements," M.I.T. Lincoln Laboratory Technical Report 854, NTIS Publ. AD-A213-404, 1989.
5. M. W. Farn and W. B. Veldkamp, "Binary Optics: Trends and Limitations," *Conference on Binary Optics*, NASA Conference Publication, 3227, 1993, pp. 19-30.
6. S. H. Lee, "Recent Advances in Computer Generated Hologram Applications," *Opt. and Phot. News* **16**(7), 1990, pp. 18-23.
7. W. T. Welford, *Aberrations of Optical Systems*, Adam Hilber Ltd., Boston, 1986, pp. 75-78, 217-225.
8. W. C. Sweatt, "Mathematical Equivalence between a Holographic Optical Element and an Ultra-high Index Lens," *J Opt. Soc. Am.* **69**, 1979, pp. 486-487.

9. M. W. Farn, "Quantitative Comparison of the General Sweatt Model and the Grating Equation," *Appl Opt.* **1992**, pp. 5312–5316.
10. W. C. Sweatt, "Describing Holographic Optical Elements as Lenses," *J. Opt. Soc. Am.* **67**, 1977, pp. 803–808.
11. D. A. Buralli and G. M. Morris, "Design of Diffractive Singlets for Monochromatic Imaging," *Appl. Opt.* **30**, 1991, pp. 2151–2158.
12. D. A. Buralli and J. R. Rogers, "Some Fundamental Limitations of Achromatic Holographic Systems," *J. Opt. Soc. Am.* **A6**, 1989, pp. 1863–1868.
13. C. W. Chen, "Application of Diffractive Optical Elements in Visible and Infrared Optical Systems," *Proc. Soc. Photo-Opt. Instrum. Eng.* **CR41**, 1992, pp. 157–172.
14. T. Stone and N. George, "Hybrid Diffractive-Refractive Lenses and Achromats," *Appl. Opt.* **27**, 1988, pp. 2960–2971.
15. R. Kingslake, *Lens Design Fundamentals*, Academic Press Inc., New York, 1978, pp. 77–78, 112–114.
16. J. R. Leger and W. C. Goltso, "Geometrical Transformation of Linear Diode-laser Arrays for Longitudinal Pumping of Solid-State Lasers," *IEEE J. Quant. Electron* **28**, 1992, pp. 1088–1100.
17. G. J. Swanson, "Binary Optics Technology: Theoretical Limits on the Diffraction Efficiency of Multilevel Diffractive Optical Elements," M.I.T. Lincoln Laboratory Technical Report 914, 1991.
18. H. Dammann, "Spectral Characteristics of Stepped-Phase Gratings," *Optik* **53**, 1979, pp. 409–417.
19. A. Vasara, et al., "Binary Surface-Relief Gratings for Array Illumination in Digital Optics," *Appl. Opt.* **31**, 1992, pp. 3320–3336.
20. U. Krackhardt, et al., "Upper Bound on the Diffraction Efficiency of Phase-Only Farnout Elements," *Appl. Opt.* **31**, 1992, pp. 27–37.
21. D. Prongue, et al., "Optimized Kinoform Structures for Highly Efficient Fan-Out Elements," *Appl. Opt.* **31**, 1992, pp. 5706–5711.
22. U. Killat, G. Rabe, and W. Rave, "Binary Phase Gratings for Star Couplers with High Splitting Ratios," *Fiber and Integrated Optics* **4**, 1982, pp. 159–167.
23. U. Krackhardt, "Binaere Phasengitter als Vielfach-Strahlteiler," *Diplomarbeit, Universitaet Erlangen-Nuernberg*, Erlangen, Germany, 1989.
24. J. Hossfeld, et al., "Rectangular Focus Spots with Uniform Intensity Profile Formed by Computer Generated Holograms," *Proc. Soc. Photo-Opt. Instrum. Eng.* **1574**, 1991, pp. 159–166.
25. C. N. Kurtz, "Transmittance Characteristics of Surface Diffusers and the Design of Nearly Band-Limited Binary Diffusers," *J. Opt. Soc. Am.* **62**, 1972, pp. 982–989.
26. J. R. Leger, et al., "Coherent Laser Beam Addition: An Application of Binary-Optics Technology," *The Lincoln Lab Journal* **1**, 1988, pp. 225–246.
27. M. A. Flavin and J. L. Horner, "Amplitude Encoded Phase-Only Filters," *Appl. Opt.* **28**, 1989, pp. 1692–1696.
28. O. Bryngdahl, "Geometrical Transforms in Optics," *J. Opt. Soc. Am.* **64**, 1974, pp. 1092–1099.
29. T. K. Gaylord, et al., "Analysis and Applications of Optical Diffraction by Gratings," *Proc. IEEE* **73**, 1985, pp. 894–937.
30. D. H. Raguin and G. M. Morris, "Antireflection Structured Surfaces for the Infrared Spectral Region," *Appl. Opt.* **32**, 1993, pp. 1154–1167.
31. J. Logue and M. L. Chisholm, "General Approaches to Mask Design for Binary Optics," *Proc. Soc. Photo-Opt. Instrum. Eng.* **1052**, 1989, pp. 19–24.
32. A. D. Kathman, "Efficient Algorithm for Encoding and Data Fracture of Electron Beam Written Holograms," *Proc. Soc. Photo-Opt. Instrum. Eng.* **1052**, 1989, pp. 47–51.
33. S. M. Rubin, *Computer Aids for VLSI Design*, Addison-Wesley Publishing Co., Reading, MA, 1987.
34. N. G. Einspruch and R. K. Watts (eds.), *Lithography for VLSI*, VLSI Electronics Series 16, Academic Press Inc., Boston, MA, 1987.
35. N. G. Einspruch and D. M. Brown (ed.), *Plasma Processing for VLSI*, VLSI Electronics Series 8, Academic Press Inc., Boston, MA, 1984.

Duncan T. Moore

*The Institute of Optics, and
Gradient Lens Corporation
Rochester, New York*

24.1 GLOSSARY

A	constant
a, b	constants
g	constant
h_i	constants
n	refractive index
r	radius
V_{ij}	Abbe numbers
z	Cartesian coordinate (optical axis direction)
Φ	power

24.2 INTRODUCTION

Gradient index (GRIN) optics¹ refers to the field of optics in which light propagates along a curved path. This contrasts with normal homogeneous materials in which light propagates in a rectilinear fashion. Other terms that have been used to describe this field are inhomogeneous optics, index of refraction gradients, and distributed index of refraction. The most familiar example of a gradient index phenomenon is the mirage when a road appears to be wet on a hot summer day. This can be understood by the fact that the road is absorbing heat, thus slightly raising the temperature of the air relative to the temperature a few meters above the surface. By the gas law, the density decreases, and therefore the index of refraction decreases. Light entering this gradient medium follows a curved path. The ray path, as shown in Fig. 1, is such that the ray propagates downward toward the road and then gradually upward to the observer's eye. The observer sees two images. One is the normal image propagating through the homogeneous material and the second is an image that is inverted and appears below the road surface. Thus, the index of refraction gradient acts as a mirror by gradual light refraction rather than reflection.

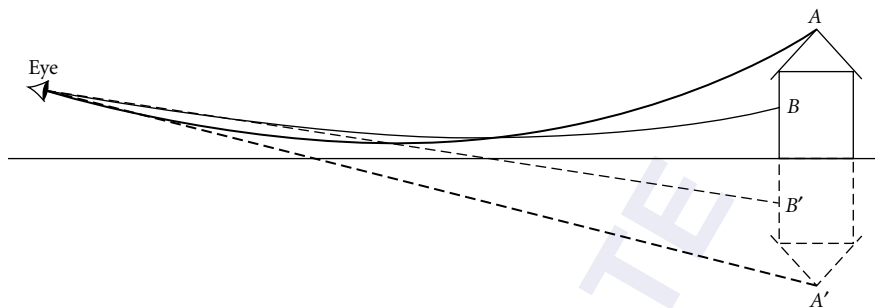


FIGURE 1 Light from point A emits or reflects in all directions. Light propagating several meters above the heated road travels in a straight line. Light passing through the lower index of refraction region near the road undergoes a bending. This light appears to have come from below the road.

24.3 ANALYTIC SOLUTIONS

Over the approximately 150 years that gradient index optics has been studied, a wealth of very interesting analytic solutions has been published.² A classic example was published by James Clerk Maxwell in 1850. Maxwell³ showed through geometrical optics that the ray paths in a spherically symmetric material whose index of refraction is given by

$$n(r) = a/(b^2 + r^2) \quad (1)$$

are circles. The object and the image lie on the surface of the sphere but, otherwise, the imaging is perfect between the conjugate points on the sphere. The medium between the object and the image is continuous with no discrete surfaces. A century later, Luneburg⁴ modified the system to allow for discontinuities of the index of refraction. While these have not been implemented in ordinary optical systems, they have, however, been shown to be useful in integrated optics.⁵

A final example of a numerical solution is that of a radial (cylindrical) gradient in which the index of refraction varies perpendicular to a line. In 1954, Fletcher⁶ showed that if the index of refraction is given by

$$n(r) = n_0 \operatorname{sech}(ar^2) \quad (2)$$

then the ray paths inside the material in the meridional plane are sinusoidal. Nearly 50 years earlier, Wood⁷ had shown experimentally that the paths appeared to be sinusoidal. This solution has several important commercial applications. It is the basis of the Selfoc lens used in arrays for facsimile and photocopying machines and in endoscopes used for medical applications.

24.4 MATHEMATICAL REPRESENTATION

Most of the gradient index profiles are represented by a polynomial expansion. While these expansions are not necessarily the most desirable from the gradient materials manufacturing standpoint, they are convenient for determining the aberrations of systems embodying GRIN materials. There are basically two major representations for gradient index materials. The first, used by the Nippon Sheet Glass, is used exclusively by representing radial gradient components. In this case, the index of refraction is written as a function of the radial coordinate r

$$N(r) = N_0(1 - Ar^2/2 + h_4r^4 + h_6r^6 + \dots) \quad (3)$$

The second method of representing index of refraction profiles is a polynomial expansion in both the radial coordinate r and the optical axis coordinate z . In this case, the representation is

$$N(r, z) = \sum_{j=0} \sum_{i=0} N_{ij} r^{2i} z^j \quad (4)$$

where the coefficients N_{ij} are the coefficients of the index of refraction polynomial. A pure axial gradient (in the z direction) has coefficients only in the form of N_{0j} and those in the radial would be of the form of N_{i0} . These representations for the index of refraction polynomial have been the basis of the aberration theory which was first developed for gradient index materials with discrete surfaces by Sands.⁸ These coefficients are wavelength dependent and are typically defined at three wavelengths. A gradient dispersion is defined by using a general Abbe number

$$V_{ij} = N_{ij,d} / (N_{ij,F} - N_{ij,C}) \quad (5)$$

except for i and j both equal to zero. In the case for $i = j = 0$, then the Abbe number becomes the standard form, namely

$$V_{00} = (N_{00,d} - 1) / (N_{00,F} - N_{00,C}) \quad (6)$$

The subscripts d , F , and C refer to the wavelengths 0.5876, 0.4861, and 0.6563 μm , respectively. Unlike the normal dispersion of glasses where V_{00} is between 20 and 90, the V_{ij} can have negative and positive values or can be infinite (implying that the gradient is the same as both the red and the blue portions of the spectrum).

24.5 AXIAL GRADIENT LENSES

When the index of refraction varies in the direction of the optical axis (the z direction), the bending of the light within the material is very small. Thus, the main feature of an axial gradient is its ability to correct aberrations rather than to add power to the lens. Sands showed that the effect of an axial gradient on monochromatic aberrations is exactly equivalent to that of an aspheric surface. In fact, one could convert any aspherical surface to an axial gradient with a spherical surface and have the same image performance to the third-order approximation. There is, however, one very important difference between aspheric surfaces and axial gradients, i.e., the variation of the index of refraction profile with wavelength. Since an aspheric is the same for all wavelengths, its effect on spherochromatism is established once the aspheric has been determined. Further, an asphere has no effect on paraxial axial or lateral chromatic aberrations. This is not the case for axial gradients. Since the index of refraction profile varies with wavelength, it is possible to significantly modify the spherochromatism of the lens and, in the case where the gradient extends from the front to the back surface, to affect the paraxial chromatic aberrations. Depending upon the dispersion of the gradient index material, the spherochromatism can be increased or decreased independent of the monochromatic correction. The effect of an axial gradient on paraxial axial chromatic aberration is best understood by placing a surface perpendicular to the optical axis in the middle of a single lens dividing it into two parts. The gradient dispersion implies that the medium will have one dispersion at the front surface and a different dispersion at the second surface. Thus, if one were to design a material in which the dispersion of the front surface is 60 and at the rear surface is 40, then the combination of a positive (convex surface on the front) and negative lens (concave surface on the back), reduces the chromatic aberration. This can only be done if the lens is meniscus. In that case, the theoretical front lens is plano-convex while the back one is plano-concave. If the negative element has the higher dispersion (lower V numbers), then it is possible to chromaticize the lens by a proper bending of the lens surfaces. This was first shown in the infrared part of the spectrum using a zinc sulfide-zinc selenide gradient material.⁹

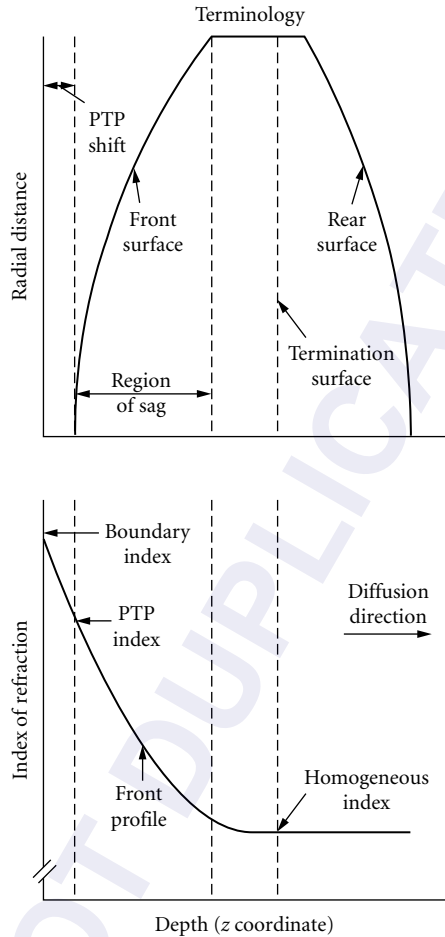


FIGURE 2 Diagram of axial gradient terminology. The effective region of the gradient is in the “region of sag.”

The simplest example of an axial gradient is the linear profile in which the index of refraction is written as

$$N(z) = N_{00} + N_{01}z \tag{7}$$

The coefficient N_{01} is an additional degree of freedom which can be used to correct any of the third-order monochromatic aberrations except Petzval curvature of field. There are two ways to approach the design of these lenses. In the case where the index of refraction profile does not continue to the rear surface (see Fig. 2), a simple formula can be used to relate the amount of index change to the F -number of the lens surface if the third-order spherical aberration and coma are to be correct to zero,¹⁰ namely,

$$\Delta n = (0.0375 / (N_{00} - 1)^2) / f \#^2 \tag{8}$$

in this formula, the important parameters are the index of refraction of the base material, N_{00} , the change in index of index of refraction, Δn , from the polar tangent to the maximum sag point, and

the F -number of the lens. One sees that if the F -number of the lens is doubled, then the amount of index change necessary to correct the spherical aberration and coma to zero increases by a factor of 4. Thus, while it is possible to correct the spherical aberration of the singlet operating at $F/4$ with an index change of only 0.0094, that same lens operating at $F/1$ will require an index change of 0.15. In most lenses, one never corrects the spherical aberration of individual elements to zero, but corrects the total amount of spherical aberration of all lens elements to zero.

Axial gradients have been used in a number of lens designs. Most of the work in this field has occurred in photographic objectives.^{11,12} In these cases, they offer a slight advantage over aspherics because of the chromatic variation of the gradient.

24.6 RADIAL GRADIENTS

In the most generalized case for radial gradients (one in which all coefficients are nonzero), it is possible not only to use the gradient for aberration correction, but also to modify the focal length of the lens. Independent of which representation is used, the coefficient of the parabolic term [Eq. (2) or Eq. (3)] dictates the amount of power that is introduced by the radial gradient component. Assuming only a radial gradient component, Eq. (4) can be expanded as

$$N(r) = N_{00} + N_{01}r^2 + N_{02}r^4 + \dots \quad (9)$$

Equating the terms in Eq. (3) and Eq. (9) gives

$$N_{00} = N_0 \quad \text{and} \quad N_{10} = -N_0 A/2 \quad (10)$$

In the most general form, the power ϕ , due to the radial gradient component, is written as

$$\phi = -N_0 A^{0.5} \sin(A^{0.5} t) \quad (11)$$

From Eq. (11), the length of the material t determines the focal length of the system. In fact, depending on the choice of length, the power can be positive, negative, or zero. See Fig. 3a. A convenient variation

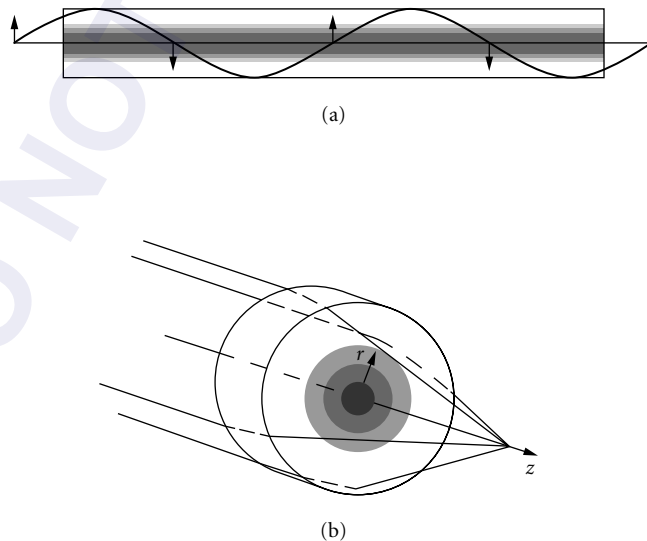


FIGURE 3 Diagram of radial gradient: (a) a long radial gradient lens illustrating period ray path and (b) wood lens.

on this formula is to determine the length at which light entering the material collimated will be focused on the rear surface. This length is called the quarter pitch length of the rod and is given by

$$P_{1/4} = (\pi/2)(-N_{00}/(2N_{10}))^{0.5} \quad (12)$$

The full period length of the rod is simply four times Eq. (12).

For the case where the focal length (the reciprocal of the power) is long compared to its thickness, this can be approximated by the formula (see Fig. 3b)

$$\phi = -2N_{10}t \quad (13)$$

This simplifying formula was derived by the entomologist Exner¹³ in 1889 while he was analyzing insect eyes and found them to have radial gradient components. Since the dispersion of a gradient material can be positive, negative, or infinity, the implication is that the paraxial axial chromatic aberration can be negative, positive, or zero. This leads to the possibility of an achromatized singlet with flat surfaces; or, by combining the dispersion of the gradient with that of the homogeneous materials, to single element lenses with curved surfaces that are color-corrected.

The radial gradient lens with flat surface is a very important example, both from a theoretical and a commercial standpoint. Consider such a lens with an object of infinity where the lens is thin relative to its focal length. As has already been shown, the value of N_{10} and the thickness determine the focal length of such a lens. According to third-order aberration theory,⁸ the only other term that can influence the third-order monochromatic aberrations is the coefficient N_{20} . This term can be used to correct any one of the third-order aberrations except Petzval curvature of field. The coefficient N_{20} is normally used to correct the spherical aberration; however, once this choice is made, there are no other degrees of freedom to reduce other aberrations such as coma. It can be shown that the coma in such a single element lens is very large if the lens is used at infinite conjugates. Of course, if such a lens is used at unit magnification in a system which is symmetric about the aperture stop, the coma (as well as the distortion and paraxial lateral color) is zero. As the length of the rod increases, the approximation for the focal length becomes inaccurate and the more rigorous formula given by Eq. (11) is appropriate. However, the rules governing the aberration correction remain the same. That is, the choice of the value of N_{20} , or in the Nippon sheet glass representation, in h_4 coefficient, corrects the third-order spherical aberration to zero. In Fletcher's original paper, he showed that rays propagating in a material whose index of refraction is given by Eq. (2) would focus light in the meridional plane periodically with no aberration along the length of such a rod. If one expands a hyperbolic secant in a polynomial expansion, one obtains

$$N_{20} = 5N_{10}^2/6N_{00} \quad (14)$$

The implication is that if N_{20} is chosen according to Eq. (14), then not only is the spherical aberration corrected, but so is the tangential field (that is, the sum of three times the astigmatism plus the field curvature). Rawson¹⁴ showed that a more appropriate value for N_{20} was $3N_{10}^2/2N_{00}$. This is a compromise for the correction of sagittal and tangential fields.

The second limiting case is to use these rods with arbitrary length but at unit magnification. This has important commercial applications in photocopying and fax machines, for couplers for single-mode fibers, and in relays used in endoscopes. In all of these systems, the magnification is ± 1 and thus there is no need to correct the coma, the distortion, or the lateral color. Thus, the choice of N_{20} can be used to either correct the spherical aberration or to achieve a compromise between the tangential and sagittal fields.

In one of the most common applications, a series of lenses is assembled to form an array (see Fig. 4). In this case, the magnification between the object and the image must be a +1 with an inverted image halfway through the gradient index rods. Light from an object point is imaged through multiple GRIN rods depending on the numerical aperture of each of the rods. The effective numerical aperture of the array is significantly higher than that of a single rod. Theoretically, a full two-dimensional array can be constructed to image an entire two-dimensional object. In practice, to reduce costs the object is scanned by moving the object across the fixed lens array with either a reduced couple device or a transfer drum used to record the image.

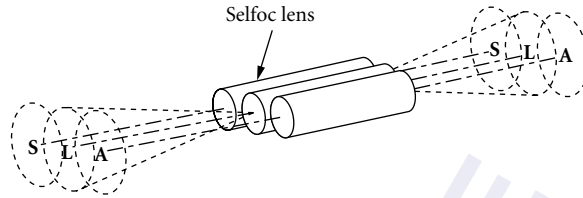


FIGURE 4 An array of radial gradient lenses (only three shown) can be used to form an image of an extended object. This principle is used in photocopying and fax machines.

24.7 RADIAL GRADIENTS WITH CURVED SURFACES

While the radial gradient with flat surfaces offers tremendous commercial applications today, it has limited applications because of the large amount of coma that is introduced unless the lens system is used at unit magnification. Thus, it is often desirable to introduce other degrees of freedom that may improve the imagery. The simplest way to do this is to make one or both of the end cases curved. The ability to chromitize such a lens is not lost so long as the power resulting from the curved surfaces and that of the radial gradient maintain the same ratio (but with opposite sign) as that of the Abbe number of homogeneous material and the Abbe number of N_{10} . Thus, the lens shape can be determined to reduce the coma to zero and the value of the N_{20} coefficient is chosen to eliminate the spherical aberration. An example of a curved lens with a radial gradient was developed by Nippon Sheet Glass for a compact disc player.¹⁵ In that case, it is not necessary to achromatize the lens since the source is a monochromatic laser diode, but it was necessary to extend the field and reduce the amount of spherical aberration simultaneously. It is also often desirable to place part of the power on the curvature rather than using the gradient to refract all of the light. This reduces the magnitude of the index change and makes the lens easier to manufacture.

In a radial gradient material with curved surfaces, it is possible to eliminate four out of five monochromatic aberrations,¹⁶ and any four can be chosen. However, these lenses tend to be very sensitive to slight manufacturing errors, as they require a very delicate balance between the coefficients of the gradient profile and typically have very large amounts of higher-order aberrations.

24.8 SHALLOW RADIAL GRADIENTS

An interesting compromise between an axial gradient and a radial gradient with power is the shallow radial gradient (SRGRIN). In this type of gradient, there is no power generated by the gradient (i.e., $N_{10} = 0$). Like the axial gradient, it has no effect on Petzval curvature of field, but its aberration correction is significantly different than that of axial gradients. Sands⁸ showed that in the case of an axial gradient, the important parameter is the differential refraction of the ray at the surface which causes an additional surface contribution. In the shallow radial gradient there is no surface contribution, since the N_{10} coefficient is zero. All of the aberration correction is from the transfer contribution through the material.⁸ The implication of this fact is that the thickness of the shallow radial gradient is very important and, in fact, the most important parameter is the product of the thickness and the N_{20} coefficient. Thus, if only a small index change can be manufactured, the same amount of aberration correction can be achieved by increasing the thickness of the element. The other significant difference between this gradient and a normal radial gradient is the sign of the index change. In most lenses designed to date, the index of refraction of a conventional radial gradient should be lower at the periphery than it is at the center, thus creating a positive lens. However, in the shallow radial gradient, the index of refraction should be higher at the periphery than at the center. This has also normally

been the case in the axial gradient in which the index of refraction should be higher at the polar tangent plane than at the maximum sag point. This has important implications for the manufacturing process. Furthermore, the amount of index change necessary for shallow gradient correction is usually very small compared to the amount of index change needed in a regular radial gradient.

24.9 MATERIALS

While several materials systems have been proposed for forming gradient index materials, gradients have only been made for commercial applications in glasses and polymers. However, research has been conducted in zinc selenide-zinc sulfide,¹⁷ and germanium-silicon¹⁸ for the infrared portion of the spectrum, and in fluoride materials for the ultraviolet.¹⁹ However, none of these have reached the stage, at this writing, which can be commercialized. For glasses, several processes have been proposed. The most common method of making gradient index materials is by the ion exchange process. In this case, a glass containing a single valence ion (such as sodium, lithium, or potassium) is placed in a molten salt bath at temperatures between 400 and 600°C. The molten salt bath contains a different ion than that in the glass. The ions from the salt diffuse into the glass and exchange for an ion of equal valence in the glass. The variation in composition leads to a variation in index of refraction. The variation in index of refraction occurs due to the change of polarizability between the two ions and the slight change in the density of the material. In some cases, these two phenomena can cancel one another, producing a composition variation, but no corresponding change in index of refraction. A model for predicting the index refraction change as well as the chromatic variation of the gradient has been developed.²⁰ In this system, it is clear that the maximum index change is limited by the changes in the properties of single valence ions. While very large index changes have been made (approaching 0.27), these gradients suffer from large amounts of chromatic aberration. In axial gradients, a large amount of chromatic aberration is desirable, as it normally improves the spherochromatism. In the case of radial gradients, however, it creates large paraxial axial chromatic aberration which is normally not desirable.

The manufacturing method is quite simple. If one wishes to make axial gradients, a sheet of glass is placed in a molten salt bath. Typical times for diffusion are a few days for diffusion depths of 3 to 7 mm at temperatures around 500°C. The higher the temperature, the faster the diffusion; however, at high temperature the glass will begin to deform. Lower temperatures increase the diffusion times. For radial gradients, one simply starts with glass with cylindrical symmetry and places the rods inside an ion exchange bath. In order to form good parabolic profiles, it is necessary for the ions to diffuse through the center.

Two other methods have been proposed for making gradients in glass. In the first, the gradient is formed by leaching or by stuffing in a sol-gel formed glass. This system has only shown to be applicable to radial gradients. After the glass is formed by the sol-gel (solution gelatin process), the glass is in a porous state where one of the components can be dissolved out in an acid bath²¹ or molecules can be stuffed into the glass to form the index of refraction gradient.²² By the leaching method, gradients have been formed in either titanium or zirconium. Index changes of up to 0.03 have been formed by this method. Alternatively, the glass can be stuffed with ions such as lead. The lead precipitates on the walls of the porous material whereupon it is included in the glass during the sintering step. While it is possible to get much larger changes using the method based on lead, both of these techniques suffer from large amounts of chromatic aberration.

A new method shown to be very useful for axial gradients is based on the fusion of glass slabs.²³ The index of refraction of each slab is slightly different than its adjacent slab. Very large index of refraction changes can be formed by this technique ($\Delta n = 0.4$). Further, these materials can be made in apertures up to 100 mm.

Two basic methods for manufacturing of polymers for gradient index have been demonstrated. In the first, an exchange of one monomer for a monomer in a partially polymerized material forms a profile in the same way as the ion exchange method.²⁴ In the second, ultraviolet light is used to induce photocopolymerization to form an index of refraction in the material.²⁵

24.10 REFERENCES

1. For a source of over 100 articles on Gradient-Index (GRIN) Optics, the reader is referred to a series of special issues in *Applied Optics*, GRIN I (April 1, 1980), GRIN II (March 15, 1982), GRIN III (Feb. 1, 1983), GRIN IV (June 1, 1984), GRIN V (December 15, 1985), GRIN VI (October 1, 1986), GRIN VII (February 1, 1988), GRIN VIII (October 1, 1990 and December 1, 1990), and GRIN IX (September 1, 1992).
2. S. Cornbleet, *Microwave Optics*, Academic Press, New York, 1976, pp. 108–187.
3. J. C. Maxwell, *The Scientific Papers of James Clerk Maxwell*, W. D. Niven (ed.), New York, 1965, pp. 76–78.
4. R. K. Luneburg, *Mathematical Theory of Optics*, University of California Press, 1966, pp. 182–195.
5. W. H. Southwell, “Planar Optical Waveguide Lens Design,” *Appl. Opt.* **21**, 1982, p. 1985.
6. A. Fletcher, T. Murphy, and A. Young, “Solutions of Two Optical Problems,” *Proc. R. Soc. London A* **223**, 1954, pp. 216–225.
7. R. W. Wood, *Physical Optics*, Macmillan, New York, 1905, pp. 86–91.
8. P. J. Sands, “Third-Order Aberrations of Inhomogeneous Lenses,” *J. Opt. Soc. Am.* **60**, 1970, pp. 1436–1443.
9. J. W. Howard and D. P. Ryan-Howard, “Optical Design of Thermal Imaging Systems Utilizing Gradient-Index Optical Materials,” *Opt. Eng.* **24**, 1985, p. 263.
10. D. S. Kindred, “Development of New Gradient Index Glasses for Optical Imaging Systems,” Ph.D. thesis, University of Rochester, 1990, pp. 207–210.
11. D. S. Kindred and D. T. Moore, “Design, Fabrication, and Testing of a Gradient-Index Binocular Objective,” *Appl. Opt.* **27**, 1988, pp. 492–495.
12. L. G. Atkinson III, et al., “Design of a Gradient-Index Photographic Objective,” *Appl. Opt.* **21**, 1984 p. 1735.
13. S. Exner, “The Retinal Image of Insect Eyes,” (in Ger.), *Sb. Akad. Wiis Wien* **98**, 1889, p. 13.
14. E. G. Rawson, D. R. Herriott, and J. McKenna, “Analysis of Refractive Index Distributions in Cylindrical Graded-Index Glass Rods Used as Image Relays,” *Appl. Opt.* **9**, pp. 753–759.
15. H. Nishi, H. Ichikawa, M. Toyama, and I. Kitano, “Gradient-Index Objective for the Compact Disk System,” *Appl. Opt.* **25**, 1986, p. 3340.
16. D. T. Moore and R. T. Salvage, “Radial Gradient-Index Lenses with Zero Petzval Aberration,” *Appl. Opt.* **19**, 1980, pp. 1081–1086.
17. M. A. Pickering, R. L. Taylor, and D. T. Moore, “Gradient Infrared Optical Material Chemical Vapor Deposition Process,” *Appl. Opt.* **25**, 1986, pp. 3364–3372.
18. J. J. Miceli, “Infrared Gradient-Index Optics: Materials, Fabrication and Testing,” Ph.D. thesis, University of Rochester, New York, 1982.
19. M. T. Houk, “Fabrication and Testing of Index Gradients in Fluoride Materials,” Ph.D. thesis, University of Rochester, New York, 1990.
20. S. D. Fantone, “Refractive Index and Spectral Models for Gradient-Index Materials,” *Appl. Opt.* **22**, 1983, pp. 432–440.
21. T. M. Che, J. B. Caldwell, and R. M. Mininni, “Sol-gel Derived Gradient-Index Optical Materials,” *SPIE* **1328**, 1990, p. 145.
22. M. Yamane, H. Kawazoe, A. Yasumori, and T. Takahashi, “Gradient-Index Glass Rods of PbO-K₂O-B₂O₃-SiO₂ System Prepared by the Sol-Gel Process,” *J. Non-Cryst. Solids* **100** (1–3), 1988, pp. 506–510.
23. J. J. Hagerty, “Glass Plate Fusion for Macro-Gradient-Index Materials,” U.S. Patent 4,929,065, 1990.
24. Y. Ohtsuka and T. Sugano, “GRIN Rod of CR 39-Trifluoroethyl Methacrylate Copolymer by Vapor-Phase Transfer Process,” *Appl. Opt.* **22**, 1983, pp. 413–417.
25. Y. Koike and Y. Ohtsuka, “Studies on the Light Focusing Plastic Rods: Control of Refractive-Index Distribution of Plastic Radial Gradient-Index Rod by Photocopolymerization,” *Appl. Opt.* **24**, 1985, pp. 4316–4320.

This page intentionally left blank.

DO NOT DUPLICATE

PART

5

INSTRUMENTS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

Norman Goldberg

Madison, Wisconsin

25.1 GLOSSARY

AF	auto focus
CCD	charge coupled device
CRT	cathode ray tube
HDTV	high definition television
IR	infrared
IREC	infrared emitting diode
FP	film plane
LCD	liquid crystal display
LP	line pairs
SLR	single lens reflex
MT	mount

25.2 INTRODUCTION

Thanks to technical progress and vigorous competition, the camera buyer faces a difficult challenge in making a choice. This chapter will attempt to reduce the difficulty by asking the buyer to consider the final image; its purpose, its audience, and its appearance.

Next, some of the more recent technical features are discussed. These include the intriguing ability to select objects in a scene for focus and/or exposure measurement by tracking the position of the user's eye. Finally, various types of cameras and their accessories are described.

In terms of technical sophistication, a moderately priced 35-mm snapshot camera made today would astonish a photographer who was suddenly time shifted from the 1950s. Consider the automation of exposure, focus, film loading, winding, rewinding, plus flash exposures from a tiny integral electronic flash unit no bigger than a spare roll of film.

The net result, for the snapshotter, is a higher percentage of "good" pictures per roll of film than ever before. The specialist also profits, particularly when the basis and limits of the feature are understood.

A good share of these technical features have been incorporated in the more advanced cameras; sometimes just because it can be done. Looking beyond this, the most basic technical camera ever made, the view camera, remains virtually unchanged for the past century. It is to photography what the wooden match is to fire making.

Portions of this chapter are adapted from the author's book, *Camera Technology: The Dark Side of The Lens* (Academic Press, 1992). The author acknowledges, with thanks, the permission granted by Academic Press to use certain material from that book in this chapter.

25.3 BACKGROUND

Imagine the first camera as nothing more than a tent with a small hole in the side casting an image upon the opposite wall. From this accidental version of a “pinhole” camera to today’s “smart” cameras, we find a cornucopia of ingenuity embracing optics, mechanics, electronics, and chemistry.

The variety of cameras ranges from one tiny enough to be concealed in a man’s ring to one large enough for several people to walk around it without obscuring the image. The price range of cameras stretches from few dollars for a disposable model (complete with film) to several thousand dollars (without film).

Cameras have recorded images of the deepest ocean trenches and the surface features of Jupiter’s moons. There are cameras that can freeze a bullet in midair or compress the germination of an acorn into a few minutes. From intimate portraits of bacteria to a 360° panoramic view of the Grand Canyon, there’s a camera for any task.

Nonetheless, there is a common denominator: all cameras produce an image. This image may be the end product, or it may be converted in some way to the final image intended for viewing, as shown in Fig. 1. To choose the best camera for a given task, the properties of this final image should be determined first.

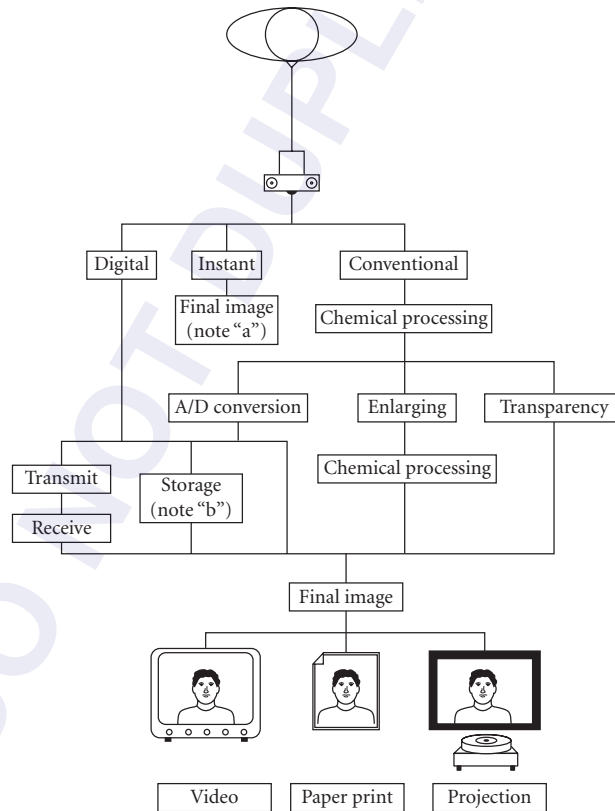


FIGURE 1 Final image flow chart. (a) Many instant photos can be manipulated just as the digital and conventional types, but are treated here in their primary use. (b) Storage means include magnetic tape and disks, optical disks, etc.

25.4 PROPERTIES OF THE FINAL IMAGE

1. Appearance
 - a.* Black-and-white
 - b.* Color
 - c.* High contrast
 - d.* Continuous tone
2. Smallest detail to be resolved
3. Type of display
 - a.* Audience population
 - b.* Viewing conditions
 - (1) Viewing distance
 - (a)* Minimum
 - (b)* Maximum
 - (2) Ambient illumination
 - c.* Display choices
 - (1) Print
 - (2) Projection
 - (3) Self-luminous
4. Distribution

By considering the properties listed, we're obliged to visualize the final image through the viewer's eyes. Esthetics aside, we'll assume that the prime purpose of the final image is to convey information to the viewer.

25.5 FILM CHOICE

The appearance of the final image affects the choice of a camera by the kind of film required to produce that appearance. There are some films that are not available in all sizes. Other films are available in certain sizes only by special order. The availability of some films in some sizes changes over time, so check with your supplier before you select a camera for which film may be scarce.

Most film makers will be glad to send you their latest data on their current films, but be prepared for changes, because this is a very competitive field. New 35-mm color films in particular seem to come out with every change in the seasons.

25.6 RESOLVING FINE DETAIL

If the information in the final image is to be of any use, it must be legible to its detector, which we'll assume to be the human eye. Figure 2 shows that for high-contrast detail viewed under at least 50 foot-candles (office lighting), the eye has an angular resolution of about 1 minute of arc. This means that we can resolve about seven line-pairs per millimeter (LP/mm) at a distance of 250 mm. Since most photographic images exhibit moderate contrast and are viewed in moderate light, a more conservative limit of resolution would be 3.4 minutes of arc, which is good enough to resolve a pattern of 2 LP/mm at 250 mm.

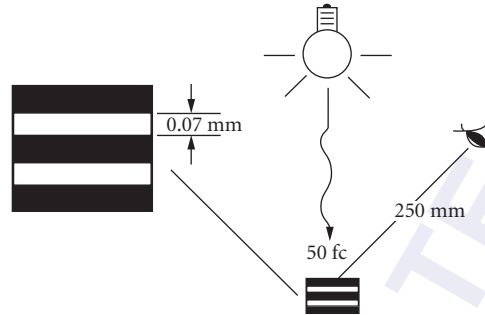


FIGURE 2 Visual resolution. Under ideal viewing conditions, we can resolve seven line-pairs per millimeter.

In most cases, the final image is a magnification of the primary image formed in the camera. All else being equal, there is a practical limit to the extent of this magnification, after which the structure of the film, residual lens aberrations, focus inaccuracy, and/or diffraction effects begin to obscure fine image details.

Suppose then, that for some film we set a practical limit of magnification at 10 \times . Based on the visual resolution limit given previously, the smallest detail in the primary image could be 20 LP/mm, each line 0.025 mm wide.

Looking at it another way, if you want to photograph fine details and display the image legibly at a distance of 250 mm from the viewer, choose a film that will clearly resolve at least 20 LP/mm and is capable of being enlarged 10 diameters without its grain or other structure obscuring the image. Most films in common use today easily satisfy this criterion.

25.7 FILM SIZES

In terms of the widest variety of films available, 35-mm ranks number one. The most common format for this film is 24 \times 36 mm. Although seldom used today, other 35-mm formats include 18 \times 24 mm and 24 \times 24 mm.

Next in line for a broad choice of film types is known as *medium-format* and is sold in 61.5-mm-wide rolls. The shortest rolls are paper-backed and are called 120. Many cameras that accept 120 film will also accept 220 film, which has an opaque paper leader and trailer, but no paper backing over the film. This permits a longer strip of film (more exposures per roll) and better film flatness than 120.

Common formats include (nominal dimensions) 45 \times 60 mm, 60 \times 60 mm, 60 \times 70 mm, and 60 \times 90 mm. Some medium-format cameras also accept 70-mm film that has a row of sprocket holes along each edge and may be loaded in special cassettes for use in the camera's large capacity, motorized, interchangeable film magazine.

The large formats, commonly referred to by their sheet film sizes in inches include 4 \times 5, 5 \times 7, and 8 \times 10, to name the most well known. They may not offer as broad a choice of film as the smaller formats, but the most essential films are available for them.

25.8 DISPLAY

Choosing the best type of display for the final image should start with the number of people in the viewing audience. For large groups, a projected transparency has the advantage of being visible to the entire audience simultaneously. This is especially important if you want to use a pointer to single out detail in the image. Image detail should be clearly resolved by everyone in the audience, from

the front row (image not too grainy) to the last row (image detail within the visual limits). In some cases, the best display is both a projected transparency that the lecturer can refer to with a pointer and a print for each viewer to examine closely, regardless of his position in the audience.

For the best viewing of projected images, the only light striking the screen should be that coming through the transparency. In other words, the room should be pitch black. Unless this condition is met, there is no possibility of reproducing the full tonal range, from deepest black to sparkling white, that the image could contain.

When this condition is difficult to satisfy, a self-luminous display may be best. One or more video monitors located at strategic points can provide good image contrast even under office illumination. The type of monitor may be the conventional cathode ray tube (CRT) or liquid crystal display (LCD). Of the two, the CRT produces a brighter image and is the least expensive. But it is bulky and fragile. The LCD has the virtue of minimal thickness; it's a flat screen display that can be hung on a wall like a framed photo.

At present, neither type can equal the fine detail and subtle color reproduction of a high-grade projected transparency viewed under the proper conditions. However, the gap in image quality is closing, especially now that high-definition television (HDTV) are widely available.

25.9 DISTRIBUTING THE IMAGE

For many applications, the ease and speed with which an image can be distributed is crucial. Thanks to scanners, fax machines, modems, color photocopy machines, rapid photofinishing plants, self-processing "instant" films, etc., we can send practically any image to practically anyone who wants it in a matter of minutes. At present, our ability to do this depends on transforming the analog information in the subject into digital information for transmission, reception, manipulation, analysis, storage, and/or display, as indicated in Fig. 1.

25.10 VIDEO CAMERAS

If speed of acquisition and distribution is most important, we can capture the image on the charge coupled device (CCD) of the widely available camcorder, whose video and audio output signals are available in real time. These video cameras are versatile and moderately priced.

The still-picture counterpart to the camcorder has come to a fork in the road. One path goes to a complete camera system, designed from scratch around the CCD chip and incorporating a miniature magnetic disk drive. The second path leads to a special video back, designed to replace the standard back of a conventional (film) camera. The video back contains a CCD chip and associated circuitry. In some cases the video back and the recorder, in which hundreds of images can be stored, require an "umbilical" cord between them. Some of the newer designs have integrated the back and recorder into a single (cordless) unit. Some of these backs can store up to 50 images internally.

As the capacity for image storage and/or manipulation grows, we see the emergence of systems within systems, where black box A converts black box B to communicate with computer C as long as you have the right adapter cables D, E, and F. This is typical of many rapidly expanding technologies.

Users of a video back on a conventional film camera will notice an unusually narrow angle of view for the lens in use if the light-sensitive area of the CCD chip is smaller than that of the film normally used in the camera. The reason is that only the central region of the camera's format is used. The result is that the camera's lenses perform as though their focal lengths have been "stretched" compared to their performance with conventional film that covers the whole format.

For example, Kodak's DCS 200 replaces the back of an unmodified 35-mm SLR camera, the Nikon N8008s. The "normal" lens for this camera's 24×36 mm film format has a 50-mm focal length, producing a (diagonal) angle of view of about 47°. The same lens used with the video back produces an angle of view of 37° because the CCD measures only 9.3×14 mm. To duplicate the 47° angle of view for this size CCD, a 19.3-mm focal length lens should be used.

Concerning the resolution from CCD images, Kodak's data for the DCS 200 gives a count of 1.54 million (square) pixels, arranged in a 1012×1524 pixel array that measures 9.3×14 mm. This gives a pixel spacing of 0.018 mm, which theoretically can resolve 54.4 monochromatic LP/mm. The color version uses a checkerboard pattern of red, green, and blue filters over the array, so divide the monochrome figure by three to come up with a color resolution of 18.1 LP/mm.

This is quite close to the criterion, discussed earlier, of 2 LP/mm for a 10× enlargement viewed at 250 mm. A 10× enlargement of the CCD image just described would measure 93×140 mm, about the size of a typical snapshot.

25.11 INSTANT PICTURES

For many applications, instant, self-processing film is the best choice. A familiar example is the oscilloscope camera loaded with high-speed film. With minimum, moderately priced equipment, a transient waveform on the scope screen can be captured on the film. Seconds later the print can be examined.

Polaroid dominates this field, which they spawned in 1948. Their range of camera models goes from snapshot to trucksize. They also have special backs which can be used on various cameras to adapt them for use with Polaroid films.

These films range from 35-mm color transparency to 8×10 in (and larger) color print. Included in this variety are black-and-white sheet films that yield both a positive print and a negative. The negative must be stabilized, then washed and dried before being placed in an enlarger or contact printer.

25.12 CRITICAL FEATURES

In many cases, the availability of an accessory such as a Polaroid and/or digital image back is important enough to dictate the choice of a camera. Other factors that may tip the scales in favor of one camera over another might not be discovered until the chosen camera is used for some time.

For example, it may be very useful to have the kind of exposure automation that measures the light reflected from the film plane, before and during the exposure, thus being capable of responding instantly to any change in the scene luminance. There are some cameras that have this capability, yet they lack another feature that may be more valuable for some kinds of photography: the ability to observe the image through the viewfinder of an SLR not just before, but during the exposure.

Most SLRs employ a mirror that swings out of the way just before the exposure begins. This allows the image-forming light to reach the film, but it also blacks out the viewfinder, so that during the crucial instant of the exposure the photographer is momentarily blind. Figure 3 illustrates that by using a beam splitter instead of a conventional mirror in an SLR, the problem is eliminated.

When the advantages of a beam-splitting system are considered, it seems strange that the feature isn't used more widely. Eliminating the swinging mirror reduces the noise and vibration generated each time an exposure is made. This can be crucial when the camera is attached to a microscope or telescope. Some SLRs provide for the mirror to be locked in its raised (shooting) position when desired.

25.13 TIME LAG

Even more important for some types of photography, substituting a beam splitter for a moving mirror in an SLR should reduce the camera's time lag. This is the interval between pressing the camera's trip button and the beginning of the exposure. It's a characteristic shared by all cameras and is rarely mentioned in a manufacturer's specifications of a camera. With few exceptions, time lag has increased in step with camera automation.

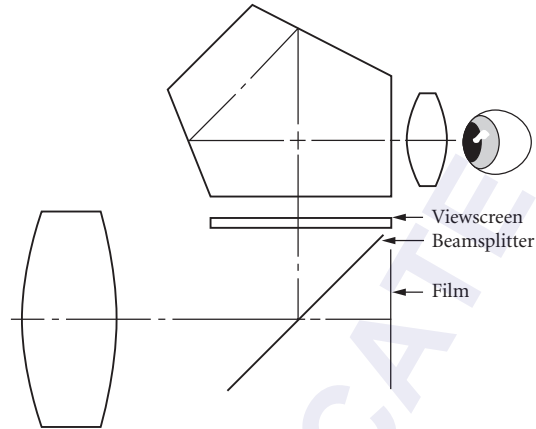
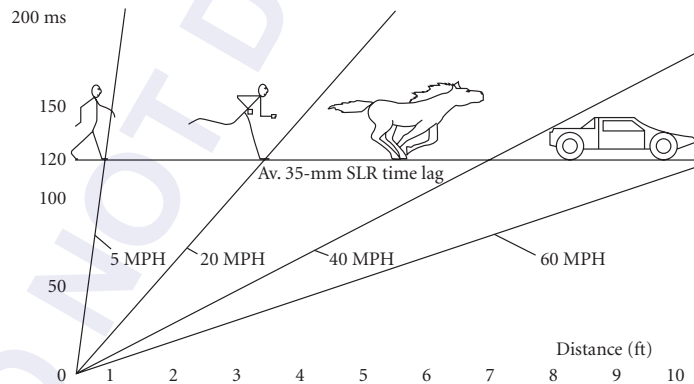


FIGURE 3 Beam-splitter SLR. The beam splitter eliminates the moving mirror, resulting in shorter time lag, reduced noise and vibration, plus the ability to monitor exposure and other image properties in real time.

Testing 40 different 35-mm SLRs for their time lag resulted in a broad range, with the minimum of 46 ms and the maximum of 230 ms. The average was 120 ms. Figure 4 shows that during this interval, a walker moves about 0.8 ft, a runner about twice as far, a galloping horse about 7.0 ft, and a car going 60 mph moves 10.6 ft.



		Distance traveled (ft)									
Time ms		20	40	60	80	100	120	140	160	180	200
Walker	5 MPH	0.1	0.3	0.4	0.6	0.7	0.8	1.0	1.1	1.3	1.4
Runner	20	0.6	1.2	1.7	2.3	2.9	3.5	4.1	4.6	5.2	5.8
Horse	40	1.2	2.4	3.5	4.6	5.8	7.0	8.1	9.3	10.4	11.6
Car	60	1.8	3.5	5.3	7.0	8.8	10.6	12.3	14.1	15.4	17.6
Car	80	2.3	4.7	7.0	9.4	11.7	14.0	16.4	18.7	21.1	23.4

FIGURE 4 Time and motion.

Various other cameras were also tested for their time lag, with these results:

Minox 35 EL (35-mm ultracompact): 8 ms

Leica M3 (35-mm coupled range finder classic): 17 ms

Hasselblad 500C (6×6 cm SLR classic): 82 ms

Kodak Disk 4000 (subminiature snapshot): 270 ms

Polaroid SX-70 Sonar (autofocus instant SLR): 600 ms

25.14 AUTOMATION

Camera automation has taken full advantage of the miniaturization and economy of electronic devices, making two features, autoexposure and autofocus, available in all but the least expensive cameras. This increases the percentage of (technically) good photos per roll of film exposed by the typical amateur.

It's the amateur photographer that is first served when it comes to most of the significant camera automation features. Curious as this may seem, camera makers prefer to introduce a new concept by offering it first in a model intended for the casual snapshotter. This generally means large numbers will be produced. If problems with the feature show up, improvements are made and a "new, improved" model follows. Typically, the feature will be scoffed at by the more seasoned photographer who has learned to overcome the difficulties of making a technically good photograph with the most basic equipment. In time, the new feature is mature enough to be included in the camera maker's premier model. Eventually, even those that scoffed at the feature in its infancy learn to love it, but only after they discover how to recognize and compensate for its weaknesses, if any.

Autoexposure

Early autoexposure systems measured the average luminance of a scene with a selenium photocell, then regulated the shutter speed and/or f-stop based on the deflection of a galvanometer connected to the photocell. These were known as trapped needle systems and were successful in their prime mission: to produce acceptable exposures in snapshot cameras with the just-available color films, whose exposure error tolerance is much smaller than that of black-and-white film.

Most of the first generation autoexposure cameras using the trapped needle system relied on brute force, requiring a long, hard push to trip the camera. This caused camera motion, resulting in a (correctly exposed) smeared image. Nonetheless, many resourceful photographers used these early autoexposure cameras, bolted together with an intervalometer and electromagnetic tripping system, to create an unmanned camera for surveillance, traffic studies, etc.

Amateur movie cameras eagerly adopted autoexposure systems, which proved to be at least as much, if not more, of an improvement for them as they were in still cameras. The movie camera autoexposure systems work by regulating the lens opening (the f-stop), either with a galvanometer or a servomotor. With autoexposure, the movie maker can follow the subject as it moves from bright sunshine to deep shade without the distraction of manually adjusting the f-stop.

This same freedom to follow action without the distraction of manually resetting camera and/or lens controls explains the need for autofocus, a feature whose introduction enjoyed greater enthusiasm from amateur movie makers than from still photographers. Once again, the amateur models were the first to incorporate the feature, but in far less time than it took for autoexposure's acceptance, autofocus became a standard feature in both the amateur and front-line models from most of the makers of 35-mm cameras.

There are similarities between the automation of exposure and focusing. Both have become increasingly sophisticated as user expectations increase. Paradoxically, in the effort to perfect the making of a routine snapshot, some of the more sophisticated automation intrudes on the process by offering the

user certain choices. Instead of simplifying photography, these technological marvels require the user to select a mode of operation from several available modes. For example, many cameras with autoexposure offer factory-programmed combinations of shutter speed and f-stop that favor

- *Action:* fast shutter speed, wide f-stop
- *Maximum depth of field:* small f-stop, slow shutter speed
- *Average scenes:* midway between the first two
- *Fill-flash:* to illuminate portraits made against the light (backlit)

It comes down to this: if you know enough about photographic principles to choose the best autoexposure program, you will rarely need any of them. But when an unexpected change in the subject occurs, such as a cloud moving across the sun, some form of autoexposure can be valuable.

One of the more helpful refinements of autoexposure is the automatic shift of shutter speed with the focal length setting of a zoom lens. This is based on the time-honored guide that gives the slowest shutter speed that may be used without objectionable image motion from normal body tremor. The rule of thumb is to use the shutter speed given by the reciprocal of the lens's focal length. For example, if you're using a 35- to 105-mm zoom lens, the slowest shutter speed for arresting body tremor will shift as you zoom, from 1/35 to 1/105 s (nominal). If the focal length's reciprocal doesn't coincide with a marked shutter speed, use the next faster speed. This guide applies to a handheld camera, not for a camera mounted on a tripod.

Another autoexposure refinement combines a segmented silicon or gallium photocell with a microprocessor to automatically select the best exposure based on the distribution of light reflected from the subject. It amounts to making a series of narrow-angle "spot" readings of the subject, then assigning weighting factors to the different readings according to their relative importance. The weighting factors are determined by the camera maker based on the analysis of thousands of photographs.

Reduced to its most spartan form, a segmented photocell could have a very small central region, surrounded by a broad field. The user can flip a switch to select the desired reading—the center segment for spot readings, the broad segment for full field readings, or both segments for center-weighted full field readings.

To ensure optimum exposure for a subject, seasoned photographers "bracket" exposure settings by making at least three exposures of the subject. The first exposure obeys the meter's reading. The next two are one exposure step less and one greater than the first. This exposure bracketing, with some variations, has been incorporated as an on-demand automatic feature in some cameras.

Autofocus

Autofocus (AF), in one form or another, has become a standard feature in camcorders and in most 35-mm cameras. The latter can be divided into two main types: (1) the snapshot "point-and-shoot," also known as "PHD" (press here, dummy) and (2) the SLR, spanning a wide range in price and sophistication. In between, there are several models which can be thought of as "PHDs on steroids." They have zoom lenses and elaborate viewfinders, making them too bulky to fit easily into a shirt pocket.

There are two main types of autofocus systems, the active and the passive. The active type emits a signal toward the subject and determines the subject's distance by measuring some property of the reflected signal. The passive type measures subject distance by analyzing the subject's image.

Active Autofocus Systems Nearly every active system uses two windows, spaced some distance apart. The user centers the subject in the viewfinder's aiming circle and presses the shutter trip button. Figure 5 shows how a narrow infrared (IR) beam is projected from one of the windows, strikes the subject, and is reflected back to the second window. A photocell behind this window detects the reflected beam. The photocell is sensitive to the position of the beam on its surface and relays this information to its associated circuitry to regulate the camera's focus setting.

Initially, this was a straight-forward triangulation system, using a single infrared beam. But too many users were getting out-of-focus pictures of the main subject when it wasn't in the center of

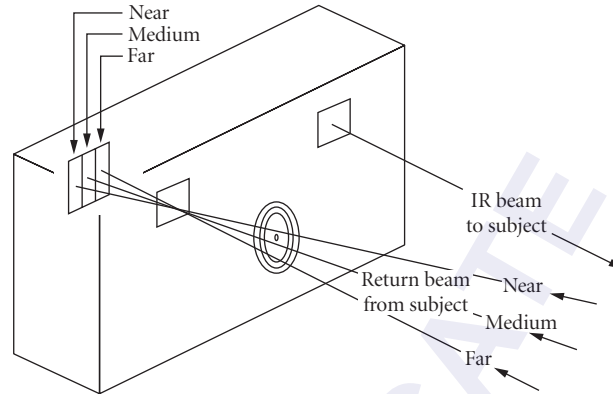


FIGURE 5 Active autofocus. Subject distance determines the angle of the reflected IR beam. The segmented photocell detects this angle, the AF system translates the angle to distance, moves the lens accordingly.

the picture. The camera's instruction book gives the solution: center the main subject in the finder's aiming circle, press the trip button halfway down, and hold it there, then recompose the scene and press the trip button all the way to make the exposure. This requires a fair amount of concentration and discipline, so it contradicted the purpose of having an automatic camera—to be free of cumbersome details, relying on the camera to make properly exposed, sharp photos.

A big improvement was made by projecting three beams from the camera, instead of one. The beams are divergent and the center beam coincides with the finder's aiming circle. Focus is set on the object closest to the camera.

A very different type of active autofocus is the ultrasonic system used by Polaroid in several models. Basically, it's a time-of-flight device that's been compared to sonar and bats. It uses an electrostatic transducer to emit an ultrasonic "chirp" toward the subject. Based on a round-trip travel time of about 5.9 ms/m, the time it takes for the chirp to reach the subject and be reflected back to the camera is translated into subject distance and a servomotor sets the focus accordingly.

A significant advantage of the active autofocus systems just described is their ability to work in total darkness. On the minus side is their inability to focus through a pane of glass or on a subject with an oblique glossy surface that reflects the signal away from the camera.

Passive Autofocus Systems Passive autofocus systems can be broadly characterized as acquiring two views of the subject, each view coming from a slightly different position, then focusing the lens to make the two views match. In this sense, the system operates just like a coincidence-type of optical range finder, but there are important differences.

With an optical range finder we rely on our ability to see when the two images are perfectly superimposed, so our focusing accuracy depends on our visual acuity. In a passive autofocus system, we relieve our eye of this burden and let the tireless electro-optical technology take over.

For the point-and-shoot camera, a passive autofocus system uses two windows, one whose line of sight coincides with that of the viewfinder's, and a second window, spaced some distance from the first. A simple, symmetrical optical system behind the windows includes a CCD for each window. The signal from the first CCD is taken as the reference against which the second CCD's signal is compared. Differences in the light distribution and/or differences in the relative location of the waveforms causes the control circuit to change the focus setting.

Autofocus SLRs Instead of the two windows just described, autofocus SLRs use two bean-shaped segments on opposite sides of the camera lens's exit pupil. Figure 6 shows how this is done. Two small lenslets are located a short distance behind the geometric equivalent of the camera's film

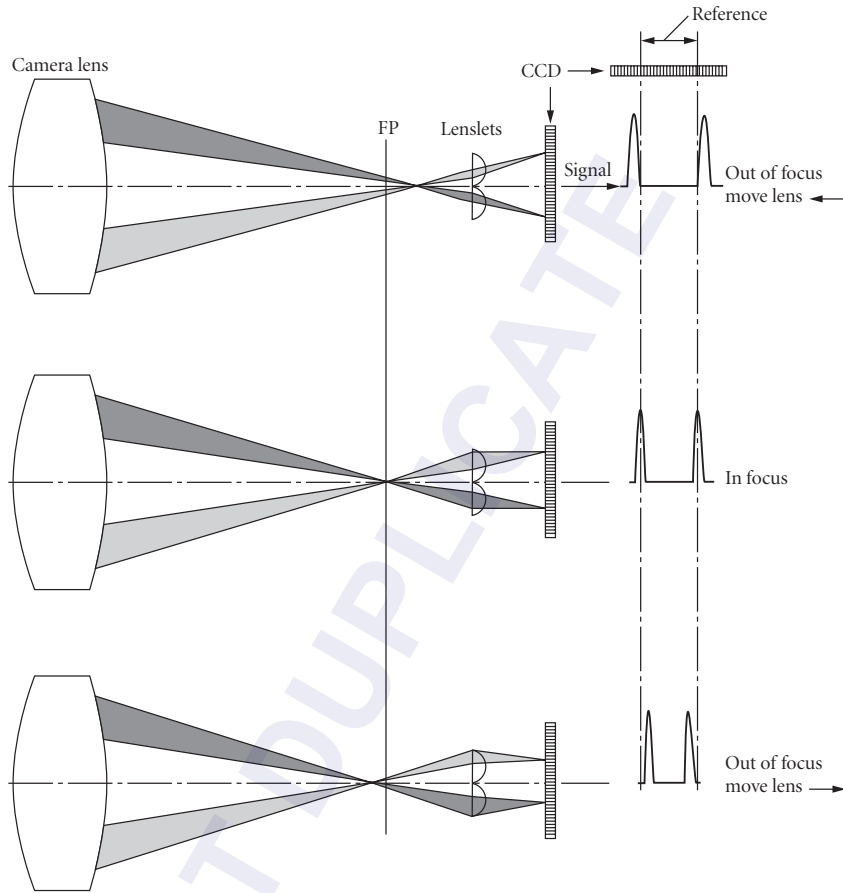


FIGURE 6 Autofocus SLR.

plane. Each lenslet receives light only from its side of the exit pupil and projects it onto a CCD line array, one for each lenslet. The relative position of each image on its CCD strip is analyzed by the system's microcomputer which is programmed to recognize the focus condition as a function of the CCD's signals. If the signals deviate from the programmed values, the microcomputer issues the appropriate command to the focus motor.

For off-center subjects, it's necessary to prefocus on them by pressing the trip button halfway, holding it there as you recompose the scene, then pressing all the way on the trip button to make the exposure. This is asking too much of a photographer shooting any sort of action, and many of them mistrusted their autofocus SLRs. In response, camera makers offered new models with broader CCD arrays to provide a larger central region of autofocus sensitivity. Some of these can be switched between narrow and broad sensitivity regions.

Other refinements to SLR autofocusing include

- Optimization of camera settings to maximize depth of field
- Prediction of moving subject's distance at instant of exposure
- Accommodation for horizontal and vertical subject detail
- Focus priority according to position of user's eye

To optimize depth of field, the user aims the camera at the near point and presses the trip button halfway. This is repeated for the far point. Then the scene is recomposed in the viewfinder and the exposure is made with the actual focus set automatically to some midpoint calculated by the camera's microcomputer.

For predicting the distance of a moving subject, the subject's motion should be constant, both in direction and velocity. Under these conditions, the autofocus sensor's signals can be used to calculate where the subject will be when the exposure is made. The calculation must consider the camera's inherent time lag.

Early AFSLRs used focus sensors that were shaped to respond to vertical image detail, with diminishing response as the detail approached the horizontal, where they were unable to respond. One solution incorporates three sets of lenslets and their CCD detector arrays. One set is laid out horizontally to respond to vertical detail, while the other two sets are vertical and straddle the first to form the letter "H." The two vertical sets respond to horizontal detail. Another solution has the individual, rectangular, detector segments (pixels) slanted to respond to both horizontal and vertical details.

By combining information from the autofocus detector and the focal length tracer in a zoom lens, some AFSLRs can maintain the image size (within the limits of the zoom range) chosen by the user, even as the subject distance changes.

Eye Tracking Figure 7 shows how Canon's model EOS A2E overcomes the need for the subject to be centered in the viewfinder in order to be in focus. Canon devised an eye tracking system that detects what portion of the viewscreen the user is looking at. Using low-power infrared emitting diodes (IREDs) to illuminate the eye, the system is matched to the user by having him/her look at the extremities of the five autofocus aiming patches in the viewfinder's center. The reflections from the eye are detected by a 60×100 pixel CCD array and the resulting signals are stored in the camera's

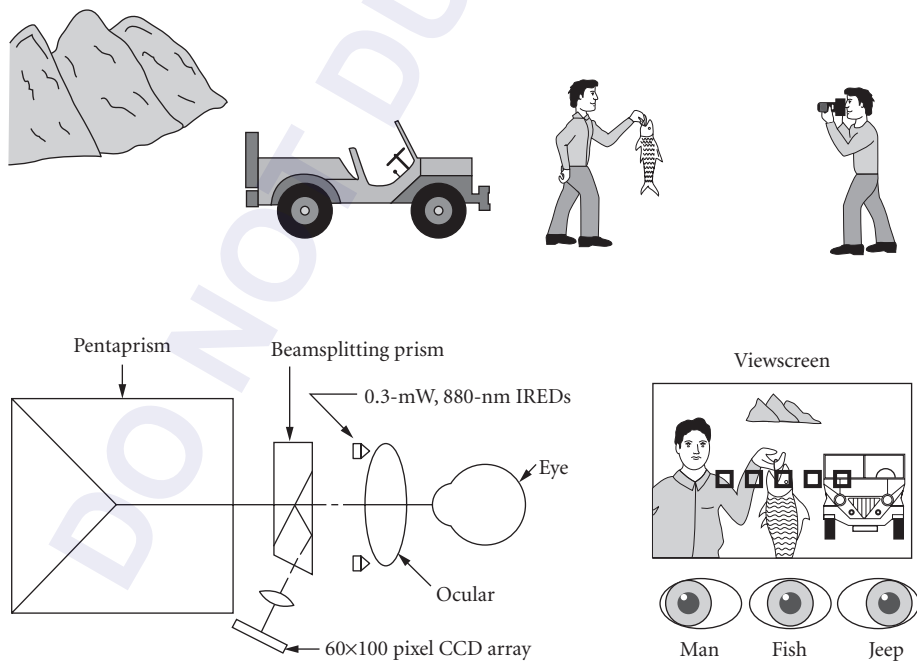


FIGURE 7 Canon's Eye Tracking SLR. The user's eye, illuminated by IREDS, is imaged on a CCD array. The resulting signal shifts in step with eye movements, causing a corresponding focus patch on the viewscreen to glow, indicating where the camera should focus.

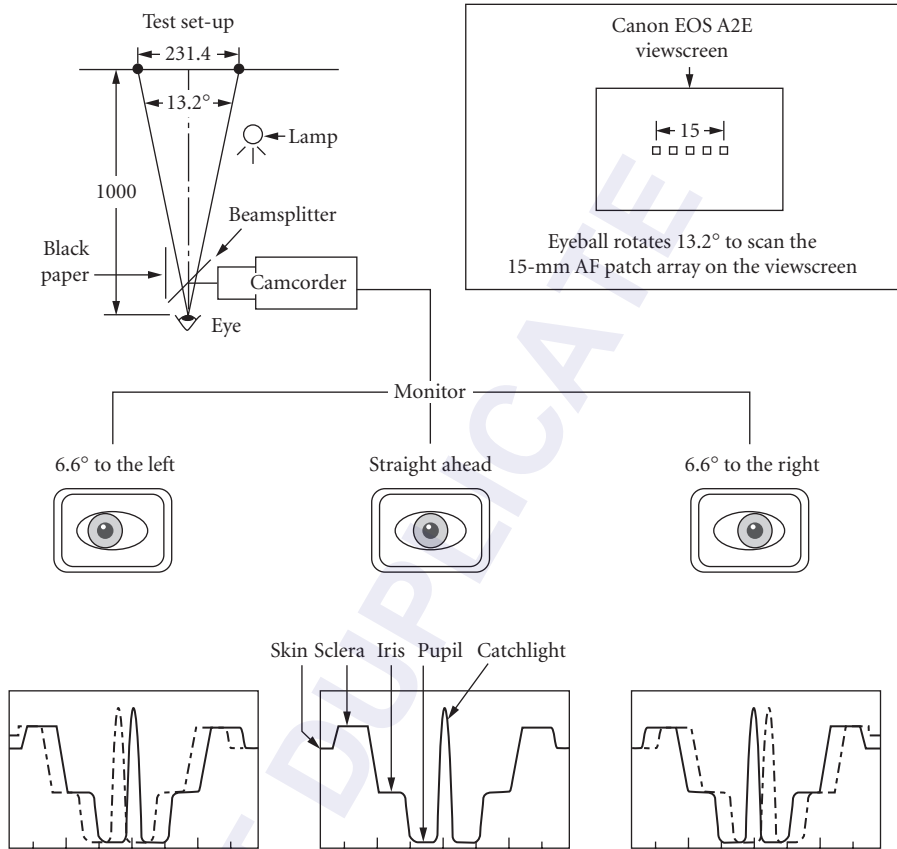


FIGURE 8 Eye tracking experiment. Line scan of the monitor's image at half screen height. (Dashed lines indicate image shift.)

memory. The five aiming patches occupy a 15-mm horizontal strip at the center of the finder. As it is, the camera's 16 user-selectable operational modes include one in which both the autofocus and autoexposure systems are commanded by the eye-tracking feature.

If the user wants to preview the depth of field, all that's necessary is to look at a small patch near the finder's upper left corner (not shown here). This brief glance causes the lens to close down to the f-stop chosen by the autoexposure system.

Because this eye-tracking feature is in an SLR, the user can see if it's working as expected just by looking at the viewscreen image. This indicates if, but not how, it works. To see how it works, I set up a simple experiment to measure the distribution of the light reflected from my eye as I shifted my gaze between two marks on a wall. The separation between the marks and their distance from my eye were chosen to duplicate the angle swept by the eye when looking from one side to the other of the 15-mm focus patch array on the Canon EOS A2E viewfinder. As indicated by Fig. 8, the format was nearly filled with the image of my eye. Consistent eye placement was assured with a chin and head rest. Once the image of my eye was recorded on tape, I could play back and pause at any point, then select a line at half screen height and store its waveform in a storage oscilloscope. By superimposing line scan waveforms from the frames showing my gaze from one side to the other, I could easily see the difference and dismissed my skepticism. This novel feature has intriguing possibilities.

25.15 FLASH

Many 35-mm cameras feature a built-in electronic flash unit. Some are designed to flash every time the shutter is tripped, unless the user switches off the flash. Others fire only when the combination of scene luminance and film speed calls for flash. In some of the more advanced models with zoom lenses, the beam angle emitted by the flash changes in step with the focal length setting of the lens.

Red Eye

In the interest of compactness, the majority of cameras with built-in flash units have the flash close to the lens. The resulting flash photos of people frequently exhibit what is commonly known as “red eye,” which describes the eerie red glow in the image of the pupils of a subject’s eyes. The red glow is the light reflected from the retina, which is laced with fine blood vessels. Young, blue-eyed subjects photographed in dim light seem to produce the most intense red-eye images.

The effect is reduced by (1) increasing the angle subtended to the subject’s eye by the separation between the centers of the lens and the flash; (2) reducing the subject’s pupil diameter by increasing the ambient brightness or having the subject look at a bright light for a few seconds before making the exposure.

Examples of how some camera makers fight red eye include Kodak’s Cobra Flash, used on several of their point-and-shoot models, and the “preflash,” used on many different camera makes and models. The Cobra Flash describes a flash unit whose flashlamp/reflector unit is hinged at the camera’s top. When the camera is not in use, the flash is folded down, covering the lens. To use the camera, the flash is swung up, positioning it further from the lens than would be possible if it had been contained in the camera’s main body. One of their most compact cameras featuring the Cobra Flash is the Cameo motordrive model, which slips easily into a dress shirt pocket when the flash is folded down. When opened for use, the flash is 72 mm above the lens. Test shots were free of red eye when the subject was no more than 7 ft away.

Another Kodak approach to the elimination of red eye is their single-use Fun Saver Portrait 35, whose integral electronic flash unit points upward, instead of forward. To use the camera, a simple white plastic panel, hinged at the camera’s top rear edge above the flash is pulled open. It latches at a 45° angle to switch the flash circuit on and direct the light from the flash forward. The result is a diffused beam that appears to originate from a point 100 mm above the lens.

Other makes and models have integral flash units that pop up at short distance when put into play. This may only gain several millimeters of lens-to-flash separation, but my experiments indicate that, as sketched in Fig. 9, for every extra millimeter of separation between the lens and the flash, the (red-eye-free) subject distance can be increased about 30 mm.

Several 35-mm cameras use the *preflash* method to reduce red eye by emitting a brief, rapid burst of low intensity flashes just before the main flash goes off for the exposure. A variation uses a steady beam from an incandescent lamp in the flash unit. The beam switches on shortly before the flashlamp fires for the exposure. The purpose in both methods is to make the subject’s pupils close down, reducing the light reflected from the eye during the exposure.

The preflash approach has two drawbacks: (1) it drains energy from the camera’s battery, reducing the number of pictures per battery; (2) many times the subject reacts to the preflash and blinks, just in time for the exposure.

25.16 FLEXIBILITY THROUGH FEATURES AND ACCESSORIES

The seemingly endless combinations of operating modes with a camera like the Canon EOS A2E might be taken as an attempt to be all things to all photographers. Another way to look at it is to see it as a 3-lb Swiss army knife: you’ll never use all of the tools all of the time, but if there’s the need for some tool, even just once, it might be nice to know you have it.

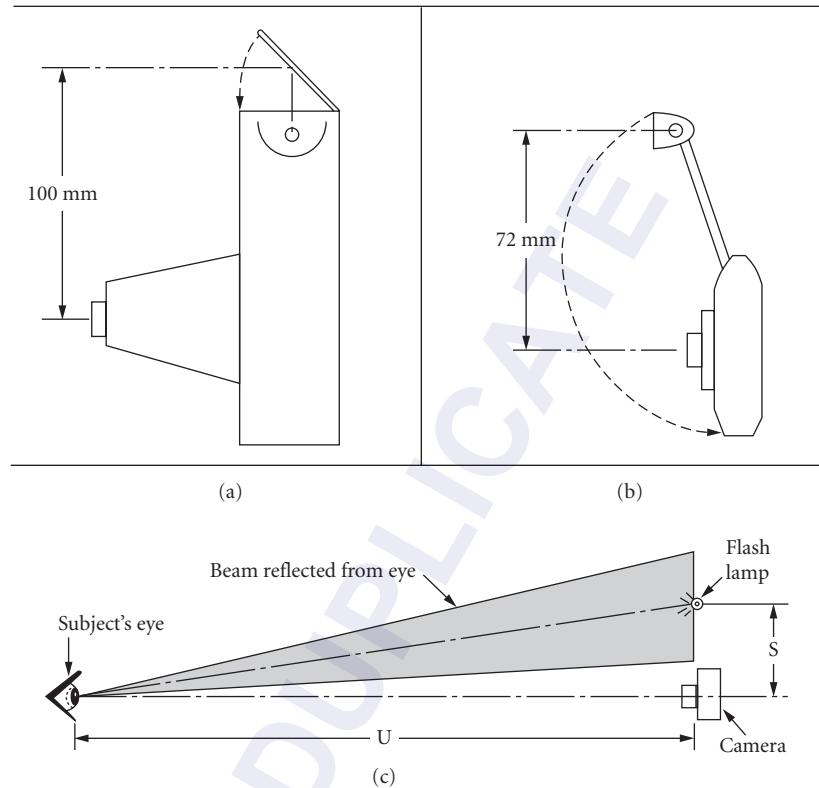


FIGURE 9 (a) Kodak's "Fun Saver Portrait 35" bounces its flash from a folding reflector. (b) Kodak's "Cameo Motordrive" uses the folding "Cobra Flash." (c) A beam reflected from the subject's eye misses the lens when the subject's distance "U" is not more than 20 S.

Many cameras have long lists of accessories. A typical camera system can be thought of as a box with an open front, top, and rear. For the front, the user may choose from as many as 40 different lenses. For the top, there may be three or more viewfinder hoods. For the back, choose one of perhaps five image receptacles.

Then there are the other groups, shown in Fig. 10: flash units, motor drives, close-up hardware, carrying cases, neck straps, lens hoods, filters, remote control cables, transmitters and receivers, mounting brackets, eyepiece magnifiers, corrective eyepiece lenses, cold weather heavy-duty battery packs, and more.

No matter how varied your photographic needs may be, the camera maker wants you to find everything you need in his or her catalog. Possibly, the availability of just one accessory, such as a wide-angle lens with tilt-shift controls for perspective correction, can decide which camera you choose.

25.17 ADVANTAGES OF VARIOUS FORMATS

In terms of versatility through a broad range of accessories plus the camera's intrinsic capabilities, it's hard to beat one of the major brands of 35-mm SLRs. No other type of camera has had as much ingenuity and as many refinements lavished on it for so many years. It's one of the most highly evolved consumer-oriented products.

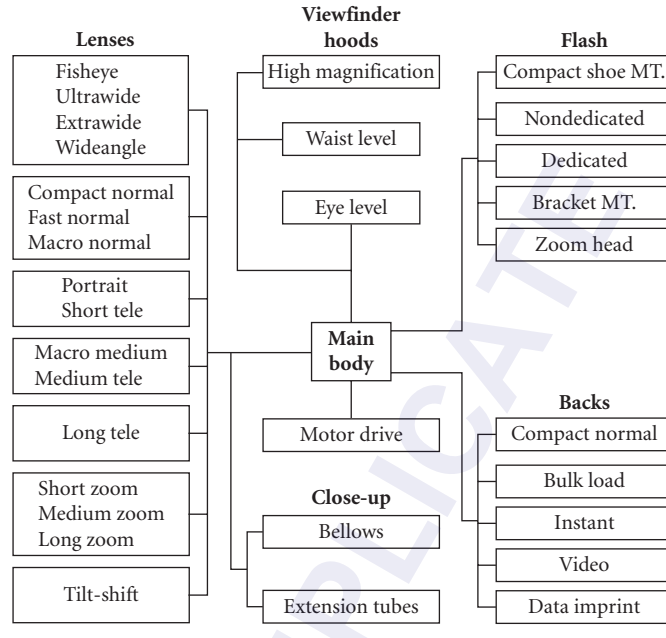


FIGURE 10 Camera system.

Accompanying the evolution in optics, mechanics, and electronics, film emulsions have improved over the years, making the 35-mm format just as able as the larger formats for most applications. Even so, all else being equal, there is no substitute for “real estate”—the precious additional square millimeters of emulsion offered by the many 120-size medium formats. As the data in Fig. 11 shows, some of these are SLRs with systems as extensive as their 35-mm counterparts.

25.18 LARGE FORMAT: A DIFFERENT WORLD

When you make the jump from medium-format to large-format, you’re in a different world. You use individual sheets of film, not rolls. Your camera will be used on a tripod or copy stand most of the time. Your photography will be contemplative, careful, and unhurried—perhaps better.

Scene composition and focusing are done with the lens at full aperture. Then the lens is stopped down, the shutter closed, the film holder inserted, its dark slide pulled, the shutter tripped, the dark slide replaced, and the film holder removed.

In a short time you’ll realize that the large-format (view) camera can be thought of as a compact optical bench. As such, it lends itself to special applications that could be difficult for the smaller formats.

View Camera Versatility

To illustrate, suppose you need a picture of a picket fence at some obliquity, with every picket board, from near to far, in sharp focus and with the lens wide open. This calls for the use of the “Scheimpflug condition,” shown in Fig. 12. It requires that the planes containing the lensboard, film, and subject all intersect on a common line. When this condition is satisfied, the entire surface of the subject plane will be in focus, even with the lens wide open.

Camera type	5 MM	Medium format	Large format
Format sizes	18×24 mm	45×60 mm	2-1/4×3-1/4"
	24×24 mm	60×60 mm	3-1/4×4-1/4"
	24×36 mm	60×70 mm	4×5", 5×7"
		60×90 mm	8×10"
Viewfinder	Galilean, SLR	Galilean, TLR, SLR	Galilean, FP
Focusing	CRF, SLR, AF	CRF, TLR, SLR	CRF, FP
Exposure	Manual, AE, TTL	Manual, AE, TTL	Manual
Options	Mot, Bulk, Pol, Dig, Data, Hood, Lens, Pgm, DEF	Mot, Bulk, Pol, Dig, Data, Hood Lens, DEF	Pol, Lens

Abbreviations:

AE = autoexposure, AF = autofocus, Bulk = bulk film back, CRF = coupled rangefinder, Data = data recording, DEF = dedicated electronic flash, Dig = digital image detector, FP = film plane, Hood = interchangeable viewfinder hoods, Lens = interchangeable lenses, Mot = motor drive, Pgm = programable AE, AF, other functions, Pol = polaroid film back, SLR = single lens reflex, TLR = twin lens reflex, TTL = through the lens metering.

FIGURE 11 Major camera features.

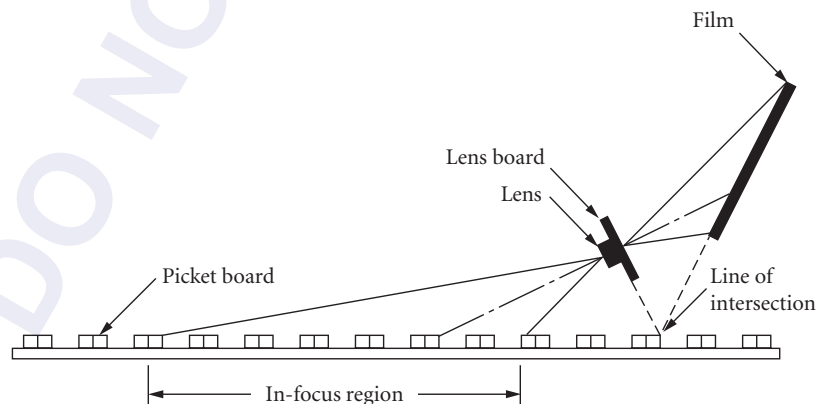


FIGURE 12 Scheimpflug condition. All of the picket boards within the field of view will be in focus when the planes of the lens board, film, and picket boards intersect on a common line.

The necessary camera movements, involving lensboard and film plane, are standard features of even the most Spartan view cameras. These movements are known as swings and tilts. They take just a few seconds to adjust on a view camera and the job doesn't require a special lens. You can do it with a smaller format camera too, but you'll need one of their special (expensive) tilt-shift lenses or a bellows unit with articulated front and rear panels, plus a lens with a large enough image circle. The resulting combination may not retain all of the small-format camera's features, such as exposure metering, autofocus, etc.

The view camera's fully articulated front and rear provide for swing, tilt, rise, fall, and left-right shift. Thanks to this flexibility, objects such as boxes and buildings can be photographed without distortion, and distracting detail near the image borders can be omitted.

It takes first-time users a while to get used to the inverted and reversed image seen on the view camera's groundglass screen. This can be annoying when shooting a portrait, since an upside-down smile looks like a frown until you accept the fact that even though you understand the basic camera optics, it doesn't mean you have to enjoy coping with it. Worse, you'll need to drape a dark cloth over the back of the camera and over your head in order to see the image if you're working in bright light. If you're claustrophobic, this may bother you.

On the plus side, large-format negatives are frequently contact-printed or only slightly enlarged for the final image. Because the image is large, depth of field and other image properties can be examined easily on the groundglass viewscreen with a small magnifier of modest power—a 4× loupe works well. The large negative has another attribute: it lends itself to retouching, masking, and other image manipulations, but these may be lost arts now that clever computer programs are available for doing the same things, provided your image is in digital form.

25.19 SPECIAL CAMERAS

Some photographic tasks call for cameras with special features, such as the ability to form images in near-total darkness or inside of a crowded mechanism. Among the long list of special cameras, we find

- Aerial
- Clandestine
- Endoscopic
- High-speed
- Periphery
- Sewer
- Stereo (3-D)
- Streak
- Thermal imaging
- Underwater
- Wide-angle

Aerial Cameras

Aerial cameras come in a variety of sizes and features. Among the more common features are image motion compensation, where focal length, speed, and altitude are factored into the movement of the film during the exposure; a vacuum back to hold the film flat during the exposure; and a calibrated lens so that any rectilinear distortion can be factored into the measurements made of the image.

Clandestine Cameras

Clandestine or “spy” cameras have been with us since photography was invented. In the broadest sense, any camera that is not recognized as such by the subject being photographed might be considered a successful spy camera. Many early box cameras were dubbed “detective” cameras because they were much smaller and more drab than a “real” camera with its prominent bellows and sturdy stand.

Cameras have been disguised as books, rings, binoculars, cigarette packs and lighters, matchboxes, portable radios, briefcases, canes, cravats, hats, even revolvers. Among them, the classic Minox is probably the best known. It can be concealed in an adult’s fist, focuses down to eight inches, and is nearly silent. Its smooth exterior and gently rounded corners have inspired the belief among many that it was designed to be concealed in a body cavity with minimal discomfort.

Endoscopic Cameras

Endoscopic cameras use a tiny, short-focal-length lens to form an image that’s transferred by a coherent, flexible fiber-optic bundle to a relay system that forms the image on the detector (film or CCD) in the camera. To illuminate the subject, the coherent bundle may be surrounded by an incoherent ring of fibers optically coupled to a light source at its free end, close to the camera.

Often fitted with a 90° prism on its tip, these cameras are used to photograph inside humans and machines. Another application is shown in Fig. 13: getting close-up views of architectural models from “ground” level. Variations include those without illumination optics but having a very small diameter image bundle to fit inconspicuously in some object for surveillance photography.

High-Speed Cameras

High-speed cameras were once defined as being able to make exposures of less than 1/1000 s. Today this would include many 35-mm SLRs which have a top speed of 1/10,000 s, a speed equaled by several consumer-grade camcorders. When shorter exposures are called for, a common, low-cost electronic flash unit can give flash durations as short as 1/32,000 s.

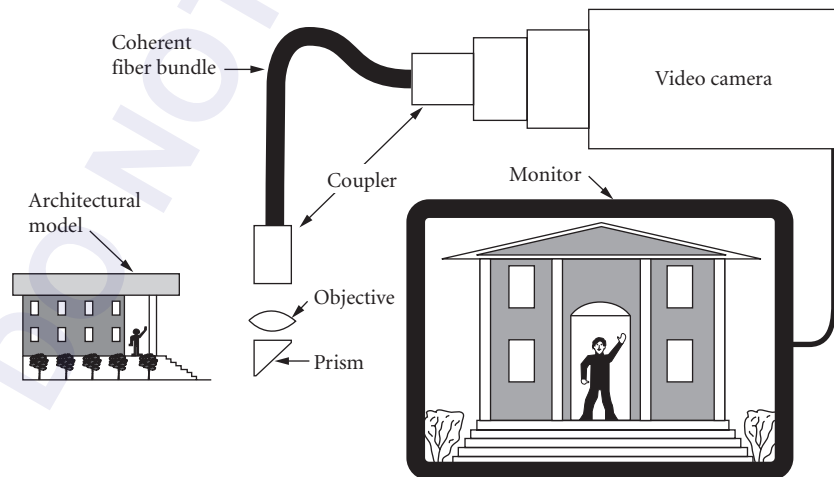


FIGURE 13 Endoscopic camera. While most often used for medical purposes, the endoscopic camera’s properties make it valuable for photographing miniature scenes from the perspective of a miniature photographer.

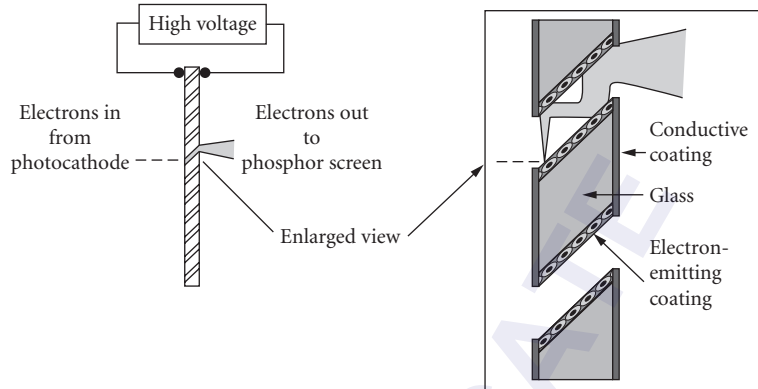


FIGURE 14 Microchannel plate.

The next step includes the Kerr cell and Faraday shutters, both of which work by discharging a high-voltage capacitor across a medium located between crossed polarizers. This produces a momentary rotation of the plane of polarization within the medium, permitting light to pass through to the detector. Exposure times are in the nanosecond range for these electro-optical/magneto-optical devices.

For exposures in the picosecond range accompanied by image intensification, there's the electronic image tube. When a lens forms an image on the photocathode at the front of this tube, electrons are emitted. Their speed and direction are controlled by electrodes within the tube. A secondary image is formed by the electrons as they strike the phosphor screen at the rear of the tube. This image may be photographed, or, if the tube has a fiber-optic faceplate behind the screen, the image can be directly transferred to a film held against the faceplate.

By placing a microchannel plate in front of the phosphor screen, the image can be intensified by a factor of 10,000 or more. A microchannel plate is a thin glass disk riddled with microscopic holes that pierce the disk at an angle. In Fig. 14 the wall surface of each hole is coated with a substance that reacts to the impact of an electron by emitting more electrons. A high voltage across the disk accelerates the stream of electrons. For every electron that enters one of the angled holes, about 100,000 electrons emerge to strike the phosphor screen.

Periphery Cameras

A periphery camera is used to make photos of objects like gas engine pistons, bullets, and other cylindrical objects whose surface detail must be imaged as though the surface was "unrolled" and laid out flat before the camera. Depending on the size of the subject, either it or the camera is rotated about its longitudinal axis at a constant angular velocity. The image strikes the film moving behind a slit that's parallel to the axis of rotation. The film's velocity matches that of the image unless deliberate image compression or elongation is desired.

Sewer Cameras

A sewer camera is designed to photograph the inside of pipes, tunnels, etc. It may be thought of as a small underwater camera on a sled. The camera's lens is encircled by an electronic flashtube and reflector to illuminate the scene. Pictures are made at regular intervals, as judged by distance marks on the cable attached to the sled. Other cables attached to the camera convey signals to and from

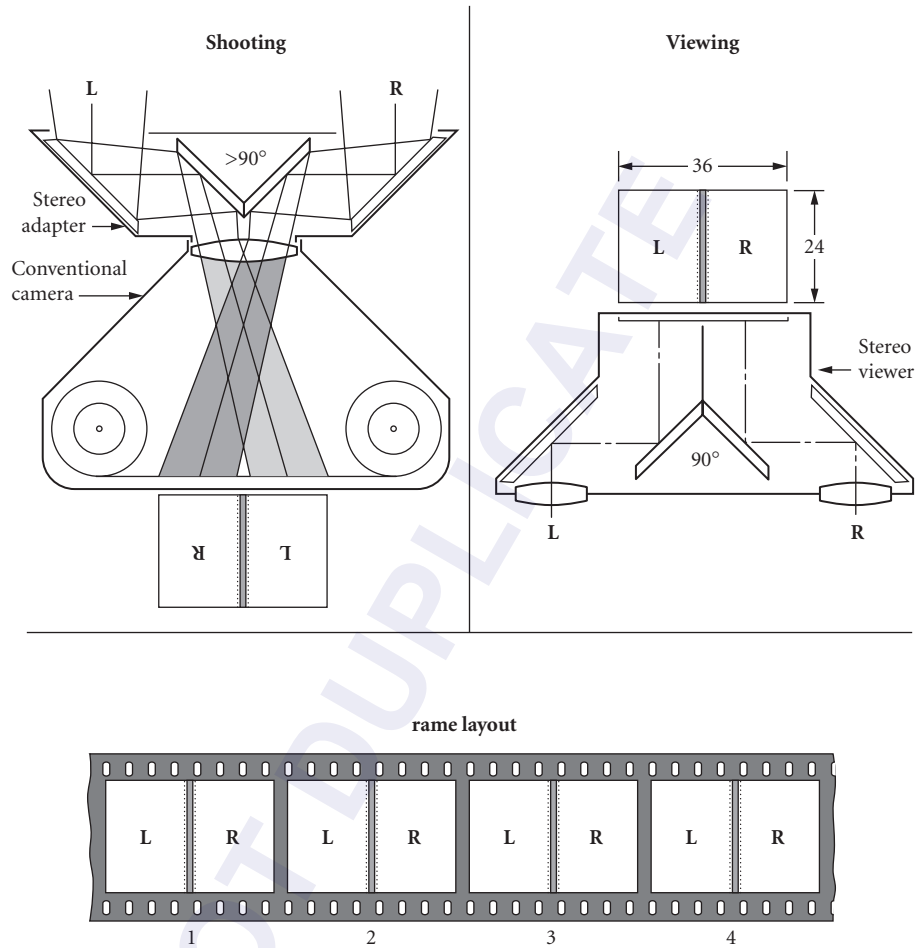


FIGURE 15 Stereo adapter set.

the camera. With the miniaturization of video cameras, they have taken over this task, except where maximum resolution is required. This is where film cameras have excelled.

Stereo Cameras

Stereo cameras seem to come in and out of vogue with some mysterious rhythmic cycle. The root idea has been around since the dawn of photography and is based on the parallax difference between the views of our left and right eyes. The classic stereo camera mimics nature by using two lenses spaced about 65 mm apart to form two images of the subject.

The two images can be made in other ways. A simple reflection system using four small mirrors or an equivalent prism system placed in front of a normal camera's lens will form two images of the subject, as shown in Fig. 15. Another method requires that the subject be stationary because two separate exposures are made, with the camera being shifted 65 mm between exposures. In aerial stereo photography, two views are made of the ground, some seconds apart.

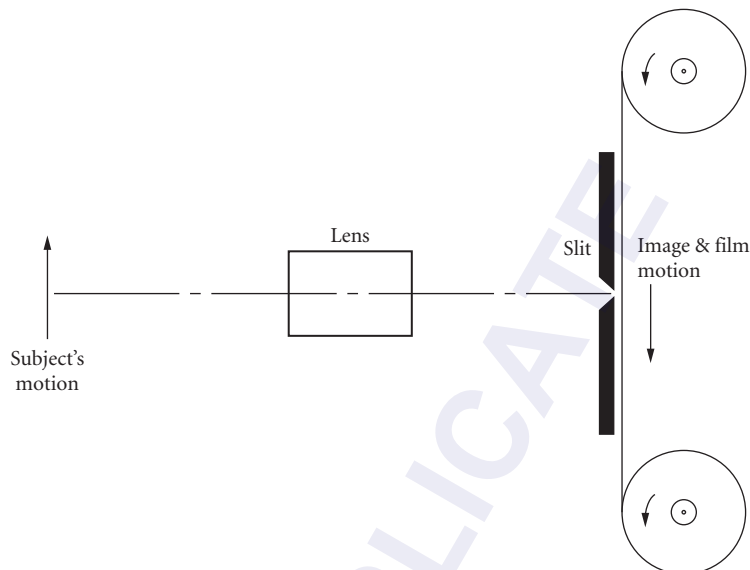


FIGURE 16 Streak camera. When film motion matches image motion, image will be free of distortion. If the film moves too fast, the image will be stretched. If the film moves too slow, the image will be compressed.

When the images are viewed in a manner that restricts the left and right images to their respective eye, the stereo effect is achieved. Various methods for viewing stereo pairs include projection, where the left and right views are polarized at 90° to each other. The viewer wears glasses with polarizing filters oriented to let each eye see the view intended for it.

Another viewing system is called a *parallax stereogram*. It (optically) slices the left and right images into narrow, interlaced strips. When viewed through a series of vertical lenticular prisms with a matching pitch, the 3-D effect is seen.

Streak Cameras

Streak cameras are useful for studying relative motion between the subject and camera. They share certain characteristics with the periphery camera described previously, insofar as they match the movement of the film to that of the image coming through a slit at the film plane. Exposure time is determined by how long it takes for a point on the film's surface to travel across the slit's width.

The basics of the streak camera are shown in Fig. 16. It would be pointless to use a streak camera without some relative motion between the image and film. Some photographers use a streak camera for creative effects, such as depicting motion tack-sharp at its beginning, then gradually elongating or compressing it, and ending in a smear. This is done by varying the relative velocity between the image and the film during the exposure, either by moving the camera, the subject, or the film. These motions may be made singly or in combination. Varying the focal-length setting of a zoom lens with the film moving also produces unusual images.

A streak camera's format has a width defined by the film it uses, but each picture has its own length, limited only by the length of the roll of film. One of the more critical factors to look for in a streak camera is freedom from *cogging*, a local density variation in exposure while the film is moving at a fixed velocity. The result of periodic or intermittent speed variations, the cause may be improperly meshed gears, a bad bearing, poor fit between the film drive sprocket teeth and the film's sprocket perforations, or the magnetic pole effects of the drive motor.

Thermal Image Cameras

Thermal imaging cameras convert the intrinsic heat of a subject into a visible image. Among their many applications are detection of heat losses from buildings, blood circulation disorders, and surveillance. Some of these cameras produce false color images, in which each color represents a different temperature.

Among the various methods to form visible images of temperature variations, the most direct way is to use a normal camera loaded with film that's sensitive to the infrared portion of the electromagnetic spectrum.

In a more elaborate system, a moving mirror scans the subject and, line-by-line, projects its image onto a heat-sensitive semiconductor device whose output is proportional to the IR intensity. The output is used to modulate a beam of light focused on the surface of a conventional film.

Another version uses the semiconductor to modulate a stream of electrons striking a phosphor screen in an image converter tube. The image can then be photographed.

Underwater Cameras

Underwater cameras come in a wide variety of sophistication, from the inexpensive disposable to the expensive high-tech versions. In between, there are dozens of underwater housings designed for specific cameras. Typically, these housings permit the user to change the camera's settings through watertight couplings. Most of the cameras used in such housings have motorized film advance, autoexposure, and autofocus, so the only external control needed is a pushbutton at one end of a simple electrical switch.

External attachments include flashguns, viewfinders, and ballast weights. The flashgun connections should be carefully examined because they are one of the leading sources of problems. In general, the simpler the connector, the better.

Wide-Angle Photography

There are several 35-mm and medium-format cameras designed specifically for wide-angle photography. These include straightforward types which use lenses designed for wide-angle views on larger-format cameras. Essentially these cameras use only a rather long horizontal strip of the broad image circle the lens produces. This type of camera is uncomplicated and rugged.

Panoramic Cameras A special kind of wide-angle camera is known as a panoramic camera, and there are two main types: one where the entire camera rotates; the other, where just the lens rotates.

The rotating camera type is capable of a full 360° vista. As the camera turns on its vertical axis, the film is moved past a narrow, stationary slit at the center of the film plane. The motion of the film is matched to that of the image. Because these cameras rotate slowly, a common prank in photos of large groups is for the prankster to stand at the edge of the group that's exposed first, then dash behind the group to the opposite edge in time for its exposure, with the result that the same person appears twice in the same photo, once at either edge of the group.

The rotating lens type shown in Fig. 17 produces images of about 140°. It works by rotating its lens on a vertical axis coinciding with its real nodal point. The image is swept across the film through a tubular image tunnel at the rear of the lens. The tunnel extends almost to the film surface and has a narrow slit at its end. The slit is parallel to the axis of rotation and extends over the width of the film. During the exposure the film is held stationary against a cylindrical film gate whose radius equals the focal length of the lens. The slit width, the rotating speed, and the lens opening may be adjusted for exposure control.

The panoramic cameras described here regulate their speed of rotation with precision governing systems to ensure edge-to-edge uniformity of exposure, so they should be kept as clean as possible. Also, to avoid unpleasant distortion, use care in leveling them and always use the best single camera accessory money can buy: a good, solid tripod.

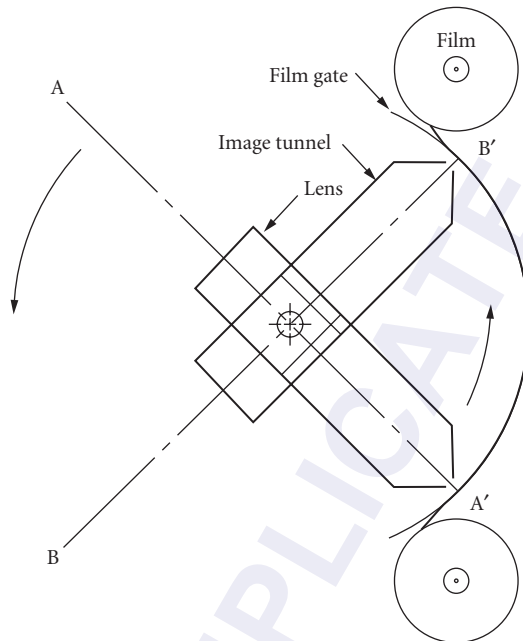


FIGURE 17 Panoramic camera. The lens rotates about its rear nodal point from A to B. Image-forming light reaches the film from A' to B' through a slit at the end of the image tunnel.

25.20 FURTHER READING

- Clerc, *Photography*, vol. 1 and 2, Focal Press, 1970.
 Edgerton, *Electronic Flash Strobe*, MIT Press, 1979.
 Feynman, *The Feynman Lectures on Physics*, vol. 1 Addison, 1975.
 Goldberg, *Camera Technology: The Dark Side of The Lens*, Academic Press, 1992.
 Habell and Cox, *Engineering Optics*, Pitman, 1966.
 Hyzer, *Engineering and Sci. Hi-Speed Photography*, McMillan, 1962.
 Jenkins and White, *Fundamentals of Optics*, McGraw-Hill, 1957.
 Kingslake, *Applied Optics and Optical Engineering*, vol. IV, Academic Press, 1967.
 Kingslake, *Optical System Design*, Academic Press, 1983.
 Kingslake, *Optics in Photography*, SPIE Optical Engineering Press, 1992.
 Morton, *Photography for the Scientist*, Academic Press, 1984.
 Ray, *Applied Photographic Optics*, Focal Press, 1988.
 Ray, *The Photographic Lens*, Focal Press, 1979.
 Smith, *Modern Optical Engineering*, McGraw-Hill, 1966, 1990.
 Spencer, *The Focal Dictionary of Photo Technologies*, Focal Press, 1973.
 Stimson, *Photometry and Radiometry for Engineers*, Wiley, 1974.
 Williamson and Cummins, *Light and Color in Nature and Art*, Wiley, 1983.

Gerald C. Holst

JCD Publishing
Winter Park, Florida

26.1 GLOSSARY

Quantity	Definition	Unit
A_D	photosensitive area of a detector	m^2
c	speed of light	$3 \times 10^8 \text{ m/s}$
C	sense node capacitance	F
$E_{e\text{-faceplate}}(\lambda)$	spectral radiant incidence	$\text{W } \mu\text{m}^{-1} \text{ m}^{-2}$
$E_{q\text{-faceplate}}(\lambda)$	spectral photon incidence	$\text{photons s}^{-1} \mu\text{m}^{-1} \text{ m}^{-2}$
d	detector size	mm
d_{CCH}	detector pitch	mm
D	aperture diameter	m
DR_{array}	array dynamic range	numeric
$\text{DR}_{\text{camera}}$	camera dynamic range	numeric
f_l	focal length	m
F	f -number	numeric
G	source follower gain	numeric
h	Planck's constant	$6.626 \times 10^{-34} \text{ J s}^{-1}$
M_e	radiant exitance	W/m^2
$M_q(\lambda)$	spectral photon exitance	$\text{photons s}^{-1} \mu\text{m}^{-1} \text{ m}^{-2}$
m_{optics}	optical magnification	numeric
M_v	photometric exitance	lumen m^2
n_{dark}	number of dark current electrons	numeric
n_{PE}	number of photoelectrons	numeric
n_{well}	charge well capacity	numeric
q	electronic charge	$1.6 \times 10^{-19} \text{ C}$
R_{eq}	equivalent resolution	mm

$\mathfrak{R}_e(\lambda)$	detector responsivity	A/W
\mathfrak{R}_{ave}	average detector responsivity	V/(J cm ⁻²)
t_{int}	integration time	s
U	nonuniformity	numeric
V_{max}	maximum output voltage	V
V_{signal}	voltage created by photoelectrons	V
$\langle n_{floor} \rangle$	noise created by on-chip amplifier	rms electrons
$\langle n_{pattern} \rangle$	pattern noise	rms electrons
$\langle n_{PRNU} \rangle$	photoresponse nonuniformity noise	rms electrons
$\langle n_{shot} \rangle$	shot noise	rms electrons
$\langle n_{sys} \rangle$	total array noise	rms electrons
η	image space horizontal spatial frequency	cycles/mm
η_C	optical cutoff in image space	cycles/mm
η_N	Nyquist frequency in image space	cycles/mm
η_s	sampling frequency in image space	cycles/mm
$\eta(\lambda)$	spectral quantum efficiency	numeric
λ	wavelength	μm

26.2 INTRODUCTION

The heart of the solid-state camera is the solid-state array. It provides the conversion of light intensity into measurable voltage signals. With appropriate timing signals, the temporal voltage signal represents spatial light intensities. When the array output is amplified and formatted into a standard video format, a solid-state camera is created. Because charge-coupled devices (CCDs) were the first solid-state detector arrays, cameras are popularly called CCD cameras even though they may contain charge injection devices (CIDs) or complementary metal-oxide semiconductors (CMOS) as detectors.

Boyle and Smith¹ and Amelis et al.² invented CCDs in 1970. Then came considerable literature³⁻¹¹ on CCD physics, fabrication, and operation. A CCD refers to a semiconductor architecture in which charge is transferred through storage areas. This architecture has three basic functions: (1) charge collection, (2) charge transfer, and (3) the conversion of charge into a measurable voltage. The basic building block of the CCD is the metal-insulator semiconductor (MIS) capacitor. The most important MIS is the metal-oxide semiconductor (MOS). Because the oxide of silicon is an insulator, it is a natural choice.

Charge generation is often considered as the initial function of the CCD. With silicon photo-detectors, each absorbed photon creates an electron-hole pair. Either the electrons or holes can be stored and transferred. For frame transfer devices, charge generation occurs under an MOS capacitor (also called a photogate). For some devices (notably interline transfer devices), photodiodes create the charge.

The CID does not use a CCD for charge transfer. Rather, two overlapping silicon MOS capacitors share the same row and column electrode. Column capacitors are typically used to integrate charge while the row capacitors sense the charge after integration. With the CID architecture, each pixel is addressable (i.e., it is a matrix-addressable device).

Active pixel sensors (APSs) are fabricated with CMOS technology. The advantage is that one or more active transistors can be integrated into the pixel. As such, they become fully addressable (to read selected pixels) and can perform on-chip image processing.

Devices may be described functionally according to their architecture (frame transfer, interline transfer, etc.) or by application. To minimize cost, array complexity, and electronic processing, the architecture is typically designed for a specific application. For example, astronomical cameras typically use full-frame arrays, whereas video systems generally use interline transfer devices.

The separation between general imagery, machine vision, scientific devices, and military devices becomes fuzzy as technology advances.

26.3 IMAGING SYSTEM APPLICATIONS

Cameras for the professional broadcast television and consumer camcorder markets are designed to operate in real time with an output that is consistent with a standard broadcast format. The resolution, in terms of array size, is matched to the bandwidth that is recommended by the standard. An array that provides an output of 768 horizontal by 484 vertical pixels creates a satisfactory image for conventional (U.S.) television. Eight bits (256 gray levels) provide an acceptable image in the broadcast and camcorder industry.

In its simplest version, a machine vision system consists of a light source, camera, and computer software that rapidly analyzes digitized images with respect to location, size, flaws, and other pre-programmed data. Unlike other types of image analysis, a machine vision system also includes a mechanism that immediately reacts to an image that does not conform to parameters stored in the computer. For example, defective parts are taken off a production line conveyor belt.

For scientific applications, low noise, high responsivity, large dynamic range, and high resolution are dominant considerations. To exploit a large dynamic range, scientific cameras may digitize the signal into 12, 14, or 16 bits. Scientific arrays may have 5000×5000 detector elements.¹² Theoretically, the array can be any size, but manufacturing considerations ultimately limit it.

Although low-light-level cameras have many applications, they tend to be used for scientific applications. There is no industrywide definition of a low-light-level imaging system. To some, it is simply a solid-state camera that can provide a usable image when the lighting conditions are less than 1 lux (lx). To others, it refers to an intensified camera and is sometimes called a low-light-level television (LLTV) system. An image intensifier amplifies a low-light-level image that can be sensed by a solid-state camera. The image intensifier/CCD camera combination is called an intensified CCD (ICCD). The image intensifier provides tremendous light amplification but also introduces additional noise. The spectral response of the ICCD is governed by the image intensifier.

The military is interested in detecting, recognizing, and identifying targets at long distances. This requires high-resolution, low-noise sensors. Target detection is a perceptible act. A human determines if the target is present. The military uses the minimum resolvable contrast (MRC) as a figure of merit.¹³

26.4 CHARGE-COUPLED DEVICE ARRAY ARCHITECTURE

Array architecture is driven by its application. Full-frame and frame transfer devices tend to be used for scientific applications. Interline transfer devices are used in consumer camcorders and professional television systems. Linear arrays, progressive scan, and time delay and integration (TDI) are used for industrial applications. Despite an ever-increasing demand for color cameras, black-and-white cameras are widely used for many scientific and industrial applications.

The basic operation of linear, full-frame, frame transfer, and interline transfer devices is described in Chap. 32, "Visible Array Detectors," by Timothy J. Tredwell in Vol. II. This section describes some additional features.

Full-Frame Arrays

In full-frame arrays, the number of pixels is often based upon powers of 2 (e.g., 512×512 or 1024×1024) to simplify memory mapping. Scientific arrays have square pixels and this simplifies image-processing algorithms.

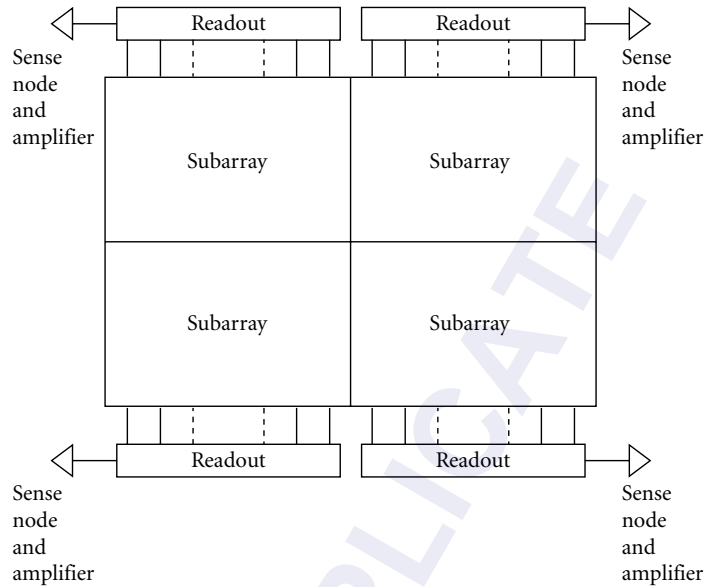


FIGURE 1 A large array divided into four subarrays. Each subarray is read out simultaneously to increase the effective data rate. Very large arrays may have up to 32 parallel outputs.

Data rates are limited by the amplifier bandwidth and, if present, the conversion capability of the analog-to-digital converter. To increase the effective readout rate, the array can be divided into subarrays that are read out simultaneously. In Fig. 1, the array is divided into four subarrays. Because they are all read out simultaneously, the effective clock rate increases by a factor of 4. Software then reconstructs the original image. This is done in a video processor that is external to the CCD device where the serial data are decoded and reformatted.

Interline Transfer

The interline transfer array consists of photodiodes separated by vertical transfer registers that are covered by an opaque metal shield (Fig. 2). Although photogates could be used, photodiodes offer higher quantum efficiency. After integration, the charge that is generated by the photodiodes is transferred to the vertical CCD registers in about $1 \mu\text{s}$. The main advantage of interline transfer is that the transfer from the active sensors to the shielded storage is quick. There is no need to shutter the incoming light. The shields act like a venetian blind that obscures half the information that is available in the scene. The area fill factor may be as low as 20 percent. Because the detector area is only 20 percent of the pixel area, the output voltage is only 20 percent of a detector that would completely fill the pixel area. A microlens can optically increase the fill factor.

Because interline devices are most often found in general imagery products, most transfer register designs are based upon standard video timing. Figure 3 illustrates a four-phase transfer register that stores charge under two gates. With 2:1 interlace, both fields are collected simultaneously but are read out alternately. This is called frame integration. With EIA 170 (formerly called RS 170), each field is read every $1/60$ s. Because the fields alternate, the maximum integration time is $1/30$ s for each field.

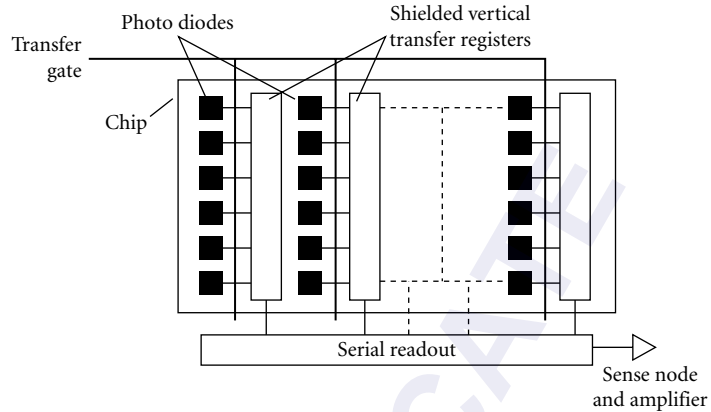


FIGURE 2 Interline transfer architecture. The charge is rapidly transferred to transfer registers via the transfer gate. Interline transfer devices can also have a split architecture similar to that shown in Fig. 1.

Pseudointerlacing (sometimes called field integration) is shown in Fig. 4. Changing the gate voltage shifts the image centroid by one-half pixel in the vertical direction. This creates 50 percent overlap between the two fields. The pixels have twice the vertical extent of standard interline transfer devices and therefore have twice the sensitivity. An array that appears to have 240 elements in the vertical direction is clocked so that it creates 480 lines. However, this reduces the vertical modulation transfer function (MTF). With some devices, the pseudointerlace device can also operate in a standard interlace mode.

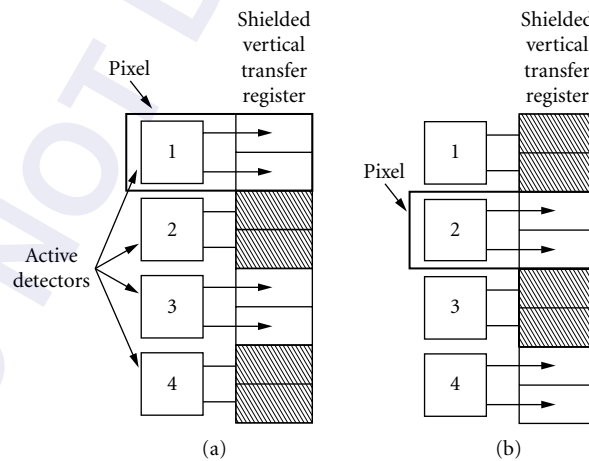


FIGURE 3 Detailed layout of the 2:1 interlaced array. (a) The odd field is clocked into the vertical transfer register and (b) the even field is transferred. The vertical transfer register has four gates and charge is stored under two wells. The pixel is defined by the detector center-to-center spacing, and it includes the shielded register area. The transfer gate is not shown.

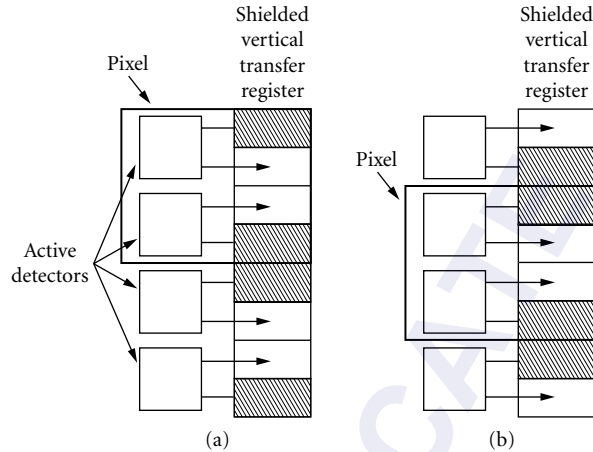


FIGURE 4 Pseudointerlace. By collecting charge from alternating active detector sites, the pixel centroid is shifted by one-half pixel. (a) Odd field and (b) even field.

26.5 CHARGE INJECTION DEVICE

A CID consists of two overlapping MOS capacitors sharing the same row and column electrode. Figure 5 illustrates the pixel architecture. The nearly contiguous pixel layout provides a fill factor of 80 percent or greater. Charge injection device readout is accomplished by transferring the integrated charge from the column capacitors to the row capacitors. After this nondestructive signal readout, the charge moves back to the columns for more integration or is injected (discarded) back into the silicon substrate.

Although the capacitors are physically orthogonal, it is easier to understand their operation by placing them side by side. Figure 6 illustrates a functional diagram of an array, and Fig. 7 illustrates the pixel operation. In Fig. 7a, a large voltage is applied to the columns, and photogenerated carriers (usually holes) are stored under the column gate. If the column voltage is brought to zero, the charge transfers to the row gate (Fig. 7b). The change in charge causes a change in the row gate potential that is then amplified and outputted. If V_1 is reapplied to the columns, the charge transfers back to the column gate. This is nondestructive readout and no charge is lost. By suspending charge injection,

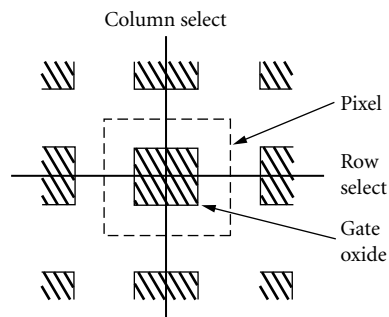


FIGURE 5 CID pixel architecture.

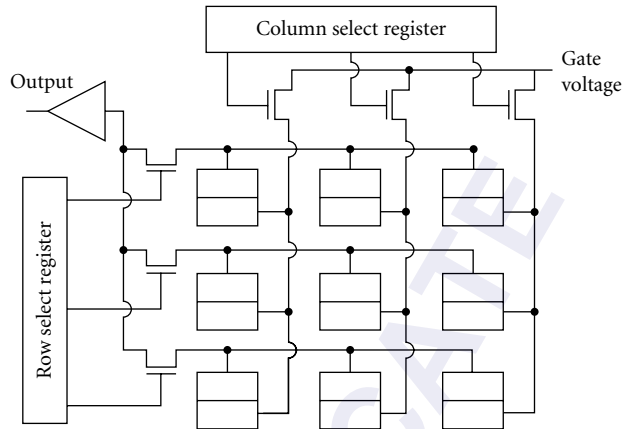


FIGURE 6 Functional layout of a 3×3 CID array. The row and column select registers are also called decoders. The select registers and readout electronics can be fabricated with CMOS technology.

multiple-frame integration (time-lapse exposure) is created. In this mode, the observer can view the image on a display as the optimum exposure develops. Integration may proceed for up to several hours. Reset occurs by momentarily setting the row and column electrodes to ground. This *injects* the charge into the substrate (Fig. 7c).

The conversion of charge into voltage depends upon the pixel, readout line, and amplifier capacitance. The capacitance is high because all the pixels on a given row are tied in parallel. Therefore, when compared with a CCD, the charge conversion is small, yielding a small signal-to-noise ratio (SNR). Because the readout is nondestructive, it can be repeated numerous times. The multiple reads are averaged together to improve the SNR.

Because each pixel sees a different capacitance, CIDs tend to have higher pattern noise compared with CCDs. However, off-chip algorithms can reduce the amount of pattern noise. With no charge transfer, CIDs are not sensitive to charge transfer efficiency effects. Without multiple gates, CIDs have larger well capacities than comparably sized CCDs. Charge injection devices inherently have antibloom capability. Because charge is limited to a single pixel, it cannot overflow into neighboring pixels. Perhaps the greatest advantage of CIDs is random access to any pixel or pixel cluster. Subframes and binned pixels can be read out at high frame rates.

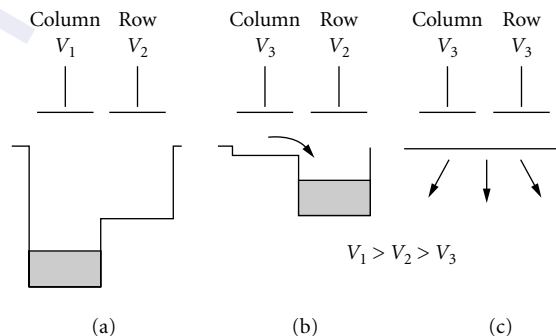


FIGURE 7 CID pixel operation. (a) Integration, (b) readout, and (c) injection.

26.6 COMPLEMENTARY METAL-OXIDE SEMICONDUCTOR

With the APS approach, highly integrated image sensors are possible.¹⁴ By placing processing on the chip, a CMOS camera can be physically smaller than a CCD camera that requires clocks, image reformatting, and signal processing in separate hardware. A sophisticated APS array can create¹⁵ a “camera on a chip.” It is possible to build a frame transfer device where pixels can be binned to enhance the SNR and provide variable resolution imaging.¹⁶

In 1993, Fossum¹⁷ described state-of-the-art active pixel concepts such as the double-gate floating surface transistor, charge modulation device, bulk charge modulation device, based-stored image sensor, and the static induction transistor. See also Ref. 14.

The APS contains a photodiode, row select transistor, and a reset transistor. As shown in Fig. 8, by activating a row, the data from the pixels in that row are simultaneously copied into the columns. Each column will have a load transistor, a column select switch, and a sampling switch. The chip may also have an analog-to-digital converter and provide correlated double sampling. The pixels are then reset and a new integration is started. The APS does not rely upon charge transfer. Rather, the photodiode drives a capacitance line. The total capacitance limits the chip speed. APSs can be fully addressable and subarrays can be read out at high frame rates just like CIDs.

Complementary metal-oxide semiconductor devices tend to have higher dark currents due to the highly doped silicon used. Therefore, CMOS sensors will not replace CCDs for low-noise scientific applications. Because the active devices often take up real estate, the area for the photosensor is reduced. This leads to reduced sensitivity that can be partially offset by a microlens. Because each pixel has its own amplifier, pattern noise is larger. However, more logic can be added to each pixel

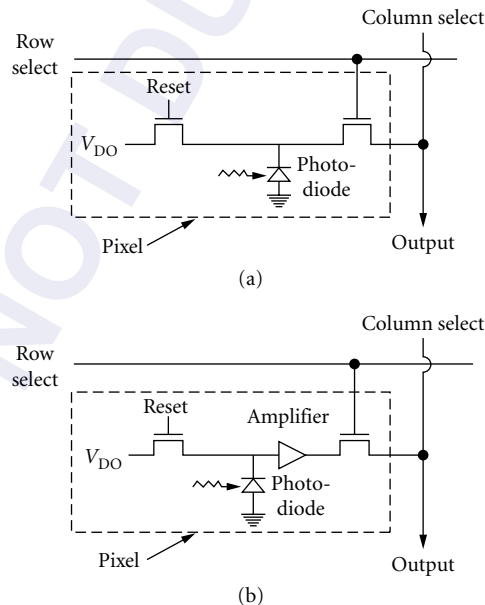


FIGURE 8 (a) Passive pixel device and (b) APS. Charge-coupled devices and CIDs are regarded as passive pixel sensors. Charge injection devices use photo-gates; CCDs use either photogates or photodiodes; CMOS devices typically use photodiodes.

for on-chip signal processing that suppresses pattern noise.¹⁸ With charge limited to a single pixel, it cannot overflow into neighboring pixels and create blooming as seen with CCDs.

A further advantage of APS is its low power consumption. It can operate from a 5-V battery or less (compared with 10 to 15 V for a CCD). CMOS may compete with CCDs in the general video marketplace where weight, size, and power consumption are deciding factors.

26.7 ARRAY PERFORMANCE

The most common array performance measures are responsivity, read noise, and charge well capacity. From these the minimum signal, maximum signal, SNR, and dynamic range can be calculated. Full characterization includes quantifying the various noise sources, charge transfer efficiency, spectral quantum efficiency, linearity, and pixel nonuniformity.¹⁹ These additional metrics are necessary for the most critical scientific applications.

Signal

Device specifications depend, in part, upon the application. Arrays for general video applications may have responsivity expressed in units of volts per lux. For scientific applications, the units may be in $V/(J\text{ cm}^{-2})$ or, if a digital output is available, $DN/(J\text{ cm}^{-2})$ where DN refers to a digital number. For example, in an 8-bit system, the digital numbers range from 0 to 255. These units are incomplete descriptors unless the device spectral response and source spectral characteristics are furnished.

The number of photoelectrons created by a detector is

$$n_{\text{PE}} = A_D \int_{\lambda_1}^{\lambda_2} E_{q\text{-faceplate}}(\lambda) \eta(\lambda) t_{\text{int}} d\lambda \quad (1)$$

where A_D is the effective photosensitive area, t_{INT} is the integration time, $\eta(\lambda)$ is the spectral quantum efficiency, and $E_{q\text{-faceplate}}(\lambda)$ is the spectral photon incidence in unit of photons $\text{s}^{-1} \mu\text{m}^{-1} \text{m}^{-2}$ (discussed further in Sec. 26.8). Some arrays have a transparent window protecting the array. The faceplate is the front surface of that window. With this approach, the array quantum efficiency includes the transmittance of the window.

For CCDs, charge is converted to a voltage by a floating diode or floating diffusion. The diode, acting as a capacitor, is precharged at a reference level. The capacitance, or sense node, is partially discharged by the amount of negative charge transferred. The difference in voltage between the final status of the diode and its precharged value (reset level) is linearly proportional to the number of photoelectrons. The signal voltage after the source follower is

$$V_{\text{signal}} = V_{\text{reset}} - V_{\text{out}} = n_{\text{PE}} \frac{qG}{C} \quad (2)$$

The gain, G , of a source follower amplifier is approximately unity, and q is the electronic charge (1.6×10^{-19} C). The charge conversion is q/C . The output gain conversion is qG/C . It typically ranges from 0.1 to 10 $\mu\text{V}/e^-$. The signal is then amplified and processed by electronics external to the CCD sensor.

Responsivity The spectral quantum efficiency is important to the scientific and military communities. When the array is placed into a general video or industrial camera, it is convenient to specify the output as a function of incident flux density or energy density averaged over the spectral response of the array.

The faceplate spectral photon incidence can be converted into the spectral radiant incidence by

$$E_{e\text{-faceplate}}(\lambda) = \frac{hc}{\lambda} E_{q\text{-faceplate}}(\lambda) \frac{W}{\mu\text{m m}^2} \quad (3)$$

where h is Planck's constant (6.626×10^{-34} J s⁻¹) and c is the speed of light (3×10^8 m s⁻¹). Quantities associated with power (watts) have the subscript e , and those associated with photons have the subscript q . The quantum efficiency can be converted to amperes per watt by

$$\mathfrak{R}_e(\lambda) = \frac{q\lambda}{hc} \eta(\lambda) \quad (4)$$

Then, the array output voltage (after the source follower amplifier) is

$$V_{\text{signal}} = \frac{G}{C} A_D \int_{\lambda_1}^{\lambda_2} E_{e\text{-faceplate}}(\lambda) \mathfrak{R}_e(\lambda) t_{\text{int}} d\lambda \quad (5)$$

It is desirable to express the responsivity in the form

$$V_{\text{signal}} = \mathfrak{R}_{\text{ave}} \left[\int_{\lambda_1}^{\lambda_2} E_{e\text{-faceplate}}(\lambda) t_{\text{int}} d\lambda \right] \quad (6)$$

The value $\mathfrak{R}_{\text{ave}}$ is an average responsivity that has units of V/J cm² and the quantity in the brackets has units of J/cm². Combining the two equations provides

$$\mathfrak{R}_{\text{ave}} = \frac{G}{C} A_D \frac{\int_{\lambda_1}^{\lambda_2} E_{e\text{-faceplate}}(\lambda) \mathfrak{R}_e(\lambda) d\lambda}{\int_{\lambda_1}^{\lambda_2} E_{e\text{-faceplate}}(\lambda) d\lambda} \text{ V/(J cm}^{-2}\text{)} \quad (7)$$

The value $\mathfrak{R}_{\text{ave}}$ is an average responsivity that depends upon the source characteristics and the spectral quantum efficiency. While the source can be standardized (e.g., CIE illuminant A or illuminant D₆₅₀₀), the spectral quantum efficiency varies by device. Therefore, extreme care must be exercised when comparing devices solely by the average responsivity.

If a very small wavelength increment is selected, $E_{e\text{-faceplate}}(\lambda)$ and $\mathfrak{R}_e(\lambda)$ may be considered as constants and Eq. (7) can be approximated as

$$\mathfrak{R}_{\text{ave}} = \frac{G}{C} A_D \mathfrak{R}_e(\lambda_0) \quad (8)$$

Minimum Signal The noise equivalent exposure (NEE) is an excellent diagnostic tool for production testing to verify noise performance. NEE is a poor array-to-array comparison parameter and should be used cautiously when comparing arrays with different architectures. This is so because it depends on array spectral responsivity and noise. The NEE is the exposure that produces a SNR of one. If the measured root-mean-square (rms) noise on the analog output is V_{noise} , then the NEE is calculated from the radiometric calibration:

$$\text{NEE} = \frac{V_{\text{noise}}}{\mathfrak{R}_{\text{ave}}} \text{ J/cm}^2 \text{ rms} \quad (9)$$

When expressed in electrons, NEE is simply the noise value in rms electrons. The absolute minimum noise level is the noise floor, and this value is used most often to calculate the NEE. Although noise is an rms value, the notation "rms" is often omitted.

Maximum Signal The maximum signal is that input signal that saturates the charge well and is called the saturation equivalent exposure (SEE). It is

$$\text{SEE} = \frac{V_{\text{max}}}{\mathfrak{R}_{\text{ave}}} \text{ J/cm}^2 \quad (10)$$

The total number of electrons that can be stored is the well capacity, n_{well} . The well size varies with architecture, number of phases, and pixel size. It is approximately proportional to pixel area. Small

pixels have small wells. If an antibloom drain is present, the maximum level is taken as the white clip level. The maximum signal is

$$V_{\max} = \frac{qG}{C}(n_{\text{well}} - n_{\text{dark}}) \quad (11)$$

For back-of-the-envelope calculations, the dark current is often considered negligible ($n_{\text{dark}} \approx 0$).

Dynamic Range The array dynamic range is

$$\text{DR}_{\text{array}} = \frac{n_{\text{well}} - n_{\text{dark}}}{\langle n_{\text{sys}} \rangle} \quad (12)$$

where $\langle n_{\text{sys}} \rangle$ is the overall rms noise. The array noise consists of dark current, shot, pattern, and readout noise (noise floor). In the absence of light and with negligible dark current, the dynamic range is most often quoted as

$$\text{DR}_{\text{array}} = \frac{n_{\text{well}}}{\langle n_{\text{floor}} \rangle} \quad (13)$$

Noise

Many books^{4–8} and articles^{20–23} have been written on noise sources. The level of detail used in noise modeling depends on the application. The noise sources include shot noise, reset noise, pattern noise, on-chip amplifier noise, and quantization noise. It is customary to specify all noise sources in units of equivalent rms electrons at the detector output.

Reset noise can be reduced to a negligible level with correlated double sampling (CDS). CDS also reduces the source follower $1/f$ noise. The off-chip amplifier is usually a low-noise amplifier such that its noise is small compared with the on-chip amplifier noise. The use of an analog-to-digital converter that has more bits reduces quantization noise.

On-chip amplifier noise may be called readout noise, mux noise, noise-equivalent electrons, or the noise floor. The value varies by device and manufacturer. For most system analyses, it is sufficient to consider

$$\langle n_{\text{sys}} \rangle = \sqrt{\langle n_{\text{shot}}^2 \rangle + \langle n_{\text{floor}}^2 \rangle + \langle n_{\text{pattern}}^2 \rangle} \quad (14)$$

where $\langle n^2 \rangle$ is the noise variance and $\langle n \rangle$ is the standard deviation measured in rms electrons.

Shot Noise Both photoelectrons and dark current contribute to shot noise. These follow Poisson statistics so that the variance is equal to the mean:

$$\langle n_{\text{shot}}^2 \rangle = n_{\text{PE}} + n_{\text{dark}} \quad (15)$$

While the dark current average value can be subtracted from the output to provide only the signal due to photoelectrons, the dark current noise cannot. Cooling the array can reduce the dark current to a negligible value and thereby reduce dark current noise to a negligible level.

Pattern Noise Pattern noise refers to any spatial pattern that does not change significantly from frame to frame. Pattern noise is not noise in the usual sense. This variation appears as spatial noise to the observer. Fixed-pattern noise (FPN) is caused by pixel-to-pixel variations in dark current.^{24,25} As a signal-independent noise, it is additive to the other noise powers. Fixed-pattern noise is due to differences in detector size, doping density, and foreign matter getting trapped during fabrication.

Photoresponse nonuniformity (PRNU) is the variation in pixel-to-pixel responsivities and, as such, is a signal-dependent noise. This noise is due to differences in detector size, spectral response, and thickness in coatings. Photoresponse nonuniformity can be specified as a peak-to-peak value or an rms value referenced to an average value. This average value may either be full well or one-half

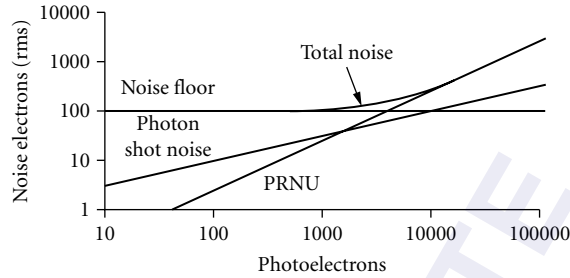


FIGURE 9 Photon transfer curve when $U = 2.5$ percent. The charge well capacity is 100,000 electrons and the noise floor is $100 e^-$ rms to produce a dynamic range of 60 dB. The noise floor, photon shot noise, and PRNU have slopes of 0, 0.5, and 1, respectively. Dark noise is considered negligible.

full well value. That is, the array is uniformly illuminated and a histogram of responses is created. The PRNU can be the rms of the histogram divided by the average value or the peak-to-peak value divided by the average value. The definition varies by manufacturer so that the test conditions must be understood when comparing arrays.

Because dark current becomes negligible when the array is sufficiently cooled, PRNU is the dominant pattern component for most arrays. As a multiplicative noise, PRNU is traditionally expressed as a fraction of the total number of charge carriers. If U is the fixed pattern ratio or nonuniformity, then

$$\langle n_{\text{pattern}} \rangle \approx \langle n_{\text{PRNU}} \rangle = U n_{\text{PE}} \quad (16)$$

Frame averaging will reduce all the noise sources except FPN and PRNU. Although FPN and PRNU are different, they are sometimes collectively called scene noise, pixel noise, pixel nonuniformity, or simply pattern noise.

Photon Transfer For many applications it is sufficient to consider photon shot noise, noise floor, and PRNU. The simplified noise model provides

$$\langle n_{\text{sys}} \rangle = \sqrt{n_{\text{PE}}^2 + \langle n_{\text{floor}}^2 \rangle + (U n_{\text{PE}})^2} \quad (17)$$

Recall that the mean square photon fluctuation is equal to the mean photon rate. Either the rms noise or noise variance can be plotted as a function of signal level. The graphs are called the photon transfer curve and the mean-variance curve, respectively. Both graphs convey the same information.

For very low photon fluxes, the noise floor dominates. As the incident flux increases, the photon shot noise dominates. For very high flux levels, the noise may be dominated by PRNU. Figure 9 illustrates the rms noise as a function of photoelectrons when the dynamic range ($n_{\text{well}} / \langle n_{\text{floor}} \rangle$) is 60 dB. With large signals and small PRNU, the total noise is dominated by photon shot noise. When PRNU is large, U dominates the array noise at high signal levels. General video and industrial cameras tend to operate in high-signal environments, and cooling will have little effect on performance. A full SNR analysis is required before selecting a cooled camera.

26.8 CAMERA PERFORMANCE

Camera performance metrics are conceptually the same as array metrics. The camera is limited by the array noise and charge well capacity. The camera's FPN and PRNU may be better than the array pattern noise when correction algorithms are present. Frame averaging and binning can reduce the

random noise floor and thereby appear to increase the camera's dynamic range. If the camera does not introduce any additional noise, modify the noise bandwidth, or minimize the maximum output, the camera SNR and DR_{camera} will be identical to the array values.

Camera Formula

The number of photoelectrons created by an object is

$$n_{\text{PE}} = \int_{\lambda_1}^{\lambda_2} \eta(\lambda) \frac{M_q(\lambda)}{4F^2(1+m_{\text{optics}})^2} A_D t_{\text{int}} \tau_{\text{optics}}(\lambda) d\lambda \quad (18)$$

where $M_q(\lambda)$ is the object's spectral photon exitance in photons/s μm^2 and $\tau_{\text{optics}}(\lambda)$ is the lens system transmittance. The f -number has the usual definition ($F = f/D$), and the optical magnification is $m_{\text{optics}} = R_2/R_1$. Here, R_1 and R_2 are related to the system's effective focal length, f , by

$$\frac{1}{R_1} + \frac{1}{R_2} = \frac{1}{f} \quad (19)$$

As the target moves to infinity ($R_1 \rightarrow \infty$), m_{optics} approaches zero.

Electronic still cameras are matched to conventional photographic cameras. In photography, shutter speeds (exposure times) vary approximately by a factor of two (e.g., 1/30, 1/60, 1/125, 1/250, etc.). Thus, changing the shutter speed by one setting changes n_{PE} approximately by a factor of 2. F-stops have been standardized to 1, 1.4, 2, 2.8, 4, 5.6, 8, The ratio of adjacent F-stops is $\sqrt{2}$. Changing the lens speed by one F-stop changes the f -number by a factor of $\sqrt{2}$. Here, also, the n_{PE} changes by a factor of 2.

The measurement of faceplate illumination [see Eq. (1)] is usually performed with a calibrated lens. The value M_q is measured with a calibrated radiometer or photometer, and then $E_{q\text{-faceplate}}$ is calculated according to

$$E_{q\text{-faceplate}} = \frac{M_q}{4F^2(1+m_{\text{optics}})^2} \tau_{\text{optics}} \quad (20)$$

Minimum Signal

The maximum and minimum signals depend on the spectral output of the source and the spectral response of the detector. The source color temperature is not always listed but is a critical parameter for comparing systems. Although the CIE illuminant A is used most often, the user should not assume that this was the source used by the camera manufacturer.

Based on signal detection theory, the minimum illumination would imply that the SNR is one. However, the definition of minimum illumination is manufacturer dependent. Its value may be (a) when the video signal is, for example, 30 IRE units, (b) when the SNR is one, or (c) when an observer just perceives a test pattern. Because of its incredible temporal and spatial integration capability, the human visual system can perceive SNRs as low as 0.05. Therefore, comparing cameras based on "minimum" illumination should be approached with care.

The voltage signal [Eq. (5)] exists at the output of the array. This voltage must be amplified by the camera electronics (gain = G_{camera}) to a value that is consistent with video standards. The minimum signal provided with gain "on" (G_{camera} greater than one) is usually calculated due to the difficulty of performing accurate, low-level radiometric and photometric measurements. These values may be provided at 30, 50, or 80 percent video levels. That is, the illumination that is given produces an output video that gives 30, 50, or 80 IRE units, respectively. Although a higher-gain amplifier could provide 100 percent video, the user can optimize the image by adjusting the gain and level of the display. That is, the display's internal amplifiers can be used for additional gain. In this context, the camera system consists of the camera and display.

If G_{camera} is expressed in decibels, it must be converted to a ratio. The scene illumination that creates a 30 percent video signal is

$$M_v(30\% \text{ video}) = \frac{0.3M_v(\text{max video})}{G_{\text{camera}}} \quad (21)$$

where $M_v(\text{max video})$ is the scene illumination that produces the maximum output. The subscript v indicates photometric units, which are usually lux. For 50 and 80 percent video, the factor becomes 0.5 and 0.8, respectively. The input signal that produces a SNR of one is

$$M_v(\text{SNR}=1) = M_v(\text{max video}) 10^{-\frac{DR_{\text{camera-dB}}}{20}} \quad (22)$$

where the camera's dynamic range is expressed in decibels (dB). Although photometric units are used most often, radiometric quantities can be used in Eqs. (21) and (22).

Dynamic range depends on integration time, binning, and frame integration. In addition, the spectral content of the source and the array spectral responsivity affect the camera output voltage. Thus, the camera dynamic range can be quite variable depending on test conditions.

Although the theoretical maximum signal just fills the charge wells, the manufacturer may limit the maximum output voltage to a lower stored-charge value. Because the dynamic range is often expressed in decibels,

$$DR_{\text{camera-dB}} = 20 \log \left(\frac{V_{\text{max}}}{V_{\text{noise}}} \right) \quad (23)$$

Many camera manufacturers list the dynamic range as the signal-to-noise ratio. This value should not be confused with the actual SNR. With most cameras, the noise level is approximately equivalent to the array read noise. The electronics may increase the camera noise and image-processing techniques may reduce it somewhat. Usually the dynamic range is calculated from the measured signal and noise. Because the read noise is amplifier white noise, it is appropriate to include the bandwidth as part of the measurements. For standard video-compatible cameras, the bandwidth is equal to the video format bandwidth.

26.9 MODULATION TRANSFER FUNCTION

The MTF of the camera is the product of all the subsystem MTFs. Usually, the electronics MTF does not significantly affect the overall MTF. Therefore, it is often adequate to consider only the optics and detector MTFs.

The MTF of a single rectangular detector in image space is

$$\text{MTF}_{\text{detector}}(u) = \left| \frac{\sin(\pi d \eta)}{\pi d \eta} \right| \quad (24)$$

where d is the horizontal extent of the photosensitive surface and η is the image-space horizontal spatial frequency variable (in cycles/mm). The detector size may be different in the horizontal and vertical directions, resulting in different MTFs in the two directions.

The MTF for a circular, clear-aperture, diffraction-limited lens is

$$\text{MTF}_{\text{optics}}(\eta) = \frac{2}{\pi} \left[\cos^{-1} \left(\frac{\eta}{\eta_c} \right) - \frac{\eta}{\eta_c} \sqrt{1 - \left(\frac{\eta}{\eta_c} \right)^2} \right] \quad (25)$$

The optical cutoff is $\eta_c = D/(\lambda_{\text{ave}} \text{fl})$, where D is the aperture diameter and λ_{ave} is the average wavelength.

26.10 RESOLUTION

Detector arrays are specified by the number of pixels, detector size, and detector pitch. These are not meaningful until an optical system is placed in front of the array. The most popular detector resolution measure is the detector angular subtense. In object space, it is

$$\text{DAS} = \frac{d}{f_l} \text{ mrad} \quad (26)$$

and in image space, the detector size, d , becomes the measure of resolution.

Perhaps the most popular measure of optical resolution is Airy disk size. It is the bright central spot of the diffraction pattern produced by an ideal optical system. In the focal plane of the lens, the Airy disk diameter is

$$d_{\text{airy}} = 2.44 \frac{\lambda}{D} f_l = 2.44 \lambda F \quad (27)$$

While often treated independently, the camera resolution depends upon both the optical and detector resolutions.

Shade created a metric for system performance. As reported by Lloyd,²⁶ Sendall modified Shade's equivalent resolution such that

$$R_{\text{eq}} = \frac{1}{2 \int_0^\infty |\text{MTF}_{\text{sys}}(\eta)|^2 d\eta} \quad (28)$$

As a summary metric, R_{eq} provides a better indication of system performance than a single metric, such as the detector size or blur diameter. R_{eq} cannot be directly measured. It is a mathematical construct simply used to express overall performance. As the MTF increases, R_{eq} decreases and the resolution "improves" (smaller is better). As an approximation, the system resolution may be estimated from the subsystem equivalent resolutions by

$$R_{\text{eq}} \approx \sqrt{R_{\text{optics}}^2 + R_{\text{detector}}^2} \quad (29)$$

Substituting the diffraction-limited MTF into Eq. (28) provides

$$R_{\text{optics}} = 1.845 \lambda F \quad (30)$$

Note that Shade's approach provides a value that is smaller than the Airy disk diameter. Recall that R_{eq} is only a mathematical construct used to analyze system performance. When R_{optics} dominates R_{eq} , we say the system is optics-limited.

Substituting the detector MTF into Eq. (28) provides

$$R_{\text{detector}} = d \quad (31)$$

Here, Shade's resolution matches the common method of describing detector performance: The smallest target that can be discerned is limited by the detector size. When R_{detector} dominates R_{eq} , we say the system is detector-limited.

Using Eq. (29) to estimate the composite resolution in image space,

$$R_{\text{eq}} \approx d \sqrt{\left(\frac{1.845 \lambda F}{d}\right)^2 + 1} \quad (32)$$

As $\lambda F/d$ decreases, R_{eq} approaches d . For large values of $\lambda F/d$, the system becomes optics-limited and the equivalent resolution increases.

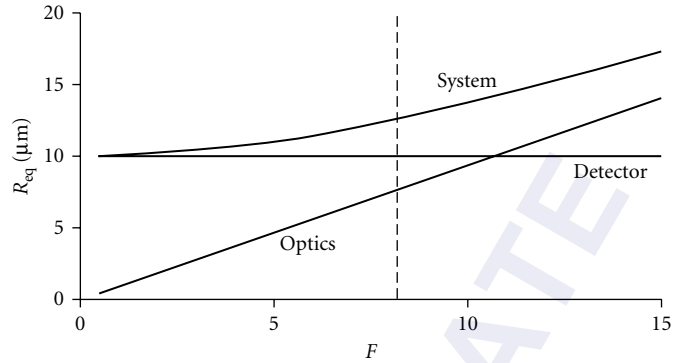


FIGURE 10 Resolution of a typical 1/2-inch-format CCD camera as a function of f -number. $\lambda = 0.5 \mu\text{m}$. The vertical line at $F = 8.20$ separates optics-limited from detector-limited operation. It occurs when the Airy disk size is equal to the detector size.

The more common 1/2-inch-format CCD arrays have detectors that are about $10 \mu\text{m}$ in size. Figure 10 illustrates R_{eq} as a function of f -number. Reducing the f -number below 5 does not improve resolution because the system is in the detector-limited region.

For most CCD camera applications, it is assumed that the camera is operating in the detector-limited region. This is only valid if $\lambda F/d$ is small. If $\lambda F/d$ is large, then the minimum discernable target size is definitely affected by the optics resolution.

We live in a world where “smaller is better.” Detector sizes are shrinking. This allows the system designer to create physically smaller cameras. Replacing a 1/2-in-format array with a 1/4-in-format array (typical detector size is $5 \mu\text{m}$) implies a $2\times$ improvement in resolution. However, this is only true if the system is operating in the detector-limited region. As d decreases, the f -number must also decrease to stay within the detector-limited region. Further, the f -number must also decrease to maintain the same signal intensity [Eq. (18)]. Reducing the f -number can place a burden on the optical designer.

26.11 SAMPLING

Sampling is an inherent feature of all electronic imaging systems. The scene is spatially sampled in both directions due to the discrete locations of the detector elements. The horizontal sample rate is

$$\eta_s = \frac{1}{d_{\text{CCH}}} \text{ cycles/mm} \quad (33)$$

Staring arrays can faithfully reproduce signals up to the Nyquist frequency:

$$\eta_N = \frac{\eta_s}{2} = \frac{1}{2d_{\text{CCH}}} \text{ cycles/mm} \quad (34)$$

The pitch in the horizontal and vertical directions, d_{CCH} and d_{CCV} , respectively, may be different and, therefore, the sampling rates will be different. Any input frequency above the Nyquist frequency will be aliased to a lower frequency (Fig. 11). After aliasing, the original signal can never be recovered. Diagonal lines appear to have jagged edges, or “jaggies,” and periodic patterns create moiré patterns. Periodic structures are rare in nature and aliasing is seldom reported when viewing natural scenery, although aliasing is always present. Aliasing may become apparent when viewing periodic targets such as test patterns, picket fences, plowed fields, railroad tracks, and Venetian blinds. It becomes bothersome when the scene geometric properties must be maintained as with mapping. It

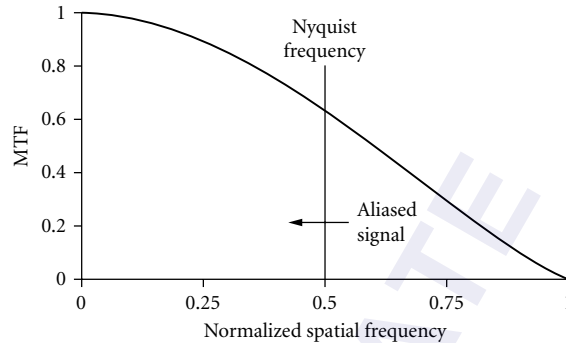


FIGURE 11 Signals above Nyquist frequency are aliased down to the base band. $\eta_s = 1$ and $\eta_N = 0.5$. The area bounded by the MTF above η_N may be considered as an “aliasing” metric.²⁷

affects the performance of most image-processing algorithms. While this is a concern for scientific and military applications, it typically is of little consequence to the average professional television broadcast and consumer markets.

We have become accustomed to the aliasing in commercial televisions. Periodic horizontal lines are distorted due to the raster. Cloth patterns, such as herringbones and stripes, produce moiré patterns. Cross-color effects occur in color imagery (red stripes may appear green or blue). Many videotape recordings are undersampled to keep the price modest, and yet the imagery is considered acceptable when observed at normal viewing distances.

Because the aliasing occurs at the detector, the signal must be band limited by the optical system to prevent it. Optical band limiting can be achieved by using small-diameter optics, blurring the image, or by inserting a birefringent crystal between the lens and array. The birefringent crystal changes the effective detector size and is found in almost all single-chip color cameras. Unfortunately, these approaches also degrade the MTF (reduce image sharpness) and are considered unacceptable for scientific applications.

If a system is Nyquist frequency-limited, then the Nyquist frequency is used as a measure of resolution. Because no frequency can exist above the Nyquist frequency, many researchers represent the MTF as zero above the Nyquist frequency (Fig. 12). This representation may be too restrictive for modeling purposes.

The fill factor can vary from 20 percent for interline transfer devices to nearly 100 percent for frame transfer devices. In the detector-limited region, the detector MTF determines the potential

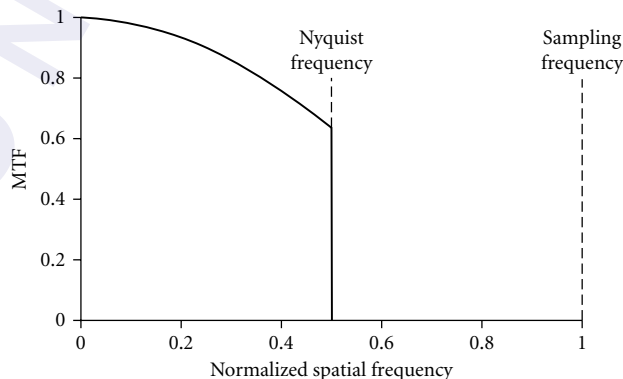


FIGURE 12 MTF representation of an undersampled system. $\eta_s = 1$ and $\eta_N = 0.5$. This restrictive representation erroneously suggests that there is no response above η_N .

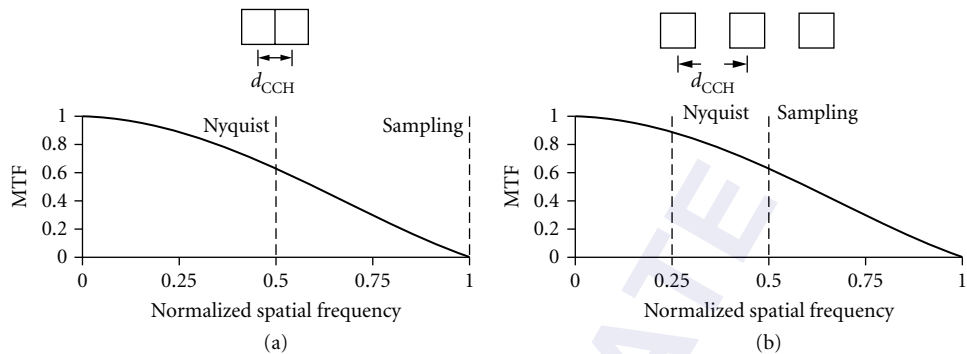


FIGURE 13 Two arrays with different horizontal center-to-center spacings. The detector size, d , is the same for both. (a) $d/d_{CCH} = 1$. This typifies frame transfer devices. (b) $d/d_{CCH} = 0.5$. This is representative of an ideal interline transfer CCD array. The detector MTF is plotted as a function of ηd .

spatial frequencies that can be reproduced. The center-to-center spacing uniquely determines the Nyquist frequency (Fig. 13). A microlens will increase the effective detector size, but it does not affect the center-to-center spacing. The absolute value of the Nyquist frequency ($1/2d_{CCH}$) does not change. By increasing the detector size, the detector cutoff decreases. The *relative* locations of the sampling and Nyquist frequencies change. The figures in this text use relative (normalized) frequency scales.

Less-expensive color cameras contain only a single CCD chip. A color filter array (CFA) is placed over the chip to create red, green, and blue pixels. Figure 20 in Chap. 32, “Visible Array Detectors,” in Vol. II, provides a variety of CFA patterns. In many sensors, the number of detectors that are devoted to each color is different. The basic reason is that the human visual system (HVS) derives its detail information primarily from the green portion of the spectrum. That is, luminance differences are associated with green, whereas color perception is associated with red and blue. The HVS requires only moderate amounts of red and blue to perceive color. Thus, many sensors have twice as many green as either red or blue detector elements. An array that has 768 horizontal elements may devote 384 to green, 192 to red, and 192 to blue. This results in an unequal sampling of the colors (Fig. 14).

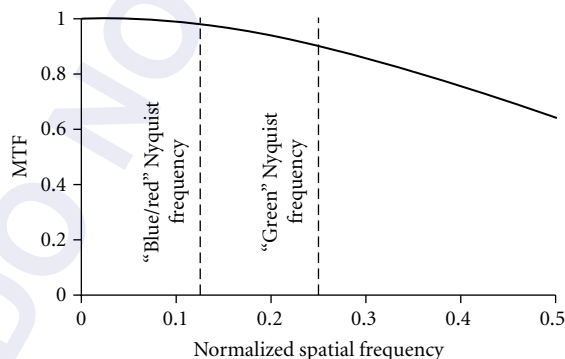


FIGURE 14 Horizontal detector MTF normalized to ηd . Unequally spaced red-, green-, and blue-sensitive detectors create different array Nyquist frequencies. This creates different amounts of aliasing and black-and-white scenes may break into color. The spacing between the red and blue detectors is twice the green detector spacing.

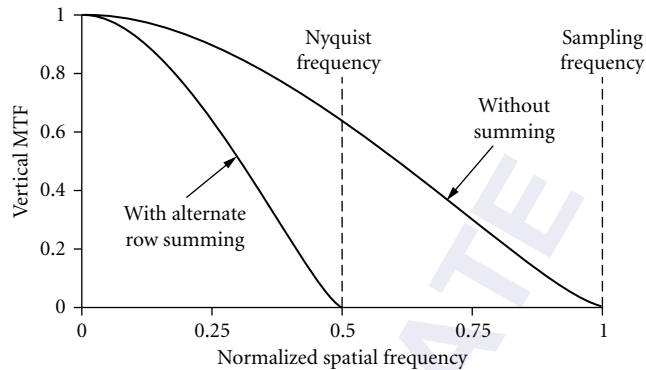


FIGURE 15 Alternate detector summing (pseudointerlacing) can reduce aliasing in the vertical direction. Vertical MTF with and without summing normalized to ηd .

The output R, G, and B signals are created by interpolation of the sparse data (sparse pixels). The output signals *appear* as if there are 768 red, green, and blue pixels. This interpolation does not change the Nyquist frequency of each color. A birefringent crystal²⁸ inserted between the lens and the array effectively increases the detector sizes. A larger detector will have reduced MTF and this reduces aliasing. It also reduces edge sharpness.

Pseudointerlacing (Fig. 4) doubles the size of the detector. This reduces the detector cutoff and makes it equal to the Nyquist frequency (Fig. 15). From an aliasing point of view, aliasing has been significantly reduced. However, the MTF has also been reduced, and this results in reduced-edge sharpness.

26.12 STORAGE, ANALYSIS, AND DISPLAY

Chapter 32, “Visible Array Detectors,” in Vol. II describes the physics of CCD detectors. This chapter has described some additional features of CCDs and introduced CID and CMOS detectors. These solid-state devices are the basic building blocks of the solid-state array. Certain array architectures lend themselves to specific applications. Once the camera is fabricated, the user selects an optical system. The minimum signal, maximum signal, and resolution discussed in this chapter include the lens f -number. It would appear that with all this knowledge, it would be easy to select a camera. A camera only becomes operational when its output is analyzed. Analysis is performed by an observer (general video) or a computer (machine vision).

For general video, the camera output must be formatted into a data stream consistent with the display device. The monochrome standard is often called EIA 170 (originally called RS 170) and the color format is simply known as NTSC (originally called EIA 170A or RS 170A). Worldwide, three color broadcast standards exist: NTSC, PAL, and SECAM. For higher vertical resolution, more lines are required. EIA 343A (originally RS 343A) is a high-resolution monochrome standard used for closed-circuit television cameras (CCTV). Although the standard encompasses equipment that operates from 675 to 1023 lines, the recommended values are 675, 729, 875, 945, and 1023 lines per frame.

These standards are commonplace. Monitors that display these formats are readily available. For computer analysis of this imagery, frame grabbers that accept multiple standards are easily obtained. The output of most general video cameras is an analog signal, and the frame grabber digitizes this signal for computer processing.

In principle, the clock rate of an analog-to-digital converter within the frame grabber can be set at any rate. However, to conserve on memory requirements and minimize clock rates, some frame grabbers tend to just satisfy the Nyquist frequency of a standard video signal. That is, if the highest frequency of the video bandwidth is f_{BW} , then the frame grabber sampling clock operates at $2f_{BW}$.

Some frame grabbers have an internal antialias filter. This filter ensures that the frame grabber does not produce any additional aliasing. The filter cutoff is linked to the frame grabber clock and is not related to the camera output. Once aliasing has occurred in the camera, the frame grabber antialias filter cannot remove it. If the filter does not exist, the frame grabber may create additional aliasing.²⁹ On the other hand, some frame grabbers have antialiasing filters that significantly limit the analog signal bandwidth and thereby reduce resolution. The number of digital samples is simply related to the frame grabber clock rate and is not necessarily equal to the number of detector elements. Even if the number of digital samples matches the number of detector elements, phasing effects will corrupt the resolution. Image-processing algorithms operate on the digitized signal and the image-processing specialist must be aware of the overall resolution of the system. The frame grabber is an integral part of that system. These issues are mitigated when the camera output is in digital form and the frame grabber can accept digital signals.

As society moves toward digital television [high-definition television (HDTV) or advanced television system (ATS)], new demands are placed upon displays and frame grabbers. New standards also require new video-recording devices. This creates storage, frame grabber, and display problems. New standards do not impose any difficulties on camera design. The solid-state array can be made any size and appropriately designed electronics can support any video format.

Industrial and scientific applications require some forethought. The camera output is no longer a conventional format. Web inspection cameras may contain a linear array with as many as 4096 elements. Split architecture (Fig. 1) may have as many as 32 parallel outputs. Charge injection devices and CMOS cameras may offer variable-sized subframes at variable frame rates. Finding suitable frame grabbers is a challenge for these applications. Nonstandard formats (even after digitization) are not easy to display.

26.13 REFERENCES

1. W. S. Boyle and G. E. Smith, "Charge Coupled Semiconductor Devices," *Bell Systems Technical Journal* **49**:587–593 (1970).
2. G. F. Amelio, M. F. Tompsett, and G. E. Smith, "Experimental Verification of the Charge Coupled Concept," *Bell Systems Technical Journal* **49**:593–600 (1970).
3. M. J. Howes and D. V. Morgan (eds.), *Charge-Coupled Devices and Systems*, John Wiley & Sons, New York, 1979.
4. C. H. Sequin and M. F. Tompsett, *Charge Transfer Devices*, Academic Press, New York, 1975.
5. E. S. Yang, *Microelectronic Devices*, McGraw-Hill, New York, 1988.
6. E. L. Dereniak and D. G. Crowe, *Optical Radiation Detectors*, John Wiley & Sons, New York, 1984.
7. A. J. P. Theuwissen, *Solid-State Imaging with Charge-Coupled Devices*, Kluwer Academic Publishers, Dordrecht, The Netherlands, 1995.
8. G. C. Holst, *CCD Cameras, Arrays, and Displays*, second edition, JCD Publishing, Winter Park, FL, 1998.
9. J. Janesick, T. Elliott, R. Winzenread, J. Pinter, and R. Dyck, "Sandbox CCDs," *Charge-Coupled Devices and Solid State Optical Sensors V*, M. M. Blouke, ed., *Proc. SPIE* Vol. 2415, 1995, pp. 2–42.
10. J. Janesick and T. Elliott, "History and Advancement of Large Area Array Scientific CCD Imagers," *Astronomical Society of the Pacific Conference Series, Vol. 23, Astronomical CCD Observing and Reduction*, BookCrafters, 1992, pp. 1–39.
11. M. Kimata and N. Tubouchi, "Charge Transfer Devices," *Infrared Photon Detectors*, A. Rogalski, ed., SPIE Press, Bellingham, WA, 1995, pp. 99–144.
12. S. G. Chamberlain, S. R. Kamasz, F. Ma, W. D. Washkurak, M. Farrier, and P. T. Jenkins, "A 26.3 Million Pixel CCD Image Sensor," *IEEE Proceedings of the International Conference on Electron Devices*, pp. 151–155, Washington, D.C., December 10, 1995.

13. G. C. Holst, *CCD Cameras, Arrays, and Displays*, second edition, JCD Publishing, Winter Park, FL, 1998, pp. 338–370.
14. S. K. Mendis, S. E. Kemeny, and E. R. Fossum, “A 128×128 CMOS Active Pixel Sensor for Highly Integrated Imaging Systems,” *IEEE IEDM Technical Digest* 583–586 (1993).
15. E. R. Fossum, “CMOS Image Sensors: Electronic Camera-on-a-Chip,” *IEEE Transactions on Electron Devices* **44**(10):1689–1698 (1997).
16. Z. Zhou, B. Pain, and E. R. Fossum, “Frame-Transfer CMOS Active Pixel Sensor with Pixel Binning,” *IEEE Transactions on Electron Devices* **44**(10):1764–1768 (1997).
17. E. R. Fossum, “Active Pixel Sensors: Are CCD’s Dinosaurs?,” *Charge-Coupled Devices and Solid State Optical Sensors III*, M. M. Blouke, ed., *Proc. SPIE* Vol. 1900, 1993, pp. 2–14.
18. J. Nakamura, T. Nomoto, T. Nakamura, and E. R. Fossum, “On-Focal-Plane Signal Processing for Current-Mode Active Pixel Sensors,” *IEEE Transactions on Electron Devices* **44**(10):1747–1757(1997).
19. G. C. Holst, *CCD Cameras, Arrays, and Displays*, second edition, JCD Publishing, Winter Park, FL, 1998, pp. 102–144.
20. D. G. Crowe, P. R. Norton, T. Limperis, and J. Mudar, “Detectors,” *Electro-Optical Components*, W. D. Rogatto, pp. 175–283. Volume 3 of *The Infrared & Electro-Optical Systems Handbook*, J. S. Accetta and D. L. Shumaker, eds., copublished by Environmental Research Institute of Michigan, Ann Arbor, MI, and SPIE Press, Bellingham, WA, 1993.
21. J. R. Janesick, T. Elliott, S. Collins, M. M. Blouke, and J. Freeman, “Scientific Charge-Coupled Devices,” *Optical Engineering* **26**(8):692–714 (1987).
22. T. W. McCurnin, L. C. Schooley, and G. R. Sims, “Charge-Coupled Device Signal Processing Models and Comparisons,” *Journal of Electronic Imaging* **2**(2):100–107 (1994).
23. M. D. Nelson, J. F. Johnson, and T. S. Lomheim, “General Noise Process in Hybrid Infrared Focal Plane Arrays,” *Optical Engineering* **30**(11):1682–1700 (1991).
24. J. M. Mooney, “Effect of Spatial Noise on the Minimum Resolvable Temperature of a Staring Array,” *Applied Optics* **30**(23):3324–3332 (1991).
25. J. M. Mooney, F. D. Shepherd, W. S. Ewing, J. E. Murguia, and J. Silverman, “Responsivity Nonuniformity Limited Performance of Infrared Staring Cameras,” *Optical Engineering* **28**(11):1151–1161(1989).
26. J. M. Lloyd, *Thermal Imaging*, Plenum Press, New York, 1975, p. 109.
27. G. C. Holst, *Sampling, Aliasing, and Data Fidelity*, JCD Publishing, Winter Park, FL, 1998, pp. 293–320.
28. J. E. Greivenkamp, “Color Dependent Optical Prefilter for the Suppression of Aliasing Artifacts,” *Applied Optics* **29**(5):676–684 (1990).
29. G. C. Holst, *Sampling, Aliasing, and Data Fidelity*, JCD Publishing, Winter Park, FL, 1998, pp. 154–156.

This page intentionally left blank.

DO NOT DUPLICATE

CAMERA LENSES

Ellis Betensky

*Opcon Associates, Inc.
West Redding, Connecticut*

Melvin H. Kreitzer and Jacob Moskovich

*Opcon Associates, Inc.
Cincinnati, Ohio*

27.1 INTRODUCTION

Camera lenses have been discussed in a large number of books and articles. The approach in this chapter is to concentrate on modern types and to describe imaging performance in detail both in terms of digital applications and in terms of the optical transfer function. By modern types, we mean lens forms that were found on cameras in 1992. The chapter deals almost entirely with lenses for the 35-mm (24×36 mm) format. This limitation is unfortunate but not really inappropriate, given the widespread use of this format. Moreover, the different lens types that are described are used for applications ranging from 8-mm video to 6×9 cm roll film.

We have not included any specific design examples of lenses for large-format cameras, but the imaging capabilities of these lenses are described in terms of digital applications. By digital applications we mean the comparison of different lens types in terms of total pixels and pixels per unit solid angle. It is hoped that this feature will make comparisons between radically different imaging systems possible and also help to classify lenses in terms of information capability. See “Further Reading” at the end of this chapter for related information about photographic lenses, particularly with respect to older design types.

27.2 IMPOSED DESIGN LIMITATIONS

There are some limitations that are imposed on the design of camera lenses. The most significant ones are listed as follows.

Microprism focusing in single lens reflex cameras (SLRs) is difficult at apertures smaller than about $F/4.5$. Recent advances permit the use of microprisms at apertures down to $F/5.6$ and this is usually the smallest maximum aperture permitted in the specification of a lens for the SLR camera.

Depending on the camera type, there is a maximum rear lens opening allowable at the flange on SLR lenses. The limitation is approximately 33 to 36-mm diameter at flange to film plane distances of 40.5 to 46 mm. This affects the maximum possible aperture on normal lenses (typically to $F/1.2$) and also requires appropriate design of the exit pupil location on long-focal-length and high-speed retrofocus lenses in order to avoid excessive vignetting.

The minimum back focal length (BFL) allowable on SLR lenses (because of the swinging mirror) is about 38.5 mm. The BFL cannot be too short on non-SLRs because of in-focus dust or cosmetic problems on optical surfaces close to the film plane. The actual limitation depends on the minimum relative aperture that would be used but is rarely less than 4 mm and usually more than 8 mm.

Since most lens accessories such as filters and lens-shades are mounted on the front of a lens, there is a practical limitation to the allowable front diameter of most lenses. Filter sizes larger than 72 mm are not desirable, and smaller is always preferred. The actual clear aperture at the front of a lens is considerably smaller than the filter size, depending on the angular field and the mounting details of the filter. Obviously there are lenses such as 600-mm F/4 telephotos for which the 72-mm limitation is not possible. In these cases, the lens can be designed to use internal filters that are incorporated into the design.

Mechanical cams are still in widespread use for the practical realization of the required motions in zoom lenses. This technology requires that the motions themselves be controlled at the design stage to be reasonably monotonic and often to have certain mutual relationships. These requirements are particularly severe for the so-called “one-touch” zoom and focus manual control found on many SLR zoom lenses.

In general, size and weight restrictions pose the biggest problems for the designer of most camera lenses. Almost any lens can be designed if there are no physical limitations. These limitations are sometimes a consequence of ergonomic considerations but can equally be an effort to achieve a marketing advantage. Size restrictions almost always adversely affect the design, and exceptionally small lenses (for a given specification) should be regarded with suspicion.

27.3 MODERN LENS TYPES

Normal (with Aspherics) and Variations

Thirty-five-mm SLR normal lenses are invariably Double-Gauss types. Refer to Fig. 1. This lens form is characterized by symmetry about a central stop to facilitate the correction of coma, distortion, and lateral color. These lenses are relatively easy to manufacture and a user can expect good quality in a production lens. Total angular coverage of about 45° is typical, and speeds as fast as F/1 are achievable. Extremely good optical performance is possible, particularly if the angular field and speed are reduced somewhat. Image quality generally deteriorates monotonically from axis to corner and improves dramatically as the lens aperture is reduced by about two F-numbers. With the addition of a fixed rear group, conjugate stability can be achieved over a wide range. Refer to Fig. 2.

Wide-Angle

An interesting new wide-angle lens type is a four-component form found commonly on the so-called compact 35-mm cameras. This lens is characterized by a triplet construction followed by a rear element that is strongly meniscus-shaped, convex to the image plane. This lens has much less astigmatism than either conventional triplets or Tessars and can cover total fields of up to 75° at speeds of around F/4. Faster speeds are possible if the angular field is reduced. Most importantly, the rear meniscus component takes the burden of field flattening away from the triplet front part. This results in considerably lower individual element powers and correspondingly lower sensitivities to tilts and decentrations of the elements. It is this problem that makes conventional triplets extremely difficult to manufacture. Refer to Fig. 3.

Inverted Telephoto (Retrofocus)

These lens types, characterized by a long back focal length, are typically used for wide-angle applications for single lens reflex cameras having a swinging viewing mirror behind the lens. Inverted telephoto implies a front negative group followed by a rear positive group, just the reverse of a

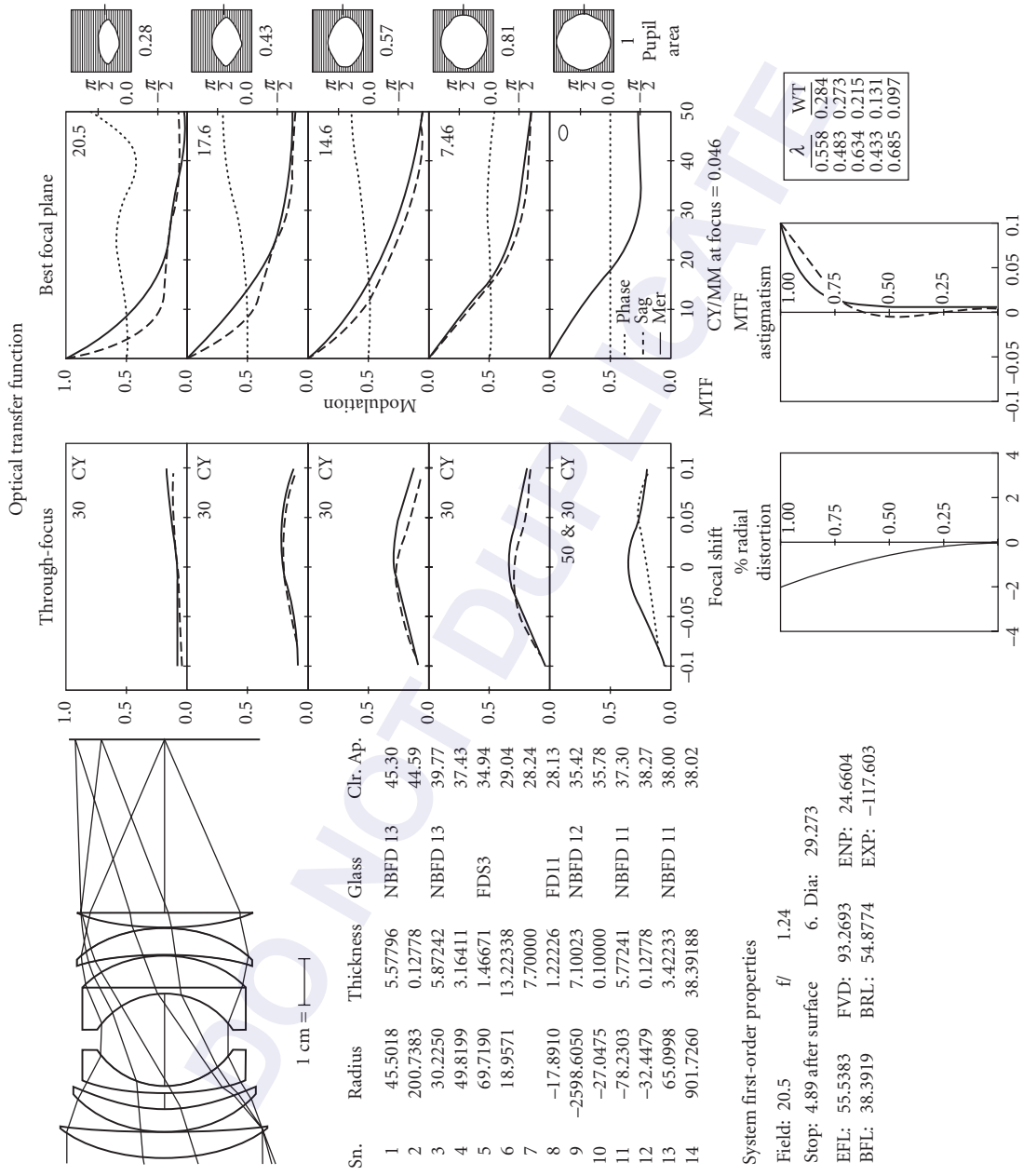
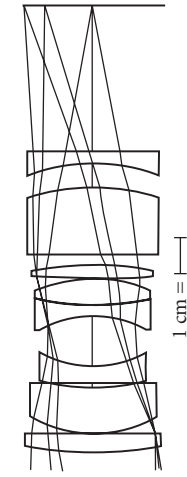


FIGURE 1 55-mm F/1.2 for 35-mm SLR.



Sn.	Radius	Thickness	Glass	Clr. Ap.
1	86.5052	5.00000	LAC12	36.91
2	-4873.6372	0.23400	LAC12	35.60
3	39.3067	12.52362	LAC12	34.24
4	81.7172	1.40000	FD5	28.81
5	258.5978	4.20000	FD5	28.46
6	29.4732	13.78955	FL7	25.49
7	-30.5568	3.00000	LAC9	25.48
8	87.9238	7.00000	LAC9	29.40
9	-42.1099	0.23400	LAC12	31.02
10	225.5319	3.60500	LAC12	32.70
11	-96.2427	3.02000	focus	33.02
12	-810.0813	17.82698	BSC7	34.35
13	-62.7744	6.66758	BSC7	35.48
14	-50.4393	3.60000	BSC7	34.25
15	-362.5221	35.45429	BSC7	34.98

Focus thicknesses:

position	τ_{11}
1	3.0200
2	45.9538

System first-order properties

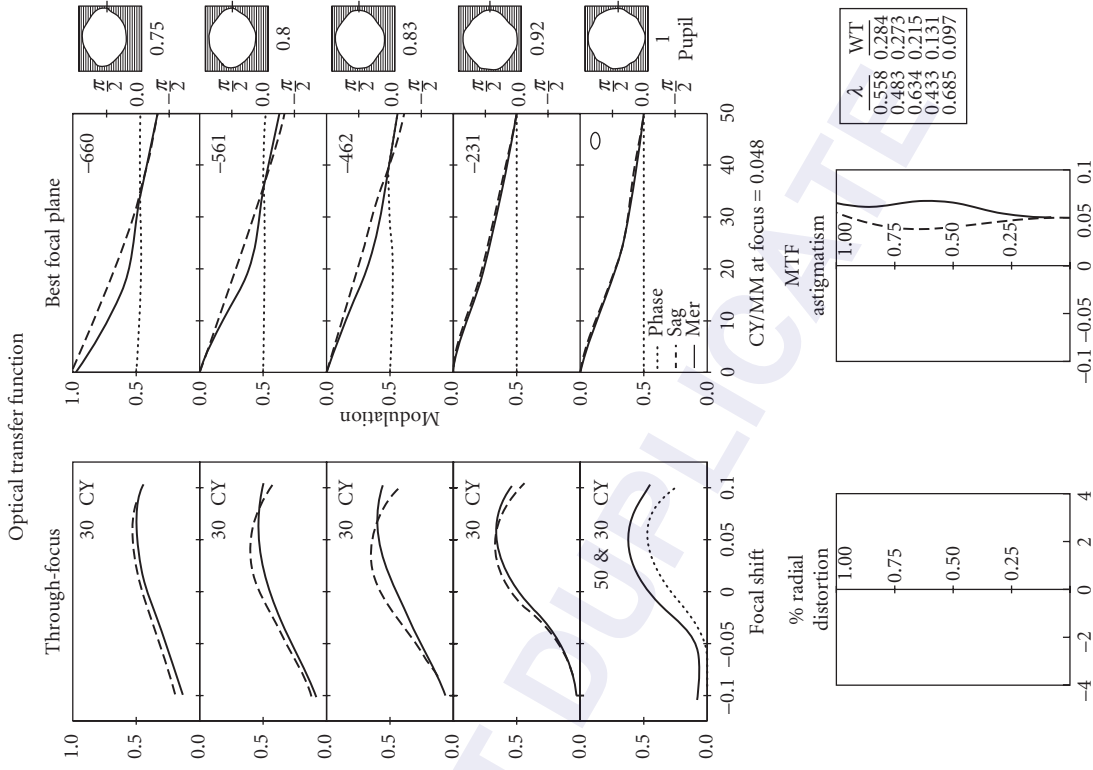
Stop: 6.89 after surface 7. Dia: 24.858

Pos 1 (Evaluation shown)

OBJ. HT:	-660.00	f	2.56	MAG:	-0.0300
EFL:	89.8128	FVD:	126.519	ENP:	39.0067
IMD:	35.4542	BRL:	91.0651	EXP:	-44.5326
OBD:	-3059.03	OVL:	3185.55		

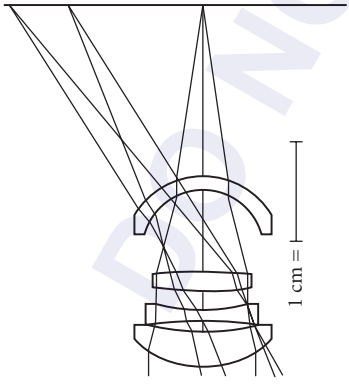
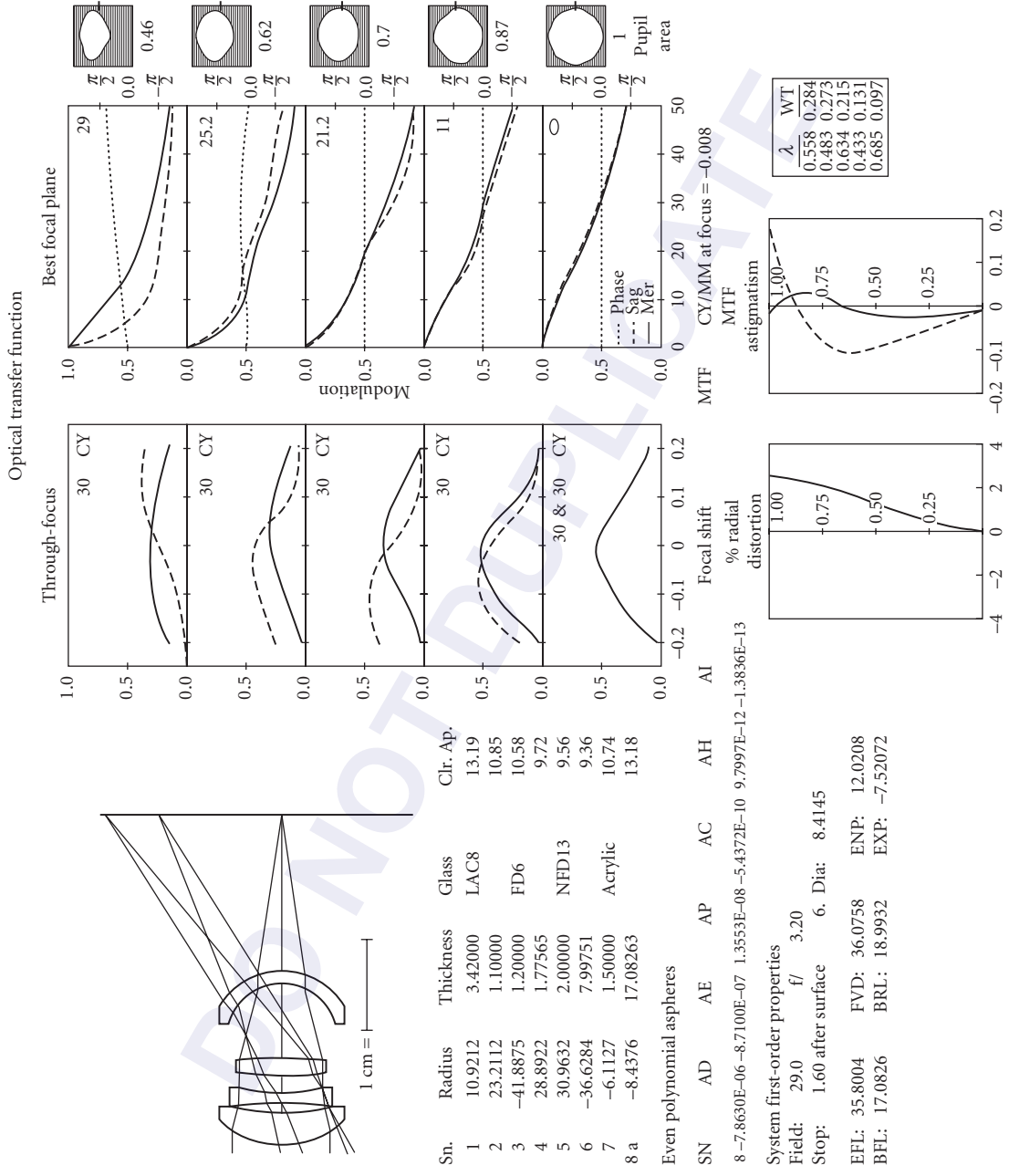
Pos 2

OBJ. HT:	-40.000	f	3.80	MAG:	-0.5000
EFL :	86.9535	FVD:	169.449	ENP:	39.0087
IMD:	35.4505	BRL:	133.999	EXP:	-80.5517
OBD:	-238.980	OVL:	408.429		



λ	WT
0.558	0.284
0.483	0.273
0.634	0.215
0.433	0.131
0.685	0.097

FIGURE 2 90-mm F/2.5 macro for 35-mm SLR.



Sn.	Radius	Thickness	Glass	Clr. Ap.
1	10.9212	3.42000	LAC8	13.19
2	23.2112	1.10000		10.85
3	-41.8875	1.20000	FD6	10.58
4	28.8922	1.77565		9.72
5	30.9632	2.00000	NFD13	9.56
6	-36.6284	7.99751		9.36
7	-6.1127	1.50000	Acrylic	10.74
8 a	-8.4376	17.08263		13.18

Even polynomial aspheres

SN AD AE AP AC AH AI
 8 -7.8630E-06 -8.7100E-07 1.3553E-08 -5.4372E-10 9.7997E-12 -1.3836E-13

System first-order properties

Field: 29.0 f/ 3.20
 Stop: 1.60 after surface 6. Dia: 8.4145
 EFL: 35.8004 FVD: 36.0758 ENP: 12.0208
 BFL: 17.0826 BRL: 18.9932 EXP: -7.52072

FIGURE 3 35-mm F/3.2 for point-and-shoot.

telephoto construction. This type of construction tends to result in relatively large front aperture sizes, and it is not easy to design small lenses without compromising on image quality. Retrofocus designs sometimes have a zone of poorer image quality in a field area between the axis and the corner. This zone is a by-product of the struggle to balance lower- and higher-order aberrations so that the outer parts of the field have acceptable image quality. These lenses have particularly good relative illumination both because the basic construction results in an exit pupil quite far from the image plane and also because it is possible for the size of the pupil to increase with field angle. In order to achieve conjugate stability, it is necessary to employ the use of so-called “floating elements” or variable airspaces that change with focusing. However, this feature does result in additional optomechanical complexity.

The newer forms of this lens type fall into four broad subcategories.

Very Compact Moderate Speed These include six-element 35-mm F/2.8 with a front negative element and seven-element 28-mm F/2.8 with a leading positive element. Refer to Figs. 4 and 5, respectively. These relatively simple constructions are suitable for speeds of F/2.8 or slower and total angular coverages of up to 75°.

Highly Complex Extreme Speed As the complexity of both the front and rear groups is increased, the inverted telephoto form can be designed to achieve speeds of F/1.4 and angular fields of 90°. The use of aspherical surfaces is essential in order to achieve these specifications. Refer to Figs. 6, 7, 8, and 9.

Highly Complex Extreme Wide-Angle with Rectilinear Distortion Correction

These are inverted telephoto designs covering total fields of up to 120°, often with speeds as fast as F/2.8. Distortion correction is rectilinear. The chromatic variations of distortion, astigmatism, and coma are usually the limiting aberrations and are virtually impossible to correct beyond a certain point. Refer to Figs. 10 and 11.

Extreme Wide-Angle with Nonrectilinear Distortion (“Fish-Eye Lenses”) Without the requirement of rectilinear correction of distortion, inverted telephoto designs can be achieved quite readily with total angular fields exceeding 180°. For these lenses, the image height h and focal length f are often related by $h = f \cdot \theta$, where θ is the semifield angle. See, for example, USP 4,412,726.

Telephoto Lenses

The term *telephoto* strictly applies to lenses having a front vertex length less than the focal length (telephoto ratio less than one). The classic telephoto construction has a front positive group followed by a rear negative group. This can lead to telephoto ratios that are as short as 0.7 or less. The term telephoto is often loosely used to refer to any long-focal-length lens and one sometimes sees references made to the telephoto ratio of a wide-angle lens.

Two significant advances characterize the newer types of telephoto lenses, particularly those used for 35-mm SLR cameras. The first is the use of small internal groups for focusing, sometimes in conjunction with the front group. This feature has also led to significant improvement in the performance of these lenses with change of conjugate. This has been a problem with telephoto lenses, particularly with respect to attaining close focus with good optical quality. Internal focusing of a long-focal-length lens also has considerable advantages in terms of mechanical simplicity because a smaller mass is being moved over a significantly shorter distance.

The second advance is the employment of optical glasses having anomalous dispersion for the correction of chromatic aberrations. These newer glasses have anomalous dispersion characteristics similar to those of calcium fluorite, but with physical and chemical properties that make their use practical. These glasses are still expensive and more difficult to use than ordinary ones, but they do

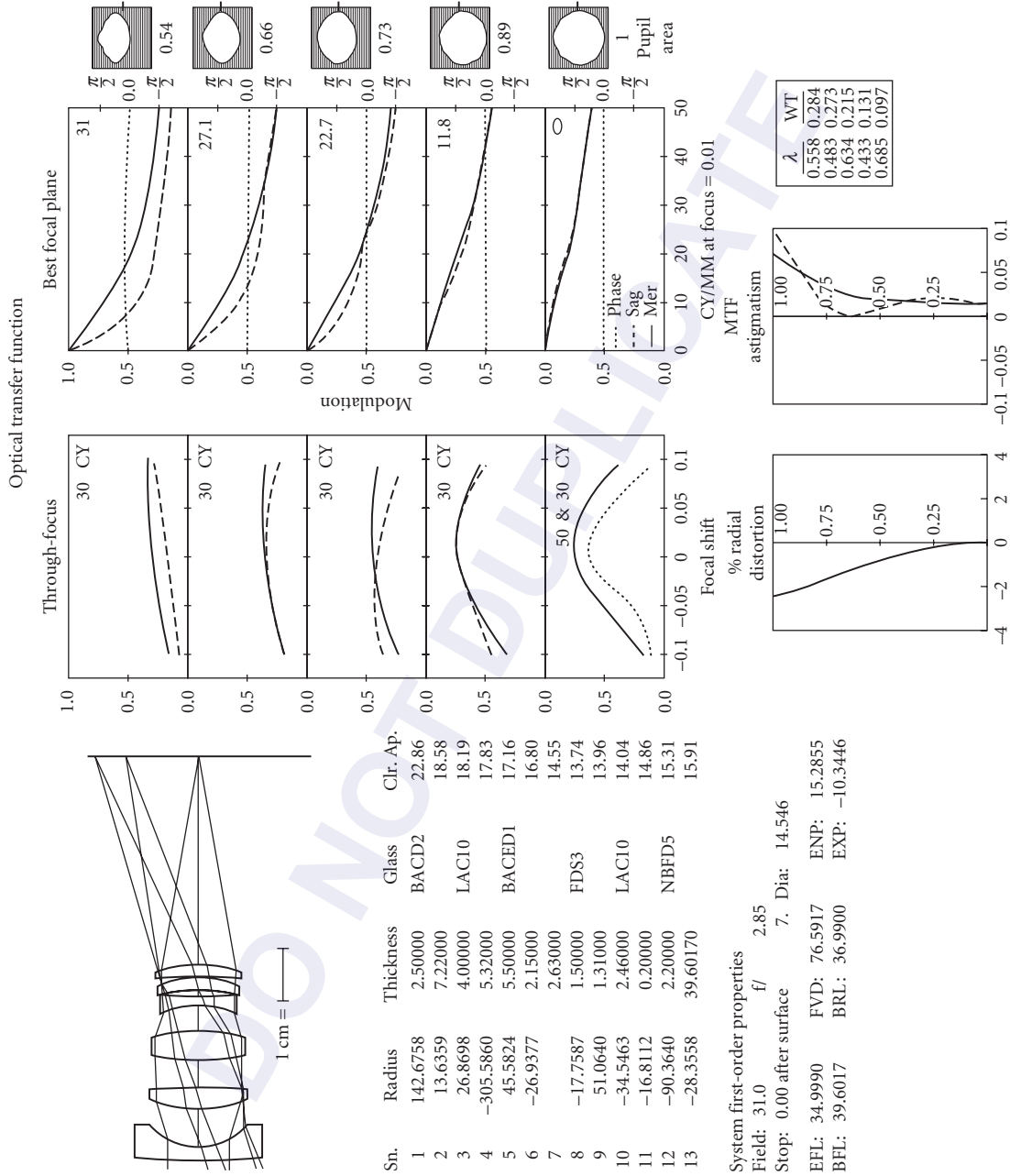


FIGURE 4 35-mm F/2.8 for 35-mm SLR.

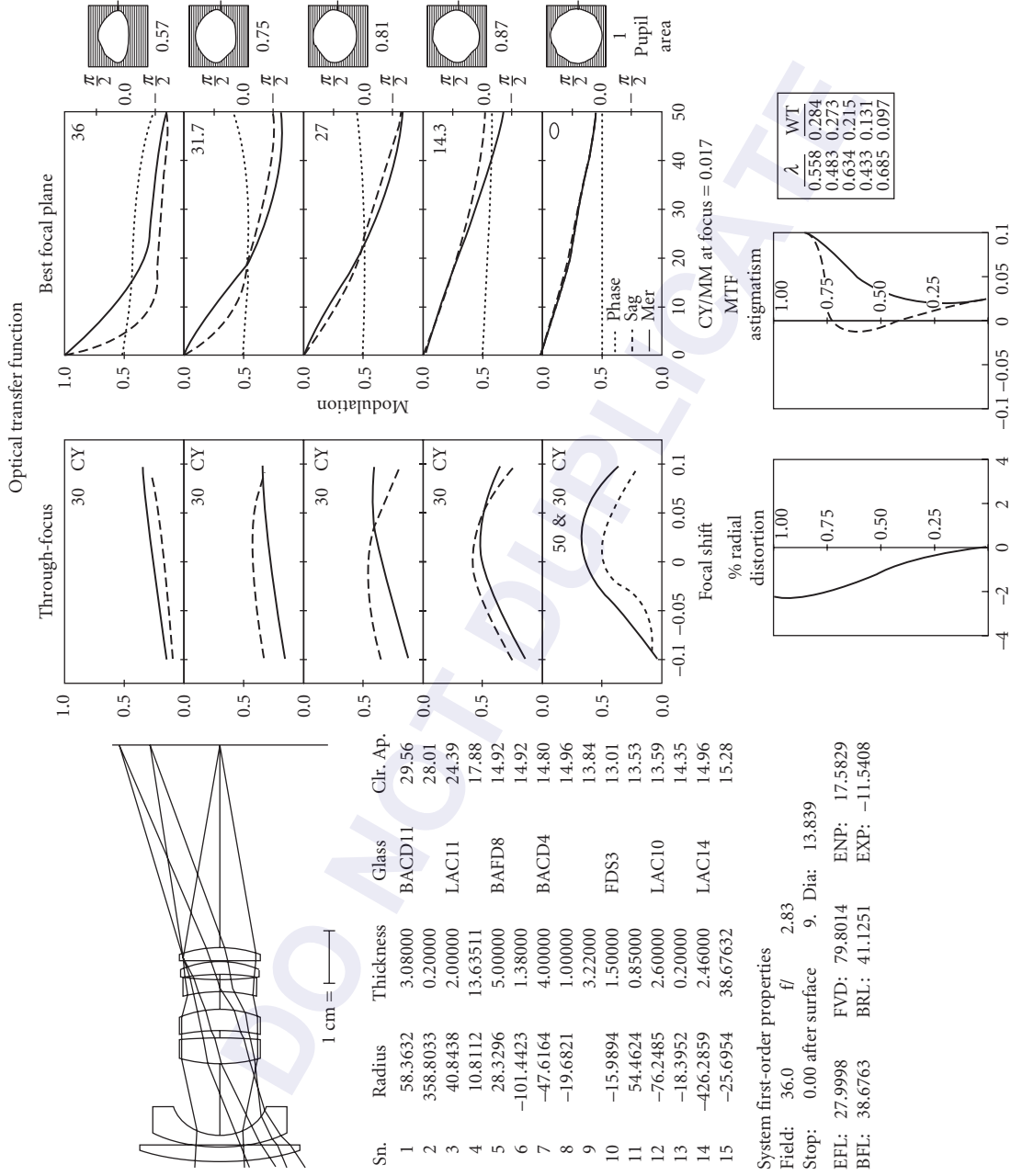
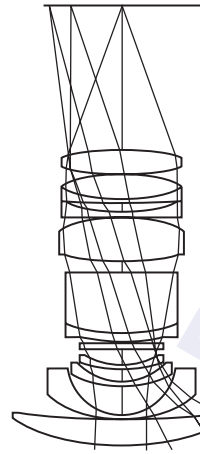
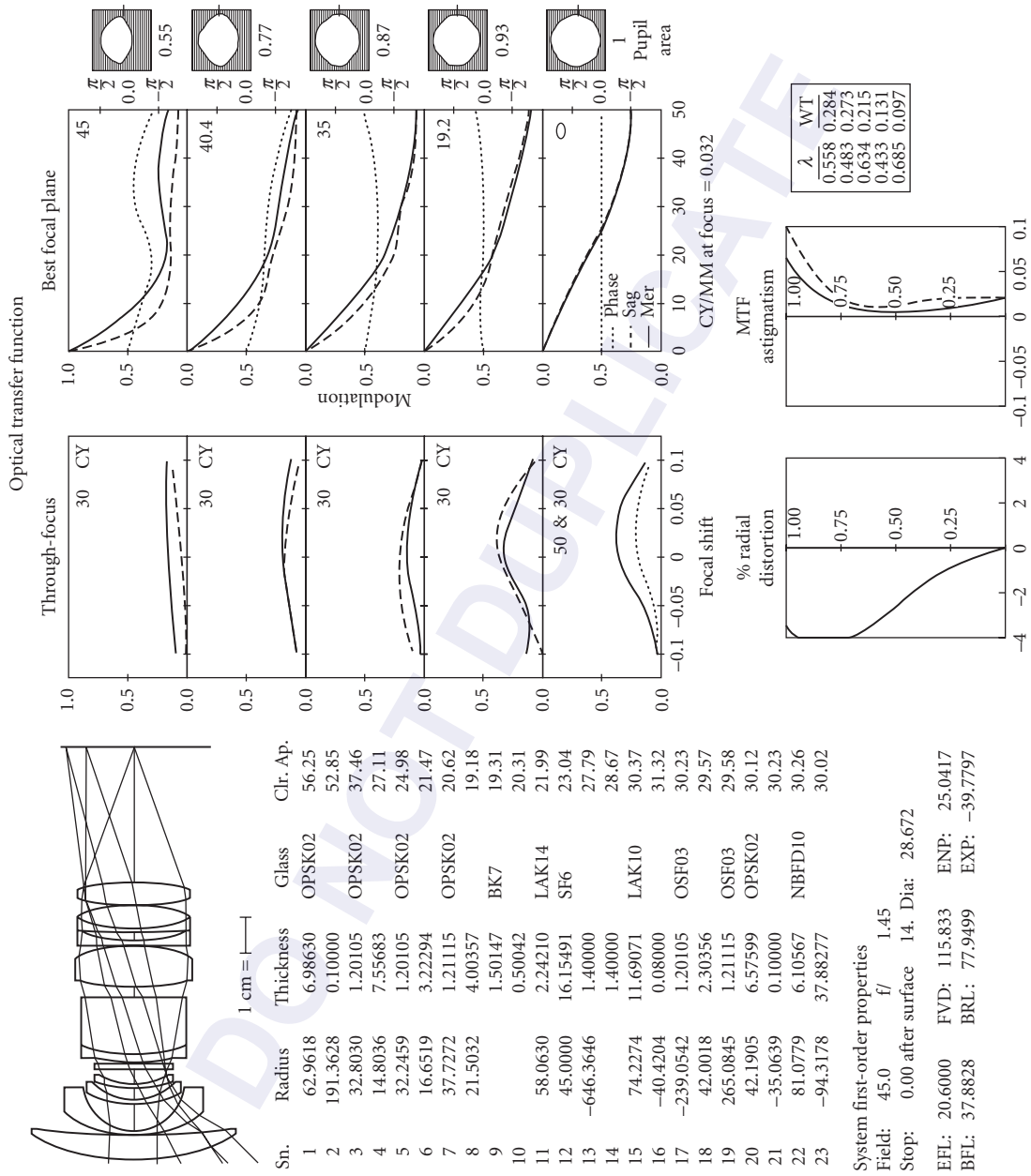


FIGURE 5 28-mm F/2.8 for 35-mm SLR.



Sn.	Radius	Thickness	Glass	Chr. Ap.
1	62.9618	6.98630	OPSK02	56.25
2	191.3628	0.10000	OPSK02	52.85
3	32.8030	1.20105	OPSK02	37.46
4	14.8036	7.55683	OPSK02	27.11
5	32.2459	1.20105	OPSK02	24.98
6	16.6519	3.22294	OPSK02	21.47
7	37.7272	1.21115	OPSK02	20.62
8	21.5032	4.00357	OPSK02	19.18
9		1.50147	BK7	19.31
10		0.50042		20.31
11	58.0630	2.24210	LAK14	21.99
12	45.0000	16.15491	SF6	23.04
13	-646.3646	1.40000		27.79
14		1.40000		28.67
15	74.2274	11.69071	LAK10	30.37
16	-40.4204	0.08000		31.32
17	-239.0542	1.20105	OSF03	30.23
18	42.0018	2.30356		29.57
19	265.0845	1.21115	OSF03	29.58
20	42.1905	6.57599	OPSK02	30.12
21	-35.0639	0.10000		30.23
22	81.0779	6.10567	NBFD10	30.26
23	-94.3178	37.88277		30.02

System first-order properties
 Field: 45.0 f/ 1.45
 Stop: 0.00 after surface 14. Dia: 28.672
 EFL: 20.6000 FVD: 115.833 ENP: 25.0417
 BFL: 37.8828 BRL: 77.9499 EXP: -39.7797

FIGURE 6 20-mm F/1.4 for 35-mm SLR.

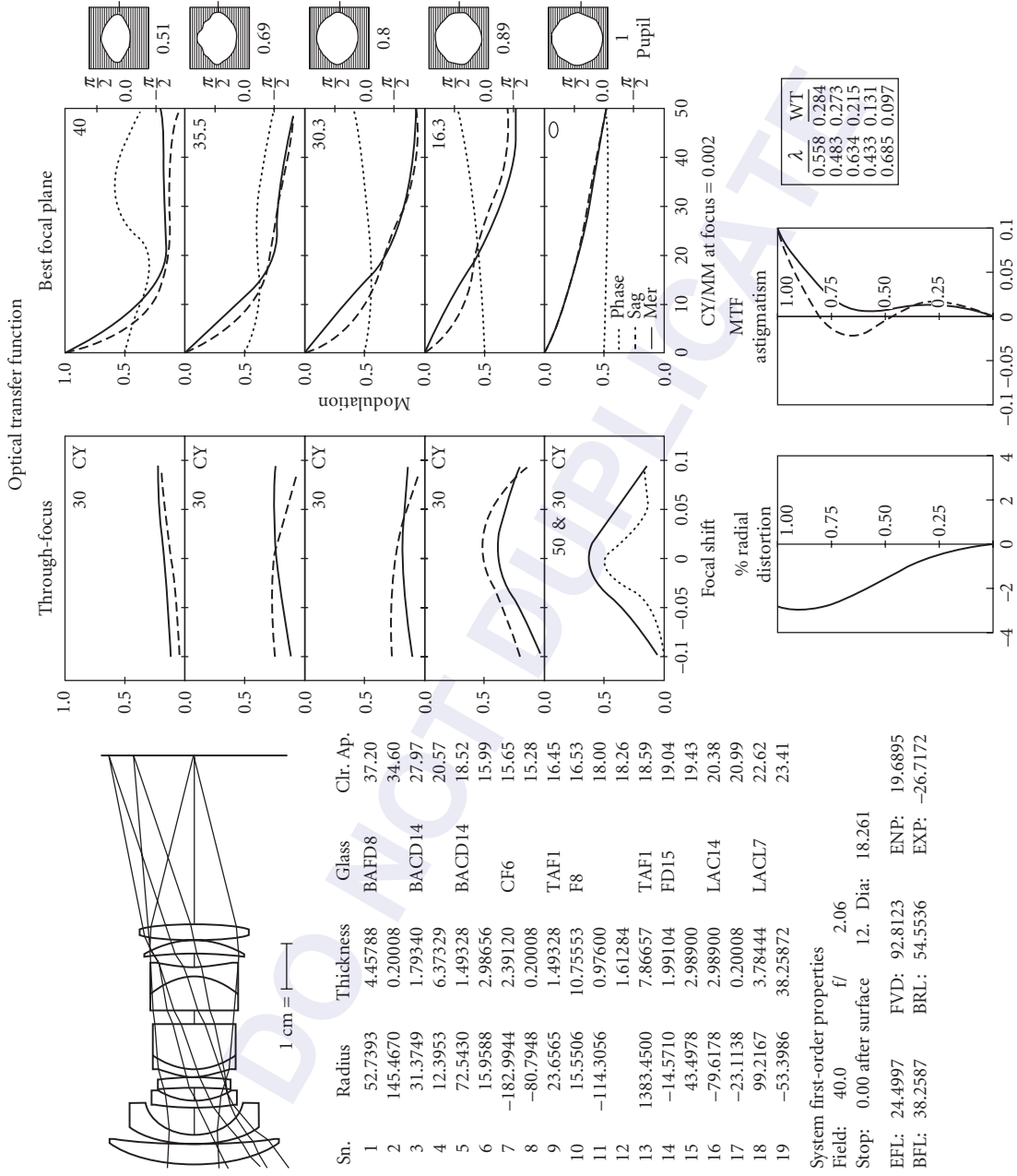


FIGURE 7 24-mm F/2 for 35-mm SLR.

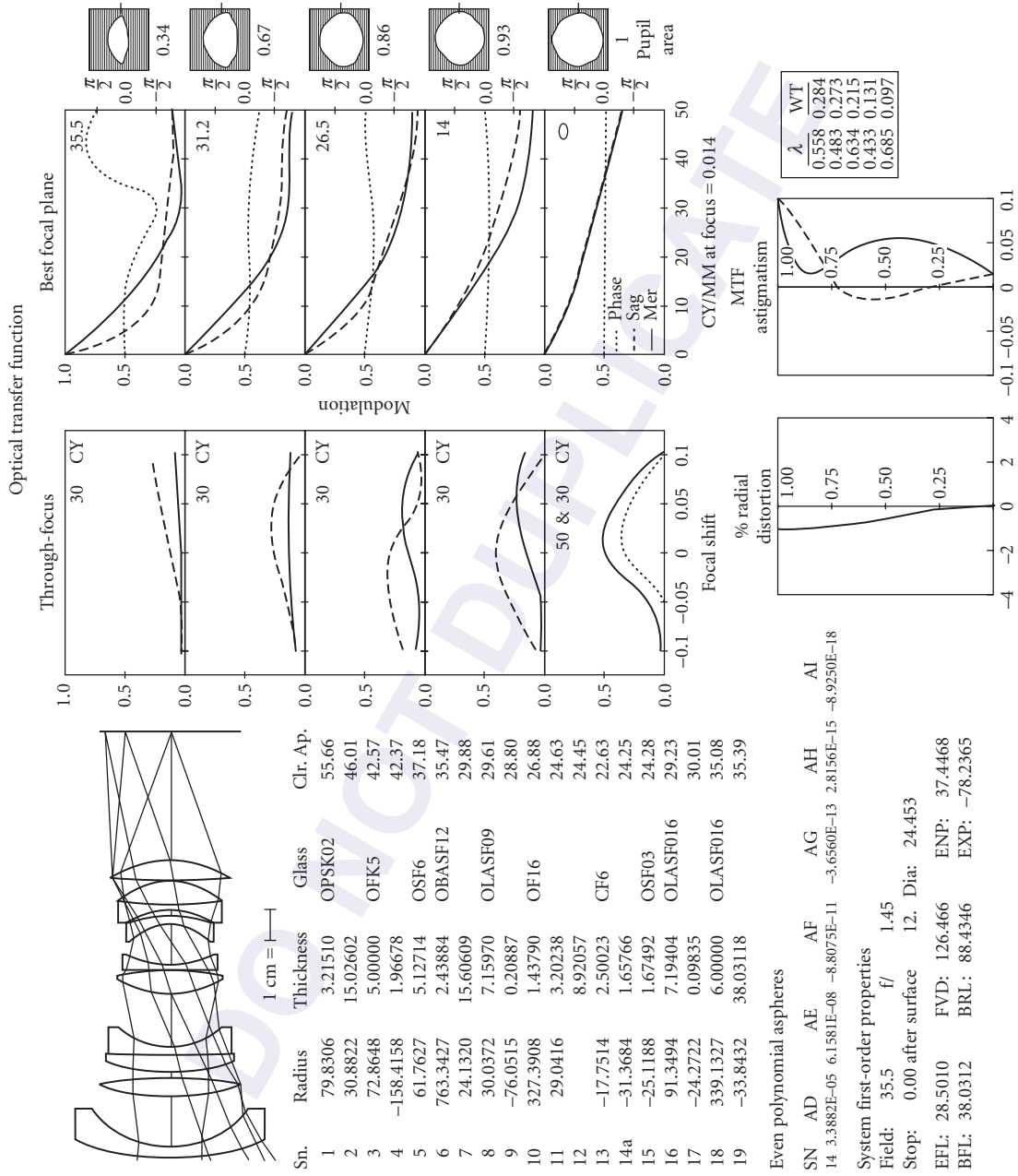


FIGURE 8 28-mm F/1.4 aspheric for 35-mm SLR.

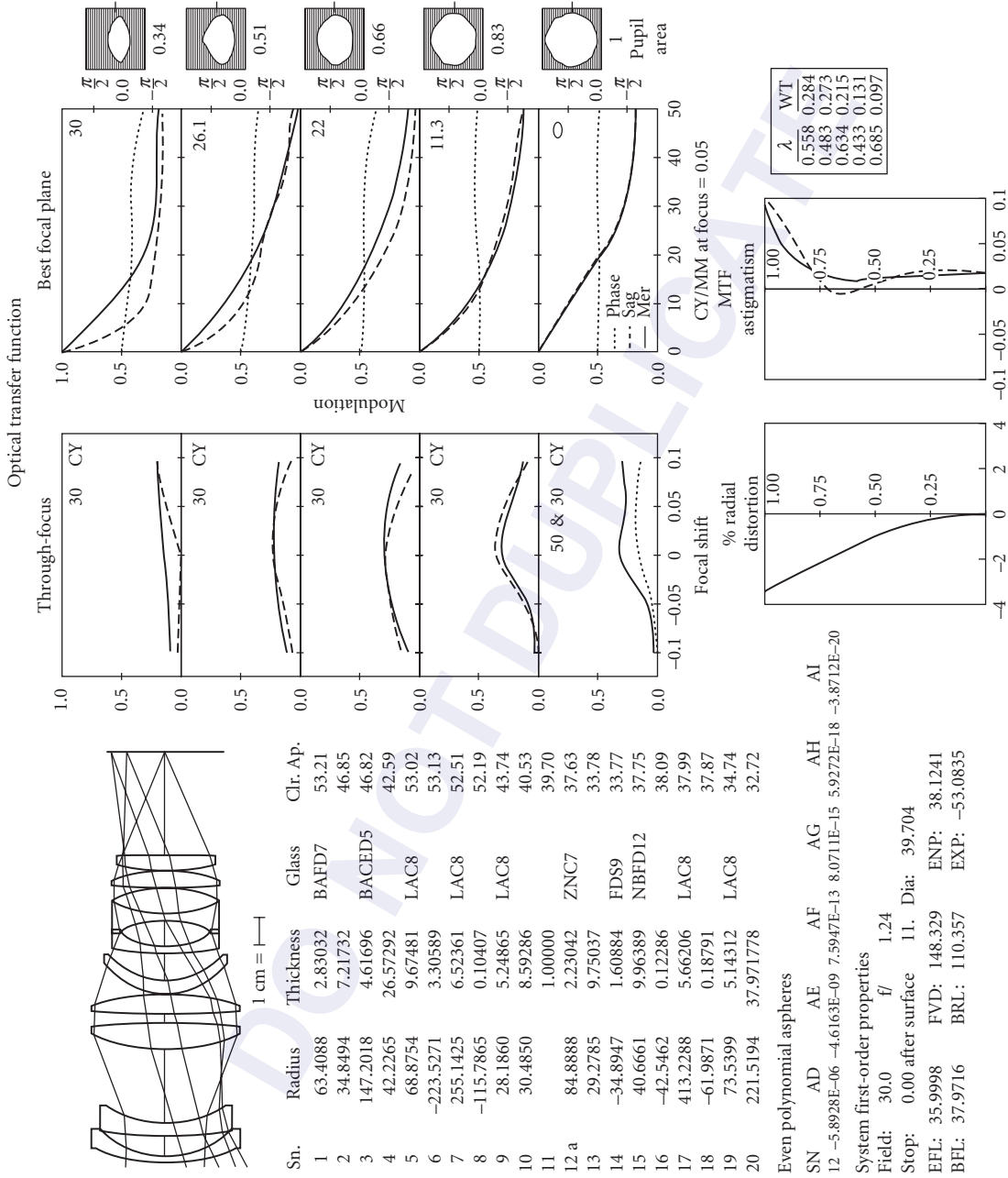


FIGURE 9 35-mm F/1.2 aspheric for 35-mm.

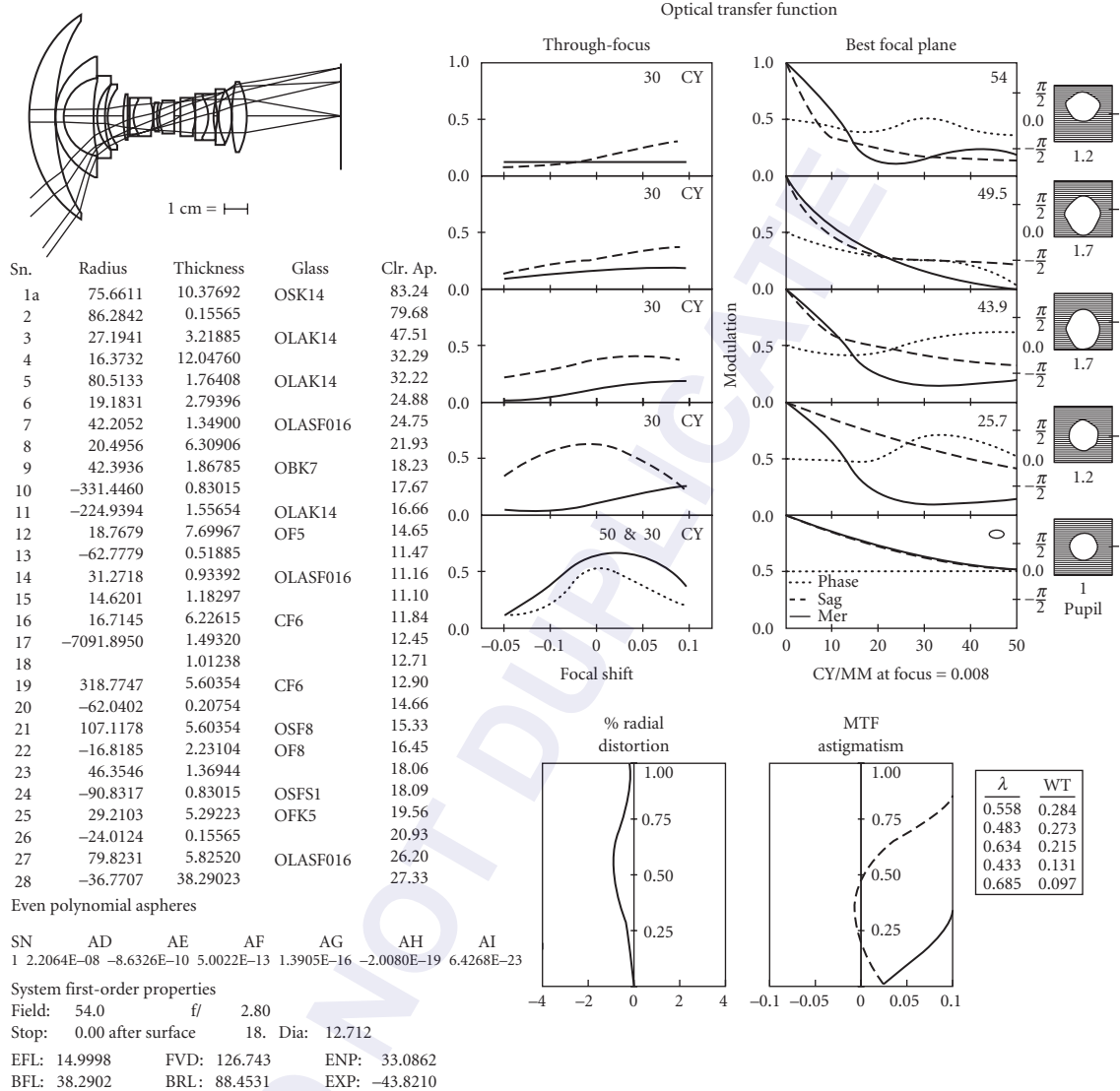


FIGURE 10 15-mm F/2.8 for 35-mm SLR.

offer significant advantages in terms of reducing the chromatic aberrations that otherwise severely limit the imaging potential of all long-focal-length refracting optics. Typical available versions of these glass types are the FK Schott, FCD Hoya, FPL Ohara, and PFC Corning series of glasses.

These design types offer outstanding optical correction together with remarkable specifications, resulting in considerable size and cost. Commercial embodiments include 300-mm F/2 and 400-mm F/2.8 for 35-mm. They are widely used for sports and wildlife photography. Since secondary color increases as the front vertex distance is reduced, it is advisable to regard excessively short all-refractive telephotos with some caution. Refer to Figs. 12 and 13. The letter *z* in the thickness column in Fig. 12 (later in Figs. 14, 16, and 17) represents the zooming or variable space between the lenses.

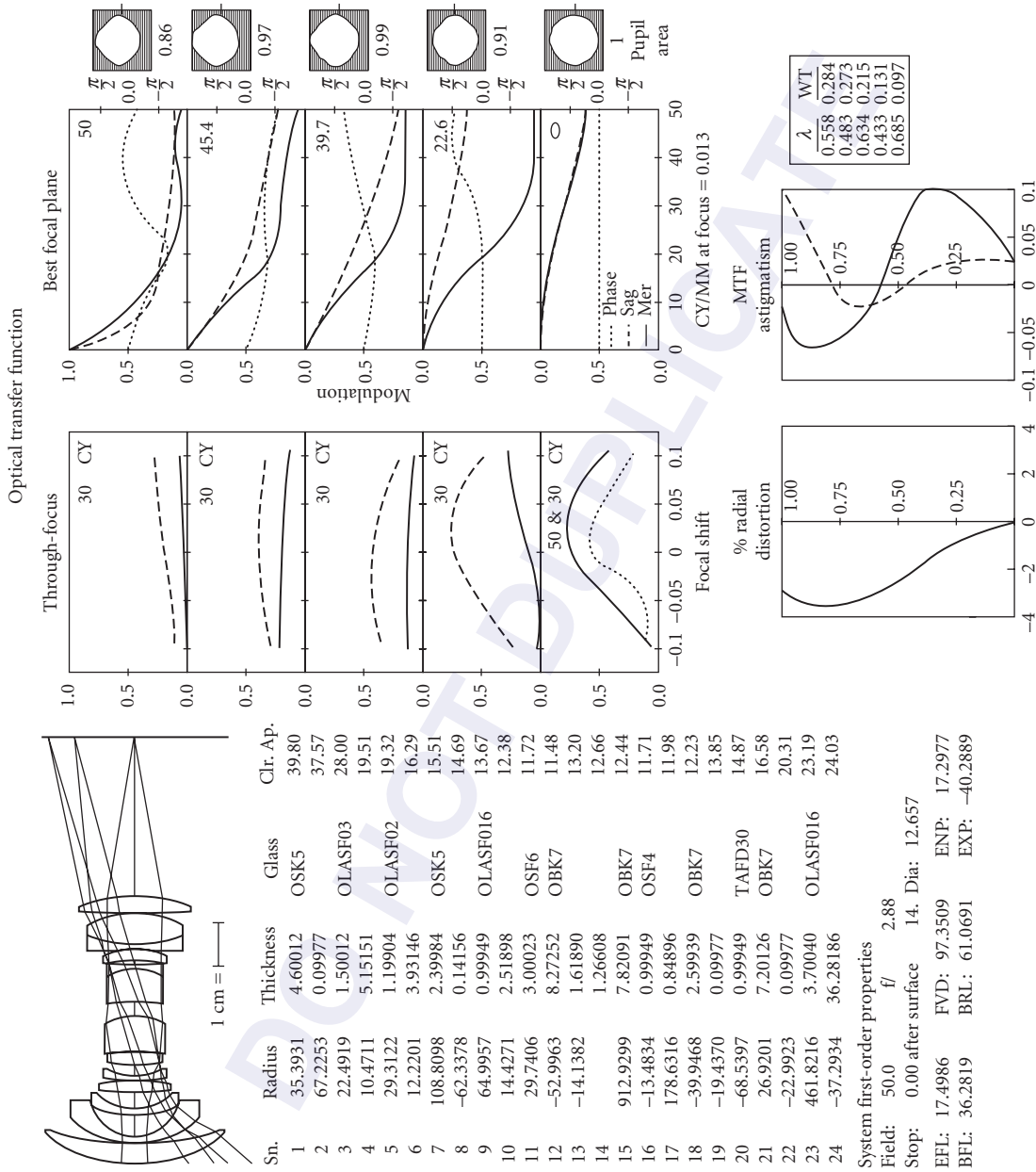


FIGURE 11 17-mm F/2.8 for 35-mm SLR.

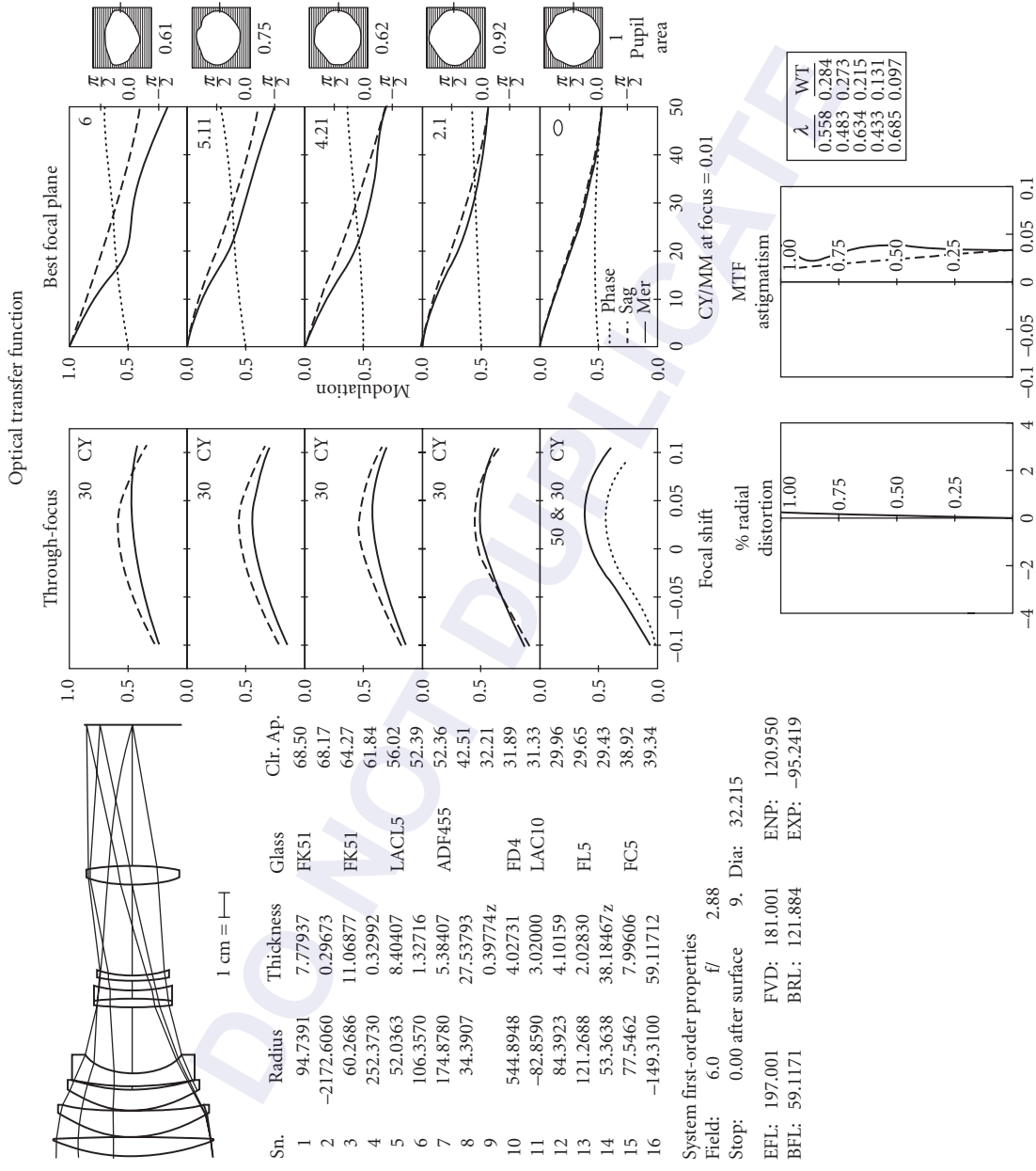


FIGURE 12 200-mm F/2.8 for 35-mm SLR.

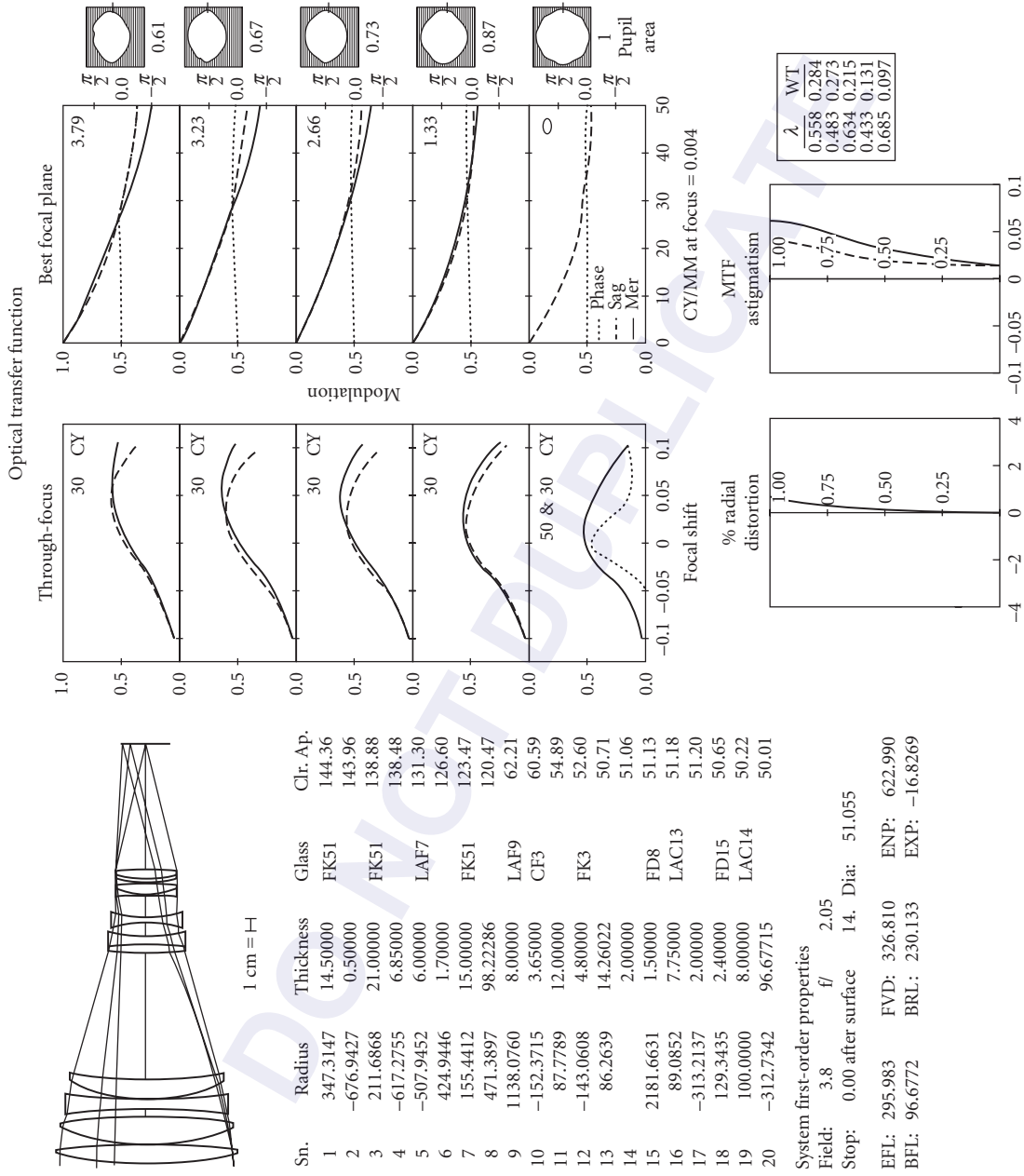


FIGURE 13 300-mm F/2 for 35-mm SLR.

Zoom Lenses

Zoom lenses have evolved significantly in the past 20 years. In the early 1970s, there was basically only the classic four-group type of zoom lens. This four-group zoom has two moving groups between a front group used only for focusing and a stationary rear (*master*) group. This type is still found on consumer video cameras. Figures 14 and 15 show a variation of this form with the rear group also moving for zooming. The master group could often be changed to yield a different zoom with the same ratio over a different range.

The second basic form, originating in the mid 1970s, was the two-group wide-angle zoom, typically 24 to 48 mm and 35 to 70 mm for the 35-mm format. Both the front negative group and the rear positive group move for zooming, and the front group is also used for focusing. This lens type has an inherently long back focal length, making it eminently suitable for the SLR camera. See, for example, USP 4, 844, 599. The maximum zoom range is about 3 : 1.

In order to achieve lens types such as a 28- to 200-mm zoom for 35-mm, new ideas had to be employed. The resulting lenses have up to five independent motions, including that of the diaphragm. These degrees of freedom allow for the location of the entrance pupil to be near the front of the lens at the short-focal-length position and also for the exit pupil to be located near the rear, particularly at the long-focal-length setting. These conditions result in acceptably small size. The extra zooming motions permit a large focal-length range to be achieved without any one motion being excessively long. There is a constant struggle in the design of these zooms to minimize the diameter of the front of the lens. This is not only to reduce size and weight, but also to permit the use of acceptably small filters. Some designs do have problems with relative illumination at the wide-angle end.

In the past, these lenses have been focused either by moving the front group or by moving the entire lens, the latter option leading to the so-called varifocal zoom. However, more recent developments in miniature electromechanical and autofocus systems have led to the evolution of extended range zooms in which the distinction between a focusing and a compensating group has become academic. As a result, small internal zooming groups can serve a dual function as focusing groups under the control of an autofocus system. Refer to Figs. 16, 17, and 18.

A recent new development in zooms is one for the so-called compact 35-mm camera. In its most basic form, this type can have as few as three elements and is characterized by having a front positive component and a rear negative component. This lens has an inherently short back focal length at the wide end, making it not suitable for SLR cameras with swinging mirrors. In more complicated versions, this idea can be extended to 28 to 160 mm or further, the main limitation being a small relative aperture at the long-focal-length end. A recent practical embodiment is a four-element 38- to 90-mm F/3.5 to 7.7 having three aspherical surfaces. See, for example, USP 4, 936, 661.

Zoom lenses are also found on most consumer video cameras. The classic fixed front-and-rear-group type (with the aperture stop in the rear group) is still commonly used because the very small format sizes can permit acceptably small lenses. This lens form is also used for motion picture and television zooms. In many of these applications, it is desirable to have an exit pupil position that does not change with zooming. Telecentricity of the exit pupil is also sometimes required. In addition, the motion picture industry still prefers zoom lenses that have conventional front-group focusing in order to easily calibrate tape-measure focus measurements.

Very long range television zooms (often 30 : 1 or more) are also of the fixed front and rear type, with a succession of cascading zooming groups in between.

27.4 CLASSIFICATION SYSTEM

A wide variety of camera lenses has been classified in Table 1 in terms of total pixel capability P and pixels per steradian AD. Pixels are defined as digital resolution elements relative to a specified modulation level and are calculated as follows:

The polychromatic optical transfer function of each lens is calculated and the spatial frequencies at which the modulation falls to 0.5 and 0.2 is noted at each of five field points. The lower of the meridional and sagittal values is used.

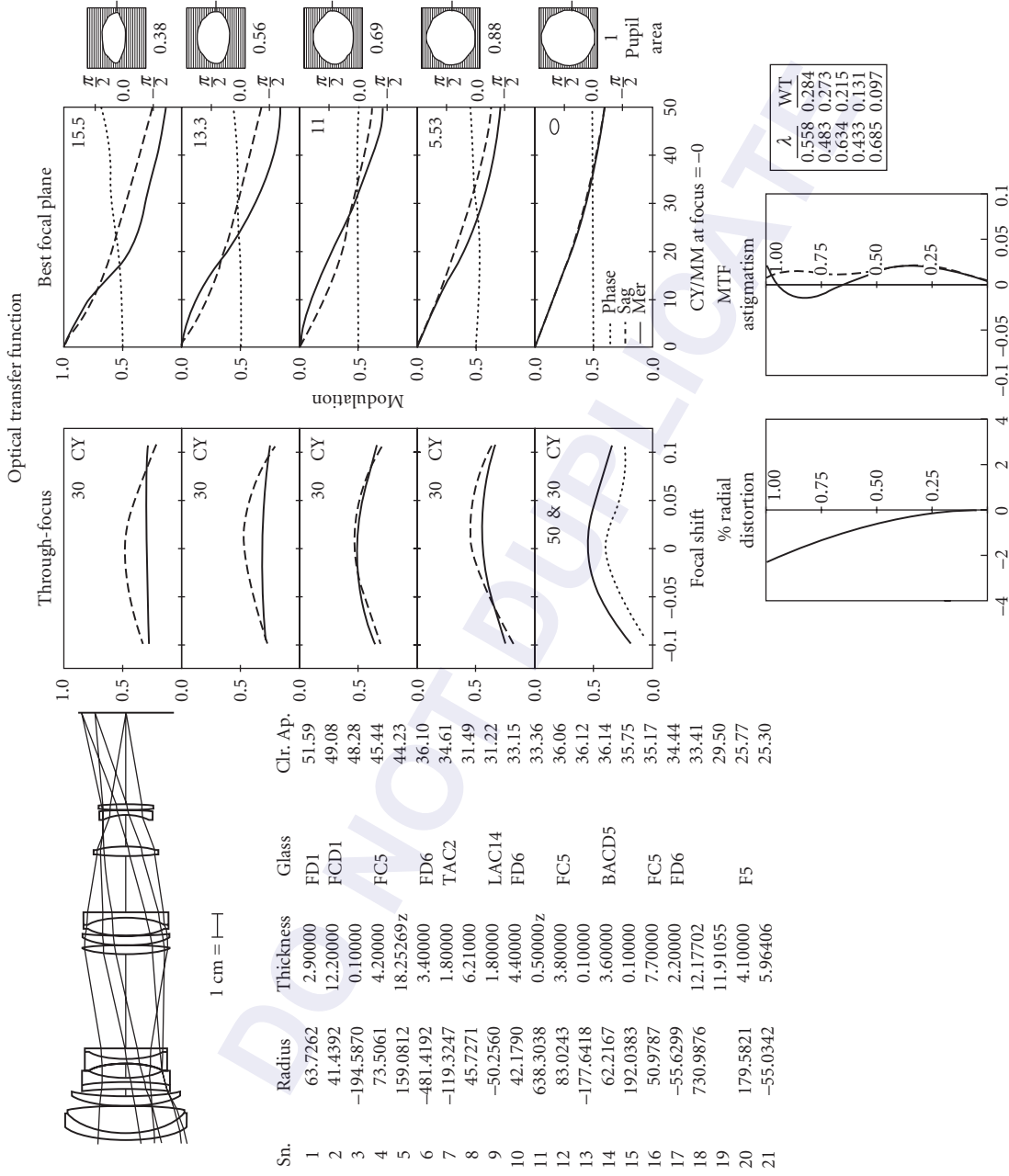


FIGURE 14 70-210-mm F/2.8-4 at $f = 70$.

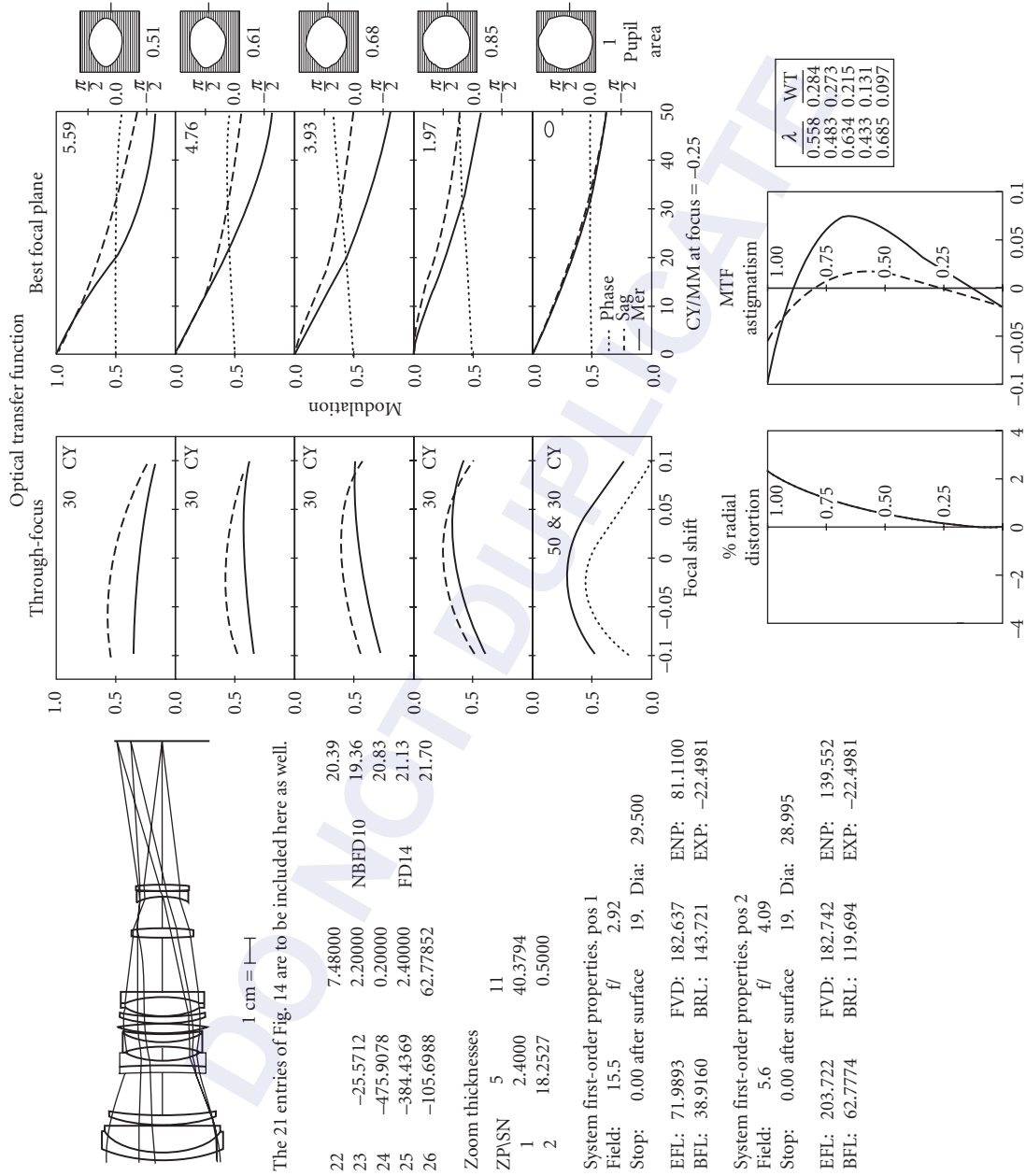


FIGURE 15 70-210-mm F/2.8-4 at $f=210$.

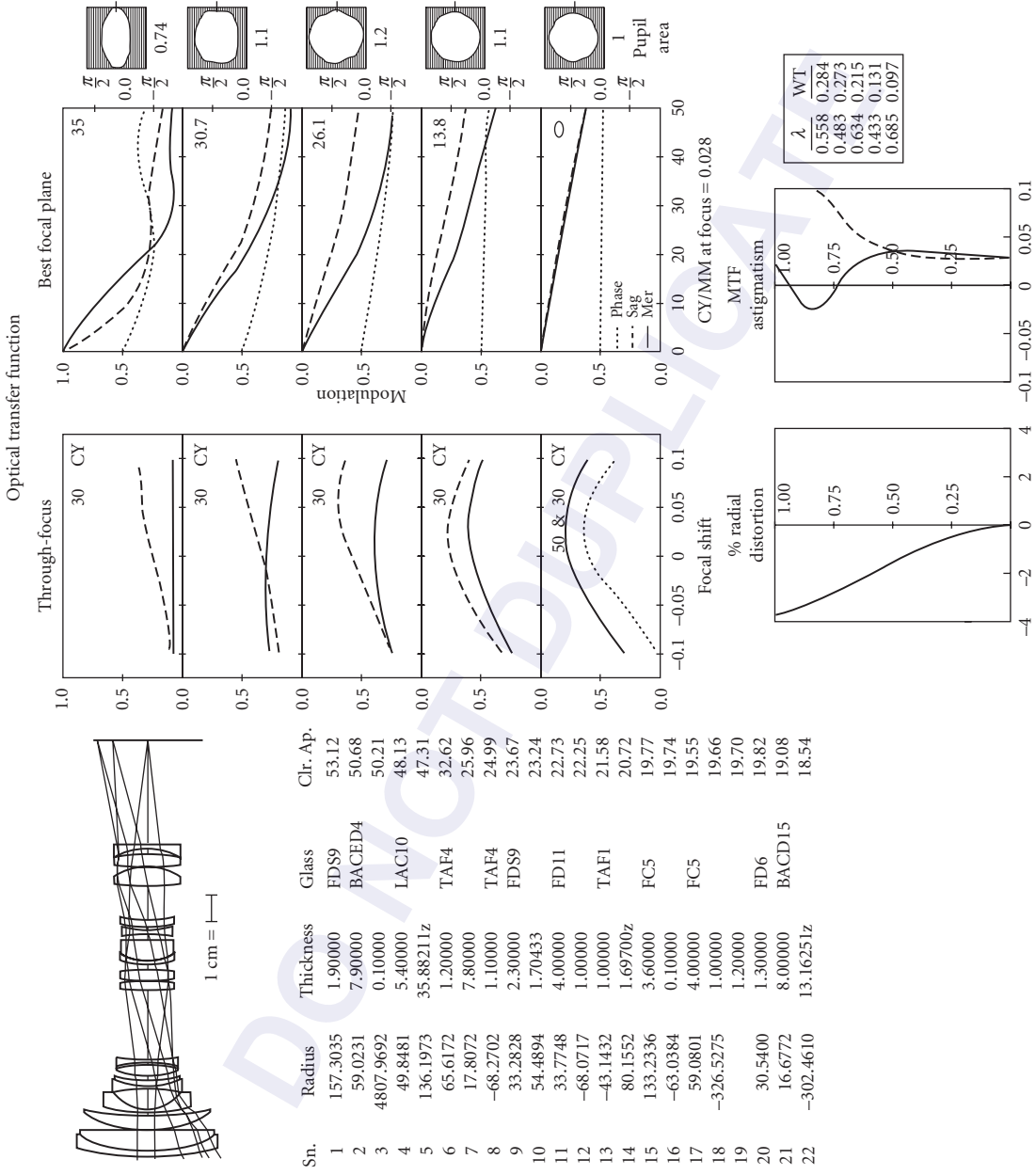


FIGURE 16 28-150-mm F/4.1-5.7 at $f = 28$.

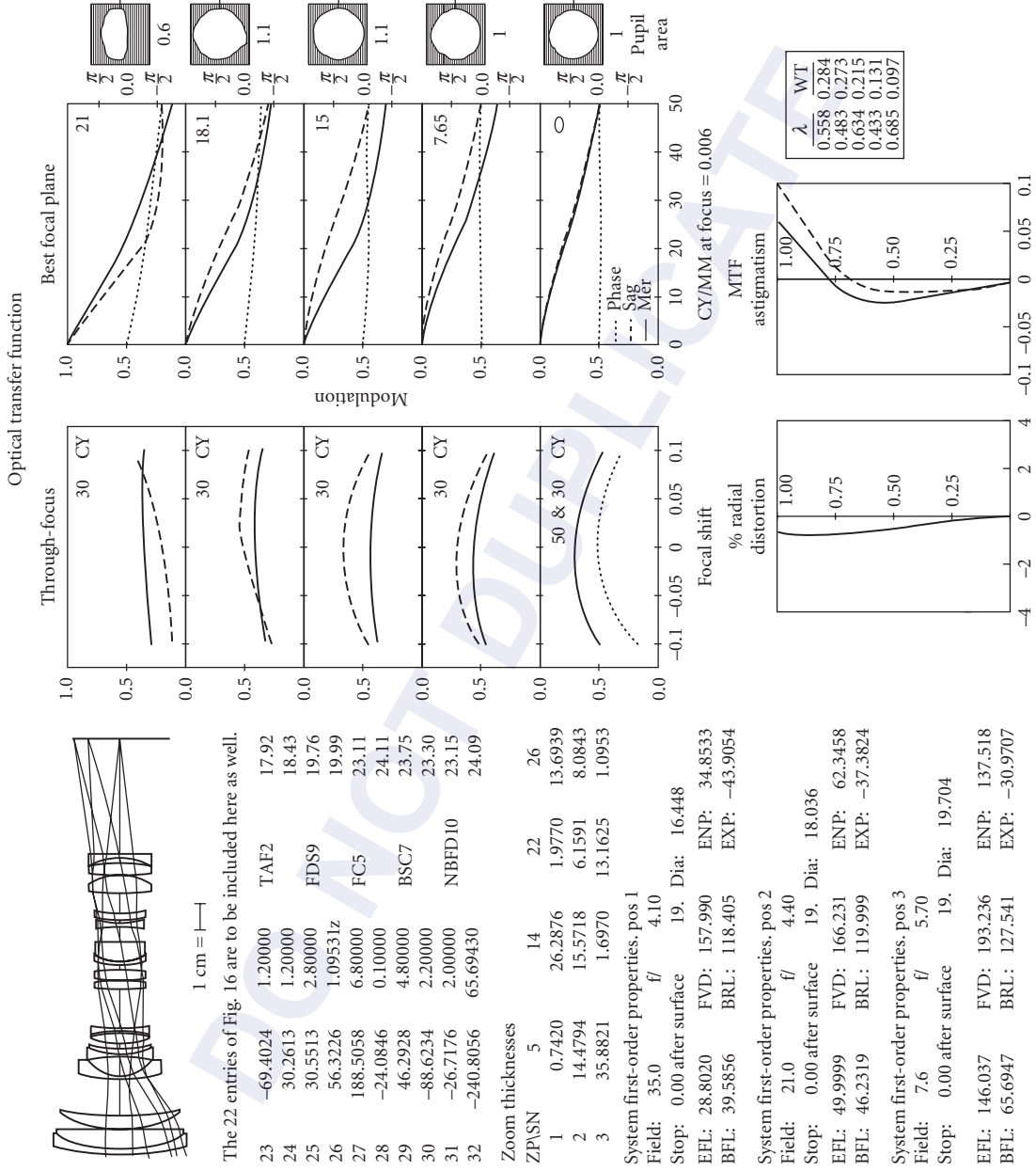
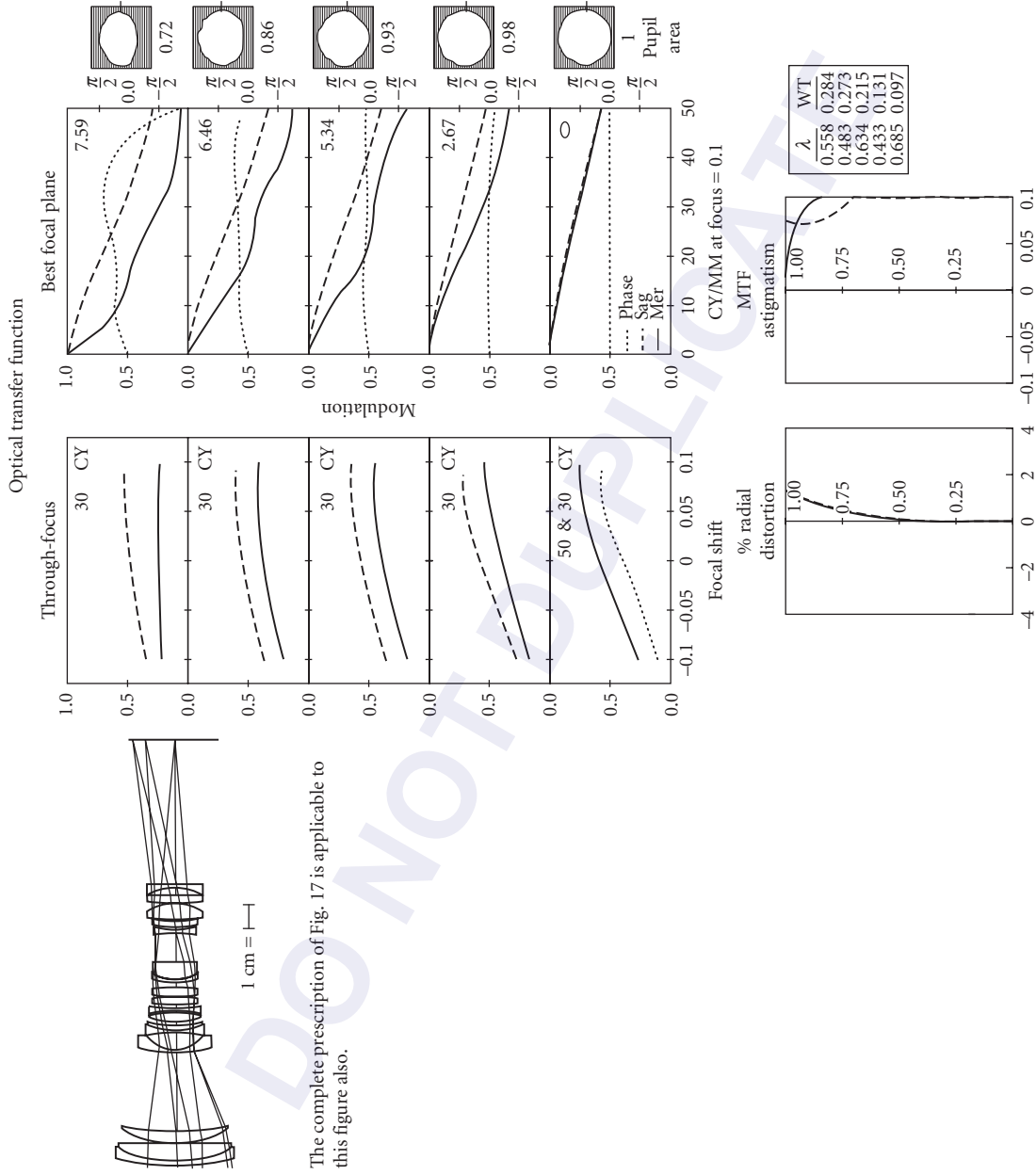


FIGURE 17 28–150-mm F/4.1–5.7 at $f = 50$.



The complete prescription of Fig. 17 is applicable to this figure also.

FIGURE 18 28-150-mm F/4.1-5.7 at $f=150$.

TABLE 1 Pixel Imaging Characteristics of Camera Lenses

EFL	D	$2W$	F-no	$P_{0.5}$	$P_{0.2}$	$AD_{0.5}$	$AD_{0.2}$	Comments
5.1	11	94	1.8	0.27	0.93	0.14	0.47	Cinegon 2/3" video
6.7	21	114	1.8	0.36	1.06	0.13	0.37	Xenoplan
10.2	16	76	1.8	0.66	1.59	0.50	1.19	Cinegon 1" video
15.0	41	108	2.8	0.95	3.39	0.37	1.31	Ultrawide 35-mm SLR
17.1	11	36	1.0	0.25	0.66	0.81	2.15	Xenar 2/3" video
17.5	42	100	2.9	1.48	3.14	0.66	1.40	Ultrawide 35-mm SLR
17.6	11	35	1.4	0.69	1.93	2.37	6.64	Xenon 2/3" video
20.6	41	90	1.5	0.52	4.05	0.28	2.20	Ultrawide 35-mm SLR
24.5	41	80	2.1	1.24	4.86	0.84	3.30	Very wide 35-mm SLR
25.5	27	56	2.8	1.61	2.56	2.19	3.48	Panavision Primo zoom
28.0	30	56	2.8	1.32	8.74	1.79	11.88	Xenar 35-mm cine
28.0	41	72	2.8	2.32	7.28	1.93	6.06	Wide 35-mm SLR
28.5	41	71	1.5	0.82	3.14	0.70	2.69	Wide 35-mm SLR
28.8	40	70	4.1	3.35	10.31	2.95	9.08	Wide-tele 35-mm SLR
35.0	42	62	2.8	3.43	17.54	3.82	19.54	Wide 35-mm SLR
35.8	40	58	3.2	1.45	10.43	1.84	13.24	Snapshot 35-mm
36.0	42	60	1.2	2.63	5.27	3.12	6.26	Wide 35-mm SLR
50.0	38	42	4.4	3.23	14.28	7.74	34.21	Wide-tele 35-mm SLR
51.0	40	43	1.4	0.76	7.02	1.74	16.06	Normal 35-mm SLR
55.5	41	41	1.2	0.87	3.41	2.17	8.58	High-speed 35-mm SLR
55.5	41	41	1.2	0.54	4.22	1.35	10.60	Asph normal 35-mm SLR
57.5	40	37	1.4	0.98	3.67	3.05	11.41	Close focus 35-mm SLR
58.0	133	98	5.6	13.77	60.12	6.37	27.82	Super-Angulon large-format
72.0	40	31	2.9	3.22	10.44	14.11	45.70	Telezoom 35-mm SLR
73.8	16	12	2.8	0.31	1.26	9.01	36.61	Tele-xenar 1" video
75.0	30	23	2.0	1.87	14.12	14.83	111.95	Xenar 35-mm cine
80.0	80	53	2.8	13.32	53.45	20.18	80.97	Xenotar medium-format
90.0	173	88	8.0	19.42	61.93	11.01	35.12	Super-Angulon large-format
90.0	207	98	5.6	30.07	75.59	13.92	34.98	Super-Angulon large-format
100.0	27	15	2.8	3.11	6.65	57.86	123.72	Panavision Primo zoom
100.0	116	60	5.6	33.11	85.72	39.34	101.84	APO-Symmar large-format
120.0	143	62	5.6	30.31	100.43	33.78	111.91	APO-Symmar large-format
120.0	168	70	5.6	23.78	113.14	20.93	99.57	Super-Symmar large-format
146.0	39	15	5.7	2.67	12.85	48.44	232.76	Wide-tele 35-mm SLR
150.0	80	30	4.0	12.35	87.89	57.69	410.55	Tele-Xenar medium-format
150.0	138	49	5.6	12.12	28.95	21.42	51.18	Xenar large-format
150.0	202	68	5.6	17.87	117.45	16.64	109.35	Super-Symmar large-format
180.0	300	45	5.6	100.80	368.45	210.77	770.41	Makro-Symmar 1 : 1
197.0	41	12	2.9	4.17	16.37	121.10	475.55	Telephoto 35-mm SLR
204.0	40	11	4.1	4.18	12.15	139.28	405.17	Telezoom 35-mm SLR
210.0	274	66	5.6	18.37	126.58	18.12	124.88	Super-Symmar large-format
210.0	400	87	8.0	22.62	130.62	13.11	75.70	Super-Angulon large-format
240.0	164	19	9.0	44.94	111.55	521.56	1294.60	Artar 1 : 1
250.0	126	28	5.6	38.05	89.71	203.88	480.70	Tele-Artar large-format
268.0	27	6	2.8	3.11	7.64	361.19	887.30	Panavision Primo zoom
296.0	39	8	2.1	5.36	20.36	388.37	1473.60	Telephoto 35-mm SLR
360.0	392	57	6.8	45.09	378.07	59.22	496.57	APO-Symmar large-format
400.0	250	35	5.6	33.77	160.01	116.13	550.27	APO-TXR large-format
480.0	327	19	11.0	76.29	256.34	885.40	2975.00	APO-Artar 1 : 1
480.0	400	45	8.4	86.48	509.51	180.83	1065.30	APO-Symmar large-format
800.0	500	35	12.0	108.62	460.37	373.54	1583.10	APO-TXR large-format

D —image diameter in mm
 P_m —pixels $\times 10^6$ at modulation level m
 F -no—F-number of the lens

W —semifield angle in degrees
 AD_m —pixels $\times 10^6$ per steradian at modulation level m
EFL—effective focal length of the lens in mm

The image field of the lens, assumed to be circular with diameter D , is divided into four annular regions. The outer boundaries of each region correspond, respectively, to $0.35H$, $0.7H$, $0.85H$, and $1.0H$, where H is the maximum field height. The area of each region is computed.

The average of the inner and outer boundary-limiting spatial frequency values is assigned to each region. This is done for both the 0.5 and 0.2 modulation levels.

The area of each annular region, in square millimeters, is multiplied by the square of the spatial frequency values from the previous step to yield regional pixel counts for both 0.5 and 0.2 modulation levels.

The pixel counts are summed over all regions to yield the D data in Table 1.

The AD data in Table 1 are obtained by dividing the total pixel values by the solid angle of the lens in object space. The solid angle S is given by the following formula:

$$S = 2\pi(1 - \cos W)$$

where W is the semifield angle of the lens in degrees.

In general, for a given image diameter D , a larger P implies higher image quality or greater information-gathering capability. A lens designed for a smaller D will have a lower P than a lens of similar quality designed for a larger D . These same generalizations hold for AD except that, in addition, a lens designed for a smaller field angle and a given D will have a larger AD than a lens of similar image quality designed to cover a wider field for the same D . In other words, for the same image-quality level and format size, wide-angle lenses have lower AD values than do narrow angle lenses.

27.5 LENS PERFORMANCE DATA

A wide variety of camera lenses has been selected to show typical performance characteristics. In most cases, the data have been derived from the referenced published U.S. patents. The authors have taken the liberty of reoptimizing most of the data to arrive at what would, in our judgment, correspond to production-level designs. All performance data have been shown at maximum aperture. It is important to realize that photographic lenses are invariably designed so that optimum performance is achieved at F -numbers at least 2 stops slower than maximum. A general explanation of the data page follows.

The lens drawing shows the marginal axial rays together with the upper and lower meridional rays for seven-tenths and full field.

The lens prescription and all other data are in millimeters. Glass catalogs are Hoya, Ohara, and Schott. Distances to the right of a surface are positive. A positive radius means that the center of curvature is to the right of the surface. The thickness and glass data indicate the distance and medium immediately following the particular surface.

The optical transfer function (OTF) plots show the through-focus modulation transfer function (MTF) on the left and the OTF at best axial focus on the right. The data are shown for five field points, viz., the axis, $0.35H$, $0.70H$, $0.85H$, and $1.0H$, where H is the maximum field angle in object space. The actual field angles are indicated in the upper-right-hand corner of each best-focus OTF block and are in degrees. The through-focus data are at the indicated spatial frequency in cycles per millimeter with an additional frequency on-axis (dotted curve). Both the through-focus and best-focus data indicate meridional (solid curves) and sagittal (dashed curves) MTF. The modulus scale is on the left of each block and runs from zero to one. The phase of the OTF is shown as a dotted curve in the best-focus plots. The scale for the phase is indicated on the right of each best-focus block and is in radian measure. All the OTF data are polychromatic. The relative weights and wavelengths used appear in the lower-right-hand corner of each page. The wavelengths are in micrometers and the weights sum to one. The axial focus shift indicated beneath the best-focus plots is relative to the zero position of the through-focus plots. The best-focus plane is at the peak of the additional axial through-focus plot (dotted curve).

Vignetting for each field angle is illustrated by the relative pupil area plots on the right-hand side of each page. The distortion plots shows the percentage of radial distortion as a function of fractional field height. The MTF astigmatism plot shows the loci of the through-focus MTF peaks as a function of fractional field height. The data can be readily determined directly from the through-focus MTF plots.

Certain acronyms are used in the system first-order properties:

Effective focal length (EFL)
Back focal length (BFL)
Front vertex distance (FVD)
Barrel length (BRL)
Entrance pupil distance (ENP)
Exit pupil distance (EXP)

The ENP and EXP data are measured from the front and rear vertices of the lens, respectively. A positive distance indicates that the pupil is to the right of the appropriate vertex.

27.6 ACKNOWLEDGMENTS

The authors would like to acknowledge data provided by R. Mühlschlag of Jos. Schneider Optische Werke and C. Marcin of Schneider Corporation of America. We also appreciate permission granted by Panavision Corporation to use data pertaining to the Primo zoom lens.

27.7 FURTHER READING

Betensky, E. I., "Photographic Lenses," in R. Shannon and J. Wyant (eds.), *Applied Optics and Optical Engineering*, vol. 8, Academic Press, New York, 1980.

Cook, G. H., "Photographic Objectives," in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. 3, Academic Press, New York, 1965.

Figures 1–18 contain data that was originally derived from the following United States Patents. The patents are listed in order corresponding to the figure order. USP 4, 431, 273; 4, 381, 888; 4, 095, 873; 3, 830, 554; 4, 770, 512; 4, 333, 714; 4, 136, 931; 4, 303, 315; 4, 792, 216; 4, 110, 007; 3, 942, 875; 4, 786, 152; 4, 732, 459; 5, 018, 843 (Figs. 14–16); 4, 758, 073 (Figs. 17–18).

Kingslake, R., *A History of the Photographic Lens*, Academic Press Inc., San Diego, Calif., 1989.

Kingslake, R., *Optics in Photography*, SPIE Press, Bellingham, Wash., 1992.

This page intentionally left blank.

DO NOT DUPLICATE

MICROSCOPES

Rudolf Oldenbourg

*Marine Biological Laboratory
Woods Hole, Massachusetts, and
Physics Department
Brown University
Providence, Rhode Island*

Michael Shribak

*Marine Biological Laboratory
Woods Hole, Massachusetts*

28.1 GLOSSARY

f	focal length
M	magnification
n	refractive index
NA	numerical aperture
z	distance along optical axis
λ	wavelength of light
I	irradiance, sometimes called intensity

28.2 INTRODUCTION

The optical principles and basic lens design needed to generate a diffraction-limited, highly magnified image with the light microscope were already essentially perfected a century ago. Ernst Abbe demonstrated how a minimum of two successive orders of diffracted light had to be captured in order for a particular spacing to be resolved (see historical sketch about Abbe principle¹). Thus, he explained and demonstrated with beautiful experiments the role of the wavelength of the imaging light and the numerical aperture ($NA = n \sin \Theta$, Fig. 1)² of the objective and condenser lenses on the resolving power of the microscope. In general, the minimum spacing δ for line gratings that can just be resolved cannot be smaller than

$$\delta = \frac{\lambda}{2NA} \quad (1)$$

when the NA of the condenser is equal to the NA of the objective.

For generating an image, contrast is just as important as resolution. Much of the early use of the light microscope depended on the relatively high image contrast that could be generated by differential absorption, scattering, reflection, birefringence, and the like due to specimen composition

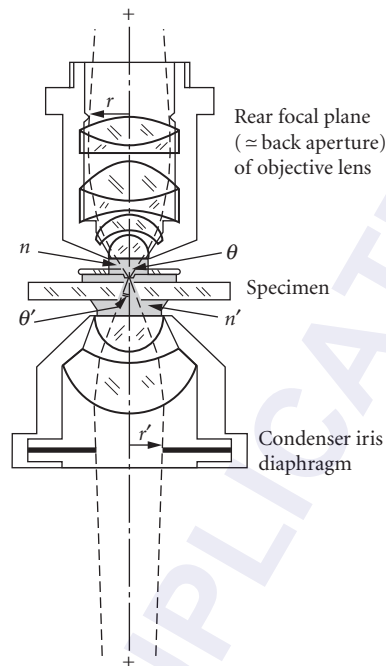


FIGURE 1 Definition of numerical aperture of objective ($NA_{\text{obj}} = n \sin\theta$) and condenser ($NA_{\text{cond}} = n' \sin\theta'$).²

or structure. Specimens, such as unstained living cells and other transparent objects introducing small optical path differences, were generally not amenable to direct microscopic observation for they would not produce detectable image contrast when brought to exact focus.

These impediments were removed by Zernike who showed how contrast in the microscope image is generated by interference between the light waves that make up the direct rays (that are undeviated by the specimen) and those that were scattered and suffered a phase difference by the presence of the specimen. Using this principle, Zernike invented the phase-contrast microscope.³ For the first time it became possible to see, in focus, the image of small, nonabsorbing objects. Zernike's revelations, together with Gabor's further contributions,⁴ not only opened up opportunities for the design of various types of interference-dependent image-forming devices but, even more importantly, improved our understanding of the basic wave optics involved in microscope image formation.

About the same time as Zernike's contributions, perfection of the electron microscope made it possible to image objects down to the nanometer range, albeit necessitating use of a high-vacuum environment and other conditions compatible with electron imaging. Thus, for four decades following World War II, the light microscope in many fields took a back seat to the electron microscope.

During the last decades, however, the light microscope has reemerged as an indispensable, powerful tool for investigating the submicron world in many fields of application. In biology and medicine, appropriate tags, such as fluorescent tags, are used to signal the presence and location of selected molecular species with exceptionally high sensitivity. Dynamic behaviors of objects far below the limit of resolution are visualized by digitally enhanced video microscopy directly in their natural (e.g., aqueous) environment. Very thin optical sections are imaged by video microscopy, and even more effectively with confocal optics. Quantitative measurements are made rapidly with the aid of digital image analysis.

At the same time, computer chips and related information-processing and storage devices, whose availability in part has spurred the new developments in light microscopy, are themselves miniaturized to microscopic dimensions and packaged with increasingly higher density. These electronic and photonic devices in turn call for improved means for mass manufacturing and inspection, both of which require advanced microscope optics.

Driven by the new needs and aided in part by computerized ray tracing and the introduction of new optical materials, we see today another epochal advance in the quality of lens design. The precision and remote control capabilities of mechanical components are also steadily improving. Furthermore, we may expect another surge of progress, hand-in-hand with development of improved electro-optical and electromechanical devices, in regulated image filtration, contrast-generating schemes, as well as in optical manipulation of the specimen employing microscope optics.

There are a number of excellent review articles and books discussing the optical principles of light microscopes^{1,5,6} and microscopic techniques,^{2,7-10} and their applications.¹¹⁻¹⁴ Among the many resources on microscopy available on the Internet, the Molecular Expressions website (<http://www.microscopy.fsu.edu/index.html>) stands out for its comprehensive treatment, beautiful illustrations, and interactive tutorials on the subject.

The present chapter is intended in part to bridge the territories of the manufacturer and the user of the microscope, including those who incorporate microscope optics into other equipment or apply them in unconventional ways. In this revision for the third edition of the *Handbook of Optics*, we reorganized the material, expanded the description of techniques that are typically covered only in passing by recent reviews and books on microscopy (e.g. interference and polarization microscopy), and added brief descriptions of imaging modes that are based on new optical concepts or new approaches to extract quantitative information from traditional imaging modes.

Many of the optical concepts and techniques, which are introduced here in the context of microscopy, are discussed in more detail in other chapters of this *Handbook*. On general optical considerations consult the *Handbook* chapters in this volume, “General Principles of Geometrical Optics” (Chap. 1) and on optical elements, such as “Lenses” (Chap. 17), “Polarizers” (Chap. 13), as well as chapters on physical optics for wave phenomena such as “Interference” (Chap. 2), “Diffraction” (Chap. 3), “Coherence Theory” (Chaps. 5 and 6), and “Polarization” (Chap. 12) which, as phenomena, are essential to the workings of the various contrast modes of the microscope. Material on image detection and processing can be found in *Handbook* chapters on vision in Vol. III, imaging detectors in Vol. II, and optical information and image processing in Chap. 11 of this volume.

28.3 OPTICAL ARRANGEMENTS, LENSES, AND RESOLUTION

Optical Arrangements

Geometric Optical Train, Magnification, Conjugate Planes In the optical train of a compound microscope (Fig. 2) invented by Galileo around 1610, the objective lens L_{ob} projects an inverted, real, magnified image O' of the specimen O (or object plane) into the intermediate image plane (or primary image plane). The intermediate image plane is located at a fixed distance $f' + z'$ behind L_{ob} , where f' is the back focal length of L_{ob} and z' is the optical tube length of the microscope. In general, O' is an aerial image for which an ocular L_{oc} (or the eyepiece) acts as a magnifier in front of the eye. Since L_{oc} , coupled with the corneal surface and lens of the eye, produces an erect image O'' of O on the retina, the object appears inverted to the observer. The ocular may also be used to project the image onto a screen. The aerial image at O' can also be exposed directly onto conventional film or an electronic sensor.

Continuing with the schematic diagram in Fig. 2, using thin-lens approximations, O is placed at a short distance z just outside of the front focal plane of L_{ob} , such that $z + f = a$, where f is the front focal length of L_{ob} and a is the distance between O and L_{ob} . O' is formed at a distance $b = (z' + f')$ behind L_{ob} . For a height y of O , the image height $y' = y \times b/a$. Thus, L_{ob} magnifies O by $M_{ob} = b/a$.

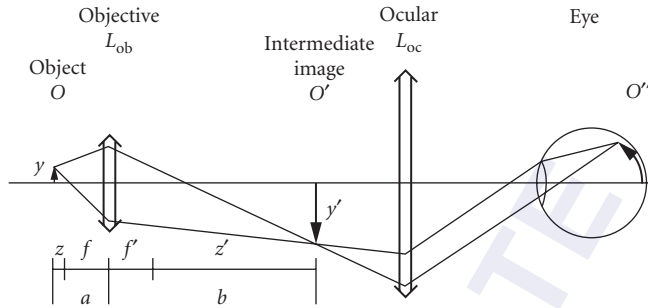


FIGURE 2 Ray path in the microscope from object to observer's eye (see text).

Also, $M_{ob} = f/z = z'/f'$. M_{ob} is the transverse or lateral magnification of L_{ob} . In the case of an infinity-corrected objective (Fig. 3), M_{ob} is the ratio f_{tb}/f , with f_{tb} the focal length of the specific tube lens L_{tb} . In turn, y' is magnified by L_{oc} by a factor $M_{oc} = 25 \text{ cm}/f_{oc}$, where f_{oc} is the focal length of the ocular (in cm) and 25 cm is the so-called near distance from the observer's eye (see Vol. II of this *Handbook*). Thus, the total transverse magnification of the microscope $M_{tot} = M_{ob} \times M_{oc}$.

Note that most microscope objectives are corrected for use only within a narrow range of image distances, and, in case of older style objectives, only in conjunction with specific groups of oculars. M_{ob} , which is the magnification inscribed on the barrel of the objective lens, is defined for its specified tube length (for high-power objectives, $M_{ob} = z'/f$) or, in case of infinity-corrected objectives, when used together with its specified tube lens. These factors, as well as those mentioned under "Microscope Lenses, Aberrations," must be kept in mind when a microscope objective is used as a magnifying lens, or in reverse as a high-numerical-aperture reducing lens, to form a truly diffraction-limited image.

Continuing the optical train back to the light source in a transilluminating microscope, Fig. 4a shows the ray paths and foci of the waves that focus on an on-axis point in the specimen. In Köhler illumination, the distance between the specimen and the condenser are adjusted so that the image of the field diaphragm in the illuminator is superimposed with the focused region of the specimen, and the lamp collector lens is adjusted so that the source image is focused in the plane of the condenser aperture diaphragm. Thus, \bar{O} , O , O' and O'' all lie in image planes that are conjugate with each other.

Tracing the rays emitted from a point in the light source (Fig. 4b), the rays are parallel between the condenser and the objective lenses. This situation arises because in Köhler illumination the light

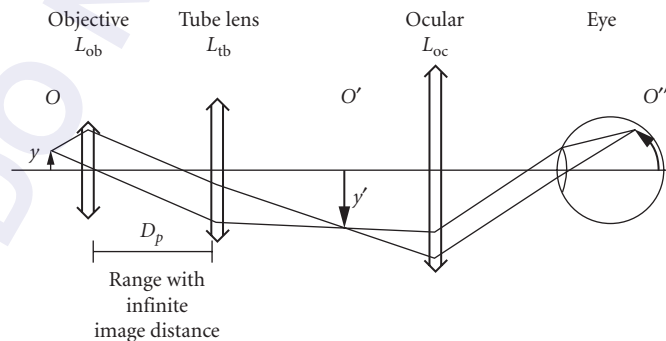


FIGURE 3 Ray path in microscope with infinity-corrected objective and tube lens.

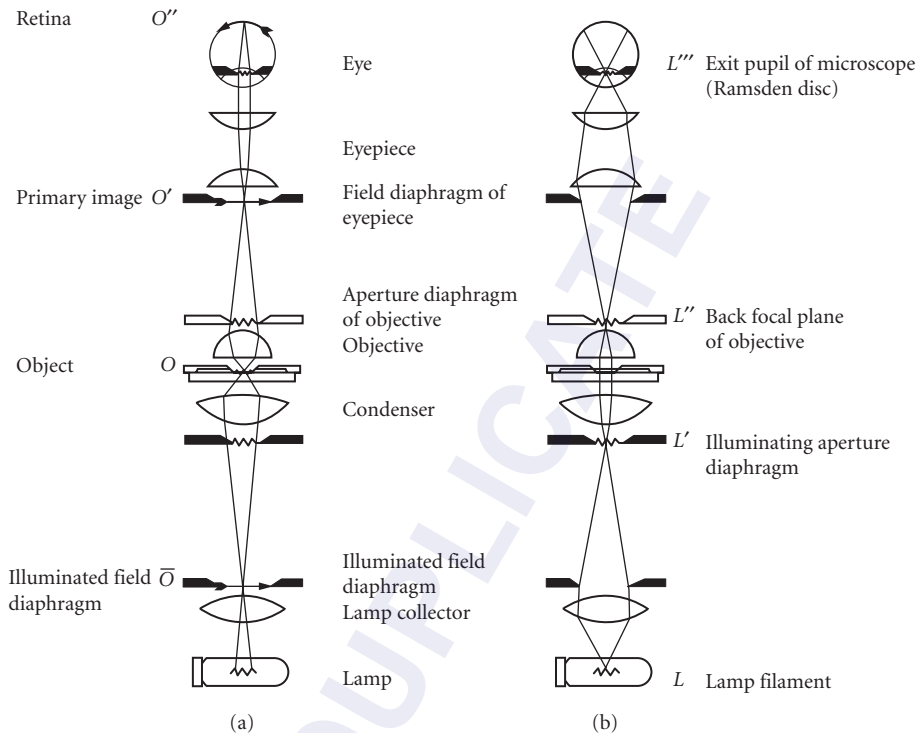


FIGURE 4 Ray paths in a transmitted light microscope adjusted for Köhler illumination. Two sets of conjugate planes are shown: set O in (a) is conjugate with the object O and with the field diaphragm planes; set L in (b) is conjugate with the lamp filament L and with aperture diaphragm planes.¹⁵

source (the filament of an incandescent bulb or the bright arc of a discharge lamp) is projected into the front focal plane of the condenser. Also, since the pupil of an (experienced) observer's eye is placed at the eyepoint or back focal plane of the ocular, the four aperture planes L , L' , L'' , and L''' are again conjugate to each other.

As inspection of Fig. 4a and b¹⁵ shows, the field planes and aperture planes are in reciprocal space relative to each other throughout the whole optical system. This reciprocal relationship explains how the various diaphragms and stops affect the cone angles, paths, and obliquity of the illuminating and image-forming rays, and the brightness, uniformity, and size of the microscope field. More fundamentally, a thorough grasp of these reciprocal relationships is needed to understand the wave optics of microscope image formation and for designing various contrast-generating devices and other microscope optical systems.

Transillumination The full impact of the illumination system on the final quality of the microscope image is often not appreciated by the microscope user or designer. Undoubtedly, part of this neglect arises from a lack of understanding of the roles played by these components, in particular the condenser, and the common practice of closing down the condenser iris diaphragm to adjust image contrast for comfortable viewing. Regardless of the conventional view, critical examination of the microscope image or point spread function reveals the importance of the alignment, focus, tilt, NA, and effective aperture function of the condenser. The effects are especially noticeable when contrast is enhanced, e.g., by video microscopy. A further illustration of the importance of the illumination on the resolving power of the light microscope can be found in the section on "Structured Illumination."

Ernst Abbe was the first to systematically analyze the resolving power of microscope optics by fabricating precision line gratings and imaging them in the microscope. As indicated earlier, a grating is resolved if the objective lens captures at least two successive diffraction orders which are typically the zero- and first-order diffraction. Abbe summarized his results in a simple expression, relating the minimum resolvable pitch δ to the numerical aperture of the objective and condenser lens:

$$\delta = \frac{\lambda}{\text{NA}_{\text{obj}} + \text{NA}_{\text{cond}}} \quad (2)$$

with λ the wavelength of light used. This formula can be derived by considering the diffraction of linear gratings that are illuminated obliquely. In the limiting case of zero condenser NA, the grating is illuminated coherently by a collimated beam of light that is parallel to the microscope's optical axis. The minimum resolvable pitch is proportional to the wavelength and inversely proportional to the objective NA. By increasing the condenser NA, oblique rays are added to the illuminating light, increasing the angular span between diffraction orders captured by the same objective lens, and thus decreasing the minimum resolvable pitch. By making the condenser and objective NA equal, the grating is effectively illuminated incoherently and Eq. (2) reduces to Eq. (1).

The influence of the condenser NA on resolving two nearby point objects was considered by Hopkins and Barham.¹⁶ They applied the Rayleigh criterion for resolving two pinholes that are equally bright and illuminated incoherently ($\text{NA}_{\text{cond}} = \text{NA}_{\text{obj}}$) and found a minimally resolved distance $d = 0.61\lambda/\text{NA}_{\text{obj}}$ (Fig. 5, $m = 1$).¹⁷ Distance d is a factor 1.22 larger than the limiting pitch of a grating illuminated and imaged by the same condenser and objective lens [Eq. (1)]. However, for the case of coherent illumination ($\text{NA}_{\text{cond}} = 0$), the minimal distance of two resolved points only increases by 40 percent instead of 100 percent, as is the case for gratings. Hopkins and Barham calculated a maximum resolution (minimal d) for $\text{NA}_{\text{cond}} = 1.5 \times \text{NA}_{\text{obj}}$. Such high NA_{cond} is usually not achievable for high-NA objective lenses, and, in addition, with most objectives, flare due to internal reflection would reduce image contrast to an extent possibly unsalvageable even with video contrast enhancement. Again, reduction of NA_{cond} , generally achieved by closing down the condenser iris diaphragm, tends

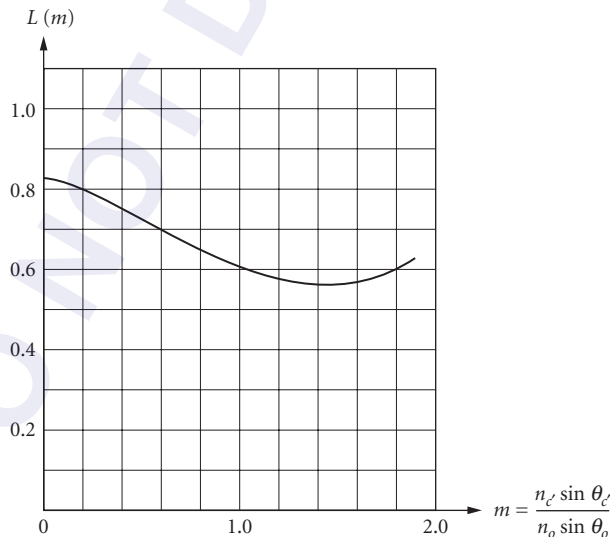


FIGURE 5 Effect of the condenser aperture on the resolution of two pinholes of equal brightness. m is the ratio of the numerical apertures of condenser to objective. L is the minimum resolved distance between the pinholes (Rayleigh criterion) in units of the wavelength divided by the objective aperture.¹⁷

to raise image contrast so that even experienced microscopists tend to use an $NA_{\text{cond}} \approx (0.3, \dots, 0.5) \times NA_{\text{obj}}$ to obtain a compromise between resolution and visibility. With video and other modes of electronic enhancement, the loss of contrast can be reversed so that improved lateral, and especially axial, resolution is achieved by using an NA_{cond} that equals, or nearly equals, the NA_{obj} .

Under optimum circumstances, the light source and condenser should be focused for Köhler illumination (Fig. 4) to minimize flare and to improve the homogeneity of field illumination. Alternately, image brightness, especially in the middle of the field, can be maximized by *critical illumination* where the condenser is somewhat defocused from Köhler illumination to produce an image of the source rather than the field diaphragm superimposed on the specimen. Either mode of illumination can yield resolution approximately as given by Eq. (2).

The aperture function of the microscope can become nonuniform, or limited, for a number of reasons. These include misalignment between the objective and condenser lenses; misalignment of the condenser iris (relative to the condenser lens elements); misalignment of the illuminator and condenser axes; tilted objective or condenser lenses or lens elements; nonuniform illumination of the condenser aperture; limited source size; nonuniform intensity distribution in the source; and improper choice, or focusing, of the condenser or source collector. Whether intentional or accidental, these conditions can reduce the effective NA_{cond} and/or induce oblique illumination, thus sacrificing resolution and image quality. An improvement, using a single optical fiber light scrambler, which allows the filling of the full condenser aperture with uniform illumination and little loss of field brightness (especially when using concentrated arc lamps) was introduced by Ellis¹⁸ (also see Figs. 3-13, 3-14 in Ref. 2).

Epi-Illumination In the epi-illumination mode, a beam splitter, part-aperture-filling mirror, or wavelength-discriminating dichromatic (unfortunately often called dichroic) mirror, placed behind the objective lens diverts the illuminating beam (originating from a light source placed in the side arm of the microscope) into the objective lens, which also acts as the condenser (Fig. 6).¹⁹ Alternatively, a second set of lenses and a beam-diverting mirror (both of whose centers are bored out and are arranged coaxially around the objective lens) can provide a larger NA-illuminating beam, much as in dark field illumination in the transillumination mode.

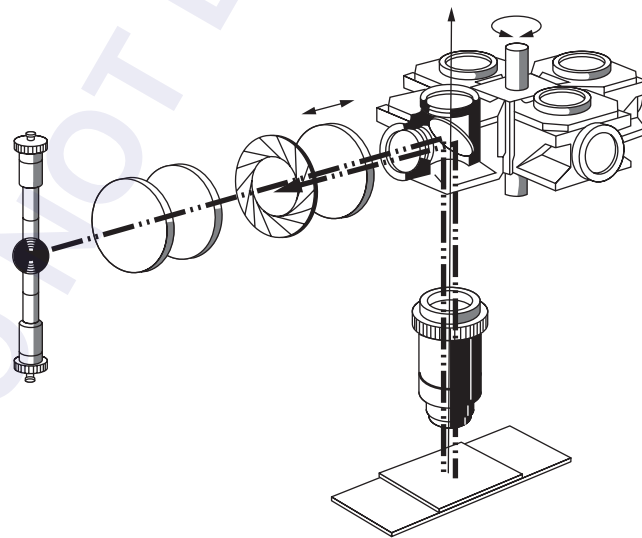


FIGURE 6 Schematic of epi-illuminating light path. The rotatable set of filter cubes with excitation filters, dichromatic mirrors, and barrier filters matched to specific fluorochromes are used in epifluorescence microscopy.¹⁹

This latter approach limits the maximum NA of the objective lens to around 1.25, but has the advantage that the illuminating beam traverses a path completely isolated from the image-forming beam. When the two beams do pass through the same objective lens, as is the case with most epi-illuminating systems, the lens elements must be carefully designed (by appropriate choice of curvature and use of highly efficient antireflection coating) to reduce hot spots and flare introduced by (multiple) reflection at the lens surfaces. Modern microscope objectives for metallurgical and industrial epi-illuminating systems in particular are designed to meet these qualities. In addition, circular polarizers (linear polarizer plus $\lambda/4$ wave plate) and appropriate stops are used to further exclude light reflected from the surfaces of lens elements, cover glass, and the like. For epi-illumination fluorescence microscopy, dichromatic beam splitters, and barrier filters can reduce background contamination that arises from the exciting beam to less than one part in 10^4 .

Orthoscopic versus Conoscopic Imaging The common mode of observation through a microscope is by orthoscopic observation of the focused image. For certain specific applications, particularly with polarizing microscopes, examination of the aperture plane, or conoscopic observation, sheds valuable complementary information.

Conoscopic observation can be made either by replacing the regular ocular with a telescope that brings the aperture plane into focus or by inserting a Bertrand lens (that serves as a telescope objective) in front of a regular ocular. Conversely, one can observe the aperture plane simply by removing the ocular and looking down the microscope body tube (in the absence of a Bertrand lens) or by examining the Ramsden disk above the ocular with a magnifier. Levoy and Oldenbourg used a microlens array for generating a hybrid image that consists of an array of small conoscopic images, each sampling a different object area.^{20,21}

The polar coordinates of each point in the aperture plane, that is the radius r and azimuth angle α are related to the rays traversing the specimen by: $r = \sin \theta$ and $\alpha =$ azimuth orientation of the ray projected onto the aperture plane (Fig. 7). Thus, conoscopic observation provides a plane projection of all of the rays traversing the specimen in three-dimensional space. For specimens, such as single crystal flakes or polished mineral sections in which a single crystal is illuminated (optically isolated)

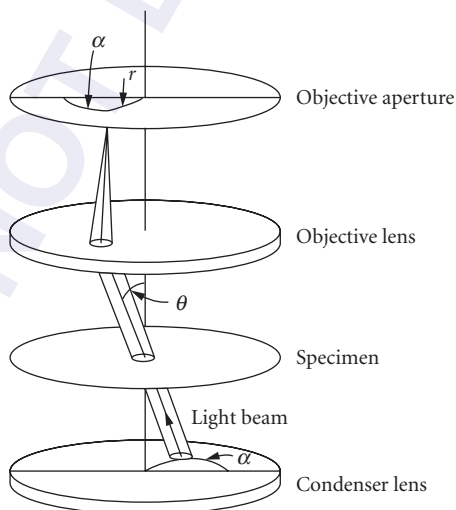


FIGURE 7 Parallel rays with inclination θ and azimuth orientation α traversing the specimen plane, and focused by the objective lens at a point with radius r and same azimuth angle α in the aperture plane.

by closing down the field diaphragm, the conoscopic image reveals whether the crystal is uniaxial or biaxial, its optic axis angle and directions, as well as sign and strength of birefringence and other anisotropic or optically active properties of the crystal.²²

Conoscopic observation also reveals several attributes of the condenser aperture plane and its conjugate planes (e.g., in Köhler illumination, the plane of the condenser iris diaphragm and the illuminating source). Thus, conoscopic observation can be used for checking the size, homogeneity, and alignment of the illuminating light source as well as the size and alignment of the condenser iris diaphragm and phase-contrast annulus (located at the front focal plane of the condenser) relative to the objective exit pupil or the phase ring (located at the back focal plane of the objective). It also reveals the state of extinction in polarized light and interference-contrast microscopy and provides a visual estimate of the aperture transfer function for the particular optical components and settings that are used.

The aperture plane of the microscope is also the Fourier plane of the image, so that diffraction introduced by periodic textures in the specimen can be visualized in the aperture plane by conoscopic observation. Depending on the NA of the objective and the spatial period in the specimen, the pattern of diffraction up to many higher orders can be visualized in the aperture plane when the condenser iris is closed down to illuminate the specimen with a parallel beam of light. Closing down the condenser iris restricts the zero-order light to a small area in the aperture plane and higher-order diffraction maxima produce additional images of the diaphragm displaced in the directions of the periodic texture in the specimen.

Microscope Lenses, Aberrations

Objective Lenses With few exceptions, microscope objective lenses are designed to form a diffraction-limited image in a specific image plane that is located at a fixed distance from the objective lens (or from the tube lens in the case of an infinity-focus system). The field of view is often quite limited, and the front element of the objective is placed close to the specimen with which it must lie in optical contact through a medium of defined refractive index n , usually air ($n = 1$, dry objectives), water ($n = 1.33$, water immersion objectives), oil ($n = 1.52$, oil immersion objectives) or other high refractive index media.

Depending on the degree of correction, objectives are generally classified into achromats, fluorites, and apochromats with a plan designation added to lenses with low curvature of field and distortion (Table 1). Some of these characteristics are inscribed on the objective lens barrel, such as Plan Apo 60/1.40 oil 160/0.17, meaning 60 power/1.40 NA Plan Apochromatic objective lens designed to be used with oil immersion between the objective front element and the specimen, covered by an 0.17-mm-thick coverslip, and used at a 160-mm mechanical tube length. Another example might be Epiplan-Neofluar 50 \times /0.85 ∞ /0, which translates to Plan “Fluorite” objective designed for epi-illumination (i.e., surface illumination of specimen through the objective lens rather than through a separate condenser) with a 50 \times magnification and 0.85 NA to be used in air (i.e., without added immersion medium between the objective front element and coverslip or specimen), with no coverslip, and an (optical) tube length of infinity. “Infinity-corrected” objectives require the use of

TABLE 1 Objective Lens Types and Corrections

Type	Spherical	Chromatic	Flatness
Achromat	*	2λ	No
F-achromat	*	2λ	Improved
Neofluar	3λ	$< 3\lambda$	No
Plan-neofluar	3λ	$< 3\lambda$	Yes
Plan apochromat	4λ	$> 4\lambda$	Yes

* = corrected for two wavelengths at two specific aperture angles.

2λ = corrected for blue and red (broad range of visible spectrum).

3λ = corrected for blue, green, and red (full range of visible spectrum).

4λ = corrected for dark blue, blue, green, and red.

Source: Zeiss publication #41-9048/83.

TABLE 2 Common Abbreviations Designating Objective Lens Types

DIC, NIC	Differential (Nomarski) interference contrast
L, LL, LD, LWD, ELWD, ULWD	Long working distance (extra-) (ultra-)
FL, FLUOR, NEOFLUOR, FLUOTAR	With corrections as with "fluorite" objectives but no longer implies the inclusion of fluorite elements
PHASE, PHACO, PC, PH 1, 2, 3, etc.	Phase contrast, using phase condenser annulus 1, 2, 3, etc.
DL, DM, PLL, PL, PM, PH, NL, NM, NH	Phase contrast: dark low, dark medium, positive low low, low, medium, high contrast (regions with higher refractive index appear darker); negative low, medium, high contrast (regions with higher refractive index appear lighter)
PL, PLAN; EF	Flat field; extended field (larger field of view but not as high as with PLAN, achromats unless otherwise designated)
PLAN APO	Flat field apochromat
NPL	Normal field of view plan
P, PO, POL	Low birefringence, for polarized light
UV	UV transmitting (down to approx. 340 nm), for UV-excited epifluorescence
ULTRAFLUAR	Fluorite objective for imaging down to approx. 250 nm in UV as well as in the visible range
CORR, W/CORR	With correction collar
I, IRIS, W/IRIS	Adjustable NA, with iris diaphragm built into back focal plane
M	Metallographic
NC, NCG	No coverslip
EPI	Surface illumination (specimen illuminated through objective lens), as contrasted to dia- or transillumination
BD, HD	For use in bright or darkfield (hell, dunkel)
CF	Chrome-free (Nikon: objective independently corrected longitudinal chromatic aberrations at specified tube length)
ICS	Infinity color-corrected system (Carl Zeiss: objective lens designed for infinity focus with lateral and longitudinal chromatic aberrations corrected in conjunction with a specified tube lens)
OIL, HI, H; WATER, W; GLY	Oil immersion, Homogeneous immersion, water immersion, glycerol immersion
U, UT	Designed to be used with universal stage (magnification/NA applies for use with glass hemisphere; divide both values by 1.51 when hemisphere is not used)
DI; MI; TI Michelson	Interferometry: noncontact; multiple-beam (Tollanski)
ICT; ICR	Interference contrast: in transillumination; in reflected light

a designated tube lens to eliminate residual aberrations and to bring the rays to focus into the image plane. Several other codes are inscribed or color-coded on microscope objectives (Tables 2 and 3).

Older style objective lenses are designed to be used with a specified group of oculars or tube lenses that are placed at specific distances in order to remove residual errors. For example, compensation oculars were used in conjunction with apochromatic and other high-NA objectives to eliminate lateral chromatic aberration and improve flatness of field. However, modern style objectives together with their tube lenses are typically fully corrected so as not to require additional chromatic or other type corrections.

Coverslip Correction For objective lenses with large NAs, the optical properties and thicknesses of the media lying between its front element and the specimen critically affect the calculations

TABLE 3 Color-Coded Rings on Microscope Objectives

Color code (narrow colored ring located near the specimen end of objective)	
Black	Oil immersion
Orange	Glycerol immersion
White	Water immersion
Red	Special
Magnification color code (narrow band located further away from specimen than immersion code)	
Color	Magnification
Black	1, 1.25, 1.5
Brown	2, 2.5
Red	4, 5
Yellow	10
Green	16, 20
Turquoise blue	25, 32
Light blue	40, 50
Cobalt (dark) blue	60, 63
White (cream)	100 and higher

needed to satisfy the aplanatic and sine conditions and otherwise to correct for image aberrations. For homogeneous immersion objectives (that are designed to be used with the refractive indices and dispersion of the immersion oil, coverslip, and medium imbibing the specimen, all matched to that of the objective lens front element), the calculation is straightforward since all the media can be considered an extension of the front lens element.

However, with nonimmersion objectives, the cover glass can become a source of chromatic aberration, which is worse the larger the dispersion and the greater the thickness of the cover glass. The spherical aberration is also proportional to the thickness of the cover glass. In designing objectives not to be used with homogeneous immersion, one assumes the presence of a standard cover glass and other specific optical media between the front lens element and the specimen. As one departs from these designated conditions, spherical aberration (and also coma) increases with the NA of the lens, since the difference between the tangent and sine of the angle of incidence is responsible for departure from the needed sine condition.

It should also be noted that oil immersion objectives fail to provide full correction, or full NA, when the specimen is mounted in an imbibing medium with a different refractive index, for example aqueous media, even with the objective and cover glass properly oil-contacted to each other. With such an arrangement, the diffraction image can degrade noticeably as one focuses into the specimen by as little as a few micrometers.²³ Special water immersion objectives (e.g., Nikon Plan Apo 60×/1.2 NA and short-wavelength transmitting Fluor 40×/1.0 NA, both with collar to correct coverslip thickness deviation from 0.17 mm) overcome such aberrations, even when the specimen is imaged through an aqueous medium of 200- μm thickness.

For lenses that are designed to be used with a standard coverslip of 0.17-mm thickness (and $n_D = 1.515$), departure from standard thickness is not overly critical for objectives with NA of 0.4 or less. However, for high-NA, nonhomogeneous immersion lenses, the problem becomes especially critical so that even a few micrometers' departure of the cover glass thickness degrades the image with *high-dry objectives* (i.e., nonimmersion objectives with high NA) of NA above 0.8 (Fig. 8).²⁴ To compensate for such error, well-corrected, high-dry objectives are equipped with correction collars that adjust the spacing of their intermediate lens elements according to the thickness of the cover glass. Likewise, objective lenses that are made to be viewed through layers of silicon or plastic, or of different immersion media (e.g., water/glycerol/oil immersion lenses), are equipped with correction collars.

The use of objective lenses with correction collars does, however, demand that the observer is experienced and alert enough to reset the collar using appropriate image criteria. Also, the focus tends to shift, and the image may wander, during adjustment of the correction collar. Figure 9 shows an

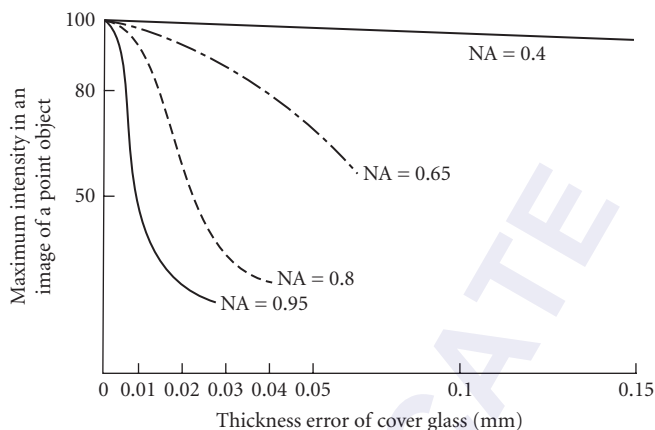


FIGURE 8 Calculated maximum intensity in the image of a point object versus the deviation of the coverglass thickness from the ideal thickness.²⁴

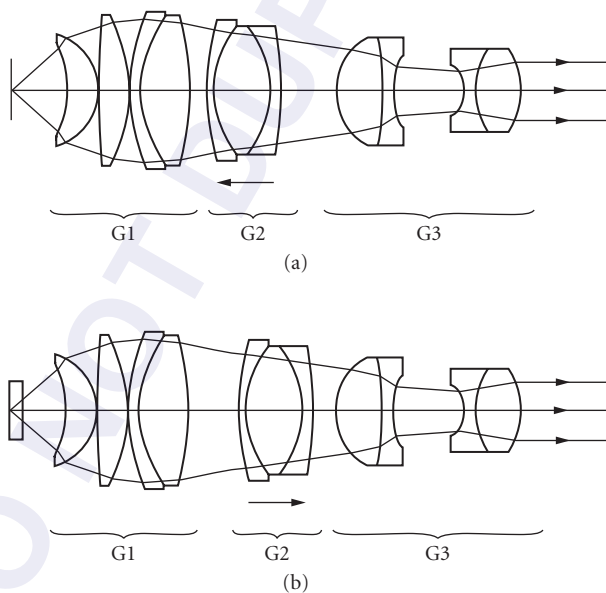


FIGURE 9 High-dry objective lens ($60\times/0.7$ NA) equipped with a correction collar for (a) focusing at the surface or (b) through plane glass of up to 1.5-mm thickness. The lens group G_2 is moved forward to enhance the spherical and chromatic correction by G_1 and G_2 when focused on the surface, while it is moving backward to compensate for the presence of the glass layer when focusing deeper through the glass.²⁴ (U.S. Patent 4666256.)

TABLE 4 Reference Focal Lengths for Infinity-Focused Objective Lenses

Leica	200 mm	B, M
Olympus	180 mm	B, M
Carl Zeiss	164.5 mm	B, M
Nikon	200 mm	B, M

B = biological, M = metallurgical.

example of a 60/0.7 objective lens equipped with a correction collar for focusing at the surface or through a cover glass of up to 1.5-mm thickness without altering the focal setting of the lens.

Tube Lengths and Tube Lenses for Which Microscope Objectives Are Corrected For finite-focused “biological” objective lenses, most manufacturers had standardized the mechanical tube length to 160 mm. More recently most manufacturers have switched to infinity focus for their biomedical and metallurgical microscopes.

For infinity-focused objective lenses, the rays emanating from a given object point are parallel between the objective and tube lens. Since the physical distance (D_p , Fig. 3) and optical path length between the objective and tube lens are not critical, optical plane-parallel components, such as compensators, analyzers, and beam splitters, can be inserted in this space without altering the objective’s corrections. The tube lens focuses the parallel rays onto the intermediate image plane.

The magnification of an infinity-focused objective lens is calculated by dividing the focal length of the tube lens (also called reference focal length) by the focal length of the objective lens. The reference focal lengths adopted by several manufacturers are listed in Table 4.

Working Distance Microscope objectives are generally designed with a short free working distance, that is the distance from the front element of the objective lens to the surface of the cover glass or, in the case of lenses that are designed to be used without cover glass, to the specimen surface. For some applications, however, a long free working distance is indispensable, and special objectives are designed for such use despite the difficulty involved in achieving large numerical apertures and the needed degree of correction.

Field Size, Distortion The diameter of the field in a microscope is expressed by the field-of-view number, or simply field number, which is the diameter of the field in millimeters measured in the intermediate image plane. The field size in the object plane is obviously the field number divided by the magnification of the objective. While the field number is often limited by the magnification and field stop of the ocular, there is clearly a limit that is also imposed by the design of the objective lens. In early microscope objectives, the maximum usable field diameter tended to be about 18 mm or considerably less, but with modern plan apochromats and other special flat field objectives, the maximum usable field can be as large as 28 mm or more. The maximum useful field number of objective lenses, while available from the manufacturers, is unfortunately not commonly listed in microscope catalogs. Acknowledging that these figures depend on proper combination with specific tube lenses and oculars, we should encourage listing of such data together with, for example, UV transmission characteristics (e.g., as the wavelength at which the transmission drops to 50 percent, or some other agreed upon fraction).

Design of Modern Microscope Objectives Unlike earlier objective lenses in which the reduction of secondary chromatic aberration or curvature of field were not stressed, modern microscope objectives that do correct for these errors over a wide field tend to be very complex. Here we shall examine two examples, the first a 60/1.40 Plan Apochromat oil-immersion lens from Nikon (Fig. 10).²⁴

Starting with the hyperhemisphere at the front end (left side of Fig. 10) of the objective, this aplanatic element is designed to fulfill Abbe’s sine condition in order to minimize off-axis spherical aberration and coma, while providing approximately half the total magnifying power of the objective (Fig. 11). In earlier designs, the hyperhemisphere has been made with as small a radius

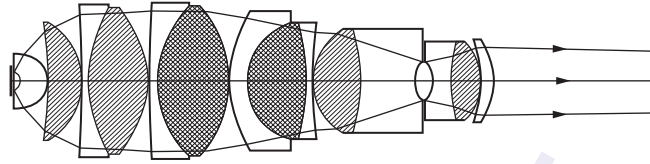


FIGURE 10 Design of Nikon Plan Apochromat oil-immersion objective with 60 \times magnification and 1.40 NA.²⁴

as possible in order to maximize its magnifying power and to minimize its spherical and chromatic aberrations, since these aberrations increase proportionally with the focal length of the lens. Modern demands for larger field size and reduced curvature of field, however, introduce a conflicting requirement, namely, the need to maintain as large a radius as practical in order to minimize the hyperhemisphere's contribution to the Petzval sum (the algebraic sum of the positive and negative curvatures multiplied by the refractive indices of the lens elements).²⁵ The hyperhemisphere in these Plan Apochromats is made with a high-index, low-dispersion material to compensate for the greater radius. Additionally, a negative meniscus is generated in the front surface of the hyperhemisphere to which is cemented a minute, plano-convex lens. The negative curvature in the hyperhemisphere contributes to the reduction of the Petzval sum. At the same time the minute plano-convex lens protects the material of the hyperhemisphere which is less resistant to weathering. Index matching between the minute plano-convex lens and immersion oil eliminates or minimizes the refraction and reflection at the lens-oil interface and provides maximum transmission of the all-important high-NA rays into the objective lens. The index matching also reduces the influence of manufacturing errors of this minute lens element on the performance of the objective.

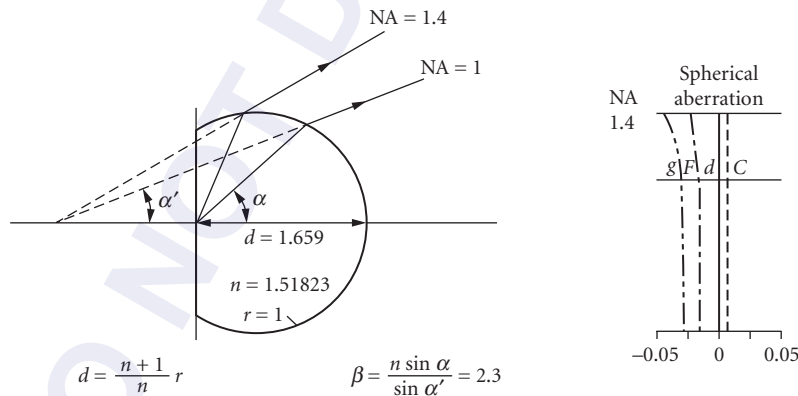


FIGURE 11 Aplanatic condition of the hyperhemisphere placed at the front end of an oil-immersion objective. The front lens has the same refractive index as the coverglass and immersion oil. The aplanatic condition describes the necessary relationship between refractive index n , distance d between object and spherical surface, and radius r of the spherical surface, in order to make all rays emanating from an object point on the axis leave the hemispherical surface after refraction without introducing spherical aberration. According to the sine condition, the magnification β has to be constant for all angles α . On the right, the small amount of longitudinal spherical aberration and chromatic deviation due to dispersion from the ideal focus point of the hyperhemisphere is shown for different wavelengths ($\lambda_c = 656$ nm, $\lambda_d = 588$ nm, $\lambda_f = 486$ nm, $\lambda_g = 436$ nm). Abscissa: longitudinal deviation on lens axis. Ordinate: numerical aperture from 0 (lens axis) to 1.4 NA.²⁴

The low-dispersion-glass singlet behind the aplanatic hyperhemisphere further reduces the cone angle of the rays entering the doublets that follow, allowing these and the subsequent lenses to concentrate on correcting axial and lateral chromatic aberration as well as curvature of field. These errors, as well as residual spherical aberrations, are corrected by inclusion of low-dispersion positive and high-dispersion negative lens elements, use of thick-lens elements, appropriate placement of positive and negative lens curvatures, and through extensive ray tracing. Near the exit pupil, the height of the ray paths through the concave surfaces is reduced in order to generate additional negative values that minimize the Petzval sum (to complement the inadequate negative contribution made by the concave surface in the hyperhemisphere), so that field flatness can be improved without overly reducing the objective lens' magnifying power or adding to its spherical aberration.

In reality, the Petzval sum of the objective as a whole is made somewhat negative in order to compensate for the inevitable positive Petzval sum contributed by the ocular. Thus, the image at the intermediate image surface, especially the sagittal surface of modern objectives, bows away from the object. Unless the image area is relatively small, one needs to use specified oculars in order to attain maximum field flatness combined with optimum correction otherwise.

Unlike earlier objective lenses whose design did not appreciably vary from one manufacturer to another, the design of lenses in modern microscope objectives can vary considerably. For example, compare the Nikon Chrome Free 60/1.4 Plan Apo objective discussed above and the Zeiss Infinity Color-Corrected Systems 63/1.4 Plan Apo objective in Fig. 12. Both are excellent, state-of-the-art lenses. But in addition to general design philosophy, including the decision to avoid or to use tube lenses to achieve full chromatic corrections, other factors such as choice of optical elements with special dispersion characteristics; degrees of UV transmission; freedom from fluorescence, birefringence, aging loss of transmittance, and the like all affect the arrangement of choice.

While a modern research-grade microscope is corrected to keep the aberrations from spreading the image of a point source beyond the Airy disk, geometrical distortion of the image formed by microscope objectives tends not to be as well-corrected (e.g., compared to photographic objectives at the same picture angle). Thus, in objectives for biological use, pincushion distortions of up to 1 percent may be present. However, in objectives that are designed for imaging semiconductors, the distortion may be as low as 0.1 percent and they can be considered nearly distortion-free. To reduce stray light and flare, modern microscope objectives contain lens elements with carefully tuned, antireflection coatings, and lens curvatures are selected to minimize ghost images arising from multiple reflections.

Given the sophisticated design to provide a wide flat field, with spherical aberrations corrected over a broad wavelength range, and with low longitudinal as well as chromatic aberrations corrected at high NA, the aberration curves of these modern microscope objectives no longer remain simple cubic curves, but turn into complex combinations of higher-order curves (Fig. 13).

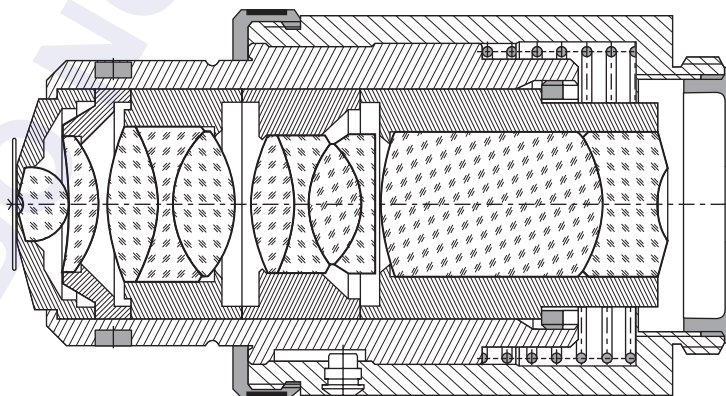


FIGURE 12 Carl Zeiss Infinity Color-Corrected 63/1.4 Plan Apo objective. (Courtesy of E. Keller, Carl Zeiss, N.Y.)

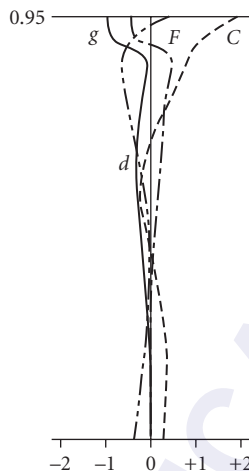


FIGURE 13 Spherical aberration curves for spectral lines (*C*, *d*, *F*, and *g*) of a highly corrected modern microscope objective with a high numerical aperture.²⁴ Ordinate: numerical aperture from 0 (lens axis) to $NA = 0.95$. Abscissa: longitudinal deviation of focal distance on lens axis indicated in millimeters. (The depth of focus for a $40\times/0.95$ NA apochromatic objective is approximately 1 mm which corresponds to a depth of field of 0.6 μm in specimen space.)

Oculars As conventionally illustrated, the ocular in a light microscope further magnifies the primary (intermediate) image formed by the objective lens (Figs. 2 and 3). The ocular can also be viewed as the front elements of a macro (relay) lens system made up of the ocular plus the refractive elements of the viewer's eye (Fig. 4a) or a video or photographic camera lens. Special video and photo oculars combine these functions of the ocular plus the video or photo lenses into single units.

The intermediate image plane (that lies between the lenses in many ocular types or precedes the lens elements in the Ramsden-type oculars), or its conjugate plane is used to place field-limiting stops, iris diaphragms, reticles, micrometer scales, comparator beam splitters, and the like that need to appear in the same focal plane as the specimen.

The Ramsden disk, the exit pupil of the objective lens imaged by the ocular, generally appears a short distance above the ocular (Fig. 4b). Since the Ramsden disk should lie in the observer's pupil, special high-eye-point oculars are provided for the benefit of observers wearing corrective eye glasses (especially those for astigmatism). High-eye-point oculars are also used for inserting beam-deviating devices (such as the scanning mirrors in laser scanning confocal microscopes) or aperture-modifying devices (such as aperture occluders for stereo viewing through single objective binocular microscopes²).

The magnification of an ocular is defined as 25 cm divided by the ocular's focal length. On the ocular, the magnification and field number are inscribed (e.g., as $10\times/20$, meaning 10-power or 25-mm focal length with a field of view of 20-mm diameter), together with manufacturer's name and special attributes of the ocular such as chromatic-aberration-free (CF), wide-field (W, WF, EWF), plan (P, Pl), compensation (Comp, C, K), high-eye-point (H, picture of glasses), with cross hair and orientation stub for crystallography (pol), projection (pro), photographic (photo), video (TV), and the like. Also, special oculars provide larger and flatter fields of view (designated wide field, extra wide field, plan, periplan, hyperplan, etc., some with field numbers ranging up to 28 mm).

Compared to microscope objective lenses, fewer design standards have been adopted and fewer standard abbreviations are used to designate the performance or function of the oculars. Two physical parameters of the oculars have, however, become more or less standardized. The outside diameter of the ocular is either 23.2 mm or 30.0 mm, and the reference distance, or the parafocalizing distance of the ocular (i.e., the location of the intermediate image plane below the flange of the ocular) is now generally set to 10 mm.

In the past, oculars with wide ranges of incremental magnifications were provided to adjust the total image magnification of the microscope, but this practice is now replaced by the use of much fewer, better-corrected oculars coupled with a telan magnification changer in the microscope's body tube, or a zoom projection ocular.

Factors affecting choice of ocular focal length and magnification include optimizing the microscope total magnification and image resolution to match the MTF characteristics of the detector and to adjust the available field coverage. In video-enhanced fluorescence, differential interference contrast, polarizing, dark field, and the like microscopy, the total magnification often needs to be raised beyond the classical "empty magnification" limit, in order to be able to visualize minute objects whose diameters lie well below the limit of microscope resolution.² However, depending on the MTF characteristics, sensitivity, and total pixels available in the sensor, conflicts may arise between the need for greater magnification, image brightness, and field coverage. To optimize the total image magnification, fine trimming of the ocular magnification may be needed, in addition to choosing an objective with the appropriate magnification and NA-to-magnification ratio. Zoom oculars are especially suited for fine-tuning the magnification to optimize S/N ratio and image integration time in video microscopy. For very low light level images, for example in photon-counting imaging, ocular magnifications of less than one may be needed in order to sufficiently elevate the S/N ratio, albeit at a sacrifice to spatial resolution.

In addition to adjusting image magnification and placing the microscope's exit pupil at a convenient location, the ocular compensates for the aberrations that have not been adequately corrected in the objective and tube lens. Huygens oculars combined with lower-power achromatic objectives, and compensating oculars combined with higher-NA achromatic and apochromatic objectives, correct for lateral chromatic aberration. Some higher-NA achromatic objectives are purposely designed to provide residual aberrations (including field curvature) that are similar to those in the apochromats, so that the same compensation oculars can be used to compensate for both types of objectives.

Modern objectives used with the appropriate tube lens are sufficiently well corrected to require minimum or no compensatory correction by the oculars. In research-grade microscopes, the image projected by the objective and tube lens is often recorded directly by placing an electronic image sensor into the intermediate image plane. With objectives that are designed to produce well-corrected intermediate images, oculars themselves are made independently free of lateral and longitudinal chromatic and some spherical aberrations. Regardless of the degree of correction relegated to the ocular, modern microscopes provide images with color corrections, fields of view, and flatness of field much superior to earlier models.

Resolution

Airy Disk and Lateral Resolution Given a perfect objective lens and an infinitely small point of light residing in the specimen plane, the image formed in the intermediate image plane by the objective lens is not another infinitely small point, but a diffraction image with a finite spread (Fig. 14a). This Airy diffraction image is the Fraunhofer diffraction pattern formed by the exit pupil of the objective lens from which spherical waves converge to the focal point. The distribution of irradiance of the diffraction image (Fig. 14b)²⁶ is given by an expression containing the first-order Bessel function $J_1(v)$:

$$I(v) = I_0 \left(\frac{2J_1(v)}{v} \right)^2 \quad (3)$$

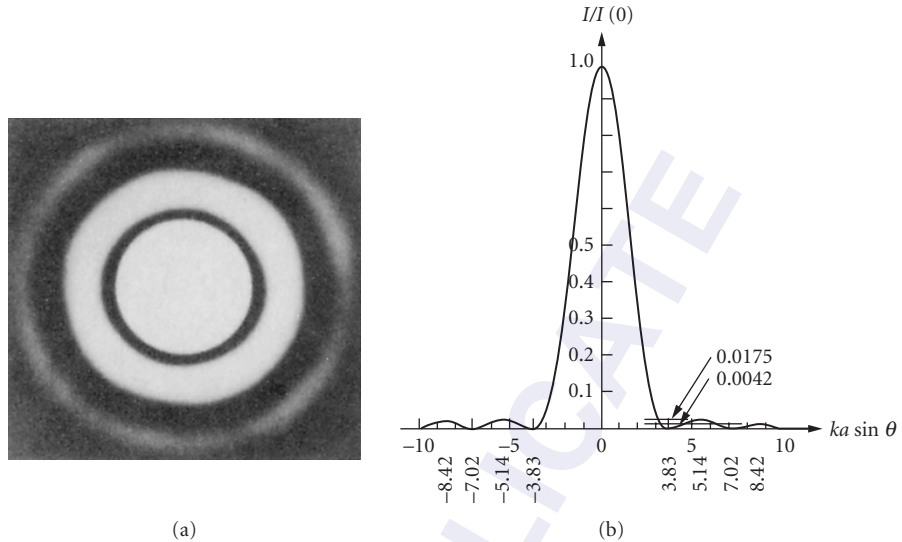


FIGURE 14 Airy pattern of circular aperture: image (a) of central Airy disk, first dark ring and subsidiary maximum and graph (b) of radial intensity distribution.²⁶

with v proportional to the diffraction angle. If the irradiance is calculated as a function of radius measured from the center of the Airy diffraction pattern located in the intermediate image plane, v takes on the form

$$v = 2\pi \frac{\text{NA}}{M\lambda} r_i \quad (4)$$

where NA is the numerical aperture and M the magnification of the objective lens, λ the wavelength of light, and r_i the radial distances measured in the intermediate image plane. If we express r_i as a distance r_o in the object plane, with $r_i = M r_o$, we obtain the more familiar relationship:

$$v = 2\pi \frac{\text{NA}}{\lambda} r_o \quad (5)$$

The central bright disk of the diffraction image is known as the Airy disk, and its radius (the radius from the central peak to the first minimum of the diffraction image) in object plane units is given by

$$r_{\text{Airy}} = 0.61 \frac{\lambda}{\text{NA}} \quad (6)$$

When there exist two equally bright, self-luminous points of light separated by a small distance d in object space, that is the specimen plane, their diffraction images lie side by side in the image plane. The sum of the two diffraction images, assuming the two points of light were mutually incoherent, appears as in Fig. 15a. As d becomes smaller so that the first minimum of one diffraction image overlaps with the central maximum of the neighboring diffraction image ($d = r_{\text{Airy}}$, Fig. 15b),²⁶ their sum (measured along the axis joining the two maxima) still contains a dip of 26.5 percent of the peak intensities that signals the twoness of the source points (the Rayleigh criterion). Once d becomes less than this distance, the two diffraction images rapidly pass a stage where instead of a small dip, their sum shows a flat peak (the Sparrow criterion) at $d = 0.78 r_{\text{Airy}}$, and thereafter the sum of the diffraction images appears essentially indistinguishable from one arising from a single point source instead of two. In other words, we can no longer resolve the image of the two points once they are closer than the Rayleigh criterion, and we lose all cues of the twoness at spacings below

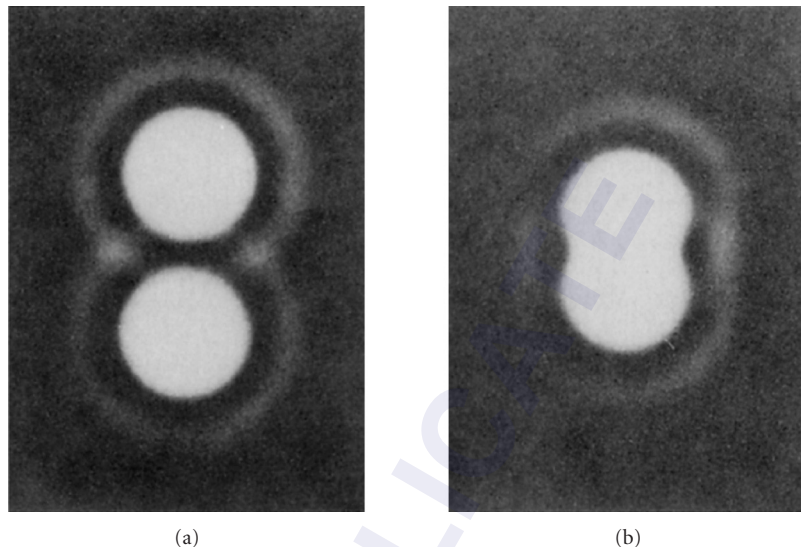


FIGURE 15 Overlapping Airy patterns: (a) clearly resolved and (b) center of Airy patterns separated by $d = r_{\text{Airy}}$, Rayleigh criterion.²⁶

the Sparrow criterion. Since the diameter of the Airy diffraction image is governed by NA_{obj} and the wavelength of the image-forming light λ , this resolution limit normally cannot be exceeded (for exceptions, see the section “Beyond the Diffraction Limit” later in this chapter).

The consideration given here for two-point sources of light applies equally well to two absorbing dots, assuming that they were illuminated incoherently. (Note, however, that it may, in fact, be difficult or impossible to illuminate the two dots totally incoherently since their spacing may approach the diameter of the diffraction image of the illuminating wave. For the influence of the condenser NA on resolution in transillumination, refer to the section on “Transillumination” earlier in this chapter. Also, the contrast of the diffraction images of the individual absorbing dots diminishes rapidly as their diameters are decreased, since the geometrical size of such small dots would occupy a decreasing fraction of the diameter of their diffraction images. For further detail see.)²⁷

The image of an infinitely small point or line thus acquires a diameter equal to that of the Airy disk when the total magnification of the image becomes sufficiently large so that we can actually perceive the diameter of the Airy disk. In classical microscopy, such a large magnification was deemed useless and defined as empty magnification. The situation is, however, quite different when one is visualizing objects smaller than the limit of resolution with video microscopy. The location of the Airy disk can, in fact, be established with very high precision. Distances between lines that are clearly isolated from each other can, therefore be measured to a precision much greater than the resolution limit of the microscope. Also, minute movements of nanometer or even Ångstrom steps have been measured with video-enhanced light microscopy using the center of gravity of the highly magnified diffraction image of marker particles (see “Beyond the Diffraction Limit” later in this chapter).

Three-Dimensional Diffraction Pattern, Axial Resolution, Depth of Focus, Depth of Field The two-dimensional Airy pattern that is formed in the image plane of a **point object** is, in fact, a cross section of a three-dimensional pattern that extends along the optical axis of the microscope. As one focuses an objective lens for short distances above and below exact focus, the brightness of the central spot periodically oscillates between bright and dark as its absolute intensity also diminishes. Simultaneously, the diameters of the outer rings expand, both events taking place symmetrically above and below the plane of focus in an aberration-free system (Fig. 16).²⁸

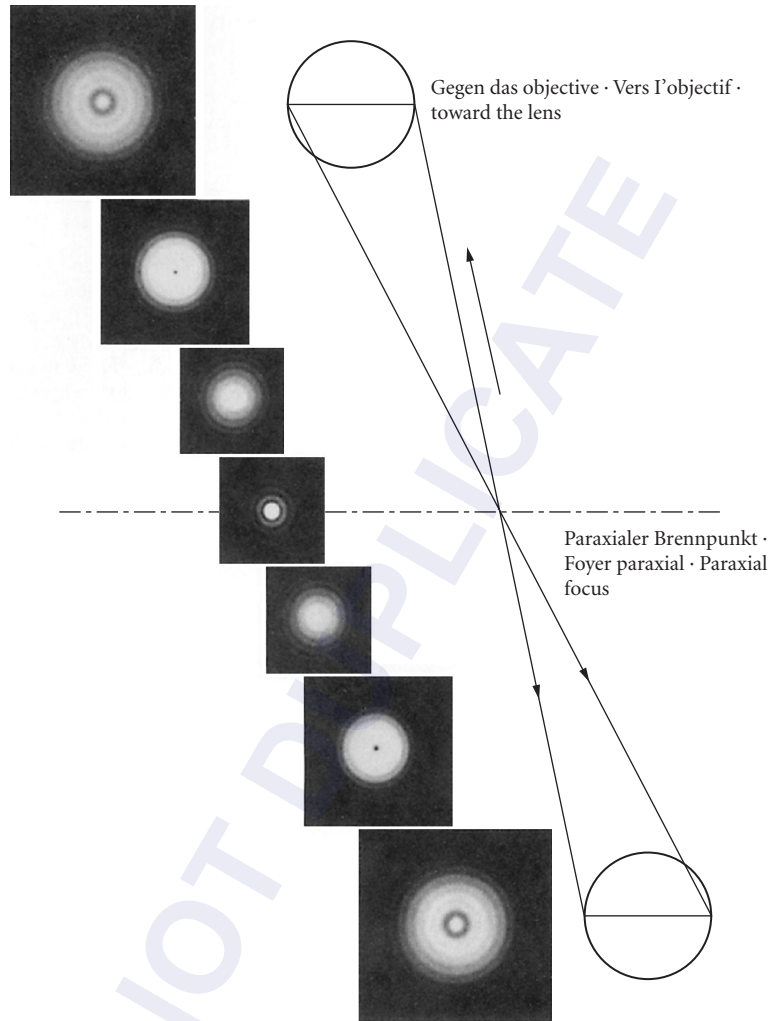


FIGURE 16 The evolution of the diffraction image of a circular aperture with differing planes of focus in an aberration-free system.²⁸

Figure 17 shows an isophote (lines of equal brightness) of the longitudinal section of this three-dimensional diffraction image. The relationship between v and the lateral distance r_i is given by Eq. (4). The axial distance z_p , oriented perpendicular to the image plane, is related to u by

$$u = 2\pi \frac{NA^2}{M^2 \lambda} z_i \quad (7)$$

In the graph we recognize at $v = 1.22\pi$ (and $u = 0$, focal plane) the first minimum of the Airy pattern which we discussed in the preceding section. The intensity distribution along u perpendicular to the focal plane has its first minima at $u = \pm 4\pi$ and $v = 0(\pm z_i$ in Fig. 17a). To find the actual extent of the three-dimensional diffraction pattern near the intermediate plane of the microscope, we express the dimensionless variables v and u of Fig. 17c as actual distances in image space.

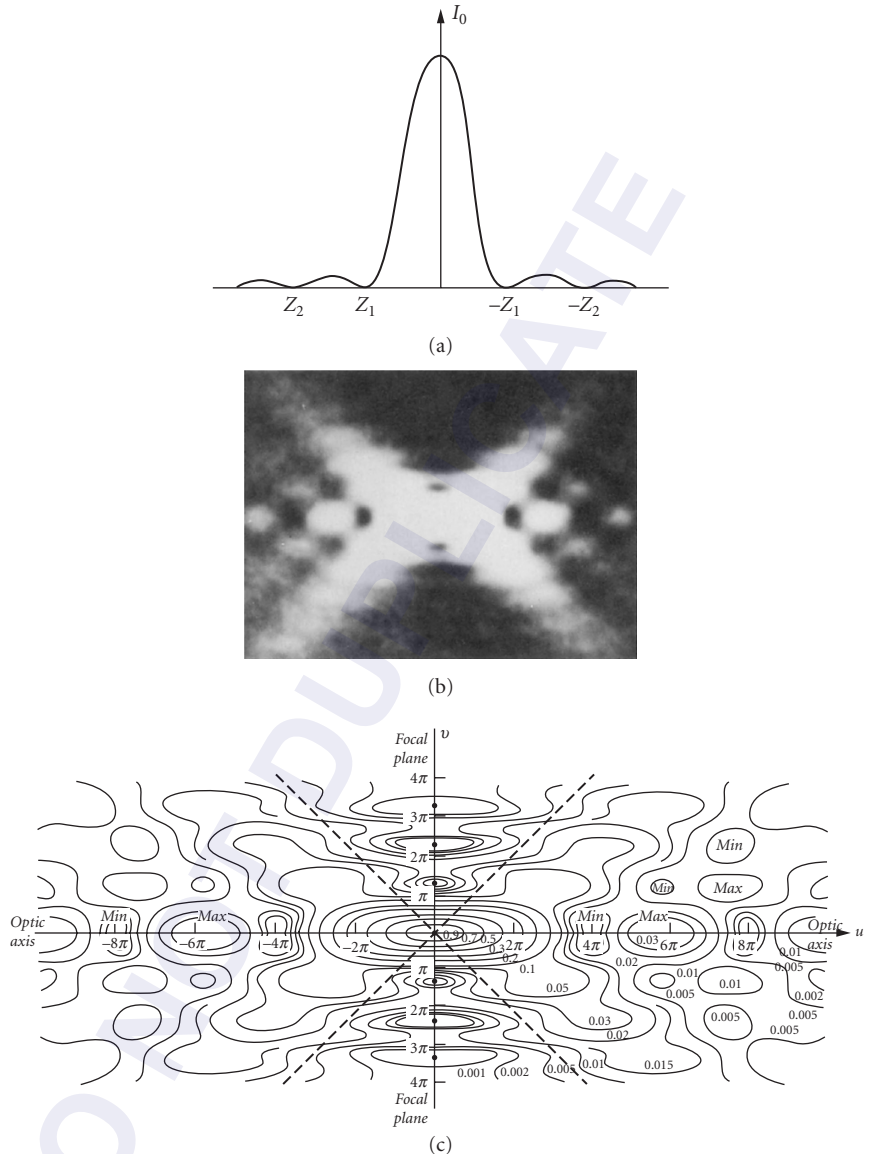


FIGURE 17 (a) Axial intensity distribution of irradiance near focal point;²³ (b) meridional section through diffraction pattern near focal point of a point source of light focused by lens with a uniform circular aperture;²³ and (c) contour plot (isophote) of the same cross section as in (b).^{17,23,29} The three-dimensional diffraction pattern is obtained by rotating the meridional section around the optic axis. The three-dimensional diffraction pattern is also called the intensity point spread function.

The first minimum ($u = 4\pi$) is at a distance $z_1 = (2M^2\lambda)/\text{NA}^2$. To transfer distance z_i in image space to distance z_o in object space, we use the relationship $z_i = z_o M^2/n$. (Note that for small axial distances, to a close approximation, the axial magnification is the square of the lateral magnification M divided by the refractive index n of the object medium.) The distance from the center of the three-dimensional diffraction pattern to the first axial minimum in object space is then given by:

$$z_{\min} = 2 \frac{\lambda n}{\text{NA}^2} \quad (8)$$

z_{\min} corresponds to the distance by which we have to raise the microscope objective in order to change from exact focus of a small pinhole to the first intensity minimum in the center of the observed diffraction pattern (see Fig. 16).

In correspondence to the lateral resolution limit, which is taken as the Airy disk radius r_{Airy} [Eq. (6)], we can use z_{\min} as a measure of the limit of *axial resolution* of microscope optics. Note that the ratio of axial to lateral resolution ($z_{\min}/r_{\text{Airy}} = 3.28 n/\text{NA}$) is inversely proportional to the numerical aperture of the objective lens.

The axial resolution of the microscope is closely related to the depth of focus, which is the axial depth on both sides of the *image* plane within which the image remains acceptably sharp (e.g., when a focusing screen at the image plane is displaced axially without moving the object or objective). The *depth of focus* D is usually defined as 1/4 of the axial distance between the first minima above and below focus of the diffraction image of a small pinhole. In the intermediate image plane, this distance is equal to $z_1/2$, with z_1 defined earlier. The depth of focus defined by z_1 is the *diffraction-limited*, or physical, depth of focus.

A second and sometimes dominating contribution to the total depth of focus derives from the lateral resolution of the detector used to capture the image. This geometric depth of focus depends on the detector resolution and the geometric shape of the light cone converging to the image point. If the detector is placed in the intermediate image plane of an objective with magnification M and numerical aperture NA, the geometrical depth of focus D is given by

$$D = \frac{M}{\text{NA}} e \quad (9)$$

with e the smallest distance resolved by the detector (e is measured on the detector's face plate).

The depth in *specimen* space that appears to be in focus within the image, without readjustment of the microscope focus, is the depth of field (unfortunately often also called the depth of focus). To derive expressions for the depth of field, we can apply the same arguments as outlined above for the depth of focus. Instead of moving the image plane in and out of focus, we keep the image plane in the ideal focus position and move the small pinhole in object space. Axial distances in object space, however, are a factor n/M^2 smaller than corresponding distances in image space. Therefore, we apply this factor to the expression for the geometrical depth of focus [Eq. (9)] and add the physical depth of field [derived from Eq. (8)] for the total depth of field d_{tot} :

$$d_{\text{tot}} = \frac{\lambda n}{\text{NA}^2} + \frac{n}{M\text{NA}} e \quad (10)$$

Notice that the diffraction-limited depth of field shrinks inversely proportionally with the square of the NA, while the lateral limit of resolution is reduced inversely proportionally to the first power of the NA. Thus, the axial resolution and thinness of optical sections that can be attained are affected by the system NA much more so than is the lateral resolution of the microscope.

These values for the depth of field, and the distribution of intensities in the three-dimensional diffraction pattern, are calculated for incoherently illuminated (or emitting) point sources (i.e., $\text{NA}_{\text{cond}} \geq \text{NA}_{\text{obj}}$). In general, the depth of field increases, up to a factor of 2, as the coherence of illumination increases (i.e., as $\text{NA}_{\text{cond}} \rightarrow 0$). However, the three-dimensional point spread function with partially coherent illumination can depart in complex ways from that so far discussed when

the aperture function is not uniform. In a number of phase-based, contrast-generating modes of microscopy, the depth of field may turn out to be unexpectedly shallower than that predicted from Eq. (9) and may yield extremely thin optical sections.³⁰

Beyond the Diffraction Limit In recent years the microscope's limit of resolution, as stated in Eq. (1), has been exceeded by different means, relying either on optical, photophysical, photochemical, or computational methods, or a combination thereof. Here we briefly refer to some of the schemes that rely on photonic properties of the specimen, while later in this chapter we will touch on schemes that rely on far-field optical methods, such as structured illumination and confocal microscopy.

Driven by the success of fluorescence microscopy in biomedical research and the need for higher resolution to understand the molecular machinery of the living cell, several methods were devised that exploit the photophysical and photochemical nature of fluorescent molecules. Most of these "super-resolving" methods take advantage of the fact that the position of a single fluorescent molecule (or point of light) can be determined to a much higher precision than the optical resolution of an imaging system. While the resolution of a traditional microscope, as described by Eq. (1), typically does not exceed 200 nm, the same microscope can be used to determine the position of a single fluorophore to 20 nm or better, depending on the number of photons captured and the mobility of the fluorescent molecule.^{31,32} Here we briefly describe those methods that have become prominent and are recognized by their acronyms. For a more detailed discussion we refer to a number of excellent reviews^{33–36} and to the original publications cited below.

Fluorescence imaging with one nanometer accuracy (FIONA) was introduced to measure the detailed stepping motion of a molecular motor (myosin V) along an immobilized track (filamentous actin).³⁷ The detailed, hand-over-hand motion was determined by measuring the location of a single fluorophore, attached to the motor-protein, with a spatial resolution of 1.5 nm and a temporal resolution of 0.5 s. The challenge here included the recording of a sufficient number of photons, within the 0.5 s time window, to localize a single fluorophore that also needed to be photostable enough to allow its observation over several minutes.

Photo-activated localization microscopy (PALM) was introduced to localize immobilized fluorophores at nanometer spatial resolution.³⁸ To this end, fluorophores are used that have to be photoactivated to become fluorescent. A low dose of typically short wavelength light activates a small, random subset of fluorophores that are spaced far enough for their point spread functions to not overlap. The locations of activated fluorophores are measured at nanometer precision and during the measurement process fluorophores become irreversibly bleached. The cycle of low-dose activation and subsequent position measurements is repeated many times and the aggregate position information from all cycles is assembled into a single, super-resolution image.

Stochastic optical reconstruction microscopy (STORM) uses similar principles as PALM but exploits photo-switchable fluorophores that can be turned on and off by exposing them to light pulses of differing wavelengths.³⁹

Single molecule high-resolution colocalization (SHREC) takes advantage of separating the fluorescence of two or more single fluorophores by their spectral characteristics.⁴⁰ By using chromatically differing fluorescent molecules as probes, the probes can approach each other closer than the Rayleigh limit and still be distinguished. The technique is typically used to measure intramolecular distances of 10 nm or more in doubly labeled macromolecules or molecular complexes.

Fluorescence resonance energy transfer or Förster resonance energy transfer (FRET) refers to a photophysical effect that transfers the excitation energy of a fluorescent donor molecule to a nearby fluorescent acceptor molecule. The appropriately chosen donor and acceptor molecules have to be less than 10 nm apart for the radiationless transfer to be effective. For example, FRET can be used to analyze the conformational change of a protein that brings two molecular subunits closer together or farther apart, resulting in enhanced or reduced acceptor fluorescence, respectively. Hence, FRET is a ratiometric method that allows measurement of the internal distance in the molecular frame rather than in the laboratory frame, which makes it largely immune to instrumental noise and drift. While regular FRET reveals the population distributions of interdye distances, single molecule FRET is used to monitor single molecules for long stretches of time.^{41,42}

Stimulated emission depletion (STED) provides a means of point spread function engineering to improve the optical resolution beyond the diffraction limit. A typical single-point scanning STED microscope uses a regularly focused excitation beam that is superimposed by a doughnut-shaped STED beam that instantly quenches excited molecules at the periphery of the excitation spot, thus confining fluorescence emission to the doughnut zero. Saturated quenching results in a fluorescent spot far below diffraction whose scanning across the sample yields a subdiffraction-resolution image.^{34,43}

All the above methods rely on fluorescence microscopy. A general approach to improve resolution was proposed by Harris⁴⁴ who argued that the diffraction pattern in the Fourier plane can be extrapolated beyond the spatial frequency that is cut off by the NA of the objective lens—in other words, that the limit of resolution can be exceeded by computational extrapolation of the diffraction orders as long as the specimen is illuminated in a narrowly limited field.

The field of illumination can be reduced beyond that defined by diffraction by placing the minute exit aperture of a tapered light guide or a minute pinhole closely adjacent to the specimen. By scanning such an aperture relative to the specimen, one obtains a proximity-scanned image whose resolution is no longer limited by the diffraction orders captured by the objective lens. Instead, only the size of the scanning pinhole and its proximity to the specimen limit the resolution.⁴⁵

For nonoptical microscopes, for example in scanning tunneling, force, and other proximity-scanning microscopes, resolution down to atomic dimensions can be obtained on images that reflect topological, electronic, ionic, and mechanical properties of the specimen surface.⁴⁶ In these types of proximity-scanning microscopes, a fine-tipped probe, mounted on a piezoelectric transducer that provides finely controlled x , y , and z displacements of the probe, interacts with specific properties of the specimen surface (alternatively, the probe may be fixed and the sample mounted to the transducer). The resulting interaction signal is detected and fed back to the z -axis transducer, which generally induces the probe tip to rise and fall with the surface contour (that reflects the particular electrical or mechanical property of the surface) as the probe is scanned in a raster fashion along the x and y directions over an area several tens of angstroms to several tens of micrometers wide. A highly magnified contour image of the atomic or molecular lattices is generated on a monitor that displays the z signal as a function of the x , y position.

28.4 CONTRAST AND IMAGING MODES

In microscopy, the generation of adequate and meaningful contrast is as important as providing the needed resolution. Many specimens are practically transparent and differ from their surroundings only by slight changes in refractive index, absorbance, reflectance, or optical anisotropy such as birefringence and dichroism. Most objects that are black or show clear color when reasonably thick become transparent or colorless when their thickness is reduced to a few tenths of a micrometer (since absorption varies exponentially with thickness). Additionally, in microscopy the specimen is often illuminated using a highly convergent beam to maximize resolution, thus reducing shadows and other contrast cues that aid detection of objects in macroscopic imaging. Furthermore, contrast is reduced at high spatial frequency because of an inherent fall-off of the contrast transfer function.

Many modes of contrast generation are used in microscopy partly to overcome these limitations and partly to measure, or detect, selected optical characteristics of the specimen. Thus, in addition to simply raising contrast to make an object visible, the introduction of contrast that reflects a specific physical or chemical characteristic of the specimen may impart particularly important information.

As a quantitative measure of expected contrast generation as functions of spatial frequencies, the modulation transfer functions (MTFs, of sinusoidal gratings) can be calculated theoretically for various contrast-generating modes assuming ideal lenses (Figs. 18 and 19),^{47,48} or on the basis of measured point or line spread functions.⁴⁹ Alternatively, the contrast transfer function (CTF, of square wave gratings) can be measured directly using test targets made by electron lithography (Fig. 20).⁵⁰

The rapid advance of electronic imaging and digital image processing in recent years made the quantitative evaluation of microscope images much more practical. Many computerized image-processing platforms provide standard functions to characterize the morphology and geometric

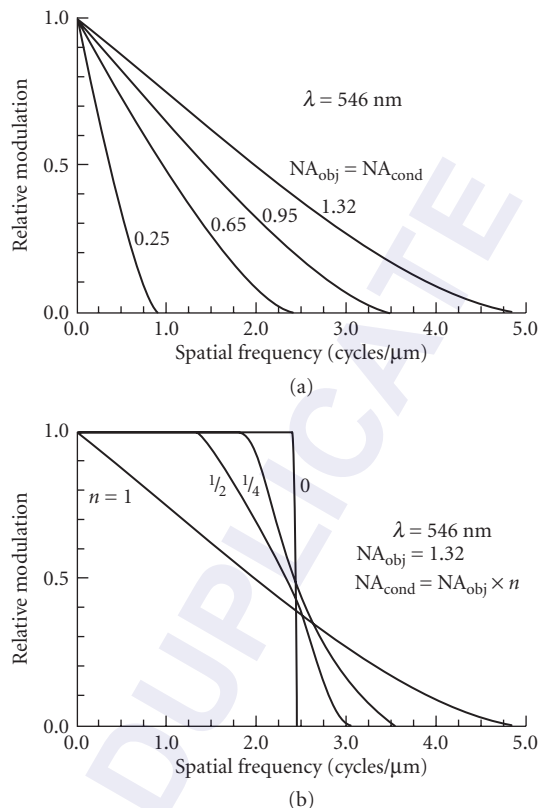


FIGURE 18 Modulation transfer function (MTF) curves for microscope lenses, calculated for periodic specimens in focus: (a) each curve represents a different numerical aperture (NA), which is the same for the objective and condenser lens in these curves. (b) These MTF curves all represent an objective lens of 1.32 NA, but with different condenser NAs; the conditions are otherwise the same as in (a). (Courtesy of Dr. G. W. Ellis.)²

relationship between image features. In addition, specialized systems that provide computer control of microscope components and settings in conjunction with quantitative image analysis provide advanced imaging modalities and new contrast modes that can no longer be viewed through the ocular, but can only be displayed on a computer screen. These hybrid contrast modes usually build on a traditional imaging mode and extend it through exact control and quantitation of image content. Therefore, in the following section we will present traditional imaging modes and give brief descriptions of related hybrid contrast modes.

Bright Field

Whether on an upright or inverted microscope, bright field is the prototypic illumination mode in microscopy (Fig. 4). In transmission bright-field illumination, image contrast commonly arises from absorption by stained objects, pigments, metal particles, etc., that possess exceptionally high extinction

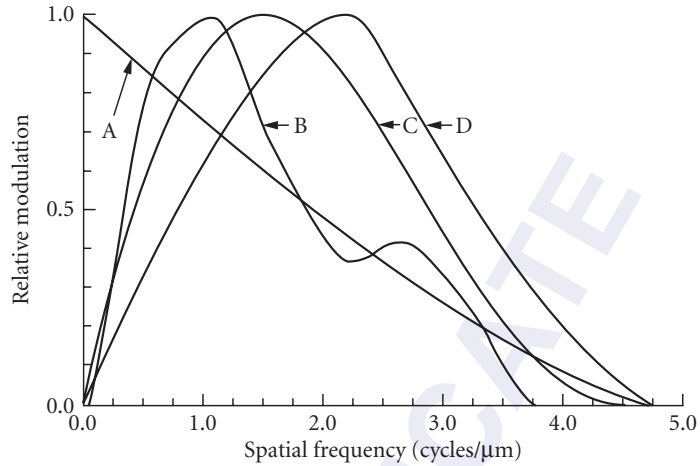


FIGURE 19 Modulation transfer function curves calculated for different modes of microscope contrast generation. A = bright field, B = phase contrast, C = differential interference contrast, and D = single-sideband edge enhancement. The curves are plotted with their peak modulation normalized to 1.0. (Courtesy of Dr. G. W. Ellis.)²

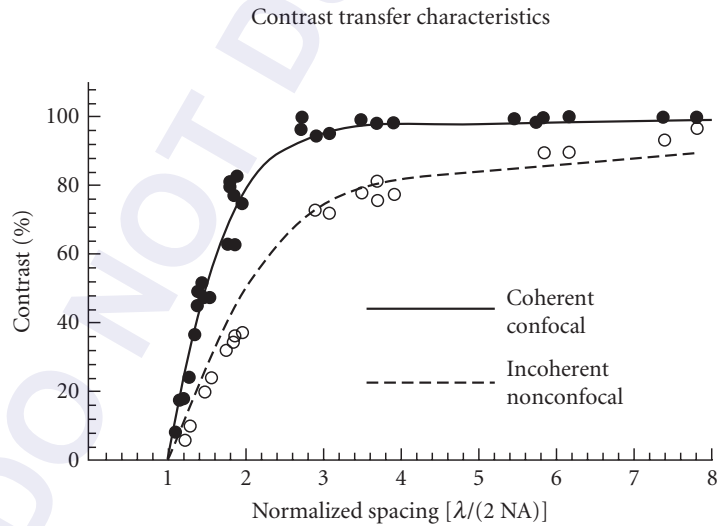


FIGURE 20 Measured contrast transfer values plotted as a function of spatial period in Airy disk diameter units, to normalize the values measured with different lenses and wavelengths. Data points were obtained with a laser spot scan microscope operating in the confocal reflection mode (*solid points*) and the nonconfocal transmission mode (*circles*). Curves are calculated contrast transfer values for the coherent confocal and the incoherent nonconfocal imaging mode.⁵⁰

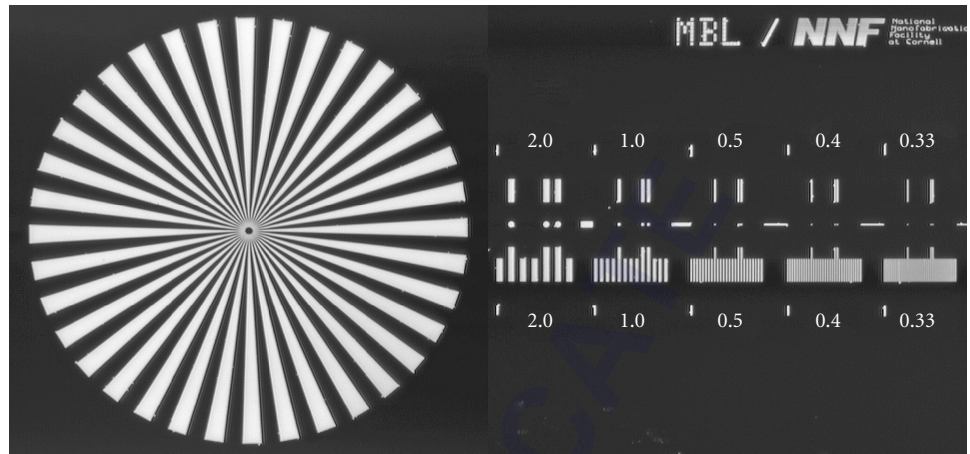


FIGURE 21 Siemens star, line and dot patterns that are part of the MBL/NNF test target imaged in bright field using transmitted light and a 60 \times /1.4 NA Plan Apo oil immersion objective lens (Nikon Inc.) and matching condenser. The dark background is due to the low transmittance of the 50-nm-thick aluminum film. Bright features were etched into the film using electron lithography. Numbers above and below bar gratings show period in microns. The Siemens star consists of 36 wedge pairs, with an outer diameter of 75 μm . The period near the outer edge is 6.5 μm , decreasing continuously toward the center. The smallest period is 0.1 μm near the inner black disk, which has a diameter of 1.2 μm .

coefficients (Fig. 21). Transparent objects only generate very weak contrast based on Becke lines introduced by refraction at object boundaries that are slightly out of focus. (The dark Becke line, which is used for immersion determination of refractive index of particles,⁵¹ surrounds, or lies just inside, a boundary with a sharp gradient of refractive index when the boundary is slightly above or below focus. The Becke line disappears altogether when a thin boundary is exactly in focus.)

To gain additional contrast, especially in bright-field microscopy, the condenser NA is commonly reduced by closing down its iris diaphragm. This practice results in loss of resolution and superimposition of diffraction rings, Becke lines, and other undesirable optical effects originating from regions of the specimen that are not exactly in focus. The various modes of optical contrast enhancement discussed in following sections obviate this limitation and provide images with improved lateral and axial resolution as well as improved contrast.

Before the advent of phase-contrast and differential interference-contrast (DIC) microscopy, oblique illumination (that can be attained by off-centering a partially closed condenser iris diaphragm) was used to generate contrast of transparent objects. While this particular approach suffered from the problems listed in the previous paragraph, combination of oblique illumination at large condenser NA with video contrast enhancement proves to be an effective method for generating DIC-like thin optical sections.⁵²

Recently, the optical phenomenon that leads to the formation of the Becke line has been explored more thoroughly, from a theoretical and an experimental point of view.^{53–56} The goal is to retrieve phase information from images of objects that affect the phase of transmitted or reflected light, but not necessarily its amplitude. Usually, phase information is gained from specially designed setups that enhance interference effects between light waves that have different optical paths through the specimen. The following sections on phase-contrast, polarized light, and interference microscopy give examples of these specialized imaging modes. Streibl,⁵³ on the other hand, proposed to use a regular bright field microscope and the phenomenon of the Becke line to retrieve phase information of weakly scattering objects. He presented a theoretical framework based on the intensity transport equation and demonstrated the enhancement of phase objects based on images that were recorded at slightly different focus positions. Nugent and collaborators^{55–57} have refined the theory and developed a practical implementation called *quantitative phase microscopy*.

In reflection bright-field microscopy, the image is formed by the reflected or backscattered light of the specimen, which is illuminated through the objective (see the section “Epi-illumination”). Reflection contrast is used primarily for opaque and thick samples, especially for metals and semiconductors. Reflection contrast is also finding increasing applications in autoradiography and in correlative light and electron microscopy for detecting the distribution of colloidal gold particles that are conjugated to antibodies and other selective indicators.

Total frustrated reflection microscopy⁵⁸ generates contrast due to objects that are present in a low-refractive-index medium located within the evanescent wave that extends over a distance only a fraction of a wavelength from the microscope coverslip surface. Regions of the specimen whose refractive index differs from its milieu produce interference fringes whose contrast sensitively reflects the refractive index difference and distance from the coverslip surface.

Dark Field

In dark field microscopy the illuminating beam is prevented from entering the image-forming ray paths. The background of the field is dark, and only light scattered by optical discontinuities in the specimen is designed to appear in the image as bright lines or dots. Thus, contrast can become extremely high, and diffraction images can be detected as bright points or lines even when the diameter of the scattering object becomes vanishingly small compared to the microscope’s limit of resolution.^{8,27,30,59}

For small objects that are not obscured by other light-scattering particles (a condition rather difficult to achieve) and are free in a fluid substrate, Brownian motion of the object and the time constant and sensitivity of the detector, rather than the object’s absolute size, are more likely to set a lower limit to the size of the object that can be clearly visualized with dark field microscopy.

Phase-Contrast and Other Aperture-Modifying Contrast Modes

Microscopic objects, distinguished from their surround only by a difference of refractive index, lose their Becke line and disappear altogether when brought exactly into focus. Nevertheless, light diffracted by the small object still suffers a $\lambda/4$ phase shift relative to the undeviated background wave by the very act of being scattered (by a nonabsorbing object; the phase shift upon scattering by an absorbing object is $\lambda/2$).⁶⁰ As shown in Fig. 22, light s scattered by the small object and the undeviated light u , both originating from a common small point A of the condenser aperture, traverse different regions of the objective lens aperture. At the objective aperture, the undeviated light traverses only point B that is conjugate to A , while the scattered light passes those regions of the aperture defined by the spatial periods of the object.

Since light waves s and u arise from the same points in object space but traverse regions that are spatially separated in the objective aperture plane, a *phase plate* introduced in that plane can be used to modify the relative phase and amplitudes of those two waves. The phase plate is configured to subtract (or add) a $\lambda/4$ phase to u relative to s so as to introduce a $\lambda/2$ (or zero) phase difference between the two and, in addition, to reduce the amplitude of the u wave so that it approximates that of the s wave. Thus, when the two waves come to focus together in the image plane, they interfere destructively or constructively to produce a darker or brighter in-focus image of the small, transparent object against a dark gray background (positive and negative phase contrast).

As generally implemented, an annulus replaces the pinhole in the condenser aperture, and a complementary phase ring in the objective aperture plane or its conjugates (covering a somewhat larger area than the undisturbed image of the annulus in order to handle the u waves displaced by out-of-focus irregularities in the specimen) replaces the simple phase disk. Figure 23 shows an example of a phase object that was imaged using phase-contrast optics as described above. The object is a Siemens star that was etched into a thin layer of silica and imaged using a Olympus $100\times/1.3\text{NA}$ Plan Apo objective and condenser with complementary phase rings.

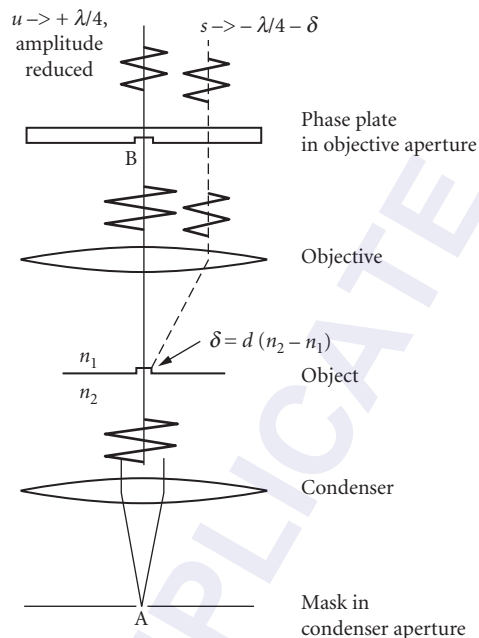


FIGURE 22 Optical principle of phase-contrast microscopy illustrating the phase relationships between waves of the light s scattered by the specimen and the undeviated light u (see text).

In the Polanret system, the phase retardation and effective absorbance of the phase ring can be modified by use of polarization optical components so that the optical path difference of a moderately small object can be measured by seeking the darkest setting of the object.^{61,62} Similarly, the Polanret system can be used to accentuate color or low contrast due to slight absorption by the object.

Several modes of microscopy, including phase contrast, take advantage of the facts that (1) the front condenser and back objective lens apertures are conjugate planes, (2) the illuminating beam arising out of each point of the condenser aperture is variously deviated by the specimen structure according to its spatial frequency, and (3) the back objective aperture is the Fourier plane of the specimen plane.

In Hoffman modulation contrast microscopy, the condenser aperture contains a slit mask with the slit placed toward the edge of the aperture. The objective aperture holds a second, complementary mask, called a modulator, which consists of two parts (Fig. 24).⁶³ The dark part covers the smaller sector to one side of the projected slit and the gray part covers the slit area. The objective mask thus attenuates the zero-order light undeviated by the specimen and removes the light diffracted by the specimen to one side of the zero-order beam. The light deviated by specimen structure away from the dark sector of the mask passes unchanged, while the light deviated toward the dark sector is blocked. Thus, the image becomes shadow-cast, similar in appearance to DIC that reflects gradients of refractive indices or of optical path differences in the specimen.

Developed by Gordon W. Ellis in 1978⁶⁴ single-sideband edge enhancement microscopy (SSEE) generates directional image contrast of phase objects, with greater modulation transfer than by phase-contrast or DIC microscopy at high spatial frequencies (Fig. 19). SSEE is also capable of generating exceptionally thin optical sections (Fig. 25). (In 1988, Ellis also devised aperture-scanning phase-contrast microscopy, a method which generates full resolution phase-contrast images with virtually no halos;⁶⁵ see Fig. 2-47 in Ref. 2.)

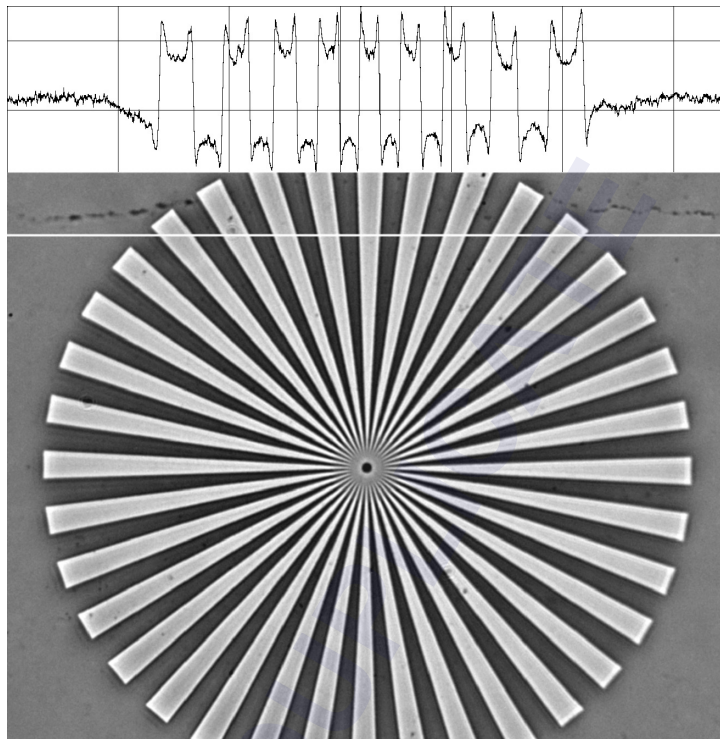


FIGURE 23 Siemens star etched into 90-nm-thick SiO_2 layer and imaged with phase contrast. The dimensions of the star pattern are the same as the one described in Fig. 21. The wedges that were etched away appear bright in this image. Light that has passed through the etched wedges is phase-advanced with respect to light that has passed through the rest of the pattern. (The SiO_2 layer was deposited on a 170- μm -thick coverglass, etched using electron lithography, and mounted on a microscope slide, leaving an air gap between slide and silica layer; 100 \times /1.3NA oil immersion Plan Apo objective.) The intensity profile along a horizontal line near the top illustrates deviations from the step function of the corresponding optical path difference. The accentuation of the edge contrast in this profile is an artifact of the phase-contrast method commonly implemented in form of an illumination ring in the condenser aperture and a complementary phase ring in the objective aperture.

SSEE takes advantage of the fact that illuminated by a condenser whose aperture is half masked, the two side bands (shifted by specimen diffraction to the left and right) are both phase shifted relative to the illuminating light (carrier wave) by $\lambda/4$, but with opposite signs. (As shown by Zernike,⁶⁰ the image contrast of a phase grating viewed with a bright field microscope disappears at exact focus, since the two side bands are in opposite phase.)

In the SSEE microscope, contrast is generated by interference between the attenuated carrier wave and one of the side bands (Fig. 26). Alternately, both side bands may be used with one of the side bands phase shifted by $\lambda/2$ (and appropriately attenuated) relative to the other. Interference between the attenuated carrier wave and the side band generates a high-contrast, high-resolution, in-focus image of the specimen's phase boundaries proportional to their orientation perpendicular to the straight edge of the half mask in the condenser.

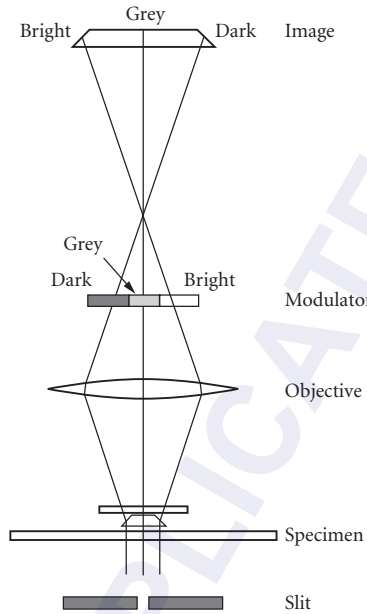


FIGURE 24 Schematic diagram indicating regions of the modulator that modify light from phase gradients in the object to enhance contrast.⁶³

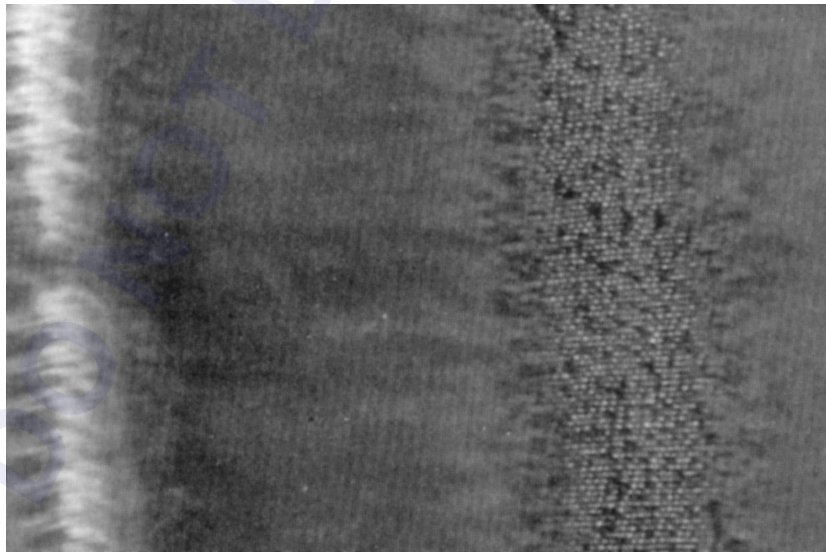


FIGURE 25 Optical section of the silicate shell of a diatom (*Surirella gemma*) observed with SSEE. The tiny pores are in focus over only a highly limited region of the shell due to the highly effective optical sectioning capability of SSEE. (Image copied and cropped from Fig. 2-50 in Ref. 2. Original image courtesy of Dr. Gordon W. Ellis, University of Pennsylvania.)

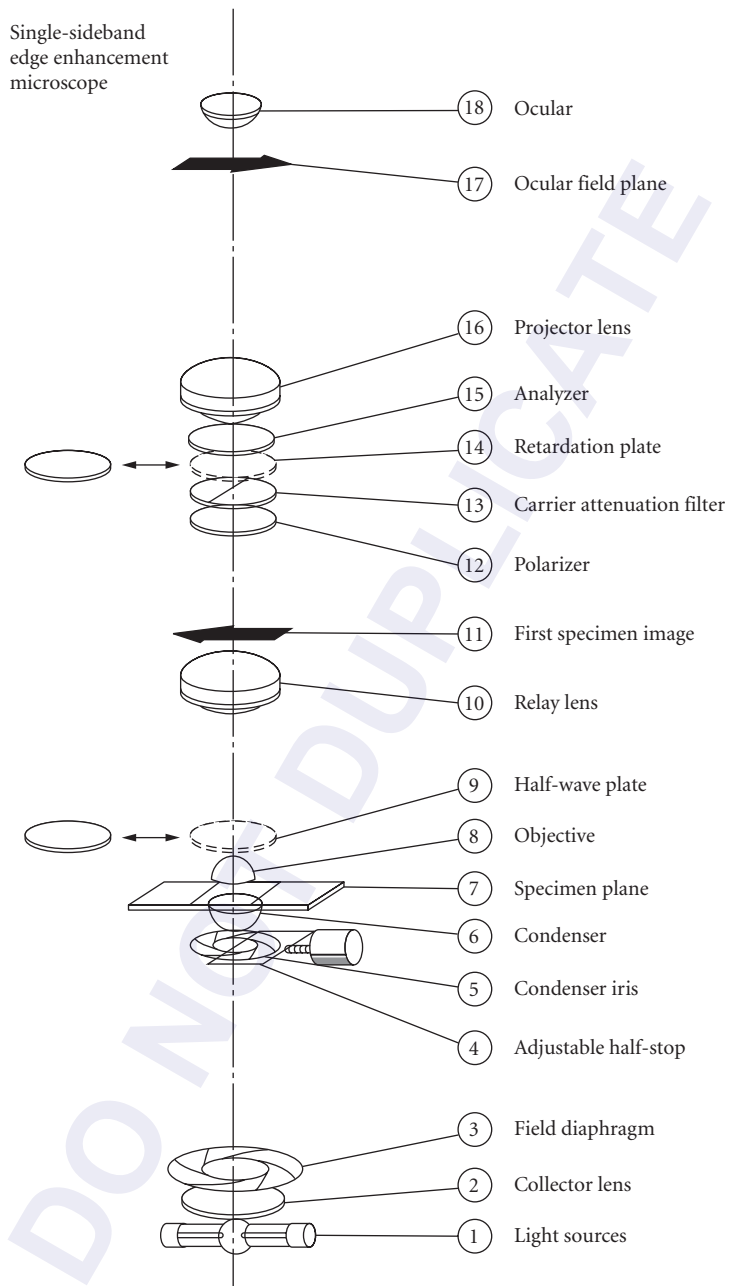


FIGURE 26 Schematic diagram of the edge enhancement single-sideband microscope (SSEE).⁶⁴

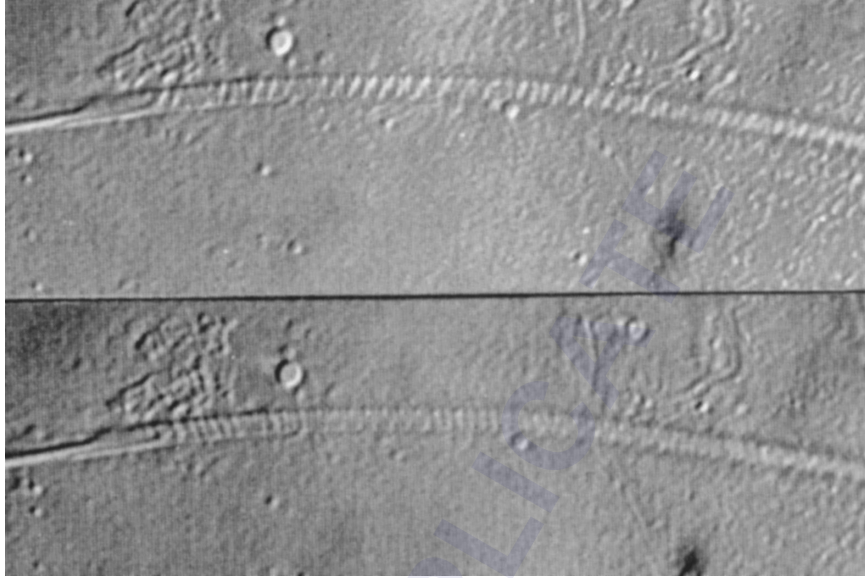


FIGURE 27 Gyres of chromosomes in live sperm head of cave cricket. The images were obtained with SSEE using selected e-vectors as described in the text. Besides rectified polarization microscopy (Fig. 36), few contrast-generating methods besides SSEE have been able to distinctly display these chromosome gyres. (Image copied from Fig. 2-51 in Ref. 2. Original image courtesy of Dr. Gordon W. Ellis, University of Pennsylvania.)

In SSEE, polarizing elements placed *after* the specimen attenuate and phase shift the carrier wave relative to the side bands (Fig. 26). Thus by adjusting the azimuth of the polarizer immediately following the specimen, one can capture exceptionally high resolution images reflecting the birefringence distribution and axes in the specimen (Fig. 27). Also, since the specimen is not sandwiched between crossed polarizers, image contrast in SSEE is not affected by birefringence of the specimen chamber as is the case with polarization and DIC microscopy.

Interference

While all modes of contrast generation in light microscopy in fact depend on interference phenomena, a group of instruments is nevertheless known separately as interference microscopes. These microscopes form part of an interferometer, or contain an interferometer, that allows direct measurements of optical path difference (or generation of contrast) based on interference between the waves passing the specimen and a reference wave. The interferometric and polarization microscopy techniques, which are considered below, generate complementary phase images of the specimen: distribution of refractive index and distribution of refractive index anisotropy, respectively.

Many interference microscopes employ amplitude-dividing beam splitters for setting up the two-beam interference scheme. Instead of amplitude division, division of wavefront can also be used to create both beam paths, especially when using a laser light source. Among the many designs that have been proposed and manufactured, amplitude division interference microscopes can be classified into three major groups: (1) the two-arm type with two separate beam paths, one containing the sample, the other for controlling the reference beam, with separate microscope optics in both arms or microscope optics only in the sample arm; (2) the beam-shearing type in which the reference wave is generated by displacing a beam laterally within the field of a single microscope; and

(3) the dual focus type in which the reference wave is focused to a different level than the specimen plane, again in a single microscope. All schemes can be implemented in transmission or reflection mode.^{66,67}

The image in an interferometric microscope is created by the superposition of a probe and a reference beam. We denote the intensities in the probe and reference beam as I_p and I_r , and their respective phases as ϕ_p and ϕ_r . The intensity that results from superimposing the probe and reference beam can be expressed as described in^{17,26}

$$I = I_p + I_r + 2|\gamma_{pr}|\sqrt{I_p I_r} \cos(\phi_p - \phi_r) \quad (11)$$

where $|\gamma_{pr}|$ is the modulus of the normalized mutual coherence function or the degree of coherence between the probe and reference image. This equation does not include polarization effects and assumes that both interfering beams have the same polarization. For quasi-monochromatic light the optical path difference (OPD) that is associated with the phase angle difference is given by

$$\text{OPD} = \frac{\bar{\lambda}}{2\pi}(\phi_p - \phi_r) \quad (12)$$

where $\bar{\lambda}$ is the center wavelength.

We note that in Eq. (11) I can stand for an array of intensity values representing the pixels of a digital image that was recorded with an appropriate camera attached to an interference microscope.

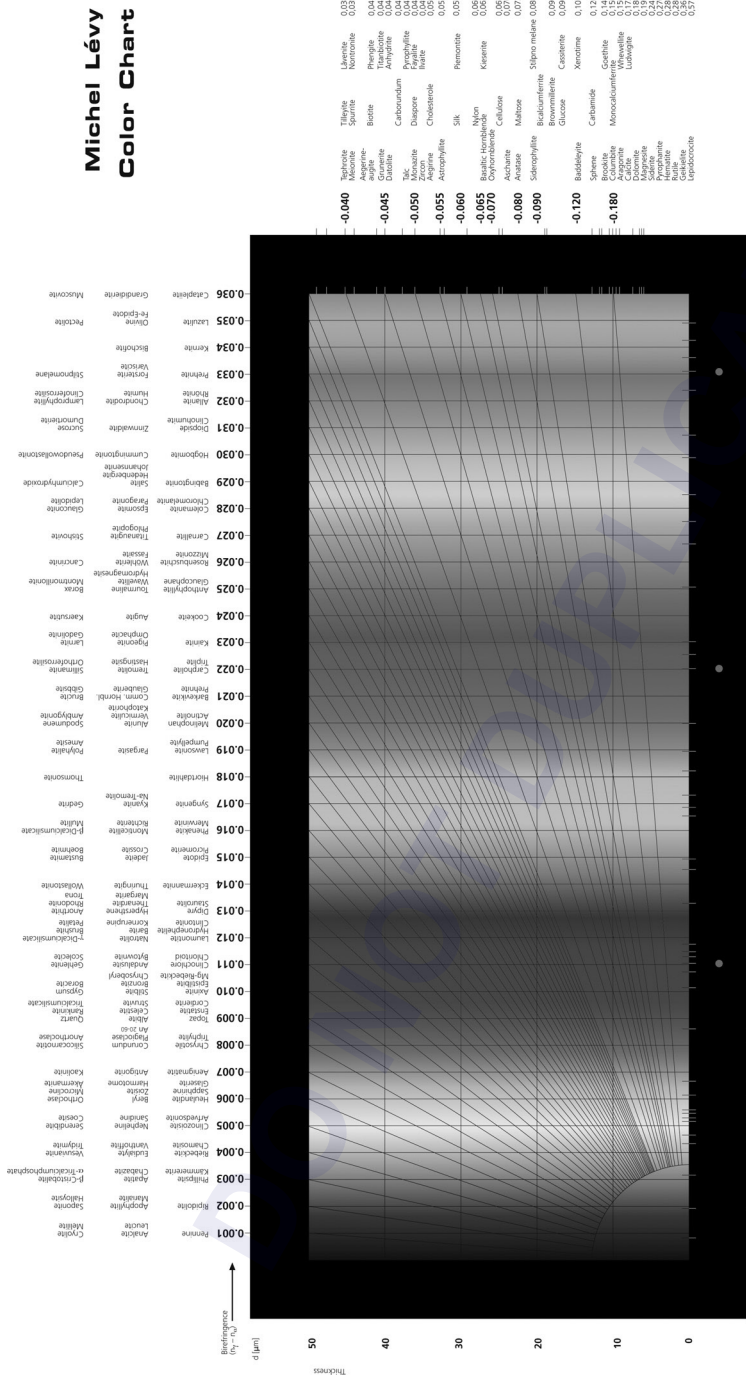
When using white light, each wavelength produces its own interference picture. White light interference pictures are only observed when the optical path difference between the probe and the reference beam is less than a few wavelengths. Let's assume that in a uniform image region the OPD is zero, hence the interference of each wavelength is constructive and the recorded spectrum in that region is white. However, if the OPD is finite, the wavelength that is twice the OPD is suppressed due to destructive interference and therefore that wavelength is missing from the spectrum recorded in the region. When systematically increasing the OPD from 0 to 2000 nm, for example, a characteristic change in spectrum is observed in the region, transitioning from white (OPD = 0), to blue (OPD = 300 nm), to yellow (OPD = 600 nm), to indigo (OPD = 900 nm), to a greenish yellow (OPD = 1600 nm), and bluish grey (OPD = 2000 nm).⁶⁸ As the OPD increases above 1000 nm, colors become less saturated and approach white again for OPDs of several thousand nanometers.

In some interferometric schemes there is an additional achromatic half-wave phase shift, for instance, due to polarization transformation, reflection, and the like. In this case, a zero optical path difference produces destructive interference at all wavelengths and a uniform image region with zero OPD appears black. For small OPDs (< 200 nm) the destructive interference is relaxed for all wavelengths simultaneously and the brightness of the region increases, first with a white spectral composition. With increasing OPD, the region becomes colored due to constructive and destructive interference of specific wavelengths leading to the following color sequence: light yellow (OPD = 300 nm), indigo (OPD = 600 nm), yellow (OPD = 900 nm), grey blue (OPD = 1600 nm), and whitish grey (OPD = 2000 nm). This sequence of interference colors is reproduced in the Michel-Lévy chart (see Fig. 28), which is used to rapidly estimate the OPD based on the observed color of a uniform region. When the OPD increases above 2000 nm, the interference colors turn white and can no longer be used to reliably determine the OPD.

Both interference schemes are implemented. The scheme with destructive interference at 0 OPD is more sensitive (higher signal to noise ratio) for measuring small OPDs, because the background of 0 OPD is black (ideally) and doesn't carry any shot noise, while the white light intensity of constructively interfering beams is subject to shot noise.

When using monochromatic light, the optical path difference between the probe and reference beam can be determined precisely by measuring the intensity in a uniform image region [see Eq. (11)]. However, measurements of OPDs that are larger than the wavelength λ of the monochromatic light result in an apparent OPD that is between 0 and λ . This ambiguity is often referred to as the order of the OPD. The order can be determined by making measurements with two or more wavelengths.

Michel Lévy Color Chart



We make it visible. 4

FIGURE 28 Michel Lévy Color Chart of interference colors. The horizontal sequence of colors is associated with the interference of two beams of white light whose mutual path difference increases from 0 (left, black for destructive interference) to more than 1700 nm (right, pale green). By comparing an experimentally observed color with the colors of this chart, one can estimate the path difference caused, for example, by the partial reflection off a thin dielectric film (e.g., soap bubble) or transmission through a thin birefringent sheet (e.g., mica) sandwiched between crossed polarizers. For birefringent materials, the path difference is the product of the birefringence (indicated along the upper and right edge) and the thickness (left edge) of the material. The diagonal lines assist in estimating one of the quantities (birefringence or thickness) from the observed color; if the other quantity is known. This brightness/color sequence is specific to interference phenomena that result in destructive interference for zero path difference (e.g., birefringent sheet between crossed polarizers). A complementary color sequence applies to interference phenomena that result in constructive interference at zero path difference (e.g., birefringent sheet between parallel polarizers). (See also color insert.) (The chart was generously provided by Rudi Rotterfusser and Becky Hohman of Carl Zeiss Microimaging Inc.)

Further improvement in measuring the OPD can be achieved by controlling the phase of the reference beam. For example, we can measure the probe plus reference image four times, each time changing the phase of the reference image by a quarter wavelength. According to Eqs. (11) and (12) we find

$$\begin{aligned}
 I &= I_p + I_r + 2|\gamma_{pr}|\sqrt{I_p I_r} \cos\left(\frac{2\pi}{\lambda} \text{OPD}\right) \\
 I &= I_p + I_r - 2|\gamma_{pr}|\sqrt{I_p I_r} \sin\left(\frac{2\pi}{\lambda} \text{OPD}\right) \\
 I &= I_p + I_r - 2|\gamma_{pr}|\sqrt{I_p I_r} \cos\left(\frac{2\pi}{\lambda} \text{OPD}\right) \\
 I &= I_p + I_r + 2|\gamma_{pr}|\sqrt{I_p I_r} \sin\left(\frac{2\pi}{\lambda} \text{OPD}\right)
 \end{aligned} \tag{13}$$

Assuming all other factors constant we can compute the OPD based on the four intensity measurements:

$$\text{OPD} = \frac{\lambda}{2\pi} \arctan\left(\frac{I_4 - I_2}{I_1 - I_3}\right) \tag{14}$$

The last expression relates the OPD to a ratio of intensity differences. Hence, the OPD is measured independent of an intensity offset (because only intensity differences are entered) and independent of a gain factor that is common to all four intensity values (because only an intensity ratio is entered).

As noted earlier, I_1, \dots, I_4 can be interpreted as arrays of intensity values representing the pixels of four digital images. In this case, the expression for OPD represents an image arithmetic operation that generates a map of the spatial variations of the measured optical path differences. This or similar image-processing schemes can be implemented using various interference microscope designs after adding appropriate equipment for electronic imaging and phase control.

Mach-Zehnder Interference Microscope The classical two-arm interference microscope with identical optics in both arms is the Mach-Zehnder interference microscope as designed by Horn (Leitz of Wetzlar) in the 1950s (Fig. 29). The intricate and sturdy design earned it the nickname “Rolls Royce of the microscopes,”⁶⁹ including its cost, which was comparable to that of an electron microscope of the time. The microscope, while straightforward in principle, requires close matching of the optics in the two interferometer arms and a mechanical design that provides exceptional precision and stability. Thus, in addition to using matched pairs of objectives and condensers and inserting a blank slide (that is similar to the specimen-containing slide) into the reference arm, one needs to carefully adjust the built-in beam deviators, path equalizers, and wedge components to reduce the difference in optical path length between the two arms to less than the coherence length of the quasi-monochromatic light. (The coherence length of light with a center wavelength $\bar{\lambda}$ and a bandwidth of $\Delta\lambda$ is $\bar{\lambda}^2/\Delta\lambda = 30 \mu\text{m}$ for $\bar{\lambda} = 550 \text{ nm}$, $\Delta\lambda = 10 \text{ nm}$.) While unfortunately no longer manufactured, this type of microscope permits precise interferometric measurements of microscopic objects both in the uniform field mode and the fringe displacement mode, and can even be used to generate holograms.²

Linnik Interference Microscope In 1933, V.P. Linnik proposed a two-arm reflective-type interference microscope with two matching objectives and a single ocular.⁷⁰ The optical scheme, also called the Linnik microwinterferometer, is shown in Fig. 30a. The illumination is split and recombined by the same beam splitter before the microscope objective lens where the beam has low divergence. The probe beam and reference beam then pass through separate but matching objectives and reflect off the specimen and reference mirror, respectively. The objectives can have high NA and short working distance, but require close matching for efficient interference of the probe and reference beams in the common image plane or behind the ocular. Closely matched objectives reduce the influence

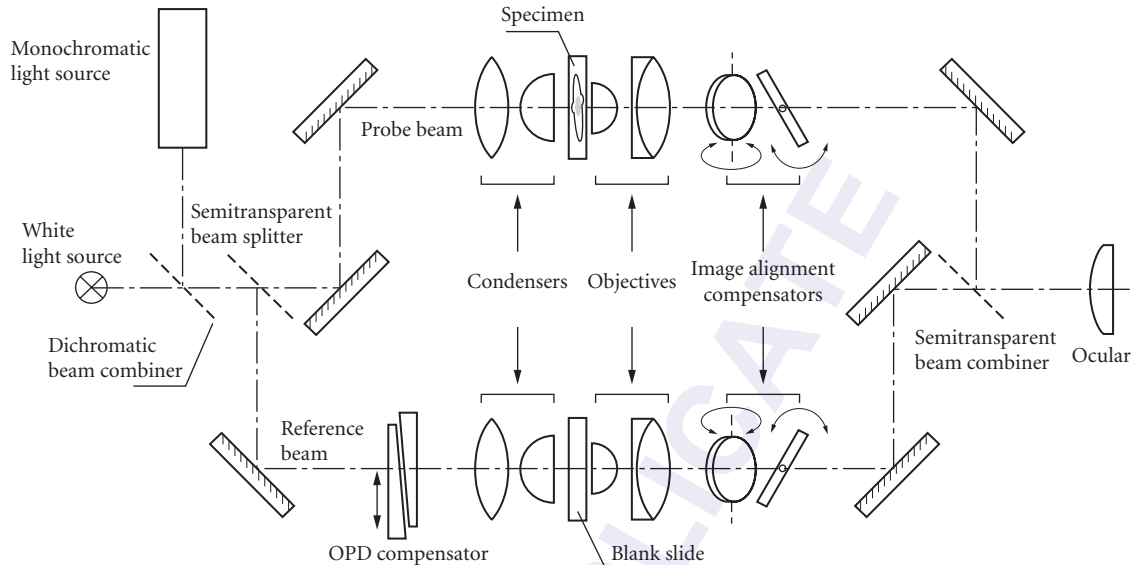


FIGURE 29 Designed by Horn the Mach-Zehnder-type interference microscope with two complete sets of microscope optics, one in each arm of a Mach-Zehnder interferometer.

of chromatic dispersion and other optical aberrations on the interference image. This is essential if a broad-band light source is used, because the dispersion and the optical path length must be closely matched across the entire useful field in each arm. Linnik type interference microscopes are still manufactured by LOMO, Russia.

The original Linnik design can be modified as proposed here by Michael Shribak and shown in Fig. 30*b*. The modification replaces the regular beam splitter with a polarizing one and adds

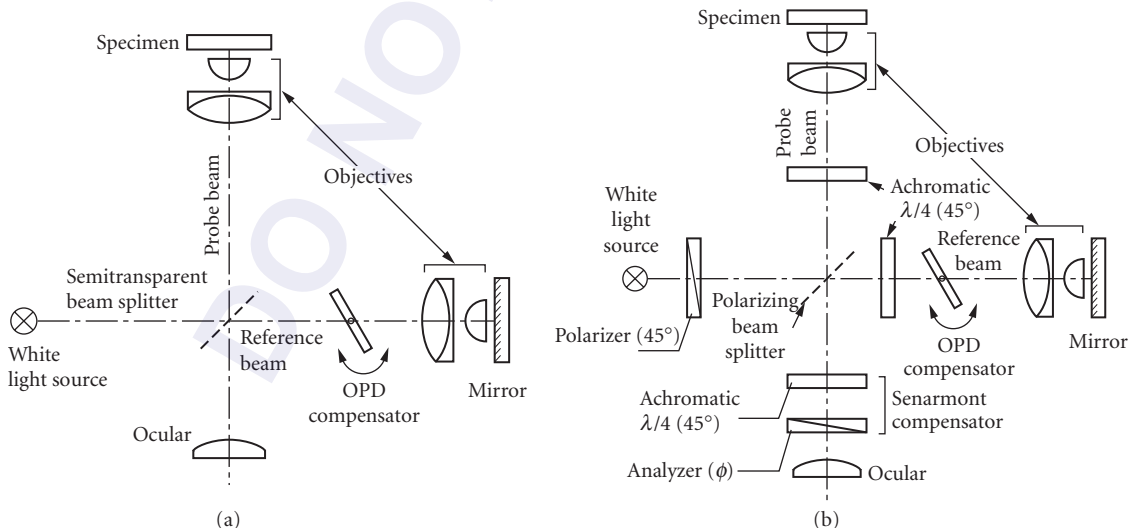


FIGURE 30 The Linnik-type interference microscope with two identical objectives, one in each arm.

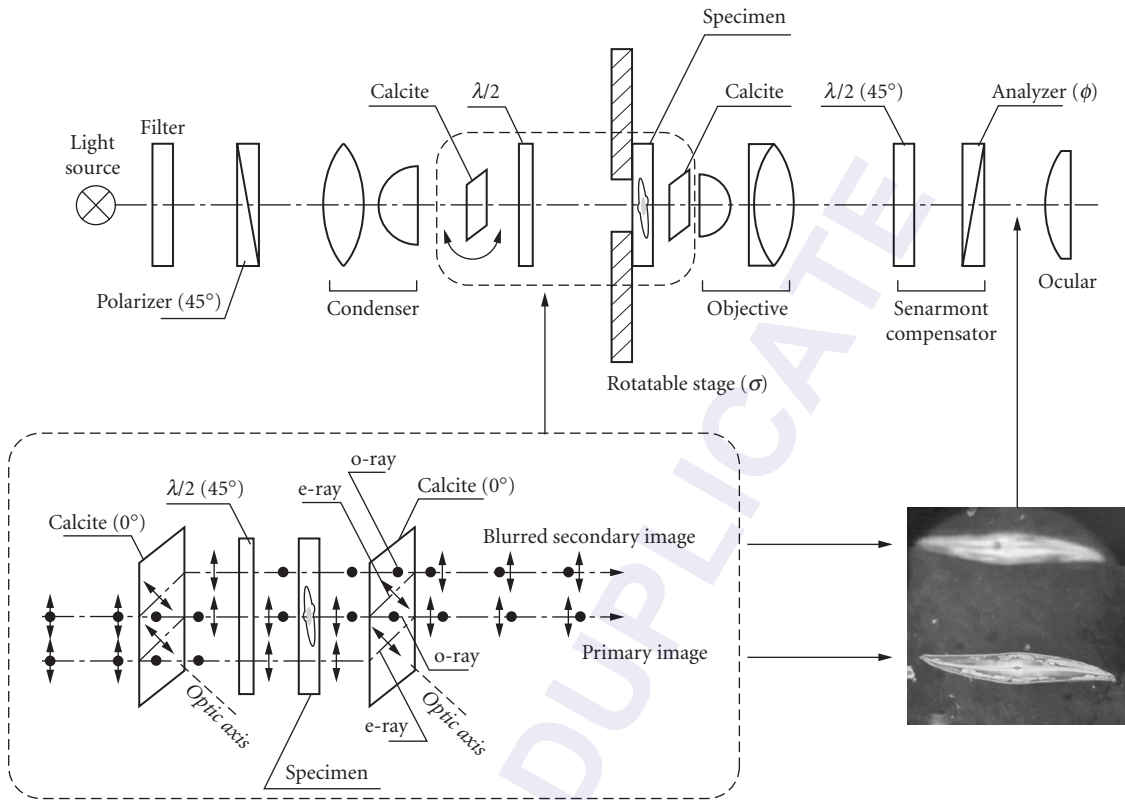


FIGURE 31 The Jamin-Lebedev type interference microscope.

quarter-wave plates to improve sensitivity and to provide a convenient way of measuring the phase. Its enhanced features include a rotatable polarizer, which is used to balance the intensities of the probe and reference beam. The quarter-wave plates following the beam splitter create circularly polarized light, which is reflected by the specimen/reference surface. The reflection induces an inversion of the circularity of the two beams, which causes them to be combined after the beam splitter in the arm with the Senarmont compensator and ocular. The compensator consists of a quarter-wave plate at azimuth 45° and a rotatable analyzer at azimuth ϕ . Image regions with different phase angles can be brought to extinction by rotating the analyzer to different angles. The phase difference Φ between two regions with extinction angles ϕ_1 and ϕ_2 is $\Phi = 2(\phi_2 - \phi_1)$. Other compensation schemes can be used, including liquid crystal devices, and a camera can be added for quantitative imaging.

Jamin-Lebedev Interference Microscope The first interference microscope was constructed by Lebedev in 1930⁷¹⁻⁷³ using a beam-shearing design based on the two-beam polarization interference scheme introduced by Jamin in 1868.⁷⁴ The optical scheme of the Jamin-Lebedev interference microscope is shown in Fig. 31.

In this instrument, a small plane-parallel plate of calcite is cut at 45° to the optic axis and cemented to the front of the objective lens. An identical calcite plate is cemented to the front of the condenser, with an additional half-wave plate facing the specimen. The axes of the two calcite plates are parallel, and at 45° to the axes of the half-wave plate. The specimen under investigation is placed between the half-wave plate and the calcite plate fixed to the objective. The plate fixed to the objective

produces the necessary lateral separation between the probe and reference beams in the intermediate image plane of the microscope. Thanks to the calcite and half-wave plate placed next to the condenser, the path difference of the interfering rays does not vary with the inclination of the rays. This compensation permits quite large openings of the substage condenser diaphragm.

In the beam-shearing Jamin-Lebedev microscope, the probe and reference beam travel a common physical path except along the short distance between the two calcite plates. Because of the common path many design criteria, including mechanical stability and duplication of optical components, can be significantly relaxed in this beam shearing microscope, compared to the dual-arm Mach-Zehnder design. The compromise lies in the lateral shear distance between probe and reference beam, which is limited by the field size and the requirement for telecentric paths for both probe and reference beam. Because both, the probe and reference beam pass through the same specimen slide, the observer has to be wary of ghost images introduced by the reference beam.

The design shown in Fig. 30 was manufactured in the 1960s by Carl Zeiss, Oberkochen, West Germany. The calcite plate next to the objective lens can be slightly rotated to align the shear planes of the two calcite plates. An additional calcite plate introduces a bias in the optical path difference adjusted by a small tilt of the plate. The microscope comes with three pairs of matched condenser and objective lenses, with the objectives designated as $10\times/0.22\text{NA}$, $40\times/0.65\text{NA}$, and $100\times/1.0\text{NA}$ Oil. Their shear distances are 500, 170, and 50 μm , respectively. Optical path differences of less than one wavelength are measured using monochromatic light and a Senarmont compensator. For measuring higher path differences, white light and a Michel-Levy chart (Fig. 28) can be used.

Differential Interference-Contrast Microscope Differential interference-contrast (DIC) microscopy is used extensively in materials research and the life sciences for observing microscopic particles and structures that are associated with refractive index and thickness changes in the specimen. A DIC microscope is a beam-shearing interferometer in which the reference beam is sheared by only a small amount, generally by less than the diameter of the Airy disk that is associated with the imaging optics. The technique produces a shadow-cast image that displays the local gradients of the optical path length. A region of the specimen where the optical path length increases along a reference direction appears brighter (or darker), while a region where the optical path length decreases appears in reverse contrast. As the gradient of the optical path grows steeper, image contrast is increased. Another important feature of the DIC technique is that it produces effective optical sectioning. This is particularly obvious when high numerical aperture (NA) objectives are used together with high NA condenser illumination. The thin optical section is a consequence of the small shear between the interfering beams, which are appreciably separated only in a thin layer around the focal plane.

The DIC technique was invented by F. H. Smith in 1947.^{75,76} He placed between a pair of polarizers one Wollaston prism at the front focal plane of the condenser and a second one in the back focal plane of the objective lens (Fig. 32). The first Wollaston prism splits the linearly polarized input beam into two orthogonally polarized beams that are separated by a small angle ε_1 . The condenser lens converts the angular split in the focal plane into a small spatial shear in the object plane. The objective lens joins the two beams again in the back focal plane where the second Wollaston prism deviates the beams to form two parallel beams again. While parallel, the two beams are orthogonally polarized and therefore cannot interfere. Therefore, a linear analyzer is needed after the second Wollaston prism to create a common polarization and to enable the beams to interfere. The interference generates the typical relief image representing the optical path gradients in the specimen (see inset in Fig. 32).

The small angular split ε_1 and ε_2 in the condenser and objective focal planes are related to the shear amount d in the object plane and the focal lengths of the condenser (f_c) and objective (f_{ob}) lenses by

$$f_c \varepsilon_1 = f_{\text{ob}} \varepsilon_2 = d \quad (15)$$

This optical configuration creates a polarizing shearing interferometer, by which one visualizes optical path gradients of the specimen under investigation.

In conventional medium- to high-NA objective lenses, the back focal plane is located inside the lens system and therefore not available for insertion of a Wollaston prism. If the Wollaston prism

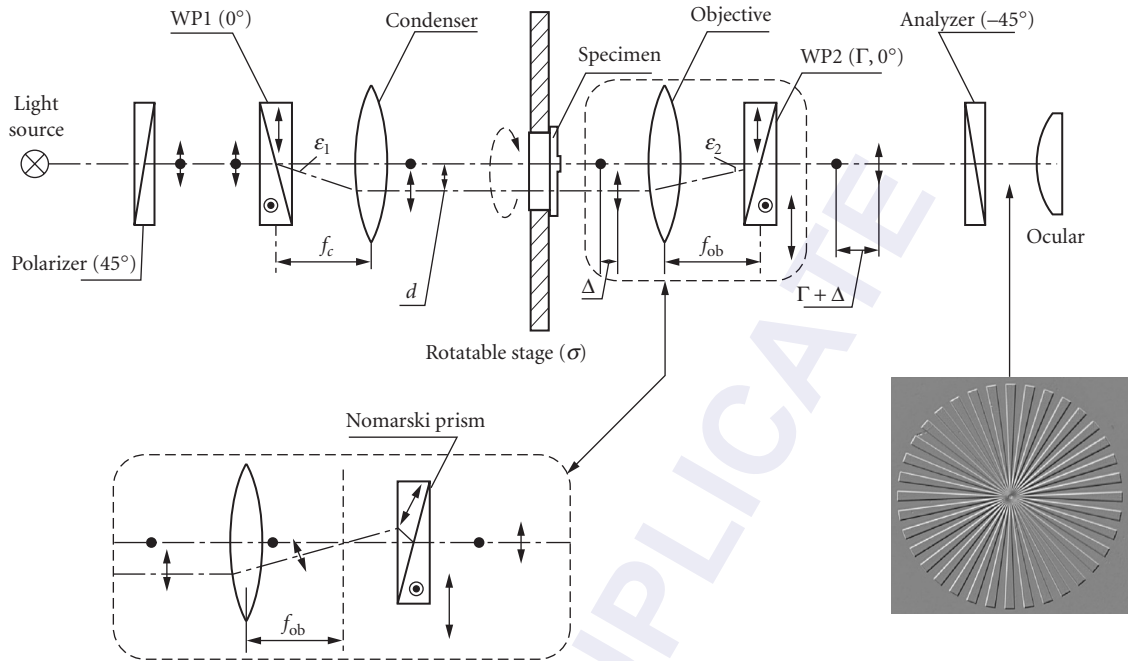


FIGURE 32 DIC microscope setup: polarizer at 45° azimuth; WP1: first Wollaston prism at 0° azimuth; ε_1 : splitting angle; f_c : condenser lens focal distance; d : shear amount; Δ : optical path difference introduced by specimen under investigation; σ : azimuth of rotatable stage; f_{ob} : objective lens focal distance; WP2: second Wollaston prism at 0° azimuth (the second prism introduces bias Γ); ε_2 : splitting angle; analyzer at -45° azimuth; Wollaston prism can be replaced by Nomarski prism.

is placed far from the back focal plane, the prism produces parallel beams, but the beams are spatially displaced and hence are not recombined. Therefore, the Smith DIC system requires specially designed objective lenses that allow the insertion of a Wollaston prism.

In 1952 G. Nomarski proposed a special prism, the Nomarski prism, which simultaneously introduced spatial displacement and angular deviation of two orthogonally polarized beams^{77,78} (see inset in Fig. 31). The prism can therefore be placed outside the objective lens. By using crystal wedges with appropriately oriented axes, the Nomarski prism recombines the two beams that were separated by the condenser Wollaston prism, as though a regular Wollaston prism were located in the back aperture plane of the objective lens. The Nomarski DIC scheme can therefore be used with regular high NA microscope objectives.

A DIC image can be modeled as the superposition of one image over an identical copy that is displaced by a small amount d and phase shifted by a bias Γ . The intensity distribution $I(x, y)$ of the combined image depends on the specimen orientation and varies proportionally with the cosine of the angle between the gradients azimuth θ and the relative direction of wavefront shear σ :⁷⁹

$$I(x, y) = \frac{1}{2} I_0 \left(1 - \cos \left(\frac{2\pi}{\lambda} (\Gamma + \gamma(x, y) d \cos(\theta(x, y) - \sigma)) \right) \right) \quad (16)$$

where I_0 is the initial beam intensity, $\gamma(x, y)$ and $\theta(x, y)$ are the gradient magnitude and azimuth. (For a theoretical framework of DIC imaging see Refs. 8, 79–81.)

Thus, regular DIC techniques show the two-dimensional distribution of optical path or phase gradients projected onto the shear direction. It is therefore prudent to examine unknown objects at several azimuth orientations.

Video-enhanced DIC (VE-DIC), in addition to providing images with improved contrast, allows the removal of unwanted background signal (such as shading and fixed image noise due to dust particles or other imperfections in the optical system) by subtraction of a reference image with no specimen.⁸² Salmon and Tran gave a comprehensive description of the VE-DIC method.⁸³ They indicate that the best optical contrast of microscopic, lowly refractile particles can be achieved with a bias of 1/15–1/20 the wavelength.

A further increase in sensitivity and sectioning capability was achieved by video-enhanced DIC microscopy with retardation modulation.^{84–87} By switching the polarization of the incident light in alternate video frames with a computer-controlled liquid crystal variable retarder, the contrast signal is increased by a factor of 2, relative to “standard” video-enhanced DIC. The modulator switches image highlights into shadows and vice versa. By subtracting alternate frames, a difference DIC image is created in which contrast is doubled while image defects and noise tend to be cancelled.

Recently, Carl Zeiss introduced a “C-DIC” technique for reflective-type microscopes, which avoids the need to rotate the specimen. Instead, the new system uses a single, mechanically rotatable Nomarski prism that is shared between the illumination and imaging path.⁸⁸

Even in a transmission-type microscope one can obtain a DIC image using only one Wollaston or Nomarski prism placed in the imaging path, if the illumination beam is made spatially coherent. Pluta described a DIC setup with a slit condenser diaphragm.⁸ A similar system is currently manufactured by Carl Zeiss called a “PlasDIC.” In the latter case the specimen is illuminated with unpolarized light using a condenser that has a slit in its aperture plane. Only a single polarizer is used and placed behind the Nomarski prism that follows the objective. The system is less sensitive to birefringence of the specimen, can be used with plastic dishes, and does not require strain-free optics. Disadvantages include a reduced illumination intensity caused by the slit (instead of a fully open aperture) and a deterioration of the optical sectioning capability.

The contrast in DIC images is proportional to the scalar product between the phase gradient in the specimen and the shear generated in the microscope’s prisms. Based on the phase gradient it is possible to restore the phase information using computational methods. The restored phase image shows the refractive index (dry mass) distribution within a thin layer of the specimen. Compared to a conventional phase-contrast image, the DIC-based phase image provides better sectioning due to the intrinsic sectioning capability of the DIC method.

The DIC phase image can be obtained by computing the line integral parallel to the shear direction.^{79,89} Other techniques use iterative phase computation,⁹⁰ noniterative Fourier phase integration,⁹¹ or nonlinear optimization using hierarchical representations.⁹² Axelrod et al. used two phase-shifted DIC images to reconstruct the phase based on linearized expressions of interference.⁹³ Biggs developed an iterative deconvolution approach for computation of phase images, based on the same principles as deconvolution techniques normally used to remove out-of-focus haze.⁹⁴

Dyson and Mirau Interference Microscopes A third group of interference microscopes, in which the reference wave is focused to a different level than the specimen plane, are represented by the Dyson and Mirau interference microscopes.

In 1950, J. Dyson designed a double-focus system for transmitted light (Fig. 33a).⁹⁵ In this system, the initial transparent beam-dividing surface is formed on the upper side of the first glass plate mounted normally to the optical axis of the microscope immediately beneath the specimen slide. The illuminating beam, convergently directed through this upper surface by the substage condenser, is then partially reflected back to the lower surface of the plate, which has a small opaquely silvered, reflecting central spot. The lower surface of the first plate therefore reflects this second beam back through the upper surface of the plate. As a result, the specimen area is illuminated by two beams, one of which is focused on the specimen after direct transmission through the plate while the other reaches the specimen in a defocused condition due to internal reflection within the plate. A similar plate between the specimen and the objective lens functions in much the same way, so that the portion of the second defocused beam, which passes directly through it becomes combined with a portion of the first focused beam internally reflected within it.

The image formed by the microscope objective consequently consists of a correctly focused image of the specimen area seen in interferometric comparison with a strongly defocused image of it.

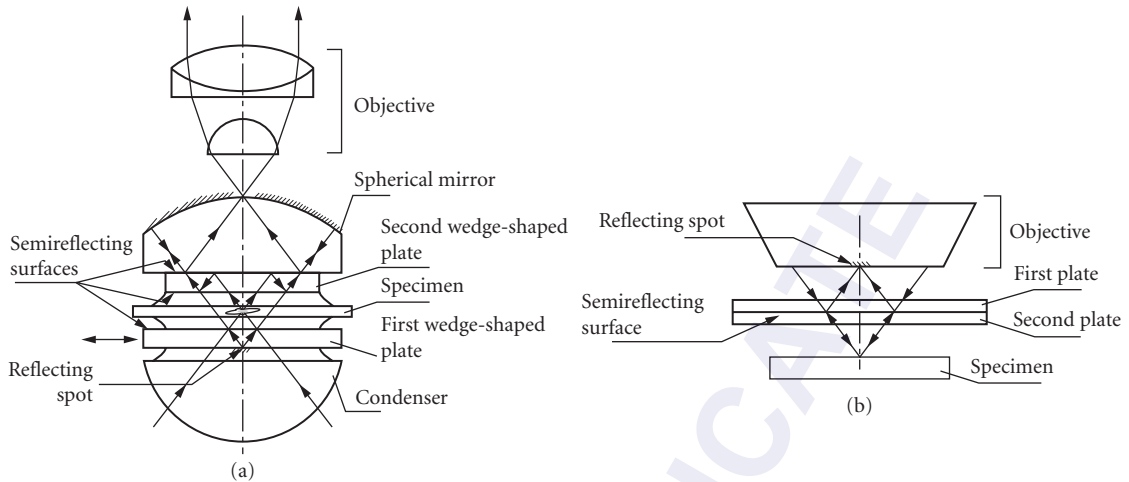


FIGURE 33 (a) Dyson and (b) Mirau's interference microscope. In (b), the incident light beam, emerging from the objective O_1 , is split in two parts in the semireflective surface. One part is transmitted to the object P and the other is reflected to the reference area R extending over a small portion of the objective front surface. The wavefronts reflected by R and P are recombined at G to produce the interference pattern.⁵

A glass block with an upper spherical surface, which is fully reflecting apart from a central totally transmitting spot is included between the second plate and the objective to allow medium- and high-power objective lenses to focus through to the specimen. The two plates are made slightly wedge-shaped so that the optical path difference can be manually adjusted by traversing the condenser plate in a direction parallel to the principal section of the wedge across the optical axis of the microscope. This operation varies the effective thickness difference between the two plates and thereby controls the optical path difference. By calibrating this movement the optical path difference can be determined.

Mirau introduced a single objective reflecting system.⁹⁶ In this design (Fig. 33b) a flat, semireflecting beam-dividing surface is placed midway between the front of the objective and the specimen surface. A small central area of the front surface of the objective is silvered to form a miniature mirror, a reflected image of which becomes superimposed on the normal image of the specimen surface by virtue of the intervening semireflecting beam divider. To maintain the required degree of optical path similarity, the beam-dividing surface is formed on the internal side of one of a pair of identical plates, which are cemented together.

Holographic Soon after its invention the laser was employed for holographic imaging in microscopy. In the early 1960s, Gordon Ellis built one of the first holographic microscopes.⁹⁷ He used a helium-neon laser as light source and photographic film for recording the hologram. After development of the film, the hologram allowed to reconstruct images using a divergent laser beam.

In digital holographic microscopy (DHM), the hologram of interfering wave fronts is recorded with an electronic sensor (e.g. CCD chip, Fig. 34)⁹⁸ and images are digitally reconstructed by a computer. A digitized hologram represents a three-dimensional record of the optical features of the specimen. Based on a single hologram, several images can be reconstructed that correspond to specific focus planes in the specimen. Furthermore, the digital reconstruction allows to simulate different contrast modes, such as phase contrast and dark field imaging.^{99–101} In addition to the specimen, the hologram can also contain information on the rest of the optical path, depending on the coherence length of the light source. For example, a hologram can provide the opportunity to correct for lens aberrations.¹⁰²

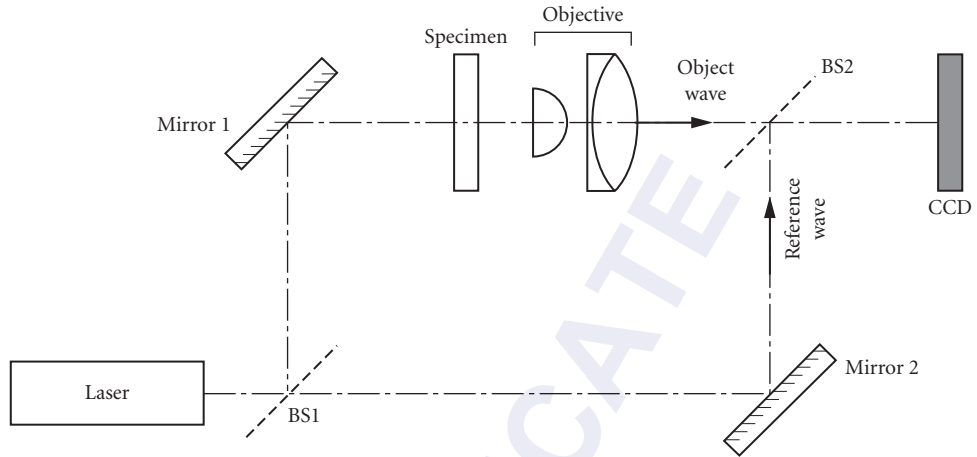


FIGURE 34 Optical principle of a holographic microscope. A collimated laser beam is divided by the beamsplitter BS1. One beam passes through the specimen and the microscope objective lens and forms the object wave. The second beam is the reference wave and is recombined on-axis with the first beam behind the objective lens. The interference pattern (hologram) of the object and reference wave is recorded by a CCD camera that is located near a conjugate plane of the backfocal plane of the objective lens. Other optical setups are possible, including for reflective-type specimens and for using an off-axis interference arrangement.⁹⁸

The digital analysis of a set of holograms, each recorded with a beam that illuminates the sample from a different direction, allows to emulate an objective with a larger numerical aperture than actually employed, leading to a corresponding enhancement in image resolution.¹⁰³

Optical Coherence Tomography Optical coherence tomography (OCT) is an imaging method that performs depth-resolved imaging of various turbid media. At the core of the OCT technique is a low-coherence, two-arm interferometer, which works in reflection mode.^{104,105} The low-coherence interferometer is used to select only a small volume named the “coherence gate” that determines the depth in the sample from where the back-reflected or back-scattered signal is processed. The depth of the coherence gate is defined and controlled by matching the optical path in the probe and reference interferometer arms. A variable delay line in one of the arms changes the gate position. In addition to its depth-selectivity feature, the low coherence interferometer is used to “amplify” very weak signals back-scattered by the sample.

In OCT, the coherence length is shortened to a distance of several micrometers, thanks to the use of a broadband light source. Light of appropriate bandwidth is typically generated by a superluminescent diode or laser with extremely short pulses. The spatial resolution of OCT in the axial direction is provided by the coherence gate, which selects signal light only from a cross-sectional volume of thickness defined by the coherence length of the illumination source. Superluminescent diodes typically provide 10- to 20- μm axial resolution. Higher resolution can be obtained with ultrashort pulsed lasers.

Interference of the light reflected by the sample and the reference mirror in the interferometer arms can occur only when the optical path lengths of the two arms match to within the coherence length of the optical source. Depth scanning can be performed in the time- or spectral domain. Time-domain OCT systems vary the reference arm path length, inducing changes in the depth from which the backscattered light of the sample is detected. In spectral-domain or Fourier-domain OCT, the axial signal intensity is calculated based on changes in the interference spectrum. The interference between probe and reference beam causes changes in the spectrum which is measured using a suitable spectrometer,¹⁰⁶ or by rapidly and repeatedly sweeping a narrow line width laser source in a mode called swept-source OCT.^{106–108}

In addition to amplitude and phase, OCT can also be used to analyze changes in polarization of the probe beam, revealing the polarization properties (birefringence, dichroism) of selected regions inside a turbid medium.¹⁰⁹

Optical coherence microscopy (OCM) combines the advantages of confocal microscopy with the principles of low-coherence interferometry.^{104,110,111} High contrast and detection sensitivity are achieved via rejection of out-of-focus light, resulting in improved optical sectioning capabilities deep within highly scattering media. Both OCT and OCM usually employ single-mode optical fibers for illuminating and collecting light from the sample. However, OCT uses a low NA objective lens with an extended depth of field, providing sectioning through coherence only. OCM, on the other hand, utilizes a high NA lens, providing sectioning through a combination of coherence and confocal effects.

Polarizing

The polarizing microscope (Fig. 35) generally differs from a standard transilluminating microscope by the addition of a polarizer before the condenser; a compensator slot and analyzer behind the objective lens; strain-free optics; a graduated, revolving stage; centrable lens mounts; cross-hairs in the ocular aligned parallel or at 45° to the polarizer axes; and a focusable Bertrand lens that can be inserted for conoscopic observation of interference patterns in the back aperture of the objective lens. In addition, the front element of the condenser can be swung into place for higher-NA conoscopic observations or swung out for low-NA orthoscopic observations of larger field areas.

The same components can be made to fit on an epi-illumination stand for observing opaque or reflective-type samples, such as in metallurgy. As outlined earlier, in epi-illumination a beam-splitting mirror separates the illumination and imaging light paths before the objective lens. In polarizing microscopy one needs to pay special attention to the beam-splitting mirror, which typically introduces polarization aberrations. The aberrations can be significantly reduced by a so-called Smith reflector replacing the regular dichromatic or half-shaded mirror. While a regular beam splitter reflects the incoming beam with a 45° angle of incidence, the Smith reflector uses two 22.5° reflections, first off a full mirror, followed by a second reflection off a 50/50 beam splitter. While the number of reflections has doubled, the steeper angle of incidence of 22.5° for both reflections reduces the overall polarization distortions compared to a regular beam splitter.

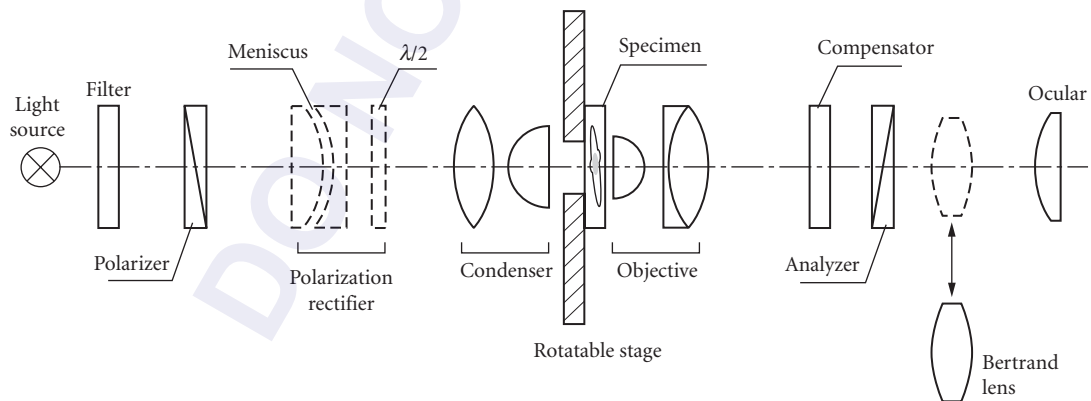


FIGURE 35 Optical train of a polarized light microscope with polarizer, analyzer, and compensator. An optional polarization rectifier can achieve improved sensitivity for low retardance measurements. With an optional Bertrand lens one examines the objective back focal plane for conoscopic interference figures.

Equipped with these standard components, the polarizing microscope represents a powerful analytical tool for the identification of crystals, fibers, and other optically anisotropic materials.^{51,68} With standard polarizing microscopes, one can image and measure polarization optical parameters on objects which are larger than a few micrometers and which introduce retardance values greater than several tens of nanometers. However, as the dimension of the object or magnitude of retardance decrease below these ranges, one needs to use special techniques or devices for detecting and measuring birefringence or even for generating a reliable image with high-NA lenses.

The basic ingredients that are needed to detect low levels of birefringence (retardance ≤ 10 nm) are high-extinction optics, use of low-retardance compensator, light source with high irradiance, and high-sensitivity detector (e.g., dark adaptation for visual observation and measurements). The need for high extinction optics applies to all components of the polarization optical train, which starts and ends with the polarizer and analyzer, respectively, and all optical components placed between them. Most manufacturers carry objective and condenser lenses that are either made or specially selected for polarized light observations. Such objectives typically carry the designation P, PO, or POL on their lens barrel and are designed to induce minimal polarization distortions (see Table 2).

Dichroic polarizing filters have replaced calcite prisms (which introduce astigmatism to all but collimated rays) as polarizer and analyzer in all but the most demanding applications. Modern dichroic polarizers are available with extinction factors better than 1000 and transmission better than 90 percent of the light that is fully polarized parallel to the transmission axis. These specs are satisfactory for most applications, in part because even POL-designated microscope lenses that are placed between the polarizer and analyzer cause polarization distortions that typically reduce the extinction of the polarization optical train significantly below 1000. The polarization distortions are typically caused by stress birefringence in the lens glass and by the differential transmission and phase shift of polarized light that passes through the peripheral regions of highly curved lens surfaces.¹¹² The latter distortions result in four bright quadrants separated by a dark cross (the Maltese cross) that is seen conoscopically for crossed linear polarizers in the absence of a birefringent specimen. These distortions also give rise to anomalous diffraction, based on a four-leaved clover pattern replacing the Airy disk or each weakly birefringent image point.^{113,114}

To counteract polarization distortions that occur at high NA lens surfaces, Inoué and colleagues have introduced polarization rectifiers^{115,116} made of a zero power lens with meniscus and a half-wave plate (Fig. 35). Using rectified optics Inoué and Sato¹¹⁷ were able to reveal the chromosome arrangement in living sperm based on high-contrast polarized light images (Fig. 36). Unfortunately, rectifiers are commercially not available for modern microscope objectives, which contain many lens elements and complex antireflection coatings, making the construction of a rectifier highly specific to each objective and condenser lens. However, some manufacturers have succeeded better than others in selecting antireflection coatings that minimize the polarization distortions. Therefore, it is advisable to carefully select microscope optics, testing the polarization performance of similar lenses from several manufacturers and even within the product range of the same manufacturer, before acquiring critical components.

The compensator is located between the polarizer and analyzer, either before or after the specimen. There are several types of compensators (often named after their inventors), which are typically made of birefringent plates or wedges that can be translated or rotated in fine increments while observing the specimen.¹¹⁸ The effect of the compensator on the polarization of the transmitted or reflected light either adds to or subtracts from (compensates) the effect caused by the specimen. While not absolutely necessary for some basic observations, the compensator (a) can significantly improve the detection and visibility of weakly birefringent objects, (b) is required to determine the slow and fast axis of specimen birefringence, and (c) is an indispensable tool for the quantitative measurement of object birefringence (see, e.g., Ref. 119); for a discussion of the Poincaré sphere as an analogue device to compute the effect of compensators, or of birefringent objects in general, on polarized light see Ref. 120).

Over the years several schemes have been proposed to automate the measurement process and exploit more fully the analytical power of the polarizing microscope. These schemes invariably involve the traditional compensator, which is either moved under computer control¹²¹ or replaced

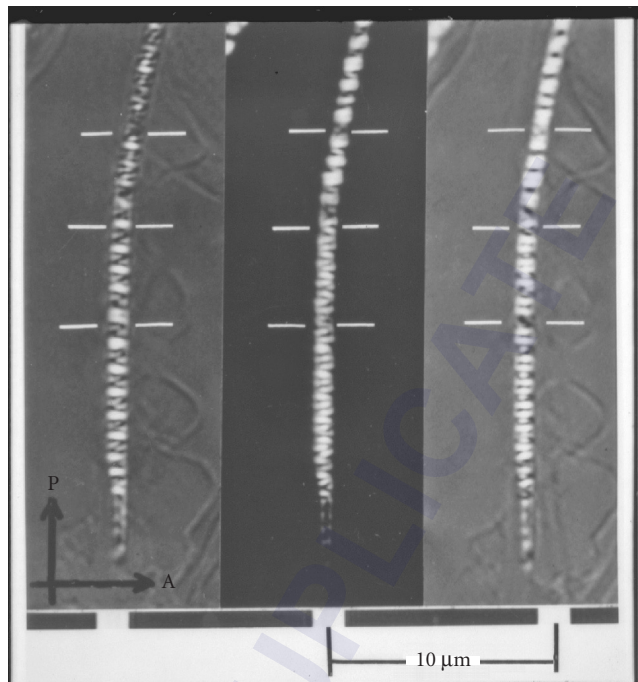


FIGURE 36 Sperm head observed with a rectified polarizing microscope at three different settings of mica compensator. Detailed distribution of birefringence in these chromosomes is shown with great clarity by immersion in dimethyl sulfoxide (refractive index 1.475). White bars: positions of chromosomal “breaks”; probably correspond to the end of each chromosome.¹¹⁷

by electro-optical modulators, such as Pockel-cells,¹²² Faraday rotators,¹²³ and liquid crystal variable retarders.¹²⁴ These schemes also involve quantitative intensity measurements using electronic light detectors, such as photomultipliers or charge-coupled device (CCD) cameras. For strictly quantitative measurements, acquisition and processing algorithms relate measured image intensities and compensator settings to optical characteristics of the specimen (see, e.g., Ref. 125). As an example of a quantitative, high-resolution birefringence map, we show in Fig. 37 the retardance image of a Siemens star that was etched into a thin silica layer.⁵⁰ The image was recorded using the LC-PolScope equipped with a liquid-crystal universal compensator.¹²⁶

Polarized light microscopy is usually practiced in two, mutually exclusive observation modes, called orthoscopy and conoscopy. In orthoscopy, the specimen is viewed directly, while in conoscopy the ocular is replaced by a telescope lens that lets one observe conoscopic interference figures formed in the back focal plane of the objective lens.⁶⁸ In conoscopy, the sample birefringence is measured as a function of the tilt angle of rays passing through the specimen. Hence, this observation mode reveals the inclination angle of the optic axis of a uniformly birefringent specimen region, in addition to the azimuth of the optic axis. In orthoscopy, the inclination angle, which is the angle between the optic axis and the plane of observation, is usually not evident. Recently, orthoscopic and conoscopic views were combined in a single, so-called polarized light field image recorded with a microlens array in the intermediate image plane of an LC-PolScope.²¹

Another approach to measuring the three-dimensional birefringence properties of small birefringent objects uses a so-called universal stage, invented by E.S. Fedorov more than 100 years ago, in

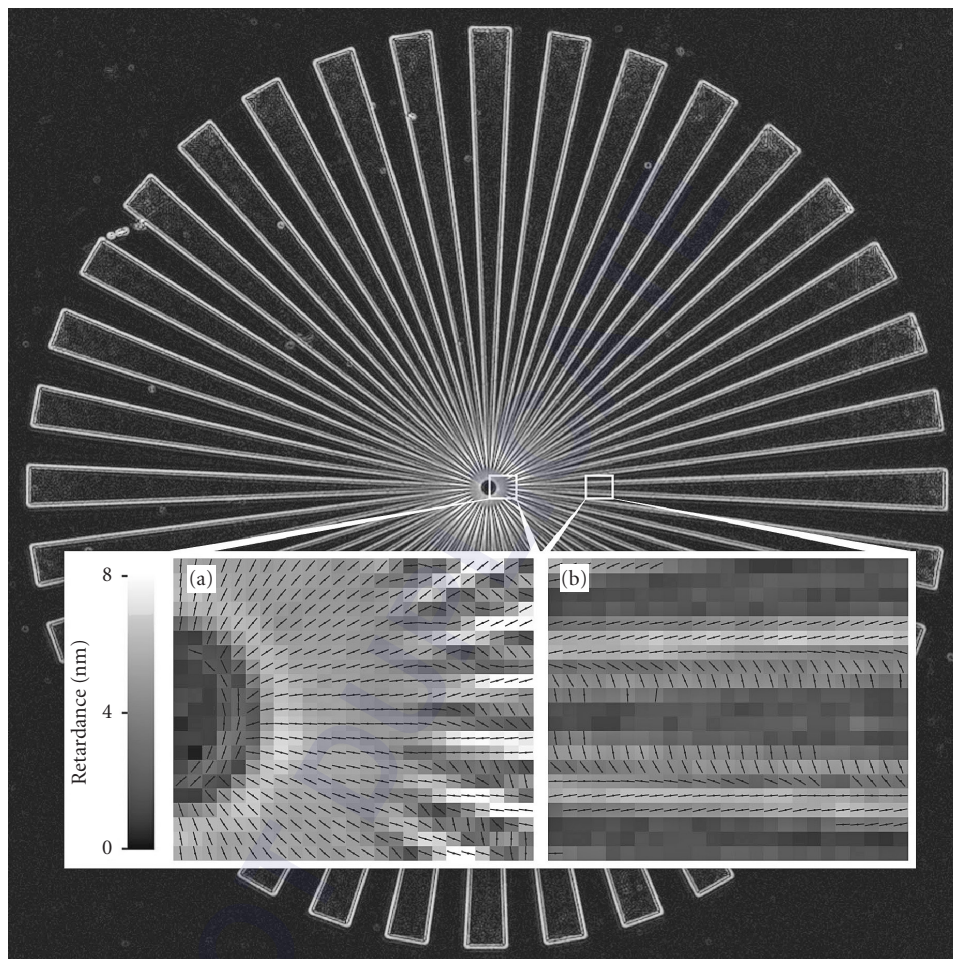


FIGURE 37 Siemens star etched into 90-nm-thick SiO_2 layer and imaged with the LC-PolScope. The dimensions of the star pattern are the same as the one described in Fig. 21. Image brightness is linearly proportional to the retardance measured at all pixel locations. Insets show magnified portions of the pattern with lines indicating the measured slow axis orientation. (a) In the central region birefringence is observed all the way to the inner black disk. Unresolved wedge tips generate form birefringence with the slow axes parallel to the wedge orientations. (b) Edges of a well-resolved wedge portion display edge birefringence, which is composed of two birefringent layers flanking each edge.¹²⁷ ($60\times/1.4$ NA PlanApo oil objective lens and oil condenser with aperture diaphragm reduced to 0.9 NA.)

which the specimen is mounted between two glass hemispheres.²² Rotation of the specimen through measured angles around two or more axes allows one to explore the three-dimensional birefringence patterns of a small specimen region that is located in the common center of rotation. Alternatives to the universal stage include the spindle stage by Bloss¹²⁸ and motorized goniometric stages by Glazer and collaborators.^{129,130}

Instead of rotating the specimen under a stationary optical system, Shribak and Oldenbourg implemented a scheme involving a high numerical aperture imaging system and oblique illumination with varying tilt angle of the illuminating beam.¹³¹ For each angle a high-resolution retardance map

is generated representing the polarization properties of the sample as projected along the tilted axis of illumination. Four such maps, each generated with a different tilt angle, are combined to produce a three-dimensional birefringence map. The system is called *Scanned Aperture LC-PolScope* and is described here in more detail in the section “Aperture Scanning.”

Fluorescence

Fluorescence microscopy is one of the few modes of microscopy in which the illuminating wavelength differs from that of the emitted. In early designs, the exciting waves were prevented from contaminating the fluorescence image by a combination of (1) special illumination (such as the use of a dark-field condenser) that prevented the direct rays from entering the objective lens, and (2) the use of a barrier filter. The barrier filter absorbs the exciting light while transmitting much of the longer fluorescence wavelengths.

Today most fluorescence microscopes (or attachments) use epi-illumination incorporating interchangeable filter cubes (after Ploem, see Fig. 6) that are matched to the fluorochrome. The filter cube is placed in the collimated beam between the objective and a tube lens, at the intersection of the microscope axis and that of the excitation illuminator located on a side arm. The objective lens serves both as the condenser and the objective. A field diaphragm, and sometimes an aperture iris, is placed in the illuminating side arm together with the source collector at appropriate conjugate planes. The illuminating beam, commonly emitted by a xenon or mercury arc lamp, is filtered through a narrow band path interference filter and reflected down into the objective by a dichromatic beam splitter. The fluorescence imaging beam originating from the specimen passes straight through the dichromatic beam splitter and associated barrier filter and reaches the ocular or camera. Each fluorescence cube contains the appropriate excitation interference filter, dichromatic beam splitter, and barrier filter so that they can be switched as a group, for example, to rapidly inspect specimens containing (or stained with) multiple fluorochromes.

For fluorochromes requiring shorter-wave UV excitation, objective lenses must be designed for greater short-wavelength transmission and low autofluorescence. While aberrations for the shorter-UV exciting wavelengths are generally not as well-corrected as for the imaging wavelengths, it should be noted that such aberrations, or lack of parfocality, directly affect the resolution and efficiency in the case of confocal fluorescence microscopes.

Also, it should be noted that, while little effort is commonly made to fill the objective aperture with the illuminating beam (presumably with the rationale that this should not affect image resolution because each fluorescent object is emitting incoherently relative to its close neighbor), one finds that in practice the fluorescent image is much improved by filling the aperture, for example, by use of an optical fiber light scrambler. While the reasons for this improvement are not fully understood, one explanation might lie in the more efficient excitation of randomly oriented fluorophores by a high-NA illumination beam, which excites even those fluorophores that have their linear transition moment aligned parallel to the microscope axis.

While most fluorescence microscopes today use epi-illumination (since epi-illumination provides advantages such as avoiding loss of excitation by self-absorption by underlying fluorochrome layers, generating an image that more closely approximates an intuitive one when reconstructed in three dimensions, etc.), improvements in interference filters open up new opportunities for fluorescence microscopy using transillumination. New interference filters are available with exceptionally high extinction ($>10^5$) and sharp cutoff of the excitation wavelengths, coupled with high transmission of the pass band. With transillumination, one can more reliably combine fluorescence with polarization-based microscopy or carry out polarized fluorescence measurements with greater accuracy, since one can avoid the use of dichromatic beam splitters, which tend to be optically anisotropic.

Fluorescence microscopy particularly benefits from electronic imaging, especially with the use of low-noise, chilled CCDs as imaging detectors, digital computers to enhance and rapidly process the signal (such as in ratio imaging), and the new fluorescence-conjugated chemical probes that provide incredible sensitivity and selectivity.^{11,12,132-135}

For imaging specimens that are labeled with more than two or three types of spectrally distinct fluorophores, a technique known as spectral imaging is becoming available. Spectral imaging combines spectroscopy and imaging, measuring the spectral composition of the light recorded at each point of the image. When spectral imaging is applied to fluorescence microscopy, the filter cube is modified as to transmit a broad range of emission wavelengths. A spectrometer placed before the detector samples the emission spectrum at appropriate resolution and intervals (channels) for wavelengths longer than the excitation wavelength. Spectral imaging systems can either be integrated into the microscope (manufacturers include Leica, Nikon, Zeiss) or can be added to an existing stand (manufacturers include Cambridge Research and Instrumentation Inc., Lightform Inc.). Datasets are typically stored as stacks of images, in which each slice corresponds to a wavelength channel. Powerful algorithms reduce an experimental dataset to indicate for each image point the weighted contributions of pure fluorophores whose spectra are stored in a library.^{136,137}

Confocal Microscopy

In confocal microscopy, the specimen is scanned point by point either by displacing the specimen (stage scanning) or by scanning a minute illuminating spot (beam scanning), generally in a TV-raster fashion. In either case, the scanning spot is an Airy disk formed by a high-NA objective lens. An exit pinhole is placed conjugate to the spot being scanned so that only the light originating from the scanned spot is transmitted through the exit pinhole. Thus, light originating from other regions of the specimen or optical system is prevented from reaching the photo detector (Fig. 38).^{138,139}

This optical arrangement reduces blurring of the image from out-of-focus light scattering, fluorescence, and the like, and yields exceptionally clear, thin optical sections. The optical sections can then be processed and assembled electronically to yield three-dimensional displays or tilted plane projections. Alternatively, the specimen itself can be scanned through a tilted plane (e.g., by implementing a series of x scans with y, z incremented) to yield a section viewed from any desired orientation, including that normal to the microscope axis.

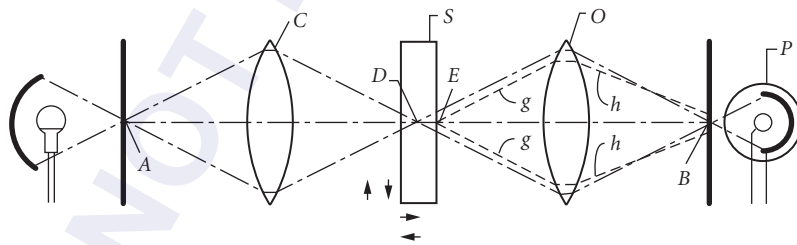


FIGURE 38 Optical path in simple confocal microscope. The condenser lens C forms an image of the first pinhole A onto a confocal spot D in the specimen S . The objective lens O forms an image of D into the second (exit) pinhole B which is confocal with D and A . Another point, such as E in the specimen, would not be focused at A or B , so that the illumination would be less and, in addition, most of the light $g-h$ scattered from E would not pass the exit pinhole. The light reaching the phototube P from E is thus greatly attenuated compared to that from the confocal point D . In addition, the exit pinhole could be made small enough to exclude the diffraction rings in the image of D , so that the resolving power of the microscope is improved. The phototube provides a signal of the light passing through points D_1, D_2, D_3 , etc. (*not shown*), as the specimen is scanned. D_1, D_2, D_3 , etc. can lie in a plane normal to the optical axis of the microscope (as in conventional microscopy), or parallel to it, or at any angle defined by the scanning pattern, so that optical sections can be made at angles tilted from the conventional image plane. Since, in the stage-scanning system, D is a small spot that lies on the axis of the microscope, lenses C and O can be considerably simpler than conventional microscope lenses.^{138,139}

The stage-scanning confocal microscope can yield vastly expanded fields of view. Here the image area is not limited by the field of view of the optics but only by the range of movement of the specimen and ability of the photo detector and processor to handle the vast information generated at high speed. Furthermore, the objective lens needs only to be corrected for a narrow field of view on axis.^{6,138} Laser disk recorders are a form of application that takes advantage of these attributes.

The beam-scanning confocal microscope is typically implemented in the reflective or epi-illumination mode. This mode has the advantage that the illuminating beam and the returning light scattered back by the sample pass through the same objective lens and beam-steering devices needed for scanning the sample. The prototype of a modern beam-scanning confocal microscope uses two galvanometric mirrors (one for each dimension of a two-dimensional image) that scan a focused laser beam across a stationary sample field. The backscattered light is collected by the objective and bounces off the same mirrors which “descan” the returning light before it passes through a stationary beamsplitter (to separate the backscattered light from the incoming beam) and a stationary exit pinhole. The exit pinhole is located in a conjugate image plane, while the scanning mirrors are located in positions that are conjugate to the back focal plane of the objective lens. By (indirectly) placing the mirrors into the objective back focal plane, the angular scan of the mirrors is translated into a positional scan of the focused laser beam in the specimen. Beam-scanning microscopes typically require additional relay optics that project the objective back focal plane into the mirror locations.

The laser-scanning, epi-illuminating confocal microscope was developed into a practical instrument in the late 1980s and immediately adopted with great enthusiasm for fluorescence imaging in the life sciences. Because laser beams are typically highly collimated, a source or entrance pinhole is commonly omitted in this instrument. The beam splitter combining and separating the illumination and imaging paths is implemented as a dichroic (also called dichromatic) mirror providing high reflectivity at short wavelengths and high transmissivity at longer wavelengths (or vice versa, depending on the particular optical design). A whole industry has evolved around designing and manufacturing dichromatic mirrors that are appropriate for specific fluorescent dyes and combination of dyes.

For direct viewing of confocal images in reflective mode a Nipkow disk is used for scanning multiple beams across a stationary sample field. The multiple beams originate in many thousands of pinholes arranged helically on a modified Nipkow disk that is located in the image plane of the objective lens. Thus, a single spinning disk can be made to provide synchronously scanning entrance and exit pinholes.^{140,141} To overcome the considerable light loss associated with the original designs by Petrán and Kino, Yokogawa Electric Corp. employed a second, coaxially aligned Nipkow disk containing microlenses in its CSU-10 disk confocal scanner (Fig. 39). Each pinhole on the first Nipkow disk has a corresponding microlens on the second Nipkow disk that focuses the laser light into the pinhole. Thus, the light efficiency is increased by a factor equal to the ratio of the microlens to pinhole area. Instead of the 1 percent or so found with conventional Nipkow disk systems, some 40 to 60 percent of the light impinging on the disk containing the microlenses becomes transmitted through the pinholes to illuminate the specimen. Accordingly, the CSU-10 provides a light efficient scan unit that permits direct visual viewing of the confocal image, a great advantage when studying moving objects such as living cells.^{142,143}

In a confocal microscope, the exit pinhole can be made smaller than the diameter of the Airy diffraction image formed by the objective lens so that the Airy disk is trimmed down to regions near its central peak. With this optical arrangement, the unit diffraction pattern that makes up the image turns out to be the square of the Airy pattern given in Eq. (2). Thus, the radius at half maximum of the central bright area (Airy disk) is reduced by a factor of 1.36. (The radial position of the first minimum in both patterns is still equal to r_{Airy} .) The shape of the unit diffraction pattern is thus sharpened so that, compared to nonconfocal imaging, two points which radiate incoherently (as in fluorescence microscopy) can be expected to approach each other by up to a factor of $\sqrt{2}$ closer to each other before their diffraction patterns encounter the Rayleigh limit. In Fig. 20 the contrast transfer characteristics of a confocal microscope in the coherent imaging mode is compared with the nonconfocal, incoherent imaging mode using the same lenses. Note that the limiting resolution is the same for both imaging modes, while the contrast transfer of the confocal mode increase more steeply for increasing grating periods.

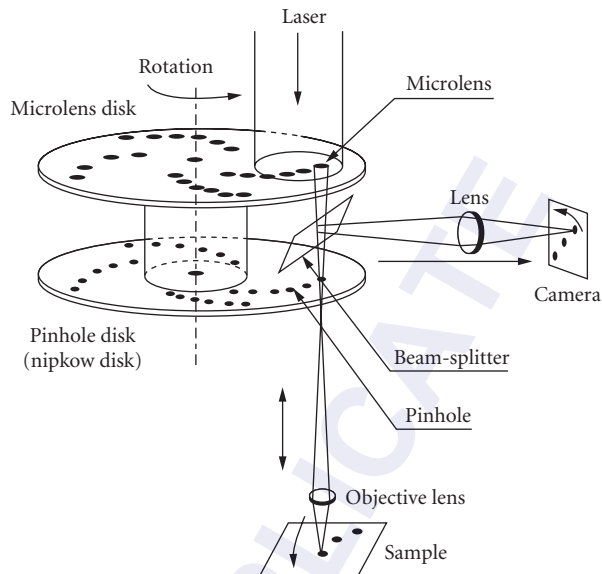


FIGURE 39 Schematic of the Yokogawa CSU-10 confocal disk scanner. The expanded and collimated laser beam illuminates the upper Nipkow disk containing some 20,000 microlenses. Each microlens focuses the laser beam onto its corresponding pinhole, thus significantly raising the fraction of the illuminating beam that is transmitted by the main Nipkow disk containing the pinhole array. The backscattered or fluorescent light is collected by the objective lens and focused back onto the Nipkow disk containing the array of pinholes, which now act as confocal exit pinholes. A beam splitter located between the first and second Nipkow disk reflects the light toward a camera. A lens projects an image of the pinhole array onto the camera, that acquires a confocal image while the Nipkow disk rotates with high speed. After carefully designing the array pattern and implementing a precise and vibration-free rotation of the Nipkow disks, the confocal disk scanner can produce clean, high-resolution images free of fixed pattern noise. In fluorescence imaging, the camera can be replaced by an ocular for direct viewing of the confocal image. (Schematic provided by Yokogawa Electric Corporation.)

Rather than using confocal optics to eliminate image blurring from out-of-focus planes, one can achieve the same end by computational deconvolution of a stack of serial optical sections obtained by wide field microscopy.^{49,144,145} While computationally intensive and time consuming, this image restoration method allows one to isolate clean optical sections from a stack of images that can be acquired at higher speed and higher efficiency than with laser-scanning confocal microscopy and in modes of contrast generation typically not accessible to confocal imaging.

Alternatively, thin optical sections can be obtained directly with digital enhanced video microscopy using high-NA condenser and objective lenses. Requiring little processing, this approach is especially convenient when many stacks of optical sections have to be acquired at high rates in succession, for example in order to record rapid, three-dimensional changes in microscopic domains over time.

Structured Illumination

The quest for improved resolution beyond the diffraction limit has led to the development of several methods that modify the illumination pattern in wide-field microscopy. In standard wide field microscopy, the specimen is illuminated using condenser optics that ideally projects a uniform field of light into the specimen. In structured illumination setups, however, a finely patterned illuminating field is projected into the specimen, providing a means for generating optical sections similar to confocal microscopy and for improving resolution.

Wilson and colleagues¹⁴⁶ first described a simple method of obtaining optical sectioning in a conventional wide-field microscope by projecting a single-spatial-frequency grid pattern onto the object. Images taken at three spatial positions of the grid were processed in real time to produce optically sectioned images that are substantially similar to those obtained with confocal microscopes. The sectioning capability is achieved by superimposing an illumination pattern that is only in focus at a thin section through the specimen, while all other sections of the specimen are illuminated with a more or less blurred version of the pattern. The specimen with the superimposed grid pattern is then imaged with regular wide field optics focused on the grid pattern inside the specimen. Hence, all image features that have the grid pattern imposed on them are located in this specimen section, while image features from other sections of the specimen appear nearly uniformly illuminated. For removing those out-of-focus features and removing the intruding effect of the illumination pattern on the specimen image, three raw sample images are recorded, each with the illumination pattern slightly shifted in position. Subsequently, the raw images are computationally combined to generate an optical section of the sample without the grid pattern noticeable in the image. The company Carl Zeiss has adopted this strategy in its ApoTome slider module for generating optical sections using epi-illumination.

Instead with a regular grid pattern, the sample can also be illuminated with a random speckle pattern to provide depth discrimination in thick, fluorescently labeled tissues.^{147,148} The technique consists of illuminating a sample with a sequence of speckle patterns and displaying the differential intensity variance of the resultant sequence of fluorescence images. The advantage of speckle illumination is that it provides diffraction-limited illumination granularity that is highly contrasted even in scattering media.

Structured illumination strategies that go beyond optical sections and provide lateral resolution that exceeds the classical diffraction limit by a factor of 2 or more have been devised by Gustafsson.¹⁴⁹ The sample is illuminated with a series of excitation light patterns, which cause normally inaccessible high-resolution information to be encoded into the observed image. The recorded images are linearly processed to extract the new information and produce a reconstruction with twice the normal resolution. Unlike confocal microscopy, the resolution improvement is achieved with no need to discard any of the light arriving from the specimen.

In addition to improving the lateral resolution this method can be applied in three dimensions to double the axial as well as the lateral resolution, with true optical sectioning.¹⁵⁰ A grating is used to generate three mutually coherent light beams, which interfere in the specimen to form an illumination pattern that varies both laterally and axially. The spatially structured excitation intensity causes normally unreachable high-resolution information to become encoded into the observed images through spatial frequency mixing. This new information is computationally extracted and used to generate a three-dimensional reconstruction with twice as high resolution, in all three dimensions, as is possible in a conventional wide-field microscope.

Structured illumination is primarily used in fluorescence microscopy, where in principle it is capable of unlimited resolution. To achieve unlimited resolution, structured illumination has to be combined with a nonlinear dependence of the fluorescence emission rate on the illumination intensity.^{151,152} As an example of this concept, Gustafsson experimentally demonstrated saturated structured-illumination microscopy, in which the nonlinearity arises from saturation of the excited state. This method can be used in a simple, wide-field (non-scanning) microscope, which uses only a single, inexpensive laser, and requires no unusual photophysical properties of the fluorophore. The practical resolving power is determined by the signal-to-noise ratio, which in turn is limited by photobleaching. Experimental results show that a two-dimensional point resolution of < 50 nm is possible on sufficiently bright and photostable samples.

Light Field

Instead of increasing resolution in a single image plane, it is sometimes desirable to trade lateral resolution for axial resolution in a three-dimensional image stack. To this end, Marc Levoy and colleagues²⁰ have replaced the regular camera on a standard, wide-field microscope by a camera with microlens array, a so-called light field camera or plenoptic camera.¹⁵³ The array consists of a hundred thousand or more microlenses arranged in a square up to the size of the microscope's field number. The array is placed in the intermediate image plane of the objective lens. Behind the array the sensor chip is located in the backfocal plane of the microlenses. In other words, the light field camera samples the specimen image on a regular grid at intervals that corresponds to the pitch of the microlens array. At each grid point the camera captures a small subimage of the objective's back focal plane. Hence, the camera captures a hybrid image of the specimen that is sampled not only in space but also along different directions through the specimen.

The raw light field image, when presented to the eye, cannot be directly interpreted since it consists of a multitude of small disk-shaped images (of the objective's back focal plane) arranged on a regular grid. However, a single light field image is used to reconstruct a multitude of conventional images of a specimen that is viewed along different directions or focused to different object planes.²⁰ These differing views are all based on a single light field image that was captured by a single camera exposure. Hence light field microscopy can be especially useful when imaging three-dimensional structures that change rapidly in time, such as living cells and tissues. Based on a single snapshot one can reconstruct a stack of optical sections that were all recorded at the same point in time, thus avoiding registration problems between sections.

However, the versatility of generating different views and optical sections from a single light field image comes at a price. The sacrifice one makes is a reduction in image size. Specifically, if each microlens subimage contains $N \times N$ pixels, then the computed images will contain N^2 fewer pixels than if the microlenses were not present. In return, we can compute N^2 unique oblique views of the specimen, and we can generate a focal stack containing N slices with nonoverlapping depths of field.²⁰

The recording of light field images is compatible with several contrast modes, including fluorescence and polarized light microscopy. One of the first areas to take advantage of simultaneous optical sections was fluorescence microscopy of functional neuronal tissues and the recording of three-dimensional excitation patterns. In polarized light field microscopy, the microlens array generates a hybrid image consisting of an array of small conoscopic images, each sampling a different object area.²¹ Analysis of the array of conoscopic images reveals the birefringence of each object area as a function of the propagation direction of transmitted light rays. Compared to traditional conoscopy and related methods, the vastly improved throughput and quantitative analysis afforded by the light field LC-PolScope, for example, make it the instrument of choice for measuring three-dimensional birefringence parameters of complex structures. Since light field microscopy was implemented only a few years ago, additional application areas of this new method are likely to emerge in the future.

Aperture Scanning

In the aperture-scanning microscope devised by Ellis for phase-contrast microscopy, the tip of a flexible signal optical fiber, illuminated by an Hg arc, makes rapid circular sweeps at the periphery of the condenser aperture.⁶⁵ This circular, scanning illumination spot replaces the conventional phase annulus in the condenser aperture plane. A quarter-wave plate and absorber, both covering only a small area conjugate to the illuminating spot, spins in synchrony with the fiber at the objective back aperture (or its projected conjugate). Thus, the specimen is illuminated by a narrow, coherent beam that enters the specimen obliquely at high NA, with the azimuth orientation of the beam swinging around and around to generate a full cone of illumination within the integration time of the detector. With this aperture-scanning approach, the specimen is illuminated by a large-NA cone of light which is temporally incoherent, with the phase disk covering only a small fraction of the area normally occupied by the phase ring in conventional phase-contrast systems. The small size of the phase

disk, while appropriately reducing the amplitude and introducing the requisite $\lambda/4$ wave phase retardation to the rays not scattered by the specimen, allows the transmission of a much larger fraction of the scattered rays that carry the high spatial frequency information. The aperture-scanning phase-contrast microscope thus provides a very thin optical section. The image is also virtually free of the phase halo that obscures image detail adjacent to refractile boundaries in conventional phase-contrast microscopy.

For polarized light microscopy an aperture scanning scheme was designed and built by Shribak and Oldenbourg using a liquid crystal device in the front focal plane of the condenser lens.^{131,154} The liquid crystal device was designed for two functions: (1) to create oblique illumination of the specimen, and (2) to measure the birefringence parameters of microscopic objects for each of four oblique tilt angles of illumination. By measuring the object retardance along each of the four tilted projections, the inclination angle of the optic axis of birefringent objects was revealed, in addition to the orientation or azimuth angle in the plane of focus. The inclination angle of the optic axis is usually not evident from traditional polarized light images (see section on polarized light).

Extending this concept, modulation of the transfer functions of the condenser and objective apertures with electro-optical devices should open up intriguing new opportunities. Such modulation eliminates the need for mechanical scanning devices, the spatial distribution of the modulation function can be altered at will, and the amplitude and phase of light passing each point in the aperture can be adjusted rapidly, even coupled dynamically to the image signal through a feedback loop to generate dynamic spatial filters that enhance or select desired features in the image.

28.5 MANIPULATION OF SPECIMEN

In addition to viewing microscopic specimens, the light microscope and microscope objectives are also used to project reduced high-intensity images of source patterns into the object plane in order to optically manipulate minute regions of the specimen. Photolithography and laser disk recorders are examples of important industrial applications, which have prompted the design of specially modified objective lenses for such purposes.

Microbeam Irradiation, Ablation

Many applications are also found in the biomedical field, initially using UV-transmitting, moderately high NA objectives that are parfocalized for visible light and UV down to approximately 250 nm (Zeiss Ultrafluor, also quartz monochromats from Leitz). In its extreme form, a concentrated image of a circular- or slit-shaped UV or laser source of selected wavelengths is imaged onto a biological specimen to locally ablate a small targeted organelle; for example, a part of a chromosome, the microtubules attached thereto, or tiny segments of cross-striated muscle, are irradiated with the microbeam in order to sever their mechanical connections and, for example, to analyze force transduction mechanisms.^{155,156} In other cases, oriented chromophores can be selectively altered at the submolecular level, for example, by polarized UV microbeam irradiation. The stacking arrangement of the DNA nucleotide bases (which exhibit a strong UV dichroism, as well as birefringence in visible light) can be selectively altered and disclose the coiling arrangement of DNA molecules within each diffraction-limited spot in the nucleus of living sperm.¹¹⁷ Brief microirradiation of slit- or grid-shaped patterns of UV are used to bleach fluorescent dyes incorporated into membranes of living cells. The time course of recovery of fluorescence into the bleached zone measures the rate of diffusion of the fluorescently tagged molecules in the membrane and reveals unexpected mobility patterns of cell membrane components.^{157,158}

Lasers have become the dominant source for microbeam irradiation experiments in cell and developmental biology and in other application areas. Laser sources can have a wide range of tunable wavelengths (217 to 800 nm), energies, and exposure durations (down to 25×10^{-12}).¹⁵⁹ They are often used together with sensitizing dyes or fluorescent markers to target specific organelles.¹⁶⁰ They can be used in conjunction with versatile beam-shaping optics such as spatial light modulators.¹⁶¹

Photosensitive and Caged Compounds

Selected target molecules within minute regions in living cells can be modified, tagged, or activated by focused beams of light. The target molecules can be naturally photosensitive species such as chlorophyll (which produces oxygen where illuminated with the appropriate visible wavelengths), rhodopsin (which isomerizes and triggers the release of calcium ions and action potentials in retinal cells), or artificially introduced photosensitive reagents such as the drug colchicine (whose antimetabolic activity is abolished locally with 366-nm irradiation).

Of the photosensitive compounds, the *caged compounds* have a far-reaching potential. These are compounds that are synthesized so as to “cage” and hide the active chemical group until a photosensitive part of the compound is altered (e.g., by long-wavelength UV irradiation) and unmasks the hidden active group. Thus, by preloading with the appropriate caged compound and irradiating the cell selectively in the region of interest, one can test the role of the uncaged compound. For example, the role of ATP can be tested using caged ATP and ATP analogs; response to subtle increase in calcium ions can be seen using caged calcium or caged calcium chelators.^{162,163} Likewise, caged fluorescent dyes are irradiated to locally label and follow the transport of subunits within macromolecular filaments in a dividing cell.¹⁶⁴ Caged glutamate in brain slices was photolyzed using a holographically generated illumination pattern for simultaneous multispot activation of different dendrites.¹⁶¹

Optical Tweezers

Intense laser beams concentrated into a diffraction spot can generate a photon-driven force great enough to capture and suspend small particles whose refractive index differs from its surrounding.^{165,166} Applied to microscopy, a single swimming bacterium or micrometer-sized organelles in live cells can be trapped and moved about at will at the focus of a near-infrared laser beam focused by an objective lens of high NA. While the energy density concentrated at the laser focus is very high, the temperature of the trapped object remains within a degree or so of its environment; biological targets typically exhibit low absorbance at near-infrared wavelengths and thermal diffusion through water from such minute bodies turns out to be highly effective. Thus, the bacterium continues to multiply while still trapped in the focused spot, and it swims away freely when the laser beam is interrupted.

The ability to use “optical tweezers,” not only to capture and move about minute objects but to be able to instantly release the object, provides the microscopist with a unique form of noninvasive, quick-release micromanipulator.¹⁶⁷

Optical tweezers are now being used in the investigation of an increasing number of biochemical and biophysical processes, from the basic mechanical properties of biological polymers to the multitude of molecular machines that drive the internal dynamics of the cell. Innovation continues in all areas of instrumentation and technique, with much of this work focusing on the refinement of established methods and on the integration of this tool with other forms of single-molecule manipulation or detection. These developments have important implications for the expanded use of optical tweezers in biochemical research.¹⁶⁸

28.6 ACKNOWLEDGMENTS

Shinya Inoué was the principal author of the chapter on microscopes in the previous edition of the *Handbook of Optics*. The authors wish to thank Dr. Inoué for his many suggestions regarding the reorganization of the content, for a new contribution on the single-sideband edge enhancement (SSEE) technique, and for his consent to use much of the previous text and figures for the new edition. In acknowledging contributions to both editions, the authors also wish to thank Gordon W. Ellis of the University of Pennsylvania, late Katsuji Rikukawa, Yoshiyuki Shimizu, and Hiroshi Takenaka from Nikon K. K., Japan, Ernst Keller and Rudi Rottenfusser of Carl Zeiss, Inc., Jan Hinsch

of Leica, Inc., Mortimer Abramovitz of Olympus, Inc., and Lee Shuett and Mel Brenner of Nikon, Inc., who all provided helpful information and insights into microscope design.

The preparation of this chapter was supported in part by the National Institutes of Health grants R01 EB002583 awarded to R.O. and R01 EB005710 awarded to M.S.

28.7 REFERENCES

1. L. C. Martin, (1966) *The Theory of the Microscope*, Blackie, London.
2. S. Inoué and K. R. Spring, (1997) *Video Microscopy the Fundamentals*, Plenum Press, New York.
3. F. Zernike, (1942) "Phase Contrast, a New Method for the Microscopic Observation of Transparent Objects," *Physica* **9**:686–693.
4. D. Gabor, (1949) "Microscopy by Reconstructed Wavefronts," *Proc. Roy. Soc. London ser A* **197**:454–487.
5. M. Françon, (1961) *Progress in Microscopy*, Row, Peterson, Evanston, Ill.
6. T. Wilson and C. Sheppard, (1984) *Theory and Practice of Scanning Optical Microscopy*, Academic Press, London.
7. M. Pluta, (1988) *Advanced Light Microscopy Vol. I: Principles and Basic Properties*, Elsevier Science Publishing Co., Amsterdam.
8. M. Pluta, (1989a) *Advanced Light Microscopy Vol. II: Specialized Methods*, Elsevier Science Publishing Co., Amsterdam.
9. M. Pluta, (1993) *Advanced Light Microscopy Vol. III: Measuring Techniques*, Elsevier Science Publishing Co., Amsterdam.
10. G. Sluder and D. E. Wolf, (1998) Video Microscopy. In: *Methods Cell Biol.* (eds. L. Wilson and P. Matsudaira). Academic Press, San Diego.
11. Y.-L. Wang and D. L. Taylor, (1989) *Fluorescence Microscopy of Living Cells in Culture. Part A*, Academic Press, San Diego.
12. D. L. Taylor and Y.-L. Wang, (1989) *Fluorescence Microscopy of Living Cells in Culture. Part B*, Academic Press, San Diego.
13. J. B. Pawley, (2006) *Handbook of Biological Confocal Microscopy*, 3d ed. Springer, New York.
14. B. R. Masters and P. T. C. So, (2008) *Handbook of Biomedical Nonlinear Optical Microscopy*, Oxford University Press, Oxford.
15. S. Bradbury, P. J. Evennett, H. Haselmann and H. Piller, (1989) *Dictionary of Light Microscopy*, Oxford University Press, Oxford.
16. H. H. Hopkins and P. M. Barham, (1950) "The Influence of the Condenser on Microscopic Resolution," *Proc. Phys. Soc. London* **63 B**:737–744.
17. M. Born and E. Wolf, (2002) *Principles of Optics: Electromagnetic Theory of Propagation, Interference and Diffraction of Light*, Cambridge University Press, Cambridge.
18. G. W. Ellis, (1985) "Microscope Illuminator with Fiber Optic Source Integrator," *J. Cell Biol.* **101**:83a.
19. E. Becker, (1985) *Fluorescence Microscopy*, Ernst Leitz Wetzlar GmbH, Wetzlar, Germany.
20. M. Levoy, R. Ng, A. Adams, M. Footer and M. Horowitz, (2006) "Light Field Microscopy," *Acm. T. Graphic* **25**:924–934.
21. R. Oldenbourg, (2008) "Polarized Light Field Microscopy: An Analytical Method Using a Microlens Array to Simultaneously Capture both Conoscopic and Orthoscopic Views of Birefringent Objects," *J. Microsc.* **231**:419–432.
22. N. H. Hartshorne and A. Stuart, (1964) *Practical Optical Crystallography*, American Elsevier Publishing Co., Inc., New York, NY.
23. H. E. Keller, (1995) Objective Lenses for Confocal Microscopy. In: *Handbook of Biological Confocal Microscopy* (ed. J. B. Pawley). 2d ed. Plenum Publ. Corp., New York.
24. Y. Shimizu and H. Takenaka, (1994) "Microscope Objective Design," *Adv. Opt. Electron Microsc.* **14**:249–334.
25. R. Kingslake, (1978) *Lens Design Fundamentals*, Academic Press, San Diego.

26. E. Hecht, (2002) *Optics*, Pearson/Addison-Wesley, San Francisco, CA.
27. L. W. Smith and H. Osterberg, (1961) "Diffraction Images of Circular Self-Radiant Disks," *J. Opt. Soc. Am.* **51**:412–414.
28. M. Cagnet, M. Françon and J. C. Thierri, (1962) *Atlas of Optical Phenomena*, Springer Verlag, Berlin.
29. E. H. Linfoot and E. Wolf (1956) "Phase Distribution Near Focus in an Aberration-Free Diffraction Image," *Proc. Phys. Soc.* **B69**:823–832.
30. S. Inoué, (1989) Imaging of Unresolved Objects, Superresolution and Precision of Distance Measurement with Video Microscopy. In: *Methods Cell Biol.* (eds. D. L. Taylor and Y.-L. Wang). Academic Press, New York.
31. R. E. Thompson, D. R. Larson and W. W. Webb, (2002) "Precise Nanometer Localization Analysis for Individual Fluorescent Probes," *Biophys. J.* **82**:2775–2783.
32. R. J. Ober, S. Ram and E. S. Ward, (2004) "Localization Accuracy in Single-Molecule Microscopy," *Biophys. J.* **86**:1185–1200.
33. Y. Garini, B. J. Vermolen and I. T. Young, (2005) "From Micro to Nano: Recent Advances in High-Resolution Microscopy," *Curr. Opin. Biotechnol.* **16**:3–12.
34. S. W. Hell, (2007) "Far-Field Optical Nanoscopy," *Science* **316**:1153–1158.
35. H. Park, E. Toprak and P. R. Selvin, (2007) "Single-Molecule Fluorescence to Study Molecular Motors," *Q. Rev. Biophys.* **40**:87–111.
36. R. Heintzmann and G. Ficz, (2007) "Breaking the Resolution Limit in Light Microscopy," *Methods Cell Biol.* **81**:561–580.
37. A. Yildiz, J. N. Forkey, S. A. McKinney, T. Ha, Y. E. Goldman and P. R. Selvin, (2003) "Myosin V Walks Hand-Over-Hand: Single Fluorophore Imaging with 1.5-nm Localization," *Science* **300**:2061–2065.
38. E. Betzig, G. H. Patterson, R. Sougrat, O. W. Lindwasser, S. Olenych, J. S. Bonifacino, M. W. Davidson, J. Lippincott-Schwartz and H. F. Hess, (2006) "Imaging Intracellular Fluorescent Proteins at Nanometer Resolution," *Science* **313**:1642–1645.
39. M. J. Rust, M. Bates and X. Zhuang, (2006) "Sub-Diffraction-Limit Imaging by Stochastic Optical Reconstruction Microscopy (STORM)," *Nat. Methods* **3**:793–795.
40. L. S. Churchman, Z. Okten, R. S. Rock, J. F. Dawson and J. A. Spudich, (2005) "Single Molecule High-Resolution Colocalization of Cy3 and Cy5 Attached to Macromolecules Measures Intramolecular Distances through Time," *Proc. Natl. Acad. Sci. U S A*, **102**:1419–1423.
41. R. Roy, S. Hohng and T. Ha, (2008) "A Practical Guide to Single-Molecule FRET," *Nat. Methods* **5**:507–516.
42. D. W. Piston and G. J. Kremers, (2007) "Fluorescent Protein FRET: The Good, the Bad and the Ugly," *Trends Biochem. Sci.* **32**:407–414.
43. S. W. Hell and J. Wichmann, (1994) "Breaking the Diffraction Resolution Limit by Stimulated-Emission—Stimulated-Emission-Depletion Fluorescence Microscopy," *Opt. Lett.* **19**:780–782.
44. J. L. Harris, (1964) "Diffraction and Resolving Power," *J. Opt. Soc. Am.* **54**:931–936.
45. E. Betzig, A. Lewis, A. Harootunian, M. Isaacson and E. Kratschmer, (1986) "Near-Field Scanning Optical Microscopy (NSOM); Development and Biophysical Applications," *Biophys. J.* **49**:269–279.
46. H. K. Wickramasinghe, (1989) "Scanned-Probe Microscopes," *Sci. Am.* **261**:98–105.
47. I. T. Young, (1989) Image Fidelity: Characterizing the Imaging Transfer Function. In: *Methods Cell Biol.* (eds. D. L. Taylor and Y.-L. Wang). Academic Press, New York.
48. C. J. Sheppard, (2004) "Defocused Transfer Function for a Partially Coherent Microscope and Application to Phase Retrieval," *J. Opt. Soc. Am. A. Opt. Image Sci. Vis.* **21**:828–831.
49. D. A. Agard, Y. Hiraoka, P. Shaw and J. W. Sedat, (1989) "Fluorescence Microscopy in Three Dimensions," *Methods Cell Biol.* **30**:353–377.
50. R. Oldenbourg, S. Inoué, R. Tiberio, A. Stemmer, G. Mei and M. Skvarla, (1996) Standard Test Targets for High Resolution Light Microscopy. In: *Nanofabrication and Biosystems: Integrating Material Science, Engineering and Biology* (eds. H. C. Hoch, L. W. Jelinsky and H. Craighead). Cambridge University Press, Cambridge, England.
51. E. M. Chamot and C. W. Mason, (1958) *Handbook of Chemical Microscopy*, John Wiley & Sons, New York.
52. B. Kachar, (1985) "Asymmetric Illumination Contrast: A Method of Image Formation for Video Microscopy," *Science* **227**:766–768.
53. N. Streibl, (1984) "Phase Imaging by the Transport-Equation of Intensity," *Opt. Commun.* **49**:6–10.

54. C. J. R. Sheppard, (2002) "Three-Dimensional Phase Imaging with the Intensity Transport Equation," *Appl. Opt.* **41**:5951–5955.
55. E. D. Barone-Nugent, A. Barty and K. A. Nugent, (2002) "Quantitative Phase-Amplitude Microscopy I: Optical Microscopy," *J. Microsc.* **206**:194–203.
56. D. Paganin, A. Barty, P. J. McMahon and K. A. Nugent, (2004) "Quantitative Phase-Amplitude Microscopy. III. The Effects of Noise," *J. Microsc.* **214**:51–61.
57. K. A. Nugent, D. Paganin and A. Barty, (2006) Phase Determination of a Radiation Wave Field. In: *U.S. Patent, Number 7,039,553*. University of Melbourne, Australia, USA.
58. C. S. Izzard and L. R. Lochner, (1980) "Formation of Cell-to-Substrate Contacts During Fibroblast Motility: An Interference-Reflection Study," *J. Cell Sci.* **42**:81–116.
59. N. Noda and S. Kamimura, (2008) "A New Microscope Optics for Laser Dark-Field Illumination Applied to High Precision Two Dimensional Measurement of Specimen Displacement," *Rev. Sci. Instrum.* **79**:023704.
60. F. Zernike, (1958) The Wave Theory of Microscopic Image Formation. In: *Concepts of Classical Optics* (ed. J. Strong). W.H. Freeman and Co., San Francisco.
61. A. H. Bennett, H. Osterberg, H. Jupnik and O. W. Richards, (1951) *Phase Microscopy, Principles and Applications*, John Wiley & Sons, Inc., New York.
62. M. Pluta, (1989b) "Simplified Polanret System for Microscopy," *Appl. Opt.* **28**:1453–1466.
63. R. Hoffman and L. Gross, (1975) "Modulation Contrast Microscopy," *Appl. Opt.* **14**:1169–1176.
64. G. W. Ellis, (1978) Advances in Visualization of Mitosis In Vivo. In: *Cell Reproduction: In Honor of Daniel Mazia*, (eds. E. Dirksen, D. Prescott and C. F. Fox). Academic Press, New York, pp. 465–476.
65. G. W. Ellis, (1988) Scanned Aperture Light Microscopy," *Proc. 46th Annual Meeting of the Electron Microscopy Society of America*, 48–49.
66. J. Dyson, (1961) "Instrument Classification and Applications," In: *The Encyclopedia of Microscopy* (ed. G.L. Clark) Reinhold Pub. Corp., New York, 412–420
67. C. J. Koester, (1961) "Theory and Techniques," In: *The Encyclopedia of Microscopy* (ed. G.L. Clark) Reinhold Pub. Corp., New York, 420–434
68. N. H. Hartshorne and A. Stuart, (1960) *Crystals and the Polarising Microscope: A Handbook for Chemists and Others*, Arnold, London.
69. G. A. Dunn, (1998) "Transmitted-Light Interference Microscopy: A Technique Born before Its Time," *Proc. RMS* **33**:189–196.
70. W. P. Linnik, (1933) "Apparatus for Interferometric Study of Reflective Specimens with a Microscope ("Microinterferometer")," *C. R. Acad. Sci. USSR* **1**:18–23 (in Russian and German).
71. A. A. Lebedeff, (1930) "Polarization Interferometer and Its Applications," *Rev. Opt.* **9**:385.
72. G. V. Rozenberg, (1953) "Interference Microscopy," *Uspekhi Fizicheskikh Nauk* **50**:271–302 (in Russian).
73. M. Françon, (1964) "Polarization Interference Microscopes," *Appl. Opt.* **3**:1033–1036.
74. M. J. Jamin, (1868) "Sur un réfracteur différentiel pour la lumière polarisée," *C. R. Acad. Sci. (Paris)* **67**:814.
75. F. H. Smith, (1952) Interference Microscope. In: *U.S. Patent, Number 2,601,175*. Francis Hugh Smith, USA.
76. F. H. Smith, (1955) "Microscopic Interferometry," *Research (London)*, **8**:385–395.
77. G. Nomarski, (1960) Interferential Polarizing Device for Study of Phase Objects. In: *U.S. Patent, Number 2,924,142*. Centre National de la Recherche Scientifique, Paris, France, USA.
78. R. D. Allen, G. B. David and G. Nomarski, (1969) "The Zeiss-Nomarski Differential Interference Equipment for Transmitted Light Microscopy," *Zeitschrift für wissenschaftliche Mikroskopie* **69**:193–221.
79. M. Shribak and S. Inoue, (2006) "Orientation-Independent Differential Interference Contrast Microscopy," *Appl. Opt.* **45**:460–469.
80. C. J. Cogswell and C. J. R. Sheppard, (1992) "Confocal Differential Interference Contrast (DIC) Microscopy: Including a Theoretical Analysis of Conventional and Confocal DIC Imaging," *J. Microscopy* **165**:81–101.
81. C. Preza, D. L. Snyder and J. A. Conchello, (1999) "Theoretical Development and Experimental Evaluation of Imaging Models for Differential-Interference-Contrast Microscopy," *J. Opt. Soc. Am. A. Opt. Image Sci. Vis.* **16**:185–2199.

82. R. D. Allen, J. L. Travis, N. S. Allen and H. Yilmaz, (1981) "Video-Enhanced Contrast Polarization (AVEC-POL) Microscopy: A New Method Applied to the Detection of Birefringence in the Motile Reticulopodial Network *Allogromia laticollaris*," *Cell Motil.* **1**:275–289.
83. E. D. Salmon and P. T. Tran, (1998) "High-Resolution Video-Enhanced Differential Interference Contrast (VE-DIC) Light Microscopy," *Methods Cell Biol.*, **56**:153–184.
84. H. Ishiwata, M. Itoh and T. Yatagai, (1996) "Retardation-Modulated Differential Interference Microscope and its Application to 3D Shape Measurement," *Proc.* **2873**:21–24.
85. G. Holzwarth, S. C. Webb, D. J. Kubinski and N. S. Allen, (1997) "Improving DIC Microscopy with Polarization Modulation," *J. Microsc.* **188**:249–254.
86. G. M. Holzwarth, D. B. Hill and E. B. McLaughlin, (2000) "Polarization-Modulated Differential-Interference Contrast Microscopy with a Variable Retarder," *Appl. Opt.* **39**:6288–6294.
87. H. Ishiwata, M. Itoh and T. Yatagai, (2006) "A New Method of Three-Dimensional Measurement by Differential Interference Contrast Microscope," *Opt. Commun.* **260**:117–126.
88. R. P. Danz, A. Dietrich, C. Soell, Hoyer and M. Wagener, (2006) Arrangement and Method for Polarization-optical Interference Contrast. In: *U.S. Patent, Number 7,046,436*. Carl Zeiss Jena GmbH (Jena, DE).
89. B. Heise, A. Sonnleitner and E. P. Klement, (2005) "DIC Image Reconstruction on Large Cell Scans," *Microsc. Res. Tech.* **66**:312–320.
90. C. Preza, (2000) "Rotational-Diversity Phase Estimation from Differential-Interference—Contrast Microscopy Images," *J. Opt. Soc. Am. A-Opt. Image Sci. Vis.* **17**:415–424.
91. M. R. Arnison, K. G. Larkin, C. J. R. Sheppard, N. I. Smith and C. J. Cogswell, (2004) "Linear Phase Imaging Using Differential Interference Contrast Microscopy," *J. Microsc.-Oxford*, **214**:7–12.
92. F. Kagalwala and T. Kanade, (2003) "Reconstructing Specimens Using DIC Microscope Images," *IEEE Transactions on Systems Man and Cybernetics Part B-Cybernetics* **33**:728–737.
93. N. Axelrod, A. Radko, A. Lewis and N. Ben-Yosef, (2004) "Topographic Profiling and Refractive-Index Analysis by Use of Differential Interference Contrast with Bright-Field Intensity and Atomic Force Imaging," *Appl. Opt.* **43**:2272–2284.
94. M. Shribak, J. LaFountain, D. Biggs and S. Inoue, (2008) "Orientation-Independent Differential Interference Contrast Microscopy and Its Combination with an Orientation-Independent Polarization System," *J Biomed. Opt.* **13**:014011.
95. J. Dyson, (1950) "An Interferometer Microscope," *Proc. Roy. Soc. London. Series A, Mathematical and Physical Sciences* **204**:170–187.
96. G. Delaunay, (1953) "Microscope interférentiel A. Mirau pour la mesure du fini des surfaces," *Rev. Optique* **32**:610–614.
97. G. W. Ellis, (1966) "Holomicrography: Transformation of Image During Reconstruction of a Posteriori," *Science* **154**:1195–1197.
98. P. Marquet, B. Rappaz, P. J. Magistretti, E. Cuche, Y. Emery, T. Colomb and C. Depeursing, (2005) "Digital Holographic Microscopy: A Noninvasive Contrast Imaging Technique Allowing Quantitative Visualization of Living Cells with Subwavelength Axial Accuracy," *Opt. Lett.* **30**:468–470.
99. B. Kemper, D. Carl, J. Schneckeburger, I. Bredebusch, M. Schafer, W. Domschke and G. von Bally, (2006) "Investigation of Living Pancreas Tumor Cells by Digital Holographic Microscopy," *J. Biomed. Opt.* **11**.
100. A. Stern and B. Lavid, (2007) "Theoretical Analysis of Three-Dimensional Imaging and Recognition of Micro-Organisms with a Single-Exposure On-line Holographic Microscope," *J. Opt. Soc. Am. A-Opt. Image Sci. Vis.* **24**:163–168.
101. A. Stadelmaier and J. H. Massig, (2000) "Compensation of Lens Aberrations in Digital Holography," *Opt. Lett.* **25**:1630–1632.
102. S. S. Kou and C. J. R. Sheppard, (2007) "Imaging in Digital Holographic Microscopy," *Opt. Exp.* **15**:13640–13648.
103. Y. Kuznetsova, A. Neumann and S. R. J. Brueck, (2007) "Imaging Interferometric Microscopy—Approaching the Linear Systems Limits of Optical Resolution," *Opt. Exp.* **15**:6651–6663.
104. C. J. R. Sheppard and M. Roy, (2007) Low-Coherence Interference Microscopy. In: *Optical Imaging and Microscopy* (eds. P. Török and F. J. Kao). 2d ed. Springer-Verlag, Berlin, Heidelberg, New York.

105. A. M. Zysk and S. A. Boppart, (2007) Optical Coherence Tomography. In: *Optical Imaging and Microscopy* (eds. P. Török and F. J. Kao). 2d ed. Springer-Verlag, Berlin, Heidelberg, New York.
106. M. Wojtkowski, R. Leitgeb, A. Kowalczyk, T. Bajraszewski and A. F. Fercher, (2002) "In Vivo Human Retinal Imaging by Fourier Domain Optical Coherence Tomography," *J. Biomed. Opt.* **7**:457–463.
107. S. R. Chinn, E. A. Swanson and J. G. Fujimoto, (1997) "Optical Coherence Tomography Using a Frequency-Tunable Optical Source," *Opt. Lett.* **22**:340–342.
108. F. Lexer, C. K. Hitznerberger, A. F. Fercher and M. Kulhavy, (1997) "Wavelength-Tuning Interferometry of Intraocular Distances," *Appl. Opt.* **36**:6548–6553.
109. V. V. Tuchin, L. V. Wang and D. A. Zimnyakov, (2006) *Optical Polarization in Biomedical Applications*, Springer-Verlag, Berlin, Heidelberg, New York.
110. J. A. Izatt, M. R. Hee, G. M. Owen, E. A. Swanson and J. G. Fujimoto, (1994) "Optical Coherence Microscopy in Scattering Media," *Opt. Lett.* **19**:590–592.
111. B. M. Hoeling, A. D. Fernandez, R. C. Haskell, E. Huang, W. R. Myers, D. C. Petersen, S. E. Ungersma, R. Y. Wang, M. E. Williams and S. E. Fraser, (2000) "An Optical Coherence Microscope for 3-Dimensional Imaging in Developmental Biology," *Opt. Exp.* **6**:136–146.
112. S. Inoué, (1952) "Studies on Depolarization of Light at Microscope Lens Surfaces. I. The Origin of Stray Light by Rotation at the Lens Surfaces," *Exp. Cell Res.* **3**:199–208.
113. S. Inoué and H. Kubota, (1958) "Diffraction Anomaly in Polarizing Microscopes," *Nature* **182**:1725–1726.
114. H. Kubota and S. Inoué, (1959) "Diffraction Images in the Polarizing Microscope," *J. Opt. Soc. Am.* **49**:191–198.
115. S. Inoué and W. L. Hyde, (1957) "Studies on Depolarization of Light at Microscope Lens Surfaces II. The Simultaneous Realization of High Resolution and High Sensitivity with the Polarizing Microscope," *J. Biophys. Biochem. Cytol.* **3**:831–838.
116. M. Shribak, S. Inoué and R. Oldenbourg, (2002) "Polarization Aberrations Caused by Differential Transmission and Phase Shift in High NA Lenses: Theory, Measurement and Rectification," *Opt. Eng.* **41**:943–954.
117. S. Inoué and H. Sato, (1966) Deoxyribonucleic acid Arrangement in Living Sperm. In: *Molecular Architecture in Cell Physiology* (eds. T. Hayashi and A. G. Szent-Gyorgyi). Prentice Hall, Englewood Cliffs, NJ.
118. H. G. Jerrard, (1948) "Optical Compensators for Measurement of Elliptical Polarization," *J. Opt. Soc. Am.* **38**:35–59.
119. R. Oldenbourg, (1999) "Polarized Light Microscopy of Spindles," *Methods Cell Biol.* **61**:175–208.
120. H. G. Jerrard, (1954) Transmission of Light through Birefringent and Optically Active Media: The Poincaré sphere," *J. Opt. Soc. Am.* **44**:634–640.
121. A. M. Glazer, J. G. Lewis and W. Kaminsky, (1996) "An Automatic Optical Imaging System for Birefringent Media," *Proc. R. Soc. London A* **452**:2751–2765.
122. R. D. Allen, J. Brault and R. D. Moore, (1963) "A New Method of Polarization Microscopic Analysis I. Scanning with a Birefringence Detection System," *J. Cell Biol.* **18**:223–235.
123. J. R. Kuhn, Z. Wu and M. Poenie, (2001) "Modulated Polarization Microscopy: A Promising New Approach to Visualizing Cytoskeletal Dynamics in Living Cells," *Biophys. J.* **80**:972–985.
124. R. Oldenbourg and G. Mei, (1995) "New Polarized Light Microscope with Precision Universal Compensator," *J. Microsc.* **180**:140–147.
125. M. Shribak and R. Oldenbourg, (2003a) "Techniques for Fast and Sensitive Measurements of Two-Dimensional Birefringence Distributions," *Appl. Opt.* **42**:3009–3017.
126. R. Oldenbourg, (2005) Polarization Microscopy with the LC-PolScope. In: *Live Cell Imaging: A Laboratory Manual* (eds. R. D. Goldman and D. L. Spector). Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
127. R. Oldenbourg, (1991) "Analysis of Edge Birefringence," *Biophys. J.* **60**:629–641.
128. F. D. Bloss, (1981) *The Spindle Stage Principles and Practice*, Cambridge University Press, Cambridge.
129. W. Kaminsky and A. M. Glazer, (1996) "Measurement of Optical Rotation in Crystals," *Ferroelectrics* **183**:133–141.
130. L. A. Pajdzik and A. M. Glazer, (2006) "Three-Dimensional Birefringence Imaging with a Microscope Tilting-Stage. I. Uniaxial Crystals," *J. Appl. Crystallogr.* **39**:326–337.

131. M. Shribak and R. Oldenbourg, (2004) "Mapping Polymer Birefringence in Three Dimensions Using a Polarizing Microscope with Oblique Illumination," *SPIE Proc. (Proc. Biophotonics Micro- and Nano-Imaging)* **5462**:57–67.
132. D. L. Taylor, M. Nederlof, F. Lanni and A. S. Waggoner, (1992) "The New Vision of Light Microscopy," *Am. Sci.* **80**:322–335.
133. R. Y. Tsien, (1989) Fluorescent Probes of Cell Signaling. In: *Annu. Rev. Neurosci.* (eds. W. M. Cowas, E. M. Shooter, C. F. Stevens, and R. F. Thompson). Annual Reviews, Inc., Palo Alto.
134. R. Y. Tsien, (1998) "The Green Fluorescent Protein," *Annu. Rev. Biochem.* **67**:509–544.
135. J. W. Lichtman and J. A. Conchello, (2005) "Fluorescence Microscopy," *Nat. Methods* **2**:910–919.
136. R. Neher and E. Neher, (2004) "Optimizing Imaging Parameters for the Separation of Multiple Labels in a Fluorescence Image," *J. Microsc.* **213**:46–62.
137. Y. Garini, I. T. Young and G. McNamara, (2006) Spectral Imaging: Principles and Applications," *Cytometry Part A* **69**:735–747.
138. M. Minsky, (1957) Microscopy Apparatus. In: *U.S. Patent #3013467* (ed. U. S. P. a. T. Office). USA.
139. S. Inoué, (1990) Foundations of Confocal Scanned Imaging in Light Microscopy. In: *Handbook of Biological Confocal Microscopy* (ed. J. B. Pawley). Plenum Publ. Corp., New York.
140. G. S. Kino, (1990) Intermediate Optics in Nipkow Disk Microscopes. In: *Handbook of Biological Confocal Microscopy* (ed. J. B. Pawley). Plenum Press, New York.
141. M. Petráň, M. Hadravsky, D. Egger and R. Galambos, (1968) "Tandem-Scanning Reflected-Light Microscope," *J. Opt. Soc. Am.* **58**:661–664.
142. S. Inoue and T. Inoue, (2002) "Direct-View High-Speed Confocal Scanner: The CSU-10," *Methods Cell Biol.* **70**:87–127.
143. E. Wang, C. M. Babbey and K. W. Dunn, (2005) "Performance Comparison between the High-Speed Yokogawa Spinning Disc Confocal System and Single-Point Scanning Confocal Systems," *J. Microsc.* **218**:148–159.
144. W. A. Carrington, R. M. Lynch, E. D. W. Moore, G. Isenberg, K. E. Fogarty and F. S. Fay, (1995) "Superresolution Three-Dimensional Images of Fluorescence in Cells with Minimal Light Exposure," *Science* **268**:1483–1487.
145. J. A. Conchello and J. W. Lichtman, (2005) "Optical Sectioning Microscopy," *Nat. Methods* **2**:920–931.
146. M. A. Neil, R. Juskaitis and T. Wilson, (1997) "Method of Obtaining Optical Sectioning by Using Structured Light in a Conventional Microscope," *Opt. Lett.* **22**:1905–1907.
147. C. Ventalon and J. Mertz, (2005) "Quasi-Confocal Fluorescence Sectioning with Dynamic Speckle Illumination," *Opt. Lett.* **30**:3350–3352.
148. C. Ventalon and J. Mertz, (2006) "Dynamic Speckle Illumination Microscopy with Translated versus Randomized Speckle Patterns," *Opt. Exp.* **14**:7198–7209.
149. M. G. Gustafsson, (2000) "Surpassing the Lateral Resolution Limit by a Factor of Two Using Structured Illumination Microscopy," *J. Microsc.* **198**:82–87.
150. M. G. Gustafsson, L. Shao, P. M. Carlton, C. J. Wang, I. N. Golubovskaya, W. Z. Cande, D. A. Agard and J. W. Sedat, (2008) "Three-Dimensional Resolution Doubling in Wide-Field Fluorescence Microscopy by Structured Illumination," *Biophys. J.* **94**:4957–4970.
151. R. Heintzmann, T. M. Jovin and C. Cremer, (2002) "Saturated Patterned Excitation Microscopy—a Concept for Optical Resolution Improvement," *J. Opt. Soc. Am. A-Opt. Image Sci. Vis.* **19**:1599–1609.
152. M. G. Gustafsson, (2005) "Nonlinear Structured-Illumination Microscopy: Wide-Field Fluorescence Imaging with Theoretically Unlimited Resolution," *Proc. Natl. Acad. Sci. U S A* **102**:13081–13086.
153. R. Ng, M. Levoy, M. Bredif, G. Duval, M. Horowitz and P. Hanrahan, (2005) Light Field Photography with a Hand-Held Plenoptic Camera. *Stanford Tech. Report, CTSR 2005-02*.
154. M. Shribak and R. Oldenbourg, (2003b) "Three-Dimensional Birefringence Distribution in Reconstituted Asters of *Spisula* Oocytes Revealed by Scanned Aperture Polarized Light Microscopy," *Biol. Bull.* **205**:194–195.
155. R. E. Stephens, (1965) "Analysis of Muscle Contraction by Ultraviolet Microbeam Disruption of Sarcomere Structure," *J. Cell Biol.* **25**:129–139.
156. M. W. Berns, (1974) *Biological Microirradiation, Classical and Laser Sources*, Prentice-Hall, Englewood Cliffs.

157. D. E. Koppel, D. Axelrod, J. Schlessinger, E. L. Elson and W. W. Webb, (1976) "Dynamics of Fluorescence Marker Concentration as a Probe of Mobility," *Biophys. J.* **16**:1315–1329.
158. H. G. Kapitza, G. McGregor and K. A. Jacobson, (1985) "Direct Measurement of Lateral Transport in Membranes by Using Time-Resolved Spatial Photometry," *Proc. Natl. Acad. Sci. (USA)* **82**:4122–4126.
159. M. W. Berns, J. Aist, J. Edwards, K. Strahs, J. Girton, P. McNeill, J. B. Rattner, et al., (1981) "Laser Microsurgery in Cell and Developmental Biology," *Science* **213**:505–513.
160. V. Magidson, J. Loncarek, P. Hergert, C. L. Rieder and A. Khodjakov, (2007) Laser Microsurgery in the GFP Era: A Cell Biologist's Perspective," *Methods Cell Biol.* **82**:239–266.
161. C. Lutz, T. S. Otis, V. Desars, S. Charpak, D. A. Digregorio and V. Emiliani, (2008) "Holographic Photolysis of Caged Neurotransmitters," *Nat. Methods* **5**:821–827.
162. J. H. Kaplan and G. C. Ellis-Davies, (1988) "Photolabile Chelators for the Rapid Photorelease of Divalent Cations," *Proc. Natl. Acad. Sci. (USA)* **85**:6571–6575.
163. J. A. Dantzig, M. G. Hibberd, D. R. Trentham and Y. E. Goldman, (1991) Cross-Bridge Kinetics in the Presence of MgADP Investigated by Photolysis of Caged ATP in Rabbit Psoas Muscle Fibres," *J. Physiol. (London)* **432**:639–680.
164. T. J. Mitchison and E. D. Salmon, (1992) "Poleward Kinetochore Fiber Movement Occurs during Both Metaphase and Anaphase—A in Newt Lung Cell Mitosis," *J. Cell Biol.* **119**:569–582.
165. A. Ashkin, J. M. Dziedzic, J. E. Bjorkholm and S. Chu, (1986) "Observation of a Single-Beam Gradient Force Optical Trap for Dielectric Particles," *Opt. Lett.* **11**:288–290.
166. K. Svoboda and S. M. Block, (1994) "Biological Applications of Optical Forces," *Annu. Rev. Biophys. Biomol. Struct.* **23**:247–285.
167. S. M. Block, (1990) "Optical Tweezers: A New Tool for Biophysics," In: *Noninvasive Techniques in Cell Biology* (eds. J. K. Foskett and S. Grinstein). Wiley-Liss, New York.
168. J. R. Moffitt, Y. R. Chemla, S. B. Smith and C. Bustamante, (2008) "Recent Advances in Optical Tweezers," *Annu. Rev. Biochem.* **77**:205–228.

REFLECTIVE AND CATADIOPTRIC OBJECTIVES

Lloyd Jones

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

29.1 GLOSSARY

A	4th-order aspheric deformation coefficient
AN	4th-order nonsymmetric deformation coefficient
B	6th-order aspheric deformation coefficient
C	8th-order aspheric deformation coefficient
c	surface base curvature
CON	conic constant
D	10th-order aspheric deformation coefficient
FN	focal ratio
GLA	glass type
h	radial surface height
INF	infinite radius of curvature
k	conic constant
n	index of refraction
R	radius of curvature
RDX	radius of curvature in the x dimension
RDY	radius of curvature in the y dimension
STO	stop surface
SUR	surface number
t	element thickness
THI	thickness of element or distance to next surface or element
Z	surface sag

29.2 INTRODUCTION

During the initial stages of an optical design, many optical engineers take advantage of existing configurations that exhibit useful properties. This chapter is a compilation of reflective and *catadioptric* objective designs that should help inform the reader of available choices and provide reasonable starting solutions.

The chapter also includes a cursory introduction to some of the more important topics in system analysis, such as angular and linear blur size, *image irradiance*, scaling, and stray light control.

An extensive list of referenced reading material and brief definitions of terms italicized throughout the text are included.

29.3 GLASS VARIETIES

Glasses used in the designs are represented in terms of index of refraction and Abbe number or V number, below. The V number indicates glass dispersion. Most glasses can be obtained from a number of vendors.

Glass	Index of Refraction	V Number
BK7	1.516	64.2
F2	1.620	36.3
F9	1.620	38.1
FK51	1.487	84.5
FN11	1.621	36.2
Germanium	4.037	117.4
LLF1	1.548	45.8
LAK21	1.640	60.1
PSK2	1.569	63.2
Silica	1.445	27.7
Silicon	3.434	147.4
Sapphire	1.735	15.5
SK1	1.610	56.5
SK2	1.607	56.8
SK3	1.609	58.9
SK16	1.620	60.3
SF5	1.673	32.1
SF10	1.728	28.5
UBK7	1.517	64.3

29.4 INTRODUCTION TO CATADIOPTRIC AND REFLECTIVE OBJECTIVES

The variety of objectives presented in this chapter is large. Most of the intricate detail relating to each design is therefore presented with the design itself. In the following paragraphs, analysis of the general features of the catadioptric and reflective objectives is undertaken.

Conic Mirrors

It is apparent after a brief perusal of the designs that there are many surface types. Among these are the sphere, paraboloid, hyperboloid, prolate ellipsoid, and oblate ellipsoid. The oblate

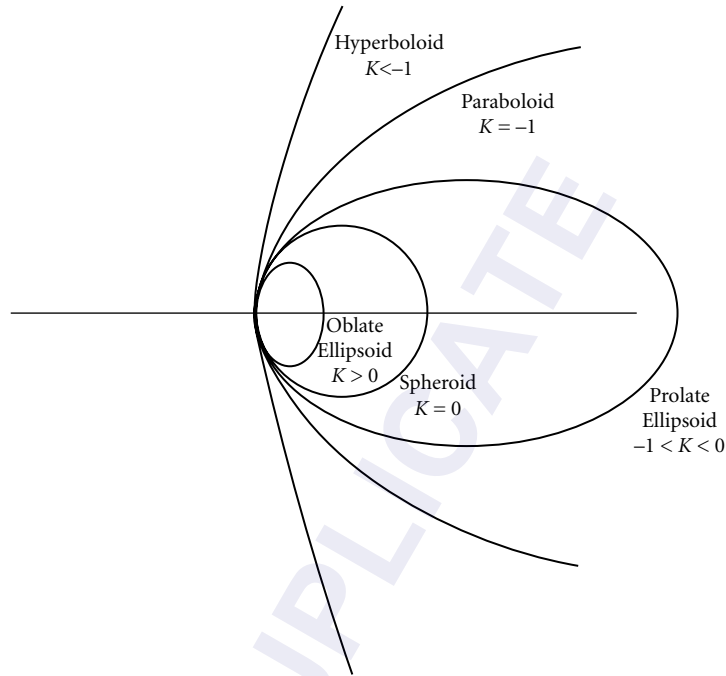


FIGURE 1 Relative shapes of conic surfaces in two dimensions.

ellipsoid is a prolate ellipsoid turned on its side. The equation of a conic is given by the expression

$$Z = \frac{ch^2}{1 + \sqrt{1 - (1+k)c^2h^2}} \quad (1)$$

where Z is the surface sag, k is the *conic constant*, c is the surface base curvature, and h is the radial height on the surface. The relative shapes of these surfaces are illustrated in Fig. 1.

Conic mirrors give perfect geometric imagery when an axial point object is located at one conic focus and the axial point image is located at the other conic focus. Figure 2 illustrates these ray paths.

General Aspheres

General aspheres are surfaces with fourth- and higher-order surface deformation on top of a flat or curved surface.¹ The surface deformation of a rotationally symmetric general asphere is given by the relation

$$Z = \frac{ch^2}{1 + \sqrt{1 - (1+k)c^2h^2}} + Ah^4 + Bh^6 + Ch^8 + Dh^{10} \quad (2)$$

where A , B , C , and D are 4th-, 6th-, 8th-, and 10th-order coefficients that determine the sign and magnitude of the deformation produced by that order. Although general aspheres allow correction of *third-* and *higher-order aberrations* and may reduce the number of elements in an optical system, general aspheres are more expensive than spheres or conics. If aspheric deformation is required, conic surfaces should be tried first, especially since a conic offers higher-order correction.²

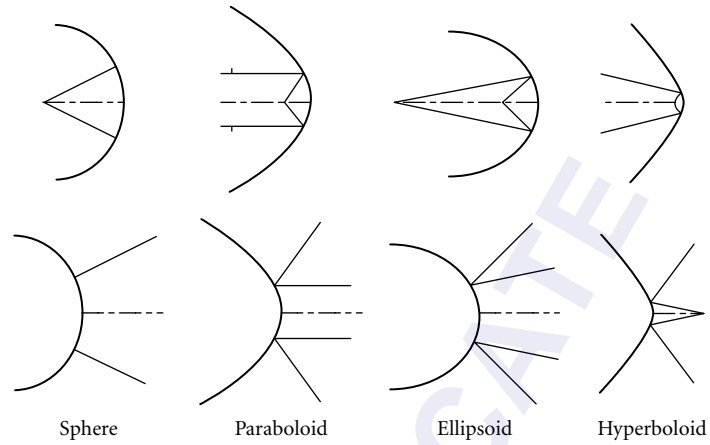


FIGURE 2 Ray paths for perfect axial imagery.

Obscurations

Obscurations that block portions of the entering beam reduce image irradiance and image contrast^{3,4} in reflective and catadioptric systems. Several methods are used to reduce or eliminate completely the effects of an obscuration (see Fig. 3).

Figure 3a illustrates a commonly employed technique for reducing large-mirror obscuration: a small secondary mirror close to the intermediate image moves the larger tertiary mirror out of the beam path.

Figure 3b is an illustration of an eccentric pupil system. All elements are symmetric about the same axis and the aperture stop is decentered for a clear light path.

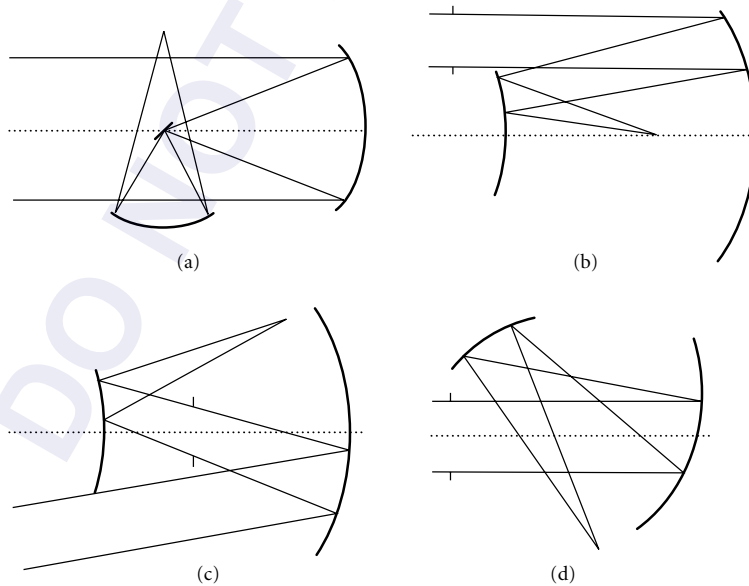


FIGURE 3 Reducing the size of or eliminating an obscuration.

Figure 3c is an example of an off-axis objective with the field of view biased to direct the center of view away from any intervening elements. All elements and apertures are symmetric about the *optical axis*.

Figure 3d is an illustration of a tilted and decentered-component objective. Each element is rotationally symmetric about its own unique optical axis which may be tilted and/or decentered. The imaging behavior of this system is more complicated to deal with than the imaging behavior of eccentric pupil and off-axis systems. Vector aberration theory^{5,6} has been developed to properly model the imaging behavior of these systems.

Stray Light Suppression

Suppression of light diffracted from apertures and obscurations is facilitated with intermediate images and a real and accessible *Lyot stop*. Figure 4a illustrates a generic refractive configuration with an intermediate image and Lyot stop. Figure 4b illustrates where the diffracted light (shaded region) originates and terminates (at one edge of each aperture, for clarity).

A *field stop* is placed at the focus of the first lens to block diffracted light produced by the front light baffle. To block unwanted objects within the field of view, an occulting disc may be inserted at the focus of the first lens, as is done with a Lyot coronagraph in order to block the sun. By oversizing the field stop slightly, the light diffracted at the field stop falls just outside of the detector area.

Following the field stop is a second lens that reimages the intermediate image to the final image and the *entrance pupil* to the Lyot stop (the shaded region in Fig. 4a illustrates how the entrance pupil is imaged). Undersizing the Lyot stop blocks the light diffracted at the entrance pupil. In this way the Lyot stop becomes the *aperture stop* of the system.

Another application of the Lyot stop in the infrared (assuming the Lyot stop is located exterior to the objective optics) is as a cold stop.⁷ The cold stop (Fig. 4a) is a baffle that prevents stray infrared light, radiated from the housing, from impinging upon the detector from outside its intended field.

Reflective and Catadioptric Objective Designs

The objectives to follow are listed according to *focal ratio* and design type. Objectives have a 20-cm diameter and catadioptric systems are optimized for a wavelength range from 480 to

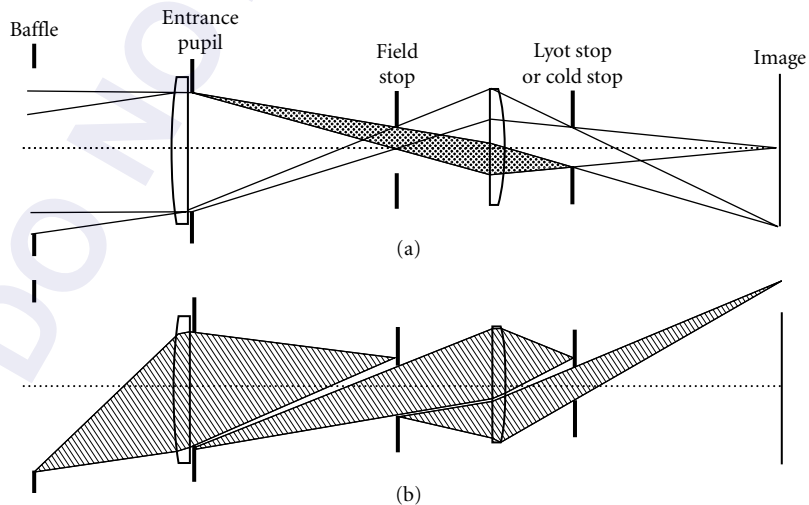


FIGURE 4 Generic objectives with apertures.

680 nm unless otherwise stated. The angles of surface and element tilts are with respect to the horizontal optical axis. Decenters are also with respect to the optical axis. Since many of the designs are *aplanatic*, *anastigmatic*, and free of chromatic aberrations, the position of the stop does not affect third-order aberration correction and may be chosen to minimize *vignetting*, *distortion*, element size, or stray light contamination. All aberrations mentioned in this section are third-order unless otherwise stated.

Definitions of the abbreviated terminology used in the lens data are as follows:

SUR Surface number.

RDY Surface radius of curvature. A positive value means the center of curvature lies to the right of the surface; negative to the left.

THI Thickness of element or distance to next element. The value is positive if the next surface lies to the right of the surface.

GLA Glass-type or mirror surface, the latter referred to by the abbreviation REFL.

CON Conic constant k .

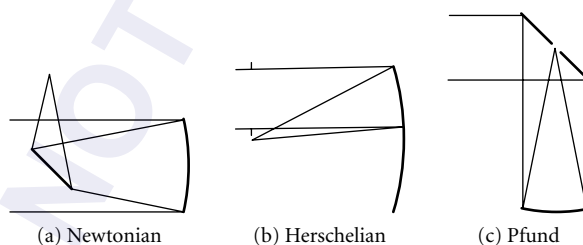
STO Stop surface.

INF A surface with an infinite radius of curvature; that is, a flat surface.

A, B, C, D The 4th-, 6th-, 8th-, and 10th-order aspheric deformation coefficients in Eq. (2).

A potential source of confusion is the terminology used to describe Mangin elements; that is, refractive elements with reflective back surfaces. This is illustrated in design 2 (F/4 Mangin): a ray enters the second refractive surface of the element (surface 2) and travels to the right where it intersects the mirror surface (surface 3). The thickness of surface 2 is therefore positive. The ray is reflected by the mirror surface (surface 3) and travels back through the glass element to surface 2; hence, the notation F9/REFL and the negative surface 3 thickness. Since surface 2 and 4 represent the same surface, the radii are the same.

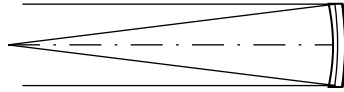
F/4 Paraboloid Objective



Comments A single parabolic mirror objective can be arranged in a variety of forms, the best known being the Newtonian. Here a mirror flat diverts the image away from the light path. A tipped-mirror configuration is the Herschelian; a modern version is untipped and employs an eccentric-pupil to give an accessible image. A “backwards” Newtonian, the Pfund has a large flat-mirror primary. The Pfund has a smaller obscuration than the Newtonian and requires no diffraction-inducing support structure for the folding flat.

As has been mentioned, the on-axis performance of a paraboloid objective is perfect. Off-axis, *coma* quickly degrades image quality. For objectives slower than F/11, the easy-to-fabricate spherical mirror gives the same performance as a paraboloid when diffraction is also considered.

The paraboloid objective has image quality as good as a Cassegrain (design 3) of equivalent FN and aperture diameter, and is easier to align. The Cassegrain has the advantage of being compact.

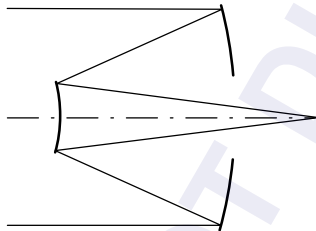
F/4 Mangin

SUR	RDY	THI	GLA
1	-75.15	1.0	BK7
2	-307.1	1.4	F9
3	-123.63	-1.4	F9/REFL
4	-307.1	-1.0	BK7
5	-75.15	-80.48	

Comments The Mangin⁸ was invented by a French engineer of the same name to replace difficult-to-fabricate paraboloids in light houses. The objective relies upon the overcorrected spherical aberration produced by the negative first surface to cancel the undercorrected spherical aberration produced by the focusing, reflective surface. The chromatic aberration of the Mangin is reduced by achromatizing with two glasses of different dispersions. *Secondary spectrum* limits on-axis performance, and coma, one-half that of a spherical mirror, is the primary field-limiting aberration. Kingslake⁹ takes the reader through the design of a Mangin mirror.

Mangin mirrors are susceptible to ghost reflections from the refractive surfaces. Antireflection coatings are usually needed.

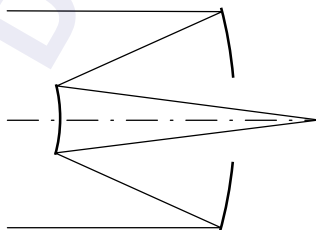
In some cases the overcorrected chromatic aberration of a Mangin is used to cancel undercorrected chromatic aberration produced by a refractive element. The Schupmann or medial objective^{10,11} has a positive refractive element with undercorrected chromatic aberration which is annulled by a Mangin element.

F/4 Cassegrain

SUR	RDY	THI	GLA	CON
STO	-45.72	-16	REFL	-1
2	-19.2	24.035	REFL	-3.236

Comments The ubiquitous Cassegrain is predominant in situations where a small field of view, high resolution, compact size, long effective focal length, and accessible image are required. The classical Cassegrain is composed of a paraboloid primary and hyperboloid secondary, giving perfect imagery on-axis whenever the primary image coincides with the hyperboloidal focus. Coma and *field curvature* limit off-axis performance.

Many books discuss the first- and third-order properties of Cassegrain objectives. The Rutten,¹² Schroeder,¹³ Korsch,¹⁴ and Smith³ texts are among these.

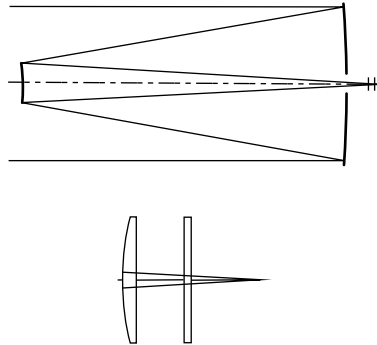
F/4 Ritchey-Chretien

SUR	RDY	THI	GLA	CON
STO	-45.72	-16	REFL	-1.072
2	-19.2	24.034	REFL	-3.898

Comments The aplanatic Cassegrain or Ritchey-Chretien¹⁵ is also corrected for coma, leaving astigmatism and field curvature uncorrected. Both mirrors of the Ritchey-Chretien are hyperboloids.

Numerous modern telescope objectives are of Ritchey-Chretien form; among these are the Hubble space telescope and the infrared astronomical satellite IRAS.

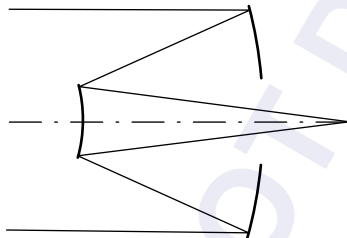
F/9 Ritchey-Chretien Telescope with Two-Lens Corrector



SUR	RDY	THI	GLA	CON
STO	-2139.7	-794.0	REFL	-1.0778
2	-802.83	853.96	REFL	-4.579
3	67.73	2.54	BK7	
4	90.39	9.9	BK7	
5	-1925.6	1.27	BK7	
6	129.1	14.39		

Comments This is a design by Wynne¹⁶ for the correction of the Cassegrain focus of a large (350-cm) Ritchey-Chretien. The corrector removes the inherent astigmatism and field curvature of the Ritchey-Chretien. Other Cassegrain focus correctors are discussed by Schulte,¹⁷ Rosin,¹⁸ and Wilson.¹⁹

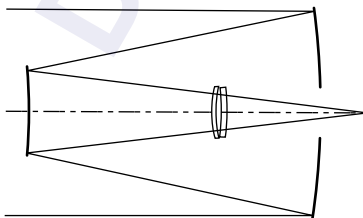
F/4 Dall-Kirkham



SUR	RDY	THI	GLA	CON
STO	-45.72	-16	REFL	-0.6456
Sixth-order term:			0.593E-10	
2	-19.2	24.035	REFL	

Comments The Dall-Kirkham is another Cassegrain corrected for spherical aberration. The primary is an ellipsoid with sixth-order aspheric deformation and the secondary is spherical. An inverse Dall-Kirkham, or Carlisle, is just the reverse, with a spherical primary. There is *zonal spherical aberration* without the sixth-order deformation. Five times more coma is produced by the Dall-Kirkham than the classical Cassegrain, seriously limiting the field of view.

F/4 Cassegrain with Field Corrector and Spherical Secondary

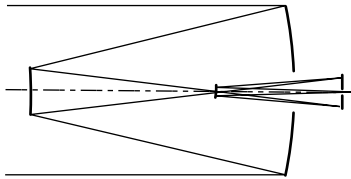


SUR	RDY	THI	GLA	CON
1	-94.21	-27.937	REFL	-1
STO	-94.29	17.72	REFL	
3	17.59	0.35	Silica	
4	8.76	0.491	Silica	
5	-64.15	0.6	Silica	
6	-13.41	13.67		

Comments By adding zero-power refractive correctors, the performance of a reflective objective is substantially enhanced. Zero power is maintained to prevent *axial color*. Such is the case with this objective similar to one designed by Rosin.²⁰ All third-order aberrations, with the exception of distortion, are corrected. The surfaces, with the exception of the primary, are spherical. One of the most attractive features of this design, in comparison to the Schmidt which will be discussed shortly, is the small size of the refractive elements. Add to this the capability of eliminating any remaining spherical aberration in an assembled objective by adjusting the axial positions of the lenses.

Zero *Petzval sum* and, hence, a flat image (in the absence of astigmatism) is ensured by giving the mirrors the same curvature and the lens elements equal and opposite power.

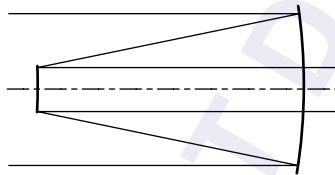
F/15 Spherical-Primary Cassegrain with Reflective Field Corrector



SUR	RDY	THI	GLA	CON
STO	-84.03	-30.69	REFL	
2	-46.56	36.83	REFL	20.97
3	-17.39	-14.77	REFL	-0.8745
4	-20.87	16.26	REFL	-96.62

Comments This well-corrected design from Korsch¹⁴ has an easily manufactured spherical primary and is intended for use as a large-aperture telescope objective. Another all-reflective corrector of Gregorian form has been developed for a fast (F/0.6) spherical primary.²¹

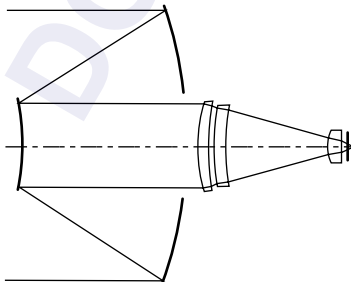
Afocal Cassegrain-Mersenne Telescope



SUR	RDY	THI	GLA	CON
STO	-100	-35	REFL	-1
2	-30	40	REFL	-1

Comments The Mersenne contains two confocal paraboloids working at infinite conjugates. It is aplanatic, anastigmatic, and can be made distortion-free by choosing an appropriate stop location. The utility of confocal mirrors has been emphasized by Baker²² and Brueggeman,²³ and is illustrated in the following design.

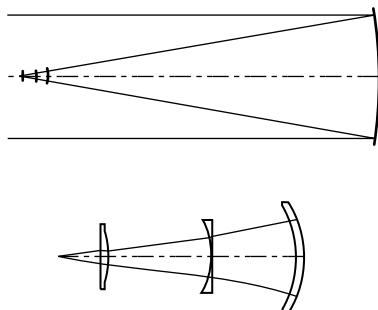
Dual-Magnification Cassegrain



SUR	RDY	THI	GLA	CON
STO	-33.99	-11.69	REFL	-1
2	-10.61	12.76	REFL	-1
3	10.486	0.877	Silicon	
4	25.673	0.462		
5	48.33	0.798	Germanium	
6	22.68	7.57		
7	3.52	1.0	Silicon	
8	4.22	0.377		
9	INF	0.16	Sapphire	
10	INF	0.396		

Comments This IR design is related to one introduced by Fjeidsted.²⁴ The system offers two magnifications and fields of view. The high-magnification configuration is with the afocal Mersenne in the optical path. Removing the secondary lets light pass directly to the refractive assembly and a larger field of view is observed. The spectral range is from 3.3 to 4.2 μm .

F/3.2 Three-Lens Prime Focus Corrector



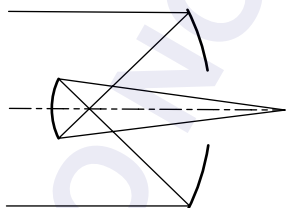
SUR	RDY	THI	GLA	CON
STO	-1494.57	-684.08	REFL	-1
2	-26.98	-2.6	UBK7	
3	-31.3	-22.43		
4	-53.96	-0.586	UBK7	
5	-19.0	-28.87		
6	-33.36	-2.042	UBK7	
7	236.7	-11.65		

Comments This is a three-lens corrector for a 250-cm parabolic mirror. The corrector was developed by Wynne^{16,25} for the region of the spectrum extending from 365 to 1014 nm. It is used to extend the field of a parabolic mirror. Versions for a Ritchey-Chretien primary also exist. The corrector is able to correct spherical aberration, coma, astigmatism, and field curvature while keeping chromatic aberrations under control. The field of view can be extended considerably for smaller apertures.

The three-spherical lens corrector is one of the best large-optics prime-focus correctors to come along, both in terms of image quality and ease of fabrication. Other designs have either not performed as well or were heavily dependent on aspheric figuring.

This and other prime-focus correctors are surveyed in articles by Gascoigne,²⁶ Ross,²⁷ Meinel,²⁸ Schulte,²⁹ Baker,³⁰ and Wynne.³¹

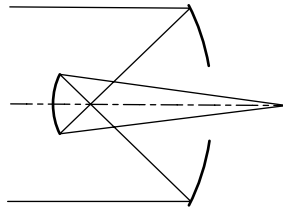
F/4 Gregorian



SUR	RDY	THI	GLA	CON
STO	-24.62	-16	REFL	-1
2	6.4	24.1	REFL	-0.5394

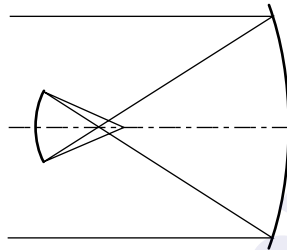
Comments The classical Gregorian is aberration-free on-axis when the paraboloidal mirror image coincides with one of the ellipsoidal-mirror foci; the other focus coincides with the final image. Like the Cassegrain, off-axis image quality is limited by coma and field curvature. The ellipsoidal secondary reimaging the entrance pupil to a location between the secondary and final image. Thus, there exists the possibility of unwanted-light suppression at the primary-mirror image and exit pupil.

The Gregorian is longer than the Cassegrain and thus more expensive to support and house, but it produces an erect image and the concave secondary is easier to produce. In eccentric-pupil versions it has an accessible prime focus.

F/4 Aplanatic Gregorian

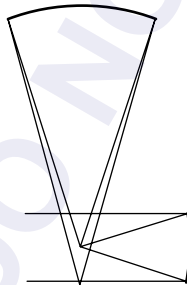
SUR	RDY	THI	GLA	CON
STO	-24.62	-16	REFL	-0.989
2	6.4	24.1	REFL	-0.5633

Comments The aplanatic Gregorian is corrected for spherical aberration and coma. Both mirrors are ellipsoids. Astigmatism and field curvature limit off-axis imagery.

F/1.25 Flat-Medial-Field Aplanatic Gregorian

SUR	RDY	THI	GLA	CON
STO	-34.68	-22.806	REFL	-0.767
2	6.47	7.924	REFL	-0.1837

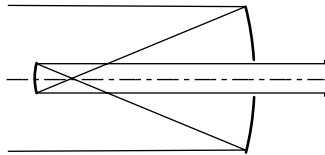
Comments The Gregorian's field performance is enhanced if image accessibility is sacrificed. This version of the Gregorian¹⁴ is aplanatic. A *flat medial image* is achieved by balancing Petzval curvature with astigmatism, which remains uncorrected.

F/1.25 Flat-Medial-Field Aplanatic Gregorian with Spherical Primary

SUR	RDY	THI	GLA	CON
STO	-42.59	-21	REFL	
2	INF	46.51	REFL	
Tilt: 45°				
3	-49.84	-54.08	REFL	-0.078

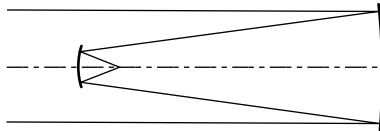
Comments The field of this objective¹⁴ is larger than its cousins, the classical and aplanatic Gregorians, even with the spherical primary. Spherical aberration and coma are corrected, and the medial image is flat. The design has a real intermediate image and exit pupil. The obvious drawback is the size of the secondary in relation to the size of the entrance pupil, which is 15 cm in diameter.

Korsch¹⁴ analyzes two other designs that are loosely referred to as Gregorians.

Afocal Gregorian-Mersenne Telescope

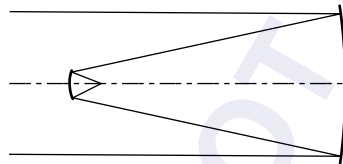
SUR	RDY	THI	GLA	CON
STO	-50	-30	REFL	-1
2	10	40	REFL	-1

Comments The Gregorian Mersenne, also composed of confocal paraboloids, is aplanatic, anastigmatic, and can be corrected for distortion. The Gregorian-Mersenne has an intermediate image and an accessible exit pupil.

F/1.25 Couder

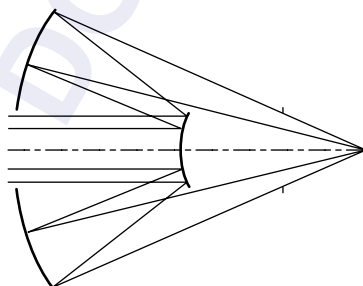
SUR	RDY	THI	GLA	CON
STO	-142.86	-52.9	REFL	-6.285
2	23.08	7.1142	REFL	-0.707

Comments The Couder,³² composed of two conic mirrors, is corrected for third-order spherical aberration, coma, and astigmatism. Unfortunately, the Couder is long for its focal length and the image is not readily accessible.

F/1.25 Aplanatic, Flat-Medial-Image Schwarzschild

SUR	RDY	THI	GLA	CON
STO	-91.57	-38.17	REFL	-2.156
2	23.67	4.637	REFL	5.256

Comments The aplanatic, flat-medial-image Schwarzschild³³ is similar in appearance to the Couder but the secondary mirror and image locations are different for identical *secondary magnifications*.

F/1.25 Aplanatic, Anastigmatic Schwarzschild

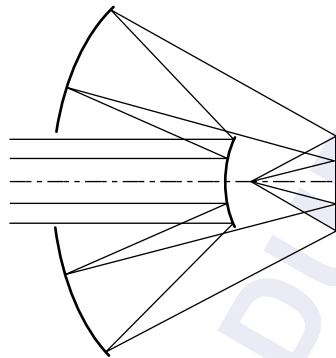
SUR	RDY	THI	GLA
1	30.62	-49.44	REFL
2	80.14	80.26	REFL
STO	INF	24.864	

Comments The spherical-mirror Schwarzschild³³ is aplanatic, *anastigmatic*, and distortion-free.³⁴ The Schwarzschild relies on the principle of symmetry for its high level of aberration correction and a large field of view. All surfaces have the same center of curvature at the aperture stop. Hence, there are no off-axis aberrations. Spherical aberration is produced but each mirror produces an equal and opposite amount, thus canceling the effect of the aberration. Some higher-order aberrations are also corrected.³⁴ Eccentric portions of this design—above and below the optical axis in the picture—form well-corrected, unobscured designs. Zonal spherical aberration from the mix of third- and higher-order terms limits on- and off-axis performance.

An aspheric plate positioned at the center-of-curvature of the mirrors removes this aberration as illustrated in the next design.

Wetherell and Rimmer,³⁵ Korsch,¹⁴ Schroeder,¹³ Linfoot,³⁶ and Gascoigne²⁶ offer a general third-order analysis of two-mirror systems. The closed-form solutions described provide insight into third-order aberration theory of reflective systems.

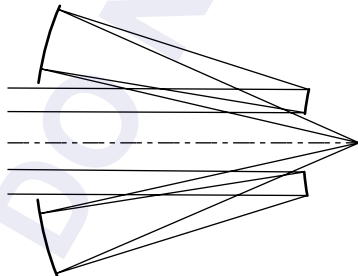
F/1 Aplanatic, Anastigmatic Schwarzschild with Aspheric Corrector Plate



SUR	RDY	THI	GLA
1	24.547	-39.456	REFL
2	63.92	64.528	REFL
STO	INF	-19.098	REFL
A: -0.9998E-7			
B: -0.1269E-9			

Comments With an aspheric plate at the aperture stop, spherical aberration is eliminated. The only aberrations still remaining are of higher order. To correct these, the mirrors must also be aspherized. Linfoot³⁶ and Abel³⁴ describe this design.

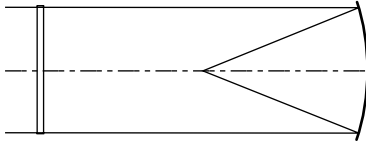
F/1.25 Anastigmatic, Flat-Image Schwarzschild



SUR	RDY	THI	GLA	CON
1	69.7	-50.56	REFL	5.47
STO	71.35	61.26	REFL	0.171

Comments With just two conics, this design type³³ achieves aplanatic and anastigmatic performance on a flat image surface. The flat field is attained by making the curvatures of the mirrors equal. Eccentric portions above or below the optical axis form unobscured versions; the design may alternatively be used off-axis. Sasian^{6,37} and Shafer³⁸ have explored many of this design's features.

F/1.25 Schmidt



SUR	RDY	THI	GLA	CON
STO	1554	1	PSK2	
		A: -0.2825E-5		
		B: -0.1716E-8		
2	INF	52.33		
3	-52.95	-26.215	REFL	

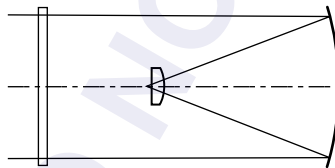
Comments The Schmidt³⁹ also relies on the principle of symmetry; that is, the aperture stop is located at the center of curvature of the spherical mirror and hence the mirror produces no off-axis aberrations.

The Schmidt corrector is flat with aspheric deformation figured in to correct the spherical aberration produced by the mirror. It is positioned at the aperture stop because off-axis aberrations are independent of aspheric deformation when an aspheric surface coincides with the stop. Hence the Schmidt plate has no effect on off-axis aberrations, and the benefits of concentricity are preserved.

The corrector introduces chromatic variation of spherical aberration (spherochromatism). A small amount of positive power in the corrector introduces undercorrected axial color to reduce the effects of this aberration. Further improvement is obtained by achromatizing the corrector with two glasses of different dispersions.

Higher-order aberrations degrade image quality at low focal ratios and large field angles. Kingslake,⁹ Schroeder,¹³ Maxwell,⁴⁰ and Linfoot³⁶ provide additional details of this and other catadioptric objectives.

F/1.25 Field-Flattened Schmidt



SUR	RDY	THI	GLA
STO	598.7	1.155	PSK2
		A: -0.273E-5	
		B: -0.129E-8	
2	INF	40.38	
3	-52.95	-24.06	REFL
4	-10.35	-1.49	PSK2
5	INF	-0.637	

Comments As is known from third-order aberration theory, a thin element will be nearly aberration-free, except for Petzval curvature, and distortion when it is placed in close proximity to an image. Therefore, by properly choosing the lens power and index to give a Petzval curvature of equal and opposite sign to the Petzval curvature introduced by the other optics, the image is flattened.

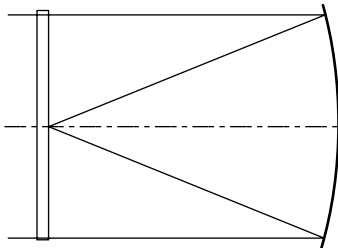
The image in the Schmidt above has been flattened with the lens near the image plane. The only aberrations introduced by the lens are spherochromatism and *lateral color*, lateral color being the most noticeable aberration; this can be removed by achromatizing the field-flattening lens. The close proximity of the lens to the image can cause problems with light scattered from areas

on the lens surfaces contaminated by dirt and dust particles. Clean optics are a must under these circumstances.

The field-flattening lens provides two more positive results. First, the lens introduces a small amount of coma which is compensated by moving the Schmidt corrector toward the mirror somewhat, thus reducing the overall length of the objective. Second, *oblique spherical aberration*, one of the primary field-limiting, higher-order aberrations of the Schmidt, is substantially reduced.

Besides its usual function as a telescope or photographic objective, the field-flattened Schmidt has also been used as a spectrograph camera.⁴¹

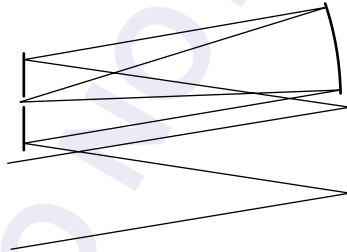
F/1.25 Wright



SUR	RDY	THI	GLA	CON
1	INF	1.0	BK7	
2	-699.8	26.1	BK7	
A: 0.6168E-5				
B: 0.5287E-8				
3	-53.24	-26.094	REFL	1.026

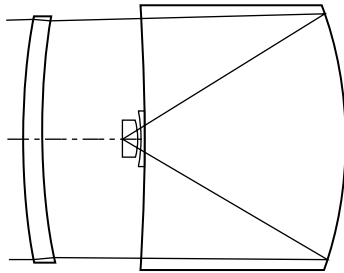
Comments The Wright⁴² is one-half the length of the Schmidt. It also relies on aspheric deformation of the corrector plate for the elimination of spherical aberration. Coma, introduced as the corrector is removed from the center of curvature of the mirror, is cancelled with conic deformation of the mirror; the surface figure is that of an oblate ellipsoid. The remaining astigmatism and Petzval curvature are balanced for a flat medial image. The only on-axis aberration, spherochromatism, is corrected by achromatizing the corrector.

F/4 Reflective Schmidt



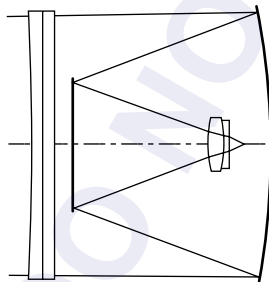
SUR	RDY	THI	GLA
STO	-66752	-67.37	REFL
A: 0.5083E-7			
2	INF	66.6	REFL
3	-133.97	-66.85	REFL

Comments Another way of defeating chromatic aberration is to eliminate it altogether with a reflective corrector.⁴³ The elements are removed from the light path with a field bias (9°), and hence the objective is off-axis. Spherical aberration, coma, and astigmatism are all corrected. At large field angles, beyond about 12° half-field angle, oblique spherical aberration becomes evident, but otherwise this design provides excellent performance on a curved image over a 24° field of view. In order to avoid severe obstruction of the light path, the full 24° can be taken advantage of only in the plane that extends perpendicular to the picture of the design above. A considerably reduced field of view is allowed in the plane of the picture.

F/0.6 Solid Schmidt

SUR	RDY	THI	GLA	CON
1	62.69	1.69	BK7	
2	103.38	8.39		
		A: 0.1492E-4		
		B: 0.1988E-7		
		C: 0.1013E-10		
		D: 0.35E-12		
3	-169.36	16.47	BK7	
STO	-37.05	-16.47	REFL\BK7	
5	INF	-0.3952	BK7	
6	-45.86	-0.245		
7	-10.295	-1.136	BK7	
8	INF	-0.026		

Comments All monochromatic aberrations, with the exception of distortion, are corrected by the appropriately-named solid Schmidt, a system used mostly as a spectrograph camera.⁴⁴ All chromatic aberrations, with the exception of lateral color, are corrected. The imaging theory behind the solid Schmidt is expounded by Baker.⁴⁵ With a refractive index n , the solid Schmidt is n^2 times faster than the conventional Schmidt. Focal ratios of F/0.3 have been obtained. Schulte⁴⁴ offers a design similar to the one given here.

F/1.25 Schmidt-Cassegrain

SUR	RDY	THI	GLA	CON
STO	INF	0.8	489.574	
		A: -0.1928E-4		
		B: 0.298E-7		
2	2269.1	1.0	583.303	
3	INF	16.49		
		A: -0.1085E-4		
		B: 0.2806E-7		
4	-55.9	-15.0	REFL	1.077
5	INF	10.267	REFL	
6	9.1	1.2	489.574	
7	-8.577	0.018		
8	-8.59	0.3	583.303	
9	-87.44	1.317		

Comments The Schmidt-Cassegrain⁴⁶ represents a successful attempt to resolve the difficulties related to the curved image, considerable length, and awkwardly located image of the Schmidt objective, without destroying the positive attributes of the design.

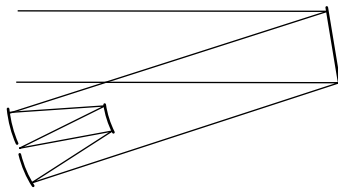
The Schmidt-Cassegrain comes in a wide variety of forms—too many to go into here. Linfoot³⁶ performs an extensive exploration of the design, with one and two aspheric plate correctors. Warmisham⁴⁷ has gone as far as three. Wayman⁴⁸ has analyzed a *monocentric* Schmidt-Cassegrain.

In this fast version of the Schmidt-Cassegrain, the corrector is close to the flat secondary. Usually one or both mirrors are aspherics. An achromatized image-flattening lens has been introduced.

An image-flattening lens is not usually required with a Schmidt-Cassegrain since enough degrees of freedom exist for aberration correction and a flat image. In this case, the secondary mirror is flat and one degree of freedom is lost. Additionally, the primary mirror power introduces a strong Petzval contribution which necessitates a field-flattening lens.

The first three digits of the six-digit code in the glass column identify the indices of the materials in the design, in this case plastics. These are 1.489 and 1.583; the *Abbe-numbers* are given by the last three digits and they are 57.4 and 30.3, respectively. The two plastics have nearly identical thermal coefficients and are very light. Buchroeder⁴⁹ analyzes designs of this variety with two aspheric correctors. Shafer⁵⁰ offers a Schmidt-Cassegrain with an aspheric meniscus corrector. Only two elements are used since the secondary mirror surface is on the corrector. Rutten¹² has examples of Schmidt-Cassegrains in a number of configurations.

F/3.4 Reflective Schmidt-Cassegrain



SUR	RDY	THI	GLA	CON
STO	INF	-92.16	REFL	
A:		0.13344E-6	AN: -0.1255E-1	
2	84	26	REFL	
3	84	-25.848	REFL	-0.3318

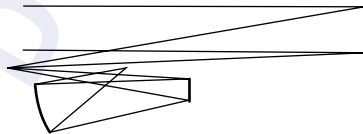
Comments The reflective Schmidt-Cassegrain exhibits all the nice properties of the Schmidt-Cassegrain and, in addition, is achromatic. Schroeder⁵¹ points out that, because the corrector is tilted (9° here), adequate aberration correction requires a nonrotationally symmetric corrector plate. The nonaxially symmetric surface deformation in this design is given by

$$Z = A[(1 - AN)X^2 + (1 + AN)Y^2]^2 \quad (3)$$

where A is the fourth-order symmetric coefficient and AN is the fourth-order nonsymmetric coefficient. The y dimension is in the plane of the picture; the x dimension is perpendicular to it.

Since the corrector is tilted by 9°, the reflected rays are deviated by twice that amount. The element spacings (THI) are no longer measured along the horizontal optical axis after reflection off the corrector, but along the line of deviation. The secondary and tertiary are tilted by 18°.

F/2 Shafer Relayed Virtual Schmidt



SUR	RDY	THI	GLA
STO	-320	-159.82	REFL
2	106.7	80.0	REFL
3	INF	-68.51	REFL
A:		0.1882E-5	
B:		0.1273E-8	
C:		-0.1757E-12	
D:		0.1766E-14	
4	63.967	40.774	REFL

Comments Shafer⁵² has introduced an eccentric-pupil (18-cm stop decenter), virtual Schmidt objective similar to this but with a decentered quaternary mirror. The center of curvature of the spherical primary is imaged by the concave secondary onto the flat Schmidt tertiary mirror.

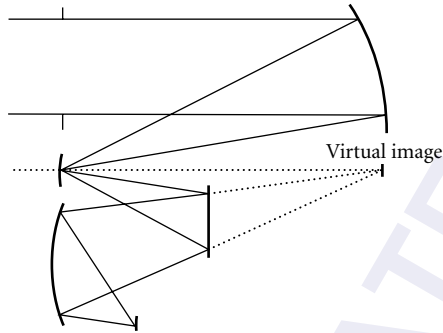


FIGURE 5 Picture of virtual Schmidt with decentered quaternary.

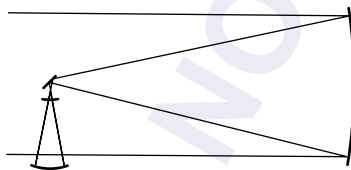
Now the Schmidt plate, which appears to be at the primary center of curvature, is aspherized to produce a well-corrected virtual image, hence the name (see Fig. 5). In this configuration, the Schmidt plate is one-half the size of the primary.

The Schmidt plate and the spherical quaternary mirror form a finite conjugate Schmidt system. Thus, the spherical aberration of this mirror is also corrected.

Figure 5 shows a pictorial representation of the Shafer design with the last mirror decentered to provide a more accessible image. Since the primary and quaternary mirrors no longer share the same axis of symmetry, a two-axis Schmidt corrector is required to remove the aberrations of both mirrors. The shape of this surface is described by Shafer, for an $F/1$, unobscured, wide-field design with an intermediate image and Lyot stop.

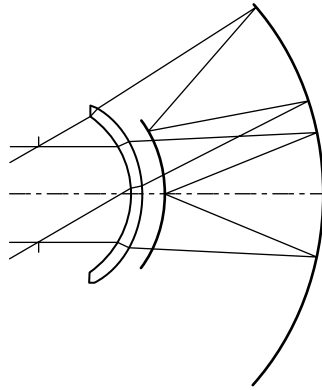
F/2.2 Spherical-Primary Objective that Employs the Schmidt Principle of Correction

SUR	RDY	THI	GLA	CON
STO	-88.07	-42.517	REFL	
2	INF	2.2	REFL	
Tilt: 45°				
3	-4.433	0.33	FK51	
4	-2.527	9.217		
5	-10.21	-9.217	REFL	-0.8631
6	-2.527	-0.33	FK51	
7	-4.433	-0.64		



Comments Baker²² reports on a system where the center of curvature of a large, spherical primary is imaged by a positive lens onto a much smaller mirror where aspheric correction of spherical aberration occurs. A small field offset (0.25°) is required so that the one-to-one relay doesn't reimaging the primary image back onto itself. To avoid overlap, this design is best used over a small or strip field of view.

Because of the geometry of the design, coma, astigmatism, image curvature, chromatic aberrations, and distortion are eliminated in addition to the spherical aberration correction from aspheric figuring of the tertiary mirror. Baker²² offers several other interesting designs in his article, including an $F/0.8$, 10.6- μm , 180° field-of-view Roesch,⁵³ a design that incorporates a Schmidt with a strong negative lens or lens group before the aspheric corrector. The strong divergence produced by this lens reduces the amount of light blocked by the image plane but increases the size of the spherical mirror.

F/2 Maksutov

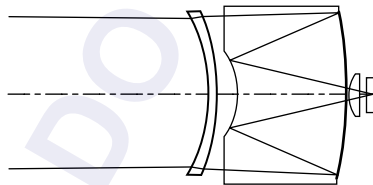
SUR	RDY	THI	GLA	CON
STO	INF	31.788		
2	-23.06	3.5	BK7	
3	-25.53	51.09		
4	-83.9	-43.8	REFL	

Comments The all-spherical Maksutov⁵⁴ was intended as an inexpensive alternative to the Schmidt at slower speeds. In small sizes it is indeed less expensive. The meniscus corrector is “self-achromatic” when the following relationship is satisfied:

$$t = \frac{n^2}{n^2 - 1} (R_2 - R_1) \quad (4)$$

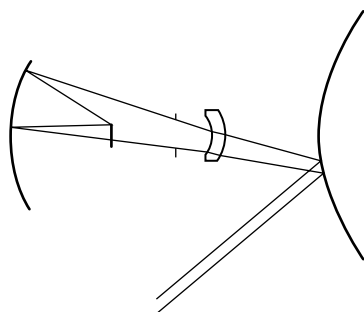
where R_1 and R_2 are the radii, t is the thickness, and n is the refractive index of the corrector.

Bouwers⁵⁵ also developed a meniscus corrector. All elements of the Bouwers are concentric about the aperture stop. This ensures correction of third-order, off-axis aberrations over a nearly unlimited field of view. In exchange for the wide field, axial color is not well-corrected.

F/1.25 Solid Maksutov-Cassegrain

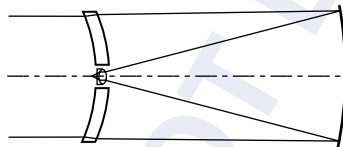
SUR	RDY	THI	GLA	CON
STO	INF	20.5		
2	-20.5	0.955	Silica	
3	-25.92	0.0313		
4	138.58	15.3	Silica	
5	-45.61	-12.973	REFL\Silica	
6	-51.89	13.41	REFL\Silica	
7	INF	0.0475		
8	12.026	1.68	Silica	
9	16.07	0.545		
10	INF	0.394	Silica	
11	INF	0.155		

Comments The solid Maksutov-Cassegrain shown here and the solid Schmidt-Cassegrains have been studied extensively by Wynne.^{56,57} Lateral color is the most consequential aberration left uncorrected.

F/1.2 Wide-Field Objective with Maksutov Correction

SUR	RDY	THI	GLA	CON
1	12.467	-9.684	REFL	-3.243
2	-4.81	-1.267	FK51	
3	-3.762	-3.679		
STO	INF	-17.189		
5	15.98	10.64	REFL	

Comments This very wide field imaging system similar to one in Courtes⁵⁸ is essentially a Maksutov focused on the virtual image of the object produced by the hyperboloidal mirror. Both speed (F/1) and a very wide field of view ($80^\circ \times 120^\circ$) can be achieved with this design form on a flat image but only for small apertures—1.25 cm in this case. Courtes et al.⁵⁹ describes similar systems with refractive and reflective Schmidt plates instead of a Maksutov corrector.

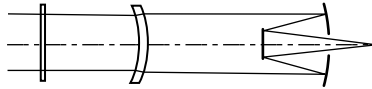
F/1 Gabor

SUR	RDY	THI	GLA
STO	-23.3	2	SKI
2	-25.468	39.5	
3	-83.33	-40	REFL
4	-1.67	-1	BK7
5	9.85	-0.5	SF5
6	-7.71	-0.942	

Comments Another meniscus design was invented by Gabor.⁶⁰ The Gabor is more compact than the Maksutov or Bouwers, and has a smaller focal ratio and field of view.

The design shown here began without the field lens. The lens was introduced into the design with the surface closest to the image being concentric about the *chief ray* and the other surface being aplanatic.⁶¹ A surface concentric about the chief ray is free of coma, astigmatism, distortion, and lateral color. The aplanatic surface is free of spherical aberration, coma, and astigmatism with the result that the lens is coma- and astigmatism-free. The spherical aberration produced by the lens is balanced against the spherical aberration produced by the two other elements. The chromatic aberrations were corrected by achromatizing the lens.

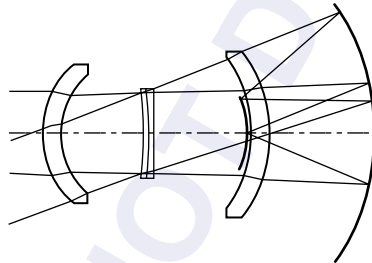
Shafer^{62,63} offers interesting suggestions for design with aplanatic and concentric surfaces. Several varieties of field-flattening lens are described. Kingslake⁹ runs through the design procedure for a Gabor.

F/4 Schmidt-Meniscus Cassegrain

SUR	RDY	THI	GLA
STO	787.7	1.4	BK7
2	INF	32.69	
3	-32.69	2.62	BK7
4	-35.6	63.446	
5	-81.97	-21.78	REFL
6	-79.5	38.65	REFL

Comments This system, originally by Bowers, uses a slightly positive plate to compensate the overcorrected chromatic aberration produced by the meniscus. The Bowers produces very good quality on a flat image, over a large field of view.

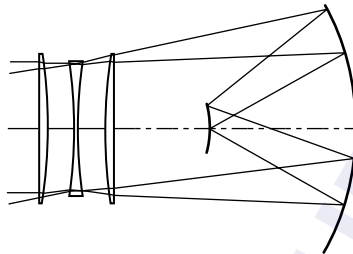
Fourth- and sixth-order deformation added to the plate eliminates any residual spherical aberration. Lateral color and oblique spherical aberration affect performance, although both are small.

F/1.2 Baker Super-Schmidt

SUR	RDY	THI	GLA
1	22.5	4.28	BK7
2	19.13	19.75	
STO	-11,783	1.4	F2
	A:	-0.558E-6	
	B:	0.244E-8	
4	-135.5	1.3	SK2
5	INF	23.92	
6	-29.31	4.28	BK7
7	-32.42	25.1	
8	-55	-25.1	REFL
9	-32.42	-4.28	BK7
10	-29.31	-1.45	

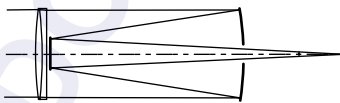
Comments The Baker⁶⁴ super-Schmidt, a design that incorporates both meniscus and Schmidt correction, achieves excellent performance over a wide field of view. The field-limiting aberration of a fast Schmidt, *oblique spherical aberration*, is controlled by adding a concentric meniscus lens which also introduces overcorrected spherical aberration, thus reducing the amount of overcorrection needed from the Schmidt plate. Since oblique spherical is proportional to the amount of overcorrection in the Schmidt plate, the effect of this aberration is reduced.

The most apparent aberration produced by the meniscus is axial color. This is minimized by achromatizing the Schmidt corrector. Spherochromatism is reduced since the magnitudes produced by the Schmidt corrector and meniscus are nearly equal and have opposite signs. Another meniscus element is added to further reduce aberrations.

F/1 Baker-Nunn

SUR	RDY	THI	GLA
1	-491.9	1.06	LLF1
2	-115.6	4.23	
A: -0.8243E-5			
B: 0.1348E-8			
STO	-125.78	0.64	SK3
A: -0.1158E-4			
B: -0.427E-8			
C: -0.7304E-11			
4	125.78	4.23	
A: 0.1158E-4			
B: 0.427E-8			
C: 0.7304E-11			
5	115.6	1.06	LLF1
A: 0.8243E-5			
B: -0.1348E-8			
6	491.87	36.77	
7	-42.03	-21.961	REFL

Comments The Baker-Nunn⁶⁵ was born of work by Houghton⁶⁶ during World War II. Houghton wished to find a less expensive alternative to the Schmidt. The result was a zero-power, three-lens combination with easy-to-make spherical surfaces. Spherical aberration and coma can be eliminated for any position of the corrector. The surfaces have equal radii so they can be tested interferometrically against one another using the Newton ring method. Residual spherical aberration that remains after assembly is removed by altering the spacing between the lenses.

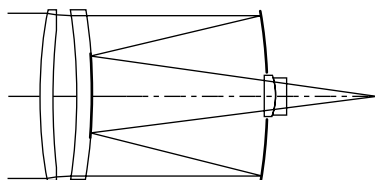
F/10 Houghton-Cassegrain

SUR	RDY	THI	GLA
STO	145	1.2	BK7
2	-172.1	0.164	
3	-111.9	0.639	BK7
4	264.7	44.61	
5	-129.7	-43.16	REFL
6	-63.94	66.84	REFL

Comments A two-lens, afocal corrector developed by Houghton and Sonnefeld⁶⁷ is used here as a corrector for a Cassegrain system. Sigler⁶⁸ has written on the subject of Cassegrains with

Houghton, Schmidt, and Maksutov refractive correctors. This Houghton-Cassegrain gives well-corrected imagery on a curved image surface. An afocal achromatized doublet corrector has also been tried.⁶⁹

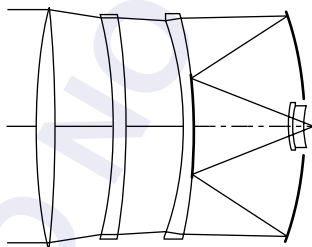
F/3.6 Houghton-Cassegrain



SUR	RDY	THI	GLA
STO	69.64	1.607	UBK7
2	148.71	3.045	
3	-61.43	1.607	LAK21
4	-97.53	21.733	
5	-85.11	-21.733	REFL
6	-97.53	21.733	REFL
7	70.44	1.2	UBK7
8	-15.47	0.18	
9	-15.23	1.3136	SK16
10	-517.29	11.03	

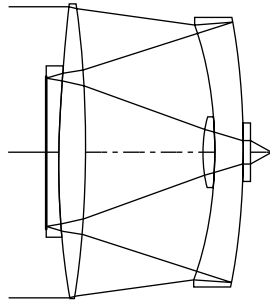
Comments Another Houghton corrector, with meniscus elements, is utilized in this design by D. Rao.⁷⁰ The spectral range is 550 to 850 nm. The design is similar to one introduced by Mandler.⁷¹ Examples of other Houghton-Cassegrains of this form are studied by Gelles.⁷²

F/1.25 Shenker



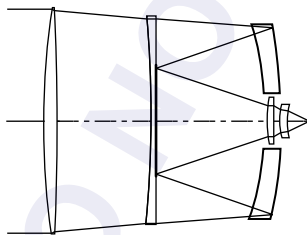
SUR	RDY	THI	GLA
STO	49.42	1.5	BK7
2	-203.6	5.4	
3	-34.7	0.863	BK7
4	-79.25	5.08	
5	-27	0.98	BK7
6	-38.87	9.32	
7	-31.96	-9.32	REFL
8	-38.87	8.1	REFL
9	13.73	0.39	BK7
10	21.8	0.05	
11	7.925	0.895	BK7
12	8.56	0.856	

Comments Shenker has studied a large number of variations on the theme of three-element correctors for a Cassegrain. This is related to one of the configurations developed by Shenker.⁷³ Note that the third corrector is also the secondary mirror. Zonal spherical aberration limits performance on-axis. This may be removed by aspherizing one or more surfaces. All elements are of the same glass. Laiken⁷⁴ has a similar version of this design as well as other catadioptric objectives. Maxwell⁴⁰ has design examples and catadioptric imaging theory.

F/1.25 Mangin-Cassegrain with Correctors

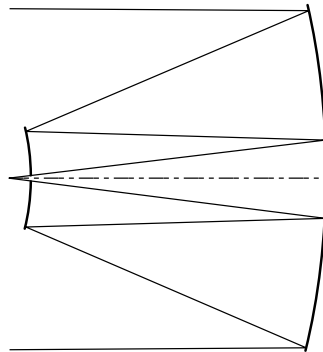
SUR	RDY	THI	GLA
STO	80.62	1.64	BK7
2	-102.3	9.07	
3	-30.43	2.02	BK7
4	-54.52	-2.02	BK7/REFL
5	-30.43	-9.07	
6	-102.3	-1.64	BK7
7	80.62	-1.01	BK7
8	-526.4	1.01	BK7/REFL
9	80.62	1.64	BK7
10	-102.3	8.32	
11	11.06	0.75	BK7
12	-30.43	2.02	BK7
13	-54.52	0.5	SF10
14	52.92	1.445	

Comments Mangin mirrors are evident in this design by Canzek⁷⁵ and two elements are used twice. The design has exceptionally good on-axis performance. Lateral color and higher-order aberrations limit the field.

F/1.25 Mangin-Cassegrain with Correctors

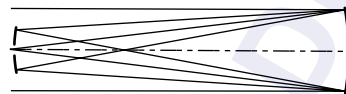
SUR	RDY	THI	GLA
STO	80.83	1.09	FN11
2	-325.9	8.5	
3	-191.4	0.728	FN11
4	-440.3	9.69	
5	-31.44	1.456	FN11
6	-46.13	-1.456	FN11/REFL
7	-31.44	-9.69	
8	-440.3	10	REFL
9	26.97	0.582	FN11
10	38.33	0.544	
11	8.44	0.728	FN11
12	40.87	2.025	

Comments Another short and fast catadioptric by Amon⁷⁶ is shown here. The second corrector is also the secondary mirror.

F/4 Eisenburg and Pearson Two-Mirror, Three-Reflection Objective

SUR	RDY	THI	GLA	CON
STO	-48.0	-17.289	REFL	-1.05
2	-14.472	17.289	REFL	-1.539
3	-48	-18.195	REFL	-1.05

Comments This aplanatic, two-mirror, three-reflection configuration was first introduced by Rumsey.⁷⁷ The design presented here comes from Eisenburg and Pearson.⁷⁸ The first and third surface represent the same surface.

F/4 Shafer Two-Mirror, Three-Reflection Objective

SUR	RDY	THI	GLA	CON
STO	-106.7	-80.01	REFL	-0.4066
2	80.01	80.01	REFL	-5.959
3	-106.7	-80.05	REFL	-0.4066

Comments Shafer⁷⁹ has documented numerous versions of multiple-reflection objectives. This is an aplanatic, anastigmatic design with field curvature. For optimum aberration correction, the primary is at the center of curvature of the secondary mirror. Shafer⁸⁰ suggests a ring field for a flat, accessible image on an annular surface, and a Lyot stop.

A simple ring field design is depicted in Fig. 6. Only one field angle θ is required, easing the difficulties associated with off-axis aberration correction. The single viewing direction is rotated about the optical axis, forming, in this case, a ring image. In reality, less than half the ring image is used to avoid interference of the image with the entering beam.

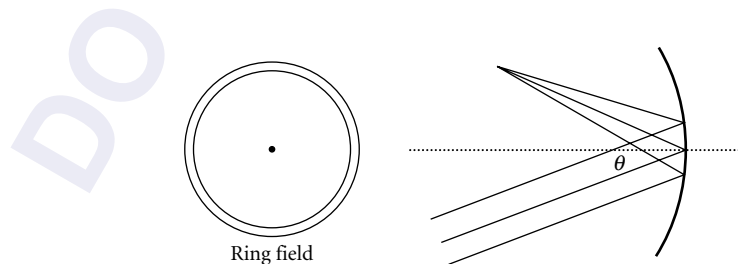
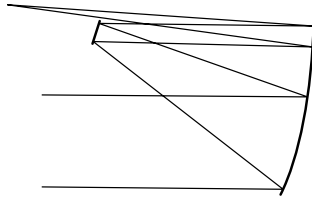


FIGURE 6 The ring field system.

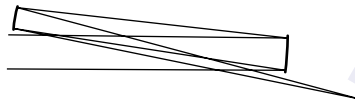
F/15 Two-Mirror, Three-Reflection Objective



SUR	RDY	THI	GLA	CON
STO	-116.33	-46.53	REFL	-1.024
2	-22.6	46.53	REFL	-1.0037
3	-116.33	-67.05	REFL	-1.024

Comments This is another aplanatic, anastigmatic, eccentric-pupil design which gives well-corrected imagery on a curved image. It has a 30-cm stop decenter.

F/15 Yolo



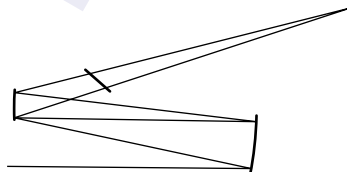
SUR	RDY	THI	GLA	CON
STO	-1015	-160.36	REFL	-4.278
Tilt: -3.5°				
2	1045.72	208.19	REFL	
RDX: 1035.0				
Tilt: -9.82°				
Image tilt: -11.7°				

Comments Leonard^{81,82} invented the Yolo (named after a scenic county in California) so that he could have an achromatic system without obscurations. The result is a tilted and decentered component objective that gives the high contrast of an unobscured refractive objective without the chromatic effects.

Spherical aberration is corrected by the hyperboloidal figuring of the first surface. The anamorphism introduced into the secondary (by a warping harness) corrects astigmatism; RDX is the surface radius of curvature perpendicular to the picture. Coma is eliminated by adjusting the curvatures and tilting the secondary.

Relatives of the two-mirror Yolo are the Solano, an in-line three-mirror Yolo, or the three-dimensional, three-mirror Yolo.⁸³ As in design 28, thickness (THI) is measured along the deviated ray paths. With the angle of reflection known, element decenter may be easily determined.

F/15 Schiefspiegler



SUR	RDY	THI	GLA	CON
STO	-397.2	-101.4	REFL	-0.607
Tilt angle: -4.5°				
2	-552.5	35.84	REFL	
Tilt angle: 3.64°				
3	3411	0.52	BK7	
Tilt angle: 50°				
4	INF	111.11		
Tilt angle: 50.0529°				
Image tilt angle: 22.80°				

Comments The Schiefspiegler (“oblique reflector” in German) was introduced about a century ago, and at the time was called a brachyt (or bent). The motivation behind the Schiefspiegler’s design is essentially the same as the Yolo’s. Like the Yolo, elements are tilted and decentered. Coma and astigmatism are corrected by tilting the secondary and corrector lens. The lens is thin and slightly wedged to minimize chromatic effects. Spherical aberration is corrected with the aspheric deformation of the primary.

A three-mirror Schiefspiegler, or Trischiefspiegler, has been developed by Kutter.⁸⁴ This design is all-reflective and completely achromatic. Like the Schiefspiegler, aspheric deformation of the primary corrects spherical aberration; coma and astigmatism are corrected with element tilts.

A four-mirror Schiefspiegler was recently introduced by Brunn.⁸⁵ For more on unusual telescope objectives, see Manly.⁸⁶

F/8 Catadioptric Herschelian Objective

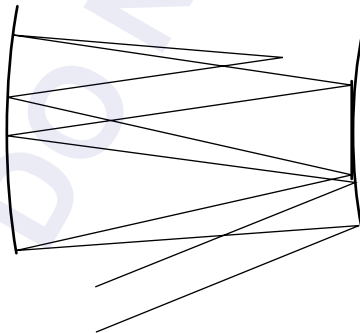


SUR	RDY	THI	GLA
STO	269.61	1.487	BK7
2	INF	2.147	
Element tilt:		0.35°	
3	-269.61	1.321	BK7
4	INF	151.97	
Element tilt:		5.38°	
5	-317.26	-158.69	REFL
Tilt angle:		3.0°	
Image tilt angle:		0.174°	

Comments Several centuries ago, Herschel tilted his large parabolic mirror to give him access to the image. A spherical mirror in this design by D. Shafer⁸⁷ has been tilted for the same reason. Element tilts in the Houghton corrector control the astigmatism introduced by tilting the mirror. The Houghton corrector also eliminates the spherical aberration of the mirror with lens bending. Note the smaller focal ratio of this design compared to either the Yolo or the Schiefspiegler.

Other catadioptric Herschelians, as well as Schiefspiegler and Yolos, have been studied by Buchroeder⁸⁸ and Leonard.⁸² Tilted, decentered, and unobscured Cassegrains are discussed by Gelles.⁸⁹

F/4 SEAL

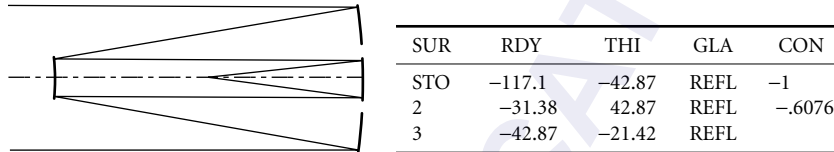


SUR	RDY	THI	GLA	CON
1	181.2	-147.8	REFL	
2	350.9	147.8	REFL	-0.404
STO	INF	-147.8	REFL	
4	350.9	119	REFL	-0.404

Comments For an all-reflective objective, this flat-image design provides an exceptionally wide, unobscured field of view—greater than 90° with a ring field. Referred to as the SEAL,⁹⁰ it is

derived from its cousin the WALRUS;⁹¹ a related design has been described by Shafer.⁹² The SEAL is another Mersenne-Schmidt hybrid: primary and secondary form an inverse-Mersenne; tertiary and quaternary (also the secondary) form a reflective Schmidt. Residual spherical aberration limits the performance, but by aspherizing the flat, this residual aberration is corrected as well. Clearing all obscurations requires at least a 22° field offset. The SEAL shown here is optimized for a 20° strip field although a square, rectangular, annular, or almost any shape field is possible.

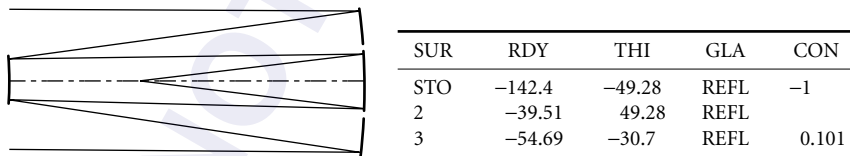
F/4 Paul Three-Mirror Objective



Comments This design is based on work by Paul⁹³ and later by Baker,⁹⁴ who was looking for an achromatic field corrector for a parabolic primary. Their efforts culminated in a design similar to this one, which combines the properties of an afocal Cassegrain-Mersenne in the first two elements with the properties of an all-reflective Schmidt in the secondary and tertiary elements. Since both modules are corrected for spherical aberration, coma, and astigmatism to third order, the complete system is aplanatic and anastigmatic. Petzval curvatures are equal and opposite so a flat image is achieved. The conic deformation of the secondary is modified to give it an ellipsoidal shape. This gives the required Schmidt aspherization needed to correct the spherical aberration of the tertiary mirror.

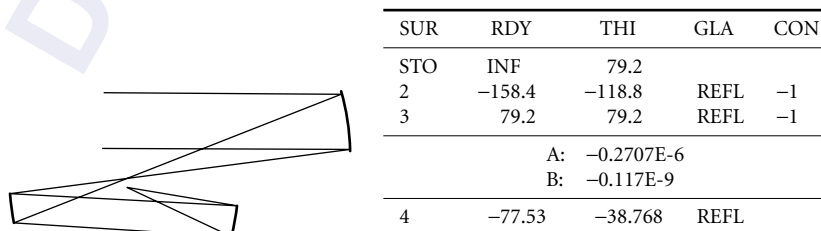
Other all-reflective designs have been proposed by Meinel^{21,95} and Baker.²² The Meinel-Shack objective⁹⁶ exhibits similar performance and offers a more accessible image.

F/4 Alternative Paul Three-Mirror Objective



Comments This Paul objective has an aspheric tertiary mirror, instead of an aspheric secondary.

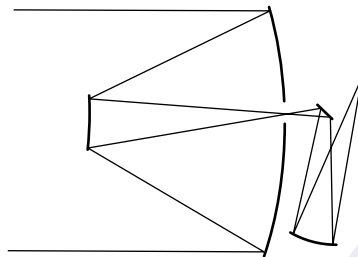
F/4 Off-Axis, Eccentric-Pupil Paul-Gregorian



Comments This eccentric-pupil (22-cm), off-axis (1°) design utilizes a Gregorian-Mersenne module in the primary and secondary mirrors. Spherical aberration produced by the spherical tertiary mirror is corrected by superimposing aspheric deformation on the paraboloid secondary, located at the tertiary mirror center of curvature. With all concave surfaces, field curvature is uncorrected. A real and accessible exit pupil and intermediate image offer possibilities for excellent stray-light suppression.

As is the case with the virtual Schmidt system, the tertiary mirror may be decentered to provide a more convenient image location. This requires two-axis aspheric deformation of the secondary mirror.⁹⁷

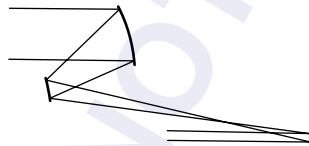
F/4 Three-Mirror Cassegrain



SUR	RDY	THI	GLA	CON
STO	-39.67	-15.814	REFL	-0.9315
2	-10.66	21	REFL	-2.04
3	INF	-9.05	REFL	
Tilt: 45°				
4	13.66	13.651	REFL	-0.4479

Comments A design similar to the aplanatic, anastigmatic, flat-image design shown here was conceived by Korsch⁹⁸ and is described by Williams,⁹⁹ Korsch,¹⁰⁰ and Abel.⁴³ The exit pupil is accessible and an intermediate image exists. A 1° field offset is needed to displace the image from the folding flat. Residual coma limits field performance. Small element tilts and decenters will improve the performance of this design.

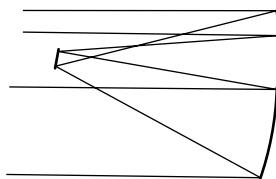
Three-Mirror Afocal Telescope



SUR	RDY	THI	GLA	CON
STO	-100.725	-33.514	REFL	-1
2	-46.109	100	REFL	-3.016
3	-74.819	-55.56	REFL	-1

Comments This $5\times$ afocal design from Smith² is an eccentric-pupil Cassegrain and a parabolic tertiary combined. The design is aplanatic and anastigmatic. The entrance pupil is decentered by 32 cm.

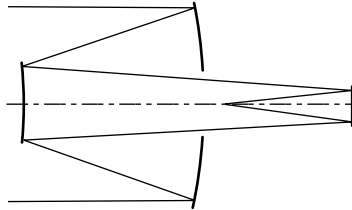
Three-Mirror Afocal Telescope



SUR	RDY	THI	GLA	CON
STO	-240	-200	REFL	-1
2	-160	200	REFL	-9
3	-480	-250	REFL	-1

Comments A similar design by Korsch¹⁴ is also aplanatic and anastigmatic. The entrance pupil is decentered by 20 cm. Other afocal designs are described by Gelles¹⁰¹ and King.¹⁰²

F/4 Three-Mirror Cassegrain



SUR	RDY	THI	GLA	CON
STO	-59.64	-18.29	REFL	-1.134
2	-28.63	33.74	REFL	-2.841
3	-55.05	-13.244	REFL	-5.938

Comments Robb¹⁰³ has introduced another aplanatic, anastigmatic, flat-image, three-mirror Cassegrain without an intermediate image.

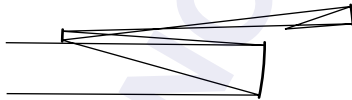
F/6.7 Spherical Primary Three-Mirror Objective



SUR	RDY	THI	GLA	CON
STO	-429.67	-149.87	REFL	
2	-104.16	211.14	REFL	3.617
3	-126.49	-73.0	REFL	-0.179

Comments Making the largest element in an objective a spherical mirror reduces cost and may enhance performance. This aplanatic, anastigmatic, flat-image, eccentric-pupil design (-35 cm stop decenter) with an unobscured light path is similar to one described by Korsch¹⁴ and another developed for use as an astrometric camera by Richardson and Morbey.¹⁰⁴

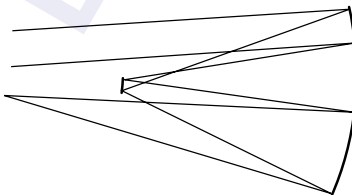
F/4 Spherical Primary Three-Mirror Objective



SUR	RDY	THI	GLA	CON
STO	-194.58	-79.13	REFL	
2	-64.42	113.68	REFL	12.952
3	-38.47	-26.24	REFL	-0.4826

Comments Here is another aplanatic, anastigmatic, flat-field, eccentric-pupil design with a 17-cm stop decenter and large spherical primary. There is an intermediate image and an accessible exit pupil.

F/4 Three-Mirror Korsch Objective

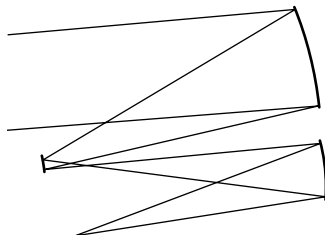


SUR	RDY	THI	GLA	CON
1	-201.67	-133.36	REFL	-0.689
STO	-96.5	131.8	REFL	-1.729
3	-172.54	-200.83	REFL	

Comments This off-axis (5°) design by Korsch¹⁰⁵ is aplanatic, anastigmatic, and has a flat image. The same configuration has been employed by Pollock¹⁰⁶ as a collimator. Characteristics include a large field of view, low pupil magnification, accessible pupils, and an intermediate image.

The tertiary in this design is spherical. With reoptimization, the secondary may also be spherical.

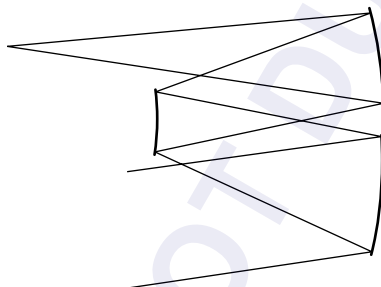
F/4 Three-Mirror Cook Objective



SUR	RDY	THI	GLA	CON
1	-123.2	-57.38	REFL	-0.7114
2	-37.46	57.45	REFL	-3.824
3	-51.89	-35.87	REFL	-0.1185
STO	INF	-15.92		

Comments This objective was introduced by Cook.¹⁰⁷⁻¹⁰⁹ The aplanatic, anastigmatic, flat-image design shown here has a larger pupil magnification and a smaller field than the previous design. The eccentric-pupil, off-axis design has a -3.2 -cm stop decenter and a 5° field bias. A space-based surveillance objective in this configuration has been developed and built by Wang et al.¹¹⁰

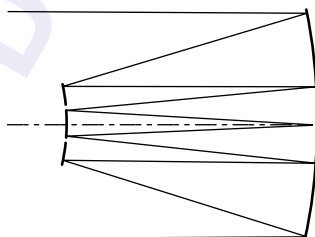
F/4 Three-Mirror Wetherell and Womble Objective



SUR	RDY	THI	GLA	CON
1	-166.19	-38.78	REFL	-2.542
STO	-55.19	38.78	REFL	-0.428
3	-82.46	-65.24	REFL	0.133

Comments Another aplanatic, anastigmatic, flat-image, off-axis (9°) design has been introduced by Wetherell and Womble.¹¹¹ Figosky¹¹² describes a variant of this form to be sent into orbit. The aperture stop is located at the secondary mirror; hence, this mirror is symmetric with respect to the optical axis.

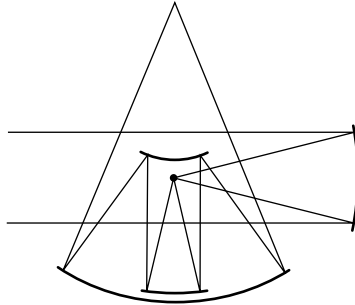
F/10 Korsch Three-Mirror, Four-Reflection Objective



SUR	RDY	THI	GLA	CON
STO	-66.44	-22.15	REFL	-1.092
2	-22.15	22.15	REFL	-1.295
3	-66.44	-22.15	REFL	-1.092
4	-44.29	21.96	REFL	0.8684

Comments The three-mirror, four-reflection design shown here from Korsch¹¹³ is extremely compact for its 200-cm focal length, and the image is accessible.

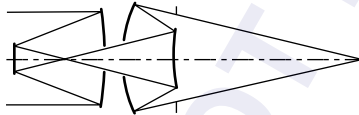
F/1.25 McCarthy



SUR	RDY	THI	GLA	CON
STO	-81.57	-40.21	REFL	-1
2	INF	25.09	REFL	
Tilt: 45°				
3	-48.68	-29	REFL	-1
4	-19.15	30.64	REFL	
5	-49.85	-65.483	REFL	

Comments McCarthy¹¹⁴ intended this design, which combines a Cassegrain-Mersenne primary and tertiary mirror with a quaternary and quintenary Schwarzschild arrangement, as a wide strip-field imager. Both the Mersenne and Schwarzschild groups are separately corrected for spherical aberration, coma, and astigmatism. The Petzval curvature of the Mersenne is equal and opposite in sign to the Petzval curvature of the Schwarzschild and hence there is no net Petzval curvature. The quaternary mirror may be moved out from the entering beam with only a slight reduction in performance.

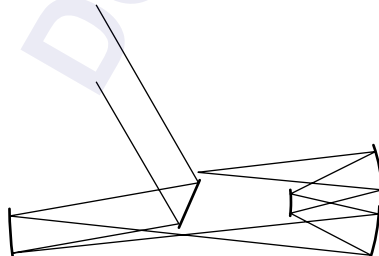
F/2.2 Cassegrain Objective with Schwarzschild Relay



SUR	RDY	THI	GLA	CON
1	-51.49	-19.01	REFL	-1.048
2	-37.37	34.19	REFL	-20.35
3	38.18	-10.493	REFL	-1.358
4	29.94	11.27	REFL	
STO	INF	40.484		

Comments Williams¹¹⁵ describes a technique for optimizing a high-resolution system similar to this one while maintaining proper clearances, obscuration sizes, and packaging requirements. An all-reflective zoom system of the above configuration, developed by Johnson et al.,¹¹⁶ gives a 4× zoom range and a field-of-view range of 1.5 to 6.0°. The Schwarzschild module and image position change with zoom position, while the front Cassegrain module remains fixed.

F/4 Altenhof Objective

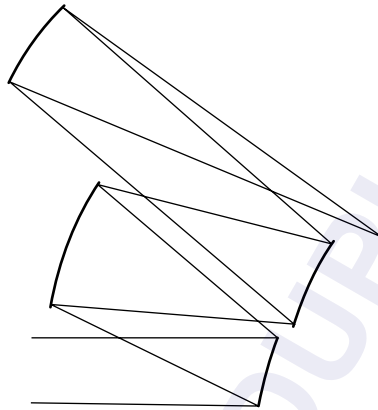


SUR	RDY	THI	GLA	CON
STO	INF	-80	REFL	
Tilt: 25°				
2	155.64	165.53	REFL	
3	-77.26	-38.57	REFL	-0.0897
4	-34.146	40.367	REFL	
5	-80.65	-82.539	REFL	

Comments This objective is similar to one designed by Altenhof.¹¹⁷ Intended as a ring field system, the flat primary couples the light incident from a large azimuthal angle (60°) into the system where the spherical primary mirror focuses the light to a poorly corrected image.

The three-mirror Offner relay,^{9,118} a unit magnification relay which is corrected for spherical aberration, coma, and astigmatism, improves the degraded image in the process of reimaging it to a flat focal surface in the form of an annulus. A two-dimensional scene is imaged by rotating the flat mirror about an axis perpendicular to the picture so as to scan the other dimension. A two-dimensional mosaic image can also be produced by building up a series of one-dimensional annular strip images as the imaging system is moved along its trajectory.

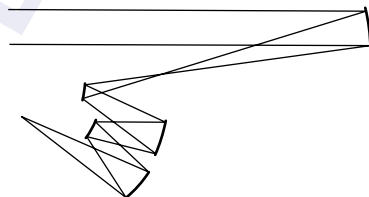
F/4.5 Shafer Four-Mirror, Unobscured Objective



SUR	RDY	THI	GLA
1	158.1	-71.21	REFL
		Tilt angle:	-16.44°
STO	186.8	74.25	REFL
		Tilt angle:	-20.88°
3	337.4	-111.4	REFL
		Tilt angle:	-24.82°
4	239.1	121.4	REFL
		Tilt angle:	-34.76°
		Image tilt angle:	-24.29°

Comments This is a tilted and decentered-component infrared imaging system by David Shafer. Mirror tilts provide an unobscured path and help correct the aberrations. Thickness is measured along the deviated ray paths. With the reflection angle known, element decenter may be easily determined.

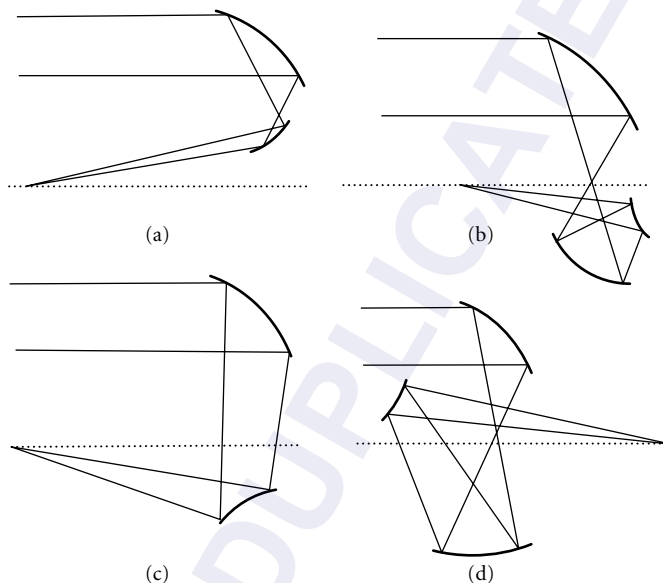
F/4.5 Shafer Five-mirror, Unobscured Objective



SUR	RDY	THI	GLA
1	-239.5	-160.2	REFL
		Tilt angle:	6.4°
2	-228.9	48.69	REFL
		Tilt angle:	-9.2°
3	-75.94	-37.24	REFL
		Tilt angle:	-19.01°
STO	-39.81	39.24	REFL
		Tilt angle:	-28.82°
5	-78.72	-74.5	REFL
		Tilt angle:	-40.55°
		Image tilt angle:	-11.28°

Comments Another all-spherical, tilted, and decentered-component infrared imager by Shafer is presented here. The entrance pupil is accessible and there is an intermediate image. A number of variations on this arrangement are described by Shafer.⁸⁰

Korsch Two- and Three-Mirror Objectives



Comments A new class of eccentric-pupil objectives has been introduced by Korsch.¹⁴ Unlike most systems, which are conceived using third-order aberration theory, these systems are based upon the fulfillment of *axial stigmatism*, the *Abbe sine condition*, and the *Herschel condition*; meeting these three conditions guarantees a perfect axial point image, axially perpendicular image area, and axial line element, respectively.

Design examples are not given for two reasons. First, rays strike some mirror surfaces at angles greater than 90° , which can cause ray-trace errors. Second, some of the surface shapes are particularly complex and must be entered in design software as a user-defined surface.

Design (c) gives perfect imagery on-axis and less than one milliradian resolution at all other points over a 6° field of view, for an aperture diameter equal to $F/6.0$ (F is focal length).

29.5 FIELD-OF-VIEW PLOTS

The plots that follow give rms spot size and angular resolution as a function of half-field of view. The curves have been generated by calculating the resolution for a number of field angles and connecting them with smooth curves. The dashed horizontal line is the Airy disc diameter for $0.55\text{-}\mu\text{m}$ radiation.

The numbers in the plots correspond to the designs presented in the previous section. The aperture of each design is 20 cm and the spectral range 480 to 680 nm, unless stated otherwise in the previous section.

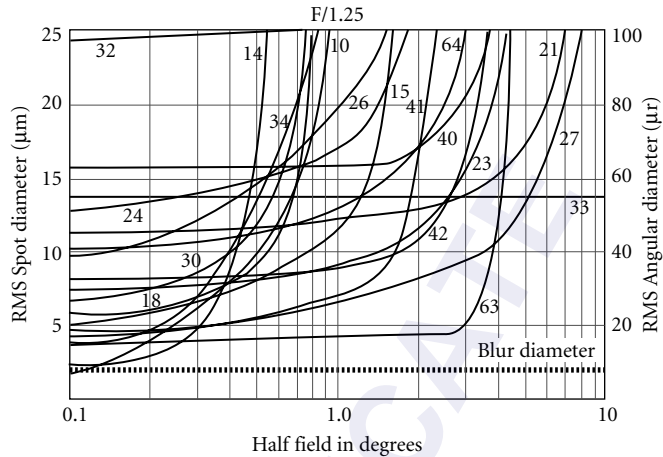


FIGURE 7 Field-of-view plots: F/1.25 on a flat image.

It should be kept in mind that these are representative designs: they have usually been optimized for a specific purpose; they are meant to be a starting place for the design of a new system that may have entirely different requirements.

Flat-field designs show consistent performance out to the field angle for which the objective is optimized. Beyond this point, the graph leaps upward. Reoptimization is needed if the field of view is to be extended further; a considerable increase in the average rms spot size may occur if this is attempted. The curved-image designs show a quadratic dependence with field angle.

Off-axis and eccentric-pupil designs have rectangular fields with most of the field of view in one dimension only. Data plotted for these designs are representative of the larger field.

In Figs. 8 and 10, plots for the curved image designs are provided. The curvature of the image is adjusted to give optimum performance. Figures 7 and 9 are for the flat image designs.

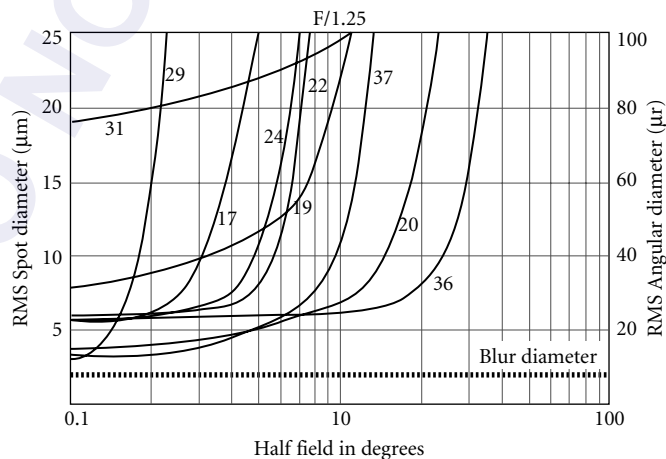


FIGURE 8 Field-of-view plots: F/1.25 on a curved image.

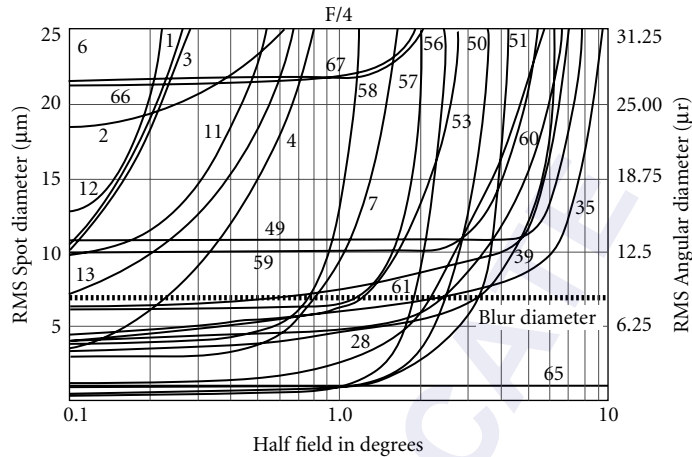


FIGURE 9 Field-of-view plots: F/4 on a flat image.

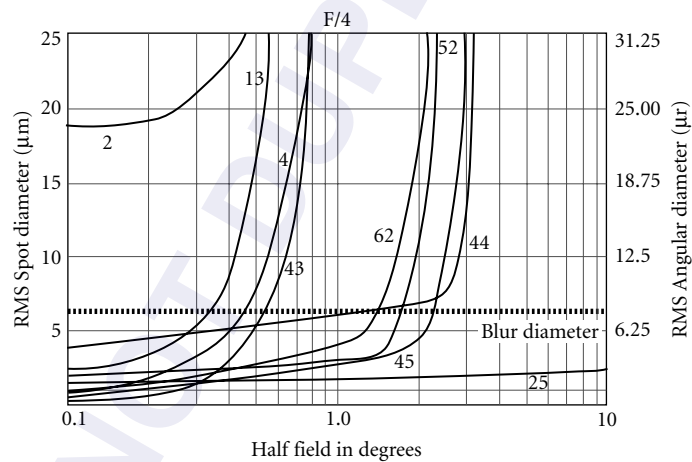


FIGURE 10 Field-of-view plots: F/4 on a curved image.

29.6 DEFINITIONS

Abbe number: A number that indicates the dispersion of a glass. Low dispersion glasses have a high Abbe number.

Abbe sine condition: A condition for zero coma, based on the requirement of equal marginal and paraxial magnifications. See Welford,⁶¹ Kingslake,⁹ or Korsch.¹⁴

anastigmatic: A surface or system free of astigmatism. Also stigmatic.

aperture stop: The aperture that limits the size of the axial beam passing through the system; the *chief ray* always passes through its center.

aplanatic: A surface or system that is corrected for spherical aberration and coma.

astigmatism: An aberration that generates two different focal positions for rays in two perpendicular planes centered on the optical axis. These are called the sagittal and tangential planes.

axial color: The variation in focal position with wavelength.

axial stigmatism: A characteristic of a surface which is able to produce a perfect image of a single point object on-axis.

catadioptric: An optical system composed of refractive and reflective elements: catoptric, reflective and dioptric, refractive.

chief ray: A ray that passes through the center of the aperture stop and the edge of the image or object.

coma: An aberration resulting from the change in magnification with ray height at the aperture, so that rays near the edge of the aperture are focused further from rays near the axis, for the same field point.

conic constant: A constant defined by

$$k = -\epsilon^2$$

where ϵ is the eccentricity of the conic.

distortion: The variation in magnification with field angle.

entrance pupil: The image of the aperture stop in object space. The chief ray passes or appears to pass through the center of the entrance pupil.

exit pupil: The image of the aperture stop in image space. The chief ray passes or appears to pass through the center of the exit pupil.

focal ratio: The effective focal length of an objective divided by its entrance-pupil diameter. Focal ratio is also referred to as the FN, F-number, and speed.

field curvature: Image curvature produced by the combined effects of astigmatism and Petzval curvature. When astigmatism is absent, the image surface coincides with the Petzval surface.

field stop: An aperture that limits the size of an intermediate or final image.

Herschel condition: A condition for invariance of aberrations with change in axial conjugates. See Welford⁶¹ and Korsch.¹⁴

higher-order aberrations: Aberrations defined by the higher-order terms in the aberration power series expansion. See Welford⁶¹ and Schulz.¹

lateral color: An aberration that produces a dependence of image size on wavelength; also called chromatic difference of magnification.

Lyot stop: A real and accessible image of the aperture stop; used to block stray light.

marginal ray: A ray that passes through the center of the object or image and past the edge of the aperture stop.

medial image: The image halfway between the sagittal and tangential images. See Welford.⁶¹

monocentric system: An optical system in which all surfaces are concentric about the chief ray.

oblique spherical aberration: A *higher-order aberration* that is the variation of spherical aberration with field angle.

optical axis: The axis about which all optical elements are symmetric. In tilted and decentered systems, each element has a unique optical axis.

Petzval sum: The sum defined by

$$p = \sum \frac{\phi}{n}$$

where ϕ is element power and n is the index of refraction. The reciprocal of the Petzval sum is the image radius of curvature.

secondary magnification: System focal length divided by primary-mirror focal length.

secondary spectrum: The difference in focal position between two wavelengths corrected for axial color and one other wavelength which is not. For example, the blue and red focus coincide and the yellow focus is axially displaced.

spherical aberration: The only on-axis monochromatic aberration, spherical aberration results from rays at different heights coming to focus at different points along the optical axis. Smith,³ Rutten,¹² Kingslake,⁹ Mackintosh,⁸³ and Welford⁶¹ discuss aberrations. Welford specifically addresses aberrations.

third-order aberrations: Any of the Seidel aberrations: spherical aberration, coma, astigmatism, Petzval curvature, and distortion. See Welford.⁶¹

vignetting: The off-axis clipping of light by apertures in an optical system.

virtual image: A real image is visible when a screen is placed at its location. The image is visible because rays from the object converge at the image. A virtual image is not visible when a screen is placed at its location since real rays do not converge.

zonal spherical aberration: The incomplete correction of spherical aberration at radial zones in the aperture. For example, spherical aberration could be corrected for rays close to the center and edge of the aperture, but not corrected at other ray heights in the aperture.

29.7 REFERENCES

1. G. Schulz, "Aspheric Surfaces," *Progress in Optics* **XXV**:351 (1988).
2. W. Smith, *Modern Lens Design*, McGraw-Hill, Inc., New York, 1992.
3. W. Smith, *Modern Optical Engineering*, McGraw-Hill, Inc., New York, 1990.
4. E. Everhart, and J. Kantorski, "Diffraction Patterns Produced by Obstructions in Reflective Telescopes of Modest Size," *Astronomical Journal* **64**:455 (1959).
5. K. Thompson, "Aberration Fields in Tilted and Decentered Optical Systems," Ph.D. dissertation, Optical Sciences Center, University of Arizona, 1980.
6. J. Sasian, "Review of Methods for the Design of Unsymmetrical Optical Systems," *SPIE* **1396**:453 (1990).
7. R. Fischer, "What's So Different About IR Lens Design?" *SPIE Proceedings* **CR41**:117 (1992).
8. A. Mangin, "Memorial de L'officier du Genie (Paris)," **25**(2):10, 211 (1876).
9. R. Kingslake, *Lens Design Fundamentals*, Academic Press, Inc., New York, 1978.
10. L. Schupmann, "Die Medial-Fernrohre," *Leipzig* (1899).
11. E. Olsen, "A Schupmann for Amateurs," in *Advanced Telescope Making Techniques*, Willmann-Bell, Inc., Richmond, Va., 1986, p. 223.
12. H. Rutten, and M. Venrooij, *Telescope Optics*, Willmann-Bell, Inc., Richmond, Va., 1988.
13. D. Schroeder, *Astronomical Optics*, Academic Press, Inc., San Diego, Calif., 1987.
14. D. Korsch, *Reflective Optics*, Academic Press, San Diego, Calif. 1991.
15. M. Chretien, "Le Telescope de Newton et le Telescope Aplanatique," *Rev. d'Optique* (2):49 (1922).
16. C. Wynne, "Field Correctors for Large Telescopes," *Applied Optics* **4**(9):1185 (1965).
17. D. Schulte, "Anastigmatic Cassegrain Type Telescope," *Applied Optics* **5**(2):309 (1966a).
18. S. Rosin, "Ritchey-Chretien Corrector System," *Applied Optics* **5**(4):675 (1966).
19. R. Wilson, "Corrector Systems for Cassegrain Telescopes," *Applied Optics* **7**(2):253 (1968).
20. S. Rosin, "Corrected Cassegrain System," *Applied Optics* **3**(1):151 (1964).
21. A. Meinel, M. Meinel, D. Su, and Ya-Nan Wang, "Four-Mirror Spherical Primary Submillimeter Telescope Design," *Applied Optics* **23**(17):3020 (1984).
22. J. Baker, "Explorations in Design," *SPIE* **147**:102 (1978).
23. H. Brueggeman, *Conic Mirrors*, Focal Press, London, 1968.

24. T. Fjeidsted, "Selectable Field-of-View Infrared Lens," U.S. Patent 4,453,800, 1984.
25. C. Wynne, "A New Wide-Field Triple Lens Paraboloid Field Corrector," *Monthly Notices Royal Astronomical Society* **167**:189 (1974).
26. S. Gascoigne, "Recent Advances in Astronomical Optics," *Applied Optics* **12**:1419 (1973).
27. F. Ross, "Lens Systems for Correcting Coma or Mirrors," *Astrophysics Journal*, **81**:156 (1935).
28. A. Meinel, "Aspheric Field Correctors for Large Telescopes," *Astrophysical Journal* **118**:335 (1953).
29. D. Schulte, "Prime Focus Correctors Involving Aspherics," *Applied Optics* **5**(2):313 (1966).
30. J. Baker, *Amateur Telescope Making, Book Three*, Scientific American, Inc., 1953, p. 1.
31. C. Wynne, "Field Correctors for Astronomical Telescopes," *Progress in Optics*, **X**:139 (1972).
32. A. Couder, *Compt. Rend. Acad. Sci. Paris*, 183, II, 1276, p. 45 (1926).
33. K. Schwarzschild, "Untersuchungen zur Geometrischen Optik, II; Theorie der Spiegelteleskope," *Abh. der Konigl. Ges. der Wiss. zu Gottingen, Math.-phys. Klasse*, 9, Folge, Bd. IV, No. 2 (1905).
34. I. Abel, and M. Hatch, "The Pursuit of Symmetry in Wide-Angle Reflective Optical Designs," *SPIE* **237**:271 (1980).
35. W. Wetherell, and P. Rimmer, "General Analysis of Aplanatic Cassegrain, Gregorian and Schwarzschild Telescopes," *Applied Optics* **11**(12):2817 (1972).
36. E. Linfoot, *Recent Advances in Optics*, Oxford University Press, London, 1955.
37. J. Sasian, "Design of a Schwarzschild Flat-Field, Anastigmatic, Unobstructed, Wide-Field Telescope," *Optical Engineering* **29**(1) January:1 (1990).
38. D. Shafer, "An All-Reflective Infrared Target Simulator Design," *October OSA Workshop on Optical Fabrication and Testing*, 1988.
39. B. Schmidt, *Mitt. Hamburg. Sternwart* **7**(36) (1932).
40. J. Maxwell, *Catadioptric Imaging Systems*, Elsevier, New York, 1972.
41. C. Wynne, "Shorter than a Schmidt," *Monthly Notices Royal Astronomical Society* **180**:485 (1977).
42. F. Wright, "An Aplanatic Reflector with a Flat-Field Related to the Schmidt Telescope," *Astronomical Society of the Pacific* **47**:300 (1935).
43. I. Abel, "Mirror Systems: Engineering Features, Benefits, Limitations and Applications," *SPIE* **531**:121 (1985).
44. D. Schulte, "Auxiliary Optical Systems for the Kitt Peak Telescopes," *Applied Optics* **2**(2):141 (1963).
45. J. Baker, "The Solid Glass Schmidt Camera and a New Type of Nebular Spectrograph," *Proceeding American Philosophical Society* **82**:323 (1940).
46. J. Baker, "A Family of Flat-Field Cameras, Equivalent in Performance to the Schmidt Camera," *Proceeding American Philosophical Society* **82**(3):339 (1940).
47. A. Warmisham, British Patent 551,112, 1941.
48. R. Wayman, "The Monocentric Schmidt-Cassegrain Cameras," *Proceeding Physics Society (London)*, **63**:553 (1944).
49. R. Buchroeder, "Catadioptric Designs," OSC, University of Arizona, Technical Report 68, May, 1971.
50. D. Shafer, "Well-Corrected Two-Element Telescope with a Flat Image," *Optica Acta* **28**(11):1477 (1981).
51. D. Schroeder, "All-Reflecting Baker-Schmidt Flat-Field Telescopes," *Applied Optics* **17**(1):141 (1978).
52. D. Shafer, "Design with Two-Axis Aspheric Surfaces," *SPIE* **147**:171 (1978).
53. M. Roesch, *Transactions of the International Astronomical Union* **VII**:103 (1950).
54. D. Maksutov, *Journal of the Optical Society of America*, **34**:270 (1944).
55. A. Bouwers, *Achievements in Optics*, Elsevier, Amsterdam, 1950.
56. C. Wynne, "Maksutov Spectrograph Cameras," *Monthly Notices Royal Astronomical Society* **153**:261 (1971).
57. C. Wynne, "Five Spectrograph Camera Designs," *Monthly Notices Royal Astronomical Society* **157**:403 (1972).
58. G. Courtes, "Optical Systems for UV Space Researches," From *New Techniques in Space Astronomy*, D. Reidel Publishing Company, Dordrecht, Holland, 1971.
59. G. Courtes, P. Cruvellier, M. Detaille, and M. Saisse, *Progress in Optics* vol. XX, 1983, p. 3.

60. D. Gabor, British Patent 544,694, 1941.
61. W. Welford, *Aberrations of Optical Systems*, Adam Hilger Ltd., Bristol, England, 1986.
62. D. Shafer, "Simple Method for Designing Lenses," *SPIE* 237:234 (1980).
63. D. Shafer, "Optical Design with Only Two Surfaces," *SPIE* 237:256 (1980).
64. J. Baker, "Schmidt Image Former with Spherical Abberation Corrector," U.S. Patent 2,458,132, 1945.
65. J. Baker, "Correcting Optical System," U.S. Patent 3,022,708, 1962.
66. J. Houghton, "Lens System," U.S. Patent 2,350,112 (1944).
67. A. Sonnefeld, "Photographic Objective," U.S. Patent 2,141,884, 1938.
68. R. Sigler, "Compound Catadioptric Telescopes with All-Spherical Surfaces," *Applied Optics* 17(10):1519 (1978).
69. D. Hawkins, and E. Linfoot, "An Improved Type of Schmidt Camera," *Monthly Notices Royal Astronomical Society* 105:334 (1945).
70. D. Rao, "Design of a Near Diffraction-Limited Catadioptric Lens," *SPIE* 766:220 (1987).
71. W. Mandler, "Reflecting Mirror and Lens System of the Cassegrain Type," U.S. Patent 2,726,574, 1951.
72. R. Gelles, "A New Family of Flat-Field Cameras," *Applied Optics* 2(10):1081 (1963).
73. M. Shenker, "High Speed Catadioptric Objective in Which Three Corrector Elements Define Two Power Balanced Air Lenses," U.S. Patent 3,252,373, 1966.
74. M. Laiken, *Lens Design*, Marcel Dekker, Inc., New York, 1991.
75. L. Canzek, "High Speed Catadioptric Objective Lens System," U.S. Patent 4,547,045, 1985.
76. M. Amon, "Large Catadioptric Objective," U.S. Patent 3,711,184, 1973.
77. N. Rumsey, "Telescopic System Utilizing Three Axially Aligned Substantially Hyperbolic Mirrors," U.S. Patent 3,460,886, 1969.
78. S. Eisenburg, and E. Pearson, "Two-Mirror Three-Surface Telescope," *SPIE* 751:24 (1987).
79. D. Shafer, "Well Baffled Two-Mirror Corrector for a Parabola," *Applied Optics* 16(5):1175 (1977).
80. D. Shafer, "Four-Mirror Unobscured Anastigmatic Telescopes with All-Spherical Surfaces," *Applied Optics* 17(7):1072 (1978).
81. A. Leonard, "New Horizons for Tilted-Component Telescopes," in *Advanced Telescope Making Techniques*, Willmann-Bell, Inc., Richmond, Va., 1986, p. 110.
82. A. Leonard, "T.C.T.'s (Tilted Component Telescopes)," in *Advanced Telescope Making Techniques*, Willmann-Bell, Inc., Richmond, Va., 1986, p. 131.
83. A. Mackintosh, *Advanced Telescope Making Techniques*, Willmann-Bell, Inc., Richmond, Va., 1986.
84. A. Kutter, "A New Three-Mirror Unobstructed Reflector," *Sky and Telescope*, Jan. 1975, p. 46; Feb. 1975, p. 115.
85. M. Brunn, "Unobstructed All-reflecting Telescopes of the Schiefspiegler Type," U.S. Patent 5,142,417, 1992.
86. P. Manly, *Unusual Telescopes*, Cambridge University Press, Cambridge, England, 1991.
87. D. Shafer, "A Simple Unobstructed Telescope with High Performance," *Telescope Making* 41:4 (TM 41).
88. R. Buchroeder, "Fundamentals of the TST," OSC, University of Arizona, Technical Report 68, May, 1971.
89. R. Gelles, "Unobscured-Aperture Two-Mirror Systems," *Journal of the Optical Society of America* 65(10):1141 (1975).
90. R. Owen, "Easily Fabricated Wide-Angle Telescope," *SPIE* 1354:430 (1990).
91. K. Hallam, B. Howell, and M. Wilson, "Wide-Angle Flat Field Telescope," U.S. Patent 4,598,981, 1986.
92. D. Shafer, "30 Degree F/1.0 Flat Field Three Mirror Telescope," presented Oct. 31 in San Francisco at *Optical Society of America Annual Meeting*, 1978.
93. M. Paul, *Rev. Opt.* 14:169 (1935).
94. J. Baker, "On Improving the Effectiveness of Large Telescopes," *IEEE Transactions AES-5*, (2):261 (1969).
95. A. Meinel, and M. Meinel, *SPIE* 332:178 (1982).
96. A. Meinel, and R. Shack, "A Wide-Angle All-Mirror UV Camera," Optical Sciences Center Tech. Rep. No. 6, 1966.

97. D. Shafer, "20 Degree Field-of-View Strip Field Unobscured Telescope with Stray Light Rejection," presented Oct. 31 in San Francisco at *Optical Society of America Annual Meeting*, 1978.
98. D. Korsch, "Anastigmatic Three-Mirror Telescope," U.S. Patent 4,101,195 (1978).
99. S. Williams, "On-Axis Three-Mirror Anastigmat with an Offset Field of View," *SPIE* **183**:212 (1979).
100. D. Korsch, "Anastigmatic Three-Mirror Telescope," *Applied Optics* **16**(8):2074 (1977).
101. R. Gelles, "Unobscured Aperture Stigmatic Telescopes," *Optical Engineering* **13**(6):534 (1974).
102. W. King, "Unobscured Laser-Beam-Expander Pointing System with Tilted Spherical Mirrors," *Applied Optics* **13**(1):21 (1974).
103. P. Robb, "Three-mirror Telescopes: Design and Optimization," *Applied Optics* **17**(17):2677 (1978).
104. E. Richardson, and C. Morbey, "Optical Design of an Astrometric Space Telescope," *SPIE* **628**:197 (1986).
105. D. Korsch, "Wide-Field Three-Mirror Collimator," U.S. Patent 4,737,021, 1988.
106. D. Pollock, "An All-Reflective, Wide Field of View Collimator," *SPIE* **751**:135 (1987).
107. L. Cook, "Three-Mirror Anastigmatic Optical Systems," U.S. Patent 4,265,510, 1981.
108. L. Cook, "Three-Mirror Anastigmat Used Off-Axis in Aperture and Field," *SPIE* **183**:207 (1979).
109. L. Cook, "The Last Three-Mirror Anastigmat (TMA)?" *SPIE Proceedings* vol. CR41:310 (1992).
110. D. Wang, L. Gardner, W. Wong, and P. Hadfield, "Space Based Visible All-Reflective Stray Light Telescope," *SPIE* **1479**:57 (1991).
111. W. Wetherell, and D. Womble, "All-Reflective Three Element Objective," U.S. Patent 4,240,707, 1980.
112. J. Figosky, "Design and Tolerance Specification of a Wide-Field, Three-Mirror, Unobscured, High-Resolution Sensor," *SPIE* **1049**:157 (1989).
113. D. Korsch, "Two Well-Corrected Four-Mirror Telescopes," *Applied Optics* **13**(8):1767 (1974).
114. E. McCarthy, "Anastigmatic Catoptric Systems," U.S. Patent 3,062,101, 1959.
115. S. Williams, "Design Techniques for WFOV High-Resolution Reflecting Optics," *SPIE* **147**:94 (1978).
116. R. Johnson, J. Hadaway, and T. Bursleson, "All-Reflective Four-Element Zoom Telescope: Design and Analysis," *SPIE* **1354**:669 (1990).
117. R. Altenhof, "The Design of a Large Aperture Infrared Optical System," *SPIE* **62**:129 (1975).
118. A. Offner, "New Concepts in Projection Mask Aligners," *Optical Engineering* **14**(2):130 (1975).

This page intentionally left blank.

DO NOT DUPLICATE

Leo Beiser

*Consultant
Flushing, New York*

R. Barry Johnson

*Consultant
Huntsville, Alabama*

30.1 GLOSSARY

a	aperture shape factor; calibrates diffraction angle as function of aperture intensity distribution
A	area
C	capacitance, electrical
d	grating spacing, array element spacing
D	useful beam width in scan direction (see W); D_m = enlarged beam width due to a
f_e	data bandwidth
f	focal length
F	F-number (f/D)
FL	field lens
FOV	field of view
FWHM	full width (of δ) measured at half maximum intensity
H	vehicle height
I	resolution invariant; adaptation of Lagrange invariant, $I = \Theta D = \Theta' D'$, normalized intensity
k	scanning constant (see m)
m	scan magnification ($d\Theta/d\Phi$) (for $m = \text{constant}$, $m = k = \Theta/\Phi$)
m'	composite magnification
M	optical magnification (image/object)
n	number of facets, refractive index, diffractive order
N	number of resolution elements subtended by Θ or S
\mathcal{N}	total number of cells in a phased array
P	radiant power (watts)
PSF	point spread function (intensity distribution of focused spot)
q	number of cells in array period
Q	q-factor (electromechanical)

r	radius
R	reciprocity failure factor (≥ 1), scanner rpm
s	sensitivity (recording medium)
S	format width, no. of scans per revolution of scanner
t	time
T	time period, optical transfer factor (≤ 1)
V	vehicle velocity, electrical voltage
W	full aperture width, of which D is utilized (see ρ)
w	width at $1/e^2$ intensity
x	along-scan direction
y	cross-scan direction
α	angular departure from normal beam landing, spectral power transfer factor, radiation intensity parameter
γ	phosphor utilization factor (≤ 1)
δ	spot size; if Gaussian, measured across $1/e^2$ intensity width or at FWHM
Δ	shift distance
ϵ	chromatic error, ellipticity
η	duty cycle, conversion efficiency (≤ 1), efficiency
Θ, θ	optical scan angle, angular beamwidth
Φ	mechanical scan angle, along-track optical scan angle, optical field angle
λ	wavelength
Λ	wavelength (acoustic grating), array period
ρ	truncation ratio (W/D)
σ	Gaussian standard deviation
τ	transit time (acoustic), retrace time, dwell time

30.2 INTRODUCTION

This chapter provides an overview of optical scanning techniques in context with their operational requirements. System objectives determine the characteristics of the scanner which, in turn, influence adjacent system elements. For example, the desired resolution, format, and data rate determine the scanner aperture size, scan angle, and speed, which then influence the associated optics. The purpose of this chapter is to review the diverse options for optical scanning and to provide insight to associated topics, such as scanned resolution and the reduction of spatial errors. This broad perspective is, however, limited to those factors which bear directly on the scanner. Referencing is provided for related system relationships, such as image processing and data display. Topics are introduced with brief expressions of the fundamentals. And, where appropriate, historical and technical origins are referenced.

The subject of scanning is often viewed quite differently in two communities. One is classified as *remote sensing* and the other, *input/output scanning*. Associated component nomenclature and jargon are, in many cases, different. While their characteristics are expanded in subsequent sections, it is useful to introduce some of their distinctions here. Remote sensing detects objects from a distance, as by a space-borne observation platform. An example is infrared imaging of terrain. Sensing is usually *passive* and the radiation *incoherent* and often multispectral. Input/output scanning, on the other hand, is *local*. A familiar example is document reading (input) or writing (output). Intensive use of the laser makes the scanning *active* and the radiation *coherent*. The scanned point is focused via finite-conjugate optics from a local fixed source.

While the scanning components may appear interchangeable, special characteristics and operational modes often preclude this option. This is most apparent for diffractive devices such as

acousto-optic and holographic deflectors. It is not so apparent regarding the differently filled scanning apertures, imparting important distinctions in resolution and duty cycle. The unification of some of the historically separated parameters and nomenclature is considered an opportunity for this writing.

A recent concentration of R&D in the field of agile beam steering is presented in Sec. 30.8. Intensive work has yielded encouragement in the long quest for achieving, for example, the performance of articulated large mirrors while avoiding some of the systemic burdens of size, weight, and inertia. Recently, remarkable work has been done using microelectromechanical systems (MEMS) to make scanners on a small scale. Two particularly interesting devices are presented and are the digital micromirror device (DMD) incorporated into digital light processing (DLP) projectors and gimbal-less two-axis scanning-micromirror devices (GSMD). The DMD is essentially a two-dimensional array binary-state scanner or light switches, while the GSMD is a fully analog dual-axis scanner.

System Classifications

The following sections introduce the two principal disciplines of optical scanning, remote sensing, and input/output scanning, in preparation for discussion of their characteristics and techniques.

Remote Sensing The applications for passive (noninvasive) remote sensing scanners are varied and cover many important aspects of our lives. A signature representative of the target is obtained to form a signal for subsequent recording or display. This process is operationally distinct from active scanning, as expressed further in this chapter. Table 1 lists typical applications of these techniques. Clearly, remote scanning sensors can be hand held to satellite-borne.

A variety of scanning methods has been developed to accomplish the formation of image (or imagelike) data for remote sensing. These methods may be roughly divided into framing, pushbroom, and mechanical. Generally stated, frame scanning requires no physical scan motion and implies that the sensor has a two-dimensional array of detectors which are read out by use of electronic means (e.g., CCD), electron beam, or light beam. Such an array requires an optical system that has two-dimensional wide-angle capability. Pushbroom methods typically employ some external means to move the image of a linear array of detectors along the area to be imaged. Mechanical methods generally include one- and two-dimensional scanning techniques incorporating as few as one detector to multiple- detector arrays. As is the case for pushbroom methods, image formation by one-dimensional mechanical scanning requires that the platform containing the sensor (or in

TABLE 1 Representative Applications of Passive Scanning Sensors

Medical	Government
Cancer	Forest fires
Arthritis	Police
Whiplash	Smuggling
Industrial	Search and rescue
Energy management	Military
Thermal fault detection	Gun sights
Electronic circuit	Night vision
detection	Tactical
Nondestructive testing	Navigation
Scientific	Missiles
Earth resources	Strategic
Weather	Aircraft
Astronomy	ICBM
	Surveillance

some cases the object) be moved to create the second dimension of the image. The latter two methods are discussed further in later sections of this chapter.

Mechanical scanners can be configured to perform either one- or two-dimensional scan patterns. In the case of a one-dimensional scanner, the second dimension needed to form an image is most often generated by movement of the sensor platform.

A number of optical scanning systems have been invented to satisfy the wide variety of applications. In the early years of passive scanning systems, the entire optical system was often moved in order to produce either a one- or two-dimensional scan pattern. The advent of airborne mapping or reconnaissance electro-optical systems began during the 1950s. Typically, the scanner performed a one-dimensional scan in object-space, as defined under the section “Object-Space and Image-Space Scanners,” orthogonal to the flight direction of the aircraft, while the motion of the aircraft generated the second dimension of the image. The resultant video information is stored on a recording medium such as photographic film, digital tape, and the like. The design and resultant performance of a scanning system are governed by a variety of parameters that are related by trade-off equations and considerations. The selection of the scanner type typically has a strong influence upon the ultimate system performance. In subsequent discussion, the more important parameters related to the scanner selection will be covered. The complexities of the total system design and optimization are not within the scope of this chapter.

Input/Output Scanning In contrast to remote sensing, which captures passive radiation, active input/output scanning illuminates an object or medium with a “flying spot,” derived typically from a laser source. Some examples appear in Table 2, divided into two principal functions: *input* (detecting radiation scattered from the scanning spot) and *output* (recording or display). Input is modulated by the target to form a signal; output is modulated by a signal.

Some merit clarification. Under input is laser radar—a special case of active remote sensing, using the same coherent and flying-spot scanning disciplines as the balance of those exemplified. Earth resources imaging is the *recording* of remotely sensed image signals. Finally, data/image display denotes the general presentation of information, which could include “hard copy” and/or actively projected and displayed images.

Active scanning is synonymous with flying-spot scanning, the discipline most identified with the ubiquitous cathode-ray tube (CRT). While the utilized devices and their performance differ significantly, the distinctions between CRT and laser radiation are primarily their degrees of monochromaticity and coherence, as addressed later in this chapter.

Thus, most high-resolution and high-speed flying-spot scanning are now conducted using the laser as a light source. This work in input/output scanning concentrates on the control of laser radiation and the unique challenges encountered in deflecting photons, devoid as they are of the electric and magnetic fields accompanying the electron beam. Reference is provided¹⁻³ for pursuit of the CRT scanning discipline.

TABLE 2 Examples of Input/Output Scanning

Input	Output
Image scanning/digitizing	Image recording/printing
Bar-code reading	Color image reproduction
Optical inspection	Medical image outputs
Optical character recognition	Data marking and engraving
Optical data readout	Microimage recording
Graphic arts camera	Reconnaissance recording
Scanning confocal microscopy	Optical data storage
Color separation	Phototypesetting
Robot vision	Graphic arts platemaking
Laser radar	Earth resources imaging
Mensuration	Data/Image display

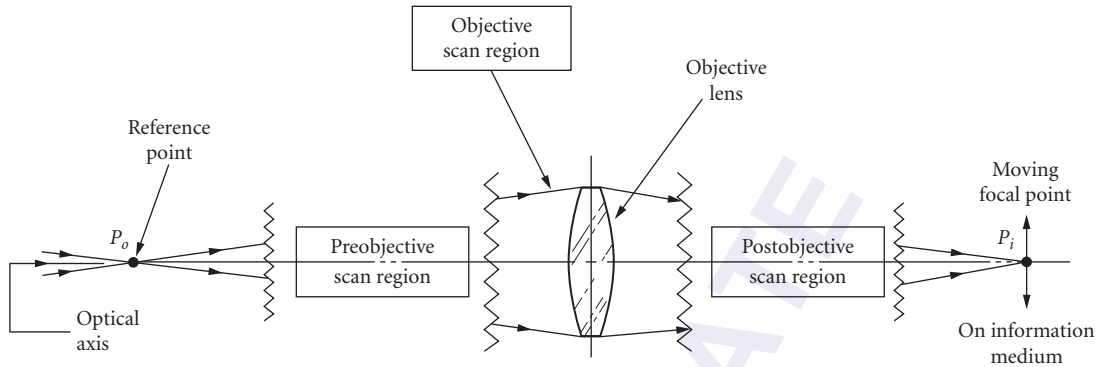


FIGURE 1 Conjugate imaging system, showing scan regions, as determined by position of scanning member relative to the objective lens. Translational (objective) scan and angular (pre/postobjective) scan can occur individually or simultaneously.⁴

Scanner classification Following the nomenclature introduced in the early 1970s,^{4,5} laser scanners are designated as *preobjective*, *objective*, and *postobjective*. Figure 1 indicates the scan regions within a general conjugate optical transfer of a fixed reference (object) point P_o to a moving focal (image) point P_i . The component which provides principal focusing of the wavefront identifies the objective lens.

The scanner can perform two functions (see section “Objective, Preobjective, and Postobjective Scanning” later in this chapter): one is *translation* of the aperture with respect to the information medium. This includes translation of the lens element(s) or translation of the object, or both, and is identified as an *objective* scan. The other is *angular change* of the optical beam with respect to the information medium. Angular scanners are exemplified by plane mirrors on rotating substrates. Although lenses can be added to an angular scanner, it is seldom so configured. The scanner is either preobjective or postobjective. In holographic scanning, however, the hologram can serve as an objective lens and scanner simultaneously.

Radial symmetry and scan magnification A basic characteristic of some angular scanners is identified as *radial symmetry*. When an illuminating beam converges to or diverges from the nodal or rotating axis of an angular scanner, it is said to exhibit radial symmetry.⁶ The collimated beam which is parallel to the rotating axis is a special case of radial symmetry, in which the illuminating beam propagates to or from a very distant point on the axis. Scanners exhibiting radial symmetry provide unity angular optical change for unity mechanical change. That is, $m = d\Theta/d\Phi = 1$, where Θ is the optical scan angle and is the mechanical change. The parameter m is called the *scan magnification*, discussed later under “Augmented Resolution, the Displaced Deflector” for Eq. (19). It ranges typically between 1 and approximately 2, depending on the scanner-illumination configuration, per Table 3. In remote sensing, $m = \Theta/\Phi = k$. (See “Compound Mirror Optics Configurations.”)

The prismatic polygon (see “Monogon and Polygon Scanners”) exhibits a variable m , depending on the degree of collimation or focusing of the output beam. When collimated, $m = 2$. When focusing, the value of m shifts from 2 according to the composite magnification

$$m' = 2 + r/f \quad (1)$$

where f and r are according to Fig. 4 and Eq. (19). This is similar to the ratio of angular velocities of the scanned focal point along the arc of a limaçon,⁵

$$\dot{\Theta}/\dot{\Phi} = 2 \left(1 + \frac{\cos \Phi}{1 + f/r} \right) \quad (2)$$

TABLE 3 Typical Features of Pyramidal and Prismatic Polygon Scanners

Item	Description	Pyramidal	Prismatic
1	Input beam direction ^a	Radially symmetric ^b (typically parallel to axis)	Perpendicular to axis ^c
2	Output beam direction ^a	Arbitrary angle to axis (typically perpendicular)	Perpendicular to axis ^c
3	Scan magnification ^b (scanning constant)	1	2 ^b
3a	Along-scan error magnification	1	2 ^b
3b	Max. scan angle, Θ_{\max} (n = no. of facets)	$2\pi/n$	$4\pi/n^b$
4	Output beam rotation about its axis ^d	Yes	No
5	Aperture shape ^e (overilluminated)	Triangular/keystone	Rectangular
6	Enlargement of along- scan beam width	No	Yes ^f $D_m = D/\cos \alpha$
7	Error due to axial polygon shift ^g	Yes	No
8	Error due to radial polygon shift ^g	Yes	Yes
9	Fabrication cost	Greater	Lower

^aWith respect to rotating axis.

^bSee sections “Radial Symmetry and Scan Magnification” and “Augmenting and Scan Magnification.” Output beam assumed collimated.

^cAll beams typically in same plane perpendicular to axis. See Figs. 26 and 28.

^dObservable when beam is nonisotropic; e.g., elliptic, polarized, etc. Rotation of isotropic beam normally not perceived. See “Image Rotation and Derotation.”

^eSee Table 4.

^f α = angular departure from normal landing. See “Scanner-Lens Relationships.”

^gShift of image focal point in noncollimated light. No error in collimated light.

Note that when $r \rightarrow 0$ or when $f \rightarrow \infty$, $\dot{\Theta}/\dot{\Phi} \rightarrow 2$. In holographic scanners which are not radially symmetric, m depends on the angles of incidence and diffraction of the input and first-order output beam,

$$m = \sin\theta_i + \sin\theta_o = \lambda/d \quad (3)$$

where θ_i and θ_o are the input and diffracted angles (with respect to the grating normal) and d is the grating spacing. For example, when $\theta_i = \theta_o = 30^\circ$, $m = 1$; when $\theta_i = \theta_o = 45^\circ$, $m = \sqrt{2}$.

30.3 SCANNED RESOLUTION

Remote Sensing Resolution and Data Rates

Figure 2 illustrates the scanning geometry for an airborne line-scanning system where the aircraft is flying along a “track” at a height H and velocity V . The total scanned field of view is θ_{\max} and the cross-track and along-track instantaneous fields of view are $\Delta\theta$ and $\Delta\phi$, respectively.* The direction directly below the aircraft and normal to the scanned surface is called the *nadir*. The instantaneous field

*Cross-track and along-track in remote sensing correspond to along-scan and cross-scan, respectively, in input/output scanning.

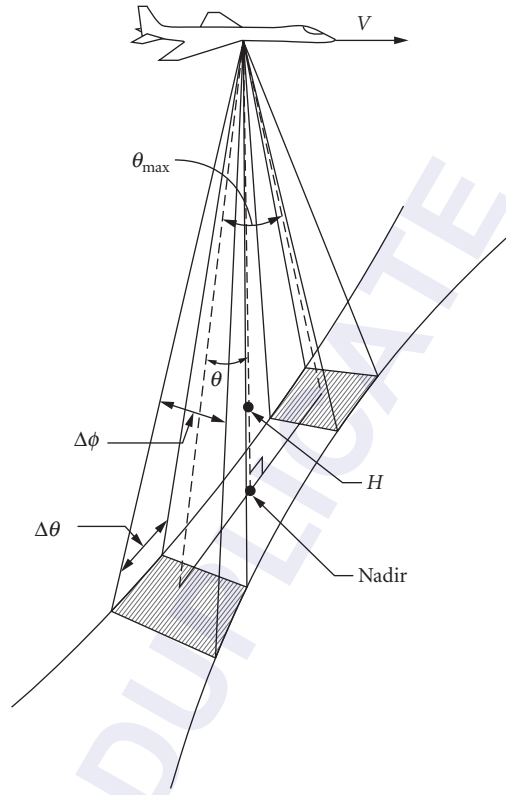


FIGURE 2 Scanning geometry for an airborne line-scanning system with a total scanned FOV of θ_{\max} . The aircraft is flying at height H and velocity V . The cross-track and across-track instantaneous FOVs are $\Delta\theta$ and $\Delta\phi$, respectively. (After Wolfe, Proc. IRE, 1958.)

of view is defined as the geometrical projection of the detector with spatial dimensions d_{ct} and d_{at} by the optics having a focal length of F . Therefore, $\Delta\theta = d_{ct}/F$ and $\Delta\phi = d_{at}/F$. Figure 2 shows the “bow-tie” distortion of the scanning geometry which will be discussed further under “Image Consequences.”

The basic equation relating the aircraft velocity to the angular velocity of the scanning system to produce contiguous scan lines at the nadir is $V/H = \dot{s} \cdot \Delta\phi$, where \dot{s} is the scanning system’s scan rate in scans per second. For a system with n detector elements aligned in the flight direction, $V/H = n\dot{s} \cdot \Delta\phi$.

The number of resolution elements or pixels in a single scan line is

$$N = \frac{\theta_{\max}}{\Delta\theta} \quad (4a)$$

$$= \frac{2\pi\theta_{\max}}{360^\circ \cdot \Delta\theta} \quad (4b)$$

where $\Delta\theta$ is in radians, θ_{\max} is the total field of view measured in radians in Eq. (4a) and in degrees in Eq. (4b), for the scanning means employed, taking due regard for the duty cycle given

in Eq. (23). The scan rate, in scans per second, may be expressed as a function of the scan mirror speed by

$$\dot{s} = \frac{RS}{60} \quad (5)$$

where R is the scan mirror rpm and S is the number of scans produced by the scanning mechanism per revolution of the optics. It follows that the number of resolution elements or pixels per second per scan line is

$$\dot{N} = \frac{2\pi\theta_{\max}RS}{60 \cdot 360 \cdot \Delta\theta} \quad (6)$$

The angle θ_{\max} (in degrees) is determined by the configuration of the scan mirror and is $\theta_{\max} = 360 \cdot k/S$, where k is the scanning constant or scan magnification* and can have values ranging from 1 to 2. The specific value is dependent upon the optical arrangement of the scanner as exemplified in Table 3. The pixel rate may now be written as

$$\dot{N} = \frac{2\pi kR}{60 \Delta\theta} \quad (7)$$

The information retrieval rate of a system is often expressed in terms of the dwell time τ or the data bandwidth f_e as

$$f_e = \frac{1}{2\tau} = \frac{\dot{N}}{2} \quad (8)$$

By combining the preceding equations, the data bandwidth for a multiple-detector system can be expressed as

$$f_e = \frac{\pi k(V/H)}{nS\Delta\theta\Delta\phi} \quad (9)$$

which illustrates clearly the relationship between important system parameters such as f_e being inversely proportional to instantaneous field-of-view solid angle ($\Delta\theta\Delta\phi$).

Input/Output Scanning

Resolution Criteria, Aperture Shape Factor The resolution of an optical scanner is expressed^{5,7} by the number N of spots or elements that can be conveyed along a contiguous spatial path. The path is usually (nearly) linear and traversed with uniform velocity. Although the elements δ are analogous to the familiar descriptors *pixels* or *pels* (picture elements), such identification is avoided, for pixels often denote spatially digitized scan, where each pixel is uniform in intensity and/or color. Active optical scan, on the other hand, is typically contiguous, except as it may be subjected to modulation. Normally, therefore, the scanned spots align and convolve to form a continuous spatial function that can be divided into elements by modulation of their intensity. To avoid perturbation of the elemental point spread function (PSF) by the modulating (or sampling) process, we assume that the scan function is modulated in intensity with a series of (Dirac) pulses of infinitesimal width, separated by a time t such that the spatial separation between spot centers is $w = vt$, where v is the velocity of the

*See "Radial Symmetry and Scan Magnification" regarding scan magnification m which represents a more general form of the scanning constant k .

scanned beam. It is often assumed that the size δ of the thus-established elemental spot corresponds to w ; that is, the width of the imaged spot is equal to the spacing from its neighbor.

To quantify the number N of such spots, those which exhibit a Gaussian intensity distribution are usually considered overlapping at one of two widths; at their $1/e^2$ intensity points, or at their 50 percent intensity points [the latter denoted as full width at half maximum (FWHM)]. Their relationship is

$$\delta_{\text{FWHM}} = 0.589\delta_{1/e^2} \quad (10)$$

The resolution N is identified with its measurement criterion, for the same system will convey different apparent N , per Eq. (10). That is, it will count approximately 1.7 times as many spots at FWHM than at $1/e^2$ intensity.

These distinctions are accommodated⁸ by their aperture shape factors a . For example, the above Gaussian aperture distribution is represented by shape factors

$$a_{1/e^2} = \frac{4}{\pi} = 1.27 \quad (11a)$$

$$a_{\text{FWHM}} = 0.589a_{1/e^2} = 0.75 \quad (11b)$$

When adapted to the applicable equation for spot size

$$\delta = aF\lambda \quad (12)$$

where $F = f/D$ is the F-number of the cone converging over the distance f from beam width D , and λ is the radiation wavelength, the resulting Gaussian spot size becomes

$$\delta_{1/e^2} = \frac{4}{\pi} \frac{f}{D} \lambda = 1.27F\lambda \quad (13a)$$

when measured across the $1/e^2$ intensity points, and

$$\delta_{\text{FWHM}} = 0.75F\lambda \quad (13b)$$

when measured across FWHM.

The factor a further accommodates the resolution changes due to changes in aperture shape, as for apodized and truncated Gaussians. Ultimate truncation is that manifest when the illuminating spatial distribution is much larger than the limiting aperture (overillumination or overfilling), forming the uniformly illuminated aperture.* Familiar analytic examples are the rectangular and round (or elliptic) apertures, which generate (for the variable x) the normalized intensity distributions $[\sin(x)/x]^2$ and $[2J_1(x)/x]^2$, respectively, where $J_1(x)$ is the first-order Bessel function of the first kind.

Figure 3 illustrates[†] the MTFs⁹ of several uniformly illuminated apertures. Their intersections with the 0.5 MTF value identifies the spatial frequency at which their modulation is 50 percent. With the rectangular aperture as a reference (its straight line intersects 0.5 MTF at 50 percent of the limit frequency, forming $a = 1$), the intersections of the others with MTF = 0.5 yield corresponding spatial frequencies and relative a -values. Since the spatial frequency bandpass is proportional to $D/f = 1/F$, the apertures of the others must be widened by their a -values (effectively lowering their F-numbers) to render equivalent response midrange.

*Although the illumination and resulting PSFs are of a coherent wave, scanning forms a sequence of incoherently related intensity measurements of the space-shifting function, yielding an incoherent MTF.

[†]See Figs. 23 and 24 and related discussion for reduced power throughput due to approaching uniform illumination within the aperture.

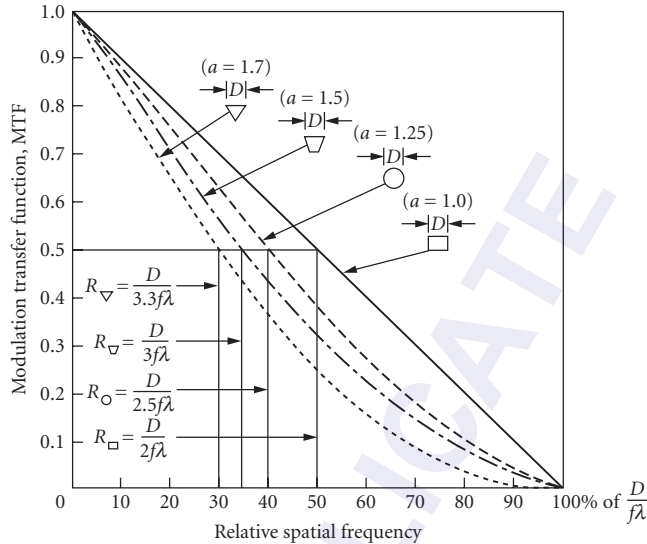


FIGURE 3 Modulation transfer function versus relative spatial frequencies for uniformly illuminated rectangular, round, keystone, and triangular apertures. Spatial frequency at 50 percent modulation (relative to that of rectangular aperture) determines the a value.⁵

Table 4 summarizes the aperture shape factors (a) for several useful distributions.⁸ Truncated, when applied, is two dimensional. Noteworthy characteristics are

1. Scanning is in the direction of the width D or W .
2. The a -value of 1.25 for the uniformly illuminated round/elliptic aperture corresponds closely to the Rayleigh radius value of 1.22.
3. The Gaussian-illuminated data requires that the width D , measured at the $1/e^2$ intensity points be centered within the available aperture W . Two conditions are tabulated: untruncated¹⁰ ($W \geq 1.7D$) and truncation at $W = D$.
4. The Gaussian-illuminated data also provides the a -values for 50 percent MTF, allowing direct comparison with performance of the uniformly illuminated apertures.

This data relates to apertures which, if apodized, are truncated two dimensionally. However, one-dimensional truncation of a Gaussian beam by parallel boundaries is not uncommon, typical of that for acousto-optic scanning. There, the limiting aperture width W is constant, as determined by the device, while the Gaussian width D is variable.^{4,5,11,12} Table 5 tabulates the shape factor a for such conditions.

TABLE 4 Aperture Shape Factor a

Uniformly Illuminated			Gaussian Illuminated		
Shape	$\rightarrow D \leftarrow$	a	δ (Spot Overlap)	a (Untruncated) $W \geq 1.7 D$	a (Truncated) $W = D$
Rectangular		1.0	@ $1/e^2$ Intensity	1.27	1.83
Round/elliptic		1.25	@ $1/2$ -Intensity	0.75	1.13
Keystone		1.5	for 50% MTF	0.85	1.38
Triangular		1.7			
Width D for 50% MTF		1.7		Beam width D @ $1/e^2$ intensity centered within aperture width W	

TABLE 5 Aperture Shape Factor a for One-Dimensional Truncation of a Gaussian Intensity Distribution

Truncation Ratio $\rho = W/D$	Shape Factor a for 50% MTF
0	1.0
0.5	1.05
1.0	1.15
1.5	1.35
2.0	1.75

W = width of aperture.

D = width of Gaussian beam at $1/e^2$ intensity points.

W and D measured in scan direction.

To relate to data in Table 4, the case of $\rho = W/D = 0$ represents illumination through a narrow slit. This corresponds to the uniformly illuminated rectangular aperture, whence $a = 1$. When $\rho = 1$, then $W = D$ and the parallel barriers truncate the Gaussian beam at its $1/e^2$ intensity points. Compared to symmetric truncation, this allows more of the Gaussian skirts to contribute to the aperture width, providing $a = 1.15$ versus 1.38. When $\rho = 2$, the Gaussian beam of half the width of the boundaries is effectively untruncated, halving the resolution ($a = 1.75$ vs. 0.85), but maximizing radiometric throughput. (See Fig. 24, observing nomenclature, where $D = 2w_x$ and $W = 2r_o$.)

Fundamental Scanned Resolution The section on “Input/Output Scanning” introduced the two forms of optical scan: translation and angular deflection. Beam translation is conducted by objective scan, while angular deflection is either preobjective or postobjective. Examples of each are provided later in this chapter.

The resolution N_s of translational scan, by a beam focused to spot size δ executing a scanned path distance S , is simply,

$$N_s = \frac{S}{\delta} \quad (14)$$

Extremely high resolutions are practical, but are often limited to moderate speeds and bandwidths. Common implementations provide $N_s = 3000$ to 100,000.*

The resolution N_θ of angular scan,[†] represented schematically in Fig. 4, capable of much higher speeds, is given by^{4,5,7}

$$N_\theta = \frac{\Theta D_o}{a\lambda} \quad (15)$$

where Θ is the useful deflected optical angle and D_o is the effective aperture width at its nodal center, discussed in the next section. Common implementations provide $N_\theta = 2000$ to 30,000. Equation (15) is independent of spot size δ and dependent only on the aperture characteristics of D_o and a , and the wavelength λ . The beam could be converging, collimated, or diverging. When collimated, $D_o = D$, the actual width of the illuminated portion of the aperture. When converging or diverging, resolution augmentation occurs (see next section).

The numerator of Eq. (9) is a form of the Lagrange invariant,¹³ expressed in this nomenclature as

$$n\Theta D = n'\Theta'D' \quad (16)$$

*A high-resolution laser printer provides $N = 3000$ to 10,000, and a high-resolution graphic arts imager $N = 10,000$ to 100,000.

[†]Derived from Eq. (4a) with $\sin \Delta\Theta \approx \Delta\Theta = a\lambda/D_o$.

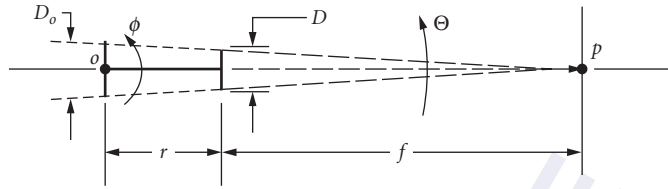


FIGURE 4 Deflecting element of width D (e.g., mirror of polygon) displaced by radius r from axis o , propagating a converging beam to p over focal distance f . Effective larger aperture D_o appears at rotating axis.⁴

where the primed terms are the refractive index, (small) angular deviation, and aperture width, respectively, in the final image space. For the common condition of $n = n'$ in air, the ΘD product and resolution N are conserved, invariant with centered optics following the deflector.

Augmented Resolution, the Displaced Deflector In general, a scanning system can accumulate resolution N by adding the two processes described previously, augmentation of angular scan with linear translation, forming

$$N = N_{\theta} + N_s \quad (17)$$

Augmentation occurs, for example, with conventional multielement scanners (such as polygons) having deflecting elements (facets) which are displaced from the rotating axis by a distance r , and whose output beam is noncollimated. One active element (of width D) and its focused output beam is illustrated in Fig. 4. For convenient analysis,⁶ the deflecting element appears as overilluminated with an incident beam. The resulting resolution equations and focal spot positions are independent of over- or underillumination (see “Duty Cycle”).

Augmentation for increased resolution is apparent in Fig. 4, in which the output beam is derived effectively from a larger aperture D_o which is located at o . By similar triangles, $D_o = D(1 + r/f)$, which yields from Eq. (15),

$$N = \frac{\Theta D}{a\lambda} \left(1 + \frac{r}{f}\right) \quad (18)$$

This corresponds to Eq. (17), for in the N_s term the aperture D executes a displacement component $S \approx r\Theta$, which, with Eq. (12) forms Eq. (14).

Following are some noteworthy observations regarding the parenthetic augmentation term:

1. Augmentation goes to zero when $r = 0$ (deflector on nodal axis) or when $f = \infty$ (output beam collimated).
2. Augmentation adds when output beam is convergent (f positive) and subtracts when output beam is divergent (f negative).
3. Augmentation adds when r is positive and subtracts when r is negative (output derived from opposite side of axis o).

The fundamental or nonaugmented portion of Eq. (18), $N = \Theta D/a\lambda$, has been transformed to a nomograph, Fig. 5, in which the angle Θ is represented directly in degrees. $D/a\lambda$ is plotted as a radius, particularly useful when $a\lambda = 1 \mu\text{m}$, whereupon $D/a\lambda$ becomes the aperture size D , directly in mm. The set of almost straight bold lines is the resolution N . Multiples of either scale yield corresponding multiples of resolution.

Augmenting and scan magnification Equation (18) develops from Fig. 4, assuming that the optimal scan angle Θ is equal to the mechanical angle Θ . This occurs only when the scanner exhibits

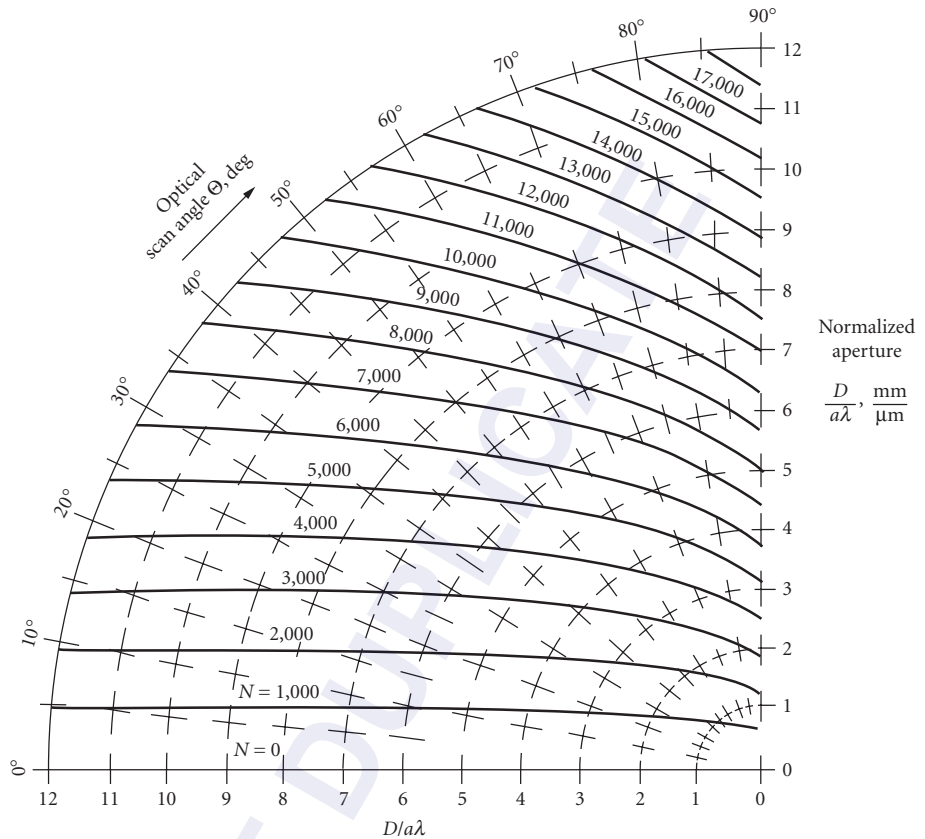


FIGURE 5 Nomograph of resolution equation $n = \frac{\Theta D}{a\lambda}$.

radial symmetry (see “Radial Symmetry and Scan Magnification”). When, however, $m = d\Theta/d\phi \neq 1$, as for configurations represented in the section “Objective, Preobjective, and Postobjective Scanning,” account must be taken of scan magnification m . Thus, the more complete resolution equation is represented by^{6,14}

$$N = \frac{\Theta D}{a\lambda} \left(1 + \frac{r}{mf} \right) \quad (19)$$

where per Fig. 4, Θ = optical scan angle (active)
 D = scan aperture width
 λ = wavelength (same units as D)
 a = aperture shape factor
 m = scan magnification ($= d\Theta/d\phi$)
 Φ = mechanical scan angle about o
 r = distance from o to D
 f = distance from D to p

[see variations for r and f under Eq. (18)].

Considering $m = \Theta/\Phi$ as a constant, another useful form is

$$N = \frac{\Phi D}{a\lambda} \left(m + \frac{r}{f} \right) \quad (20)$$

whose augmenting term shows a composite magnification

$$m' = m + r/f \quad (21a)$$

which, for the typical prismatic polygon becomes

$$m' = 2 + r/f \quad (21b)$$

Duty Cycle The foregoing resolution equations refer to the active portion of a scan cycle. The full scan period almost always includes a blanking or retrace interval. The ratio of the active portion to the full scan period is termed the duty cycle η . The blanking interval can include short overscan portions (straddling the active format), which are used typically for radiometric and timing calibration. The duty cycle is then expressed as

$$\eta = 1 - \tau/T \quad (22)$$

where τ is the blanking interval time and T is the full scan period. A reduced duty cycle increases instantaneous bandwidth for a given average data rate. In terms of the scan angle of polygons, for example, it limits the useful component to

$$\theta = \eta \theta_{\max} \quad (23)$$

where θ_{\max} is the full available scan angle (see Table 3).

Over- and Underillumination (Over- and Underfilling) In overillumination, the light flux encompasses the entire useful aperture. This is usually implemented by illuminating at least two adjacent apertures (e.g., polygon facets) such that the active one is always filled with light flux. This not only allows unity duty cycle, but provides for resolution to be maximized for two reasons: (1) blanking or retrace may be reduced to zero; and (2) the full available aperture width is operative as D throughout scan. The trade-off is the loss of illuminating flux beyond the aperture edges (truncation) and attendant reduction in optical power throughput (see “Coherent Source” under Sec. 30.5, p. 30.26). An alternative is to prescan¹⁵ the light flux synchronously with the path of the scanning aperture such that it is filled with illumination during its entire transit.

In underillumination, the light flux is incident on a portion of the available aperture, such that this subtense delimits the useful portion D . A finite and often substantive blanking or retrace interval results, thereby depleting the duty cycle, but maximizing the transfer of incident flux to the scanned output.

30.4 SCANNERS FOR REMOTE SENSING

Early Single-Mirror Scanners

Early scanning systems comprised an object-space mirror followed by focusing optics and a detector element (or array). The first scanners were simple rotating mirrors oriented typically at 45° to the axis as illustrated in Fig. 6. The rotational axis of the scan mirror lies parallel to the flight direction. In Fig. 6a, the scan efficiency and duty cycle of the oblique or single ax-blade scanner (see monogon under “Monogon and Polygon Scanners”) is quite low since only one scan per revolution ($S = 1$) is generated. The scan efficiency of the wedge or double ax-blade scanner shown in Fig. 6b is twice as great ($S = 2$), although the effective optical aperture is less than half that of the oblique scanner for the same mirror diameter. The scanning constant is $k = 1$ for both types (see “Remote Sensing Resolution and Data Rates” in Sec. 30.3).

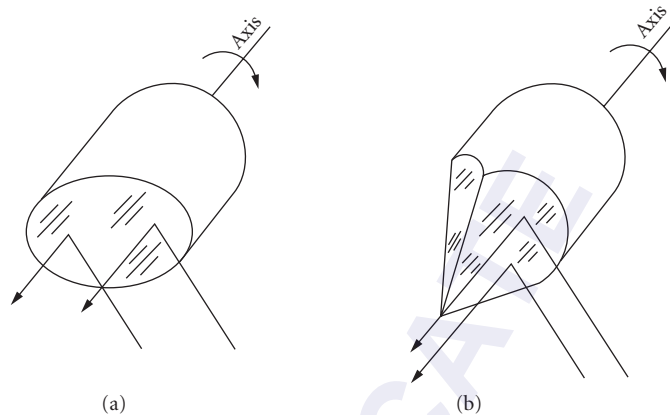


FIGURE 6 Early forms of scanners for remote sensing. The oblique or single ax-blade scanner is shown in (a) and the wedge or double ax-blade is shown in (b).

Compound Mirror Optics Configurations

The aforementioned scanners suffered from a varying optical aperture as a function of view angle. To overcome this difficulty that causes severe variation in the video resolution during a scan line, several new line scanner configurations were developed. Most notable among these was the rotating prism scanner invented by Howard Kennedy¹⁶ in the early 1960s and which forms the basis for most of the produced wide-field-of-view line scanners. Figures 7 and 8 illustrate two

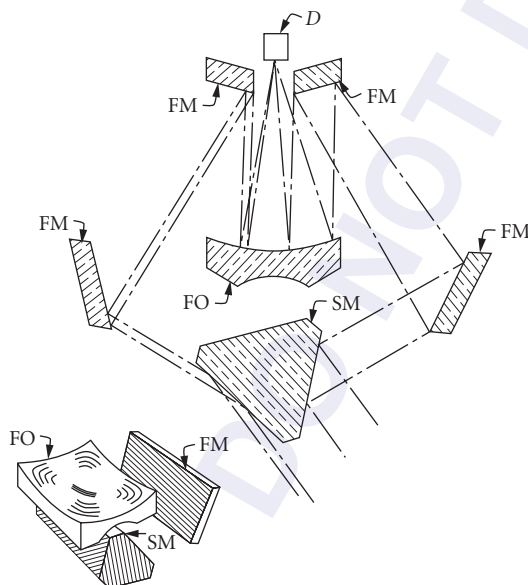


FIGURE 7 Basic split-aperture scanner with a three-sided scan mirror developed in the early 1960s. This scanner features wide FOV, constant optical aperture versus scan angle, and compact mechanical configuration.

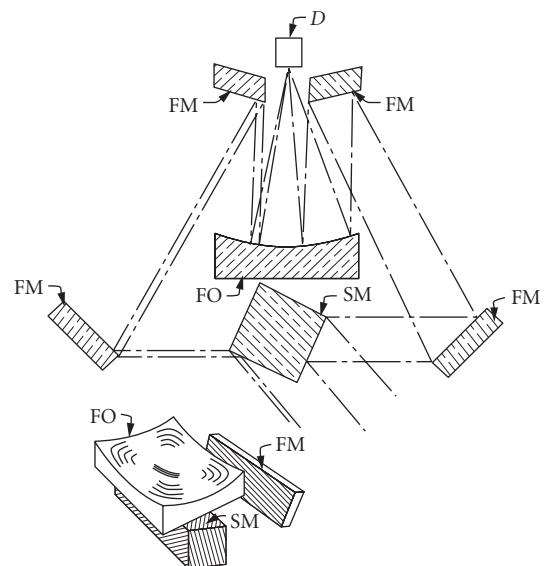


FIGURE 8 Basic split-aperture scanner with a four-sided scan mirror developed in the early 1960s. This scanner features wide FOV, constant optical aperture versus scan angle, and compact mechanical configuration.

configurations of this scanner. The three-sided scan mirror (SM) shown in Fig. 7 rotates about its longitudinal axis at a high rate and the concomitant folding mirrors (FMs) are arranged such that the scanned flux is directed onto suitable focusing optics (FO) which focuses the flux at the detector D . As may be seen in the drawing of a four-sided scanner shown in Fig. 8, the effective optical aperture is split into two portions such that their sum is a constant value as a function of view angle. The width of each portion varies as the view angle is changed, with the portions being of equal value at the nadir position. The isometric view in Fig. 8 shows a portion of the scanner comprising the scan mirror, one folding mirror, and the focusing mirror. For this design, the number of scans per rotation of the scan mirror is equal to the number of faces on the scan mirror, and the scanning constant is $k = 2$, which is also known as optical doubling (see item 3 of the prismatic polygon in Table 3). Also, two faces of the scan mirror are always used to form the total optical aperture. Another advantage of this scanner configuration is that it produces a compact design for the total scanner system, a major reason for its popularity for over a quarter of a century.

Image Consequences

In airborne sensing, it is reasonable to assume that the earth is flat beneath the aircraft. When viewing along the nadir, the detector spatial footprint on the ground is $H \Delta\theta$ and $H \Delta\phi$ in the across- and along-track directions, respectively. As the view angle (θ) moves away from the nadir, the geometric resolution on the ground changes as illustrated in Fig. 2, which creates the bow-tie pattern. In the cross-track direction, it is easily shown that the footprint dimension is $H \Delta\theta \cdot \sec^2\theta$, while in the along-track direction, the footprint dimension is $H \Delta\phi \cdot \sec\theta$. The change in footprint as a function of view angle can be significant. For example, if $\theta_{\max} = 120^\circ$, then the footprint area at the extremes of the scan line is about eight times greater than at the nadir.

Image Relation and Overlap

When a linear array of n detectors is used, it is easily seen that the image of the detector array rotates by exactly the same amount as the view angle if the scanner is pyramidal as shown in Fig. 6. No such rotation occurs for the prismatic polygon, as in the Kennedy scanner, for which each scan comprises n adjacent detector footprints on the ground that form a segmented bow tie. The next scan footprint has significant overlap with the preceding scan(s) for $\theta \neq 0$. A means to compensate for the radiometric difficulties caused by the overlap of scans has been developed.¹⁷ In a single detector system, this artifact is easily compensated by electronic means.

Rotating Wedge Scanner

Figure 9 shows a simple rotating wedge scanner that allows the generation of a wide variety of scan patterns, including a line scan. By controlling the rotational rates and phasing of the wedges, such patterns as included in Fig. 10 can be realized.¹⁸

Circular Scan

In some cases, a circular scan pattern has found utility. Typically, the entire optical system is rotated about the nadir with the optical axis inclined at an angle ψ to the nadir. Figure 11 depicts an object-plane scanner showing how the aircraft or satellite motion creates contiguous scans. Although the duty cycle is limited, an advantage of such a scanner is that, at a given altitude, the footprint has the same spatial size over the scanned arc.

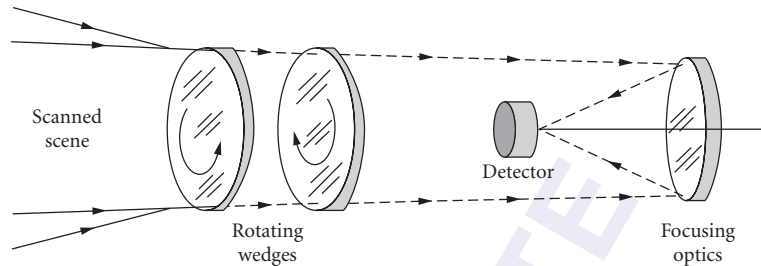


FIGURE 9 Basic geometry of a simple rotating wedge scanner. A wide variety of scan patterns can be produced by controlling the rotational rates and phasing of the wedges. In principal, the detector can view any point within the circular scanned FOV of the scanner. Figure 10 presents typical scan patterns for constant rotational rates. Two-dimensional raster scans can be generated if general positional control of each prism is allowed.

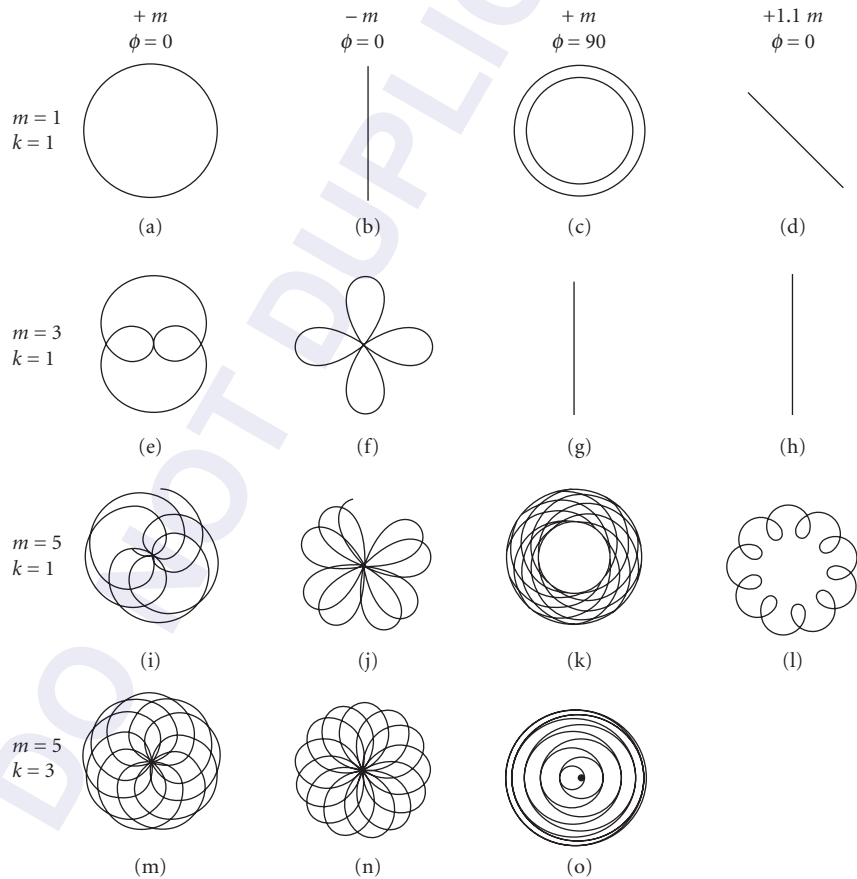


FIGURE 10 Typical scan patterns for a rotating wedge scanner. The ratio of the rotational frequencies of the two prisms is m , the ratio of the prism angles is k , and the phase relation at time zero is ϕ ($\phi = 0$ implies the prism apices are oriented in the same direction). A negative value of m indicates that the prisms are counter-rotating. (After Ref. 18, Fig. 12, p. 12.)

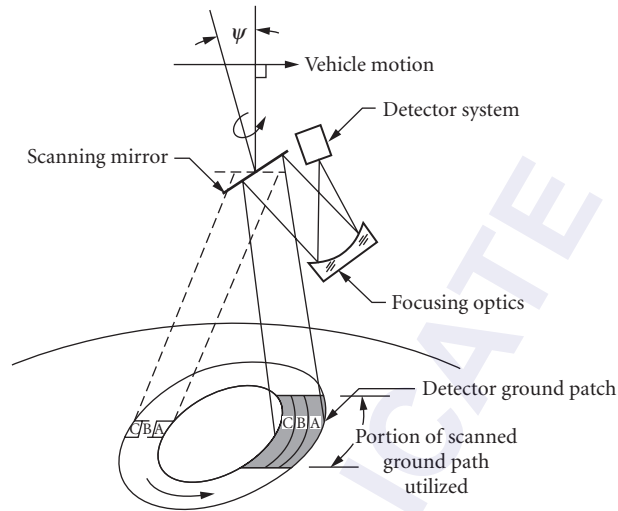


FIGURE 11 Basic configuration of a circular or conical scanner. The normal of the scanning mirror makes an angle ψ with the nadir. The scan pattern of the ground forms a series of arcs illustrated by scans A, B, and C. (After Ref. 19, Fig. 8-3, p. 340.)

Pushbroom Scan

A pushbroom scanner comprises typically an optical system that images onto the ground a linear array of detectors aligned in the cross-track direction or orthogonal to the flight direction. The entire array of detectors is read out every along-track dwell time which is $\tau_{at} = \Delta\phi/(V/H)$. Often, when a serial read-out array is employed, the array is rotated slightly such that the read-out time delay between detectors creates an image that is properly aligned to the direction of motion. Some state-of-the-art arrays can transfer the image data in parallel to a storage register for further processing. The principal advantage of the pushbroom scanner is that no moving parts are required other than the moving platform upon which it is located.

Two-Dimensional Scanners

Two-dimensional scanners have become the workhorses of the infrared community during the past two decades even though line scanners still find many applications, particularly in the area of earth resources. Scanners of this category can be classified into three basic groups, namely, object-space scanner, convergent-beam or image-space scanner, and parallel-beam or intermediate space scanner. Figure 12 depicts the generic form of each group.

Object-Space and Image-Space Scanners

The earliest two-dimensional scanners utilized an object-space scan mechanism. The simplest optical configuration is a single flat-mirror (see Fig. 12a) that is articulated in such a manner as to form a raster scan. The difficulty with this scan mechanism is that movement of a large mirror with the necessary accuracy is challenging. The size of the mirror aperture when in object space must be greater than that of the focusing optics. By using two mirrors rotating about orthogonal axes, the scan can be generated by using smaller mirrors, although

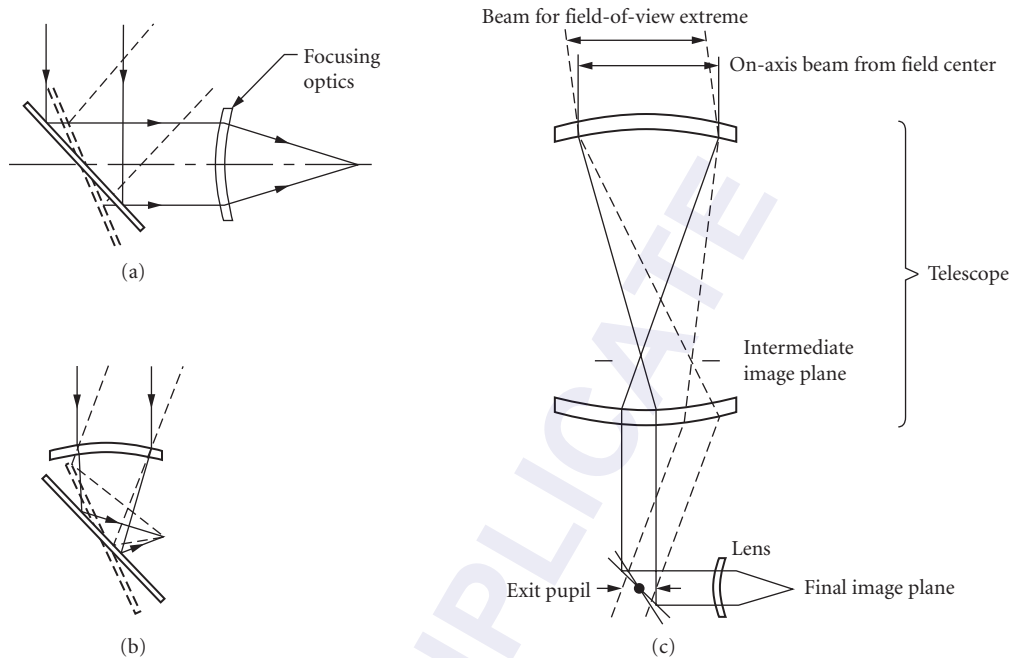


FIGURE 12 The three basic forms of two-dimensional scanners are shown in (a), (b), and (c). The object-space scanner in (a) comprises a scan mirror located in object space where the mirror may be moved in orthogonal angular directions. Image-space or convergent-beam scanner in (b) forms a spherical image surface due to motion of the scan mirror (unless special compensation motion of the scan mirror is provided). The parallel-beam or intermediate-space scanner is illustrated in (c). It is similar to the scanner in (a) except that the scan mirror is preceded by an afocal telescope. By proper selection of the afocal ratio and FOV, the scan mirror and focusing lens become of reasonable size. The scan mirror forms the effective exit pupil. (After Ref. 20, Figs. 7.1 and 7.10.)

the objective optics must have the capability to cover the entire field of view rather than the FOV of the detector. Figure 13 illustrates such a scanner²¹ where mirror SM1 moves the beam in the vertical direction at a slow rate while mirror SM2 generates the high-speed horizontal scan. Although the focusing optics is shown preceding the image-space scan mirrors, the optics could be placed following the mirrors which would then be in object space. Although the high-F-number or low-numerical-aperture focusing lens before the mirrors must accommodate the FOV, it allows the use of smaller mirror facets. The left-hand side of Fig. 13 shows an integral recording mechanism that is automatically synchronized to the infrared receptor side. This feature is one of the more notable aspects of the configuration and sets the stage for other scanner designs incorporating the integrated scene and display scanner. A disadvantage of this scanner is the large size and weight of the vertical scan mirror, in part, to accommodate both scene and display scan.

A variation of the two-mirror object-space scanner is known as the discoid scanner, which produces a raster scan at TV rates.²² Figure 14 depicts the scanner configuration which uses a high-speed, multiple-facet scan mirror SM1 to generate the horizontal scan and a small, oscillating flat mirror SM2 to produce the vertical scan. An advantage of this scanner is that only a single detector is needed to cover the FOV, although a linear array oriented in the scan direction is sometimes used, with time-delay integration, to improve sensitivity. A feature of the “paddle” mirror scanner is the maintenance of a relatively stable aperture on the second deflector without the use of relay optics (see Figs. 21 and 32 and the section on the “Parallel-Beam Scanner”).

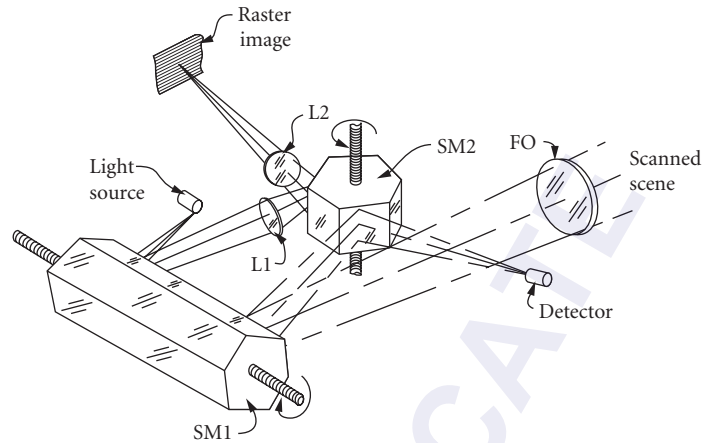


FIGURE 13 Early slow-scan-rate, image-space scanner where the flux from the scanned scene is imaged by the focusing objective lens FO onto the detector. The scene is scanned by mirrors SM1 (vertical) and SM2 (horizontal). A raster image of the scene is formed by imaging the light source using lenses L1 and L2. The display image and the scanned scene are synchronized by using the same pair of mirrors. The light source is modulated by the output of the detector.

Figure 15 depicts a reflective polygon scanner that generates the high-speed horizontal scan (per facet) by rotation of mirror SM about its rotational axis and the vertical movement of the scan pattern by tilting the spinning mirror about pivots P1 and P2 using cam C and its follower CF.²³ The path of the flux from the object reflects from the active facet A of the scan mirror to the folding mirror FM to the focusing mirror FO back through a hole in mirror FM to the detector located in dewar D. Almost

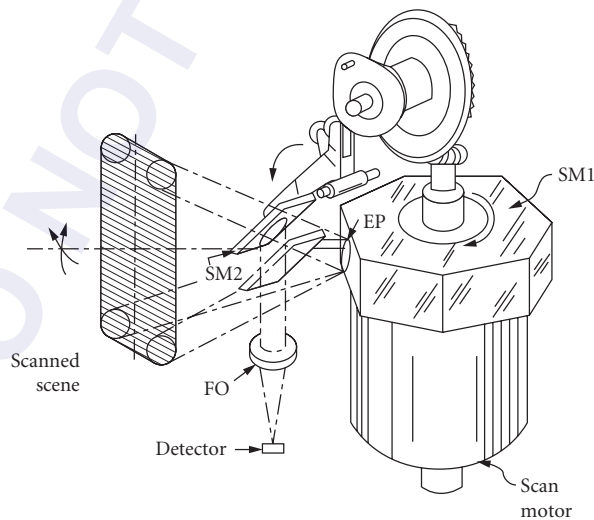


FIGURE 14 Real-time, object-space scanner that has a compact geometry. The exit pupil EP is located on the active facet of the horizontal scan mirror SM1.

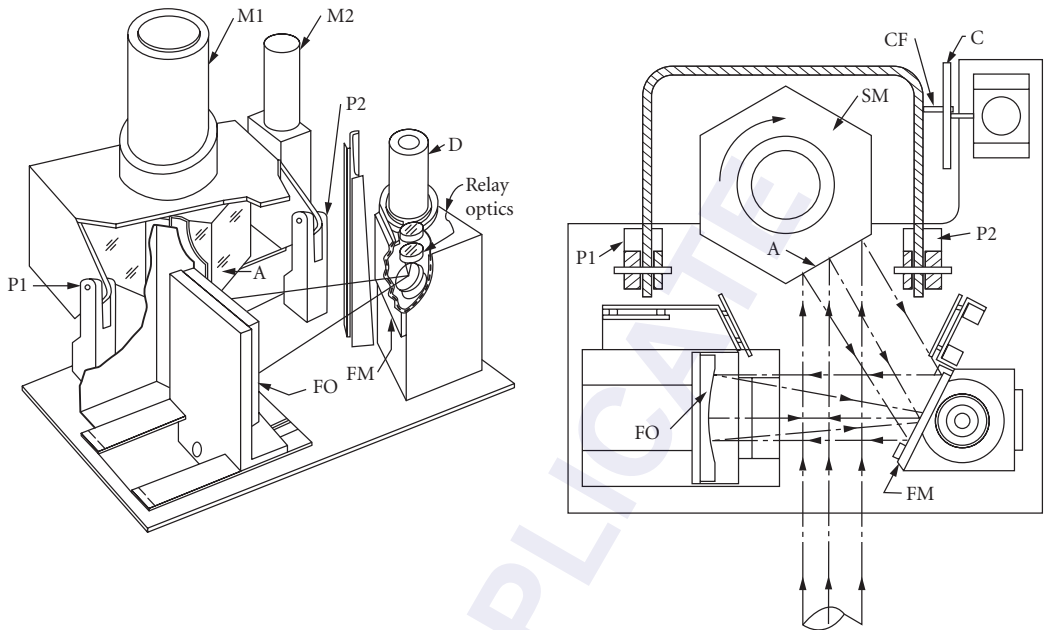


FIGURE 15 Object-space scanner that generates a raster scan using a single mirror SM which is driven by motor M1. The vertical motion of the mirror is accomplished by tilting the housing containing the scan mirror SM about pivots P1 and P2 using the drive mechanism comprising motor M2, cam C, and CF. The FOV of scanners of this type can exceed 30°.

all scanners of this type exhibit scanned-field distortion; that is the mapping of object to image space is nonrectilinear (e.g., see the target distortion discussion in the section “Image Consequences” on p. 30.16).

In general, convergent-beam, image-space scanners suffer from severe defocus over the scanned field of view due to the field curvature produced by the rotation of the scan mirror in a convergent beam. The use of this type scanner is therefore rather limited unless some form of focus correction or curved detector array is employed. A clever invention by Lindberg^{24,25} uses a high-speed refractive prism and a low-speed refractive prism to create the scanned frame. Figure 16 shows the basic

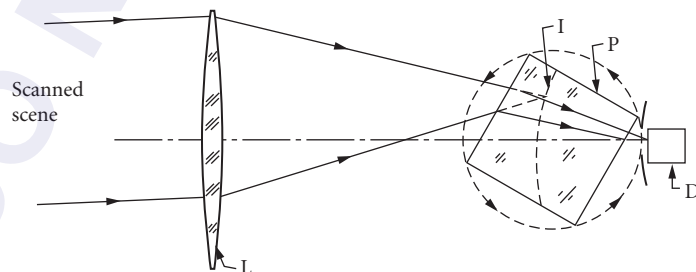


FIGURE 16 Basic configuration of a refractive prism scanner. The scan is generated by rotation of the prism. As shown, four scans are produced per complete rotation of the prism. By proper selection of the refractive index of the prism, reasonably wide FOV can be realized. Although only one prism is shown, a second prism can be included to produce the orthogonal scan direction.

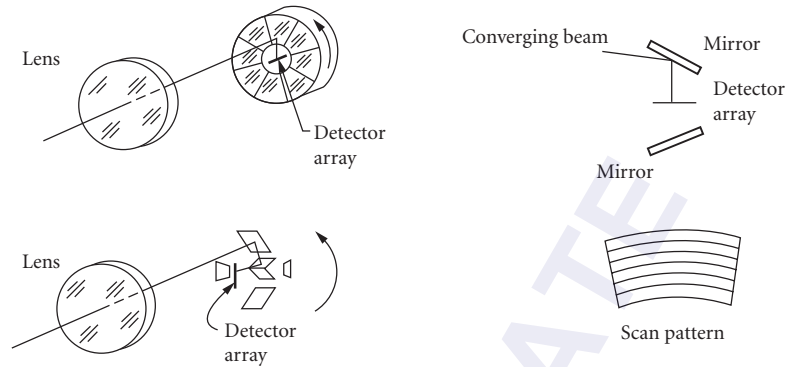


FIGURE 17 Rotating reflective or “soupbowl” scanner. (After Ref. 20, Fig. 7.21, p. 309.)

configuration for a one-dimensional scanner where the cube P is rotated about the axis orthogonal to the page. By proper selection of the refractive index and the geometry of the prism, focus is maintained over a significant and useful field of view. As can be seen from the figure, flux from the object at a given view angle is focused by lens L onto surface I which is then directed to the detector D by the refraction caused by the rotated prism. Numerous commercial and military thermographic systems have utilized this principle for passive scanning. Since the field of view, maximum numerical aperture, optical throughput, and scan and frame rates are tightly coupled together, such scanners have a reasonably constrained design region.

Other image-space scanners used in a convergent beam are the “soupbowl” and carousel scanners.²⁰ The soupbowl scanner shown in Fig. 17 uses a rotating array of mirrors to produce a circularly segmented raster scan. The mirror facets may be at the same angle to generate more frames per rotation, given a detector array that has adequate extent to cover the field of view. The facets could also be tilted appropriately with respect to one another to produce contiguous segments of the field of view if a small detector array is employed. Figure 18 illustrates the configuration of the carousel scanner which uses an array of mirrors arranged such that they create essentially a rectangular scan of the field of view. Another scanning means that has been used for certain forward-looking infrared systems (FLIRs) was to mechanically rotate a detector array of angular extent Φ about the

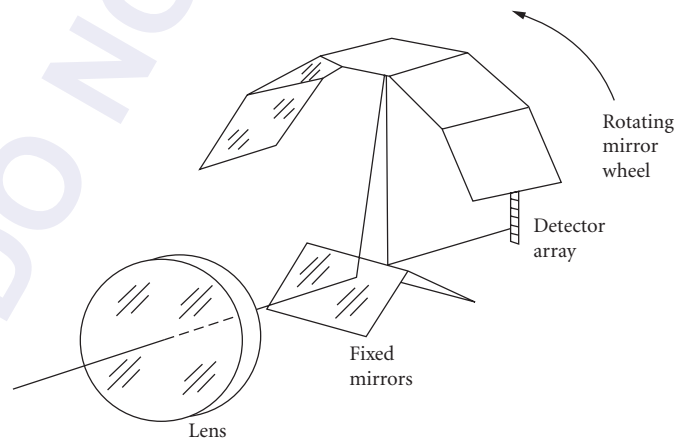


FIGURE 18 Rotating reflective carousel scanner. (After Ref. 20, Fig. 7.22, p. 309.)

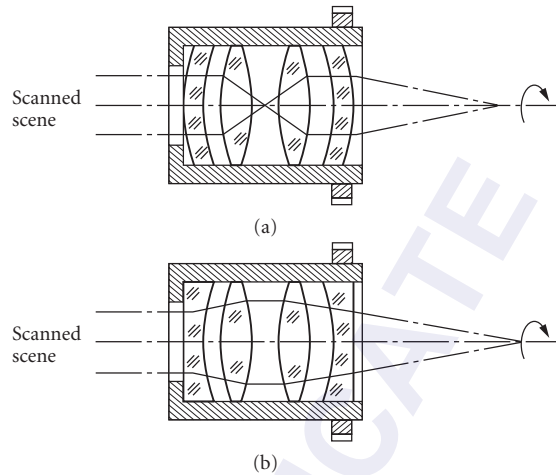


FIGURE 19 “Windshield wiper” scanner. The circular scan is generated by rotating the anamorphic optics about its optical axis. The detector array is located radially and offset from the optical axis. Two scans are produced for each complete rotation of the optics. The lens is shown in (a) at rotation angle θ and in (b) at rotation angle $\theta + 90^\circ$. Although the lens shown is focal, the lens could be afocal and followed by a focusing lens.

optical axis of the focusing optics such that one end of the array was located an angular distance Φ_{os} from the optical axis. The rotating action generated a circular arc scan pattern similar to that of a windshield wiper. The inner radius of the scan pattern is Φ_{os} and the outer radius is $\Phi_{os} + \Phi$. Clearly, the scan efficiency is rather poor and the necessity to use slip rings or the equivalent to maintain electrical connections to the detector array complicated acceptance of this scanner. The windshield wiper scan can also be generated by rotating only the optics if the optics incorporates anamorphic elements. A pair of cylindrical lenses placed in an afocal arrangement, as illustrated in Fig. 19 at rotational angles θ and $\theta + 90^\circ$, will rotate the beam passing through it at twice the rotational rate of the optics.²⁶ See “Image Rotation in Derotation” in Sec. 30.6.

Multiplexed Image Scanning

With the advent of detector arrays comprising significant numbers of elements, the use of a single scan mirror became attractive. Figure 20 presents the basic parallel-beam scanner configuration used for the common module FLIR and thermal night sights. The flat scan mirror SM is oscillated with either a sawtooth or a triangular waveform such that the detector array D (comprising 60, 120, or 180 elements) is scanned over the FOV in the azimuthal direction while the extent of the detector covers the elevation FOV. Since the detectors are spaced on centers two detector widths apart, the scan mirror is tilted slightly in elevation every other scan to produce a 2:1 interlaced scan of the FOV. As shown in Fig. 20, the back side of the scan mirror is used to produce a display of the scanned scene by coupling the outputs of the detectors to a corresponding array of LEDs which are projected to the user’s eye by lenses L1, L2, L3, and L4.

Parallel-Beam Scanner

A more complex two-dimensional, parallel-beam scanner configuration of the type shown in Fig. 12c has been developed by Barr & Stroud and is illustrated in Fig. 21, which incorporates an oscillating mirror SM1, a high-speed polygon mirror SM2 driven by motor M2, and relay optics L1. (See discussion

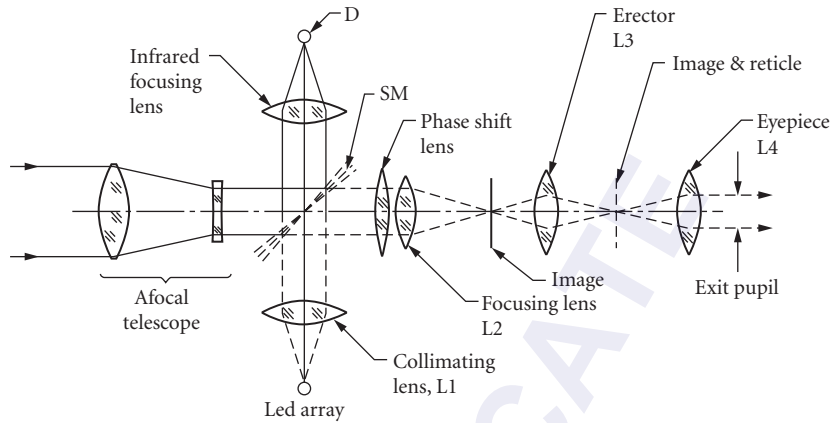


FIGURE 20 Basic configuration of the common module scanner. The front side of the flat scan mirror SM is used to direct the flux from the scanned scene to the detector D through the focusing lens. The outputs from the infrared detector array are used to modulate a corresponding LED array. The LED array is collimated by L1 and scanned over image space by the back side of SM. Lenses L2–L4 are used to project the image of the LED array to the observer’s eye.

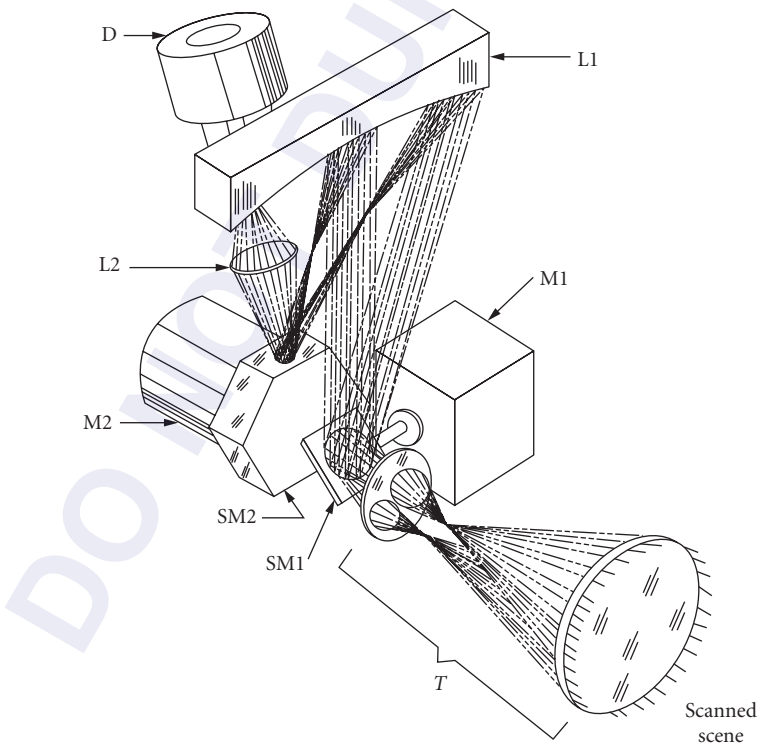


FIGURE 21 Compact real-time scanner. The horizontal scan mirror SM2 is shown in two positions to illustrate how the field of view is related to location on mirror L1. Mirror L1 serves as a relay mirror of the pupil on mirror SM1 to SM2.

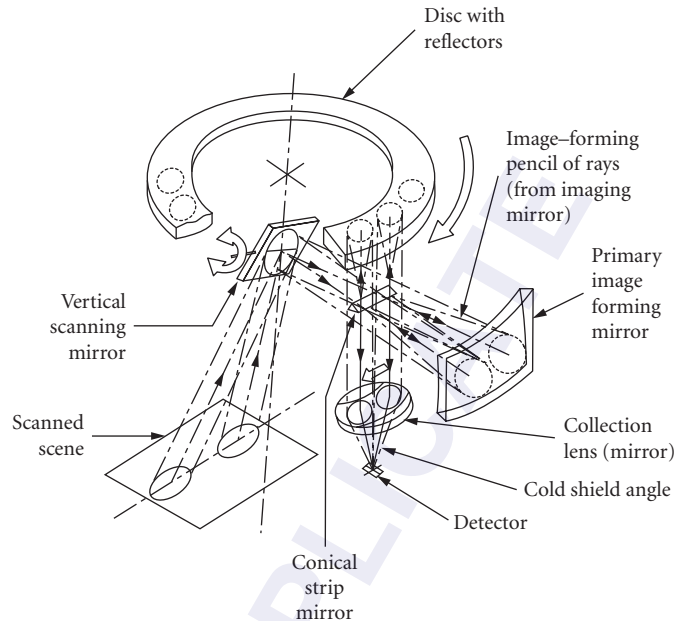


FIGURE 22 Extremely compact, real-time scanner. Diamond-turned optics are used throughout this configuration. (After Ref. 28, Fig. 1.)

at end of Sec. 30.5.)²⁷ An afocal telescope is located before the scanner to change the FOV, as is typical of parallel-beam scanners. Another innovative and compact two-dimensional scanner design by Kollmorgen is depicted in Fig. 22 and features diamond-turned fabrication technology for ease of manufacture and alignment of the mirrors and mounts.²⁸ Another parallel-beam scanner that uses a simple scan mirror has been developed.²⁹ The scan mirror is multifaceted with each facet tilted at an angle that positions the detector array in a contiguous manner in elevation. By having the nominal tilt angle of the facets be 45° to the rotation axis, minimal scanned-field distortion is realized.

30.5 SCANNING FOR INPUT/OUTPUT IMAGING

Power Density and Power Transfer

Incoherent Source This topic merits introduction as the predecessor to laser scanning—cathode-ray tube (CRT), flying-spot scanning and recording.^{1,2,3} Adaptation to other forms of incoherent sources, such as light-emitting diodes (LEDs) will be apparent. Similarities and contrasts with the handling of coherent sources are expressed.

In a CRT, the electron beam power P (accelerating voltage beam current) excites a phosphor of conversion efficiency η and utilization factor γ . The resulting radiant power is transferred through an imaging system of optical transmission efficiency T and spectral power transfer α to a photosensitive medium of area a during a time t . The resulting actinic energy density is given by³

$$E = \frac{\eta \alpha T \gamma P t}{A} \text{ J/cm}^2 \quad (24)$$

[1 J (joule) = 1 W-s (watt-sec) = 10^7 ergs].

The first four terms are transfer factors (≤ 1) relating to the CRT, but adaptable to other radiant sources. They are determined³ for a principal group of CRT recording phosphors having varying processes of deposition and aluminizing, and for two typical (silver halide) photosensitive spectral responses: noncolor sensitized and orthochromatic. The spectral transfer term α is determined from the relatively broad and nonanalytic spectral characteristics of the CRT phosphors and the photosensors,

$$\alpha \cong \frac{\sum_{i=1}^n \frac{P_i}{P_{\max}} \cdot \frac{S_i}{S_{\max}} \Delta \lambda_i}{\sum_{j=1}^m \frac{P_j}{P_{\max}} \Delta \lambda_j} \quad (j \geq i) \quad (25)$$

where the P s and the S s are the radiant power and medium sensitivity, respectively, taken at significant equal wavelength increments $\Delta \lambda$.

The optical transfer term T is composed of three principal factors, $T = T_r T_f T_v$, where T_r is the fixed transmission which survives losses due to, for example, reflection and scatter, T_f is the fixed geometric transfer, and T_v is the spectrally variable transmission of, for example, different glass types. The fixed geometric transfer is given by³

$$T_f = \frac{\cos^4 \Phi V_\Phi}{1 + 4F^2(M+1)^2} \quad (26)$$

The numerator (≤ 1) is a transfer factor due to field angle Φ and vignetting³⁰ losses, F is the lens F-number, and M is the magnification, image/object. The variable component T_v requires evaluation in a manner similar to that conducted for the α . The resulting available energy density E is determined from Eq. (24) and compared to that required for satisfactory exposure of the selected storage material.

Coherent Source Determination of power transfer is much simplified by utilization of a monochromatic (single-line laser) source. Even if it radiates several useful lines (as a multispectral source), power transfer is established with a finite number of relatively simple determinations. Laser lines are sufficiently narrow, compared to the spectral characteristics of most transmission and detection media, so that single point evaluations at the wavelengths of interest are usually adequate. The complexity due to spectral and spatial distributions per Eqs. (25) and (26) are effectively eliminated.

In contrast to the incoherent imaging system described above, which suffers a significant geometric power loss represented by T_f of Eq. (26), essentially all the radiant power from the laser (under controlled conditions discussed subsequently) can be transferred to the focal spot. Further, in contrast to the typical increase in radiating spot size with increased electron beam power of a CRT, the radiating source size of the laser remains essentially constant with power variation. The focused spot size is determined (per the section "Resolution Criteria, Aperture Shape Factor") by the converging beam angle or corresponding numerical aperture or F-number, allowing for extremely high power densities. Thus, a more useful form of Eq. (24) expresses directly the laser power required to irradiate a photosensitive material as

$$P = \frac{sR}{T} \left(\frac{A}{t} \right) \text{watts} \quad (27)$$

where s = material sensitivity, J/cm²
 R = reciprocity failure factor, ≥ 1
 T = optical throughput efficiency, ≤ 1
 A = exposed area, cm²
 t = time over area, A

The reciprocity failure factor R appears in Eq. (27), since the exposure interval t (by laser) can be sufficiently short to elicit a loss in sensitivity of the photosensitive medium (usually registered by silver halide media). If the A/t value is taken as, for example, an entire frame of assumed uniform

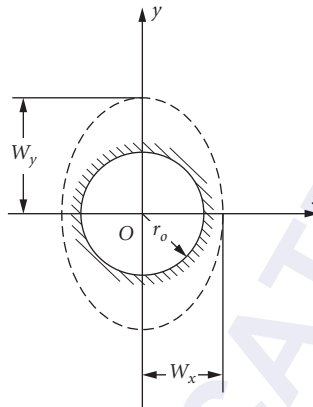


FIGURE 23 Irradiance of a single-mode laser beam, generalized to elliptical, centered within a circular aperture of radius r_o .³¹ Glossary as published; D (as used here) = $2w_x$ and w (as used here) = $2r_o$.

exposure interval (including blanking), then the two-dimensional values of η must appear in the denominator, for they could represent a significant combined loss of exposure time.

The optical throughput efficiency T is a result of loss factors including those due to absorption, reflection, scatter, diffraction, polarization, diffraction inefficiency in acousto-optic and holographic elements, and beam truncation or vignetting. Each requires disciplined attention. While the radiation from (fundamental mode) laser sources is essentially conserved in traversing sufficiently large apertures, practical implementation can be burdensome in scanners. To evaluate the aperture size consistent with throughput power transfer, Figs. 23 and 24 are useful. The data is generalized to

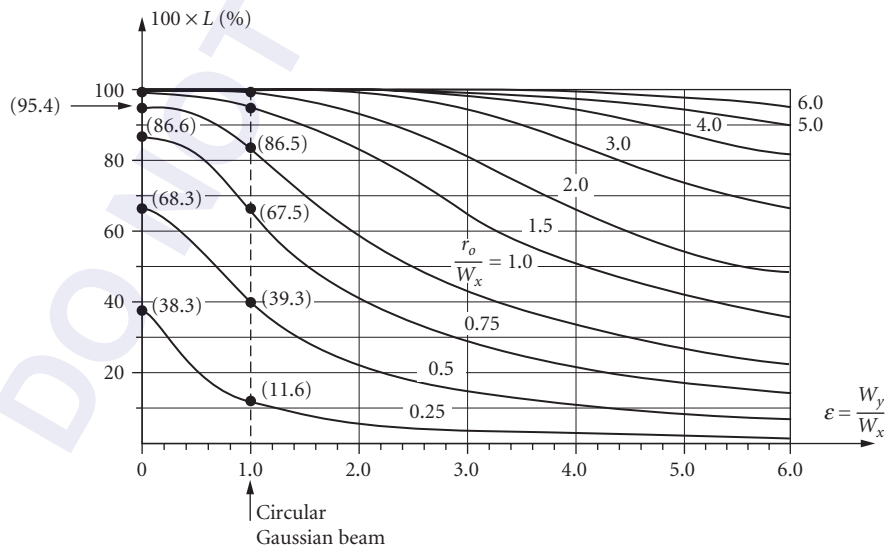


FIGURE 24 Variations of the encircled energy $100 \times L$ (%) versus the ellipticity ϵ and the ratio r_o/w_x as a parameter.³¹ Glossary as published; D (as used here) = $2w_o$ and w (as used here) = $2r_o$.

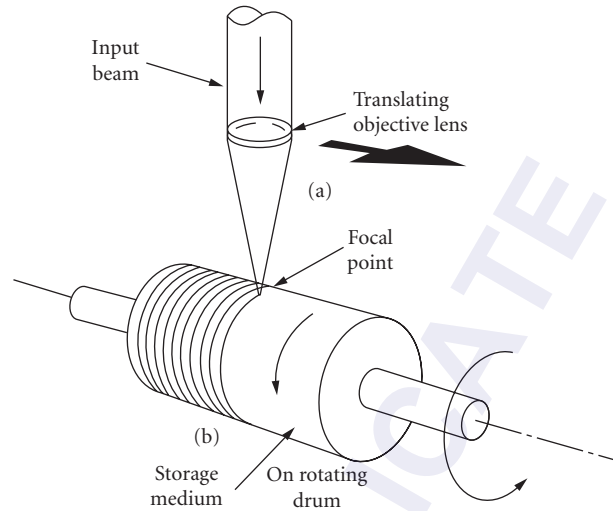


FIGURE 25 Drum configuration executing two forms of objective scan: (a) lens and its focal point translate with respect to storage medium and (b) storage medium translates with respect to lens during drum rotation.⁴

elliptic, accommodating the irradiance of typical laser diodes.^{31,32} Figure 23 shows an irradiance distribution having ellipticity $\varepsilon = w_x/w_y$ ($w @ 1/e^2$ intensity) apertured by a circle of radius r_o . Figure 24 plots the encircled power (percent) versus the ellipticity, with the ratio r_o/w_x as a parameter. When $\varepsilon = 1$, it represents the circular Gaussian beam. Another parameter closely related to this efficiency is the aperture shape factor (discussed previously) affecting scanned resolution. (Note: $D = 2w_x$ and $w = 2r_o$.)

Objective, Preobjective, and Postobjective Scanning

Classification Characteristics The scanner classifications designated as preobjective, objective, and postobjective were introduced previously and represented in Fig. 1 as a general conjugate optical transfer. This section expresses their characteristics.

Objective scan (transverse translational) Translation of an objective lens transverse to its axis translates the imaged focal point on the information surface. (Axial lens translation which optimizes focus is not normally considered scanning.) Translation of the information medium (or object) with respect to the objective lens forms the same effect, both termed objective scan.* The two forms of objective scan appear in Fig. 25, the configuration of a drum scanner.

Preobjective scan (angular) Preobjective scan can provide a flat image field.† This is exemplified by angularly scanning a laser beam into a flat-field or $f-\theta$ lens,³³ as illustrated in Fig. 26, an important technique discussed further under “Pyramidal and Prismatic Facets” and “Flat-Field Objective Optics.”

*Objective scan is limited in speed because the translating lens elements or storage medium must execute the desired final scan velocity, limited by the articulation of relatively massive components.

†Applies beyond the small scan angle $\theta \approx \arctan \theta$, which may be considered linear. Also, no separate dynamic focus to aid forming a flat field.

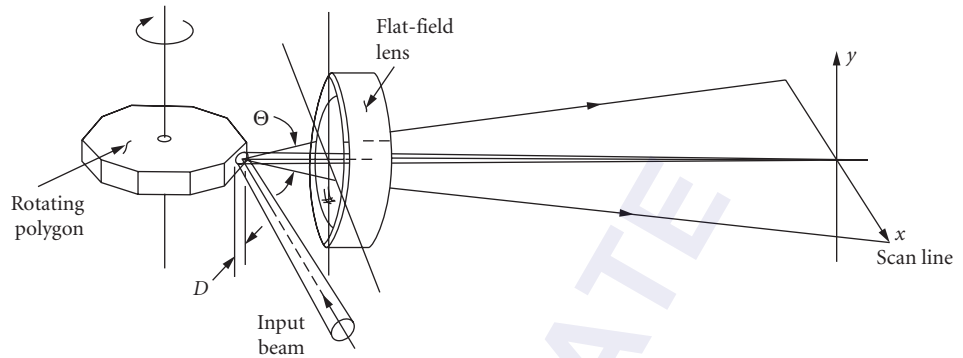


FIGURE 26 Polygon preobjective scan. The rotating polygon reflects and scans the input beam through the angle Θ . The flat-field lens transforms this Θ change to (nominally) linear x displacement along a straight scan line. The input beam and the scanned beam reside nominally in the same plane.⁴

Postobjective scan (angular) Postobjective scan which is radially symmetric per Fig. 27 generates a perfectly circular scan locus.* Departure from radial symmetry (e.g., focal point not on the axis of Fig. 27) generates noncircular (e.g., limaççon⁵) scan, except for when a postobjective mirror with its surface on its axis generates a perfectly circular scan locus, illustrated in Fig. 28. The input beam is focused beyond the axis at point o . Scan magnification $m = 2$.

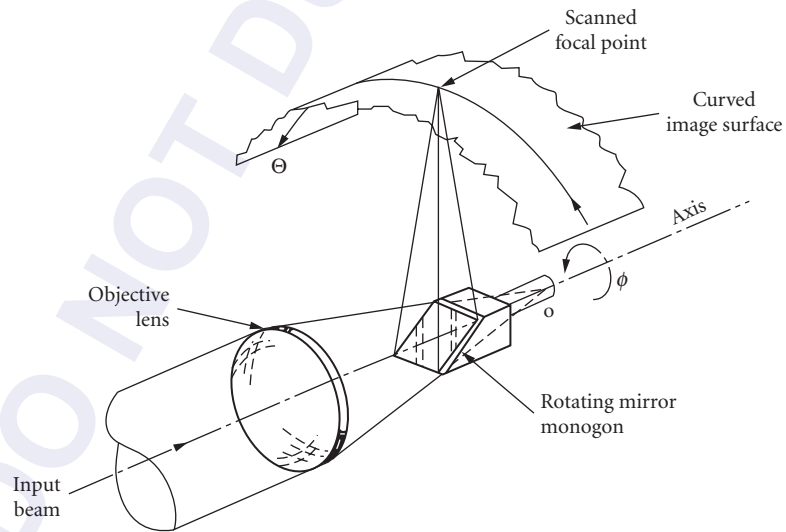


FIGURE 27 Monogon postobjective scan. Generates a curved image. When input beam is focused on the axis (at point o), then system becomes radially symmetric, and image locus forms section of a perfect circle.⁴

*Placing the objective lens in the output beam and coupling it rigidly to the scanner (e.g., of Fig. 27) maintains the same characteristic. This is identified as objective scan (angular). The scanner and lens may be combined, as in a hologram.

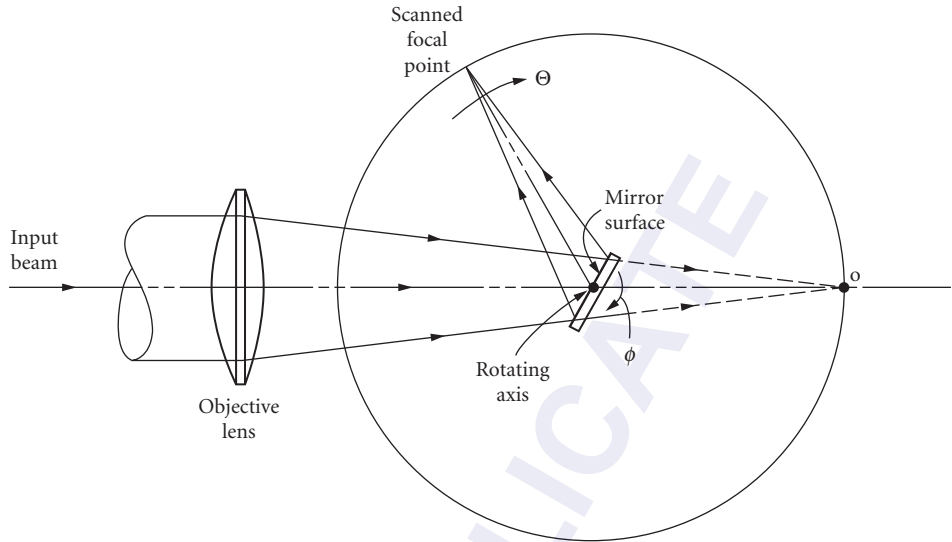


FIGURE 28 Postobjective scan with mirror surface on rotating axis. Input beam is focused by objective lens to point o , intercepted by mirror and reflected to scanning circular locus. Scan magnification $m = d\Theta/d\phi = 2^4$.

Objective Optics

The objective lens converges a scanned laser beam to a moving focal point. The deflector can appear before, at, or after the lens, forming preobjective, objective, and postobjective scanning, respectively (see previous discussion).

On-Axis Objective Optics The simplest objective lens arrangement is one which appears before the deflector, as in Fig. 27, where it is required only to focus a monochromatic beam on-axis. The (postobjective) deflector intercepts the converging beam to scan its focal point. Ideally, the process is conducted almost aberrationlessly, approaching diffraction-limited performance. Since the lens operates on-axis only (accommodates no field angle), if the F-number of the converging cone is sufficiently high (see “Resolution Criteria, Aperture Shape Factor”), it can be composed of a single lens element. This simple arrangement can scan a perfectly circular arc [see “Postobjective Scan (Angular)”], the basis for the elegance of the internal drum scanner and the requirement for adapting the information medium to a curved surface.

Flat-Field Objective Optics Almost all other lens configurations are required to form a flat field by transforming the angular scan to a straight line.³³ The deflector appears before the lens—preobjective. The most common configuration is similar to that of Fig. 26, as detailed further in “Design Considerations” under “Monogon and Polygon Scanners,” in which the scanned beam is shown collimated. Application is not limited to polygon scanners. Similar lenses adapt to a variety of techniques, including galvanometer, acousto-optic, electro-optic, and holographic scanners. The lens must accept the scanned angle θ from the aperture D and converge the beam throughout the scanned field to a best-focus along a straight-line locus. Depending on the magnitudes of θ and D , the F-number of the converging cone and the desired perfection of straight-line focus and linearity, the lens assembly can be composed of from 2 to 7 (or more) elements, with an equal number of choices of index of refraction, 4 to 14 (or more) surfaces, and 3 to 8 (or more) lens spacings, all representing the degrees of freedom for the lens designer to accommodate performance. A typical arrangement of three elements is illustrated in Fig. 34.

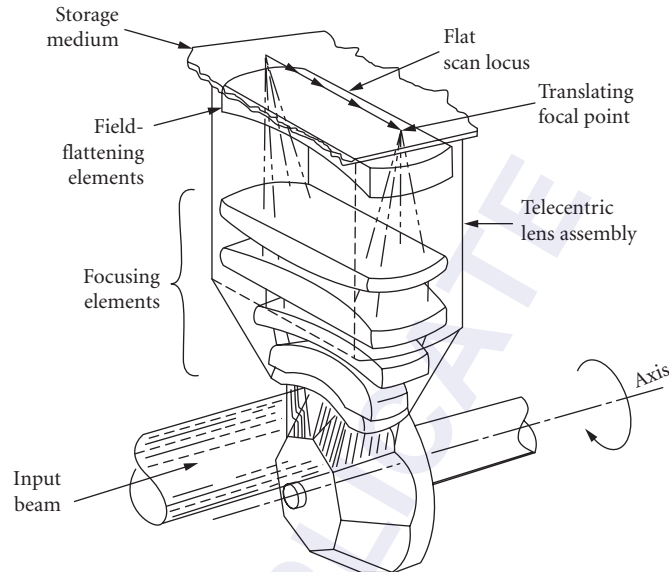


FIGURE 29 High-performance telecentric lens at output of pyramidal polygon scanner, see Fig. 34. (Lens elements shown cut for illustration.)⁴

Telecentricity A more demanding arrangement is illustrated in Fig. 29, showing six elements forming a high-performance scan lens in a telecentric configuration.³⁰ Telecentricity is represented schematically in Fig. 30, in which an ideal thin-lens element depicts the actual arrangement of Fig. 29. Interposed one focal length f between the scanning aperture D (entrance pupil) and the flat image surface, the ideal lens transforms the angular change at the input to a translation of the output cone. The chief ray of the ideal output beam lands normal to the image surface. The degree of telecentricity is expressed by the angular departure from normal landing. Telecentricity is applied typically to restrict the spread of the landing beam and/or to retroreflect the probing beam efficiently

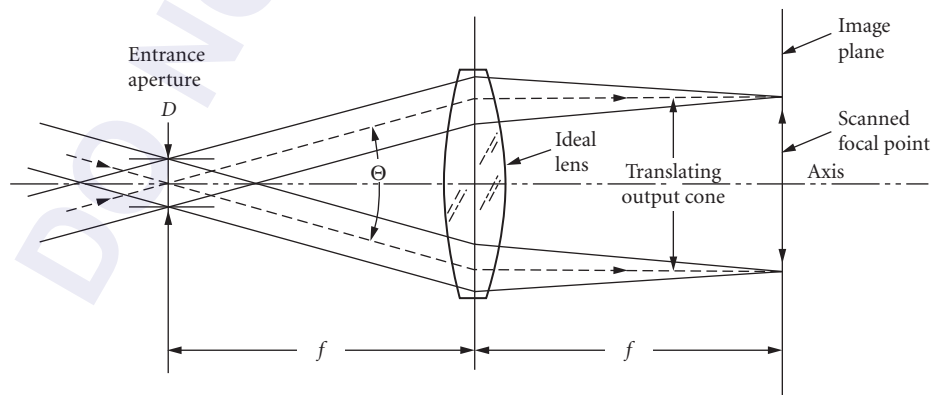


FIGURE 30 Telecentric optical system. Schematic illustration of ideal lens transforming angular scan Θ from aperture D to translational scan landing normal to the image plane.⁴

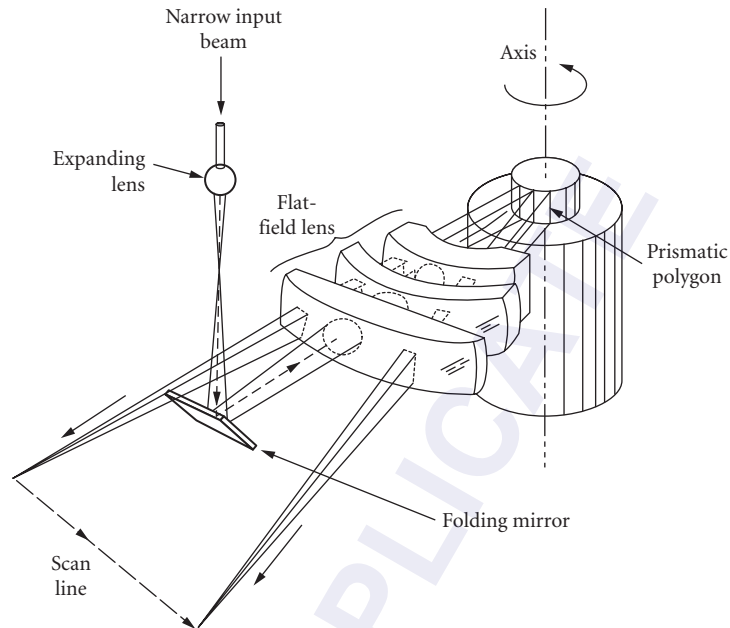


FIGURE 31 Prismatic polygon in double-pass configuration. Narrow input beam is focused by positive lens, expanded, and picked-off by folding mirror to be launched through flat-field lens (in reverse direction). Beam emerges collimated, illuminates facets and is reflected and scanned by rotating polygon to be reconverged to focus on scan line. Input and output beams are slightly skewed above and below lens axis to allow clear separation by folding mirror. (Flat-field lens elements shown cut for illustration.)⁴

for internal system calibration. This facility comes dearly, however, for the final lens elements must be at least as wide as the desired scan format. A further requirement is the need to correct the non-linearity of the simple system of Fig. 30, in which the spot displacement is proportional to the tangent of the scan angle, rather than to the angle directly. As in all scan lenses, compensation to make displacement proportional to scan angle is termed the f - θ correction.

Double-pass and beam expansion Another variation of the objective lens is its adaptation to double-pass,^{4,5,33} as depicted in Fig. 31. The lens assembly serves two functions: first, as the collimating portion of a lenticular beam expander³³ and second, upon reflection by the scanner, as a conventional flat-field lens. This not only provides compaction, but since the illuminating beam is normal to the undeflected facet, the beam and facet undergo minimum enlargement, conserving the size of the deflector. A slight skew of the input and output planes, per Fig. 31, avoids obstruction of the input and scanned beams at the folding mirror. An alternate input method is to make the lens array wide enough to include injection of the input beam (via a small mirror) from the side; at an angle sufficiently off-axis to avoid obstruction of the reflected scanned beam.³³ This method imposes an off-axis angle and consequential facet enlargement and beam aberration, but allows all beams to remain in the same plane normal to the axis, avoiding the (typically) minor

⁴Beam expansion/compression can also be achieved nonlenticularly with prisms.³⁴ Introduction in 1964 of the phrase "beam expander" by Leo Beiser, and its dissemination to generic form, is summarized in App.1 of Ref. 6.

scanned bow which develops in the aforementioned center-skewed method. Other factors relating to increased surface scatter and reflection need be considered.

The requirement for beam expansion noted here is fundamental to the formation of the aperture width D which provides a desired scanned resolution. Since most gas lasers radiate a collimated beam which is narrower than that required, the beam is broadened by propagating it through an inverted telescope beam expander, that is, an afocal lens group having the shorter focal length followed by the longer focal length.* Operation may be reversed, forming *beam compression*, as required. In the previously described double-pass system (Fig. 31), the objective lens provides the collimating portion (long-focal-length group) of a beam expander.

Conservation of Resolution A most significant role of objective optics following the scanner is its determination of the integrity of scanned format, *not of scanned resolution*, as discussed under "Input/Output Scanning". Denoting N as the total number of scanned elements of resolution to be conveyed over a full format width, in first analysis, N is *invariant* with intervening ideal optics. In reasonably stigmatic systems, the lens determines the *size* of the spots, *not their total number*. The number of spots is determined *at the deflector*, whether it be galvanometer, acousto-optic, electro-optic, polygonal, holographic, phased array, or any other angular scanner. This invariance is expressed as

$$I = \Theta D = \Theta' D' \quad (28)$$

an adaptation of the Lagrange invariant [see "Fundamental Scanned Resolution," Eq.(16)], which is illustrated effectively with telescopic operation. If the scanned beam is directed through a telescope (beam compression), as in Fig. 32, the demagnification of f_2/f_1 reduces D to D' , but also expands θ to θ' by the same paraxial factor, sustaining resolution invariance. If D were the deflecting aperture width and L_1 were its objective lens (telecentric in this case), then the image along surface S would exhibit the same number of N spots as would appear if the output beam from D' were focused by another objective lens to another image plane. This schematic represents, effectively, a pupil-transferring optical relay.^{4,5} (See Chap. 17, "Lenses" and Chap. 18, "Afocal Systems" in this volume.)

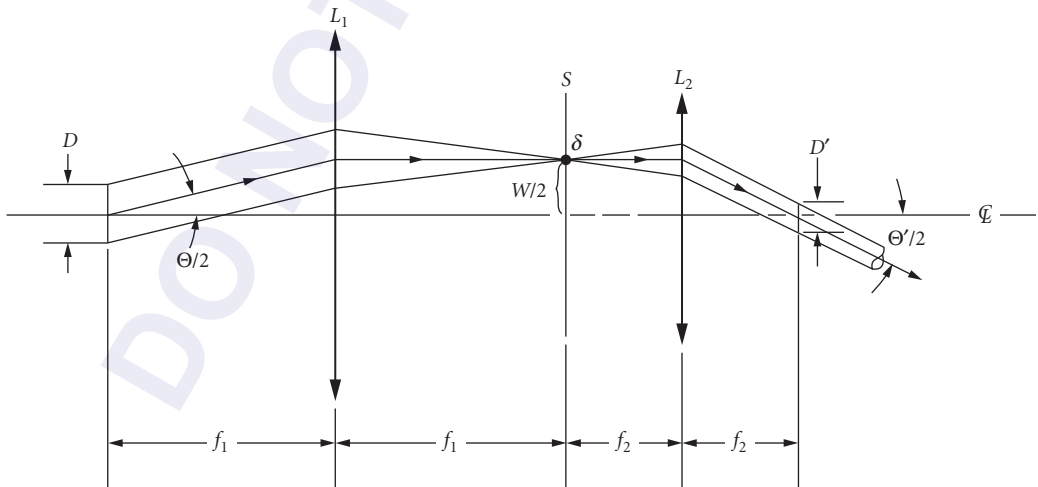


FIGURE 32 Illustration of invariance $I = \Theta D = \Theta' D'$ with telescopic transfer of scanned angle Θ from aperture D .⁴

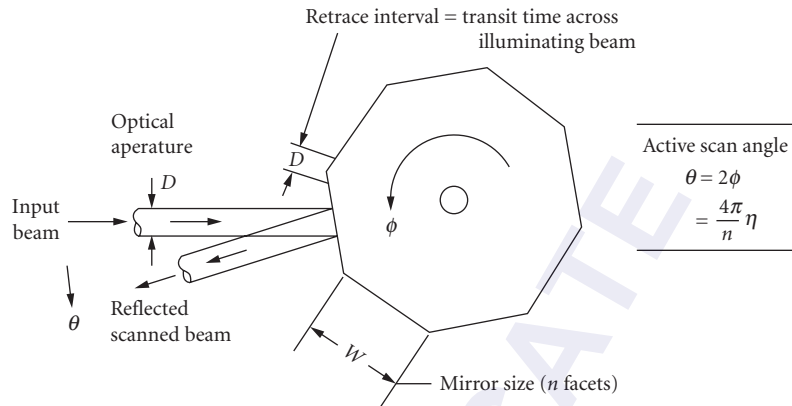


FIGURE 33 Prismatic polygon (underilluminated). Input beam perpendicular to axis. Its width D illuminates a portion of facet width W . Rotation through angle Φ yields scanned angle $\Theta = 2\Phi$, till facet corner encounters beam. Scan remains inactive for fraction D/W , yielding duty cycle $\eta = 1 - D_m/W^4$

30.6 SCANNER DEVICES AND TECHNIQUES

Many of the techniques addressed here for input/output imaging apply equally to remote sensing. Their reciprocal characteristic can be applied effectively by reversing the positions (and ray directions) of the light source(s) and detector(s). Preobjective and postobjective scanning have their counterparts in object-space and image-space scanning. A notable distinction, however, is in the option of underillumination or overillumination of the deflecting aperture in input/output imaging, while the aperture is most often fully subtended in collecting flux from a remote source. This leads to the required attention to aperture shape factor in input/output imaging, which is less of an issue in remote sensing. Another is the need to accommodate a relatively broad spectral range in remote sensing, while input/output operation can be monochromatic, or at least polychromatic. The frequent use of reflective optics in both disciplines tends to normalize this distinction.

Monogon and Polygon Scanners

The rotating mirrored polygon is noted for its capacity to render high data rate at high resolution. It is characterized by a multiplicity of facets which are usually plane and disposed in a regular array on a shaft which is rotatable about an axis. When the number of facets reduces to one, it is identified as a monogon scanner.

Pyramidal and Prismatic Facets Principal arrangements of facets are termed *prismatic* (Fig. 33) or *pyramidal* (see Fig. 34 and “Scanner-Lens Relationship”). Figure 27 is a single-facet pyramidal equivalent, while Fig. 28 is a single-facet prismatic equivalent (common galvanometer mount).

The prismatic polygon of Fig. 33 is oriented typically with respect to its objective optics in a manner shown in Fig. 26, while Fig. 34 shows the relationship of the pyramidal polygon to its flat-field lens. The pyramidal arrangement allows the lens to be oriented close to the polygon, while, as in Fig. 26, the prismatic configuration requires space for clear passage of the input beam.* Design consideration for this most popular arrangement is provided later in this chapter.

*In remote sensing, applying reciprocity, this is the *detected beam*.

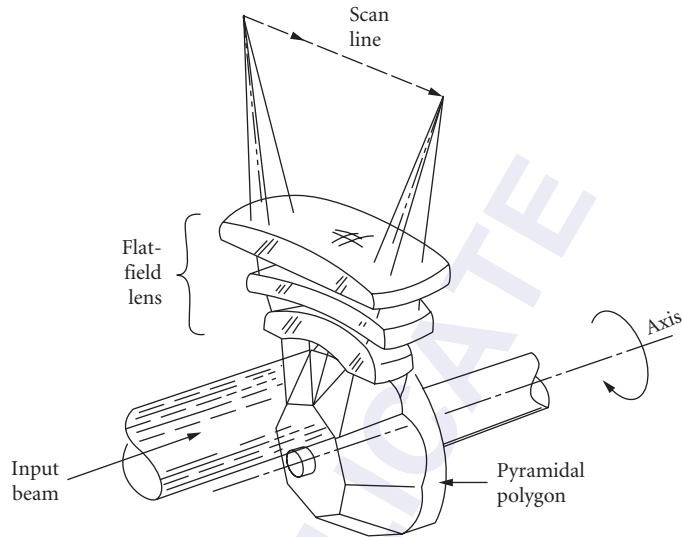


FIGURE 34 Pyramidal polygon (overilluminated). Input beam, parallel to axis, is scanned (by 45° pyramidal angle) in plane perpendicular to axis. When input beam illuminates two facets (as shown), one facet is active at all times, yielding (up to) 100 percent duty cycle. Facet width is full optical aperture D , minimizing polygon size for a given resolution, but wasting illumination around unused facet regions. Can operate underilluminated (per Fig. 33) to conserve throughput efficiency, but requires increased facet width (and polygon size) to attain high-duty cycle. (Flat-field lens elements shown cut for illustration.⁴)

Table 3 lists significant features and distinctions of typical polygon scanners. Consequences of item 3, for example, are that the scan angle of the prismatic polygon is twice that of the pyramidal one for a given rotation. To obtain equal scan angles θ of equal beam width D (equal resolutions N) and to provide equal duty cycle (see “Augmenting and Scan Magnification”) at equal scan rates, the prismatic polygon requires twice the number of facets, is almost twice the diameter, and rotates at half the speed of the pyramidal polygon. The actual diameter is determined with regard for the aperture shape factor (previously discussed) and the propagation of the beam cross section (pupil) across the facet during its rotation (see “Design Considerations”).

Image Rotation and Derotation When a beam having a round cross section is focused to an isotropic point spread function (psf), the rotation of this distribution about its axis is typically undetectable. If, however, the psf is nonisotropic (asymmetric or polarized), or if an array of 2 or more points is scanned to provide beam multiplexing,⁶ certain scanning techniques can cause an undesired rotation of the point and the array of points on the image surface.

Consider a monogon scanner, per Fig. 27. As shown, the input beam overilluminates the rotating mirror. Thus, the mirror delimits the beam, establishing a rectangular cross section which maintains its relative orientation with respect to the image surface. Thus, if uniformly illuminated, the focal point on the image surface (having in this case a $\text{sinc}^2 x \cdot \text{sinc}^2 y$ psf, x = along-scan and y = cross-scan) maintains the same orientation along its scanned line. If, however, the input beam is polarized, the axis of polarization of the imaged point will rotate directly with mirror rotation within the rectangular psf. Similarly will be rotation for any radial asymmetry (e.g., intensity or ellipticity) within the aperture, resulting in rotation of the psf.

Consider, therefore, the same scanner underilluminated with, for example, an elliptical Gaussian beam (with major axis horizontal). The axis of the imaged elliptic spot (intended major axis vertical)

will rotate directly with the mirror. Similarly, if the scanner is illuminated with multiple beams displaced slightly angularly in a plane (to generate an in-line array of spots), the axis of the imaged array will rotate directly with mirror rotation.

This effect is transferrable directly to the pyramidal polygon which is illuminated per Fig. 34. It may be considered as an array of mirrors, each exhibiting the same rotation characteristics as the monogon of Fig. 27. The mirrors of Fig. 34 are also overilluminated, maintaining a stationary geometric psf during scan (if uniformly illuminated), but subject to rotation of, for example, polarization within the psf. Similarly, it is subject to rotation of an elliptical beam within the aperture, or of a multiple-beam array.

Not so, however, for the mirror mounted per Fig. 28 (galvanometer mount), or for the prismatic polygon of Figs. 26 and 33, which may be considered a multifacet extension of Fig. 28. When the illuminating beam and the scanned beam form a plane which is normal to the axis of mirror rotation, execution of scan does not alter the characteristics of the psf, except for the possible vignetting of the optical aperture and possible alteration of reflection characteristics (e.g., polarization) with variation in incident angle. It is noteworthy that in the prior examples, the angles of incidence remained constant, while the image is subject to rotation; and here, the angles of incidence change, while the image develops no rotation.

The distinction is in the symmetry of the scanning system with respect to the illumination. The prior examples (maintaining constant incidence while exhibiting image rotation) are radially symmetric. The latter examples (which vary incidence but execute no image rotation) represent the limit of radial asymmetry. While mirrored optical scanners seldom operate in regions between these two extremes, holographic scanners can, creating possible complications with, for example, polarization rotation. This is manifest in the variation in diffraction efficiency of gratings for variation in p and polarizations during rotation. (See "Operation in the Bragg Regime.")

Image Derotation Image derotation can be implemented by interposing another image-rotating component in the optical path to cancel that caused by the scanner. The characteristic of an image rotator is that it inverts an image.¹³ Thus, with continuous rotation, it rotates the image, developing two complete rotations per rotation of the component. It must, therefore, be rotated at half the angular velocity of the scanner.

While the Dove prism¹³ is one of the most familiar components used for image rotation, other coaxial image inverters include¹³

- Three-mirror arrangement, which simulates the optical path of the Dove prism
- Cylindrical/spherical lens optical relay
- Pechan prism, which allows operation in converging or diverging beams

Design Considerations A commonly encountered scanner configuration is the prismatic polygon feeding a flat-field lens in preobjective scan, illustrated schematically in Fig. 26. The size and cost of the flat-field lens (given resolution and accuracy constraints) is determined primarily by its proximity to the scanner and the demand on its field angle. A larger distance from the scanner (pupil relief distance) imposes a wider acceptance aperture for a given scan angle, and a wider scan angle imposes more complex correction for off-axis aberration and field flattening. The pupil relief distance is determined primarily by the need for the input beam (Fig. 26) to clear the edge of the flat-field lens. Also, a wider scan angle reduces the accuracy requirement for pixel placement. Since the scan angle θ subtends the desired number N of resolution elements, a wider angle provides a larger angular subtense per element and correspondingly larger allowed error in angle $\Delta\theta$ for a desired elemental placement accuracy ΔN . This applies in both along-scan and cross-scan directions, $\Delta\theta_x$ and $\Delta\theta_y$, respectively (see Sec. 30.7).

Subsequent consideration of the scanner-lens relationships requires a preliminary estimate of the polygon facet count, in light of its diameter and speed. Its speed is determined by the desired data rates and entails considerations which transcend the optogeometric ones developed here. Diffraction-limited relationships are used throughout, requiring adjustment for anticipated aberration in real systems. The wavelength λ is a convenient parameter for buffering the design to accommodate

aberration. For example, an anticipated fractional spot growth of 15 percent due to systematic aberration is accommodated by using $\lambda_+ = 1.15\lambda$.

Performance characteristics which are usually predisposed are the resolution N (elements per scan), the optical scan angle θ , and the duty cycle η . Their interrelationships are presented under "Input/Output Scanning," notably by Eqs. (15) and (23). The values of N and θ for a desired image format width must be deemed practical for the flat-field lens.

Following these preliminary judgments, the collimated input beam width D is determined from [see "Fundamental Scanned Resolution" Eq. (15)]

$$D = Na\lambda/\theta \quad (29)$$

where a is the aperture shape factor and λ is the wavelength. For noncollimated beams, see "Augmented Resolution, the Displaced Deflector," notably Eq. (19). The number of facets is determined from Table 3 and Eq. (23),

$$n = 4\pi\eta/\theta \quad (30)$$

whereupon it is adjusted to an integer.

Scanner-lens relationships The polygon size and related scan geometry into the flat-field lens may now be determined.³⁵ Figure 35 illustrates a typical prismatic polygon and its input and output beams, all in the same plane. One of n facets of width W is shown in three positions: undeflected and in its limit-rotated positions. The optical beams are shown in corresponding undeflected and limit positions, deflected by $\pm \theta/2$. A lens housing edge denotes the input surface of a flat-field lens. Angle γ provides clear separation between the input beam and the down-deflected beam or lens housing. The pupil relief distance P (distance ac) and its slant distance P_e (distance bc) are system parameters which establish angle α such that $\cos \alpha = P/P_e$. Angle α represents the off-axis illumination on the

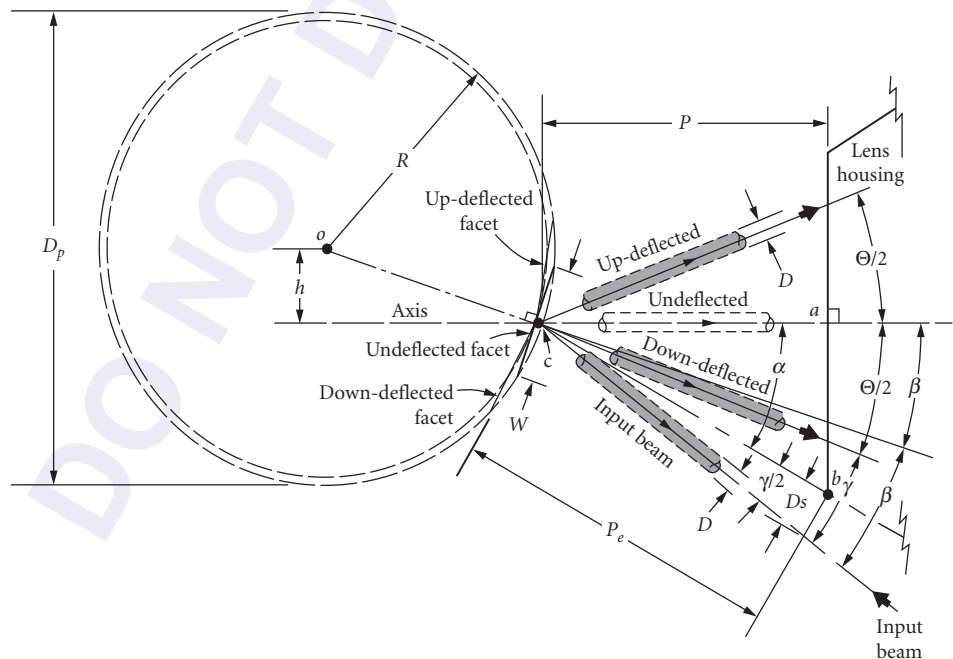


FIGURE 35 Polygon, beam, and lens relationships. Showing undeflected and limit facet and beam positions.⁴

polygon which broadens the input beam on the facet. The beam width D_m on the facet is widened due to α and due to an additional safety factor t ($1 \leq t \leq 1.4$) which limits one-sided truncation of the beam by the edge of the facet at the end of scan. Applying these factors, the beam width becomes

$$D_m = \frac{Dt}{\cos\alpha} \quad (31)$$

Following Eq. (15), the duty cycle is represented by $\eta = 1 - D_m/W$, yielding the facet width

$$W = D_m / (1 - \eta) \quad (32)$$

from which the outer (circumscribed) polygon diameter is developed;³⁵ expressed by

$$D_p = \frac{Dt}{(1 - \eta)\sin\pi/n\cos\alpha} \quad (33)$$

Solution of Eq. (33) or expressions of similar form³⁶ entails determination of α , the angle of off-axis illumination on the facet. This usually requires a detailed layout, similar to that of Fig. 35. Series approximation of $\cos\alpha$ allows transformation of Eq. (33) to replace α with more direct dependence on the important lens parameter P (pupil relief distance), yielding,

$$D_p = \frac{Dt}{(1 - \eta)\sin\pi/n} \cdot \frac{1 + \theta Ds/2P}{1 - \theta^2/8} \quad (34)$$

where, per Fig. 35, $s \approx 2$ is a safety multiplier on D for secure input/output beam separation and clearance.

Orientation of the scanner and lens also requires the height h , the normal distance from the lens axis to the polygon center. This is developed³⁵ as

$$h = R_c \sin(\gamma/2 + \theta/4) \quad (35)$$

where R_c is the radial distance oc , slightly shorter than the outer radius R , approximated to be

$$R_c = R \left[1 - \frac{1}{4}(\pi/n)^2 \right] \quad (36)$$

Holographic Scanners

General Characteristics Almost all holographic scanners comprise a substrate which is rotated about an axis, and utilize many of the concepts representative of polygons. An array of holographic elements disposed about the substrate serves as facets, to transfer a fixed incident beam to one which scans. As with polygons, the number of facets is determined by the optical scan angle and duty cycle (see "Duty Cycle"), and the elemental resolution is determined by the incident beam width and the scan angle [see Eq. (15)]. In radially symmetric systems, scan functions can be identical to those of the pyramidal polygon. While there are many similarities to polygons, there are significant advantages and limitations.⁶ The most attractive features of holographic scanners are

1. Reduced aerodynamic loading and windage with elimination of radial discontinuities of the substrate
2. Reduced inertial deformation with elimination of radial variations
3. Reduced optical-beam wobble when operated near the Bragg transmission angle

Additional favorable factors are

1. Operation in transmission, allowing efficient beam transfer and lens-size accommodation
2. Provision for disk-scanner configuration, with facets disposed on the periphery of a flat surface, designable for replication

3. No physical contact during exposure; precision shaft indexing between exposures allows for high accuracy in facet orientation
4. Filtering in retrocollection, allowing spatial and spectral selection by rediffraction
5. Adjustability of focus, size, and orientation of individual facets

Some limiting factors are

1. Need for stringent design and fabrication procedures, with special expertise and facilities in diffractive optics, instrumentation, metrology, and processing chemistry.
2. Accommodation of wavelength shift: exposure at one wavelength (of high photosensitivity) and reconstruction at another (for system operation). Per the grating equation⁶ for first-order diffraction,

$$\sin\theta_i + \sin\theta_o = \lambda/d \quad (37)$$

where θ_i and θ_o are the input and diffracted output angles with respect to the grating normal and d is the grating spacing, a plane linear grating reconstructs a collimated beam of a shifted wavelength at a shifted angle. Since wavefront purity is maintained, it is commonly employed,⁶ although it requires separate focusing optics (as does a polygon). When optical power is added to the hologram (to provide self-focusing), its wavelength shift requires compensation for aberration.⁶ Further complications arise when intended for multicolor operation, even if plane linear gratings. Further, even small wavelength shifts, as from laser diodes, can cause unacceptable beam misplacements, requiring corrective action.^{6,37}

3. Departure from radial symmetry develops complex interactions which require critical balancing to achieve good scan linearity, scan-angle range, wobble correction, radiometric uniformity, and insensitivity to input beam polarization.^{6,36} This is especially demanding in systems having optical power in the holograms.
4. Systems which retain radial symmetry to maintain scan uniformity may be limited in Bragg angle wobble reduction, and can require auxiliary compensation, such as anamorphic error correction.

Holographic Scanner Configurations A scanner which embodies some of the characteristics expressed above is represented in Fig. 36.³⁸ A cylindrical glass substrate supports an array of equally spaced *hololenses* which image the input beam incident at o to the output point at P . Since point o intersects the axis, the scanner is radially symmetric, whereupon P executes a circular (arced) scan concentric with the axis, maintaining magnification $m = \theta/\Phi = 1$. A portion of the radiation incident on the image surface is backscattered and intercepted by the hololens, and reflected to a detector which is located at the mirror image o' of point o . The resolution of this configuration is shown to be analogous to that of the pyramidal polygon.⁶

An even closer analogy is provided by an earlier reflective form illustrated in Fig. 37, emulating the pyramidal polygon, Fig. 34. It scans a collimated beam which is transformed by a conventional flat-field lens to a scanned focused line. This is one of a family of holofacet scanners,⁶ the most prominent of which tested to the highest performance yet achieved in combined resolution and speed—20,000 elements per scan at 200 Mpixels/s. This apparatus is now in the permanent collection of the Smithsonian Institution.

Operation in the Bragg Regime The aforementioned systems are radially symmetric and utilize substrates which allow derivation of the output beam normal to the rotating axis. While operation with radial asymmetry was anticipated in 1967,³⁹ it was not until operation in the Bragg regime was introduced^{6,40} that major progress developed in disk configurations. Referring to Fig. 38, the input and output beams I and O appear as principal rays directed to and diffracted from the holographic sector HS, forming angles θ_i and θ_o with respect to the grating surface normal.

For the tilt-error reduction in the vicinity of Bragg operation, the differential in output angle $d\theta_o$ for a differential in hologram tilt angle $d\alpha$ during tilt error $\Delta\alpha$ is given by

$$d\theta_o = \left[1 - \frac{\cos(\theta_i + \Delta\alpha)}{\cos(\theta_o - \Delta\alpha)} \right] d\alpha \quad (38)$$

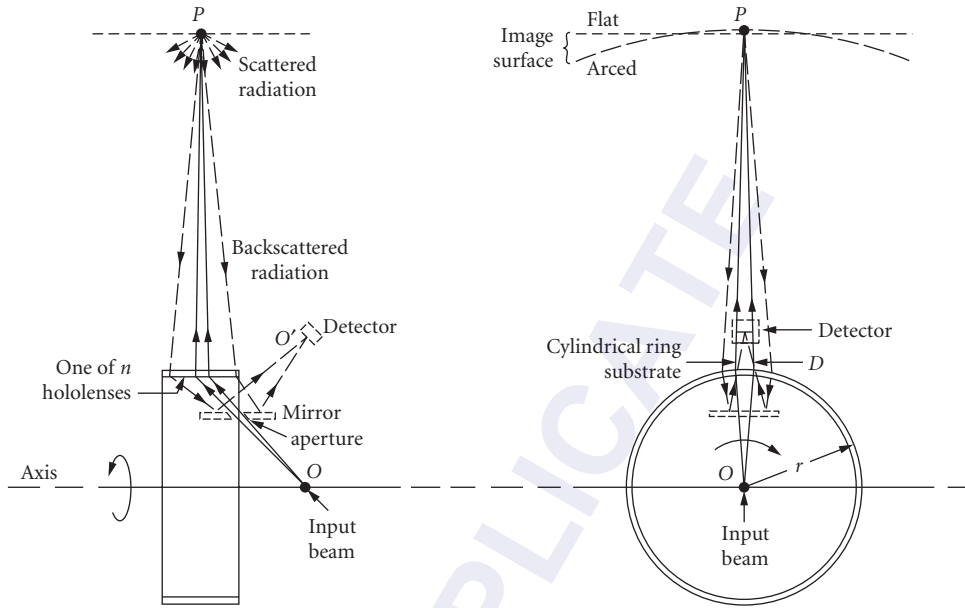


FIGURE 36 Transmissive cylindrical holographic scanner. Input beam underilluminates hololens which focuses diffracted beam to image surface. Dashed lines designate optional collection of backscattered radiation for document scanning.⁶

Hence, when $\theta_i = \theta_o$, a small $\Delta\alpha$ is effectively nulled. While the θ_i and θ_o depart from perfect Bragg symmetry during hologram rotation and scan, the reduction in error remains significant. An analogy of this important property is developed for the tilting of a refractive wedge operating at minimum deviation.⁶ When $\theta_i = \theta_o \approx 45^\circ$, another property develops in the unbowing of the output

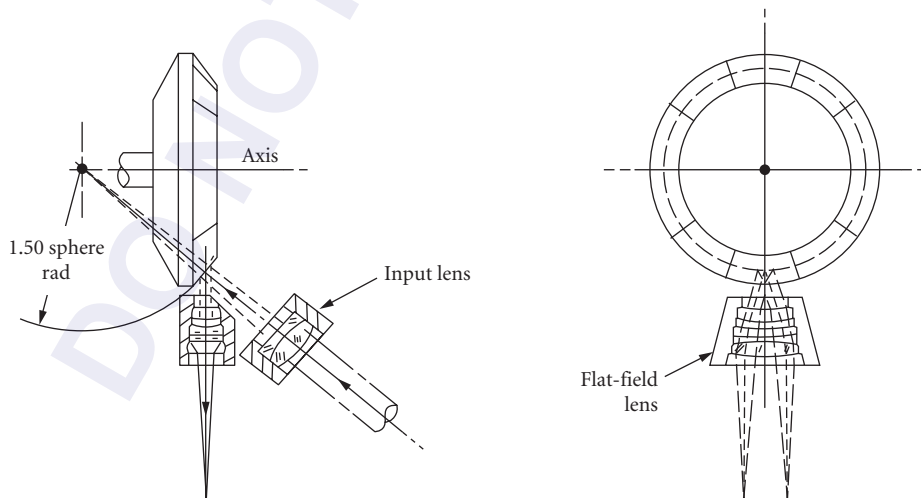


FIGURE 37 Reflective holofacet scanner, underilluminated. Flat-field microimage scanner (100 lp/mm over 11-mm format).

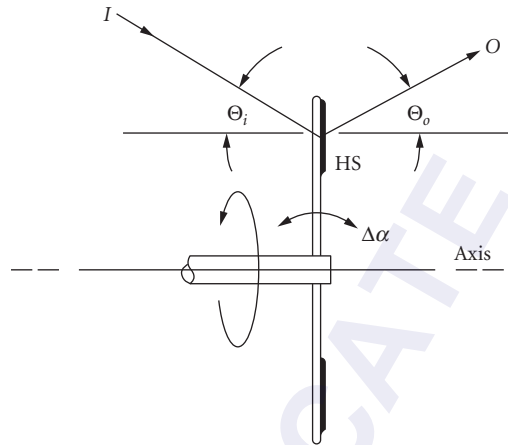


FIGURE 38 Holographic scanner in Bragg regime; $\Theta_i = \Theta_o$, in which output angle θ_o is stabilized against tilt error $\Delta\alpha$ of holographic segment HS.

scanned beam: the locus of the output beam resides (almost) in a plane normal to that of the paper over a limited but useful range.^{6,37} Further, the incremental angular scan for incremental disk rotation becomes almost uniform: their ratio m at small scan angles is shown to be equal to the ratio λ/d of the grating equation [see the section on “Radial Symmetry and Scan Magnification” and Eq. (3)].⁶ At $\theta_i = \theta_o = 45^\circ$, $m = \lambda/d = \sqrt{2}$. This results in the output-scan angle to be $\sqrt{2}$ larger (in its plane) than the disk-rotation angle. While such operation provides the above attributes, it imposes two practical restrictions.

1. For high diffraction efficiency from relief gratings (e.g., photoresist), the depth-to-spacing ratio of the gratings must be extremely high, while the spacing $d = \lambda/\sqrt{2}$ must be extremely narrow. This is difficult to achieve and maintain, and difficult to replicate gratings which provide efficient diffraction.
2. Such gratings exhibit a high polarization selectivity, imposing a significant variation in diffraction efficiency with grating rotation (see “Image Rotation and Derotation”).

Accommodation of these limitations is provided by reducing the Bragg angle and introducing a bow correction element to straighten the scan line. This is represented in Fig. 39; a high-performance scanner intended for application to the graphic arts. The Bragg angle is reduced to 30° . This reduces the magnification to $m = 1 = \lambda/d$ (as in radially symmetric systems), increases d to equal λ for more realizable deep-groove gratings, and reduces significantly the angular polarization sensitivity of the grating.

The elegance of the 45° Bragg configuration has been adapted⁶ to achieve self-focusing in less demanding tasks (e.g., laser printing). This is exemplified in Fig. 40, which includes a holographic lens to balance the wavelength shift of the laser diode,^{41,42,43} to shape the laser output for proper illumination of the scanner and to accommodate wavelength shift reconstruction. However, such multifunction systems are compounded by more critical centration requirements⁶ and balancing of characteristics for achievement of a discrete set of objectives.

Galvanometer and Resonant Scanners

To avoid the scan nonuniformities which can arise from facet variations (see Sec. 30.7) of polygons or holographic deflectors, one might avoid multifacets. Reducing the number to one, the polygon becomes a monogon. This adapts well to the internal drum scanner (Fig. 27), which achieves a high

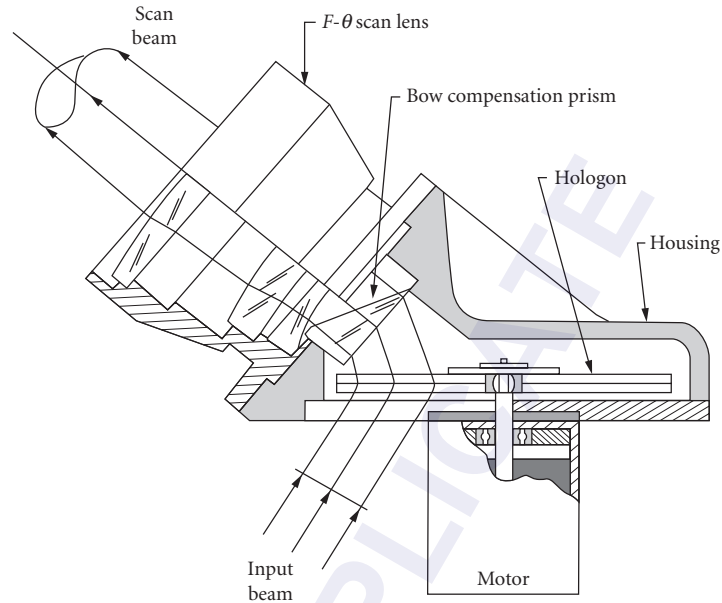


FIGURE 39 Plane linear grating (hologon) holographic disk scanner. Bragg angle of 30° provides fabricatable and polarization-insensitive grating structure, but requires bow compensation prism. Useful scan beam is in and out of plane of paper. (After Holotek Ltd, Rochester, NY product data.)

duty cycle, executing a very large angular scan within a cylindrical image surface. Flat-field scanning, however, as projected through a flat-field lens, allows limited optical scan angle, resulting in a limited duty cycle from a rotating monogon. If the mirror is vibrated rather than rotated completely, the wasted scan interval may be reduced. Such components must, however, satisfy system speed, resolution, and linearity. Vibrational scanners include the familiar galvanometer and resonant devices^{4,5,44,45} and the less commonly encountered piezoelectrically driven mirror transducer.^{5,45}

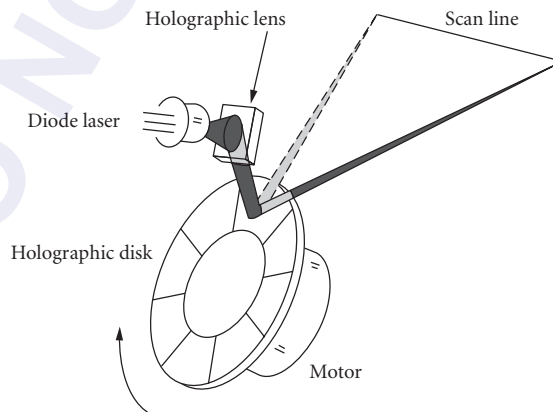


FIGURE 40 Holographic disk scanner with corrective holographic lens, both operating in approximately Bragg regime, providing complex error balancing.⁴²

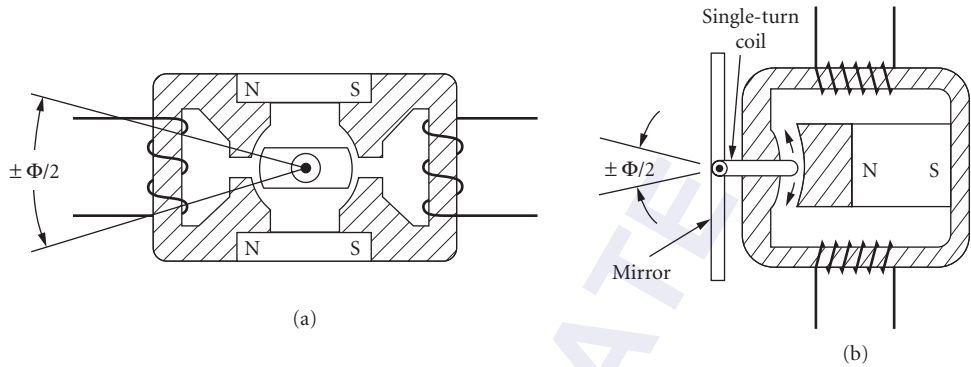


FIGURE 41 Examples of galvanometer and resonant scanner transducers. Fixed field of permanent magnet(s) is augmented by variable field from current through stator coils. (a) Galvanometer: torque rotates iron or magnetic core. Mirror surface (not shown) on extended shaft axis. (b) Resonant scanner: torque from field induced into single-turn armature coil (in plane perpendicular to paper) rotates mirror suspended between torsion bars. One stator coil may be nondriven and used for velocity pick-off.

The Galvanometer Referring to Fig. 41a, a typical galvanometer driver is similar to a torque motor. Permanent magnets provide a fixed field which is augmented (\pm) by the variable field developed from an adjustable current through the stator coils. Seeking a new balanced field, the rotor* executes a limited angular excursion ($\pm\Phi/2$). With the mirror and principal ray per Fig. 28, the reflected beam scans through $\pm\theta/2$, twice that of the rotor.

The galvanometer is a broadband device, damped sufficiently to scan within a wide range of frequencies, from zero to an upper value close to its mechanical resonance. Thus, it can provide the sawtooth waveform with a longer active linearized portion and shorter retrace time τ . This is represented in Fig. 42 (solid lines) showing rotation angle Φ versus time. As a broadband device, it can also serve for random access, positioning to an arbitrary location within its access-time limitations. For this feature of waveform shaping, the galvanometer was categorized as a *low inertia scanner*.⁵

The Resonant Scanner When damping is removed almost completely, large vibrations can be sustained only very near the resonant frequency of the oscillating system. The resonant scanner is thus characterized by larger angular excursions at a fixed and usually higher frequency, executing near-perfect sinusoidal oscillations. A typical driver configuration is illustrated in Fig. 41b. Figure 42 (dashed lines) shows a sinusoid with the same zero-crossings as those of the sawtooth waveform. Contrary to its popular designation as “low-inertia,” the resonant scanner provides rigid time increments, as though it exhibits a high inertia. While the rotary inertia of the suspension system is low to allow high repetition rates, it permits no random access and no scan waveform shaping, as do the galvanometer, acousto-optic, electro-optic, and other wideband scanners designated as low-inertia devices.⁵

Suspension Systems In the vibrational scanners, the bearings and suspension systems are the principal determinants of scan uniformity. The galvanometer shaft must be sufficiently stiff and long to inhibit cross-scan wobble. However, to maximize the oscillating frequency, the armature is restricted in size and mass. Fortunately, its reciprocating motion tends to retrace its path (and its perturbations) faithfully over many cycles, making adjacent scans more uniform than if the same shaft rotated completely within the same bearings, as in a motor.

*Rotor types include moving iron, moving magnet, or moving coil. Figure 41a illustrates the first two and Fig. 41b exemplifies the moving coil type. Moving magnet types (having NdFeB high-energy magnetic material) can exhibit some advantage in lower inertia and higher torque.

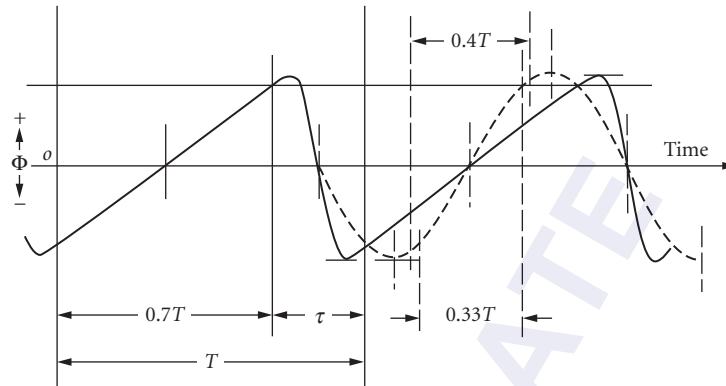


FIGURE 42 Waveforms (Φ vs. time) of vibrational scanners having same period and zero crossings. *Solid line*: galvanometer with linearized scan, providing 70 percent duty cycle. *Dashed line*: Resonant scanner providing 33.3 percent duty cycle (unidirectional) with 2:1 slope change, or 40 percent duty cycle with 3.24:1 slope change. Ratio of maximum slopes: resonant/galvanometer 2.6/1.⁴

Some bearings are flexure, torsion, or taut-band devices which insert almost no along-scan perturbations.⁴⁴ Because of their low damping, these suspensions are most often applied to the resonant scanner. When damped, they can serve for the galvanometer, suffering a small sacrifice in bandwidth and maximum excursion, but gaining more uniform scan with very low noise and almost unlimited life. Some considerations are their low radial stiffness and possible coupling perturbation from torsion, shift of the axis of rotation with scan angle, and possible appearance of spurious modes when lightly damped. Most of these factors can be well-controlled in commercial instrument designs.

Adaptations and Comparisons Because the resonant scanner oscillates sinusoidally, and we seek typically a linearized scan, some significant adaptations are often required. As illustrated in Fig. 42 (dashed lines), we must select from the sine function a central portion which is sufficiently linear to be linearized further by timing the pixels or extracting them out of memory at a corresponding rate.⁴⁶ To limit the variation in pixel rate to 2:1 (i.e., velocity at zero crossover will be twice that at the same limit), then the useful excursion must be restricted to $60^\circ/90^\circ$ or 66.7 percent of its peak angle. When scanning with only one slope of the sinusoid (as for generation of a uniformly spaced raster), this represents a duty cycle of only 33.3 percent. To raise the duty cycle, one must accommodate a greater variation in data rate. If, for example, the useful scan is 80 percent of its full excursion (40 percent when using one slope), then the velocity variation rises to 3.24 \times . That is, the data rate or bandwidth at crossover is 3.24 times that at the scan limit. Also, its bandwidth at crossover is approximately 2 1/2 times that of the galvanometer, as represented by their relative slopes in Fig. 42.

There is a corresponding variation in the dwell time of the pixels, resulting in predictable but significant variation in image exposure or detectivity: 2:1 for 33.3 percent duty cycle and 3 1/4: 1 for 40 percent duty cycle. This may require compensation over the full scan interval, using position sensing and control.⁴⁵⁻⁴⁸ In contrast, the broadband galvanometer with feedback can provide a highly linearized scan⁴⁴ at a duty cycle of approximately 70 percent.⁵

Acousto-Optic Scanners

Acousto-optic diffraction serves effectively for high-speed low-inertia optical deflection. It can provide random beam positioning within extremely short access times, or generate repetitive linear scans at very high rates, or divide a single beam into multiple beams for multiplexing applications. The trade-off is, however, relatively low resolution, seldom providing more than $N = 1000$ elements per scan.

The principles of acousto-optics were formulated in 1932⁴⁹ and its attributes were applied only 5 years later to the Scophony TV projection system.⁵⁰ Its potential for laser scanning was explored in the mid-1960s.^{51,52} While various acousto-optic interactions exist, laser scanning is dominated by operation in the Bragg regime.^{5,53}

Fundamental Characteristics Diffraction from a structure having a periodic spacing Λ is expressed as $\sin \theta_i + \sin \theta_o = n\lambda/\Lambda$, where θ_i and θ_o are the input and output beam angles respectively, n is the diffractive order, and λ is the wavelength. Bragg operation requires that $\theta_i = \theta_o = \theta_B$. In a “thick” diffractor, length $L \geq \Lambda^2/\lambda$, wherein all the orders are transferred efficiently to the first, and the Bragg angle reduces to

$$\theta_B = \frac{1}{2} \frac{\lambda}{\Lambda} \quad (39)$$

Per Fig. 43, the grating spacing is synthesized by the wavefront spacing formed by an acoustic wave traveling through an elastic medium. An acoustic transducer at one end converts an electrical drive signal to a corresponding pressure wave which traverses the medium at the velocity v_s , whereupon it is absorbed at the far end to suppress standing waves. The varying pressure wave in the medium forms a corresponding variation in its refractive index. An incident light beam of width D is introduced at the Bragg angle (angle shown exaggerated). An electrical drive signal at the center frequency f_o develops a variable index grating of spacing Λ which diffracts the output beam at θ_B into position b . The drive signal magnitude is adjusted to maximize beam intensity at position b , minimizing intensity of the zero-order beam at position a . When f_o is increased to $f_s = f_o + \Delta f$, the grating spacing is decreased, diffracting the output beam through a larger angle, to position c . The small scan angle θ is effectively proportional to the change in frequency Δf .

The scan angle is $\theta = \lambda/\Delta\Lambda = (\lambda/v_s)\Delta f$. The beam width, traversed by the acoustic wave over the transit time τ is $D = v_s\tau$. Substituting into Eq. (15) and accounting for duty cycle per Eq. (22), the resolution of the acousto-optic scanner (total N elements for total Δf) is

$$N = \frac{\tau\Delta f}{a}(1 - \tau/T) \quad (40)$$

The $\tau\Delta f$ component represents the familiar time-bandwidth product, a measure of information-handling capacity.

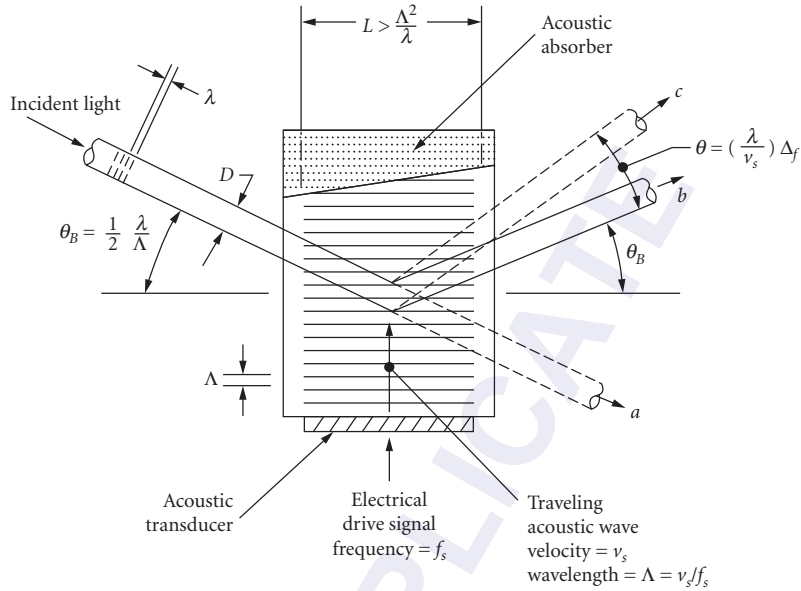
Deflection Techniques Because the clear aperture width W of the device is fixed, anamorphic optics is often used to illuminate W with an adjusted beam width D —encountering selective truncation by the parallel boundaries of W . The beam height (in quadrature to D) can be arbitrarily narrow to avoid apodization by the aperture. This one-dimensional truncation of the Gaussian beam requires assignment of an appropriate aperture shape factor a , summarized in Table 5.

Additional topics in acousto-optic deflection are cylindrical lensing due to linearly swept f_s ,⁵³ correction for decollimation in random access operation,⁵ Scophony operation,⁵⁴ traveling lens or chirp operation,⁵⁵ correction for color dispersion,⁵⁶ polarization effects,⁵⁷ and multibeam operation.⁵⁸

Electro-Optic (Gradient) Scanners

The gradient deflector is a generalized form of beam scanner^{4,5,59} in which the propagating wavefronts undergo increasing retardation transverse to the beam, thereby changing the wavefront spacing (wavelength) transverse to the beam. To maintain wavefront continuity, the rays (orthogonal trajectories of the wavefronts) bend in the direction of the shorter wavelength. Referring to Fig. 44a, this bend angle θ through such a deflection cell may be expressed as

$$\theta = k_o (dn/dy)l \quad (41)$$



Output beam position	Beam condition	Acoustic frequency
a	Undiffracted zero order	$f_s = \text{zero (or amplitude} = 0)$
b	Nominal Bragg angle Θ_B	$f_s = f_o$ (center frequency)
c	Added diffraction, scanned through angle Θ	$f_s = f_o + \Delta f$

FIGURE 43 Bragg acousto-optic deflector (angles exaggerated for illustration). Electrical drive signal generates traveling acoustic wave in medium which simulates a thick optical grating. Relationship between the output beam position and the electrical drive frequency is tabulated.⁴

where n is taken as the number of wavelengths per unit axial length l , y is the transverse distance, and k_o is a cell system constant. For the refractive material form in which the wavefront traverses a change Δn in index of refraction and the light rays traverse the change in index over the full beam aperture D in a cell of length L , then the relatively small deflection angle becomes^{4,5}

$$\theta = (\Delta n/n_f)L/D \tag{42}$$

where n_f is the refractive index of the final medium (applying Snell's law and assuming $\sin \theta = \theta$). Following Eq. (15), the corresponding resolution in elements per scan angle is expressed as

$$N = (\Delta n/n_f)L/a\lambda \tag{43}$$

The Δn is given by

$$\text{(for class I materials)} \quad \Delta n = n_o^3 r_{ij} E_z \tag{44a}$$

$$\text{(for class II materials)} \quad \Delta n = n_e^3 r_{ij} E_z \tag{44b}$$

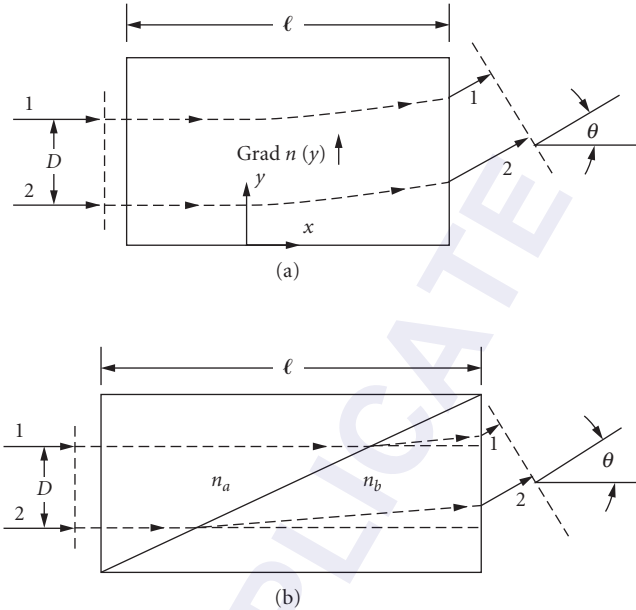


FIGURE 44 Equivalent gradient deflectors: (a) basic deflector cell composed of material having grad $n(y)$. Ray 1, propagating through a higher refractive index, is retarded more than Ray 2, tipping the wavefront through angle Θ (including boundary effect). (b) Analogous prismatic cell, in which $n_a > n_b$, such that Ray 1 is retarded more than Ray 2, tipping the wavefront through angle Θ .⁴

where $n_{o,e}$ is the (ordinary, extraordinary) index of refraction, r_{ij} is the electro-optic coefficient, and $E_z = V/Z$ is the electric field in the z direction (see Fig. 45).

Methods of Implementation An electroacoustic method of developing a time-dependent index gradient was proposed in 1963⁶⁰ utilizing the (harmonic) pressure variations in an acoustically driven cell (of a transparent elastic material). Although this appears similar to acousto-optic deflection (see “Acousto-Optic Scanners”), it differs fundamentally in that the cell is terminated reflectively rather than absorptively (to support a standing wave). Also, the acoustic wavelength is much longer than the beam width, rather than much shorter for Bragg operation. A method of approaching a linearly varying index gradient utilizes a quadrupolar array of electrodes bounding an electro-optic material;^{61,62} and is available commercially.⁶³

A continuous index gradient can be simulated by the use of alternating electro-optic prisms.^{59,64} A single stage biprism is illustrated in Fig. 44b and an iterated array for practical implementation appears in Fig. 45. Each interface imparts a cumulative differential in retardation across the beam. The direction and speed of retardation is controlled by the index changes in the electro-optic material. While resolution is limited primarily by available materials, significant experiment and test is reported for this form of deflector.⁵

Drive Power Considerations The electrical power dissipated within the electro-optic material is given by⁵

$$P = \frac{1}{4} \pi V^2 C f / Q \quad (45)$$

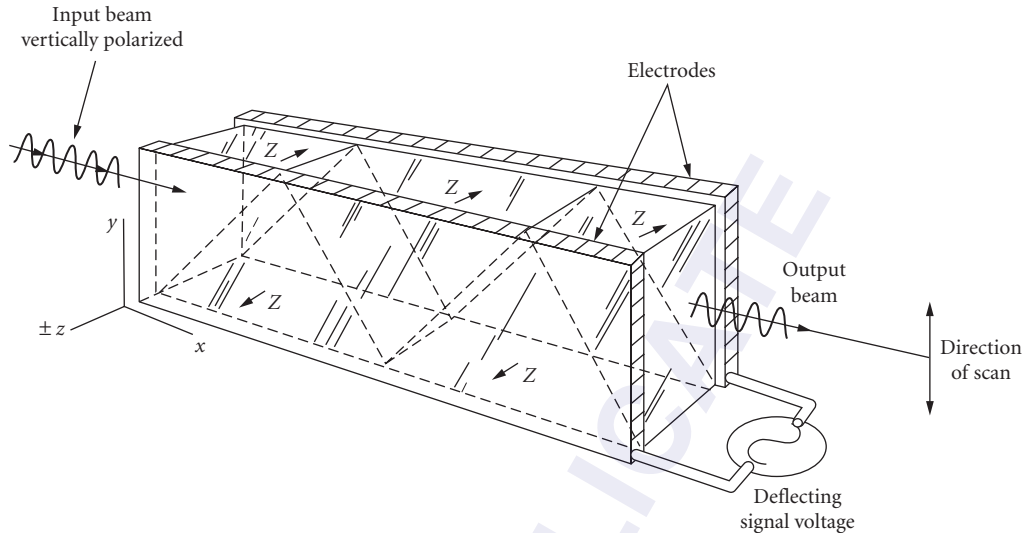


FIGURE 45 Iterated electro-optic prism deflector. Indicating alternating crystallographic (z) axes. Input beam polarization for class 1 electro-optic coefficient (r_{63}) materials.⁴

where V is the applied (p-p sinusoidal) voltage in volts, C is the deflector capacitance in farads, f is the drive frequency in hertz and Q is the material Q factor [$Q = 1/\text{loss tangent } (\tan \delta) \approx 1/\text{power factor}$, ($Q > 5$)].

The capacitance C for transverse electroded deflectors is approximately that for a parallel-plate capacitor of (rectangular) length L , width Y , and dielectric thickness Z (per Fig. 45)

$$C = 0.09\kappa LY/Z \text{ picofarads} \quad (46)$$

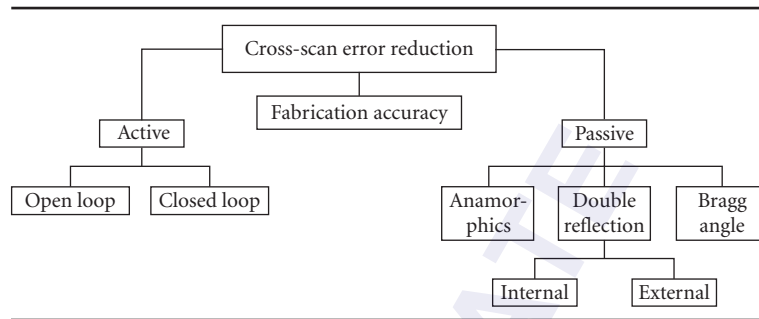
where k is the dielectric constant of the material (L, Y, Z in centimeters).

The loss characteristics of materials which determine their operating Q are often a strong function of frequency beyond 10^5 Hz. The dissipation characteristics of some electro-optic materials are provided,^{5,65,66} and a resolution-speed-power figure of merit has been proposed.⁶⁷

Unique Characteristics Most electro-optic coefficients are extremely low, requiring high drive voltages to achieve even moderate resolutions (to $N \approx 100$). However, these devices can scan to very high speeds (to $10^5/\text{s}$) and suffer effectively no time delay (as do acousto-optic devices), allowing use of broadband feedback for position control.

30.7 SCAN-ERROR REDUCTION

High-resolution scanners often depend on precise rotation of a shaft about its axis, said shaft supporting a multiplicity of deflecting elements (facets, mirrors, holograms). The control of angular uniformity of these multielements with respect to the axis, and of the axis with respect to its frame, can create an imposing demand on fabrication procedures and consequential cost. Since uniformity of beam position in the cross-scan direction may not be approached by phasing and timing of the data (as can be the along-scan errors), several noteworthy techniques have been developed to alleviate this burden.

TABLE 6 Techniques for Cross-Scan Error Reduction

Available Methods

The general field of cross-scan error reduction is represented in Table 6. *Fabrication accuracy* may be selected as the only discipline, or it may be augmented by any of the auxiliary methods. The *active* ones utilize high-speed low-inertia (A-O or E-O) or piezoelectric deflectors^{5,68} or lower-speed (galvanometer) deflectors which are programmed to rectify the beam-position errors. While open-loop programming is straightforward (while accounting for angular magnification/demagnification as a function of the accessed beam size), elegant closed-loop methods may be required to rectify pseudo-random perturbations. This must, however, be cost-effective when compared to the alternatives of increased fabrication accuracy and of the *passive* techniques.

Passive Methods

Passive techniques require no programming. They incorporate optical principles in novel configurations to reduce beam misplacement due to angular error in reflection or diffraction. Bragg-angle error reduction of tilted holographic deflectors is discussed in the section, "Operation in the Bragg Regime."

Anamorphic Error Control Anamorphic control, the most prominent treatment, may be applied to any deflector. The basics and operational characteristics⁶ are summarized here.

Separating the nonaugmented portion of the resolution equation [Eq. (19)] into quadrature components and denoting the cross-scan direction as y , then the error, expressed in the number of resolvable elements, is

$$N_y = \frac{\theta_y D_y}{a\lambda} \quad (47)$$

where $a\lambda$ is assumed constant, θ_y is the angular error of the output beam direction, and D_y is the height of the beam illuminating the deflector. The objective is to make $N_y \rightarrow 0$. Mechanical accuracies determine θ_y , while anamorphics are introduced to reduce D_y ; usually accomplished with a cylindrical lens focusing the illuminating beam in the y direction upon the deflector. [The quadrature (along-scan) resolution is retained by the unmodified D_x and scan angle θ_x .] As D_y is reduced, the y displacement error is reduced. Following deflection, the y -direction scanned spot distribution is restored by additional anamorphics—restoring the nominal converging beam angle (via F_y , the F-number forming the scanning spot in the y direction).

The error reduction ratio is

$$R = D'_y/D_y \quad (48)$$

where D'_y is the compressed beam height and D_y is the original beam height on the deflector.

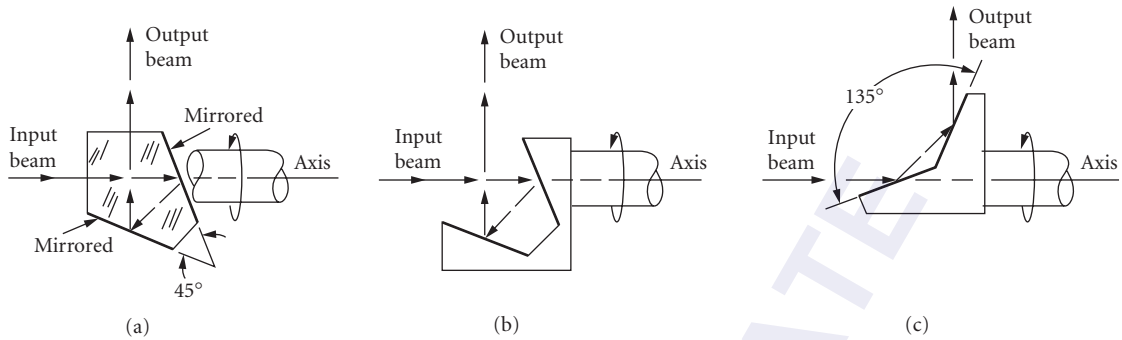


FIGURE 46 Monogon scanners employing double-reflection: (a) pentaprism; (b) pentamirror; and (c) open mirror.⁴

A variety of anamorphic configurations has been instituted, with principal variations in the output region, in consort with the objective lens, to reestablish the nominal F_y while maintaining focused spot quality and uniformity.

Double-Reflection Error Control In double-reflection (Table 6), the deflector which creates a cross-scan error is reilluminated by the scanned beam in such phase as to tend to null the error. This can be conducted in two forms: internal and external.

An internal double-reflection scanner is exemplified by the pentaprism monogon⁶⁹ in Fig. 46a; a (glass) substrate having two of its five surfaces mirrored. This is an optically stabilized alternate to the 45° monogon of Fig. 27, operating preobjective in collimated light. Tipping the pentaprism cross-scan (in the plane of the paper) leaves the 90° output beam unaffected. A minor translation of the beam is nulled when focused by the objective lens. The pentamirror⁶⁹ per Fig. 46b, requires, however, significant balancing and support of the mirrors, since any shift in the nominal 45° included angle causes twice the error in the output beam. A stable double-reflector is the open mirror monogon⁷⁰ of Fig. 46c. Its nominal 135° angle serves identically to maintain the output beam angle at 90° from the axis, independently of cross-scan wobble. With a rigid included angle and simple balancing, it can provide high-speed operation.

Two variations which double the duty cycle, as would a two-faceted pyramidal polygon or ax-blade scanner (see “Early Single-Mirror Scanners”) appear in Fig. 47. Figure 47a is effectively two pentamirrors

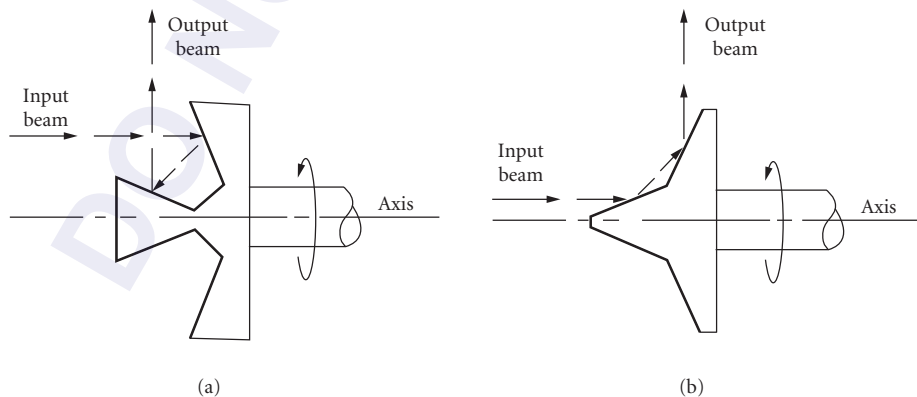


FIGURE 47 Paired scanners employing double-reflection: (a) paired pentamirror “butterfly” scanner and (b) paired open mirror scanner.⁴

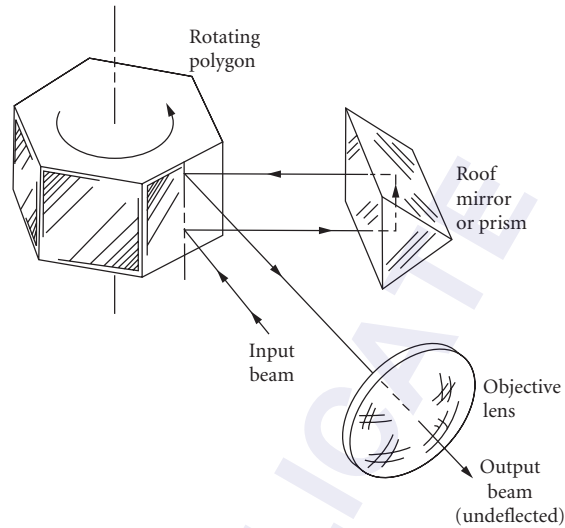


FIGURE 48 Method of external double-reflection—shown in undeflected position. Components and distances not to scale. Only principal rays shown.⁴

forming a *butterfly scanner*⁷¹ and Fig. 47b is effectively a pair of open mirrors.⁷² The absolute angles of each half-section must maintain equality to within half of the allowed error in the output beam. Also, the center section of Fig. 47a must be angularly stable to within one-quarter of the allowed error, because an increased included angle on one side forms a corresponding decrease on the other. Other dynamic considerations involve inertial deformation, and the beam displacements and mirror widths (not shown) to accommodate the distance of the input beam from the axis during rotation.

The need for near-perfect symmetry of the multiple double-reflectors can be avoided by transferring the accuracy requirement to an *external* element that redirects recurrent beam scans. One such form⁷³ is illustrated in Fig. 48. A prismatic polygon illuminated with a collimated beam of required width (only principal rays shown) deflects the beam first to a roof mirror, which returns the beam to the same facet for a second deflection toward the flat-field lens. The roof mirror phases the returned beam such as to null the cross-scan error upon the second reflection. Several characteristics are noteworthy:

1. The along-scan angle is doubled. That is, scan magnification $m = 4$ rather than 2.
2. This normally requires increasing the number of facets to provide the same angle with the same duty cycle.
3. However, during polygon rotation, the point of second reflection shifts significantly along the facet and sacrifices duty cycle.
4. The pupil distance from the flat-field lens is effectively extended by the extra reflections, requiring a larger lens to avoid vignetting.
5. The roof mirror and flat-field lens must be sized and positioned to minimize obstruction of the input and scanned beams. Allow for finite beam widths (see “Scanner-Lens Relationships”).

30.8 AGILE BEAM STEERING

A class of low-inertia scanning, called *agile beam steering*^{74–77} was developed initially for such challenging tasks as laser radar (LIDAR) and forward-looking infrared (FLIR) systems. Further advancement may allow its extension to more general application. The motivation for this

work is to achieve the performance of the scanned mirror while avoiding some of the concomitant burdens of size, weight, and inertia. This has been a long-envisioned goal of many earlier researchers in work having important similarities⁵ to the more current activity in agile beam steering. Recent research has harvested new resources such as liquid crystal E-O phase shifters micromachined devices, and microlens techniques assembled in novel configurations. Two principal approaches have dominated investigation and development, viz., the phased array and the decentered microlens array; along with some of their principal variations. Although the basic operation of these two array types differ, they both develop the same form of the steered output wavefronts

Phased Array Beam Steering

The directing of radiation in the radio and microwave regions by driving antenna arrays with controlled relative phase was especially familiar to the early radar specialist.^{78,79} Its adaptation to the optical spectrum, notably as radiated by lasers, was investigated⁸⁰ in 1964 following the invention of the laser and continued⁸¹⁻⁸³ through the early 1970s. The prospect of altering the direction of a laser beam with small adjustments on a group of radiators appeared *very* attractive. With the introduction of electrostatically actuated membrane mirror arrays⁸² in 1971, and the programming of electro-optic crystal arrays⁸³ in 1972, operational utility was affirmed. Beam steering with arrays of mirrors was investigated⁸¹ in 1967. Further work was conducted in mirror array beam steering in the infrared region, where mirror reflectance exceeds the transmittance of even the exotic infrared materials, and where the longer wavelength imposes lower requirements on mirror flatness. With the current use of faster acting electro-optic materials and novel design variations, substantive advances have been achieved.

The steering of an optical wavefront by phase variation is introduced with the effect of refractive prisms.^{84,87} Figure 49a illustrates a plane wavefront in air incident parallel to the plane surface of one dielectric wedge. Within the material of refractive index >1 , the wavelength is compressed proportionately, while the fronts remain parallel to the incident wavefront. A linearly increasing local phase delay results from the progressive retardation of the wavefronts across the enlarging wedge thickness. Traversing the tilted boundary, the wavelength in air is reexpanded and its angle of propagation is refracted as illustrated. This is exemplified in the iterated prismatic deflector of Fig. 45 given the dynamics of the electro-optic material.

To provide a wide aperture, the single prism of Fig. 49a can exhibit substantive bulk. To relieve this, the long wedge profile is divided into an array of smaller wedge increments where each causes a linear phase delay of from 0 to 2π . As illustrated in Fig. 49b, when implemented as described below, it provides the same deflection as the continuous single wedge. Along with the need to accommodate both the slope and refractive index of the material, one must dimension the periods of the wedges such that they form 2π phase differentials (or multiples thereof, i.e., modulo 2π) at the operating wavelength to assemble continuous nonstaircased wavefronts in the near field. This is functionally analogous to the reflective blazed grating, in which high efficiency is achieved when the angle of specular reflection (from the sawtooth slopes of the grating surfaces) coincides with the angle of diffraction *at the selected wavelength*. Thus, an array formed of such 2π phase differentials exhibits dispersion which limits efficient performance to narrow spectral-band (near-monochromatic) operation. Work conducted to alleviate this limitation is discussed in if "Development of Phased Arrays" subsection.

The above examples provide continuous phase retardation by virtue of their linear surface slopes; incremental or continuous. Incremental phase retardation can be controlled in transmission by an array of small electro-optic cells, and in reflection by precise actuation of individual mirrored pistons. An array of refractive phase retarders and the assembly of the radiated wavelets into contiguous wavefronts is represented in Fig. 49c. Operation is similar with pistons, except that the reflective piston requires displacement of only half of the 2π phase-retardation distance. A thin electro-optic retarder having a reflective ground plane requires not only half thickness, but attains a

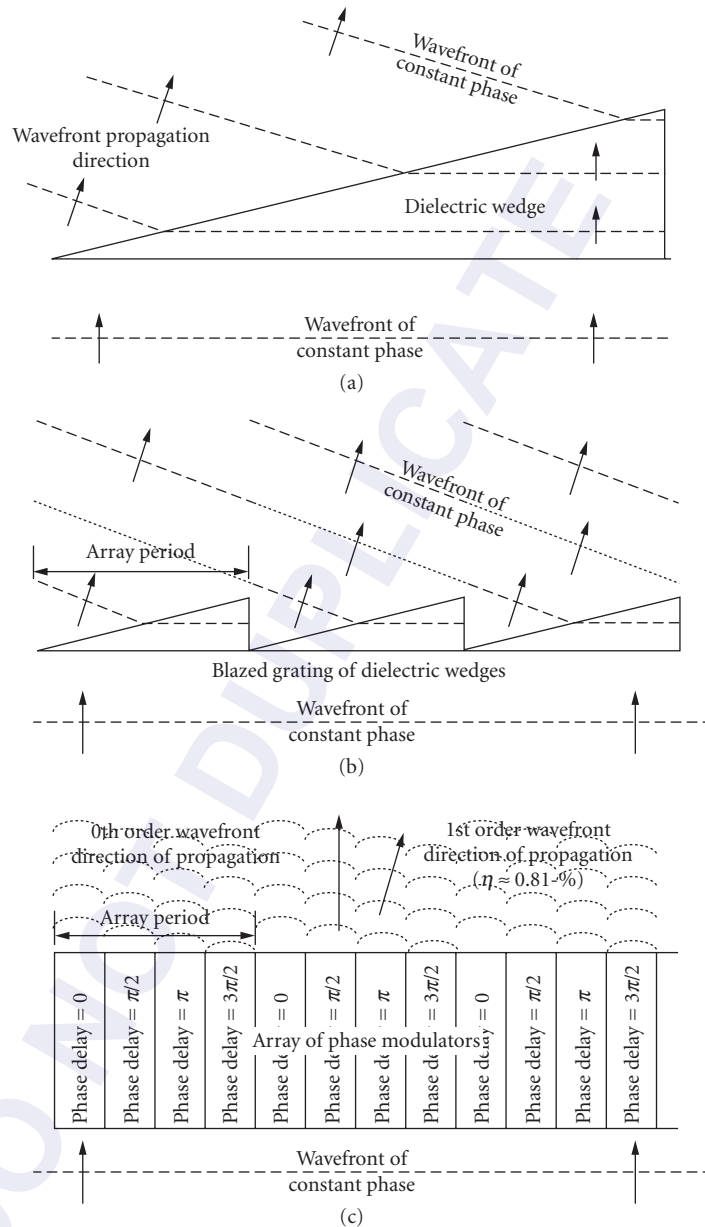


FIGURE 49 Illustrative steps (a, b, and c) toward optical-phased array beam steering. (Wavefront enters parallel to each bottom surface.) (a) Single dielectric wedge, illustrating familiar refraction of plane wavefront. (b) Synthesis of (a) with array of wedges. Each wedge imparts a 2π phase delay over each array period. (c) Synthesis of (b) with multiple delay elements (4 per 2π period). Output wavelets superpose into wavefronts having idealized efficiency of 81 percent. Greater multiplicity provides higher efficiency. (After Refs. 84 and 87.)

fourfold increase in switching speed which can be useful for a nominally slow liquid crystal retarder. An alternative option, shown in this figure, shows division of the full 2π phase change into four substeps per full cycle, which diffracts 81 percent of the energy into the first order. If more steps are available, then the diffraction efficiency is greater. For example, eight-step cycle attains an otherwise lossless efficiency of 95 percent. Figure 49 is a combination of three separated illustrations^{84,87} to unify their progression.^{74,85,86} The heuristic observations rendered above are affirmed by the considerations which follow.

The Analytic Base

The angular beam relationships of the phased array are expressed by the diffraction grating equation,⁸⁷

$$\sin(\theta_i) + \sin(\theta_o) = n\lambda/\Lambda \quad (49a)$$

where θ_i and θ_o are the input and output beam angles with respect to the grating normal (boresight), n is the diffraction order, λ is the free space wavelength, and Λ is the grating (array) period, per Fig. 49b and c. As in Fig. 49c showing four delay elements per array period, for q delay elements, each separated by a fixed distance d , $\Lambda = qd$. Since the number of elements in each period is $q = 2\pi/\phi$, where ϕ is the phase shift between elements, then $\Lambda = (2\pi/\phi)d$ is the distance required to assemble a one-wave phase difference. When $\theta_i = 0$, the angle of first-order wave propagation (Fig. 49c) is given by

$$\sin\theta_o = \lambda/\Lambda \quad (49b)$$

$$= \lambda/qd = \lambda\phi/2\pi d \quad (49c)$$

The normalized intensity I of the radiation pattern follows the analogy of the one-dimensional microwave phased array⁷⁹ expressed as

$$I = (\sin N\alpha / N \sin \alpha)^2 \quad (50)$$

with

$$\alpha = \pi d / \lambda (\sin\theta - \sin\theta_o) \quad (51)$$

where θ is the angle with respect to the grating normal at which the field in free space is measured, and N is the number of phase shifters in the array. The elemental spacing d provides uniform phase difference ϕ between elements.

The efficiency η_q of a linear array having the nominal (blazed) 2π phase resets, illustrated in Fig. 49b and c is expressed by

$$\eta_q = \left(\frac{\sin \frac{\pi}{q}}{\pi/q} \right)^2 = \text{sinc}^2 \frac{\pi}{q} \quad (52)$$

where q is number of elements per 2π array period. This may be recognized as similar to the Fourier transform¹³ of a uniformly illuminated linear aperture.⁸⁷ Inserting values of q , i.e., 4 and 8, Eq. (52) yields $\eta_4 = 0.81$ and $\eta_8 = 0.95$, respectively, as indicated earlier. With a reduced q , lower efficiency results from depletion of the main lobe to the sidelobes due to wavefront staircasing.

Typical liquid crystal phase retarder elements exhibit a unique loss factor established by the minimum space required to relax its orientation from a 2π phase shift to zero. This “flyback” transition is analogous to the flyback time τ of many conventional scanners as expressed in Eq. (22) as $\eta = 1 - \tau/T$, where T is the full scan period. This represents a time loss burdening high-speed operation. The duty cycle for liquid crystal elements is given by

$$\eta_{\Lambda} = (1 - \tau/\Lambda)^2 \quad (53)$$

where the time terms are replaced with τ representing the flyback width and Λ the full 2π width. The expression is squared to denote the radiated intensity rather than the time in Eq. (22), as illustrated by the solid line in Fig. 42. In addition, fill factor accounts for a cell having its operating portion occupy less than its full allotted area while the vignetting factor accounts for the loss of input illumination beyond the boundary of the array.

The far-field angular beamwidth θ_B is expressed as a minor variation to the familiar diffraction relation

$$\theta_B = a\lambda/D \quad (54a)$$

where a is the aperture shape factor⁸⁷ modifying the beamwidth, discussed in Sec. 30.3. With the full aperture width $D = Nd$,

$$\theta_B = a\lambda/Nd \quad (54b)$$

The Resolution of Phased Arrays

Equation (54) denotes the output beamwidth, that is the breadth of the principal lobe of radiation. For scanned resolution (Sec. 30.3), the number N of these adjacent lobes which fill the field-of-view along a linear track represents the number of phase-shifting calls in a linear phased array. Analogous to Eq. (4a),

$$\mathcal{N} = \theta/\theta_B \quad (55)$$

Half of the full deflected field angle θ is represented by the (positive) first diffracted order of grating [i.e., $n = +1$ in Eq. (49)]. When the array is addressed in complementary phase sequence, the same deflection magnitude results in the opposite ($n = -1$) direction. Thus, for typically small θ_0 in Eq. (49) and with $D = \mathcal{N}d$, the numerator for Eq. (55) becomes

$$\theta = \frac{2\lambda}{qd} = \frac{2\lambda}{qD} \mathcal{N} \quad (56)$$

With $\theta_B = a\lambda/D$ providing the denominator of Eq. (55), and accounting for the central boresight position, the steered resolution reduces to

$$\mathcal{N} = 1 + \frac{2}{aq} \mathcal{N} \quad (57a)$$

which is independent of wavelength. Although Eq. (57) appears to differ from the fundamental Eq. (15) for scanned resolution, substituting Eq. (56) into the fundamental equation and adding one for the boresight beam yields Eq. (57).

With a as a relative constant, the ratio of the two variables \mathcal{N} and q dominates the total number of elements divided by the number of elements per phase reset. Thus, the number of phase resets is the principal variable which determines the steered resolution and $\mathcal{N}/q = D/\Lambda$, where Λ is the array period. Since q is a parameter of the system, its adjustment also affects⁸⁸ the closeness of the adjacent steering states.

When a is unity, it denotes uniform illumination upon a rectangular aperture. This yields a far-field intensity distribution⁸⁹ of the $\text{sinc}^2(x)$ function having a main lobe within equispaced null

intervals. Rayleigh resolution requires this uniform illumination upon a rectangular aperture, and that the adjacent spots in the far field overlap such that the maximum of each main lobe coincides with the first null of each adjacent one. Further delimiting Eq. (57a) is that it is impractical to form a modulo 2π array in which q is less than three cells in view of the resulting disruption of the ramp wavefronts and the loss in efficiency. Letting $a = 1$ and setting $q_{\min} = 3$, the steered resolution is often expressed as an assumed Rayleigh resolution yielding

$$\mathcal{N}_{\max} = 1 + \frac{2}{3} \mathcal{N} \quad (57b)$$

A common illumination (laser) is the Gaussian function, with adjustment of the aperture overfill (and/or with complementary Gaussian filtering) to control the intensity distributed across the full aperture width W . The degree of overfill is, however, moderated by the reduction in light throughput due to aperture vignetting. As W is illuminated more uniformly, it may be limited by the appearance of fine structure beyond the main lobe when approaching the appearance of the sinc² function.

To quantify this value of a , different conditions can be considered⁸⁷ and tabulated, summarizing its value for the Gaussian beam of width D , either falling substantially within the aperture W (untruncated beam), or when the $1/e^2$ intensity of the input beam occurs at the aperture boundary (truncated beam). These are two typical illumination conditions of most conventional deflectors. For the rectangular aperture of width $W >$ height, illumination with a Gaussian beam primarily in the W or scan direction is further evaluated and tabulated (Table 5) providing data of current interest. A variable beam width D (at $1/e^2$ intensity) illuminates the full width W of a linear array. Assigning a parameter $\rho = W/D$, when $\rho = 1$, the $1/e^2$ beamwidth matches the full aperture width W . At $\rho = 1.5$, the array is 1.5 times wider than that of the $1/e^2$ beamwidth. This terminates the Gaussian function at $\pm 3\sigma$, where its intensity tapers to a small fraction of its maximum value, representing a practical limit on the narrowness of the input beam. At the other extreme, when $\rho \rightarrow 0$, the input beamwidth $D \gg W$, extracting near-uniform illumination from the center of the beam, and imposing extreme light loss beyond the aperture. This is the case of $a \rightarrow 1$. The aperture shape factors for the other two cases are determined; at $\rho = 1$, $a = 1.15$ and at $\rho = 1.5$, $a = 1.35$. Related to the topic of resolution is the *finesse* which is the smallest addressable increment of beam position. Consideration of this factor⁸⁷ involves (a potentially nonuniform) adjustment of the values of q (number of delay elements per array period).

Development of Phased Arrays

Work using nematic-phase liquid crystal electro-optic retarders is detailed comprehensively in a 1993 Air Force document.⁸⁶ The materials are known as types E7 and PTPP-33 liquid crystals, having birefringence $\Delta = (n_e - n_o) \approx 0.2$ in the infrared, requiring a cell be only 5 optical waves thick for a full-wave phase shift in transmission and only 2.5 waves thick in reflection. The thinner the cell, the shorter is its reorientation time. Switching speeds in the millisecond range with high-efficiency, diffraction-limited steering have been demonstrated at 10.6 μm with CO₂ lasers, and at 1.06 μm and 0.53 μm with Nd:YAG lasers. The cascading of tandem scanners by optical relaying⁸⁵ is a means for adding the contributions of two or more deflectors with each operating optimally. One-dimensional arrays may be compounded having one for azimuth and one for elevation. Also, individual deflectors requiring excessive spatial separation may be cascaded using relay optics to avoid walk-off of the beam from the second aperture by the action of the first deflector.⁸⁹

Another approach to tandem arrays,⁹⁰ named discrete/offset cascade, reduces potential “noise” (beam artifacts) n in the instances of large quantization mismatches when cascading phase-delayed groups. Experiments have demonstrated improved overall diffraction efficiency, along with requiring a reduced number of control lines. (A similar approach was demonstrated with microlens arrays.⁹¹) Also significant is the use of an electro-optic phase retarder other than liquid crystal. The material selected was PLZT (lead lanthanum zirconate titanate) which exhibits a large electro-optic coefficient, broadband optical transmission, very fast switching, and good thermal stability,⁹² and is a well-documented ceramic material, familiar in electro-optic modulation and deflection. Mirrored

piston-like phase adjustment is also reviewed,⁷⁶ and later work⁹³ describes both continuous phase change and binary phase shift.

Problems in broadband operation of phased arrays are reviewed⁷⁴ and early work was directed toward their solution.^{94,95} A wavelength-independent phase shift is achieved by polarization modulation of chiral smectic liquid crystals (CSLC), providing action similar to the mechanical rotation of a waveplate. However, grating dispersion remains due to wavelength deviation from nominal 2π phase resets, rendering a variation in efficiency η_d similar to Eq. (52)

$$\eta_d = \text{sinc}^2 \varepsilon \quad (58)$$

where ε is the chromatic error due to mismatch of the nominal 2π phase reset. Not only is energy lost, but side-lobe amplitudes increase and nondiffracted components result in image blurring and interference from sources outside the desired acceptance angle. This dispersion is reduced with the application of achromatic Fourier transform optics,⁹⁶ as investigated in more recent comprehensive work,⁹⁷ yielding precautions regarding the ability to reduce dispersion completely and the difficulty in implementation of the technology. More conventional achromatic optics has been applied⁹⁸ to the decentered lens; the second of the two major techniques for agile beam steering which is discussed next.

The Decentered Lens and Microlens Arrays

The decentering of one lens with respect to a matching afocal lens is an alternative to the phased array described above. Although its basic action differs from the phased systems, when smaller lenses are formed into a mating periodic array, the assembly can exhibit some of the characteristics of phased arrays, including functioning as blazed gratings.⁷⁵ However, single lens-group operation can avoid some of the image faults of array steering as later discussed.

Consider Fig. 50a illustrating a pair of afocal lenses (denoted 1 and 2) oriented originally on a common axis, now with lens 2 shifted “downward” through a distance Δ (dotted axes). Beyond the focus of lens (left lens), the diverging beam continues into lens 2 shifted angularly off its axis, resulting in deflecting the recollimated output beam through an angle θ_o . Thus, a transverse shift of one lens with respect to the other affects beam steering. The vignetting of the output beam and the related diversion of its residual output flux outside the lens is discussed subsequently. Constraining this simple two-lens technique is its limitation on the width of the lens aperture, consistent with the energy required for rapid Δ shift within a reasonable burden on acceleration of its more massive components.

Consider combining many lens pairs like lenses 1 and 2 (maintaining the F) to form an array of microlenses, as illustrated in Fig. 50b, and illuminating the group from the left by a single broad beam. The steered waves sum into the total field in a manner similar to those of the prior phased arrays and similar to Fig. 49b. This results in a significant decrease in mass for a given full aperture size and a decrease in shift distance Δ for the same steered angle. The effect is a reduction in Δ and in the acceleration/deceleration forces required for rapid beam steering. Although the composite wavefront is discontinuous, the segments are tipped at the same slopes such that the output exhibits the characteristics of a blazed grating. When, at the operating wavelength their junctions exhibit 2π phase differentials, a sawtooth pattern is formed typified by a high-diffraction efficiency blazed grating. This technique is satisfactory for small steered angles, where high fill factors remain at the second lens array, and the spurious components are a small residue. At wider steering angles, however, when the vignetting and the disruptive effects of the spurious components become significant, remedies are required.

A classic method for the control of such vignetting is to include a field lens⁹⁹ into the microlens array.⁷⁵ Fig. 50c illustrates this as a variation of Fig. 50a with a field lens (FL) inserted at the common focal plane of the original two lenses. The bar over the output pair of lenses represents physical connection for simultaneous shift. With equal focal lengths for all lenses, the expanding light cone completely fills the lens pair during Δ shift. This technique for realizing output efficiency and spectral quality is directly transferrable to the microlens array of Fig. 50b with an added plane of field

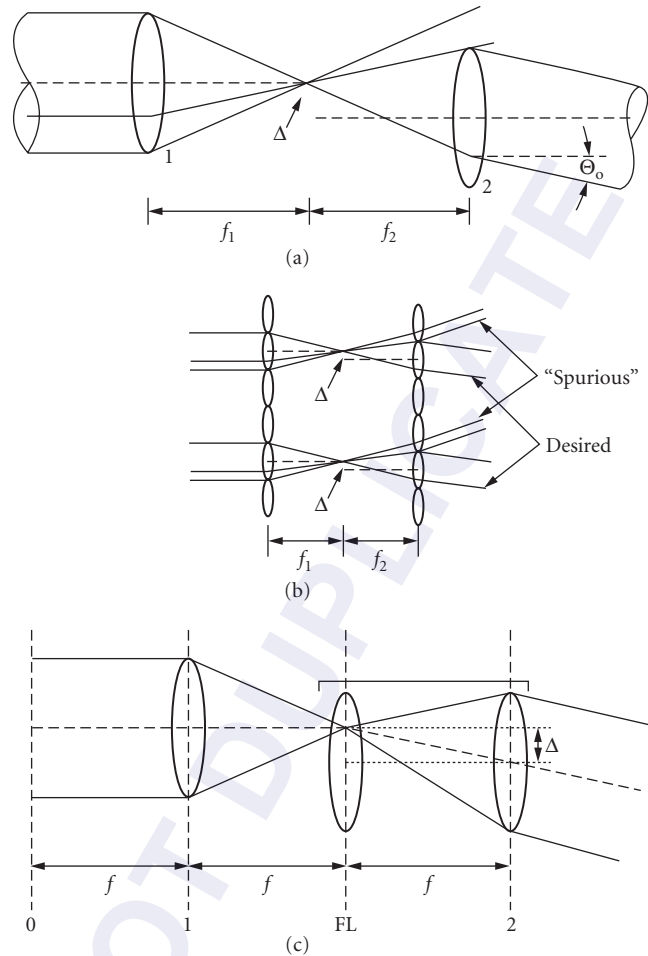


FIGURE 50 Beam steering with decentered lens of afocal pair shifted through distance Δ . (a) Single lens pair, showing Δ -shift deflecting major portion of output beam through angle Θ_o while upper portion of the beam by-passes lens 2. (b) Microlens arrays operating as in (a), but lighter and with smaller Δ -shift. The desired outputs accumulate, while the by-pass portions become spurious. (c) Field lens (FL) added to (a) provides constant filling of lens 2. When added to (b), the FLs maintain the wavefront synthesis of a blazed grating. (After Refs. 75 and 87.)

lenses affixed to the output array. The inertia can be accommodated by the force of piezoelectric or electrodynamic drive transducers. Alternatively, the single-element group may be shifted instead. A microlens-field lens design was fabricated and tested¹⁰⁰ over a $\pm 1.6^\circ$ field. Larger angles ($\pm 17^\circ$) have been demonstrated,¹⁰¹ but with loss of beam quality.

Further consideration for reduction of the spurious beams during Δ shift is represented in Fig. 51. The method of Fig. 51a provides⁷⁵ some tolerance for beam displacement on lens 2 by changing the ratio of focal lengths. The initial condition of $f_1 = f_2$ is adjusted to $f_1/f_2 > 1$. This forms a beam compressor (see Fig. 32) with a compression ratio $\approx 2:1$. A similar approach is investigated¹⁰² using

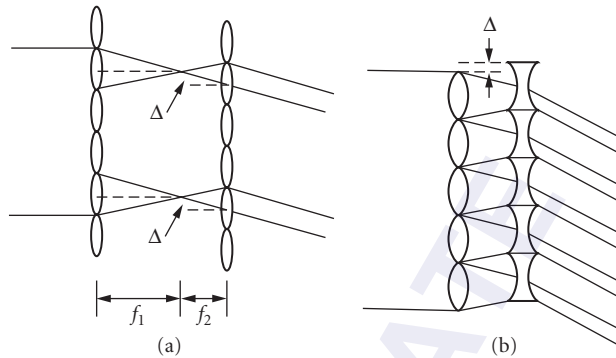


FIGURE 51 Elimination of spurious radiation over a range of operation with increased focal length ratio f_1/f_2 . (a) Similar to Fig. 50b with $f_1/f_2 \approx 2$. (b) Analogous galilean form with $f_1/f_2 \approx 2.5$ and the beam energy is conserved. However, the spaced output beams represent diffraction from a discontinuous blazed grating. (After Refs. 75, 102 and 87).

a positive-negative lens combination. Method in Fig. 51a employs the equivalent of a Keplerian telescope and method in Fig. 51b that of a Galilean telescope. While the spurious components of Fig. 50b may be abated over its initial range of operation, the fill factor at the second array is reduced. Although the energy is conserved in this reduced light cone, the ideal sawtooth pattern of the blazed grating is disrupted by the truncated sawtooth function. This, in turn, causes its own spurious noise^{75,102} which limits operation to a small range of Δ shift. It is proposed¹⁰² that the second array be maximally filled ideally by reducing the lens separation in Fig. 51b toward zero. Practically, this is approached with the development of thin binary optics microlens arrays. Binary forms of Fresnel zone patterns are fabricated utilizing high-resolution etching and transfer techniques formed on substrate materials. Hundred percent fill factors of lenslet arrays are attainable, approximating a continuous phase profile in a stepwise manner, to allow achievement of high diffraction efficiency. As presented earlier for a phased array composed of q elements per 2π phase reset [Eq. (52)], the efficiency η_b of a multilevel binary optic of m levels within one width of a Fresnel feature is given by¹⁰³ $\eta_b = \text{sinc}^2(1/m)$. An experimental system¹⁰² utilized such arrays of F/5 microlenses; each having a 0.2 mm diameter. The second lenslet array was spaced from the first by 10 μm , allowing relative translation. This system steered a 6-mm HeNe beam over an 11.5° field using ± 0.1 -mm travel at a 35-Hz sweep rate. Practical mask alignment, etch errors, and transfer errors during fabrication reduced the 95 percent theoretical efficiency to 84 percent and 72 percent for the positive and negative lens arrays, respectively. Overall efficiency of the unsteered beam measured approximately 50 percent. The F/5 system exhibited low efficiency when steered. This is expected to improve with more precise fabrication and operation at lower F-numbers.

A variation to the above work was conducted using “phased-arraylike” binary optics.⁷⁷ A continuous quadratic phase function was sampled at equal intervals of Δx , forming a stepwise matching of the continuous phase profile. The Δ shifts are conducted in integral increments of Δx , where the formerly disruptive region is shown analytically to render a continuous linear phase profile across the full aperture. Experimental binary micro-optics were designed to compare phased-arraylike and microlens arrays, fabricated simultaneously and adjacent on the same quartz substrate. Tests confirmed that the phased-arraylike structure provided approximately 50 percent increased intensity in the steered mode, and less than 1 percent leakage into the immediate (local) sidelobes. While stronger distant sidelobes developed, they were well separated from the steered mode.

As mentioned previously, achromatic optics applied to a single decentered lens system⁹⁸ reduces dispersion and allows operation in the 2 to 5 μm IR band. The prototype design of Fig. 50c was altered to avoid the high power intensity focus and its potential damage within the field lens FL.

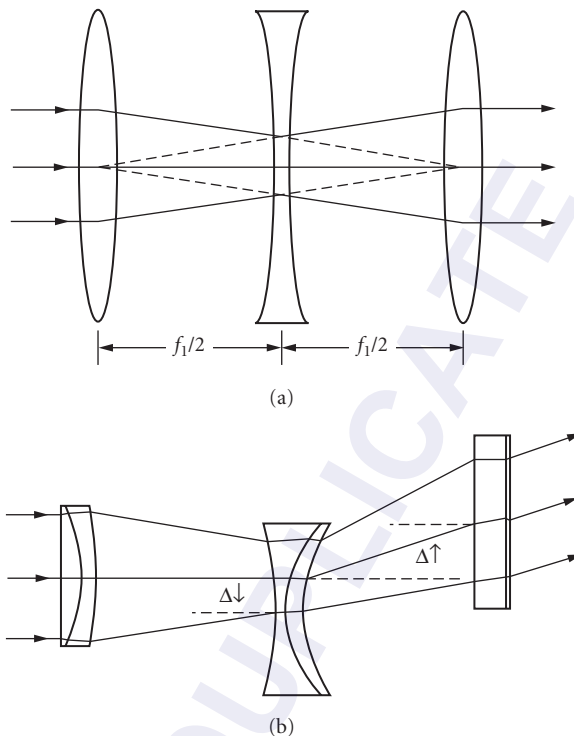


FIGURE 52 Adaptation for high power and achromatic operation. (a) Intense power at focus within the (positive) field lens FL of Fig. 50c is avoided by replacing it with a negative field lens (having $1/2$ the focal length of the positive lenses). (b) System (to scale) achromatized and optimized for beam deflection of 22.5° . The field lens is shifted downward through distance $\Delta\downarrow$, while the output lens is shifted upward through the greater distance $\Delta\uparrow$. (After Ref. 98.)

Thus, per Fig. 52a, the positive field lens of Fig. 50c was replaced with a negative one, requiring a focal length one-fourth that of the positive lens. The system was achromatized and optimized for $\pm 22.5^\circ$ maximum deflection, as represented in Fig. 52b. The negative field lens and the output lens now require unequal and *opposite* shift directions to implement scan. The 60-mm diameter output lens requires ≈ 30 mm of Δ shift. This work demonstrated that for some systems, a single group of three cemented achromatic doublets enables eliminating several problems associated with microlens arrays, such as spurious diffraction, multiple beam orders, blind spots, and large dispersion while accommodating relatively wide scan angles.

Digital Micromirror Device

In 1987, the first digital micromirror device (DMD) was created at Texas Instruments.^{104,105} Larry Hornbeck was granted the first patent for the DMD design in 1991.¹⁰⁶ This initial design was the basis for current digital micromirror device chips which were incorporated into digital light processing (DLP) projectors for both visible and infrared applications.^{107–109,110} Although the DMD is an all-digital spatial light modulator (SLM) rather than the more traditional optical scanners of this

chapter, it is briefly discussed as it can be considered a binary-state scanner. Its impact on the display industry has been most significant.

A DMD is an array of “light switches” having a MEMS structure that is fabricated by using micromachining CMOS processes over a CMOS static random access memory (SRAM) integrated circuit.^{110,111} Each light switch has an aluminum mirror that can reflect light in one of two directions depending on the state of the underlying memory cell. With the memory cell in the On state, the mirror rotates to $+10^\circ$. With the memory cell in the Off state, the mirror rotates to -10° . By combining the DMD with a suitable light source and projection optics, the mirror reflects incident light either into or out of the pupil of the projection lens. Thus, the On state of the mirror appears bright and the Off state of the mirror appears dark. Gray scale is achieved by binary pulsewidth modulation of the incident light while color is achieved by using color filters, either stationary or rotating, in combination with one, two, or three DMD chips.

The DMD light switch is a MEMS structure consisting of a mirror that is rigidly connected to an underlying yoke. The yoke in turn is connected by two thin, mechanically compliant torsion hinges to support posts that are attached to the underlying substrate. Electrostatic fields developed between the underlying memory cell and the yoke and mirror cause rotation in the positive or negative rotation direction. The rotation is limited by mechanical stops to typically $\pm 10^\circ$. The use of semiconductor and MEMS processing technologies allow production of very large arrays of these individually controllable micromirrors having minimal defects and high reliability.^{112–114}

Gimbal-Less Two-Axis Scanning Micromirrors

Gimbal-less two-axis scanning-micromirror devices (GSMD) have been recently developed by Mirrorcle Technologies, Inc. and are based on multilevel beam silicon-on-insulator micro-electromechanical (SOI-MEMS) fabrication technology.¹¹⁵ Due to their small scale and electrostatic actuation, these devices require ultralow power and can provide fast optical beam scanning in two axes when compared to the large-scale galvanometer-based optical scanners.¹¹⁶ Laser beams can be deflected to optical scanning angles of up to 32° at very high speeds in both axes. Continuous full-speed operation of the electrostatic actuators that drive the GSMD dissipates less than 1 mW of power. These devices are made entirely of monolithic single-crystal silicon, resulting in excellent repeatability and reliability. The flat, smooth mirror surfaces can be coated with a thin film of metal with desired reflectivity. Larger mirrors can be bonded onto actuators for custom aperture sizes. In contrast to the two-state mirror movement of the DMD, the GSMD mirror movement is fully analog and can maintain a selected tilt angle or move dynamically upon command. At this time, huge arrays of micromirrors comprise typical DMD arrays, while GSMD arrays are presently limited to a few micromirrors; however, GSMD mirrors can be physically much larger than the DMD mirrors.

The GSMD are designed and optimized for point-to-point optical beam scanning mode of operation. A steady-state analog actuation voltage results in a steady-state analog angle of rotation of the micromirror. Specifically, there is a one-to-one correspondence of actuation voltages and resulting angles that is highly repeatable. Positional precision in open-loop driving of the micromirrors is at least 14 bits, that is within $10\ \mu\text{rad}$. For a particular high-speed 3D tracking and position-measurement application, 16-bit precision has been demonstrated.¹¹⁷ Devices can be operated over a very wide bandwidth from DC (they maintain position at constant voltage) to several kilohertz. Such fast and broadband capability allows nearly arbitrary waveforms such as vector graphics, constant velocity scanning, point-to-point step scanning, and the like.¹¹⁸ The major advantage of the gimbal-less design is the capability to scan optical beams at equally high speeds in both axes. A typical GSMD with a 0.8-mm diameter micromirror achieves angular beam scanning of up to 600 rad/s and has first resonant frequency in both axes above 5 kHz. Devices with larger-diameter micromirrors are correspondingly slower due to the increased inertia.

The gimbal-less design combines one-axis electrostatic combdrive-based rotators,¹¹⁵ and allows their operation to be nearly independent of the other axis' operation without the added inertia of a gimbal frame. A schematic diagram of the conceptual operation of the gimbal-less two-dimensional

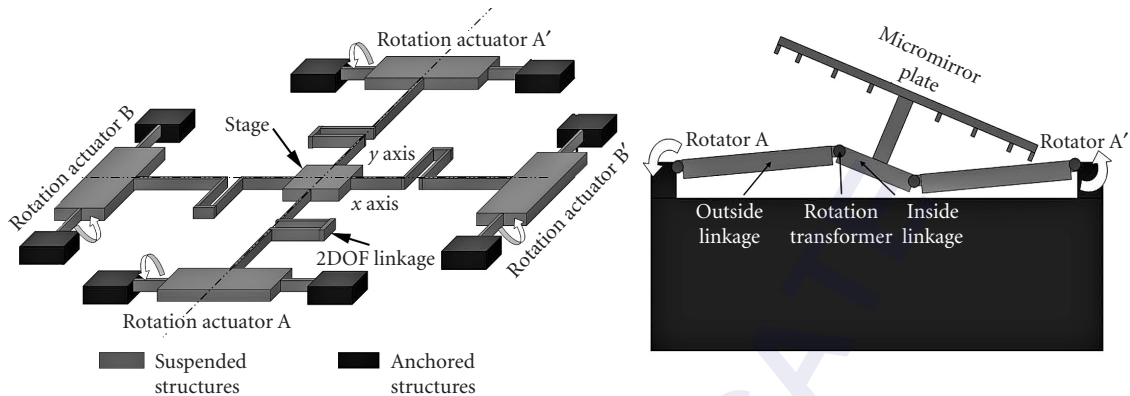


FIGURE 53 Schematic diagram of a gimbal-less two-axis scanning actuator based on four high aspect ratio rotators connected to the central pedestal by two degrees-of-freedom (2 DOF) linkages. Cross-sectional depiction of device operation.¹¹⁶ (Diagram courtesy of Mirroracle Technologies, Inc.).

designs is shown in Fig. 53. Two one-axis rotators are utilized for each axis of the overall two-dimensional scanner. For the x axis, actuators A and A' are utilized, and for the y axis, actuators B and B'. The inside linkages are designed such that they allow torsion on the axis, specifically during the operation of the orthogonal axis. In other words, each linkage that connects a rotator to the central micromirror is actually designed to be a two degree-of-freedom component.

GSMD can also operate in the dynamic, resonant mode. When operated near the resonant frequency, devices give significantly more angle at lower operating voltages and sinusoidal motion. Resonant frequencies are in the range of several kilohertz, although in some cases one of the axes is made exceptionally stiff to achieve 16 kHz and faster actuation for video projection applications.¹¹⁹

The gimbal-less design lends itself inherently to a modular design approach, hence, several types of dual-axis actuator designs are available.¹²⁰ Each actuator can utilize rotators of arbitrary length, arbitrarily stiff linkages, and arbitrarily positioned mechanical rotation transformers. In addition, the GSMD can have an arbitrarily large mirror diameter. Because of modularity, these devices can be customized for the requirements of a particular application.

Silicon mirrors of up to 1.2 mm diameter can be fabricated as an integral (monolithic) part of some GSMDs. Due to the limitations of the fabrication steps of the actuator, the standard mirrors are relatively thick (24 μm .) The inherent properties of the single crystal silicon substrate yield a polished surface with nearly perfect flatness. Larger and customizable mirror sizes and shapes can be utilized by fabricating those separately and assembling them on top of gimbal-less actuators. A SEM image of a 2-mm diameter bonded mirror on an actuator is shown in Fig. 54. Sets of electrostatic actuators optimized for speed, angle, area footprint, or resonant driving are designed and realized in a self-aligned deep reactive ion etching (DRIE) fabrication process.¹²¹ Metalized, ultralow inertia, single crystal mirrors stiffened by a backbone of thicker silicon beams (see Fig. 53) are created in a separate fabrication process. The diameter, as well as geometry, of the mirror is selected by customers, in order to optimize the trade-offs between speed, beam size, and scan angle for each individual application. The mirrors are subsequently bonded to the actuators. The modular approach allows either the absolute optimization of a device prior to fabrication, or the ability to economically adapt a small set of fabricated devices for a wide range of applications. Larger sizes up to 3.6 mm are regularly assembled in experiments and applications.

Summary of Agile Beam Steering

Two basic methods for low-inertia laser beam steering are presented: the phased array and the decentered lens. Although their operating principles differ, as arrays, they both form diffraction

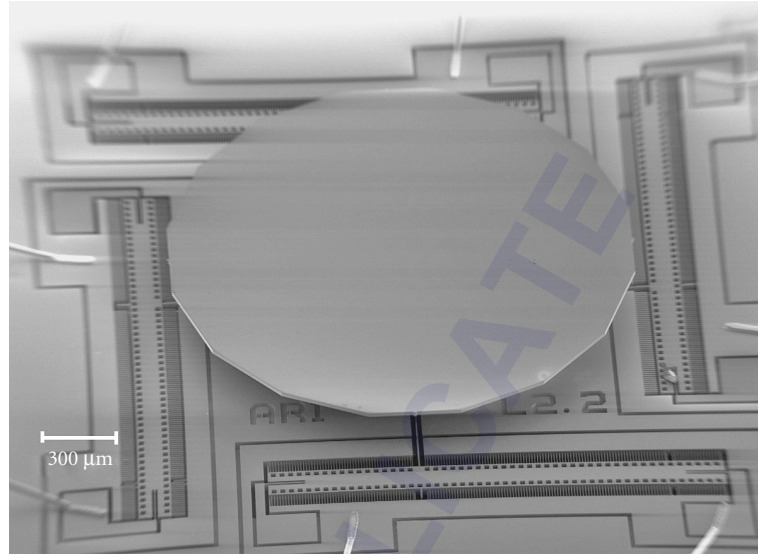


FIGURE 54 SEM image of a 2.0-mm diameter micromirror bonded to dual-axis actuators. (Photograph courtesy of Mirrorcle Technologies, Inc.)

gratings yielding output wavefronts having the properties of a blazed grating. While this results in high diffraction-efficiency at a selected wavelength, typical grating dispersion limits broadband operation. Complicating corrective measures have been applied to both array systems to approach broadband operation. The single decentered lens group avoids array dispersion. And, if its increased size and inertia can be tolerated, lenticular dispersion remains controllable with familiar achromatizing techniques. The dominant phased arrays are completely electro-optic, while the microlens array requires very small translations of reasonably low inertia assemblies. An alternative phased array utilizes individual micromirrors requiring minute ($\lambda/4$) axial displacements. However, difficulty may be encountered in fabrication of high optical-integrity and in providing the high optical-fill-factor of the electro-optic types. Another alternative microlens array is formed of Fresnel-lens-type binary optics. For low F-number lenticules, whose theoretical steering efficiency is high, their minimum feature sizes become miniscule and presently are difficult to fabricate.

Auxiliary control facilities can impose burdens of mass, volume, and cost. Phased arrays utilize complex multielement electrical programming, while the lens arrays require small but very precise positioning of their assemblies. Though such additional requirements are generally not detailed in the literature, a comparative analysis¹²¹ provided some related observations. The authors preferred a microlens array over the liquid crystal phased array, thereby avoiding the “heavy burden” on electronic control of the many phase-delay elements. An x - y microlens array system was designed and built for test, and was compared to a two-galvanometer x - y system assembled of commercial components. Evaluations confirmed that the microlens system steered faster, consumed lower power, and packaged smaller and lighter. However, no comment appeared on design to minimize mirror inertia and to reduce the bulk of the components and their assembly. Nor was x - y relay optics⁵ considered to allow minimum-sized mirrors serve to reduce inertia. Also, meriting evaluation is the *single mirror* suspended and actuated in x - y for precise two-dimensional scan.^{122,123} Such diverse considerations are invaluable for rating design alternatives for their relative compliance to system requirements.

30.9 REFERENCES

1. S. Sherr, *Fundamentals of Display System Design*, John Wiley & Sons, New York, 1970.
2. S. Sherr, *Electronic Displays*, John Wiley & Sons, New York, 1979.
3. L. Beiser, "Unified Approach to Photographic Recording from the Cathode-Ray Tube," *Photo. Sci. Engr.* **7**(3):196–204 (1963).
4. L. Beiser, *Laser Scanning Notebook*, SPIE Press, Washington, vol. PM13, 1992.
5. L. Beiser, "Laser Scanning Systems," in *Laser Applications*, vol. 2. Academic Press, New York, 1974, pp. 53–159.
6. L. Beiser, *Holographic Scanning*, John Wiley & Sons, New York, 1988.
7. J. D. Zook, "Light Beam Deflector Performance: A Comparative Analysis," *Appl. Opt.* **13**(4):875–887 (1974).
8. L. Beiser (ed.), "A Guide and Nomograph for Laser Scanned Resolution," *Proc. SPIE* **1079**:2–5 (1989).
9. G. Boreman, "Modulation Transfer Techniques," Chap. 32, *Handbook of Optics*, vol. II, McGraw Hill, New York, 1995.
10. L. R. Dickson, "Characteristics of Propagating Gaussian Beams," *Appl. Opt.* **9**(8):1854–1861 (1970).
11. J. Randolph and J. Morrison, "Rayleigh-Equivalent Resolution of Acoustooptic Deflection Cells," *Appl. Opt.* **10**(6):1383–1385 (1971).
12. L. R. Dickson, "Optical Considerations for Acoustooptic Deflectors," *Appl. Opt.* **11**(10):2196–2202 (1972).
13. L. Levy, *Applied Optics*, vol. 1, John Wiley & Sons, New York, 1968.
14. L. Beiser, "Generalized Equations for the Resolution of Laser Scanners," *Appl. Opt.* **22**(20):3149–3151 (1983).
15. J. C. Urbach, T. S. Fisli, and G. K. Starkweather, "Laser Scanning for Electronic Printing," *Proc. IEEE* **70**(6):597–618 (1982).
16. H. V. Kennedy, Split Image, "High Scanning Rate Optical System with Constant Aperture," U.S. Patent 3,211,046, 1965.
17. H. B. Henderson, "Scanner System," U.S. Patent 3,632,870, 1972.
18. W. L. Wolfe, "Optical-Mechanical Scanning Techniques and Devices," Chap. 10, in W. L. Wolfe and G. Zissis (eds.), *The Infrared Handbook*, ERIM, Ann Arbor, Mich., 1989; U.S. Office of Naval Research, Washington, Fig. 10-12, p. 10–12.
19. *Manual of Remote Sensing*, 2nd ed., Amer. Soc. of Photogrammetry, Falls Church, 1983, Figs. 3–8, p. 340.
20. J. M. Lloyd, *Thermal Imaging Systems*, Plenum Press, New York, 1975, pp. 308–309.
21. R. B. Barnes, "Infrared Thermogram Camera and Scanning Means Theorem," U.S. Patent 3,287,559, 1966.
22. B. A. Wheeler, "Optical Raster Scan Generator," U.S. Patent 3,764,192, 1973.
23. R. B. Johnson, "Target-Scanning Camera Comprising a Constant Temperature Source for Providing a Calibrated Signal," U.S. Patent 3,631,248, 1971.
24. P. J. Lindberg, "A Prism Line-Scanner for High Speed Thermography," *Optica Acta* **15**:305–316 (1966).
25. P. J. Lindberg, "Scanning Mechanism for Electromagnetic Radiation," U.S. Patent 3,253,498, 1966.
26. R. B. Johnson, "Image Rotation Device for an Infrared Scanning System or the Like," U.S. Patent 3,813,552, 1974.
27. P. J. Berry and H. M. Runciman, "Radiation Scanning System," U.S. Patent 4,210,810, 1980.
28. T. H. Jamieson, "Optical Design of Compact Thermal Scanner," *Proc. SPIE* **518**:15–21 (1984).
29. J. E. Modisette and R. B. Johnson, "High Speed Infrared Imaging System," U.S. Patent 4,419,692, 1983.
30. D. C. O'Shea, *Elements of Modern Optical Design*, John Wiley & Sons, New York, 1985.
31. K. S. Petche, *Handbook of Optics*, OSA, W. G. Driscoll and W. Vaughn (eds.), McGraw Hill, New York, 1978, pp. 5–10.
32. Y. Li and J. Katz, "Encircled Energy of Laser Diode Beams," *Appl. Opt.* **30**(30):4283 (1991).
33. R. E. Hopkins and D. Stephenson, "Optical Systems for Laser Scanners," in G. F. Marshall (ed.), *Optical Scanning*, Marcel Dekker, New York, 1991.
34. A. B. Marchant, U.S. Patent 4,759,616, 1988; and A. W. Lohman and W. Stork, "Modified Brewster Telescope," *Appl. Opt.* **28**(7):1318–1319 (1989).

35. L. Beiser, "Design Equations for a Polygon Laser Scanner," in G. F. Marshall and L. Beiser (eds.), *Proc. SPIE* **1454**:60–66 (1991).
36. D. Kessler, "High Resolution Laser Writer," in L. Beiser (ed.), *Proc. SPIE* **1079**:27–35 (1989).
37. C. J. Kramer, "Holographic Deflectors for Graphic Arts Systems," in G. F. Marshall (ed.), *Optical Scanning*, Marcel Dekker, New York, 1991.
38. R. V. Pole and H. P. Wollenmann, "Holographic Laser Beam Deflector," *Appl. Opt.* **14**(4):976–980 (1975).
39. I. Cindrich, "Image Scanning by Rotation of a Hologram," *Appl. Opt.* **6**(9):1531–1534 (1967).
40. C. J. Kramer, "Holo-Scanner for Reconstructing a Scanning Light Spot Insensitive to Mechanical Wobble," U.S. Patent 4,239,326, 1980.
41. D. B. Kay, "Optical Scanning System with Wavelength Shift Correction," U.S. Patent 4,428,643, 1984.
42. H. Ikeda, et al., "Hologram Scanner," *Fujitsu Sci. Tech J.* **23**:3 (1987).
43. F. Yamagishi, S. Hasegawa, H. Ikeda, and T. Inagaki, "Lensless Holographic Line Scanner," *Proc. SPIE* **615**:128–132 (1986).
44. J. Montagu, "Galvanometric and Resonant Low Inertia Scanners," in G. F. Marshall (ed.), *Optical Scanning*, Marcel Dekker, New York, 1991.
45. S. Reich, "Use of Electro-Mechanical Mirror Scanning Devices," in L. Beiser (ed.), *SPIE Milestone Series* **378**:229–238 (1985).
46. D. G. Tweed, "Resonant Scanner Linearization Techniques," *Opt. Eng.* **24**(6):1018–1022 (1985).
47. F. Blais, "Control of Galvanometers for High Precision Laser Scanning Systems," *Opt. Eng.* **27**(2):104–110 (1988).
48. G. F. Marshall and J. S. Gadhok, "Resonant and Galvanometer Scanners: Integral Position Sensing," *Photonics Spectra*, 155–160 (1991).
49. P. Debye and F. W. Sears, "On the Scattering of Light by Supersonic Waves," *Proc. Natl. Acad. Sci. (USA)* **18**:409 (1932).
50. F. Okolicsanyi, "The Waveslot, an Optical TV System," *Wireless Eng.* **14**:536–572 (1937).
51. R. Adler, "Interaction between Light and Sound," *IEEE Spectrum* **4**(5):42–54 (1967).
52. E. I. Gordon, "Review of Acoustooptical Deflection and Modulation Devices," *Proc. IEEE* **54**(10):1391–1401 (1966).
53. A. Korpel, "Acousto-Optics," in R. Kingslake and B. J. Thompson (eds.), *Appl. Opt. and Opt. Eng.* VI, Academic Press, New York, 1980, pp. 89–141.
54. R. V. Johnson, "Scophony Light Valve," *Appl. Opt.* **18**(23):4030–4038 (1979).
55. L. Bedamian, "Acousto-Optic Laser Recording," *Opt. Eng.* **20**(1):143–149 (1981).
56. W. H. Watson and A. Korpel, "Equalization of Acoustooptical Deflection Cells in a Laser Color TV System," *Appl. Opt.* **7**(5):1176–1179 (1970).
57. M. Gottlieb, "Acousto-Optical Scanners and Modulators," in G. Marshall (ed.), *Optical Scanning*, Marcel Dekker, New York, 1991.
58. D. L. Hecht, "Multibeam Acoustooptical and Electrooptic Modulators," in L. Beiser (ed.), *Proc. SPIE* **396**:2–9 (1983).
59. L. Beiser, "Generalized Gradient Deflector and Consequences of Scan of Convergent Light," *J. Opt. Soc. Am.* **57**:923–931 (1967).
60. A. J. Giarola and T. R. Billeter, "Electroacoustic Deflection of a Coherent Light Beam," *Proc. IEEE* **51**:1150 (1963).
61. V. J. Fowler, et al., "Electro-Optical Light Beam Deflection," *Proc. IEEE* **52**:193 (1964).
62. J. F. Lospeich, "Electro-Optical Light Beam Deflections," *IEEE Spectrum* **5**(2):45–52 (1968).
63. R. J. Pizzo and R. J. Kocka, Conoptics Inc., Danbury, CT 06810. Available at <http://www.conoptics.com/>.
64. T. C. Lee and J. D. Zook, "Light Beam Deflection with Electrooptic Prisms," *IEEE J. Quant. Electron.* **QE-4**(7):442–454 (1968).
65. C. S. Tsai and J. M. Kraushaar, "Electro-Optical Multiplexers and Demultiplexers for Time-Multiplexed PCM Laser Communication Systems," *Proc. Electro-Opt. Syst. Des. Conf.* 176–182 (1972).
66. A. S. Vasilevskaya, *Soviet Phys. Crystallogr.* **12**(2):308 (1967).

67. J. D. Zook and T. C. Lee, "Design of Analog Electrooptic Deflectors," *Proc. SPIE* (of 1969) **281** (1970).
68. J. P. Donohue, "Laser Pattern Generator for Printed Circuit Board Artwork Generation," *Proc. SPIE* **200**:179–186 (1979); also in L. Beiser (ed.), *SPIE Milestone Series* **378**:421–428 (1985).
69. G. K. Starkweather, "Single Facet Wobble Free Scanner," U.S. Patent 4,475,787, 1984.
70. L. Beiser, "Light Scanner," U.S. Patent 4,936,643, 1990.
71. G. F. Marshall, T. Vettese, and J. H. Carosella, "Butterfly Line Scanner," G. F. Marshall and L. Beiser (eds.), *Proc. SPIE* **1454** (1991).
72. L. Beiser, "Double Reflection Light Scanner," U.S. Patent 5,114,217, 1992.
73. D. F. Hanson and R. J. Sherman, "Method and Apparatus for Generating Optical Scans," U.S. Patent 4,433,894, Feb. 28, 1984; and R. J. Garwin, "Optical Alignment Compensation," U.S. Patent 4,429,948, Feb. 7, 1984.
74. P. F. McManamon, T. A. Dorschner, D. L. Corkum, L. J. Friedman, D. S. Hobbs, M. Holz, S. Liberman, et al., "Optical Phased Array Technology," *Proc. IEEE* **84**(2):268–298 (Feb. 1996).
75. E. A. Watson, "Analysis of Beam Steering with Decentered Microlens Arrays," *Opt. Eng.* **32**(11):2665–2670 (Nov. 1993).
76. E. A. Watson and A. R. Miller, "Analysis of Optical Beam Steering Using Phased Micromirror Arrays," *SPIE* **2687**:60–67 (1966).
77. M. W. Farn, "Agile Beam Steering Using Phased Array-Like Binary Optics," *Appl. Opt.* **33**(22):5151–5158 (Aug. 1994).
78. T. C. Cheston and J. Frank, *Radar Handbook*, M. I. Skolnick (ed.) McGraw Hill, New York, 1970.
79. M. I. Skolnick, *Introduction to Radar Systems*, McGraw Hill, New York, 1962.
80. L. W. Procopio, F. A. Jessen, and L. J. Brown., "Laser Phased Arrays," *Proc. 8th Int. Conf. Mil. Electron. Group* p. 67; RADC contract AF30(602)-2901 (1964).
81. L. J. Alpet, et al., "Laser Array Techniques," *Quart. Rep. 3A*, Air Force contract AF33(615)3918, AD No. 821235 (1967).
82. P. G. Grant, R. A. Meyer, and D. N. Qualkinbush, "An Optical Phased Array Beam Steering Technique," *Proc. Electr.-Opt. Syst. Des. Conf.* 259–264 (1971).
83. R. A. Meyer, "Optical Beam Steering Using a Multichannel Lithium Tantalate Crystal," *Appl. Opt.* **11**(3):613 (Mar. 1972).
84. A. S. Keys, R. L. Fork, T. R. Nelson, and J. P. Loehr, "Resonant Transmissive Modulator Construction for use in Beam Steering Arrays", in L. Beiser, S. Sagan and G. F. Marshall (eds.), *Optical Scanning: Design and Applications*, *SPIE* **3787**:115–125 (Jul. 1999).
85. P. F. McManamon and E. A. Watson, "Optical Beam Steering Using Phased Array Technology," in L. Beiser and S. Sagan (eds.), *SPIE* **3131**:90–98 (1992).
86. T. A. Dorschner, R. C. Sharp, D. P. Resler, L. J. Friedman, and D. C. Hobbs, "Basic Laser Beam Agility Techniques," *Wright Patterson Air Force Base*, WL-TR-93-1020 (AD-B-175-883) (1993).
87. L. Beiser, *Unified Optical Scanning Technology*, John Wiley & Sons, New York, 2003.
88. E. A. Watson, D. T. Miller, and P. F. McManamon, "Applications and Requirements for Nonmechanical Beam Steering in Active Electro-Optic Modulators," in I. Cindrich, S. H. Lee, and R. L. Sutherland (eds.), *SPIE* **3633**:216–225 (1999).
89. L. Beiser, "Laser Scanning Systems," *Laser Applications*, vol. 2, M. Ross (ed.), Academic Press, New York, pp. 55–159, 1974.
90. J. A. Thomas, M. E. Lasher, Y. Fainman, and P. Soltan, "A PLZT-Based Dynamic Diffractive Optical Element for High Speed Random Access Beam Steering," in L. Beiser and S. Sagan (eds.), *SPIE* **3131**:124–132 (1997).
91. K. M. Flood, W. J. Cassarly, C. Sigg, and M. J. Finlan, "Continuous Wide Angle Beam Steering Using Translation of Binary Microlens Arrays and a Liquid Crystal Phased Array," *SPIE* **1211**:296–304 (1990).
92. G. H. Haertling, "PLZT Electrooptic Materials and Applications—A Review," *Ferroelectrics* **75**:25–55 (1987).
93. D. M. Burns, V. M. Bright, S. C. Gustafson and E. A. Watson, "Optical Beam Steering Using Surface Machined Gratings and Optical Phased Arrays," in L. Beiser and S. Sagan (eds.), *SPIE* **3131**:99–110 (1997).
94. E. A. Watson, P. R. McManamon, L. J. Barnes, and A. J. Carney, "Applications of Dynamic Gratings to Broad Spectral Band Beam Steering," *SPIE* **2120** (1994).

95. J. E. Stockley, S. A. Serati, G. D. Sharp, P. Wang, K. F. Walsh, and K. M. Johnson, "Broad-Band Beam Steering," in *Optical Scanning Systems: Design and Applications*, L. Beiser and S. Sagan (eds.), *Proc. SPIE* **3131**:111–123 (1997).
96. G. M. Morris and D. A. Zweig, "White Light Fourier Transforms," in J. J. Hovner (ed.), *Optical Signal Processing*, Chap. 1.2, Academic Press, New York, 1987.
97. P. F. McManamon, E. A. Watson, S. Jianru, and P. J. Bos, "Nonmechanical Steering of the Field of View of Broad Spectral Band Optical Systems," in *Optical Scanning 2005*, S. F. Sagan and G. F. Marshall (eds.), *SPIE* **5873**:26–37 (2005).
98. J. L. Gibson, B. D. Duncan, E. A. Watson, and J. S. Loomis, "Wide Angle Decentered Lens Beam Steering for Infrared Countermeasures Applications," *Opt. Eng.* **43**(10):2312–2321 (Oct. 2004).
99. L. Levi, *Applied Optics*, vol. 1, John Wiley & Sons, New York 1968.
100. T. D. Milster and J. N. Wang, "Modeling and Measurement of a Micro-Optic Beam Deflector," in Y. Kohanzadeh, G. W. Lawrence, J. G. McCoy, and H. Weichel (eds.), *SPIE* **1625**:78–83 (1992).
101. E. A. Watson, W. E. Whitaker, C. D. Brewer, and S. R. Harris, "Implementing Optical Phased Array Beam Steering with Cascaded Microlens Arrays," *IEEE Aerospace Conf.* Paper No. 5.036, Big Sky, MT (March 2002).
102. W. Goltos and M. Holz, "Agile Beam Steering Using Binary Optics Microlens Arrays," *Opt. Eng.* **29**(11):1392–1397 (Nov. 1990).
103. G. J. Swanson, "Binary Optics Technology: Theory and Design of Multi-Level Diffractive Optics Elements," *Lincoln Labs Tech. Rep.* **854** (Aug. 14, 1989).
104. L. J. Hornbeck, "Deformable-Mirror Spatial Light Modulators," *Spatial Light Modulators and Applications I*, *SPIE Critical Reviews* **1150**:86 (Aug. 1989).
105. J. Younse, "Mirrors on a Chip," *IEEE Spectrum* **30**(11):27 (Nov. 1993).
106. L. J. Hornbeck, "Spatial Light Modulator and Method," U.S. Patent 5,061,049 (1991).
107. P. F. Van Kessel, L. J. Hornbeck, R. E. Meier, and M. R. Douglass, "A MEMS-Based Projection Display," *Proc. IEEE*, **86**(8):1687–1704 (Aug. 1998).
108. J. M. Younse, "Projection Display Systems Based on the Digital Micromirror Device (DMD)," *SPIE Conference on Microelectronic Structures and Microelectromechanical Devices for Optical Processing and Multimedia Applications*, Austin, Texas, *SPIE Proc.* **2641**:64–75 (Oct. 1995).
109. J. B. Sampsel, "An Overview of the Digital Micromirror Device (DMD) and Its Application to Projection Displays," *1993 SID International Symposium Digest of Technical Papers*, **24**:1012 (1993).
110. L. J. Hornbeck, "Digital Light Processing™ and MEMS: Timely Convergence for a Bright Future," Plenary Session, *SPIE Micromachining and Microfabrication '95*, Austin, TX, October 24, 1995, p. 2.
111. P. F. Van Kessel, L. J. Hornbeck, R. E. Meier, M. R. Douglass, "A MEMS-Based Projection Display," *Proc. IEEE*, **86**(8):1687–1704, (Aug. 1998).
112. L. J. Hornbeck, "From Cathode Rays to Digital Micromirrors: A History of Electronic Projection Display Technology," Texas Instruments White Paper, <http://www.dlp.com/tech/research.aspx> (Sept. 1998).
113. L. J. Hornbeck, "Digital Light Processing™: A New MEMS-Based Display Technology," Texas Instruments White Paper, <http://www.dlp.com/tech/research.aspx> (2008).
114. M. R. Douglass, "Development of the Digital Micromirror Device™ (DMD™) Microsystem," Texas Instruments White Paper, <http://www.dlp.com/tech/research.aspx>, Accessed on Sept. 2005.
115. V. Milanović, "Multilevel-Beam SOI-MEMS Fabrication and Applications," *IEEE/ASME J. Microelectromechanical Systems*, **13**(1):19–30, (Feb. 2004).
116. V. Milanović, D. T. McCormick, and G. Matus, "Gimbal-Less Monolithic Silicon Actuators For Tip-Tilt-Piston Micromirror Applications," *IEEE J. Select Topics in Quantum Electron.* **10**(3):462–471, (May-Jun. 2004).
117. V. Milanović, and W. K. Lo, "Fast and High-Precision 3D Tracking and Position Measurement with MEMS Micromirrors," *2008 IEEE/LEOS Optical MEMS and Their Applications Conf.*, Freiburg, Germany, Aug. 12, 2008.
118. V. Milanović, K. Castelino, and D. McCormick, "Highly Adaptable MEMS-based Display with Wide Projection Angle," *2007 IEEE Int. Conf. on Microelectromechanical Systems (MEMS'07)*, Kobe, Japan, Jan. 25, 2007.
119. V. Milanović, "Improved Control of the Vertical Axis Scan for MEMS Projection Displays," *2007 IEEE/LEOS Optical MEMS and Their Applications Conf.*, Hualien, Taiwan, Aug. 12, 2007.

120. D. T. McCormick, V. Milanović, and K. Castelino, "A Modular Custom Aperture Technology for Optimization of MEMS Scanners," *2005 IEEE/LEOS Optical MEMS and Their Applications Conf.*, Oulu, Finland, Aug. 2005, pp. 29–30.
121. Devoe, Donald Lad, and Tsai, Lung-Wen, "Process for Fabrication of 3-Dimensional Micromechanisms," U.S. Patent 6,664,126 (2003).
122. G. F. McDearmon, K. M. Flood, and J. M. Finlan, "Comparison of Conventional and Microlens-Array Agile Beam Steerers," in L. Beiser and M. Modamedi (eds.), *Proc. SPIE 2383*:167–178 (Feb. 1995).
123. A. Berta, "Development of High Performance Fast Steering Mirrors," in *Optical Scanning: Design and Applications*, L. Beiser, S. Sagan, and G. F. Marshall (eds.), *Proc. SPIE 3787*:115–125 (Jul. 1999).

30.10 FURTHER READING

The following listing augments the text and its references with a select group of publications, many arcane and of historic value, representing substantive work in the field.

- Barnes, C. W., "Direct Optical Data Readout and Transmission at High Data Rates," AFAL-TR-65-45, Wright Pat. A. F. B., AD-460486, 1965.
- Becker, C. H., "Coherent Light Recording Techniques," RADC-TR-65-130, AFSC, Griffiss AFB, New York, 1965.
- Beiser, L., "Laser Beam Scanning for High Density Data Extraction and Storage," *Proc. SPSE Symp. on Photography in Information Storage and Retrieval*, Washington, D.C., 1965.
- Beiser, L., "Laser Beam and Electron Beam Extremely Wideband Information Storage and Retrieval," *Photo. Sci. Eng.* **10**(4):222 (1966).
- Beiser, L., "Laser Scanning Systems," in M. Ross (ed.), *Laser Applications*, vol. 2, Academic Press, New York, 1974.
- Beiser, L. (ed.), "Selected Papers on Laser Scanning and Recording," *SPIE Milestone Series*, vol. 378, 1985.
- Beiser, L., *Holographic Scanning*, John Wiley & Sons, New York, 1988.
- Bousky, S., "Scanning Techniques for Light Modulation Recording," Wright Pat. AFB, Ohio, Contr. No. AF33(615)–2632, 1965.
- Cindrlich, I., et al., "Wide Swath Wide Bandwidth Recording," AFAL-TR-11-361, Wright Pat. AFB, Ohio, DDC No. 13217, AD 513171, 1971.
- Dubovic, A. S., *Bases of Mirror Scanning Theory* (trans from Russian), AD-703569, Wright Pat. AFB, F. T. D. HT-23-654-69, 1969.
- Dubovic, A. S., *The Photographic Recording of High Speed Processes*, John Wiley & Sons, New York, 1981.
- Fowler, V. J., "Survey of Laser Beam Deflection Techniques," *Appl. Opt.* **5**(10):1675 (1966).
- Johnson, R. B. and W. L. Wolfe (eds.), "Selected Papers on Infrared Design," *SPIE Milestone Series*, vol. 513, 1985.
- Katys, G. P., *Automatic Scanning* (trans. from Russian, 1969), National Tech Info. Service, JPRS 51512, 1970.
- Korpel, A., et al., "A TV Display Having Acoustic Deflection and Modulation of Coherent Light," *Appl. Opt.* **5**(10):667 (1966).
- Klein, D. A., "Laser Imaging Techniques," Naval Air Syst. Comm. TR-1615 (AIR-52022), Naval Avionics Facility, Indiana 46218, 1970.
- Laemmel, A. E., "The Scanning Process in Picture Transmission," Polytech. Inst. of B'klyn, RADC-TN-57-215, DDC AD-131112, 1957.
- Marshall, G. F. (ed.), *Laser Beam Scanning*, Marcel Dekker, New York, 1985.
- Marshall, G. F. (ed.), *Optical Scanning*, Marcel Dekker, New York, 1991.
- Newgard, P. M. and Brain, E. E., "Optical Transmission at High Data Rates," AFAL-TR-67-217, Wright. Pat. AFB, Ohio, 1967.
- Townsend, C. V. R., "Study of Advanced Scanning Methods," ASD-TDR-63-595, DDC 412794, Wright Pat. AFB, Ohio, 1963.
- Wolfe, W. L. and G. J. Zissis (eds.), *The Infrared Handbook*, ERIM, Ann Arbor, Mich., 1989.
- Zook, J. D., "Light Beam Deflector Performance," *Appl. Opt.* **13**:875 (1974).

Brian Henderson

*Department of Physics and Applied Physics
University of Strathclyde
Glasgow, United Kingdom*

31.1 GLOSSARY

A_{ba}	Einstein coefficient for spontaneous emission
a_o	Bohr radius
B_{if}	Einstein coefficient for transition between initial state $ i\rangle$ and final state $ f\rangle$
e	electron charge
ED	electric dipole term
E_{DC}	Dirac Coulomb term
E_{hf}	hyperfine energy
E_n	eigenvalues of quantum state n
EQ	electric quadrupole term
$\mathbf{E}(t)$	electric field at time t
$\mathbf{E}(\omega)$	electric field at frequency ω
g_a	degeneracy of ground level
g_b	degeneracy of excited level
g_N	gyromagnetic ratio of nucleus
h	Planck's constant
H_{so}	spin-orbit interaction Hamiltonian
I	nuclear spin
$I(t)$	emission intensity at time t
\mathbf{j}	total angular momentum vector given by $\mathbf{j} = l \pm 1/2$
l_i	orbital state
m	electron mass
MD	magnetic dipole term
$n_\omega(T)$	equilibrium number of photons in a blackbody cavity radiator at angular frequency ω and temperature T

QED	quantum electrodynamics
$R_{nl}^{(r)}$	radial wave function
R_{∞}	Rydberg constant for an infinitely heavy nucleus
s	spin quantum number with the value 1/2
s_i	electronic spin
T	absolute temperature
W_{ab}	transition rate in absorption transition between states a and b
W_{ba}	transition rate in emission transition from state b to state a
Z	charge on the nucleus
$\alpha = e^2/4\pi\epsilon_0\hbar c$	fine structure constant
$\Delta\omega$	natural linewidth of transition
$\Delta\omega_D$	Doppler width of transition
ϵ_0	permittivity of free space
$\zeta(r)$	spin-orbit parameter
μ_B	Bohr magneton
$\rho(\omega)$	energy-density at frequency ω
τ_R	radiative lifetime
ω	angular frequency
ω_k	mode k with angular frequency ω
$\langle f V' i\rangle$	matrix element of perturbation V

31.2 INTRODUCTION

This chapter outlines the physical basis of optical measurements in the wavelength/frequency and time domains. From the multiplicity of different apparatus, only simple examples are given of spectrometers designed for optical absorption, photoluminescence, and radiative decay measurements. Rather more detailed expositions are given of modern developments in laser spectrometers especially where high resolution is possible in both frequency and time domains. Included are specific developments for linewidth measurements down to tens of kilohertz using saturated absorption techniques as well as temporal decay characteristics in the sub-picosecond domain. A description is also given of a multiple resonance spectrometer including optically detected electron spin resonance and optically detected electron nuclear double resonance.

31.3 OPTICAL ABSORPTION SPECTROMETERS

General Principles

In optical absorption spectroscopy, electromagnetic radiation in the near-ultraviolet, visible, or near-infrared regions is used to excite transitions between the electronic states. Whereas atoms in low-pressure gas discharges exhibit very sharp lines, electronic centers in molecules and condensed matter display a variety of different bandshapes. In consequence, the absorbed intensity is a function of the photon wavelength (or energy). The most desirable experimental format plots the absorption coefficient α as a function of the radiation frequency ν , because ν is directly proportional to the energy separation between the states involved in the transition. Nevertheless, optical

spectroscopists quote peak positions and linewidths in energy units E (eV, meV), in wave numbers $\bar{\nu}$ (in cm^{-1}), in frequency units (ν or ω), or in wavelength λ [nanometers (nm) or micrometers (μm)]. The following approximate relationships exist: $1 \text{ cm}^{-1} = 1.24 \times 10^{-4} \text{ eV}$; $1 \text{ eV} = 8066 \text{ cm}^{-1}$; and $E(\text{eV}) = 1.24/\lambda (\mu\text{m})$.

Very often the spectrometer output is given in terms of the specimen transmission, $T = I(\nu)/I_o(\nu)$ expressed as a percentage, or the optical density (or absorbance), $\text{OD} = \log_{10}(1/T)$, which are related to the absorption coefficient α by

$$\text{OD} = \log_{10}(1/T) = \alpha(\nu)l/2.303 \quad (1)$$

where l is the thickness of the sample. Typically one measures the absorption coefficient α over the wavelength range 185 to 3000 nm. Since α may be a function of both frequency ν and polarization $\hat{\epsilon}$, we may use the designation $\alpha(\nu, \hat{\epsilon})$. For a center containing N noninteracting absorbing centers per unit volume, each absorbing radiation at a frequency ν and polarization $\hat{\epsilon}$, the attenuation of a beam of intensity $I_o(\nu, \hat{\epsilon})$ by a solid of thickness l is given by

$$I(\nu, \hat{\epsilon}) = I_o(\nu, \hat{\epsilon}) \exp[-\alpha(\nu, \hat{\epsilon})l] \quad (2)$$

Experimentally $I_o(\nu, \hat{\epsilon})$ represents the transmission of the system in the absence of an absorbing specimen. In practice $I_o(\nu, \hat{\epsilon})$ and $I(\nu, \hat{\epsilon})$ are measured and the value of the absorption coefficient $\alpha(\nu, \hat{\epsilon})$ at a particular frequency is obtained using the formula

$$\alpha(\nu, \hat{\epsilon}) = \frac{1}{l} \ln \frac{I_o(\nu, \hat{\epsilon})}{I(\nu, \hat{\epsilon})} \quad (3)$$

$\alpha(\nu, \hat{\epsilon})$ has units of cm^{-1} or m^{-1} . The variation of the absorption coefficient with frequency is difficult to predict. In general, the absorption transition has a finite width, and the absorption strength, $\int \alpha(\nu, \hat{\epsilon}) d\nu$, is related to the density of absorbing centers and to the transition probability.

The value of the absorption coefficient in an isotropic material is related to the Einstein A coefficient for spontaneous emission by¹

$$\alpha(\nu) = \left(N_a \frac{g_b}{g_a} - N_b \right) A_{ba} \frac{c^2}{8\pi\nu^2} \frac{1}{n^2} G(\nu) \quad (4)$$

where g_a and g_b are the statistical weights of the states, $G(\nu)$ is the lineshape function [defined such that $\int G(\nu) d\nu = 1$], c/n is the velocity of light in the medium, and n is the refractive index. In Eq. (4), the population densities in the ground and excited states, N_a and N_b , respectively, have been assumed to be invariant with time and unaffected by the absorption process. Under conditions of weak excitation we can ignore the small value of N_b , and replace N_a by N so that

$$\alpha(\nu) = NA_{ba} \frac{c^2}{8\pi\nu^2} \frac{1}{n^2} \frac{g_b}{g_a} G(\nu) = N\sigma G(\nu) \quad (5)$$

where σ is the *absorption cross section per center*. The absorption strength, i.e., the area under the absorption band, is related to σ by

$$\int \alpha(\nu) d\nu = N\sigma \quad (6)$$

If we ignore the refractive index and local field correction factors and assume a gaussian-shaped absorption band, then

$$Nf_{ab}^c = 0.87 \times 10^{17} \alpha(\nu_o) \Delta\nu \quad (7)$$

where $\alpha(\nu_o)$ is measured in cm^{-1} , $\Delta\nu$, the full-width at half maximum absorption, is measured in eV , and N is the number of centers cm^{-3} . Equation (7) is often referred to as *Smakula's formula*. To obtain the oscillator strength from the area under the absorption band, one needs an independent determination of the density of centers. For impurity ions in solids, N may be determined by chemical assay or by electron spin resonance.

The Double-Beam Spectrophotometer

The first essential of an absorption spectrophotometer is a broadband source: deuterium, hydrogen, xenon, and tungsten lamps are commonly used. Their outputs cover different wavelength ranges: a hydrogen lamp is suitable for wavelengths in the range 150 to 350 nm whereas high-pressure xenon lamps have usable light outputs in the range 270 to 1100 nm. For Xe arc lamps the output is relatively continuous in the wavelength range 270 to 800 nm apart from some sharp lines near 450 nm. In the infrared region 800 to 1100 nm, much of the most intense part of the output is in the form of sharp lines. In the arc lamp, radiation is due to the collision of Xe atoms with electrons which flow across the arc. Complete separation of the excited electrons from the atoms leads to ionization and the continuum output. The formation of Xe atoms in excited states leads to the sharp lines in the output from Xe arc lamps. Tungsten filament lamps may also be used in absorption spectrophotometers. The spectral output from such a heated filament is approximately that of a blackbody radiator at a temperature of 2000 K. In consequence, the emission intensity is a smooth function of wavelength with peak at 1500 nm, with the detailed curve following Planck's thermal radiancy law. Accordingly, the peak in the distribution of light varies with filament temperature (and therefore current through the filament), being determined by $\lambda_{\text{max}} T = 2.898 \times 10^{-3} \text{ mK}$. This relationship expresses Wien's wavelength displacement law. Although containing all wavelengths from the ultraviolet into the infrared region, the total output is fairly modest compared with a high-pressure mercury lamp.

Accurate measurements of the absorption coefficient at different wavelengths are best made using a double-beam spectrophotometer: a schematic is shown in Fig. 1. The exciting beam from the broadband source passes through a grating monochromator: the resulting narrow band radiation is divided by a beam-splitting chopper into separate monochromatic beams which traverse the sample and a reference channel. Thus the light incident on the sample and that which passes through the reference channel have the same wavelengths and is square-wave modulated (on/off) at some frequency in the range 1 to 5 kHz. The sample and reference beams are recombined at the phototube, and the magnitude and phase of sample and reference signals are amplified and compared by the lock-in detector. Chopping at a preselected frequency permits narrowband amplification of the detected signal. Thus any noise components in the signal are limited to a narrowband centered at the chopping frequency. The dc output from the lock-in detector is plotted as a function of wavelength using a pen recorder. Alternatively, the signal may be processed using a microcomputer, so that the absorbed intensity may be signaled as the transmission, the optical density [Eq. (1)], or the absorption coefficient [Eq. (3)] of the sample as a function of wavelength λ , wave number $\bar{\nu}$, or photon energy ($E = h\nu$).

Ensuring a high light throughput in both sample and reference channels usually limits the resolution of the monochromator used in the spectrophotometer (Fig. 1). In consequence, very narrow absorption lines, $\Delta\lambda < 0.1 \text{ nm}$, are normally broadened instrumentally. Note that because in an absorption spectrophotometer one measures the light transmitted by the sample relative to that transmitted by the reference chamber [Eqs. (2) and (3)], the absorption coefficient is independent of the spectral dependencies of the lamp, the monochromator, and the detection system.

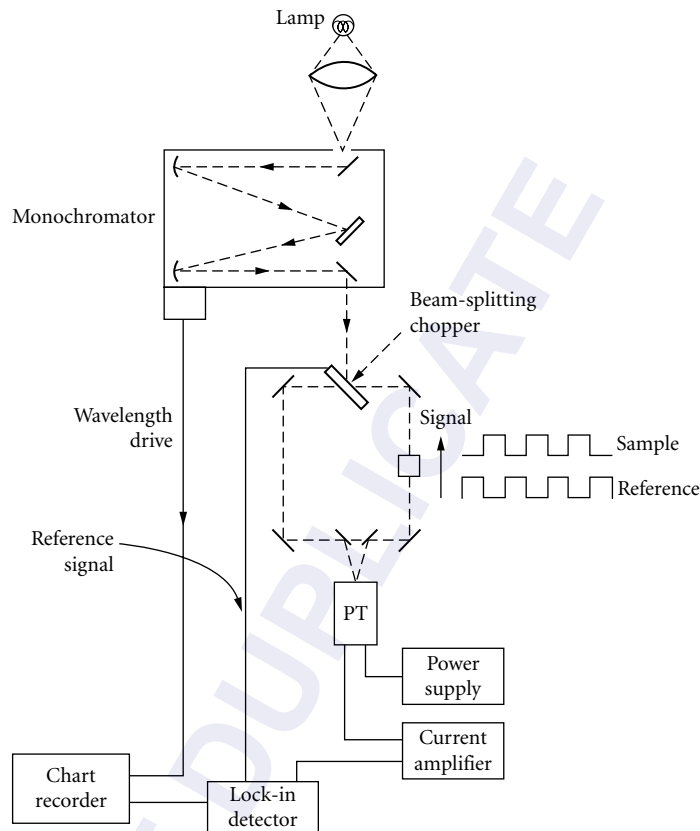


FIGURE 1 Block diagram of a dual-beam absorption spectrometer.

The measurement is also independent of the polarization properties of the monochromator system. By taking ratios, many nonideal behaviors of the components cancel.

31.4 LUMINESCENCE SPECTROMETERS

General Principles

To study luminescence it is necessary to optically pump into the absorption spectrum using high-intensity sources. Typical sources used in luminescence spectroscopy, which have broadbands in near ultraviolet and blue regions, include hydrogen and xenon arc lamps. The xenon arc lamp is particularly useful for exciting luminescence in the yellow-red region of the spectrum since xenon does not show interfering sharp line emission in this region. In general, high-pressure mercury (Hg) arc lamps have higher intensities than Xe arc lamps. However, the intensity is concentrated in sharp lines. Consequently, such lamps are utilized mainly with broadband absorbers or in situations that permit the individual lines to suit the absorption lines of the particular sample. In addition, a variety of lasers may be used, including Ar⁺, Kr⁺, He-Ne, and He-Cd lasers which have emissions at fixed wavelengths. Tunable dye lasers can be selected to closely match the absorption bands of particular

materials. Because of their low intensity, tungsten filament lamps are not normally used in luminescence spectrometers.

The light emitted is resolved into its component lines/bands using a monochromator. For medium resolution a 1-m Czerny-Turner monochromator will give a spectral resolution of about 0.02 nm. An order of magnitude lower resolution can be achieved using a grating spectrometer with focal length 0.25 m. The light emerging from the monochromator is detected using an electron multiplier phototube with associated high-voltage power supplies. Gallium arsenide phototubes operate with good quantum efficiency in the range 280 to 860 nm. For measurements in the near-infrared, a lead sulphide cell, cooled germanium photodetector, or special III-V compound photodiode may be used. Under steady-state optical pumping, a steady-state luminescence output is obtained and detected as a photocurrent which is amplified and converted to a voltage signal to be displayed on a pen recorder. Luminescence detection is inherently more sensitive than absorption measurements and sensitivities of 10^{11} centers cm^{-3} are routine.

Ideally, the excitation source should yield a constant light output at all wavelengths, the monochromator must pass all wavelengths with equal efficiency, and be independent of polarization. In addition, the detector should detect all wavelengths with equal efficiency. Unfortunately, such ideal light sources, monochromators, and phototubes are not available and it is necessary to compromise on the selection of components and to correct for the nonideal response of the luminescence spectrometer. Generally, luminescence spectra are recorded by selecting the excitation wavelength which results in the most intense emission and then scanning the wavelength of the emission monochromator. In consequence, techniques must be developed to allow for the wavelength-dependent efficiency of the emission monochromator and photomultiplier tube. This is not required in absorption spectrophotometers where the ratio of $I(\nu, \hat{\epsilon})/I_0(\nu, \hat{\epsilon})$ are used to compute the values of $\alpha(\nu, \hat{\epsilon})$ from Eq. (3).

Modern spectrometers use diffraction gratings in monochromators rather than prisms. This results in less interference from stray light and in greater dispersion. Stray light may also be reduced using narrow entrance and exit slits as well as double monochromators (i.e., monochromators incorporating two gratings). Nevertheless, the transmission efficiency of the grating monochromator is a strong function of wavelength, which can be maximized at any given wavelength by choice of the blaze angle: the efficiency is less at other wavelengths as Fig. 2 shows. The stray light levels are to some extent controlled by exit and entrance slits. Smaller slit widths also yield higher resolution as do gratings with greater numbers of grooves per unit area. The efficiency of a grating monochromator also depends upon the polarization of the light. For this reason, the observed fluorescence intensities can be dependent upon the polarization of the emitted radiation. A typical plot of the

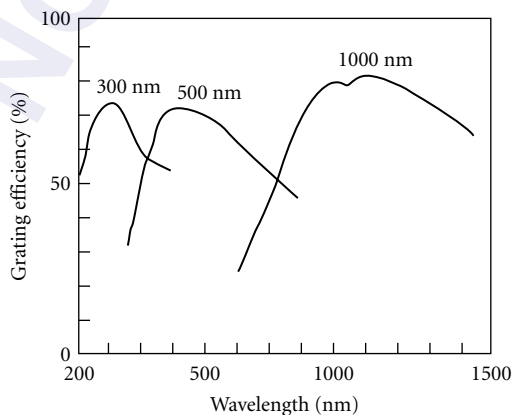


FIGURE 2 Showing how the grating efficiency varies with wavelength for gratings blazed at 300, 500, and 1000 nm.

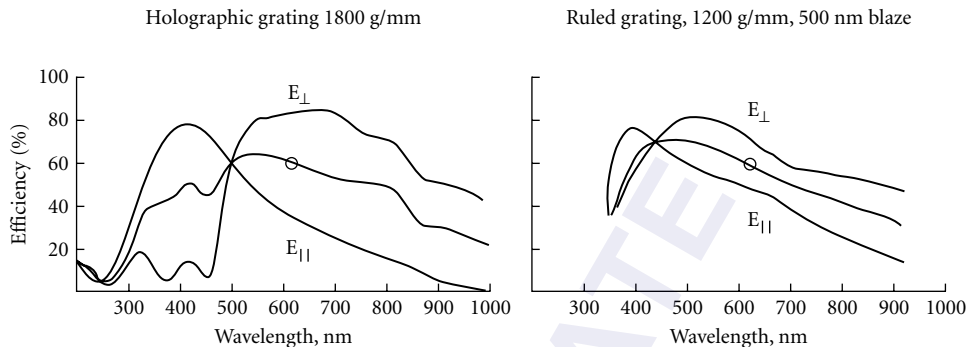


FIGURE 3 Showing the effect of polarization on the efficiency of a ruled grating with 1200 grooves/mm and blazed at 500 nm.

wavelength dependence of the efficiency of a ruled grating as a function of polarization is shown in Fig. 3. As a consequence, the emission spectrum of a sample can be shifted in wavelength and altered in shape by the polarization properties of the monochromator. In modern spectrometers the monochromators can be calibrated using standard lamps and polarizers, the information stored in the memory of the control computer, and the detected intensities corrected at the data processing stage of the experiment. Most manufacturers also provide data sheets describing monochromator performance, and use can be made of such data for approximate corrections to the measured spectra.

Care must be taken with polarization anisotropy measurements. Thin-film polarizers have absorption properties which are strongly wavelength-dependent. Precise corrections can be made using computer-controlled facilities with provision for computerized data processing. However, it is preferable to use a Glan-Thompson prism made from quartz or calcite which has good transparency from the ultraviolet into the infrared. Furthermore, the polarization properties are not wavelength-dependent.

In general terms, the light signal is detected using a photomultiplier tube in which the photon flux produces an electrical current that is proportional to the light intensity. The basis of the device is the photoelectric effect. Incident photons cause photoelectrons to be emitted from a photocathode with an efficiency dependent upon the incident wavelength. The photocathode is held at a high negative potential of 1000 to 2000 V. The photoelectrons are incident upon a series of dynodes which are also held at negative potentials in order to accelerate electrons toward the next dynode. Each photoelectron arriving at the first dynode chain causes the ejection of a further 10 to 20 electrons, depending on the voltage difference between photocathode and first dynode. This process of electron multiplication and consequent current amplification continues down the dynode chain until a current pulse arrives at the anode. Although the photomultiplier tube responds to individual photons, the individual current pulses are generally detected as an average signal.

The anode current must be directly proportional to the light intensity. However, at wavelengths longer than the work function of the photocathode, the photomultiplier tube is no longer sensitive to the incident photons. Thus, different photocathodes are used in different wavelength ranges. For phototubes used in the ultraviolet region, quartz windows are used. For the ultraviolet-visible region (200 to 550 nm) a K-Cs alkali photocathode may be used; such devices have high quantum efficiency, up to 25 percent between 350 to 500 nm, high gain, and low dark current. Typically, the operating anode currents are of the order of a few microamps, whereas the dark current is in the nanoamp range. A somewhat more useful device, in that the quantum efficiency is almost constant from 300 to 860 nm, uses a GaAs photocathode. For longer wavelength operation, 800 to 1250 nm, a germanium photodiode may be used. In other words, spectroscopic studies over a wide wavelength range may require several different photodetectors to be used. Techniques for correcting for the nonideal wavelength-dependent properties of the monochromator, polarizers, and photomultiplier tubes have been described at length by Lackowicz.²

Luminescence Spectrometers Using Phase-Sensitive Detection

Where phase-sensitive detection techniques are used, the excitation intensity is switched on and off at a certain reference frequency so that the luminescence intensity is modulated at this same frequency. The detection system is then set to record signals at the reference frequency only. This effectively eliminates all noise signals except those closely centered on the modulation frequency. A typical luminescence spectrometer is shown in Fig. 4. The pumping light is modulated by a mechanical light chopper operating at frequencies up to 5 kHz. A reference signal is taken from the chopper to one channel of a lock-in detector. The magnitude and phase of the luminescence signal is then compared with the reference signal. Because of the finite radiative lifetime of the emission and phase changes within the electronics, the luminescence signal is not in phase with the reference signal. Hence, to maximize the output from the lock-in detector, the phase control of the reference signal is adjusted until input (luminescence) and reference signals to the lock-in detector are in phase. Of course, the phase of the reference signal may also be adjusted so that reference and luminescence signals are in quadrature giving zero output from the lock-in. This method of phase adjusting may enable one to separate the overlapping luminescence bands from different centers. In such experiments, the chopping frequency is adjusted so that there is an appreciable reduction in the luminescence intensity during the “off” half-cycle. This effectively puts an upper limit on the rate at which the lock-in system can operate.

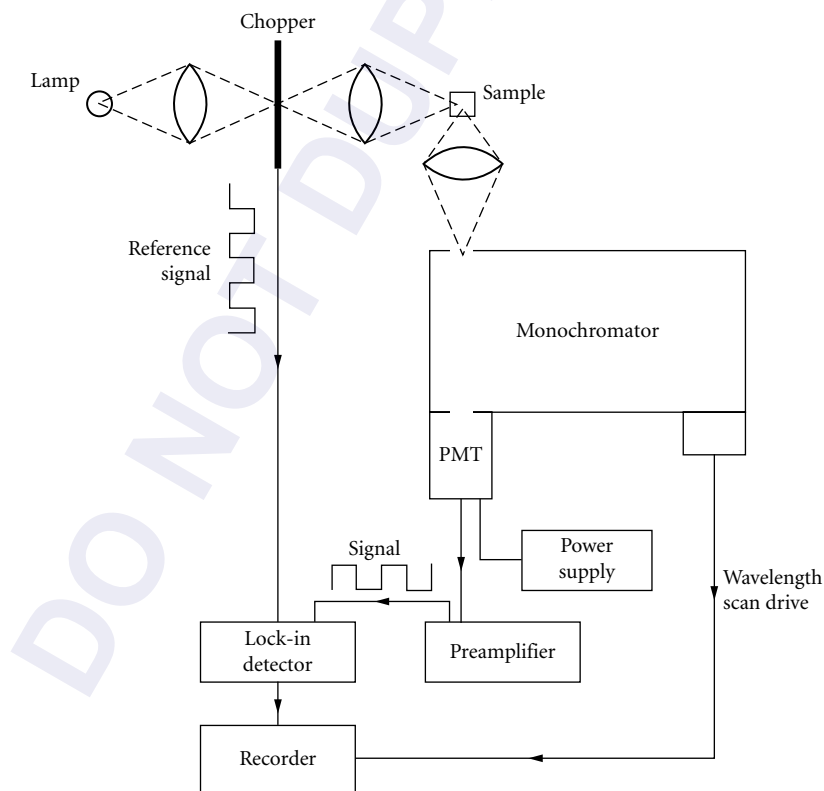


FIGURE 4 Schematic of a spectrometer for measuring luminescence spectra by phase-sensitive detection techniques.

The use of a mechanical chopper restricts the maximum modulation frequency to 5 kHz. Essentially, the mechanical chopper consists of a rotating blade of metal into which slots are cut at regular angular intervals. When the excitation beam is incident on the metal section, the excitation intensity at the sample is zero, and when on the slot, the sample receives the full intensity in the excitation beam. If the blade is rotated at a frequency of 1 Hz and there are n slots cut in the blade, then the excitation beam is essentially switched on/off at a rate of n Hz. Obviously, the modulation rate can be increased either by increasing the number of slots cut in the blade and/or by increasing the revolution rate. If the excitation is well-focused onto the chopper, then the modulation is in the form of a square wave with maximum modulation index $M = 1$.

Other modulators impose sinusoidal variations in the excitation intensity, at frequencies up to 50 MHz.² There are various means for providing sinusoidal intensity variations, including Kerr cells and Pockels cells. Both types require high electric fields to obtain the desired modulation and such high driver voltages may interfere with the detection of weak signals. Kerr cells do not transmit ultraviolet light and so may only be used in the visible/near-infrared region. The Pockels cells may be used in the ultraviolet region as well as at visible and infrared wavelengths. They may also be operated at variable frequencies. However, since they require highly collimated light sources for efficient operation, they require a laser for excitation. The ultrasonic Debye-Sears modulator overcomes the experimental difficulties associated with both Pockel cells and Kerr cells. A vibrating quartz crystal is used to set up standing waves in a tank containing an ethanol-water mixture. (The crystal restricts the device to operate at the fundamental and one or two harmonic frequencies only.) The standing waves act as a closely spaced refractive index diffraction grating normal to the incident exciting radiation. A slit permits only the undiffracted light to pass to the sample. The result is a sinusoidally varying light intensity with about 50 percent modulation index.

The emission signal is forced to respond to the modulated excitation at the same circular frequency ω as the excitation. However, the detected emission signal is delayed in phase by an angle ϕ relative to the excitation, and with reduced modulation depth. The radiative lifetime may be calculated from the measured phase angle ϕ and demodulation factor m . For a single exponential decay the appropriate relations are³

$$\tan\phi = \omega\tau_R \quad (8)$$

and

$$m = [1 + \omega^2\tau_R^2]^{-1/2} \quad (9)$$

Even with more complex processes, where several decaying species are present, phase angles and demodulation factors can be measured and used to calculate actual lifetimes.^{3,4}

Phase-sensitive detection techniques may also be used to "time-resolve" overlapping absorption/luminescence spectra with different decay characteristics. The phase-sensitive detector (PSD) yields a direct-current signal proportional to the modulated amplitude and to the cosine of the phase difference between the detector phase ϕ_D and the signal phase ϕ , i.e.,

$$I(\lambda, \phi_D) = m_s I_o(\lambda) \cos(\phi_D - \phi) \quad (10)$$

where λ is the wavelength, $I_o(\lambda)$ is the steady-state excitation intensity, and m_s is the source modulation index. Now suppose that there are two components A and B with lifetimes $\tau_A < \tau_B$. The modulated emission measured with the PSD results in an unmodulated signal given by

$$I(\lambda, \phi_D) = m_s^A I_o^A(\lambda) \cos(\phi_D - \phi_A) + m_s^B I_o^B(\lambda) \cos(\phi_D - \phi_B) \quad (11)$$

If the phase-control of the PSD is adjusted so that $\phi_D = \phi_B + 90^\circ$, then the second term in Eq. (11) is zero, and the output intensity is given by

$$I(\lambda, \phi_D) = m_s^A I_o^A(\lambda) \sin(\phi_B - \phi_A) \quad (12)$$

In other words, the emission output from species *B* has been suppressed. Species *A* can be suppressed at the detector phase angle $\phi_D = \phi_A + 90^\circ$. If we now scan the wavelength, then the consequence of Eq. (12) is that the steady-state spectrum of species *A* is recorded, i.e., $I_o^A(\lambda)$, and conversely for species *B*.

In the example given in Fig. 5*a*, the steady-state fluorescence of a mixture of indole and dimethylindole dissolved in dodecane is shown.⁵ With the detector phase angle set to $90^\circ + 9.7$ and using a modulation frequency of 10 MHz in Fig. 5*b*, we resolve the indole emission with wavelength maximum at 306 nm. The phase angle of 9.7° corresponds to a radiative lifetime close to the isolated

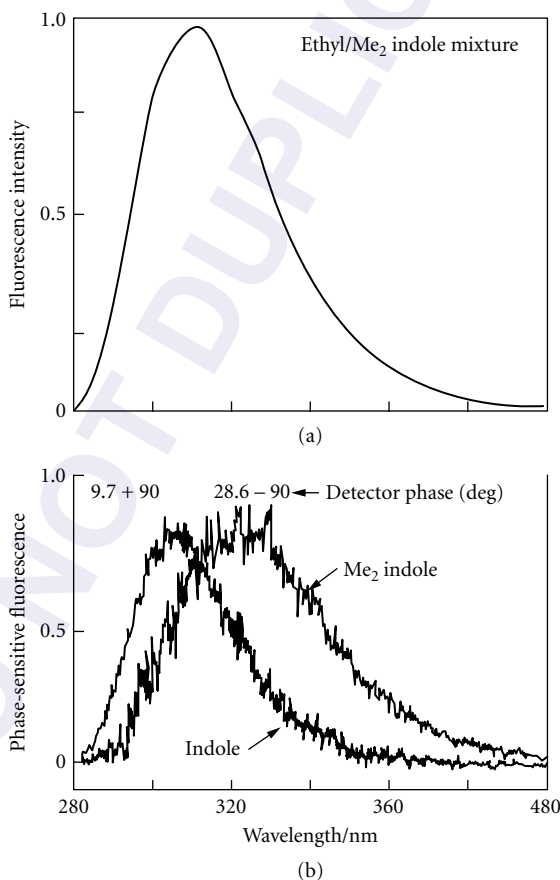


FIGURE 5 (a) Steady-state emission spectra of a mixture of indole and dimethylindole in dodecane and (b) shows the phase-resolved spectra of the indole and dimethylindole.⁵

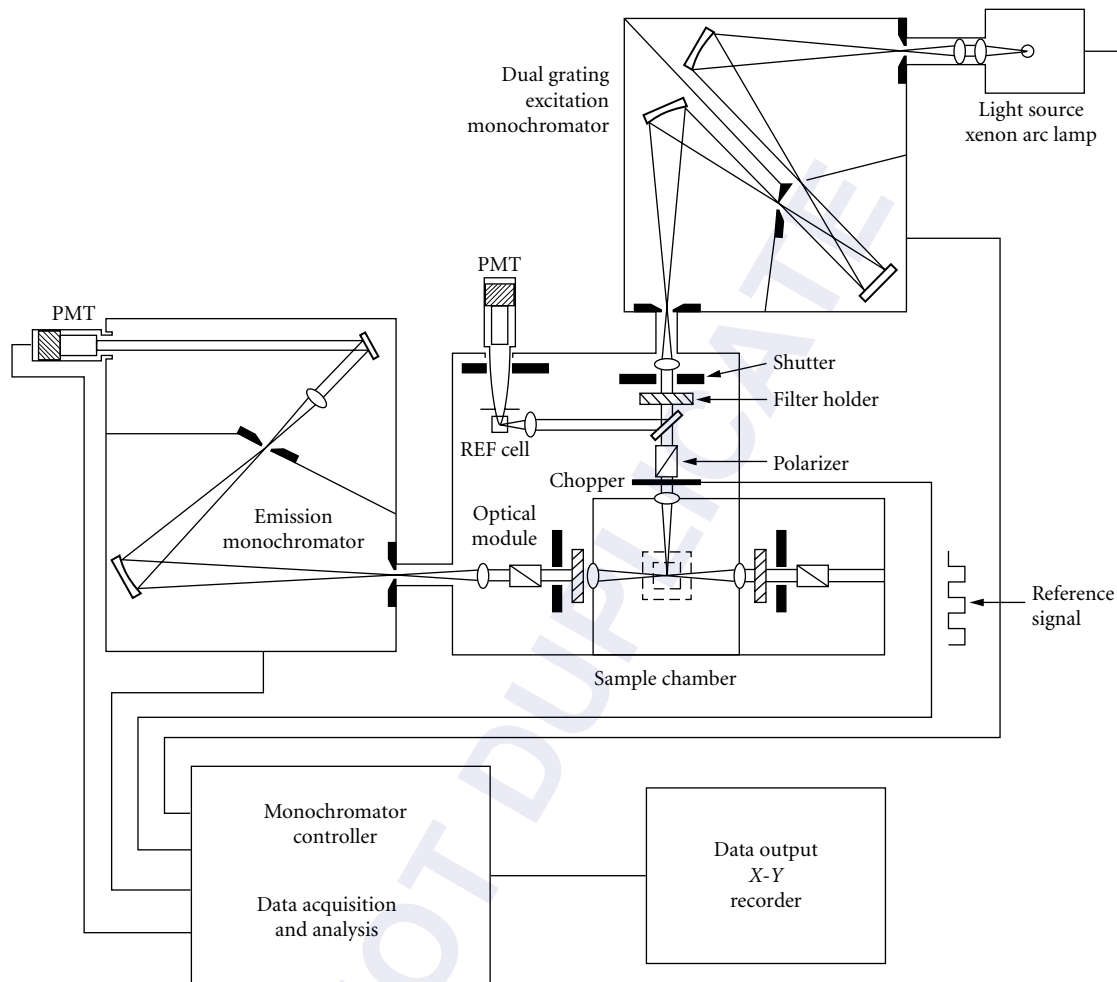


FIGURE 6 Schematic representation of a luminescence excitation spectrometer.

methylindole molecules in dodecane ($\tau_R = 5$ ns). The suppression of the indole signal gives the dimethylindole spectrum with peak at 323 nm at a phase angle of $28.6 - 90^\circ$, giving the τ_R value of indole as 9.0 ns.

Luminescence Excitation Spectrometers

Some inorganic solids have strong overlapping absorption bands due to nonluminescent centers, which completely overwhelm the absorption spectrum related to a particular luminescence center. These difficulties are overcome by excitation spectroscopy, Fig. 6, in which the intensity of the luminescence output is recorded as a function of the wavelength of the excitation beam. Strong emission at a particular excitation wavelength signals that the emitting center absorbs strongly at that wavelength. In this way it is possible to determine the shape and position of the

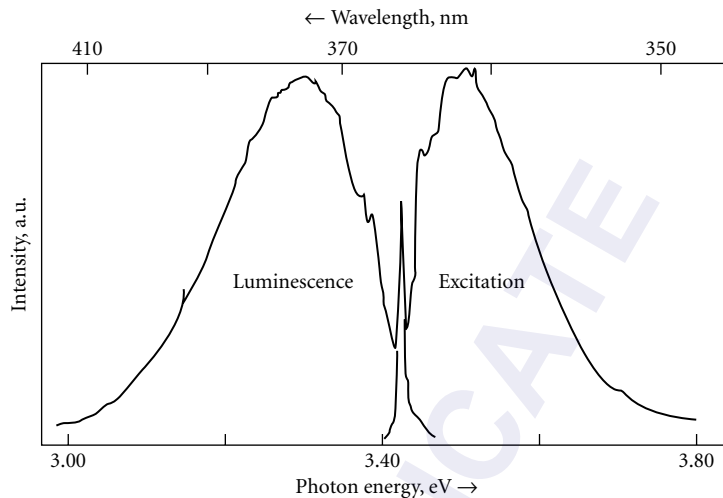


FIGURE 7 Luminescence and excitation luminescence spectrum of F_2 centers in MgO .

absorption bands which excite the emission process. A low-resolution scanning monochromator is placed immediately after the chopper, and light from its exit slit is then focused onto the sample. This monochromator may be of focal length only 250 mm and have a grating of area $5\text{ cm} \times 5\text{ cm}$, ruled with only 600 lines per mm. Alternatively, it may be a double monochromator chosen to reduce stray light levels. In either case, an optical band pass filter may be used in front of the monochromator. Generally, the grating blaze is chosen to give high efficiency in the ultraviolet/blue/green regions for the excitation monochromator (e.g., gratings blazed at 300 nm or 500 nm), whereas the emission monochromator is chosen to give high efficiency at visible and near-infrared wavelength (i.e., gratings blazed at 500 nm or 750 nm). With such an apparatus, it is possible to distinguish absorption transitions from several centers whose absorption bands partially or completely overlap. The example given in Fig. 7 shows the luminescence pattern emitted by F_2 centers in magnesium oxide and the excitation spectrum associated with this emission. Other strong absorption bands due to Fe^{3+} ions and F centers which overlap the F_2 -absorption bands are strongly discriminated against by selective detection of the F_2 -center luminescence.

31.5 PHOTOLUMINESCENCE DECAY TIME

Radiative Lifetime

In order to measure the radiative lifetime of a transition it is necessary to use a sharp intense pulse of excitation in the absorption band together with some means of recording the temporal evolution of the luminescence signal. Suitable excitation sources include pulsed lasers, flash lamps, or stroboscopes. Laser systems may produce pulses of duration 0.1 to 100 ps; flash lamps and stroboscopes will produce pulses of order 10^{-8} s and 10^{-5} s, respectively. A possible spectrometer system is shown in Fig. 8. Usually the luminescence yield following a single excitation pulse is too small for good signal-to-noise throughout the decay period. In consequence, repetitive pulsing techniques are used together with signal averaging to obtain good decay statistics. The pulse reproducibility of the stroboscope is advantageous in the signal averaging process in which the output from the detector is sampled at equally spaced time intervals after each excitation pulse. If the pulse is repeated N times

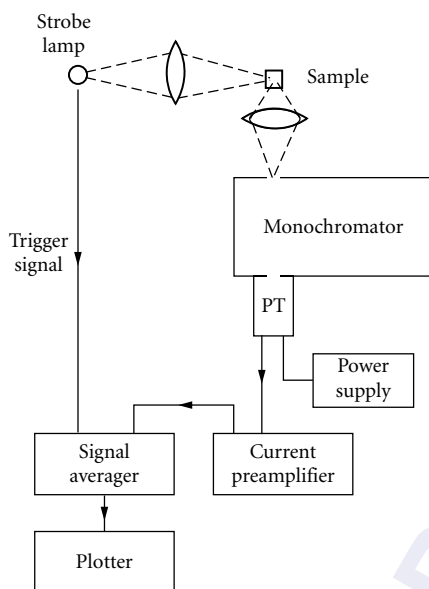


FIGURE 8 Spectrometer for measuring luminescence decay times.

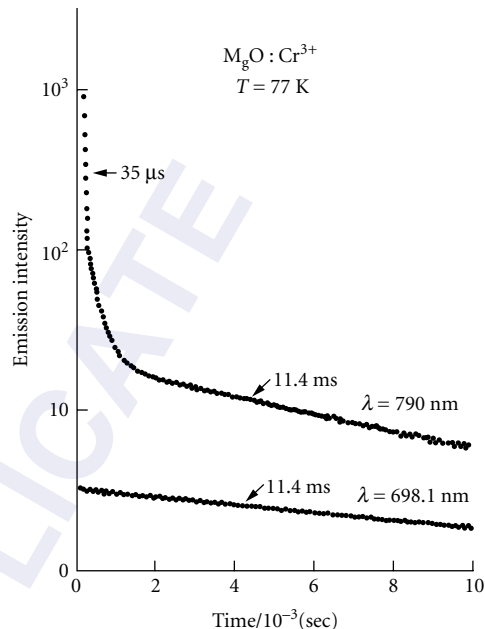


FIGURE 9 Decay in the intensity of Cr^{3+} luminescence in MgO . Detection of the broadband luminescence at 790 nm shows two components, one fast ($\tau_R = 35 \mu\text{s}$) and one slow ($\tau_R = 11.4 \text{ ms}$). Detection of the R-line at 698.1 nm shows a single component with $\tau_R = 11.4 \text{ ms}$.

then there is an $N^{1/2}$ improvement in the signal-to-noise ratio. If a multichannel analyzer is used, the excitation pulse is used to trigger the analyzer, and hence the time between pulses need not be constant. Of course the phase sensitive detection spectrometer may also be used to measure lifetimes, but only down to about 100 μs .

An illustration of the data obtainable using the stroboscope technique is shown in Fig. 9. The luminescence signal detected is the broadband emission with peak at 780 nm from Cr^{3+} ions in orthorhombic symmetry sites in magnesium oxide measured at 77 K. At low Cr^{3+} ion concentration, the radiative lifetime of this luminescence center is 35 μs . These data show that the evolution of the intensity during the pulse-decay cycle is not necessarily in the form of a single exponential decay. On sampling the emission at times long relative to τ_R there is a component with characteristic decay time of 11.4 ms, which is the lifetime of Cr^{3+} ions occupying octahedral symmetry sites in magnesium oxide and which emit a characteristic R-line emission at 698.1 nm. This result implies that excitation is being transferred from excited Cr^{3+} ions in octahedral sites to Cr^{3+} in orthorhombic sites.

For rather faster decay processes (10^{-10} – 10^{-8} s), fast flashlamps are used to excite the luminescence. The gated flashlamps have extremely reproducible pulses, down to 0.8-ns width with repetition rates of up to 50 kHz.² The usual gases for such lamps are hydrogen, deuterium, and air. Hydrogen has several advantages, not the least being the continuum output in the ultraviolet and visible ranges, with pulse profiles which are independent of wavelength. The combination of pulse-sampling techniques and computer deconvolution of the decaying luminescence enables decay times to be measured down to 20 ps. However, judicious choice of photomultiplier tube and careful design of the photomultiplier dynode chain is necessary to eliminate signal noise. It is usual to use coincidence single-photon counting techniques to obtain good decay data.²

Picosecond and Sub-Picosecond Relaxation

During the past two decades there have been quite remarkable developments in techniques for generating and measuring ultrashort pulses into the femtosecond domain. In semiconductors, a very wide range of ultrafast phenomena are being studied—electronhole plasma formation, exciton and biexciton formation dynamics, hot electron effects, phase-conjugate self-defocusing, and degenerate four-wave mixing. However, one very general optical phenomenon that may be addressed using ultrashort pulses involves nonradiative decay times in nonresonant fluorescence spectra. Such processes include ionic relaxations around a center in decaying from an excited state, sometimes including reorientations of anisotropic centers. Many picosecond phenomena, especially nonradiative decay processes, are studied by excite-and-probe techniques in which light pulses at wavelength λ_1 are used to excite a phenomenon of interest, and then a delayed optical pulse at wavelength λ_2 interrogates a change of some optical property of this phenomenon. Ideally, two sources of picosecond pulses at different, independently tunable wavelengths are required, which must be synchronized on the picosecond timescale.

A convenient experimental system for studying vibrational relaxation at color centers and transition metal ions in ionic crystals is shown in Fig. 10.⁶ A mode-locked dye laser producing sub-picosecond pulses at wavelength λ_1 is used both to pump in the absorption band and to provide the timing beam. Such pumping leads to optical gain in the luminescence band and prepares the centers in their relaxed state. The CW probe beam, collinear with the pump beam, operates at a longer wavelength, λ_2 . The probe beam and gated pulses from the pump laser are mixed in a nonlinear optical crystal and a filter allows only the sum frequency of the pump and probe beams, which is detected by a phototube. The photomultiplier tube actually measures the rise in intensity of the probe beam which signals the appearance of gain when the $F_A(\text{Li})$ -centers have reached the relaxed excited state. The pump beam is chopped at low frequency to permit phase-sensitive detection. The temporal evolution gain signal is measured by varying the time delay between pump and gating pulses. Although specifically used by Mollenauer et al.⁶ to probe the relaxation dynamics of color centers, the spectrometer system shown in Fig. 10 is readily adapted to other four-level systems, including transition metal ions.

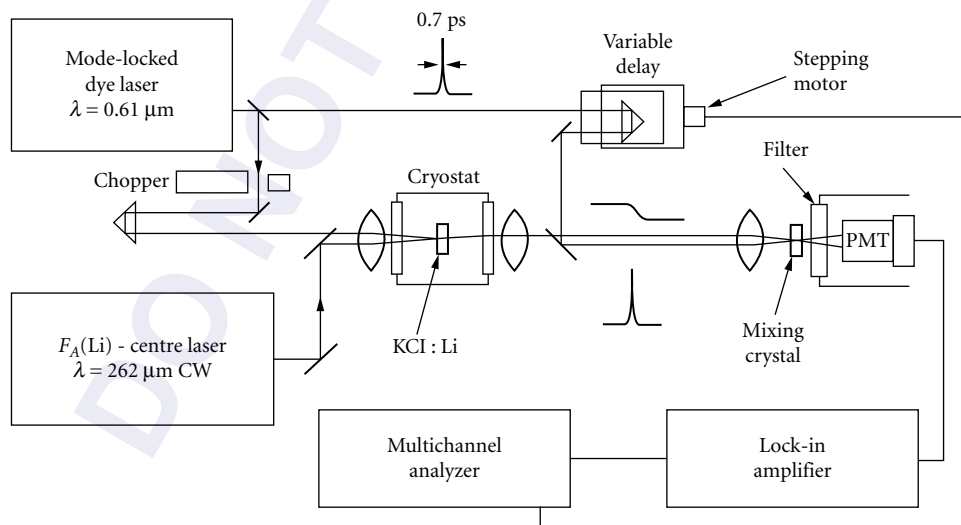


FIGURE 10 A sub-picosecond pump and probe spectrometer for measuring vibrational relaxation times in excited defects and transition metal ions.

31.6 POLARIZATION SPECTROMETERS

General Principles

The absorbed intensity is sometimes dependent on the optical anisotropy of the sample. Whether or not a transition is allowed in a particular polarization is determined by examining the value of the square of the matrix element $\langle b | \hat{\mu} \cdot \hat{\epsilon}_j | a \rangle$, see "Optical Spectroscopy and Spectroscopic Lineshapes" (Chap. 10 in this volume), Eqs. (12) and (13), where a and b are the states involved in the transition, $\hat{\mu} \cdot \hat{\epsilon}_j$ is the appropriate component of the dipole operator summed over all j electrons involved in the transition. Optical transitions may be linearly or circularly polarized. For an electronic dipole transition, the dipole operator is $\hat{\mu}_e \cdot \hat{\epsilon}_E$ where $\hat{\mu}_e = \sum_j e \mathbf{r}_j$ is summed over the j electrons and $\hat{\epsilon}_E$ is the unit electric polarization vector parallel to the \hat{E} -field of the radiation. The matrix element is evaluated using group theory, which shows how the symmetry properties of the states affect the transition rate.¹ From this matrix element the selection rules of the transition are determined. The polarization of the radiation is defined in Fig. 11 by reference to the \hat{z} direction of the system, which itself is assumed to be parallel to an external perturbation (static electric or magnetic fields) or to unique symmetry direction in a crystal. For the π - and σ -senses of linear polarization, the radiation travels in a direction perpendicular to \hat{z} with its electric field $\hat{\epsilon}_E$ either parallel to \hat{z} (π -polarization) or perpendicular to \hat{z} (σ -polarization). The electric dipole operators are then given by $\sum_j e \mathbf{r}_j \cdot \hat{z} = \sum_j e z_j$ for π -polarization and $\sum_j e x_j$ or $\sum_j e y_j$ for σ -polarization. The \hat{x} and \hat{y} directions have been assumed equivalent. In α -polarization the radiation propagates along the unique symmetry axis, \hat{z} , with $\hat{\epsilon}_E$ anywhere in the x - y plane: in this case the electric dipole operator is also $\sum_j e x_j$ or $\sum_j e y_j$. We define right circularly polarized (RCP) radiation as having electric (and magnetic) polarization vectors which rotate clockwise when viewed from behind the direction of propagation. For electric dipole absorption transitions, the electric dipole operator for RCP light propagating in the z direction is $\sum_j e(x + jy)_j / \sqrt{2}$. Accordingly, in the case of LCP light, where the sense of rotation is anticlockwise, the electric dipole operator is $\sum_j e(x - jy)_j / \sqrt{2}$.

Polarized Absorption

Although the selection rules of dipole transitions provide for polarized spectra, the optical spectra of atoms in the absence of external fields and of electronic centers with octahedral symmetry in solids are isotropic. Since the unit polarization vector, $\hat{\epsilon}_E$, has direction cosines $\cos \alpha$, $\cos \beta$, and

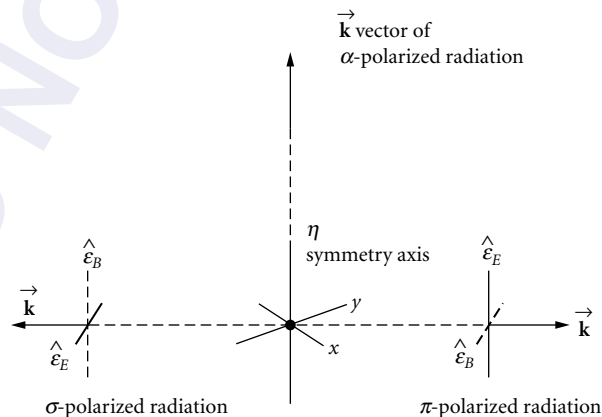


FIGURE 11 Definitions of the senses of α -, π - and σ -polarized light beams relative to a local symmetry axis.¹

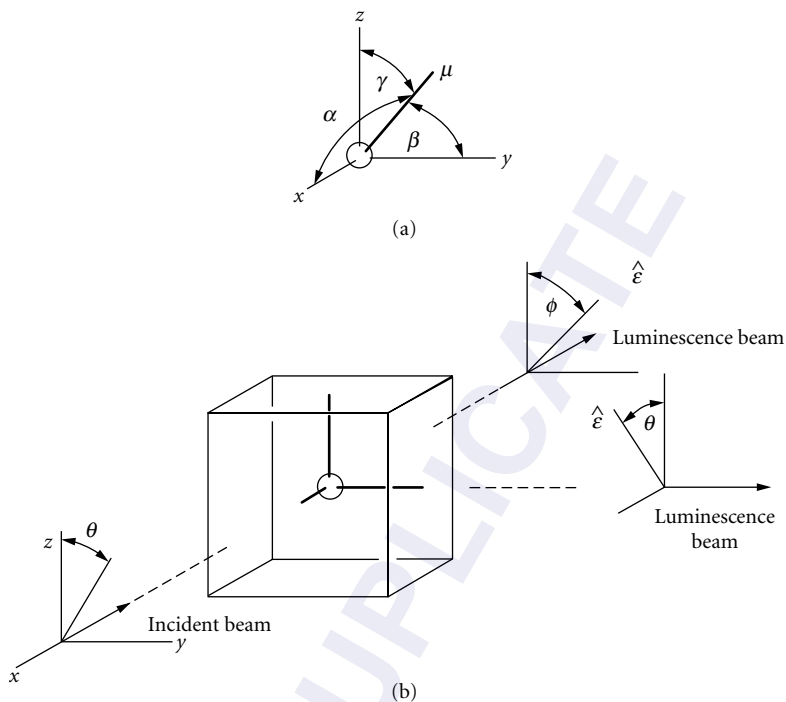


FIGURE 12 Showing (a) the orientation of the dipole moment μ relative to the x , y , and z axes and (b) two different geometrical arrangements used for polarized excitation luminescence spectroscopy.

cos γ , where the angles α , β , and γ are defined in Fig. 12a, the square of the electric dipole matrix element is

$$| \langle b | \mu_x \cos \alpha + \mu_y \cos \beta + \mu_z \cos \gamma | a \rangle |^2 \quad (13)$$

When the center has octahedral symmetry the cross terms in Eq. (13) are zero so that the squared matrix element becomes

$$\langle \mu_x \rangle^2 \cos^2 \alpha + \langle \mu_y \rangle^2 \cos^2 \beta + \langle \mu_z \rangle^2 \cos^2 \gamma \quad (14)$$

using $\langle \mu_x \rangle = \langle b | \mu_x | a \rangle$ with similar expressions for $\langle \mu_y \rangle$ and $\langle \mu_z \rangle$. Since in octahedral symmetry

$$\langle \mu_x \rangle^2 = \langle \mu_y \rangle^2 = \langle \mu_z \rangle^2 \quad (15)$$

$|\langle \mu \cdot \hat{\epsilon} \rangle|^2$ becomes $\langle \mu_x \rangle^2$ and the strength of the transition is independent of the direction of the polarization of the incident radiation and the direction of propagation.

In octahedral solids, the local symmetry of an electronic center may be reduced by the application of an external perturbation or internally through the presence of a nearby defect. In tetragonal symmetry with the z axis parallel to the symmetry axis, the transition probability is again given by

Eq. (14) but with $\langle \mu_x \rangle^2 = \langle \mu_y \rangle^2 \neq \langle \mu_z \rangle^2$. Since the transition probability for radiation at polarization $\hat{\mathbf{e}}_{\beta}$ is then proportional to

$$\langle \mu_x \rangle^2 (\cos^2 \alpha + \cos^2 \beta) + \langle \mu_z \rangle^2 \cos^2 \gamma = A + B \cos^2 \gamma \quad (16)$$

where A and B are constants, the spectroscopic properties of the center are anisotropic. In terms of the experimental situation referred to in Fig. 12b, in α -polarization the angle γ is always $\pi/4$ radians, and the intensity is proportional to A . For π -polarization, the angles are $\alpha = \pi/4$, $\beta = \pi/4 - \alpha$, and $\gamma = 0$, and the intensity is proportional to B . Similarly, for σ -polarization, the intensity is proportional to A . This shows that in tetragonal symmetry a rotation of the polarizer from the $\gamma = 0$ to $\gamma = \pi/4$ in, for example, the y - z plane determines the magnitudes of A and B . The linear dichroism D is then given by $D = (B - A)/(A + B)$.

To illustrate these ideas, Fig. 13 shows the polarization of the ${}^2S_{1/2} \rightarrow {}^2P_{1/2}, {}^2P_{3/2}$ lines of atomic sodium, i.e., the D_1 and D_2 absorption lines, in the presence of an applied magnetic field. The Zeeman splittings of energy levels are much smaller than the spin-orbit splitting between the ${}^2P_{1/2}$ and ${}^2P_{3/2}$ levels. The wave functions are labeled in Fig. 11 by the M_j -values: the relevant Clebsch-Gordan coefficients and theoretical intensities of the transitions for linear and circular polarizations are shown in Fig. 13, as are the theoretical intensities of the ${}^2S_{1/2} \rightarrow {}^2P_{1/2}, {}^2P_{3/2}$ right circularly polarized (RCP) and left circularly polarized (LCP) absorption transitions. The experimental pattern of lines for π - and σ -polarizations are in excellent agreement with the predicted Zeeman pattern.

An analysis of the polarization properties of the sample absorption requires that a polarizer be inserted in the light path immediately prior to the sample chamber. For accurate measurements of the absorption anisotropy, the polarizers must be accurately positioned relative to the beam and rotatable. The angle of rotation about the beam must be accurately indexed so that the orientation-dependence of the anisotropy may be determined. The polarizer should also be removable since it is unnecessary for measurements with optically isotropic solids. A sample which has different absorption coefficients in different crystallographic directions is said to be *dichroic*. The dichroism is defined as

$$D = \frac{\alpha(\pi) - \alpha(\sigma)}{\alpha(\pi) + \alpha(\sigma)} = \frac{1}{I} \left(\frac{I(\pi) - I(\sigma)}{I(\pi) + I(\sigma)} \right) \quad (17a)$$

in the limit of small absorption coefficients.

Although discussion has focused on radiative absorption transitions via electric dipole transitions, a similar analysis can be made for magnetic dipole transitions. In this case, the phase relationships between the magnetic fields \mathbf{B}_x and \mathbf{B}_y are exactly the same as those between \mathbf{E}_x and \mathbf{E}_y , and the magnetic dipole operator is $\boldsymbol{\mu}_m \cdot \boldsymbol{\epsilon}_B$; where $\boldsymbol{\mu}_m = \sum_j (e/2m) \times (l + 2s)_j$ and $\boldsymbol{\epsilon}_B$ is the unit vector along the direction of the magnetic field of the radiation. If the absorption transitions used to excite the luminescence are unpolarized, so too will be the resulting luminescence spectrum. However, as discussed above, the absorption spectrum of an atomic system may be made anisotropic by the application of an external field or by using polarized exciting radiation. The resulting emission spectrum will be correspondingly polarized. Absorption and luminescence spectra from optically isotropic solids can also be made anisotropic using similar techniques.

Polarized Absorption/Luminescence

Just as the absorption spectra of free atoms and isotropic solids are unpolarized, so too are the luminescence spectra, at least when exciting with unpolarized radiation. This is shown by simple extensions to the arguments leading to Eq. (15) in which the electric dipole operators for luminescence are the complex conjugates of those for the appropriate absorption transitions. In practice, both absorption and emission properties are anisotropic. Although the host crystal may possess a cubic unit cell in which the electronic centers are anisotropic, a regular distribution of equivalent sites will

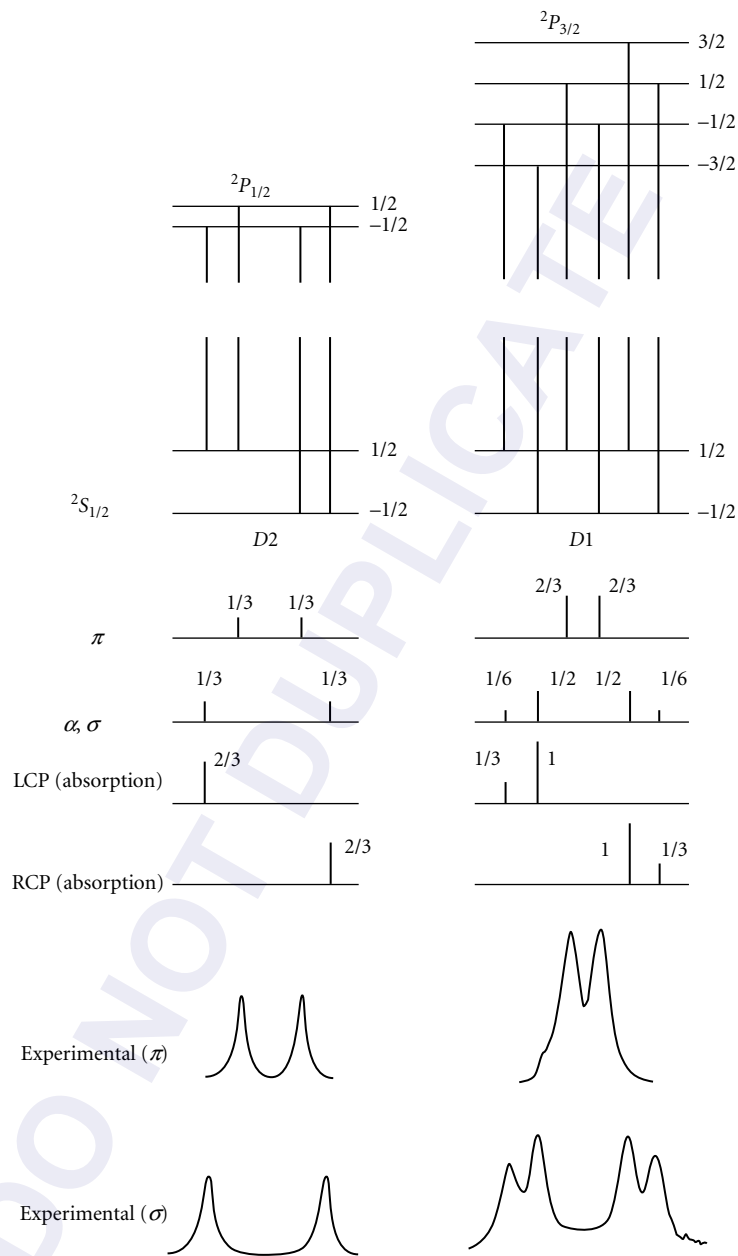


FIGURE 13 Zeeman splittings of $^2S_{1/2} \rightarrow ^2P_{1/2}, ^2P_{3/2}$ levels of sodium. The electric dipole matrix elements and relative intensities of linearly and circularly polarized absorption transitions are compared with some experimental spectra.¹

still result in isotropic spectra. The use of polarized absorption/luminescence techniques can reveal details of the local site anisotropy. Such methods have been discussed by Henderson and Imbusch¹ and in more detail by Feofilov.⁷

To measure the effects of polarization on the absorption coefficient [i.e., $\alpha(\nu, \hat{\epsilon})$] it is necessary to place a polarizer immediately before the beam splitter in the double-beam spectrophotometer, Fig. 1. In polarized luminescence measurements, linear polarizers are placed immediately before the sample in the absorption channel and just after the sample in the emission channel of a luminescence excitation spectrometer such as that shown in Fig. 6. The spectrometer may then operate in the "straight through" configuration or the emitted light may be collected in a direction at 90° to the direction of the excitation light, as illustrated in Fig. 13. Note that provision is made for rotatable channels in both excitation (θ) and detection channels (ϕ), and the measured emission signal will be a function of both θ and ϕ .

The circular dichroism may be defined in an analogous manner to the linear dichroism, i.e., Eq. (17a). Since circular dichroism has a specific relevance to the Zeeman effect, we use Fig. 14a and consider circularly polarized absorption transitions which are excited between two Kramers doublets. With light propagating along the direction of the magnetic field, the selection rule is that

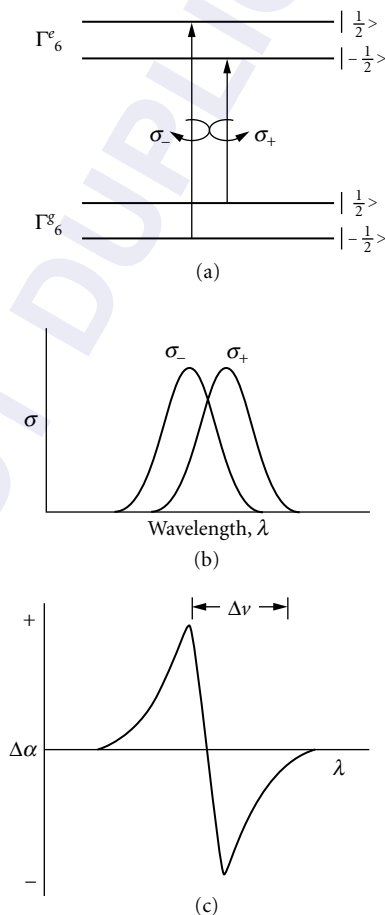


FIGURE 14 Circularly polarized dipole transitions excited between Kramers states.

σ -polarized light induces $\Delta M_s = +1$ absorption transitions and σ_+ -polarized light induces absorption transitions in which $\Delta M_s = -1$. As a result of the Zeeman effect, the absorption peak splits into two overlapping bands (Fig. 14b) centered at different wavelengths with absorption coefficients α_{\pm} in σ_+ - and σ_- -polarizations that are different at particular frequencies. The peaks in the two oppositely polarized bands are separated in energy by $(g_e + g_g)\mu_B B$, where the g -values refer to the excited e and ground g states. The energy difference corresponds to a frequency splitting $\Delta\nu = (g_e + g_g)eB/4\pi m$, which for $g_e = g_g = 2.0$ and $B = 1 T$ gives a separation between band peaks of ≈ 0.04 nm for a band centered at 500 nm. In a magnetic field, the difference $[\Delta\alpha(\nu)]$ in the absorption coefficients for σ_+ and σ_- circularly polarized light is referred to as *magnetic circular dichroism* (MCD). In the limit of small absorption coefficient the circular dichroism is

$$\Delta\alpha(\nu) \approx -\frac{2(I_+(\nu) - I_-(\nu))}{l(I_+(\nu) + I_-(\nu))} \quad (17b)$$

where l is the sample thickness and $I_+(\nu)$, $I_-(\nu)$ refer to the transmitted intensities of the σ_+ and σ_- circularly polarized light at frequency ν .

In most cases, a splitting of only 0.04 nm would be hard to resolve directly by Zeeman effect measurements on a broadband. However, this Zeeman structure may be resolved by measuring $\Delta\alpha(\nu)$ as a function of magnetic field, as can be seen from a simple estimate. We approximate the MCD signal, $\Delta\alpha(\nu)$ for a sample of thickness l , as the product of the magnetic splitting $\Delta\nu$ with the rate of change of the absorption coefficient with frequency which is given by $d\alpha(\nu)/d\nu \approx \alpha(\nu_o)/\Gamma$, for a symmetrical, structureless band. Hence

$$\Delta\nu = \frac{\Delta\alpha(\nu)}{l} \times \frac{\Gamma}{\alpha(\nu_o)} \quad (18)$$

In a typical experiment, $\Delta\alpha(\nu)l \approx 10^{-5}$ and $\alpha(\nu_o)l \sim 1$, hence $\Delta\nu \approx 10^{-5}\Gamma$. For a typical broadband, $\Gamma \approx 0.25$ eV ≈ 2000 cm $^{-1}$ and Eq. (18) yields $\Delta\nu \approx 0.02$ cm $^{-1}$ (i.e., $\Delta\lambda \sim 0.05$ nm) which is of the same order of magnitude as the Zeeman splitting calculated above. Although the intensity changes, which determine the magnitude of $\Delta\alpha(\nu)$, may be quite small, they may be assumed very precisely using lock-in techniques. This is done very efficiently by replacing the circular polarizer in the excitation system by a stress-modulated quarter-wave plate, a device which transmits circularly polarized light, the polarization of which is switched between σ_+ and σ_- at the vibration frequency of the plate, usually ≈ 50 kHz. Using this piezo-optic modulation, MCD signals as low as 10^{-6} can be measured.⁸

The MCD signal is strongly dependent on both frequency and temperature. Since at low temperatures the populations N_{\pm} of the $M_s = \pm 1/2$ levels of the spin 1/2 ground state are different for a system in thermal equilibrium, the MCD signal [Eq. (17b)] is given by

$$\Delta\alpha(\nu)l = \alpha_o(\nu)G(\Delta\nu) \tan h \left[\frac{g\mu_B B}{2kT} \right] \quad (19)$$

In this expression $\alpha_o(\nu)$ and the sample thickness, l , are experimental constants and, in consequence, the MCD signal only varies through the Brillouin function for the $s=1/2$ ground state [i.e., $\tan h(g\mu_B B/2kT)$]. This MCD signal is paramagnetic, being strongest at high field and low temperature, and measurement of its magnitude probes the ground-state magnetization. In order to test Eq. (19) experimentally, it is best to work at either the positive or negative peak in Fig. 14 and so maximize the MCD signal. Having thus obtained a suitable MCD signal, its variation with temperature and magnetic field can then be measured. Excitation of the Kramers' system in Fig. 14 with circularly polarized radiation of appropriate frequency results in the circularly polarized emission. The electric dipole operators for RCP and LCP emission are the complex conjugates of those

for absorption. The consequent change in emission intensity, Eq. (17), is referred to as the magnetic circular polarization (MCP).

Optically Detected Magnetic Resonance

In optical absorption spectroscopy, electronic transitions (usually) out of the ground state may result in one of a rich tapestry of possible bandshapes, depending upon the strength of the electron-phonon coupling. Photoluminescence measurements involve transitions which originate on an excited electronic state and frequently the terminal state is the electronic ground state. Overlapping absorption and luminescence bands can cause difficulty in assigning individual optical absorption and luminescence bands to particular optical centers. Since the lifetimes of excited states are in the range 10^{-3} to 10^{-8} s, it is no trivial matter to measure excited-state electron spin resonance using the microwave detection techniques pioneered in ground-state studies. Geschwind et al.⁹ developed techniques in which the excited-state ESR was detected optically. In favorable cases this method enables one to correlate in a single experiment ESR spectra in the ground state and in the excited state with particular optical absorption and luminescence spectra. The technique of measuring the effect of resonant microwave absorption on the MCD and/or MCP signal may be termed optically detected magnetic resonance (ODMR). In ODMR measurements involving the MCD signal, microwave-induced transitions between the Zeeman levels of the ground state are detected by a change in intensity of the absorption (i.e., MCD) spectrum. Electron spin resonance transitions in the excited state are signaled by microwave-induced changes in the MCP signal.

Figure 15 is a schematic drawing of an ODMR spectrometer. There are three necessary channels: a microwave channel and channels for optical excitation and detection. The microwave system is relatively simple, comprising a klystron or Gunn diode operating at some frequency in the range 8.5 to 50 GHz, followed by an isolator to protect the microwave source from unwanted reflected signals in the waveguide path. The microwave power is then square-wave modulated at frequencies up to 10 kHz, using a PIN diode. A variable attenuator determines the power incident upon the resonant cavity, although for high-power operation a traveling-wave amplifier might be added to the waveguide system. The sample is contained in the microwave cavity, which is designed to allow optical access of the sample by linearly or circularly polarized light traveling either parallel or perpendicular to the magnetic field direction. The cavity is submerged in liquid helium to achieve as large a population difference as possible between the Zeeman levels. The magnetic field is provided either by an electromagnet ($B \approx 0-2.0$ T) or a superconducting solenoid ($B \approx 0-6.5$ T). Radiation from the sample is focused onto the detection system, which in its simplest form consists of suitable filters, a polarizer, and photomultiplier tube. A high-resolution monochromator may be used instead of the filters to resolve sharp features in the optical spectrum. The signal from the phototube is processed using a phase-sensitive detector, or alternatively using computer data collection with a multichannel analyzer or transient recorder. The recorded spectrum is plotted out using a pen recorder as a function of either magnetic field or photon energy (or wavelength). With such an experimental arrangement one may examine the spectral dependence of the ODMR signal on the wavelength of the optical excitation or on the wavelength of the detected luminescence by use of one of the two scanning monochromators.

In order to carry out ODMR, microwave radiation of fixed frequency ν is introduced while the optical wavelength is kept at the positive or negative peak in Fig. 14c. The magnetic field is then adjusted until the ESR condition, $h\nu = g\mu_B B$, is satisfied. Since ESR transitions tend to equalize the populations N_+ and N_- , resonance is observed as a decrease in $\Delta\alpha(\nu)$, and as the microwave power is increased, the MCD gradually tends to zero. In certain circumstances the ground-state spin polarization may be used to monitor *excited-state* ESR transitions because of the selectivity of the transitions induced by circularly polarized radiation. This measurement technique is an example of *trigger detection* where one microwave photon in absorption triggers the detection of one optical photon emitted. The resulting enhancement in sensitivity relative to the normal ESR technique is approximately in the ratio of optical to microwave frequency (i.e., $10^{15}/10^{10} = 10^5$). At *x*-band (≈ 10 GHz), the ESR sensitivity is about 10^{10} spins per gauss linewidth so that ODMR sensitivity is of order 10^5

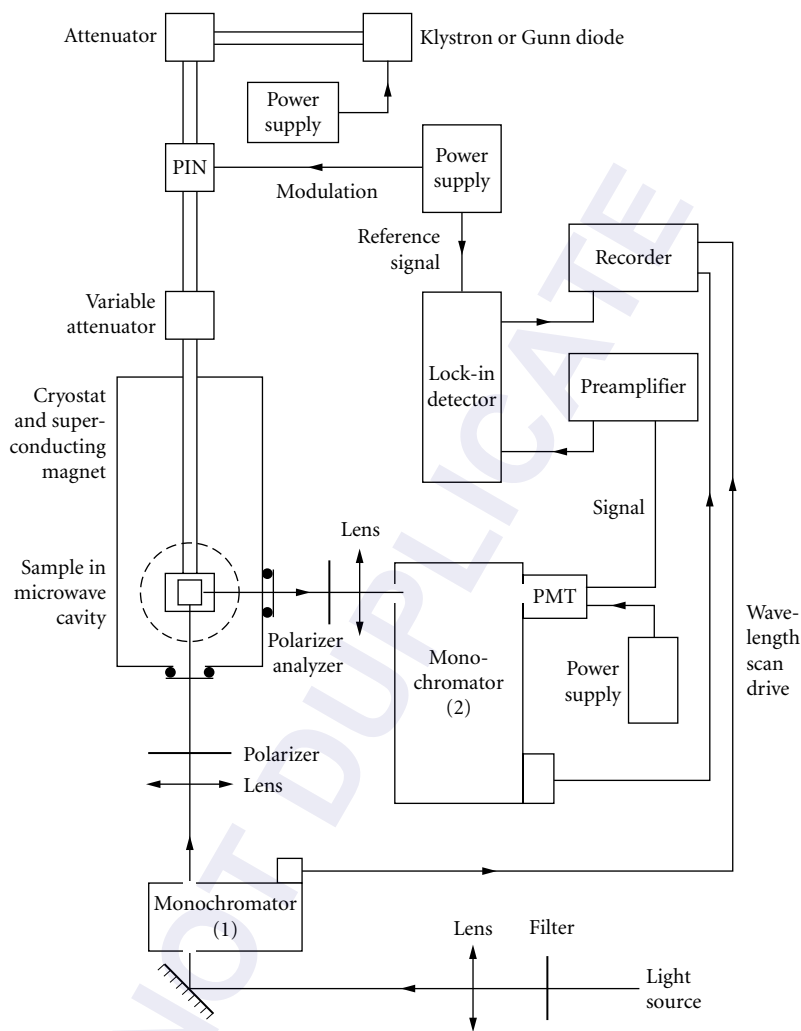


FIGURE 15 A schematic representation of a spectrometer for measuring optically detected magnetic resonance spectra via circularly polarized absorption or emission transitions.

atoms in the excited state. With the ODMR technique, one may gather information on a wide range of important solid-state processes including spin-lattice and cross relaxation, spin memory, energy transfer, electron-hole recombination, phonon bottlenecks, and spin coherence effects.

A major attribute of the ODMR technique is illustrated in Fig. 16, showing the optical characteristics of the ODMR spectrum of F centers in calcium oxide.¹⁰ These spectra were measured at 18.7 GHz and 1.6 K with the magnetic field along a crystal $\langle 100 \rangle$ direction. A high-pressure xenon discharge lamp and monochromator (M_1 in Fig. 15) set at 400 nm was used to excite the fluorescence, which was detected through monochromator M_2 . The spectrum consists of four equally spaced lines due to an $S = 1$ state of a center with tetragonal symmetry. Then with the magnetic field set at the strongest ODMR line, the excitation wavelength is scanned using monochromator M_1 (Fig. 15) over the visible and near-ultraviolet region. A single broad structureless excitation peak is observed

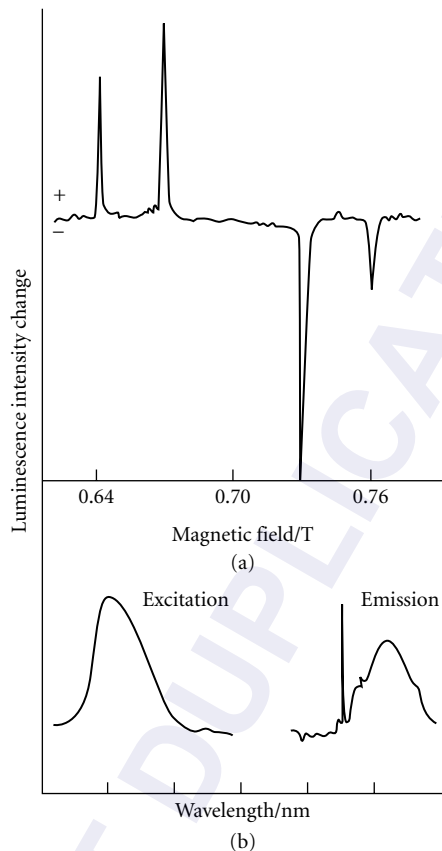


FIGURE 16 The ODMR spectrum of triplet state of F -centers in CaO .

at 400 nm corresponding to the ${}^1A_{1g} \rightarrow {}^1T_{1u}$ absorption band of the F center (Fig. 16). Subsequently, the excitation monochromator is set at the peak of this excitation band and the same magnetic field while the detecting monochromator (M_2 in Fig. 15) is scanned over the fluorescence spectrum. This spectral dependence (Fig. 16) shows a sharp zero-phonon line at a wavelength of 574 nm with an accompanying broad vibronic sideband with peak at 602 nm. In a single experiment, a unique and unambiguous relationship is established between the ESR spectrum, absorption, and fluorescence bands of an intrinsic lattice defect.

31.7 HIGH-RESOLUTION TECHNIQUES

Inhomogeneous broadening arises when individual atoms are distinguished by the frequency at which they absorb light. The absorption profile is then the sum of separate absorption lines. In atomic spectroscopy, the major source of the spectral linewidth is Doppler broadening; the frequency shift is $(\Delta\nu/\nu) = \pm(v_z/c)$ due to an atom moving with velocity component v_z towards (+) or away from (-) the observer. At thermal equilibrium, a gaussian lineshape is observed because of the

Maxwell-Boltzmann velocity distribution. In solids, the distribution of internal strains is a source of inhomogeneous broadening. Because crystals contain imperfections, electronic centers experience crystal fields which vary slightly from site to site in the crystal; in consequence, zero-phonon lines may have linewidths of order 0.1 to 50 cm^{-1} . The use of narrow-band laser excitation makes it possible to eliminate inhomogeneous broadening and to realize a resolution limited only by the homogeneous width of the transition, which in crystals can vary from kilohertz to gigahertz. This factor of 10^3 to 10^4 improvement in resolution enables the spectroscopist to carry out high-resolution studies of the physical properties and electronic structures of centers and of the mechanisms responsible for homogeneous broadening. Contributions to homogeneous width come from population dynamics and random modulation of the optical frequency by phonons and nuclear spins.

Saturated Absorption and Optical Holeburning

The experimental basis of recovering the homogeneous width of an inhomogeneously broadened optical spectrum, so-called saturated absorption or optical holeburning (OHB) spectroscopy, is illustrated in Fig. 17. An inhomogeneously broadened line of width Γ_{inh} is produced by many narrow components of homogeneous width $\Gamma_{\text{hom}} \ll \Gamma_{\text{inh}}$. Each component is centered at a different

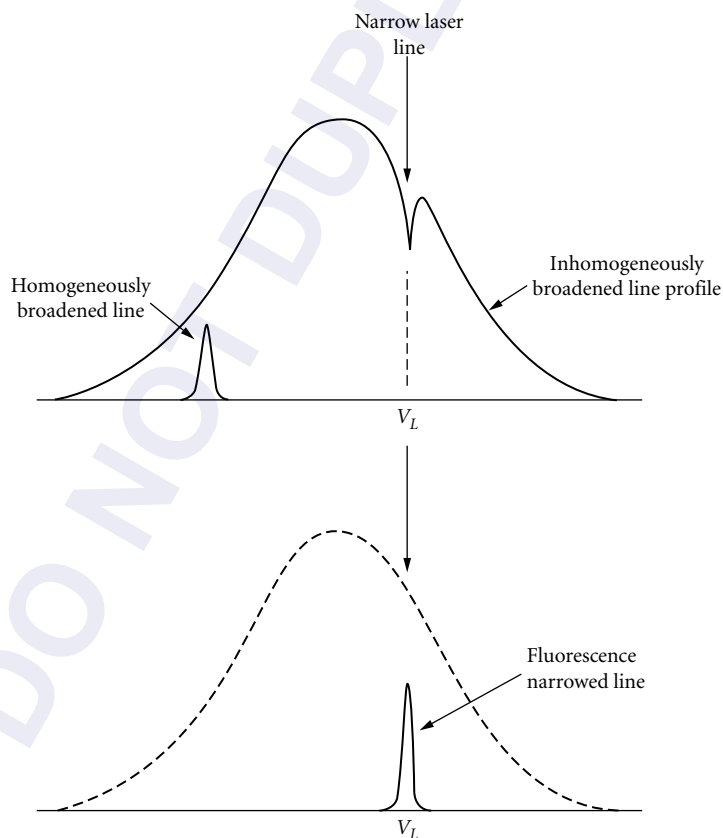


FIGURE 17 Optical holeburning (OHB) and fluorescence line narrowing (FLN) of an inhomogeneously broadened spectroscopic line.¹

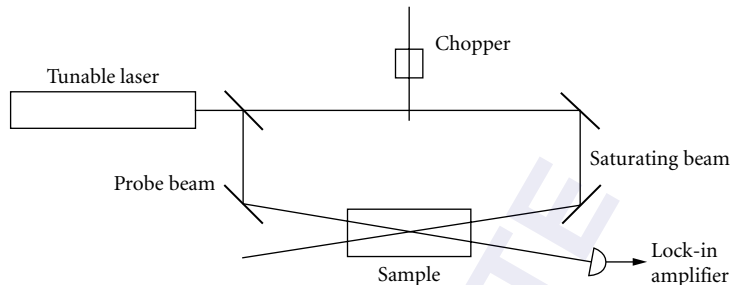


FIGURE 18 A spectrometer system for Doppler-free saturation spectroscopy.¹²

frequency within the inhomogeneous line profile. If a narrow laser line of frequency ν_L and bandwidth $\Gamma_L < \Gamma_{\text{hom}}$ is incident upon an atomic assembly having an inhomogeneously broadened linewidth Γ_{inh} , the resulting absorption of laser radiation depletes only that subassembly of excited centers whose energies are within Γ_{hom} of the laser line frequency ν_L . Consequently, a “hole” is burned in the lineshape in the neighborhood of ν_L . Resolution of the homogeneous width requires that $\Gamma_L < \Gamma_{\text{hom}} \ll \Gamma_{\text{inh}}$. In holeburning spectroscopy, the narrow laser linewidth and high power make it possible to maintain a significant fraction of those atoms with transition frequency ν_L in the excited state, where they no longer contribute to the absorption at this frequency. To observe holeburning experimentally requires that ~5 percent of those centers within the pump laser bandwidth be transferred to the excited state.

The first measurements of optical holeburning or saturated absorption spectroscopy in atoms were made by the Stanford group on the H_{α} -line ($n=2 \rightarrow n=3$) in hydrogen using a pulsed dye laser.¹¹ A schematic diagram of an appropriate absorption spectrometer is shown in Fig. 18. A strong laser beam interacts with those atoms that are moving with the right velocity to Doppler and shift them into resonance. If the laser beam is intense enough, it tends to equalize the population in the two levels, thereby reducing the intensity. The hole burned in the absorption profile, which extends over the natural width of the transition, is probed by a second beam at the same frequency but lower intensity and traveling in the opposite direction. This beam interacts with atoms having the same velocity but in the opposite direction to the saturating beam. When the laser is tuned to line center, both pump and probe beams interact with atoms moving with zero longitudinal velocity. The probe beam then measures the reduced absorption caused by the saturating beam. In experiments using pulsed lasers, very high intensity is required to achieve saturation and hence there must be very tight focusing and overlap of pump and probe beam. In consequence, CW lasers are preferred in both gas-phase and solid-state spectroscopy. Saturated absorption measurements on atomic hydrogen have been reviewed by Ferguson and Tolchard¹² and OHB in solids by Selzer¹³ and by Yen.¹⁴

To burn a hole in a narrow absorption line in crystals requires that the laser be focused onto the sample for periods of order 10^2 to 10^3 s, depending upon the specific system. When the laser excitation is switched off, the holes recover on some timescale characteristic of the physical process responsible for holeburning. For short-lived holes the exciting beam is divided using a beam splitter into pump and probe beams. The weaker probe beam passes through an optoacoustic modulator which scans it backward and forward over the hole.¹³ To observe long-lived holes, the sample is irradiated for a short time in the zero-phonon line with a few hundred milliwatts of single-mode dye laser light with a width of a few megahertz. The shape of the hole is then displayed by reducing the laser intensity to a few milliwatts and scanning the laser over the inhomogeneous line profile. Figure 19 shows an example of holeburning in the 601.28-nm line of $\text{Pr}^{3+}:\text{LaCl}_3$. The homogeneous width measured in this holeburning experiment is $\Gamma_{\text{hom}} = 10$ MHz, which corresponds to a lifetime of 10 ns. There have been many reports of holeburning spectroscopy on transition metal ions, rare-earth ions, and color centers in inorganic materials. For rare-earth ions, holeburning with lifetimes determined by the nuclear spin relaxation processes have been reported to vary from 10 to 10^3 s. Many measurements are aimed at the mechanisms leading to the homogeneous width of optical transitions.

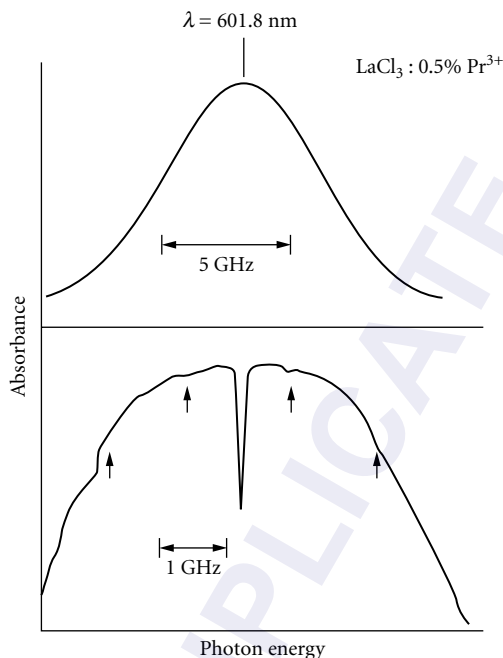


FIGURE 19 Optical holeburning in the 601.28-nm line of Pr^{3+} ions in LaCl_3 . The hole was burned using some 200 W cm^{-2} of single-frequency laser light with a bandwidth of 2 MHz. The zero-phonon line has an inhomogeneous width of 7.5 GHz. (After Harley and Macfarlane, 1986, unpublished.)

In these cases, techniques have been developed for the detection of coherent transients (e.g., photon echo or free induction decay) because the measurements are made on the timescale of the dephasing and are not affected by spectral diffusion and other such processes.

Polarized Absorption Spectrometers

Polarized absorption spectroscopy is a technique related to sub-Doppler absorption spectroscopy. However, in this case use is made of the circularly polarized absorption properties of atomic transitions. In the Wieman-Hansch experiment,¹¹ the probe beam passes through crossed polarizers immediately before and after the sample. If the pump beam is unpolarized, the sample is optically isotropic and no light falls on the detector. However, if the $s \rightarrow p$ transitions are excited using RCP light, the pump beam induces optical anisotropy in the medium with which it interacts. In consequence, as pump and probe beams are tuned to line center so that both interact with the same class of atoms, the weak probe beam becomes slightly elliptically polarized and light is transmitted through the crossed polarizers. The advantage of the method is a factor of about 10^3 enhancement in sensitivity relative to saturation absorption spectroscopy. Sub-Doppler two-photon absorption spectroscopy is also much used in atomic physics.¹⁵ The selection rule for two-photon absorption is that $\Delta l = 0$ or 2. In consequence, for $l = 1$ electron atoms $S \rightarrow S$ and $S \rightarrow D$ transitions are allowed.

Laser Stark Spectroscopy of Molecules

Sub-Doppler resolution enhancement is also used in studying the heterogeneously broadened rotational/vibrational spectra of molecules. Such spectra are generally observed in the mid-IR region and are studied using a wide variety of gas lasers (e.g., N_2O , CO, and CO_2). Such laser lines are not usually in exact resonance with the particular molecular transition: laser and molecular transition are brought into register using a variable electric field to tune the molecular system into resonance. In general, parallel-plate Stark cells are used in which free-space propagation of the short-wavelength infrared radiation occurs. This makes it easy to use both perpendicular and parallel polarization configurations in the electric resonance experiments so that both $\Delta M_J = 0$ and $\Delta M_J = \pm 1$ transitions are observed. The subject of laser Stark spectroscopy has been discussed at length by Duxbury.¹⁶

A schematic intracavity laser Stark spectrometer is shown in Fig. 20; the same basic principles are obtained as with optical holeburning spectroscopy. The effects of the saturating laser field are confined to a narrow frequency region centered on the velocity component of those molecules whose absorption is Doppler-shifted into resonance. In a standing wave field, two holes are burned, one on either side of the line center, corresponding to molecules moving toward or away from the detector. The applied electric field is used to tune the two holes to line center where they coalesce to give a sharp dip in the absorption coefficient at line center. Since the resonance method relies on the use of an electric field for tuning, it is necessary both to generate high uniform fields and to study molecules with appreciable Stark tuning coefficients. In order to generate high electric fields, which may approach 90 kV cm^{-1} , narrow electrode spacings from 1 to 4 mm are commonly used. With such narrow gaps, the plates must be flat to one or two fringes of visible light, and must be held accurately parallel. The gas pressure used must also be restricted to the low-pressure region below 100 mtorr. A potential difference of roughly 3000 V may be sustained without electrical breakdown across any gas at a pressure of 100 mtorr and below.

The electric field is then modulated at some convenient frequency to permit the use of phase-sensitive detection techniques. In order to get above the principal noise region of the electric discharge lasers used in the 5- and 10- μm regions and as pumps for the FIR lasers, it is necessary to use electric field modulation frequencies in the range from 5 to 100 kHz. The amplitude of the electric field modulation used to detect the signals is usually small compared to the equivalent electric field linewidth of the transitions. The most common modulation waveform is sinusoidal.

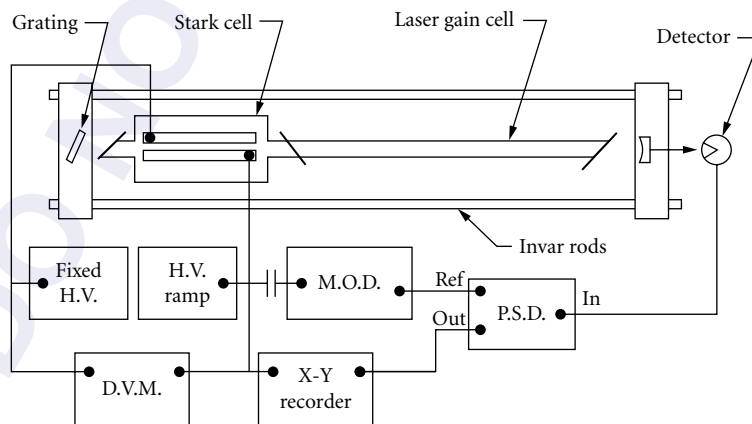


FIGURE 20 Schematic diagram of an intracavity laser Stark spectrometer. PSD stands for phase sensitive detector, DVM for digital voltmeter, HV for high voltage, and MOD for modulation source.¹⁶

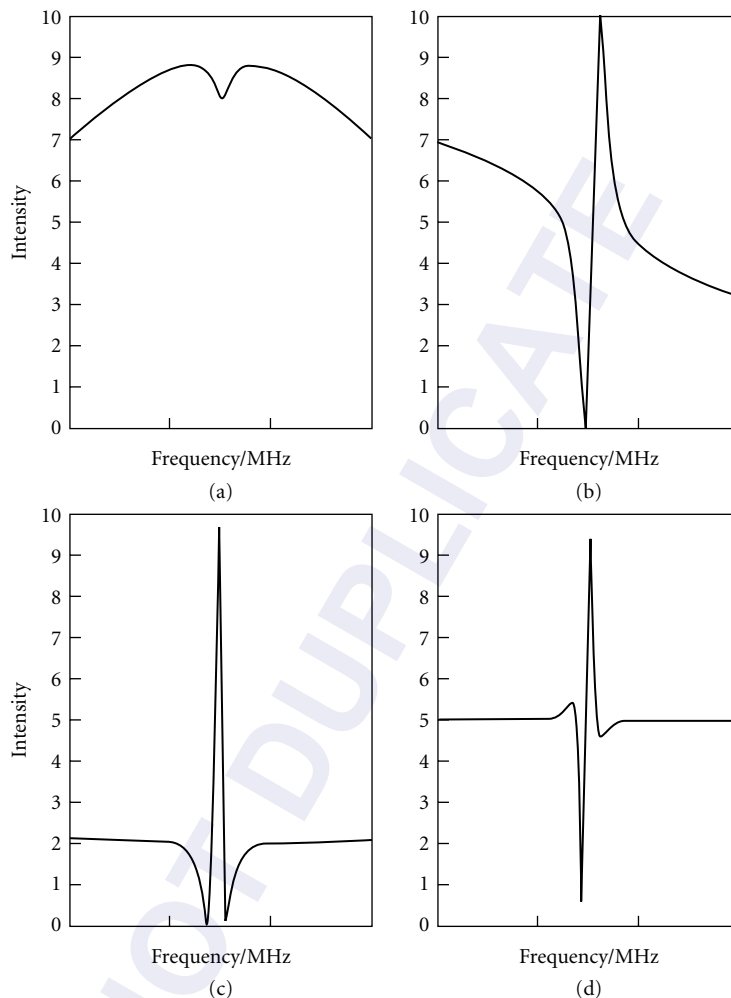


FIGURE 21 Lineshapes which occur when detecting at harmonics of the modulation frequency when small-amplitude field modulation is used. Doppler-broadened line showing the Lamb dip. (a) 30-MHz scan of a partially saturated Doppler-broadened line showing the Lamb dip; (b) 30-MHz scan with first-derivative detection; (c) 30-MHz scan with second-harmonic, second-derivative detection. The gain is increased by a factor of four from (a); and (d) 30-MHz scan with third-harmonic, third-derivative detection.

If the modulation amplitude is much smaller than the linewidth, detection at the fundamental modulation frequency results in a first derivative lineshape as in analogous electron spin resonance spectra. In order to remove the effects of sloping baselines produced by transitions with a slow Stark effect, it is common to use detection at either the second or third harmonic of the modulation frequency. Second-harmonic detection produces a second-derivative signal resembling a sharpened absorption line but with negative side lobes. Third-harmonic detection produces a third-derivative signal which resembles a sharpened first derivative, but which again possesses side lobes. Theoretical lineshapes are illustrated in Fig. 21. Second- and third-harmonic detection are particularly useful for

the observation of narrow saturation features free from background effects. The detectors used are quantum-limited liquid nitrogen cooled devices, PbSnTe or CdHgTe in the 10- μm region and InSb or Au doped Ge in the 5- μm region. In the far infrared, Goly cells have been used but in order to achieve a better signal-to-noise ratio it is necessary to use detectors cooled by liquid helium.

Just as in atomic spectroscopy one may use atomic beam spectroscopy as an alternative to absorption saturations, so too may one use molecular beam systems in high-resolution studies of the rotational-vibrational spectra of molecules.

Fluorescence Line Narrowing

Fluorescence line narrowing (FLN) is a technique complementary to that of OHB. It may also be understood by referring to Fig. 17. A narrow laser line is used to pump within the inhomogeneous linewidth Γ_{inh} . The laser interacts only with the subset of levels spanning the bandwidth of the laser Γ_L . These centers reradiate to some lower lying level, with a fluorescence linewidth much narrower than the inhomogeneous width. The fluorescence linewidth approaches the homogeneous width. In fact, for centers involved in a resonance fluorescence transition, the total FLN lineshape is a convolution of the laser lineshape and twice the homogeneous lineshape (once for the pump bandwidth and once for the fluorescence). The FLN linewidth Γ is then usually written as $\Gamma = \Gamma_L + 2\Gamma_h$. Experimentally, FLN requires a little more sophistication than does holeburning spectroscopy. Of course, one still requires a stable, high-resolution laser. Care must be used in extracting the true homogeneous linewidth, especially for nonresonant fluorescence. Many of the experimental problems relative to solid samples are discussed in the review by Selzer,¹⁵ and numerous examples are given by Yen and Selzer.¹⁷ The CW FLN spectrum shown in Fig. 22 is for the Cr^{3+} ion in aluminum oxide.¹⁸ The fluorescence lifetime is 3.4 ms at 4.2 K. Hence the homogeneous width is of the order 0.3 kHz. A direct-phonon relaxation process between the two ${}^2\text{E}$ levels, $2\bar{\text{A}}$ and $\bar{\text{E}}$, separated in energy by 29 cm^{-1} , broadens the homogeneous width to ≈ 130 kHz. In CW measurements, a homogeneous width in excess of 100 MHz was reported.¹⁸ The problem is relaxations due to super-hyperfine interactions with neighboring aluminum nuclei. The application of a dc magnetic field of only 40 mT has the effect of inhibiting relaxation due to local fields at the Cr^{3+} ions due to the ${}^{27}\text{Al}$ nuclear moments. A very considerable narrowing of the Cr^{3+} FLN spectrum is then achieved.

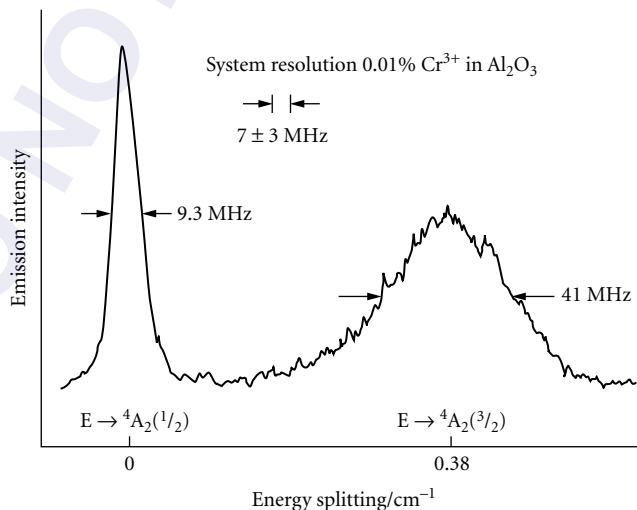


FIGURE 22 FLN in the R_1 transition of ruby.¹⁸

31.8 LIGHT SCATTERING

Light-scattering experiments are now a routine feature in many optical laboratories. The first observations of light scattering by small particles were reported by Tyndall.¹⁹ Subsequently, theoretical work by Lord Rayleigh²⁰ showed both that the scattered intensity varied as the fourth power of the frequency and that the scattering was due to molecules rather than dust particles. Many of the early studies were concerned with the depolarization of the light after being scattered by the dust-free atmosphere. Of course, in the pre-laser era, sufficient light intensity could only be achieved by use of strongly condensing lenses to focus light onto the gas cell. Very great care was then necessary to obtain reliable depolarization measurements. Even in the laser era it is still essential to avoid any effects due to parasitic light which often plague light-scattering experiments.

A significant early result from scattering of light by gases was that the scattered light intensity varied with the density of the gas being used as the sample. However, Lord Rayleigh discovered that the intensity scattered per molecule decreased by a factor of order 10 on condensation to the liquid phase. There is a somewhat smaller decrease in going from the liquid phase to the solid. Obviously, some scattering experiments become rather difficult in the solid state. The classical experimental geometry for studying Rayleigh scattering is in the 90° orientation for the scattered radiation. This is also the most useful orientation for Raman scattering in solids.²¹

One important feature of the structure of solids is the periodic disturbance of the crystal structure by the propagation of quantized elastic waves (i.e., phonons). Those elastic waves which travel at the velocity of sound (i.e., sonic waves) are essentially thermal density fluctuations in the elastic medium. Brillouin predicted that such fluctuations should give rise to fine structure in the Rayleigh scattered light when the Bragg coherence condition $\lambda_1 = 2\lambda_p \sin(\phi/2)$ is obeyed. Here λ_1 is the wavelength of light, λ_p is the wavelength of those phonons responsible for scattering the light, and ϕ is the scattering angle. Because the scattering centers are in motion, the scattered light is frequency shifted by the Doppler effect. It is an easy matter to show that the Doppler shift, $\Delta\nu$, is given by

$$\Delta\nu = \pm v_p = \pm 2v_1(v/c) \sin(\phi/2) \quad (20)$$

where v_p is the frequency of the density fluctuations in the medium and v is the velocity of sound in the medium. For light in the visible region then, that part of the phonon spectrum probed by the Brillouin scattering is in the gigahertz frequency region. In addition, the Brillouin components are completely polarized for 90° scattering. Before the advent of lasers, the study of Brillouin scattering effects in solids was exceedingly difficult. It remains a technique more used in gases than in condensed media.

Raman was one of numerous scientists engaged in research into light scattering during the decade 1920 to 1930. Much of his work was carried out using sunlight as a source. However, in experiments using monochromatic light, he observed in the spectrum of light scattered at 90° by liquid samples, new lines at wavelengths not present in the original light.²¹ The frequency displacement of these new lines from source frequency was found to be independent of the wavelength of the incident light. This was contrary both to fluorescence excitation and Brillouin scattering [Eq. (20)]; hence was born a new scattering phenomenon for which Raman was awarded the Nobel prize and which now bears his name. The frequency shifts in the Raman spectrum of a particular substance are related to but not identical to infrared absorption frequencies. In general, infrared transitions occur when there is a change in the electric dipole moment of a center as a consequence of the local atomic vibrations. The Raman lines occur when a change in polarizability is involved during atomic vibrations. This usually means that infrared transitions occur only between states of opposite parity whereas Raman transitions occur between states of the same parity. Thus the infrared and Raman spectra give complementary information about the vibrational spectra of spectroscopic centers.

Raman scattering measurements have found wide application in condensed matter physics. The spectrometer systems have much in common with fluorescence spectrometers, although lasers provide the excitation source almost without exception. Single-frequency lasers (He-Ne, Ar⁺, Ke⁺) and tunable dye lasers and solid-state lasers have all been used. Most lasers provide a polarized output and it is necessary to modify this to allow the excitation polarization to be varied. The scattered

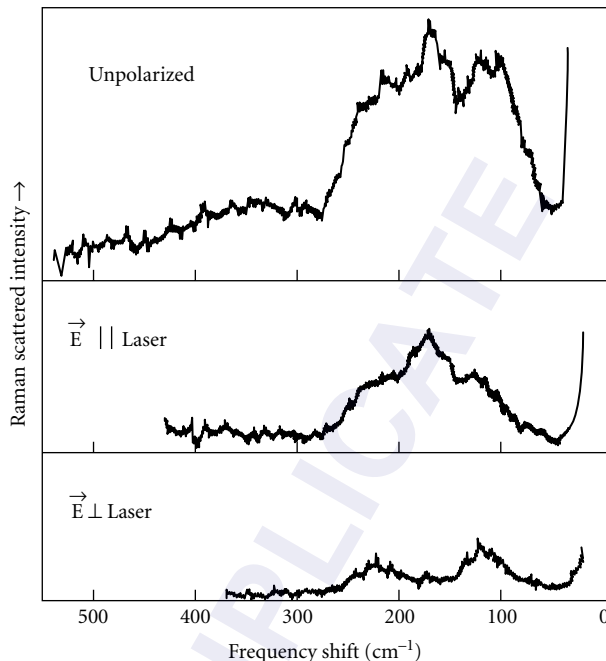


FIGURE 23 Raman spectra of F -centers in NaCl for different polarizations of the detected radiation. The shifts are measured relative to the 514.5 nm wavelength of the Ar⁺ laser.²⁰

light is observed via a monochromator in a direction normal to the laser beam. Again, provision is made for the polarization of the scattered radiation to be analyzed. To permit observation closer to the laser line, double or triple monochromators are used to eliminate all traces of stray light. Furthermore, one must take trouble to separate out the Raman-scattered light from any fluorescence signal. Since the Raman signal is instantaneous, it is comparatively straightforward to recover the desired signal from the decaying fluorescence signal using time-resolution techniques.

An example of the application of Raman spectroscopy in color center physics is shown in Fig. 23. The intensity of scattering versus wavelength shift from the Ar⁺ laser excitation is shown for F -centers in NaCl²² for which the longitudinal optic frequency is 270 cm⁻¹. The major Raman-shifted spectrum occurs below 200 cm⁻¹, showing that the vibrational interaction is due to ionic displacements close to the defect. These local modes have broad peak centers near $\hbar\omega = 175$ cm⁻¹. A comparison of the polarized and unpolarized excitation spectra shows that the local mode scattering is supplemented by a lattice vibrational contribution covering much of the 0 to 500 cm⁻¹ frequency shift.

31.9 REFERENCES

1. B. Henderson and G. F. Imbusch, *Optical Spectroscopy of Inorganic Solids*, Oxford University Press, Oxford, England, 1989.
2. J. R. Lakowicz, *Principles of Fluorescence Spectroscopy*, Plenum Press, New York, 1983.
3. H. Engstrom and L. F. Mollenauer, *Phys. Rev. B* **7**:1616 (1973).
4. M. O. Henry, J. P. Larkin, and G. F. Imbusch, *Phys. Rev. B* **13**:1893 (1976).

5. J. R. Lakowicz and H. Charek, *J. Biol. Chem.* **256**:6348 (1981).
6. L. F. Mollenauer, J. M. Wiesenfeld, and E. P. Ippen, *Radiation Effects* **72**:73 (1983).
7. P. P. Feofilov, *The Physical Basis of Polarized Emission*, Consultants Bureau, New York, 1961.
8. J. C. Kemp, *J. Opt. Soc. Amer.* **59**:915 (1966).
9. S. Geschwind, R. J. Collins, and A. L. Schawlow, *Phys. Rev. Lett.* **3**:545 (1959).
10. P. Edel, C. Hennies, Y. Merle d'Aubigné, R. Romestain, and Y. Twarowski, *Phys. Rev. Lett.*, **28**:1268 (1972).
11. The Stanford group made a number of experimental improvements; see for example T. W. Hänsch, I. S. Shakin, and A. L. Schawlow, *Nature (Lond.)* **235**:63 (1972); T. W. Hänsch, M. H. Nayfeh, S. A. Lee, S. M. Curry, and I. S. Shakin, *Phys. Rev. Lett.* **32**:1336 (1974); and C. E. Wieman and T. W. Hänsch, *Phys. Rev. Lett.* **36**:1170 (1976).
12. A. I. Ferguson and J. M. Tolchard, *Contemp. Phys.* **28**:383 (1987).
13. P. M. Selzer, in W. M. Yen and P. M. Selzer (eds.), *Laser Spectroscopy of Solids I*, Springer-Verlag, Berlin, 1981, p. 113.
14. W. Yen, in W. M. Yen (ed.), *Laser Spectroscopy of Solids II*, Springer-Verlag, Berlin (1988).
15. B. Cagnac, G. Grynberg, and F. Biraben, *J. Phys. Paris* **34**:845 (1973).
16. G. Duxbury, *International Reviews in Physical Chemistry* **4**:237 (1985).
17. W. M. Yen and P. M. Selzer, in W. M. Yen and P. M. Selzer (eds.), *Laser Spectroscopy of Solids I*, Springer-Verlag, Berlin, 1981.
18. P. E. Jessop, T. Muramoto, and A. Szabo, *Phys. Rev. B* **21**:926 (1980).
19. J. W. Tyndall, *Notes of a Course of Nine Lectures on Light*, Royal Institution of Great Britain, Longmans, London, 1869.
20. Lord Rayleigh (Hon. J. W. Strutt), *Phil. Mag.* **x/i**:107 (1871).
21. C. V. Raman, *Indian J. Phys.* **2**:1 (1928).
22. J. M. Worlock and S. P. S. Porto, *Phys. Rev. Lett.* **15**:697 (1965).

INTERFEROMETERS

Parameswaran Hariharan

*School of Physics
University of Sydney
Sydney, Australia*

32.1 GLOSSARY

A	area
a	amplitude
C	ratio of peaks to valleys
d	distance
E	electric field
F	finesse
FSR	free spectral range
I	intensity
$J_i(\)$	Bessel function
L	length
m	integer
N	number of fringes
n	refractive index
p	optical path difference
R	reflectance
r	radius
T	transmittance
v	velocity
λ	wavelength
θ	angle
ν	frequency
φ	phase
ψ	phase difference
ω	angular velocity

32.2 INTRODUCTION

Optical interferometers have made possible a variety of precision measurements using the interference phenomena produced by light waves.^{1,2} This chapter presents a brief survey of the basic types of interferometers and discusses some of their applications.

32.3 BASIC TYPES OF INTERFEROMETERS

Interferometric measurements require an optical arrangement in which two or more beams, derived from the same source but traveling along separate paths, are made to interfere. Interferometers can be classified as *two-beam* interferometers or *multiple-beam* interferometers according to the number of interfering beams; they can also be grouped according to the methods used to obtain these beams.

The Fizeau Interferometer

In the Fizeau interferometer, as shown in Fig. 1, interference fringes of equal thickness are formed between two flat surfaces separated by an air gap and illuminated with a collimated beam. If one of the surfaces is a standard reference flat surface, the fringe pattern is a contour map of the errors of the test surface. Absolute measurements of deviations from flatness can be made by an intercomparison of three surfaces. Modified forms of the Fizeau interferometer are also used to test convex and concave surfaces by using a converging or diverging beam.³

The Michelson Interferometer

The Michelson interferometer, shown schematically in Fig. 2, uses a beam splitter to divide and recombine the beams. As can be seen, one of the beams traverses the beam splitter three times, while the other traverses it only once. Accordingly, a compensating plate of the same thickness as the beam splitter is introduced in the second beam to equalize the optical paths in glass. With an extended source, the interference pattern is similar to that produced in a layer of air bounded by the mirror M_1 and M_2' , the image of the other mirror in the beam splitter. With collimated light, fringes of equal thickness are obtained. The Michelson interferometer modified to use collimated light (the Twyman-Green interferometer) is used extensively in optical testing.⁴

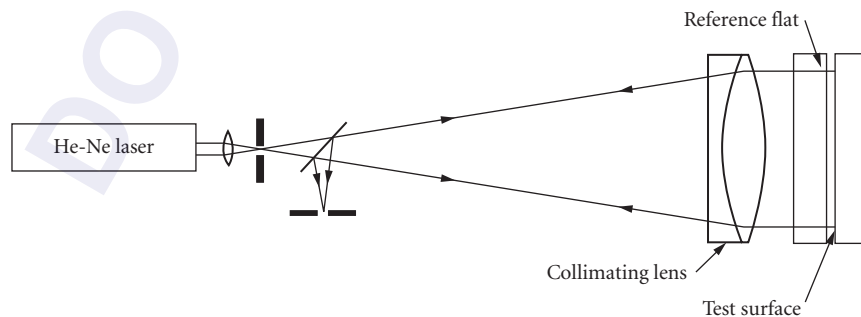


FIGURE 1 The Fizeau interferometer.

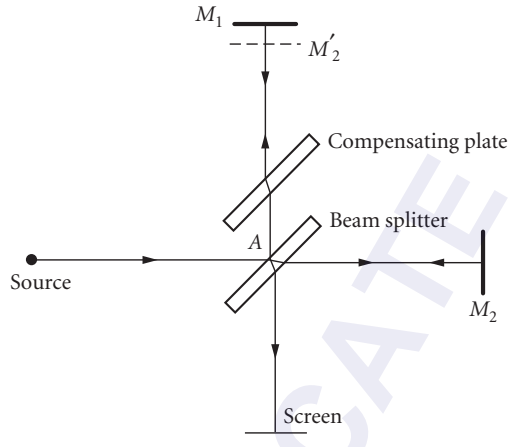


FIGURE 2 The Michelson interferometer.

The Mach-Zehnder Interferometer

The Mach-Zehnder interferometer uses two beam splitters and two mirrors to divide and recombine the beams. As shown in Fig. 3, the fringe spacing and the plane of localization of the fringes obtained with an extended source can be controlled by varying the angle between the beams and their lateral separation when they emerge from the interferometer. The Mach-Zehnder interferometer has been used for studies of gas flows and plasmas.

The Sagnac Interferometer

In the Sagnac interferometer, as shown in Fig. 4, the two beams traverse the same closed path in opposite directions. Because of this, the interferometer is extremely stable and easy to align, even with an extended broadband light source.

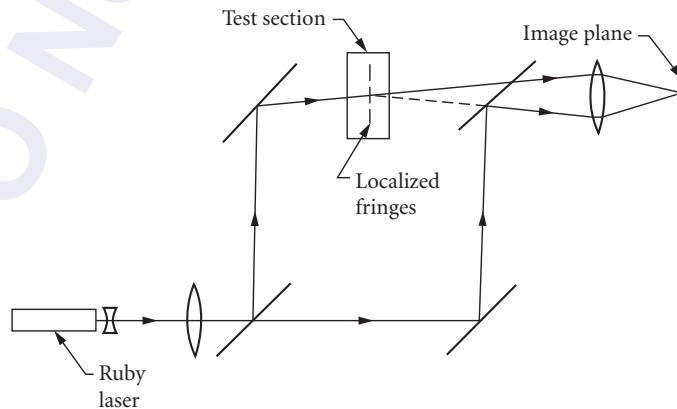


FIGURE 3 The Mach-Zehnder interferometer.

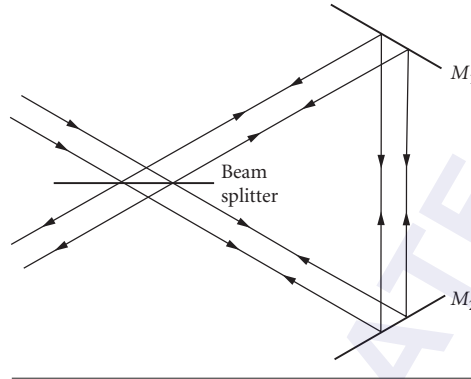


FIGURE 4 The Sagnac interferometer.

The Sagnac interferometer has been used for rotation sensing. When the interferometer is rotated with an angular velocity ω about an axis making an angle θ with the normal to the plane of the interferometer, a phase difference φ is introduced between the beams given by the relation

$$\varphi = (8\pi\omega A \cos \theta) / \lambda c \quad (1)$$

where A is the area enclosed by the light path, λ is the wavelength, and c is the speed of light.

Polarization Interferometers

Polarization interferometers are used in interference microscopy.⁵ The Nomarski interferometer, shown schematically in Fig. 5, uses two Wollaston (polarizing) prisms to split and recombine the beams. If the separation of the beams in the object plane (the lateral shear) is small compared to the dimensions of the object, the optical path difference corresponds to the phase gradients in the test object.

Grating Interferometers

Gratings can be used as beam splitters in the Michelson and Mach-Zender interferometers. Such an arrangement is very stable, since the angle between the beams is affected only to a small extent by the orientation of the gratings. Figure 6 is a schematic of an interferometer that has been used to test fine-ground surfaces at grazing incidence utilizing two diffraction gratings to split and recombine the beams.⁶

Shearing Interferometers

Shearing interferometers are widely used for optical testing, since they eliminate the need for a reference surface. As shown in Fig. 7, in a lateral shearing interferometer two images of the test wavefront are superimposed with a mutual lateral displacement, while in a radial shearing interferometer one of the images is contracted or expanded with respect to the other.^{7,8}

The Fabry-Perot Interferometer

The Fabry-Perot interferometer⁹ is used widely in high-resolution spectroscopy. It consists of two flat, parallel surfaces with highly reflecting, semitransparent coatings. If the surfaces are separated by a distance d and the medium between them has a refractive index n , the normalized value of the transmitted intensity at a wavelength λ for rays traversing the interferometer at an angle θ is

$$I_T(\lambda) = T^2 / (1 + R^2 - 2R \cos \varphi) \quad (2)$$

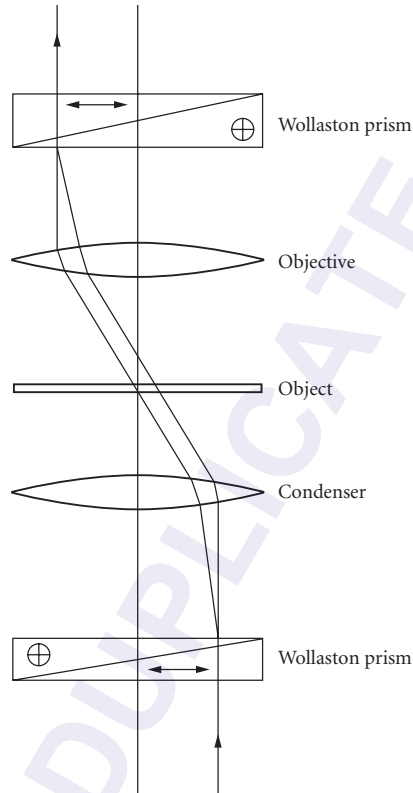


FIGURE 5 The Nomarski interferometer.

where T and R are, respectively, the transmittance and reflectance of the surfaces and $\varphi = (4\pi/\lambda)nd \cos \theta$. With an extended source of monochromatic light, the fringes seen by transmission are narrow, concentric rings. The free spectral range (FSR), which corresponds to the range of wavelengths that can be handled without successive orders overlapping, is given by the relation

$$\text{FSR}_\lambda = \lambda^2/2nd \quad (3)$$

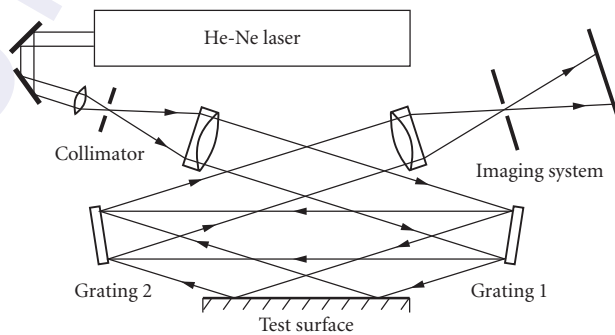


FIGURE 6 Grating interferometer used to test fine-ground surfaces at grazing incidence. (From Ref. 6.)

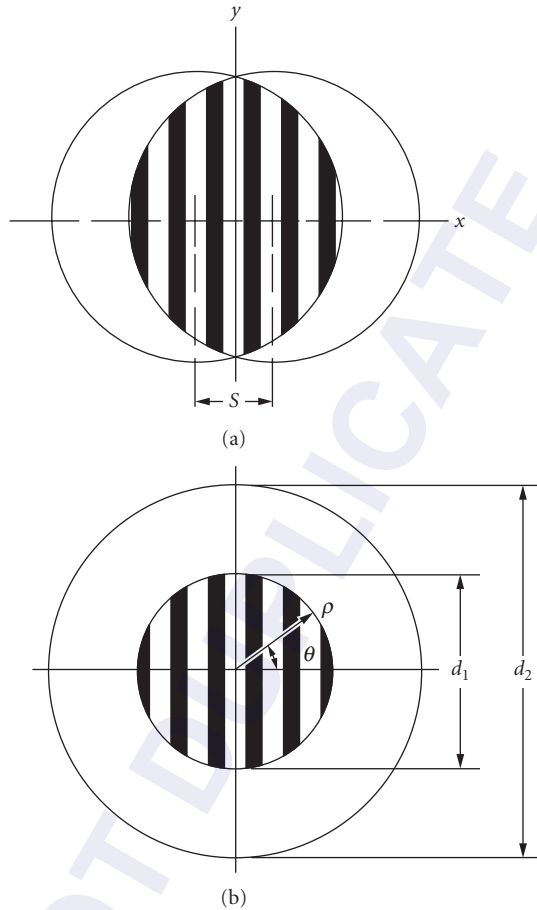


FIGURE 7 Fields of view in (a) lateral and (b) radial shearing interferometers.

while the width of the peaks at half the maximum intensity corresponds to a change in φ given by the relation

$$\Delta\varphi = 2(1-R)/R^{1/2} \tag{4}$$

The ratio of the free spectral range to the width of the fringes at half maximum intensity is known as the *finesse* F , and is given by the relation

$$F = \pi R^{1/2} / (1-R) \tag{5}$$

Two useful variants of the Fabry-Perot interferometer are the multiple-passed Fabry-Perot interferometer and the confocal Fabry-Perot interferometer. With the conventional Fabry-Perot interferometer, the ratio of the intensity at the maxima to that at the minima between them is

$$C = [(1+R)/(1-R)]^2 \tag{6}$$

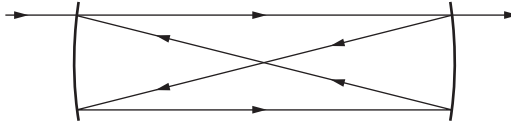


FIGURE 8 Ray paths in a confocal Fabry-Perot interferometer.

and for typical values of reflectance ($R \approx 0.95$), the background due to a strong spectral line may mask a neighboring weak satellite. A much higher contrast factor may be obtained by double- or multiple-passing the interferometer.^{10,11}

The confocal Fabry-Perot interferometer uses two spherical mirrors whose spacing is chosen, as shown in Fig. 8, so that their foci coincide. Any ray, after traversing the interferometer four times, then emerges along its original path.¹² The confocal Fabry-Perot interferometer has a higher throughput than the plane Fabry-Perot interferometer and produces a uniform output field. It is, therefore, the preferred form for operation in a scanning mode by using piezoelectric spacers to vary the separation of the mirrors.

32.4 THREE-BEAM AND DOUBLE-PASSED TWO-BEAM INTERFEROMETERS

Because of the sinusoidal intensity distribution in two-beam interference fringes, it is difficult to estimate their position visually to better than $1/20$ of their spacing. However, it is possible to detect much smaller optical path variations using the intensity changes in a system of interference fringes.

Three-Beam Interferometers

Zernike's three-beam interferometer, shown schematically in Fig. 9, uses three beams produced by division of a wavefront at a screen containing three parallel, equidistant slits.¹³ In this arrangement, the optical paths of all three beams are equal at a point in the back focal plane of the lens L_2 . The two outer slits provide the reference beams, while the beam from the middle slit, which is twice as broad, is used for measurements. The intensity at any point in the interference pattern is then given by the relation

$$I = I_0 [3 + 2 \cos 2\psi + 4 \cos \psi \cos \varphi] \quad (7)$$

where ψ is the phase difference between the two outer beams, and φ is the phase difference between the middle beam and the two outer beams at the center of the field. The intensities at adjacent maxima are equal only when φ is an odd multiple of $\pi/2$. Two positions of the plane of observation can

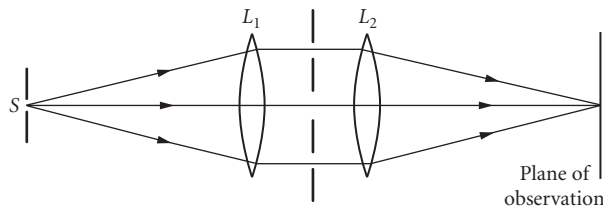


FIGURE 9 Zernike's three-beam interferometer.

be found that satisfy this condition, one inside and the other outside the focus, and any small change in the optical path of the middle beam can be measured from the shift in these positions.

Three-beam fringes can also be produced with an optical system similar to that in the Jamin interferometer.¹⁴ Settings are made by means of a compensator in the middle beam and can be repeated to $\lambda/200$ by visual observation, and to better than $\lambda/1000$ with a photoelectric detector.¹⁵

Double-Passed Two-Beam Interferometers

Fringes whose intensity is modulated in the same manner as three-beam fringes can be produced by reflecting the beams emerging from a two-beam interferometer back through the interferometer.¹⁶ In this case also, the intensity of the adjacent fringes is equal when the phase difference between the single-passed beams is

$$\varphi = (2m + 1)\pi/2 \quad (8)$$

where m is an integer. Measurements can be made with a precision of $\lambda/1000$.

32.5 FRINGE-COUNTING INTERFEROMETERS

One of the main applications of interferometry has been in accurate measurements of length using the wavelengths of stabilized lasers. Electronic fringe counting has become a practical technique for such measurements.¹⁷

The very narrow spectral line widths of lasers make it possible to use a heterodyne system. In one implementation of this technique, a He-Ne laser is forced to oscillate simultaneously at two frequencies, ν_1 and ν_2 , separated by a constant frequency difference of about 2 MHz, by applying an axial magnetic field.¹⁸ These two waves, which are circularly polarized in opposite senses, are converted to orthogonal linear polarizations by a $\lambda/4$ plate.

As shown in Fig. 10, a polarizing beam splitter reflects one beam to a fixed reflector, while the other is transmitted to a movable reflector. A differential counter receives the beat frequencies from the photodetector D_S and a reference photodetector D_R . If the two reflectors are stationary, the two beat frequencies are the same, and the net count is zero. However, if one of the reflectors is moved, the change in the optical path is given by the net count.

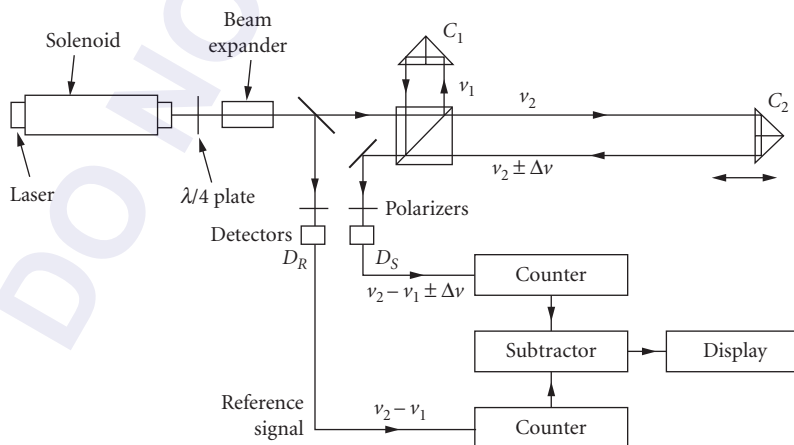


FIGURE 10 Heterodyne fringe-counting interferometer. (After Ref. 18 © Copyright Hewlett-Packard Company. Reproduced with permission.)

32.6 TWO-WAVELENGTH INTERFEROMETRY

If a length is known within certain limits, the use of a wavelength longer than the separation of these limits permits its exact value to be determined unambiguously by a single interferometric measurement. One way to synthesize such a long wavelength is by illuminating the interferometer simultaneously with two wavelengths λ_1 and λ_2 . The envelope of the fringes then corresponds to the interference pattern that would be obtained with a synthetic wavelength

$$\lambda_s = \lambda_1 \lambda_2 / |\lambda_1 - \lambda_2| \quad (9)$$

This technique can be implemented very effectively with a carbon dioxide laser, since it can operate at a number of wavelengths that are known very accurately, yielding a wide range of synthetic wavelengths.¹⁹

Two-wavelength interferometry and fringe-counting can be combined to measure lengths up to 100 m by switching the laser rapidly between two wavelengths as one of the mirrors of a Twyman-Green interferometer is moved over the distance to be measured.²⁰

32.7 FREQUENCY-MODULATION INTERFEROMETERS

New interferometric techniques are possible with laser diodes which can be tuned electrically over a range of wavelengths.²¹ One of these is frequency-modulation interferometry.

Figure 11, shows a frequency-modulation interferometer that can be used to measure absolute distances, as well as relative displacements, with high accuracy.²² In this arrangement, the signal beam reflected from the movable mirror returns as a circularly polarized beam, since it traverses the $\lambda/8$ plate twice. The reference beam reflected from the front surface of the $\lambda/8$ plate interferes with the two orthogonally polarized components of the signal beam at the two detectors to produce outputs that vary in quadrature and can be fed to a counter to determine the magnitude and sign of any displacement of the movable mirror.

To make direct measurements of the optical path difference, the frequency of the laser is ramped linearly with time by using a function generator to vary the injection current of the laser. An optical

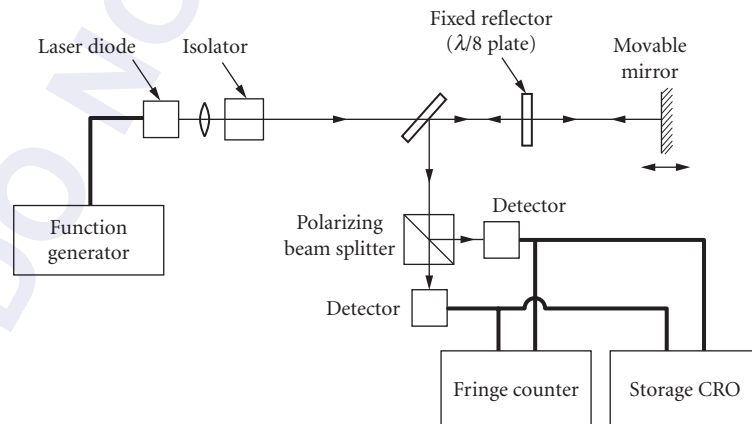


FIGURE 11 Frequency-modulation interferometer for measurements of distance. (From Ref. 22.)

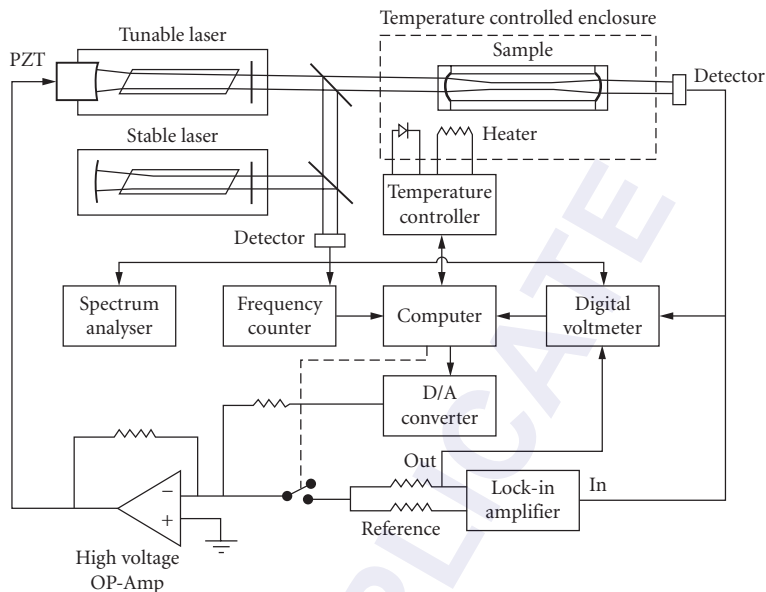


FIGURE 12 Heterodyne interferometer for measurements of thermal expansion. (From Ref. 24.)

path difference p introduces a time delay p/c between the two beams, so that they produce a beat signal with a frequency

$$f = (p/c)(dv/dt) \quad (10)$$

where dv/dt is the rate at which the laser frequency is varying with time.

32.8 HETERODYNE INTERFEROMETERS

In heterodyne interferometers, a frequency difference is introduced between the two beams by means of two acousto-optic modulators operated at slightly different frequencies. The output signal from a square-law detector then contains an ac component at the difference frequency whose phase corresponds to the phase difference between the interfering light waves.²³

Heterodyne techniques can also be used for measurements of very small changes in length.^{24,25} In the setup shown in Fig. 12, the frequency of a laser is locked to a transmission peak of a Fabry-Perot interferometer formed by attaching two mirrors to the ends of the sample. The beam from this slave laser is mixed at a photodetector with the beam from a stable reference laser. Changes in the separation of the mirrors can be evaluated from the changes in the beat frequency.

A simple arrangement for measuring small displacements uses two diode lasers with external cavities. A displacement of the reflecting mirror of one cavity results in a change in the beat frequency.²⁶

32.9 PHASE-SHIFTING INTERFEROMETERS

In phase-shifting interferometers, the phase difference between the two beams in the interferometer is varied linearly with time and the values of intensity at any point in the interference pattern are integrated over a number of equal segments covering one period of the sinusoidal signal.

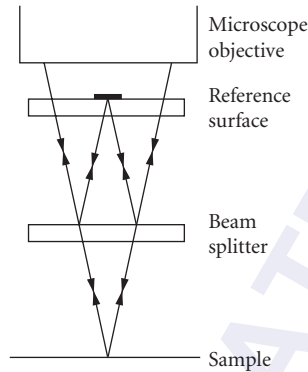


FIGURE 13 Schematic of a compact optical system (the Mirau interferometer) used for phase-stepping interference microscopy.

Alternatively the phase difference between the two beams can be changed in a number of equal steps, and the corresponding values of intensity at each data point are measured and stored. In both cases, the values obtained can be represented by a Fourier series, whose coefficients can be evaluated to obtain the original phase difference between the interfering beams at each point.^{27,28} Typically, four measurements are made at each point, corresponding to phase intervals of 90° . If I_1 , I_2 , I_3 , and I_4 are the values of intensity obtained, the phase difference between the interfering beams is given by the relation

$$\tan \varphi(x, y) = (I_1 - I_3) / (I_2 - I_4) \quad (11)$$

Phase-shifting interferometers are used widely in optical testing, since a detector array can be used in conjunction with a microcomputer to make measurements simultaneously at a large number of points covering the interference pattern.

Figure 13 is a schematic of a compact optical system (the Mirau interferometer) used for phase-shifting interference microscopy. In this setup, the phase-steps are introduced by mounting the sample on a piezoelectric transducer (PZT) to which an appropriately varying voltage is applied. In a Fizeau interferometer, it is possible to use a laser diode as the light source and vary its output frequency.²⁹ If the initial optical path difference between the beams in the interferometer is p , a frequency shift $\Delta\nu$ in the output of the laser diode introduces an additional phase difference between the beams

$$\Delta\varphi = (2\pi p/\nu) \Delta\nu \quad (12)$$

Another way of shifting the phase of a beam of light is by a cyclic change in its state of polarization. Since the resulting phase shift (the Pancharatnam phase) is very nearly achromatic, measurements can be made with white light, so that phase ambiguities at steps are eliminated.³⁰

32.10 PHASE-LOCKED INTERFEROMETERS

The output intensity from an interferometer depends on the phase difference between the beams. In phase-locked interferometers, any variation in the output intensity is detected and fed back to a phase modulator in the measurement path so as to hold the output intensity constant. The changes in the optical path can then be estimated from the changes in the drive signal to the phase modulator.³¹

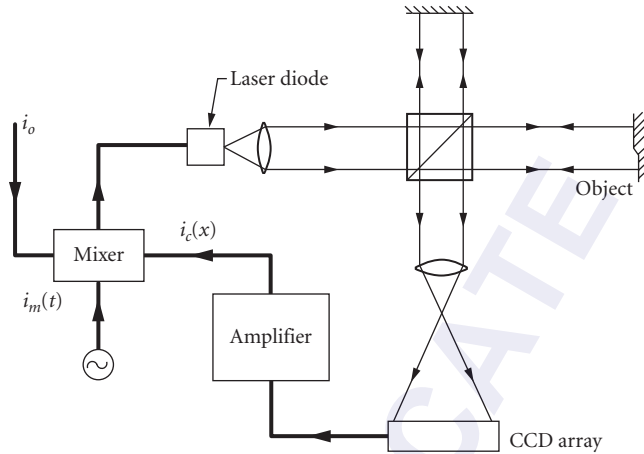


FIGURE 14 Schematic of a phase-locked interferometer using a laser diode source. (From Ref. 32.)

Drifts can be eliminated by using an ac amplifier. If the phase of one beam in the interferometer is subjected to a sinusoidal modulation

$$\Delta\phi(t) = \Delta\phi \sin \omega t \quad (13)$$

with an amplitude $\Delta\phi \ll \pi$, the output signal at the modulation frequency has an amplitude

$$I_{\omega}(t) = 4(I_1 I_2)^{1/2} J_1(\Delta\phi) \sin \phi \quad (14)$$

and drops to zero when $\phi = m\pi$, where m is an integer. Since, at this point, both the magnitude and the sign of this signal change, it can be used as the input to a servo system that locks the phase difference between the beams at this point.

With a laser diode, it is possible to compensate for changes in the optical path difference by a change in the illuminating wavelength. A typical setup is shown in Fig. 14. The injection current of the laser then consists of a dc bias current i_o , a control current i_c , and a sinusoidal modulation current $i_m(t) = i_m \cos \omega t$ whose amplitude is chosen to produce the required phase modulation.³²

Direct measurements of changes in the optical path are possible by sinusoidal phase-modulating interferometry, which uses a similar setup, except that in this case the amplitude of the phase modulation is much larger (typically around π radians). The modulation amplitude is determined from the amplitudes of the components in the detector output corresponding to the modulation frequency and its third harmonic. The average phase difference between the beams can then be determined from the amplitudes of the components at the modulation frequency and its second harmonic.³³

32.11 LASER-DOPPLER INTERFEROMETERS

Light scattered from a moving particle undergoes a frequency shift, due to the Doppler effect, that is proportional to the component of its velocity in a direction determined by the directions of illumination and viewing. With laser light, this frequency shift can be evaluated by measuring the frequency of the beats produced by the scattered light and a reference beam, or by the scattered light from two illuminating beams incident at different angles.^{34,35}

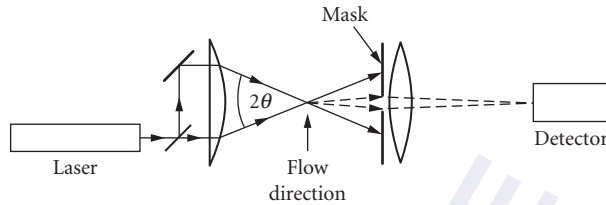


FIGURE 15 Optical arrangement used for laser-Doppler velocimetry.

Laser-Doppler interferometry can be used for measurements of the velocity of moving materials,³⁶ as well as for measurements, at a given point, of the instantaneous flow velocity of a moving fluid to which suitable tracer particles have been added.³⁷ A typical optical system for measurements on fluids is shown in Fig. 15. If the two illuminating beams in this arrangement make equal but opposite angles $\pm\theta$ with the viewing direction, the frequency of the beat signal is given by the relation

$$f = (2|v| \sin \theta) / \lambda \quad (15)$$

where v is the component of the velocity of the particle in the plane of the beams at right angles to the direction of observation. To distinguish between positive and negative flow directions, the frequency of one of the beams is offset by a known amount by means of an acousto-optical modulator. Simultaneous measurements of the velocity components along two orthogonal directions can be made by using two pairs of beams in orthogonal planes. Interactions between the two pairs of beams are avoided by using different laser wavelengths.

Laser diodes and optical fibers can be used to build very compact laser-Doppler interferometers.^{38,39} A frequency offset can be introduced between the beams either by using a piezoelectric fiber-stretcher driven by a sawtooth waveform in one path, or by ramping the injection current of the laser diode linearly.

Laser-Doppler interferometry can also be used to measure vibration amplitudes. Typically, one of the beams in an interferometer is reflected from a point on the vibrating specimen, while the other, whose frequency is offset, is reflected from a fixed reference mirror. The output from a photodetector then consists of a component at the offset frequency (the carrier) and two sidebands. The amplitude of the vibration can be determined from a comparison of the amplitudes of the carrier and the sidebands.⁴⁰ This technique can measure vibration amplitudes down to a few thousandths of a nanometer.⁴¹

32.12 LASER-FEEDBACK INTERFEROMETERS

Laser-feedback interferometers use the fact that the output of a laser is strongly affected if, as shown in Fig. 16, a fraction of the output beam is reflected back into the laser cavity by an external mirror M_3 . The output of the laser then varies cyclically with the position of M_3 , one cycle of modulation corresponding to a displacement of M_3 by half a wavelength.⁴²

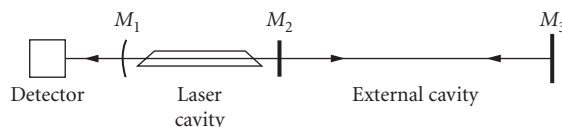


FIGURE 16 Schematic of a laser-feedback interferometer.

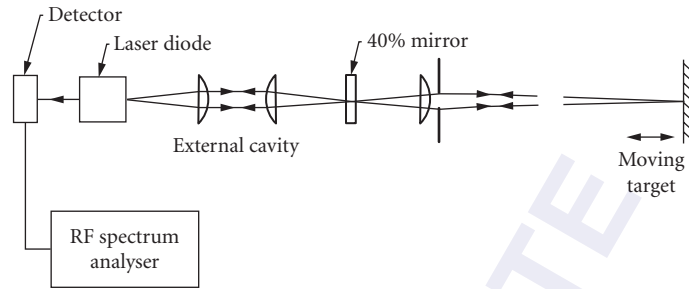


FIGURE 17 Feedback interferometer using a diode laser for velocimetry. (From Ref. 47.)

The operation of such an interferometer can be analyzed by considering the two mirrors M_3 and M_2 as a Fabry-Perot interferometer that replaces the output mirror of the laser. A variation in the spacing of M_3 and M_2 results in a variation in the reflectivity of this interferometer for the laser wavelength and, hence, in the gain of the laser.

A very compact laser-feedback interferometer can be set up with a single-mode laser diode.⁴³ Small displacements can be detected by measuring the changes in the laser output when the laser current is held constant. Measurements can be made over a larger range by mounting the laser on a piezoelectric transducer and using an active feedback loop to stabilize the length of the optical path from the laser to the mirror.⁴⁴

Laser-feedback interferometers can also be used for velocimetry. If the light reflected from the moving object is mixed with the original oscillating wave inside the laser cavity, the beat signal can be observed in the beam leaving the rear end of the laser.^{45,46} Very high sensitivity can be obtained with a laser diode operated near threshold.⁴⁷ If a separate external cavity is used, as shown in Fig. 17, to ensure single-mode operation, measurements can be made at distances up to 50 m.

32.13 FIBER INTERFEROMETERS

Analogues of conventional two-beam interferometers can be built with single-mode optical fibers. High sensitivity can be obtained with fiber interferometers because it is possible to have very long optical paths in a small space. In addition, because of the extremely low noise level, sophisticated detection techniques can be used.

Fiber-Interferometer Rotation Sensors

Fiber interferometers were first used for rotation sensing, by replacing the ring cavity in a conventional Sagnac interferometer with a closed, multiturn loop made of a single-mode fiber.⁴⁸ For a loop rotating with an angular velocity ω about an axis making an angle θ with the plane of the loop, the phase difference introduced between the two counterpropagating beams is

$$\Delta\phi = (4\pi\omega Lr \cos\theta)/\lambda c \quad (16)$$

where L is the length of the fiber, r is the radius of the loop, λ is the wavelength, and c is the speed of light. High sensitivity can be obtained by increasing the length of the fiber in the loop. In addition, very small phase shifts can be measured, and the sense of rotation determined, by introducing a nonreciprocal phase modulation in the beams and using a phase-sensitive detector.⁴⁹

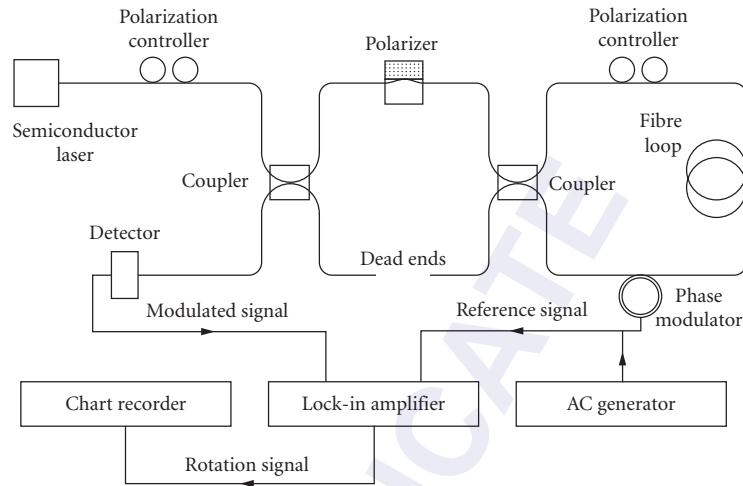


FIGURE 18 Fiber-interferometer for rotation sensing. (From Ref. 50.)

Figure 18 is a schematic of a typical all-fiber interferometric rotation sensor.⁵⁰ In this arrangement, the beam splitters are replaced by optical couplers, and a phase modulator consisting of a few turns of the fiber wound around a piezoelectric cylinder is located near one end of the optical fiber coil.

Fiber-interferometer rotation sensors have the advantages of small size and low cost. If care is taken to minimize noise due to back scattering and nonreciprocal effects due to fiber birefringence, performance close to the limit set by photon noise can be obtained.⁵¹

Generalized Fiber-Interferometer Sensors

The optical path length in a fiber is affected by its temperature and also changes when the fiber is stretched, or when the pressure changes. Accordingly, an optical fiber can be used in an interferometer to sense changes in these parameters.⁵²

Figure 19 is a schematic of an all-fiber interferometer that can be used for such measurements.⁵³ A layout analogous to a Mach-Zehnder interferometer avoids optical feedback to the laser. Optical

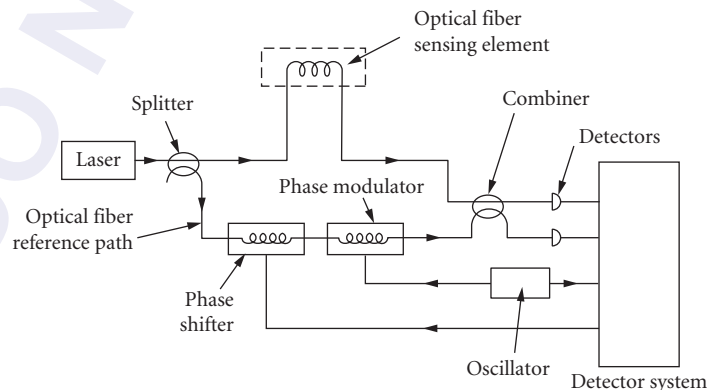


FIGURE 19 Schematic of a typical fiber-interferometer sensor. (From Ref. 53.)

fiber couplers are used to divide and recombine the beams, and measurements can be made with either a heterodyne system or a phase-tracking system. Detection schemes involving either laser-frequency switching or a modulated laser source can also be used. Optical phase shifts as small as 10^{-6} radian can be detected.

Fiber interferometers can also be used for measurements of magnetic or electric fields with a fiber sensor bonded to a magnetostrictive element⁵⁴ or jacketed with a piezoelectric polymer.⁵⁵ Phase ambiguities can be overcome by using a birefringent fiber⁵⁶ or by fiber-optic low-coherence interferometry, using a broad-band source.⁵⁷

Multiplexed Fiber-Interferometer Sensors

Fiber-interferometer sensors can be multiplexed to measure different quantities at different locations with a single light source and detector and the same set of transmission lines. Techniques developed for this purpose include frequency-division multiplexing, time-division multiplexing, and coherence multiplexing.⁵⁸⁻⁶²

32.14 INTERFEROMETRIC WAVE METERS

Tunable lasers have created a need for instruments that can measure their output wavelengths with an accuracy commensurate with their narrow line width. Dynamic wave meters have greater accuracy but can be used only with continuous wave (cw) sources; static wave meters can also be used with pulsed lasers.

Dynamic Wave Meters

A dynamic wave meter typically consists of a two-beam interferometer in which the number of fringes crossing the field is counted as the optical path is changed by a known amount. In one form, shown in Fig. 20, two beams, one from the laser whose wavelength is to be determined and another from a frequency stabilized He-Ne laser, traverse the same two paths in opposite directions.⁶³ The fringe systems formed by these two lasers are imaged on the two detectors D_1 and D_2 , respectively. If, then, the end reflector is moved through a distance d , we have

$$\lambda_1/\lambda_2 = N_2/N_1 \quad (17)$$

where N_1 and N_2 are the numbers of fringes seen by D_1 and D_2 , respectively, and λ_1 and λ_2 are the wavelengths in air. To obtain the highest precision, it is also necessary to measure the fractional order numbers. This can be done by phase-locking an oscillator to an exact multiple of the frequency

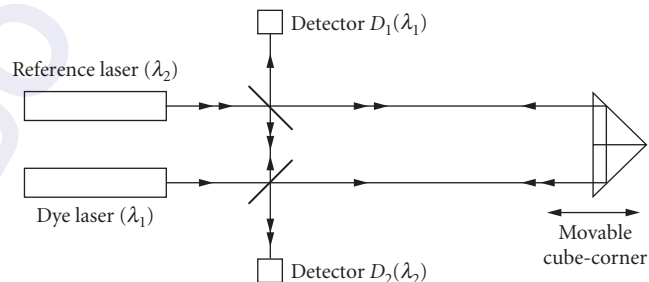


FIGURE 20 Optical system of a dynamic interferometric wave meter. (From Ref. 63.)

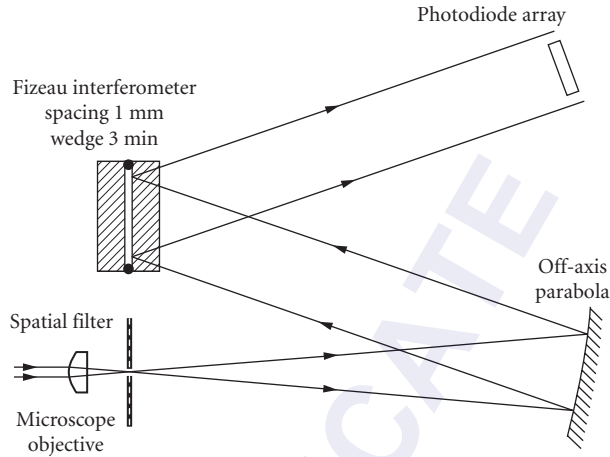


FIGURE 21 Schematic of a static interferometric wave meter. (From Ref. 67.)

of the ac signal from the reference channel, or by digitally averaging the two signal frequencies.⁶⁴ It is also possible to use a vernier method in which the counting cycle starts and stops when the phases of the two signals coincide.⁶⁵ With these techniques, a precision of 1 part in 10^9 can be obtained.

Another type of dynamic wave meter uses a scanning Fabry-Perot interferometer in which the separation of the mirrors is changed slowly. If this interferometer is illuminated with the two wavelengths to be compared, peak transmission will be obtained for both wavelengths at intervals given by the condition

$$m_1 \lambda_1 = m_2 \lambda_2 = p \quad (18)$$

where m_1 and m_2 are the changes in the integer order and p is the change in the optical path difference.⁶⁶ A precision of 1 part in 10^7 can be obtained with a range of movement of only 25 μm , because the Fabry-Perot fringes are much sharper than two-beam fringes.

Static Wave Meters

The simplest type of static wave meter is based on the Fizeau interferometer.⁶⁷ As shown in Fig. 21, a collimated beam from the laser is incident on two uncoated fused-silica flats separated by about 1 mm and making an angle of about 3 min of arc with each other. The intensity distribution in the fringe pattern formed in the region in which the shear between the two reflected beams is zero is recorded by a linear detector array.⁶⁸ In the first step, the integral interference order is calculated from the spatial period of the interference pattern; the exact value of the wavelength is then calculated from the positions of the maxima and minima.

32.15 SECOND-HARMONIC AND PHASE-CONJUGATE INTERFEROMETERS

Nonlinear optical elements are used in second-harmonic and phase-conjugate interferometers.⁶⁹

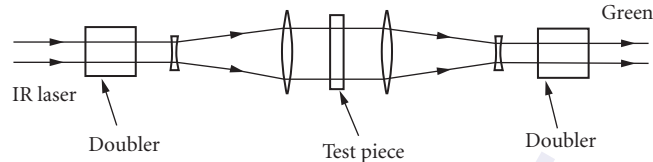


FIGURE 22 Second-harmonic interferometer: analog of the Mach-Zehnder interferometer. (From Ref. 70.)

Second-Harmonic Interferometers

One type of second-harmonic interferometer, shown in Fig. 22, is an analog of the Mach-Zehnder interferometer.⁷⁰ In this interferometer, the infrared beam from a Q-switched Nd:YAG laser ($\lambda_1 = 1.06 \mu\text{m}$) is incident on a frequency-doubling crystal. The green ($\lambda_2 = 0.53 \mu\text{m}$) and infrared beams emerging from this crystal traverse the test piece and are then incident on another frequency-doubling crystal.

The fringe number at any point in this interferometer is

$$N = (n_1 - n_2)d/\lambda_2 \quad (19)$$

where n_1 and n_2 are the refractive indices of the test specimen at 1.06 and $0.53 \mu\text{m}$, respectively, and d is its thickness.

Phase-Conjugate Interferometers

In a phase-conjugate interferometer, the wavefront that is being studied is made to interfere with its conjugate.⁷¹ Such an interferometer has the advantage that a reference wavefront is not required; in addition, the sensitivity of the interferometer is doubled.

Figure 23 is a schematic of a phase-conjugate interferometer that is an analog of the Fizeau interferometer.⁷² In this interferometer, the signal beam is incident on a conventional, partially reflecting mirror placed in front of a single crystal of barium titanate which functions as an internally self-pumped phase-conjugate mirror.

An interferometer in which both mirrors have been replaced by phase-conjugating mirrors is unaffected by misalignment of the mirrors and the field of view is normally completely dark. However, because of the delay in the response of the phase conjugator, dynamic changes in the optical path difference are displayed.^{73,74}

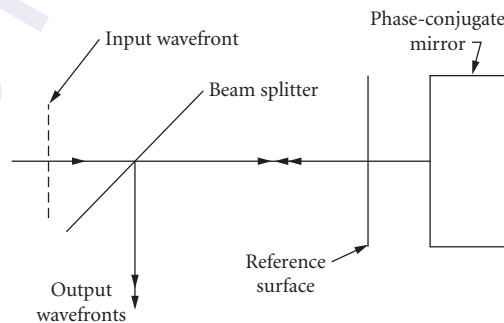


FIGURE 23 Schematic of a phase-conjugate Fizeau interferometer. (From Ref. 72.)

Interferometric Optical Switches

Nonlinear optical effects have also been exploited to develop high-speed interferometric optical switches.⁷⁵

32.16 STELLAR INTERFEROMETERS

A star can be considered as an incoherent light source whose dimensions are small compared to its distance from the earth. Accordingly, the complex degree of coherence between the fields at two points on the earth's surface is given by the normalized Fourier transform of the intensity distribution over the stellar disc.

Michelson's Stellar Interferometer

Michelson used the interferometer shown schematically in Fig. 24 to make observations of the visibility of the fringes formed by light from a star, for different separations of the mirrors. The separation at which the fringes disappeared was used to determine the angular diameter of the star. The problems encountered by Michelson in making measurements at mirror separations greater than 6 m have been overcome in new versions of this interferometer.⁷⁶

The Intensity Interferometer

The intensity interferometer⁷⁷ uses measurements of the correlation between the fluctuations in the intensity at two photodetectors separated by a suitable distance, which is proportional to the square of the modulus of the degree of coherence of the fields. Atmospheric turbulence only affects the phase of the incident waves and has no effect on the measured correlation. In addition, since the spectral bandwidth is limited by the electronics, it is only necessary to equalize the optical paths to within a few centimeters. It was therefore possible to use light collectors separated by distances up to 188 m.

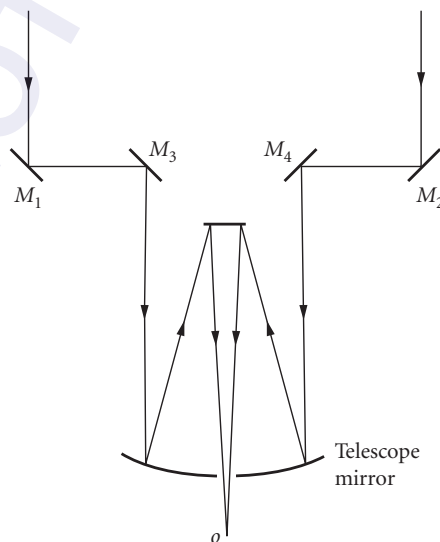


FIGURE 24 Michelson's stellar interferometer.

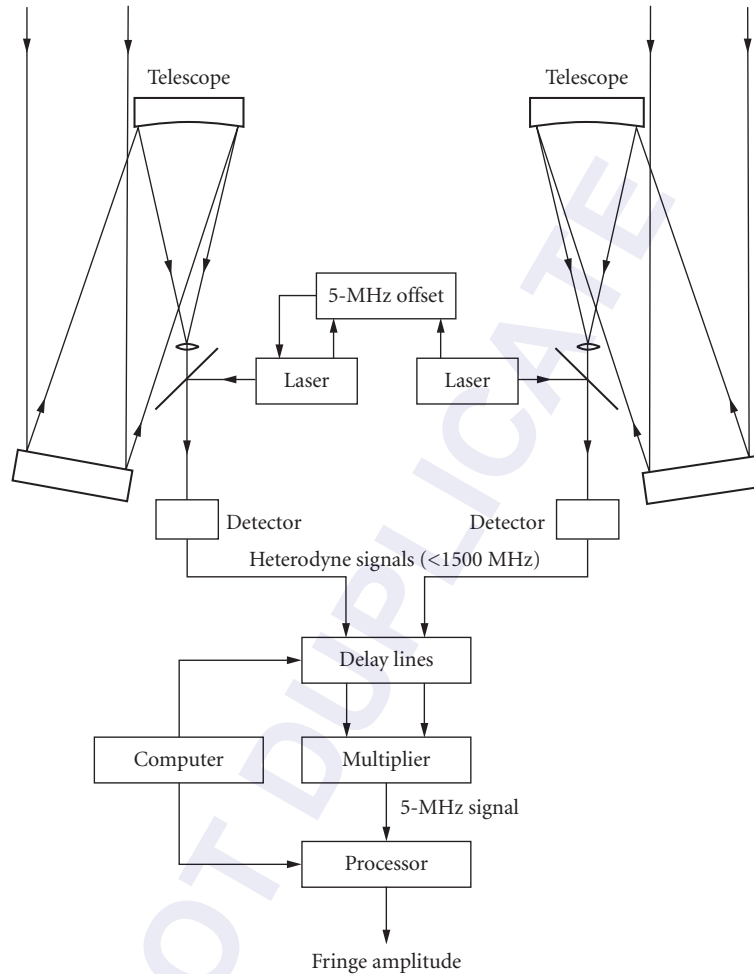


FIGURE 25 Schematic of an infrared heterodyne stellar interferometer. (From Ref. 78.)

Heterodyne Stellar Interferometers

In heterodyne stellar interferometers, as shown in Fig. 25, light from the star is mixed with light from two CO_2 lasers, whose frequencies are offset by 5 MHz with respect to each other, at two photodetectors, and the resulting heterodyne signals are multiplied in a correlator. The output signal from the correlator is a measure of the degree of coherence of the wave fields at the two photodetectors.^{78–80}

As in the intensity interferometer, it is only necessary to equalize the two paths to within a few centimeters. However, higher sensitivity is obtained, because the output is proportional to the product of the intensities of the laser and the star.

Nulling Interferometers and Interferometric Arrays

Problems arise when trying to detect a planet near a star. Nulling interferometers reduce the flux from the star, relative to its surroundings, by making the light from the star interfere with itself.⁸¹

Another advance is the application of multielement interferometric arrays to obtain high-resolution images of stellar objects.^{82,83}

32.17 GRAVITATIONAL-WAVE INTERFEROMETERS

Gravitational waves produced by cosmic sources, such as binary systems of neutron stars, collapsing supernovas and black holes, can be thought of as an alternating strain that propagates through space, affecting the dimensions and spacing of all material objects.

Since gravitational waves are transverse quadrupole waves, the effect of a gravitational wave on a Michelson interferometer would be a change in the difference of the lengths of the two arms.⁸⁴ However, to obtain the required sensitivity to strains, of the order of 1 part in 10^{21} , unrealistically long arms (>100 km) would be needed. In the LIGO project, higher sensitivity is obtained by using, as shown in Fig. 26, two identical Fabry-Perot interferometers ($d = 4$ km) at right angles to each other, with their mirrors mounted on freely suspended masses.⁸⁵ The separations of the mirrors are compared by locking the frequency of a laser to a transmission peak of one interferometer and using a servo system to adjust the length of the other interferometer continuously, so that its peak transmittance is also at the same frequency.

Even higher sensitivity is obtained by making use of the fact that, to avoid overloading the detector, the interferometer is normally adjusted so that observations are made on a dark fringe. Most of the light is then returned to the source. This light is recycled by using an extra mirror to reflect it back into the interferometer with the right phase.

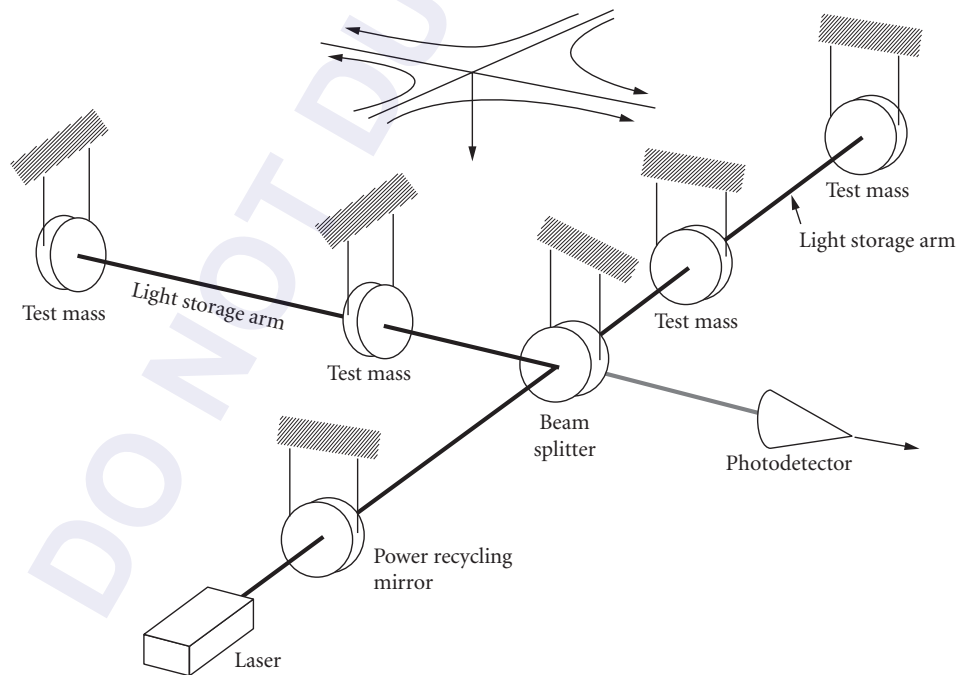


FIGURE 26 Gravitational-wave detector using two Fabry-Perot interferometers. (See also color insert.) (From Ref. 85.)

32.18 REFERENCES

1. W. H. Steel, *Interferometry*, Cambridge University Press, Cambridge, 1983.
2. P. Hariharan, *Optical Interferometry*, Academic Press, San Diego, 2003.
3. M. V. Mantravadi and D. Malacara, In: D. Malacara (ed.), *Optical Shop Testing*, John Wiley, Hoboken, 2007, pp. 1–45.
4. D. Malacara, In: D. Malacara (ed.), *Optical Shop Testing*, John Wiley, Hoboken, 2007, pp. 46–96.
5. M. Francon and S. Mallick, *Polarization Interferometers: Applications in Microscopy and Macroscopy*, Wiley-Interscience, London, 1971.
6. P. Hariharan, *Opt. Eng.* **14**:257–258 (1975).
7. M. Strojnik, G. Paez, and M. V. Mantravadi, In: D. Malacara (ed.), *Optical Shop Testing*, John Wiley, Hoboken, 2007, pp. 122–184.
8. D. Malacara, In: D. Malacara (ed.), *Optical Shop Testing*, John Wiley, Hoboken, 2007, pp. 185–218.
9. J. M. Vaughan, *The Fabry-Perot Interferometer*, Adam Hilger, Bristol, 1989.
10. P. Hariharan and D. Sen, *J. Opt. Soc. Am.* **51**:398–399 (1961).
11. J. R. Sandercock, *Opt. Commun.* **2**:73–76 (1970).
12. M. Hercher, *Appl. Opt.* **7**:951–966 (1968).
13. F. Zernike, *J. Opt. Soc. Am.* **40**:326–328 (1950).
14. P. Hariharan and D. Sen, *J. Sci. Instrum.* **36**:70–72 (1959).
15. P. Hariharan, D. Sen, and M. S. Bhalla, *J. Sci. Instrum.* **36**:72–75 (1959).
16. P. Hariharan and D. Sen, *J. Opt. Soc. Am.* **50**:357–361 (1960).
17. G. R. Hopkinson, *J. Opt. (Paris)* **9**:151–155 (1978).
18. J. N. Dukes and G. B. Gordon, *J. Hewlett-Packard*, **21**(12):2–8 (1970).
19. C. W. Gillard and N. E. Buholz, *Opt. Eng.* **22**:348–353 (1983).
20. H. Matsumoto, *Appl. Opt.* **25**: 493–498 (1986).
21. P. Hariharan, *Proc. SPIE* **1400**:2–10 (1991).
22. T. Kubota, M. Nara, and T. Yoshino, *Opt. Lett.* **12**:310–312 (1987).
23. R. Crane, *Appl. Opt.* **8**:538–542 (1969).
24. S. F. Jacobs and D. Shough, *Appl. Opt.* **20**:3461–3463 (1981).
25. S. F. Jacobs, D. Shough, and C. Connors, *Appl. Opt.* **23**:4237–4244 (1984).
26. N. Takahashi, S. Kakuma, and R. Ohba, *Opt. Eng.* **35**:802–807 (1996).
27. J. H. Bruning, D. R. Herriott, J. E. Gallagher, D. P. Rosenfeld, A. D. White, and D. J. Brangaccio, *Appl. Opt.* **13**:2693–2703 (1974).
28. K. Creath, In: E. Wolf (ed.), *Progress in Optics*, vol. XXVI, Elsevier, Amsterdam, 1988, pp. 349–393.
29. Y. Ishii, J. Chen, and K. Murata, *Opt. Lett.* **12**:233–235 (1987).
30. P. Hariharan, In: E. Wolf (ed.), *Progress in Optics*, vol. XLVIII, Elsevier, Amsterdam, 2005, pp. 149–201.
31. G. W. Johnson, D. C. Leiner, and D. T. Moore, *Proc. SPIE* **126**:152–160 (1977).
32. T. Suzuki, O. Sasaki, and T. Maruyama, *Appl. Opt.* **28**:4407–4410 (1989).
33. O. Sasaki and H. Okazaki, *Appl. Opt.* **25**:3137–3140 (1986).
34. Y. Yeh and H. Z. Cummins, *Appl. Phys. Lett.* **4**:176–178 (1964).
35. F. Durst and J. H. Whitelaw, *J. Phys. E: Sci. Instrum.* **4**:804–808 (1971).
36. B. E. Truax, F. C. Demarest, and G. E. Sommargren, *Appl. Opt.* **23**:67–73 (1984).
37. F. Durst, A. Melling, and J. H. Whitelaw, *Principles and Practice of Laser-Doppler Anemometry*, Academic Press, London, 1976.
38. O. Sasaki, T. Sato, T. Abe, T. Mizuguchi, and M. Niwayama, *Appl. Opt.* **19**:1306–1308 (1980).

39. J. D. C. Jones, M. Corke, A. D. Kersey, and D. A. Jackson, *Electron. Lett.* **18**:967–969 (1982).
40. W. Puschert, *Opt. Commun.* **10**:357–361 (1974).
41. P. Hariharan, in E. Wolf (ed.), *Progress in Optics*, vol. XXIV, Elsevier, Amsterdam, 1987, pp. 123–125.
42. D. E. T. F. Ashby and D. F. Jephcott, *Appl. Phys. Lett.* **3**:13–16 (1963).
43. A. Dandridge, R. O. Miles, and T. G. Giallorenzi, *Electron. Lett.* **16**:943–949 (1980).
44. T. Yoshino, M. Nara, S. Mnatzakanian, B. S. Lee, and T. C. Strand, *Appl. Opt.* **26**:892–897 (1987).
45. L. H. Churnside, *Appl. Opt.* **23**:61–66 (1984).
46. S. Shinohara, A. Mochizuki, H. Yoshida, and M. Sumio, *Appl. Opt.* **25**:1417–1419 (1986).
47. P. J. de Groot and G. M. Gallatin, *Opt. Lett.* **14**:165–167 (1989).
48. V. Vali and R. W. Shorthill, *Appl. Opt.* **15**:1099–1100 (1976).
49. S. Ezekiel, *Proc. SPIE.* **487**:13–20 (1984).
50. R. A. Bergh, H. C. Lefevre, and H. J. Shaw, *Opt. Lett.* **6**:502–504 (1981).
51. R. A. Bergh, H. C. Lefevre, and H. J. Shaw, *IEEE J. Lightwave Technol.* **LT-2**:91–107 (1984).
52. B. Culshaw, *Optical Fibre Sensing and Signal Processing*, Peregrinus, London, 1984.
53. T. G. Giallorenzi, J. A. Bucaro, A. Dandridge, G. H. Sigel Jr., J. H. Cole, S. C. Rashleigh, and R. G. Priest, *IEEE J. Quantum Electron.* **QE-18**:626–665 (1982).
54. J. P. Willson and R. E. Jones, *Opt. Lett.* **8**:333–335 (1983).
55. P. D. De Souza and M. D. Mermelstein, *Appl. Opt.* **21**:4214–4218 (1982).
56. P. A. Leilabady, J. D. C. Jones, M. Corke, and D. A. Jackson, *J. Phys. E: Sci. Instrum.* **19**:143–146 (1986).
57. Y.-J. Rao and D. A. Jackson, *Meas. Sci. Technol.* **7**:981–999 (1996).
58. I. P. Gilles, D. Uttam, B. Culshaw, and D. E. N. Davies, *Electron. Lett.* **19**:14–15 (1983).
59. J. L. Brooks, R. H. Wentworth, R. C. Youngquist, M. Tur, B. Y. Kim, and H. J. Shaw, *J. Lightwave Technol.* **LT-3**:1062–1071 (1985).
60. I. Sakai, G. Parry, and R. C. Youngquist, *Opt. Lett.* **11**:183–185 (1986).
61. J. L. Brooks, B. Moslehi, B. Y. Kim, and H. J. Shaw, *J. Lightwave Technol.* **LT-5**:1014–1023 (1987).
62. F. Farahi, J. D. C. Jones, and D. A. Jackson, *Electron. Lett.* **24**:409–410 (1988).
63. F. V. Kowalski, R. T. Hawkins, and A. L. Schawlow, *J. Opt. Soc. Am.* **66**:965–966 (1976).
64. J. L. Hall and S. A. Lee, *Appl. Phys. Lett.* **29**:367–369 (1976).
65. A. Kahane, M. S. O’Sullivan, N. M. Sanford, and B. P. Stoicheff, *Rev. Sci. Instrum.* **54**:1138–1142 (1983).
66. R. Salimbeni and R. V. Pole, *Opt. Lett.* **5**:39–41 (1980).
67. J. J. Snyder, *Proc. SPIE.* **288**:258–262 (1981).
68. J. L. Gardner, *Opt. Lett.* **8**:91–93 (1983).
69. P. Hariharan, In: E. Wolf (ed.), *Progress in Optics*, vol. XXIV, Elsevier, Amsterdam, 1987, pp. 144–151.
70. F. A. Hopf, A. Tomita, and G. Al-Jumaily, *Opt. Lett.* **5**:386–388 (1980).
71. F. A. Hopf, *J. Opt. Soc. Am.* **70**:1320–1322 (1980).
72. D. J. Gauthier, R. W. Boyd, R. K. Jungquist, J. B. Lisson, and L. L. Voci, *Opt. Lett.* **14**:323–325 (1989).
73. J. Feinberg, *Opt. Lett.* **8**:569–571 (1983).
74. D. Z. Anderson, D. M. Lininger, and J. Feinberg, *Opt. Lett.* **12**:123–125 (1987).
75. N. S. Patel, K. L. Hall, and K. A. Rauschenbach, *Appl. Opt.* **37**:2831–2842 (1998).
76. J. Davis, W. J. Tango, A. J. Booth, R. A. Minard, S. M. Owens, and R. R. Shobbrook, *Proc. SPIE.* **2200**:231–241 (1994).
77. R. Hanbury Brown, *The Intensity Interferometer*, Taylor and Francis, London, 1974.
78. M. A. Johnson, A. L. Betz, and C. H. Townes, *Phys. Rev. Lett.* **33**:1617–1620 (1974).
79. C. H. Townes, *J. Astrophys. Astron.* **5**:111–130 (1984).

80. P. D. S. Hale, M. Bester, W. C. Danchi, W. Fitelson, S. Hoss, E. A. Lipman, J. D. Monnier, P. G. Tuthill, and C. H. Townes, *Astrophys. J.* **537**:998–1012 (2000).
81. E. Serabyn and M. M. Colavita, *Appl. Opt.* **40**:1668–1671 (2001).
82. J. E. Baldwin, C. A. Haniff, C. D. Mackay, and J. P. Warner, *Nature* **320**:595–597 (1986).
83. W. Traub, (ed.), *Interferometry for Optical Astronomy II*, *Proc. SPIE* **4838**, SPIE, Bellingham, 2002.
84. G. E. Moss, L. R. Miller, and R. L. Forward, *Appl. Opt.* **10**:2495–2498 (1971).
85. R. Weiss, *Rev. Mod. Phys.* **71**:S187–S196 (1999).

DO NOT DUPLICATE

HOLOGRAPHY AND HOLOGRAPHIC INSTRUMENTS

Lloyd Huff

*Research Institute
University of Dayton
Dayton, Ohio*

33.1 GLOSSARY

A	wave amplitude
a	diameter of viewing lens in speckle imaging system
d	fringe spacing
d_{sp}	characteristic speckle diameter
\mathbf{E}	electric field vector
$\hat{\mathbf{e}}$	polarization unit vector
f	wave frequency
I	field irradiance
I_H	irradiance of the field in plane of hologram
K	proportionality constant
k	propagation constant
r	radial position coordinate
T	transmittance of the hologram
v	distance from lens to image plane in speckle imaging system
θ_1, θ_2	object and reference beam angles
λ	wavelength
ϕ	wave phase
Ψ	complex field amplitude
Ψ_H	complex field amplitude in plane of hologram
Ψ_O	complex amplitude of object wave field
Ψ_R	complex amplitude of reference wave field
Ψ_T	complex amplitude of field transmitted by hologram

33.2 INTRODUCTION

The three-dimensional imagery produced by holography accounts for much of the popular interest in this technique. Conceptual applications, such as the holodeck seen on the television series *Star Trek: The Next Generation*, and actual applications, like the widespread use of embossed holograms on book and magazine covers, gift wrapping, product packaging, and credit cards, have fascinated and captured the imagination of millions. Holography was discovered in 1947 by Gabor and revived in the early 1960s through the work of Leith and Upatnieks. Since that time, most practitioners in the field believe that technical applications, rather than imaging, have represented the utility of holography in a more significant way. This chapter is a brief overview of some of the more important technical applications, particularly as they relate to a variety of instrumentation problems. The discussion addresses several ways holography has been used to observe, detect, inspect, measure, or record numerous physical phenomena. The second section presents a brief review of the basic principles of wavefront reconstruction. The third section (Sec. 33.4) addresses one of the more important applications of holography—holographic interferometry. Included in this section is a review of electronic or television holography which takes this powerful interferometric technique into the real-time domain. Section 33.5 addresses several instrumental applications of holographic optical elements (HOEs). Sections 33.6 and 33.7 discuss ways in which holography has been applied in the semiconductor industry. Section 33.8 briefly addresses the holographic storage of information.

33.3 BACKGROUND AND BASIC PRINCIPLES

Holography is a method of recording and reconstructing wavefronts residing anywhere in the electromagnetic spectrum or acoustic spectrum. This chapter addresses optical holography as practiced in or near the visible region of the electromagnetic spectrum. The principals of wavefront reconstruction were discovered by Gabor¹⁻³ in an attempt to improve the resolving power of the electron microscope. The original purpose was never accomplished, but this basic discovery evolved into one of the most significant new fields of study in the twentieth century. Gabor's early work received little attention because the lack of a light source with sufficient coherence severely limited the quality of the images produced. However, the invention of the laser in the early 1960s heralded a holographic renaissance. During this period, Leith and Upatnieks^{4,5} recognized the parallels between their work in coherent radar and Gabor's wavefront reconstruction concepts. Their experiments in the optical region of the spectrum with the newly available HeNe laser produced the first high-quality, three-dimensional images. The publication of this work created an explosive interest in the field as well as many unrealistic predictions about what might be accomplished with holographic three-dimensional imagery. The work of numerous researchers established the medium's true capabilities and limitations; consequently, many successful applications ensued. Progress continues to be made in the development of new materials and techniques sustaining a high level of interest in holography and its technical, commercial, and artistic applications.

Holography is most often associated with its ability to produce striking three-dimensional images. Therefore, a logical place to start in understanding holography is to compare this imaging science with its two-dimensional predecessor—photography. A light wavefront is characterized by several parameters; the two most important of these are its intensity (or irradiance) and its local direction of propagation. Photography records only one of these parameters—intensity—in the plane of the recording medium or photographic film. The intensity distribution of a light wave emanating from an object may be recorded by simply exposing a film plate placed in proximity to the object; however, this will not produce a discernible image. Recording a photograph is accomplished by imaging the object onto the film with a lens, thereby establishing a correspondence between points on the object and points in the film plane.

Holography also records the intensity distribution of a wavefront; in addition, the local propagation direction (or phase) is recorded through the process of optical interference. The process in its simplest form is illustrated in Fig. 1. The light from a laser is split into two parts, expanded with a short-focal-length lens (usually a microscope objective), and spatially filtered with a pinhole to

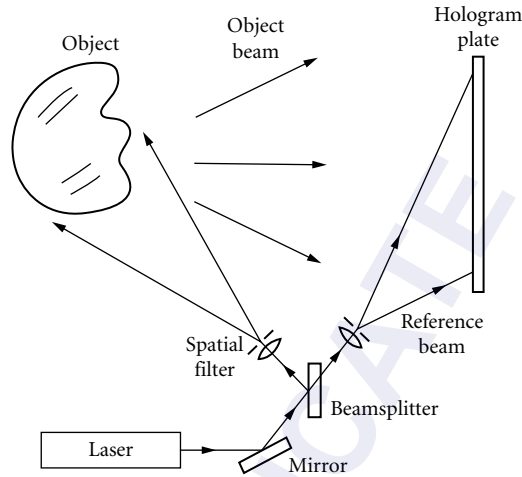


FIGURE 1 Typical optical arrangement for making a simple laser transmission hologram.

remove intensity variations caused primarily by nonuniformities in the lens. One of the split beams (object beam) is directed to the object; the other (reference beam) is incident directly on the recording medium (such as high-resolution silver halide film). The light reflected from and scattered by the object combines with the reference beam at the plate to form an interference fringe field. These fringes are recorded by the film. The spacing of these fringes d is given by the grating equation

$$d = \frac{\lambda}{(\sin\theta_1 + \sin\theta_2)} \quad (1)$$

where λ is the wavelength of the light, and θ_1 and θ_2 are the angles made by the object and reference beams relative to the normal. For visible light and common recording geometries, the fringe frequency ($1/d$) can exceed 2000 fringes (or line-pairs) per millimeter. Therefore, the recording material must be of very high resolution relative to conventional photographic film which is usually in the range of 50 line-pairs per millimeter. The stability of the fringe is extremely sensitive to the mechanical motion of the object and optical components. To record holograms with good fringe stability, the optical system must be stable enough to prevent motions greater than a fraction of a wavelength. For this reason, the common practice is to use rigid optical components placed on a stable, vibration-isolated table.

Illumination of the developed hologram by the reference beam alone reveals a three-dimensional image which is essentially identical to the original object as viewed in laser light. Observing the holographic image of the object is exactly like looking at the object through the window formed by the plate with full parallax and look-around capability. The object wave is reconstructed when the illumination (reference) wave is diffracted by the grating formed in the recording medium. This grating is formed by variation of the optical transmittance or optical thickness of the material along the fringe lines. The amplitude hologram formed with silver halide film may be converted to a phase hologram by bleaching; this results in a significant increase in diffraction efficiency. Other materials (such as photopolymer film) produce phase holograms directly.

The holographic recording and reconstruction process may be described in general mathematical terms as follows. The object and reference fields satisfy the Helmholtz equation

$$\nabla^2 \mathbf{E} + k^2 \mathbf{E} = 0 \quad (2)$$

where \mathbf{E} is the electric field vector and $k = 2\pi/\lambda$ is the propagation constant.

A spherical wave solution of this equation may be expressed in the form

$$\mathbf{E} = A e^{ikr} \mathbf{e}^{i2\pi ft} \hat{\mathbf{e}} \quad (3)$$

where A is the amplitude of the wave, f is the frequency, and $\hat{\mathbf{e}}$ is the polarization unit vector.

The complex field amplitude Ψ is defined as

$$\Psi = A e^{i\phi} \quad (4)$$

where $\phi = kr$ is the phase. The irradiance of the field is given by

$$\begin{aligned} I &= \mathbf{E} \cdot \mathbf{E}^* = \Psi \Psi^* \hat{\mathbf{e}} \cdot \hat{\mathbf{e}} = \Psi \Psi^* \\ &= A^2 = |\Psi|^2 \end{aligned} \quad (5)$$

where $*$ denotes the complex conjugate.

The field in the plane of the hologram Ψ_H is the sum of the object and reference fields:

$$\Psi_H = \Psi_O + \Psi_R \quad (6)$$

and the irradiance of the field at the hologram is given by (assuming parallel polarization of the waves):

$$\begin{aligned} I_H &= \Psi_H \Psi_H^* = (\Psi_O + \Psi_R)(\Psi_O + \Psi_R)^* \\ &= |\Psi_O|^2 + |\Psi_R|^2 + \Psi_O \Psi_R^* + \Psi_R \Psi_O^* \end{aligned} \quad (7)$$

After processing, and apart from a constant term, the transmittance T of the hologram is proportional to the irradiance of the field at the hologram;

$$T = K[|\Psi_O|^2 + |\Psi_R|^2 + \Psi_O \Psi_R^* + \Psi_R \Psi_O^*] \quad (8)$$

When illuminated by the reference wave, the field transmitted by the hologram Ψ_T is given by the hologram transmittance multiplied by the reference wave field:

$$\Psi_T = K[\Psi_R(|\Psi_O|^2 + |\Psi_R|^2) + |\Psi_R|^2 \Psi_O + \Psi_R^2 \Psi_O^*] \quad (9)$$

The first term in this equation for Ψ_T is simply the transmitted wave altered by an attenuation factor. The second term is the original object wave multiplied by an amplitude factor; this term represents the virtual holographic image of the object. The third term is the conjugate object wave. In off-axis holography, the real image formed by this wave is weak, lies out of the field of view, and does not make a significant contribution to the imaging process. However, for Gabor's original in-line holography, this term represented an objectionable twin image which overlapped and obstructed viewing of the desired image. An important contribution of the Leith and Upatnieks off-axis reference scheme was the elimination of this twin image.

Many different types of holograms can be made by varying the location of the object relative to the recording medium, the directions and relative angles of the object and reference beams, and the wavefront curvature of these beams. The properties of these hologram types vary greatly; much research has been performed to characterize and successfully apply the different formats. Vigorous work is still being pursued in both areas of imaging and technical applications. For a thorough explanation of holography and its many applications, the reader may consult any of several standard texts on the subject (e.g., Refs. 6–11).

33.4 HOLOGRAPHIC INTERFEROMETRY

Interferometry provides a means of measuring optical path differences through the analysis of fringe patterns formed by the interference of coherent light waves. Optical path differences of interest may be produced by mechanical displacements, variations in the contour of one surface relative

to another, and variations in the refractive index of a material volume. Classical interferometry involves the interference of two relatively simple optical wavefronts which are formed and directed by optical components. These components must be of sufficient quality that they do not introduce random phase variations across the field that compete with or totally mask the optical path length differences of interest. Typical examples of classical interferometry include the use of configurations such as the Michelson or Twyman Green and Mach-Zehnder interferometers to determine the surface figure of optical components, study the refractive index variation in optical materials, and visualize the properties of flowing gases. The need for high-quality optical surfaces in classical interferometry is a consequence of the difficulty, using classical optical methods, of generating two separate but identical optical wavefronts of arbitrary shape. Although it was not immediately recognized by early holography researchers, the ability to holographically record then replay an arbitrary wavefront in a predictable fashion obviated this basic limitation of classical interferometry. With holographic interferometry, polished optical surfaces are not required and diffusely reflecting objects of any shape may be studied.

In holographic interferometry, a wavefront is stored in the hologram and later compared interferometrically with another wavefront. Phase differences between these two wavefronts produce fringes that can be analyzed to yield a wide range of both qualitative and quantitative information about the system originating these two wavefronts. Several researchers working independently made experimental observations related to this fact.¹²⁻¹⁷ Once the full implication of this discovery was realized, a period of intense research activity began to develop a solid theoretical understanding of this powerful new technique. Holographic interferometry quickly became the most important application of the relatively young science of holography. Although other branches of holography have successfully matured, most notably HOEs, holographic interferometry remains today the area in which holography has probably made the greatest impact.

As stated earlier, holographic interferometry involves the interferometric comparison of two wavefronts separated in time. This comparison can be made in a variety of ways which constitute the basic methods of holographic interferometry: real-time, double-exposure, and time-average. Real-time interferometry is realized by the interference of a holographically reconstructed object wave with the wave emanating from the actual object. This is accomplished as follows. The holographic plate is exposed, developed, then replaced in its holder in its original position. Reference and object-beam intensities are adjusted so that the illuminated object and its holographic image are of approximately equal brightness. Since the reconstructed object wavefront is 180° out of phase with the object wavefront, the object should be dark when viewed through the holographic plate. In practice, one or two broad fringes usually appear across the object due to emulsion shrinkage effects and lack of complete mechanical precision in returning the holographic plate to its original position. Any disturbance of the object which results in a mechanical displacement of its surface will now produce a fringe system which can be viewed in real time. The structure and periodicity of the fringes are related to the surface displacement. The mechanical surface deformation can result from an applied force, change in pressure, change in temperature, or any combination of the three. The quantitative details of this surface deformation can be derived from an analysis of this fringe system.

In double-exposure holographic interferometry, the two wavefronts to be compared are stored in the same hologram. This is done by holographically recording the image of the object under study in two exposures separated in time in the same holographic plate. If nothing is done to alter the object wavefront between these two exposures, the resulting image will appear as for a single-exposure hologram. However, if the object is perturbed in some way between these two exposures, an interference fringe system will appear in the final image. Again, this fringe system is related to the mechanical deformations of the object surface caused by the disturbance. Real-time holographic interferometry allows one to study the effects of object perturbation of varying types and degrees over any desired length of time and in real time. In contrast, double-exposure holographic interferometry examines a particular change of the state of the object between two particular points in time. A double-exposure holographic interferogram then might be thought of as a single data point record. The interferogram might be a record of the change of the object from one stable state to another, which might be recorded with a continuous wave laser. Or the interest might be to compare two states of a rapidly varying system most effectively recorded using a pulsed laser.

In some respects, the double-exposure holographic interferogram involves less experimental complexity, because both exposures are made in a single hologram plate held in a fixed position. Mechanical registration of the plate to a baseline position is not required. Emulsion shrinkage due to wet process development of silver halide film affects the holographic fringe systems for both exposures in the same way; therefore, it is not a problem for double-exposure interferometry. Another feature which may or may not be of benefit, depending on the parameters of the experiment, is that the superimposed images for the two exposures are in phase, thereby producing a bright baseline image. Double-exposure interferometry has been used very effectively to record fast events such as the flight of a bullet passing through a chamber.¹⁸ The first exposure is made of the chamber alone before passage of the bullet; the second exposure is made with the bullet in midflight. The interferogram is an interference recording of the refractive index variations in the chamber created by the disturbance of the bullet.

Vibrating surfaces may be studied using either real-time or time-average holographic interferometry.^{16,19} A vibrating object presents a continuum of surface configurations or surface deformations to an interferometric system. A unique interferometric fringe system is associated with each state of the surface for any particular point in time during the vibrating cycle. In real-time interferometry, the interference fringes are formed by the addition of wavefronts from the surface of the object at rest and from the vibrating surface at some point during its vibrating cycle. The fringe pattern observed is a visual time average of this continuum set of interferogram fringe systems. A time-average holographic interferogram is made by exposing the holographic plate while the object is vibrating; the exposure time is usually multiple vibration periods. This hologram may be thought of as a continuum set of exposures, each recording the interference of the object wave with its temporal counterparts in the rest of the continuum over one cycle. The end result is a fringe pattern, directly related to the surface vibration pattern, in which the fringe lines represent contours of constant vibration amplitude. A more physical view of the process derives from observing that holographic fringe movement in the recording medium due to object movement nullifies the holographic recording. Thus, regions of a vibrating surface in motion (antinodes) will appear dark while regions at rest (nodes) will appear bright.

Stroboscopic illumination has long been used to study objects in motion. Coupling this technique with holographic interferometry produces interferograms of vibrating objects with enhanced fringe visibility and greater information content. The technique may be used to make real-time observations or to record double-exposure holograms of vibrating surfaces. In the real-time configuration, the hologram of the test object at rest is made in the usual fashion. The object is then vibrated and strobed with a laser set to flash at a particular point in the vibrating cycle. The fringe system formed is produced by interference of the wavefront from the object at rest and the wavefront from the object at this particular point in the cycle. The timing of the laser flash may be varied to observe the evolution of the surface vibration throughout its entire cycle. A double-exposure interferogram is formed by exposing the hologram plate to two flashes from the strobed laser. In this manner, any two states of the vibrating surface may be compared during its cycle by varying the timing of the laser flashes and their separation. The actual exposure may extend over several vibration cycles. This method results in fringes of much higher contrast than those yielded by time-average holography.

An interesting variant of hologram interferometry is contour generation.²⁰⁻²² This is accomplished by making two exposures of an object with a refractive index change in the medium surrounding the object or a change in the laser wavelength between exposures. Either method yields fringes on the surface of the object. The fringe positions are related to the height of points on the object relative to a fixed plane. The two-wavelength method may be implemented in real-time by first making a hologram of an object in the usual manner at one wavelength, then illuminating both the object and the processed hologram (carefully placed in its original position) at a different wavelength. Both methods have been successfully applied in a variety of situations.

In many cases, the interpretation of the fringe pattern produced by a holographic interferogram is a simple matter of qualitative assessment. For example, defects or flaws may be identified by anomalous local variations in a background fringe pattern. The presence or absence of these

anomalies may provide all the information required in a nondestructive evaluation experiment. However, a detailed quantitative assessment of the mechanical deformation of the surface may be desired—this can involve a complex mathematical analysis of the fringe system. Quantitative analysis of the fringe pattern is often complicated by the fact that the fringes are not necessarily localized on the surface of the object. The interpretation of holographic interferograms and analysis of the fringe data have been the subject of considerable study.^{23–27} Numerous methods such as sandwich holographic interferometry,^{28,29} fringe linearization interferometry,³⁰ difference holographic interferometry,³¹ and fringe carrier techniques³² have been developed to facilitate this interpretation. In addition, the development of automatic fringe reading systems and data reduction software have greatly aided this process. However, fringe data reduction remains a challenge for many situations despite the progress that has been made.

The basic methods of holographic interferometry (real-time, double-exposure, and time-average) are in widespread use and continue to be the mainstay of this technique. However, important refinements have been made which have greatly added to the power of holographic interferometry. These advances include the use of real-time recording media^{33,34} and heterodyne holographic interferometry.³⁵ Real-time holographic recording materials (such as photorefractive crystals) provide an adaptive feature that makes the interferometer less sensitive to vibration, air currents, and other instabilities. The reliability of the interferometric process in a hostile environment is thus improved. Heterodyne techniques using two separate reference waves and a frequency shift between these two waves upon reconstruction has greatly improved the accuracy of holographic interferometry. Measurements with accuracies as high as $\lambda/1000$ can be made using heterodyne methods. Holographic moiré,³⁶ infrared holographic interferometry,³⁷ and the use of optical fibers³⁸ have also significantly extended the capabilities of holographic interferometry.

In the laboratory, where conditions are well-controlled, silver halide film has been the recording material of choice for holographic interferometry due to its relatively high sensitivity, low cost, and reliability. However, in field applications such as the factory floor, the wet processing requirements of silver halide film make this material much less attractive and, for some time, inhibited the use of holographic interferometry in many situations. Other materials that do not require wet processing (such as photopolymer films) are available, but these materials have very low sensitivity. The development of the thermoplastic recording material,³⁹ which does not require wet processing but retains the high sensitivity of silver halide film, made possible the much more convenient application of holographic interferometry in industrial situations. Several companies have commercially marketed holocamera systems using this material. A holographic interferogram made using one of these holocamera systems is shown in Fig. 2.

Holographic interferometry has been applied to an enormous range of problems; this is a simple testimony to its utility and versatility. The classical interferometric testing of the figure of optical components during fabrication can be augmented with holographic interferometry to test for figure during the grinding process since the surface of the test object does not need to be polished.²¹ The ability of holographic interferometry to make precise measurements of very small mechanical displacements has enabled it to be used in stress-strain measurements in materials, components, and systems. Mechanical displacements observed with holographic interferometry are often the result of thermal disturbances. Measurement of these thermally induced mechanical displacements with holographic interferometry can provide an accurate determination of the heat transfer properties of the material or system under study.⁴⁰ Similarly, diffusion coefficients in liquids can be determined using holographic interferometry.^{41,42} Flow visualization and the accurate determination of fluid-flow properties using holographic interferometry has been an intense area of study.^{43–45} As noted earlier, the technique can also be used to study high-speed events using short-pulse lasers in the double-pulse mode. Vibration analysis is one of the more powerful applications of holographic interferometry. In this area, the technique has been applied to a diverse array of problems including studies of the vibration properties of musical instruments,⁴⁶ vibration patterns in the human eardrum,⁴⁷ and vibration properties of mechanical parts such as turbine blades.⁴⁸ The application of holographic interferometry to turbine blade mechanics is illustrated in Fig. 3. One of the great virtues of holographic interferometry is that a tremendous wealth of

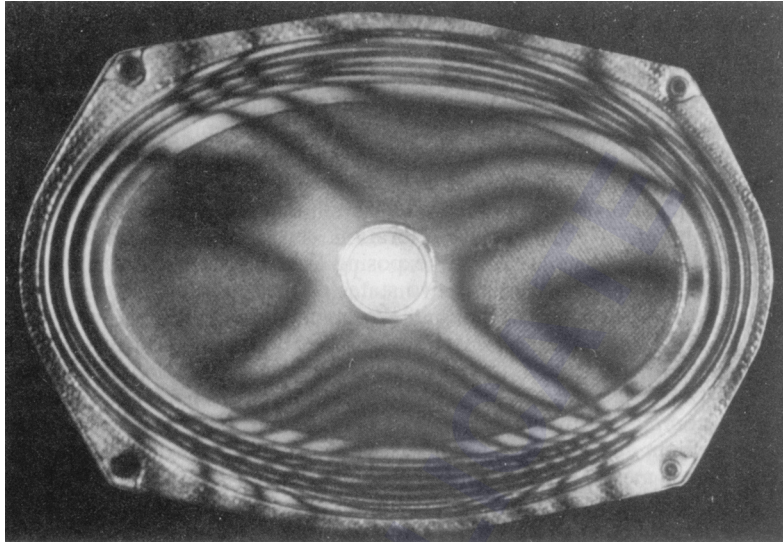


FIGURE 2 Time-average holographic interferogram of a loudspeaker vibrating at resonance. (Photo courtesy of Newport Corporation.)

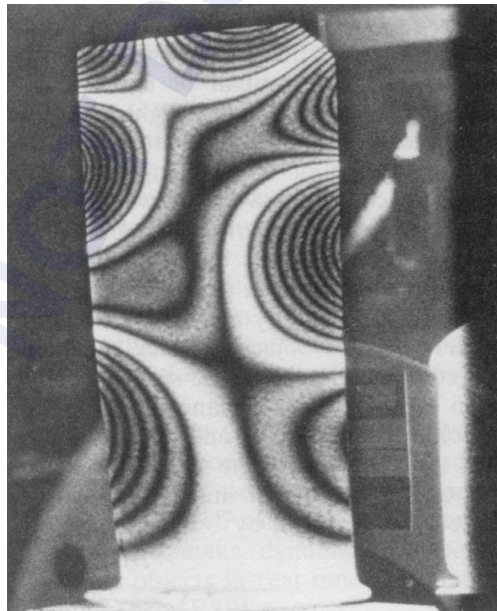


FIGURE 3 Time-average holographic interferogram displaying one of the vibration modes of a turbine blade. (Photo courtesy of Karl Stetson, United Technologies Research Center.)

information can be garnered from its application without destroying the test object or system. As a result, nondestructive evaluation or nondestructive testing has been one of the most important areas of application for this technique. As an example, holographic interferometry has been successfully used to observe subsurface defects in solid opaque objects. Even though the interference pattern is produced strictly by mechanical surface deformations, these surface variations are often indicative of subsurface changes (e.g., ply separations in automobile tires and interlayer delamination in composite materials).⁴⁹ Subsurface defects are usually manifest in local anomalies of the fringe pattern and a qualitative examination of the interferogram will often discern the effect. The literature is replete with articles describing these and many other applications of holographic interferometry. For the reader interested in an in-depth discussion of the theory and application of holographic interferometry, numerous textbooks and review articles are available (see, for example, Refs. 6–8 and 50–57).

Electronic Holography

Even with the use of thermoplastic recording media, holographic interferometry remains a challenge and, in many cases, unacceptable technique for industrial applications—particularly those that involve on-line quality testing in a production environment. Speckle-pattern interferometry,⁵⁸ another technique closely associated with holographic interferometry, alleviates many shortcomings of the traditional holographic approach when combined with electronic image recording and processing equipment.

Speckle is the coarse granular or mottled intensity pattern observed when a diffuse surface is illuminated with coherent light. Wavelets reflected from the randomly oriented facets of an optically rough surface interfere to produce this effect when the size of the facets is on the order of a wavelength or larger. Although this interference occurs throughout the space occupied by the wave scattered by the surface, the interference that produces the observable pattern takes place in the plane of the detector or recording medium (i.e., the retina of the eye or the film plane of a camera). The speckle-pattern recorded by an imaging system (eye or camera) is known as *subjective speckle*, while the intensity variation detected by a scanning detector above a coherently illuminated diffuse surface is referred to as *objective speckle*. Objective speckle is the resultant sum of the waves scattered from all parts of the surface to a point in space; in subjective speckle, wave summation in the observation plane is limited to the resolution cell of the system. The objective-speckle scale depends only on the plane in space where it is viewed, not on the image system used to view it. The size of the image plane or subjective speckle depends on the aperture of the viewing or imaging system. For subjective speckle, the characteristic speckle diameter, d_{sp} is given by⁵⁹

$$d_{sp} \approx \frac{2.4\lambda v}{a} \quad (10)$$

where λ is the wavelength, v is the distance from the lens to the image plane, and a is the diameter of the viewing lens aperture.

A fringe pattern is formed when the speckle patterns of the diffuse surface in its original and displaced positions are properly combined. The formation of this fringe pattern is known as speckle-pattern interferometry, of which there are two basic types: speckle-pattern photography and speckle-pattern correlation interferometry. Both techniques, which form the basis for electronic speckle-pattern interferometry (ESPI), and other forms of electronic holography will be discussed in this section. The remarks made here are derived from Ref. 58, which contains a thorough discussion of the topic.

By varying the recording and viewing configurations, speckle-pattern interference fringes can be made sensitive to in-plane and out-of-plane displacements, displacement gradients, and the first derivative of displacement gradients. Speckle-pattern interferometry has two distinct advantages over holographic interferometry: (1) the direction of the magnitude sensitivity of the fringes can be varied over a larger range, and (2) the resolution of the recording medium

for speckle-pattern interferometry does not need to be nearly as high. Therefore, speckle-pattern interferometry is a much more flexible technique, although the fringe definition is not nearly as good as with holographic interferometry due to the degradation of the images by the speckle pattern.

In speckle-pattern photography, the object is illuminated by a single light beam; no reference beam is involved. Some of the light scattered by the object is collected by a lens and recorded on photographic film. The film plane may be an image plane (in-focus) or some other plane (out-of-focus). The location of the film plane determines whether the resulting interferometric fringes are sensitive to in-plane or out-of-plane motion. Two exposures of the film are made: one with the object in its original position, the second with the object deformed or displaced. Proper illumination of the film negative with coherent light produces a fringe pattern in the observation plane which is related to the object motion. With the use of appropriate recording and viewing geometries, the fringes may be made to superimpose an image of the object. If the object is illuminated by a plane wavefront and the film is in the focal plane of a lens, the fringes are related to out-of-plane displacement gradients. Illumination of the object by a diverging wavefront in the proper geometry yields fringes related to the tilt of the object. Speckle-pattern photography can be used to make time-average stroboscopic and double-pulse measurements just as in holographic interferometry. In speckle-pattern correlation interferometry, a reference beam (either specular or diffuse) is incident upon the observation or recording plane in addition to the light scattered by the object. Interferometric fringes are produced by the correlation of the speckle patterns in the observation plane for the displaced and undisplaced object. Real-time or live-correlation fringes may be produced as follows. The object and reference beams are recorded with the object in its original position using photographic film. The film is developed and the film negative is replaced in its original position. The negative is illuminated with object and reference light, and the object is displaced. Correlation fringes are produced by the process of intensity multiplication. Because of the contrast reversal of the film negative, minimum transmission is found in areas of maximum correlation between the pattern recorded and the pattern produced by the displaced object. Unfortunately, the correlation fringes produced using this technique are of low contrast.

The variation in the correlation of the two speckle patterns which produces the fringes may be made sensitive to different components of surface displacement by using different object and reference beam geometries. One of the most important configurations uses a specular in-line reference beam introduced with a beam splitter or mirror with a pin hole. This configuration may be used to make dynamic displacement measurements or to observe the behavior of vibrating objects in real time. This particular optical geometry is also the most popular arrangement for ESPI.

In ESPI, the photographic film processing methods used for speckle-pattern photography and speckle-pattern correlation interferometry are replaced by video recording and display technology. The concept of using video equipment for this purpose was originated by several researchers working independently during the same period.⁶⁰⁻⁶³ For speckle-pattern interferometry, the minimum speckle size is usually in the range of 5 to 100 μm so that standard television (TV) cameras can be used to record the pattern. The main advantage of using TV equipment is the high data rate. Real-time correlation fringes may be produced and displayed on the TV monitor at 30 frames per second. In addition, the full array of modern video image processing technology is available to further manipulate the image once it is recorded and stored in the system. Another advantage is the relatively high light sensitivity of TV cameras which operate at low light levels, thus enabling satisfactory ESPI measurements with relatively low power lasers.

The ready availability of advanced video recording and processing equipment, its ease of use, and its flexible adaption to various applications make ESPI a near-ideal measurement system in many instances—particularly, industrial situations (such as an assembly line) where rapid data generation and retrieval, and high throughput are required. ESPI overcomes many of the objections of holographic interferometry and has been used extensively for industrial measurements.

Intensity correlation fringes in ESPI are produced by a process of video subtraction or addition. In the subtraction process, the image of a displaced object is subtracted from an electronically

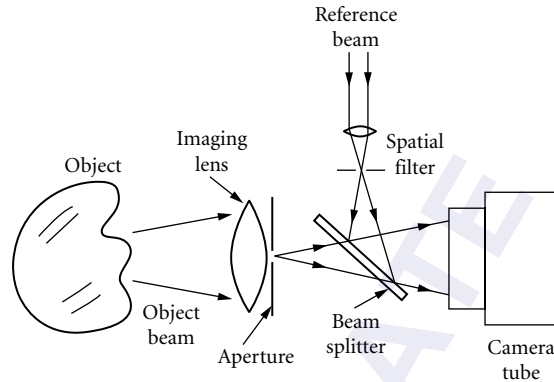
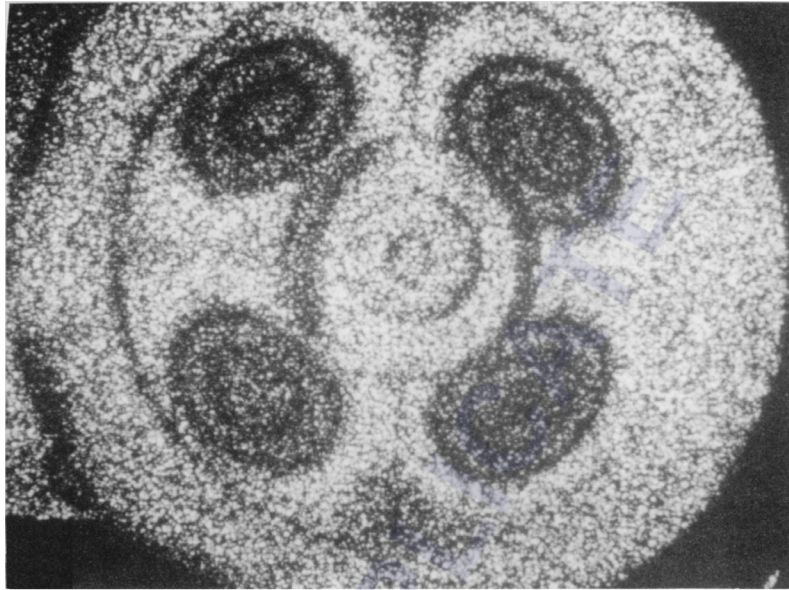


FIGURE 4 Typical ESPI optical arrangement with in-line reference beam.

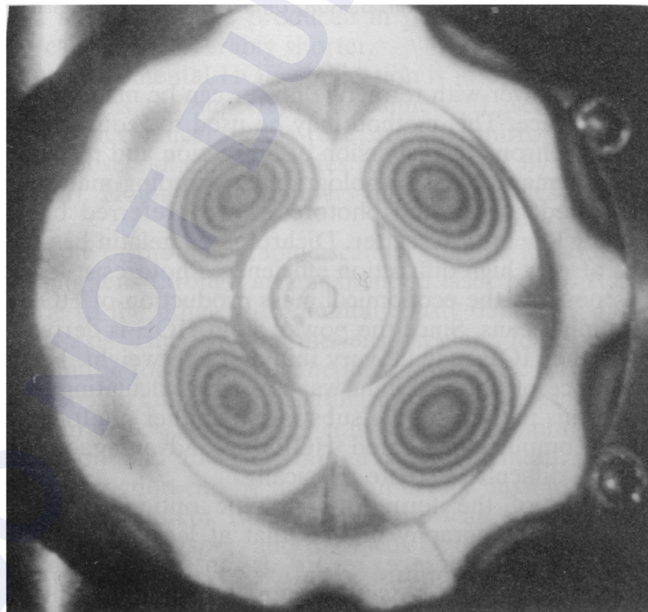
stored image of the object in its original position to produce the correlation fringes. To observe these fringes, the subtracted video signal is rectified and high-pass filtered, then displayed on a video monitor. This video processing is analogous to the reconstruction step in holography.

For the addition method, both original and displaced images are added optically at the photo cathode of the TV camera. The TV camera detects the light intensity and, again, the signal is full-wave rectified, high-pass filtered, and displayed on the TV monitor. Because of the persistence of the TV tube, the two images need not be recorded simultaneously; however, they must be presented to the camera within its persistence time, usually on the order of 100 ms. The various optical configurations used in speckle-pattern correlation interferometry, employing both specular and diffuse reference beams, may be used in ESPI as well. The most popular of these configurations uses a specular in-line reference beam and may be used to make real-time vibration studies. ESPI has been used for this purpose more than any other application. This optical configuration is shown in Fig. 4. The object and reference beams in Fig. 4 are derived from the same laser with the use of a beam splitter, not shown for simplicity. In holographic interferometry, a high-resolution film is used and the reference beam angle may be any practical value desired. In ESPI, the recording medium (TV camera) has a resolution two orders of magnitude lower than holographic film (on the order of 30 line-pairs per millimeter). Therefore, in ESPI, an in-line reference beam must be used. Furthermore, the aperture of the system must be small enough to keep the interference angle below 1° . All the usual modes of holographic interferometry (real-time, time-average, stroboscopic, and double-exposure) may be performed with ESPI. A time-average ESPI interferogram of an object made with the system operating in the subtraction mode is shown in Fig. 5a. For comparison, a holographic interferogram of the same object is shown in Fig. 5b. ESPI has been applied to a wide range of measurement problems. These applications are discussed in numerous books, technical papers, and review articles.⁶⁴⁻⁷¹

Despite the flexibility of ESPI and its ease of use, fringe definition is poor compared to holographic interferometry—this has somewhat limited its use. Speckle-averaging and video-processing techniques have provided some improvement in interferogram quality, but very fine interference fringes are still difficult to discern with ESPI. A significant improvement in interferogram quality has been achieved with a newer technique: electro-optic holography (EOH).⁷²⁻⁷⁴ This technique uses the same speckle interferometer optical configuration as ESPI, but processes the video images in a different manner. In EOH, a phase-stepping mirror is added to the reference leg to advance the phase of the reference beam by 90° between successive video frames. Subsequent processing of these phase-stepped images combined with frame and speckle averaging provides interferograms in real-time with nearly the same



(a)



(b)

FIGURE 5 Time-average interferograms of a pulley made with (a) ESPI and (b) holographic interferometry.

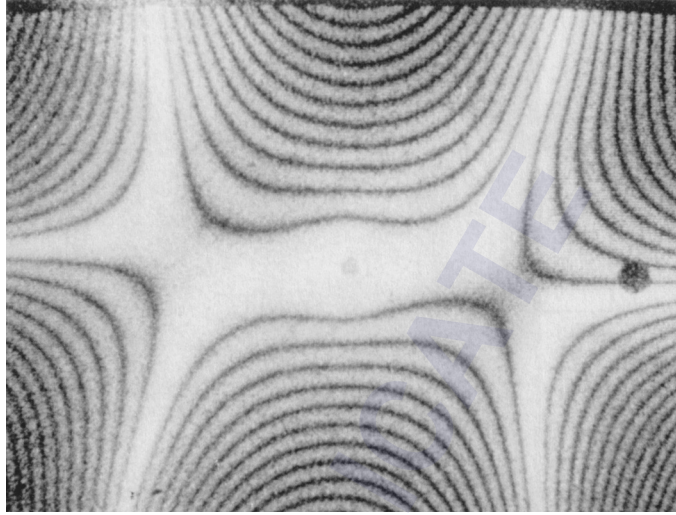


FIGURE 6 An interferogram illustrating a vibrating mode of a center-mounted rectangular plate made with an electronic holography system. (Photo courtesy of Karl Stetson, United Technologies Research Center.)

resolution and clarity of the traditional film-based holographic interferogram. An interferogram made with an EOH system is shown in Fig. 6.

33.5 HOLOGRAPHIC OPTICAL ELEMENTS

An optical element with the power to direct and/or focus a light wave can be made by recording the fringe pattern formed by two interfering light beams. HOEs or diffractive optical elements, thus produced, have several features that distinguish them from conventional optics. A compilation of articles on diffractive optics may be found in Ref. 75. References 76 and 77 provide past reviews of the field.

Different types of elements may be produced by varying the curvature of the interfering wavefronts, interbeam angle, and configuration of the recording surface. Lenses with mild or very strong focusing power and plane or focusing mirrors are easily produced. The recording substrate may be plane or any arbitrarily curved shape. An element combining diffractive power with refractive power can be made by placing the recording medium on a curved surface. This method may be used to reduce the aberration of the combined optic since the achromatic dispersion of diffraction and refraction are of opposite sign. Although any of the many types of holographic recording materials may be used to make an HOE, dichromated gelatin and photoresist are preferred because of their high resolution and extremely low optical scatter. Dichromated gelatin has the additional advantage of forming HOEs of very high diffraction efficiency. Photoresist forms a surface relief hologram which makes possible the economical mass production of HOEs with straightforward mechanical replication means. Since the power of an HOE is derived from diffraction at the element's surface, the HOE may be very thin and lightweight. HOEs may be produced by the direct interference of physical light waves, or by calculating the desired interference pattern and printing this pattern onto a substrate by either photographic or electron beam lithographic means. Computer-generated HOEs are advantageous when the required optical wavefronts are difficult to create physically.

The use of the term "holographic optics" is not technically correct because the definition of holography implies that at least one of the wavefronts being combined to produce an interference record is an information carrier. Consequently, the term "diffractive optics" has gained popularity

and, when applied to gratings, the term “interference gratings” is certainly more appropriate. In this brief discussion, however, we will continue to use the “holographic” terminology.

Certainly one of the most common and successful applications of holographic optics is as a grating in spectrographic instruments.^{78–80} The main advantage of holographic gratings* are that they can be made free of the random and periodic groove variation found in even the finest-ruled gratings, and they have low light scatter. This latter property is especially important when even a small amount of stray light is objectionable (such as in the study of Raman spectra of solid samples). To produce high-quality holographic gratings, extreme care must be used in the fabrication process. Photoresist is the preferred recording material for reasons mentioned above; however, it is very insensitive and requires long exposure times, often many hours. Therefore, a highly stable optical system is essential. The recording room must be free of air currents, the air must be filtered and dust-free, and the photoresist coating must be defect-free. Stray light scatter from optical mounts and other objects in the recording setup must be eliminated by proper baffling and masking. The interfering beams must be appropriately conditioned by spatial filtering to ensure diffraction-limited performance. If beam-forming optics are used, for example, to produce collimated beams for making plane holographic gratings, these optics must be aberration-free and of diffraction-limited quality.

After exposure and chemical development, the surface relief pattern is metalized and the holographic grating is replicated as conventional master-ruled gratings are replicated. Both positive and negative photoresists are available for making holographic gratings, however, the negative resist is seldom used.

Grating-diffraction efficiency in the various orders is determined by the groove profile. In a ruled grating, the groove is profiled by appropriately shaping the diamond tool. Holographic gratings have a sinusoidal profile; blazing, in this case, is accomplished by ion etching. A wide range of groove spacing is possible with ruled gratings, however, holographic gratings offer more flexibility with respect to the groove pattern. For example, groove curvature may be used in holographic gratings to reduce aberrations in the spectrum, thereby improving the throughput and resolution of imaging spectrometers. The grooves in ruled gratings are produced by the traveling diamond tool, one after another. In holographic gratings, all grooves are produced in parallel; thus, the fabrication time for holographic gratings can be considerably shorter.

Another important instrumental HOE application is in optical beam scanning; the most common example of this is the supermarket scanner.⁸¹ Laser beam scanning is usually accomplished by either mechanical means (e.g., rotating mirror) or with the use of some transparent medium whose optical properties are changed by some sort of stimulation (e.g., acousto-optic cell). Holographic scanners offer advantages over both of these more conventional methods.

The working principle of the holographic scanner may be illustrated by considering the translation of a focusing lens through an unexpanded laser beam. As the beam intercepts the lens from one side to the other, it is simultaneously deflected and focused at the lens focal point on the lens axis. Thus, moving the lens back and forth causes the laser beam to sweep back and forth in the focal plane of the lens. In its basic form, a holographic scanner is simply an HOE lens or mirror translating through the scan beam. The principal advantage of the holographic scanner over conventional scanning means is the ability to combine beam deflection and focusing into a single element. The form of both the deflection and focusing function can be tailored in a very flexible way by proper design of the HOE formation optical system. For example, the scan element may be easily made a line segment rather than a focal spot, and the locus of the scanned focal spot or line may be placed on either plane or curved surfaces. Multiple scan beams with multiple focal points may be generated by a scanner in the form of a segmented rotating disk. Each segment or facet of this disk has different deflection and focal properties. As the disk rotates through the beam, a multiplicity of scan beams is produced which can densely fill the desired scanned volume. This feature is especially important in a supermarket application because it allows products of varying sizes and shapes to be rapidly scanned. A further advantage of the HOE scanner is that the scanner disk can be small, thin, and lightweight, thereby greatly reducing the demands on the drive system. The disk format also produces little air movement due to windage and is very quiet in operation.

*The discussion here on holographic gratings is taken largely from Ref. 80. The author is indebted to Christopher Palmer of the Milton Roy Company for making an advance copy of this material available.

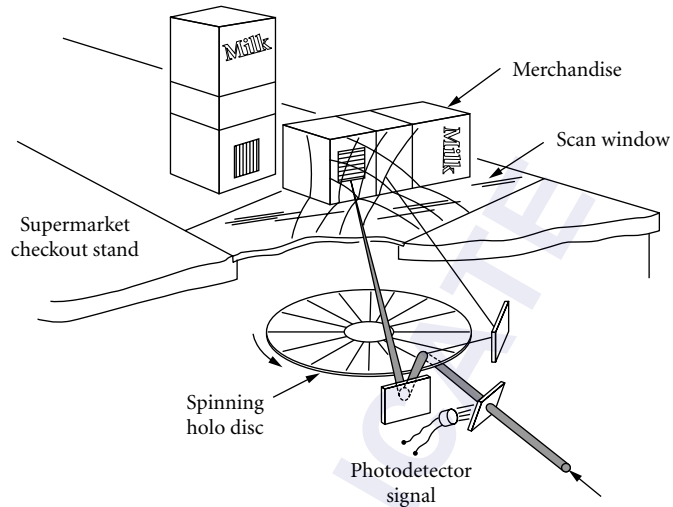


FIGURE 7 Schematic of a holographic supermarket scanner. (Reprinted from Ref. 81, p. 9; courtesy of Marcel Dekker, Inc.)

In addition to serving as the beam deflector, the holographic scanner also collects the light scattered from the laser spot on the product and directs it to the optical detector. This scattered light illuminates the holographic scanner along the conjugate object beam path and is diffracted into the fixed conjugate reference or primary scan beam path where it is accessed by a beam splitter. The light scattered from all points in the scan volume is thus directed to a single, fixed detector position. An optical schematic for a point-of-sales scanner is shown in Fig. 7.

The ability to combine several optical functions into a single HOE makes this device attractive in many situations. Significant savings in space, weight, and cost can often be realized by replacing several conventional elements with a single HOE device. This feature has been incorporated into an optical head for compact disk applications with the use of a multifunctional HOE.⁸² An optical diagram of the device is shown in Fig. 8. The objective lens images the light-emitting point

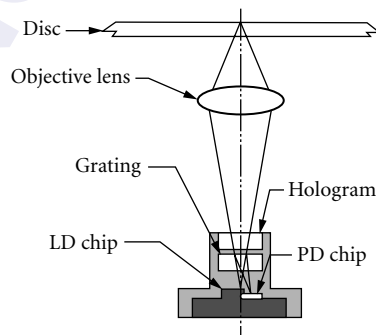


FIGURE 8 Diagram of compact disk optical head employing a holographic optical element. (Diagram courtesy of Wai-Hon Lee, Hoetron, Inc.)

of the laser diode (LD) to the compact disk. The light scattered from this focal point on the disk is reimaged by the objective lens through the HOE to the photodetector (PD). In this application, the HOE serves as a spherical lens, beam splitter, aberration-correcting lens, and cylindrical lens. In addition to simplifying the optical system, the HOE provides a better means of aligning the optical system.

Holography has even been applied to one of the oldest instrumental functions known—the keeping of time. This has been accomplished by using an HOE as a holographic sundial.⁸³

33.6 HOLOGRAPHIC INSPECTION

Quality assurance inspection and testing, important functions in any industrial manufacturing process, have also benefited from advances in holography. Holographic methods have been applied to the quality control problem in several areas, such as identifying and locating subsurface mechanical defects and determining the presence or absence of certain surface features. Holographic interferometry has been successfully applied to some of these problems,^{84–86} and optical processing methods have also yielded good results in many cases.⁸⁷ Matched and spatial filtering in the Fourier transform plane of an optical processing system have proved to be especially powerful means of identifying features and determining surface detail. In this section, we describe a unique combination of holography and classical optical processing methods that made possible a very successful means of rapidly detecting defects in devices with highly regular and repetitive patterns.

Defects in integrated circuit photomasks and wafers at various stages of processing can greatly diminish the final device yield. Since the economics of wafer production is strongly influenced by yield, there has been an ever-present incentive to increase this yield by minimizing the number of defects introduced at various points in the production process. One way to increase yield is to identify these defects early, and eliminate them before they cascade in a multiplicative manner through the various stages of the process.

For years, the inspection of integrated circuit photomasks and wafers was performed by either manual microscopic examination or automated serial scanning using an optical detector. The latter method involves comparison of the detector signal from a magnified portion of the test pattern to a similar portion of a reference pattern, adjacent pattern, or digital design database. In either case, the inspection process was long and tedious, often requiring many hours or days to inspect a single photomask or wafer. These methods are very slow because of the large number of pixels involved and the serial pixel-by-pixel nature of the inspection. Clearly, a great advantage would be afforded by the ability to examine all the pixels in parallel rather than serial format.

This observation prompted a number of researchers to consider optical processing methods for integrated circuit inspection and for addressing other types of problems involving highly repetitive patterns such as cathode-ray tube masks and TV camera tube array targets.^{88–90} The concept in all cases was to eliminate perfect pattern information and highlight defect areas of the image by spatial filtering in the Fourier transform plane. These methods met with only limited success because of the difficulty in fashioning effective blocking filters and the need for extremely high quality, large-aperture, low F-number Fourier transform lenses. Despite considerable work in this area, none of these efforts resulted in the development, production, and in-process use of integrated circuit inspection systems using Fourier optical processing concepts.

This situation was reversed by the adaption of a holographic documentation system used to document the surface microstructure of high-energy laser optical component test samples.^{91,92} A schematic diagram of the holographic documentation optical system is shown in Fig. 9. An F/3.42 lens was used to image the test target onto the hologram film and an argon laser operating at a wavelength of 514.5 nm was the illumination source. A polarizing beam-splitting cube and a half-wave plate were used to split the beam into object and reference beams, and to adjust the beam ratios. A second polarizing beam-splitting cube and quarter-wave plate were used to efficiently illuminate the test target and direct the reflected light to the holographic plate. The holographic image was reconstructed

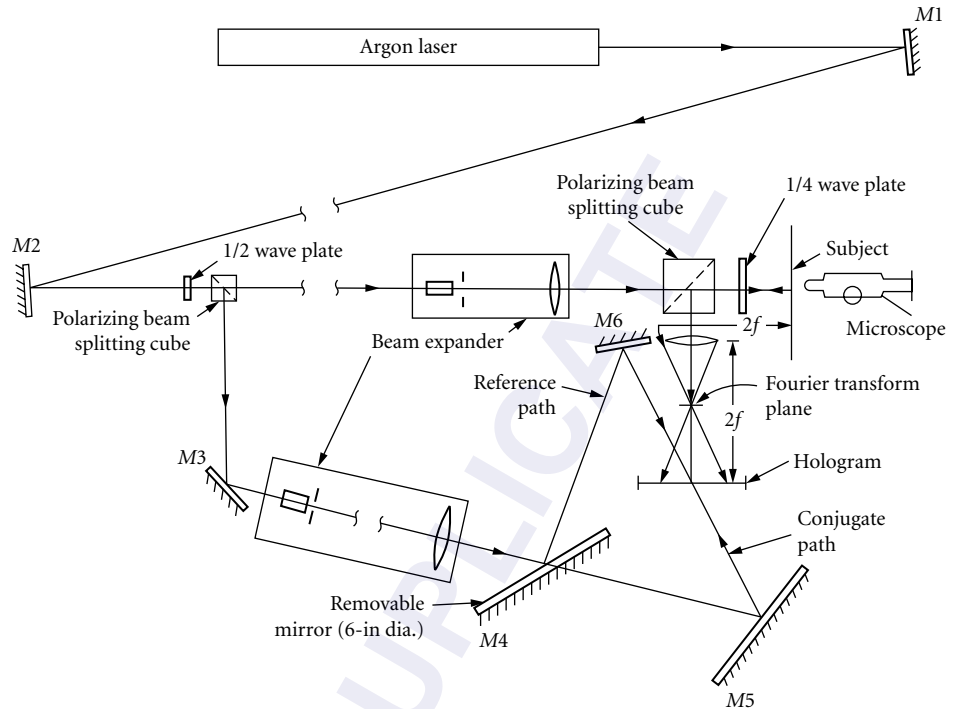


FIGURE 9 Holographic documentation optical system. (Reprinted from Ref. 92, p. 87; courtesy of Oxford University Press.)

with the conjugate reference beam by removing mirror *M4*. The test target was removed and the holographic image of the sample was examined with the aid of a microscope. The calculated resolution of this system (classic Rayleigh resolution limit) was $4.0\ \mu\text{m}$. A photograph of the holographic reconstruction of the standard Air Force resolution target made with this system is shown in Fig. 10. The smallest bars in this target are on $4.4\text{-}\mu\text{m}$ centers. Thus, the resolution observed with this system was comparable to the calculated value.

In considering how the documentation system might be adapted to other applications (including the integrated circuit inspection problem), Fusek et al. observed that because the real image of the test target was being examined by conjugate reconstruction, both functions of the classic Fourier optical system (i.e., transform and inverse transform) were performed in the documentation system.^{93,94} Because of reverse ray tracing, a high-quality matched pair of specially designed Fourier transform lenses is no longer required to produce a diffraction-limited output image. As with previous work, the objective was to attenuate the image area where the pattern is defect-free and highlight the defects. This is done simply by placing the appropriate blocking filter in the Fourier transform plane. The method works effectively only if the filter efficiently blocks the light associated with defect-free areas of the image and efficiently transmits the defect light. Fortunately, this is the case for integrated circuit masks and wafers which consist of regular patterns of circuit elements repeated many times over the area of the wafer. For such patterns, the intensity distribution in the Fourier transform plane is a series of sharp spikes or bright points of light. Low spatial frequencies associated with slowly varying features (such as line-spacing) are represented by light points near the optical axis of the imaging lens. High spatial frequencies, representing such features as edge and corner detail, lie farther out in the Fourier transform plane. The Fourier transform plane intensity pattern for a production-integrated circuit photomask is shown in Fig. 11. Since defects are usually

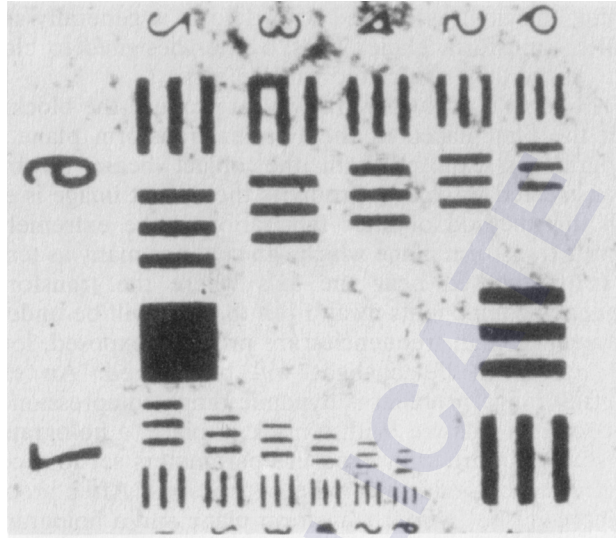


FIGURE 10 Magnified image of the holographic reconstruction of the Air Force resolution target. (Reprinted from Ref. 92, p. 88; courtesy of Oxford University Press.)

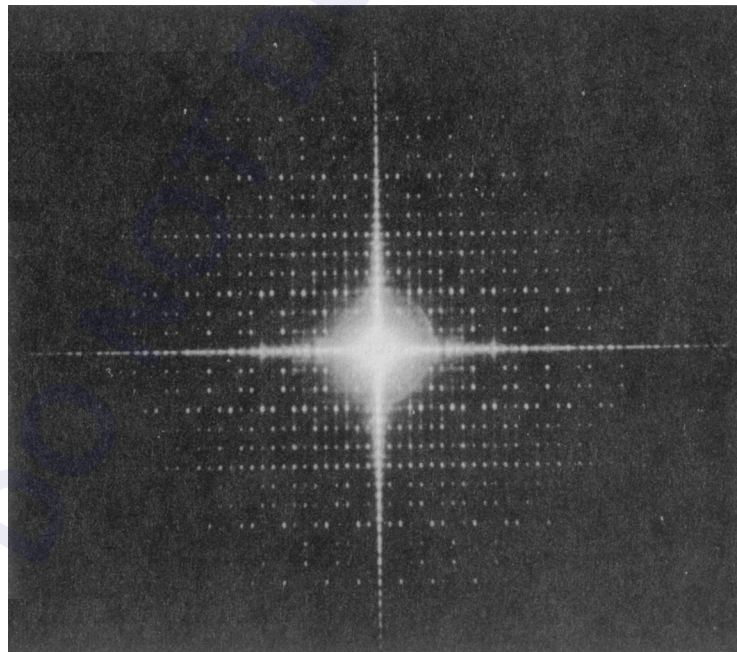


FIGURE 11 Optical Fourier transform of a production-integrated circuit photo-mask. (Reprinted from Ref. 92, p. 98; courtesy of Oxford University Press.)

of a nonregular nature, the light associated with defects is generally spread out fairly uniformly over the Fourier transform plane. Thus, a filter designed to block the regular-pattern light passes most of the light associated with defects.

The most straightforward way to produce the blocking filter is by photographic means, with the film placed in the Fourier transform plane. The filter is made first, then the hologram is exposed with the object beam passing through this filter. Thus, the reconstructed wave that produces the output image is effectively filtered twice. A problem with this method of filter fabrication is the extremely large range of intensities in the Fourier transform plane which can cover as many as ten orders of magnitude. If the filter is properly exposed near the axis where the transform light is brightest, high-spatial-frequency components away from the axis will be underexposed and inadequately filtered. However, if high frequencies are properly exposed, low frequencies will be overexposed, and too much defect light will be filtered. An effective means of alleviating this intensity-range problem is dynamic range compression by multistep filter generation. The process is as follows. With a mask in place, a holographic plate is exposed to the Fourier transform pattern with exposure parameters set to record the intensity distribution in the low-frequency region near the optical axis. After processing, this plate (Stage 1 filter) is replaced in the Fourier transform plane and a hologram of the mask is made through this filter. The hologram is illuminated by the conjugate reference beam and a second filter plate (Stage 2 filter) is exposed to the resulting Fourier transform intensity pattern. This pattern now has its low-frequency components attenuated because of the action of the Stage 1 filter. The Stage 2 filter can now be used to record a Stage 3 filter. The process may be repeated as many times as necessary to produce a filter with the desired attenuation properties. Because the defect light is of much lower intensity in the Fourier transform plane than the light corresponding to nondefect areas, defect light does not contribute significantly to the exposure of the filter, and the object under test (containing defects) may be used to generate the filter. Stage 1 and 2 filters for a defect calibration test mask (Master Images VeriMask™) are shown in Fig. 12.

Performance of the breadboard documentation system using the VeriMask is shown in Fig. 13. Figure 13*a* is a photo of a magnified region of the VeriMask containing a pinhole defect. Figure 13*b* shows the Stage 2 filter image of this same defect which is clearly enhanced. In addition, dimensional variations in the mask pattern from die to die are also highlighted.

The breadboard holographic inspection technology developed by Fusek and coworkers was further advanced and placed into production by Insystems of San Jose, California.⁹⁵⁻⁹⁸ This company produced a series of mask and wafer inspection machines based on the holographic optical processing technique. The optical configuration of the Insystems Model 8800 Wafer Inspection System is shown in Fig. 14. The commercial instrument used an argon ion laser as the light source, and the system functioned in an optical manner identical to that of the original breadboard device. However, many refinements were incorporated into the commercial system which yielded substantially improved performance over the breadboard system. These refinements, which included a sophisticated Fourier transform lens design, made possible adequate performance without using the multistep filter generation technique. The wafer test piece and the hologram were placed on a rotating stage under a microscope and a video camera so that the filtered image of the defect and the microscopic image of the defect on the actual wafer could be viewed simultaneously. Figure 15 illustrates the advanced filtering capability of the commercial instrument. This instrument was sensitive to defects as small as 0.35 μm .

Disadvantages of the holographic defect detection system are the inconvenience and time delays associated with the wet processing of the silver halide holographic recording material. The use of photorefractive crystals, which operate in real time and do not require wet processing, has been studied as a means of eliminating these disadvantages.⁹⁷ The dual functions of image recording and spatial filtering are combined by placing the crystal in the Fourier transform plane and adjusting the reference beam intensity to the level of the defect light intensity. Since the light in the Fourier transform plane associated with defects is much lower in intensity than nondefect light, only the defects will be recorded with high diffraction efficiency. Practical use of photorefractive crystals in this application has not been realized, however, because these crystals have relatively low sensitivity and are not available in large sizes with good optical quality.

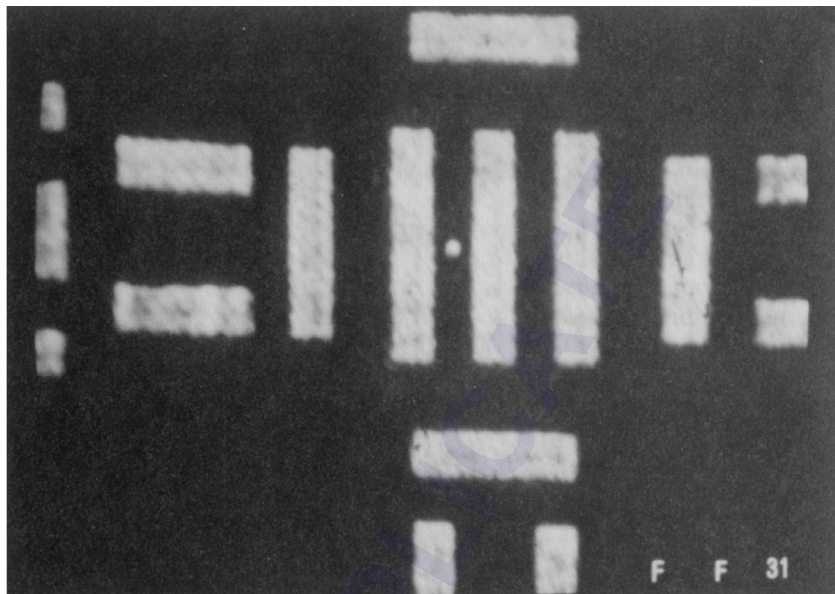


(a)

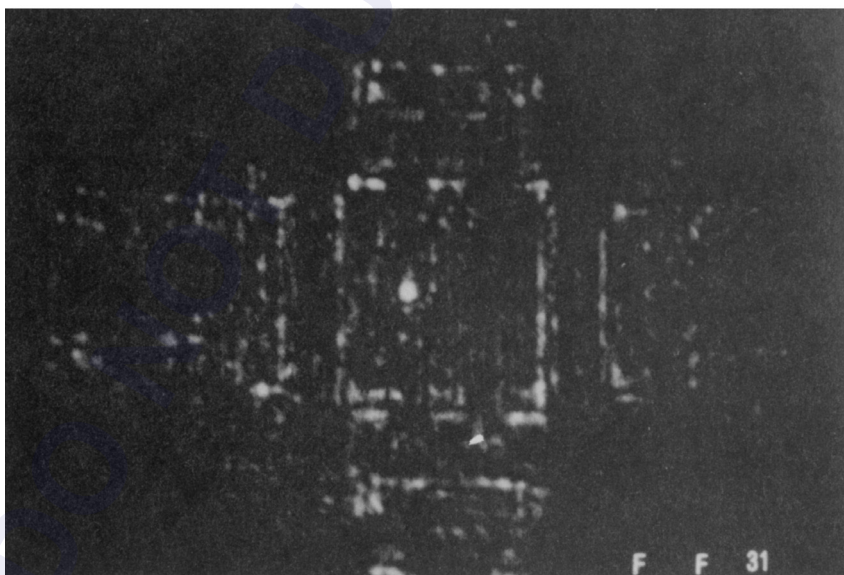


(b)

FIGURE 12 Fourier transform plane blocking filters for a defect calibration test mask: (a) first-stage initial filter and (b) second-stage dynamic range compressed filter. (Reprinted from Ref. 92, p. 101; courtesy of Oxford University Press.)



(a)



(b)

FIGURE 13 Images of the holographic reconstruction of a 2.01- μm pinhole on a calibration test mask: (a) unfiltered and (b) filtered. (Reprinted from Ref. 92, pp. 102–103; courtesy of Oxford University Press.)

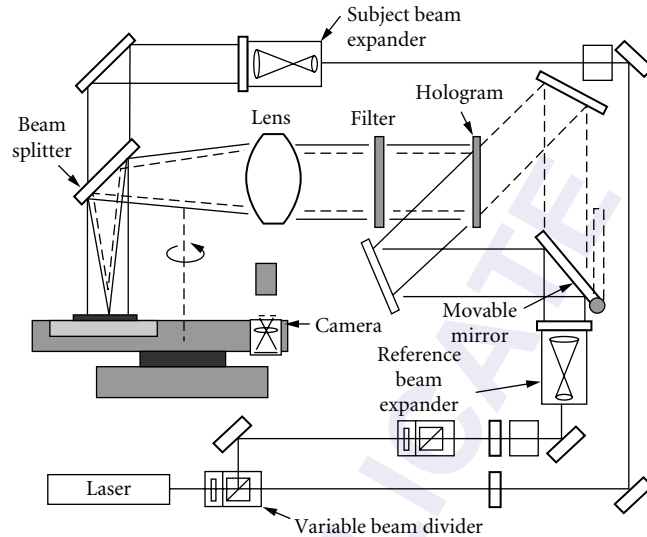


FIGURE 14 Optical schematic for the Insystems Model 8800 holographic wafer inspection system. (Diagram courtesy of Insystems.)

33.7 HOLOGRAPHIC LITHOGRAPHY

The lithographic transfer of an integrated circuit photomask pattern to a resist-coated integrated circuit wafer has been accomplished by several methods, including contact printing, proximity printing, and step-and-repeat imaging. Each method has advantages and disadvantages. Contact printing is a simple, straightforward method suitable for printing large wafer areas, but damage to

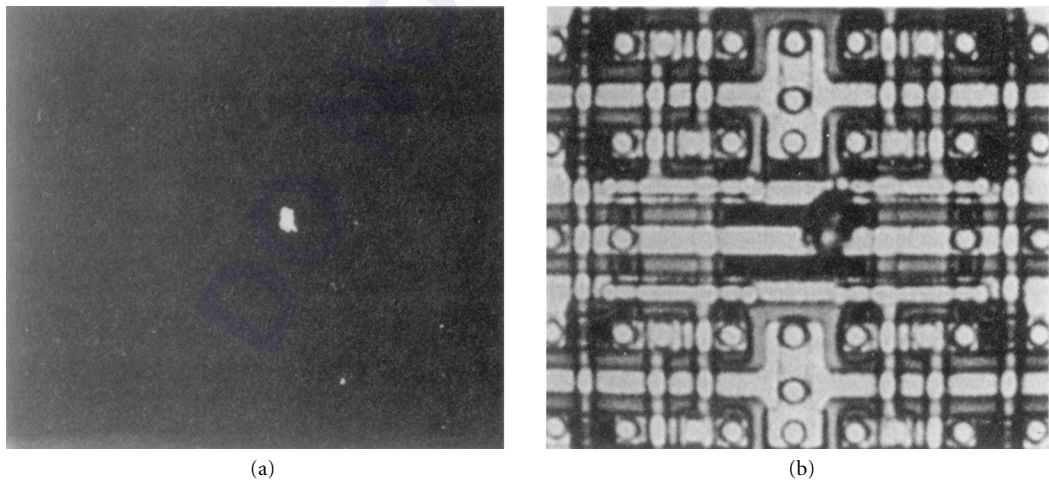


FIGURE 15 Metal layer defect in an integrated circuit wafer highlighted by the Insystems Model 8800 wafer inspection system: (a) filtered image and (b) unfiltered image. (Photo courtesy of Insystems.)

the mask and contamination of the wafer are common problems. Proximity printing is resolution-limited because of near-field diffraction. The diffraction problem can be eliminated by imaging the mask onto the wafer, but full-field imaging systems do not provide the resolution required over the full area of the wafer. Stepper systems, which image only a small area of the mask at a time, provide the required resolution; however, they are complex and expensive because of the resolution demands placed on the optical imaging system and the mechanical difficulty in accurately stitching together the multiple image patterns over the full area of the wafer.

Clearly, a full-field method of printing the image onto the full area of the wafer with the required resolution is desirable. A holographic system capable of accomplishing this task was developed by Holtronic Technologies Limited.⁹⁸⁻¹⁰⁰ Rather than using a lens system to image the mask on to the wafer, the Holtronic holographic system used real-image projection by conjugate illumination to overlay the mask image onto the wafer. Near-field holography was used to record the mask image by placing the mask in close proximity to the recording medium (100- μm separation) and illuminating the mask from the back with a collimated laser beam (364-nm line from an argon ion laser). To allow the introduction of an off-axis reference beam, an image of the mask could be relayed to the hologram plane with a lens. The approach taken was to eliminate the need for this lens with the use of the total internal reflection holography scheme of Stetson.¹⁰¹ A diagram illustrating the optical principle is shown in Fig. 16.

Light transmitted and diffracted by the mask is incident directly on the holographic recording material which is deposited onto the opposing surface of a prism. Reference light is introduced through the diagonal surface of the prism and reflected at the holographic coating/air interface by total internal reflection. The recording medium was Dupont photopolymer. The holographic exposure was made with an expanded reference beam illuminating the entire hologram area.

Since three beams pass through the photosensitive material, three gratings are formed: (1) a reflection grating formed by the incident and reflected reference beams, (2) a reflection grating formed by the object beam and the incident reference beam, and (3) a transmission grating formed

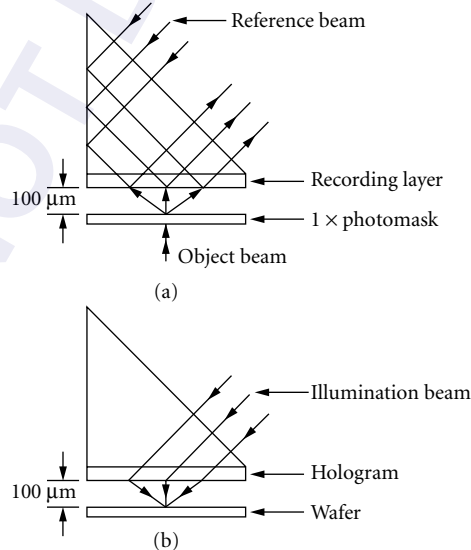


FIGURE 16 Basic optical arrangement for total internal reflection holographic lithography: (a) hologram exposure and (b) reconstruction onto a resist-coated substrate. (Reprinted from Ref. 99; courtesy of PennWell Publications.)

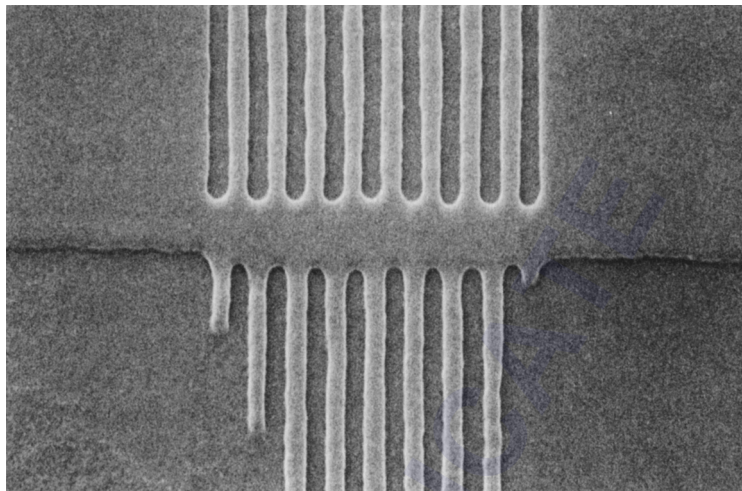


FIGURE 17 Scanning electron micrograph of 0.3- μm lines and spaces printed in photoresist by holographic lithography. (Photo courtesy of Holtronic Technologies Ltd.)

by the object beam and the reflected reference beam. To reconstruct the mask wavefront, an illumination beam conjugate to the reflected reference beam is introduced through the prism. This incident illumination beam interacts with the transmission object grating to form the conjugate real image of the mask. The totally internally reflected and the Lippmann-reflected illumination beams interact with the reflection object grating to reinforce the mask image.

During reconstruction and exposure of the photoresist-coated wafer, a small-area illumination beam is scanned over the hologram surface. Dynamic focusing of this scanned exposure beam eliminates the requirement that the mask and wafer substrates be ultraflat. Figure 17 illustrates the 0.3- μm resolving capability of the holographic lithography process in 0.3-mm thick resist. The lines shown were printed on a silicon wafer coated with Olin Hunt HPR 204 i-line photoresist and developed using the Hunt HPRD 429 developer. The effective numerical aperture of the system was greater than 0.7.

33.8 HOLOGRAPHIC MEMORY

The information storage capability of holograms has been the subject of considerable study over the years with several applications in mind. With the parallel information storage and processing capability of holograms, and the promise of shorter access times, computer memory has received particular attention. Two review articles provide good summaries of this field of research through 1990.^{102,103} However, little mention is made of the work of Russian scientists who have also been very active in this field (see, for example, Refs. 104–116).

Of all the holographic material recording possibilities available, volume storage in photorefractive crystals has received the most attention. There are two main reasons for this emphasis: (1) the large information storage capacity of these crystals, and (2) their capability to meet the write-read-erase requirement in real time with no wet-chemical or other material-processing delays. However, no commercial memory systems using holographic storage have been developed. One reason for this is the relatively large volume of space occupied by the laser beam-steering equipment and other optical components required in such a system, even though the actual holographic storage element may occupy a volume of less than a few cubic centimeters. However, the primary reason

is the limitations of the recording material. Despite their distinct advantages over other recording material candidates, photorefractive crystals have some significant limitations. It is difficult to grow large crystals with good optical quality and to achieve stable, long-term storage without destructive readout.

Work by Redfield and Hesselink¹¹⁷⁻¹¹⁹ was directed toward overcoming these previous limitations. Rather than concentrating on developing large, high-quality crystals, their approach was to form a large-volume memory element by using an array of small crystallites of strontium barium niobate in the form of small cubes or crystalline fibers. Techniques have also been developed for accessing the holographic information stored in these crystallites without destructive readout. Information is stored in the memory structure by recording Fourier holograms of checkerboard patterns (pages) of digital information. Access times were projected to be 100 to 1000 times faster than with conventional magnetic disk drives.

33.9 CONCLUSION

This chapter has briefly reviewed some of the more important instrumental applications of holography and demonstrated how holographic methods have been used to creatively solve a variety of measurement and recording problems. These successful applications should pave the way for additional advances in this field. Consequently, we anticipate that the list of technical applications of holography will expand significantly in the future.

33.10 REFERENCES

1. D. Gabor, "A New Microscope Principle," *Nature* **161** (1948).
2. D. Gabor, "Microscopy by Reconstructed Wavefronts," *Proc. Roy. Soc.* **A197** (1949).
3. D. Gabor, "Microscopy by Reconstructed Wavefronts: II," *Proc. Phys. Soc.* **B64** (1951).
4. E. N. Leith and J. Upatnieks, "Reconstructed Wavefronts and Communication Theory," *J. Opt. Soc. Am.* **52**(10) (1962).
5. E. N. Leith and J. Upatnieks, "Wavefront Reconstruction with Diffused Illumination and Three-Dimensional Objects," *J. Opt. Soc. Am.* **54**(11) (1964).
6. R. J. Collier, C. B. Burchhardt, and L. H. Lin, *Optical Holography*, Academic Press, New York, 1971.
7. H. M. Smith, *Principles of Holography*, 2d ed., John Wiley & Sons, New York, 1975.
8. H. J. Caulfield (ed.), *Handbook of Optical Holography*, Academic Press, New York, 1979.
9. N. Abramson, *The Making and Evaluation of Holograms*, Academic Press, New York, 1981.
10. P. Hariharan, *Optical Holography*, Cambridge University Press, Cambridge, 1984.
11. G. Saxby, *Practical Holography*, Prentice-Hall, New York, 1988.
12. J. M. Burch, "The Application of Lasers in Production Engineering," *Prod. Eng. (London)* **44**:431 (1965).
13. K. A. Haines and B. P. Hildebrand, "Contour Generation by Wavefront Reconstruction," *Phys. Lett.* **19**:10 (1965).
14. R. J. Collier, E. T. Doherty, and K. S. Pennington, "Application of Moire Techniques to Holography," *Appl. Phys. Lett.* **7**:223 (1965).
15. R. E. Brooks, L. O. Heflinger, and R. F. Wuerker, "Interferometry with a Holographically Reconstructed Comparison Beam," *Appl. Phys. Lett.* **7**:248 (1965).
16. R. L. Powell and K. A. Stetson, "Interferometric Vibration Analysis by Wavefront Reconstruction," *J. Opt. Soc. Am.* **55**:1593 (1965).
17. K. A. Stetson and R. L. Powell, "Interferometric Hologram Evaluation and Real-Time Vibration Analysis of Diffuse Objects," *J. Opt. Soc. Am.* **55**:1694 (1965).

18. L. O. Heflinger, R. F. Wuerker, and R. E. Brooks, "Holographic Interferometry," *J. Appl. Phys.* **37**:642 (1966).
19. M. A. Monahan and K. Bromley, "Vibration Analysis by Holographic Interferometry," *J. Acoust. Soc. Am.* **44**:1225 (1968).
20. B. P. Hildebrand and K. A. Haines, "Multiple-Wavelength and Multiple-Source Holography Applied to Contour Generation," *J. Opt. Soc. Am.* **57**:155 (1967).
21. T. Tsuruta, N. Shiotake, J. Tsujiuchi, and K. Matsuda, "Holographic Generation of Contour Map of Diffusely Reflecting Surface by Using Immersion Method," *Jpn. J. Appl. Phys.* **6**:661 (1967).
22. N. Shiotake, T. Tsuruta, Y. Itoh, J. Tsujiuchi, N. Takeya, and K. Matsuda, "Holographic Generation of Contour Map of Diffusely Reflecting Surface by Using Immersion Method," *Jpn. J. Appl. Phys.* **7**:904 (1968).
23. J. D. Trolinger, "Automated Data Reduction in Holographic Interferometry," *Opt. Eng.* **24**(5) (1985).
24. R. J. Pryputniewicz, "Time Average Holography in Vibration Analysis," *Opt. Eng.* **24**(5) (1985).
25. R. J. Pryputniewicz, "Quantification of Holographic Interferograms: State of the Art Methods," *Topical Meeting on Holography Technical Digest* **86**(5), *Opt. Soc. Am.* Washington, D.C. (1986).
26. R. J. Pryputniewicz, "Review of Methods for Automatic Analysis of Fringes in Hologram Interferometry," *SPIE Proc.* **816** (1987).
27. R. J. Pryputniewicz, "Automated Systems for Quantitative Analysis of Holograms," *SPIE Institute Series*, vol. **IS 8** (1990).
28. N. Abramson, "Sandwich Hologram Interferometry: A New Dimension in Holographic Comparison," *Appl. Opt.* **13**(9) (1974).
29. H. Bjelkhagen, "Sandwich Holography for Compensation of Rigid Body Motion and Reposition of Large Objects," *SPIE Proc.* **215** (1980).
30. G.O. Reynolds, D. A. Servaes, L. Ramos-Izquierdo, J. B. DeVelis, D. C. Peirce, P. D. Hilton, and R. A. Mayville, "Holographic Fringe Linearization Interferometry for Defect Detection," *Opt. Eng.* **24**(5) (1985).
31. Z. Fuzessy and F. Gyimesi, "Difference Holographic Interferometry: An Overview," *SPIE Institute Series*, vol. **IS 8** (1990).
32. P. D. Plotkowski, Y. Y. Hung, J. D. Hovanesian, and G. Gerhart, "Improved Fringe Carrier Technique for Unambiguously Determination of Holographically Recorded Displacements," *Opt. Eng.* **24**(5) (1985).
33. A. A. Kamshilin, E. V. Mokrushina, and M. P. Petrov, "Adaptive Holographic Interferometers Operating Through Self-Diffraction of Recording Beams in Photorefractive Crystals," *Opt. Eng.* **28**(6) (1989).
34. V. I. Vlad, D. Popa, M. P. Petrov, and A. A. Kamshilin, "Optical Testing by Dynamic Holographic Interferometry with Photorefractive Crystals and Computer Image Processing," *SPIE Proc.* **1332** (1990).
35. R. Dandliker and R. Thalmann, "Heterodyne and Quasi-Heterodyne Holographic Interferometry," *Opt. Eng.* **24**(5) (1985).
36. X. Youren, C. M. Vest, and E. J. Delp, "Optical and Digital Moiré Detection of Flaws Applied to Holographic Nondestructive Testing," *Appl. Opt.* **8**:452-454 (1983).
37. M. Cormier, J. Lewandowski, B. Mongeau, F. Ledoyen, and J. Lapiere, "Infrared Holographic Interferometry," *Topical Meeting on Holography Technical Digest* **86**(5), *Opt. Soc. Am.* Washington, D.C. (1986).
38. J. A. Gilbert and T. D. Dudderar, "The Use of Fiber Optics to Enhance and Extend the Capabilities of Holographic Interferometry," *SPIE Institute Series*, **IS 8** (1990).
39. T. C. Lee, "Holographic Recording on Thermoplastic Films," *Appl. Opt.* **13**(4) (1974).
40. N. G. Patil, C. R. Prasad, and V. H. Arakeri, "Holographic Interferometric Study of Heat Transfer in Rectangular Cavities," *Topical Meeting on Holography Technical Digest* **86**(5), *Opt. Soc. Am.* Washington, D.C. (1986).
41. H. Fenichel and M. Lin, "Application of Holographic Interferometry to Investigations of Diffusion Processes in Liquid Solutions," *SPIE Proc.* **523** (1985).
42. H. Fenichel, G. E. Lohman, and D. Will, "Measurements of Diffusion Coefficients in Liquids Using Holographic Interferometry," *Topical Meeting on Holography Technical Digest* **86**(5), *Opt. Soc. Am.* Washington, D.C. (1986).
43. R. L. Perry and G. Lee, "Holographic Interferometry Applied to Symmetric Aerodynamic Models in a Wind Tunnel," *SPIE Proc.* **523** (1985).

44. V. A. Deason, L. D. Reynolds, and M. E. McIlwain, "Velocities of Gases and Plasmas from Real-Time Holographic Interferograms," *Opt. Eng.* **24**(5) (1985).
45. P. J. Bryanston-Cross, "Holographic Flow Visualization," *J. Phot. Sci.* **37**(1) (1989).
46. C. Agren and K. A. Stetson, "Measuring the Wood Resonance of Treble-Viol Plates by Hologram Interferometry," *J. Acoust. Soc. Am.* **46**(1) (1969).
47. G. von Bally, "Otological Investigations in Living Man Using Holographic Interferometry," in G. von Bally (ed.), *Holography in Medicine and Biology*, Springer Series in Optical Sciences, vol. 18, Springer-Verlag, Berlin, 1979.
48. K. A. Stetson, "Holography as a Tool in the Gas Turbine Industry," *Topical Meeting on Holography Technical Digest* **86**(5), *Opt. Soc. Am.* Washington, D.C. (1986).
49. Y. Y. Hung, "Shearography versus Holography in Nondestructive Evaluation of Tyres and Composites," *SPIE Proc.* **814** (1987).
50. G. M. Brown, R. M. Grant, and G. W. Stroke, "Theory of Holographic Interferometry," *J. Acoust. Soc. Am.* **45**(5) (1969).
51. C. M. Vest, *Holographic Interferometry*, John Wiley & Sons, New York, 1979.
52. K. A. Stetson, "A Critical Review of Hologram Interferometry," *SPIE Proc.* **532** (1985).
53. C. M. Vest, "Holographic Metrology and Nondestructive Testing—Past and Future," *Proc. of the NATO Advanced Study Institute*, Martinus Nijhoff, Dordrecht, Netherlands, 1987.
54. P. Hariharan, "Interferometric Metrology: Current Trends and Future Prospects," *SPIE Proc.* **816** (1987).
55. R. J. Parker and D. G. Jones, "Holography in an Industrial Environment," *Opt. Eng.* **27**(1) (1988).
56. B. Ovryn, "Holographic Interferometry," *CRC Critical Reviews in Biomedical Engineering* **16**(4) (1989).
57. H. Rottenkolber and W. Juptner, "Holographic Interferometry in the Next Decade," *SPIE Proc.* **1162** (1990).
58. R. Jones and C. Wykes, *Holographic and Speckle Interferometry*, Cambridge University Press, Cambridge, 1983.
59. R. Jones and C. Wykes, *Holographic and Speckle Interferometry*, Cambridge University Press, Cambridge, 1983, p. 57.
60. J. N. Butters and J. A. Leendertz, "Holographic and Video Techniques Applied to Engineering Measurements," *J. Meas. Control* **4** (1971).
61. A. Macovski, D. Ramsey, and L. F. Schaefer, "Time Lapse Interferometry and Contouring Using Television Systems," *Appl. Opt.* **10**(12) (1971).
62. O. Schwomma, Österreichisches, Patent No. 298830, 1972.
63. U. Kopf, in *Messtechnik* (in German), vol. 4, 1972, p. 105.
64. O. J. Lokberg, "Advances and Application of Electronic Speckle Pattern Interferometry (ESPI)," *SPIE Proc.* **215** (1980).
65. B. D. Bergquist, P. C. Montgomery, F. Mendoza-Santoyo, P. Henry, and J. Tyrer, "The Present Status of Electronic Speckle Pattern Interferometry (ESPI) With Respect to Automatic Inspection and Measurement," *SPIE Proc.* **654** (1986).
66. O. J. Lokberg, "The Present and Future Importance of ESPI," *SPIE Proc.* **746** (1987).
67. O. J. Lokberg, "Electronic Speckle Pattern Interferometry," *Proc. of the NATO Advanced Study Institute*, Martinus Nijhoff, Dordrecht, Netherlands, 1987, pp. 542–572.
68. O. J. Lokberg and G. A. Slettermoen, "Basic Electronic Speckle Pattern Interferometry," chap. 8, in R. R. Shannon and J. C. Wyant (eds.), *Applied Optics and Optical Engineering*, vol. X, Academic Press, Inc., 1987.
69. D. W. Robinson, "Holographic and Speckle Interferometry in the UK: A Review of Recent Developments," *SPIE Proc.* **814** (1988).
70. D. E. Parker, "Introductory Overview of Holography and Speckle," *SPIE Proc.* **1375** (1990).
71. O. J. Lokberg and S. Ellingsrud, "TV-Holography (ESPI) and Image Processing in Practical Use," *SPIE Proc.* **1332** (1990).
72. K. A. Stetson, W. R. Brohinsky, J. Wahid, and T. Bushman, "An Electro-Optic Holography System with Real-Time Arithmetic Processing," *J. Nondest. Eval.* **8**(2) (1989).
73. T. Bushman, "Development of a Holographic Computing System," *SPIE Proc.* **1162** (1989).

74. R. J. Pryputniewicz and K. A. Stetson, "Measurement of Vibration Patterns Using Electro-Optic Holography," *SPIE Proc.* **1162** (1989).
75. T. W. Stone and B. J. Thompson (eds.), "Selected Papers on Holographic and Diffractive Lenses and Mirrors," *SPIE Milestone Series*, vol. MS 34, 1991.
76. S. V. Pappu, "Holographic Optical Elements: State-of-the-Art Review Part 2," *Opt. Laser Technol.* **21**(6) (1989).
77. S. V. Pappu, "Holographic Optical Elements: State-of-the-Art Review Part 1," *Opt. Laser Technol.* **21**(5) (1989).
78. J. M. Lerner, J. Flamand, J. P. Laude, G. Passereau, and A. Thevenon, "Diffraction Gratings, Ruled and Holographic: A Review," *SPIE Proc.* **240** (1980).
79. E. G. Loewen, "Diffraction Gratings, Ruled and Holographic," chap. 2, in *Appl. Opt. and Opt. Eng.* **IX** (1983).
80. "Interference (Holographic) Gratings," chap. 5, in C. Palmer and E. Loewen (eds.), *Diffraction Grating Handbook*, Milton Roy Company, 1991.
81. G. T. Sincerbox, "Holographic Scanners," chap. 1, in G. F. Marshall (ed.), *Laser Beam Scanning*, Marcel Dekker, Inc., New York, 1985.
82. W. Lee, "Holographic Optical Head for Compact Disk Applications," *Opt. Eng.* **28**(6) (1989).
83. K. M. Johnson, B. Cormack, A. Strasser, K. Dixon, and J. Carsten, "The Digital Holographic Sundial," *Topical Meeting on Holography Technical Digest* **86**(5), *Opt. Soc. Am.* Washington, D.C. (1986).
84. K. A. Arunkumar, J. D. Trolinger, S. Hall, and D. Cooper, "Holographic Inspection of Printed Circuit Board," *SPIE Proc.* **693** (1986).
85. Y. Lu, L. Jiang, L. Zou, X. Zhao, and J. Sun, "The Non-Destructive Testing of Printed Circuit Board by Phase Shifting Interferometry," *SPIE Proc.* **1332** (1990).
86. C. P. Wood and J. D. Trolinger, "The Application of Real-Time Holographic Interferometry in the Nondestructive Inspection of Electronic Parts and Assemblies," *SPIE Proc.* **1332** (1990).
87. D. Casasent, "Computer Generated Holograms in Pattern Recognition: A Review," *SPIE Proc.* **532** (1985).
88. L. S. Watkins, *Proc. IEEE* **57**:1634 (1969).
89. N. N. Axelrod, *Proc. IEEE* **60**:447 (1972).
90. R. A. Heinz, R. L. Odenweller, Jr., R. C. Oehrle, and L. S. Watkins, *Western Elect. Eng.* **17**:39 (1973).
91. R. L. Fusek, J. S. Harris, J. Murphy, and K. G. Harding, "Holographic Documentation Camera for Component Study Evaluation," in *High Power Lasers and Applications: Proceedings of the Meeting, SPIE Proc.*, Los Angeles, Calif., February 11–13, 1981.
92. L. Huff, "Holographic Documentation and Inspection," chap. 8, in J. Robillard and H. J. Caulfield (eds.), *Industrial Application of Holography*, Oxford University Press, 1990.
93. R. L. Fusek, L. H. Linn, K. Harding, and S. Gustafson, "Holographic Optical Processing for Submicrometer Defect Detection," *Opt. Eng.* **24**(5) (1985).
94. R. L. Fusek, J. S. Harris, and K. G. Harding, U.S. Patent 4,566,757, 28 January 1986.
95. L. H. Din, D. L. Cavan, R. B. Howe, and R. E. Graves, "A Holographic Photomask Defect Inspection System," *SPIE Proc.* **538**:110–116 (1985).
96. D. L. Cavan, L. H. Lin, R. B. Howe, R. E. Graves, and R. L. Fusek, "Patterned Wafer Inspection Using Laser Holography and Spatial Frequency Filtering," *J. Vac. Sci. Technol.* **B6**(6) (1988).
97. E. Ochoa, J. W. Goodman, and L. Hesselink, "Real-Time Enhancement of Defects in a Periodic Mask Using Photorefractive $\text{Bi}_{12}\text{SiO}_{20}$," *Opt. Lett.* **10**(9) (1985).
98. J. Brook and R. Dandliker, "Submicrometer Holographic Photolithography," *Solid State Technology* (November 1989).
99. B. A. Omar, F. Clube, M. Hamidi, D. Struchen, and S. Gray, "Advances in Holographic Lithography," *Solid State Technology* (September 1991).
100. S. Gray and M. Hamidi, "Holographic Microlithography for Flat-Panel Displays," *SID 91 Digest* (1991).
101. K. Stetson, "Holography with Totally Internally Reflected Light," *Appl. Phys. Lett.* **11**:225 (1967).
102. B. Hill, "Holographic Memories and Their Future," in N. H. Farhat (ed.), *Advances in Holography*, vol. 3, Marcel Dekker, New York, 1976, pp. 1–251.

103. S. V. Pappu, "Holographic Memories: A Critical Review," *Int. J. Optoelectronics* 5:3 (1990).
104. G. A. Voskoboinik, I. S. Gibin, V. P. Koronkevich, E. S. Nezhevenko, P. E. Tverdokhlebl, and Y. V. Ghugui, "Holographic Memory Device for Identifying Substances from Their Infrared Spectra," *Optika i Spektroskopiya* 30(6) (1971), translated in *Optics and Spectroscopy* 30(6) (1971).
105. I. S. Gibin, A. Gofman, S. K. Kibirev, E. F. Pen, and P. E. Tverdokhlebl, "Holographic Memory Devices with Data Search Functions," *Avtometriya* 5:37–51 (1977) translated in *Optoelectronics, Instrumentation and Data Processing* (1977).
106. E. F. Pen, P. E. Tverdokhlebl, Y. N. Tishchenko, and A. V. Trubetskoi, "Acoustooptical Deflector for a Holographic Memory," *Optika i Spektroskopiya* 55(1):148–155 (1983); translated in *Opt. Spectrosc.* 55(1):86–90 (1983).
107. V. A. Dombrovskii, S. A. Dombrovskii, and E. F. Pen, "Investigation of Noise Stability of Holograms in a Holographic Memory," *Avtometriya* 4 (1985), translated in *Optoelectronics, Instrumentation and Data Processing* (1985).
108. A. A. Verbovetskii, A. P. Grammatin, V. N. Ivanov, V. G. Mityakov, A. A. Novikov, N. N. Rukavitsin, Y. S. Skvortsov, V. B. Fedorov, and V. V. Tsvetkov, "Holographic Memory for Archival Storage of Binary Information," *Optiko-Mekhanicheskaya Promyshlennost* 55:5 (1988), translated in *Sov. J. Opt. Technol.* 55:5 (1988).
109. V. A. Dombrovskii, S. A. Dombrovskii, and E. F. Pen, "Reliability of Data Readout in a Holographic Memory Channel with Constant Characteristics," *Avtometriya* 6 (1988), translated in *Optoelectronics, Instrumentation and Data Processing* 6 (1988).
110. Y. V. Vovk, L. V. Vydrin, N. N. V'yukhina, V. N. Zatolokin, P. E. Tverdokhlebl, I. S. Shteinberg, and Y. A. Shchepetkin, "Fast Storage Device for Digital Data Based on a Pack of Optical Disks," *Avtometriya* 3:82–94 (1989), translated in *Optoelectronics, Instrumentation and Data Processing* 3:78–90 (1989).
111. Y. V. Vovk, L. V. Vydrin, P. E. Tverdokhlebl, and Y. A. Shchepetkin, "Method for Multichannel Recording of Binary Data on Optical Disk," *Avtometriya* 2:77–87 (1989), translated in *Optoelectronics, Instrumentation and Data Processing* 2:79–89 (1989).
112. B. V. Vanyushev, N. N. V'yukhina, I. S. Gibin, A. P. Litvintseva, T. N. Mantush, B. N. Pankov, E. F. Pen, A. N. Potapov, I. B. Tatarnikova, and P. E. Tverdokhlebl, "Architecture of Data System Based on Large Capacity Holographic Memory," *Avtometriya* 3:74–82 (1989), translated in *Optoelectronics, Instrumentation and Data Processing* 3:70–77 (1989).
113. A. A. Blok, R. S. Kucheruk, and E. F. Pen, "Diffraction Efficiency of Partially Superimposed Holograms," *Avtometriya* 3 (1989), translated in *Optoelectronics, Instrumentation and Data Processing* 3 (1989).
114. A. A. Blok, "Effect of Data Coding Methods in Holographic Memory on Characteristics of Reconstructed Images of Data Pages," *Avtometriya* 5 (1989), translated in *Optoelectronics Instrumentation and Data Processing* 5 (1989).
115. V. A. Dombrovskii, S. A. Dombrovskii, and E. F. Pen, "Noise Immunity of Holographic Memory with Paraphrase Data Coding," *Avtometriya* 2 (1989), translated in *Optoelectronics, Instrumentation and Data Processing* 2 (1989).
116. P. E. Tverdokhlebl and B. N. Pankov, "Parallel Associative VLSI Processor with Optical Input," *SPIE Proc.* 1230 (1990).
117. S. Redfield and L. Hesselink, "Data Storage in Photorefractives Revisited," *SPIE, Optical Computing* 88 963 (1988).
118. S. Redfield and L. Hesselink, "Enhanced Nondestructive Holographic Readout in Strontium Barium Niobate," *Opt. Lett.* 13(10) (1988).
119. L. Hesselink and S. Redfield, "Photorefractive Holographic Recording in Strontium Barium Niobate Fibers," *Opt. Lett.* 13(10) (1988).

This page intentionally left blank.

DO NOT DUPLICATE

Howard Stark*

*Xerox Corporation
Corporate Research and Technology
Rochester, New York*

34.1 INTRODUCTION AND OVERVIEW

The xerographic process was invented in the 1930s by Chester Carlson, who was looking for a simple process to copy office documents. The process consists of the creation of an electrostatic image on an *image receptor*, development of the image with dyed or pigmented charged particles referred to as *toner*, transfer of the toner from the image receptor to the paper, fusing the toner to the paper, cleaning the residual toner from the image receptor, and finally, erasing whatever is left of the original electrostatic image. The process is then repeated on the cleaned, electrostatically uniform image receptor. In the most common embodiment of the process, the electrostatic image is created optically from either a digital or a light lens imaging system on a charged *photoreceptor*, a material that conducts electric charge in the light and is an insulator in the dark. These steps are shown schematically in Fig. 1, in which the photoreceptor drum is shown to be rotating clockwise. In this review we summarize the more common ways in which these steps are carried out.

The process just outlined is the heart of copying and digital printing systems whose speeds can range from a few to 180 copies per minute. Repeating the process several times (once for each color and black if needed) can produce full-color images. Often the system contains means for either input or output collation and stapling or binding. The cost of these systems can range from hundreds of dollars to several hundred thousand dollars.

This review will not attempt to give complete references to the technical literature. There are several excellent books that do this.¹⁻⁴ In addition, there are older books that give an interesting historical perspective on the development of the technology.^{5,6}

34.2 CREATION OF THE LATENT IMAGE

This section covers the creation of the electrostatic field image. First the more common optical systems are considered. Here, exposing a charged photoconductor to the optical image creates the latent electrostatic image. Then, ion writing systems, in which an insulator is charged imagewise with an ion writing head or bar to create the latent electrostatic image, are briefly discussed.

*Retired.

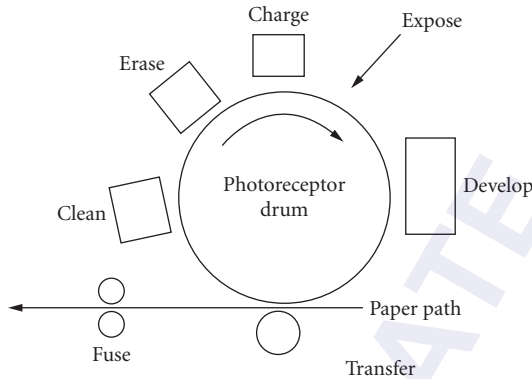


FIGURE 1 Schematic of xerographic process.

Optical Systems

We consider here ways in which the photoreceptor is charged, the required physical properties of the photoreceptor, and common exposure systems.

Charging Figure 2 schematically shows the charging and exposure of the photoreceptor. In this case the charging is shown to be positive. The two devices commonly used to charge the photoreceptor, the *corotron* and the *scorotron*, are shown schematically in Fig. 3. The corotron approximates a constant-current device, the scorotron a constant-voltage device.

The operational difference between the two devices is that the scorotron has a control screen. In both cases, a high potential, shown negative here, is applied to the corotron wires, creating a cloud of negative ions around the wires. In the case of the corotron, the negative ions drift under the influence of the electric field between the wires and the photoreceptor. Since the charging voltage of the photoreceptor is significantly less than that of the corona wires, the electric field and the resulting photoreceptor-charging current remain roughly constant. The charge voltage of the photoreceptor is then simply determined from the current per unit length of the corotron, the photoreceptor velocity under the corotron, and the capacitance per unit area of the photoreceptor.

In the case of the scorotron, the photoreceptor voltage and the voltage on the control grid determine the charging field. Thus, when the photoreceptor reaches the grid voltage, the field and the current go to zero. Hence the constant-voltage-like behavior.

Photoreceptor The discharge of the photoreceptor is accomplished by charge transport through the photoconductive medium. There are many materials that have been used as photoconductors.

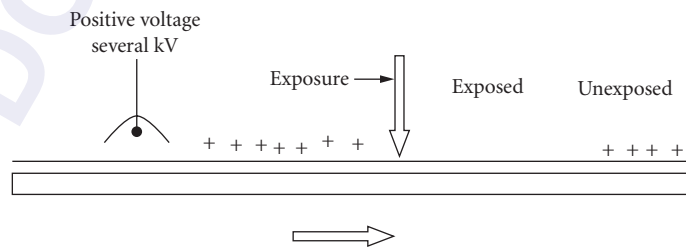


FIGURE 2 Charging and exposure of photoreceptor.

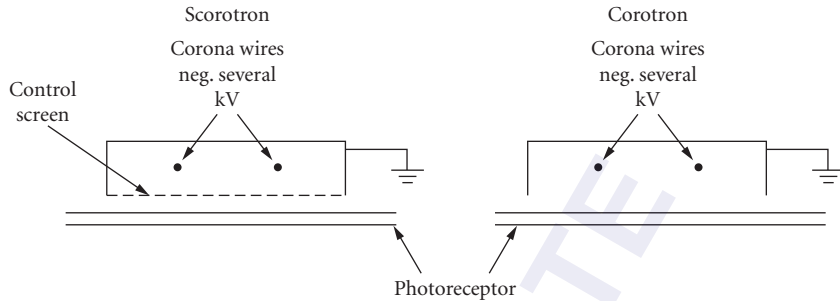


FIGURE 3 Scorotron and corotron.

The first photoreceptors were films 50 or 60 μm thick of amorphous selenium on a metallic substrate. These were followed by amorphous films of various selenium alloys, which were panchromatic and in some cases more robust. Other materials that have been used include amorphous silicon, which is quite robust as well as being panchromatic. Organic photoreceptors are used in most of the recent designs. Here the photoreceptor consists of a photogeneration layer on the order of 1 μm and a transport layer on the order of 20 μm thick.

The photoreceptor discharge process is shown in Fig. 4. The photoconductor is negatively charged, creating an electric field between the deposited charge and the ground plane. Light is shown to be incident on the generator layer. A hole is released that drifts upward under the influence of the electric field. Ideally the hole reaches the surface and neutralizes the applied surface charge. The electron that remains in the generator layer neutralizes the positive charge in the ground plane. Important characteristics of this process include the *dark decay* of the photoreceptor—how well the photoreceptor holds its charge in the dark—the quantum efficiency of the generation process, the transit time of the hole across the transport layer, whether or not it gets trapped in the process, and whether or not there are any residual fields remaining across the generator layer.

In order for the photoreceptor to hold its charge in the dark (Ref. 2, pp. 104–112), the charge on the surface must not be injected into the transport layer and drift to the substrate. There must be no bulk generation of charge in the transport layer. Finally, there must be no injection and transport of charge from the conductive ground plane into the transport layer. Modern photoreceptors dark-decay at rates of less than a few volts per second.

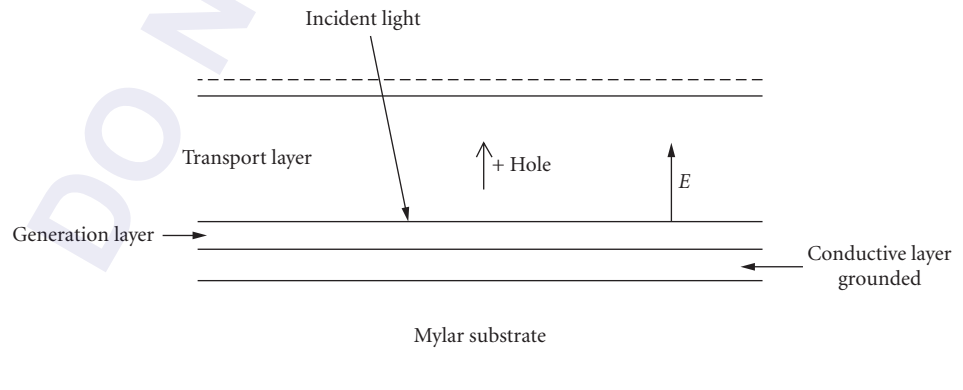


FIGURE 4 Photoreceptor discharge process.

The transport time of the photogenerated charge through the transport layer determines the rate at which the electrostatic latent image builds up. This limits the shortest time between exposure and development.

Charge trapped within the bulk of the photoreceptor can cause electrostatic ghosts of earlier images that may be developed. Proper erase procedures as well as careful photoreceptor processing are required to eliminate ghosts.

Exposure At present, both conventional light lens and digital exposure systems are in use. Conventional systems include both full-frame flash and slit scanning systems. Belt photoreceptors allow for a flat focal plane that permits a quite conventional exposure system; a full-frame flash exposure is used from a flat platen. More interesting is the system that is shown in Fig. 1. Here a slit is exposed at a fixed position on the rotating drum. To accommodate the movement of the drum, the image must move with it within the exposure slit. This is done with a moving platen for the document or a fixed platen with a moving exposure system and pivoting mirrors. Often a selfoc lens is used to conserve space.

The “original” in a digital imaging system is a document stored in computer memory. The idea includes both computer printers and digital copiers. The two most common means of optically writing the image on the photoreceptor are scanning lasers and image bars. In its simplest form an image bar exposes a line at a time across the photoreceptor. It consists of a full-width array of adjacent light-emitting diodes, one for each pixel. As the photoreceptor rotates under the image bar the diodes are turned on and off to write the image.

Laser scanning systems, also known as raster output scanning (ROS) systems, in their simplest embodiment use a laser diode that is focused and scanned across the photoreceptor by a rotating polygon. A so-called f - θ is used to achieve constant linear velocity of the spot across the photoreceptor. Often two or more diodes are focused several raster lines apart in order to write several lines at the same polygon speed. The laser diodes are modulated appropriately with the image information.

Prior to the development of laser diodes, HeNe lasers were used with acousto-optical modulators. In order to accommodate the slow response time of the acousto-optical crystal, the Scophony⁷ system developed in the 1930s was used. The acousto-optic modulator consists of a piezoelectric transducer, which launches an acoustical wave in a transparent medium whose index of refraction is pressure sensitive. The acoustic wave, which is modulated with the image information, creates a phase-modulated diffraction pattern. The laser beam is expanded, passed through the crystal and by a stop that blocks the zeroth order of the diffraction pattern. The image of the acousto-optic modulator is then focused on the photoreceptor. Because of the phased imaging system, the resulting image is intensity modulated. However, the diffraction pattern is moving and thus the pixels are moving on the photoreceptor surface. To compensate for this motion, the image of the modulator is scanned in the opposite direction by the polygon at precisely the same speed at which the pixels are moving.

Ion Writing Systems

In an ion writing system the electrostatic image is created by depositing ions on a dielectric receiver in an imagewise fashion. It is typically used in high-speed applications. The requirements for the dielectric receiver are that it be mechanically robust and that it hold the charge through the development process. Transfer, fusing, and cleaning are essentially the same as in conventional xerography. Since (in principle at least) the photoreceptor can be replaced by a more durable dielectric receiver and since charging is eliminated, the process promises to be cheaper and more robust.

At least two techniques have been used commercially for writing the image: stylus writing and ion writing heads. In both cases the limitation appears to be resolution. A stylus writing head consists of an array of styli, one for each pixel. The dielectric receiver is moved under the array and a voltage greater than air breakdown is applied to each stylus as appropriate. The ion writing heads are an array of ion guns, which uses an electrode to control the ion flow to the receiver.

34.3 DEVELOPMENT

The role of the developer system is to apply toner to the appropriate areas of the photoreceptor. In the case of conventional exposure, these areas are the charged areas of the photoreceptor. This system is referred to as *charged area development* (CAD). For digital systems where lasers or image bars are employed, the designer has a choice; the regions to be developed can be left charged as in the conventional system. Alternatively, the photoreceptor can be discharged in the regions to be toned. This is referred to as *discharged area development* (DAD). Image quality and reliability drive the choice. For either CAD or DAD, charged toner is electrostatically attracted to the photoreceptor.

There are many different techniques for developing the electrostatic latent image. We consider first two-component magnetic brush development and, in that context, outline many of the more general considerations for all development systems. Other interesting systems will then be described.

Two-Component Magnetic Brush Development

The developer in two-component magnetic brush development consists of magnetized carrier beads and toner. Here the toner is typically $10\ \mu\text{m}$ and the carrier 200 to $300\ \mu\text{m}$. The two components are mixed together and, by means of triboelectric charging, charge is exchanged between the toner and carrier. The much smaller toner particles remain attached to the carrier beads so that in a properly mixed developer there is little or no free toner. The role of the carrier is thus twofold: to charge the toner and, because of its magnetic properties, to enable the transport of the two-component developer. As will be seen, the conductivity of the carrier plays an important role in development. The carrier often consists of a ferrite core coated with a polymer chosen principally to control the charging characteristics of the developer.

The toner is a polymer containing pigment particles. For black systems the pigment is carbon black; for full color the subtractive primaries (cyan, magenta, and yellow) are used. In highlight color systems (black plus a highlight color) the pigment is the highlight color or perhaps a combination of pigments yielding the desired color. The choice of the polymer and, to some degree, the colorants is also constrained by the charging properties against the carrier and by the softening temperature, which is set by the fusing requirements. The covering power of the toner is determined by the concentration of the pigment. Typically a density of 1 is achieved on the order of $1\ \text{mg}/\text{cm}^2$ of toner.

A typical magnetic brush development system is shown schematically in Fig. 5.

The developer roll transports the developer (beads and carrier) from the sump to the nip between the developer roll and the photoreceptor where development takes place. The magnetic fields hold the developer on the roll and the material is moved along by friction. A carefully spaced doctor blade is used to control and limit the amount of developer on the roll.

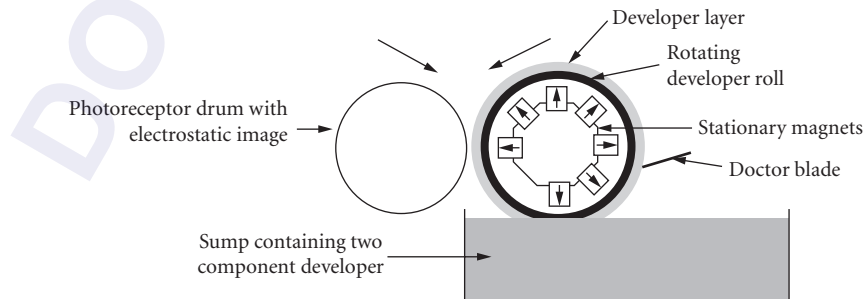


FIGURE 5 Magnetic brush development system.

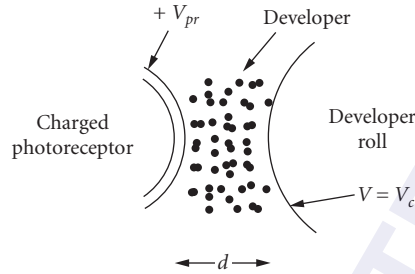


FIGURE 6 Details of nip between charged photoreceptor and developer roll.

Key to the development process is the triboelectric charge on the toner particle. The charge exchange between the toner and carrier can be thought of in terms of the alignment of the Fermi levels of the two materials in order to reach thermodynamic equilibrium. Thus a knowledge of the work functions (the energy required to lift an electron from the Fermi level to vacuum) gives a first-order estimate of the toner charge. Impurities and the pigments also play a large role in the triboelectric charging, as do the manufacturing processes. Often charge control agents are used to control the toner charging. The toner adheres to the carrier bead because of electrostatic attraction (they are of opposite polarity) and whatever other adhesion forces there are. Figure 6 schematically shows the details of the nip between the charged photoreceptor and the development roll.

The Development of Solid Areas The principal driving force for development is the electric field $E = (V_{pr} - V_c)/d$ where it is assumed that the developer is an insulator. (The conductive case will be considered presently.) V_{pr} is the voltage on the photoreceptor and V_c is a small voltage used to suppress development in background regions by reversing the field. The toner, however, is attached to the carrier beads and must be detached before it can be deposited on the photoreceptor. The electric field plays a role in this, as do the mechanical forces that result from impaction with the photoreceptor and the mixing of the developer within the nip. The density of the developer in the nip, the toner concentration, the magnetic field, and the electric field thus all affect the rate of toner deposition. Developed toner may also be scavenged from the photoreceptor by, say, an oppositely charged carrier bead impacting on the developed area. Development proceeds until the two rates are equal or until the photoreceptor emerges from the nip.

The voltage V_c is used to provide a reverse electric field to prevent toner deposition in what should be toner-free regions. This is required to prevent toner from adhering to the photoreceptor for nonelectrostatic reasons. Developers may contain a small amount of wrong-sign toner for which this field is a development field. This requires careful formulation of the developer and as well as judicious sizing of the cleaning field.

As development proceeds and toner is deposited on the photoreceptor, current flows from the developer roll to the photoreceptor, neutralizing the charge on the photoreceptor. This process may be viewed as the discharging of a capacitor—the photoreceptor—through a resistor—the developer. Thus, as a first-order approximation, the time constant for development is simply determined from the capacitance per unit area of the photoreceptor and the resistivity of the developer. The nip design must be such that the photoreceptor is in the nip on the order of a time constant or more. Typically development takes place to 50 percent or more of complete neutralization of the photoreceptor and is roughly a linear function of the development field until saturation is reached.

The resistivity of the developer plays a large role in the rate of development. Two cases may be considered: the insulating magnetic brush (IMB) and the conductive magnetic brush (CMB). In the case of conductive development (CMB) the effective spacing to the development electrode or roller is smaller than d (see Fig. 6), thereby increasing the apparent electric field and the rate of development. Ideally development proceeds to neutralization for CMB. If the resistivity of the developer is

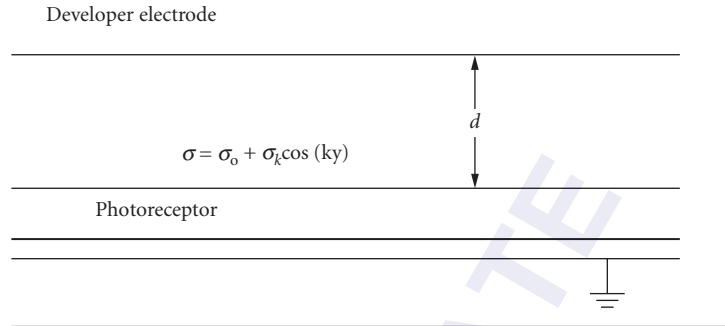


FIGURE 7 Calculating electric fields above a sinusoidal charge distribution.

large, the spacing is larger and the development is slower. In addition, space charge may develop in the nip, further slowing down development.

The Development of Lines For solid-area development with a highly resistive developer, the electric field that controls development is given by $E = (V_{pr} - V_c)/d$ (Fig. 6). For lines this is no longer true. It was recognized early on that lines develop much faster than solids due to the fact that the electric field at the edge of a large solid area is quite a bit stronger than at the center. This edge-enhancing effect was quite prominent in early xerographic development systems. A general approach to understanding line development is to calculate the modulation transfer function (MTF) or sine-wave response. It is relatively straightforward to calculate the electric fields above a sinusoidal charge distribution (Ref. 2, pp. 25–37), as shown in Fig. 7. The question is what field to use and whether or not a linear analysis is appropriate in what would appear to be a very nonlinear system.

As development proceeds, the fields decrease due to the neutralization of the charge on the photoreceptor. Furthermore, the fields fall off approximately exponentially in distance from the surface of the photoreceptor. Finally, space charge can accumulate in the developer nip; thus the assumption of a well-defined dielectric constant is questionable. Shown in Fig. 8 (taken from Ref. 2) is the normal component of the initial electric field 27 μm above the photoreceptor surface. The photoreceptor is 60 μm thick. The dielectric constant is 6.6, and the photoreceptor is charged to an average field of 15 $\text{V}/\mu\text{m}$. The dielectric constant of the nip is assumed to be 21 and the thickness of the nip is assumed to be 1700 μm . Here it is seen that, at least initially, lines with a spatial frequency of, say, 5 lines per mm develop at a rate 4 times faster than a solid area. If development is designed to go close to completion, this ratio can be much reduced.

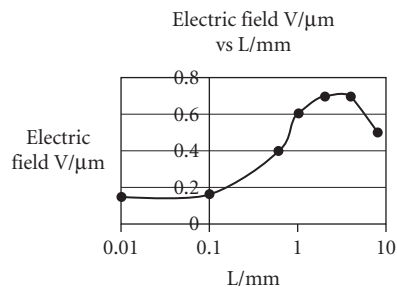


FIGURE 8 The normal electric field as a function of spatial frequency.

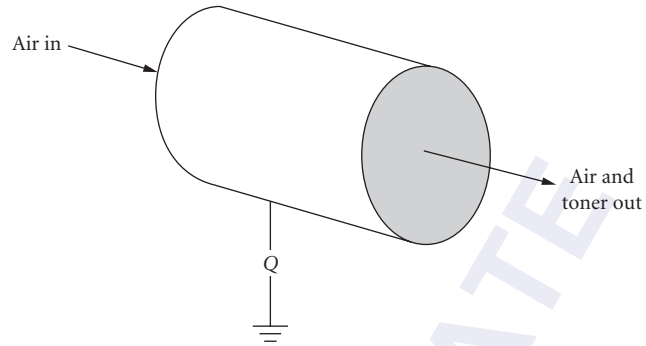


FIGURE 9 Faraday cage used to measure charge on toner.

Measuring Toner Charge

The measurement of the charge on the toner is fundamental to characterization of a developer, that is, a toner and carrier bead mix. A Faraday cage is used with screens on either end (Fig. 9). The screen mesh is such that toner can pass through and the carrier cannot.

The developer is loaded into the cage. A jet of air is blown through the cage, removing the toner. Both the charge and the weight of the removed toner are measured. The quotient is referred to as the *tribo* and is measured in units of microcoulombs per gram. The results depend on how the developer is mixed. Useful, properly mixed developers have tribos ranging between 10 and, say, $30 \mu\text{C/g}$.

The distribution of the charge can be obtained from what is called a *charge spectrograph*.⁸ (See Fig. 10.) The charged toner is blown off the developer mixture and inserted into the laminar air stream flowing in the tube. An electric field is applied normal to the toner flow. Within the tube the toner is entrained in the air and drifts transversely in the direction of the electric field. It is collected on the filter. The displacement of the toner d can be calculated from the charge on the toner Q , the

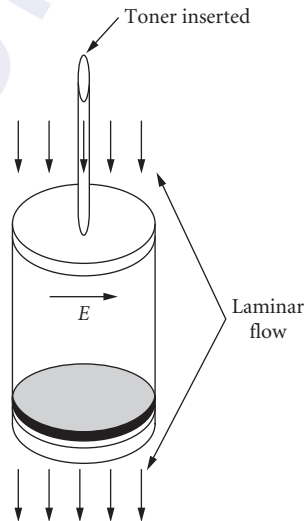


FIGURE 10 Charge spectrograph.

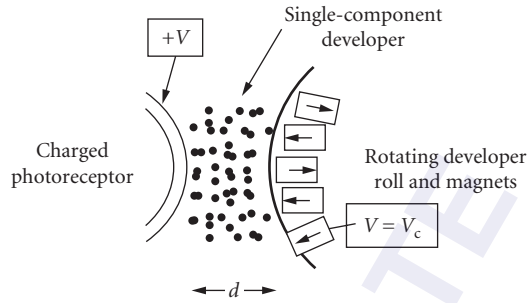


FIGURE 11 Single-component development.

electric field E , and the viscous drag, which is proportional to the radius r_t of the particle and the viscosity of the air η . Thus

$$d = (Q_t / r_t)(E / 6\pi\eta)$$

Using a computerized microscope to measure the number of toner particles as well as the radius and displacement, it is possible to obtain the charge and size distribution of the toner. This technique is particularly important as it yields the amount of wrong-sign toner in the developer.

Other Development Systems

Among the other useful development systems are inductively charged single-component development, powder cloud development, and electrophoretic or liquid immersion development (LID).

Single-Component Development In inductively charged single-component development the toner is both conductive and magnetic. The system is shown schematically in Fig. 11. The toner is charged inductively in response to the electric field with charge injected from the developer roll. The charge is transferred to the toner closest to the developer roll. The toner then is attracted to the photoreceptor. The materials issues are to ensure charge injection from the developer roll while at the same time preventing charge transfer from the toner to the photoreceptor.

Powder Cloud Development A powder cloud development system is shown in Fig. 12. Here toner is injected above the grid, drifts through the grid, acquires a negative charge, and is transported to the photoreceptor by the electric field. As is seen, at the edge of a charged area, the fields are such as to prevent toner from entering into the region near the edge, thus diminishing edge development.

Liquid Immersion Development An electrophoretic developer consists of toner particles suspended in a nonconducting fluid. Charge exchange takes place between the fluid and the toner particles.

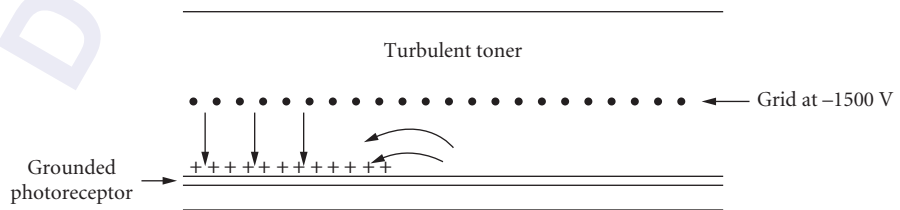


FIGURE 12 Powder cloud development system.

The charged toner then follows the field lines to the photoreceptor. The development rates can be determined from the toner charge and the image-driven electric field. In commercial systems, the developer fluid is pumped between the photoreceptor and a development roll. Controlling the toner charge is a major materials problem. The details of preparation play an important role. Surface-active agents also play an important role in the charging of the toner.

34.4 TRANSFER

After development, the toner is held on the photoreceptor by electrostatic forces. Electrostatic transfer is accomplished by applying an electric field either by means of a corotron or a bias transfer roll to attract the toner to the paper (Fig. 1). The paper is brought into contact with the image on the photoreceptor if nonelectrostatic forces are assumed negligible; it is possible to calculate⁹ where the toner splits as a function of the applied field and the thickness of the photoreceptor, toner layer, and paper. The nonelectrostatic forces, say, Van der Waals forces, between toner particles and between the toner and the photoreceptor can play an important role. Transfer efficiencies can run over 90 percent.

The difficult engineering problem is to bring the possibly charged paper into and out of contact with the photoreceptor without causing toner disturbances due possibly to air breakdown or premature transfer resulting in a loss of resolution. Bias transfer rolls with a semiconductive coating having carefully controlled relaxation times are required.

34.5 FUSING

After transfer the toner must be permanently fixed to the paper. This can be accomplished with the application of heat and possibly pressure. The idea is to have the toner flow together as well as into the paper. Surface tension and pressure play important roles. Many different types of fusing systems exist, the simplest of which is to pass the paper under a radiant heater. Here the optical absorption of the toner must be matched to the output of the lamp. This usually requires black toner.

Roll fusing systems are quite common. Here the paper is passed between two rolls, with the heated roll on the toner side. The important parameters are the roll temperature and dwell time of the image in the nip of the rollers. Release agents are used to assist the release of the paper from the rollers.

The fusing system imposes material constraints on the toner. Low-melt toners are preferred for fusing, but they cause the developer to age faster.

34.6 CLEANING AND ERASING

There are many ways of removing the untransferred toner from the photoreceptor in preparation for the next imaging cycle. Vacuum-aided brushes are common. Here a fur brush is rotated against the photoreceptor; the toner is removed from the photoreceptor by the brush and from the brush by the vacuum. These systems tend to be noisy because of the vacuum assist. Electrostatic brushes have also been used. Here a biased conductive brush removes the toner from the photoreceptor and then “develops” it onto a conductive roller, which in turn is cleaned with a blade. A development system biased to remove toner from the photoreceptor has also been used. The simplest of the cleaning systems is a blade cleaner; it is compact, quiet, and inexpensive.

Along with the removal of untransferred toner, the photoreceptor must be returned to a uniform and preferably uncharged state. Ghosting from the previous image may result from trapped charge within the photoreceptor. Erasing is accomplished using strong uniform exposure.

34.7 CONTROL SYSTEMS

Proper operation of the xerographic system requires system feedback control to maintain excellent image quality. Among the things to be controlled are charging, exposure, and toner concentration and development. At a minimum, toner gets used and must be replaced. The simplest control system counts pages, assumes area coverage, and replenishes the toner appropriately. The photoreceptor potential after charging and after exposure can be measured with an electrostatic voltmeter. Changes to the charging and exposures can then be appropriately made. The toner optical density of a developed patch on the photoreceptor of known voltage contrast can be measured with a densitometer. In digital systems the actual area coverage can be determined by counting pixels. These two measurements allow the control of toner concentration.

34.8 COLOR

There are many full-color and highlight-color copiers and printers available commercially. All of the recent designs are digital in that the image is either created on a computer or read in from an original document and stored in a digital format. Thus in a full-color process the half-toned cyan, magenta, yellow, and black separations are written on the photoreceptor by the ROS or image bar. The process is then essentially repeated for each separation.

Full Color

The xerographic processor can be configured in several different ways for full color. In a cyclic process a system like the one shown in Fig. 13 is used. The paper is attached to the bias transfer roll and the cycle is repeated four times for full color. The significant new issues are the need for four developer housings, registration, and the fusing of the thicker toner layers. The cyan image must be developed with just cyan. It is also important not to disturb the cyan image on the photoreceptor as it passes through the other housings. The developer housings are therefore mechanically or electrically switched in or out depending on the image written.

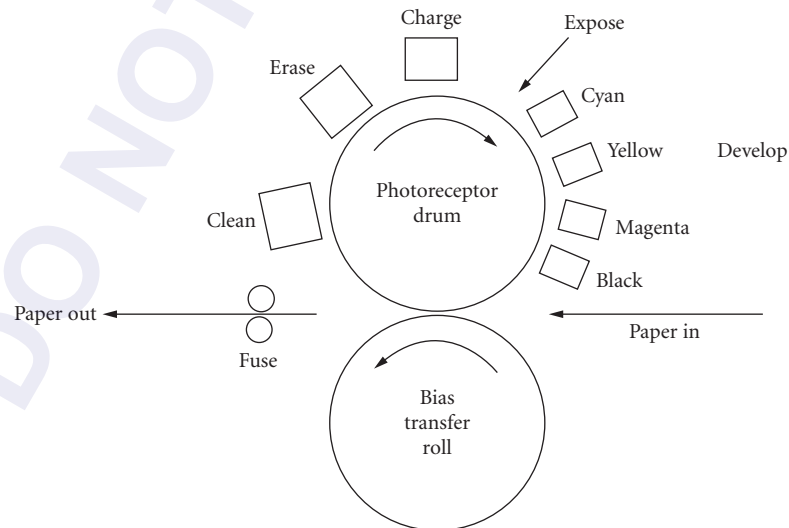


FIGURE 13 Cyclic full-color process.

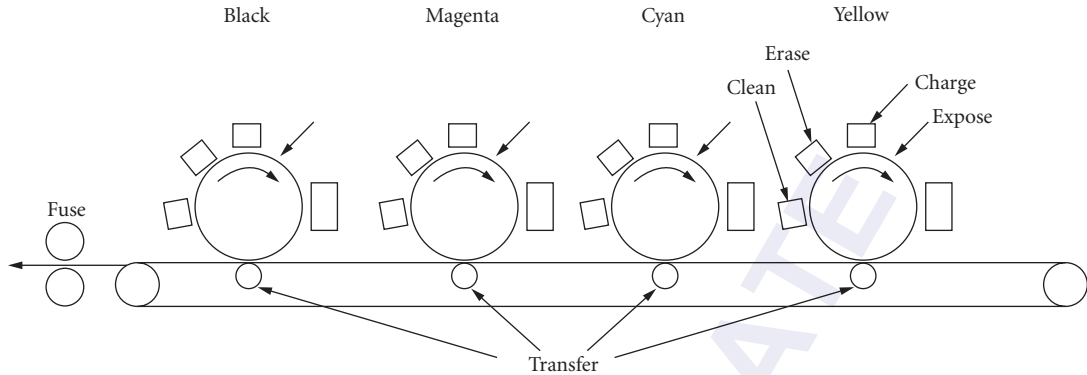


FIGURE 14 Tandem full-color process.

In Fig. 13 the transfer roll and the photoreceptor are the same diameter. Registration is accomplished by writing the image on the same place on the photoreceptor in each cycle. The temperature and dwell time in the fuser are controlled to achieve a good fuse.

Tandem configurations are also used. Here four separate processors are used as shown in Fig. 14. A belt transfer system is shown. The paper is attached to the belt and moved from station to station. The rollers shown behind the belt are used to apply the bias required for transfer. The order of development, shown here with yellow first, is chosen to minimize the effects of contamination. Other transfer systems are also possible. The image can be transferred directly to the belt and then, after the last station, to the paper.

Comparison of these two systems is interesting. The cyclic system has fewer parts and is therefore less expensive and likely more reliable. The tandem system has more parts but is a factor of 4 faster for full color than the cyclic system.

Highlight Color

Highlight color is black plus one color, say red, which is used to accentuate or highlight important parts of the image. It can be achieved with the full-color systems just discussed. Also, the full-color system can be modified to contain just the black plus the highlight color. A single-pass highlight color system was developed at Xerox that retains much of the simplicity and throughput of a single-color system. Referred to as *Tri-Level*,^{10,11} it encodes the black and highlight color on the photoreceptor as different voltages. The electrostatic arrangement is shown in Fig. 15. The photoreceptor is discharged to three levels. In this case full charge is black, the highlight is discharged to a low level, and white is maintained at some intermediate level.

Two developer housings are used, one for each color. In this case the black toner is negative and the color toner is positive. The housings are biased as shown. The image is passed sequentially through the two housings. The black region appears as a development field for black. Both the background and highlight color regions are cleaning fields for black and no development takes place. The same considerations apply to the highlight color. On the photoreceptor the resulting image contains opposite-polarity toner. The images are passed under a pretransfer corotron to reverse the sign of one of the toners. Thus, at the cost of an additional housing and a pretransfer charging device and no cost in throughput, a black-plus-one-color system can be designed.

The difficulties of the system are under consideration. Since the maximum charging voltage is limited, the voltage contrast is reduced by more than a factor of 2. The first image must not be disturbed when passing through the second developer housing. Both of these constraints require sophisticated developer housing design.

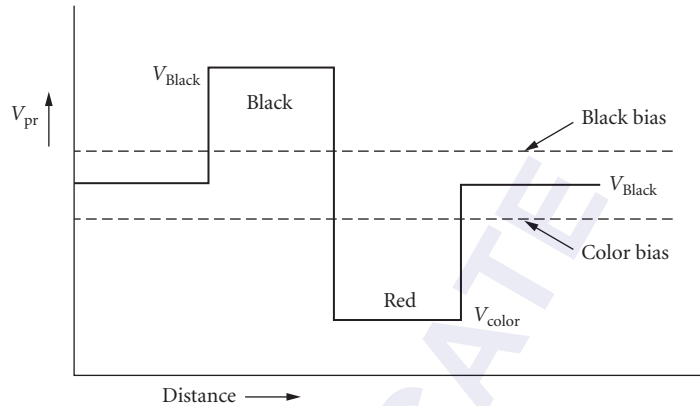


FIGURE 15 Tri-Level process.

34.9 REFERENCES

1. L. B. Shein, *Electrophotography and Development Physics*, 2d ed., Springer-Verlag, 1992.
2. M. Scharfe, *Electrophotography Principles and Optimization*, Research Studies Press, 1983.
3. E. M. Williams, *The Physics and Technology of the Xerographic Process*, Wiley, 1984.
4. D. M. Pai and B. E. Springett, "Physics of Electrophotography," *Rev. Modern Phys.* **65**:163–211 (January 1983).
5. J. Dessauer and H. Clark (eds.), *Xerography and Related Processes*, Focal Press Limited, 1965.
6. R. M. Schaffert, *Electrophotography*, Focal Press Limited, 1975.
7. L. M. Myers, "The Scophony System: An Analysis of Its Possibilities," *TV and Shortwave World*, 201–294 (April 1936).
8. R. B. Lewis, E. W. Connors, and R. F. Koehler, U.S. Patent 4,375,673, 1983.
9. C. C. Yang, and G. C. Hartman, *IEEE Trans. ED* **23**:308 (1976).
10. W. E. Haas, D. G. Parker, and H. M. Stark, "The Start of a New Machine, IS&T's Seventh International Congress on Advances in Non-Impact Printing," Portland, Oreg., October 6–11, 1991.
11. W. E. Haas, D. G. Parker, and H. M. Stark, "Highlight Color Printing: The Start of a New Machine," *J. Imag. Sci. Technol.* **36**(4):366 (July/August 1992).

This page intentionally left blank.

DO NOT DUPLICATE

PRINCIPLES OF OPTICAL DISK DATA STORAGE

Masud Mansuripur

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

35.1 INTRODUCTION

Since the early 1940s magnetic recording has been the mainstay of electronic information storage worldwide. Audiotapes provided the first major application for the storage of information on magnetic media. Magnetic tape has been used extensively in consumer products such as audiotapes and videocassette recorders (VCR); it has also found application in backup/archival storage of computer files, satellite images, medical records, etc. Large volumetric capacity and low cost are the hallmarks of tape data storage, although sequential access to the recorded information is perhaps the main drawback of this technology. Magnetic hard disk drives have been used as mass-storage devices in the computer industry ever since their inception in 1957. With an areal density that has doubled roughly every 2 years, hard disks have been and remain the medium of choice for secondary storage in computers.* Another magnetic storage device, the floppy disk, has been successful in areas where compactness, removability, and rapid access to the recorded information have been of primary concern. In addition to providing backup and safe storage, inexpensive floppies with their moderate capacities (2 Mbytes on a 3.5-in-diameter platter is typical) and reasonable transfer rates have provided the crucial function of file/data transfer between isolated machines. All in all, it has been a great half-century of progress and market dominance for magnetic storage which is only now beginning to face a serious challenge from the technology of optical recording.

Like magnetic recording, a major application of optical data storage is the secondary storage of information for computers and computerized systems. Like the high-end magnetic media, optical disks can provide recording densities in the range of 10^7 bits/cm² and beyond. The added advantage of optical recording is that, like floppies, these disks can be removed from the drive and stored on the shelf. Thus the functions of the hard disk (i.e., high capacity, high data transfer rate, rapid access) may be combined with those of the floppy (i.e., backup storage, removable media) in a single optical disk drive. Applications of optical recording are not confined to computer data storage. The enormously successful compact audio disk (CD) which was introduced in 1983 and has since

*Achievable densities on hard disks are presently in the range of 10^7 bits/cm²; random access to arbitrary blocks of data in these devices can take on the order of 10 ms, and individual read-write heads can transfer data at the rate of several megabits per second.

become the de facto standard of the music industry, is but one example of the tremendous potentials of the optical disk technology.

A strength of optical recording is that, unlike its magnetic counterpart, it can support read-only, write-once, and erasable/rewritable modes of data storage. Consider, for example, the technology of optical audio/video disks. Here the information is recorded on a master disk which is then used as a stamper to transfer the embossed patterns to a plastic substrate for rapid, accurate, and inexpensive reproduction. The same process is employed in the mass production of read-only files (CD-ROM, O-ROM) which are now being used to distribute software, catalogs, and other large databases. Or consider the write-once-read-many (WORM) technology, where one can permanently store massive amounts of information on a given medium and have rapid, random access to them afterward. The optical drive can be designed to handle read-only, WORM, and erasable media all in one unit, thus combining their useful features without sacrificing performance and ease of use. Moreover, the media can contain regions with prerecorded information as well as regions for read/write/erase operations on the same platter, thus offering opportunities for applications that have heretofore been unthinkable.

This chapter presents the conceptual basis for optical storage systems, with emphasis on disk technology in general and magneto-optical (MO) disk in particular. Section 35.2 is devoted to a discussion of some elementary aspects of disk data storage including the concept of track, definition of the access time, and the physical layout of data. Section 35.3 describes the function of the optical path; included are properties of the semiconductor laser diode, characteristics of the beam-shaping optics, and features of the focusing (objective) lens. The limited depth of focus of the objective lens and the eccentricity of tracks dictate that optical disk systems utilize closed-loop feedback mechanisms for maintaining the focused light spot on the right track at all times. Automatic focusing and automatic track-following schemes are described in Secs. 35.4 and 35.5. The physical process of thermomagnetic recording is the subject of Sec. 35.6, followed by a discussion of MO readout in Sec. 35.7. Certain important characteristics of MO media are summarized in Sec. 35.8. Concluding remarks and an examination of trends for future optical recording devices are the subject of Sec. 35.9.

Alternative methods of optical data storage such as reversible phase-change, photochemical spectral hole burning, three-dimensional volume holographic storage, photon echo, photon trapping, etc., will not be discussed in this chapter. The interested reader may consult the following references for information concerning these alternative storage schemes:

Proceedings of the International Symposium on Optical Memory, ISOM'89, published as supplement 28-3 of the *Japanese Journal of Applied Physics*, vol. 28 (1989).

Proceedings of the Optical Data Storage Conference, SPIE, vol. 1316 (1990).

Proceedings of the Optical Data Storage Conference, SPIE, vol. 1499 (1991).

Proceedings of the Optical Data Storage Conference, SPIE, vol. 1663 (1992).

R. G. Zech, "Volume Hologram Optical Memories: Mass Storage Future Perfect," *Optics and Photonics News*, vol. 3, no. 8, pp. 16-25 (1992).

35.2 PRELIMINARIES AND BASIC DEFINITIONS

The format and physical layout of recorded data on the storage medium as well as certain operational aspects of disk drive mechanism will be described in the present section.

The Concept of Track

The information on magnetic and optical disks is recorded along tracks. Typically, a track is a narrow annulus at some distance r from the disk center, as shown in Fig. 1. The width of the annulus is denoted by W_t , while the width of the guard band, if any, between adjacent tracks is denoted by W_g .

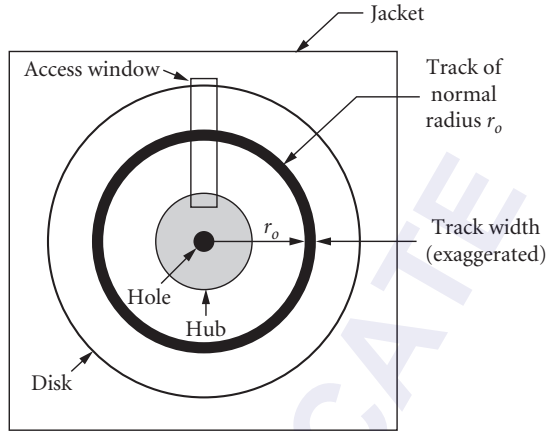


FIGURE 1 Physical appearance and general features of an optical disk. The read-write head gains access to the disk through a window in the jacket; the jacket itself is for protection purposes only. The hub is the mechanical interface with the drive for mounting and centering the disk on the spindle. The track shown here is of the concentric-ring type, with radius r_o and width W_t .

The track-pitch is the center-to-center distance between neighboring tracks and is therefore equal to $W_t + W_g$. A major difference between the magnetic floppy disk, the magnetic hard disk, and the optical disk is that their respective track-pitches are presently of the order of 100, 10, and $1 \mu\text{m}$. Tracks may be fictitious entities, in the sense that no independent existence outside the pattern of recorded marks may be ascribed to them. This is the case, for example, with the compact audio disk format where prerecorded marks simply define their own tracks and help guide the laser beam during read-out. In the other extreme are tracks that are physically engraved on the disk surface before any data is ever recorded. Examples of this type of track are provided by pregrooved WORM and magneto-optical disks. Figure 2 shows micrographs from several recorded optical disk surfaces. The tracks along which data is written are clearly visible in these pictures.

It is generally desired to keep the read-write head stationary while the disk spins and a given track is being read from or written onto. Thus, in an ideal situation, not only should the track be perfectly circular, but also the disk must be precisely centered on the spindle axis. In practical systems, however, tracks are neither precisely circular, nor are they concentric with the spindle axis. These eccentricity problems are solved in low-performance floppy drives by making tracks wide enough to provide tolerance for misregistrations and misalignments. Thus the head moves blindly to a radius where the track center is nominally expected to be, and stays put until the reading or writing is over. By making the head narrower than the track-pitch, the track center is allowed to wobble around its nominal position without significantly degrading the performance during read-write operations. This kind of wobble, however, is unacceptable in optical disk systems which have a very narrow track, about the same size as the focused beam spot. In a typical situation arising in practice the eccentricity of a given track may be as much as $50 \mu\text{m}$, while the track-pitch is only about $1 \mu\text{m}$, thus requiring active track-following procedures.

A popular method of defining tracks on an optical disk is by means of pregrooves, which are either etched, stamped, or molded onto the substrate. The space between neighboring grooves is called *land* (see Fig. 3a). Data may be written in the groove with the land acting as a guard band. Alternatively, the land may be used for recording while the grooves separate adjacent tracks. The groove depth is optimized for generating an optical signal sensitive to the radial position of the

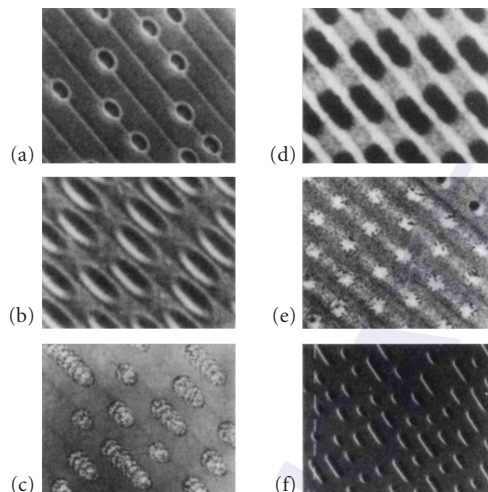


FIGURE 2 Micrographs of several types of optical storage media. The tracks are straight and narrow with a $1.6\text{-}\mu\text{m}$ pitch, and are diagonally oriented in each frame. (a) Ablative, write-once tellurium alloy. (b) Ablative, write-once organic dye. (c) Amorphous-to-crystalline, write-once phase-change alloy GaSb. (d) Erasable, amorphous magneto-optic alloy GdTbFe. (e) Erasable, crystalline-to-amorphous phase-change tellurium alloy. (f) Read-only CD-Audio, injection-molded from polycarbonate with a nickel stamper. (From Ullmann's "Encyclopedia of Industrial Chemistry," Verlagsgesellschaft, mbH, Weinheim, 1989.)

read-write laser beam. For the push-pull method of track-error detection (described in Sec. 35.5) the groove depth is in the neighborhood of $\lambda/8$, where λ is the wavelength of the light beam.

In digital data storage each track is divided into small segments called *sectors*. A sector is intended for the storage of a single block of data which is typically either 512 or 1024 bytes. The physical length of a sector is thus several millimeters. Each sector is preceded by header information such as the identity of the sector, identity of the corresponding track, synchronization marks, etc. The header information may be preformatted onto the substrate, or it may be written directly on the storage layer. Pregrooved tracks may be "carved" on the optical disk either as concentric rings or as a single continuous spiral. There are certain advantages to each format. A spiral track contains a succession of sectors without interruption, whereas concentric rings may each end up with some empty space that is too small to become a sector. Also, large files may be written onto (and read from) spiral tracks without jumping to the next track, which is something that occurs when concentric tracks are used. On the other hand, multiple-path operations such as write-and-verify or erase-and-write which require two paths each for a given sector, or still-frame video are more conveniently handled on concentric-ring tracks.

Another suggested track format is based on the idea of a sampling servo. Here the tracks are identified by occasional marks placed permanently on the substrate at regular intervals, as shown in Fig. 3b. Details of track-following by the sampled-servo scheme will follow shortly (see Sec. 35.5), suffice it to say at this point that servo marks help the system identify the position of the focused spot relative to the track center. Once the position is determined it is fairly simple to steer the beam and adjust its position on the track.

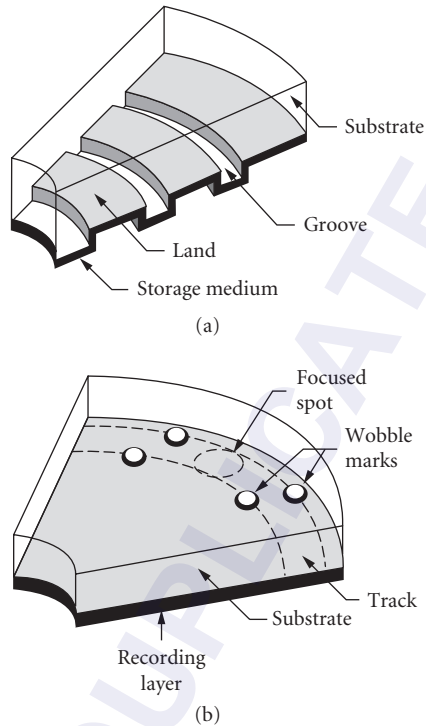


FIGURE 3 (a) Lands and grooves in an optical disk. The substrate is transparent, and the laser beam must pass through it before reaching the storage medium. (b) Sampled-servo marks in an optical disk. These marks which are offset from the track center provide information regarding the position of focused spot.

Disk Rotation Speed

When a disk rotates at a constant angular velocity ω , a track of radius r moves with the constant linear velocity $V = r\omega$. Ideally, one would like to have the same linear velocity for all the tracks, but this is impractical except in a limited number of situations. For instance, when the desired mode of access to the various tracks is sequential, such as in audio- and video-disk applications, it is possible to place the head in the beginning at the inner radius and move outward from the center thereafter while continuously decreasing the angular velocity. By keeping the product of r and ω constant, one can achieve constant linear velocity for all tracks.* Sequential access mode, however, is the exception rather than the norm in data storage systems. In most applications, the tracks are accessed randomly with such rapidity that it becomes impossible to adjust the rotation speed for constant linear velocity. Under these circumstances the angular velocity is kept constant during normal operation. Typical

*In compact audio disk players the linear velocity is kept constant at 1.2 m/s. The starting position of the head is at the inner radius $r_{\min} = 25$ mm, where the disk spins at 460 rpm. The spiral track ends at the outer radius $r_{\max} = 58$ mm, where the disk's angular velocity is 200 rpm.

rotation rates are 1200 and 1800 rpm for slower drives, and 3600 rpm for the high-end systems. Higher rotation rates (5000 rpm and beyond) are certainly feasible and will likely appear in future generations of optical storage devices.

Access Time

The direct access storage device used in computer systems for the mass storage of digital information is a disk drive capable of storing large quantities of data and accessing blocks of this data rapidly and in random order. In read-write operations it is often necessary to move the head to new locations in search of sectors containing specific data items. Such random relocations are usually time-consuming and can become the factor that limits performance in certain applications. The access time τ_a is defined as the average time spent in going from one randomly selected spot on the disk to another. The access time τ_a can be considered the sum of seek time τ_s , which is the average time needed to acquire the target track, and a latency τ_l , which is the average time spent on the target track waiting for the desired sector; thus $\tau_a = \tau_s + \tau_l$. The latency is half the revolution period of the disk, since a randomly selected sector is, on the average, halfway along the track from the point where the head initially lands. Thus, for a disk rotating at 1200 rpm $\tau_l = 25$ ms, while at 3600 rpm $\tau_l \approx 8.3$ ms. The seek time, on the other hand, is independent of the rotation speed, but is determined by the travel distance of the head during an average seek, as well as by the mechanism of head actuation. (It can be shown that the average length of travel in a random seek is one-third of the full stroke.) In magnetic disk drives where the head/actuator assembly is relatively lightweight, (a typical Winchester head weighs about 5 g) the acceleration and deceleration periods are short, and seek times are typically around 10 ms. In optical disk systems, on the other hand, the head, being an assembly of discrete elements, is fairly large and heavy (typical weight ≈ 50 to 100 g), resulting in values of τ_s that are several times greater than those obtained in magnetic recording. The seek times reported for commercially available optical drives presently range from 20 msec in high-performance 3.5-in drives to 100 ms in larger drives. One must emphasize, however, that the optical disk technology is still in its infancy; with the passage of time the integration and miniaturization of the elements within the optical head will surely produce lightweight devices capable of achieving seek times in the range of several milliseconds.

Organization of Data on Disk

For applications involving computer files and data, each track is divided into a number of sectors where each sector can store a fixed-length block of binary data. The size of the block varies among the various disk/drive manufacturers, but typically it is either 512 or 1024 bytes. As long as the disk is dedicated to a particular drive (such as in magnetic hard drives) the sector size is of little importance to the outside world. However, with removable media the sector size (among other things) must be standardized, since now various drives need to read from and write onto the same disk.

A block of user data cannot be directly recorded on a sector. First, it must be coded for protection against errors (error-correction coding) and for the satisfaction of channel requirements (modulation coding). Also, it may be necessary to add synchronization bits or other kinds of information to the data before recording. Thus a sector's capacity must be somewhat greater than the amount of raw data assigned to it. A sector also must have room for "header" information. The header is either recorded during the first use of the disk by the user, as in formatting a floppy disk, or is written by the manufacturer before shipping the disk. The header typically contains the address of the sector plus synchronization and servo bits. In magnetic disks the header is recorded magnetically, which makes it erasable and provides the option of reformatting at later times. On the negative side, formatting is time-consuming and the information is subject to accidental erasure. In contrast, the optical disk's sector headers may be mass-produced from a master at the time of manufacture, thus eliminating the slow process of soft formatting. The additional space used by the codes and by the header information constitutes the overhead. Depending on the quality of the disk, the degree of

sophistication of the drive, and the particular needs of a given application, the overhead may take as little as 10 percent and as much as 30 percent of a disk's raw capacity.

35.3 THE OPTICAL PATH

The optical path begins at the light source which, in all laser disk systems in use today, is a semiconductor GaAs diode laser. Several unique features of the laser diode have made it indispensable for optical recording applications: its small size ($\approx 300 \times 50 \times 10 \mu\text{m}$) makes possible the construction of compact head assemblies, its coherence properties allow diffraction-limited focusing to extremely small spots, and its direct modulation capability eliminates the need for external intensity modulators. The operating wavelength of the laser diode can be selected within a limited range by proper choice of material composition; presently, the shortest wavelength available from the III-V class of semiconductor materials is 670 nm.

Figure 4a shows a typical plot of laser power output versus input current for a GaAs-based laser diode. The lasing starts at the threshold current, and the output power rapidly increases beyond that

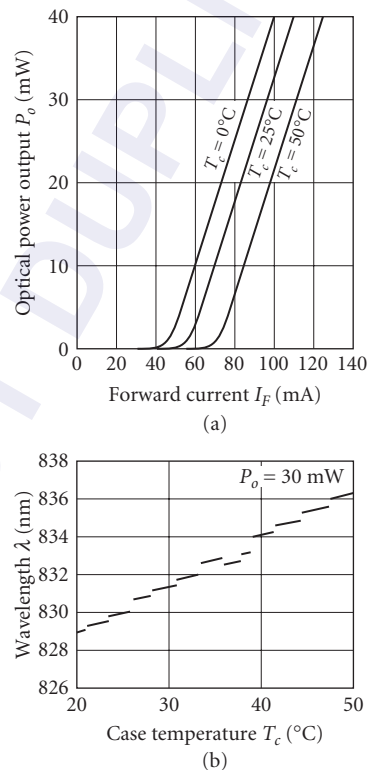


FIGURE 4 (a) Optical output power versus forward-bias current for a typical diode laser. Different curves were obtained at different ambient temperatures. (b) Variations of wavelength as function of case temperature for typical diode laser. The output power is fixed at $P_o = 30 \text{ mW}$. (From Sharp Laser Diode User's Manual.)

point. Below threshold, the diode operates in the spontaneous emission mode and its output is incoherent. After threshold, stimulated emission takes place, yielding coherent radiation. Of course, the output power cannot increase indefinitely and beyond a certain point the laser fails catastrophically. Fortunately, the required optical power levels for the read/write/erase operations in present-day data storage systems are well below the failure levels of these lasers. Available lasers for data storage applications have threshold currents around 40 mA, maximum allowable currents of about 100 mA, and peak output powers [CW (continuous wave) mode] around 50 mW. The relationship between the injection current and the output light power is very sensitive to the operating temperature of the laser, as evidenced by the various plots in Fig. 4*a*. Also, because the semiconductor material's band-gap is a function of the ambient temperature, there is a small shift in the operating wavelength of the device when the temperature fluctuates (see Fig. 4*b*). For best performance it is usually necessary to mount the laser on a good heat-sink, or try to steady its temperature by closed-loop feedback.

The output optical power of the laser can be modulated by controlling the injection current. One can apply pulses of variable duration to turn the laser on and off during the recording process. The pulse duration can be as short as a few nanoseconds, with rise and fall times which are typically less than 1 ns. This direct-modulation capability of the laser diode is particularly welcome in optical disk systems, considering that most other sources of coherent light (such as gas lasers) require bulky and expensive devices for external modulation. Although readout of optical disks can be accomplished at constant power level in CW mode, it is customary (for noise reduction purposes) to modulate the laser at a high frequency in the range of several hundred MHz.

Collimation and Beam Shaping

Since the cross-sectional area of the active region in a laser diode is only about $1 \mu\text{m}^2$, diffraction effects cause the emerging beam to diverge rapidly. This phenomenon is depicted schematically in Fig. 5*a*. In practical applications of the laser diode, the expansion of the emerging beam is arrested by a collimating lens, such as that shown in Fig. 5*b*. If the beam happens to have aberrations (astigmatism is particularly severe in diode lasers), then the collimating lens must be designed to correct this defect as well.

In optical recording it is most desirable to have a beam with circular cross section. The need for beam shaping arises from the special geometry of the laser cavity with its rectangular cross section. Since the emerging beam has different dimensions in the directions parallel and perpendicular to the junction, its cross section at the collimator becomes elliptical, with the initially narrow dimension expanding more rapidly to become the major axis of the ellipse. The collimating lens thus produces a beam with elliptical cross section. Circularization may be achieved by bending various rays of the beam at a prism, as shown in Fig. 5*c*. The bending changes the beam's diameter in the plane of incidence, but leaves its diameter in the perpendicular direction intact.

The output of the laser diode is linearly polarized in the plane of the junction. In some applications (such as readout of compact disks or read-write on WORM media) the polarization state is immaterial as far as interaction with the storage medium is concerned. In such applications one usually passes the beam through a polarizing beam splitter (PBS) and a quarter-wave plate, as in Fig. 6, and converts its polarization to circular. Upon reflection from the disk, the beam passes through the quarter-wave plate once again, but this time emerges as linearly polarized in a direction perpendicular to the original direction of polarization. The returning beam is thus directed away from the laser and toward the detection module, where its data content is extracted and its phase/amplitude pattern is used to generate error signals for automatic focusing and tracking. By thoroughly separating the returning beam from the incident beam, one not only achieves efficiency in the use of the optical power, but also succeeds in preventing the beam from going back to the laser where it causes instabilities in the laser cavity and, subsequently, increases the noise level. Unfortunately, there are situations where a specific polarization state is required for interaction with the disk; magneto-optical readout which requires linear polarization is a case in point. In such instances the simple combination of PBS and quarter-wave plate becomes inadequate and one must resort to other (less efficient) means of separating the beams.

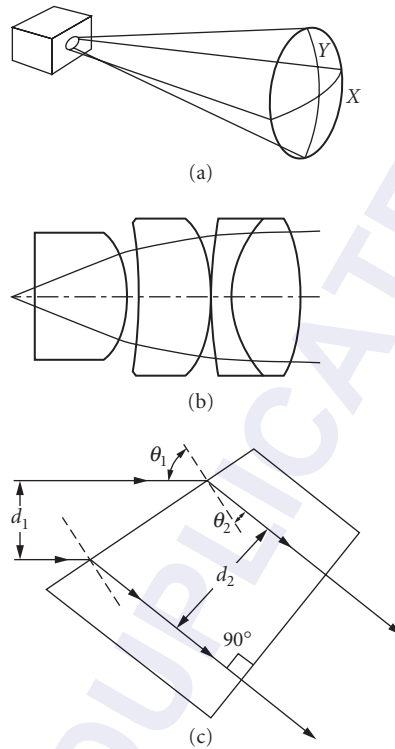


FIGURE 5 (a) Away from the facet, the output beam of a diode laser diverges rapidly. In general, the beam diameter along X is different from that along Y , which makes the cross section of the beam elliptical. Also, the radii of curvature R_x and R_y are not the same, thus creating a certain amount of astigmatism in the beam. (b) Multi-element collimator lens for laser diode applications. Aside from collimating, this lens also corrects astigmatic aberrations of the beam. (c) Beam-shaping by deflection at a prism surface. Θ_1 and Θ_2 are related by the Snell's law, and the ratio d_2/d_1 is the same as $\cos \Theta_2 / \cos \Theta_1$. Passage through the prism circularizes the elliptical cross section of the beam.

Focusing

The collimated and circularized beam of the laser is focused on the surface of the disk using an objective lens. The objective is designed to be aberration-free, so that its focused spot size is limited only by the effects of diffraction. Figure 7a shows the design of a typical objective made from spherical optics. According to the classical theory of diffraction, the diameter of the beam, d , at the objective's focal plane is

$$d \approx \frac{\lambda}{\text{NA}} \quad (1)$$

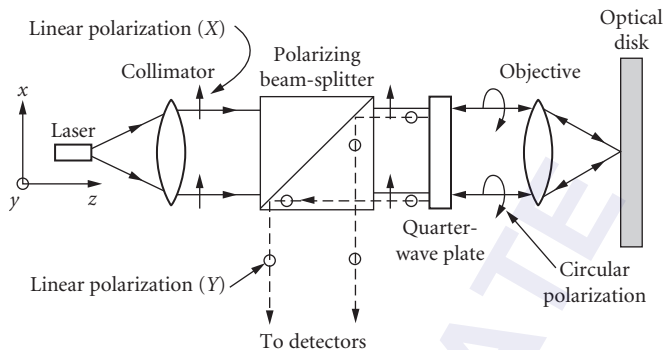


FIGURE 6 Separation of incident and reflected beams at the polarizing beam splitter (PBS). The quarter-wave plate converts the linearly polarized incident beam into one with circular polarization and converts the returning beam back to linear, but with its polarization vector orthogonal to that of the incident beam. This 90° rotation of polarization is responsible for the diversion of the reflected beam toward the detection channel.

where λ is the wavelength of light and NA is the numerical aperture of the objective. In optical recording it is desired to achieve the smallest possible spot, since the size of the spot is directly related to the size of marks recorded on the medium. Also, in readout, the spot size determines the resolution of the system. According to Eq. (1) there are two ways to achieve a small spot: reducing the wavelength and increasing the numerical aperture. The wavelengths currently available from GaAs lasers are in the range of 670 to 840 nm. It is possible to use a nonlinear optical device to double the frequency of these lasers, thus achieving blue light. Good efficiencies have been demonstrated by frequency doubling. Also recent developments in II-VI materials have improved the prospects for obtaining green and blue light directly from semiconductor lasers. Consequently, there is hope that in the near future optical storage systems will operate in the wavelength range of 400 to 500 nm. As for the numerical aperture, current practice is to use a lens with $NA \approx 0.5\text{--}0.6$. Although this value might increase slightly in the coming years, much higher numerical apertures are unlikely, since they put strict constraints on the other characteristics of the system and limit the tolerances. For instance, the working distance at high NA is relatively short, making access to the recording layer through the substrate more difficult. The smaller depth of focus of a high-NA lens will make attaining/maintaining proper focus more difficult, while the limited field of view might restrict automatic track-following procedures. A small field of view also places constraints on the possibility of read/write/erase operations involving multiple beams.

The depth of focus of a lens, δ , is the distance away from the focal plane over which tight focus can be maintained (see Fig. 7b). According to the classical theory of diffraction,

$$\delta \approx \frac{\lambda}{NA^2} \quad (2)$$

Thus for $\lambda = 700$ nm and $NA = 0.6$ the depth of focus is about ± 1 μm . As the disk spins under the optical head at the rate of several thousand rpm, the objective must stay within a distance of $f \pm \delta$ from the active layer if proper focus is to be maintained. Given the conditions under which drives usually operate, it is impossible to make rigid enough mechanical systems to yield the required positioning tolerances. On the other hand, it is fairly simple to mount the objective lens in an actuator capable of adjusting its position with the aid of closed-loop feedback control. We emphasize that by going to shorter wavelengths and/or larger numerical apertures (as is required for attaining higher data densities) one will have to face a much stricter regime as far as automatic focusing is concerned. Increasing the numerical aperture is particularly worrisome, since δ drops with the square of NA.

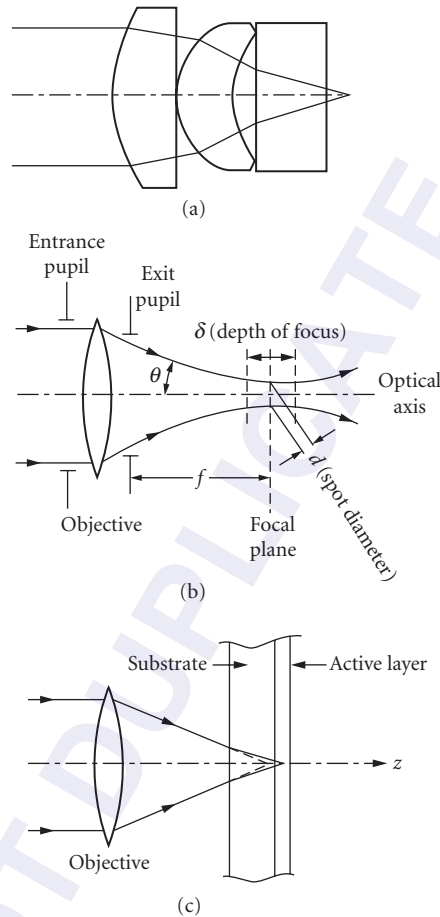


FIGURE 7 (a) Multielement lens design for a high-NA videodisc objective. (After D. Kuntz, "Specifying Laser Diode Optics," *Laser Focus*, March 1984.) (b) Various parameters of the objective lens. The numerical aperture is $NA = \sin \Theta$. The spot diameter d and the depth of focus δ are given by Eqs. (1) and (2), respectively. (c) Focusing through the substrate can cause spherical aberration at the active layer. The problem is corrected by a proper design for the objective lens, which takes the substrate into account.

A source of spherical aberrations in optical disk systems is the substrate through which the light must pass in order to reach the active layer. Figure 7c shows the bending of the rays at the surface of the disk, which causes the aberration. This problem can be solved by taking into account the effects of the substrate in the design of the objective, so that the lens is corrected for all aberrations, including those arising at the substrate. Recent developments in molding of aspheric glass lenses have gone a long way in simplifying the lens design problem. Figure 8 shows a pair of molded glass aspherics

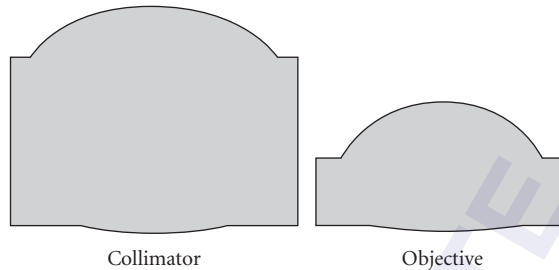


FIGURE 8 Molded glass aspheric lens pair for optical disk application. These singlets can replace the multielement spherical lenses shown in Figs. 5*b* and 7*a*.

designed for optical storage applications; both the collimator and the objective are single-element lenses and are corrected for axial aberrations.

Laser Noise

Compared to other sources of coherent light such as gas lasers, laser diodes are noisy and unstable. Typically, within a diode laser's cavity several modes compete for dominance. Under these circumstances, small variations in the environment can cause mode-hopping which results in unpredictable power-level fluctuations and wavelength shifts. Unwanted optical feedback is specially troublesome, as even a small fraction of light returning to the cavity can cause a significant rise in the noise level. Fortunately, it has been found that high-frequency modulation of the injection current can be used to instigate power sharing among the modes and thereby reduce fluctuations of the output optical power. In general, a combination of efforts such as temperature stabilization of the laser, antireflection coating of the various surfaces within the system, optical isolation of the laser, and high-frequency modulation of the injection current can yield acceptable levels of noise for practical operation of the device.

35.4 AUTOMATIC FOCUSING

Since the objective lens has a large numerical aperture ($NA \geq 0.5$) its depth of focus δ is shallow ($\delta \approx \pm 1 \mu\text{m}$ at $\lambda = 780 \text{ nm}$). During all read/write/erase operations, therefore, the disk must remain within a fraction of a micrometer from the focal plane of the objective. In practice, however, the disks are not flat and are not always mounted rigidly parallel to the focal plane, so that during any given revolution movements away from focus (by as much as $\pm 50 \mu\text{m}$) may occur. Without automatic adjustment of the objective along the optic axis, this runout (or disk flutter) will be detrimental to the operation of the system. In practice, the objective is mounted on a small actuator (usually a voice coil) and allowed to move back and forth to keep its distance from the disk within an acceptable range. Since the spindle turns at a few thousand rpm, if the disk moves in and out of focus a few times during each revolution, then the voice coil must be fast enough to follow these movements in real time; in other words, its frequency response must extend from DC to several kHz.

The signal that controls the voice coil is obtained from the light reflected from the disk. There are several techniques for deriving the focus error signal (FES), one of which is depicted in Fig. 9*a*. In this so-called obscuration method a secondary lens with one-half of its aperture covered is placed in the path of the reflected light, and a split-detector is placed at the focal plane of this secondary lens. When the disk is in focus, the returning beam is collimated and the secondary lens will focus the

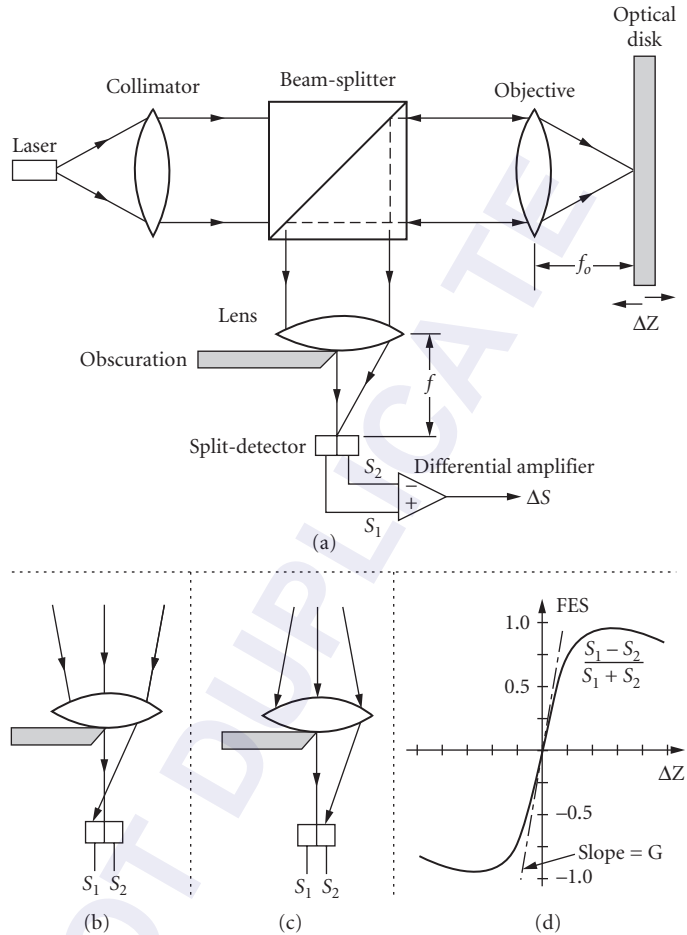


FIGURE 9 Focus error detection by the obscuration method. In (a) the disk is in focus, and the two halves of the split detector receive equal amounts of light. When the disk is too far from the objective (b) or too close to it (c), the balance of detector signals shifts to one side or the other. A plot of the focus error signal versus defocus is shown in (d), and its slope near the origin is identified as the FES gain, G .

beam at the center of the split-detector, giving a difference signal ΔS equal to zero. If the disk now moves away from the objective, the returning beam will become converging, as in Fig. 9b, sending all the light to detector 1. In this case ΔS will be positive and the voice coil will push the lens toward the disk. On the other hand, when the disk moves close to the objective, the returning beam becomes diverging and detector 2 receives the light (see Fig. 9c). This results in a negative ΔS which forces the voice coil to pull back and return ΔS to zero.

A given focus error detection scheme is generally characterized by the shape of its focus error signal ΔS versus the amount of defocus ΔZ . One such curve is shown in Fig. 9d. The slope of the FES curve near the origin is of particular importance, since it determines the overall performance and stability of the servo loop. In general, schemes with a large slope are preferred, although certain

other aspects of system performance should also be taken into consideration. For instance, variations of the FES during seek operations (where multiple track-crossings occur) should be kept at a minimum, or else the resulting “feedthrough” might destabilize the focus servo. Also, it is important for a focus-error-detection scheme to be insensitive to slight imperfections of the optical elements, as well as to the positioning and mechanical misalignments; otherwise, the manufacturing cost of the device may become prohibitive. Finally, the focusing scheme must have a reasonable acquisition range, so that at start-up (or in those occasions where focus is lost and needs to be acquired again) the system can move in the proper direction to establish focus.

35.5 AUTOMATIC TRACKING

Consider a circular track with a certain radius, say, r_0 , and imagine viewing a portion of it through the access window (see Fig. 1). It is through this window that the read-write head gains access to the disk and, by moving in the radial direction, reaches the various tracks. To a viewer looking through the window, a perfectly circular track centered on the spindle axis will look stationary, irrespective of the rotational speed of the disk. However, any track eccentricity will cause an apparent motion toward or away from the center. The peak-to-peak radial distance traveled by a track (as seen through the window) might depend on a number of factors, including centering accuracy of the hub, deformability of the disk substrate, mechanical vibrations, manufacturing tolerances, etc. For a 3.5-in plastic disk, for example, this peak-to-peak motion can be as much as 100 μm . Assuming a rotation rate of 3600 rpm, the apparent radial velocity of the track will be a few millimeter per second. Now, if the focused spot (which is only about 1 μm) remains stationary while trying to read or write on this track (whose width is also about 1 μm), it is clear that the beam will miss the track for a good fraction of every revolution cycle.

Practical solutions to the above problem are provided by automatic track-following techniques. Here the objective lens is placed in a fine actuator, typically a voice coil, which is capable of moving the necessary radial distances and maintaining a lock on the desired track. The signal that controls the movement of this actuator is derived from the reflected light itself, which carries information about the position of the focused spot relative to the track. There exist several mechanisms for extracting the track-error signal (TES) from the reflected light. All these methods require some sort of structure on the disk surface to identify the position of the track. In the case of read-only disks (CD, CD-ROM, and video disk) the embossed pattern of data provides ample information for tracking purposes. In the case of write-once and erasable disks, tracking guides are impressed on the substrate during the manufacturing process. The two major formats for these tracking guides are pregrooves (for continuous tracking) and sampled-servo marks (for discrete tracking). A combination of the two schemes, known as continuous/composite format, is often used in practice. This format is depicted schematically in Fig. 10 which shows a small section containing five tracks, each consisting of the tail end of a groove, synchronization marks, a mirror area for adjusting offsets, a pair of wobble marks for sampled tracking, and header information for sector identification.

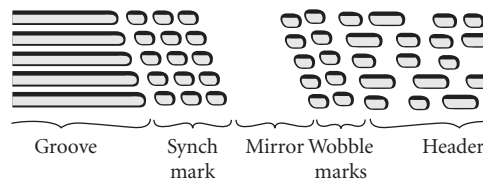


FIGURE 10 Servo offset fields in continuous/composite format contain a mirror area and offset marks for tracking. (Marchant, 1990.)

Tracking on Grooved Regions

As shown in Fig. 3a, grooves are continuous depressions that are embossed, etched, or molded onto the substrate prior to deposition of the storage medium. If the data is recorded on the grooves, then the lands are not used except for providing a guard band between neighboring grooves. Conversely, the land regions may be used to record the information, in which case grooves provide the guard band. Typical track widths are about one wavelength of the light. The guard bands are somewhat narrower, their exact shape and dimensions depending on the beam size, required track-servo accuracy, and the acceptable levels of crosstalk between adjacent tracks. The groove depth is typically around one-eighth of one wavelength ($\lambda/8$) which gives the largest TES in the push-pull method. The geometrical shape of the groove's cross section might be rectangular, trapezoidal, triangular, or some smooth version of these curves.

When the focused spot is centered on a given track, it is diffracted symmetrically from the two edges, resulting in a balanced far-field pattern. As soon as the spot moves away from the center, the symmetry breaks down and the far-field distribution tends to shift to one side or the other. A split photodetector placed in the path of the reflected light can therefore sense the relative position of the spot and provide the appropriate feedback signal (see Fig. 11). This is the essence of the push-pull

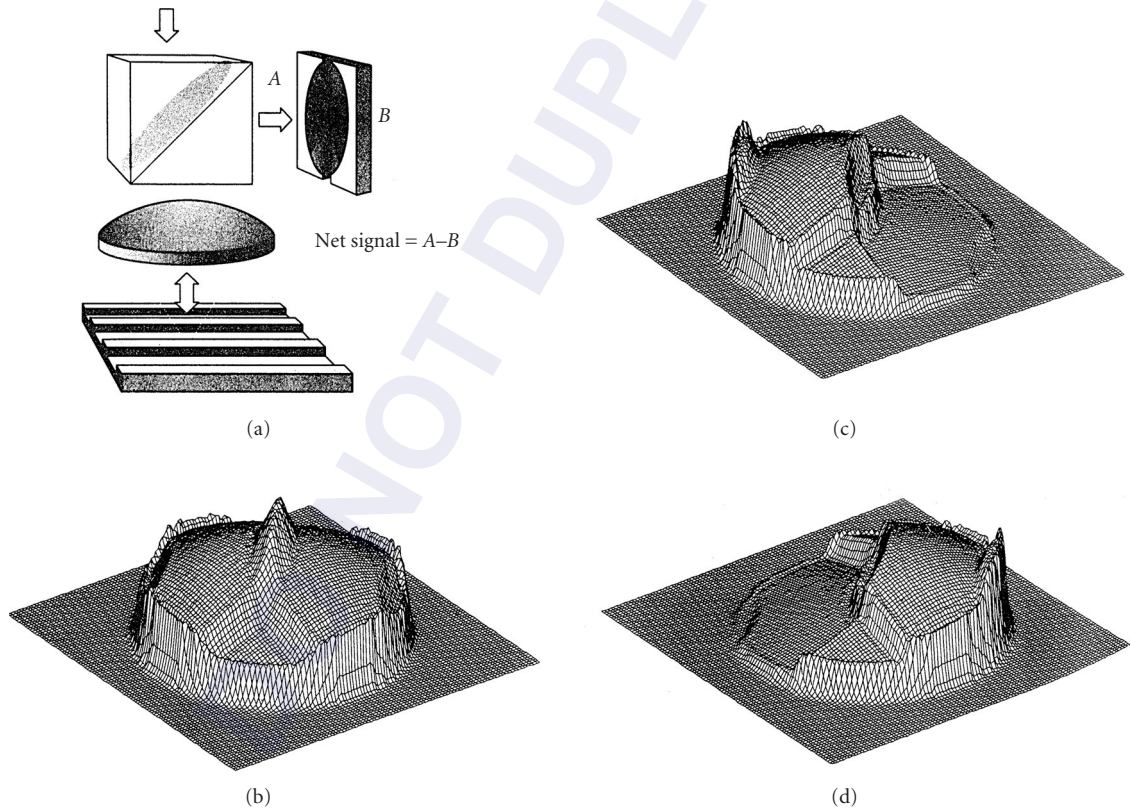


FIGURE 11 (a) Push-pull sensor for tracking on grooves. (Marchant, 1990.) (b) Light intensity distribution at the detector plane when the disk is in focus and the beam is centered on the track. (c) Light intensity distribution at the detector plane when the disk is in focus and the beam is centered on the groove edge. (d) Same as (c) except for the spot being on the opposite edge of the groove.

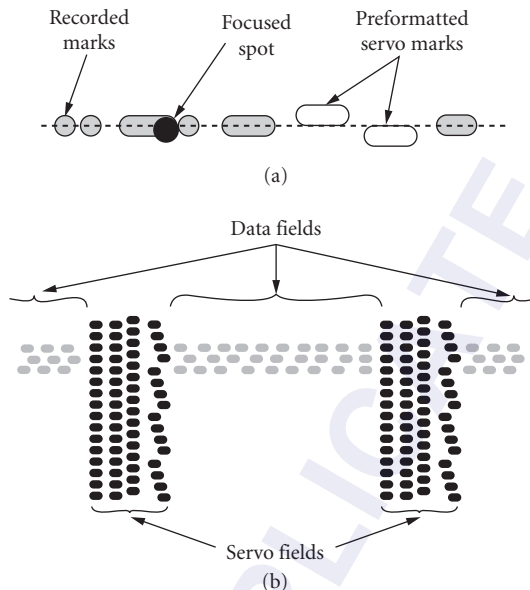


FIGURE 12 (a) In sampled tracking a pair of preformatted servo marks helps locate the position of the focused spot relative to the track center. (b) servo fields occur frequently and at regular intervals in sampled servo format. The data area shown here has data recorded on three tracks. (Marchant, 1990.)

method. Figure 11 also shows intensity plots at the detector plane after reflection from various locations on the grooved surface. Note how the intensity shifts to one side or the other depending on whether the spot moves to the right edge or to the left edge of the groove.

Sampled Tracking

Since dynamic track runout is usually a slow and gradual process, there is actually no need for continuous tracking as done on grooved media. A pair of embedded marks, offset from the track center as in Fig. 12a, can provide the necessary information for correcting the relative position of the focused spot. The reflected intensity will indicate the positions of the two servo marks as two successive short pulses. If the beam happens to be on track, the two pulses will have equal magnitudes and there shall be no need for correction. If, on the other hand, the beam is off-track, one of the pulses will be stronger than the other. Depending on which pulse is the stronger, the system will recognize the direction in which it has to move and will correct the error accordingly. Sampled-servo mark pairs must be provided frequently enough to ensure proper track-following. In a typical application, the track might be divided into groups of 18 bytes, 2 bytes dedicated as servo offset areas and 16 bytes filled with other format information or left blank for user data. Figure 12b shows a small section from a sampled-servo disk containing a number of tracks, three of which are recorded with user data. The track-servo marks in this case are preceded by synch marks (also prerecorded on the servo offset area). Note in Fig. 12b that the format marks repeat a certain pattern every four tracks. This pattern is known as a “gray code,” and allows the system to recognize and correct minor track-counting errors during the seek operation.

Track Counting During the Seek Operation

In the seek operation the coarse actuator moves the head assembly across the disk to a new location where the desired track is located. In order to avoid landing on a nearby track and being forced to perform a second (fine) seek, most systems in use today count the tracks as they are being crossed. In this way the head can land on the correct track and thereby minimize the overall seek time. The sampled-servo format is not suitable for this purpose, since the servo marks do not occur frequently enough to allow uninterrupted counting. In contrast, grooved media provide the necessary information for track-counting.

During a seek operation the focus servo loop remains closed, maintaining focus as the head crosses the tracks. The tracking loop, on the other hand, must be opened. The zero crossings of the TES then provide the track count. Complications may arise in this process, however, due to eccentricities of tracks. As was mentioned earlier, to an observer looking through the access window, an eccentric track moves in and out radially with a small (but not insignificant) velocity. As the head approaches the desired track and slows down to capture it, its velocity might fall just short of the apparent track velocity. Under these circumstances, a track which has already been counted may catch up with the head and be counted once again. Intelligence must be built into the system to recognize and avoid such problems. Also, through the use of gray codes and similar schemes, the system can be made to correct its occasional miscounts before finally locking onto the destination track.

35.6 THERMOMAGNETIC RECORDING PROCESS

Recording and erasure of information on a magneto-optical disk are both achieved by the thermomagnetic process. The essence of thermomagnetic recording is shown in Fig. 13. At the ambient temperature the film has a high magnetic coercivity* and therefore does not respond to the externally applied field. When a focused laser beam raises the local temperature of the film, the hot spot becomes magnetically soft (i.e., its coercivity drops). As the temperature rises, coercivity drops continuously until such time as the field of the electromagnet finally overcomes the material's resistance to reversal and switches its magnetization. Turning the laser off brings the temperatures back to normal, but the reverse-magnetized domain remains frozen in the film. In a typical situation in practice, the film thickness may be around 300 Å, laser power at the disk ≈ 10 mW, diameter of the focused spot ≈ 1 μm , laser pulse duration ~ 50 ns, linear velocity of the track ≈ 10 m/s, and the magnetic field strength ≈ 200 gauss. The temperature may reach a peak of 500 K at the center of the spot, which is certainly sufficient for magnetization reversal, but is not nearly high enough to melt or crystallize or in any other way modify the structure of the material.

The materials of MO recording have strong perpendicular magnetic anisotropy. This type of anisotropy favors the "up" and "down" directions of magnetization over all other orientations. The disk is initialized in one of these two directions, say, up, and the recording takes place when small regions are selectively reverse-magnetized by the thermomagnetic process. The resulting magnetization distribution then represents the pattern of recorded information. For instance, binary sequences may be represented by a mapping of zeros to up-magnetized regions and ones to down-magnetized regions [non-return to zero (NRZ) scheme]. Alternatively, the non-return to zero inverted (NRZI) scheme might be used, whereby transitions (up-to-down and down-to-up) are used to represent the ones in the bit sequence.

Recording by Laser Power Modulation

In this traditional approach to thermomagnetic recording, the electromagnet produces a constant field, while the information signal is used to modulate the power of the laser beam. As the disk rotates under the focused spot, the pulsed laser beam creates a sequence of up/down domains along

*Coercivity of a magnetic medium is a measure of its resistance to magnetization reversal. For example, consider a thin film with perpendicular magnetic moment saturated in the +Z direction, as in Fig. 13a. A magnetic field applied along -Z will succeed in reversing the direction of magnetization only if the field is stronger than the coercivity of the film.

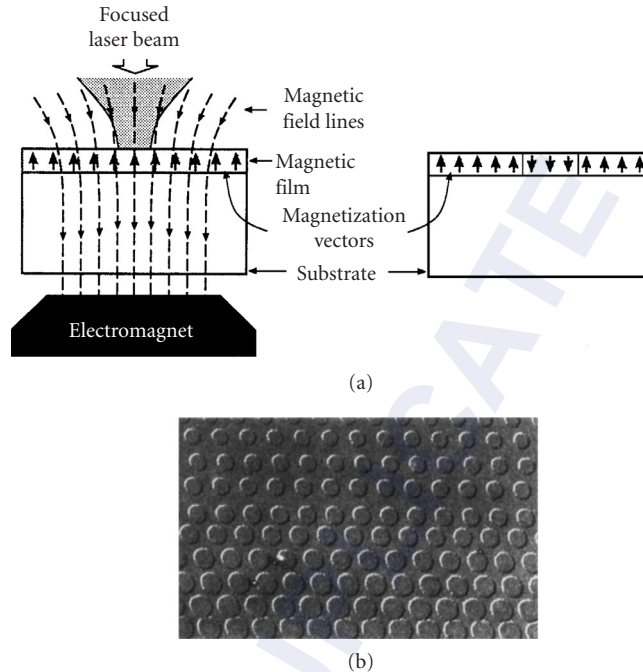


FIGURE 13 (a) Thermomagnetic recording process. The field of the electromagnet helps reverse the direction of magnetization in the area heated by the focused laser beam. (b) Lorentz micrograph of domains written thermomagnetically. The various tracks shown here were written at different laser powers, with power level decreasing from top to bottom.

the track. The Lorentz electron micrograph in Fig. 13*b* shows a number of domains recorded by laser power modulation (LPM). The domains are highly stable and may be read over and over again without significant degradation. If, however, the user decides to discard a recorded block and to use the space for new data, the LPM scheme does not allow direct overwrite; the system must erase the old data during one revolution and record the new data in a subsequent revolution cycle.

During erasure, the direction of the external field is reversed, so that up-magnetized domains in Fig. 13*a* now become the favored ones. Whereas writing is achieved with a modulated laser beam, in erasure the laser stays on for a relatively long period of time, erasing an entire sector. Selective erasure of individual domains is not practical, nor is it desired, since mass data storage systems generally deal with data at the level of blocks, which are recorded onto and read from individual sectors. Note that at least one revolution cycle elapses between the erasure of an old block and its replacement by a new block. The electromagnet therefore need not be capable of rapid switchings. (When the disk rotates at 3600 rpm, for example, there is a period of 16 ms or so between successive switchings.) This kind of slow reversal allows the magnet to be large enough to cover all the tracks simultaneously, thereby eliminating the need for a moving magnet and an actuator. It also affords a relatively large gap between the disk and the magnet tip, which enables the use of double-sided disks and relaxes the mechanical tolerances of the system without overburdening the magnet's power supply.

The obvious disadvantage of LPM is its lack of direct overwrite capability. A more subtle concern is that it is perhaps unsuitable for the pulse width modulation (PWM) scheme of representing binary waveforms. Due to fluctuations in the laser power, spatial variations of material properties,

lack of perfect focusing and track-following, etc., the length of a recorded domain along the track may fluctuate in small but unpredictable ways. If the information is to be encoded in the distance between adjacent domain walls (i.e., PWM), then the LPM scheme of thermomagnetic writing may suffer from excessive domain-wall jitter. Laser power modulation works well, however, when the information is encoded in the position of domain centers [i.e., pulse position modulation (PPM)]. In general, PWM is superior to PPM in terms of the recording density, and methods that allow PWM are therefore preferred.

Recording by Magnetic Field Modulation

Another method of thermomagnetic recording is based on magnetic field modulation (MFM), and is depicted schematically in Fig. 14a. Here the laser power may be kept constant while the information signal is used to modulate the direction of the magnetic field. Photomicrographs of typical domain patterns recorded in the MFM scheme are shown in Fig. 14b. Crescent-shaped domains are the hallmark of the field modulation technique. If one assumes (using a much simplified model) that the magnetization aligns itself with the applied field within a region whose temperature has passed a certain critical value, T_{crit} , then one can explain the crescent shape of these domains in the following way: with the laser operating in the CW mode and the disk moving at constant velocity, temperature distribution in the magnetic medium assumes a steady-state profile, such as that in Fig. 14c. Of course, relative to the laser beam, the temperature profile is stationary, but in the frame

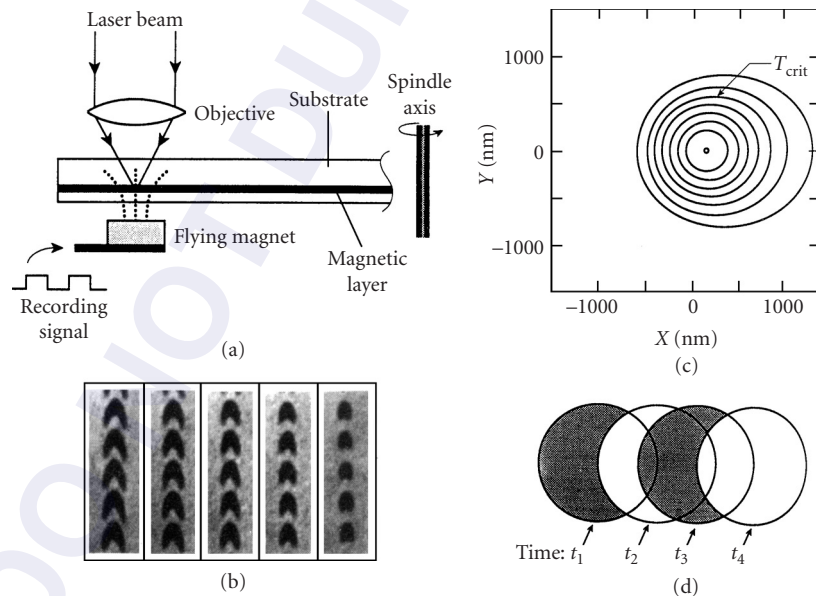


FIGURE 14 (a) Thermomagnetic recording by magnetic field modulation. The power of the beam is kept constant, while the magnetic field direction is switched by the data signal. (b) Polarized-light microphotograph of recorded domains. (c) Computed isotherms produced by a CW laser beam, focused on the magnetic layer of a disk. The disk moves with constant velocity under the beam. The region inside the isotherm marked as T_{crit} is above the critical temperature for writing, thus its magnetization aligns itself with the direction of the applied magnetic field. (d) Magnetization within the heated region (above T_{crit}) follows the direction of the applied magnetic field, whose switchings occur at times t_n . The resulting domains are crescent-shaped.

of reference of the disk the profile moves along the track with the linear track velocity. The isotherm corresponding to T_{crit} is identified as such in the figure; within this isotherm the magnetization always aligns itself with the applied field. A succession of critical isotherms along the track, each obtained at the particular instant of time when the magnetic field switches direction, is shown in Fig. 14*d*. From this picture it is not difficult to see how the crescent-shaped domains form, and also to understand the relation between the waveform that controls the magnet and the resulting domain pattern.

The advantages of magnetic field modulation recording are that (1) direct overwriting is possible, and (2) domain wall positions along the track, being rather insensitive to defocus and laser power fluctuations, are fairly accurately controlled by the timing of the magnetic field switchings. On the negative side, the magnet must now be small and fly close to the disk surface if it is to produce rapidly switched fields with a magnitude of a few hundred gauss. Systems that utilize magnetic field modulation often fly a small electromagnet on the opposite side of the disk from the optical stylus. Since mechanical tolerances are tight, this might compromise the removability of the disk in such systems. Moreover, the requirement of close proximity between the magnet and the storage medium dictates the use of single-sided disks in practice.

Thermal Optimization of the Media—Multilayer Structures

The thermal behavior of an optical disk can be modified and improved if the active layer is incorporated into a properly designed multilayer structure, such as that shown in Fig. 15. In addition to thermal engineering, multilayers allow protective mechanisms to be built around the active layer; they also enable the enhancement of the signal-to-noise ratio in readout. (This latter feature is further explored in Sec. 35.7.) Multilayers are generally designed to optimize the absorption of light by creating an antireflection structure, whereby a good fraction of the incident optical power is absorbed in the active layer. Whereas the reflectivity of bare metal films is typically over 50 percent, a quadrilayer structure can easily reduce that to 20 percent or even less, if so desired. Multilayers can also be designed to control the flow of heat generated by the absorbed light. The aluminum reflecting layer in the quadrilayer of Fig. 15, for instance, may be used as a heat sink for the magnetic layer, thus minimizing the undesirable effects of lateral heat diffusion within the magnetic medium.

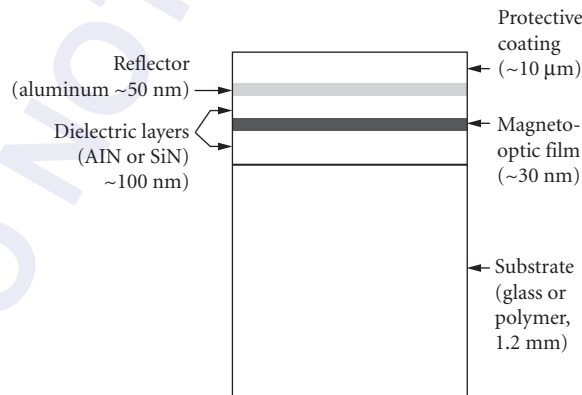


FIGURE 15 Quadrilayer magneto-optical disk structure. This particular design is for use in the substrate incident mode, where the light goes through the substrate before reaching the MO layer. The thicknesses of the various layers can be optimized for enhancing the read signal, increasing the absorbed laser power, and controlling the thermal profile. Note in particular that the aluminum layer can play the dual roles of light reflector and heat sink.

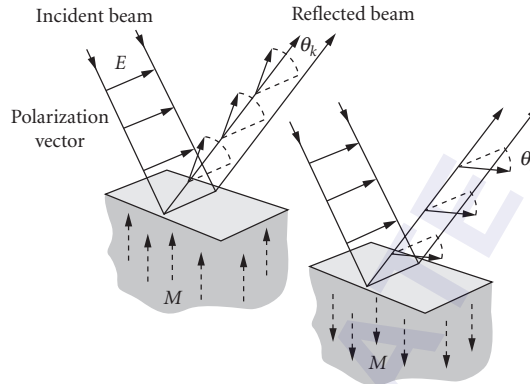


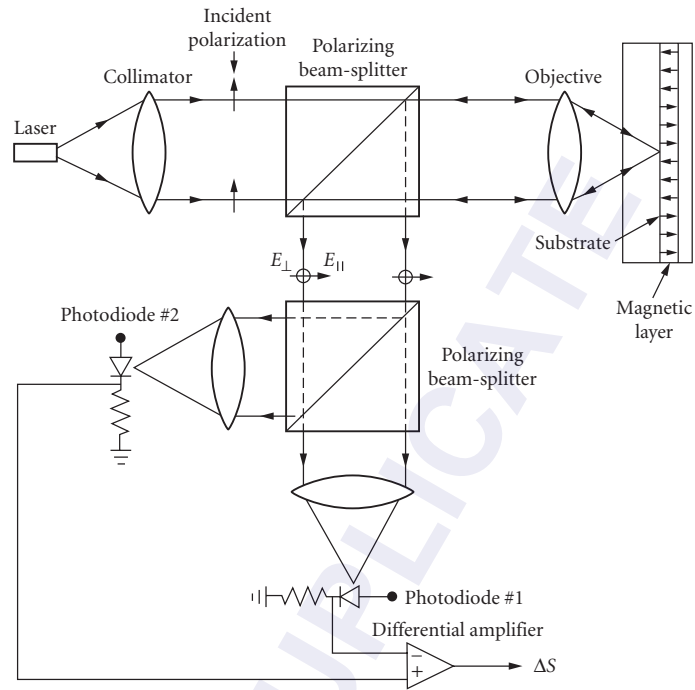
FIGURE 16 Schematic diagram describing the polar magneto-optical Kerr effect. Upon reflection from the surface of a perpendicularly magnetized medium, the polarization vector undergoes a rotation. The sense of rotation depends on the direction of magnetization \mathbf{M} , and switches sign when \mathbf{M} is reversed.

35.7 MAGNETO-OPTICAL READOUT

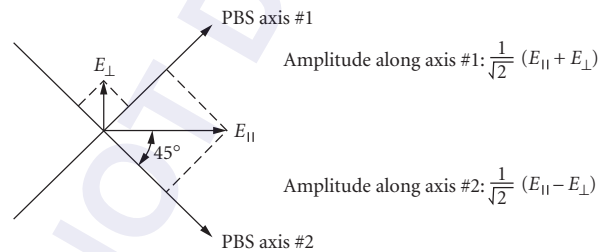
The information recorded on a perpendicularly magnetized medium may be read with the aid of the polar magneto-optical Kerr effect. When linearly polarized light is normally incident on a perpendicular magnetic medium, its plane of polarization undergoes a slight rotation upon reflection. This rotation of the plane of polarization, whose sense depends on the direction of magnetization in the medium, is known as the polar Kerr effect. The schematic representation of this phenomenon in Fig. 16 shows that if the polarization vector suffers a counterclockwise rotation upon reflection from an up-magnetized region, then the same vector will rotate clockwise when the magnetization is down. A magneto-optical medium is characterized in terms of its reflectivity R and its Kerr rotation angle θ_k .^{*} In MO readout, it is the sign of the rotation angle that carries the information about the state of magnetization of the medium, i.e., the recorded bit pattern.

The laser used for readout is usually the same as that used for recording, but its output power level is substantially reduced in order to avoid erasing (or otherwise obliterating) the previously recorded information. For instance, if the power of the write/erase beam is 20 mW, then for the read operation the beam is attenuated to about 3 or 4 mW. The same objective lens that focuses the write beam is now used to focus the read beam, creating a diffraction-limited spot for resolving the recorded marks. Whereas in writing the laser was pulsed to selectively reverse-magnetize small regions along the track, in readout it operates with constant power, i.e., in CW mode. Both up- and down-magnetized regions are read as the track passes under the focused light spot. The reflected beam, which is now polarization-modulated, goes back through the objective and becomes collimated once again; its information content is subsequently decoded by polarization-sensitive optics, and the scanned pattern of magnetization is reproduced as an electronic signal.

^{*}In reality, the reflected state of polarization is not linear, but has a certain degree of ellipticity. One may consider the reflected polarization as consisting of two linear components: E_{\parallel} which is parallel to the direction of incident polarization, and E_{\perp} which is perpendicular to it. Now, if E_{\parallel} is in phase with E_{\perp} , the net magneto-optic effect will be a pure rotation of the polarization vector. On the other hand, if E_{\parallel} and E_{\perp} are 90° out of phase, then the reflected polarization will be elliptical, with no rotation whatsoever. In practice, the phase difference between E_{\parallel} and E_{\perp} is somewhere between 0° and 90°, resulting in a reflected beam which has some degree of ellipticity ϵ_k , with the major axis of the polarization ellipse rotated by an angle θ_k (relative to the incident \mathbf{E} vector). By inserting a Soleil-Babinet compensator in the reflected beam's path, one can change the phase relationship between E_{\parallel} and E_{\perp} in such a way as to eliminate the beam's ellipticity; the emerging polarization then will become linear with an enhanced rotation angle. In this chapter, reference to Kerr angle implies the effective angle which includes the above correction for ellipticity.



(a)



(b)

FIGURE 17 Differential detection scheme utilizes a polarizing beam splitter and two photodetectors in order to convert the rotation of polarization to an electronic signal. E_{\parallel} and E_{\perp} are the reflected components of polarization; they are, respectively, parallel and perpendicular to the direction of incident polarization. The diagram in (b) shows the orientation of the PBS axes relative to the polarization vectors.

Differential Detection

Figure 17 shows the differential detection system that is the basis of magneto-optical readout in practically all erasable optical storage systems in use today. The beam splitter (BS) diverts half of the reflected beam away from the laser and into the detection module. The polarizing beam splitter (PBS) splits the beam into two parts, each carrying the projection of the incident polarization along

one axis of the PBS, as shown in Fig. 17b. The component of polarization along one of the axes goes straight through, while the component along the other axis splits off to the side. If, upon reflection from the disk, the polarization did not undergo any rotations whatsoever, then the beam entering the PBS would be polarized at 45° to the PBS axes, in which case it would split equally between the two branches. Under this condition, the two detectors generate identical signals and the differential signal ΔS will be zero. Now, if the beam returns from the disk with its polarization rotated clockwise (rotation angle = θ_k), then detector 1 will receive more light than detector 2, and the differential signal will be positive. Similarly, a counterclockwise rotated beam entering the PBS will generate a negative ΔS . The electronic signal ΔS thus reproduces the pattern of magnetization along the scanned track.

Enhancement of the Signal-to-Noise Ratio by Multilayering

The materials suitable for optical recording presently have very small Kerr angles (typically $\theta_k \approx 0.5^\circ$), with the result that the signal ΔS is correspondingly small. Multilayering schemes designed for the enhancement of the MO signal increase the interaction between the light and the magnetic medium by encapsulating a thin film of the MO material in an antireflection-type structure. By providing a better index match between the MO film and its surroundings, and also by circulating the light through the MO film, multilayered structures manage to trap a large fraction of the incident light within the magnetized medium, and thus increase the Kerr rotation angle. These efforts inevitably result in a reduced reflectivity, but since the important parameter is the magneto-optically generated component of polarization, $E_\perp = \sqrt{R} \sin \theta_k$, it turns out that a net gain in the signal-to-noise ratio can be achieved by adopting the multilayering schemes. Reported enhancements of E_\perp have been as large as a factor of 5. The popular quadrilayer structure depicted in Fig. 15 consists of a thin film of the MO material, sandwiched between two transparent dielectric layers, and capped off with a reflecting metallic layer. The entire structure, which is grown on a transparent substrate (through which light must travel to reach the MO film), is protected by a lacquer layer on the top. Numbers shown in Fig. 15 for the various layer thicknesses are representative of currently designed quadrilayers.

The advantage of sending the light through the substrate is that the front facet of the disk stays out of focus during operation. In this way, small dust particles, finger prints, and scratches will not block the passage of light, and their deteriorating effects on the quality of the focused spot (which affects the integrity of both writing and readout) will be minimized. Any optical storage medium designed for removability ought to have the kind of protection that illumination through the substrate provides. The note of caution with substrate-side illumination is that, if the objective is simply designed for focusing in the air, then the oblique rays will bend upon entering the substrate and deviate from nominal focus, causing severe aberrations (see Fig. 7c). Therefore, the substrate thickness and refractive index must be taken into account in the objective's design.

Sources of Noise in Readout

The read signal is always accompanied by random noise. The effective noise amplitude (relative to the strength of the signal) ultimately limits the performance of any readout system. Part of the noise is thermal in nature, arising from random fluctuations of charge carriers within the photodiodes, resistors, and amplifiers. In principle, this source of noise can be mitigated by reducing the operating temperature of the device. However, since operating below the normal room temperature is not very practical for data storage systems, one must accept some of the limitations brought about by the thermal noise.

Another source of readout noise is shot noise which, in classical language, is due to random arrival of photons at the photodetector(s). This noise is a permanent companion of the read signal and cannot be eliminated, but the system parameters may be adjusted to minimize its effect. One

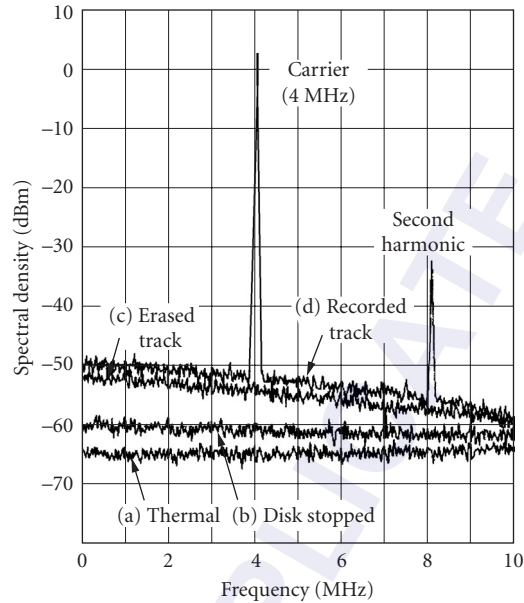


FIGURE 18 Spectra of the various noise components in magneto-optical readout.

property of the shot noise is that its rms amplitude is proportional to the square root of the available optical power P_o . Since the signal strength is directly proportional to P_o , it is clear that by increasing the read power of the laser one can enhance the ratio of signal-to-shot noise. There is, however, an upper limit on the laser read power, since the rise in the temperature of the medium will force the decline of its magneto-optical response.

Other sources of noise in magneto-optical readout include the laser noise, the media noise, and the data noise. Laser noise is caused by amplitude/phase fluctuations of the electromagnetic radiation that comprises the optical beam. Media noise arises from variations in the reflectivity/magneto-optic activity of the medium across its surface. The presence of grooves with rough and nonuniform edges can be a source of media noise as well. The term *data noise* refers to the unpredictable variations of the read signal arising from the imperfect shape/position of the recorded marks.

Figure 18 shows the various components of noise in a typical MO readout system, as detected by a spectrum analyzer. In (a) the light beam is blocked and the trace on the analyzer screen is solely due to the thermal noise. The trace in (b) where the beam reaches the detectors but the disk is stationary shows the combined effect of thermal, shot, and laser noise. Trace (c) corresponds to reading an erased track on a spinning disk; the noise here includes all of the above plus the media noise. When a single-frequency tone was recorded on the track and the read-back signal was fed to the spectrum analyzer, trace (d) was obtained. The narrow pulse at frequency f_0 is the first harmonic of the recorded signal; the corresponding second harmonic appears at $2f_0$. The noise level in this case is somewhat greater than that from the same track before the data was recorded. This difference is due to “data noise” and arises from jitter and nonuniformity of the recorded marks.

A commonly used measure of performance for optical recording media is the carrier-to-noise ratio (CNR). This is the ratio of the signal amplitude at the carrier frequency f_0 to the average level of noise. On a logarithmic scale the ratio is simply the difference between the two levels; in Fig. 18 the CNR is 53 decibels (dB).

35.8 MATERIALS OF MAGNETO-OPTICAL RECORDING

Amorphous rare earth transition metal alloys are presently the media of choice for erasable optical data storage applications. The general formula for the composition of the alloy may be written $(\text{Tb}_y\text{Gd}_{1-y})_x(\text{Fe}_x\text{Co}_{1-x})_{1-x}$ where terbium and gadolinium are the rare earth (RE) elements, while iron and cobalt are the transition metals (TM). In practice, the transition metals constitute roughly 80 atomic percent of the alloy (i.e., $x \approx 0.2$). In the transition metal subnetwork the fraction of cobalt is usually small, typically around 10 percent, and iron is the dominant element ($z \approx 0.9$). Similarly, in the rare earth subnetwork Tb is the main element ($y \approx 0.9$) while the Gd content is small or it may even be absent in some cases. Since the rare earth elements are highly reactive, RE-TM films tend to have poor corrosion resistance and, therefore, require protective coatings. In a disk structure such as that shown in Fig. 15, the dielectric layers that enable optimization of the medium for the best optical/thermal behavior also perform the crucial task of protecting the MO layer from the environment.

The amorphous nature of the material allows its composition to be continuously varied until a number of desirable properties are achieved. (In other words, the fractions x , y , z of the various elements are not constrained by the rules of stoichiometry.) Disks with large surface areas are coated uniformly with thin films of these media, and, in contrast to polycrystalline films whose grains and grain boundaries scatter the light beam and cause noise, these amorphous films are smooth and substantially noise-free. The films are deposited either by sputtering from an alloy target, or by cosputtering from multiple elemental targets. In the latter case, the substrate moves under the various targets and the fraction of a given element in the alloy film is determined by the time spent under the target as well as the power applied to that target. Substrates are usually kept at a low temperature (by water cooling, for instance) in order to reduce the mobility of deposited atoms and to inhibit crystal growth. Factors that affect the composition and short-range order of the deposited films include the type of the sputtering gas (argon, krypton, xenon, etc.) and its pressure during sputtering, the bias voltage applied to the substrate, deposition rate, nature of the substrate and its pretreatment, temperature of the substrate, etc.

Ferrimagnetism

The RE-TM alloys of interest in MO recording are ferrimagnetic, in the sense that the magnetization of the TM subnetwork is antiparallel to that of the RE subnetwork. The net magnetic moment exhibited by the material is the vector sum of the two subnetwork magnetizations. Figure 19 shows a typical temperature dependence of RE and TM magnetic moments, as well as the net saturation moment of the material. The exchange coupling between the two magnetic subnetworks is strong enough to give them the same critical temperature T_c . At $T = 0$ K the rare earth moment is stronger than that of the transition metal, giving the material a net moment along the direction of the RE magnetization. As the temperature rises, thermal disorder competes with interatomic exchange forces that tend to align the individual atomic dipole moments. The decrease of M_{RE} with the increasing temperature is faster than that of M_{TM} , and the net moment M_s begins to approach zero. At the compensation point temperature T_{comp} , the net moment vanishes. Between T_{comp} and T_c the net moment is dominated by the TM subnetwork and the material is said to exhibit TM-rich behavior (as opposed to when $T < T_{\text{comp}}$, where it exhibits RE-rich behavior). At the Curie temperature, thermal agitations finally break the hold of the exchange forces on magnetic dipoles, and the magnetic order disappears. Beyond T_c the material is in the paramagnetic state.

The composition of the materials of magneto-optical storage is chosen so that T_{comp} appears near the ambient temperature of $T_a \approx 300$ K. Thus, under normal conditions, the net magnetization of the material is close to zero. Figure 20 shows a schematic drawing of the magnetization pattern in the cross section of a recorded track. Note that, although the net magnetization is nearly zero everywhere, the subnetwork moments have opposite orientations in adjacent domains. During readout the light from the GaAs laser interacts mainly with the transition metal subnetwork; thus, the MO Kerr signal is strong even though the net magnetization of the storage layer may be small. The magnetic

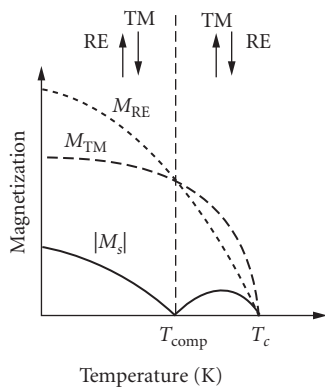


FIGURE 19 Temperature dependence of magnetization in amorphous RE-TM films. The moments of RE and TM subnetworks decrease monotonically, until they both vanish at the critical (Curie) temperature T_c . The net magnetization is the difference between the two subnetwork moments, and goes through zero at the compensation point T_{comp} .

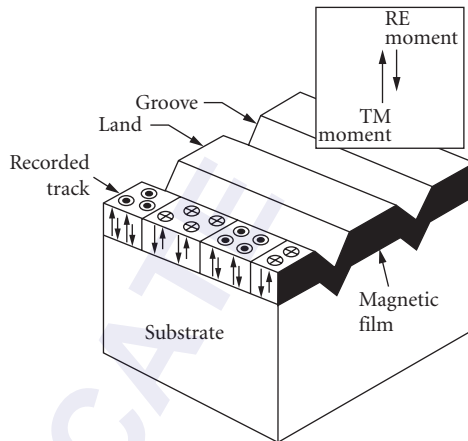


FIGURE 20 Schematic diagram showing the pattern of magnetization along a recorded track. The rare earth and the transition metal moments couple antiferromagnetically, so that the net magnetization everywhere is small. However, since the read beam interacts mainly with the TM subnetwork, the read-out signal is not necessarily small.

electrons of iron and cobalt are in the $3d$ electronic shell, which forms the outer layer of the ion once the $4s$ electrons have escaped into the sea of conduction electrons. The magnetic electrons of Tb and Gd, in contrast, are within the $4f$ shell, concealed by the $5s$, $5p$, and $5d$ shells, even after the $6s$ electrons escape to the conduction band. A red or near-infrared photon is not energetic enough to penetrate the outer shell and interact with the magnetic $4f$ electrons, but it readily interacts with the exposed $3d$ electrons that constitute the magnetic moment of the TM subnetwork. It is for this reason that the MO Kerr signal in the visible and in the infrared is a probe of the state of magnetization of the TM subnetwork.

Perpendicular Magnetic Anisotropy

An important property of amorphous RE-TM alloy films is that they possess perpendicular magnetic anisotropy. The magnetization in these films favors perpendicular orientation even though there is no discernible crystallinity or microstructure that might obviously be responsible for this behavior. It is generally believed that atomic short-range order, established in the deposition process and aided by the symmetry-breaking at the surface of the film, gives preference to perpendicular orientation. Unequivocal proof of this assertion, however, is not presently available due to a lack of high-resolution observation instruments.

The perpendicular magnetization of MO media is in sharp contrast to the in-plane orientation of the magnetization vector in ordinary magnetic recording. In magnetic recording, the neighboring domains are magnetized in head-to-head fashion, which is an energetically unfavorable situation, since the domain walls are charged and highly unstable. The boundary between neighboring domains in fact breaks down into zigzags, vortices, and all manner of jagged, uneven structure in an attempt to reduce the magnetostatic energy. In contrast, adjacent domains in MO media are highly stable, since the pattern of magnetization causes flux closure, which reduces the magnetostatic energy.

Coercivity and the Hysteresis Loop

Typical hysteresis loops of an amorphous RE-TM thin film at various temperatures are shown in Fig. 21a. These loops, obtained with a vibrating sample magnetometer (VSM), show several characteristics of the MO media. (The VSM applies a slowly varying magnetic field to the sample and measures its net magnetic moment as a function of the field.) The horizontal axis in Fig. 21a is the applied field, which varies from -12 to $+12$ kOe, while the vertical axis is the measured magnetic moment per unit volume (in CGS units of emu/cm^3). The high degree of squareness of the loops signifies the following:

1. The remanent magnetization M_r is the same as the saturation magnetization M_s . Thus, once the sample is saturated with the help of an applied field, removing that field does not cause a reduction of the magnetic moment.
2. Transitions of the magnetization from up to down (or from down to up) are very sharp. The reverse field does not affect the magnetization until the critical value of H_c , the coercive field, is reached. At the coercive field the magnetization suddenly reverses direction, and saturation in the opposite direction is almost immediate.

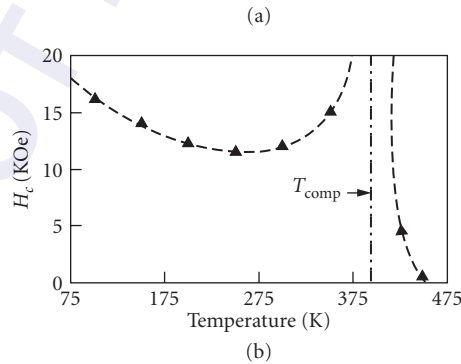
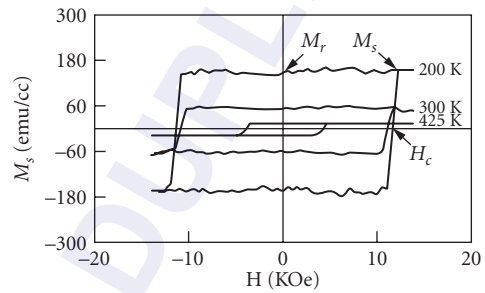


FIGURE 21 (a) Hysteresis loops of an amorphous $\text{Tb}_{27}(\text{FeCo})_{73}$ film, measured by VSM at three different temperatures. The saturation moment M_s , the remanent moment M_r , and the coercive field H_c are identified for the loop measured at $T = 200$ K. (b) Coercivity as function of temperature for the above sample. At the compensation temperature, $T_{\text{comp}} = 400$ K, the coercivity is infinite; it drops to zero at the Curie point $T_c = 450$ K.

The fact that M_r is very nearly equal to M_s in MO media is significant, since it means that the recorded domains remain fully saturated and exhibit maximum signal during readout. The coercivity H_c , in addition to being responsible for the stability of recorded domains, plays an important role in the processes of thermomagnetic recording and erasure. The coercivity at room temperature, being of the order of several thousand oersteds, prevents fields weaker than H_c from destroying (or disturbing) any recorded data. With the increasing temperature, the coercivity decreases and drops to zero at the Curie point, T_c . Figure 21b is the plot of H_c versus T for the same sample as in (a). Note that at the compensation point the coercivity goes to infinity, simply because the magnetization vanishes, and the external field does not see any magnetic moments to interact with. Above T_{comp} the coercive field decreases monotonically, which explains the process of magnetization reversal during thermomagnetic recording: \mathbf{M} switches sign once the coercivity drops below the level of the applied field.

35.9 CONCLUDING REMARKS

In this chapter we have reviewed the basic characteristics of optical disk data storage systems, with emphasis on magneto-optical recording. The goal has been to convey the important concepts without getting distracted by secondary issues and less significant details. As a result, we have glossed over several interesting developments that have played a role in the technological evolution of optical data storage. In this final section some of these developments are briefly described.

Multiple-Track Read-Write with Diode Laser Arrays

It is possible in an optical disk system to use an array of lasers instead of just one, focus all the lasers simultaneously through the same objective lens, and perform parallel read/write/erase operations on multiple tracks. Since the individual lasers of an array can be modulated independently, the parallel channels thus obtained are totally independent of each other. In going from a single-channel drive to a multiple-channel one, the optics of the system (i.e., lenses, beam splitters, polarization-sensitive elements, focus and track servos, etc.) remain essentially the same; only the lasers and detectors proliferate in number. Parallel track operations boost the sustainable data rates in proportion to the number of channels used.

Diffractive Optics

The use of holographic optical elements (HOEs) to replace individual refractive optics is a promising new development. Diffractive optical elements are relatively easy to manufacture, they are lightweight and inexpensive, and can combine the functions of several elements on a single plate. These devices are therefore expected to help reduce the cost, size, and weight of optical heads, making optical drives more competitive in terms of price and performance.

An example of the application of HOEs in MO systems is shown in Fig. 22, which shows a reflection-type element consisting of four separate holograms. The light incident on the HOE at an angle of 60° has a p component which is the original polarization of the laser beam, and an s component (parallel to the hologram's grooves) which is the magneto-optically generated polarization at the disk. Nearly 90 percent of the s and 70 percent of the p polarization in this case are reflected from the HOE without suffering any diffraction (i.e., in the zero-order beam); they are captured in the differential detection module and yield the MO read signal. The four holograms deflect 20 percent of the incident p -polarized light in the form of first-order diffracted beams and bring them to focus at four different spots on a multielement detector. The two small holograms in the middle, H_3 and H_4 , focus their first-order beams on detectors P_3 and P_4 , to generate the push-pull

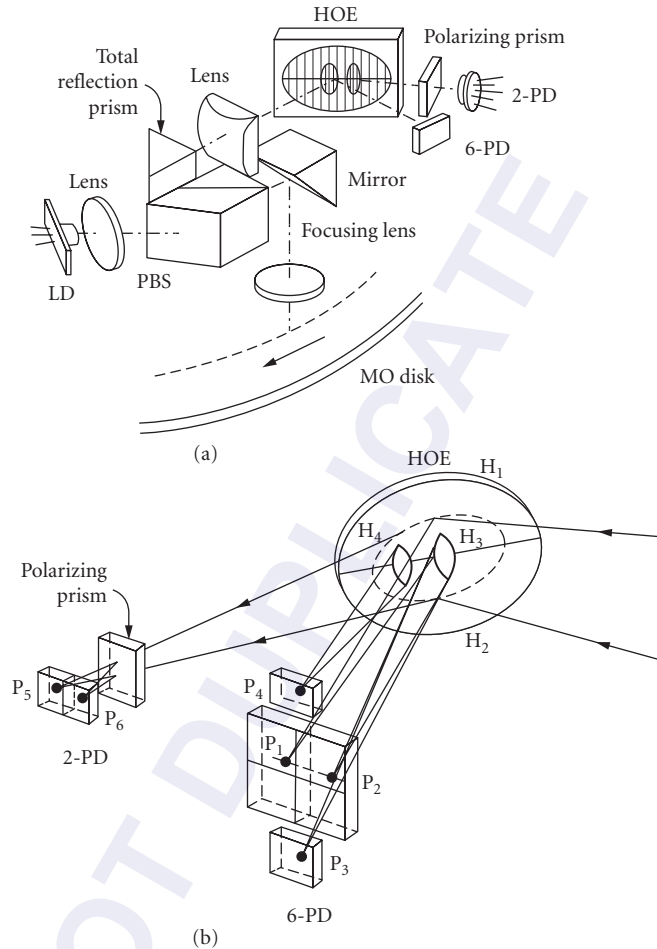


FIGURE 22 Application of holographic optical element (HOE) in optical recording. (a) Configuration of MO head using a polarization-sensitive reflection HOE. (b) Geometrical relation between holograms and detectors. (After A. Ohba et al., *SPIE Proceedings*, Vol. 1078, 1989.)

tracking-error signal. The other two holograms, H_1 and H_2 , send the diffracted beams to a four-element detector in order to generate a focus-error signal based on the double knife-edge scheme. This HOE, therefore, combines the functions of beam splitting, masking, and focusing all in one compact unit.

Alternative Storage Media

The GaAlAs lasers of the present optical disk technology will likely be replaced in the future by light sources that emit in the blue end of the spectrum. Shorter wavelengths allow smaller marks to be recorded, and also enable the resolving of those marks in readout. Aside from the imposition of tighter tolerances on focusing and tracking servos, operation in the blue will require storage materials

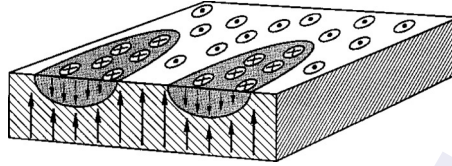


FIGURE 23 Direct overwrite in exchange-coupled magnetic multilayers involves the formation of domains that do not extend through the entire thickness of the magnetic medium.

that are sensitive to the short wavelengths. The current favorites for erasable optical recording media, the amorphous RE-TM alloys, may not be suitable for readout in the blue, since their magneto-optic Kerr signal drops at short wavelengths. A new class of magnetic materials which holds promise for future-generation device applications is the class of TM/TM superlattice-type media. The best-known material in this class is the Co/Pt-layered structure which consists of very thin layers of cobalt (typically one or two atomic layers), separated by several atomic layers of platinum. These polycrystalline films which have very small crystallites (grain diameter $\approx 200 \text{ \AA}$) are prepared either by electron beam evaporation or by sputtering. Co/Pt films have large perpendicular magnetic anisotropy, good signal-to-noise ratios in the blue, sufficient sensitivity for write/erase operations, and are environmentally more stable than the RE-TM media.

Direct Overwrite in Magneto-Optical Systems

The problem of direct overwrite (DOW) on MO media has been the subject of extensive research in recent years. Some of the most promising solutions have been based on exchange-coupled magnetic multilayered structures. The basic idea of recording on exchange-coupled bilayers (or trilayers) is simple and involves the writing of reverse domains that do not extend through the entire film thickness, such as those shown schematically in Fig. 23. Such domains are under pressure from their excessive wall energies to collapse and can readily be erased with a moderate-power laser pulse. DOW on exchange-coupled media is thus achieved by writing (i.e., creating reverse domains) with a high-power pulse, and erasing (i.e., eliminating domains) with a moderate-power pulse. An external magnetic field is usually required for writing on such media, but neither the direction nor the magnitude of this field needs to change during erasure.

Optical recording is an evolving technology, which will undoubtedly see many innovations and improvements in the coming years. Some of the ideas and concepts described here will hopefully remain useful for some time to come, while others may have a shorter lifetime and limited usefulness. It is the author's hope, however, that they all serve as a stepping-stone to profound new ideas.

35.10 FURTHER INFORMATION

Proceedings of the *Optical Data Storage Conference* are published annually by SPIE, the International Society for Optical Engineering. These proceedings document each year the latest developments in the field of optical recording.

Two other conferences in this field are the *International Symposium on Optical Memory* (ISOM) whose proceedings are published as a special issue of the *Japanese Journal of Applied Physics*, and the *Magneto-Optical Recording International Symposium* (MORIS) whose proceedings appear in a special issue of the *Journal of the Magnetics Society of Japan*.

35.11 BIBLIOGRAPHY

1. A. B. Marchant, *Optical Recording*, Addison-Wesley, Boston, Massachusetts, 1990.
2. P. Hansen and H. Heitmann, "Media for Erasable Magneto-Optic Recording," *IEEE Trans. Mag.* **25**:4390–4404 (1989).
3. M. H. Kryder, "Data-Storage Technologies for Advanced Computing," *Scientific American* **257**:116–125 (1987).
4. M. H. Kryder, "Data Storage in 2000—Trends in Data Storage Technologies," *IEEE Trans. Mag.* **25**:4358–4363 (1989).
5. A. E. Bell, "Optical Data Storage Technology: Status and Prospects," *Computer Design*, pp. 133–146 (1983).
6. D. S. Bloomberg and G. A. N. Connell, "Magneto-Optical Recording," *Magnetic Recording*, vol. II, C. D. Meo and E. D. Daniel (eds.), McGraw-Hill, New York, 1988, chap. 6.
7. S. Miyaoka, "Digital Audio is Compact and Rugged," *IEEE Spectrum*, pp. 35–39 (1984).
8. G. Bouwhuis and J. J. M. Braat, "Recording and Reading of Information on Optical Disks," *Applied Optics and Applied Engineering*, vol. IX, R. R. Shannon and J. C. Wyant (eds.), Academic Press, New York, 1983, chap. 3.
9. G. Bouwhuis, J. Braat, A. Huijser, J. Pasman, G. Van Rosmalen, and K. S. Immink, *Principles of Optical Disk Systems*, Adam Hilger Ltd., Bristol and Boston, 1985, chaps. 2 and 3.
10. H. H. Hopkins, "Diffraction Theory of Laser Read-Out Systems for Optical Video Discs," *J. Opt. Soc. Am.* **69**:4–24 (1979).
11. Special issue of *Applied Optics* on video disks, July 1, 1978.
12. M. Mansuripur, "Certain Computational Aspects of Vector Diffraction Problems," *J. Opt. Soc. Am.* **A6**:786–805 (1989).
13. P. Sheng, "Theoretical Considerations of Optical Diffraction from RCA Video Disk Signals," *RCA Review* **39**:513–555 (1978).
14. D. O. Smith, "Magneto-Optical Scattering from Multilayer Magnetic and Dielectric Films," *Opt. Acta* **12**:13 (1985).
15. P. S. Pershan, "Magneto-Optic Effects," *J. Appl. Phys.* **38**:1482–1490 (1967).
16. K. Balasubramanian, A. S. Marathay, and H. A. Macleod, "Modeling Magneto-Optical Thin Film Media for Optical Data Storage," *Thin Solid Films* **164**:341–403 (1988).
17. Y. Tomita and T. Yoshini, "Optimum Design of Multilayer-Medium Structures in a Magneto-Optical Readout System," *J. Opt. Soc. Am.* **A1**:809–817 (1984).
18. K. Egashira and T. Yamada, "Kerr Effect Enhancement and Improvement of Readout Characteristics in MnBi Film Memory," *J. Appl. Phys.* **45**:3643–3648 (1974).
19. G. A. N. Connell, "Interference Enhanced Kerr Spectroscopy for Very Thin Absorbing Films," *Appl. Phys. Lett.* **40**:212 (1982).
20. K. Y. Ahn and G. J. Fan, "Kerr Effect Enhancement in Ferromagnetic Films," *IEEE Magnet.* **2**:678 (1966).
21. H. Wieder and R. A. Burn, "Direct Comparison of Thermal and Magnetic Profiles in Curie Point Writing on MnGaGe Films," *J. Appl. Phys.* **44**:1774 (1973).
22. P. Kivits, R. deBont, and P. Zalm, "Superheating of Thin Films for Optical Recording," *Appl. Phys.* **24**:273–278 (1981).
23. M. Mansuripur and G. A. N. Connell, "Laser-Induced Local Heating of Moving Multilayer Media," *Appl. Opt.* **22**:666 (1983).
24. J. P. J. Heemskerck, "Noise in a Video Disk System: Experiments with an (AlGa) As Laser," *Appl. Opt.* **17**:2007 (1978).
25. G. A. Acket, D. Lenstra, A. J. DenBoef, and B. H. Verbeek, "Influence of Feedback Intensity on Longitudinal Mode Properties and Optical Noise in Index-Guided Semiconductor Lasers," *IEEE J. Quant. Electron.* **QE-20**:1163 (1984).
26. A. Arimoto, M. Ojima, N. Chinone, A. Oishi, T. Gotoh, and N. Ohnuki, "Optimum Conditions for the High Frequency Noise Reduction Method in Optical Video Disk Players," *Appl. Opt.* **25**:1398 (1986).
27. M. Ojima, A. Arimoto, N. Chinone, T. Gotoh, and K. Aiki, "Diode Laser Noise at Video Frequencies in Optical Video Disk Players," *Appl. Opt.* **25**:1404 (1986).

28. J. W. Beck, "Noise Consideration of Optical Beam Recording," *Appl. Opt.* **9**:2559 (1970).
29. D. Treves and D. S. Bloomberg, "Signal, Noise, and Codes in Optical Memories," *Opt. Eng.* **25**:881 (1986).
30. D. Treves, "Magneto-Optic Detection of High-Density Recordings," *J. Appl. Phys.* **38**:1192 (1967).
31. B. R. Brown, "Readout Performance Analysis of a Cryogenic Magneto-Optical Data Storage System," *IBM J. Res. Dev.* pp. 19–26 (1972).
32. R. L. Aagard, "Signal to Noise Ratio for Magneto-Optical Readout from MnBi Films," *IEEE Trans. Mag.* **9**:705 (1973).
33. M. Mansuripur, G. A. N. Connell, and J. W. Goodman, "Signal and Noise in Magneto-Optical Readout," *J. Appl. Phys.* **53**:4485 (1982).
34. B. G. Huth, "Calculation of Stable Domain Radii Produced by Thermomagnetic Writing," *IBM J. Res. Dev.*, pp. 100–109 (1974).
35. D. H. Howe, "Signal-to-Noise Ratio for Reliable Data Recording," *Proc. SPIE* **695**:255–261 (1986).
36. K. A. S. Immink, "Coding Methods for High-Density Optical Recording," *Philips J. Res.* **41**:410–430 (1986).
37. R. Hasegawa (ed.), *Glassy Metals: Magnetic, Chemical and Structural Properties*, CRC Press, Boca Raton, Florida, 1983.

INDEX

Index note: The *f* after a page number refers to a figure, the *n* to a note, and the *t* to a table.

- Abbe numbers, **29.36**
and axial chromatic aberrations, **17.22**
of binary optics, **23.6, 23.6f**
in gradient index optics, **24.3, 24.7**
of molded microlenses, **22.12, 22.12t**
of reflective and catadioptric objectives, **29.17**
- Abbe sine condition, **29.36**
and reflective and catadioptric objectives, **29.34**
of stigmatic conditioning, **1.30–1.31, 17.10**
- Abbe-Porter experiments, **11.1**
- Abbe's prisms, **19.3t, 19.7, 19.7f–19.8f**
- ABC model, of surface finish, **8.14–8.15**
- Abel transform, **8.13**
- Abelès method, **12.11–12.12**
- Aberrated wavefronts, **2.12, 2.13**
- Aberration(s):
in binary optics, **23.4–23.7, 23.5f, 23.6f**
chromatic, **1.91–1.92**
defined, **1.28**
defocus as, **1.85–1.86**
and general aspheres, **29.3**
higher-order, **29.37**
in instrumental optics, **1.85**
of point images, **1.85–1.92, 1.88f**
polarization, **15.35–15.37, 15.35f–15.37f**
pupil, **1.76**
ray, **1.87–1.88, 1.88f**
Seidel, **29.38**
spherical, **1.90, 29.7, 29.8, 29.15, 29.21, 29.37, 29.38**
and stop position, **1.92**
and stop size, **1.92**
in systems with rotational symmetry, **1.89–1.90**
third-order, **1.90–1.91, 29.38**
transverse primary chromatic, **17.22**
wavefront, **1.86–1.88, 1.86f**
- Aberration theory, in gradient index optics, **24.3**
- Ablation, microscopes, **28.54**
- Absolute instruments, **1.29**
- Absorption cross section, **7.5**
- Absorption cross section, **31.3**
- Absorption index, **12.6**
- Acceptance (étendue), **1.22, 1.81, 13.7**
- Access time, for optical disk data, **35.6**
- Accessories, for cameras, **25.16–25.17, 25.18f**
- Achromatic doublets (lenses), **17.22–17.25, 17.23f–17.25f, 17.24t**
- Achromatic retardation plates, **13.48–13.52, 13.50f, 13.53t**
- Acousto-optic cells, **11.8–11.9, 11.9f**
- Acousto-optic correlators, **11.10–11.12, 11.11f**
- Acousto-optic scanners, **30.44–30.45**
- Active autofocus systems, for cameras, **25.11–25.12, 25.12f**
- Active devices:
fabrication of, **21.14**
for integrated optics, **21.25–21.31, 21.26f–21.31f**
- Active layer removal, in PIC manufacturing, **21.19**
- Active pixel sensors, **26.2, 26.8–26.9, 26.8f**
- Active scanning, **30.4**
- Active-passive transitions, in PICs, **21.19–21.20**
- Addition, as analog operation, **11.2**
- Aerial cameras, **25.20**
- Aerial images, **1.26**
- Aerosol Polarimeter Sensor, **15.38**
- Affine transformations, **1.57**
- Afocal Cassegrain-Mersenne telescope objective, **29.9**
- Afocal Gregorian-Mersenne telescope objective, **29.12**
- Afocal lenses, **18.1–18.22**
for binoculars, **18.13–18.14, 18.14f**
catadioptric, **18.21–18.22, 18.22f**
Galilean, **18.15–18.17, 18.15f, 18.16f**
Gaussian, **1.45, 1.46f, 1.53–1.54, 1.53f, 1.54f**

- Afocal lenses (*Cont.*):
- Gaussian analysis, **18.4–18.6**
 - and focusing lenses, **18.2–18.4, 18.3f, 18.5f**
 - and optical invariant, **18.7**
 - subjective aspects of, **18.6, 18.7f**
 - Keplerian, **18.7–18.14**
 - and eye relief manipulation, **18.8–18.10, 18.9f, 18.10f**
 - field-of-view limitations in, **18.11**
 - finite conjugate afocal relays, **18.11–18.12, 18.12f**
 - thin-lens model of, **18.7–18.8, 18.8f**
 - paraxial matrix methods, **1.70**
 - for periscopes, **18.19, 18.19f**
 - reflecting, **18.19–18.21, 18.20f, 18.21f**
 - in relay trains, **18.17–18.19, 18.17f, 18.18f**
 - for scanners, **18.13, 18.13f**
 - for telescopes, **18.10–18.11, 18.11f**
- Afocal magnification, **18.5–18.6**
- Afocal objectives, **29.9, 29.12, 29.29–29.30**
- Agile beam steering, **30.51–30.63**
- with decentered-lens and microlens arrays, **30.57–30.60, 30.58f–30.60f, 30.62–30.63**
 - with digital micromirror devices, **30.60–30.61**
 - with gimbal-less two-axis scanning micromirrors, **30.61–30.62, 30.62f**
 - phased-array, **30.52–30.57, 30.53f, 30.62–30.63**
- Ahrens method of spar cutting, **13.12**
- Ahrens Nicol prisms, **13.16f, 13.17**
- Ahrens prisms, **13.14**
- Aircraft, synthetic aperture radar for, **11.6–11.7, 11.7f**
- Airy disks:
- in confocal microscopy, **28.50**
 - defined, **3.26**
 - of DIC microscopes, **28.39**
 - of microscopes, **28.17–28.19, 28.18f, 28.19f**
 - of solid-state cameras, **26.15**
 - and vector diffraction, **3.33**
- Airy equation, **12.10**
- Airy pattern, **17.38**
- Airy-Drude formula, **16.5**
- Altenhof objectives, **29.32–29.33**
- Alternative Paul three-mirror objective, **29.28**
- Alvarez plates, **22.16**
- Alvarez-Humphrey plates, **22.37**
- Amici lenses, **17.10, 17.10f**
- Amici prisms, **19.3t, 19.11, 19.12f, 20.6f**
- Amplitude, of waves, **2.4, 2.5, 12.5**
- Amplitude division, interference by, **2.14, 2.19–2.28**
- and extended sources, **2.20**
 - and Fizeau interferometers, **2.24–2.26, 2.25f**
 - and fringes of equal inclination, **2.20–2.22, 2.21f, 2.22f**
 - and fringes of equal thickness, **2.22–2.24, 2.23f**
 - and Michelson interferometers, **2.26–2.28, 2.26f–2.27f**
 - plane-parallel plate, **2.19, 2.20f, 2.30–2.33, 2.30f, 2.32f, 2.33f**
 - thin films, **2.24**
- Amplitude penetration depth, **12.5**
- Amplitude reflection coefficients, Fresnel, **12.7–12.8, 12.10**
- Amplitude scattering matrix, **7.10, 7.13**
(*See also* Jones matrix)
- Amplitude transmission coefficients, Fresnel, **12.8**
- Amplitude-shift-keyed (ASK) transmission, **21.30**
- Analog optical signal and image processing, **11.1–11.20**
- Fourier transforms in, **11.3–11.5, 11.3f, 11.5f**
 - and fundamental analog operations, **11.2–11.3**
 - incoherent processing, **11.17–11.20, 11.18f, 11.19f**
 - and spatial filtering, **11.5–11.6, 11.6f**
 - of synthetic aperture radar data, **11.6–11.8, 11.7f–11.8f**
 - of temporal signals, **11.8–11.12, 11.9f–11.11f**
 - of two-dimensional images, **11.12–11.17, 11.13f**
- Analog transmission, **21.32–21.34, 21.33f, 21.34f**
- Analytical signal representation, in coherence theory, **5.2–5.3**
- Analyzed states, of polarizers, **15.19**
- Analyzer vectors, **15.11**
- Anamorphic afocal attachments (anamorphosers), **18.16–18.17, 18.16f**
- Anamorphic error control, **30.49, 30.50**
- Anastigmatic (term), **29.36**
- Anastigmatic objectives, **29.12–29.13**
- Angle characteristic function, **1.14–1.15, 1.15f, 1.17**
- Angle of incidence, **1.23, 1.39, 13.48**
- Angle-point characteristic function, **1.16, 1.17**
- Angular change, of optical beams, **30.5**
- Angular correction function, **5.6, 5.7f**
- Angular magnification, **1.52, 1.78, 18.4**
- Angular scan, **30.28–30.29**
- Angular spectrum representation, **5.14–5.15, 5.14f**

- Anisotropy, perpendicular magnetic, 35.26
- Annealed proton exchange (APE) process:
 for fiber optic gyroscopes, 21.35, 21.36, 21.36*t*
 for LiNbO₃ waveguides, 21.16–21.17
- Anomalous diffraction, 7.5, 7.6*f*
- Antinodal points, of lens systems, 17.7
- Antiprincipal points, of lens systems, 17.7
- Aperture(s), 1.74–1.76, 1.75*f*
 circular
 diffraction of light from, 3.6–3.7, 3.6*f*, 3.7*f*, 3.9–3.11
 Fraunhofer patterns for, 3.25, 3.26, 3.27*f*
 double-slit, 3.26–3.28, 3.27*f*, 3.28*f*
 image space numerical, 1.79
 linear, 30.54
 rectangular, 3.19–3.20
 Fraunhofer patterns for, 3.25, 3.26
 Fresnel diffraction from, 3.19–3.20, 3.20*f*
 uniformly illuminated, 30.9, 30.10, 30.10*f*
- Aperture, numerical, 1.79, 17.9
- Aperture stops, 1.74, 1.75*f*, 17.8, 29.5, 29.36
- Aperture-scanning microscopy, 28.53–28.54
- Aplanatic (term), 29.36
- Aplanatic lenses, 17.5, 17.11–17.12, 17.11*f*
- Aplanatic objectives, 29.11–29.13
- Arbitrary phase profiles, 23.8
- Arbitrary systems, paraxial matrix methods for, 1.67
- Area-solid-angle-product, 1.22
 (See also Étendue)
- Arrayed waveguide gratings (AWGs), 21.24
- Aspherical surfaces:
 and axial gradients, 24.3
 and reflective/catadioptric objectives, 29.3
 in systems of revolution, 1.35
- Astigmatic difference (term), 1.43
- Astigmatism, 1.90, 1.91, 29.34, 29.37
- Astronomical telescopes, 18.10
- Astronomy, radio, 5.23
- Atmospheric particles, scattering by, 7.2
- Atoms:
 electronic structure of, 10.12–10.16, 10.13*f*–10.15*f*
 multielectron, 10.10–10.11
 one-electron, 10.7–10.9, 10.8*f*, 10.9*f*
 spectra of, 10.3
- Augmented resolution, 30.12–30.14, 30.13*f*
- Autofocus, of cameras, 25.11–25.15, 25.12*f*–25.15*f*
- Autofocus SLRs, 25.12–25.14, 25.13*f*
- Automatic focusing, optical disks and, 35.12–35.14, 35.13*f*
- Automatic tracking, on optical disks, 35.14–35.17, 35.14*f*–35.16*f*
- Average degree of polarization (average DoP), 14.32–14.33
- Axial astigmatism, 29.34, 29.37
- Axial color, 1.91, 29.9, 29.37
- Axial gradient lenses, 24.3–24.5, 24.4*f*
- Axial image point, 1.27
- Axial (longitudinal) magnification, 1.28, 1.52, 17.5
- Axicons, 11.7, 11.7*f*
- Axis wander, of prisms, 13.15
- AxoScan Mueller matrix polarimeters, 15.33
- Azimuth, 15.10, 16.16
- Babinet compensators, 13.53–13.55, 13.54*f*
- Babinet principle, 3.9–3.11, 3.10*f*, 3.11*f*, 3.13
- Babinet-Soleil compensators, 13.55–13.56, 13.55*f*
- Back focal length (BFL), of camera lenses, 27.2, 27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
- Backscatter and backscattering:
 coherent, 9.14–9.15, 9.14*f*
 enhanced, 6.5–6.7, 6.5*f*
- Baker relays, 18.18, 18.18*f*
- Baker super-Schmidt objective, 29.21
- Baker-Nunn objective, 29.22
- Balmer α -spectra, 10.7–10.8, 10.8*f*
- Balmer β -transition, 10.9
- Barrel distortion, 1.91
- Barrel length (BRL), of camera lenses, 27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
- Beam deviation and displacement, 19.2
- Beam diffusers, 23.13
- Beam propagation method (BPM), 21.8
- Beam shapers and shaping:
 in binary optics, 23.13
 for optical disks, 35.8, 35.9*f*, 35.10*f*
- Beam splitters, polarizing, 13.41–13.42
- Beam steering (see Agile beam steering)
- Beam-splitter gratings, 23.11, 23.12
 Dammann approach, for binary gratings, 23.12
- Beam-splitter prisms, 13.18–13.22
 Foster, 13.7, 13.18*f*, 13.21–13.22
 Glan-Thompson, 13.18*f*, 13.22
 Rochon, 13.7, 13.18–13.21, 13.18*f*, 13.24
 Sénarmont, 13.7, 13.18, 13.18*f*, 13.21
 Wollaston, 13.7, 13.18, 13.18*f*, 13.21, 13.24

- Becke line, **28.27**
- Beilby-layer polarizers, **13.28**
- Benes architecture, for switches, **21.34–21.35**
- Berremann calculus, **15.32**
- Bertrand lenses, **28.8, 28.44f**
- Bertrand-type Feussner prisms, **13.23**
- Bessel functions, **7.12, 7.14–7.15, 17.38, 28.17**
- Bidirectional reflectance distribution function (BRDF), **8.4, 15.38–15.39, 15.38f–15.40f**
- Bifocal lenses, **23.12, 23.12f**
- Billet's split lens, **2.16, 2.17f**
- Binary optics, **23.1–23.17**
- fabrication of
 - mask layout, **23.14–23.16, 23.15f, 23.15t**
 - micromachining techniques, **23.16, 23.16f, 23.17**
 - and geometrical optics, **23.2–23.9**
 - aberration correction, **23.4–23.7, 23.5f, 23.6f**
 - analytical models, **23.2–23.4, 23.6**
 - micro-optics, **23.7–23.8, 23.7f–23.8f**
 - optical performance, **23.8–23.9, 23.9f, 23.10t**
 - and scalar diffraction theory, **23.10–23.13, 23.11t, 23.12f, 23.13f**
 - and vector diffraction theory, **23.13–23.14, 23.14f**
- Binoculars, **18.13–18.14, 18.14f**
- Binomial vectors, of space curves, **1.19**
- Biotar lenses, **17.28**
- Biplates, **13.56**
- Birefringence, **15.6, 15.41**
- Bistatic radar cross-section (RCS), in surface scattering, **8.3, 8.4**
- Blocks (optical disk data), **35.6–35.7**
- Bohr frequency condition, **10.4**
- Bohr's theory of hydrogen, **10.3**
- Born series, **9.4**
- Born-Oppenheimer approximation, **10.19, 10.20, 10.22**
- Bound modes, of optical waveguides, **21.3**
- Brace half-shade plates, **13.57**
- Bragg cell spectrum analyzers, **11.9–11.10, 11.10f**
- Bragg cells, **11.11–11.12, 11.19**
- Bragg diffraction, **11.9, 11.9f**
- Bragg reflection filters, **21.30**
- Bragg reflectors, **21.8, 21.30**
- Bragg regime, **30.39–30.41, 30.41f, 30.42f**
- Brashear-Hastings prisms, **19.3t, 19.25, 19.25f**
- Bravais biplates, **13.56**
- Brewster angle:
 - and Airy equation, **12.10**
 - defined, **12.12**
 - and extinction ratio, **12.18, 12.22, 12.22f, 12.23f, 12.24**
 - and polarization, **12.15**
- Brewster angle prisms, **13.13**
- Brewster angle reflection polarizers, **12.16–12.18, 12.16f, 12.17f, 13.34–13.37, 13.34t–13.36t**
- Brewster angle transmission polarizers, **12.18–12.24, 12.19t–12.20t, 12.21f, 13.37–13.39, 13.38t–13.39t**
- Bright field microscopy, **28.25, 28.27–28.28, 28.27f**
- Brillouin scattering, **31.30**
- Broadening, of lineshapes, **10.7**
- Brownian fractals, **8.9, 8.17**
- Brownian movement, **28.28**
- Buckbee Mears wire-grid polarizers, **13.31**
- Bunsen-Kirchhoff spectrometers, **20.5f**
- Butterfly scanners, **30.50f, 30.51**
- Cable television (CATV), **21.2, 21.32–21.34**
- Caged compounds, in microscopy, **28.55**
- Calcite:
 - double refraction in, **13.2–13.6, 13.2f–13.3f, 13.4t–13.5t**
 - Feussner prisms of, **13.23**
 - Rochon prisms of, **13.20**
- Camera formula, for solid-state cameras, **26.13**
- Camera lenses:
 - classification system for, **27.17, 27.23t, 27.24**
 - design limitations of, **27.1–27.2**
 - fish-eye, **27.6**
 - inverted telephoto, **27.2, 27.6, 27.7f–27.14f**
 - extreme wide-angle, **27.6, 27.13f, 27.14f**
 - highly complex extreme speed, **27.6, 27.9f–27.12f**
 - very compact moderate speed, **27.6, 27.7f–27.8f**
 - SLR normal lenses, **27.2, 27.3f–27.4f**
 - telephoto lenses, **27.6, 27.13, 27.15f–27.16f**
 - wide-angle lenses, **27.2, 27.5f**
 - zoom lenses, **27.17, 27.20f–27.22f**
- Cameras, **25.3–25.26**
- accessories for, **25.16, 25.17, 25.17f**
 - aerial, **25.20**
 - and autoexposure, **25.10–25.11**
 - and autofocus, **25.11–25.15, 25.12f–25.15f**
 - characteristics of, **25.3–25.4**
 - clandestine, **25.21**

- Cameras (*Cont.*):
 critical features of, 25.8, 25.9f
 display types, 25.6–25.7
 endoscopic, 25.21, 25.21f
 features of, 25.8, 25.18, 25.19f
 film for, 25.5, 25.6
 and flash, 25.16, 25.17f
 formats for, 25.18
 high-speed, 25.21–25.22, 25.22f
 and image, 25.5
 images from, 25.7
 and instant pictures, 25.8
 large-format, 25.18–25.20
 and red eye, 25.16
 and resolution of fine detail, 25.5–25.6, 25.6f
 sewer, 25.22–25.23
 solid-state (*see* Solid-state cameras)
 stereo, 25.23–25.24, 25.23f
 streak, 25.24, 25.24f
 thermal imaging, 25.25
 and time lag, 25.8–25.9, 25.9f
 underwater, 25.25
 video, 25.7–25.8
 view, 25.18–25.20, 25.19f
 for wide-angle photography, 25.25, 25.26f
 Canon EOS A2E camera, 25.14f, 25.15, 25.16, 25.15f
 Cantor sets, 8.9
 Cardinal points, of lenses, 1.44, 17.7
 Carl Zeiss prism system, 19.3f, 19.16, 19.16f
 Carlisle objectives, 29.8
 Carlson, Chester, 34.1
 Carrier effects, in integrated optics, 21.10–21.12
 Carrier-to-noise ratio (CNR), 35.24
 Cartesian coordinates, 1.20, 1.21
 Cassegrain objectives, 29.6, 29.7
 afocal Cassegrain-Mersenne telescope, 29.9
 dual magnification, 29.9–29.10
 with field corrector and spherical secondary, 29.8–29.9
 Houghton-Cassegrain, 29.22–29.23
 Mangin-Cassegrain with correctors, 29.24
 reflective Schmidt-Cassegrain, 29.17
 Schmidt-Cassegrain, 29.16–29.17
 Schmidt-meniscus Cassegrain, 29.21
 with Schwarzschild relay, 29.32
 solid Makutsov-Cassegrain, 29.19
 spherical-primary, with reflective field corrector, 29.9
 three-mirror, 29.30
 Cassegrainian telescopes, 18.21
 Catadioptric (term), 29.37
 Catadioptric Herschelian objective, 29.27
 Catadioptric lenses:
 afocal, 18.21–18.22, 18.22f
 systems of, 1.9
 Catadioptric objectives (*see* Reflective and catadioptric objectives)
 Cathode ray tube (CRT) monitors, 25.6–25.7
 Cathode-ray tubes (CRTs), 30.4, 30.25–30.26
 Catoptric systems, 1.9
 Caustics, ray densities and, 1.88
 Center, of afocal lens, 1.54
 Channels, for wave propagation, 9.16
 Characteristic functions (geometrical optics), 1.13–1.18
 angle characteristic function, 1.14–1.15, 1.15f, 1.17
 angle-point characteristic function, 1.16, 1.17
 and expansions about rays, 1.16
 and expansions about the axis, 1.16–1.17
 ideal, 1.17–1.18
 mixed, 1.13
 paraxial forms of, 1.17
 and paraxial matrices, 1.74
 point characteristic function, 1.14
 point eikonal, 1.14, 1.17
 point-angle characteristic function, 1.15–1.17
 Charge injection devices (CIDs), 26.6–26.7, 26.6f–26.8f
 Charge spectrographs, 34.8, 34.8f
 Charge-coupled devices (CCDs), 25.7, 26.3–26.5, 26.4f–26.6f
 Charged area development (CAD), in xerographic systems, 34.4
 Chebyshev polynomials, 7.15
 Chemical beam epitaxy (CBE), 21.17, 21.18
 Chief rays, 1.75, 17.8, 29.20, 29.37
 Chiolite, 13.42
 Chiral particles, scattering by, 7.2
 Cholesky decomposition, 14.41, 14.42
 Chromatic aberration correction, 23.5–23.6, 23.6f
 Chromatic aberrations, 1.91–1.92
 Chromatism, of axial gradients, 24.3–24.6
 Circle (Zernike) polynomials, 1.90, 23.3
 Circuits, in integrated optics, 21.21–21.31
 for active devices, 21.25–21.31, 21.26f–21.31f
 for passive devices, 21.21–21.25, 21.22f–21.25f
 Circular analyzers, of polarized light, 15.18, 15.19

- Circular apertures:
 diffraction of light from, 3.6–3.7, 3.6f, 3.7f, 3.9–3.11
 Fraunhofer patterns for, 3.25, 3.26, 3.27f
- Circular polarizers, 15.17–15.19
- Circular scan, 30.16, 30.18f
- Clandestine cameras, 25.21
- Clausius-Mosotti theory, 7.16
- Cleaning, in xerographic systems, 34.10
- Clebsch-Gordon coefficients, 31.17
- Coated spheres, scattering by, 7.14
- Coddington's equations, 1.44
- Coefficient of finesse (interference), 2.31
- Coercivity, of optical disk data, 35.17n, 35.27, 35.27f, 35.28
- Cogging, of streak cameras, 25.24
- Coherence, 5.1–5.23, 6.2–6.13
 analytical signal representation, 5.2–5.3
 applications of, 5.22–5.23
 in binary optics, 23.7, 23.8f
 classical, 5.1–5.2
 coherence area, 5.3
 and coherence functions, 5.4–5.9
 angular correction function, 5.6, 5.7f
 complex degree of coherence, 5.4
 complex degree of spectral coherence, 5.5
 cross-spectral density function, 5.5
 efficient sampling of, 6.10–6.12, 6.11f
 higher-order functions, 5.8–5.9
 intensity, 5.7
 mutual coherence function, 5.4
 radiance, 5.8
 radiant emittance, 5.7–5.8
 radiant intensity, 5.8
 spectrum and normalized spectrum, 5.5–5.6
 coherence time, 5.3
 complex degree of, 2.37, 5.4
 and enhanced backscatter, 6.5–6.7, 6.5f
 and general linear systems, 6.3–6.4
 and image formation, 6.9–6.10, 6.9f
 and interference, 2.13, 2.36–2.42
 laser sources, 2.41–2.42, 2.42f
 Michelson stellar interferometers, 2.40–2.41, 2.40f
 mutual coherence function, 2.36–2.38, 2.36f
 spatial coherence, 2.38–2.40, 2.38f–2.39f
 temporal coherence, 2.41
 and Koehler-illumination, 6.12–6.13, 6.12f
 and laser modes, 5.23
 and Lau effect, 6.7–6.8, 6.8f
 of light sources, 5.9–5.13
- Coherence (*Cont.*):
 and Lukosz-type super-resolving systems, 6.9–6.10, 6.9f
 measurements of coherence, 5.3–5.4
 and noncosmological red shift, 5.23
 and optical image enhancement, 11.14–11.17
 partial, 2.38
 and polarization effects, 5.22
 propagation in, 5.13–5.19, 5.14f–5.16f, 6.4
 and radio astronomy, 5.23
 scalar field amplitude, 5.3
 spatial coherence, 5.3
 and speckle, 5.22
 and spectral representation, 5.22
 and spectrum of light, 5.19–5.22, 5.20f
 and statistical radiometry, 5.22
 temporal, 2.41, 5.3
 time averages in, 6.4–6.5
- Coherence area, 5.3
- Coherence length, 2.19
- Coherence time, 5.3
- Coherence volume, 5.3
- Coherency matrix, 12.29–12.30, 14.41
- Coherent arrays, scattering by, 7.2–7.3
- Coherent backscattering, 9.14–9.15, 9.14f
- Coherent mode representation (spectrum of light), 5.20–5.21
- Coherent optical image enhancement, 11.14–11.17
- Coherent radiation, 30.2, 30.25–30.26
- Coherent scattering, 7.3, 9.2, 9.3, 9.5–9.7, 9.6f
- Collimation, 35.8, 35.9f, 35.10f
- Collineation, 1.56–1.63
 of conjugate lines, 1.59
 of conjugate planes, 1.58–1.59
 coordinate systems and degrees of freedom for, 1.57
 equations of, 1.57–1.58
 general properties of, 1.62–1.63
 matrix representation of, 1.59–1.60
 of rotationally symmetric lenses, 1.60–1.62, 1.62f
- Colocalization, single molecule high-resolution, 28.23
- Color:
 axial, 1.91, 29.9, 29.37
 in human visual system, 26.18
 lateral, 1.91–1.92
 in xerographic systems, 34.11–34.12, 34.11f–34.13f

- Color filter arrays (CFAs), in solid-state cameras, 26.18
- Coma, 1.31, 24.7, 29.37
- Coming's glass molding process, 22.9, 22.9f
- Compact disks (CDs), 35.1–35.2, 35.5n
- Compensators, 13.53–13.56, 13.54f, 13.55f, 28.38
- Complementary aperture screens, 3.9–3.11
- Complementary metal-oxide semiconductors (CMOSs), 26.8–26.9
- Complex amplitude, 2.4, 2.5, 11.2
- Complex degree of coherence, 2.37, 5.4
- Complex degree of spectral coherence, 5.5
- Complex refractive index, 7.12–7.13, 12.5, 12.6
- Composite retardation plates, 13.52, 13.53
- Compound microscopes, 17.10
- Compound mirror optics configurations, 30.15–30.16, 30.15f
- Conductance channels, in speckle patterns, 9.16
- Conductive magnetic brush (CMB), 34.6
- Conductivity, diffraction and, 3.32–3.33
- Configurational coordinate model, of lineshapes, 10.22, 10.23f
- Confocal Cassegrainians (telescopes), 18.21
- Confocal microscopy, 28.49–28.51, 28.49f, 28.51f
- Confocal parabolas, 18.20, 18.20f, 18.21
- Conic constant (term), 1.34, 29.37
- Conic mirrors, 29.3f, 29.4f
- Conical surfaces, in systems of revolution, 1.34–1.35
- Conjugate lines, collineation of, 1.59
- Conjugate matrices, 1.68–1.71, 1.73
- Conjugate planes: collineation of, 1.58–1.59 in microscopes, 28.4–28.5
- Conoscopic imaging, 28.8–28.9
- Conservation of étendue, law of, 1.22
- Constructive interference, 2.7
- Continuous wave (CW) dye lasers, 10.8
- Contrast, in microscopy: bright field microscopy, 28.25, 28.27–28.28, 28.27f dark field microscopy, 28.28 Hoffman modulation contrast, 28.29 interference microscopy, 28.33–28.44, 28.35f, 28.37f, 28.38f, 28.40f, 28.42f, 28.43f and modulation transfer function, 28.24–28.25, 28.25f, 28.26f phase contrast, 28.28–28.29, 28.29f SSEE microscopy, 28.29, 28.30, 28.30f–28.33f, 28.33
- Contrast transfer function (CTF), 4.7–4.8, 4.8f, 28.24
- Cook objectives, 29.31
- Coordinate systems: for aberrations of point images, 1.86 for collineation, 1.57 for Fresnel equations, 12.6–12.7, 12.7f left-handed, 12.6 for Mueller matrices, 14.19–14.20
- Cornu's spiral, 3.16–3.19, 3.18f
- Corotron, in xerographic systems, 34.2, 34.3f
- Correctors, for reflective and catadioptric objectives: aplanatic, anastigmatic Schwarzschild with aspheric corrector plate, 29.13 Cassegrain with spherical secondary and field corrector, 29.8–29.9 Mangin-Cassegrain with correctors, 29.24 Ritchey-Chretien telescope with two-lens corrector, 29.8 spherical-primary Cassegrain with reflective field corrector, 29.9 three-lens prime focus corrector, 29.10
- Correlated double sampling (CDS), 26.11
- Correlators, acousto-optic, 11.10–11.12, 11.11f
- Cosine condition, of stigmatic imaging, 1.30, 1.30f
- Cosine-to-the-fourth approximation, 1.81
- Cotton polarizers, 13.21
- Couder objective, 29.12
- Coulomb repulsion, 10.10, 10.12, 10.16
- Coupled-dipole method, 7.15
- Couple-mode theory, in integrated optics, 21.8
- Coverslip correction, 28.10–28.11, 28.12f, 28.13
- Critical illumination, 28.7
- Cross-spectral density function, 5.5, 5.9, 5.10, 5.16
- Crystalline-quartz retardation plates, 13.46–13.48
- Curie temperature, 35.25
- Curvature: of space curves, 1.18–1.19 vertex (paraxial), 1.32–1.33
- Curved surfaces, radial gradients with, 24.7
- Cylinders, scattering by, 7.14
- Cylindrical lenses, 22.45, 22.46f
- Cylindrical wavefronts, 3.13–3.21, 3.14f and Cornu's spiral, 3.16–3.19, 3.18f, 3.19t and opaque strip construction, 3.20–3.21 from rectangular apertures, 3.19–3.20 from straight edges, 3.14–3.16

- Czerny-Turner monochromators, 31.6
 Czerny-Turner mounts, 20.8*f*, 20.14*t*
- Dall-Kirkham objective, 29.8
 Dammann approach, for binary gratings, 23.12
 Dark decay, of xerographic photoreceptors, 34.3
 Dark field microscopy, 28.28
 Data noise, 35.24
 Data rates, 30.6–30.8
 Data-reduction equations, polarimetric, 15.14–15.15
 Decentered lens arrays, for agile beam steering, 30.57–30.60, 30.58*f*–30.60*f*, 30.62–30.63
 Decomposition:
 of conjugate matrices, 1.71–1.72
 of Mueller matrices, 14.33
 Cholesky decomposition, 14.41, 14.42
 polar decomposition, 14.39–14.40
 SVD, 15.25–15.27
 Weyl's plane-wave, 3.23
 Deep reactive ion etching (DRIE), 22.7, 22.23, 30.62
 Deep x-ray lithography (DXRL), 22.7
 Defocus, 1.82, 1.83, 1.83*f*, 1.85–1.86, 1.90, 1.91
 Degenerate integrated structures, 21.12
 Degree of circular polarization (DoCP), 15.10
 Degree of coherence, 2.37
 Degree of linear polarization (DoLP), 15.10
 Degree of polarization (DoP), 12.14–12.15, 15.9
 Degree of polarization (DoP) surfaces and maps, 14.31–14.32, 14.32*f*
 Density (coherency) matrix, 12.29–12.30, 14.41
 Depolarization:
 defined, 15.7
 and diagonal depolarizers, 14.30
 Mueller matrices for, 14.30–14.39
 depolarization index, 14.32
 generators of, 14.33–14.39, 14.36*f*, 14.37*f*
 and volume scattering, 9.16–9.17, 9.17*f*
 depolarization index, 14.32
 Depth of field, 1.84, 17.37, 28.22–28.23
 Depth of focus, 1.84, 17.37, 17.37*f*, 28.22
 Derivative matrices, 1.73
 Derotation, of polygon scanners, 30.35–30.36
 Destructive interference, 2.7
 Development, in xerographic systems, 34.5–34.10, 34.5*f*–34.9*f*
- Diagonal depolarizers, 14.30
 Diallytes (lenses), 17.25, 17.25*f*
 Diamagnification, 1.23
 Diamond turning, 22.15–22.18, 22.16*t*, 22.17*f*, 22.18*f*
 Diattenuation and diattenuators, 14.6, 15.7
 linear, 14.8, 14.17
 Mueller matrices for, 14.16–14.19, 14.16*f*
 Dichroic polarizers, 13.24–13.33, 13.26*f*, 13.27*f*
 coatings as, 13.28
 measuring polarization of, 13.33
 pyrolytic-graphite polarizers, 13.28–13.29, 13.29*f*
 sheet polarizers, 13.25–13.28
 Dichroism, 15.19, 15.41, 31.17, 31.20, 31.21
 Dielectric impermeability, 21.9
 Differential detection, of optical disk data, 35.22–35.23, 35.22*f*
 Differential geometry, of rays, 1.19–1.21
 Differential scattering cross-sections (DSCs), 7.8, 8.3, 8.4
 Differential-interference contrast (DIC) microscopy, 28.27, 28.39–28.41, 28.40*f*
 Differential-phase-shift-keyed (DPSK) transmission, 21.30, 21.32
 Diffraction, 3.2–3.3, 3.3*f*
 anomalous, 7.5, 7.6*f*
 Bragg, 11.9, 11.9*f*
 from circular apertures, 3.6–3.7, 3.6*f*, 3.7*f*
 and Cornu's spiral, 3.16–3.17
 of cylindrical wavefronts, 3.13–3.21, 3.14*f*
 Cornu's spiral, 3.16–3.19, 3.18*f*, 3.19*t*
 opaque strip construction, 3.20–3.21, 3.21*f*
 from rectangular apertures, 3.19–3.20, 3.20*f*
 from straight edge, 3.14–3.16, 3.14*f*, 3.15*f*
 definition of, 3.6
 from disks, 3.7–3.8
 Fraunhofer, 3.24–3.28, 3.24*f*–3.26*f*
 Airy diffraction as, 28.17
 conducting screens for, 3.33*f*
 and gratings, 20.3
 Fresnel-Kirchhoff formula, 3.21, 3.22, 3.32
 Green's function, 3.22–3.23
 Huygens-Fresnel construction, 3.4–3.13
 Babinet principle, 3.9–3.11, 3.10*f*, 3.11*f*
 circular apertures and disks, light from, 3.6–3.9, 3.6*f*–3.9*f*
 Fresnel zones, 3.4–3.6
 zone plates, 3.11–3.13

- Diffraction (*Cont.*):
 mathematical theory of, 3.21–3.29
 diffraction grating, 3.28–3.29, 3.29f, 3.30f
 Fraunhofer diffraction, 3.24–3.28
 Fresnel and Fraunhofer approximations, 3.23–3.24
 Rayleigh-Sommerfeld, 3.9, 3.10, 3.23, 3.29
 and resolution of microscopes, 28.19–28.22, 28.20f, 28.21f
 scalar diffraction theory, 23.10–23.13, 23.11t, 23.12f, 23.13f
 by spheres, 7.4
 stationary phase approximation for, 3.29, 3.31–3.32
 vector, 3.32–3.37, 3.32f–3.37f, 23.13–23.14, 23.14f
- Diffraction efficiency, 23.9, 23.9f, 23.10t
- Diffraction gratings, 3.28–3.29, 3.29f, 3.30f, 20.4
 arrayed waveguide, 21.24
 beam-splitter, 23.11, 23.12
 binary, 23.13
 dispersive gratings vs., 20.3–20.4
 interference, 33.14
 multiple beam, 2.28–2.29, 2.29f, 2.30f
 for PICs, 21.19
 as polarizers, 13.30–13.33
- Diffraction patterns, three-dimensional, 28.19–28.22
- Diffraction-limited depth of focus, 28.22
- Diffraction-limited lenses, 17.37–17.39, 17.41f–17.42f
- Diffraction-type polarizers, 13.30–13.33
- Diffractional optics, 33.13, 35.28, 35.28f, 35.29
- Diffusers, 6.5, 6.5f, 6.7, 23.13
- Diffusion approximation, of radiative transfer, 9.11–9.12
- Digital holographic microscopy (DHM), 28.42, 28.43
- Digital light processing (DLP), 30.3, 30.60
- Digital micromirror devices (DMDs), 22.23, 30.60–30.61
- Digital transmission, in integrated optics, 21.31–21.32
- Digitized Green's function, 7.15, 7.16
- Dioptric systems, 1.9
- Dipole approximation, discrete, 7.16
- Dipole model of light, 3.33–3.36, 3.37f
- Dirac delta function, 5.9, 5.12, 6.8
- Dirac equation, for one-electron atom, 10.10
- Dirac series, 30.8
- Direct overwrite (DOW), of optical disks, 35.30, 35.30f
- Direct-vision prisms, 19.2
- Discharged area development (DAD), 34.4
- Discrete dipole approximation, 7.16
- Discrete signals, incoherent processing of, 11.17–11.20, 11.18f, 11.19f
- Disk rotation speed, of optical disks, 35.5–35.6
- Disks, as aperture screens, 3.7–3.11, 3.8f
- Dispersion, 20.1
- Dispersive prisms and gratings, 20.1–20.15
 configurations of, 20.4–20.15
 diffraction gratings vs., 20.3–20.4
 Eagle configuration, 20.7, 20.11f
 Ebert-Fastie configuration, 20.8, 20.12f
 Littrow configuration, 20.7f, 20.10
 Paschen-Runge configuration, 20.7, 20.11f
 Pfund configuration, 20.8f, 20.10, 20.13f
 in spectrometers, 20.2–20.3, 20.3f
 in spectroradiometers, 20.1, 20.2f, 20.14t
 Wadsworth configuration, 20.5f, 20.8, 20.12f
- Displays, of cameras, 25.6–25.7
- Distances:
 in Gaussian lenses, 1.51–1.53
 hyperfocal, 1.85
- Distortion, 29.37
 barrel, 1.91
 nonrectilinear, 27.6
 of objectives, 29.6
 pincushion, 1.91
 pupil, 1.78
 rectilinear, 27.6, 27.13f, 27.14f
- Distortion-free focusing lenses, 18.4
- Distributed Bragg reflector (DBR) lasers, 21.25, 21.30, 21.32, 21.37, 21.37f
- Distributed feedback (DFB) lasers, 21.25, 21.29, 21.30, 21.32, 21.38, 21.42
- Distributed index of refraction, 24.1 [*See also* Gradient index (GRIN) optics]
- Distributed-index planar microlenses, 22.26–22.31, 22.27f–22.30f, 22.27t, 22.31t
- Disturbance of wavefront, 3.5–3.6, 3.14–3.15, 3.14f
- Division-of-amplitude photopolarimeter (DOAP), 16.15, 16.15f, 16.16
- Division-of-amplitude polarimeters, 15.5–15.6
- Division-of-aperture polarimeters, 15.5
- Division-of-wavefront photopolarimeter (DOWP), 16.14, 16.14f
- Doppler broadening, 10.2, 10.7, 31.23

- Doppler shift, 2.13, 5.23, 11.6, 31.30
- Double Dove prisms, 19.3*t*, 19.10, 19.10*f*
- Double refraction, in calcite, 13.2–13.6, 13.4*t*–13.5*t*
- Double-beam spectrophotometers, 31.4–31.5, 31.5*f*
- Double-Gauss lenses, 17.27–17.28, 17.28*f*, 27.2
- Double-pass monochromators, 20.9*f*
- Double-pass objective optics, 30.32–30.33, 30.32*f*
- Double-passed two-beam interferometers, 32.8
- Double-reflection error control, 30.50–30.51, 30.50*f*, 30.51*f*
- Double-slit apertures, 3.26–3.28, 3.27*f*, 3.28*f*
- Doublets, achromatic (lenses), 17.22–17.25, 17.23*f*–17.25*f*, 17.24*t*
- Dove prisms, 19.3*t*, 19.9, 19.9*f*, 19.10, 19.10*f*
- Drude model, 21.10
- Dual magnification Cassegrain objective, 29.9–29.10
- Dual rotating retarder polarimeters, 15.16, 15.16*f*
- Dummond, D. G., 13.47
- Duty cycle, for input/output scanning, 30.14
- Dynamic range, of solid-state cameras, 26.11, 26.14
- Dynamic scattering, 9.7–9.8, 9.7*f*
- Dyson interference microscopes, 28.41, 28.42, 28.42*f*
- Dyson lenses, 18.21, 18.22, 18.22*f*
- Eagle configuration, of dispersive prisms, 20.7, 20.11*f*
- Ebert-Fastie configuration, 20.14*t*
- Ebert-Fastie configuration, of dispersive prisms, 20.8, 20.12*f*
- Eccentric field arrangement, of lenses, 18.22
- Eccentricity, 1.34, 15.10
- Edge response, OTF and, 4.7
- Effective focal length (EFL):
 of camera lenses, 27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
 of Gaussian lenses, 1.48
 of microlenses, 22.10
- Effective index method, 12.11, 12.12, 21.4
- Effective medium theories (EMTs), 16.4, 16.9
- Effective-medium representation, of volume scattering, 9.8
- Efficiency factors, for scattering by particles, 7.5, 7.6*f*
- E-folding distance, 7.13
- Eigenpolarization, 15.7
- Eikonals, 1.12–1.14
- Einstein coefficients:
 for spontaneous emission, 31.3
 for stimulated absorption, 10.6
- Eisenburg and Pearson two-mirror, three reflection objective, 29.25
- Elastic scattering, 7.3
- Electric fields, 2.3–2.4, 3.2, 3.3
- Electromagnetic dipole model of light, 3.33–3.36, 3.37*f*
- Electromagnetic theory, 12.4
- Electronic holography, 33.9–33.14, 33.11*f*–33.13*f*
- Electronic structure, of atoms, 10.12–10.16, 10.13*f*–10.15*f*
- Electronic-speckle pattern interferometry (ESPI), 33.9–33.13, 33.11*f*–33.13*f*
- Electrons:
 lifetimes of, 10.6
 in multielectron atoms, 10.10–10.11
 in one-electron atoms, 10.7–10.9, 10.8*f*, 10.9*f*
 and optical spectra, 10.12–10.16, 10.13*f*–10.15*f*
- Electro-optic holography (EOH), 33.11–33.13, 33.12*f*, 33.13*f*
- Electro-optic (gradient) scanners, 30.45–30.48, 30.46*f*–30.48*f*
- Electro-optic tensors, 21.10
- Electro-optical modulators, 15.23
- Ellipsometers, 16.10–16.18, 16.10*f*–16.12*f*
 for azimuth measurements, 16.16
 four-detector photopolarimeters, 16.14–16.16, 16.14*f*–16.16*f*
 interferometric arrangements of, 16.18
 normal-incidence rotating-sample, 16.18
 null, 16.11, 16.12
 perpendicular-incidence, 16.17–16.18, 16.18*f*
 photometric, 16.12–16.14, 16.13*f*, 16.14*f*
 return-path, 16.16–16.17, 16.17*f*
 rotating-analyzer, 16.13, 16.13*f*, 16.14
 rotating-detector, 16.14, 16.14*f*
- Ellipsometric angles, 16.3
- Ellipsometry, 16.1–16.21
 about, 16.2–16.3, 16.2*f*
 applications, 16.21
 conventions, 16.3–16.4, 16.3*f*
 defined, 15.7
 generalized, 16.19
 instrumentation for (*see* Ellipsometers)
 Jones-matrix generalized, 16.19
 modeling and inversion, 16.4–16.9, 16.6*f*–16.8*f*, 16.10*f*

- Ellipsometry (*Cont.*):
 Mueller-matrix generalized, 16.19–16.21,
 16.20f, 16.20t, 16.21f
 multiple-angle-of-incidence, 16.3
 and polarimetry, 15.30–15.32, 15.31f, 15.32f
 spectroscopic, 16.3
 transmission, 16.10
 variable-angle spectroscopic, 16.3
- Elliptical polarizers, 14.10, 15.17–15.18
- Elliptical retarders, Mueller matrices for, 14.14
- Ellipticity, of polarization elements, 15.10
- Empty magnification limit, 28.17
- Endoscopic cameras, 25.21, 25.21f
- Enhanced backscatter (EBS), 6.5–6.7, 6.5f
- Entrance pupil (lens), 1.76, 17.8, 18.4–18.6, 29.37
- Entrance pupil distance (ENP) of camera
 lenses, 27.3f–27.5f, 27.7f–27.16f,
 27.18f–27.22f, 27.25
- Entrance window (lens), 17.9
- Epi-illumination, in microscopes, 28.7–28.9,
 28.7f, 28.8f
- Epitaxy, 21.17–21.20
- Equivalent lenses, 17.20
- Equivalent particles, in volume scattering, 9.6
- Erasing, in xerographic systems, 34.10
- Etching, for PICs, 21.18–21.19
- Étendue, 1.22, 1.81, 13.7
- Euler equations, 1.20
- Excitonic effects, in integrated optics, 21.11
- Exit pupil (lens), 1.76, 17.8, 18.6, 29.37
- Exit pupil distance (EXP), of camera lenses,
 27.3f–27.5f, 27.7f–27.16f,
 27.18f–27.22f, 27.25
- Exit window (lens), 17.9
- Exiting beam, of polarimeters, 15.4
- Exposure, in xerographic systems, 34.3
- Extended boundary condition method
 (EBCM), 7.15
- Extended objects, images of, 1.27
- Extended sources, interference by, 2.20
- Extinction cross section, 7.5, 7.8
- Extinction paradox, 7.8
- Extinction ratio:
 about, 12.14–12.15, 12.17
 of polarizers, 12.21–12.24, 12.22f, 12.23f, 14.6
- Extreme wide-angle lenses, 27.6, 27.13f, 27.14f
- Eye loupes, 17.9–17.10
- Eye relief (ER), 18.8–18.10, 18.9f, 18.10f
- Eye space, of afocal lenses, 18.4
- Eye space pupil diameter, 18.6
- Eye tracking, in cameras, 25.14–25.15,
 25.14f, 25.15f
- Eyepieces, in afocal systems, 18.7
- Fabry-Perot etalon (cavity), 2.33, 2.33f
 Fabry-Perot interferometers, 32.4–32.7,
 32.7f, 32.14
 in dynamic wave meters, 32.17
 in gravitational wave interferometers,
 32.21, 32.21f
 as heterodyne interferometers, 32.10
 and multiple beam interference, 2.33–2.36,
 2.34f, 2.35f
 and wire-grid polarizers, 13.31
- False polarization, 15.38
- Faraday cages, 34.8, 34.8f
- Faraday rotators, 28.46
- Faraday shutters, 25.22
- Fast axis, 12.25, 15.7
- Fax machines, 24.6
- Feedthrough, of optical disk data, 35.14
- Fermat's principle, 1.11–1.13, 1.24
- Ferrimagnetism, 35.25, 35.26, 35.26f
- Feussner prisms, 13.6, 13.7, 13.22–13.23, 13.22f
- Fiber interferometers, 32.14–32.16, 32.15f
- Fiber optic gyroscopes (FOG), 21.2,
 21.35–21.37, 21.36f, 21.36t
- Fiber squeezers, 15.24
- Fiber-to-fiber excess loss, 21.13–21.14
- Fick's law, 9.12
- Field (lens), 1.74
- Fields, of rays, 1.13
- Field angles:
 of apertures, 1.75, 1.75f, 1.76
 of Glan-Thompson type prisms, 13.12
- Field curvature, 1.91, 29.7, 29.37
 (*See also* Petzval curvature)
- Field flatness (aberration), 1.91
- Field flattener lenses, 17.28
- Field intensities, of waves, 2.5–2.6
- Field lenses, 1.82, 1.82f, 17.10
- Field of view (FOV), 1.74
 in Keplerian afocal lenses, 18.11
 for reflective and catadioptric objectives,
 29.34–29.35, 29.35f, 29.36f
 in telescopes, 18.15–18.16, 18.15f
- Field size (lens), 28.13
- Field stop, 1.74, 17.9, 29.5, 29.37
- Field-effect-transistors (FETs), 21.38
- Field-flattened Schmidt objective, 29.14–29.15

- Fill factor, of binary optics, **23.8**
- Film (camera), **25.5, 25.6**
- Filters:
- for coherent optical image enhancement, **11.14–11.17**
 - of Mach-Zehnder interference, **21.23, 21.24f, 21.25**
 - narrowband, **3.3**
 - for pattern recognition, **11.12–11.14**
 - spatial, **11.5–11.6, 11.6f**
- Finesse:
- coefficient of, **2.31**
 - of Fabry-Perot etalon, **2.34, 2.35**
 - of interferometers, **32.6**
- Finish models, for surface scattering, **8.14–8.15**
- Finite conjugate afocal relays, **18.11–18.12, 18.12f**
- Finite rays, **1.35**
- Finite-difference time-domain (FDTD) technique, **7.16–7.17**
- First-order optics, **1.29, 1.37**
- First-order retardation plates, **13.46**
- Fish-eye lenses, **27.6**
- Fixed-pattern noise (FPN), **26.11**
- Fizeau interferometers, **2.24–2.26, 2.25f, 32.2, 32.2f, 32.17**
- Flame hydrolysis (FHD), **21.13–21.14**
- Flash (camera), **25.16, 25.17f**
- Flat-field objective optics, **30.30–30.33**
- Flat-medial-field objectives, **29.11**
- Fluorescence imaging with one nanometer accuracy (FIONA), **28.23**
- Fluorescence line narrowing (FLN), **10.18, 31.29, 31.29f**
- Fluorescent microscopy, **28.48–28.49**
- Flux:
- of polarization elements, **15.9**
 - and radiative transfer, **9.12–9.13, 9.13f**
- F-number, **1.79, 17.9**
- Focal Gaussian lenses, **1.45–1.53, 1.47f**
- conjugate equations of, **1.49–1.50, 1.49f**
 - magnifications and distances in, **1.50–1.53**
 - nodal points of, **1.48, 1.48f, 1.49**
 - principal focal points of, **1.47**
 - principal planes of, **1.47, 1.47f, 1.48**
 - reduced coordinates of, **1.53**
- Focal length:
- effective
 - of camera lenses, **27.3f–27.5f, 27.7f–27.16f, 27.18f–27.22f, 27.25**
 - of Gaussian lenses, **1.48**
 - of microlenses, **22.10**
- Focal length (*Cont.*):
- in gradient index optics, **24.5–24.6**
 - primary, **3.12**
 - of surfaces, **1.39–1.40**
 - in systems of revolution, **1.40**
- Focal lines, **1.58**
- Focal planes, **1.57, 1.70**
- Focal plane-to-focal plane conjugate matrices, **1.69**
- Focal points:
- front and rear, **1.40, 1.47**
 - of lens systems, **17.7**
 - principal, **1.47, 1.58**
- Focal ratio, **29.5, 29.37**
- Focal-plane-to-focal-plane geometry, **11.3–11.4, 11.3f**
- Focus, range of, **1.85**
- Focus error signal (FES), **35.12–35.14, 35.13f**
- Focusing, of optical disks, **35.9–35.12, 35.11f, 35.12f**
- Form birefringence (term), **15.41**
- Form dichroism (term), **15.41**
- Förster resonance energy transfer (FRET), **28.23**
- 45° half-wave linear retarders, Mueller matrices for, **14.12t**
- 45° linear polarizers, **14.10t**
- 45° quarter-wave linear retarders, Mueller matrices for, **14.12t**
- Forward-looking infrared (FLIR) systems, **30.22, 30.23, 30.51**
- Foster prisms, **13.7, 13.18f, 13.21–13.22**
- Foucault prisms, **13.7, 13.17**
- Four-detector photopolarimeters (FDPs), **16.14–16.16, 16.14f–16.16f**
- Fourier analysis, for radiative transfer, **9.11**
- Fourier domain filters, **11.13**
- Fourier transform lenses, **18.12**
- Fourier transform plane, **33.16–33.19, 33.18f, 33.20f**
- Fourier transforms:
- in analog optical and image processing, **11.3–11.5, 11.5f**
 - and coherence theory, **5.20**
 - and diffraction, **3.2, 3.24**
 - for focal plane-to-focal plane matrices, **1.69**
 - and transfer functions, **4.2**
 - of uniformly illuminated linear aperture, **30.54**
- Four-powered-mirror lenses, **18.21, 18.21f**
- Fractal model of surface finish, **8.14**

- Fractals:
 brownian, 8.9, 8.17
 Fresnel-Kirchhoff approximation for, 8.8–8.9
- Frankford Arsenal prisms, 19.3*t*, 19.18–19.24, 19.18*f*–19.24*f*
- Frank-Ritter-type prisms, 13.6, 13.6*f*, 13.13–13.14
- Franz-Keldysh effect, 21.11, 21.11*f*
- Fraunhofer approximation, 5.14, 5.16
- Fraunhofer diffraction, 3.24–3.28, 3.24*f*–3.26*f*
 Airy diffraction as, 28.17
 conducting screens for, 3.33*f*
 and gratings, 20.3
- Fraunhofer theory, 7.11
- Fredericksz cells, 14.31, 15.21
- Free spectral range (FSR), of interferometers, 2.34, 2.35, 2.35*f*, 32.5
- Frenet equation, 1.19
- Frequency domain, 7.17
- Frequency-modulation interferometers, 32.9, 32.9*f*, 32.10
- Fresnel amplitude reflection coefficients, 12.7–12.8
- Fresnel amplitude transmission coefficients, 12.8
- Fresnel diffraction, 1.24
- Fresnel equations, 12.6–12.13, 12.15
 for absorbing materials, 12.10–12.13, 12.13*f*
 coordinate system for, 12.6–12.7, 12.7*f*
 for nonabsorbing materials, 12.8–12.10, 12.9*f*
- Fresnel intensity reflectivities, 8.6
- Fresnel lenses, micro-, 22.31–22.37, 22.31*f*–22.33*f*, 22.35*f*–22.37*f*
- Fresnel losses, 22.10
- Fresnel number, 3.25
- Fresnel propagation kernels, 6.6
- Fresnel reflection coefficients, 8.10
- Fresnel rhombs, 13.45, 13.50*f*
- Fresnel Risley prisms, 19.27*f*
- Fresnel zones, 3.4–3.6, 3.4*f*
 for cylindrical wavefronts, 3.6*f*, 3.7*f*, 3.13–3.14, 3.14*f*
 and Fraunhofer diffraction, 3.25
 opaque strip obstruction, 3.20–3.21, 3.21*f*
- Fresnel-Kirchhoff approximation, 8.5–8.9
- Fresnel-Kirchhoff diffraction formula, 3.21, 3.22, 3.32
- Fresnel's biprism, 2.16, 2.16*f*
- Fresnel's mirror, 2.16, 2.16*f*
- Fringe localization, 2.14
- Fringe visibility (contrast), 2.7–2.8
- Fringe-counting interferometers, 32.8, 32.8*f*
- Fringes of equal inclination, 2.20–2.22, 2.21*f*, 2.22*f*
- Fringes of equal thickness, 2.22–2.24, 2.23*f*
- Front focal lengths, 1.40
- Front focal points, 1.40, 1.47
- Front principal planes, 1.48
- Front vertex distance (FVD), camera lens performance and, 27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
- Full width at half maximum (FWHM) points, 30.9
- Full-frame arrays, of CCDs, 26.3–26.4, 26.4*f*
- Fusing, in xerographic systems, 34.10
- Gabor objective, 29.20
- Galilean lenses, 18.7, 18.15–18.17, 18.15*f*, 18.16*f*
- Galilean telescopes, 18.15
- Galvanometer scanners, 30.41–30.44, 30.43*f*, 30.44*f*
- Gauss points, of lenses, 1.44
- Gaussian analyses, of lenses:
 afocal, 18.4–18.7, 18.7*f*
 focusing, 18.2–18.4, 18.3*f*, 18.5*f*
- Gaussian lenses, 1.44–1.55
 afocal, 1.45, 1.46*f*, 1.53–1.54, 1.53*f*, 1.54*f*
 focal, 1.45–1.53, 1.47*f*
 conjugate equations of, 1.49–1.50, 1.49*f*
 magnifications and distances in, 1.50–1.53
 nodal points of, 1.48, 1.48*f*, 1.49
 principal focal points of, 1.47
 principal planes of, 1.47, 1.47*f*, 1.48
 reduced coordinates of, 1.53
 notation for, 1.45*t*
 properties of, 1.54
 systems of, 1.54–1.55
- Gaussian optics, 1.29, 1.44
- GDx Nerve Fiber Analyzer, 15.41
- Gédamine, 13.9, 13.10, 13.11*f*, 13.20
- General deviation prisms, 19.3*t*, 19.28, 19.28*f*–19.29*f*
- General Photonics, 15.24
- Generalized ellipsometry (GE), 16.19
- Generalized Lagrange invariant (étendue), 1.22, 1.81, 13.7
- Geometrical optics:
 for aberrations of point images, 1.86–1.92, 1.86*f*
 and binary optics, 23.5–23.9, 23.5*f*–23.9*f*, 23.10*t*
 characteristic functions of, 1.13–1.18, 1.15*f*

- Geometrical optics (*Cont.*):
 for collineation, 1.56–1.65, 1.65*f*
 and conservation of étendue, 1.22
 defined, 1.8
 for Gaussian lenses, 1.44–1.55
 afocal, 1.45, 1.46*f*, 1.53–1.54, 1.53*f*, 1.54*f*
 focal, 1.45–1.55, 1.47*f*–1.49*f*
 notation for, 1.45*t*
 properties of, 1.54
 systems of, 1.54–1.55
 and images about known rays,
 1.43–1.44, 1.44*f*
 for imaging, 1.26–1.31, 1.30*f*
 at interfaces of homogeneous media, 1.23–1.26
 lens sizes and fields in, 1.74–1.85
 apertures, 1.74–1.77, 1.75*f*
 cosine-to-the-fourth approximation, 1.81
 field lenses, 1.82, 1.82*f*
 fields, 1.74, 1.77, 1.84
 F-number, 1.79
 focus and defocus, 1.82–1.85, 1.83*f*
 irradiance, 1.80
 power per pixel, 1.80
 pupils, 1.76–1.79, 1.76*f*, 1.78*f*
 telecentricity, 1.83–1.84
 total lens étendue, 1.81
 vignetting, 1.81, 1.81*f*, 1.82
 paraxial matrix methods, 1.65–1.74
 and rays, 1.8–1.13
 in heterogeneous media, 1.18–1.22
 paths of, 1.10–1.13
 and skew invariant, 1.23
 of systems of revolution, 1.32–1.43
 paraxial optics of, 1.37–1.43
 ray tracing in, 1.35–1.37, 1.36*f*
 surfaces, 1.32–1.35
 unfolded reflections, 1.32
 Geometrical path length, 1.11
 Geometrical wavefronts, 1.12–1.13
 Gimbal-less two-axis scanning-micromirror
 devices (GSMDs), 30.61–30.62, 30.62*f*
 Glan-Foucault prisms, 13.7, 13.9, 13.11*f*,
 13.12–13.14
 Glan-Taylor prisms, 13.7, 13.9*n*, 13.10–13.14,
 13.10*f*, 13.11*f*
 Glan-Thompson prisms, 13.6, 13.6*f*,
 13.9–13.12, 13.10*f*, 13.18*f*, 13.22
 field angle of, 13.12
 and optical spectrometers, 31.7
 sheet polarizers vs., 13.27
 transmission by, 13.9–13.10, 13.11*f*
 Glan-type prisms, 13.6, 13.6*f*, 13.8–13.15
 defects and testing of, 13.14–13.15
 Frank-Ritter, 13.6, 13.6*f*, 13.13–13.14
 Glan-Foucault, 13.7, 13.9, 13.11*f*,
 13.12–13.14
 Glan-Taylor, 13.7, 13.9*n*, 13.10–13.14,
 13.10*f*, 13.11*f*
 Glan-Thompson, 13.6*f*, 13.9–13.12, 13.10*f*,
 13.11*f*, 13.27
 Lippich, 13.6, 13.6*f*, 13.12–13.13
 Marple-Hess, 13.12, 13.13
 precautions with, 13.14
 Glass:
 for microlenses, 22.9–22.10, 22.10*t*, 22.11*f*,
 22.12*t*, 22.13*t*, 22.14*f*
 for objectives, 29.2, 29.2*t*
 sol-gel formed, 24.8
 Glass-calcite Rochon prisms, 13.20
 Glazebrook prisms, 13.6, 13.9 (*See also*
 Glan-Thompson prisms)
 Goerz prism system, 19.3*t*, 19.17, 19.17*f*
 Goos-Hanchen shift, 21.4
 Graded index profile, of optical waveguides, 21.4
 Graded-index (GRIN) films, 16.9
 Gradient dispersion, 24.3
 Gradient index (GRIN) optics, 24.1–24.8
 analytic solutions of, 24.2
 axial gradient lenses, 24.3–24.5, 24.4*f*
 materials, 24.8
 mathematical representations of, 24.2–24.3
 radial gradients, 24.5–24.8, 24.5*f*, 24.7*f*
 Grating equation, 23.4
 Grating interferometers, 32.4, 32.5*f*
 Grating multiplexers, 23.11, 23.12
 Grating polarizers, 13.33
 Gratings (*see* Diffraction gratings; Dispersive
 prisms and gratings)
 Gravitational-wave interferometers,
 32.21, 32.21*f*
 Gray code, on optical disks, 35.16
 Gray-scale masks, 22.23–22.25,
 22.23*f*–22.25*f*
 Green's function, 3.22–3.23, 3.22*f*
 in approximations of multiple
 scattering, 9.10
 digitized, 7.15, 7.16
 in scattering, 9.3
 Gregorian objective, 29.10
 Grooved regions, on optical disks, 35.15,
 35.15*f*, 35.16
 Gunn diodes, 31.21

- Haidinger fringes, 2.22, 2.22f, 2.27
- Half-period zones, 3.5
- Half-shade devices, 13.56–13.57
- Half-wave linear retarders, Mueller matrices
for, 14.12t, 14.13
- Half-wave plates, 12.25, 12.27
- Half-wave retarders, Mueller matrices
for, 14.14
- Halle prisms, 13.16f, 13.17
- Hamiltonian optics, 1.13, 1.43
- Hamilton's equations for rays, 1.21
- Harting-Dove (Dove) prisms, 19.3t, 19.9,
19.9f, 19.10, 19.10f
- Hartman testing, 22.26
- Hartnack-Prazmowski prisms, 13.16f, 13.17
- Hartree-Fock variational approach, 10.11
- Hastings (Brashear-Hastings) prisms, 19.3t,
19.25, 19.25f
- Hausdorff-Besicovitch dimension, 8.8
- Header information, on optical disks, 35.6
- Heisenberg uncertainty principle, 21.11
- Heisenberg's matrix mechanics, 10.3
- Helmholtz equation, 3.2, 5.9, 5.10, 33.3
- Helmholtz invariant, 1.77 (*See also* Two-ray
paraxial invariant)
- Hermitian matrices, 14.19, 14.41
- Herschel condition, 1.31, 29.34, 29.37
- Herschelian objectives, 29.6, 29.27
- Heterodyne interferometers, 32.10, 32.10f,
32.20, 32.20f
- Heterogeneous media, rays in, 1.9, 1.18–1.22
- High-dry objectives, 28.11, 28.12f
- Higher-order aberrations, 29.37
- Highlight color, in xerographic systems,
34.12, 34.13f
- High-performance miniature systems,
22.5–22.8
- High-speed cameras, 25.21–25.22, 25.22f
- Hilbert space formulation, 5.2
- Hoffman modulation contrast
microscopy, 28.29
- Holographic inspection, 33.16–33.19,
33.17f–33.18f, 33.20f–33.22f
- Holographic memory, 33.24–33.25
- Holographic microscopes, 28.42,
28.43, 28.43f
- Holographic optical elements (HOEs),
33.13–33.16, 33.15f, 35.28–35.29, 35.29f
- Holographic optics, 33.13
- Holographic scanners, 30.38–30.41,
30.40f–30.42f
- Holography, 33.1–33.25
electronic, 33.9–33.14, 33.11f–33.13f
and holographic inspection, 33.16–33.19,
33.17f–33.18f, 33.20f–33.22f
and holographic memory, 33.24–33.25
and interferometry, 33.4–33.9, 33.8f
and lithography, 33.22–33.24, 33.22f–33.24f
optical elements for, 33.13–33.16, 33.15f
principles of, 33.2–33.4, 33.3f
- Holtronic Technologies holographic system, 33.23
- Homogeneous broadening, of lineshapes, 10.7
- Homogeneous coordinates, of collineation
matrix, 1.59
- Homogeneous media, 1.9, 1.23–1.26
- Homogeneous polarization elements:
defined, 15.7
in Mueller matrices, 14.25–14.26
- Homogeneous sources of light, 5.11–5.12, 5.19
- Hopkin's formula, 5.13
- Horizontal half-wave linear retarders, Mueller
matrices for, 14.12t
- Horizontal linear polarizers, 14.9, 14.10t
- Horizontal quarter-wave linear retarders,
Mueller matrices for, 14.12t
- Houghton objectives, 29.22–29.23
- Huang-Rhys parameter, 10.22
- Hubble space telescope, 29.8
- Human visual system (HVS), color in, 26.18
- Hund's rule, 10.12
- Hurst dimension, 8.8
- Huygens wavelets, 5.16
- Huygens-Fresnel construction (diffraction),
3.4–3.13
Babinet principle, 3.9–3.11
for double refraction in calcite, 13.5
Fresnel zones in, 3.4–3.6
and light from circular apertures
and disks, 3.6–3.9
zone plates in, 3.11–3.13, 3.12f
- Hydrogen, Bohr's theory of, 10.3
- Hyperfocal distance, 1.85
- Hysteresis loops, of optical disks, 35.27,
35.27f, 35.28
- Ideal imaging (term), 1.28, 1.38
- Identity matrix, 14.8
- Illumination:
critical, 28.7
epi-, 28.7–28.9, 28.7f, 28.8f
for input/output scanning, 30.14
trans-, 28.5–28.7, 28.6f

- Illumination area, surface scattering and, **8.12**
- Image(s), **1.26**
 about known rays, **1.43–1.44, 1.44f**
 aerial, **1.26**
 from cameras, **25.5, 25.7**
 of extended objects, **1.27**
 in focal Gaussian lenses, **1.51**
 formation of, **6.9–6.10, 6.9f, 17.5–17.8, 17.5f, 17.6f**
 latent, **34.1–34.4, 34.2f–34.3f**
 medial, **29.11, 29.37**
 point, **1.85–1.92, 1.88f**
 of points, **1.27**
 received, **1.26**
 recorded, **1.26**
 in systems of revolution, **1.42**
 two-dimensional, **11.12–11.17, 11.13f**
 virtual, **29.38**
- Image derotation, **30.36**
- Image distance, **18.3**
- Image distribution, for cameras, **25.7**
- Image enhancement, coherent optical, **11.14–11.17**
- Image erectors, **18.10**
- Image height, **1.27, 18.4**
- Image inversion, **19.2**
- Image plane, **1.27**
- Image point, **1.27**
- Image processing, analog (*see* Analog optical signal and image processing)
- Image receptors, in xerographic systems, **34.1**
- Image reversion, **19.2**
- Image rotation, **30.35–30.36**
- Image space, **1.26, 1.83–1.84, 1.83f**
- Image space numerical aperture, **1.79**
- Image-forming cone (bundle), **1.74**
- Image-forming rays, **1.74**
- Imaging:
 conspicious vs. orthoscopic, **28.8–28.9**
 fluorescence, **28.23**
 in geometrical optics, **1.26–1.31**
 ideal, **1.17, 1.28–1.29, 1.38**
 in microscopes, **28.44–28.54**
 aperture-scanning microscopy, **28.53–28.54**
 confocal microscopy, **28.49–28.51, 28.49f, 28.51f**
 fluorescent microscopy, **28.48–28.49**
 light field microscopy, **28.53**
 polarizing microscopes, **28.44–28.48, 28.44f, 28.46f, 28.47f**
 structured illumination, **28.52**
- Imaging (*Cont.*):
 Newtonian equation for, **17.8**
 in polarimeters, **15.6**
 pupil, **1.76**
 stigmatic, **1.29–1.31, 1.30f**
- Image-space scanners, **30.18–30.23**
- Immersion lenses, **17.11**
- Imperfect polarizers, **13.33**
- Impermeability, dielectric, **21.9**
- Inclination factor, **3.5**
- Incoherent arrays, scattering by, **7.2–7.3**
- Incoherent light sources, **5.12, 5.18–5.19**
- Incoherent processing, of discrete signals, **11.17–11.20, 11.18f, 11.19f**
- Incoherent radiation, **30.2, 30.26, 30.27**
- Incoherent scattering, **9.2–9.5**
- Incomplete polarimeters, **15.4**
- Incomplete sample-measuring polarimeters, **15.16–15.17**
- Index ellipsoid equation, **21.9–21.10**
- Index of refraction (refractive index), **1.9**
 for Brewster angle transmission polarizers, **12.21–12.22**
 complex, **7.12–7.13, 12.5, 12.6**
 distributed, **24.1**
 in gradient index optics, **24.2–24.3**
 in integrated optics, **21.8–21.9**
 of polarizers, **12.16, 12.18**
 for rays in heterogeneous media, **1.21–1.22**
 of shallow radical gradients, **24.7–24.8**
- Infrared emitting diodes (IREDs), **25.14, 25.14f**
- Infrared radiation, single-order plates and, **13.47**
- Inhomogeneous (heterogeneous) media, rays in, **1.9, 1.18–1.22**
- Inhomogeneous optics, **24.1** (*See also* Gradient index (GRIN) optics)
- Inhomogeneous polarization elements, **14.25–14.26, 14.26f, 15.7, 15.20**
- Input planes, translations of, **1.68**
- Input/output scanning, **30.2, 30.4–30.6, 30.25–30.34**
 objective, preobjective, and postobjective, **30.28–30.29, 30.29f, 30.30f**
 objective optics, **30.30–30.33, 30.32f–30.33f**
 power density and power transfer of, **30.25–30.28, 30.27f, 30.28f**
 resolution of, **30.8–30.14, 30.10f, 30.10t, 30.11t, 30.12f–30.13f**
- Inspection, holographic, **33.16–33.19, 33.17f–33.18f, 33.20f–33.22f**
- Instant pictures, **25.8**

- Instrumental polarization, 12.15
- Insulating magnetic brush (IMB), 34.6
- Integrated circuits:
- photonic
 - of III-V materials, 21.17–21.20
 - in integrated optics, 21.2
 - in WDM systems, 21.37, 21.38
 - polarizers for, 13.57
- Integrated optic circuits (IOCs), 21.2
- Integrated optics (IO), 21.1–21.41
- applications of, 21.31–21.39
 - analog transmission, 21.32–21.34, 21.33*f*, 21.34*f*
 - digital transmission, 21.31–21.32
 - fiber optic gyroscopes, 21.35–21.37, 21.36*f*, 21.36*t*
 - silicon photonics transmission, 21.38, 21.39
 - switching, 21.34–21.35, 21.34*f*
 - WDM systems, 21.37–21.38, 21.37*f*–21.39*f*
 - circuit elements of, 21.21–21.31
 - active devices, 21.25–21.31, 21.26*f*–21.31*f*
 - passive devices, 21.21–21.25, 21.22*f*–21.25*f*
 - future trends in, 21.39–21.41
 - advanced integration, 21.40–21.41
 - shift from R&D to manufacturing, 21.39–21.40
 - materials and fabrication techniques for, 21.13–21.21
 - ion-exchanged glass waveguides, 21.13
 - LiNbO₃ and LiTaO₃, 21.16–21.17, 21.28, 21.33
 - PICs of III-V materials, 21.17–21.20
 - silicon photonics, 21.14–21.16, 21.15*f*
 - thin film oxides, 21.13–21.14
 - physics of, 21.3–21.12
 - carrier effects, 21.10–21.12
 - index of refraction, 21.8–21.9
 - linear electro-optical effect, 21.9–21.10
 - nonlinear effects, 21.12
 - optical waveguides, 21.3–21.8, 21.3*f*–21.5*f*, 21.7*f*
 - thermal effects, 21.12
- Intensified CCD (ICCD), 26.3
- Intensity, 5.7, 11.2
- Intensity interferometers, 32.19
- Intensity reflection coefficients, 12.8–12.12
- Intensity transmission coefficients, 12.9–12.10
- Interference, 2.3–2.42
- by amplitude division, 2.19–2.28
 - extended source, 2.20
 - Fizeau interferometers, 2.24–2.26, 2.25*f*
 - fringes of equal inclination, 2.20–2.22, 2.21*f*, 2.22*f*
 - fringes of equal thickness, 2.22–2.24, 2.23*f*
 - Michelson interferometer, 2.26–2.28, 2.26*f*–2.27*f*
 - plane-parallel plate, 2.19, 2.20*f*
 - thin films, 2.24
 - applications of, 2.42
 - and coherence, 2.36–2.42, 2.36*f*, 2.38*f*–2.40*f*, 2.42*f*
 - constructive and destructive, 2.7
 - effects of, 2.5–2.14
 - aberrated wavefronts, 2.12, 2.13
 - coherence, 2.13
 - interference fringes, 2.6–2.8, 2.6*f*, 2.7*t*
 - plane wave and spherical wave, 2.9–2.11, 2.10*f*
 - temporal beats, 2.13
 - two plane waves, 2.8–2.9, 2.9*f*
 - two spherical waves, 2.11–2.12, 2.12*f*, 2.13*f*
 - multiple beam, 2.28–2.36
 - diffraction gratings, 2.28–2.29, 2.29*f*, 2.30*f*
 - Fabry-Perot interferometer, 2.33–2.36, 2.34*f*, 2.35*f*
 - plane-parallel plates, 2.30–2.33, 2.30*f*, 2.32*f*, 2.33*f*
 - order of, 2.7
 - by wavefront division, 2.14–2.19, 2.15*f*–2.18*f*
 - and wavefronts, 2.4–2.5, 2.5*f*
 - and waves, 2.3–2.4, 2.5*f*
- Interference filters, Mach-Zehnder, 21.23, 21.24*f*, 21.25
- Interference fringes, 2.6–2.8, 2.6*f*, 2.7*t*
- Interference gratings, 33.14
- Interference microscopy, 28.33–28.44
- differential-interference contrast
 - microscopes, 28.39–28.41, 28.40*f*
 - Dyson microscopes, 28.41, 28.42, 28.42*f*
 - holographic, 28.42, 28.43, 28.43*f*
 - Jamin-Lebedev microscopes, 28.38–28.39, 28.38*f*
 - Linnik microscopes, 28.36–28.38, 28.37*f*
 - Mach-Zehnder microscopes, 28.36, 28.37*f*
 - Mirau microscopes, 28.41, 28.42, 28.42*f*
 - and optical coherence tomography, 28.43–28.44
 - optical path difference (OPD) in, 28.33–28.34, 28.35*f*, 28.36

- Interference polarizers, 13.39–13.41, 13.40f
- Interferometers (interferometry), 32.1–32.21
- double-passed two-beam, 32.8
 - Fabry-Perot, 32.4–32.7, 32.7f, 32.14
 - in dynamic wave meters, 32.17
 - in gravitational wave interferometers, 32.21, 32.21f
 - as heterodyne interferometers, 32.10
 - and multiple beam interference, 2.33–2.36, 2.34f, 2.35f
 - and wire-grid polarizers, 13.31
- fiber, 32.14–32.16, 32.15f
- finesse of, 32.6
- Fizeau, 2.24–2.26, 2.25f, 32.2, 32.2f, 32.17
- free spectral range of, 2.34, 2.35, 2.35f, 32.5
- frequency-modulation, 32.9, 32.9f, 32.10
- fringe-counting, 32.8, 32.8f
- grating, 32.4, 32.5f
- gravitational-wave, 32.21, 32.21f
- heterodyne, 32.10, 32.10f, 32.20, 32.20f
- and holography, 33.4–33.9, 33.8f
- intensity, 32.19
- and interferometric optical switches, 32.19
 - and interferometric wave meters, 32.16–32.17, 32.16f, 32.17f
- laser-Doppler, 32.12–32.13, 32.13f
- laser-feedback, 32.13–32.14, 32.14f
- Mach-Zehnder, 21.10, 21.12, 21.14, 21.16, 21.32, 21.40, 32.3, 32.3f, 33.5
- Michelson, 2.26–2.28, 2.26f–2.27f, 32.2, 32.3f, 32.21, 33.4
- Michelson stellar, 2.40–2.41, 2.40f, 32.19, 32.19f
- Newton, 2.25
- Nomarski, 32.4, 32.5f
- nulling, 32.20–32.21
- phase-conjugate, 32.17, 32.18, 32.18f
- phase-locked, 32.11–32.12, 32.12f
- phase-shifting, 32.10–32.11, 32.11f
- polarization, 32.4, 32.5f
- Sagnac, 21.35, 21.36, 21.36f, 32.3–32.4, 32.4f
- second-harmonic, 32.17–32.18, 32.18f
- shearing, 32.4, 32.6f
- stellar, 32.19–32.21, 32.19f, 32.20f
- three-beam, 32.7–32.8, 32.7f
- and two-wavelength interferometry, 32.9
- Twyman-Green, 2.28, 32.2, 32.9, 33.5
- Young's two pinhole, 6.3
- Interferometric arrays, 32.20–32.21
- Interferometric ellipsometry, 16.18
- Interferometric optical switches, 32.19
- Interferometric wave meters, 32.16–32.17, 32.16f, 32.17f
- Interline transfer, of CCDs, 26.4–26.5, 26.5f, 26.6f
- Intersection points, ray tracing for, 1.36, 1.36f
- Invariance properties, of rays, 1.10
- Invariants, optical, 18.7
- Inverse filters, for coherent optical image enhancement, 11.14–11.17
- Inverse Galilean telescopes, 18.15, 18.16, 18.16f
- Inverse systems, conjugate matrices for, 1.71
- Inverted telephoto lenses, 27.2, 27.6, 27.7f–27.14f
- highly complex extreme speed, 27.6, 27.9f–27.12f
 - with nonrectilinear distortion, 27.6
 - with rectilinear distortion correction, 27.6, 27.13f, 27.14f
 - very compact moderate speed, 27.6, 27.7f–27.8f
- Ion exchange process, 24.8
- Ion-exchanged glass waveguides, 21.13
- Ions, tri-positive rare earth, 10.16–10.18, 10.16t, 10.17f
- Iris (lens), 17.8
- Irradiance, 3.3
- of circular patterns and disks, 3.7–3.8
 - of complementary aperture screen patterns, 3.10–3.11
 - and diffraction gratings, 3.28, 3.29
 - of double-slit patterns, 3.27, 3.27f
 - of lambertian objects, 1.80
 - of straight-edge patterns, 3.15–3.17, 3.18f
 - as vector, 3.33, 3.37f
 - and zone plates, 3.11, 3.12f
- Irradiation, microbeam, 28.54
- Isotropic homogenous spheres, scattering by, 7.11–7.14
- Jamin-Lebedev microscopes, 28.38–28.39, 28.38f
- Jellett-Cornu prisms, 13.56
- J-K model, of surface finish, 8.13, 8.15
- Jones calculus, 12.29–12.30
- Jones matrix, 7.10
- in ellipsometry, 15.30, 16.19
 - and Mueller matrices, 14.3, 14.22–14.24, 14.27, 14.33
 - tensor product of, 14.23

- Keplerian afocal lenses, **18.7–18.14**
 in binoculars, **18.13–18.14, 18.14f**
 and eye relief manipulation, **18.8–18.10, 18.9f, 18.10f**
 field-of-view limitations in, **18.11**
 finite conjugate afocal relays, **18.11–18.12, 18.12f**
 for scanners, **18.13, 18.13f**
 in terrestrial telescopes, **18.10–18.11, 18.11f**
 thin-lens model, **18.7–18.8, 18.8f**
- Keplerian telescopes, **18.10**
- Kerr cell shutters, **25.22**
- Kerr cells, **31.9**
- Kerr effect, **15.23, 35.21, 35.23, 35.25**
- Kinematic lens mounts, **22.7, 22.8f**
- Klein bottle, **14.15**
- Kodak Cameo Motordrive camera, **25.16, 25.16f**
- Kodak Cobra flash, **25.16, 25.17f**
- Kodak DCS 200 camera, **25.7**
- Kodak Fun Saver Portait 35, **25.16, 25.17f**
- Kodak glass molding process, **22.9, 22.10**
- Koehler (Köhler) illumination:
 and coherence theory, **6.12–6.13, 6.12f**
 in microscopes, **28.4–28.5, 28.5f, 28.7, 28.9**
- Korsch objectives, **29.30–29.32, 29.34**
- Kramers states, **31.19, 31.19f, 31.20**
- Kramers-Kronig relation, **21.11**
- Kronecker delta function, **5.20**
- Kubelka-Munk theory, **9.13**
- Lagrange invariant, **1.22, 1.41, 1.77, 1.81, 13.7, 17.5, 30.11**
- Lamb shift, **10.4**
- Lambertian objects, irradiance for, **1.80**
- Lambert's law, **5.12**
- Lamipol structures, **13.57**
- Landé interval formula, **10.17**
- Landolt fringe, **13.14, 13.17, 13.18**
- Lands, of optical disks, **35.3, 35.5f**
- Landscape lenses, **17.17–17.20, 17.17f–17.20f, 17.19t**
- Laporte selection rule, **10.10**
- Large-format cameras, **25.18–25.20**
- Large-format film, **25.6**
- Laser direct write (LDW) fabrication, **22.19–22.23, 22.20f–22.22f**
- Laser modes, coherence theory and, **5.23**
- Laser noise, **35.12, 35.24**
- Laser power modulation (LPM), **35.17–35.19**
- Laser radar (LIDAR) systems, **30.51**
- Laser speckle, **9.14**
- Laser Stark spectroscopy, **31.27–31.29, 31.27f, 31.28f**
- Laser-assisted chemical etching (LACE), **22.45**
- Laser-Doppler interferometers, **32.12–32.13, 32.13f**
- Laser-feedback interferometers, **32.13–32.14, 32.14f**
- Lasers:
 continuous wave (CW) dye, **10.8**
 distributed Bragg reflector, **21.25, 21.30, 21.32, 21.37, 21.37f**
 distributed feedback, **21.25, 21.29, 21.30, 21.32, 21.38, 21.40**
 for interferometry, **2.41–2.42**
- Latent images, in xerographic systems, **34.1–34.4, 34.2f–34.3f**
- Lateral color, **1.91–1.92, 29.14, 29.37**
- Lateral magnification, **17.5**
- Lateral resolution, of microscopes, **28.17–28.19**
- Lau effect, **6.7–6.8, 6.8f**
- Laurent half shades, **13.56**
- Leaky waveguides, **21.3**
- Left circular polarizers, Mueller matrices for, **14.10t**
- Left half-wave circular retarders, Mueller matrices for, **14.12t**
- Left-circularly polarized light, **12.27**
- Left-handed coordinate systems, **12.6**
- Legendre transformations, **1.13, 1.15, 1.16**
- Leica Summitar lenses, **17.28, 17.28f**
- Leman prisms, **19.3t, 19.13, 19.13f**
- Length-to-aperture (L/A) ratio, for prisms, **13.7**
- Lens axis, **1.32** (*See also* Optical axis)
- Lens law, **17.8**
- Lenses, **17.3–17.39**
 axial separations in, **1.52**
 cardinal points of, **1.44, 17.7**
 conjugate matrices for, **1.71**
 defined, **1.9**
 effective focal length of, **1.48, 22.10, 27.3f–27.5f, 27.7f–27.16f, 27.18f–27.22f, 27.25**
 entrance pupil of, **1.76, 17.8, 18.4–18.6, 29.37**
 entrance window of, **17.9**
 equivalent, **17.20**
 exit pupil of, **1.76, 17.8, 18.6, 29.37**
 exit window of, **17.9**
 field size of, **28.13**
 fields of, **1.74**

Lenses (*Cont.*):

- F-number and numerical aperture of, 17.9
- Gauss points of, 1.44
- Gaussian analyses of, 18.2–18.7, 18.3*f*, 18.5*f*, 18.7*f*
- geometrical optics for, 1.74–1.85
 - apertures, 1.74–1.77, 1.75*f*
 - cosine-to-the-fourth approximation, 1.81
 - field lenses, 1.82, 1.82*f*
 - fields, 1.74, 1.77, 1.84
 - F-number, 1.79
 - focus and defocus, 1.82–1.85, 1.83*f*
 - irradiance, 1.80
 - power per pixel, 1.80
 - pupils, 1.76–1.79, 1.76*f*, 1.78*f*
 - telecentricity, 1.83–1.84
 - total lens étendue, 1.81
 - vignetting, 1.81, 1.81*f*, 1.82
- image formation in, 17.5–17.8, 17.5*f*, 17.6*f*
- inverses of, 1.71
- for magnifiers, 17.9–17.10
- for microscopes, 28.9–28.17 (*See also* Reflective and catadioptric lenses)
 - compound, 17.10
 - objective lenses, 28.9–28.15, 28.10*t*, 28.11*t*, 28.12*f*, 28.13*t*, 28.14*f*–28.16*f*
 - oculars, 28.16–28.17
- natural stop shift of, 22.3
- nodal points of, 1.48, 1.49
- optical center point of, 17.16, 17.17
- performance of, 17.29–17.36, 17.30*f*–17.35*f*
- for periscopes, 18.19, 18.19*f*
- pupils of, 17.8–17.9
- rays in, 1.35
- for scanners, 18.13, 30.57–30.60, 30.58*f*–30.60*f*
- shape factor of, 17.12–17.13, 17.13*f*
- single element, 17.12–17.17, 17.13*f*–17.16*f*, 17.17*t*
- stops of, 17.8–17.9
- systems of, 17.7, 17.20–17.22, 17.21*f*–17.22*f*
- See also specific types of lenses*
- Lenslets, monolithic, 22.25–22.26, 22.25*f*
- Lifetimes, of electrons, 10.6
- Light-emitting diodes (LEDs):
 - in integrated optics, 21.2
 - and parallel matrix-vector multipliers, 11.18
 - and serial incoherent matrix-vector multipliers, 11.17–11.18
 - transmission by, 21.32
- Light field microscopy, 28.53
- Light grasp (étendue), 1.22, 1.81, 13.7
- Light sources, 5.6, 5.9–5.13
- Light-gathering power (étendue), 1.22, 1.81, 13.7
- Light-measuring polarimeters, 15.3–15.4, 15.11–15.13
- Linear diattenuation and diattenuators, 14.8, 14.17
- Linear electro-optic effect, 21.9–21.10
- Linear magnification, 18.4
- Linear polarization sensitivity, 14.17
- Linear polarizers, 14.9, 14.10*t*, 15.7
- Linear systems, coherence theory for, 6.3–6.4
- Linearity:
 - in paraxial matrix methods, 1.66
 - of systems of revolution, 1.41
 - and transfer functions, 4.2
- Linnik interference microscopes, 28.36–28.38, 28.37*f*
- Liouville's theorem, 1.22
- Lippich-type prisms, 13.6, 13.12–13.13
 - Glan-Taylor, 13.7, 13.9*n*, 13.10, 13.10*f*, 13.11*f*, 13.12–13.14
 - half-shade, 13.12*n*, 13.56
 - Marple-Hess, 13.12, 13.13
- Liquid crystal (LC) cells, 15.32*f*, 15.33–15.35, 15.33*f*, 15.34*t*
- Liquid crystal displays (LCDs), 25.7
- Liquid crystal (LC) lenses, 22.40–22.41, 22.42*t*
- Liquid crystal on silicon (LCOS) panels, 15.28
- Liquid crystal retarders, 15.21–15.23, 15.22*f*
- Liquid crystal variable retarders (LCVRs), 15.21–15.23, 15.22*f*
- Liquid immersion development, in xerographic systems, 34.9, 34.10
- Liquid lenses, 22.37–22.41, 22.38*f*–22.41*f*, 22.42*t*
- Liquid-phase epitaxial (LPE) growth, 21.17
- Lissajous-type figures, from Risley prisms, 19.25, 19.26*f*
- Lithium niobate (LiNbO₃), 21.16–21.17, 21.28, 21.33, 21.39
- Lithium tantalate (LiTaO₃), 21.17
- Lithography:
 - deep x-ray, 22.7
 - and holography, 33.22–33.24, 33.22*f*–33.24*f*
 - and miniature and micro-optics, 22.18–22.25
 - with gray-scale masks, 22.23–22.25, 22.23*f*–22.25*f*
 - laser direct write fabrication, 22.19–22.23, 22.20*f*–22.22*f*

- Littrow configuration, of dispersive prisms, 20.7f, 20.10
- Littrow mirrors, 20.4
- Lloyd's mirror, 2.16, 2.17f, 2.18
- Localization, in volume scattering, 9.13–9.17, 9.14f
- Longitudinal aberrations, 1.87
- Longitudinal magnification, 1.28, 1.52, 17.5
- Lorentzian distribution of frequencies, 10.7
- Loupes, eye, 17.9–17.10
- Low inertia scanners, 30.43
- Low-light-level television (LLTV) systems, 26.3
- Low-order flux models, of radiative transfer, 9.12–9.13, 9.13f
- Lu-Chipman polar decomposition, 14.39–14.40
- Lukosz-type super-resolving systems, 6.9–6.10, 6.9f
- Luminescence, 31.17, 31.19–31.21, 31.19f
- Luminescence excitation spectrometers, 31.11–31.12, 31.11f
- Luminescence spectrometers, 31.5–31.12, 31.8f, 31.11f
- Luminosity (étendue), 1.22, 1.81, 13.7
- Luneburg lenses, 1.21, 22.26 (*See also* Distributed-index planar microlenses)
- Lyot coronagraphs, 29.5
- Lyot stop, 29.5, 29.5f, 29.37
- Mach-Zehnder (MZ) interference filters, 21.23, 21.24f, 21.25
- Mach-Zehnder (MZ) interference microscopes, 28.36, 28.37f, 28.39
- Mach-Zehnder (MZ) interferometers, 21.10, 21.12, 21.14, 21.16, 21.32, 21.40, 32.3, 32.3f, 33.5
- Mach-Zehnder (MZ) modulators, 21.26, 21.27f, 21.28, 21.32, 21.34
- Magnetic brush development, in xerographic systems, 34.5–34.7, 34.5f–34.7f
- Magnetic circular dichroism (MCD), 31.20, 31.21
- Magnetic circular polarization (MCP), 31.21
- Magnetic field modulation (MFM), 35.19, 35.19f, 35.20
- Magnetic resonance, optically detected, 31.21–31.23, 31.22f, 31.23f
- Magneto-optical (MO) disks, 35.2
- Magneto-optical (MO) modulators, 15.23
- Magneto-optical (MO) readout, 35.21–35.24, 35.22f, 35.24f
- Magneto-optical (MO) recording, 35.25–35.28, 35.26f, 35.27f (*See also* Optical disk data storage)
- Magnification, 1.28
- afocal, 18.5–18.6
- angular, 1.52, 1.78, 18.4
- axial (longitudinal), 1.28, 1.52, 17.5
- dia-, 1.23
- by focal Gaussian lenses, 1.52–1.53
- lateral, 17.5
- linear, 18.4
- longitudinal, 1.28
- and magnifiers, 17.9–17.10
- in microscopes, 28.3–28.4
- pupil, 1.76
- pupil angular, 1.78
- scan, 30.5, 30.12–30.14
- secondary, 29.12, 29.38
- in systems of revolution, 1.42
- transverse, 1.28, 1.50–1.51
- visual, 1.28
- Maksutov objectives, 29.19, 29.20
- Malus-Dupin principle, 1.12
- Mangin elements, in objective design, 29.6
- Mangin objectives, 29.7, 29.24
- Mapping, of object and image space, 1.27
- Marginal rays, 1.75, 17.8, 29.37
- Marple-Hess prisms, 13.12, 13.13
- Mask layout, for binary optics, 23.14–23.16, 23.15f, 23.15t
- Maskless lithography tool (MLT), 22.23
- Mass-transport process, for miniature and micro-optics, 22.45, 22.46f
- Master groups, of zoom lenses, 27.17
- Master-oscillator/power-amplifiers (MOPAs), 21.30, 21.30f
- Matched filters, for pattern recognition, 11.12–11.13
- Matrices:
- amplitude scattering, 7.10
- for collineation, 1.59–1.60
- computing polarization with, 12.27–12.30
- density (coherency), 12.29–12.30, 14.41
- identity, 14.8
- Jones, 7.10
- in ellipsometry, 15.30, 16.19
- and Mueller matrices, 14.3, 14.22–14.24, 14.27, 14.33
- tensor product of, 14.23
- paraxial, and geometrical optics, 1.65–1.74

- Matrices (*Cont.*):
 Pauli spin, 14.24, 14.41
 point spread, 15.36, 15.36f
 power, 1.67
 for radiative transfer, 9.11
 for single scattering, 9.16, 9.17
 Stokes, 14.4
 T-matrix method, 7.15
See also Mueller matrices
- Matrix generators, 14.27
- Maxwell equations, 7.3
- Maxwell fisheye (lens), 1.21
- Maxwell-Boltzmann velocity distribution, 31.24
- Maxwell-Garnett mixing formula, 9.8
- Maxwellian ideal imaging, 1.17, 1.28, 1.38
- Maxwell's electromagnetic theory, 10.3
- Maxwell's equations:
 and binary gratings, 23.13
 and coherence theory, 5.2, 5.3
 and diffraction, 3.1, 3.2, 3.4
 and laws of reflection and refraction, 1.24
 and optical waveguides, 21.3, 21.5
 and surface scattering, 8.4
- McCarthy objective, 29.32
- McKinley relays, 18.18, 18.18f, 18.19
- Mean-field approximation, 7.16
- Mechanical distances, in focal Gaussian lenses, 1.53
- Medial images, 29.11, 29.37
- Medium-format film, 25.6
- Meinel-Shack objective, 29.28
- Mellin transforms, 11.14
- Melted-resin arrays, 22.42–22.45, 22.42f–22.44f
- Memory, holographic, 33.24–33.25
- Meridians (meridional planes), 1.27, 1.32
- Meridional rays, 1.35, 1.37
- Mersenne objectives, 29.9, 29.12
- Mersenne telescopes, 18.20
- Metal-insulator semiconductor (MIS) capacitors, 26.2
- Metal-organic molecular beam epitaxy (MOMBE), 21.17, 21.18
- Metal-organic vapor-phase epitaxy (MOVPE), 21.17–21.18
- Metal-oxide semiconductor (MOS) capacitors, 26.2, 26.6
- Mica retardation plates, 13.45–13.46
- Michel Lévy Color Chart, 28.35f
- Michelson interferometers, 2.26–2.28, 2.26f–2.27f, 32.2, 32.3f, 32.21, 33.4
- Michelson stellar interferometers, 2.40–2.41, 2.40f, 32.19, 32.19f
- Microbeam irradiation, 28.54
- Microelectromechanical systems (MEMS), 30.3
- Micro-Fresnel lenses (MFLs), 22.31–22.37, 22.31f–22.33f, 22.35f–22.37f
- Microlenses:
 distributed-index planar, 22.26–22.31, 22.27f–22.30f, 22.27t, 22.31t
 micro-Fresnel lenses (MFLs), 22.31–22.37, 22.31f–22.33f, 22.35f–22.37f
 molded glass, 22.9–22.10, 22.10t, 22.11f, 22.12t, 22.13t, 22.14f
 molded plastic, 22.10, 22.12–22.15
- Micromachining techniques, for binary optics, 23.16, 23.16f, 23.17
- Micromirrors, 22.23, 30.61–30.62, 30.62f
- Micro-optical table (MOT) techniques, 22.6, 22.7f
- Micro-optics (*see* Miniature and micro-optics)
- Microscopes, 28.1–28.56
 aperture-scanning microscopy, 28.53–28.54
 bright field microscopy, 28.25, 28.27–28.28, 28.27f
 compound, 17.10
 confocal microscopy, 28.49–28.51, 28.49f, 28.51f
 contrast in
 bright field microscopy, 28.25, 28.27–28.28, 28.27f
 dark field microscopy, 28.28
 Hoffman modulation contrast, 28.29
 interference microscopy, 28.33–28.44, 28.35f, 28.37f, 28.38f, 28.40f, 28.42f, 28.43f
 modulation transfer function, 28.24–28.25, 28.25f, 28.26f
 phase contrast, 28.28–28.29, 28.29f
 SSEE microscopy, 28.29–28.33, 28.30f–28.33f
 dark field microscopy, 28.28
 differential-interference contrast, 28.39–28.41, 28.40f
 Dyson, 28.41, 28.42, 28.42f
 fluorescent microscopy, 28.48–28.49
 history of, 28.1–28.3
 holographic, 28.42, 28.43, 28.43f
 imaging modes of, 28.44–28.54, 28.44f, 28.46f, 28.47f, 28.49f, 28.51f
 interference microscopy, 28.35f, 28.37f, 28.38f, 28.42f, 28.43f

- Microscopes (*Cont.*):
 Jamin-Lebedev, 28.38–28.39, 28.38f
 lenses in, 28.9–28.17, 28.10t, 28.11t,
 28.12f–28.16f
 light field microscopy, 28.53
 Linnik, 28.36–28.38, 28.37f
 Mach-Zehnder, 28.36, 28.37f
 Mirau, 28.41, 28.42, 28.42f
 optical arrangements in, 28.3–28.9,
 28.4f–28.8f
 optical path difference (OPD) in,
 28.33–28.34, 28.35f, 28.36
 resolution, 28.17–28.24, 28.18f–28.21f
 specimen manipulation for, 28.54–28.55
 SSEE microscopy, 28.29–28.33, 28.30f–28.33f
- Mie scattering, 7.11, 7.12, 9.17
- Miller algorithm, 7.15
- Miniature and micro-optics, 22.1–22.46, 22.46f
 and binary optics, 23.7–23.8, 23.7f–23.8f
 design considerations, 22.2–22.8
 diamond turning, 22.15–22.18, 22.16t,
 22.17f, 22.18f
 distributed-index planar microlenses,
 22.26–22.31, 22.27f–22.30f, 22.27t, 22.31t
 drawn preform cylindrical lenses, 22.45, 22.46f
 high-performance miniature systems,
 22.5–22.8
 laser-assisted chemical etching (LACE), 22.45
 liquid lenses, 22.37–22.41, 22.38f–22.41f,
 22.42t
 and lithography, 22.18–22.25, 22.20f–22.25f
 mass-transport process, 22.45, 22.46f
 melted-resin arrays, 22.42–22.45,
 22.42f–22.44f
 micro-Fresnel lenses, 22.31–22.37,
 22.31f–22.33f, 22.35f–22.37f
 molded microlenses, 22.8–22.15
 molded glass, 22.9–22.10, 22.10t, 22.11f,
 22.12t, 22.13t, 22.14f
 molded plastic, 22.10, 22.12–22.15
 monolithic lenslet modules,
 22.25–22.26, 22.25f
- Minimum signal, for solid-state cameras,
 26.13–26.14
- Minox camera, 25.21
- Mirages, 24.1, 24.2f
- Mirau interference microscopes, 28.41,
 28.42, 28.42f
- Mircolens arrays, for agile beam steering,
 30.57–30.60
- Mirror-image effect, 12.6
- Mirrors:
 compound, 30.15–30.16, 30.15f
 conic, 29.3f, 29.4f
 Fresnel's, 2.16, 2.16f
 Littrow, 20.4
 Lloyd's, 2.16, 2.17f, 2.18
 micromirrors, 22.23, 30.61–30.62, 30.62f
 plane, 1.25
 in reflecting afocal lenses, 18.20, 18.20f,
 18.21, 18.21f
 reflection from, 1.25
 for scanners, 30.14–30.16, 30.15f,
 30.60–30.62, 30.62f
 as thin lenses, 1.55
See also Reflective and catadioptric objectives
- Misfocus, 1.82
- Mixed characteristic functions, 1.13
- Möbius strip, 14.15
- Modulation transfer function (MTF), 4.3
 calculations, 4.3–4.6, 4.4f, 4.5f
 and camera lens performance, 27.3f–27.5f,
 27.7f–27.16f, 27.18f–27.22f, 27.24, 27.25
 and characteristics of objective detectors, 28.16
 and contrast of microscopes, 28.24–28.25,
 28.25f, 28.26f
 and development in xerographic
 systems, 34.7
 diffraction-limited, 4.4–4.5, 4.4f, 4.5f
 measurements of, 4.6–4.8, 4.8f
 for microscopes, 28.24, 28.25f
 for scanners, 4.6
 for solid-state cameras, 26.14
 at specific wavelengths, 17.38, 17.39
 of uniformly illuminated apertures, 30.9,
 30.10, 30.10f
- Modulators:
 electro-optical modulators, 15.23
 Mach-Zehnder, 21.26–21.28, 21.27f,
 21.32, 21.34
 magneto-optical modulators, 15.23
 photo-elastic, 15.21, 16.13
 polarization (retardance), 15.20–15.24
 traveling wave, 21.26
- Molecular beam epitaxy (MBE), 21.17, 21.18
- Molecular Expressions (website), 28.3
- Molecular scattering, 7.11
- Monocentric Schmidt-Cassegrain
 objectives, 29.16
- Monocentric systems, 29.37
- Monochromatic aberration correction,
 23.6–23.7

- Monochromatic sources of light, 5.11
- Monochromators:
- Czerny-Turner, 31.6
 - double-pass, 20.9*f*
 - Perkin-Elmer Model 99, 20.9*f*
 - Uvicam double-, 20.10, 20.13, 20.15*f*
- Monogon scanners, 30.34–30.36
- Monolithic lenslet modules (MLMs), 22.25–22.26, 22.25*f*
- Moving electron beam exposure system (MEBES), 23.15, 23.16, 23.16*f*
- Mueller calculus, 12.28–12.30
- Mueller matrices, 7.10, 14.1–14.42
- about, 14.3–14.4
 - coordinate system, 14.19–14.20
 - and depolarization, 14.30–14.31
 - depolarization index, 14.32
 - generators of, 14.33–14.39, 14.36*f*, 14.37*f*
 - nondepolarizing matrices, 14.24–14.25, 14.27–14.30
 - diattenuators/diattenuation, 14.16–14.19, 14.16*f*
 - in ellipsometry, 15.30, 16.19–16.21, 16.20*f*, 16.20*t*, 16.21*f*
 - and Jones matrices, 14.22–14.24
 - normalization of, 14.19
 - physically realizable, 14.40–14.42
 - polar decomposition of matrices, 14.39–14.40
 - and polarimetry, 15.8–15.9, 15.11
 - elements of, 15.13–15.14
 - in error analysis, 15.28
 - singular value decomposition, 15.25–15.27
 - polarizance, 14.18
 - and polarization, 14.7, 14.8, 14.25–14.27, 14.33
 - average degree of polarization, 14.32–14.33
 - degree of polarization surfaces and maps, 14.31–14.32, 14.32*f*
 - ideal polarizers, 14.8–14.10, 14.10*t*
 - nonpolarizing, 14.8
 - for radiative transfer, 9.11
 - for refraction and reflection, 14.20–14.22
 - retarder, 14.11, 14.12*t*, 14.13–14.15, 14.15*f*
 - for single scattering, 9.16, 9.17
 - and Stokes parameters, 14.4–14.6
 - transmittance, 14.16–14.17
- Mueller matrix bidirectional reflectance distribution function (MMBRDF), 15.39, 15.39*f*, 15.40*f*
- Mueller matrix polarimeters, 15.26–15.27
- Mueller polarimeters, 15.4
- Mueller vectors, 15.15
- Mueller-Jones matrices, 14.24, 14.27–14.29
- (See also Nondepolarizing Mueller matrices)
- Muller convention, 12.6
- Multielectron atoms, 10.10–10.11
- Multifocal lenses, 23.12, 23.12*f*, 23.13*f*
- Multiple beam interference, 2.28–2.36
- diffraction gratings, 2.28–2.29, 2.29*f*, 2.30*f*
 - Fabry-Perot interferometers, 2.33–2.36, 2.34*f*, 2.35*f*
 - plane-parallel plates, 2.30–2.33, 2.30*f*, 2.32*f*, 2.33*f*
- Multiple (volume) scattering, 9.2, 9.3, 9.8–9.17
- analytical theory of, 9.9–9.10, 9.9*f*
 - depolarization, 9.16–9.17, 9.17*f*
 - effective-medium representation, 9.8
 - radiative transfer, 9.10–9.13, 9.11*f*, 9.13*f*
 - speckle patterns, 9.15–9.16, 9.15*f*
 - weak localization, 9.13–9.17, 9.14*f*
- Multiple-angle-of-incidence ellipsometry (MAIE), 16.3
- Multiple-order retardation plates, 13.47–13.48
- Multiple-track read-write with diode laser arrays, 35.28
- Multiplexed image scanning, 30.23, 30.24*f*
- Multiplexed sensors, for fiber interferometers, 32.16
- Multiplexing, spatial, 23.8
- Multiplication, 11.3
- Multipliers, serial incoherent matrix-vector, 11.17–11.18, 11.18*f*
- Mutual coherence function, 2.36–2.38, 2.36*f*, 5.4, 5.10, 6.2–6.3
- Mutual intensity, 5.11, 6.3–6.4, 6.7, 6.8
- Nakamura biplates, 13.56
- Narrowband filters, 3.3
- Natural broadening, of lineshapes, 10.7
- Natural stop shift, of lenses, 22.3
- Net complex amplitude, 2.5
- Neutron scattering, 9.6
- Newton interferometers, 2.25
- Newton ring method, 29.22
- Newtonian form, of Gaussian equations, 18.4
- Newtonian imaging equation, 17.8
- Newtonian objectives, 29.6
- Newton's equation, for Gaussian focal lenses, 1.49
- Newton's ring pattern, 2.10, 2.19, 2.25–2.26
- Nicol curtate prisms, 13.16*f*, 13.17

- Nicol-type polarizers, 13.6, 13.6f, 13.10f, 13.15–13.18
 conventional, 13.6f, 13.15–13.16
 Glan-type vs., 13.8–13.9
 trimmed, 13.16–13.18, 13.16f
- Nikon N8008s camera, 25.7
- Nikon Plan Apochromat, 28.13, 28.14f
- Nipkow disks, 28.50, 28.51f
- Nippon Sheet Glass, 24.2, 24.6, 24.7
- Nodal planes, of Gaussian lenses, 1.48–1.49
- Nodal plane-to-nodal plane conjugate matrices, 1.69
- Nodal points, of lenses, 1.48, 1.48f, 1.49, 17.7
- Nodal rays, 17.16
- Noether's theorem, 1.21
- Noise:
 data, 35.24
 fixed-pattern, 26.11
 laser, 35.12, 35.24
 and optical disk data, 35.12, 35.23–35.24, 35.24f
 pattern, 26.11–26.12, 26.12f
 shot, 26.11
 in solid-state cameras, 26.11–26.12, 26.12f
- Noise equivalent exposure (NEE), 26.10
- Nomarski interferometers, 32.4, 32.5f
- Nomarski prisms, 28.40, 28.41
- Nonafocal lenses, 1.46 (*See also* Focal lenses)
- Noncalcite prisms, 13.23–13.24
- Noncosmological red shift, 5.23
- Nondepolarizing Mueller matrices, 14.24–14.25, 14.27–14.30
- Nondispersive prisms, 19.1–19.29
 Abbe's, 19.3t, 19.7, 19.7f–19.8f
 Amici (roof), 19.3t, 19.11, 19.12f
 and beam deviation, 19.2
 and beam displacement, 19.2
 Brashear-Hastings, 19.3t, 19.25, 19.25f
 Carl Zeiss, 19.3t, 19.16, 19.16f
 Dove, 19.3t, 19.9, 19.9f, 19.10, 19.10f
 Frankford Arsenal, 19.3t, 19.18–19.24, 19.18f–19.24f
 general deviation, 19.3t, 19.28, 19.28f–19.29f
 Goetz, 19.3t, 19.17, 19.17f
 and image inversion/reversion, 19.2
 Leman, 19.3t, 19.13, 19.13f
 Pechan, 19.3t, 19.11, 19.11f
 penta, 19.3t, 19.13, 19.14f
 Porro, 19.3, 19.3t, 19.5f, 19.6, 19.6f
 retroreflectors, 19.3t, 19.28, 19.28f
- Nondispersive prisms (*Cont.*):
 reversion, 19.3t, 19.14, 19.14f
 rhomboidal, 19.3t, 19.25, 19.25f
 right-angle, 19.3, 19.3t, 19.4f
 Risley, 19.3t, 19.25, 19.25f–19.27f, 19.27
 Schmidt, 19.3t, 19.12, 19.12f
 Wollaston, 19.3t, 19.15, 19.15f
- Nonhomogeneous polarization elements, 14.25–14.26, 14.26f, 15.7, 15.20
- Nonhomogeneous polarization elements (Mueller matrices), 14.25–14.26, 14.26f
- Nonlinear effects, in integrated optics, 21.12
- Nonnormal-incidence reflection:
 in Brewster angle reflection polarizers, 13.34–13.37, 13.34t–13.36t
 in pile-of-plates polarizers, 12.15–12.18, 12.16f, 12.17f
 in polarizing beam splitters, 13.41–13.42
- Nonnormal-incidence transmission:
 in Brewster angle transmission polarizers, 13.37–13.39, 13.38t–13.39t
 in interference polarizers, 13.39–13.41, 13.40f
 in pile-of-plates polarizers, 12.18–12.24, 12.19t–12.20t, 12.21f
 in polarizing beam splitters, 13.41–13.42
- Nonpolarizing elements, 15.7
- Nonpolarizing Mueller matrices, 14.8
- Nonradiating sources of light, 5.18
- Nonrectilinear distortion, 27.6
- Non-return to zero inverted (NRZI) scheme, 35.17
- Non-return to zero (NRZ) scheme, 35.17
- Nonrotationally symmetric systems, 1.74
- Nonspherical particles, scattering by, 7.15–7.17
- Normal congruence rays, 1.10
- Normal vectors, 1.18, 1.19
- Normal-incidence rotating-sample ellipsometers (NIRSE), 16.18
- Normalization, of Mueller matrices, 14.19
- Normalized spectrum, coherence functions for, 5.5–5.6
- Null ellipsometers, 16.11, 16.12
- Nulling interferometers, 32.20–32.21
- Numerical aperture, 1.78, 1.79, 17.9
- Nyquist frequency, 26.16–26.20, 26.17f
- Object relief distance, 18.11–18.12, 18.12f
- Object space, 1.26, 1.83
- Object space numerical aperture, 1.78
- Object space pupil diameter, 18.6

- Object transparencies, 11.4–11.5, 11.5*f*
- Objective optics, 30.30–30.33, 30.32*f*–30.33*f*
- Objective scanning, 30.5, 30.28
- Objective speckle, 33.9
- Objectives:
- for microscopes, 28.9–28.15, 28.10*t*, 28.11*t*
 - in afocal systems, 18.7
 - corrections for tube length, 28.13, 28.13*t*
 - coverslip correction, 28.10–28.11, 28.12*f*, 28.13
 - design of, 28.13–28.15, 28.14*f*–28.16*f*
 - field size, 28.13
 - working distance, 28.13
 - reflective and catadioptric (*see* Reflective and catadioptric objectives)
 - for telescopes
 - afocal Cassegrain-Mersenne, 29.9
 - afocal Gregorian-Mersenne, 29.12
 - Ritchey-Chretien, 29.8
 - three-mirror afocal, 29.29–29.30
- Object-space scanners, 30.18–30.23, 30.19*f*–30.23*f*
- Object-to-image distance, 1.51–1.52
- Oblique spherical aberrations, 1.90, 29.15, 29.21, 29.37
- Obliquity factor, 3.5
- Obscurations, of reflective and catadioptric objectives, 29.4–29.5, 29.4*f*
- Oculars, for microscopes, 28.16–28.17
- Off-axis, eccentric-pupil Paul-Gregorian objective, 29.28–29.29
- Off-axis double-pass grating spectrograph, 20.10, 20.13*f*
- Offner relay, 29.33
- Offset quantum wells, 21.19
- On-axis objective optics, 30.30
- One-electron atoms, 10.7–10.9, 10.8*f*, 10.9*f*
- 135° linear polarizers, Mueller matrices for, 14.10*t*
- 135° half-wave linear retarders, Mueller matrices for, 14.12*t*
- 135° quarter-wave linear retarders, Mueller matrices for, 14.12*t*
- Ophthalmic polarimetry, 15.39, 15.41
- Optic axis, of calcite crystals, 13.2, 13.2*f*, 13.3*f*, 13.3*n*
- Optical absorption spectrometers, 31.2–31.5, 31.5*f*
- Optical axes, 1.32, 14.8, 18.2, 29.5, 29.37
- Optical center point, of lenses, 17.16, 17.17
- Optical coherence microscopy (OCM), 28.44
- Optical coherence tomography (OCT), 22.2, 22.39, 22.40*f*, 28.43–28.44
- Optical constants, 7.12, 12.4–12.6, 16.5
- Optical disk data storage, 35.1–35.30
 - alternative storage media, 35.29, 35.30
 - automatic focusing, 35.12–35.14, 35.13*f*
 - automatic tracking, 35.14–35.17, 35.14*f*–35.16*f*
 - data format and layout for, 35.2–35.7, 35.3*f*–35.5*f*
 - developments in, 35.28–35.30
 - diffractive optics, 35.28, 35.28*f*, 35.29
 - direct overwrite, 35.30, 35.30*f*
 - materials for recording, 35.25–35.28, 35.26*f*, 35.27*f*
 - multiple-track read-write with diode laser arrays, 35.28
 - and optical path, 35.7–35.12, 35.7*f*, 35.9*f*–35.11*f*
 - readout, 35.21–35.24, 35.22*f*, 35.24*f*
 - thermomagnetic recording process, 35.17–35.20, 35.18*f*–35.20*f*
- Optical extent (*étendue*), 1.22, 1.81, 13.7
- Optical fibers, polarizers for, 13.57
- Optical hole burning (OHB), 10.18
- Optical holeburning (OHB) spectroscopy, 31.24–31.26, 31.24*f*–31.26*f*
- Optical invariants, 18.7
- Optical matched filtering, for pattern recognition, 11.12–11.14, 11.13*f*
- Optical metrology, 15.35
- Optical path, of optical disks, 35.7–35.12, 35.7*f*, 35.9*f*–35.12*f*
- Optical path difference (OPD), 2.7, 28.33–28.34, 28.35*f*, 28.36
- Optical path length (OPL), 1.11, 2.5
- Optical processing systems, for synthetic aperture radar data, 11.7–11.8
- Optical sine theorem, 17.5
- Optical spectrometers (*see* Spectrometers)
- Optical spectroscopy (*see* Spectroscopy)
- Optical theorem, 7.8
- Optical train, in microscopes, 28.3–28.5, 28.5*f*
- Optical transfer function (OTF):
 - calculations of, 4.3, 4.5
 - and camera lens performance, 27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.24
 - measurements of, 4.6–4.7
- Optical tube length, 17.10

- Optical tweezers, **28.55**
- Optically detected magnetic resonance
(ODMR), **31.21–31.23, 31.22f, 31.23f**
- Optoelectronic integrated circuit (OEIC), **21.2**
- Orthogonal matrices, **14.11**
- Orthoscopic imaging, **28.8**
- Orthotomic systems, **1.10, 1.12**
- Oscillations, wavelength interval between, **12.10**
- Osculating planes, of space curves, **1.18, 1.19**
- Outer product processors, **11.19**
- Output planes, conjugate matrices for, **1.68**
- Overillumination, **30.14**
- Oxides, thin film, **21.13–21.14**
- Pancharatnam phase, **32.11**
- Panoramic cameras, **25.25, 25.26f**
- Parabasal optics, **1.43**
- Paraboloid objective, **29.6**
- Parallax stereogram, **25.24**
- Parallel matrix-vector multipliers,
11.18–11.19, 11.18f
- Parallel-beam scanners, **30.23–30.25,**
30.24f–30.25f
- Paraxial chief rays, **1.75**
- Paraxial curvature, **1.32–1.33**
- Paraxial invariant, **1.41**
- Paraxial invariants, **1.77**
- Paraxial limit, of systems of revolution, **1.38**
- Paraxial matrices, for geometrical optics,
1.65–1.74
angle instead of reduced angle, **1.72**
arbitrary systems, **1.67**
and characteristic functions, **1.74**
conjugate matrices, **1.68–1.71, 1.73**
linearity, **1.66**
nonrotationally symmetric systems, **1.74**
operation on two rays, **1.68**
possible zeros, **1.68**
power matrix, **1.67**
skew rays, **1.73**
transfer matrices, **1.66**
two-ray specification, **1.72**
unit determinants, **1.67**
- Paraxial optics, **1.29, 1.37**
- Paraxial optics, of systems of revolution,
1.37–1.43
angle of incidence at a surface, **1.39**
axial object and image locations, **1.40**
image location and magnification, **1.42**
linearity of, **1.41**
- Paraxial optics, of systems of revolution (*Cont.*):
paraxial limit, **1.38**
principal focal lengths of surfaces, **1.39–1.40**
ray tracing, **1.40**
reflection and refraction, **1.38**
switching axial objects and viewing
positions, **1.43**
three-ray rule, **1.42**
transfer, **1.38**
two-ray paraxial invariant, **1.41**
- Paraxial pupils, **1.77**
- Paraxial rays, **1.35**
- Partially polarized light, **15.7**
- Particles, scattering by, **7.1–7.17**
coherent vs. incoherent arrays, **7.2–7.3**
concepts of, **7.4–7.5, 7.6f–7.7f, 7.8–7.10,**
7.9f, 7.10f
isotropic homogenous spheres, **7.11–7.14**
Mie, **7.11, 7.12**
nonspherical particles, **7.15–7.17**
regular particles, **7.14–7.15**
single particles, **7.2–7.3**
theories of, **7.3–7.4**
- Paschen-Runge configuration, **20.7, 20.11f, 20.14t**
- Passive autofocus systems, for cameras, **25.12**
- Passive devices, for integrated optics,
21.21–21.25, 21.22f–21.25f
- Pattern noise, **26.11–26.12, 26.12f**
- Pattern recognition, optical matched filtering
for, **11.12–11.14, 11.13f**
- Paul objectives, **29.28–29.29**
- Pauli spin matrices, **14.24, 14.41**
- Pechan prisms, **19.3t, 19.11, 19.11f**
- Pellin-Broca prisms, **20.6f**
- Penta prisms, **19.3t, 19.13, 19.14f**
- Percus-Yevick approximation, **9.5**
- Perfectly reflecting (PEC) surfaces, **8.10**
- Periphery cameras, **25.22**
- Periscopes, lenses in, **18.19, 18.19f**
- Periscopic lenses, **17.27, 17.27f**
- Perpendicular magnetic anisotropy, **35.26**
- Perpendicular-incidence ellipsometers (PIEs),
16.17–16.18, 16.18f
- Petzval (field) curvature:
in gradient index optics, **24.4, 24.6, 24.7**
of reflective and catadioptric objectives, **29.7,**
29.11, 29.14, 29.15, 29.32
as wavefront aberration, **1.91**
- Petzval lenses, **17.10, 17.28, 17.28f, 17.35f**
- Petzval sum, **28.15, 29.37**

- Pfund configuration, of dispersive prisms, 20.8*f*, 20.10, 20.13*f*
- Pfund objectives, 29.6
- Phase contrast microscopy, 28.28–28.29, 28.29*f*
- Phase plates, 28.28
- Phase retardation, 12.24
- Phase transfer functions (PTF), 4.3, 4.7
- Phase-conjugate interferometers, 32.17, 32.18, 32.18*f*
- Phased-arrays, for agile beam steering, 30.52–30.57, 30.53*f*, 30.62–30.63
- Phase-locked interferometers, 32.11–32.12, 32.12*f*
- Phase-matching, in integrated optics, 21.12
- Phase-sensitive detection (PSD), 31.8–31.11, 31.8*f*
- Phase-shifting interferometers, 32.10–32.11, 32.11*f*
- Phasors, 5.2
- Photo-activated localization microscopy (PALM), 28.23
- Photocopolymerization, 24.8
- Photocopying, 24.6
- Photo-elastic modulators (PEMs), 15.21, 16.13
- Photographic plates, 3.3
- Photography, wide-angle, 25.25
- Photoluminescence decay time, 31.12–31.15, 31.13*f*, 31.14*f*
- Photometric ellipsometers, 16.12–16.14, 16.13*f*, 16.14*f*
- Photon correlation spectroscopy (PCS), 9.8
- Photon migration approach, to radiative transfer, 9.12
- Photon transfer, 26.12
- Photonic integrated circuits (PICs):
of III-V materials, 21.17–21.20
in integrated optics, 21.2
in WDM systems, 21.37, 21.38
- Photopolarimeters, 16.13–16.16, 16.14*f*–16.16*f*
- Photoreceptors, in xerographic systems, 34.1, 34.2–34.4, 34.2*f*, 34.3*f*
- Photoresponse nonuniformity (PRNU), 26.11–26.12, 26.12*f*
- Photosensitive compounds, in microscopy, 28.55
- Physically realizable Mueller matrices, 14.40–14.42
- Picosecond and sub-picosecond relaxation, 31.14, 31.14*f*
- Pictures, instant, 25.8
- Piezoelectric transducer (PZT), 15.21, 15.24
- Pile-of-plates polarizers:
nonnormal-incidence reflection, 12.15–12.18, 12.16*f*, 12.17*f*
nonnormal-incidence transmission, 12.18–12.24
- Pincushion distortion, 1.91
- Piston error, 1.91
- Pixels, 30.8
- Planar lenses, 17.28
- Planar objects, transmissive, 6.3
- Planar secondary source of light, 5.9
- Plane mirrors, 1.25
- Plane of incidence, 1.23, 12.6
- Plane of polarization, 12.6*n*
- Plane waves, 2.4, 2.5*f*, 3.3, 3.17
decomposition of, 3.23
interference of, 2.8–2.9, 2.9*f*
and spherical waves, 2.9–2.11, 2.10*f*
- Plane-parallel plates, 2.19, 2.20*f*, 2.30–2.33, 2.30*f*, 2.32*f*, 2.33*f*
- Planes of vibration, 12.6
- Pockels cells, 28.45, 31.9
- Pockels effect, 15.23, 21.9 (*See also* Linear electro-optic effect)
- Poincaré spheres, 12.27–12.29, 14.4–14.6, 14.5*f*, 14.26*f*, 28.45
- Point characteristic function, 1.11, 1.14
- Point eikonal, 1.14, 1.17
- Point images, aberrations of, 1.85–1.92, 1.86*f*
- Point objects, image planes of, 28.19
- Point spread matrix, 15.36, 15.36*f*
- Point-angle characteristic function, 1.15–1.17
- Points, images of, 1.27
- Polacoat dichroic polarizers, 13.25, 13.26, 13.28
- Polanret system, 28.29
- Polar decomposition, of Mueller matrices, 14.39–14.40
- Polarimeters, 15.3–15.6
AxoScan Mueller matrix, 15.33
classes of, 15.5
complete and incomplete, 15.4
defined, 15.7
design metrics for, 15.24–15.25
division-of-amplitude, 15.5–15.6
division-of-aperture, 15.5
dual rotating retarder, 15.16
dual rotating retarder polarimeters, 15.16, 15.16*f*
imaging, 15.6
light-measuring, 15.3–15.4, 15.11–15.13

- Polarimeters (*Cont.*):
 Mueller, 15.4
 Mueller matrix, 15.26–15.27
 polarization modulation, 15.5
 sample-measuring, 15.4
 incomplete, 15.16–15.17
 for Mueller matrix elements,
 15.13–15.14, 15.13*f*
 spectropolarimeters, 15.6
 Stokes, 15.4, 15.5, 15.25
 time-sequential, 15.5
See also Photopolarimeters
- Polarimetric data-reduction equations,
 15.14–15.15
- Polarimetric measurement equation, 15.12,
 15.14–15.15
- Polarimetry, 15.3–15.41
 applications of, 15.29–15.41
 ellipsometry, 15.30–15.32, 15.31*f*, 15.32*f*
 liquid crystal cell and system testing,
 15.32–15.35, 15.32*f*, 15.33*f*, 15.34*t*
 ophthalmic polarimetry, 15.39, 15.41
 polarization aberrations, 15.35–15.37,
 15.35*f*–15.37*f*
 polarization light scattering, 15.38–15.39,
 15.38*f*–15.40*f*
 remote sensing, 15.37–15.38
 error analysis in, 15.27–15.29, 15.29*f*
 instruments for (*see* Polarimeters)
 Mueller matrices in, 15.8–15.9, 15.11
 elements of, 15.13–15.14
 in error analysis, 15.28
 singular value decomposition, 15.25–15.27
 polarimetric data-reduction equations,
 15.14–15.15
 polarimetric measurement equation,
 15.14–15.15
 and polarization elements, 15.17, 15.19–15.20
 polarization generators and analyzers,
 15.4–15.5
 polarization (retardance) modulators,
 15.20–15.24, 15.22*f*
 Stokes vectors in, 15.8–15.10
 terms in, 15.6–15.7
- Polariscopes, Sénarmont, 12.30
- Polarizance, 12.14*n*, 14.18
- Polarization, 12.3–12.30
 average degree of, 14.32–14.33
 and coherence theory, 5.22
 concepts and conventions, 12.4–12.6
- Polarization (*Cont.*):
 defined, 15.8
 degree of, 12.14–12.15
 of dichroic polarizers, 13.33
 false, 15.38
 Fresnel equations for, 12.6–12.13
 for absorbing materials, 12.10–12.13, 12.13*f*
 coordinate system for, 12.6–12.7, 12.7*f*
 for nonabsorbing materials,
 12.8–12.10, 12.9*f*
 generators and analyzers of, 15.4–15.5
 instrumental, 12.15
 magnetic circular, 31.21
 matrix methods for computing, 12.27–12.30
 and Mueller matrices, 14.7, 14.8,
 14.25–14.27, 14.33
 pile-of-plates polarizers, 12.15–12.24
 nonnormal-incidence reflection,
 12.15–12.18, 12.16*f*, 12.17*f*
 nonnormal-incidence transmission,
 12.18–12.24, 12.19*t*–12.20*t*, 12.21*f*
 plane of, 12.6*n*
 relations for polarizers, 12.14–12.15
 retardation plates, 12.24–12.27, 12.25*f*, 12.26*f*
See also related topics
- Polarization aberration function (PAF), 15.35
- Polarization aberrations, 15.35–15.37,
 15.35*f*–15.37*f*
- Polarization analyzer, 15.11
- Polarization and Directionality of Earth's
 Reflectances (POLDER) instrument, 15.37
- Polarization artifacts, 15.38
- Polarization coupling, 15.8
- Polarization critical region, 15.28
- Polarization instruments, 12.29
- Polarization interferometers, 32.4, 32.5*f*
- Polarization light scattering, 15.38–15.39,
 15.38*f*–15.40*f*
- Polarization modulation polarimeters, 15.5
- Polarization (retardance) modulators,
 15.20–15.24, 15.22*f*
- Polarization spectrometers, 31.15–31.23
 optically detected magnetic resonance,
 31.21–31.23, 31.22*f*, 31.23*f*
 polarized absorption by, 31.15–31.17,
 31.15*f*, 31.18*f*
 polarized absorption/luminescence
 techniques, 31.17, 31.19–31.21, 31.19*f*
 principles of, 31.15, 31.15*f*
- Polarization state detectors (PSDs), 16.10

- Polarization state generators (PSGs), 16.10
Polarization-dependent loss (PDL), 14.17
Polarized absorption, 31.15–31.21, 31.15f, 31.18f, 31.19f
Polarized light, 15.8
Polarizer-compensator-sample analyzer (PCSA) ellipsometer arrangement, 16.10–16.14, 16.11f
Polarizers, 13.1–13.57
 beam-splitter prisms as, 13.6, 13.18–13.22
 Foster, 13.7, 13.18f, 13.21–13.22
 Glan-Thompson, 13.18f, 13.22
 Rochon, 13.7, 13.18–13.21, 13.18f, 13.19f, 13.24
 Sénarmont, 13.7, 13.18, 13.18f, 13.21
 Wollaston, 13.7, 13.18, 13.18f, 13.21, 13.24
 circular, 15.17–15.19
 compensators, 13.53–13.56, 13.54f, 13.55f
 Cotton, 13.21
 defined, 15.8
 dichroic and diffraction-type, 13.24–13.33, 13.26f, 13.27f
 dichroic polarizing coatings, 13.28
 measuring polarization of, 13.33
 pyrolytic-graphite polarizers, 13.28–13.29, 13.29f
 sheet polarizers, 13.25–13.28
 wire-grid and grating polarizers, 13.30–13.33, 13.31f, 13.32t
 elliptical, 15.17–15.18
 Feussner prisms, 13.6, 13.7, 13.22–13.23, 13.22f
 Glan-Foucault prisms, 13.7, 13.9, 13.11f, 13.12–13.14
 Glan-type prisms, 13.6, 13.6f, 13.8–13.15
 Frank-Ritter-type, 13.6, 13.6f, 13.13–13.14
 Glan-Foucault, 13.7, 13.9, 13.11f, 13.12–13.14
 Glan-Thompson type, 13.6f, 13.9–13.12, 13.10f, 13.11f, 13.27
 Lippich-type, 13.6, 13.6f, 13.7, 13.9n, 13.10–13.14, 13.10f, 13.11f
 half-shade devices, 13.56–13.57
 ideal, 14.8–14.10, 14.10t, 15.7
 imperfect, 13.33
 miniature, 13.57
 Nicol-type, 13.6, 13.6f, 13.10f, 13.15–13.18
 conventional, 13.6f, 13.15–13.16
 Glan-type vs., 13.8–13.9
 trimmed, 13.6f, 13.7, 13.16–13.18, 13.16f
Polarizers (*Cont.*):
 noncalcite prisms as, 13.23–13.24
 nonnormal-incidence reflection by
 Brewster angle reflection polarizers, 13.34–13.37, 13.34t–13.36t
 pile-of-plates polarizers, 12.15–12.18, 12.16f, 12.17f
 polarizing beam splitters, 13.41–13.42
 nonnormal-incidence transmission by
 Brewster angle transmission polarizers, 13.37–13.39, 13.38t–13.39t
 interference polarizers, 13.39–13.41, 13.40f
 pile-of plates polarizers, 12.18–12.24, 12.19t–12.20t, 12.21f
 polarizing beam splitters, 13.41–13.42
 prism, 13.2–13.8, 13.2f–13.3f, 13.4t–13.5t, 13.6f
 relations for, 12.14–12.15
 retardation plates as, 13.43–13.53, 13.43t–13.44t, 13.50f, 13.53t
Polarizing angle, 12.12, 12.15
Polarizing beam splitter (PBS) prisms, 13.6, 13.41–13.42
Polarizing beam splitters (PBSs), 13.41–13.42, 35.22–35.23, 35.22f
Polarizing coatings, dichroic, 13.28
Polaroid dichroic polarizers, 13.25–13.28, 13.26f
Polygon scanners, 30.34–30.38, 30.34f, 30.35f
Porro prisms, 19.3, 19.3t, 19.5f, 19.6, 19.6f
Postobjective scanning, 30.5, 30.29, 30.29f, 30.30f
Powder cloud development, in xerographic systems, 34.9, 34.9f
Power (Gaussian lenses), 1.46–1.47, 1.47f
Power density, 30.25
Power exponential (PEX) model, of surface finish, 8.15
Power matrix, 1.67
Power per pixel, 1.80
Power spectra, for surface scattering, 8.12–8.13, 8.13f
Power transfer, 30.25–30.28, 30.27f, 30.28f
Poynting vectors, 1.8, 3.3
Preflash, of cameras, 25.16
Preobjective scanning, 30.5, 30.28, 30.29f
Primary focal length (term), 3.12
Principal plane-to-principal plane conjugate matrices, 1.69
Principal rays (term), 1.75, 17.8
Principal transmittance (term), 12.14–12.16
Prism polarizers, 13.2–13.8, 13.2f–13.3f, 13.4t–13.5t, 13.8f (*See also specific types*)

- Prism spectrometers, 20.2–20.3, 20.3f
 Prismatic facets, 30.34–30.35
 Prisms (*See also* Dispersive prisms and gratings;
 Nondispersive prisms)
 axis wander of, 13.15
 beam-splitter, 13.7, 13.18–13.22, 13.18f
 Bertrand-type Feussner, 13.23
 Brewster angle, 13.13
 of calcite, 13.20, 13.23
 Feussner, 13.6, 13.7, 13.22–13.23, 13.22f
 Foster, 13.7, 13.18f, 13.21–13.22
 Foucault, 13.7, 13.17
 Frank-Ritter-type, 13.6, 13.6f, 13.13–13.14
 Fresnel's biprism, 2.16, 2.16f
 Glan-Foucault, 13.7, 13.9, 13.11f, 13.12–13.14
 Glan-Taylor, 13.7, 13.9n, 13.10–13.14,
 13.10f, 13.11f
 Glan-Thompson, 13.6, 13.6f, 13.9–13.12,
 13.10f, 13.18f, 13.22
 field angle of, 13.12
 sheet polarizers vs., 13.27
 transmission by, 13.9–13.10, 13.11f
 Glan-type, 13.6, 13.6f, 13.8–13.15
 Glazebrook, 13.6, 13.9
 Halle, 13.16f, 13.17
 Hartnack-Prazmowski, 13.16f, 13.17
 Jellett-Cornu, 13.56
 length-to-aperture (L/A) ratio, 13.7
 Lippich-type, 13.6, 13.9n, 13.10f, 13.11f,
 13.12–13.14, 13.12n, 13.56
 Marple-Hess, 13.12, 13.13
 Nicol curtate, 13.16f, 13.17
 Nomarski, 28.40, 28.41
 noncalcite, 13.23–13.24
 polarizing beam splitter, 13.6, 13.41–13.42
 Rochon, 13.7, 13.18–13.21, 13.18f,
 13.19–13.20, 13.19f, 13.24
 semifield angle of, 13.7
 Sénarmont, 13.7, 13.18, 13.18f, 13.21
 Steeg and Reuter Nicol, 13.17
 Wollaston, 13.7, 13.18, 13.18f, 13.21, 13.24,
 28.39, 28.40, 32.4
 Projective transformation, 1.56 (*See also*
 Collineation)
 Propagation of light, coherence theory and,
 5.13–5.19, 5.14f–5.16f
 Propagation of mutual intensity, 6.4
 Pseudo-Brewster angle, 12.13
 Pulse width modulation (PWM), in
 thermomagnetic recording, 35.18–35.19
 Pupil aberrations, 1.76
 Pupil angular magnification, 1.78
 Pupil distortion, 1.78
 Pupil imaging, 1.76
 Pupil magnification, 1.76
 Pupils, of lenses, 1.76–1.79, 1.76f, 1.78f,
 17.8–17.9
 Purcell, M., 7.16
 Purcell-Pennypacker method, 7.15
 Pure diattenuators, 15.8
 Pure retarders, 15.8
 Pushbroom scan, 30.18
 Pyramidal facets, 30.34–30.35, 30.34f, 30.35f
 Pyrolytic-graphite polarizers, 13.28–13.29, 13.29f
 Quantitative phase microscopy, 28.27
 Quantum coherence theory, 5.2
 Quantum electrodynamic (QED) shifts, 10.4
 Quantum-confined Stark effect (QCSE),
 21.11–21.12, 21.32
 Quantum-well intermixing, 21.19, 21.20
 Quarter pitch length, of the rod, 24.6
 Quarter-wave circular retarders, Mueller
 matrices for, 14.12t
 Quarter-wave linear retarders, Mueller matrices
 for, 14.11, 14.12t
 Quarter-wave plates, 12.25–12.27, 12.26f
 Quartz retardation plates, 13.46–13.48
 Quasi-homogeneous sources of light,
 5.11–5.12, 5.19
 Quasi-monochromatic sources of light, 5.11
 Racah parameters, for ion energy levels, 10.12
 Radar, synthetic aperture, 11.6–11.8,
 11.7f–11.8f
 Radial gradients, 24.5–24.8, 24.5f, 24.7f
 Radial symmetry, of scanners, 30.5
 Radiance, 5.8
 Radiant emittance, 5.7–5.8
 Radiant intensity, 5.8
 Radiation:
 coherent, 30.2, 30.25–30.26
 incoherent, 30.2, 30.26, 30.27
 infrared, 13.47
 Radiation fields, coherence theory and,
 5.15–5.16, 5.15f, 5.16f
 Radiation modes, of optical waveguides, 21.4
 Radiative lifetime, 31.12–31.13, 31.13f
 Radiative transfer, in volume scattering,
 9.10–9.13, 9.11f, 9.13f

- Radio astronomy, 5.23
Radiometry, statistical, 5.22
Radius of torsion, for space curves, 1.19
Radius of torsion, of space curves, 1.19
Ramachandran, G. N., 12.28, 13.53
Raman scattering, 31.30–31.31, 31.31*f*
Raman-Nath diffraction, 11.9, 11.9*f*
Ramsden disk, 28.8, 28.16
Range of focus, 1.85
Rapid Rectilinear lenses, 17.27
Rare earth ions, tri-positive, 10.16–10.18, 10.16*t*, 10.17*f*
Raster output scanning (ROS) systems, 34.4
Ray aberrations, 1.87–1.88
Ray densities, 1.88
Ray equation, 1.20
Ray fans, 1.35
Ray intercept diagrams, 1.87
Ray optics, 1.8
Ray paths, 1.10–1.13, 24.2
Ray tracing:
 for binary optics, 23.4, 23.6
 in systems of revolution, 1.35–1.37, 1.36*f*, 1.40
Rayleigh criterion, 17.37, 28.6, 28.18
Rayleigh criterion of resolving power, 3.26
Rayleigh index, 8.6
Rayleigh radius value, 30.10
Rayleigh range of origin, 5.14, 5.16
Rayleigh resolution, 30.56, 33.17
Rayleigh scattering, 7.11, 9.17, 31.30
Rayleigh-Gans approximation, 7.9, 7.9*f*
Rayleigh-Rice (RR) approximation, 8.4, 8.9–8.12
Rayleigh's diffraction integral, 5.13
Rayleigh-Sommerfeld diffraction, 3.9, 3.10, 3.23, 3.29
Rays, 1.8–1.13
 chief, 1.75, 17.8, 29.20, 29.37
 for collineation, 1.61
 defined, 1.8–1.9
 differential geometry of, 1.19–1.21
 direction of, 1.10
 expansions about, 1.16
 fields of, 1.13
 finite, 1.35
 groups of, 1.10
 Hamilton's equations for, 1.21
 in heterogeneous media, 1.9, 1.18–1.22
 image-forming, 1.74
 images about known rays, 1.43–1.44, 1.44*f*
Rays (*Cont.*):
 invariance properties of, 1.10
 in lenses, 1.35
 marginal, 1.75
 meridional, 1.37
 meridional, 1.35
 nodal, 17.16
 normal congruence, 1.10
 paraxial, 1.35, 1.75
 paraxial matrices for, 1.68, 1.73
 paths of, 1.10–1.13
 principal, 1.75, 17.8
 principal index of, 13.3
 real and virtual, 1.10, 1.35
 reversibility of, 1.9
 skew, 1.35, 1.73
 variational integral of, 1.19
Reactive ion etching (RIE), 21.18–21.19
Readout, from optical disk data storage, 35.21–35.24
Real pupils, 1.76
Real rays, 1.10, 1.35
Rear focal lengths, 1.40
Rear focal points, 1.40, 1.47
Rear principal plane, of Gaussian focal lenses, 1.48
Received images, 1.26
Receiving surfaces, in imaging, 1.26
Reciprocity theorem, coherence theory and, 5.17–5.18
Recorded images, 1.26
Recording, of optical disk data, 35.25–35.28, 35.26*f*, 35.27*f*
Rectangular apertures, 3.19–3.20, 3.20*f*, 3.25, 3.26
Rectilinear distortion correction, 27.6, 27.13*f*, 27.14*f*
Red eye, cameras and, 25.16
Red shift, noncosmological, 5.23
Rediagonalization, of index ellipsoid equation, 21.10
Reference spheres, for wavefronts, 1.86
Reflectance, 12.17
Reflecting afocal lenses, 18.19–18.21, 18.20*f*
Reflection(s):
 in Gaussian lens systems, 1.55
 in homogeneous media, 1.25
 Mueller matrices for, 14.21–14.22
 nonnormal-incidence, 12.15–12.18
 and phase changes, 12.12–12.13
 and ray tracing, 1.37

- Reflection(s) (*Cont.*):
of systems of revolution, 1.38, 1.39
total internal, 13.20, 21.3
unfolded, 1.32
volumetric, 7.7f
- Reflective and catadioptric objectives, 29.1–29.38
afocal telescope designs
 Cassegrain-Mersenne, 29.9
 Gregorian-Mersenne, 29.12
 three-mirror, 29.29–29.30
Altenhof, 29.32–29.33
anastigmatic designs, 29.12–29.13
aplanatic designs, 29.11–29.13
Baker-Nunn, 29.22
Cassegrain designs, 29.6, 29.7
 afocal Cassegrain-Mersenne telescope, 29.9
 dual magnification, 29.9–29.10
 with field corrector and spherical secondary, 29.8–29.9
 Houghton-Cassegrain, 29.22–29.23
 Mangin-Cassegrain with correctors, 29.24
 reflective Schmidt-Cassegrain, 29.17
 Schmidt-Cassegrain, 29.16–29.17
 Schmidt-meniscus Cassegrain, 29.21
 with Schwarzschild relay, 29.32
 solid Makutsov-Cassegrain, 29.19
 spherical-primary, with reflective field corrector, 29.9
 three-mirror, 29.30
Cook three-mirror, 29.31
correctors, in designs
 aplanatic, anastigmatic Schwarzschild
 with aspheric corrector plate, 29.13
 Cassegrain with spherical secondary and field corrector, 29.8–29.9
 Mangin-Cassegrain with correctors, 29.24
 Ritchey-Chretien telescope with two-lens corrector, 29.8
 spherical-primary Cassegrain with reflective field corrector, 29.9
 three-lens prime focus corrector, 29.10
Couder, 29.12
Dall-Kirkham, 29.8
Eisenburg and Pearson two-mirror, three reflection, 29.25
features of, 29.2–29.5, 29.3f–29.5f
field-of-view plots, 29.34–29.35, 29.35f, 29.36f
flat-medial-field designs, 29.11
Gabor, 29.20
glass varieties for, 29.2, 29.2t
- Reflective and catadioptric objectives (*Cont.*):
Herschelian catadioptric, 29.27
Houghton designs, 29.22–29.23
Korsch designs, 29.30–29.32, 29.34
Maksutov designs, 29.19, 29.20
Mangin designs, 29.7, 29.24
Mersenne designs, 29.9, 29.12
Paul designs, 29.28–29.29
Ritchey-Chretien with two-lens corrector, 29.8
Schiefspiegler, 29.26–29.27
Schmidt designs, 29.14
 Baker super-Schmidt, 29.21
 field-flattened, 29.14–29.15
 reflective, 29.15
 reflective Schmidt-Cassegrain, 29.17
 Schmidt-Cassegrain, 29.16–29.17
 Schmidt-meniscus Cassegrain, 29.21
 Shafer-relayed-virtual, 29.17–29.18, 29.18f
 solid, 29.16
Schwarzschild designs, 29.12–29.13, 29.32
SEAL, 29.27–29.28
Shafer designs
 five mirror unobscured, 29.33–29.34
 four mirror unobscured, 29.33
 Shafer relayed virtual Schmidt, 29.17–29.18, 29.18f
 two-mirror three reflection, 29.25
Shenker, 29.23
spherical primaries in designs, 29.9, 29.11, 29.18, 29.30
for telescopes
 afocal Cassegrain-Mersenne, 29.9
 afocal Gregorian-Mersenne, 29.12
 Ritchey-Chretien, 29.8
 three-mirror afocal, 29.29–29.30
terminology, 29.36–29.38
three-mirror designs, 29.28–29.32, 29.34
Wetherell and Womble three-mirror, 29.31
Wright, 29.15
Yolo, 29.26
- Reflective Schmidt objective, 29.15
Reflective Schmidt-Cassegrain objective, 29.17
Reflective systems, 1.9
- Refraction:
in calcite, 13.2–13.6, 13.4t–13.5t
double, 13.2–13.6, 13.2f–13.3f, 13.4t–13.5t
in Gaussian lens systems, 1.54
in homogeneous media, 1.24–1.25
Mueller matrices for, 14.20–14.21
ray tracing, 1.37
in systems of revolution, 1.38, 1.39

- Refraction gradients, index of, **24.1** (*See also* Gradient index (GRIN) optics)
- Refractive index (index of refraction), **1.9**
for Brewster angle transmission polarizers, **12.21–12.22**
complex, **7.12–7.13, 12.5, 12.6**
distributed, **24.1**
in gradient index optics, **24.2–24.3**
in integrated optics, **21.8–21.9**
of polarizers, **12.16, 12.18**
for rays in heterogeneous media, **1.21–1.22**
of shallow radical gradients, **24.7–24.8**
- Refractive optics, **23.7, 23.8**
- Refractive systems, **1.9**
- Region of sag (axial gradients), **24.4f**
- Relativity, of conjugate matrices, **1.70**
- Relay lenses, **17.10**
- Relay trains, in afocal lenses, **18.17–18.19, 18.17f, 18.18f**
- Remote sensing, polarimetry and, **15.37–15.38**
- Remote sensing scanners, **30.2–30.4, 30.14–30.25**
circular scan, **30.16, 30.18f**
compound mirror optics configurations for, **30.15–30.16, 30.15f**
multiplexed image scanning by, **30.23, 30.24f**
object- and image-space, **30.18–30.23**
parallel-beam, **30.23–30.25, 30.24f–30.25f**
pushbroom scan, **30.18**
resolution of, **30.6–30.8, 30.7f**
rotating wedge, **30.16, 30.17f**
single-mirror, **30.14, 30.15f**
two-dimensional, **30.18, 30.19f**
- Resolution:
of cameras, **25.5–25.6, 25.6f**
of microscopes, **28.17–28.24**
Airy disk and lateral resolution, **28.17–28.19, 28.18f, 28.19f**
depth of field, **28.22–28.23**
depth of focus, **28.22**
three-dimensional diffraction pattern, **28.19–28.22, 28.20f, 28.21f**
and objective optics, **30.33**
- Rayleigh, **30.56, 33.17**
- of scanners, **30.6–30.14**
data rates and remote sensing, **30.6–30.8, 30.7f**
input/output scanning, **30.8–30.14, 30.10f, 30.10t, 30.11t, 30.12f–30.13f**
of solid-state cameras, **26.15–26.16, 26.16f**
- Resolution limit, **1.80**
- Resonant scanners, **30.41–30.44, 30.43f, 30.44f**
- Responsivity, of solid-state cameras, **26.9–26.10**
- Retardance, **14.6, 15.8**
- Retardance modulators, **15.20**
- Retardance space, **14.6**
- Retardation plates, **12.24–12.27, 12.25f, 12.26f, 13.43–13.53, 13.43t–13.44t**
achromatic, **13.48–13.52, 13.50f, 13.53t**
composite, **13.52, 13.53**
crystalline-quartz, **13.46–13.48**
defined, **15.8**
mica, **13.45–13.46**
quarter-wave and half-wave, **12.24–12.27, 12.25f, 12.26f**
rhomb-type, **13.52, 13.53t**
variable, **13.53**
- Retarder space, **14.14–14.15, 14.15f**
- Retarders:
defined, **15.8**
Mueller matrices for, **14.11–14.15, 14.12t, 14.15f**
- Retrofocus lenses, **27.2** (*See also* Inverted telephoto camera lenses)
- Retro-reflection testing and correction, **15.28–15.29, 15.29f**
- Retroreflectors, **19.3t, 19.28, 19.28f**
- Return-path ellipsometers, **16.16–16.17, 16.17f**
- Return-path ellipsometers (RPEs), **16.16–16.17, 16.17f**
- Reverse telephoto lenses, **17.29, 17.34f**
- Reversibility, of rays, **1.9**
- Reversion prisms, **19.3t, 19.14, 19.14f**
- Revolution, systems of, **1.32–1.43**
paraxial optics of, **1.37–1.43**
ray tracing in, **1.35–1.37**
surfaces, **1.32–1.35**
unfolded reflections, **1.32**
- Rhomboidal prisms, **19.3t, 19.25, 19.25f**
- Rhomb-type retardation plates, **13.52, 13.53t**
- Riccati-Bessel functions, **7.12**
- Right circular polarizers, Mueller matrices for, **14.10t**
- Right half-wave circular retarders, Mueller matrices for, **14.12t**
- Right-angle prisms, **19.3, 19.3t, 19.4f**
- Right-circularly polarized light, **12.27, 12.28n**
- Ring field lens design, **18.22**
- Risley prisms, **19.3t, 19.25–19.27, 19.25f–19.27f**
- Ritchey-Chretien objectives, **29.7–29.8**

- Ritchey-Chretien primaries, 29.10
 Rochon prisms, 13.7, 13.18–13.21, 13.18f, 13.19f, 13.24
 Roof prisms, 19.11, 19.12f
 (See also Amici prisms)
 Rotating retarders, 15.20
 Rotating wedge scanners, 30.16, 30.17f
 Rotating-analyzer ellipsometer (RAE), 16.13, 16.13f, 16.14
 Rotating-compensator fixed analyzer (RCFA) photopolarimeter, 16.14
 Rotating-detector ellipsometer (RODE), 16.14, 16.14f
 Rotating-element photopolarimeters (REPs), 16.13
 Rotation sensors, for fiber interferometers, 32.14–32.15, 32.15f
 Rotational spectra, 10.20–10.22
 Rotationally symmetric lenses, 1.27, 1.60–1.62, 1.62f
 Rotationally symmetric systems, 1.17, 1.89–1.90
 Routers, waveguide grating, 21.24
 Rowland circle, 20.5, 20.7, 20.8, 20.10f
 Rydberg constant, 10.3
 Rytov's series of exponential approximations, 9.4
- Sag, of surfaces, 1.32, 1.33f
 Sagittal fans and foci, 1.35
 Sagnac interferometers, 21.35, 21.36, 21.36f, 32.3–32.4, 32.4f
 Sampled tracking, on optical disks, 35.16, 35.16f
 Sample-measuring polarimeters, 15.4, 15.13–15.14, 15.13f, 15.16–15.17
 Sampling, with solid-state cameras, 26.16–26.19, 26.17f–26.19f
 SAOBIC processor, 11.20
 Saturated absorption spectroscopy, 31.24–31.26, 31.24f–31.26f
 Saturation equivalent exposure (SEE), 26.10–26.11
 Savart plates, 13.56
 Scalar diffraction theory, for binary optics, 23.10–23.13, 23.11t, 23.12f, 23.13f
 Scalar field amplitude, 5.3
 Scaling law, of spectrum of light, 5.21
 Scan error reduction, 30.48–30.51, 30.49t, 30.50f, 30.51f
 Scan magnification, 30.5, 30.12–30.14
- Scanners, 30.1–30.63
 acousto-optic, 30.44–30.45
 agile beam steering, 30.51–30.63
 decentered lens and microlens arrays, 30.57–30.60, 30.58f–30.60f, 30.62–30.63
 digital micromirror devices, 30.60–30.61
 gimbal-less two-axis scanning micromirrors, 30.61–30.62, 30.62f
 phased-array, 30.52–30.57, 30.53f, 30.62–30.63
 electro-optic (gradient), 30.45–30.48, 30.46f–30.48f
 error reduction in, 30.48–30.51, 30.49t, 30.50f, 30.51f
 galvanometer and resonant, 30.41–30.44, 30.43f, 30.44f
 holographic, 30.38–30.41, 30.40f–30.42f
 input/output scanning, 30.2, 30.4–30.6, 30.4t, 30.25–30.34
 objective, preobjective, and postobjective, 30.28–30.29, 30.29f, 30.30f
 objective optics, 30.30–30.33, 30.32f–30.33f
 power density and power transfer of, 30.25–30.28, 30.27f, 30.28f
 resolution of, 30.8–30.14, 30.10f, 30.10t, 30.11t, 30.12f–30.13f
 Keplerian afocal lenses for, 18.13, 18.13f
 modulation transfer function (MTF) for, 4.6
 monogon and polygon, 30.34–30.38, 30.34f, 30.35f
 remote sensing, 30.2–30.4, 30.14–30.25
 circular scan, 30.16, 30.18f
 compound mirror optics configurations, 30.15–30.16, 30.15f
 multiplexed image scanning, 30.23, 30.24f
 object- and image-space, 30.18–30.23, 30.19f–30.23f
 parallel-beam, 30.23–30.25, 30.24f–30.25f
 pushbroom scan, 30.18
 rotating wedge, 30.16, 30.17f
 single-mirror, 30.14, 30.15f
 two-dimensional, 30.18, 30.19f
 resolution of, 30.6–30.14
 data rates and remote sensing, 30.6–30.8
 input/output scanning, 30.8–30.14, 30.10f, 30.10t, 30.11t, 30.12f–30.13f
- Scanning, active, 30.4

- Scattering:
- and backscattering, 6.5–6.7, 6.5f, 9.14–9.15, 9.14f
 - Brillouin, 31.30
 - by coated spheres, 7.14
 - coherent and incoherent, 9.2, 9.3
 - by cylinders, 7.14
 - Mie, 7.11, 7.12, 9.17
 - molecular, 7.11
 - neutron, 9.6
 - in optical spectrometers, 31.30–31.31, 31.31f
 - and polarization, 15.38–15.39, 15.38f–15.40f
 - Raman, 31.30–31.31, 31.31f
 - Rayleigh, 7.11, 9.17, 31.30
 - theory of, 9.3–9.4, 9.4f
 - x-ray, 9.6
 - See also related topics, e.g.:* Volume (multiple) scattering
- Scattering, by particles, 7.1–7.17
- in coherent vs. incoherent arrays, 7.2–7.3
 - isotropic homogenous spheres, 7.11–7.14
 - Mie scattering, 7.11, 7.12
 - nonspherical particles, 7.15–7.17
 - regular particles, 7.14–7.15
 - single particles, 7.2–7.3
 - theories of, 7.3–7.4
 - volume scattering vs., 9.2–9.3
- Scattering cross section, 7.4
- Scattering length, 9.6
- Scattering matrices, 7.10, 16.8–16.9
(*See also* Mueller matrices)
- Scattering planes, 7.9
- Scattering potentials, 9.6
- Scheimpflug condition, 1.61, 17.6f, 17.7, 18.4, 25.18, 25.19f
- Scheimpflug rule, 18.4, 18.8
- Schell model sources (of light), 5.11
- Schiefspiegler objectives, 29.26–29.27
- Schmidt objectives, 29.14
- Baker super-Schmidt, 29.21
 - field-flattened, 29.14–29.15
 - reflective, 29.15
 - reflective Schmidt-Cassegrain, 29.17
 - Schmidt-Cassegrain, 29.16–29.17
 - Schmidt-meniscus Cassegrain, 29.21
 - Shafer-relayed-virtual, 29.17–29.18, 29.18f
 - solid, 29.16
- Schmidt prisms, 19.3t, 19.12, 19.12f
- Schrödinger equation, 10.4
- Schwarzschild arrangement, for McCarthy objective, 29.33
- Schwarzschild objectives, 29.12–29.13, 29.32
- Scophony TV projection system, 30.45
- Scorotron, in xerographic systems, 34.2, 34.3f
- SEAL objective, 29.27–29.28
- Second Brewster angle, 12.13
- Secondary magnification, 29.12, 29.38
- Secondary sources of light, 5.9–5.10
- Secondary spectrum, 29.7, 29.38
- Second-harmonic interferometers, 32.17–32.18, 32.18f
- Sectors (optical disk data), 35.4, 35.6–35.7
- Seek operation, on optical disks, 35.17
- Seidel, Philipp Ludwig von, 17.27
- Seidel aberrations, 1.90, 29.38
- Selective area epitaxy, 21.19–21.20
- Self-centering lens springs, 22.8, 22.8f
- Self-coherence function, 2.41
- Selfoc lenses, 24.2, 24.7f
- Semiconductors, complementary metal-oxide, 26.8–26.9
- Semifield angles, of prisms, 13.7
- Sénarmont compensators, 13.53, 28.38
- Sénarmont polariscopes, 12.30
- Sénarmont prisms, 13.7, 13.18, 13.18f, 13.21
- Sensors:
- active pixel, 26.2, 26.8–26.9, 26.8f
 - generalized, 32.15–32.16, 32.15f
 - multiplexed, 32.16
 - rotation, 32.14–32.15, 32.15f
- Serial incoherent matrix-vector multipliers, 11.17–11.18, 11.18f
- Sewer cameras, 25.22–25.23
- Shafer objectives:
- five mirror unobscured, 29.33–29.34
 - four mirror unobscured, 29.33
 - Shafer relayed virtual Schmidt, 29.17–29.18, 29.18f
 - two-mirror three reflection, 29.25
- Shallow radial gradient index (SRGRIN), 24.7
- Shape factor, of lenses, 17.12–17.13, 17.13f
- Shearing interferometers, 32.4, 32.6f
- Sheet polarizers, 13.25–13.28
- Shenker objective, 29.23
- Shenker objectives, 29.23
- Shift invariance, 4.2
- Shot noise, of solid-state cameras, 26.11
- Siemens star, 28.30f, 28.47f
- Sierpinski Gasket, 8.9
- Signal, for solid-state cameras, 26.9–26.11, 26.13–26.14

- Signal-to-noise ratio (SNR):
of optical disk data, 35.23
of solid-state cameras, 26.13–26.14
- Silicon photonics transmission, 21.14–21.16,
21.15*f*, 21.38–21.40
- Silicon-on-insulator (SOI) technology,
21.14–21.15, 21.15*f*
- Sine condition, for stigmatic imaging, 1.30–1.31
- Singham, Shermila Brito, 7.16
- Single element lenses, 17.12–17.17,
17.13*f*–17.16*f*, 17.17*t*
- Single lens reflex (SLR) cameras:
autofocus, 25.12–25.14, 25.13*f*
features of, 25.8, 25.9*f*
formats of, 25.18
lenses for, 27.1–27.2, 27.3*f*–27.4*f*
normal lenses for, 27.2, 27.3*f*–27.4*f*
and time lag, 25.8–25.9
- Single molecule high-resolution colocalization
(SHREC), 28.23
- Single scattering, 7.2–7.3
coherent and incoherent, 9.4–9.7, 9.6*f*
dynamic, 9.7–9.8, 9.7*f*
and volume scattering, 9.2–9.8
- Single sideband edge enhancement (SSEE)
microscopy, 28.29–28.33, 28.30*f*–28.33*f*
- Single speckle, 8.17
- Single-component development,
in xerographic systems, 34.9, 34.9*f*
- Single-mirror scanners, 30.14, 30.15*f*
- Single-mode waveguides, 21.4
- Single-order plates, 13.47
- Singlet lenses, 17.37–17.38
- Singular value decomposition (SVD),
15.25–15.27
- Sinusoidal ray paths, 24.2
- Skew invariant, 1.21, 1.23
- Skew rays, 1.35, 1.73
- Skewness, 1.23
- Skin depth, 7.13
- Slater parameters, 10.12
- Slow axis, 12.25, 15.8
- Small-perturbation approximation, for surface
scattering, 8.9–8.12
- SMARTCUT technique, 21.14, 21.15
- Smith invariant, 1.77 (*See also* Two-ray
paraxial invariant)
- Smith reflectors, 28.44
- Snell's law, 1.24, 1.38, 12.16, 13.3, 13.5, 13.19
- Solano objectives, 29.26
- Soleil compensators, 13.55–13.56, 13.55*f*
- Soleil-Babinet compensators, 35.21*n*
- Sol-gel formed glass, 24.8
- Solid Makutsov-Cassegrain objective, 29.19
- Solid Schmidt objective, 29.16
- Solid state spectroscopy, 10.22–10.26,
10.23*f*–10.27*f*
- Solid-state cameras, 26.1–26.20
applications, 26.3
array performance in, 26.9–26.12, 26.12*f*
and charge injection devices, 26.6–26.7,
26.6*f*–26.8*f*
and charge-coupled devices, 26.3–26.5,
26.4*f*–26.6*f*
complementary metal-oxide semiconductor
(CMOS), 26.8–26.9
modulation transfer function (MTF)
for, 26.14
performance metrics for, 26.12–26.16, 26.16*f*
sampling with, 26.16–26.19, 26.17*f*–26.19*f*
- Space curves, differential geometry of,
1.18–1.19
- Space-bandwidth product, 6.9
- Space-integrating correlator, 11.11
- Spar cutting, Ahrens method of, 13.12
- Sparrow criterion, 28.18, 28.19
- Spatial coherence, 2.38–2.40, 2.38*f*–2.39*f*,
5.3, 5.5
- Spatial filtering, 11.5–11.6, 11.6*f*
- Spatial multiplexing, 23.8
- Spatial-frequency content, 4.8
- Specimen, for microscopes, 28.22, 28.54–28.55
- Speckle:
and coherence theory, 5.22
laser, 9.14
objective and subjective, 33.9
single, 8.17
in volume scattering, 9.15–9.16, 9.15*f*
- Spectral coherence, 5.5
- Spectral transitions, 10.6–10.7
- Spectrographs, charge, 34.8
- Spectrometers, 31.1–31.31
Bunsen-Kirchhoff, 20.5*f*
dispersive prisms and gratings for, 20.2–20.3,
20.3*f*
high-resolution techniques of, 31.23–31.29
fluorescence line narrowing, 31.29, 31.29*f*
laser Stark spectroscopy of molecules,
31.27–31.29, 31.27*f*, 31.28*f*
polarized absorption spectrometers, 31.26
saturated absorption, 31.24–31.26,
31.24*f*–31.26*f*

- Spectrometers (*Cont.*):
 and light scattering, 31.30–31.31, 31.31f
 luminescence, 31.5–31.12, 31.8f, 31.11f
 optical absorption, 31.2–31.5, 31.5f
 and photoluminescence decay time,
 31.12–31.15, 31.13f, 31.14f
 polarization, 31.15–31.23
 and optically detected magnetic resonance,
 31.21–31.23, 31.22f, 31.23f
 polarized absorption by, 31.15–31.17,
 31.15f, 31.18f
 polarized absorption/luminescence
 techniques for, 31.17,
 31.19–31.21, 31.19f
 prism, 20.2–20.3, 20.3f
 Unicam double-monochromator, 20.10,
 20.13, 20.15f
- Spectrophotometric measurements, 12.15
- Spectropolarimetry and spectropolarimeters,
 15.6, 15.8
- Spectroradiometers, 20.1, 20.2f, 20.14t
- Spectroscopic ellipsometry (SE), 16.3
- Spectroscopic lineshapes, 10.6–10.7, 10.22–10.27
 in solid state spectroscopy, 10.22–10.26,
 10.23f–10.27f
 of spectral transitions, 10.6–10.7
- Spectroscopic transition, rates of, 10.4–10.6
- Spectroscopy, 3.29, 10.1–10.22
 of multielectron atoms, 10.10–10.11
 of one-electron atoms, 10.7–10.9, 10.8f, 10.9f
 and outer electronic structure, 10.12–10.16,
 10.13f–10.15f
 photon correlation, 9.8
 rates of spectroscopic transition, 10.4–10.6
 solid state, 10.22–10.26
 theoretical basis, 10.3–10.4
 of tri-positive rare earth ions, 10.16–10.18,
 10.16t, 10.17f
 and vibrational and rotational spectra,
 10.18–10.22, 10.19f, 10.21f
- Spectrum:
 power, 8.12–8.13
 of primary light source, 5.6
 rotational, 10.20–10.22
 secondary, 29.7, 29.38
 vibrational, 10.18–10.20
 and wavefront division, 2.17–2.18, 2.18f
- Spectrum of light, 5.19–5.22
 coherence functions for, 5.5–5.6
 coherent mode representation of, 5.20–5.21
- Spectrum of light (*Cont.*):
 limitations, 5.19, 5.20f
 for primary sources, 5.6
 scaling law, 5.21
 Wolf shift, 5.21
- Spheres, scattering by, 7.11–7.14
- Spherical aberrations, 1.90, 29.7, 29.38
 oblique, 29.15, 29.21, 29.37
 zonal, 29.8, 29.38
- Spherical lenses, 6.3
- Spherical primaries, in objective designs, 29.9,
 29.11, 29.18, 29.30
- Spherical surfaces, in systems of revolution, 1.34
- Spherical waves, 2.4, 2.5f, 3.2–3.3
 interference from, 2.11–2.12, 2.12f, 2.13f
 and plane waves, 2.9–2.11, 2.10f
- Spherochromatism, 24.3–24.6, 29.14
- Split-aperture scanners, 30.15–30.16, 30.15f
- Spontaneous decay rate, 10.6
- Square-ended Nicol prisms, 13.16f, 13.17
- Squirm, 13.14
- Stark effect, 21.11–21.12
- Stark spectroscopy of molecules, laser,
 31.27–31.29
- Stationary phase approximation, in diffraction,
 3.29, 3.31–3.32
- Stationary surfaces, Fresnel-Kirchhoff
 approximation for, 8.6–8.8
- Statistical radiometry, 5.22
- Steeg and Reuter Nicol prisms, 13.17
- Stellar interferometers, 32.19–32.21, 32.19f, 32.20f
- Stereo cameras, 25.23–25.24, 25.23f
- Stigmatic imaging, 1.29–1.31, 1.30f
- Stimulated emission depletion (STED), 28.24
- Stochastic optical reconstruction microscopy
 (STORM), 28.23
- Stokes matrix, 14.4 (*See also* Mueller matrix)
- Stokes parameters:
 and Mueller matrices, 14.4–14.6
 and Poincaré sphere, 14.4–14.6, 14.5f,
 14.26f, 14.33
- Stokes polarimeters, 15.4, 15.5, 15.25
- Stokes vectors, 12.14n, 12.28
 and Mueller matrices, 14.15–14.17,
 14.19–14.21
 for nonhomogeneous polarization
 elements, 14.26f
 in polarimetry, 15.3, 15.8–15.13
 for radiative transfer, 9.11
 for speckle patterns, 9.17

- Stop shift, 1.92
- Stops:
 - aperture, 1.74, 1.75*f*, 17.8, 29.5, 29.36
 - field, 1.74, 17.9, 29.5, 29.37
 - of lenses, 17.8–17.9
- S-trace formula, 1.44
- Straight edges, cylindrical wavefronts and, 3.14–3.16, 3.14*f*, 3.15*f*
- Stratified-medium model (SMM),
 - in ellipsometry, 16.4
- Stray light suppression, 29.5, 29.5*f*
- Streak cameras, 25.24, 25.24*f*
- Strehl index, 8.6
- Sub-Doppler absorption spectroscopy, 31.26–31.27
- Subjective speckle, 33.9
- Superzone construction, for micro-Fresnel lenses, 22.36
- Surface acoustic wave (SAW) devices, 21.16
- Surface scattering, 8.1–8.17
 - finish models for, 8.14–8.15
 - and finite illumination area, 8.12
 - Fresnel-Kirchhoff approximation for, 8.5–8.9
 - fractal surfaces, 8.8–8.9
 - statistically stationary surfaces, 8.6–8.8
 - notation for, 8.2–8.4
 - and power spectra, 8.12–8.13, 8.13*f*
 - Rayleigh-Rice approximation, 8.9–8.12
 - second-order statistical functions for, 8.12–8.13
 - statistics for, 8.12–8.15
 - and surface finish specifications, 8.15–8.17
- Surfaces, 1.32
 - aspherical, 1.35
 - conical, 1.34–1.35
 - spherical, 1.34
- Suspension systems, of resonant scanners, 30.43, 30.43*f*, 30.44
- Sweatt lenses, 23.4
- Sweatt model, 23.4, 23.6
- Switches and switching:
 - in integrated optics, 21.34–21.35, 21.34*f*
 - interferometric optical switches, 32.19
- Symmetrical lenses, 1.71, 17.26–17.27, 17.27*f*
- Synthetic aperture radar data, 11.6–11.8, 11.7*f*–11.8*f*
- System-response (SR) function, 8.12
- Tangent vectors, of space curves, 1.18, 1.19
- Tangential fans and foci, 1.35
- Telecentric lenses, 18.12
- Telecentric stop, 17.9
- Telecentricity, 1.83–1.84, 30.31, 30.31*f*, 30.32
- Telephoto lenses, 17.29, 27.2, 27.6, 27.7*f*–27.16*f*, 27.13
- Telescopes:
 - astronomical, 18.10
 - Cassegrainian, 18.21
 - field of view in, 18.15–18.16, 18.15*f*
 - Galilean, 18.15
 - Keplerian, 18.10
 - Mersenne, 18.19
 - objectives for, 29.8, 29.9, 29.12, 29.29–29.30
 - terrestrial, 18.10–18.11, 18.11*f*
- Telescopic lenses, 1.46 (*See also* Afocal lenses)
- Telescopic transformations, 1.57
- Television:
 - integrated optics and cable, 21.2, 21.32–21.34
 - Scophony TV projection system, 30.45
- Temperature:
 - and crystalline-quartz retardation plates, 13.48
 - Curie, 35.25
 - and integrated optics, 21.12
- Temporal beats, interference and, 2.13
- Temporal coherence, 2.41, 5.3
- Temporal signals, analog processing of, 11.8–11.12, 11.9*f*–11.11*f*
- Tensor, electro-optic, 21.10
- Tensor product, of Jones matrices, 14.23
- Terrestrial telescopes, 18.10–18.11, 18.11*f*
- Tessar lenses, 17.26, 17.26*f*
- Tewarson, 12.30
- Thermal imaging cameras, 25.25
- Thermal optimization, of media, 35.20, 35.20*f*
- Thermal (Lambertian) sources of light, 5.12–5.13
- Thermomagnetic recording process, 35.17–35.20, 35.18*f*–35.20*f*
- Thick lens systems, 1.55
- Thin film oxides, 21.13–21.14
- Thin films, interference and, 2.24
- Thin lens systems, 1.55
- Thin-lens model:
 - of Galilean afocal lenses, 18.15, 18.15*f*
 - of Keplerian afocal lenses, 18.7–18.8, 18.8*f*
- Third-order aberrations, 1.90–1.91, 29.38
- Thompson reversed Nicol prisms, 13.16*f*, 13.17
- Thomson CSE, 21.35
- Three-beam interferometers, 32.7–32.8, 32.7*f*
- Three-dimensional diffraction patterns, 28.19–28.22, 28.20*f*, 28.21*f*

- Three-lens prime focus corrector objective, **29.10**
- Three-mirror objectives, **29.28–29.32, 29.34**
- Three-phase model (ellipsometry), **16.5–16.8, 16.6f–16.8f**
- Three-powered-mirror lenses, **18.20, 18.20f, 18.21**
- Three-ray rule, **1.42**
- Throughput, **1.22, 13.7** (*See also Étendue*)
- Tilted planes, for collineation, **1.61, 1.62, 1.62f**
- Tilted-plane processors, **11.8, 11.8f**
- Time, coherence, **5.3**
- Time averages, in coherence theory, **6.2n, 6.4–6.5**
- Time domain, **7.17**
- Time lag, of cameras, **25.8–25.9, 25.9f**
- Time-integrating correlators, **11.11–11.12, 11.11f**
- Time-sequential polarimeters, **15.5**
- T*-matrix method, **7.15**
- T*-number, **1.79**
- Toner, in xerographic systems, **34.1, 34.8–34.9, 34.8f**
- Topothesis, of fractals, **8.8**
- Torsion, radius of, **1.19**
- Total internal reflection:
 - of optical waveguides, **21.3**
 - of Rochon prisms, **13.20**
- Track-error signal (TES), **35.14–35.15, 35.15f, 35.17**
- Tracks, on optical disks, **35.2–35.5, 35.3f–35.5f**
- Transfer, in xerographic systems, **34.10**
- Transfer functions (*see specific functions, e.g.:*
 - Modulation transfer function (MTF))
- Transfer matrices, **1.66**
- Transillumination, in microscopes, **28.5–28.7, 28.6f**
- Translation, by scanners, **30.5**
- Transmission:
 - amplitude-shift-keyed, **21.30**
 - analog, **21.32–21.34, 21.33f, 21.34f**
 - differential-phase-shift-keyed, **21.30, 21.32**
 - digital, **21.31–21.32**
 - by Glan-Thompson-type prisms, **13.9–13.10, 13.10f, 13.11f**
 - by LEDs, **21.32**
 - in multilayer systems, **16.8–16.9**
 - nonnormal-incidence, **12.18–12.24**
 - silicon photonics, **21.14–21.16, 21.15f**
 - by silicon photonics, **21.38, 21.39**
- Transmission ellipsometry, **16.10**
- Transmissive planar objects, **6.3**
- Transmittance:
 - of Brewster angle transmission polarizers, **12.22**
 - of Mueller matrices, **14.16–14.17**
 - of pile-of-plates polarizers, **12.15–12.17**
 - principal, **12.14–12.16**
 - of spherical lenses, **6.3**
- Transmitted state, of polarizers, **15.19**
- Transverse electric (TE) modes, of optical waveguides, **21.6, 21.7**
- Transverse magnetic (TM) modes, of optical waveguides, **21.6, 21.7**
- Transverse magnification, **1.28, 1.50–1.51**
- Transverse primary chromatic aberration (TPAC), **17.22**
- Transverse ray aberrations, **1.87**
- Transverse translational scan, **30.28**
- Traveling wave modulators, **21.26**
- Trigger detection, with ODMR, **31.21**
- Tri-Level highlight color process, **34.12, 34.13f**
- Trim retarders, for LC panels, **15.33**
- Trimmed Nicol-type polarizers, **13.16–13.18, 13.16f**
- Triplet lenses, **17.26, 17.26f**
- Tri-positive rare earth ions, **10.16–10.18, 10.16t, 10.17f**
- Trischiefspiegler objective, **29.27**
- T*-trace formula, **1.44**
- Tube length (objective lenses), **28.13, 28.13t**
- Tutton's test, **13.45**
- Two-component magnetic brush
 - development, in xerographic systems, **34.5–34.7, 34.5f–34.7f**
- Two-dimensional images, analog processing
 - for, **11.12–11.17, 11.13f**
- Two-dimensional scanners, **30.18, 30.19f**
- Two-lens systems, **17.20–17.22, 17.21f–17.22f**
- Two-mirror, three reflection objective, **29.26**
- Two-phase model (ellipsometry), **16.5, 16.6f**
- Two-powered-mirror lenses, **18.20, 18.20f**
- Two-ray paraxial invariant, **1.41**
- Two-wavelength interferometry, **32.9**
- Twyman-Green interferometers, **2.28, 32.2, 32.9, 33.5**
- Unblazed gratings, **20.3–20.4**
- Underillumination, for input/output scanning, **30.14**
- Underwater cameras, **25.25**
- Unfolded reflections, in systems
 - of revolution, **1.32**

- Unicam double-monochromator spectrometer, 20.10, 20.13, 20.15f
- Uniformity, of binary optics, 23.7
- Unit determinants, paraxial matrix methods for, 1.67
- Upatnieks, J., 33.2
- Van Cittert-Zernike theorem, 2.38, 2.39, 5.13–5.14, 5.17–5.19, 6.4, 6.8
- Vander Lugt filters, 11.13–11.14
- Vapor-phase epitaxy (VPE), 21.17, 21.18
- Variable retardation plates, 13.53
- Variable-angle spectroscopic ellipsometry (VASE), 16.3
- Variational integral, of rays, 1.19
- Vector aberration theory, 29.5
- Vector diffraction, 3.32–3.37, 3.32f–3.37f, 23.13–23.14, 23.14f
- Vector Huygens secondary source (unit), 3.33–3.36, 3.37f
- Vectors:
 - electric field, 2.3–2.4
 - of space curves, 1.18, 1.19
- Vector-scattering amplitude, 7.8
- VeriMask, 33.19
- Vertex, for figure of revolution, 1.32
- Vertex curvature, 1.32–1.33
- Vertical half-wave linear retarders, Mueller matrices for, 14.12t
- Vertical linear polarizers, 14.10t
- Vertical quarter-wave linear retarders, Mueller matrices for, 14.12t
- Vibration, planes of, 12.6
- Vibrational relaxation, 31.14, 31.14f
- Vibrational sidebands, 10.24
- Vibrational spectra, 10.18–10.20
- Video cameras, 25.7–25.8
- Videocassette recorders (VCRs), 35.1
- Video-enhanced differential-interference contrast (VE-DIC), 28.41
- View cameras, 25.18–25.20, 25.19f
- Vignetting, 1.81, 1.81f, 1.82, 17.8, 29.6, 29.38
- Virtual images, 29.38
- Virtual pupils, 1.76
- Virtual rays, 1.10
- Visual magnification, 1.28
- VLSI-CMOS technology, silicon photonics for, 21.15, 21.16
- Volume imaging ideal, 1.29
- Volume scattering, 9.1–9.17
 - multiple scattering, 9.8–9.17
 - analytical theory of, 9.9–9.10, 9.9f
 - depolarization, 9.16–9.17, 9.17f
 - effective-medium representation, 9.8
 - radiative transfer, 9.10–9.13, 9.11f, 9.13f
 - speckle patterns, 9.15–9.16, 9.15f
 - weak localization, 9.13–9.17, 9.14f
 - single particle scattering vs., 9.2–9.3
 - and single scattering, 9.4–9.8, 9.6f, 9.7f
 - theory of, 9.3–9.4, 9.4f
- Volumetric reflection, 7.7f
- Volumetric scattering cross section, 7.7f
- Wadsworth configuration, 20.5f, 20.8, 20.12f, 20.14t
- WALRUS objective, 29.28
- Wave equation, 12.4–12.6
- Wave normals, 13.5
- Wavefront aberration coefficients, 1.90
- Wavefront aberrations, 1.86–1.88
- Wavefront division, 2.14
- Wavefront division, interference by, 2.14–2.19, 2.15f–2.18f
- Wavefront multiplexers, 23.11, 23.11t, 23.12
- Wavefront quality, of binary optics, 23.8
- Wavefronts, 3.4
 - aberrated, 2.12, 2.13
 - cylindrical, 3.13–3.21, 3.14f
 - Cornu's spiral, 3.16–3.19
 - opaque strip construction, 3.20–3.21
 - from rectangular apertures, 3.19–3.20
 - from straight edge, 3.14–3.16
 - disturbance of, 3.5–3.6
 - for cylindrical wavefronts, 3.13–3.14, 3.14f
 - and straight edges, 3.14–3.15
 - geometrical, 1.12–1.13
 - and interference, 2.4–2.5
- Waveguide grating routers (WGRs), 21.24
- Waveguides:
 - integrated optics, 21.3–21.8, 21.3f–21.5f, 21.7f
 - leaky, 21.3
- Wavelength:
 - and modulation transfer function, 17.38, 17.39
 - of plane waves, 2.4
- Wavelength division multiplexes (WDM) systems:
 - fabrication of, 21.14
 - filters for, 21.23
 - in integrated optics, 21.37–21.38, 21.37f–21.39f

- Wavelength interval (between oscillations), 12.10
- Waveplates, 15.8
- Waves, 2.3–2.6
- amplitudes of, 2.4, 2.5, 12.5
 - diffraction of, 3.2–3.3, 3.3f
 - interference of, 2.5–2.6
 - plane, 2.4, 2.5f, 3.3, 3.17
 - decomposition of, 3.23
 - interference of, 2.8–2.9, 2.9f
 - and spherical waves, 2.9–2.11, 2.10f
 - spherical, 2.4, 2.5f, 3.2–3.3
 - interference of, 2.11–2.12, 2.12f, 2.13f
 - and plane waves, 2.9–2.11, 2.10f
- Weak polarization elements (Mueller matrices), 14.26–14.27
- Wernicke prisms, 20.6f
- Wetherell and Womble objectives, 29.31
- Weyl's integral, 5.17
- Weyl's plane-wave decomposition, 3.23
- Wide-angle photography:
- cameras for, 25.25, 25.26f
 - lenses for, 27.2, 27.5f
 - with nonrectilinear distortion, 27.6
 - with rectilinear distortion correction, 27.6, 27.13f, 27.14f
- Wide-field objective with Maksutov correction, 29.20
- Wiemann-Hansch experiment, 31.26
- Wiener filter, 11.15–11.17
- Wiener-Khinchine theorem, 5.5, 5.21
- Wien's wavelength displacement law, 31.4
- Wigner distribution function, 5.8
- Winchester heads, 35.6
- Wire grids, 13.30n
- Wire-grid polarizers, 13.30–13.33, 13.31f, 13.32t
- Wiscombe, W., 7.12, 7.15
- Wizinowich, P. L., 12.15, 13.52
- Wolf shift, 5.21
- Wollaston, W. H., 17.17, 17.18
- Wollaston prisms, 13.7, 13.18, 13.18f, 13.21, 13.24
 - for DIC microscopes, 28.39, 28.40
 - in Nomarski interferometers, 32.4
 - and nondispersive prisms, 19.3t, 19.15, 19.15f
- Working distance, of objective lenses, 28.13
- Wright objective, 29.15
- Wright objectives, 29.15
- Write-once-read-many (WORM) technology, 35.2
- Xerographic systems, 34.1–34.13, 34.2f
 - cleaning and erasing in, 34.10
 - color in, 34.11–34.12, 34.11f–34.13f
 - control of, 34.11
 - development in, 34.5–34.10, 34.5f–34.9f
 - fusing in, 34.10
 - and latent image, 34.1–34.4, 34.2f–34.3f
 - transfer in, 34.10
- X-ray scattering, 9.6
- Yolo objectives, 29.26
- Young's astigmatic formulae, 1.44
- Young's double slit experiment, 2.14–2.15, 2.15f
- Young's fringes, 13.45
- Young's modulus, for molded microlenses, 22.12t
- Young's two pinhole interferometer, 6.3
- Young-Thollon half prisms, 20.7f
- Zeeman effect, 31.17, 31.18f, 31.19–31.21
- Zeiss Infinity Color-Corrected Systems Plan Apo, 28.15, 28.15f
- Zeiss prism system, 19.3t, 19.16, 19.16f
- Zeiss sheet polarizers, 13.26
- Zenger prisms, 20.6f
- Zernike polynomials, 1.90, 23.3
- Zero-phonon transitions, 10.24–10.26, 10.24f
- Zeros, in paraxial matrices, 1.68
- Zonal spherical aberrations, 29.8, 29.38
- Zone plates, 3.11–3.13, 3.12f
- Zoom lenses, 27.17, 27.20f–27.22f

COLOR PLATES

DO NOT DUPLICATE

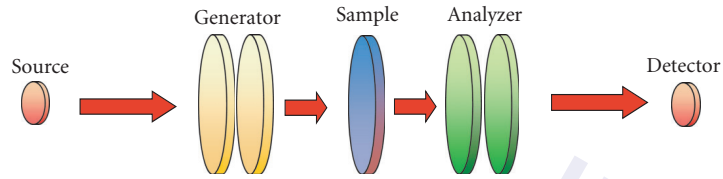


FIGURE 15.2 A sample-measuring polarimeter consists of a source, polarization state generator (PSG), the sample, a polarization state analyzer (PSA), and the detector.

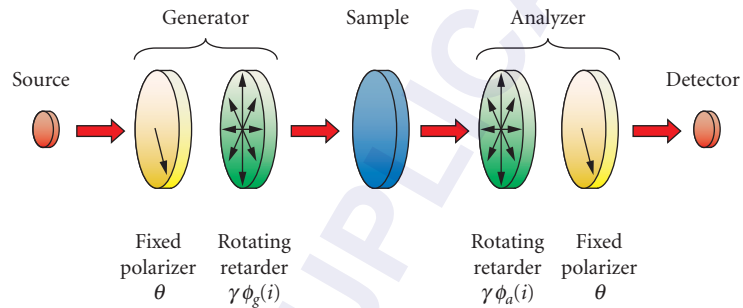


FIGURE 15.3 The dual rotating retarder polarimeter consists of a source, a fixed linear polarizer, a retarder which rotates in steps, the sample, a second retarder which rotates in steps, a fixed linear polarizer, and the detector.

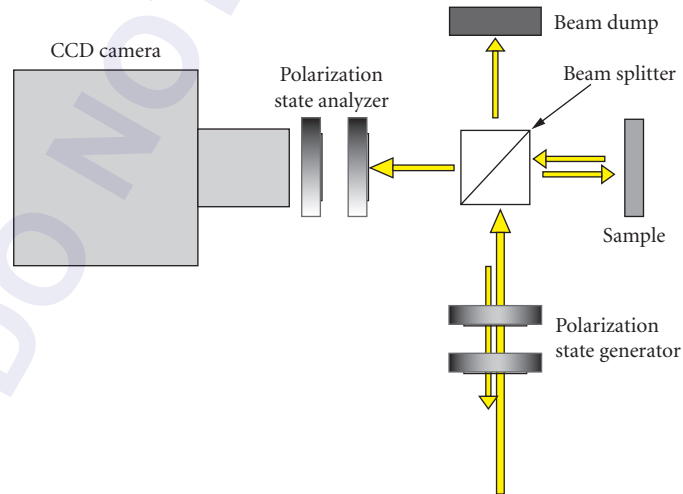


FIGURE 15.6 Imaging polarimeter configured for retro reflection testing using a non polarizing beam splitter and beam dump.

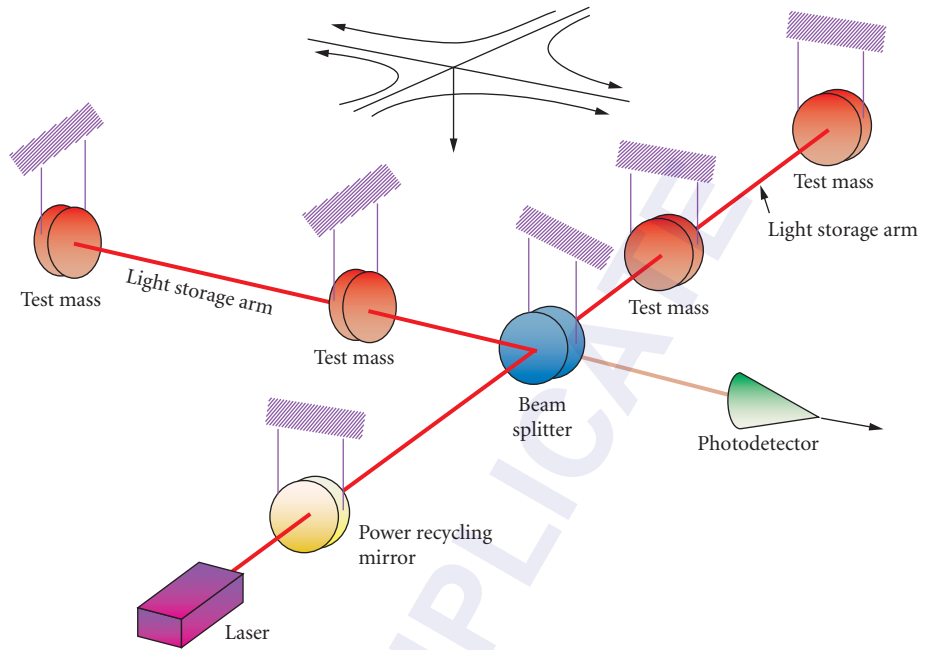


FIGURE 32.26 Gravitational-wave detector using two Fabry-Perot interferometers.

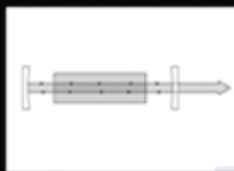
Third Edition

Sponsored by the Optical Society of America

HANDBOOK OF OPTICS

Volume II

*Design, Fabrication, and Testing; Sources
and Detectors; Radiometry and Photometry*



Editor-in-Chief:
Michael Bass

Associate Editors:
Casimer M. DeCusatis
Jay M. Enoch
Vasudevan Lakshminarayanan
Guifang Li
Carolyn MacDonald
Virendra N. Mahajan
Eric Van Stryland

OSA[®]

HANDBOOK OF OPTICS

DO NOT DUPLICATE

ABOUT THE EDITORS

Editor-in-Chief: Dr. Michael Bass is professor emeritus at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Associate Editors:

Dr. Casimer M. DeCusatis is a distinguished engineer and technical executive with IBM Corporation.

Dr. Jay M. Enoch is dean emeritus and professor at the School of Optometry at the University of California, Berkeley.

Dr. Vasudevan Lakshminarayanan is professor of Optometry, Physics, and Electrical Engineering at the University of Waterloo, Ontario, Canada.

Dr. Guifang Li is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Dr. Carolyn MacDonald is a professor at the University at Albany, and director of the Center for X-Ray Optics.

Dr. Virendra N. Mahajan is a distinguished scientist at The Aerospace Corporation.

Dr. Eric Van Stryland is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

HANDBOOK OF OPTICS

Volume II
Design, Fabrication, and Testing;
Sources and Detectors;
Radiometry and Photometry

THIRD EDITION

Sponsored by the
OPTICAL SOCIETY OF AMERICA

Michael Bass Editor-in-Chief
*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

Virendra N. Mahajan Associate Editor
*The Aerospace Corporation
El Segundo, California*

Eric Van Stryland Associate Editor
*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*



New York Chicago San Francisco Lisbon London Madrid
Mexico City Milan New Delhi San Juan Seoul
Singapore Sydney Toronto

This page intentionally left blank.

DO NOT DUPLICATE

COVER ILLUSTRATIONS

Left: Telescope such as used by Galileo to discover lunar craters and Jupiter's moons. The basic design is still used in opera and sports glasses. See Chap. 1.

Middle: Simplified schematic of a laser showing the gain medium which amplifies the light, and the resonator which defines the light's direction and spatial distribution. The third critical part, the source to excite the gain medium, is not shown. See Chap. 16.

Right: Zernike circle polynomial representing balanced astigmatism with a standard deviation of one wave illustrated as an isometric plot on the top, interferogram on the left, and point-spread function on the right. See Chap. 11.

This page intentionally left blank.

DO NOT DUPLICATE

CONTENTS

Contributors	xvii
Brief Contents of All Volumes	xix
Editors' Preface	xxv
Preface to Volume II	xxvii
Glossary and Fundamental Constants	xxix

Part 1. Design

Chapter 1. Techniques of First-Order Layout	<i>Warren J. Smith</i>	1.3
<hr/>		
1.1	Glossary / 1.3	
1.2	First-Order Layout / 1.4	
1.3	Ray-Tracing / 1.4	
1.4	Two-Component Systems / 1.5	
1.5	Afocal Systems / 1.7	
1.6	Magnifiers and Microscopes / 1.8	
1.7	Afocal Attachments / 1.8	
1.8	Field Lenses / 1.8	
1.9	Condensers / 1.10	
1.10	Zoom or Varifocal Systems / 1.11	
1.11	Additional Rays / 1.12	
1.12	Minimizing Component Power / 1.13	
1.13	Is It a Reasonable Layout? / 1.13	
1.14	Achromatism / 1.14	
1.15	Athermalization / 1.15	
Chapter 2. Aberration Curves in Lens Design	<i>Donald C. O'Shea and Michael E. Harrigan</i>	2.1
<hr/>		
2.1	Glossary / 2.1	
2.2	Introduction / 2.1	
2.3	Transverse Ray Plots / 2.2	
2.4	Field Plots / 2.4	
2.5	Additional Considerations / 2.5	
2.6	Summary / 2.6	
2.7	References / 2.6	
Chapter 3. Optical Design Software	<i>Douglas C. Sinclair</i>	3.1
<hr/>		
3.1	Glossary / 3.1	
3.2	Introduction / 3.2	
3.3	Lens Entry / 3.2	
3.4	Evaluation / 3.8	
3.5	Optimization / 3.16	
3.6	Other Topics / 3.21	
3.7	Buying Optical Design Software / 3.22	
3.8	Summary / 3.24	
3.9	References / 3.24	

Chapter 4. Optical Specifications *Robert R. Shannon* 4.1

- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.1
- 4.3 Preparation of Optical Specifications / 4.5
- 4.4 Image Specifications / 4.6
- 4.5 Element Description / 4.8
- 4.6 Environmental Specifications / 4.10
- 4.7 Presentation of Specifications / 4.10
- 4.8 Problems with Specification Writing / 4.11
- 4.9 References / 4.12

Chapter 5. Tolerancing Techniques *Robert R. Shannon* 5.1

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.1
- 5.3 Wavefront Tolerances / 5.3
- 5.4 Other Tolerances / 5.7
- 5.5 Starting Points / 5.8
- 5.6 Material Properties / 5.9
- 5.7 Tolerancing Procedures / 5.9
- 5.8 Problems in Tolerancing / 5.11
- 5.9 References / 5.11

Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.* 6.1

- 6.1 Glossary / 6.1
- 6.2 Introduction and Summary / 6.1
- 6.3 Mounting Individual Rotationally Symmetric Optics / 6.2
- 6.4 Multicomponent Lens Assemblies / 6.5
- 6.5 Mounting Windows and Domes / 6.11
- 6.6 Mounting Small Mirrors and Prisms / 6.11
- 6.7 Mounting Moderate-Sized Mirrors / 6.17
- 6.8 Contact Stresses in Optics / 6.21
- 6.9 Temperature Effects on Mounted Optics / 6.21
- 6.10 References / 6.25

Chapter 7. Control of Stray Light *Robert P. Breault* 7.1

- 7.1 Glossary / 7.1
- 7.2 Introduction / 7.1
- 7.3 Concepts / 7.2
- 7.4 Optical Software for Stray Light Analysis / 7.24
- 7.5 Methods / 7.27
- 7.6 Conclusion / 7.30
- 7.7 Sources of Information on Stray Light and Scattered Light / 7.31
- 7.8 References / 7.32

Chapter 8. Thermal Compensation Techniques
Philip J. Rogers and Michael Roberts 8.1

- 8.1 Glossary / 8.1
- 8.2 Introduction / 8.2
- 8.3 Homogeneous Thermal Effects / 8.2
- 8.4 Tolerable Homogeneous Temperature Change (No Compensation) / 8.5
- 8.5 Effect of Thermal Gradients / 8.6
- 8.6 Intrinsic Athermalization / 8.7
- 8.7 Mechanical Athermalization / 8.8
- 8.8 Optical Athermalization / 8.12
- 8.9 References / 8.15

Part 2. Fabrication

Chapter 9. Optical Fabrication *Michael P. Mandina* 9.3

- 9.1 Introduction / 9.3
- 9.2 Material Forms of Supply / 9.3
- 9.3 Basic Steps in Spherical Optics Fabrication / 9.4
- 9.4 Plano Optics Fabrication / 9.7
- 9.5 Asphere Optics Fabrication / 9.7
- 9.6 Crystalline Optics / 9.8
- 9.7 Purchasing Optics / 9.9
- 9.8 Conclusion / 9.9
- 9.9 References / 9.9

Chapter 10. Fabrication of Optics by Diamond Turning *Richard L. Rhorer and Chris J. Evans* 10.1

- 10.1 Glossary / 10.1
- 10.2 Introduction / 10.1
- 10.3 The Diamond-Turning Process / 10.2
- 10.4 The Advantages of Diamond Turning / 10.2
- 10.5 Diamond-Turnable Materials / 10.4
- 10.6 Comparison of Diamond Turning and Traditional Optical Fabrication / 10.6
- 10.7 Machine Tools for Diamond Turning / 10.6
- 10.8 Basic Steps in Diamond Turning / 10.8
- 10.9 Surface Finish of Diamond-Turned Optics / 10.9
- 10.10 Metrology of Diamond-Turned Optics / 10.12
- 10.11 Conclusions / 10.13
- 10.12 References / 10.14

Part 3. Testing

Chapter 11. Orthonormal Polynomials in Wavefront Analysis *Virendra N. Mahajan* 11.3

- Abstract / 11.3
- 11.1 Glossary / 11.3
- 11.2 Introduction / 11.4
- 11.3 Orthonormal Polynomials / 11.5
- 11.4 Zernike Circle Polynomials / 11.6
- 11.5 Zernike Annular Polynomials / 11.13
- 11.6 Hexagonal Polynomials / 11.21
- 11.7 Elliptical Polynomials / 11.21
- 11.8 Rectangular Polynomials / 11.27
- 11.9 Square Polynomials / 11.30
- 11.10 Slit Polynomials / 11.30
- 11.11 Aberration Balancing and Tolerancing, and Diffraction Focus / 11.30
- 11.12 Isometric, Interferometric, and PSF Plots for Orthonormal Aberrations / 11.36
- 11.13 Use of Circle Polynomials for Noncircular Pupils / 11.37
- 11.14 Discussion and Conclusions / 11.39
- 11.15 References / 11.40

Chapter 12. Optical Metrology *Zacarias Malacara and Daniel Malacara-Hernández* 12.1

- 12.1 Glossary / 12.1
- 12.2 Introduction and Definitions / 12.2

- 12.3 Length and Straightness Measurements / 12.2
- 12.4 Angle Measurements / 12.10
- 12.5 Curvature and Focal Length Measurements / 12.17
- 12.6 References / 12.25

Chapter 13. Optical Testing *Daniel Malacara-Hernández* **13.1**

- 13.1 Glossary / 13.1
- 13.2 Introduction / 13.1
- 13.3 Classical Noninterferometric Tests / 13.1
- 13.4 Interferometric Tests / 13.7
- 13.5 Increasing the Sensitivity of Interferometers / 13.13
- 13.6 Interferogram Evaluation / 13.14
- 13.7 Phase-Shifting Interferometry / 13.18
- 13.8 Measuring Aspherical Wavefronts / 13.23
- 13.9 References / 13.28

Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant* **14.1**

- 14.1 Glossary / 14.1
- 14.2 Introduction / 14.1
- 14.3 Plotting CGHs / 14.3
- 14.4 Interferometers Using Computer-Generated Holograms / 14.4
- 14.5 Accuracy Limitations / 14.6
- 14.6 Experimental Results / 14.7
- 14.7 Discussion / 14.9
- 14.8 References / 14.9

Part 4. Sources

Chapter 15. Artificial Sources *Anthony LaRocca* **15.3**

- 15.1 Glossary / 15.3
- 15.2 Introduction / 15.3
- 15.3 Radiation Law / 15.4
- 15.4 Laboratory Sources / 15.7
- 15.5 Commercial Sources / 15.13
- 15.6 References / 15.53

Chapter 16. Lasers *William T. Silfvast* **16.1**

- 16.1 Glossary / 16.1
- 16.2 Introduction / 16.2
- 16.3 Laser Properties Associated with the Laser Gain Medium / 16.4
- 16.4 Laser Properties Associated with Optical Cavities or Resonators / 16.19
- 16.5 Special Laser Cavities / 16.25
- 16.6 Specific Types of Lasers / 16.29
- 16.7 References / 16.37

Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman* **17.1**

- 17.1 Glossary / 17.1
- 17.2 Introduction / 17.2
- 17.3 Light-Generation Processes / 17.2
- 17.4 Light Extraction / 17.6
- 17.5 Device Structures / 17.8

- 17.6 Material Systems / 17.15
- 17.7 Substrate Technology / 17.20
- 17.8 Epitaxial Technology / 17.21
- 17.9 Wafer Processing / 17.23
- 17.10 Led Quality and Reliability / 17.25
- 17.11 Led-Based Products / 17.29
- 17.12 References / 17.35

Chapter 18. High-Brightness Visible LEDs **18.1**
Winston V. Schoenfeld

- 18.1 The Materials Systems / 18.1
- 18.2 Substrates and Epitaxial Growth / 18.2
- 18.3 Processing / 18.3
- 18.4 Solid-State Lighting / 18.4
- 18.5 Packaging / 18.5

Chapter 19. Semiconductor Lasers **19.1**
*Pamela L. Derry,
Luis Figueroa, and Chi-Shain Hong*

- 19.1 Glossary / 19.1
- 19.2 Introduction / 19.3
- 19.3 Applications for Semiconductor Lasers / 19.3
- 19.4 Basic Operation / 19.4
- 19.5 Fabrication and Configurations / 19.6
- 19.6 Quantum Well Lasers / 19.9
- 19.7 High-Power Semiconductor Lasers / 19.18
- 19.8 High-Speed Modulation / 19.30
- 19.9 Spectral Properties / 19.36
- 19.10 Surface-Emitting Lasers / 19.39
- 19.11 Conclusion / 19.41
- 19.12 References / 19.43

Chapter 20. Ultrashort Optical Sources and Applications **20.1**
Jean-Claude Diels and Ladan Arissian

- 20.1 Introduction / 20.1
- 20.2 Description of Optical Pulses and Pulse Trains / 20.2
- 20.3 Pulse Evolution toward Steady State / 20.9
- 20.4 Coupling Circulating Pulses Inside a Cavity / 20.12
- 20.5 Designs of Cavities with Two Circulating Pulses / 20.15
- 20.6 Analogy of a Two-Level System / 20.22
- 20.7 Conclusion / 20.28
- 20.8 References / 20.28

Chapter 21. Attosecond Optics **21.1**
Zenghu Chang

- 21.1 Glossary / 21.1
- 21.2 Introduction / 21.2
- 21.3 The Driving Laser / 21.4
- 21.4 Attosecond Pulse Generation / 21.6
- 21.5 Attosecond Pulse Characterization / 21.8
- 21.6 Acknowledgments / 21.10
- 21.7 References / 21.10

Chapter 22. Laser Stabilization **22.1**
*John L. Hall,
Matthew S. Taubman, and Jun Ye*

- 22.1 Introduction and Overview / 22.1
- 22.2 Servo Principles and Issues / 22.5

- 22.3 Practical Issues / 22.12
- 22.4 Summary and Outlook / 22.23
- 22.5 Conclusions and Recommendations / 22.24
- 22.6 Acknowledgments / 22.24
- 22.7 References / 22.24

Chapter 23. Quantum Theory of the Laser *János A. Bergou, Berthold-Georg Englert, Melvin Lax, Marian O. Scully, Herbert Walther, and M. Suhail Zubairy* 23.1

- 23.1 Glossary / 23.1
- 23.2 Introduction / 23.5
- 23.3 Some History of the Photon Concept / 23.6
- 23.4 Quantum Theory of the Laser / 23.14
- 23.5 The Laser Phase-Transition Analogy / 23.35
- 23.6 Exotic Masers and Lasers / 23.40
- 23.7 Acknowledgments / 23.45
- 23.8 References / 23.46

Part 5. Detectors

Chapter 24. Photodetectors *Paul R. Norton* 24.3

- 24.1 Scope / 24.3
- 24.2 Thermal Detectors / 24.4
- 24.3 Quantum Detectors / 24.6
- 24.4 Definitions / 24.10
- 24.5 Detector Performance and Sensitivity / 24.13
- 24.6 Other Performance Parameters / 24.18
- 24.7 Detector Performance / 24.21
- 24.8 References / 24.101
- 24.9 Suggested Readings / 24.102

Chapter 25. Photodetection *Abhay M. Joshi and Gregory H. Olsen* 25.1

- 25.1 Glossary / 25.1
- 25.2 Introduction / 25.2
- 25.3 Principle of Operation / 25.3
- 25.4 Applications / 25.11
- 25.5 Reliability / 25.13
- 25.6 Future Photodetectors / 25.15
- 25.7 Acknowledgment / 25.17
- 25.8 References / 25.18
- 25.9 Additional Reading / 25.19

Chapter 26. High-Speed Photodetectors *J. E. Bowers and Y. G. Wey* 26.1

- 26.1 Glossary / 26.1
- 26.2 Introduction / 26.3
- 26.3 Photodetector Structures / 26.3
- 26.4 Speed Limitations / 26.5
- 26.5 *p-i-n* Photodetectors / 26.10
- 26.6 Schottky Photodiode / 26.16
- 26.7 Avalanche Photodetectors / 26.17
- 26.8 Photoconductors / 26.20

- 26.9 Summary / 26.24
26.10 References / 26.24

Chapter 27. Signal Detection and Analysis*John R. Willison***27.1**

-
- 27.1 Glossary / 27.1
27.2 Introduction / 27.1
27.3 Prototype Experiment / 27.2
27.4 Noise Sources / 27.3
27.5 Applications Using Photomultipliers / 27.6
27.6 Amplifiers / 27.10
27.7 Signal Analysis / 27.12
27.8 References / 27.15

Chapter 28. Thermal Detectors *William L. Wolfe and**Paul W. Kruse***28.1**

-
- 28.1 Glossary / 28.1
28.2 Thermal Detector Elements / 28.1
28.3 Arrays / 28.7
28.4 References / 28.13

Part 6. Imaging Detectors

Chapter 29. Photographic Films *Joseph H. Altman***29.3**

-
- 29.1 Glossary / 29.3
29.2 Structure of Silver Halide Photographic Layers / 29.4
29.3 Grains / 29.5
29.4 Processing / 29.5
29.5 Exposure / 29.5
29.6 Optical Density / 29.6
29.7 The D-Log H Curve / 29.8
29.8 Spectral Sensitivity / 29.11
29.9 Reciprocity Failure / 29.11
29.10 Development Effects / 29.12
29.11 Color Photography / 29.12
29.12 Microdensitometers / 29.15
29.13 Performance of Photographic Systems / 29.16
29.14 Image Structure / 29.17
29.15 Acutance / 29.17
29.16 Graininess / 29.19
29.17 Sharpness and Graininess Considered Together / 29.22
29.18 Signal-to-Noise Ratio and Detective Quantum Efficiency / 29.22
29.19 Resolving Power / 29.24
29.20 Information Capacity / 29.24
29.21 List of Photographic Manufacturers / 29.25
29.22 References / 29.25

Chapter 30. Photographic Materials *John D. Baloga***30.1**

-
- 30.1 Introduction / 30.1
30.2 The Optics of Photographic Films and Papers / 30.2
30.3 The Photophysics of Silver Halide Light Detectors / 30.7
30.4 The Stability of Photographic Image Dyes toward Light Fade / 30.10
30.5 Photographic Spectral Sensitizers / 30.13

- 30.6 General Characteristics of Photographic Films / 30.18
- 30.7 References / 30.28

Chapter 31. Image Tube Intensified Electronic Imaging **31.1**
C. Bruce Johnson and Larry D. Owen

- 31.1 Glossary / 31.1
- 31.2 Introduction / 31.2
- 31.3 The Optical Interface / 31.3
- 31.4 Image Intensifiers / 31.7
- 31.5 Image Intensified Self-Scanned Arrays / 31.19
- 31.6 Applications / 31.27
- 31.7 References / 31.30

Chapter 32. Visible Array Detectors **32.1**
Timothy J. Tredwell

- 32.1 Glossary / 32.1
- 32.2 Introduction / 32.2
- 32.3 Image Sensing Elements / 32.2
- 32.4 Readout Elements / 32.12
- 32.5 Sensor Architectures / 32.21
- 32.6 References / 32.35

Chapter 33. Infrared Detector Arrays **33.1**
Lester J. Kozlowski and Walter F. Kosonocky

- 33.1 Glossary / 33.1
- 33.2 Introduction / 33.3
- 33.3 Monolithic FPAs / 33.10
- 33.4 Hybrid FPAs / 33.14
- 33.5 Performance: Figures of Merit / 33.23
- 33.6 Current Status and Future Trends / 33.28
- 33.7 References / 33.31

Part 7. Radiometry and Photometry

Chapter 34. Radiometry and Photometry **34.3**
Edward F. Zalewski

- 34.1 Glossary / 34.3
- 34.2 Introduction / 34.5
- 34.3 Radiometric Definitions and Basic Concepts / 34.7
- 34.4 Radiant Transfer Approximations / 34.13
- 34.5 Absolute Measurements / 34.20
- 34.6 Photometry / 34.37
- 34.7 References / 34.44

Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection **35.1**
James M. Palmer

- 35.1 Glossary / 35.1
- 35.2 Introduction and Terminology / 35.2
- 35.3 Transmittance / 35.3
- 35.4 Absorptance / 35.4
- 35.5 Reflectance / 35.4
- 35.6 Emittance / 35.7
- 35.7 Kirchhoff's Law / 35.7
- 35.8 Relationship between Transmittance, Reflectance, and Absorptance / 35.7
- 35.9 Measurement of Transmittance / 35.8

- 35.10 Measurement of Absorptance / 35.10
 35.11 Measurement of Reflectance / 35.10
 35.12 Measurement of Emittance / 35.14
 35.13 References / 35.16
 35.14 Further Reading / 35.23

Chapter 36. Radiometry and Photometry: Units and Conversions *James M. Palmer* **36.1**

- 36.1 Glossary / 36.1
 36.2 Introduction and Background / 36.2
 36.3 Symbols, Units, and Nomenclature in Radiometry / 36.4
 36.4 Symbols, Units, and Nomenclature in Photometry / 36.5
 36.5 Conversion of Radiometric Quantities to Photometric Quantities / 36.11
 36.6 Conversion of Photometric Quantities to Radiometric Quantities / 36.12
 36.7 Radiometric/Photometric Normalization / 36.14
 36.8 Other Weighting Functions and Conversions / 36.17
 36.9 References / 36.17
 36.10 Further Reading / 36.18

Chapter 37. Radiometry and Photometry for Vision Optics *Yoshi Ohno* **37.1**

- 37.1 Introduction / 37.1
 37.2 Basis of Physical Photometry / 37.1
 37.3 Photometric Base Unit—the Candela / 37.3
 37.4 Quantities and Units in Photometry and Radiometry / 37.3
 37.5 Principles in Photometry and Radiometry / 37.8
 37.6 Practice in Photometry and Radiometry / 37.11
 37.7 References / 37.12

Chapter 38. Spectroradiometry *Carolyn J. Sher DeCusatis* **38.1**

- 38.1 Introduction / 38.1
 38.2 Definitions, Calculations, and Figures of Merit / 38.1
 38.3 General Features of Spectroradiometry Systems / 38.7
 38.4 Typical Spectroradiometry System Designs / 38.13
 38.5 References / 38.19

Chapter 39. Nonimaging Optics: Concentration and Illumination *William Cassarly* **39.1**

- 39.1 Introduction / 39.1
 39.2 Basic Calculations / 39.2
 39.3 Software Modeling of Nonimaging Systems / 39.6
 39.4 Basic Building Blocks / 39.8
 39.5 Concentration / 39.12
 39.6 Uniformity and Illumination / 39.22
 39.7 Acknowledgments / 39.41
 39.8 References / 39.41

Chapter 40. Lighting and Applications *Anurag Gupta and R. John Koshel* **40.1**

- 40.1 Glossary / 40.1
 40.2 Introduction / 40.1
 40.3 Vision Biology and Perception / 40.3
 40.4 The Science of Lighting Design / 40.6

40.5	Luminaires	/	40.24
40.6	Lighting Measurements	/	40.51
40.7	Lighting Application Areas	/	40.54
40.8	Acknowledgments	/	40.71
40.9	References	/	40.72

Index	I.1
--------------	------------

DO NOT DUPLICATE

CONTRIBUTORS

- Joseph H. Altman** *Institute of Optics, University of Rochester, Rochester, New York* (CHAP. 29)
- Ladan Arissian** *Texas A&M University, College Station, Texas, and National Research Council of Canada, Ottawa, Ontario, Canada* (CHAP. 20)
- John D. Baloga** *Imaging Materials and Media, Eastman Kodak Company, Rochester, New York* (CHAP. 30)
- János A. Bergou** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Department of Physics and Astronomy, Hunter College of the City University of New York, New York, New York* (CHAP. 23)
- John E. Bowers** *Department of Electrical and Computer Engineering, University of California, Santa Barbara, California* (CHAP. 26)
- Robert P. Breault** *Breault Research Organization, Tucson, Arizona* (CHAP. 7)
- William Cassarly** *Optical Research Associates, Pasadena, California* (CHAP. 39)
- Zenghu Chang** *Department of Physics, Kansas State University, Cardwell Hall, Manhattan, Kansas* (CHAP. 21)
- M. George Craford** *Hewlett-Packard Co., San Jose, California* (CHAP. 17)
- Katherine Creath** *Optineering, Tucson, Arizona, and College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 14)
- Pamela L. Derry** *Boeing Defense & Space Group, Seattle, Washington* (CHAP. 19)
- Jean-Claude Diels** *Departments of Physics and Electrical Engineering, University of New Mexico, Albuquerque, New Mexico* (CHAP. 20)
- Berthold-Georg Englert** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Max-Planck-Institut für Quantenoptik, Garching bei München, Germany, and Abteilung Quantenphysik der Universität Ulm, Ulm, Germany* (CHAP. 23)
- Chris J. Evans** *Zygo Corporation, Middlefield, Connecticut* (CHAP. 10)
- Luis Figueroa** *Boeing Defense & Space Group, Seattle, Washington* (CHAP. 19)
- Anurag Gupta** *Optical Research Associates, Tucson, Arizona* (CHAP. 40)
- Roland H. Haitz** *Hewlett-Packard Co., San Jose, California* (CHAP. 17)
- John L. Hall** *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (CHAP. 22)
- Michael E. Harrigan** *Harrigan Optical Design, Victor, New York* (CHAP. 2)
- Chi-Shain Hong** *Boeing Defense & Space Group, Seattle, Washington* (CHAP. 19)
- C. Bruce Johnson** *Johnson Scientific Group, Inc., Phoenix, Arizona* (CHAP. 31)
- Abhay M. Joshi** *Discovery Semiconductors, Inc., Cranbury, New Jersey* (CHAP. 25)
- R. John Koschel** *Photon Engineering LLC, and College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 40)
- Walter F. Kosonocky*** *New Jersey Institute of Technology, University Heights, Newark, New Jersey* (CHAP. 33)
- Lester J. Kozlowski** *Altasens, Inc., Westlake Village, California* (CHAP. 33)
- Paul W. Kruse** *Consultant, Edina, Minnesota* (CHAP. 28)
- Anthony LaRocca†** *General Dynamics, Advanced Information Systems, Ypsilanti, Michigan* (CHAP. 15)

*Deceased.

†Retired.

- Melvin Lax*** *Department of Physics, City College of the City University of New York, New York, New York* (CHAP. 23)
- Virendra N. Mahajan** *The Aerospace Corporation, El Segundo, California* (CHAP. 11)
- Zacarias Malacara** *Centro de Investigaciones en Óptica, A. C., León, Gto., México* (CHAP. 12)
- Daniel Malacara-Hernández** *Centro de Investigaciones en Óptica, A. C., León, Gto., México* (CHAPS. 12, 13)
- Michael P. Mandina** *Brandon Light, Optimax Systems, Inc., Ontario, New York* (CHAP. 9)
- Paul R. Norton** *U.S. Army Night Vision and Electronics Directorate, Fort Belvoir, Virginia* (CHAP. 24)
- Donald C. O'Shea** *Georgia Institute of Technology, School of Physics, Atlanta, Georgia* (CHAP. 2)
- Yoshi Ohno** *Optical Technology Division, National Institute of Standards and Technology, Gaithersburg, Maryland* (CHAP. 37)
- Gregory H. Olsen** *Sensors Unlimited, Inc., Princeton, New Jersey* (CHAP. 25)
- Larry D. Owen** *NuOptics International, Phoenix, Arizona* (CHAP. 31)
- James M. Palmer*** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAPS. 35, 36)
- Richard L. Rhorer** *National Institute of Standards and Technology, Gaithersburg, Maryland* (CHAP. 10)
- Michael Roberts** *Pilkington Optronics, Wales, United Kingdom* (CHAP. 8)
- Philip J. Rogers** *Pilkington Optronics, Wales, United Kingdom* (CHAP. 8)
- Winston V. Schoenfeld** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 18)
- Marian O. Scully** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Max-Planck-Institut für Quantenoptik, Garching bei München, Germany* (CHAP. 23)
- Robert R. Shannon†** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAPS. 4, 5)
- Carolyn J. Sher DeCusatis** *Pace University, White Plains, New York* (CHAP. 38)
- William T. Silfvast** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 16)
- Douglas C. Sinclair** *Sinclair Optics, Inc., Fairport, New York* (CHAP. 3)
- Warren J. Smith*** *Kaiser Electro-Optics, Inc., Carlsbad, California* (CHAP. 1)
- Matthew S. Taubman** *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (CHAP. 22)
- Timothy J. Tredwell** *Sensor Systems Division, Imager Systems Development Laboratory, Eastman Kodak Company, Rochester, New York* (CHAP. 32)
- Herbert Walther*** *Max-Planck-Institut für Quantenoptik, Garching bei München, Germany, and Sektion Physik der Universität München, Garching bei München, Germany* (CHAP. 23)
- Robert H. Weissman** *Hewlett-Packard Co., San Jose, California* (CHAP. 17)
- Yih G. Wey** *Department of Electrical and Computer Engineering, University of California, Santa Barbara, California* (CHAP. 26)
- John R. Willison** *Stanford Research Systems, Inc., Sunnyvale, California* (CHAP. 27)
- William L. Wolfe** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 28)
- James C. Wyant** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 14)
- Jun Ye** *JILA, University of Colorado and National Institute of Standards and Technology, Boulder, Colorado* (CHAP. 22)
- Paul R. Yoder, Jr.** *Consultant in Optical Engineering, Norwalk, Connecticut* (CHAP. 6)
- Edward F. Zalewski** *College of Optical Sciences, University of Arizona, Tucson, Arizona* (CHAP. 34)
- M. Suhail Zubairy** *Institute for Quantum Studies and Department of Physics, Texas A&M University, College Station, Texas, and Department of Electronics, Quaid-i-Azam University, Islamabad, Pakistan* (CHAP. 23)

*Deceased.

†Retired.

BRIEF CONTENTS OF ALL VOLUMES

VOLUME I. GEOMETRICAL AND PHYSICAL OPTICS, POLARIZED LIGHT, COMPONENT AND INSTRUMENTS

PART 1. GEOMETRICAL OPTICS

Chapter 1. General Principles of Geometrical Optics *Douglas S. Goodman*

PART 2. PHYSICAL OPTICS

Chapter 2. Interference *John E. Greivenkamp*

Chapter 3. Diffraction *Arvind S. Marathay and John F. McCalmont*

Chapter 4. Transfer Function Techniques *Glenn D. Boreman*

Chapter 5. Coherence Theory *William H. Carter*

Chapter 6. Coherence Theory: Tools and Applications *Gisele Bennett, William T. Rhodes, and J. Christopher James*

Chapter 7. Scattering by Particles *Craig F. Bohren*

Chapter 8. Surface Scattering *Eugene L. Church and Peter Z. Takacs*

Chapter 9. Volume Scattering in Random Media *Aristide Dogariu and Jeremy Ellis*

Chapter 10. Optical Spectroscopy and Spectroscopic Lineshapes *Brian Henderson*

Chapter 11. Analog Optical Signal and Image Processing *Joseph W. Goodman*

PART 3. POLARIZED LIGHT

Chapter 12. Polarization *Jean M. Bennett*

Chapter 13. Polarizers *Jean M. Bennett*

Chapter 14. Mueller Matrices *Russell A. Chipman*

Chapter 15. Polarimetry *Russell A. Chipman*

Chapter 16. Ellipsometry *Rasheed M. A. Azzam*

PART 4. COMPONENTS

Chapter 17. Lenses *R. Barry Johnson*

Chapter 18. Afocal Systems *William B. Wetherell*

Chapter 19. Nondispersive Prisms *William L. Wolfe*

Chapter 20. Dispersive Prisms and Gratings *George J. Zissis*

Chapter 21. Integrated Optics *Thomas L. Koch, Frederick J. Leonberger, and Paul G. Suchoski*

Chapter 22. Miniature and Micro-Optics *Tom D. Milster and Tomasz S. Tkaczyk*

Chapter 23. Binary Optics *Michael W. Farn and Wilfrid B. Veldkamp*

Chapter 24. Gradient Index Optics *Duncan T. Moore*

PART 5. INSTRUMENTS

Chapter 25. Cameras *Norman Goldberg*

Chapter 26. Solid-State Cameras *Gerald C. Holst*

Chapter 27. Camera Lenses *Ellis Betensky, Melvin H. Kreitzer, and Jacob Moskovich*

Chapter 28. Microscopes *Rudolf Oldenbourg and Michael Shribak*

Chapter 29. Reflective and Catadioptric Objectives *Lloyd Jones*

- Chapter 30. Scanners *Leo Beiser and R. Barry Johnson*
- Chapter 31. Optical Spectrometers *Brian Henderson*
- Chapter 32. Interferometers *Parameswaran Hariharan*
- Chapter 33. Holography and Holographic Instruments *Lloyd Huff*
- Chapter 34. Xerographic Systems *Howard Stark*
- Chapter 35. Principles of Optical Disk Data Storage *Masud Mansuripur*

VOLUME II. DESIGN, FABRICATION, AND TESTING; SOURCES AND DETECTORS; RADIOMETRY AND PHOTOMETRY

PART 1. DESIGN

- Chapter 1. Techniques of First-Order Layout *Warren J. Smith*
- Chapter 2. Aberration Curves in Lens Design *Donald C. O'Shea and Michael E. Harrigan*
- Chapter 3. Optical Design Software *Douglas C. Sinclair*
- Chapter 4. Optical Specifications *Robert R. Shannon*
- Chapter 5. Tolerancing Techniques *Robert R. Shannon*
- Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.*
- Chapter 7. Control of Stray Light *Robert P. Breault*
- Chapter 8. Thermal Compensation Techniques *Philip J. Rogers and Michael Roberts*

PART 2. FABRICATION

- Chapter 9. Optical Fabrication *Michael P. Mandina*
- Chapter 10. Fabrication of Optics by Diamond Turning *Richard L. Rhorer and Chris J. Evans*

PART 3. TESTING

- Chapter 11. Orthonormal Polynomials in Wavefront Analysis *Virendra N. Mahajan*
- Chapter 12. Optical Metrology *Zacarias Malacara and Daniel Malacara-Hernández*
- Chapter 13. Optical Testing *Daniel Malacara-Hernández*
- Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant*

PART 4. SOURCES

- Chapter 15. Artificial Sources *Anthony LaRocca*
- Chapter 16. Lasers *William T. Silfvast*
- Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman*
- Chapter 18. High-Brightness Visible LEDs *Winston V. Schoenfeld*
- Chapter 19. Semiconductor Lasers *Pamela L. Derry, Luis Figueroa, and Chi-shain Hong*
- Chapter 20. Ultrashort Optical Sources and Applications *Jean-Claude Diels and Ladan Arissian*
- Chapter 21. Attosecond Optics *Zenghu Chang*
- Chapter 22. Laser Stabilization *John L. Hall, Matthew S. Taubman, and Jun Ye*
- Chapter 23. Quantum Theory of the Laser *János A. Bergou, Berthold-Georg Englert, Melvin Lax, Marian O. Scully, Herbert Walther, and M. Suhail Zubairy*

PART 5. DETECTORS

- Chapter 24. Photodetectors *Paul R. Norton*
- Chapter 25. Photodetection *Abhay M. Joshi and Gregory H. Olsen*
- Chapter 26. High-Speed Photodetectors *John E. Bowers and Yih G. Wey*
- Chapter 27. Signal Detection and Analysis *John R. Willison*
- Chapter 28. Thermal Detectors *William L. Wolfe and Paul W. Kruse*

PART 6. IMAGING DETECTORS

- Chapter 29. Photographic Films *Joseph H. Altman*
- Chapter 30. Photographic Materials *John D. Baloga*

- Chapter 31. Image Tube Intensified Electronic Imaging *C. Bruce Johnson and Larry D. Owen*
 Chapter 32. Visible Array Detectors *Timothy J. Tredwell*
 Chapter 33. Infrared Detector Arrays *Lester J. Kozlowski and Walter F. Kosonocky*

PART 7. RADIOMETRY AND PHOTOMETRY

- Chapter 34. Radiometry and Photometry *Edward F. Zalewski*
 Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection *James M. Palmer*
 Chapter 36. Radiometry and Photometry: Units and Conversions *James M. Palmer*
 Chapter 37. Radiometry and Photometry for Vision Optics *Yoshi Ohno*
 Chapter 38. Spectroradiometry *Carolyn J. Sher DeCusatis*
 Chapter 39. Nonimaging Optics: Concentration and Illumination *William Cassarly*
 Chapter 40. Lighting and Applications *Anurag Gupta and R. John Koshel*

VOLUME III. VISION AND VISION OPTICS

- Chapter 1. Optics of the Eye *Neil Charman*
 Chapter 2. Visual Performance *Wilson S. Geisler and Martin S. Banks*
 Chapter 3. Psychophysical Methods *Denis G. Pelli and Bart Farell*
 Chapter 4. Visual Acuity and Hyperacuity *Gerald Westheimer*
 Chapter 5. Optical Generation of the Visual Stimulus *Stephen A. Burns and Robert H. Webb*
 Chapter 6. The Maxwellian View with an Addendum on Apodization *Gerald Westheimer*
 Chapter 7. Ocular Radiation Hazards *David H. Sliney*
 Chapter 8. Biological Waveguides *Vasudevan Lakshminarayanan and Jay M. Enoch*
 Chapter 9. The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution *Jay M. Enoch and Vasudevan Lakshminarayanan*
 Chapter 10. Colorimetry *David H. Brainard and Andrew Stockman*
 Chapter 11. Color Vision Mechanisms *Andrew Stockman and David H. Brainard*
 Chapter 12. Assessment of Refraction and Refractive Errors and Their Influence on Optical Design *B. Ralph Chou*
 Chapter 13. Binocular Vision Factors That Influence Optical Design *Clifton Schor*
 Chapter 14. Optics and Vision of the Aging Eye *John S. Werner, Brooke E. Scheffrin, and Arthur Bradley*
 Chapter 15. Adaptive Optics in Retinal Microscopy and Vision *Donald T. Miller and Austin Roorda*
 Chapter 16. Refractive Surgery, Correction of Vision, PRK and LASIK *L. Diaz-Santana and Harilaos Ginis*
 Chapter 17. Three-Dimensional Confocal Microscopy of the Living Human Cornea *Barry R. Masters*
 Chapter 18. Diagnostic Use of Optical Coherence Tomography in the Eye *Johannes F. de Boer*
 Chapter 19. Gradient Index Optics in the Eye *Barbara K. Pierscionek*
 Chapter 20. Optics of Contact Lenses *Edward S. Bennett*
 Chapter 21. Intraocular Lenses *Jim Schwiegerling*
 Chapter 22. Displays for Vision Research *William Cowan*
 Chapter 23. Vision Problems at Computers *Jeffrey Anshel and James E. Sheedy*
 Chapter 24. Human Vision and Electronic Imaging *Bernice E. Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allebach*
 Chapter 25. Visual Factors Associated with Head-Mounted Displays *Brian H. Tsou and Martin Shenker*

VOLUME IV. OPTICAL PROPERTIES OF MATERIALS, NONLINEAR OPTICS, QUANTUM OPTICS

PART 1. PROPERTIES

- Chapter 1. Optical Properties of Water *Curtis D. Mobley*
 Chapter 2. Properties of Crystals and Glasses *William J. Tropf, Michael E. Thomas, and Eric W. Rogala*
 Chapter 3. Polymeric Optics *John D. Lytle*
 Chapter 4. Properties of Metals *Roger A. Paquin*

- Chapter 5. Optical Properties of Semiconductors *David G. Seiler, Stefan Zollner, Alain C. Diebold, and Paul M. Amirtharaj*
- Chapter 6. Characterization and Use of Black Surfaces for Optical Systems *Stephen M. Pompea and Robert P. Breault*
- Chapter 7. Optical Properties of Films and Coatings *Jerzy A. Dobrowolski*
- Chapter 8. Fundamental Optical Properties of Solids *Alan Miller*
- Chapter 9. Photonic Bandgap Materials *Pierre R. Villeneuve*

PART 2. NONLINEAR OPTICS

- Chapter 10. Nonlinear Optics *Chung L. Tang*
- Chapter 11. Coherent Optical Transients *Paul R. Berman and D. G. Steel*
- Chapter 12. Photorefractive Materials and Devices *Mark Cronin-Golomb and Marvin Klein*
- Chapter 13. Optical Limiting *David J. Hagan*
- Chapter 14. Electromagnetically Induced Transparency *Jonathan P. Marangos and Thomas Halfmann*
- Chapter 15. Stimulated Raman and Brillouin Scattering *John Reintjes and M. Bashkansky*
- Chapter 16. Third-Order Optical Nonlinearities *Mansoor Sheik-Bahae and Michael P. Hasselbeck*
- Chapter 17. Continuous-Wave Optical Parametric Oscillators *M. Ebrahim-Zadeh*
- Chapter 18. Nonlinear Optical Processes for Ultrashort Pulse Generation *Uwe Siegner and Ursula Keller*
- Chapter 19. Laser-Induced Damage to Optical Materials *Marion J. Soileau*

PART 3. QUANTUM AND MOLECULAR OPTICS

- Chapter 20. Laser Cooling and Trapping of Atoms *Harold J. Metcalf and Peter van der Straten*
- Chapter 21. Strong Field Physics *Todd Ditmire*
- Chapter 22. Slow Light Propagation in Atomic and Photonic Media *Jacob B. Khurgin*
- Chapter 23. Quantum Entanglement in Optical Interferometry *Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver, and Jonathan P. Dowling*

VOLUME V. ATMOSPHERIC OPTICS, MODULATORS, FIBER OPTICS, X-RAY AND NEUTRON OPTICS

PART 1. MEASUREMENTS

- Chapter 1. Scatterometers *John C. Stover*
- Chapter 2. Spectroscopic Measurements *Brian Henderson*

PART 2. ATMOSPHERIC OPTICS

- Chapter 3. Atmospheric Optics *Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman*
- Chapter 4. Imaging through Atmospheric Turbulence *Virendra N. Mahajan and Guang-ming Dai*
- Chapter 5. Adaptive Optics *Robert Q. Fugate*

PART 3. MODULATORS

- Chapter 6. Acousto-Optic Devices *I-Cheng Chang*
- Chapter 7. Electro-Optic Modulators *Georgeanne M. Purvinis and Theresa A. Maldonado*
- Chapter 8. Liquid Crystals *Sebastian Gauza and Shin-Tson Wu*

PART 4. FIBER OPTICS

- Chapter 9. Optical Fiber Communication Technology and System Overview *Ira Jacobs*
- Chapter 10. Nonlinear Effects in Optical Fibers *John A. Buck*
- Chapter 11. Photonic Crystal Fibers *Philip St. J. Russell and G. J. Pearce*
- Chapter 12. Infrared Fibers *James A. Harrington*
- Chapter 13. Sources, Modulators, and Detectors for Fiber Optic Communication Systems *Elsa Garmire*
- Chapter 14. Optical Fiber Amplifiers *John A. Buck*
- Chapter 15. Fiber Optic Communication Links (Telecom, Datacom, and Analog) *Casimer DeCusatis and Guifang Li*

- Chapter 16. Fiber-Based Couplers *Daniel Nolan*
 Chapter 17. Fiber Bragg Gratings *Kenneth O. Hill*
 Chapter 18. Micro-Optics-Based Components for Networking *Joseph C. Palais*
 Chapter 19. Semiconductor Optical Amplifiers *Jay M. Wiesenfeld and Leo H. Spiekman*
 Chapter 20. Optical Time-Division Multiplexed Communication Networks *Peter J. Delfyett*
 Chapter 21. WDM Fiber-Optic Communication Networks *Alan E. Willner, Changyuan Yu, Zhongqi Pan, and Yong Xie*
 Chapter 22. Solitons in Optical Fiber Communication Systems *Pavel V. Mamyshev*
 Chapter 23. Fiber-Optic Communication Standards *Casimer DeCusatis*
 Chapter 24. Optical Fiber Sensors *Richard O. Claus, Ignacio Matias, and Francisco Arregui*
 Chapter 25. High-Power Fiber Lasers and Amplifiers *Timothy S. McComb, Martin C. Richardson, and Michael Bass*

PART 5. X-RAY AND NEUTRON OPTICS

Subpart 5.1. Introduction and Applications

- Chapter 26. An Introduction to X-Ray and Neutron Optics *Carolyn A. MacDonald*
 Chapter 27. Coherent X-Ray Optics and Microscopy *Qun Shen*
 Chapter 28. Requirements for X-Ray diffraction *Scott T. Misture*
 Chapter 29. Requirements for X-Ray Fluorescence *George J. Havrilla*
 Chapter 30. Requirements for X-Ray Spectroscopy *Dirk Lützenkirchen-Hecht and Ronald Frahm*
 Chapter 31. Requirements for Medical Imaging and X-Ray Inspection *Douglas Pfeiffer*
 Chapter 32. Requirements for Nuclear Medicine *Lars R. Furenlid*
 Chapter 33. Requirements for X-Ray Astronomy *Scott O. Rohrbach*
 Chapter 34. Extreme Ultraviolet Lithography *Franco Cerrina and Fan Jiang*
 Chapter 35. Ray Tracing of X-Ray Optical Systems *Franco Cerrina and M. Sanchez del Rio*
 Chapter 36. X-Ray Properties of Materials *Eric M. Gullikson*

Subpart 5.2. Refractive and Interference Optics

- Chapter 37. Refractive X-Ray Lenses *Bruno Lengeler and Christian G. Schroer*
 Chapter 38. Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region *Malcolm R. Howells*
 Chapter 39. Crystal Monochromators and Bent Crystals *Peter Siddons*
 Chapter 40. Zone Plates *Alan Michette*
 Chapter 41. Multilayers *Eberhard Spiller*
 Chapter 42. Nanofocusing of Hard X-Rays with Multilayer Laue Lenses *Albert T. Macrander, Hanfei Yan, Hyon Chol Kang, Jörg Maser, Chian Liu, Ray Conley, and G. Brian Stephenson*
 Chapter 43. Polarizing Crystal Optics *Qun Shen*

Subpart 5.3. Reflective Optics

- Chapter 44. Reflective Optics *James Harvey*
 Chapter 45. Aberrations for Grazing Incidence Optics *Timo T. Saha*
 Chapter 46. X-Ray Mirror Metrology *Peter Z. Takacs*
 Chapter 47. Astronomical X-Ray Optics *Marshall K. Joy and Brian D. Ramsey*
 Chapter 48. Multifoil X-Ray Optics *Ladislav Pina*
 Chapter 49. Pore Optics *Marco Beijersbergen*
 Chapter 50. Adaptive X-Ray Optics *Ali Khounsary*
 Chapter 51. The Schwarzschild Objective *Franco Cerrina*
 Chapter 52. Single Capillaries *Donald H. Bilderback and Sterling W. Cornaby*
 Chapter 53. Polycapillary X-Ray Optics *Carolyn MacDonald and Walter M. Gibson*

Subpart 5.4. X-Ray Sources

- Chapter 54. X-Ray Tube Sources *Susanne M. Lee and Carolyn MacDonald*
 Chapter 55. Synchrotron Sources *Steven L. Hulbert and Gwyn P. Williams*
 Chapter 56. Laser Generated Plasmas *Alan Michette*

- Chapter 57. Pinch Plasma Sources *Victor Kantsyrev*
Chapter 58. X-Ray Lasers *Greg Tallents*
Chapter 59. Inverse Compton X-Ray Sources *Frank Carroll*

Subpart 5.5. X-Ray Detectors

- Chapter 60. Introduction to X-Ray Detectors *Walter M. Gibson and Peter Siddons*
Chapter 61. Advances in Imaging Detectors *Aaron Couture*
Chapter 62. X-Ray Spectral Detection and Imaging *Eric Lifshin*

Subpart 5.6. Neutron Optics and Applications

- Chapter 63. Neutron Optics *David Mildner*
Chapter 64. Grazing-Incidence Neutron Optics *Mikhail Gubarev and Brian Ramsey*

DO NOT DUPLICATE

EDITORS' PREFACE

The third edition of the *Handbook of Optics* is designed to pull together the dramatic developments in both the basic and applied aspects of the field while retaining the archival, reference book value of a handbook. This means that it is much more extensive than either the first edition, published in 1978, or the second edition, with Volumes I and II appearing in 1995 and Volumes III and IV in 2001. To cover the greatly expanded field of optics, the *Handbook* now appears in five volumes. Over 100 authors or author teams have contributed to this work.

Volume I is devoted to the fundamentals, components, and instruments that make optics possible. Volume II contains chapters on design, fabrication, testing, sources of light, detection, and a new section devoted to radiometry and photometry. Volume III concerns vision optics only and is printed entirely in color. In Volume IV there are chapters on the optical properties of materials, nonlinear, quantum and molecular optics. Volume V has extensive sections on fiber optics and x ray and neutron optics, along with shorter sections on measurements, modulators, and atmospheric optical properties and turbulence. Several pages of color inserts are provided where appropriate to aid the reader. A purchaser of the print version of any volume of the *Handbook* will be able to download a digital version containing all of the material in that volume in PDF format to one computer (see download instructions on bound-in card). The combined index for all five volumes can be downloaded from www.HandbookofOpticsOnline.com.

It is possible by careful selection of what and how to present that the third edition of the *Handbook* could serve as a text for a comprehensive course in optics. In addition, students who take such a course would have the *Handbook* as a career-long reference.

Topics were selected by the editors so that the *Handbook* could be a desktop (bookshelf) general reference for the parts of optics that had matured enough to warrant archival presentation. New chapters were included on topics that had reached this stage since the second edition, and existing chapters from the second edition were updated where necessary to provide this compendium. In selecting subjects to include, we also had to select which subjects to leave out. The criteria we applied were: (1) was it a specific application of optics rather than a core science or technology and (2) was it a subject in which the role of optics was peripheral to the central issue addressed. Thus, such topics as medical optics, laser surgery, and laser materials processing were not included. While applications of optics are mentioned in the chapters there is no space in the *Handbook* to include separate chapters devoted to all of the myriad uses of optics in today's world. If we had, the third edition would be much longer than it is and much of it would soon be outdated. We designed the third edition of the *Handbook of Optics* so that it concentrates on the principles of optics that make applications possible.

Authors were asked to try to achieve the dual purpose of preparing a chapter that was a worthwhile reference for someone working in the field and that could be used as a starting point to become acquainted with that aspect of optics. They did that and we thank them for the outstanding results seen throughout the *Handbook*. We also thank Mr. Taisuke Soda of McGraw-Hill for his help in putting this complex project together and Mr. Alan Tourtlotte and Ms. Susannah Lehman of the Optical Society of America for logistical help that made this effort possible.

We dedicate the third edition of the *Handbook of Optics* to all of the OSA volunteers who, since OSA's founding in 1916, give their time and energy to promoting the generation, application, archiving, and worldwide dissemination of knowledge in optics and photonics.

Michael Bass, Editor-in-Chief

Associate Editors:

Casimer M. DeCusatis

Jay M. Enoch

Vasudevan Lakshminarayanan

Guifang Li

Carolyn MacDonald

Virendra N. Mahajan

Eric Van Stryland

This page intentionally left blank.

DO NOT DUPLICATE

PREFACE TO VOLUME II

Volume II of the *Handbook of Optics* is a continuation of Volume I. It starts with optical system design and covers first-order layout, aberration curves, design software, specifications and tolerances, component mounting, stray light control, and thermal compensation techniques. Optical fabrication and testing are discussed next. A new chapter on the use of orthonormal polynomials in optical design and testing has been added. Such a polynomial representing balanced astigmatism is illustrated on the cover. The section on sources includes different types of lasers, laser stabilization, laser theory, and a discussion of ultrashort laser sources. Light-emitting diodes including the new “high-brightness” LEDs are presented. Artificial sources of light for both the laboratory and field are described along with a discussion of light standards calibration. The section on detectors includes high-speed and thermal detectors along with an analysis of signal detection. Imaging using film, detector arrays, and image tubes is discussed. This volume ends with a section on radiometry and photometry. Two new chapters have been added in this area. One is on spectroradiometry and the other is on lighting and applications.

Every effort was made to contact all the authors of chapters in the second edition that would appear in this edition so that they could update their chapters. However, the authors of several chapters could not be located or were not available. Their chapters are reproduced without update. Every effort has been made to ensure that such chapters have been correctly reproduced.

There are many other chapters in this edition of the *Handbook* that could have been included in Volumes I and II. However, page limitations prevented that. For example, in Volume V there is a section on Atmospheric Optics. It consists of three chapters, one on transmission through the atmosphere, another on imaging through atmospheric turbulence, and a third on adaptive optics to overcome some of the deleterious effects of turbulence.

The chapters are generally aimed at the graduate students, though practicing scientists and engineers will find them equally suitable as references on the topics discussed. Each chapter has sufficient references for additional and/or further study.

Virendra N. Mahajan
The Aerospace Corporation
Eric Van Stryland
CREOL, The College of Optics and Photonics
Associate Editors

This page intentionally left blank.

DO NOT DUPLICATE

GLOSSARY AND FUNDAMENTAL CONSTANTS

Introduction

This glossary of the terms used in the *Handbook* represents to a large extent the language of optics. The symbols are representations of numbers, variables, and concepts. Although the basic list was compiled by the author of this section, all the editors have contributed and agreed to this set of symbols and definitions. Every attempt has been made to use the same symbols for the same concepts throughout the entire *Handbook*, although there are exceptions. Some symbols seem to be used for many concepts. The symbol α is a prime example, as it is used for absorptivity, absorption coefficient, coefficient of linear thermal expansion, and more. Although we have tried to limit this kind of redundancy, we have also bowed deeply to custom.

Units

The abbreviations for the most common units are given first. They are consistent with most of the established lists of symbols, such as given by the International Standards Organization ISO¹ and the International Union of Pure and Applied Physics, IUPAP.²

Prefixes

Similarly, a list of the numerical prefixes¹ that are most frequently used is given, along with both the common names (where they exist) and the multiples of ten that they represent.

Fundamental Constants

The values of the fundamental constants³ are listed following the sections on SI units.

Symbols

The most commonly used symbols are then given. Most chapters of the *Handbook* also have a glossary of the terms and symbols specific to them for the convenience of the reader. In the following list, the symbol is given, its meaning is next, and the most customary unit of measure for the quantity is presented in brackets. A bracket with a dash in it indicates that the quantity is unitless. Note that there is a difference between units and dimensions. An angle has units of degrees or radians and a solid angle square degrees or steradians, but both are pure ratios and are dimensionless. The unit symbols as recommended in the SI system are used, but decimal multiples of some of the dimensions are sometimes given. The symbols chosen, with some cited exceptions, are also those of the first two references.

RATIONALE FOR SOME DISPUTED SYMBOLS

The choice of symbols is a personal decision, but commonality improves communication. This section explains why the editors have chosen the preferred symbols for the *Handbook*. We hope that this will encourage more agreement.

Fundamental Constants

It is encouraging that there is almost universal agreement for the symbols for the fundamental constants. We have taken one small exception by adding a subscript B to the k for Boltzmann's constant.

Mathematics

We have chosen i as the imaginary unit arbitrarily. IUPAP lists both i and j , while ISO does not report on these.

Spectral Variables

These include expressions for the wavelength λ , frequency ν , wave number σ , ω for circular or radian frequency, k for circular or radian wave number and dimensionless frequency x . Although some use f for frequency, it can be easily confused with electronic or spatial frequency. Some use $\tilde{\nu}$ for wave number, but, because of typography problems and agreement with ISO and IUPAP, we have chosen σ ; it should not be confused with the Stefan-Boltzmann constant. For spatial frequencies we have chosen ξ and η , although f_x and f_y are sometimes used. ISO and IUPAP do not report on these.

Radiometry

Radiometric terms are contentious. The most recent set of recommendations by ISO and IUPAP are L for radiance [$\text{Wcm}^{-2}\text{sr}^{-1}$], M for radiant emittance or exitance [Wcm^{-2}], E for irradiance or incidence [Wcm^{-2}], and I for intensity [Wsr^{-2}]. The previous terms, W , H , N , and J , respectively, are still in many texts, notably Smith⁴ and Lloyd⁵ but we have used the revised set, although there are still shortcomings. We have tried to deal with the vexatious term *intensity* by using *specific intensity* when the units are $\text{Wcm}^{-2}\text{sr}^{-1}$, *field intensity* when they are Wcm^{-2} , and *radiometric intensity* when they are Wsr^{-1} .

There are two sets of terms for these radiometric quantities, which arise in part from the terms for different types of reflection, transmission, absorption, and emission. It has been proposed that the *ion* ending indicate a process, that the *ance* ending indicate a value associated with a particular sample, and that the *ivity* ending indicate a generic value for a "pure" substance. Then one also has reflectance, transmittance, absorptance, and emittance as well as reflectivity, transmissivity, absorptivity, and emissivity. There are now two different uses of the word emissivity. Thus the words *exitance*, *incidence*, and *sterance* were coined to be used in place of emittance, irradiance, and radiance. It is interesting that ISO uses radiance, exitance, and irradiance whereas IUPAP uses radiance, exitance [*sic*], and irradiance. We have chosen to use them both, i.e., emittance, irradiance, and radiance will be followed in square brackets by exitance, incidence, and sterance (or vice versa). Individual authors will use the different endings for transmission, reflection, absorption, and emission as they see fit.

We are still troubled by the use of the symbol E for irradiance, as it is so close in meaning to electric field, but we have maintained that accepted use. The spectral concentrations of these quantities, indicated by a wavelength, wave number, or frequency subscript (e.g., L_λ) represent partial differentiations; a subscript q represents a photon quantity; and a subscript ν indicates a quantity normalized to the response of the eye. Thereby, L_ν is luminance, E_ν illuminance, and M_ν and I_ν luminous emittance and luminous intensity. The symbols we have chosen are consistent with ISO and IUPAP.

The refractive index may be considered a radiometric quantity. It is generally complex and is indicated by $\tilde{n} = n - ik$. The real part is the relative refractive index and k is the extinction coefficient. These are consistent with ISO and IUPAP, but they do not address the complex index or extinction coefficient.

Optical Design

For the most part ISO and IUPAP do not address the symbols that are important in this area.

There were at least 20 different ways to indicate focal ratio; we have chosen FN as symmetrical with NA; we chose f and efl to indicate the effective focal length. Object and image distance, although given many different symbols, were finally called s_o and s_i since s is an almost universal symbol for distance. Field angles are θ and ϕ ; angles that measure the slope of a ray to the optical axis are u ; u can also be $\sin u$. Wave aberrations are indicated by W_{ijk} , while third-order ray aberrations are indicated by σ_i and more mnemonic symbols.

Electromagnetic Fields

There is no argument about \mathbf{E} and \mathbf{H} for the electric and magnetic field strengths, Q for quantity of charge, ρ for volume charge density, σ for surface charge density, etc. There is no guidance from Refs. 1 and 2 on polarization indication. We chose \perp and \parallel rather than p and s , partly because s is sometimes also used to indicate scattered light.

There are several sets of symbols used for reflection transmission, and (sometimes) absorption, each with good logic. The versions of these quantities dealing with field amplitudes are usually specified with lower case symbols: r , t , and a . The versions dealing with power are alternately given by the uppercase symbols or the corresponding Greek symbols: R and T versus ρ and τ . We have chosen to use the Greek, mainly because these quantities are also closely associated with Kirchhoff's law that is usually stated symbolically as $\alpha = \epsilon$. The law of conservation of energy for light on a surface is also usually written as $\alpha + \rho + \tau = 1$.

Base SI Quantities

length	m	meter
time	s	second
mass	kg	kilogram
electric current	A	ampere
temperature	K	kelvin
amount of substance	mol	mole
luminous intensity	cd	candela

Derived SI Quantities

energy	J	joule
electric charge	C	coulomb
electric potential	V	volt
electric capacitance	F	farad
electric resistance	Ω	ohm
electric conductance	S	siemens
magnetic flux	Wb	weber
inductance	H	henry
pressure	Pa	pascal
magnetic flux density	T	tesla
frequency	Hz	hertz
power	W	watt
force	N	newton
angle	rad	radian
angle	sr	steradian

Prefixes

Symbol	Name	Common name	Exponent of ten
F	exa		18
P	peta		15
T	tera	trillion	12
G	giga	billion	9
M	mega	million	6
k	kilo	thousand	3
h	hecto	hundred	2
da	deca	ten	1
d	deci	tenth	-1
c	centi	hundredth	-2
m	milli	thousandth	-3
μ	micro	millionth	-6
n	nano	billionth	-9
p	pico	trillionth	-12
f	femto		-15
a	atto		-18

Constants

c	speed of light vacuo [299792458 ms ⁻¹]
c_1	first radiation constant = $2\pi^2 h = 3.7417749 \times 10^{-16}$ [Wm ²]
c_2	second radiation constant = $hc/k = 0.014838769$ [mK]
e	elementary charge [$1.60217733 \times 10^{-19}$ C]
g_n	free fall constant [9.80665 ms ⁻²]
h	Planck's constant [$6.6260755 \times 10^{-34}$ Ws]
k_B	Boltzmann constant [1.380658×10^{-23} JK ⁻¹]
m_e	mass of the electron [$9.1093897 \times 10^{-31}$ kg]
N_A	Avogadro constant [6.0221367×10^{23} mol ⁻¹]
R_∞	Rydberg constant [10973731.534 m ⁻¹]
ϵ_0	vacuum permittivity [$\mu_0^{-1}c^{-2}$]
σ	Stefan-Boltzmann constant [5.67051×10^{-8} Wm ⁻¹ K ⁻⁴]
μ_0	vacuum permeability [$4\pi \times 10^{-7}$ NA ⁻²]
μ_B	Bohr magneton [$9.2740154 \times 10^{-24}$ JT ⁻¹]

General

B	magnetic induction [Wbm ⁻² , kgs ⁻¹ C ⁻¹]
C	capacitance [f, C ² s ² m ⁻² kg ⁻¹]
C	curvature [m ⁻¹]
c	speed of light in vacuo [ms ⁻¹]
c_1	first radiation constant [Wm ²]
c_2	second radiation constant [mK]
D	electric displacement [Cm ⁻²]
E	incidence [irradiance] [Wm ⁻²]
e	electronic charge [coulomb]
E_v	illuminance [lux, lmm ⁻²]
E	electrical field strength [Vm ⁻¹]
E	transition energy [J]
E_g	band-gap energy [eV]
f^g	focal length [m]
f_f	Fermi occupation function, conduction band
f_v	Fermi occupation function, valence band

FN	focal ratio (<i>f</i> /number) [—]
<i>g</i>	gain per unit length [m^{-1}]
g_{th}	gain threshold per unit length [m^{-1}]
H	magnetic field strength [Am^{-1} , $\text{Cs}^{-1} \text{m}^{-1}$]
<i>h</i>	height [m]
<i>I</i>	irradiance (see also <i>E</i>) [Wm^{-2}]
<i>I</i>	radiant intensity [Wsr^{-1}]
<i>I</i>	nuclear spin quantum number [—]
<i>I</i>	current [A]
<i>i</i>	$\sqrt{-1}$
Im()	imaginary part of
<i>J</i>	current density [Am^{-2}]
j	total angular momentum [$\text{kg m}^2 \text{s}^{-1}$]
$J_1()$	Bessel function of the first kind [—]
<i>k</i>	radian wave number $=2\pi/\lambda$ [rad cm^{-1}]
k	wave vector [rad cm^{-1}]
<i>k</i>	extinction coefficient [—]
<i>L</i>	sterance [radiance] [$\text{Wm}^{-2} \text{sr}^{-1}$]
L_v	luminance [cdm^{-2}]
<i>L</i>	inductance [h, $\text{m}^2 \text{kg C}^2$]
<i>L</i>	laser cavity length
<i>L, M, N</i>	direction cosines [—]
<i>M</i>	angular magnification [—]
<i>M</i>	radiant exitance [radiant emittance] [Wm^{-2}]
<i>m</i>	linear magnification [—]
<i>m</i>	effective mass [kg]
MTF	modulation transfer function [—]
<i>N</i>	photon flux [s^{-1}]
<i>N</i>	carrier (number) density [m^{-3}]
<i>n</i>	real part of the relative refractive index [—]
\tilde{n}	complex index of refraction [—]
NA	numerical aperture [—]
OPD	optical path difference [m]
<i>P</i>	macroscopic polarization [C m^{-2}]
Re()	real part of [—]
<i>R</i>	resistance [Ω]
r	position vector [m]
<i>S</i>	Seebeck coefficient [VK^{-1}]
<i>s</i>	spin quantum number [—]
<i>s</i>	path length [m]
S_o	object distance [m]
S_i	image distance [m]
T	temperature [K, C]
<i>t</i>	time [s]
<i>t</i>	thickness [m]
<i>u</i>	slope of ray with the optical axis [rad]
<i>V</i>	Abbe reciprocal dispersion [—]
<i>V</i>	voltage [V , $\text{m}^2 \text{kg s}^{-2} \text{C}^{-1}$]
<i>x, y, z</i>	rectangular coordinates [m]
<i>Z</i>	atomic number [—]

Greek Symbols

α	absorption coefficient [cm^{-1}]
α	(power) absorptance (absorptivity)

ϵ	dielectric coefficient (constant) [—]
ϵ	emittance (emissivity) [—]
ϵ	eccentricity [—]
ϵ_1	Re (ϵ)
ϵ_2	Im (ϵ)
τ	(power) transmittance (transmissivity) [—]
ν	radiation frequency [Hz]
ω	circular frequency = $2\pi\nu$ [rads ⁻¹]
ω	plasma frequency [Hz]
λ	wavelength [μm , nm]
σ	wave number = $1/\lambda$ [cm ⁻¹]
σ	Stefan Boltzmann constant [Wm ⁻² K ⁻¹]
ρ	reflectance (reflectivity) [—]
θ, ϕ	angular coordinates [rad, °]
ξ, η	rectangular spatial frequencies [m ⁻¹ , r ⁻¹]
ϕ	phase [rad, °]
ϕ	lens power [m ⁻²]
Φ	flux [W]
χ	electric susceptibility tensor [—]
Ω	solid angle [sr]

Other

\Re	responsivity
$\exp(x)$	e^x
$\log_a(x)$	log to the base a of x
$\ln(x)$	natural log of x
$\log(x)$	standard log of x : $\log_{10}(x)$
Σ	summation
Π	product
Δ	finite difference
δx	variation in x
dx	total differential
∂x	partial derivative of x
$\delta(x)$	Dirac delta function of x
δ_{ij}	Kronecker delta

REFERENCES

1. Anonymous, *ISO Standards Handbook 2: Units of Measurement*, 2nd ed., International Organization for Standardization, 1982.
2. Anonymous, *Symbols, Units and Nomenclature in Physics*, Document U.I.P. 20, International Union of Pure and Applied Physics, 1978.
3. E. Cohen and B. Taylor, "The Fundamental Physical Constants," *Physics Today*, 9 August 1990.
4. W. J. Smith, *Modern Optical Engineering*, 2nd ed., McGraw-Hill, 1990.
5. J. M. Lloyd, *Thermal Imaging Systems*, Plenum Press, 1972.

William L. Wolfe
 College of Optical Sciences
 University of Arizona
 Tucson, Arizona

PART

1

DESIGN

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

1

TECHNIQUES OF FIRST-ORDER LAYOUT

Warren J. Smith*

*Kaiser Electro-Optics, Inc.
Carlsbad, California*

1.1 GLOSSARY

A, B	scaling constants
d	distance between components
f	focal length
h	image height
I	invariant
j, k	indices
l	axial intercept distance
M	angular magnification
m	linear, lateral magnification
n	refractive index
P	partial dispersion, projection lens diameter
r	radius
S	source or detector linear dimension
SS	secondary spectrum
s	object distance
s'	image distance
t	temperature
u	ray slope
V	Abbe number
y	height above optical axis
α	radiometer field of view, projector field of view
ϕ	component power ($= 1/f$)

*Deceased.

1.2 FIRST-ORDER LAYOUT

First-order layout is the determination of the arrangement of the components of an optical system in order to satisfy the first-order requirements imposed on the system. The term “first-order” means the paraxial image properties: the size of the image, its orientations, its location, and the illumination or brightness of the image. This also implies apertures, f -numbers, fields of view, physical size limitations, and the like. It does not ordinarily include considerations of aberration correction; these are usually third- and higher-order matters, not first-order. However, ordinary chromatic aberration and secondary spectrum are first-order aberrations. Additionally, the first-order layout can have an effect on the Petzval curvature of field, the cost of the optics, the sensitivity to misalignment, and the defocusing effects of temperature changes.

The primary task of first-order layout is to determine the powers and spacings of the system components so that the image is located in the right place and has the right size and orientation. It is not necessary to deal with surface-by-surface ray-tracing here; the concern is with components. “Components” may mean single elements, cemented doublets, or even complex assemblies of many elements. The first-order properties of a component can be described by its Gauss points: the focal points and principal points. For layout purposes, however, the initial work can be done assuming that each component is of zero thickness; then only the component location and its power (or focal length) need be defined.

1.3 RAY-TRACING

The most general way to determine the characteristics of an image is by ray-tracing. As shown in Fig. 1, if an “axial (marginal)” ray is started at the foot (axial intercept) of the object, then an image is located at each place that this ray crosses the axis. The size of the image can be determined by tracing a second, “principal (chief),” ray from the top of the object and passing through the center of the limiting aperture of the system, the “aperture stop;” the intersection height of this ray at the image plane indicates the image size. This size can also be determined from the ratio of the ray slopes of the axial ray at the object and at the image; this yields the magnification $m = u_0/u'_k$; object height times magnification yields the image height.

The ray-tracing equations are

$$y_1 = -l_1 u_1 \quad (1)$$

$$u'_j = u_j - y_j \phi_j \quad (2)$$

$$y_{j+1} = y_j + d_j u'_j \quad (3)$$

$$l'_k = -y_k / u'_k \quad (4)$$

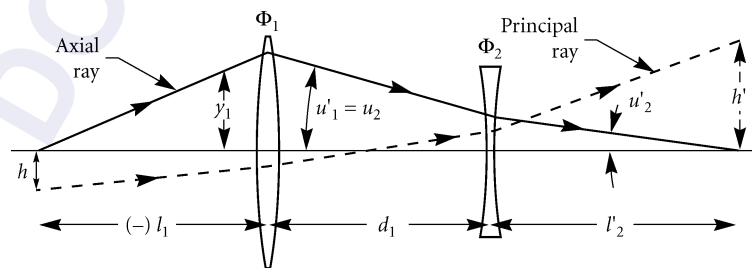


FIGURE 1

where l and l' are the axial intersection distances of the ray before and after refraction by the component, u and u' are the ray slopes before and after refraction, ϕ is the component power ($\phi = 1/f$), y_j is the height at which the ray strikes the j th component, and d_j is the distance from the j th to the $(j + 1)$ th component. Equations (2) and (3) are applied sequentially to the components, from object to image.

These equations can be used in two different ways. When the components and spacings are known, the image characteristics can readily be calculated. In the inverse application, the (unknown) powers and spaces can be represented by symbols, and the ray can be traced symbolically through the postulated number of components. The results of this symbolic ray-tracing can be equated to the required characteristics of the system; these equations can then be solved for the unknowns, which are the component powers and spacings.

As an example, given the starting ray data, y_1 and u_1 , we get

$$\begin{aligned} u'_1 &= u_1 - y_1\phi_1 \\ y_2 &= y_1 + d_1u'_1 = y_1 + d_1(u_1 - y_1\phi_1) \\ u'_2 &= u'_1 - y_2\phi_2 \\ &= u_1 - y_1\phi_1 - [y_1 + d_1(u_1 - y_1\phi_1)]\phi_2 \\ y_3 &= y_2 + d_2u'_2 = \text{etc.} \end{aligned}$$

Obviously the equations can become rather complex in very short order. However, because of the linear characteristics of the paraxial ray equations, they can be simplified by setting either y_1 or u_1 equal to one (1.0) without any loss of generality. But the algebra can still be daunting.

1.4 TWO-COMPONENT SYSTEMS

Many systems are either limited to two components or can be separated into two-component segments. There are relatively simple expressions for solving two-component systems.

Although the figures in this chapter show thick lenses with appropriate principal planes, “thin” lenses (whose thickness is zero and whose principal planes are coincident with the two coincident lens surfaces) may be used.

For systems with infinitely distant objects, as shown in Fig. 2, the following equations for the focal length and focus distance are useful:

$$f_{AB} = f_A f_B / (f_A + f_B - d) \quad (5)$$

$$\phi_{AB} = \phi_A + \phi_B - d\phi_A\phi_B \quad (6)$$

$$B = f_{AB}(f_A - d)/f_A \quad (7)$$

$$F = f_{AB}(f_B - d)/f_B \quad (8)$$

$$h' = f_{AB} \tan u_p \quad (9)$$

where f_{AB} is the focal length of the combination, ϕ_{AB} is its power, f_A and f_B are the focal lengths of the components, ϕ_A and ϕ_B are their powers, d is the spacing between the components, B is the “back

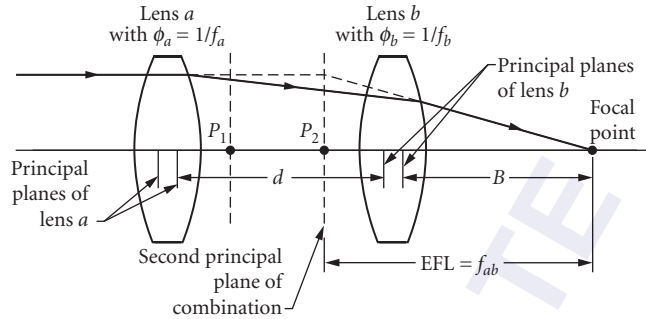


FIGURE 2

focus” distance from the B component, F is the “front focus” distance, u_p is the angle subtended by the object, and h' is the image height.

If f_{AB} , d , and B (or F) are known, the component focal lengths can be found from

$$f_A = df_{AB} / (f_{AB} - B) \tag{10}$$

$$f_B = -dB / (f_{AB} - B - d) \tag{11}$$

These simple expressions are probably the most widely used equations in optical layout work.

If a two-component system operates at *finite* conjugates, as shown in Fig. 3, the following equations can be used to determine the layout. When the required system magnification and the component locations are known, the powers of the components are given by

$$\phi_A = (ms - md - s') / msd \tag{12}$$

$$\phi_B = (d - ms + s') / ds' \tag{13}$$

where $m = h' / h$ is the magnification, s and s' are the object and image distances.

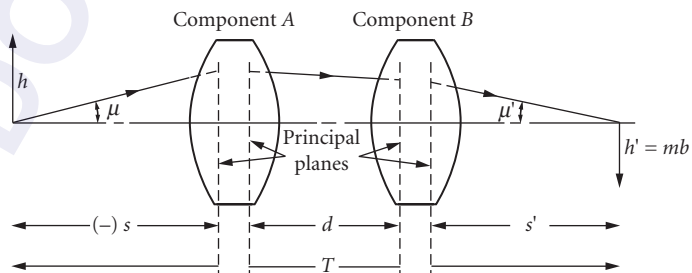


FIGURE 3

In different circumstances, the component powers, the object-to-image distance, and the magnification may be known and the component locations are to be determined. The following quadratic equation [Eq. (14)] in d (the spacing) is solved for d :

$$0 = d^2 - dT + T(f_A + f_B) + (m - 1)^2 f_A f_B / m \quad (14)$$

and then

$$s = [(m - 1)d + T] / [(m - 1) - md\phi_A] \quad (15)$$

$$s' = T + s - d \quad (16)$$

1.5 AFOCAL SYSTEMS

If the system is afocal, then the following relations will apply:

$$MP = -(f_O / f_E) = (u_E / u_O) = (d_O / d_E) \quad (17)$$

and, if the components are “thin,”

$$L = f_O + f_E \quad (18)$$

$$f_O = -L \cdot MP / (1 - MP) \quad (19)$$

$$f_E = L / (1 - MP) \quad (20)$$

where MP is the angular magnification, f_O and f_E are the objective and eyepiece focal lengths, u_E and u_O are the apparent (image) and real (object) angular fields, d_O and d_E are the entrance and exit pupil diameters, and L is the length of the telescope as indicated in Fig. 4.

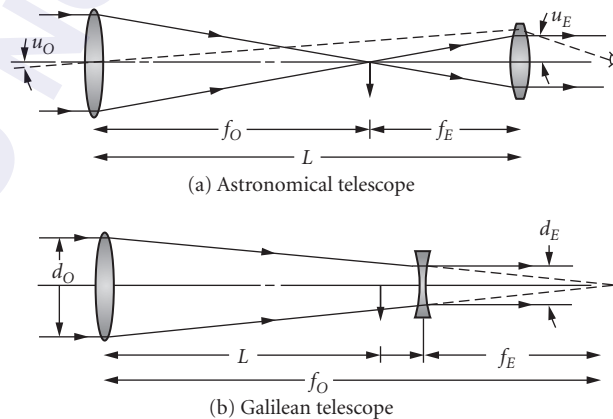


FIGURE 4

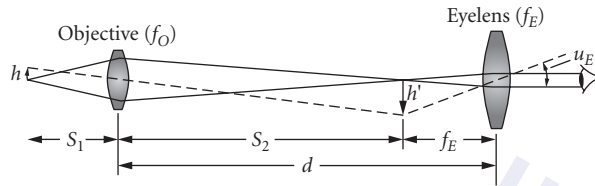


FIGURE 5

1.6 MAGNIFIERS AND MICROSCOPES

The conventional definition of magnifying power for either a magnifier or microscope compares the angular size of the image with the angular size of the object *when the object is viewed from a (conventional) distance of 10 inches*. Thus the magnification can be found from

$$MP = 10''/f \quad (21)$$

for either a simple microscope (i.e., magnifier) or a compound microscope, where f is the focal length of the system. Using the symbols of Fig. 5, we can also write the following for the compound microscope

$$MP = (f_E + f_O - d)10''/f_E f_O \quad (22)$$

$$\begin{aligned} MP &= m_O \times m_E \\ &= (S_2/S_1)(10''/f_E) \end{aligned} \quad (23)$$

1.7 AFOCAL ATTACHMENTS

In addition to functioning as a telescope, beam expander, etc., an afocal system can be used to modify the characteristics of another system. It can change the focal length, power, or field of the “prime” system. Figure 6 shows several examples of an afocal device placed (in these examples) before an imaging system. The combination has a focal length equal to the focal length of the prime system multiplied by the angular magnification of the afocal device. Note that in Fig. 6a and b the same afocal attachment has been reversed to provide two different focal lengths. If the size of the film or detector is kept constant, the angular field is changed by a factor equal to the inverse of the afocal magnification.

1.8 FIELD LENSES

Figure 7 illustrates the function of the field lens in a telescope. It is placed near (but rarely exactly at) an internal image; its power is chosen so that it converges the oblique ray bundle toward the axis sufficiently so that the rays pass through the subsequent component. A field lens is useful to keep the component diameters at reasonable sizes. It acts to relay the pupil image to a more acceptable location.

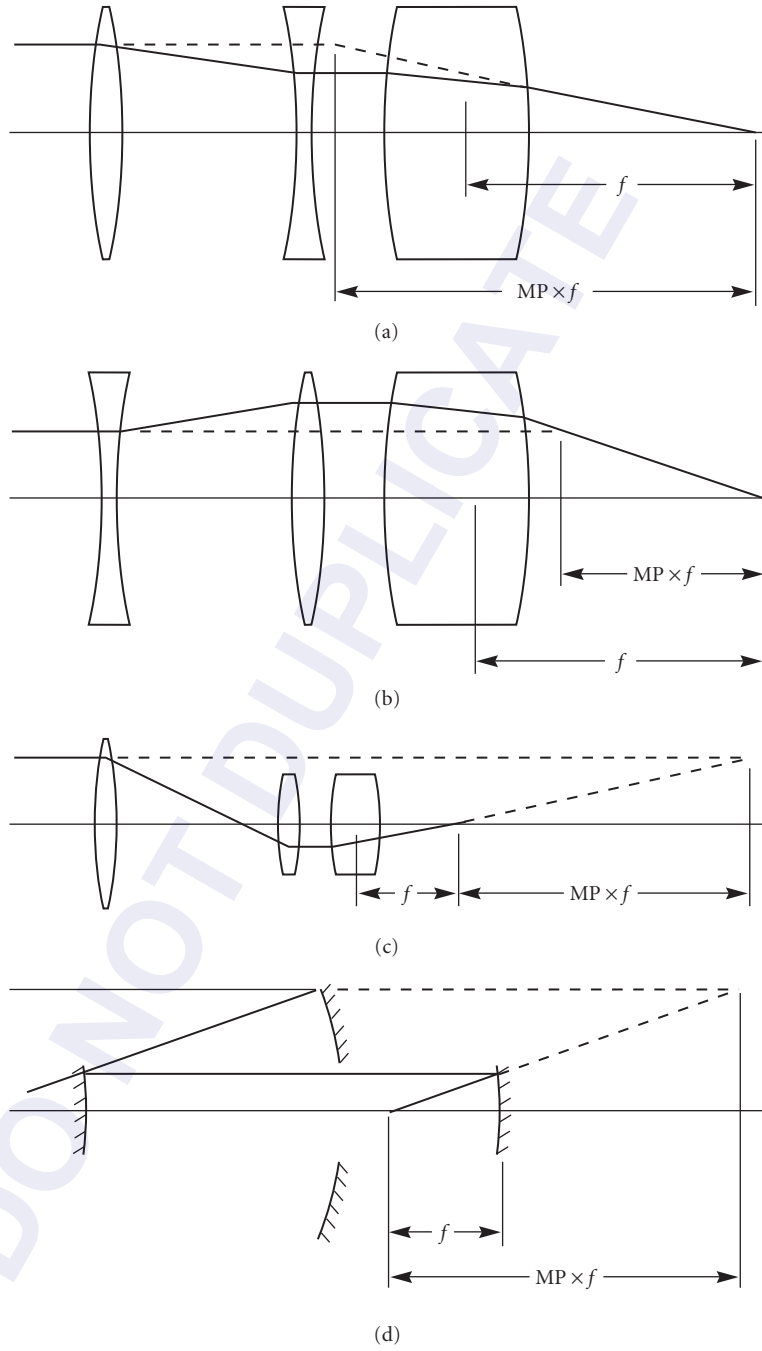


FIGURE 6

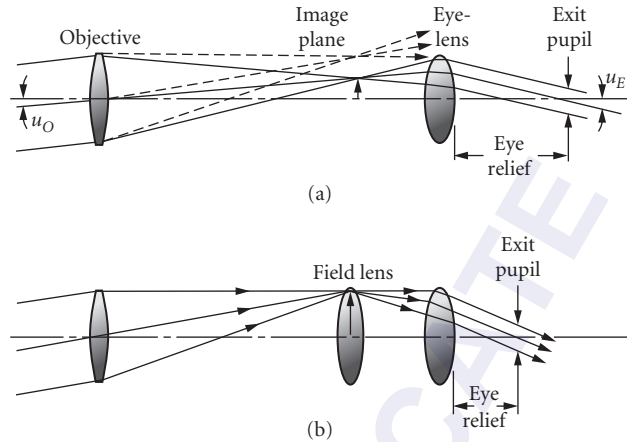


FIGURE 7

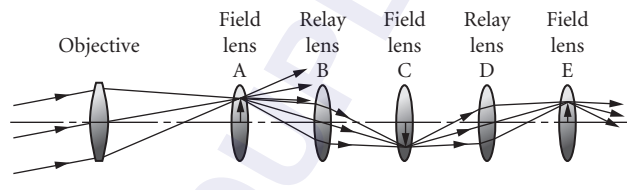


FIGURE 8

The required field lens power is easily determined. In Fig. 7 the most troublesome ray is that from the bottom of the objective aperture; its slope (u) is simply the height that it climbs divided by the distance that it travels. The required slope (u') for the ray after refraction by the field lens is defined by the image height (y), the “eyelens” semidiameter, and the spacing between them. Then Eq. (2) can be solved for the field lens power,

$$\phi = (u - u')/y \quad (24)$$

A periscope is used to carry an image through a long, small-diameter space. As shown in Fig. 8, the elements of a periscope are alternating field lenses and relay lenses. An optimum arrangement occurs when the images at the field lenses and the apertures of the relay lenses are as large as the available space allows. This arrangement has the fewest number of relay stages and the lowest power components. For a space of uniform diameter, both the field lenses and the relay lenses operate at unit magnification.

1.9 CONDENSERS

The projection/illumination condenser and the field lens of a radiation measuring system operate in exactly the same way. The condenser (Fig. 9) forms an image of the light source in the aperture of the projection lens, thereby producing even illumination from a nonuniform source. If the source

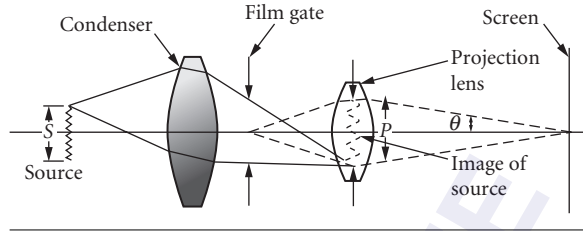


FIGURE 9

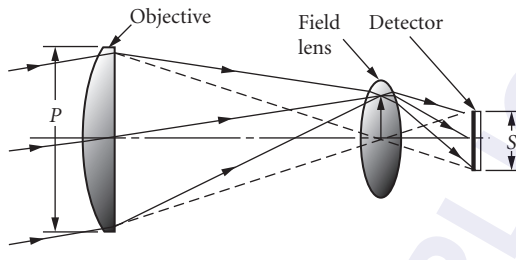


FIGURE 10

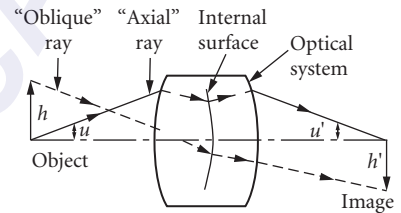


FIGURE 11

image fills the projection lens aperture, this will produce the maximum illumination that the source brightness and the projection lens aperture diameter will allow. This is often called Köhler illumination. In a radiometer type of application (Fig. 10), the field lens images the objective lens aperture on the detector, uniformly illuminating its surface and permitting the use of a smaller detector. Often, the smallest possible source or detector is desired in order to minimize power or maximize signal-to-noise. The smallest possible size is given by

$$S = P\alpha/2n \quad (25)$$

where S is the source or detector size, P is the projection lens or objective aperture diameter, α is the field angle of projection or the radiometer field of view, and n is the index in which the source or detector is immersed. This value for S corresponds to an (impractical) system speed of $F/0.5$. A source or detector size twice as large is a far more realistic limit, corresponding to a speed of $F/1.0$.

The invariant, $I = n(y_2u_1 - y_1u_2)$, where y_1 , u_1 , y_2 , and u_2 are the ray heights and slopes of two different rays, is an expression which has the same value everywhere in an optical system. If the two rays used are an axial ray and a principal (or chief) ray as shown in Fig. 11, and if the invariant is evaluated at the object and image surfaces, the result is

$$hnu = h'n'u' \quad (26)$$

1.10 ZOOM OR VARIFOCAL SYSTEMS

If the spacing between two components is changed, the effective focal length and the back focus are changed in accord with Eqs. (5) through (9). If the motions of the two components are arranged so that the image location is constant, this is a mechanically compensated zoom lens, so called because

the component motions are usually effected with a mechanical cam. A zoom system may consist of just the two basic components or it may include one or more additional members. Usually the two basic components have opposite-signed powers.

If a component is working at unit magnification, it can be moved in one direction or the other to increase or decrease the magnification. There are pairs of positions where the magnifications are m and $1/m$ and for which the object-to-image distance is the same. This is the basis of what is called a “bang-bang” zoom; this is a simple way to provide two different focal lengths (or powers, or fields of view, or magnifications) for a system.

1.11 ADDITIONAL RAYS

When the system layout has been determined, an “axial” ray at full aperture and a “principal” ray at full field can be traced through the system. Because of the linearity of the paraxial equations, we can determine the ray-trace data (i.e., y and u) of *any* third ray from the data of these two traced rays by

$$y_3 = Ay_1 + By_2 \quad (27)$$

$$u_3 = Au_1 + Bu_2 \quad (28)$$

where A and B are scaling constants which can be determined from

$$A = (y_3u_1 - u_3y_1)/(u_1y_2 - y_1u_2) \quad (29)$$

$$B = (u_3y_2 - y_3u_2)/(u_1y_2 - y_1u_2) \quad (30)$$

where y_1 , u_1 , y_2 , and u_2 are the ray heights and slopes of the axial and principal rays and y_3 and u_3 are the data of the third ray; these data are determined at any component of the system where the specifications for all three rays are known. These equations can, for example, be used to determine the necessary component diameters to pass a bundle of rays which are A times the diameter of the axial bundle at a field angle B times the full-field angle. In Fig. 12, for the dashed rays $A = +0.5$ and -0.5 and $B = 1.0$. Another application of Eqs. (27) through (30) is to locate either a pupil or an aperture stop when the other is known.

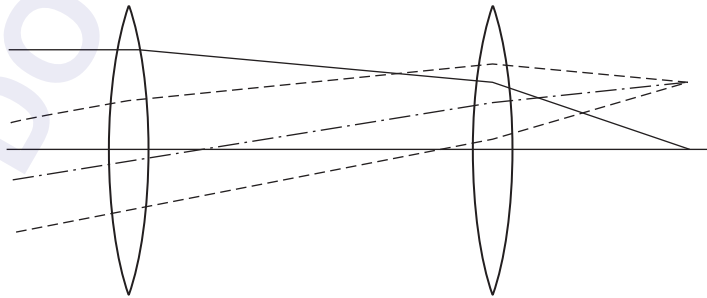


FIGURE 12

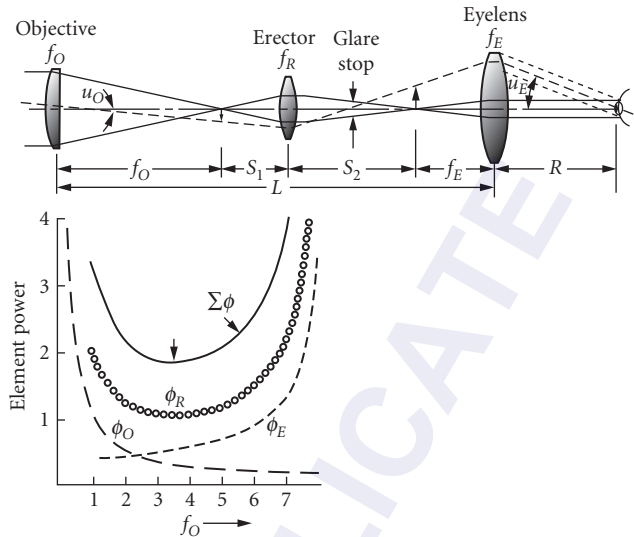


FIGURE 13

1.12 MINIMIZING COMPONENT POWER

The first-order layout may in fact determine the ultimate quality, cost, and manufacturability of the system. The residual aberrations in a system are a function of the component powers, relative apertures, and angular fields. The relationships are complex, but a good choice for a system layout is one which minimizes the sum of the (absolute) component powers, or possibly the sum of the (absolute) $y\phi$ product for all the components.

For example, in Fig. 13 the length, magnification, and the eye relief of the rifle scope are specified. There are five variables: three component powers and two spaces. This is one more variable than is necessary to achieve the specified characteristics. If we take the focal length of the objective component as the free variable, the component powers which satisfy the specifications can be plotted against the objective focal length, as in Fig. 13, and the minimum power arrangement is easily determined.

Minimizing the component powers will strongly tend to minimize the aberrations and also the sensitivity of the system to fabrication errors and misalignments. The *cost* of an optical element will vary with its diameter (or perhaps the square of the diameter) and also with the product of the diameter and the power. Thus, while first-order layout deals only with components, these relationships still apply reasonably well even when applied to components rather than elements. Minimizing the component powers does tend to reduce the cost on these grounds (and also because it tends to reduce the complexity of the components).

1.13 IS IT A REASONABLE LAYOUT?

A simple way to get a feel for the reasonableness of a layout is to make a rough scale drawing showing each component as single element. An element can be drawn as an equiconvex lens with radii which are approximately $r = 2(n - 1)f$; for an element with an index of 1.5 the radii equal the focal length. The elements should be drawn to the diameter necessary to pass the (suitably vignettted) off-axis

bundle of rays as well as the axial bundle. The on-axis and off-axis ray bundles should be sketched in. This will very quickly indicate which elements or components are the difficult ones. If the design is being started from scratch (as opposed to simply combining existing components), each component can be drawn as an achromat. The following section describes achromat layout, but for visual-spectrum systems it is often sufficient to assume that the positive (crown) element has twice the power of the achromat and the (negative) flint element has a power equal to that of the achromat. Thus an achromat may be sketched to the simplified, approximate prescription: $r_1 = -r_2 = f/2$ and $r_3 = \text{plano}$.

Any elements which are too fat must then be divided or “split” until they look “reasonable.” This yields a reasonable estimate of the required complexity of the system, even before the lens design process is begun.

If more or less standard design types are to be utilized for the components, it is useful to tabulate the focal lengths and diameters to get the (infinity) f -number of each component, and also its angular field coverage. The field coverage should be expressed both in terms of the angle that the object and image subtend from the component, and also the angle that the smaller of these two heights subtends as a function of the focal length (rather than as a function of that conjugate distance). This latter angle is useful because the coverage capability of a given design form is usually known in these terms, that is, h/f , rather than in finite conjugate terms. With this information at hand, a reasonable decision can be made as to the design type necessary to perform the function required of the component.

1.14 ACHROMATISM

The powers of the elements of an achromat can be determined from

$$\phi_A = \phi_{AB} V_A / (V_A - V_B) \quad (31)$$

$$\begin{aligned} \phi_B &= \phi_{AB} V_B / (V_B - V_A) \quad (32) \\ &= \phi_{AB} - \phi_A \end{aligned}$$

where ϕ_{AB} is the power of the achromatic doublet and V_A is the Abbe V -value for the element whose power is ϕ_A , etc. For the visible spectral region $V = (n_d - 1)/(n_F - n_C)$; this can be extended to any spectral region by substituting the indices at middle, short, and long wavelengths for n_d , n_F , and n_C .

If the elements are to be spaced apart, and the back focus is B , then the powers and the spacing are given by

$$\phi_A = \phi_{AB} B V_A / (V_A B - V_B / \phi_{AB}) \quad (33)$$

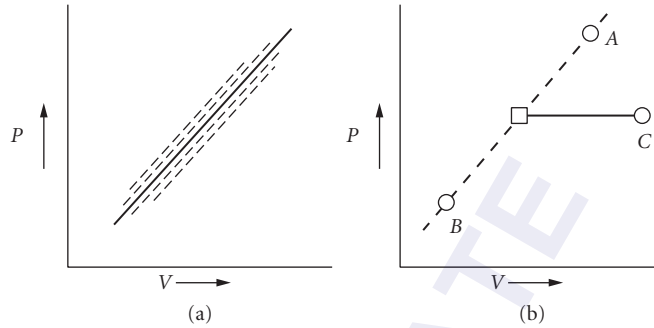
$$\phi_B = -\phi_{AB} V_B / B (V_A B - V_B / \phi_{AB}) \quad (34)$$

$$D = (1 - B \phi_{AB}) / \phi_A \quad (35)$$

For a complete system, the transverse *axial chromatic* aberration is the sum of $y^2 \phi / V u'_k$ for all the elements, where y is the height of the axial ray at the element and u'_k is the ray slope at the image. The *lateral color* is the sum of $y y_p \phi / V u'_k$, where y_p is the principal ray height.

The *secondary spectrum* is the sum of $y^2 \phi P / V u'_k$, where P is the partial dispersion, $P = (n_d - n_c) / (n_F - n_c)$. Summed over two elements, this leads to an expression for the longitudinal secondary spectrum of an achromatic doublet

$$\begin{aligned} \text{SS} &= f(P_B - P_A) / (V_A - V_B) \\ &= -f(\Delta P / \Delta V) \end{aligned} \quad (36)$$


FIGURE 14

This indicates that in order to eliminate secondary spectrum for a doublet, two glasses with identical partial dispersions [so that $(P_A - P_B)$ is zero] are required. A large difference in V -value is desired so that $(V_A - V_B)$ in the denominator of Eqs. (31) and (32) will produce reasonably low element powers. As indicated in the schematic and simplified plot of P versus V in Fig. 14a, most glasses fall into a nearly linear array, and $(\Delta P/\Delta V)$ is nearly a constant for the vast majority of glasses. The few glasses which are away from the "normal" line can be used for apochromats, but the ΔV for glass pairs with a small ΔP tends to be quite small. In order to get an exact match for the partial dispersions so that ΔP is equal to zero, two glasses can be combined to simulate a third, as indicated in Fig. 14b. For a unit power ($\phi = 1$) apochromatic triplet, the element powers can be found from

$$X = [V_A(P_B - P_C) + V_B(P_C - P_A)] / (P_B - P_A) \quad (37)$$

$$\phi_C = V_C / (V_C - X) \quad (38)$$

$$\phi_B = (1 - \phi_C)(P_C - P_A) V_B / [V_B(P_C - P_A) + V_A(P_B - P_C)] \quad (39)$$

$$\phi_A = 1 - \phi_B - \phi_C \quad (40)$$

1.15 ATHERMALIZATION

When the temperature of a lens element is changed, two factors affect its focus or focal length. As the temperature rises, all dimensions of the element are increased; this, by itself, would lengthen the focal length. However, the index of refraction of the lens material also changes with temperature. For many glasses the index rises with temperature; this effect tends to shorten the focal length.

The thermal change in the power of a thin element is given by

$$d\phi/dt = -\phi[a - (dn/dt)/(n-1)] \quad (41)$$

where dn/dt is the differential of index with temperature and a is the thermal expansion coefficient of the lens material. Then for a thin doublet

$$d\phi/dt = \phi_A T_A + \phi_B T_B \quad (42)$$

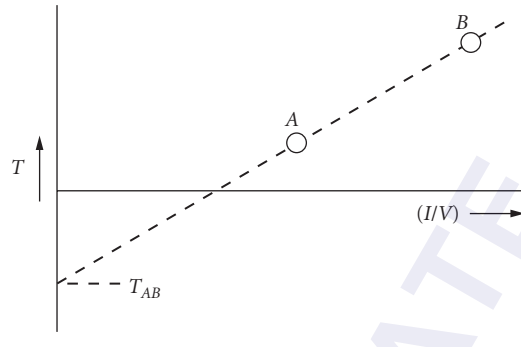


FIGURE 15

where

$$T = [-a + (dn/dt)/(n-1)] \quad (43)$$

and ϕ is the doublet power.

For an athermalized doublet (or one with some desired $d\phi/dt$) the element powers are given by

$$\phi_A = [(d\phi/dt) - \phi T_B] / (T_A - T_B) \quad (44)$$

$$\phi_B = \phi - \phi_A \quad (45)$$

To get an athermalized *achromatic* doublet, a plot of T against $(1/V)$ for all the glasses/materials under consideration is made. A line drawn between two glass points is extended to intersect the T axis as indicated in Fig. 15. Then the value of the $d\phi/dt$ for the achromatic doublet is equal to the doublet power times the value of T at which the line intersects the T axis. A pair of glasses with a large V -value difference and a small or zero T axis intersection is desirable.

An athermal achromatic triplet can be made with three glasses as follows:

$$\phi_A = \phi V_A (T_B V_B - T_C V_C) / D \quad (46)$$

$$\phi_B = \phi V_B (T_C V_C - T_A V_A) / D \quad (47)$$

$$\phi_C = \phi V_C (T_A V_A - T_B V_B) / D \quad (48)$$

$$D = V_A (T_B V_B - T_C V_C) + V_B (T_C V_C - T_A V_A) + V_C (T_A V_A - T_B V_B) \quad (49)$$

See also Chap. 8, "Thermal Compensation Techniques," by Philip J. Rogers and Michael Roberts.

NOTE: Figures 2, 3, 4, 5, 7, 8, 9, 10, 11, and 13 are adapted from W. Smith, *Modern Optical Engineering*, 2nd ed., McGraw-Hill, New York, 1990. The remaining figures are adapted from *Critical Reviews of Lens Design*, W. Smith (Ed.), SPIE, vol. CR41, 1992.

ABERRATION CURVES IN LENS DESIGN

Donald C. O'Shea

*Georgia Institute of Technology
School of Physics
Atlanta, Georgia*

Michael E. Harrigan

*Harrigan Optical Design
Victor, New York*

2.1 GLOSSARY

H	ray height
NA	numerical aperture
OPD	optical path difference
P	petzval
S	sagittal
T	tangential
$\tan U$	slope

2.2 INTRODUCTION

Many optical designers use aberration curves to summarize the state of correction of an optical system, primarily because these curves give a designer important details about the relative contributions of individual aberrations to lens performance. Because a certain design technique may affect only one particular aberration type, these curves are more helpful to the lens designer than a single-value merit function. When a design is finished, the aberration curves serve as a summary of the lens performance and a record for future efforts. For applications such as photography, they are most useful because they provide a quick estimate of the effective blur circle diameter.

The aberration curves can be divided into two types: those that are expressed in terms of ray errors and those in terms of the optical path difference (OPD). OPD plots are usually plotted against the relative ray height in the entrance pupil. Ray errors can be displayed in a number of ways. Either the transverse or longitudinal error of a particular ray relative to the chief ray can be plotted as a function of the ray height in the entrance pupil. Depending upon the amount and type of aberration present, it is sometimes more appropriate to plot the longitudinal aberration as a function of field angle.

For example, astigmatism or field curvature is more easily estimated from field plots, described below. Frequently, the curves are also plotted for several wavelengths to characterize chromatic performance. Because ray error plots are the most commonly used format, this entry will concentrate on them.

2.3 TRANSVERSE RAY PLOTS

These curves can take several different forms, depending on the particular application of the optical system. The most common form is the transverse ray aberration curve. It is also called lateral aberration, or ray intercept curve (also referred to by the misleading term “rim ray plots”). These plots are generated by tracing fans of rays from a specific object point for finite object distances (or a specific field angle for an object at infinity) to a linear array of points across the entrance pupil of the lens. The curves are plots of the ray error at an evaluation plane measured from the chief ray as a function of the relative ray height in the entrance pupil (Fig. 1). For afocal systems, one generally plots angular aberrations, the differences between the tangents of exiting rays and their chief ray in image space.

If the evaluation plane is in the image of a perfect image, there would be no ray error and the curve would be a straight line coincident with the abscissa of the plot. If the curve were plotted for a different evaluation plane parallel to the image plane, the curve would remain a straight line but it would be rotated about the origin. Usually the aberration is plotted along the vertical axis, although some designers plot it along the horizontal axis.

The curves in Fig. 1 indicate a lens with substantial undercorrected spherical aberration as evidenced by the characteristic S-shaped curve. Since a change of the evaluation plane serves only to rotate the curve about the origin, a quick estimate of the aberrations of a lens can be made by reading the scale of the ray error axis (y axis) and mentally rotating the plot. For example, the blur spot can be estimated from the extent of a band that would enclose the curve a in Fig. 1, but a similar estimate could be made from the curves b or c , also.

The simplest form of chromatic aberration is axial color. It is shown in Fig. 2 in the presence of spherical aberration. Axial color is the variation of paraxial focus with wavelength and is seen as a difference in slope of the aberration curves at the origin as a function of wavelength. If the slopes of the curves at the origin for the end wavelengths are different, primary axial color is present. If primary axial color is corrected, then the curves for the end wavelengths will have the same slope at the origin. But if that slope differs from the slope of the curve for the center wavelength, then secondary axial color is present.

A more complex chromatic aberration occurs when the aberrations themselves vary with wavelength. Spherochromatism, the change of spherical aberration with wavelength, manifests itself as a difference in the shapes of the curves for different colors. Another curve that provides a measure of lateral color, an off-axis chromatic aberration, is described below.

For a point on the axis of the optical system, all ray fans lie in the meridional plane and only one plot is needed to evaluate the system. For off-axis object points, a second plot is added to evaluate a fan of skew rays traced in a sagittal plane. Because a skew ray fan is symmetrical across the meridional

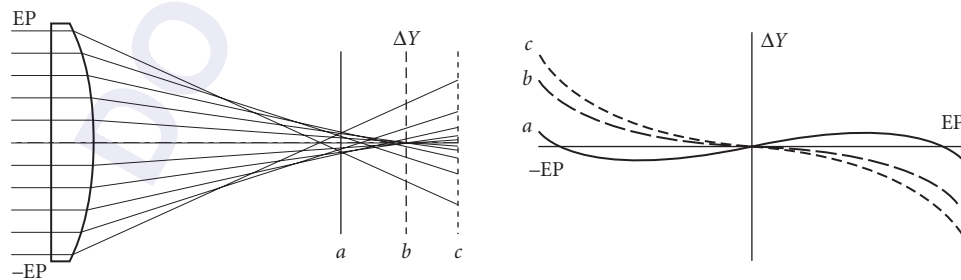


FIGURE 1 (Left) Rays exiting a lens are intercepted at three evaluation planes. (Right) Ray intercept curves plotted for the evaluation planes: (a) at the point of minimum ray error (circle of least confusion); (b) at the paraxial image plane; and (c) outside the paraxial image plane. (See also color insert.)

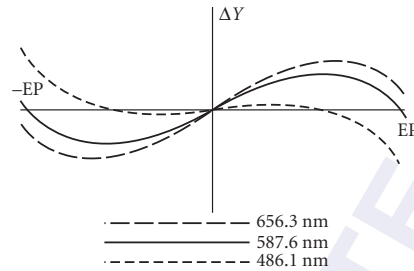


FIGURE 2 Meridional ray intercept curves of a lens with spherical aberration plotted for three colors. (See also color insert.)

plane, only one side of the curve is usually plotted. For all curves the plots are departures from the chief ray location in the evaluation plane (Fig. 3). (In the case of the on-axis point, the chief ray is coincident with the optical axis.) For systems of small-field coverage only two or three object points need to be analyzed, but for wide-angle systems, four or more field points may be necessary.

What can be determined most easily from a comparison between the meridional and sagittal fans is the amount of astigmatism in the image for that field point. When astigmatism is present, the image planes for the tangential and sagittal fans are located at different distances along the chief ray. This is manifested in the ray intercept curves by different slopes at the origin for the tangential and sagittal curves. In Fig. 3 the slopes at the origins of the two curves are different at both 70 percent and full field, indicating astigmatism at both field points. The fact that the difference in the slopes of these two curves has changed sign between the two field points indicates that at some field angle between 70 percent and full field, the slopes are equal and there is no astigmatism there. In addition, the variation of slopes for each curve as a function of field angle is evidence of field curvature.

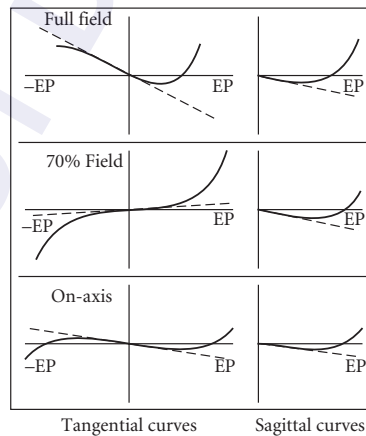


FIGURE 3 Evaluation of a lens on-axis and at two off-axis points. The reduction of the length of the curve with higher field indicates that the lens is vignetting these angles. The differences in slopes (dashed lines) at the origin between the meridional and skew curves indicate that the lens has astigmatism at these field angles. The variation in the slopes with field indicates the presence of field curvature. (See also color insert.)

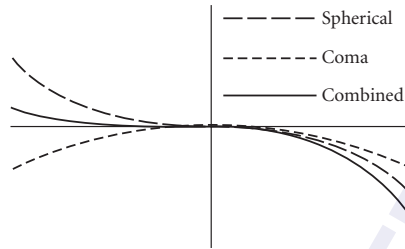


FIGURE 4 Ray intercept curve showing coma combined with spherical aberration. (See also color insert.)

The off-axis aberration of pure primary coma would be evident on these plots as a U-shaped curves for the meridional fan and sagittal fans, the tangential curve being three times larger than the sagittal curve. The “U” will be either upright or upside down depending on the sign of the coma. In almost all cases coma is combined with spherical to produce an S-shaped curve that elongates one of the arms of the “S” and shortens the other (Fig. 4).

The amount of vignetting can be determined from the ray intercept curves also. When it is present, the meridional curves get progressively shorter as the field angle is increased (Fig. 3), since rays at the edges of the entrance pupil are not transmitted. Taken from another perspective, ray intercept curves can also provide the designer with an estimate of how far a system must be stopped down to provide a required degree of correction.

2.4 FIELD PLOTS

The ray intercept curves provide evaluation for a limited number of object points—usually a point on the optical axis and several field points. The field plots present information on certain aberrations across the entire field. In these plots, the independent variable is usually the field angle and is plotted vertically and the aberration is plotted horizontally. The three field plots most often used are: distortion, field curvature, and lateral color. The first of these shows percentage distortion as a function of field angle (Fig. 5).

The second type of plot, field curvature, displays the tangential and sagittal foci as a function of object point or field angle (Fig. 6a). In some plots the Petzval surface, the surface to which the image would



FIGURE 5 Field curve: distortion plot. The percentage distortion is plotted as a function of field angle. Note that the axis of the dependent variable is the horizontal axis. (See also color insert.)

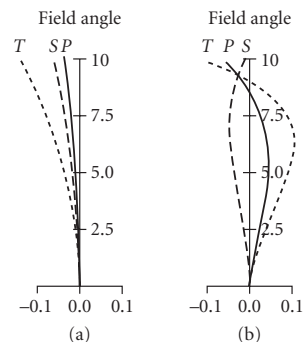


FIGURE 6 Field curve: field curvature plot. The locations of the tangential *T* and sagittal *S* foci are plotted for a full range of field angles. The Petzval surface *P* is also plotted. The tangential surface is always three times farther from the Petzval surface than from the sagittal surface: (a) an uncorrected system and (b) a corrected system. (See also color insert.)

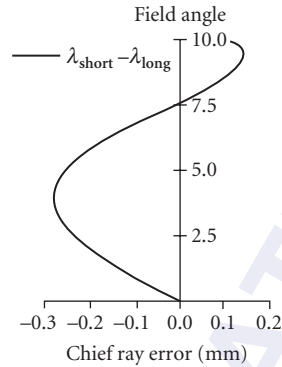


FIGURE 7 Field curve: lateral color plot. A plot of the transverse ray error between red and blue chief ray heights in the image plane for a full range of field angles. Here the distance along the horizontal axis is the color error in the image plane. (See also color insert.)

collapse if there were no astigmatism, is also plotted. This plot shows the amount of curvature in the image plane and amount of astigmatism over the entire field. In cases of corrected field curvature (Fig. 6*b*), this plot provides an estimate of the residual astigmatism between the axis and the corrected zone and an estimate of the maximum field angle at which the image possesses reasonable correction.

The last of the field curves provides information on color error as a function of field angle (Fig. 7). Lateral color, the variation of magnification with wavelength, is plotted as the difference between the chief ray heights at the red and blue wavelengths as a function of field angle. This provides the designer with an estimate of the amount of color separation in the image at various points in the field. In the transverse ray error curves, lateral color is seen as a vertical displacement of the end wavelength curves from the central wavelength curve at the origin.

Although there are other plots that can describe aberrations of optical systems (e.g., plot of longitudinal error as a function of entrance pupil height), the ones described here represent the ensemble that is used in most ray evaluation presentations.

2.5 ADDITIONAL CONSIDERATIONS

In many ray intercept curves the independent variable is the relative entrance pupil coordinate of the ray. However, for systems with high NA or large field of view, where the principal surface cannot be approximated by a plane, it is better to plot the difference between the tangent of the convergence angle of the chosen ray and the tangent of the convergence angle of the chief ray. This is because the curve for a corrected image will remain a straight line in any evaluation plane.¹ When plotted this way, the curves are called *H-tan U* curves.

Shifting the stop of an optical system has no effect on the on-axis curves. However, it causes the origin of the meridional curves of off-axis points to be shifted along the curve. In Fig. 8, the off-axis meridional curves are plotted for three stop positions of a double Gauss lens. The center curve (Fig. 8*b*) is plotted for a symmetrically located stop; the outer curves are plots when the stop is located at lens surfaces before and after the central stop.

It is usually sufficient to make a plot of the aberrations in the meridional and sagittal sections of the beam. The meridional section, defined for an optical system with rotational symmetry, is any plane containing the optical axis. It is sometimes called the tangential section. The sagittal section is a plane perpendicular to the meridional plane containing the chief ray. There are some forms of higher-order

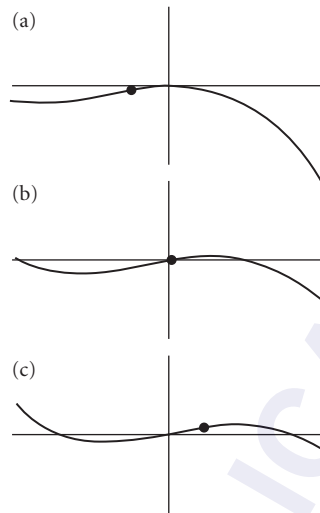


FIGURE 8 The effect of stop shifting on the meridional ray intercept curves of a double Gauss lens. (a) Stop located in front of the normal centrally located stop. (b) Stop at the normal stop position. (c) Stop behind the normal stop position. The dot locates the point on the curve where the origin is located for case (b). (See also color insert.)

coma that do not show in these sections.² In those cases where this aberration is suspected to be a problem, it may be helpful to look at a spot diagram generated from rays in all sections of the bundle.

For a rotationally symmetric system, only objects in a meridional plane need to be analyzed. Also for such systems, only meridional ray errors are possible for purely meridional rays. To observe certain coma types, it is a good idea to plot both the meridional and sagittal ray errors for sagittal rays. It is possible for the meridional section to show no coma and have it show only in the meridional error component of the sagittal fan,² but this aberration is normally small.

In addition to plots of the ray error in an evaluation plane, another aberration plot is one that expresses wavefront aberrations as an optical path difference from a spherical wavefront centered about the image point. These OPD plots are particularly useful for applications where the lens must be close to diffraction-limited.

2.6 SUMMARY

Aberration curves provide experienced designers with the information needed to enable them to correct different types of aberrations. Chromatic effects are much more easily classified from aberration curves also. In comparison to spot diagrams and modulation transfer function curves, the types of aberrations can be more easily seen and quantified. In the case of diffraction-limited systems, modulation transfer functions may provide better estimates of system performance.

2.7 REFERENCES

1. R. Kingslake, *Lens Design Fundamentals*, Academic Press, San Diego, 1978, p. 144.
2. F. D. Cruickshank and G. A. Hills, "Use of Optical Aberration Coefficients in Optical Design," *J. Opt. Soc. Am.* 50:379 (1960).

Douglas C. Sinclair

*Sinclair Optics, Inc.
Fairport, New York*

3.1 GLOSSARY

a	axial ray
b	chief ray
c	curvature
d, e, f, g	aspheric coefficients
efl	effective focal length
FN	focal ratio
$f()$	function
h	ray height
m	linear, lateral magnification
n	refractive index
PIV	paraxial (Lagrange) invariant
r	entrance pupil radius
t	thickness
u	ray slope
y	coordinate
z	coordinate
α	tilt about x (Euler angles)
β	tilt about y (Euler angles)
γ	tilt about z (Euler angles)
ϵ	displacement of a ray from the chief ray
μ	Buchdahl coefficients
κ	conic constant
ρ	radial coordinate
σ_1	spherical aberration

σ_2	coma
σ_3	astigmatism
σ_4	Petzval blur
σ_5	distortion

3.2 INTRODUCTION

The primary function of optical design software is to produce a mathematical description, or *prescription*, describing the shapes, locations, materials, etc., of an optical system that satisfies a given set of specifications. A typical optical design program contains three principal sections: *data entry*, *evaluation*, and *optimization*. The optical design programs considered here are to be distinguished from *ray-trace* programs, which are mainly concerned with evaluation, and *CAD* programs, which are mainly concerned with drawings. The essence of an optical design program is its optimization section, which takes a starting design and produces a new design that minimizes an *error function* that characterizes the system performance.

The first practical computer software for optical design was developed in the 1950s and 1960s.¹⁻⁴ Several commercially available programs were introduced during the 1970s, and development of these programs has continued through the 1980s to the present time. Although decades have passed since the introduction of optical design software, developments continue in optimization algorithms, evaluation methods, and user interfaces.

This chapter attempts to describe a typical optical design program. It is intended for readers that have a general background in optics, but who are not familiar with the capabilities of optical design software. We present a brief description of some of the most important mathematical concepts, but make no attempt to give a detailed development. We hope that this approach will give readers enough understanding to know whether an optical design program will be a useful tool for their own work.

Of course, many different programs are available, each with its own advantages and disadvantages. Our purpose is not to review or explain specific programs, but to concentrate on the basic capabilities. Some programs work better than others, but we make no quality judgment. In fact, we avoid reference, either explicit or implicit, to any particular program. The features and benefits of particular optical design programs are more than adequately described by software vendors, who are listed in optical industry buyer's guides.⁵

Figure 1 is a flowchart of a typical optical design project. Usually, the designer not only must enter the starting design and initial optimization data, but also must continually monitor the progress of the computer, modifying either the lens data or the optimization data as the design progresses to achieve the best solution. Even when the performance requirements are tightly specified, it is often necessary to change the error function during the design process. This occurs when the error function does not correlate with the desired performance sufficiently well, either because it is ill-conceived, or because the designer has purposefully chosen a simple error function to achieve improved speed.

The fact that there are alternate choices of action to be taken when the design is not good enough has led to two schools of thought concerning the design of an optical design program. The first school tries to make the interface between the designer and the program as smooth as possible, emphasizing the interactive side of the process. The second school tries to make the error function comprehensive, and the iteration procedure powerful, minimizing the need for the designer to intervene.

3.3 LENS ENTRY

In early lens design programs, lens entry was a "phase" in which the lens data for a starting design was read into the computer from a deck of cards. At that time, the numerical aspects of optical design on a computer were so amazing that scant attention was paid to the lens entry process. However, as the use of optical design software became more widespread, it was found that a great deal of a designer's time was spent punching cards and submitting new jobs, often to correct past mistakes. Many times, it turned out that the hardest part of a design job was preparing a "correct" lens deck!

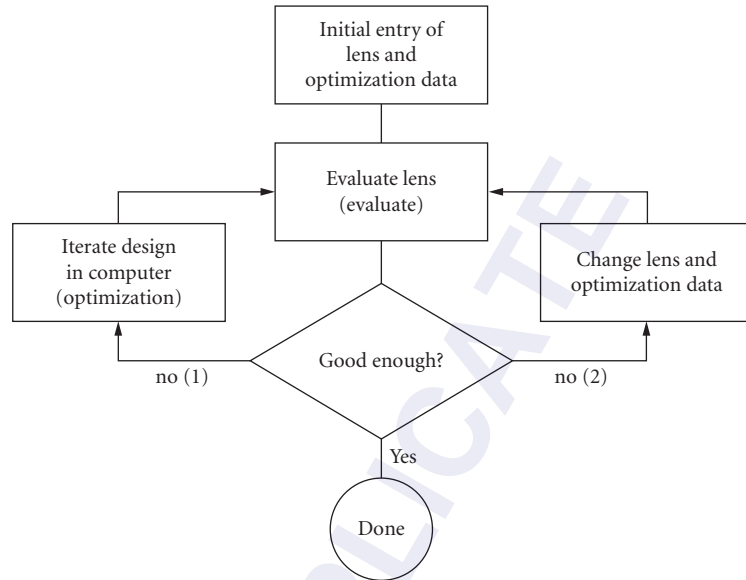


FIGURE 1 Flowchart for the lens design process. The action taken when a design is not satisfactory depends on how bad it is. The designer (or a design program) may change the lens data, or redefine the targets to ones that can be achieved.

Over the years, optical design programs have been expanded to improve the lens entry process, changing the function of this part of the program from simple lens entry to what might be called lens database management. A typical contemporary program provides on-line access to a library of hundreds of lenses, interactive editing, automatic lens drawings, and many features designed to simplify this aspect of optical design.

The lens database contains all items needed to describe the optical system under study, including not only the physical data needed to construct the system (curvatures, thicknesses, etc.), but also data that describe the conditions of use (object and image location, field of view, etc.). Some programs also incorporate optimization data in the lens database, while others provide separate routines for handling such data. In any case, the lens database is often the largest part of an optical design program.

The management of lens data in an optical design program is complicated by two factors. One is that there is a tremendous range of complexity in the types of systems that can be accommodated, so there are many different data items. The other is that the data are often described indirectly. A surface curvature may be specified, for example, by the required slope of a ray that emerges from the surface, rather than the actual curvature itself. Such a specification is called a *solve*, and is based on the fact that paraxial ray tracing is incorporated in the lens entry portion of most optical design programs.

It might seem curious that paraxial ray tracing is still used in contemporary optical design programs that can trace *exact* rays in microseconds. (The term *exact* ray is used in this chapter to mean a real skew ray. Meridional rays are treated as a special case of skew rays in contemporary software; there is not sufficient computational advantage to warrant a separate ray trace for them.) In fact, paraxial rays have some important properties that account for their incorporation in the lens database.

First, paraxial rays provide a linear system model that leads to analysis of optical systems in terms of simple bilinear transforms. Second, paraxial ray tracing does not fail. Exact rays can miss surfaces or undergo total internal reflection. Finally, paraxial rays determine the ideal performance of a lens. In a well-corrected lens, the aberrations are balanced so that the exact rays come to the image points defined by the paraxial rays, not the other way around.

Two basic types of data are used to describe optical systems. The first are the *general* data that are used to describe the system as a whole, and the other are the *surface* data that describe the individual surfaces and their locations. Usually, an optical system is described as an ordered set of surfaces,

beginning with an *object* surface and ending with an *image* surface (where there may or may not be an actual image). It is assumed that the designer knows the order in which rays strike the various surfaces. Systems for which this is not the case are said to contain *nonsequential* surfaces, which are discussed later.

General System Data

The general data used to describe a system include the aperture and field of view, the wavelengths at which the system is to be evaluated, and perhaps other data that specify evaluation modes, vignetting conditions, etc.

Aperture and Field of View The aperture and field of view of a system determine its conditions of use. The aperture is specified by the *axial ray*, which emerges from the vertex of the object surface and passes through the edge of the entrance pupil. The field of view is specified by the *chief ray*, which emerges from the edge of the object and passes through the center of the entrance pupil.

There are various common ways to provide data for the axial and chief rays. If the object is at an infinite distance, the entrance pupil radius and semifield angle form a convenient way to specify the axial and chief rays. For finite conjugates, the numerical aperture in object space and the object height are usually more convenient.

Some programs permit the specification of paraxial ray data by image-space quantities such as the *f*-number and the image height, but such a specification is less desirable from a computational point of view because it requires an iterative process to determine initial ray-aiming data.

Wavelengths It is necessary to specify the wavelengths to be used to evaluate polychromatic systems. Three wavelengths are needed to enable the calculation of primary and secondary chromatic aberrations. More than three wavelengths are required to provide an accurate evaluation of a typical system, and many programs provide additional wavelengths for this reason. There has been little standardization of wavelength specification. Some programs assume that the first wavelength is the central wavelength, while others assume that it is one of the extreme wavelengths; some require wavelengths in micrometers, while others in nanometers.

Other General Data Several other items of general data are needed to furnish a complete lens description, but there is little consistency between programs on how these items are treated, or even what they are. The only one that warrants mention here is the aperture stop. The *aperture stop* is usually defined to be the surface whose aperture limits the angle of the axial ray. Once the aperture stop surface is given, the positions of the paraxial pupils are determined by the imaging properties of the system. Since the aperture and field of view are determined formally by the paraxial pupils, the apertures are not associated with the exact ray behavior.

The “vignetting factor” is used to account for the differences between paraxial and exact off-axis ray heights at apertures. In particular, the vignetting factor provides, in terms of fractional (paraxial) coordinates, the data for an exact ray that grazes the apertures of a system. Typically, there is an upper, lower, and skew vignetting factor. The details of how such factors are defined and handled are program dependent.

Surface Data

Surface Location There are two basic ways to specify the location of surfaces that make up a lens. One is to specify the position of a surface relative to the immediately preceding surface. The other is to specify its position relative to some fixed surface (for example, the first surface). The two ways lead to what are called *local* and *global* coordinates, respectively. For ordinary lenses consisting of a series of rotationally symmetric surfaces centered on an optical axis, local coordinates are more convenient, but for systems that include reflectors, tilted, and/or decentered surfaces, etc., global

coordinates are simpler. Internally, optical design programs convert global surface data to local coordinates for speed in ray tracing.

Most optical design programs use a standard coordinate system and standard sign conventions, although there are exceptions.⁶ Each surface defines a local right-handed coordinate system in which the z axis is the symmetry axis and the yz plane is the meridional plane. The local coordinate system is used to describe the surface under consideration and also the origin of the next coordinate system. Tilted elements are described by an Euler-angle system in which α is a tilt around the x axis, β is a tilt around the y axis, and γ is a tilt around the z axis. Since tilting and decentering operations are not commutative; some data item must be provided to indicate which comes first.

Surface Profile Of the various surfaces used in optical systems, the most common by far is the rotationally symmetric surface, which can be written as⁷

$$z = \frac{cr^2}{1 + \sqrt{1 - c^2(\kappa + 1)r^2}} + dr^4 + er^6 + fr^8 + gr^{10}$$

$$r = \sqrt{x^2 + y^2}$$

c is the curvature of the surface; κ is the conic constant; and d , e , f , and g are aspheric constants. The use of the above equation is almost universal in optical design programs. The description of conic surfaces in terms of a conic constant κ instead of the eccentricity e used in the standard mathematical literature allows spherical surfaces to be specified as those with no conic constant. (The conic constant is minus the square of the eccentricity.)

Although aspheric surfaces include all surfaces that are not spherical, from a design standpoint there is a demarcation between “conic” aspheres and “polynomial” aspheres described using the coefficients d , e , f , and g . Rays can be traced analytically through the former, while the latter require numerical iterative methods.

Many optical design programs can handle surface profiles that are more complicated than the above, including cylinders, torics, splines, and even general aspheres of the form $z = f(x, y)$, where $f(x, y)$ is an arbitrary polynomial. The general operation of an optical design program, however, can be understood by considering only the rotationally symmetric surfaces described here.

As mentioned above, the importance of paraxial rays in optical system design has led to the indirect specification of lens data, using *solves*, as they are called, which permit a designer to incorporate the basic paraxial data describing a lens with the lens itself, rather than having to compute and optimize the paraxial performance of a lens as a separate task. Considering the j th surface of an optical system, let

- y_j = ray height on surface
- u_j = ray slope on image side
- c_j = curvature of surface
- n_j = refractive index on image side
- t_j = thickness on image side

The paraxial ray trace equations can then be written as⁸

$$y_j = y_{j-1} + t_{j-1}u_{j-1}$$

$$n_j u_j = n_{j-1} u_{j-1} - y_j c_j (n_j - n_{j-1})$$

These equations can be inverted to give the curvatures and thicknesses in terms of the ray data. We have

$$c_j = \frac{n_{j-1}u_{j-1} - n_j u_j}{y_j(n_j - n_{j-1})}$$

$$t_j = \frac{y_{j+1} - y_j}{u_j}$$

The specification of curvatures and thicknesses by solves is considered to be on an equal basis with the direct specification of these items. The terminology used to specify solves is that the solves used to determine thickness are called *height solves*, and the solves used to determine curvature are called *angle solves*. Often, an axial ray height solve on the last surface is used to automatically locate the paraxial image plane, a chief ray height solve on the same surface to locate the exit pupil, and an axial ray angle solve is used to maintain a given focal length (if the entrance pupil radius is fixed). In some programs, additional types of solves are allowed, such as center of curvature solves, or aperture solves.

Of course, specifying lens data in terms of paraxial ray data means that whenever any lens data is changed, two paraxial rays must be traced through the system to resolve any following data that are determined by a solve. In an optical design program, this function is performed by a *lens setup* routine, which must be efficiently coded, since it is executed thousands of times in even a small design project.

Other functions of the lens setup routine are to precalculate values that are needed for repetitive calculations, such as refractive indices, rotation and translation matrices, etc. Many programs have the capability of specifying certain data items to be equal to (\pm) the value of the corresponding item on a previous surface. These are called *pickups*, and are needed for optimization of systems containing mirrors, as well as maintaining special geometrical relationships. Programs that lack pickups usually have an alternate means for maintaining the required linking between data items. Like solves, pickups are resolved by the lens setup routine, although they do not use paraxial data.

Other Surface Data A variety of other data is required to specify surfaces. Most important are apertures, discussed below, and refractive indices. Refractive indices are usually given by specifying the name of a catalog glass. In the lens setup routine, the actual refractive indices are calculated using an index interpolation formula and coefficient supplied by the glass manufacturer, together with the design wavelengths stored with the lens data. Other surface-related items include phase data for diffractive surfaces, gradient-index data, holographic construction data, and coatings.

Apertures have a somewhat obscure status in many optical design programs. Although apertures have a major role to play in determining the performance of a typical system, they do not usually appear directly in optimization functions. Instead, apertures are usually controlled in optimization by targets on the heights of rays that define their edges. If an aperture is specified directly, it will block rays that pass outside of it and cause typical optimization procedures to become unstable. Accordingly, some programs ignore apertures during optimization. Other programs allow the apertures to be determined by a set of exact “reference rays” that graze their extremities.

Nonsequential Surfaces In some optical systems, it is not possible to specify the order in which a ray will intersect the surfaces as it progresses through the system. The most common examples of such systems are prisms such as the corner-cube reflector, where the ordering of surfaces depends on the entering ray coordinates. Other examples of nonsequential surfaces include light pipes and a variety of nonimaging concentrators. Nonsequential surfaces can be accommodated by many optical design programs, but for the most part they are not “designed” using the program, but rather are included as a subsystem used in conjunction with another part of the system that is the actual

system being designed. Data specification of nonsequential surfaces is more complicated than ordinary systems, and ray tracing is much slower, since several surfaces must be investigated to see which surface is the one actually traversed by a given ray.

Lens Setup

Whenever the lens entry process is completed, the lens must be “set up.” Pickup constraints must be resolved. If the system contains an internal aperture stop, the position of the entrance pupil must be determined. Then paraxial axial and chief rays must be traced through the system so that surface data specified by solves can be computed. Depending on the program, a variety of other data may be precomputed for later use, including aperture radii, refractive indices, and various paraxial constants.

The lens setup routine must be very efficient, since it is the most heavily used code in an optical design program. In addition to running whenever explicit data entry is complete, the code is also executed whenever the lens is modified internally by the program, such as when derivatives are computed during optimization, or when configurations are changed in a multiconfiguration system. Typically, lens setup takes milliseconds (at most), so it is not noticed by the user, other than through its effects.

Programming Considerations

In writing an optical design program, the programmer must make a number of compromises between speed, size, accuracy, and ease of use. These compromises affect the usefulness of a particular program for a particular application. For example, a simple, fast, small program may be well suited to a casual user with a simple problem to solve, but this same program may not be suited for an experienced designer who routinely tackles state-of-the-art problems.

The lens entry portion of an optical design program shows, more than any other part, the difference in programming models that occurred during the 1980s. Before the 1980s, most application programs were of a type called *procedural* programs. When such a program needs data, it requests it, perhaps from a file or by issuing a prompt for keyboard input. The type of data needed is known, and the program is only prepared to accept that kind of data at any given point. Although the program may branch through different paths in response to the data it receives, the program is responsible for its own execution.

With the popularization in the 1980s of computer systems that use a mouse for input the model for an application program changed from the procedural model described above to what is called an *event-driven* model. An event-driven program has a very simple top-level structure consisting of an initialization section followed by an infinite loop usually called something like the *main event loop*. The function of the main event loop is to react to user-initiated interrupts (such as pressing a key, or clicking a mouse button), dispatching such events to appropriate processing functions. In such a program, the user controls the execution, unlike a procedural program, where the execution controls the user.

An event-driven program usually provides a better user interface than a procedural program. Unfortunately, most optical design programs were originally written as procedural programs, and it is difficult to convert a procedural program into an event-driven program by “patching” it. Usually it is easier to start over. In addition, it is harder to write an event-driven program than a procedural program, because the event-driven program must be set up to handle a wide variety of unpredictable requests received at random times. Of course, it is this very fact that makes the user interface better. There is an aphorism sometimes called the “conservation of complexity,” which states that the simpler a program is to use, the more complicated the program itself must be.

The data structures used to define lens data in an optical design program may have a major impact on its capabilities. For example, for various reasons it is usually desirable to represent a lens in a computer program as an array of surfaces. If the maximum size of the array is determined at compile time, then the maximum size lens that can be accommodated is built into the program.

As more data items are added to the surface data, the space required for storage can become unwieldy. For example, it takes about 10 items of real data to specify a holographic surface. If every surface were allowed to be a hologram, then 10 array elements would have to be reserved for each surface's holographic data. On the other hand, in most systems, the elements would never be used, so the data structure would be very inefficient. To avoid this, a more complicated data structure could be implemented in which only one element would be devoted to holograms, and this item would be used as an index into a separate array containing the actual holographic data. Such an array might have a sufficient number of elements to accommodate up to, say, five holograms, the maximum number expected in any one system.

The preceding is a simple example of how the data structure in an optical design program can grow in complexity. In fact, in a large optical design program the data structure may contain all sorts of indices, pointers, flags, etc., used to implement special data types and control their use. Managing this data while maintaining its integrity is a programming task of a magnitude often greater than the numerical processing that takes place during optical design.

Consider, for example, the task of deleting a surface from a lens. To do this, the surface data must of course be deleted, and all of the higher-numbered surfaces renumbered. But, in addition, the surface must be checked to see whether it is a hologram and, if so, the holographic data must also be deleted and that data structure "cleaned up." All other possible special data items must be tested and handled similarly. Then all the renumbered surfaces must be checked to see if any of the "pick up" data from a surface that has been renumbered, and the reference adjusted accordingly. Then other data structures such as the optimization files must be checked to see if they refer to any of the renumbered surfaces, and appropriate adjustments made. There may be several other checks and adjustments that must also be carried out.

Related to the lens entry process is the method used to store lens data on disc. Of course, lens data are originally provided to a program in the form of text (e.g., "TH 1.0"). The program *parses* this data to identify its type (a thickness) and value (1.0). The results of the parsing process (the binary values) are stored in appropriate memory locations (arrays). To store lens data on disc, early optical design programs used the binary data to avoid having to reparse it when it was recovered. However, the use of binary files has decreased markedly as computers have become fast enough that parsing lens input does not take long. The disadvantages of binary files are that they tend to be quite large, and usually have a structure that makes them obsolete when the internal data structure of the program is changed. The alternative is to store lens data as text files, similar in form to ordinary keyboard input files.

3.4 EVALUATION

Paraxial Analysis

Although the lens setup routine contains a paraxial ray trace, a separate paraxial ray trace routine is used to compute data for display to the user. At a minimum, the paraxial ray heights and slopes of the axial and chief ray are shown for each surface, in each color, and in each configuration.

The equations used for paraxial ray tracing were described in the previous section. Although such equations become exact only for "true" paraxial rays that are infinitesimally displaced from the optical axis, it is customary to consider paraxial ray data to describe "formal" paraxial rays that refract at the tangent planes to surfaces, as shown in Fig. 2. Here, the ray ABC is a paraxial ray that provides a first-order approximation to the exact ray ADE. Not only does the paraxial ray refract at the (imaginary) tangent plane BVP, but also it bends a different amount from the exact ray.

In addition to the computation of ray heights and slopes for the axial and chief ray, various paraxial constants that characterize the overall system are computed. The particular values computed depend on whether the system is *focal* (finite image distance) or *afocal* (image at infinity). For focal systems, the quantities of interest are (at a minimum) the focal length f , the f -number FN, the paraxial (Lagrange) invariant PIV, and the transverse magnification m . It is desirable to compute such

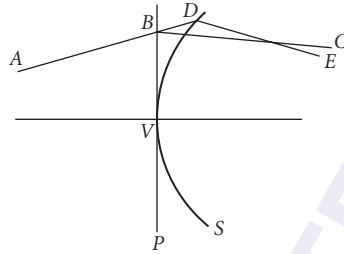


FIGURE 2 Showing the difference between a paraxial ray and a real ray. The paraxial ray propagates along ABC , while the real ray propagates along ADE .

quantities in a way that does not depend on the position of the final image surface. Let the object height be h , the entrance pupil radius be r , the axial ray data in object and image spaces be y , u and y' , u' , the chief ray data be \bar{y} , \bar{u} and \bar{y}' , \bar{u}' , and the refractive indices be n and n' .

The above-mentioned paraxial constants are then given by

$$\text{efl} = \frac{-rh}{\bar{u}'r + u'h}$$

$$\text{FN} = -\frac{n}{2n'u'}$$

$$\text{PIV} = n(\bar{y}u - y\bar{u})$$

$$m = \frac{nu}{n'u'}$$

In addition to the paraxial constants, most programs display the locations of the entrance and exit pupils, which are easily determined using chief-ray data. Surprisingly, most optical design programs do not explicitly show the locations of the principal planes. In addition, although most programs have the capability to display “ $y - \bar{y}$ ” plots, few have integrated this method into the main lens entry routine.

Aberrations

Although most optical designs are based on exact ray data, virtually all programs have the capability to compute and display first-order chromatic aberrations and third-order monochromatic (Seidel) aberrations. Many programs can compute fifth-order aberrations as well. The form in which aberrations are displayed depends on the program and the type of system under study, but as a general rule, for focal systems aberrations are displayed as equivalent ray displacements in the paraxial image plane.

In the case of the chromatic aberrations, the primary and secondary chromatic aberration of the axial and chief rays are computed. In a system for which three wavelengths are defined, the primary aberration is usually taken between the two outer wavelengths, and the secondary aberration between the central and short wavelengths.

The Seidel aberrations are computed according to the usual aberration polynomial. If we let ϵ be the displacement of a ray from the chief ray, then

$$\epsilon_y = \epsilon_{3y} + \epsilon_{5y} + \dots$$

$$\epsilon_x = \epsilon_{3x} + \epsilon_{5x} + \dots$$

For a relative field height h and normalized entrance pupil coordinates r and θ , the third-order terms are

$$\epsilon_{3y} = \sigma_1 \cos \theta r^3 + \sigma_2 (2 + \cos 2\theta) r^2 h + (3\sigma_3 + \sigma_4) \cos \theta r h^2 + \sigma_5 h^3$$

$$\epsilon_{3x} = \sigma_1 \sin \theta r^3 + \sigma_2 \sin 2\theta r^2 h + (\sigma_3 + \sigma_4) \sin \theta r h^2$$

The interpretation of the coefficients is generally as follows, but several optical design programs display tangential coma, rather than the sagittal coma indicated in the table.

σ_1	Spherical aberration
σ_2	Coma
σ_3	Astigmatism
σ_4	Petzval blur
σ_5	Distortion

The fifth-order terms are

$$\begin{aligned} \epsilon_{5y} = & \mu_1 \cos \theta r^5 + (\mu_2 + \mu_3 \cos 2\theta) r^4 h + (\mu_4 + \mu_6 \cos^2 \theta) \cos \theta r^2 h^2 \\ & + (\mu_7 + \mu_8 \cos 2\theta) r^2 h^3 + \mu_{10} \cos \theta r h^4 + \mu_{12} h^5 \end{aligned}$$

$$\begin{aligned} \epsilon_{5x} = & \mu_1 \sin \theta r^5 + \mu_3 \sin 2\theta r^4 h + (\mu_5 + \mu_6 \cos^2 \theta) \sin \theta r^3 h^2 \\ & + \mu_9 \sin 2\theta r^2 h^3 + \mu_{11} \sin \theta r h^4 \end{aligned}$$

These equations express the fifth-order aberration in terms of the Buchdahl μ coefficients. In systems for which the third-order aberrations are corrected, the following identities exist:

$$\mu_2 = \frac{3}{2} \mu_3$$

$$\mu_4 = \mu_5 + \mu_6$$

$$\mu_7 = \mu_8 + \mu_9$$

μ_1	Spherical aberration
μ_3	Coma
$(\mu_{10} - \mu_{11})/4$	Astigmatism
$(5\mu_{11} - \mu_{10})/4$	Petzval blur

$\mu_4 + \mu_6$	Tangential oblique spherical aberration
μ_5	Sagittal oblique spherical aberration
$\mu_7 + \mu_8$	Tangential elliptical coma
μ_9	Sagittal elliptical coma
μ_{12}	Distortion

Some programs display only the aberrations that have corresponding third-order coefficients, omitting oblique spherical aberration and elliptical coma.

The formulas needed to calculate the chromatic and third-order aberrations are given in the *U.S. Military Handbook of Optical Design*. The formulas for calculating the fifth-order aberrations are given in Buchdahl's book.⁹

Aberration coefficients are useful in optical design because they characterize the system in terms of its symmetries, allow the overall performance to be expressed as a sum of surface contributions, and are calculated quickly. On the negative side, aberration coefficients are not valid for systems that have tilted and decentered elements for systems that cover an appreciable field of view, and the accuracy of aberration coefficients in predicting performance is usually inadequate. Moreover, for systems that include unusual elements like diffractive surfaces and gradient index materials, the computation of aberration coefficients is cumbersome at best.

Ray Tracing

Exact ray tracing is the foundation of an optical design program, serving as a base for both evaluation and optimization. From the programmer's standpoint, the exact ray-trace routines must be accurate and efficient. From the user's viewpoint, the data produced by the ray-trace routines must be accurate and comprehensible. Misunderstanding the meaning of ray-trace results can be the source of costly errors in design.

To trace rays in an optical design program, it is necessary to understand how exact rays are specified. Although the details may vary from one program to the next, many programs define a ray by a two-step process. In the first step, an object point is specified. Once this has been done, all rays are assumed to originate from this point until a new object point is specified. The rays themselves are then specified by aperture coordinates and wavelength.

Exact ray starting data is usually normalized to the object and pupil coordinates specified by the axial and chief rays. That is, the aperture coordinates of a ray are specified as a fractional number, with 0.0 representing a point on the vertex of the entrance pupil, and 1.0 representing the edge of the pupil. Field angles or object heights are similarly described, with 0.0 being a point on the axis, and 1.0 being a point at the edge of the field of view.

Although the above normalization is useful when the object plane is at infinity, it is not so good when the object is at a finite distance and the numerical aperture in object space is appreciable. Then, fractional aperture coordinates should be chosen proportional to the direction cosines of rays leaving an object point. There are two reasons for this. One is that it allows an object point to be considered a point source, so that the amount of energy is proportional to the "area" on the entrance pupil. The other is that for systems without pupil aberrations, the fractional coordinates on the second principal surface should be the same as those on the first principal surface. Notwithstanding these requirements, many optical design programs do not define fractional coordinates proportional to direction cosines.

It is sometimes a point of confusion that the aperture and field of view of a system are specified by paraxial quantities, when the actual performance is determined by exact rays. In fact, the paraxial specifications merely establish a normalization for exact ray data. For example, in a real system the field of view is determined not by the angle of the paraxial chief ray, but by the angle at which exact rays blocked by actual apertures just fail to pass through the system. Using an iterative procedure, it is not too hard to find this angle, but because of the nonlinear behavior of Snell's law, it does not provide a convenient reference point.

There are two types of exact rays: *ordinary* or *lagrangian* rays, and *iterated* or *hamiltonian* rays. The designation of rays as lagrangian or hamiltonian comes from the analogy to the equations of motion of a particle in classical mechanics. Here we use the more common designation as ordinary or iterated rays. An ordinary ray is a ray that starts from a known object point in a known direction. An iterated ray also starts from a known object point, but its direction is not known at the start. Instead, it is known that the ray passes through some known (nonconjugate) point inside the system, and the initial ray direction is determined by an iterative procedure.

Iterated rays have several applications in optical design programs. For example, whenever a new object point is specified, it is common to trace an iterated ray through the center of the aperture stop (or some other point) to serve as a reference ray, or to trace several iterated rays through the edges of limiting apertures to serve as reference rays. In fact, many programs use the term *reference ray* to mean iterated ray (although in others, reference rays are ordinary rays). Iterated rays are traced using differentially displaced rays to compute corrections to the initial ray directions. Because of this, they are traced slower than ordinary rays. On the other hand, they carry more information in the form of the differentials, which is useful for computing ancillary data like field sags.

Reference rays are used as base rays in the interpretation of ordinary ray data. For example, the term *ray displacement* often refers to the difference in coordinates on the image surface of a ray from those of the reference ray. Similarly, the *optical path difference* of a ray may compare its phase length to that of the corresponding reference ray. The qualifications expressed in the preceding sentences indicate that the definitions are not universal. Indeed, although the terms *ray displacement* and *optical path difference* are very commonly used in optical design, they are not precisely defined, nor can they be. Let us consider, for example, the optical path difference.

Imagine a monochromatic wavefront from a specified object point that passes through an optical system. Figure 3 shows the wavefront PE emerging in image space, where it is labeled “actual wavefront.” Because of aberrations, an ordinary ray perpendicular to the actual wavefront will not intersect the final image surface at the ideal image point I , but at some other point Q . The optical path difference (OPD) may be defined as the optical path measured along the actual ray between the actual wavefront and a reference sphere centered on the ideal image point.

Unfortunately, the ideal image point is not precisely defined. In the figure, it is shown as the intersection of the reference ray with the image surface, but the reference ray itself may not be precisely defined.

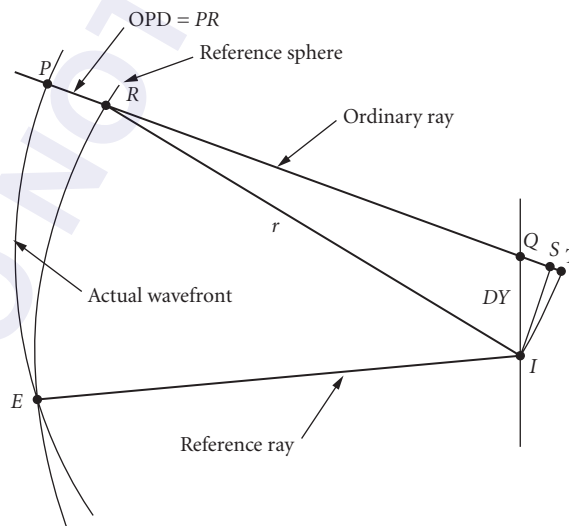


FIGURE 3 The relation between ray trajectories and optical path difference (OPD). See text.

Is it the ray through the center of the aperture stop, or perhaps the ray through the center of the actual vignetted aperture? These two definitions will result in different reference rays, and correspondingly different values for the optical path difference. In fact, in many practical applications neither definition is used, and the actual ideal image point is defined to be the one that minimizes the variance of the optical path difference (and hence maximizes the peak intensity of the diffraction image).

Moreover, the figure shows that even if the ideal image point is precisely defined, the value of the optical path difference depends on the point E where the actual wavefront intersects the reference sphere. For the particular point shown, the optical path difference is the optical length along the ordinary ray from the object point to the point T , less the optical length along the reference ray from the object point to the point I . As the radius of the reference sphere is increased, the point T merges with the point S , where a perpendicular from the ideal image point intersects the ordinary ray.

The above somewhat extended discussion is meant to demonstrate that even “well-known” optical terms are not always precisely defined. Not surprisingly, various optical design programs in common use produce different values for such quantities. There has been little effort to standardize the definitions of many terms, possibly because one cannot legislate physics. In any case, it is important for the user of an optical design program to understand precisely what the program is computing.

Virtually all optical design programs can trace single rays and display the ray heights and direction cosines on each surface. Other data, such as the path length, angles of incidence and refraction, and direction of the normal vector, are also commonly computed. Another type of ray-data display that is nearly universal is the ray-intercept curve, which shows ray displacement on the final image surface versus (fractional) pupil coordinates. A variation plots optical path difference versus pupil coordinates.

In addition to the uncertainty concerning the definition of ray displacement and optical path difference, there are different methods for handling the pupil coordinates. Some programs use entrance pupil coordinates, while others use exit pupil coordinates. In most cases, there is not a significant difference, but in the case of systems containing cylindrical lenses, for example, there are major differences.

Another consideration relating to ray-intercept curves is the way in which vignetting is handled. This is coupled to the way the program handles apertures. As mentioned before, apertures have a special status in many optical design programs. Rays can be blocked by apertures, but this must be handled as a special case by the program, because there is nothing inherent in the ray-trace equations that prevents a blocked ray from being traced, in contrast to a ray that misses a surface or undergoes total internal reflection.

Even though a surface may have a blocking aperture, it may be desirable to let the ray trace proceed anyway. As mentioned before, blocking rays in optimization can produce instabilities that prevent convergence to a solution even though all the rays in the final solution are contained within the allowed apertures. Another situation where blocking can be a problem concerns central obstructions. In such systems, the reference ray may be blocked by an obstruction, so its data are not available to compute the displacement or optical path difference of an ordinary ray (which is not blocked). The programmer must anticipate such situations and build in the proper code to handle them.

In the case of ray-intercept curves, it is not unusual for programs to display data for rays that are actually blocked by apertures. The user is expected to know which rays get through, and ignore the others, a somewhat unreasonable expectation. The justification for allowing it is that the designer can see what would happen to the rays if the apertures were increased.

In addition to ray-intercept curves, optical design programs usually display field sag plots showing the locations of the tangential and sagittal foci as a function of field angle and distortion curves. In the case of distortion, there is the question of what to choose as a reference height. It is generally easiest to refer distortion to the paraxial chief ray height in the final image surface, but in many cases it is more meaningful to refer it to the centroid height of a bundle of exact rays from the same object point. Again, it is important for the user to know what the program is computing.

Spot-Diagram Analysis

Spot diagrams provide the basis for realistic modeling of optical systems in an optical design program. In contrast to simple ray-trace evaluation, which shows data from one or a few rays, spot

diagrams average data from hundreds or thousands of rays to evaluate the image of a point source. Notwithstanding this, it should be understood that the principal purpose of an optical design program is to design a system, not to simulate its performance. It is generally up to the designer to understand whether or not the evaluation model of a system is adequate to characterize its real performance, and the prudent designer will view unexpected results with suspicion.

From a programmer's point of view, the most difficult task in spot-diagram analysis is to accurately locate the aperture of the system. For systems that have rotational symmetry, this is not difficult, but for off-axis systems with vignetted apertures it can be a challenging exercise. However, the results of image evaluation routines are often critically dependent on effects that occur near the edges of apertures, so particular care must be paid to this problem in writing optical design software. Like many other aspects of an optical design program, there is a trade-off between efficiency and accuracy.

A spot diagram is an assemblage of data describing the image-space coordinates of a large number of rays traced from a single object point. The data may be either monochromatic or polychromatic. Each ray is assigned a weight proportional to the fractional energy that it carries. Usually, the data saved for each ray include its xyz coordinates on the image surface, the direction cosines klm , and the optical path length or optical path difference from the reference ray. The ray coordinates are treated statistically to calculate root-mean-square spot sizes. The optical path lengths yield a measure of the wavefront quality, expressed through its variance and peak-to-valley error.

To obtain a spot diagram, the entrance pupil must be divided into cells, usually of equal area. Although for many purposes the arrangement of the cells does not matter, for some computations (e.g., transfer functions) it is advantageous to have the cells arranged on a rectangular grid. To make the computations have the proper symmetry, the grid should be symmetrical about the x and y axis. The size of the grid cells determines the total number of rays in the spot diagram.

In computing spot diagrams, the same considerations concerning the reference point appear as for ray fans. That is, it is possible to define ray displacements with respect to the chief-ray, the paraxial ray height, or the centroid of the spot diagram. However, for spot diagrams it is most common to use the centroid as the reference point, both because many image evaluation computations require this definition, and also because the value for the centroid is readily available from the computed ray data.

$$a_i = \epsilon_{xi} = \text{ray displacement in the } x \text{ direction}$$

$$b_i = \epsilon_{yi} = \text{ray displacement in the } y \text{ direction}$$

$$c_i = k_i/m_i = \text{ray slope in the } x \text{ direction}$$

$$d_i = l_i/m_i = \text{ray slope in the } y \text{ direction}$$

$$w_i = \text{weight assigned to ray}$$

The displacements of rays on a plane shifted in the z direction from the nominal image plane by an amount Δz are given by

$$\delta x_i = a_i + b_i \Delta z$$

$$\delta y_i = c_i + d_i \Delta z$$

If there are n rays, the coordinates of the centroid of the spot diagram are

$$\delta \bar{x} = \frac{1}{W} \sum_{i=1}^n w_i \delta x_i = A + B \Delta z$$

$$\delta \bar{y} = \frac{1}{W} \sum_{i=1}^n w_i \delta y_i = C + D \Delta z$$

where W is a normalizing constant that ensures that the total energy in the image adds up to 100 percent, and

$$A = \frac{1}{W} \sum_{i=1}^n w_i a_i$$

$$B = \frac{1}{W} \sum_{i=1}^n w_i b_i$$

$$C = \frac{1}{W} \sum_{i=1}^n w_i c_i$$

$$D = \frac{1}{W} \sum_{i=1}^n w_i d_i$$

The mean-square spot size can then be written as

$$\text{MSS} = \frac{1}{W} \sum_{i=1}^n w_i \{(\delta x_i - \delta \bar{x})^2 + (\delta y_i - \delta \bar{y})^2\}$$

Usually, the root-mean-square (rms) spot size, which is the square root of this quantity, is reported. Since the MSS has a quadratic form, it can be written explicitly as a function of the focus shift by

$$\text{MSS} = P + 2Q \Delta z + R(\Delta z)^2$$

where

$$P = \frac{1}{W} \sum_{i=1}^n w_i \{(a_i^2 + c_i^2) - (A^2 + C^2)\}$$

$$Q = \frac{1}{W} \sum_{i=1}^n w_i \{(a_i b_i + c_i d_i) - (AB + CD)\}$$

$$R = \frac{1}{W} \sum_{i=1}^n w_i \{(b_i^2 + d_i^2) - (B^2 + D^2)\}$$

Differentiating this expression for the MSS with respect to focus shift, then setting the derivative to zero, determines the focus shift at which the rms spot size has its minimum value:

$$\Delta z_{\text{opt}} = -Q/R$$

Although the above equations determine the rms spot size in two dimensions, similar one-dimensional equations can be written for x and y separately, allowing the ready computation of the tangential and sagittal foci from spot-diagram data. In addition, it is straightforward to carry out the preceding type of analysis using optical path data, which leads to the determination of the center of the reference sphere that minimizes the variance of the wavefront.

Beyond the computation of the statistical rms spot size and the wavefront variance, most optical design programs include a variety of image evaluation routines that are based on spot diagram data. It is useful to characterize them as belonging to geometrical optics or physical optics, according to whether they are based on ray displacements or wavefronts, although, of course, all are based on the results of geometrical ray tracing.

Geometrical Optics Most optical design programs provide routines for computing radial diagrams and knife-edge scans. To compute a radial energy diagram, the spot-diagram data are sorted according to increasing ray displacement from the centroid of the spot. The fractional energy is then plotted as a function of spot radius. The knife-edge scan involves a similar computation, except that the spot-diagram data are sorted according to x or y coordinates, instead of total ray displacement.

Another type of geometrical image evaluation based on spot-diagram data is the so-called *geometrical optical transfer function* (GOTF). This function can be developed as the limiting case, as the wavelength approaches zero, of the actual diffraction MTF, or, alternately, in a more heuristic way as the Fourier transform of a line spread function found directly from spot-diagram ray displacements (see, for example, Smith's book¹⁰). From a programming standpoint, computation of the GOTF involves multiplying the ray displacements by 2π times the spatial frequency under consideration, forming cosine and sine terms, and summing over all the rays in the spot diagram. The computation is quick, flexible, and if there are more than a few waves of aberration, accurate. The results of the GOTF computation are typically shown as either plots of the magnitude of the GOTF as a function of frequency, or alternately in the form of what is called a "through-focus" MTF, in which the GOTF at a chosen frequency is plotted as a function of focus shift from the nominal image surface.

Physical Optics The principal physical optics calculations based on spot-diagram data are the modulation transfer function, sometimes called the "diffraction" MTF, and the point spread function (PSF). Both are based on the wavefront derived from the optical path length data in the spot diagram. There are various ways to compute the MTF and PSF, and not all programs use the same method. The PSF, for example, can be computed from the pupil function using the fast Fourier transform algorithm or, alternately, using direct evaluation of the Fraunhofer diffraction integral. The MTF can be computed either as the Fourier transfer of the PSF or, alternately, using the convolution of the pupil function.¹¹ The decision as to which method to use involves speed, accuracy, flexibility, and ease of coding.

In physical-optics-based image evaluation, accuracy can be a problem of substantial magnitude. In many optical design programs, diffraction-based computations are only accurate for systems in which diffraction plays an important role in limiting the performance. Systems that are limited primarily by geometrical aberrations are difficult to evaluate using physical optics, because the wavefront changes so much across the pupil that it may be difficult to sample it sufficiently using a reasonable number of rays. If the actual wavefront in the exit pupil is compared to a reference sphere, the resultant fringe spacing defines the size required for the spot diagram grid, since there must be several sample points per fringe to obtain accurate diffraction calculations. To obtain a small grid spacing, one can either trace many rays, or trace fewer rays but interpolate the resulting data to obtain intermediate data.

Diffraction calculations are necessarily restricted to one wavelength. To obtain polychromatic diffraction results it is necessary to repeat the calculations in each color, adding the results while keeping track of the phase shifts caused by the chromatic aberration.

3.5 OPTIMIZATION

The function of the optimization part of the program is to take a *starting design* and modify its construction so that it meets a given set of specifications. The starting design may be the result of a previous design task, a lens from the library, or a new design based on general optical principles and the designer's intuition.

The performance of the design must be measured by a single number, often known in optics as the *merit function*, although the term *error function* is more descriptive and will be used here. The error function is the sum of squares of quantities called *operands* that characterize the desired attributes. Examples of typical operands include paraxial constants, aberration coefficients, and exact ray displacements. Sometimes, the operands are broken into two groups: those that must be satisfied exactly, which may be called *constraints*, and others that must be minimized. Examples of constraints might include paraxial conditions such as the focal length or numerical aperture.

The constructional parameters to be adjusted are called *variables*, which include lens curvatures, thicknesses, refractive indices, etc. Often the allowed values of the variables are restricted, either by requirements of physical reality (e.g., positive thickness) or the given specifications (e.g., lens diameters less than a prescribed value). These restrictions are called *boundary conditions*, and represent another form of constraint.

Usually, both the operands and constraints are nonlinear functions of the variables, so optical design involves nonlinear optimization with nonlinear constraints, the most difficult type of problem from a mathematical point of view. A great deal of work has been carried out to develop efficient, general methods to solve such problems. Detailed consideration of these methods is beyond the scope of this chapter, and the reader is referred to a paper by Hayford.¹²

In a typical optical design task, there are more operands than variables. This means that there is, in general, no solution that makes all of the operands equal to their target values. However, there is a well-defined solution called the *least-squares* solution, which is the state of the system for which the operands are collectively as close to their targets as is possible. This is the solution for which the error function is a minimum.

The Damped Least-Squares Method

Most optical design programs utilize some form of the *damped least-squares* (DLS) method, sometimes in combination with other techniques. DLS was introduced to optics in about 1960, so it has a history of 50 years of (usually) successful application. It is an example of what is known as a *downhill* optimizer, meaning that in a system with multiple minima, it is supposed to find the nearest local minimum. In practice, it sometimes suffers from *stagnation*, yielding slow convergence. On the other hand, many designers over the years have learned to manipulate the damping factor to overcome this deficiency, and even in some cases to find solutions beyond the local minimum.

We consider first the case of unconstrained optimization. Let the system have M operands f_i and N variables x_j . The error function ϕ is given by

$$\phi = f_1^2 + f_2^2 + \cdots + f_M^2$$

Define the following:

$$\mathbf{A} = \text{derivative matrix, } A_{ij} \equiv \frac{\partial f_i}{\partial x_j}$$

$$\mathbf{G} = \text{gradient vector, } G_k \equiv \frac{1}{2} \frac{\partial \phi}{\partial x_k}$$

\mathbf{x} = change vector

\mathbf{f} = error vector

With these definitions, we have

$$\mathbf{G} = \mathbf{A}^T \mathbf{f}$$

If we assume that the changes in the operands are linearly proportional to the changes in the variables, we have

$$\mathbf{f} = \mathbf{A}\mathbf{x} + \mathbf{f}_0$$

$$\mathbf{G} = \mathbf{A}^T \mathbf{A}\mathbf{x} + \mathbf{G}_0$$

At the solution point, the gradient vector is zero, since the error function is at a minimum. The change vector is thus

$$\mathbf{x} = -(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{G}_0$$

These are called the least-squares *normal equations*, and are the basis for linear least-squares analysis. When nonlinear effects are involved, repeated use of these equations to iterate to a minimum often leads to a diverging solution. To prevent such divergence, it is common to add another term to the error function, and this limits the magnitude of the change vector \mathbf{x} . In the DLS method, this is accomplished by defining a new error function

$$\varphi = \phi + p\mathbf{x}^T \mathbf{x}$$

A key property of DLS is that the minimum of φ is the same as the minimum of ϕ since, at the minimum, the change vector \mathbf{x} is zero. By differentiating and setting the derivative equal to zero at the minimum, we arrive at the damped least-squares equations

$$\mathbf{x} = -(\mathbf{A}^T \mathbf{A} + p\mathbf{I})^{-1} \mathbf{G}_0$$

which look like the normal equations with terms added along the diagonal. These terms provide the damping, and the factor p is called the *damping factor*. This particular choice of damping is called *additive damping* but, more generally, it is possible to add any terms to the diagonal and still maintain the same minimum. Some optical design programs multiply the diagonal elements of the $\mathbf{A}^T \mathbf{A}$ matrix by a damping factor, while others make them proportional to the second derivative terms. Although theoretical arguments are sometimes advanced to support the choice of a particular method of damping, in practice the choice of damping factor is an ad hoc way to accelerate convergence to a solution by limiting the magnitude (and changing the direction) of the change vector found from the normal equations.

In practical optical design work, it has been found that no single method for choosing the damping factor works best in all cases. In a particular problem, one method may be dramatically better than another, but in a different problem, the situation may be completely reversed. Every optical design program has its unique way of choosing the optimum damping, which makes each program different from the others, and gives it a *raison d'être*.

Although the principal use of the damping factor is to accelerate convergence by limiting the magnitude of the change vector, the damping factor has also been used routinely to increase the magnitude of the change vector to escape a local minimum. During the course of a minimization task, if the solution stagnates, or does not converge to what the designer believes to be an acceptable configuration, it may be possible to force the solution into another region by running one or more iterations with reduced damping in which the error function increases.

Constraints and Boundary Conditions There are two general methods used in optical design programs for handling constraints and boundary conditions. The first is to add a term (called a *penalty function*) to the error function that targets the constraint to its desired value. In the case of boundary violations, “one-sided” terms can be added, or special weighting functions can be constructed that increase in magnitude as a violation goes farther into a forbidden region. The other method augments

the number of equations by the number of constraints and solves the resulting equations using the Lagrange multiplier method. This produces a minimum that satisfies the constraints exactly and minimizes the remaining error function.

The penalty function method is more flexible and faster (since there are fewer equations) than the Lagrange multiplier method. On the other hand, the Lagrange multiplier method gives more precise control over the constraints. Both are commonly used in optical design software.

Other Methods

Although DLS is used in the vast majority of optical design applications, other methods are occasionally used,¹² and two warrant mention. These are *orthonormalization*, which has been used to overcome stagnation in some DLS problems, and *simulated annealing*, which has been used for global optimization.

Orthonormalization The technique of orthonormalization for the solution of optical design problems was introduced by Grey.² Although it solves the same problem as DLS does, it proceeds in a very different fashion. Instead of forming the least-squares normal equations, Grey works directly with the operand equations

$$\mathbf{Ax} = -\mathbf{f}$$

To understand Grey's method, it is best to forget about optics and consider the solution of these equations strictly from a mathematical point of view. The point of view that Grey uses is that \mathbf{f} represents a vector in m -dimensional space. The columns of \mathbf{A} can be regarded as basis vectors in this m -dimensional space. Since there are only n columns, the basis vectors do not span the space. The change vector \mathbf{x} represents a projection of \mathbf{f} on the basis vectors defined by \mathbf{A} . At the solution point, the residual part of \mathbf{f} will be orthogonal to its projection on the basis vectors.

In Grey's orthonormalization method, the solution of the equations is found by a technique similar to Gram-Schmidt orthogonalization, but during the solution process, the actual error function is evaluated several times in an effort to use the best variables to maximum advantage. Because of this, the method is computationally intensive compared to DLS. However, the extra computation is justified by a more accurate solution. The common wisdom is that orthogonalization is superior to DLS near a solution point, and inferior to DLS when the solution is far removed from the starting point.

Simulated Annealing Simulated annealing has been applied to optical design optimization, chiefly in problems where the task is to find a global minimum. The method varies drastically from other techniques. It makes no use of derivative information, and takes random steps to form trial solutions. If a trial solution has a lower error function than the current system, the new system replaces the old. If a trial solution has a higher error function than the current system, it may be accepted, depending on how much worse it is. The probability of acceptance is taken to be $\exp(-\Delta\phi/T)$, where T is an experimentally determined quantity. In general simulated annealing, T is provided by the user. In adaptive simulated annealing, T is reduced automatically according to algorithms that hold the system near statistical equilibrium.

Error Functions

Obviously, the choice of an error function has a major impact on the success of an optical design task. There are a number of requirements that an error function should meet. Most importantly, the error function should accurately characterize the desired properties of the system under design. There is little chance of success if the program is optimizing the wrong thing. Yet this is an area of great difficulty in computer-aided optical design, because it is at odds with efficiency. In order to obtain more accuracy, more extensive computations should be carried out, but this takes time.

There are two schools of thought concerning the implementation of error functions in optical design programs. The first holds that the designer should have complete control over the items included in the error function, while the second holds that the program itself should set up the basic error function, allowing the designer some degree of control through weighting functions. Neither school has demonstrated superiority, but the approach to error function construction taken by various optical design programs accounts for user allegiances that are sometimes remarkably strong.

The different ways that optical design programs handle error functions makes it difficult to discuss the topic here in anything other than broad detail. At one extreme are programs that provide practically no capability for the user to insert operands, displaying only the value of the overall error function, while at the other extreme are programs that make the user enter every operand individually. Regardless of the user interface, however, there are some general concepts that are universally relevant.

Error functions can be based on either aberration coefficients or exact-ray data (or both). In the early stages of design, aberration coefficients are sometimes favored because they provide insight into the nature of the design, and do not suffer ray failures. However, the accuracy of aberration coefficients for evaluating complex systems is not very good, and exact-ray data are used in virtually all final optimization work.

So far as exact-ray error functions are concerned, there is the question of whether to use ray displacements or optical path difference (or both). This is a matter of user (or programmer) preference. The use of ray displacements leads to minimizing geometrical spot sizes, while the use of optical path difference leads to minimizing the wavefront variance.

For exact-ray error functions, a suitable pattern of rays must be set up. This is often called a *ray set*. There are three common methods for setting up a ray set. The first is to allow the designer to specify the coordinates (object, pupil, wavelength, etc.) for a desired set of rays. This gives great flexibility, but demands considerable skill from the user to ensure that the resulting error function accurately characterizes performance.

The other two methods for setting up ray sets are more automatic. The first is to allow the user to specify object points, and have the program define a rectangular grid of rays in the aperture for each point. The second uses a Gaussian integration scheme proposed by Forbes to compute the rms spot size, averaged over field, aperture, and wavelength.¹³ The Forbes method, which is restricted to systems having plane symmetry, leads to dividing the aperture into *rings* and *spokes*. For systems having circular pupils, the Forbes method has both superior accuracy and efficiency, but for vignetted pupils, there is little difference between the two.

Multiconfiguration Optimization

Multiconfiguration optimization refers to a process in which several systems having some common elements are optimized jointly, so that none of the individual systems but the *ensemble* of all of the systems is optimized. The archetype of multiconfiguration systems is the *zoom* system in which the focal length is changed by changing the separation between certain elements. The system is optimized simultaneously at high, medium, and low magnifications to produce the best overall performance.

Most of the larger optical design programs have the capability to carry out multiconfiguration optimization, and this capability is probably used more for non-zoom systems than for zoom systems. A common use of this feature is to optimize a focal system for through-focus performance in order to minimize sensitivity to image plane shifts. In fact, multiconfiguration optimization is used routinely to control tolerances.

Tolerancing

Beyond the task of desensitizing a given design, considerations of manufacturing tolerances become increasingly important as the complexity of optical designs increases. It is quite easy to

design optical systems that cannot be built because the fabrication tolerances are beyond the capability of optical manufacturing technology. In any case, specifying tolerances is an integral part of optical design, and a design project cannot be considered finished until appropriate tolerances are established.

Tolerancing is closely related to optimization. The basic tolerance computation is to calculate how much the error function changes for a small change in a construction parameter, which is the same type of computation carried out when computing a derivative matrix. Even more relevant, however, is the use of *compensators*, which requires reoptimization. A compensator is a construction parameter that can be adjusted to compensate for an error introduced by another construction parameter. For example, a typical compensator would be the image distance, which could be adjusted to compensate for power changes introduced by curvature errors.

There is considerable variation in how different optical design programs handle tolerancing. Some use the reoptimization method described here, while others use Monte Carlo techniques. Some stress interaction with the designer, while others use defaults for more automatic operation.

3.6 OTHER TOPICS

Of course, many other topics would be included in a full discussion of optical design software. Space limitations and our intended purpose prevents any detailed consideration, but a few of the areas where there is still considerable interest are the following.

Simulation

There is increased interest in using optical design programs to simulate the performance of actual systems. The goal is often to be able to calculate radiometric throughput of a system used in conjunction with a real extended source. It is difficult to provide software to do this with much generality, because brute force methods are very inefficient and hard to specify, while elegant methods tend to have restricted scope, and demand good judgment by the person modeling the physical situation. Nevertheless, with the increase in the speed of computers, there is bound to be an increasing use of optical design software for evaluating real systems.

Global Optimization

After several years during which there was little interest in optimization methodology, the tremendous increase in the speed of new computers has spawned a renewal of efforts to find global, rather than local, solutions to optical design problems. Global optimization is a much more difficult problem than local optimization. In the absence of an analytic solution, one never knows whether a global optimum has been achieved. All solution criteria must specify a region of interest and a time limit, and the method cannot depend on the starting point. The simulated annealing method described above is one area of continuing interest. Several methods for what might be called *pseudo-global* optimization have been used in commercial optical design programs, combining DLS with algorithms that allow the solution to move away from the current local minimum.

Computing Environment

Increasingly, optical design programs are used in conjunction with other software. Drawing programs, manufacturing inventory software, and intelligent databases are all relevant to optical design. While the conventional optical design program has been a *stand-alone* application, there is increasing demand for integrating optical design into more general design tasks.

3.7 BUYING OPTICAL DESIGN SOFTWARE

The complexity of the optical design process, together with the breadth of applications of optics, has created an ongoing market for commercial optical design software. For people new to optical design, however, the abundance of advertisements, feature lists, and even technical data sheets doesn't make purchasing decisions easy. The following commentary, adapted from an article, may be helpful in selecting an optical design program.¹⁴ It considers five key factors: hardware, features, user interface, cost, and support.

Hardware

It used to be that the choice of an optical design program was governed by the computer hardware available to the designer. Of course, when the hardware cost was many times higher than the software cost, this made a great deal of sense. Today, however, the software often costs more than the hardware, and many programs can be run on several different computer platforms, so the choice of computer hardware is less important. The hardware currently used for optical design is principally IBM-PC compatible.

To run optical design software, the fastest computer that can be obtained easily is recommended. The iterative nature of optical design makes the process interminable. There is a rule, sometimes called the Hyde maxim, that states that an optical design is finished when the time or money runs out. Notwithstanding this, the speed of computers has ceased to be a significant impediment to ordinary optical design. Even low-cost computers now trace more than 1000 ray-surfaces/s, a speed considered the minimum for ordinary design work, and create the potential for solving new types of problems formerly beyond the range of optical design software.

Before desktop computers, optical design software was usually run on time-shared central computers accessed by terminals, and some programs are still in that mode. There seems to be general agreement, however, that the memory-mapped display found on PCs provide a superior working environment and dedicated desktop computer systems are currently most popular.

Features

If you need a particular feature to carry out your optical design task, then it is obviously important that your optical design program have that feature. But using the number of features as a way to select an optical design program is probably a mistake. There are more important factors, such as cost, ease of use, and scope. Moreover, you might assume that all the features listed for a program work simultaneously, which may not be true. For example, if a vendor states that its program handles holograms and toric surfaces, you might assume that you can work with holographic toroids, but this may not be true.

The continuing growth of optics and the power of desktop computers has put heavy demands on software vendors to keep up with the development of new technology. Moreover, since the customer base is small and most vendors now support the same computer hardware, the market has become highly competitive. These factors have led to a "feature" contest in which software suppliers vie to outdo each other. While this is generally good for the consumer, the introduction of a highly visible new feature can overshadow an equally important but less obvious improvement (for example, fewer bugs or better documentation). In addition, the presence of a number of extra features is no guarantee that the underlying program is structurally sound.

User Interface

There is very little in common between the user interfaces used by various optical design programs. Each seems to have its own personality. The older programs, originally designed to run in batch

mode on a large computer, are usually less interactive than ones that were written specifically for desktop computers. Batch programs tend to be built more around default actions than interactive programs, which require more user input. It would be hard to put any of today's major optical design programs in a box classified as either batch or interactive, but the look and feel of a program has a strong influence on its usefulness.

Many people don't realize that the most important benefit of using an optical design program is often the understanding that it provides the user about how a particular design works. It's often tempting to think that if the computer could just come up with a satisfactory solution, the design would be finished. In practice, it is important to know the trade-offs that are made during a design project. This is where the judgment of the optical engineer comes in, knowing whether to make changes in mechanical or electrical specifications to achieve the optimum balance in the overall system. Lens designers often say that the easiest lens to design is one that has to be diffraction-limited, because it is clear when to stop. If the question of how to fit the optics together with other system components is important, then the ability of the user to work interactively with the design program can be a big help.

Cost

In today's market, there is a wide range of prices for optical design software. This can be very confusing for the first-time buyer, who often can't see much difference in the specifications. The pricing of optical design software is influenced by (at least) three factors.

First, the range of tasks that can be carried out using an optical design program is enormous. The difference in complexity between the job of designing a singlet lens for a simple camera, and that of designing a contemporary objective for a microlithographic masking camera is somewhat akin to the difference between a firecracker and a hydrogen bomb.

Second, all software is governed by the factors originally studied in F. P. Brooks' famous essay *The Mythical Man-Month*.¹⁵ Brooks was director of the group that developed the operating system for the IBM 360, a mainframe computer introduced in the 1960s. Despite its provocative title, Brooks' essay is a serious work that has become a standard reference for software developers. In it, he notes that if the task of developing a program to be used on a single computer by its author has a difficulty of 1, then the overall difficulty of producing integrated software written by a group of people and usable by anyone on a wide range of computers may be as high as 10. In recent years, the scope of the major optical design programs has grown too big for a single programmer to develop and maintain, which raises costs.

Third, there are structural differences in the way optical design software is sold. The original mainframe programs were rented, not sold. If the user did not want to continue monthly payments, the software had to be returned. PC programs, on the other hand, are usually sold with a one-time fee. In the optical design software business, several vendors offer a compromise policy, combining a permanent license with an optional ongoing support fee.

It would be nice if the buyer could feel comfortable that "you get what you pay for," but unfortunately this view is too simplistic. One program may lack essential capabilities, another may contain several unnecessary features when evaluated for a particular installation. Buying on the basis of cost, like features, is probably not a good idea.

Support

Support is an important aspect to consider in selecting an optical design program, and it is often difficult to know what is included in support. Minimal support consists of fixing outright bugs in the program. More commonly, support includes software updates and phone or email assistance in working around problems.

Optical design programs are typically not bug-free. Unlike simple programs like word processors, optimization programs cannot be fully tested, because they generate their own data. One result of this is that software vendors are generally reluctant to offer any warranty beyond a "best-effort"

attempt to fix reported problems. Unfortunately, there is no good way for buyers to know whether and when their particular problems may be fixed; the best approach is probably to assess the track record of the vendor by talking to other users.

Coupled with support is user training. Although it should be possible to use a program by studying the documentation, the major optical design software vendors offer regular seminars, often covering not only the mechanics of using their program, but also general instruction in optical design. For new users, this can be a valuable experience.

3.8 SUMMARY

As stated in the introduction, this chapter is intended as a survey for readers who are not regular users of optical design software. The form of an optical design program described here, consisting of lens entry, evaluation, and optimization sections, is used in many different programs. There has been little standardization in this field, so the “look and feel,” performance features and extent of various programs are quite different. Nonetheless, it is hoped that with a knowledge of the basic features described here, the reader will be in a good position to judge whether an optical design program is of use, and to make an informed decision about whether one particular program is better than another.

3.9 REFERENCES

1. D. P. Feder, “Automatic Optical Design,” *Appl. Opt.* **2**:1209–1226 (1963).
2. D. S. Grey, “Aberration Theories for Semiautomatic Lens Design by Electronic Computers,” *J. Opt. Soc. Am.* **53**:672–680 (1963).
3. G. H. Spencer, “A Flexible Automatic Lens Correction Procedure,” *Appl. Opt.* **2**:1257–1264 (1963).
4. C. G. Wynne and P. Wormell, “Lens Design by Computer,” *Appl. Opt.* **2**:1233–1238 (1963).
5. See, for example, *The Photonic Industry Buyer’s Guide*, Laurin Publishing, Pittsfield, MA 01202.
6. *U.S. Military Handbook for Optical Design*, republished by Sinclair Optics, Fairport, NY 14450 (1987).
7. G. H. Spencer and M. V. R. K. Murty, “Generalized Ray-Tracing Procedure,” *J. Opt. Soc. Am.* **52**:672–678 (1962).
8. D. C. O’Shea, *Elements of Modern Optical Design*, John Wiley & Sons, New York (1985).
9. H. A. Buchdahl, *Optical Aberration Coefficients*, Oxford Press, London (1954).
10. W. J. Smith, *Modern Optical Engineering*, McGraw-Hill, New York (1990).
11. H. H. Hopkins, “Numerical Evaluation of the Frequency Response of Optical Systems,” *Proc. Phys. Soc. B* **70**:1002–1005 (1957).
12. M. J. Hayford, “Optimization Methodology,” *Proc. SPIE* **531**: 68–81 (1985).
13. G. W. Forbes, “Optical System Assessment for Design: Numerical Ray Tracing in the Gaussian Pupil,” *J. Opt. Soc. Am. A* **5**:1943–1956 (1988).
14. D. C. Sinclair, “Optical Design Software: What to Look For in a Program,” *Photonics Spectra*, Nov. 1991.
15. F. P. Brooks, Jr., *The Mythical Man-Month*, Addison-Wesley, Reading, MA (1975).

OPTICAL SPECIFICATIONS

Robert R. Shannon*

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

4.1 GLOSSARY

ATF	approximate transfer factor
DTF	diffraction transfer function
MTF	modulation transfer function
W	wavefront error in units of wavelengths
W_{rms}	root-mean-square wavefront error
t -number	f -number adjusted for lens transmission
ν	normalized spatial frequency

4.2 INTRODUCTION

Setting the specifications for an optical instrument or system is an essential part of engineering, designing, or purchasing an optical system. Since the optics usually serve as a portion of a larger system, the specifications are frequently set by project managers who do not have specific knowledge in the basics of optical systems. This can at times lead to unrealistic requirements being established; this can profoundly affect the probability of success for the system. Properly drafted specifications can make the entire project successful and cost effective. Poorly written specifications can lead to excess cost and ultimately project failure.

One of the difficulties with setting optical specifications is that the ultimate result of a beam of light passing through a complex assembly of components is affected by each of those components, which in turn need to be specified and tolerances placed upon the fabrication and assembly of those components. In the case of an imaging system, the problem is compounded by the need to describe an optical system which passes many bundles of light across a wide field of view. Even in the case of single beam, optical communications components, indirect issues such as scattered light and environmental stability may prove to be major issues.

*Retired.

In the worst case, the specifications may be set so high that the system is not capable of being manufactured. In most cases, the specifications interact with other devices, such as detector arrays, and matching the quality of the optic to the limits of the sensor is required. In this section some of the principles involved in setting the specifications will be discussed, and guidelines provided for carrying out the process of specification setting. The reader will have to extend these principles to the device or system that is being considered. In this chapter, the stress will be placed on imaging systems.

Specifications for optical systems cover a wide range of needs. Functional specifications of the image quality or other optical characteristics are required for the satisfactory operation of a system. These functional specifications serve as the goal for the design and construction of the optical system. In addition, these specifications are a basis for tolerances placed upon the components of the optical system and lead to detailed component specifications used for procurement of the optical elements of the system. Assembly specifications and detailed specifications of optical parts to be produced by a shop can be written based upon these component specifications. The detail and extent of information required is different at each step. Over- or underspecification can contribute significantly to the cost or feasibility of design of an optical system.

Functional specifications are also used to describe the characteristics that an instrument must demonstrate in order to meet the needs of the user. This may include top-level requirements such as size, weight, image scale, image format, power levels, spectral range, and so on. Component specifications are developed after design of the system and describe the optical components, surface, and materials used in the system to the detail necessary to permit fabrication of the components. Assembly specifications are another derivative of the design and system specifications. These include the statement of tolerances upon location of the components, as well as the procedure to be used in assembling and testing the system.

The development and writing of these specifications is important both for initiating and for tracking the course of development of an optical instrument. In a business or legal sense, specifications are used to establish responsibility for a contractor or subcontractor, as well as to define the basis for bidding on the job. Thus the technical specifications can have business importance as well as engineering significance. "Meeting the customer's specifications" is an essential part of any design and fabrication task. Identifying areas where the specifications could be altered with benefit to all parties is an important business and engineering responsibility.

Specifications are usually communicated as a written document following some logical format. Although there are some international standards that may cover the details of drawings of components, there is no established uniform set of standards for stating the specifications on a system or component. The detailed or component specifications are usually added as explanatory notes to drawings of the components to be fabricated. In modern production facilities, the specifications and tolerances are often part of a digital database that is accessed as part of the production of the components of the system.

The detail and the intent of each of these classes of specifications are different. Optical specifications differ from many mechanical or other sets of specifications in that numbers are applied to surfaces and dimensions that control the cumulative effect of errors imposed on a wavefront passing through the total system. Each of the specifications must be verifiable during fabrication, and the overall result must be testable after completion.

Mechanical versus Optical Specifications

There are two types of specifications that are applied to an optical system or assembly. One set of these includes mechanical tolerances on the shape or location of the components that indirectly affect the optical quality of the image produced by the system. Examples of this include the overall size or weight of the system. The other set consists of specialized descriptions that directly affect the image quality. Examples of this latter type of specification are modulation transfer function (MTF), illumination level, and location of the focal plane relative to the system.

System versus Components Specifications

Some specifications have meaning only with respect to the behavior of the entire optical system. Others apply to the individual components, but may affect the ability of the entire system to function.

An example of a system specification is a set of numbers limiting the range of acceptable values of the MTF that are required for the system. Another system specification is the desired total light transmission of the system.

Examples of component specifications are tolerances upon surface irregularity, sphericity, and scattering. The related component specification based upon the system light transmission specification might provide detailed statements about the nature and properties of the antireflective coatings to be applied to the surface of each element.

Image Specifications

The specifications that are applied to the image usually deal with image quality. Examples are modulation transfer function, fraction of scattered light, resolution, or distortion. In some cases, these specifications can be quite general, referring to the ability of the lens to deliver an image suitable for a given purpose, such as the identification of serial numbers on specific products that are to be read by an automated scanner. In other cases, the requirements will be given in a physically meaningful manner, such as “the MTF will be greater than 40 percent at 50 lines per millimeter throughout the field of view.”

Other criteria may be used for the image specifications. One example is the energy concentration. This approach specifies the concentration of light from a point object on the image surface. For example, the specification might read “75 percent of the light shall fall within a 25- μm -diameter circle on the image.” This quantity is obviously measurable by a photometer with appropriate-size apertures. The function may be computed from the design data by a method of numerical integration similar to that providing the point spread function or modulation transfer function.

Wavefront Specifications

Wavefront specifications describe the extent to which the wavefront leaving the lens or components conforms to the ideal or desired shape. Usually the true requirement for an optical system is the specification of image quality, such as MTF, but there is a relation between the image quality and the wavefront error introduced by the optical system. The wavefront error may be left to be derived from the functional image quality specification, or it may be defined by the intended user of the system.

For example, a wavefront leaving a lens would ideally conform to a sphere centered on the chosen focal location. The departure of the actual wavefront from this ideal, would be expressed either as a matrix or map of departure of the wavefront from the ideal sphere, as a set of functional forms representing the deviation, or as an average [usually root-mean-square (rms)] departure from the ideal surface. By convention, these departures are expressed in units of wavelength, although there is a growing tendency to use micrometers as the unit of measure.

The rms wavefront error is a specific average over the wavefront phase errors in the exit pupil. The basic definition is found by defining the n th power average of the wavefront $W(x, y)$ over the area A of the pupil and then specifically defining W_{rms} or, in words, the rms wavefront error is the square root of the mean square error minus the square of the mean wavefront error:

$$\bar{W}^n = \frac{1}{A} \int W(x, y)^n dx dy$$

$$W_{\text{rms}} = \sqrt{[\bar{W}^2] - [\bar{W}]^2}$$

The ability to conveniently obtain a complete specification of image quality by a single number describing the wavefront shape has proven to be questionable in many cases. Addition of a correlation

length, sometimes expressed as a phase difference between separated points, has become common. In other cases, the relative magnitude of the error when represented by various orders of Zernike polynomials is used.

There is, of course, a specific relationship between the wavefront error produced by a lens and the resulting image quality. In the lens, this is established by the process of diffraction image formation. In establishing specifications, the image quality can be determined by computation of the modulation transfer function from the known wavefront aberrations. This computation is quite detailed and, while rapidly done using present day computer techniques, is quite complex for general specification setting. An approximation which provides an average MTF or guide to acceptable values relating wavefront error and MTF is of great aid.

A perfect lens is one that produces a wavefront with no aberration, or zero rms wavefront error. By convention, any wavefront with less than 0.07 wave, rms, of aberration is considered to be essentially perfect. It is referred to as *diffraction-limited*, since the image produced by such a lens is deemed to be essentially indistinguishable from a perfect image.

The definition of image quality depends upon the intended application for the lens. In general, nearly perfect image quality is produced by lenses with wavefront errors of less than 0.15 wave, rms. Somewhat poorer image quality is found with lenses that have greater than about 0.15 wave of error. The vast majority of imaging systems operate with wavefront errors in the range of 0.1 to 0.25 wave, rms.

There are several different methods that can be used to establish this relationship. The most useful comparison is with the MTF for a lens with varying amounts of aberration. The larger the wavefront error, the lower will be the contrast at specific spatial frequencies. For rms error levels of less than 0.25 or so, the relation is generally monotonic. For larger aberrations, the MTF becomes rather complex, and the relation between rms wavefront error and MTF value can be multiple valued. Nevertheless, an approximate relation between MTF and rms wavefront error would be useful in setting reasonable specifications for a lens.

There are several possible approximate relations, but one useful one is the empirical formula relating root-mean-square wavefront error and MTF given by

$$\text{MTF}(\nu) = \text{DTF}(\nu) \times \text{ATF}(\nu)$$

The functional forms for these values are

$$\text{DTF}(\nu) = \frac{2}{\pi} [\arccos(\nu) - \nu\sqrt{1-\nu^2}]$$

$$\text{ATF}(\nu) = \left[1 - \left(\frac{W_{\text{rms}}}{0.18} \right)^2 \right] (1 - 4(\nu - 0.5)^2)$$

$$\nu = \frac{\text{spatial frequency}}{\text{spatial frequency cutoff}} = \frac{N}{\left(\frac{1}{\lambda f\text{-number}} \right)}$$

These look quite complicated, but are relatively simple, as is shown in Fig. 1. This is an approximation, however, and it becomes progressively less accurate as the amount of the rms wavefront error W_{rms} exceeds about 0.18 wavelength. The approximation remains reasonably valid for lower spatial frequencies, less than about 25 percent of the diffraction limited cutoff frequency. The majority of imaging systems fall into this category.

Figure 1 shows a plot of several values for the MTF of an optical system using this approximate method of computation. The system designer can use this information to determine the appropriate level of residual rms wavefront error that will be acceptable for the system of interest. It is important to note that this is an empirical attempt to provide a link between the wavefront error and the

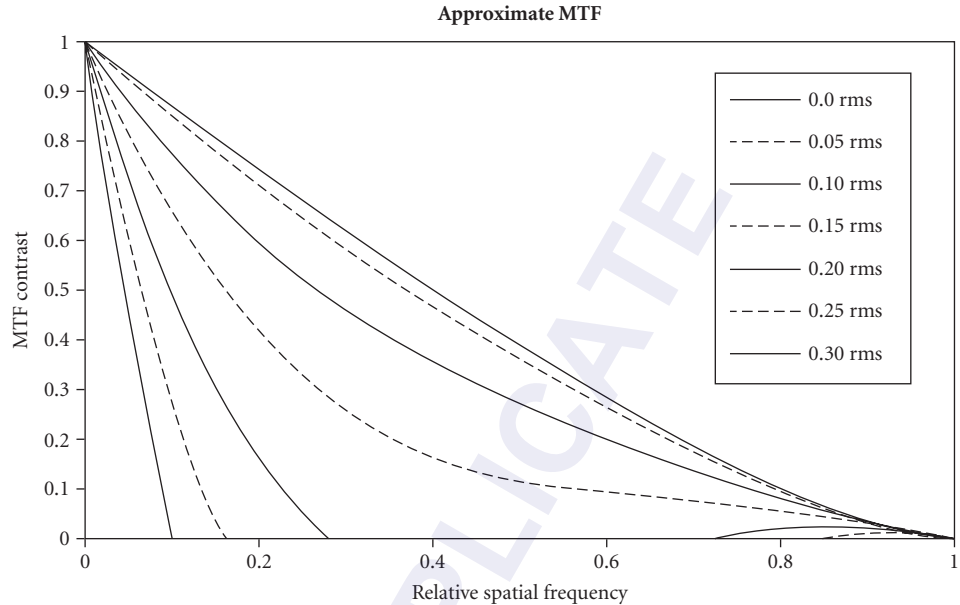


FIGURE 1 Approximate MTF curves from formula.

MTF as a single-value description of the state of correction of a system. Examination of the curves provides a method of communicating the specification to the system designer and fabricator. More detail on applying the rms wavefront error can be found in Chap.5, “Tolerancing Techniques.”

In addition, it must be pointed out that most imaging systems operate over a finite wavelength range. Thus the specification of “wavefront error” can be a bit fuzzy, but is usually meant to mean either the wavefront error at a specified wavelength, or a weighted average over the wavelength band. This should be mentioned when writing the specifications. In either case, the stated wavefront error contains a measure that communicates the extent of perfection required of the optical system performance.

4.3 PREPARATION OF OPTICAL SPECIFICATIONS

Gaussian Parameters

The gaussian parameters determine the basic imaging properties of the lens. They are the starting point for setting the specifications for a lens system. In principle these numbers can be specified precisely as desired. In reality, overly tight specifications can greatly increase the cost of the lens. Some of the important parameters are shown in Table 1.

Table 1 is a sample of reasonable values that may be placed upon a lens. A specific case may vary from these nominal values. The image location, radiometry, and scale are fixed by these numbers. A specific application will require some adjustment of these nominal values. In general, specifications that are tighter than these values will likely result in increased cost and difficulty of manufacture.

There is an interaction between these numbers. For example, the tight specification of magnification and overall conjugate distance will require a very closely held specification upon the focal length. The interaction between these numbers should be considered by the user to avoid accidentally producing an undue difficulty for the fabricator. It may be appropriate to specify a looser tolerance on some of these quantities for the prototype lens, and later design a manufacturing process

TABLE 1 Gaussian Parameters

Parameter	Precision Target	Importance	How Verified
Focal length	1–2%	Determines focal position and image size	Lens bench
<i>f</i> -number	<±5%	Determines irradiance at image plane	Geometrical measurement
Field angle	<±2%	Determines extent of image	Lens bench
Magnification	<±2%	Determines overall conjugate distances	Trial setup of lens
Back focus	±5%	Image location	Lens bench
Wavelength range	As needed; set by detector and source	Describes spectral range covered by lens	Image measurement
Transmission	Usually specified as >0.98 ⁿ for <i>n</i> surfaces	Total energy through lens	Imaging test, radiometric test of lens
Vignetting	Usually by requiring transmission to drop by less than 20% or so at the edge of the field	Uniformity of irradiance in the image	Imaging test, radiometric test of lens

to bring the production values within a smaller tolerance. However, it is appropriate that this be investigated fully at the design stage. The designer should be encouraged to consider the possibility of leaving adjustment possibilities in the lens design, so that a final assembly adjustment can bring the Gaussian parameters into the required tolerance range.

4.4 IMAGE SPECIFICATIONS

Image Quality

The rms wavefront error and the MTF for a lens have been discussed earlier as useful items to specify for a lens. Frequently, the user desires to apply a detection criterion to the image. This is always related to the application for the lens.

The most familiar functional specification that is widely used for system image quality is the resolution of the system. This is usually stated as the number of line pairs per millimeter that need to be visually distinguished or recognized by the user of the system. Since this involves both the physics of image formation as well as the psychophysics of vision, this is an interesting goal, but needs to be specified clearly to be of use to a designer. The reading of the resolution by a human observer is subjective, and the values obtained may differ between observers. Therefore, it is necessary to specify the conditions under which the test is to be carried out.

The type of target and its contrast need to be stated. The default standard in this case is the “standard” U.S. Air Force three-bar target, with high contrast, and a 6:1 ratio of bar length to width. This is usually selected as it will give the highest numerical values, certainly politically desirable. However, studies have shown that there is a better correlation between the resolution produced by a low-contrast target, say 2:1 contrast ratio, or 0.33 modulation contrast and the general acceptability of an image.

The resolution is, of course, related to the value of the MTF in the spatial frequency region of the resolution, as well as the threshold of detection or recognition for the observer viewing the target. If the thresholds are available, the above-described empirical relation between the rms wavefront error and the MTF can be used to estimate the allowable aberration that can be left in the system after design or fabrication.

In the case of a system not intended to produce an image to be viewed by a human, a specific definition of the required image contrast or energy concentration is usually possible. The signal-to-noise ratio of the data transmitted to some electronic device that is to make a decision can be calculated once a model for operation of the detector is assumed. The specification writer can then work backward through the required MTF to establish an acceptable level of image quality. The process is similar to that for the visual system above, except that the threshold is fully calculable.

In some cases, the fractional amount of energy collected by the aperture of the lens from a small angular source, such as a star, falling within the dimensions of a detector of a given size is desired. Such a requirement can be given directly to the designer.

Image Irradiance

The radiometry of the image is usually of importance. With an optical system containing the source, such as a viewer, projector, or printer, the usual specification is of the irradiance of the image in some appropriate units. Specifying the screen irradiance in watts per square centimeter or, more commonly, foot-lamberts, implies a number of optical properties. The radiance of the source, the transmission of the system, and the apertures of the lenses are derived from this requirement.

In the more usual imaging situation, the f -number and the transmission of the lens are specified. If the lens covers a reasonable field, the allowable reduction in image irradiance over the field of view must also be specified. This leads directly to the level of vignetting that can be allowed by the designer in carrying out the setup and design of the lens system.

There is an interaction between these irradiance specifications and the image quality that can be obtained. The requirement of a large numerical aperture leads to a more difficult design problem, as the high-order aberration content is increased in lenses of high numerical aperture.

An attempt to separate the geometrical aperture effects from the transmission of the components of the lens is accomplished through the t -number specification. Since the relative amount of irradiance falling on the focal plane is inversely proportional to the square of the f -number of a lens, the effect of transmission of the lens can be included by dividing the f -number by the transmission of the lens.

$$t\text{-number} = \frac{f\text{-number}}{\sqrt{t}}$$

where t is the transmission factor for the lens. The transmission factor for the lens is the product of the bulk transmission of the glass and the transmission factor for each of the surfaces. When this is specified, the designer must provide a combination of lens transmission and relative aperture that meets or exceeds a stated value.

Depth of Focus

The definition of the depth of focus is usually the result of a tolerance investigation. The allowable focal depth is obtained by determining when an unacceptable level of image quality is obtained. There is an obvious relation between the geometry of the lens numerical aperture and the aberrations that establishes the change of MTF with focal position. This effect can be computed for specific cases, or estimated by recognizing that the relation between rms wavefront error and focus shift is

$$W_{\text{rms}_{\text{def}}} = \frac{\delta l}{8\lambda(f\text{-number})^2}$$

which can be used in the above approximate MTF to provide an estimate of the likely MTF over a focal range.

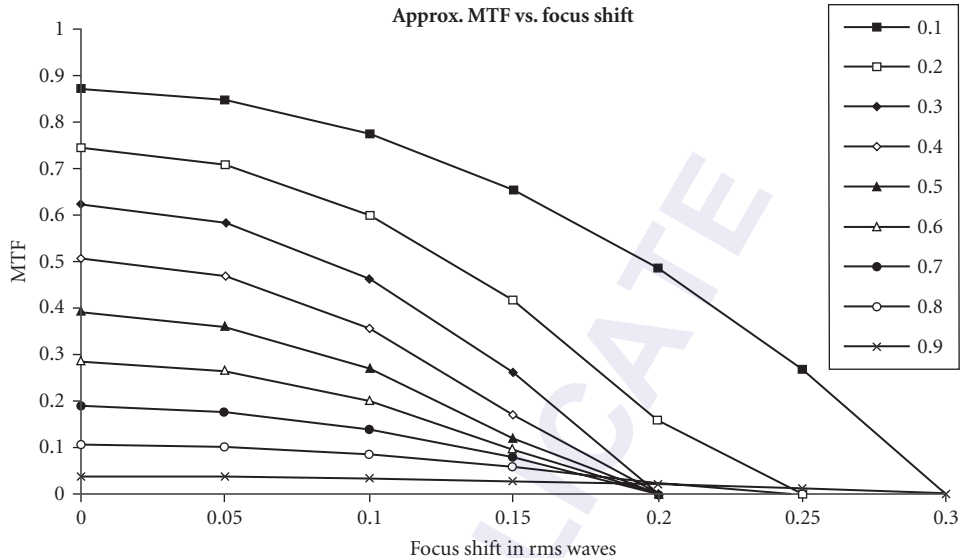


FIGURE 2 Approximate MTF as a function of focal position for various spatial frequencies.

If there is a basic amount of aberration present in the lens, then the approximation is that

$$W_{\text{rms, total}} = \sqrt{W_{\text{rms, def}}^2 + W_{\text{rms, lens}}^2}$$

leading to a calculation of the estimated MTF value for the given spatial frequency and focal position. As an example, Fig. 2 provides a plot of the focal position change of the MTF. The interpretation of this curve is straightforward. Each of the curves provides the MTF versus focus for a spatial frequency that is the stated fraction of the cutoff spatial frequency. The defocus is given in units of rms wavelength error, which can be obtained from reference to the appropriate formula. Using this approximate data, a specification writer can determine whether the requirements for image quality, f -number, and focal depth are realistic.

An additional consideration regarding depth of focus is that the field of the lens must be considered. The approximate model presented here is used at an individual field point. An actual lens must show the expected depth of field across the entire image surface, which places some limits upon the allowable field curvature. In general, it is the responsibility of the specification writer to establish the goal. It is the responsibility of the optical designer to determine whether the goal is realizable, and to design a system to meet the needs. In a sensible project, there will be some discussion between the designer and the engineer writing the specifications in order to avoid an unrealizable set of goals being set.

4.5 ELEMENT DESCRIPTION

Each element of the lens to be fabricated must be described in detail, usually through a drawing. All of the dimensions will require tolerances, or plus and minus values that, if met, lead to a high probability that the specified image quality goals will be met.

Mechanical Dimensions

The mechanical dimensions are specified to ensure that the element will fit into the cell sufficiently closely that the lens elements are held in alignment. This will be a result of tolerance evaluation, and must include allowances for assembly, thermal changes, and so on.

An important item for any lens is the interface specification, which describes the method of mounting the lens to the optical device used with it. For some items, such as cameras and microscopes, there are standard sizes and screw threads that should be used. In other cases, the specification needs to describe a method for coupling or mounting the optics in which there is a strain-free transfer of load between the lens and the mounting.

Optical Parameters

The optical parameters of the lens element relate to the surfaces that are part of the image-forming process. The radius of curvature of the spherical refracting (or reflecting) surfaces needs to be specified, as well as a plus or minus value providing the allowable tolerances. When tested using a test glass or an interferometer, the important radius specification is usually expressed in terms of fringes of spherical departure from the nominal radius. In addition, the shape of the surface is usually specified in terms of the fringes of irregularity that may be permitted.

When specifying a surface that will be measured on an interferometer, adjustment of focus during the test can be made. In this case, the spherical component of the surface, that is, the fringes of radius error, can be specified independently of the irregularity fringes that are applicable to the surface. When test glasses are to be used, the spherical component must be fabricated to within a small level of error to permit accurate reading of the irregularity component of the surface.

The cosmetic characteristics of the surface also need to be stated. The specification for this is as yet a bit imperfect, with the use of a scratch-and-dig number. This is actually intended to be a comparison of surface scratches with a visual standard, but is generally accepted to be in terms of a ratio, such as 20:10, which means, more or less, scratches of less than 2- μm width and digs of less than 100- μm diameter. This specification is described in MIL-O-13830, and is referred to a set of standards that are used for visual comparison to the defects on the surface. There have been several attempts to quantify this specification in detail, but no generally accepted standard has yet been achieved. A broader description of these specifications is found in International Standards 10110 and 9211, discussed later in this chapter.

Material Specifications

The usual material for a lens is optical glass, although plastics are becoming more commonly used in optics. The specification of a material requires identification of the type, as, for example, BK7 glass from Schott. Additional data upon the homogeneity class and the birefringence needs to be stated in ordering the glass. The homogeneity is usually specified by class, currently P1 through P4 with the higher number representing the highest homogeneity, or lowest variation of index of refraction throughout the glass. The method of specifying glass varies with the manufacturer, and with the catalog date. It is necessary to refer to a current catalog to ensure that the correct specification is being used.

Similar data should be provided regarding plastics. Additional data about transmission is usually not necessary, as the type of material is selected from a catalog which provides the physical description of the material. Usually, the manufacturer of the plastic will be noted to ensure that the proper material is obtained.

Materials for reflective components similarly have catalog data describing the class and properties of the material. In specifying such materials, it is usually necessary to add a description of the form and the final shape required for the blank from which the components will be made.

Coating Specifications

The thin film coatings that are applied to the optical surfaces require some careful specification writing. In general, the spectral characteristics need to be spelled out, such as passband and maximum reflectivity for an antireflection coating. Requirements for the environmental stability also need to be described, with reference to tests for film adhesion and durability. Generally, the coating supplier will have a set of “in-house” specifications that will guarantee a specific result that can be used as the basis for the coating specification.

4.6 ENVIRONMENTAL SPECIFICATIONS

Temperature and Humidity

Specification setting should also include a description of the temperature range that will be experienced in use or storage. This greatly affects the choice of materials that can be used. The humidity and such militarily favorite specifications as salt spray tests are very important in material selection and design.

Shock and Vibration

The ruggedness of an instrument is determined by the extent to which it survives bad handling. A requirement that the lens shall survive some specified drop test can be used. In other cases, stating the audio frequency power spectrum that is likely to be encountered by the lens is a method of specifying ruggedness in environments such as spacecraft and aircraft. In most cases, the delivery and storage environment is far more stressing than the usage environment. Any specification written in this respect should be careful to state the limits under which the instrument is actually supposed to operate, and the range over which it is merely meant to survive storage.

4.7 PRESENTATION OF SPECIFICATIONS

Published Standards

There are published standards from various sources. The most frequently referred to are those from the U.S. Department of Defense, but a number of standards are being proposed by the International Standards Organization.

Format for Specifications

The format used in conveying specifications for an optical system is sometimes constrained by the governmental or industrial policy of the purchaser. Most often, there is no specific format for expressing the specifications.

The best approach is to precede the specifications with a brief statement as to the goals for the use of the instrument being specified. Following this, the most important optical parameters, such as focal length, f -number, and field size (object and image) should be stated. In some cases, magnification and overall object-to-image distance along with object dimension will be the defining quantities.

Following this, the wavelength range, detector specifications, and a statement regarding the required image quality should be given. The transmission of the lens is also important at this stage.

Following the optical specifications, the mechanical and environmental requirements should be stated. The temperature and humidity relations under which the optical system needs to operate as

well as a statement of storage environment are needed. Descriptions of the mechanical environment, such as shock and vibration, are also important, even if expressed generally.

Other important pieces of information, such as a desired cost target, can then be included. Any special conditions, such as the need to be exposed to rapid temperature changes or a radiation environment, should be clearly stated. Finally, some statement of the finish quality for the optical system should be given.

In many cases, a list of applicable governmental specifications will be listed. In each case it is appropriate to ensure that these referenced documents are actually available to the individual who has to respond to the specification.

Use of International Standards

An important tool in writing specifications is the ability to refer to an established set of standards that may be applicable to the system being designed. In some cases, the development of specifications is simplified by the specification writer being able to refer to a set of codified statements about the environment or other characteristics the system must meet. In other cases, the established standards can be used as a reference to interpret parts of the specifications being written. For example, there may be a set of standards regarding interpretation of items included in a drawing.

Standards are an aid, not an end to specification. If the instrument must be interchangeable with parts from other sources, then the standards must be adhered to carefully. In other cases, the standards can serve as an indicator of accepted good practice in design or fabrication. It is the responsibility of the specification writer to ensure that the standard is applicable and meaningful in any particular situation.

At the present time there is growing activity in the preparation of standards for drawings, interfaces and dimensions, MTF, and other properties of optical and electro-optical systems. The efforts in this direction are coordinated by the International Standards Organization, but there are a number of individual standards published by national standards organizations in Germany, England, Japan, and the United States. The first major publications are ISO 10110, detailing preparation of drawings for optical elements and systems and ISO 9211, on optical coatings. Other standards on optical testing and environmental requirements are in draft form.

The ISO standards are expected to provide significant detail on various standards issues, and should become the principal guiding documents. At present, the standards documents that are most used in the United States are the various military specifications, or MIL-SPEC documents, that cover many different aspects of optical systems.

Information on published standards is available from the American National Standards Institute (ANSI), 11 West 42d Street, New York, NY 11036, or may be downloaded at www.webstore.ansi.org. A recent (2008) review of this website showed over 350 individual documents dealing with these issues, the majority of which deal with issues regarding fiber optical systems. Information on U.S. Department of Defense standards can be searched for through the National Search Engine for Standards at www.nssn.org. Additional information about worldwide standards is available at www.worldwidestandards.com.

4.8 PROBLEMS WITH SPECIFICATION WRITING

Underspecification

Failure to specify all of the conditions leaves the user vulnerable to having an instrument that will not operate properly in the real world. In many cases, the designer may not be aware of situations that may arise in operation that may affect the proper choice of design methods. Therefore, the design may not meet the actual needs.

The engineer developing a specification should examine all aspects of the problem to be solved, and carefully set the boundaries for the requirements on an optical system to meet the needs. All of

the pertinent information about the image quality, environment, and relation to other systems that may interact with the lens being specified should be considered. The specifying engineer should also review the physical limits on the image quality and ensure that these are translated into realistic values.

Overspecification

Specifying image quality and focal position requirements too tightly can lead to problems. Overspecification would seem to ensure that the needs will be met, but difficulties in meeting these requirements can lead to designs that are difficult and expensive to build. Achieving the goals can be costly and may fail. In such cases, the penalty for not quite meeting very tight specifications can be serious, both economically and technically.

Boundary-Limit Specification

In most cases, the statement of goals or boundaries within which a lens must operate is better than stating specific values. This leaves the designer with some room to maneuver to find an economic solution to the design. Obviously, some fixed values are needed, such as focal length, f -number, and field angle. However, too-tight specifications upon such items as weight, space, and materials can force the design engineer into a corner where a less desirable solution is achieved.

Negotiation of Specifications

Finally it is important to note that unless there is an existing closely defined set of established specifications for an specific optical device (such as a fiber optics coupler, for example) each specification is the product of a single individual or group and reflects the experience and understanding of that individual. The procuring official should be prepared in some cases to act as a negotiator between the engineer and the supplier to ensure that a reasonable and successful set of verifiable specifications has been stated.

4.9 REFERENCES

There are many useful references on optical specifications that deal with specific topics not directly covered by the general discussion in this chapter. The most useful suggestion is that the users having the task of setting specifications on a specific product or system use the massive capabilities of Internet search engines to look for specific data applicable to that task. A general Google search on "Optical Specifications" provided 360,000 hits, of which probably less than 10 will be applicable to any specific problem.

TOLERANCING TECHNIQUES

Robert R. Shannon*

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

5.1 GLOSSARY

a	relative tolerance error
BK7, SF2	types of optical glass
C to F	spectral region 0.486 to 0.656 μm
f -number	relative aperture as in $F/2.8$
n	refractive index
r_1, r_2, r_3, r_4	radii of curvature of surfaces
t_2	airspace thickness in sample lens
V	Abbe number or reciprocal dispersion
W	wavefront or pupil function
x	change factor
Δ	finite change in a parameter
δ	small change in a parameter
λ	wavelength

5.2 INTRODUCTION

Determination of the tolerances on an optical system is one of the most important parts of carrying out an optical design. No component can be made perfectly; thus, stating a reasonable acceptable range for the dimensions or characteristics is important to ensure that an economical, functioning instrument results. The tolerances attached to the dimensions describing the parts of the lens system are an important communication by the designer to the fabrication shop of the precision required in making the components and assembling them into a final lens.

*Retired.

The tolerances are related to but are not the same as the specifications. Setting specifications is discussed in Chap. 4. The tolerances are responsive to the requested system specifications and are intended to ensure that the final, assembled instrument meets the requested performance. The specifications placed on the individual lens elements or components are derived from the tolerances. Thus there is an interactive relation between the tolerancing activity and the setting of specifications. The system specifications drive the tolerances that need to be determined, and the tolerances are used in setting the specifications for the components of the system. The reality is that neither of these processes can be done fully independently.

At this point it is important to note that most optical design programs include a tolerancing utility that can be used to generate and distribute tolerances automatically once a few questions about goals have been answered by the designer. This appears to be a seductively simple process that is usually quite useful, but can be very disastrous if used uncritically by anyone who does not understand the basics of the process being carried out.

There are three principal issues in optical tolerancing. The first is the setting of an appropriate goal for the image quality or transmitted wavefront to be expected from the system. The second is the translation of this goal into allowable changes introduced by errors occurring on each component of the system. The third is the distribution of these allowable errors against all of the components of the system, in which some components of the optical system may partially or completely compensate for errors introduced by other components.

In this chapter, some basic approaches to distributing tolerances within an optical assembly are discussed. The examples will deal with tolerancing to meet a specified wave-front error and level of image quality. Similar principles apply to nonimaging optical systems, once the procedures necessary for relating errors in components or alignment to the specified operating requirements are established. The user will obviously have to adapt these approaches to the specific system being tolerated.

Optical versus Mechanical Tolerances

The tolerances on mechanical parts, in which a dimension may be stated as a specific value, plus or minus some allowable error, are familiar to any engineer. For example, the diameter of a rotating shaft may be expressed as 20.00 mm + 0.01/−0.02 mm. This dimension ensures that the shaft will fit into another component, such as the inner part of a bearing, and that fabricating the shaft to within the specified range will ensure that proper operational fit occurs. These tolerances may include the effect of environmental effects, such as operating temperature or lubrication needs, on the mechanical assembly.

Optical tolerances are more complicated, as they are generally stated as a mechanical error in a dimension, but the allowable error is determined by the effect upon an entire set of wavefronts passing through the lens. For example, the radius of curvature of a surface may be specified as 27.00 mm ± 0.05 mm. The interpretation of this is that the shape of the optical surface should conform to a specific spherical form, but remain within a range of allowable curvatures. Meeting this criterion indicates that the surface will perform properly in producing a focused wavefront, along with other surfaces in the optical system. Verification that the specific component tolerance is met is usually carried out by an optical test, such as examining the fit to a test plate. Verification that the entire system operates properly is accomplished by an assembled system test in which a specified image quality criterion is measured.

Basis for Tolerances

The process involved in setting tolerances begins with setting of the minimum level of acceptable image quality. This is usually expressed as the desired level of contrast at a specific spatial frequency as expressed by the modulation transfer function. Each parameter of the system, such as a radius of curvature of a surface, is individually varied to determine how large an error in each component is allowed before the contrast is reduced to the specified level. This differential change is then used to set the allowable range of error in each component.

In most cases, direct computation of the change in the contrast is a lengthy procedure, so that a more direct function, such as the rms wavefront error, is used as the quality-defining criterion. In other cases, the quantity of importance will be the focal length, image position, or distortion.

In nonimaging systems, the beam divergence or the uniformity of illumination after passage through the system may be the criterion of interest.

Relating the computed individual errors in the system to the tolerances to be specified is not always a simple matter. If there are several components, some errors may compensate other errors. Thus, it would be easy to assign too tight a tolerance for each surface unless these compensating effects, as well as the probability of a specific distribution of errors, are used in assigning the final tolerances.

Tolerance Budgeting

The method of incorporating compensation of one error by another, as well as the likelihood of obtaining a certain level of error in a defined fabrication process, is called *tolerance budgeting*. As an example, in a lens system it may be found that maintaining the thickness of a component may be easier than keeping the surfaces of the component at the right spherical form. The designer may choose to allow a looser tolerance for the thickness and use some of the distributed error to tighten the tolerance on the radius of curvature. In other cases, the shop carrying out the fabrication may be known to be able to measure surfaces well, but has difficulty with the centering of the lenses. The designer may choose to trade a tight tolerance on the irregularity of the surfaces for a looser tolerance on the wedge in the lens components.

Finally, the effect of a plus error on one surface may be partially compensated by a minus error on another surface. If the probability of errors is considered, the designer may choose to budget a looser tolerance to both surfaces.

This budgeting of tolerances is one of the most difficult parts of a tolerancing process, since judgment, rather than hard numbers, is very much a part of the budget decisions. It is advisable for the engineer carrying out the tolerance budgeting to do some modeling of the system performance using trial sets of parameter variations based on the tolerances that have been obtained. This verification serves as a method of ensuring that the tolerances are indeed reasonable and justified.

Tolerance Verification

Simply stating a set of allowable errors does not complete the integrated process of design and fabrication. The errors must be measurable. Measurement of length can be gauged, but has to be within the capabilities of the shop fabricating the optics. Measurement of error in radius of curvature requires the use of an interferometer or test plates to determine the shape of the surface. Measurement of the nonspherical component of the surface, or the irregularity, requires either an estimate from the test plate, or a computation of the lack of fit to a spherical surface based on measured fringes.

Finally, the quality of a completed lens must be measurable. Use of a criterion that cannot be measured or controlled by the shop or by the user is not acceptable. The contrast mentioned is not always measurable by the optical shop. The surface errors as measured by an interferometer or by test plates are common.

As shown in the Chap. 4 on specifications, the average or rms wavefront error can be related to the level of contrast, or modulation transfer function (MTF), that can be expected in the image. In addition, measurement of the final wavefront from an assembled optical system is most frequently obtained by an optical shop in a summary method by using an interferometer. For this reason, wavefront tolerancing methods have become the most commonly used methods of defining and verifying tolerances.

5.3 WAVEFRONT TOLERANCES

The rms wavefront error tolerancing method will now be used as an example of the approach to evaluating the tolerances required to fabricate a lens. An example which discusses the axial image tolerances for a doublet will be used to provide insight into the tolerancing of a relatively simple system. Most optical tolerancing problems are far more complex, but this example provides an insight into the methods applied. The specific example selected for this chapter is an airspaced achromatic

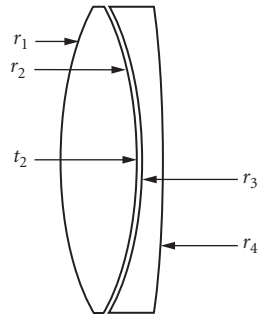


FIGURE 1 Drawing of the doublet lens used for tolerancing.

doublet using BK7 and SF2 glasses, $F/2.8$, 100-mm focal length. The lens design is nominally of moderately good quality, and is optimized over the usual C to F spectral range, with balanced spherical aberration, and is corrected for coma.

Figure 1 is a drawing of the sample doublet used in this chapter. The locations of the four radii of curvature and the airspace to be toleranced are indicated. The number of possible errors that actually can occur in such a simple lens is surprisingly large. There are four curvatures, two thicknesses, one airspace, and two materials that may have refractive index or dispersion errors. In addition to these seven quantities, there are two element wedge angles, four possible tilts, and four decenter possibilities, plus the irregularity on four surfaces and the homogeneity of two materials. So far, there are 21 possible tolerances that are required in order to completely define the lens. For interfacing to the lens mount, the element diameters, roundness after edging, and cone angle on the edges must be considered as well. More complex systems have far more possible sources of error. For the example here, only the four radii and the airspace separation will be considered.

Parameter Error Quality Definitions

The starting point for a tolerance calculation is the definition of a set of levels that may be used to define the initial allowable range of variation of the parameters in the lens. The magnitude of these classes of errors is determined by experience, and usually depends upon the type of fabrication facility being used. Table 1 presents some realistic values for different levels of shop capability.

These values are based on the type of work that can be expected from different shops, and serves as a guide for initiating the tolerancing process. It is obvious that the degree of difficulty in meeting the quality goals becomes more expensive as the required image quality increases.

Computation of Individual Tolerances

The individual tolerances to be applied to the parameters are obtained by computing the effect of some arbitrary but reasonable parameter changes upon the image-quality function. For the example doublet, if made perfectly with no errors in the individual components or assembly, the nominal amount of rms wavefront error at the central wavelength is 0.116 waves, rms. It is determined by the user from consideration of the needs for the application that the maximum amount of error that is acceptable is 0.15 waves, rms. Thus a distribution in allowable errors that results in no greater than about a 0.15 wavelength rms wavefront error would produce an acceptable system. The tolerancing task is to specify the tolerances on the radii of curvature and the separation between the component surfaces such that the goal is met.

TABLE 1 Reasonable Starting Points for Tolerancing a Lens System

Parameter	Commercial	Precision	High Precision
Wavefront residual	0.25 wave rms 2-wave peak	0.1 wave rms 0.5 wave peak	<0.07 wave rms <0.25 wave peak
Thickness	0.1 mm	0.01 mm	0.001 mm
Radius	1.0%	0.1%	0.01%
Index	0.001	0.0001	0.00001
V-number	1.0%	0.1%	0.01%
Homogeneity	0.0001	0.00001	0.000002
Decenter	0.1 mm	0.01 mm	0.001 mm
Tilt	1 arc min	10 arc sec	1 arc sec
Sphericity	2 rings	1 ring	0.25 ring
Irregularity	1 ring	0.25 ring	<0.1 ring

The reason for the choice of 0.15 wave rms is indicated by Fig. 2. The designer experimented with several choices of focus position to obtain a set of plots of the MTF for different amounts of error. This is not a completely general conclusion since the source of error produces an rms error which may not be the same for every source of error. But the samples permit the intelligent selection of an upper bound to the required error. In a lens with more sources of error, and with larger amounts of aberration in the basic design, setting up an example such as this is extremely important to avoid an error in the goal for the final image quality.

The first computation of the effect of nominal changes in the parameters on the rms wavefront error leads to the results in Table 2. (Only the radii and thickness are considered for this example.) The two columns for rms wavefront effect are first, for the aberration if no adjustment for best focus is made, and second, permitting the establishment of best focus after assembly of the lens.

Table 2 shows that the effect of a change in the parameter will have an effect proportional to the change, but that the factor relating the change to the resulting wavefront error is different for the various parameters. The amount of change permitted if the parameter is not compensated by allowing an adjustment for best focus is quite small. Any one of the parameters would have to be maintained within a range far less than the delta used in computation. The allowance of a compensating focal shift does greatly loosen the tolerance.

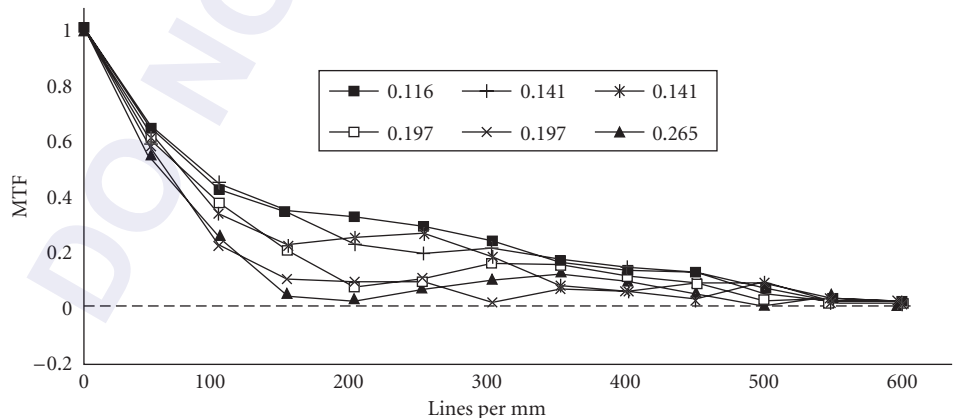
**FIGURE 2** Some examples of the effect of various rms errors on the MTF of the sample doublet. (The rms error is stated in wavelengths.)

TABLE 2 Finite Differentials for Computing Tolerances

Parameter	Delta	Rms Uncompensated	Rms Compensated
$r1$	0.1%	0.740	0.117
$r2$	0.1%	1.187	0.171
$r3$	0.1%	1.456	0.157
$r4$	0.1%	0.346	0.110
$t2$	0.025 mm	1.155	0.152

Since the acceptable goal is 0.15 waves, rms, the amount of change of an individual parameter to attain the acceptable level is about 50 times that of $r1$ and $r4$, but the change of 0.1 percent would be excessive for $r2$ and $r3$. The allowable change for $t2$ alone would be just about the delta of 0.025 mm.

Combination of Tolerances

No parameter in a lens lives alone. The effect upon the image will be the result of combining the effect of all of the errors. If the errors are uncorrelated, then the usual statistical summing of errors can be used. This states that the total amount of aberration produced by the errors can be found by using

$$W_{\text{rms}} = \sqrt{\sum_i W_i^2} = \sqrt{\sum_i \left(a_i x_i \frac{\partial W_{\text{rms}}}{\partial x_i} \right)^2}$$

where the sum is taken over the i parameters of interest. The x factors are the amount of change used in computing the change of wavefront error and the a factors are the relative amount of tolerance error allotted to each parameter in units of the delta used in the computation.

There are implicit assumptions in the application of this method to distributing tolerances. The principal assumption is that the fabrication errors will follow normal gaussian statistics. For many fabrication processes, this is not true, and modification of the approach is required.

For the example of the doublet, Table 3 can be generated to evaluate the different possibilities in assigning tolerances. The allowable change in rms wavefront error is 0.033 waves; thus the root sum square of all of the contributors must not exceed that amount.

In Table 3 the first column identifies the parameter, the second states the delta used in the computation, the third states the amount of wavefront error caused by a delta amount. The final two columns show different budgeting of the allowable error. Distribution 1 loosens the outer radii and the thickness at the cost of maintaining the inner radii very tightly. Distribution 2 tightens the outer radii and spacing tolerances, but loosens the inner radii tolerances. Depending upon the capabilities of the shop selected to make the optics, one of these may be preferable.

The interpretation of these statistical summations is that they are the sum of a number of different random processes. Thus, if the interpretation of each of the values given is the width of a normal distribution, which implies that 67 percent of the samples lie within that value, then 67 percent of the resulting combinations will lie within that range. If the interpretation is a two- or three-sigma value, the interpretation of the result follows similarly.

TABLE 3 Two Possible Tolerance Distributions for the Doublet

Parameter	Delta	Coefficient	Distbn. 1	Distbn. 2
$r1$	0.1%	0.00071	10.0	0.1
$r2$	0.1%	0.05440	0.25	0.4
$r3$	0.1%	0.04069	0.25	0.4
$r4$	0.1%	0.00069	10.0	0.1
$t2$	0.025 mm	0.03589	0.75	0.5
		rms change	0.033	0.033

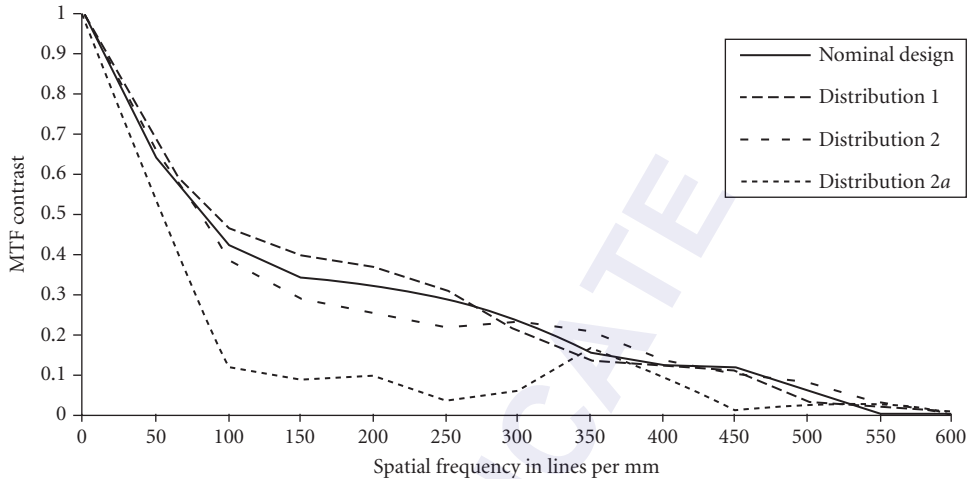


FIGURE 3 The resulting effect of different tolerance budgets for the sample doublet.

Figure 3 shows the effect of the various choices of allowable error distribution on the modulation transfer function of the lens. Either of the distributions stated in Table 3 provides an acceptable lens. For comparison, the allowable tolerances were doubled to provide distribution 2a, which is clearly not an acceptable lens. A spot check of some sample distributions is always relevant when doing tolerancing just to ensure that a reasonable relation between the tolerated system and acceptable image quality exists.

In any manufacturing process, the individual statistics will not necessarily follow a random rule. The interpretation is somewhat modified, but the principle still remains. In some cases, parameters may not be independent. For example, in the doublet there is a linking between the values of r_2 and r_3 that would loosen the tolerances if both are in error in the same direction. This could be used to advantage if the manufacturing process is carefully defined.

Use of Compensators

The use of compensators to loosen the tolerances was indicated above. An example of compensation for aberrations can be seen from a single lens element. If the first curvature is varied, both the element power and the spherical aberration from the element will change. However, a specific change of the second radius can restore the focal position and reduce the change of spherical aberration. Thus the tolerance allowed to the first curvature needs to take into account the possibility of a correlated or deliberate change in the second curvature. It is evident that the proper use of compensators can greatly loosen the tolerances applied to a surface.

A compensation that is frequently employed is the establishment of the correct focal position after assembly of the lens. If this procedure is followed, the individual tolerances on the surface of the elements can be loosened. It is obvious that the tolerancing and the development of a plan for fabrication and assembly must be coordinated.

5.4 OTHER TOLERANCES

Often, a particular optical parameter for the lens must be specified and maintained. Sometimes, for example, the focal length or back focus must be obtained within some tolerance. The computation of these paraxial constants for the lens can be made in the usual manner, and tolerances

obtained by using differentials relating each of the parameters to the quantity, such as the focal length, and then distributing the tolerances in a manner similar to that shown above for the doublet in Table 3.

Boresight

The pointing direction, or boresight, for a lens is sometimes of interest. Errors in boresight are usually due to asymmetric fabrication or mounting errors for the lens. In the simplest manner, one needs simply to trace an axial ray through the lens, and evaluate the direction of this ray as a result of introducing tilts and decenters of the surfaces, or of entire components. Tolerances on the lens parameters can be obtained by the procedure described for the doublet above, substituting the boresight error for the wavefront error.

Distortion

Distortion is the failure of the lens to provide a constant mapping from object space to image space. There are alternate interpretations for this error, which can have radial components due to symmetric errors in the lens, as well as tangential components from tilt and decenter of the lens components. This can be toleranced in the usual manner, but may be related to some general properties of the lens, such as the overall glass thickness of the components. In the simplest case, the tolerances upon distortion may be obtained from direct aberration computation. In complex cases, it may be necessary to compute the actual location of the centroid of the image as a function of image position and in the best image location.

Assembly

Assembly tolerances are related to the tolerances on image quality. The elements must be located and held in position so that the resulting image-quality goals are met. There are additional questions of allowing sufficient clearance between the elements and the lens barrel so that the elements can be inserted into the barrel without breaking or being strained by the mountings. These must be considered in stating the allowable dimensional range in the diameter, wedge, and concentricity of the edge of the lens.

5.5 STARTING POINTS

Shop Practice

Table 3, given as part of the sample tolerancing of the doublet, provides some estimates of the accuracies to which an optical shop may operate. These generic levels of error convey what is likely to be possible. The designer carrying out a tolerance evaluation should consult with probable fabrication shops for modification to this table. The tolerances that are ultimately assigned relate errors in the system to acceptable errors in the image. However, an understanding of shop practice is of great assistance in intelligent budgeting of tolerances.

Measurement Practice

Contemporary practice in optical fabrication and testing is to use interferometry to define wavefronts and surfaces. A convention that has become common in recent years is the use of polynomials

fitted to the wavefront as a method of describing the wavefront. There are several representations used, the most frequent of which is a limited set of Zernike polynomials. These ideally serve as an orthonormal set describing the wavefront or surface up to a specific order or symmetry. In tolerancing, the principal use for the coefficients of the Zernike set has been the easy computation of the rms wavefront error fitted to a given order. Thus the residual error in the system can be described after removal of low-order error such as focus or coma, which can sometimes be attributed to properties of the test setup.

5.6 MATERIAL PROPERTIES

The most important material is optical glass. Specification of the material usually includes some expected level of error in index of refraction or dispersion. In addition, glass is offered having several different levels of homogeneity of index of refraction. The usual range allows for grades of glass having index of refraction inhomogeneity ranging from ± 0.00001 to less than ± 0.0000001 within a single glass blank. It is usually assumed that this variation will be random, but the process of glass manufacturing does not guarantee this.

To place a tolerance upon the required glass homogeneity variation, the concept of wavefront tolerancing can be used. In general, the amount of wavefront error that can be expected along a glass path of length t through the glass is

$$\delta W = \frac{\delta N \times t}{\lambda}$$

For example, if a lens has a glass path of 5 cm but an error of 0.01 wavelengths is assigned to glass homogeneity, then the allowable glass homogeneity is about 0.0000013 within the glass. Thus precision-quality glass is needed for this application. For less glass path or looser tolerance assignment to glass homogeneity error, the required glass precision can be loosened. For a prism, the light path may be folded within the glass, so that an effective longer glass path occurs.

5.7 TOLERANCING PROCEDURES

The example of the doublet serves to illustrate the basic principles involved in determining the tolerances on a lens system. Most lens design programs contain routines that carry out tolerancing to various degrees of sophistication. Some programs are capable of presenting a set of tolerances automatically with only limited input from the designer. The output is a neat table of parameters and allowable ranges in the parameters that can be handed to the shop. This appears to be a quite painless method of carrying out a complex procedure, but it must be remembered that the process is based on application of a set of principles defined by the program writer, and the result is limited by the algorithms and specific logic used. In most cases, some trials of samples of the suggested tolerance distribution will suggest changes that can be made to simplify production of the lens system.

Direct Calculation

The preceding discussion describes methods used in calculating the tolerance distribution for a lens system. Frequently, tolerance determinations for special optical systems are required that either do not require the formal calculation described above, or may be sufficiently unusual that the use of a lens-design program is not possible. In that case, application of the principles is best accomplished directly.

The procedure is first to decide on a meaningful measure of the image quality required. In fact, the term "image quality" may require some broader interpretation. For example, the problem may be to optimize the amount of energy that is collected by a sensor in an optical communication system;

or the goal may be to scan a specific pattern with specific goals on the straightness of the projected spot or line during the scan.

The next step is to express the desired image quality in a numeric form. Usually, the rms wavefront error is the useful quantity. In some cases, other values such as the focus location of a beam waist or the size of the beam waist may be pertinent. In radiometric cases, the amount of flux within a specified area on the image surface (or within a specified angular diameter when projected to the object space) may be the pertinent value.

Once this is accomplished, the third step is to determine the relation between small changes in the parameters of the optical system and changes in the desired image quality function. This is usually accomplished by making small changes in each of the parameters, and computing the value of each differential as

$$\frac{dW}{dx_j} = \frac{\Delta W}{\Delta x_j}$$

where the right side is a finite differential. On occasion, the magnitude of this relation is nonlinear, and may require verification by using different magnitudes of change in the parameter.

This computation provides relations for independent, individual changes of the parameters. The possibility of compensation by joint changes in two or more parameters also has to be investigated. The best approach is to compute changes in the image-quality parameter in which specific coupling of parameters is included. For example, the differentials for variation of the curvatures of the two surfaces of a lens independently will be significantly different than the coupled changes of both surfaces simultaneously, either in the same or different directions.

Spreadsheet Calculation

The differentials must be combined in some manner to provide insight into the tolerance distribution. The best way to accomplish this is to develop a spreadsheet which allows simultaneous evaluation of the combination of errors using the equation for rms summation stated earlier. The use of spreadsheets for calculation is so common today that details need not be covered in this chapter.

Lens-Design Programs

The use of lens-design programs for tolerance calculation has become very widespread because of the proliferation of programs for use on the PC-level computer. The status of lens-design programs changes rapidly, so that any specific comments regarding the use of any program is sure to be out of date by the time this book appears in print. Suffice it to say that all of the principal programs have sections devoted to establishing tolerances. Usually the approach follows the procedures illustrated earlier in this chapter, with finite changes, or sometimes true computed analytical derivatives, used to establish a change table relating parameters to changes in the state of correction of the lens. In this case, the tolerances would be a listing of the allowed changes in the parameter to remain within some specified distance from the design values in aberration space. Some programs use a more complex approach where the allowable change in such quantities as the contrast value of the modulation transfer function at specified spatial frequencies is computed.

The distribution of tolerances is usually established according to the statistical addition rules given above. Some programs permit the user to specify the type of distribution of errors to be expected for various types of parameters.

As recommended above, it is strongly suggested that the user or designer not accept blindly the results of any tolerancing run but, rather, do some spot checking to verify the validity of the range of numbers computed. It is frequently found that alterations in the specified tolerances will occur as a result of such an investigation.

5.8 PROBLEMS IN TOLERANCING

Finally, it is useful to recite some of the problems remaining in establishing tolerances for a lens system. Even though the computational approach has reached a high level of sophistication for some lens-design programs, there are aspects of tolerancing that more closely approach an art than a science. The judgment of the designer or user of a tolerance program is of importance in obtaining a successful conclusion to a project.

Use of Computer Techniques

The use of a computer program mandates the application of rules that have been established by the writer of the program. These rules, of necessity, are general and designed to cover as many cases as possible. As such, they are not likely to be optimum for any specific problem. User modifications of the weighting, aberration goals, and tolerance image-quality requirements are almost always necessary.

Overtightening

The safe thing for a designer to do is to require very tight tolerances. This overtightening may ensure that the fabricated system comes close to the designed system, but the cost of production will likely be significantly higher. In some cases, the added cost generated by the overtight tolerances can raise the cost of the lens to the point where the entire project is abandoned.

The designer should consult with the fabricators of the optical system to develop an approach to assembly and testing that will allow the use of more compensating spacings or alignments to permit loosening of some of the tight tolerances.

Overloosening

A similar set of comments can be made about too-generous tolerances. In many schemes for production, these loose tolerances are justified by inclusion of an alignment step that corrects or compensates for cumulative system error. Too casual an approach to developing tolerances that require specific assembly processes, which are not fully communicated to the project, can result in a lens which is initially inexpensive to build, but becomes expensive after significant rework required to correct the errors.

Judgment factors

The preceding two sections really state that judgment is required. There is no completely “cookbook” approach to tolerancing any but the very simplest cases. The principles stated in this chapter need to be applied with a full knowledge of the relation between a change in a system parameter and the effect upon the image quality. In some cases, a completely novel relationship needs to be developed, which may include, for example, the connection between the alignment of a laser cavity and a nonlinear component included within the cavity. Finite difference calculations to obtain the output level can be developed using whatever computation techniques are appropriate. These values can be combined in a spreadsheet to examine the consequence of various distributions of the allowable errors.

5.9 REFERENCES

There are many useful references on optical tolerances that deal with specific topics not directly covered by the general discussion in this chapter. The most useful suggestion is that the users having the task of setting specifications on a specific product or system use the massive capabilities of internet search engines to look for specific data applicable to that task. A general Google search on “Optical Tolerances” provided 1,600,000 hits, of which probably less than 10 will be applicable to any specific problem.

This page intentionally left blank.

DO NOT DUPLICATE

MOUNTING OPTICAL COMPONENTS

Paul R. Yoder, Jr.

*Consultant in Optical Engineering
Norwalk, Connecticut*

6.1 GLOSSARY

a_G	acceleration factor
D_G	diameter of optic
E	Young's modulus
ID	internal diameter
K	constant factor
m	mass
OD	outer diameter
P	total preload
S_Y	yield stress
SPDT	single point diamond turning
t	thickness
Δ	deflection of spring or flange
ν	Poisson's ratio

6.2 INTRODUCTION AND SUMMARY

This chapter summarizes the techniques most commonly used to mount lenses, windows, small mirrors, and similar optical components as well as moderate-sized mirrors, and prisms within their mechanical surrounds to form optical instruments. Because of space limitations, mountings suitable for large (i.e., >85-cm diameter) optics are not discussed here. Two basic approaches for mounting optical components are considered: those in which the optic is held firmly against mechanical reference surfaces by applied forces (hard mounting) or those supported by benign means that do not inherently apply force (soft mounting). Descriptions of hard mountings include ones using threaded retaining rings, flanges, or springs while descriptions of soft

mountings include ones using flexures, elastomeric encapsulation, or bonding to mechanical pads. With either type of mounting, the required location and orientation of the optic relative to other portions of the instrument, that is, its alignment, is established during assembly in order to maximize performance. An important aspect of mounting design considered here is how the adverse influences of shock, vibration, temperature change, and moisture on alignment and system performance can be minimized. References cited here provide equations for designing and analyzing a large variety of mountings. Although we speak here of optics as if they are always made of glass and to mechanical parts (housings, cells, spacers, retainers, etc.) as if they are always made of metal, it should be understood that many of these mounting considerations also apply to other materials such as crystals, plastics, and composites.

6.3 MOUNTING INDIVIDUAL ROTATIONALLY SYMMETRIC OPTICS

Hard Mounting Techniques

In order to constrain an optic and preserve its alignment relative to other critical components of an optical instrument, hard mountings apply compressive forces to the glass at discrete locations or along line contacts. These forces, called preloads, are established during assembly and generally are of sufficient magnitude to hold the optic against appropriately located mechanical reference surfaces in the mount under all environmental conditions, including shock, vibration, and temperature changes. The magnitude of the preload P in N applied along any axis should be at least $9.81ma_G$, where m is the mass of the optic and any related components to be held by a single constraining means and a_G is the worst case acceleration expected to be encountered by the subassembly. This term a_G is understood to be a multiple of ambient gravity. If the optic is rotationally symmetric, glass-to-metal contact can be provided at the optic's cylindrical rim, at its ground bevels, or at its polished surfaces. Five degrees of freedom (three translations and two tilts) must be controlled. The sixth degree of freedom (rotation about the optic axis) is also adjusted and controlled in some cases to improve performance in the presence of residual optical wedge or if nonsymmetrical aspheric surfaces are involved. All six degrees of freedom must be constrained for noncircular optics, such as prisms.

The forces applied at the interfaces as well as those from gravity or imposed accelerations may distort the optical surfaces (thereby affecting performance) and introduce stress into the glass. Stress is known to cause birefringence, or, in extreme cases, damage to the optic—especially at low temperatures where shrinkage of the metal exerts maximum force on the glass. To minimize these adverse effects, forces must be kept within acceptable limits. Very few closed-form equations are available for predicting refracting or reflecting surface deformations due to applied forces. Finite element analyses are most frequently used for this purpose.^{1,2} Explanation of how this is done is beyond the scope of this presentation.

Relatively simple analytical means for estimating compressive and tensile stresses introduced by mounting forces are detailed in the literature.^{3,4} Most of these techniques are based on adaptations of standard formulations by Roark⁵ and Timoshenko and Goodier.⁶ The magnitude of the stress generated by a force depends not only upon the magnitude of that force, but also on the shapes of the surfaces in contact and Young's modulus and Poisson's ratio values for the glass and metal involved.

Statistical analyses backed by experimentation indicate that an optical component made by conventional high-quality grinding and polishing methods can usually withstand tensile stress as large as ~ 6.9 MPa without failure. This value is generally accepted as a "rule-of-thumb" tolerance for survival of the optic under stress.⁷ Optics made by "controlled grinding" techniques,⁸ polished and assembled with great care, and not scratched or otherwise damaged during use, might well survive long-term stress about 1.7 times greater.⁹

Under the more benign conditions of the operating environment (wherein the instrument must perform to specifications), survival is not a concern, but distortions of optical surfaces

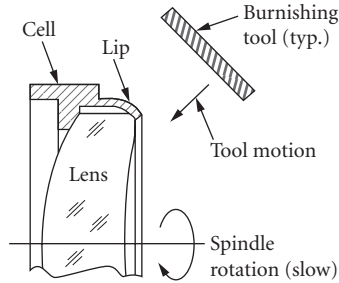


FIGURE 1 A small lens burnished into a cell made of malleable metal.

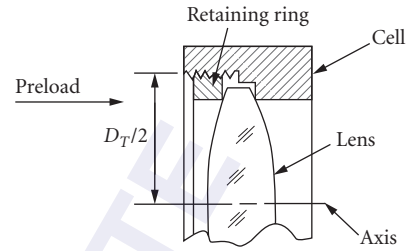


FIGURE 2 A lens preloaded in its cell with a threaded retaining ring.

due to mounting forces may degrade performance. High-performance optical systems and those using polarized light may also be especially sensitive to stress-induced birefringence. A commonly applied tolerance for stress in the glass in such cases is ~ 3.4 MPa. Analytical methods outlined in this chapter allow mounting stresses to be estimated to predict the potential success of a given design.

Burnished Mountings Figure 1 shows a very simple way to mount a lens element in a tubular cell. This cell has an internal shoulder against which the lens is to be held. The cell is mounted on a spindle, the lens is inserted and held in place gently, and the assembly is slowly rotated. The cell has a lip that extends beyond the rim of the lens. That lip is burnished with one or more hardened rod-shaped tool(s) over the edge of the lens as indicated in the right-hand view. The cell material must be malleable so it can be bent easily. Brass or annealed aluminum are common choices. The magnitude of the force, if any, introduced by the mounting cannot be quantified in this case because the bent metal tends to spring slightly away from the glass once tool pressure is removed. This type mounting is most suitable for use with small elements used in some endoscopes, simple microscope objectives, or low-cost cameras.

Mounts Using Threaded Retaining Rings This mounting, shown schematically in Fig. 2, is the type most frequently used to secure a lens element in its mount. Torque Q in N-mm applied to the ring with a wrench creates axial preload P to hold the lens against the shoulder very approximately as $5Q/D_T$, where D_T is the pitch diameter of the thread in millimeters.

The fit of the mating threads in the cell and retainer should be loose enough for the retainer to align itself to the centered lens surface; otherwise, lens alignment may be altered when the retainer is tightened. Such a fit may be specified as Class-1 or -2 per ANSI/ASME B1.1-2003.* During assembly, the lens should first be aligned in the cell and then held in place as the retainer is tightened to the required torque.

Mounts Using Annular Flanges Figure 3 shows a lens element preloaded against a shoulder by an annular flange that is deflected axially by a distance Δ from its nominal flat shape. Adapting an equation from Roark,⁵ the deflection Δ required to produce a given preload P in N equals $(K_A - K_B)P/t^3$ where t is the flange thickness in millimeters and the constants K_A and K_B are determined by the material properties and the dimensions a and b indicated in the figure.

For a given design, the required deflection may be obtained by customizing the thickness of the spacer located under the flange and should be at least 10 times larger than the resolution capability of the device to be used to measure the flange deflection at the time of assembly. As the flange is bent, stress is developed within that component. To prevent damage to the flange, its thickness should equal $K_C Pf_S/S_Y$ where the constant K_C depends upon the dimensions a and b and the flange material

*Unified Inch Screw Threads (UN and UNR Thread Form).

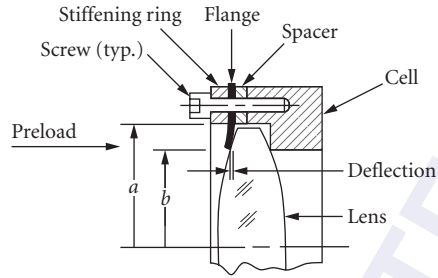


FIGURE 3 A lens preloaded in its cell with a deflected continuous ring flange.

properties. The quantity S_y is the yield stress of the flange material and f_s is the desired safety factor. The stiffening ring shown next to the flange maintains uniform flange deflection between the attaching screws.

A significant advantage of the annular flange as compared to the threaded retaining ring is that the flange can be calibrated before installation by measuring the actual preload developed as a function of deflection. Then, one can be quite confident that the preload on the lens is as stated by the above relationship when that flange is deflected by the specific distance Δ . This level of confidence cannot be achieved with a threaded ring.

Soft Mounting Techniques

Elastomeric Mountings A convenient technique for mounting a lens in a cell is to inject a continuous annular ring of an elastomeric material such as room temperature vulcanizing (RTV) sealing compound between the lens rim and the inside surface of the cell (see Fig. 4). This is sometimes called an “elastomeric ring mounting.” The thickness t_e equals $K_D D_G$, where D_G is the lens diameter and the constant K_D is determined from the material properties using a relationship attributed by Herbert¹⁰ to R. Vanbezooijen. The lens is then virtually free of radial stress at all temperatures. This is because the elastomer expands or contracts with temperature changes just enough to always fill the radial gap between the glass and the metal.

Some designs using the elastomeric ring approach also benefit from the fact that a continuous ring of this material effectively seals the lens to its mount so, if this subassembly forms part of the exterior skin of an optical instrument, leakage of gases and moisture through that interface is prevented.

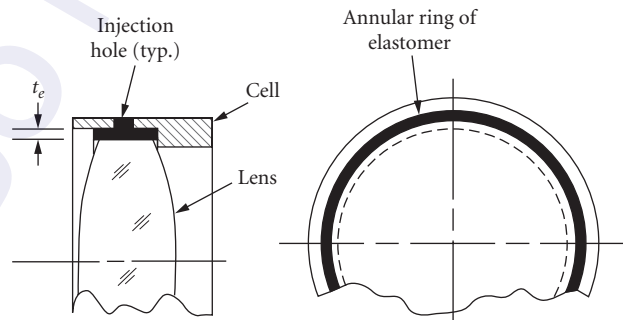


FIGURE 4 A lens supported in its cell by a continuous annular ring of elastomer.

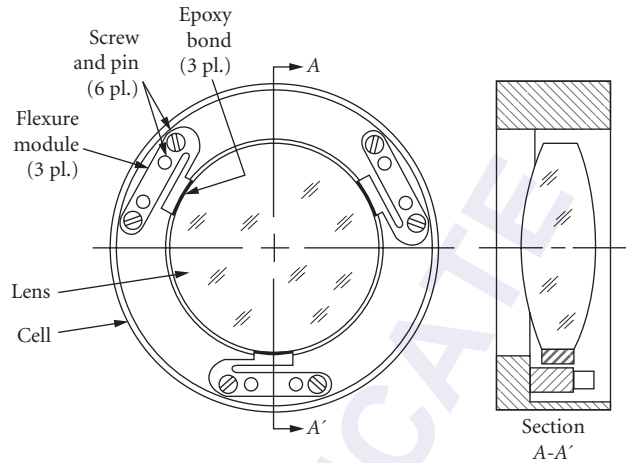


FIGURE 5 A lens bonded to three flexure modules attached to a cell.

Another type of elastomeric mounting for lenses uses discrete pads of elastomer located between the lens rim and the cell's inside surface. At least three such pads are needed to fully constrain the optic. They should be symmetrically distributed around the lens. The radial thicknesses of these pads can be sized as described above for the continuous ring.

Flexure Mountings High-performance optical systems require the optical axes of their lenses to be precisely centered mechanically with respect to some mechanical reference and to remain in that condition when the temperature changes. Because metal and glass components expand or contract with temperature at different rates, the optics may become decentered, tilted, or stressed when temperature changes occur in the above-described hard mountings. A properly designed support configuration using three or more symmetrically located identical flexures between the mount and lens rim will ensure that the lens stays as originally aligned and free of stress in spite of such changes.

Figure 5 shows a concept for a simple flexure mount design suggested by Ahmad and Huse.¹¹ Three identical flexure modules are made with narrow slots cut into them (by an electric discharge machining method) to form cantilevered flexure blades. Each blade has, at its free end, a curved pad shaped to interface with the lens rim. These modules are attached to the lens cell with screws passing through slightly oversized holes. In an alignment fixture, the optical axis of the lens is centered with respect to the axis of the cell. The modules are then adjusted to provide specific gaps between the pads and the lens rim and pinned in place. Epoxy is injected into those gaps and cured. Because the flexures are separate from the cell, they can be made from a material (such as titanium) with a higher yield stress than the cell. The cell is typically made of less expensive yet dimensionally stable material (such as stainless steel). More complex flexure designs and ones featuring a larger number of radial flexures have also been described.¹²⁻¹⁶ Because of their inherent flexibility, flexure mountings should be analyzed to determine their responses to externally imposed shock and vibration.

6.4 MULTICOMPONENT LENS ASSEMBLIES

Groups of lenses used in optomechanical assemblies typically are individually mounted and constrained in seats machined into a housing or are separated axially by spacer rings within a common cylindrical bore in the housing. It is important for those lenses to have a common optical axis and the correct axial airspaces within allowable tolerances. We here consider several ways in which such assemblies can be designed.

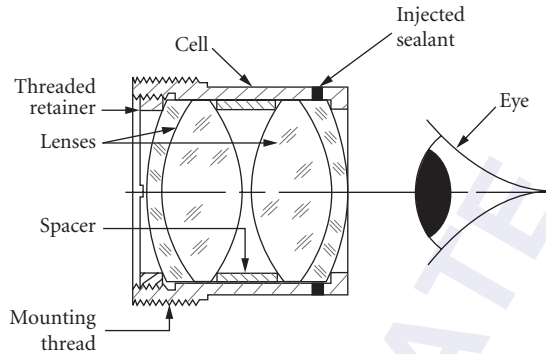


FIGURE 6 A telescope eyepiece with two lenses and a spacer assembled by the “drop-in” method.

Drop-In Assembly

If alignment requirements are not too demanding, lenses, housings, and spacers can be machined to reasonable tolerances and simply assembled without further machining or adjustment other than tightening a retainer to provide prescribed preload. Figure 6 shows an eyepiece for a low-power telescope assembled in this way. Radial clearances are typically ~ 0.075 mm, so individual lenses can easily be inserted into their seats. The eyepiece is configured to thread into a cylindrical opening in the telescope housing and to be focused by rotating the entire eyepiece. In some cases, the lenses are sealed to the cell and the threaded joint with the telescope also is sealed.

Many all-plastic lens assemblies used in consumer products are designed for swift drop-in assembly. The example shown in Fig. 7 is the objective for a rear-projection television system. Flexible tabs molded into the inside walls of both halves of the plastic housing project inward to form pockets for insertion of the three injection-molded plastic lens elements. Grooves (not shown) molded into the inside walls of the housings reduce stray light that otherwise could reduce contrast of the image. The housings are fastened together by self-tapping screws passing through flanges along each side, as indicated in the end view. Optical alignment relies on accuracy of the molding processes and is adequate for the application.

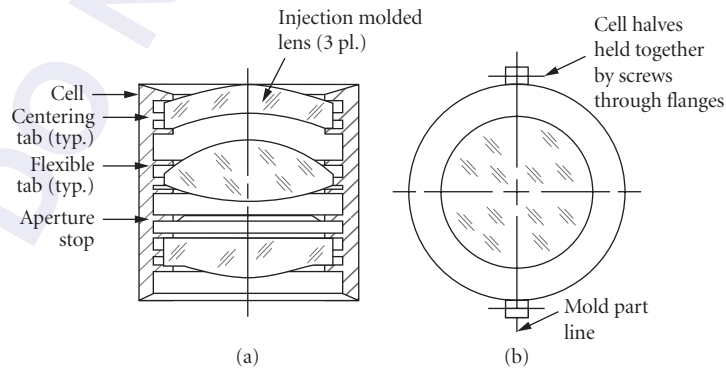


FIGURE 7 An all-plastic projection lens assembled by the “drop-in” method.

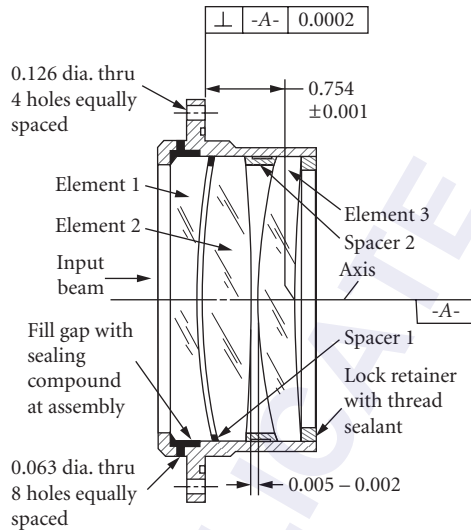


FIGURE 8 A telescope objective assembly with alignment resulting from tightly toleranced dimensions. Dimensions are in inches. (From Yoder.¹⁷)

Tightly Toleranced Assemblies

When higher performance is required, the dimensional tolerances are tightened and better optical alignment is achieved. Figure 8 shows the objective lens assembly for a military telescope.¹⁷ The three lenses are edged to fit the cell inside diameter with nominal radial clearances of 0.012 mm. All metal parts are made of stainless steel. The first spacer is made of sheet metal stock 0.025 ± 0.005 mm thick. It conforms to the spherical shapes of the adjacent lens surfaces under preload. The axial thicknesses of the lenses are toleranced to ± 0.005 mm. Residual optical wedge tolerances for the lenses are 12 arcsec. The beam deviation from these wedges is minimized at assembly by rotating (i.e., clocking) two lenses about their axes relative to the third lens to obtain maximum symmetry of the image of an on-axis artificial star. This image is observed with a microscope during alignment.

Lathe Assembly

A technique that is frequently used to obtain lens centration by minimizing radial clearances between lens ODs and cell IDs is called “lathe assembly” because it is done on a machinist’s lathe. The diameters and thicknesses of a selected set of lenses are measured and recorded. The required central air spaces and their tolerances are obtained from the optical system design. Actual lens surface radii are obtained from interferometric measurements made during lens manufacture. This data accompanies the lenses to the machine shop where a partially machined cell or housing is customized to provide conical or toroidal interfaces with the polished surfaces of that particular set of lenses and to provide all other required dimensions for the optomechanical assembly within the required tolerances. Radial clearances of ~ 0.005 mm can be achieved by this method. This clearance is adequate for careful assembly of the lenses into the cell.

Figure 9¹⁷ shows an air-spaced doublet lens subassembly created by this process. The individual lens seats are finish-machined at the time of assembly to fit those lenses. The length of the spacer (dimension E) and the location of the mounting flange (Datum B) relative to the front lens vertex

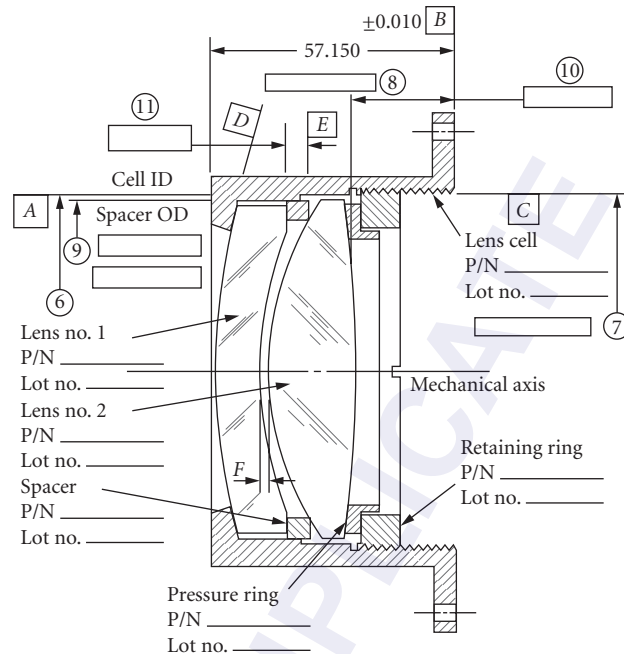


FIGURE 9 An air-spaced doublet assembled by the “lathe assembly” process. (From Yoder.¹⁷)

are also machined to produce the proper air space and overall length. Both lenses are secured by the retainer. The actual values of the numbered dimensions are recorded in the boxes. If required, a complete pedigree of that particular assembly can be established for future reference from the measured data and inspection reports. This type of construction is especially suited for applications involving high accelerations. Bayar described an aerial camera lens assembled by this method.¹⁸

“Poker Chip” Assembly

Figure 10 is a partial section view through a lens assembly that features seven lenses: four doublets and three singlets. Each lens, except the largest, was centered interferometrically to the mechanical axis of its cell OD and held in place in that cell with annular rings of epoxy nominally 0.381 mm thick. After the adhesive was cured, the axial thicknesses of the cells were final machined so all axial air spaces would be within design tolerances. The cell subassemblies Numbers 6 through 2, which had been machined to the same ODs within tight tolerances, were then inserted into the stainless steel housing and secured by Cell No. 1 that was threaded into the housing to act as a retainer. The largest lens (No. 12) was held directly in the housing by its own retainer. Accuracy of internal alignment was built into the assembly by the fabrication process.¹⁹ This type of construction is frequently referred to as “poker chip” assembly because the individual lens/cell subassemblies are stacked on top of each other inside the housing.

Lenses Adjusted at Assembly

Many complex lens assemblies to be used in very high-performance applications such as micro-lithographic projection systems need positional adjustment of a few carefully selected elements at the final stage of assembly. This is because application of the best possible optical and mechanical

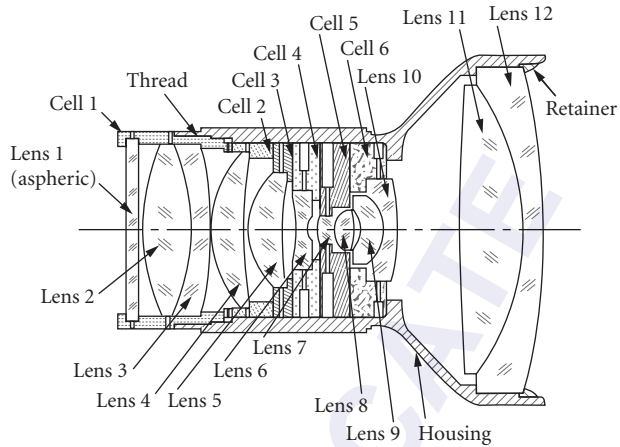


FIGURE 10 A projection lens comprising a stack of “poker chip” lens/cell subassemblies inserted into the bore of the mount. (From Fischer.¹⁹)

manufacturing processes and extremely tight dimensional tolerances cannot make the lenses and mechanical parts accurately enough to obtain the full level of performance required by the application. An example is shown schematically in Fig. 11.²⁰

This optomechanical system comprises twelve air-spaced “poker chip” subassemblies, stacked on top of each other with custom-made spacers placed between lapped coplanar pads on the cell faces

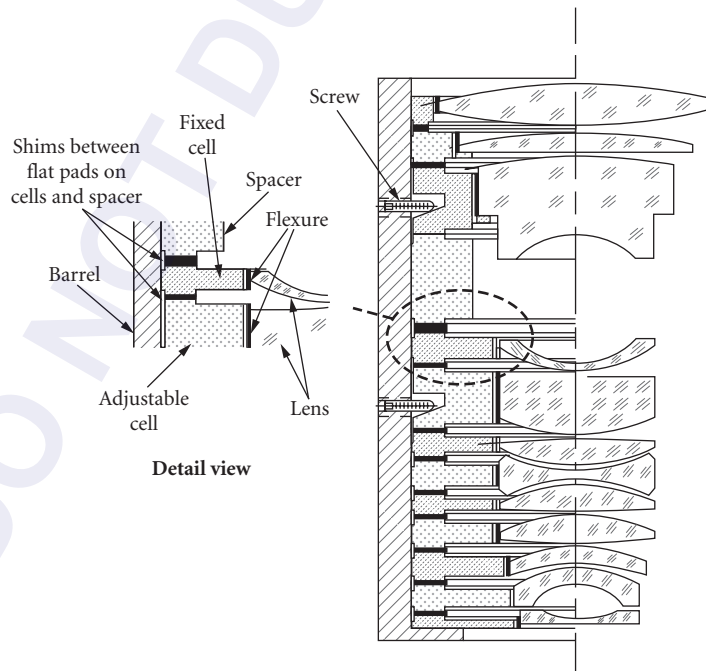


FIGURE 11 Partial section view of a “poker chip” lens assembly with two lenses adjusted after assembly to optimize performance. (From Yoder.⁴)

to control axial air spaces. The lenses are mounted on flexures machined directly into the interior surfaces the cells. Ten of the subassemblies fit closely inside the stainless steel barrel. Two subassemblies are adjustable laterally in orthogonal directions from the outside. Optical performance of the system is measured interferometrically in near real time, while the adjustable lenses are moved very slightly until optimal performance is achieved. The adjustment mechanisms are then locked and the lens is installed into the microlithography system.

Determination of which lens elements to move in a given optical system to correct residual aberrations is a job for the lens designer working with mechanical engineers and metrology experts who help decide how to incorporate the needed mechanisms and to conduct the necessary tests. The sensitivities of spherical, coma, astigmatism, and distortion aberration contributions from each lens to lateral and axial displacements are determined by raytracing. The ideal candidates for correcting each aberration are lens shifts that modify that aberration significantly, but do not excessively affect the other aberrations. The results are reviewed to determine which lens movements are best to minimize each aberration. Williamson²⁰ outlined a procedure in which the aberration contributions of the optical system shown in Fig. 12a for specific axial and lateral displacements of each element are plotted as shown in Fig. 12b and c, respectively. Using phase-measuring ultraviolet interferometry

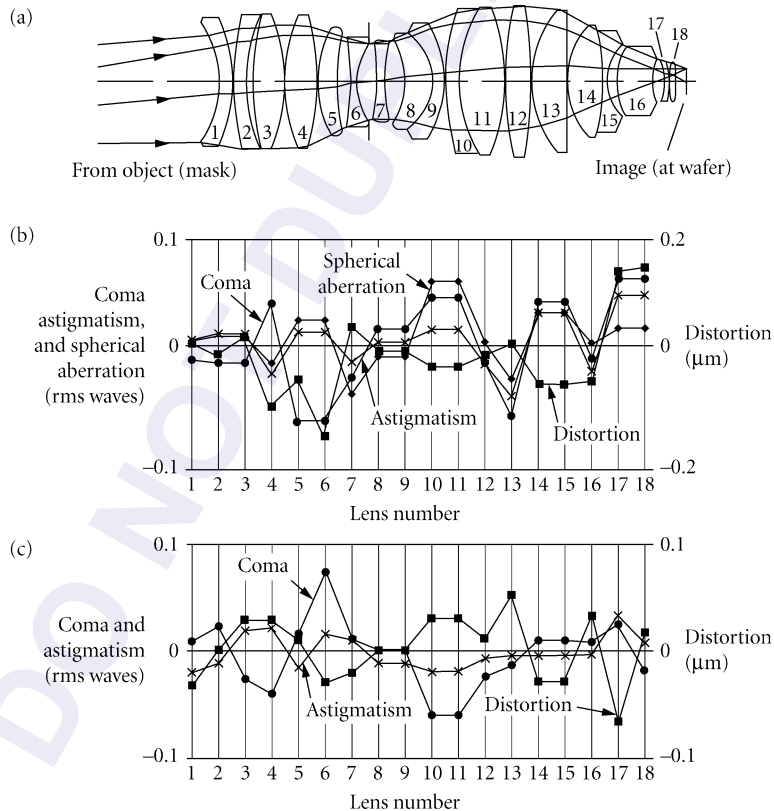


FIGURE 12 (a) Optical schematic of an 18-element lithographic projection lens. (b) The effects on aberrations of individually displacing each element axially by $25\ \mu\text{m}$. (c) Similar effects of displacing each element laterally by $5\ \mu\text{m}$. The best lenses to adjust for optimum system performance can be determined. (From Williamson.²⁰)

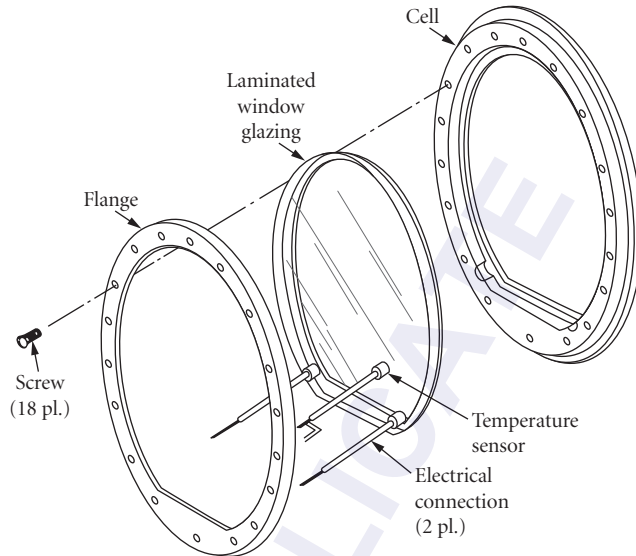


FIGURE 13 Exploded view of a heated window assembly used in a military application. (Courtesy of Goodrich Corporation, Danbury, CT)

as the image quality monitor, all these aberrations can be significantly reduced by iteration of lens movements and the final system performance of production lens assemblies greatly enhanced.

6.5 MOUNTING WINDOWS AND DOMES

Small circular windows usually are secured in a mount with a threaded retainer or by an elastomeric ring. Noncircular ones are best held in place with elastomer. The continuous flange-mounting method can be used to advantage with larger windows. The one shown in Fig. 13 has an elliptical aperture of 20.32×30.5 cm.⁴ The electrical connections shown provide current to a conductive coating on a buried surface, which keeps the window free of fog in high-humidity situations. The flange preloads the window into an aluminum cell. An elastomeric sealant is injected into a groove around the window's rim to seal it to the cell. The cell is sealed to the instrument housing with a gasket or an O-ring.

Figure 14 shows typical mountings for deeply curved spherical windows, called *shells* or *domes*. That in Fig. 14*a* is sealed and secured with a Neoprene gasket clamped in place by a flange²¹ while that in Fig. 14*b* is secured and sealed with a continuous ring of elastomer.⁴ In some more elaborate designs, an elliptically shaped sapphire dome is brazed with special metallic alloys to a titanium mount.²²

6.6 MOUNTING SMALL MIRRORS AND PRISMS

General Considerations

The appropriateness of designs for mechanical mountings for small mirrors and prisms depends upon a variety of factors including: tolerable rigid body movement of the optic and distortion of its reflecting and/or refracting surface(s); the magnitudes, application locations, and directions

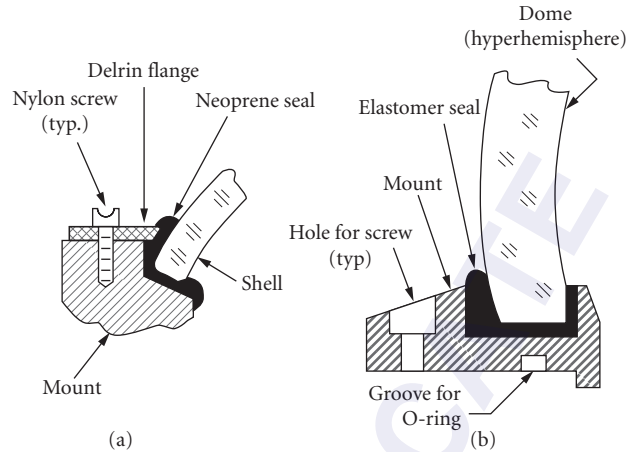


FIGURE 14 Typical mountings for (a) a thin shell and (b) a hyperhemispherical dome. [(a) From Vukobratovich.²¹ (b) from Yoder.⁴]

of forces tending to move the optic with respect to its mount; steady-state and transient thermal effects (including gradients); the sizes and kinematic compatibility of interfacing optomechanical surfaces; and the rigidity and long-term stability of the structure supporting the optic. In addition, the designs must be compatible with assembly, maintenance, package size, weight, and configuration constraints, as well as being cost effective. The representative mounting designs described in the following sections illustrate proven mounting techniques.

Mechanically Clamped Mountings

Figure 15 shows a simple means for attaching a first surface flat mirror to a mechanical bracket.²³ Three cantilevered springs press the reflecting surface against three pads that have been lapped coplanar. The contacts between the springs and the mirror's back face are directly opposite the pads to minimize bending moments. This design constrains one translation and two tilts. Translations in the plane of the reflecting surface can most easily be constrained by dimensioning the spacers supporting the springs so as to just clear the rim of the mirror at minimum temperature. Rotation in that plane usually does not need to be constrained. Given the number of springs N , the spring material's Young's modulus E_M , its yield stress S_Y , and its Poisson's ratio ν_M , the spring lengths L and widths b , and an appropriate safety factor f_S , the spring thickness t that will provide a total preload P to the mirror is determined as $[K_{S1}PLf_S/(bS_YN)]^{1/2}$. The length of the spacer located under each spring is chosen to cause that spring to be deflected from its flat condition by Δ equal to $(K_{S2}L^3)/(1 - \nu_M^2)(E_Mbt^3N)$. In these relationships, K_{S1} is 4 and K_{S2} is 0.75.

Elastomeric Mountings for Mirrors

Small mirrors can often be mounted in the manner illustrated by Fig. 4 for a lens. In applications where the optic does not need to be sealed in place with a continuous ring of elastomer, three or more discrete pads located between the lens rim and the cell ID can support it. Vanbezoijen's equation is again used to determine the pad thicknesses.¹⁰ The lateral dimensions of the pads have, in some designs, been determined by finite element analysis that predicts the dynamic response of the subassembly to vibration inputs from the environment.²⁴

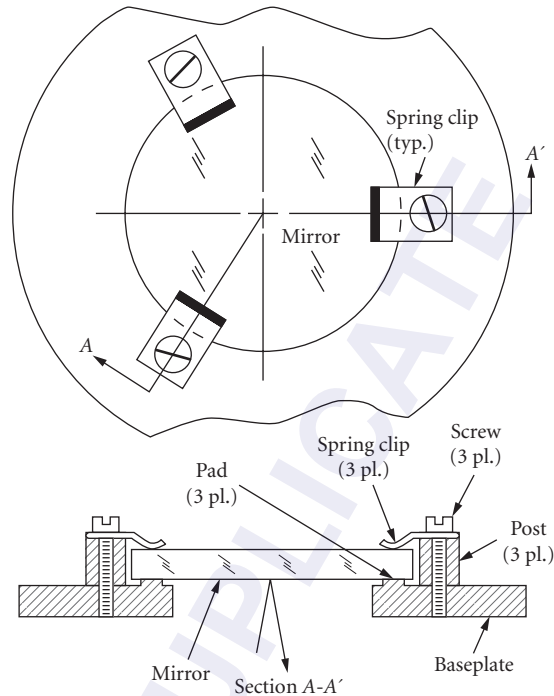


FIGURE 15 A simple mounting for a flat mirror preloaded with deflected cantilevered springs. (From Yoder.²³)

Spring-Loaded Mountings for Prisms

A spring-loaded mounting for a prism is illustrated in Fig. 16.²⁵ Here, a penta prism is preloaded against three coplanar pads on the baseplate by three cantilevered springs supported by posts with spacers machined to produce the necessary spring deflection and resultant preload, as described earlier for a mirror mounting. Constraint in the plane parallel to the pad surfaces is provided by a single spring (called a straddling spring) that is supported at each end and presses against the end of the prism. The dimensions and deflection of this spring are chosen to preload the prism against three locating pins that are pressed or threaded into strategically located holes in the baseplate. The relationships for t and Δ given in Sec. “Mechanically Clamped Mountings” apply also to the straddling spring, but K_{s1} equals 0.75, $N = 1$, and $K_{s2} = 0.0625$. How the direction of the force exerted by the straddling spring can be optimized to nearly equalize the stresses created in the prism at the interfaces with the pins has been explained in the literature.²⁶

Bonded Mountings for Small Mirrors and Prisms

A widely used and successful technique for mounting small mirrors and prisms is to bond them directly to a plate or bracket with an adhesive such as epoxy. Any alignment adjustments that are needed should be built into the mount rather than into the glass-to-metal joint. Figure 17 shows a typical mirror mount of this type.²³ It has proven satisfactory for cases where the diameter-to-thickness ratio for the mirror substrate is at least 6:1. The mirror should then be stiff enough not to be excessively

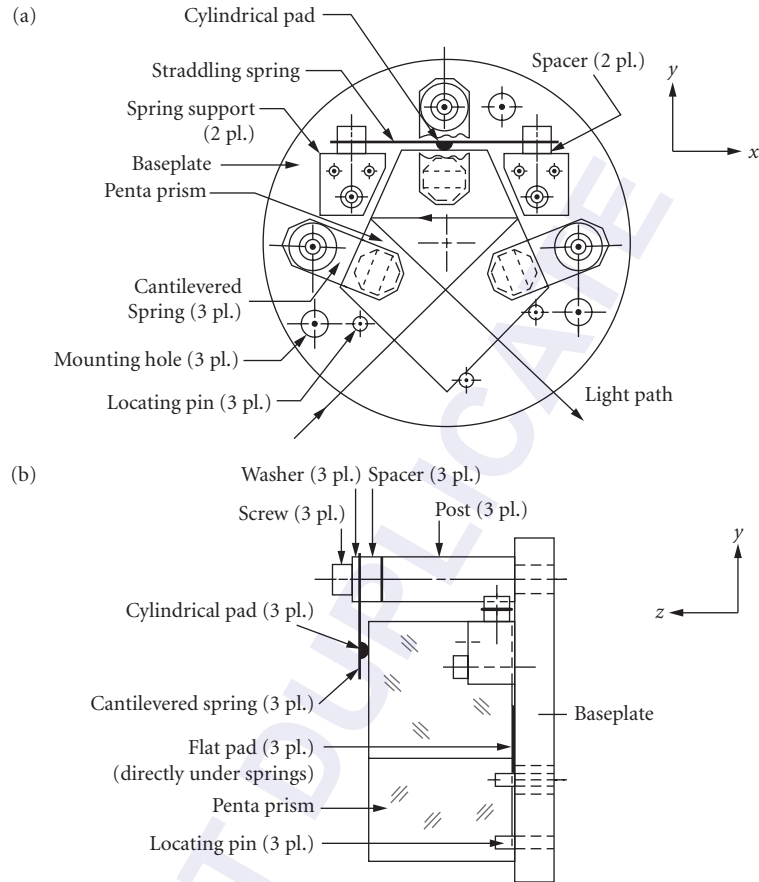


FIGURE 16 A penta prism preloaded against lapped pads on a baseplate with cantilevered and straddling springs. (From Yoder.²⁵)

distorted by shrinkage of the adhesive as it cures. Most prisms are thick enough that they are not distorted by this shrinkage.

All adhesive bonds need to have sufficient area for the joint to be strong enough to support the optic under all anticipated levels of acceleration. The minimum bond area Q_{MIN} is $9.81ma_g f_s / J$, where J is the strength of the cured adhesive joint and all other terms are as previously defined.²⁷ Bonding should be done on a fine ground surface of the optic for maximum joint strength. A typical value for J for a two-part epoxy such as 3M 2216B/A with bond thickness of 0.100 ± 0.025 mm is ~ 17.2 MPa. For conservative design, the factor f_s should be ~ 4 . Successful bonding requires careful cleaning of the surfaces to be bonded and adequate curing time. The adhesive manufacturer's recommendations should be followed unless tests indicate otherwise for a specific application.

The 29-mm aperture roof penta prism in Fig. 18 is bonded in cantilevered fashion to a bracket nominally oriented vertically. The circular bond area is adequate to withstand a severe military shock and vibration environment. Some designs work better if the prism is supported from both sides. Figure 19 shows one way to do this. It was adapted from Beckmann.²⁸ The mount is designed with two arms, one of which has a hole bored through it. The prism is supported by a fixture in the proper location and orientation relative to the mount and epoxy bonded to the flat pad on the left arm. A plug made of the same

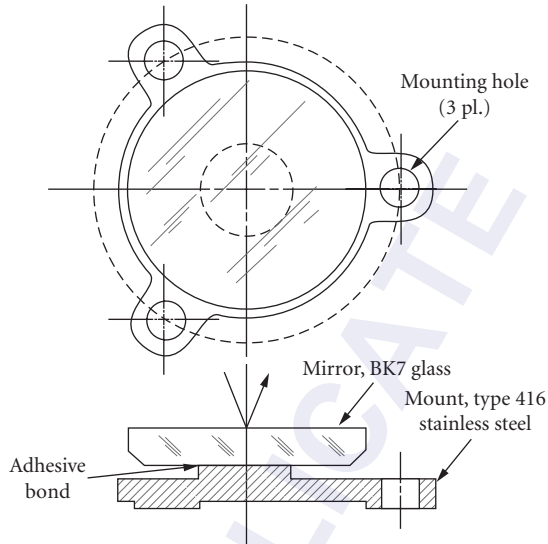


FIGURE 17 A flat mirror bonded on its back to a pad on its mount. (From Yoder.²³)

metal as the mount is then centered in the hole in the right arm and bonded to the right surface of the prism. When those bonds have cured, the plug is bonded into the right arm.

Flexure Mountings for Small Mirrors and Prisms

Circular mirrors as large as ~15-cm diameter have been successfully mounted on flexures in the general manner shown in Fig. 5 for a lens. Usually, these are image-forming mirrors, perhaps aspheric, that need to have constant centration relative to a system axis.

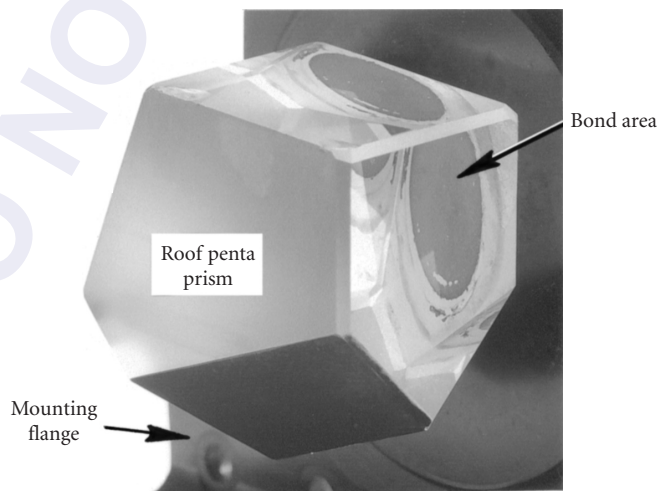


FIGURE 18 A roof penta prism bonded to a pad on a mounting flange.

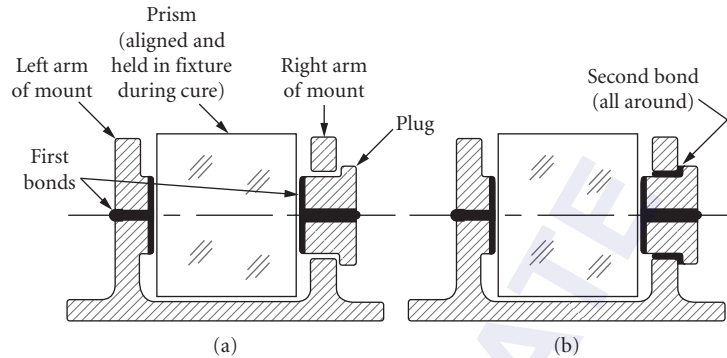


FIGURE 19 Concept for supporting a prism from both sides. (a) Prism bonded to left arm of mount and plug bonded to prism. (b) Plug bonded into right arm of mount. (Adapted from Beckmann.²⁸)

Prisms intended for use in relatively benign environments can be mounted on flexures. One way to do this is by attaching the prism to instrument structure through three posts with integral flexures at each end. Figure 20 shows a large multiple component Zerodur prism mounted in this manner. It has two wing prisms optically contacted to a third (base) prism, to which the flexures are bonded. The wing prism surfaces are perpendicular to each other and form a 15.2-cm-wide roof mirror that is inclined at 45° to the vertical. The reflected image is inverted horizontally as the optic turns the incident beam axis 90° . The orientations of three of the flexure joints are

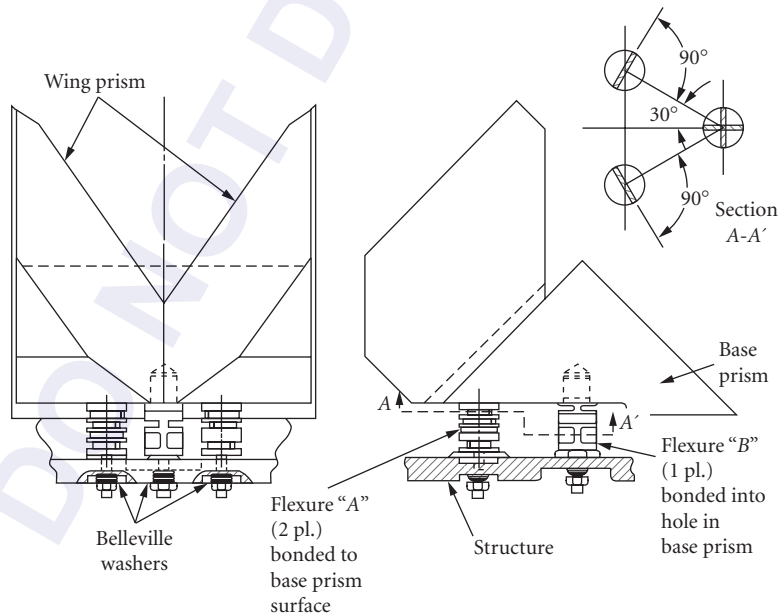


FIGURE 20 Optomechanical configuration of a large prism assembly with three flexure mounting posts to isolate the optic from dimensional changes under temperature changes. (Courtesy of ASML Lithography, Wilton, CT.)

as indicated in the section view A-A'. If attached to a structure that expands or contracts more than the prism as the temperature changes, the flexures simply bend very slightly and prevent the introduction of mounting forces that could distort the reflecting surfaces and interfere with performance of the optical system.

6.7 MOUNTING MODERATE-SIZED MIRRORS

General Considerations

The simple mirror mountings described earlier are not satisfactory for mirrors larger than about 15-cm diameter because they are too flexible to be treated as rigid bodies. The important criteria for selecting a suitable mounting are orientation with respect to gravity, performance level required, substrate material stiffness, and weight limitations. The mounting for a mirror to be used in a fixed horizontal- or vertical-axis orientation can be figured during polishing to compensate for gravity effects. Variable orientation applications require mounts that change their force distribution with inclination to keep surface deflections within tolerance. Both axial and radial supports are required. Mirrors to be used in space have the added requirement of release of gravitational force after being fabricated, tested, and installed into the instrument in a normal gravity environment. Choice of mounting depends strongly on the substrate configuration. Weight constraints generally lead to solid substrates with shaped back surfaces or ones built-up from multiple parts that are attached together. We here describe a few typical ways to support mirrors of various shapes as large as ~85 cm. Designs appropriate to both nonmetallic and metallic mirrors are considered.

Substrate Configurations

Figures 21*b* through *e* shows half-section views of four first-surface mirror solid substrates of the same diameter and material with concave surfaces of the same radius. Their back surface shapes differ and reduce the mirror weight as compared to a flat-back baseline design (Fig. 21*a*). All these mirrors, except one, can be supported within the telescope housing on a hub passing through the mirror's central perforation. For example, see Fig. 22. Here, the hub has a toroidal-shaped locating land that supports the 41-cm-diameter meniscus-shaped mirror radially and a shoulder that locates it axially. The radial support lies in the mirror's neutral plane where fore and aft bending moments are balanced. A threaded retaining ring provides axial preload. To focus, the locating ring is moved on the hub and secured with the clamping ring. The substrate configuration from Fig. 21 that cannot be hub mounted is the double arch configuration (Fig. 21*e*). It is best supported on flexures at three or more points spaced equally around the zone of greatest thickness.

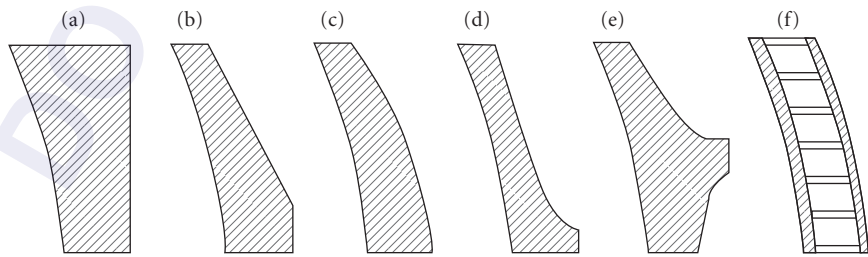


FIGURE 21 Sectional views of baseline concave-plane (a) and lightweighted mirror substrates (b) through (e) with contoured backs. Figure (f) shows a built-up substrate configuration.

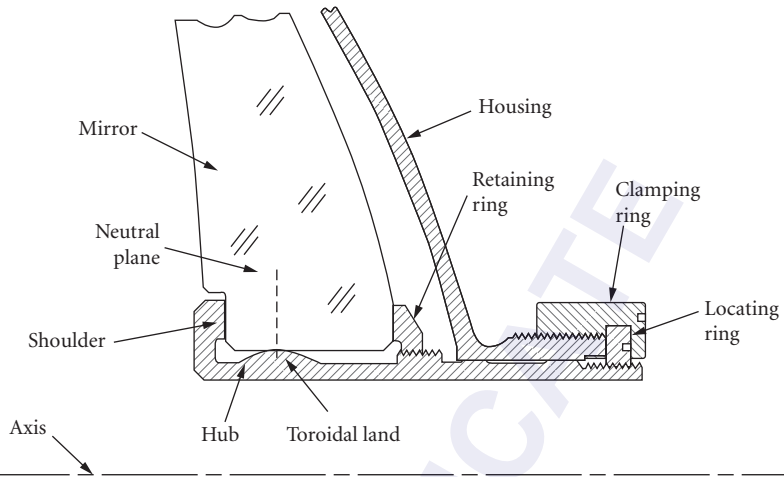


FIGURE 22 Hub mounting for a meniscus-shaped telescope mirror. Focus adjustment means is illustrated.

Lighter-weight mirror constructions typically employ built-up substrates such as that shown in Fig. 21*f*. A very successful type is the monolithic meniscus construction illustrated by Fig. 23. Such mirrors are usually made of Corning ULE. Strips of the material form the webs of a core to which front and back facesheets are fused. All joints in the core also are fused together. The spacing of the webs is large except at locations where axial and radial supports attach to the substrate. There, the spacing is considerably smaller to increase strength.

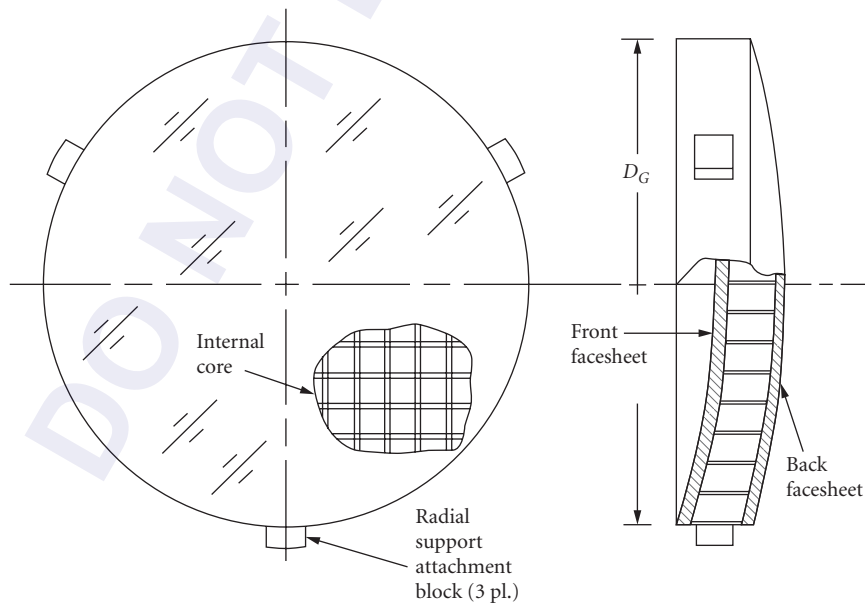


FIGURE 23 A completely fused (monolithic) built-up lightweight mirror substrate.

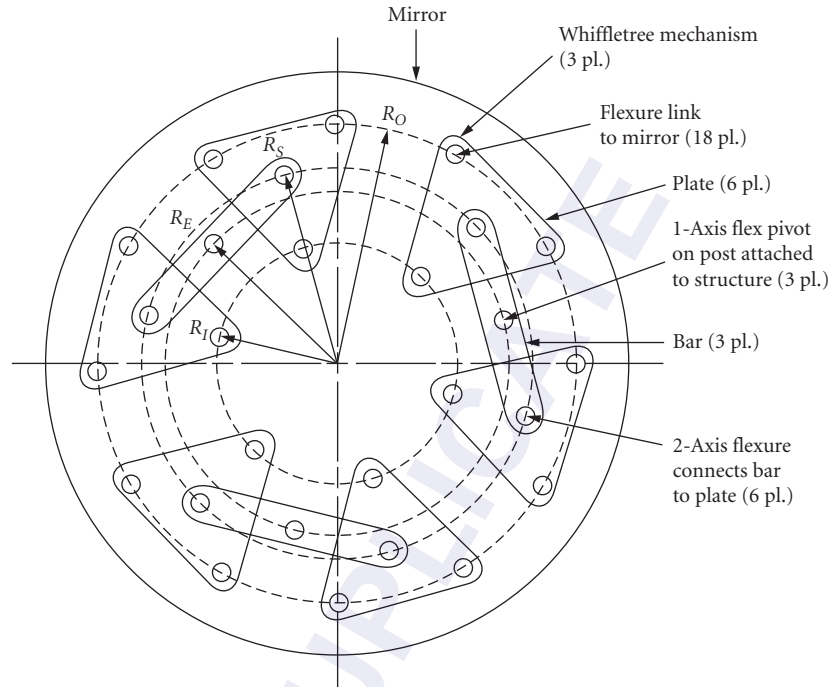


FIGURE 24 An 18-support Hindle-type mirror mount supporting the optic at multiple points on rings of radii R_O and R_I from three posts attached to structure at radius R_E . The whiffletree plates are centered on the ring of radius R_S .

Lever Mechanism Mountings

Because of the flexibility of lightweighted mirrors, axial support is frequently provided at many points on the back of the substrate. Hindle mounts²⁹ using multiple lever mechanisms (called “whiffletrees”) are commonly used. Figure 24 shows such a mount with 18 supports for the mirror. The number of supports needed is the minimum number that keeps the gravitational sag of the reflecting surface between support points smaller than the deflection tolerance when the mirror axis is vertical.⁴ To avoid friction effects, flexures (sometimes called “Flex-Pivots”) are typically used as single-axis bearings in these mounts. Dual-axis bearings are usually necked-down posts that serve as flexures.

A mirror on a Hindle mount also needs radial support if it is to be used in any orientation other than axis vertical. This might be in the form of three or more mechanical links with universal-joint flexures at each end that are oriented tangent to the rim of the mirror and connect the mirror rim to the surrounding structure. Provision for such a support is shown in the mirror of Fig. 23. Multiple-point whiffletree radial supports have also been used for this purpose.²¹

Mountings for Metallic Mirrors

Metallic mirrors are generally easier to support than nonmetallic ones because attachments can be made directly to the substrate through, for example, threaded holes for screws. The metallic substrate may also be stiffer than the glass counterpart. An example is the aluminum mirror shown

by section and back views in Fig. 25a.³⁰ Here, a single-point diamond-turning (SPDT) method is employed to machine the optical surface and the axial and radial mounting interface surfaces on the mirror's back. In this method, many extremely fine cuts are made with a precision diamond tool as the substrate rotates about a common mechanical and optical axis. The tool moves on a prescribed path under interferometric control. This results in very accurate surface shapes and surface interrelationships, as well as smooth surfaces and very low residual stresses in the parts. The mating surfaces on the mount are also created by diamond turning. The mirror is shown installed in its mount in Fig. 25b. Optical surface distortions due to mounting forces are minimal because the contacting surfaces on the optic and its mount are parallel when drawn together.³¹ When the mirror and its mount are made of the same material, the effects of temperature changes are minimized.

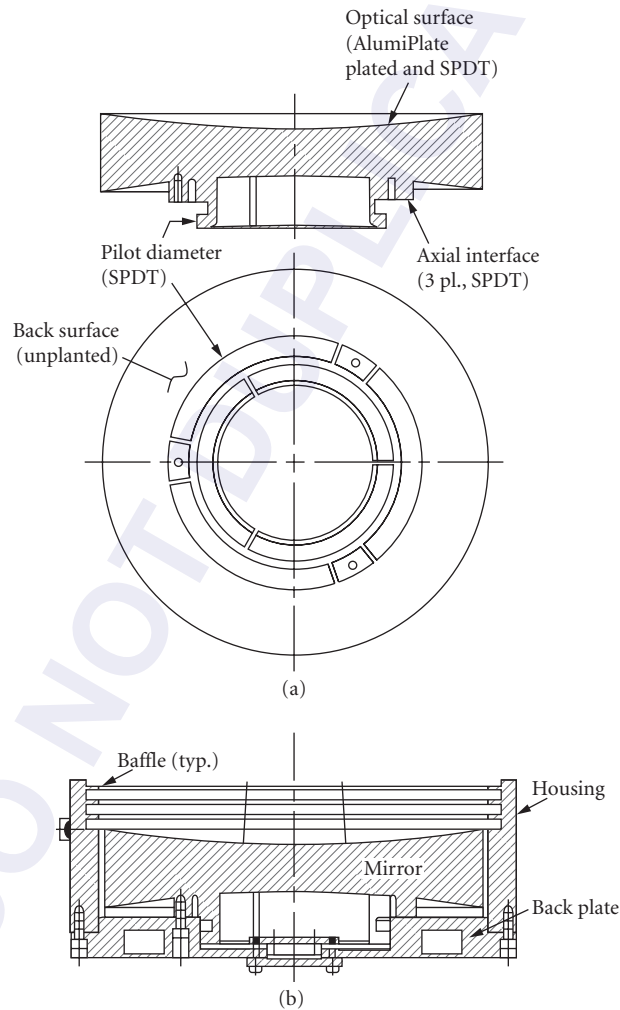


FIGURE 25 Optomechanical configuration (b) of an aluminum mirror (a) with optical and mounting surfaces machined by SPDT methods. The radial and axial interfaces are shown. (From Vukobratovich et al.³⁰)

6.8 CONTACT STRESSES IN OPTICS

The shape of the mechanical surface touching a lens, mirror, or prism surface is typically spherical, cylindrical, conical, or flat. Point contacts occur at spherical pads attached to the ends of springs while short line contacts occur if cylindrical pads are used on springs or if pins are used to locate the optic. Lenses, windows, and mirrors preloaded against mechanical constraints by a retaining ring or flange typically have circular line contacts with the metal around the edges of their apertures. The metal surface typically is conical for a convex glass surface and toroidal for a concave glass surface. The preloads applied through all of these interfaces cause elastic deformations of the glass and metal parts. Associated with these deformations are compressive and tensile stresses in those materials. Up-to-date analytical methods for estimation of these stresses have been presented in detail elsewhere.⁴ Space constraints preclude discussion of those methods here. Once the tensile stress to be expected in a given optomechanical design has been quantified, it can be compared to the aforementioned rule-of-thumb tolerance to predict success or failure of the optomechanical design. Should the stress appear to be too large, certain design changes that can be made to reduce it are suggested in the referenced publication.

6.9 TEMPERATURE EFFECTS ON MOUNTED OPTICS

General Considerations

Because the temperature environment of any optical instrument is seldom constant, we should anticipate changes in dimensions of all parts, in refractive indices, and in material parameters [such as coefficient of thermal expansion (CTE) and Young's modulus] to occur throughout the lifetime of the device. These changes may defocus the system, change aberration balance, or degrade alignment. Athermalized designs are created in a manner to reduce the magnitudes of these effects to tolerable levels.

Prevention of Axial Gaps

Differential expansions and contractions of all types of materials with temperature changes may change the axial and/or radial relationships, that is, alignment between optics and their mechanical reference surfaces. Optomechanical assemblies that are adequately preloaded at assembly will tend to maintain optic-to-mount contact, but this preload will change as the temperature changes. It may disappear completely at elevated temperatures. Then the optics may be free to move if externally disturbed, as by vibration or shock. These component shifts may become permanent if the optic is decentered or tilted when the temperature drops and the mount reapplies forces to the optics.

To reduce this effect, each optical assembly might be designed to compensate for axial dimensional changes so axial preload changes are reduced to insignificance.³² For example, the air-spaced triplet assembly of Fig. 26a is constructed of three optical glasses, an aluminum cell, and two aluminum spacers. The scale of the figure is as indicated. At maximum temperature, the physical separation of the interfacing points *A* and *B* in this particular assembly changes by 0.015 mm if computed for a path through the lenses and spacers, but changes by 0.030 mm if computed for a path through the cell. One or more axial air gaps totaling 0.015 mm would then exist somewhere within the assembly and the lenses might move or tilt within that space. If the design were to be modified by changing the metals in the cell and in one spacer, lengthening that spacer, and providing space for the larger spacer by adding a step bevel to the second lens—as indicated in Fig. 26b—the chosen materials and component dimensions would make the *A*-to-*B* separation remain equal for both paths for all temperature changes. Preload would then remain unchanged and misalignment would not occur.

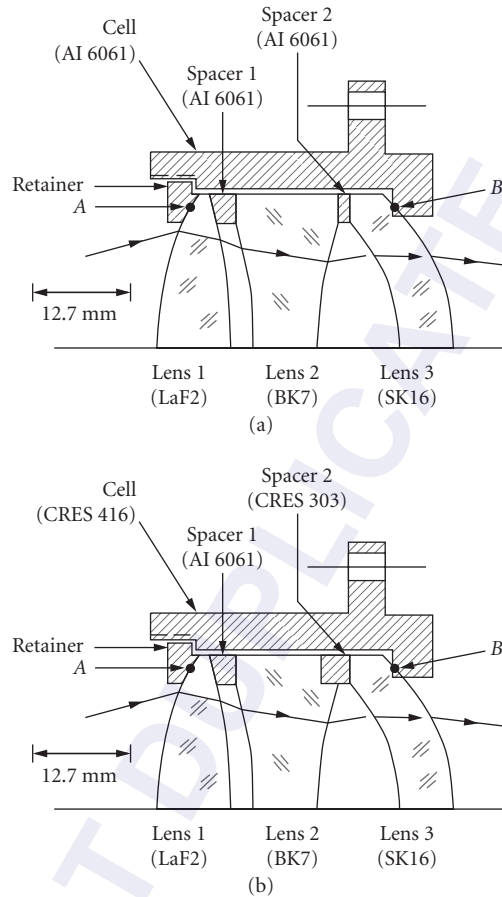


FIGURE 26 An air-spaced triplet lens assembly (a) in which an axial gap between glass and metal parts exists at high temperature, possibly allowing the lenses to become misaligned. Modified design (b) is athermalized to maintain registry of the optics in the mount. (Adapted from Yoder and Hatheway.³²)

Focus Athermalization Techniques

Single Material Designs Figure 27 shows a reflecting telescope made of a single material, in this case, aluminum.³³ All dimensions change, but the assembly remains in alignment and the optical performance is unchanged (other than a small change in image scale) as the temperature changes. This telescope is an example of the use of single-point diamond-machining methods as all optical and mounting surfaces are precisely made in the proper geometric relationships so alignment accuracy is built-in.

Passive Athermalization Figure 28a illustrates the use of materials with dissimilar CTEs and carefully chosen axial dimensions so the axial distance between optical components (in this case, the two mirrors) remains constant when the temperature changes.³⁴ This keeps the optical performance within required limits. Control of the mirror separation of this telescope is modeled schematically in Fig. 28b. Positive signs associated with lengths of low and high CTE materials indicate how the

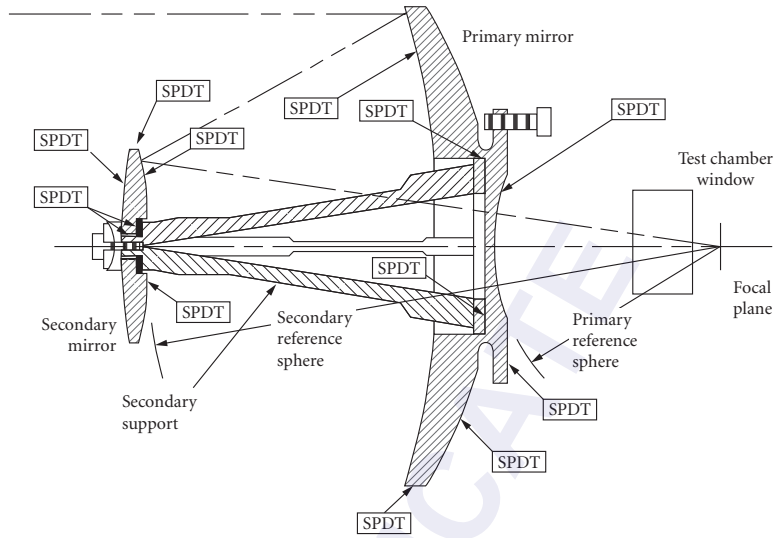
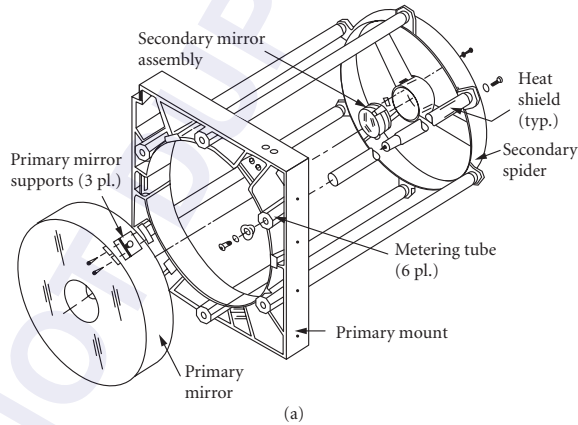
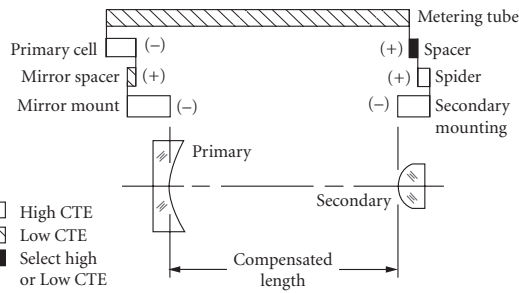


FIGURE 27 Schematic of an all-aluminum (athermalized) telescope objective with optical and mechanical interface surfaces finished by SPDT for ease of assembly without alignment. (From Erickson et al.³³)



(a)

(+)



(b)

FIGURE 28 A passively athermalized telescope structure using Invar metering tubes to connect the primary and secondary mirror mounts. (a) Exploded view of the telescope. (b) Model of the compensation system. (From Zurmehly and Hookman³⁴)

mirror separation changes as the temperature rises. Proper choices of materials for their coefficients of thermal expansion and dimensions make the mirror separation remain constant as the temperature changes.

Active Athermalization When a source of power is available, components in an optical system can be physically moved to compensate for the effects of temperature changes. For example, Fig. 29a shows a concept for a zoom lens system in which locations of the moveable lenses are varied by motors as commanded by an internal microprocessor that monitors the temperature of the system.³⁵ As indicated in Fig. 29b, desired magnification inputs from the operator are automatically converted into the lens shifts required to focus properly on the object at the measured temperature.

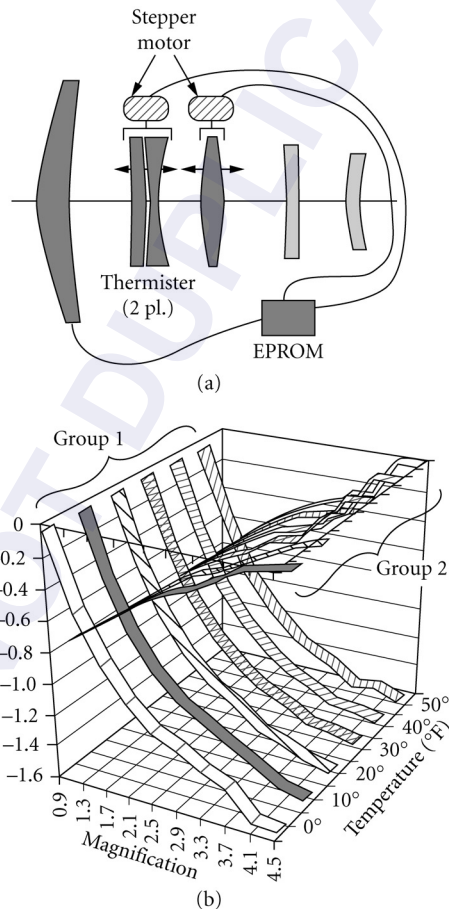


FIGURE 29 An actively athermalized zoom lens system that drives two lens groups to maintain focus at selectable magnification settings in spite of temperature changes. (From Fischer and Kampe.³⁵)

6.10 REFERENCES

1. V. L. Genberg, "Structural Analysis of Optics," Chapter 8 in *Handbook of Optomechanical Engineering*, A. Ahmad, (ed.), CRC Press, Boca Raton, 1997.
2. K. B. Doyle, V. L. Genberg, and G. J. Michels, *Integrated Optomechanical Analysis*, SPIE Press, Bellingham, 2002.
3. P. R. Yoder, Jr., *Opto-Mechanical Systems Design*, 3rd ed., CRC Press, Boca Raton, 2005.
4. P. R. Yoder, Jr., *Mounting Optics in Optical Instruments*, 2nd ed., SPIE Press, Bellingham, 2008.
5. R. J. Roark, *Formulas for Stress and Strain*, 3rd ed., McGraw-Hill, New York, 1954.
6. S. P. Timoshenko and J. N. Goodier, *Theory of Elasticity*, 3rd ed., McGraw-Hill, New York, 1970.
7. K. B. Doyle and M. Kahan, "Design Strength of Optical Glass," *Proc. SPIE*, 5176, 14, 2003.
8. R. Stoll, P. F. Forman, and J. Edelman, "The Effect of Different Grinding Procedures on the Strength of Scratched Fused Silica," *Proc. of the Symposium on the Strength of Glass and Ways to Improve It*, Union Scientifique du Verre, Florence, 1961.
9. K. B. Doyle, private communication, 2008.
10. J. J. Herbert, "Techniques for Deriving Optimal Bondlines for Athermal Bonded Mounts," *Proc. SPIE* **6288OJ-1**, 2006.
11. A. Ahmad and R. L. Huse, "Mounting for High Resolution Projection Lenses," *U.S. Patent 4,929,054*, 1990.
12. J. J. Bacich, "Precision Lens Mounting," *U.S. Patent 4,733,945*, 1988.
13. J. H. Bruning, F. A. DeWitt, and K. E. Hanford, "Decoupled Mount for Optical Element and Stacked Annuli Assembly," *U.S. Patent 5,428,482*, 1995.
14. E. T. Kvamme, D. Trevias, R. Simonson, and L. Sokolsky, "A Low Stress Cryogenic Mount for Space-Borne Lithium Fluoride Optics," *Proc. of SPIE* **58770T**, 2005.
15. E. T. Kvamme and Michael Jacoby, "A Second Generation Low Stress Cryogenic Mount for Space-Borne Lithium Fluoride Optics," *Proc. SPIE* **66920I**, 2007.
16. E. T. Kvamme and M. Jacoby, "Opto-Mechanical Testing Results for the Near Infra-red Camera on the James Webb Space Telescope," *Proc. SPIE* **7010**, 2008.
17. P. R. Yoder, Jr., "Lens Mounting Techniques," *Proc. SPIE* **389**: 2, 1983.
18. M. Bayar, "Lens Barrel Optomechanical Design Principles," *Opt. Eng.* **20**:181, 1981.
19. R. E. Fischer, "Case Study of Elastomeric Lens Mounts," *Proc. SPIE* **1533**: 27, 1991.
20. D. M. Williamson, "Compensator Selection in the Tolerancing of a Microlithography Lens," *Proc. SPIE* **1049**: 178, 1989.
21. D. Vukobratovich, *Introduction to Opto-Mechanical Design*, SPIE Short Course SC014, 2003.
22. W. Sunne, "Dome Attachment with Brazing for Increased Aperture and Strength," *Proc. SPIE* **5078**: 121, 2003.
23. P. R. Yoder, Jr., "Non-Image-Forming Optical Components," *Proc. SPIE* **531**: 206, 1985.
24. P. Mammini, B. Holmes, A. Nordt, and D. Stubbs, "Sensitivity Evaluation of Mounting Optics Using Elastomer and Bipod Flexures," *Proc. SPIE* **5176**: 26, 2003.
25. P. R. Yoder, Jr., "Mounting-Induced Contact Stresses in Prisms," *Proc. SPIE* **3429**: 71, 1998.
26. P. R. Yoder, Jr., "Improved Semikinematic Mounting for Prisms," *Proc. SPIE* **4771**:173, 2002.
27. P. R. Yoder, Jr., "Design Guidelines for Bonding Prisms to Mounts," *Proc. SPIE* **1013**: 112, 1988.
28. L. H. J. F. Beckmann, private communication, 1990.
29. J. H. Hindle, "Mechanical Floatation of Mirrors," *Amateur Telescope Making, Advanced*, A.G., Ingalls, ed., Scientific American, New York, 1945: 229. (Reprinted 1996 as Chapter B.8 in *Amateur Telescope Making, 2*, William-Bell, Inc., Richmond.)
30. D. Vukobratovich, A. Gerzoff, and M. K. Cho, "Therm-Optic Analysis of Bi-Metallic Mirrors," *Proc. SPIE* **3132**: 12, 1997.
31. R. L. Rhorer and C. J. Evans, "Fabrication of Optics by Diamond Turning," Chapter 41 in *Optical Society of America Handbook of Optics*, 2nd ed., Vol. I, Bass, M., Van Stryland, E. W., and Wolfe, W. L., eds., McGraw-Hill, New York, 1995.

32. P. R. Yoder, Jr. and A. E. Hatheway, "Further Considerations of Axial Preload Variations with Temperature and the Resultant Effects on Contact Stresses in Simple Lens Mountings," *Proc. SPIE* **587705**, 2005.
33. D. J. Erickson, R. A. Johnston, and A. B. Hull, "Optimization of the Opto-Mechanical Interface Employing Diamond Machining in a Concurrent Engineering Environment," *Proc. SPIE* **CR43**: 329, 1992.
34. G. E. Zurmehly and R. Hookman, "Thermal/Optical Test Setup for the Geostationary Operational Environmental Satellite Telescope," *Proc. SPIE* **1167**: 360, 1989.
35. R. E. Fischer and T. U. Kampe, "Actively Controlled 5:1 Afocal Zoom Attachment for Common Module FLIR," *Proc. SPIE* **1690**: 137, 1992.

DO NOT DUPLICATE

CONTROL OF STRAY LIGHT

Robert P. Breault

*Breault Research Organization
Tucson, Arizona*

7.1 GLOSSARY

A	area
BRDF	bidirectional reflectance distribution function
GCF	geometric configuration factor
L	radiance
R	distance
θ, ϕ	angles
Φ	power
Ω	solid angle

7.2 INTRODUCTION

The analysis of stray light suppression is the study of all unwanted sources that reduce contrast or image quality. The control of stray light encompasses several very specialized fields of both experimental and theoretical research. Its basic input must consider (1) the optical design of the system; (2) the mechanical design, size, and shape of the objects in the system; (3) the thermal emittance characteristics for some systems; and (4) the scattering and reflectance characteristics of each surface for all input and output angles. It may also include spectral characteristics, spatial distribution, and polarization. Each of these areas may be concentrated on individually, but ultimately the analysis culminates in the merging of the various inputs.

Developments in detector technology, optical design software, diffraction-limited optical designs, fabrication techniques, and metrology testing have created a demand for sensors with lower levels of stray radiation. Ways to control stray light to meet these demands must be considered during the “preliminary” conceptual design. Decisions made at this time are, more often than not, irrevocable. This is because parallel studies based upon the initially accepted starting design are often very expensive. The task of minimizing the stray radiation that reaches the detector after the system has been

designed by “adding on” a suppression system is very difficult. Therefore, every effort should be made to start off with a sound stray light design. To ensure a sound design, some stray light analysis should be incorporated in the earliest stages of a preliminary design study.

This chapter presents some basic concepts, tools, and methods that you, the optical or mechanical designer, can consider when creating a sensor system. You do not need to be very experienced in stray light suppression to design basic features into the system, or to consider alternative designs that may significantly enhance the sensor’s performance. The concepts are applicable to all sizes of optical instrumentation and to virtually all wavelengths. In some cases, you can use the concepts to rescue a design when experimental test results indicate a major design flaw.

7.3 CONCEPTS

This section outlines some concepts that you can use to reduce stray radiation in any optical system. The section also contains some experimental and computer-calculated data as examples that should give you some idea of the magnitude of the enhancement that is possible.

The power on a collector depends on the following:

1. The power from the stray light source.
2. The surface scatter characteristics of the source; these characteristics are defined by the bidirectional scatter distribution function (BSDF).
3. The geometrical relationship between the source and collector. This relationship is called the geometrical configuration factor (discussed later in this section).

To reduce the power on the detector, we can try to reduce the contributions from these elements:

$$\Phi_{\text{collector power}} = \Phi_{\text{source power}} \times \text{BRDF}_{\text{source}} \times \text{GCF}_{\text{source collector}} \times \pi \quad (1)$$

Ways to reduce each of these factors are discussed below. The creative use of aperture stops, Lyot stops, and field stops is an important part of any attempt to reduce the GCF term of the power transfer equation.

For the discussion that follows, examples from a two-mirror Cassegrain design, with the aperture stop at the primary, will be used to stimulate thoughts about stray light reduction possibilities for other sensors.¹ The system is shown in Fig. 1.

Critical Objects

The most fundamental concept is to start the stray light analysis from the detector plane of the proposed designs. The most critical surfaces in a system are those that can be seen from the detector position or focal surface. These structures are the only ones that contribute power to the detector. For this reason, direct your initial attention toward minimizing their power contributions by removing them from the field of view of the detector.

The basic idea is to visualize what would be seen if you were to look out of the system from the image plane. Unlike most users of optical instruments, the stray light designer’s primary concern is seldom the object field, but rather all the interior surfaces that scatter light. It is necessary to look beyond the radii of the imaging apertures to find the sources of unwanted energy. Removing these sources from the field of the detector is a real possibility, and will result in a significant improvement in the system.

Real-Space Critical Objects

I will start out by identifying a particular critical real-space object that can be seen by the detector in our example; it is the inside of the secondary baffle. The direct view discussed here is different than the image of the same baffle reflected by the secondary which is discussed in the next section.

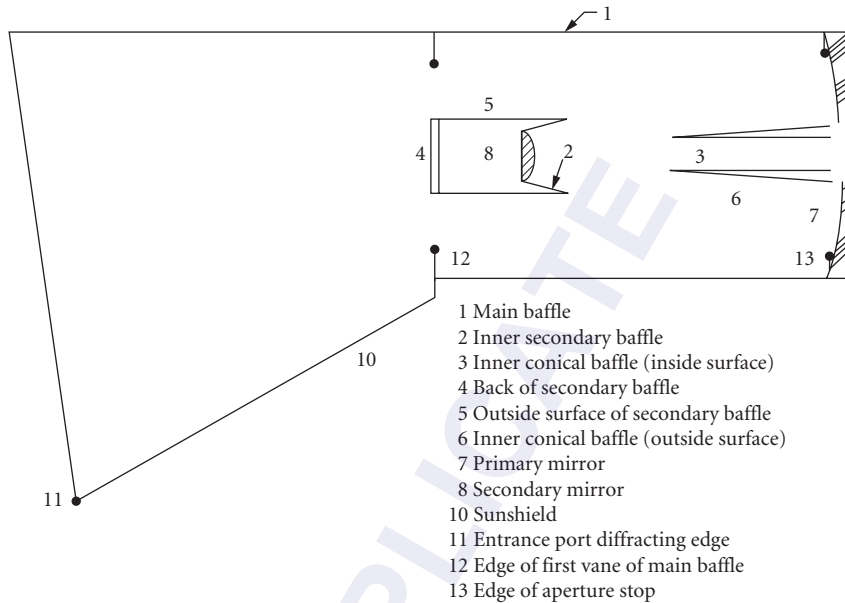


FIGURE 1 Typical Cassegrain design with the aperture stop at the primary. (Ref. 2, p. 52.)

Many Cassegrain secondary baffles have been designed to be cone-shaped (Fig. 2), usually approximating the converging cone of light from the primary. From the detector, portions of this secondary cone are seen directly as a critical surface. Since most of the unwanted energy is incident on this baffle from nearly the same direction as this surface is seen from the detector, the addition of vane structures would be of little help, assuming an optimum coating is used on the simple baffle. If the cone is made *more* cylindrical, the amount of critical cone area is reduced, and the angle at which the surface is seen gives a smaller projected area (Fig. 3).

Avoid making the baffle cylindrical because the outside of it would be seen. Since the detector is of a finite size there is a fan of rays off the primary representing the field of view of the telescope.

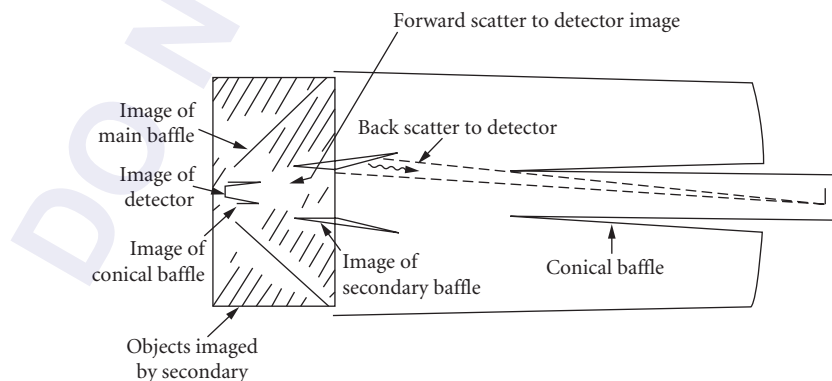


FIGURE 2 Direct and reflected scatter from the cone-shaped secondary baffle. (Ref. 1, p. 4.)

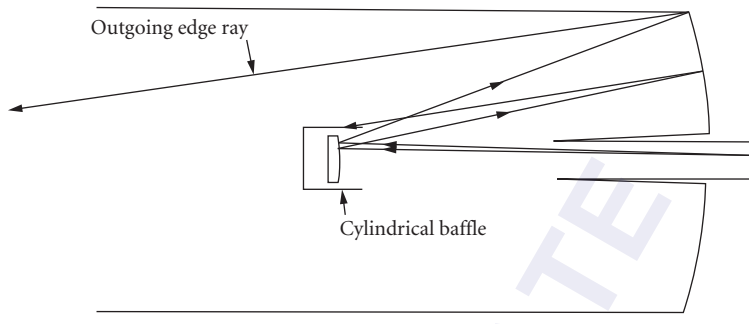


FIGURE 3 Reduced scatter from an almost cylindrical-shaped secondary baffle.

Although collimated for any point on the detector, any point not on axis would have its ray bundle at some angle to the optical axis, hence a cylindrical secondary baffle would be seen from off-axis positions on the detector.

Imaged Critical Objects

Imaged objects are often critical objects. They too can be seen from the detector. Determining which of the imaged objects are critical requires a bit of imagination and usually some calculations; stray light software can help you make the calculations. The Y-Y bar diagram can help you to conceptualize and determine the relative image distances and sizes with a minimum number of calculations.³ The same could be done with other first-order imaging techniques (see Chap. 1 in this volume). Using the Cassegrain example again (Fig. 2), you can see that reflected off the secondary mirror are the images of the detector and the inside of the inner conical baffle (object 3 in Fig. 1). In some designs the outside of the conical baffle will be seen in reflection. These are *imaged critical objects*. If you wish, you can eliminate some of these images with a central obscuration on the secondary, or for the conical baffle, with a spherical mirror concentric about the image plane. The direct view to the inner conical baffle will remain, but the path from the image of it is removed.

The cone-shaped secondary baffle is also seen in reflection (Fig. 2). For the incident angles of radiation on this surface, the near specular (forward-scattering) characteristics will often be one of the most important stray light paths because the image of the detector is in that direction. This is an extension of *starting from the detector*. There is an image of the detector at the prime focus of the primary mirror. Often, as in this case, one location may be easier than another for you to determine what could be seen. By making the baffle more cylindrical, part of the *image* of the baffle is removed from the detector's view; as a result, the power that can scatter to the detector is reduced. Furthermore, it is sometimes possible to baffle most of this power from the field of view with one or two vanes (Fig. 4).

Continue the process of removing critical surfaces until all the critical surfaces have been considered for all points in the image plane. The power contributions from these surfaces will either go to zero, or at least be lessened after you reduce the area of the sections seen.

There is still more that can be done, since only the GCF term in the power transfer equation to the detector has been reduced. It is also possible to minimize the power onto the critical sections, which will become the source of power, Φ , at the next level of scatter. This approach can be very similar to the approach used for minimizing the power scattered to the image plane. The viewing is now forward from the critical surfaces instead of the image plane. By minimizing the BRDF and GCF factors of the surface scattering to the critical sections, the power incident on the critical surfaces will be reduced. Hence, the power to the detector is also reduced.

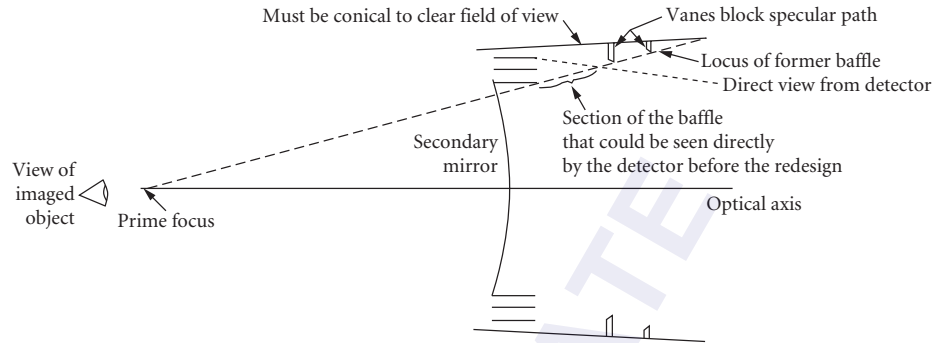


FIGURE 4 A cylindrical secondary baffle can be seen from off-axis positions on the detector.

Illuminated Objects

Minimizing the GCFs and BRDFs for the specific input and output angle is sometimes easier if you look into the system from the position of the stray light source in object space. By doing this, you can identify the surfaces that directly receive the unwanted energy. I will call these the *illuminated objects*. If any of these illuminated surfaces contain sections that the detector can see, then you should direct your initial efforts toward eliminating these paths. These paths will usually dominate all other stray light paths because there is only a single scatter before the stray light reaches the detector. An example of such a path that is often encountered is from the source onto the inner conical baffle of multimirror systems (Fig. 1). Some of the ways that the direct radiation can be eliminated is by extending the main baffle tube, increasing the obscuration ratio by increasing the diameter of the secondary baffle (Fig. 5), or by narrowing the field of view, which will allow you to extend the secondary baffle and the inner conical baffle toward each other.

The effect of eliminating this path is shown in a composite Point Source Transmittance (PST) plot in Fig. 6.⁴ The PST plot is defined as the reference plane (detector plane in most cases) irradiance divided by the input irradiance along the line of sight. (See the section called “Point Source Transmittance Definitions” for a more detailed definition of PST.) For the case shown, the unwanted irradiance on the detector is reduced by over an order of magnitude.

Aperture Placement

I will now focus on the optical design aspects of a stray light suppression system, and give a qualitative discussion of some general aspects that you might consider. All optical systems will have at least

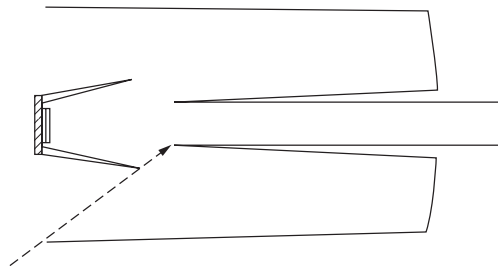


FIGURE 5 Increased obscuration ratio blocks direct path to inner conical baffle. (Ref. 1, p. 6.)

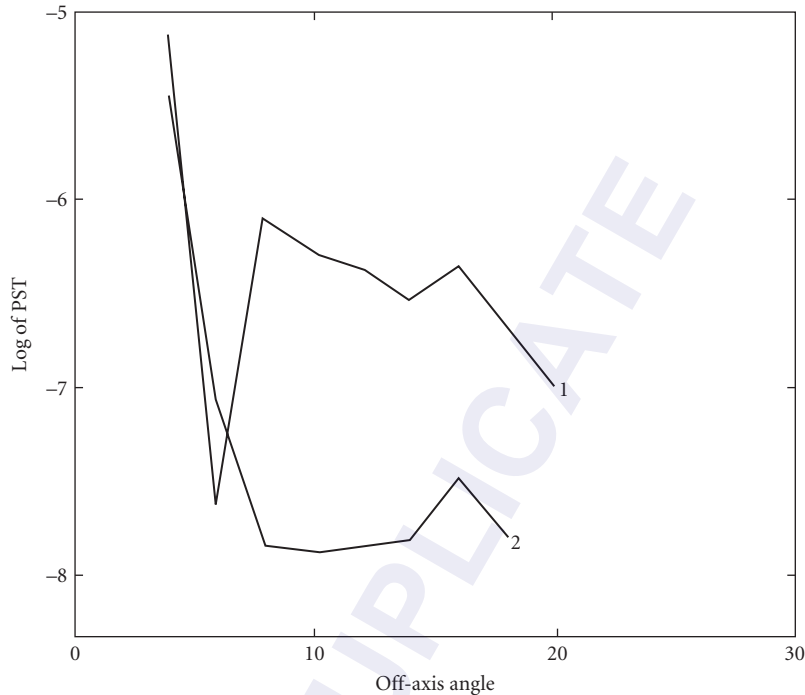


FIGURE 6 Point source transmittance with obscurations of (1) 0.333 and (2) 0.4. The 0.4 obscuration removed the direct path from the source to the inner conical baffle. (Ref. 1, p. 6.)

one aperture, called an *aperture stop*, that limits the size of the bundle of the incoming signal rays. Some systems will have field stops and/or Lyot stops. Each type of stop has a clearly defined role in stray radiation suppression, which is discussed in the following sections.

In many cases stop placement will have a much more noticeable effect on system performance than any vane structure, coating, or baffle redesign. Probably the only factor with more effect on the PST curve is the off-axis position of the source. Therefore, the benefits of any of the stops cannot be overemphasized.

Aperture Stops The aperture stop is the aperture that limits the size of the cone of radiation that will reach a point on the image plane. Sometimes shifting this stop allows the optical designer to better balance the aberrations. In a stray radiation suppression design, it plays a similar important role. All objects in the spaces preceding the stop in the optical path will not be seen unless they are imaging elements, central obscurations, or objects that vignette the field of view. Only a limited number of critical objects is possible before the aperture stop. In the intervening spaces from the stop to the image plane it is likely that many of the baffle surfaces will be seen. Figure 7 represents a two-mirror design, and Fig. 8 represents a three-element refracting system; both have the stop at the first element. In both cases the second element is oversized to accommodate the field of view from a point in the field stop; the amount depends on the full field of view of the design. Because the elements are oversized, the main baffle following the first element will be seen. This baffle will be a critical object, a direct path of unwanted energy. The “overviewing” is characteristic of all of the optical elements past the aperture stop.

If you move the stop along the optical path toward the detector plane, its performance as a stray radiation baffle will improve. If you shift the stop to the second element, the intermediate baffle will

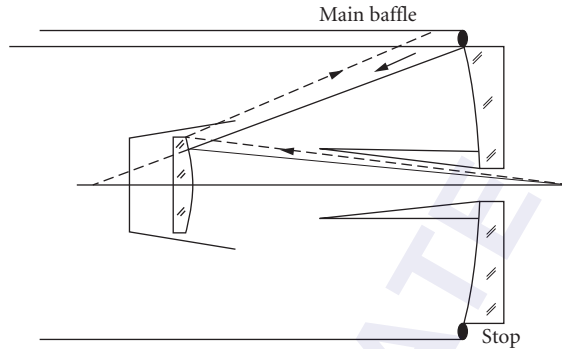


FIGURE 7 The oversized secondary allows the main baffle to be seen in reflection. (Ref. 1, p. 8.)

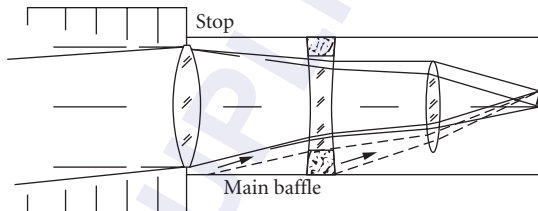


FIGURE 8 The main baffle is seen through the oversized secondary and tertiary. (Ref. 1, p. 8.)

not be seen. It is removed from the field of view of the detector, since the stop now eliminates direct paths from baffles in all spaces that precede it. Figure 9 shows the improvement in the PST curve for a two-mirror system. By moving the stop you have reduced the PST by a factor of 10. This is a desirable feature to consider for stray radiation reduction.

Direct paths from central obscurations can be blocked by a central disk located at some location deeper into the system; however, because of the parallax involved between the central obscuration disk and this central disk, the central disk obscuration will usually be a larger obstruction to imaging rays. In a reimaging design it is often possible to locate a central disk conjugate to the actual central obscuration.

Field Stops An aperture can be placed at intermediate *images* in a system to limit the field of view. Such an aperture will usually prevent any stray light from outside of the field of view from being directly imaged into the system beyond this field stop aperture. In a sense, its operation is just opposite that of an aperture stop. Baffle surfaces following a field stop cannot be seen from outside the field of view in the object plane, unless they are central obscurations. Note that with just a field stop, succeeding optical elements may allow out-of-field critical sections to be seen *through* the field stop, from within the field in the image plane (Fig. 11). Aperture stops are necessary to block such paths. Figures 10 and 11 show two such cases. Although for some designs the field stop is not 100 percent effective because of optical aberrations, its small size limits most of the unwanted stray light. Field stops therefore do not remove critical sections, but rather limit the propagation of power to illuminated objects. In reflecting systems, take care that the object side of the field stops does not become a critical area, which can be seen directly or in reflection from the image plane because unwanted energy is being focused onto them.⁵

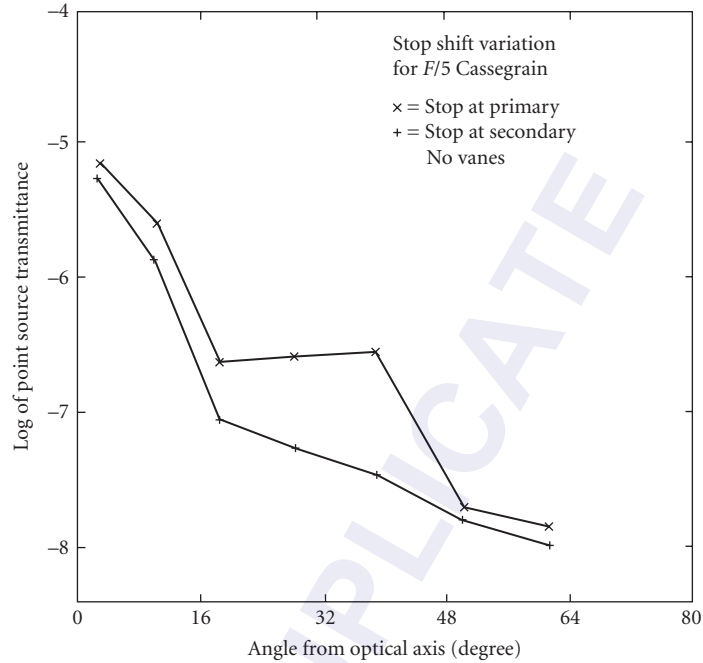


FIGURE 9 PST improvement with stop shift for the two-mirror system. (Ref. 1, p. 8.)

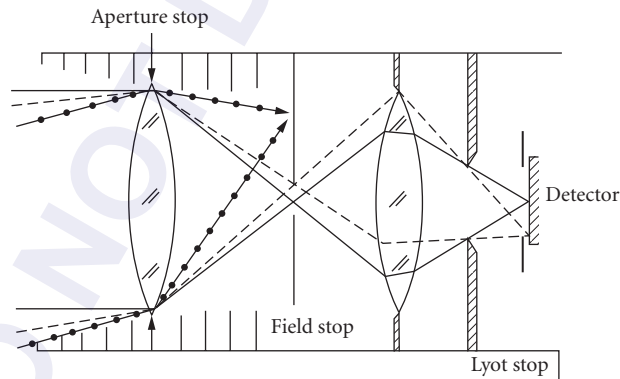


FIGURE 10 The addition of a Lyot stop prevents preceding baffles from being seen from the image.

Lyot Stops A limiting aperture placed at the location of the image of an aperture stop, sometimes called a glare stop or Lyot stop, has the same property as described for aperture stops. It should be slightly smaller than the image of the aperture stop. It limits the critical sections which are out of the field of view to those objects in succeeding spaces only. Since Lyot stops are by definition further along the optical path to the detector, the number of critical surfaces seen by the detector will be reduced. Usually, these stops are incorporated into the design to block the diffracted energy from an aperture stop and field stop pair, so that only secondary or tertiary diffracted energy reaches the image. Nevertheless, both diffracted and scattered energy are removed from the direct

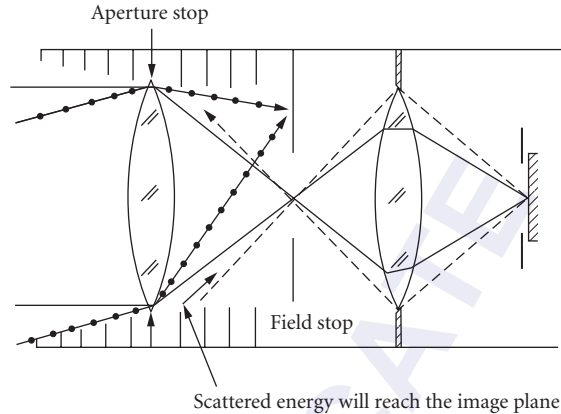


FIGURE 11 Out-of-field energy in object plane will not be imaged beyond field stop. Out-of-field energy elsewhere may be seen. (Ref. 1, p. 9.)

view of the image, re imaging the largest optical element as the stop takes full advantage of both the light-gathering power of the optics and the stray radiation suppression features provided by the stop. Figure 10 shows a system with a Lyot stop.

On space-based telescopes the image plane is often *shared* by one or more instruments. Each instrument reimages the telescope's image through some optical train, and eventually onto the detector. In that optical train there could be a logical place to use a Lyot stop to improve the stray light performance of the viewing instrument well beyond that of the telescope.

It is the combination of these different stops or apertures that helps minimize the propagation of unwanted energy by limiting the number of critical objects seen by the detector, and the objects illuminated by the stray light source.

When all direct paths have been eliminated, the next step is to determine the relationship between the sections that received power (illuminated objects) and the critical surfaces. This relationship takes the form of scattering *paths*; that is, stray light can scatter from the illuminated objects to the critical objects. To start reducing the stray light contributions from these paths, you can start at the critical surfaces as described above. But now you have more knowledge about where the direct incident power is being distributed throughout the system, since you can also look into the system from the source side to find the surfaces receiving direct power. With this information you can identify the possible paths between the illuminated and critical surfaces.

Design considerations for extended baffle shields (Figs. 12 and 13) provide a good example of starting from the source side to identify possible paths. In the examples, object 2 (an optical surface) is the largest contributor of scattered radiation and is the best superpolished mirror available

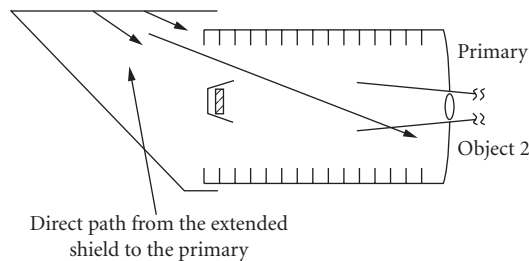


FIGURE 12 There is a direct path from the baffle to the primary mirror. (Ref. 1, p. 7.)

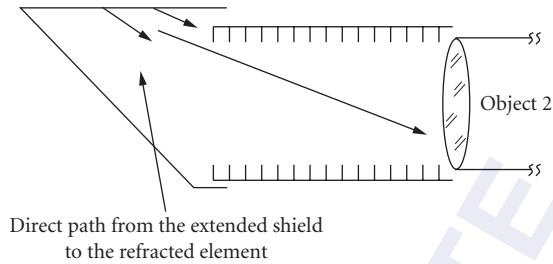


FIGURE 13 There is a direct path from the baffle to the refracting element. (Ref. 1, p. 7.)

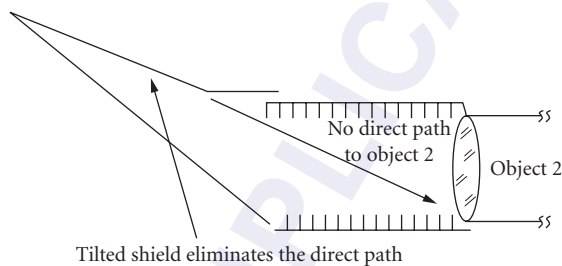


FIGURE 14 Direct paths are removed by properly tilting the shield. (Ref. 1, p. 7.)

(Fig. 12), or if it is a lens as in Fig. 13, it has the lowest possible scattering characteristics. It cannot be removed from the view of the image plane. If the initial power incident on object 2 is only from the extended shield, then by tilting the shield (Fig. 14), the power on the shield must first scatter to the main tube and then to the optical element. The combination is then referred to as a two-stage baffle. If vanes are added to the main baffle, the scattered radiation incident on the optical element will be reduced by many orders of magnitude when all the scattering solid angles and the number of absorbing surfaces are considered. Note that without the tilt to the shield, vanes on the main baffle are worthless because there is a direct path to the objective.

Figure 15 is an abstract representation of the process of reducing stray light in a sensor system. Start at the detector, then work from its conjugate image locations. Starting from the detector simplifies the analysis and directs your attention to the most productive solutions, because you can identify all the possible sources of stray light to the detector. You can then work at decreasing their number by slightly redesigning the baffles and stops. Next, identify which objects are illuminated. Discover how energy may propagate between them and you have identified the paths of stray light propagation. From then on the process of moving objects or blocking paths is quite simple, although the quantitative calculations might get difficult and may require some analysis software.

Baffles and Vanes

A few definitions are required to define baffles and vanes. Other authors have used their own different definitions. In this section the term *baffle* is used to describe conical structures (including cylindrical) that can also be described as tubelike structures. Their function is to shade, or occult, stray light from the source to one or more system components. The main baffle shields the primary

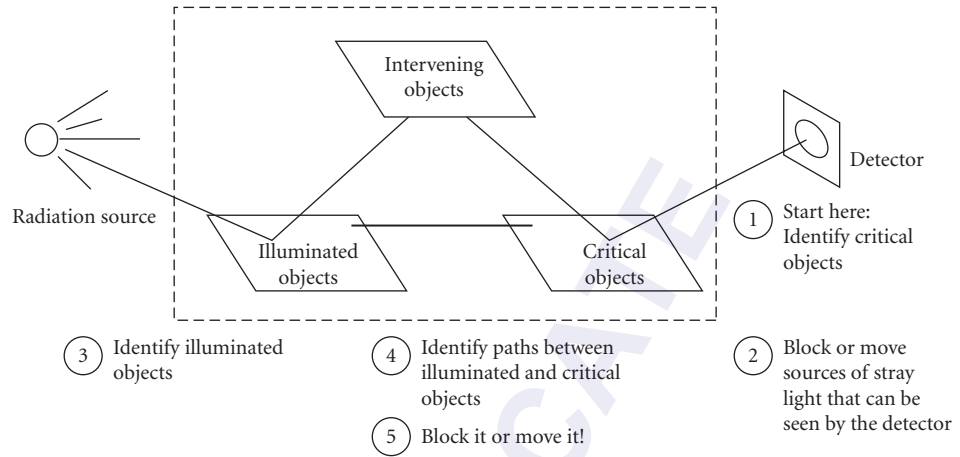


FIGURE 15 The first step in a stray light analysis begins from the detector plane, not from the source.

mirror from direct radiation at the larger off-axis angles. *Vanes* are structures put *on* baffles to affect the scatter characteristics of the surface. Other authors have used the term “baffles,” or “glare stops,” to describe these vanes.

Baffles In a well-designed system vanes play an important role only at large off-axis angles. For example, when one-tenth of the stray light falls on the primary of the Cassegrain design, then the main baffle receives the remaining 90 percent of the stray light. When the main baffle has properly designed vanes on it, light that falls on the baffles is attenuated by five orders of magnitude before it reaches the primary mirror. The resultant power on the primary mirror is then about 9.0×10^{-5} compared to the direct 10 percent that fell on the mirror. This results in less than 0.1 percent of the total on the primary. In addition, most of the subsequent scatter off the primary will be at much higher scatter angles. This will cause the scattered energy to have much lower BRDFs off the primary mirror when scattered in the direction of the detector, further reducing the scatter contribution from the baffle.

Only when no power illuminates the objective will the baffles play a significant role in the propagation paths of the stray light. Usually the system’s performance merit function is then very good. Only if the stray light source has a tremendous amount of energy, like the sun, does the stray light become measurable.

Vanes The depth, separation, angle, and bevel of vanes are variables that need to be evaluated for every design. In the following paragraphs stray light analysis results are presented for both a centrally obscured system (Cassegrain) and an unobscured eccentric pupil design (Z-system).⁶ Profiles of these systems are given in Figs. 2 and 16. Of the two designs, only the eccentric pupil design has a reimager that would allow for the placement of an intermediate field stop and an accessible Lyot stop, as discussed above.

The APART stray light analysis program was used to analyze the two designs. The APART program was a substantial software package that performed deterministic calculations of stray light propagation in optical systems.^{2,7,8}

As an example of vane design considerations, the design of vanes on a main baffle tube will be explained. With minor differences, the design steps are the same for the Cassegrain and the eccentric pupil designs. In a reimaging system, vane structure deep in the system is usually not necessary, but there are exceptions. Figure 17 shows a collecting optical element that has some small field of view (FOV). The optical element could represent a primary mirror or a refractive element.

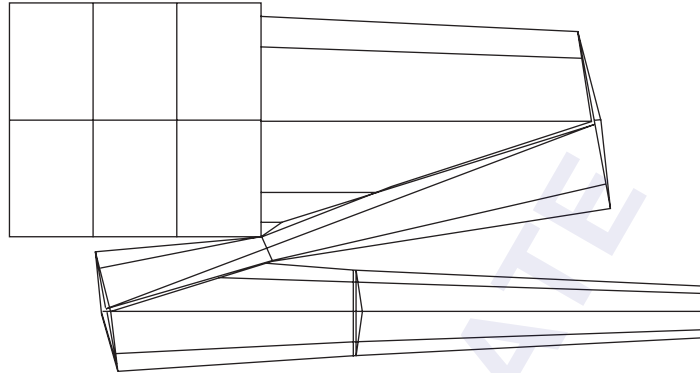


FIGURE 16 Confocal mirror system, eccentric pupil, no obscuration, low-scatter system. (Ref. 6, p. 91.)

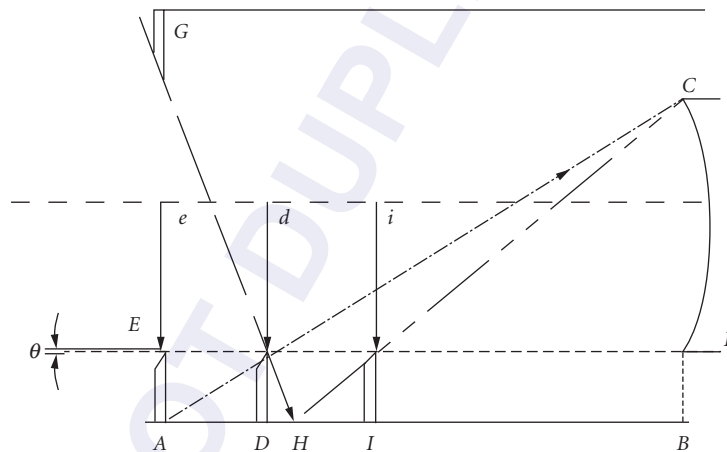


FIGURE 17 Vane placement design, lowercase letters are radii (measured from the optical axis), uppercase are z locations. (Ref. 6, p. 94.)

The placement of a straight, diffusely coated cylindrical tube would block the direct radiation from an external source, such as the sun, from reaching the optical element for a certain range of off-axis angles. If it were at a large off-axis angle, the forward scatter off the inside of the tube would be so high that it would normally not be acceptable. The solution is to add vanes to block this path.

Figures 18 and 19 depict the two cases that could represent the scatter from a baffle. In one case there are no vanes; in the other case there are vanes. This example shows how a propagation path is blocked by vanes. Vanes are useful, but a better approach is to make the solid angle (Ω_c) from the baffle (not the vanes) to the collectors of the scattered light go to zero, so that there is no path from the baffle and vane structure to the collecting object. By moving the baffle out of the field of view of the collector, the baffle's contribution goes to zero. There is no edge scatter, and no edge diffraction effects. That topic is in the realm of baffle design, which has already been touched on, and is well covered in the literature.^{9,10}

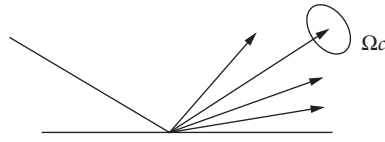


FIGURE 18 High forward scatter path. (Ref. 6, p. 92.)

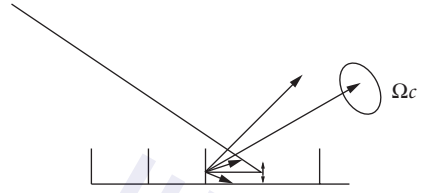


FIGURE 19 Forward scatter path highly attenuated by the vane structure. (Ref. 6, p. 92.)

Designing the Vane Spacing and Depth¹¹

A first vane is most often placed at the entrance of the baffle and an external ray is brought in from object space at a maximum off-axis position. If there is no forebaffle the angle is 90° off axis. The depth of the vane cavity is normally dictated by space and weight requirements. Too little depth will dictate the requirement for many vanes. Then vane edge scatter eventually becomes the major source of scatter instead of multiple vane scatter.

The initial ray will strike the side wall at the base of the first vane (point *A* in Fig. 17). From this point, a design line is drawn/calculated (*AC*) from the wall to the edge of the optical element on the opposite side. This line (*AC*) intersects the edge ray (*EF*), at *z* position *D*. At this point a vane could be placed. Mathematically, this assures that any point below *C*, including those on the optical element, would not see any directly illuminated side wall. However, practicality dictates that some offset of point *D* to a point *D'* (not shown) is required to allow for tolerance errors in fabrication of the vane, thermal effects, assembly errors, and for stray light edge scatter and diffraction effects. The tolerance allowance is company-, material-, and design-dependent. Acceptable numbers are often about 0.125 mm for fabrication and assembly tolerances. For the rest of this analysis, assume that this is accounted for.

Continue the design process by constructing another line from the edge of the entrance aperture to the tip of the second vane to the wall (line *GH*). Draw a new *HC* line to the area near the objective and determine the placement of the third vane (at *I*); once tolerances are considered, iterate the process to reach a final design. In some cases you may have to consider more than just the scatter path to the objective. In the Cassegrain design you may also have to consider the inner conical baffle opening. It is beyond the present scope to go into further detail.¹²

Bevel Placement on Vanes In this short discussion on baffle-vane design and placement, I did not mention the placement of a slanted surface, or bevel, sometimes placed on a vane edge as shown in Fig. 20. Which side should the bevel go on? The answer is usually dictated by first-order scatter principles.^{13,14} Near the front of the tube, direct radiation from a source at large off-axis angles will strike this bevel. If it is placed on the right side (Fig. 21), then the illuminated bevel will scatter its radiation all the way down into the tube to some optical surface. If placed on the left side, as depicted in Fig. 20, then it will go only 16° deeper into the system to the opposing vanes, a much better solution. For vanes deeper into the system, the bevel is placed on the right side. The point at which this is done is determined by the angle of the bevel and the diameter of the baffle tube. At some point, external radiation will not be able to directly strike the beveled edge if it is on the right-hand side of the vane. Only the nonbeveled, straight side will be illuminated. Therefore, the vane can rescatter only in the left side of the hemisphere, which is in the direction of going out of the system. If the bevel is placed on the left side, it can scatter 16° (in the example) deeper into the sensor; this is usually a needless design shortcoming that could be a significant error.

Vane Angle Considerations Another variation on the design feature of vanes that has sometimes been incorporated onto baffles in an optical system is angled vane structure. These vanes are non-planar objects. This makes them quite tedious to cut out of sheet metal, fabricate, and install.

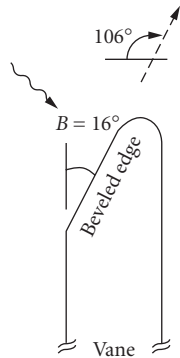


FIGURE 20 Placement of the bevel on the left side of the vane structure. (Ref. 6, p. 96.)

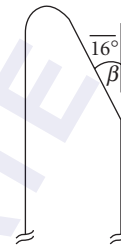


FIGURE 21 Placement of the bevel on the right side of the vane structure. (Ref. 6, p. 96.)



FIGURE 22 Vane structure angled at 90, 70, and 45°, respectively. (Ref. 6, p. 97.)

The next few paragraphs will present computer analysis results from two designs to show the effect of vanes on the propagated stray light. The vane angles used were 90, 70, and 45°, as depicted in Fig. 22.

The comparative stray light results for the Cassegrain system (Fig. 1) with a Martin Black coating on the vanes are shown in Fig. 23; in this system the vanes are on the main baffle, but not on the sunshade. There is no difference in the performance as the vane angle is varied from 45° to 90° (all three curves lie one on top of the other).

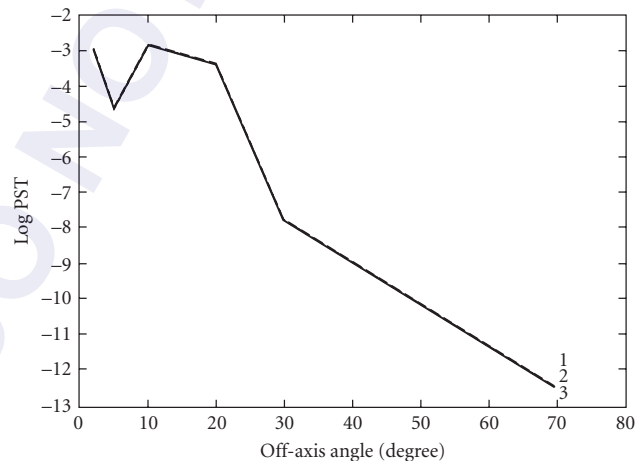


FIGURE 23 Cassegrain with Martin Black. Vane angles 1 = 90°, 2 = 70°, 3 = 45°. Log PST = detector irradiance divided by input irradiance. (Ref. 6, p. 97.)

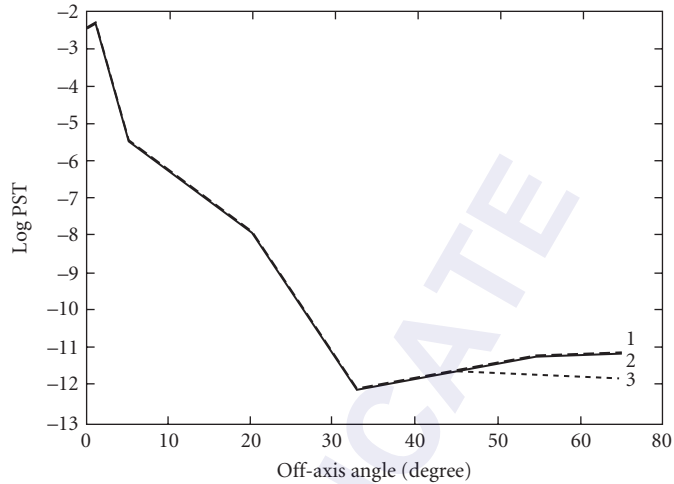


FIGURE 24 Z-system with Martin Black. Vane angles 1 = 45°, 2 = 70°, 3 = 90°. Vanes are on the sunshield. (Ref. 6, p. 114.)

The comparative results for the Z-system (Fig. 16) with vanes on the sunshield are shown in Fig. 24. The results differ from the Cassegrain results for source located at angles greater than 45° off-axis. This is because the primary side of the baffle is illuminated and scatters light directly to the primary mirror. The 70° baffles would fail for sources beyond 70° off-axis. The Cassegrain system has vanes on the main baffle (not the sunshade) and the sunshade occulted the direct illumination of the primary side of the 45° vanes. This accounts for the subtle but important difference in the results.

Usually the first-order scattering properties of the vane structures are more important than whether the vanes are angled or not. The results presented above confirm this statement. There are occasions where angled vanes would be beneficial, but to fully understand those cases a much longer explanation of diffuse vane baffle scatter is necessary. These results are detailed elsewhere.^{15,16}

There are special situations where angled vanes will have a significant advantage over annular vanes. One example is a bright source at a fixed offset angle. I have seen such a feature on a spaceborne telescope on a platform where there was nearby a brightly sunlit rocket-thruster casing at a fixed angle outside the field of view. Figure 25 shows the design where the vanes were aimed at the thruster at an angle where the primary mirror side (right side in Fig. 25) could not be directly illuminated by the sunlight scattered off the thruster. Under those circumstances most of the stray light

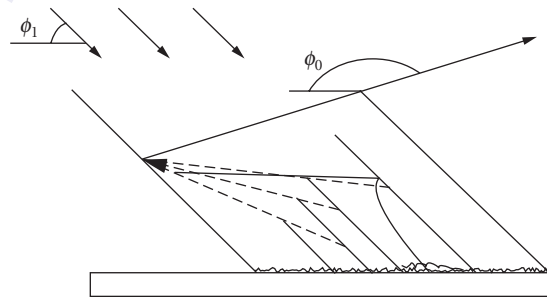


FIGURE 25 Angle-staggered vanes for fixed input angle. (Ref. 6, p. 104.)

had to make three scatters before exiting the vane cavity. In general, as soon as the position of the bright object is moved over a range of angles, the advantage of the angled vanes is lost. Nevertheless, there are many occasions within a sensor where the relative positions of a scattering source and a collecting object are fixed along a major stray light path. The front parts of the main-barrel baffle and the opening of the inner conical baffle in the Cassegrain design is an example. Many more examples could be cited. But the point is that you, as a designer, should first consider the first-order, single scatter paths off the baffle wall, each side of the vanes, and the bevel, for the full range of input values. Based on that information you can make the decision to use planar or angular vanes.

Vane Depth Considerations By varying the vane depth in the example analysis we can evaluate how the vane spacing-to-depth ratio affects system performance. Figure 26 gives the results of an analysis of the Cassegrain system with varying vane depths on the main baffle of 0.2, 0.4, and 0.8 inches. Figure 27 gives similar output from the Z-system analysis results. The performance of the system gets better as the vane depth increases from 0.2 to 0.4 inches, but there is little performance difference between the 0.4- and 0.8-inch baffle depths. The latter is the normal case. The 0.2-inch vane depth allows for a single path from the walls of the baffle tube, which increases the stray light propagation. Once that path is blocked by a greater vane depth, no further improvement should be expected due to further increases in vane depth.

The intent of presenting the two different optical designs was not to trade off one optical design against another. It needs to be made clear that the two optical sensors being used as examples are intentionally not equivalent from stray light design considerations. This is why the changes in performance are design-dependent. The nominal design of the eccentric pupil has a reimager, and the Cassegrain does not. The Cassegrain could have a reimager, in which case the stray light performance of both could be made essentially equal. It would depend on the optical design characteristics, $F/\#$, field of view, obscuration ratio, etc. The Cassegrain design has a specular sunshield and the Z-system has a vaned diffuse baffle structure. Which would perform better could only be determined after all of these features are considered.

To summarize, the general points being made in this section are

1. Usually, angled vane structure has little, if any, additional benefit over straight, annular vanes, and the annular vanes are much easier to fabricate and assemble.
2. Once the depth of a diffuse black vane structure is deep enough to block the single scatter path, further increases will not improve performance.

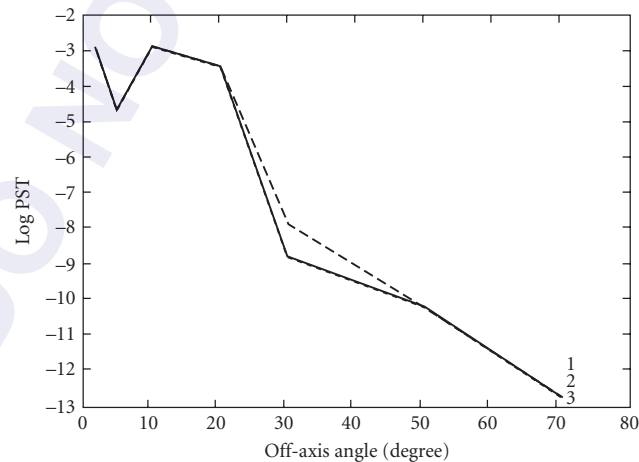


FIGURE 26 Cassegrain 90° baffles, coated with Martin Black, at varying depths; 1 = 0.2-inch, 2 = 0.4-inch, 3 = 0.8-inch depth. (Scatter is dominated by baffles.) (Ref. 6, p. 98.)

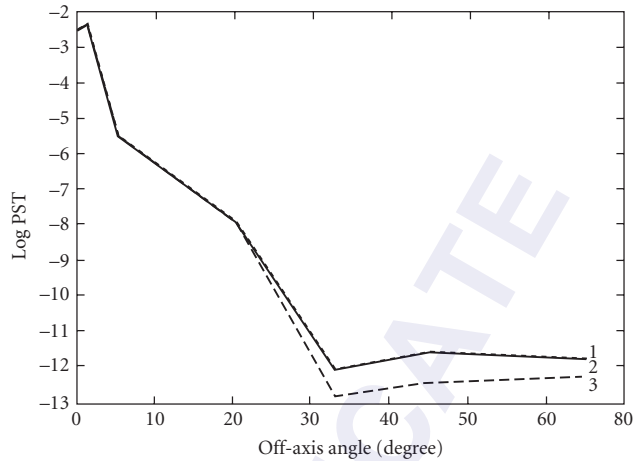


FIGURE 27 Z-system with varying vane depths. 1 = 0.2-inch, 2 = 0.4-inch, and 3 = 0.8-inch depth. (Ref. 6, p. 114.)

Specular Vanes Another aspect about vane structure that has been explored, but only in a limited way, is the specular vane cavity. Previous studies indicated that specular vanes have a problem with the aberrated rays and near specular angle scatter; this problem is severe enough to degrade the performance significantly.^{17,18} In another study by Freniere this was not always true.¹⁹ The ASAP²⁰ stray light software was used to evaluate the Z-system (Fig. 28) with (1) no vane structure, but with the main barrel baffle coated with Martin Black; (2) with Martin-Black-coated vanes; and (3) with a specular vane structure. The results show a dramatic degradation in the stray light performance without the coating on the main baffle tube. A subsequent specular baffle design developed by Nick Stavroudis has been shown to be a major improvement over previous concepts.²¹

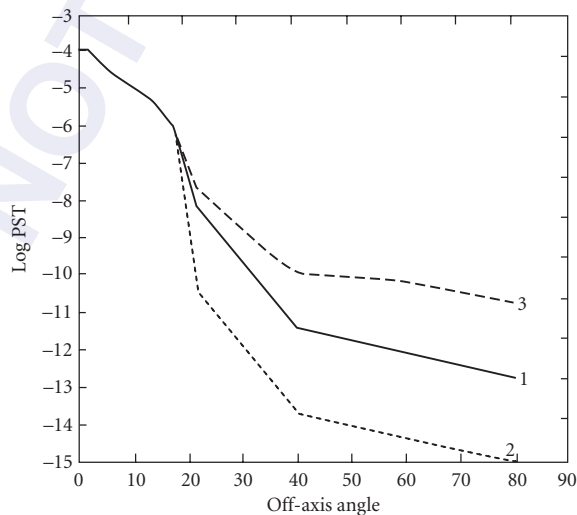


FIGURE 28 PST for unobscured pupil design without vane structure, with diffuse vane structure, and with specular vanes. 1 (solid) = no vanes, diffuse black coating; 2 (dotted) = diffuse vanes on main tube; 3 (dashed) = specular vanes on main tube. (Ref. 6, p. 115.)

Contamination Levels

Light scattered from a particulate-contaminated surface can have a pronounced effect on the stray light performance of a system.²²

I will now relate the performance of both designs (the centrally obscured Cassegrain, and the unobscured eccentric pupil) as a function of the level of scatter, per MIL-STD 1246A.²³ This analysis evaluates the sensor for different amounts of contamination on the optics only. The levels of contamination as defined in IEST-STD-CC1246D are for a distribution of particles with a specified range in particle sizes.

Ray Young used Mie scattering theory to predict the BRDF of a mirror covered with such MIL-STD distributions.²⁴ Table 1 was generated from Young's work for the 10 μm radiation. This table shows the base BRDF value and the BRDF slope that would be used in a typical stray light analysis program for input. The base value is the BRDF value at $(\beta - \beta_0) = 0.01$ and the second term is the slope of the BRDF in a (log-log) plot of BRDF versus $(\beta - \beta_0)$. β_0 is the sine of the angle of incidence and β is the sine of the observation angle.²⁵ The terms work equally well for out-of-plane values, but the above definitions, for simplicity, assume in-plane scattering data. See also the works of Spyak.²²

Spyak and Wolfe²⁶ did a series of experiments and calculations that relate BRDF to particulate contamination. They counted and sized particles on a mirror surface, and then measured the contribution of these particles to the mirror's BRDF at both visible (633 nm) and infrared (10.6 μm) wavelengths. They also performed Mie theory calculations and compared their calculations with the measured BRDFs. At both visible and infrared wavelengths, Mie theory calculations were a reasonable estimate for contribution of particulates to a mirror's BRDF. In most cases agreement between Mie calculations and their measurements were within a factor of two. Spyak and Wolfe also published Mie calculations of the BRDF expected from the MIL-STD-1246B (now IEST-STD-CC1246D) standard at 633 nm and 10.6 μm .*

Michael Dittman²⁷ has published a series of Mie calculations at five wavelengths. His calculations were done for the IEST 1246D distributions, but he considered an additional distribution in which the "particle slope" on the distribution was reduced from 0.926 to 0.383. The latter slope results in more large particles, and is commonly thought to be a more realistic distribution on surfaces that are exposed in a cleanroom environment.*

There is a problem with specifying the optics with this standard because it is difficult to reliably relate a level of contamination by particles to a BRDF performance. Two equal sizes and distributions of particulates may not give the same BRDF, because the index of refraction, the reflectivity, and the roughness of the particulates enter the calculations. In general, few people go to the trouble to determine these other factors. These factors will vary from one distribution to another. BRDF is the most usable value when performing a stray light analysis, so it should be the stray light specification. For manufacturing specifications, other parameters may be more appropriate, but they are not as good as BRDF for a stray light specification.

The level of scatter is also given in Table 1 along with the BRDF. The BRDF data from particulate scatter for the 5-, 10-, and 20- μm wavelengths for the 100, 300, and 500 contamination level have

TABLE 1 Mirror Scatter Relationships [Wavelength = 10 μm , BRDF Slope in Log $(\beta - \beta_0)$]

BRDF at $(\beta - \beta_0) = 0.01$	BRDF Slope	Cleanliness Level
0.02	-1.17	500
0.01	-1.17	454
0.001	-1.17	300
0.0001	-1.17	204
0.00001	-1.17	100

*Personal communication on partial contamination paragraphs, contributed by Dr. Gary Peterson, Brealut Research Organization, Inc.

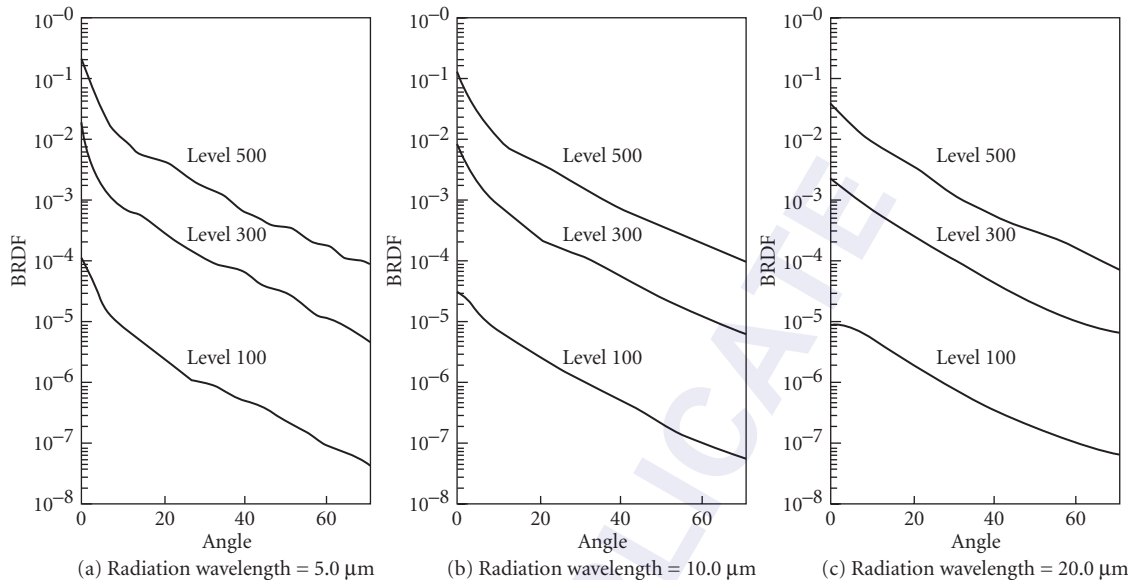


FIGURE 29 Predicted BRDFs on particles deposited on low-scatter mirror for cleanliness levels of 100, 300, and 500 at radiation wavelengths of 5.0, 10.0, and 20.0 μm .²⁴ (Ref. 6, p. 105.)

been plotted in Fig. 29. Consensus, not factually documented, indicates that the current state of the art of contamination control is at the cleanliness level of 300 to 500 for the 10- μm -wavelength region. Measured BRDFs below level 200 are achievable in the lab for short periods of time. A stray light analyst is strongly advised not to predict a system's performance with values below $1.0\text{E-}3 \text{ sr}^{-1}$ in the 10- μm region. Based on historical performance, mirrors in the IR (10- μm region) consistently degrade to this value, usually because of particulate scatter. Research work performed under Rome Air Development Center contract for the detection, prevention, and removal of contamination from the ground and in space could greatly reduce the degradation currently experienced by IR sensors.²⁸

Hal Bennett presents the significance of particulate scatter, as shown in Fig. 30.²⁹ This figure shows an agreement between measured data and theoretical data, and illustrates why IR sensors are usually more sensitive to particulate scatter than RMS scatter; the opposite is true in the visible. Figure 30 also indicates why the wavelength scaling law does not usually relate visible BRDF measurements to BRDF measurements in the IR. The physical process is different.

Figures 31 and 32 are the representative point source transmittances (defined as the irradiance on the detector divided by the incident irradiance) for the cleanliness levels of 100 through 500 for each design. The Cassegrain is much less affected by changes in contamination level, because the scatter from the black-coated surfaces dominates all other scatters. If the system had a reimager its performance would be better because these black surfaces would be blocked from the field of view of the detector, and the stray light performance would be due to the cleanliness level of the optics. The eccentric pupil design is sensitive to changes in the mirror coatings because it does have a reimager, and the major source of scatter is from the mirror surfaces.

In summary, the impact of particle contamination on the performance of a system will depend on how well the system is designed to suppress stray light. The goal is to be limited by a single optical element, such as the collecting lens or mirror, which is the objective of the system. The eccentric pupil design (Z-system) has this design feature. The better the optical design from a stray light point of view, the more the system's performance will be degraded by particle contamination. The more the system performance is determined by the black coatings, the more it will be sensitive to degradations in the coatings on the baffles.

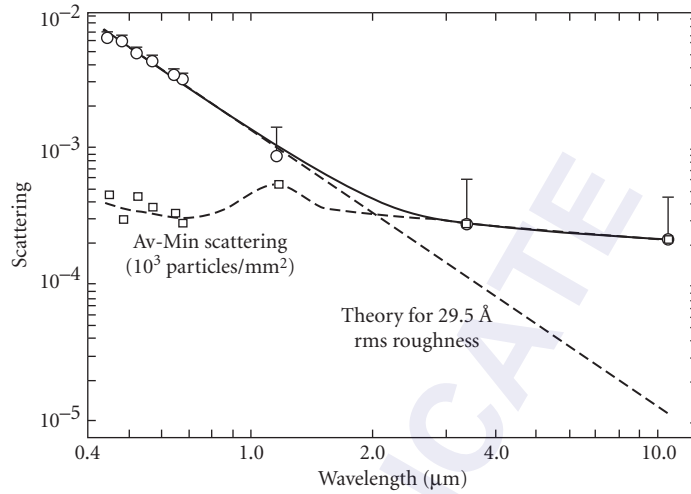


FIGURE 30 Scattering from polished dense flint glass. The diagonal line gives the contribution predicted for microirregularity scattering by a 29.5 Å rough surface. Circles indicate the minimum scattering observed, and the bars and squares the difference between the average and minimum scattering observed at several points on the surface. This difference may be related to particulate scattering. (Ref. 29, p. 32.)

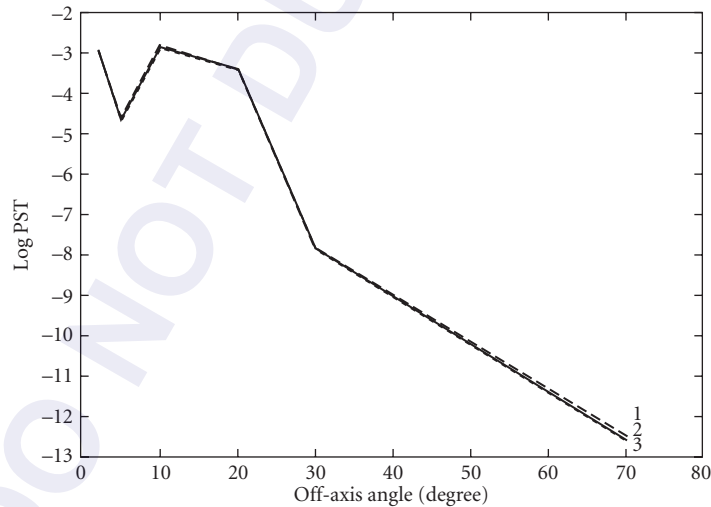


FIGURE 31 Cassegrain system with mirrors at all five contamination levels. 1 = 100, 2 = 204, 3 = 300, 4 = 454, 5 = 500. (Ref. 5, p. 113.)

Strut Design

In a centrally obscured system the central obscuration must be supported. In some designs (Schmidt-Cassegrains) the obscuration can be supported by a refractive element, but in most designs some form of struts are used. The most common error in strut design is to specify manufacture

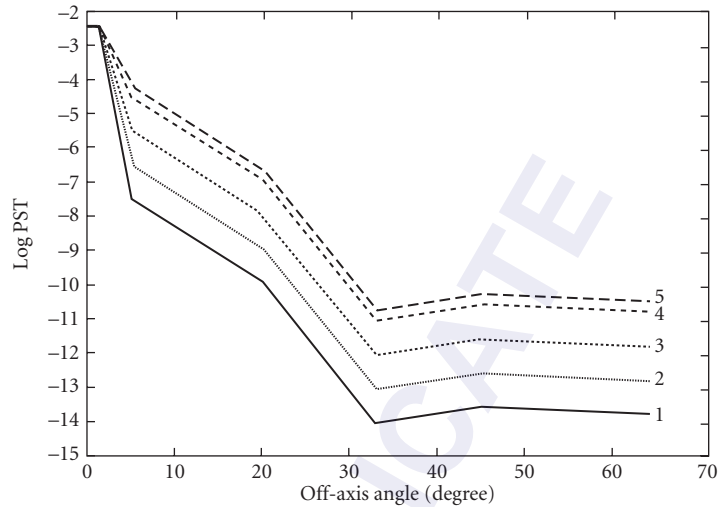


FIGURE 32 Z-system with mirrors at all five contamination levels. 1 = 100, 2 = 204, 3 = 300, 4 = 454, 5 = 500. (Ref. 6, p. 113.)

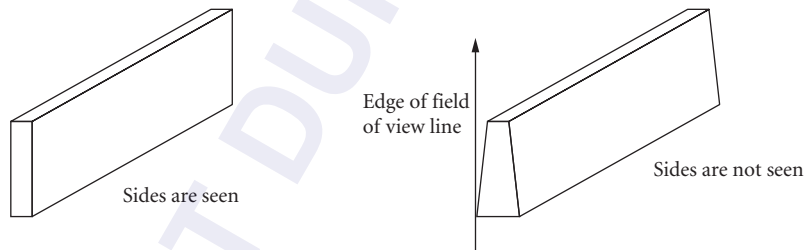


FIGURE 33 Angled strut design does not allow the detector to see the sides of the strut.

from a slab or plate of coated metal. Because all detectors have some finite field of view, the scatter from the sides of the struts can be seen from the image plane. Usually the struts are out “front” and exposed to more stray light sources than the objects deeper into the system. The near off-axis angles of incidence of scattered light off of the strut sides make for very high scattering toward the detector.

The proper strut design will preclude this path by making the object end of the strut narrower than the side nearest the objective (primary). This shape, shown in profile in Fig. 33, does not allow the detector to see the sides of the struts. The angle of the taper depends upon the object space field of view of the detector. It requires only a small change in design to remove this stray light path.

Basic Equation of Radiation Transfer

This section briefly discusses the most fundamental equation needed to perform the quantitative calculations of a stray light analysis. It reinforces the concept of first identifying what the detector can see and working on the geometry of the system to limit the stray light propagation, and not the BRDF term.

The fundamental equation relating to power transfer from one section to another is:

$$d\Phi_c = L_s(\theta_c, \phi_c) dA_s \frac{\cos(\phi_s) dA_c \cos(\phi_c)}{R_{sc}^2} \quad (2)$$

where $d\Phi_c$ is the differential power transferred, $L_s(\theta_c, \phi_c)$ is the radiance of the source section, dA_s and dA_c are the elemental areas of the source and collector, and ϕ_s and ϕ_c are the angles that the line of sight from the source to the collector makes with their respective normals. This equation can be rewritten as three factors that help clarify the reduction of scattered radiation.

$$d\Phi_c = \frac{L_s(\theta_c, \phi_c)}{E(\theta_i, \phi_i)} E(\theta_i, \phi_i) dA_s \frac{\cos(\phi_s) dA_c \cos(\phi_c)}{R_{sc}^2} \quad (3)$$

$$d\Phi_c = \text{BRDF}(\theta_i, \phi_i; \theta_c, \phi_c) d\Phi_s(\theta_i, \phi_i) d\Omega_{sc} \cos(\phi_s) \quad (4)$$

$$d\Phi_c = \text{BRDF}(\theta_i, \phi_i; \theta_c, \phi_c) d\Phi_s(\theta_i, \phi_i) \text{GCF}_{sc} \pi \quad (5)$$

$E(\theta_i, \phi_i)$ is the incident irradiance on the source section dA_s . GCF_{sc} is the projected solid angle from the source to the collector divided by π .

The GCF is independent of the first two terms and solely determined by the geometry of the system, including obscurations. The first term, $\text{BRDF}(\theta_i, \phi_i; \theta_c, \phi_c)$, is the bidirectional reflectance distribution function. It is usually considered independent of the second term, the incident power, and is therefore a function of the surface characteristics only. When reducing stray radiation propagation, one or more of these terms must be reduced. If any one of these terms is reduced to zero, no power will be transferred between the source and collector.

Stray Radiation Paths

Since the third term (GCF) in Eq. (4) is the *only* term that can be reduced to *zero*, it should receive attention first. This is a crucial point in a stray light analysis. Therefore, the logical starting place for stray light reduction is with the critical objects, since it is the GCF terms for these transfers which can be reduced to zero. Most novice analysts make the mistake of working on the BRDF term first.

$$\text{GCF} = \frac{\cos(\phi_s) dA_c \cos(\phi_c)}{\pi R_{sc}^2}$$

The apparent possibilities for decreasing the GCF are to increase R_{sc} , ϕ_s , ϕ_c or to reduce the area dA_c . Not readily apparent is that the GCF is limited by apertures and obstructions. These features will, in some cases, block out the entire view of the source section from the collector so that there is no direct path. This is the mathematical basis for the logical approach, discussed at the beginning of the chapter. First block off as many direct paths of unwanted energy to the detector as possible, and then minimize the GCF for the remaining paths.

Point Source Transmittance Definitions

There are five common ways to define the merit function of the stray light in an optical sensor. The most common and preferred method is to define it as the output irradiance divided by the input irradiance, in terms of the *normalized detector irradiance* (NDI),³⁰ or in terms of the *point*

source normalized irradiance transmittance (PSNIT).³¹ This merit function is appropriate because it describes an irradiance transmittance, and it is relatively independent of the detector size.

A term often used in the past was the *off-axis rejection* (OAR), defined as the detector power divided by the input power from the same source *on-axis*. The term *rejection* is a misnomer because by definition the term describes a power transmittance, which can have little correlation with the rejected stray light. The second objection is that as a merit function it varies significantly with the detector size. If you double the area of the detector, the OAR will increase by about the same factor even though the system hasn't performed significantly worse in any way.

Another term commonly used is the system's stray light *point source power transmittance* (PSPT), or its reciprocal, the *attenuation* of the system. The PSPT is the detector power divided by the input power into the sensor from the specified *off-axis angle*. Again, this term varies with the detector size. Sometimes there is no well-defined entrance port so the denominator is impossible to define. Note that the magnitude of attenuation would normally be expressed in terms of a positive exponential. Beware that attenuations are often incorrectly called out with negative exponents.

A final PST definition that is sometimes specified is the *point source irradiance transmittance* (PSIT), defined as the output irradiance divided by the entrance port input irradiance. This definition becomes inappropriate when there is no clearly defined entrance port.

Surface Scattering Characteristics

Of the three potentially important factors in scattered radiation analysis cited above (the radiance of the undesirable source or sources, the geometry of the scattered radiation paths (GCF), and the surface scattering characteristics, (BRDF)), usually the first possibility considered is to improve the surface coatings or the addition of vane structure. In concept it *appears* to be the right place to start and that it is straightforward. Neither is the case; the BRDF never goes to zero as does the GCF, and the BRDF varies with input and output angles. However, with accurate *bidirectional reflectance distribution function* (BRDF) data and knowledge about the variations with applications, time, wavelength, and other factors, BRDF problems can be dealt with. The scattering characteristics of surfaces are discussed by Church, and the scattering characteristics of black coatings by Pompea and Breault elsewhere in this *Handbook*. The addition of vane sections on baffles can usually be considered as a specialized "coating" with its own specialized BRDF.

BRDF Characteristics

Usually, BRDF data that are presented represent only one profile of the BRDF, and many such profiles for various angles of incidence are necessary for understanding the scattering characteristics. However, studies have shown that a single profile of a mirror's surface scattering characteristics can be used, with some approximations, to define the BRDF for all angles of incidence.³² This is a significant achievement. It reduces the amount of data that must be taken, and it makes it easier to calculate, or estimate, the BRDF value for any set of input and output angles. The BRDF can also be reconstructed for cases where only a single profile of the function has been presented, which has been the usual practice.

The approximation has its limitations, as clearly detailed by Stover.³³ The approximation is quite good for nominal angles of incidence (see Fig. 34).³⁴ However, it breaks down for very high θ_i and high observation angles θ_o .

It is important to understand qualitatively the scattering characteristics of diffuse black coatings. Figure 35 shows the BRDF profile of Martin Black at 10.6 μm for several angles of incidence.³⁵ At near-normal angles of incidence the BRDF values are bowl-shaped; the values increase at large observation angles from the normal. At high angles of incidence the BRDF values in the near specular direction have increased by 2.5 orders of magnitude. There is a good discussion of the qualitative characteristics of diffuse black surfaces by Pompea and Breault elsewhere in this *Handbook*.

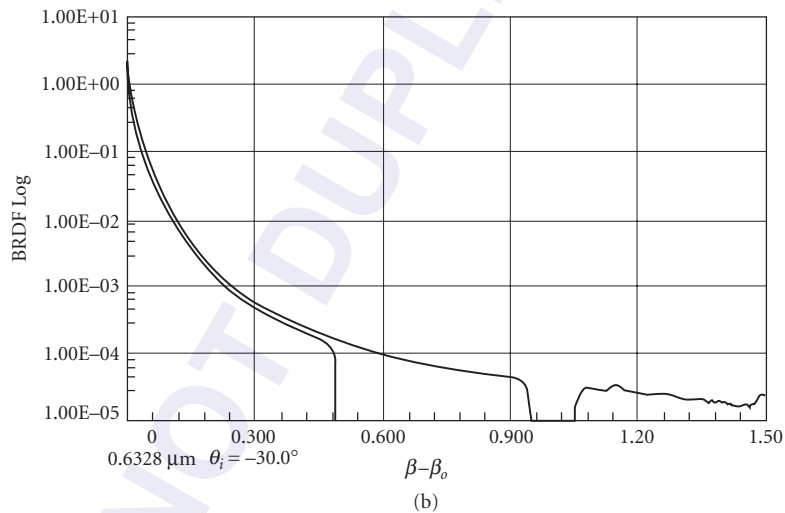
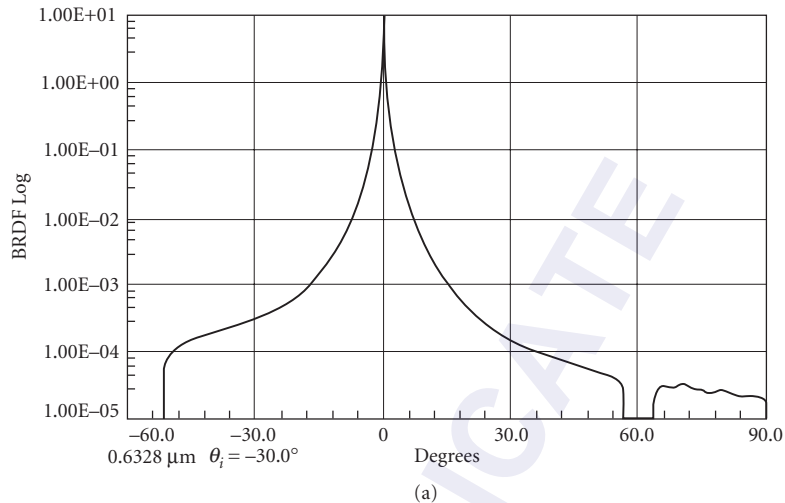


FIGURE 34 (a) The BRDF is asymmetrical when plotted against $\theta_s - \theta_i$, (b) The data in (a) exhibits near symmetry when plotted against $|\beta - \beta_o|$. The slight deviation from symmetry is due to the factor $(\cos \theta_s Q)$, where Q is a polarization factor. (Ref. 34, p. 69, reprinted with permission.)

7.4 OPTICAL SOFTWARE FOR STRAY LIGHT ANALYSIS

There is a small bevy of commercial optical software programs on the market that perform stray light analyses or some aspects of the stray light problems that are typically encountered. One must be knowledgeable of what each package can accomplish so the best thing to do is to ask for a demonstration of what each optical software manufacturer has to offer that is relevant to your design challenge. Here's a summary of a few commercially available optical software codes for stray light analysis. Note also that the programs' capabilities are always in a state of flux.

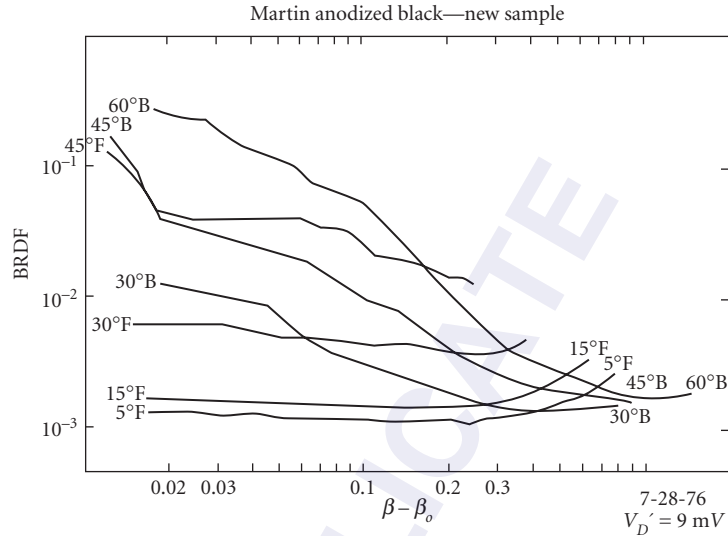


FIGURE 35 BRDF profile of Martin Black at $10\ \mu\text{m}$. (F. O. Bartell *et al.*, "A Study Leading to Improvements in Radiation Focusing and Control in Infrared Sensors," Final Report, Army Materials and Mechanics Research Center, December 1976.)

ASAP, by Breault Research Organization, Inc.

ASAP Optical Software was developed to meet design and analysis criteria of imaging and illumination systems and the unique challenges of stray light analysis with CAD interoperability. ASAP is powered by the ASAP nonsequential ray-tracing engine—known throughout the optics industry for its accuracy and efficiency. Rays can encounter surfaces in any order and any number of times, with automatic ray splitting. Optimized for speed, ASAP will trace millions of rays in minutes. The standard edition of ASAP includes a license of the SolidWorks Parts Only 3D Modeling Engine—an intuitive 3D-design environment optimized for use with ASAP. The user can write ASAP geometry files from within SolidWorks, import XML files, or use BRO's proprietary smart-IGES system to import system models from any CAD package while maintaining fast, efficient ray trace speed.

Use ASAP to model complex imaging systems, illumination systems, and light-concentrating devices. Create highly accurate source models using source images, point sources, ray grids, and fans. Model incandescent bulbs, LEDs, CCFLs, and HID arc lamps, or import from the BRO Light Source Library. Perform the analyses necessary to validate your designs without experimental prototyping.

ASAP includes a distributed-processing capability allowing the user to complete big design and analysis jobs effectively in short-time span—spawn up to 5 additional ASAP sessions on other local area network (LAN), without leaving your desk. Web site: www.breault.com

FRED, by Photon Engineering

FRED is an optical engineering software package that uses a statistical ray sampling approach to analyzing incoherent stray light mechanisms in optical systems. The user can assign one or more of many different BSDF scatter functions to a surface. When a ray is incident on the surface, a specified number of scatter rays are generated in random directions into the hemisphere according to a uniformly random angular distribution, although a Monte-Carlo technique can be used

to generate scatter rays with directional ray density proportional to the BSDF function. The power assigned to any particular scatter ray is proportional to the incident ray power and value of the BSDF function evaluated in the direction of the scatter ray. Uniform sampling of angle space has the advantage that lower power scatter paths will be realized with higher probability. Ray-density directional sampling has the advantage that more rays are directed in the higher power directions which decreases the statistical noise in those directions.

When analyzing scatter from a surface, the user is often interested in only a subset of the entire hemispherical angular range. In these cases a technique called importance sampling can often be used to great advantage. The user can specify that the same number of scatter rays be generated into the importance sample direction as was originally directed into the hemisphere to decrease the statistical noise. An alternative is to decrease the number of rays but still achieve the same noise level. The desired angular range can be specified in a variety of ways and FRED will then direct scatter rays randomly into the specified range. The power assigned to each scatter ray is adjusted to account for the fact that it can be directed only into a subset of the entire hemisphere. Web site: www.photonengr.com

LightTools and CodeV, by Optical Research Associates

LightTools is a complete illumination design and analysis software product. It combines full optical accuracy, powerful optical and illumination analysis, and an intuitive graphical user interface in a 3D solid modeling environment where models interact with rays to produce virtual prototypes of manufacturable systems. A fully integrated illumination optimization capability automatically improves model performance. *LightTools'* Monte Carlo ray tracing facilitates accurate spectral modeling of the illuminance, luminance, intensity distributions, and CIE colorimetric data anywhere in the optomechanical model.

In addition to illumination system design, *LightTools* supports many other applications, from packaging studies to stray light analysis. For example, its ray path visualization collects and displays information about ray-surface interactions to identify system elements that are contributing to light loss, scatter, unintentional reflections, or ghost images. A unique point-and-shoot ray tracing capability allows rapid, interactive evaluation of optical behavior.

CODE V is used for the optimization, analysis, and tolerancing of image-forming optical systems and free-space photonic devices. Its many capabilities include powerful local and global optimization for optics, fast wavefront differential tolerancing that allows as-built considerations to be evaluated throughout the design process, and highly accurate diffraction beam propagation analysis. For stray light applications, CODE V can be used to analyze ghosts in imaging systems due to Fresnel reflections. Web site: www.opticalres.com

ZEMAX, by ZEMAX Development Corporation

The ZEMAX program, from ZEMAX Development Corporation, has two modes of use. Its primary use is as a sequential optical design (optimization) program. In this mode it has tools to help identify location of ghost pupils and images resulting from Fresnel reflections. A separate nonsequential mode has many capabilities necessary for stray light analysis, including scatter modeling with importance sampling.

ZEMAX's Nonsequential ray-tracing capabilities can further be extended to finding rays which have specific characteristics or properties. For example, imagine you are studying the stray light in a telescope:

How significant are rays which "ghost" reflected off of various surfaces (both mechanical and optical)?

Rays which are experience multiple reflections may be important, but how significant are those which experience more than four reflections?

How effective is a strategically placed baffle in terms of limiting the amount of stray light on the detector?

Website: www.zemax.com

TracePRO, by Lambda Research Corporation

TracePro, from Lambda Research Corporation, is a 3D Computer Aided Design (CAD) program for simulating the performance of illumination and optical systems. TracePro can model the propagation of light in imaging and nonimaging optomechanical systems. Models are created by combining imported lens designs, imported CAD geometry (IGES, STEP, SolidWorks, Pro/E, CATIA, or Inventor files), and geometrical objects created using TracePro's user interface. Optical properties are then assigned to each solid and surface using the TracePro interface or through the TracePro Bridge for SolidWorks. Source models are added by specifying grids, surface emitters, ray file data or by using the surface source utility. Rays are ray traced through the model, while keeping track of absorption, specular reflection and transmission, fluorescence and scatter at each intersected surface or volume scatter site.

From TracePro models, the user may ray trace and analyze:

- Light distributions in illumination and imaging systems
- Stray light, scattered light, and aperture diffraction
- Throughput, loss, or system transmittance
- Flux or power absorbed by surfaces and bulk media
- Light scattering in biological tissue
- Polarization effects
- Fluorescence effects
- Birefringence effects
- Lit Appearance

Website: www.lambdare.com

SPEOS, by Optis

OPTIS' simulation software family, SPEOS and OptisWorks. It manages and optimizes many of the optical aspects of a broad range of sensors: reflection, refraction, scatter from surfaces, diffraction, absorption, polarization, and Gaussian beam propagation. It calculates stray light, illumination, and realistic optical simulations. Any product that needs to manage interactions of light and surfaces is calculated. It deals with the various types of light sources also. The simulations limit the need for costly prototyping of systems.

OPTIS simulation software allows the designer to "see" and realistically render products to depict what the final performance of the illuminator will look like in its applied application with stunning similarity. Its software produces a unique and accurate physiological human vision model of the final lit product for comfort, safety, and performance.

7.5 METHODS

There are two distinct methods that have been used to evaluate a system for stray radiation. You can either build the system and test it, or you can model the system and try to predict its performance. Both methods have advantages and disadvantages. Taken *together* the two methods provide the means to ensure that the system will perform as desired.

Build-and-Test Approach

A common approach is to make the system and either use it or test it for stray radiation rejection. Certainly if the system consistently performs satisfactorily *in its operational environment*, it has passed the ultimate test. But what if it does not meet the desired or expected level of performance? Making more systems to test becomes expensive rapidly. In fact, for very large systems, usually only modifications (“fixes”) can be contemplated because of the high cost. This is not the only argument against the build-and-test approach. The tests are rarely designed to determine *how* the scattered radiation is propagating through the system and which surfaces contribute most of the undesired radiation.

It is this information, and a thorough knowledge of the surface scattering characteristics, that is necessary to make measurable improvements to the system. Such a test, when determining the propagation paths, should yield information about how the system is reacting to its *test* environment, including the test equipment. Unless the tests are being conducted in the environment for which the system was designed, it is imperative to determine that the *test* environment is not giving erroneous results (either better or worse). Without analyzing the test configuration, you should expect that the environment *will* affect the system stray light measurements. It is also incorrect to assume that the test environment can only add to the stray light background. It is sometimes assumed that if the system passes the stray light tests in the lab it will only perform better in space or wherever its intended environment is. This is not necessarily true.

Now that several points have been made about the difficulty of making valid experimental tests, it must be stated that valid tests can and should be made. The measurement costs need not be prohibitive. Even relatively large optical systems have been fabricated and then modestly redesigned. Changes to the system can be made until the desired information and stray radiation rejection is attained. In some cases it will be less expensive to test an existing system and modify it if necessary than to analyze the system with computer software.

The system-level test need not be extensive; it is not necessary to have an all-encompassing measurement from on-axis to 90° off-axis. Indeed, few facilities are capable of making such tests when the attenuation gets even modestly high. An important point to recognize is that the most important paths to check are those at the nearer off-axis angles where the attenuation is not so high. These can usually be measured reliably.

At small off-axis angles the stray light noise is more often much higher than the detector background noise, while at the higher off-axis angles the stray light noise is well below the electronic/detector noise. From a performance point of view, at the higher off-axis angles there is usually only one additional scattering object (scatter from the main baffle) before these same near off-axis angle paths are encountered or are reinvolved. The validation of the analysis will only be susceptible to the scatter from this one object that can't be fully tested at the system level, but most of the scatter paths, and usually all the most important ones, will have been validated by the near off-axis measurements.

This one additional surface scatter most often (especially on space-based sensors) involves the vanes on the main baffle that shields the primary objective. It will normally reduce the optical noise incident upon it by four to five orders of magnitude. That's why the optical noise goes dramatically below the electronic noise. Its most important role is to occult the direct illumination of the objective which is usually part of the most significant direct scatter path. The performance of this baffle and its vane structure could be analyzed separately and then measured independently to confirm that it too will perform as predicted.

NOTE: Contrary to some published papers you cannot, in general, multiply the stray light transmittance of two parts of a sensor and determine the system's overall performance. Although the main baffle system can be analyzed (or measured) independently from the rest of the system, it is not correct to take its performance and multiply it times the stray light performance of the rest of the system. The stray light propagation paths are far more important than the magnitudes of the two parts. In the above analysis where it was proposed that the main baffle could be measured independently it was to confirm its performance alone. A full-system stray light analysis was assumed.

Computer Analysis

As with the experimental tests, computerized analyses are also subject to errors. The three most significant ones are software limitations, scatter data of samples (not the real system), and user error. No software is capable of putting in every detail of a complex design, yet the computer model must faithfully represent the actual performance of the system. On the other hand, the software can put in “parts” with far greater mathematical precision than these parts can actually be assembled. Unless special studies are made the analyst does not usually account for assembly errors that might affect the actual system. The scatter characteristics of the surfaces, usually defined in terms of the bi-directional reflectance distribution function (BRDF), are usually measured on sample substrates, and controls must be exercised to ensure that the samples tested represent the sensor’s actual coatings, and that they do not change with time. The stray light analysis programs are also subject to errors in determining the significant paths. The experimental test is for the actual design, with real coatings, and will include any extraneous unintentional paths due to misalignment or other causes.

On the positive side, a software program can point out many flaws in the system that contribute stray radiation by considering the input BRDF characteristics of the coatings. A program can also do trade-off studies, parametric analysis, and in many other ways aid in the study of alternate designs. The analysis of the paths of scatter will suggest meaningful modifications and help to discard impossible designs. These analyses allow designers to test designs and make modifications before the design goes into production. This is very useful, since rejecting a sensor design is much easier when it is on paper than after it has already been built. It is usually much more cost-efficient than cutting new hardware, redesigning the system, or making fixes on the built system.

If you are in a field related to the optical design of a sensor, be it at the design level or the system level, you know that it would be preposterous to perform the optical design analysis and then put the system together without testing it for its image quality. Yet that is how far the pendulum has swung in favor of performing a stray light analysis over making a system-level stray light test. It reflects a major change in attitude since the early 1970s. It has been stated by stray light analysts that the reliability of a stray light analysis is now much higher than experimental test results, so some people avoid the latter. While there is a degree of truth in this statement, it is wrong to omit the stray light test at the system level.

The advantages and disadvantages of the two methods are summarized in Fig. 36. The disadvantages of the build-and-test approach are the strengths of the analysis method, whereas the strengths of the build-and-test approach cover the weaknesses of an analysis. Taken together these two methods give the greatest amount of reliable information which you can use to create the optimal system and have confidence in its performance. Jointly, they indicate the reliability of the analysis and test results.

Strength	Inexpensive	High	Easy	Real performance complete with manufacturing error	No missed paths	Real BRDF
Weakness	Expensive	Limited	Hard	Models only	Programmer error	Sampled BRDF measurements
	Cost	Information level	Changes	Real-world performance	Error	BRDF
	<div style="display: inline-block; border: 1px solid black; width: 15px; height: 15px; margin-right: 5px;"></div> Build and test		<div style="display: inline-block; border: 1px solid black; width: 15px; height: 15px; background-color: #cccccc; margin-right: 5px;"></div> Analysis			

FIGURE 36 Build-and-test and analysis methods complement each other.

7.6 CONCLUSION

In summary, the issues involved in designing a system with stray light suppression in mind are

- I. System design concepts
 - A. Critical objects seen by the detector
 - B. Illuminated objects
 - C. Lyot stops
 - D. Field stops
 - E. Optical designs
- II. Baffle and vane design
 - A. Diffuse and specular vane cavities
 - B. Vane edge scatter
- III. Diffraction
- IV. Strut design
- V. Scattering theory
- VI. BRDF data
 - A. Log BRDF versus θ
 - B. Log BRDF versus $\log(\beta - \beta_0)$
 - C. Polar plots
 - D. Isometric projections (3-D characteristics)
- VII. Coatings
 - A. Paints and anodized surfaces
 - B. AR coatings and other thin films
 - C. Mirror coatings
- VIII. Thermal emission
- IX. Ghost images
- X. Software
- XI. Detection, prevention, and removal of contamination

A step by step procedure that can help you to improve your system is:

- I. Start from the detector and identify what objects, called “critical objects,” can be seen from various positions on the detector. Be sure to include a point near the edge of the detector.
- II. Work to remove the number of critical objects that the detector can see.
- III. Determine what objects the source of unwanted radiation can see, called the “illuminated objects.”
- IV. If possible, reduce the number of illuminated objects seen.
- V. If there are illuminated objects that are also critical objects, work very hard on these paths. Orders of magnitude in improvement will be your reward.
- VI. If task V is not possible, then the computations are quite easy.
 - A. Calculate the power incident on the illuminated/critical objects.
 - B. Use Eq. (1) to calculate the transfer of power from the critical objects to the detector. Remember to properly account for the input and output angles when calculating the BRDF. *Do not* use a straight lambertian scatter distribution; there is no such distribution in reality.
- VII. Find all the paths connecting the illuminated objects to the critical objects.
- VIII. Evaluate the corresponding input and output angles at the illuminated and critical objects.

- IX. Determine if vane structure will help, or if some other redesign will effectively block these paths.
- X. For the calculated input and output angles, evaluate which coating would be lowest.
- XI. Perform the stray light calculation using Eq. (1) in an iterative fashion. This should determine the most significant stray light path and quantify the amount of stray light on the detector
- XII. Perform the above tasks for a series of off-axis positions of the point source.

7.7 SOURCES OF INFORMATION ON STRAY LIGHT AND SCATTERED LIGHT

- J. D. Lytle and H. Morrow (eds.), "Stray Light Problems in Optical Systems," *Proc. SPIE*, vol. 107, April 18–21, 1977 (22 papers).
- M. Kahan (ed.), "Optics in Adverse Environments," *Proc. SPIE*, vol. 216, Feb. 4–5, 1980 (30 papers).
- G. H. Hunt (ed.), "Radiation Scattering in Optical Systems," *Proc. SPIE*, vol. 257, Sept. 30–Oct. 1, 1980 (28 papers).
- S. Musikant (ed.), "Scattering in Optical Materials," *Proc. SPIE*, vol. 362, Aug. 25–27, 1982 (28 papers).
- R. P. Breault (ed.), "Generation, Measurement, and Control of Stray Radiation III," *Proc. SPIE*, vol. 384, Jan. 18–19, 1983 (15 papers).
- R. P. Breault (ed.), "Stray Radiation IV," *Proc. SPIE*, vol. 511, Aug. 23, 1984 (14 papers).
- R. P. Breault (ed.), "Stray Radiation V," *Proc. SPIE*, vol. 675, Aug. 18, 1986 (46 papers).
- R. P. Breault (ed.), "Stray Light and Contamination in Optical Systems," *Proc. SPIE*, vol. 967, Aug. 17–19, 1988 (33 papers).
- J. C. Stover (ed.), "Scatter from Optical Components," *Proc. SPIE*, vol. 1165, Aug. 8–10, 1989 (42 papers).
- R. P. Breault (ed.), "Stray Radiation in Optical Systems," *Proc. SPIE*, vol. 1331, July 12–13, 1990 (29 papers).
- J. C. Stover (ed.), "Optical Scatter: Applications, Measurement, and Theory," *Proc. SPIE*, vol. 1530, July 21–27, 1991.
- R. P. Breault (ed.), "Stray Light and Contamination in Optical Systems II," *Proc. SPIE*, vol. 1753, July 21–23, 1992.
- F. O. Bartell et al., "A Study Leading to Improvements in Radiation Focusing and Control in Infrared Sensors," *Final Report Prepared for Army Materials and Mechanics Research Center*, December 1976.
- J. A. Gunderson, "Goniometric Reflection Scattering Measurements and Techniques at 10.6 Micrometers," M.S. thesis, University of Arizona, 1977.
- P. J. Peters, "Stray Light Control, Evaluation, and Suppression," *Proc. SPIE*, vol. 531, January 1985.
- T. W. Stuhlinger, "Bidirectional Reflectance Distribution Function (BRDF) of Gold-Plated Sandpaper," M.S. thesis, University of Arizona, 1981.
- A. W. Greynolds, "Computer-Assisted Design of Well-Baffled Axially Symmetric Optical Systems," M.S. thesis, University of Arizona, 1981.
- J. W. Figoski, "Interferometric Technique for the Reduction of Scattered Light," M.S. thesis, University of Arizona, 1977.
- J. S. Fender, "An Investigation of Computer-Assisted Stray Radiation Analysis Programs," Ph.D. dissertation, University of Arizona, 1981.
- D. A. Thomas, "Light Scattering from Reflecting Optical Surfaces," Ph.D. dissertation, University of Arizona, 1980.
- R. P. Breault, "Suppression of Scattered Light," Ph.D. dissertation, University of Arizona, 1979.
- P. R. Spyak, "A Cryogenic Scatterometer and Scatter from Particulate Contaminants on Mirrors," Ph.D. dissertation, University of Arizona, 1990.
- F. O. Bartell, "Blackbody Simulator Cavity Radiation Theory," Ph.D. dissertation, University of Arizona, 1978.

- L. D. Brooks, "Microprocessor-based Instrumentation for BSDF Measurements from Visible to FIR," Ph.D. dissertation, University of Arizona, 1982.
- A. G. Lusk, "Measurements of the Light Scattering Profile of Small Size Parameter Fibers," M.S. thesis, University of Arizona, 1987.
- K. Nahm, "Light Scattering by Polystyrene Spheres on a Conducting Plane," Ph.D. dissertation, University of Arizona, 1985.
- G. W. Videen, "Light Scattering from a Sphere on or Near an Interface," Ph.D. dissertation, University of Arizona, 1992.
- Y. Wang, "Comparisons of BRDF Theories with Experiment," Ph.D. dissertation, University of Arizona, 1983.
- S. J. Wein, "Small-Angle Scatter Measurement," Ph.D. dissertation, University of Arizona, 1989.
- J. M. Bennett and L. Mattsson, *Introduction to Surface Roughness and Scattering*, Optical Society of America, Washington, D.C., 1989.
- J. C. Stover, *Optical Scattering Measurement and Analysis*, McGraw-Hill, Inc., New York, 1990.

7.8 REFERENCES

1. R. P. Breault, "Problems and Techniques in Stray Radiation Suppression," *Stray Light Problems in Optical Systems*, J. D. Lytle and Howard Morrow (eds.), *Proc. SPIE* **107**, 1977, pp. 2–23.
2. R. P. Breault, A. W. Greynolds, and S. R. Lange, "APART/PADE Version 7: A Deterministic Computer Program Used to Calculate Scattered and Diffracted Energy," *Radiation Scattering in Optical Systems*, G. Hunt (ed.), *Proc. SPIE* **257**, 1980, pp. 50–63.
3. R. V. Shack, "Analytic System Design with a Pencil and Ruler—The Advantages of the $y-\bar{y}$ Diagram," *Applications of Geometrical Optics*, *Proc. SPIE* **39**, 1973.
4. S. R. Lange, R. P. Breault, and A. W. Greynolds, "APART, A First-Order Deterministic Stray Radiation Analysis Program," *Stray-Light Problems in Optical Systems*, J. D. Lytle and H. Morrow (eds.), *Proc. SPIE* **107**, 1977, pp. 89–97.
5. Ibid., R. P. Breault, "Problems and Techniques in Stray Radiation Suppression."
6. R. P. Breault, "Vane Structure Design Trade-Off and Performance Analysis," *Stray Light and Contamination in Optical Systems*, *Proc. SPIE* **967**, 1988.
7. Ibid., S. R. Lange, R. P. Breault, and A. W. Greynolds, "APART, A First-Order Deterministic Stray Radiation Analysis Program."
8. Ibid., R. P. Breault, A. W. Greynolds, and M. A. Gauvin, "Stray Light Analysis with APART/PADE, Version 8.7."
9. Ibid., R. P. Breault, "Problems and Techniques in Stray Radiation Suppression."
10. Ibid., S. R. Lange, R. P. Breault, and A. W. Greynolds, "APART, A First-Order Deterministic Stray Radiation Analysis Program."
11. This section is a part of a more complete mathematical description of the process which can be found in R. P. Breault, "Vane Structure Design Trade-Off and Performance Analysis." (Ibid.)
12. R. P. Breault, "Vane Structure Design Trade-Off and Performance Analysis" (Ibid.) contains a more detailed description of the process.
13. Ibid., R. P. Breault, "Problems and Techniques in Stray Radiation Suppression."
14. R. P. Breault, "Suppression of Scattered Light," Ph.D. dissertation, University of Arizona, 1979.
15. Ibid., R. P. Breault, "Problems and Techniques in Stray Radiation Suppression."
16. J. Gunderson, "Goniometric Reflection Scattering Measurements and Techniques at 10.6 Micrometers," M. S. thesis, University of Arizona, 1977.
17. Ibid., R. P. Breault, A. W. Greynolds, and S. R. Lange, "APART/PADE Version 7: A Deterministic Computer Program Used to Calculate Scattered and Diffracted Energy."
18. Ibid., R. P. Breault, A. W. Greynolds, and M. A. Gauvin, "Stray Light Analysis with APART/PADE, Version 8.7."
19. E. R. Friere and D. L. Skelton, "Use of Specular Black Coatings in Well-Baffled Optical Systems," *Stray Radiation V*, Robert Breault (ed.), *Proc. SPIE* **675**, 1986, pp. 126–132.

20. A. W. Greynolds, computer code ASAP, Breault Research Organization, Inc., 1992.
21. G. L. Peterson, S. C. Johnston, and J. Thomas, "Specular Baffles," *Stray Radiation in Optical Systems II*, Robert P. Breault (ed.), *Proc. SPIE* **1753**, 1992.
22. P. R. Spyak, "A Cryogenic Scatterometer and Scatter from Particulate Contaminants on Mirrors," Ph.D. dissertation, University of Arizona, 1990.
23. "Product Cleanliness Levels and Contamination Control Program," *MIL-STD-1246A*, Dept. of Defense, Global Engineering Documents, Santa Ana, Calif., 18 Aug. 1967.
24. R. Young, "Low-Scatter Mirror Degradation by Particle Contamination," *Optical Engineering* **15**, no. 6, Nov.–Dec. 1976.
25. J. E. Harvey, "Light-Scattering Characteristics of Optical Surfaces," Ph.D. dissertation, University of Arizona, 1976.
26. P. R. Spyak and W. L. Wolfe, "Scatter from Particulate Contaminated Mirrors," *Optical Engineering* **31**, no. 8, Aug. 1992, pp. 1746–1784.
27. M. G. Dittman, "Contamination Scatter Functions for Stray Light Analysis," *Optical System Contamination: Effects, Measurements, and Control VII*, P. T. Chen and O. M. Uy (eds.), *Proc. SPIE*, **4774**, 2002, pp. 99–110.
28. Contact Dierdre Dykeman at Rome Air Development Center for the best way to access the results of this work.
29. H. E. Bennett, "Reduction of Stray Light from Optical Components," *Stray Light Problems in Optical Systems*, J. D. Lytle and H. Morrow (eds.), *Proc. SPIE* **107**, 1977, pp. 24–33.
30. Name coined by D. Rock, Hughes Aircraft Co., El Segundo, Calif.
31. Name coined by R. P. Breault; this is the name most often used in an APART stray light analysis.
32. *Ibid.*, J. E. Harvey, "Light-Scattering Characteristics of Optical Surfaces."
33. J. C. Stover, *Optical Scattering: Measurement and Analysis*, New York: McGraw-Hill, Inc., New York, 1990.
34. *Ibid.*, J. C. Stover, *Optical Scattering: Measurement and Analysis*, p. 69.
35. *Ibid.*, J. Gunderson, "Goniometric Reflection Scattering Measurements and Techniques at 10.6 Micrometers."

This page intentionally left blank.

DO NOT DUPLICATE

THERMAL COMPENSATION TECHNIQUES

Philip J. Rogers and Michael Roberts

*Pilkington Optronics
Wales, United Kingdom*

8.1 GLOSSARY

c	surface curvature of an optical element
D	diameter
FN	F-number or focal ratio
f	paraxial focal length
G	thermo-optical constant (normalized thermal change of OPD)
h	(subscript) signifies “pertaining to the optic housing”
i	(subscript) number of a specific optical element
j	signifies a number of optical elements
K	Kelvin
k	thermal conductivity
n	refractive index
OPD	optical path difference
T	temperature
t	thickness
V	Abbe number of a refracting optical material
ν	spatial frequency
α	linear coefficient of thermal expansion
γ	thermal glass constant (normalized thermal change of optical power)
Δ	small, finite change
λ	wavelength
δ	infinitely small change of a parameter
ϕ	optical power (reciprocal of focal length)

8.2 INTRODUCTION

In the following, the thermal effects for which compensation is required are taken to be those that affect the focus and image scale of an optical system. Methods for quantifying and offsetting these effects were described some time ago,¹ similar information being provided by several other authorities.^{2,3,4} The thermal compensation techniques described in this chapter, with the exception of intrinsic athermalization, involve either mechanical movement of one or more parts of the optical system, or compensation achieved solely by choice of optical materials. Except in Sec. 8.5 titled “Effect of Thermal Gradients,” a homogeneous temperature change of all parts of the optical system is assumed.

Most optical materials undergo a change of refractive index n with temperature T , conveniently quoted as a rate of change $\delta n/\delta T$. The usual values of n and $\delta n/\delta T$ given for a material (and assumed in this chapter unless stated otherwise) are those relative to the surrounding air rather than the absolute values with respect to vacuum. Air has a $\delta n/\delta T$ of -1×10^{-6} at $T = 288$ K and 1 atmosphere air pressure for wavelengths between 0.25 and 20 μm .⁵ allowance for this must be made when a lens operates in vacuum or in an enclosed space where the number of air molecules per unit volume does not change with temperature. The absolute $\delta n/\delta T$ of an optical material can be found from

$$\left(\frac{\delta n}{\delta T}\right)_{\text{abs}} = n_{\text{air}} \left(\frac{\delta n}{\delta T}\right) + n \left(\frac{\delta n}{\delta T}\right)_{\text{air}} \quad (1)$$

where the value of n_{air} is approximately 1.0.

8.3 HOMOGENEOUS THERMAL EFFECTS

Thermal Focus Shift of a Simple Lens

The rate of change of the power ϕ (reciprocal of the focal length f) of an optical element with temperature T can be obtained by differentiating the thin lens power equation $\phi = c(n - 1)$, where c is the total surface curvature of the element. For a linear thermal expansion coefficient α of the material from which the element is formed this gives

$$f = \frac{1}{\phi} \quad \frac{\delta f}{\delta T} = -\frac{1}{\phi^2} \frac{\delta \phi}{\delta T}$$

$$\frac{\delta \phi}{\delta T} = +\phi \left(\frac{\delta n/\delta T}{n-1} - \alpha \right) \quad (2)$$

Therefore

$$\frac{\delta f}{\delta T} = -f \left(\frac{\delta n/\delta T}{n-1} - \alpha \right) \quad (3)$$

The material-dependent factor inside the parenthesis in Eqs. (2) and (3) is known as the thermal “glass” constant (γ) and represents the thermal power change due to an optical material normalized to unit ϕ and unit change of T . Tables 1 to 3 give γ values for a selected number of visual and

TABLE 1 Optical and Thermal Data for a Number of Visual Waveband Materials

Schott Glass	Optical Plastic Type	Refractive Index, n_e^*	Abbe Number, V_e^\dagger	Thermal Constant, $\gamma(\times 10^6)^\ddagger$	Thermo-Optical Constant, $G(\times 10^6)^\ddagger$	Thermal Conductivity, $k(W \cdot m^{-1} \cdot K^{-1})$
	FK52	1.487	81.4	-27	+1	0.9
	FK5	1.489	70.2	-11	+4	0.9
	BK7	1.519	64.0	-1	+7	1.1
	PSK53A	1.622	63.2	-13	+4	—
	SK5	1.591	61.0	+1	+7	1.0
	BaLKN3	1.521	60.0	-3	+7	1.0
	BaK2	1.542	59.4	-5	+6	—
	SK4	1.615	58.4	-2	+7	0.9
	LaK9	1.694	54.5	-1	+8	0.9
	KzFSN4	1.617	44.1	+4	+8	0.8
	LF5	1.585	40.6	-6	+7	0.9
	BaSF51	1.728	37.9	+8	+14	0.7
	LaFN7	1.755	34.7	+6	+12	0.8
	SF5	1.678	32.0	0	+11	—
	SFN64	1.711	30.1	-4	+9	—
	SF6	1.813	25.2	+6	+18	0.7
	Acrylic [§]	1.497	57	-279	-71	0.2
	Polycarbonate [§]	1.590	30	-247	-68	0.2

*At $\lambda = 546$ nm.†Defined as $(n_{546} - 1)/(n_{480} - n_{644})$.‡At $\lambda = 546$ nm and $T = 20$ °C.§Values (except conductivity) from Waxier et al. *Appl. Opt.* 18:102 (1979).

infrared materials along with the relevant V value (Abbe number) and other data. The much higher level of γ for infrared as opposed to glass optical materials indicates that thermal defocus (focus shift) is generally a much more serious problem in the infrared wavebands. The actual value of γ varies with both wavelength and temperature due to variations in the value of $\delta n/\delta T$ and α . In general, this is unlikely to cause major problems unless a wide wavelength or temperature range is

TABLE 2 Optical and Thermal Data for Selected 3- to 5- μ m Waveband Infrared Materials

Optical Material	Refractive Index, $n_{4\mu}$	Abbe Number, $V_{3-5\mu}$	Thermal "Glass" Constant, γ	Thermal Thermo-Optical Constant, G	Conductivity, $k(W \cdot m^{-1} \cdot K^{-1})$
Silicon	3.43	2.4×10^2	$+6.3 \times 10^{-5}$	$+1.7 \times 10^{-4}$	1.5×10^2
KRS5*	2.38	2.3×10^2	-2.3×10^{-4}	-1.5×10^{-4}	5.0×10^{-1}
AMTIR1†	2.51	1.9×10^2	$+3.9 \times 10^{-5}$	$+9.5 \times 10^{-5}$	3.0×10^{-1}
Zinc selenide	2.43	1.8×10^2	$+3.6 \times 10^{-5}$	$+7.3 \times 10^{-5}$	1.8×10^1
Arsenic trisulfide	2.41	1.6×10^2	-1.9×10^{-5}	$+3.4 \times 10^{-5}$	1.7×10^{-1}
Zinc sulfide	2.25	1.1×10^2	$+2.6 \times 10^{-5}$	$+5.2 \times 10^{-5}$	1.7×10^1
Germanium	4.02	1.0×10^2	$+1.3 \times 10^{-4}$	$+4.2 \times 10^{-4}$	5.9×10^1
Calcium fluoride	1.41	2.2×10^1	-5.1×10^{-5}	-1×10^{-6}	9
Magnesium oxide	1.67	1.2×10^1	$+1.9 \times 10^{-5}$	$+2.6 \times 10^{-5}$	4.4×10^1

*Thallium bromo-iodide.

†Ge/As/Se chalcogenide from Amorphous Materials Inc.

TABLE 3 Optical and Thermal Data for Selected 8- to 12- μm Waveband Infrared Materials

Optical Material	Refractive Index, $n_{10\mu}$	Abbe Number, $V_{8-12\mu}$	Thermal "Glass" Constant, γ	Thermo-Optical Constant, G	Thermal Conductivity, $k(\text{W} \cdot \text{m}^{-1} \cdot \text{K}^{-1})$
Germanium	4.00	8.6×10^2	$+1.2 \times 10^{-4}$	$+4.1 \times 10^{-4}$	5.9×10^1
Cesium iodide	1.74	2.3×10^2	-1.7×10^{-4}	-5.3×10^{-5}	1
Cadmium telluride	2.68	1.7×10^2	$+5.3 \times 10^{-5}$	$+1.1 \times 10^{-4}$	6
KRS5	2.37	1.7×10^2	-2.3×10^{-4}	-1.6×10^{-4}	5.0×10^{-1}
AMTIR1	2.50	1.1×10^2	$+3.6 \times 10^{-5}$	$+9.0 \times 10^{-5}$	3.0×10^{-1}
Gallium arsenide	3.28	1.1×10^2	$+7.6 \times 10^{-5}$	$+2.0 \times 10^{-4}$	4.8×10^1
Zinc selenide	2.41	5.8×10^1	$+3.6 \times 10^{-5}$	$+7.2 \times 10^{-5}$	1.8×10^1
Zinc sulfide	2.20	2.3×10^1	$+2.6 \times 10^{-5}$	$+5.0 \times 10^{-5}$	1.7×10^1
Sodium chloride	1.49	1.9×10^1	-9.5×10^{-5}	-3×10^{-6}	6

being considered.⁶ Thermal defocus results not only from a change of optical power but also from the thermal expansion coefficient α_h of the housing. Equation (3) can be modified to allow for the effect of the latter:

$$\text{Single thin lens:} \quad \Delta f = -f(\gamma + \alpha_h)\Delta T \quad (4)$$

$$j \text{ thin lenses in contact:} \quad \Delta f = -f \left[f \sum_{i=1}^j (\gamma_i \phi_i) + \alpha_h \right] \Delta T \quad (5)$$

Thermal Defocus of a Compound Optical Construction

Consider a homogeneous temperature change in an optical system that comprises two thin-lens groups separated from each other, the normalized thermal power change being the same in each lens group. Taking the thermal defocus calculated from Eq. (4) as unity, then that due to a compound optic comprising two separated components and of the same overall power can be estimated from Fig. 1.⁷ The latter shows scaling of thermal defocus with respect to a simple thin lens, relative to front lens/image plane distance (overall length) for three different positions of the second lens group. The graph is divided into three basic lens constructions distinguished from each other by overall optical length and/or the sign of the power of the front lens group.

Figure 1 assumes germanium optics in an aluminum housing but change of either material, while altering values slightly, has no effect on the following two conclusions:

1. Telephoto/inverted telephoto constructions always give more (and Petzval lenses always give less) thermal defocus than an equivalent simple lens.
2. Thermal defocus reduces as the second lens group is moved toward the image plane, irrespective of lens construction: the efficacy of this procedure is limited, however, by the increased imbalance of optical powers between the groups.

The thermal defocus scaling technique could be extended to cover an optic comprising more than two lens groups. This extension has been carried out for the Cooke triplet construction⁸ and for a series of separated thin lenses,⁶ although only for the case of a zero-expansion housing.

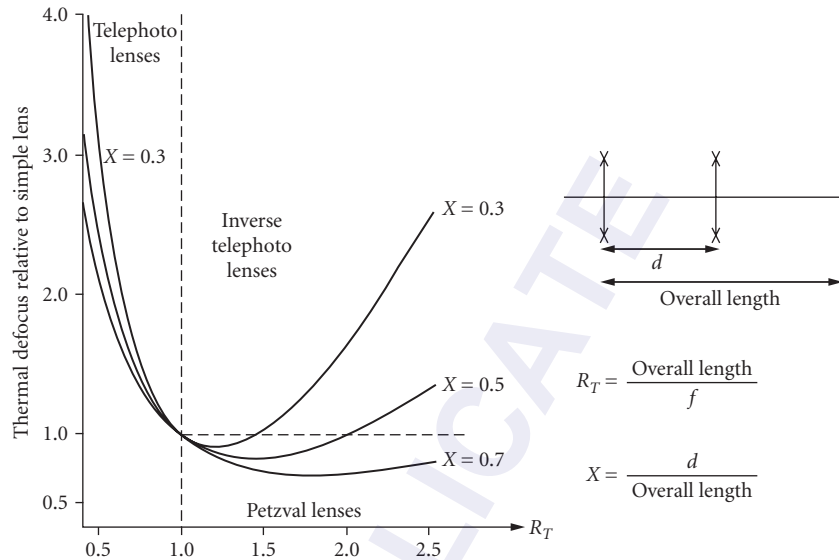


FIGURE 1 Effect of compound lens construction on thermal defocus. (From Rogers.⁷)

8.4 TOLERABLE HOMOGENEOUS TEMPERATURE CHANGE (NO COMPENSATION)

Diffraction-Limited Optic

Equation (5) can be used to establish the temperature change ΔT that will result in a quarter-wave of longitudinal thermal defocus, a reasonable limit for a simple optic that is nominally diffraction-limited. Given an optic of diameter D and focal ratio FN imaging at a mean wavelength of λ :

$$\text{Diffraction-based depth of focus:} \quad \Delta f = \pm 2\lambda(\text{FN})^2 \quad (6)$$

$$\text{Combining Eqs. (5) and (6):} \quad \Delta T = \pm \frac{2\lambda(\text{FN})}{D \left[f \sum_{i=1}^i (\gamma_i \phi_i) + \alpha_h \right]} \quad (7)$$

Figure 2 gives ΔT against D results for a simple 8- to 12- μm bandwidth germanium optic in an aluminum housing;⁹ the curves illustrate the small temperature change that can be tolerated in germanium optics before focus compensation is required. Partial avoidance of this particular problem may be achieved by the replacement of germanium by other infrared optical materials having lower values of γ : this may also be desirable to reduce high-temperature absorption but generally leads to much greater optical complexity.

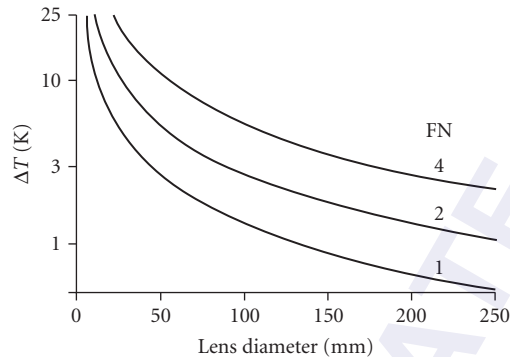


FIGURE 2 Tolerable temperature change for a simple germanium infrared lens. (From Rogers.⁹)

Nondiffraction-Limited Optic

The depth of focus of an optic having a nominal performance far from the diffraction limit is a function of the residual aberration level and balance in the optic as well as its first-order parameters. An estimate related to a cutoff spatial frequency ν that gives a reasonable approximation in many cases can be obtained¹⁰ from

$$\text{Approximate depth of focus:} \quad \Delta f = \pm \frac{(\text{FN})}{\nu} \quad (8)$$

$$\text{Combining Eqs. (5) and (8):} \quad \Delta T = \pm \left\{ D\nu \left[f \sum_{i=1}^j (\gamma_i \phi_i) + \alpha_h \right] \right\}^{-1} \quad (9)$$

Notice that, given the approximation of this method, the value of ν can be determined by extending a straight line MTF from 1.0 response at zero spatial frequency, through the MTF point of interest, to the intersection of the line with the spatial frequency axis.

8.5 EFFECT OF THERMAL GRADIENTS

The previous sections assume a homogeneous temperature change in all parts of the optical system: in situations where steady-state or transient temperature gradients exist, further consideration is required.¹

Allowance for the effect of steady-state longitudinal gradients can be made by applying a different value of T to each lens group and an average local temperature to each portion of the housing that separates two adjacent lens groups. Transient longitudinal gradients are a more difficult problem and, if severe, may require individual athermalization of each lens group in its own housing domain.

Steady-state or transient radial thermal gradients cause at least a shift of focus position, with the possible addition of a change of aberration correction. A localized radial temperature difference of

ΔT through the thickness t of a plane-parallel plate will cause a deviation of a ray of light¹¹ that can be quantified as an optical path difference (OPD):

$$\text{OPD} = t[\alpha(n-1) + \delta n / \delta T] \Delta T \quad (10)$$

The expression in the square bracket is often referred to as the thermo-optical constant G and is an approximate measure of the sensitivity of an optical material to radial gradients. More thorough analysis of the effects produced by radial thermal gradients includes computation of thermally induced stress and consequent anisotropic change of refractive index: in some cases, this may be a significant factor in image degradation.¹²⁻¹⁴

Tables 1 to 3 give values of G for the selected optical materials. Also tabulated is the thermal conductivity k , as in many cases G/k is a more appropriate measure of sensitivity given the greater ability of high-conductivity materials to achieve thermal equilibrium.

8.6 INTRINSIC ATHERMALIZATION

The need for athermalization can be avoided or minimized for some applications by employing optical power and mounting techniques that are inherently insensitive to temperature change. A concave spherical mirror fabricated from the same material that separates the mirror from its focal plane (e.g., an aluminum mirror in an aluminum housing) is in effect “self-athermalized” for a homogeneous distribution of temperature. The optical performance of a single spherical mirror is limited, but the above principle applies for more complex all-reflective optical constructions employing conic or other aspheric surface forms. A glass spherical mirror, although not thermally matched to an aluminum mounting, may be used as part of a self-athermalized catadioptric afocal in the infrared, a germanium Mangin being used in this case as a secondary mirror lens.¹⁵ The high thermal power change of the negatively powered lens in the germanium Mangin, used in double-pass, compensates for the thermal defocus due to the glass primary, the housing, and the remaining germanium optics in the afocal—Fig. 3.¹⁶

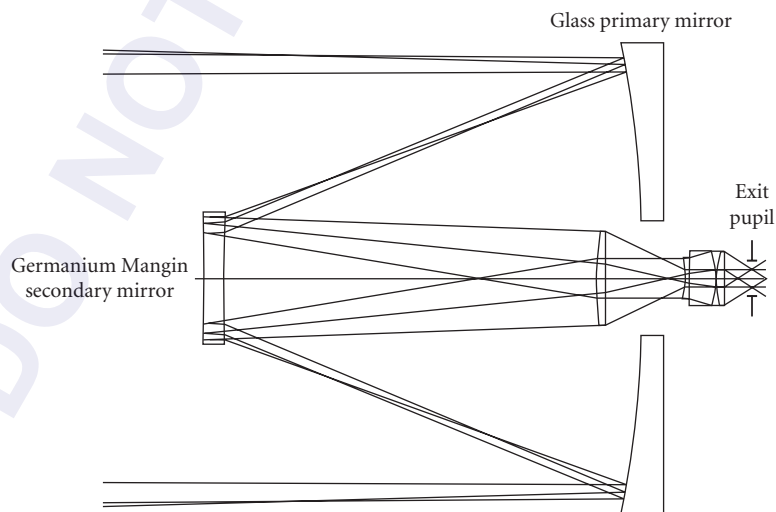


FIGURE 3 High magnification self-athermalized catadioptric afocal. (From Norrie.¹⁶)

An alternative approach to the above is to use glass-ceramic mirrors within a nickel-iron alloy housing, since they can have thermal expansion coefficients approaching zero. A major advantage of this approach is its insensitivity to thermal gradients.

8.7 MECHANICAL ATHERMALIZATION

General

Mechanical athermalization essentially involves some agency moving one or more lens elements by an amount that compensates for thermal defocus—a simple manual option being to use an existing focus mechanism. Automatic methods are, however, preferable in many cases and can be divided into passive or active. Passive athermalization employs an agency, often involving materials (including liquids) with abnormal thermal expansion coefficients, to maintain focus without any powered drive mechanism being required. Automatic active athermalization involves the computation of focus compensation algorithms that are stored (usually electronically) and implemented by a powered device such as an electric motor. The following sections refer to a number of passive and active athermalization methods, although the list is by no means exhaustive.

Passive Mechanical Athermalization

The principal advantages of passive thermal compensation methods are their relative simplicity and potential reliability. Disadvantages are their inadequate response to transient temperature gradients and, generally, lack of adjustment to allow for errors or unforeseen circumstances. Passive methods are ideal in glass optics¹⁷ where thermal effects are low, although here it is not too difficult to achieve optical athermalization (see later under “Optical Athermalization”) except where very low secondary spectrum is required. In the infrared wavebands, where thermal effects are much greater due to the nature of the optical materials, it is difficult to achieve simple passive mechanical athermalization due to the large refocusing movement required, typically 1.5×10^{-4} per unit focal length per Kelvin for an aluminum-housed germanium optic. An exception to the above is the combination of silicon and germanium in 3- to 5- μm optics, where thermal defocus results largely from the expansion of the housing. In this case, use of more than one nonmetallic housing material can result in an athermalized optic, even one having two fields of view¹⁸—Fig. 4.

For infrared cases other than the above, the options are either to provide a mechanism that modifies mechanical expansion effects or to reduce the required refocusing movement by optical means.

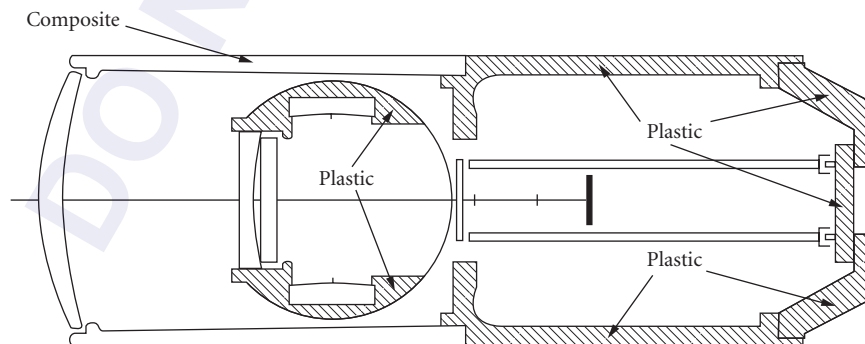


FIGURE 4 Part composite/part plastic mounting structure used for athermalization of a 3- to 5- μm IR optic. (From Garcia-Nuñez and Michika.¹⁸)

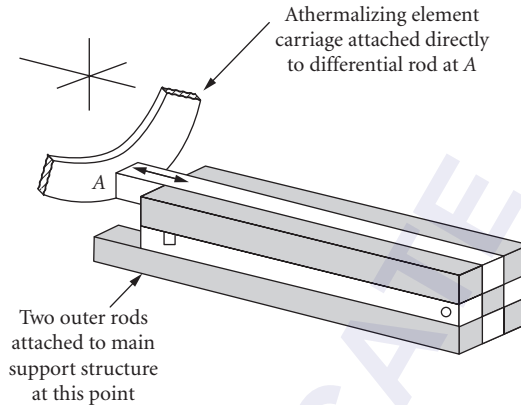


FIGURE 5 Passive mechanical thermal compensation using differential expansion rods. (From Povey.¹⁹)

Examples of the former include¹⁹ a series of linked rods of alternatively high and low expansion coefficient—Fig. 5—and a hydraulic method where the fluid contained in a large-volume reservoir expands into a small-bore cylinder—Fig. 6.

An interesting alternative employs shape-memory metal²⁰ to provide a large movement over a relatively small temperature range.¹⁹ Another alternative is to employ a geodetic arrangement: in this method¹⁹ an athermalizing adjustment of, for example, the separation between primary and secondary mirrors in a catadioptric, is produced by differing expansion coefficients of the primary mirror

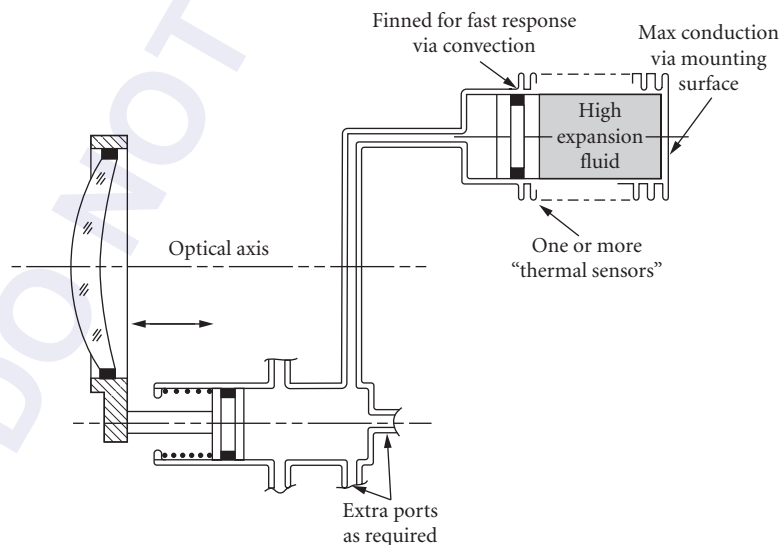


FIGURE 6 Passive mechanical thermal compensation using high-expansion-fluid thermal sensors. (From Povey.¹⁹)

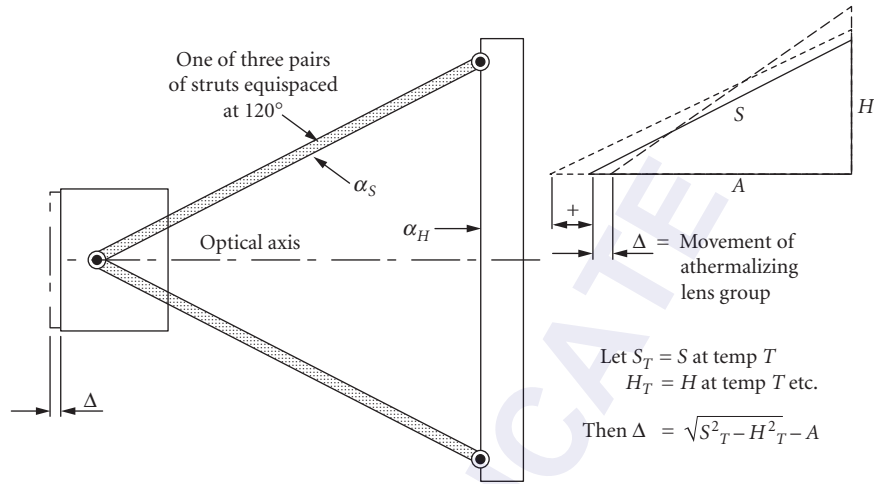


FIGURE 7 Geodetic support structure for positive or negative thermal compensation movement. (From Povey.¹⁹)

mount and the secondary mirror struts—Fig. 7. Where none of the above methods are desirable, the option to reduce the necessary movement may be the only alternative. This may be achieved by an optical layout⁹ configured such that the required athermalizing movement is reduced typically by a factor of four, but at the expense of somewhat greater optical complexity—Fig. 8.

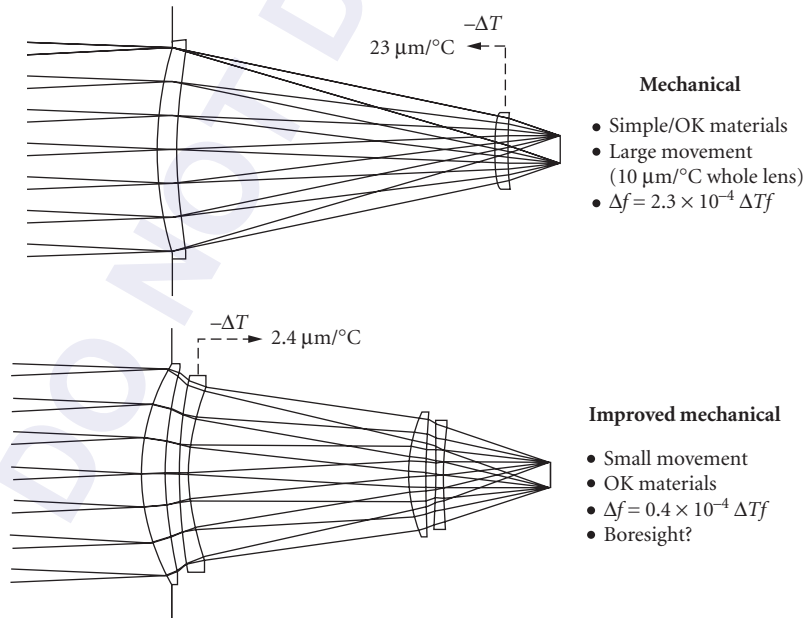


FIGURE 8 Alternative optical configurations for mechanically athermalized forward looking infrared (FLIR) systems. (From Rogers.⁹)

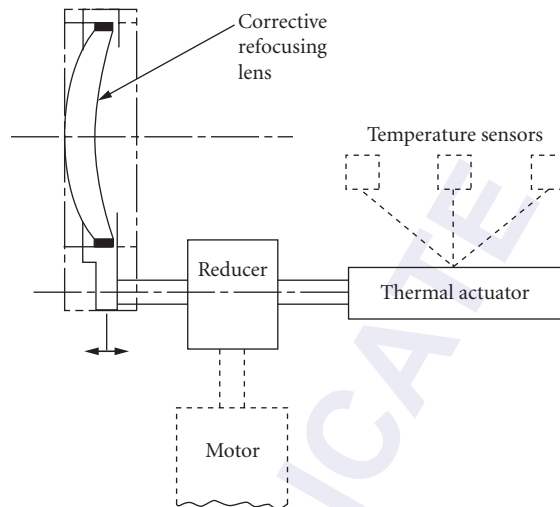


FIGURE 9 Active electromechanical athermalization—schematic.

Active Mechanical Athermalization

Active mechanical athermalization in its simplest form can be manual adjustment of a lens element or group for refocusing. For more complex optics, such as multi-field-of-view, a procedure can be specified for manual (or motorized) adjustment of several lens elements to maintain focus over a range of magnifications and temperatures.^{21,22} Where automatic athermalization is required, a method can be employed that uses a combination of electronics and mechanics—Fig. 9. One or more temperature sensors located along the body of the optic feed their signals into an algorithm that calculates the required movement of a compensating lens and then initiates the motion. For simplicity, the compensating lens may be that which already provides close-distance focusing, thus requiring only an increase in the range of movement for athermalization. The location of sensors is especially important for infrared optics and should be dependent on the thermal sensitivity variations within the optical system.

Active electromechanical thermal compensation is particularly suitable where transient longitudinal temperature gradients are expected and for multi-field-of-view optics where thermal defocus is dependent on field-of-view setting. The algorithm required for elimination of the effects of a combination of both of the above is complex, but compensation may be accomplished by a single mechanical motion.²³

A single motion does not, however, guarantee athermalization of image scale, in which case more than one compensatory movement may be required. Two motion athermalization in a zoom or dual-field-of-view infrared telescope can take advantage of the existing mechanisms required for field-of-view change. Also, by utilizing internal lens elements, problems associated with hermetically sealing an external focusing lens element can be avoided.^{24–28} In order to maintain stability of aberration correction in infrared zoom telescopes, particularly those having a large zoom range, three-motion athermalization has been proposed.^{29,30}

Active/Passive Athermalization

An improvement over simple manual active athermalization is to include partial passive athermalization. This is best suited to systems that already contain axially moving lens components, for example, a dual-field-of-view infrared telescope³¹—Fig. 10.³² Here the majority of the athermalization is provided

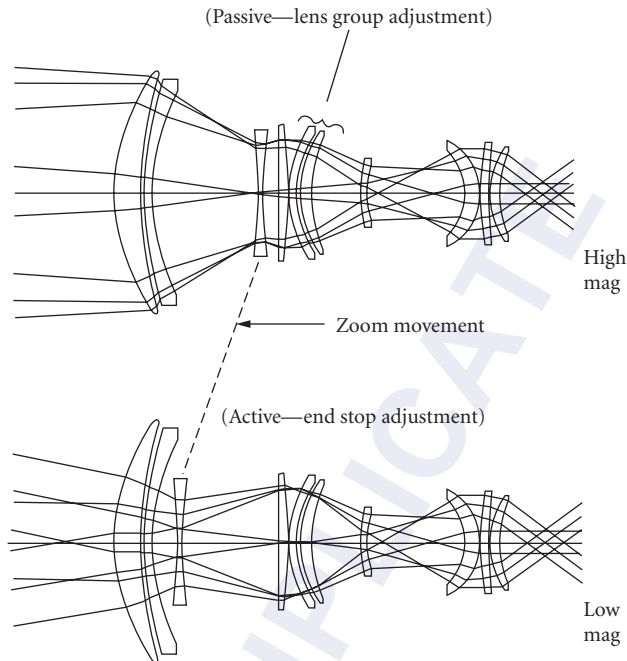


FIGURE 10 Part active, part passive mechanical athermalization. (From Roberts.³²)

by a mechanically passive device that adjusts the position of the rear lens group in the objective. The residual focus error is then corrected by small manual adjustments to the magnification change element. This technique can minimize the change of image scale and aberrations with temperature. A potential problem, however, is the subjective nature of best-focus determination.

Athermalization by Image Processing

Athermalization by image processing is suitable for some applications. A range of automatic focusing techniques exists but, while this approach has the advantage of not requiring temperature sensors, it does suffer the potential disadvantage of misinterpretation of image information.

8.8 OPTICAL ATHERMALIZATION

General

Athermalization of the focus position of an optical system by choice of refractive materials has been described extensively in the literature.^{33–42} The requirements of overall optical power, achromatism, and athermalism demand that three conditions be satisfied for j thin lens elements in contact:

Power:
$$\sum_{i=1}^j \phi_i = \phi \quad (11a)$$

TABLE 4 Unity Focal Length Athermal Two-Material Achromatic Combinations

Material Type	Material Combination	Total Curvatures	Secondary Spectrum [‡]	Petzval Sum	Normalized Mass
Optical glasses	BaLKN3 + KzFSN4	+7.24/−4.49	3.6×10^{-4}	0.77	2.1
	BaK2 + LaFN7	+4.43/−1.86	4.8×10^{-4}	0.76	1.4
	FK5 + LF5	+4.85/−2.34	4.8×10^{-4}	0.73	1.3
	PSK53A + SFN64	+3.06/−1.28	5.0×10^{-4}	0.65	1.0
	BaLKN3 + BaSF51	+5.21/−2.35	5.2×10^{-4}	0.79	1.5
Stabilized optical glasses	SK4 + KzFSN4	+7.30/−5.64	1.8×10^{-4}	0.62	2.0 [‡]
	SK5 + SF5 [*]	+3.97/−1.99	10.6×10^{-4}	0.67	1.0 [‡]
3- to 5- μm Materials	As ₂ S ₃ + MgO	+0.77/−0.12	8.6×10^{-4}	0.40	0.8 [§]

^{*}Thermally invariant housing, all others aluminum.

[†]Over wavebands of 480 to 644 nm, 546 to 852 nm, and 3 to 5 μm respectively.

[‡]Relative to SK5/SF5 solution.

[§]Relative to lowest value in Table 5.

Source: Rogers.⁴³

$$\text{Achromatism: } \sum_{i=1}^j \frac{\phi_i}{V_i} = 0 \quad (11b)$$

$$\text{Athermalism: } \sum_{i=1}^j (\gamma_i \phi_i) + \phi \alpha_h = 0 \quad (11c)$$

The presence of three conditions implies the need for three different materials in order to obtain an exact solution. It is possible, however, to find achromatic combinations of two materials that are also athermal, provided that a simple condition is satisfied:⁴¹

$$V_1(\gamma_1 + \alpha_h) = V_2(\gamma_2 + \alpha_h) \quad (12)$$

Suitable combinations for thin-lens athermal achromats can be found by plotting a range of materials on a graph of γV against V ; the slope of the line joining a chosen pair representing the required thermal expansion coefficient of the housing.⁶

A number of approximately athermal optical glass achromats exist of which those listed in Table 4⁴³—with the exception of the last entry—represent examples with low to moderate secondary spectrum over the visible to near infrared waveband. The data given for these achromats are lens element total curvatures for unity focal length; secondary spectrum (second-order color); thin-lens Petzval sum; and an approximate indication of mass, normalized to the lowest value. The pairing of radiation-stabilized versions of SK5 and SF5, both of which have a low value of γ makes a good choice for athermalized space optics in a temperature-invariant mount.⁶

In the infrared wavebands the options are far more limited: at least one 3- to 5- μm waveband two-material athermal combination exists, namely, arsenic trisulfide and magnesium oxide, but there is currently no realistic pairing of materials in the 8- to 12- μm band.

Athermal Laser Beam Expanders

Many more two-material athermal combinations exist if the requirement for achromatism [Eq. (11b)] is removed. This is the situation that occurs with a (preferably) galilean laser beam expander, although here the two lens materials are separated.⁴⁴ From Eq. (4), making the thermal defocus Δf values equal

and opposite for the two lenses leads to a value for magnification at which two given materials in a specific housing material will provide an athermal beam expander (for a homogeneous temperature distribution):

$$\text{Magnification} = \frac{\gamma_1 + \alpha_h}{\gamma_2 + \alpha_h} \quad (13)$$

Three-Material Athermal Solutions

Graphical methods have been described that allow investigation of preferred three-material athermalized achromatic solutions.⁴¹ An alternative method is the systematic evaluation of all possible combinations of three materials selected from a short list, each combination being allocated a risk factor dependent on material characteristics and solution sensitivity.⁹ The optical powers of the three in-contact thin-lens elements are determined by solving Eq. (11a to c) which give for a unity focal length:

$$a = \frac{V_1 V_2 - V_2 V_3}{V_1 V_3 - V_2 V_3}, \phi_3 = \frac{(1-b)\gamma_1 + b\gamma_2 + \alpha_h}{(1-a)\gamma_1 + a\gamma_2 + \gamma_3} \quad (14a)$$

$$b = \frac{V_2}{V_2 - V_1}, \phi_2 = b - a\phi_3 \quad (14b)$$

$$\phi_1 = 1 - (\phi_2 + \phi_3) \quad (14c)$$

Tables 5 and 6 give a selection of lower-risk three-material solutions, in approximate order of increasing risk, for 3- to 5- and 8 to 12- μm infrared combinations, respectively. The data given are similar to those in Table 4, but the housing is assumed to be aluminum in all cases. Note that these tables are intended as a guide only and are based on currently available material data.

Athermalization of Separated Components

In many ways, thermal defocus and thermal change of focal length are analogous to longitudinal and lateral chromatic aberration, having the same first-order dependencies. For this reason it has been suggested that a thermal Abbe number, defined as $\gamma^{-1.3}$ be used to replace the chromatic

TABLE 5 Unity Focal Length Athermal Three-Material Achromatic Combinations for the 3- to 5- μm Waveband

Material Combination	Total Curvatures	Petzval Sum	Normalized Mass
Si + Ge + ZnS	+0.72/-0.36/+0.27	0.39	1.3
ZnSe + Ge + MgO	+1.16/-0.21/-0.06	0.51	1.8
Si + Ge + KRS5	+0.69/-0.26/+0.08	0.34	1.0
[ZnS + MgO + Ge]	+1.28/-0.17/-0.16	0.52	1.5
AMTIR1 + Ge + Si	+0.56/-0.32/+0.46	0.42	1.4
Si + MgO + KRS5	+0.31/-0.08/+0.22	0.31	1.1
ZnSe + ZnS + Ge	+1.80/-0.69/-0.23	0.50	2.8
Si + CaF ₂ + KRS5	+0.32/-0.25/+0.24	0.29	1.1

[] Low residual high-order chromatic aberration.

TABLE 6 Unity Focal Length Athermal Three-Material Achromatic Combinations for the 8- to 12- μm Waveband

Material Combination	Total Curvatures	Petzval Sum	Normalized Mass
KRS5 + ZnSe + Ge	+0.34/−0.15/+0.25	0.30	1.0
ZnSe + ZnS + Ge	+2.05/−0.92/−0.26	0.50	2.5
GaAs + ZnS + KRS5	+0.38/−0.20/+0.26	0.31	1.0
AMTIR1 + ZnS + Ge	+1.42/−0.35/−0.24	0.48	1.3
{CdTe + ZnSe + KRS5}	+0.72/−0.37/+0.22	0.37	1.5
GaAs + ZnSe + KRS5	+0.68/−0.71/+0.33	0.25	1.8
[AMTIR1 + ZnSe + KRS5]	+1.19/−0.72/+0.16	0.39	1.9
[CsI + NaCl + GaAs]	+0.68/−0.32/+0.29	0.38	1.1

{ } Very high transmission.

[] Low residual high-order chromatic aberration.

Abbe number (V value) in the usual chromatic aberration equations. Thermal expansion of the housing—obviously not present in chromatic calculations—does, however, complicate the situation a little.

In equations thus far, Δf has meant both thermal defocus and focal length change, as numerically these are the same for a thin lens. For separated components, rules similar to those for chromatic aberration apply, for example, two separated thin-lens groups—such as those described by Fig. 1—must be individually athermalized if both types of thermal “aberration” are to be corrected simultaneously. More complex optics (for example, multistage) may have transfer of thermal aberration between constituent lens groups but may still be corrected simultaneously for thermal focus shift and focal length change as a whole. This procedure can, however, lead to one lens group requiring excessive optical powers in order to achieve full overall correction—transient longitudinal thermal gradients may also cause problems.

Use of Diffractive Optics in Optical Athermalization

The term “hybrid optic” is generally used to signify a combination of refractive and diffractive means in an optical element. The diffractive part of the hybrid is usually a transmission hologram which for high efficiency would be of surface relief form, the surface structure being machined or etched onto the refractive surface.⁴⁵ The diffractive surface acts as a powered diffraction grating, producing large amounts of chromatic aberration which could be employed in an optic where a lightweight optically athermalized combination of two materials could be chosen without regard to achromatism: residual chromatic aberration could then be corrected by the hologram.⁴⁶

8.9 REFERENCES

1. J. W. Perry, *Proc. Phys. Soc.*, vol. 55, 1943, pp. 257–285.
2. J. Johnson and J. H. Jeffree, U.K. Patent No. 561 503, 1942 U.K. priority.
3. D. S. Grey, *J. Opt. Soc. Am.*, vol. 38, 1948, pp. 542–546.
4. D. S. Volosov, *Opt. Spectrosc.*, U.S.S.R., vol. 4, pp. 663–669 and pp. 772–778, vol. 5, 1958, pp. 191–199.
5. R. Penndorf, *J. Opt. Soc. Am.*, vol. 47, 1957, pp. 176–182.
6. H. Köhler and F. Strähle, *Space Optics*, B. J. Thompson and R. R. Shannon (eds.), National Academy of Sciences, 1974, pp. 116–153.
7. P. J. Rogers, *SPIE Critical Reviews*, vol. CR38, 1991, pp. 69–94.
8. L. R. Estelle, *SPIE*, vol. 237, 1980, pp. 392–401.

9. P. J. Rogers, *SPIE*, vol. 1354, 1990, pp. 742–751.
10. M. Laikin, *Lens Design*, Marcel Dekker, New York, 1991, p. 28.
11. G. G. Slyusarev, *Opt. Spectrosc.*, U.S.S.R., vol. 6, 1959, pp. 134–138.
12. W. H. Turner, *Optical Sciences Center*, vol. 4, University of Arizona, Tucson, Arizona, 1970, pp. 123–125.
13. V. M. Mit'kin and O. S. Shchavlev, *Sov. J. Opt. Tech.*, vol. 40, 1973, pp. 558–561.
14. F. Reitmayer and H. Schroeder, *Appl. Opt.*, vol. 14, 1975, pp. 716–720.
15. P. J. Rogers, *SPIE*, vol. 147, 1978, pp. 141–148; and U.K. Patent No. 2,030,315.
16. D. G. Norrie, *Opt. Eng.*, vol. 25, 1986, pp. 319–322.
17. J. Angénieux et al., *SPIE*, vol. 399, 1983, pp. 446–448.
18. D. S. Garcia-Núñez and D. Michika, *SPIE*, vol. 1049, 1989, pp. 82–85.
19. V. Povey, *SPIE*, vol. 655, 1986, pp. 142–153.
20. A. D. Michael and W. B. Hart, *Metallurgist and Materials Technologist*, August 1980, pp. 434–440.
21. G. V. Thompson, U.S. Patent No. 4,148,548, 1976 U.K. priority.
22. P. J. Rogers and G. N. Andrews, *SPIE*, vol. 99, 1976, pp. 163–175.
23. P. M. Parr-Burman and A. Gardam, *SPIE*, vol. 590, 1985, pp. 11–17.
24. I. A. Neil and W. McCreath, U.K. Patent No. 2,141,260, 1983 U.K. priority.
25. I. A. Neil, U.S. Patent No. 4,659,171, 1984 U.K. priority.
26. I. A. Neil and M. Y. Turnbull, *SPIE*, vol. 590, 1985, pp. 18–29.
27. P. Nory, *SPIE*, vol. 590, 1985, pp. 30–39.
28. P. M. Parr-Burman and P. Madgwick, *SPIE*, vol. 1013, 1988, pp. 92–99.
29. R. C. Simmons and P. A. Blaine, *SPIE*, vol. 916, 1988, pp. 19–26.
30. M. Shechterman, *SPIE*, vol. 1442, 1990, pp. 276–285.
31. M. Roberts and P. R. Crew, U.K. Patent No. 2,201,011, 1987 U.K. priority.
32. M. Roberts, *SPIE*, vol. 1049, 1989, pp. 72–81.
33. C. B. Estes, U.S. Patent No. 3,205,774, 1961 U.S. priority.
34. R. C. Gibbons, *ERIM Report No. 120200-I-X*, 1976, p. 71.
35. K. Straw, *SPIE*, vol. 237, 1980, pp. 386–391.
36. M. O. Lidwell, U.S. Patent No. 4,494,819, 1980 U.K. priority.
37. T. H. Jamieson, *Opt. Eng.*, vol. 20, 1981, pp. 156–160.
38. I. A. Neil, U.S. Patent No. 4,505,535, 1982 U.K. priority.
39. M. Roberts, U.S. Patent No. 4,679,891, 1984 U.K. priority.
40. P. J. Rogers, Thermal Imaging course notes, Institute of Optics, Summer School on Lens Design, Rochester, 1988 (unpublished).
41. J. L. Rayces and L. Lebich, *SPIE*, vol. 1354, 1990, pp. 752–759.
42. M. Yatsu et al., *SPIE*, vol. 1354, 1990, pp. 663–668.
43. P. J. Rogers, *SPIE*, vols. 1780 and 1781, 1992, pp. 36–48.
44. J. M. Palmer, U.K. Patent No. 2,194,072, 1986 U.K. priority.
45. G. J. Swanson and W. B. Veldkamp, *Opt. Eng.*, vol. 28, 1989, pp. 605–608.
46. P. J. Rogers, *SPIE*, vol. 1573, 1992, pp. 13–18.

PART

2

FABRICATION

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

OPTICAL FABRICATION

Michael P. Mandina

*Brandon Light
Optimax Systems, Inc.
Ontario, New York*

9.1 INTRODUCTION

The novel creations of optical designers have been limited by the fabricator's ability to manufacture and measure the elements of the optical prescription. A solution to a design criteria often existed only on paper as the required elements were not physically realizable. Optics manufacturing technology innovations continually expand the possibilities for optical components. Increasingly, manufacturing is tethered to metrology. Creation of optics metrology instruments with accuracy equal to that of optics manufacturing equipment and vice versa has driven process development. It is this developmental symbiosis that has brought determinism to the art of precision optics manufacturing. Metrology and machine innovations offer optics of higher quality and complexity in predictable timeframes. The requirement for skilled technicians is still vital in the manufacturing process; however, the skill set is increasingly one of craft in combination with science. Artisan opticians of yesteryear still provide value; however, the future of optics manufacturing is in the hands of the 21st century optics technicians.

The methods described below are the most common for typical optical components used in industrial, aerospace, and defense applications. For the spherical lens section, a brief overview of the traditional process is described first and then the latest methods. The remaining sections will provide general overviews. The focus will be exclusively on brittle materials. For our purposes, brittle materials are defined as those where the removal process is achieved by applying mechanical forces that fracture the surface, releasing fragmented particles in a controlled manner. Much has been documented on fine finishing of brittle materials such as optical glasses, ceramics, and crystals. Works by Preston¹, Silvernail², Izumitani³, Buijs⁴, Bach⁵, Kaller⁶, Lambropoulos⁷, Golini⁸, Cook⁹, Jacobs¹⁰ and DeGroot¹¹ have contributed greatly to the understanding of optics finishing processes.

9.2 MATERIAL FORMS OF SUPPLY

Optical glass is available in boule, slab, and gob forms. Boules are formed in special disposable pots that yield a batch or glass melt of a specific glass type, such as the borosilicate glass BK7, but whose detailed characteristics are unique to that batch. Slab is yielded from a continuous flow process.

Materials are homogeneously mixed and heated, and a continuous ribbon of glass is produced. These ribbons are cut into slabs that are generally 250 mm wide, 25 mm thick, and 350 mm long, although sizes vary significantly based on supplier and material. Gobs are also yielded from a continuous flow process; however, the molten glass flows through an orifice and is sliced like cookie dough at a predetermined frequency that ensures the desired volume for the application. Gobs are almost always made to customer specifications for glass type and volume. They are used as the preblank to produce near net shape molded blanks for high-volume lens systems. Many glass suppliers also provide polished preforms, usually balls. These are used for glass molding finished optics components.

Manufacturers will order the form that best suits their purpose. The closer to final form the material, the less waste and time consumed in bulk removal operations. When rough shaping material blanks from boule or slab forms, most manufacturers use diamond impregnated saw blades and core drills to yield a part appropriate to yield the finished optic, generally called disks. Typically blanks or disks are several millimeters oversized from the final part's critical dimensions.

9.3 BASIC STEPS IN SPHERICAL OPTICS FABRICATION

Generating

This is a bulk material removal operation that starts with a near net shape molded blank or a disk.

Generating—Traditional The removal is accomplished through the application of diamonds embedded in a matrix on the cutting surface of a cup-shaped ring tool. The material is ground away as the diamonds create cracks in the surface, sweeping away glass particles where the fractures intersect.¹² The machine accuracy is generally akin to manually set, mechanically based control production equipment used in the machine tool industry. The operator continually monitors results and modifies machine settings as the cupped ring tools wear. Machine precision is adequate to control thicknesses to ± 0.025 mm and radius to ± 0.010 -mm sagittal height, but the machine is only capable of coarse removal. Subsequent lapping operations are required in order to reduce subsurface damage¹³ to a level where polishing is possible.

Generating—Modern The advent of deterministic microgrinding processes spearheaded by the work at the University of Rochester's Center for Optics Manufacturing¹⁴ in the 1990s shifted the paradigm for finishing expectations from the generating operation. As a result, machines used in the generating operation have evolved to precision machine tool status. Generators manufactured by mainstream optics manufacturing equipment providers such as OptiPro,¹⁵ Satisloh,¹⁶ Schneider¹⁷ and others, have created grinding solutions that enable the generating operation to predictably yield surfaces that are ready for polishing operations. This modern equipment makes use of CNC (computer numeric control) systems, robust motion and motion control systems such as precision linear ball slides, advanced machine base materials, structure design, and improved positioning feedback through optical encoders and other submicron feedback systems. Additionally, most of the machine builders provide in situ metrology options that enable operator assisted or completely automated parameter adjustment optimization. This is an important feature as the tool consumes itself during the process.

Complementing the advent of advanced generating machine tools for optics generating has been the increased understanding of fixed abrasive grinding mechanisms. Deterministic microgrinding is typically preferred to loose abrasive lapping when fabricators have a choice. The residual damage from microgrinding can be estimated based on glass properties.¹⁸ This aids in determining prepolish finish requirements so the overall process time for the optics can be optimized. Even with recent advancements in understanding the microgrinding process, the industry is far from offering a directory of ring tools optimized for the array of optical materials. Therefore the industry continues to rely heavily on empirical results to determine optimal setups.

Lapping

This process reduces subsurface damage left from generating to a manageable level in preparation for polishing.

Lapping—Traditional Lapping is the application of loose abrasive particles applied as slurry and pressed into the work surface by nominally constant applied pressure.¹⁹ The process typically consists of applying the abrasive slurry between a cast-iron-rotating lapping tool and the optic. Both surfaces abrade away as they remain in random dynamic contact. The fabricator controls the material removal so the operation yields the desired surface radius and smoothness. The abrasive material, often aluminum oxide, is typically between 30- and 5- μm particle size. The operator steps through particle sizes, using progressively smaller abrasives. Removal amounts account for the subsurface damage from the prior generating or lapping operation, ideally completely removing it.

Lapping—Modern The use of diamond particles embedded in a resin or metal matrix have been popular for some time. Initially, these matrices were fabricated in pellets, fastened as desired on metal backing plates, and used as laps. Abrasive work is done by the diamonds and coolant serves as lubricant and carries the glass particulate away. Unlike loose abrasive lapping, the slurry is not the abrasive. Diamond tool manufacturers also make diamond-sheet material for ready application to tools, and more recent products such as resin-bonded sheet materials from abrasive manufacturers such as 3M²⁰ can be used the same way.

Polishing

Polishing converts the finely fractured surface from the lapping or deterministic grinding operation [typical roughness of about 1- μm rms (root mean square)] into a specular surface of a surface roughness typically 1 to 3 nm rms. Polishing is a chemical-mechanical process. Water attacks the surface creating a chemically softened layer, and then the mechanical action of the abrasive in the polishing slurry, usually ceria based for optical glasses, removes the chemically softened outer layer of glass.³

Polishing—Traditional The polishing process is expected to remove the damage left from preceding operations, typically 5 to 20 μm of material. The intimate contact between the polishing tool and the optic, working with the slurry, slowly enhances the surface finish. The process is feedback based, and the fabricator works the part for a while and checks the outcome. Reacting to the results, the experienced fabricator controls various parameters to yield the desired form and finish of the polished surface.

The most basic polishing tool is a pitch polisher. Optical polishing pitch is a viscoelastic material. To form a pitch polisher, a metal tool of proper radius is coated with a 4 to 5-mm layer of polishing pitch. The pitch is warmed and formed to the optic. Once cool the brittle pitch will be cut to allow irrigation grooves for slurry access. When performed by artisans, pitch polishing can yield form errors equal to fractional wavelength of visible light routinely. Less capable fabricators may be limited to commercial quality, multiple wavelength form error outcomes.

By the 1980s, high-speed polishing had become very popular. One of the key innovations was the use of polyurethane polishing pads as a replacement for pitch. Polyurethane pads are a viscoelastic thermoplastic material with a higher viscosity than pitch. Polyurethane remains a polishing material staple and is the polishing material of choice for the fast removal seen in high-volume optics manufacturing.

Polishing—Modern Advances in deterministic polishing are dramatically changing the demands placed on optics manufacturers. Deterministic polishing is a feed-forward process, where the outcome is reasonably certain. Industry leaders in deterministic polishing development are QED/Schneider consisting of Magnetorheological Finishing (MRF)²¹ and Zeeko/SATISLOH²² who promote

a precessions polishing and air bladder solution. Each has created opportunities for manufacturers to produce optics at more predictable cost. Their application of CNC machining systems to the polishing process is revolutionizing the precision optics industry.²³

All these new solutions rely on subaperture small “pad” polishing with a known removal rate, where the “pad” may be in the form of variable stiffness polishing fluid or compliant tool made from a variety of materials and consistencies. Originally used to finish large astronomical telescope optics, small pad methods have advanced in recent years to scale cost and size down so these technologies are available to the broader population of precision optics manufactures.

Deterministic subaperture polishing solutions combine a tool’s known removal rate with an error map of the optic to produce a removal schedule. This feed-forward process relies completely on the accuracy of actual surface form information. In most cases this information is acquired from a variety of instruments such as coordinate measurement machines (CMM) or surface profilometers like the Taylor-Hobson Form TalySurf.²⁴ These instruments themselves or the software of the polishing tool convert points of data into an error map for a continuous surface. The removal profile dictates the dwell time for the small aperture polishing pad, and in general form error decreases by a factor of five per iteration. For example, if the form error is 1 wave, it is reasonable to expect that after a deterministic polishing iteration the form error will be $\sim 1/5$ wave, and after another iteration would be $\sim 1/25$ wave.

Newer technology that is also under development at a number of equipment manufacturers including QED and ZEEKO²⁵ incorporates fluid jet technology. Surfaces are corrected by directing a jet of abrasive/fluid mixture at a surface, the flow generates sufficient surface shear stress that chemical-mechanical polishing occurs.^{26,27} The jet-polishing technology is especially promising for difficult to reach areas seen in asphere and conformal optical surfacing.

Edging

Most applications of lenses require mounting into a lens housing. Lens system performance is maximized when the centers of curvature reside on the cylindrical axis of the housing. The edging operation simultaneously creates a precise (± 0.025 mm or less) diameter for mounting and aligns the centers of curvature on the mechanical centerline of the lens.

Edging—Traditional²⁸ Earlier pitch-based methods consisted of using a precision spindle where a brass cup was trued using a cutting tool. This was basically a lathe-type operation and required a skilled combination of heat, pitch, spindle velocity, timing, and consistent axial force applied by a skilled artisan in order to set the lens in a “trued” position. Once the lens was blocked, a diamond wheel ground the diameter to final dimension. Lenses are typically brought to final polished state with the diameter of the lens 1- to 3-mm oversized.

This pitch method was almost entirely replaced with mechanical bell-clamping edging machines. Bell clamping employs two opposed coaxial synchronized precision spindles and is a pitch-free process. Each spindle is affixed with a precision cup of the appropriate size to capture the lens and allow auto alignment by virtue of the mechanical forces on the variably sloped surfaces. Once the lens is “clamped” into true position, an operator mechanically defines and initiates an automated grinding sequence.

Edging—Modern In recent years, the use of CNC edging equipment is enabling a single setup for multiple grinding operations. For example, it is fairly routine to process the diameter, sagitta with a step, bevel and a fiduciary flat, all in one operation. The CNC controller interface shows a series of cross-sections, and the operator fills in inputs for what is the starting point and what features are needed in the end. Simultaneous creation ensures the features will all run true relative to one another. In addition to facilitating grinding of more complex features, optional features such as micropositioning air blasts for automated alignment optimization and measurement enable precision placement of the optic. Lenses are still mechanically bell clamped.

9.4 PLANO OPTICS FABRICATION

A plano surface has a radius equal to infinity. Typically plano form specification does not differentiate between spherical power and irregularity, specifying lump sum reflected errors as flatness. Therefore, maintaining perfect flatness is critical during plano surface finishing. The process steps for plano surfaces are exactly the same as for spheres. Planos have the advantage of fixed radius, so often, companies, departments within companies, personnel and equipment will be plano specific. This specificity allows economies of scale and development of plano-specific solutions. An example of this are continuous polishers (CPs), in which a large (40–60 inches in diameter) annular lap is “conditioned” to maintain lap flatness independent of the workpiece size. The lap is forced by a large glass (or similar material) “conditioner” to stay flat. This persuasion by the conditioner imprints onto the work piece and maintains flatness as a result.

Double-sided CPs polish both sides of a window simultaneously. Much of the recent technology used in the precision plano window manufacturing has been taken from semiconductor industry’s work-optimizing silicon wafer processes.

9.5 ASPHERE OPTICS FABRICATION

Aspheric lenses contain at least one optically active surface of nonconstant curvature. This is the primary differentiator from a spherical lens. Rotationally symmetric aspheric lenses are solids of revolution, where a general equation describes the cross section to be revolved (Fig. 1). Lenses of this style are capable of higher aberration order correction than spherical lenses. While the forms and their promise have been known to optical designers for centuries, for most of that time only the mildest forms have been physically realizable. The methods, machinery, and metrology are specific to asphere manufacturing.¹⁴

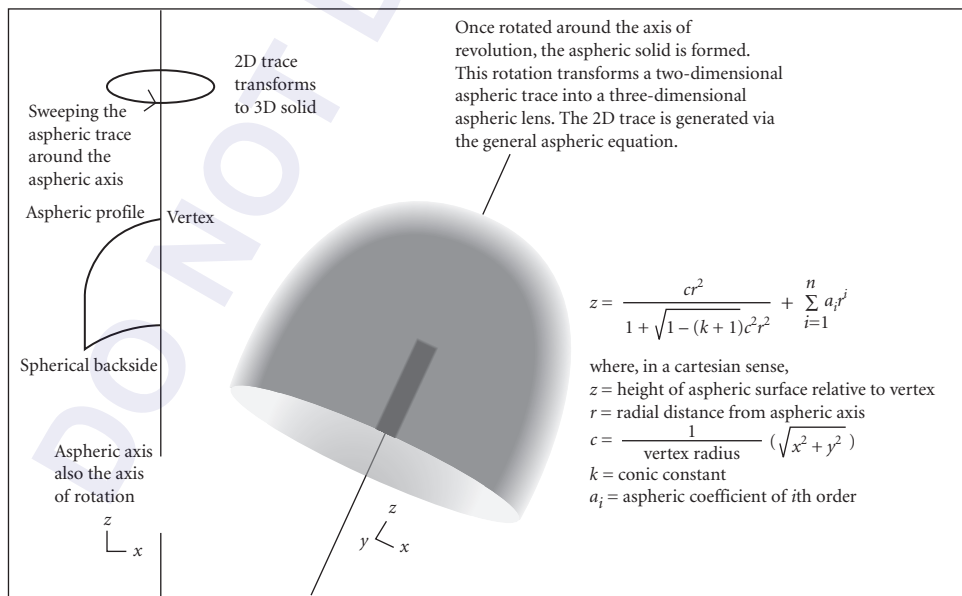


FIGURE 1 Sample general asphere form. (Brandon Light, Optimax Systems, Inc.)

Traditional full-aperture fabrication methods are not capable of manufacturing aspheric surfaces due to their nonconstant curvature. By changing the amount of contact from full to a region where change in local curvature approaches zero, some portions of traditional spherical lens-manufacturing techniques can be applied. Brittle removal by high-speed diamond grinding followed by ductile removal using a polishing slurry (ceria, alumina, etc.) can be used to prepare aspheric surfaces. Instead of full contact, curvature-insensitive local contact is used in grinding and polishing.

In aspheric grinding, a peripheral diamond wheel on a CNC platform traces the surface to generate the aspheric profile. In grinding, machine accuracy determines profile accuracy. A more accurate ground profile makes a more accurate polished profile more likely, since there's less correction needed. Particular attention must be paid to wheel wear, wheel balance, positional accuracy, and overall stiffness of the grinding platform. Imperfection in any of these grinding parameters will leave signatures in the ground surface.

The surface is then polished by working only a small area at a time. All of this must be done while maintaining location of the aspheric axis, the axis around which the solid of revolution was formed. Each iteration has an error inducement associated with it, so making as few correction runs as possible is a primary focus. Typically, asphere polishing is a feed-forward, deterministic process. While the local curvature may be constant, globally it is not. Polishing requires an adaptive tool and knowledge of what's ahead. The polishing tool needs to change to suit local curvature at a suitable rate of change. This requires knowledge of how the tool will evolve and how much removal is needed in which region. Deterministic processes provided by Zeeko/Satisloh and QED machinery, discussed earlier, are examples of such tools. These processes characterize the removal rate as a function of curvature for a given tool and combine that with an error map of the surface to be worked. The resulting removal schedule accommodates for volumes to be removed and tool performance at that local curvature.

Conventional interferometric techniques do not translate to aspheric manufacturing either. Since local curvature is nonconstant, interferometric techniques for aspheres are lock and key. The setup and equipment can be unique for a given aspheric form, so time and money demands are large. For example, form-specific computer-generated holograms (CGHs) may be required to provide feedback to the closed-loop deterministic polishing process. For a more cost-effective solution profilometry is the main two-dimensional compromise, and it is the current industry standard. Although more generalized interferometric solutions are beginning to be offered by QED, Zygo, and others.

Errors in centration are unrecoverable. In centering a spherical lens, errors can be removed. With sufficient diameter overage both centers of curvature could be positioned on the same axis and that axis could be made concurrent and coincidental with the mechanical axis. Since an aspheric surface is centered about an axis and not a point such realignment is not possible. Therefore, centration must be conserved throughout processing.

9.6 CRYSTALLINE OPTICS

As more optical work occurs outside the visible spectrum, use of optics made from nonglass brittle materials will grow. Single crystalline and polycrystalline materials are transparent far outside the usual spectral transmission range of glass. In many cases, the surfaces of these materials have differing hardness values depending on the orientation of the crystal boundaries. Soft laps tend to accentuate the grain boundaries of these materials, and that can lead to wavefront errors, mottled surfaces, and scattering. The traditional optical fabrication process can be adapted to crystal materials if some substitutions are made. The lapping process may substitute finely graded diamond for alumina and tin or zinc laps in place of the typical cast iron. Similarly, diamond suspensions are often used in polishing in place of ceria. Polishing laps may consist of synthetic materials like polyurethane or beeswax instead of optical pitches.^{29,30}

9.7 PURCHASING OPTICS

There are a number of companies who offer lines of standard optical components. These suppliers can provide off-the-shelf optics in a variety of sizes, shapes, and quality levels. Most optics providers have areas of specialization, and the informed optics buyer will select vendors that match their specific optics requirements. When custom optics are required, it is best to understand the capabilities of prospective suppliers. Most optics companies promote a broad range of capabilities, but many tend to specialize in some manner. Professionals who are engaged in optics purchasing on a regular basis learn where to go for specific optics requirements. Often this education is paid for by awarding of numerous contracts across a broad array of parts and suppliers and experiencing the consequences of the decisions. Much is learned in the contract's postmortem review.

Optics purchasing is further complicated with the predominance of the internet as a research tool. Web sites and promotional materials often do not reflect a supplier's true capability and know-how.

Whether buying off-the-shelf or custom optics, it is always best to engage potential suppliers in dialog, preferably addressing tolerances and other manufacturing cost drivers. The buyer should be satisfied the supplier has the ability to meet and measure all critical criteria. For optics that approach a manufacturer's limits, it is especially important to understand the test and acceptance process, as there can be quite a divergence of metrology equipment and methods available for testing various parameters.³¹ This is especially true for aspheres, where full format phase measuring interferometry or transmitted wavefront testing is not always within a supplier's capability.

9.8 CONCLUSION

Optics fabrication requires serial application of relatively simple steps. In the past, these steps were carefully carried out by artisans using traditional techniques. Modern approaches incorporate scientific research into the manufacturing process. Artisan skills integrate with the scientific know-how yielding a new breed of technology workers of the twenty-first century. Nevertheless, the basic process steps of grinding followed by polishing have remained. Introduction of new optical materials, more complex shapes and more narrow tolerance budgets will enable designers to develop improved solutions to old problems over an expanded spectrum, and modern manufacturing methods can make the optics physically realizable. Finally, the tendency for specialization among optics supplier requires open dialog between supplier and designer as a means to optimize successful relationships.

9.9 REFERENCES

1. F. W. Preston, "Structure of Abraded Glass Surfaces," *Trans. Opt. Soc.* **23**:141 (1922).
2. W. W. Silvernail, "Role of Cerium Oxide in Glass Polishing," in *The Science of Polishing*, Duncan Moore (ed.), OSA: Washington, D.C., 1984.
3. T. S. Izumitani, *Optical Glass*, American Institute of Physics: New York, 1986.
4. M. Buijs and K. Korpel-van Houten, "A Model for Lapping of Glass," *J. Mater. Sci.* **28**(11):3014–3020 (1993).
5. H. Bach, "Analysis of Subsurface Layers and Spots and the Reactivity of Glass Components," in *The Science of Polishing*, Duncan Moore (ed.), OSA: Washington, D.C., 1984.
6. A. Kaller, "Properties of Polishing Media for Polishing Optics," *Glastechnische Berichte—Glass Sci. Technol.* **71**(6):174–183 (1998).
7. J. C. Lambropoulos, S. Xu, and T. Fang, "Loose Abrasive Lapping Hardness of Optical Glasses and Its Interpretation," *Appl. Opt.* **36**(7):1501–1516 (1997).
8. D. Golini and S. D. Jacobs, "The Physics of Loose Abrasive Microgrinding," *Appl. Opt.* **30**:2761–2777 (1991).
9. L. M. Cook, "Chemical Processes in Glass Polishing," *J. Non-Cryst. Solids* **120**:152–171 (1990).

10. S. D. Jacobs, D. Golini, Y. Hsu, et al., "Magnetorheological Finishing: A Deterministic Process for Optics Manufacturing," *Opt. Fabrication and Testing*, T. Kasai (ed.), *SPIE* 2576:372–383 (1995).
11. J. E. DeGroote, A. E. Marino, J. P. Wilson, et al., "Removal Rate Model for Magnetorheological Finishing of Glass," *Appl. Opt.* 46:7927–7941 (2007).
12. R. E. Parks, "Optical Fabrication," in *Handbook of Optics*, 2d ed., M. Bass, E. W. Van Stryland, D. R. Williams, and W. L. Wolfe (eds.), McGraw-Hill: New York, 1995, Vol. 1. Chap. 41.
13. P. Hed and D. F. Edwards, "Relationship between Subsurface Damage Depth and Surface Roughness during Grinding of Optical Glass with Diamond Tools," *Appl. Opt.* 26(13):2491 (1987).
14. H. Pollicove and D. Golini, "Computer Numerically Controlled Fabrication," Chap. 7: *International Trends in Applied Optics*, A. H. Guenther (ed.), *SPIE*: Bellingham, Wash., Vol. PM119 (2002).
15. www.Optipro.com, OptiPro Systems, Optical Fabrication Equipment, April 21, 2009.
16. www.Satisloh.com, Optical Manufacturing Solutions, Products, Precision Optics, April 21, 2009.
17. www.schneider-om.com/home.html, Schneider GmbH & Co. KG—Fascination for Innovation: Products, April 21, 2009.
18. J. C. Lambropoulos, "Surface Microroughness of Optical Glasses under Deterministic Microgrinding," *Appl. Opt.* 35:4448–4462.
19. J. C. Lambropoulos, "Using the Grinding Merit Function (GMF): What Quality of Grind Can You Expect in the Shop?" *Convergence, Newsletter of the Center for Optics Manufacturing*, Sept./Oct. 1998.
20. www.3M.com, April 21, 2009.
21. D. Golini, G. Schneider, P. Flug, M. Demarco, "Magnetorheological Finishing," *Optics and Photonics News*, October 2001, pp. 20–24.
22. D. D. Walker, D. Brooks, A. King, et al., "The 'Precessions' Tooling for Polishing and Figuring Flat, Spherical and Aspheric Surfaces," *OSA*, 21 April, 2003; *Opt. Exp.* 11(8):958–964.
23. S. D. Jacobs, "Innovations in Polishing of Precision Optics," *Convergence, Newsletter of the Center for Optics Manufacturing*, 1st/2nd qtr. 2003.
24. www.taylor-hobson.com, Taylor Hobson—Surface Profilers, April 21, 2009.
25. www.zeeko.co.uk, Zeeko Ltd. Ultra-Precision Polishing Solutions for Optics and Other Complex Surfaces, April 21, 2009.
26. W. I. Kordonsk, A. Shorey, and M. Tricard, "Magnetorheological Jet (MRJet™) Finishing Technology," *ASME*, 128:20–26, Jan. 2006.
27. S. M. Booij, O. W. Föhnle, and J. J. M. Braat, "Shaping with Fluid Jet Polishing by Footprint Optimization," *Appl. Opt.* 43:67–69.
28. D. F. Horne, *Optical Production Technology*, Crane, Russack & Co.: New York, 1988.
29. G. W. Fynn and W. J. A. Powell, *Cutting and Polishing Optical and Electronic Materials*, 2d ed., Adam Hilger: Philadelphia, Pa., 1988.
30. R. Sumner, "Polishing of IR Materials," in *The Infrared Handbook*, W. Wolfe and G. Zissis (eds.), ERIM: Ann Arbor, Mich., 1978.
31. Y. A. Carts, "How to Buy Custom Optics That Meet Your Specifications," *Laser Focus World*, Aug. 1992, pp. 91–100.

FABRICATION OF OPTICS BY DIAMOND TURNING*†

Richard L. Rhorer

*National Institute of Standards and Technology
Gaithersburg, Maryland*

Chris J. Evans

*Zygo Corporation
Middlefield, Connecticut*

10.1 GLOSSARY

f	feed rate
h	peak-to-valley height
R	tip radius of diamond tool

10.2 INTRODUCTION

The use of special machine tools with single-crystal diamond-cutting tools to produce optical surfaces on some metals and a limited range of other materials is called *diamond turning*. Over the last 50 years or so, diamond turning has matured to become the method of choice for producing some optical surfaces; in other applications, diamond turning provides a critical process step with radically different characteristics from most other optical fabrication methods.

In terms of geometry and motions required, the diamond-turning process is much like the step of “generating the optical surface” in traditional optical fabrication. However, the diamond-turning machine is a more sophisticated piece of equipment that produces the final surface, which frequently does not need the traditional polishing operation. The surface quality produced by the “best” diamond turning does not yet match the best produced by conventional polishing practice. The limits of diamond turning for both figure and surface-finish accuracy have not yet been reached—and diamond turning can be combined with postpolishing to improve surface finish and reduce scatter.¹ Also subaperture processing with small polishing tools or magnetorheological finishing (MRF) can be used to improve figure.

There are several important advantages of using diamond turning, including the ability to produce good optical surfaces to the edge of the element, to fabricate soft ductile materials that are difficult to polish, to eliminate alignment adjustments in some systems, and to fabricate shapes difficult to produce by other methods. The latest generation of diamond-turning machines incorporates up to five axes

*Official contribution of the National Institute of Standards and Technology; not subject to copyright in the United States.

†Certain commercial equipment, instruments, or materials are identified in this chapter. Such identification does not imply recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply that the products identified are necessarily the best available for the purpose.

of computer-controlled motion, allowing for production of anamorphic optics. Use of form tools on multiaxis machines enables production of “structured” optical surfaces² ranging from subwavelength structures through diffractive/refractive infrared (IR) elements to optical component molds.

If the advantages of diamond turning suggest this fabrication method, then it is important to determine early in the design phase of a project whether the material specified is appropriate for diamond turning and whether slideway travels and linear and rotary axis controls are available on the diamond-turning machine to support fabrication of complex structures.

Sections in this chapter highlight the following:

- The diamond-turning process
- The advantages of diamond turning
- Diamond-turnable materials
- Comparison of diamond turning and traditional optical fabrication
- Machine tools for diamond turning
- Basic steps in diamond turning
- Surface finish of diamond-turned optics
- Metrology of diamond-turned optics
- Conclusions

10.3 THE DIAMOND-TURNING PROCESS

The diamond-turning process produces finished surfaces by very accurately cutting away a thin chip or layer of the surface. Thus, it is generally applicable to ductile materials that machine well rather than to hard brittle materials traditionally used for optical elements. However, by using a grinding head on a diamond-turning machine in place of the tool, hard brittle materials can be finished. At very small effective depths of cut, brittle materials behave in an apparently ductile manner. This attribute allows fracture-free grinding of glasses and ceramics as well as diamond turning of optical surfaces on materials such as germanium, zinc selenide, and potassium dihydrogen phosphate (KDP).

In diamond turning, both the figure and surface finish are largely determined by the machine tool and the cutting process. Note, however, that material characteristics such as grain size and inclusion size limit the ultimate surface finish achievable. The tool has to be very accurately moved with respect to the optical element to generate a good optical surface, and the edge of the diamond tool has to be extremely sharp and free of defects.

10.4 THE ADVANTAGES OF DIAMOND TURNING

Diamond turning fits within a broad spectrum of optics fabrication processes. When compared with traditional optical fabrication methods of lapping and polishing (see, for example, Chap. 9, “Optical Fabrication,” by Michael P. Mandina) diamond turning has several advantages.

- It can produce good optical surfaces clear to the edge of the optical element. This is important, for example, in making scanners, polygons, special shaped flats, and when producing parts with interrupted cuts.
- It can produce optical surfaces on soft ductile materials that are extremely difficult to polish.
- It can easily produce off-axis parabolas and other difficult-to-lap aspherical shapes.
- It can produce optical elements with a significant cost advantage over conventional lapping and polishing where the relationship of the mounting surface—or other feature—to the optical surface is very critical. Expressed differently, this feature of diamond turning offers the opportunity to eliminate alignment adjustments in some systems.

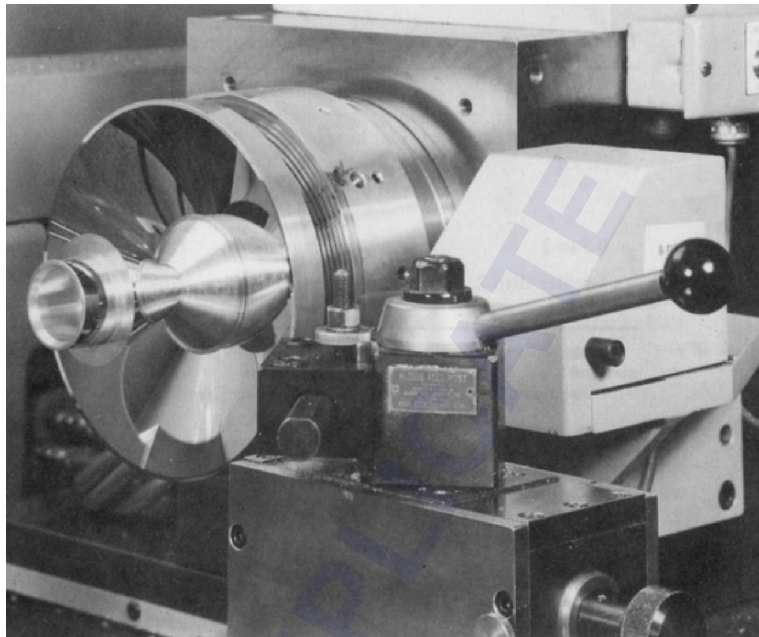


FIGURE 1 An axicon optical element being diamond turned. (Courtesy of Rank Taylor Hobson, Keene, New Hampshire.)

- It can fabricate optical shapes such as axicons, faceted optics, and grazing incidence X-ray optics that would be extremely difficult to fabricate by methods other than diamond turning (see Fig. 1).

Conflicts between optical requirements and diamond turnability on the one hand, and mechanical considerations on the other, often lead to the use of platings. Plating deficiencies, however, can cause as much trouble as poor bulk materials. For example, small changes in the composition of plated electroless nickel may cause dramatic changes in tool wear.³

Residual stress in the mirror blank, whether plated or not, can lead to changes in mirror shape with time. It is essential to pay careful attention to stress-relief prior to final diamond turning.

A decision to diamond turn an optical element, rather than fabricate it by the conventional polishing techniques, might be based on several different considerations such as type of element, size, and material. A general guide to different considerations in selecting diamond turning as a fabrication technique is presented in Table 1.

TABLE 1 General Guide to Optical Fabrication Methods

Size, m	Shape	Material	Preferred Method
Less than 0.5	Flat or sphere	Glass/ceramic	Polish
		Ductile metal	Diamond turn
	Asphere	Glass/ceramic	Grind/polish*
		Ductile metal	Diamond turn
0.5 to 2.0	Any axisymmetric	Ductile metal	Diamond turn [†]
Greater than 2.0	Any	Any	Large polishing machine

*Can generate shape or figure on a diamond-turning machine with a grinding head replacing the diamond tool. Subaperture polishing techniques, including techniques such as MRF, may be applied to advantage.

[†]Diamond-turning machines up to 2-m diameter have been built.

As indicated above, diamond turning has some unique characteristics. In some IR (and even visible) imaging systems, considerable improvements in optical performance have been obtained by combining a refractive aspheric surface and a diffractive surface in a single element. For IR applications, it is hard to produce such a component by any other fabrication process; for visible applications, such optics can be produced in volume from diamond-turned molds.

Another unique capability of diamond turning is to provide datums or alignment features machined in the same setup as the optical surface. “Snap-together” optical systems requiring no alignment adjustments after assembly are very attractive in some applications.

Over the last decade, there have been considerable advances in the ability to produce aspheric optics using computer-controlled generators and pad polishers. These technologies, combined with ion polishing, magnetorheological finishing, and computer-controlled polishing have enabled a new generation of aspheric optics. Ultimately the choice of manufacturing process requires a careful analysis of the options and the system requirements.

10.5 DIAMOND-TURNABLE MATERIALS

Selection of appropriate materials is, necessarily, a trade-off between application-specific requirements and optimization of the manufacturing process. This trade-off may drive the selection of a plated surface, for example, or the choice of fabrication steps.

Historically materials have been described as either “diamond turnable,” or not, as if this were an inherent material property. This shorthand covers two different situations. One is that, in practice, some materials cause very rapid wear of the diamond; for example, it is widely known that ferrous materials cause rapid tool wear. The other is that, particularly for certain plastics, tool-workpiece interactions produce unacceptable optical surfaces.

A number of listings of diamond-turnable materials, such as the one included in Table 2, have been published. Such listings should be treated with caution. Typically, they are incomplete and do not provide sufficient information on the materials that are listed. For example, good optical surfaces are not generally produced on all aluminum alloys: Aluminum Alloy 6061 (Aluminum Association, Inc. designation) is the most commonly used alloy, although certain 5000 series and 7000 series alloys have their proponents, and 2024 aluminum has been used but, in general, does not produce the best surfaces.

TABLE 2 Diamond-Turnable Materials

Metals	Nonmetals	Plastics
Aluminum	Calcium fluoride	Polymethacrylate
Brass	Magnesium fluoride	Polycarbonates
Copper	Cadmium telluride	Polyimide
Beryllium copper	Zinc selenide	
Bronze	Zinc sulphide	
Gold	Gallium arsenide	
Silver	Sodium chloride	
Lead	Calcium chloride	
Platinum	Germanium	
Tin	Strontium fluoride	
Zinc	Sodium fluoride	
Electroless nickel	KDP	
	KTP	
	Silicon	

Similarly, gold is considered diamond-turnable, but problems have been reported machining large gold-plated optics. Conventional electroplated nickels (and bulk nickel) give rapid tool wear, but electroless nickel with phosphorous contents above about 10 percent, if appropriately heat treated, can be machined effectively.⁴ Higher phosphorous contents (up to 15 percent) are obtainable in electroplated nickel^{5,6} which also machines extremely well. Both materials are metastable and will transform—with exposure to the necessary time/temperature conditions—to a mixture of crystalline nickel with hard phosphides. This transformation is accompanied by a volume change, a degradation of optical characteristics of the surface, and a dramatic increase in tool wear.

Platings may also be optimized to give low ductility and hence minimum burr formation when machining Fresnels or the molds for micro-optics arrays such as arrays of retroreflectors. Platings over a diamond-turned sacrificial mandrel allow production of otherwise unobtainable forms. Plated surfaces, however, have characteristics which can adversely affect both fabrication and application; at some level the resulting structure is a temperature-sensitive bimetal strip. Pits and inclusions can cause significant fabrication issues.⁷

Silicon, although included in the listing given here, should be considered marginal as tool wear can be high. Reasonably large areas of amorphous silicon cladding are reported to have been successfully machined.

Over the last decade or so there have been significant advances in understanding of diamond tool wear. Mechanisms associated with abrasion and chipping typically provide one limit to diamond tool life. For example, when machining bulk aluminums the interactions between hard inclusions and the diamond tool clearly lead to wear. Such mechanisms, however, do not explain the very rapid wear observed when machining soft, high-purity iron.

Paul et al.⁸ showed that machining metallic elements containing unpaired *d*-shell electrons results in catalyzed reactions between diamond and the work material. The same mechanism explains the role of phosphorous in electroless nickel and led to a recent breakthrough by Brinksmeier et al.⁹ They showed that, like phosphorous, nitrogen in a nitride surface layer on steel combines with the unpaired *d*-shell electrons from the iron. The result is a dramatic reduction in tool wear, suggesting that diamond-turned steel molds (e.g., for plastic optics) may become practical in the near future. Previous approaches—such as diamond turning at cryogenic temperatures¹⁰ or in methane or acetylene environments¹¹—provided evidence of the mechanisms at work but were not widely adopted (and in the case of cryogenic machining was not intended by the original researchers to be practical).

For a number of years, Moriwaki¹² has been developing ultrasonic-assisted machining, including cutting when the tool is vibrated with an elliptical motion. Significant reductions in tool wear have been demonstrated, although the mechanism remains controversial. The amplitude and frequency of oscillation in the cutting direction are generally selected such that separation between the rake face of the tool and the chip is expected. In this case, one might postulate poisoning of the catalytic process by, for example, hydrocarbon-based cutting fluids. Elliptical motion would also move the clearance face out of contact. Other measurements show significant reductions in cutting forces,¹³ suggesting a reduction of tool temperatures.

In case of some plastics, recent work by Gubbels et al.¹⁴ shows that the chemical explanations of Paul et al.⁸ do not apply, but that triboelectric effects dominate. In general, there are some suggestions that parameters, such as surface speed, are more important for successful diamond machining of plastics than for metal and crystalline substrates. Some plastics are diamond turned in volume production.

Other material characteristics, in addition to the material composition, are important. For example, large grain size results in a more pronounced “orange peel” as tools become dull and the variation in modulus of the different grain orientations leads to different deflections due to cutting forces. Residual stresses can relax over time and cause changes in figure. Because of these types of problems it is important to involve experienced personnel early in the design phase¹⁵ to ensure that the material specified is appropriate. In some projects, the part is so valuable and/or so difficult to produce by other techniques, it is worth consuming tools more rapidly than would normally be acceptable. However, such a decision should be taken consciously, not by default late in a project.

10.6 COMPARISON OF DIAMOND TURNING AND TRADITIONAL OPTICAL FABRICATION

In diamond turning, the final shape and surface of the optic produced depend on the machine tool accuracy, whereas, in traditional optical fabrication, the final shape and surface of the optical element depend on the process variables involved with using an abrasive-loaded lap. The differences between diamond turning and traditional optical fabrication can be summarized by describing diamond turning as a displacement-controlled process versus a force-controlled process for traditional optical fabrication.¹⁶ The goal in diamond turning is to have a machine tool that produces an extremely accurate path with the diamond tool, hence a displacement-controlled process. A traditional polishing machine used for optical fabrication depends on the force being constant over the area where the abrasive-loaded lap—or tool—touches the surface being worked. Selective removal of material can be produced by increasing the lap pressure in selected areas or by use of a zone lap. The stiffness of a diamond-turning machine is important because, to control the displacement, it is important that cutting forces and other influences do not cause unwanted displacements. Feeds, speeds, and depth of cut are typically much lower in diamond turning than conventional machining, thus giving lower forces. However, the displacements of concern are also much lower. Thus the stiffness required is as much, or more, of a concern than conventional machining even though the total force capability may be lower for diamond turning.

10.7 MACHINE TOOLS FOR DIAMOND TURNING

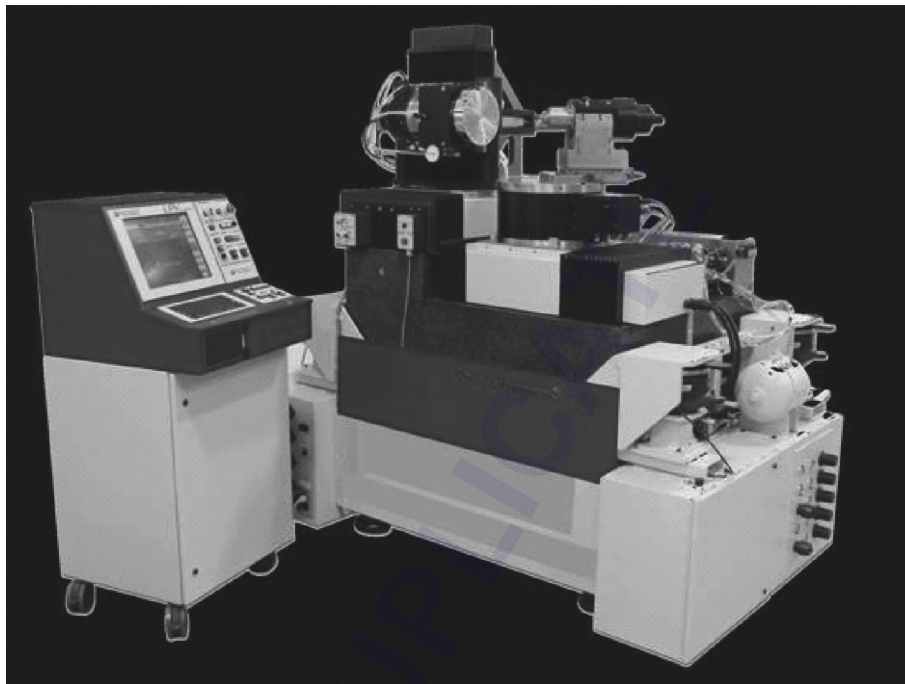
In general, the machine tools used for diamond turning are very expensive compared to the equipment needed for traditional optical fabrication. The positioning accuracy required for diamond turning is beyond the capability of conventional machine tools, thus some of the first widely adopted diamond-turning machines for fabricating optics were modified Moore measuring machines.¹⁷

Although there are some records of machine tools being used to generate optical surfaces as early as the seventeenth century, most of the effort is modern, accelerated in the 1960s and 1970s with the advent of computer-based machine tool controls and laser interferometer systems used as positional feedback devices. Evans¹⁸ has documented much of the history of diamond turning and provides an extensive reference list. Ikawa¹⁹ summarizes some of the research in metal cutting related to diamond turning and associated machine tools.

The early diamond-turning machines were two-axis lathes that could produce axisymmetric optical elements. With recent advances in computer-based control systems, and improved motion control and feedback sensors, multiaxis diamond-turning machines have become readily available. Two commercial diamond-turning machines are shown in Fig. 2. Both machines can be configured with five-axis motion control combining both linear and rotary axes. Measuring scales have replaced the laser interferometers in many diamond-turning machines and give a very reliable positioning feedback system at lower cost.

Programming of these multiaxis machines draws on the technology developed in precision machine shops for large five-axis machine tools used to make complicated parts. By adapting the multiaxis control to diamond-turning machines, a great variety of shapes can now be diamond turned which opens up the process to many new optical applications. Before judging an optical element shape to be unsuitable for diamond turning, a manufacturer of modern diamond-turning machines should be consulted.

Producing nonaxisymmetric parts—such as an off-axis parabola machined while centered on the rotating axis—has become possible with fast tool servos. These systems can rapidly move a cutting tool a short distance coordinated with the rotation of the spindle.²⁰ There are also cases where the machine's slideways or rotary motions can be used to produce nonaxisymmetric parts.



(a)



(b)

FIGURE 2 (a) Diamond-turning machine configurable for five-axis machining. (Courtesy of Precitech, Inc., Keene, New Hampshire.) (b) Diamond-turning machine configurable for five-axis machining. (Courtesy of Moore Nanotechnology Systems, LLC., Keene, New Hampshire.)

Many diamond-turning machines are used in the traditional turning lathe mode where the workpiece turns and the tool is held stationary in the tool post. Most diamond-turning machines can also be configured such that the tool rotates about the spindle axis—commonly called fly cutting—to produce components such as long flat mirror surfaces or other milled surfaces.

10.8 BASIC STEPS IN DIAMOND TURNING

Much like the traditional optical-fabrication process, the diamond-turning process can be described as a series of steps used to make an optical element. The steps used in diamond turning are

1. *Preparing the blank* with all the required features of the element with an extra thickness of material (generally 0.1-mm extra material or plating is adequate) on the surface to be diamond turned
2. *Mounting the blank* in an appropriate fixture or chuck on the diamond-turning machine
3. *Selecting the diamond tool* appropriate for the material and shape of the optical component
4. *Mounting and adjusting the diamond tool* on the machine
5. *Machining the optical surface* to final shape and surface quality
6. *Cleaning the optical surface* to remove cutting oils or solvents

Mounting the optical element blank on a diamond-turning machine is extremely important. If a blank is slightly distorted in the holding fixture, and then machined to a perfect shape on the machine, it will be a distorted mirror when released from the fixture. Therefore, fixtures and chucks to hold mirrors during diamond turning need to be carefully designed to prevent distortion. Often the best way to hold a mirror during machining is to use the same mounting method that will be used to hold the mirror in service.

It is advantageous in many applications to machine a substrate of aluminum or copper and then add a plating to be diamond turned. The design and application of platings is part science and part art. Many aspects of the platings as related to diamond turning were covered at the ASPE Spring 1991 Topical Meeting.⁷

Tool setting—the mounting and adjusting of the diamond-tipped cutting tool—is often accomplished by cutting a test surface, either on the actual mirror blank to be later machined over, or by placing a test piece on the machine just for tool setting. If the cutting tool is too high or too low, a defect at the center of a mirror is produced. It is possible, using reasonable care and patience, to set the tool height within about 0.1 μm of the exact center. Setting the tool in the feed direction after the height is correct is somewhat more difficult. For example, an error in setting will produce an ogive shape rather than a sphere which is not obvious until the figure is measured. Gerchman²¹ describes these types of defects.

The selection of the tool for diamond turning is important. Large cutting tip radii (2 mm or greater) are often used when producing flats, convex, or concave mirrors with large radius of curvature. However, small-radii diamond tools are available (in the range of 0.1 mm) for making small deep mirrors or molds. Tools with special geometries, including so-called “dead sharp” tools, can be obtained for such applications as Fresnel lenses or retroreflector arrays. In general, approximately 0° rake tools, with about 5° or 6° front clearance, are used for diamond-turning ductile metals. Negative rake tools are often good for crystalline materials and positive rakes may be beneficial when machining some plastics. The cutting edge has to be chip free to produce a good diamond-turned surface. A normal specification for edge quality is “chip free when examined at 1000 \times .” The edge sharpness is a concern for very small depths of cut—especially where the depth of cut is close to the cutting edge sharpness—because the cutting forces increase and more of a plowing than a cutting process occurs. The effect of cutting edge sharpness has been investigated by researchers, for example Lucca, et al,²² however, there is currently no convenient way to specify and inspect tools for edge sharpness.

The orientation of the diamond itself on the shank is of concern because the single-crystal diamond is anisotropic. The orientation of diamond tools has been studied, for example, by Wilks,²³ Decker,²⁴ and Hurt.²⁵ It is necessary for the tool manufacturer to mount the diamond so that it can be shaped to the required radius and produce a good cutting edge. The usual orientation for diamond tools is with the cleavage plane parallel to the rake face.

The actual diamond turning, or machining to final size and surface finish, is often the fastest part of the process. The machine-tool controller has to be programmed to move the tool along the correct path, the chip-removal system has to be positioned, and the cutting-fluid applicator needs to be adjusted to provide consistent clean cutting.

For machining of flats and spherical surfaces, the part programs that define the machine motion are straightforward. But when cutting aspherical surfaces, caution has to be exercised so that the radius of the tool is properly handled in calculating the tool path. Modern computer-aided design (CAD) systems perform the necessary calculations, but tests should be performed prior to cutting a difficult or expensive component.

In general, the cutting speeds for diamond turning are similar to those used for conventional machining: less than 1 m/min to more than 100 m/min. However, the slower cutting speeds produced by facing to the center of a workpiece do not affect the surface finish in diamond turning as is often the case with nondiamond tools. Thus, varying the spindle speed to keep the cutting speed constant is not necessary in diamond turning. The upper speed for diamond turning is often limited by the distortion of the optical element due to inertial forces, especially for larger elements. The upper spindle speed can also be limited due to any unbalance of the workpiece and fixture. The feed rate in diamond turning is usually adjusted to give a good theoretical surface finish. (See the following section.)

Cleaning of diamond-turned optics has a lot in common with cleaning conventionally polished optics. But because many of the diamond-turned elements are of soft metals, caution has to be exercised to prevent scratching. In general, a degreaser is used (soap or solvent), followed by a rinse in pure ethyl alcohol. The drag-wiping technique traditionally used on some glass optics can be used on some diamond-turned elements. Care must be taken to ensure that the lens tissue is very clean and remains wet. Some work has been done to study the best solvents to use for cleaning diamond-turned optics from an environmental-impact standpoint.²⁶

10.9 SURFACE FINISH OF DIAMOND-TURNED OPTICS

The surface structure is different for diamond-turned surfaces as compared with conventionally polished surfaces. A diamond-turned surface is produced by moving a cutting tool across the surface of the turning component, such as the facing operation illustrated in Fig. 3. Therefore, diamond-turned elements always have some periodic surface roughness, which can produce a diffraction-grating effect, whereas polished optical surfaces have a random roughness pattern. The traditional “scratch and dig” approach to describing surfaces is not meaningful for diamond-turned surfaces.

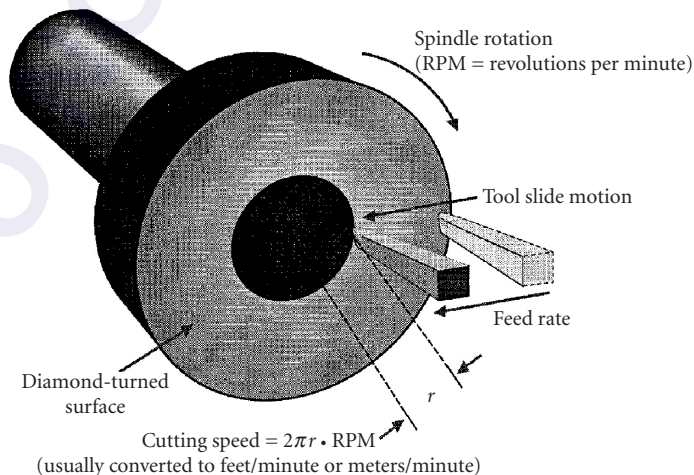


FIGURE 3 Diamond turning an optical element.

The machining process produces a periodic surface structure directly related to the tool tip radius and feed rate. The theoretical diamond-turned surface is illustrated in Fig. 4. The formula displayed in the figure for calculating the height of the cusps is

$$h = \frac{f^2}{8 \cdot R} \tag{1}$$

where h = peak-to-valley height of the periodic surface defect
 f = feed per revolution
 R = tool tip radius

For example, if a surface is diamond turned using a spindle speed of 31.4 rad/s (300 rpm), a feed of 7.5 mm/min, and a 5.0-mm tool tip radius:

$$h = \frac{(7.5/300)^2}{8 \times 5} = 1.56 \times 10^{-5} \text{ mm}$$

$$h = 15.6 \text{ nm} \tag{2}$$

In addition to the “theoretical finish” based on cusp structure, the measured surface finish on diamond-turned parts is influenced by other factors.

- Waviness within the long-wavelength cut-off for surface measurement may be correlated, for example, with slide straightness errors.
- Asynchronous error motions of the spindle can cause surface defects. If, for a given angular spindle position, there is nonrepeatability in axial, radial, or tilt directions, these errors will transfer into surface structure. Details of spindle errors are important in diamond turning. Further information can be found in the “Axis of Rotation Standard.”²⁷
- External and self-induced vibration, not at the spindle frequency or at one of its harmonics, can have the same effect on finish—measured across the lay—as asynchronous spindle motions.
- Materials effects such as the differential elastic recovery of adjacent grains can cause steps in the machined surface that have an appearance commonly referred to as “orange peel.” Impurities in the material can also degrade surface finish.
- Within each cusp, there can be a repeated surface structure related to chips in the edge of the tool.

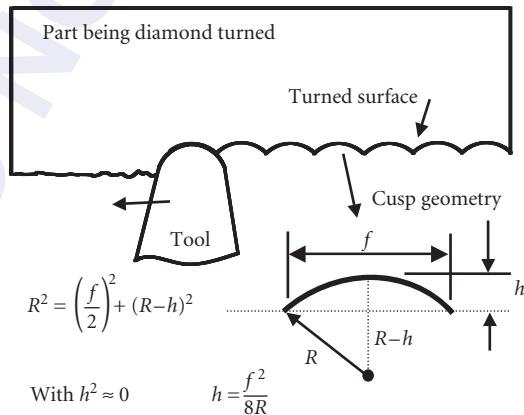
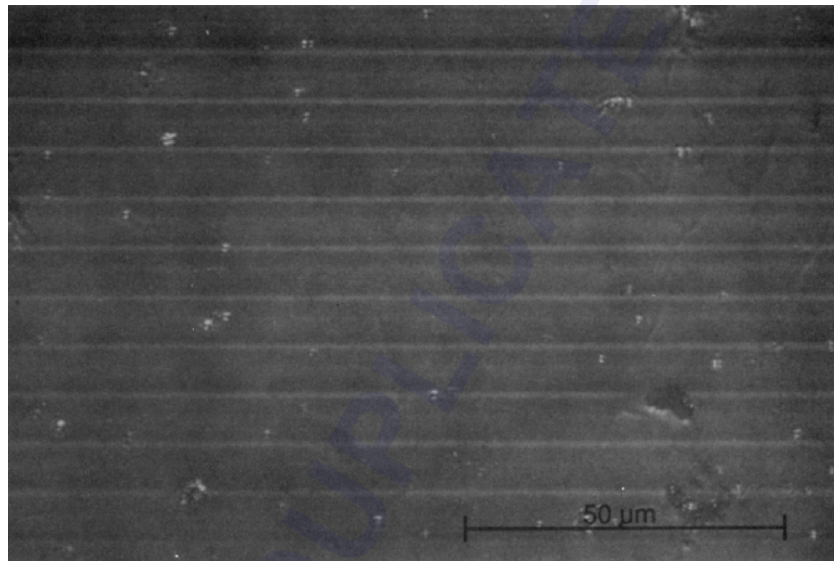
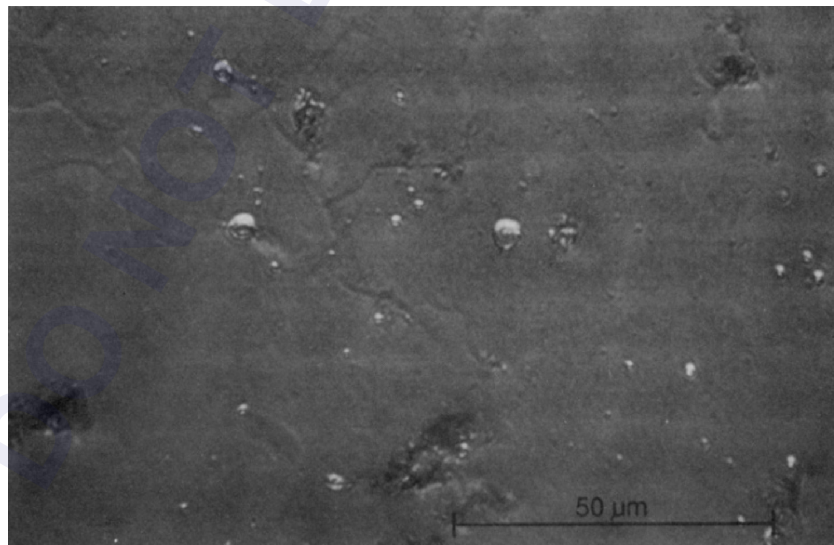


FIGURE 4 “Cusp” surface of diamond-turned optical element.

The Nomarski microscope is an excellent means of qualitatively evaluating diamond-turned surfaces. The Nomarski photo in Fig. 5a illustrates the periodic structure of a diamond-turned surface. The feed rate used in producing the surface causes the wavelength of the periodic structure to be about $8\ \mu\text{m}$. Figure 5b illustrates other defects in the diamond-turned surface when the Nomarski microscope is adjusted such that the periodic cusps are not seen.²⁸



(a)



(b)

FIGURE 5 Nomarski micrograph of a diamond-turned aluminum alloy (a) aligned so that the grooves can be seen and (b) aligned so that the grooves are canceled. (From Bennett, p. 84.²⁸)

10.10 METROLOGY OF DIAMOND-TURNED OPTICS

In general, measurement of diamond-turned optics is similar to the measurement of any other optic; figure, midspatial frequency errors, transmitted wavefront, and surface roughness may all need characterization, depending on the specification. As with other optics, the choice of figure metrology is driven by the optical surface itself. Classical null tests—especially autocollimation tests for parabolae and the related tests for other conics—are widely used. Over the last decade or so, use of in-cavity holograms in Fizeau tests has increased. What remains elusive is a general test. Over a limited range of surfaces, subaperture stitching²⁹ may be viable or, for circularly symmetric aspheres, zonal stitching.³⁰ Kuechel³¹ described a zonal technique that uses only the zone of null data and, hence, is free of retrace errors and applicable to a range of aspheres without the need for null optics. An instrument based on this technique is shown in Fig. 6.

One area in which diamond turning differs from conventional optics production is that the machine itself can be used as a measuring machine. The diamond tool can be replaced with an appropriate sensor (such as a capacitance sensor, air bearing linear variable differential transformer (LVDT), optical triangulation sensor, etc.) or the sensor can be built into an auxiliary mount. With sufficient care,³² the geometric errors of the machine can be mapped so that the limits in the metrology are the uncertainties associated with probing and with the environment. This approach is particularly advantageous when making (and measuring) radical aspheres or discontinuous, structured surfaces² such as molds for faceted automotive lighting. On multi-axis machines, it is sometimes more useful to use a different combination of axes for metrology than for machining to better decouple machine geometry errors from measurement uncertainty. For example, on a diamond-turning machine with a B axis (rotary table), near hemispheres and some aspheres can conveniently be machined using only the x and y axes, with measurement of the departure from a best-fit sphere performed using a separate probe mounted on the rotary table.



FIGURE 6 Aspheric measuring system. (Courtesy Zygo Corporation, Middlefield, CT.)



FIGURE 7 Microinterferometer. (Courtesy of Zygo Corporation).

There is little practical difference between measuring optical surfaces produced using traditional methods and by diamond turning. It is worth bearing in mind, however, that during diamond turning there is usually a monotonic progression in cutting from outside diameter to inside diameter or vice versa; hence, diamond tool wear or small edge nicks will cause a degradation in finish that depends on position on the part. The surface finish measurement sampling strategy should be adjusted accordingly. Surfaces produced using traditional 2-axis or 3-axis diamond-turning have significantly different characteristics along and transverse to the lay; scattering is isotropic, a characteristic that should be considered in both the specification and metrology of diamond-turned optics. Four-axis and 5-axis machining using methods akin to milling produce cusp structures usually at different spatial wavelengths in both directions.

Microinterferometers (Fig. 7) have become the tool of choice for characterizing optical surfaces at spatial wavelengths down to the limits posed by the instrument transfer function.³³ Microinterferometers—particularly those using scanning coherence techniques frequently referred to as scanning white light interferometry (SWLI)—can be useful, provided the surface slopes and lateral extent are compatible with the available numerical aperture of the objective and the field of view. Replication—for example, using dental replica materials, silicone-based caulks, two-part epoxies, and the like—allows sampling of large surfaces, although there is inevitably some increase in “noise” due to the replication process.

Higher spatial frequency structured surfaces, such as retroreflectors or other micro-optic arrays,³⁴ often pose metrology challenges for which there is no general solution.

10.11 CONCLUSIONS

Diamond turning has been used for many years to commercially produce infrared optics. Some visible and ultraviolet applications are now possible. Moreover, the limits of diamond turning for both figure and surface finish accuracy have not yet been reached. Taniguchi³⁵ and others have shown that precision in both conventional machining and ultraprecision machining, such as diamond turning, has steadily improved for many decades, with roughly a factor of three improvements possible every 10 years. If this trend continues, we could expect diamond-turning machines with accuracies below 10 nm and even approaching 1 nm by the year 2020. Yet, it is important to remember that it becomes increasingly difficult to push the capabilities in this regime—nor is it clear that it is cost effective to do so. Other manufacturing techniques may be more appropriate for production of the highest quality optics.

The technology developed for diamond-turning optics in some industries is now beginning to impact the precision machining of nonoptical components. In the future, the improvement of all machine tools will likely be driven by both optical and nonoptical applications, with diamond-turning machines possibly reaching the accuracy level that will allow visible and ultraviolet optics to be fabricated by machining or grinding without postpolishing.

10.12 REFERENCES

1. R. E. Parks and C. J. Evans, "Rapid Post-Polishing of Diamond-Turned Optics," *Precision Engineering* **16**: 223–227 (1994).
2. C. J. Evans and J. B. Bryan, "Structured, Engineered and Textured Surfaces," *CIRP Annals* **48/2**:541–546 (1999).
3. C. K. Syn, J. S. Taylor, and R. R. Donaldson, "Diamond Tool Wear vs. Cutting Distance on Electroless Nickel Mirrors," *Proc. SPIE* **676**:128–140 (1986).
4. J. S. Taylor, C. K. Syn, T. T. Saito, and R. R. Donaldson, "Surface Finish Measurements of Diamond Turned Electroless Nickel Plated Mirrors," *Optical Engineering* **25**(9):1013–1020 (1986).
5. A. Mayer, et al., "Electrodeposited Coatings for Diamond Turning Applications," *Proc. ASPE Spring Topical Meeting, Metal Platings for Precision Finishing Operations*, 1991.
6. C. J. Evans, R. S. Polvani, and A. Mayer, "Diamond Turned Electrodeposited Nickel Alloys," *OSA Technical Digest Series* **9**:110 (1990).
7. ASPE, "Metal Platings for Precision Finishing Operations," *Spring Topical Meeting*, Raleigh, N.C., 1991.
8. E. Paul, C. J. Evans, A. Mangamelli, M. L. McGlauffin, and R. S. Polvani, "Chemical Aspects of Tool Wear in Single Point Diamond Turning," *Precision Engineering* **18**:4–19 (1996).
9. E. Brinksmeier, R. Glabe, and J. Osmer, "Ultra-Precision Diamond Cutting of Steel Molds," *CIRP Annals* **55**:551–554 (2006).
10. C. J. Evans, "Cryogenic Diamond Turning of Stainless Steel," *CIRP Annals* **40**:571–575 (1991).
11. J. Casstevens, "Diamond Turning of Steel in Carbon Saturated Environments," *Precision Engineering* **5**:9–15 (1983).
12. T. Moriwaki, "Ultraprecision Diamond Turning of Stainless Steel by Applying Ultrasonic Vibration," *CIRP Annals* **40**:559–562 (1991).
13. E. Shamoto and T. Moriwaki, "Ultraprecision Diamond Cutting of Hardened Steel by Applying Elliptical Vibration Cutting," *CIRP Annals* **48**:441–444 (1999).
14. G. P. H. Gubbels, G. J. F. T. van der Beek, A. L. Hoep, F. L. M. Delbressine, and H. Halewijn, "Diamond Tool Wear When Cutting Amorphous Polymers," *CIRP Annals* **53**:447–450 (2004).
15. J. S. Taylor and C. J. Evans, "Fabrication of a Metal Plated Mirror, Beginning from a Performance Specification," *Proc. ASPE Spring Topical Meeting, Metal Platings for Precision Finishing Operations*, 1991.
16. A. Gee, Cranfield Institute of Technology, in a private communication to the authors used the description of "displacement" controlled and "force" controlled.
17. W. R. Moore, *Foundations of Mechanical Accuracy*, Moore Special Tool Co., Bridgeport, CT., 1970.
18. C. J. Evans, *Precision Engineering: An Evolutionary View*, Cranfield Press, Bedford, UK, 1989, pp. 135–155.
19. N. Ikawa, et al., "Ultraprecision Metal Cutting : the Past, the Present, and the Future," *CIRP Annals* **40**: 587–594 (1991).
20. S. Patterson and E. Magrab, "Design and Testing of a Fast Tool Servo for Diamond Turning," *Precision Engineering* **7**:131–136 (1985).
21. M. C. Gerchman, "Optical Tolerancing for Diamond Turning Ogive Error," *Proc. SPIE (Reflective Optics II)*: 224–229 (1989).
22. D. A. Lucca, Y. W. Seo, R. L. Rhorer, and R. R. Donaldson, "Aspects of Surface Generation in Orthogonal Ultraprecision Machining," *CIRP Annals* **43**:43–46 (1994).
23. J. Wilks, "Performance of Diamonds as Cutting Tools for Precision Machining," *Precision Engineering* **2**: 57–71 (1980).

24. D. C. Decker, H. H. Hurt, J. H. Dancy, and C. W. Fountain, "Preselection of Diamond Single Point Tools," *Proc. SPIE* **508** (1986).
25. H. H. Hurt and D. L. Decker, "Tribological Considerations of the Diamond Single Point Tool," *Proc. SPIE* **508** (1986).
26. L. A. Theye and R. D. Day, "Evaluation of Environmentally Safe Cleaning Agents for Diamond Turned Optics," *Proc. ASPE*, ASPE, Raleigh, N.C., 1991.
27. ANSI/ASME B89.3.4M-2006 Standard: "Axes of Rotation: Methods for Specifying and Testing," ASME, New York, 2006.
28. J. M. Bennett and L. Mattsson, *Introduction to Surface Roughness and Scattering*, Optical Society of America, Washington, D.C., 1989.
29. P. Murphy, J. Fleig, G. Forbes, D. Miladinovic, G. DeVries, and S. O'Donohue, "Subaperture Stitching Interferometry for Testing Mild Aspheres," *Proc. SPIE* **6293** (2006).
30. M. J. Tronolone, J. F. Fleig, C. Huang, and J. H. Bruning, US Patent 5,416,586 (1995).
31. M. Kuechel, US Patent 6,781,700, (2004).
32. W. T. Estler and E. B. Magrab, "Validation Metrology of the Large Optics Diamond Turning Machine," NBSIR 85-3182(R), U.S. Department of Commerce, National Bureau of Standards (1985).
33. P. de Groot and X. Colonna de Lega, "Interpreting Interferometric Height Measurements Using the Instrument Transfer Function," *Proc. Fringe 2005: 5th International Workshop on Automatic Processing of Fringe Patterns*, 2005.
34. M. A. Davies, C. J. Evans, R. R. Vohra, B. C. Bergner, and S. R. Patterson, "Application of Precision Diamond Machining to the Manufacture of Microphotonics Components," *Proc. SPIE* **5183**:94–108 (2003).
35. N. Taniguchi, "Nanotechnology for Ultraprecision Instruments," *Precision Engineering* **16**:5–24 (1994).

This page intentionally left blank.

DO NOT DUPLICATE

PART

3

TESTING

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

ORTHONORMAL POLYNOMIALS IN WAVEFRONT ANALYSIS

Virendra N. Mahajan*

*The Aerospace Corporation
El Segundo, California*

ABSTRACT

Zernike circle polynomials are in widespread use for wavefront analysis because they are orthogonal over a unit circle and represent balanced classical aberrations for imaging systems with circular pupils. However, they are not suitable for systems with noncircular pupils. Examples of such pupils are annular as in astronomical telescopes, elliptical as in the off-axis pupil of an otherwise rotationally symmetric system with a circular on-axis pupil, hexagonal as in the hexagonal segments of a large telescope, for example, Keck, and rectangular and square as in high-power laser beams. In this chapter, we list the orthonormal circle, annular, elliptical, hexagonal, rectangular, and square polynomials. The polynomials for a noncircular pupil can be obtained by orthogonalizing the circle polynomials over the pupil using the recursive Gram-Schmidt process or a nonrecursive matrix approach. These polynomials are unique in that they are not only orthogonal across such pupils, but also represent balanced classical aberrations for such pupils, just as the Zernike circle polynomials are unique in these respects for circular pupils. The polynomials are given in terms of the circle polynomials as well as in polar and Cartesian coordinates. The orthonormal polynomials for a one-dimensional slit pupil are given as a limiting case of a rectangular pupil. The polynomials corresponding to Seidel aberrations are illustrated isometrically, interferometrically, and with the corresponding point-spread functions (PSFs).

11.1 GLOSSARY

- a half width of a unit rectangular pupil
- a_j j th expansion coefficient
- A area of pupil
- b aspect ratio of a unit elliptical pupil

*The author is also an adjunct professor at the College of Optical Sciences, University of Arizona, Tucson, Arizona and Department of Optics and Photonics, National Central University, Chung Li, Taiwan. He gratefully acknowledges helpful discussions with Drs. Guang-ming Dai and Bill Swantner.

$E_j(x, y)$	orthonormal elliptical polynomial in Cartesian coordinates (x, y)
F	focal ratio of the image-forming light cone
$F_j(x, y)$	j th orthonormal polynomial
$H_j(x, y)$	orthonormal hexagonal polynomial
j	polynomial number
N_n	number of polynomials through an order n
$P_j(x)$	orthonormal slit polynomial along the x axis
$P_n(\cdot)$	Legendre polynomial of order n
$R_j(x, y)$	orthonormal rectangular polynomial
$R_n^m(\rho)$	Zernike circle radial polynomial
$R_n^m(\rho; \epsilon)$	Zernike annular radial polynomial
$S_j(x, y)$	orthonormal square polynomial
$W(x, y)$	wave aberration at a point (x, y)
$Z_j(\rho, \theta)$	orthonormal Zernike circle polynomial in polar coordinates (ρ, θ)
$Z_j(\rho, \theta; \epsilon)$	orthonormal Zernike annular polynomial
σ	standard deviation
σ^2	variance
ϵ	obscuration ratio of an annular pupil

11.2 INTRODUCTION

Optical systems generally have a circular pupil. The imaging elements of such systems have a circular boundary. Hence they also represent circular pupils in fabrication and testing. As a result, the Zernike circle polynomials have been in widespread use since Zernike introduced them in his phase contrast method for testing circular mirrors.¹ They are used in optical design and testing to understand the aberration content of a wavefront. They have also been used for analyzing the wavefront aberration introduced by atmospheric turbulence on a wave propagating through it.² Their utility stems from the fact that they are orthogonal over a unit circle and they represent balanced classical aberrations yielding minimum variance over a circular pupil.³⁻⁶ They are unique in this respect since no other polynomials have these properties. Because of their orthogonality, when a wavefront is expanded in terms of them, the value of an expansion coefficient is independent of the number of polynomials used in the expansion. Hence, one or more polynomial terms can be added or subtracted without affecting the other coefficients. The piston coefficient represents the mean value of the aberration function and the variance of the function is given simply by the sum of the squares of the other expansion coefficients.⁷

For systems with noncircular pupils, the Zernike circle polynomials are neither orthogonal over such pupils nor do they represent balanced aberrations. Hence their special utility is lost. However, since they form a complete set, an aberration function over a noncircular wavefront can be expanded in terms of them. The expansion coefficients are no longer independent of each other and their values change as the number of polynomials used in the expansion changes. The piston coefficient does not represent the mean value of the aberration function, and the sum of the squares of the other coefficients does not yield the aberration variance.

The reflecting telescopes, such as the Hubble, have annular pupils and require polynomials that are orthogonal across an annulus to describe their aberrations.⁸⁻¹¹ The primary mirrors of large telescopes, such as the Keck, consist of hexagonal segments.¹² The wavefront analysis of such segments requires polynomials that are orthogonal over a hexagon. The pupil for off-axis imaging by a system with an axial circular pupil is vignetted, but can be approximated by an ellipse.¹³ When a flat mirror is tested by shining a circular beam on it at some angle (other than normal incidence), the illuminated spot is elliptical. Similarly, the overlap region of two circular wavefronts that are

displaced from each other, as in lateral shearing interferometry¹⁴ or in the calculation of the optical transfer function of a system,¹⁵ can also be approximated by an ellipse. In such cases we need polynomials that are orthogonal over an ellipse. In Refs. 14 and 15, the polynomials that are orthogonal over an elliptical region were obtained simply by scaling the Cartesian coordinates by its aspect ratio. However, such orthogonal polynomials cannot represent classical aberrations. For example, defocus, which varies as ρ^2 , has the same scale for both the x and y coordinates. Similarly, they cannot represent balanced classical aberrations, for example, coma balanced with tilt. High-power laser beams have rectangular or square cross sections¹⁶ and require polynomials that are orthogonal over a rectangle or a square, respectively.

The polynomials orthonormal over a unit annulus, hexagon, ellipse, rectangle, and a square inscribed inside a unit circle may be obtained from the circle polynomials by the recursive Gram-Schmidt orthogonalization process^{17,18} or a nonrecursive matrix approach.¹⁹ The orthonormal polynomials representing balanced aberrations for a slit pupil can be obtained as a limiting case of the rectangular polynomials, where one dimension of the rectangle approaches zero. They are the Legendre polynomials.²⁰ We use the circle polynomials as the basis functions for the orthogonalization process, so that the relationship of a noncircle polynomial to the circle polynomials is evident, since the former is a linear combination of the latter. We give the orthonormal form of the polynomials so that when an aberration function is expanded in terms of them, each expansion coefficient (with the exception of piston) represents the standard deviation of the corresponding expansion term. The noncircle polynomials are given not only in terms of the circle polynomials, but in polar and Cartesian coordinates as well. The circle, annular, hexagonal, and square polynomials are given up to the eighth order, and the elliptical and rectangular polynomials are given up to the fourth order. Just as the Zernike circle polynomials uniquely represent the orthogonal and balanced aberrations across circular pupils, similarly, the orthonormal polynomials for the noncircular pupils given in this chapter also uniquely represent the orthogonal and balanced aberrations across such pupils.

Orthogonal square polynomials were obtained by Bray by orthogonalizing the circle polynomials, but he chose a circle inscribed inside a square instead of the other way around.²¹ Thus his square with a full width of unity has regions that fall outside the unit circle. Defining a unit square in this manner has the disadvantage that the coefficient of a term in a certain polynomial does not represent its peak value. Products of x and y Legendre polynomials,¹⁷ which are orthogonal over a square pupil, have been suggested for analysis of square wavefronts.²² But they do not represent classical or balanced aberrations. For example, defocus is represented by a term in $x^2 + y^2$. While it can be expanded in terms of a complete set of Legendre polynomials, it cannot be represented by a single two-dimensional Legendre polynomial (i.e., as a product of x and y Legendre polynomial). The same difficulty holds for spherical aberration and coma, and the like.

Although in many imaging applications, the amplitude across the pupil is uniform, such is not always the case, for example, a system with an apodized pupil. An example of such a pupil is the Gaussian pupil, where the amplitude has the form of a Gaussian due either to an amplitude filter placed at the pupil or to the wave incident on the pupil being Gaussian, as in the case of a Gaussian laser beam. Again, the balanced aberrations for a Gaussian pupil have a form that is different from the corresponding balanced aberrations for a uniform pupil due to the amplitude weighting of the pupil.^{23–25} The amount of defocus to optimally balance spherical aberration, or the amount of wavefront tilt to optimally balance coma, for example, is different for a Gaussian pupil than its corresponding value for a uniform pupil.

11.3 ORTHONORMAL POLYNOMIALS

In Cartesian coordinates (x, y) , the aberration function $W(x, y)$ for a certain pupil may be expanded in terms of J polynomials $F_j(x, y)$ that are orthonormal over the pupil:²⁶

$$W(x, y) = \sum_{j=1}^J a_j F_j(x, y) \quad (1)$$

where a_j is an expansion or the aberration coefficient of the polynomial $F_j(x, y)$. The orthonormality of the polynomials is represented by

$$\frac{1}{A} \int_{\text{pupil}} F_j(x, y) F_{j'}(x, y) dx dy = \delta_{jj'} \quad (2)$$

where A is the area of the pupil inscribed inside a unit circle, the integration is carried out over the area of the pupil, and $\delta_{jj'}$ is a Kronecker delta. If $F_1 = 1$, then the mean value of each polynomial, except for $j = 1$, is zero, that is,

$$\frac{1}{A} \int_{\text{pupil}} F_j(x, y) dx dy = 0 \quad \text{for } j \neq 1 \quad (3)$$

as may be seen by letting $j' = 1$ in Eq. (2). The aberration coefficients are given by

$$a_j = \frac{1}{A} \int_{\text{pupil}} W(x, y) F_j(x, y) dx dy \quad (4)$$

as may be seen by substituting Eq. (1) into Eq. (4) and using the orthonormality Eq. (2).

The mean and the mean square values of the aberration function are given by

$$\langle W(x, y) \rangle = a_1 \quad (5)$$

and

$$\langle W^2(x, y) \rangle = \sum_{j=1}^J a_j^2 \quad (6)$$

Accordingly, the variance σ^2 of the aberration function is given by

$$\sigma^2 = \langle W^2(x, y) \rangle - \langle W(x, y) \rangle^2 = \sum_{j=2}^J a_j^2 \quad (7)$$

where σ is the standard deviation of the aberration function. The number of polynomials J used in the expansion is a sufficiently large that the variance obtained from Eq. (6) equals the actual value obtained from the function $W(x, y)$ within some prescribed tolerance.

11.4 ZERNIKE CIRCLE POLYNOMIALS

An aberration function $W(\rho, \theta)$, across a *unit circle* can be expanded in terms of the orthonormal *Zernike circle polynomials* $Z_j(\rho, \theta)$ in the form^{2,5}

$$W(\rho, \theta) = \sum_j a_j Z_j(\rho, \theta) \quad (8)$$

where (ρ, θ) are the polar coordinates of a point on the circle, $0 \leq \rho \leq 1$, $0 \leq \theta < 2\pi$, and a_j are the expansion coefficients. The polynomials may be written in the form

$$Z_{\text{even } j}(\rho, \theta) = \sqrt{2(n+1)} R_n^m(\rho) \cos m\theta, m \neq 0 \quad (9a)$$

$$Z_{\text{odd } j}(\rho, \theta) = \sqrt{2(n+1)} R_n^m(\rho) \sin m\theta, m \neq 0 \quad (9b)$$

$$Z_j(\rho, \theta) = \sqrt{n+1} R_n^0(\rho), m = 0 \quad (9c)$$

where n and m are positive integers (including zero) and $n - m \geq 0$ and even. It is evident from Eqs. (9) that the circle polynomials are separable in the polar coordinates ρ and θ . A radial polynomial $R_n^m(\rho)$ is given by

$$R_n^m(\rho) = \sum_{s=0}^{(n-m)/2} \frac{(-1)^s (n-s)!}{s! \left(\frac{n+m}{2} - s\right)! \left(\frac{n-m}{2} - s\right)!} \rho^{n-2s} \quad (10)$$

with a degree n in ρ containing terms in $\rho^n, \rho^{n-2}, \dots$, and ρ^m . It is even or odd in ρ depending on whether n (or m) is even or odd. Also, $R_n^n(\rho) = \rho^n$, $R_n^n(1) = 1$, and $R_n^m(0) = \delta_{m0}$ for even $n/2$ and $-\delta_{m0}$ for odd $n/2$. The polynomials $R_n^m(\rho)$ obey the orthogonality relation

$$\int_0^1 R_n^m(\rho) R_{n'}^{m'}(\rho) \rho d\rho = \frac{1}{2(n+1)} \delta_{nn'} \quad (11)$$

The orthogonality of the angular functions yields

$$\int_0^{2\pi} d\theta \begin{cases} \cos m\theta \cos m'\theta, & j \text{ and } j' \text{ are both even} \\ \cos m\theta \sin m'\theta, & j \text{ is even and } j' \text{ is odd} \\ \sin m\theta \cos m'\theta, & j \text{ is odd and } j' \text{ is even} \\ \sin m\theta \sin m'\theta, & j \text{ and } j' \text{ are both odd} \end{cases} \\ = \begin{cases} \pi(1 + \delta_{m0})\delta_{mm'}, & j \text{ and } j' \text{ are both even} \\ \pi\delta_{mm'}, & j \text{ and } j' \text{ are both odd} \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

Therefore, the Zernike polynomials are orthonormal according to

$$\int_0^1 \int_0^{2\pi} Z_j(\rho, \theta) Z_{j'}(\rho, \theta) \rho d\rho d\theta \bigg/ \int_0^1 \int_0^{2\pi} \rho d\rho d\theta = \delta_{jj'} \quad (13)$$

The expansion coefficients are given by

$$a_j = \frac{1}{\pi} \int_0^1 \int_0^{2\pi} W(\rho, \theta) Z_j(\rho, \theta) \rho d\rho d\theta \quad (14)$$

as may be seen by substituting Eq. (8) into Eq. (14) and using the orthonormality Eq. (13).

While the index n represents the radial *degree* or the *order* of a polynomial, since it represents the highest power of ρ in the polynomial, m is referred to as its *azimuthal frequency*. The index j is a *polynomial-ordering number* and is a function of both n and m . The polynomials are ordered such that an even j corresponds to a symmetric polynomial varying as $\cos m\theta$, while an odd j corresponds to an antisymmetric polynomial varying as $\sin m\theta$. A polynomial with a lower value of n is ordered first, and for a given value of n , a polynomial with a lower value of m is ordered first.

The Zernike circle polynomials are unique in that they are the only polynomials in two variables ρ and θ , which (a) are orthogonal over a circle, (b) are invariant in form with respect to rotation of the coordinate axes about the origin, and (c) include a polynomial for each permissible pair of n and m values.^{4,27}

The orthonormal Zernike circle polynomials and the names associated with some of them when identified with classical aberrations are listed in Table 1a for $n \leq 8$. The polynomials independent of θ are the spherical aberrations, those varying as $\cos\theta$ are the coma aberrations, and those varying as $\cos 2\theta$ are the astigmatism aberrations. The variation of several radial polynomials $R_n^m(\rho)$ with ρ is illustrated in Fig. 1.

TABLE 1a Orthonormal Zernike Circle Polynomials $Z_j(\rho, \theta)$ Ordered Such That an Even j Corresponds to a Symmetric Polynomial Varying as $\cos m\theta$, While an Odd j Corresponds to an Antisymmetric Polynomial Varying as $\sin m\theta$

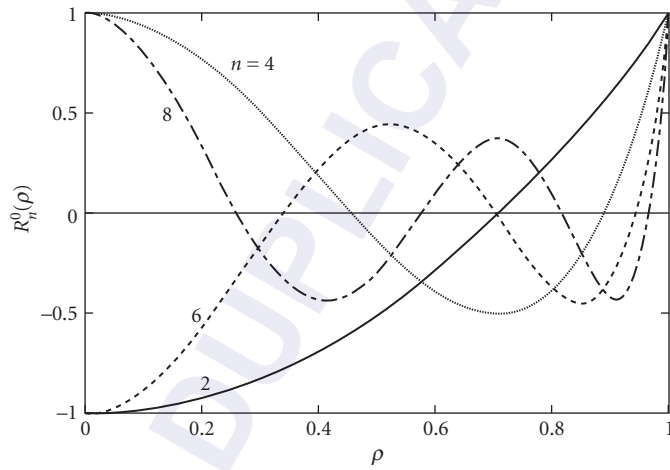
j	n	m	$Z_j(\rho, \theta)$	Aberration Name ^a
1	0	0	1	Piston
2	1	1	$2\rho\cos\theta$	x tilt
3	1	1	$2\rho\sin\theta$	y tilt
4	2	0	$\sqrt{3}(2\rho^2-1)$	Defocus
5	2	2	$\sqrt{6}\rho^2\sin 2\theta$	Primary astigmatism at 45°
6	2	2	$\sqrt{6}\rho^2\cos 2\theta$	Primary astigmatism at 0°
7	3	1	$\sqrt{8}(3\rho^3-2\rho)\sin\theta$	Primary y coma
8	3	1	$\sqrt{8}(3\rho^3-2\rho)\cos\theta$	Primary x coma
9	3	3	$\sqrt{8}\rho^3\sin 3\theta$	
10	3	3	$\sqrt{8}\rho^3\cos 3\theta$	
11	4	0	$\sqrt{5}(6\rho^4-6\rho^2+1)$	Primary spherical aberration
12	4	2	$\sqrt{10}(4\rho^4-3\rho^2)\cos 2\theta$	Secondary astigmatism at 0°
13	4	2	$\sqrt{10}(4\rho^4-3\rho^2)\sin 2\theta$	Secondary astigmatism at 45°
14	4	4	$\sqrt{10}\rho^4\cos 4\theta$	
15	4	4	$\sqrt{10}\rho^4\sin 4\theta$	
16	5	1	$\sqrt{12}(10\rho^5-12\rho^3+3\rho)\cos\theta$	Secondary x coma
17	5	1	$\sqrt{12}(10\rho^5-12\rho^3+3\rho)\sin\theta$	Secondary y coma
18	5	3	$\sqrt{12}(5\rho^5-4\rho^3)\cos 3\theta$	
19	5	3	$\sqrt{12}(5\rho^5-4\rho^3)\sin 3\theta$	
20	5	5	$\sqrt{12}\rho^5\cos 5\theta$	
21	5	5	$\sqrt{12}\rho^5\sin 5\theta$	
22	6	0	$\sqrt{7}(20\rho^6-30\rho^4+12\rho^2-1)$	Secondary spherical aberration
23	6	2	$\sqrt{14}(15\rho^6-20\rho^4+6\rho^2)\sin 2\theta$	Tertiary astigmatism at 45°
24	6	2	$\sqrt{14}(15\rho^6-20\rho^4+6\rho^2)\cos 2\theta$	Tertiary astigmatism at 0°
25	6	4	$\sqrt{14}(6\rho^6-5\rho^4)\sin 4\theta$	
26	6	4	$\sqrt{14}(6\rho^6-5\rho^4)\cos 4\theta$	
27	6	6	$\sqrt{14}6\rho^6\sin 6\theta$	
28	6	6	$\sqrt{14}6\rho^6\cos 6\theta$	
29	7	1	$4(35\rho^7-60\rho^5+30\rho^3-4\rho)\sin\theta$	Tertiary y coma
30	7	1	$4(35\rho^7-60\rho^5+30\rho^3-4\rho)\cos\theta$	Tertiary x coma
31	7	3	$4(21\rho^7-30\rho^5+10\rho^3)\sin 3\theta$	
32	7	3	$4(21\rho^7-30\rho^5+10\rho^3)\cos 3\theta$	
33	7	5	$4(7\rho^7-6\rho^5)\sin 5\theta$	
34	7	5	$4(7\rho^7-6\rho^5)\cos 5\theta$	
35	7	7	$4\rho^7\sin 7\theta$	
36	7	7	$4\rho^7\cos 7\theta$	
37	8	0	$3(70\rho^8-140\rho^6+90\rho^4-20\rho^2+1)$	Tertiary spherical aberration
38	8	2	$\sqrt{18}(56\rho^8-105\rho^6+60\rho^4-10\rho^2)\cos 2\theta$	Quaternary astigmatism at 0°
39	8	2	$\sqrt{18}(56\rho^8-105\rho^6+60\rho^4-10\rho^2)\sin 2\theta$	Quaternary astigmatism at 45°
40	8	4	$\sqrt{18}(28\rho^8-42\rho^6+15\rho^4)\cos 4\theta$	

(Continued)

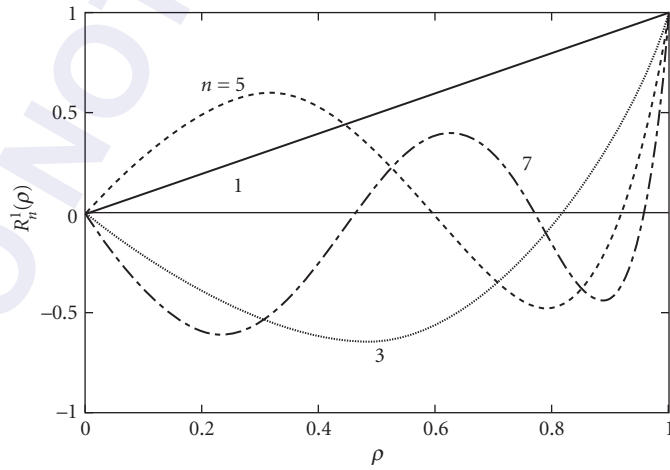
TABLE 1a Orthonormal Zernike Circle Polynomials $Z_j(\rho, \theta)$ Ordered Such That an Even j Corresponds to a Symmetric Polynomial Varying as $\cos m\theta$, While an Odd j Corresponds to an Antisymmetric Polynomial Varying as $\sin m\theta$ (Continued)

j	n	m	$Z_j(\rho, \theta)$	Aberration Name*
41	8	4	$\sqrt{18}(28\rho^8 - 42\rho^6 + 15\rho^4)\sin 4\theta$	
42	8	6	$\sqrt{18}(8\rho^8 - 7\rho^6)\cos 6\theta$	
43	8	6	$\sqrt{18}(8\rho^8 - 7\rho^6)\sin 6\theta$	
44	8	8	$\sqrt{18}\rho^8 \cos 8\theta$	
45	8	8	$\sqrt{18}\rho^8 \sin 8\theta$	

*The words *orthonormal Zernike circle* are to be associated with these names, e.g., *orthonormal Zernike circle primary astigmatism at 0°*.



(a)



(b)

FIGURE 1 Variation of a Zernike circle radial polynomial $R_n^m(\rho)$ with ρ : (a) defocus and spherical aberrations; (b) tilt and coma; and (c) astigmatism.

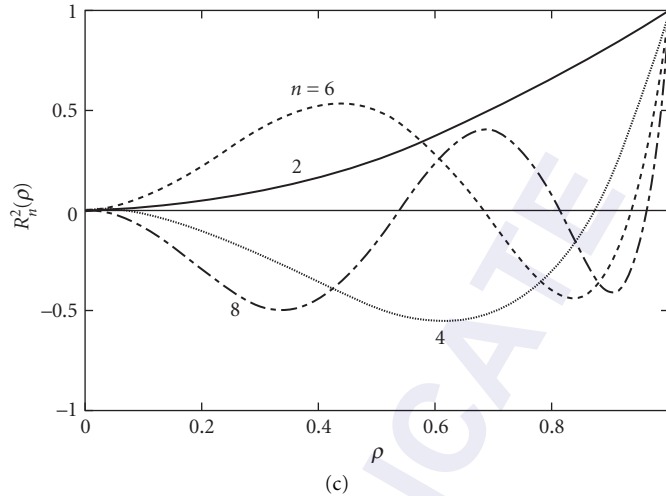


FIGURE 1 (Continued)

The number of polynomials of a given order n is $n + 1$. Their number through a certain order n is given by

$$N_n = (n+1)(n+2)/2 \quad (15)$$

For a rotationally symmetric imaging system, each of the $\sin m\theta$ terms is zero.^{4,28-32} Accordingly the number of polynomials of an even order is $(n/2) + 1$ and $(n + 1)/2$ for an odd order. Their number through an order n is given by

$$N_n = \left(\frac{n}{2} + 1\right)^2 \quad \text{for even } n \quad (16a)$$

$$= (n+1)(n+3)/4 \quad \text{for odd } n \quad (16b)$$

Relationships among the Indices n , m , and j

The number of polynomials N_n through a certain order n represents the largest value of j . Since the number of terms with the same value of n but different values of m is equal to $n + 1$, the smallest value of j for a given value of n is $N_n - n$. For a given value of n and m , there are two j values, $N_n - n + m - 1$ and $N_n - n + m$. The even value of j represents the $\cos m\theta$ term and the odd value of j represents the $\sin m\theta$ term. The value of j with $m = 0$ is $N_n - n$. For example, for $n = 5$, $N_n = 21$, and $j = 21$ represents the $\sin 5\theta$ term. The number of the corresponding $\cos 5\theta$ term is $j = 20$. The two terms with $m = 3$, for example, have j values of 18 and 19 representing the $\cos 3\theta$ and the $\sin 3\theta$ terms, respectively.

For a given value of j , n is given by

$$n = [(2j-1)^{1/2} + 0.5]_{\text{integer}} - 1 \quad (17)$$

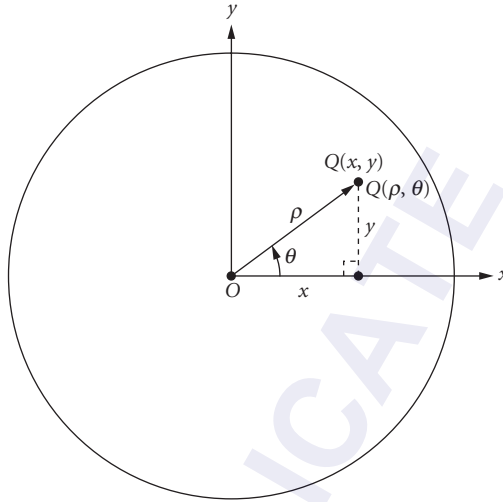


FIGURE 2 Cartesian and polar coordinates (x, y) and (ρ, θ) , respectively, of a point Q in the plane of a unit circle representing the circular exit pupil of an imaging system.

where the subscript integer implies the integer value of the number in brackets. Once n is known, the value of m is given by

$$m = \begin{cases} 2\{[2j+1-n(n+1)]/4\}_{\text{integer}} & \text{when } n \text{ is even} \\ 2\{[2(j+1)-n(n+1)]/4\}_{\text{integer}} - 1 & \text{when } n \text{ is odd} \end{cases} \quad (18a)$$

$$m = \begin{cases} 2\{[2j+1-n(n+1)]/4\}_{\text{integer}} & \text{when } n \text{ is even} \\ 2\{[2(j+1)-n(n+1)]/4\}_{\text{integer}} - 1 & \text{when } n \text{ is odd} \end{cases} \quad (18b)$$

For example, suppose we want to know the values of n and m for the term $j = 10$. From Eq. (17), $n = 3$ and from Eq. (18b), $m = 3$. Hence, it is a $\cos 3\theta$ term.

The polar coordinates (ρ, θ) and the Cartesian coordinates (x, y) of a pupil point Q , as illustrated in Fig. 2, are related to each other according to

$$(x, y) = \rho(\cos\theta, \sin\theta) \quad (19)$$

The circle polynomials in the Cartesian coordinates (x, y) of a pupil point are listed in Table 1b. It is quite common in the optics literature to consider a point object lying along the y axis when imaged by a rotationally symmetric optical system, thus making the yz plane the tangential plane.^{4,28-32} To maintain symmetry of the aberration function about this plane, the polar angle θ of a pupil point is accordingly defined as the angle made by its position vector OQ with the y axis, contrary to the standard convention as the angle with the x axis. We choose a point object along the x axis so that, for example, the coma aberration is expressed as $x(x^2 + y^2)$ and not as $y(x^2 + y^2)$. A positive value of our coma aberration yields a diffraction point spread function that is symmetric about the x axis (or symmetric in y) with its peak and centroid shifted to a positive value of x with respect to the Gaussian image point.

TABLE 1b Orthonormal Zernike Circle Polynomials $Z_j(x, y)$ in Cartesian Coordinates (x, y) , Where $x = \rho \cos \theta$, $y = \rho \sin \theta$, and $0 \leq \rho = \sqrt{x^2 + y^2} \leq 1$

Polynomial	$Z_j(x, y)$
Z_1	1
Z_2	$2x$
Z_3	$2y$
Z_4	$\sqrt{3}(2\rho^2 - 1)$
Z_5	$2\sqrt{6}xy$
Z_6	$\sqrt{6}(x^2 - y^2)$
Z_7	$\sqrt{8}y(3\rho^2 - 2)$
Z_8	$\sqrt{8}x(3\rho^2 - 2)$
Z_9	$\sqrt{8}y(3x^2 - y^2)$
Z_{10}	$\sqrt{8}x(x^2 - 3y^2)$
Z_{11}	$\sqrt{5}(6\rho^4 - 6\rho^2 + 1)$
Z_{12}	$\sqrt{10}(x^2 - y^2)(4\rho^2 - 3)$
Z_{13}	$2\sqrt{10}xy(4\rho^2 - 3)$
Z_{14}	$\sqrt{10}(\rho^4 - 8x^2y^2)$
Z_{15}	$4\sqrt{10}xy(x^2 - y^2)$
Z_{16}	$\sqrt{12}x(10\rho^4 - 12\rho^2 + 3)$
Z_{17}	$\sqrt{12}y(10\rho^4 - 12\rho^2 + 3)$
Z_{18}	$\sqrt{12}x(x^2 - 3y^2)(5\rho^2 - 4)$
Z_{19}	$\sqrt{12}y(3x^2 - y^2)(5\rho^2 - 4)$
Z_{20}	$\sqrt{12}x(16x^4 - 20x^2\rho^2 + 5\rho^4)$
Z_{21}	$\sqrt{12}y(16y^4 - 20y^2\rho^2 + 5\rho^4)$
Z_{22}	$\sqrt{7}(20\rho^6 - 30\rho^4 + 12\rho^2 - 1)$
Z_{23}	$2\sqrt{14}xy(15\rho^2 - 20\rho^2 + 6)$
Z_{24}	$\sqrt{14}(x^2 - y^2)(15\rho^4 - 20\rho^2 + 6)$
Z_{25}	$4\sqrt{14}xy(x^2 - y^2)(6\rho^2 - 5)$
Z_{26}	$\sqrt{14}(8x^4 - 8x^2\rho^2 + \rho^4)(6\rho^2 - 5)$
Z_{27}	$\sqrt{14}xy(32x^4 - 32x^2\rho^2 + 6\rho^4)$
Z_{28}	$\sqrt{14}(32x^6 - 48x^4\rho^2 + 18x^2\rho^4 - \rho^6)$
Z_{29}	$4y(35\rho^6 - 60\rho^4 + 30\rho^2 - 4)$
Z_{30}	$4x(35\rho^6 - 60\rho^4 + 30\rho^2 - 4)$
Z_{31}	$4y(3x^2 - y^2)(21\rho^4 - 30\rho^2 + 10)$
Z_{32}	$4x(x^2 - 3y^2)(21\rho^4 - 30\rho^2 + 10)$
Z_{33}	$4(7\rho^2 - 6)[4x^2y(x^2 - y^2) + y(\rho^4 - 8x^2y^2)]$
Z_{34}	$4(7\rho^2 - 6)[x(\rho^4 - 8x^2y^2) - 4xy^2(x^2 - y^2)]$
Z_{35}	$8x^2y(3\rho^4 - 16x^2y^2) + 4y(x^2 - y^2)(\rho^4 - 16x^2y^2)$
Z_{36}	$4x(x^2 - y^2)(\rho^4 - 16x^2y^2) - 8xy^2(3\rho^4 - 16x^2y^2)$
Z_{37}	$3(70\rho^8 - 140\rho^6 + 90\rho^4 - 20\rho^2 + 1)$
Z_{38}	$\sqrt{18}(56\rho^6 - 105\rho^4 + 60\rho^2 - 10)(x^2 - y^2)$
Z_{39}	$2\sqrt{18}xy(56\rho^6 - 105\rho^4 + 60\rho^2 - 10)$
Z_{40}	$\sqrt{18}(28\rho^4 - 42\rho^2 + 15)(\rho^4 - 8x^2y^2)$
Z_{41}	$4\sqrt{18}xy(28\rho^4 - 42\rho^2 + 15)(x^2 - y^2)$
Z_{42}	$\sqrt{18}(x^2 - y^2)(\rho^4 - 16x^2y^2)(8\rho^2 - 7)$
Z_{43}	$2\sqrt{18}xy(3\rho^4 - 16x^2y^2)$
Z_{44}	$2\sqrt{18}(\rho^4 - 8x^2y^2)^2 - \rho^8$
Z_{45}	$8\sqrt{18}xy(x^2 - y^2)(\rho^4 - 8x^2y^2)$

11.5 ZERNIKE ANNULAR POLYNOMIALS

The aberration function $W(\rho, \theta; \epsilon)$ across a *unit annulus* with an obscuration ratio ϵ , representing the ratio of its inner and outer radii, as illustrated in Fig. 3a, can be expanded in terms of a complete set of *Zernike annular polynomials* $Z_j(\rho, \theta; \epsilon)$ that are orthonormal over the unit annulus in the form⁸⁻¹¹

$$W(\rho, \theta; \epsilon) = \sum_j a_j Z_j(\rho, \theta; \epsilon) \quad (20)$$

where a_j is an expansion coefficient of the polynomial, $\epsilon \leq \rho \leq 1$ and $0 \leq \theta < 2\pi$. The annular polynomials are written in a manner similar to the circle polynomials. Thus

$$Z_{\text{even } j}(\rho, \theta; \epsilon) = \sqrt{2(n+1)} R_n^m(\rho; \epsilon) \cos m\theta, m \neq 0 \quad (21a)$$

$$Z_{\text{odd } j}(\rho, \theta; \epsilon) = \sqrt{2(n+1)} R_n^m(\rho; \epsilon) \sin m\theta, m \neq 0 \quad (21b)$$

$$Z_j(\rho, \theta; \epsilon) = \sqrt{n+1} R_n^0(\rho; \epsilon), m = 0 \quad (21c)$$

where n and m are positive integers (including zero) and $n - m \geq 0$ and even. The radial annular polynomials $R_n^m(\rho; \epsilon)$ obey the orthogonality relation

$$\int_{\epsilon}^1 R_n^m(\rho; \epsilon) R_n^m(\rho; \epsilon) \rho d\rho = \frac{1-\epsilon^2}{2(n+1)} \delta_{nn'} \quad (22)$$

Accordingly, the annular polynomials obey the orthonormality condition

$$\int_{\epsilon}^1 \int_0^{2\pi} Z_j(\rho, \theta; \epsilon) Z_{j'}(\rho, \theta; \epsilon) \rho d\rho d\theta \bigg/ \int_{\epsilon}^1 \int_0^{2\pi} \rho d\rho d\theta = \delta_{jj'} \quad (23)$$

The Zernike expansion coefficients are given by

$$a_j = \frac{1}{\pi(1-\epsilon)^2} \int_{\epsilon}^1 \int_0^{2\pi} W(\rho, \theta; \epsilon) Z_j(\rho, \theta; \epsilon) \rho d\rho d\theta \quad (24)$$

as may be seen by substituting Eq. (20) into Eq. (24) and using Eq. (23) for the orthonormality of the polynomials.

The annular polynomials are similar to the circle polynomials, except that they are orthogonal over an annular pupil. They can be obtained from the circle polynomials by the Gram-Schmidt orthogonalization process.¹⁷ The radial polynomials are accordingly given by

$$R_n^m(\rho; \epsilon) = N_n^m \left[R_n^m(\rho) - \sum_{i \geq 1}^{(n-m)/2} (n-2i+1) \langle R_n^m(\rho) R_{n-2i}^m(\rho; \epsilon) \rangle R_{n-2i}^m(\rho; \epsilon) \right] \quad (25)$$

where

$$\langle R_n^m(\rho) R_n^m(\rho; \epsilon) \rangle = \frac{2}{1-\epsilon^2} \int_{\epsilon}^1 R_n^m(\rho) R_n^m(\rho; \epsilon) \rho d\rho \quad (26)$$

and N_n^m is a normalization constant such that the radial polynomials satisfy the orthogonality Eq. (22). Thus, $R_n^m(\rho; \epsilon)$ is a radial polynomial of degree n in ρ containing terms in $\rho^n, \rho^{n-2}, \dots$, and ρ^m with coefficients that depend on ϵ . The radial polynomials are even or odd in ρ depending on whether n

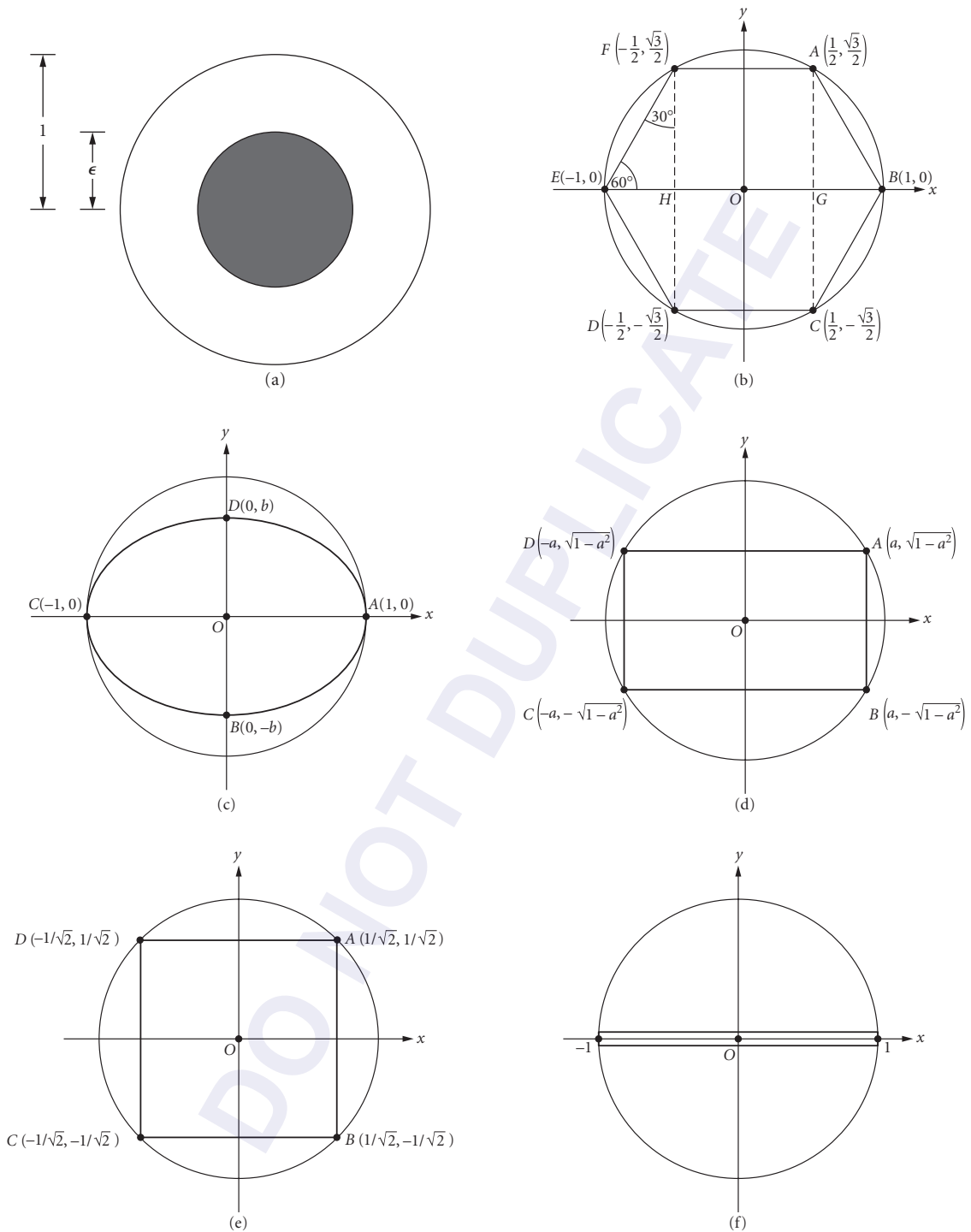


FIGURE 3 Unit pupils inscribed inside a unit circle: (a) annulus of obscuration ratio ϵ ; (b) hexagon; (c) ellipse of aspect ratio b ; (d) rectangle of half width a ; (e) square of half width $1/\sqrt{2}$; and (f) slit.

(or m) is even or odd. For $m = 0$, the radial polynomials are equal to the Legendre polynomials $P_n(\cdot)$ according to

$$R_{2n}^0(\rho; \epsilon) = P_n \left[\frac{2(\rho^2 - \epsilon^2)}{1 - \epsilon^2} - 1 \right] \quad (27)$$

Thus, they can be obtained from the circle radial polynomials $R_{2n}^0(\rho)$ by replacing ρ by $[(\rho^2 - \epsilon^2)/(1 - \epsilon^2)]^{1/2}$, that is,

$$R_{2n}^0(\rho; \epsilon) = R_{2n}^0 \left[\left(\frac{\rho^2 - \epsilon^2}{1 - \epsilon^2} \right)^{1/2} \right] \quad (28)$$

It can be seen from Eqs. (22) and (25) that

$$R_n^n(\rho; \epsilon) = \rho^n / \left(\sum_{i=0}^n \epsilon^{2i} \right)^{1/2} \quad (29)$$

$$= \rho^n \{ (1 - \epsilon^2) / [1 - \epsilon^{2(n+1)}] \}^{1/2} \quad (30)$$

Moreover,

$$R_2^{n-2}(\rho; \epsilon) = \frac{n\rho^n - (n-1)[(1 - \epsilon^{2n})/(1 - \epsilon^{2(n-1)})]\rho^{n-2}}{\{(1 - \epsilon^2)^{-1}[n^2(1 - \epsilon^{2(n+1)}) - (n^2 - 1)(1 - \epsilon^{2n})^2 / (1 - \epsilon^{2(n-1)})]\}^{1/2}} \quad (31)$$

It is evident that the radial polynomial $R_n^m(\rho; \epsilon)$ differs from the corresponding circle polynomial $R_n^m(\rho)$ only in its normalization. We also note that

$$\begin{aligned} R_n^m(1; \epsilon) &= 1, \quad m=0 \\ &\neq 1, \quad m \neq 0 \end{aligned} \quad (32)$$

The variation of several Zernike annular radial polynomials with ρ is shown in Fig. 4 for $\epsilon = 0.5$.

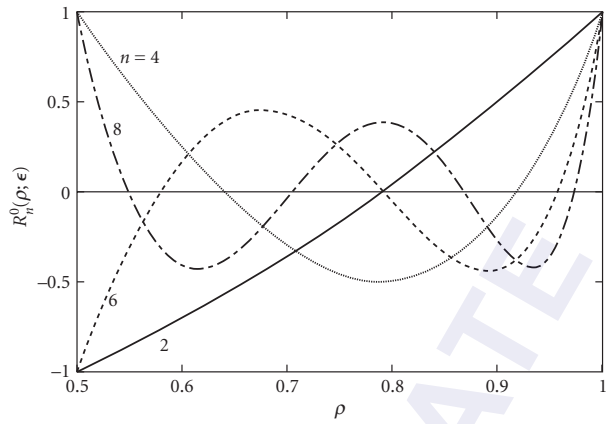
It is evident from Eqs. (21) that the annular polynomials, like the circle polynomials, are separable in the polar coordinates ρ and θ . This is a consequence of the radial symmetry of the annular pupil. As may be evident from the Gram-Schmidt orthogonalization process, each annular polynomial is a linear combination of the circle polynomials.³³ Accordingly, each radial polynomial $R_n^m(\rho; \epsilon)$ can be written as a linear combination of the polynomials $R_n^m(\rho)$, $R_{n-2}^m(\rho)$, \dots , and $R_m^m(\rho)$. For example,

$$R_3^1(\rho; \epsilon) = \frac{1}{(1 - \epsilon^2)(1 + 5\epsilon^2 + 5\epsilon^4 + \epsilon^6)^{1/2}} [(1 + \epsilon^2)R_3^1(\rho) - 2\epsilon^4 R_1^1(\rho)] \quad (33a)$$

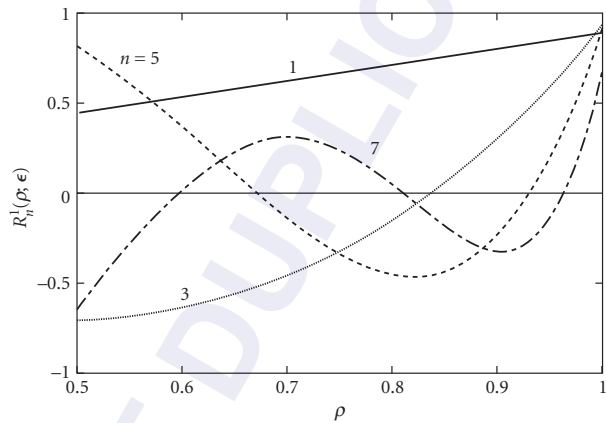
and

$$R_4^0(\rho; \epsilon) = \frac{1}{(1 - \epsilon^2)^2} [R_4^0(\rho) - 3\epsilon^2 R_2^0(\rho) + \epsilon^2(1 + \epsilon^2)R_0^0(\rho)] \quad (33b)$$

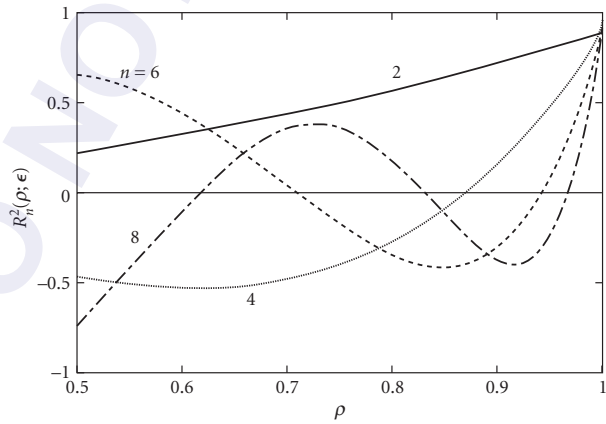
The Zernike annular radial polynomials for $n \leq 8$ are listed in Table 2a. The number polynomials of a certain order or through a certain order n is given by the same expressions as in the case of Zernike circle polynomials. Table 2b lists the full annular polynomials illustrating their ordering. In Table 2c, they are given in the Cartesian coordinates.



(a)



(b)



(c)

FIGURE 4 Variation of a Zernike *annular* radial polynomial $R_n^m(\rho; \epsilon)$ with ρ for $\epsilon = 0.5$: (a) defocus and spherical aberrations; (b) tilt and coma; and (c) astigmatism.

TABLE 2a Zernike Annular Radial Polynomials $R_n^m(\rho; \epsilon)$, Where ϵ Is the Obscuration Ratio of Annular Pupil and $\epsilon \leq \rho \leq 1$

n	m	$R_n^m(\rho; \epsilon)$
0	0	1
1	1	$\rho/(1+\epsilon^2)^{1/2}$
2	0	$(2\rho^2-1-\epsilon^2)/(1-\epsilon^2)$
2	2	$\rho^2/(1+\epsilon^2+\epsilon^2)^{1/2}$
3	1	$\frac{3(1+\epsilon^2)\rho^3-2(1+\epsilon^2+\epsilon^4)\rho}{(1-\epsilon^2)[(1+\epsilon^2)(1+4\epsilon^2+\epsilon^4)]^{1/2}}$
3	3	$\rho^3/(1+\epsilon^2+\epsilon^4+\epsilon^6)^{1/2}$
4	0	$[6\rho^4-6(1+\epsilon^2)\rho^2+1+4\epsilon^2+\epsilon^4]/(1-\epsilon^2)^2$
4	2	$\frac{4\rho^4-3[(1-\epsilon^8)/(1-\epsilon^6)]\rho^2}{\{(1-\epsilon^2)^{-1}[16(1-\epsilon^{10})-15(1-\epsilon^8)^2/(1-\epsilon^6)]^{1/2}\}}$
4	4	$\rho^4/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8)^{1/2}$
5	1	$\frac{10(1+4\epsilon^2+\epsilon^4)\rho^5-12(1+4\epsilon^2+4\epsilon^4+\epsilon^6)\rho^3+3(1+4\epsilon^2+10\epsilon^4+4\epsilon^6+\epsilon^8)\rho}{(1-\epsilon^2)^2[(1+4\epsilon^2+\epsilon^4)(1+9\epsilon^2+9\epsilon^4+\epsilon^6)]^{1/2}}$
5	3	$\frac{5\rho^5-4[(1-\epsilon^{10})/(1-\epsilon^8)]\rho^3}{\{(1-\epsilon^2)^{-1}[25(1-\epsilon^{12})-24(1-\epsilon^{10})^2/(1-\epsilon^8)]^{1/2}\}}$
5	5	$\rho^5/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10})^{1/2}$
6	0	$[20\rho^6-30(1+\epsilon^2)\rho^4+12(1+3\epsilon^2+\epsilon^4)\rho^2-(1+9\epsilon^2+9\epsilon^4+\epsilon^6)]/(1-\epsilon^2)^3$ $15(1+4\epsilon^2+10\epsilon^4+4\epsilon^6+\epsilon^8)\rho^6-20(1+4\epsilon^2+10\epsilon^4+10\epsilon^6+4\epsilon^8+\epsilon^{10})\rho^4$
6	2	$\frac{+6(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+10\epsilon^8+4\epsilon^{10}+\epsilon^{12})\rho^2}{(1+\epsilon^2)^2[(1+4\epsilon^2+10\epsilon^4+4\epsilon^6+\epsilon^8)(1+9\epsilon^2+45\epsilon^4+65\epsilon^6+45\epsilon^8+9\epsilon^{10}+\epsilon^{12})]^{1/2}}$
6	4	$\frac{6\rho^6-5[(1-\epsilon^{12})/(1-\epsilon^{10})]\rho^4}{\{(1-\epsilon^2)^{-1}[36(1-\epsilon^{14})-35(1-\epsilon^{12})^2/(1-\epsilon^{10})]^{1/2}\}}$
6	6	$\rho^6/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12})^{1/2}$
7	1	$a_7^1\rho^7+b_7^1\rho^5+c_7^1\rho^3+d_7^1\rho$
7	3	$a_7^3\rho^7+b_7^3\rho^5+c_7^3\rho^3$
7	5	$\frac{7\rho^7-6[(1-\epsilon^{14})/(1-\epsilon^{12})]\rho^5}{\{(1-\epsilon^2)^{-1}[49(1-\epsilon^{16})-48(1-\epsilon^{14})^2/(1-\epsilon^{12})]^{1/2}\}}$
7	7	$\rho^7/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12}+\epsilon^{14})^{1/2}$
8	0	$\frac{70\rho^8-140(1+\epsilon^2)\rho^6+30(3+8\epsilon^2+3\epsilon^4)\rho^4-20(1+6\epsilon^2+6\epsilon^4+\epsilon^6)\rho^2+\epsilon_8^0}{(1-\epsilon^2)^4}$
8	2	$a_8^2\rho^8+b_8^2\rho^6+c_8^2\rho^4+d_8^2\rho^2$
8	4	$a_8^4\rho^8+b_8^4\rho^6+c_8^4\rho^4$
8	6	$a_8^6\rho^8+b_8^6\rho^6$
8	8	$\rho^8/(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12}+\epsilon^{14}+\epsilon^{16})^{1/2}$

(Continued)

TABLE 2a Zernike Annular Radial Polynomials $R_n^m(\rho; \epsilon)$, Where ϵ Is the Obscuration Ratio of Annular Pupil and $\epsilon \leq \rho \leq 1$ (Continued)

$$\begin{aligned}
 a_7^1 &= 35(1+9\epsilon^2+9\epsilon^4+\epsilon^6)/A_7^1 \\
 b_7^1 &= -60(1+9\epsilon^2+15\epsilon^4+9\epsilon^6+\epsilon^8)/A_7^1 \\
 c_7^1 &= 30(1+9\epsilon^2+25\epsilon^4+25\epsilon^6+9\epsilon^8+\epsilon^{10})/A_7^1 \\
 d_7^1 &= -4(1+9\epsilon^2+45\epsilon^4+65\epsilon^6+45\epsilon^8+9\epsilon^{10}+\epsilon^{12})/A_7^1 \\
 A_7^1 &= (1-\epsilon^2)^3(1+9\epsilon^2+9\epsilon^4+\epsilon^6)^{1/2}(1+16\epsilon^2+36\epsilon^4+16\epsilon^6+\epsilon^8)^{1/2} \\
 a_7^3 &= 21(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+10\epsilon^8+4\epsilon^{10}+\epsilon^{12})/A_7^3 \\
 b_7^3 &= -30(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+20\epsilon^8+10\epsilon^{10}+4\epsilon^{12}+\epsilon^{14})/A_7^3 \\
 c_7^3 &= 10(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+20\epsilon^{10}+10\epsilon^{12}+4\epsilon^{14}+\epsilon^{16})/A_7^3 \\
 A_7^3 &= (1-\epsilon^2)^2(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+10\epsilon^8+4\epsilon^{10}+\epsilon^{12})^{1/2} \\
 &\quad \times (1+9\epsilon^2+45\epsilon^4+165\epsilon^6+270\epsilon^8+270\epsilon^{10}+165\epsilon^{12}+45\epsilon^{14}+9\epsilon^{16}+\epsilon^{18})^{1/2} \\
 e_8^0 &= 1+16\epsilon^2+36\epsilon^4+16\epsilon^6+\epsilon^8 \\
 a_8^2 &= 56(1+9\epsilon^2+45\epsilon^4+65\epsilon^6+45\epsilon^8+9\epsilon^{10}+\epsilon^{12})/A_8^2 \\
 b_8^2 &= -105(1+9\epsilon^2+45\epsilon^4+85\epsilon^6+85\epsilon^8+45\epsilon^{10}+9\epsilon^{12}+\epsilon^{14})/A_8^2 \\
 c_8^2 &= 60(1+9\epsilon^2+45\epsilon^4+115\epsilon^6+150\epsilon^8+115\epsilon^{10}+45\epsilon^{12}+9\epsilon^{14}+\epsilon^{16})/A_8^2 \\
 d_8^2 &= -10(1+9\epsilon^2+45\epsilon^4+165\epsilon^6+270\epsilon^8+270\epsilon^{10}+165\epsilon^{12}+45\epsilon^{14}+9\epsilon^{16}+\epsilon^{18})/A_8^2 \\
 A_8^2 &= (1-\epsilon^2)^3(1+9\epsilon^2+45\epsilon^4+65\epsilon^6+45\epsilon^8+9\epsilon^{10}+\epsilon^{12})^{1/2} \\
 &\quad \times (1+16\epsilon^2+136\epsilon^4+416\epsilon^6+626\epsilon^8+416\epsilon^{10}+136\epsilon^{12}+16\epsilon^{14}+\epsilon^{16})^{1/2} \\
 a_8^4 &= 28(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+20\epsilon^{10}+10\epsilon^{12}+4\epsilon^{14}+\epsilon^{16})/A_8^4 \\
 b_8^4 &= -42(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+35\epsilon^{10}+20\epsilon^{12}+10\epsilon^{14}+4\epsilon^{16}+\epsilon^{18})/A_8^4 \\
 c_8^4 &= 15(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+56\epsilon^{10}+35\epsilon^{12}+20\epsilon^{14}+10\epsilon^{16}+4\epsilon^{18}+\epsilon^{20})/A_8^4 \\
 A_8^4 &= (1-\epsilon^2)^2(1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+20\epsilon^{10}+10\epsilon^{12}+4\epsilon^{14}+\epsilon^{16})^{1/2} \\
 &\quad \times (1+9\epsilon^2+45\epsilon^4+165\epsilon^6+495\epsilon^8+846\epsilon^{10}+994\epsilon^{12}+846\epsilon^{14}+495\epsilon^{16}+165\epsilon^{18}+45\epsilon^{20}+9\epsilon^{22}+\epsilon^{24})^{1/2} \\
 a_8^6 &= 8(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12})/A_8^6 \\
 b_8^6 &= -7(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12}+\epsilon^{14})/A_8^6 \\
 A_8^6 &= (1-\epsilon^2)(1+\epsilon^2+\epsilon^4+\epsilon^6+\epsilon^8+\epsilon^{10}+\epsilon^{12})^{1/2} \\
 &\quad \times (1+4\epsilon^2+10\epsilon^4+20\epsilon^6+35\epsilon^8+56\epsilon^{10}+84\epsilon^{12}+56\epsilon^{14}+35\epsilon^{16}+20\epsilon^{18}+10\epsilon^{20}+4\epsilon^{22}+\epsilon^{24})^{1/2}
 \end{aligned}$$

TABLE 2b Orthonormal Zernike *Annular* Polynomials $Z_j(\rho, \theta; \epsilon)$, Ordered in the Same Manner as the Zernike Circle Polynomials in Table 1a

j	n	m	$Z_j(\rho, \theta; \epsilon)^*$	Aberration Name*
1	0	0	$R_0^0(\rho; \epsilon) = 1$	Piston
2	1	1	$2R_1^1(\rho; \epsilon)\cos\theta$	x tilt
3	1	1	$2R_1^1(\rho; \epsilon)\sin\theta$	y tilt
4	2	0	$\sqrt{3}R_2^0(\rho; \epsilon)$	Defocus
5	2	2	$\sqrt{6}R_2^2(\rho; \epsilon)\sin 2\theta$	Primary astigmatism at 45°
6	2	2	$\sqrt{6}R_2^2(\rho; \epsilon)\cos 2\theta$	Primary astigmatism at 0°
7	3	1	$\sqrt{8}R_3^1(\rho; \epsilon)\sin\theta$	Primary y coma
8	3	1	$\sqrt{8}R_3^1(\rho; \epsilon)\cos\theta$	Primary x coma
9	3	3	$\sqrt{8}R_3^3(\rho; \epsilon)\sin 3\theta$	
10	3	3	$\sqrt{8}R_3^3(\rho; \epsilon)\cos 3\theta$	
11	4	0	$\sqrt{5}R_4^0(\rho; \epsilon)$	Primary spherical aberration
12	4	2	$\sqrt{10}R_4^2(\rho; \epsilon)\cos 2\theta$	Secondary astigmatism at 0°
13	4	2	$\sqrt{10}R_4^2(\rho; \epsilon)\sin 2\theta$	Secondary astigmatism at 45°
14	4	4	$\sqrt{10}R_4^4(\rho; \epsilon)\cos 4\theta$	
15	4	4	$\sqrt{10}R_4^4(\rho; \epsilon)\sin 4\theta$	
16	5	1	$\sqrt{12}R_5^1(\rho; \epsilon)\cos\theta$	Secondary x coma
17	5	1	$\sqrt{12}R_5^1(\rho; \epsilon)\sin\theta$	Secondary y coma
18	5	3	$\sqrt{12}R_5^3(\rho; \epsilon)\cos 3\theta$	
19	5	3	$\sqrt{12}R_5^3(\rho; \epsilon)\sin 3\theta$	
20	5	5	$\sqrt{12}R_5^5(\rho; \epsilon)\cos 5\theta$	
21	5	5	$\sqrt{12}R_5^5(\rho; \epsilon)\sin 5\theta$	
22	6	0	$\sqrt{7}R_6^0(\rho; \epsilon)$	Secondary spherical aberration
23	6	2	$\sqrt{14}R_6^2(\rho; \epsilon)\sin 2\theta$	Tertiary astigmatism at 45°
24	6	2	$\sqrt{14}R_6^2(\rho; \epsilon)\cos 2\theta$	Tertiary astigmatism at 0°
25	6	4	$\sqrt{14}R_6^4(\rho; \epsilon)\cos 4\theta$	
26	6	4	$\sqrt{14}R_6^4(\rho; \epsilon)\sin 4\theta$	
27	6	6	$\sqrt{14}R_6^6(\rho; \epsilon)\sin 6\theta$	
28	6	6	$\sqrt{14}R_6^6(\rho; \epsilon)\cos 6\theta$	
29	7	1	$4R_7^1(\rho; \epsilon)\sin\theta$	
30	7	1	$4R_7^1(\rho; \epsilon)\cos\theta$	
31	7	3	$4R_7^3(\rho; \epsilon)\sin 3\theta$	
32	7	3	$4R_7^3(\rho; \epsilon)\cos 3\theta$	
33	7	5	$4R_7^5(\rho; \epsilon)\sin 5\theta$	
34	7	5	$4R_7^5(\rho; \epsilon)\cos 5\theta$	
35	7	7	$4R_7^7(\rho; \epsilon)\sin 7\theta$	
36	7	7	$4R_7^7(\rho; \epsilon)\cos 7\theta$	
37	8	0	$3R_8^0(\rho; \epsilon)$	Tertiary spherical aberration
38	8	2	$\sqrt{18}R_8^2(\rho; \epsilon)\cos 2\theta$	Quaternary astigmatism at 0°
39	8	2	$\sqrt{18}R_8^2(\rho; \epsilon)\sin 2\theta$	Quaternary astigmatism at 45°
40	8	4	$\sqrt{18}R_8^4(\rho; \epsilon)\cos 4\theta$	
41	8	4	$\sqrt{18}R_8^4(\rho; \epsilon)\sin 4\theta$	
42	8	6	$\sqrt{18}R_8^6(\rho; \epsilon)\cos 6\theta$	
43	8	6	$\sqrt{18}R_8^6(\rho; \epsilon)\sin 6\theta$	
44	8	8	$\sqrt{18}R_8^8(\rho; \epsilon)\cos 8\theta$	
45	8	8	$\sqrt{18}R_8^8(\rho; \epsilon)\sin 8\theta$	

*The words “orthonormal Zernike annular” should be added to the name, e.g., orthonormal Zernike annular primary spherical aberration.

TABLE 2c Orthonormal Zernike *Annular* Polynomials $Z_j(x, y; \epsilon)$ in Cartesian Coordinates (x, y) , Where $x = \rho \cos \theta$, $y = \rho \sin \theta$, and $\epsilon \leq \rho = \sqrt{x^2 + y^2} \leq 1$

Polynomial	$Z_j(x, y; \epsilon)$
Z_1	1
Z_2	$2x/(1 + \epsilon^2)^{1/2}$
Z_3	$2y/(1 + \epsilon^2)^{1/2}$
Z_4	$\sqrt{3}(2\rho^2 - 1 - \epsilon^2)/(1 - \epsilon^2)$
Z_5	$2\sqrt{6}xy/(1 + \epsilon^2 + \epsilon^4)^{1/2}$
Z_6	$\sqrt{6}(x^2 - y^2)/(1 + \epsilon^2 + \epsilon^4)^{1/2}$
Z_7	$\frac{\sqrt{8}y[3(1 + \epsilon^2)\rho^2 - 2(1 + \epsilon^2 + \epsilon^4)]}{(1 - \epsilon^2)[1 + \epsilon^2(1 + 4\epsilon^2 + \epsilon^4)]^{1/2}}$
Z_8	$\frac{\sqrt{8}x[3(1 + \epsilon^2)\rho^2 - 2(1 + \epsilon^2 + \epsilon^4)]}{(1 - \epsilon^2)[1 + \epsilon^2(1 + 4\epsilon^2 + \epsilon^4)]^{1/2}}$
Z_9	$\sqrt{8}y(3x^2 - y^2)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6)^{1/2}$
Z_{10}	$\sqrt{8}x(x^2 - 3y^2)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6)^{1/2}$
Z_{11}	$\sqrt{5}[6\rho^4 - 6(1 + \epsilon^2)\rho^2 + (1 + 4\epsilon^2 + \epsilon^4)]/(1 - \epsilon^2)^2$
Z_{12}	$\frac{\sqrt{10}(x^2 - y^2)[4\rho^2 - 3(1 - \epsilon^8)]/(1 - \epsilon^6)}{\{(1 - \epsilon^2)^{-1}[16(1 - \epsilon^{10}) - 15(1 - \epsilon^8)^2/(1 - \epsilon^{12})]\}^{1/2}}$
Z_{13}	$\frac{2\sqrt{10}xy[4\rho^2 - 3(1 - \epsilon^8)]/(1 - \epsilon^6)}{\{(1 - \epsilon^2)^{-1}[16(1 - \epsilon^{10}) - 15(1 - \epsilon^8)^2/(1 - \epsilon^6)]\}^{1/2}}$
Z_{14}	$\sqrt{10}(\rho^4 - 8x^2y^2)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8)^{1/2}$
Z_{15}	$4\sqrt{10}xy(x^2 - y^2)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8)^{1/2}$
Z_{16}	$\frac{\sqrt{12}x[10(1 + 4\epsilon^2 + \epsilon^4)\rho^4 - 12(1 + 4\epsilon^2 + 4\epsilon^4 + \epsilon^6)\rho^2 + 3(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)]}{(1 - \epsilon^2)^2[(1 + 4\epsilon^2 + \epsilon^4)(1 + 9\epsilon^2 + 9\epsilon^4 + \epsilon^6)]^{1/2}}$
Z_{17}	$\frac{\sqrt{12}y[10(1 + 4\epsilon^2 + \epsilon^4)\rho^4 - 12(1 + 4\epsilon^2 + 4\epsilon^4 + \epsilon^6)\rho^2 + 3(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)]}{(1 - \epsilon^2)^2[(1 + 4\epsilon^2 + \epsilon^4)(1 + 9\epsilon^2 + 9\epsilon^4 + \epsilon^6)]^{1/2}}$
Z_{18}	$\frac{\sqrt{12}x(x^2 - 3y^2)[5\rho^2 - 4(1 - \epsilon^{10})/(1 - \epsilon^8)]}{\{(1 - \epsilon^2)^{-1}[25(1 - \epsilon^{12}) - 24(1 - \epsilon^{10})^2/(1 - \epsilon^8)]\}^{1/2}}$
Z_{19}	$\frac{\sqrt{12}y(3x^2 - y^2)[5\rho^2 - 4(1 - \epsilon^{10})/(1 - \epsilon^8)]}{\{(1 - \epsilon^2)^{-1}[25(1 - \epsilon^{12}) - 24(1 - \epsilon^{10})^2/(1 - \epsilon^8)]\}^{1/2}}$
Z_{20}	$\sqrt{12}x(16x^4 - 20x^2\rho^2 + 5\rho^4)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8 + \epsilon^{10})^{1/2}$
Z_{21}	$\sqrt{12}y(16y^4 - 20y^2\rho^2 + 5\rho^4)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8 + \epsilon^{10})^{1/2}$
Z_{22}	$\sqrt{7}[20\rho^6 - 30(1 + \epsilon^2)\rho^4 + 12(1 + 3\epsilon^2 + \epsilon^4)\rho^2 - (1 + 9\epsilon^2 + 9\epsilon^4 + \epsilon^6)]/(1 - \epsilon^2)^3$
Z_{23}	$\frac{2\sqrt{14}xy[15(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)\rho^4 - 20(1 + 4\epsilon^2 + 10\epsilon^4 + 10\epsilon^6 + 4\epsilon^8 + \epsilon^{10})\rho^2 + 6(1 + 4\epsilon^2 + 10\epsilon^4 + 20\epsilon^6 + 10\epsilon^8 + 4\epsilon^{10} + \epsilon^{12})]}{(1 - \epsilon^2)^2[(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)(1 + 9\epsilon^2 + 45\epsilon^4 + 65\epsilon^6 + 45\epsilon^8 + 9\epsilon^{10} + \epsilon^{12})]^{1/2}}$
Z_{24}	$\frac{\sqrt{14}(x^2 - y^2)[15(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)\rho^4 - 20(1 + 4\epsilon^2 + 10\epsilon^4 + 10\epsilon^6 + 4\epsilon^8 + \epsilon^{10})\rho^2 + 6(1 + 4\epsilon^2 + 10\epsilon^4 + 20\epsilon^6 + 10\epsilon^8 + 4\epsilon^{10} + \epsilon^{12})]}{(1 - \epsilon^2)^2(1 + 4\epsilon^2 + 10\epsilon^4 + 4\epsilon^6 + \epsilon^8)(1 + 9\epsilon^2 + 45\epsilon^4 + 65\epsilon^6 + 45\epsilon^8 + 9\epsilon^{10} + \epsilon^{12})^{1/2}}$

(Continued)

TABLE 2c Orthonormal Zernike Annular Polynomials $Z_j(x, y; \epsilon)$ in Cartesian Coordinates (x, y) , Where $x = \rho \cos \theta$, $y = \rho \sin \theta$, and $\epsilon \leq \rho = \sqrt{x^2 + y^2} \leq 1$ (Continued)

Polynomial	$Z_j(x, y; \epsilon)$
Z_{25}	$\frac{4\sqrt{14}xy(x^2 - y^2)[6\rho^2 - 5(1 - \epsilon^{12})/(1 - \epsilon^{10})]}{\{(1 - \epsilon^2)^{-1}[36(1 - \epsilon^{14}) - 35(1 - \epsilon^{12})^2/(1 - \epsilon^{10})]\}^{1/2}}$
Z_{26}	$\frac{\sqrt{14}(8x^4 - 8x^2\rho^2 + \rho^4)[6\rho^2 - 5(1 - \epsilon^{12})/(1 - \epsilon^{10})]}{\{(1 - \epsilon^2)^{-1}[36(1 - \epsilon^{14}) - 35(1 - \epsilon^{12})^2/(1 - \epsilon^{10})]\}^{1/2}}$
Z_{27}	$\sqrt{14}xy(32x^4 - 32x^2\rho^2 + 6\rho^4)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8 + \epsilon^{10} + \epsilon^{12})^{1/2}$
Z_{28}	$\sqrt{14}(32x^6 - 48x^4\rho^2 + 18x^2\rho^4 - \rho^6)/(1 + \epsilon^2 + \epsilon^4 + \epsilon^6 + \epsilon^8 + \epsilon^{10} + \epsilon^{12})^{1/2}$

11.6 HEXAGONAL POLYNOMIALS

Figure 3b shows a *unit hexagon* inscribed inside a unit circle. Each side of the hexagon has a length of unity. The area of the hexagon is $A = 3\sqrt{3}/2$. The orthonormality of the hexagonal polynomials $H_j(x, y)$ implies that²⁶

$$\frac{2}{3\sqrt{3}} \int_{\text{hexagon}} H_j(x, y) H_{j'}(x, y) dx dy = \delta_{jj'} \quad (34)$$

The orthonormal hexagonal polynomials are given in Tables 3 up to the eighth order in three different but equivalent forms. In Table 3a, each hexagonal polynomial is written in terms of the circle polynomials, thus illustrating the relationship between the two. In particular, it helps determine the potential error made when a hexagonal aberration function is expanded in terms of the circle polynomials.³⁴ The polynomials up to H_{19} are given in their analytical form, but those with $j > 19$ are written in a numerical form because of the increasing complexity of the coefficients of the circle polynomials. In Table 3b, the hexagonal polynomials are given in polar coordinates, showing one-to-one correspondence with the circle polynomials, but illustrating the difference from them. This form is convenient for analytical calculations because of the integration of trigonometric functions over symmetric limits. Finally, in Table 3c, they are given in Cartesian coordinates, as they would be used for any quantitative numerical analysis of, say, an interferogram.

From Table 3a, we note that each hexagonal polynomial consists of cosine or sine terms, but not both. Unlike the circle,³⁻⁶ annular,⁸⁻¹¹ or Gauss^{23,24} polynomials, the hexagonal polynomials are generally not separable in ρ and θ due to the lack of radial symmetry of the hexagonal pupil. The first 13 polynomials, that is, up to H_{13} , are separable, but H_{14} and H_{15} are not; H_{16} through H_{19} are separable, but H_{20} and H_{21} are not. Accordingly, the notion of two indices n and m with dependence on m in the form of $\cos m\theta$ or $\sin m\theta$, as in the case of circle polynomials, loses significance. For example, the Zernike polynomial Z_{14} for $n = 4$ and $m = 4$ varies as $\cos 4\theta$, but H_{14} has a term in $\cos 2\theta$ also. Hence, the hexagonal polynomials can be ordered by a single index only. While the polynomials H_{11} and H_{22} representing the balanced primary and secondary spherical aberrations are radially symmetric, the polynomial H_{37} representing the balanced tertiary spherical aberration is not, since it consists of an angle-dependent term in Z_{28} or $\cos 6\theta$ also. If this term is not included in the polynomial H_{37} , the standard deviation of the aberration increases from a value of unity to 1.3339.

11.7 ELLIPTICAL POLYNOMIALS

Figure 3c shows a *unit ellipse* of an aspect ratio b inscribed inside a unit circle. The semimajor and semiminor axes of the ellipse have lengths of unity and b , respectively. Of course, a unit ellipse is not unique, since b can have any value between 0 and 1. It is represented by an equation

$$x^2 + y^2/b^2 = 1 \quad (35a)$$

TABLE 3a Orthonormal Hexagonal Polynomials H_j in Terms of Zernike Circle Polynomials Z_j

$$\begin{aligned}
H_1 &= Z_1 \\
H_2 &= \sqrt{6/5}Z_2 \\
H_3 &= \sqrt{6/5}Z_3 \\
H_4 &= \sqrt{5/43}Z_1 + (2\sqrt{15/43})Z_4 \\
H_5 &= \sqrt{10/7}Z_5 \\
H_6 &= \sqrt{10/7}Z_6 \\
H_7 &= 16\sqrt{14/11055}Z_3 + 10\sqrt{35/2211}Z_7 \\
H_8 &= 16\sqrt{14/11055}Z_2 + 10\sqrt{35/2211}Z_8 \\
H_9 &= (2\sqrt{5/3})Z_9 \\
H_{10} &= 2\sqrt{35/103}Z_{10} \\
H_{11} &= (521/\sqrt{1072205})Z_1 + 88\sqrt{15/214441}Z_4 + 14\sqrt{43/4987}Z_{11} \\
H_{12} &= 225\sqrt{6/492583}Z_6 + 42\sqrt{70/70369}Z_{12} \\
H_{13} &= 225\sqrt{6/492583}Z_5 + 42\sqrt{70/70369}Z_{13} \\
H_{14} &= -2525\sqrt{14/297774543}Z_6 - (1495\sqrt{70/99258181/3})Z_{12} + (\sqrt{378910/18337/3})Z_{14} \\
H_{15} &= 2525\sqrt{14/297774543}Z_5 + (1495\sqrt{70/99258181/3})Z_{13} + (\sqrt{378910/18337/3})Z_{15} \\
H_{16} &= 30857\sqrt{2/3268147641}Z_2 + (49168/\sqrt{3268147641})Z_8 + 42\sqrt{1474/1478131}Z_{16} \\
H_{17} &= 30857\sqrt{2/3268147641}Z_3 + (49168/\sqrt{3268147641})Z_7 + 42\sqrt{1474/1478131}Z_{17} \\
H_{18} &= 386\sqrt{770/295894589}Z_{10} + 6\sqrt{118965/2872763}Z_{18} \\
H_{19} &= 6\sqrt{10/97}Z_9 + 14\sqrt{5/291}Z_{19} \\
H_{20} &= -0.71499593Z_2 - 0.72488884Z_8 - 0.46636441Z_{16} + 1.72029850Z_{20} \\
H_{21} &= 0.71499594Z_3 + 0.72488884Z_7 + 0.46636441Z_{17} + 1.72029850Z_{21} \\
H_{22} &= 0.58113135Z_1 + 0.89024136Z_4 + 0.89044507Z_{11} + 1.32320623Z_{22} \\
H_{23} &= 1.15667686Z_5 + 1.10775599Z_{13} + 0.43375081Z_{15} + 1.39889072Z_{23} \\
H_{24} &= 1.15667686Z_6 + 1.10775599Z_{12} - 0.43375081Z_{14} + 1.39889072Z_{24} \\
H_{25} &= 1.31832566Z_5 + 1.14465174Z_{13} + 1.94724032Z_{15} + 0.67629133Z_{23} + 1.75496998Z_{25} \\
H_{26} &= -1.31832566Z_6 - 1.14465174Z_{12} + 1.94724032Z_{14} - 0.67629133Z_{24} + 1.75496998Z_{26} \\
H_{27} &= 2\sqrt{77/93}Z_{27} \\
H_{28} &= -1.07362889Z_1 - 1.52546162Z_4 - 1.28216588Z_{11} - 0.70446308Z_{22} + 2.09532473Z_{28} \\
H_{29} &= 0.97998834Z_3 + 1.16162002Z_7 + 1.04573775Z_{17} + 0.40808953Z_{21} + 1.36410394Z_{29} \\
H_{30} &= 0.97998834Z_2 + 1.16162002Z_8 + 1.04573775Z_{16} - 0.40808953Z_{20} + 1.36410394Z_{30} \\
H_{31} &= 3.63513758Z_9 + 2.92084414Z_{19} + 2.11189625Z_{31} \\
H_{32} &= 0.69734874Z_{10} + 0.67589740Z_{18} + 1.22484055Z_{32} \\
H_{33} &= 1.56189763Z_3 + 1.69985309Z_7 + 1.29338869Z_{17} + 2.57680871Z_{21} + 0.67653220Z_{29} \\
&\quad + 1.95719339Z_{33} \\
H_{34} &= -1.56189763Z_2 - 1.69985309Z_8 - 1.29338869Z_{16} + 2.57680871Z_{20} - 0.67653220Z_{30} \\
&\quad + 1.95719339Z_{34} \\
H_{35} &= -1.63832594Z_3 - 1.74759886Z_7 - 1.27572528Z_{17} - 0.77446421Z_{21} - 0.60947360Z_{29} \\
&\quad - 0.36228537Z_{33} + 2.24453237Z_{35} \\
H_{36} &= -1.63832594Z_2 - 1.74759886Z_8 - 1.27572528Z_{16} + 0.77446421Z_{20} - 0.60947360Z_{30} \\
&\quad + 0.36228537Z_{34} + 2.24453237Z_{36} \\
H_{37} &= 0.82154671Z_1 + 1.27988084Z_4 + 1.32912377Z_{11} + 1.11636637Z_{22} - 0.54097038Z_{28} \\
&\quad + 1.37406534Z_{37} \\
H_{38} &= 1.54526522Z_6 + 1.57785242Z_{12} - 0.89280081Z_{14} + 1.28876176Z_{24} - 0.60514082Z_{26} \\
&\quad + 1.43097780Z_{38} \\
H_{39} &= 1.54526522Z_5 + 1.57785242Z_{13} + 0.89280081Z_{15} + 1.28876176Z_{23} + 0.60514082Z_{25} \\
&\quad + 1.43097780Z_{39} \\
H_{40} &= -2.51783502Z_6 - 2.38279377Z_{12} + 3.42458933Z_{14} - 1.69296616Z_{24} + 2.56612920Z_{26} \\
&\quad - 0.85703819Z_{38} + 1.89468756Z_{40} \\
H_{41} &= 2.51783502Z_5 + 2.38279377Z_{13} + 3.42458933Z_{15} + 1.69296616Z_{23} + 2.56612920Z_{25} \\
&\quad + 0.85703819Z_{39} + 1.89468756Z_{41}
\end{aligned}$$

(Continued)

TABLE 3a Orthonormal Hexagonal Polynomials H_j in Terms of Zernike Circle Polynomials Z_j (Continued)

$$H_{42} = -2.72919646Z_1 - 4.02313214Z_4 - 3.69899239Z_{11} - 2.49229315Z_{22} + 4.36717121Z_{28} - 1.13485132Z_{37} + 2.52330106Z_{42}$$

$$H_{43} = 1362\sqrt{77/20334667}Z_{27} + (260/3)\sqrt{341/655957}Z_{43}$$

$$H_{44} = -2.76678413Z_6 - 2.50005278Z_{12} + 1.48041348Z_{14} - 1.62947374Z_{24} + 0.95864121Z_{26} - 0.69034812Z_{38} + 0.40743941Z_{40} + 2.56965299Z_{44}$$

$$H_{45} = -2.76678413Z_5 - 2.50005278Z_{13} - 1.48041348Z_{15} - 1.62947374Z_{23} - 0.95864121Z_{25} - 0.69034812Z_{39} - 0.40743941Z_{41} + 2.56965299Z_{45}$$

TABLE 3b Orthonormal Hexagonal Polynomials $H_j(\rho, \theta)$ in Polar Coordinates

$$H_1 = 1$$

$$H_2 = 2\sqrt{6/5}\rho \cos \theta$$

$$H_3 = 2\sqrt{6/5}\rho \sin \theta$$

$$H_4 = \sqrt{5/43}(-5 + 12\rho^2)$$

$$H_5 = 2\sqrt{15/7}\rho^2 \sin 2\theta$$

$$H_6 = 2\sqrt{15/7}\rho^2 \cos 2\theta$$

$$H_7 = 4\sqrt{42/3685}(-14\rho + 25\rho^3) \sin \theta$$

$$H_8 = 4\sqrt{42/3685}(-14\rho + 25\rho^3) \cos \theta$$

$$H_9 = (4\sqrt{10/3})\rho^3 \sin 3\theta$$

$$H_{10} = 4\sqrt{70/103}\rho^3 \cos 3\theta$$

$$H_{11} = (3/\sqrt{1072205})(737 - 5140\rho^2 + 6020\rho^4)$$

$$H_{12} = (30/\sqrt{492583})(-249\rho^2 + 392\rho^4) \cos 2\theta$$

$$H_{13} = (30/\sqrt{492583})(-249\rho^2 + 392\rho^4) \sin 2\theta$$

$$H_{14} = (10/3)\sqrt{7/99258181}[10(297 - 598\rho^2)\rho^2 \cos 2\theta + 5413\rho^4 \cos 4\theta]$$

$$H_{15} = (10/3)\sqrt{7/99258181}[-10(297 - 598\rho^2)\rho^2 \sin 2\theta + 5413\rho^4 \sin 4\theta]$$

$$H_{16} = 2\sqrt{6/1089382547}(70369\rho - 322280\rho^3 + 309540\rho^5) \cos \theta$$

$$H_{17} = 2\sqrt{6/1089382547}(70369\rho - 322280\rho^3 + 309540\rho^5) \sin \theta$$

$$H_{18} = 4\sqrt{385/295894589}(-3322\rho^3 + 4635\rho^5) \cos 3\theta$$

$$H_{19} = 4\sqrt{5/97}(-22\rho^3 + 35\rho^5) \sin 3\theta$$

$$H_{20} = (-2.17600248\rho + 13.23551876\rho^3 - 16.15533716\rho^5) \cos \theta + 5.95928883\rho^5 \cos 5\theta$$

$$H_{21} = (2.17600248\rho - 13.23551876\rho^3 + 16.15533716\rho^5) \sin \theta + 5.95928883\rho^5 \sin 5\theta$$

$$H_{22} = -2.47059083 + 33.14780774\rho^2 - 93.07966445\rho^4 + 70.01749250\rho^6$$

$$H_{23} = (23.72919095\rho^2 - 90.67126833\rho^4 + 78.51254738\rho^6) \sin 2\theta + 1.37164051\rho^4 \sin 4\theta$$

$$H_{24} = (23.72919095\rho^2 - 90.67126833\rho^4 + 78.51254738\rho^6) \cos 2\theta - 1.37164051\rho^4 \cos 4\theta$$

$$H_{25} = (7.55280798\rho^2 - 36.13018255\rho^4 + 37.95675688\rho^6) \sin 2\theta + (-26.67476754\rho^4 + 39.39897852\rho^6) \sin 4\theta$$

$$H_{26} = (-7.55280798\rho^2 + 36.13018255\rho^4 - 37.95675688\rho^6) \cos 2\theta + (-26.67476754\rho^4 + 39.39897852\rho^6) \cos 4\theta$$

$$H_{27} = 14\sqrt{22/93}\rho^6 \sin 6\theta$$

$$H_{28} = 0.56537219 - 10.44830313\rho^2 + 38.71296332\rho^4 - 37.27668254\rho^6 + 7.83998727\rho^6 \cos 6\theta$$

$$H_{29} = (-15.56917599 + 130.07864353\rho^2 - 288.33220017\rho^4 + 190.97455178\rho^6)\rho \sin \theta + 2.82732724\rho^5 \sin 3\theta + 1.41366362\rho^5 \sin 5\theta$$

$$H_{30} = (-15.56917599 + 130.07864353\rho^2 - 288.33220017\rho^4 + 190.97455178\rho^6)\rho \cos \theta + 2.82732724\rho^5 \cos 3\theta + 1.41366362\rho^5 \cos 5\theta$$

$$H_{31} = (54.28516840 - 202.83704634\rho^2 + 177.39928561\rho^4)\rho^3 \sin 3\theta$$

$$H_{32} = (41.60051295 - 135.27397959\rho^2 + 102.88660624\rho^4)\rho^3 \cos 3\theta$$

$$H_{33} = (-3.87525156 + 41.84243767\rho^2 - 193.65605837\rho^4 + 204.31733848\rho^6)\rho \sin \theta + (76.09262860 - 109.60283027\rho^2)\rho^5 \sin 3\theta + (38.04631430 - 54.80141514\rho^2)\rho^5 \sin 5\theta$$

$$H_{34} = (3.87525156 - 41.84243767\rho^2 + 117.56342977\rho^4 - 94.71450820\rho^6)\rho \cos \theta + (-76.09262860 + 109.60283027\rho^2)\rho^5 \cos 3\theta + (38.04631430 - 54.80141514\rho^2)\rho^5 \cos 5\theta$$

(Continued)

TABLE 3b Orthonormal Hexagonal Polynomials $H_j(\rho, \theta)$ in Polar Coordinates (Continued)

$$H_{35} = (3.10311187 - 34.93479698\rho^2 + 114.10529848\rho^4 - 87.65802721\rho^6)\rho \sin \theta + (12.02405243 - 2.33172188\rho^2)\rho^5 \sin 3\theta + (12.02405243 + 3.68030434\rho^2)\rho^5 \sin 5\theta + 6.01202622\rho^7 \sin 7\theta$$

$$H_{36} = (3.10311187 - 34.93479698\rho^2 + 114.10529848\rho^4 - 87.65802721\rho^6)\rho \cos \theta + (12.02405243 - 2.33172188\rho^2)\rho^5 \cos 3\theta + (12.02405243 + 3.68030434\rho^2)\rho^5 \sin 5\theta + 6.01202622\rho^7 \cos 7\theta$$

$$H_{37} = 2.74530738 - 60.39881618\rho^2 + 300.22087475\rho^4 - 518.03488742\rho^6 + 288.55372176\rho^8 - 2.02412582\rho^6 \cos 6\theta$$

$$H_{38} = (-42.96232789 + 287.78381063\rho^2 - 565.13651608\rho^4 + 339.98298180\rho^4)\rho^2 \cos 2\theta + (8.49786414 - 13.58537785\rho^2)\rho^4 \cos 4\theta$$

$$H_{39} = (-42.96232789 + 287.78381063\rho^2 - 565.13651608\rho^4 + 339.98298180\rho^4)\rho^2 \sin 2\theta + (8.49786414 - 13.58537785\rho^2)\rho^4 \sin 4\theta$$

$$H_{40} = (14.79181046 - 121.61654135\rho^2 + 286.77354559\rho^4 - 203.62188574\rho^6)\rho^2 \cos 2\theta + (83.39879886 - 280.00664075\rho^2 + 225.07739907\rho^4)\rho^4 \cos 4\theta$$

$$H_{41} = (-14.79181046 + 121.61654135\rho^2 - 286.77354559\rho^4 + 203.62188574\rho^6)\rho^2 \sin 2\theta + (83.39879886 - 280.00664075\rho^2 + 225.07739907\rho^4)\rho^4 \sin 4\theta$$

$$H_{42} = -0.84269170 + 24.65387703\rho^2 - 158.21741244\rho^4 + 344.75780000\rho^6 - 238.31877895\rho^8 + (-58.59775991 + 85.64367812\rho^2)\rho^6 \cos 6\theta$$

$$H_{43} = 2\sqrt{22/20334667(-23443 + 32240\rho^2)}\rho^6 \sin 6\theta$$

$$H_{44} = (9.64776957 - 85.41873843\rho^2 + 216.08041438\rho^4 - 164.01834750\rho^6)\rho^2 \cos 2\theta + (12.67622930 - 51.08055822\rho^2 + 48.40133344\rho^4)\rho^4 \cos 4\theta + 10.90211434\rho^8 \cos 8\theta$$

$$H_{45} = (9.64776957 - 85.41873843\rho^2 + 216.08041438\rho^4 - 164.01834750\rho^6)\rho^2 \sin 2\theta - (12.67622930 - 51.08055822\rho^2 + 48.40133344\rho^4)\rho^4 \sin 4\theta + 10.90211434\rho^8 \sin 8\theta$$

TABLE 3c Orthonormal Hexagonal Polynomials $H_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$

$$H_1 = 1$$

$$H_2 = 2\sqrt{6/5}x$$

$$H_3 = 2\sqrt{6/5}y$$

$$H_4 = \sqrt{5/43}(-5 + 12\rho^2)$$

$$H_5 = 4\sqrt{15/7}xy$$

$$H_6 = 2\sqrt{15/7}(x^2 - y^2)$$

$$H_7 = 4\sqrt{42/3685}(-14 + 25\rho^2)y$$

$$H_8 = 4\sqrt{42/3685}(-14 + 25\rho^2)x$$

$$H_9 = (4/3)\sqrt{10}(3x^2y - y^3)$$

$$H_{10} = 4\sqrt{70/103}(x^3 - 3xy^2)$$

$$H_{11} = (3/\sqrt{1072205})(737 - 5140\rho^2 + 6020\rho^4)$$

$$H_{12} = (30/\sqrt{492583})(392\rho^2 - 249)(x^2 - y^2)$$

$$H_{13} = (60/\sqrt{492583})(392\rho^2 - 249)xy$$

$$H_{14} = -(10/3)\sqrt{7/99258181}[567x^4 + 32478x^2y^2 - 11393y^4 - 2970(x^2 - y^2)]$$

$$H_{15} = (40/3)\sqrt{7/99258181}(-1485 + 8403x^2 - 2423y^2)xy$$

$$H_{16} = 2\sqrt{2/3268147641}(211107 - 966840\rho^2 + 928620\rho^4)x$$

$$H_{17} = 2\sqrt{2/3268147641}(211107 - 966840\rho^2 + 928620\rho^4)y$$

$$H_{18} = 4\sqrt{385/295894589}(-3322 + 4635\rho^2)(x^3 - 3xy^2)$$

$$H_{19} = 4\sqrt{5/97}(-22 + 35\rho^2)(3x^2y - y^3)$$

$$H_{20} = (-2.17600247 + 13.23551876\rho^2 - 10.19604832x^4 - 91.90356268x^2y^2 + 13.64110702y^4)x$$

$$H_{21} = (2.17600247 - 13.23551876\rho^2 + 45.95178134x^4 - 27.28221405x^2y^2 + 22.11462599y^4)y$$

$$H_{22} = -2.47059083 + 33.14780774\rho^2 - 93.07966445\rho^4 + 70.01749250\rho^6$$

$$H_{23} = (47.45838189 - 175.85597460x^2 - 186.82909872y^2 + 157.02509476x^4 + 314.05018953x^2y^2 + 157.02509476y^4)xy$$

$$H_{24} = (23.72919094 - 92.04290884x^2 + 78.51254738x^4 + (-23.72919094 + 8.22984309x^2 + 89.29962781y^2 + 78.51254738x^4 - 78.51254738x^2y^2 - 78.51254738y^4)y^2$$

TABLE 3c Orthonormal Hexagonal Polynomials $H_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$ (Continued)

$$\begin{aligned}
 H_{25} &= (15.10561596 - 178.95943525x^2 + 34.43870505y^2 + 233.50942786x^4 \\
 &\quad + 151.82702751x^2y^2 - 81.68240034y^4)xy \\
 H_{26} &= (-7.55280798 + 9.45541501x^2 + 1.44222164x^4)x^2 + (7.55280798 + 160.04860523x^2 \\
 &\quad - 62.80495008y^2 - 234.95164950x^4 - 159.03813574x^2y^2 + 77.35573540y^4)y^2 \\
 H_{27} &= (40.85537039x^4 - 136.18456799x^2y^2 + 40.85537039y^4)xy \\
 H_{28} &= 0.56537219 - 10.44830312\rho^2 + 38.71296332x^4 + 77.42592664x^2y^2 + 38.71296332y^4 \\
 &\quad - 29.43669525x^6 - 229.42985678x^4y^2 + 5.76976155x^2y^4 - 45.11666981y^6 \\
 H_{29} &= (-15.56917599 + 7.06831810x^4 - 14.13663621x^2y^2 + 1.41366362y^4 + 130.07864353\rho^2 \\
 &\quad - 291.15952741\rho^4 + 190.97455178\rho^6)y \\
 H_{30} &= (-15.56917599 - 1.41366362x^4 + 14.13663621x^2y^2 - 7.06831810y^4 + 130.07864353\rho^2 \\
 &\quad - 291.15952741\rho^4 + 190.97455178\rho^6)x \\
 H_{31} &= 162.85550520x^2 - 54.28516840y^2 - 608.51113904x^2\rho^2 + 202.83704634y^2\rho^2 \\
 &\quad + 532.19785685x^2\rho^4 - 177.39928561y^2\rho^4)y \\
 H_{32} &= [(41.60051295 - 135.27397959x^2 + 102.88660624x^4)x^2 + (-124.80153887 + 270.54795919x^2 \\
 &\quad + 405.82193879y^2 - 102.88660624x^4 - 514.43303123x^2y^2 - 308.65981874y^4)y^2]x \\
 H_{33} &= [-3.87525156 + (41.84243767 - 307.79500129x^2 + 368.72158389x^4)x^2 + (41.84243767 \\
 &\quad + 145.33628349x^2 - 155.60974407y^2 + 10.13644892x^4 - 209.06921162x^2y^2 + 149.51592334y^4)y^2]y \\
 H_{34} &= [3.87525156 + (-41.84243767 + 79.51711547x^2 - 39.91309306x^4)x^2 + (-41.84243767 \\
 &\quad + 615.59000259x^2 - 72.66814174y^2 - 777.35626084x^4 - 558.15060029x^2y^2 + 179.29256748y^4)y^2]x \\
 H_{35} &= [3.10311187 + (-34.93479698 + 132.14137712x^2 - 73.19935100x^4)x^2 + (-34.93479698 \\
 &\quad + 144.04222993x^2 + 108.09327226y^2 - 519.49349681x^4 + 23.85771799x^2y^2 - 104.44842531y^4)y^2]y \\
 H_{36} &= [3.10311187 + (-34.93479698 + 96.06921983x^2 - 66.20418535x^4)x^2 + (-34.93479698 \\
 &\quad + 264.28275425x^2 + 72.02111496y^2 - 535.81555000x^4 + 7.53566481x^2y^2 - 97.45325965y^4)y^2]x \\
 H_{37} &= 2.74530738 - 60.39881618\rho^2 + 300.22087475\rho^4 + 288.55372176\rho^8 - 520.05901324x^6 \\
 &\quad - 1523.74277487x^4y^2 - 1584.46654966x^2y^4 - 516.01076159y^6 \\
 H_{38} &= (-42.96232789 + 296.28167478x^2 - 578.72189394x^4 + 339.98298180x^6)x^2 + (42.96232789 \\
 &\quad - 50.98718488x^2 - 279.28594648y^2 - 497.20962679x^4 + 633.06340537x^2y^2 + 551.55113822y^4 \\
 &\quad + 679.96596360x^6 - 679.96596360x^2y^4 - 339.98298180y^6)y^2 \\
 H_{39} &= [-85.92465579 + (541.57616468 - 1075.93152073x^2 + 679.96596360x^4)x^2 + (609.55907786 \\
 &\quad - 2260.54606433x^2 - 1184.61454360y^2 + 2039.89789081x^4 + 2039.89789081x^2y^2 + 679.96596360y^4)y^2]xy \\
 H_{40} &= (14.79181046 - 38.21774249x^2 + 6.76690483x^4 + 21.45551332x^6)x^2 + (-14.79181046 \\
 &\quad - 500.39279319x^2 + 205.01534022y^2 + 1686.80674937x^4 + 1113.25965819x^2y^2 - 566.78018634y^4 \\
 &\quad - 1307.55336779x^6 - 2250.77399075x^4y^2 - 493.06582480x^2y^4 + 428.69928482y^6)y^2 \\
 H_{41} &= [-29.58362093 + (576.82827818 - 1693.57365421x^2 + 1307.55336779x^4)x^2 + (-90.36211274 \\
 &\quad - 1147.09418236x^2 + 546.47947184y^2 + 2122.04091078x^4 + 321.42171817x^2y^2 - 493.06582480y^4)y^2]xy \\
 H_{42} &= -0.84269170 + (24.65387703 - 158.21741244x^2 + 286.16004008x^4 - 152.67510082x^6)x^2 \\
 &\quad + (24.65387703 - 316.43482489x^2 - 158.21741244y^2 + 1913.23979875x^4 + 155.30700127x^2y^2 \\
 &\quad + 403.3555992y^4 - 2152.28660953x^6 - 1429.91267370x^4y^2 + 245.73637792x^2y^4 - 323.96245707y^6)y^2 \\
 H_{43} &= 2\sqrt{22/20334667}(6x^5y - 20x^3y^3 + 6xy^5)(-23443 + 32240\rho^2) \\
 H_{44} &= (9.64776957 - 72.74250912x^2 + 164.99985615x^4 - 104.71489971x^6)x^2 + (-9.64776957 \\
 &\quad - 76.05737585x^2 + 98.09496774y^2 + 471.48320551x^4 + 39.32237674x^2y^2 - 267.16097261y^4 \\
 &\quad - 826.90123032x^6 + 279.13466933x^4y^2 - 170.82784030x^2y^4 + 223.32179529y^6)y^2 \\
 H_{45} &= [19.29553915 + (-221.54239411 + 636.48306167x^2 - 434.42511407x^4)x^2 + (-120.13255963 \\
 &\quad + 864.32165754x^2 + 227.83859586y^2 - 1788.23382186x^4 - 179.98634818x^2y^2 - 221.64827593y^4)y^2]xy
 \end{aligned}$$

or

$$y = \pm b\sqrt{1-x^2} \quad (35b)$$

 Its area is equal to πb . The orthonormality of the elliptical polynomials $E_j(x, y)$ is represented by²⁶

$$\frac{1}{\pi b} \int_{-1}^1 dx \int_{-b\sqrt{1-x^2}}^{b\sqrt{1-x^2}} E_j(x, y) E_{j'}(x, y) dy = \delta_{jj'} \quad (36)$$

The orthonormal elliptical polynomials up to the fourth order are given in Tables 4 in three different but equivalent forms, as in the case of hexagonal polynomials. As in the case of a hexagonal

TABLE 4a Orthonormal *Elliptical* Polynomials E_j in terms of Zernike Circle Polynomials Z_j , Which Reduce to the Corresponding Circle Polynomials as the Aspect Ratio $b \rightarrow 1$

$$\begin{aligned}
E_1 &= Z_1 \\
E_2 &= Z_2 \\
E_3 &= Z_3/b \\
E_4 &= (1/\sqrt{3-2b^2+3b^4})[\sqrt{3}(1-b^2)Z_1+2Z_4] \\
E_5 &= Z_5/b \\
E_6 &= [1/(2\sqrt{2}b^2\sqrt{3-2b^2+3b^4})][-\sqrt{3}(3-4b^2+b^4)Z_1-3(1-b^4)Z_4+\sqrt{2}(3-2b^2+3b^4)Z_6] \\
E_7 &= [1/(b\sqrt{5-6b^2+9b^4})][6(1-b^2)Z_3+2\sqrt{2}Z_7] \\
E_8 &= (2/\sqrt{9-6b^2+5b^4})[(1-b^2)Z_2+\sqrt{2}Z_8] \\
E_9 &= [1/(2\sqrt{2}b^3\sqrt{5-6b^2+9b^4})][2\sqrt{2}(5-8b^2+3b^4)Z_3-(5-2b^2-3b^4)Z_7+(5-6b^2+9b^4)Z_9] \\
E_{10} &= [1/(2\sqrt{2}b^3\sqrt{9-6b^2+5b^4})][2\sqrt{2}(3-4b^2+b^4)Z_2-(3+2b^2-5b^4)Z_8+(9-6b^2+5b^4)Z_{10}] \\
E_{11} &= (1/\alpha)[\sqrt{5}(7-10b^2+3b^4)Z_1+4\sqrt{15}(1-b^2)Z_4-2\sqrt{30}(1-b^2)Z_6+8Z_{11}] \\
E_{12} &= -\sqrt{5/8}b^{-2}(195-475b^2+558b^4-422b^6+159b^8-15b^{10})\beta^{-1}Z_1 \\
&\quad -\sqrt{15/8}b^{-2}(105-205b^2+194b^4-114b^6+5b^8+15b^{10})\beta^{-1}Z_4 \\
&\quad +\sqrt{15/4}(75-155b^2+174b^4-134b^6+55b^8-15b^{10})\beta^{-1}Z_6 \\
&\quad -10\sqrt{2}b^{-2}(3-2b^2+2b^6-3b^8)\beta^{-1}Z_{11}+b^{-2}\alpha\gamma^{-1}Z_{12} \\
E_{13} &= [1/(b\sqrt{5-6b^2+5b^4})][\sqrt{15}(1-b^2)Z_5+2Z_{13}] \\
E_{14} &= (\sqrt{5/2}/4)(1-b^2)^2b^{-4}(35-10b^2-b^4)\gamma^{-1}Z_1+(5\sqrt{15/2}/8)(1-b^2)^2b^{-4}(7+2b^2-b^4)\gamma^{-1}Z_4 \\
&\quad -(\sqrt{15}/8)(35-70b^2+56b^4-26b^6+5b^8)\gamma^{-1}Z_6+[5/(8\sqrt{2})](1-b^2)^2b^{-4}(7+10b^2+7b^4)\gamma^{-1}Z_{11} \\
&\quad -(5/8)b^{-4}(7-6b^2+6b^6-7b^8)\gamma^{-1}Z_{12}+[\gamma/(8b^4)]Z_{14} \\
E_{15} &= -(\sqrt{15/4})b^{-3}(5-8b^2+3b^4)\delta^{-1}Z_5-(5/4)(1-b^4)b^{-3}\beta^{-1}Z_{13}+[\delta/(2b^3)]Z_{15} \\
\alpha &= (45-60b^2+94b^4-60b^6+45b^8)^{1/2} \\
\beta &= (1575-4800b^2+12020b^4-17280b^6+21066b^8-17280b^{10}+12020b^{12}-4800b^{14}+1575b^{16})^{1/2} \\
\gamma &= (35-60b^2+114b^4-60b^6+35b^8)^{1/2} \\
\delta &= (5-6b^2+5b^4)^{1/2}
\end{aligned}$$

TABLE 4b Orthonormal *Elliptical* Polynomials $E_j(\rho, \theta)$ in Polar Coordinates

$$\begin{aligned}
E_1 &= 1 \\
E_2 &= 2\rho \cos \theta \\
E_3 &= (2\rho \sin \theta)/b \\
E_4 &= \sqrt{3}/(3-2b^2+3b^4)(-1-b^2+4\rho^2) \\
E_5 &= (\sqrt{6}/b)\rho^2 \sin 2\theta \\
E_6 &= [1/(2b^2)]\sqrt{6/(3-2b^2+3b^4)}[2b^2(1-b^2)-3(1-b^4)\rho^2+(3-2b^2+3b^4)\rho^2 \cos 2\theta] \\
E_7 &= [4/(b\sqrt{5-6b^2+9b^4})][-(1+3b^2)\rho+6\rho^3] \sin \theta \\
E_8 &= (4/\sqrt{9-6b^2+5b^4})[-(3+b^2)\rho+6\rho^3] \cos \theta \\
E_9 &= [1/(b^3\sqrt{5-6b^2+9b^4})]\{3[4b^2(1-b^2)\rho-(5-2b^2-3b^4)\rho^3] \sin \theta+(5-6b^2+9b^4)\rho^3 \sin 3\theta\} \\
E_{10} &= [1/(b^2\sqrt{9-6b^2+5b^4})]\{3[4b^2(1-b^2)\rho-(3+2b^2-5b^4)\rho^3] \cos \theta+(9-6b^2+5b^4)\rho^3 \cos 3\theta\} \\
E_{11} &= \sqrt{5}[3+2b^2+3b^4-24(1+b^2)\rho^2+48\rho^4-12(1-b^2)\rho^2 \cos 2\theta]\alpha \\
E_{12} &= [\sqrt{10}\alpha/(\gamma b^2)](-3\rho^2+4\rho^4) \cos 2\theta+[\sqrt{5/2}/(2b^2\beta)][-12b^2(5-2b^2+2b^6-5b^8) \\
&\quad +6(15+125b^2-194b^4+194b^6-125b^8-15b^{10})\rho^2+240(-3+2b^2-2b^6+3b^8)\rho^4 \\
&\quad +6(75-155b^2+174b^4-134b^6+55b^8-15b^{10})\rho^2 \cos 2\theta] \\
E_{13} &= (\sqrt{10}/b)\delta^{-1}[-3(1+b^2)\rho^2+8\rho^4] \sin 2\theta \\
E_{14} &= [\sqrt{10}/(8b^4\gamma)]\{3(1-b^2)^2[8b^4-40b^2(1+b^2)\rho^2+5(7+10b^2+7b^4)\rho^4] \\
&\quad +4[6b^2(5-7b^2+7b^4-5b^6)-5(7-6b^2+6b^6-7b^8)\rho^2]\rho^2 \cos 2\theta+(35-60b^2+114b^4 \\
&\quad -60b^6+35b^8)\rho^4 \cos 4\theta\} \\
E_{15} &= (\sqrt{10}/b^3)\delta^{-1}\{[6b^2(1-b^2)-5(1-b^4)\rho^2]\rho^2 \sin 2\theta+[(5-6b^2+5b^4)/2]\rho^4 \sin 4\theta\}
\end{aligned}$$

TABLE 4c Orthonormal *Elliptical* Polynomials $E_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$, $-1 \leq x \leq 1$, and $-\sqrt{1-b^2x^2} \leq y \leq \sqrt{1-b^2x^2}$

$$\begin{aligned}
 E_1 &= 1 \\
 E_2 &= 2x \\
 E_3 &= 2y/b \\
 E_4 &= (\sqrt{3}/\sqrt{3-2b^2+3b^4})(-1-b^2+4\rho^2) \\
 E_5 &= (2\sqrt{6}/b)xy \\
 E_6 &= [\sqrt{3}/(b^2\sqrt{6-4b^2+6b^4})][b^2(1-b^2)+b^2(3b^2-1)x^2-(3-b^2)y^2] \\
 E_7 &= [4/(b\sqrt{5-6b^2+9b^4})][-(1+3b^2)+6\rho^2]y \\
 E_8 &= (4/\sqrt{9-6b^2+5b^4})[-(3+b^2)+6\rho^2]x \\
 E_9 &= [4/(b^3\sqrt{5-6b^2+9b^4})][3b^2(3b^2-1)x^2-(5-3b^2)y^2+3b^2(1-b^2)]y \\
 E_{10} &= [4/(b^2\sqrt{9-6b^2+5b^4})][b^2(5b^2-3)x^2-3(3-b^2)y^2+3b^2(1-b^2)]x \\
 E_{11} &= (\sqrt{5}/\alpha)[48\rho^4-12(3+b^2)x^2-12(1+3b^2)y^2+3+2b^2+3b^4] \\
 E_{12} &= [\sqrt{10}\alpha/(b^2\gamma)][(x^2-y^2)(4\rho^2-3)+[\sqrt{5}/(2\sqrt{2}b^2\beta)][240(-3+2b^2-2b^6+3b^8)\rho^4 \\
 &\quad -60(-9+3b^2+2b^4-6b^6+7b^8+3b^{10})x^2-24(15-70b^2+92b^4-82b^6+45b^8)y^2 \\
 &\quad +12b^2(-5+2b^2-2b^6+5b^8)]] \\
 E_{13} &= [2\sqrt{10}/(b\delta)](8\rho^2-3-3b^2)xy \\
 E_{14} &= [\sqrt{10}/(b^4\gamma)][b^4(3-30b^2+35b^4)x^4+6b^2(5-18b^2+5b^4)x^2y^2+(35-30b^2+3b^4)y^4 \\
 &\quad -6b^4(1-6b^2+5b^4)x^2-6b^2(5-6b^2+b^4)y^2+3b^4(1-b^2)^2] \\
 E_{15} &= [4\sqrt{10}/(b^3\delta)][b^2(5b^2-3)x^2-(5-3b^2)y^2+3b^2(1-b^2)]xy
 \end{aligned}$$

pupil, each elliptical polynomial consists of either cosine or sine terms, but not both. For example, E_6 is a linear combination of Z_6 , Z_4 , and Z_1 . It also shows that the balancing defocus for (zero-degree) Seidel astigmatism is different for an elliptical pupil compared to that for a circular,³⁻⁶ annular,⁸⁻¹¹ or a Gaussian pupil.²³⁻²⁵ Moreover, E_{11} is a linear combination of Z_{11} , Z_6 , Z_4 , and Z_1 . Thus, spherical aberration ρ^4 is balanced with not only defocus ρ^2 but astigmatism $\rho^2 \cos^2\theta$ as well. The elliptical polynomials are generally more complex in that they are made up of a larger number of circle polynomials. These results are a consequence of the fact that the x and y dimensions of the elliptical pupil are not equal. As expected, the elliptical polynomials reduce to the circle polynomials as $b \rightarrow 1$, that is, as the unit ellipse approaches a unit circle.

11.8 RECTANGULAR POLYNOMIALS

Figure 3d shows a *unit rectangle* inscribed inside a unit circle. While the distance of a corner point of the rectangle, such as A , from its center O is unity, the half widths of the rectangle along the x and y axes are a and $\sqrt{1-a^2}$, respectively. Accordingly, the aspect ratio of the rectangle is $\sqrt{1-a^2}/a$, and its area is $4a\sqrt{1-a^2}$. As in the case of a unit ellipse, a unit rectangle is also not unique, since a can have any value between 0 and 1. The orthonormality of the rectangular polynomials $R_j(x, y)$ is represented by²⁶

$$\frac{1}{4a\sqrt{1-a^2}} \int_{-\sqrt{1-a^2}}^{\sqrt{1-a^2}} dy \int_{-a}^a R_j(x, y) R_j(x, y) dx = \delta_{jj} \quad (37)$$

The rectangular polynomials thus obtained up to the fourth order are given in Tables 5 in the same manner as the hexagonal and elliptical polynomials. As in the case of hexagonal and elliptical polynomials, each rectangular polynomial also consists of either cosine or sine terms, but not both. Like the elliptical polynomials, the rectangular polynomials also consist of a larger number of circle polynomials. The rectangular polynomial R_{11} , like the elliptical polynomials E_{11} , representing a balanced primary spherical aberration is not radially symmetric, since it consists of a term in astigmatism Z_6 or $\cos 2\theta$. As discussed below, the rectangular polynomials reduce to the square polynomials as $a \rightarrow 1/\sqrt{2}$, and the slit polynomials for a slit pupil parallel to the x axis as $a \rightarrow 1$.

TABLE 5a Orthonormal Rectangular Polynomials R_j in Terms of Zernike Circle Polynomials Z_j , Which Reduce to the Corresponding Square Polynomials as $a \rightarrow 1/\sqrt{2}$

$$\begin{aligned}
 R_1 &= Z_1 \\
 R_2 &= [\sqrt{3}/(2a)]Z_2 \\
 R_3 &= [\sqrt{3}/(2\sqrt{1-a^2})]Z_3 \\
 R_4 &= [\sqrt{5}/(4\sqrt{1-2a^2+2a^4})](Z_1 + \sqrt{3}Z_4) \\
 R_5 &= [\sqrt{3/2}/(2a\sqrt{1-a^2})]Z_5 \\
 R_6 &= \{\sqrt{5}/[8a^2(1-a^2)\sqrt{1-2a^2+2a^4}]\}[(3-10a^2+12a^4-8a^6)Z_1 + \sqrt{3}(1-2a^2)Z_4 \\
 &\quad + \sqrt{6}(1-2a^2+2a^4)Z_6] \\
 R_7 &= [\sqrt{21}/(4\sqrt{2}\sqrt{27-81a^2+116a^4-62a^6})][\sqrt{2}(1+4a^2)Z_3 + 5Z_7] \\
 R_8 &= [\sqrt{21}/(4\sqrt{2a}\sqrt{35-70a^2+62a^4})][\sqrt{2}(5-4a^2)Z_2 + 5Z_8] \\
 R_9 &= \{\sqrt{5/2}\sqrt{(27-54a^2+62a^4)/(1-a^2)}/[16a^2(27-81a^2+116a^4-62a^6)]\}[2\sqrt{2}(9-36a^2 \\
 &\quad + 52a^4-60a^6)Z_3 + (9-18a^2-26a^4)Z_7 + (27-54a^2+62a^4)Z_9] \\
 R_{10} &= \{\sqrt{5/2}/[16a^3(1-a^2)\sqrt{35-70a^2+62a^4}]\}[2\sqrt{2}(35-112a^2+128a^4-60a^6)Z_2 \\
 &\quad + (35-70a^2+26a^4)Z_8 + (35-70a^2+62a^4)Z_{10}] \\
 R_{11} &= [1/(16\mu)][8(3+4a^2-4a^4)Z_1 + 25\sqrt{3}Z_4 + 10\sqrt{6}(1-2a^2)Z_6 + 21\sqrt{5}Z_{11}] \\
 R_{12} &= \{3\mu/[16a^2\nu\eta]\}[(105-550a^2+1559a^4-2836a^6+2695a^8-1078a^{10})Z_1 \\
 &\quad + 5\sqrt{3}(14-74a^2+205a^4-360a^6+335a^8-134a^{10})Z_4 + (5\sqrt{3}/2)(35-156a^2 \\
 &\quad + 421a^4-530a^6+265a^8)Z_6 + 21\sqrt{5}(1-4a^2+6a^4-4a^6)Z_{11} + [(7/2)\sqrt{5/2}\eta/(1-a^2)]Z_{12}] \\
 R_{13} &= [\sqrt{21}/(16\sqrt{2a}\sqrt{1-3a^2+4a^4-2a^6})](\sqrt{3}Z_5 + \sqrt{5}Z_{13}) \\
 R_{14} &= \tau[6(245-1400a^2+3378a^4-4452a^6+3466a^8-1488a^{10}+496a^{12})Z_1 \\
 &\quad + 15\sqrt{3}(49-252a^2+522a^4-540a^6+270a^8)Z_4 + 15\sqrt{6}(49-252a^2+534a^4-596a^6 \\
 &\quad + 360a^8-144a^{10})Z_6 + 3\sqrt{5}(49-196a^2+282a^4-172a^6+86a^8)Z_{11} \\
 &\quad + 147\sqrt{10}(1-4a^2+6a^4-4a^6)Z_{12} + 3\sqrt{10}\nu^2Z_{14}] \\
 R_{15} &= \{1/[32a^3(1-a^2)(1-3a^2+4a^4-2a^6)^{1/2}]\}[3\sqrt{7/2}(5-18a^2+24a^4-16a^6)Z_5 \\
 &\quad + \sqrt{105/2}(1-2a^2)Z_{13} + \sqrt{210}(1-2a^2+2a^4)Z_{15}] \\
 \mu &= (9-36a^2+103a^4-134a^6+67a^8)^{1/2} \\
 \nu &= (49-196a^2+330a^4-268a^6+134a^8)^{1/2} \\
 \tau &= 1/[128\nu a^4(1-a^2)^2] \\
 \eta &= 9-45a^2+139a^4-237a^6+210a^8-67a^{10}
 \end{aligned}$$

TABLE 5b Orthonormal *Rectangular* Polynomials $R_j(\rho, \theta)$ in Polar Coordinates

$$\begin{aligned}
 R_1 &= 1 \\
 R_2 &= (\sqrt{3}/a)\rho \cos \theta \\
 R_3 &= \sqrt{3/(1-a^2)}\rho \sin \theta \\
 R_4 &= [\sqrt{5}/(2\sqrt{1-2a^2+2a^4})](3\rho^2-1) \\
 R_5 &= [3/(2a\sqrt{1-a^2})]\rho^2 \sin 2\theta \\
 R_6 &= \{\sqrt{5}/[4a^2(1-a^2)\sqrt{1-2a^2+2a^4}]\}[3(1-2a^2+2a^4)\rho^2 \cos 2\theta + 3(1-2a^2)\rho^2 \\
 &\quad - 2a^2(1-a^2)(1-2a^2)] \\
 R_7 &= [\sqrt{21}/(2\sqrt{27-81a^2+116a^4-62a^6})](15\rho^2-9+4a^2)\rho \sin \theta \\
 R_8 &= [\sqrt{21}/(2a\sqrt{35-70a^2+62a^4})](15\rho^2-5-4a^2)\rho \cos \theta \\
 R_9 &= \{\sqrt{5}\sqrt{(27-54a^2+62a^4)/(1-a^2)}/[8a^2(27-81a^2+116a^4-62a^6)]\}\{(27-54a^2+62a^4) \\
 &\quad \times \rho^3 \sin 3\theta - 3[4a^2(3-13a^2+10a^4) - (9-18a^2-26a^4)\rho^2]\rho \sin \theta\} \\
 R_{10} &= \{\sqrt{5}/[8a^3(1-a^2)\sqrt{35-70a^2+62a^4}]\}\{(35-70a^2+62a^4)\rho^3 \cos 3\theta \\
 &\quad - 3[4a^2(7-17a^2+10a^4) - (35-70a^2+26a^4)\rho^2]\rho \cos \theta\} \\
 R_{11} &= [1/(8\mu)][315\rho^4+30(1-2a^2)\rho^2 \cos 2\theta - 240\rho^2+27+16a^2-16a^4] \\
 R_{12} &= [3\mu/(8a^2\nu\eta)][315(1-2a^2)(1-2a^2+2a^4)\rho^4+5(7\mu^2\rho^2-21+72a^2-225a^4+306a^6 \\
 &\quad -153a^8)\rho^2 \cos 2\theta - 15(1-2a^2)(7+4a^2-71a^4+134a^6-67a^8)\rho^2 \\
 &\quad + a^2(1-a^2)(1-2a^2)(70-233a^2+233a^4)] \\
 R_{13} &= [\sqrt{21}/(4a\sqrt{1-3a^2+4a^4-2a^6})](5\rho^2-3)\rho^2 \sin 2\theta \\
 R_{14} &= 6\tau\{5\nu^2\rho^4 \cos 4\theta - 20(1-2a^2)[6a^2(7-16a^2+18a^4-9a^6) - 49(1-2a^2+2a^4)\rho^2]\rho^2 \cos 2\theta \\
 &\quad + 8a^4(1-a^2)^2(21-62a^2+62a^4) - 120a^2(7-30a^2+46a^4-23a^6)\rho^2 \\
 &\quad + 15(49-196a^2+282a^4-172a^6+86a^8)\rho^4\} \\
 R_{15} &= \{\sqrt{21}/[8a^3(1-a^2)^{3/2}(1-2a^2+2a^4)^{1/2}]\}[-(1-2a^2)(6a^2-6a^4-5\rho^2)\rho^2 \sin 2\theta \\
 &\quad + (5/2)(1-2a^2+2a^4)\rho^4 \sin 4\theta]
 \end{aligned}$$

TABLE 5c Orthonormal *Rectangular* Polynomials $R_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$, $-a \leq x \leq a$, and $-\sqrt{1-a^2} \leq y \leq \sqrt{1-a^2}$

$$\begin{aligned}
 R_1 &= 1 \\
 R_2 &= (\sqrt{3}/a)x \\
 R_3 &= \sqrt{3/(1-a^2)}y \\
 R_4 &= [\sqrt{5}/(2\sqrt{1-2a^2+2a^4})](3\rho^2-1) \\
 R_5 &= [3/(a\sqrt{1-a^2})]xy \\
 R_6 &= \{\sqrt{5}/[2a^2(1-a^2)\sqrt{1-2a^2+2a^4}]\}[3(1-a^2)^2x^2 - 3a^4y^2 - a^2(1-3a^2+2a^4)] \\
 R_7 &= [\sqrt{21}/(2\sqrt{27-81a^2+116a^4-62a^6})](15\rho^2-9+4a^2)y \\
 R_8 &= [\sqrt{21}/(2a\sqrt{35-70a^2+62a^4})](15\rho^2-5-4a^2)x \\
 R_9 &= \{\sqrt{5}\sqrt{(27-54a^2+62a^4)/(1-a^2)}/[2a^2(27-81a^2+116a^4-62a^6)]\}[27(1-a^2)^2x^2 \\
 &\quad - 35a^4y^2 - a^2(9-39a^2+30a^4)]y \\
 R_{10} &= \{\sqrt{5}/[2a^3(1-a^2)\sqrt{35-70a^2+62a^4}]\}[35(1-a^2)^2x^2 - 27a^4y^2 - a^2(21-51a^2+30a^4)]x \\
 R_{11} &= [1/(8\mu)][315\rho^4 - 30(7+2a^2)x^2 - 30(9-2a^2)y^2 + 27+16a^2-16a^4] \\
 R_{12} &= [3\mu/(8a^2\nu\eta)][35(1-a^2)^2(18-36a^2+67a^4)x^4 + 630(1-2a^2)(1-2a^2+2a^4)x^2y^2 \\
 &\quad - 35a^4(49-98a^2+67a^4)y^4 - 30(1-a^2)(7-10a^2-12a^4+75a^6-67a^8)x^2 \\
 &\quad - 30a^2(7-77a^2+189a^4-193a^6+67a^8)y^2 + a^2(1-a^2)(1-2a^2)(70-233a^2+233a^4)] \\
 R_{13} &= [\sqrt{21}/(2a\sqrt{1-3a^2+4a^4-2a^6})](5\rho^2-3)xy \\
 R_{14} &= 16\tau\{735(1-a^2)^4x^4 - 540a^4(1-a^2)^2x^2y^2 + 735a^8y^4 - 90a^2(1-a^2)^3(7-9a^2)x^2 \\
 &\quad + 90a^6(1-a^2)(2-9a^2)y^2 + 3a^4(1-a^2)^2(21-62a^2+62a^4)\} \\
 R_{15} &= \{\sqrt{21}/[2a^3(1-a^2)\sqrt{1-3a^2+4a^4-2a^6}]\}[5(1-a^2)^2x^2 - 5a^4y^2 - a^2(3-9a^2+6a^4)]xy
 \end{aligned}$$

11.9 SQUARE POLYNOMIALS

Figure 3e shows a *unit square* inscribed inside a unit circle, as in the case of a rectangle. The distance of a corner point of the square, such as A , from its center O is unity, but each of its sides has a length of $\sqrt{2}$ and its area is 2. The orthonormality of the square polynomials $S_j(x, y)$ is represented by²⁶

$$\frac{1}{2} \int_{-1/\sqrt{2}}^{1/\sqrt{2}} dy \int_{-1/\sqrt{2}}^{1/\sqrt{2}} S_j(x, y) S_{j'}(x, y) dx = \delta_{jj'} \quad (38)$$

The square polynomials through the eighth order are given in terms of the Zernike polynomials in Table 6a. The first 15 polynomials are given in their analytical form, but those with $j > 15$ are written in a numerical form because of the increasing complexity of the coefficients of the circle polynomials. The corresponding polynomials in polar and Cartesian coordinates are given in Tables 6b and 6c, respectively. Of course, up to the fourth order, they can be obtained simply from the rectangular polynomials $R_j(x, y)$ given in Tables 5 by letting $a = 1/\sqrt{2}$. The square polynomial S_{11} representing the balanced primary spherical aberration is radially symmetric, but the polynomial S_{22} representing the balanced secondary spherical aberration is not, since it consists of a term in Z_{14} or $\cos 4\theta$ also. Similarly, the polynomial S_{37} representing the balanced tertiary spherical aberration is also not radially symmetric, since it consists of terms in Z_{14} and Z_{26} both varying as $\cos 4\theta$.

11.10 SLIT POLYNOMIALS

By letting $a \rightarrow 1$ in the rectangular pupil, we obtain a *unit slit* pupil that is parallel to the x axis, as illustrated in Figure 3f. The corresponding orthonormal polynomials representing balanced aberrations for such pupils can be obtained from the rectangular polynomials $R_j(x, y)$ given in Table 5c by letting $y \rightarrow 0$ and $a \rightarrow 1$. Half of the rectangular polynomials thus reduce to zero. Some of the other polynomials are redundant. For example, the one-dimensional defocus and astigmatism can not be distinguished from each other. The slit polynomials are orthonormal according to²⁶

$$\frac{1}{2} \int_{-1}^1 P_j(x) P_{j'}(x) dx = \delta_{jj'} \quad (39)$$

The relevant orthonormal slit polynomials are listed in Table 7. They are the Legendre polynomials,¹⁷ which represent balanced aberrations uniquely.^{20,26} Since the pupil is one dimensional along the x axis, the aberrations vary with x only.

11.11 ABERRATION BALANCING AND TOLERANCING, AND DIFFRACTION FOCUS

For small aberrations, the Strehl ratio of the image of a point object is approximately given by $1 - \sigma^2$ or $\exp(-\sigma^2)$ when the standard deviation σ of the aberration is in units of radians.^{4,5,35} The Zernike circle and annular polynomials are separable in ρ and θ . The balanced spherical aberrations for these radially symmetric pupils are radially symmetric, and the balanced primary astigmatism for them has the same form. This is also true of a Gaussian circular or annular pupil, again because of the radial symmetry of the pupil and the amplitude across it.²³⁻²⁵ From the orthonormal form H_4 of defocus for a hexagonal pupil, the sigma of the defocus aberration ρ^2 is given by $(1/12)\sqrt{43/5}$. The hexagonal polynomials H_5 and H_6 show that the balanced astigmatism has the same form as the circle polynomials Z_5 and Z_6 , respectively. Thus the relative amount of defocus ρ^2 that balances classical or Seidel astigmatism $\rho^2 \cos^2 \theta$ is the same for a hexagonal pupil as for a circular pupil. Hence, for

TABLE 6a Orthonormal Square Polynomials S_j in Terms of Zernike Circle Polynomials Z_j

$$\begin{aligned}
 S_1 &= Z_1 \\
 S_2 &= \sqrt{3/2}Z_2 \\
 S_3 &= \sqrt{3/2}Z_3 \\
 S_4 &= (\sqrt{5/2/2})Z_1 + (\sqrt{15/2/2})Z_4 \\
 S_5 &= \sqrt{3/2}Z_5 \\
 S_6 &= (\sqrt{15/2})Z_6 \\
 S_7 &= (3\sqrt{21/31/2})Z_3 + (5\sqrt{21/62/2})Z_7 \\
 S_8 &= (3\sqrt{21/31/2})Z_2 + (5\sqrt{21/62/2})Z_8 \\
 S_9 &= -7(\sqrt{5/31/2})Z_3 - (13\sqrt{5/62/4})Z_7 + (\sqrt{155/2/4})Z_9 \\
 S_{10} &= (7\sqrt{5/31/2})Z_2 + (13\sqrt{5/62/4})Z_8 + (\sqrt{155/2/4})Z_{10} \\
 S_{11} &= (8/\sqrt{67})Z_1 + (25\sqrt{3/67/4})Z_4 + (21\sqrt{5/67/4})Z_{11} \\
 S_{12} &= (45\sqrt{3/16})Z_6 + (21\sqrt{5/16})Z_{12} \\
 S_{13} &= (3\sqrt{7/8})Z_5 + (\sqrt{105/8})Z_{13} \\
 S_{14} &= 261/(8\sqrt{134})Z_1 + (345\sqrt{3/134/16})Z_4 + (129\sqrt{5/134/16})Z_{11} + (3\sqrt{335/16})Z_{14} \\
 S_{15} &= (\sqrt{105/4})Z_{15} \\
 S_{16} &= 1.71440511Z_2 + 1.71491497Z_8 + 0.65048499Z_{10} + 1.52093102Z_{16} \\
 S_{17} &= 1.71440511Z_3 + 1.71491497Z_7 - 0.65048449Z_9 + 1.52093102Z_{17} \\
 S_{18} &= 4.10471345Z_2 + 3.45884077Z_8 + 5.34411808Z_{10} + 1.51830574Z_{16} + 2.80808005Z_{18} \\
 S_{19} &= -4.10471345Z_3 - 3.45884078Z_7 + 5.34411808Z_9 - 1.51830575Z_{17} + 2.80808005Z_{19} \\
 S_{20} &= 5.57146696Z_2 + 4.44429264Z_8 + 3.00807599Z_{10} + 1.70525179Z_{16} + 1.16777987Z_{18} \\
 &\quad + 4.19716701Z_{20} \\
 S_{21} &= 5.57146696Z_3 + 4.44429264Z_7 - 3.00807599Z_9 + 1.70525179Z_{17} - 1.16777988Z_{19} \\
 &\quad + 4.19716701Z_{21} \\
 S_{22} &= 1.33159935Z_1 + 1.94695912Z_4 + 1.74012467Z_{11} + 0.65624211Z_{14} + 1.50989174Z_{22} \\
 S_{23} &= 0.95479991Z_5 + 1.01511643Z_{13} + 1.28689496Z_{23} \\
 S_{24} &= 9.87992565Z_6 + 7.28853095Z_{12} + 3.38796312Z_{24} \\
 S_{25} &= 5.61978925Z_{15} + 2.84975327Z_{25} \\
 S_{26} &= 11.00650275Z_1 + 14.00366597Z_4 + 9.22698484Z_{11} + 13.55765720Z_{14} + 3.18799971Z_{22} \\
 &\quad + 5.11045000Z_{26} \\
 S_{27} &= 4.24396143Z_5 + 2.70990074Z_{13} + 0.84615108Z_{23} + 5.17855026Z_{27} \\
 S_{28} &= 17.58672314Z_6 + 11.15913268Z_{12} + 3.57668869Z_{24} + 6.44185987Z_{28} \\
 S_{29} &= 2.42764289Z_3 + 2.69721906Z_7 - 1.56598064Z_9 + 2.12208902Z_{17} - 0.93135653Z_{19} \\
 &\quad + 0.25252773Z_{21} + 1.59017528Z_{29} \\
 S_{30} &= 2.42764289Z_2 + 2.69721906Z_8 + 1.56598064Z_{10} + 2.12208902Z_{16} + 0.93135653Z_{18} \\
 &\quad + 0.25252773Z_{20} + 1.59017528Z_{30} \\
 S_{31} &= -9.10300982Z_3 - 8.79978208Z_7 + 10.69381427Z_9 - 5.37383385Z_{17} + 7.01044701Z_{19} \\
 &\quad - 1.26347272Z_{21} - 1.90131756Z_{29} + 3.07960207Z_{31} \\
 S_{32} &= 9.10300982Z_2 + 8.79978208Z_8 + 10.69381427Z_{10} + 5.37383385Z_{16} + 7.01044701Z_{18} \\
 &\quad + 1.26347272Z_{20} + 1.90131756Z_{30} + 3.07960207Z_{32} \\
 S_{33} &= 21.39630883Z_3 + 19.76696884Z_7 - 12.70550260Z_9 + 11.05819453Z_{17} - 7.02178756Z_{19} \\
 &\quad + 15.80286172Z_{21} + 3.29259996Z_{29} - 2.07602718Z_{31} + 5.40902889Z_{33} \\
 S_{34} &= 21.39630883Z_2 + 19.76696884Z_8 + 12.70550260Z_{10} + 11.05819453Z_{16} + 7.02178756Z_{18} \\
 &\quad + 15.80286172Z_{20} + 3.29259996Z_{30} + 2.07602718Z_{32} + 5.40902889Z_{34} \\
 S_{35} &= -16.54454462Z_3 - 14.89205549Z_7 + 22.18054997Z_9 - 7.94524849Z_{17} + 11.85458952Z_{19} \\
 &\quad - 6.18963457Z_{21} - 2.19431441Z_{29} + 3.24324400Z_{31} - 1.72001172Z_{33} + 8.16384008Z_{35} \\
 S_{36} &= 16.54454462Z_2 + 14.89205549Z_8 + 22.18054997Z_{10} + 7.94524849Z_{16} + 11.85458952Z_{18} \\
 &\quad + 6.18963457Z_{20} + 2.19431441Z_{30} + 3.24324400Z_{32} + 1.72001172Z_{34} + 8.16384008Z_{36}
 \end{aligned}$$

(Continued)

TABLE 6a Orthonormal Square Polynomials $S_j(\rho, \theta)$ in Terms of Zernike Circle Polynomials
(Continued)

$$S_{37} = 1.75238960Z_1 + 2.72870567Z_4 + 2.76530671Z_{11} + 1.43647360Z_{14} + 2.12459170Z_{22}$$

$$+ 0.92450043Z_{26} + 1.58545010Z_{37}$$

$$S_{38} = 19.24848143Z_6 + 16.41468913Z_{12} + 9.76776798Z_{24} + 1.47438007Z_{28} + 3.83118509Z_{38}$$

$$S_{39} = 0.46604820Z_5 + 0.84124290Z_{13} + 1.00986774Z_{23} - 0.42520747Z_{27} + 1.30579570Z_{39}$$

$$S_{40} = 28.18104531Z_1 + 38.52219208Z_4 + 30.18363661Z_{11} + 36.44278147Z_{14} + 15.52577202Z_{22}$$

$$+ 19.21524879Z_{26} + 4.44731721Z_{37} + 6.00189814Z_{40}$$

$$S_{41} = (369/4)\sqrt{35/3574}Z_{15} + [11781/(32\sqrt{3574})]Z_{25} + (2145/32)\sqrt{7/3574}Z_{41}$$

$$S_{42} = 85.33469748Z_6 + 64.01249391Z_{12} + 30.59874671Z_{24} + 34.09158819Z_{28} + 7.75796322Z_{38}$$

$$+ 9.37150432Z_{42}$$

$$S_{43} = 14.30642479Z_5 + 11.17404702Z_{13} + 5.68231935Z_{23} + 18.15306055Z_{27} + 1.54919583Z_{39}$$

$$+ 5.90178984Z_{43}$$

$$S_{44} = 36.12567424Z_1 + 47.95305224Z_4 + 35.30691679Z_{11} + 56.72014548Z_{14} + 16.36470429Z_{22}$$

$$+ 26.32636277Z_{26} + 3.95466397Z_{37} + 6.33853092Z_{40} + 12.38056785Z_{44}$$

$$S_{45} = 21.45429746Z_{15} + 9.94633083Z_{25} + 2.34632890Z_{41} + 10.39130049Z_{45}$$

TABLE 6b Orthonormal Square Polynomials $S_j(\rho, \theta)$ in Polar Coordinates

$$S_1 = 1$$

$$S_2 = \sqrt{6}\rho \cos \theta$$

$$S_3 = \sqrt{6}\rho \sin \theta$$

$$S_4 = \sqrt{5/2}(3\rho^2 - 1)$$

$$S_5 = 3\rho^2 \sin 2\theta$$

$$S_6 = 3\sqrt{5/2}\rho^2 \cos 2\theta$$

$$S_7 = \sqrt{21/31}(15\rho^2 - 7)\rho \sin \theta$$

$$S_8 = \sqrt{21/31}(15\rho^2 - 7)\rho \cos \theta$$

$$S_9 = (\sqrt{5/31}/2)[31\rho^3 \sin 3\theta - 3(13\rho^2 - 4)\rho \sin \theta]$$

$$S_{10} = (\sqrt{5/31}/2)[31\rho^3 \cos 3\theta + 3(13\rho^2 - 4)\rho \cos \theta]$$

$$S_{11} = [1/(2\sqrt{67})](315\rho^4 - 240\rho^2 + 31)$$

$$S_{12} = [15/(2\sqrt{2})](7\rho^2 - 3)\rho^2 \cos 2\theta$$

$$S_{13} = \sqrt{21/2}(5\rho^2 - 3)\rho^2 \sin 2\theta$$

$$S_{14} = [3/(8\sqrt{134})](335\rho^4 \cos 4\theta + 645\rho^4 - 300\rho^2 + 22)$$

$$S_{15} = (5/2)\sqrt{21/2}\rho^4 \sin 4\theta$$

$$S_{16} = \sqrt{55/1966}[11\rho^3 \cos 3\theta + 3(19 - 97\rho^2 + 105\rho^4)\rho \cos \theta]$$

$$S_{17} = \sqrt{55/1966}[-11\rho^3 \sin 3\theta + 3(19 - 97\rho^2 + 105\rho^4)\rho \sin \theta]$$

$$S_{18} = (1/4)\sqrt{3/844397}[5(-10099 + 20643\rho^2)\rho^3 \cos 3\theta + 3(3128 - 23885\rho^2 + 37205\rho^4)\rho \cos \theta]$$

$$S_{19} = (1/4)\sqrt{3/844397}[5(-10099 + 20643\rho^2)\rho^3 \sin 3\theta - 3(3128 - 23885\rho^2 + 37205\rho^4)\rho \sin \theta]$$

$$S_{20} = (1/16)\sqrt{7/859}[2577\rho^5 \cos 5\theta - 5(272 - 717\rho^2)\rho^3 \cos 3\theta + 30(22 - 196\rho^2 + 349\rho^4)\rho \cos \theta]$$

$$S_{21} = (1/16)\sqrt{7/859}[2577\rho^5 \sin 5\theta + 5(272 - 717\rho^2)\rho^3 \sin 3\theta + 30(22 - 196\rho^2 + 349\rho^4)\rho \sin \theta]$$

$$S_{22} = (1/4)\sqrt{65/849}[1155\rho^6 + 30\rho^4 \cos 4\theta - 1395\rho^4 + 453\rho^2 - 31]$$

$$S_{23} = (1/2)\sqrt{33/3923}[471 - 1820\rho^2 + 1575\rho^4]\rho^2 \sin 2\theta$$

$$S_{24} = (21/4)\sqrt{65/1349}(27 - 140\rho^2 + 165\rho^4)\rho^2 \cos 2\theta$$

$$S_{25} = (7/4)\sqrt{33/2}(9\rho^2 - 5)\rho^4 \sin 4\theta$$

$$S_{26} = [1/(16\sqrt{849})][5(-98 + 2418\rho^2 - 12051\rho^4 + 15729\rho^6) + 3(-8195 + 17829\rho^2)\rho^4 \cos 4\theta]$$

$$S_{27} = [1/(16\sqrt{7846})][27461\rho^6 \sin 6\theta + 15(348 - 2744\rho^2 + 4487\rho^4)\rho^2 \sin 2\theta]$$

$$S_{28} = [21/(32\sqrt{1349})][1349\rho^6 \cos 6\theta + 5(196 - 1416\rho^2 + 2247\rho^4)\rho^2 \cos 2\theta]$$

(Continued)

TABLE 6b Orthonormal Square Polynomials $S_j(\rho, \theta)$ in Polar Coordinates (Continued)

$$S_{29} = (-13.79189793\rho + 125.49411319\rho^3 - 308.13074909\rho^5 + 222.62454035\rho^7) \sin \theta$$

$$+ (8.47599260\rho^3 - 16.13156842\rho^5) \sin 3\theta + 0.87478174\rho^5 \sin 5\theta$$

$$S_{30} = (-13.79189793\rho + 125.49411319\rho^3 - 308.13074909\rho^5 + 222.62454035\rho^7) \cos \theta$$

$$+ (-8.47599260\rho^3 + 16.13156842\rho^5) \cos 3\theta + 0.87478174\rho^5 \cos 5\theta$$

$$S_{31} = (6.14762642\rho - 79.44065626\rho^3 + 270.16115026\rho^5 - 266.18445920\rho^7) \sin \theta$$

$$+ (56.29115383\rho^3 - 248.12774426\rho^5 + 258.68657393\rho^7) \sin 3\theta - 4.37679791\rho^5 \sin 5\theta$$

$$S_{32} = (-6.14762642\rho + 79.44065626\rho^3 - 270.16115026\rho^5 + 266.18445920\rho^7) \cos \theta$$

$$+ (56.29115383\rho^3 - 248.12774426\rho^5 + 258.68657393\rho^7) \cos 3\theta + 4.37679791\rho^5 \cos 5\theta$$

$$S_{33} = (-6.78771487\rho + 103.15977419\rho^3 - 407.15689696\rho^5 + 460.96399558\rho^7) \sin \theta$$

$$+ (-21.68093294\rho^3 + 127.50233381\rho^5 - 174.38628345\rho^7) \sin 3\theta$$

$$+ (-75.07397471\rho^5 + 151.45280913\rho^7) \sin 5\theta$$

$$S_{34} = (-6.78771487\rho + 103.15977419\rho^3 - 407.15689696\rho^5 + 460.96399558\rho^7) \cos \theta$$

$$+ (21.68093294\rho^3 - 127.50233381\rho^5 + 174.38628345\rho^7) \cos 3\theta$$

$$+ \rho^5(-75.07397471 + 151.45280913\rho^2) \cos 5\theta$$

$$S_{35} = (3.69268433\rho - 59.40323317\rho^3 + 251.40397826\rho^5 - 307.20401818\rho^7) \sin \theta$$

$$+ (28.20381860\rho^3 - 183.86176738\rho^5 + 272.43249673\rho^7) \sin 3\theta$$

$$+ (19.83875817\rho^5 - 48.16032819\rho^7) \sin 5\theta + 32.65536033\rho^7 \sin 7\theta$$

$$S_{36} = (-3.69268433\rho + 59.40323317\rho^3 - 251.40397826\rho^5 + 307.20401818\rho^7) \cos \theta$$

$$+ (28.20381860\rho^3 - 183.86176738\rho^5 + 272.43249673\rho^7) \cos 3\theta$$

$$+ (-19.83875817\rho^5 + 48.16032819\rho^7) \cos 5\theta + 32.65536033\rho^7 \cos 7\theta$$

$$S_{37} = 2.34475558 - 55.32128002\rho^2 + 296.53777290\rho^4 - 553.46621887\rho^6 + 332.94452229\rho^8$$

$$+ (-12.75329096\rho^4 + 20.75498320\rho^6) \cos 4\theta$$

$$S_{38} = (-51.83202694\rho^2 + 451.93890159\rho^4 - 1158.49126888\rho^6 + 910.24313983\rho^8) \cos 2\theta$$

$$+ 5.51662508\rho^6 \cos 6\theta$$

$$S_{39} = (-39.56789598\rho^2 + 267.47071204\rho^4 - 525.02362247\rho^6 + 310.24123146\rho^8) \sin 2\theta$$

$$- 1.59098067\rho^6 \sin 6\theta$$

$$S_{40} = 1.21593465 - 45.42224477\rho^2 + 373.41167834\rho^4 - 1046.32659847\rho^6 + 933.93661610\rho^8$$

$$+ (137.71626496\rho^4 - 638.10242034\rho^6 + 712.98912399\rho^8) \cos 4\theta$$

$$S_{41} = (9/8)\sqrt{7/1787}(1455 - 5544\rho^2 + 5005\rho^4)\rho^4 \sin 4\theta$$

$$S_{42} = (-40.45171657\rho^2 + 494.75561036\rho^4 - 1738.64589491\rho^6 + 1843.19802390\rho^8) \cos 2\theta$$

$$+ (-150.76043598\rho^6 + 318.07940431\rho^8) \cos 6\theta$$

$$S_{43} = (-9.12193686\rho^2 + 110.47679089\rho^4 - 371.21215287\rho^6 + 368.07015240\rho^8) \sin 2\theta$$

$$+ (-107.35168289\rho^6 + 200.31338972\rho^8) \sin 6\theta$$

$$S_{44} = 0.58427150 - 25.29433513\rho^2 + 242.54313549\rho^4 - 795.02011474\rho^6 + 830.47943579\rho^8$$

$$+ (90.22533813\rho^4 - 538.44320774\rho^6 + 752.97905752\rho^8) \cos 4\theta + 52.52630092\rho^8 \cos 8\theta$$

$$S_{45} = (31.08509142\rho^4 - 194.79990628\rho^6 + 278.72965314\rho^8) \sin 4\theta + 44.08655427\rho^8 \sin 8\theta$$

TABLE 6c Orthonormal Square Polynomials $S_j(x, y)$ in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2, -1/\sqrt{2} \leq x \leq 1/\sqrt{2}$, and, $-1/\sqrt{2} \leq y \leq 1/\sqrt{2}$

$$S_1 = 1$$

$$S_2 = \sqrt{6}x$$

$$S_3 = \sqrt{6}y$$

$$S_4 = \sqrt{5/2}(3\rho^2 - 1)$$

$$S_5 = 6xy$$

$$S_6 = 3\sqrt{5/2}(x^2 - y^2)$$

$$S_7 = \sqrt{21/31}(15\rho^2 - 7)y$$

$$S_8 = \sqrt{21/31}(15\rho^2 - 7)x$$

$$S_9 = \sqrt{5/31}(27x^2 - 35y^2 + 6)y$$

$$S_{10} = \sqrt{5/31}(35x^2 - 27y^2 - 6)x$$

$$S_{11} = [1/(2\sqrt{67})](315\rho^4 - 240\rho^2 + 31)$$

$$S_{12} = [15/(2\sqrt{2})](x^2 - y^2)(7\rho^2 - 3)$$

$$S_{13} = \sqrt{42}(5\rho^2 - 3)xy$$

(Continued)

TABLE 6c Orthonormal Square Polynomials $S_j(x, y)$, in Cartesian Coordinates, Where $\rho^2 = x^2 + y^2$, $-1/\sqrt{2} \leq x \leq 1/\sqrt{2}$, and, $-1/\sqrt{2} \leq y \leq 1/\sqrt{2}$ (Continued)

$$S_{14} = [3/(4\sqrt{134})][10(49x^4 - 36x^2y^2 + 49y^4) - 150\rho^2 + 11]$$

$$S_{15} = 5\sqrt{42}(x^2 - y^2)xy$$

$$S_{16} = \sqrt{55/1966}(315\rho^4 - 280x^2 - 324y^2 + 57)x$$

$$S_{17} = \sqrt{55/1966}(315\rho^4 - 324x^2 - 280y^2 + 57)y$$

$$S_{18} = (1/2)\sqrt{3/844397}[105(1023x^4 + 80x^2y^2 - 943y^4) - 61075x^2 + 39915y^2 + 4692]x$$

$$S_{19} = (1/2)\sqrt{3/844397}[105(943x^4 - 80x^2y^2 - 1023y^4) - 39915x^2 + 61075y^2 - 4692]y$$

$$S_{20} = (1/4)\sqrt{7/859}[6(693x^4 - 500x^2y^2 + 525y^4) - 1810x^2 - 450y^2 + 165]x$$

$$S_{21} = (1/4)\sqrt{7/859}[6(525x^4 - 500x^2y^2 + 693y^4) - 450x^2 - 1810y^2 + 165]y$$

$$S_{22} = (1/4)\sqrt{65/849}[1155\rho^6 - 15(91x^4 + 198x^2y^2 + 91y^4) + 453\rho^2 - 31]$$

$$S_{23} = \sqrt{33/3923}(1575\rho^4 - 1820\rho^2 + 471)xy$$

$$S_{24} = (21/4)\sqrt{65/1349}(165\rho^4 - 140\rho^2 + 27)(x^2 - y^2)$$

$$S_{25} = 7\sqrt{33/2}(9\rho^2 - 5)xy(x^2 - y^2)$$

$$S_{26} = [1/(8\sqrt{849})][42(1573x^6 - 375x^4y^2 - 375x^2y^4 + 1573y^6) - 60(707x^4 - 225x^2y^2 + 707y^4) + 6045\rho^2 - 245]$$

$$S_{27} = [1/(2\sqrt{7846})][14(2673x^4 - 2500x^2y^2 + 2673y^4) - 10290\rho^2 + 1305]xy$$

$$S_{28} = [21/(8\sqrt{1349})][3146x^6 - 2250x^4y^2 + 2250x^2y^4 - 3146y^6 - 1770(x^4 - y^4) + 245(x^2 - y^2)]$$

$$S_{29} = (-13.79189793 + 150.92209099x^2 + 117.01812058y^2 - 352.15154565x^4 - 657.27245247x^2y^2 - 291.12439892y^4 + 222.62454035x^6 + 667.87362106x^4y^2 + 667.87362106x^2y^4 + 222.62454035y^6)y$$

$$S_{30} = (-13.79189793 + 117.01812058x^2 + 150.92209099y^2 - 291.12439892x^4 - 657.27245247x^2y^2 - 352.15154565y^4 + 222.62454035x^6 + 667.87362106x^4y^2 + 667.87362106x^2y^4 + 222.62454035y^6)x$$

$$S_{31} = (6.14762642 + 89.43280522x^2 - 135.73181009y^2 - 496.10607212x^4 + 87.83479115x^2y^2 + 513.91209661y^4 + 509.87526260x^6 + 494.87949207x^4y^2 - 539.86680367x^2y^4 - 524.87103314y^6)y$$

$$S_{32} = (-6.14762642 + 135.73181009x^2 - 89.43280522y^2 - 513.91209661x^4 - 87.83479115x^2y^2 + 496.10607212y^4 + 524.87103314x^6 + 539.86680367x^4y^2 - 494.87949207x^2y^4 - 509.87526260y^6)x$$

$$S_{33} = (-6.78771487 + 38.11697536x^2 + 124.84070714y^2 - 400.01976911x^4 + 191.43062089x^2y^2 - 609.73320550y^4 + 695.06919087x^6 - 246.30347616x^4y^2 - 154.56957886x^2y^4 + 786.80308817y^6)y$$

$$S_{34} = (-6.78771487 + 124.84070714x^2 + 38.11697536y^2 - 609.73320550x^4 + 191.43062089x^2y^2 - 400.01976911y^4 + 786.80308817x^6 - 154.56957886x^4y^2 - 246.30347616x^2y^4 + 695.06919087y^6)x$$

$$S_{35} = (3.69268433 + 25.20822264x^2 - 87.60705178y^2 - 200.98753298x^4 - 63.30315999x^2y^2 + 455.10450382y^4 + 497.87935336x^6 - 461.58554163x^4y^2 + 470.02596297x^2y^4 - 660.45220344y^6)y$$

$$S_{36} = (-3.69268433 + 87.60705178x^2 - 25.20822264y^2 - 455.10450382x^4 + 63.30315999x^2y^2 + 200.98753298y^4 + 660.45220344x^6 - 470.02596297x^4y^2 + 461.58554163x^2y^4 - 497.87935336y^6)x$$

$$S_{37} = 9.37902233 - 221.28512011\rho^2 + 1186.15109160\rho^4 - 2213.86487550\rho^6 + 1331.77808917\rho^8 + 0.0190064(x^4 - 6x^2y^2 + y^4)(-671 + 1092\rho^2)$$

$$S_{38} = (-51.83202694 + 451.93890159x^2 - 1152.97464379x^4 + 910.24313983x^6)x^2 + (51.83202694 - 451.93890159y^2 - 1241.24064523x^4 + 1241.24064523x^2y^2 + 1152.97464379y^4 + 1820.48627967x^6 - 1820.48627967x^2y^2 - 910.24313983y^6)y^2$$

$$S_{39} = (-79.13579197 + 534.94142408x^2 + 534.94142408y^2 - 1059.59312899x^4 - 2068.27487642x^2y^2 - 1059.59312899y^4 + 620.48246292x^6 + 1861.44738877x^4y^2 + 1861.44738877x^2y^4 + 620.48246292y^6)xy$$

$$S_{40} = 1.21593465 + (-45.42224477 + 511.12794331x^2 - 1684.42901882x^4 + 1646.92574009x^6)x^2 + (-45.42224477 - 79.47423312x^2 + 511.12794331y^2 + 51.53230630x^4 + 51.53230630x^2y^2 - 1684.42901882y^4 + 883.78996844x^6 - 1526.27154329x^4y^2 + 883.78996844x^2y^4 + 1646.92574009y^6)y^2$$

$$S_{41} = (409.79084415x^2 - 409.79084415y^2 - 1561.42985567x^4 + 1561.42985567y^4 + 1409.62417525x^6 + 1409.62417525x^4y^2 - 1409.62417525x^2y^4 - 1409.62417525y^6)xy$$

$$S_{42} = (-40.45171657 + 494.75561036x^2 - 1889.40633090x^4 + 2161.27742821x^6)x^2 + (40.45171657 - 494.75561036y^2 + 522.76064491x^4 - 522.76064491x^2y^2 + 1889.40633090y^4 - 766.71561254x^6 + 766.71561254x^4y^2 - 2161.27742821y^6)y^2$$

$$S_{43} = (-18.24387372 + 220.95358178x^2 + 220.95358178y^2 - 1386.53440310x^4 + 662.18504631x^2y^2 - 1386.53440310y^4 + 1938.02064313x^6 - 595.96654168x^4y^2 - 595.96654168x^2y^4 + 1938.02064313y^6)xy$$

$$S_{44} = 0.58427150 + (-25.29433513 + 332.76847363x^2 - 1333.46332249x^4 + 1635.98479424x^6)x^2 + (-25.29433513 - 56.26575785x^2 + 332.76847363y^2 + 307.15569451x^4 + 307.15569451x^2y^2 - 1333.46332249y^4 - 1160.73491284x^6 + 1129.92710444x^4y^2 - 1160.73491284x^2y^4 + 1635.98479424y^6)y^2$$

$$S_{45} = (124.34036571x^2 - 124.34036571y^2 - 779.19962514x^4 + 779.19962514y^4 + 1467.61104674x^6 - 1353.92842666x^4y^2 + 1353.92842666x^2y^4 - 1467.61104674y^6)xy$$

TABLE 7 Orthonormal Polynomials for a Unit *Slit* Pupil

j	Aberration	Orthonormal Polynomials
1	Piston	1
2	Tilt	$\sqrt{3}x$
3	Defocus	$(\sqrt{5}/2)(3x^2 - 1)$
4	Coma	$(\sqrt{7}/2)(5x^3 - 3x)$
5	Spherical aberration	$(3/8)(35x^4 - 30x^2 + 3)$
6	Secondary coma	$(\sqrt{11}/8)(63x^5 - 70x^3 + 15x)$
7	Secondary spherical aberration	$(\sqrt{13}/16)(231x^6 - 315x^4 + 105x^2 - 5)$

a small amount of astigmatism, the diffraction focus for an inscribed hexagonal pupil is the same as for a circular pupil.^{4,5} For an image with a focal ratio of F , it lies along the z axis at a distance of $-8F^2$ times the amount of the balancing defocus from the Gaussian image point. However, the hexagonal polynomials H_7 and H_8 show that the relative amount of tilt $\rho \cos \theta$ that optimally balances classical or Seidel coma $\rho^3 \cos \theta$ is $-14/25 \approx -0.56$ compared to $-2/3$ for a circular pupil. The diffraction focus in this case lies along the x axis at a distance of $-2F$ times the amount of tilt from the Gaussian image point. Similarly, the hexagonal polynomial H_{11} shows that the relative amount of defocus that optimally balances classical primary or Seidel spherical aberration ρ^4 is $-257/301 \approx -0.85$ compared to a value of -1 for a circular pupil. It has the consequence that the diffraction focus lies closer to the Gaussian image point in the case of coma, and closer to the Gaussian image plane in the case of spherical aberration, compared to their corresponding locations for a circular pupil. While the balanced primary and secondary spherical aberrations H_{11} and H_{22} are radially symmetric, the balanced tertiary spherical aberration H_{37} is not. The tertiary spherical aberration ρ^8 is balanced not only by defocus ρ^2 and primary and secondary spherical aberrations ρ^4 and ρ^6 , but by a term in Z_{28} or $\rho^6 \cos 6\theta$ as well.

In the case of an elliptical pupil, the sigma of Seidel astigmatism $\rho^2 \cos \theta$ is given by $\sigma_a = 1/4$, independent of its aspect ratio b , and thus equal to that for a circular pupil. Since Seidel astigmatism x^2 varies only along the x axis for which the unit ellipse has the same length as a unit circle, the sigma is independent of b . The amount of balancing defocus ρ^2 for astigmatism is different in the case of an elliptical or a rectangular pupil from the value of $-1/2$ for a circular pupil. Moreover, for these pupils, spherical aberration ρ^4 is balanced not only by defocus ρ^2 but astigmatism $\rho^2 \cos^2 \theta$ as well. This is a consequence of the fact that the x and y dimensions of these pupils are not equal.

A square pupil is a special case of a rectangular pupil for which $a = 1/\sqrt{2}$. It is evident from the square polynomials S_5 and S_6 that they have the same form as the corresponding circle polynomials. Thus there is no additional defocus for balancing astigmatism, as may be seen by the absence of a Z_4 term in the expression for S_6 . Hence, the diffraction focus of a system does not change when its circular pupil is replaced by an inscribed square pupil. Unlike an elliptical or a rectangular pupil, the primary spherical aberration in a square pupil is balanced by defocus only, as is evident from the radially symmetric expression for S_{11} . However, the balanced secondary and tertiary spherical aberrations are not radially symmetric, since they contain angle-dependent terms varying as $\cos 4\theta$. From the polynomials S_7 , S_8 , and S_{11} , the diffraction foci in the case of coma and spherical aberration are closer to the Gaussian image point and the Gaussian image plane, respectively, compared to their corresponding locations for a circular pupil.

The sigma of Seidel aberrations with and without balancing are listed in Table 8 for elliptical and rectangular pupils. The corresponding values for a circular, hexagonal, square, and a slit pupil are listed in Table 9.²⁶ As expected, the results for an elliptical pupil reduce to those for a circular pupil as $b \rightarrow 1$, and the results for a rectangular pupil reduce to those for a square pupil as $a \rightarrow 1/\sqrt{2}$. As the area of a unit pupil decreases in going from a circular to a hexagonal to a square pupil (from π to $3\sqrt{3}/2 \approx 2.6$ to 2), the sigma of an aberration decreases and its tolerance for a certain Strehl ratio

TABLE 8 Standard Deviation or Sigma of a Primary and a Balanced Primary Aberration for Elliptical and Rectangular Pupils

Sigma	Elliptical	Rectangular
σ_d	$(1/4)[(3 - 2b^2 + 3b^4)/3]^{1/2}$	$(2/3)[(1 - 2a^2 + 2a^4)/5]^{1/2}$
σ_a	1/4	$2a^2/(3\sqrt{5})$
σ_{ba}	$b^2/[6(3 - 2b^2 + 3b^4)]^{1/2}$	$2a^2(1 - a^2)/\{3[5(1 - 2a^2 + 2a^4)]^{1/2}\}$
σ_c	$(5 + 2b^2 + b^4)^{1/2}/8$	$a[(7 + 8a^4)/105]^{1/2}$
σ_{bc}	$(9 - 6b^2 + 5b^4)^{1/2}/24$	$2a(35 - 70a^2 + 62a^4)^{1/2}/(15\sqrt{21})$
σ_s	$(225 + 60b^2 - 58b^4 + 60b^6 + 225b^8)^{1/2}/(24\sqrt{10})$	$4(63 - 162a^2 + 206a^4 - 88a^6 + 44a^8)^{1/2}/(45\sqrt{7})$
σ_{bc}	$(45 - 60b^2 + 94b^4 - 60b^6 + 45b^8)^{1/2}/(48\sqrt{5})$	$(8/315)(9 - 36a^2 + 103a^4 - 134a^6 + 67a^8)^{1/2}$

TABLE 9 Standard Deviation or Sigma of a Primary and a Balanced Primary Aberration for Circular, Hexagonal, Square, and Slit Pupils

Sigma	Circle	Hexagon	Square	Slit
σ_d	$1/(2\sqrt{3})$ = 1/3.464	$(1/12)\sqrt{43/5}$ = 1/4.092	$(1/3)\sqrt{2/5}$ = 1/4.743	$2/(3\sqrt{5})$ = 1/3.354
σ_a	1/4	$(1/24)\sqrt{127/5}$ = 1/4.762	$1/(3\sqrt{5})$ = 1/6.708	-
σ_{ba}	$1/(2\sqrt{6})$ = 1/4.899	$(1/4)\sqrt{7/15}$ = 1/5.855	$1/(3\sqrt{10})$ = 1/9.487	-
σ_c	$1/(2\sqrt{2})$ = 1/2.828	$(1/4)\sqrt{83/70}$ = 1/3.673	$\sqrt{3/70}$ = 1/4.831	$1\sqrt{7}$ = 1/2.646
σ_{bc}	$1/(6\sqrt{2})$ = 1/8.485	$(1/20)\sqrt{737/210}$ = 1/10.676	$(1/15)\sqrt{31/21}$ = 1/12.346	$2/(5\sqrt{7})$ = 1/6.614
σ_s	$2/(3\sqrt{5})$ = 1/3.354	$(1/6)\sqrt{59/35}$ = 1/4.621	$(2/45)\sqrt{101/7}$ = 1/5.923	4/15 = 1/3.750
σ_{bs}	$1/(6\sqrt{5})$ = 1/13.416	$(1/84)\sqrt{4987/215}$ = 1/17.441	$(2/315)\sqrt{67}$ = 1/19.242	8/105 = 1/13.125

increases. The slit pupil is more sensitive compared to a circular pupil, except for spherical aberration for which it is slightly less sensitive. To obtain the Seidel coefficients from the orthonormal coefficients of a noncircular wavefront, all significant coefficients that contain a Seidel term must be taken into account, just as in the case of Zernike coefficients.³⁶

11.12 ISOMETRIC, INTERFEROMETRIC, AND PSF PLOTS FOR ORTHONORMAL ABERRATIONS

The aberration-free point-spread functions (PSFs) for unit pupils considered in this chapter are shown in Fig. 5, illustrating their symmetry, for example, 6-fold symmetry for a hexagonal pupil. Their linear scale is such that the first zero of the PSF, for example, for a square pupil occurs at unity in units of λF (corresponding to 1.22 for a circular pupil). Here λ is the wavelength of the object radiation and F is the focal ratio of the image-forming light cone. These PSFs are the ultimate goal of fabrication and testing. The obscuration ratio of the annular pupil in Fig. 5b is $\epsilon = 0.5$; the aspect ratio of the elliptical pupil in Fig. 5d is $b = 0.85$; and the half width of the rectangular pupil in Fig. 5e

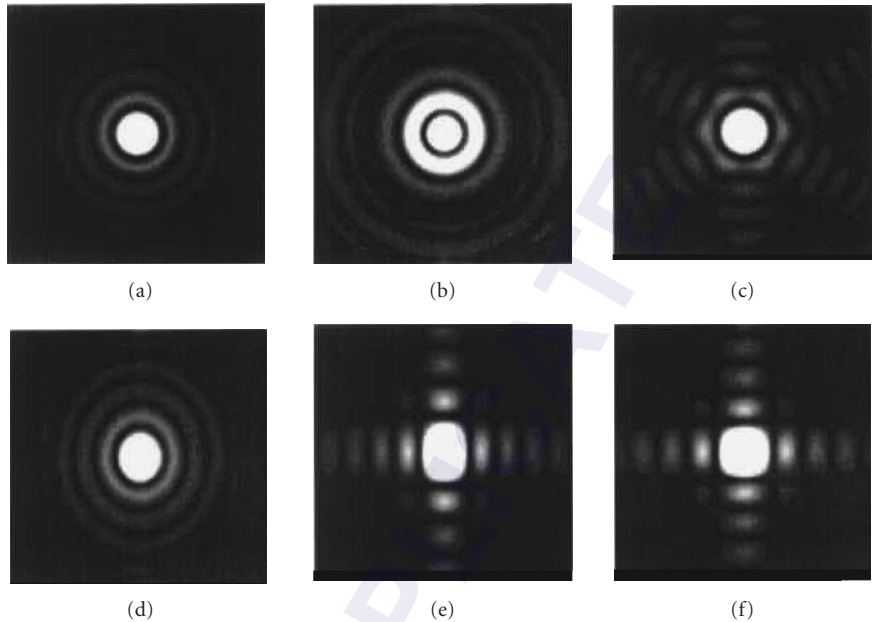


FIGURE 5 Aberration-free PSFs for different unit pupils: (a) Circular; (b) annular with obscuration ratio $\epsilon = 0.5$; (c) hexagonal; (d) elliptical with aspect ratio $b = 0.85$; (e) rectangular with half width $a = 0.8$; and (f) square.

is $a = 0.8$. The orthonormal polynomials corresponding to a Seidel aberration for a hexagonal, elliptical, rectangular, and square pupils are illustrated in three different but equivalent ways in Fig. 6.⁷ In Fig. 6d, as in Fig. 5a the aspect ratio of the elliptical pupil is $b = 0.85$. In Fig. 6e, as in Fig. 5e, the half width of the rectangular pupil is $a = 0.8$. For each polynomial, the isometric plot at the top illustrates its shape as produced, for example, in a deformable mirror. The standard deviation of each polynomial aberration in the figure is one wave. An interferogram, as in optical testing, is shown on the left. The number of fringes, which is equal to the number of times the aberration changes by one wave as we move from the center to the edge of a pupil, is different for the different polynomials. Each fringe represents a contour of constant phase or aberration. The fringe is dark when the phase is an odd multiple of π or the aberration is an odd multiple of $\lambda/2$. On the right for each polynomial are shown the PSFs, which represent the images of a point object in the presence of a polynomial aberration.

11.13 USE OF CIRCLE POLYNOMIALS FOR NONCIRCULAR PUPILS

Since the Zernike circle polynomials form a complete set, any wavefront, regardless of the shape of the pupil (which defines the perimeter of the wavefront) can be expanded in terms of them.³⁴ However, unless the pupil is circular, advantages of orthogonality and aberration balancing are lost. For example, the mean value of a Zernike circle polynomial across a noncircular pupil is not zero, the Zernike piston coefficient does not represent the mean value of the aberration, the other Zernike coefficients do not represent the standard deviation of the corresponding aberration terms, and the variance of the aberration is not equal to the sum of the squares of these other coefficients. Moreover, the value of a Zernike coefficient changes as the number of polynomials used in the

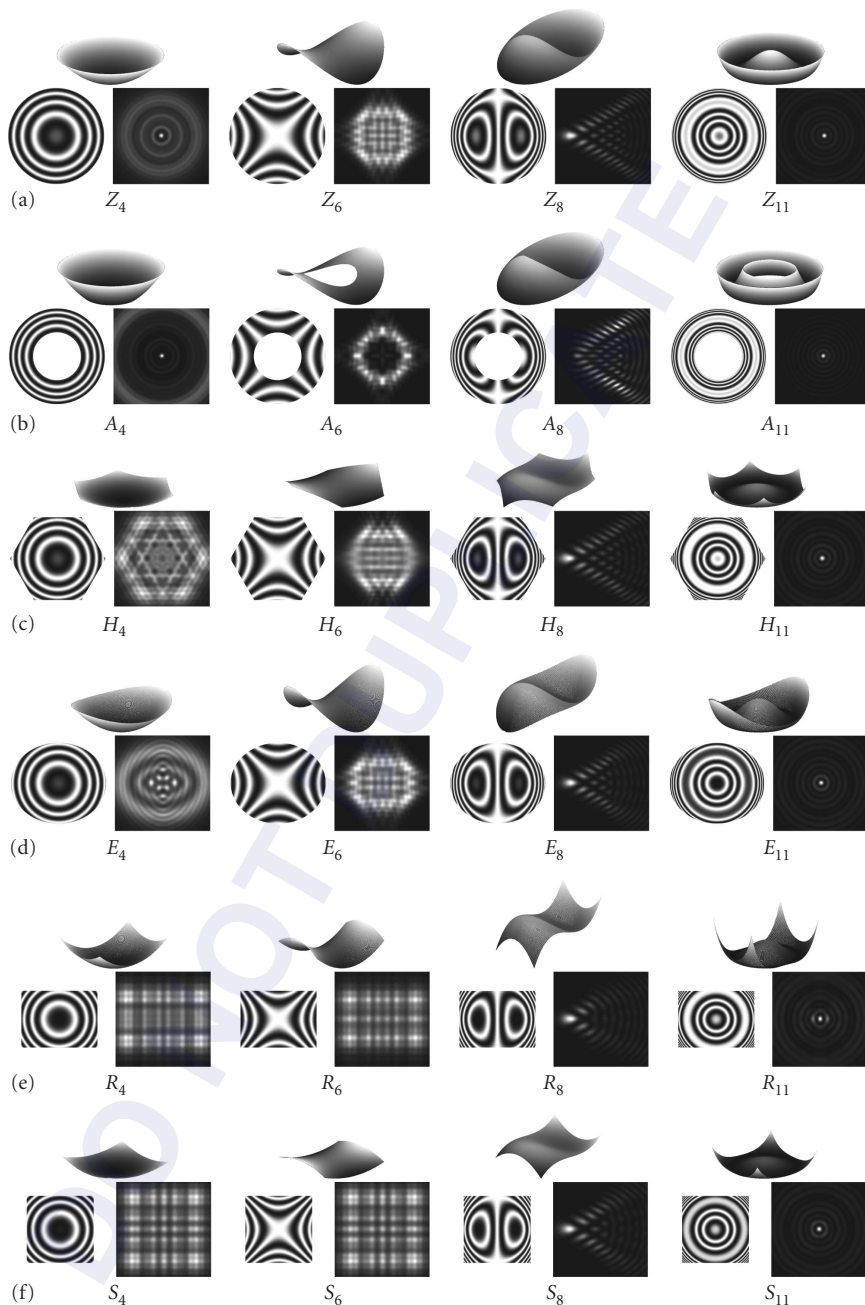


FIGURE 6 Isometric plots, interferograms, and PSFs for defocus ($j = 4$), astigmatism ($j = 6$), coma ($j = 8$), and spherical aberration ($j = 11$) in unit pupils. (a) Circular; (b) annular with $\epsilon = 0.5$; (c) hexagonal; (d) elliptical with aspect ratio $b = 0.85$; (e) rectangular with half width $a = 0.8$; and (f) square.

expansion of an aberration function changes. Hence, the circle polynomials are not appropriate for analysis of noncircular wavefronts. The polynomials given in this chapter for various pupils uniquely represent balanced classical aberrations that are also orthogonal across those pupils, just like the Zernike circle polynomials are for a circular pupil. Since each orthonormal polynomial is a linear combination of the Zernike circle polynomials, the wavefront fitting is as complete with the latter as it is with the former. However, since the circle polynomials do not represent the balanced classical aberrations for a noncircular pupil, the Zernike coefficients do not have the physical significance of their orthonormal counterparts. But the tip/tilt and defocus values in an interferometrically obtained aberration function, representing the lateral and longitudinal errors of an interferometer setting, obtained from the corresponding Zernike circle coefficients when the function is approximated with only the first four circle polynomials in a least square sense are identically the same as those obtained from the corresponding orthonormal coefficients. Accordingly, the aberration function obtained by subtracting the tip/tilt and defocus values from the measured aberration function is independent of the nature of the polynomials used in the expansion, regardless of the domain of the function or the shape of the pupil, so long as the nonorthogonal expansion is in terms of only the first four circle polynomials. The difference function is what is provided to the optician to zero out from the surface under fabrication by polishing.

11.14 DISCUSSION AND CONCLUSIONS

The Zernike circle polynomials are in widespread use for wavefront analysis in optical design and testing, because they are orthogonal over a unit circle and represent balanced aberrations of systems with circular pupils. When an aberration function of a circular wavefront is expanded in terms of them, the value of an expansion coefficient is independent of the number of polynomials used in the expansion. Accordingly, one or more terms can be added or subtracted without affecting the other coefficients. The piston coefficient represents the mean value of the aberration function and the other coefficients represent the standard deviation of the corresponding terms. The variance of the aberration is given simply by the sum of the squares of those other aberration coefficients.

We have also listed the orthonormal polynomials for analyzing the wavefronts across noncircular pupils, such as annular, hexagonal, elliptical, rectangular, and square. These polynomials are for unit pupils inscribed inside a unit circle. Such a choice keeps the maximum value of the distance of a point on the pupil from its center to be unity, thus easily identifying the peak of value of a classical aberration across it. Each orthonormal polynomial for the pupils considered consists of either the cosine or the sine terms, but not both due to the biaxial symmetry of the pupils. Whereas the circle and annular polynomials are separable in their dependence on the polar coordinates ρ and θ of a pupil point due to the radial symmetry of the pupils, only some of the polynomials for other pupils are separable. Hence polynomial numbering with two indices n and m , as for circular and annular polynomials, loses significance, and must be numbered with a single index j . The hexagonal polynomials H_{11} and H_{22} representing the balanced primary and secondary spherical aberrations are radially symmetric, but the polynomial H_{37} representing the balanced tertiary spherical aberration is not, since it contains an angle-dependent term in Z_{28} or $\cos 6\theta$ also. A hexagonal pupil has two distinct configurations where the hexagon in one is rotated by 30° with respect to that in the other. Only some of the polynomials are the same for the two configurations.²⁶ While the balancing defocus to optimally balance Seidel astigmatism for a hexagonal or a square pupil is the same as that for circular and annular pupils, it is different for the elliptical and rectangular pupils. For the elliptical and rectangular pupils, the Seidel or primary aberration ρ^4 is balanced not only by defocus but astigmatism as well. The square polynomial S_{11} representing the balanced primary spherical aberration is radially symmetric, but the square polynomial S_{22} representing the balanced secondary spherical aberration is not, since it contains a term in Z_{14} or $\cos 4\theta$ also. Similarly, the polynomial S_{37} representing the balanced tertiary spherical aberration is also not radially symmetric, since it consists of terms in Z_{14} and Z_{26} both varying as $\cos 4\theta$. We have illustrated orthonormal polynomials in three different but equivalent ways: isometrically, interferometrically, and by the corresponding aberrated PSFs.

The sigma of a Seidel aberration with and without balancing decreases as the area of a unit pupil decreases in going from a circular to a hexagonal to a square pupil. The sigma for Seidel astigmatism $\rho^2 \cos \theta$ for an elliptical pupil is independent of its aspect ratio and, therefore, is the same as for a circular pupil. This is due to the fact that the aberration is one dimensional along the dimension for which the unit ellipse has the same length as the unit circle. Since a slit pupil is one dimensional, there is no distinction between defocus and astigmatism. It is more sensitive to a Seidel aberration with or without balancing compared to a circular pupil, except for spherical aberration for which it is slightly less sensitive.

When the aberration function is known only at a discrete set of points, as in a digitized interferogram, the integral for determining the aberration coefficients reduces to a sum and the orthonormal coefficients thus obtained may be in error, since the polynomials are not orthonormal over the discrete points of the aberration data set. The magnitude of the error decreases as the number of points increases. This is not a serious problem when the wavefront errors are determined by, say, phase-shifting interferometry,³⁷ since the number of points can be very large. However, when the number of data points is small, or the pupil is irregular in shape due to vignetting, then ray tracing or testing of the system yields wavefront error data at an array of points across a region for which closed-form orthonormal polynomials are not available. In such cases, we can determine the coefficients of an expansion in terms of numerical polynomials that are orthogonal over the data set, obtained by the Gram-Schmidt orthogonalization process.^{7,38} However, if we just want to determine the values of tip/tilt and defocus terms, yielding the errors in interferometer settings, they can be obtained by least squares fitting the aberration function data with only these terms.

11.15 REFERENCES

1. F. Zernike, "Diffraction Theory of Knife-Edge Test and Its Improved Form, the Phase Contrast Method," *Mon. Not. R. Astron. Soc.* **94**: 377–384 (1934).
2. R. J. Noll, "Zernike Polynomials and Atmospheric Turbulence," *J. Opt. Soc. Am.* **66**: 207–211 (1976).
3. B. R. A. Nijboer, "The Diffraction Theory of Optical Aberrations. Part II: Diffraction Pattern in the Presence of Small Aberrations," *Physica*. **13**: 605–620 (1947).
4. M. Born and E. Wolf, *Principles of Optics*, 7th ed., Oxford, New York (1999).
5. V. N. Mahajan, *Optical Imaging and Aberrations*, Part II: Wave Diffraction Optics, SPIE Press, Bellingham, Washington (Second Printing 2004).
6. V. N. Mahajan, "Zernike Polynomials and Aberration Balancing," *SPIE Proc.* **5173**: 1–17 (2003).
7. V. N. Mahajan, "Zernike Polynomials and Wavefront Fitting," in *Optical Shop Testing*, 3rd ed., D. Malacara, ed., Wiley, New York, pp. 498–546 (2007).
8. V. N. Mahajan, "Zernike Annular Polynomials for Imaging Systems with Annular Pupils," *J. Opt. Soc. Am.* **71**: 75–85 (1981).
9. V. N. Mahajan, "Zernike Annular Polynomials for Imaging Systems with Annular Pupils," *J. Opt. Soc. Am.* **71**: 1408 (1981).
10. V. N. Mahajan, "Zernike Annular Polynomials for Imaging Systems with Annular Pupils," *J. Opt. Soc. Am. A* **1**: 685 (1984).
11. V. N. Mahajan, "Zernike Annular Polynomials and Optical Aberrations of Systems with Annular Pupils," *Appl. Opt.* **33**: 8125–8127 (1994).
12. <http://scikits.com/KFacts.html>
13. W. B. King, "The Approximation of Vignetted Pupil Shape by an Ellipse," *Appl. Opt.* **7**: 197–201 (1968).
14. G. Harbers, P. J. Kunst, and G. W. R. Leibbrandt, "Analysis of Lateral Shearing Interferograms by Use of Zernike Polynomials," *Appl. Opt.* **35**: 6162–6172 (1996).
15. H. Sumita, "Orthogonal Expansion of the Aberration Difference Function and Its Application to Image Evaluation," *Jpn. J. Appl. Phys.* **8**: 1027–1036 (1969).

16. K. N. LaFortune, R. L. Hurd, S. N. Fochs, M. D. Rotter, P. H. Pax, R. L. Combs, S. S. Olivier, J. M. Brase, and R. M. Yamamoto, "Technical Challenges for the Future of High Energy Lasers," *SPIE Proc.* **6454**: 1–11 (2007).
17. G. A. Korn and T. M. Korn, *Mathematical Handbook for Scientists and Engineers*, McGraw–Hill, New York, (1968).
18. V. N. Mahajan and G.-m. Dai, "Orthonormal Polynomials for Hexagonal Pupils," *Opt. Lett.* **31**: 2462–2465 (2006).
19. G.-m. Dai and V. N. Mahajan, "Nonrecursive Orthonormal Polynomials with Matrix Formulation," *Opt. Lett.* **32**: 74–76 (2007).
20. R. Barakat and L. Riseberg, "Diffraction Theory of the Aberrations of a Slit Aperture," *J. Opt. Soc. Am.* **55**: 878–881 (1965). There is an error in their polynomial S_2 , which should read as $x^2 - 1/3$.
21. M. Bray, "Orthogonal Polynomials: A Set for Square Areas," *SPIE Proc.* **5252**: 314–320 (2004).
22. J. L. Rayces, "Least-Squares Fitting of Orthogonal Polynomials to the Wave-Aberration Function," *Appl. Opt.* **31**: 2223–2228 (1992).
23. V. N. Mahajan, "Uniform Versus Gaussian Beams: a Comparison of the Effects of Diffraction, Obscuration, and Aberrations," *J. Opt. Soc. Am.* **A3**: 470–485 (1986).
24. V. N. Mahajan, "Zernike-Gauss Polynomials and Optical Aberrations of Systems with Gaussian Pupils," *Appl. Opt.* **34**: 8057–8059 (1995).
25. S. Szapiel, "Aberration Balancing Techniques for Radially Symmetric Amplitude Distributions; a Generalization of the Maréchal Approach," *J. Opt. Soc. Am.* **72**: 947–956 (1982).
26. V. N. Mahajan and G.-m. Dai, "Orthonormal Polynomials in Wavefront Analysis: Analytical Solution," *J. Opt. Soc. Am.* **A24**: 2994–3016 (2007).
27. A. B. Bhatia and E. Wolf, "On the Circle Polynomials of Zernike and Related Orthogonal Sets," in *Proc. Camb. Phil. Soc.* **50**: 40–48 (1954).
28. V. N. Mahajan, *Optical Imaging and Aberrations, Part I: Ray Geometrical Optics*, SPIE Press, Bellingham, Washington (Second Printing 2001).
29. W. T. Welford, *Aberrations of the Symmetrical Optical System*, Academic Press, New York (1974).
30. R. R. Shannon, *The Art and Science of Optical Design*, Cambridge University Press, New York (1997).
31. P. Mouroulis and J. Macdonald, *Geometrical Optics and Optical Design*, Oxford, New York (1997).
32. D. Malacara and Z. Malacara, *Handbook of Lens Design*, Dekkar, New York (1994).
33. G.-m. Dai and V. N. Mahajan, "Zernike Annular Polynomials and Atmospheric Turbulence," *J. Opt. Soc. Am.* **A24**: 139–155 (2007).
34. G.-m. Dai and V. N. Mahajan, "Orthonormal Polynomials in Wavefront Analysis: Error Analysis," *Appl. Opt.* **47**: 3433–3445 (2008).
35. V. N. Mahajan, "Strehl Ratio for Primary Aberrations in Terms of Their Aberration Variance," *J. Opt. Soc. Am.* **73**: 860–861 (1983).
36. V. N. Mahajan and W. H. Swantner, "Seidel Coefficients in Optical Testing," *Asian J. Phys.* **15**: 203–209 (2006).
37. K. Creath, "Phase-Measurement Interferometry Techniques," *Progress in Optics*, E. Wolf, ed., Elsevier, New York, **26**: 349–393 (1988).
38. D. Malacara, J. M. Carpio-Valdéz, and J. Javier Sánchez-Mondragón, "Wavefront Fitting with Discrete Orthogonal Polynomials in a Unit Radius Circle," *Opt. Eng.* **29**: 672–675 (1990).

This page intentionally left blank.

DO NOT DUPLICATE

Zacarías Malacara and Daniel Malacara-Hernández

*Centro de Investigaciones en Óptica, A. C.
León, Gto., México*

12.1 GLOSSARY

B	baseline length
f	focal length
f	signal frequency
f_b	back focal length
I	irradiance
N	average group refractive index
R	radius of curvature of an optical surface
r	radius of curvature of a spherometer ball
R	range
α	attenuation coefficient
λ	wavelength of light
Λ	synthetic wavelength
τ	delay time

In the optical shop, the measuring process has the purpose to obtain a comparison of physical variables using optical means. In this chapter we describe the most common procedures for the measurements of length, angle, curvature, and focal length of lenses and mirrors. The reader may obtain some more details about these procedures in the book *Optical Shop Testing* by D. Malacara,¹ or in Chap. 32, “Interferometers,” in Vol. I.

*Note: Figures 27, 30 and 31 are from *Optical Shop Testing*, 3d ed., edited by D. Malacara. (Reprinted with permission John Wiley and Sons, Inc. New York, 2007.)

12.2 INTRODUCTION AND DEFINITIONS

Lens parameter measuring in the optical shop has been a permanent problem in respect to the unit's agreement. Usually, rather than using SI units for length specifications, wavelength units is still a common reference for precision length measurements. Although the meter is obtained practically from the wavelength of a krypton source, helium-neon lasers are a common reference wavelength. The generally accepted measurement system of units is the International System or *Système International* (SI).

Time is the fundamental standard and is defined as follows: "The second is the duration of 9,192,631,770 periods of the radiation corresponding to the transition between the two hyperfine levels of the ground state of the cesium 133 atom." Formerly, the meter was a fundamental standard defined as 1,650,763.73 wavelengths in vacuum of the orange-red spectral line from the $2p_{10}$ and $5d_5$ levels of the krypton 86 atom. Shortly after the invention of the laser, it was proposed to use a laser line as a length standard.² In 1986, the speed of light was defined as 299,792,458 m/s, thus the meter is now a derived unit defined as "the distance traveled by light during $1/299,792,458$ of a second." The advantages of this definition, compared with the former meter definition, lie in the fact that it uses a relativistic constant, not subjected to physical influences, and is accessible and invariant. To avoid a previous definition of a time standard, the meter could be defined as "The length equal to $9,192,631,770/299,792,458$ wavelengths in vacuum of the radiation corresponding to the transition between the two hyperfine levels of the ground state of the cesium 133 atom."^{3,4} Modern description of the time and distance standards are described by Cundiff et al.⁵

The measurement process is prone to errors. They may be systematic or stochastic. A systematic error occurs in a poorly calibrated instrument. An instrument low in systematic error is said to be accurate. Accuracy is a measure of the amount of systematic errors. The accuracy is improved by adequate tracing to a primary standard. Stochastic errors appear due to random noise and other time-dependent fluctuations that affect the instrument. Stochastic errors may be reduced by taking several measurements and averaging them. A measurement from an instrument is said to be reproducible when the magnitude of stochastic errors is low. Reproducibility is a term used to define the repeatability for the measurements of an instrument. In measurements, a method for data acquisition and analysis has to be developed. Techniques for experimentation, planning, and data reduction can be found in the references.^{6,7}

12.3 LENGTH AND STRAIGHTNESS MEASUREMENTS

Length measurements may be performed by optical methods, since the definition of the meter is in terms of a light wavelength. Most of the length measurements are, in fact, comparisons to a secondary standard. Optical length measurements are made by comparisons to an external or internal scale a light time of flight, or by interferometric fringe counting.

Stadia and Range Finders

A *stadia* is an optical device used to determine distances. The principle of the measurement is a bar of known length W set at one end of the distance to be measured (Fig. 1). At the other end, a prism superimposes two images, one coming directly and the other after a reflection from a mirror, which are then observed through a telescope. At this point, the image of one end of the bar is brought in coincidence with that of the other end by rotating the mirror an angle θ . The mirror rotator is calibrated in such a way that for a given bar length W , a range R can be read directly in a dial, according to the equation:⁸

$$R = \frac{W}{\theta} \quad (1)$$

where θ is small and expressed in radians.

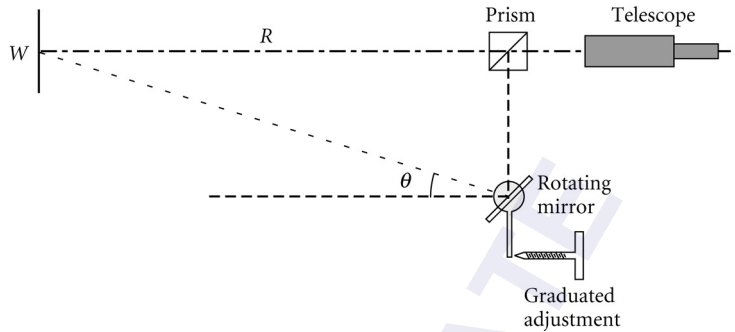


FIGURE 1 A stadia range meter. (From Patrick.⁸)

Another stadia method uses a graduated bar and a calibrated reticle in the telescope. For a known bar length W imaged on a telescope with focal length f , the bar on the focal plane will have a size i , and the range can be calculated approximately by:⁹

$$R = \left(\frac{f}{i} \right) W \quad (2)$$

Most surveying instruments have stadia markings usually in a 1:100 ratio, to measure distances from the so-called anallatic point, typically the instrument's center. A theodolite may be used for range measurements using the subtense bar method. In this method, a bar of known length is placed at the distance to be measured from the theodolite. By measuring the angle from one end to another end of the bar, the range can be easily measured.

A range finder is different from the stadia, in that the reference line is self-contained within the instrument. A general layout for a range finder is shown in Fig. 2. Two pentaprisms are separated a known baseline B ; two telescopes form an image, each through a coincidence prism. The images from the same reference point are superimposed in a split field formed by the coincidence prism. In one branch, a range compensator is adjusted to permit an effective coincidence from the reference images.

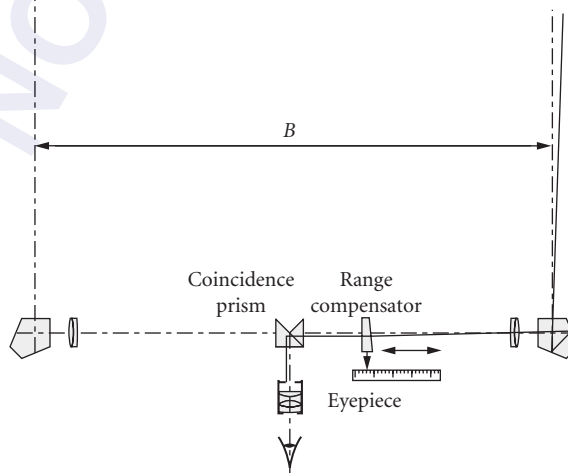


FIGURE 2 A range finder.

Assuming a baseline B and a range R , for small angles (large distances compared with the baseline), the range is

$$R = \frac{B}{\theta} \quad (3)$$

For an error $\Delta\theta$ in the angle measurement, the corresponding error ΔR in the range determination would be

$$\Delta R = -B\theta^{-2} \Delta\theta \quad (4)$$

and by substituting in Eq. (3),

$$\Delta R = -\frac{R^2}{B} \Delta\theta \quad (5)$$

From this last equation, it can be seen that the range error increases with the square of the range. Also, it is inversely proportional to the baseline. The angle error $\Delta\theta$ is a function of the eye's angular acuity, about 10 arcsec or 0.00005 rad.¹⁰

Pentaprisms permit a precise 90° deflection, independent of the alignment, and are like mirror systems with an angle of 45° between them. The focal length for the two telescopes in a range finder must be closely matched to avoid any difference in magnification. There are several versions of the range compensator. In one design, a sliding prism makes a variable deviation angle (Fig. 3a); in another system (Fig. 3b), a sliding prism displaces the image, without deviating it. A deviation can be also made with a rotating glass block (Fig. 3c). The Risley prisms (Fig. 3d) are a pair of counter rotating prisms located on the entrance pupil for one of the arms. Since the light beam is collimated, there is no astigmatism contribution from the prisms.¹⁰ A unit magnification Galilean telescope (Fig. 3e) is made with two weak lenses. A sliding lens is used to form a variable wedge to deviate the image path.⁸

Time-Based and Optical Radar

Distance measurements can also be done by the time-of-flight method. An application of the laser, obvious at the time of its advent, is the measurement of range. Distance determination by precise timing is known as optical radar. Optical radar has been used to measure the distance to the moon.

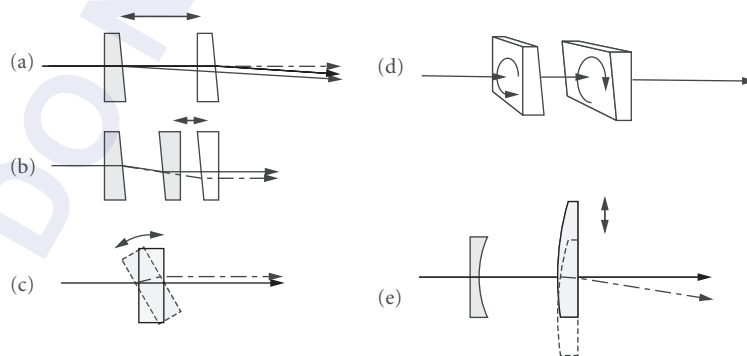


FIGURE 3 Range compensators for range finders (a) and (b) sliding prisms; (c) rotating glass block; (d) counterrotating prisms; and (e) sliding lens.

Since a small timing is involved, optical radar is applicable for distances from about 10 km. For distances from about a meter up to 50 km, modulated beams are used.

Laser radars measure the time of flight for a pulsed laser. Since accuracy depends on the temporal response of the electronic and detection system, optical radars are limited to distances larger than 1 km. Whenever possible, a cat's eye retro reflector is set at the range distance, making possible the use of a low power-pulsed laser. High-power lasers can be used over small distances without the need of a reflector. Accuracies of the order of 10^{-6} can be obtained.¹¹

The beam modulation method requires a high-frequency signal, about 10 to 30 MHz (modulating wavelength between 30 and 10 m) to modulate a laser beam carrier. The amplitude modulation is usually applied, but polarization modulation may also be used. With beam modulation distance measurements, the phase for the returning beam is compared with that of the output beam. The following description is for an amplitude modulation distance meter, although the same applies for polarizing modulation. Assuming a sinusoidally modulated beam, the returning beam can be described by

$$I_R = \alpha I_0 [1 + A \sin \omega(t - \tau)] \quad (6)$$

where α is the attenuation coefficient for the propagation media, I_0 is the output intensity for the exit beam; ω is 2π times the modulated beam frequency, and τ is the delay time. By measuring the relative phase between the outgoing and the returning beam, the quantity ω is measured. In most electronic systems, the delay time τ is measured, so, the length in multiples of the modulating wavelength is

$$L = \frac{c}{2N_g} \tau \quad (7)$$

where N_g is the average group refractive index for the modulating frequency.¹¹

Since the measured length is a multiple of the modulating wavelength, one is limited in range to one-half of the modulating wavelength. To measure with acceptable precision, and at the same time measure over large distances, several modulating frequencies are used, sometimes at multiples of ten to one another. The purpose is to obtain a synthetic wavelength Λ obtained from the mixing of two wavelengths λ_1, λ_2 as follows:^{12,13}

$$\Lambda = \frac{\lambda_1 \lambda_2}{\lambda_1 - \lambda_2} \quad (8)$$

To increment the range, several wavelengths are used to have several synthetic wavelengths. A frequency comb can increase the accuracy up to 8 nm in a range of 800 mm.¹⁴ The use of a femto-second mode locked laser produces the desired frequency comb.¹⁵

The traveling time is measured by comparing the phase of the modulating signal for the exiting and the returning beams. This phase comparison is sometimes made using a null method. In this method, a delay is introduced in the exiting signal, until it has the same phase as the returning beam, as shown in Fig. 4. Since the introduced delay is known, the light traveling time is thus determined.

Another method uses the down conversion of the frequency of both signals. These signals, with frequency f_s , are mixed with an electrical signal with frequency f_o , in order to obtain a signal with lower frequency f_L , in the range of a few kHz. The phase difference between the two low-frequency signals is the same as that between the two original signals. The lowering of the frequencies permits us to use conventional electronic methods to determine the phase difference between the two signals. A broad study on range finders has been done by Stitch.¹⁶

Interferometric Measurement of Small Distances

Interferometric methods may be used to measure small distances with a high degree of accuracy. This is usually done with a Michelson interferometer by comparing the thickness of the lens or glass

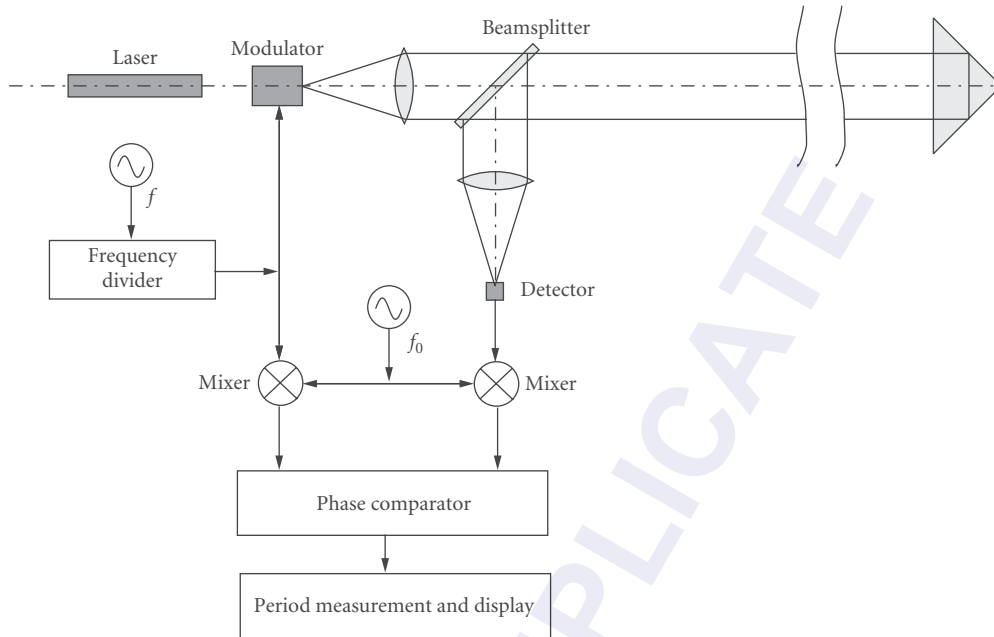


FIGURE 4 A wave modulation distance meter.

plate with that of a calibrated reference glass plate. Both plates must have approximately the same thickness and should be made with the same material.¹⁷ The two mirrors of a dispersion-compensated Michelson interferometer are replaced by the glass plate to be measured and by a reference plane-parallel plate of the same material as the lens. (The next step is to adjust the interferometer, to produce white-light Newton rings with the front surfaces of the lens and the plate.) Then, the plate is translated along the arm of the interferometer until the rear surface produces white-light rings. The displacement is the optical thickness Nt of the measured plate.

Interferometric Measurement of Medium Distances

Long and medium distances may also be measured by interferometric methods.^{18,19,20} Basically, the method counts the fringes in an interferometer while increasing or decreasing the optical path difference. The low temporal coherence or monochromaticity of most light sources limits this procedure to short distances. Lasers, however, have a much longer coherence length, due to their higher monochromaticity. Using lasers, it has been possible to make interferometric distance measurements over several meters.

In these interferometers, three things should be considered during their design. The first is that the laser light illuminating the interferometer should not be reflected back to the laser because that would cause instabilities in the laser, resulting in large irradiance fluctuations. As the optical path difference is changed by moving one of the mirrors, an irradiance detector in the pattern will detect a sinusoidally varying signal, but the direction of this change cannot be determined. Therefore, the second thing to be considered is that there should be a way to determine if the fringe count is going up or down; that is, if the distance is increasing or decreasing. There are two basic approaches to satisfy this last requirement. One is by producing two sinusoidal signals in phase quadrature (phase difference of 90° between them). The direction of motion of the moving prism may be sensed by determining the phase of which signal leads or lags the phase of the other signal. This information

is used to make the fringe counter increase or decrease. The alternative method uses a laser with two different frequencies. Finally, the third thing to consider in the interferometer design is that the number of fringes across its aperture should remain low and constant while moving the reflector in one of the two interferometer arms. This last condition is easily satisfied by using retroreflectors instead of mirrors. Then the two interfering wavefronts will always be almost flat and parallel. A typical retroreflector is a cube corner prism of reasonable quality to keep the number of fringes over the aperture low.

One method of producing two signals in quadrature is to have only a small number of fringes over the interferogram aperture. One possible method is by deliberately using an imperfect retroreflector. Then, the two desired signals are two small slits parallel to the fringes and separated by one-fourth of the distance between the fringes.¹¹ To avoid illuminating back the laser, the light from this laser should be linearly polarized. Then, a $\lambda/4$ phase plate is inserted in front of the beam, with its slow axis set at 45° to the interferometer plane to transform it into a circularly polarized beam.

Another method used to produce the two signals in quadrature phase is to take the signals from the two interference patterns that are produced in the interferometer. If the beam splitters are dielectric (no energy losses), the interference patterns will be complements of each other and, thus, the signals will be 90° apart. By introducing appropriate phase shifts in the beam splitter using metal coatings, the phase difference between the two patterns may be made 90° , as desired.²¹ This method was used by Rowley,²² as illustrated in Fig. 5. In order to separate the two patterns from the incident light beam, a large beam splitter and large retroreflectors are used. This configuration has the advantage that the laser beam is not reflected back. This is called a nonreacting interferometer configuration.

One more method, illustrated in Fig. 6, is the nonreacting interferometer designed at the Perkin-Elmer Corporation by Minkowitz and Vanir.²³ A circularly polarized light beam, produced by a linearly polarized laser and a $\lambda/4$ phase plate, illuminates the interferometer. This beam is divided by a beam splitter into two beams going to both arms of the interferometer. Upon reflection by the retro reflector, one of the beams changes its state of polarization from right to left circularly polarized. The two beams with opposite circular polarization are recombined at the beam splitter, thus producing linearly polarized light. The angle of the plane of polarization is determined by the phase difference between the two beams. The plane of polarization rotates 360° if the optical path difference

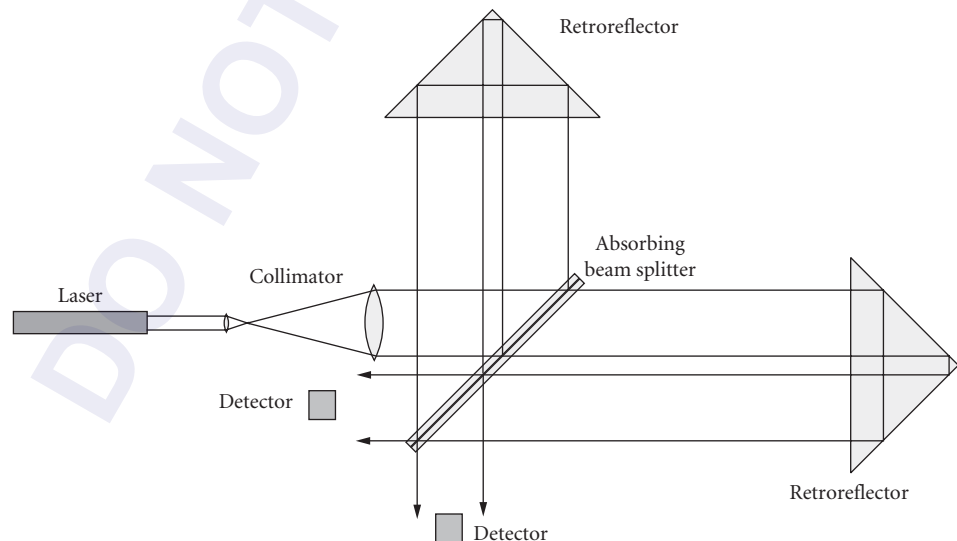


FIGURE 5 Two-interference pattern distance-measuring interferometer.

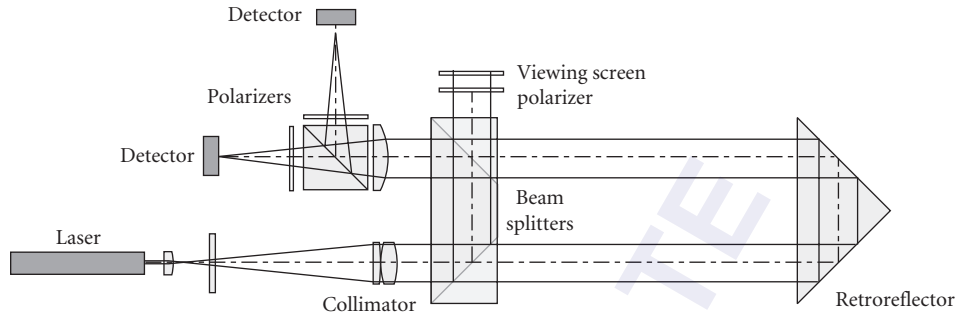


FIGURE 6 Minkowitz distance-measuring interferometer.

is changed by $\lambda/2$, the direction of rotation being given by the direction of the displacement. This linearly polarized beam is divided into two beams by a beam splitter. On each of the exiting beams, a linear polarizer is placed, one at an angle of $+45^\circ$ and the other at an angle of -45° . Then the two beams are in quadrature to each other.

In still another method, shown in Fig. 7, a beam of light, linearly polarized at 45° (or circularly polarized), is divided at a beam splitter, the p and s components. Then, both beams are converted to circular polarization with a $\lambda/4$ phase plate in front of each of them, with their axis at 45° . Upon reflection on the retro-reflectors, the handedness of the polarization is reversed. Thus, the linearly polarized beams exiting from the phase plates on the return to the beam splitter will have a plane of polarization orthogonal to that of the incoming beams. It is easily seen that no light returns to the laser. Here, the nonreacting configuration is not necessary but it may be used for additional protection. After recombination on the beam splitter, two orthogonal polarizations are present. Each plane of polarization contains information about the phase from one arm only so that no interference between the two beams has occurred. The two desired signals in quadrature are then generated by producing two identical beams with a nonpolarizing beam splitter with a polarizer on each exiting beam with their axes at $+45^\circ$ and -45° with respect to the vertical plane. The desired phase difference between the two beams is obtained by introducing a

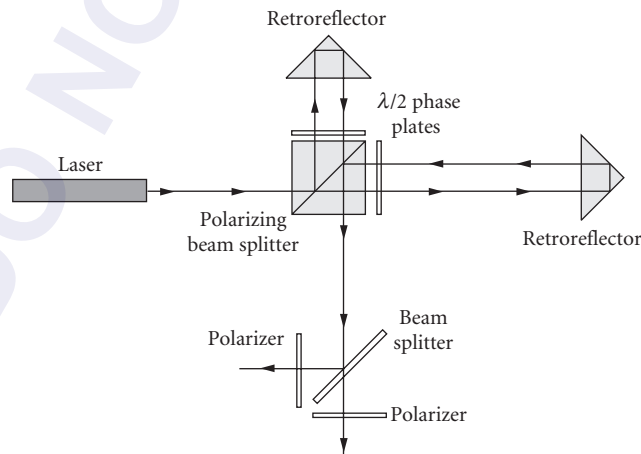


FIGURE 7 Brunning distance-measuring interferometer.

$\lambda/4$ phase plate after the beam splitter but before one of the polarizers with its slow axis vertical or horizontal. These two polarizers may be slightly rotated to make the two irradiances equal, at the same time preserving the 90° angle between their axes. If the prism is shifted a distance x , the fringe count is

$$\Delta_{\text{count}} = \pm \left[\frac{2x}{\lambda} \right] \quad (9)$$

where $[]$ denotes the integer part of the argument.

These interferometers are problematic in that any change in the irradiance may be easily interpreted as a fringe monitoring the light source. A more serious problem is the requirement that the static interference pattern be free of, or with very few, fringes. Fringes may appear because of multiple reflections or turbulence.

A completely different method uses two frequencies. It was developed by the Hewlett Packard Co.^{24,25} and is illustrated in Fig. 8. The light source is a frequency-stabilized He-Ne laser whose light beam is Zeeman split into two frequencies f_1 and f_2 by application of an axial magnetic field. The frequency difference is several megahertz and both beams have circular polarization, but with opposite sense. A $\lambda/4$ phase plate transforms the signals f_1 and f_2 into two orthogonal linearly polarized beams, one in the horizontal and the other in the vertical plane. A sample of this mixed signal is deviated by a beam splitter and detected at photo detector A, by using a polarizer at 45° . The amplitude modulation of this signal, with frequency $f_1 - f_2$, is detected and passed to a counter. Then, the two orthogonally polarized beams with frequencies f_1 and f_2 are separated at a polarizing beam splitter. Each is transformed into a circularly polarized beam by means of $\lambda/4$ phase plates. After reflection by the prisms, the handedness of these polarizations is changed. Then they go through the same phase plates where they are converted again to orthogonal linearly polarized beams. There is no light reflecting back to the laser. After recombination at the beam splitter, a polaroid at 45° will take the components of both beams in this plane. This signal is detected at the photo detector B. As with the other signal, the modulation with frequency $f_1 - f_2 + \Delta f$ is extracted from the carrier and sent to

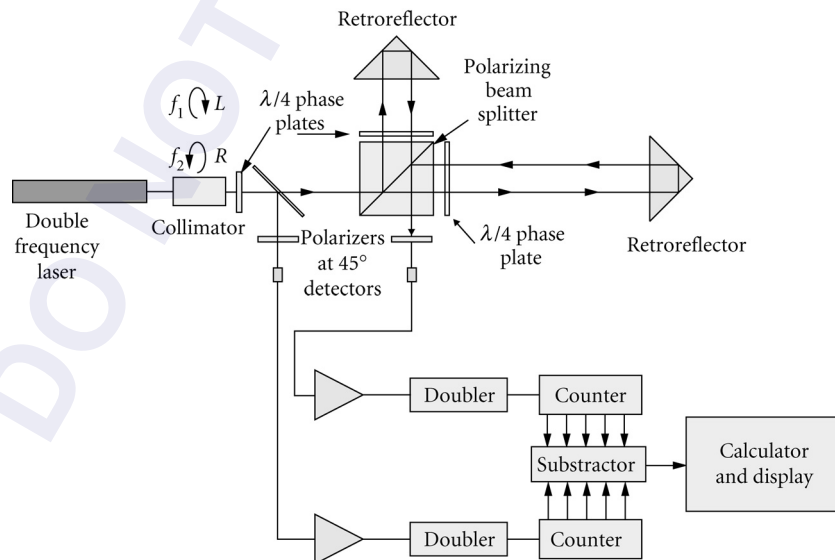


FIGURE 8 Hewlett Packard double-frequency distance-measuring interferometer.

another counter. The shift Δf in the frequency of this signal comes from a Doppler shift due to the movement of one of the retroreflectors,

$$\Delta f = \frac{2}{\lambda_2} \frac{dx}{dt} \quad (10)$$

where dx/dt is the cube corner prism velocity and λ_1 is the wavelength corresponding to the frequency f_1 . The difference between the results of the two counters is produced by the displacement of the retroreflector. If the prism moves a distance x , the number of pulses detected is given by

$$\Delta_{\text{count}} = \pm \left[\frac{2x}{\lambda_2} \right] \quad (11)$$

The advantage of this method compared with the first is that fringe counting is not subject to drift. These signals may be processed to obtain a better signal-to-noise ratio and higher resolution.²⁵

Straightness Measurements

Light propagation is assumed to be rectilinear in a homogeneous medium. This permits the use of a propagating light beam as a straightness reference. Besides the homogeneous medium, it is necessary to get a truly narrow pencil of light to improve accuracy. Laser light is an obvious application because of the high degree of spatial coherence. Beam divergence is usually less than 1 mrad for a He-Ne laser. One method uses a position-sensing detector to measure the centroid of the light spot despite its shape. In front of the laser, McLeod²⁶ used an axicon as an aligning device. When a collimated light beam is incident on an axicon, it produces a bright spot on a circular field. An axicon can give as much as 0.01 arcsec.

Another method to measure the deviation from an ideal reference line uses an autocollimator. A light beam leaving an autocollimator is reflected by a mirror. The surface slope is measured at the mirror. By knowing the distance to the mirror, one can determine the surface's profile by integration, as in a curvature measurement (see "Optical Methods for Measuring Curvature," later in this chapter). A method can be designed for measuring flatness for tables on lathe beds.²⁷

12.4 ANGLE MEASUREMENTS

Angle measurements, as well as distance measurements, require different levels of accuracy. For cutting glass, the required accuracy can be as high as several degrees, while for test plates, an error of less than a second of arc may be required. For each case, different measurement methods are developed.

Mechanical Methods

The easiest way to measure angles with medium accuracy is by means of mechanical nonoptical methods. These are

Sine Plate Essentially it is a table with one end higher than the other by a fixed amount, as shown in Fig. 9. Accuracy close to 30 arcmin may be obtained.

Goniometer This is a precision spectrometer. It has a fixed collimator and a moving telescope pointing to the center of a divided circle. Accuracies close to 20 arcsec may be obtained.

Bevel Gauge Another nonoptical method is by the use of a bevel gauge. This is made of two straight bars hinged at their edges by a pivot, as shown in Fig. 10. This device may be used to measure angle

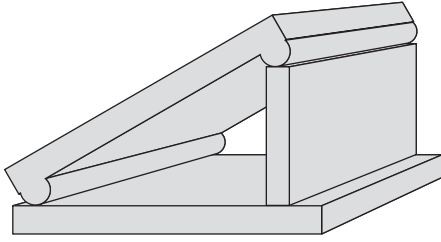


FIGURE 9 Sine plate.

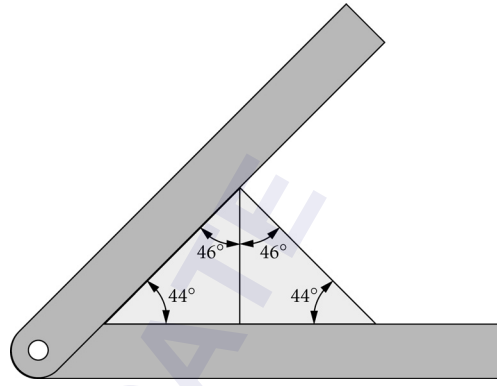


FIGURE 10 Bevel gauge.

prisms whose angle accuracies are from 45 to about 20 arcsec.²⁸ For example, if the measured prism has a 50-mm hypotenuse, a space of 5 μm at one end represents an angle of 0.0001 rad or 20 arcsec.

Numerically Controlled Machines (CNC) The advent of digitally controlled machines has brought to the optical shop machines to work prisms and angles with a high level of accuracy. Most machines can be adjusted within 0.5 arcmin of error.

Autocollimators

As shown in Fig. 11, an autocollimator is essentially a telescope focused at infinity with an illuminated reticle located at the focal plane. A complete description of autocollimators is found in Hume.²⁹ A flat reflecting surface, perpendicular to the exiting light beam, forms an image of the reticle on the same plane as the original reticle. Then, both the reticle and its image are observed through the eyepiece. When the reflecting surface is not exactly perpendicular to the exiting light beam, the reticle image is laterally displaced in the focal plane with respect to the object reticle. The magnitude of this displacement d is

$$d = 2\alpha f \quad (12)$$

where α is the tilt angle for the mirror in radians and f is the focal length of the telescope.

Autocollimator objective lenses are usually corrected doublets. Sometimes a negative lens is included to form a telephoto lens to increase the effective focal length while maintaining compactness. The collimating lens adjustment is critical for the final accuracy. Talbot interferometry can be used for a precise focus adjustment.³⁰

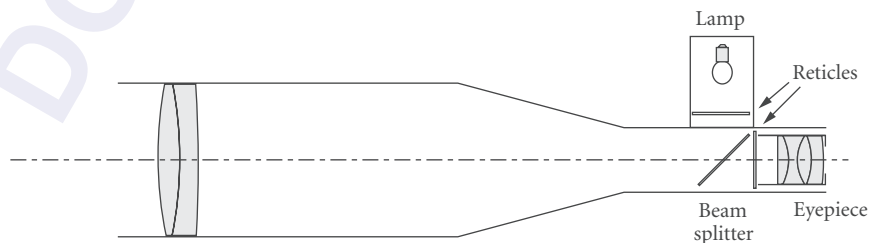


FIGURE 11 An autocollimator.

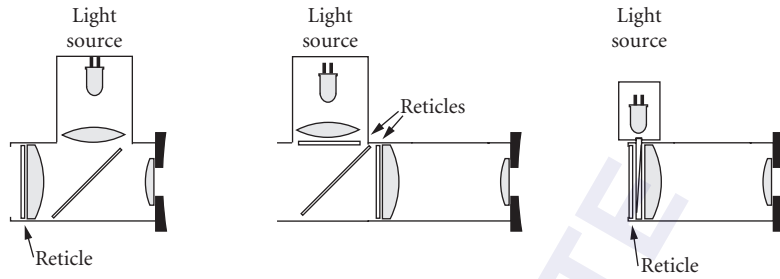


FIGURE 12 Illuminated eyepieces for autocollimators and microscopes. (a) Gauss; (b) bright line; and (c) Abbe.

The focal plane is observed through an eyepiece. Several types of illuminated reticles and eyepieces have been developed. Figure 12³¹ illustrates some illuminated eyepieces, in all of which the reticle is calibrated to measure the displacement. Gauss and Abbe illuminators show a dark reticle on a bright field. A bright field³² may be more appropriate for low reflectance surfaces. Rank³³ modified a Gauss eyepiece to produce a dark field. In other systems, a drum micrometer displaces a reticle to position it at the image plane of the first reticle. To increase sensitivity some systems, called microoptic autocollimators, use a microscope to observe the image.

Direct-reading autocollimators have a field of view of about 1° . Precision in an autocollimator is limited by the method for measuring the centroid of the image. In a diffraction-limited visual system, the diffraction image size sets the limit of precision. In a precision electronic measuring system, the accuracy of the centroid measurement is limited by the electronic detector, independent of the diffraction image itself, and can exceed the diffraction limit. In some photoelectric systems, the precision is improved by more than an order of magnitude.

Autocollimators are used for angle measurements in prisms and glass polygons. But they also have other applications; for example, to evaluate the parallelism between faces in optical flats or to manufacture divided circles.³⁴ By integrating measured slope values with an autocollimator, flatness deviations for a machine tool for an optical bed can also be evaluated.²⁷

The reflecting surface in autocollimation measurements must be kept close to the objective in order to make the alignment easier and to be sure that all of the reflected beam enters the system. The reflecting surface must be of high quality. A curved surface is equivalent to introducing another lens in the system with a change in the effective focal length.²⁷

Several accessories for autocollimators have been designed. For single-axis angle measurement, a pentaprism is used. An optical square permits angle measurements for surfaces at right angles. Perpendicularity is measured with a pentaprism and a mirror, as shown in Fig. 13. A handy horizontal reference can be produced with an oil pool, but great care must be taken with the surface stability.

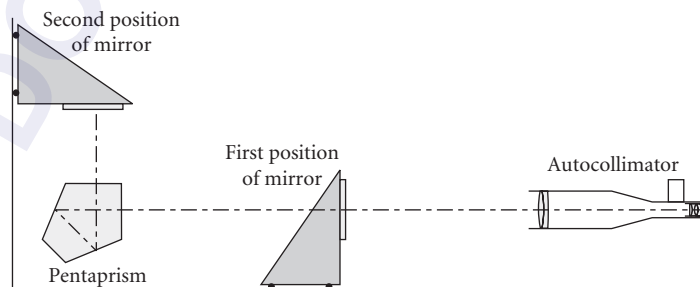


FIGURE 13 Perpendicularity measurement with an autocollimator.

Theodolites

Theodolites are surveying instruments to measure vertical and horizontal angles. A telescope with reticle is the reference to direct the instrument axis. In some theodolites, the telescope has an inverting optical system to produce an erect image. The reticle is composed of a crosswire and has a couple of parallel horizontal lines, called stadia lines, used for range measurements (see “Stadia and Range Finders” earlier in this chapter). The telescope has two focus adjustments: one to sharply focus the reticle according to the personal setting for each observer, and the other to focus the objective on the reticle. This later focus adjustment is performed by moving the whole eyepiece.

The theodolite telescope has a tree-screw base altitude-azimuth mounting on a base made horizontal with a spirit level. Divided circles attached to both telescope axes measure vertical and horizontal angles. In older instruments, an accuracy of 20 arcmin was standard. Modern instruments are accurate to within 20 arcsec, the most expensive of which can reach an accuracy of 1 arcsec. To remove errors derived from eccentric scales as well as orthogonality errors, both axes are rotated 180° and the measurement repeated. Older instruments had a provision for reading at opposite points of the scale. Scales for theodolites can be graduated in sexagesimal degrees or may use a centesimal system that divides a circle into 400 grades.

Some of the accessories available for theodolites include

1. An illuminated reticle that can be used as an autocollimator when directed to a remote retroreflector. The observer adjusts the angles until both the reflected and the instrument's reticle are superposed. This increases the pointing accuracy.
2. A solar filter, which can be attached to the eyepiece or objective side of the telescope. This is used mainly for geographic determination in surveying.
3. An electronic range meter, which is superposed to the theodolite to measure the distance. Additionally, some instruments have electronic position encoders that allow a computer to be used as an immediate data-gathering and reducing device.
4. A transverse micrometer for measuring angular separation in the image plane.

Accuracy in a theodolite depends on several factors in its construction. Several of these errors can be removed by a careful measuring routine. Some of the systematic or accuracy limiting errors are

1. Perpendicularity—deviation between vertical and horizontal scales. This error can be nulled by plunging and rotating the telescope, then averaging.
2. Concentricity deviation of scales. When scales are not concentric, they are read at opposite ends to reduce this error. Further accuracy can be obtained by rotating the instrument 90° , 180° , and 270° and averaging measurements.

Level

Levels are surveying instruments for measuring the deviation from a horizontal reference line. A level is a telescope with an attached spirit level. The angle between the telescope axis and the one defined by the spirit level must be minimized. It can be adjusted by a series of measurements to a pair of surveying staves (Fig. 14). Once the bubble is centered in the tube, two measurements are taken on

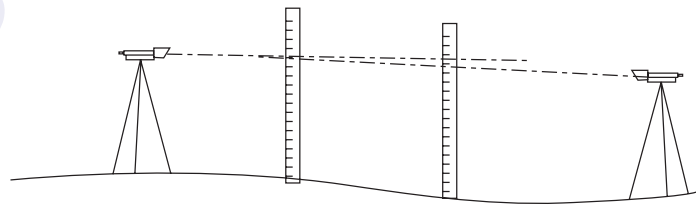


FIGURE 14 Level adjustment. (After Kingslake.³⁵)

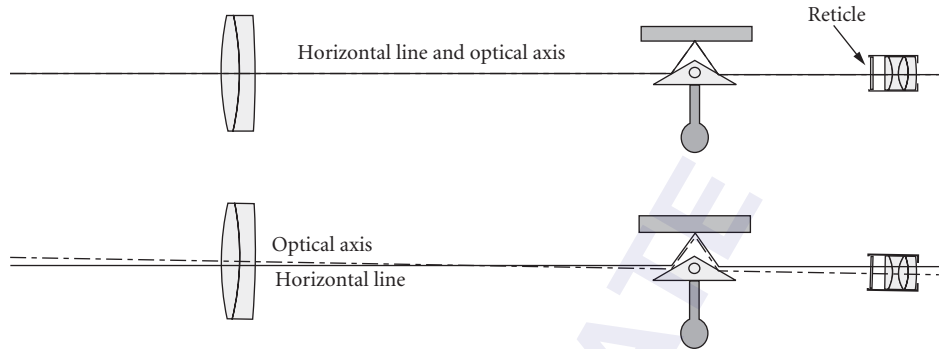


FIGURE 15 The autoset level.

each side of the staves.³⁵ The level differences between the two staves must be equal if the telescope axis is parallel with the level horizontal axis.

The autoset level (Fig. 15) uses a suspended prism and a fixed mirror on the telescope tube. The moving prism maintains the line aimed at the horizon and passing through the center of the reticle, despite the tube orientation, as long as it is within about 15 arcmin. Typical precision for this automatic level can go up to 1 arcsec.²⁷

Interferometric Measurements

Interferometric methods find their main applications in measuring very small wedge angles in glass slabs^{36,37} and in parallelism evaluation³⁸ by means of the Fizeau or Haidinger interferometers.³⁹

Interferometric measurements of large angles may also be performed. In one method, a collimated laser beam is reflected from the surfaces by a rotating glass slab. The resulting fringes can be considered as coming from a Murty lateral shear interferometer.⁴⁰ This device can be used as a secondary standard to produce angles from 0° to 360° with accuracy within a second of arc. Further analysis of this method has been done by Tentori and Celaya.⁴¹ In another system, a Michelson interferometer is used with an electronic counter to measure over a range of $\pm 5^\circ$ with a resolution of 10° .^{42,43} An interferometric optical sine bar for angles in the milliseconds of arc was built by Chapman.⁴⁴

Angle Measurements in Prisms

A problem frequently encountered in the manufacture of prisms is the precise measurement of angle. In most cases, prism angles are 90° , 45° , and 30° . These angles are easily measured by comparison with a standard but it is not always necessary.

An important aspect of measuring angles in a prism is to determine if the prism is free of pyramidal error. Consider a prism with angles A, B, and C (Fig. 16a). Let OA be perpendicular to plane ABC. If line AP is perpendicular to segment BC, then the angle AOP is a measurement of the pyramidal error. In a prism with pyramidal error, the angles between the faces, as measured in planes perpendicular to the edges between these faces, add up to over 180° . To simply detect pyramidal error in a prism, Johnson⁴⁵ and Martin⁴⁶ suggest that both the refracted and the reflected images from a straight line be examined (Fig. 16b). When pyramidal error is present, the line appears to be broken. A remote target could be graduated to measure directly in minutes of arc. A sensitivity of up to 3 arcmin may be obtained.

During the milling process in the production of a prism, a glass blank is mounted in a jig collinear with a master prism (Fig. 17). An autocollimator aimed at the master prism accurately sets the

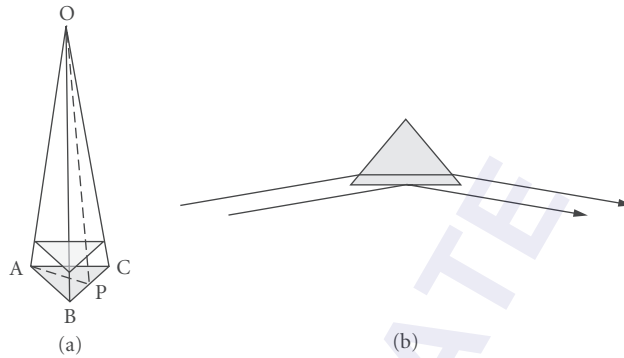


FIGURE 16 Pyramidal error in a prism (a) nature of the error and (b) test of the error.

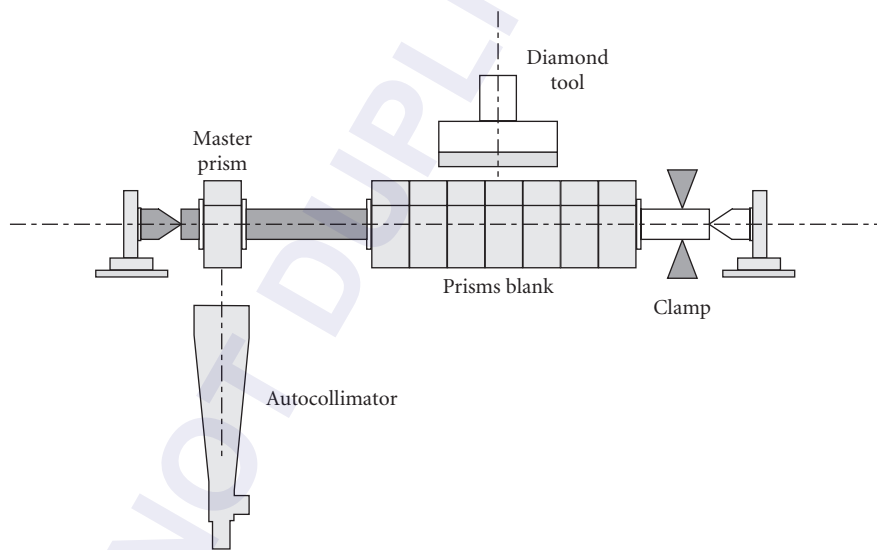


FIGURE 17 Milling prisms for replication.

position for each prism face.^{47,48} With a carefully set diamond lap, pyramidal error is minimized. In a short run, angles can be checked with a bevel gauge. Visual tests for a prism in a bevel gauge can measure an error smaller than a minute of arc.³¹

A 90° angle in a prism can be measured by internal reflection, as shown in Fig. 18a. At the autocollimator image plane, two images are seen with an angular separation of $2N\alpha$, where α is the magnitude of the prism angle error, and its sign is unknown. Since the hypotenuse face has to be polished and the glass must be homogeneous, the measurement of the external angle with respect to a reference flat is preferred (Fig. 18b). In this case, the sign of the angle error is determined by a change in the angle by tilting the prism. If the external angle is decreased and the images separate further, then the external angle is less than 90° . Conversely, if the images separate by tilting in such way that the external angle increases, then the external angle is larger than 90° .

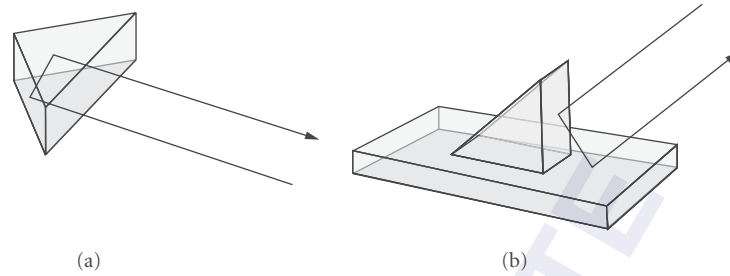


FIGURE 18 Right angle measurement in prisms: (a) internal measurement and (b) external measurement.

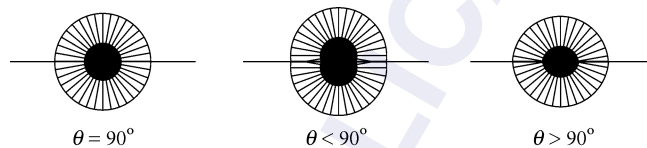


FIGURE 19 Retroreflected images of the observer's pupil in a 90° prism.

To determine the sign of the error, several other methods have been proposed. DeVany⁴⁹ suggested that when looking at the double image from the autocollimator, the image should be defocused inward. If the images tend to separate, then the angle in the prism is greater than 90° . Conversely, an outward defocusing will move the images closer to each other for an angle greater than 90° . Another way to eliminate the sign of the error in the angle is by introducing, between the autocollimator and the prism, a glass plate with a small wedge whose orientation is known. The wedge should cover only one-half of the prism aperture. Ratajczyk and Bodner⁵⁰ suggested a different method using polarized light.

Right-angle prisms can be measured using an autocollimator with acceptable precision.⁵¹ With some practice, perfect cubes with angles more accurate than 2 arcsec can be obtained.⁴⁹ An extremely simple test for the 90° angle in prisms⁴⁵ is performed by looking to the retroreflected image of the observer's pupil without any instrument. The shape of the image of the pupil determines the error, as shown in Fig. 19. The sensitivity of this test is not very great and may be used only as a coarse qualitative test. As shown by Malacara and Flores,⁵² a small improvement in the sensitivity of this test may be obtained if a screen with a small hole is placed in front of the eye, as in Fig. 20a. A cross centered on the small hole is painted on the front face of the screen. The observed images are as shown in the same Fig. 20b. As opposed to the collimator test, there is no uncertainty in the sign of the error in the tests just described, since the observed plane is located where the two prism surfaces intersect. An improvement described by Malacara and Flores,⁵² combining these simple tests with an autocollimator, is obtained with the instrument in Fig. 21. In this system, the line defining the intersection between the two surfaces is out of focus and barely visible while the reticle is in perfect focus at the eyepiece.

Corner cube prisms are a real challenge to manufacture, since besides the large precision required in the angles, all surfaces should be exempt of any curvature. The dihedral angle in pentaprisms is tested usually with an interferometer. An error in the prism alignment results in an error in angle determination.⁵³ A simple geometric method for angle measurement in corner cube reflector has been described by Rao.⁵⁴ Also, the calculations for the electric field in a corner cube are performed by Schöll.⁵⁵ These calculations include the effect by nonhomogeneities, angle of incidence, and errors in the surface finish.

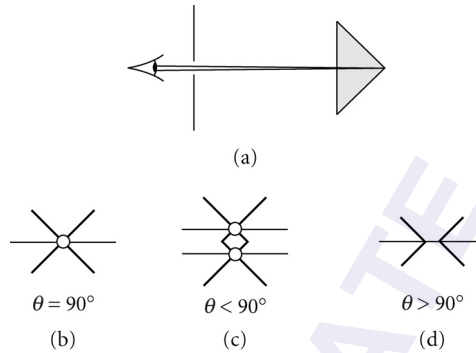


FIGURE 20 Testing a right-angle prism: (a) screen in front of the eye and (b) to (d) its observed images.

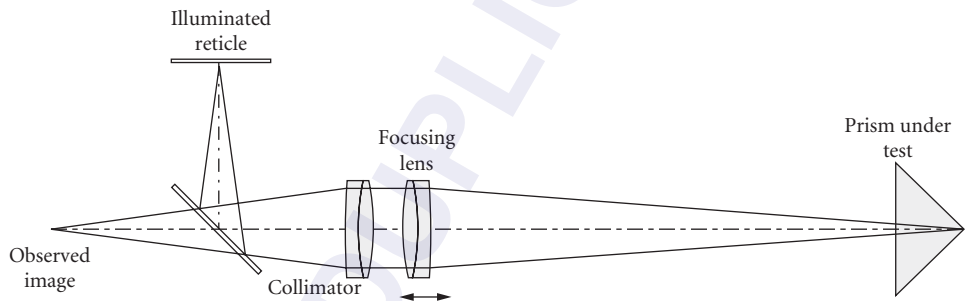


FIGURE 21 Modified autocollimator for testing the right angle in prisms without sign uncertainty in measured error.

12.5 CURVATURE AND FOCAL LENGTH MEASUREMENTS

The curvature of a spherical optical surface or the local curvature of an aspherical surface may be measured by means of mechanical or optical methods. Some methods measure the sagitta, some the surface slope, and some others directly locate the position of the center of curvature.

Mechanical Methods for Measuring Curvature

Templates The simplest and most common way to measure the radius of curvature is by comparing it with metal templates with known radii of curvature until a good fit is obtained. The template is held in contact with the optical surface with a bright light source behind the template and the optical surface. If the surface is polished, gaps between the template and the surface may be detected to an accuracy of one wavelength. If the opening is very narrow, the light passing through the gap becomes blue due to diffraction.

Test Plates This method uses a glass test plate with a curvature opposite to that of the glass surface to be measured. The accuracy is much higher than in the template method, but the surface must be polished.

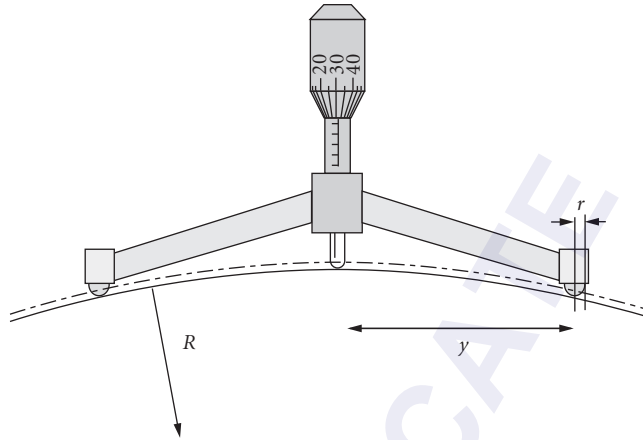


FIGURE 22 Three-leg spherometer.

Spherometers This is probably the most popular mechanical device used to measure radii of curvature. It consists of three equally spaced feet with a central moving plunger. The value of the radius of curvature is calculated after measuring the sagitta, as shown in Fig. 22. The spherometer must first be calibrated by placing it on top of a flat surface. Then it is placed on the surface to be measured. The difference in the position of the central plunger for these two measurements is the sagitta of the spherical surface being measured. Frequently, a steel ball is placed at the end of the legs as well as at the end of the plunger to avoid the possibility of scratching the surface with sharp points. In this case, if the measured sagitta is z , the radius of curvature R of the surface is given by

$$R = \frac{z}{2} + \frac{y^2}{2z} \pm r \quad (13)$$

where r is the radius of curvature of the balls. The plus sign is used for concave surfaces and the minus sign for convex surfaces. The precision of this instrument in the measurement of the radius of curvature for a given uncertainty in the measured sagitta may be obtained by differentiating Eq. (13)

$$\frac{dR}{dz} = \frac{1}{2} - \frac{y^2}{2z^2} \quad (14)$$

obtaining

$$\Delta R = \frac{\Delta z}{2} \left(1 - \frac{y^2}{z^2} \right) \quad (15)$$

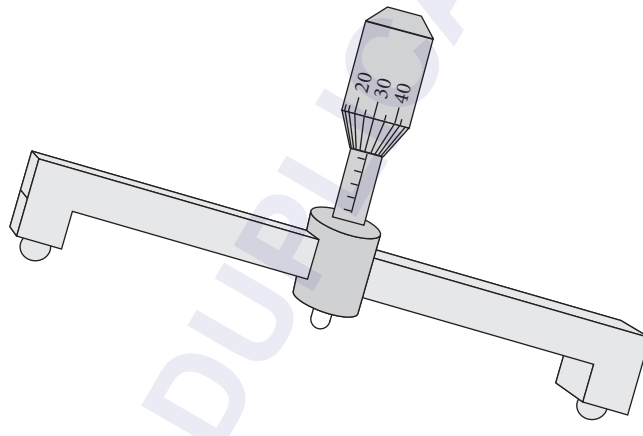
This accuracy assumes that the spherometer is perfectly built and that its dimensional parameters y and r are well known. The uncertainty comes only from human or instrumental errors in the measurement of the sagitta. Noble³¹ has evaluated the repeatability for a spherometer with $y = 50$ mm and an uncertainty in the sagitta reading equal to $5 \mu\text{m}$, and has reported the results in Table 1 where it can be seen that the precision is better than 2 percent. An extensive analysis of the precision and accuracy of several types of spherometers is given by Jurek.⁵⁶

A ring may be used instead of the three legs in a mechanical spherometer. A concave surface contacts the external edge of the cup, and a convex surface is contacted by the internal edge of the ring. Thus, Eq. (5) may be used if a different value of y is used for concave and convex surfaces, and r is taken as zero. Frequently in spherometers of this type, the cups are interchangeable in order to have different diameters for different surface diameters and radii of curvature. The main advantage of the

TABLE 1 Spherometer precision*

Radius of Sphere R (mm)	Sagitta Z (mm)	Fractional Precision ΔR (mm)	Precision $\Delta R/R$
10,000	0.125	-400	-0.040
5,000	0.250	-100	-0.020
2,000	0.625	-16	-0.008
1,000	1.251	-4	-0.004
500	2.506	-1	-0.002
200	6.351	-0.15	-0.0008

* $y = 50$ mm; $\Delta z = 5$ μm .
Source: From Noble.³¹

**FIGURE 23** Bar spherometer.

use of a ring instead of three legs is that an astigmatic deformation of the surface is easily detected, although it cannot be measured.

A spherometer that permits the evaluation of astigmatism is the bar spherometer, shown in Fig. 23. It can measure the curvature along any diameter. A commercial version of a small bar spherometer for the specific application in optometric work is the Geneva gauge, where the scale is directly calibrated in diopters assuming that the refractive index of the glass is 1.53.

Automatic spherometers use a differential transformer as a transducer to measure the plunger displacement. This transformer is coupled to an electronic circuit and produces a voltage that is linear with respect to the plunger displacement. This voltage is fed to a microprocessor which calculates the radius of curvature or power in any desired units and displays it.

Optical Methods for Measuring Curvature

Foucault Test Probably the oldest and easiest method to measure the radius of curvature of a concave surface is the knife-edge test. In this method, the center of curvature is first located by means of the knife edge. Then, the distance from the center of curvature to the optical surface is measured.

Autocollimator The radius of curvature may also be determined through measurements of the slopes of the optical surface with an autocollimator as described by Horne.⁵⁷ A pentaprism producing a 90° deflection of a light beam independent of small errors in its orientation is used in this technique,

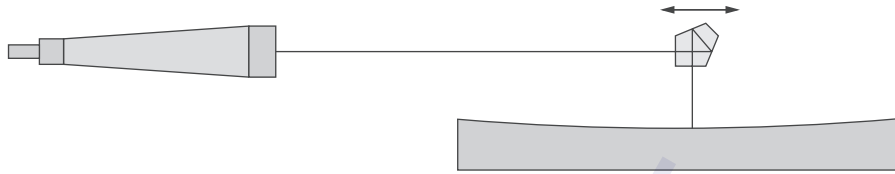


FIGURE 24 Autocollimator and pentaprism used to determine radius of curvature by measuring surface slopes.

as illustrated in Fig. 24, where the pentaprism travels over the optical surface to be measured along one diameter. First the light on the reticle of the autocollimator is centered on the vertex of the surface being examined. Then the pentaprism is moved toward the edge of the surface in order to measure any slope variations. From these slope measurements, the radius of curvature may be calculated. This method is useful only for large radii of curvature for either concave or convex surfaces.

Confocal Cavity Technique Gerchman and Hunter^{58,59} have described the so-called optical cavity technique that permits the interferometric measurement of very long radii of curvature with an accuracy of 0.1 percent. The cavity of a Fizeau interferometer is formed, as illustrated in Fig. 25. This is a confocal cavity of n th order, where n is the number of times the path is folded. The radius of curvature is equal to approximately $2n$ times the cavity length Z .

Traveling Microscope This instrument is used to measure the radius of curvature of small concave optical surfaces with short radius of curvature. As illustrated in Fig. 26, a point light source is produced at the front focus of a microscope objective. This light source illuminates the concave optical surface to be measured near its center of curvature. Then this concave surface forms an image which is also close to its center of curvature. This image is observed with the same microscope used to illuminate the surface. During this procedure, the microscope is focused both at the center of curvature and at the surface to be measured. A sharp image of the light source is observed at both places. The radius of curvature is the distance between these two positions for the microscope.

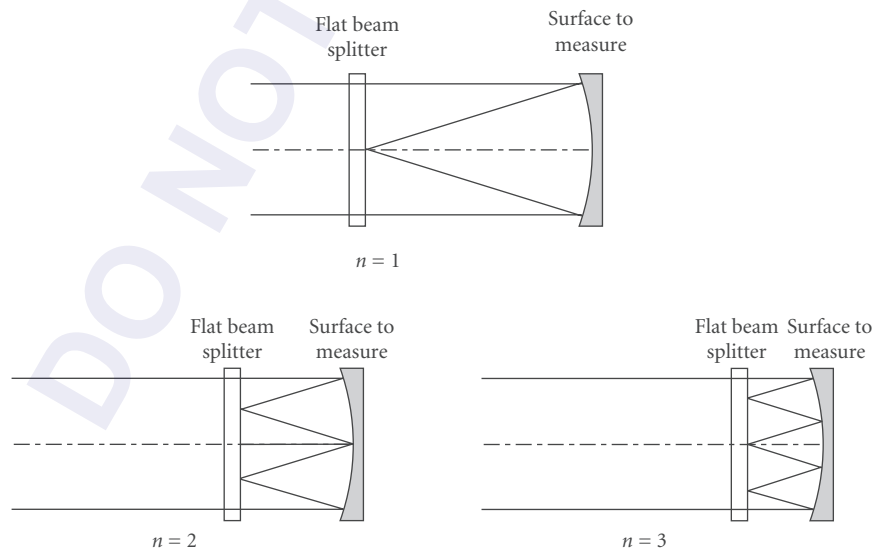


FIGURE 25 Confocal cavity arrangements used to measure radius of curvature.

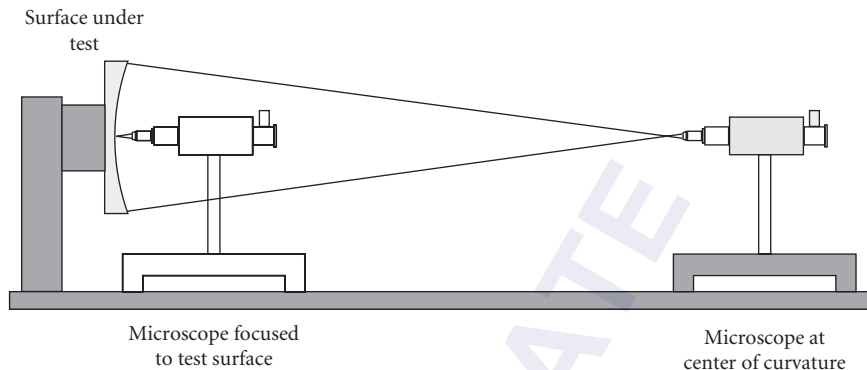


FIGURE 26 Traveling microscope to measure radii of curvature.

This distance traveled by the microscope may be measured on a vernier scale, obtaining a precision of about 0.1 mm. If a bar micrometer is used, the precision may be increased by an order of magnitude. In this case, two small convex buttons are required: one fixed to the microscope carriage and the other to the stationary part of the bench. They must face each other when the microscope carriage is close to the optical bench fixed component.

Carnell and Welford³² describe a method that requires only one measurement. The microscope is focused only at the center of curvature. Then the radius of curvature is measured by inserting a bar micrometer with one end touching the vertex of the optical surface. The other end is adjusted until it is observed to be in focus on the microscope. Accuracies of a few microns are obtained with this method.

In order to focus the microscope properly, the image of an illuminated reticle must fall, after reflection, back on itself, as in the Gauss eyepiece shown in Fig. 12. The reticle and its image appear as dark lines in a bright field. The focusing accuracy may be increased with a dark field. Carnell and Welford obtained a dark field with two reticles, as in Fig. 12, one illuminated with bright lines and the other with dark lines.

A convex surface may also be measured with this method if a well-corrected lens with a conjugate longer than the radius of curvature of the surface under test is used. Another alternative for measuring convex surfaces is by inserting an optical device with prisms in front of the microscope, as described by Jurek.⁵⁶

Some practical aspects of the traveling microscope are examined by Rank,³³ who obtained a dark field at focus with an Abbe eyepiece which introduces the illumination with a small prism. This method has been implemented using a laser light source by O'Shea and Tilstra.⁶⁰

Additional optical methods to measure the radius of curvature of a spherical surface have been described. Evans^{61,62,63} determines the radius by measuring the lateral displacements on a screen of a laser beam reflected on the optical surface when this optical surface is laterally displaced. Cornejo-Rodriguez and Cordero-Dávila,⁶⁴ Klingsporn,⁶⁵ and Diaz-Urbe et al.⁶⁶ rotate the surface about its center of curvature on a nodal bench.

Focal Length Measurements

There are two focal lengths in an optical system: the back focal length and the effective focal length. The back focal length is the distance from the last surface of the system to the focus. The effective focal length is the distance from the principal plane to the focus. The back focal length is easily measured, following the same procedure used for measuring the radius of curvature, using a microscope and the lens bench. On the other hand, the effective focal length requires the previous location of the principal plane.

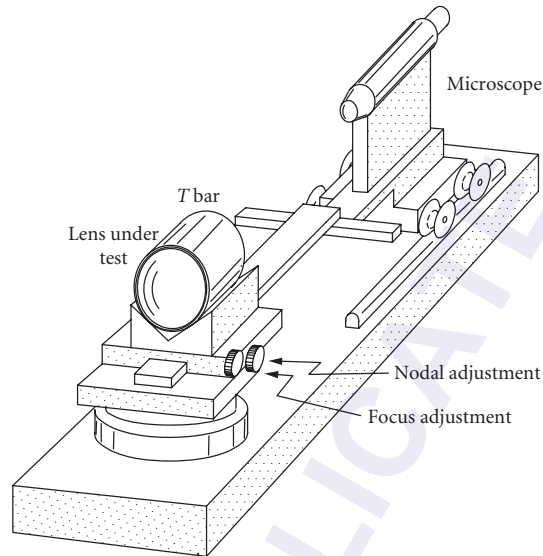


FIGURE 27 Nodal slide bench. (From Malacara.¹)

Nodal Slide Bench In an optical system in air, the principal points (intersection of the principal plane and the optical axis) coincide with the nodal points. Thus, to locate this point we may use the well-known property that small rotations of the lens about an axis perpendicular to the optical axis and passing through the nodal point do not produce any lateral shift of the image. The instrument used to perform this procedure, shown in Fig. 27, is called an optical nodal slide bench.⁶⁷ This bench has a provision for slowly moving the lens under test longitudinally in order to find the nodal point.

The bench is illuminated with a collimated light source and the image produced by the lens under test is examined with a microscope. The lens is then displaced slightly about a vertical axis as it is being displaced longitudinally. This procedure is stopped until a point is found in which the image does not move laterally while rotating the lens. This axis of rotation is the nodal point. Then, the distance from the nodal point to the image is the effective focal length.

Focimeters A focimeter is an instrument designed to measure the focal length of lenses in a simple manner. The optical scheme for the classical version of this instrument is shown in Fig. 28. A light source illuminates a reticle and a convergent lens, with focal length f , displaced at a distance x from the reticle. The lens to be measured is placed at a distance d from the convergent lens. The magnitude

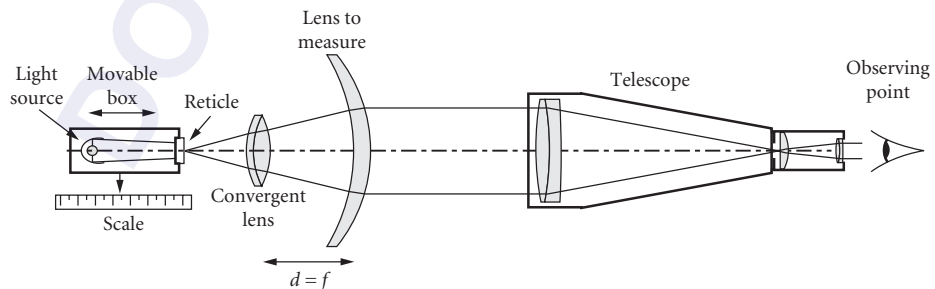


FIGURE 28 Focimeter schematics.

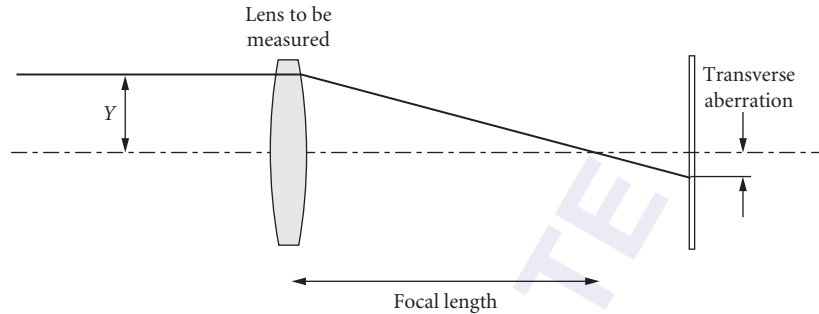


FIGURE 29 Focal length determination by transverse aberration measurements.

of x is variable and is adjusted until the light beam going out from the lens under test becomes collimated. This collimation is verified by means of a small telescope in front of this lens focused at infinity. The values of d and the focal length f are set to be equal. Then, the back focal length f_b of the lens under test is given by

$$\frac{1}{f_b} = \frac{1}{d} - \frac{x}{d^2} = P_v \quad (16)$$

where its inverse P_v is the vertex power. As can be seen, the power of the lens being measured is linear with respect to the distance x . There are many variations of this instrument. Some modern focimeters measure the lateral deviation of a light ray from the optical axis (transverse aberration), as in Fig. 29, when a defocus is introduced.^{61,62,63} This method is mainly used in some modern automatic focimeters for optometric applications. To measure the transverse aberration, a position-sensing detector is frequently used. The power error of a focimeter can be obtained by derivation of Eq. (16)

$$\delta P_v = -\frac{\delta x}{f_c^2} \quad (17)$$

Thus, the power error is a linear function of the target position error and decreases with the square of the collimating lens focal distance.⁷⁰

Other Focal Length Measurements

Moiré Deflectometry Moiré deflectometry method for focal length determination sends a collimated beam over a pair of Ronchi rulings (Fig. 30) depending on the convergence or divergence of the beam, the resulting moiré pattern rotates according to the convergence (Fig. 31). The rotation has an angle α that is related to the focal distance by^{71,72}

$$f \approx \frac{d}{\theta \tan \alpha} \quad (18)$$

d being the ruling's pitch, θ is the angle between the ruling's lines, and α is the rotation angle of the moiré pattern.

Talbot Autoimages Talbot autoimages method for focal length determination is performed by sending a coherent beam of light into a Ronchi ruling. An image of the grating will be produced periodically and evenly spaced along the light beam. Every Talbot autoimage is an object for the lens

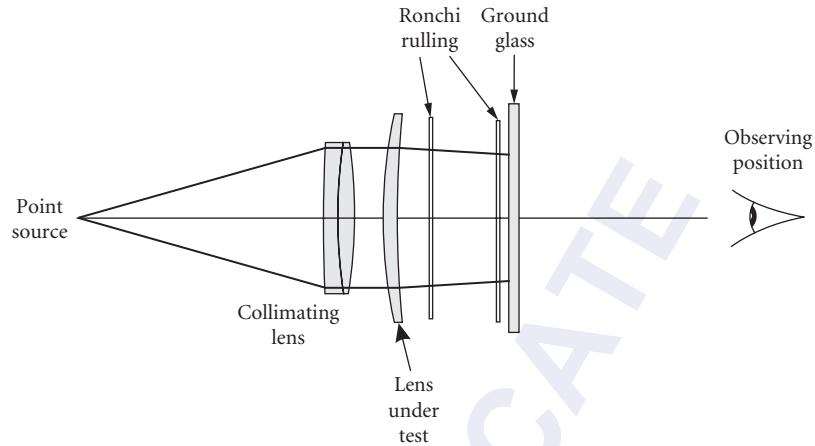


FIGURE 30 Moiré deflectometry lens power measurement. (From Malacara.¹)

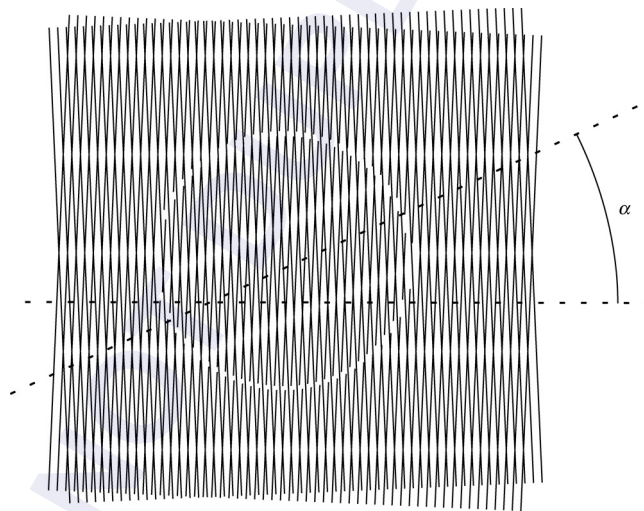


FIGURE 31 Moiré pattern as produced in a moiré deflectometer. (From Malacara.¹)

under test and produces a set of autoimages at the image side of the lens. To determine the lens focal length, another ruling coincident to the autoimages at the image plane is used.^{73,74}

Fourier Transforms The Fourier-transforming property of a lens can be used in the focal-length determination. Horner⁷⁵ measured the diffraction pattern at the focal plane produced by a slit. For this method, the light beam does not have to be perfectly collimated.

Microlenses Micro-lenses applications require new methods for small focal length determination. The propagated Gaussian beam of a laser can be analyzed.^{76,77} In this method, a lens is placed at the laser beam waist, then the propagating beam is measured to determine the focal length.

Fiber Optics A clever method used to automatically find the position of the focus has been described by Howland and Proll.⁷⁸ They used optical fibers to illuminate the lens in an autocollimating configuration, and the location of the image was also determined using optical fibers.

12.6 REFERENCES

1. D. Malacara, (ed.) *Optical Shop Testing*, 3d ed., John Wiley and Sons, New York, 2007.
2. A. L. Bloom, "Gas Lasers and Their Application to Precise Length Measurements," in E. Wolf (ed.), *Progress in Optics*, vol. IX, North Holland, Amsterdam, 1971.
3. D. T. Goldman, "Proposed New Definition of the Meter," *J. Opt. Soc. Am.* **70**: 1640–1641 (1980).
4. P. Giacomo, "Metrology and Fundamental Constants." *Proc. Int. School of Phys. "Enrico Fermi," course 68*, North-Holland, Amsterdam, 1980.
5. S. Cundiff, J. Ye, and J. Hall, "Rulers of Light," *Sci. Am.* **298**(4): 52–59 (2008).
6. D. C. Baird, *Experimentation*, Prentice-Hall, New Jersey, 1962.
7. J. C. Gibbins, *The Systematic Experiment*, Cambridge Univ. Press, Cambridge, 1986.
8. F. B. Patrick, "Military Optical Instruments," in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. V, Academic Press, New York, 1969, chap. 7.
9. M. S. Dickson and D. Harkness, "Surveying and Tracking Instruments," in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. V, Academic Press, New York, 1969, chap. 8.
10. W. J. Smith, *Modern Optical Engineering*, 2d ed., McGraw-Hill, New York, 1990.
11. A. Sona, "Lasers in Metrology," in F. T. Arecchi and E. O. Schulz-Dubois (eds.), *Laser Handbook*, vol. 2, North Holland, Amsterdam, 1972.
12. R. Dändliker, R. Thalmann, and D. Prongué, "Two-Wavelength Laser Interferometry Using Superheterodyne Detection," *Opt. Lett.* **13**: 339–342 (1988).
13. R. Dändliker, Y. Salvadé, and E. Zimmermann, "Distance Measurement by Multiple-Wavelength Interferometry," *J. Opt.* **29**: 105–114 (1998).
14. Y. Salvadé, N. Schuler, S. Lévêque, and S. Le Floch, "High-Accuracy Absolute Distance Measurement Using Frequency Comb Referenced Multiwavelength Source," *Appl. Opt.* **47**(14): 2715–2720 (2008).
15. K. Minoshima and H. Matsumoto, "High-Accuracy Measurement of 240-m Distance in an Optical Tunnel by Use of a Compact Femtosecond Laser," *Appl. Opt.* **39**: 5512–5517 (2000).
16. M. L. Stitch, "Laser Rangefinding," in F. T. Arecchi and E. O. Schulz-Dubois (eds.), *Laser Handbook*, vol. 2, North Holland, Amsterdam, 1972.
17. T. Tsuruta and Y. Ichihara, "Accurate Measurement of Lens Thickness by Using White-Light Fringes," *Jpn. J. Appl. Phys. Suppl.* **14-1**: 369–372 (1975).
18. J. Bruning, "Fringe Scanning," in D. Malacara (ed.), *Optical Shop Testing*, 1st ed., John Wiley and Sons, New York, 1978.
19. C. Steinmetz, R. Burgoon, and J. Herris, "Accuracy Analysis and Improvements for the Hewlett-Packard Laser Interferometer System," *Proc. SPIE* **816**: 79–94 (1987).
20. N. A. Massie, and J. Caulfield, "Absolute Distance Interferometry," *Proc. SPIE* **816**: 149–157 (1987).
21. E. R. Peck and S. W. Obetz, "Wavelength or Length Measurement by Reversible Fringe Counting," *J. Opt. Soc. Am.* **43**: 505–507 (1953).
22. W. R. C. Rowley, "Some Aspects of Fringe Counting in Loser Interferometers," *IEEE Trans. on Instr. and Measur.* **15**(4): 146–149 (1966).
23. S. Minkowitz and W. Reid Smith Vanir, "Laser Interferometer," *Proc. 1st Congress on Laser Applications (Paris)*, *J. Quantum Electronics* **3**: 237 (1967).
24. G. M. Burgwald and W. P. Kruger, "An Instant-On Laser for Length Measurements," *HPJ* **21**: 2 (1970).
25. J. N. Dukes and G. B. Gordon, "A Two-Hundred-Foot Yardstick with Graduations Every Microinch," *HPJ* **21**: 2 (1970).
26. J. H. McLeod, "The Axicon : A New Type of Optical Element," *J. Opt. Soc. Am.* **44**: 592–597 (1954).

27. A. W. Young, "Optical Workshop Instruments," in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. 4, Academic Press, New York, 1967, chap. 7.
28. C. Deve, *Optical Workshop Principles*, T. L. Tippell (transl.), Hilger and Watts, London, 1945.
29. K. J. Hume, *Metrology with Autocollimators*, Hilger and Watts, London, 1965.
30. M. P. Kothiyal and R. S. Sirohi, "Improved Collimation Testing Using Talbot Interferometry," *Appl. Opt.* **26**: 4056–4057 (1987).
31. R. E. Noble, "Some Parameter Measurements," in D. Malacara (ed.), *Optical Shop Testing*, 1st ed., John Wiley and Sons, New York, 1978.
32. K. H. Carnell and W. T. Welford, "A Method for Precision Spherometry of Concave Surfaces," *J. Phys. E.* **4**: 1060–1062 (1971).
33. D. H. Rank, "Measurement of the Radius of Curvature of Concave Spheres," *J. Opt. Soc. Am.* **36**: 108–110 (1946).
34. D. F. Horne, *Dividing, Ruling and Mask Making*, Adam Hilger, London, 1974, chap. VII.
35. R. Kingslake, *Optical System Design*, Academic Press, New York, 1983, chap. 13.
36. V. Met, "Determination of Small Wedge Angles Using a Gas Laser," *Appl. Opt.* **5**: 1242–1244 (1966).
37. G. W. Leppelmeier and D. J. Mullenhoff, "A Technique to Measure the Wedge Angle of Optical Flats," *Appl. Opt.* **9**: 509–510 (1970).
38. J. H. Wasilik, T. V. Blomquist, and C. S. Willett, "Measurement of Parallelism of the Surfaces of a Transparent Sample Using Two-Beam Non-Localized Fringes Produced by a Laser," *Appl. Opt.* **10**: 2107–2112 (1971).
39. D. Malacara, (ed.), *Optical Shop Testing*, 2d ed., John Wiley and Sons, New York, 1992.
40. D. Malacara and O. Harris, "Interferometric Measurement of Angles," *Appl. Opt.* **9**: 1630–1633 (1970).
41. D. Tentori and M. Celaya, "Continuous Angle Measurement with a Jamin Interferometer," *Appl. Opt.* **25**: 215–220 (1986).
42. E. Stijns, "Measuring Small Rotation Rates with a Modified Michelson Interferometer," *Proc. SPIE* **661**: 264–266 (1986).
43. P. Shi and E. Stijns, "New Optical Method for Measuring Small Angle Rotations," *Appl. Opt.* **27**: 4342–4344 (1988).
44. G. D. Chapman, "Interferometric Angular Measurement," *Appl. Opt.* **13**: 1646–1651 (1974).
45. B. K. Johnson, *Optics and Optical Instruments*, Dover, New York, 1947, chaps. II and VIII.
46. L. C. Martin, *Optical Measuring Instruments*, Blackie and Sons Ltd., London, 1924.
47. F. Twyman, *Prisms and Lens Making*, 2d ed., Hilger and Watts, London, 1957.
48. A. S. DeVany, "Reduplication of a Penta-Prism Angle Using Master Angle Prisms and Plano Interferometer," *Appl. Opt.* **10**: 1371–1375 (1971).
49. A. S. DeVany, "Testing Glass Reflecting-Angles of Prisms," *Appl. Opt.* **17**: 1661–1662 (1978).
50. F. Ratajczyk and Z. Bodner, "An Autocollimation Measurement of the Right Angle Error with the Help of Polarized Light," *Appl. Opt.* **5**: 755–758 (1966).
51. A. M. Tareev, "Testing the Angles of High-Precision Prisms by Means of an Autocollimator and a Mirror Unit," *Sov. J. Opt. Technol.* **52**: 50–52 (1985).
52. D. Malacara and R. Flores, "A Simple Test for the 90 Degrees Angle in Prisms," *Proc. SPIE* **1332**: 678 (1990).
53. C. Ai and K. L. Smith, "Accurate Measurement of the Dihedral Angle of a Corner Cube," *Appl. Opt.* **31**: 519–527 (1992).
54. S. M. Rao, "Method for the Measurement of the Angles of a Tetragonal or Corner Cube Prism," *Opt. Eng.* **41**: 1612–1614 (2002).
55. M. S. Scholl, "Ray Trace through a Corner-Cube Retroreflector with Complex Reflection Coefficients," *J. Opt. Soc. Am. A.* **12**(7): 1589–1592 (1995).
56. B. Jurek, *Optical Surfaces*, Elsevier Scient. Pub. Co., New York, 1977.
57. D. F. Horne, *Optical Production Technology*, Adam Hilger, London, and Crane Russak, New York, 1972, chap. XI.
58. M. C. Gerchman and G. C. Hunter, "Differential Technique for Accurately Measuring the Radius of Curvature of Long Radius Concave Optical Surfaces," *Proc. SPIE* **192**: 75–84 (1979).

59. M. C. Gerchman and G. C. Hunter, "Differential Technique for Accurately Measuring the Radius of Curvature of Long Radius Concave Optical Surfaces," *Opt. Eng.* **19**: 843–848 (1980).
60. D. C. O'Shea and S. A. Tilstra, "Non-Contact Measurements of Refractive Index and Surface Curvature," *Proc. SPIE* **966**: 172–176 (1988).
61. J. D. Evans, "Method for Approximating the Radius of Curvature of Small Concave Spherical Mirrors Using a He-Ne Laser," *Appl. Opt.* **10**: 995–996 (1971).
62. J. D. Evans, "Equations for Determining the Focal Length of On-Axis Parabolic Mirrors by He-Ne Laser Reflection," *Appl. Opt.* **11**: 712–714 (1972).
63. J. D. Evans, "Error Analysis to: Method for Approximating the Radius of Curvature of Small Concave Spherical Mirrors Using a He-Ne Laser," *Appl. Opt.* **11**: 945–946 (1972).
64. A. Cornejo-Rodriguez and A. Cordero-Dávila, "Measurement of Radii of Curvature of Convex and Concave Surfaces Using a Nodal Bench and a He-Ne Laser," *Appl. Opt.* **19**: 1743–1745 (1980).
65. P. E. Klingsporn, "Use of a Laser Interferometric Displacement-Measuring System for Noncontact Positioning of a Sphere on a Rotation Axis through Its Center and for Measuring the Spherical Contour," *Appl. Opt.* **18**: 2881–2890 (1979).
66. R. Díaz-Urribe, J. Pedraza-Contreras, O. Cardona-Nuñez, A. Cordero-Dávila, and A. Cornejo Rodriguez, "Cylindrical Lenses: Testing and Radius of Curvature Measurement," *Appl. Opt.* **25**: 1707–1709 (1986).
67. R. Kingslake, "A New Bench for Testing Photographic Lenses," *J. Opt. Soc. Am.* **22**: 207–222 (1932).
68. P. Bouchaud, and J. A. Cogno, "Automatic Method for Measuring Simple Lens Power," *Appl. Opt.* **21**: 3068 (1982).
69. D. Malacara and Z. Malacara, "Testing and Centering of Lenses by Means of Hartmann Test with Four Holes," *Opt. Eng.* **31**: 1551–1555 (1996).
70. M. Martínez-Corral, W. D. Furlan, A. Pons, and G. Saavedra, *Instrumentos Ópticos y Optométricos. Teoría y Prácticas*, Universitat de Valencia, Valencia, (1998).
71. O. Kafri and I. Glatt, *The Physics of Moiré Metrology*, Wiley Interscience, New York, 1990.
72. I. Glatt and O. Kafri, "Determination of the Focal Length of Non-Paraxial Lenses by Moiré Deflectometry," *Appl. Opt.* **26**: 2507–2508, (1987).
73. D. Malacara-Doblado and D. Malacara-Hernández, "Measuring Convergence or Divergence Power with Moiré fringes," *Proc. SPIE* **2860**: 390–393 (1996).
74. Y. Nakano, and K. Murata, "Talbot Interferometry for Measuring the Focal Length of a Lens," *Appl. Opt.* **24**:19 (1985).
75. J. L. Horner, Collimation Invariant Technique for Measuring the Focal Length of a Lens," *Appl. Opt.* **28**: 1047–1047 (1989).
76. A. A. Camacho, C. Solano, M. Cywiak, G. Martínez-Ponce, and R. Baltazar, "Method for the Determination of the Focal Length of a Microlens," *Opt. Eng.* **39**: 2149–2152 (2000).
77. A. A. Camacho, C. Solano, G. Martínez-Ponce, and R. Baltazar, "Simple Method to Measure the Focal Length of a Lens," *Opt. Eng.* **41**: 2899–2902 (2002).
78. B. Howland and A. F. Proll, "Apparatus for the Accurate Determination of Flange Focal Distance," *Appl. Opt.* **11**: 1247–1251 (1970).

This page intentionally left blank.

DO NOT DUPLICATE

OPTICAL TESTING

Daniel Malacara-Hernández

*Centro de Investigaciones en Optica, A. C.
León, Gto., Mexico*

13.1 GLOSSARY

E	electric field strength
k	radian wave number
r	position
t	time
λ	wavelength
φ	phase
ω	radian frequency

13.2 INTRODUCTION

The requirements for high-quality optical surfaces are more demanding every day. They should be tested in an easier, faster, and more accurate manner. Optical surfaces usually have a flat or a spherical *shape*, but they also may be toroidal or generally aspheric. Frequently, an aspherical surface is a conic of revolution, but an aspherical surface can only be made as good as it can be tested. Here, the field of optical testing will be reviewed. There are some references that the reader may consult for further details.¹

13.3 CLASSICAL NONINTERFEROMETRIC TESTS

Some classical tests will never be obsolete, because they are cheap, simple, and provide qualitative results about the shape of the optical surface or wavefront almost instantly. These are the Foucault or knife-edge test, the Ronchi test, and the Hartmann test. They will be described next.

Foucault Test

The Foucault or knife-edge test was invented by Leon Foucault² in France, to evaluate the quality of spherical surfaces. This test described by Ojeda-Castañeda³ detects the presence of transverse aberrations by intercepting the reflected rays deviated from their ideal trajectory, as Fig. 1 shows. The observer is behind the knife, looking at the illuminated optical surface, with the reflected rays entering the eye. The regions corresponding to the intercepted rays will appear dark, as in Fig. 2.

This test is extremely sensitive. If the wavefront is nearly spherical, irregularities as small as a fraction of the wavelength of the light may be easily detected. This is the simplest and most powerful qualitative test for observing small irregularities and evaluating the general smoothness of the spherical

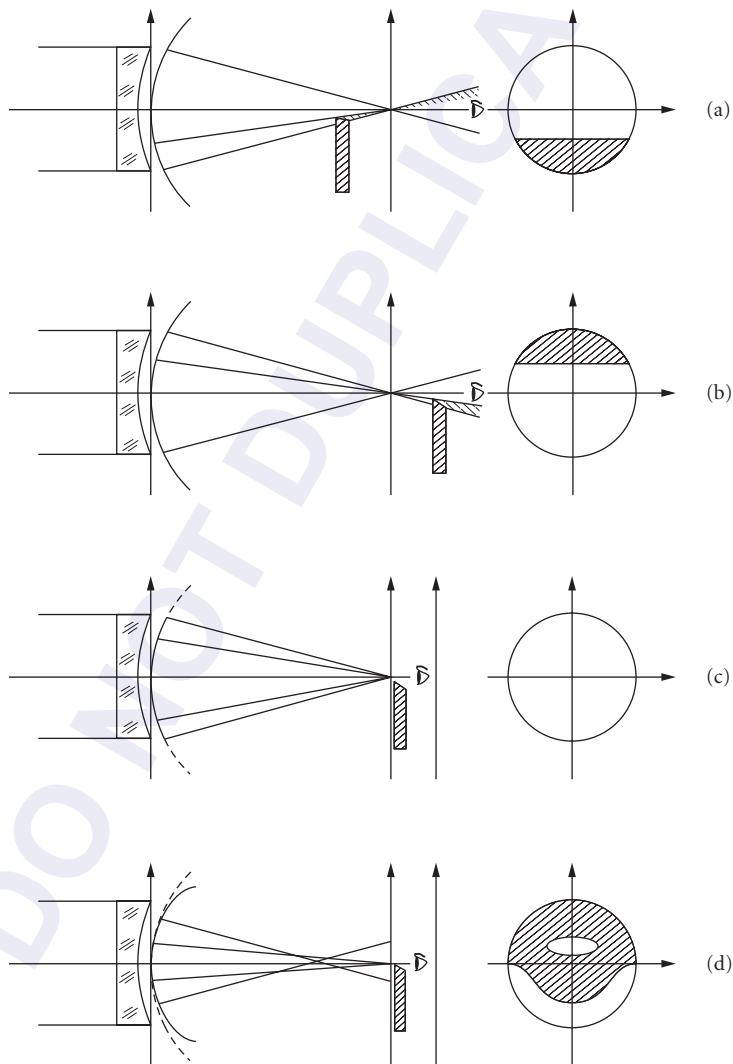


FIGURE 1 Optical schematics for the Foucault test of a spherical mirror at several positions of the knife edge.

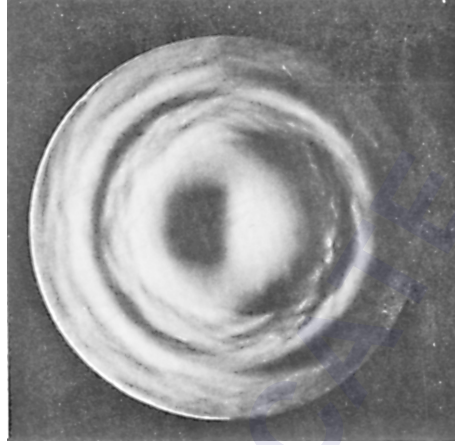


FIGURE 2 An optical surface being examined by the Foucault test. (From Ojeda-Castañeda.³)

surface under test. Any other surface or lens may be tested, as long as it produces an almost spherical wavefront, otherwise, an aberration compensator must be used, as will be described later. Very often a razor blade makes a good, straight, sharp edge that is large enough to cover the focal region.

Ronchi Test

Vasco Ronchi⁴ invented his famous test in Italy in 1923. A coarse ruling (50–100 lines per inch) is placed in the convergent light beam reflected from the surface under test, near its focus. The observer is behind the ruling, as Fig. 3 shows, with the light entering the eye. The dark bands in the ruling intercept light, forming shadows on the illuminated optical surface. These shadows will be straight and parallel only if the reflected wavefront is perfectly spherical. Otherwise, the fringes will be curves whose shape and separation depends on the wavefront deformations. The Ronchi test

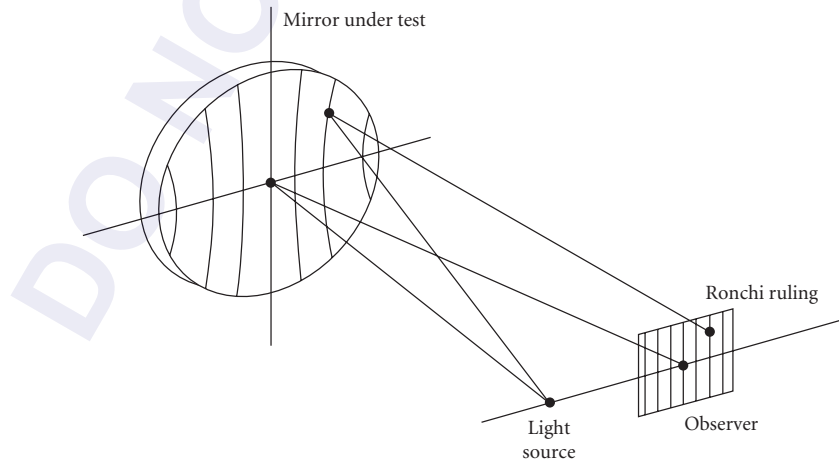


FIGURE 3 Testing a concave surface by means of the Ronchi test.

measures the transverse aberrations in the direction perpendicular to the slits on the grating. The wavefront deformations $W(x, y)$ are related to the transverse aberrations $TA_x(x, y)$ and $TA_y(x, y)$ by the following well-known relations:

$$TA_x(x, y) = -r \frac{\partial W(x, y)}{\partial x} \quad (1)$$

and

$$TA_y(x, y) = -r \frac{\partial W(x, y)}{\partial y} \quad (2)$$

where r is the radius of curvature of the wavefront $W(x, y)$. Thus, if we assume a ruling with period d , the expression describing the m th fringe on the optical surface is given by

$$\frac{\partial W(x, y)}{\partial x} = -\frac{md}{r} \quad (3)$$

Each type of aberration wavefront has a characteristic Ronchi pattern, as shown in Fig. 4; thus, the aberrations in the optical system may be easily identified, and their magnitude estimated. We may interpret the Ronchi fringes not only as geometrical shadows, but also as interferometric fringes, identical with those produced by a lateral shear interferometer.

Hartmann Test

J. Hartmann⁵ invented his test in Germany. It is one of the most powerful methods to determine the figure of a concave spherical or aspherical mirror. Figure 5 shows the optical configuration used

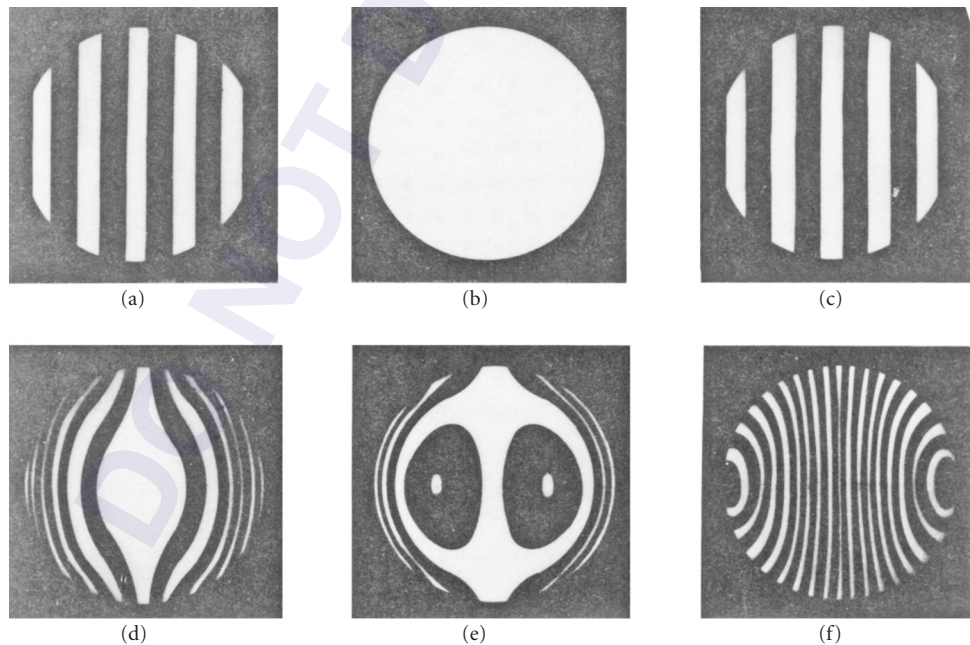


FIGURE 4 Typical Ronchi patterns for a spherical and a parabolic mirror for different positions of the Ronchi ruling.

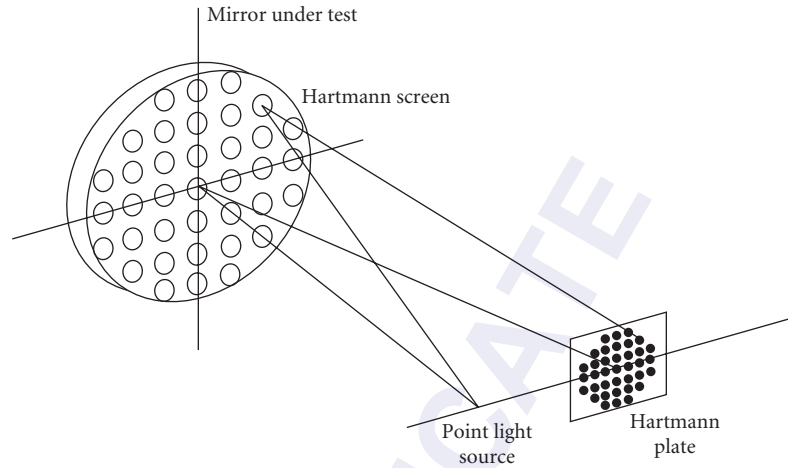


FIGURE 5 Optical arrangement to perform the Hartmann test.

in this test, where a point light source illuminates the optical surface, with its Hartmann screen in front of it. The light beams reflected through each hole on the screen are intercepted on a photographic plate near the focus. Then, the position of the recorded spots is measured to find the value of the transverse aberration on each point. If the screen has a rectangular array of holes, the typical Hartmann plate image for a parabolic mirror looks like that in Fig. 6. The wavefront $W(x, y)$ may be obtained from integration of Eqs. (1) and (2) as follows:

$$W(x, y) = -\frac{1}{r} \int_0^x TA_x(x, y) dx \quad (4)$$

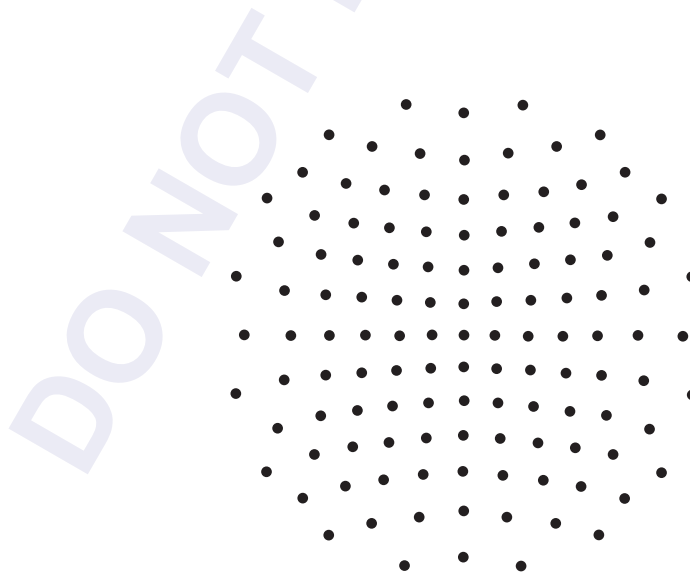


FIGURE 6 Array of spots in a Hartmann plate of a parabolic mirror.

and

$$W(x, y) = -\frac{1}{r} \int_0^y TA_y(x, y) dy \quad (5)$$

After numerical integration of the values of the transverse aberrations, this test provides the concave surface shape with very high accuracy. If the surface is not spherical, the transverse aberrations to be integrated are the difference between the measured values and the ideal values for a perfect surface. Extended, localized errors, as well as asymmetric errors like astigmatism, are detected with this test. The two main problems of this test are that small, localized defects are not detected if they are not covered by the holes on the screen. Not only is this information lost, but the integration results will be false if the localized errors are large. The second important problem of the Hartmann test is that it is very time consuming, due to the time used in measuring all the data points on the Hartmann plate. These problems are avoided by complementing this test with the Foucault test, using an Offner compensator, in order to be sure about the smoothness of the surface (discussed under “Measuring Aspherical Wavefronts”). Various stratagems are available to speed the process. These include modulating the light at different frequencies at each of the holes. Variations also include measuring in front of, behind, or at the focus to get slope information. This technique can be considered an experimental ray trace.

Hartmann-Shack Test

Platt and Shack,⁶ proposed using a lenticular screen, instead of a screen with an array of holes, as illustrated in Fig. 7. This is a simple but important modification from the classic Hartmann. Some differences are

- (a) In the Hartmann test the pattern is obtained in focused convergent light beam, near the focus. On the other hand, in the Hartmann-Shack the test is made in a nearly collimated beam of light.
- (b) A practical advantage of the Hartmann-Shack method is that any positive or negative power can be easily detected and measured.

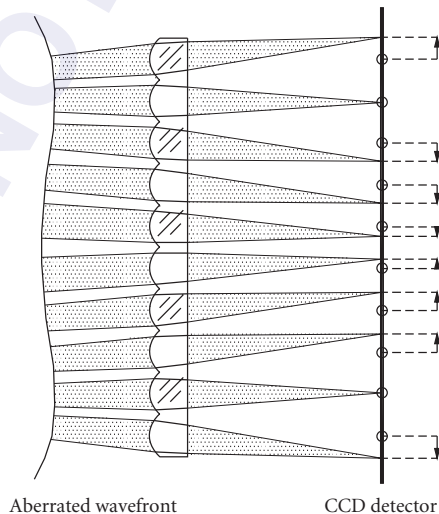


FIGURE 7 Hartmann-Shack Test.

- (c) A second advantage is that each of the spots is individually focused on the detector, making the light energy density of the spot higher than in the Hartmann test.
- (d) The Hartmann-Shack lenticular screen can be made with two identical layers of cylindrical lenses perpendicular to each other, or with a lenslet arrays in molded plastic, glass, or fused silica.

If the wavefront is flat, the light beam passing through each lens is focused close to the optical axis of each lenslet. Since the lens array is not perfect, the lenticular array must be previously calibrated with a reference well-known flat wavefront.

The spot displacement on the detector is equal to the wavefront slope multiplied by the focal length of the lenslet, thus, a shorter focal length will give a greater dynamic range but a reduced angular sensitivity. The optimum focal length depends on the application.

13.4 INTERFEROMETRIC TESTS

Classical geometrical tests are very simple, but they do not provide the accuracy of the interferometric tests. Quite generally, an interferometric test produces an interferogram by producing the interference between two wavefronts. One of these two wavefronts is the wavefront under test. The other wavefront is either a perfectly spherical or flat wavefront, or a copy of the wavefront under test.

When the second wavefront is perfectly spherical or flat, this wavefront acts as a reference. The separation between the two wavefronts, or optical path difference $OPD(x, y)$, is a direct indication of the deformations $W(x, y)$ of the wavefront under test. Then, we may simply write $W(x, y) = OPD(x, y)$. There are many types of interferometers producing interferograms of these type of interferograms, for example, the Twyman-Green and the Fizeau interferometers. Some other interferometers can be considered as modifications of these two basic interferometers, such as the Point Diffraction and the Burch interferometers, and many others that will not be described.

Twyman-Green Interferometer

The Twyman-Green interferometer is illustrated in Fig. 8. The light from a monochromatic point light source is collimated to produce a flat wavefront. Then, the two interfering wavefronts are generated by means of a partially reflective and partially transmitting glass plate, called beam splitter.

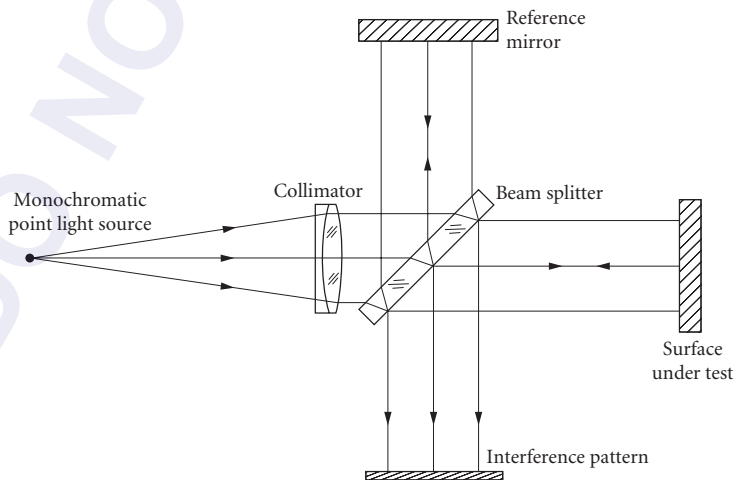


FIGURE 8 Twyman-Green interferometer.

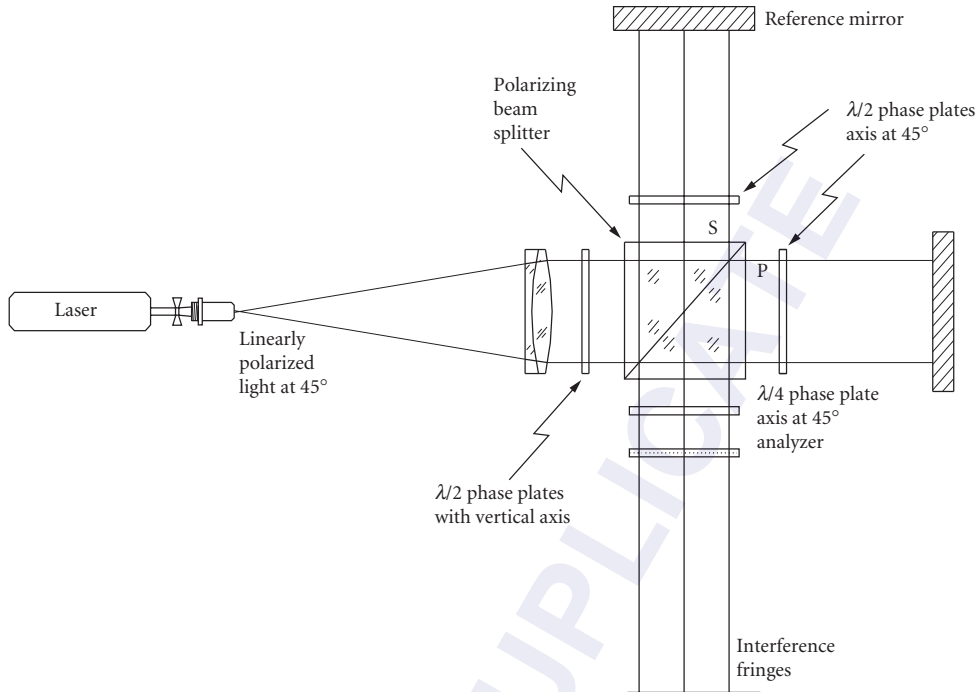


FIGURE 9 Twyman-Green interferometer with polarized beam splitter.

After the beam splitter two flat wavefronts travel in orthogonal directions, one of them to the flat reference mirror and the other to the surface or optical element under test. After returning to the beam splitter, the two wavefronts are recombined to produce an interference pattern. The beam splitter can be oriented at 45° as in Fig. 8, but sometimes a different angle is chosen, for example, a Brewster angle to avoid the reflected beam from the second face on the beam splitter.

Instead of a plane parallel beam splitter, sometimes a polarizing cube beam splitter is used, as in Fig. 9. In this system the plane wavefront entering the beam splitter is linearly polarized at an angle of 45° with respect to the plane of the interferometer. Then, the two wavefronts exiting the cube in orthogonal directions will also be linearly polarized but one of them in the vertical plane and the other in the horizontal plane. A $\lambda/4$ phase plate is located at each of the two exiting faces on the cube. These phase plates produce circularly polarized beams, one going to the reference mirror and the other to the surface under test, but one is right handed and the other is left handed. When arriving back to the cube, after passing twice through the phase plates, the two beams will be again linearly polarized, one in the vertical plane and the other in the horizontal plane. However, these planes of polarization will be orthogonal to the planes of polarization when the wavefront exited the cube. Thus, the two wavefronts are sent to the observation plane and not back to the light source.

Two important facts should be noticed: (a) that the two wavefronts are orthogonally polarized and hence they can not interfere and (b) that there are not any light beams going back to the light source to produce a complementary pattern as in the classic Twyman-Green interferometer. In order to produce observable fringes a linear polarizer should be placed just before the observing plane.

This interferometer with a polarizing beam splitter has many practical advantages.

Fizeau Interferometer

A Fizeau interferometer, illustrated in Fig. 10, is also a two-beam system like the Twyman-Green interferometer. The main difference is that the plate beam splitter is not at 45° with the illuminated

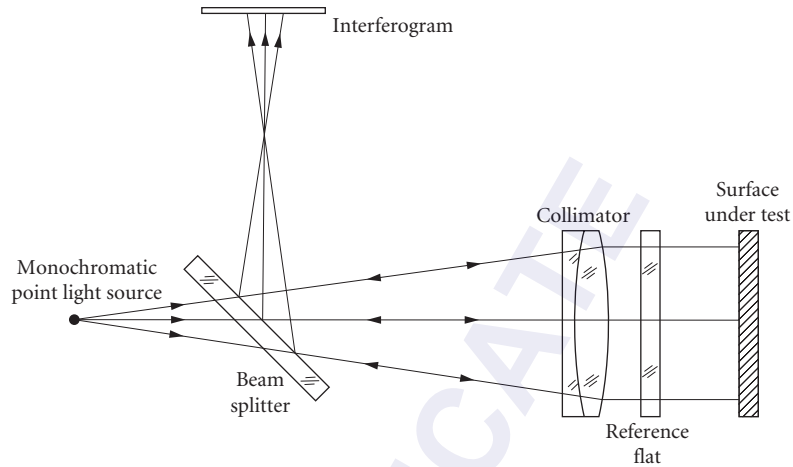


FIGURE 10 Fizeau interferometer.

collimated beam, but perpendicular to its incidence path. Some important practical advantages are that (a) the beam splitter can be smaller for the same aperture, (b) it is more compact, and (c) it is easier to align.

Figure 11 shows some typical interferograms for the Seidel primary aberrations obtained with a Twyman-Green or a Fizeau interferometer.¹ The mathematical analysis of interferograms for wavefronts with arbitrary deformations is a research topic of great interest that has been described in a large number of publications.⁷

Common Path Interferometer

A common path interferometer is one for which the two interfering beams travel the same paths. The optical path difference at the center of the optical axis is zero and cannot be modified. Thus interference with a white light source can be easily achieved. An example of this kind of interferometer is the point-diffraction interferometer first described by Linnik in 1933 and later rediscovered by Smart and Strong.⁸ The lens of optical element under test focuses the light at the center of a small diffracting plate as illustrated in Fig. 12. This diffracting plate is coated with a thin partially transmitting film with a small uncoated disk at its center. The diameter of the central clear disk is about the size of the diffraction air disk produced by a perfect optical system. If the wavefront from the system is not spherical but distorted, the focused spot would be larger than the central disk. Then, two wavefronts are produced. One is a reference spherical wavefront arising from the light diffracted at the central disk. The undiffracted light passing outside of the disk is the wavefront under test.

The fringe patterns in the point diffraction interferometer are identical to those produced by the Twyman-Green interferometer.

Lateral Shearing Interferometers

When the second wavefront is not perfectly flat or spherical, but a copy of the wavefront under test, its relative dimensions or orientation must be changed (sheared) in some way with respect to the wavefront under test. Otherwise, no information about the wavefront deformations is obtained, because the fringes will always be straight and parallel independent of any aberrations. There are several kinds of shearing interferometers, depending on the kind of transformation applied to the reference wavefront.

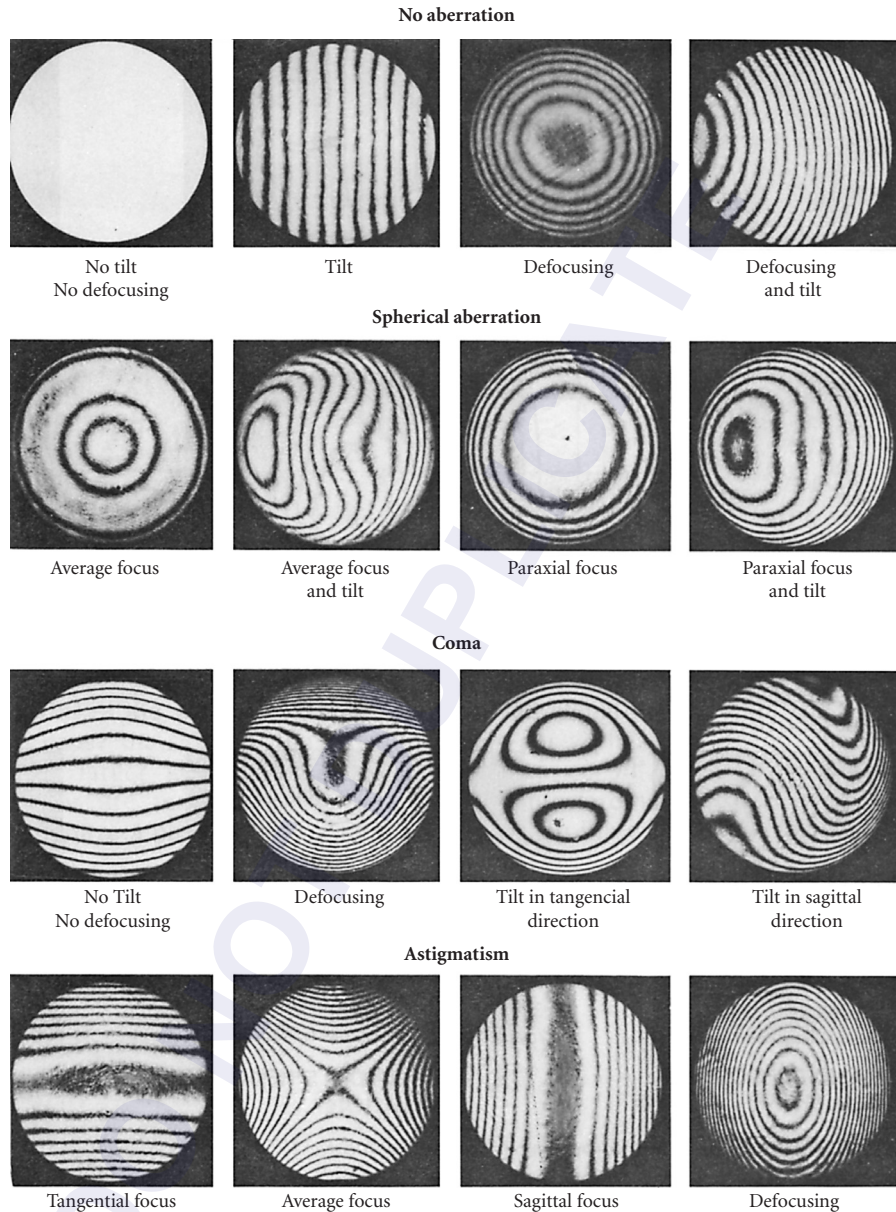


FIGURE 11 Twyman-Green interferograms. (From Malacara.¹)

The most popular of these instruments is the lateral shearing interferometer, with the reference wavefront laterally displaced with respect to the other, as in the interferograms in Fig. 13 shows. The optical path difference $OPD(x, y)$ and the wavefront deformations $W(x, y)$ are related by

$$OPD(x, y) = W(x, y) - W(x - S, y) \quad (6)$$

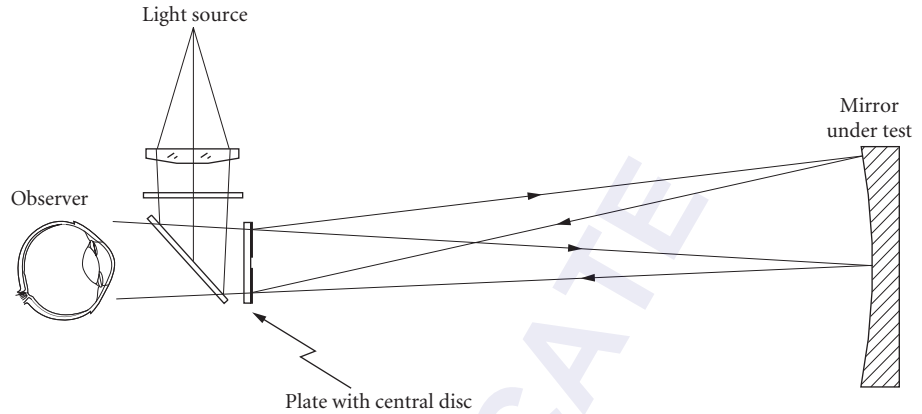


FIGURE 12 Point diffraction interferometer.

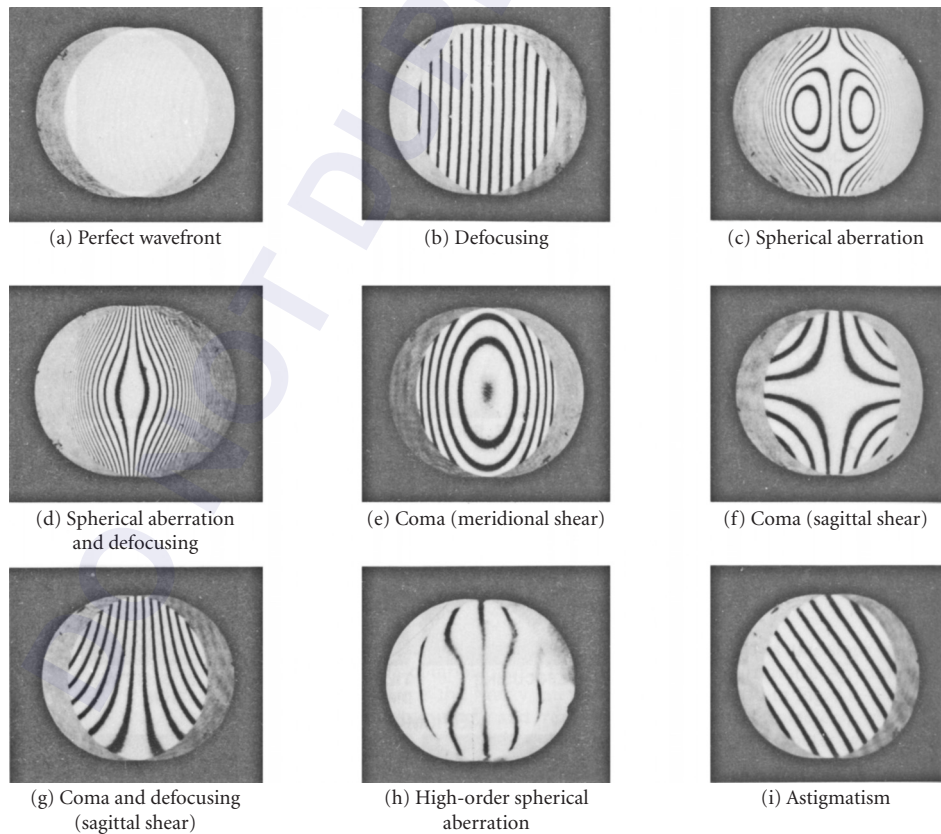


FIGURE 13 Laterally sheared interferograms. (From Malacara.¹)

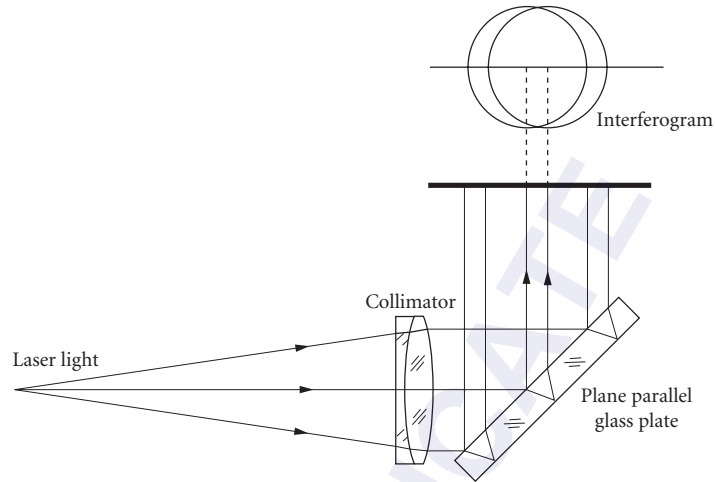


FIGURE 14 Murty's lateral shear interferometer.

where S is the lateral shear of one wavefront with respect to the other. If the shear is small with respect to the diameter of the wavefront, this expression may be approximated by

$$\text{OPD}(x, y) = -S \frac{\partial W(x, y)}{\partial x} = -\frac{S}{r} T A_x(x, y) \quad (7)$$

This relation suggests that the parameter being directly measured is the slope in the x direction of the wavefront (x component $T A_x$ of the transverse aberration). An example of a lateral shear interferometer is the Murty interferometer, illustrated in Fig. 14.

Radial, Rotational, and Reversal Shearing Interferometers

There are also radial, rotational, and reversal shearing interferometers, where the interfering wavefronts are as illustrated in Fig. 15. A radial shear interferometer with a large shear approaches an interferometer with a perfect reference wavefront. These interferometers procedure fringe patterns by the interference of two wavefronts with the same aberrations and deformations. The difference between the two interfering wavefronts is their size or orientation.

The optical path in the radial shear interferometer can be represented by

$$\text{OPD}(x, y) = W(x, y) - W(\rho x, \rho y) \quad (8)$$

A typical radial shear interferogram is in Fig. 16.

In the rotational shear interferometer. The two interfering wavefronts have the same size, but one of those is rotated with respect to the other. In the particular case on a 180° rotation the sensitivity of the interferometer is zero for symmetrical (even) aberrations, like spherical aberration. However, the sensitivity is doubled for antisymmetrical (odd) aberrations, like coma.

In a reversing shear interferometer one of the two wavefronts is reversed with respect to the other about any diameter on the wavefront's pupil. As in the rotational shear interferometer, the sensitivity to aberrations is symmetric with respect to the axis of reversion. Also, the sensitivity to aberration antisymmetric about the axis of reversion is doubled.

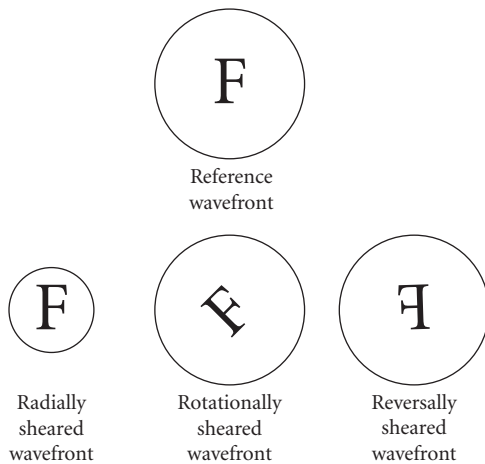


FIGURE 15 Wavefronts in radial, rotational, and reversal shear interferometers.

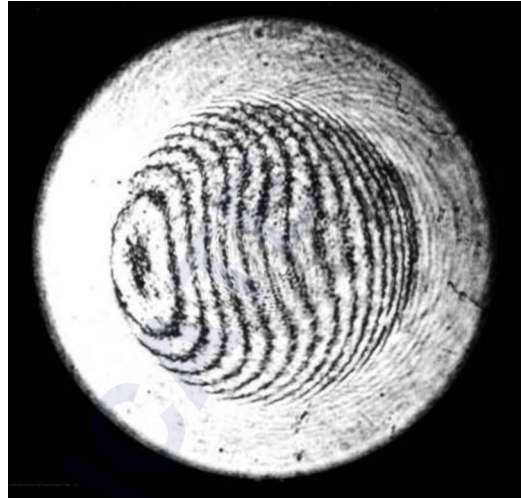


FIGURE 16 A radial shear interferogram.

13.5 INCREASING THE SENSITIVITY OF INTERFEROMETERS

The sensitivity of interferometers is a small fraction of the wavelength being used (about $\lambda/20$). There are several methods to increase this sensitivity, but the most common methods will now be described.

Multiple-Reflection Interferometers

A method to increase the sensitivity of interferometric methods is to use multiple reflections, as in the Fabry-Perot interferometer. The Newton as well as the Fizeau interferometers can be made multiple-reflection interferometers by coating the reference surface and the surface under test with a high-reflection film. Then, the fringes are greatly narrowed and their deviations from straightness are more accurately measured.⁹

Multiple-Pass Interferometers

Another method to increase the sensitivity of interferometers is by double, or even multiple, pass. An additional advantage of double-pass interferometry is that the symmetrical and antisymmetrical parts of the wavefront aberration may be separated. This makes their identification easier, as Hariharan and Sen¹⁰ have proved. Several arrangements have been devised to use multiple pass.¹¹

Zernike Tests

The Zernike phase-contrast method is another way to improve the sensitivity of an interferometer to small aberrations. It was suggested by Zernike as a way to improve the knife-edge test.¹² There are several versions of this test. The basic principle in all of them is the introduction of a

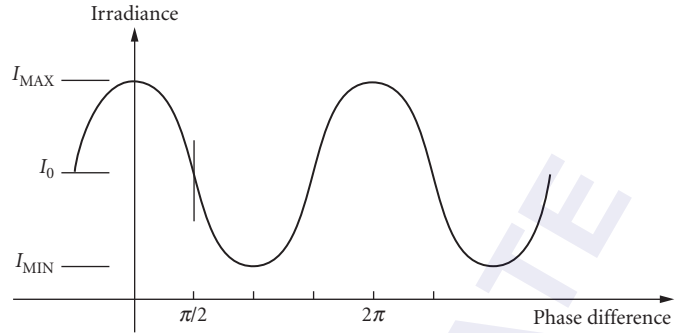


FIGURE 17 Irradiance in an interference pattern, as a function of the phase difference between the two interfering waves.

phase difference equal to $\lambda/2$ between the wavefront under test and the reference wavefront. To understand why this phase difference is convenient, let us consider two interfering beams and irradiances $I_1(x, y)$ and $I_2(x, y)$ and a phase $\phi(x, y)$ between them. The final irradiance $I(x, y)$ in the interferogram is given by

$$I(x, y) = I_1(x, y) + I_2(x, y) + 2\sqrt{I_1(x, y)I_2(x, y)} \cos\phi(x, y) \quad (9)$$

Thus, the irradiance $I(x, y)$ of the combination would be a sinusoidal function of the phase, as illustrated in Fig. 17. If the phase difference is zero for a perfect wavefront, deformations of the wavefront smaller than the wavelength of the light will not be easy to detect, because the slope of the function is zero for a phase near zero. The slope of this function is larger and linear for a phase value of 90° . Thus, the small wavefront deformations are more easily detected if the interferometer is adjusted, so that the wavefronts have a phase difference equal to 90° when the wavefront under test is perfect.

13.6 INTERFEROGRAM EVALUATION

An interferogram may be analyzed in several manners. One way begins by measuring several points on the interferogram, on top of the fringes. Then, the wavefront values between the fringes are interpolated. Another way uses a Fourier analysis of the interferogram. A third method interprets the fringe deformations as a phase modulation.

Fixed Interferogram Evaluation

Once the interferogram has been formed, a quantitative evaluation of it is a convenient method to find the wavefront deformations. The fixed interferogram evaluation by fringe measurements is done by measuring the position of several data points located on top of the fringes. These measurements are made in many ways, for example, with a measuring microscope, with a digitizing tablet, or with a video camera connected to a computer.

The fringe centers can be located either manually, using a digitizing tablet, or automatically, with the computer directly examining a single fringe image that has been captured using a digital frame grabber. After locating the fringe centers, fringe order numbers must be assigned to each point. The wavefront can then be characterized by direct analysis of the fringe centers. If desired, instead of global interpolation, a local interpolation procedure may be used.

To analyze the fringes by a computer, they must first be digitized by locating the fringe centers, and assigning fringe order numbers to them. The optical path difference (OPD) at the center of any fringe is a multiple m of the wavelength ($OPD = m\lambda$), where m is the fringe order. To obtain the

wavefront deformation, only the relative values of the fringe order are important. So any value of the fringe order may be assigned to the first fringe being measured. However, for the second fringe, it may be increased or decreased by one. This choice affects only the sign of the OPD. An important disadvantage of the fixed interferogram analysis is that the sign of the OPD cannot be obtained from the interferogram alone. This information can be retrieved if the sign of any term in the wavefront deformation expression, like defocusing or tilt, is previously determined when taking the interferogram.

Fringes have been digitized using scanners,¹³ television cameras,¹⁴ photoelectric scanners, and digitizing tablets. Review articles by Reid^{15,16} give useful references for fringe digitization using television cameras.

Global and Local Interpolation of Interferograms

After the measurements are made, the wavefront is computed with the measured points. The data density depends on the density of fringes in the interferogram. Given a wavefront deformation, the ratio of the fringe deviations from straightness to the separation between the fringes remains a constant, independently of the number of fringes introduced by tilting of the reference wavefront. If the number of fringes is large due to a large tilt, the fringes look more straight than if the number of fringes is small. Thus, the fringe deviations may more accurately be measured if there are few fringes in the interferogram. Thus, information about many large zones is lost. A way to overcome this problem is to interpolate intermediate values by any of several existing methods. One method is to fit the wavefront data to a two-dimensional polynomial with a least-squares fitting, as described by Malacara et al.¹⁷ or by using splines as described by Hayslett and Swantner¹⁸ and Becker et al.¹⁹ Unfortunately, this procedure has many problems if the wavefront is very irregular. The values obtained with the polynomial may be wrong, especially near the edge, or between fringes if the wavefront is too irregular.

The main disadvantage of global fits is that they smooth the measured surface more than desired. Depending on the degree of the polynomial, there will only be a few degrees of freedom to fit many data points. It is even possible that the fitted surface will pass through none of the measured points. If the surface contains irregular features that are not well described by the chosen polynomial, such as steps or small bumps, the polynomial fit will smooth these features. Then, they will not be visible in the fitted surface.

Global interpolation is done by least-squares fitting the measured data to a two-dimensional polynomial in polar coordinates. The procedure to make the least-squares fitting begins by defining the variance of the discrete wavefront fitting as follows:

$$\sigma = \frac{1}{N} \sum_{i=1}^N [W_i' - W(\rho_i, \theta_i)]^2 \quad (10)$$

where N is the number of data points, W_i is the measured wavefront deviation for data point i , and $W(\rho_i, \theta_i)$ is the functional wavefront deviation after the polynomial fitting. The only requirement is that this variance or fit error is minimized. It is well known that the normal least-squares procedure leads to the inversion of an almost singular matrix. Then, the round-off errors will be so large that the results will be useless. To avoid this problem, the normal approach is to fit the measured points to a linear combination of polynomials that are orthogonal over the discrete set of data points. Thus, the wavefront is represented by

$$W(\rho_i, \theta_i) = \sum_{n=1}^L B_n V_n(\rho_i, \theta_i) \quad (11)$$

$V(\rho, \theta)$ are polynomials of degree r and not the monomials x . These polynomials satisfy the orthogonality condition

$$\sum_{i=1}^N V_n(\rho_i, \theta_i) V_m(\rho_i, \theta_i) = F_n \rho_{nm} \quad (12)$$

where $F_n = \sum V_n^2$.

The advantage of using these orthogonal polynomials is that the matrix of the system becomes diagonal and there is no need to invert it.

The only problem that remains is to obtain the orthogonal polynomials by means of the Gram-Schmidt orthogonalization procedure. It is important to notice that the set of orthogonal polynomials is different for every set of data points. If only one data point is removed or added, the orthogonal polynomials are modified. If the number of data points tends to infinity and they are uniformly distributed over a circular pupil with unit radius, these polynomials $V_r(\rho, \theta)$ approach the Zernike polynomials.²⁰

Several properties of orthogonal polynomials make them ideal for representing wavefronts, but the most important of them is that we may add or subtract one or more polynomial terms without affecting the fit coefficients of the other terms. Thus, we can subtract one or more fitted terms—defocus, for example—without having to recalculate the least-squares fit. In an interferometric optical testing procedure the main objective is to determine the shape of the wavefront measured with respect to a best-fit sphere. Nearly always it will be necessary to add or subtract some terms.

The only problem with these orthogonal polynomials over the discrete set of data points is that they are different for every set of data points. A better choice for the wavefront representation is the set of Zernike polynomials, which are orthogonal on the circle with unit radius, as follows:

$$\int_0^1 \int_0^{2\pi} U_n(\rho, \theta) U_m(\rho, \theta) \rho d\rho d\theta = F_{nm} \delta_{nm} \quad (13)$$

These polynomials are not exactly orthogonal on the set of data points, but they are close to satisfying this condition. Therefore, it is common to transform the wavefront representation in terms of the polynomials V_n to another similar representation in terms of Zernike polynomials $U_n(\rho, \theta)$, as

$$W(\rho, \theta) = \sum_{n=1}^L A_n U_n(\rho, \theta) \quad (14)$$

Fourier Analysis of Interferograms

A completely different way to analyze an interferogram without having to make any interpolation between the fringes is by a Fourier analysis of the interferogram. An interpolation procedure is not needed because the irradiance at a fine two-dimensional array of points is measured and not only at the top of the fringes. The irradiance should be measured directly on the interferogram with a two-dimensional detector or television camera, and not on a photographic picture. Womack,²¹ Macy,²² Takeda et al.,²³ and Roddier and Roddier²⁴ have studied in detail the Fourier analysis of interferograms to obtain the wavefront deformations.

Consider an interferogram produced by the interference of the wavefront under test and a flat reference wavefront, with a large tilt between them, as in the interferogram in Fig. 18. The tilt is about the y axis, increasing the distance between the wavefronts in the x direction. The picture of this interferogram may be thought of as a hologram reconstructing the wavefront. Thus, three wavefronts (images) are generated when this hologram is illuminated with a flat wavefront. In order to have complete separation between these images, the tilt between the wavefronts must be large enough, so that the angle between them is not zero at any point over the interferogram. This is equivalent to saying that the fringes must be open, and never cross any line parallel to the x axis more than once. One image is the wavefront under test and another is the conjugate of this wavefront.

If the tilt between the wavefront is θ and the wavefront shape is $W(x, y)$, the irradiance is given by

$$I(x, y) = I_1(x, y) + I_2(x, y) + 2\sqrt{I_1(x, y)I_2(x, y)} \cos(\phi_0 + kx \sin\theta + kW(x, y)) \quad (15)$$

where $k = 2\pi/\lambda$. This expression may be rewritten as

$$I = [I_1 + I_2] + \sqrt{I_1 I_2} e^{i(kx \sin\theta + kW)} + \sqrt{I_1 I_2} e^{-i(kx \sin\theta + kW)} \quad (16)$$

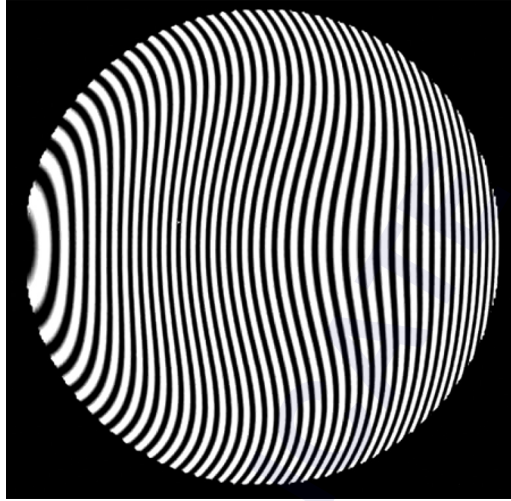


FIGURE 18 Interferogram with a large tilt (linear carrier) to avoid closed fringes.

The first term represents the zero order, the second is the real image, and the third is the virtual image. We also may say that the Fourier transform of the interferogram is formed by a Dirac impulse $\delta(f)$ at the origin and two terms shifted from the origin, at frequencies $+f_0$ and $-f_0$. The quantity f is the spatial frequency, defined by the tilt between the reference wavefront and the wavefront under test ($f = \sin \theta / \lambda$). These terms may be found by taking the Fourier transform of the interferogram. The term at f_0 is due to the wavefront under test. This wavefront may be obtained by taking the Fourier transform of this term, mathematically isolated from the others. This method is performed in a computer by using the fast Fourier transform. The undesired terms are simply eliminated before taking the second fast Fourier transform in order to obtain the wavefront.

Direct Interferometry

This is another method to obtain the wavefront from an interferogram without the need of any interpolation. As in the Fourier method, the image of the interferogram is directly measured with a two-dimensional detector or television camera. The interferogram must have many fringes, produced with a large tilt between the wavefronts. The requirements for the magnitude of this tilt are the same as in the Fourier method.

Consider the irradiance in the interferogram in Fig. 18 along a line parallel to the x axis. This irradiance plotted versus the coordinate x is a perfectly sinusoidal function only if the wavefront is perfect, that is, if the fringes are straight, parallel, and equidistant. Otherwise, this function appears as a wave with a phase modulation. The phase-modulating function is the wavefront shape $W(x, y)$. If the tilt between the wavefronts is θ , the irradiance function is described by Eq. (15). If φ_0 is a multiple of 2π , this expression may be rewritten as

$$I(x, y) = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos(kx \sin \theta + kW) \quad (17)$$

Multiplying this phase-modulated function by a sinusoidal signal with the same frequency as the carrier $\sin(kx \sin \theta)$ a new signal S is obtained. Similarly, multiplying by a cosinusoidal signal

$\cos(kx \sin \theta)$ a new signal C is obtained. If all terms in the signals S and C with frequencies equal to or greater than the carrier frequency are removed with a low-pass filter, they become

$$S(x, y) = -\sqrt{I_1 I_2} \sin kW(x, y) \quad (18)$$

$$C(x, y) = \sqrt{I_1 I_2} \cos kW(x, y) \quad (19)$$

then, the wavefront $W(x, y)$ is given by

$$W(x, y) = -\frac{1}{k} \tan^{-1} \left[\frac{S(x, y)}{C(x, y)} \right] \quad (20)$$

which is our desired result.

13.7 PHASE-SHIFTING INTERFEROMETRY

All the methods just described are based on the analysis of a single static interferogram. Static fringe analysis is generally less precise than phase-shifting interferometry, by more than one order of magnitude. However, fringe analysis has the advantage that a single image of the fringes is needed. On the other hand, phase-shifting interferometry requires several images, acquired over a long time span during which the fringes must be stable. This is the main reason why phase-shifting interferometry has seldom been used for the testing of astronomical optics.

Phase-shifting interferometry^{25,26} is possible, thanks to modern tools like array detectors and microprocessors. Figure 19 shows a Twyman-Green interferometer adapted to perform phase-shifting interferometry. Most conventional interferometers, like the Fizeau and the Twyman-Green, have been used to do phase shifting. A good review about these techniques may be found in the review article by Creath.²⁷

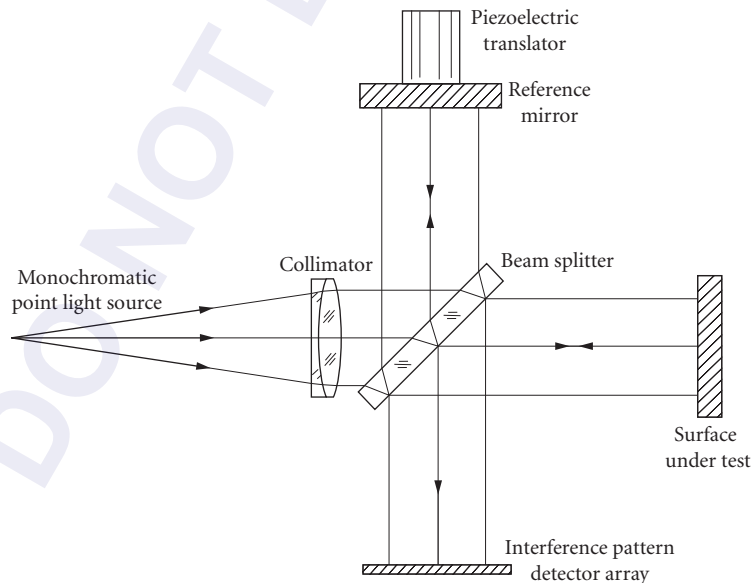


FIGURE 19 Twyman-Green interferogram adapted to do phase shifting.

In phase-shifting interferometers, the reference wavefront is moved along the direction of propagation, with respect to the wavefront under test, changing in this manner their phase difference. This phase shifting is made in steps or in a continuous manner. Of course, this relative displacement of one wavefront with respect to the other may only be achieved through a momentary or continuous change in the frequency of one of the beams, for example, by Doppler shift, moving one of the mirrors in the interferometer. In other words, this change in the phase is accomplished when the frequency of one of the beams is modified in order to form beats.

By measuring the irradiance changes for different values of the phase shifts, it is possible to determine the initial difference in phase between the wavefront under test and the reference wavefront, for that measured point over the wavefront. By obtaining this initial phase difference for many points over the wavefront, the complete wavefront shape is thus determined. If we consider any fixed point in the interferogram, the initial phase difference between the two wavefronts has to be changed in order to make several measurements.

One method that can be used to shift this phase is by moving the mirror for the reference beam along the light trajectory, as in Fig. 19. This can be done in many ways, for example, with a piezoelectric crystal or with a coil in a magnetic field. If the mirror moves with a speed V , the frequency of the reflected light is shifted by an amount equal to $\Delta\nu = 2V/\lambda$.

Another method to shift the phase is by inserting a plane parallel glass plate in the light beam (see Fig. 20). Then the plate is rotated about an axis perpendicular to the optical axis. The phase may also be shifted by means of the device shown in Fig. 21. The first quarter-wave retarding plate is stationary, with its slow axis at 45° with respect to the plane of polarization of the incident linearly polarized light. This plate also transforms the returning circularly polarized light back to linearly

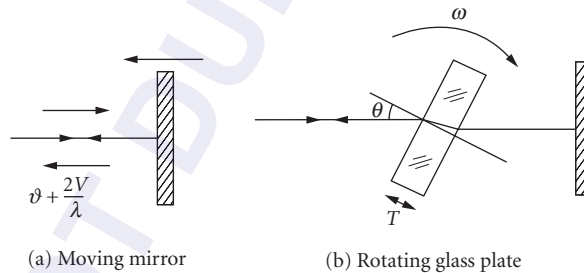


FIGURE 20 Obtaining the phase shift by means of a moving mirror or a rotating glass plate.

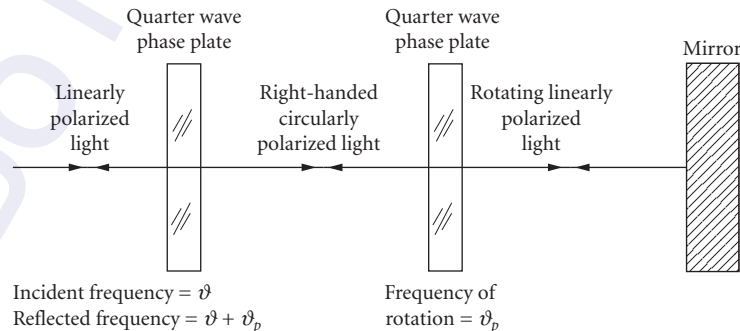


FIGURE 21 Obtaining the phase shift by means of phase plates and polarized light, with a double pass of the light beam.

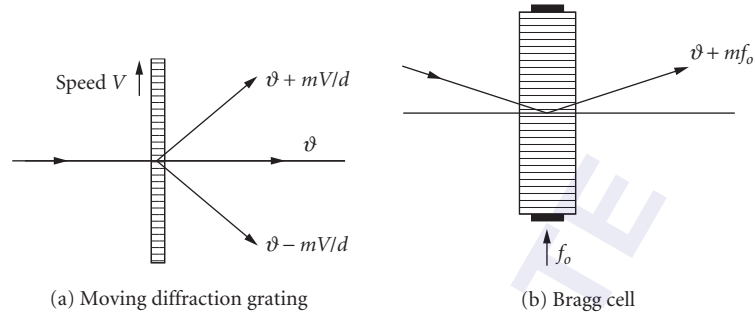


FIGURE 22 Obtaining the phase shift by means of diffraction: (a) with a diffraction grating and (b) with an acousto-optic Bragg cell.

polarized. The second phase retarder is also a quarter-wave plate, it is rotating, and the light goes through it twice; therefore, it is acting as a half-wave plate.

Still another manner to obtain the shift of the phase is by a diffraction grating moving perpendicularly to the light beam, as shown in Fig. 22a, or with an acousto-optic Bragg cell, as shown in Fig. 22b. The change in the frequency is equal to the frequency f of the ultrasonic wave times the order of diffraction m . Thus: $\Delta\nu = mf$.

The nonshifted relative phase of the two interfering wavefronts is found by measuring the irradiance with several predefined and known phase shifts. Let us assume that the irradiance of each of the two interfering light beams at the point x, y in the interference patterns are $I_1(x, y)$ and $I_2(x, y)$ and that their phase difference is $\phi(x, y)$. It was shown before, in Eq. (9), that the resultant irradiance $I(x, y)$ is a sinusoidal function describing the phase difference between the two waves. The basic problem is to determine the nonshifted phase difference between the two waves, with the highest possible precision. This may be done by any of several different procedures.

Phase-Stepping and Four Steps Algorithms

This method²⁷ consists of measuring the irradiance values for several known increments of the phase. There are several versions of this method, which will be described later. The measurement of the irradiance for any given phase takes some time, since there is a time response for the detector. Therefore, the phase has to be stationary during a short time in order to take the measurement. Between two consecutive measurements, the phase is changed by an increment α_i . For those values of the phase, the irradiance becomes

$$I(x, y) = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos(\phi + \alpha_i) \quad (21)$$

There are many different algorithms, with many different phase steps, as shown in Fig. 23. The minimum number of steps needed to reconstruct this sinusoidal function is three. As an example with four steps,

$$I_A = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos\phi \quad (22)$$

$$I_B = I_1 + I_2 - 2\sqrt{I_1 I_2} \sin\phi \quad (23)$$

$$I_C = I_1 + I_2 - 2\sqrt{I_1 I_2} \cos\phi \quad (24)$$

$$I_D = I_1 + I_2 + 2\sqrt{I_1 I_2} \sin\phi \quad (25)$$

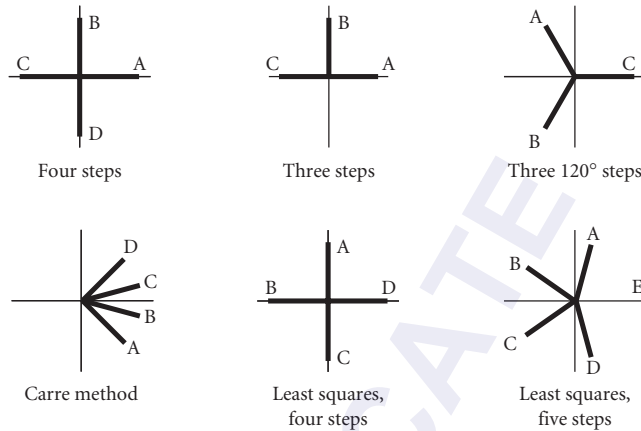


FIGURE 23 Six different ways to shift the phase using phase steps.

Integrating Bucket

In the integrating phase-shifting method the detector continuously measures the irradiance during a fixed time interval, without stopping the phase. Since the phase changes continuously, the average value of the irradiance during the measuring time interval is measured. Thus, the integrating phase-stepping method may be mathematically considered a particular case of the phase-stepping method if the detector has an infinitely short time response. Then, the measurement time interval is reduced to zero. If the measurement is taken, as in Fig. 24, from $\alpha_i + \Delta/2$ to $\alpha_i - \Delta/2$ with center at α_i , then

$$I = \frac{1}{\Delta} \int_{\alpha_i - \Delta/2}^{\alpha_i + \Delta/2} [I_1 + I_2 + 2\sqrt{I_1 I_2} \cos(\phi + \alpha_i)] d\alpha \quad (26)$$

$$I = I_1 + I_2 + 2\sqrt{I_1 I_2} \text{sinc}(\Delta/2) \cos(\phi + \alpha_i) \quad (27)$$

In general, in the phase-stepping as well as in the integrating phase-shifting methods, the irradiance is measured at several different values of the phase α_p , and then the phase is calculated.

Two Steps Plus One Method

As pointed out before, phase-shifting interferometry is not useful for testing systems with vibrations or turbulence because the three or four interferograms are taken at different times. An attempt

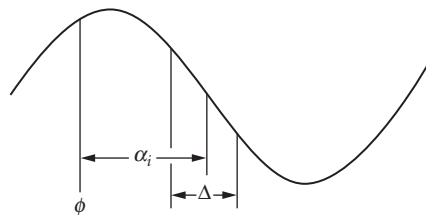


FIGURE 24 Averaged signal measurements with the integrating phase-shifting method.

to reduce this time is the so-called two steps plus one method, in which only two measurements separated by 90° are taken.²⁸ A third reading is taken any time later, of the sum of the irradiance of the beams, independently of their relative phase. This last reading may be taken using an integrating interval $\Delta = 2\pi$. Thus,

$$I_A = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos\phi \quad (28)$$

$$I_B = I_1 + I_2 + 2\sqrt{I_1 I_2} \sin\phi \quad (29)$$

$$I_C = I_1 + I_2 \quad (30)$$

Therefore,

$$\phi = \tan^{-1} \left\{ \frac{I_B - I_C}{I_A - I_C} \right\} \quad (31)$$

Other Phase-Stepping Algorithms

When calculating the phase with several different three or more phase steps, several sources of error affect the accuracy of the result, for example:

- (a) A line miscalibration of the phase shifter. Then the introduced phase steps will have a linear error, which can be interpreted as a change of the reference frequency f_r , making it different from the signed frequency as it should be.
- (b) A nonlinearity in the phase shifter or on the detector that introduces light order harmonics in the detected signed.

These phase errors can be reduced or minimized by properly selecting the number phase steps and their phase increments. Using Fourier theory as shown by Freischlad and Koliopolus,²⁹ a large number of different algorithms with different contributors and number of phase steps have been described in the literature. Depending on the kind of source error and the maximum number of phase steps desired, the proper algorithm can be selected.

Simultaneous Measurement

It has been said several times that the great disadvantage of phase-shifting interferometry is its great sensitivity to vibrations and atmospheric turbulence. To eliminate this problem, it has been proposed that the different interferograms corresponding to different phases be taken simultaneously.^{30,31} To obtain the phase-shifted interferogram, they have used polarization-based interferometers. The great disadvantage of these interferometers is their complexity. To measure the images these interferometers have to use several television cameras.

Heterodyne Interferometer

When the phase shift is made in a continuous manner rather than in steps, the frequency of the shifting beam is permanently modified, and a beating between the two interferometer beams is formed.³²

The phase of the modulated or beating wave may be determined in many ways. One way is by electronic analog techniques, for example, using leading-edge detectors. Another way is by detecting when the irradiance passes through zero, that is, through the axis of symmetry of the irradiance function.

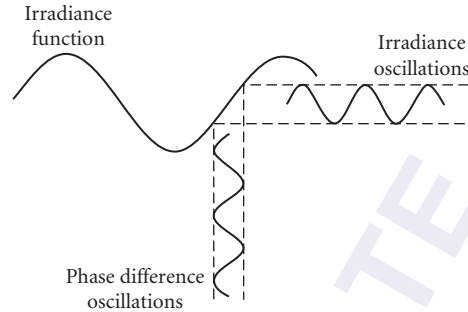


FIGURE 25 Phase-lock method to find the phase with a small sinusoidal modulation of the phase.

Phase Lock

The phase-lock method^{31–35} can be explained with the help of Fig. 25. Assume that an additional phase difference is added to the initial phase $\phi(x, y)$. The additional phase being added has two components: one of them with a fixed value and the other with a sinusoidal time shape. Both components can have any predetermined desired value. Thus, the resultant phase ϕ_r is given by

$$\phi_r = \phi(x, y) + \delta(x, y) + a \sin \omega t \quad (32)$$

then, the irradiance $I(x, y)$ would be given by

$$I = I_1 + I_2 + 2\sqrt{I_1 I_2} \cos[\phi + \delta + a \sin \omega t] \quad (33)$$

The amplitude of the phase oscillations $a \sin t$ is much smaller than π . We may now adjust the fixed phase to a value such that $\phi + \delta = \pi/2 + n\pi$. Then the value of $\cos(\phi + \delta)$ is zero. The curve is antisymmetric at this point; hence, only odd harmonics remain on the irradiance signal. This is done in practice by slowly changing the value of the phase δ , while maintaining the oscillation $a \sin \omega t$, until the maximum amplitude of the first harmonic, or fundamental frequency, is obtained. At this point, then, we have $\delta + \phi = \pi/2 + n\pi$, and since the value of δ is known, the value ϕ has been determined.

13.8 MEASURING ASPHERICAL WAVEFRONTS

The most common type of interferometer, with the exception of lateral or rotational shearing interferometers, produces interference patterns in which the fringes are straight, equidistant, and parallel, when the wavefront under test is perfect and spherical with the same radius of curvature as the reference wavefront.

If the surface under test does not have a perfect shape, the fringes will not be straight and their separations will be variable. The deformations of the wavefront may be determined by a mathematical examination of the shape of the fringes. By introducing a small spherical curvature on the reference wavefront (focus shift) or by changing its angle with respect to the wavefront under test (tilt), the number of fringes in the interferogram may be changed. This is done to reduce the number of fringes as much as possible, since the greater the number of fringes, the smaller the sensitivity of the test. However, for aspherical surfaces this number of fringes cannot be smaller than a certain minimum. The larger the asphericity is, the greater is this minimum number of fringes. Since the fringe

separations are not constant, in some places the fringes will be widely spaced, but in some others the fringes will be too close together.

The sensitivity of the test depends on the separation between the fringes, because an error of one wavelength in the wavefront distorts the fringe shape by an amount equal to the separation between the fringes. Thus, the sensitivity is directly proportional to the fringe separation. When the fringes are widely separated, the sampled points will be quite separated from each other, leaving many zones without any information. On the other hand, where the fringes are very close to each other, there is a high density of sampled data points, but the sensitivity is low.

Then, it is desirable that the spherical aberration of the wavefront under test is compensated in some way, so that the fringes appear straight, parallel, and equidistant, for a perfect wavefront. This is called a null test and may be accomplished by means of some special configurations. These special configurations may be used to conduct a null test of a conic surface. These are described in several books.¹ Almost all of these surfaces have rotational symmetry.

If no testing configuration can be found to get rid of the spherical aberration, additional optical components, called null compensators, have to be used. Many different types of compensators have been invented. The compensators may be refractive (lenses), reflective (mirrors), or diffractive (real or computer-generated holograms).

Refractive or Reflective Compensators

The simplest way to compensate the spherical aberration of a paraboloid or a hyperboloid tested at the center of curvature is a single convergent lens placed near the point of convergence of the rays, as Fig. 26 shows. This lens is called a Dall compensator. Unfortunately, the correction due to a single lens is not complete, so a system of two lenses must be used to obtain a better compensation. This system is called an Offner compensator and is shown in Fig. 27. The field lens L is used to image the surface under test on the plane of the compensating lens L . Mirrors may also be used to design a null compensator.

As the sad experience of the Hubble space telescope proves, the construction parameters in a lens compensator have to be very carefully measured and adjusted, otherwise an imperfect correction is

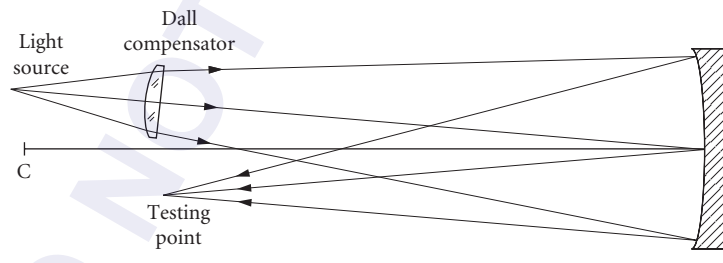


FIGURE 26 The Dall compensator.

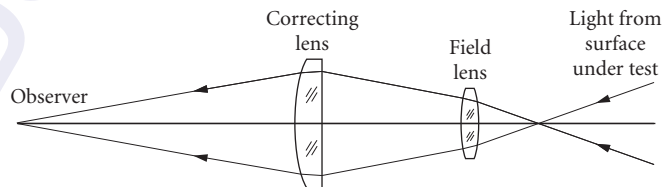


FIGURE 27 The Offner compensator. Only the reflected beam is shown.

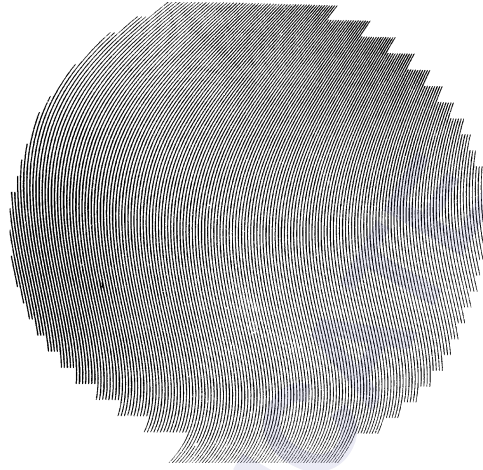


FIGURE 28 Computer-generated hologram for testing an aspherical wavefront. (From Wyant.³⁷)

obtained either by undercorrection or overcorrection. The distance from the compensator to the surface under test is one of those parameters to be carefully measured. A way around this problem would be to assume that the compensator detects smoothness imperfections but not the exact degree of asphericity. This degree of asphericity may then be measured with some independent measurement like the Hartmann test.

Holographic Compensators

Diffractive holographic elements also may be used to compensate the spherical aberration of the system and to obtain a null test. The hologram may be real, produced by photographing an interferometric pattern. This pattern has to be formed by superimposing on the screen a wavefront like the one we have to test and a perfectly flat or spherical wavefront. The only problem with this procedure is that a perfect wavefront with the same shape as the wavefront to be tested has first to be produced. This is not always easy.

A better approach is to simulate the holographic interference pattern in a computer,³⁶ as in Fig. 28. Then this image is transferred to a small photographic plate, with the desired dimensions. There are many experimental arrangements to compensate the aspherical wavefront aberration with a hologram. One of these is illustrated in Fig. 29.

Infrared Interferometry

Another simple approach to reduce the number of fringes in the interferogram is to use a long infrared wavelength. Light from a CO₂ laser has been used with this purpose. It can also be used when the surface is still quite rough.

Two-Wavelength Interferometry

In phase-shifting interferometry, each detector must have a phase difference smaller than π from the closest neighboring detector, in order to avoid 2π phase ambiguities and ensure phase continuity. In

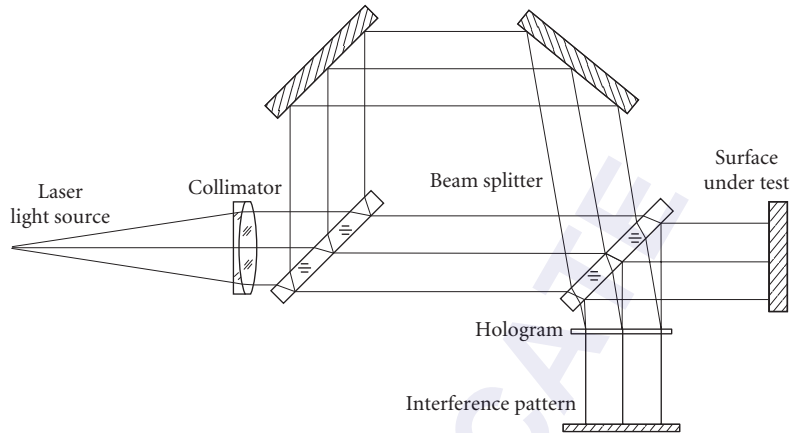


FIGURE 29 An optical arrangement for testing an aspherical wavefront with a computer-generated hologram.

other words, there should be at least two detector elements for each fringe. If the slope of the wavefront is very large, the fringes will be too close together and the number of detector elements would be extremely large.³⁷

A solution to this problem is to use two different wavelengths λ_1 and λ_2 simultaneously. The group wavelength or equivalent wavelength λ_{eq} is longer than any of the two components and is given by

$$\lambda_{\text{eq}} = \frac{\lambda_1 \lambda_2}{|\lambda_1 - \lambda_2|} \quad (34)$$

Under these conditions, the requirement in order to avoid phase uncertainties is that there should be at least two detectors for each fringe produced if the wavelength is λ_{eq} . The great advantage of this method is that we may test wavefronts with large asphericities, limited in asphericity by the group wavelength, and accuracy limited by the shortest wavelength of the two components.

Moiré Tests

An interferogram in which a large amount of tilt has been introduced is an ideal periodic structure to form moiré patterns. A moiré pattern represents the difference between two periodic structures. Thus, a moiré formed by two interferograms represents the difference between the two interferograms. There are several possibilities for the use in optical testing of this technique, as shown by Patorski.³⁸

Let us assume that the two interferograms are taken from the same optical system producing an aspherical wavefront, but with two different wavelengths λ_1 and λ_2 . The moiré obtained represents the interferogram that would be obtained with an equivalent wavelength λ_{eq} given by Eq. (31). If the tilt is of different magnitude in the two interferograms, the difference appears as a tilt in the moiré between them. Strong aspheric wavefronts may be tested with this method.

A second possibility is to produce the moiré between the ideal interferogram for an aspheric wavefront and the actual wavefront. Any differences between both would be easily detected.

Another possibility of application is for eliminating the wavefront imperfections in a low-quality interferometer. One interferogram is taken with the interferometer alone, without any optical piece under test. The second interferogram is taken with the optical component being tested. The moiré

represents the wavefront deformations due to the piece being tested, without the interferometer imperfections.

Sub-Nyquist Interferometry

It was pointed out before that in phase-shifting interferometry each detector must have a phase difference smaller than π from the closest neighboring detector, in order to avoid 2π phase ambiguities and to ensure phase continuity. In other words, there should be at least two detector elements for each fringe. This condition is known as the Nyquist condition.

Since there is a minimum practical distance between detectors, the maximum asphericity in a surface to be tested by phase-shifting interferometry is only a few wavelengths. This condition may be relaxed³⁹ if the wavefront and its slope are assumed to be continuous on the whole aperture. Then, optical surfaces with larger asphericities may be tested.

Wavefront Stitching

When an aspheric and a flat reference wavefronts interfere, the fringe spacing is minimum where the angle between the two wavefronts is larger. When the aspheric wavefront has rotational symmetry and there is no angle between them (no tilt) near the optical axis, the minimum fringe spacing occurs at the edge of the pupil. The maximum fringe spacing is at the center, where the two wavefronts are parallel to each other.

Most times the fringe pattern is imaged on a CCD detector with a rectangular array of small square pixels. According to the Nyquist sampling condition, the fringes can be detected only if the detector has more than two pixels per fringe spacing. With strongly aspheric wavefronts this condition can not be satisfied near the pupil edge.

If the flat reference wavefront is tilted or made slightly spherical, the zone of the interferogram with maximum fringe spacing can be located at any desired. Then the interferogram can be measured at a zone around this point with maximum fringe spacing where the Nyquist condition is satisfied. By moving this point the whole interferogram can thus be measured in small pieces. Then, all pieces should be joined together in a process called wavefront stitching,⁴⁰ as illustrated in Fig. 30.

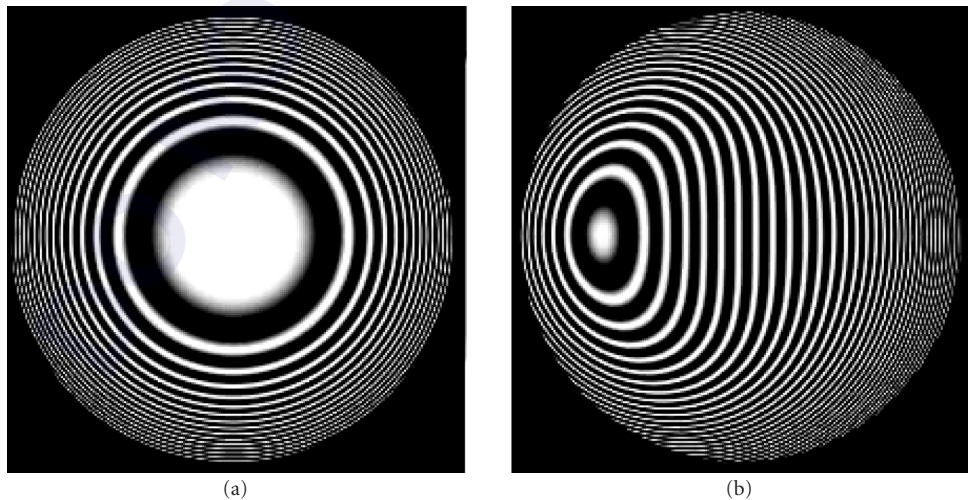


FIGURE 30 Two different interferograms of the same wavefront with different tilts to do wavefront stitching.

13.9 REFERENCES

1. D. Malacara, *Optical Shop Testing*, 3d ed., John Wiley and Sons, New York, 2007.
2. L. M. Foucault, "Description des Procédes Employes pour Reconnaître la Configuration des Surfaces Optiques," *C.R. Acad. Sci. Paris* **47**:958 (1852); reprinted in Armand Colin, *Classiques de al Science*, vol. II.
3. J. Ojeda-Castañeda, "Foucault, Wire and Phase Modulation Tests," in D. Malacara (ed.), *Optical Shop Testing*, 3d ed., John Wiley and Sons, New York, 2007.
4. V. Ronchi, "Le Franque di Combinazione Nello Studio Delle Superficie e Dei Sistemi Ottici," *Ri. Ottica mecc. Precis.* **2**:9 (1923).
5. J. Hartmann, "Bemerkungen über den Bann und die Justirung von Spektrographen," *Zt. Instrumentenk.* **20**:47 (1900).
6. B. C. Platt and R. V. Shack, "Lenticular Hartmann Screen," *Opt. Sci. Newsl.* **5**: 15–16 (1971).
7. D. Malacara, *Interferogram Analysis for Optical Testing*, 2d ed., CRC Press, Taylor and Francis Group, Boca Raton, FL, 2005.
8. R. N. Smart and J. Strong, "Point Diffraction Interferometer," (Abstract only) *J. Opt. Soc. Am.* **62**:737 (1972).
9. C. Roychoudhuri, "Multiple-Beam Interferometers," in D. Malacara (ed.), *Optical Shop Testing*, 3d ed., John Wiley and Sons, New York, 2007.
10. P. Hariharan and D. Sen, "The Separation of Symmetrical and Asymmetrical Wave-Front Aberrations in the Twyman Interferometer," *Proc. Phys. Soc.* **77**:328 (1961).
11. P. Hariharan, "Multiple-Pass Interferometers," in D. Malacara (ed.), *Optical Shop Testing*, 2d ed., John Wiley and Sons, New York, 1991.
12. F. Zernike, "Diffraction Theory of Knife Edge Test and Its Improved Form: The Phase Contrast," *Mon. Not. R. Astron. Soc.* **94**:371 (1934).
13. D. Rozenzweig and B. Alte, "A Facility for the Analysis of Interferograms," in A. H. Guenther, D. H. Liedbergh (eds), *Optical Interferograms—Reduction and Interpretation*, ASTM Symposium, *Am. Soc. for Test and Mat. Tech.* Publ. 666, West Conshohocken, PA, 1978.
14. K. H. Womack, J. A. Jonas, C. L. Koliopoulos, K. L. Underwood, J. C. Wyant, J. S. Loomis, and C. R. Hayslett, "Microprocessor-Based Instrument for Analysis of Video Interferograms," *Proc. SPIE* **192**:134 (1979).
15. G. T. Reid, "Automatic Fringe Pattern Analysis: A Review," *Opt. and Lasers in Eng.* **7**:37 (1986).
16. G. T. Reid, "Image Processing Techniques for Fringe Pattern Analysis," *Proc. SPIE* **954**:468 (1988).
17. D. Malacara, J. M. Carpio-Valadéz, and J. J. Sánchez-Mondragón, "Wavefront Fitting with Discrete Orthogonal Polynomials in a Unit Radius Circle," *Opt. Eng.* **29**:672 (1990).
18. C. R. Hayslett and W. Swantner, "Wave Front Derivation from Interferograms by Three Computer Programs," *Appl. Opt.* **19**:3401 (1980).
19. F. Becker, G. E. A. Maier, and H. Wegner, "Automatic Evaluation of Interferograms," *Proc. SPIE* **359**:386 (1982).
20. F. Zernike, "Begünsttheorie des Schneidener-Fahrens und Seiner Verbasserten Form, der Phasenkontrastmethode," *Physica* **1**:689 (1934).
21. K. H. Womack, "Frequency Domain Description of Interferogram Analysis," *Opt. Eng.* **23**:396 (1984).
22. W. W. Macy, Jr., "Two Dimensional Fringe Pattern Analysis," *Appl. Opt.* **22**:3898 (1983).
23. M. Takeda, H. Ina, and S. Kobayashi, "Fourier Transform Method of Fringe-Pattern Analysis for Computer-Based Topography and Interferometry," *J. Opt. Soc. Am.* **72**:156 (1982).
24. C. Roddier and F. Roddier, "Interferogram Analysis Using Fourier Transform Techniques," *Appl. Opt.* **26**:1668 (1987).
25. J. H. Bruning, D. J. Herriott, J. E. Gallagher, D. P. Rosenfeld, A. D. White, and D. J. Brangaccio, "Digital Wavefront Measurement Interferometer," *Appl. Opt.* **13**:2693 (1974).
26. J. Greivenkamp and J. H. Bruning, "Phase Shifting Interferometers," in D. Malacara (ed.), *Optical Shop Testing*, 2d ed., John Wiley and Sons, New York, 1991.
27. K. Creath, "Phase-Measurement Interferometry Techniques," in E. Wolf (ed.), *Progress in Optics*, vol. XXVI, Elsevier Science Publishers, Amsterdam, 1988.

28. P. L. Wizinowich, "Systems for Phase Shifting Interferometry in the Presence of Vibration: A New Algorithm and System," *Appl. Opt.* **29**:3271–3279 (1990).
29. K. Freischlad and C. L. Koliopoulos, "Fourier Description of Digital Phase Measuring Interferometry," *J. Opt. Soc. Am. A.* **7**:542–551 (1990).
30. N. Bareket, "Three-Channel Phase Detector for Pulsed Wavefront Sensing," *Proc. SPIE* **551**:12 (1985).
31. C. L. Koliopoulos, "Simultaneous Phase Shift Interferometer," *Proc. SPIE* **1531**:119–127 (1991).
32. N. A. Massie, "Digital Heterodyne Interferometry," *Proc. SPIE* **816**:40 (1987).
33. G. W. Johnson, D. C. Leiner, and D. T. Moore, "Phase Locked Interferometry," *Proc. SPIE* **126**:152 (1977).
34. G. W. Johnson, D. C. Leiner, and D. T. Moore, "Phase Locked Interferometry," *Opt. Eng.* **18**:46 (1979).
35. D. T. Moore, "Phase-Locked Moire Fringe Analysis for Automated Contouring of Diffuse Surfaces," *Appl. Opt.* **18**:91 (1979).
36. J. C. Wyant, "Holographic and Moire Techniques," in D. Malacara (ed.), *Optical Shop Testing*, John Wiley and Sons, New York, 1978.
37. J. C. Wyant, B. F. Oreb, and P. Hariharan, "Testing Aspherics Using Two-Wavelength Holography: Use of Digital Electronic Techniques," *Appl. Opt.* **23**:4020 (1984).
38. K. Patorski, "Moire' Methods in Interferometry," *Opt. and Lasers in Eng.* **8**:147 (1988).
39. J. E. Greivenkamp, "Sub-Nyquist Interferometry," *Appl. Opt.* **26**:5245 (1987).
40. J. Liesener and H. Tiziani, "Interferometer with Dynamic Reference," *Proc. SPIE* **5252**:264–271 (2004).

This page intentionally left blank.

DO NOT DUPLICATE

USE OF COMPUTER-GENERATED HOLOGRAMS IN OPTICAL TESTING

Katherine Creath

*Optineering
Tucson, Arizona, and
College of Optical Sciences
University of Arizona
Tucson, Arizona*

James C. Wyant

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

14.1 GLOSSARY

CGH	computer-generated hologram
M	linear, lateral magnification
N	diffracted order number
n	integers
P	number of distortion-free resolution points
r	radius
S	maximum wavefront slope (waves/radius)
$x, \Delta x$	distance
$\Delta\theta$	rotational angle error
$\Delta\phi$	wavefront phase error
θ	rotational angle
λ	wavelength
$\phi(\)$	wavefront phase described by hologram

14.2 INTRODUCTION

Holography is extremely useful for the testing of optical components and systems. If a master optical component or optical system is available, a hologram can be made of the wavefront produced by the component or system and this stored wavefront can be used to perform null tests of similar

optical systems. If a master optical system is not available for making a hologram, a synthetic or a computer-generated hologram (CGH) can be made to provide the reference wavefront.¹⁻⁸ When an aspheric optical element with a large departure from a sphere is tested, a CGH can be combined with null optics to perform a null test.

There are several ways of thinking about CGHs. For the testing of aspheric surfaces, it is easiest to think of a CGH as a binary representation of the ideal interferogram that would be produced by interfering the reference wavefront with the wavefront produced by a perfect sphere. In the making of the CGH the entire interferometer should be ray traced to determine the so-called perfect aspheric wavefront at the hologram plane. This ray trace is essential because the aspheric wavefront will change as it propagates, and the interferometer components may change the shape of the perfect aspheric wavefront.

Figure 1 shows an example of a CGH. Since the amplitude of the aspheric wavefront is constant across the wavefront, best results are obtained if the lines making up the hologram have approximately one-half the spacing of the lines (i.e., fringe spacing) at the location of the lines. Thus, the line width will vary across the hologram. The major difference between the binary synthetic hologram

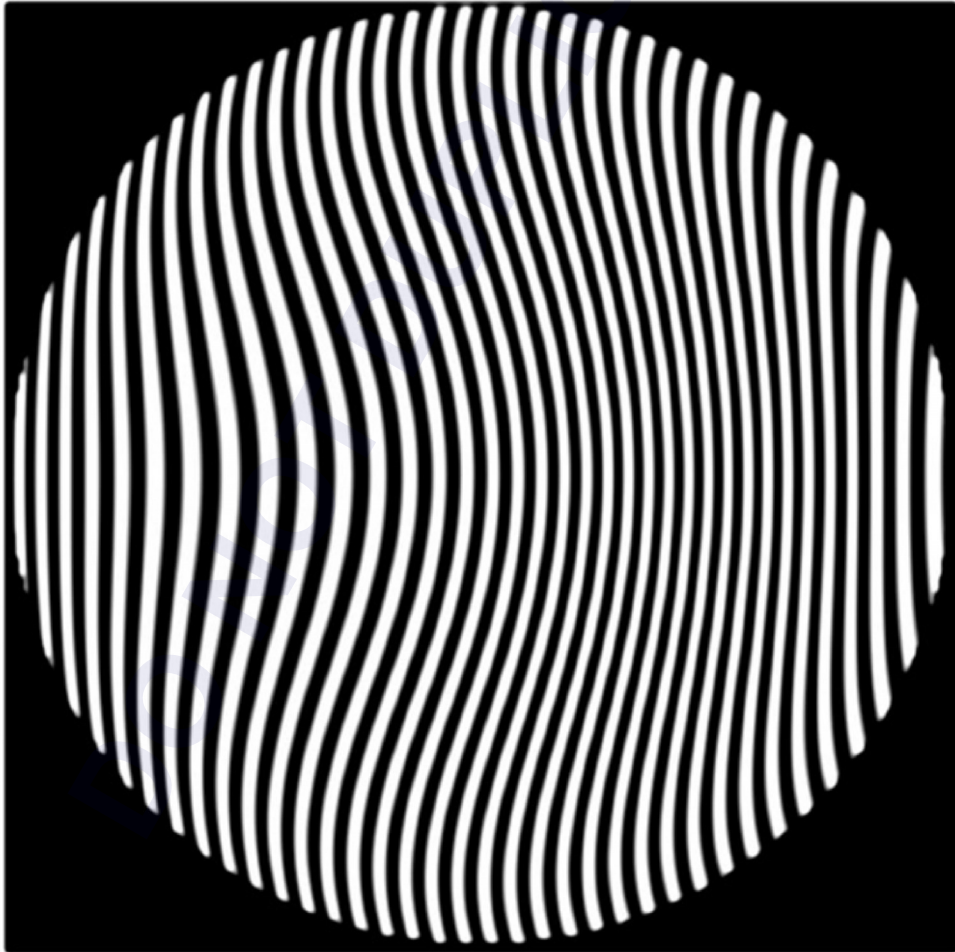


FIGURE 1 Sample computer-generated hologram (CGH).

and the real grayscale hologram that would be produced by interfering a reference wavefront and the aspheric wavefront is that additional diffraction orders are produced. These additional diffraction orders can be eliminated by spatial filtering.

14.3 PLOTTING CGHs

The largest problem in making CGHs is the plotting. The accuracy of the plot determines the accuracy of the wavefront. It is easier to see the plotting accuracy by comparing a binary synthetic hologram with an interferogram. In an interferogram, a wavefront error of $1/n$ waves causes a fringe to deviate from the ideal position by $1/n$ the fringe spacing. The same is true for CGHs. A plotting error of $1/n$ the fringe spacing will cause an error in the produced aspheric wavefront of $1/n$ wave. As an example, assume the error in drawing a line is $0.1\ \mu\text{m}$ and the fringe spacing is $20\ \mu\text{m}$, then the wavefront produced by the CGH will have an error in units of wave of $0.1/20$, or $1/200$ wave.

To minimize wavefront error due to the plotter, the fringe spacing in the CGH should be as large as possible. The minimum fringe spacing is set by the slope difference between the aspheric wavefront and the reference wavefront used in the making of the synthetic hologram. While it is not mandatory, the interferogram is cleaner if the slope difference is large enough to separate the diffraction orders so spatial filtering can be used to select out only the first order. Figure 2 shows a photograph of the diffracted orders. As shown in Fig. 2, to ensure no overlapping of the first and second orders in the Fourier plane, the tilt angle of the reference beam needs to be greater than three times the maximum slope of the aberrated wave.⁹ This means that, in general, the maximum slope difference between the reference and test beams is four times the maximum slope of the test beam. Thus, the error produced by plotter distortion is proportional to the slope of the aspheric wavefront being produced.

Many plotters have been used to plot holograms, but the best holograms are now made using either laser-beam recorders or more commonly electron-beam (e-beam) recorders of the type used for producing masks in the semiconductor industry.¹⁰ The e-beam recorders write onto photoresist

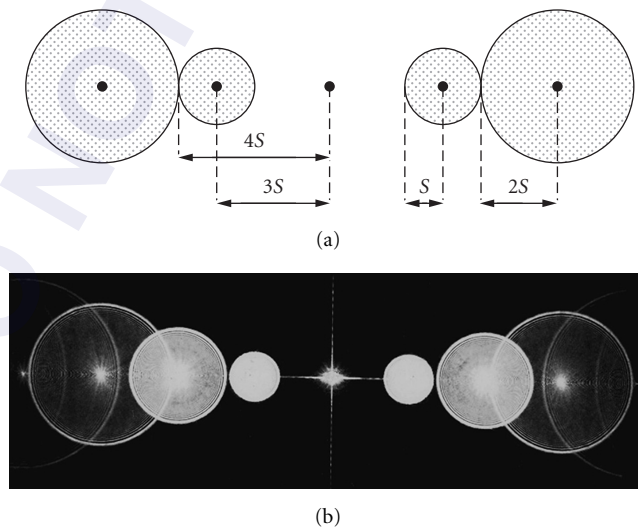


FIGURE 2 Diffracted orders in Fourier plane of CGH: (a) drawing and (b) photograph.

deposited on an optical-quality glass plate and currently produce the highest-quality CGHs. Typical e-beam recorders will write areas larger than $100 \text{ mm} \times 100 \text{ mm}$ with positional accuracies of less than 100 nm .¹¹

If needed, plotter distortion can be measured and calibrated out in the making of the hologram.^{12,13} The easiest way of determining plotter distortion is to draw straight lines and then treat this plot as a diffraction grating. If the computer-generated grating is illuminated with two plane waves, and the $-N$ order of beam 1 is interfaced with the $+N$ order of beam 2, the resulting interferogram gives us the plotter distortion. If the lines drawn by the plotter are spaced a distance Δx , a fringe error in the interferogram corresponds to a distortion error of $\Delta x/2N$ in the plot.

14.4 INTERFEROMETERS USING COMPUTER-GENERATED HOLOGRAMS

Many different experimental setups can be used for the holographic testing of optical elements. Figure 3 shows one common setup. The setup must be ray traced so the aberration in the hologram plane is known. While in theory there are many locations where the hologram can be placed, it is convenient to place the hologram in a plane conjugate to the asphere under test so the intensity across the image of the asphere is uniform. The longitudinal positional sensitivity for the hologram is reduced if the hologram is made in a region where the beams are collimated. Another advantage of this setup is that both the test and the reference beams pass through the hologram so errors resulting from hologram substrate thickness variations are eliminated without requiring the hologram be made on a good optical flat.

Another common setup for using a CGH to test aspheres is shown in Fig. 4.¹⁴ The largest advantage of this setup is that it works well with commercial Fizeau interferometers. The only addition to the commercial interferometer is a mount to hold the CGH between the transmission sphere and the optics under test. Since the light is diffracted by the CGH twice, the CGH must be a phase

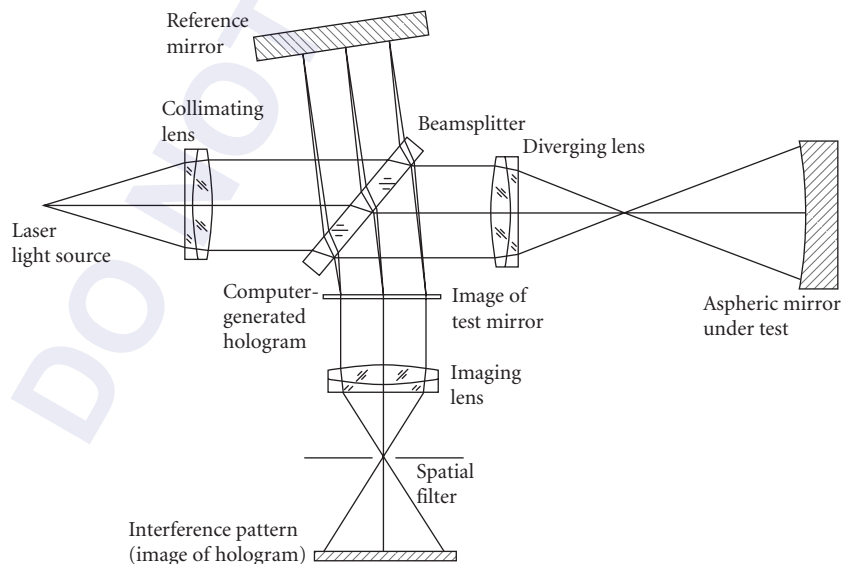


FIGURE 3 Interferometer setup using CGHs to test aspheric wavefronts.

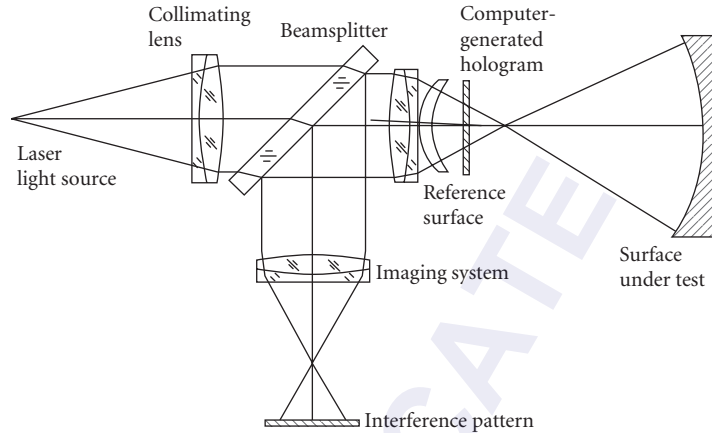


FIGURE 4 Use of CGH with Fizeau interferometer.

hologram so the diffraction efficiency is good, and since only the test beam is transmitted through the CGH, the substrate must either be high quality or thickness variations in the substrate must be measured and subtracted from the test results.

Figure 5 shows a setup for testing convex surfaces. In this case an on-axis CGH is used and the CGH is made on the concave reference surface.¹⁵ The light waves are perpendicular to the concave reference surface and then after diffraction they become perpendicular to the surface under test. The CGH pattern may be drawn exposing photoresist, ablating a metallic coating, or by creating a thin oxidation layer by heating a metal coating with a focused laser beam.¹⁶

CGHs can also be combined with partial null optics to test much more complicated aspherics than can be practically tested with either a CGH or null optics. This combination gives the real power of computer-generated holograms.¹⁷

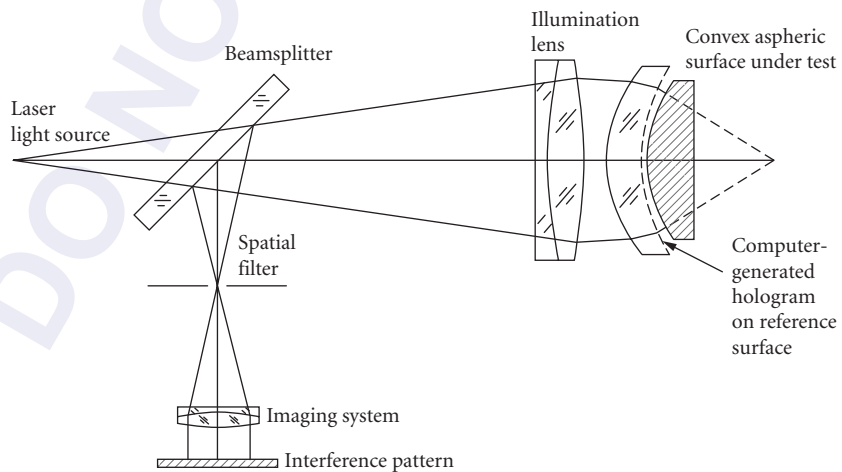


FIGURE 5 Using CGH to test convex surface.

14.5 ACCURACY LIMITATIONS

The largest source of error is the error due to plotter distortion as discussed previously. The other large sources of error are improper positioning of the hologram in the interferometer, and incorrect hologram size.

Any translation or rotation of the hologram produces error.² If the hologram is made conjugate to the exit pupil of the master optical system, the exit pupil of the system under test must coincide with the hologram. If the test wavefront in the hologram plane is described by the function $\phi(x, y)$, a displacement of the hologram a distance Δx in the x direction produces an error

$$\Delta\phi(x, y) \approx \frac{\partial\phi(x, y)}{\partial x} \Delta x \quad (1)$$

where $\partial\phi/\partial x$ is the slope of the wavefront in the x direction. Similarly, for a wavefront described by $\phi(r, \theta)$, the rotational error $\Delta\theta$ is given by

$$\Delta\phi(r, \theta) \approx \frac{\partial\phi(r, \theta)}{\partial\theta} \Delta\theta \quad (2)$$

Another source of error is incorrect hologram size. If the aberrated test wavefront in the plane of the hologram is given by $\phi(r, \theta)$, a hologram of incorrect size will be given by $\phi(r/M, \theta)$, where M is a magnification factor. The error due to incorrect hologram size will be given by the difference $\phi(r/M, \theta) - \phi(r, \theta)$, and can be written in terms of a Taylor expansion as

$$\begin{aligned} \phi\left(\frac{r}{M}, \theta\right) - \phi(r, \theta) &= \phi\left[r + \left(\frac{1}{M} - 1\right)r, \theta\right] - \phi(r, \theta) \\ &= \left[\frac{\partial\phi(r, \theta)}{\partial r}\right] \left(\frac{1}{M} - 1\right)r + \dots \end{aligned} \quad (3)$$

where terms higher than first order can be neglected if M is sufficiently close to 1, and a small region is examined. Note that this error is similar to a radial shear. When the CGH is plotted, alignment aids, which can help in obtaining the proper hologram size, can be drawn on the hologram plot. Figure 6 shows a CGH where the hologram is made in the center of the substrate and alignment aids are placed on the outer portion of the CGH.¹¹ Not only can the alignment aids help in putting the CGH in the proper position, but they can be used to help position the optics being tested. Figure 7

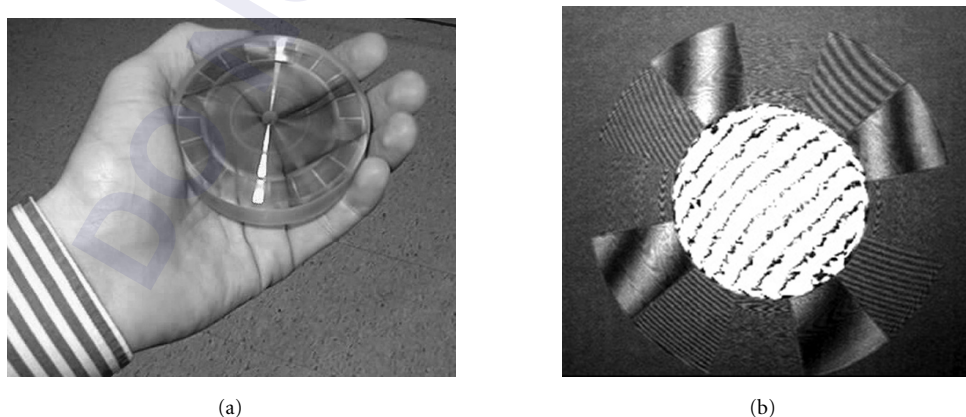


FIGURE 6 Use of CGH for alignment: (a) note structure in CGH and (b) interferogram produced with this CGH.

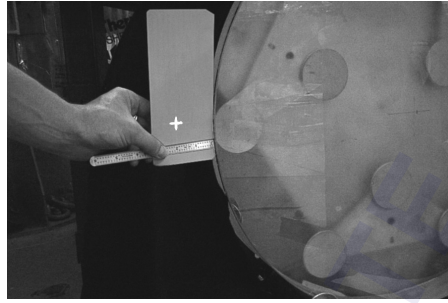


FIGURE 7 Use of crosshair produced by CGH to aid in the alignment of an off-axis parabola mirror.

shows a crosshair produced by a CGH that aids in the alignment of an off-axis parabolic mirror. The same CGH used to produce the crosshair produces the aspheric wavefront required to provide a null test of the off-axis parabola.¹¹

14.6 EXPERIMENTAL RESULTS

Figure 8 shows the results of using the setup shown in Fig. 3 to measure a 10-cm-diameter F/2 parabola using a CGH generated with an e-beam recorder. The fringes obtained in a Twyman-Green interferometer using a helium-neon source without the CGH present are shown in Fig. 8*a*. After the CGH is placed in the interferometer, a much less complicated interferogram is obtained as shown in Fig. 8*b*. The CGH corrects for about 80 fringes of spherical aberration, and makes the test much easier to perform.

To illustrate the potential of a combined CGH/null-lens test, results for a CGH/null-lens test of the primary mirror of an eccentric Cassegrain system with a departure of approximately 455 waves (at 514.5 nm) and a maximum slope of approximately 1500 waves per radius are shown.¹⁷ The mirror was a 69-cm diameter off-axis segment whose center lies 81 cm from the axis of symmetry of the parent aspheric surface. The null optics was a Maksutov sphere (as illustrated in Fig. 9), which reduces the departure and slope of the aspheric wavefront from 910 to 45 waves, and 300 to 70 waves per radius, respectively. A hologram was then used to remove the remaining asphericity.

Figure 10*a* shows interferograms of the mirror under test obtained using the CGH Maksutov test. Figure 10*b* shows the results when the same test was performed using a rather expensive refractive

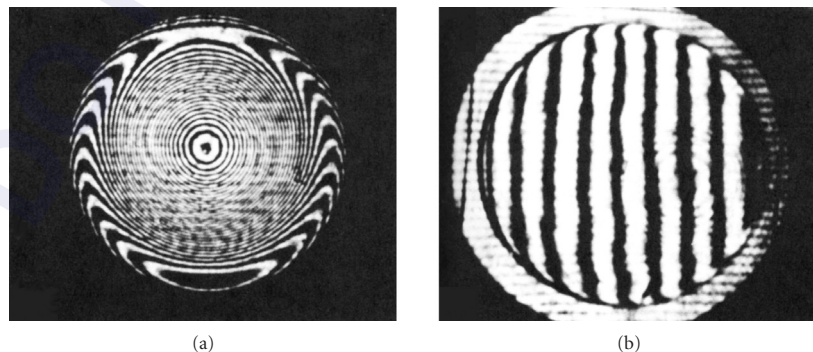


FIGURE 8 Results obtained testing a 10-cm-diameter F/2 parabola: (a) without using CGH and (b) using CGH made using an e-beam recorder.

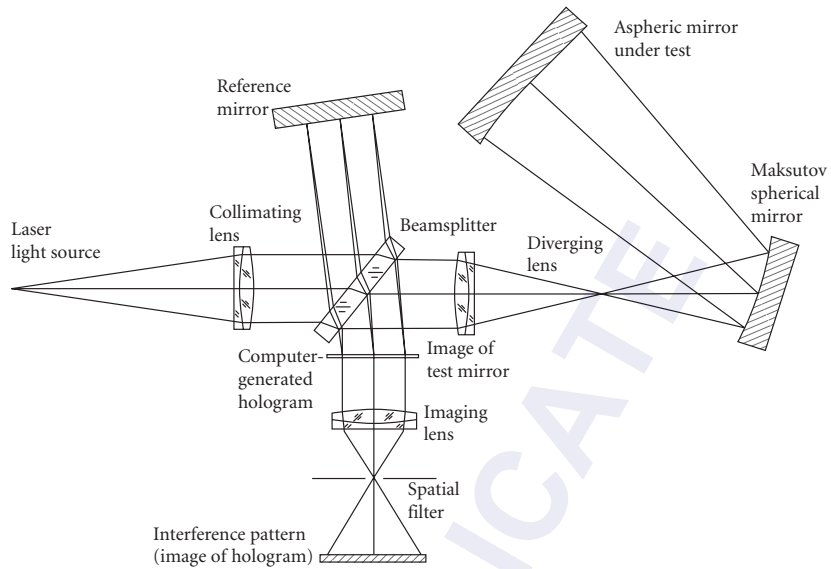


FIGURE 9 Setup to test the primary mirror of a Cassegrain telescope using a Maksutov sphere as a partial null and a CGH.

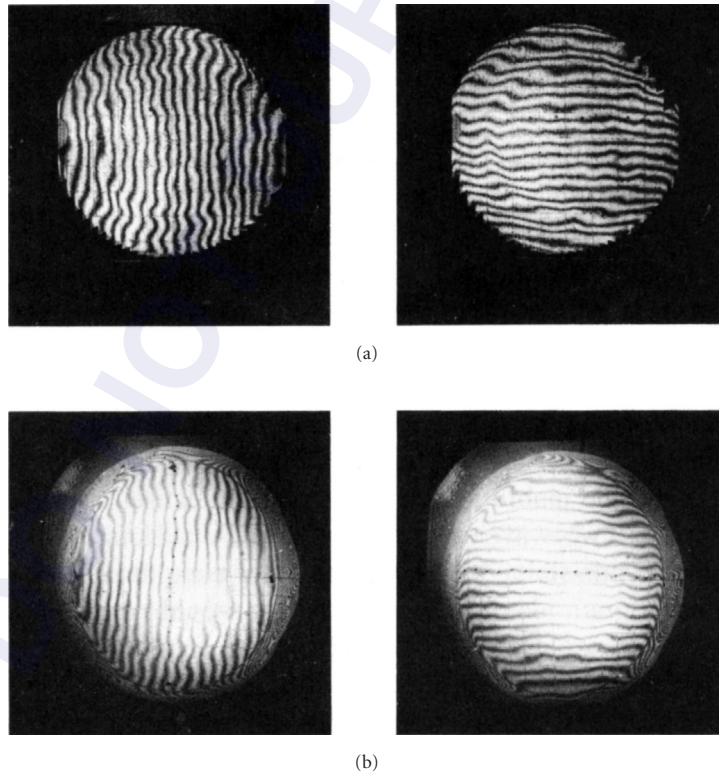


FIGURE 10 Results obtained using Fig. 9: (a) CGH-Maksutov test ($\lambda = 514.5 \text{ nm}$) and (b) using null lens ($\lambda = 632.8 \text{ nm}$).

null lens. When allowance is made for the fact that the interferogram obtained with the null lens has much more distortion than the CGH Maksutov interferogram, and for the difference in sensitivity ($\lambda = 632.8$ nm for the null-lens test and 514.5 nm for the CGH-Maksutov test), the results for the two tests are seen to be very similar. The “hills” and “valleys” on the mirror surface appear the same for both tests, as expected. The peak-to-valley surface error measured using the null lens was 0.46 waves (632.8 nm), while for the CGH-Maksutov test it was 0.39 waves (514.5 nm). The rms surface error measured was 0.06 waves (632.8 nm) for the null lens, while the CGH Maksutov test gave 0.07 wave (514.5 nm). These results certainly demonstrate that expensive null optics can be replaced by a combination of relatively inexpensive null optics and a CGH.

14.7 DISCUSSION

The difficult problem of testing aspheric surfaces, which are becoming increasingly popular in optical design, is made easier by the use of CGHs. The technology has reached the point that commercial interferometers using computer-generated holograms are now available. The main problem with testing aspheric optical elements is reducing the aberration sufficiently to ensure that light gets back through the interferometer. Combinations of simple null optics with a CGH to perform a test enable the measurement of a wide variety of optical surfaces. The making and use of a CGH are analogous to using an interferometer setup that yields a large number of interference fringes, and measuring the interferogram at a large number of data points. Difficulties involved in recording and analyzing a high-density interferogram and making a CGH are very similar. In both cases, a large number of data points are necessary, and the interferometer must be ray traced so that the aberrations due to the interferometer are well known. The advantage of the CGH technique is that once the CGH is made, it can be used for testing a single piece of optics many times or for testing several identical optical components. Additional alignment aids can be placed on the CGH to aid in the alignment of the CGH and the optics under test.

14.8 REFERENCES

1. A. J. MacGovern and J. C. Wyant, “Computer Generated Holograms for Testing Optical Elements,” *Appl. Opt.* **10**(3):619–624 (1971).
2. J. C. Wyant and V. P. Bennett, “Using Computer Generated Holograms to Test Aspheric Wavefronts,” *Appl. Opt.* **11**(12):2833–2839 (1972).
3. A. F. Fercher and M. Kriese, “Binäre Synthetische Hologramme zur Prüfung Asphärischer Optischer Elemente,” *Optik*. **35**(2):168–179 (1972).
4. Y. Ichioka and A. W. Lohmann, “Interferometric Testing of Large Optical Components with Circular Computer Holograms,” *Appl. Opt.* **11**(11):2597–2602 (1972).
5. J. Schwider and R. Burrow, “The Testing of Aspherics by Means of Rotational-Symmetric Synthetic Holograms,” *Optica Applicata* **6**:83 (1976).
6. T. Yatagai and H. Saito, “Interferometric Testing with Computer-Generated Holograms: Aberration Balancing Method and Error Analysis,” *Appl. Opt.* **17**(4):558–565 (1978).
7. J. Schwider, R. Burrow, and J. Grzanna, “CGH—Testing of Rotational Symmetric Aspheric in Compensated Interferometers,” *Optica Applicata* **9**:39 (1979).
8. C.S. Pruss, S. Reichelt, H.J. Tiziani, and W. Osten, “Computer-Generated Holograms in Interferometric Testing,” *Opt. Eng.* **43**:2534–2540 (2004).
9. J. W. Goodman, *Introduction to Fourier Optics*, 3d ed. Roberts & Company: Greenwood Village, Colorado, 2004.
10. Y. C. Chang and J. H. Burge, “Error Analysis for CGH Optical Testing,” *Proc. SPIE* **3872**:358–366 (1999).
11. J. H. Burge, R. Zehnder, and Chunyu Zhao, “Optical Alignment with Computer Generated Holograms,” *Proc. SPIE* **6676**:66760C (2007).

12. J. C. Wyant, P. K. O'Neill, and A. J. MacGovern, "Interferometric Method of Measuring Plotter Distortion," *Appl. Opt.* **13**(7):1549–1551 (1974).
13. A. F. Fercher, "Computer Generated Holograms for Testing Optical Elements: Error Analysis and Error Compensation," *Optica Applicata* **23**(5):347–365 (1976).
14. H. J. Tiziani, J. S. Reichlet, C. Pruss, M. Rocktachel, and U. Hofbauer, "Testing of Aspheric Surfaces," *Proc. SPIE*, **4440**:109–119 (2001).
15. J. H. Burge and D. S. Anderson, "Full-Aperture Interferometric Test of Convex Secondary Mirrors Using Holographic Test Plates," *Proc. SPIE* **2199**:181–192 (1994).
16. J. H. Burge, M. J. Fehniger, and G. C. Cole, "Demonstration of Accuracy and Flexibility of Using CGH Test Plates for Measuring Aspheric Surfaces," *Proc. SPIE* **3134**:379–389 (1997).
17. J. C. Wyant and P. K. O'Neill, "Computer Generated Hologram; Null Lens Test of Aspheric Wavefronts," *Appl. Opt.* **13**(12):2762–2765 (1974).

PART

4

SOURCES

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

ARTIFICIAL SOURCES

Anthony LaRocca*

*General Dynamics
Advanced Information Systems
Ypsilanti, Michigan*

15.1 GLOSSARY

λ	wavelength
$d\lambda$	differential wavelength
M_λ	spectral radiant exitance
T	absolute temperature
c_1	first radiation constant
c_2	second radiation constant
h	Planck's constant
c	velocity of light
k	Boltzmann constant
λ_{\max}	wavelength at peak of radiant exitance
K	factor
R	radius of the interior surface of the cavity
r	radius of the aperture
S	interior surface area of the cavity
s	aperture area
ϵ	emissivity
ϵ_0	uncorrected emissivity

15.2 INTRODUCTION

Whereas most of the sources described in this chapter can be used for any purpose for which one can justify their use, the emphasis is on the production of the appropriate radiation for the calibration of measurement instrumentation. This implies that the basis for their use is supported by their

*Retired

traceabilities to calibrated standards of radiation from an internationally known and respected standards laboratory such as the National Institute of Standards and Technology (NIST) in the United States or, say, the National Physical Laboratory (NPL) in the United Kingdom. Because calibration implies a high degree of accuracy, the chapter initially contains a short exposition on the so-called Planckian, or blackbody radiation standard, and the equation which describes the Planck radiation.

This chapter deals with artificial sources of radiation as subdivided into two classes: laboratory and field sources. Much of the information on commercial sources is taken from a chapter previously written by the author.¹ Where it was feasible, similar information here has been updated. When vendors failed to comply to requests for information, the older data were retained to maintain completeness, but the reader should be aware that some sources cited here may no longer exist, or perhaps may not exist in the specification presented. Normally, laboratory sources are used in some standard capacity and field sources are used as targets. Both varieties appear to be limitless. Only laboratory sources are covered here.

The sources in this chapter were chosen arbitrarily, often depending on manufacturer response to requests for information. The purpose of this chapter is to consolidate much of this information to assist the optical-systems designer in making reasonable choices. To attempt to include the hundreds of types of lasers, however, and the thousands of varieties, would be useless for several reasons, but particularly because they change often.

Complementing the material in the following chapter on Lasers, a fairly comprehensive source of information on lasers can be found in the *CRC Handbook of Laser Science and Technology*, Supplement 1: published by the CRC Press, Inc., 2000 Corporate Blvd., N.W., Boca Raton, Florida, 33431. Of course, the literature is laden with material on lasers, including the chapter on Lasers in this *Handbook*, and the reader would be wise to consult the Internet from which compilations such as the one cited above or a host of others can be obtained from companies like Amazon.com, or, better yet, by accessing a literature-rich source such as Google.

Regarding the selection of a source, Worthing² suggests that one ask the following questions:

1. Does it supply energy at such a rate or in such an amount as to make measurements possible?
2. Does it yield an irradiation that is generally constant or that may be varied with time as desired?
3. Is it reproducible?
4. Does it yield irradiations of the desired magnitudes over areas of the desired extent?
5. Has it the desired spectral distribution?
6. Has it the necessary operating life?
7. Has it sufficient ruggedness for the proposed problem?
8. Is it sufficiently easy to obtain and replace, or is its purchase price or its construction cost reasonable?

15.3 RADIATION LAW

All of the radiation sources described in this chapter span the region of the electromagnetic spectrum mainly from the visible region (starting from about 400) through the infrared (around 400 μm and beyond). Given that they have a demonstrable temperature, they relate in their own peculiar ways, depending on material properties, to the radiation called “blackbody” radiation, which is described by the Planck radiation law. Many attempts were made in the latter part of the nineteenth century and the early twentieth century to describe blackbody radiation mathematically, all doomed to failure before the recognition of quantum concepts, in particular, by Max Planck. Any attempt to describe the mechanisms surrounding the Planck theory would be superfluous here. Suffice it to say that the basis for the theory can be explained from an examination of the experimental curve shown in Fig. 1 (borrowed from Richtmeyer and Kennard³) determined from the examination of the radiation from a “blackbody” at several different temperatures. By plotting the points one concludes, as shown on the graph, that the spectral radiant exitance, M_λ , is equal to the product of the negative fifth power of λ times some function of the product, λT , where λ is the wavelength (in micrometers) of the radiation and T is the absolute temperature of the radiator (in Kelvins). Confirmation of this fact is shown in Table 1. The radiant exitance values in the table were calculated

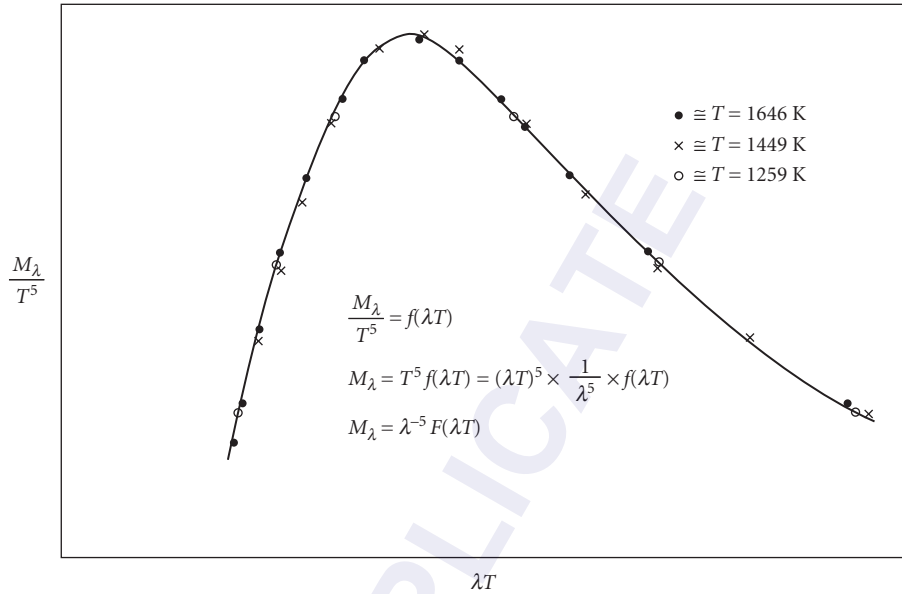


FIGURE 1 Experimental verification of the blackbody displacement law.

TABLE 1 Numerical Values to Support Fig. 1

T	λ	λT	M_λ	M_λ/T^5	
1,646	1.761	2,899	4.95	4.09×10^{-16}	
	1	1,646	1.91	1.56×10^{-16}	
	0.8	1,317	0.653	5.40×10^{-17}	
	2	3,291	4.77	3.90×10^{-16}	
	3	4,938	2.81	2.33×10^{-16}	
	5	8,230	0.803	6.65×10^{-17}	
	7	11,522	0.285	2.36×10^{-17}	
	10	16,460	0.085	7.04×10^{-18}	
	1,449	2	2,898	2.62	4.10×10^{-16}
		1.14	1,646	1.02	1.60×10^{-16}
0.909		1,317	0.346	5.42×10^{-17}	
2.27		3,291	2.52	3.95×10^{-16}	
3.41		4,938	1.49	2.32×10^{-16}	
5.68		8,230	0.425	6.65×10^{-17}	
7.95		11,522	0.151	2.36×10^{-17}	
11.36		16,460	0.045	7.06×10^{-18}	
1,259	2.3	2,896	1.3	4.11×10^{-16}	
	1.31	1,646	0.502	1.59×10^{-16}	
	1.046	1,317	0.171	5.41×10^{-17}	
	2.61	3,291	1.25	3.95×10^{-16}	
	3.91	4,938	0.741	2.34×10^{-16}	
	6.54	8,230	0.21	6.64×10^{-17}	
	9.15	11,522	0.0749	2.37×10^{-17}	
	13.07	16,460	0.0223	7.05×10^{-18}	

T = deg Kelvin; λ = μm ; M_λ = radiant exitance; $\text{W}\cdot\text{cm}^{-2}\cdot\text{ster}^{-1}\cdot\mu\text{m}^{-1}$

using the Infrared Radiance Calculator created by the author and found by choosing the term “Calculators” from the Military Sensing Information Analysis Center (SENSIAC) in a search of the Internet under www.sensiac.gatech.edu.

Postulating the quantum nature of the radiation, Planck, in a clever demonstration of the entropies resulting from small and large values of λT , was able to establish an expression for blackbody radiation as

$$M_{\lambda}(\lambda)d\lambda = c_1\lambda^{-5}(e^{c_2/\lambda T} - 1)^{-1}d\lambda$$

where

$$c_1 = 2\pi hc^2 = 3.7413 \times 10^4 \text{ W-cm}^{-2}\text{-}\mu\text{m}^4 \text{ (first radiation constant)}$$

and

$$c_2 = hc/k = 14388 \text{ }\mu\text{m-K (second radiation constant)}$$

h = Planck's constant = $6.6252 \times 10^{-34} \text{ W-s}^2$

k = Boltzmann constant = $1.38042 \times 10^{-23} \text{ W-s-K}^{-1}$

c = Velocity of light = $2.99793 \times 10^{10}\text{-s}^{-1}$

He later established the same equation from first principles.

When the Planck function is plotted on log-log paper the graph of Fig. 2 results. The special feature of this type of plot is that, regardless of the temperature of the blackbody, the shape of the curve is constant. It merely moves up and to the left (i.e., toward shorter wavelengths) as the temperature increases. The straight line, with a slope of -5 , drawn through the set of curves of Fig. 2, depicts the

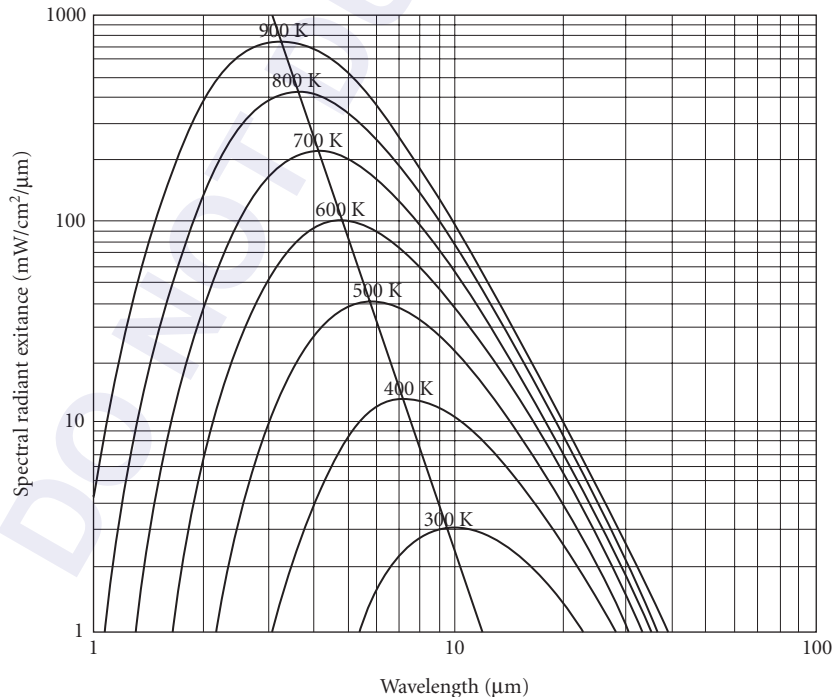


FIGURE 2 Spectral radiant exitance versus wavelength.

wavelength at which each one peaks, resulting in what is known as the Wein displacement law given for a specific temperature by

$$\lambda_{\max} T = 2897.9 \mu\text{m} - \text{K}^{-1}$$

Thus the peak of any Planck curve can be determined, given the temperature of the blackbody.

15.4 LABORATORY SOURCES

Standard Sources

The reader may be interested in the exposition by Quinn⁴ on the calculation of the emissivity of cylindrical cavities in which the method of DeVos⁵ is used. In a more recent paper Irani⁶ refers to the method of Gouffé⁷ for the construction of blackbody calibration sources. Quinn states that for certain constructions there are errors in the method of Gouffé. However, for a well-constructed source the shape of the construction is least at fault, since any heat-resistant material with a reasonably high surface emissivity will produce a resultant emissivity of better than 0.99. However, the accuracy of the value of the radiation for a given temperature depends not only on the emissivity but on generally high numerical powers of the temperature especially for high-temperature blackbodies. Therefore, very small variations of temperature over the inner surface of the source can cause relatively large errors in the radiation accuracy. Thus, great caution is used in creating a uniform temperature, resulting in the use of the fixed-point temperatures of various metals for the most basic and accurate calibration standards.

Blackbody Cavity Theory Radiation levels can be standardized by the use of a source that will emit a quantity of radiation that is both reproducible and predictable. Cavity configurations can be produced to yield radiation theoretically sufficiently close to Planckian that it is necessary only to determine what the imprecision is. Several theories have been expounded over the years to calculate the quality of a blackbody simulator.*

*The Method of Gouffé.*⁷ For the total emissivity of the cavity forming a blackbody (disregarding temperature variations) Gouffé gives

$$\epsilon_0 = \epsilon'_0(1 + K) \quad (1)$$

where

$$\epsilon'_0 = \frac{\epsilon}{\epsilon \left(1 - \frac{s}{S}\right) + \frac{s}{S}} \quad (2)$$

and $K = (1 - \epsilon) \left[\left(\frac{s}{S} \right) - \left(\frac{s}{S_0} \right) \right]$, and is always nearly zero—it can be either positive or negative.

ϵ = emissivity of materials forming the blackbody surface

s = area of aperture

S = area of interior surface

S_0 = surface of a sphere of the same depth as the cavity in the direction normal to the aperture

Figure 3 is a graph for determining the emissivities of cavities with simple geometric shapes. In the lower section, the value of the ratio s/S is given as a function of the ratio $1/r$. (Note the scale

*Generically used to describe those sources designed to produce radiation that is nearly Planckian.

change at the value for $1/r = 5$.) The values of ϵ'_0 is found by reading up from this value of the intrinsic emissivity of the cavity material. The emissivity of the cavity is found by multiplying ϵ'_0 by the factor $(1 + K)$.

When the aperture diameter is smaller than the interior diameter of the cylindrical cavity, or the base diameter of a conical cavity, it is necessary to multiply the value of s/S determined from the graph by $(r/R)^2$, which is the ratio of the squares of the aperture and cavity radii (Fig. 3).

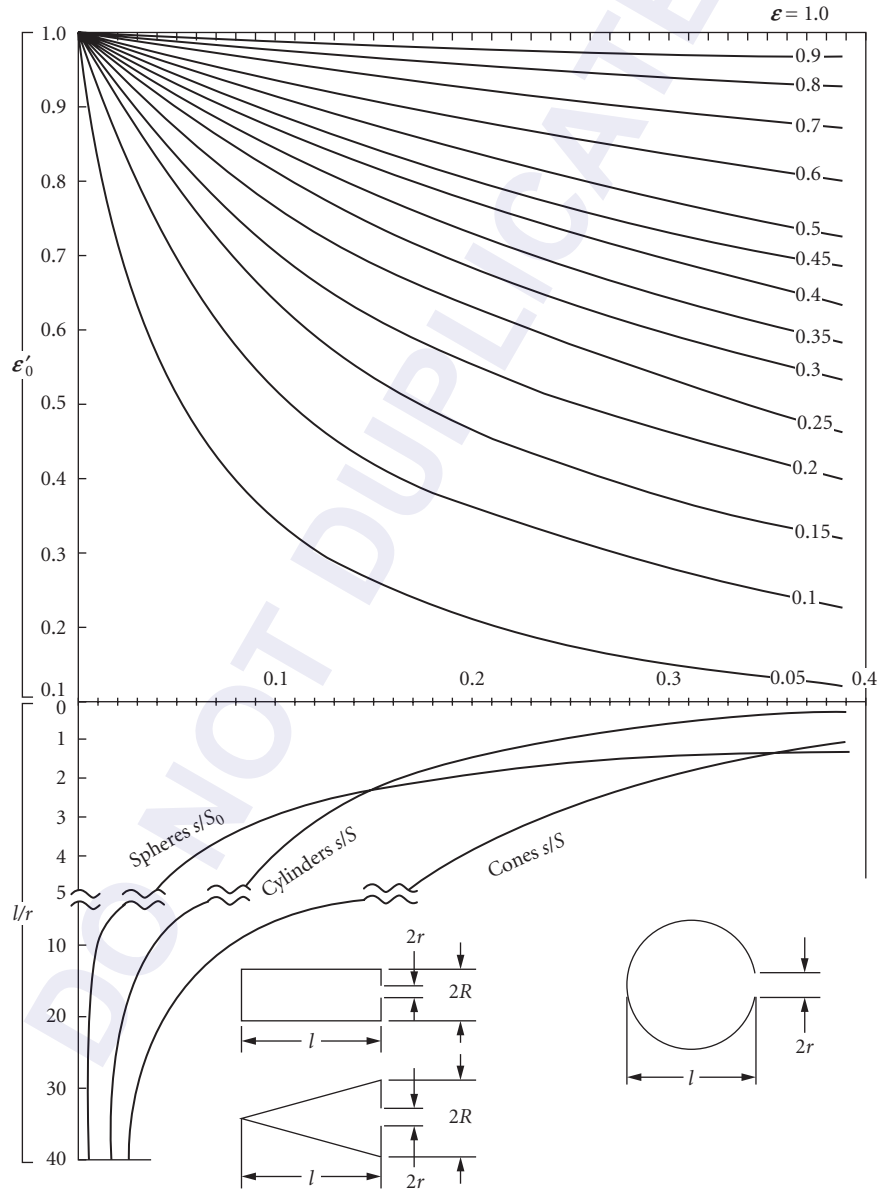


FIGURE 3 Emissivities of conical, spherical, and cylindrical cavities.

It is important to be aware of the effect of temperature gradients in a cavity. This factor is the most important in determining the quality of a blackbody, since it is not very difficult to achieve emissivities as near to unity as desired.

Manufacturers of blackbody simulators strive to achieve uniform heating of the cavity because it is only under this condition that the radiation is Planckian. The ultimate determination of a radiator that is to be used as the standard is the quality of the radiation that it emits.

A recent investigation on comparison of IR radiators is presented by Leupin et al.⁸ There has been, incidentally, a division historically between the standards of photometry and those used to establish thermal radiation and the thermodynamic temperature scale. Thus, in photometry the standard has changed from the use of candles, the Carcel lamp, the Harcourt pentane lamp, and the Hefner lamp⁹ to more modern radiators.

Baseline Standard of Radiation Although there is no internationally accepted standard of radiation, the National Institute of Standards and Technology (NIST) uses as its substitute standard the goldpoint blackbody (see Fig. 4),¹⁰ which fixes one point on the international temperature scale, now reported to be 1337.33 ± 0.34 K. Starting from this point, NIST is able to transfer fixed radiation values to working standards of radiation through an accurately constructed variable-temperature radiator as shown in Fig. 5.¹¹

The goldpoint blackbody is shown mainly for information. It is quite feasible to build a replica of the variable-temperature radiator, especially in the laboratory equipped to do fundamental radiation measurements.

Working Standards of Radiation For the calibration of instruments in the ordinary laboratory, the user is likely to use a source which is traceable to NIST, and generally supplied by NIST or one of the recognized vendors of calibrated sources, mainly in the form of a heated filament, a gaseous arc enclosed in an envelope of glass or quartz (or fused silica), or in glass with a quartz or sapphire window.

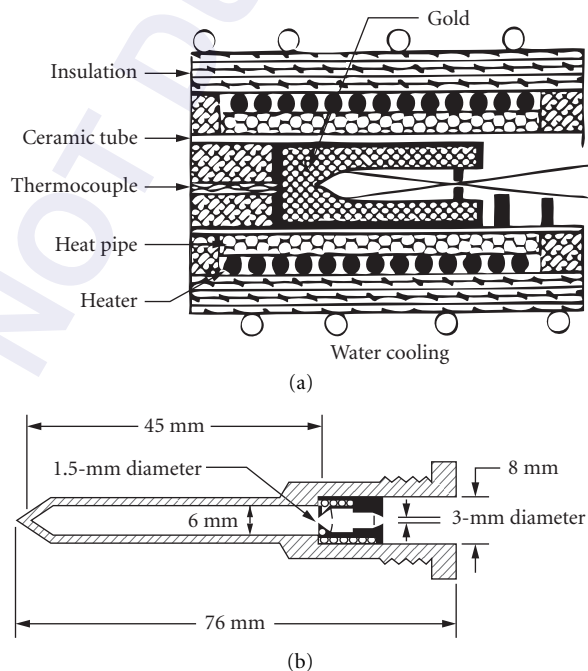


FIGURE 4 (a) Cross section of heat-pipe blackbody furnace. (b) Blackbody inner cavity dimensions.

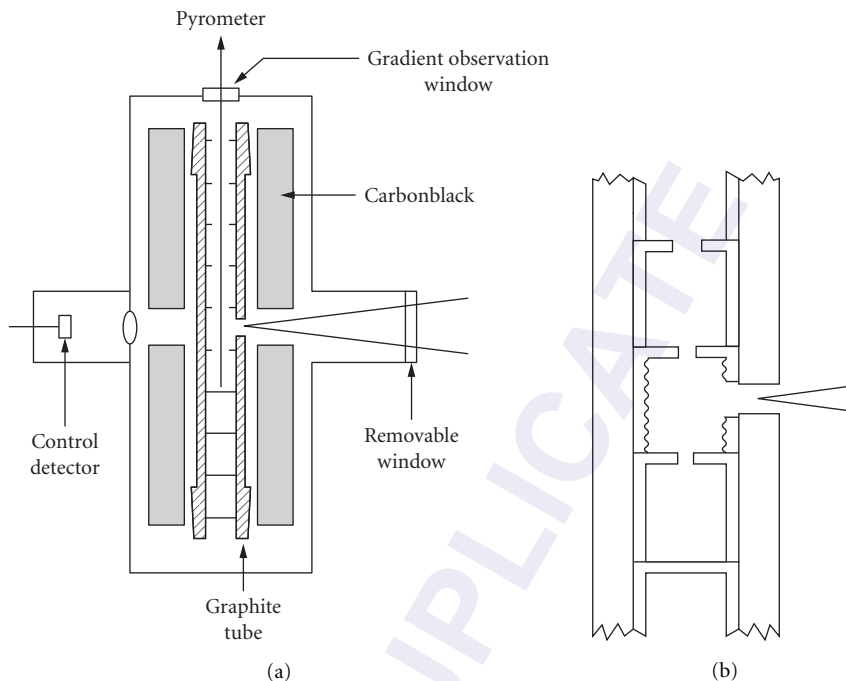


FIGURE 5 (a) Variable-temperature blackbody schematic. (b) Central section of variable-temperature blackbody.

Any source whose radiation deviates from that described by Planck's law is nonblackbody. Even the sources previously described are not strictly blackbodies, but can come as close as the user desires within the constraints of bulk and price. Any other source has an emissivity less than unity and can, and usually does, have a highly variable spectral emissivity. The lamps used by NIST, for example (see the following), fit into this category, but they differ in one large respect. They are transfer standards which have been carefully determined to emit specified radiation within certain specific spectral regions.

The following discussion of these types of sources is reproduced (in some cases with slight modifications), with permission, from the NIST Special Publication 250.¹² For specific details of calibration, and for the exact source designations, the user should contact NIST at

*U.S. Department of Commerce
National Institute of Standards and Technology
Office of Physical Measurement Services
Rm. B362, Physics Bldg.
Gaithersburg, MD 20899
Photometric Standards*

The following text on working standards is left untouched from the presentation of the earlier edition of the *Handbook* because there do not appear to be significant changes since the publication of that material. However, to the extent that information on the Internet is current, a reasonable complement to the information published in this chapter would be a search of the Internet at the location designated www.physics.nist.gov. The reader will find numerous features of the National Institute of Science and Technology from which to choose the service or other information required.

1. Sources/Lamps

Luminous Intensity Standard (100-W Frosted Tungsten Lamp, 90 cd)
 Luminous Intensity Standard (100-W Frosted Tungsten Lamp, color temp., 2700 K)
 Luminous Intensity Standard (100-W Frosted Tungsten Lamp, color temp., 2856 K)
 Luminous Intensity Standard (500-W Frosted Tungsten Lamp, 700 cd)
 Luminous Intensity Standard (1000-W Frosted Tungsten Lamp, 1400 cd)
 Luminous Intensity Standard (1000-W Frosted Tungsten Lamp, color temp., 2856 K)
 Luminous Flux Standard (25-W Vacuum Lamp about 270 lm)
 Luminous Flux Standard (60-W Gas-filled Lamp about 870 lm)
 Luminous Flux Standard (100-W Gas-Filled Lamp about 1600 lm)
 Luminous Flux Standard (200-W Gas-Filled Lamp about 3300 lm)
 Luminous Flux Standard (500-W Gas-Filled Lamp about 10,000 lm)
 Luminous Flux Standard (Miniature Lamps 7 sizes 6 to 400 lm)
 Airway Beacon Lamps for Color Temperature (500-W, 1 point in range, 2000 to 3000 K)

2. General Information

Calibration services provide access to the photometric scales realized and maintained at NIST. Lamp standards of luminous intensity, luminous flux, and color temperature, as described next, are calibrated on a routine basis.

a. Luminous Intensity Standards

Luminous intensity standard lamps supplied by NIST [100-W (90–140 cd), 500-W (approximately 700 cd), and 1000-W (approximately 1400 cd) tungsten filament lamps with C-13B filaments in inside-frosted bulbs and having medium bipost bases] are calibrated at either a set current or a specified color temperature in the range 2700 to 3000 K. Approximate 3-sigma uncertainties are 1 percent relative to the SI unit of luminous intensity and 0.8 percent relative to NIST standards.

b. Luminous Flux Standards

Vacuum tungsten lamps of 25 W and 60-, 100-, 200-, and 500-W gas-filled tungsten lamps that are submitted by customers are calibrated. Lamps must be base-up burning and rated at 120 V. Approximate 3-sigma uncertainties are 1.4 percent relative to SI units and 1.2 percent relative to NIST standards. Luminous flux standards for miniature lamps producing 6 to 400 lm are calibrated with uncertainties of about 2 percent.

c. Airway Beacon Lamps

Color temperature standard lamps supplied by NIST (airway beacon 500-W medium bipost lamps) are calibrated for color temperature in the range 2000 to 3000 K with 3-sigma uncertainties ranging from 10 to 15°.

IR Radiometric Standards

General Information

a. Spectral Radiance Ribbon Filament Lamps

These spectral radiance standards are supplied by NIST. Tungsten, ribbon filament lamps (30A/T24/13) are provided as lamp standards of spectral radiance. The lamps are calibrated at 34 wavelengths from 225 to 2400 nm, with a target area 0.6 mm wide by 0.8 mm high. Radiance temperature ranges from 2650 K at 225 nm and 2475 K at 650 nm to 1610 K at 2400 nm, with corresponding uncertainties of 2, 0.6, and 0.4 percent. For spectral radiance lamps, errors are stated as the quadrature sum of individual uncertainties at the three standard deviation level.

Figure 6 summarizes the measurement uncertainty for NIST spectral radiance calibrations.

b. Spectral Irradiance Lamps

These spectral irradiance standards are supplied by NIST. Lamp standards of spectral irradiance are provided in two forms. Tungsten filament, 1000 W quartz halogen-type FEL lamps are calibrated at 31 wavelengths in the range 250 to 2400 nm. At the working distance of 50 cm, the

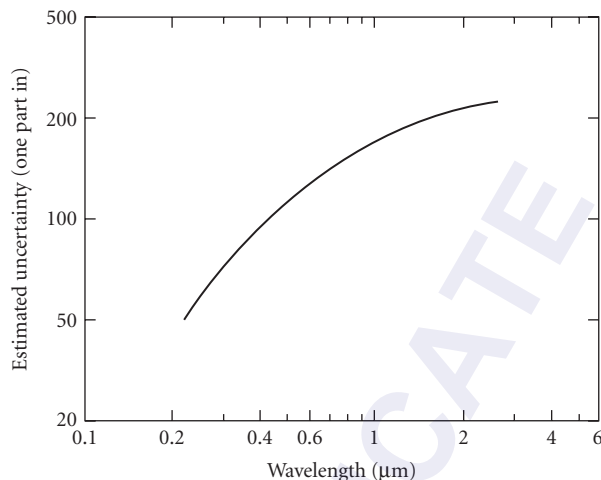


FIGURE 6 Uncertainties for NIST spectral radiance calibrations.

lamps produce 0.2 W/cm²/cm at 250 nm, 220 W/cm²/cm at 900 nm, 115 W/cm²/cm at 1600 nm, and 40 W/cm²/cm at 2400 nm, with corresponding uncertainties of 2.2, 1.3, 1.9, and 6.5 percent. For spectral irradiance lamps, errors are stated as the quadrature sum of individual uncertainties at the three standard deviation level. Deuterium lamp standards of spectral irradiance are also provided and are calibrated at 16 wavelengths from 200 to 350 nm. At the working distance of 50 cm, the spectral irradiance produced by the lamp ranges from about 0.5 W/cm²/cm at 200 nm and 0.3 W/cm²/cm at 250 nm to 0.07 W/cm²/cm at 350 nm. The deuterium lamps are intended primarily for the spectral region 200 to 250 nm. The approximate uncertainty relative to SI units is 7.5 percent at 200 nm and 5 percent at 250 nm. The approximate uncertainty in relative spectral distribution is 3 percent. It is strongly recommended that the deuterium standards be compared to an FEL tungsten standard over the range 250 to 300 nm each time the deuterium lamp is lighted to take advantage of the accuracy of the relative spectral distribution.

Figure 7 summarizes the measurement uncertainty for NIST spectral irradiance calibrations of type FEL lamps.

Radiometric Sources in the Far Ultraviolet

1. Sources

Spectral Irradiance Standard, Argon Mini-Arc (140 to 330 nm)

Spectral Radiance Standard, Argon Mini-Arc (115 to 330 nm)

Spectral Irradiance Standard, Deuterium Arc Lamp (165 to 200 nm)

2. General Information

a. Source Calibrations in the Ultraviolet

NIST maintains a collection of secondary sources such as argon maxi-arcs, argon mini-arcs, and deuterium arc lamps in the near and vacuum ultraviolet radiometric standards program to provide calibrations for user-supplied sources. The calibrations of these sources are traceable to a hydrogen arc whose radiance is calculable and which NIST maintains as a primary standard. The collection also includes tungsten strip lamps and tungsten halogen lamps whose calibrations are based on a blackbody rather than a hydrogen arc. Customer-supplied sources are calibrated in both radiance and irradiance by comparing them with NIST secondary standards.

Argon arcs are used to calibrate other sources in the wavelength range 115 to 330 nm for radiance and 140 to 330 nm for irradiance. The lower wavelength limit is determined in radiance by the cutoff of the magnesium fluoride windows used in the arcs, and in irradiance by the decrease

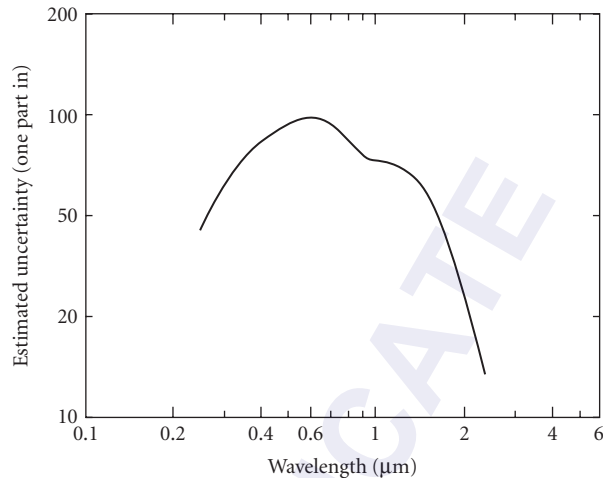


FIGURE 7 Uncertainties for NIST spectral irradiance calibrations of type FEL lamps.

in signal produced by the addition of a diffuser. Deuterium arc lamps are used in the range 165 to 200 nm, with the low wavelength cutoff due to the onset of blended molecular lines.

The high wavelength limit is the starting point of the range for the Radiometric Standards group. The tungsten lamps are used at 250 nm and above, since their signals are too weak at shorter wavelengths. It should be noted that the wavelength range of the NIST arcs partially overlap the range of tungsten lamps, thus providing an independent check on calibrations.

An argon mini-arc lamp supplied by the customer is calibrated for spectral irradiance at 10-nm intervals in the wavelength region 140 to 300 nm. Absolute values are obtained by comparison of the radiative output with laboratory standards of both spectral irradiance and spectral radiance. The spectral irradiance measurement is made at a distance of 50 cm from the field stop. Uncertainties are estimated to be less than ± 10 percent in the wavelength region 140 to 200 nm and within ± 5 percent in the wavelength region 200 to 330 nm. A measurement of the spectral transmission of the lamp window is included in order that the calibration be independent of possible window deterioration or damage. The uncertainties are taken to be two standard deviations.

The spectral radiance of argon mini-arc radiation sources is determined to within an uncertainty of less than 7 percent over the wavelength range 140 to 330 nm and 20 percent over the wavelength range 115 to 140 nm. The calibrated area of the 4-mm diameter radiation source is the central 0.3-mm diameter region. Typical values of the spectral radiance are: at 250 nm, $L(\lambda) = 30 \text{ mW/cm}^2/\text{nm/sr}$; and at 150 nm, $L(\lambda) = 3 \text{ mW/cm}^2/\text{nm/sr}$. The transmission of the demountable lamp window and that of an additional MgF_2 window are determined individually so that the user may check periodically for possible long-term variations.

The deuterium arc lamp is calibrated at 10 wavelengths from 165 to 200 nm, at a distance of 50 cm, at a spectral irradiance of about $0.5 \text{ W/cm}^2/\text{cm}$ at 165 nm, $0.5 \text{ W/cm}^2/\text{cm}$ at 170 nm, and $0.5 \text{ W/cm}^2/\text{cm}$ at 200 nm. The approximate uncertainty relative to SI units is estimated to be less than 10 percent. The lamp is normally supplied by NIST and requires 300 mA at about 100 V.

15.5 COMMERCIAL SOURCES

The commercial sources described here are derived from the 1995 edition of *The Handbook of Optics*, which were taken from catalogs available at the time and from the literature of the day and prior thereto, providing choices that have obviously been available for years. Evidently changes in

the makeup of these products are slower than those of other areas of technology, making it reasonable to retain the same sources in this chapter, mainly as examples of the types that are available. With the universality of the Internet, accessibility of information on various sources of radiation, far beyond what is attainable from a limited collection of company catalogs, is at one's fingertips. Thus, it is recommended that, in seeking information on various sources, one use the examples in the text as a reference to what can be found currently on the Internet. Experience demonstrates that, in many cases, due to the stability of the lamp industry, there will be few changes between what is found currently on the Internet and what appears in the text of this chapter.

Obviously the choice of a source is dependent on the application. Many, if not most of the sources described are multipurpose ones, although most of them have been selected specifically for scientific study, tailored in a way to produce an image amenable to different optical systems. For basic research it is usually essential to have a blackbody source, especially for infrared research, that is traceable to a Standards Laboratory, along with traceable secondary standards for calibrating research instrumentation. Many of the sources can be used to produce spectra, which are capable of calibrating spectral measuring instrumentation. Other sources are included mostly to provide the user with an array of choices.

Blackbody Simulators

Virtually any cavity can be used to produce radiation of high quality, but practicality limits the shapes to a few. The most popular shapes are cones and cylinders, the former being more popular. Spheres, combinations of shapes, and even flat-plate radiators are used occasionally. Blackbodies can be bought rather inexpensively, but there is a fairly direct correlation between cost and quality (i.e., the higher the cost the better the quality).

Few manufacturers specialize in blackbody construction. Some, whose products are specifically described here, have been specializing in blackbody construction for many years. Other companies of this description may be found, for example, in the latest *Lasers and Optronics Buying Guide*¹³ or the latest *Photonics Directory of Optical Industries*.¹⁴ These references are the latest as of the writing of this work. It is expected that they will continue in succeeding years.

A large selection of standard (or blackbody) radiators is offered by Electro-Optical Industries, Inc. (EOI), Santa Barbara, California.* Most blackbodies can be characterized as one of the following: primary, secondary, or working standard. The output of the primary must, of course, be checked with those standards retained at NIST. Figure 8 pictures an EOI blackbody and its controller. Figure 9 pictures a similar blackbody from Mikron, Inc. and its controller. All of the companies sell separate apertures (some of which are water cooled) for controlling the radiation output of the radiators. Another piece of auxiliary equipment which can be purchased is a multispeed chopper. It is impossible to cite all of the companies that sell these kinds of sources; therefore, the reader is referred to one of the buyers' guides already referenced for a relatively complete list. It is prudent to shop around for the source that suits one's own purpose.

Figure 10 demonstrates a less conventional working standard manufactured by EOI. Its grooves-and-honeycomb structure is designed to improve the absorptance of such a large and open structure. A coating with a good absorbing paint increases its absorptance further.

Incandescent Nongaseous Sources (Exclusive of High-Temperature Blackbodies)

Nernst Glower[†] The Nernst glower is usually constructed in the form of a cylindrical rod or tube from refractory materials (usually zirconia, yttria, beria, and thoria) in various sizes. Platinum leads at the ends of the tube conduct power to the glower from the source. Since the resistivity of

*Many of the sources in the text are portrayed using certain specific company products, only for the sake of demonstration. This does not necessarily imply an endorsement of these products by the author. The reader is encouraged in all cases to consult the *Photonics Directory of Optical Industries*¹⁴ or a similar directory for competitive products.

[†]Since Nernst glower is probably obsolete, this section is retained only for historical purpose.

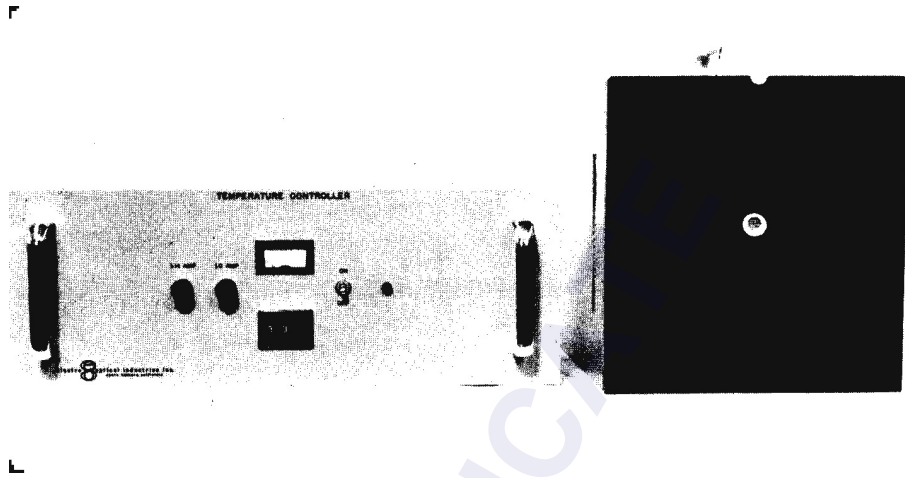


FIGURE 8 EOI blackbody.

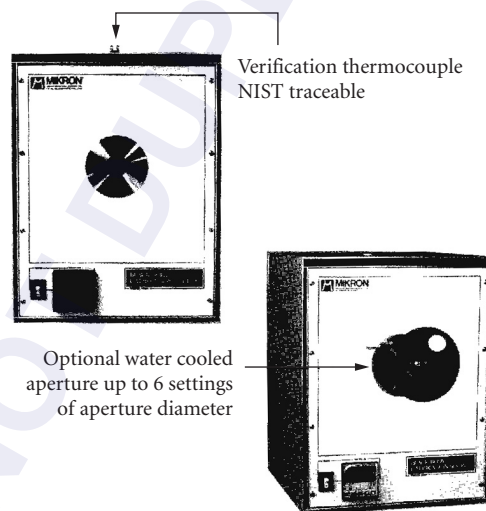


FIGURE 9 Mikron blackbody.

the material at room temperature is quite high, the working voltage is insufficient to get the glower started. Once started, its negative temperature coefficient-of-resistance tends to increase current, which would cause its destruction, so that a ballast is required in the circuit. Starting is effected by applying external heat, either with a flame or an adjacent electrically heated wire, until the glower begins to radiate.

Data from a typical glower are as follows:

1. Power requirements: 117 V, 50 to 60 A, 200 W
2. Color temperature range: 1500 to 1950 K
3. Dimension: 0.05-in. diameter by 0.3 in.

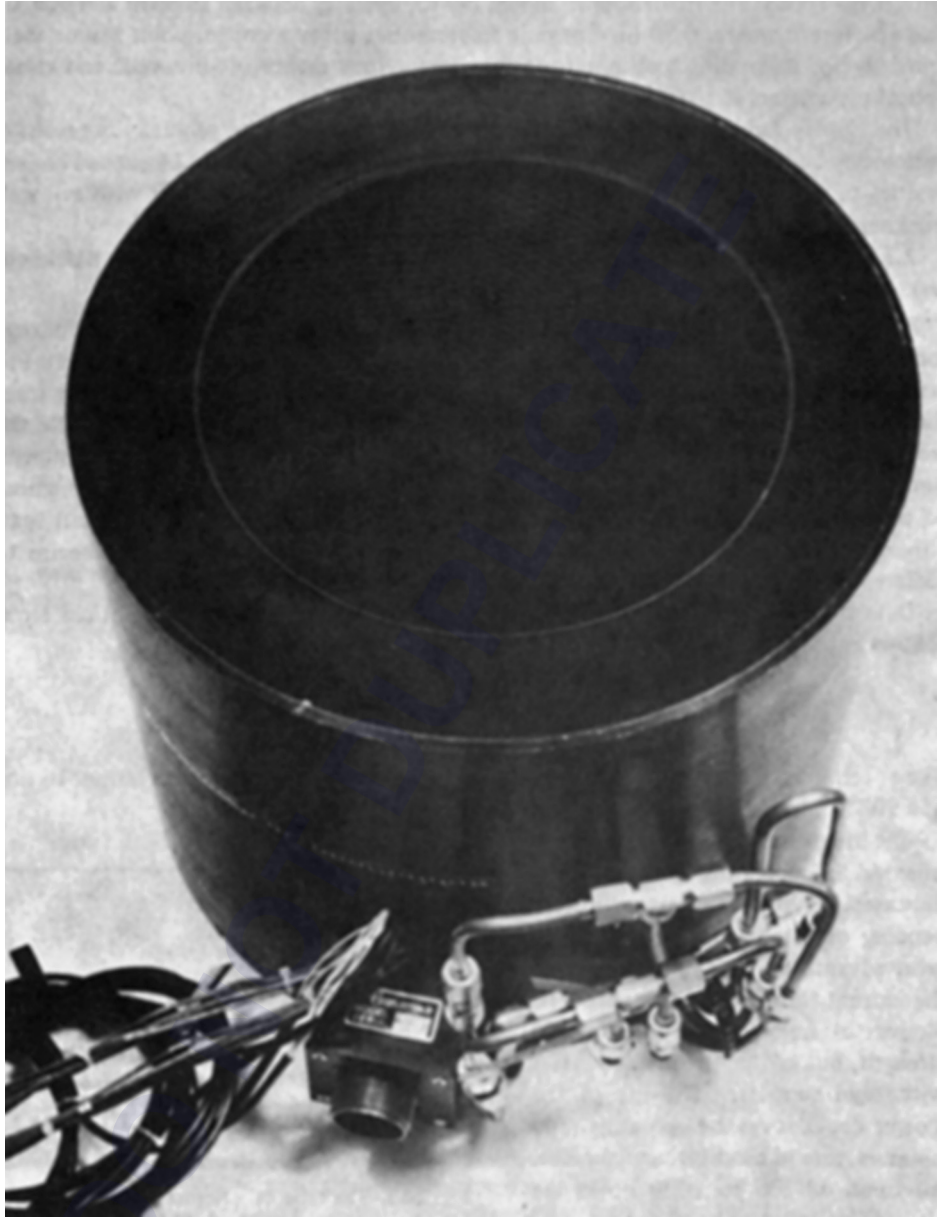


FIGURE 10 EOI model 1965. This model is 12 in. in diameter and 9 in. deep. The base is an array of intersecting conical cavities. The walls are hex-honeycomb and the temperature range is 175 to 340 K.

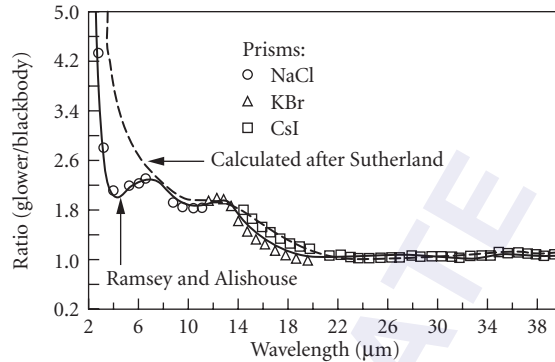


FIGURE 11 The ratio of a Nernst glower to a 900°C blackbody versus wavelength.

The spectral characteristics of a Nernst glower in terms of the ratio of its output to that of a 900°C blackbody are shown in Fig. 11.

The life of the Nernst glower diminishes as the operating temperature is increased. Beyond a certain point, depending on the particular glower, no great advantage is gained by increasing the current through the element. The glower is fragile, with low tensile strength, but can be maintained intact with rigid support. The life of the glower depends on the operating temperature, care in handling, and the like. Lifetimes of 200 to 1000 hours are claimed by various manufacturers.

Since the Nernst glower is made in the form of a long thin cylinder, it is particularly useful for illuminating spectrometer slits. Its useful spectral range is from the visible region to almost 30 μm , although its usefulness compared with other sources diminishes beyond about 15 μm . As a rough estimate, the radiance of a glower is nearly that of a graybody at the operating temperature with an emissivity in excess of 75 percent, especially below about 15 μm . The relatively low cost of the glower makes it a desirable source of moderate radiant power for optical uses in the laboratory. The makers of spectroscopic equipment constitute the usual source of supply of glowers (or of information about suppliers).

Globar The globar is a rod of bonded silicon carbide usually capped with metallic caps which serve as electrodes for the conduction of current through the globar from the power source. The passage of current causes the globar to heat, yielding radiation at a temperature above 1000°C. A flow of water through the housing that contains the rod is needed to cool the electrodes (usually silver). This complexity makes the globar less convenient to use than the Nernst glower and necessarily more expensive. This source can be obtained already mounted, from a number of manufacturers of spectroscopic equipment. Feedback in the controlled power source makes it possible to obtain high radiation output.

Ramsey and Alishouse¹⁵ provide information on a particular sample globar as follows:

1. Power consumption: 200 W, 6 A
2. Color temperature: 1470 K

They also provide the spectral characteristics of the globar in terms of the ratio of its output to that of a 900°C blackbody. This ratio is plotted as a function of wavelength in Fig. 12. Figure 13¹⁶ is a representation of the spectral emissivity of a globar as a function of wavelength. The emissivity values are only representative and can be expected to change considerably with use.

Gas Mantle The Welsbach mantle is typified by the kind found in high-intensity gasoline lamps used where electricity is not available. The mantle is composed of thorium oxide with some additive

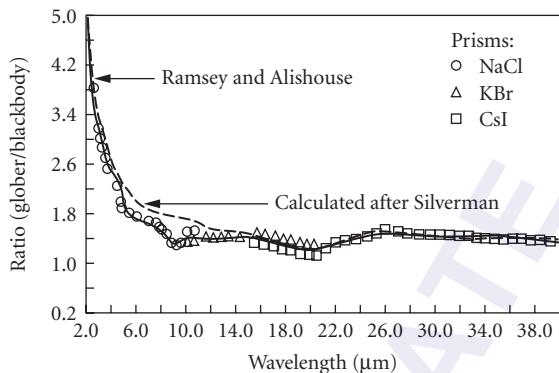


FIGURE 12 The ratio of a globar to a 900°C blackbody versus wavelength.

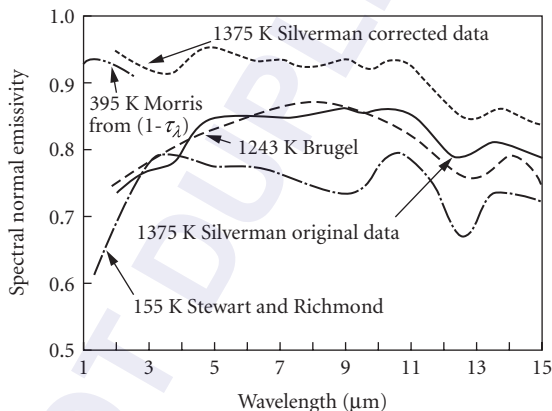


FIGURE 13 The spectral emissivity of a globar.

to increase its efficiency in the visible region. Its near-infrared emissivity is quite small, except for regions exemplified by gaseous emission, but increases considerably beyond 10 μm .

Ramsey and Alishouse¹⁵ provide information on a propane-heated sample from an experiment in which a comparison of several sources is made:

1. Color temperature: 1670 K
2. Dimensions: 25.4 by 38.1 mm

The spectral characteristics of the mantle in terms of the ratio of its output to that of a 900°C blackbody are shown in Fig. 14.

Pfund modified the gas mantle so that it became more a laboratory experimental source than an ordinary radiator. By playing a gas flame on an electrically heated mantle, he was able to increase its radiation over that from the gas mantle itself.¹⁷ Figure 15 shows a comparison of the gas mantle and the electrically heated gas mantle, with a Nernst glower. Strong¹⁸ points out that playing a flame against the mantle at an angle produces an elongated area of intense radiation useful for illuminating the slits of a spectrometer.

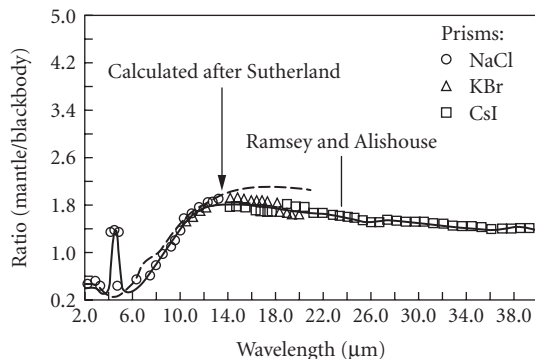


FIGURE 14 The ratio of the gas mantle to a 900°C blackbody versus wavelength.

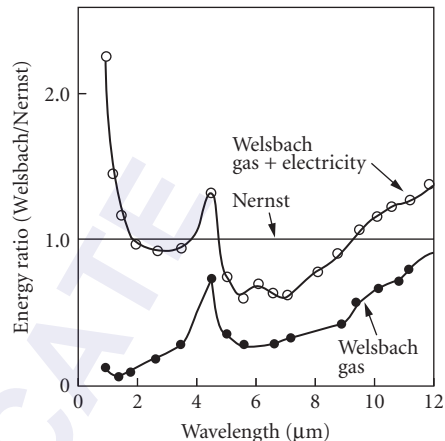


FIGURE 15 Emission relative to that of a Nernst glower (2240 K) of the gas-heated mantle (lower curve) and that of the mantle heated by gap plus electricity (upper curve).

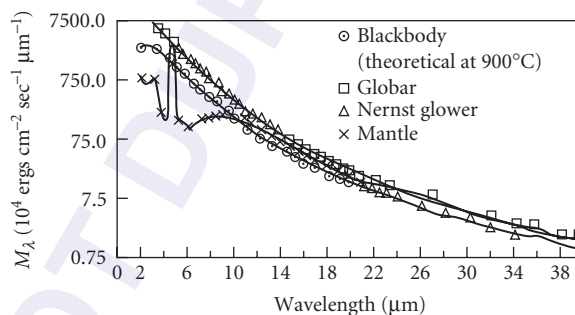


FIGURE 16 The spectral radiant emittances of a globar, Nernst glower, 900°C blackbody, and gas mantle versus wavelength.

Comparison of Nernst Glower, Globar, and Gas Mantle Figure 16 compares these three types of sources, omitting a consideration of differences in the instrumentation used in making measurements of the radiation from the sources.

Availability, convenience, and cost usually influence a choice of sources. At the very long wavelength regions in the infrared, the gas mantle and the globar have a slight edge over the Nernst glower because the Nernst glower (a convenient, small, and inexpensive source) does not have the power of the gas mantle and globar.

Tungsten-Filament Lamps A comprehensive discussion of tungsten-filament lamps is given by Carlson and Clark.¹⁹ Figures 17 to 19 show the configurations of lamp housings and filaments. The types and variations of lamps are too numerous to be meaningfully included in this chapter. The reader is referred to one of the buyer's guides for a comprehensive delineation of manufacturers from whom unlimited literature can be obtained.

Tungsten lamps have been designed for a variety of applications; few lamps are directed toward scientific research, but some bear directly or indirectly on scientific pursuits insofar as they can

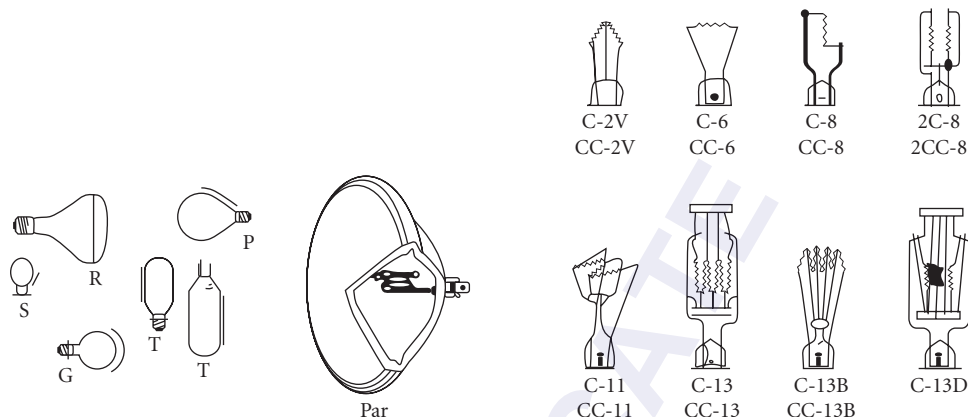


FIGURE 17 Bulk shapes most frequently used for lamps in optical devices. Letter designations are for particular shapes.

FIGURE 18 Most commonly used filament forms. Letters designate the type of filament.

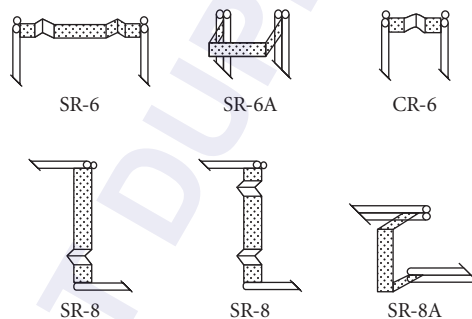
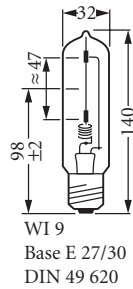


FIGURE 19 Ribbon-type tungsten filaments. Type designations are by number.

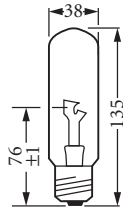
provide steady sources of numerous types of radiation. One set of sources cited here, particularly for what the manufacturer calls their scientific usefulness, is described in *Lamps for Scientific Purposes*.²⁰ Their filament structures are similar to those already described, but their designs reduce extraneous radiation and ensure the quality and stability of the desired radiation. The lamps can be obtained with a certification of their calibration values.

The physical descriptions of some of these sources are given in Fig. 20. Applications (according to the manufacturer, Osram) are photometry, pyrometry, optical radiometry, sensitometry, spectroscopy, spectrometry, polarimetry, saccharimetry, spectrophotometry, colorimetry, microscopy, microphotography, microprojection, and stroboscopy.

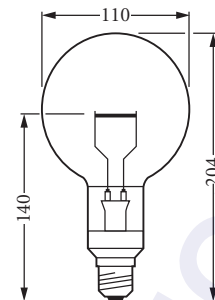
Quartz Envelope Lamps These are particularly useful as standards because they are longer lasting (due to action of iodine in the quartz-iodine series), can be heated to higher temperatures, are sturdier, and can transmit radiation to longer wavelengths in the infrared than glass-envelope lamps. Studer and Van Beers²¹ have shown the spectral deviation to be expected of lamps containing no iodine. The deviation, when known, is readily acceptable in lieu of the degradation in the lamp caused by the absence of iodine. The particular tungsten-quartz-iodine lamps used in accordance



WI 9
Base E 27/30
DIN 49 620



WI 14
Base E 27/30
DIN 49 620



WI 16/G
Base E 27/51 x 39

Order reference	Upper limits for electric data (V) (A)		Color temperature T_f max.*	Luminance T_s max.	Dimensions of luminous width (mm)	Area height (mm)	Burning position [†]	Base
-----------------	---	--	-------------------------------	----------------------	-----------------------------------	------------------	-------------------------------	------

Lamps for scientific purposes

WI 9	8.5	6	2856 K	—	0.2	47	s	E 27
WI 14	5	16	—	2400 K	1.6	8	s	E 27
WI 16/G	9	16	—	2600 K	21	1.6	s+h	E 27
WI 17/G	9	16	—	2600 K	1.6	20	s	E 27
WI 40/G	31	6	2856 K	—	18	18	s+h	E 27
WI 41/G	31	6	2856 K	—	18	18	s+h	E 27

Lamps for scientific purposes are gas-filled, incandescent lamps for calibration of luminous intensity, luminous flux, luminance (spectral radiant temperature), color temperature (luminance temperature and spectral radiance distribution). A test certificate can be issued for these types of lamps.

Also for other types of lamp with sufficiently constant electric and photometric data, a test certificate can be issued. To order a test certificate, the order reference of the lamp, the type of measurement, and the desired burning position have to be given. Example: Lamp 41/G, measurement of the electric data and the luminous intensity for $T_f = 2856$ K (light type A), burning position vertical, base up.

Variables for which test certificates can be issued are shown in the following table by +. The sign (+) indicates that certificates can be issued for variables although the lamps were not designed for such measurements.

Type of lamp	Light intensity	Luminous flux	Luminance	Color temperature	Spectral radiance distribution [‡]
WI 9	+	—	—	+	—
WI 14	(+)	—	+	+	+ 300–800 mm
WI 16/G	(+)	—	+	+	+ 300–800 mm
WI 17/G	—	—	+	(+)	+ 250–800 mm 250–2500 mm
WI 40/G	+	+	—	(+)	—
WI 42/G	+	—	—	+	—

Description

WI 9:
Lamp with uncoiled straight filament.

WI 14:
Tungsten ribbon lamp with tubular bulb. The portion of the tungsten ribbon to be utilized for measurement is mounted parallel to the lamp axis and positioned approx. 8 mm off-axis in the measuring direction.

* The color temperature of 2856 K corresponds with light-type A (DIN 5035).

† s = vertical (base down); h = vertical (base up).

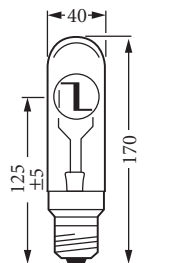
‡ Only for additional measurement of luminance or color temperature.

FIGURE 20a Lamps for scientific purposes. (Note dimensions in mm.)

with the NIST are described earlier in this chapter. Others can be obtained in a variety of sizes and wattages from General Electric, Sylvania, and a variety of other lamp manufacturers and secondary sources.

Carbon Arc

The carbon arc has been passed down from early lighting applications in three forms: low-intensity arc, flame, and high-intensity arc. The low- and high-intensity arcs are usually operated on direct



WI 16/G:

Tungsten ribbon lamp with spherical bulb. Horizontal tungsten ribbon with a small notch to indicate the measuring point. The ribbon is positioned approx. 3 mm off-axis.

WI 17/G:

Tungsten ribbon lamp with horn-shaped bulb. The bulb has a tubular extension with a sealed-on quartz glass window (homogenized ultrasil). Vertical tungsten ribbon with a small notch to indicate the measuring point.

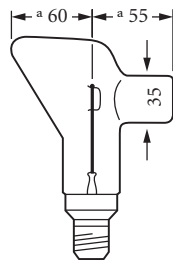
WI 40/G:

Standard lamp for total radiation, luminous flux, and color temperatures with conic bulb. The bulb shape prevents reflections in the direction of the plane normal of the luminous area, which is formed by the meandrous-shaped filament.

WI 41/G:

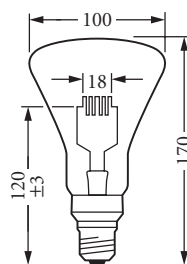
Standard lamp for light intensity and color temperature with conic bulb. Differs from the WI 40/G lamp by a black, opaque coating which covers one side of the bulb.

A window is left open in the coating opposite the filament, through which over an angle of approx. $\pm 3^\circ$ a constant light intensity is emitted. The black coating prevents stray light being reflected in the measuring direction.



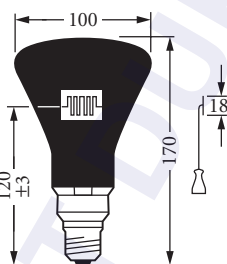
WI 17/G

Base E 27/51 × 39



WI 40/G

Base E 27/51 × 39



WI 41/G

Base E 27/51 × 39

FIGURE 20b (Continued)

current; the flame type adapts to either direct or alternating current. In all cases, a ballast must be used. In the alternating current arc, the combined radiation from the two terminals is less than that from the positive crater of the direct-current arc of the same wattage.²²

Spatial variation in the amount of light energy across the crater of dc arcs for different currents is shown in Fig. 21.

The carbon arc is a good example of an open arc, widely used because of its very high radiation and color temperatures (from approximately 3800 to 6500 K, or higher). The rate at which the material is consumed and expended during burning (5 to 30 cm/h) depends on the intensity of the arc. The arc is discharged between two electrodes that are moved to compensate for the rate of consumption of the material. The anode forms a crater of decomposing material which provides a center of very high luminosity. Some electrodes are hollowed out and filled with a softer carbon material which helps keep the arc fixed in the anode and prevents it from wandering on the anode surface.

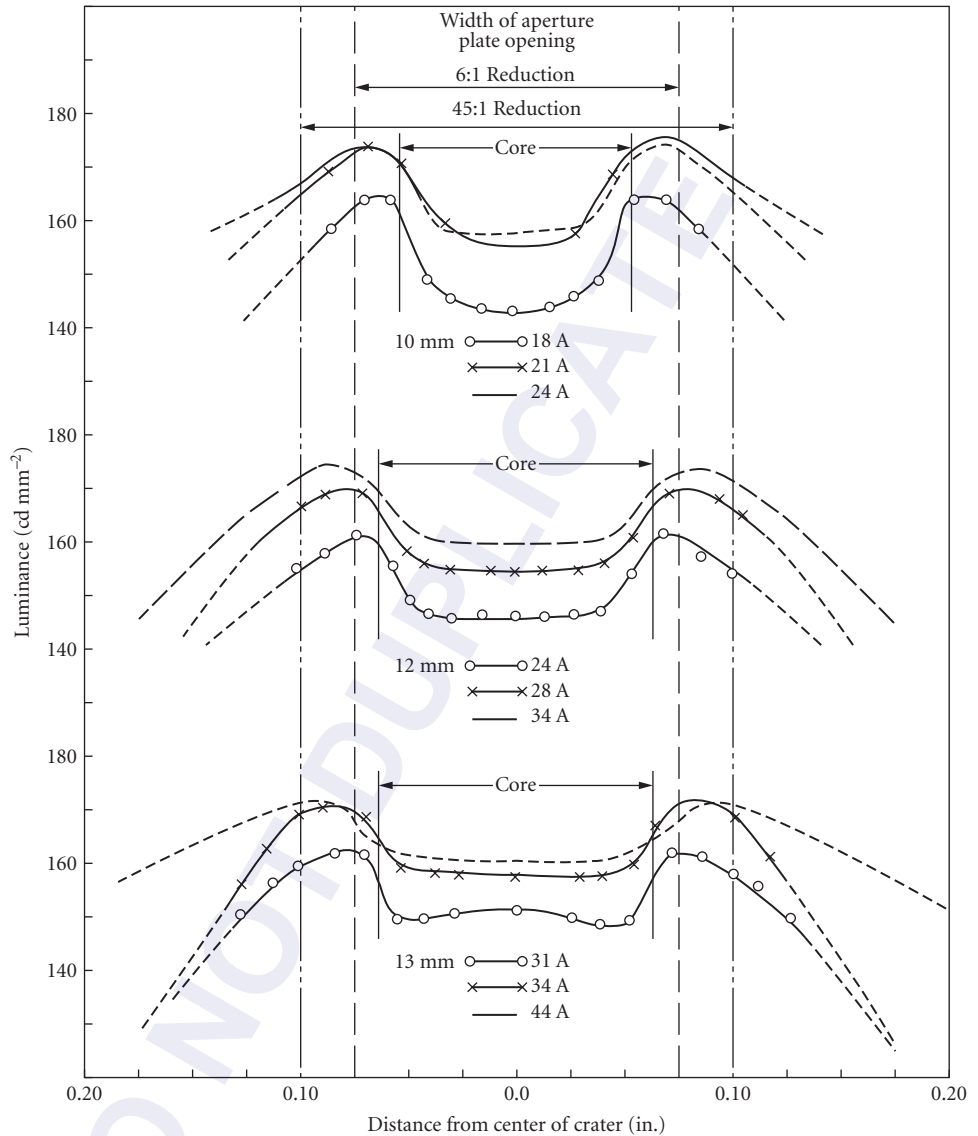


FIGURE 21 Variations in brightness across the craters of 10-, 12-, and 13-mm positive carbons of dc plain arcs operated at different currents in the regions of recommended operation.

In some cored electrodes, the center is filled with whatever material is needed to produce desired spectral characteristics in the arc. In such devices, the flame between the electrodes becomes the important center of luminosity, and color temperatures reach values as high as 8000 K.²² An example of this so-called flaming arc is shown in Fig. 22a. Figure 22b and c shows the low-intensity dc carbon arc and the high-intensity dc carbon arc with rotating positive electrodes. Tables 2 and 3 give characteristics of dc high-intensity and flame carbon arcs.

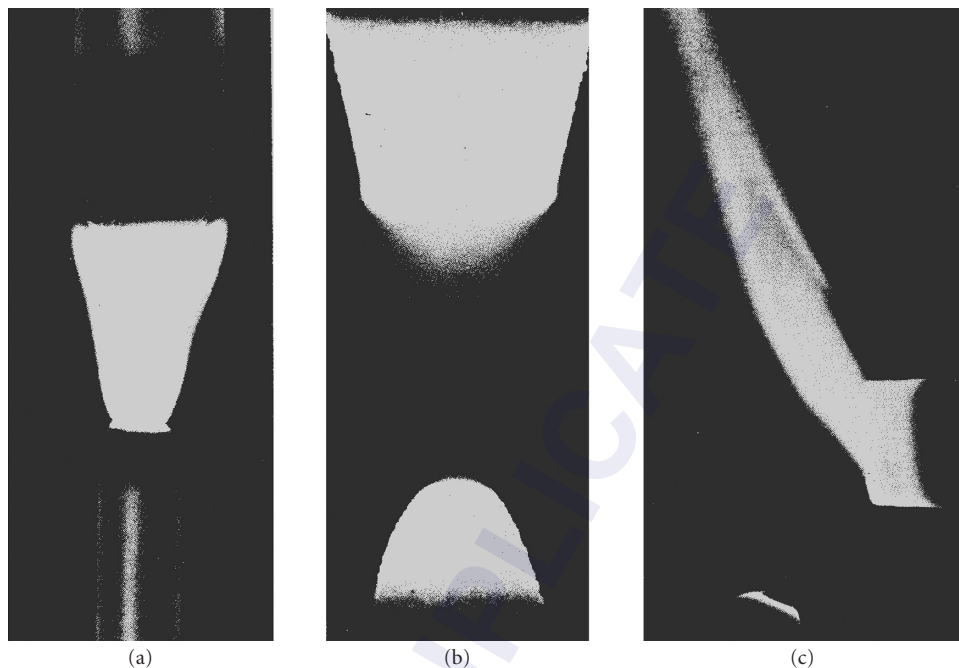


FIGURE 22 Various types of carbon arc: (a) flame type; (b) low-intensity dc arc; and (c) high-intensity dc arc with rotating positive carbon.

A spectrum of low-intensity arc (Fig. 23) shows the similarity between the radiation from it and a 3800-K blackbody, except for the band structure at 0.25 and 0.39 μm . In Koller²² an assortment of spectra are given for cored carbons containing different materials. Those for a core of soft carbon and for a polymetallic core are shown in Figs. 24 and 25. Because radiation emitted from the carbon arc is very intense, this arc supplants, for many applications, sources which radiate at lower temperatures. Among the disadvantages in using the carbon arc are its inconvenience relative to the use of other sources (e.g., lamps) and its relative instability. However, Null and Lozier²³ have studied the properties of the low-intensity carbon arc extensively and have found that under the proper operating conditions the carbon arc can be made quite stable; in fact, in their treatise they recommend its use as a standard of radiation at high temperatures.

Enclosed Arc and Discharge Sources (High-Pressure)

Koller²² states that the carbon arc is generally desired if a high intensity is required from a single unit but that it is less efficient than the mercury arc. Other disadvantages are the short life of the carbon with respect to mercury, and combustion products which may be undesirable. Worthing² describes a number of the older, enclosed, metallic arc sources, many of which can be built in the laboratory for laboratory use. Today, however, it is rarely necessary to build one's own source unless it is highly specialized.

*The Infrared Handbook*¹ compiles a large number of these types of sources, some of which will be repeated here, in case that publication would not be currently available to the reader. However, the reader should take caution that many changes might have occurred in the characteristics of these sources and in the supplier whose product is preferred. Consultation with the Photonics Directory (see preceding) is usually a good procedure. In some cases a certain type of source described previously

TABLE 2 DC Carbon Arcs

	Low Intensity			Nonrotating High Intensity			Rotating High Intensity			
	1	2	3	4	5	6	7	8	9	10
Type of carbon	Microscope	Projector	Projector	Projector	Projector	Projector	Projector	Searchlight	Studio	
Positive carbon:										
Diameter (mm)	5	7	8	10	11	13.6	13.6	16	16	16
Length (in.)	8	12-14	12-14	20	20	22	22	22	22	22-30
Negative carbon:										
Diameter	6 mm	6 mm	7 mm	11/32 in.	3/8 in.	0.5 in.	0.5 in.	11 mm	17/32 in.	7/16 in.
Length (in.)	4.5	9	9	9	9	9	9	12	9	12-48
Arc current (A)	5	50	70	105	120	160	180	150	225	400
Arc volts (dc)	59	40	42	59	57	66	74	78	70	80
Arc power (W)	295	2,000	2,940	6,200	6,840	10,600	13,300	11,700	15,800	32,000
Burning rate (in. h ⁻¹)										
Positive carbon	4.5	11.6	13.6	21.5	16.5	17	21.5	8.9	20.2	55
Negative carbon	2.1	4.3	4.3	2.9	2.4	2.2	2.5	3.9	2.2	3.5
Approximate crater diameter (in.)	0.12	0.23	0.28	0.36	0.39	0.5	0.5	0.55	0.59	0.59
Maximum luminance of crater (cd cm ⁻²)	15,000	55,000	83,000	90,000	85,000	96,000	95,000	65,000	68,000	45,000
Forward crater candlepower	975	10,500	22,000	36,000	44,000	63,000	78,000	68,000	99,000	185,000
Crater lumens [†]	3,100	36,800	77,000	126,000	154,000	221,000	273,000	250,000	347,000	660,000
Total lumens [‡]	3,100	55,000	115,000	189,000	231,000	368,000	410,000	374,000	521,000	999,000
Total lumens per arc watt	10.4	29.7	39.1	30.5	33.8	34.7	30.8	32	33	30.9
Color temperature (K) [§]	3,600	5,950	5,500-6,500	5,500-6,500	5,500-6,500	5,500-6,500	5,500-6,500	5,400	4,100	5,800-6,100

[†]Typical applications: 1, microscope illumination and projection; 2 to 7, motion-picture projection; 8, searchlight projection; 9, motion-picture-set lighting and motion-picture and television background projection.

[‡]Includes light radiated in forward hemisphere.

[§]Includes light from crater and arc flame in forward hemisphere.

[¶]Crater radiation only.

TABLE 3 Flame-Type Carbon Arcs

Type of carbon Flame materials Burning position ^g Upper carbon ^d	Application Number ^g									
	1	3	3	4	5	6	7 ^b	8 ^{c,d}	9 ^{d,e}	10 ^f
C										
Polymetallic Vertical										
Diameter	22 mm	22 mm	22 mm	22 mm	22 mm	1/2 in.	1/2 in.	6 mm	9 mm	8 mm
Length (in.)	12	12	12	12	12	3-16	12	6.5	8	12
Lower carbon ^d										
Diameter	13 mm	13 mm	13 mm	13 mm	13 mm	1/2 in.	1/2 in.	6 mm	9 mm	7 mm
Length (in.)	12	12	12	12	12	3-16	12	6.5	8	9
Arc current (A)	60	60	60	80	80	16	38	40	95	40
Arc voltage (ac) ^h	50	50	50	50	50	138	50	24	30	37 dc
Arc power (kW)	3	3	3	4	4	2.2	1.9	1	2.85	1.5
Candlepower ⁱ	2,100	6,300	9,100	10,000	8,400	1,170	6,700	4,830	14,200	11,000
Lumens	23,000	69,000	100,000	110,000	92,000	13,000	74,000	53,000	156,000	110,000
Lumens per arc watt	7.6	23	33.3	27.5	23	5.9	39.8	53	54.8	73.5
Color temperature (K)			12,800 ^j	24,000 ^j			7,420 ^j	6,590	8,150	4,700
Spectral intensity ($\mu\text{W cm}^{-2}$)										
1 m from arc axis:										
Below 270 nm	540.0	180.0	102	140	1,020		95	11	100	12
270-320 nm	540.0	150.0	186	244	1,860		76	49	100	48
320-400 nm	1,800.0	1,200.0	2,046	2,816	3,120	1,700	684	415	1,590	464
400-450 nm	300.0	1,100.0	1,704	2,306	1,480	177	722	405	844	726
450-700 nm	600.0	4,050.0	3,210	3,520	2,600	442	2,223	1,602	3,671	3,965
700-1125 nm	1,580.0	2,480.0	3,032	3,500	3,220	1,681	1,264	1,368	5,632	2,123
Above 1125 nm	9,480.0	10,290.0	9,820	11,420	14,500	6,600	5,189	3,290	8,763	4,593
Total	14,930	19,460	20,100	24,000	27,800	10,600	10,253	7,140	20,600	11,930

Spectral radiation (percent of input power):										
Below 270 nm	1.8	0.6	0.34	0.35	2.55	0.5	0.11	0.32	0.08	
270–320 nm	1.8	0.5	0.62	0.61	4.65	0.4	0.49	0.35		
320–400 nm	6.0	4.0	6.82	7.04	7.80	7.7	4.15	5.59	3.09	
400–450 nm	1.3	3.7	5.68	5.90	3.70	0.8	4.05	2.96	4.84	
450–700 nm	2.0	13.5	10.7	8.80	6.50	11.7	16.02	12.86	26.43	
700–1125 nm	5.27	8.27	10.1	8.75	8.05	7.6	13.68	10.75	14.15	
Above 1125 nm	31.6	34.3	32.7	28.55	36.25	29.9	32.90	30.60	30.62	
Total	49.77	64.87	67.00	60.00	69.50	48.0	71.40	72.20	79.53	

^a Typical applications: 1 to 5 and 8, photochemical, therapeutic, accelerated exposure testing, or accelerated plant growth; 6, 7, and 9 blue-printing diazo printing, photo copying, and graphic arts; 10, motion-picture and television studio lighting.

^b Photographic white-flame carbons.

^c High-intensity copper-coated sunshine carbons.

^d Both carbons are same in horizontal, coaxial arcs.

^e High-intensity photo carbons.

^f Motion-picture-studio carbons

^g All combinations shown are operated coaxially.

^h All operated on alternating current except item 10.

ⁱ Horizontal candlepower, transverse to arc axis.

^j Deviates enough from blackbody colors to make color temperature of doubtful meaning.

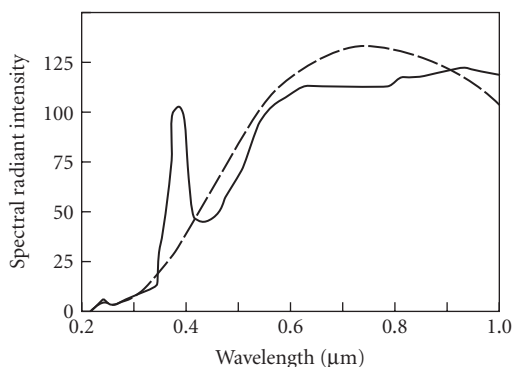


FIGURE 23 Spectral distribution of radiant flux from 30-A, 55-V dc low-intensity arc with 12-mm positive carbon (solid line) and a 3800-K blackbody radiator (broken line).

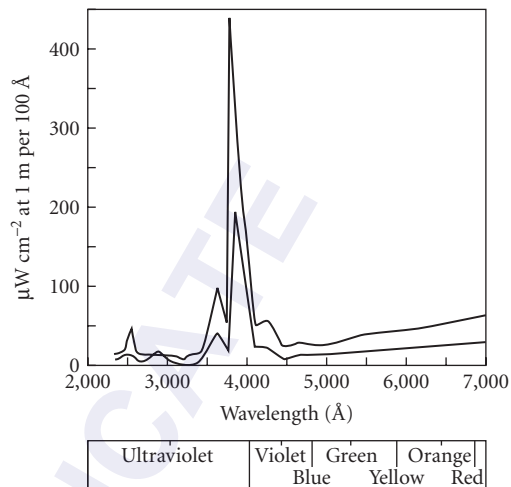


FIGURE 24 Spectral energy distribution of carbon arc with core of soft carbon. Upper curve: 60-A ac 50-V across the arc; lower curve: 30-A ac 50-V across the arc.

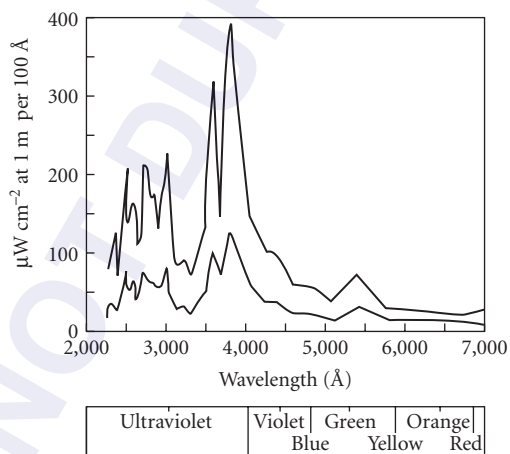


FIGURE 25 Spectral energy distribution of carbon arc with polymetallic-cored carbons. Upper curve: 60-A ac 50-V across the arc; lower curve: 30-A ac 50-V across the arc.

may not still exist. Thus, whereas some manufacturers were less compliant in providing data, they should be expected to respond more readily to a potential customer.

*Uviarc** This lamp is an efficient radiator of ultraviolet radiation. The energy distribution of one type is given in Fig. 26. Since the pressure of this mercury-vapor lamp is intermediate between the usual high- and the low-pressure lamps, little background (or continuum) radiation

*Registered trademark of General Electric.

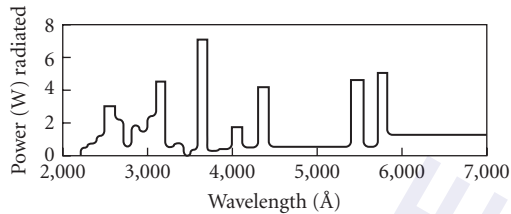


FIGURE 26 Intensity distribution of UA-2 intermediate-pressure lamp.

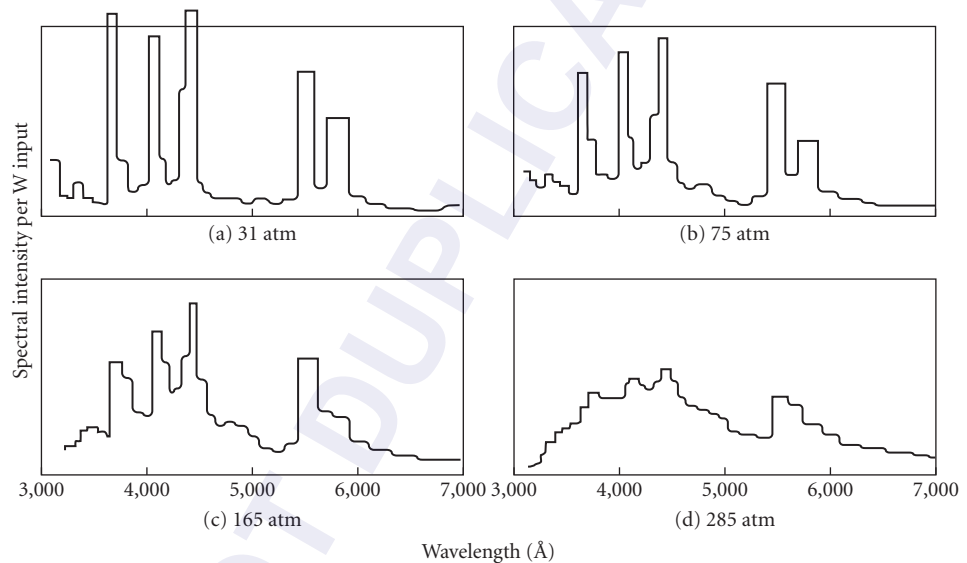


FIGURE 27 Emission spectrum of high-pressure mercury-arc lamps showing continuum background.

is present. In the truly high-pressure lamp, considerable continuum radiation results from greater molecular interaction. Figure 27²⁴ shows the dependence on pressure of the amount of continuum in mercury lamps of differing pressure. Bulb shapes and sizes are shown in Fig. 28.

Mercury Arcs A widely used type of high-pressure, mercury-arc lamp and the components necessary for its successful operation are shown in Fig. 29. The coiled tungsten cathode is coated with a rare-earth material (e.g., thorium). The auxiliary electrode is used to help in starting. A high resistance limits the starting current. Once the arc is started, the operating current is limited by ballast supplied by the high reactance of the power transformer. Spectral data for clear, 400-W mercury lamps of this type are given in Fig. 30.

Multivapor Arcs In these lamps, argon and mercury provide the starting action. Then sodium iodide, thallium iodide, and indium iodide vaporize and dissociate to yield the bulk of the lamp radiation. The physical appearance is like that of mercury lamps of the same general nature. Ballasts are similar to their counterparts for the mercury lamp. Up-to-date information on these sources should be obtained from the General Electric Corporation Lamp Division in Nela Park near Cleveland, Ohio. Spectral features of these sources are given in Fig. 31.

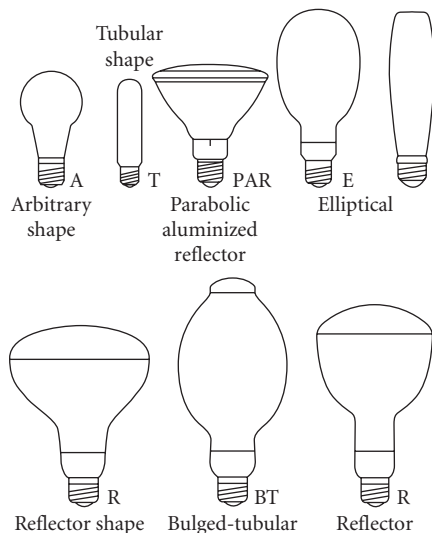


FIGURE 28 Bulb shapes and sizes (not to scale).

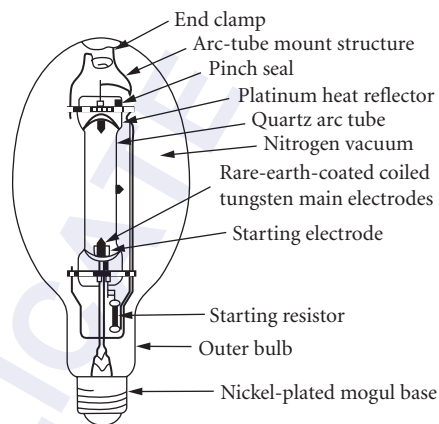


FIGURE 29 High-pressure mercury lamp showing various components.

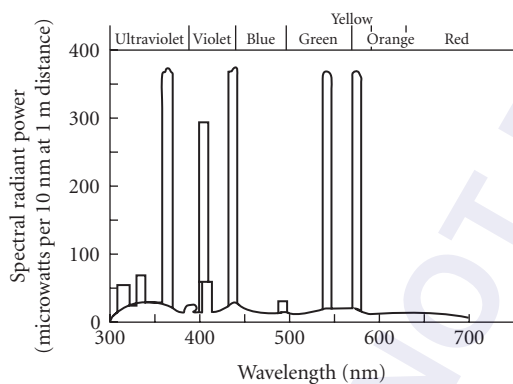


FIGURE 30 Spectral energy distribution for clear mercury-arc lamp.

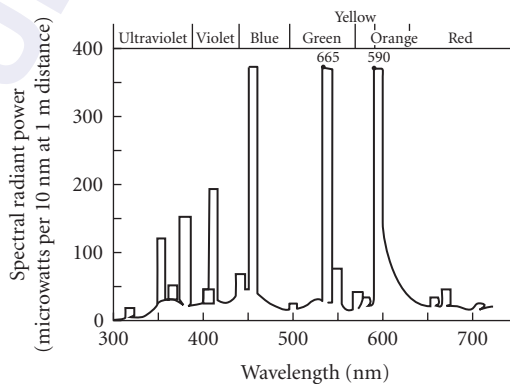


FIGURE 31 Spectral energy distribution of multivapor-arc lamp.

Lucalox[®] Lamps The chief characteristics of this lamp are high-pressure sodium discharge and a high temperature withstanding ceramic, Lucalox (translucent aluminum oxide), to yield performance typified in the spectral output of the 400-W Lucalox lamp shown in Fig. 32. Ballasts for this lamp are described in the General Electric *Bulletin TP-109R*.²⁵

Capillary Mercury-Arc Lamps²² As the pressure of the arc increases, cooling is required to avoid catastrophic effects on the tube. The AH6 tube (Fig. 33) is constructed with a quartz bulb wall and a quartz outer jacket, to allow 2800 K radiation to pass, or a Pyrex^{®†} outer jacket to eliminate ultraviolet. Pure water is forced through at a rapid rate, while the tube is maintained at a potential of 840 V.

[®]Registered trademark of General Electric.

[†]Registered trademark of Corning Glass Works.

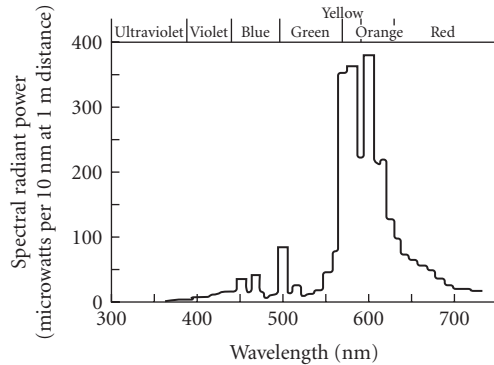


FIGURE 32 Spectral output of 400-W Lucalox lamp.

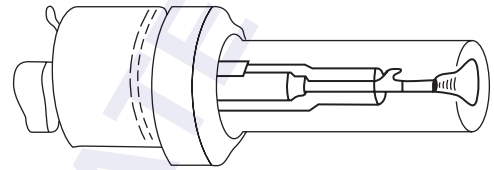


FIGURE 33 Water-cooled high-pressure (110 atm) mercury-arc lamp showing lamp in water jacket.

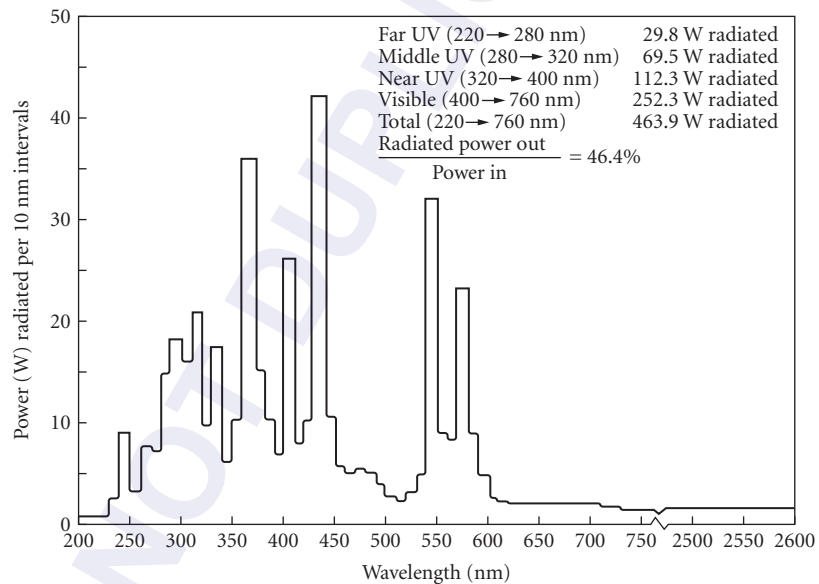


FIGURE 34 Spectral energy distribution of type BH6-1 mercury capillary lamp.

The spectral characteristics of certain tubes²⁶ are shown in Fig. 34. This company does not appear in the *Photonics Guide of 1989*, so the catalog referenced in the figure may not be current.

Compact-Source Arcs^{19,27} Some common characteristics of currently available compact-source arc lamps are as follows:

1. A clear quartz bulb of roughly spherical shape with extensions at opposite ends constituting the electrode terminals. In some cases, the quartz bulb is then sealed within a larger glass bulb, which is filled with an inert gas.
2. A pair of electrodes with relatively close spacing (from less than 1 mm to about 1 cm); hence the sometimes-used term short-arc lamps.

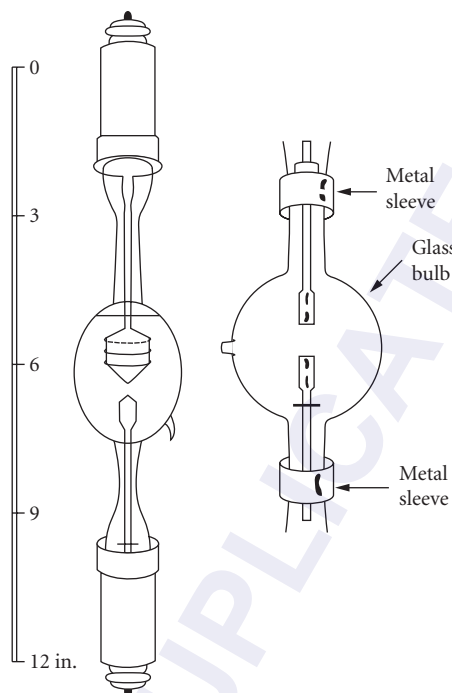


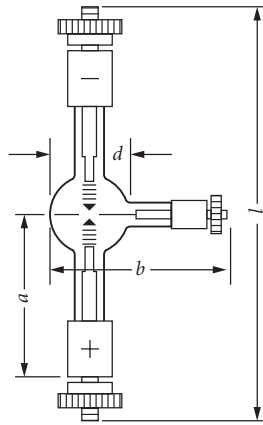
FIGURE 35 Construction of different lamps showing differences in relative sizes of electrodes for dc (*left*) and ac (*right*) operation.

3. A filling of gas or vapor through which the arc discharge takes place.
4. Extreme electrical loading of the arc gap, which results in very high luminance, internal pressures of many atmospheres, and bulb temperatures as high as 900°C . Precautions are necessary to protect people and equipment in case the lamps should fail violently.
5. The need for a momentary high-voltage ignition pulse, and a ballast or other auxiliary equipment to limit current during operation.
6. Clean, attention-free operation for long periods of time.

These lamps are designated by the chief radiating gases enclosed as mercury, mercury xenon, and xenon lamps.

Figure 35 shows a compact-source construction for a 1000-W lamp. Since starting may be a problem, some lamps (Fig. 36) are constructed with a third (i.e., a starting) electrode, to which a momentary high voltage is applied for starting (and especially restarting) while hot. The usual ballast is required for compact-source arcs. For stability, these arcs, particularly mercury and mercury-xenon, should be operated near rated power on a well-regulated power supply.²⁷

The spatial distribution of luminance from these lamps is reported in the literature already cited, and typical contours are shown in Fig. 37. Polar distributions are similar to those shown in Fig. 38. Spectral distributions are given in Figs. 39 through 41 for a 1000-W ac mercury lamp, a 5-kW dc xenon lamp, and 1000-W dc mercury-xenon lamp. Lamps are available at considerably less wattage.



Lamp (order reference)	HBO 200	
Type of current	DC	AC
Lamp supply voltage	V	>105 220
Operating voltage of lamp	$V \frac{L_1}{L_2}$	65...47 $\frac{61 \pm 4}{53 \pm 4}$
Operating current at operating voltage range	$A \frac{L_1}{L_2}$	3.1...4.2 $\frac{3.6}{4.2}$
Rated power of lamp	W	200
Luminous flux	lm	9,500
Luminous efficacy	lm/W	47.5
Light intensity	cd	1,000
Average luminance	cd / cm ²	40,000
Arc (width $w \times$ height h^4)	mm	0.6 \times 2.2
Average lamp life	h	200
Diameter d	mm	18
Length l_{max}	mm	108
Distance a	mm	41 \pm 2
Width b	mm	45
Burning position with stamped base down		s^{45}

FIGURE 36 Construction of a lamp with a third, starting electrode.

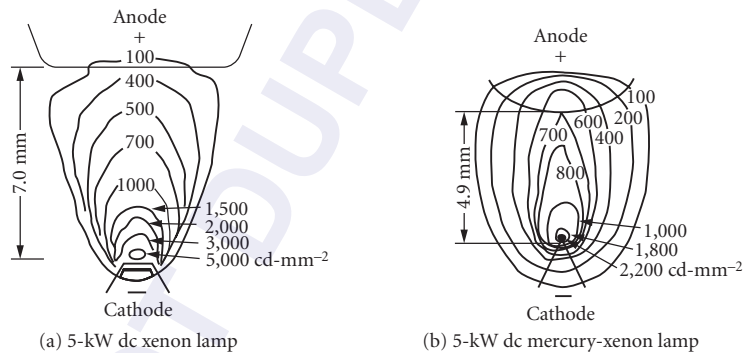


FIGURE 37 Spatial luminance distribution of compact-arc lamps.

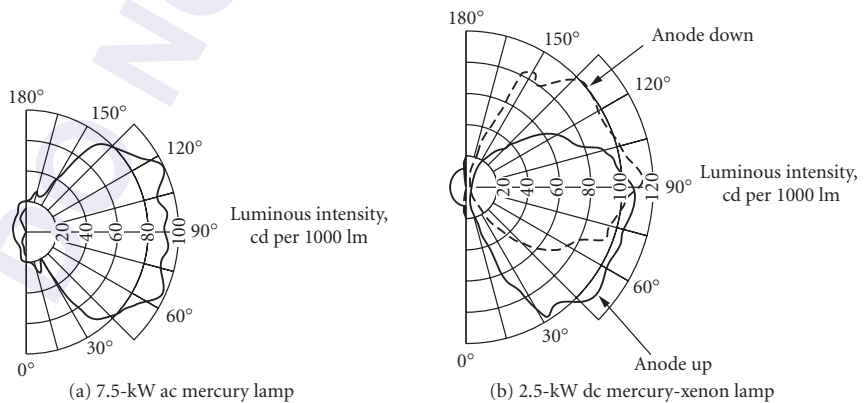


FIGURE 38 Polar distribution of radiation in planes that include arc axis. Asymmetry in (b) is due to unequal size of electrodes.

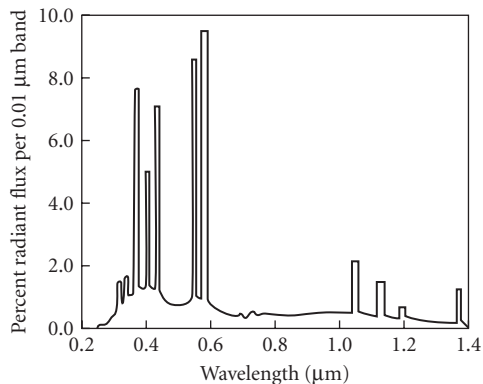


FIGURE 39 Spectral distribution of radiant intensity from a 1000-W ac mercury lamp perpendicular to the lamp axis.

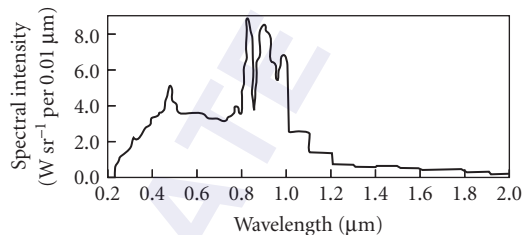


FIGURE 40 Spectral distribution of radiant intensity from a 5-kW dc xenon lamp perpendicular to the lamp axis with electrode and bulb radiation excluded.

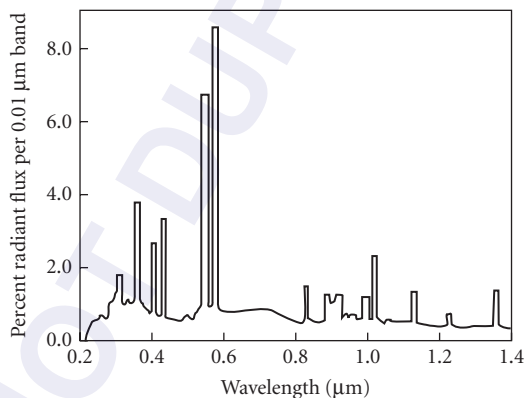


FIGURE 41 Spectral distribution of radiant flux from a 1000-W mercury-xenon lamp.

Cann²⁷ reports on some interesting special lamps tested by Jet Propulsion Laboratories for the purpose of obtaining a good spectral match to the solar distribution. The types of lamps tested were Xe, Xe-Zn, Xe-Cd, Hg-Xe-Zn, Hg-Xe-Cd, Kr, Kr-Zn, Kr-Cd, Hg-Kr-Zn, Hg-Kr-Cd, Ar, Ne, and Hg-Xe with variable mercury-vapor pressure. For details, the reader should consult the literature.

A special design of a short-arc lamp manufactured by Varian²⁸ is shown in Fig. 42. Aside from its compactness and parabolic sector, it has a sapphire window which allows a greater amount of IR energy to be emitted. It is operated either dc or pulsed, but the user should obtain complete specifications, because the reflector can become contaminated, with a resultant decrease in output.



FIGURE 42 High-pressure, short-arc xenon illuminators with sapphire windows. Low starting voltage, 150 through 800 W; VIX150, VIX300, VIX500, VIX800.

Enclosed Arc and Discharge Sources (Low-Pressure)²²

With pressure reduction in a tube filled with mercury vapor, the 2537 Å line becomes predominant so that low-pressure mercury tubes are usually selected for their ability to emit ultraviolet radiation.

Germicidal Lamps These are hot-cathode lamps which operate at relatively low voltages. They differ from ordinary fluorescent lamps which are used in lighting in that they are designed to transmit ultraviolet, whereas the wall of the fluorescent lamp is coated with a material that absorbs ultraviolet and reemits visible light. The germicidal lamp is constructed of glass of 1-mm thickness which transmits about 65 percent of the 2537 Å radiation and virtually cuts off shorter wavelength ultraviolet radiation.

Sterilamp[®] Types These cold cathode lamps start and operate at higher voltages than the hot-cathode type and can be obtained in relatively small sizes as shown in Fig. 43. Operating characteristics of the Sterilamps should be obtained from the manufacturer.

Black-Light Fluorescent Lamps This fluorescent lamp is coated with a phosphor efficient in the absorption of 2537Å radiation, emitting ultraviolet radiation in a broadband around 3650 Å. The phosphor is a cerium-activated calcium phosphate, and the glass bulb is impervious to shorter wavelength ultraviolet radiation. Characteristics of one type are given in Table 4.

Hollow Cathode Lamps A device described early in this century and used for many years by spectroscopists is the hollow-cathode tube. The one used by Paschen² consisted of a hollow metal cylinder and contained a small quantity of inert gas, yielding an intense cathode-glow characteristic of the cathode constituents. Materials that vaporize easily can be incorporated into the tube so that their spectral characteristics predominate.²

Several companies sell hollow-cathode lamps which do not differ significantly from those constructed in early laboratories. The external appearance of these modern tubes shows the marks of mass production and emphasis on convenience. They come with a large number of vaporizable elements, singly or in multiples, and with Pyrex[®] or quartz windows. A partial list of the characteristics or the lamps available from two manufacturers is given in Table 5. Their physical appearance is shown in Fig. 44a. A schematic of the different elements obtainable in various lamps is shown in Fig. 44b.²⁹

²⁹Registered trademark of Westinghouse Electric.

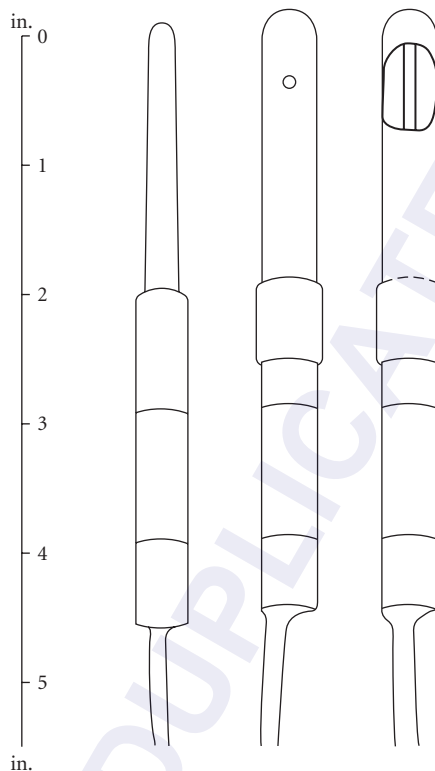


FIGURE 43 Pen-Ray low-pressure lamp. Pen-Ray is a registered trademark of Ultraviolet Products, Inc.

TABLE 4 Spectral Energy Distribution for Black-Light (360 BL) Lamps

(W)	Length (in.)	3200–3800 Å		Total Ultraviolet below 3800 Å		Total Visible (W)	3800–7600 Å %*	Erythemat Flux
		(W)	%*	(W)	%*			
6	9	0.55	9.1	0.56	9.4	0.1	1.7	250
15	18	2.10	14.0	2.20	14.6	0.4	2.7	950
30	36	4.60	15.3	4.70	15.8	0.9	3.0	2100
40	48	6.70	16.8	6.90	17.3	1.5	3.8	3000

* Percentage of input power.

Electrodeless Discharge Lamps^{30–32} The electrodeless lamp gained popularity when Meggers used it in his attempt to produce a highly precise standard of radiation. Simplicity of design makes laboratory construction of this type of lamp easy. Some of the simplest lamps consist of a tube, containing the radiation-producing element, and a microwave generator, for producing the electric field (within the tube) which in turn excites the elemental spectra. Lamps of this type can be purchased with specially designed microwave cavities for greater efficiency in coupling. Those made of fused quartz can transmit from ultraviolet to near infrared. The electrodeless lamp is better able than the

TABLE 5 Single-Element and Multiple-Element Hollow-Cathode Lamps*

Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ	Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ
Aluminum	Q	A	B	3092	WL22804	Copper	P	N	A	3247	JA-45-458
	P	N	A	3092	JA-45-452		Q	A	B	3247	WL22606
	Q	A	A	3092	WL22870		Q	A	A	3247	WL22879
Antimony	Q	N	B	3092	WL22929	Dysprosium	Q	N	A	3247	JA-45-490
	Q	N	A	3092	WL22954		Q	N	A	3247	WL23042
	Q	N	A	3092	WL22840		Q	N	B	4212	JA-45-595
	Q	A	B	2311	WL22840		Q	N	A	4212	WL22880
Arsenic	Q	A	A	2311	WL22872	Erbium	Q	N	B	4008	JA-45-571
	Q	N	B	2311	JA-45-461		Q	N	A	4008	WL22881
	Q	N	A	2311	WL22956		Q	N	B	4594	JA-45-572
Barium	Q	N	A	1937	JA-45-315	Europium	Q	N	A	4594	WL22882
	Q	N	B	1937	WL22873		Q	N	A	4079	WL22975
	Q	N	A	1937	JA-45-315		Q	N	B	4079	JA-45-573
Beryllium	P	N	A	5536	JA-45-480	Gadolinium	Q	N	A	4079	WL22986
	Q	N	B	2349	WL23407		Q	N	A	4172	JA-45-470
	Q	A	B	3068	WL22841		Q	N	A	4172	WL22884
Bismuth	Q	A	A	3068	WL22874	Germanium	Q	N	B	2651	JA-45-575
	Q	N	B	3068	JA-45-469		Q	N	A	2651	JA-45-313
	Q	N	A	3068	WL22957		Q	A	B	2676	WL22839
Boron	Q	A	B	2497	JA-45-568	Gold	Q	A	B	2676	WL22883
	Q	A	A	2497	WL22917		Q	A	A	2676	JA-45-467
	Q	A	A	2497	WL22816		Q	N	A	2676	WL22960
Cadmium	Q	A	B	3261	WL22816	Hafnium	Q	N	B	3072	JA-45-303
	Q	A	A	3261	WL22875		Q	N	A	4104	WL22885
	Q	N	B	3261	JA-45-462		Q	N	B	4104	JA-45-576
Calcium	Q	N	A	3261	WL22958	Indium	Q	N	B	3040	WL22867
	P	N	A	4227	JA-45-440		Q	A	B	3040	WL22915
	Q	N	B	-	JA-45-569		Q	A	A	3040	JA-45-471
Cerium	Q	N	A	4556	WL22978	Iridium	Q	N	B	2850	JA-45-577
	P	A	A	4556	WL22817		P	A	A	3270	WL22602
	P	N	A	4566	JA-45-141		Q	A	B	3720	WL22611

(Continued)

TABLE 5 Single-Element and Multiple-Element Hollow-Cathode Lamps* (Continued)

Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ	Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ
Chromium	P	A	A	3579	WL22812	Iron, high-purity	Q	N	B	3720	JA-45-155
	Q	A	B	3579	WL22521		P	N	A	3720	WL22820
	Q	A	A	3579	WL22877		Q	A	A	3720	WL22886
Cobalt	Q	N	B	3579	JA-45-454	Lanthanum	Q	N	A	3720	WL22887
	Q	N	A	3579	WL22959		Q	N	B	3720	WL22837
	P	A	A	3454	WL22813		Q	N	A	3720	WL22888
	Q	A	B	3454	WL22814		Q	A	B	5501	WL22846
	Q	A	A	3454	WL22878		Q	A	A	5501	WL22889
	Q	N	B	3454	JA-45-456		Q	N	B	5501	JA-45-495
Lead	Q	N	A	3454	WL22953	Palladium	Q	A	B	3404	WL22857
	Q	A	B	2833	WL22838		Q	A	A	3404	WL22911
	Q	A	A	2833	WL22890		Q	N	B	3404	JA-45-475
Lithium 6	Q	N	B	2833	JA-45-468	Phosphorus	Q	N	A	3404	WL22970
	Q	N	A	2833	WL22952		Q	N	B	2136	JA-45-449
	P	N	A	6708	JA-45-579		Q	N	A	2136	WL22990
Lithium 7	P	N	A	6708	WL22925	Platinum	Q	N	A	2659	WL22851
	P	A	A	6708	JA-45-580		Q	A	A	2659	WL22896
	P	A	A	6708	WL22926		Q	N	B	2659	JA-45-466
Lithium, natural	P	A	A	6708	WL22825	Potassium	Q	N	B	4044	JA-45-484
	P	N	A	6708	JA-45-444		Q	N	A	4951	JA-45-585
	Q	A	B	6708	WL23115		Q	N	A	4951	WL22982
Lutetium	Q	N	B	3282	JA-45-581	Praseodymium	Q	N	B	3460	JA-45-489
	Q	N	A	3282	WL23010		Q	N	B	3460	WL22967
	Q	A	B	2852	WL22609		Q	N	A	3435	WL22850
Magnesium	Q	A	A	2852	WL22891	Rhodium	Q	A	B	3435	WL22850
	Q	N	A	2852	WL22951		Q	A	A	3435	WL22897
	Q	N	B	2852	JA-45-451		Q	N	B	3435	JA-45-476
Manganese	Q	A	B	2795	WL22608	Rubidium	P	N	A	7800	JA-45-443
	P	A	A	2795	WL22815		Q	N	B	7800	WL23046
	Q	N	B	2795	JA-45-472		Q	A	B	3499	JA-45-586
	Q	N	A	2795	WL22961	Ruthenium	Q	N	B	4760	JA-45-587
	Q	A	A	2795	WL22876		Q	N	A	4760	WL22899
	Q	A	A	2795							

Mercury	Q	A	B	2537	JA-45-493	Scandium	Q	N	B	3912	JA-45-309
	Q	A	A	2537	WL22892	Selenium	Q	A	B	1960	WL22843
Molybdenum	Q	A	B	3133	WL22805		Q	A	A	1960	WL22898
	Q	A	A	3133	WL22893		Q	N	B	1960	JA-45-477
	Q	N	B	3133	JA-45-460		Q	N	A	1960	WL22963
Neodymium	Q	N	A	3133	WL22962	Silicon	Q	A	B	2516	WL22832
	Q	N	B	4925	JA-45-582		Q	A	A	2516	WL22900
	Q	N	A	4925	WL22980		Q	N	B	2516	JA-45-479
Nickel	P	A	A	3415	WL22605		Q	N	A	2516	WL22964
	Q	A	B	3415	WL22663	Silver	Q	A	B	3281	JA-45-483
	Q	N	B	3415	JA-45-457		Q	A	A	3281	WL22901
	Q	A	A	3415	WL22894	Sodium	P	A	A	5890	WL22864
Niobium	Q	N	A	3415	WL22895		P	N	A	5890	JA-45-485
	Q	N	B	4059	JA-45-486	Strontium	P	N	A	4607	JA-45-481
Osmium	Q	N	A	4059	WL22912	Sulphur	Q	N	B	-	JA-45-588
Tantalum	Q	A	B	2909	JA-45-584						
	Q	A	B	2714	JA-45-488	Zirconium	Q	A	B	3601	JA-45-482
	Q	A	A	2714	WL22913		Q	A	A	3601	WL22914
	Q	N	B	2714	WL22971		Q	N	B	3601	WL22998
	Q	N	A	2714	WL22972						
Tellurium	Q	A	B	2143	WL22842						
	Q	A	A	2143	WL22902						
	Q	N	B	2143	JA-45-473	Aluminum	P	N	C	3092	JA-45-36009
	Q	N	A	2143	WL22965	Antimony	Q	N	C	2311	JA-4-36010
Terbium	Q	N	B	4326	JA-45-589	Arsenic	Q	A	C	1937	JA-45-36011
	Q	N	A	4326	WL22903	Barium	P	N	C	5536	JA-45-36012
Thallium	Q	N	B	3776	WL23408	Beryllium	Q	N	C	2349	JA-45-36013
Thorium	Q	N	A	3245	WL23028	Bismuth	Q	N	C	3068	JA-45-36014
	Q	N	B	3245	JA-45-590	Boron	Q	A	C	2497	JA-45-36015
Thulium	Q	N	B	4105	JA-45-591	Cadmium	Q	N	C	3261	JA-45-36016
	Q	N	A	4105	WL23008	Calcium	P	N	C	4227	JA-45-36017
Tin	Q	A	B	2863	WL22822	Cerium	Q	N	C	-	JA-45-36019
	Q	A	A	2863	WL22904	Cesium	P	N	C	4556	JA-45-36020
	Q	N	B	2863	JA-45-463	Chromium	P	N	C	3579	JA-45-36021
	Q	N	A	2863	WL22966	Cobalt	Q	N	C	3454	JA-45-36022
Titanium	Q	N	B	3643	JA-45-592	Copper	P	N	C	3247	JA-45-36024
	Q	N	A	3643	WL22992	Dysprosium	P	N	C	4212	JA-45-36025

Single-Element 36000 Series

(Continued)

TABLE 5 Single-Element and Multiple-Element Hollow-Cathode Lamps* (Continued)

Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ	Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ
Tungsten	Q	N	B	4009	JA-45-465	Erbium	P	N	C	4008	JA-45-36026
	Q	A	B	4009	WL22849	Europium	P	N	C	4594	JA-45-36027
	Q	N	A	4009	WL22905	Gadolinium	P	N	C	4079	JA-45-36028
Uranium	Q	A	A	4009	WL22906	Gallium	Q	N	C	4172	JA-45-36029
	Q	N	B	5027	JA-45-447	Germanium	Q	N	C	2651	JA-45-36030
	Q	N	A	5027	WL22907	Gold	Q	N	C	2676	JA-45-36031
Vanadium	Q	A	B	3184	WL22856	Hafnium	Q	N	C	3072	JA-45-36032
	Q	A	A	3184	WL22910	Holmium	P	N	C	4104	JA-45-36033
	Q	N	B	3184	JA-45-453	Indium	Q	N	C	3040	JA-45-36034
Ytterbium	Q	N	A	3184	WL22974	Iridium	Q	N	C	2850	JA-45-36036
	Q	A	B	3988	JA-45-593	Iron	Q	N	C	3720	JA-45-36037
	Q	A	A	3988	WL22984	Lanthanum	P	N	C	5501	JA-45-36038
Yttrium	Q	N	A	4102	WL22976	Lead	Q	N	C	2833	JA-45-36039
	Q	N	B	4102	JA-45-594	Lithium 6	P	N	C	6708	JA-45-36090
	Q	N	A	4102	WL22988	Lithium 7	P	N	C	6708	JA-45-36091
Zinc	Q	A	B	2139	WL22607	Lithium, natural	P	N	C	6708	JA-45-36040
	Q	N	B	2139	JA-45-459	Luettium	P	N	C	3282	JA-45-36041
	Q	A	A	2139	WL22908	Magnesium	Q	N	C	2852	JA-45-36042
Mercury	Q	N	A	2139	WL22909	Manganese	Q	N	C	2795	JA-45-36043
	Q	A	C	2537	JA-45-36044	Multiple-Element 22000 Series					
	Q	N	C	3133	JA-45-36045	Aluminum, calcium	Q	N	B	-	WL23246
Neodymium	P	N	C	4925	JA-45-36046	Aluminum, calcium, magnesium	Q	A	B	-	WL22604
Nickel	Q	N	C	3415	JA-45-36047						
Niobium	P	N	C	4059	JA-45-36023	Aluminum, calcium, magnesium	Q	A	A	-	WL22871
Osmium	Q	A	C	2909	JA-45-36048						
Palladium	Q	N	C	3404	JA-45-36049	Aluminum, calcium, magnesium	Q	N	B	-	JA-45-450
Phosphorus	Q	N	C	2136	JA-45-36050						
Platinum	Q	N	C	2659	JA-45-36051	Aluminum, calcium, magnesium	Q	N	A	-	WL22955
Potassium	P	N	C	4044	JA-45-36052						

Praseodymium	P	N	C	4951	JA-45-36053	Aluminum, calcium, magnesium, iron	Q	N	B	-	JA-45-310
Rhenium	P	N	C	3460	JA-45-36056						
Rhodium	P	N	C	3435	JA-45-36057						
Rubidium	P	N	C	7800	JA-45-36058	Aluminum, calcium, magnesium, lithium	Q	N	B	-	JA-45-436
Ruthenium	P	A	C	3499	JA-45-36059						
Samarium	P	N	C	4760	JA-45-36060	Aluminum, calcium, magnesium, lithium	Q	A	A	-	WL23036
Scandium	P	N	C	3912	JA-45-36061						
Selenium	Q	N	C	1960	JA-45-36062	Aluminum, calcium, strontium	P	N	A	-	WL23403
Silicon	Q	N	C	2516	JA-45-36063						
Silver	P	A	C	3281	JA-45-36064	Antimony, arsenic, bismuth	Q	N	B	-	WL23147
Sodium	P	N	C	5890	JA-45-36065						
Strontium	P	N	C	4607	JA-45-36066	Arsenic, nickel	Q	N	B	-	JA-45-434
Sulphur	Q	N	C	-	JA-45-36067	Arsenic, selenium, tellurium	Q	N	B	-	JA-45-598
Tantalum	Q	A	C	2714	JA-45-36068	Barium, calcium, strontium	P	N	A	-	JA-45-437
Tellurium	Q	N	C	2143	JA-45-36069	Barium, calcium, silicon, magnesium	Q	N	B	-	JA-45-478
Terbium	P	N	C	4326	JA-45-36070						
Thallium	Q	N	C	3776	JA-45-36071	Cadmium, copper	Q	N	B	-	JA-45-597
Thorium	Q	N	C	3245	JA-45-36072,	zinc, lead	Q	N	B	-	JA-45-308
Thulium	P	N	C	4105	JA-45-36073	Cadmium, silver, zinc, lead	Q	N	B	-	WL23605
Tin	Q	N	C	2863	JA-45-36074	Calcium, magnesium, strontium	Q	N	B	-	JA-45-311
Titanium	P	N	C	3643	JA-45-36075	Calcium, magnesium, zinc	Q	N	B	-	
Tungsten	Q	N	C	4009	JA-45-36076	Calcium, magnesium, aluminum, lithium	Q	N	B	-	WL23158
Uranium	P	N	C	5027	JA-45-36077	Calcium, zinc	Q	N	B	-	JA-45-304
Vanadium	Q	N	C	3184	JA-45-36078	Chromium, iron, manganese, nickel	Q	N	B	-	JA-45-442
Ytterbium	P	A	C	3988	JA-45-36079	Chromium, cobalt, nickel	Q	N	B	-	WL23174

(Continued)

TABLE 5 Single-Element and Multiple-Element Hollow-Cathode Lamps* (Continued)

Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ	Element	Window ^f	Gas Fill ^g	Size ^h	Analytical Line (Å)	Catalog Number ⁱ
Yttrium	P	N	C	4102	JA-45-36080	Chromium, copper	Q	N	B	-	JA-45-306
Zinc	Q	N	C	2139	JA-45-36081						
Zirconium	P	A	C	3601	JA-45-36082						
Chromium, manganese	Q	N	B	-	WL23499	Antimony, arsenic, bismuth	Q	N	C	-	JA-45-36203
Chromium, cobalt, copper, manganese, nickel	Q	N	B	-	WL23601	Barium, calcium, strontium, magnesium	Q	N	C	-	JA-45-36228
Chromium, cobalt, copper, iron,	Q	N	B	-	JA-45-599	Cadmium, silver, zinc, lead	Q	N	C	-	JA-45-36205
Cobalt, copper, manganese, nickel	Q	N	B	-	JA-45-599						
Cobalt, copper, gold, nickel	Q	N	B	-	JA-45-305	Cadmium, copper, zinc, lead	Q	N	C	-	JA-45-36227
Cobalt, copper, zinc, molybdenum	Q	N	B	-	WL23295	Calcium, magnesium	Q	N	C	-	JA-45-36092
Cobalt, iron	Q	N	B	-	WL23291						
Cobalt, nickel	Q	N	B	-	WL23426	Calcium, magnesium, zinc	Q	N	C	-	JA-45-36097
Copper, gallium	Q	N	B	-	JA-45-431	Calcium, zinc	Q	N	C	-	JA-45-36093
						Chromium, iron, manganese, nickel	Q	N	C	-	JA-45-36201
Copper, iron	Q	N	B	-	JA-45-312	Chromium, cobalt, copper	Q	N	C	-	JA-45-36094
Copper, iron, manganese	Q	N	B	-	JA-45-435	manganese, nickel	Q	N	C	-	JA-45-36103
Copper, iron, molybdenum	Q	N	B	-	JA-45-301	Chromium, copper, nickel, silver	Q	N	C	-	JA-45-36096
Copper, iron, gold, nickel	Q	N	B	-	JA-45-307	Chromium, copper, iron, nickel, silver	Q	N	C	-	JA-45-36108
Copper, iron, manganese, zinc	Q	N	B	-	JA-45-492	Cobalt, copper, iron, manganese, molybdenum	Q	N	C	-	JA-45-36102
Copper, manganese	Q	N	B	-	JA-45-491	Copper, zinc, lead, tin	Q	N	C	-	JA-45-36202

Copper, nickel	Q	N	B	-	WL23441A	Copper, iron	Q	N	C	JA-45-36200
Copper, nickel, zinc	Q	N	B	-	WL23405	Copper, iron, nickel	Q	N	C	JA-45-36101
Copper, zinc, molybdenum	Q	N	B	-	JA-45-496	Copper, iron, lead, nickel, zinc	Q	N	C	JA-45-36204
Copper, zinc, lead, silver	Q	N	B	-	JA-45-448	Copper, iron, manganese, zinc	Q	N	C	JA-45-36105
Copper, zinc, lead, tin	Q	N	B	-	JA-45-438	Sodium, potassium	P	A	C	JA-45-36095
Gold, nickel	Q	N	B	-	JA-4S-433					
Gold, silver	Q	N	B	-	WL23269					
Indium, silver	Q	N	B	-	WL23294					
Lead, silver, zinc	Q	N	B	-	WL23171					
Magnesium, zinc	Q	N	B	-	WL23455					
Sodium, potassium	P	N	A	-	JA-45-439					
Sodium, potassium	P	A	A	-	WL23230					
Zinc, lead, tin	Q	N	B	-	WL23404					
Multiple-Element 36000 Series										
Aluminum, calcium, magnesium	Q	N	C		JA-45-36099					
Aluminum, calcium, magnesium, lithium	Q	N	C		JA-45-36250					

*Tubes listed in this table are issued by Fisher Scientific and produced by Westinghouse Electric.

¹P = Pyrex, Q = quartz, ²N = neon, A = argon, ³A = 1 1/2-in. diameter, B = 1-in. diameter, C = 2-in. diameter, ⁴WL = Westinghouse, JA = Jarrell-Ash.



FIGURE 44a Hollow-cathode spectral tubes described in Table 5.

		Transition elements															
Li	Be											B					
Na	Mg	Group 8										Al	Si	P	S		
K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Gn	Ge	As	Se		
Rb	Sr	Y	Zr	Nb	Mo		Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te		
Cs	Ba	La	Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi			
Lanthanides	Ce	Pr	Nd		Sm	Eu	Gd	Tb	Dy	Ho	Er	Tm	Yb	Lu			
Actinides	Th		U														

FIGURE 44b Periodic table showing the prevalence of elements obtainable in hollow-cathode tubes.

arc lamp to produce stable radiation of sharp spectral lines; this makes it useful in spectroscopy and interferometry. The Hg 198 lamp makes a suitable secondary standard of radiation.

Spectral Lamps³¹ Some manufacturers produce groups of arc sources, which are similar in construction and filled with different elements and rare gases, and which yield discontinuous or monochromatic radiation throughout most of the ultraviolet and visible spectrum. They are called

spectral lamps. The envelopes of these lamps are constructed of glass or quartz, depending on the part of the spectrum desired. Thus, discrete radiation can be obtained from around 2300 Å into the near infrared. Figure 45³¹ represents the various atomic lines observable from Osram spectral lamps. Figure 46³³ gives a physical description of various spectral lamps obtainable from Philips. Table 6 lists the characteristics of the various types of lamps obtainable from Philips.

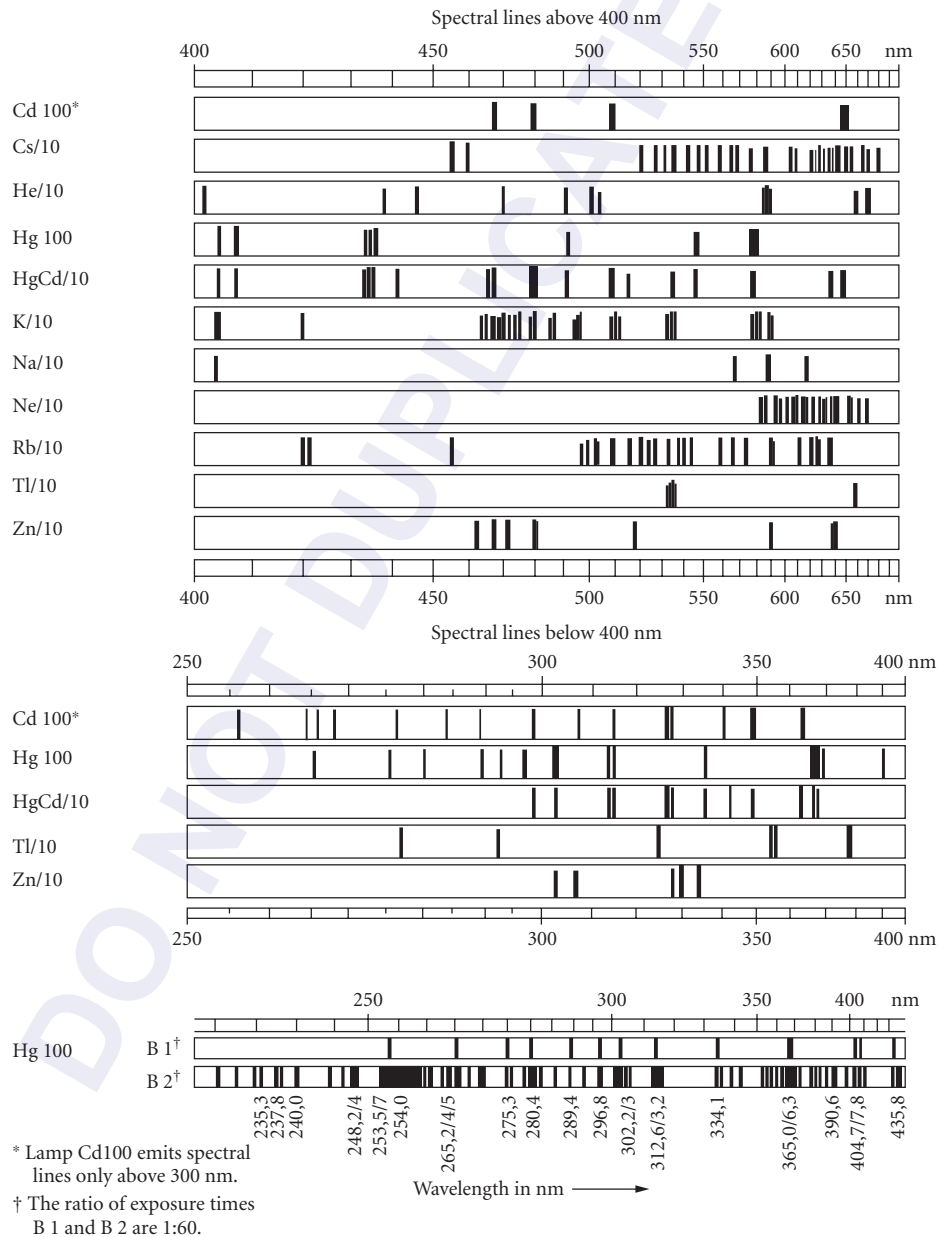


FIGURE 45 Spectral lamps.

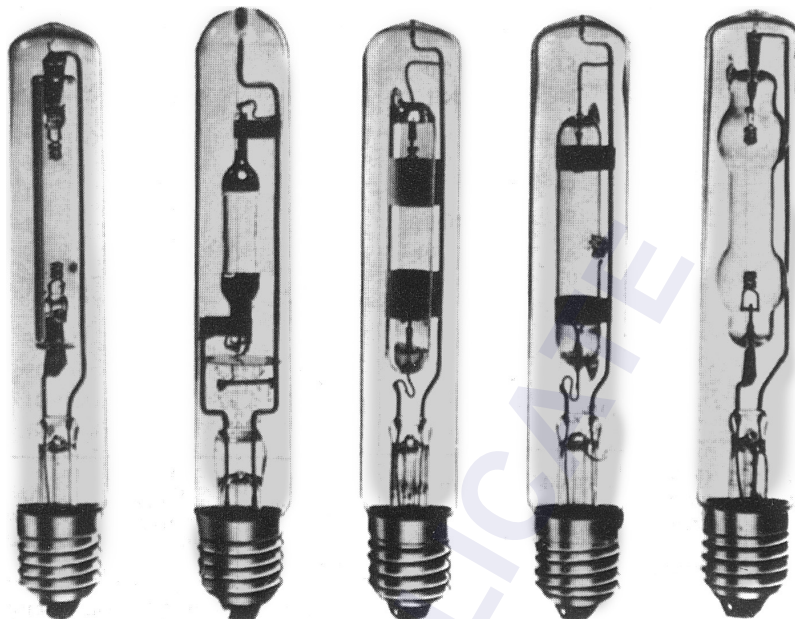


FIGURE 46 Examples of Philips spectral lamps.

TABLE 6 Specifications of Philips Spectral Lamps

Catalog Number	Symbols	Type	Material		Operating Current (A)	Wattage	Arc Length (mm)
			Burner	Envelope			
26-2709	Hg	Mercury (low-pressure)	Quartz	Glass	0.9	15	40
26-2717	Hg	Mercury (high-pressure)	Quartz	Glass	0.9	90	30
26-2725	Cd	Cadmium	Quartz	Glass	0.9	25	30
26-2733	Zn	Zinc	Quartz	Glass	0.9	25	30
26-2741	Hg, Cd, Zn	Mercury, cadmium, and zinc	Quartz	Glass	0.9	90	30
26-2758	He	Helium	Glass	Glass	0.9	45	32
26-2766	Ne	Neon	Glass	Glass	0.9	25	40
26-2774	A	Argon	Glass	Glass	0.9	15	40
26-2782	Kr	Krypton	Glass	Glass	0.9	15	40
26-2790	Xe	Xenon	Glass	Glass	0.9	10	40
26-2808	Na	Sodium	Glass	Glass	0.9	15	40
26-2816	Rb	Rubidium	Glass	Glass	0.9	15	40
26-2824	Cs	Caesium	Glass	Glass	0.9	10	40
26-2832	K	Potassium	Glass	Glass	0.9	10	40
26-2857	Hg	Mercury (low-pressure)	Quartz	Quartz	0.9	15	40
26-2865	Hg	Mercury (high-pressure)	Quartz	Quartz	0.9	90	30
26-2873	Cd	Cadmium	Quartz	Quartz	0.9	25	30
26-2881	Zn	Zinc	Quartz	Quartz	0.9	25	30
26-2899	Hg, Cd, Zn	Mercury, cadmium, and zinc	Quartz	Quartz	0.9	90	30
26-2907	In	Indium*	Quartz	Quartz	0.9	25	25
26-2915	Tl	Thallium	Quartz	Quartz	0.9	20	30
26-2923	Ga	Gallium	Quartz	Quartz	0.9	20	30

* Requires a Tesla coil to cause it to strike initially.



FIGURE 47 Physical construction of Pluecker spectrum tubes.

TABLE 7 Gas Fills in Pluecker Tubes*

Cenco Number	Type
87210	Argon Gas
87215	Helium Gas
87220	Neon Gas
87225	Carbonic Acid Gas
87230	Chlorine Gas
87235	Nitrogen Gas
87240	Nitrogen Gas
87242	Air
87245	Oxygen Gas
87255	Iodine Vapor
87256	Krypton Gas
87258	Xenon
87260	Mercury Vapor
87265	Water Vapor

*Consists of glass tube with overall length of 25 cm with capillary portion about 8.5 to 10 cm long. Glass-to-metal seal wires are welded in metal caps with loops for wire connection are firmly sealed to the ends. Power supply no. 87208 is recommended as a source of excitation.

*Pluecker Spectrum Tubes*³⁴ These are inexpensive tubes made of glass (Fig. 47) with an overall length of 25 cm and capillary portion of 8.5 to 10 cm long. They operate from an ordinary supply with a special transformer which supports the tubes in a vertical position and maintains the voltage and current values adequate to operate the discharge and regulate the spectral intensity. Table 7 lists the various gases in available tubes.

Concentrated Arc Lamps

*Zirconium Arc*²⁵ The cathodes of these lamps are made of a hollow refractory metal containing zirconium oxide. The anode, a disk of metal with an aperture, resides directly above the cathode with the normal to the aperture coincident with the longitudinal axis of the cathode. Argon gas fills the tube. The arc discharge causes the zirconium to heat (to about 3000 K) and produce an intense, very small source of light. These lamps have been demonstrated in older catalogs from the Cenco Company in a number of wattages (from 2 to 300). The end of the bulk through which the radiation passes comes with ordinary curvature or (for a slight increase in price) flat. Examples are shown in Fig. 48.

*Tungsten-Arc (Photomicrographic) Lamp*²⁵ The essential elements of this discharge-type lamp (see Fig. 49) are a ring electrode and a pellet electrode, both made of tungsten. The arc forms between these electrodes, causing the pellet to heat incandescently. The ring also incandesces, but to a lesser extent. Thus, the hot pellet (approximately 3100 K) provides an intense source of small-area

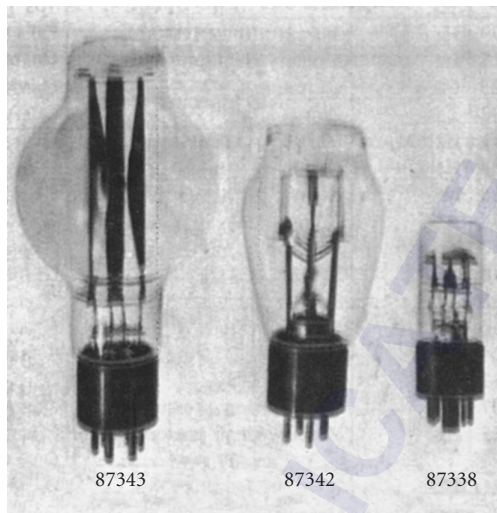


FIGURE 48 Physical construction of some zirconium arc lamps. Two 2-W lamps are available but not shown here.

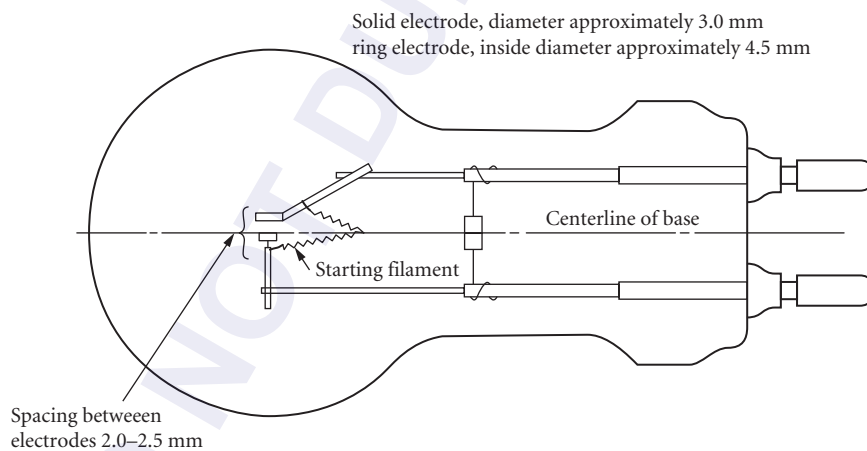


FIGURE 49 Construction of tungsten-arc lamp. The lamp must be operated baseup on a well-ventilated housing and using a special high-current socket which does not distort the position of the posts.

radiation. A plot of the spectral variation of this radiation is given in Fig. 50. As with all tungsten sources, evaporation causes a steady erosion of the pellet surface with the introduction of gradients, which is not serious if the pellet is used as a point source.

General Electric, manufacturer of the 30A/PS22 photomicrographic lamp, which uses a 30 Å operating current, states that this lamp requires a special heavy-duty socket obtainable through certain manufacturers suggested in its brochure, which may now be out of print in the original, but obtainable presumably as a copy from GE.

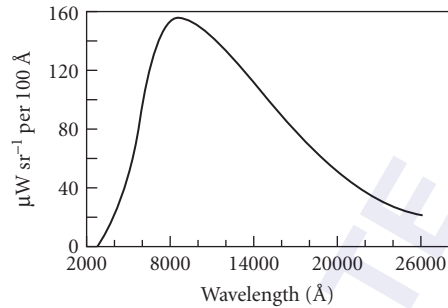


FIGURE 50 Spectral distribution of a 30/PS-22 photomicrographic lamp (superceded by the 330 watt, 30 amp PS-70).

Glow Modulator Tubes³⁵

According to technical data supplied by Sylvania, these are cold-cathode light sources uniquely adaptable to high-frequency modulation. (These tubes are now manufactured by The English Electric Valve Company, Elmsford, New York.) Pictures of two types are shown in Fig. 51. The cathode is a small hollow cylinder, and the high ionization density in the region of the cathode provides an intense source of radiation. Figure 52 is a graph of the light output as a function of tube current. Figure 53 is a graph depicting the response of the tube to a modulating input. The spectral outputs of a variety of tubes are shown in Fig. 54. Table 8 gives some of the glow-modulator specifications.

Hydrogen and Deuterium Arcs

For applications requiring a strong continuum in the ultraviolet region, the hydrogen arc at a few millimeters pressure provides a useful source. It can be operated with a cold or hot cathode. One hot-cathode type is shown in Fig. 55. Koller²² plots a distribution for this lamp down to about 200 Å.

Deuterium lamps (Fig. 56) provide a continuum in the ultraviolet with increased intensity over the hydrogen arc. Both lamps have quartz envelopes. The one on the left is designed for operation down to 2000 Å; the one on the right is provided with a Suprasil[®] window to increase the ultraviolet range down to 1650 Å. NIST is offering a deuterium lamp standard of spectral irradiance between 200 and 350 nm. The lamp output at 50 cm from its medium bipost base is about 0.7 W cm^{-3} at 200 nm and drops off smoothly to 0.3 W cm^{-3} at 250 nm and 0.07 W cm^{-3} at 350 nm. A working standard of the deuterium lamp can be obtained also, for example, from Optronic Laboratories, Incorporated, Orlando, Florida.

Other Commercial Sources

Activated-Phosphor Sources Of particular importance and convenience in the use of photometers are sources composed of a phosphor activated by radioactive substances. Readily available, and not subject to licensing with small quantities of radioactive material, are the ^{14}C -activated phosphor light sources. These are relatively stable sources of low intensity, losing about 0.02 percent per year due to the half-life of ^{14}C and the destruction of phosphor centers.

[®]Registered trademark of Heraeus-Amersil.

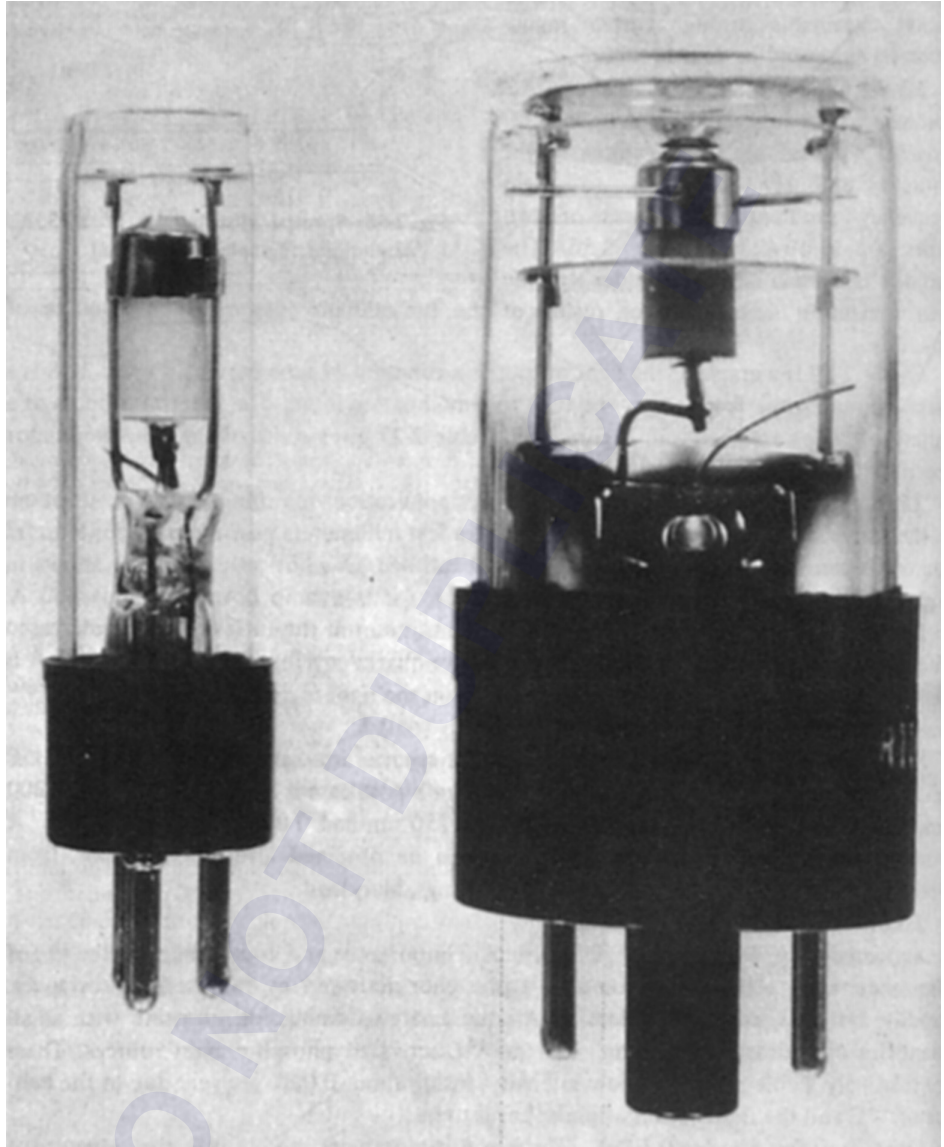


FIGURE 51 Construction of two glow modulator tubes.

Other (High-Energy) Sources Radiation at very high powers can be produced. Sources are synchrotrons, plasmotrons, arcs, sparks, exploding wires, shock tubes, and atomic and molecular beams, to name but a few. Among these, one can purchase in convenient, usable form precisely controlled spark-sources for yielding many joules of energy in a time interval of the order of microseconds. The number of vendors will be few, but check the directories.

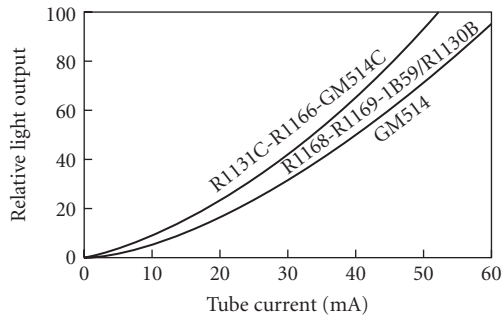


FIGURE 52 Variations of the light output from a glow modulator tube as a function of tube current.

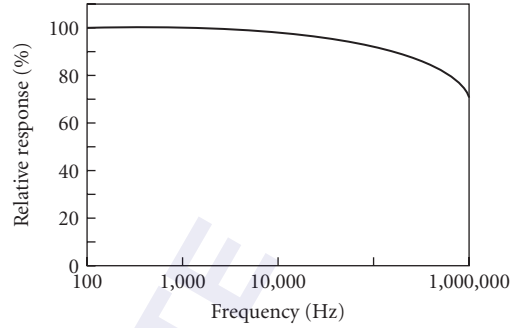
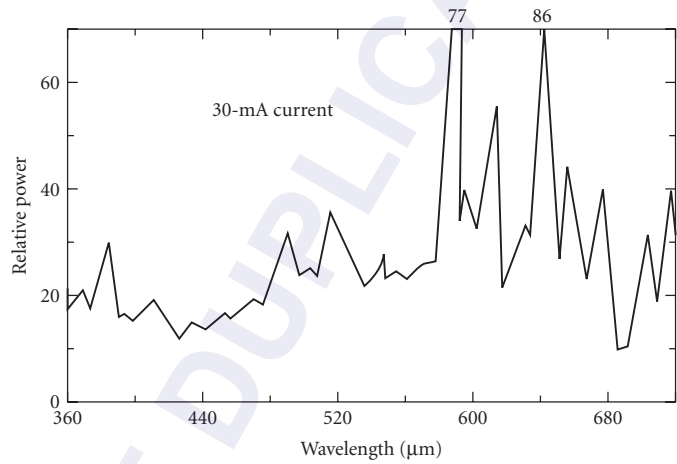
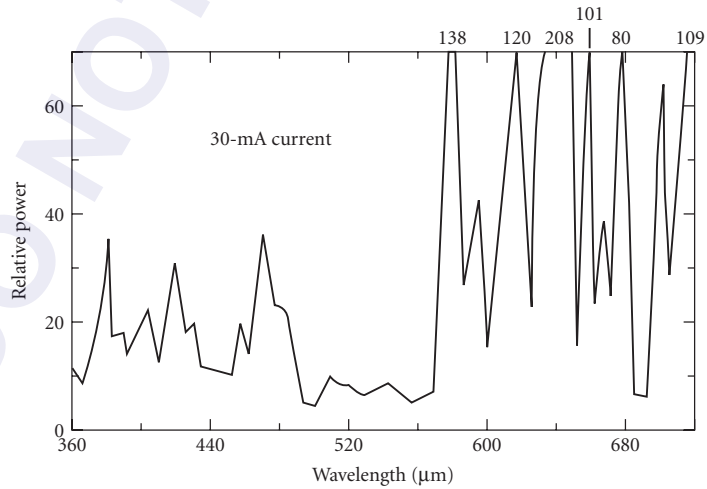


FIGURE 53 Response of the glow modulator tube to a modulating input.



(a) GM514C-R1166-R1131C



(b) R1168-R1169-1B59/R1130D-GM514

FIGURE 54 Spectral variation of the output of glow modulator tubes.

TABLE 8 Glow-Modulator Specifications

No.*	Maximum Operating Voltage	Current (mA)		Minimum Starting Voltage (V)	Crater Diameter (in.)	Approximate Light Center Length (in.)	Light Output (cd)	Brightness (cd in. ⁻²)	Rated Life (h)	Base Type	Bulb Type	Maximum Overall Length (in.)	Maximum Diameter (in.)	Color of Discharge
		Average	Peak											
GM-514	160	5-25	55	240	0.056	1-3/4	0.1 at 25 mA	41 at 25 mA	100 at 15 mA	3-pin miniature [†]	T-4 1/2 2-5/8	4 1/64	4 1/64	Blue-red
GM-514C	160	5-15	35	240	0.093	1-3/4	0.1 at 15 mA	15 at 15 mA	25 at 10 mA	3-pin miniature [‡]	T-4 1/2 2-5/8	4 1/64	4 1/64	White
IB59/														
R-1130D	150	5-35	75	225	0.056	2	0.13 at 30 mA	43 at 30 mA	250 at 20 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	Blue-red
R-1131C	150	3-25	55	225	0.093	2	0.2 at 25 mA	29 at 25 mA	150 at 15 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	White
R-1166	150	3-25	55	225	0.093	2	0.2 at 25 mA	29 at 25 mA	150 at 15 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	White
R-1168	150	5-15	30	225	0.015	2	0.023 at 15 mA	132 at 15 mA	150 at 15 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	Blue-red
R-1169	150	5-25	45	225	0.025	2	0.036 at 15 mA	72 at 15 mA	250 at 15 mA	Intermediate shell oct. [‡]	T-9	3-1/16	1-9/32	Blue-red

[†]Type R-1166 is opaque-coated with the exception of a circle 3/8 in. in diameter at end of lamp. All other types have clear-finish bulb.

[‡]Pins 1 and 3 arc anode; pin 2 cathode.

[§]Pin 7 anode; pin 3 cathode.

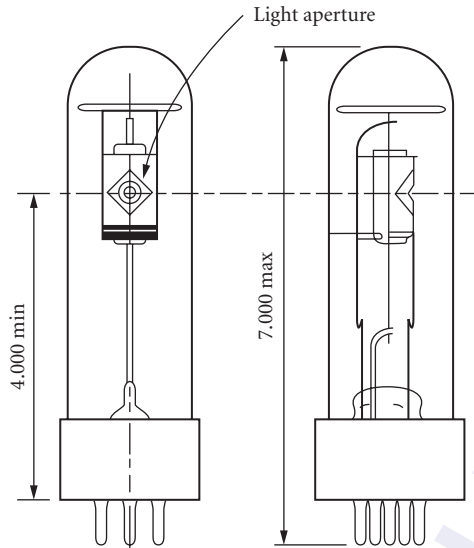


FIGURE 55 Hydrogen-arc lamp.

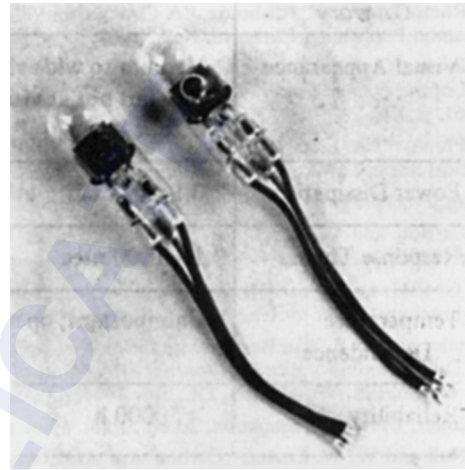


FIGURE 56 Two types of deuterium arc lamps.

Other Special Sources An enormous number of special-purpose sources are obtainable from manufacturers and scientific instrument suppliers. One source that remains to be mentioned is the so-called miniature, sub- and microminiature lamps. These are small, even tiny, incandescent bulbs of glass or quartz, containing tungsten filaments. They serve excellently in certain applications where small, intense radiators of visible and near-infrared radiation are needed. Second-source vendors advertise in the trade magazines.

15.6 REFERENCES

1. A. J. LaRocca, "Artificial Sources," in *The Infrared Handbook*, rev. ed., ONR, Washington, D.C., 1985, chap. 2.
2. A. G. Worthing, "Sources of Radiant Energy," in W. E. Forsythe (ed.), *Measurement of Radiant Energy*, McGraw-Hill, New York, 1937, chap. 2.
3. F. K. Richtmeyer and E. H. Kennard, *Introduction to Modern Physics*, McGraw-Hill Book Company, New York, 1947, p. 159.
4. T. J. Quinn, "The Calculation of the Emissivity of Cylindrical Cavities Giving Near Back-body Radiation," *British Journal of Applied Physics*, vol. 18, 1967, p. 1105.
5. J. C. DeVos, "Evaluation of the Quality of a Conical Blackbody," *Physica*, North Holland Publishing, Amsterdam, Netherlands, vol. 20, 1954, p. 669.
6. K. Irani, "Theory and Construction of Blackbody Calibration Sources," *Proceedings of SPIE*, vol. 4360, March 2001, p. 347.
7. Andre Gouffé, "Corrections d'Ouverture des Corps-noir Artificiels Compte Tenu des Diffusions Multiples Internes," in *Revue d'Optique*, vol. 24, Masson, Paris, 1945, p. 1.
8. A. Leupin, H. Vetsch, and F. K. Kneibuhl, "Investigation, Comparison, and Improvement of Technical Infrared Radiators," *Infrared Physics*, vol. 30, no. 3, 1990, pp. 199–258.
9. J. W. T. Walsh, *Photometry*, 3d ed., Dover, New York, 1965.
10. K. D. Mielenz, R. D. Saunders, and J. B. Shumaker, "Spectroradiometric Determination of the Freezing Temperature of Gold," *Journal of Research of the National Institute of Standards and Technology*, vol. 95, no. 1, Jan.–Feb. 1990, pp. 49–67.

11. *Spectral Radiance Calibrations*. NBS Special Publication, Jan. 1987, p. 250–251.*
12. *NIST Calibration Services Users Guide 1989*, NIST Special Publication, 1989, p. 250.*
13. *Lasers and Optronics 1990 Buying Guide*, Elsevier Communications, Morris Plains, N.J., 1990.*
14. *Photonics Directory of Optical Industries*, Laurin Publishing Co., Pittsfield, Mass., 1994.*
15. W. Y. Ramsey and J. C. Alishouse, “A Comparison of Infrared Sources,” *Infrared Physics*, vol. 8, Pergamon Publishing, Elmsford, N.Y., 1968, p. 143.
16. J. C. Morris, “Comments on the Measurements of Emittance of the Globar Radiation Source,” *Journal of the Optical Society of America*, vol. 51, Optical Society of America, Washington, D.C., July 1961, p. 798.
17. A. H. Pfund, “The Electric Welsbach Lamp,” *Journal of the Optical Society of America*, vol. 26, Optical Society of America, Washington, D.C., Dec. 1936, p. 439.
18. J. Strong, *Procedures in Experimental Physics*, Prentice-Hall, New York, 1938, p. 346.
19. F. E. Carlson and C. N. Clark, “Light Sources for Optical Devices,” in R. Kingslake (ed.), *Applied Optics and Optical Engineering*, vol. 1, Academic Press, New York, 1965, p. 80.
20. G. M. B. H. Osram, *Lamps for Scientific Purposes*, Munchen, West Germany, 1966. (See also the later Osram brochure, *Light for Cine Projection, Technology and Science*, 1987.)*
21. F. J. Studer and R. F. Van Beers, “Modification of Spectrum of Tungsten Filament Quartz-Iodine Lamps due to Iodine Vapor,” *Journal of the Optical Society of America*, vol. 54, no. 7, Optical Society of America, Washington, D.C., July 1964, p. 945.
22. L. R. Koller, *Ultraviolet Radiation*, 2d ed., John Wiley and Sons, New York, 1965.
23. M. R. Null and W. W. Lozier, “Carbon Arc as a Radiation Standard,” *Journal of the Optical Society of America*, vol. 52, no. 10, Optical Society of America, Washington, D.C., Oct. 1962, pp. 1156–1162.
24. E. B. Noel, “Radiation from High Pressure Mercury Arcs,” *Illuminating Engineering*, vol. 36, 1941, p. 243.
25. General Electric, *Bulletin TP-109R*, Cleveland, 1975.*
26. Illumination Industries, *Catalog No. 108-672-3M*, Sunnyvale, Calif., 1972.*
27. M. W. P. Cann, *Light Sources for Remote Sensing Systems*, in NASA-CR-854, Aug. 1967.*
28. Varian Associates, Palo Alto, Calif., 1969.*
29. Fisher Scientific, *Special Catalog to Spectrophotometer Users*, Pittsburgh, Mass., 1972.*
30. E. F. Worden, R. G. Gutmacher, and J. F. Conway, “Use of Electrodeless Discharge Lamps in the Analysis of Atomic Spectra,” *Applied Optics*, vol. 2, no. 7, Optical Society of America, Washington, D.C., July 1963, pp. 707–713.
31. W. F. Meggers and F. O. Westfall, “Lamps and Wavelengths of Mercury 198,” *NBS Journal of Research*, vol. 44, National Bureau of Standards, Washington, D.C., 1950, pp. 447–55.
32. W. F. Meggers, “Present Experimental Status of Rare Earth Spectra,” *Journal of the Optical Society of America*, vol. 50, Optical Society of America, Washington, D.C., 1960, p. 405.
33. Ealing Corporation, *Ealing Catalog, Optical Components Section*, South Natick, Mass., 1976–77. (See also *Optical Services Supplement*, no. 1, 1969–70, p. 26; See also *Ealing Electro-optics Product Guide*, 1990.)*
34. Central Scientific, *Cenco Scientific Education Catalog*, Physics-Light Section, Chicago, 1975, p. 602.*
35. GTE Sylvania, *Special Purpose Lamps*, Lighting Products Group, Danvers, Mass., in TR-29R, May 1966.*

*Many of these references are likely to be inaccessible as shown, either because they are out-of-date or perhaps otherwise obsolete. They are retained here because there are up-to-date versions of many of them, and the user is advised to use them as starting points in a search of the Internet for current information. Most of the companies still exist and have more recent catalogs that are available.

William T. Silfvast

CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida

16.1 GLOSSARY

A_2	radiative transition probability from level 2 to all other possible lower-lying levels
A_{21}	radiative transition probability from level 2 to level 1
a_L	scattering losses within a laser cavity for a single pass through the cavity
B_{12}	Einstein B coefficient associated with absorption
B_{21}	Einstein B coefficient associated with stimulated emission
E_1, E_2	energies of levels 1 and 2 above the ground state energy for that species
g_1, g_2	stability parameters for laser modes when describing the laser optical cavity
g_1, g_2	statistical weights of energy levels 1 and 2 that indicate the degeneracy of the levels
g_{21}	gain coefficient for amplification of radiation within a medium at a wavelength of λ_{21}
I_{sat}	saturation intensity of a beam in a medium; intensity at which exponential growth will cease to occur even though the medium has uniform gain (energy/time-area)
N_1, N_2	population densities (number of species per unit volume) in energy levels 1 and 2
r_c	radius of curvature of the expanding wavefront of a gaussian beam
R_1, R_2	reflectivities of mirrors 1 and 2 at the desired wavelength
T_1	lifetime of a level when dominated by collisional decay
T_2	average time between phase-interrupting collisions of a species in a specific excited state
t_{opt}	optimum mirror transmission for a laser of a given gain and loss
$w(z)$	beam waist radius at a distance z from the minimum beam waist for a gaussian beam
w_0	minimum beam waist radius for a gaussian mode
α_{12}	absorption coefficient for absorption of radiation within a medium at wavelength λ_{21}
γ_{21}	angular frequency bandwidth of an emission or absorption line
Δt_p	pulse duration of a mode-locked laser pulse
$\Delta\nu$	frequency bandwidth over which emission, absorption, or amplification can occur

$\Delta\nu_D$	frequency bandwidth (FWHM) when the dominant broadening process is Doppler or motional broadening
η	index of refraction of the laser medium at the desired wavelength
λ_{21}	wavelength of a radiative transition occurring between energy levels 2 and 1
ν_{21}	frequency of a radiative transition occurring between energy levels 2 and 1
σ_{21}	stimulated emission cross section (area) associated with levels 2 and 1
σ_{21}^D	stimulated emission cross section at line center when Doppler broadening dominates (area)
σ_{21}^H	stimulated emission cross section at line center when homogeneous broadening dominates (area)
τ_2	lifetime of energy level 2
τ_{21}	lifetime of energy level 2 if it can only decay to level 1

16.2 INTRODUCTION

A laser is a device that amplifies light and produces a highly directional, high-intensity beam that typically has a very pure frequency or wavelength. It comes in sizes ranging from approximately one-tenth the diameter of a human hair to the size of a very large building, in powers ranging from 10^{-9} to 10^{20} W and in wavelengths ranging from the microwave to the soft-x-ray spectral regions with corresponding frequencies from 10^{11} to 10^{17} Hz. Lasers have pulse energies as high as 10^4 J and pulse durations as short as 6×10^{-15} seconds. They can easily drill holes in the most durable of materials and can weld detached retinas within the human eye.

Lasers are a key component of some of our most modern communication systems and are the “phonograph needle” of compact disc players. They are used for heat treatment of high-strength materials, such as the pistons of automobile engines, and provide a special surgical knife for many types of medical procedures. They act as target designators for military weapons and are used in the checkout scanners we see everyday at the supermarket.

The word *laser* is an acronym for *Light Amplification by Stimulated Emission of Radiation*. The laser makes use of processes that increase or amplify light signals after those signals have been generated by other means. These processes include (1) stimulated emission, a natural effect that arises out of considerations relating to thermodynamic equilibrium, and (2) optical feedback (present in most lasers) that is usually provided by mirrors. Thus, in its simplest form, a laser consists of a gain or amplifying medium (where stimulated emission occurs) and a set of mirrors to feed the light back into the amplifier for continued growth of the developing beam (Fig. 1).

The entire spectrum of electromagnetic radiation is shown in Fig. 2, along with the region covered by currently existing lasers. Such lasers span the wavelength range from the far infrared part of

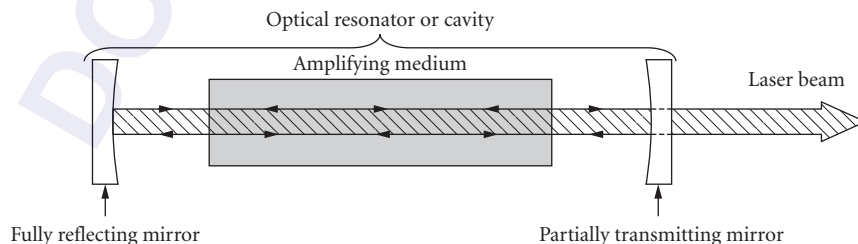


FIGURE 1 Simplified diagram of a laser, including the amplifying medium and the optical resonator.

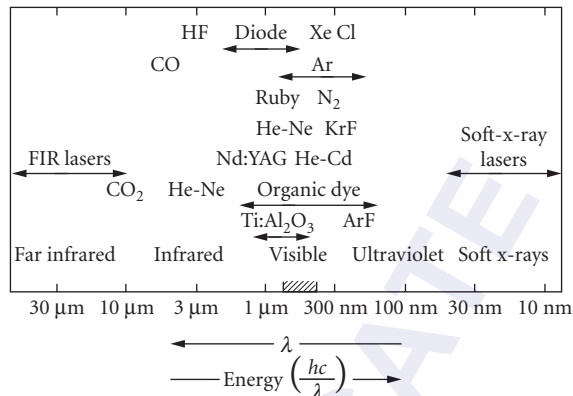


FIGURE 2 The portion of the electromagnetic spectrum that involves lasers, along with the general wavelengths of operation of most of the common lasers.

the spectrum ($\lambda = 1000 \mu\text{m}$) to the soft-x-ray region ($\lambda = 3 \text{ nm}$), thereby covering a range of almost six orders of magnitude! There are several types of units that are used to define laser wavelengths. These range from micrometers (μm) in the infrared to nanometers (nm) and angstroms (\AA) in the visible, ultraviolet (UV), vacuum ultraviolet (VUV), extreme ultraviolet (EUV or XUV), and soft-x-ray (SXR) spectral regions.

This chapter provides a brief overview of how a laser operates. It considers a laser as having two primary components: (1) a region where light amplification occurs which is referred to as a *gain medium* or an *amplifier*, and (2) a *cavity*, which generally consists of two mirrors placed at either end of the amplifier.

Properties of the amplifier include the concept of discrete excited energy levels and their associated finite lifetimes. The broadening of these energy levels will be associated with the emission linewidth which is related to decay of the population in these levels.

Stimulated emission will be described, and the formulas for calculating the amount of gain that can occur via stimulated emission will be given in terms of the radiative properties of the medium. The concept of the saturation intensity will be introduced and related to the amount of gain that is necessary for laser output. The addition of mirrors at the ends of the amplifier will be used to increase the gain length and to reduce the divergence of the amplified beam. The threshold conditions for laser output will be described in terms of the amplifier properties and the mirror reflectivities. This section will conclude with a review of excitation or pumping processes that are used to produce the necessary population density in the upper laser level.

Cavity properties will begin with a discussion of both longitudinal and transverse cavity modes which provide the laser beam with a gaussian-shaped transverse profile. The properties of those gaussian beams will be reviewed. The types of optical cavities that allow stable operation of laser modes will then be described. A number of special types of laser cavity arrangements and techniques will be reviewed, including unstable resonators, Q-switching, mode-locking, and ring lasers. A brief review will then be given of the various common types of gaseous, liquid, and solid-state lasers.

Additional information related to spectral lineshape and the mechanisms of spectral broadening can be found in Chap. 10 (Vol. I), "Optical Spectroscopy and Spectroscopic Lineshapes." Other related material can be found in Chap. 8 (Vol. IV), "Fundamental Optical Properties of Solids," and Chap. 33 (Vol. I), "Holography and Holographic Instruments." As lasers are widely used in many of the devices and techniques discussed in other chapters in this *Handbook*, the reader is directed to those topics for information on specific lasers.

16.3 LASER PROPERTIES ASSOCIATED WITH THE LASER GAIN MEDIUM

Energy Levels and Radiation^{1,2}

Nearly all lasers involve electronic charge distributions of atoms, molecules, organic dye solutions, or solids that make transitions from one energy state or level E_2 to another lower-lying level E_1 . The loss of energy resulting from this transition is given off in the form of electromagnetic radiation. The relationship between the energy difference between the levels, $E_2 - E_1$ or ΔE_{21} , and the frequency ν_{21} of radiation occurring as a result of the transition, is determined by the Einstein relationship $E_{21} = h\nu_{21}$ where h is Planck's constant. It was first shown by Bohr in 1913 that the discrete set of emission wavelengths from a hydrogen discharge could be explained by the occurrence of discrete energy levels in the hydrogen atom that have a fixed relationship. This discrete arrangement of energy levels was later shown to occur in other atoms, in molecules, and also in liquids and solids. In atoms these energy levels are very precisely defined and narrow in width ($\approx 10^9$ Hz) and can be accurately calculated with sophisticated atomic physics codes. In molecules and high-density materials the locations of the levels are more difficult to calculate and they tend to be much broader in width, the largest widths occurring in liquids and solids (up to 5×10^{13} Hz).

The lowest energy level of a species is referred to as the *ground state* and is usually the most stable state of the species. There are some exceptions to this, for example, ground states of ionized species or unstable ground states of some molecular species such as excimer molecules. Energy levels above the ground state are inherently unstable and have lifetimes that are precisely determined by the arrangement of the atoms and electrons associated with any particular level as well as to the particular species or material. Thus, when an excited state is produced by applying energy to the system, that state will eventually decay by emitting radiation over a time period ranging from 10^{-15} seconds or less to times as long as seconds or more, depending upon the particular state or level involved. For *strongly allowed* transitions that involve the electron charge cloud changing from an atomic energy level of energy E_2 to a lower-lying level of energy E_1 , the radiative decay time τ_{21} can be approximated by $\tau_{21} \approx 10^4 \lambda_{21}^2$ where τ_{21} is in seconds and λ_{21} is the wavelength of the emitted radiation in meters. λ_{21} is related to ν_{21} by the relationship $\lambda_{21}\nu_{21} = c/\eta$ where c is the velocity of light (3×10^8 m/s) and η is the index of refraction of the material. For most gases, η is near unity and for solids and liquids it ranges between 1 and 10, with most values ranging from ≈ 1.3 to 2.0.

Using the expression suggested above for the approximate value of the lifetime of an excited energy level, one obtains a decay time of several ns for green light ($\lambda_{21} = 5 \times 10^{-7}$ m). This represents a minimum radiative lifetime since most excited energy levels have a weaker radiative decay probability than mentioned above and would therefore have radiative lifetimes one or two orders of magnitude longer. Other laser materials such as molecules, organic dye solutions, and semiconductor lasers have similar radiative lifetimes. The one exception is the class of dielectric solid-state laser materials (both crystalline and glass) such as ruby and Nd:YAG, in which the lifetimes are of the order of 1 μ s to 3 ms. This much longer radiative lifetime in solid-state laser materials is due to the nature of the particular state of the laser species and to the crystal matrix in which it is contained. This is a very desirable property for a laser medium since it allows excitation and energy storage within the laser medium over a relatively long period of time.

Emission Linewidth and Line Broadening of Radiating Species^{1,3}

Assume that population in energy level 2 decays to energy level 1 with an exponential decay time of τ_{21} and emits radiation at frequency ν_{21} during that decay. It can be shown by Fourier analysis that the exponential decay of that radiation requires the frequency width of the emission to be of the order of $\Delta\nu \approx 1/2\pi\tau_{21}$. This suggests that the energy width ΔE_2 of level 2 is of the order of $\Delta E_2 = h/2\pi\tau_{21}$. If the energy level 2 can decay to more levels than level 1, with a corresponding decay time of τ_2 , then its energy is broadened by an amount $\Delta E_2 = h/2\pi\tau_2$. If the decay is due primarily to

radiation at a rate A_{2i} to one or more individual lower-lying levels i , then $1/\tau_2 = A_2 = \sum_i A_{2i}$. A_2 represents the total radiative decay rate of level 2, whereas A_{2i} is the specific radiative decay rate from level 2 to a lower-lying level i .

If population in level 2 decays radiatively at a radiative rate A_2 and population in level 1 decays radiatively at a rate A_1 , then the emission linewidth of radiation from level 2 to 1 is given by

$$\Delta\nu_{21} = \frac{\sum_i A_{2i} + \sum_j A_{1j}}{2\pi} \quad (1)$$

which is referred to as the *natural linewidth* of the transition and represents the sum of the widths of levels 2 and 1 in frequency units. If, in the above example, level 1 is a ground state with infinite lifetime or a long-lived metastable level, then the natural linewidth of the emission from level 2 to level 1 would be represented by

$$\Delta\nu_{21} = \frac{\sum_i A_{2i}}{2\pi} \quad (2)$$

since the ground state would have an infinite lifetime and would therefore not contribute to the broadening. This type of linewidth or broadening is known as *natural broadening* since it results specifically from the radiative decay of a species. Thus the natural linewidth associated with a specific transition between two levels has an inherent value determined only by the factors associated with specific atomic and electronic characteristics of those levels.

The emission-line broadening or natural broadening described above is the minimum line broadening that can occur for a specific radiative transition. There are a number of mechanisms that can increase the emission linewidth. These include collisional broadening, phase-interruption broadening, Doppler broadening, and isotope broadening. The first two of these, along with natural broadening, are all referred to as *homogeneous broadening*. Homogeneous broadening is a type of emission broadening in which all of the atoms radiating from the specific level under consideration participate in the same way. In other words, all of the atoms have the identical opportunity to radiate with equal probability.

The type of broadening associated with either Doppler or isotope broadening is referred to as *inhomogeneous broadening*. For this type of broadening, only certain atoms radiating from that level that have a specific property such as a specific velocity, or are of a specific isotope, participate in radiation at a certain frequency within the emission bandwidth.

Collisional broadening is a type of broadening that is produced when surrounding atoms, molecules, solvents (in the case of dye lasers), or crystal structures interact with the radiating level and cause the population to decay before it has a chance to decay by its normal radiative processes. The emission broadening is then associated with the faster decay time T_1 , or $\Delta\nu = 1/2\pi T_1$.

Phase-interruption broadening or *phonon broadening* is a type of broadening that does not increase the decay rate of the level, but it does interrupt the phase of the rotating electron cloud on average over a time interval T_2 which is much shorter than the radiative decay time τ_2 which includes all possible radiative decay channels from level 2. The result of this phase interruption is to increase the emission linewidth beyond that of both natural broadening and T_1 broadening (if it exists) to an amount $\Delta\nu = 1/2\pi T_2$.

Doppler broadening is a type of inhomogeneous broadening in which the Doppler effect shifts the frequencies of radiating atoms moving toward the observer to a higher value and the frequencies of atoms moving away from the observer to a lower value. This effect occurs only in gases since they are the only species that are moving fast enough to produce such broadening. Doppler broadening is the dominant broadening process in most visible gas lasers. The expression for the Doppler linewidth (FWHM) is given by

$$\Delta\nu_D = 7.16 \cdot 10^{-7} \nu_o \sqrt{\frac{T}{M}} \quad (3)$$

in which ν_o is the center frequency associated with atoms that are not moving either toward or away from the observer, T is the gas temperature in kelvin and M is the atomic or molecular weight (number of nucleons/atom or molecule) of the gas atoms or molecules.

Isotope broadening also occurs in some gas lasers. It becomes the dominant broadening process if the specific gas consists of several isotopes of the species and if the isotope shifts for the specific radiative transition are broader than the Doppler width of the transition. The helium-cadmium laser is dominated more by this effect than any other laser since the naturally occurring cadmium isotopic mixture contains eight different isotopes and the isotope shift between adjacent isotopes (adjacent neutron numbers) is approximately equal to the Doppler width of the individual radiating isotopes. This broadening effect can be eliminated by the use of isotopically pure individual isotopes, but the cost for such isotopes is often prohibitive.

All homogeneous broadening processes have a frequency distribution that is described by a Lorentzian mathematical function

$$I_{21}(\nu) = I_o \frac{\gamma_{21}/4\pi^2}{(\nu - \nu_o)^2 + (\gamma_{21}/4\pi)^2} \quad (4)$$

in which γ_{21} represents the decay rate of level 2, I_o is the total emission intensity of the transition over the entire linewidth, and ν_o is the center frequency of the emission line. In Eq. (4), γ_{21} is determined by the relationship $\gamma_{21} = 2\pi\Delta\nu_{21}$. For natural broadening, $\Delta\nu_{21}$ is given by either Eq. (1) or Eq. (2), whichever is applicable. For T_1 -dominated broadening, $\Delta\nu_{21} = 1/2\pi T_1$, and for T_2 -dominated broadening, $\gamma_{21} = 1/2\pi T_2$.

The frequency distribution for Doppler broadening is described by a gaussian function

$$I(\nu) = \frac{2(\ln 2)^{1/2}}{\pi^{1/2}\Delta\nu_D} I_o \exp\left[-\frac{4\ln 2(\nu - \nu_o)^2}{\Delta\nu_D^2}\right] \quad (5)$$

Both of these lineshape functions are indicated in Fig. 3. In this figure, both the total emission intensity integrated over all frequencies and the emission linewidth (full width at half maximum or FWHM) for both functions are identical.

Isotope broadening involves the superposition of a series of either Lorentzian shapes or Gaussian shapes for each isotope of the species, separated by the frequencies associated with the isotope shifts of that particular transition.

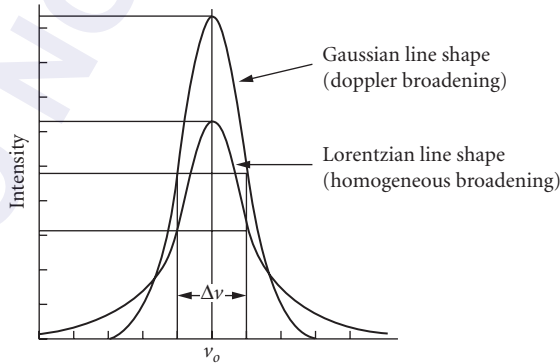


FIGURE 3 Lineshape functions for both homogeneous broadening, with a Lorentzian shape, and Doppler broadening (inhomogeneous), with a Gaussian shape. Both lines are arranged with equal linewidths (FWHM) and equal total intensities.

TABLE 1 Amplifier Parameters for a Wide Range of Lasers

Type of Laser	λ_{21} (nm)	τ_2 (s)	$\Delta\nu_{21}$ (Hz)	$\Delta\lambda_{21}$ (nm)	σ_{21} (m ²)	ΔN_{21} (m ⁻³)	L (m)	g_{21} (m ⁻¹)
Helium-Neon	632.8	3×10^{-7}	2×10^9	2×10^{-3}	3×10^{-17}	5×10^{15}	0.2	0.15
Argon	488.0	1×10^{-8}	2×10^9	1.6×10^{-3}	5×10^{-16}	1×10^{15}	0.2–1.0	0.5
He-Cadmium	441.6	7×10^{-7}	2×10^9	1.3×10^{-3}	8×10^{-18}	4×10^{16}	0.2–1.0	0.3
Copper vapor	510.5	5×10^{-7}	2×10^9	1.3×10^{-3}	8×10^{-18}	6×10^{17}	1.0–2.0	5
CO ₂	10,600	4	6×10^7	2.2×10^{-2}	1.6×10^{-20}	5×10^{19}	0.2–2.0	0.8
Excimer	248.0	9×10^{-9}	1×10^{13}	2	2.6×10^{-20}	1×10^{20}	0.5–1.0	2.6
Dye (Rh6G)	577.0	5×10^{-9}	5×10^{13}	60	1.2×10^{-20}	2×10^{22}	0.01	240
Semiconductor	800	1×10^{-9}	1×10^{13}	20	1×10^{-19}	10^{24}	0.00025	100,000
Nd:Yag	1064.1	2.3×10^{-4}	1.2×10^{11}	0.4	6.5×10^{-23}	1.6×10^{23}	0.1	10
Nd:Glass	1054	3.0×10^{-4}	7.5×10^{12}	26	4.0×10^{-24}	8×10^{23}	0.1	3
Cr:LiSAF	840	6.7×10^{-5}	9.0×10^{13}	250	5.0×10^{-24}	2×10^{24}	0.1	10
Ti:Al ₂ O ₃	760	3.2×10^{-6}	1.5×10^{14}	400	4.1×10^{-23}	5×10^{23}	0.1	20

Table 1 gives examples of the dominant broadening process and the value of the broadening for most of the common commercial lasers.

Bandwidths of laser transitions in semiconductors are actually made narrower by making the active regions of the material extremely thin in one or more dimensions. In doing so the energy levels become quantized and thus behave more like single atom atomic levels. Quantizing in one dimension, by making the thickness of the order of 50 to 100 nm, leads to a quantum well laser. Narrowing and thereby quantizing in two dimensions is referred to as a quantum wire and in three dimensions, a quantum dot. The advantages of quantizing the dimensions is that the reduced thickness leads to a significantly reduced heat loss during the excitation of the semiconductor as well as a narrower laser emission linewidth because of the smaller size of the electron energy distribution in the upper laser level. In the case of the quantum dot, the material takes on atom-like properties because the energies are quantized in all three dimensions and the lasing threshold is reduced much more so than even with the quantum well laser. Of course there is less laser gain medium produced in such materials per unit volume, because of the reduced gain volume, and thereby less laser power per unit volume. This can be made up by having many such gain media in parallel and/or in series, taking into account the heat removal requirements of the pumping process.

Stimulated Radiative Processes— Absorption and Emission^{1,2}

Two types of stimulated radiative processes, absorption and stimulated emission, occur between energy levels 1 and 2 of a gain medium when light of frequency ν_{21} corresponding to an energy difference $\Delta E_{21} = (E_2 - E_1) = h\nu_{21}$ passes through the medium. These processes are proportional to the light intensity I as indicated in Fig. 4 for a two-level system as well as to the stimulated absorption and emission coefficients B_{12} and B_{21} , respectively. These coefficients are related to the frequency ν_{21} and the spontaneous emission probability A_{21} associated with the two levels. A_{21} has units of (1/s).

Absorption results in the loss of light of intensity I when the light interacts with the medium. The energy is transferred from the beam to the medium by raising population from level 1 to the higher-energy level 2. In this situation, the species within the medium can either reradiate the energy and return to its initial level 1, it can reradiate a different energy and decay to a different level, or it can lose the energy to the surrounding medium via collisions, which results in the heating of the medium, and return to the lower level. The absorption probability is proportional to the intensity I which has units of energy/s-m² times B_{12} , which is the absorption probability coefficient for that transition. B_{12} is one of the Einstein B coefficients and has the units of m³/energy-s².

Stimulated emission results in the increase in the light intensity I when light of the appropriate frequency ν_{21} interacts with population occupying level 2 of the gain medium. The energy is given

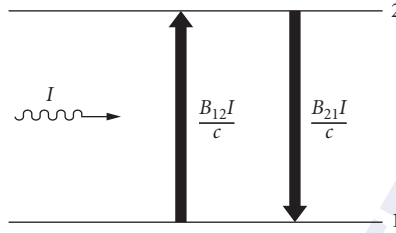


FIGURE 4 Stimulated emission and absorption processes that can occur between two energy levels, 1 and 2, and can significantly alter the population densities of the levels compared to when a beam of intensity I is not present.

up by the species to the radiation field. In the case of stimulated emission, the emitted photons or bundles of light have exactly the same frequency ν_{21} and direction as the incident photons of intensity I that produce the stimulation. B_{21} is the associated stimulated emission coefficient and has the same units as B_{12} . It is known as the other Einstein B coefficient.

Einstein showed the relationship between B_{12} , B_{21} , and A_{21} as

$$g_2 B_{21} = g_1 B_{12} \quad (6)$$

and

$$A_{21} = \frac{8\pi h \nu_{21}^3}{c^3} B_{21} \quad (7)$$

where g_2 and g_1 are the statistical weights of levels 2 and 1 and h is Planck's constant. Since $A_{21} = 1/\tau_{21}$ for the case where radiative decay dominates and where there is only one decay path from level 2, B_{21} can be determined from lifetime measurements or from absorption measurements on the transition at frequency ν_{21} .

Population Inversions^{1,2}

The two processes of absorption and stimulated emission are the principal interactions involved in a laser amplifier. Assume that a collection of atoms of a particular species is energized to populate two excited states 1 and 2 with population densities N_1 and N_2 (number of species/m³) and state 2 is at a higher energy than state 1 by an amount ΔE_{21} as described in the previous section. If a photon beam of energy ΔE_{21} , with an intensity I_0 and a corresponding wavelength $\lambda_{21} = c/\nu_{21} = hc/\Delta E_{21}$, passes through this collection of atoms, then the intensity I after the beam emerges from the medium can be expressed as

$$I = I_0 e^{\sigma_{21}(N_2 - (g_2/g_1)N_1)L} = I_0 e^{\sigma_{21}\Delta N_{21}L} \quad (8)$$

where σ_{21} is referred to as the *stimulated emission cross section* with dimensions of m², L is the thickness of the medium (in meters) through which the beam passes, and $N_2 - (g_2/g_1)N_1 = \Delta N_{21}$ is known as the *population inversion density*. The exponents in Eq. (8) are dimensionless quantities that can be either greater or less than unity, depending upon whether N_2 is greater than or less than $(g_2/g_1)N_1$.

The general form of the stimulated emission cross section per unit frequency is given as

$$\sigma_{21} = \frac{\lambda_{21}^2 A_{21}}{8\pi \Delta\nu} \quad (9)$$

in which $\Delta\nu$ represents the linewidth over which the stimulated emission or absorption occurs.

For the case of homogeneous broadening, at the *center* of the emission line, σ_{21} is expressed as

$$\sigma_{21}^H = \frac{\lambda_{21}^2 A_{21}}{4\pi^2 \Delta\nu_{21}^H} \quad (10)$$

where $\Delta\nu_{21}^H$ is the homogeneous emission linewidth (FWHM) which was described earlier for several different situations.

For the case of Doppler broadening, σ_{21}^D can be expressed as

$$\sigma_{21}^D = \sqrt{\frac{\ln 2}{16\pi^3}} \frac{\lambda_{21}^2 A_{21}}{\Delta\nu_{21}^D} \quad (11)$$

at the *center* of the emission line and $\Delta\nu_{21}^D$ is the Doppler emission linewidth expressed earlier in Eq. (3).

For all types of matter, the population density ratio of levels 1 and 2 would normally be such that $N_2 \ll N_1$. This can be shown by the Boltzmann relationship for the population ratio in thermal equilibrium which provides the ratio of N_2/N_1 to be

$$\frac{N_2}{N_1} = \frac{g_2}{g_1} e^{-\Delta E_{21}/kT} \quad (12)$$

in which T is the temperature and k is Boltzmann's constant. Thus, for energy levels separated by energies corresponding to visible transitions, in a medium at or near room temperature, the ratio of $N_2/N_1 \cong e^{-100} = 10^{-44}$. For most situations in everyday life we can ignore the population N_2 in the upper level and rewrite Eq. (8) as

$$I = I_o e^{-\sigma_{12} N_1 L} = I_o e^{-\alpha_{12} L} \quad (13)$$

where $\sigma_{21} = (g_2/g_1)\sigma_{12}$ and $\alpha_{12} = \sigma_{12} N_1$. Equation (13) is known as Beer's law, which is used to describe the absorption of light within a medium. α_{12} is referred to as the absorption coefficient with units of m^{-1} .

In laser amplifiers we cannot ignore the population in level 2. In fact, the condition for amplification and laser action is

$$N_2 - (g_2/g_1)N_1 > 0 \quad \text{or} \quad g_1 N_2 / g_2 N_1 > 1 \quad (14)$$

since, if Eq. (14) is satisfied, the exponent of Eq. (8) will be greater than 1 and I will emerge from the medium with a greater value than I_o , or amplification will occur. The condition of Eq. (14) is a necessary condition for laser amplification and is known as a population inversion since N_2 is greater than $(g_2/g_1)N_1$. In most cases, (g_2/g_1) is either unity or close to unity.

Considering that the ratio of $g_1 N_2 / g_2 N_1$ [Eq. (14)] could be greater than unity does not follow from normal thermodynamic equilibrium considerations [Eq. (12)] since it would represent a population ratio that could never exist in thermal equilibrium. When Eq. (14) is satisfied and amplification of the beam occurs, the medium is said to have *gain* or *amplification*. The factor $\sigma_{21}(N_2 - (g_2/g_1)N_1)$ or $\sigma_{21}\Delta N_{21}$ is often referred to as the *gain coefficient* g_{21} and is given in units of m^{-1} such that $g_{21} = \sigma_{21}\Delta N_{21}$. Typical values of σ_{21} , ΔN_{21} , and g_{21} are given in Table 1 for a variety

of lasers. The term g_{21} is also referred to as the *small-signal gain coefficient* since it is the gain coefficient determined when the laser beam intensity within the laser gain medium is small enough that stimulated emission does not significantly alter the populations in the laser levels.

Gain Saturation²

It was stated in the previous section that Eq. (14) is a necessary condition for making a laser but it is not a sufficient condition. For example, a medium might satisfy Eq. (14) by having a gain of $e^{g_{21}L} \approx 10^{-10}$, but this would not be sufficient to allow any reasonable beam to develop. Lasers generally start by having a pumping process that produces enough population in level 2 to create a population inversion with respect to level 1. As the population decays from level 2, radiation occurs spontaneously on the transition from level 2 to level 1 equally in all directions within the gain medium. In most of the directions very little gain or enhancement of the spontaneous emission occurs, since the length is not sufficient to cause significant growth according to Eq. (8). It is only in the elongated direction of the amplifier, with a much greater length, that significant gain exists and, consequently, the spontaneous emission is significantly enhanced. The requirement for a laser beam to develop in the elongated direction is that the exponent of Eq. (8) be large enough for the beam to grow to the point where it begins to significantly reduce the population in level 2 by stimulated emission. The beam will eventually grow, according to Eq. (8), to an intensity such that the stimulated emission rate is equal to the spontaneous emission rate. At that point the beam is said to reach its saturation intensity I_{sat} , which is given by

$$I_{\text{sat}} = \frac{h\nu_{21}}{\sigma_{21}^H \tau_{21}} \quad (15)$$

The saturation intensity is that value at which the beam can no longer grow exponentially according to Eq. (8) because there are no longer enough atoms in level 2 to provide the additional gain. When the beam grows above I_{sat} it begins to extract significant energy since at this point the stimulated emission rate exceeds the spontaneous emission rate for that transition. The beam essentially takes energy that would normally be radiated in all directions spontaneously and redirects it via stimulated emission, thereby increasing the beam intensity.

Threshold Conditions with No Mirrors

I_{sat} can be achieved by having any combination of values of the three parameters σ_{21}^H , ΔN_{21} , and L large enough that their product provides sufficient gain. It turns out that the requirement to reach I_{sat} can be given by making the exponent of Eq. (8) have the following range of values:

$$\sigma_{21}^H \Delta N_{21} L \approx 10-20 \quad \text{or} \quad g_{21}^H L \approx 10-20 \quad (16)$$

where the specific value between 10 and 20 is determined by the geometry of the laser cavity. Equation (16) suggests that the beam grows to a value of $I/I_0 = e^{10-20} = 2 \times 10^4 - 5 \times 10^8$, which is a very large amplification.

Threshold Conditions with Mirrors

One could conceivably make L sufficiently long to always satisfy Eq. (16), but this is not practical. Some lasers can reach the saturation intensity over a length L of a few centimeters, but most require much longer lengths. Since one cannot readily extend the lasing medium to be long enough to achieve I_{sat} , the same result is obtained by putting mirrors around the gain medium. This effectively increases the path length by having the beam pass many times through the amplifier.

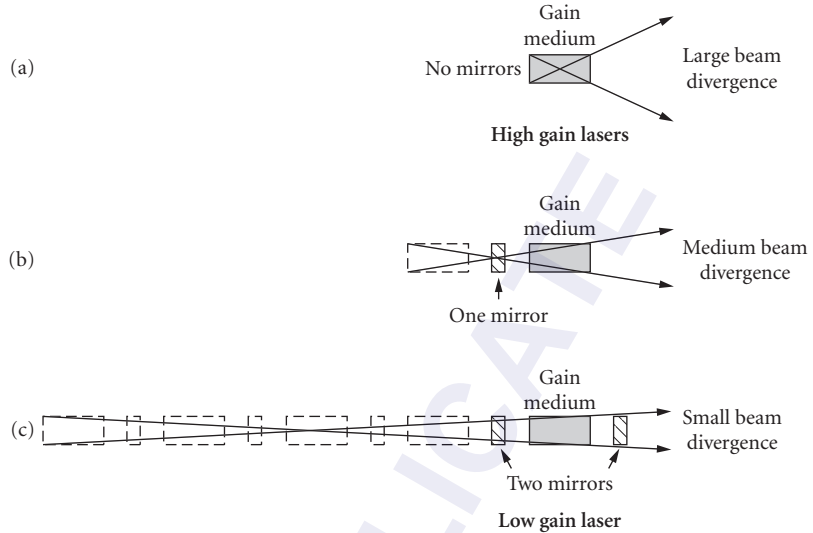


FIGURE 5 Laser beam divergence for amplifier configurations having high gain and (a) no mirrors or (b) one mirror, and high or low gain and (c) two mirrors.

A simple understanding of how the mirrors affect the beam is shown in Fig. 5. The diagram effectively shows how multiple passes through the amplifier can be considered for the situation where flat mirrors are used at the ends of the amplifier. The dashed outlines that are the same size as the gain media represent the images of those gain media produced by their reflections in the mirrors. It can be seen that as the beam makes multiple passes, its divergence narrows up significantly in addition to the large increase in intensity that occurs due to amplification. For high-gain lasers, such as excimer or organic dye lasers, the beam only need pass through the amplifier a few times to reach saturation. For low-gain lasers, such as the helium-neon laser, it might take 500 passes through the amplifier in order to reach I_{sat} .

The laser beam that emerges from the laser is usually coupled out of the amplifier by having a partially transmitting mirror at one end of the amplifier which typically reflects most of the beam back into the medium for more growth. To ensure that the beam develops, the transmission of the output coupling mirror (as it is usually referred to) must be lower than the gain incurred by the beam during a round-trip pass through the amplifier. If the transmission is higher than the round-trip gain, the beam undergoes no net amplification. It simply never develops. Thus a relationship that describes the threshold for laser oscillation balances the laser gain with the cavity losses. In the most simplified form, those losses are due to the mirrors having a reflectivity less than unity. Thus, for a round-trip pass-through the laser cavity, the threshold for inversion can be expressed as

$$R_1 R_2 e^{\sigma_{21} \Delta N_{21} 2L} > 1 \quad \text{or} \quad R_1 R_2 e^{g_{21} 2L} > 1 \quad (17a)$$

which is a similar requirement to that of Eq. (16) since it defines the minimum gain requirements for a laser.

A more general version of Eq. (17a) can be expressed as

$$R_1 R_2 (1 - a_L)^2 e^{(\sigma_{21} \Delta N_{21} - \alpha) 2L} = 1 \quad \text{or} \quad R_1 R_2 (1 - a_L)^2 e^{(g_{21} - \alpha) 2L} = 1 \quad (17b)$$

in which we have included a distributed absorption α throughout the length of the gain medium at the laser wavelength, as well as the total scattering losses a_L per pass through the cavity (excluding the gain medium and the mirror surfaces). The absorption loss α is essentially a separate absorbing transition within the gain medium that could be a separate molecule as in an excimer laser,

absorption from either the ground state or from the triplet state in a dye laser, or absorption from the ground state in the broadband tunable solid-state lasers and in semiconductor lasers or from the upper laser state in most solid-state lasers. The scattering losses a_L , per pass, would include scattering at the windows of the gain medium, such as Brewster angle windows, or scattering losses from other elements that are inserted within the cavity. These are typically of the order of one or two percent or less.

Laser Operation above Threshold

Significant power output is achieved by operating the laser at a gain greater than the threshold value defined above in Eqs. (16) and (17). For such a situation, the higher gain that would normally be produced by increased pumping is reduced to the threshold value by stimulated emission. The additional energy obtained from the reduced population inversion is transferred to the laser beam in the form of increased laser power. If the laser has low gain, as most cw (continuously operating) gas lasers do, the gain and also the power output tend to stabilize rather readily.

For solid-state lasers, which tend to have higher gain and also a long upper-laser-level lifetime, a phenomenon known as *relaxation oscillations*⁷ occurs in the laser output. For pulsed (non-Q-switched) lasers in which the gain lasts for many microseconds, these oscillations occur in the form of a regularly repeated spiked laser output superimposed on a lower steady-state value. For cw lasers it takes the form of a sinusoidal oscillation of the output. The phenomenon is caused by an oscillation of the gain due to the interchange of pumped energy between the upper laser level and the laser field in the cavity. This effect can be controlled by using an active feedback mechanism, in the form of an intensity-dependent loss, in the laser cavity.

Laser mirrors not only provide the additional length required for the laser beam to reach I_{sat} , but they also provide very important resonant cavity effects that will be discussed in a later section. Using mirrors at the ends of the laser gain medium (or amplifier) is referred to as having the gain medium located within an optical cavity.

How Population Inversions Are Achieved⁴

It was mentioned earlier that population inversions are not easily achieved in normal situations. All types of matter tend to be driven toward thermal equilibrium. From an energy-level standpoint, to be in thermal equilibrium implies that the ratio of the populations of two excited states of a particular material, whether it be a gaseous, liquid, or solid material, is described by Eq. (12). For any finite value of the temperature this leads to a value of $N_2/(g_2/g_1)N_1$ that is always less than unity and therefore Eq. (14) can never be satisfied under conditions of thermal equilibrium. Population inversions are therefore produced in either one of two ways: (1) selective excitation (pumping) of the upper-laser-level 2, or (2) more rapid decay of the population of the lower-laser-level 1 than of the upper-laser-level 2, even if they are both populated by the same pumping process.

The first requirement mentioned above was met in producing the very first laser, the ruby laser.⁵ In this laser the flash lamp selectively pumped chromium atoms to the upper laser level (through an intermediate level) until the ground state (lower laser level) was depleted enough to produce the inversion (Fig. 6). Another laser that uses this selective pumping process is the copper vapor laser⁶ (CVL). In this case, electrons in a gaseous discharge containing the copper vapor have a much preferred probability of pumping the upper laser level than the lower laser level (Fig. 7). Both of these lasers involve essentially three levels.

The second type of excitation is used for most solid-state lasers, such as the Nd³⁺ doped yttrium aluminum garnet laser⁷ (commonly referred to as the Nd:YAG laser), for organic dye lasers,⁸ and many others. It is probably the most common mechanism used to achieve the necessary population inversion. This process involves four level² (although it can include more) and generally occurs via excitation from the ground state 0 to an excited state 3 which energetically lies above the upper-laser-level 2. The population then decays from level 3 to level 2 by nonradiative processes (such as

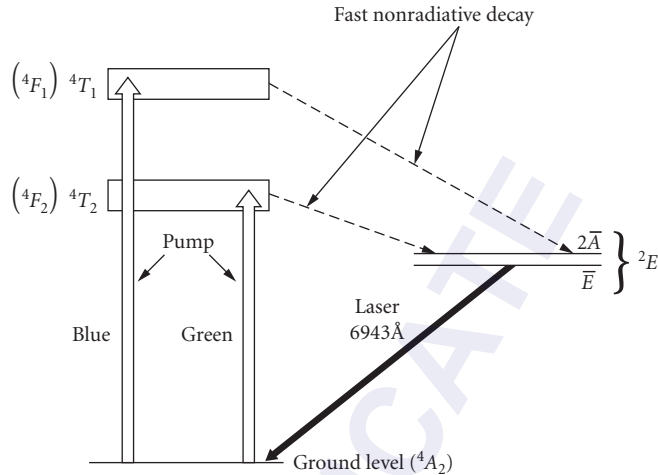


FIGURE 6 Energy-level diagram for a ruby laser showing the pump wavelength bands and the laser transition. The symbols 4T_1 and 4T_2 are shown as the appropriate designations for the pumping levels in ruby along with the more traditional designations in parenthesis.

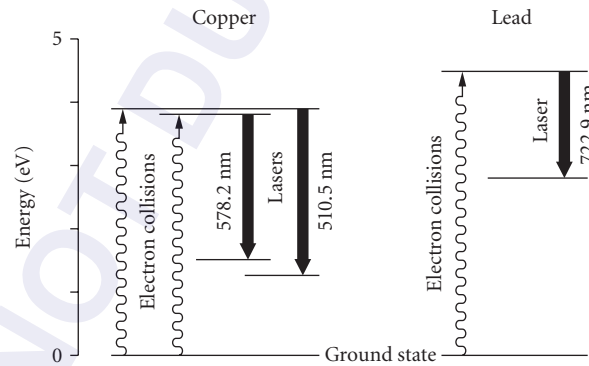


FIGURE 7 Energy-level diagrams of the transient three-level copper and lead vapor lasers showing the pump processes as well as the laser transitions.

collisional processes in gases), but can also decay radiatively to level 2. If the proper choice of materials has been made, the lower-laser-level 1 in some systems will decay rapidly to the ground state (level 0) which allows the condition $N_2 > (g_2/g_1)N_1$ to be satisfied. This situation is shown in Fig. 8 for an Nd:YAG laser crystal.

Optimization of the Output Coupling from a Laser Cavity⁹

A laser will operate with any combination of mirror reflectivities subject to the constraints of the threshold condition of Eq. (17a or b). However, since lasers are devices that are designed to use the

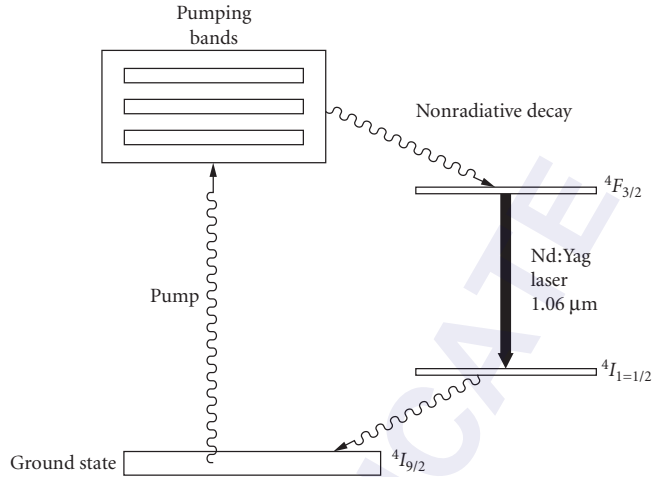


FIGURE 8 Energy-level diagram of the Nd:YAG laser indicating the four-level laser excitation process.

laser power in various applications, it is desirable to extract the most power from the laser in the most efficient manner. A simple expression for the optimum laser output coupling is given as

$$t_{\text{opt}} = (a_L g_{21} L)^{1/2} - a_L \quad (18)$$

in which a_L is the absorption and scattering losses per pass through the amplifier (the same as in Eq. (17a and b)), g_{21} is the small signal gain per pass through the amplifier, and L is the length of the gain medium. This value of t_{opt} is obtained by assuming that equal output couplings are used for both mirrors at the ends of the cavity. To obtain all of the power from one end of the laser, the output transmission must be doubled for one mirror and the other mirror is made to be a high reflector.

The intensity of the beam that would be emitted from the output mirror can be estimated to be

$$I_{t_{\text{max}}} = \frac{I_{\text{sat}} t_{\text{opt}}^2}{2a_L} \quad (19)$$

in which I_{sat} is the saturation intensity as obtained from Eq. (15). If all of the power is desired from one end of the laser, as discussed above, then $I_{t_{\text{max}}}$ would be doubled in the above expression.

Pumping Techniques to Produce Inversions

Excitation or pumping of the upper laser level generally occurs by two techniques: (1) particle pumping and (2) optical or photon pumping. No matter which process is used, the goal is to achieve sufficient pumping flux and, consequently enough, population in the upper-laser-level 2 to exceed the requirements of either Eq. (16) or Eq. (17).

Particle Pumping²⁰ Particle pumping occurs when a high-speed particle collides with a laser species and converts its kinetic energy to internal energy of the laser species. Particle pumping occurs mostly with electrons as the pumping particles. This is especially common in a gas discharge where a voltage is applied across a low-pressure gas and the electrons flow through the tube in the form of a discharge current that can range from a few milliamps to tens of amperes, depending upon the particular laser and the power level desired. This type of excitation process is used for lasers such as the argon (Fig. 9) and krypton ion lasers, the copper vapor laser, excimer lasers, and the molecular nitrogen laser.

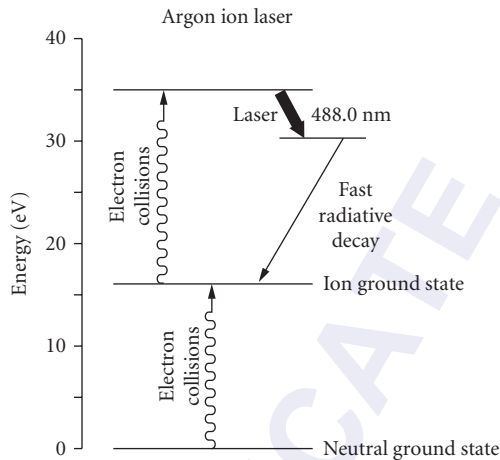


FIGURE 9 Energy-level diagram of the argon ion laser indicating the two-step excitation process.

Two other well-known gas lasers, the helium-neon laser and the helium-cadmium laser, operate in a gas discharge containing a mixture of helium gas and the laser species (neon gas or cadmium vapor). When an electric current is produced within the discharge, the high-speed electrons first pump an excited metastable state in helium (essentially a storage reservoir). The energy is then transferred from this reservoir to the upper laser levels of neon or cadmium by collisions of the helium metastable level with the neon or cadmium ground-state atoms as shown in Fig. 10. Electron collisions with the cadmium ion ground state have also been shown to produce excitation in the case of the helium-cadmium laser.

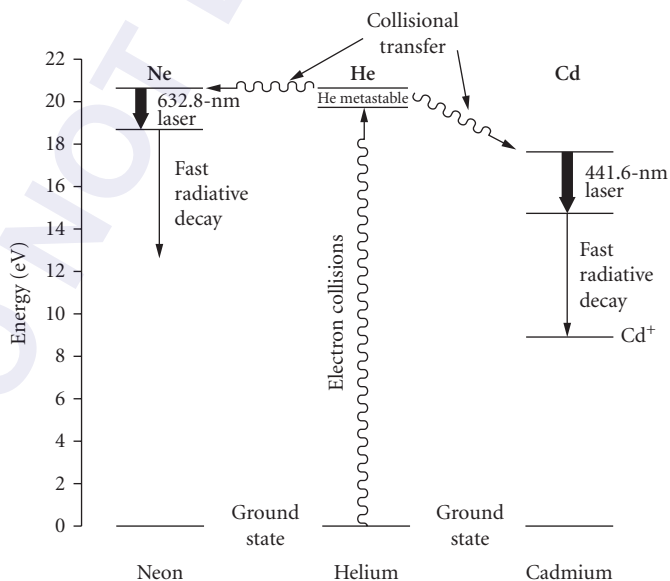


FIGURE 10 Energy-level diagrams of the helium-neon (He-Ne) laser and the helium-cadmium (He-Cd) laser that also include the helium metastable energy levels that transfer their energy to the upper laser levels by collisions.

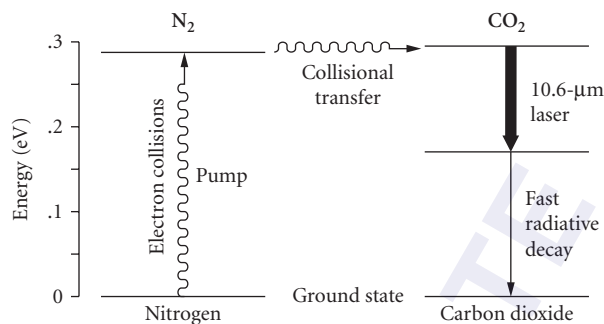


FIGURE 11 Energy-level diagram of the carbon dioxide (CO_2) laser along with the energy level of molecular nitrogen that collisionally transfers its energy to the CO_2 upper laser level.

In an energy-transfer process similar to that of helium with neon or cadmium, the CO_2 laser operates by using electrons of the gas discharge to produce excitation of molecular nitrogen vibrational levels that subsequently transfer their energy to the vibrational upper laser levels of CO_2 as indicated in Fig. 11. Helium is used in the CO_2 laser to control the electron temperature and also to cool (reduce) the population of the lower laser level via collisions of helium atoms with CO_2 atoms in the lower laser level.

High-energy electron beams and even nuclear reactor particles have also been used for particle pumping of lasers, but such techniques are not normally used in commercial laser devices.

Optical Pumping⁷ Optical pumping involves the process of focusing light into the gain medium at the appropriate wavelength such that the gain medium will absorb most (or all) of the light and thereby pump that energy into the upper laser level as shown in Fig. 12. The selectivity in pumping the laser level with an optical pumping process is determined by choosing a gain medium having significant absorption at a wavelength at which a suitable pump light source is available. This of course implies that the absorbing wavelengths provide efficient pumping pathways to the upper laser level. Optical pumping requires very intense pumping light sources, including flash lamps and other lasers. Lasers that are produced by optical pumping include organic dye lasers and solid-state lasers. The two types of energy level arrangements for producing lasers via optical pumping were described in detail in the section “How Population Inversions Are Achieved” and shown in Figs. 6 and 8.

Flash lamps used in optically pumped laser systems are typically long, cylindrically shaped, fused quartz structures of a few millimeters to a few centimeters in diameter and 10 to 50 cm in

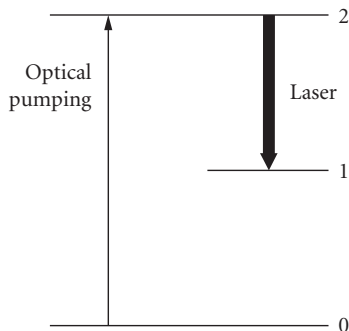


FIGURE 12 A general diagram showing optical pumping of the upper laser level.

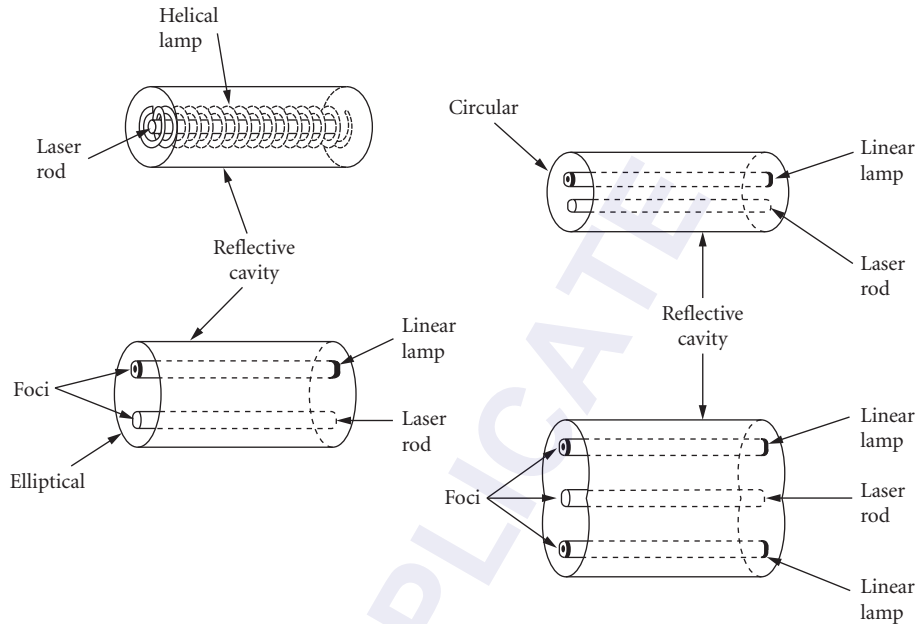


FIGURE 13 Flash lamp pumping arrangements for solid-state laser rods showing the use of helical lamps as well as linear lamps in a circular cavity, a single elliptical cavity, and a double elliptical cavity.

length. The lamps are filled with gases such as xenon and are initiated by running an electrical current through the gas. The light is concentrated into the lasing medium by using elliptically shaped reflecting cavities that surround both the laser medium and the flash lamps as shown in Fig. 13. These cavities efficiently collect and transfer to the laser rod most of the lamp energy in wavelengths within the pump absorbing band of the rod. The most common flash lamp-pumped lasers are Nd:YAG, Nd:glass, and organic dye lasers. New crystals such as Cr:LiSAF and HoTm:YAG are also amenable to flash lamp pumping.

Solid-state lasers also use energy-transfer processes as part of the pumping sequence in a way similar to that of the He-Ne and He-Cd gas lasers. For example, Cr^{3+} ions are added into neodymium-doped crystals to improve the absorption of the pumping light. The energy is subsequently transferred to the Nd^{3+} laser species. Such a process of adding desirable impurities is known as *sensitizing*.

Lasers are used as optical pumping sources in situations where (1) it is desirable to be able to concentrate the pump energy into a small-gain region or (2) it is useful to have a narrow spectral output of the pump source in contrast with the broadband spectral distribution of a flash lamp.

Laser pumping is achieved by either transverse pumping (a direction perpendicular to the direction of the laser beam) or longitudinal pumping (a direction in the same direction as the emerging laser beam). Frequency doubled and tripled pulsed Nd:YAG lasers are used to transversely pump organic dye lasers⁸ that provide continuously tunable laser radiation over the near-ultraviolet, visible, and near-infrared spectral regions (by changing dyes at appropriate wavelength intervals). For transverse pumping, the pump lasers are typically focused into the dye medium with a cylindrical lens to provide a 1- to 2-cm-long (but very narrow) gain medium in the liquid dye solution. The dye concentration is adjusted to absorb the pump light within a millimeter or so into the dye cell to provide the very high concentration of gain near the surface of the cell.

Both cw and mode-locked argon ion lasers and Nd:YAG lasers are typically used for longitudinal or end pumping of cw and/or mode-locked organic dye lasers and also of solid-state gain media. In this pumping arrangement the pump laser is focused into a very thin gain region, which is provided by either a thin jet stream of flowing dye solution (Fig. 14) or a solid-state crystal such as $\text{Ti:Al}_2\text{O}_3$,

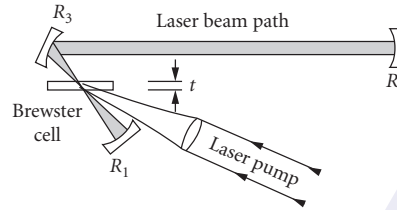


FIGURE 14 End-pumping cavity arrangement for either organic dye lasers or solid-state lasers.

The thin gain medium is used, in the case of the generation of ultrashort mode-locked pulses, so as to allow precise timing of the short-duration pump pulses with the ultrashort laser pulses that develop within the gain medium as they travel within the optical cavity.

Thin-disk lasers are diode-pumped solid-state lasers that efficiently produce high output power with good beam quality. Such lasers have gain media with a very short axial dimension of several hundred microns and wide transverse dimensions of several centimeters. The disks are antireflection coated for both the lasing and pumping wavelengths on the front side and a high reflection coating on the rear side. The laser crystal is mounted on a heat sink to efficiently remove the wasted heat of the pumping process. An output coupling mirror is mounted in front of the disk to provide multiple passes of the beam through the gain medium, for maximum power extraction, and to provide good mode quality. Yb:YAG is the most successful laser material for this type of laser.

Gallium arsenide semiconductor diode lasers, operating at wavelengths around $0.8 \mu\text{m}$, can be effectively used to pump Nd:YAG lasers because the laser wavelength is near that of the strongest absorption feature of the pump band of the Nd:YAG laser crystal, thereby minimizing excess heating of the laser medium. Also, the diode lasers are very efficient light sources that can be precisely focused into the desired mode volume of a small Nd:YAG crystal (Fig. 15a) which results in minimal waste of the pump light. Figure 15b shows how close-coupling of the pump laser and the Nd:YAG crystal can be used to provide compact efficient diode laser pumping. The infrared output of the

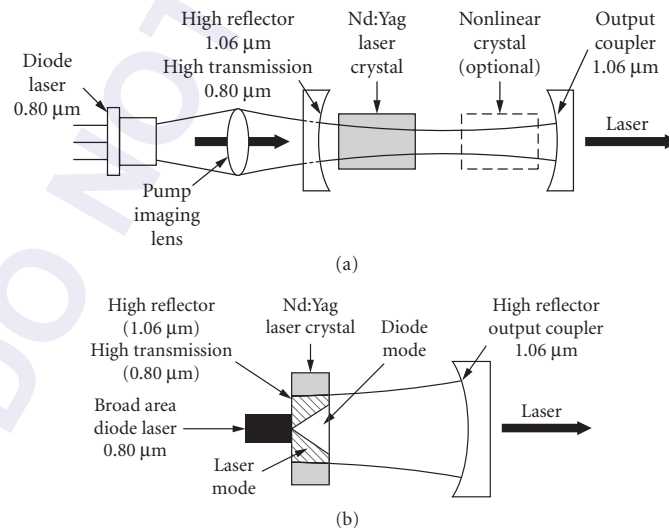


FIGURE 15 Pumping arrangements for diode-pumped Nd:YAG lasers showing both (a) standard pumping using an imaging lens and (b) close-coupled pumping in which the diode is located adjacent to the laser crystal.

Nd:YAG laser can then also be frequency doubled as indicated in Fig. 15a, using nonlinear optical techniques to produce a green laser beam in a relatively compact package.

Semiconductor Diode Laser Pumping Semiconductor laser pumping occurs when electrons are made to flow from an n -type semiconductor to a p -type semiconductor. In this case, as opposed to particle pumping described previously, it is not the kinetic energy of the electrons that does the excitation. Instead, it is the electrons themselves flowing into a p -doped material that produces the inversion. An analogy might be the water in a mountain region approaching a waterfall. The water is already at the upper energy site and loses its energy when it cascades down the waterfall. In the same sense, the electrons already have sufficient potential energy when they are pulled into the p -type material via an external electric field. Once they arrive, a population inversion exists where they recombine with the holes and cascade downward to produce the recombination radiation.

16.4 LASER PROPERTIES ASSOCIATED WITH OPTICAL CAVITIES OR RESONATORS

Longitudinal Laser Modes^{1,2}

When a collimated optical beam of infinite lateral extent (a plane wave) passes through two reflecting surfaces of reflectivity R and also of infinite extent that are placed normal (or nearly normal) to the beam and separated by a distance d , as shown in Fig. 16a, the plot of transmission versus wavelength for the light as it emerges from the second reflecting surface is shown in Fig. 16b. The transmission reaches a maximum of 100 percent (if there are no absorption losses at the reflecting surfaces) at frequency spacings of $\Delta\nu = c/2\eta d$ where c is the speed of light in a vacuum, and η is the index of refraction of the medium between the mirrors. This frequency-selective optical device is known as a Fabry-Perot interferometer and has many useful applications in optics.

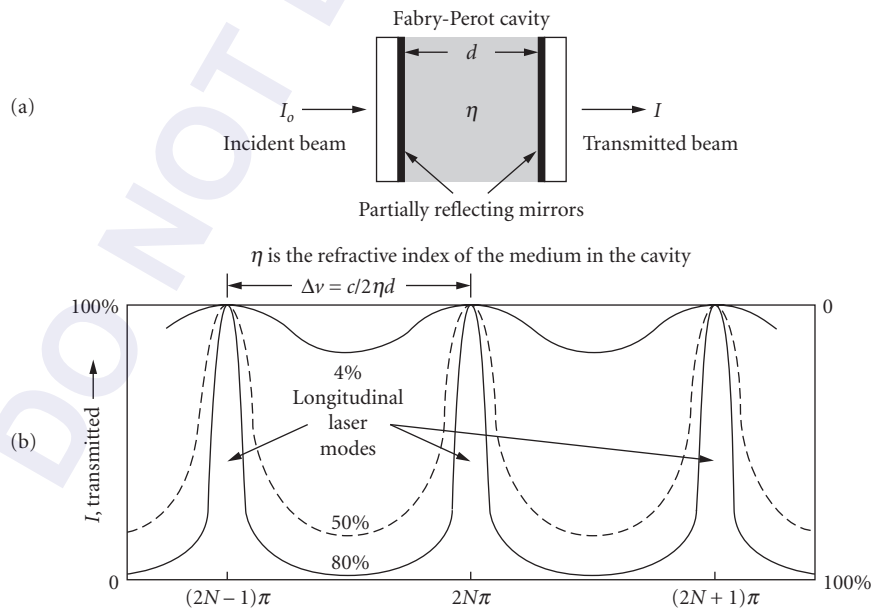


FIGURE 16 Fabry-Perot optical cavity consisting of two plane-parallel mirrors with a specific reflectivity separated by a distance d indicating the frequency spacing of longitudinal modes.

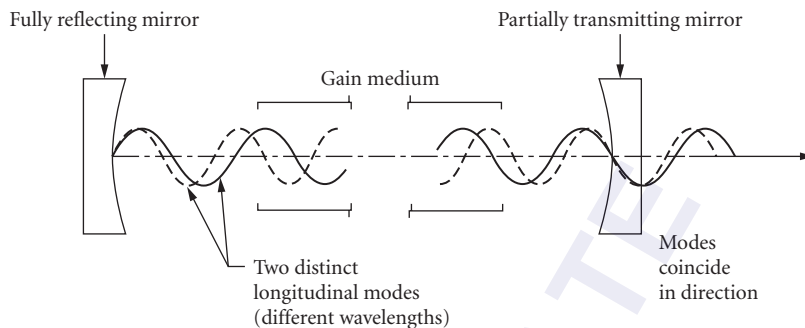


FIGURE 17 Laser resonator showing two distinct longitudinal modes traveling in the same direction but with slightly different frequencies.

The transmission through this device, as shown in Fig. 16*b* is enhanced at regular frequency or wavelength intervals due to the development of standing waves that resonate within the optical cavity. The enhancement occurs at frequencies (or wavelengths) at which complete sinusoidal half-cycles of the electromagnetic wave exactly “fit” between the mirrors such that the value of the electric field of the wave is zero at the mirror surfaces.

If a laser amplifier is placed between two mirrors in the same arrangement as described above, the same standing waves tend to be enhanced at frequency intervals of

$$\nu = n(c/2\eta d) \quad (20)$$

where n is an integer that indicates the number of half wavelengths of the laser light that fit within the spacing d of the two mirrors. In a typical laser operating in the visible spectral region, n would be of the order of 30,000 to 40,000. In such a laser, the output of the laser beam emerging from the cavity is very strongly enhanced at the resonant wavelengths as shown in Fig. 17, since these are the wavelengths that have the lowest loss within the cavity. The widths of the resonances shown in Fig. 17 are those of a passive Fabry-Perot cavity. When an active gain medium is placed within the cavity, the linewidth of the beam that is continually amplified as it reflects back and forth between the mirrors is narrowed even further.

These enhanced regions of very narrow frequency output are known as *longitudinal modes* of the laser. They are referred to as modes since they represent discrete values of frequency associated with the integral values of n at which laser output occurs. Lasers operating on a single longitudinal mode with ultrastable cavities and ultrahigh reflectivity mirrors have generated linewidths as narrow as a few hundred hertz or less. Since the longitudinal or temporal coherence length of a beam of light is determined by $c/\Delta\nu$, a very narrow laser linewidth can provide an extremely long coherence length and thus a very coherent beam of light.

For a typical gas laser (not including excimer lasers), the laser gain bandwidth is of the order of 10^9 to 10^{10} Hz. Thus, for a laser mirror cavity length of 0.5 m, the mode spacing would be of the order of 300 MHz and there would be anywhere from 3 to 30 longitudinal modes operating within the laser cavity. For an organic dye laser or a broadband solid-state laser, such as a Ti:Al₂O₃ laser, there could be as many as one million distinct longitudinal modes, each of a slightly different frequency than the next one, oscillating at the same time. However, if mode-locking is not present, typically only one or a few modes will dominate the laser output of a homogeneously broadened laser gain medium.

Transverse Laser Modes¹¹⁻¹⁴

The previous section considered the implications of having a collimated or parallel beam of light of infinite lateral extent pass through two infinite reflecting surfaces that are arranged normal to the direction of propagation of the beam and separated by a specific distance d . We must now

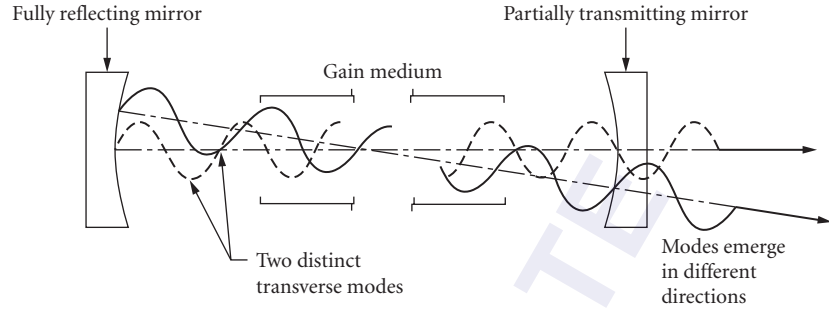


FIGURE 18 Laser resonator showing two distinct transverse modes traveling in different directions with slightly different frequencies.

consider the consequences of having the light originate within the space between the two mirrors in an amplifier that has only a very narrow lateral extent limited by either the diameter of the mirrors or by the diameter of the amplifying medium. The beam evolves from the spontaneous emission within the gain medium and eventually becomes a nearly collimated beam when it reaches I_{sat} since only rays traveling in a very limited range of directions normal to the laser mirrors will experience enough reflections to reach I_{sat} . The fact that the beam has a restricted aperture in the direction transverse to the direction of propagation causes it to evolve with a slight transverse component due to diffraction, which effectively causes the beam to diverge. This slight divergence actually consists of one or more distinctly separate beams that can operate individually or in combination.

These separate beams that propagate in the z direction are referred to as *transverse modes* as shown in Fig. 18. They are characterized by the various lateral spatial distributions of the electric field vector $\mathbf{E}(x, y)$ in the x - y directions as they emerge from the laser. These transverse amplitude distributions for the waves can be described by the relationship.

$$E_{pq}(x, y) = H_p\left(\frac{\sqrt{2}x}{w}\right) H_q\left(\frac{\sqrt{2}y}{w}\right) e^{-(x^2+y^2)/w^2} \quad (21)$$

In this solution, p and q are positive integers ranging from zero to infinity that designate the different modes which are associated with the order of the Hermite polynomials. Thus every set of p, q represents a specific distribution of wave amplitude at one of the mirrors, or a specific transverse mode of the open-walled cavity. We can list several Hermite polynomials as follows:

$$\begin{aligned} H_0(u) &= 1 & H_1(u) &= 2u \\ H_2(u) &= 2(2u^2 - 1) \\ H_m(u) &= (-1)^m e^{u^2} \frac{d^m(e^{-u^2})}{du^m} \end{aligned} \quad (22)$$

The spatial intensity distribution would be obtained by squaring the amplitude distribution function of Eq. (21). The transverse modes are designated TEM for transverse electromagnetic. The lowest order mode is given by TEM_{00} . It could also be written as TEM_{n00} in which n would designate the longitudinal mode number [Eq. (20)]. Since this number is generally very large for optical frequencies, it is not normally given.

The lowest order TEM_{00} mode has a circular distribution with a gaussian shape (often referred to as the *gaussian mode*) and has the smallest divergence of any of the transverse modes. Such a mode can be focused to a spot size with dimensions of the order of the wavelength of the beam. It has a minimum width or waist $2w_0$ that is typically located between the laser mirrors (determined

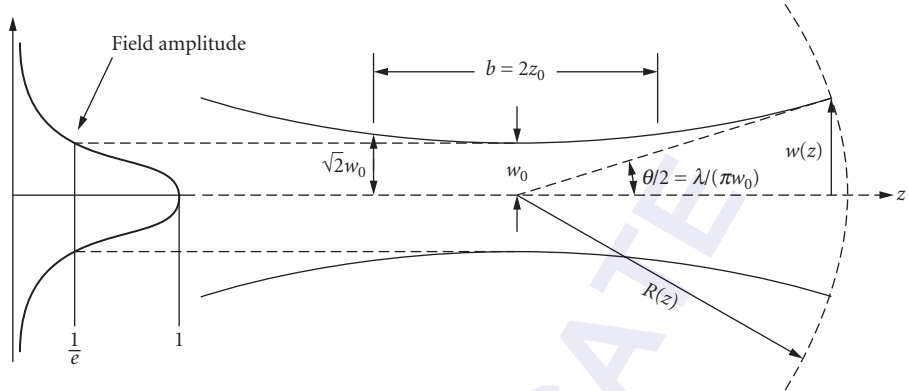


FIGURE 19 Parameters of a gaussian-shaped beam which are the features of a TEM_{00} transverse laser mode.

by the mirror curvatures and separation) and expands symmetrically in opposite directions from that minimum waist according to the following equation:

$$w(z) = w_0 \left[1 + \left(\frac{\lambda z}{\eta \pi w_0^2} \right)^2 \right]^{1/2} = w_0 \left[1 + \left(\frac{z}{z_0} \right)^2 \right]^{1/2} \quad (23)$$

where $w(z)$ is the beam waist at any location z measured from w_0 , η is the index of refraction of the medium, and $z_0 = \eta \pi w_0^2 / \lambda$ is the distance over which the beam waist expands to a value of $\sqrt{2}w_0$.

The waist $w(z)$ at any location, as shown in Fig. 19, describes the transverse dimension within which the electric field distribution of the beam decreases to a value of 37 percent ($1/e$) of its maximum on the beam axis and within which 86.5 percent of the beam energy is contained. The TEM_{00} mode would have an intensity distribution that is proportional to the square of Eq. (21) for $p=q=0$ and, since it is symmetrical around the axis of propagation, it would have a cylindrically symmetric distribution of the form

$$I(r, z) = I_0 e^{-2r^2/w^2(z)} \quad (24)$$

where I_0 is the intensity on the beam axis.

The beam also has a wavefront curvature given by

$$R(z) = z \left[1 + \left(\frac{\eta \pi w_0^2}{\lambda z} \right)^2 \right] \quad (25)$$

which is indicated in Fig. 19, and a far-field angular divergence given by

$$\theta = \lim_{z \rightarrow \infty} \frac{2w(z)}{z} = \frac{2\lambda}{\pi w_0} = 0.64 \frac{\lambda}{w_0} \quad (26)$$

For a symmetrical cavity formed by two mirrors, each of radius of curvature R , separated by a distance d and in a medium in which $\eta=1$ the minimum beam waist w_0 is given by

$$w_0^2 = \frac{\lambda}{2\pi} [d(2R-d)]^{1/2} \quad (27)$$

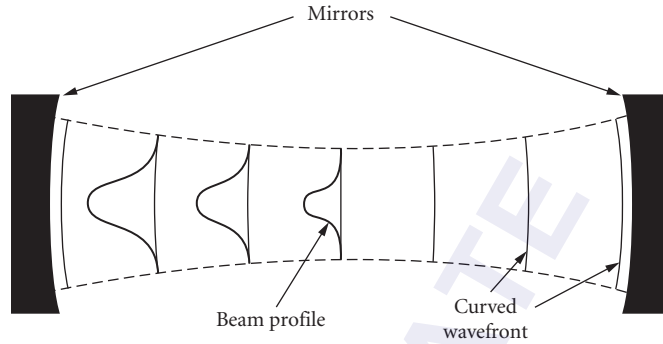


FIGURE 20 A stable laser resonator indicating the beam profile at various locations along the beam axis as well as the wavefront at the mirrors that matches the curvature of the mirrors.

and the radius of curvature r_c of the wavefront is

$$r_c = z + \frac{d(2R-d)}{4z} \quad (28)$$

For a confocal resonator in which $R=d$, w_o is given by

$$w_o = \sqrt{\frac{\lambda d}{2\pi}} \quad (29)$$

and the beam waist (spot size) at each mirror located a distance $d/2$ from the minimum is

$$w = \sqrt{\frac{\lambda d}{\pi}} \quad (30)$$

Thus, for a confocal resonator, $w(d/2)$ at each of the mirrors is equal to $\sqrt{2}w_o$ and thus at that location, $z = z_o$. The distance between mirrors for such a cavity configuration is referred to as the *confocal parameter* b such that $b = 2z_o$. In a stable resonator, the curvature of the wavefront at the mirrors, according to Eqs. (25) and (28), exactly matches the curvature of the mirrors as shown in Fig. 20.

Each individual transverse mode of the beam is produced by traveling a specific path between the laser mirrors such that, as it passes from one mirror to the other and returns, the gain it receives from the amplifier is at least as great as the total losses of the mirror, as indicated from Eq. (17), plus the additional diffraction losses produced by either the finite lateral extent of the laser mirrors or the finite diameter of the laser amplifier or some other optical aperture placed in the system, whichever is smaller. Thus the TEM_{00} mode is produced by a beam passing straight down the axis of the resonator as indicated in Fig. 19.

Laser Resonator Configurations and Cavity Stability^{4,14}

There are a variety of resonator configurations that can be used for lasers. The use of slightly curved mirrors leads to much lower diffraction losses of the transverse modes than do plane parallel mirrors, and they also have much less stringent alignment tolerances. Therefore, most lasers use curved mirrors for the optical cavity. For a cavity with two mirrors of curvature R_1 and R_2 , and a separation distance d , a number of possible cavity configurations are shown in Fig. 21.

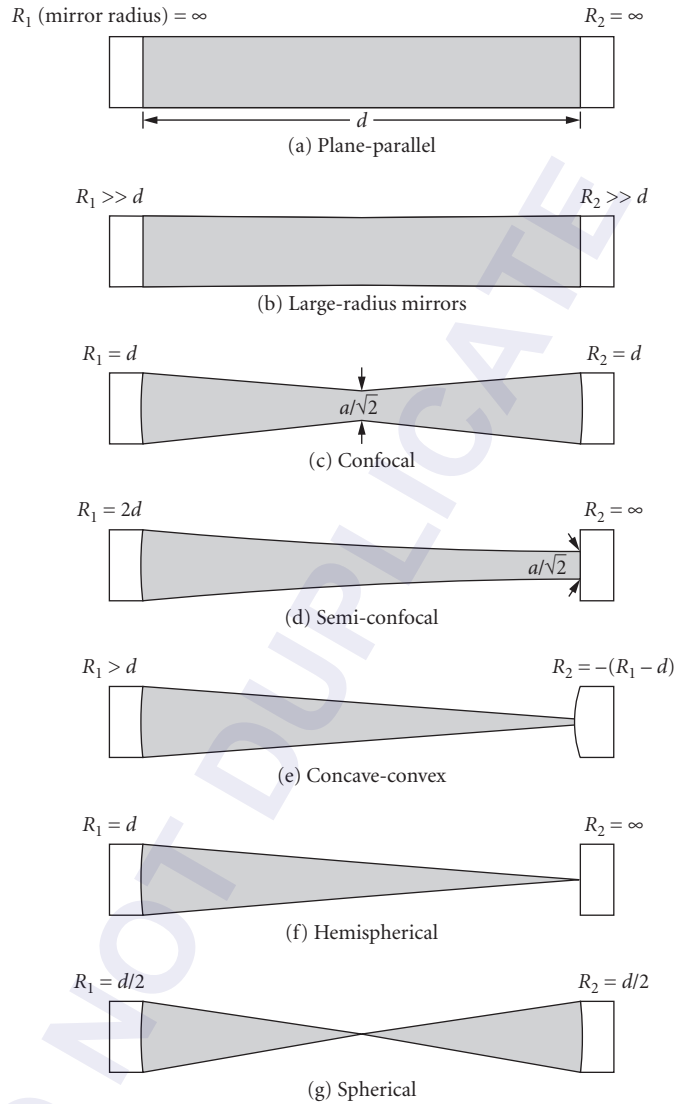


FIGURE 21 Possible two-mirror laser cavity configurations indicating the relationship of the radii of curvature of the mirrors with respect to the separation between mirrors.

A relationship between the radii of curvature and the separation between mirrors can be defined as

$$g_1 = 1 - \frac{d}{R_1} \quad \text{and} \quad g_2 = 1 - \frac{d}{R_2} \quad (31)$$

such that the condition for stable transverse modes is given by

$$0 < g_1 g_2 < 1 \quad (32)$$

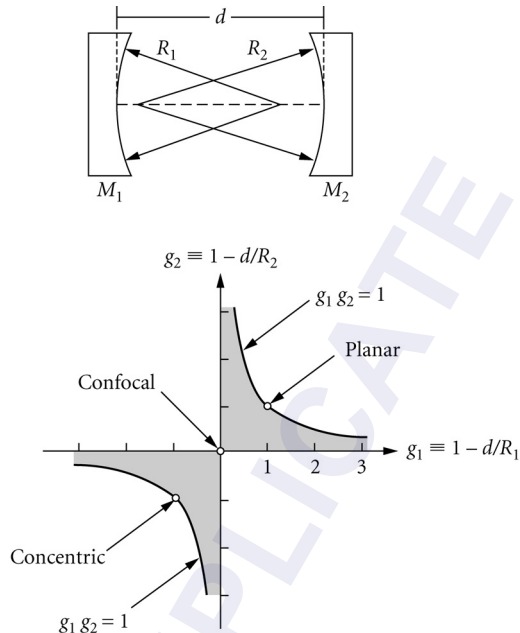


FIGURE 22 Stability diagram for two-mirror laser cavities indicating the shaded regions where stable cavities exist.

A stable mode is a beam that can be maintained with a steady output and profile over a relatively long period of time. It results from a cavity configuration that concentrates the beam toward the resonator axis in a regular pattern as it traverses back and forth within the cavity, rather than allowing it to diverge and escape from the resonator. In considering the various possible combinations of curved mirror cavities, one must keep the relation between the curvatures and mirror separation d within the stable regions of the graph, as shown in Fig. 22, in order to produce stable modes. Thus it can be seen from Fig. 22 that not all configurations shown in Fig. 21 are stable. For example, the planar, confocal, and concentric arrangements are just on the edge of stability according to Fig. 22.

There may be several transverse modes oscillating simultaneously “within” a single longitudinal mode. Each transverse mode can have the same value of n [Eq. (20)] but will have a slightly different value of d as it travels a different optical path between the resonator mirrors, thereby generating a slight frequency shift from that of an adjacent transverse mode. For most optical cavities, the mode that operates most easily is the TEM_{00} mode, since it travels a direct path along the axis of the gain medium.

16.5 SPECIAL LASER CAVITIES

Unstable Resonators¹⁴

A laser that is operating in a TEM_{00} mode, as outlined above, typically has a beam within the laser cavity that is relatively narrow in width compared to the cavity length. Thus, if a laser with a relatively wide gain region is used, to obtain more energy in the output beam, it is not possible to extract the energy from that entire region from a typically very narrow low-order gaussian mode.

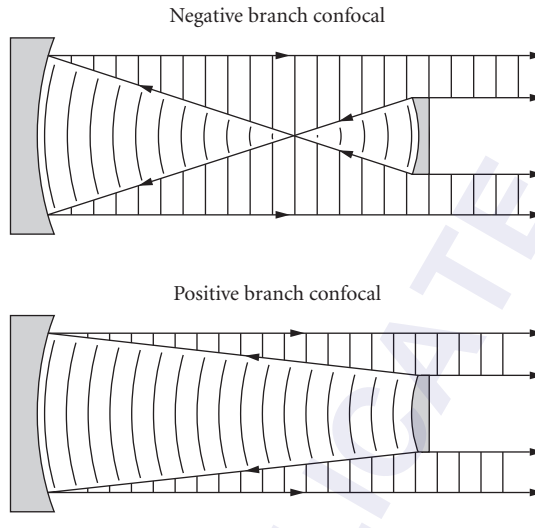


FIGURE 23 Unstable resonator cavity configurations showing both the negative and positive branch confocal cavities.

A class of resonators has been developed that can extract the energy from such wide laser volumes and also produce a beam with a nearly gaussian profile that makes it easily focusable. These resonators do not meet the criteria for stability, as outlined above, but still provide a good beam quality for some types of lasers. This class of resonators is referred to as *unstable resonators*.

Unstable resonators are typically used with high-gain laser media under conditions such that only a few passes through the amplifier will allow the beam to reach the saturation intensity and thus extract useful energy. A diagram of two unstable resonator cavity configurations is shown in Fig. 23. Figure 23b is the positive branch confocal geometry and is one of the most common unstable resonator configurations. In this arrangement, the small mirror has a convex shape and the large mirror a concave shape with a separation of length d such that $R_2 - R_1 = 2d$. With this configuration, any ray traveling parallel to the axis from left to right that intercepts the small convex mirror will diverge toward the large convex mirror as though it came from the focus of that mirror. The beam then reflects off of the larger mirror and continues to the right as a beam parallel to the axis. It emerges as a reasonably well-collimated beam with a hole in the center (due to the obscuration of the small mirror). The beam is designed to reach the saturation intensity when it arrives at the large mirror and will therefore proceed to extract energy as it makes its final pass through the amplifier. In the far field, the beam is near gaussian in shape, which allows it to be propagated and focused according to the equations described in the previous section. A number of different unstable resonator configurations can be found in the literature for specialized applications.

Q-Switching¹⁵

A typical laser, after the pumping or excitation is first applied, will reach the saturation intensity in a time period ranging from approximately 10 ns to 1 μ s, depending upon the value of the gain in the medium. For lasers such as solid-state lasers the upper-laser-level lifetime is considerably longer than this time (typically 50 to 200 μ s). It is possible to store and accumulate energy in the form of population in the upper laser level for a time duration of the order of that upper-level lifetime. If the laser cavity could be obscured during this pumping time and then suddenly switched into the system at a time order of the upper-laser-level lifetime, it would be possible for the gain, as well as

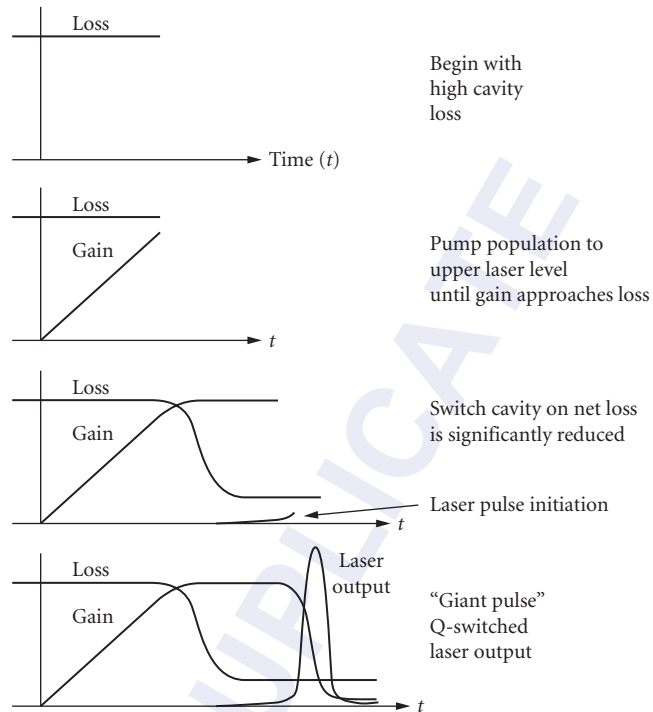


FIGURE 24 Schematic diagram of a Q-switched laser indicating how the loss is switched out of the cavity and a giant laser pulse is produced as the gain builds up.

the laser energy, to reach a much larger value than it normally would under steady-state conditions. This would produce a "giant" laser pulse of energy from the system.

Such a technique can in fact be realized and is referred to as *Q-switching* to suggest that the cavity *Q* is changed or switched into place. The cavity is switched on by using either a rapidly rotating mirror or an electro-optic shutter such as a Pockel cell or a Kerr cell. Nd:YAG and Nd:glass lasers are the most common Q-switched lasers. A diagram of the sequence of events involved in Q-switching is shown in Fig. 24.

Another technique that is similar to Q-switching is referred to as *cavity dumping*. With this technique, the intense laser beam inside of a normal laser cavity is rapidly switched out of the cavity by a device such as an acousto-optic modulator. Such a device is inserted at Brewster's angle inside the cavity and is normally transparent to the laser beam. When the device is activated, it rapidly inserts a high-reflecting surface into the cavity and reflects the beam out of the cavity. Since the beam can be as much as two orders of magnitude higher in intensity within the cavity than that which leaks through an output mirror, it is possible to extract high power on a pulsed basis with such a technique.

Mode-Locking¹⁴

In the discussion of longitudinal modes it was indicated that such modes of intense laser output occur at regularly spaced frequency intervals of $\Delta\nu = c/2\eta d$ over the gain bandwidth of the laser medium. For laser cavity lengths of the order of 10 to 100 cm these frequency intervals range from approximately 10^8 to 10^9 Hz. Under normal laser operation, specific modes with the highest gain

tend to dominate and quench other modes (especially if the gain medium is homogeneously broadened). However, under certain conditions it is possible to obtain all of the longitudinal modes lasing simultaneously, as shown in Fig. 25a. If this occurs, and the modes are all phased together so that they can act in concert by constructively and destructively interfering with each other, it is possible to produce a series of giant pulses separated in a time Δt of

$$\Delta t = 2\eta d/c \tag{33}$$

or approximately 1 to 10 ns for the cavity lengths mentioned above (see Fig. 25b). The pulse duration is approximately the reciprocal of the separation between the two extreme longitudinal laser modes or

$$\Delta t_p = \frac{1}{n\Delta\nu} \tag{34}$$

as can also be seen in Fig. 25b. This pulse duration is approximately the reciprocal of the laser gain bandwidth. However, if the index of refraction varies significantly over the gain bandwidth, then all of the frequencies are not equally spaced and the mode-locked pulse duration will occur only within the frequency width over which the frequency separations are approximately the same.

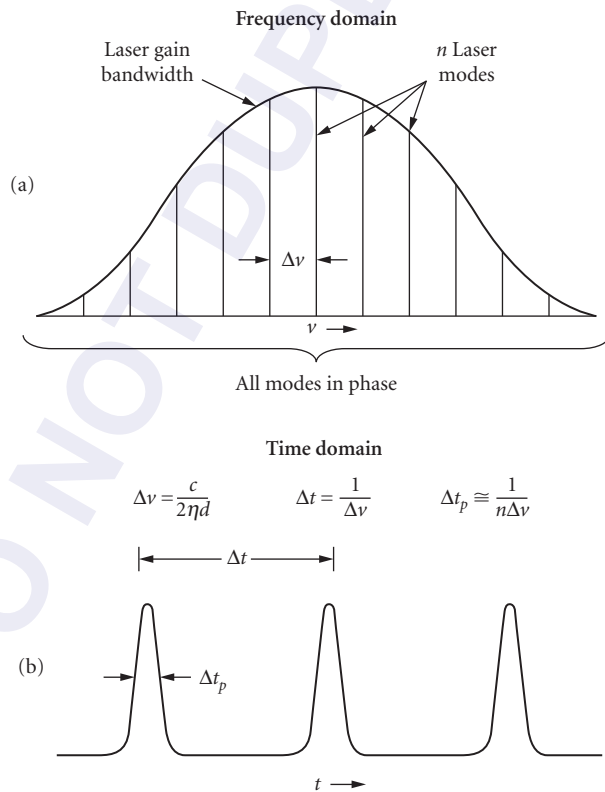


FIGURE 25 Diagrams in both the frequency and time domains of how mode-locking is produced by phasing n longitudinal modes together to produce an ultrashort laser pulse.

The narrowest pulses are produced from lasers having the largest gain bandwidth such as organic dye lasers, and solid-state lasers including Ti:Al₂O₃, Cr:Al₂O₃, and Cr:LiSAIF. The shortest mode-locked pulses to date, 4.4×10^{-15} seconds, have been produced in Ti:sapphire at a center wavelength of 820 nm at a rep rate of 80 MHz. Using pulse compression techniques, such pulses have been made as short as 3.5×10^{-15} seconds.

Such short-pulse operation is referred to as *mode-locking* and is achieved by inserting a very fast shutter within the cavity which is opened and closed at the intervals of the round-trip time of the short laser pulse within the cavity. This shutter coordinates the time at which all of the modes arrive at the mirror and thus brings them all into phase at that location. Electro-optic shutters, short duration gain pumping by another mode-locked laser, or passive saturable absorbers are techniques that can serve as the fast shutter. The second technique, short-duration gain pumping, is referred to as *synchronous pumping*. Three fast saturable absorber shutter techniques for solid-state lasers include colliding pulse mode-locking,¹⁶ additive pulse mode-locking,¹⁷ and Kerr lens mode-locking.¹⁸

Extremely short soft-x-ray pulses have also been produced via the interaction of intense laser beams with atoms. These pulses have been made as short as 70 as (70×10^{-18} seconds).

Distributed Feedback Lasers¹⁹

The typical method of obtaining feedback into the laser gain medium is to have the mirrors located at the ends of the amplifier as discussed previously. It is also possible, however, to provide the reflectivity within the amplifying medium in the form of a periodic variation of either the gain or the index of refraction throughout the medium. This process is referred to as *distributed feedback* (DFB). Such feedback methods are particularly effective in semiconductor lasers in which the gain is high and the fabrication of periodic variations is not difficult. The reader is referred to the reference section at the end of this chapter for further information concerning this type of feedback.

Ring Lasers¹⁴

Ring lasers are lasers that have an optical path within the cavity that involves the beam circulating in a loop rather than passing back and forth over the same path. This requires optical cavities that have more than two mirrors. The laser beam within the cavity consists of two waves traveling in opposite directions with separate and independent resonances within the cavity. In some instances an optical device is placed within the cavity that provides a unidirectional loss. This loss suppresses one of the beams, allowing the beam propagating in the other direction to become dominant. The laser output then consists of a traveling wave instead of a standing wave and therefore there are no longitudinal modes. Such an arrangement also eliminates the variation of the gain due to the standing waves in the cavity (spatial hole burning), and thus the beam tends to be more homogeneous than that of a normal standing-wave cavity. Ring lasers are useful for producing ultrashort mode-locked pulses and also for use in laser gyroscopes as stable reference sources.

16.6 SPECIFIC TYPES OF LASERS

Lasers can be categorized in several different ways including wavelength, material type, and applications. In this section we will summarize them by material type such as gas, liquid, solid-state, and semiconductor lasers. We will include only lasers that are available commercially since such lasers now provide a very wide range of available wavelengths and powers without having to consider special laboratory lasers.

Gaseous Laser Gain Media

Helium-Neon Laser¹⁰ The helium-neon laser was the first gas laser. The most widely used laser output wavelength is a red beam at 632.8 nm with a cw output ranging from 1 to 100 mW and sizes ranging from 10 to 100 cm in length. It can also be operated at a wavelength of 543.5 nm for some specialized applications. The gain medium is produced by passing a relatively low electrical current (10 mA) through a low pressure gaseous discharge tube containing a mixture of helium and neon. In this mixture, helium metastable atoms are first excited by electron collisions as shown in Fig. 10. The energy is then collisionally transferred to neon atom excited states which serve as upper laser levels.

Argon Ion Laser¹⁰ The argon ion laser and the similar krypton ion laser operate over a wide range of wavelengths in the visible and near-ultraviolet regions of the spectrum. The wavelengths in argon that have the highest power are at 488.0 and 514.5 nm. Power outputs on these laser transitions are available up to 20 W cw in sizes ranging from 50 to 200 cm in length. The gain medium is produced by running a high electric current (many amperes) through a very low pressure argon or krypton gas. The argon atoms must be ionized to the second and third ionization stages (Fig. 9) in order to produce the population inversions. As a result, these lasers are inherently inefficient devices.

Helium-Cadmium Laser¹⁰ The helium-cadmium laser operates cw in the blue at 441.6 nm, and in the ultraviolet at 353.6 and 325.0 nm with powers ranging from 20 to 200 mW in lasers ranging from 40 to 100 cm in length. The gain medium is produced by heating cadmium metal and evaporating it into a gaseous discharge of helium where the laser gain is produced. The excitation mechanisms include Penning ionization (helium metastables collide with cadmium atoms) and electron collisional ionization within the discharge as indicated in Fig. 10. The laser uses an effect known as *cataphoresis* to transport the cadmium through the discharge and provide the uniform gain medium.

Copper Vapor Laser¹⁰ This pulsed laser provides high-average powers of up to 100 W at wavelengths of 510.5 and 578.2 nm. The copper laser and other metal vapor lasers of this class, including gold and lead lasers, typically operate at a repetition rate of up to 20 kHz with a current pulse duration of 10 to 50 ns and a laser output of 1 to 10 mJ/pulse. The copper lasers operate at temperatures in the range of 1600°C in 2- to 10-cm-diameter temperature-resistant tubes typically 100 to 150 cm in length. The lasers are self-heated such that all of the energy losses from the discharge current provide heat to bring the plasma tube to the required operating temperature. Excitation occurs by electron collisions with copper atoms vaporized in the plasma tube as indicated in Fig. 7.

Carbon-Dioxide Laser¹⁰ The CO₂ laser, operating primarily at a wavelength of 10.6 μm, is one of the most powerful lasers in the world, producing cw powers of over 100 kW and pulsed energies of up to 10 kJ. It is also available in small versions with powers of up to 100 W from a laser the size of a shoe box. CO₂ lasers typically operate in a mixture of carbon dioxide, nitrogen, and helium gases. Electron collisions excite the metastable levels in nitrogen molecules with subsequent transfer of that energy to carbon dioxide laser levels as shown in Fig. 11. The helium gas acts to keep the average electron energy high in the gas discharge region and to cool or depopulate the lower laser level. This laser is one of the most efficient lasers, with conversion from electrical energy to laser energy of up to 30 percent.

Excimer Laser²⁰ The rare gas-halide excimer lasers operate with a pulsed output primarily in the ultraviolet spectral region at 351 nm in xenon fluoride, 308 nm in xenon chloride, 248 nm in krypton fluoride, and 193 nm in argon fluoride. The laser output, with pulse durations of 10 to 50 ns, is typically of the order of 0.2 to 1.0 J/pulse at repetition rates up to several hundred hertz. The lasers are relatively efficient (1 to 5 percent) and are of a size that would fit on a desktop. The excitation occurs via electrons within the discharge colliding with and ionizing the rare gas molecules and at the same time disassociating the halogen molecules to form negative halogen ions. These two species

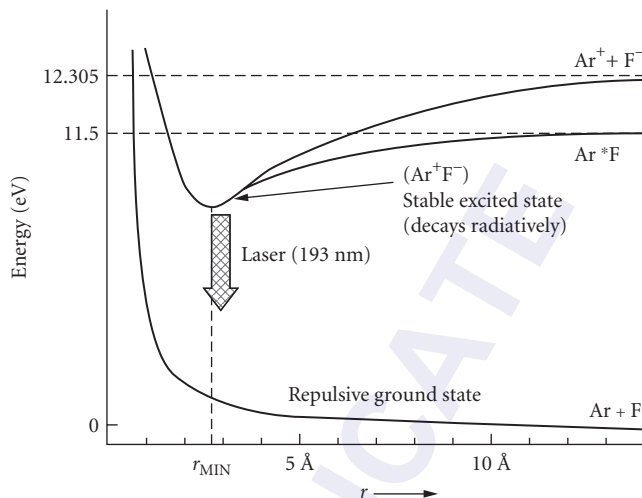


FIGURE 26 Energy-level diagram of an argon fluoride (ArF) excimer laser showing the stable excited state (upper laser level) and the unstable or repulsive ground state.

then combine to form an excited molecular state of the rare gas-halogen molecule which serves as the upper laser state. The molecule then radiates at the laser transition and the lower level advantageously disassociates since it is unstable, as shown in Fig. 26, for the ArF excimer molecule. The excited laser state is an excited state dimer which is referred to as an *excimer state*.

X-Ray Laser²¹ Laser output in the soft-x-ray spectral region has been produced in plasmas of highly ionized ions of a number of atomic species. The highly ionized ions are produced by the absorption of powerful solid-state lasers focused onto solid material of the desired atomic species. Since mirrors are not available for most of the soft-x-ray laser wavelengths (4 to 30 nm), the gain has to be high enough to obtain laser output in a single pass through the laser amplifier [Eq. (16)]. The lasers with the highest gain are the selenium laser (Se^{24+}) at 20.6 and 20.9 nm and the germanium laser (Ge^{22+}) at 23.2, 23.6, and 28.6 nm.

Liquid Laser Gain Media

Organic Dye Lasers⁸ A dye laser consists of a host or solvent material, such as alcohol or water, into which is mixed a laser species in the form of an organic dye molecule, typically in the proportion of one part in ten thousand. A large number of different dye molecules are used to make lasers covering a wavelength range of from 320 to 1500 nm with each dye having a laser bandwidth of the order of 30 to 50 nm. The wide, homogeneously broadened gain spectrum for each dye allows laser tunability over a wide spectrum in the ultraviolet, visible, and near-infrared. Combining the broad gain spectrum (Fig. 27) with a diffraction-grating or prism-tuning element allows tunable laser output to be obtained over the entire dye emission spectrum with a laser linewidth of 10 GHz or less. Dye lasers are available in either pulsed (up to 50 to 100 mJ/pulse) or continuous output (up to a few watts) in tabletop systems that are pumped by either flash lamps or by other lasers such as frequency doubled or tripled YAG lasers or argon ion lasers. Most dye lasers are arranged to have the dye mixture circulated by a pump into the gain region from a much larger reservoir since the dyes degrade at a slow rate during the excitation process.

Dye lasers, with their broad gain spectrum, are particularly attractive for producing ultrashort mode-locked pulses. Some of the shortest light pulses ever generated, of the order of 6×10^{-15} seconds,

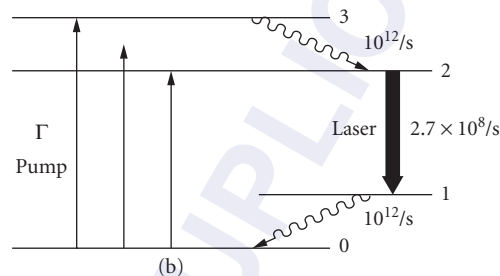
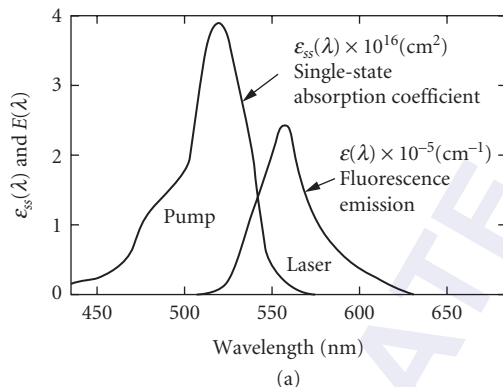


FIGURE 27 Absorption and emission spectra along with the energy-level diagram of an Rh6G organic dye laser showing both (a) the broad pump and emission bandwidths and (b) the fast decay of the lower laser level.

were produced with mode-locked dye lasers. A mode-locked dye laser cavity is shown in Fig. 14 with a thin dye gain region located within an astigmatically compensated laser cavity.

Dielectric Solid-State Laser Gain Media

Ruby Laser^{5,7} The ruby laser, with an output at 694.3 nm, was the first laser ever developed. It consisted of a sapphire (Al_2O_3) host material into which was implanted a chromium laser species in the form of Cr^{3+} ions at a concentration of 0.05 percent as the amplifying medium. The ruby laser involves a three-level optical pumping scheme, with the excitation provided by flash lamps, and operates either in a pulsed or cw mode. The three-level scheme for the ruby laser (Fig. 6) requires a large fraction of the population to be pumped out of the ground state before an inversion occurs. Therefore, the ruby laser is not as efficient as other solid-state lasers such as the Nd:YAG laser which employs a four-level pumping scheme (Fig. 8). It is therefore no longer used as much as it was in the early days.

Nd:YAG and Nd:Glass Lasers⁷ Neodymium atoms, in the form of Nd^{3+} ions, are doped into host materials, including crystals such as yttrium-aluminum-garnet (YAG) and various forms of glass, in concentrations of approximately one part per hundred. The pumping scheme is a four-level system as shown in Fig. 8. When implanted in YAG crystals to produce what is referred to as an Nd:YAG laser, the laser emits primarily at 1.06 μm with continuous powers of up to 250 W and with pulsed

powers as high as several mega watts. Difficulties in growing Nd:YAG crystals limit the size of commonly used laser rods to a maximum of 1 cm in diameter and 10 cm in length. Although this size limitation is somewhat restrictive, a YAG crystal has the advantage of high thermal conductivity allowing the rapid removal of wasted heat due to inefficient excitation. Efforts are being made to use Nd:YAG powders to be able to construct larger Nd:YAG laser rods, thereby providing the advantages of Nd:YAG laser without their usual size and laser power limitations. Slab geometries have recently been developed to compensate for focusing due to thermal gradients and gradient-induced stresses within the amplifier medium, thereby allowing higher average powers to be achieved. Efficient GaAs laser diodes are also being used to pump Nd:YAG amplifiers (see Fig. 15) since the diode pump wavelength matches a very strong pump absorption wavelength for the Nd^{3+} ion. This absorption wavelength is near the threshold pump energy, thereby minimizing the generation of wasted heat. Also, the diode pump laser can be made to more efficiently match the Nd laser mode volume than can a flash lamp.

Under long-pulse, flash lamp-pumped operation, Nd:YAG lasers can produce 5 J/pulse at 100 Hz from a single rod. The pulses are trains of relaxation oscillation spikes lasting 3 to 4 ms which is obtained by powerful, long-pulse flash lamp pumping. Q-switched Nd:YAG lasers, using a single Nd:YAG laser rod, can provide pulse energies of up to 1.5 J/pulse at repetition rates up to 200 Hz (75 W of average power).

Currently power levels of solid state Nd:YAG lasers are approaching 2 kW with repetition rates of nearly 10 kHz. These lasers are presently being developed for use in EUV microlithography as well as in applications such as liquid crystal display production, micromachining, and silicon processing (where carbon dioxide lasers are presently being used.)

Nd:glass lasers have several advantages over Nd:YAG lasers. They can be made in much larger sizes, allowing the construction of very large amplifiers. Nd:glass has a wider gain bandwidth which makes possible the production of shorter mode-locked pulses and a lower stimulated emission cross section. The latter property allows a larger population inversion to be created and thus more energy to be stored before energy is extracted from the laser. Nd:glass amplifiers operate at a wavelength near that of Nd:YAG, at 1.054 μm , for example, with phosphate glass, with a gain bandwidth 10 to 15 times broader than that of Nd:YAG. This larger gain bandwidth lowers the stimulated emission cross section compared to that of Nd:YAG, which allows increased energy storage when used as an amplifier, as mentioned above.

The largest laser system in the world is located at the National Ignition Facility (NIF) at Lawrence Livermore National Laboratories in Livermore, California. It consists of 192 separate laser beams traveling about 1000 ft from their origination (one of two master laser oscillators), to the center of a target chamber. The amplifiers are Nd-doped phosphate glass disc-shaped amplifiers with diameters of over a meter, installed at Brewster's angle within the laser beam path to reduce reflection losses. The beam of each of the 192 amplifiers is capable of producing an energy of 20,000 J with pulse durations ranging from 100 fs to 25 ps. One of the objectives of the laser facility is to test the concept of laser fusion as a source of useful commercial power production.

Neodymium-YLF Lasers Nd:YLF has relatively recently become an attractive laser with the successful implementation of diode laser pumping of solid state laser materials. It has the advantage over Nd:YAG in that the upper-laser-level lifetime is twice as long, thereby allowing twice the energy storage. Its lower stimulated emission cross section than that of Nd:YAG also increases the energy storage and thus the power output per unit volume of the crystal. It has a relatively large thermal conductivity, similar to Nd:YAG. Also the output is polarized and the crystal exhibits low thermal birefringence. Because its emission wavelength matches that of phosphate glass, the laser makes an ideal laser oscillator for seeding large Nd:glass amplifiers for laser fusion studies.

Neodymium:Yttrium Vanadate Lasers The Nd:VO_4 laser has also come into prominence with the use of diode pumping. Its advantages over Nd:YAG include a 5 times larger stimulated emission cross section (and hence higher gain) as well as a four times stronger absorption cross section with a 6 times wider pump absorption cross section centered at 809 nm. Therefore it can be effectively pumped in crystals of only a few millimeters in length and is therefore attractive for use in producing

small diode-pumped lasers. Typically this laser is frequency doubled or tripled intracavity to produce several watts of power at either 532 or 355 nm. It can be operated either cw or Q-switched at repetition rates of up to 100 KHz. It can also be operated at 1.342 μm .

Fiber Lasers The first major application of fiber lasers was in their use as amplifiers in the field of optical communications. The erbium-doped fiber is installed directly in the fiber-optic transmission line and pumped through the fiber itself. The useful wavelengths are in the 1.53 μm region where optical fibers have their lowest loss. A more recent application is in the production of high power output from large mode-area fibers. These lasers utilize single-emitter semiconductor diodes as the light source to pump the cladding of rare-earth doped optical fibers. Pulsed Ytterbium fiber lasers have demonstrated high average powers and peak powers of up to 25 kW with variable repetition rates from a few kilohertz to up to 400 kHz. Fiber lasers are being considered as the replacement technology for conventional solid state lasers due to their compactness, high wall-plug efficiency, high average and peak powers, stability, close to diffraction-limited beam quality and lack of thermal effects.

Ti:Al₂O₃ Laser and Other Broad Bandwidth Solid-State Lasers⁷

Another class of solid-state lasers provides emission and gain over a bandwidth of the order of 100 to 400 nm in the near-infrared, as indicated in Fig. 28. The pump absorption band is in the visible spectrum, thus allowing such pump sources as flash lamps and other lasers. The Ti:Al₂O₃ laser is perhaps the most well-known laser of this category in that it has the widest bandwidth, covering a wavelength range from 0.67 μm to greater than 1.07 μm . Other lasers of this type include alexandrite (Cr:BeAl₂O₄), lasing from 0.7 to 0.8 μm , and Cr:LiSAF which lases from 0.8 to 1.05 μm . These lasers are used in applications where either wide tunability or short-pulse production are desired. Pulses as short as 4 fs have been produced with mode-locked versions of these lasers. Ti:Al₂O₃ lasers, although offering very wide gain bandwidth, have relatively short upper-laser-level lifetimes (3 μs), thereby making them less efficient when pumping with conventional flash lamps. Cr:LiSAF lasers

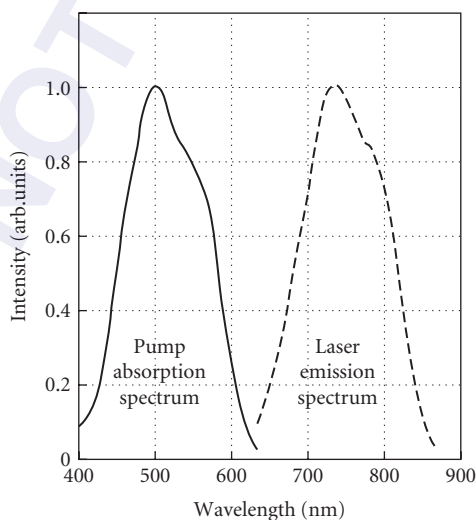


FIGURE 28 Absorption and emission spectra of a Ti:Al₂O₃ laser crystal showing the extremely broad emission (gain) spectrum for this material.

have longer upper-level lifetimes (67 μs), closer to that of Nd:YAG or Nd:glass and have demonstrated efficient laser operation with pumping technologies developed for Nd:YAG lasers.

Color-Center Laser²² Color-center laser gain media are produced by a different form of impurity species than most solid-state lasers. Special defect centers (F-centers) are produced within alkali-halide crystals at a density of 1 part in 10,000 by irradiation with x rays. These defect centers have absorbing regions in the visible portion of the spectrum and emission (and gain) in the near-infrared. A variety of crystals are used to span the laser wavelength spectrum from 0.8 to 4.0 μm . Disadvantages of color-center lasers include operation at temperatures well below room temperature and the necessity to re-form the color centers at intervals of weeks or months in most cases.

Semiconductor Laser Gain Media

Semiconductor Lasers²³ Semiconductor or diode lasers, typically about the size of a grain of salt, are the smallest lasers yet devised. They consist of a p - n junction formed in semiconductor crystal such as GaAs or InP in which the p -type material has an excess of holes (vacancies due to missing electrons) and the n -type material has an excess of electrons. When these two types of materials are brought together to form a junction, and an electric field in the form of a voltage is applied across the junction in the appropriate direction, the electrons and holes are brought together and recombine to produce recombination radiation at or near the wavelength associated with the bandgap energy of the material. The population of electrons and holes within the junction provides the upper-laser-level population, and the recombination radiation spectrum is the gain bandwidth $\Delta\nu$ of the laser, typically of the order of 0.5 to 1.0 nm.

The extended gain length required for these lasers is generally provided by partially reflecting cleaved parallel faces at the ends of the crystals which serve as an optical cavity. Because the cavity is so short, the longitudinal modes are spaced far apart in frequency ($\Delta\nu \cong 1-5 \times 10^{11}$ Hz or several tenths of a nanometer), and thus it is possible to obtain single longitudinal mode operation in such lasers. They require a few volts to operate with milliamperes of current.

Heterostructure semiconductor lasers include additional layers of different materials of similar electronic configurations, such as aluminum, indium, and phosphorous, grown adjacent to the junction to help confine the electron current to the junction region in order to minimize current and heat dissipation requirements (see Fig. 29). The laser mode in the transverse direction is either

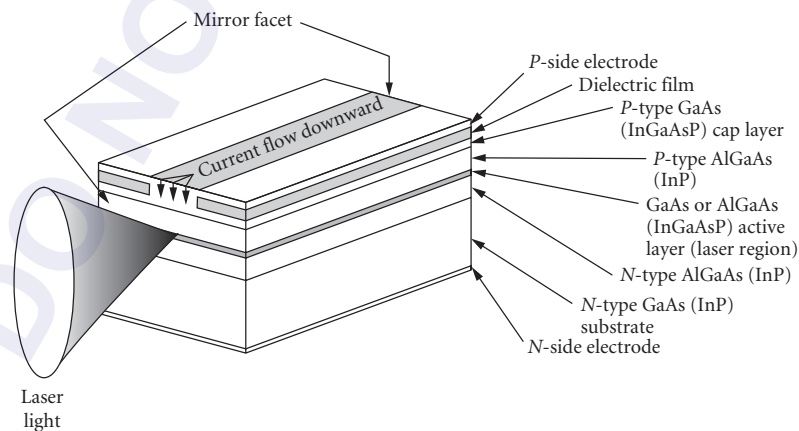


FIGURE 29 A typical heterostructure semiconductor laser showing the various layers of differential materials and the narrow region where current flows to produce the gain region in the active layer.

controlled by gain guiding, in which the gain is produced over a specific narrow lateral extent determined by fabrication techniques, or by index guiding, in which the index of refraction in the transverse direction is varied to provide total internal reflection of a guided mode. Quantum-well lasers have a smaller gain region (cross section), which confines the excitation current and thus the laser mode to an even smaller lateral region, thereby significantly reducing the threshold current and also the heat dissipation requirements. Because of these low threshold requirements, quantum-well semiconductor lasers are used almost exclusively for most semiconductor laser applications.

Multiple semiconductor lasers fabricated within the same bulk structure, known as *semiconductor arrays*, can be operated simultaneously to produce combined cw power outputs of up to 10 W from a laser crystal of dimensions of the order of 1 mm or less. Semiconductor lasers have also been fabricated in multiple arrays mounted vertically on a chip with the mirrors in the plane of the chip. They are known as *vertical cavity semiconductor lasers* (see below). Each of the individual lasers has dimensions of the order of 5 μm and can be separately accessed and excited.

Semiconductor lasers operate over wavelengths ranging from 400 nm to 2.2 μm by using special doping materials to provide the expanded or contracted bandgap energies that provide the varied wavelengths. The newest additions to this class of lasers are based upon the GaN laser materials with the active region consisting of various combinations of InGaN dopings that provide laser wavelengths in the green, blue, and violet portions of the spectrum.

Quantum Cascade Lasers Laser action in semiconductor lasers typically occurs when recombination radiation occurs across the band gap of the semiconductor. Quantum cascade lasers are different in that the radiation occurs from transitions between quantized conduction band states. Such transitions are inherently low-energy transitions and hence the laser output occurs in the middle infrared at wavelengths ranging from 3 to 24 μm . In the laser process, electrons are fed into the injector region of each stage via an electric field and transition across a mini-band from $n = 3$ to $n = 2$ quantum well levels thereby emitting the near infrared radiation. Several of these stages are stacked together in series to produce high power output. The active region is typically made up of aluminum indium arsenide and gallium indium arsenide. The lasers are particularly useful for operation in the two atmospheric windows at 3 to 5 μm and 8 to 13 μm .

Vertical Cavity Surface-Emitting Lasers (VCSELs) The vertical cavity laser is a different type of semiconductor laser than the typical edge-emitting lasers in that the emission occurs from the top surface of the laser and the cavity mirrors are comprised of multilayer dielectric coatings on the top and bottom surfaces of the very thin gain medium. These lasers can be made, using lithographic techniques, in large arrays on a microchip and can also be tested on the chip before being cleaved into individual lasers. Applications include optical fiber data transmission, absorption spectroscopy, and laser printers.

Laser Gain Media in Vacuum

Free Electron Laser²⁴ Free-electron lasers are significantly different than any other type of laser in that the laser output does not result from transitions between discrete energy levels in specific materials. Instead, a high-energy beam of electrons, such as that produced by a synchrotron, traveling in a vacuum with kinetic energies of the order of 1 MeV, are directed to pass through a spatially varying magnetic field produced by two regular arrays of alternating magnet poles located on opposite sides of the beam as shown in Fig. 30. The alternating magnetic field causes the electrons to oscillate back and forth in a direction transverse to the beam direction. The frequency of oscillation is determined by the electron beam energy, the longitudinal spacing of the alternating poles (the magnet period), and the separation between the magnet arrays on opposite sides of the beam. The transverse oscillation of the electrons causes them to radiate at the oscillation frequency and to thereby stimulate other electrons to oscillate and thereby radiate at the same frequency, in phase with the originally oscillating electrons. The result is an intense tunable beam of light emerging from the end

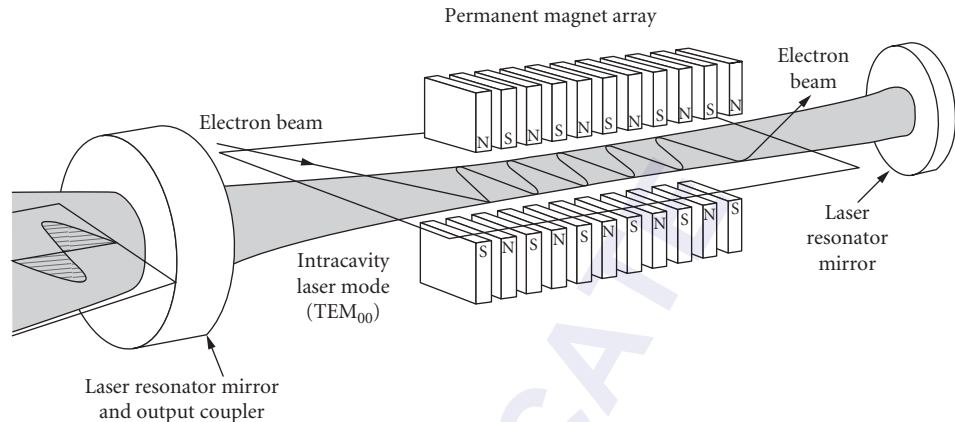


FIGURE 30 A general diagram of a free-electron laser showing how the electron beam is introduced into the cavity and how the alternating magnets cause the beam to oscillate to produce laser radiation.

of the device. Mirrors can be placed at the ends of the radiated beam to produce feedback and extra amplification.

Free-electron lasers have operated at wavelengths ranging from the near-ultraviolet ($0.25\ \mu\text{m}$) to the far-infrared (6 mm) spectral regions. They are efficient devices and also offer the potential of very high average output powers.

16.7 REFERENCES

1. A. Corney, *Atomic and Laser Spectroscopy*, Clarendon Press, Oxford, 1977.
2. P. W. Milonni and J. H. Everly, *Lasers*, John Wiley & Sons, New York, 1988.
3. H. G. Kuhn, *Atomic Spectra*, 2d ed., Academic Press, New York, 1969.
4. J. T. Verdeyen, *Laser Electronics*, 2d ed., Prentice-Hall, Englewood Cliffs, New Jersey, 1989.
5. T. H. Maiman, *Nature* **187**:493–494 (1960).
6. W. T. Walter, N. Solimene, M. Piltch, and G. Gould, *IEEE J. of Quant. Elect.* **QE-2**:474–479 (1966).
7. W. Koehner, *Solid-State Laser Engineering*, 3d ed., Springer-Verlag, New York, 1992.
8. F. J. Duarte and L. W. Hillman (eds.), *Dye Laser Principles*, Academic Press, New York, 1990.
9. W. W. Rigrod, *J. Appl. Phys.* **34**:2602–2609 (1963); **36**:2487–2490 (1965).
10. C. S. Willett, *An Introduction to Gas Lasers: Population Inversion Mechanisms*, Pergamon Press, Oxford, 1974.
11. A. G. Fox and T. Li, *Bell Syst. Tech. J.* **40**:453–488 (1961).
12. G. D. Boyd and J. P. Gordon, *Bell Syst. Tech. J.* **40**:489–508 (1961).
13. H. Kogelnik and T. Li, *Proc. IEEE* **54**:1312–1329 (1966).
14. A. E. Siegman, *Lasers*, University Science Books, Mill Valley, Calif., 1986.
15. W. Wagner and B. Lengyel, *J. Appl. Phys.* **34**:2040–2046 (1963).
16. R. L. Fork, B. I. Greene, and C. V. Shank, *Appl. Phys. Lett.* **38**:671–672 (1981).
17. J. Mark, L. Y. Liu, K. L. Hall, H. A. Haus, and E. P. Ippen, *Opt. Lett.* **14**:48–50 (1989).
18. D. K. Negus, L. Spinelli, N. Goldblatt, and G. Feuget, *Proc. on Advanced Solid State Lasers* (Optical Soc. of Am.) **10**:120–124 (1991).
19. H. Kogelnik and C. V. Shank, *J. Appl. Phys.* **43**:2327–2335 (1972).

20. M. Rokni, J. A. Mangano, J. H. Jacobs, and J. C Hsia, *IEEE Journ. of Quant. Elect.* **QE-14**:464–481 (1978).
21. R. C. Elton, *X-Ray Lasers*, Academic Press, New York, 1990.
22. L. F. Mollenauer and D. H. Olson, *J. Appl. Phys.* **46**:3109–3118 (1975).
23. A. Yariv, *Quantum Electronics*, John Wiley & Sons, New York, 1989.
24. L. R. Elias, W. M. Fairbank, J. M. Madey, H. A. Schewttman, and T. J. Smith, *Phys. Rev. Lett.* **36**:717–720 (1976).

DO NOT DUPLICATE

LIGHT-EMITTING DIODES

Roland H. Haitz, M. George Craford, and Robert H. Weissman

*Hewlett-Packard Co.
San Jose, California*

17.1 GLOSSARY

c	velocity of light
E_g	semiconductor energy bandgap
h	Planck's constant
I_T	total LED current
J	LED current density
k	Boltzmann's constant
M	magnification
n_0	low index of refraction medium
n_1	high index of refraction medium
q	electron charge
T	temperature
V	applied voltage
η_i	internal quantum efficiency
θ_c	critical angle
λ	emission wavelength
τ	total minority carrier lifetime
τ_n	nonradiative minority carrier lifetime
τ_r	radiative minority carrier lifetime

17.2 INTRODUCTION

Over the past 25 years the light-emitting diode (LED) has grown from a laboratory curiosity to a broadly used light source for signaling applications. In 1992 LED production reached a level of approximately 25 billion chips, and \$2.5 billion worth of LED-based components were shipped to original equipment manufacturers.

This chapter covers light-emitting diodes from the basic light-generation processes to descriptions of LED products. First, we will deal with light-generation mechanisms and light extraction. Four major types of device structures—from simple grown or diffused homojunctions to complex double heterojunction devices are discussed next, followed by a description of the commercially important semiconductors used for LEDs, from the pioneering GaAsP system to the AlGaInP system that is currently revolutionizing LED technology. Then processes used to fabricate LED chips are explained—the growth of GaAs and GaP substrates, the major techniques used for growing the epitaxial material in which the light-generation processes occur, and the steps required to create LED chips up to the point of assembly. Next the important topics of quality and reliability—in particular, chip degradation and package-related failure mechanisms—will be addressed. Finally, LED-based products, such as indicator lamps, numeric and alphanumeric displays, optocouplers, fiber-optic transmitters, and sensors, are described.

This chapter covers the mainstream structures, materials, processes, and applications in use today. It does not cover certain advanced structures, such as quantum well or strained layer devices, a discussion of which can be found in Chap. 19, “Semiconductor Lasers.” The reader is also referred to Chap. 19 for current information on edge-emitting LEDs, whose fabrication and use are similar to lasers.

For further information on the physics of light generation, the reader should consult Refs. 1 to 11. Semiconductor material systems for LEDs are discussed in Refs. 13 to 24. Crystal growth, epitaxial, and wafer fabrication processes are discussed in detail in Refs. 25 to 29.

17.3 LIGHT-GENERATION PROCESSES

When a p - n junction is biased in the forward direction, the resulting current flow across the boundary layer between the p and n regions has two components: holes are injected from the p region into the n region and electrons are injected from the n region into the p region. This so-called minority-carrier injection disturbs the carrier distribution from its equilibrium condition. The injected minority carriers recombine with majority carriers until thermal equilibrium is reestablished. As long as the current continues to flow, minority-carrier injection continues. On both sides of the junction, a new steady-state carrier distribution is established such that the recombination rate equals the injection rate.^{1,2}

Minority-carrier recombination is not instantaneous. The injected minority carriers have to find proper conditions before the recombination process can take place. Both energy and momentum conservation have to be met. Energy conservation can be readily met since a photon can take up the energy of the electron-hole pair, but the photon doesn't contribute much to the conservation of momentum. Therefore, an electron can only combine with a hole of practically identical and opposite momentum. Such proper conditions are not readily met, resulting in a delay. In other words, the injected minority carrier has a finite lifetime τ_r before it combines radiatively through the emission of a photon.² This average time to recombine radiatively through the emission of light can be visualized as the average time it takes an injected minority carrier to find a majority carrier with the right momentum to allow radiative recombination without violating momentum conservation.

Unfortunately, radiative recombination is not the only recombination path. There are also crystalline defects, such as impurities, dislocations, surfaces, etc., that can trap the injected minority carriers. This type of recombination process may or may not generate light. Energy and momentum conservation are met through the successive emission of phonons. Again, the recombination process is not instantaneous because the minority carrier first has to diffuse to a recombination site. This nonradiative recombination process is characterized by a lifetime τ_n .²

Of primary interest in design of light-emitting diodes is the maximization of the radiative recombination relative to the nonradiative recombination. In other words, it is of interest to develop conditions where radiative recombination occurs fairly rapidly compared with nonradiative recombination. The effectiveness of the light-generation process is described by the fraction of the injected minority carriers that recombine radiatively compared to the total injection. The internal quantum efficiency η_i can be calculated from τ_r and τ_n . The combined recombination processes lead to a total minority-carrier lifetime τ given by Eq. (1):

$$\frac{1}{\tau} = \frac{1}{\tau_r} + \frac{1}{\tau_n} \quad (1)$$

η_i is simply computed from Eq. (1) as the fraction of carriers recombining radiatively:²

$$\eta_i = \frac{\tau_n}{\tau_r + \tau_n} \quad (2)$$

Of interest are two simple cases: in the case of excellent material quality (large τ_n) or efficient radiative recombination conditions (small τ_r), the internal quantum efficiency approaches 100 percent. For the opposite case ($\tau_n \ll \tau_r$), we find $\eta_i \approx \tau_n / \tau_r \ll 1$. As discussed under “Material Systems,” there are several families of III–V compounds with internal quantum efficiencies approaching 100 percent. There are also other useful semiconductor materials with internal quantum efficiencies in the 1- to 10-percent range.

To find material systems for LEDs with a high quantum efficiency, one has to understand the band structure of semiconductors. The band structure describes the allowed distribution of energy and momentum states for electrons and holes (see Fig. 1 and Ref. 2). In practically all semiconductors the lower band, also known as the *valence band*, has a fairly simple structure, a

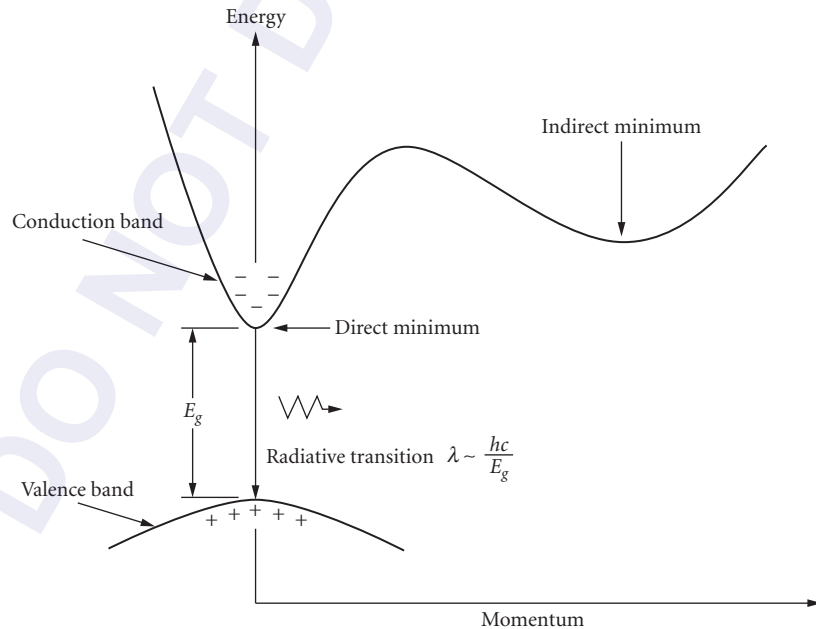


FIGURE 1 Energy band structure of a direct semiconductor showing radiative recombination of electrons in the conduction band with holes in the valence band.

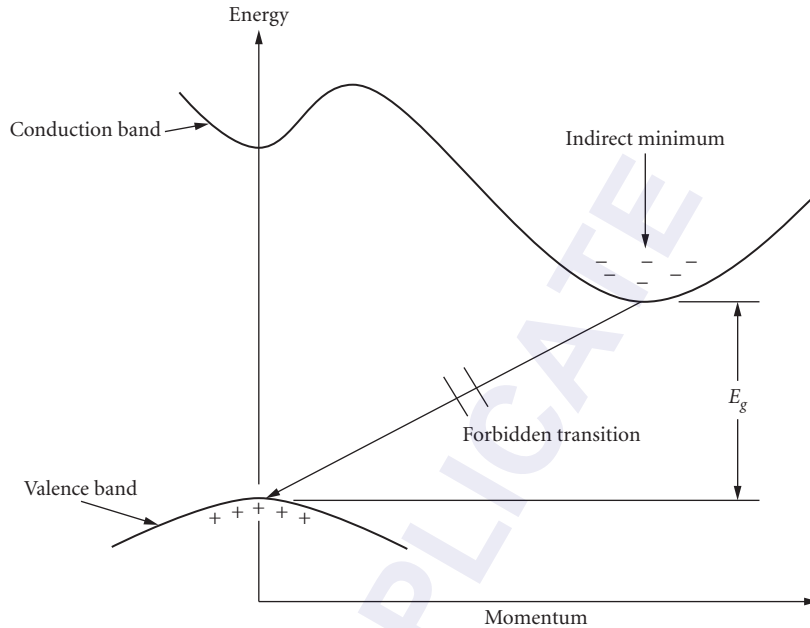


FIGURE 2 Energy band structure of an indirect semiconductor showing the conduction-band minima and valence-band maximum at different positions in momentum space. Radiative recombination of conduction-band electrons and valence-band holes is generally forbidden.

paraboloid around the $\langle 0, 0, 0 \rangle$ crystalline direction. Holes will take up a position near the apex of the paraboloid and have very small momentum. The upper band, also known as the *conduction band*, is different for various semiconductor materials. All semiconductors have multiple valleys in the conduction band. Of practical interest are the valleys with the lowest energy. Semiconductor materials are classified as either *direct* or *indirect*.¹⁻³ In a direct semiconductor, the lowest valley in the conduction band is directly above the apex of the valence-band paraboloid. In an indirect semiconductor, the lowest valleys are not at $\langle 0, 0, 0 \rangle$, but at different positions in momentum per energy space (see Fig. 2). Majority or minority carriers mostly occupy the lowest energy states, i.e., holes near the top of the valence band paraboloid and electrons near the bottom of the lowest conduction-band valley.

In the case of a direct semiconductor the electrons are positioned directly above the holes at the same momentum coordinates. It is relatively easy to match up electrons and holes with proper momentum-conserving conditions. Thus, the resulting radiative lifetime τ_r is short. On the other hand, electrons in an indirect valley will find it practically impossible to find momentum-matching holes and the resulting radiative lifetime will be long. Injected carriers in indirect material generally recombine nonradiatively through defects.

In a direct semiconductor, such as GaAs, the radiative lifetime τ_r is in the range of 1 to 100 ns, depending on doping, temperature, and other factors. It is relatively easy to grow crystals with sufficiently low defect density such that τ_n is in the same range as τ_r .

For indirect semiconductors, such as germanium or silicon, the radiative recombination process is extremely unlikely, and τ_r is in the range of seconds.¹ In this case, $\tau_r \gg \tau_n$, and practically all injected carriers recombine nonradiatively.

The wavelength of the photons emitted in a radiative recombination event is determined by the energy difference between the recombining electron-hole pair. Since carriers relax quickly to an energy level near the top of the valence band (holes), or the bottom of the conduction band

(electrons), we have the following approximation for the wavelength λ of the emitted photon (see Fig. 1):

$$\lambda \approx hc/E_g \quad (3)$$

where h = Planck's constant, c = velocity of light, and E_g = bandgap energy.

This relation is only an approximation since holes and electrons are thermally distributed at levels slightly below the valence-band maximum and above the conduction-band minimum, resulting in a finite linewidth in the energy or wavelength of the emitted light. Another modification results from a recombination between a free electron and a hole trapped in a deep acceptor state. See Ref. 2 for a discussion of various recombination processes.

To change the wavelength or energy of the emitted light, one has to change the bandgap of the semiconductor material. For example, GaAs with a bandgap of 1.4 eV has an infrared emission wavelength of 900 nm. To achieve emission in the visible red region, the bandgap has to be raised to around 1.9 eV. This increase in E_g can be achieved by mixing GaAs with another material with a wider bandgap, for instance GaP with $E_g = 2.3$ eV. By adjusting the ratio of arsenic to phosphorous the bandgap of the resulting ternary compound, GaAsP, can be tailored to any value between 1.4 and 2.3 eV.³

The resulting band structure with varying As to P ratio is illustrated in Fig. 3. Note, the two conduction-band valleys do not move upward in energy space at the same rate. The direct valley moves up faster than the indirect valley with increasing phosphorous composition. At a composition of around 40 percent GaP and 60 percent GaAs, the direct and indirect valleys are about equal in energy. When the valleys are approximately equal in energy, electrons in the conduction band can scatter from the direct valley into the indirect valley. While the direct valley electrons still undergo rapid radiative recombination, the indirect valley electrons have a long radiative lifetime and either have to be scattered back to the direct valley or they will recombine nonradiatively. In other words, near this crossover

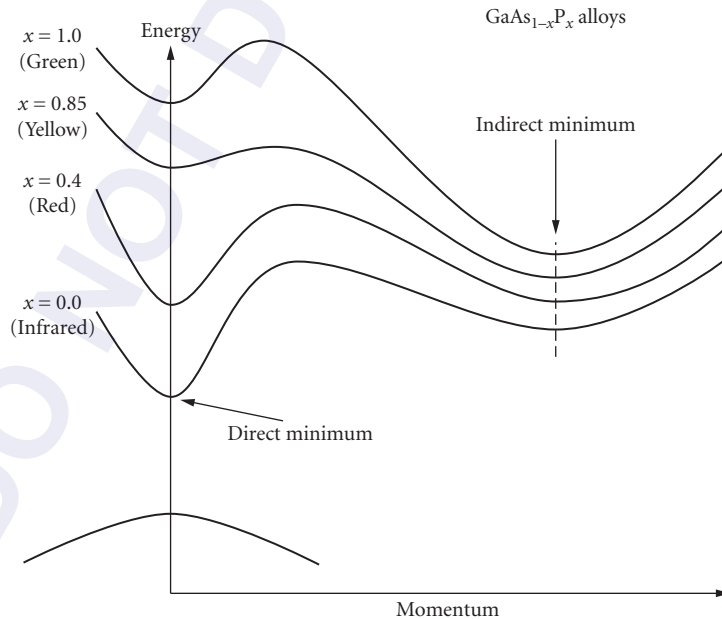


FIGURE 3 Energy band diagram for various alloys of the GaAs_{1-x}P_x material system showing the direct and indirect conduction-band minima for various alloy compositions.

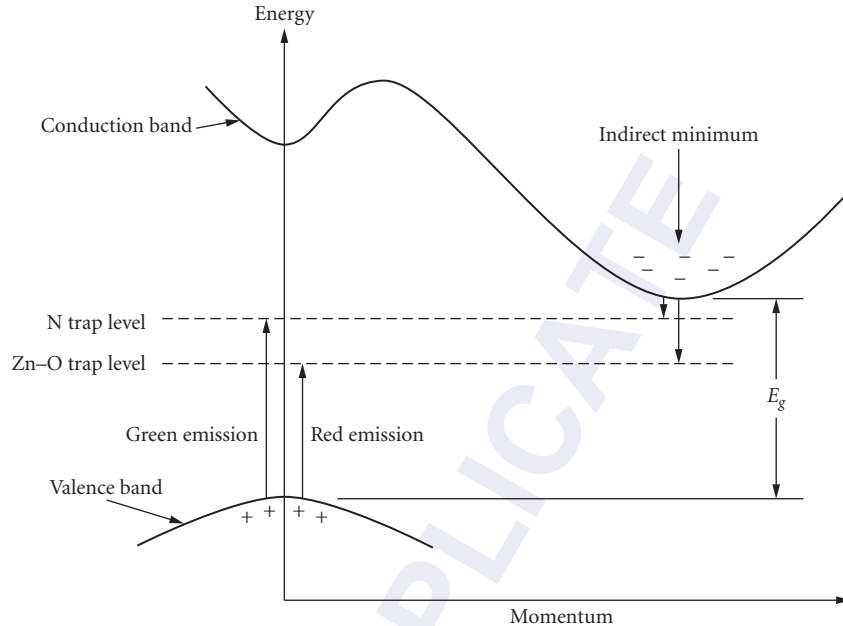


FIGURE 4 Formation of excitons (electron-hole pairs) by the addition of isoelectronic dopants N and ZnO to an indirect semiconductor. The excitons have a high probability to recombine radiatively.

between direct and indirect valleys, the radiative efficiency drops off dramatically, and for compositions with greater than 40 percent phosphorous the direct radiative recombination is practically nonexistent.^{3,4}

The above discussion indicates that indirect semiconductors are not suitable for efficient generation of light through minority-carrier recombination. Fortunately, the introduction of so-called isoelectronic impurities can circumvent this limitation and introduces a new radiative recombination process.⁵ A commonly used isoelectronic trap is generated by substituting a nitrogen atom for phosphorous in the GaAsP system^{6–11} (see Fig. 4). Since N and P both have five electrons in their outer shell, the trap is electrically neutral. However, the stronger electronegativity of N relative to P can result in capture of an electron from the conduction band. Since the electron is very tightly bound to the impurity atom, its wave function in momentum space is spread out and has reasonable magnitude at $\langle k=0 \rangle$ in momentum space.¹⁰ The negatively charged defect can attract a free hole to form a loosely bound electron-hole pair or “exciton.” This electron-hole pair has a high probability to recombine radiatively. The energy of the emitted light is less than E_g . Another isoelectronic trap in GaP is formed by ZnO pairs (Zn on a Ga site and O on a P site) (see Fig. 4). The ZnO trap is deeper than the N trap, resulting in longer wavelength emission in the red region of the spectrum.¹

The recombination process for exciton recombination is quite complex. For a detailed analysis, the reader is referred to Ref. 10. One result of this analysis is the recognition that the bound exciton has a relatively long lifetime in the range of 100 to 1000 ns. Light emission by exciton recombination is generally slower than emission due to direct band-to-band recombination.

17.4 LIGHT EXTRACTION

Generating light efficiently within a semiconductor material is only one part of the problem to build an efficient light source. The next challenge is the extraction of light from within the LED chip to the outside. The designer must consider total internal reflection.¹ According to Snell’s law, light can

escape from a medium of high index of refraction n_1 into a medium of low index refraction n_0 only if it intersects the surface between the two media at an angle from normal less than the critical angle θ_c with θ_c being defined by Eq. (4):

$$\theta_c = \arcsin n_0/n_1 \quad (4)$$

Most semiconductor LEDs have an isotropic emission pattern as seen from within the light-generating material. Assuming a cubic shape for the LED chip, because of internal reflections, only a small fraction of the isotropically emitted light can escape any of the six surfaces. As a case in point, let us calculate the emission through the top surface. For typical light-emitting semiconductors, n_1 is in the range of 2.9 to 3.6. If $n_1 = 3.3$ and $n_0 = 1.0$ (air), we find $\theta_c = 17.6^\circ$. The emission from an isotropic source into a cone with a half angle of θ_c is given by $(1 - \cos \theta_c)/2$. After correcting for Fresnel reflections, only 1.6 percent of the light generated escapes through the LED top surface into air. Depending on chip and p - n junction geometry, virtually all of the remaining light (98.4 percent) is reflected and absorbed within the LED chip.

The fraction of light coupled from chip to air is a function of the number of surfaces through which the chip can transmit light effectively. Most LED chips are called “absorbing substrate” (AS) chips. In such a chip, the starting substrate material (discussed later under “Substrate Technology”) has a narrow bandgap and absorbs all the light with energy greater than the bandgap of the substrate. Consider the case of a GaAsP LED grown on a GaAs substrate. The emitted light ($E_g > 1.9\text{eV}$) is absorbed by the GaAs substrate ($E_g = 1.4\text{ eV}$). Thus, a GaAsP-emitting layer on a GaAs substrate can transmit only through its top surface. Light transmitted toward the side surfaces or downward is absorbed.

To increase light extraction, the substrate or part of the epitaxial layers near the top of the chip has to be made of a material transparent to the emitted light. The “transparent substrate” (TS) chip is designed such that light transmitted toward the side surfaces within θ_c half-angle cones can escape. Assuming that there is negligible absorption between the point of light generation and the side walls, this increases the extraction efficiency by a factor of 5 (5 instead of 1 escape cones).

In a TS chip, additional light can be extracted if the side walls are nonplanar, i.e., if light from outside an escape cone can be scattered into an escape cone. This process increases the optical path within the chip and is very dependent on residual absorption. In a chip with low absorption and randomizing side surfaces, most of the light should escape. Unfortunately, in practical LED structures, there are several absorption mechanisms left such as front and back contacts, crystal defects, and absorption in areas where secondary radiative recombination is inefficient.

A common approach is to use a hybrid chip with properties between AS and TS chips. These chips utilize a thick, transparent window layer above the light-emitting layer. If this layer is sufficiently thick, then most of the light in the top half of the cones transmitted toward the side surfaces will reach the side of the chip before hitting the substrate. In this case of hybrid chips, the efficiency is between that of AS and TS chips as shown in Table 1.

Another important way to increase extraction efficiency is derived from a stepwise reduction in the index of refraction from chip to air. If the chip is first imbedded in a material with an intermediate index, i.e., plastic with $n_2 = 1.5$, then the critical angle θ_c between chip and plastic is increased to 27° . The extraction efficiency relative to air increases by the ratio of $(n_2/n_0)^2$ plus some additional correction for Fresnel-reflection losses. The gain from plastic encapsulation is usually around

TABLE 1 Extraction Efficiency into Air or Plastic for Three Types of Commonly Used LED Chips

Chip Type	No. Cones	Typical extraction efficiency	
		Air (%)	Plastic (%)
AS	1	1.5	4
Thick window	3	4.5	12
TS	5	7.5	20

2.7 times compared to air. Chips with multipath internal reflection will result in lower gains. It is important to note that this gain can be achieved only if the plastic/air interface can accommodate the increased angular distribution through proper lenslike surface shaping or efficient scattering optics. Table 1 illustrates the approximate extraction efficiencies achieved by the three dominant chip structures in air and in plastic. The numbers assume only first-pass extraction, limited absorption, and no multiple reflections within the chip.

17.5 DEVICE STRUCTURES

LED devices come in a broad range of structures. Each material system (see following section) requires a different optimization. The only common feature for all LED structures is the placement of the p - n junction where the light is generated. The p - n junction is practically never placed in the bulk-grown substrate material for the following reasons:

- The bulk-grown materials such as GaAs, GaP, and InP usually do not have the right energy gap for the desired wavelength of the emitted light.
- The light-generating region requires moderately low doping that is inconsistent with the need for a low series resistance.
- Bulk-grown material often has a relatively high defect density, making it difficult to achieve high efficiency.

Because of these reasons, practically all commercially important LED structures utilize a secondary growth step on top of a single-crystal bulk-grown substrate material. The secondary growth step consists of a single-crystal layer lattice matched to the substrate. This growth process is known as *epitaxial growth* and is described in a later section of this chapter.

The commonly used epitaxial structures can be classified into the following categories:

- Homojunctions
 - grown
 - diffused
- Heterojunctions
 - single confinement
 - double confinement

Grown Homojunctions

Figure 5 illustrates one of the simplest design approaches to an LED chip. An n -type GaAs layer with low to moderate doping density is grown on top of a highly doped n -type substrate by a vapor or liquid-phase epitaxial process (see “Epitaxial Technology”). After a growth of 5 to 10 μm , the doping is changed to p type for another 5 to 10 μm . A critical dimension is the thickness of the epitaxial p layer. The thickness should be larger than the diffusion length of electrons. In other words, the electrons should recombine radiatively in the epitaxially grown p layer before reaching the surface. The p layer should be of sufficiently high quality to meet the condition for efficient recombination, i.e., $\tau_n \gg \tau_r$. In addition, the side surfaces may have to be etched to remove damage. Damage and other defects where the p - n junction intercepts the chip surface can lead to a substantial leakage current that reduces efficiency, especially at low drive levels.

The structure of Fig. 5 was used in some of the earliest infrared emitters (wavelength 900 nm). Efficiency was low, typically 1 percent. Modern infrared emitters use Si-doped GaAs (see Fig. 6). The detailed recombination mechanism in GaAs:Si even today is quite controversial and goes beyond the scope of this publication. The recombination process has two important characteristics: (1) the radiative lifetime τ_r is relatively slow, i.e., in the range of 1 μs and (2) the wavelength is shifted to 940 nm.

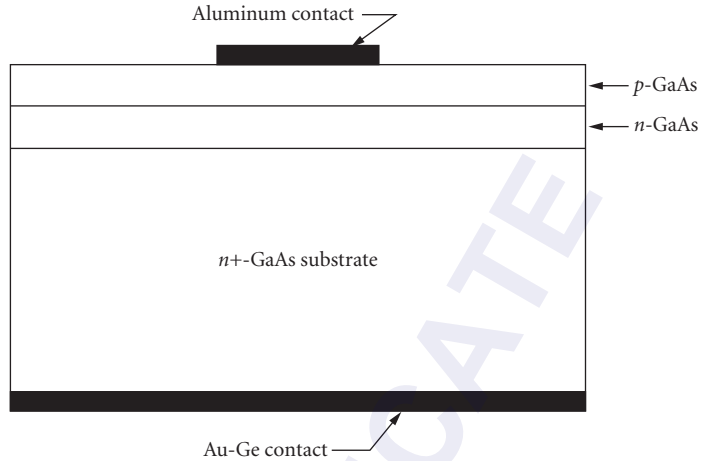


FIGURE 5 Cross section of an infrared LED chip. A p - n junction is formed by epitaxially growing n - and p -doped GaAs onto an n + -doped GaAs substrate.

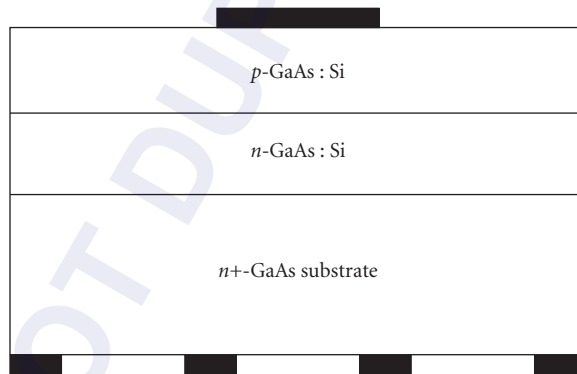


FIGURE 6 A high-efficiency IR LED made by LPE growth of GaAs that is doped with silicon on both the p and n sides of the junction. To increase external quantum efficiency, a partly reflective back contact is employed.

At this wavelength the GaAs substrate is partly transparent, making this device a quasi-TS structure with efficiencies into plastic of 5 to 10 percent.

Diffused Homojunction

The chip structure of Fig. 5 can also be produced by a zinc (Zn) diffusion into a thick n layer. The commercially most significant structure of this type is shown in Fig. 7. By replacing 40 percent of the As atoms with P atoms, the bandgap is increased to 1.92 eV to make a GaAsP LED that emits visible red light. In this case, the p - n junction is diffused selectively by using a deposited layer of silicon-nitride as a diffusion mask. This structure of Fig. 7 has several advantages over the structure of Fig. 5. Lateral diffusion of Zn moves the intersection of the junction with the chip surface

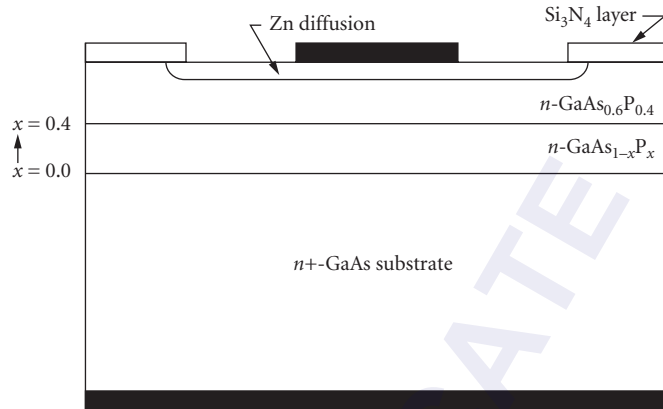


FIGURE 7 GaAsP LED which emits at 650 nm. On a GaAs substrate, a layer is grown whose composition varies linearly from GaAs to GaAs $_{0.6}$ P $_{0.4}$, followed by a layer of constant composition. Zinc is selectively diffused using an Si_3N_4 mask to form the light-emitting junction.

underneath the protecting Si_3N_4 layer. This layer protects the junction from contamination and adds to the long-term stability of the device. In addition, it is important for applications requiring more than one clearly separated light-emitting area. For instance, seven-segment displays, such as those used in the early handheld calculators, are made by diffusing seven long and narrow stripes into a single chip of GaAsP material. (See Fig. 8.) This chip consists of eight (seven segments plus decimal point) individually addressable p -regions (anodes) with a common n -type cathode. Such a chip is feasible only in an AS-type structure because the individual segments have to be optically isolated from each other. A TS-type structure results in unacceptable levels of crosstalk.

Figure 7 shows another feature of practical LED devices. The composition with 40 percent P and 60 percent As has a lattice constant (atomic spacing) that is different from the GaAs substrate. Such a lattice mismatch between adjacent layers would result in a very high density of dislocations. To reduce this problem to an acceptable level, one has to slowly increase the phosphorous composition from 0 percent at the GaAs interface to 40 percent over a 10- to 20- μm -thick buffer layer. Typically, the buffer layer is graded linearly. The phosphorous composition is increased linearly from bottom to top. The thicker the buffer layer, the lower is the resulting dislocation density. Cost constraints keep the layer in the 10- to 15- μm range. The layer of constant composition (40 percent P) has to be thick enough to accommodate the Zn diffusion plus the diffusion length of minority carriers. A thickness range of 5 to 10 μm is typical.

Another variation of a homojunction is shown in the TS-chip structure of Fig. 9. Instead of an absorbing GaAs substrate, one starts with a transparent GaP substrate. The graded layer has an inverse gradient relative to the chip shown in Fig. 7. The initial growth is 100 percent GaP phasing in As linearly over 10 to 15 μm . At 15 percent As, the emission is in the yellow range (585 nm), at 25 percent in the orange range (605 nm), and at 35 percent in the red range (635 nm). Figure 3 shows the approximate band structure of this material system. The composition range mentioned above has an indirect band structure. To obtain efficient light emission, the region of minority-carrier injection is doped with nitrogen forming an isoelectronic recombination center (see exciton recombination in the first section).

Figure 9 shows an important technique to increase extraction efficiency. In a TS-chip, the major light loss is due to free carrier absorption at the alloyed contacts. Rather than covering the entire bottom surface of the chip with contact metal, one can reduce the contact area either by depositing small contact islands (see Fig. 6) or by placing a dielectric mirror (deposited SiO_2) between the substrate and the unused areas of the back contact (Fig. 9). This dielectric mirror increases the efficiency by 20 to 50 percent at the expense of higher manufacturing cost.

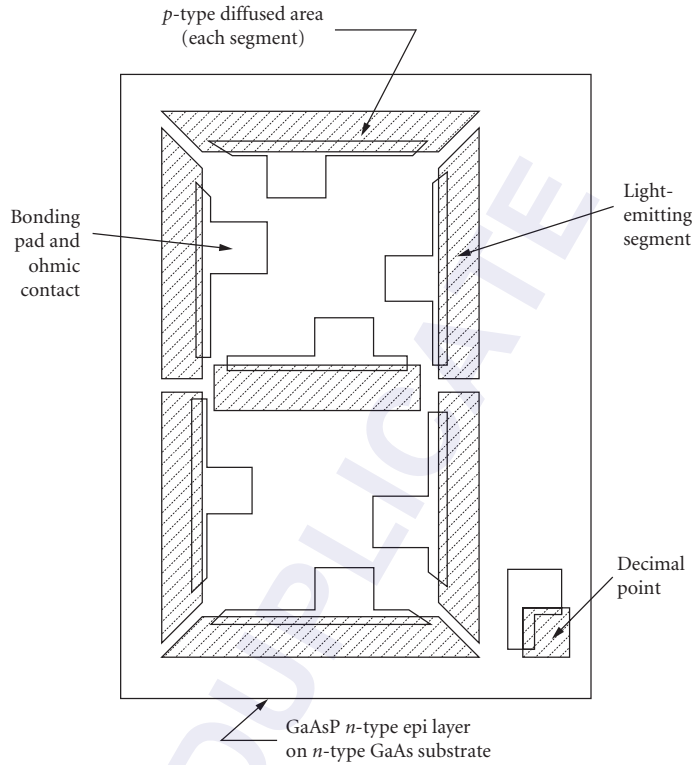


FIGURE 8 Monolithic seven-segment display chip with eight separate diffused regions (anodes) and a common cathode (the GaAs substrate).

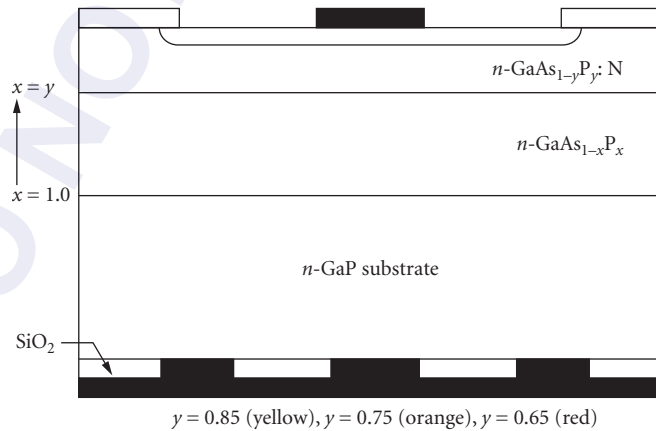


FIGURE 9 Cross section of a $\text{GaAs}_{1-x}\text{P}_x$ LED which, by changing the composition “ x ,” produces red, orange, or yellow light. The top layer is doped with nitrogen to increase the quantum efficiency. The GaP substrate is transparent to the emitted light. Also shown is a reflective back contact made by depositing the contact metallization on top of SiO_2 .

Single Heterojunctions

Heterojunctions introduce a new variable: local variation of the energy bandgap resulting in carrier confinement. Figure 10 shows a popular structure for a red LED emitter chip. A p -type layer of GaAlAs with 38 percent Al is grown on a GaAs substrate. The GaAlAs alloy system can be lattice-matched to GaAs; therefore, no graded layer is required as in the GaAsP system of Fig. 7. Next, an n -type layer is grown with 75 percent Al. The variation of E_g from substrate to top is illustrated in Fig. 11. Holes accumulate in the GaAlAs p layer with the narrower bandgap. Electrons are injected from the n layer into this p layer. The holes have insufficient energy to climb the potential barrier into the wide-bandgap material. Holes are confined to the p layer. In the p layer, the radiative recombination time is very short because of the high concentration of holes. As a result, the internal quantum efficiency is quite high. A variation of this structure is a widely used infrared emitter that emits at 880 nm.

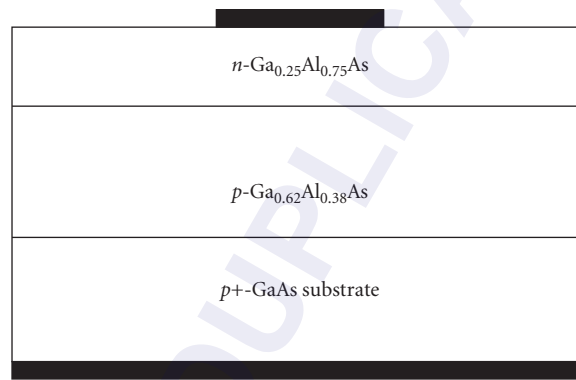


FIGURE 10 Cross section of a single heterostructure LED.

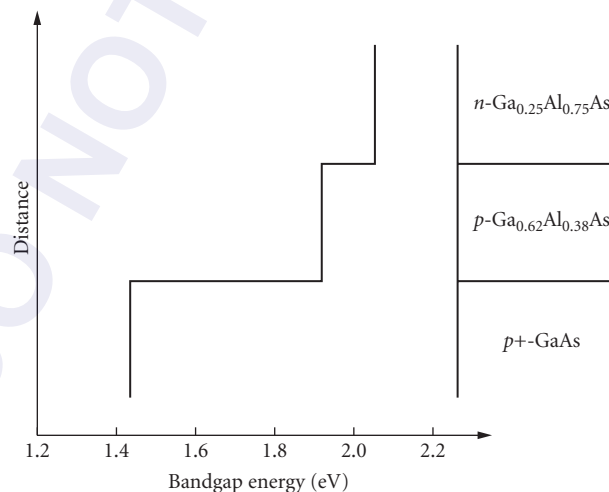


FIGURE 11 Variation in energy bandgap for the various layers in the GaAlAs LED shown in Fig. 10.

Double Heterostructures

The double heterostructure shown in Fig. 12 repeats the p -side confinement of Fig. 10 on the n side. An n -type buffer layer is grown on the GaAs substrate to create a high-quality surface onto which the first n -type GaAlAs confinement layer with 75 percent Al is grown. The active or light-generation layer is a 3- μm -thick p -type layer with 38 percent Al. The top p -type confinement layer again uses 75 percent Al. This structure with the energy-band diagram of Fig. 13 has two advantages: (1) There is no hole injection into the n -type layer with reduced efficiency and a slow hole recombination in the lowly doped n layer. (2) The high electron and hole density in the active layer reduces τ_r , thus increasing device speed and efficiency. The increased speed is quite important for LED sources in fiber-optic communication applications. (See “Fiber Optics” subsection later in this chapter.)

The double heterostructure of Fig. 12 represents a one-dimensional containment of injected carriers. Injection and light emission occurs across the entire lateral dimension of the chip. For fiber-optic applications, the light generated over such a large area cannot be effectively coupled into small-diameter fiber. A rule of thumb for fiber coupling requires that the light-emitting area be equal to or, preferably, smaller than the fiber core diameter. This rule requires lateral constraint of carrier injection. The localized diffusion of Fig. 7 is not applicable to grown structures such as those in Fig. 12. The preferred solution inserts an n layer between buffer and lower confinement layer (see Fig. 14). A hole etched into the n layer allows current flow. Outside of this hole, the p - n junction between n layer and lower confinement layer is reverse-biased, thus blocking any current flow. The disadvantage of this approach is a complication of the growth process. In a first growth process, the n layer is grown. Then a hole is etched into the n layer using standard photolithographic etching techniques. Finally, a second epitaxial growth is used for the remaining layers.

Another technique to constrain current injection utilizes a small ohmic contact.¹² It is used frequently in conjunction with InP-based fiber-optic emitters (see Fig. 15). An SiO_2 layer limits contact to a small-diameter (typically 25 μm) hole that results in a relatively small light-emitting area. The etched lens shown for this structure helps to collimate the light for more efficient coupling onto the fiber.¹²

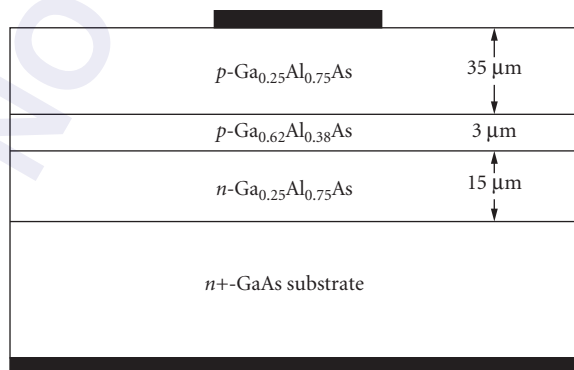


FIGURE 12 Structure of a double heterostructure (DH) GaAlAs LED. The DH is composed of a $\text{Ga}_{0.62}\text{Al}_{0.38}\text{As}$ layer surrounded on either side by a $\text{Ga}_{0.25}\text{Al}_{0.75}\text{As}$ layer. The thick top layer acts as window to increase light extraction through the side walls of the chip.

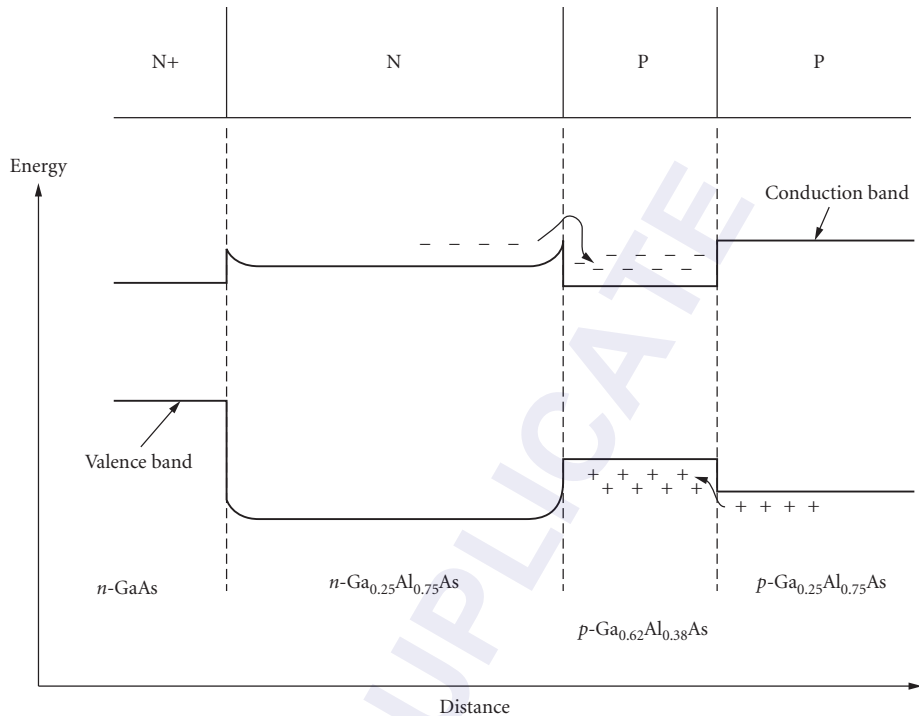


FIGURE 13 Energy band diagram of the GaAlAs LED shown in Fig. 12. The LED is forward biased. Electrons and holes are confined in the p -doped $\text{Ga}_{0.62}\text{Al}_{0.38}\text{As}$ layer, which increases the radiative recombination efficiency.

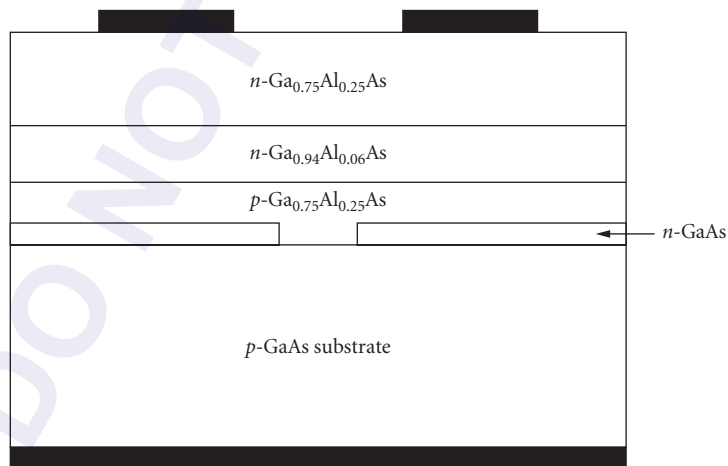


FIGURE 14 Cross section of an LED with three-dimensional carrier confinement. A DH structure is used to confine injected carriers in the $\text{Ga}_{0.94}\text{Al}_{0.06}\text{As}$ layer (direction perpendicular to the junction). The patterned, n -type GaAs layer is used to limit current flow in the lateral direction. The small emitting area and the 820-nm emission of this LED makes it ideal for fiber-optic applications.

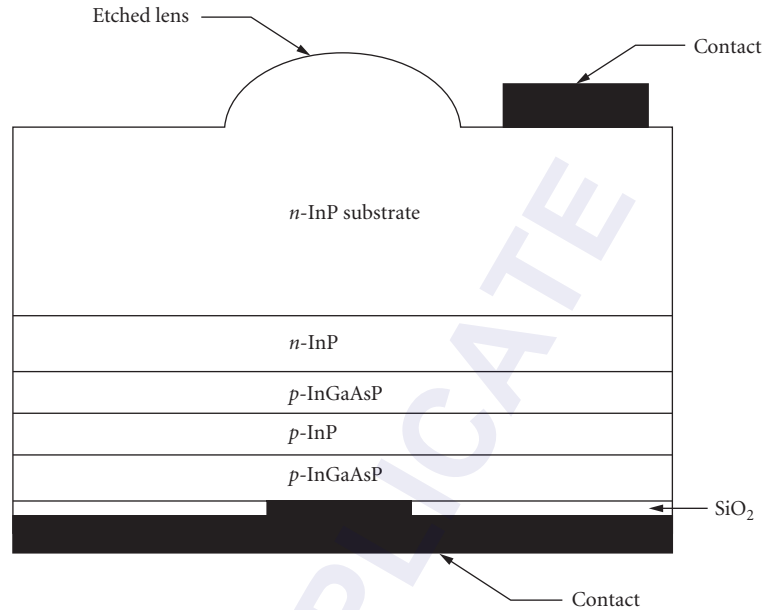


FIGURE 15 Structure of a 1300-nm LED used for optical fiber communications. The cross section shows the DH-layer configuration, limited area back contact for emission-size control, and an etched lens at the top of the chip which, magnifies ($M = 2$) the source area and collimates the light for effective coupling into a fiber.

17.6 MATERIAL SYSTEMS

The $\text{GaAs}_{1-x}\text{P}_x$ System

The most widely used alloy for LEDs is the ternary $\text{GaAs}_{1-x}\text{P}_x$ system, including its two binary components GaAs and GaP. This system is best described by the composition parameter x with $0 \leq x \leq 1$. For $x = 0$, we have GaAs and for $x = 1$ the composition is GaP. For $x \leq 0.4$, the alloy has a direct bandgap. GaAs was developed in the early 1960s as an infrared emitter with a wavelength of 910 nm and an efficiency in the range of 1 percent. This emitter was soon followed by a Si-doped variety. As discussed earlier, this leads to an emission wavelength of 940 nm, a wavelength at which the GaAs substrate is partly transparent. The resulting efficiency is increased substantially and, depending on configuration, is in the 5- to 10-percent range. However, the recombination process is quite slow, resulting in rise and fall times in the 50-ns to 1.0- μs range (see Table 2). One other drawback is caused by the low absorption coefficient of Si detectors at 940 nm. To absorb 90 percent of the light requires a detector thickness of 60 to 70 μm . Conventional photo transistors are quite suitable as detectors. Integrated photo ICs with their 5- to 7- μm -thick epitaxial layers are very inefficient as detectors at 940 nm.

To shift the wavelength toward the near-infrared or into the visible spectrum, one has to grow a ternary alloy, a mixture between GaAs and GaP. Of commercial interest are two alloys with $x = 0.3$ and $x = 0.4$ grown on a GaAs substrate (see Table 2). The $x = 0.4$ alloy was the first commercially produced material with a wavelength in the visible range of the spectrum. Grown on an absorbing substrate, it has a modest luminous efficacy of around 0.2 lm/A. [Luminous efficacy is the luminous (visible) flux output measured in lumens divided by the electrical current input.] The absorbing substrate allows the integration of multiple light sources into a single chip without crosstalk. Such

TABLE 2 Performance Summary of LED Chips in the GaAs_{1-x}P_x System

x	Growth Process*	Isoel. Dopant	Substrate†	Dominant Wavelength (nm)	Colon‡	Luminous Efficacy§ (lm/A)	Quantum Efficiency§ (%)	Speed (ns)
0	LPE	—	AS	910	IR		1	50
0	LPE	—	TS	940	IR		10	1000
0.3	VPE	—	AS	700	IR		0.5	50
0.4	VPE	—	AS	650	Red	0.2		50
0.65	VPE	N	TS	635	Red	2.5		300
0.75	VPE	N	TS	605	Orange	2.5		300
0.85	VPE	N	TS	585	Yellow	2.5		300
1.0	LPE	N	TS	572	Y/green	6		300
1.0	LPE	—	TS	565	Green	1		
1.0	LPE	ZnO	TS	640	Red	1		

*LPE = liquid phase epitaxy; VPE = vapor phase epitaxy.

†AS = absorbing substrate; TS = transparent substrate.

‡IR = infrared.

§Into plastic (index = 1.5).

monolithic seven-segment displays became the workhorse display technology for handheld calculators from 1972 to 1976.^{9,13} Today this alloy is used in LED arrays for printers.¹⁴

The $x = 0.3$ alloy with a wavelength of 700 nm became important in the mid 1970s as a light source in applications using integrated photodetectors. It has 3 to 5 times the quantum efficiency of the $x = 0.4$ alloy (see Table 2), but has a lower luminous efficacy because of the much-reduced eye sensitivity at 700 nm.

For $x > 0.4$, the GaAsP material system becomes indirect (see earlier under “Light-Generation Processes” and Fig. 3). The quantum efficiency decreases faster than the increase in eye sensitivity.⁴ The only way to achieve a meaningful efficiency is through the use of isoelectronic dopants as described earlier. The choices for isoelectronic dopants that have been successful are N for GaAsP¹¹ and N or ZnO for GaP. Nitrogen doping is used widely for alloys with $x = 0.65$ to $x = 1.0$.⁶⁻¹⁰ The resulting light sources cover the wavelength range from 635 nm to approximately 565 nm (see Table 2). Since these alloys are either GaP or very close in composition to GaP, they are all grown on GaP substrates. The resulting transparent-substrate chip structure increases the luminous efficacy.

In the case of the binary GaP compound, the dominant wavelength depends on N concentration. With low concentrations, practically all N atoms are isolated in single sites. With increasing concentration, N atoms can arrange themselves as pairs or triplets. The resulting electron traps have lower energy states, which shift the emitted light toward longer wavelength. Phonon coupling can also reduce the emission energy. Commercially significant are two compositions: (1) undoped GaP which emits at 565 nm (dominant wavelength) with a reasonably green appearance and a low efficiency of around 1 lm/A, and (2) GaP with a nitrogen concentration in the range of 10^{19} cm⁻³ with a substantially higher luminous efficiency of around 6 lm/A at 572 nm. At this wavelength, the color appearance is yellow-green, often described as chartreuse.

Three nitrogen-doped ternary alloys of GaAsP are commercially important for red, orange, and yellow. The red source with $x = 0.65$ has an efficiency in the range of 2 to 3 lm/A. With increasing bandgap or decreasing wavelength, the drop in quantum efficiency is compensated by an increase in eye sensitivity, resulting in a practically wavelength-independent luminous efficiency for the range of 635 to 585 nm.^{8,15}

ZnO-doped GaP is an interesting material. The quantum efficiency of such chips is relatively high, around 3 percent. However, the linewidth is quite broad. The quantum efficiency peaks at 700 nm, but the luminous efficiency peaks at 640 nm (dominant wavelength). In other words, most of the photons are emitted at wavelengths with low eye sensitivity. Another problem of GaP:ZnO is saturation. The deep ZnO electron trap causes very slow exciton recombination. At high injection currents, all traps are saturated and most of the injected carriers recombine

nonradiatively. At low injection levels (≤ 1 A/cm²), the efficacy is relatively high, 3 to 5 lm/A. At a more useful density of 10 to 30 A/cm², the emission saturates, resulting in an efficacy of around 1 lm/A.

The Al_xGa_{1-x}As System

The Al_xGa_{1-x}As material system has a direct bandgap for $0 \leq x \leq 0.38$. This system has one very significant advantage over the GaAsP system described earlier, the entire alloy range from $x = 0$ to $x = 1$ can be lattice-matched to GaAs. In other words, every alloy composition can be directly grown on any other alloy composition without the need for transition layers. This feature allows the growth of very abrupt heterojunctions, i.e., abrupt transition in composition and bandgap. These heterojunctions add one important property not available in the GaAsP system: carrier containment (see earlier under “Device Structures”). Carrier containment reduces the movement of injected carriers in a direction perpendicular to the junction. Thus, carrier density can be increased beyond the diffusion-limited levels. This results in increased internal quantum efficiency and higher speed. Another benefit is reduced absorption and improved extraction efficiency (under “Light Extraction”).

Of practical significance are two compositions: $x = 0.06$ and $x = 0.38$ (see Table 3). Both compositions exist in single and double heterojunction variations (see under “Device Structures”). The double heterojunctions usually have a 1.5 to 2.0 times advantage in efficiency and speed. In all cases, the efficiency strongly depends on the thickness of the window layer and, to a lesser degree, on the thickness of the transparent layer between active layer and absorbing substrate (see Fig. 12). Chips with a transparent substrate have an additional efficiency improvement of 1.5 to 3.0 times again, depending on layer thickness and contact area. The efficiency variation is best understood by counting exit cones as described in the text in conjunction with Table 1. For $x = 0.06$, the internal quantum efficiency of a double heterojunction approaches 100 percent. For $x = 0.38$, the direct and indirect valleys are practically at the same level and the internal quantum efficiency is reduced to the range of 50 percent, again dependent upon the quality of the manufacturing process.

The best compromise for efficiency and speed is the $x = 0.06$ alloy as a double heterostructure. Depending on layer thickness, substrate, and contact area, these devices have efficiencies of 5 to 20 percent and rise/fall times of 20 to 50 ns. This alloy is becoming the workhorse for all infrared applications demanding power and speed. A structural variation as shown in Fig. 14 is an important light source for fiber-optic communication.

The $x = 0.38$ alloy is optimized for applications in the visible spectrum. The highest product of quantum efficiency and eye response is achieved at $x = 0.38$ and $\lambda = 650$ nm. The single heterostructure on an absorbing substrate has an efficacy of around 4 lm/A. The equivalent double heterostructure is in the 6- to 8-lm/A range. On a transparent GaAlAs substrate, the efficacy is typically in the 15- to 20-lm/A range and results of as high as 30 lm/A have been reported in the literature.¹⁶ The major application for these red LEDs is in light-flux-intensive applications, such as message panels and automotive stoplights. A variation optimized for speed is widely used for optical communication using plastic fiber.

TABLE 3 Performance Summary of LED Chips in the Al_xGa_{1-x}As System

x	λ (nm)	Substrate*	Structure [†]	Efficiency or Efficacy	Speed (ns)
0.06	820	AS, TW	DH	8%	30
0.06	820	TS	DH	15%	30
0.38	650	AS, TW	SH	4 lm/A	
0.38	650	AS, TW	DH	8 lm/A	
0.38	650	TS	DH	16 lm/A	

*AS = absorbing substrate; TW = thick window layer; TS = transport substrate.

[†]DH = double heterostructure; SH = single heterostructure.

The AlInGaP System

The AlInGaP system has most of the advantages of the AlGaAs system with the additional advantage that it has a higher-energy direct energy gap of 2.3 eV that corresponds to green emission at 540 nm. AlInGaP can be lattice-matched to GaAs substrates. Indium occupies about half of the Group III atomic sites. The ratio of aluminum to gallium can be changed without affecting the lattice match, just as it can in the AlGaAs material system, since AlP and GaP have nearly the same lattice spacing. This enables the growth of heterostructures that have the efficiency advantages described in the previous section.

Various AlInGaP device structures have been grown. A simple DH structure with an AlInGaP active layer surrounded by higher bandgap AlInGaP confining layers has been effective for injection lasers, but has not produced efficient surface-emitting LEDs.¹⁷ The main problem has been that AlInGaP is relatively resistive and the top AlInGaP layer is not effective in distributing the current uniformly over the chip. This is not a problem with lasers since the top surface is covered with metal and the light is emitted from the edge of the chip.

The top layer must also be transparent to the light that is generated. Two window layers used are AlGaAs or GaP on top of the AlInGaP heterostructure.^{18,19} AlGaAs has the advantage that it is lattice-matched and introduces a minimum number of defects at the interface, but it has the disadvantage that it is somewhat absorptive to the yellow and green light which is generated for high-aluminum compositions. The highest AlInGaP device efficiencies have been achieved using GaP window layers.^{18,20} GaP has the advantage that it is transparent to shorter wavelengths than AlGaAs and that it is easy to grow thick GaP layers, using either VPE or LPE, on top of the AlInGaP DH that was grown by MOVPE. Both VPE and LPE have substantially higher growth rates than MOCVD. The various growth techniques are discussed later under "Epitaxial Technology."

AlInGaP devices with 45- μm -thick GaP window layers have achieved external quantum efficiencies exceeding 5 percent in the red and yellow regions of the emission spectrum.²¹ This is more than twice as bright as devices that have thinner AlGaAs window layers.

Green-emitting AlInGaP devices have also been grown which are brighter than the conventional GaP and GaP:N green emitters.^{18,20,21} Substantial further improvement in green is expected since the quantum efficiency is not as high as would be expected based on the energy position of the transition from a direct to an indirect semiconductor.

The performance of AlInGaP LEDs compared to the most important other types of visible emitters is shown in Figs. 16 and 17. GaAsP on a GaAs substrate and GaP: ZnO are not shown since their luminous efficacy are off the chart at the bottom and the lower-right-hand corner, respectively. It is clear from Fig. 17 that the luminous efficacy of AlInGaP is substantially higher than the other technologies in all color regions except for red beyond 640 nm. Since forward voltage is typically about 2 V, the lumen per watt value for a given device is about one-half the lumen per ampere value that is given in Tables 2 and 3. The quantum efficiency of AlInGaP is also better than all of the other technologies except for the highest-performance AlGaAs devices operating at about 650 nm. Because of the eye sensitivity variations (see C.I.E. curve in Fig. 16), the 620 nm (red/orange) AlInGaP devices have a higher luminous efficacy than 650 nm AlGaAs LEDs (see Fig. 17).

Blue LED Technology

Blue emitters have been commercially available for more than a decade, but have only begun to have a significant impact on the market in the last few years. SiC is the leading technology for blue emitters with a quantum efficiency of about 0.02 percent and 0.04 lm/A luminous performance. SiC devices are not much used due to their high price and relatively low performance efficiency.

Other approaches for making blue LEDs are the use of II–VI compounds such as ZnSe or the nitride system GaN, AlGaN, or AlGaInN. It has been difficult to make good *p-n* junctions in these materials. Recently improved *p-n* junctions have been demonstrated in both ZnSe^{22,23} and GaN.²⁴

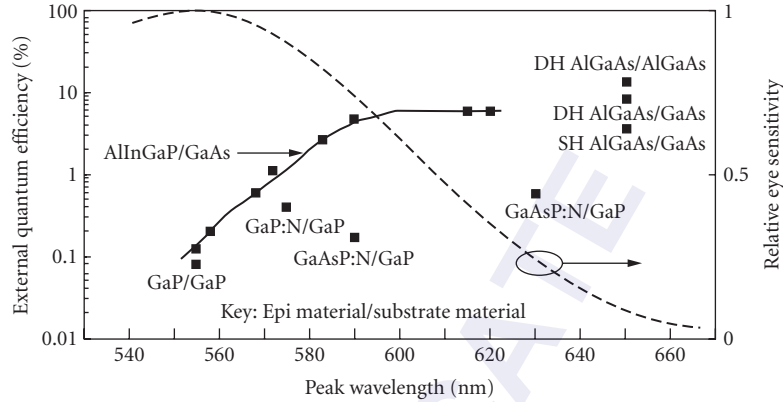


FIGURE 16 External quantum efficiency as a function of peak wavelength for various types of visible LEDs. Below 590 nm the efficiency of AlInGaP LEDs decreases due to the approaching transition from direct to indirect semiconductor. The human eye sensitivity curve is also shown. Since the eye response increases sharply from 660 to 540 nm, it partially makes up for the drop in AlInGaP LED efficiency. The resulting luminous performance is shown in Fig. 17.

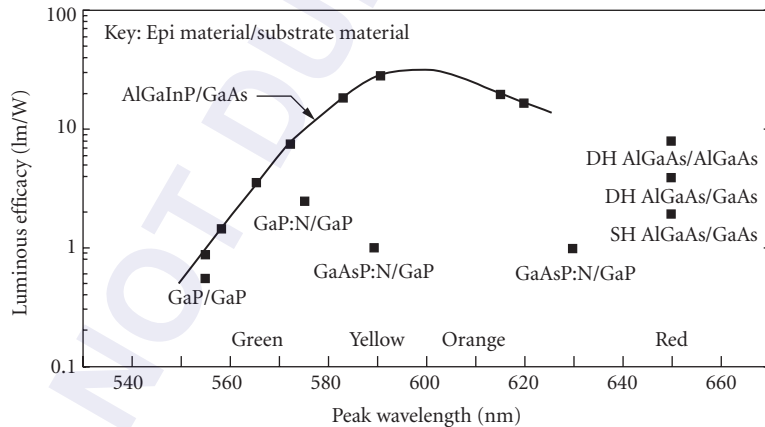


FIGURE 17 Luminous efficacy for AlInGaP LEDs versus wavelength compared to other LED technologies. AlInGaP LEDs are more than an order of magnitude brighter in the orange and yellow regions than other LEDs. AlInGaP LEDs compare favorably to the best AlGaAs red LEDs.

Device performance is still in the 0.1-lm/A range and reliability is unproven. However, this recent progress is very encouraging for blue-emission technology and could lead to a high-performance device in the next few years. Both ZnSe and the nitride system have a major advantage over SiC because they are direct bandgap semiconductors, so a much higher internal quantum efficiency is possible. However, it is difficult to find suitable lattice-matched substrates for these materials.

17.7 SUBSTRATE TECHNOLOGY

Substrate Criteria

There are several requirements for substrates for LEDs. The substrate must be as conductive as possible both thermally and electrically to minimize power loss. In order to minimize defects it should match the epitaxial layers as closely as possible in atomic lattice spacing, and in the coefficient of thermal expansion. The substrate should also have a low defect density itself. Finally, the substrate should, if possible, be transparent to the light generated by the LED structure since this will enhance the external quantum efficiency.

Substrate Choices

The substrates used for nearly all visible LEDs are GaAs and GaP. GaAs or InP is used for infrared devices, depending upon the device structure required. Substrate parameters are summarized in Table 4, along with Si and Ge for comparison. GaAs is used for the AlGaAs and AlInGaP material systems since they can be lattice-matched to it. GaAs is also used for the GaAs_{1-x}P_x system for $x \leq 0.4$ because it is more nearly lattice-matched and, since it is absorbing, it is useful in multiple-junction devices where optical crosstalk must be minimized (see under “Diffused Homojunction” and Fig. 8). GaP is used for compositions of $x > 0.6$ due to its transparency and closer lattice match. However, neither GaP nor GaAs are well matched to GaAsP, and grading layers are required to grow epitaxial layers with decent quality. InP is the choice for long wavelength emitters made using the InGaAs or InGaAsP materials systems, due to the better lattice match.

Substrate Doping

Generally, substrates are *n* type and are doped with Te, S, or Si, although sometimes Se and Sn are also used. In some cases, particularly for some AlGaAs LEDs and laser structures, *p*-type substrates are required and Zn is nearly always the dopant. The doping levels are typically in the 10¹⁸-cm⁻³ range. Basically, the substrates are doped as heavily as possible to maximize conductivity. However, the doping must be below the solubility limit to eliminate precipitates and other defects. In the case of substrates that are transparent to the light which is generated, such as GaP, with a GaAs_{1-x}P_x epitaxial layer, the doping should also remain below the level at which substantial free carrier absorption occurs.

Growth Techniques

Substrates can be grown by either the Bridgeman or Liquid Encapsulated Czochralski (LEC) technique. The LEC technique is the most widely used. Both techniques are described in detail elsewhere and will be only briefly summarized here.²⁵

TABLE 4 Properties of Common Semiconductor Substrates

Substrate	Lattice Parameter	Energy Gap @ 300 K (eV)	Melting Point (°C)
GaAs	5.653	1.428 Direct	1238
GaP	5.451	2.268 Indirect	1467
InP	5.868	1.34 Direct	1062
Si	5.431	1.11 Indirect	1415
Ge	5.646	0.664 Indirect	937

The LEC technique for GaAs consists of a crucible containing a molten GaAs solution, into which a single crystal “seed” is dipped. The temperature is carefully controlled so that the molten GaAs slowly freezes on the seed as the seed is rotated and raised out of the molten solution. By properly controlling the temperature, rotation rate, seed lift rate, etc., the seed can be grown into a single crystal weighing several kilograms and having a diameter of typically 2 to 4 in for GaAs and for 2 to 2.5 in for GaP. At the GaAs and GaP melting points As and P would rapidly evaporate from the growth crucible if they were not contained with a molten boric oxide layer covering the growth solution. This layer is the reason for the name “liquid encapsulated.” The seed is dipped through the boric oxide to grow the crystal. The growth chamber must be pressurized to keep the phosphorus and arsenic from bubbling through the boric oxide. The growth pressure for GaP is 80 atm, and for GaAs is 20 atm or less, depending upon the approach for synthesis and growth. The LEC technique, used for GaAs, GaP, and InP, is similar to the Czochralski technique used for silicon, but the silicon process is much simpler since encapsulation is not required and the growth can be done at atmospheric pressure.

The Bridgeman and the gradient-freeze technique, which is a variation of the Bridgeman technique, can also be used to grow compound semiconductors. In this technique the growth solution and seed are contained in a sealed chamber so liquid encapsulation is not required. Growth is accomplished by having a temperature gradient in the solution, with the lowest temperature at the melting point in the vicinity of the seed. Growth can be accomplished by lowering the temperature of the entire chamber (gradient freeze), or by physically moving the growth chamber relative to the furnace (Bridgeman technique) to sweep the temperature gradient through the molten solution. The growing crystal can be in either a vertical or horizontal position.

GaAs for LEDs is commonly grown using either LEC or horizontal gradient freeze, also called “boat grown,” but sometimes a vertical Bridgeman approach is used. GaP and InP are almost always grown using LEC but sometimes a vertical Bridgeman approach has also been used.

17.8 EPITAXIAL TECHNOLOGY

Growth Techniques Available

Epitaxial layers are grown using one of several techniques, depending on the material system. The most common techniques are liquid phase epitaxy (LPE), which is primary used to grow GaAs, GaP, and AlGaAs and vapor phase epitaxy (VPE), which is used to grow GaAsP. Metalorganic vapor phase epitaxy (MOVPE) is also used to grow AlGaAs, GaInAsP, and AlInGaP. Molecular beam epitaxy (MBE) is used for lasers and high-speed devices but is not used for high-volume commercial LEDs at this time. It has been used to grow blue ZnSe-based lasers and LEDs and could be important in the future. All of these epitaxial techniques have been discussed extensively elsewhere and will be only briefly described here.

LPE

LPE growth consists of a liquid growth solution, generally gallium, which is saturated with the compound to be grown.²⁶ The saturated solution is placed in contact with the substrate at the desired growth temperature, and cooled. As the solution cools, an epitaxial film is grown on the substrate. The technique has the advantage that it is relatively easy to grow high-quality epitaxial layers, and materials containing aluminum (such as AlGaAs) can be readily grown. The disadvantages are that composition control can be difficult. Also, the growth of epitaxial structures involving more than two or three layers, particularly thin layers, can be mechanically complicated since each layer requires a separate growth solution that must be carefully saturated and sequentially brought into contact with the substrate.

One important use of the LPE technique has been for the growth of GaP:ZnO and GaP:N for red and green LEDs, respectively. These devices each consist of two relatively thick layers: an *n*-type layer,

followed by the growth of a *p*-type layer. While other growth techniques can be used to grow GaP LEDs, the best results have been obtained using LPE. As a result of a high-volume, low-cost production technology has evolved, which produces more visible LED chips than any other technique. Another major use of LPE is for the growth of GaAs: Si for infrared emitters. LPE is the only technique with which it has been possible to produce the recombination center that gives rise to the 940-nm emission characteristic of this material. The GaAs: Si structures are generally grown from a single growth solution. At high temperatures the silicon is incorporated on Ga sites and the layers are *n* type. As the solution cools the silicon becomes preferentially incorporated on the As sites and the layer becomes *p* type.

AlGaAs devices for both visible (red) and IR devices are also generally grown by LPE. AlGaAs devices can also be grown by MOCVD, but LPE has the advantage that thick layers can be more easily grown. This is important for high extraction efficiency (see under "Device Structures" and Fig. 12). In the case of visible devices at 650 nm, the internal quantum efficiency is also higher using LPE than MOCVD. This is not understood, but the result is that virtually all of the visible AlGaAs LEDs are produced using LPE.

VPE

VPE is the other major commercial epitaxial technology for LEDs.^{9,27} VPE consists of a quartz chamber containing the substrate wafers at the appropriate growth temperature. The reactants are transported to the substrates in gaseous form. The technique is mainly used for the growth of GaAsP which, along with GaP, dominates the high-volume visible LED market. In this case HCl is passed over Ga metal to form gallium chlorides, and AsH₃ and PH₃ are used to provide the As and P compounds. Appropriate dopant gases are added to achieve the *n*- and *p*-type doping. NH₃ is used to achieve nitrogen doping for the growth of GaAsP:N. The VPE technique has the advantages that it is relatively easy to scale up the growth chamber size so large quantities of material can be grown and layer composition and thickness can be easily controlled by adjusting the flow conditions. A limitation of VPE is that it has not been possible to grow high-quality compounds containing aluminum because the aluminum-bearing reactants attack the quartz chamber resulting in contaminated films. Thus, AlGaAs and AlInGaP, the emerging high-brightness technologies, cannot be grown using this technique.

MOCVD

MOCVD growth, like conventional VPE, uses gases to transport the reactants to the substrates in a growth chamber.²⁸ However, in this case metallorganic compounds such as trimethylgallium (TMG) are used for one or more of the reactants. A major difference between VPE and MOCVD is that in the case of MOCVD the decomposition of the source gas (e.g., TMG) occurs as a reaction at or near the substrate surface, and the substrate is in the hottest area of the reactor such that the decomposition occurs on the substrate instead of the walls of the growth chamber. The walls of the growth chamber remain relatively cool. This is the key factor that makes MOCVD suitable for the growth of aluminum-bearing compounds which, unlike the VPE situation, do not react significantly with the cooler reactor walls. Thus AlGaAs and AlInGaP can be readily grown with MOCVD, and this technology is widely used for infrared AlGaAs LEDs and lasers, and for the emerging visible AlInGaP laser and LED technology.

MOCVD is also used for the growth of GaN and AlGaN that are candidates for blue emission, and for the growth of II–VI compounds, such as ZnSe, that is also a potential blue emitter. However, at this time the key limitation in obtaining blue ZnSe emitters is the growth of low-resistivity *p*-type ZnSe. For reasons that are not yet understood, low-resistivity *p*-type ZnSe has been grown using MBE only.

MBE

MBE is a high vacuum growth technique in which the reactants are essentially evaporated onto the substrates under very controlled conditions.²⁹ MBE, like MOCVD, can be used to deposit compounds

containing aluminum. The growth rates using MBE are generally slower than the other epitaxial techniques, so MBE is most suitable for structures requiring thin layers and precise control of layer thickness. MBE equipment is somewhat more expensive than the equipment used for the other types of epitaxial growth, so it has not been suitable for the high-volume, low-cost production that is required for most types of LEDs. MBE has generally been utilized for lasers and high-speed devices where control of complicated epitaxial structures is critical and where relatively low volumes of devices are required. One advantage that MBE has over the other growth technologies is that the reactants utilized are generally less hazardous. Consequently, MBE equipment is often cheaper and easier to install since there are less safety issues and safety-code restrictions to deal with.

17.9 WAFER PROCESSING

Wafer Processing Overview

Wafer processing of compound semiconductors for LED applications has many of the same general steps used to process silicon integrated-circuit wafers, namely passivation, diffusion, metallization, testing, and die fabrication. The LED device structures are much simpler, so fewer steps are required; but, due to the materials involved, the individual steps are generally different and sometimes more complicated.

Compound semiconductor processing has been described in detail elsewhere, so only a brief summary is discussed here.³⁰ Some types of LED structures require all of the processing steps listed here, but in many cases fewer process steps are required. An example is a GaP or AlGaAs device with a grown *p-n* junction. For these devices no passivation or diffusion is required.

Passivation

Some types of LED structures, particularly multijunction structures, require a passivation layer prior to diffusion, as shown in Fig. 7. This layer must be deposited relatively free of pinholes, be patternable with standard photolithographic techniques, and must block the diffusing element, generally zinc. In the case of silicon, a native oxide is grown which is suitable for most diffusions. Unfortunately, the compound semiconductors do not form a coherent native oxide as readily as silicon. Silicon nitride (Si_3N_4) is the most widely used passivation layer for LEDs. Si_3N_4 is grown by reacting silane (SiH_4) and ammonia at high temperature in a furnace. Si_3N_4 blocks zinc very effectively, and is easily grown, patterned, and removed. Sometimes an SiO_2 layer is used in conjunction with Si_3N_4 for applications such as protecting the surface of the compound semiconductor during high-temperature processing. Silicon oxynitride can also be used instead of or in addition to pure Si_3N_4 . Silicon oxynitride is somewhat more complicated to deposit and control, but can have superior properties, such as a better match of coefficient of expansion, resulting in lower stress at the interface.

Diffusion

Generally, only *p*-type impurities, usually Zn, are diffused in compound semiconductors. *N*-type impurities have prohibitively small diffusion coefficients for most applications. Zn is commonly used because it diffuses rapidly in most materials and because it is nontoxic in contrast to Be, which also diffuses rapidly. Mg is another reasonable *p*-type dopant, but it diffuses more slowly than zinc. Diffusions are generally done in evacuated and sealed ampoules using metallic zinc as source material. A column V element such as As is also generally added to the ampoule to provide an overpressure that helps to prevent decomposition of the semiconductor surface during diffusion. Diffusion conditions typically range from 600 to 900°C for times ranging from minutes to days, depending upon the material and device involved. Junction depths can range from a fraction of a μm to a more than 10 μm .

Open-tube diffusions have also been employed but have generally been harder to control than the sealed ampoule approach, often because of surface decomposition problems. Open-tube diffusions have the advantage that one does not have to deal with the expense and hazard of sealing, breaking, and replacing quartz ampoules.

A third type of diffusion that has been used is a “semisealed” ampoule approach in which the ampoule can be opened and reused. The diffusion is carried out at atmospheric pressure and the pressure is controlled by having a one-way pressure relief valve on the ampoule.

Contacting

The contacts must make good ohmic contact to both the *p*- and *n*-type semiconductor, and the top surface of the top contact must be well suited for high-speed wire bonding. Generally, multilayer contacts are required to meet these conditions. Evaporation, sputtering, and *e*-beam deposition are all employed in LED fabrication. The *p*-type contact generally uses an alloy of either Zn or Be to make the ohmic contact. An Au-Zn alloy is the most common due to the toxicity of Be. The Au-Zn can be covered with a layer of Al or Au to enable high-yield wire bonding. A refractory metal barrier layer may be included between the Au-Zn and top Al or Au layer to prevent intermixing of the two layers and the out-diffusion of Ga, both of which can have a deleterious effect on the bondability of the Al or Au top layer. The *n*-type contact can be similar to the *p*-type contact except that an element which acts as an *n*-type dopant, commonly Ge, is used instead of Zn. An Au-Ge alloy is probably most frequently used to form the *n*-type ohmic contact since it has a suitably low melting point. If the *n*-type contact is the top, or bonded, contact it will be covered by one or more metallic layers to enhance bondability.

Testing

The key parameters that need to be tested are light output, optical rise and fall times, emission wavelength, forward voltage, and leakage current. The equipment used is similar to that used to test other semiconductor devices except that a detector must be added to measure light output. Rise and fall times and wavelength are generally measured on only a sample basis and not for each device on a wafer. In order to test the individual LEDs, the devices must be isolated on the wafer. This occurs automatically for LEDs that are masked and diffused, but if the LEDs are sheet diffused or have a grown junction, the top layer must be processed to isolate individual junctions. This can be accomplished by etching or sawing with a dicing saw. Generally, sawing is used, followed by an etch to remove saw damage, because the layers are so thick that etching deep (>10 μm) grooves is required. It is advisable to avoid deep groove etching because undercutting and lateral etching often occur and the process becomes hard to control. In many cases LED junctions are not 100 percent tested. This is particularly true of GaP and AlGaAs red-emitting devices in which the top layer may be 30 μm thick. Wafers of this type can be sampled by “coring” through the top layer of the wafer in one or more places with an ultrasonic tool in order to verify that the wafer is generally satisfactory. Later, when chips are fully processed, chips can be selected from several regions of the wafer and fully tested to determine if the wafer should be used or rejected.

Die Fab

Die fab is the process of separating the wafer into individual dice so they are ready for assembly. Generally, the wafer is first mounted on a piece of expandable tape. Next the wafer is either scribed or sawed to form individual dice. Mechanical diamond scribing or laser scribing were the preferred technologies in the past. Mechanical scribing has zero kerf loss, but the chips tend to have jagged edges and visual inspection is required. Laser scribing provides uniform chips but the molten waste material from between the chips damages neighboring chips, and in the case of full function chips

the edges of the junction can be damaged by the laser. As a result of the limitations of scribing, sawing (using a dicing saw with a thin diamond impregnated blade) has become the technology of choice for most LEDs.

The kerf loss for sawing has been reduced to about 40 μm and the chip uniformity is excellent such that a minimum of inspection and testing is required. For most materials a “cleanup” etch is required after sawing to remove work damage at the edges of the chips, which can both affect the electrical performance and absorb light. The wafer remains on the expandable tape during the sawing process. After sawing, the tape is expanded so that the chips are separated. The tape is clamped in a ring that keeps it expanded and the chips aligned. In this form the chips are easily individually picked off the tape by the die-attach machine that places the chips in the LED package.

17.10 LED QUALITY AND RELIABILITY

LEDs offer many advantages over other types of light sources. They have long operating life, they operate over a wide temperature range, and they are unaffected by many adverse environmental conditions. LED devices also are mechanically robust, making them suitable for applications where there is high vibration, shock, or acceleration. Excellent quality and reliability are obtainable when an LED product is properly designed, fabricated, packaged, tested, and operated.

Product quality is defined as “fitness for use” in a customer’s application. Quality is measured in units of the average number of defective parts per million shipped (i.e., ppm), and is inferred from product sampling and testing. LED product quality is assured by (1) robust chip and product design, (2) high-quality piece parts, (3) well-controlled fabrication processes, (4) use of statistical process control during manufacturing, (5) careful product testing, and (6) proper handling and storage. Most III–V LEDs are comparable in quality to the best silicon devices manufactured today. Well-designed LED products have total defect levels well below 100 ppm.

Reliability measures the probability that a product will perform its intended function under defined use conditions over the useful life of the product. Probability of survival is characterized by a failure rate, which is calculated by dividing the number of failures by the total number of operating hours (number of products tested per x hours operated). Common measures for reliability are percent failures per 1000 hours (percent/khr) and number failures per 10^9 hours (FITS). LED failure rates typically are better than 0.01 percent per khr at 50°C.

The reliability of an LED product is dependent on the reliability of the LED semiconductor chip and on the robustness of the package into which the chip is placed. Interactions between the chip and package can affect product reliability as well. Aspects of LED packaging and LED chip reliability are discussed in the following paragraphs.

LED Package Reliability

The package into which the LED chip is assembled should provide mechanical stability, electrical connection, and environmental protection. To evaluate package integrity, stress tests such as temperature cycling, thermal and mechanical shock, moisture resistance, and vibration are used to establish the worst-case conditions under which a product can survive. Generally, product data sheets contain relevant information about safe conditions for product application and operation.

Plastic materials are commonly used to package LEDs (see under “LED-Based Products”). Thermal fatigue is a limiting factor in plastic-packaged LEDs. Take the case of the plastic LED lamp shown later in Fig. 21. Because of the different materials used (epoxy plastic, copper lead frame, gold wire, III–V LED, etc.) and the different coefficients of expansion of these materials, temperature changes cause internal stresses. If the package is not well designed and properly assembled, thermal changes can cause cracking, chip-attach failure, or failure of the wire bond (open circuits). Careful design can reduce these problems to negligible levels over wide temperature ranges. Today’s high-quality plastic lamps are capable of being cycled from -55 to $+100^\circ\text{C}$ for 100 cycles without failure.

Long-term exposure to water vapor can lead to moisture penetration through the plastic, subjecting the chips to humidity. High humidity can cause chip corrosion, plastic delamination, or surface leakage problems. Plastic-packaged LEDs are typically not harmed when used under normal use conditions. Accelerated moisture resistance testing can be used to test the limits of LED packages. Plastic materials have been improved to the point that LED products can withstand 1000 hours of environmental testing at elevated temperatures and high humidity (i.e., 85°C and 85 percent RH).

The thermal stability of plastic packaging materials is another important parameter. Over normal service conditions, the expansion coefficient of plastic is relatively constant. Above the so-called glass transition temperature T_g , the coefficient increases rapidly. Reliable operation of plastic-packaged LEDs generally requires operation at ambient temperatures below T_g . Failures associated with improper soldering operations, such as too high a soldering temperature or for too long a time can cause the package to fail. Excessive storage temperatures also should be avoided. When an LED product is operated, internal ohmic heating occurs; hence, the safe operating maximum temperature is generally somewhat lower than the safe storage temperature.

LED Chip Reliability

The reliability characteristics of the LED chip determine the safe limits of operation of the product. When operated at a given temperature and drive current there is some probability that the LED will fail. In general, LED failure rates can be separated into three time periods: (1) infantile failure, (2) useful life, and (3) the wearout period (see Fig. 18). During the infant mortality period, failures occur due to weak or substandard units. Typically, the failure rate decreases during the infantile period until no weak units remain. During the “useful life” period, the failure rate is relatively low and constant. The number of failures that do occur are random in nature and cannot be eliminated by more testing. The useful life of an LED is a function of the operating temperature and drive conditions. Under normal use conditions, LEDs have useful lives exceeding 100,000 hours. The wearout period is characterized by a rapidly rising failure rate. Generally, wearout for LEDs is not a concern, as the useful life far exceeds the useful life of the product that the LEDs are designed into.

The principle failure mode for LED chips is light-output degradation. In the case of a visible lamp or display, failure is typically defined as a 50 percent decrease in light output from its initial

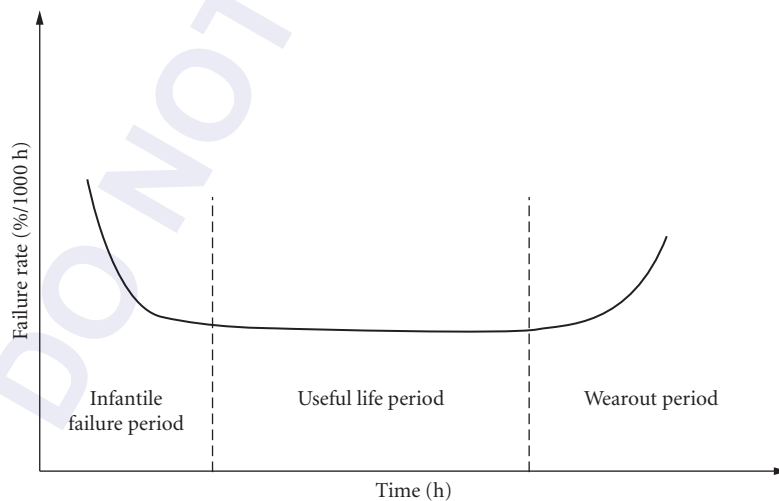


FIGURE 18 Plot of LED failure rate versus time, showing the infantile failure period (decreasing rate), the useful life period (constant, low rate), and the wear-out period (rapidly rising rate).

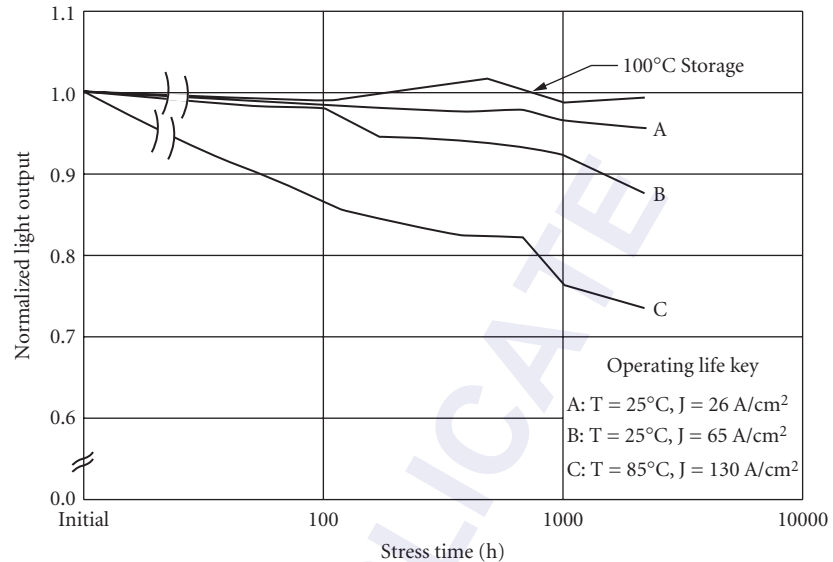


FIGURE 19 Curves of light output versus time for GaAsP indicator lamps stressed under various conditions. Light output is normalized to the initial value. Each curve shows the average degradation of 20 lamps

value, since that is the level where the human eye begins to observe a noticeable change. For infrared emitters or visible LEDs where flux is sensed by a semiconductor detector, failure is commonly defined as a 20 to 30 percent decrease in flux output.

Figure 19 shows degradation curves for a direct bandgap GaAsP LED packaged in a 5-mm plastic lamp.³¹ Current must flow for degradation to occur, as negligible change is observed after 1 khr of 100°C storage. Degradation is a function of the temperature at the p - n junction and the junction current density. As shown in Fig. 19, a larger decline in light output is observed as junction current density and/or temperature at the junction increases. The dependence of degradation on current density J is superlinear, varying as J^x with $1.5 < x < 2.5$. Hence, accelerated-aging tests typically use high currents and temperatures to shorten the time needed to observe LED degradation. The maximum stress level shown in Fig. 19 is 200 percent of the maximum allowable drive current specified in the data sheet.

Light-output degradation in GaAsP LEDs is due to an increase in the nonradiative space-charge recombination current. Total current flowing through the LED is made up of the sum of diffusion current and space-charge recombination current, as shown in the following equations:

$$I_T(V, t) = A(t)e^{qV/kT} + B(t)e^{qV/2kT} \quad (5)$$

where q is electron charge, k is Boltzmann's constant, and T is temperature.³² The first term is diffusion current that produces light output, while the second term is space-charge recombination that is nonradiative. At fixed I_T , if $B(t)$ increases, then the diffusion term must decrease and, hence, the light output decreases. The reason for the increase in space-charge recombination in GaAsP LEDs is not fully understood.

The degradation characteristics of GaAlAs LEDs differ from those of GaAsP LEDs. Typical curves of normalized light output versus time for GaAlAs LEDs are shown in Fig. 20. "Good" units show negligible decrease in light output when operated under normal service conditions. Gradual degradation may occur, but it is relatively uncommon. The predominate failure mode in GaAlAs LEDs is catastrophic degradation. The light emission decreases very rapidly over a period of less than 100 hours

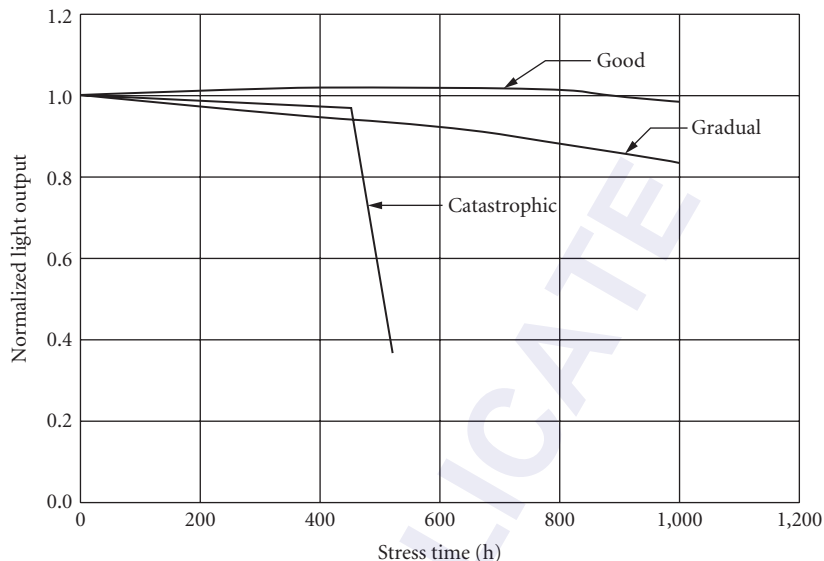


FIGURE 20 Light output degradation of AlGaAs LEDs. Three modes are shown: “good” devices with negligible light-output decrease over time, devices which degrade gradually over time, and “catastrophic” degradation devices where the flux rapidly decreases over a short time period

and simultaneously nonradiative regions (“dark-spot” or “dark-line defects”) are observed to form. The catastrophic degradation mechanism is described in detail in Refs. 33 and 34. In brief, the dark regions are caused by nonradiative recombination at dislocation networks that grow rapidly from a crystal dislocation located in the light-emitting region of the LED. Network formation depends on carrier recombination, both nonradiative, which creates mobile point defects, and radiative recombination, which enhances the movement of the point defects to the growing network.

Formation of dark-line defects is enhanced by mechanical stress either present in the LED chip or occurring during assembly. Properly designed products reduce such stress by minimizing bending caused by the different coefficients of expansion in the LED, and by stress-free die attach, wire bond, and encapsulation of the LED during assembly.

Failures due to dark-spot or dark-line defects can be effectively screened out by operating the LED at high current and temperature. Units with defects typically fail within the first few hundred hours. GaAlAs LEDs with dark-spot or dark-line defects are screened out by means of visual inspection and/or by eliminating units with large decreases in light output. GaAlAs LEDs used for fiber-optic applications have small emitting areas (see under “Fiber Optics” discussion and Fig. 14). Due to the high current densities present in such devices, high temperature and current burn-in is used extensively for these types of LEDs.

Another degradation mode in LEDs is the change in the reverse breakdown characteristics over time. The reverse characteristics become soft and the breakdown voltage may decrease to a very low value. Several mechanisms have been observed in LEDs. One cause is localized avalanche breakdown due to microplasma formation at points where electric fields are high. Microplasmas have been observed in GaAs and GaAsP LEDs.

Damage or contamination of the edges of an LED chip can cause increased surface leakage and reduced reverse breakdown. Incomplete removal of damage during die separation operations and damage induced during handling and assembly are known to cause reverse breakdown changes. Die-attach materials also can unintentionally contaminate the edges of LEDs. Copper, frequently found in LED packaging materials, can diffuse into the exposed surfaces of the LED, causing excess leakage and, in some cases, light-output degradation.³⁵ Chips whose p - n junction extends to the edges (i.e., Fig. 6) are very susceptible to damage and/or contamination.

17.11 LED-BASED PRODUCTS

Indicator Lamps

The simplest LED product is an indicator lamp or its infrared equivalent. The most frequently used lamp is shown in Fig. 21. A LED chip with a typical dimension of $250 \times 250 \mu\text{m}$ is attached with conductive silver-loaded epoxy into a reflective cavity coined into the end of a silver-plated copper or steel lead frame. The top of the LED chip is connected with a thin $25\text{-}\mu\text{m}$ gold wire to the second terminal of the lead frame. The lead frame subassembly is then embedded in epoxy. The epoxy serves several functions: (1) it holds the assembly together and protects the delicate chip and bond wire; (2) it increases the light extraction from the chip (see under "Light Extraction," discussed earlier); and (3) it determines the spatial light distribution.

There are a large number of variations of the lamp shown in Fig. 21. Besides the obvious variation of source wavelength, there are variations of size, shape, radiation pattern, etc. The cross section of the plastic body ranges from 2 to 10 mm. The radiation pattern is affected by three factors: the shape of the dome, the relative position of the chip/reflector combination, and by the presence of a diffusant in the epoxy. Figure 22*a* shows the radiation pattern of a lamp using clear plastic as an encapsulant. The rays emanating from chip and reflector are collimated into a narrow beam. For many indicator lamps a broader viewing angle is desired such as that shown in the radiation pattern of Fig. 22*b*. This effect is achieved by adding a diffusant, such as glass powder, to the clear plastic. Another variation is shape. Common shapes are round, square, rectangular, or triangular.

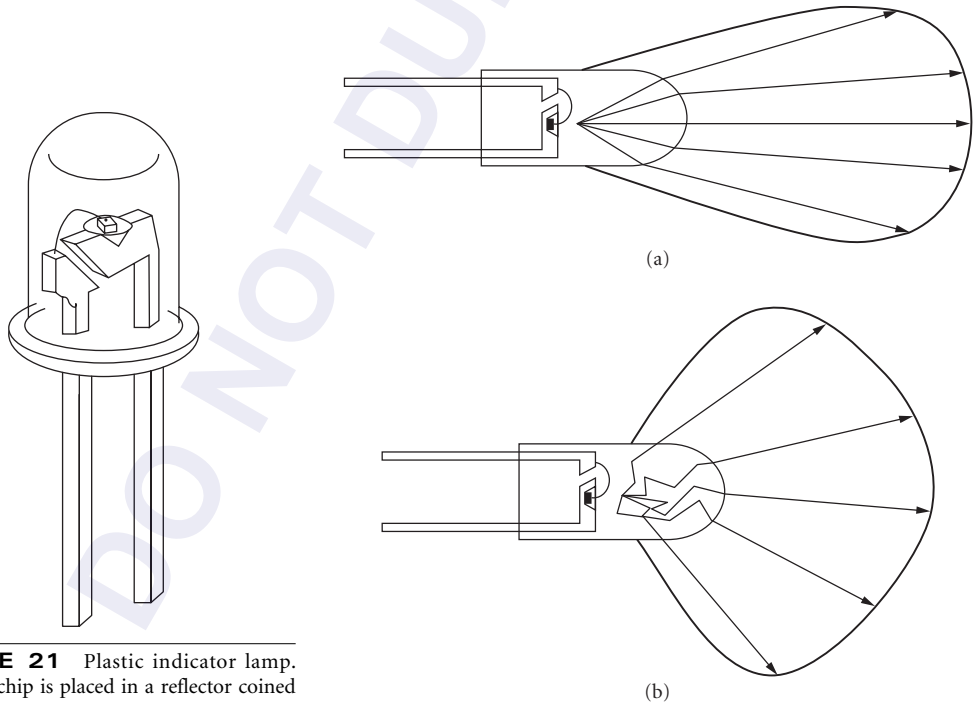


FIGURE 21 Plastic indicator lamp. The LED chip is placed in a reflector coined into the end of one electrode lead. The top of the chip is connected with a gold wire to the second electrode. The electrodes are encapsulated in plastic to form a mechanically robust package.

FIGURE 22 Radiation pattern of two types of LED indicator lamps: (a) lamp with clear plastic package with a narrow beam and (b) lamp with a diffusing plastic package (glass powder added) with a broader radiation pattern.

Another variation is achieved by placing two different chips into the reflector cup. For instance, a lamp with green and red chips connected to the second post in an antiparallel fashion can operate as a red indicator if the reflector post is biased positive, and as a green indicator if it is biased negative. Rapidly switching between the two polarities one can achieve any color between the two basic colors, i.e., yellow or orange, depending on current and duty cycle.

A number of chips can be combined in a single package to illuminate a rectangular area. These so-called annunciator assemblies range in size from 1 to several cm. They typically use 4 chips per cm^2 . By placing an aperture-limiting symbol or telltale in front of the lit area, these structures are cost-effective means to display a fixed message, such as warning lights in an automotive dashboard.

Numeric Displays

Numeric displays are usually made up of a nearly rectangular arrangement of seven elongated segments in a figure-eight pattern. Selectively switching these segments generates all ten digits from 0 to 9. Often decimal point, colon, comma, and other symbols are added.

There are two main types of LED numeric displays: (1) monolithic displays and (2) stretched segment displays. All monolithic displays are based on GaAs-GaAsP technology. Seven elongated p -doped regions and a decimal point are diffused into an n -type epitaxial layer of a single or monolithic chip (see Fig. 8). Electrically, this is a structure with eight anodes and one common cathode. This monolithic approach is relatively expensive. For arms-length viewing, a character height of 3 to 5 mm is required. Adding space for bonding pads, decimal point, and edge separation, such a display consumes around 10 mm^2 of expensive semiconductor material per digit. One way to reduce material and power consumption is optical magnification. Viewing-angle limitation and distortion limit the magnification M to $M \leq 2.0$. Power consumption is reduced by M^2 —an important feature for battery-powered applications.

For digits $>5 \text{ mm}$, a stretched segment display is most cost effective. The design of Fig. 23 utilizes a $250 \times 250 \text{ }\mu\text{m}$ chip to generate a segment with dimensions of up to $8 \times 2 \text{ mm}$ for a 20-mm digit height. This corresponds to a real magnification of $M^2 = 256$, or an equivalent linear magnification of $M = 16$. This magnification is achieved without a reduction in viewing angle by using scattering optics. An LED chip is placed at the bottom of a cavity having the desired rectangular exit shape and

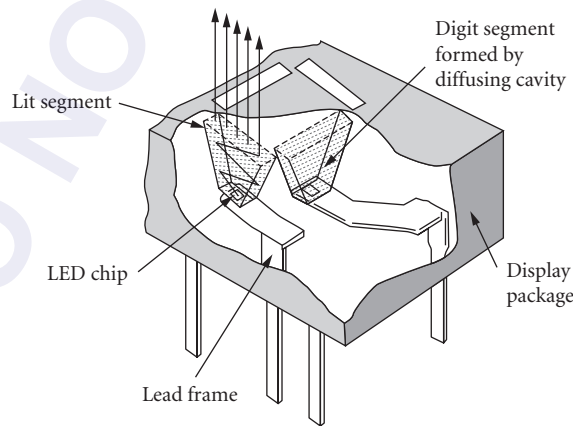


FIGURE 23 Cutaway of a seven-segment numeric LED display, showing how light from a small LED chip is stretched to a large character segment using a diffusing cavity. Characters 0 to 9 are created by turning on appropriate combinations of segments.

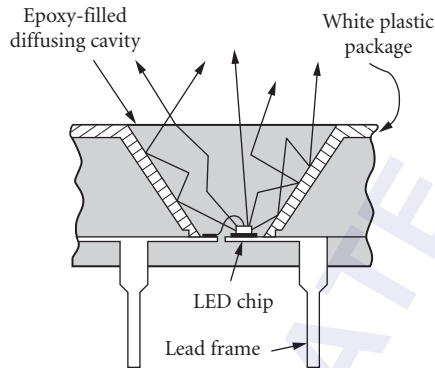


FIGURE 24 Cross section through one segment of the seven-segment numeric display shown in Fig. 23. A LED chip is placed at the bottom of a highly reflective, white plastic cavity which is filled with clear plastic containing a diffusant. Light is scattered within the cavity to produce uniform emission at the segment surface.

the cavity is filled with a diffusing plastic material (see Fig. 24). The side and bottom surfaces of the cavity are made as reflective as possible. Highly reflective white surfaces are typically used. A good white plastic surface measured in air may have a reflectivity of 94 percent compared with 98 percent for Ag and 91 percent for Al. Ag and Al achieve this reflectivity only if evaporated on a specularly smooth surface. Measured in plastic, the reflectivity of the white surface increases from 94 to 98 percent, the Ag surface remains at 98 percent, while the Al surface decreases to 86 percent. Both metallic surfaces have substantially lower reflectivities if they are evaporated onto a nonspecular surface, or if they are deposited by plating. Practically all numeric LED displays above 5-mm character heights are made using white cavity walls and diffusing epoxy within the cavity.

This case of magnification by scattering does not result in a power saving as in the case of magnification of monolithic displays. Since there is practically no reduction in emission angle, the law of energy conservation requires an increased light flux from the chip that equals the area magnification plus reflection losses in the cavity.

Alphanumeric Displays

There are two ways LEDs are used to display alphanumeric information: either by using more than seven elongated segments, i.e., 14, 16, or 24; or by using an array of LED chips in a 5×7 dot matrix format. The multiple segment products are similar in design to the monolithic or stretched segment numeric displays described above.

In the case of small monolithic characters, the number of input terminals quickly exceeds conventional pin spacing. These products are usually clusters of 4 to 16 characters combined with a decoder/driver integrated circuit within the same package. To reduce cost and power, some modest optical magnification is usually used. The segmented displays are usually larger, i.e., 12 to 25 mm and limited to 14 segments per character. At this size, there is no pin density constraint and the decoder is usually placed outside the package.

The most frequently used alphanumeric LED display is based on a 5×7 matrix per character. For small characters in the range of 5 to 8 mm, the LED chips are directly viewed and pin density limitations require an on-board decoder/driver IC. Products are offered as end-stackable clusters of 4, 8, and 16 characters.

For larger displays, the LED chip is magnified by the same optical scattering technique described earlier for numerics. Exit apertures per pixel have a diameter of 2 to 5 mm. Products are offered as 5×7 single characters or end-stackable 20×20 tiles for large message- or graphics-display panels. At this size, pin density is not a limitation.

Optocouplers

An optocoupler is a device where signal input and signal output have no galvanic connection. It is mainly used in applications as the interface between the line voltage side of a system and the low-voltage circuit functions, or in systems where the separate ground connection of interconnected subsystems causes magnetic coupling in the galvanic loop between signal and ground connections. By interrupting the galvanic loop with an optical signal path, many sources of signal interference are eliminated.

The oldest optocouplers consist of an IR LED and a photodetector facing each other in an insulating tube. The second generation utilized the so-called dual-in-line package widely used by logic ICs. In this package an IR emitter and a phototransistor are mounted face to face on two separate lead frames. The center of the package between emitter and detector is filled with a clear insulating material. The subassembly is then molded in opaque plastic to shield external light and to mechanically stabilize the assembly (Fig. 25). The second generation optocouplers have limited speed performance for two reasons: (1) the slow response time of the GaAs:Si LED and (2) the slow response of the photo-transistor detector because of the high collector-base capacitance.

The third generation of optocouplers overcomes the speed limitation. It uses an integrated photodetector and a decoupled gain element. Integration limits the thickness of the effective detection region in silicon to 5 to 7 μm . This thickness range forces a shift of the source to wavelengths shorter than the 940-nm sources used in the second-generation couplers. Third-generation couplers use GaAsP (700 nm) or AlGaAs emitters (880 nm).

Within the last decade, the optocoupler product family has seen further proliferation by adding features on the input or output side of the coupler. One proliferation resulted in couplers behaving

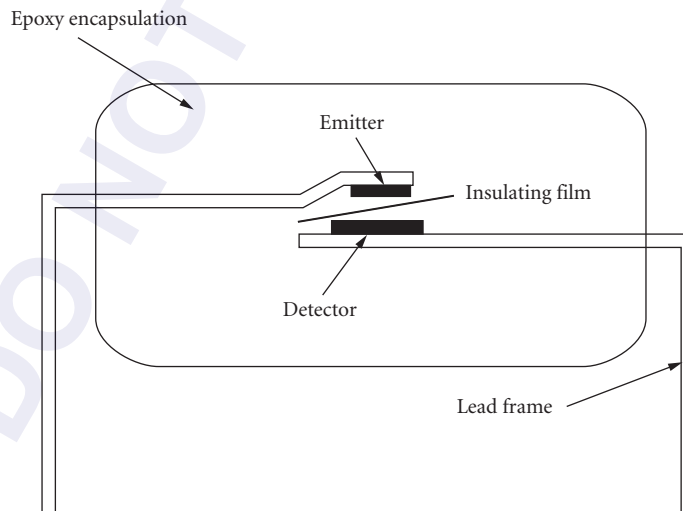


FIGURE 25 Optocoupler consisting of face-to-face emitter and detector chips. An insulating film is placed between the chips to increase the ability of the optocoupler to withstand high voltages between input and output electrodes

like logic gates. Another variation used MOS FET devices on the output side, eliminating the offset voltages of bipolar devices. These couplers are comparable to the performance of conventional relays and are classified as solid-state relays. Other types use a CMOS input driver and CMOS output circuitry to achieve data transfer rates of 50 Mb/s and CMOS interface compatibility.

Fiber Optics

LEDs are the primary light source used in fiber-optic links for speeds up to 200 Mb/s and distances up to 2 km. For higher speed and longer distances, diode lasers are the preferred source.

For fiber-optic applications, LEDs have to meet a number of requirements that go well beyond the requirements for lamps and displays. The major issues are minimum and maximum flux coupled into the fiber, optical rise and fall times, source diameter, and wavelength. Analysis of flux budget, speed, fiber dispersion, wavelength, and maximum distance is quite complex and goes far beyond the scope of this work. A simplified discussion for the popular standard, Fiber Distributed Data Interface (FDDI) will highlight the issues. For a detailed discussion, the reader is referred to Ref. 36.

Flux Budget The minimum flux that has to be coupled into the fiber is determined by receiver sensitivity (-31 dBm), fiber attenuation over the 2-km maximum distance (3 dB), connector and coupler losses (5 dB), and miscellaneous penalties for detector response variations, bandwidth limitations, jitter, etc. (3 dB). With this flux budget of 11 dB, the minimum coupled power has to be -20 dBm. Another flux constraint arises from the fact that the receiver can only handle a maximum level of power before saturation (-14 dBm). These two specifications bracket the power level coupled into the fiber at a minimum of -20 dBm (maximum fiber and connector losses) and at a maximum of -14 dBm (no fiber or connector losses for the case of very short fiber).

Speed The transmitter speed or baud rate directly translates into a maximum rise and fall time of the LED. The 125 Mdb FDDI specifications call for a maximum rise and fall time of 3.5 ns. As a rule of thumb, the sum of rise and fall time should be a little shorter than the inverse baud period (8 ns for 125 Mbd).

Source Diameter and Fiber Alignment Efficient coupling of the LED to the fiber requires a source diameter that is equal to or preferably smaller than the diameter of the fiber core. For sources smaller than the core diameter, a lens between source and fiber can magnify the source to a diameter equal to or larger than the core diameter. The magnification has two benefits: (1) It increases the coupling efficiency between source and fiber. The improvement is limited to the ratio of source area to core cross-sectional area. (2) It increases the apparent spot size to a diameter larger than the fiber core. This effect relaxes the alignment tolerance between source and core and results in substantially reduced assembly and connector costs.

Wavelength The LEDs used in fiber-optic applications operate at three narrowly defined wavelength bands 650, 820 to 870, and 1300 nm, as determined by optical fiber transmission characteristics.

650 nm This band is defined by an absorption window in acrylic plastic fiber. It is a very narrow window between two C-H resonances of the polymer material. The bottom of the window has an absorption of approximately 0.17 dB/m. However, the effective absorption is in the 0.3 to 0.4 dB/m range because the LED linewidth is comparable to the width of the absorption window and the LED wavelength changes with temperature. The 650-nm LEDs use either $\text{GaAs}_{1-x}\text{P}_x$ with $x = 0.4$ or GaAlAs_{1-x} with $x = 0.38$. Quantum efficiencies are at 0.2 and 1.5 percent, respectively. Maximum link length is in the range of 20 to 100 m, depending on source efficiency, detector sensitivity, speed, and temperature range.

820 to 870 nm This window was chosen for several reasons. GaAlAs emitters (see Fig. 14) and Si detectors are readily available at this wavelength. Early fibers had an absorption peak from water

contamination at approximately 870 nm. As fiber technology improved, the absorption peak was eliminated and the wavelength of choice moved from 820- to the 850 to 870-nm range. Fiber attenuation at 850 nm is typically 3 dB/km. Maximum link length in the 500- to 2000-m range, depending on data rate. In the 850 to 870-nm window the maximum link length is limited by chromatic dispersion. GaAlAs emitters have a half-power linewidth of approximately 35 nm. The velocity of light in the fiber is determined by the index of refraction of the fiber core. The index is wavelength-dependent resulting in dispersion of the light pulse. This dispersion grows with distance. The compounded effect of LED linewidth and fiber dispersion is a constant distance-speed product. For a typical multimode fiber and GaAlAs LED combination, this product is in the range of 100 Mbd-km.³⁶

1300 nm At this wavelength, the index of refraction as a function of wavelength reaches a minimum. At this minimum, the velocity of light is practically independent of wavelength, and chromatic dispersion is nearly eliminated. The distance-speed limitation is caused by modal dispersion. Modal dispersion can be envisioned as a different path length for rays of different entrance angles into the fiber. A ray going down the middle of the core will have a shorter path than a ray entering the fiber at the maximum acceptance angle undergoing many bends as it travels down the fiber. The resulting modal dispersion limits multimode fibers to distance bandwidth products of approximately 500 Mbd-km. LED sources used at this wavelength are GaInAsP emitters, as shown in Fig. 15. At this wavelength, fiber attenuation is typically <1.0 dB/km and maximum link length is in the range of 500 to 5000 m, depending on data rate and flux budget.

Sensors

LED/detector combinations are used in a wide range of sensor applications. They can be grouped into three classes: transmissive, reflective, and scattering sensors.

Transmissive Sensors The most widely used transmissive sensor is the slot interrupter. A U-shaped plastic holder aligns an emitter and detector face to face. It is used widely for such applications as sensing the presence or absence of paper in printers, end-of-tape in tape-recorders, erase/overwrite protection on floppy disks, and many other applications where the presence or absence of an opaque obstruction in the light path determines a system response.

A widely used slot interrupter is a two- or three-channel optical encoder. A pattern of opaque and transmissive sections moves in front of a fixed pattern with the same spatial frequency. Two optical channels positioned such that they are 90° out of phase to each other with regard to the pattern allow the measurement of both distance (number of transmissive/opaque sequences) and direction (phase of channel A with regard to channel B). Such encoders are widely used in industrial control applications, paper motion in printers, pen movement in plotters, scales, motor rotation, etc. A third channel is often used to detect an index pulse per revolution to obtain a quasi-absolute reference.

Reflective Sensors In a reflective sensor an LED, detector, and associated optical elements are positioned such that the detector senses a reflection when a reflective surface (specular or diffuse) is positioned within a narrow sensing range. A black surface or nonaligned specular surface or the absence of any reflective surface can be discriminated from a white surface or properly aligned specular surface. Applications include bar-code reading (black or white surface), object-counting on a conveyor belt (presence or no presence of a reflecting surface), and many others. Many transmissive sensor applications can be replaced by using reflective sensors and visa versa. The choice is usually determined by the optical properties of the sensing media or by cost. An emerging application for reflective sensors is blood gas analysis. The concentration of O₂ or CO₂ in blood can be determined by absorption at two different LED wavelengths, i.e., red and infrared.

Scattering Sensors One design of smoke detectors is based on light scattering. The LED light beam and the detector path are crossed. In the absence of smoke, no light from the LED can reach the detector. In the presence of smoke, light is scattered into the detector.

17.12 REFERENCES

1. A. A. Bergh and P. J. Dean, "Light-Emitting Diodes," *Proc. IEEE* vol. 60, 1972, pp. 156–224.
2. K. Gillessen and W. Shairer, *LEDs—An Introduction*, Prentice-Hall, 1987.
3. M. G. Craford, "Properties and Electroluminescence of the GaAsP Ternary System," *Progress in Solid State Chemistry*, vol. 8, 1973, pp. 127–165.
4. A. H. Herzog, W. O. Groves, and M. G. Craford, "Electroluminescence of Diffused GaAs_{1-x}P_x Diodes with Low Donor Concentrations," *J. Appl. Phys.* vol. 40, 1969, pp. 1830–1838.
5. R. A. Faulkner, "Toward a Theory of Isoelectronic Impurities in Semiconductors," *Phys. Rev.* vol. 175, 1968, pp. 991–1009.
6. W. O. Groves, A. J. Herzog, and M. G. Craford, "The Effect of Nitrogen Doping on GaAs_{1-x}P_x Electroluminescent Diodes," *Appl. Phys. Lett.* vol. 19, 1971, pp. 184–186.
7. M. G. Craford, R. W. Shaw, W. O. Groves, and A. H. Herzog, "Radiative Recombination Mechanisms in GaAsP Diodes with and without Nitrogen Doping," *J. Appl. Physics* vol. 43, 1972, pp. 4075–4083.
8. M. G. Craford, D. L. Keune, W. O. Groves, and A. H. Herzog, "The Luminescent Properties of Nitrogen Doped GaAsP Light Emitting Diodes," *J. Electron. Matls.* vol. 2, 1973, pp. 137–158.
9. M. G. Craford and W. O. Groves, "Vapor Phase Epitaxial Materials for LED Applications," *Proc. IEEE* vol. 61, 1973, pp. 862–880.
10. J. C. Campbell, N. Holonyak Jr., M. G. Craford, and D. L. Keune, "Band Structure Enhancement and Optimization of Radiative Recombination in GaAs_{1-x}P_x:N (and In_{1-x}Ga_xP:N)," *J. Appl. Phys.* vol. 45, 1974, pp. 4543–4553.
11. R. A. Logan, H. G. White, and W. Wiegmann, "Efficient Green Electroluminescence in Nitrogen-Doped GaP *p-n* Junctions," *Appl. Phys. Lett.* vol. 13, 1968, p. 139.
12. A. A. Bergh and J. A. Copeland, "Optical Sources for Fiber Transmission Systems," *Proc. IEEE*, vol. 68, 1980, pp. 1240–1247.
13. M. G. Craford, "Recent Developments in Light-Emitting Diode Technology," *IEEE Trans. Electron Devices* vol. 24, 1977, pp. 935–943.
14. H. Nather, V. Nitsche, and W. Schairer, "High Resolution Printing Capability of LED-Based Print Heads," *Proc. SPIE*, 1988, pp. 396–404.
15. M. G. Craford, "Light-Emitting Diode Displays," in *Flat-Panel Displays and CRTs*, L. E. Tannas Jr. (ed.), Van Nostrand Reinhold, 1985, pp. 289–331.
16. L. W. Cook, M. D. Camras, S. L. Rudaz, and F. M. Steranka, "High Efficiency 650 nm Aluminum Gallium Arsenide Light Emitting Diodes," *Proc. 14th International Symposium on GaAs and Related Compounds*, Institute of Physics, Bristol, 1988, pp. 777–780.
17. J. M. Dallesasse, D. W. Nam, D. G. Deppe, N. Holonyak Jr., R. M. Fletcher, C. P. Kuo, T. D. Osentowski, and M. G. Craford, "Short-Wavelength (<6400 Å) Room Temperature Continuous Operation of *p-n* InAlGaP Quantum Well Lasers," *Appl. Phys. Lett.* vol. 19, 1988, pp. 1826–1828.
18. C. P. Kuo, R. M. Fletcher, T. D. Osentowski, M. C. Lardizabel, and M. G. Craford, "High Performance AlGaInP Visible Light-Emitting Diodes," *Appl. Phys. Lett.* vol. 57, 1990, pp. 2937–2939.
19. H. Sugawara, M. Ishikawa, and G. Hatakoshi, "High-efficiency InGaAlP/GaAs Visible Light-Emitting Diodes," *Appl. Phys. Lett.* vol. 58, 1991, 1010–1012.
20. R. M. Fletcher, C. P. Kuo, T. D. Osentowski, K. H. Huang, and M. G. Craford, "The Growth and Properties of High Performance AlGaInP Emitters Using a Lattice Mismatched GaP Window Layer," *Jour. Elec. Matls.* vol. 20, 1991, pp. 1125–1130.
21. K. H. Huang, J. G. Yu, C. P. Kuo, R. M. Fletcher, T. D. Osentowski, L. J. Stinson, A. S. H. Liao, and M. G. Craford, "Twofold Efficiency Improvement in High Performance AlGaInP Light-Emitting Diodes in the 555–620 nm Spectral Region Using a Thick GaP Window Layer," *Appl. Phys. Lett.* vol. 61, 1992, pp. 1045–1047.
22. M. A. Haase, J. Qiv, J. M. DePoydt, and H. Cheng, "Blue-green Laser Diodes," *Appl. Phys. Lett.* vol. 58, 1991, pp. 1272–1275.
23. J. Jeon, J. Ding, A. V. Normikko, W. Xie, M. Kobayashi, and R. L. Gunshore, "ZnSe Based Multilayer *p/n* Junctions as Efficient Light Emitting Diodes for Display Applications," *Appl. Phys. Lett.* vol. 60, 1992, pp. 892–894.

24. S. Nakamura, M. Senoh, and T. Mukai, "Highly P-typed Mg-doped GaN Films Grown with GaN Buffer Layers," *Jpn. Jour. Appl. Phys.* vol. 30, 1991, pp. L1701–L1711.
25. A. G. Fischer, "Methods of Growing Crystals under Pressure," in *Crystal Growth*, B. R. Pamplin (ed.), Pergamon Press, Oxford, 1975, pp. 521–555.
26. R. L. Moon, "Liquid Phase Epitaxy," in *Crystal Growth*, 2d ed., B. R. Pamplin (ed.), Pergamon Press, Oxford, 1980.
27. J. W. Burd, "A Multi-Wafer Growth System for the Epitaxial Deposition of GaAs and GaAs_{1-x}P_x," *Trans. Met. Soc. AIME*, 1969, pp. 571–576.
28. G. B. Stringfellow, *Organometallic Vapor Phase Epitaxial Growth of III-V Semiconductor: Theory and Practice*, Academic Press, Oxford, 1989.
29. E. C. H. Parker (ed.), *Technology and Physics of Molecular Beam Epitaxy*, Plenum Press, New York, 1985.
30. S. K. Ghandhi, *VLSI Fabrication Principles*, John Wiley and Sons, New York, 1982.
31. Hewlett-Packard, *Optoelectronics/Fiber-Optics Application Manual*, 2d ed., McGraw-Hill, New York, 1981, p. 82.
32. A. S. Grove, *Physics and Technology of Semiconductor Devices*, sec. 6.6, John Wiley, New York, 1967.
33. O. Veda, *Material Research Society Symposium Proc.* vol. 184, 1991, p. 125.
34. M. Fukuda, *Reliability and Degradation of Semiconductor Lasers and LEDs*, Artech House, 1991.
35. A. A. Bergh, "Bulk Degradation of GaP Red LEDs," *IEEE Trans. Electron Devices* vol. 18, 1971, pp. 166–170.
36. D. C. Hanson, "Progress in Fiber Optic LAN and MAN Standards," *IEEE LCS Magazine* vol. 1, 1990, pp. 17–25.

HIGH-BRIGHTNESS VISIBLE LEDs

Winston V. Schoenfeld

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

18.1 INTRODUCTION

Over the past decade a transformation in light emitting diode (LED) application has occurred from indicator use to the more demanding solid-state lighting and illumination markets. This includes areas such as backlighting in consumer hand-held products, outdoor displays, traffic signals, and general room lighting through the replacement of incandescent and fluorescent light sources. This transformation has been driven directly by the significant increases in LED efficiency that has enabled LEDs to penetrate into such markets. These new high-efficiency LEDs are typically referred to as high-brightness LEDs (HB-LEDs), and utilize quantum well active regions to achieve their high efficiency and lumen output.

Quantum well-active regions can have either a single quantum well (SQW) or multiple quantum well (MQW) structure, reducing internal reabsorption and increasing the radiative recombination rate of the device through greater spatial overlap of electrons and holes. A discussion of quantum wells can be found in Chap. 19, “Semiconductor Lasers.” Figure 1 provides a basic structure for a modern HB-LED. It is similar to the DH structure, but utilizes a MQW active region between the n and p -type layers of the device, composed of quantum wells (QWs) that are on the order of a few nanometers in thickness. By adjusting either the composition or width of the QW, the emission wavelength of the LED can be tuned. The use of QW active regions in LEDs has not only allowed for an increase in the efficiency of the devices, it has also resulted in the ability to obtain new wavelength emissions that were not previously available due to epitaxial strain and lattice matching constraints.

18.2 THE MATERIALS SYSTEMS

There are two main semiconductor material systems currently exploited for visible HB-LEDs. These are the AlInGaP and AlInGaN systems. A detailed discussion of the AlInGaP system is in Chap. 17, “Light-Emitting Diodes.” AlInGaP is used for amber and red LEDs that fall within the 590- to 650-nm wavelength region. It is limited to this wavelength range because the AlInGaP quaternary shifts to an indirect band gap as wavelengths below 590 nm are targeted. This left a void in the blue and green spectral regions for visible HB-LEDs.

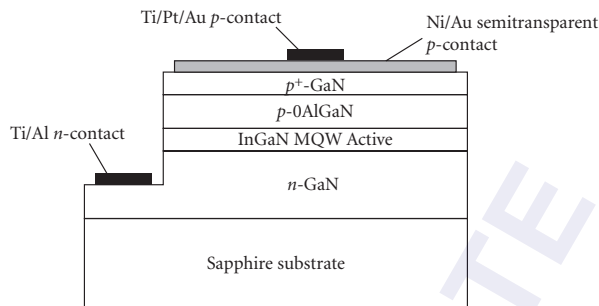


FIGURE 1 Typical structure for a modern InGaN HB-LED showing associated epitaxial layers, including the multiquantum well (MQW) active region.

In the past decade, this void has been quickly filled with MQW LEDs made from the AlInGaN material system. GaN has a wurtzite crystal structure and band gap at room temperature of about 3.2 eV (385 nm). By alloying with Al or In, the band gap can be shifted to shorter or longer wavelengths, respectively. This has led to the realization of InGaN QW active regions in HB-LEDs offering light output in the blue to green spectral regions. Although InN has a direct energy gap of 0.7 eV, phase segregation in InGaN alloys limits the ability to obtain active regions emitting at wavelengths above 530 nm. Despite this, InGaN QW active regions have successfully been used to create blue and green HB-LEDs with complete wavelength tunability across the visible 400- to 530-nm range.

18.3 SUBSTRATES AND EPITAXIAL GROWTH

Current AlInGaN HB-LEDs use one of two substrates: sapphire or silicon carbide (SiC). Despite considerable lattice mismatch between these and GaN (roughly 16 percent for sapphire and 3.5 percent for SiC), low-temperature buffer layer technologies have been developed to allow for nucleation and growth of high-quality AlInGaN HB-LEDs. Low-temperature AlN buffer layers on sapphire substrates are directly formed by MOCVD that result in dislocation entanglement just above the nucleation interface. As depicted in Fig. 2, many of the dislocations interact and annihilate. The subsequent growth of a thick (typically 3 to 4 μm) *n*-GaN buffer layer further reduces dislocation density below 10^9 per cm^2 . The use of 6H-SiC substrates offers closer lattice matching and lower defect densities, although this comes with a higher cost that has kept sapphire as a more popular substrate solution. Recent research has aimed to produce native or lattice-matched substrate solutions for GaN growth. Native GaN substrates have been produced, although the diameter of such substrates and cost remain a challenge. A sister wide band gap compound, ZnO, has strikingly similar properties

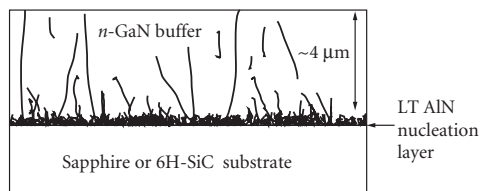


FIGURE 2 Cross-section schematic of the dislocation reduction that occurs due to use of a low temperature (LT) AlN nucleation layer. Dislocations originate from the large lattice mismatch of the nitrides with SiC or sapphire substrates, but can be reduced well below $10^9/\text{cm}^2$ when a LT-AlN nucleation layer and thick GaN buffer layer are used.

to GaN and is currently considered a good candidate for a lattice matched substrate for AlInGaN HB-LED growth. ZnO has the same wurtzite crystal structure as GaN, is nearly lattice matched to GaN, and can be easily doped n -type to form a conductive substrate. Lattice matched substrates made from ZnO have just become commercially available through hydrothermal growth.

Growth of AlInGaN HB-LEDs is accomplished commercially with MOCVD (a short discussion of MOCVD was provided in Sec. 17.8 in Chap. 17). While high quality growth of AlInGaN epitaxial films has been demonstrated by MBE, typical HB-LED structures require 4- to 5- μm -thick films that are not economically realizable in MBE systems due to the slow growth rate and limited number of substrates per growth. MOCVD reactors are able to accommodate multiple 2" substrates per growth (as many as 40 in larger systems) with the necessary uniformity and control of film thickness. One of the initial challenges in AlInGaN HB-LED growth was p -type doping of GaN. The metal-organic precursor sources of MOCVD introduce a considerable amount of hydrogen into the films that directly compensates p -type acceptors. This issue was resolved through post-growth annealing at temperatures above 800°C in which the excess hydrogen is driven out of the epitaxial films. Upon out-diffusion of the hydrogen, the p -type acceptors become active and the necessary hole injection into the structures is then possible.

18.4 PROCESSING

Many of the steps covered in Sec. 17.9 of Chap. 17 are used in the fabrication of HB-LED epitaxial wafers. Sapphire substrates are electrically insulating, imposing the need for creating both n - and p -type contacts on the front of the HB-LED surface as was shown in Fig. 1. As a result, the epitaxial structure must be etched down to the n -GaN layer in order to make electrical contact to the n -side of the HB-LED structure. AlInGaN is relatively resistant to wet chemical etching and must be dry etched using reactive ion etching (RIE) or inductively coupled plasma (ICP) methods. RIE/ICP is capable of high etch rates that are anisotropic, meaning that they are capable of etching vertically down into the structure with little to no lateral etching. SiC substrates are electrically conductive and thus allow for the n -contact to be formed on the back side of the substrate without the need for etching of the front surface.

The typical fabrication process for AlInGaN HB-LEDs utilizes four lithographic steps. These include the mesa etch, SiO₂ passivation, n -contact, and p -contact layers. A top view and associated cross section of a standard HB-LED device are given in Fig. 3, indicating these layers. Prior to the

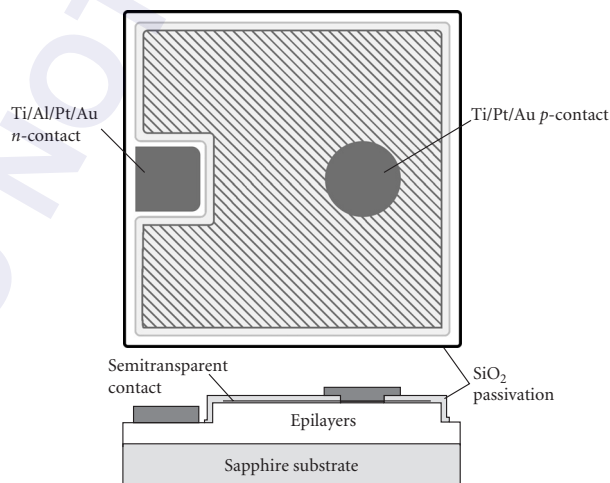


FIGURE 3 Top view and cross-section schematic of a standard nitride HB-LED. A mesa is etched to allow for contact to the n -type underlayer. A semitransparent contact on top of the mesa is used to promote uniform hole injection due to poor hole mobility in the p -GaN.

first lithographic step a semitransparent top contact is deposited on the surface (p -side) of the entire epi wafer. This is necessary in order to achieve uniform current injection across the device due to the higher resistivity of the p -GaN in comparison to the n -GaN. The most common semitransparent contacts are Ni/Au and indium tin oxide (ITO), with Ni/Au being the most widely used. Once the semitransparent contact is formed the lithographic steps are then carried out. The mesa etch is a dry etch step to access the n -GaN layers for the n -type contact as required when using sapphire substrates. The SiO_2 passivation step covers the side walls of the mesa and protects its lateral edges. Once the passivation is in place, the last two lithographic steps define the n - and p -type contacts to the device through a lift-off process. Common metallizations are Ti/Al/Pt/Au for the n -contact and Ti/Pt/Au for the p -contact to the semitransparent contact underlayer. To separate the individual LEDs on the fabricated wafer, the substrate must first be thinned from its typical 400 μm thickness to on the order of 80 to 100 μm using a multiple step wafer polishing method on automated multiwafer polishers. Once thinned, singulation of the individual die is accomplished either by a scribe and break method or by using laser separation. The sawing method typically employed for other III-V LEDs is not possible due to the substantial hardness of sapphire and SiC that greatly limits their ability to be cut using a dicing saw.

18.5 SOLID-STATE LIGHTING

The current push towards energy conservation has created a considerable interest in the replacement of conventional lighting by solid-state LED fixtures. LEDs offer the potential to considerably reduce power consumption while maintaining the necessary lumen output for lighting. Generating white light from HB-LEDs is typically accomplished by one of two methods as depicted in Fig. 4. The first method, shown in Figs. 4a and b, utilizes an AlInGaN-LED in conjunction with one or more phosphors. When a UV-LED is used, the UV light emission is absorbed and re-emitted by a mixture of red, green, and blue phosphors. As indicated in Fig. 4a the phosphors down-convert the UV light (dashed line) to visible light (solid line), and when the appropriate ratio of phosphors is used, white light is emitted. A more common LED/phosphor combination used for general illumination is the combination of a blue (~ 465 nm) AlInGaN-LED and a yttrium aluminum garnet (YAG) phosphor such as cerium doped $\text{Y}_3\text{Al}_5\text{O}_{12}$. In this approach (Fig. 4b) a portion of the blue LED emission is absorbed by the YAG phosphor and down-converted to the yellow spectral region. When the proper ratio of YAG phosphor is used, the resulting binary complimentary output of the blue LED and yellow YAG phosphor creates white light as perceived by observers. The second method for white light emission is preferred when color tuning is necessary, such as in outdoor displays. As indicated in Fig. 4c, by using red, green, and blue LEDs, one can create white light when the appropriate ratio of each is selected. This RGB approach has the added benefit of allowing the user to create any color within the associated color gamut by balancing the ratio of the individual LEDs.

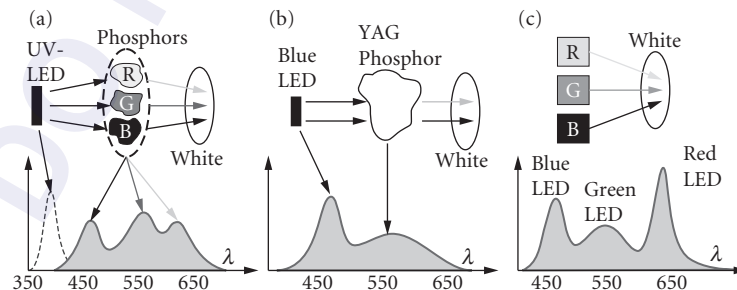


FIGURE 4 Methods for white light generation using HB-LEDs. White light is achieved by using a HB-LED and phosphors (a and b), or the combination of a red, green, and blue HB-LED (c).

This method also allows for active adjustment of the color temperature of white light emission which is not possible using the LED/phosphor approach.

The primary figure of merit for solid-state lighting is the luminous efficacy of the fixture. Luminous efficacy refers to the ratio of lumen output from the source to the power input, and has the units of lumens per Watt (lm/W). Conventional incandescent and fluorescent lights have typical luminous efficacies of 15 lm/W and 70 lm/W, respectively, varying somewhat depending on manufacturer. As of early 2008, currently available LED-based replacements have luminous efficacies as high as 80 lm/W, exceeding compact fluorescents. By comparison, lab demonstrations of white LEDs are commonly breaking 120 lm/W, with the expectation of reaching 150 lm/W in the near future. Despite these accomplishments the added cost of solid-state white lighting has slowed the market acceptance. A common cost comparison is the first cost of a light source, typically measured in terms of the cost per kilolumen (klm). Incandescent and fluorescent light sources fall in the \$0.5 to \$3.00 per klm, while current light emitting diode sources fall easily in the \$20 to \$30 per klm range. Considerable progress in efficiency continues to climb and it is expected that as the luminous efficacy and lifetime of solid-state white LEDs continues to rise that the cost per kilolumen will decrease while market acceptance increases.

18.6 PACKAGING

HB-LEDs chips for standard current use (20 to 70 mA) have lateral geometries on the order of $350 \times 350 \mu\text{m}^2$. Such LED chips are commonly packaged similarly to IR-LEDs in a 5-mm T1-3/4 format, as shown in Figs. 21 and 22 in Chap. 17. For AlInGaP HB-LEDs the substrate is conductive and silver filled epoxy is used for attaching the die to the lead frame to form one of the contacts. The other contact is formed using standard ball wire bonding to the top of the LED chip. For AlInGaN HB-LEDs, commonly grown on insulating sapphire substrates, two wire bonds are required to make the electrical connections to the lead frame. Silver-filled epoxy is still used for such LEDs since the sapphire is transparent and emission can be effectively redirected upward from the epoxy surface below the chip. In addition to the 5-mm T1-3/4 package, a considerable number of new surface mount device (SMD) packages have become available to support the introduction of HB-LEDs into the consumer hand-held device market. A schematic of a typical SMD package is provided in Fig. 5. The LED is placed on a metal lead frame that is encased in a plastic outer shell. Once the HB-LED has been die bonded into the package with silver-filled epoxy, it is then wire bonded and encapsulated with transparent epoxy for protection. Such SMD packages are typically several millimeters on a side; however, smaller versions with formats nearly identical to chip resistors are available that have

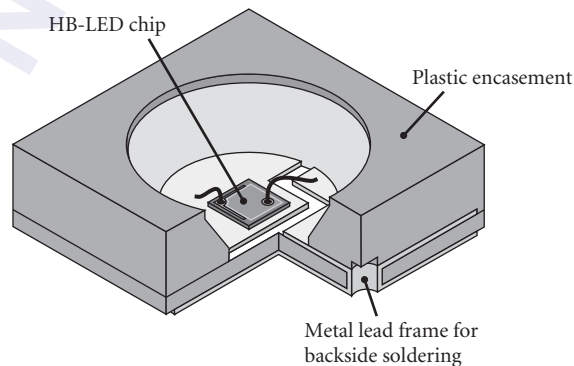


FIGURE 5 Schematic of a standard surface mount device (SMD) package for HB-LEDs.

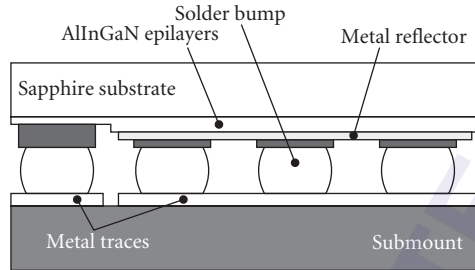


FIGURE 6 Cross-section schematic of the flip chip geometry used for high-power nitride LEDs. Flip chip packaging allows for increased heat dissipation from the LED junction, enabling them to be driven at much higher current densities.

very small form factors providing the low profile necessary for applications such as cell phone key pad backlighting.

In recent years the push for solid-state white lighting has resulted in a transition to larger HB-LED chip sizes and an increase in the operating current densities. The latter brings new constraints to packaging due to the increased need for thermal management in an effort to keep the junction temperature of the HB-LED as low as possible. While the traditional packaging formats have proven very effective for standard drive current use, they are not suitable for high drive current applications using large area HB-LEDs since they do not provide adequate thermal management. This has forced a significant change in not only the design of packaging for high drive current, large area HB-LEDs, but also the devices themselves. Most HB-LEDs use sapphire substrates that provide very poor heat dissipation from the LED junction due to the low thermal conductivity of sapphire. This has been overcome by using a flip chip approach, where the LED chip is flipped over onto a carrier and attached using solder bump methods similar to the silicon IC industry. A rough schematic of a flip chip geometry is shown in Fig. 6. The HB-LED chip is designed to support multiple solder bump attachments to the submount. A highly reflective metal layer is formed on the *p*-side of the device to redirect light emission through the backside of the device. The solder bumps affix the LED die firmly to the submount and create the electrical connections to the metal traces below, allowing for subsequent wire bonding between the submount and the external leads of the package. The solder bumps also serve as a route for efficient heat transfer from the LED junction to the submount enabling the junction temperature of the LED chip to remain low under high drive current conditions. There are a variety of external package geometries that have been developed to support flip chip HB-LEDs with no specific convention between manufacturers. Among the main considerations in the package design has been the encapsulant around the LED and providing a low thermal resistance path to the external housing. Standard clear epoxies that are used in conventional packages typically cannot withstand temperatures much above 100°C. In flip chip high drive packages this can be easily exceeded and would cause thermal damage to conventional epoxies. This has led many companies to develop a variety of new high-index encapsulants, such as silicones, that are able to withstand the higher temperature demands of high drive solid-state lighting. Low thermal resistance packaging has also followed many new routes such as packages with copper tungsten metal bases, metal core board, and high thermal conducting insulating materials (e.g., BeO and AlN).

Pamela L. Derry, Luis Figueroa, and Chi-Shain Hong

*Boeing Defense & Space Group
Seattle, Washington*

19.1 GLOSSARY

A	Constant approximating the slope of gain versus current or carrier density
C	Capacitance
c	Speed of light
D	Density of states for a transition
D_c	Density of states for the conduction band
D_v	Density of states for the valence band
d	Active layer thickness
d_{eff}	Effective beam width in the transverse direction
d_G	Guide layer thickness
dg/dN	Differential gain
E	Energy of a transition
E_c	Total energy of an electron in the conduction band
E_g	Bandgap energy
E_n	The n th quantized energy level in a quantum well
E_n^c	The n th quantized energy level in the conduction band
E_n^v	The n th quantized energy level in the valence band
E_v	Total energy of a hole in a valence band
e	Electronic charge
F_c	Quasi-Fermi level in the conduction band
F_v	Quasi-Fermi level in the valence band
f_c	Fermi occupation function for the conduction band
f_d	Damping frequency
f_o	Resonant frequency of an LRC circuit
f_p	Peak frequency

f_r	Resonance frequency
f_v	Fermi occupation function for the valence band
g	Model gain per unit length
g_{th}	Threshold modal gain per unit length
H	Heavyside function
h	Refers to heavy holes
\hbar	Plank's constant divided by 2π
I	Current
I_{off}	DC bias current before a modulation pulse
I_{on}	Bias current during a modulation pulse
I_{th}	Threshold current
J	Current density
J_o	Transparency current density
J_{th}	Threshold current density
K	Constant dependent on the distribution of spectral output function
\mathbf{k}	Wavevector
k	Boltzmann constant
L	Inductance
L	Laser cavity length
L_c	Coherence length
L_z	Quantum well thickness
l	Refers to light holes
$ M ^2$	Matrix element for a transition
m	Effective mass of a particle
m_c	Conduction band mass
m_r	Effective mass of a transition
m_v	Valence band mass
N	Carrier density
N_o	Transparency carrier density
n_{eff}	Effective index of refraction
n_r	Index of refraction
n_{sp}	Spontaneous emission factor
P	Photon density
P_{off}	Photon density before a modulation pulse
P_{on}	Photon density during a modulation pulse
R	Resistance
R_F	Front facet reflectivity
R_R	Rear facet reflectivity
T	Temperature
w	Laser stripe width
α	Absorption coefficient
α	Linewidth enhancement factor
α_i	Internal loss per unit length
β	Spontaneous emission factor
Γ	Optical confinement factor

$\Delta f_{1/2}$	Frequency spectral linewidth
$\Delta\lambda_L$	Longitudinal mode spacing
$\Delta\lambda_{1/2}$	Half-width of the spectral emission in terms of wavelength
λ	Wavelength
λ_o	Wavelength of the stimulated emission peak
τ_d	Turn-on time delay
τ_p	Photon lifetime
τ_s	Carrier lifetime

19.2 INTRODUCTION

This chapter is devoted to the performance characteristics of semiconductor lasers. In addition, some discussion is provided on fabrication and applications. In the first section we describe some of the applications being considered for semiconductor lasers. The following several sections describe the basic physics, fabrication, and operation of a variety of semiconductor laser types, including quantum well and strained layer lasers. Then we describe the operation of high-power laser diodes, including single element and arrays. A number of tables are presented which summarize the characteristics of a variety of lasers. Next we discuss the high-speed operation and provide the latest results, after which we summarize the important characteristics dealing with the spectral properties of semiconductor lasers. Finally, we discuss the properties of surface emitting lasers and summarize the latest results in this rapidly evolving field.

More than 260 references are provided for the interested reader who requires more information. In this *Handbook*, Chap. 17 (LEDs) also contains related information. For further in-depth reviews of semiconductor lasers we refer the reader to the several excellent books which have been written on the subject.¹⁻⁵

19.3 APPLICATIONS FOR SEMICONDUCTOR LASERS

The best-known application of diode lasers is in optical communication systems. However, there are many other potential applications. In particular, semiconductor lasers are being considered for high-speed optical recording,⁶ high-speed printing,⁷ single- and multimode database distribution systems,⁸ long-distance transmission,⁹ submarine cable transmission,¹⁰ free-space communications,¹¹ local area networks,¹² Doppler optical radar,¹³ optical signal processing,¹⁴ high-speed optical microwave sources,¹⁵ pump sources for other solid-state lasers,¹⁶ fiber amplifiers,¹⁷ and medical applications.¹⁸

For very high-speed optical recording systems (>100 MB/s), laser diodes operating at relatively short wavelengths ($\lambda < 0.75 \mu\text{m}$) are required. In the past few years, much progress has been made in developing short-wavelength semiconductor lasers, although the output powers are not yet as high as those of more standard semiconductor lasers.

One of the major applications for lasers with higher power and wide temperature of operation is in local area networks. Such networks will be widely used in high-speed computer networks, avionic systems, satellite networks, and high-definition TV. These systems have a large number of couplers, switches, and other lossy interfaces that determine the total system loss. In order to maximize the number of terminals, a higher-power laser diode will be required.

Wide temperature operation and high reliability are required for aerospace applications in flight control and avionics. One such application involves the use of fiber optics to directly link the flight control computer to the flight control surfaces, and is referred to as fly-by-light (FBL). A second application involves the use of a fiber-optic data network for distributing sensor and video information.

Finally, with the advent of efficient high-power laser diodes, it has become practical to replace flash lamps for the pumping of solid-state lasers such as Nd:YAG. Such an approach has the advantages

of compactness and high efficiency. In addition, the use of strained quantum well lasers operating at $0.98\ \mu\text{m}$ has opened significant applications for high-gain fiber amplifiers for communications operating in the $1.55\text{-}\mu\text{m}$ wavelength region.

19.4 BASIC OPERATION

Lasing in a semiconductor laser, as in all lasers, is made possible by the existence of a gain mechanism plus a resonant cavity. In a semiconductor laser the gain mechanism is provided by light generation from the recombination of holes and electrons (see Fig. 1). The wavelength of the light is determined by the energy bandgap of the lasing material. The recombining holes and electrons are injected, respectively, from the p and n sides of a p - n junction. The recombing carriers can be generated by optical pumping or, more commonly, by electrical pumping, i.e., forward-biasing the p - n junction. In order for the light generation to be efficient enough to result in lasing, the active region of a semiconductor laser, where the carrier recombination occurs, must be a direct bandgap semiconductor. The surrounding carrier injection layers, which are called *cladding layers*, can be indirect bandgap semiconductors. For a discussion of semiconductor band levels see any solid-state physics textbook such as that by Kittel.¹⁹ For a more detailed discussion of carrier recombination see Chap. 17 in this *Handbook*.

For a practical laser, the cladding layers have a wider bandgap and a lower index of refraction than the active region. This type of semiconductor laser is called a *double heterostructure* (DH) laser, since both cladding layers are made of a different material than the active region (see Fig. 2). The first semiconductor lasers were homojunction lasers,^{21–24} which did not operate at room temperature; it is much easier to achieve lasing in semiconductors at low temperatures. Today, however, all semiconductor lasers contain heterojunctions. The narrower bandgap of the active region confines carrier recombination to a narrow optical gain region. The sandwich of the larger refractive index active region surrounded by cladding layers forms a waveguide, which concentrates the optical modes generated by lasing in the active region. For efficient carrier recombination the active region must be fairly thin, typically on the order of $1000\ \text{\AA}$, so a significant fraction of the optical mode spreads into the cladding layers. In order to completely confine the optical mode in the semiconductor structure, the cladding layers must be fairly thick, usually about $1\ \mu\text{m}$.

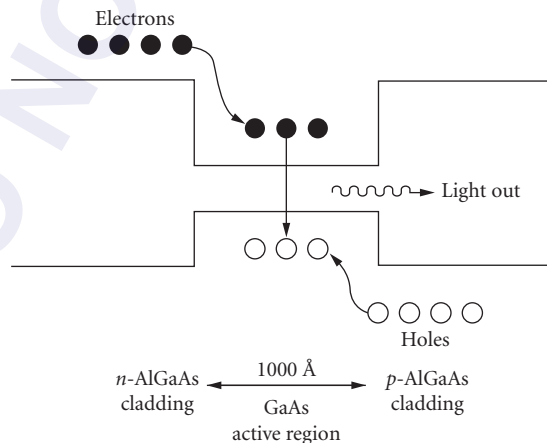


FIGURE 1 Schematic diagram of the recombination of electrons and holes in a semiconductor laser.

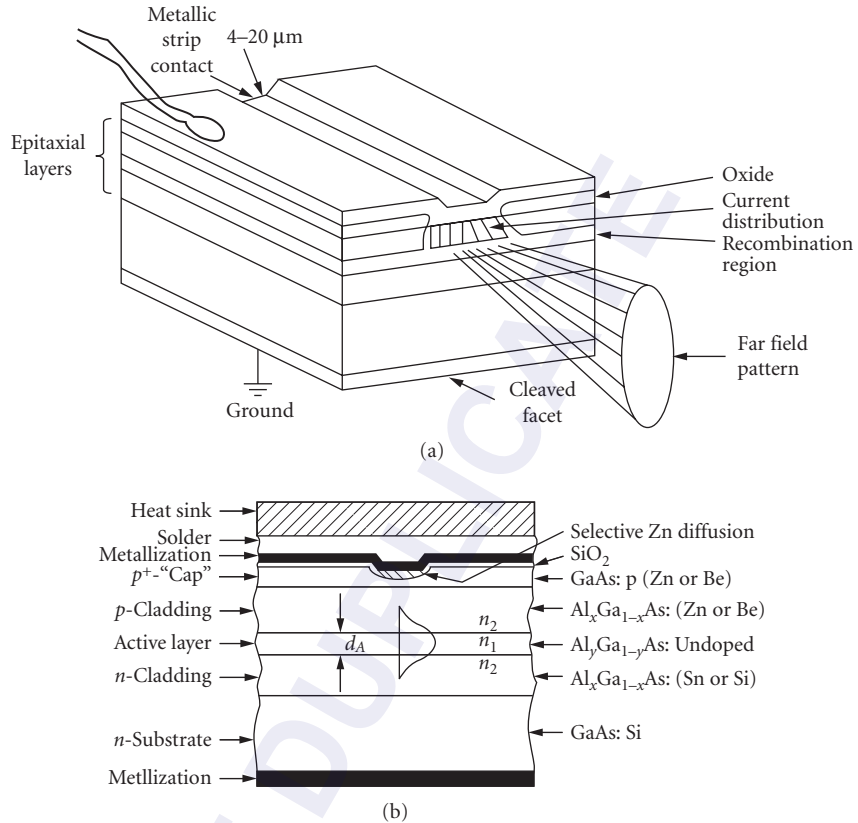


FIGURE 2 (a) Schematic diagram of a simple double heterostructure laser and (b) cross-sectional view showing the various epitaxial layers. (After Ref. 20.)

The resonant cavity of a simple semiconductor laser is formed by cleaving the ends of the structure. Lasers are fabricated with their lasing cavity oriented perpendicular to a natural cleavage plane. For typical semiconductor materials this results in mirror reflectivities of about 30 percent. If necessary, the reflectivities of the end facets can be modified by applying dielectric coatings to them.²⁵ For applications where it is not possible to cleave the laser facets, it is also possible to etch them,²⁶⁻²⁸ although this is much more difficult and usually does not work as well. Laser cavity lengths can be anywhere from about 50 to 2000 μm, although commercially available lasers are typically 200 to 1000 μm long.

Unpumped semiconductor material absorbs light of energy greater than or equal to its bandgap. When the semiconductor material is pumped optically or electrically, it reaches a point at which it stops being absorbing. This point is called *transparency*. If it is pumped beyond this point, it will have optical gain, which is the opposite of absorption. A semiconductor laser is subject to both internal and external losses. For lasing to begin, i.e., to reach threshold, the gain must be equal to these optical losses. The threshold gain per unit length is given by

$$g_{\text{th}} = \alpha_i + \frac{1}{2L} \ln \left(\frac{1}{R_F R_R} \right) \quad (1)$$

where α_i is the internal loss per unit length, L is the laser cavity length, and R_F and R_R are the front and rear facet reflectivities. (For semiconductor lasers, gain is normally quoted as gain per unit

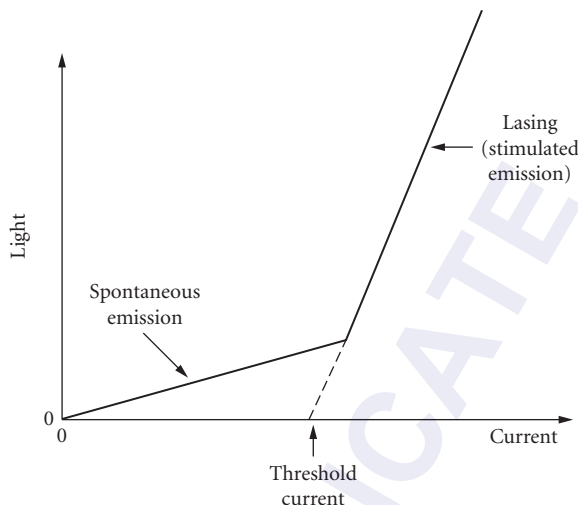


FIGURE 3 An example of the light-versus-current relationship of a semiconductor laser illustrating the definition of threshold current.

length in cm^{-1} . This turns out to be very convenient, but unfortunately is confusing for people in other fields, who are used to gain being unitless.)

The internal loss is a material parameter determined by the quality of the semiconductor layers. Mechanisms such as free-carrier absorption and scattering contribute to α_i .¹ The second term in Eq. (1) is the end loss. A long laser cavity will have reduced end loss, since the laser light reaches the cavity ends less frequently. Similarly, high facet reflectivities also decrease the end loss, since less light leaves the laser through them.

For biases below threshold, a semiconductor laser emits a small amount of incoherent light spontaneously (see Fig. 3). This is the same type of light emitted by an LED (see Chap. 17). Above threshold, stimulated emission results in lasing. The relationship between lasing emission and the bias current of a healthy semiconductor laser is linear. To find the threshold current of a laser this line is extrapolated to the point at which the stimulated emission is zero (see Fig. 3).

For further discussion of optical gain in semiconductors, see under “Quantum Well Lasers” later in this chapter. For more detail, see one of the books that has been written about semiconductor lasers.¹⁻⁵

19.5 FABRICATION AND CONFIGURATIONS

In order to fabricate a heterostructure laser, thin semiconductor layers of varying composition must be grown on a semiconductor substrate (normally GaAs or InP). There are three primary epitaxial methods for growing these layers: liquid phase epitaxy (LPE), molecular beam epitaxy (MBE), and metalorganic chemical vapor deposition (MOCVD), which is also called organometallic vapor phase deposition (OMVPE).

Most of the laser diode structures which have been developed were first grown by LPE.²⁹ For a description of LPE see Chap. 17. Most commercially available lasers are grown by LPE; however, it is not well suited for growing thin structures such as quantum well lasers, because of lack of control and uniformity, especially over large substrates.³⁰ MBE and MOCVD are better suited for growths of thin, uniform structures.

MOCVD^{31,32} is basically a specialized form of chemical vapor deposition. In MOCVD, gases reacting over the surface of a substrate form epitaxial layers; some of the gases are metalorganics. MOCVD is well suited for production environments, since epitaxial layers can be grown simultaneously on multiple large-area substrates and quickly, compared to MBE. It is expected that more commercial laser diodes will be grown by MOCVD in the future.

In the simplest terms, MBE^{30,33–35} is a form of vacuum evaporation. In MBE growth occurs through the thermal reaction of thermal beams of atoms and molecules with the substrate, which is held at an appropriate temperature in an ultrahigh vacuum. MBE is different from simple vacuum evaporation for several reasons: the growth is single crystalline; the growth is much more controlled; and the vacuum system, evaporation materials, and substrate are cleaner.

With MOCVD the sources are gases, while with MBE they are solids. There are advantages and disadvantages to both types of sources. With gaseous sources the operator must work with arsine and/or phosphine, which are extremely hazardous gases. Solid-source phosphorus, however, is very flammable. Also, with MBE balancing the ratios of arsenic and phosphorus is extremely difficult; therefore, MOCVD is the preferred method for growth of GaInAsP and InP. MBE is a slower growth process (on the order of 1 to 2 $\mu\text{m}/\text{h}$) than MOCVD. MBE therefore has the control necessary to grow very thin structures (10 Å), but MOCVD is more efficient for production. MBE has a cleaner background environment, which tends to make it better suited for growths at which background impurities must be kept at a minimum. Newer growth techniques,^{36,37} which combine some of the advantages of both MBE and MOCVD, are gas source MBE, metalorganic MBE (MOMBE), and chemical beam epitaxy (CBE). In these growth techniques the background environment is that of MBE, but some or all of the sources are gases, which makes them more practical for growth of phosphorus-based materials.

Double heterostructure (DH) semiconductor lasers can be fabricated from a variety of lattice-matched semiconductor materials. The two material systems most frequently used are GaAs/Al_xGa_{1-x}As and In_{1-x}Ga_xAsP_{1-y}/InP. All of these semiconductors are III–V alloys. The GaAs/Al_xGa_{1-x}As material system has the advantage that all compositions of Al_xGa_{1-x}As are closely lattice-matched to GaAs, which is the substrate. For GaAs-based lasers, the active region is usually GaAs or low-Al-concentration Al_xGa_{1-x}As ($x < 0.15$), which results in lasing wavelengths of 0.78 to 0.87 μm . Al_xGa_{1-x}As quantum well lasers with wavelengths as low as 0.68 μm have been fabricated,³⁸ but the performance is reduced.

In the In_{1-x}Ga_xAsP_{1-y}/InP material system, the active region is In_{1-x}Ga_xAsF_{1-y} and the cladding layers and substrate are InP. Not all compositions of In_{1-x}Ga_xAsP_{1-y} are lattice-matched to InP; x and y must be chosen appropriately to achieve both lattice match and the desired lasing wavelength.⁴ The lasing wavelength range of InP-based lasers, 1.1 to 1.65 μm , includes the wavelengths at which optical fibers have the lowest loss (1.55 μm) and material dispersion (1.3 μm). This match with fiber characteristics makes In_{1-x}Ga_xAsP_{1-y}/InP lasers the preferred laser for long-distance communication applications. InP-based lasers can also include lattice-matched In_{1-x-y}Al_xGa_yAs layers,^{39–42} but the performance is reduced.

There is a great deal of interest in developing true visible lasers for optical data storage applications. (Al_xGa_{1-x})_{0.5}In_{0.5}P lasers^{43–46} lattice-matched to GaAs have proven superior to very short wavelength GaAs/Al_xGa_{1-x}As lasers. Higher Al concentration layers are cladding layers and a low Al concentration layer or Ga_{0.5}In_{0.5}P is the active region. In this material, system lasers with a lasing wavelength as low as 0.63 μm which operate continuously at room temperature have been fabricated.⁴⁶

In order to fabricate a blue semiconductor laser, other material systems will be required. Recently, lasing at 0.49 μm at a temperature of 77 K was demonstrated in a ZnSe- (II–VI semiconductor) based laser.⁴⁷

Very long-wavelength (>2 μm) semiconductor lasers are of interest for optical communication and molecular spectroscopy. The most promising results so far have been achieved with GaInAsSb/AlGaAsSb lattice-matched to a GaSb substrate. These lasers have been demonstrated to operate continuously at 30°C with a wavelength of 2.2 μm .⁴⁸

Lead salt lasers (Pb_{1-x}Eu_xSeTe_{1-y}, Pb_{1-x}Sn_xSe, PbS_{1-x}Se, PbS_{1-x}Sn_xTe, Pb_{1-x}Sr_xS) can be fabricated for operation at even longer wavelengths,^{4,49–52} but they have not been demonstrated at room temperature. Progress has been made, however, increasing the operating temperature; currently Pb_{1-x}Eu_xSeTe_{1-y}/PbTe lasers operating continuously at 203 K with a lasing wavelength of 4.2 μm have been

demonstrated.⁵³ Other very long-wavelength lasers are possible; recently, HgCdTe lasers with pulsed operation at 90 K and a lasing wavelength of 3.4 μm have been fabricated.⁵⁴

Laser Stripe Structures

We have discussed the optical and electrical confinement provided by a double heterostructure parallel to the direction of epitaxial growth; practical laser structures also require a confinement structure in the direction parallel to the substrate.

The simplest semiconductor laser stripe structure is called an *oxide stripe laser* (see Fig. 4a). The metallic contact on the *n*-doped side of a semiconductor laser is normally applied with no definition for current confinement; current confinement is introduced on the *p* side of the device. For a wide-stripe laser, a dielectric coating (usually SiO_2 or Si_3N_4) is evaporated on the *p* side of the laser. Contact openings in the dielectric are made through photolithography combined with etching of the dielectric. The *p* metallic contact is then applied across the whole device, but makes electrical contact only at the dielectric openings. A contact stripe works very well for wide stripes, but in narrow stripes current spreading is a very significant drawback, because there is no mechanism to prevent current spreading after the current is injected. Also, since the active region extends outside of the stripe, there is no mechanism to prevent optical leakage in a contact-stripe laser. Lasers like this, which provide electrical confinement, but no optical confinement, are called *gain-guided lasers*.

Another type of gain-guide laser is an ion bombardment stripe (see Fig. 4b). The material outside the stripe is made highly resistive by ion bombardment or implantation, which produces lattice defects.⁵⁵ Implantation causes optical damage,⁵⁶ so implantation should not be heavy enough to reach the active region.

A more complicated stripe structure with electrical and optical confinement is required for an efficient narrow-stripe laser. A number of structures which accomplish the necessary confinement have been developed. These structures are called *index-guided lasers*, since optical confinement is achieved through a change in refractive index.

The buried heterostructure laser (BH) was first developed by Tsukada.⁵⁷ To form a BH stripe, a planar laser structure is first grown. Stripe mesas of the laser structure are formed by photolithography combined with etching. For a GaAs-based BH laser, AlGaAs is then regrown around the lasing stripe. Figure 4c is a schematic diagram of a buried heterostructure. Since the active region is completely surrounded by AlGaAs, a BH has tight optical confinement. If the regrown layers are doped to produce a reverse-biased junction or are semi-insulating, a BH laser can also provide good current confinement. There are many variations on the BH structure. In some cases the active region is grown in the second growth step (see Fig. 4d). The tight optical confinement of BH lasers allows practical fabrication of very narrow stripes, on the order of 1 to 2 μm .

There are many other stripe structures that provide weaker optical confinement than a buried heterostructure. One of the simplest and most widely used of these is the ridge waveguide laser (RWG) (Fig. 4e). After epitaxial growth, most of the *p*-cladding layer is etched away, leaving a mesa where the lasing stripe will be. Only this mesa is contacted, which provides electrical confinement. The change in surrounding refractive index produces an effective change in refractive index in the active region beyond the mesa and provides optical confinement. Other stripe structures are described later in this chapter under "High-Power Semiconductor Lasers."

Another type of laser stripe is one in which confinement is provided by the *p-n* junction. The best-known laser of this type is the transverse junction stripe⁵⁸⁻⁶⁰ (see Fig. 4f). In order to fabricate a TJS laser, both cladding layers are grown as *n*-AlGaAs. Zn diffusion is then used to create a *p-n* junction and contacts are applied on either side of the junction. In this laser the current flows parallel to the substrate rather than perpendicular to it. In a TJS laser the active region is limited to the small region of GaAs in which the Zn diffusion front ends.

The examples of laser stripe structures described here are GaAs/AlGaAs lasers. Long-wavelength laser structures (InP-based) are very similar,⁴ but the active region is InGaAsP and the cladding layers are InP. With an *n*-InP substrate the substrate can be used as the *n*-type cladding, which allows greater flexibility in designing structures such as that illustrated in Fig. 4d. For a more detailed discussion of GaAs-based laser stripe structures see Casey and Panish² or Thompson.³

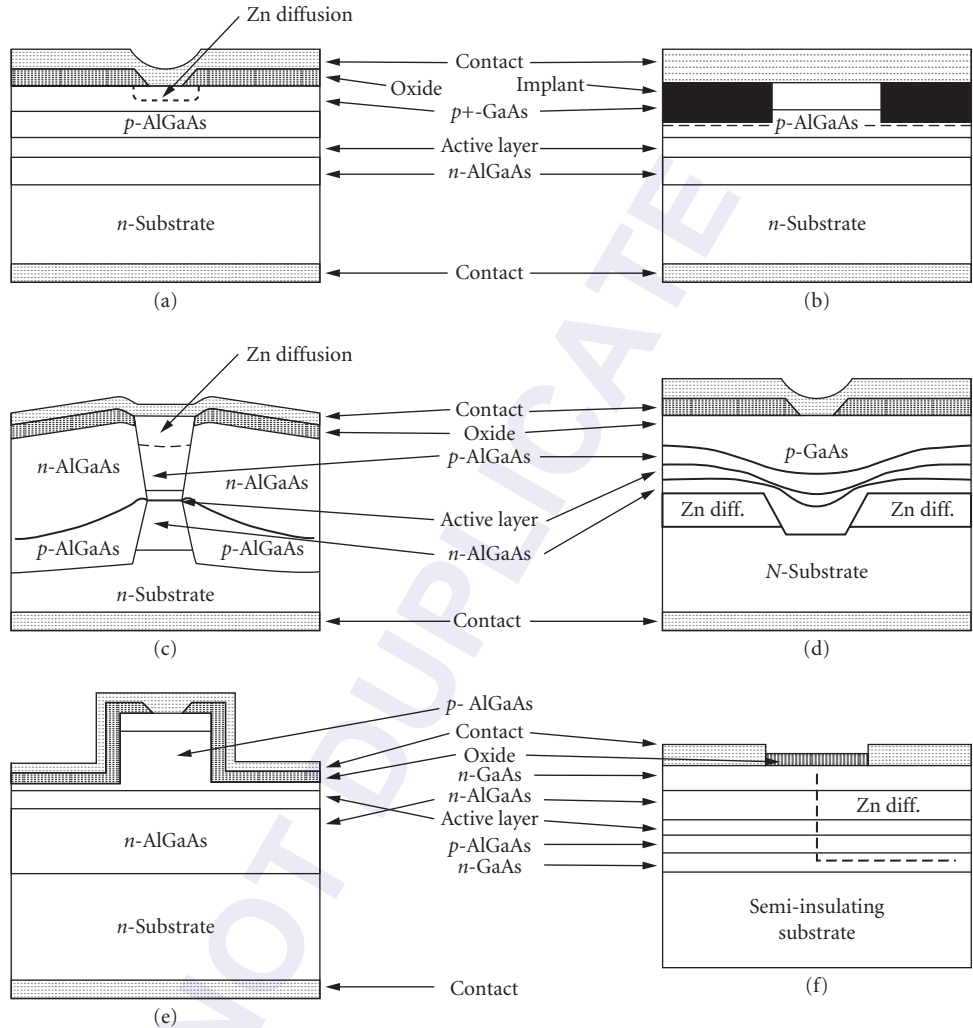


FIGURE 4 Schematic diagrams of GaAs/AlGaAs stripe laser structures. (a) Oxide stripe laser; (b) ion bombardment laser; (c) buried heterostructure (BH) laser; (d) variation on buried heterostructure laser; (e) ridge waveguide (RWG) laser; and (f) transverse junction stripe (TJS) laser.

19.6 QUANTUM WELL LASERS

The active region in a conventional DH semiconductor laser is wide enough ($\sim 1000 \text{ \AA}$) that it acts as bulk material and no quantum effects are apparent. In such a laser the conduction band and valence band are continuous (Fig. 5a). In bulk material the density of states, $D(E)$, for a transition energy E per unit volume per unit energy is²⁶

$$D(E) = \sum_{i=l,h} \frac{m_r^i}{\pi^2 \hbar^3} \sqrt{2m_r^i(E - E_g)} \quad E > E_g \quad (2)$$

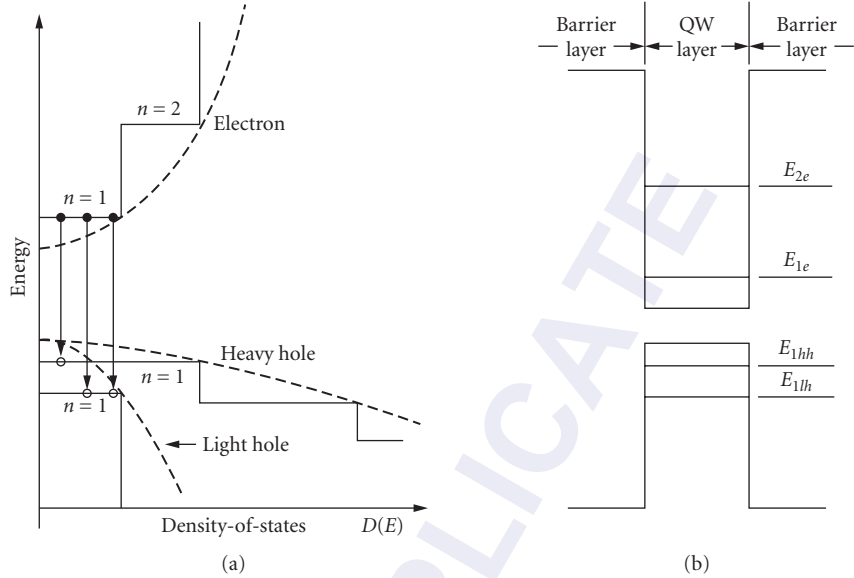


FIGURE 5 Schematic diagrams of (a) the density of states for a quantum well (solid line) and for a bulk DH (dotted line) and (b) quantized energy levels in a quantum well for $n = 1$ and 2 for the conduction band and for the light and heavy hole bands. (After Ref. 63.)

where E_g is the bandgap energy, \hbar is Planck's constant divided by 2π , l and h refer to light and heavy holes, and m_r is the effective mass of the transition which is defined as

$$\frac{1}{m_r} = \frac{1}{m_c} + \frac{1}{m_v} \quad (3)$$

where m_c is the conduction band mass and m_v is the valence band mass. (The split-off hole band and the indirect conduction bands are neglected here and have a negligible effect on most semiconductor laser calculations.)

If the active region of a semiconductor laser is very thin (on the order of the DeBroglie wavelength of an electron) quantum effects become important. When the active region is this narrow (less than ~ 200 Å) the structure is called a *quantum well* (QW). (For a review of QWs see Dingle,⁶¹ Holonyak et al.,⁶² Okamoto,⁶³ or the book edited by Zory.⁶⁴) Since the quantum effects in a QW are occurring in only one dimension they can be described by the elementary quantum mechanical problem of a particle in a one-dimensional quantum box.⁶⁵ In such a well, solution of Schrödinger's equation shows that a series of discrete energy levels (Fig. 5b) are formed instead of the continuous energy bands of the bulk material. With the approximation that the well is infinitely deep, the allowed energy levels are given by

$$E_n = \frac{(n\pi\hbar)^2}{2mL_z^2} \quad (4)$$

where $n = 1, 2, 3, \dots, m$ is the effective mass of the particle in the well, and L_z is the quantum well thickness. Setting the energy at the top of the valence band equal to zero, the allowed energies for an electron in the conduction band of a semiconductor QW become $E = E_g + E_n^c$, where E_n^c is E_n with m

equal to m_c . The allowed energies for a hole in the valence band are then $E = -E_n^v$, where E_n^v is E_n with m equal to m_v . The allowed transition energies E are limited to

$$E = E_g + E_n^c + E_n^v + \frac{\hbar^2 \mathbf{k}^2}{2m_r} \quad (5)$$

where \mathbf{k} is the wavevector, rigorous \mathbf{k} -selection is assumed, and transitions are limited to those with $\Delta n = 0$.

This quantization of energy levels will, of course, change the density of states. For a QW the density of states is given by

$$D(E) = \sum_{i=l,h} \sum_{n=1}^{\infty} \frac{m_r^i}{L_z \pi \hbar^2} H(E - E_g - E_n^c - E_n^v) \quad (6)$$

where H is the heavyside function. The difference in the density of states directly affects the modal optical gain generated by the injection of carriers. The modal gain is proportional to the stimulated emission rate:^{1,66,67}

$$g(E, N) = \alpha \frac{\Gamma D(E) |M|^2 (f_c(E, N) - f_v(E, N))}{E} \quad (7)$$

where I is the optical confinement factor, $|M|^2$ is the matrix element for the transition, N is the carrier density of either electrons or holes (the active region is undoped so they have equal densities), and $f_c(E, N)$ and $f_v(E, N)$ are the Fermi occupation factors for the conduction and valence bands. (For a detailed review of gain in semiconductor lasers see Ref. 67.)

The optical confinement factor Γ is defined as the ratio of the light intensity of the lasing mode within the active region to the total intensity over all space. Since a QW is very thin, Γ_{QW} will be much smaller than Γ_{DH} . Γ_{DH} is typically around 0.5 whereas for a single QW, Γ_{QW} is around 0.03.

The Fermi occupation functions describe the probability that the carriers necessary for stimulated emission have been excited to the states required. They are given by^{1,19}

$$f_c(E_c, N) = \frac{1}{1 + \exp((E_c + E_c - F_c)/kT)} \quad (8)$$

and

$$f_v(E_v, N) = \frac{1}{1 + \exp(-(E_v - F_v)/kT)} \quad (9)$$

where k is the Boltzmann constant, T is temperature, E_c is the energy level of the electron in the conduction band relative to the bottom of the band (including both the quantized energy level and kinetic energy), E_v is the absolute value of the energy level of the hole in a valence band, and F_c and F_v are the quasi-Fermi levels in the conduction and valence bands. Note that E_c and E_v are dependent on E , so f_c and f_v are functions of E . f_c and f_v are also functions of N through F_c and F_v . F_c and F_v are obtained by evaluating the expressions for the electron and hole densities:

$$N = \int D_c(E_c) f_c(E_c) dE_c \quad (10)$$

and

$$N = \int D_v(E_v) f_v(E_v) dE_v \quad (11)$$

where $D_c(E_c)$ and $D_v(E_v)$ are the densities of states for the conduction and valence bands and follow the same form as $D(E)$ for a transition.

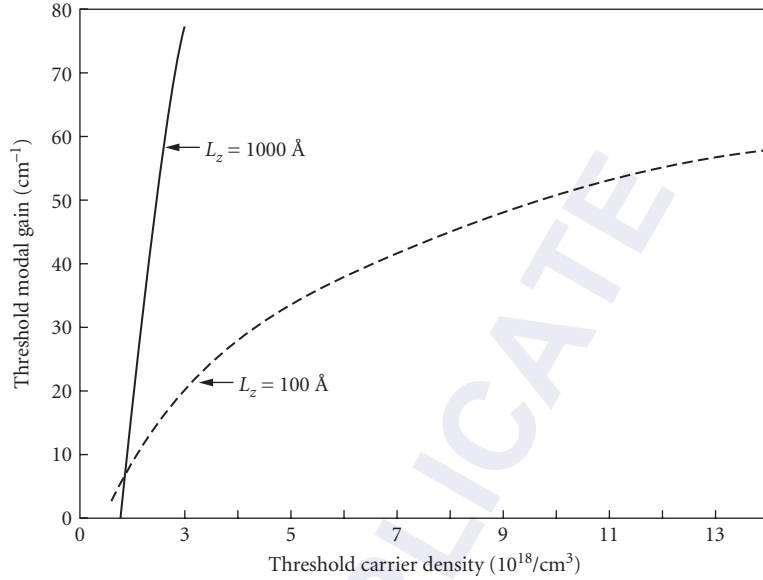


FIGURE 6 Threshold modal gain as a function of threshold carrier density for a conventional (Al, Ga)As double heterostructure with an active region thickness of 1000 Å and for a 100-Å single quantum well. (From Ref. 68.)

In Fig. 6, the results of a detailed calculation⁶⁸ based on Eq. (7) for the threshold modal gain as a function of threshold carrier density are plotted for a 100-Å single QW. The corresponding curve for a DH laser with an active region thickness of 1000 Å is also shown. The gain curves for the QWs are very nonlinear because of saturation of the first quantized state as the carrier density increases. The transparency carrier density N_0 is the carrier density at which the gain is zero. From Fig. 6 it is clear that the transparency carrier densities for QW and DH lasers are very similar and are on the order of $2 \times 10^{18} \text{ cm}^{-3}$.

The advantage of a QW over a DH laser is not immediately apparent. Consider, however, the transparency current density J_o . At transparency²⁸

$$J_o = \frac{N_0 L_z e}{\tau_s} \quad (12)$$

where L_z is the active layer thickness, e is the charge of an electron, and τ_s is the carrier lifetime near transparency. τ_s is approximately 2 to 4 ns for either a QW or a DH laser. Since N_0 is about the same for either structure, any difference in J_o will be directly proportional to L_z . But L_z is approximately 10 times smaller for a QW; therefore, J_o will be approximately 10 times lower for a QW than for a conventional DH laser. (A lower J_o will result in a lower threshold current density since the threshold current density is equal to J_o plus a term proportional to the threshold gain.) Note that this result is not determined by the quantization of energy levels; it occurs because fewer carriers are needed to reach the same carrier density in a QW as in a DH laser. In other words, this result is achieved because the QW is thin!

In this discussion we are considering current density instead of current. The threshold current density (current divided by the length and width of the stripe) is a more meaningful measure of the relative quality of the lasing material than is current. Current depends very strongly on the geometry and stripe fabrication of the device. In order to eliminate geometry-induced variations, current density is normally measured on broad-area (50 to 150 μm wide) oxide stripe lasers (see earlier under

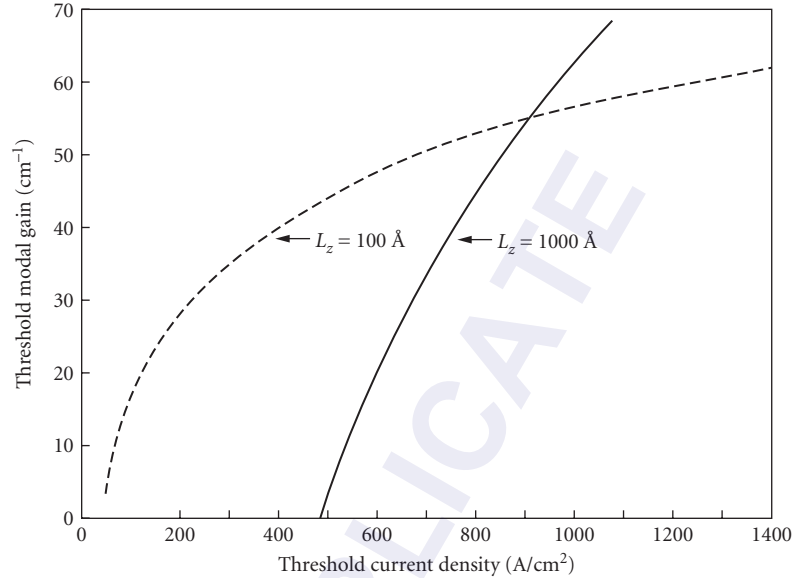


FIGURE 7 Threshold modal gain as a function of threshold current density for a conventional (Al, Ga)As double heterostructure with an active region thickness of 1000 Å and for a 100-Å single quantum well. (From Ref. 68.)

“Fabrication and Configurations”). With a narrow stripe the current spreads beyond the intended stripe width, so it is difficult to accurately measure the current density.

Figure 7 shows the results of a detailed calculation⁶⁸ of the threshold modal gain versus the threshold current density for a DH laser with an active region thickness of 1000 Å and for a 100-Å single QW. The potential for lower threshold current densities for QW lasers is clear for threshold gains less than that where the curve for the DH laser intercepts those of the QWs. With low losses, the threshold current of a QW laser will be substantially lower than that of a DH laser, since the threshold gain will be below the interception point.

To get an appreciation for how the threshold current density of a single QW will compare to that of a DH laser, consider that near transparency, the modal gain is approximately linearly dependent on the current density:

$$g(J) = A(J - J_o) \quad (13)$$

where A is a constant which should have a similar value for either a QW or a DH laser (this can be seen visually on Fig. 7). Taking Eq. (13) at threshold we can equate it to Eq. (1) and solve for J_{th} (the threshold current density):

$$J_{th} = J_o + \frac{\alpha_i}{A} + \frac{1}{2AL} \ln \left(\frac{1}{R_F R_R} \right) \quad (14)$$

α_i is related primarily to losses occurring through the interaction of the optical mode with the active region. In a QW, the optical confinement is lower, which means that the optical mode interacts less with the active region and α_i tends to be smaller. Let's substitute in the numbers in order to get an idea for the difference between a QW and a DH laser. Reasonable values are⁶⁹ $A_{QW} \sim 0.7 A^{-1} \text{ cm}$, $A_{DH} \sim 0.4 A^{-1} \text{ cm}$, $J_o^{QW} \sim 50 \text{ A/cm}^2$, $J_o^{DH} \sim 500 \text{ A/cm}^2$, $\alpha_i^{QW} \sim 2 \text{ cm}^{-1}$, $\alpha_i^{DH} \sim 15 \text{ cm}^{-1}$, $L \sim 400 \mu\text{m}$, and for uncoated facets $R_F = R_R = 0.32$. Substituting in we get $J_{th}^{QW} \sim 95 \text{ A/cm}^2$ and $J_{th}^{DH} \sim 610 \text{ A/cm}^2$.

It is clear that changes in the losses will have a more noticeable effect on threshold current for a QW than for a DH laser since losses are responsible for a more significant portion of the threshold current of a QW laser. The gain curve of a QW laser saturates due to the filling of the first quantized energy level, so operating with low losses is even more important for a QW than is illustrated by the above calculation. When the gain saturates, the simple approximation of Eq. (13) is invalid. Operating with low-end losses is also important for a QW, since they are a large fraction of the total losses. This explains why threshold current density results for QW lasers are typically quoted for long laser cavity lengths (greater than 400 μm), while DH lasers are normally cleaved to lengths on the order of 250 μm . High-quality broad-area single QW lasers (without strain) have threshold current densities lower than 200 A/cm^2 (threshold current densities as low as 93 A/cm^2 have been achieved^{69–71}), while the very best DH lasers have threshold current densities around 600 A/cm^2 .⁷² The end loss can also be reduced by the use of high-reflectivity coatings.²⁵ The combination of a single QW active region with a narrow stripe and high-reflectivity coatings has allowed the realization of submilliampere threshold current semiconductor lasers^{68,69,73} and high-temperature operation.^{74,75,76}

A disadvantage of QW lasers compared to DH lasers is the loss of optical confinement. One of the advantages of a DH laser is that the active region acts as a waveguide, but in a QW the active region is too thin to make a reasonable waveguide. Guiding layers are needed between the QW and the (Al, Ga)As cladding layers. As the bandgap diagram of Fig. 8 illustrates, a graded layer of intermediate aluminum content can be inserted between the QW and each cladding layer. The advantage of this structure, which is called a *graded-index separate-confinement heterostructure* (GRIN SCH),^{77,78} is separate optical and electrical confinement. The carriers are confined in the QW, but the optical mode is confined in the surrounding layers. The grading can be either parabolic (as illustrated in Fig. 8) or linear. Experimentally it has been found that the optimum AlAs mole fraction x for layers around a GaAs QW is approximately 0.2.⁷⁹ Typically, each additional layer is on the order of 2000 \AA thick. In order to confine the optical mode, the cladding layers need a low index of refraction compared to that of an $x = 0.2$ layer. In a simple DH laser, the cladding layers typically have x between 0.3 and 0.4, but for good confinement in an $x = 0.2$ layer, more aluminum should be incorporated into the cladding layers; x should be between 0.5 and 0.7.

In the discussion so far we have considered only single QWs. Structures in which several quantum wells are separated by thin AlGaAs barriers are called *multi-quantum wells* (MQWs) and also have useful properties. For a given carrier density, an MQW with n QWs of equal thickness, L_z has gain which is approximately n times the gain for a single QW of the same thickness L_z , but the current

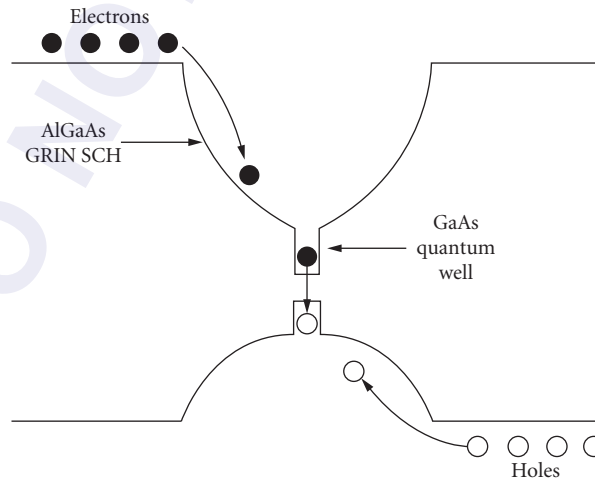


FIGURE 8 Schematic energy-band diagram for a GRIN SCH single QW.

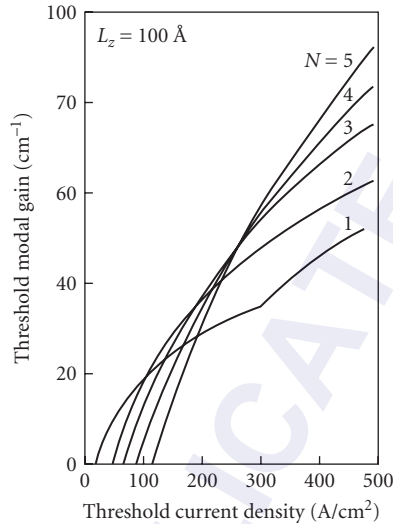


FIGURE 9 Threshold modal gain as a function of threshold current density for (Al, Ga) As single QW and MQWs with 2, 3, 4, and 5 QWs. Each QW has a thickness of 100 Å. (From Ref. 80.)

density is also approximately increased by a factor of n .⁸⁰ The transparency current density will be larger for the MQW than for the single QW since the total active region thickness is larger. Figure 9 shows that, as a function of current density, the gain in the single QW will start out higher than that in the MQW because of the lower transparency current, but the gain in the MQW increases more quickly so the MQW gain curve crosses that of the single QW at some point.⁸¹

Which QW structure has a lower threshold current will depend on how large the losses are for a particular device structure and on where the gain curves cross. The best structure for low threshold current in a GaAs-based laser is normally a single QW, but for some applications involving very large losses and requiring high gain an MQW is superior. For applications in which high output power is more important than low current an MQW is appropriate. An MQW is also preferred for high-modulation bandwidth (see “Spectral Properties” later in this chapter).

An advantage of a QW structure over a bulk DH laser is that the lasing wavelength, which is determined by the bulk bandgap plus the first quantized energy levels, can be changed by changing the quantum well thickness [see Eq. (4)]. A bulk GaAs laser has a lasing wavelength of about 0.87 μm , while a GaAs QW laser of normal thickness (60 to 120 Å) has a lasing wavelength of 0.83 to 0.86 μm . Further bandgap engineering can be introduced with a strained QW.^{67,81,82}

Strained Quantum Well Lasers

Normally, if a semiconductor layer of significantly different lattice constant is grown in an epitaxial structure, it will maintain its own lattice constant and generate misfit dislocations. If this layer is very thin, below a certain critical thickness,^{81,83–85} it will be distorted to match the lattice constant (perpendicular to the substrate) of the surrounding layers and will not generate misfit locations. A layer with thickness above the critical thickness is called “relaxed,” one below is called “strained.”

Straining a semiconductor layer changes the valence-band structure. Figure 10a shows the band structure of an unstrained III–V semiconductor, while Fig. 10b and c show the band structure under biaxial tension and compression, respectively.^{81,82} For the unstrained semiconductor, the light and

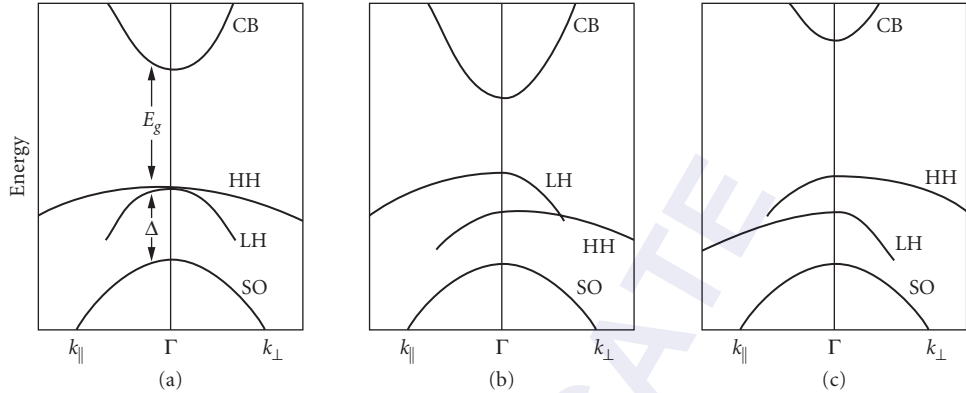


FIGURE 10 Schematic diagrams of the band structure of a III-V semiconductor: (a) unstrained; (b) under tensile strain; and (c) compressively strained. (After Ref. 81.)

heavy hole bands are degenerate at $\mathbf{k}=0(\Gamma)$. Strain lifts this degeneracy and changes the effective masses of the light and heavy holes. In the direction parallel to the substrate the heavy hole band becomes light and the light hole band becomes heavy. Under biaxial tension (Fig. 10b) the bandgap decreases and the “heavy” hole band lies below the “light” hole band. Under biaxial compression the bandgap increases and the “heavy” hole band lies above the “light” hole band.

Figure 10 is a simplification of the true band structure.^{81,82} The bands are not really exactly parabolic, especially the hole bands, and strain increases the nonparabolicity of the hole bands. The details of the band structure can be derived using $\mathbf{k}\cdot\mathbf{p}$ theory.^{86–90}

For a GaAs-based QW, strain can be introduced by adding In to the QW. Since InAs has a larger lattice constant than GaAs, this is a compressively strained QW. In the direction of quantum confinement the highest quantized hole level is the first heavy hole level. This hole level will, therefore, have the largest influence on the density of states and on the optical gain. The effective mass of holes in this level confined in the QW is that parallel to the substrate and is reduced by strain. The reduction of the hole mass within the QW results in a reduced density of hole states [see Eq. (6)].

The reduction in density of hole states is a significant improvement. In order for a semiconductor laser to have optical gain (and lase), $f_c(E, N) - f_v(E, N)$ must be greater than zero [see Eq. (7)]. In an unstrained semiconductor the electrons have a much lighter mass than the heavy holes; the holes therefore have a higher density of states than the electrons. F_c and $f_c(E, N)$ change much more quickly with the injection of carriers than F_v and $f_v(E, N)$. Since approximately equal numbers of holes and electrons are injected into the undoped active layer, reducing the mass of the holes by introducing compressive strain reduces the carrier density required to reach transparency and therefore reduces the threshold current of a semiconductor laser.⁹¹

This theoretical prediction is well supported by experimental results. Strained InGaAs single-QW lasers with record-low threshold current densities of 45 to 65 A/cm² have been demonstrated.^{92–95} These very high-quality strained QW lasers typically have lasing wavelengths from 0.98 to 1.02 μm , QW widths of 60 to 70 \AA , and In concentrations of 20 to 25 percent. InGaAs QWs with wavelengths as long as 1.1 μm have been successfully fabricated,^{96,97} but staying below the critical thickness of the InGaAs layers becomes a problem since the wavelength is increased by increasing the In concentration. (With higher In concentration the amount of strain is increased and the critical thickness is reduced.)

Strained InGaAs QWs have another advantage over GaAs QWs. Strained QW lasers are more reliable than GaAs lasers, i.e., they have longer lifetimes. Even at high temperatures (70 to 100°C), they are very reliable.^{75,76} The reasons for this are not well understood, but it has been suggested that the strain inhibits the growth of defects in the active region.^{98–100} Improving the reliability of GaAs-based lasers is of great practical significance since GaAs lasers are generally less reliable than InP-based lasers.^{4,101,102}

Up to this point our discussion of QW lasers has been limited to GaAs-based QW lasers. QW lasers can also be fabricated in other material systems. GaInP/AlGaInP visible lasers have been improved significantly with the use of a single strained QW active region.^{103–105} These are also compressively strained QWs formed by adding excess In to the active region. This is a much less developed material system than GaAs, so recent results such as 215 A/cm² for a single strained Ga_{0.43}In_{0.57}P QW¹⁰³ are very impressive.

Long-Wavelength (1.3 and 1.55 μm) Quantum Well Lasers

Long-wavelength (InGaAsP/InP) QWs generally do not perform as well as GaAs-based QWs; however, with the advent of strained QW lasers significant progress has been made. Narrow bandgap lasers are believed to be significantly affected by nonradiative recombination processes such as Auger recombination^{4,106–108} and intervalence band absorption.¹⁰⁹ In Auger recombination (illustrated in Fig. 11) the energy from the recombination of an electron and a hole is transferred to another carrier (either an electron or a hole). This newly created carrier relaxes by emitting a phonon; therefore, no photons are created. In intervalence band absorption (IVBA) a photon is emitted, but is reabsorbed to excite a hole from the split-off band to the heavy hole band. These processes reduce the performance of long-wavelength QW lasers enough to make an MQW a lower threshold device than a single QW. As illustrated by Fig. 9, this means that the threshold gain is above the point where the gain versus current density curve of a single QW crosses that of an MQW. Good threshold current density results for lasers operating at 1.5 μm are 750 A/cm² for a single QW and 450 A/cm² for an MQW.¹¹⁰

Long-wavelength QW lasers can be improved by use of a strained QW. For these narrow bandgap lasers strain has the additional benefits of suppressing Auger recombination^{111,112} and intervalence band absorption.¹¹¹ Several groups have demonstrated excellent results with compressively strained InGaAsP/InP QW lasers.^{113–115} Compressively strained single QW lasers operating at 1.5 μm have been demonstrated with threshold current densities as low as 160 A/cm².¹¹⁵

Surprisingly, tensile strained InGaAsP/InP QW lasers also show improved characteristics.^{115–117} Tensile strained QW lasers operating at 1.5 μm have been fabricated with threshold current densities as low as 197 A/cm².¹⁶ These results had not been expected (although some benefit could be expected through suppression of Auger recombination), but have since been explained in terms of TM-mode lasing^{118,119} (normally, semiconductor lasers lase in the TE mode) and suppression of spontaneous emission.¹¹⁸

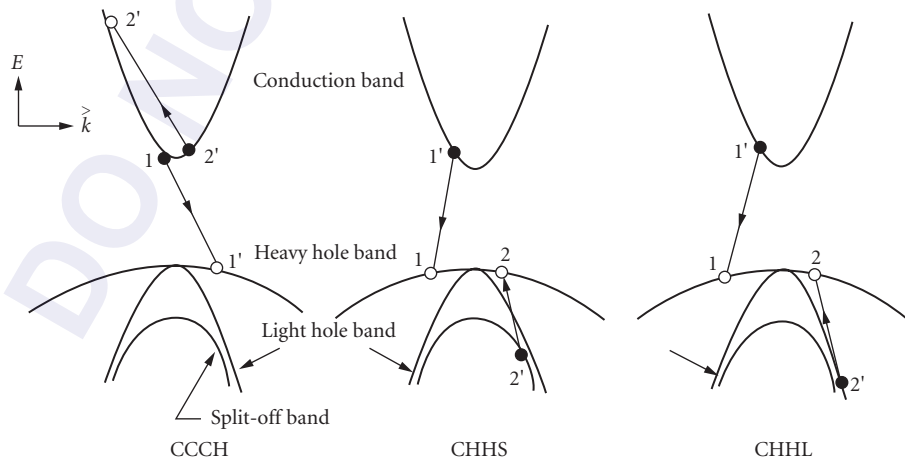


FIGURE 11 Schematic diagrams of band-to-band Auger recombination processes. (After Ref. 4.)

Long-wavelength semiconductor lasers are more sensitive to increases in operating temperature than GaAs-based semiconductor lasers. This temperature sensitivity has been attributed to the strong temperature dependence of Auger recombination^{106–108} and intervalence-band absorption.¹⁰⁹ The use of a strained QW should therefore improve the high-temperature operation of these lasers. This is in fact the case the best results are reported for tensile strained QW lasers with continuous operation at 140°C.^{115,117}

In summary, the use of QW active regions has significantly improved the performance of semiconductor lasers. In this section we have emphasized the dramatic reductions in threshold current density. Improvements have also been realized in quantum efficiency,¹¹⁹ high-temperature operation,^{74–76,115,117} modulation speed (discussed later in this chapter), and spectral linewidth (discussed later). We have limited our discussion to quantum wells. It is also possible to provide quantum confinement in two directions, which is called a *quantum wire*, or three directions, which is called a *quantum dot* or *quantum box*. It is much more difficult to fabricate a quantum wire than a QW, but quantum wire lasers have been successfully demonstrated.¹²¹ For a review of these novel structures we refer the reader to Ref. 122.

19.7 HIGH-POWER SEMICONDUCTOR LASERS

There are several useful methods for stabilizing the lateral modes of an injection laser.^{123–129} In this section, we will discuss techniques for the achievement of high-power operation in a single spatial and spectral mode. There are several physical mechanisms that limit the output power of the injection laser: spatial hole-burning effects lead to multispatial mode operation and are intimately related to multispectral mode operation, temperature increases in the active layer will eventually cause the output power to reach a maximum, and catastrophic facet damage will limit the ultimate power of the laser diode (GaAlAs/GaAs). Thus, the high-power laser designer must optimize these three physical mechanisms to achieve maximum power. In this section, we discuss the design criteria for optimizing the laser power.

High-Power Mode-Stabilized Lasers with Reduced Facet Intensity

One of the most significant concerns for achieving high-power operation and high reliability is to reduce the facet intensity while, at the same time, providing a method for stabilizing the laser lateral mode. Over the years, researchers have developed four approaches for performing this task: (1) increasing the lasing spot size, both perpendicular to and in the plane of the junction, and at the same time introducing a mechanism for providing lateral mode-dependent absorption loss to discriminate against higher-order modes; (2) modifying the facet reflectivities by providing a combination of high-reflectivity and low-reflectivity dielectric coatings; (3) eliminating or reducing the facet absorption by using structures with nonabsorbing mirrors (NAM); (4) using laser arrays and unstable resonator configurations to increase the mode volume. Techniques 1 and 2 are the commonly used techniques and will be further discussed in this section. Techniques 3 and 4 (laser arrays) will be discussed shortly.

Given the proper heat sinking, in order to increase the output power of a semiconductor GaAlAs/GaAs laser, we must increase the size of the beam and thus reduce the power density at the facets for a given power level. The first step in increasing the spot size involves the transverse direction (perpendicular to the junction). There are two approaches for accomplishing this, with the constraint of keeping threshold current low: (1) thinning the active layer in a conventional double-heterostructure (DH) laser (Fig. 12a) below 1000 Å and (2) creating a large optical cavity structure (Fig. 12b).

Thinning the active layer from a conventional value of 0.2 to 0.03 μm causes the transverse-mode spot size to triple for a constant index of refraction step, Δn_r .¹³⁰ The catastrophic power level is proportional to the effective beam width in the transverse direction, d_{eff} , the asymmetric large optical

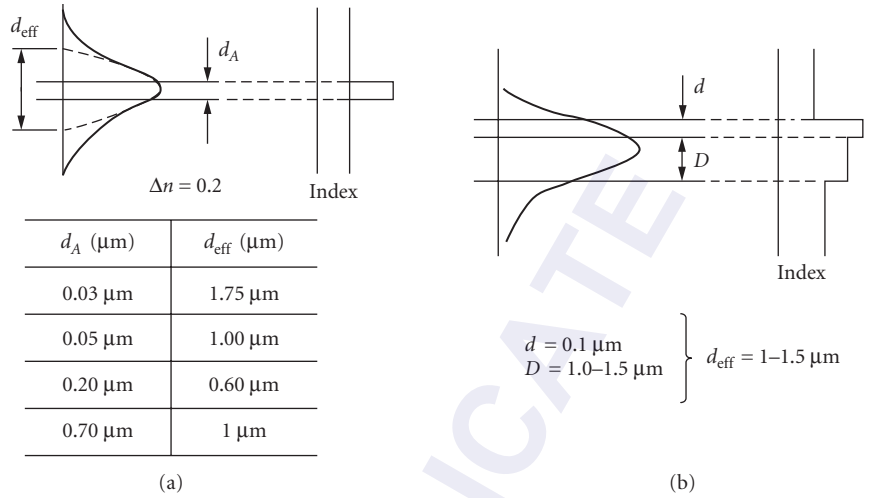


FIGURE 12 Schematic diagrams of the two most commonly used heterostructure configurations for fabricating high-power laser diodes: (a) DH structure and (b) layer large-optical cavity structure. The d_{eff} calculations are after Botez.¹³⁰

cavity (A-LOC) concept^{3,131} involves the epitaxial growth of an additional cladding layer referred to as the *guide layer* (d_G), with index of refraction intermediate between the n -AlGaAs cladding layer and the active layer. By using a relatively small index of refraction step ($\Delta n_r = 0.1$) versus 0.20 to 0.30 for DH lasers, it is possible to force the optical mode to propagate with most of its energy in the guide layer. The effective beam width for the A-LOC can be approximately expressed as

$$d_{\text{eff}} = d_A + d_G \quad (15)$$

where d_A is the active layer thickness. Mode spot sizes in the transverse direction of approximately 1.5 μm can be achieved.

Important Commercial High-Power Diode Lasers

In the last few years, several important high-power laser geometries have either become commercially available or have demonstrated impressive laboratory results. Table 1 summarizes the characteristics of the more important structures. It is evident that the structures that emit the highest cw (continuous wave) power (>100 mW), QW ridge waveguide (QWR), twin-ridge structure (TRS), buried TRS (BTRS), current-confined constricted double-heterostructure large optical cavity (CC-CDH-LOC), and buried V-groove-substrate inner stripe (BVSIS)) have several common features: (1) large spot size (CDH-LOC, TRS, BTRS, QW ridge), (2) low threshold current and high quantum efficiency, and (3) a combination of low- and high-reflectivity coatings. All the lasers with the highest powers, except for the CDH-LOC, use the thin active laser design. A recent trend has been the use of quantum well-active layers.

Figure 13 contains schematic diagrams for five of the more common DH laser designs for high cw power operation, and Fig. 14 shows plots of output power versus current for various important geometries listed in Table 1.

The CC-CDH-LOC device with improved current confinement¹³⁶ (Fig. 13a) is fabricated by one-step liquid-phase epitaxy (LPE) above a mesa separating two dovetail channels. Current confinement is provided by a deep zinc diffusion and a narrow oxide stripe geometry. The final cross-sectional

TABLE 1 Summary of Mode-Stabilized High-Power Laser Characteristics (GaAlAs/GaAs)*†

Manufacturer [Reference]	Geometry	Construction	Rated Power (mW)	Max. cw Power (mW)	Spectral Qual (cw)	Spacial Qual (cw)	I _{th} (mA)	Slope EFF (mW/mA)	Far Field
General Optonics [132]	CNS-LOC	Two-step LPE	—	60	SLM (50)	SSM (50)	50	0.67	12° × 26°
Hitachi [127]	CSP	One-step LPE (TA)	30	100	SLM (40)	SSM (40)	75	0.5	(10–12)° × 27°
MATS. [133]	TRS	One-step LPE	25	115	SLM (50)	SSM (80)	90	0.43	6° × 20°
MATS. [134]	BTRS	Two-step LPE	40	200	SLM (50)	SSM (100)	50	0.8	6° × 16°
NEC [135]	BCM	Two-step LPE	—	80	SLM (80)	SSM (80)	40	0.78	7° × 20°
RCA [136]	CC-CDH	One-step LPE	—	165	SLM (50)	SSM (50)	50	0.77	6° × 30°
RCA [137]	CSP	One-step LPE	—	190	SLM (70)	SSM (70)	50	—	6.5° × 30°
Sharp [138]	V _{SiS}	Two-step LPE	30	100	SLM (50)	SSM (50)	50	0.74	12° × 25°
Sharp [139]	B _{VSiS}	Two-step LPE	—	100	—	SSM (70)	50	0.80	12° × 25°
HP [140]	TCSM	One-step MOCVD	—	65	SLM (65)	SSM (40)	60	0.4	—
TRW [141]	ICSP	Two-step MOCVD (AH/HR)	—	100	SLM (30)	150 (50% duty cycle)	75	0.86	(8–11)° × 35°
Ortel [142]	BH/LOC (NAM)	Two-step LPE (AH/HR)	30	90	—	90	30–50	0.85	—
Spectra Diode [143]	QWR	MOCVD	—	500	SLM(100)	SSM(180)	16	1.3	8° × 22°
BN(STC) [144]	QWR	MOCVD	—	300	SLM(150)	SSM(175)	—	0.8	—

*BH Buried heterostructure

†BTRS Buried TRS

B_{VSiS} Buried V_{SiS}

CC-CDH Current-confined constricted double heterostructure

CNS Channeled narrow planar

CSP Channeled substrate planar

ICSP Inverted CSP

LOC Large optical cavity

NAM Nonabsorbing mirror

QWR Quantum well ridge

SLM Single longitudinal mode

SSM Single spatial mode

TCSM Twin-channel-substrate mesa

TRS Twin-ridge substrate

V_{SiS} V-groove-substrate inner stripe

*Approaches for achieving high-power GaAlAs lasers:

- Thin active (TA) or A-LOC layer to decrease facet power density
- Tight current confinement to produce high current utilization
- Combination of low-/high-reflectivity facet coatings (AR/HR) to produced high differential efficiency and lower facet intensity
- Quantum well design with long cavity

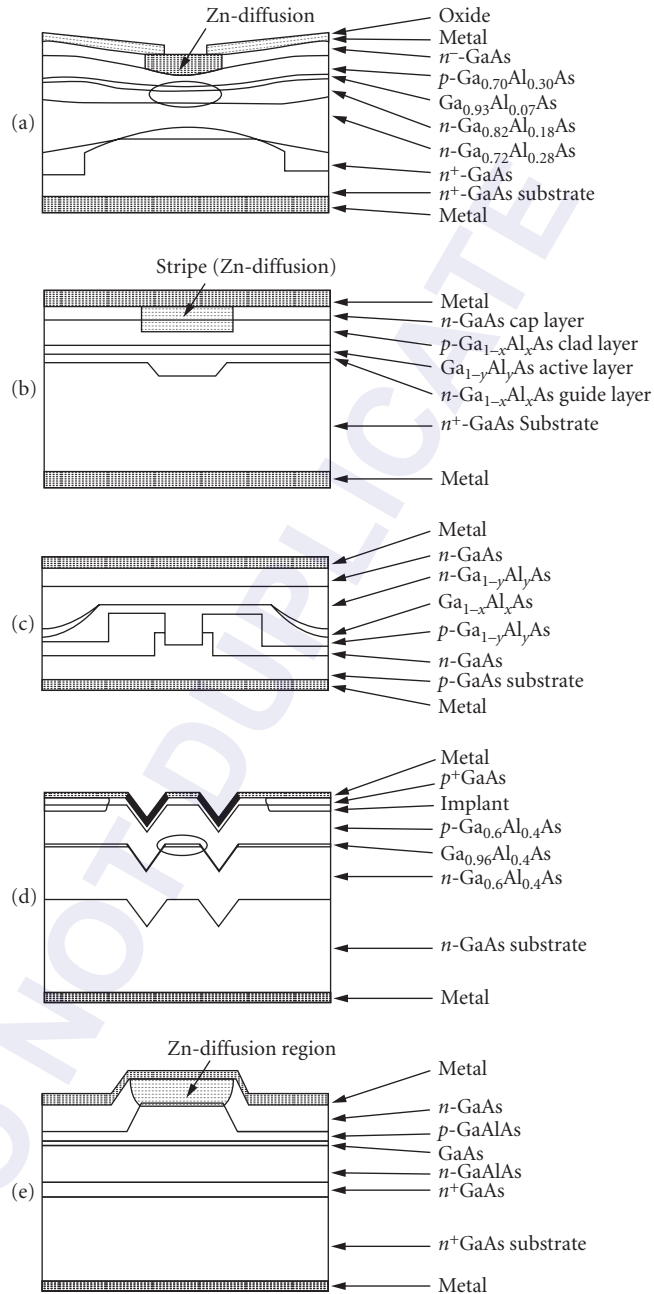


FIGURE 13 Geometries for several important high-power diode lasers. (a) Constricted double-heterostructure large-optical cavity laser (CDH-LOC);¹³⁶ (b) channel substrate planar laser (CSP);¹²⁷ (c) broad-area twin-ridge structure (BTRS);¹³⁴ (d) twin-channel substrate mesa (TCSM); and (e) inverted CSP.

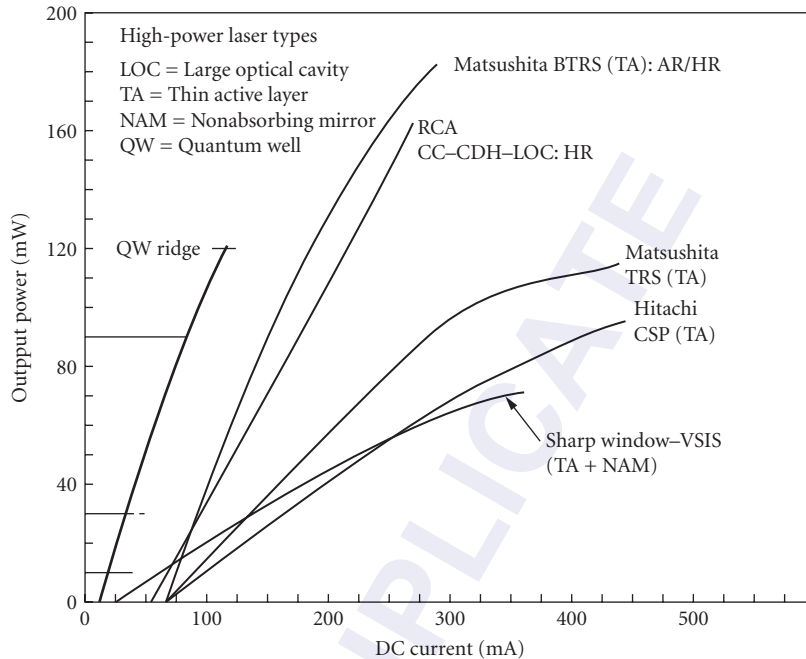


FIGURE 14 Plots showing output power versus cw current for the major high-power laser diodes. The maximum power in a single spatial mode is in the range of 100 to 150 mW and total cw power can approach 200 mW.

geometry of the device is very dependent on the exact substrate orientation and LPE growth conditions.¹⁴⁵ By properly choosing these conditions, it is possible to grow a convex-lens-shaped active layer ($\text{Al}_{0.07}\text{Ga}_{0.93}\text{As}$) on top of a concave-lens-shaped guide layer ($\text{Al}_{0.21}\text{Ga}_{0.79}\text{As}$). The combination of the two leads to a structure with antiwaveguiding properties and a large spot size. Discrimination against higher-order modes is provided by a leaky-mode waveguide. The cw threshold current is in the range of 50 to 70 mA. Single-mode operation has been obtained to 60 mW under 50 percent duty cycle, and the maximum cw power from the device is 165 mW. The power conversion efficiency at this power level is 35 percent considering only the front facet.

The channeled substrate planar (CSP) laser¹²⁷ (Fig. 13b) is fabricated by one-step LPE above a substrate channel. The current stripe is purposely made larger than the channel to ensure uniform current flow across the channel. However, this leads to some waste of current and thus a lower differential efficiency than other similar high-power laser structures (BVSIS, BTRS). Lateral mode control is very effectively obtained by the large difference in the absorption coefficient α between the center and edges of the channel and by changes in the index of refraction that result from changes in the geometry. By proper control of the active and n -cladding layer thicknesses, it is possible to obtain $\Delta\alpha \cong 1000 \text{ cm}^{-1}$ (see Ref. 146) and $\Delta n_c \cong 10^{-2}$. Threshold currents are in the range of 55 to 70 mA. The transverse far field is relatively narrow due to the very thin active layer. Researchers from RCA have obtained power levels in excess of 150 mW (cw) with a CSP-type laser.¹³⁷ A detailed study of the CSP laser has been presented by Lee et al., in a recent publication.¹⁴⁷

Matsushita¹³³ has also developed a CSP-like structure called the twin-ridge structure (TRS) that uses a 400-Å active layer thickness (Fig. 13c). The structure has demonstrated fundamental-mode cw power to 200 mW and single-longitudinal-mode cw power to 100 mW. The maximum available power for the TRS laser is 115 mW, and threshold currents are in the range of 80 to 120 mA. It appears that even though their geometry is similar to that of the CSP, lasers with ultrathin and

planar-active layers have been fabricated. It should be further pointed out that one of the keys to achieving ultrahigh power from CSP-like structures is the achievement of ultrathin (<1000 Å) active layers that are highly uniform in thickness. Small nonuniformities in the active layer thickness lead to a larger Δn_r difference, and thus a smaller lateral spot size, which will lead to lower power levels and reduced lateral mode stability.

Metalorganic chemical vapor deposition (MOCVD, discussed earlier) has been used to fabricate lasers with higher layer uniformity, which leads to a reduced spectral width and more uniform threshold characteristics. Several MOCVD laser structures with demonstrated high-power capability are schematically shown in Fig. 13, and their characteristics are summarized in Table 1. Figure 13*d* shows the twin-channel substrate mesa (TCSM) laser.¹⁴⁰ The fabrication consists of growing a DH laser structure over a chemically etched twin-channel structure using MOCVD. Optical guiding is provided by the curvature of the active layer. The TCSM laser has achieved cw powers of 40 mW in a single spatial mode and 65 mW in a single longitudinal mode. The inverted channel substrate planar (ICSP) laser¹³⁵ is schematically shown in Fig. 13*e*. This structure is one MOCVD version of the very successful CSP structure (Fig. 13*b*).¹⁴¹ The ICSP laser has achieved powers in excess of 150 mW (50 percent duty cycle) in a single spatial mode and a 100-mW (cw) catastrophic power level.

More recently, quantum well lasers using the separate carrier and optical confinement (see previous section, “Quantum Well Lasers”) and ridge waveguide geometries have been used for producing power levels in excess of 150 mW (cw) in a single spatial mode.^{143,144} The QW ridge resembles a standard RWG (Fig. 4*e*), but with a QW active region. Such laser structures have low threshold current density and low internal absorption losses, thus permitting higher-power operation.

Future Directions for High-Power Lasers

Nonabsorbing Mirror Technology The catastrophic facet damage is the ultimate limit to the power from a semiconductor laser. In order to prevent catastrophic damage, one has to create a region of higher-energy bandgap and low surface recombination at the laser facets. Thus, the concept of a laser with a nonabsorbing mirror (NAM) was developed. The first NAM structure was demonstrated by Yonezu et al.¹⁴⁸ by selectively diffusing zinc along the length of the stripe, except near the facets. This created a bandgap difference between the facet and bulk regions and permitted a three- to fourfold increase in the cw facet damage threshold and a four- to fivefold increase in pulse power operation.¹⁴⁹ More recent structures have involved several steps of liquid-phase epitaxy.^{150,151}

The incorporation of the NAM structure is strongly device dependent. For example, in the diffused device structures, such as deep-diffused stripe (DDS)¹⁴⁸ and transverse junction stripe (TJS) lasers, NAM structures have been formed by selective diffusion of zinc in the cavity direction.¹⁴⁹ The *n*-type region will have a wider bandgap than the diffused region, and thus there will be little absorption near the facets. However, most index-guided structures require an additional growth step for forming the NAM region.^{150,151} The NAM structures in the past have suffered from several problems: (1) Due to their complex fabrication, they tend to have low yields. Furthermore, cw operation has been difficult to obtain. (2) Cleaving must be carefully controlled for NAM structures having no lateral confinement, in order to avoid excessive radiation losses in the NAM region. The NAM length is a function of the spot size. (3) The effect of the NAM structure on lateral mode control has not been documented, but could lead to excessive scattering and a rough far-field pattern.

It is now becoming more clear that the use of a NAM structure will be required for the reliable operation of high-power GaAlAs/GaAs laser diodes. Experimental results¹⁵² appear to indicate that laser structures without a NAM region show a decrease in the catastrophic power level as the device degrades. However, most of the approaches currently being implemented require elaborate processing steps. A potentially more fundamental approach would involve the deposition of a coating that would reduce the surface recombination velocity and thus enhance the catastrophic intensity level.^{153,154} Such coatings have been recently used by researchers from Sharp and the University of Florida to increase the uncoated facet catastrophic power level by a factor of 2.^{155,156}

Recently, the use of NAM technology has been appearing in commercial products. The crank transverse junction stripe (TJS) laser (a TJS laser with NAM) can operate reliably at an output

power of 15 mW (cw), while the TJS laser without the NAM can operate only at 3 mW (cw).¹⁴⁷ The Ortel Corporation has developed a buried heterostructure (BH) laser with significantly improved output power characteristics compared to conventional BH lasers.¹⁴² The NAM BH laser is rated at 30 mW (cw)¹⁴² compared to 3 to 10 mW for the conventional BH/LOC device.

Last, the use of alloy disordering, whereby the bandgap of a quantum well laser can be increased by diffusion of various types of impurities (for example, Zn and Si),¹⁵⁷ can lead to a very effective technique for the fabrication of a NAM structure. Such structures have produced an enhancement of the maximum pulsed power by a factor of 3 to 4.

High-Power 1.3/1.48/1.55- μm Lasers Previous sections have discussed high cw power operation from (GaAl)As/GaAs laser devices. In the past several years there have been reports of the increasing power levels achieved with GaInAsP/InP lasers operating at $\lambda = 1.3 \mu\text{m}$. The physical mechanisms limiting high-power operation in this material system are quite different than those for GaAlAs/GaAs lasers. The surface recombination at the laser facets is significantly lower than in GaAlAs/GaAs, and thus catastrophic damage has not been observed. Maximum output power is limited by either heating or carrier leakage effects. With the advent of structures having low threshold current density and high quantum efficiency, it was just a matter of time before high-power results would become available. Furthermore, since facet damage is not a problem, the only real need for facet coatings is for improving the output power from one facet and sealing the device for improved reliability.

In Fig. 15, we schematically show the two most common long-wavelength laser structures that have demonstrated high cw power operation. In Fig. 15a, the double-channel planar buried heterostructure (DC-PBH) is systematically shown.^{158,159} The structure requires a two-step LPE growth process. The first step is the growth of the first and top cladding layers in addition to the active layer. This is followed by the etching of the structure, which is followed by a regrowth to form the blocking and contact layers. LPE growth of this material system is such that if the mesa region is narrow enough, no growth occurs on top of it during the deposition of the blocking layer, and this occurs for mesa widths of less than $\sim 5 \mu\text{m}$. Low threshold current is achieved due to the narrow mesa geometry and the good carrier and current confinement.

The DC-PBH has proved to be a laser structure with excellent output characteristics and high reliability. NEC has been able to obtain thresholds as low as 10 mA with 70 percent quantum efficiency. Degradation rates of the order of $10^{-6}/\text{h}$ for an output power of 5 mW at a temperature of 70°C have also been obtained. More recently, NEC has obtained 140-mW power in a single spatial mode.¹⁵⁸ Lasers at 50 mW and 25°C have been placed on lifetest and show relatively low degradation rates after several hundred hours. Degradation rates at 20 and 30 mW (50°C) are $1.3 \times 10^{-5}/\text{h}$ and $2.22 \times 10^{-5}/\text{h}$, respectively. TRW has also worked with DC-PBH/PBC-(planar buried crescent) type laser diodes and has obtained 100 mW(cw).¹⁵⁹ A summary of the various high-power $\lambda = 1.3\text{-}\mu\text{m}$ laser diode structures and characteristics is given in Table 2.

The other structure that has demonstrated high cw power is the buried crescent laser first investigated by Mitsubishi (Fig. 15b).¹⁶⁰ The structure is grown using a two-step LPE process and a p substrate. The final structure resembles a channel laser with an active layer that tapers to zero near the edges of the channel. The tapering provides good carrier and optical confinement. Researchers from Oki with a structure similar to the Mitsubishi structure have demonstrated maximum power levels of 200 and 140 mW in a single spatial mode.¹⁶² Lifetests¹⁶³ on these lasers demonstrated a mean time to failure of $\sim 7 \times 10^5/\text{h}$ (at 20°C) at 75 percent of the maximum cw output power (maximum = 25 to 85 mW). These results appear to indicate that $1.3\text{-}\mu\text{m}$ lasers are reliable for high-power applications.

A more recent development has been the use of multiquantum well (MQW) high-power lasers in the 1.5 to $1.55\text{-}\mu\text{m}$ wavelength band. The use of an MQW ridge waveguide structure has produced power levels of ~ 170 mW (cw).¹⁶⁴ The MQW structure consists of five wells of InGaAs, 60 Å thick, separated by four GaInAsP barriers of 100-Å thickness. The two thicker, outermost barriers of GaInAsP provide a separate confinement heterostructure (SCH) waveguide. In addition, buried heterostructure (BH) lasers¹⁶⁵ with power levels in excess of 200 mW have been achieved by incorporating strain into MQW structures. A review article by Henshall describes the state of the art in more detail.¹⁶⁶

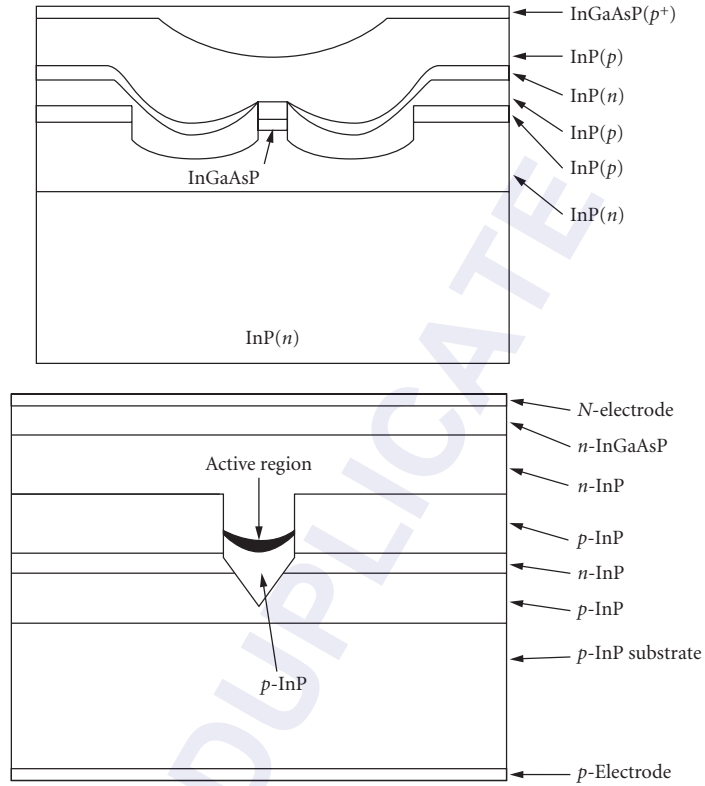


FIGURE 15 Schematic diagrams of the most prominent LPE-based 1.3- μm laser structures: (a) Double-channel planar-buried heterostructure (DC-PBH)^{158,159} and (b) buried crescent.^{160,161}

TABLE 2 Summary of Mode-Stabilized High-Power Laser Characteristics (GaInAsP/InP)^{**†}

Manufacturer [Reference]	Geometry	Construction	Max. Power (mW)	Spatial Quality	I_{th} (mA)
Mitsubishi (160)	PBC	Two-step LPE (p-subst.)	140	SSM (70)	10–30
NEC [157, 158]	DC-PBH	Two-step LPE	140	SSM (140)	10–30
OKI [161, 162]	VIPS	Two-step LPE (p-subst.)	200	SSM (200)	10–30
TRW/EORC [159]	DC-PBH type	Two-step LPE	100	SSM (70)	10–30
TRW/EORC [159]	PBC	Two-step LPE (p-subst.)	107	SSM (78)	10–30
STC [164]	MQW	MOCVD	170	—	—
ATT [165]	MQW BH	MBE	200	—	—

^{*}DC-PBH Double-channel planar buried heterostructure

MQW Multiquantum well

PBC Planar buried crescent

SSM Single spatial mode

[†]Approaches for high-power GaInAsP/InP lasers:

- Tight current confinement to reduce the threshold current
- Facet coatings (reflector/low reflecting front facet)
- Diamond heat sinking
- Long cavity length

TABLE 3 High-Power GaInAs Strained Layer Quantum Well Lasers

Laser Group [Reference]	Ridge Width (μm)	Wavelength (μm)	Threshold Current (mA)	Max. Power in Single Spatial Mode (mW)	Max. cw Power (mW)
JPL [168]	6	0.984–0.989	13	—	24
JPL [168]	3	0.978	8	116	400
NTT [169]	3	0.973–0.983	9	115	500
Spectra Diode [170]	4	0.9–0.91	~20	180	350
Boeing [76]	4	0.98	10–15	150	440

High-Power Strained Quantum Well Lasers Over the last several years, there has been extensive research in the area of strained layer quantum well high-power lasers. As with GaAs QW high-power lasers, the geometry is typically a QW ridge. Table 3 summarizes some of the latest single-spatial-mode high-power results.

Thermal Properties An important parameter in the operation of high-power laser diodes is the optimization of thermal properties of the device. In particular, optimizing the laser geometry for achieving high-power operation is an important design criterion. Arvind et al.¹⁶⁷ used a simple one-dimensional thermal model for estimating the maximum output power as a function of laser geometry (cavity length, active layer thickness substrate type, etc.). The results obtained for GaInAsP/InP narrow stripe PH lasers were as follows:

- Maximum output power is achieved for an optimum active layer thickness in the 0.15- μm region. This result applies only to nonquantum well lasers.
- Significantly higher output powers (25 to 60 percent) are obtained for lasers fabricated on p substrates compared to those on n substrates. The result is based on the lower electrical resistance of the top epitaxial layers in the p substrate compared to n substrate.
- Significantly higher output powers (~60 percent) are obtained for lasers mounted on diamond rather than silicon heat sinks as a result of the higher thermal conductivity of diamond compared to silicon, 22 versus 1.3 W/($^{\circ}\text{C}\text{-cm}$).
- Significantly higher output powers (~100 percent) are obtained for lasers having a length of 700 μm compared to the conventional 300 μm . The higher power results from the reduced threshold current density and thermal resistance for the longer laser devices.

A plot of the calculation and experimental data from Oki¹⁶² is given in Fig. 16. Note that there are no adjustable parameters in the calculation.

An important conclusion from the thermal modeling is that longer cavity semiconductor lasers (700 to 1000 μm) will be able to operate at higher heat sink temperatures when the power level is nominal (~5 mW) compared to shorter cavity devices (~100 to 300 μm). In addition, the reliability of the longer cavity devices is also expected to be better. More recent calculations and experimental results using strained layer lasers have verified this.⁹⁵

Semiconductor Laser Arrays

One of the most common methods used for increasing the power from a semiconductor laser is to increase the width of the emitting region. However, as the width is increased, the occurrence of multilateral modes, filaments, and lateral-mode instabilities becomes more significant. A far-field pattern is produced that is not diffraction-limited and has reduced brightness. The most practical method to overcome this problem is to use a monolithic array of phase-locked semiconductor lasers. Such lasers have been used to generate powers in excess of 10 W (cw)¹⁷⁰ and over 200 W¹⁷¹ from a single laser bar.

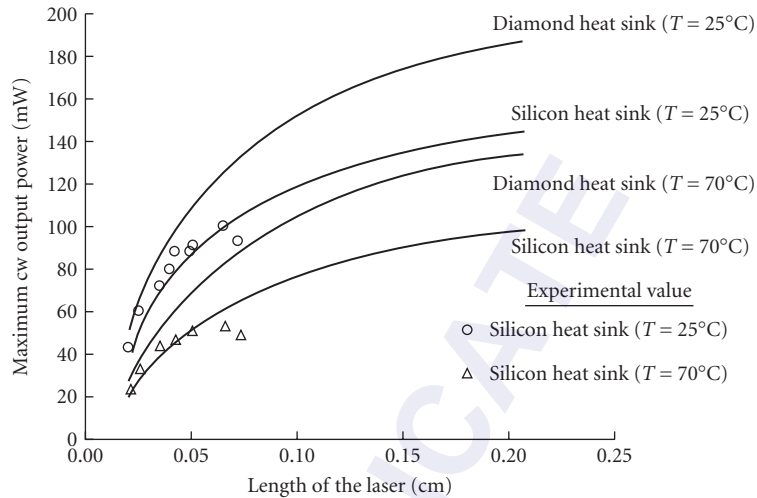


FIGURE 16 Calculation of maximum output power per facet as a function of device length for n - p -substrate-type lasers and different heat sinks.¹⁶⁷ Note the increased power level achieved for longer lasers and p -type substrates.

It was not until 1978 that Scifres and coworkers¹⁷² first reported on the phase-locked operation of a monolithic array consisting of five closely coupled proton-bombarded lasers. The original coupling scheme involved branched waveguides, but this was quickly abandoned in favor of evanescent field coupling by placing the individual elements of the array in close proximity (Fig. 17a). More recently, arrays of index-guided lasers¹⁷³ have been fabricated; one example is shown in Fig. 17b.

Recent emphasis has been on achieving higher cw power and controlling the output far-field distribution. Some of the more significant events in the development of practical semiconductor laser arrays are summarized in Table 4.

The subject of array-mode stability has become of great interest. In a series of significant papers, Butler et al.¹⁹⁰ and Kapon et al.¹⁹¹ recognized that to a first approximation, an array can be modeled as a system of n weakly coupled waveguides. The results indicate that the general solution for the field amplitudes will consist of a superposition of these array modes. The analytic results permitted, for the first time, a simple explanation for the observed far-field patterns and provided a means for designing device structures that would operate in the fundamental array mode (i.e., all elements in phase). This particular mode will provide the greatest brightness.

Many techniques have been used for improving array-mode selection^{179,188,192,197} and thus achieving a well-controlled spatial mode. Two of the more successful earlier techniques involve (1) incorporation of optical gain in the interelement regions (gain coupling of the laser array^{179-184,193}) and (2) use of interferometric techniques that involves Y -coupled junctions.^{186,188,192}

The gain-coupled arrays achieve mode selectivity by introducing optical gain in the interelement regions and thus increasing the gain of the fundamental array mode since this mode has a significant portion of its energy in the interelement regions. The first demonstration of this approach was the twin-channel laser (TCL) developed by researchers from TRW;^{180,181} since then, there have been other demonstrations.^{182,184,193}

The theoretical foundations of the Y -coupled junction were first described in a paper by Chen and Wang.¹⁹² Mode-selectivity is accomplished because the in-phase mode adds coherently at each Y junction, while the out-of-phase mode has destructive interference, since the single waveguides after the Y junction can support only the fundamental mode. Similar interferometric and mode-selective techniques have been used in the development of optical modulators.¹⁹⁸

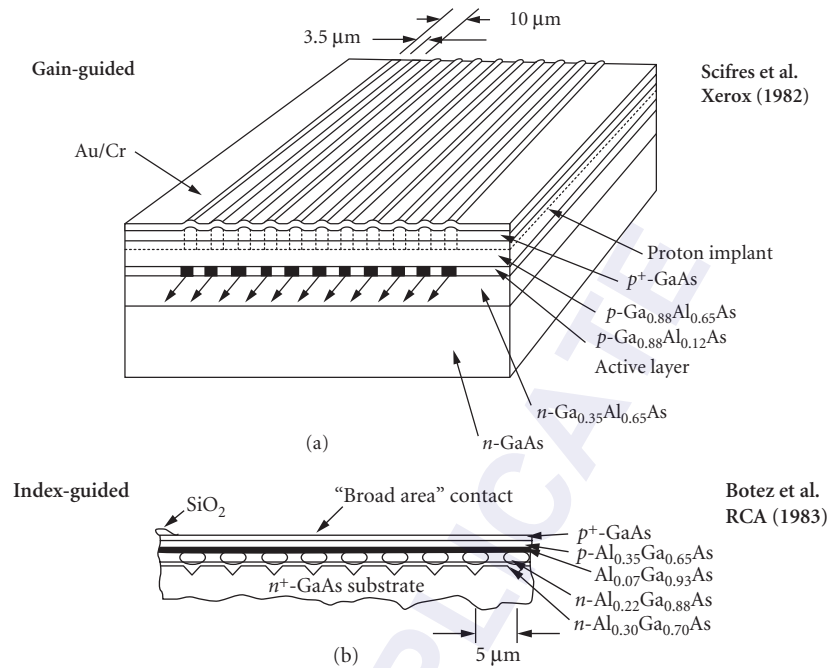


FIGURE 17 Schematic diagrams showing two types of laser array structures: (a) Gain-guided phased array using quantum well-active layers and grown by MOCVD¹⁷⁵ and (b) index-guided phased array using CSP-LOC structures grown by LPE.¹⁷³

TABLE 4 Summary of High-Power Phase-Locked Laser Arrays*

Laser Group [Reference]	No. of Elements	Material System	Type of Array*	Max. Power (mW)	Max. Power (mW)	Far Field
Xerox, 1978 [172]	5	GaAlAs-GaAs	GG	60 (P)	130 (P)	SL (2°)
HP, 1981 [174]	10	GaAlAs-GaAs	IG	1W(P)	1400 (P)	DL
Xerox, 1982 [175]	10	GaAlAs-GaAs	GG	200 (P)	270 (cw)	SL (1°)
Xerox, 1983 [176]	40	GaAlAs-GaAs	GG	800 (P)	2600 (cw)	DL
RCA, 1983 [173]	10	GaAlAs-GaAs	IG	400 (P)	1000 (P)	DL
Siemens, 1984 [177]	40	GaAlAs-GaAs	GG	—	1600 (cw)	DL
Bell Labs, 1984 [178]	10	GaAlAs-GaAs	IG	—	—	DL
TRW, 1984 [179–181]	2	GaAlAs-GaAs	IG	75 (cw)	115 (cw)	SL (4–6°)
UC Berkeley, 1982 [182]	10	GaAlAs-GaAs	IG	—	200	SL (2–7°)
Cal-Tech, 1984 [183]	5	GaAlAs-GaAs	IG	—	—	SL (3°)
Xerox/Spectra Diode, 1985 [184]	10	GaAlAs-GaAs	Offset stripe GG	575 (P)	—	SL (1.9°)
Bell Labs, 1985 [185]	10	InGaAsP-InP	GG	100 (cw)	600	SL (4°)
Sharp, 1985 [186]	2	GaAlAs-GaAs	IG; Y-C	65 (cw)	90 (cw)	SL (4.22°)
Mitsubishi, 1985 [187]	3	GaAlAs-GaAs	IG	100 (cw)	150 (cw)	SL (3.6°)
Xerox/Spectra Diode, 1986 [188]	10	GaAlAs-GaAs	IG(Y-C) stripe GG	200 (cw)	575 (P)	SL (3°)
TRW [189]	10	GaAlAs-GaAs	ROW	380 (cw) 1500 (P)	—	SL (0.7°)

*GG Gain guided
 IG Index guided
 Y-C Y-coupled
 ROW Resonant optical waveguide
 P Pulsed
 DL Double lobe
 SL Single lobe

Finally, the most recent mode control mechanism for laser arrays involves the resonant phase-locking of leaky-mode elements.¹⁸⁹ With a properly optimized geometry, the fundamental array mode has significant mode discrimination analogous to the high discrimination found in single-element leaky-mode devices. Recent results indicate power levels in excess of 360 mW (cw) in the fundamental array mode, and the beam broadens to ~ 2 times diffraction limited for output powers of ~ 500 mW (cw).

At the present time, it is not clear which technique will be most useful for achieving stable, fundamental array mode operation. The gain-coupling concept works well for two or three elements.¹⁹⁹ However, array-mode selection described by the difference in gain between the first and second array modes rapidly decreases as the number of array elements increases beyond two or three.¹⁹⁹ The Y junction and leaky-mode approaches do not appear to have the same limitations. The resonant leaky-mode arrays appear to have the most promising performance at high-power levels; however, the structures are complex and thus yield and reliability need to be more fully addressed.

Two-Dimensional, High-Power Laser Arrays

There has been a significant amount of research activity in the past few years in the area of very high-power diode lasers.^{200–205} The activity has been driven by the significant reductions in threshold current density of GaAlAs/GaAs lasers that can be achieved with metalorganic chemical vapor deposition (MOCVD), utilizing a quantum well design. Threshold current densities as low as 200 to 300 A/cm² with external efficiencies exceeding 80 percent have been achieved using GRIN-SCH quantum well lasers.^{200,201} CW powers of ~ 6 to 9 W have been achieved from single-laser bars. Such power levels correspond to a maximum of 11 W/cm from a single-laser bar. Table 5 lists some of the more recent results on very high-power diode laser arrays.

In order to increase the output power from laser array structures, researchers have investigated the use of a two-dimensional laser array. One particular configuration, referred to as the “rack-stack approach,” is schematically shown in Fig. 18. In essence, the approach involves stacking a linear array of edge emitters into a two-dimensional array. The two-dimensional arrays are fabricated²⁰³ by (1) cleaving linear arrays of laser diodes from a processed wafer, (2) mounting the bars on heat sinks, and (3) stacking the heat sinks into a two-dimensional array.

As shown in Table 5, the main players in this business are McDonnell-Douglas and Spectra Diode Labs. The largest stacked²⁰⁴ two-dimensional array has been manufactured by McDonnell-Douglas and has an active area with five laser bars 8 mm in length. The array was operated with 150- μ s pulses to the limit of the driver²⁰⁴ at pulse repetition rates of 20 to 666 Hz. Approximately 2.5 kW/cm² was obtained at 20 Hz (average power ~ 300 W) and 0.9 kW/cm² at 666 Hz (average power ~ 92 W). Higher output powers will be obtained as a result of achieving the ultimate limits in

TABLE 5 Summary of High-Power Laser Array Results*

Laboratory [Reference]	Array Type	Maximum Output Power	Power Efficiency (%)	Slope Efficiency (W/A)	Power Density (W/cm ²)
General Electric [201]	ID array	80 W (200- μ s pulse; 10–100 Hz)	20	0.9	80
McDonnell-Douglas [203, 204]	Broad stripe ($L = 1200 \mu\text{m}$)	6 W (cw) ($W = 300 \mu\text{m}$)	38	0.91	200
	2D: 4 bars, 8 mm	15 W (cw)	15	—	50
	2D: 5 bars, 8 mm	320 W (0.3% DF [†])	—	—	2560
Spectra Diode [170]	ID array	8 W (cw)	—	—	—
	ID array	134 W (150- μ s pulse)	49	1.26	134

*All high-power laser structures are fabricated using MOCVD and quantum well design.

[†]DF = Duty factor.

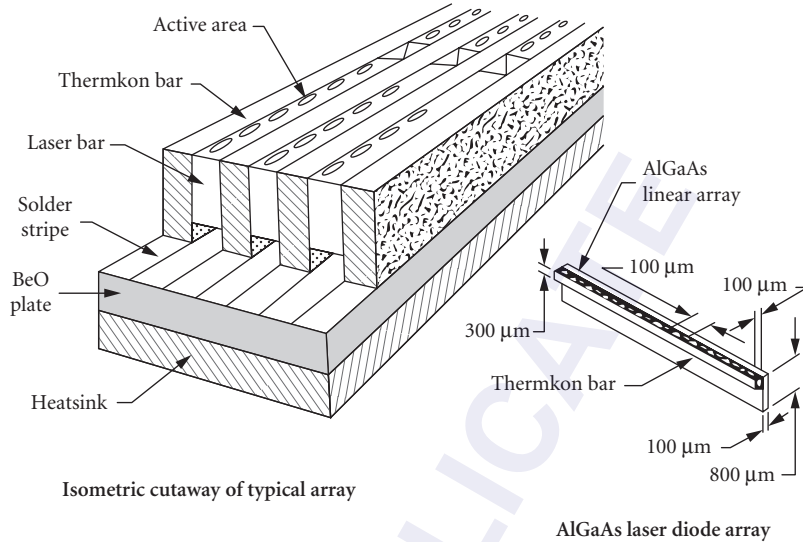


FIGURE 18 Schematic diagram of the two-dimensional rack-stack laser array architecture. (Courtesy of R. Solarz, Lawrence Livermore.)

threshold current density and optical losses in the individual with advances in nonabsorbing mirror and active cooling techniques. Some recent progress has been seen in the latter with the use of etched silicon grooves for fluid flow, which function as radiator elements to remove the heat.

19.8 HIGH-SPEED MODULATION

In many applications semiconductor lasers are modulated in order to carry information. Semiconductor laser dynamics are usually described by the rate equations for the photon and carrier densities:^{3-5,205,206}

$$\frac{dN}{dt} = \frac{I}{edLw} - \frac{cg}{n_r \Gamma} P - \frac{N}{\tau_s} \tag{16}$$

$$\frac{dP}{dt} = \frac{c}{n_r} g P - \frac{P}{\tau_p} + \Gamma \beta \frac{N}{\tau_s} \tag{17}$$

where N is the carrier density, P is the photon density, I is the current, e is the charge of an electron, d is the active layer thickness, L is the laser cavity length, w is the laser stripe width, c is the speed of light, n_r is the refractive index of the active region, g is the threshold modal gain, Γ is the optical confinement factor, τ_s is the carrier lifetime, β is the spontaneous emission factor, and τ_p is the photon lifetime of the cavity. [When these equations are written in terms of total gain instead of modal gain, Γ does not appear in Eq. (16), but multiplies the gain in Eq. (17).]

$$\tau_p = \frac{n_r}{c(\alpha_i + (1/2L)\ln(1/R_F R_R))} \tag{18}$$

where α_i is the internal loss and R_F and R_R are the front- and rear-facet reflectivities, β is the ratio of spontaneous emission power into the lasing mode to the total spontaneous emission rate.²⁰⁷ (Do not confuse β with the other spontaneous emission factor, which is used in linewidth theory and is defined as the ratio of the spontaneous emission power into the lasing mode to the stimulated emission power of the mode.)

When a semiconductor laser is modulated there is some delay before it reaches a steady state. Because it takes time for a carrier population to build up, there will be a time delay τ_d before the final photon density P_{on} is reached (see Fig. 19). Once P_{on} is reached, additional time is required for the carrier and photon populations to come into equilibrium. The output power therefore goes through relaxation oscillations before finally reaching a steady state. This type of oscillation has many parallels in other second-order systems,²⁰⁸ such as the vibration of a damped spring or an RLC circuit.

The frequency of these relaxation oscillations, f_r is called the relaxation, resonance, or corner frequency. By considering small deviations from the steady state where $N = N_{\text{th}} + \Delta N$ and $P = P_{\text{on}} + \Delta P$, we can solve Eqs. (16) and (17) for f_r with the result:^{5,205}

$$f_r \cong \frac{1}{2\pi} \sqrt{\frac{c}{n_r \Gamma} \frac{dg}{dN} \frac{P_{\text{on}}}{\tau_p}} = \frac{1}{2\pi} \sqrt{\frac{c}{n_r} \frac{dg}{dN} \frac{(I - I_{\text{th}})}{edLw}} \quad (19)$$

where dg/dN is the differential modal gain. For bulk double-heterostructure lasers the gain is linearly dependent on the carrier density and $(c/n_r)dg/dN$ is replaced by A , where A is a constant. As discussed earlier in the chapter, the gain-versus-carrier-density relationship of a quantum well is nonlinear, so A is not a constant for a QW laser.

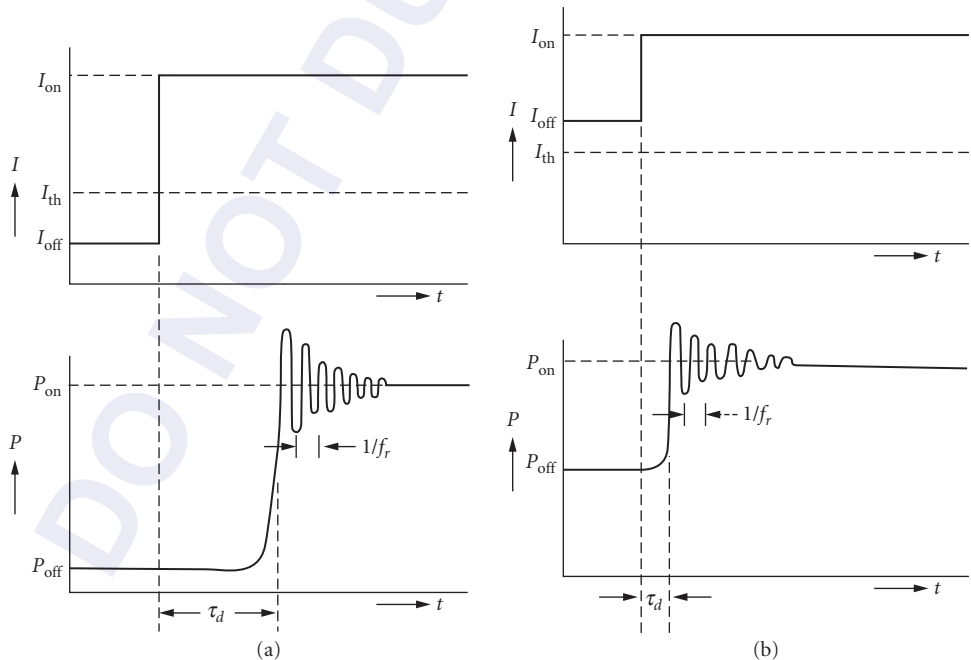


FIGURE 19 Schematic diagrams of the turn on delay and relaxation oscillations for a semiconductor laser (a) prebiased below threshold and (b) prebiased above threshold. (After Ref. 206.)

Let us return our attention to the turn on delay τ_d between a current increase and the beginning of relaxation oscillations as illustrated in Fig. 19. If the initial current I_{off} is below I_{th} , the initial photon density P_{off} can be neglected in Eq. (16). Assuming an exponential increase in carrier density we can derive²⁰⁶

$$\tau_d = \tau_s \ln \left(\frac{I_{\text{on}} - I_{\text{off}}}{I_{\text{on}} - I_{\text{th}}} \right), \quad I_{\text{off}} < I_{\text{th}} < I_{\text{on}} \quad (20)$$

Since τ_s is on the order of several nanoseconds, τ_d is usually very large for semiconductor lasers with I_{off} below I_{th} . For example, with $\tau_s = 4$ ns, $I_{\text{on}} = 20$ mA, $I_{\text{th}} = 10$ mA, and $I_{\text{off}} = 5$ mA, τ_d will be 1.6 ns. If a short current pulse is applied to a laser biased below threshold, τ_d may be so long that the laser barely responds (see Fig. 20). When a current pulse ends, it takes time for the carrier population to decay. If another identical pulse is applied before the carrier population decays fully, it will produce a larger light pulse than the first current pulse did. This phenomenon, which is illustrated in Fig. 20, is called the *pattern effect*.²⁰⁹ The pattern effect is clearly undesirable as it will distort information carried by the laser modulation. The pattern effect can be eliminated by prebiasing the laser at a current sufficient to maintain a carrier population; for most semiconductor lasers this will mean prebiasing at or above threshold. In order to modulate with a prebias below threshold, I_{on} must be much greater than I_{th} . For most lasers this will require an unpractically large I_{on} , but if I_{th} is very low, it may be possible.^{211,212}

Even a semiconductor biased above threshold will have a nonzero τ_d before it reaches its final photon density. Equations (16) and (17) may be solved for τ_d above threshold:²⁰⁶

$$\tau_d = \frac{1}{2\pi f_r} \sqrt{2 \ln \left(\frac{P_{\text{on}}}{P_{\text{off}}} \right)} = \frac{1}{2\pi f_r} \sqrt{2 \ln \left(\frac{I_{\text{on}} - I_{\text{th}}}{I_{\text{off}} - I_{\text{th}}} \right)} \quad I_{\text{th}} < I_{\text{off}} < I_{\text{on}} \quad (21)$$

τ_d will be much shorter for a prebias above threshold. For example if $f_r = 5$ GHz, $I_{\text{on}} = 40$ mA, $I_{\text{off}} = 15$ mA, and $I_{\text{th}} = 10$ mA, $\tau_d = 60$ ps. τ_d will be shortest for large P_{on} and P_{off} , so the shortest time

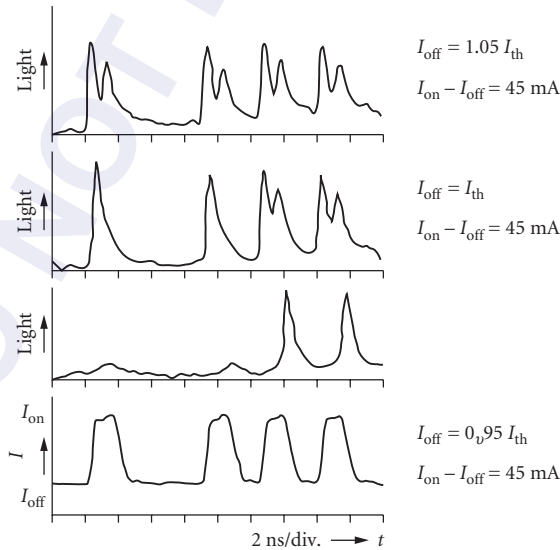


FIGURE 20 An illustration of the pattern effect for a AlGaAs laser diode with a 280 Mbit/s 10111 return to zero pattern. (From Ref. 210.)

delays will be achieved for small-scale modulation at high photon density. For digital applications in which fairly large-scale modulation is required, the maximum modulation speed of a semiconductor laser is to a large extent determined by τ_d .

For very high-speed microwave applications, lasers are prebiased at a current greater than the threshold current and modulated at high frequencies through small amplitudes about the continuous current prebias. The frequency response of a semiconductor laser has the typical shape expected from a second-order system. (For a discussion of the frequency response of a second-order system see Ref. 208.) The laser amplitude response is fairly uniform at frequencies less than the relaxation oscillation frequency. At f_r the response goes through a resonance and then drops off sharply. The relaxation oscillation frequency is therefore the primary intrinsic parameter determining the modulation bandwidth. The actual useful bandwidth is generally considered to be the frequency at which the response of the laser drops by 3 dB. Figure 21 is an example of the frequency response of a semiconductor laser under amplitude modulation.^{213,214} In this example, the maximum 3-dB bandwidth is 16 GHz. If the 3-dB bandwidth is measured in electrical dB, as in our example, it is located at approximately $1.55f_r$. Sometimes the 3-dB frequency is quoted as that at which the optical power is reduced by a factor of 2; this actually corresponds to 6 dB in electrical power and occurs at approximately $1.73f_r$. The 0-dB frequency occurs at approximately $1.41f_r$.^{213–215}

The description of the relaxation oscillation frequency given here is rather simplistic since it is based on rate equations, which consider only one type of carrier, neglect the spatial dependences of the carrier and photon distributions, and neglect the effects of carrier diffusion and nonlinear gain. In addition, the spontaneous emission term of Eq. (17) was neglected in the derivation. The neglected effects are particularly important when considering damping of the

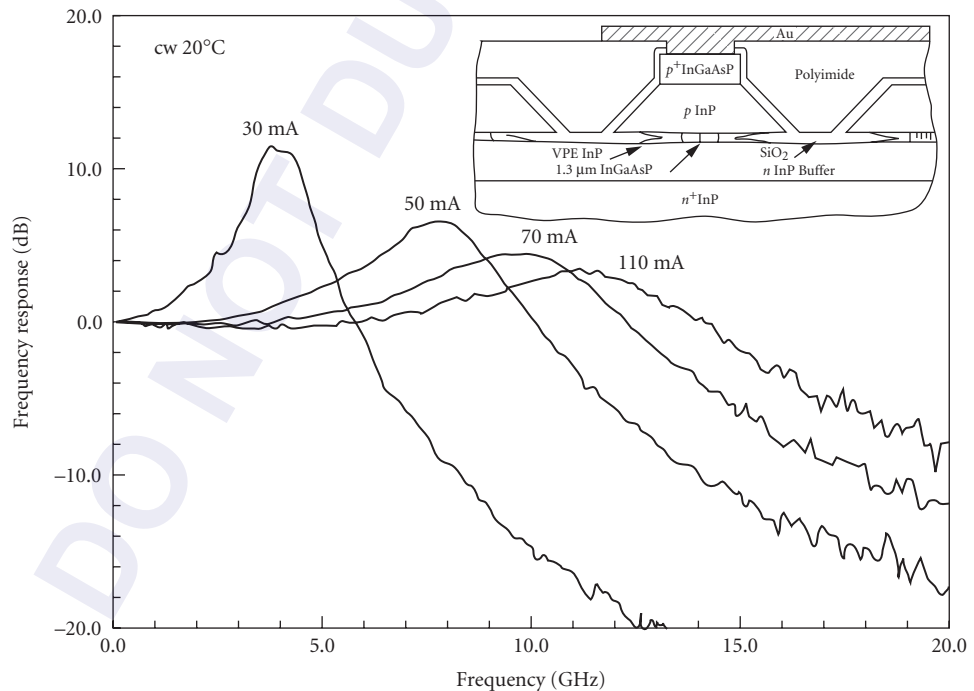


FIGURE 21 The small-signal modulation response of a 1.3- μm InGaAsP-constricted mesa laser for different bias levels. The cavity length is 170 μm and the stripe width is 1 μm . Inset: Schematic diagram of a 1.3- μm InGaAsP-constricted mesa laser. (From Refs. 213 and 214.)

relaxation oscillations.^{204–206} With significant damping²¹⁵ the measured peak frequency f_p will be more accurately determined by

$$f_p^2 = f_r^2 - \frac{f_d^2}{4} \quad (22)$$

where f_d is the damping frequency.

We have also neglected, however, the electrical parasitics of the laser and its operating circuit (bonding wires, etc.). Figure 22 is a simple equivalent circuit, which describes the parasitic elements influencing a semiconductor laser. Here L is the inductance of bond wire, R is the laser resistance including contact resistance, and C is capacitance primarily due to bonding-pad capacitance and capacitance of the current-confining structure of the laser stripe.^{213,214} The 50- Ω resistance is included to represent a 50- Ω drive. The resonant frequency of this circuit is

$$f_0 = \frac{1}{2\pi} \sqrt{\frac{R+50}{LRC}} \quad (23)$$

The circuit is strongly damped so no resonance peak occurs; instead, the response simply drops off.^{205,218} The amplitude response of a modulated semiconductor laser will begin to drop off at frequencies at which the response of the parasitics drops off even if the intrinsic peak frequency of the laser is higher. Therefore, the maximum practical modulation bandwidth may be determined by f_0 instead of f_r . Figure 23 shows the modulation response of a semiconductor laser strongly affected by parasitics.²¹⁹

The most significant parasitic limiting the performance of a semiconductor laser is normally the capacitance.²⁰⁵ In order to achieve high speeds, the laser stripe must be very narrow [see Eq. (18)]; practical narrow stripe lasers are often some form of buried heterostructure (see Fig. 4c). The substrate doping and confinement-layer doping of buried heterostructure form a parallel plate capacitor. In order to reduce the capacitance the laser can be fabricated on a semi-insulating substrate,^{205,218} the confinement layers can be semi-insulating,¹²⁴ the active area of the device can be isolated from the confinement layers by etching trenches on either side of it,^{124,125} or the confinement layers can be replaced by a thick dielectric layer such as polyimide.^{213,214} The inset on Fig. 21 is a schematic diagram of the high-speed laser stripe whose frequency response is shown in Fig. 21.

Assuming that the parasitics have been minimized, consider how a semiconductor laser can be optimized for high-speed operation. As already mentioned, minimizing the stripe width is desirable.

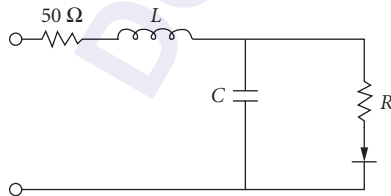


FIGURE 22 Simple equivalent circuit of the parasitics affecting the modulation of a semiconductor laser.

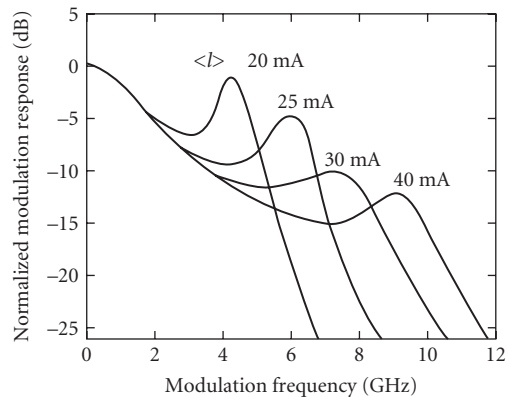


FIGURE 23 The modulation response of a 1.3- μm double-channel, planar-buried heterostructure laser with a cavity length of 80 μm and a threshold current of 18 mA. The effect of parasitics is apparent. (From Ref. 219.)

Figure 21 illustrates that increasing the photon density (or equivalently the current) increases the speed. Of course, there will be a limit as to how much the photon density can be increased; when the photon density of a semiconductor laser is increased, eventually a maximum power is reached at which the laser fails due to catastrophic facet damage. InGaAsP/InP lasers have a higher threshold for catastrophic facet damage than AlGaAs/GaAs lasers; InGaAsP/InP lasers, therefore, tend to have higher bandwidths.²¹⁴

Consideration of Eq. (19) shows that decreasing the cavity length will also increase the speed.^{205,213,214,218,219} (Decreasing the length reduces τ_p .) Increasing dg/dN will also increase the speed. With a bulk active region dg/dN is approximately constant, so use of a short cavity length will not affect it. As illustrated in Fig. 6, however, dg/dN of a QW laser will decrease with increasing threshold gain, and therefore with decreasing cavity length. A single QW laser will, therefore, make a relatively poor high-speed laser. An MQW laser, however, will have higher dg/dN than a single QW. In the past, the highest modulation bandwidths were achieved with InGaAsP/InP bulk active region InGaAsP/InP lasers^{213,214,222,223} with the best results on the order of 24 GHz²²³ for room-temperature CW measurements. The advent of strained MQW lasers has, however, recently resulted in higher bandwidths because strain increases dg/dN .^{224–226} Strained GaAs-based In_{0.3}Ga_{0.7}AsMQW lasers with bandwidths as high as 28 GHz have been demonstrated.²²⁶

So far our discussion has dealt only with the amplitude response to modulation. The phase and lasing wavelength (or optical frequency) are also affected by modulation. Figure 24 is an example of the phase response which accompanies the amplitude response of a semiconductor laser under modulation.

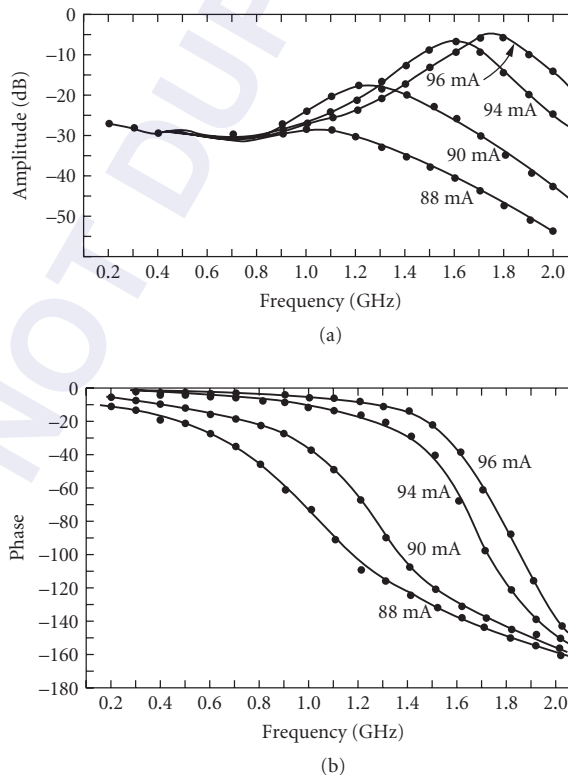


FIGURE 24 The modulation response of a proton stripe laser at various bias currents: (a) amplitude response and (b) phase. The threshold current is approximately 80 mA. (From Ref. 205.)

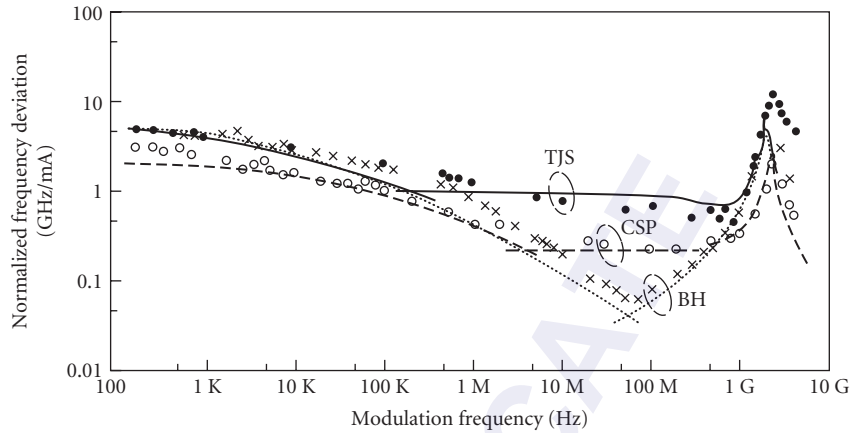


FIGURE 25 The FM response of TJS, BH, and channeled substrate planar (CSP) lasers. All the lasers are biased at 1.2 times threshold. (From Ref. 229.)

Assuming that a semiconductor laser lases in a single longitudinal mode under static conditions, high-frequency modulation can cause it to shift to another mode or to become multimode.²²⁷ The tendency to become multimode increases with the depth of modulation. For many applications single-mode operation under modulation is required. In this case a laser with built-in frequency selectivity such as a distributed feedback laser²²⁸ (DFB) (see “Spectral Properties” following) can be used to maintain single-mode operation. Even with single-mode operation under modulation the linewidth of the lasing mode will be broadened.^{206,227} This broadening which is often called *chirp*, is discussed in more detail in the next section.

While the frequency changes associated with small-scale modulation are generally undesirable in amplitude modulation (AM), they can be utilized for frequency modulation^{206,229,230} (FM). In digital systems, FM is often called frequency shift keying (FSK). FM requires only a very small amplitude modulation, so it normally refers to the effect of modulation on a single mode. In FM modulation, very fine shifts in optical frequency are detected; frequency stabilized lasers such as DFB lasers²³⁰ as well as standard semiconductor lasers will show an FM response. Typically, FM response shows a low frequency decay below f_r and a resonance at f_r ^{206,229} (see Fig. 25).

If the reader requires more in-depth information on high-speed modulation of semiconductor lasers, the book by Petermann²⁰⁶ or the review by Lau and Yariv²⁰⁵ will be particularly helpful. For a recent review of the state of the art see the tutorial by Bowers.²¹⁵ References 3, 4, and 5 also contain chapters on high-speed modulation.

19.9 SPECTRAL PROPERTIES

One of the most important features of a semiconductor laser is its high degree of spectral coherence. There are several aspects to laser coherence. First, the laser must have spatial coherence in the various transverse directions. This is usually accomplished by controlling both the geometry and the lateral-mode geometry using a structure with a built-in index as discussed earlier under “Fabrication and Configurations.” In order to achieve high spectral coherence, the semiconductor laser must operate in a single longitudinal mode. There are four technical approaches for accomplishing this:^{231,232} (1) coupled cavity, (2) frequency selective feedback, (3) injection locking, and (4) geometry control. The various techniques for achieving spectral control are described in Fig. 26.

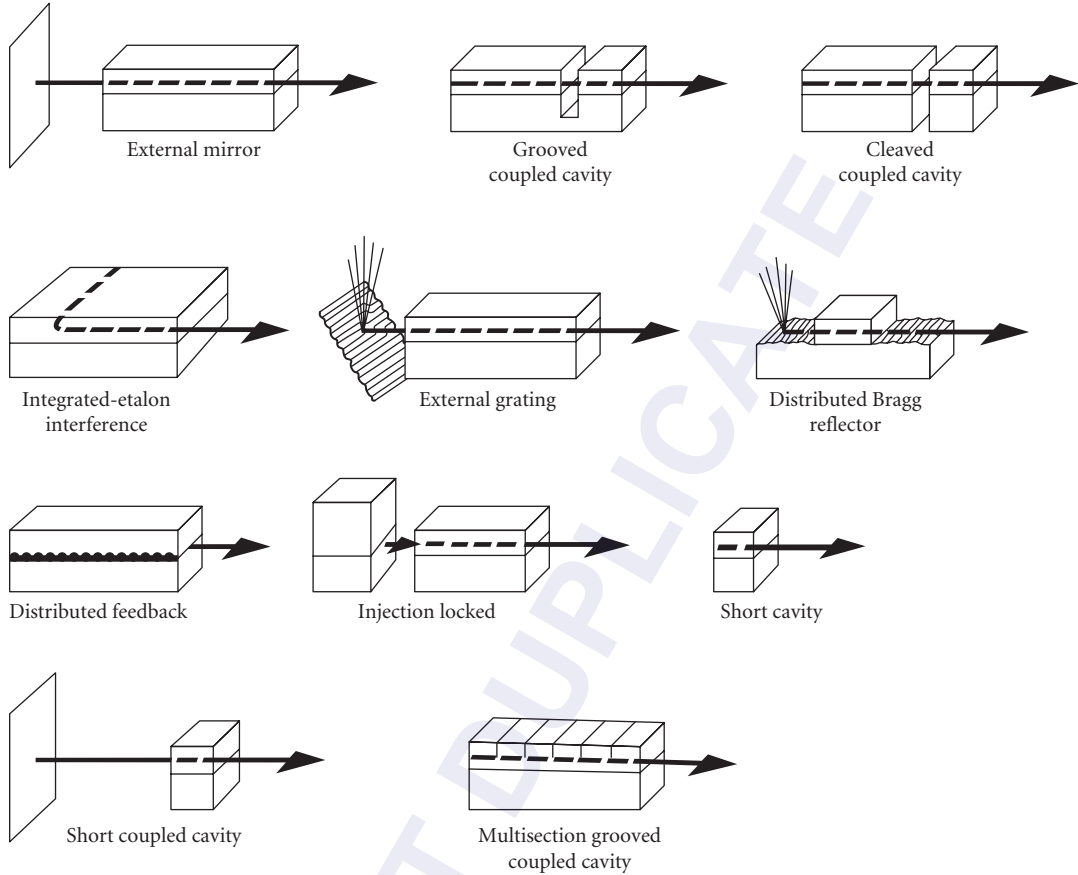


FIGURE 26 Eleven major designs for single-frequency lasers. The three in the top row and the first in the second row are coupled cavity lasers; the next three are frequency-selective-feedback lasers; the next is an injection-locked laser; the last one in the third row is a geometry-controlled laser. The left two are hybrid designs. (From Ref. 232.)

The temporal coherence of the laser is related to the spectral width of the stimulated emission spectrum by

$$L_c = K \frac{\lambda_0^2}{\Delta\lambda_{1/2}} \quad (24)$$

where K is a constant dependent on the distribution of spectral output function, L_c is coherence length, λ_0 is the wavelength of the stimulated emission peak, and $\Delta\lambda_{1/2}$ is the halfwidth of the spectral emission. $K = 1$ for rectangular, $K = 0.32$ for a lorentzian, and $K = 0.66$ for a gaussian.

The spectral linewidth, $\Delta f_{1/2}$, for a single longitudinal mode can be expressed as²³³

$$\Delta f_{1/2} = \frac{n_{sp}}{4\pi\tau_p} \left(\frac{J}{J_{th}} - 1 \right)^{-1} (1 + \alpha^2) \quad (25)$$

where n_{sp} is the spontaneous emission factor, defined as the ratio of spontaneous to stimulated emission in the lasing mode, τ_p is the cavity lifetime, and α is the linewidth enhancement factor.

Typical linewidths for a solitary single-longitudinal-mode laser are in the range of 5 to 20 MHz. Narrower linewidths can be achieved by using some of the techniques described in Fig. 26. More recently, the use of QW lasers, as described earlier, has led to a significant reduction in the linewidth enhancement factor and the corresponding laser linewidths.²³⁴ Typical linewidths in the range of 0.9 to 1.3 MHz have been achieved.

Coupled cavity lasers make up a family of devices, whereby spectral control is achieved by reinforcing certain wavelengths which resonate in several cavities.²³⁵ A typical configuration is shown in Fig. 26, whereby the long cavity is cleaved into two smaller cavities. By properly controlling the length ratios and the gap width, good longitudinal mode discrimination (better than 20-dB side-mode suppression) can be obtained.

Another important technique is the use of an external resonant optical cavity (frequency-selective feedback) as shown in Fig. 27. This technique has been used by researchers at Boeing to achieve extremely narrow linewidth ($\Delta f_{1/2} \sim 1-2$ KHz) single-longitudinal-mode operation.²³⁶

Frequency-selective feedback can also be achieved by using either a distributed feedback (DFB) or distributed Bragg reflector (DBR) laser. As shown in Fig. 26, it differs from other types of lasers in that the feedback is provided by a grating internal to the diode laser. By using a DFB/DBR in combination with a long external cavity, it is possible to achieve linewidths below 1 MHz in a monolithic diode.²³⁷

Injection-locked lasers have also been under investigation at several research labs.²³⁸ In this technique, a low-power, single-frequency laser, which does not have to be a semiconductor laser, is coupled to a single-mode semiconductor laser by injecting the continuous wave emission of a single wavelength of radiation into the laser's cavity.

The last technique for achieving single-longitudinal-mode operation involves the geometry-controlled cavity. Basically, this involves a short cavity 50 μm or less in length, since the longitudinal-mode spacing $\Delta\lambda_L$ in a semiconductor laser is given by¹²⁸

$$\Delta\lambda_L = \frac{\lambda_0^2}{2n_{\text{eff}}L} \tag{26}$$

where n_{eff} is the effective index of refraction and L is the cavity length. Then if $L < 50$ mm, $\Delta\lambda_L > 20$ Å and the gain available to modes away from the gain maximum falls rapidly, the laser operates in a

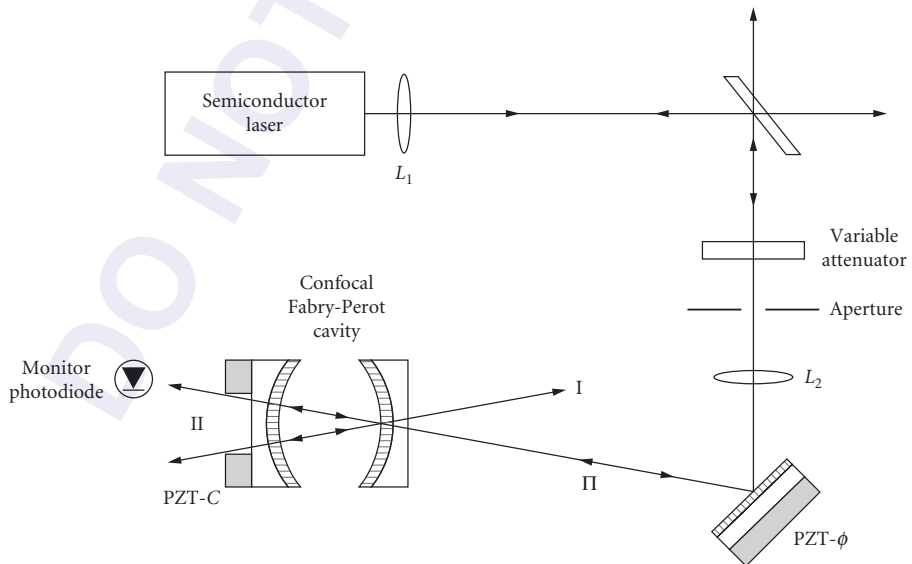


FIGURE 27 Semiconductor laser using resonant optical feedback.

single longitudinal mode. However, the width of the spectral model can still be rather large as dictated by Eq. (25), unless special precautions are made (e.g., ultrahigh mirror reflectivities).

19.10 SURFACE-EMITTING LASERS

Monolithic two-dimensional (2D) laser arrays are key to many important applications such as massive parallel data transfer, interconnect, processing, computing, and high-power, diode-pumped, solid-state lasers. Conventional lasers, as described in previous sections, require a pair of parallel crystalline facets (by cleaving) for delineating the laser cavity, thus limiting laser emission parallel to the junction plane. In this section, we describe laser structures and fabrication techniques which allow light to emit perpendicular to the junction plane, namely, surface-emitting lasers (SEL). SEL structures are compatible with monolithic 2D laser array integration and requirements.

There are three designs for the fabrication of surface-emitting lasers and arrays: (1) in-plane laser with a 45° mirror, (2) in-plane laser with a distributed grating coupler, and (3) vertical cavity laser. The main body of the first two structures is very similar to the conventional cavity design with the axis parallel to the junction plane (in-plane). Light is coupled out from the surface via an integrated mirror or grating coupler. The third structure is an ultrashort cavity ($10\ \mu\text{m}$) "microlaser" requiring no cleaving and compatible with photodiode and integrated-circuit processing techniques. High-density, surface-emitting laser arrays of this type have been demonstrated jointly by AT&T and Bellcore.²³⁹ The following subsections will summarize each of the three structures.

Integrated Laser with a 45° Mirror

The development of this SEL structure requires the wafer processing of two 90° laser mirrors as well as a 45° mirror for deflecting the laser output from the junction plane as shown in Fig. 28. Dry etching techniques such as reactive ion (beam) etching (RIE), chemical-assisted ion beam etching (CAIBE), and ion beam milling are usually used for the fabrication. In combination with a mass-transport process,²⁴⁰ a smooth parabolic sidewall has been demonstrated for the 45° mirror of InGaAsP/InP lasers.

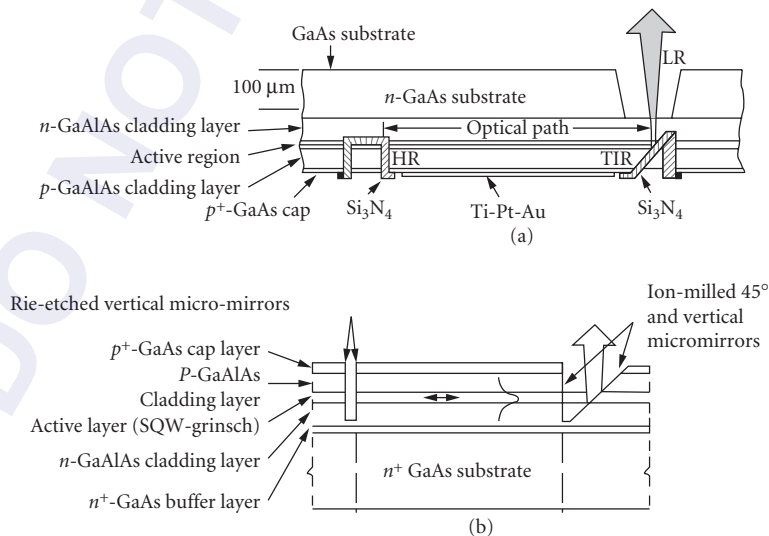


FIGURE 28 Monolithic in-plane-cavity surface-emitting lasers with 45° mirrors: (a) Junction-down and (b) junction-up configurations. (From Ref. 241.)

This approach for an SEL takes advantage of well-established layer structure growth for typical lasers. The laser performance relies on the optical quality and reliability of the facet mirrors formed. The etched-mirror lasers have been improved over the past decade to the stage very comparable with the cleaved lasers. Two-dimensional, high-power (over 1 W) laser arrays have been demonstrated by both TRW²⁴¹ and MIT Lincoln Laboratory.²⁴⁰ These structures would require injection-locking or other external optical techniques in order to achieve coherent phased array operation as mentioned in the high-power laser section.

Distributed Grating Surface-Emitting Lasers

Distributed feedback (DFB) (see under “Spectral Properties,” discussed earlier) and distributed Bragg reflector (DBR) lasers were proposed and demonstrated in the early 1970s. It is well known that for the second-order gratings fabricated in the laser, the first-order diffraction will be normal to the grating surface, as shown in Fig. 29. Since early demonstrations, it has taken over 10 years for both the applications and processing techniques to become mature. Low-threshold, high-reliability DFB lasers with true single-mode characteristics are readily fabricated. The critical issue involved in the fabrication of the laser structure is the fabrication of the gratings with a period on the order of 2000 Å. Holographic interference techniques with an Ar⁺ or He-Cd laser are generally used in many laboratories. The fabrication of large-area gratings with good throughput can be easily achieved by this technique. Another technique involves direct electron-beam writing, which is effective for design iterations.

The development of DBR structure with second-order gratings for surface-emitting lasers did not occur until it was funded by the U.S. Air Force pilot program. This type of laser does not require discrete mirrors for the laser action, so that one could link an array of the lasers with residual in-plane light injection (leaking) across neighboring lasers for coherent operation. A near-diffraction-limited array operation has been demonstrated with this type of SEL. The concept was recently used

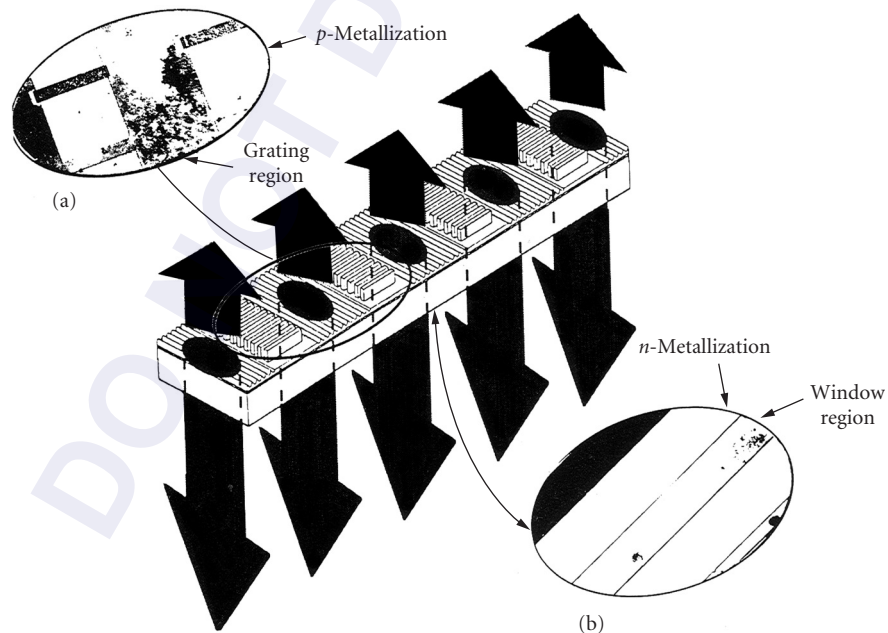


FIGURE 29 Grating surface-emitting laser array. Groups of 10 parallel ridge-guide lasers are laterally coupled in each gain section. (From Ref. 242.)

for a high-power MOPA (master oscillator power amplifier) laser amplifier demonstration with the slave lasers' grating slightly off the resonant second-order diffraction condition. High-power laser arrays of this type have been demonstrated by SRI DSRC²⁴² and Spectra Diode Laboratories.²⁴³ Coherent output powers of 3 to 5 W with an incremental quantum efficiency in excess of 30 percent have been obtained with the array.

Vertical Cavity Lasers

The term "vertical" refers to the laser axis (or cavity) perpendicular to the wafer surface when it is fabricated. Conventional lasers have a relatively long cavity, on the order of 250 μm . It is not practical to grow such a thick layer for the laser. From the analysis, if we reduce the cavity length down to 10 μm , one needs to have a pair of very high reflectivity mirrors to make it lase at room temperature. To satisfy these conditions, researchers at Tokyo Institute of Technology²⁴⁴ have used a metal thin-film or a quarter-wavelength stack of dielectric layers (Bragg reflectors) of high- and low-index material for the mirror post to the growth of laser layers. The advances in epitaxial growth techniques allow an accurate control of semiconductor layer compositions and thicknesses such that Bragg reflectors with 99.9 percent reflectivity can be attained. Therefore, a complete vertical cavity laser structure as shown in Fig. 30 consisting of a gain medium- and high-reflectivity (more than 10 periods of alternate layers due to incremental index difference) mirrors can be grown successfully in one step by MBE or MOCVD techniques.

It is important to optimize the structure for optical gain. To maximize the modal gain, one can locate the standing wave field peak at the thin quantum well-active layer(s) (quantum well lasers were discussed earlier in this chapter) to form a resonant periodic gain structure.²⁴⁵ The issue associated with the semiconductor superlattice Bragg reflectors is the built-in carrier resistance across the abrupt heterojunction discontinuity. Without modifying the structure, the series resistance is on the order of several hundred to a thousand ohms. There have been two techniques applied to lower the resistance, namely, the use of graded junctions²⁴⁶ and peripheral Zn diffusion²⁴⁷ for conducting current to the active region. Both have demonstrated improvement over the original design.

The laser size is defined by etching into a circular column that can be mode-matched to a single-mode fiber for high coupling efficiency. It is desirable that the lasers can be planarized. Proton-bombardment-defined lasers²⁴⁸ with good performance and high yield have been obtained. Meanwhile, the small size of the laser has resulted in low threshold currents close to 1 mA.²⁴⁹ The differential quantum efficiency has been improved from a few percent to more than 30 percent; the output power level, modulation frequency, and maximum operating temperature have also increased over the past several years. As mentioned previously, the advantage of this SEL structure is the potential of high packing density. Bellcore researchers²⁵⁰ have demonstrated a novel WDM (wavelength division multiplexing) laser source with a good histogram of wavelength distribution. The grading of layer thickness across the wafer during a portion of growth translates into different lasing wavelengths. Two-dimensional, individually addressable lasers in a matrix form have also been demonstrated.²⁵¹ In the future, 2D laser arrays operating at a visible wavelength will be very useful for display and optical recording/reading applications. The performance characteristics of vertical cavity SELs reported are shown in Table 6.

19.11 CONCLUSION

In this chapter we have introduced the basic properties of semiconductor lasers and reviewed some areas of the field, including high-power operation, high-speed operation, coherence, and surface-emitting lasers. We have particularly emphasized the advantages of quantum well lasers and strained quantum well lasers. Up until very recently, all the major laser diodes were fabricated using GaAs/GaAlAs and GaInAsP/InP heterostructures. However, there have been such significant advances in the use of strained quantum wells that these lasers have performance levels which exceed, in many

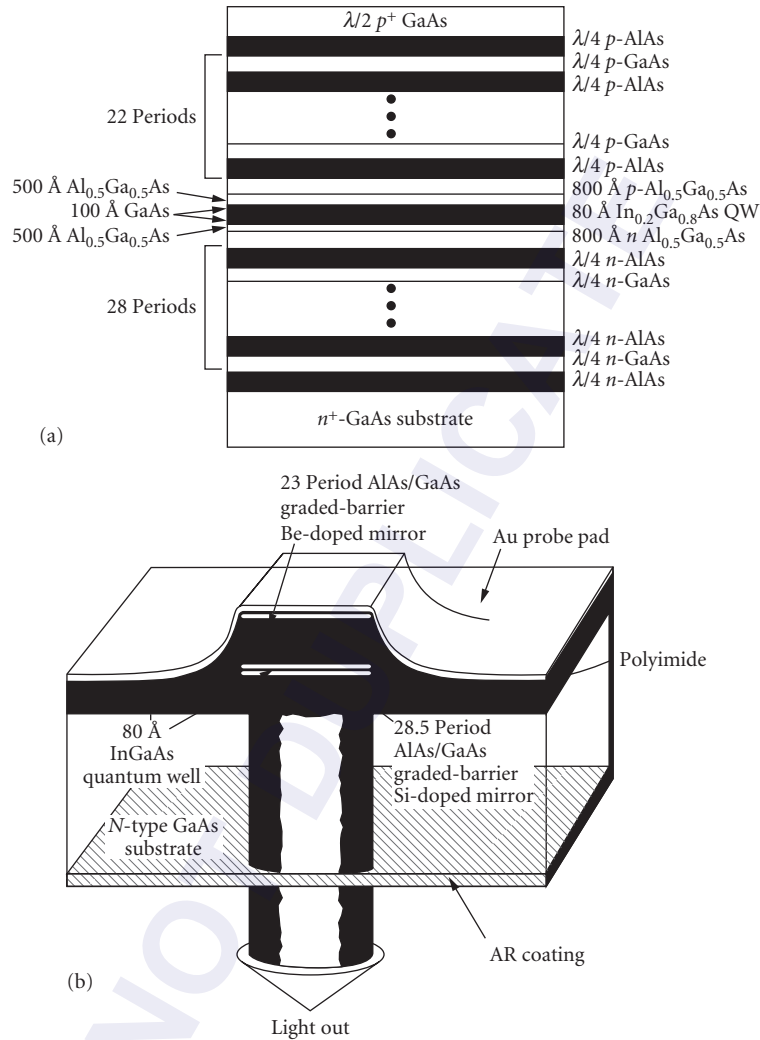


FIGURE 30 Vertical cavity surface-emitting laser with (a) layer structure and (b) device geometry. (From Ref. 249.)

cases, that found for the unstrained lasers. Coupled with the measured excellent reliability results and better output power/temperature performance, these types of lasers will experience a high demand in the future. High-power strained InGaAs/GaAs quantum well lasers are also of interest because their lasing wavelength range includes $0.98 \mu\text{m}$, which makes them useful for pumping erbium-doped fiber amplifiers.¹⁷

Another trend in the future will be the extension of commercial semiconductor lasers to a wider variety of lasing wavelengths. Just as with the standard lasers, strained quantum wells result in significant performance improvements in these more novel laser systems. Shorter wavelength ($\lambda < 0.75 \mu\text{m}$) lasers with high output powers are of interest for high-density optical recording and printing. In the last few years, much progress has been made in developing short-wavelength semiconductor lasers. In the near future, practical visible wavelength lasers will be in the red-to-yellow range, but progress has begun on even shorter wavelengths in the blue. Recent work on very

TABLE 6 Performance Characteristics of Vertical Cavity Surface-Emitting Lasers Developed in Various Research Laboratories*

J_{th} (kA/cm ²)	I_{th} (mA)	V_{th} (V)	T_0 (K)	P_{max} (mW)	η (%)	Size (μ m)	Structure	Reference
22.6	40		115	1.2	1.2	15 ϕ	0.5 mm, DH-DBR	Ref. 246, Tai, AT&T
6.6	2.5			0.3	3.9	7 ϕ	3 \times 8 nm MQW strained columnar μ laser	
6.0	1.5					5 \times 5	10 nm SQW strained columnar μ laser	Ref. 240 Jewell, AT&T
4.1	0.8					5 ϕ	Passivated 3 \times 8 nm MQS strained columnar μ laser	
3.6	3.6	3.7		0.7	4.7	10 \times 10	3 \times 8 nm MQW strained ion-implanted μ laser	Ref. 252 Lee 1990, AT&T
2.8	2.2	7.5		0.6	7.4	10 ϕ	4 \times 10 nm MQW proton- implanted μ laser	
1.4	0.7	4.0				7 \times 7	8 nm SQW strained columnar μ laser	Ref. 249 Geels, UCSB
1.2	4.8			3.0 [†]	12 [†]	20 \times 20	Strained columnar μ laser	Ref. 253 Clausen 1990 Bellcore
1.1	7.5	4.0		3.2	8.3	30 ϕ	4 \times 10 nm MQW proton- implanted μ laser	Ref. 254 Tell 1990, AT&T
0.8	1.1	4.0				12 \times 12	8 nm SQW strained columnar μ laser	Ref. 249 Geels, UCSB
				1.5	14.5	10 ϕ	4 \times 10 nm GRIN-SCH proton-implanted μ laser	Ref. 254 Tell 1990, AT&T

* J_{th} = threshold current density; I_{th} = threshold current; V_{th} = threshold voltage; P_{max} = maximum power; η = overall efficiency; T_0 = characteristic temperature;

[†]Pulsed.

long-wavelength ($\lambda > 2.0 \mu\text{m}$) GaInAsSb/AlGaAsSb lasers is also very promising. Long wavelengths ($> 1.55 \mu\text{m}$) are of interest for eye-safe laser radar, metrology, and medical instrumentation.

Currently, commercial lasers are of the edge-emitting variety, but two-dimensional surface-emitting laser arrays have advanced considerably in the past few years. When they reach maturity, they will be used for pixel interconnect and display applications.

In a limited chapter such as this it is impossible to cover all areas of the field of semiconductor lasers in depth. One of the most important areas neglected is that of tunable lasers.²⁵⁵⁻²⁵⁷ Single-mode tunable DFB and DBR lasers are of great interest for future coherent optical transmission systems. These lasers lase in a single longitudinal mode, but that mode can be tuned to a range of frequencies.

Another area not discussed in detail here is amplifiers. Amplifiers are of great interest for long-haul communication systems, for example, submarine cable systems. Amplifiers can be laser pumped fiber amplifiers¹⁷ or laser amplifiers.²⁵⁸⁻²⁶² A laser amplifier has a structure similar to that of a semiconductor laser and has some optical gain, but only enough to amplify an existing signal, not enough to lase on its own.

It is hoped that the information presented in this chapter will satisfy readers who are interested in the basics of the subject and will give readers interested in greater depth the understanding necessary to probe further in order to satisfy their specific requirements.

19.12 REFERENCES

1. H. C. Casey, Jr. and M. B. Panish, *Heterostructure Lasers, Part A: Fundamental Principles*, Academic Press, Orlando, 1978.
2. H. C. Casey, Jr. and M. B. Panish, *Heterostructure Lasers, Part B: Materials and Operating Characteristics*, Academic Press, Orlando, 1978.

3. G. H. B. Thompson, *Physics of Semiconductor Laser Devices*, John Wiley & Sons, New York, 1980.
4. G. P. Agrawal and N. K. Dutta, *Long-Wavelength Semiconductor Lasers*, Van Nostrand Reinhold, New York, 1986.
5. H. Kressel and J. K. Butler, *Semiconductor Lasers and Heterojunction LEDs*, Academic Press, New York, 1977.
6. R. A. Bartolini, A. E. Bell, and F. W. Spang, *IEEE J. Quantum Electron.* **QE-17**:69 (1981).
7. R. N. Bhargava, *J. Cryst. Growth* **117**:894 (1992).
8. C. Lin (ed.), *Optoelectronic Technology and Lightwave Communication Systems*, Van Nostrand, Reinhold, New York, 1989.
9. S. E. Miller and I. Kaminow (eds.), *Optical Fiber Telecommunications*, Academic Press, Orlando, 1988.
10. M. Katzman (ed.), *Laser Satellite Communications*, Prentice-Hall, Englewood Cliffs, N.J., 1987.
11. M. Ross, *Proc. SPIE* **885**:2 (1988).
12. J. D. McClure, *Proc. SPIE* **1219**:446 (1990).
13. G. Abbas, W. R. Babbitt, M. de La Chappelle, M. L. Fleshner, and J. D. McClure, *Proc. SPIE* **1219**:468 (1990).
14. J. W. Goodman, *International Trends in Optics*, Academic Press, Orlando, 1991.
15. R. Olshansky, V. Lanzisera, and P. Hill, *J. Lightwave Technology* **7**:1329 (1989).
16. R. L. Byer, *Proc. of the CLEO/IQEC Conf.*, Plenary Session, Baltimore, Md., 1987.
17. K. Nakagawa, S. Nishi, K. Aida, and E. Yonoda, *J. Lightwave Technology* **9**:198 (1991).
18. T. F. Deutch, J. Boll, C. A. Poliafito, K. To, *Proc. of the CLEO Conf.*, San Francisco, Calif., 1986.
19. C. Kittel, *Introduction to Solid State Physics*, John Wiley & Sons, New York, 1976.
20. L. Figueroa, "Semiconductor Lasers," *Handbook of Microwave and Optical Components*, K. Chang (ed.), J. Wiley & Sons, New York, 1990.
21. R. N. Hall, G. E. Fenner, J. D. Kingsley, T. J. Soltys, and R. O. Carlson, *Phys. Rev. Lett.* **9**:366 (1962).
22. M. I. Nathan, W. P. Dumke, G. Burns, F. H. Dill, Jr., and G. Lasher, *Appl. Phys. Lett.* **1**:62 (1962).
23. N. Holonyak, Jr. and S. F. Bevacqua, *Appl. Phys. Lett.* **1**:82 (1962).
24. T. M. Quist, R. H. Rediker, R. J. Keyes, W. E. Krag, B. Lax, A. L. McWhorter, and H. J. Zeigler, *Appl. Phys. Lett.* **1**:91 (1962).
25. T. R. Chen, Y. Zhuang, Y. J. Xu, P. Derry, N. Bar-Chaim, A. Yariv, B. Yu, Q. Z. Wang, and Y. Q. Zhou, *Optics & Laser Tech.* **22**:245 (1990).
26. G. L. Bona, P. Buchmann, R. Clauberg, H. Jaeckel, P. Vettiger, O. Voegeli, and D. J. Webb, *IEEE Photon. Tech. Lett.* **3**:412 (1991).
27. A. Behfar-Rad, S. S. Wong, J. M. Ballantyne, B. A. Stolz, and C. M. Harding, *Appl. Phys. Lett.* **54**:493 (1989).
28. N. Bouadma, J. F. Hogrel, J. Charil, and M. Carre, *IEEE J. Quantum Electron.* **QE-23**:909 (1987).
29. M. B. Panish, J. Sumski, and I. Hayashi, *Met. Trans.* **2**:795 (1971).
30. W. T. Tsang (ed.), *Semiconductors and Semimetals*, vol. 22, part A, Academic Press, New York, 1971, pp. 95–207.
31. R. D. Dupuis and P. D. Dapkus, *IEEE J. Quantum Electron.* **QE-15**:128 (1979).
32. R. D. Burnham, W. Streifer, T. L. Paoli, and N. Holonyak, Jr., *J. Cryst. Growth* **68**:370 (1984).
33. A. Y. Cho, *Thin Solid Films* **100**:291 (1983).
34. K. Ploog, *Crystal Growth, Properties and Applications*, vol. 3, H. C. Freyhardt (ed.), Springer-Verlag, Berlin, 1980, pp. 73–162.
35. B. A. Joyce, *Rep. Prog. Phys.* **48**:1637 (1985).
36. W. T. Tsang, *J. Cryst. Growth* **105**:1 (1990).
37. W. T. Tsang, *J. Cryst. Growth* **95**:121 (1989).
38. T. Hayakawa, T. Suyama, K. Takahashi, M. Kondo, S. Yamamoto, and T. Hijikata, *Appl. Phys. Lett.* **51**:707 (1987).
39. A. Kasukawa, R. Bhat, C. E. Zah, S. A. Schwarz, D. M. Hwang, M. A. Koza, and T. P. Lee, *Electron. Lett.* **27**:1063 (1991).
40. J. I. Davies, A. C. Marchall, P. J. Williams, M. D. Scott, and A. C. Carter, *Electron. Lett.* **24**:732 (1988).

41. W. T. Tsang and N. A. Olsson, *Appl. Phys. Lett.* **42**:922 (1983).
42. H. Temkin, K. Alavi, W. R. Wagner, T. P. Pearsall, and A. Y. Cho, *Appl. Phys. Lett.* **42**:845 (1983).
43. D. P. Bour, *Proc. SPIE* **1078**:60 (1989).
44. M. Ishikawa, K. Itaya, M. Okajima, and G. Hatakoshi, *Proc. SPIE* **1418**:344 (1991).
45. G. Hatakoshi, K. Itaya, M. Ishikawa, M. Okajima, and Y. Uematsu, *IEEE J. Quantum Electron.* **QE-27**:1476 (1991).
46. H. Hamada, M. Shono, S. Honda, R. Hiroshima, K. Yodoshi, and T. Yamaguchi, *IEEE J. Quantum Electron.* **QE-27**:1483 (1991).
47. M. A. Haase, J. Qiu, J. M. DePuydt, and H. Cheng, *Appl. Phys. Lett.* **59**:1272 (1991).
48. H. K. Choi and S. J. Eglash, *Appl. Phys. Lett.* **59**:1165 (1991).
49. D. L. Partin, *IEEE J. Quantum Electron.* **QE-24**:1716 (1988).
50. Z. Feit, D. Kostyk, R. J. Woods, and P. Mak, *J. Vac. Sci. Technol.* **B8**:200 (1990).
51. Y. Nishijima, *J. Appl. Phys.* **65**:935 (1989).
52. A. Ishida, K. Muramatsu, H. Takashiba, and H. Fujiasu, *Appl. Phys. Lett.* **55**:430 (1989).
53. Z. Feit, D. Kostyk, R. J. Woods, and P. Mak, *Appl. Phys. Lett.* **58**:343 (1991).
54. R. Zucca, M. Zandian, J. M. Arias, and R. V. Gill, *Proc. SPIE* **1634**:161 (1992).
55. K. Wohlleben and W. Beck, *Z. Naturforsch.* **A21**:1057 (1966).
56. J. C. Dymont, J. C. North, and L. A. D'Asaro, *J. Appl. Phys.* **44**:207 (1973).
57. T. Tsukada, *J. Appl. Phys.* **45**:4899 (1974).
58. H. Namizaki, H. Kan, M. Ishii, and A. Ito, *J. Appl. Phys.* **45**:2785 (1974).
59. H. Namizaki, *IEEE J. Quantum Electron.* **QE-11**:427 (1975).
60. C. P. Lee, S. Margalit, I. Ury, and A. Yariv, *Appl. Phys. Lett.* **32**:410 (1978).
61. R. Dingle, *Festkörper Probleme XV (Advances in Solid State Physics)*, H. Queisser (ed.), Pergamon, New York, 1975, pp. 21–48.
62. N. Holonyak, Jr., R. M. Kolbas, R. D. Dupuis, and P. D. Dapkus, *IEEE J. Quantum Electron.* **QE-16**:170 (1980).
63. N. Okamoto, *Jpn. J. Appl. Phys.* **26**:315 (1987).
64. P. Zory (ed.), *Quantum Well Lasers*, Academic Press, Orlando, 1993.
65. C. Cohen-Tannoudji, B. Diu, and F. Lalöe, *Quantum Mechanics*, vol. 1, John Wiley & Sons, New York, 1977.
66. G. Lasher and F. Stern, *Phys. Rev.* **133**:A553 (1964).
67. S. W. Corzine, R. H. Yan, and L. A. Coldren, *Quantum Well Lasers*, P. Zory (ed.), Academic Press, Orlando, 1993.
68. P. L. Derry, *Properties of Buried Heterostructure Single Quantum Well (Al, Ga)As Lasers*, thesis, Calif. Inst. of Tech., Pasadena, Calif., 1989.
69. P. L. Derry, A. Yariv, K. Y. Lau, N. Bar-Chaim, K. Lee, and J. Rosenberg, *Appl. Phys. Lett.* **50**:1773 (1987).
70. H. Z. Chen, A. Ghaffari, H. Morkoç, and A. Yariv, *Appl. Phys. Lett.* **51**:2094 (1987).
71. H. Chen, A. Ghaffari, H. Morkoç, and A. Yariv, *Electron. Lett.* **23**:1334 (1987).
72. R. Fischer, J. Klem, T. J. Drummond, W. Kopp, H. Morkoç, E. Anderson, and M. Pion, *Appl. Phys. Lett.* **44**:1 (1984).
73. P. L. Derry, T. R. Chen, Y. Zhuang, J. Paslaski, M. Middlestein, K. Vahala, A. Yariv, K. Y. Lau, and N. Bar-Chaim, *Optoelectronics—Dev. and Tech.* **3**:117 (1988).
74. P. L. Derry, R. J. Fu, C. S. Hong, E. Y. Chan, K. Chiu, H. E. Hager, and L. Figueroa, *Proc. SPIE* **1634**:374 (1992).
75. R. J. Fu, C. S. Hong, E. Y. Chan, D. J. Booher, and L. Figueroa, *IEEE Photon. Tech. Lett.* **3**:308 (1991).
76. R. J. Fu, C. S. Hong, E. Y. Chan, D. J. Booher, and L. Figueroa, *Proc. SPIE* **1418**:108 (1991).
77. W. T. Tsang, *Appl. Phys. Lett.* **40**:217 (1982).
78. W. T. Tsang, *Appl. Phys. Lett.* **39**:134 (1981).
79. W. T. Tsang, *Appl. Phys. Lett.* **39**:786 (1981).

80. Y. Arakawa and A. Yariv, *IEEE J. Quantum Electron.* **QE-21**:1666 (1985).
81. E. P. O'Reilly, *Semicond. Sci. Technol.* **4**:121 (1989).
82. E. P. O'Reilly and A. Ghiti, *Quantum Well Lasers*, P. Zory (ed.), Academic Press, Orlando, 1993.
83. J. W. Matthews and A. E. Blakeslee, *J. Cryst. Growth* **27**:118 (1974).
84. I. J. Fritz, S. T. Picraux, L. R. Dawson, T. J. Drummon, W. D. Laidig, and N. G. Anderson, *Appl. Phys. Lett.* **46**:967 (1985).
85. T. G. Andersson, Z. G. Chen, V. D. Kulakovskii, A. Uddin, and J. T. Vallin, *Appl. Phys. Lett.* **51**:752 (1987).
86. M. Altarelli, *Heterojunctions and Semiconductor Superlattices*, G. Allan, G. Bastard, N. Boccarda, M. Lannoo, and M. Voos (eds.), Springer-Verlag, Berlin, 1985, p. 12.
87. J. M. Luttinger and W. Kohn, *Phys. Rev.* **97**:869 (1955).
88. P. Lawaetz, *Phys. Rev.* **B4**:3640 (1971).
89. D. Ahn and S. L. Chuang, *IEEE J. Quantum Electron.* **QE-24**:2400 (1988).
90. S. L. Chuang, *Phys. Rev.* **B43**:9649 (1991).
91. E. Yablonovitch and E. O. Kane, *J. Lightwave Tech.* **LT-4**:504 (1986).
92. H. K. Choi and C. A. Wang, *Appl. Phys. Lett.* **57**:321 (1990).
93. N. Chand, E. E. Becker, J. P. van der Ziel, S. N. G. Chu, and N. K. Dutta, *Appl. Phys. Lett.* **58**:1704 (1991).
94. C. A. Wang and H. K. Choi, *IEEE J. Quantum Electron.* **QE-27**:681 (1991).
95. R. L. Williams, M. Dion, F. Chatenoud, and K. Dzurko, *Appl. Phys. Lett.* **58**:1816 (1991).
96. S. L. Yellen, R. G. Waters, P. K. York, K. J. Beerink, and J. J. Coleman, *Electron Lett.* **27**:552 (1991).
97. P. K. York, K. J. Beerink, J. Kim, J. J. Alwan, J. J. Coleman, and C. M. Wayman, *J. Cryst. Growth* **107**:741 (1991).
98. R. G. Waters, D. P. Bour, S. L. Yellen, and N. F. Ruggieri, *IEEE Photon. Tech. Lett.* **2**:531 (1990).
99. K. Fukagai, S. Ishikawa, K. Endo, and T. Yuasa, *Japan J. Appl. Phys.* **30**:L371 (1991).
100. J. J. Coleman, R. G. Waters, and D. P. Bour, *Proc. SPIE* **1418**:318 (1991).
101. S. Tsuji, K. Mizuishi, H. Hirao, and M. Nakamura, *Links for the Future: Science, Systems & Services for Communications*, P. Dewilde and C. A. May (eds.), IEEE/Elsevier Science, North Holland, 1984, p. 1123.
102. H. D. Wolf and K. Mettler, *Proc. SPIE* **717**:46 (1986).
103. J. Hashimoto, T. Katsyama, J. Shinkai, I. Yoshida, and H. Hayashi, *Appl. Phys. Lett.* **58**:879 (1991).
104. H. B. Serreze, Y. C. Chen, and R. G. Waters, *Appl. Phys. Lett.* **58**:2464 (1991).
105. D. F. Welch and D. R. Scifres, *Electron. Lett.* **27**:1915 (1991).
106. J. I. Pankove, *Optical Processes in Semiconductors*, Dover Publ. Inc., New York, 1971.
107. N. K. Dutta, *J. Appl. Phys.* **54**:1236 (1983).
108. N. K. Dutta and R. J. Nelson, *J. Appl. Phys.* **53**:74 (1982).
109. A. R. Adams, M. Asada, Y. Suematsu, and S. Arai, *Jpn. J. Appl. Phys.* **19**:L621 (1980).
110. T. Tanbun-Ek, R. A. Logan, H. Temkin, K. Berthold, A. F. J. Levi, and S. N. G. Chu, *Appl. Phys. Lett.* **55**:2283 (1989).
111. A. R. Adams, *Electron. Lett.* **22**:249 (1986).
112. Y. Jiang, M. C. Teich, and W. I. Wang, *Appl. Phys. Lett.* **57**:2922 (1990).
113. H. Temkin, R. A. Logan, and T. Tanbun-Ek, *Proc. SPIE* **1418**:88 (1991).
114. C. E. Zah, R. Bhat, R. J. Favire, Jr., S. G. Menocal, N. C. Andreadakis, K. W. Cheung, D. M. Hwang, M. A. Koza, and T. P. Lee, *IEEE J. Quantum Electron.* **27**:1440 (1991).
115. P. J. A. Thijs, L. F. Tiemeijer, P. I. Kuindersma, J. J. M. Binsma, and T. Van Dongen, *IEEE J. Quantum Electron.* **27**:1426 (1991).
116. C. E. Zah, R. Bhat, B. Pathak, C. Caneau, F. J. Favire, Jr., N. C. Andreadakis, D. M. Hwang, M. A. Koza, C. Y. Chen, and T. P. Lee, *Electron. Lett.* **27**:1414 (1991).
117. P. J. A. Thijs, J. J. M. Binsma, E. W. A. Young, and W. M. E. Van Gils, *Electron. Lett.* **27**:791 (1991).
118. E. P. O'Reilly, G. Jones, A. Ghiti, and A. R. Adams, *Electron. Lett.* **27**:1417 (1991).
119. S. W. Corzine and L. A. Coldren, *Appl. Phys. Lett.* **59**:588 (1991).

120. A. Larsson, M. Mittelstein, Y. Arakawa, and A. Yariv, *Electron. Lett.* **22**:79 (1986).
121. S. Simhony, E. Kapon, E. Colas, R. Bhat, N. G. Stoffel, and D. M. Hwang, *IEEE Photon. Tech. Lett.* **2**:305 (1990).
122. K. J. Vahala, J. A. Lebens, C. S. Tsai, T. F. Kuech, P. C. Sercel, M. E. Hoenk, and H. Zarem, *Proc. SPIE* **1216**:120 (1990).
123. N. Chinone, *J. Appl. Phys.* **48**:3237 (1978).
124. P. A. Kirby, A. R. Goodwin, G. H. B. Thompson, D. F. Lovelace, and S. E. Turley, *IEEE J. Quantum Electron.* **QE-13**:720 (1977).
125. R. Lang, *IEEE J. Quantum Electron.* **QE-15**:718 (1979).
126. S. Wang, C. Y. Chen, A. S. Liao, and L. Figueroa, *IEEE J. Quantum Electron.* **QE-17**:453 (1981).
127. K. Aiki, N. Nakamura, T. Kurada, and J. Umeda, *Appl. Phys. Lett.* **30**:649 (1977).
128. M. Nakamura, *IEEE Trans. Circuits Syst.* **26**:1055 (1979).
129. D. Botez, *IEEE Spectrum* **22**:43 (1985).
130. D. Botez, *RCA Rev.* **39**:577 (1978).]
131. H. C. Casey, M. B. Panish, W. O. Schlosser, and T. L. Paoli, *J. Appl. Phys.* **45**:322 (1974).
132. R. J. Fu, C. J. Hwang, C. S. Wang, and B. Lolevic, *Appl. Phys. Lett.* **45**:716 (1984).
133. M. Wada, K. Hamada, H. Himuza, T. Sugino, F. Tujiri, K. Itoh, G. Kano, and I. Teramoto, *Appl. Phys. Lett.* **42**:853 (1983).
134. K. Hamada, M. Wada, H. Shimuzu, M. Kume, A. Yoshikawa, F. Tajiri, K. Itoh, and G. Kano, *Proc. IEEE Int. Semicond. Lasers Conf.*, Rio de Janeiro, Brazil, 1984, p. 34.
135. K. Endo, H. Kawamo, M. Ueno, N. Nido, Y. Kuwamura, T. Furese, and I. Sukuma, *Proc. IEEE Int. Semicond. Laser Conf.*, Rio de Janeiro, Brazil, 1984, p. 38.
136. D. Botez, J. C. Connolly, M. Ettenberg, and D. B. Gilbert, *Electron. Lett.* **19**:882 (1983).
137. B. Goldstein, J. K. Butler, and M. Ettenberg, *Proc. CLEO Conf.*, Baltimore, Md., 1985, p. 180.
138. Y. Yamamoto, N. Miyauchi, S. Maci, T. Morimoto, O. Yamamoto, S. Yomo, and T. Hijikata, *Appl. Phys. Lett.* **46**:319 (1985).
139. S. Yamamoto, H. Hayashi, T. Hayashi, T. Hayakawa, N. Miyauchi, S. Yomo, and T. Hijikata, *Appl. Phys. Lett.* **42**:406 (1983).
140. D. Ackley, *Electron. Lett.* **20**:509 (1984).
141. J. Yang, C. S. Hong, L. Zinkiewicz, and L. Figueroa, *Electron. Lett.* **21**:751 (1985).
142. J. Ungar, N. Bar-Chaim, and I. Ury, *Electron. Lett.* **22**:280 (1986).
143. D. F. Welch, W. Streifer, D. R. Scifres, *Proc. SPIE* **1043**:54 (1989).
144. D. R. Daniel, D. Buckley, B. Garrett, *Proc. SPIE* **1043**:61 (1989).
145. D. Botez, *IEEE J. Quantum Electron.* **QE-17**:2290 (1981).
146. T. Kuroda, M. Nakamura, K. Aiki, and J. Umeda, *Appl. Opt.* **17**:3264 (1978).
147. S. J. Lee, L. Figueroa, and R. Rammaswamy, *IEEE J. Quant. Electron.* **25**:1632 (1989).
148. H. Yonezu, M. Ueno, T. Kamejima, and I. Hayashi, *IEEE J. Quantum Electron.* **15**:775 (1979).
149. H. Kumabe, T. Tumuka, S. Nita, Y. Seiwa, T. Sugo, and S. Takamija, *Jpn. J. Appl. Phys.* **21**:347 (1982).
150. H. Blauvelt, S. Margalit, and A. Yariv, *Appl. Phys. Lett.* **40**:1029 (1982).
151. D. Botez and J. C. Connolly, *Proc. IEEE Int. Semicond. Laser Conf.*, Rio de Janeiro, Brazil, 1984, p. 36.
152. H. Matsubara, K. Ishiki, H. Kumabe, H. Namazaki, and W. Susaki, *Proc. CLEO.*, Baltimore, Md., 1985, p. 180.
153. F. Capasso, and G. F. Williams, *J. Electrochem. Soc.* **129**:821 (1982).
154. H. H. Lee and L. Figueroa, *J. Electrochem. Soc.* **135**:496 (1988).
155. H. Kawanishi, H. Ohno, T. Morimoto, S. Kaneiwa, N. Miyauchi, H. Hayashi, Y. Akagi, Y. Nakajima, *Proc. SPIE* **1219**:309 (1990).
156. J. Yoo, H. Lee, and P. Zory, *IEEE Photonics Lett.* **3**:594 (1991).
157. Y. Suzuki, Y. Horikoshi, M. Kobayashi, and H. Okamoto, *Electron. Lett.* **20**:384 (1984).

158. M. Yamaguchi, H. Nishimoto, M. Kitumara, S. Yamazaki, I. Moto, and K. Kobayashi, *Proc. CLEO*, Baltimore, 1988, p. 180.
159. C. B. Morrison, D. Botez, L. M. Zinkiewicz, D. Tran, E. A. Rezek, and E. R. Anderson, *Proc. SPIE* **893**:84 (1988).
160. Y. Sakakibara, E. Oomura, H. Higuchi, H. Namazaki, K. Ikeda, and W. Susaki, *Electron. Lett.* **20**:762 (1984).
161. K. Imanaka, H. Horikawa, A. Matoba, Y. Kawai, and M. Sakuta, *Appl. Phys. Lett.* **45**:282 (1984).
162. M. Kawahara, S. Shiba, A. Matoba, Y. Kawai, and Y. Tamara, *Proc. Opt. Fiber Commun.* (OFC 1987), paper ME1, 1987.
163. S. Oshiba, A. Matoba, H. Horikawa, Y. Kawai, and M. Sakuta, *Electron. Lett.* **22**:429 (1986).
164. B. S. Bhumbra, R. W. Glew, P. D. Greene, G. D. Henshall, C. M. Lowney, and J. E. A. Whiteaway, *Electron. Lett.* **26**:1755 (1990).
165. T. Tanbun-Ek, R. A. Logan, N. A. Olsson, H. Temkin, A. M. Sergent, and K. W. Wecht, *International Semiconductor Laser Conference*, paper D-3, Davos, 1990.
166. G. D. Henshall, A. Hadjifotiou, R. A. Baker, and K. J. Warwick, *Proc. SPIE* **1418**:286 (1991).
167. M. Arvind, H. Hsing, and L. Figueroa, *J. Appl. Phys.* **63**:1009 (1988).
168. A. Larsson, S. Forouher, J. Cody, and R. J. Lang, *Proc. SPIE* **1418**:292 (1991).
169. M. Okayasu, M. Fukuda, T. Takeshita, O. Kogure, T. Hirone, and S. Uehara, *Proc. of Optical Fiber Communications Conf.*, 29, 1990.
170. D. F. Welch, C. F. Schaus, S. Sun, M. Cardinal, W. Streifer, and D. R. Scifres, *Proc. SPIE* **1219**:186 (1990).
171. G. L. Harnagel, J. M. Haden, G. S. Browder, Jr., M. Cardinal, J. G. Endriz, and D. R. Scifres, *Proc. SPIE* **1219**:186 (1990).
172. D. R. Scifres, R. D. Burnham, and W. Steifer, *Appl. Phys. Lett.* **33**:1015 (1978).
173. D. Botez and J. C. Connally, *Appl. Phys. Lett.* **43**:1096 (1983).
174. D. E. Ackley and R. G. Engelmann, *Appl. Phys. Lett.* **39**:27 (1981).
175. D. R. Scifres, R. D. Burnham, W. Streifer, and M. Bernstein, *Appl. Phys. Lett.* **41**:614 (1982).
176. D. R. Scifres, C. Lindstrom, R. D. Burnham, W. Streifer, and T. L. Paoli, *Appl. Phys. Lett.* **19**:160 (1983).
177. F. Kappeler, H. Westmeier, R. Gessner, M. Druminski, and K. H. Zschauer, *Proc. IEEE Int. Semicond. Laser Conf.*, Rio de Janeiro, Brazil, 1984, p. 90.
178. J. P. Van der Ziel, H. Temkin, and R. D. Dupuis, *Proc. IEEE Int. Semicond. Laser Conf.*, Rio de Janeiro, Brazil, 1984, p. 92.
179. L. Figueroa, C. Morrison, H. D. Law, and F. Goodwin, *Proc. Int. Electron Devices Meeting*, 1983, p. 760.
180. L. Figueroa, C. Morrison, H. D. Law, and F. Goodwin, *J. Appl. Phys.* **56**:3357 (1984).
181. C. Morrison, L. Zinkiewicz, A. Burghard, and L. Figueroa, *Electron. Lett.* **21**:337 (1985).
182. Y. Twu, A. Dienes, S. Wang, and J. R. Whinnery, *Appl. Phys. Lett.* **45**:709 (1984).
183. S. Mukai, C. Lindsey, J. Katz, E. Kapon, Z. Rav-Noy, S. Margalit, and A. Yariv, *Appl. Phys. Lett.* **45**:834 (1984).
184. D. F. Welch, D. Scifres, P. Cross, H. Kung, W. Streifer, R. D. Burnham, and J. Yaeli, *Electron. Lett.* **21**:603 (1985).
185. N. Dutta, L. A. Kozzi, S. G. Napholtz, and B. P. Seger, *Proc. Conf. Lasers Electro-Optics (CLEO)*, Baltimore, Md., 1985, p. 44.
186. M. Taneya, M. Matsumoto, S. Matsui, Y. Yano, and T. Hijikata, *Appl. Phys. Lett.* **47**:341 (1985).
187. J. Ohsawa, S. Himota, T. Aoyagi, T. Kadowaki, N. Kaneno, K. Ikeda, and W. Susaki, *Electron. Lett.* **21**:779 (1985).
188. D. F. Welch, P. S. Cross, D. R. Scifres, W. Streifer, and R. D. Burnham, *Proc. CLEO*, San Francisco, Calif., 1986, p. 66.
189. L. Mawst, D. Botez, E. R. Anderson, M. Jansen S. Ou, M. Sargent, G. L. Peterson, and T. J. Roth, *Proc. SPIE* **1418**:353 (1991).
190. J. K. Butler, D. E. Ackley, and D. Botez, *Appl. Phys. Lett.* **44**:293 (1984).
191. E. Kapon, J. Katz, and A. Yariv, *Opt. Lett.* **10**:125 (1984).
192. K. L. Chen and S. Wang, *Electron. Lett.* **21**:347 (1985).

193. W. Streifer, A. Hardy, R. D. Burnham, and D. R. Scifres, *Electron. Lett.* **21**:118 (1985).
194. S. Chinn and R. J. Spier, *IEEE J. Quantum Electron.* **20**:358 (1985).
195. J. Katz, E. Kapon, C. Lindsey, S. Margalit, U. Shreter, and A. Yariv, *Appl. Phys. Lett.* **42**:521 (1983).
196. E. Kapon, C. P. Lindsey, J. S. Smith, S. Margalit, and A. Yariv, *Appl. Phys. Lett.* **45**:1257 (1984).
197. D. Ackley, *Electron. Lett.* **20**:695 (1984).
198. T. R. Ranganath and S. Wang, *IEEE J. Quantum Electron.* **13**:290 (1977).
199. L. Figueroa, T. Holcomb, K. Burghard, D. Bullock, C. Morrison, L. Zinkiewicz, and G. Evans, *IEEE J. Quantum Electron.* **22**:241 (1986).
200. L. J. Mawst, M. E. Givens, C. A. Zmudzinski, M. A. Emanuel, and J. J. Coleman, *IEEE J. Quantum Electron.* **QE-23**:696 (1987).
201. P. S. Zory, A. R. Reisinger, R. G. Walters, L. J. Mawst, C. A. Zmudzinski, M. A. Emanuel, M. E. Givens, and J. J. Coleman, *Appl. Phys. Lett.* **49**:16 (1986).
202. R. G. Walters, P. L. Tihanyi, D. S. Hill, and B. A. Soltz, *Proc. SPIE* **893**:103 (1988).
203. M. S. Zediker, D. J. Krebs, J. L. Levy, R. R. Rice, G. M. Bender, and D. L. Begley, *Proc. SPIE* **893**:21 (1988).
204. C. Krebs and B. Vivian, *Proc. SPIE* **893**:38 (1988).
205. K. Y. Lau and A. Yariv, *Semiconductors and Semimetals Volume 22: Lightwave Communications Technology*, W. T. Tsang (ed.), Academic Press, New York, 1985, pp. 69–151.
206. K. Petermann, *Laser Diode Modulation and Noise*, Kluwer Academic Publ., Dordrecht, The Netherlands, 1988.
207. K. Petermann, *IEEE J. Quantum Electron.* **QE-15**:566 (1979).
208. G. F. Franklin, J. D. Powell, and A. Emami-Naeini, *Feedback Control of Dynamic Systems*, Addison-Wesley Publishing Company, Reading, 1986.
209. T. P. Lee and R. M. Derosier, *Proc. IEEE* **62**:1176 (1974).
210. G. Arnold, P. Russer, and K. Petermann, *Topics in Applied Physics, vol. 39: Semiconductor Devices for Optical Communication*, H. Kressel (ed.), Springer-Verlag, Berlin, 1982, p. 213.
211. K. Y. Lau, N. Bar-Chaim, P. L. Derry, and A. Yariv, *Appl. Phys. Lett.* **51**:69 (1987).
212. K. Y. Lau, P. L. Derry, and A. Yariv, *Appl. Phys. Lett.* **52**:88 (1988).
213. J. E. Bowers, B. R. Hemenway, A. H. Gnauck, and D. P. Wilt, *IEEE J. Quantum Electron.*, **QE-22**:833 (1986).
214. J. E. Bowers, *Solid-State Electron.* **30**:1 (1987).
215. J. Bowers, *Conference on Optical Fiber Communication: Tutorial Sessions*, San Jose, Calif. 1992, p. 233.
216. K. Furuya, Y. Suematsu, and T. Hong, *Appl. Optics* **17**:1949 (1978).
217. D. Wilt, K. Y. Lau, and A. Yariv, *J. Appl. Phys.* **52**:4970 (1981).
218. K. Y. Lau and A. Yariv, *IEEE J. Quantum Electron.* **QE-21**:121 (1985).
219. R. S. Tucker, C. Lin, C. A. Burrus, P. Besomi, and R. J. Nelson, *Electron. Lett.* **20**:393 (1984).
220. W. H. Cheng, A. Appelbaum, R. T. Huang, D. Renner, and K. R. Cioffi, *Proc. SPIE* **1418**:279 (1991).
221. C. B. Su and V. A. Lanzisera, *IEEE J. Quantum Electron.* **QE-22**:1568 (1986).
222. R. Olshansky, W. Powazink, P. Hill, V. Lanzisera, and R. B. Lauer, *Electron. Lett.* **23**:839 (1987).
223. E. Meland, R. Holmstrom, J. Schlafer, R. B. Lauer, and W. Powazink, *Electron. Lett.* **26**:1827 (1990).
224. S. D. Offsey, W. J. Schaff, L. F. Lester, L. F. Eastman, and S. K. McKernan, *IEEE J. Quantum Electron.* **27**:1455 (1991).
225. R. Nagarajan, T. Fukushima, J. E. Bowers, R. S. Geels, and L. A. Coldren, *Appl. Phys. Lett.* **58**:2326 (1991).
226. L. F. Lester, W. J. Schaff, X. J. Song, and L. F. Eastman, *Proc. SPIE* **1634**:127 (1992).
227. K. Y. Lau, C. Harder, and A. Yariv, *IEEE J. Quantum Electron.* **QE-20**:71 (1984).
228. Y. Sakakibara, K. Furuya, K. Utaka, and Y. Suematsu, *Electron. Lett.* **16**:456 (1980).
229. S. Kobayashi, Y. Yamamoto, M. Ito, and T. Kimura, *IEEE J. Quantum Electron.* **QE-18**:582 (1982).
230. P. Vankwikelberge, F. Buytaert, A. Franchois, R. Baets, P. I. Kuindersma, and C. W. Fredrksz, *IEEE J. Quantum Electron.* **25**:2239 (1989).

231. T. Ikegami, "Longitudinal Mode Control in Laser Diodes," *Opto-Electronics Technology and Lightwave Communication Systems*, Van Nostrand Reinhold, New York, 1989, p. 264.
232. T. E. Bell, *IEEE Spectrum* **20**(12):38 (December 1983).
233. M. Osinsky and J. Boos, *IEEE J. Quantum Electron.* **QE-23**:9 (1987).
234. S. Takano, T. Sasaki, H. Yamada, M. Kitomura, and I. Mito, *Electron Lett.* **25**:356 (1989).
235. W. T. Tsang, N. A. Olson, R. A. Linke, and R. A. Logan, *Electron Lett.* **19**:415 (1983).
236. R. Beausoleil, J. A. McGarvey, R. L. Hagman, and C. S. Hong, *Proc. of the CLEO Conference*, Baltimore, Md., 1989.
237. S. Murata, S. Yamazaki, I. Mito, and K. Koboyashi, *Electron Lett.* **22**:1197 (1986).
238. L. Goldberg and J. F. Weller, *Electron Lett.* **22**:858 (1986).
239. J. L. Jewell, A. Scherer, S. L. McCall, Y. H. Lee, S. J. Walker, J. P. Harbison, and L. T. Florez, *Electron. Lett.* **25**:1123 (1989).
240. Z. L. Liao and J. N. Walpole, *Appl. Phys. Lett.* **50**:528 (1987).
241. M. Jansen, J. J. Yang, S. S. Ou, M. Sergeant, L. Mawst, T. J. Roth, D. Botez, and J. Wilcox, *Proc. SPIE* **1582**:94 (1991).
242. G. A. Evans, D. P. Bour, N. W. Carlson, et al., *IEEE J. Quantum Electron.* **27**:1594 (1991).
243. D. Mehuys, D. Welch, R. Parke, R. Waarts, A. Hardy, and D. Scifres, *Proc. SPIE* **1418**:57 (1991).
244. K. Iga, F. Koyama, and S. Kinoshita, *IEEE J. Quantum Electron.* **24**:1845 (1988).
245. M. Y. A. Raja, S. R. J. Brueck, M. Osinski, C. F. Schaus, J. G. McInerney, T. M. Brennan, and B. E. Hammons, *IEEE J. Quantum Electron.* **25**:1500 (1989).
246. K. Tai, L. Yang, Y. H. Wang, J. D. Wynn, and A. Y. Cho, *Appl. Phys. Lett.* **56**:2496 (1990).
247. Y. J. Yang, T. G. Dziura, R. Fernandez, S. C. Wang, G. Du, and S. Wang, *Appl. Phys. Lett.* **58**:1780 (1991).
248. M. Orenstein, A. C. Von Lehmen, C. Chang-Hasnain, N. G. Stoffel, J. P. Harbison, L. T. Florez, E. Clausen, and J. E. Jewell, *Appl. Phys. Lett.* **56**:2384 (1990).
249. R. S. Geels and L. A. Coldren, *Appl. Phys. Lett.* **57**:1605 (1990).
250. C. J. Chang-Hasnain, J. R. Wullert, J. P. Harbison, L. T. Florez, N. G. Stoffel, and M. W. Maeda, *Appl. Phys. Lett.* **58**:31 (1991).
251. A. Von Lehmen, M. Orenstein, W. Chan, C. Chang-Hasnain, J. Wullert, L. Florez, J. Harbison, and N. Stoffel, *Proc. of the CLEO Conference*, Baltimore, Md., 1991, p. 46.
252. Y. H. Lee, J. L. Jewell, B. Tell, K. F. Brown-Goebeler, A. Scherer, J. P. Harbison, and L. T. Florez, *Electron. Lett.* **26**:225 (1990).
253. E. M. Clausen, Jr., A. Von Lehmen, C. Chang-Hasnain, J. P. Harbison, and L. T. Florez, *Techn. Digest Postdeadline Papers, OSA 1990 Annual Meeting*, Boston, Mass., 1990, p. 52.
254. B. Tell, Y. H. Lee, K. F. Brown-Goebeler, J. L. Jewell, R. E. Leibenguth, M. T. Asom, G. Livescu, L. Luther, and V. D. Matterna, *Appl. Phys. Lett.* **57**:1855 (1990).
255. T. P. Lee, *IEEE Proceedings* **79**:253 (1991).
256. M.-C. Amann and W. Thulke, *IEEE J. Selected Areas Comm.* **8**:1169 (1990).
257. K. Kobayashi and I. Mito, *J. Lightwave Tech.* **6**:1623 (1988).
258. T. Saitoh and T. Mukai, *IEEE Global Telecommunications Conference and Exhibition*, San Diego, Calif., 1990, p. 1274.
259. N. A. Olsson, *J. Lightwave Tech.* **7**:1071 (1989).
260. A. F. Mitchell and W. A. Stallard, *IEEE Int. Conference on Communications*, Boston, Mass., 1989, p. 1546.
261. T. Saitoh and T. Mukai, *J. Lightwave Tech.* **6**:1656 (1988).
262. M. J. O'Mahony, *J. Lightwave Tech.* **6**:531 (1988).

ULTRASHORT OPTICAL SOURCES AND APPLICATIONS

Jean-Claude Diels

*Departments of Physics and Electrical Engineering
University of New Mexico
Albuquerque, New Mexico*

Ladan Arissian

*Texas A&M University
College Station Texas, and
National Research Council of Canada
Ottawa, Ontario, Canada*

20.1 INTRODUCTION

It is considered an easy task to control waveforms down to a few cycles with electronic circuits, at frequencies in the megahertz range. Ultrafast optics has seen the development of the same capability at optical frequencies, i.e., in the peta Hertz range. Laser pulses of a few optical cycles (pulse duration of a few femtoseconds) are routinely generated, with a suboptical cycle accuracy. The high power of these ultrashort bursts of electromagnetic radiation have led to new type of high field interactions. Electrons ejected from an atom/molecule by tunnel or multiphoton ionization can be recaptured by the next half optical cycle of opposite sign. The interaction of the returning electron with the atom/molecule is rich of new physics, including high harmonic generation, generation of single attosecond pulses of attosecond pulse trains, scattering of returning electrons by the atom/molecule, etc. Generation, amplification, control, and manipulation of optical pulses is an important starting point for these high field studies.

As compared to fast electronics, ultrafast optical pulses have reached a considerable higher level of accuracy. Pulse trains can be generated, of which the spacing between pulses (of the order of nanoseconds) is a measurable number of optical cycles (one optical cycle being approximately 2 fs in the visible). The frequency spectrum of these pulse trains is a frequency comb, of which each tooth can be an absolute standard with a subhertz accuracy. These frequency combs have numerous applications in metrology and physics—for instance, determining the eventual drift of physical constants, or in astronomy, a considerable improvement in the determination of Doppler shifts of various sources. In addition to the high level of accuracy and control in time and frequency, the femtosecond sources have a remarkable amplitude stability. This stability is the result of nonlinear intracavity losses being minimum for a particular intensity.

This chapter starts with a detailed description of an optical pulse and an optical pulse train. Nonlinear mechanisms are described that can be exploited to control pulse duration, chirp, intensity

of the mode-locked lasers. In particular, a mode-locked laser with two intracavity pulses will be discussed, and its analogy with a quantum mechanical two-level system.

20.2 DESCRIPTION OF OPTICAL PULSES AND PULSE TRAINS

Single Optical Pulse

In this first section we will summarize the essential notations and definitions used throughout the chapter. Ideally, a mode-locked laser emits a continuous train of identical ultrashort pulses. To this infinite series of identical pulses corresponds, in the frequency domain, a finite (but large) number of equally spaced modes, generally referred to as a *frequency comb*. Inside the laser typically, only one pulse circulates. The shape of an intracavity pulse results from a steady-state equilibrium between various mechanisms of pulse stretching (saturable gain, dispersion), compression (saturable absorption, combination of self-phase modulation, and dispersion), amplification, and losses.

The pulse is characterized by measurable quantities which can be directly related to the electric field. A complex representation of the field amplitude is particularly convenient in dealing with propagation problems of electromagnetic pulses.

The complex spectrum of the pulse $\tilde{E}(\Omega)^*$ is defined by taking the complex Fourier transform \mathcal{F} of the real electric field $E(t) = \varepsilon(t)\cos[\omega t + \varphi(t)]$:

$$\tilde{E}(\Omega) = \mathcal{F}\{E(t)\} = \int_{-\infty}^{\infty} E(t)e^{-i\Omega t} dt = \left| \tilde{E}(\Omega) \right| e^{i\Phi(\Omega)} = \tilde{\varepsilon}(\Omega - \omega) = \tilde{\varepsilon}(\Delta\Omega) \quad (1)$$

In the definition (1), $|\tilde{E}(\Omega)|$ denotes the spectral amplitude and $\Phi(\Omega)$ is the spectral phase. Since $E(t)$ is a real function, its Fourier transform is symmetric, and its negative frequency part can be considered as redundant information. We will therefore choose to represent the light pulse by either the positive frequency function $\tilde{E}(\Omega) = E(\Omega)e^{i\Phi(\Omega)}$ (defined as being equal to zero for $\Omega < 0$) or its complex inverse Fourier transform in the time domain

$$\tilde{E}(t) = \frac{1}{2\pi} \int_0^{\infty} \tilde{E}(\Omega)e^{i\Omega t} d\Omega = \frac{1}{2} \tilde{\varepsilon}(t)e^{i\omega t} = \frac{1}{2} \varepsilon(t)e^{i[\omega t + \varphi(t)]} \quad (2)$$

The relation with the real physical measurable field $E(t)$ is

$$E(t) = \tilde{E}(t) + \text{c.c.} = \varepsilon(t)\cos[\omega t + \varphi(t)] \quad (3)$$

The latter part of Eq. (2) defines a pulse envelope function $\varepsilon(t)$, a carrier frequency ω and a phase $\varphi(t)$. The decomposition is somewhat arbitrary, since the instantaneous frequency is given by

$$\omega(t) = \omega + \frac{d}{dt}\varphi(t) \quad (4)$$

In general, the carrier frequency ω will be chosen such that the average contribution from the phase factor $\varphi(t)$ is zero:

$$\langle \varphi(t) \rangle = \frac{\int_{-\infty}^{\infty} \varepsilon(t)^2 \dot{\varphi}(t) dt}{\int_{-\infty}^{\infty} \varepsilon(t)^2 dt} = 0 \quad (5)$$

*Complex quantities related to the field will be represented with a tilde.

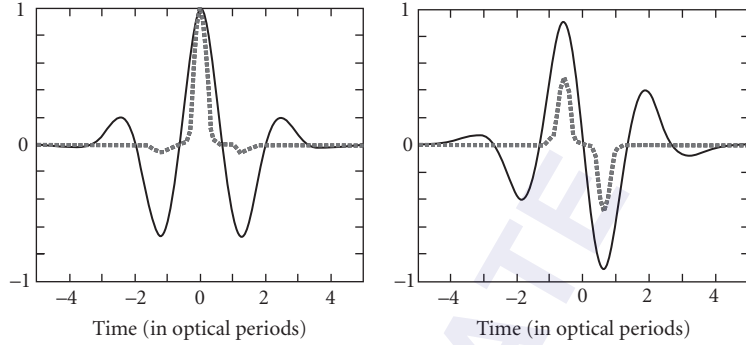


FIGURE 1 Comparison of a two-cycle pulse with $\varphi_e = 0$ (left) and $\varphi_e = \pi/2$ (right). The solid line traces the instantaneous electric field normalized to the peak value of its envelope, as a function of time in units of the optical period. The dotted lines correspond to the seventh power of the electric field, which would be driving a seven photon process.

For pulses of a few optical cycles, the variation of the phase factor can often be neglected across the pulse, and $\varphi(t) = \varphi_e$ is constant. Even for a single pulse, the phase factor φ_e is of practical significance, when a nonlinear phenomena traces the electric field under the envelope of a pulse of only a few cycle duration. If the phase φ_e is zero, the time dependence of the electric field is symmetric, with a peak in the center at $t = 0$, larger than the two opposite maxima at $t = \pm T/2$. If the phase $\varphi_e = \pi/2$, the time dependence of the electric field is antisymmetric, with equal opposite extrema at $t = \pm \pi/4$. These two pulses can give a different response in highly nonlinear phenomena. Let us consider for instance the shortest pulse that can be generated at 800 nm, which has a full width at half maximum (FWHM) of the intensity of 2.5 fs. Its complex electric field envelope can be written as $\tilde{\mathcal{E}}(t) = (\varepsilon_0/2) \exp[-(t/2T)^2 + i\varphi_e]$, which corresponds to the real electric field $E(t) = \varepsilon_0 \exp[-(t/2T)^2] \cos[2\pi(t/2T) + \varphi_e]$ which is plotted as a solid line in Fig. 1 for $\varphi_e = \pi/2$ (left) and $\varphi_e = \pi/2$ (right). If we consider that this pulse is used to excite a seven-photon process (for instance a seven-photon ionization), the driving function for that process is the seventh power of the field, which is plotted as a dotted line in Fig. 1. One can see that the different values of φ_e make a significant difference on how the process is driven. For $\varphi_e = 0$, the excitation is a single spike, as close approximation as practical to a δ -function. In the case of $\varphi_e = \pi/2$ (right), the excitation consists in a succession of positive and negative spikes.

Train of Pulses

The “ideal” mode-locked laser emits a train of identical pulses, at equal time interval. The period of the pulse train is τ_{RT} , defined as the separation between two successive envelopes. In the particular case that the pulse separation is an integer number of optical cycles $\tau_{RT} = NT = N/\nu$ (T being the light period and $\nu = \omega/(2\pi)$ the optical frequency) the successive pulses are identical. This will generally not be the case, and there will be a phase shift $\varphi_p = \omega\tau_{RT} \neq 2N\pi$ between successive pulses. The complex electric field of the total pulse train \tilde{E}_{pt} is

$$\tilde{E}_{pt}(t) = e^{i\omega t} [\tilde{\mathcal{E}}(t) + \tilde{\mathcal{E}}(t - \tau_{RT}) e^{i\varphi_p} + \tilde{\mathcal{E}}(t - 2\tau_{RT}) e^{2i\varphi_p} + \dots] \quad (6)$$

where $\tilde{\epsilon}(t) = \epsilon(t)e^{i\varphi_c}$ is the electric field of one particular pulse. The n th pulse has the phase factor $\exp[i(\varphi_e + n\varphi_p)]$, different from the previous and next pulse. To the change in phase between successive pulses φ_p , corresponds a frequency:

$$f_0 = \frac{1}{2\pi} \frac{\varphi_p}{\tau_{RT}} \quad (7)$$

This frequency is called the *carrier to envelope offset*. The *carrier to envelope offset* is an important parameter of pulse trains, where the change in phase from pulse to pulse is a measurable quantity, independent of the duration of the individual pulse in the train.

One can “idealize” to the extreme the concept of a pulse train, by considering an infinite train of δ -functions, equally spaced by the period of the train τ_{RT} , as shown in Fig. 2a. The Fourier transform of this ideal pulse train shown in Fig. 2b is an identical picture in the frequency domain: a comb of infinite extent (because the pulses were δ -function in time), with δ -function teeth (because of the infinite extent of the train).

Since the comb extends to infinity, there is no particular tooth that can be called an average frequency. Each mode ν_m of index m carries the same weight, and corresponds in the time domain to an infinite sine wave, which is a particular term of a Fourier series representation of δ -function. The first tooth at frequency $\nu_0 = f_0$ represents the carrier to envelope offset defined above. The corresponding carrier to envelope phase φ_e defined previously can be identified in the time domain, even with a train of δ -functions. The harmonic wave corresponding to the mode ν_2 is shown in Fig. 2a,* and the phase φ_e is identified as the phase at which each δ -function crosses the harmonic field. In the sketch of Fig. 2a, $\varphi_e = 0$ for the first pulse, and φ_p is then the carrier to envelope phase φ_e of the second pulse as indicated in the figure.

A somewhat more mundane train of pulses of finite duration τ^\dagger is sketched in Fig. 3a. In the frequency domain (Fig. 3b), the infinite pulse train is represented by a finite frequency comb. The

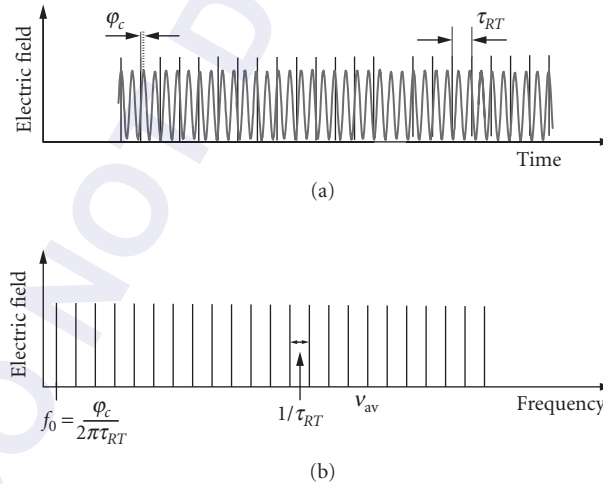


FIGURE 2 Idealized infinite train of δ -function pulses (a), and its Fourier transform (b). In (a), the carrier to envelope phase φ_e of the first pulse is assumed to be zero.

*This harmonic wave sketched is associated with ν_2 because there are two periods between pulses. In the Fourier spectrum of a train of δ -functions, any mode ν_n can be chosen as being the “average frequency”.

†When not otherwise specified, the pulse duration will be the full width at half maximum (FWHM) of the pulse intensity profile.

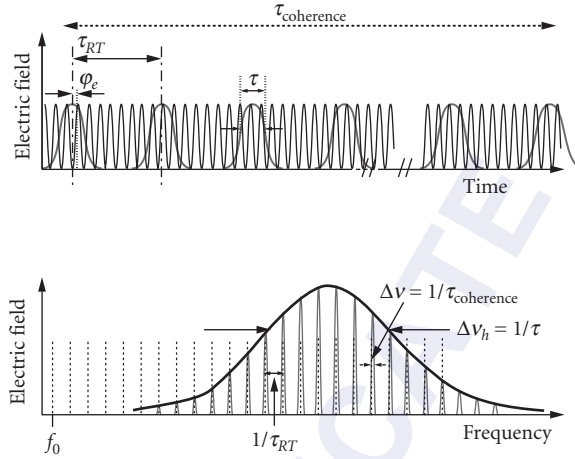


FIGURE 3 (a) Train of pulses of finite duration τ . The successive pulse envelope repeats every τ_{RT} . Within the coherence time of the train $T_{coherence}$, it is the same carrier at the optical frequency that is modulated by the successive envelopes. (b) The Fourier transform of the pulse train shown in (a).

envelope of the comb is the Fourier transform of the envelope of a single pulse of the train, thus of extension $\approx 1/\tau$. The teeth of the frequency comb are no longer δ -functions, but sharp peaks of width $1/\tau$, where τ is the coherence time of the pulse train. The carrier to envelope phase φ_p is indicated for a pulse of the time sequence in Fig. 3a. Note that this phase is changing from one pulse to the next. The rate of change $\varphi_p/(2\pi\tau_{RT})$ is the frequency f_0 , which is indicated in the frequency picture by the lowest frequency tooth of the extension (dashed lines in Fig. 3b).

The angular frequency ω_m of the m th mode of the comb is given by

$$\omega_m = 2\pi f_0 + m \frac{2\pi}{\tau_{RT}} \quad (8)$$

In the case of a train of pulses of finite duration, the frequency f_0 is no longer a real tooth of the comb, but the first mode of an extension of the frequency comb beyond the pulse bandwidth as shown by the dotted line in (Fig. 3b).

It can be seen from this definition that f_0 is indeed the change of phase per round-trip between the envelope and the carrier. By definition of τ_{RT} , the pulse envelope peaks exactly at the same locations after one round-trip. With respect to this envelope, the shift of phase of the mode m is obtained by multiplying Eq. (8) by τ_{RT} :

$$\omega_m \tau_{RT} = 2\pi f_0 \tau_{RT} + 2m\pi \quad (9)$$

which, after substitution of the definition of f_0 Eq. (7), is indeed the phase φ_p defined earlier.

Soliton Solution and Steady-State Pulse Train

As mentioned in the introduction, if a laser is to generate a pulse of well defined duration and shape, there has to be compression and broadening mechanisms that balance each other, and lead to

a stable pulse. The mechanisms that lead to an emergence of a pulse out of noise in a laser cavity are usually dissipative, i.e., the pulse that emerges dissipates a minimum amount of energy by nonlinear loss mechanisms and extracts the maximum amount of gain from the active laser medium. We will consider in this section only nondissipative interaction that plays a dominant role for the stable formation of the shortest pulses. First we consider the evolution of a single pulse as it propagates through a cavity. Next we will study the formation of a pulse train with a similar nondissipative model.

Evolution of a Single Pulse in an “Ideal” Cavity When a laser is in continuous operation, the cavity gain and losses are in equilibrium. In the case of femtosecond mode-locked lasers, the major pulse-shaping mechanism is a combination of self-phase modulation and dispersion at each round-trip. The self-phase modulation results from a nonlinear index of refraction $n_2 I$ (I being the intensity in W/cm^2 , n_2 being the nonlinear index in cm^2/W of a nonlinear element of length ℓ in the cavity). The dispersion k'' results from the frequency dependence of the average index of refraction n_{av} , defined previously, and is characterized by the second derivative of a cavity averaged k vector with respect to frequency. In what follows, for simplicity, we will neglect higher-order terms in Kerr effect and in dispersion. The evolution of a pulse in the mode-locked laser cavity can be considered as a propagation (of a nondiffracting beam) through an infinite lossless medium, with a positive Kerr nonlinearity ($n_2 > 0$) and a negative dispersion (as can be introduced with intracavity prisms¹ or chirped mirrors² in the cavity). The pulse evolution generally converges toward a steady-state solution, designated as “solitons,” which can be explained as follows. The nonlinearity is responsible for spectral broadening and up-chirp. Because of the anomalous dispersion, $k'' < 0$, the high-frequency components produced in the trailing part of the pulse, travel faster than the low-frequency components of the pulse leading edge. Therefore, the tendency of pulse broadening owing to the exclusive action of group velocity dispersion can be counterbalanced. To determine the approximate parameters of that solution, let us assume a Gaussian pulse $\mathcal{E}(t) = \varepsilon_0 \exp[-(t/\tau_{G0})^2]$, and let us state that the chirp produced in the pulse center by the nonlinearity and the dispersion are of equal magnitude (but of opposite sign). Under this equilibrium condition the pulse circulates in the cavity without developing a frequency modulation and spectral broadening. The effect of group velocity dispersion is to create a pulse broadening and a down-chirp (in a medium of negative dispersion). The change (per round-trip) of the second derivative of the phase versus time, at the center of the pulse, is given, in a first-order approximation, by Ref. 3:

$$\Delta \left(\frac{\partial^2 \varphi(t)}{\partial t^2} \Big|_{(t=0)} \right) = \frac{4k''_{\text{av}}}{\tau_{G0}^4} P \quad (10)$$

where $k''_{\text{av}} = d^2 k_{\text{av}}/d\Omega^2$ is the second-order dispersion averaged over the cavity of perimeter P , which is < 0 for an optical element with negative dispersion. Assuming that the cavity contains an element with a nonlinear index $n = n_0 + n_2 I$ of length ℓ_{Kerr} , the phase induced by self-phase modulation, near the center of the Gaussian pulse, is

$$\Delta \varphi(t) = -k_{\text{NL}} \cdot \ell_{\text{Kerr}} \Big|_{(t=0)} = \frac{2\pi n_2}{\lambda} \ell_{\text{Kerr}} I \approx \frac{4\pi n_2 I_0 \ell_{\text{Kerr}}}{\lambda} \frac{t}{\tau_{G0}^2} \quad (11)$$

where we have used a quadratic approximation for the Gaussian near $t = 0$. Taking the second derivative yields the chirp induced by phase modulation at the pulse center:

$$\Delta \left(\frac{\partial^2 \varphi(t)}{\partial t^2} \right) \approx \frac{8\pi n_2 I_0 \ell_{\text{Kerr}}}{\lambda \tau_{G0}^2} \quad (12)$$

The peak intensity of the pulse (at $t = 0$) is $I_0 = \epsilon_0^2 / 2\eta$; $\eta = \sqrt{\mu_0 / \epsilon}$ being the characteristic impedance of the medium. Expressing that the chirps induced by phase modulation [Eq. (12)] and dispersion [Eq. (10)] should cancel each other leads to

$$I_0 \tau_{G0}^2 = -\frac{\lambda k''_{av}}{2\pi n_2} \frac{P}{\ell_{Kerr}} \quad (13)$$

This expression leads to a value for the pulse duration given by

$$\tau_{G0}^2 = \frac{\lambda |k''_{av}|}{2\pi n_2 I_0} \frac{P}{\ell_{Kerr}} \quad (14)$$

In a Ti:sapphire laser, $n_2 = 10.5 \cdot 10^{16} \text{ cm}^2/\text{W}$; the crystal length is typically $\ell = 4 \text{ mm}$, and a well-designed laser can produce a train of 10 fs pulses with an intracavity average power of the order of 10 W. These numbers can be used in Eq. (14) to determine the negative cavity dispersion k''_{av} required for stable laser operation.

How Unequally Spaced Modes Lead to a Perfect Frequency Comb The physical reason for a non-zero carrier to envelope offset f_0 is the dispersion of the laser cavity in which a pulse is circulating. The components of the laser cavity impose different group velocities on the pulse. We can define an average group velocity $v_g = P / \tau_{RT}$ of the pulse envelope, where P represents twice the length of a linear cavity, or the perimeter of a ring cavity. That group velocity is different from the phase velocity c/n_{av} (n_{av} being the linear index of refraction averaged over the laser components). The two quantities are related by

$$\frac{1}{v_g} = \frac{n_{av}}{c} + \frac{\omega}{c} \frac{dn_{av}}{d\Omega} \bigg|_{\omega} \quad (15)$$

Note that these quantities n_{av} and v_g are function of the spectral frequency of the pulse. The “ideal mode-locked laser” considered in this section already poses a conceptual dilemma. Mode-locking is generally described as putting the modes of a laser cavity in phase. If the cavity has dispersion, we have seen that the mode-comb issued from the laser does not start at zero frequency but with a frequency offset f_0 . Keeping in mind that a cavity with dispersion has unequally spaced modes, is contradictory to the fact that the frequency comb has rigorously equally spaced teeth.* To resolve this apparent contradiction, we will look at the pulse train formation, and discuss how an initially irregular set of modes can lead to a perfect frequency comb.

As shown above, a minimum negative cavity dispersion k''_{av} is required for stable mode-locked operation. Such a cavity dispersion implies that the index of refraction n_{av} is frequency (wavelength) dependent, hence the spacing of the cavity modes $c/[n_{av}(\Omega)P]$ varies across the pulse spectrum.

The laser is modeled by a circulating pulse, which enters a Kerr medium of thickness ℓ , resulting in phase modulation at each passage, and a medium that represents the linear dispersive properties of the cavity. We will assume that the balance of gain and losses maintains a constant Gaussian shape for the envelope of the circulating pulse. At each passage through the cavity, the phase of the pulse is modified in the time domain through the Kerr effect, and in the frequency domain through dispersion. We consider first the modulation in the time domain:

$$\varphi(t) = -k_{NL} \ell_{Kerr} = -\frac{2\pi n_2 \ell_{Kerr}}{\lambda} I_0 e^{-2(t/\tau_G)^2} \quad (16)$$

*A fact that has been verified experimentally with millihertz accuracy.⁴

where τ_G is the $1/e$ half-width of the pulse electric field envelope (the FWHM of the intensity is $\tau_p = \sqrt{2\ln 2}\tau_G$). Ignoring at this point the influence of dispersion (which will be introduced after Fourier transformation into the frequency domain), the pulse *train* issued from the laser can be represented by

$$\sum_{q=0}^{\infty} \mathcal{E}(t - \tau_q) e^{iq\varphi(t - \tau_q)} e^{i\omega t} \quad (17)$$

where τ_q is the time of arrival of the center of gravity of the successive pulses. At this point τ_q is not set to any value. It is assumed here that at $t = 0$, the first pulse is unmodulated. Using a parabolic approximation for the Gaussian intensity profile, the time-dependent phase is

$$\varphi(t - \tau_q) \approx \frac{4\pi n_2 I_0 \ell_{\text{Kerr}}}{\lambda} \left(\frac{t - \tau_q}{\tau_G} \right)^2 = a \left(\frac{t - \tau_q}{\tau_G} \right)^2 \quad (18)$$

The Fourier transform of the pulse train given by Eq. (17) is

$$\mathcal{E}(\Delta\Omega) \left[\sum_{q=0}^{\infty} e^{i\Delta\Omega\tau_q} e^{iq\Delta\Omega^2\tau_k^2} \right] \quad (19)$$

where

$$\begin{aligned} \Delta\Omega &= \Omega - \omega \\ \mathcal{E}(\Delta\Omega) &= \frac{\epsilon_0 \sqrt{\pi} \tau_G}{\sqrt[4]{1+a^2}} \exp\left\{ -\frac{\Delta\Omega^2 \tau_G^2}{4(1+a^2)} \right\} \\ \tau_k^2 &= \frac{a \tau_G^2}{4(1+a^2)} \end{aligned} \quad (20)$$

The width of the Gaussian pulse spectrum, broadened by the Kerr effect, is the inverse of the characteristic time τ_k . Let us now take dispersion into account. The operation representing the dispersion of the cavity is a product of the spectral field by $\exp[-ik_{\text{av}}(\Delta\Omega)P]$, where $-k_{\text{av}}(\Delta\Omega)P^*$ is the phase change per round-trip. The combined Kerr effect and dispersion, in the frequency domain, leads to the output spectral field:

$$\mathcal{E}_{\text{out}}(\Delta\Omega) = \mathcal{E}(\Delta\Omega) \left[\sum_{q=0}^{\infty} e^{i\Delta\Omega\tau_q} e^{iq\Delta\Omega^2\tau_k^2} e^{-iqk_{\text{av}}(\Delta\Omega)P} \right] \quad (21)$$

Expanding the wave vector $k_{\text{av}}(\Delta\Omega)$ in series, to second order:

$$\begin{aligned} k_{\text{av}}(\Delta\Omega)P &= k_{\text{av}}(\Delta\Omega=0)P + \Delta\Omega k'_{\text{av}}P + \frac{\Delta\Omega^2}{2} k''_{\text{av}}P \\ &= k_{\text{av}}(\Delta\Omega=0)P + \Delta\Omega \tau_{RT} + \frac{k''_{\text{av}}P}{2} \Delta\Omega^2 \end{aligned} \quad (22)$$

*In the argument of k_{av} , the light frequency ω is taken as origin ($\Delta\Omega = 0$) of the frequency scale.

where the derivatives k'_{av} and k''_{av} are calculated at the light frequency $\omega(\Delta\Omega=0)$. Note that $k'_{av}=1/v_g=\tau_{RT}/P$ [cf. Eq. (15)] are material properties independent of the index q , as is the cavity perimeter P . The modes of the cavity are not equally spaced. The parameter k'' characterizes the departure from equal spacing. Substituting (22) in Eq. (21),

$$\mathcal{E}_{out}(\Delta\Omega) = \mathcal{E}(\Delta\Omega) \left[\sum_{q=0}^{\infty} e^{i\Delta\Omega(\tau_q - q\tau_{RT})} e^{iq\Delta\Omega^2(\tau_k^2 - k''_{av}P/2)} \right] \quad (23)$$

The conditions

$$\tau_k^2 = -\frac{k''_{av}P}{2} \quad (24)$$

$$\tau_q = (q+1)\tau_{RT} \quad (25)$$

leads to modes that are exactly equally spaced. The inverse Fourier transform of the frequency comb becomes then

$$\tilde{\mathcal{E}}_{out}(t) = \mathcal{E}(t) + \mathcal{E}(t - \tau_{RT})e^{-ik_{av0}P} + \mathcal{E}(t - 2\tau_{RT})e^{-2k_{av0}P} + \dots \quad (26)$$

This last equation corresponds indeed to the description of the ideal frequency comb, with equally spaced pulses in time and frequency, and a carrier to envelope phase shift of $\varphi_p = -k_{av0}P$. In the case of small Kerr modulation, $a \ll 1$, it can easily be verified that the condition in Eq. (24) is identical to the soliton Eq. (14). Indeed, substituting

$$\tau_k^2 = \frac{a\tau_G^2}{4(1+a^2)} \approx \frac{4\pi n_2 I_0 \ell_{Kerr} \tau_G^2}{4\lambda} = \frac{k''_{av}P}{2} \quad (27)$$

which is indeed equivalent to Eq. (14). One can thus conclude that the mechanism that leads to an equal spacing for the teeth of the frequency comb emitted by the laser is the same Kerr effect responsible for creating maximum intracavity pulse compression.

20.3 PULSE EVOLUTION TOWARD STEADY STATE

A Simple Model

In the previous section we have considered the dispersive mechanisms that ultimately give the final shape in amplitude and phase to the steady-state pulse. This mechanism dominates in the sub-picosecond regime, where dissipative mechanisms have reached equilibrium. Other elements play a decisive role in initiating the mode-locking, which are usually referred to as the passive mode-locking elements. The latter can most often be represented by intensity-dependent intracavity loss. Larger losses at low intensity imply that the laser has less gain—and may be below threshold—for low-intensity continuous wave (cw) radiation than for pulses with higher peak intensity. This leads to the emergence of a pulse out of the amplified spontaneous emission noise of the laser. Rather than concentrating on the primary process of formation of a precursor of a pulse from random noise, let us follow the evolution of the pulse from its birth from noise until it has blossomed into a fully shaped stable laser pulse. In this intermediate stage of the evolution toward steady state, the main shaping elements are dissipative, as opposed to the purely dispersive interaction considered in the previous section. We will look for simple evolution equations for the pulse energy $W = \int_{-\infty}^{\infty} I(t) dt$, with $I(t)$ being the pulse intensity. The element responsible for saturable losses (gain) should have

typically a *linear* loss (gain) factor at low energies, and a *constant* loss (gain) at higher energies. We thus have, at low energies: $dW/dz = \mp \alpha W$. At large energies $W \gg W_s$: $dW/dz = \mp W_s$ where W_s is the saturation energy for the chosen geometry. The simplest differential equation to combine these two limits is

$$\frac{dW}{dz} = \alpha_g W_{sg} \left[1 - e^{W/W_{sg}} \right] \quad (28)$$

Equation (28) is written for a medium with a linear gain α_g and a saturation energy W_{sg} . It can be integrated to yield the energy W_2 at the end of the amplifier of thickness d_g , as a function of the input energy W_1 :

$$W_2 = G(W_2, W_1)W_1 = W_{sg} \ln \left\{ 1 - e^{\alpha_g d_g} (1 - e^{W_1/W_{sg}}) \right\} \quad (29)$$

A similar equation applies to the saturable absorber, with a negative absorption coefficient $-\alpha_a$ and a smaller saturation energy W_{sa} :

$$W_2 = A(W_2, W_1)W_1 = W_{sa} \ln \left\{ 1 - e^{\alpha_a d_a} (1 - e^{W_1/W_{sa}}) \right\} \quad (30)$$

The dominant linear loss element is the output coupler, with (intensity) reflectivity r . The transfer function for that element is simply

$$W_2 = L(W_2, W_1)W_1 = rW_1 \quad (31)$$

and the energy of the output pulse is $(1-r)W_1$. The evolution of the pulse energy in a single round-trip can be simply calculated from the product of all three transfer functions given by Eqs. (29), (30), and (31). For instance, if we consider a ring laser with the sequence: mirror, gain, and absorber, the pulse energy W_4 after the absorber is given by the product $A(W_4, W_3)G(W_3, W_2)L(W_2, W_1)$. One can also express the relation between the energy W_4 and the pulse energy W_1 before the output mirror by the algebraic relation:

$$1 + a[e^{W_4/W_{sa}} - 1] = \left\{ 1 + g[e^{rW_1/W_{sg}} - 1] \right\}^{W_{sg}/W_{sa}} \quad (32)$$

where $a = \exp\{-\alpha_a d_a\}$ is the *linear* small-signal attenuation of the passive element and $g = \exp\{\alpha_g d_g\}$ is the *linear* small-signal amplification.

High-Gain Oscillators

Unlike laser amplifiers, where it is desirable to use a gain medium with as high a saturation energy density as possible, mode-locked oscillators will often use high-gain laser media. These are opposite requirements: the larger the amplification cross section σ_g , the larger the gain $\alpha_g = \Delta N \sigma_g$, and the smaller the saturation energy density $W_s = \hbar \omega / (2\sigma_g)$. Both numbers a and g can be large, and the reflectivity of the output coupler r can be even lower than 50 percent. Examples are dye lasers, semiconductor lasers with tapered amplifiers, and to a smaller extent the Ti:sapphire laser. As a result, the order of the elements matters in the design of the laser, and in its performances. To illustrate this point, let us assume that the passive element is totally saturated in normal operation. In full saturation, the input energy $W(0)$ is related to the output energy $W(d)$ by

$$W(d) = W(0) - \alpha_a d_a W_{sa} \quad (33)$$

Given an initial energy W_1 , the energy W_4 for a single passage through three different sequences of the same elements is given below. For the sequence mirror-absorber-gain

$$e^{W_4/W_{sg}} = 1 - g \left[1 - e^{(rW_1 - \alpha_a d_a W_{sa})/W_{sg}} \right] \quad (34)$$

For the sequence absorber-mirror-gain

$$e^{W_4/W_{sg}} = 1 - g \left[1 - e^{r(W_1 - \alpha_a d_a W_{sa})/W_{sg}} \right] \quad (35)$$

For the sequence absorber-gain-mirror

$$e^{W_4/W_{sg}} = \left\{ 1 - g \left[1 - e^{(W_1 - \alpha_a d_a W_{sa})/W_{sg}} \right] \right\}^{1/r} \quad (36)$$

Let us take a numerical example for a high-gain system such as the flash-lamp pumped Nd:glass laser, with $W_{sa}/W_{sg} = 0.1$, $\alpha_a d_a = 1$, $\alpha_g d_g = 1.5$ and output coupling of $r = 0.5$. For an initial energy $W_1/W_{sg} = 0.5$, we find for the sequence absorber-mirror-gain $W_4/W_{sg} = 0.689$, and for the sequence absorber-gain-mirror $W_4/W_{sg} = 0.582$. The order of the elements, the relative saturation of the gain to that of the passive element, as well as the output coupling influence the stability and output power of the laser, as shown in Ref. 5. The evolution of the pulse in the cavity can be calculated by repeated applications of products of operations such as $A(W_4, W_3)G(W_3, W_2)L(W_2, W_1)$ for the sequence (passive element, gain, output coupler), starting from a minimum value of W_1 above threshold for pulsed operation, and recycling at each step the value of W_4 as the new input energy W_1 . Figure 4 shows the growth of intracavity pulse energy as a function of the round-trip index j , for different orders of the elements. The initial pulse energy is 1 percent of the saturation energy W_{sg} in the gain medium. The saturation energy and optical thickness of the absorbing medium are, respectively, $W_{sa} = 0.8W_{sg}$ and $\alpha_a d_a = 1.2$. The linear gain is $\alpha_g d_g = 1.5$ and the output coupling $r = 0.8$.

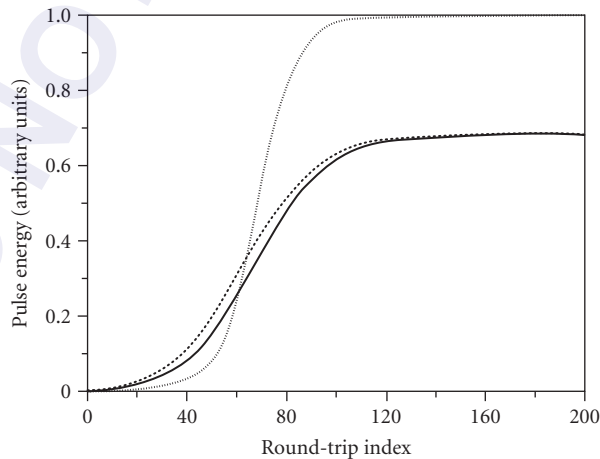


FIGURE 4 Intracavity pulse energy versus round-trip index j . The solid, dashed, and dotted lines correspond to the sequence m (irror)- a (bsorber)- g (ain), a - g - m and a - m - g , respectively.

The main point of this exercise is that the order of the elements is important, a fact confirmed by measurement on moderate-gain lasers such as $\text{Ti:Al}_2\text{O}_3$. Such dependence on the order of the elements indicates that analytical theories based on the approximation of infinitesimal change per element and per cavity round-trip are not quite adequate. Numerical codes have been developed that attempt to include all physical phenomena affecting pulse shape and duration. Unfortunately, the sheer number of these mechanisms makes it difficult to reach a physical understanding of the pulse-generation process, or even identify the essential parameters. Therefore, the most popular approach is to construct a simplified analytical model on a selected mechanism.

20.4 COUPLING CIRCULATING PULSES INSIDE A CAVITY

Pulse Train Interferometry

A pulse train combines the temporal resolution of a single pulse, and the spectral resolution of a cw beam. It is therefore not surprising that new interferometric technique can be developed exploiting these properties.

Because the ratio of pulse duration to the period of the train is generally less than $1:10^5$, interferences of pulse trains of different repetition rates will not be considered here. We will instead focus on a situation where two pulse trains of identical repetition rate, but different carrier to envelope phase, are made to interfere. It will be shown in the next section how pulse trains of identical repetition rate are generated. The experimental arrangement for pulse train interferometry is depicted in Fig. 5. The two pulse trains are combined by a beam splitter, and their relative delay adjusted in an optical delay line, in order to have superposition of the pulse envelopes on the detector. If the carrier to envelope phase is identical for both pulse trains, the detector will simply register a constant signal, with an amplitude dependent on the relative phase of the two trains at the detector.

If instead the two pulse trains have a different carrier to envelope phase, successive pulses will interfere differently. The envelope of the interfering pulse trains, as seen by the detector, will be modulated at the frequency $f_{01} - f_{02}$, where f_{01} and f_{02} are the carrier to envelope offsets of either pulse trains.

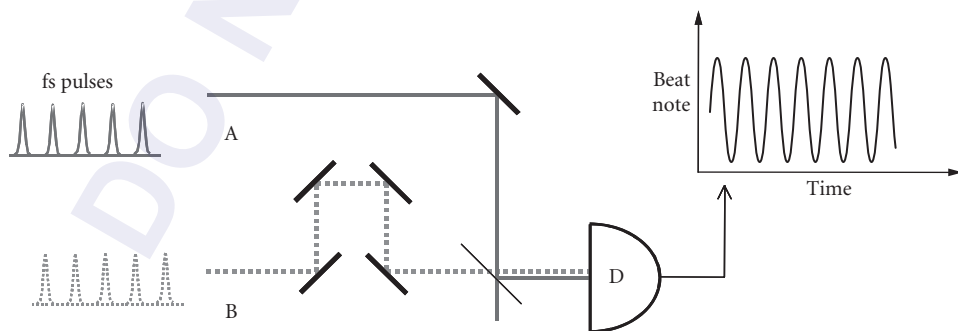


FIGURE 5 Interference of two pulse trains of the same repetition rate. An optical delay line is required to ensure temporal overlap of the pulses in either train.

Interwoven Pulse Trains Generated by Two Intracavity Pulses

We have seen the interference of pulse trains. These can be generated by a laser cavity. Inside the laser, there are then two or more pulses circulating in the cavity. The interactions of these two pulses will determine the relative properties of the pulse trains. The common picture of the mode locked laser is that of a resonator cavity with gain in which a single pulse circulates. It is interesting to consider, both from the point of view of fundamental understanding of the laser operation and applications to sensors, the situation where several pulses circulate in the cavity. Mode-locking with multiple pulses/round-trip is sometimes referred to as “harmonic mode-locking” 6–10. Techniques of harmonic mode-locking have been developed for telecommunications, where a pulse rate of over a GHz is desirable. The fundamental clock remains the round-trip time, which may typically be of the order of 100 MHz. If there are m pulses in the cavity, there will be m values of interpulse delay, which, in the frequency domain, will mean some splitting of the modes, which, is small, may be seen as a slight broadening of the tooth of the comb.

In the following we will consider only the case of two pulses circulating in the cavity, separated exactly by half the cavity round-trip time. Such a situation is encountered in bidirectional ring lasers, but also in linear cavities. Of particular interest is to determine the type of coupling that may or may not exist between the two pulses, and the resulting correlation between the two pulse trains. As will be the case in most sensor applications,¹¹ we will assume in the following that the two pulses experience a relative shift in phase δ at each round-trip.

Locking Two Pulse Trains by Backscattering

Whether in a ring or linear laser, the two circulating pulses will meet at two points of the cavity. Unless the meeting point is in vacuum, there will be some coupling introduced by the medium in which the pulses meet. The most common case is that of a medium with random scattering. Using the plane wave description of Eq. (2), the backscattering component $\tilde{r}_{ij} = r \exp(\theta_{ij})$ of scattering will couple the pulse with field envelope \tilde{e}_j into the pulse with field envelope \tilde{e}_i :

$$\begin{aligned}\frac{\partial \tilde{e}_1}{\partial t} &= i \frac{\delta}{2\tau_{RT}} \tilde{e}_1 + \frac{1}{\tau_{RT}} \tilde{r}_{12} \tilde{e}_2 \\ \frac{\partial \tilde{e}_2}{\partial t} &= \frac{1}{\tau_{RT}} \tilde{r}_{21} \tilde{e}_1 - i \frac{\delta}{2\tau_{RT}} \tilde{e}_2\end{aligned}\quad (37)$$

Of particular interest here is the impact of the coupling on the phase of the two fields, since the balance of saturable gain and losses will in general restore a steady-state value of the pulses energy. Expressing the fields in terms of amplitude and phase as in Eq. (2) in the system of Eqs. (37), and taking the difference of the imaginary parts, yields:

$$\frac{\partial(\varphi_2 - \varphi_1)}{\partial t} = \frac{\partial\psi}{\partial t} = \frac{\delta}{\tau_{RT}} + \frac{r}{\varepsilon_1 \varepsilon_2 \tau_{RT}} [\varepsilon_1^2 \sin(\theta_{21} - \psi) - \varepsilon_2^2 \sin(\theta_{12} + \psi)] \quad (38)$$

In the absence of coupling ($r=0$), the carrier frequency of the two pulses would differ by $\dot{\psi} = \delta/\tau_{RT}$, a frequency difference that can easily be detected by beating the two output pulse trains of the laser (corresponding to either pulse in the cavity) against each other on a detector.

Distributed Backscattering In presence of a sufficient coupling $r \neq 0$, there is generally a solution $\partial\psi/\partial t$ to Eq. (38), such that the two circulating pulses are identical, only differing by a phase factor. If we take for instance the particular case of $\theta_{12} = \theta_{21} = 0$, the constant phase difference is ψ_0 given by

$$\sin \psi_0 = \frac{2\varepsilon_1\varepsilon_2}{\varepsilon_1^2 + \varepsilon_2^2} \frac{\delta}{2r} \quad (39)$$

Any backscattering such that

$$r \geq \frac{\delta\varepsilon_1\varepsilon_2}{\varepsilon_1^2 + \varepsilon_2^2} \quad (40)$$

will lock the carrier frequency of the two waves to each other. This implies that the mode frequencies, the repetition rates, and the CEO of the two pulse trains are identical.

Interface Coupling A reciprocal backscattering, where $\theta_{12} = \theta_{21}$ is the norm when dealing with distributed scattering of a solid, liquid, or gaseous medium.^{12,13} The situation is different however in a short pulse laser, where the meeting points of the two pulses are localized rather than being distributed over the whole length of the laser resonator. In the case of the mode-locked laser, the backscattering can be due to an interface, in which case $\tilde{r}_{21} = -\tilde{r}_{12}^* = -\tilde{r}^*$; and $\theta_{21} = \theta_{12} + \pi = \theta + \pi$. This type of coupling does not prevent lock-in, since Eq. (38) becomes

$$\frac{\partial\psi}{\partial t} = \frac{\delta}{\tau_{RT}} - \frac{r}{\varepsilon_1\varepsilon_2\tau_{RT}} [\varepsilon_1^2 + \varepsilon_2^2] \sin(\theta + \psi) \quad (41)$$

which, for sufficient large r , still has a lock-in solution ψ_0 for which $\partial\psi/\partial t = 0$.

Phase Conjugated Coupling Not all couplings lead to identical mode frequencies of the two output pulse trains of the laser. In a phase conjugated coupling, a fraction r_c of the complex conjugate of one field is coupled into the other field. Such a phase conjugated coupling¹⁴ does preserve the phase identity of each intracavity pulse. The coupled equations for the two pulses are then

$$\begin{aligned} \frac{\partial\tilde{\varepsilon}_1}{\partial t} &= i\frac{\delta}{2\tau_{RT}}\tilde{\varepsilon}_1 + \frac{r_c}{\tau_{RT}}\tilde{\varepsilon}_2^* \\ \frac{\partial\tilde{\varepsilon}_2}{\partial t} &= \frac{r_c}{\tau_{RT}}\tilde{\varepsilon}_1^* - i\frac{\delta}{2\tau_{RT}}\tilde{\varepsilon}_2 \end{aligned} \quad (42)$$

Subtracting the imaginary parts of this equation:

$$\frac{\partial\psi}{\partial t} = \frac{\delta}{\tau_{RT}} \quad (43)$$

and there is no “lock-in” possible with this type of coupling.

Repetition Rate Coupling It has been observed that the repetition rate in both directions can be locked by a saturable absorber. The mechanism by which the average group velocity of the two pulses are locked to each other is described below. This mechanism leaves the carrier frequencies of the two pulses uncoupled.

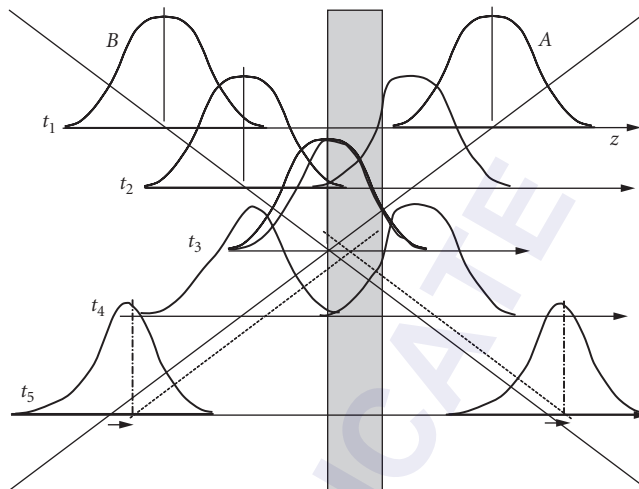


FIGURE 6 Representation of the intracavity pulses entering a saturable absorber. The pulses are plotted as a function of space (z) at successive times. The saturable absorber is initially at the left of the pulse crossing point. Because of mutual saturation, there is only significant absorption when only one of the pulses is present in the absorber. Therefore, the leading edge of pulse A is attenuated more, resulting in an apparent slowing down of the pulse. Similarly, the trailing edge of pulse B is absorbed more, resulting in an apparent acceleration of that pulse. The effect of the absorption combined with mutual saturation is to “pull” the pulse crossing point toward the center of the absorber.

To appreciate this envelope coupling, let us consider a saturable absorber of smaller longitudinal dimensions than the optical pulse, as sketched in Fig. 6. In the figure, the meeting point of the two pulses is on the left of the saturable absorber. Therefore, the pulse A entering from the left enters first the absorber, and its leading edge is attenuated. The absorption is saturated when the two pulses meet in the absorber. The pulse B coming from the left is still partly in the absorber when pulse A has left the absorber. Therefore, its tail will be more absorbed. The net effect is a shift of the center of gravity of both pulses, such that at the next round-trip they will meet closer to the middle of the absorber.

20.5 DESIGNS OF CAVITIES WITH TWO CIRCULATING PULSES

The properties of interwoven intracavity pulses have been discussed in the preceding sections. A few examples of laser cavity designs where two pulses are circulating independently are presented next.

Ring Dye Laser

In a bidirectional ring laser, two pulses circulate in opposite direction. The first realization was a dye laser,^{15,16} in which the gain is provided by a jet of Rhodamine 6G in ethylene glycol (pumped by an argon ion laser). A saturable absorber jet of DODCI* has three functions: (i) to mode-lock the laser,

*Di-oxa-di-carbo-cyanide-iodide.

(ii) to ensure bidirectional operation of the laser, and (iii) to define and maintain the pulse crossing point of the envelopes. Bidirectional operation is favored over unidirectional operation because, for the same pulse energy and duration, a single pulse will suffer more intracavity losses than two pulses creating a standing wave by crossing in the absorber.³ The pulse crossing point is set and maintained by the mechanism just described in Section “Locking Two Pulse Trains by Backscattering”. The saturable absorber has to be located approximately $1/4$ cavity perimeter ($P/4$) away from the gain jet, in order for each counter-circulating pulse to enter the gain medium at equal time interval ($P/2c$). The phase difference per round-trip between the clockwise (CW) and counter-clockwise (CCW) pulses is measured by combining the CW and CCW pulses on a detector.

Ti:Sapphire Ring Laser with Saturable Absorber

A similar cavity configuration has been used with a Ti:sapphire laser as a gain medium, and a saturable absorber jet of HITCI* for repetition rate synchronization between the two counter-circulating pulses. A sketch of such a cavity is shown in Fig. 7.

The nonlinear index of the gain medium results in a lensing effect, which can be approximated by a positive lens collocated with the gain medium. In mode-locked operation, Kerr lensing induces a positive lens at the location of the gain medium, which modifies the beam size distribution. For the empty cavity, the beam size versus position is represented by the solid line in the graph of Fig. 8. The beams size distribution modified by the self-lensing in the gain rod is indicated by the dotted line in the figure. An aperture located at the position A_1 will favor mode-locked operation, since the losses will decrease with intensity. An aperture located at A_2 will create increasing losses with intensity. This negative feedback stabilizes the mode-locked laser operation. Kerr lensing, which could be considered as an instantaneous saturable absorption,³ is the technique commonly used to generate the shortest pulses with Ti:sapphire lasers. It is however not a preferred technique for achieving stable bidirectional operation, when, as is typically the case, the active element for Kerr lensing is the gain medium. There is a competition in the gain medium between mutual Kerr lensing, favoring bidirectional operation (with the pulses crossing in the gain medium) and mutual gain saturation favoring unidirectional operation, with the latter generally dominating. An experimental study of a Kerr-lens mode-locked Ti:sapphire laser¹⁸ showed unidirectional operation, switching direction periodically (approximately every 0.1 second). The operation became bidirectional after insertion of a dilute saturable absorber jet inside the cavity.¹⁸

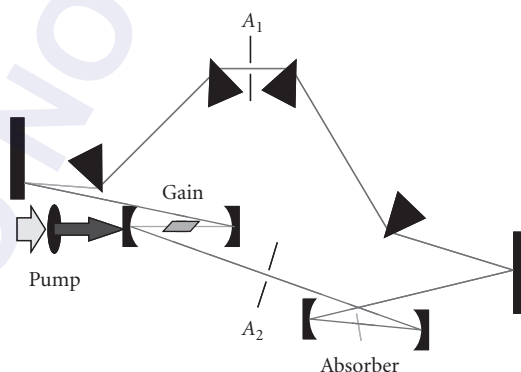


FIGURE 7 Ti:sapphire mode-locked ring laser mode-locked with a saturable absorber jet. Four prisms are used for the control of cavity dispersion.¹⁷

*Hexa-methyl-indo-carbo-cyanide-iodide.

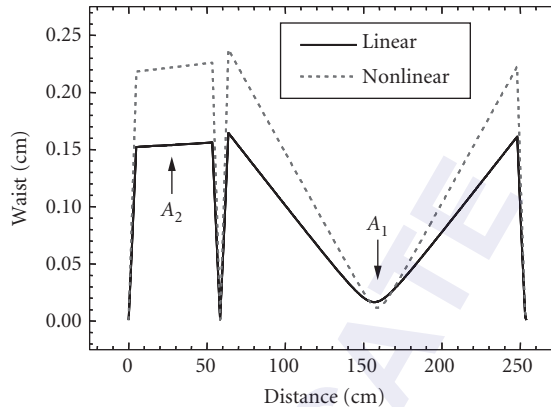


FIGURE 8 Calculation of the cavity mode as a function of position along the cavity, for the empty cavity (solid line) and the cavity modified by nonlinear lensing at the position of the gain medium (dotted line).

The use of a liquid saturable absorber flowing at high velocity (several meters per second) through a narrow nozzle (typically $100\ \mu\text{m}$ thick and $5\ \text{mm}$ wide) is essential to ensure the absence of phase coupling between the two pulses. If the saturable absorber were a nonmoving solid, the coupling by scattering would result in a mutual locking of the carrier frequency of the two pulses, as discussed in section “Locking Two Pulse Trains by Backscattering.” In the case of the moving fluid of the absorbing dye jet, θ_{12} and θ_{21} are both random functions of time, varying much faster than ψ . Over the time scale that the variation of ψ is negligible, the last terms of Eq. (38) average to zero. Therefore, the dead band has been eliminated, as has been verified experimentally.¹⁹

Ring Laser with Additional Kerr Crystal

Another technique to achieve bidirectional operation for a Kerr-lens mode-locked laser^{20,21} is to insert a nonlinear crystal (for which the nonlinear phase shift is larger than that produced in the gain medium) $1/4$ cavity perimeter away from the gain medium. The laser cavity is sketched in Fig. 9.

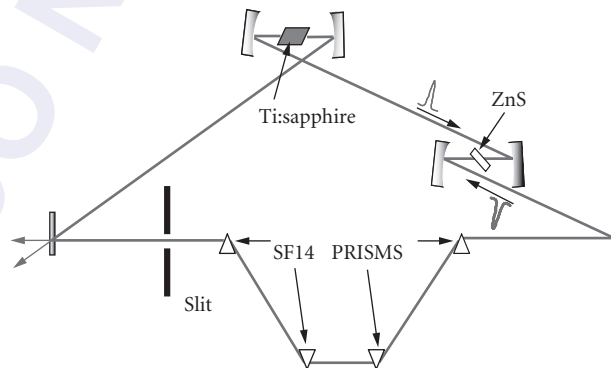


FIGURE 9 Ring laser mode-locked with a nonlinear crystal (ZnS).

The modification to the nonlinear index of refraction due to two counter-propagating fields of direction i and j is

$$n_2^i = n_2(I_i + 2I_j) \quad (44)$$

The factor 2 on the right hand side of Eq. (44) reflects the well-known act that cross-phase modulation is twice as effective as self-phase modulation for the same intensity.²² Assuming equal intensity in the counter-propagating fields and similar waists for single pulse operation versus bidirectional operation yields the following approximate relationship for the nonlinear index of refraction in the ZnS crystal for single pulse (unidirectional) versus double pulse (bidirectional) operation:

$$\Delta n_{\text{bidirectional}} \sim \frac{3}{2} \Delta n_{\text{single}} \left[\frac{\text{interaction length}}{\text{crystal length}} \right] \quad (45)$$

Because the overlap region of the pulse envelopes is approximately the pulsewidth/ c , the mutual Kerr-lens will dominate only if the length of the Kerr medium is not much longer than the pulsewidth. Depending on the pulsewidth and crystal thickness, there may be enough nonlinearity to distinguish between single-pulse operation and bidirectional operation. To enhance the mutual Kerr effect, a crystal of ZnS was chosen, because its nonlinear index is 50 times larger than that of Ti:sapphire.²¹ Pulses as short as 60 fs were generated, meeting in the ZnS crystal used as a nonlinear element. Because the two countercirculating pulses have different intensities, a differential phase shift between the two directions results in a difference of cavity modes of 60 kHz. In addition to the relative shift of the modes, this laser system has the additional complexity that the center of spectral envelope of either countercirculating pulse is shifted by 2 nm.

Imperfection in the ZnS crystal used as nonlinear element resulted in coupling of the beams by scattering. The amount of coupling can be measured through the spectrum of the beat note of the two pulse trains. Such a spectrum shows for instance a fundamental at 60 kHz and harmonics at 120 kHz and 180 kHz. Such a spectrum can lead to the lock-in frequency between the two beams,²³ which in the case of the experiment cited is 1.8 kHz.²¹

Linear Lasers

Considering the ring laser sketched on the top of Fig. 7, it is possible to visualize stretching the cavity by the two mirrors at the extreme left and right, while keeping the perimeter constant. The limit of the stretched out ring is a linear cavity in which two pulses circulate.

As an example, let us consider a linear cavity used to measure with high accuracy the electro-optic coefficient.²⁴ The laser cavity is similar to the typical linear cavity mode-locked Ti:sapphire laser, but with a saturable absorber (a jet of HITCI dye dissolved in ethylene glycol) placed in the center of the cavity, as sketched in Fig. 10. As in the case of the ring laser discussed above, Kerr-lens mode-locking does not appear to be possible with double pulse operation. Instead, a saturable absorber is positioned in the middle of the cavity by translating one of the end mirrors. The distance that the end mirror can be translated while maintaining double pulses is about 2 cm, in excess of the pulse length of approximately 0.6 mm (2-ps pulses). The 2-cm distance is a 120-ps delay and corresponds to the lifetime of the dye. The dye concentration is not a critical parameter and can be varied over a broad range without affecting the performance of the laser. The 140-MHz output from one end of the laser, detected on a fast photodiode, is filtered, its frequency divided by 2 in an ECL logic. The resulting 70-MHz signal is which yielded a 70-MHz sinusoidal signal. Finally the signal is amplified again and applied to synchronize the measurement to be performed. In the case of measurement of an electro-optic coefficient, the 70-MHz signal is applied directly to electrode on the crystal to be measured, in parallel with a 50-ohm terminator.²⁴

Whether to use a ring or linear cavity depends on the quantity to be measured. The ring laser is sensitive to rotation, and to fresnel drag in one of its arms.²⁵ Without any rotation or modulation a

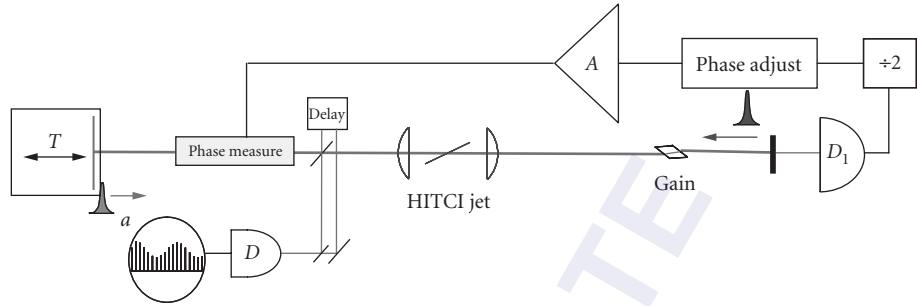


FIGURE 10 Linear laser mode-locked with a saturable absorber to produce two pulses/cavity round-trip. The end mirror on the left is on a translation stage T , in order to set the position of the saturable absorber in the middle of the cavity. The saturable absorber is a jet of HITCI dissolved in ethylene glycol, between two focusing elements. The output pulse train, recorded on a photodiode, reduced to half frequency, is used to synchronize a phase measurement. Each of the two intracavity pulses experience a different phase shift per round-trip, hence a different carrier frequency. The two intracavity pulses are extracted from the cavity with a beam splitter, and recombined after an appropriate optical delay. The two interfering pulse trains show a beat signal on the detector D .

mode-locked ring laser normally has a beat frequency offset of at least 100 Hz and often as high as 100 kHz.²⁶ This is a result of the asymmetry in the CW and CCW pulse. Because of the nonlinear intracavity elements, the order in which the pulse encounters the optical elements will affect the pulsewidth and pulse amplitude.^{3,27,28} Any variation in pulse amplitude or pulsewidth will be seen as a beat signal. Since the pulses in a linear cavity travel through the same optical elements in the same order, there is no asymmetry. Therefore, one advantage of a linear cavity versus a ring geometry is the improvement in the frequency offset.

As the electronic delay of the signal applied to the sample is varied, the beat note shows a sinusoidal dependence, as shown in Fig. 11, which is a plot of the beat frequency versus the delay. The optimum timing occurs when one pulse sees a voltage on the sample of $+V_0$ and the second pulse sees a voltage of $-V_0$ at the sample. The line plotted in Fig. 11 is not a fit, but a plot of $V_0 \left| \sin\left(\frac{2\pi c}{2L} \tau - \phi_0\right) \right|$, where the fixed phase ϕ_0 was the only free parameter.

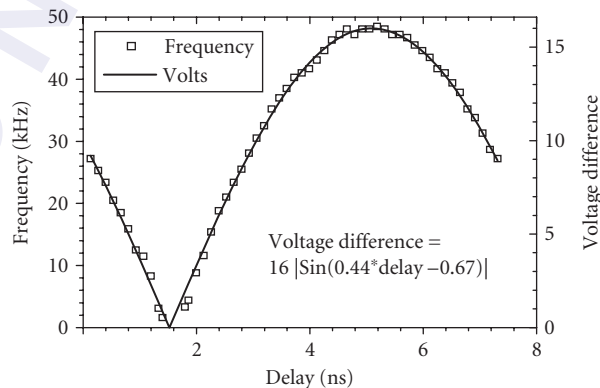


FIGURE 11 Delay dependence of the beat note frequency.

Optical Parametric Oscillators

The coupling between oppositely circulating beams in a ring laser is eliminated if a laser operates with ultrashort pulses* circulating in opposite direction and crossing in a nonscattering medium such as vacuum or air. It is however quite a challenge to find means to couple the pulse envelopes, without introducing any phase coupling. One solution described in the previous sections is to couple the pulses in amplitude in a medium that moves transversally to the beam, such as a jet of liquid saturable absorber. Another solution is the synchronously pumped optical parametric oscillator (OPO), which offers the possibility to decouple relative phase and repetition rates of the oscillating signals, without the need for any moving element.

A simple configuration is that of an OPO-pumped extracavity by a Ti:sapphire mode-locked laser, as sketched in Fig. 12. The position of the crossing point of the two circulating pulses in the OPO is simply determined by the timing of the pump pulses, rather than by a saturable absorber¹⁹ or a nonlinear crystal.²¹ The mode frequencies are still set by the cavity. Another advantage over other systems is the tunability, which is important for applications such as detecting ultra low magnetic fields, where the laser radiation has to be tuned to a narrow atomic transition.

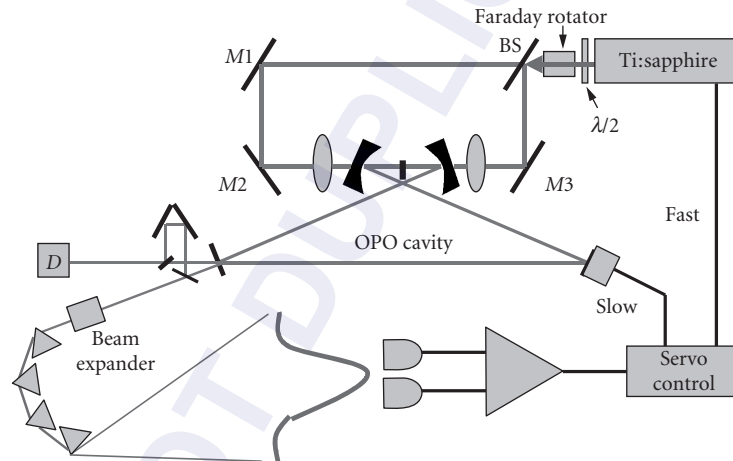


FIGURE 12 Illustration of the OPO cavity pumped by the Ti:sapphire laser. The reflected and transmitted parts of the beam splitter BS are focused into the periodically poled lithium niobate crystal via the two branches of an antiresonant ring. The beams should destructively interfere at the antiresonant ring output, which is monitored with a CCD. Since all the radiation is, in exact alignment, reflected back into the laser, a two-stage optical isolator is required to prevent disruption of the mode-locked operation. The nonlinear crystal is a 0.8-mm-long periodically poled lithium niobate (PPLN) crystal with a period of $19.75 \mu\text{m}$, temperature stabilized at 353 K to prevent photorefractive damage and achieve quasi-phase-matching condition for generation of a signal of $1.35 \mu\text{m}$ with an average power of 30 mW per direction. The difference between the two optical paths from the beam splitter to the crystal determines the crossing point of the signal pulses in the OPO cavity. The two output pulses are made to interfere on a detector *D* after an optical delay line brings them in coincidence. A four-prism sequence sends the pulse spectrum onto a pair of detectors. The difference between the two detected signals is amplified and applied to a piezo to stabilize the cavity length.

*Pulse duration τ in the range from femtoseconds to a few picoseconds, pulse length $C\tau$ much shorter than any linear dimension in the cavity.

A periodically poled LiNbO_3 crystal (PPLN) is excited by a “pump pulse” at an optical frequency ω_p to provide gain at a “signal” frequency ω_s through the process $\omega_p = \omega_s + \omega_i$ where ω_i is the “idler.” The OPO is an oscillator which uses the gated gain at ω_s . Therefore, as opposed to a conventional gain medium, in an OPO, the timing, direction, and position of the gain are determined by the time of arrival, k vector and focal spot location of the pump pulse in the parametric crystal. A Ti:sapphire laser provides a train of gating pump pulses of 200-fs duration at 789 nm, 143-MHz repetition rate, and 400-mW average power. The cavity of the pump laser is in a ring configuration, and operates unidirectionally. Such a configuration is less sensitive to feedback than a linear cavity. A double stage Faraday isolator (providing -60 -db isolation) is still required to prevent the feedback from the OPO antiresonant ring pumping arrangement (Fig. 12)^{29,30} from destroying the mode-locked operation. The sensitivity of the OPO wavelength to cavity mismatch can be exploited to stabilize the synchronously pumped OPO. Indeed, since the repetition rate of the OPO is fixed by the pump laser, the signal wavelength will adjust to a value for which the round-trip rate matches the pump rate.^{31,32} As a result, any fluctuation of OPO cavity length relative to that of the pump cavity will be translated into a change in wavelength of the OPO laser. In the arrangement sketched in Fig. 12, motions of the spectrum are detected by spectrally dispersing (with 4 prisms) an expanded output from the counter-clockwise OPO beam. The signal spectrum, centered at $1.35 \mu\text{m}$, is split into two parts and collected by a pair of lenses into two infrared photodiodes (Fig. 12). The difference signal of the two detectors monitoring two spectral components on either side of the pulse spectrum is sent through a high-gain amplifier ($\omega_{3db} = 1 \text{ kHz}$) to drive piezoelectric transducers (PZT) translating an OPO (slow servo loop) and a Ti:sapphire mirror (fast servo loop).

The beat note observed with the configuration of Fig. 12 has a bandwidth of tens of kilohertz,³³ because of the extreme (nanometer) sensitivity of the OPO to the pump spot position, as demonstrated experimentally in reference.³⁴ The beat note bandwidth is thus fundamentally due to fluctuations in the gain spot position for either circulating pulse, due to the beam pointing instability of the pump laser. The basic remedy is to make the two pump spots part of the same spatial mode of a cavity. One solution that has been implemented³⁴ is to insert the OPO crystal inside the cavity of the pump laser. Implementation of an OPO pumped intracavity by a linear Ti:sapphire laser is shown in Fig. 13. Four LaKL21 prisms are incorporated in the pump cavity to compensate the group velocity dispersion (GVD) from the Ti:sapphire crystal, the PPLN crystal, and other intracavity elements such as lenses and mirrors. This four-prisms configuration was necessitated by the desire to have large GVD compensation (needed because of the large positive GVD of LiNbO_3) and a reasonably short cavity length ($1/2$ of the perimeter of the OPO cavity). Two quantum wells (MQW) of AlGaAs on top of a mirror structure are used in the cavity as a saturable absorber to mode-lock the laser. The

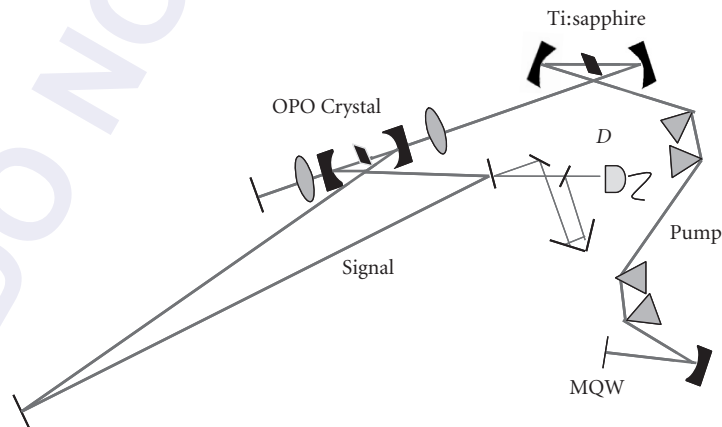


FIGURE 13 Illustration of the intracavity OPO pumped by the Ti:sapphire laser. The main control of the GVD compensation is the prism spacing L_2 .

Ti:sapphire laser radiation consists of 200-fs pulses centered at 785 nm, repetition rate of 95 MHz. The OPO crystal is a 3-mm-long Brewster cut PPLN crystal (HC Photonics, Taiwan) with a period of 19.4 μm (quasi-phase matching for signal near 1.36 μm), which is temperature stabilized at 408 K to prevent photorefractive damage. Attempts to use OPO crystals cut near normal incidence and antireflection coated failed, because the Ti:sapphire mode-locked operation was prevented by the smallest feedback from the antireflection coating. This feedback problem was completely eliminated through the use of a Brewster angle cut, but at the price of a considerably more difficult alignment procedure. Because of the Brewster angle cut of LiNbO_3 , the idler, the pump radiation and its second harmonic all exited the crystal at a different angle, and could not be used for alignment of the OPO cavity (the OPO cavity mirrors were reflecting at the signal wavelength of 1.4 μm and at the second harmonic of the pump).

20.6 ANALOGY OF A TWO-LEVEL SYSTEM

In the previous section, some methods for mode-locking a laser with two intracavity pulses were described. Inside the cavity, these two pulses can couple to each other in diverse ways which were analyzed in Sec. 20.4. The purpose of circulating two pulses in a laser cavity is that intracavity perturbations created by the quantity to be measured will alter the carrier frequency of a pulse, a shift in frequency that can easily be detected by interfering the output pulse trains. The laser itself is used as an interferometer, with the remarkable properties that the phase shift occurring in the cavity is transformed in a frequency shift. There are numerous factors influencing the accuracy as well as the sensitivity of a measurement performed intracavity. A better understanding of the two pulse per cavity laser can be reached by noting the complete analogy with a quantum mechanical two-level system. In the quantum mechanical situation an atomic or molecular system can be in one or two quantum states $|k\rangle$, with $k = 1$ and 2, of energy $\pm\omega_0/2$. Each of these states correspond to pulse $|2\rangle$ and pulse $|1\rangle$ in the laser cavity. For instance, in a ring laser, one of the states would correspond to a counterclockwise circulating pulse, the other to the clockwise circulating pulse. The interaction of a two-level system with a near resonant field is the most thoroughly studied problem in atomic and molecular physics. Techniques developed to achieve sublinewidth resolution in atomic physics may be transposed to the laser situation, and, thanks to the analogy, lead to methods to enhance the resolution of intracavity laser sensors.

Review of Coherent Interaction of Two-Level Systems

Considering the case of two level with a dipole allowed transition, in presence of a near resonant electromagnetic field $E = 1/2\tilde{E}(t)\exp(i\omega t) + \text{c.c.}$ In presence of this electric field, the state of the atomic/molecular system is described by the wavefunction ψ , a solution of the time-dependent Schrödinger equation:

$$H\psi = i\hbar \frac{\partial\psi}{\partial t} \quad (46)$$

with the total Hamiltonian given by

$$H = H_0 + H' = H_0 - p \cdot E(t) \quad (47)$$

where p is the dipole moment. In the standard technique for solving time dependent problems, the wave function ψ is written as a linear combination of the basis functions $|k\rangle$:

$$\psi(t) = \sum_k a_k(t) |k\rangle \quad (48)$$

This expression for ψ is inserted in the time dependent Schrödinger Eq. (46). Taking into account the normalization conditions for the basis functions ψ_k , one finds the coefficients a_k have to satisfy the following set of differential equations:

$$\frac{da_k}{dt} = -i\omega_k a_k + \sum_j \frac{i}{2\hbar} p_{k,j} [\tilde{\mathcal{E}} e^{i\omega_j t} + \tilde{\mathcal{E}}^* e^{-i\omega_j t}] a_j \quad (49)$$

where $p_{k,j}$ are the components of the dipole coupling matrix for the transition $k \rightarrow j$, and a_k are the amplitudes of the eigenstates. Phase and amplitude relaxation have been neglected so far and will be introduced later. It should be noted that Eq. (49) is of a quite general nature, is ideally suited to numerical integration, and is not limited to two level systems. Similarly, the laser analogy can be extended to lasers with more than two intracavity pulses. We will consider here only two levels, with a very small detuning $\Delta\omega = \omega_0 - \omega \ll \omega_0$:

$$\Delta\omega = \omega_0 - \omega \quad (50)$$

Consistent with the approximation of small detuning, we replace the set of coefficients a_k , which have temporal variations at optical frequencies, by the “slowly varying” set of coefficients c_k , using the transformation:

$$a_k = e^{-ik\omega_0 t} c_k \quad (51)$$

Inserting in the pair of Eqs. (49), leads to the pair of differential equations for the two coefficients c_k :

$$\frac{d}{dt} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} i\frac{\Delta\omega}{2} & i\frac{1}{2\hbar} p \tilde{\mathcal{E}} \\ -i\frac{1}{2\hbar} p \tilde{\mathcal{E}}^* & -i\frac{\Delta\omega}{2} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \quad (52)$$

where $\kappa|\tilde{\mathcal{E}}| = (p/\hbar)|\tilde{\mathcal{E}}|$ (p being the dipole moment of the single photon transition) is the Rabi frequency.

The Laser as a Two-Level System

The two-level system considered above is isolated, hence the total population is conserved. The analogue of the electromagnetic field coupling states $|1\rangle$ and $|2\rangle$ is a conservative intracavity coupling $\tilde{r}_{12} = -\tilde{r}_{21}^*$. Such a type of coupling can only be considered in the case of mode-locked lasers, where the localization of the radiation in the cavity enables one to select a truly conservative coupling. The coupling, localized at the crossing point of the two circulating pulses, can be produced by the back-scattering at a dielectric interface[†] between two media a and b , for which $\tilde{r}_{ab} = \tilde{r}$ and $\tilde{r}_{ba} = -\tilde{r}^*$. It can easily be verified that the total intensity change introduced by this coupling is zero, as expected for a conservative coupling. In fact, the phase relation between the two reflections at either sides of the interface is a consequence of energy conservation.

In the analogy of the laser, the coefficients $c_i(t)$ correspond to the complex field amplitudes $\tilde{\mathcal{E}}_i$ (the tilde indicating a complex quantity) of each pulse circulating in the ring cavity (round-trip

[†]In the case of a linear laser with two pulses/round-trip, the conservative coupling is only possible when there are two crossing points in the cavity, and that the interface is located at one of the crossing points.

time τ_{RT}). The state of the system is also defined by $\psi(t) = \tilde{\epsilon}_1(t)|1\rangle + \tilde{\epsilon}_2(t)|2\rangle$. The evolution equation of these fields are

$$\frac{d}{dt} \begin{pmatrix} \tilde{\epsilon}_1 \\ \tilde{\epsilon}_2 \end{pmatrix} = \frac{1}{\tau_{RT}} \begin{pmatrix} \tilde{r}_{11} & \tilde{r}_{12} \\ \tilde{r}_{21} & \tilde{r}_{22} \end{pmatrix} \begin{pmatrix} \tilde{\epsilon}_1 \\ \tilde{\epsilon}_2 \end{pmatrix} = \frac{1}{\tau_{RT}} ||R|| \cdot ||E|| \quad (53)$$

In order to have an equivalence between Eqs. (52) and (53), the matrix $||R||$ should be *Anti-Hermitian*, which, in addition to the condition $\tilde{r}_{21} = -\tilde{r}_{12}^*$, imposes and that \tilde{r}_{kk} be purely imaginary. It can also easily be verified that this is the only form of interaction matrix for which energy is conserved $d/dt(|\tilde{\epsilon}_1|^2 + |\tilde{\epsilon}_2|^2) = 0$. The real parts of the diagonal elements of the matrix $||R||$ represent gain and loss in the cavity. In steady state, the gain and loss are in equilibrium, and the real parts of \tilde{r}_{kk} are zero. A gain (or absorber) with a recovery (relaxation) time longer than $\tau_{RT}/2$ will cause transients in population. The laser equivalent to the detuning $\Delta\omega$ is a differential phase shift for the pulses $|1\rangle$ and $|2\rangle$ in the cavity. Such a differential phase shift is introduced either by rotation in a ring laser,^{15,16} or with an electro-optic modulator in a linear laser.²⁴ In the latter case, the electro-optic phase modulator imposes an opposite phase shift ($\Delta\phi/2$ and $-\Delta\phi/2$) for either pulse, thereby modifying the resonance of the cavity for the pulse $\tilde{\epsilon}_1$ by $\Delta\omega/2 = \Delta\phi/(2\tau_{RT})$, and for pulse $\tilde{\epsilon}_2$ by $-\Delta\omega/2 = -\Delta\phi/(2\tau_{RT})$. These detuning terms contribute to the diagonal terms of the matrix $||R||$: $\tilde{r}_{11} = -\tilde{r}_{22} = i\Delta\phi/2$.

The analogy between the laser and a two-level system applies also to the set of density matrix equations. These can be obtained by rewriting Eq. (53) in terms of the intensities in either sense of rotation $\rho_{22} = \tilde{\epsilon}_2 \tilde{\epsilon}_2^*$ and $\rho_{11} = \tilde{\epsilon}_1 \tilde{\epsilon}_1^*$, and the quantities $\rho_{12} = \tilde{\epsilon}_1 \tilde{\epsilon}_2^*$ and $\rho_{21} = \tilde{\epsilon}_2 \tilde{\epsilon}_1^*$:

$$\frac{d(\rho_{22} - \rho_{11})}{dt / \tau_{RT}} = -4\text{Re}(\tilde{r}_{12}\rho_{21}) - \frac{(\rho_{22} - \rho_{11})}{T_1} \quad (54)$$

$$\frac{d\rho_{21}}{dt / \tau_{RT}} = -i\Delta\omega\tau_{RT}\rho_{21} + \tilde{r}_{12}^*(\rho_{22} - \rho_{11}) - \frac{\rho_{21}}{T_2} \quad (55)$$

where, as in the case of the two-level system interacting with a near resonant field, phenomenological relaxation times T_1 and T_2 have been introduced. One recognizes here Bloch's equation for a two-level system driven off-resonance by a step function Rabi frequency of amplitude \tilde{r}/τ_{RT} .³⁵ The difference in intensities $(\rho_{22} - \rho_{11})$ is the direct analogue of the population difference between the two levels. The off-diagonal matrix element ρ_{21} is the interference signal obtained by beating the two outputs of the laser on a detector. As in the case of the two-level system, one can introduce phenomenological relaxation times T_1 for the energy relaxation (diagonal matrix element) and T_2 for the coherence relaxation (off-diagonal matrix elements). As for the quantum mechanical two-level system, $1/T_2$ is the homogeneous component of the linewidth of the beat note between the two pulse trains. There is also an "inhomogeneous" component to that linewidth, which has as physical origin the mechanical vibration of the laser components, causing random fluctuations of the beat note. Because of mechanical vibrations, each pulse sees random differences in the cavity length caused by mirror motion over a time of $\tau_{RT}/2$.

Table 1 summarizes the main points of the analogy between a laser with two pulses/cavity and the coherent interaction of a two-level system with a near resonant electromagnetic field.

Experimental Demonstration of the Analogy

The most typical manifestation of a two-level system interacting with a step function electromagnetic pulse is Rabi cycling, which is a periodic transfer of population from one state to the other. To observe such a periodic transfer, the system should be in one of the two states at $t = 0$. One method to prepare the ring laser with one state (direction) dominating, is to feedback one direction into the other outside of the cavity (Fig. 14). The output pulse from one direction is extracted, and fed back (<1 percent) with a mirror, after appropriate optical delay, into the opposite direction. By using a fast switch (turn-off time of less than the cavity round-trip time of 10 ns) at the Pockel's cell, the coupling can be turned off, to let the counter-circulating fields evolve in the cavity.

TABLE 1 Summary of the Analogy between a Two-Level System and the Laser with Two Circulating Intracavity Pulses

	Two-Level System	Laser
Basic states	$ k\rangle; k=1, 2$	$ k\rangle; k=1, 2$
corresponding to	energy level $\pm \frac{\omega_0}{2}$	Intracavity pulses 1, 2 selected by geometry
Coupling through	Near-resonant E-field at ω	Backscattering at interface
Detuning	$\Delta\omega = \omega_0 - \omega$	$\Delta\omega = \Delta\phi \tau_{RT}$
Slowly varying	$\Delta\omega \ll \omega$	$\Delta\omega \ll 1/\tau_{RT}$
Wave function	$\psi(t) = a_1(t) 1\rangle + a_2(t) 2\rangle$	$\psi(t) = \varepsilon_1(t) 1\rangle + \varepsilon_2(t) 2\rangle$
Density matrix	$\rho_{kk} = a_k a_k^*$ Populations	$\rho_{kk} = \tilde{\varepsilon}_k \varepsilon_k^*$ (Intensities)
elements	$\rho_{kk} = -a_i a_j^*$	$\rho_{ij} = \tilde{\varepsilon}_i \varepsilon_j^*$ (beat signal)

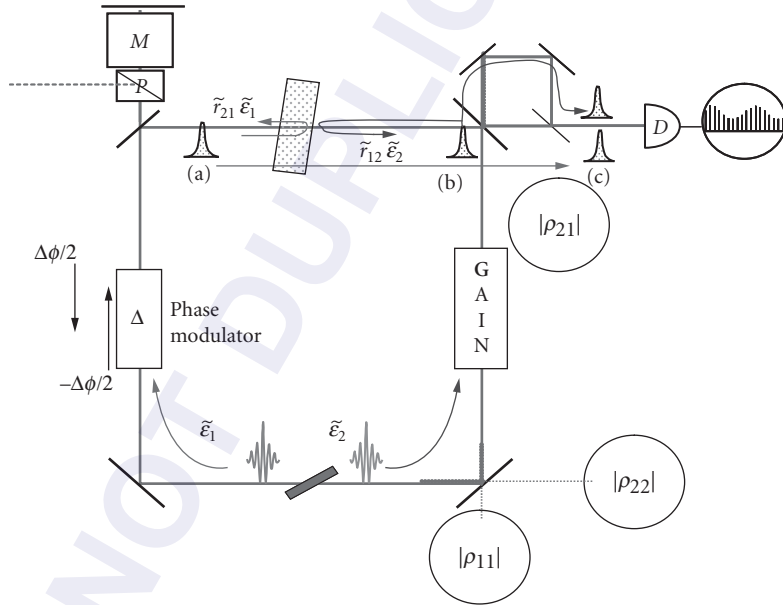


FIGURE 14 Sketch of the ring laser used to demonstrate the analogy with a two-level system. In the bidirectional mode-locked ring laser two circulating pulses meet in a saturable absorber jet. Three successive positions: (a), (b), and (c) of the two pulses are shown. An interface, positioned at or near the opposite crossing point of the two pulses, controls the amplitude of the coupling parameter \tilde{r}_{ij} . The circulating intensities in the laser, measured for each direction by quadratic detectors, are the diagonal elements (populations) of the density matrix of the equivalent two-level system. The absence of phase modulation corresponds to the two levels being on resonance, driven at the Rabi frequency $(p/\hbar)\mathcal{E}$ by a resonant field (the Rabi frequency $(p/\hbar)\mathcal{E} = \kappa\mathcal{E}$ corresponds to the frequency r_{12}/τ_{RT} in the ring laser analogy). The backscattering at the interface provides thus coherent coupling (Rabi cycling) between the two states, while other noncoherent decays tend to equalize the population in the two directions, and washes out the phase information. The detuning $\Delta\omega$ corresponds to the phase difference per round-trip $\Delta\phi/\tau_{RT}$, imposed by an electro-optic phase modulator driven exactly at the cavity round-trip time. A beat-note detector measuring the interference between the two fields, records the off-diagonal matrix element. A combination of a Pockel's cell M and polarizer P controls a feedback of the clockwise pulse into the counterclockwise one.

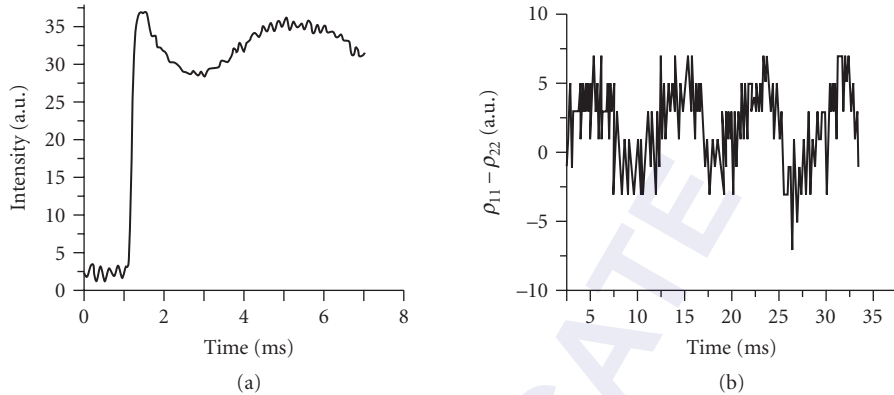


FIGURE 15 The evolution of the intensities are shown after switching the Pockels' cell. (a) The counterclockwise direction is shown—the intensity at clockwise direction is 180° out of phase with this graph, with population dropping from the maximum initial value. The fast initial transient reflects the gain and cavity dynamics associated with the sudden change in cavity loss at the switching time. Thereafter, a slow oscillation due to population transfer or Rabi oscillation between two directions is observed. (b) Population difference showing the Rabi cycling.

Rabi Cycling on Resonance In the measurements that follow, the system is “at resonance”; (i.e., $\Delta\omega=0$). An example of “Rabi cycling” is shown in Fig. 15. The counterclockwise intensity (ρ_{22}) is plotted as a function of time [(Fig. 15a)]. The clockwise intensity ρ_{11} (not shown) is complementary. The system is prepared so that the ρ_{11} is initially populated ($\rho_{11}=0.8, \rho_{22}=0.2$). As the feedback that creates the initial state is switched off at $t=1$ ms, there is a fast (approximately $10 \mu\text{s}$) transient. This risetime reflects combined dynamics of the gain and cavity, as the laser adapts to the different (now symmetrical) cavity losses. This risetime corresponds roughly to the fluorescence lifetime of the upper state of Ti:sapphire. The “Rabi cycling” of the “population difference” $\rho_{22}-\rho_{11}$ is plotted in Fig. 15b. One can also record the beat note frequency (off-diagonal element $|\rho_{12}|$) as sketched in Fig. 14. As can easily be seen from the Bloch vector model of Feynman et al.,³⁵ the oscillation of the diagonal elements and the off-diagonal elements are 90° out of phase. This property can indeed be seen in Fig. 16a. The Rabi frequency $|\tilde{r}|/\tau_{RT}$ can be varied by changing the position of the scattering surface, as shown in Fig. 16b. The maximum value measured³⁶ for this interface corresponds to a

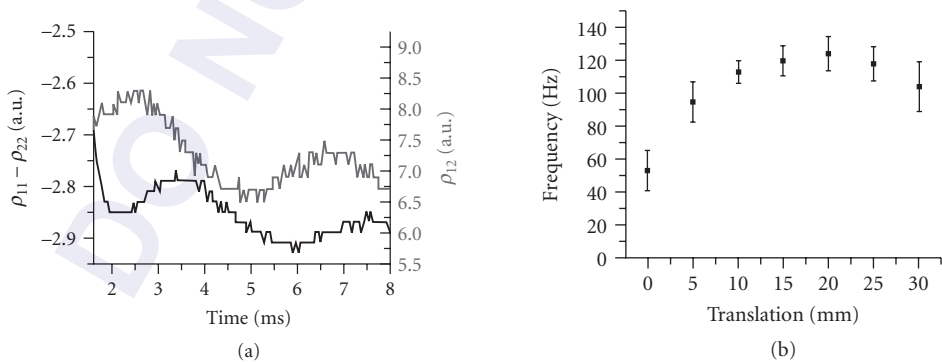


FIGURE 16 (a) Comparison of the oscillation of the population difference $\rho_{22}-\rho_{11}$ and the off-diagonal element (beat note) ρ_{12} . (b) Rabi frequency as a function of position of the glass at the meeting point of the two directions. Translation of the glass-air interface along the beam results in different values of coupling \tilde{r} .

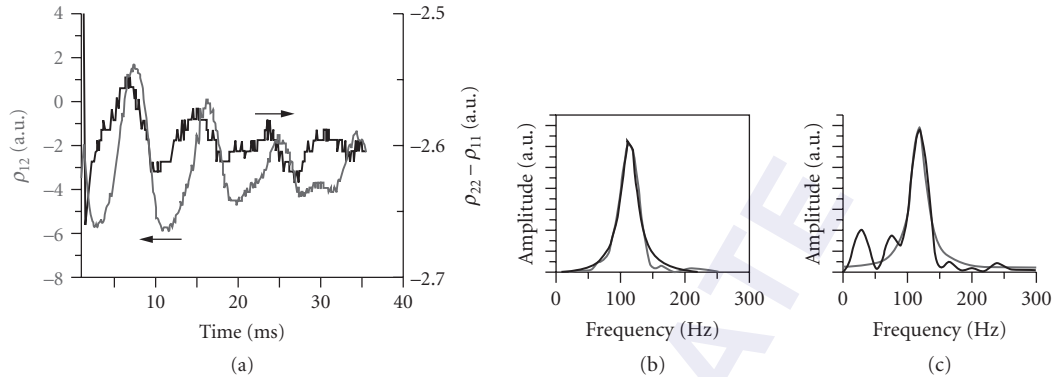


FIGURE 17 (a) Measurement of the decay of the Rabi oscillations $\rho_{22} - \rho_{11}$ and ρ_{21} . (b and c) The Fourier transforms of the relative measurements are shown on the right.

backscattering coefficient of $|\tilde{r}| \approx 1 \cdot 10^{-6}$. Note that the Rabi frequency provides a direct measurement of very minute backscattering coefficients, without the need to trace a complete gyroscopic response as in Refs. 36 and 37.

In the case of a two-level system, the phenomenological “longitudinal” and “transverse” relaxation times have been identified as energy relaxation time (fluorescence decay) and phase relaxation time (due for instance to atomic collisions). Figure 17 shows a measurement of the decay of the Rabi oscillation for the diagonal and off-diagonal elements. The decay is measured by fitting the Fourier transform of the measurement to a Lorentzian, and measuring its FWHM. The values are 27 and 30 Hz. As noted previously, there are at least two origins to the decay of the off-diagonal element: vibration of mirrors, and coupling through absorption (gain). The latter affects equally the diagonal and off-diagonal elements. The former can be seen as a type of “inhomogeneous broadening,” since it has its origin in random cavity length fluctuation, expressed as randomness in the value of $\Delta\omega$. The approximately 30-Hz bandwidth of both decays is consistent with 0.3- μm amplitude vibrations at 100 Hz of cavity components, causing differential cavity fluctuations of 0.3 pm/round-trip.

The role of the gain coupling is particularly important in the present Ti:sapphire laser because of the long-gain recovery time. A better candidate for this analogy would be a dye or semiconductor laser for which the gain lifetime is shorter than the cavity round-trip time. An OPO provides an even better situation, since the gain exists only at the time of pumping.

Rabi Cycling Off-Resonance

If the radiation of amplitude \mathcal{E} (Rabi frequency $\kappa\mathcal{E}$) is off-resonance with a two-level system by an amount $\Delta\omega$, the Rabi frequency becomes $\sqrt{\kappa^2\mathcal{E}^2 + \Delta\omega^2}$. In the case of the ring laser, we can control the off-resonance amount $\Delta\omega$ with a Pockel’s cell (Fig. 14), the initial condition is set favorable to the counter-clockwise direction as shown in Fig. 15a. The Rabi cycling is measured indeed to correspond to $\sqrt{|\tilde{r}|^2/\tau^2 + \Delta\omega^2}$. In resonance case ($\Delta\omega = 0$), measurement of ρ_{12} leads to $r/\tau_{RT} = 138 \text{ Hz} \pm 15 \text{ Hz}$. With $\Delta\omega$ of $171 \text{ Hz} \pm 12 \text{ Hz}$, the off-resonant measurement is $r/\tau_{RT} = 237 \text{ Hz} \pm 21 \text{ Hz}$, which follows the behavior of an off-resonance two-level system.

Impact of the Analogy

The analogy between two-level and laser systems may be more than just a scientific curiosity. A thorough understanding of two-level systems led to powerful spectroscopic techniques, using sophisticated

pulse sequences. Rather than using optical pulses, sublinewidth molecular spectroscopy has been successfully realized by pulsing the detuning by the Stark shifts.^{38–40} In the case of the ring laser, similar pulse sequences can be applied to the detuning. The information sought in spectroscopy is contained in the measurement of $|\rho_{12}|$, as a function of the driving field (measurement of the Rabi frequency $\kappa\epsilon$ leading to the determination of the dipole moment) or detuning $\Delta\omega$. In the case of the ring laser, the measurement of $|\rho_{12}|$ is linked to the properties of some sample inserted in the cavity.¹¹ Any resolution enhancing technique that has been devised in spectroscopy, such as the Ramsey fringes,^{41,42} could be transposed to a laser phase sensor with two intracavity pulses.

20.7 CONCLUSION

This chapter started with a mathematical description of short optical pulses and optical pulse train. Basic physics of short pulse generation in a mode-locked laser are discussed. It is shown in particular that the mechanism by which a steady-state pulse is generated inside the laser, is also responsible for creating equally spaced modes in the frequency domain. It is shown that pulse train interferometry combines the properties of temporal and spatial resolution. The laser can be used as a most sensitive interferometer, when the reference and sample pulses are circulating in the same cavity. Measurements of extreme sensitivity can be performed by interfering the two pulse trains emitted by such a laser. The exquisite sensitivity to phase results from the fact that a phase shift is transposed into a frequency shift inside the active cavity. Exploitation of such lasers as sensors requires a thorough understanding of the coupling between the two intracavity pulses. A new modeling of the laser with two intracavity pulses is introduced by making an analogy with a quantum mechanical two-level system. Beyond its physical elegance, this analogy inspires new sensitivity enhancement techniques for the use of the two-pulse per cavity laser sensor.

20.8 REFERENCES

1. Ladan Arissian and Jean-Claude Diels. "Carrier to Envelope and Dispersion Control in a Cavity with Prism Pairs." *Phys. Rev. A* **75**:013814–013824, 2007.
2. F. X. Kärtner, N. Matuschek, T. Schibli, U. Keller, H. A. Haus, C. Heine, R. Morf, V. Scheuer, M. Tilsch, and T. Tschudi. "Design and Fabrication of Double Chirped Mirrors." *Opt. Lett.* **22**:831–833, 1997.
3. J.-C. Diels and Wolfgang Rudolph. *Ultrashort Laser Pulse Phenomena*, 2d ed. Elsevier, Boston, 2006, ISBN 0-12-215492-4.
4. Th. Udem, J. Reichert, R. Holzwarth, and T.W. Hänsch. "Accurate Measurement of Large Optical Frequency Differences with a Mode-Locked Laser." *Opt. Lett.* **24**:881–883, 1999.
5. J.-C. Diels. Femtosecond Dye Lasers. In F. Duarte and L. Hillman, (eds.) *Dye Laser Principles: With Applications*, Academic Press, Boston, 1990. ISBN 0-12-215492-4, chapter 3, pages 41–132.
6. C. M. Depriest, T. Yilmaz, P. J. Delfyett, S. Etemad, A. Braun, and J. H. Abeles. "Ultralow Noise and Supermode Suppression in an Actively Mode-Locked External-Cavity Semiconductor Diode Ring Laser." *Opt. Lett.* **27**:719–721, 2002.
7. B. Resan and P. J. Delfyett. "Dispersion-Managed Breathing-Mode Semiconductor Mode-Locked Ring Laser: Experimental Characterization and Numerical Simulations." *IEEE J. of Quantum Elect.* **40**:214–220, 2004.
8. T. Yilmaz, C. M. Depriest, and P. J. Delfyett. "Complete Noise Characterisation of External Cavity Semiconductor Laser Hybridly Modelocked at 10 GHz." *Elect. Lett.* **22**:1338–1339, 2003.
9. T. Yilmaz, C. M. Depriest, A. Braun, and J. H. Abeles and P. J. Delfyett. "Residual Phase Noise and Longitudinal Mode Linewidth Measurements of Hybridly Modelocked External Linear Cavity Semiconductor Laser." *Opt. Lett.* **27**:872–874, 2002.
10. T. Yilmaz, C. M. Depriest, A. Braun, J. H. Abeles, and P. J. Delfyett. "Noise in Fundamental and Harmonic Mode-Locked Semiconductor Lasers: Experiments and Simulations." *IEEE J. of Quantum Elect.* **39**:838–849, 2003.

11. J.-C. Diels, Jason Jones, and Ladan Arissian. "Applications to Sensors of Extreme Sensitivity." In Jun Ye and Stephen Cundiff, (eds.) *Femtosecond Optical Frequency Comb: Principle, Operation and Applications*, Springer, New York, 2005, chapter 13, 333–354.
12. F. Aronowitz and R. J. Collins. "Mode Coupling due to Backscattering in a He-Ne Traveling-Wave Ring Laser." *Appl Phys. Lett.* **9**:55–58, 1966.
13. F. Aronowitz. "The Laser Gyro." In Ross, (ed.) *Laser Applications*, Academic Press, New York, 1971, 133–200.
14. J.-C. Diels, I. C. McMichael, J. J. Fontaine, and C. Y. Wang. "Subpicosecond Pulse Shape Measurement and Modeling of Passively Mode locked Dye Lasers Including Saturation and Spatial Hole Burning." In K. B. Eisenthal, R. M. Hochstrasser, W. Kaiser, and A. Laubereau, (eds.) *Picosecond Phenomena III*, Springer-Verlag, New York, 1982, 116–119.
15. M. L. Dennis, J.-C. Diels, and M. Lai. "The Femtosecond Ring Dye Laser: A Potential New Laser Gyro." *Opt. Lett.* **16**:529–531, 1991.
16. Ming Lai, Jean-Claude Diels, and Michael Dennis. "Nonreciprocal Measurements in fs Ring Lasers." *Opt. Lett.* **17**:1535–1537, 1992.
17. Ladan Arissian and Jean-Claude Diels. "Repetition Rate Spectroscopy of the Dark Line Resonance in Rubidium." *Opt. Comm.* **264**:169–173, 2006.
18. Matthew J. Bohn and Jean-Claude Diels. "Bidirectional Kerr-Lens Mode-Locked Femtosecond Ring Laser." *Opt. Comm.* **141**:53–58, 1997.
19. Scott Diddams, Briggs Atherton, and Jean-Claude Diels. "Frequency Locking and Unlocking in a Femtosecond Ring Laser with the Application to Intracavity Phase Measurements." *Appl. Phys. B* **63**:473–480, 1996.
20. Czeslaw Radzewicz, Gary W. Pearson, and Jerzy S. Krasinski. "Use of ZnS as an Additional Highly Nonlinear Intracavity Self-Focusing Element in a Ti:sapphire Self-Modelocked Laser." *Opt. Comm.* **102**:464–468, 1993.
21. M. J. Bohn, R. J. Jones, and J.-C. Diels. "Mutual Kerr-Lens Mode-Locking." *Opt. Comm.* **170**:85–92, 1999.
22. G. P. Agrawal. *Nonlinear Fiber Optics*. Academic Press, Boston, 1995, ISBN 0-12-045142-5.
23. G. E. Stedman, Z. Li, C. H. Rowe, A. D. McGregor, and H. R. Bilger. "Harmonic Analysis in a Large Ring Laser with Backscatter-Induced Pulling." *Physical Review A* **51**(6), June 1995.
24. Matthew J. Bohn, Jean-Claude Diels, and R. K. Jain. "Measuring Intracavity Phase Changes Using Double Pulses in a Linear Cavity." *Opt. Lett.* **22**:642–644, 1997.
25. Ming Lai and Jean-Claude Diels. "Wave-Particle Duality of a Photon in Emission." *J. of the Opt Soc. Am. B* **9**:2290–2294, 1992.
26. D. Gnass, N. P. Ernsting, and F. P. Schaefer. "Sagnac Effect in the Colliding-Pulse-Mode-Locked Dye Ring Laser." *Appl. Phys. B* **53**:119–120, 1991.
27. F. Krausz, Ch. Spielman, T. Brabec, E. Wintner, and A. J. Schmidt. "Generation of 33-fs Optical Pulses from a Solid-State Laser." *Opt. Lett.* **17**:204, 1992.
28. C. Spielmann, P. F. Curley, T. Brabec, and F. Krausz. "Ultrabroadband Femtosecond Lasers." *IEEE J. Quant. Elec.* **QE-30**:1100–1114, 1994.
29. A. E. Siegman. "An Antiresonant Ring Interferometer for Coupled Laser Cavities, Laser Output Coupling, Mode-Locking, and Cavity Dumping." *IEEE J. Quantum Electron.* **QE-9**:247–250, 1973.
30. N. Jamasbi, J.-C. Diels, and L. Sarger. "Study of a Linear Femtosecond Laser in Passive and Hybrid Operation." *J. of Modern Optics* **35**:1891–1906, 1988.
31. D. T. Reid, M. Padgett, C. McGowan, W. E. Sleat, and W. Sibbett. "Light-Emitting Diodes as Measurement Devices for Femtosecond Laser Pulses." *Opt. Lett.* **22**:233–235, 1997.
32. E. S. Wachman, D. C. Edelstein, and C. L. Tang. "Continuous-Wave Mode-Locked and Dispersion Compensated fs Optical Parametric Oscillator." *Opt. Lett.* **15**:136–139, 1990.
33. Xianmei Meng, Jean-Claude Diels, Dietrich Kuehlke, Robert Batchko, and Robert Byer. "Bidirectional, Synchronously Pumped, Ring Optical Parametric Oscillator." *Opt. Lett.* **26**:265–267, 2001.
34. Xianmei Meng, Raphael Quintero, and Jean-Claude Diels. "Intracavity Pumped Optical Parametric Bidirectional Ring Laser as a Differential Interferometer." *Opt. Comm.* **233**:167–172, 2004.
35. R. P. Feynman, F. L. Vernon, and R. W. Hellwarth. "Geometrical Representation of the Schroedinger Equation for Solving Maser Problems." *J. Appl. Phys.* **28**:49–52, 1957.
36. Rafael Quintero-Torres, Mark Ackerman, Martha Navarro, and Jean-Claude Diels. "Scatterometer Using a Bidirectional Ring Laser." *Opt. Comm.* **241**:179–183, 2004.

37. M. Navarro, O. Chalus, and Jean-Claude Diels. "Mode-Locked Ring Lasers for Backscattering Measurement of Mirror." *Opt. Lett.* **31**:2864–2866, 2006.
38. R. G. Brewer and R. L. Shoemaker. "Photon Echoes and Optical Nutation in Molecules." *Phys. Rev. Lett.* **27**:631–634, 1971.
39. R. G. Brewer and R. L. Shoemaker. "Optical Free Induction Decay." *Phys. Rev. A* **6**:2001–2007, 1972.
40. P. R. Berman, J. M. Levy, and R. G. Brewer. "Coherent Optical Transient Study of Molecular Collisions: Theory and Observations." *Phys. Rev.* **11**:1668–1688, 1975.
41. N. F. Ramsey. "A Molecular Beam Resonance Method with Separated Oscillating Fields." *Phys. Rev.* **78**:695–699, 1950.
42. M. M. Salour and C. Cohen-Tannoudji. "Observation of Ramsey's Interference Fringes in the Profile of Doppler-Free Two-Photon Resonances." *Phys. Rev. Lett.* **38**:757–760, 1977.

Zenghu Chang

*Department of Physics
Kansas State University
Cardwell Hall
Manhattan, Kansas*

21.1 GLOSSARY

A	electric field envelope of a laser pulse
E	kinetic energy of an electron in a laser field
ϵ_L	electric field strength of a laser pulse at a given time
E_{\max}	maximum kinetic energy of an electron
ϵ_x	electric field strength of an attosecond XUV pulse at a given time
f_0	carrier-envelope offset frequency of a frequency comb
f_{rep}	repetition frequency of a pulse train
G	temporal gate function
φ_{CE}	carrier-envelope phase (also called absolute phase) of a laser pulse
h	Planck constant
\hbar	Planck constant divided by 2π
I	laser
I_p	ionization potential of an atom
λ_0	center wavelength of a laser pulse
S	trace of the frequency-resolved optical gating
τ	time delay between a laser pulse and an attosecond XUV pulse
U_p	ponderomotive potential of an electron in a laser field
ν_c	frequency of the cutoff harmonic order
ω_0	carrier angular frequency of a laser pulse

21.2 INTRODUCTION

Since the invention of the laser in 1960, the duration of coherent optical pulses has decreased from hundreds of microseconds^{1,2} to 6 femtoseconds in the first 27 years.³ Such tremendous progress was driven by the desire to generate high peak power, study dynamics in matter, increase the speed of telecommunications, and many other applications. However, by the year 1987, the optical pulse length was approaching the limit, i.e., one optical cycle of visible light, which is a few femtoseconds. The bandwidth required to support such few-cycle pulses is generated by perturbative nonlinear interactions such as self-phase modulation.

The characteristic time scale of electron motion in atoms is one atomic unit of time, which is 24.2 attoseconds. One attosecond is 10^{-18} seconds. In the Bohr's model of the atom, the electron orbital time around the hydrogen nucleus is 152 attoseconds. The study of electron dynamics in atoms and molecules called for optical pulses with attosecond duration.⁴⁻⁶ In the frequency domain, a transform-limited Gaussian pulse with 24 attosecond full width at half maximum (FWHM) corresponds to a 73 eV FWHM power spectrum, which is much broader than the entire visible light range. In other words, attosecond pulses are inherently XUV light or x rays. The duration of such extremely short pulses was first measured in 2001.^{7,8} The required ultrabroad spectrum was obtained by using a nonperturbative nonlinear optics process called high-order harmonic generation, discovered in 1987–1988.^{9,10}

High Harmonic Generation

When a linearly polarized, short-pulse laser beam with an intensity on the order of 10^{14} W/cm² interacts with noble gases, odd harmonics of the fundamental frequency—up to tens or even hundreds in order—emerge in the output beam,^{11,12} as depicted in Fig. 1. The intensity of the first few order harmonics decreases quickly as the order increases, then the intensity remains almost unchanged over many harmonic orders, forming a plateau. Finally, the signal cuts off abruptly at the highest order. The broad width of the plateau provided the required spectral bandwidth to support attosecond pulses. The appearance of the intensity plateau is the signature of this nonperturbative laser-atom interaction, which can be described by a semiclassical model.

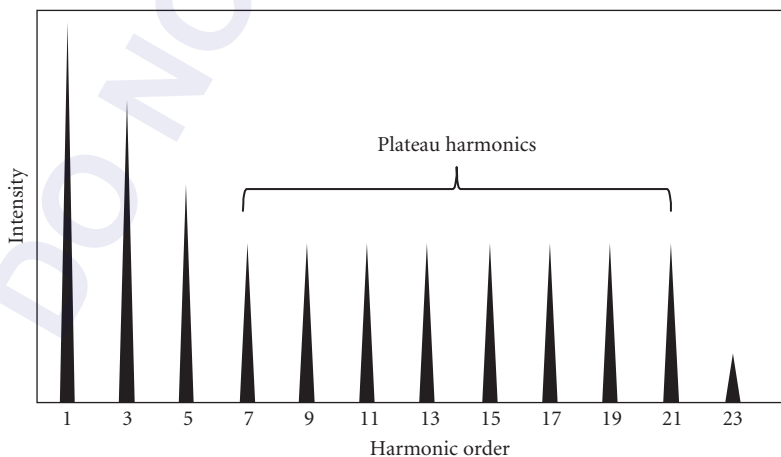


FIGURE 1 High-order harmonic spectrum.

Semiclassical Model

It is also called three-step or two-step model, rescattering model. The electric field acting on an atom changes sinusoidally within one laser cycle. As the laser intensity reaches the level of 10^{14} W/cm², the field near the peak of the oscillation is comparable to the atomic Coulomb field. The superposition of the laser field and the Coulomb field transforms the potential well that binds the electron into a potential barrier. As a result, the electron in the ground state tunnels through the barrier (the first step). The freed electron moves in the laser field like a classical particle and its trajectory can be calculated using Newton's second law. In one laser cycle, the electron first moves away from the nucleus, then is driven back when the force changes direction. During the returning journey, the electron can acquire kinetic energy up to hundreds of electron volts (the second step). Finally, the electron recombines with the parent ion with the emission of a photon (the third step).^{13,14}

When all electrons released near one peak of a laser cycle are considered, the emitted photons form an attosecond pulse. Since there are two field maxima in one laser cycle, two attosecond pulses are generated. For a laser pulse that contains many cycles, an attosecond pulse train is produced. The pulse train corresponds to discrete harmonic peaks in the frequency domain. In other words, high harmonic generation and the attosecond pulse train are two manifestations of the same nonperturbative interaction.

The photon energy of the cutoff harmonic order, $h\nu_c$, is determined by the maximum kinetic energy of the electron gained in the laser field, E_{\max} . It can be shown that $h\nu_c = I_p + E_{\max} \approx I_p + 3U_p$. Here I_p is the ionization potential of the atom and U_p is the ponderomotive potential of the electron in the laser field. Apparently, the width of the plateau and therefore the minimum attosecond pulse duration is limited by $h\nu_c$. The cutoff order is also affected by the depletion of the ground state population due to the ionization.¹⁵

Ponderomotive Potential

The ponderomotive potential is the cycle-averaged kinetic energy of an electron in a laser field, $U_p [eV] = 9.33 \times 10^{-14} I \cdot \lambda_0^2$, where I is the laser intensity in W/cm² and λ_0 is the center wavelength of the laser in micrometer. It is clear that the cutoff photon energy of the high harmonic spectrum can be extended by using longer wavelength driving lasers.¹⁶

Strong Field Approximation

A fully quantum three-step model was developed in 1994.¹⁷ It is valid when the ponderomotive potential is much larger than the ionization potential. It assumes that the harmonic emission is the result of the dipole transition between the ground state and the continuum states only, with the excitation states playing no role. Solving the Schrödinger equation results in an analytical solution of the dipole moment, from which one can obtain both the phase and the intensity of each harmonic order. The model reveals that there are two quantum trajectories that contribute to each plateau harmonic. One is called the long trajectory and the other is the short trajectory. The phase of each harmonic depends on the laser intensity. The intensity dependence of the dipole phase (also called intrinsic phase) is different for the two trajectories.

Quantum Trajectories

By solving the equation of motion, it can be shown that an electron released right at the peak of the laser field will return to the starting point one cycle later, with zero kinetic energy. As the releasing time from the field peak increases, the returning energy increases first, reaches the maximum value ($3U_p$), then decreases to zero. Therefore electrons releasing at two different moments can come back to the parent ion with the same kinetic energy, which corresponds to the same harmonic order.^{13,14}

The electron that starts the journey earlier returns later. Its path is called the long trajectory. The other one is the short trajectory. Quantum mechanically, there are many trajectories contribute to each harmonics (Feynman's path-integral), but the dominating contributions are from the two trajectories corresponding to the classical ones.¹⁷

Phase-Matching

The semiclassical model and the strong field approximation describe the single atom response. To generate a high-intensity high harmonic beam, many atoms must contribute to the output constructively.¹⁸ Ionization of the atom is unavoidable in high harmonic generation because it is the first step of the process. In highly ionized gas targets, the phase velocity of the laser field (and thus the polarization) is greater than that of the harmonic field. The resulting phase mismatch can be compensated for by several approaches. One of them utilizes the intensity dependent phase.¹⁹ In most cases, only the short trajectory is phase matched. Nevertheless, low laser to harmonic conversion efficiency is still a major problem that needs to be solved. The spatial coherence of the high harmonic/attosecond train beam is excellent when the phase matching conditions are fulfilled. The divergence angle of the XUV beam is smaller than the driving laser beam.¹²

Single Isolated Pulses

The attosecond pulse train corresponding to high order harmonics is useful for some applications. In general, however, single isolated attosecond pulses are required for performing pump-probe experiments with arbitrary delay between the pump and the probe pulses. Such pulses can be generated by suppressing all the pulses in the train except one, which can be accomplished by using single-cycle driving lasers²⁰ or pulse extraction switches with a subcycle opening time.²¹ Also the pulses from the gas target are positively chirped.²² Dispersion compensation over a broad XUV spectral range is a major challenge. By 2008, the shortest single isolated pulses, which were generated from neon gas by using 3.3-fs driving lasers centered at 720 nm, were 80 attoseconds and contained ~0.5 nJ of energy.²³ Their spectrum was centered at 80 eV.

21.3 THE DRIVING LASER

There are several basic requirements on the driving lasers for the generation of single isolated attosecond pulses. First, the intensity at the focus must be high enough, on the order of 10^{14} to 10^{15} W/cm², which is a fraction of an atomic unit of intensity (3.55×10^{16} W/cm²). The corresponding pulse energy is 100 μ J or higher. The spectral bandwidth of the attosecond pulses is proportional to the driving laser intensity. Second, the laser pulse duration must be short enough. The ionization of the target atoms by the laser field before the cycle where the attosecond pulse is generated must not deplete the ground state population completely. Depending on the generation scheme, acceptable laser pulses range from 3 to 30 fs. Third, the carrier-envelope phase needs to be stabilized. Since the single attosecond pulses are generated in a fraction of the laser cycle, a shift in the carrier-envelope phase results in shot-to-shot variations of the attosecond pulses. Finally, the repetition rate of the laser should be high, on the order of kilohertz. Many attosecond characterization and application schemes rely on photoelectron measurements. There is an upper limit on the number of electrons per shot to avoid the space charge effect. Thus the signal count rate is primarily determined by the repetition rate. The energy stability of high-repetition-rate lasers is also better than those with low repetition rates.

High power laser pulses with duration around 30 fs can be generated with chirped pulse amplification. Pulses down to ~4 fs with submillijoule energy can be obtained by spectral broadening in hollow-core fibers filled with gases, followed by dispersion compensation using chirped mirrors or phase modulators,²⁴⁻²⁷ as illustrated by the block diagram in Fig. 2.

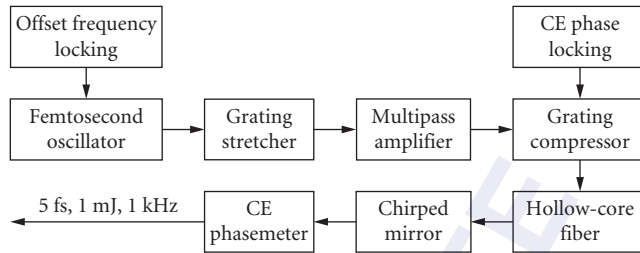


FIGURE 2 A carrier-envelope phase stabilized, few-cycle laser system.

Chirped Pulse Amplification (CPA)

Ti:Sapphire is a commonly used gain medium for femtosecond lasers primarily because of its broad gain bandwidth.²⁸ Its center wavelength is ~ 800 nm, which corresponds to a ~ 2.6 -fs optical period. Femtosecond oscillators that use Ti:Sapphire as the gain medium can generate pulses with nanojoule level energy. Direct amplification of the pulse to millijoule level may cause damage to the laser crystal. In a chirped amplifier, the pulses from the oscillator are stretched to hundreds of picoseconds to lower the peak power. Then the pulses are amplified in multipass or regenerative amplifiers. The high energy pulses are finally compressed to femtosecond duration.^{29,30} Most high energy (>5 mJ) lasers use grating pairs to stretch and compress pulses.

Carrier-Envelope Phase

The effects of the magnetic field of the laser on high harmonic generation can be ignored because the electron velocity during the three-step journey is much slower than the speed of light. The electric field of a linearly polarized laser pulse can be expressed as $\epsilon_L(t) = A(t)\cos(\omega_0 t + \varphi_{CE})$. Here $A(t)$ is the pulse envelope, ω_0 is the carrier frequency, φ_{CE} is the carrier-envelope (CE) phase that specifies the offset between the peak of the pulse envelope and the closest maximum of the oscillating field. To stabilize the carrier-envelope phase of the laser pulses from chirped pulse amplifiers followed by the hollow-core fiber compressors, the carrier-envelope offset frequency of the oscillator must first be stabilized. Furthermore, the CE phase drifts in the chirped pulse amplifier and in the hollow-core fiber compressor must be compensated.^{31,32}

Carrier-Envelope Offset Frequency

The technique for stabilizing carrier-envelope offset frequency was originally developed for frequency metrology in 2000.^{33,34} Most femtosecond oscillators used for seeding amplifiers work at a repetition rate $f_{rep} \sim 80$ MHz. The oscillator output is a femosecond pulse train that corresponds to a frequency comb. The frequency of the n th tooth of the comb is $f_0 + nf_{rep}$, where f_0 is the frequency of the zeroth tooth. The rate of carrier-envelope phase change is determined by the offset frequency f_0 , which can be stabilized by using f -to- $2f$ technology. For this to work, the laser spectrum must cover an octave. The f_0 is measured by beating the $2n$ th tooth with the frequency doubled n th tooth, that is $f_0 = 2[f_0 + nf_{rep}] - (f_0 + 2nf_{rep})$. It was found that f_0 can be stabilized by controlling the pump power to the gain medium of chirped mirror based oscillators. When f_0 is fixed to f_{rep}/m , the CE phase of every m th pulse from the oscillator is the same. It is common practice to choose $m = 4$.

Carrier-Envelope Phase of Chirped Pulse Amplifiers

The oscillator pulses with the same CE phase are switched out by a Pockels cell and sent to Ti:Sapphire amplifiers that operate at kilohertz repetition rates. When grating pairs are used to stretch and compress laser pulses for the chirped pulse amplification, a submicrometer change of separation between gratings can lead to a 2π CE phase shift.³⁵ This effect has been used to correct the slow CE phase variation introduced by the amplifier components. It was accomplished by measuring the CE phase variation after the amplifier and using the measured signal for feedback control of the grating separation.³⁶ The CE phase error of CPA systems can be controlled to <200 mrad over hours. The relative CE phase variation can be measured by a single shot f -to- $2f$ interferometer, whereas the absolute phase value can be determined by a phasemeter (discussed below), which measures electrons from the above-threshold ionization of atoms by the laser pulses.

Single Shot f -to- $2f$ Interferometer

The laser pulse from the hollow-core fiber compressor is a white-light continuum that can cover an octave spectral range. One can select a narrow range near 1000 nm and frequency double it to interfere with the light around 500 nm. The interferogram in the frequency domain is a sinusoidal fringe. The period of the fringe pattern is inversely proportional to the delay between the two interfering pulses. A CE phase shift will cause the fringes to shift. Thus by measuring the interferogram with a spectrometer, the CE phase variation can be measured.³⁷

Carrier-Envelope Phasemeter

In the three-step semiclassical model, the attosecond photon pulses are generated by the recombination of the returning electrons. A returning electron can also scatter away from the parent ion. The kinetic energy distribution of the rescattered electrons after the laser field vanishes can extend to $10U_p$. There is also a plateau in the electron spectrum similar to the high harmonic spectrum. This electron emission process is called above-threshold ionization. The angular distribution of the electrons is concentrated along the field polarization direction. When the laser pulse is only a few cycles long, the number of plateau electrons flying to one direction can be different from those to the opposite direction. The asymmetry depends on the carrier-envelope phase of the laser. Thus, by simultaneously measuring electrons in two directions, the absolute CE phase value can be determined.³⁸

21.4 ATTOSECOND PULSE GENERATION

A typical attosecond pulse generation setup consists of a kilohertz femtosecond Ti:Sapphire laser system, a vacuum chamber where the gas target is located, and an XUV spectrometer/attosecond streak camera that characterizes the pulses in the spectral domain and the time domain, as shown in Fig. 3. The attosecond pulses are XUV or soft x-ray light that cannot propagate in air because of high absorption. The gas density in the laser interaction region is on the order of 10^{17} to 10^{18} atoms/cm³. The interaction length is typically a few millimeters for gas cells or gas jets. It should be smaller than the Rayleigh range of the focusing laser beam, so that the carrier-envelope phase does not change significantly due to the Gouy phase shift inside the target. The target is located after the focal point to achieve good phase-matching.

Attosecond Pulse Train

Such pulses are generated with linearly polarized laser pulses that contain many optical cycles. When only the fundamental frequency is used, the spacing between two neighboring harmonic

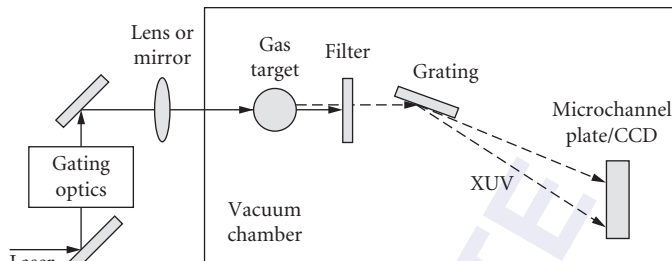


FIGURE 3 Setup for generating attosecond pulses and measuring their spectrum.

peaks is two-photon energies. In the time domain, the spacing between adjacent pulses is one-half of an optical cycle; 1.3 fs for Ti:Sapphire lasers.^{7,39} The amplitude changes from pulse to pulse. Since many-cycle lasers (>20 fs) can be generated directly from high power (terawatt) chirped pulses amplifiers, the attosecond pulse energy can be high enough to perform nonlinear physics experiments.⁴⁰

Two-Color Gating

When the many-cycle driving laser is a combination of the fundamental frequency and its second harmonic, the breaking of symmetry of the laser field leads to the generation of both odd and even high harmonics and thus the spacing between two neighboring harmonic peaks is one photon-energy. In the time domain, the spacing between adjacent pulses becomes a full optical cycle.⁴¹ Such pulse trains are useful for performing experiments using the powerful attosecond streaking technique.

Amplitude Gating

As the driving laser approaches a single optical cycle, the cycle-to-cycle field amplitude variation becomes significant. When the carrier-envelope phase of the pump laser is set to zero, the spectrum of the attosecond pulses generated near the peak of the laser pulse envelope extends to a shorter XUV wavelength range as compared to the adjacent attosecond pulses emitted when the laser field is weaker. As a result, the high-order harmonic spectrum becomes a continuum in the cutoff region.²⁰ Discrete harmonics remain in other portions of the spectrum. The shorter the driving laser is, the broader the XUV continuum becomes. A single isolated attosecond pulse as short as 80 attoseconds was obtained by selecting the continuum region of the XUV spectrum with a high-pass filter, using <4-fs pump lasers.²¹ The scaling of the attosecond pulse energy is limited by the maximum energy of the driving laser from the hollow-core fiber compressor. Combining this type of gating with the two-color gating can relax the requirement of the laser pulse duration.

Polarization Gating

In the plateau region, four attosecond pulses are produced in one laser cycle taking into account both the long and short trajectory's contributions. However, only the two pulses from the short trajectory can be phase matched on axis. Consequently, the spacing between pulses is still half of a laser cycle. Single isolated attosecond pulses can be extracted by a scheme called polarization gating.²¹ It uses a laser field with a rapid change of ellipticity. Since XUV attosecond pulses can only be efficiently generated with linearly polarized driving fields, a single attosecond pulse is emitted if the laser field

is linearly polarized in only a short time range and elliptically polarized in the other portion of the driving pulse. The time range over which the attosecond pulse is generated is called the polarization gate. So far, single isolated XUV pulses as short as 130 attoseconds were generated with this method using 5 fs pump lasers.⁴² For the same driving laser pulse duration, polarization gating has the potential to generate shorter attosecond pulses because it can create a broader continuum in the plateau region.⁴³

Double Optical Gating

The few-cycle laser pulses used in amplitude gating and polarization gating are difficult to generate daily. A method called double optical gating was proposed to allow the generation of single isolated attosecond pulses with longer pump lasers.⁴⁴ It is a combination of the two-color gating and the polarization gating. A second harmonic field is added to the fundamental field in order to break the symmetry of the field and increases the spacing between the adjacent attosecond pulses to one optical cycle. When the polarization gating is applied, the polarization gate width equals one optical cycle to select one isolated XUV pulse. The depletion of the ground state population by the leading edge of the laser pulses can be significantly reduced with this scheme; as a result, multicycle lasers can be used. This scheme has been demonstrated with laser pulses as long as 20 fs. Since such lasers do not necessarily need hollow-core fiber compressors, they are much easier to operate. The laser pulse energy can also be much higher, which is important for the energy scaling of the attosecond pulses.

21.5 ATTOSECOND PULSE CHARACTERIZATION

Measurement of the optical pulse duration requires a temporal gate. For femtosecond lasers, nonlinear optics phenomena such as second harmonic generation can serve as the gating, which is the foundation of widely implemented autocorrelation and the frequency-resolved optical gating (FROG) techniques.⁴⁵ The intensity of the attosecond pulses is not high enough to generate second harmonic light yet. Most of the methods for determining the width of the attosecond pulses require the measurements of photoelectrons or ions. The XUV beam is focused to a gas target to generate the photoelectrons/ions. The charged particles are detected by a time-of-flight spectrometer. A second beam, either an XUV or an intense laser beam is also focused to the same target, overlapping spatially and temporally with the first beam. The interaction of the two pulses in the gas serves as the temporal gate. A typical setup is shown in Fig. 4, where the attosecond XUV pulses are generated in the first gas target and are measured in the second gas target. Similar apparatus have been used for studying electron dynamics in atoms.

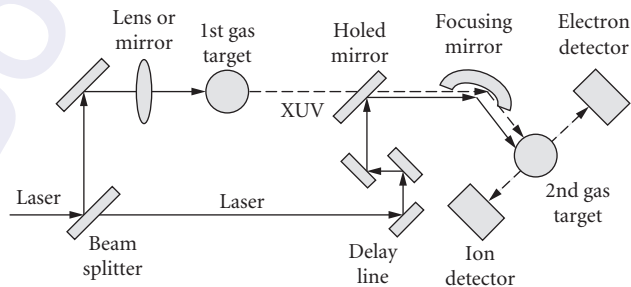


FIGURE 4 Setup for measuring the attosecond pulse duration.

Second-Order Autocorrelator or FROG

This technique resembles the second harmonic autocorrelation in femtosecond optics. Ionization of atoms (such as helium) or Coulomb explosion of molecules (such as N_2) by nonresonant two-photon absorption can serve as the nonlinearity. The ion signal as a function of the time delay between the two attosecond pulses is measured to yield the second order autocorrelation function.⁴⁰ By assuming a certain pulse shape, the pulse duration can be obtained by fitting the autocorrelation trace. The interferometric second-order autocorrelations have been used to characterize attosecond pulse trains generated with low-repetition rate femtosecond lasers with tens of millijoule pulse energy.

When the photoelectron kinetic energy spectrum is measured as a function of the delay, a two-dimensional frequency-resolved optical gating pattern is obtained. Both the phase and pulse profile of subfemtosecond pulses can be reconstructed using this method.⁴⁶

RABITT (Reconstruction of Attosecond Beating by Interference of Two-Photon Transition)

When the attosecond pulse train is generated with the fundamental wave of kilohertz lasers, the intensity of the XUV light may not be strong enough to cause measurable nonlinear effects. A cross-correlation method based on the two-color above-threshold ionization was developed to determine the duration of the pulses in the train.⁷ A high harmonic beam interacting with atomic gases alone will generate photoelectron peaks separated by two laser photon-energies. Adding a dressing laser with intensity of 10^{11} W/cm² generates an electron sideband located in the middle of two peaks. By measuring sideband intensity as a function of the delay between the XUV pulse and the dressing laser, relative phase between adjacent harmonics can be determined. Combining this with the high harmonic power spectrum, one can deduce the attosecond pulse duration. This approach assumes that the width of all the pulses in the train is the same.

Attosecond Streak Camera

The photoelectron replicas generated by attosecond XUV pulses have durations shorter than the optical cycle of the driving lasers. When the photoelectrons are released in the presence of a laser field, their momentum after the laser pulse is gone will be different from the initial value. The momentum shift is determined by the vector potential of the laser field at the time the electron is released. Thus, the leading edge of the electron pulse will gain an additional momentum that is different from the electron in the trailing edge. By measuring the momentum distribution of photoelectrons, the width of the photoelectron pulse (and thus the XUV pulse) can be determined.^{47,48} The required laser intensity is on the order of 10^{12} to 10^{13} W/cm². This approach is similar to the picosecond optical streak camera. It has been used to measure single isolated attosecond pulses and the pulse trains generated from the two-color gating. It is, however, difficult to measure pulses with half laser cycle spacing with this method.

FROG-CRAB (Frequency-Resolved Optical Gating for Complete Reconstruction of Attosecond Bursts)

The momentum streaking of the photoelectrons in a laser field can also be described as the phase shift of an electron wave packet. The phase shift of the electron wave by the laser field can be considered as a temporal phase gate, $G(t)$. When the photoelectron spectrum is measured as a function of the delay τ between the XUV field $\epsilon_x(t)$ and the laser field, a FROG trace is obtained, given by
$$S(E, \tau) = \left| \int_{-\infty}^{\infty} dt \epsilon_x(t - \tau) G(t) e^{j(E + I_p)t/\hbar} \right|^2.$$
 Here E is the energy of the photoelectron. Such a spectrogram can be processed using a FROG retrieval algorithm to fully characterize the XUV pulse as well as the electric field of the near IR laser pulse.^{49,50} This method works well for measuring both attosecond pulse trains and single isolated pulses.

21.6 ACKNOWLEDGMENTS

This material is supported by the U.S. Army Research Office under grant number W911NF-07-1-0475.

21.7 REFERENCES

1. T. H. Maiman, "Stimulated Optical Radiation in Ruby," *Nature* **187**: 493–494 (1960).
2. A. E. Siegman, *Lasers*, University Science Books, Mill Valley, California (1986), ISBN 0-935702-11-3, p. 61.
3. R. L. Fork, C. H. B. Cruz, P. C. Becker, and C. V. Shank, "Compression of Optical Pulses to Six Femtoseconds by Using Cubic Phase Compensation," *Opt. Lett.* **12**: 483–485 (1987).
4. P. Agostini, and L. F. DiMauro, "The Physics of Attosecond Light Pulses," *Reports on Progress in Physics* **67**: 813 (2004).
5. P. B. Corkum and F. Krausz, "Attosecond Science," *Nat. Phys.* **3**: 381 (2007).
6. M. F. Kling and M. J. J. Vrakking, "Attosecond Electron Dynamics," *Annual Review of Physical Chemistry* **59**: 463 (2008).
7. P. M. Paul, E. S. Toma, P. Breger, G. Mullot, F. Auge, Ph. Balcou, H. G. Muller, and P. Agostini, "Observation of a Train of Attosecond Pulses from High Harmonic Generation," *Science* **292**: 1689 (2001).
8. M. Hentschel, R. Kienberger, Ch. Spielmann, G. A. Reider, N. Milosevic, T. Brabec, P. Corkum, U. Heinzmann, M. Drescher, and F. Krausz, "Attosecond Metrology," *Nature* **414**: 509 (2001).
9. A. McPherson, G. Gibson, H. Jara, U. Johann, T. S. Luk, I. A. McIntyre, K. Boyer, and C. K. Rhodes, "Studies of Multiphoton Production of Vacuum-Ultraviolet Radiation in the Rare Gases," *J. Opt. Soc. Am. B* **4**: 595 (1987).
10. M. Ferray, A. L'Huillier, X. F. Li, L. A. Lompré, G. Mainfray, and C. Manus, "Multiple-Harmonic Conversion of 1064 nm Radiation in Rare Gases," *J. Phys. B* **21**: L31 (1988).
11. A. L'Huillier, T. Auguste, Ph. Balcou, B. Carfé, P. Monot, P. Salières, C. Altucci, et al., "High-Order Harmonics: A Coherent Source in the XUV Range," *J. Nonl. Opt. Phys. and Mat.* **4**: 647 (1995).
12. P. Salières, A. L'Huillier, P. Antoine, M. Lewenstein, "Studies of the Spatial and Temporal Coherence of High Order Harmonics," *Adv. Atom. Mol. Opt. Phys.* **41**: 83 (1999).
13. P. B. Corkum, "Plasma Perspective on Strong-Field Multiphoton Ionization," *Phys. Rev. Lett.* **71**: 1994–1997 (1993).
14. K. C. Kulander, K. J. Schafer, J. L. Krause, in *Super-Intense Laser-Atom Physics*, B. Piraux, A. L'Huillier, and K. Rzazewski (eds.) Plenum, New York (1993). NATO ASI, Ser. B, Vol. **316**: p. 95.
15. Z. Chang, A. Rundquist, H. Wang, M. M. Murnane, and H. C. Kapteyn, "Generation of Coherent Soft X Rays at 2.7 nm Using High Harmonics," *Phys. Rev. Lett.* **79**: 2967 (1997).
16. B. Shan, Z. Chang, "Dramatic Extension of the High-Order Harmonic Cutoff by Using a Long-Wavelength Pump," *Phys. Rev. A* **65**: 011804(R) (2002).
17. M. Lewenstein, Ph. Balcou, M. Yu. Ivanov, A. L'Huillier, and P. B. Corkum, "Theory of High-Harmonic Generation by Low-Frequency Laser Fields," *Phys. Rev. A* **49**: 2117–2132 (1994).
18. M. B. Gaarde, J. L. Tate, and K. J. Schafer, "Macroscopic Aspects of Attosecond Pulse Generation," *J. Phys. B: At. Mol. Opt. Phys.* **41**: 32001 (2008).
19. M. Lewenstein, P. Salières, and A. L'Huillier "Phase of the Atomic Polarization in High-Order Harmonic Generation," *Phys. Rev. A* **52**: 4747 (1995).
20. I. P. Christov, M. M. Murnane, and H. Kapteyn, "High-Harmonic Generation of Attosecond Pulses in the 'Single-Cycle' Regime," *Phys. Rev. Lett.* **78**: 1251–1254 (1997).
21. P. B. Corkum, N. H. Burnett, and M. Y. Ivanov, "Subfemtosecond Pulses," *Opt. Lett.* **19**: 1870 (1994).
22. Z. Chang, "Chirp of the Attosecond Pulses Generated by a Polarization Gating," *Phys. Rev. A* **71**: 023813 (2005).
23. E. Goulielmakis, M. Schultze, M. Hofstetter, V. S. Yakovlev, J. Gagnon, M. Uiberacker, A. L. Aquila, et al., "Single-Cycle Nonlinear Optics," *Science* **320**: 1614 (2008).
24. M. Nisoli, S. D. Silvestri, and O. Svelto, "Generation of High Energy 10 fs Pulses by a New Pulse Compression Technique," *Appl. Phys. Lett.* **68**: 2793–2975 (1996).

25. R. Szipöcs, K. Ferencz, C. Spielmann, and F. Krausz, "Chirped Multilayer Coatings for Broadband Dispersion Control in Femtosecond Lasers," *Opt. Lett.* **19**: 201–203 (1994).
26. M. Nisoli, S. D. Sverstri, O. Svelto, R. Szipöcs, K. Ferencz, Ch. Spielmann, S. Sartania, and F. Krausz, "Compression of High-Energy Laser Pulse below 5 fs," *Opt. Lett.* **22**: 522–524 (1997).
27. H. Wang, Y. Wu, C. Li, H. Mashiko, S. Gilbertson, and Z. Chang, "Generation of 0.5 mJ, Few-Cycle Laser Pulses by an Adaptive Phase Modulator," *Opt. Exp.* **16**: 14448–14455 (2008).
28. P. F. Moulton, "Spectroscopic and Laser Characteristics of $\text{Ti:Al}_2\text{O}_3$," *J. Opt. Soc. Am. B* **3**: 125 (1986).
29. D. Strickland and G. Mourou, "Compression of Amplified Chirped Optical Pulses," *Opt. Commun.* **56**: 219 (1985).
30. G. A. Mourou, T. Tajima, and S. V. Bulanov, "Optics in the Relativistic Regime," *Rev. Mod. Phys.* **78**: 309 (2006).
31. A. Baltuska, Th. Udem, M. Uiberacker, M. Hentschel, E. Goulielmakis, Ch. Gohle, R. Holzwarth, et al., "Attosecond Control of Electronic Processes by Intense Light Fields," *Nature* **421**: 611 (2003).
32. A. Baltuska, M. Uiberacker, E. Goulielmakis, R. Kienberger, V. S. Yakovlev, T. Udem, T. W. Hänsch, and F. Krausz, "Phase-Controlled Amplification of Few-Cycle Laser Pulses," *IEEE J. Sel. Topics Quantum Electron.* **9**: 972 (2003).
33. D. J. Jones, S. A. Diddams, J. K. Ranka, A. Stentz, R. S. Windeler, J. L. Hall, and S. T. Cundiff, "Carrier-Envelope Phase Control of Femtosecond Mode-Locked Lasers and Direct Optical Frequency Synthesis," *Science* **288**: 635–639 (2000).
34. A. Apolonski, A. Poppe, G. Tempea, C. Spielmann, T. Udem, R. Holzwarth, T. W. Hänsch, and F. Krausz, "Controlling the Phase Evolution of Few-Cycle Light Pulses," *Phys. Rev. Lett.* **85**: 740–743 (2000).
35. Z. Chang, "Carrier Envelope Phase Shift Caused by Grating-Based Stretchers and Compressors," *Appl. Opt.* **45**: 8350(2006).
36. C. Li, E. Moon, and Z. Chang, "Carrier-Envelope Phase Shift Caused by Variation of Grating Separation," *Opt. Lett.* **31**: 3113–3115 (2006).
37. M. Kakehata, H. Takada, Y. Kobayashi, K. Torizuka, Y. Fujihara, T. Homma, and H. Takahashi, "Single-Shot Measurement of Carrier-Envelope Phase Changes by Spectral Interferometry," *Opt. Lett.* **26**: 1436–1438 (2001).
38. G. G. Paulus, F. Grabson, H. Walther, P. Villorosti, M. Nisoli, S. Stagira, E. Priori, and S. De Silvestri, "Absolute-Phase Phenomena in Photoionization with Few-Cycle Laser Pulses," *Nature* **414**: 182–184, (2001).
39. P. Antoine, A. L'Huillier, and M. Lewenstein, "Attosecond Pulse Trains Using High-Order Harmonics," *Phys. Rev. Lett.* **77**, 1234 (1996).
40. P. Tzallas, D. Charalambidis, N. A. Papadogiannis, K. Witte, and G. D. Tsakiris, "Direct Observation of Attosecond Light Bunching," *Nature* **426**: 267–271 (2003).
41. J. Mauritsson, P. Johnsson, E. Gustafsson, A. L'Huillier, K. J. Schafer, and M. B. Gaarde, "Attosecond Pulse Trains Generated Using Two Color Laser Fields," *Phys. Rev. Lett.* **97**: 013001 (2006).
42. G. Sansone, E. Benedetti, F. Calegari, C. Vozzi, L. Avaldi, R. Flammini, L. Poletto, et al., "Isolated Single-Cycle Attosecond Pulses," *Science* **314**: 443 (2006).
43. Z. Chang, "Single Attosecond Pulse and xuv Supercontinuum in the High-Order Harmonic Plateau," *Phys. Rev. A* **70**: 043802 (2004).
44. H. Mashiko, S. Gilbertson, C. Li, S. D. Khan, M. M. Shakya, E. Moon, and Z. Chang, "Double Optical Gating of High-Order Harmonic Generation with Carrier-Envelope Phase Stabilized Lasers," *Phys. Rev. Lett.* **100**: 103906 (2008).
45. R. Trebino, D. J. Kane, "Using Phase Retrieval to Measure the Intensity and Phase of Ultrashort Pulses: Frequency-Resolved Optical Gating," *J. Opt. Soc. Am. A* **10**: 1101 (1993).
46. A. Kosuge, T. Sekikawa, X. Zhou, T. Kanai, S. Adachi, and S. Watanabe, "Frequency-Resolved Optical Gating of Isolated Attosecond Pulses in the Extreme Ultraviolet," *Phys. Rev. Lett.* **97**: 263901 (2006).
47. R. Kienberger, M. Hentschel, M. Uiberacker, Ch. Spielmann, M. Kitzler, A. Scrinzi, M. Wieland, et al., "Steering Attosecond Electron Wave Packets with Light," *Science* **297**: 1144–1148 (2002).
48. J. Itatani, F. Quéré, G. L. Yudin, M. Yu. Ivanov, F. Krausz, and P. B. Corkum, "Attosecond Streak Camera," *Phys. Rev. Lett.* **88**: 173903 (2002).
49. Y. Mairesse, and F. Quéré, "Frequency-Resolved Optical Gating for Complete Reconstruction of Attosecond Bursts," *Phys. Rev. A* **71**: 011401(R) (2005).
50. J. Gagnon, E. Goulielmakis, V. S. Yakovlev, "The Accurate FROG Characterization of Attosecond Pulses from Streaking Measurements," *App. Phys. B* **92**: 25 (2008).

This page intentionally left blank.

DO NOT DUPLICATE

LASER STABILIZATION

John L. Hall, Matthew S. Taubman*, and Jun Ye

JILA

*University of Colorado and National Institute of Standards and Technology
Boulder, Colorado*

22.1 INTRODUCTION AND OVERVIEW

For laser applications in which measurement precision is a key feature, frequency-stabilized lasers are preferred, if not essential. This observation was true in the gas laser days when the 10^{-6} fractional Doppler width set the uncertainty scale. Now we have diode-pumped solid state lasers with fractional tuning range approaching 10^{-2} or more, and laser diode systems with several percent tuning. Such tuning is useful to find the exact frequency for our locking resonance, but then stabilization will be essential. Locking to cavities and atomic references can provide excellent stability, even using a widely tunable laser source. Indeed, laser frequency stability between independent systems has been demonstrated at 1×10^{-15} in 1 s averaging time, and more than a decade better at 300 seconds. This incredible performance enhancement is possible because of a feedback system, beginning from measurement of the laser's frequency error from our chosen setpoint, suitable processing of this error signal by a filter/amplifier system, and finally application of a correction signal to an actuator on the laser itself, which changes its frequency in response. While such feedback in response to performance may be the most important principle in evolution, in machines and lasers feedback enables the design of lighter, less costly systems. The accuracy is obtained, not by great bulk and stiffness, but rather by error correction, comparing the actual output against the ideal. This continuous correction will also detect and suppress the system's internal nonlinearity and noise. The performance limitation ultimately is set by imprecision of the measurement, but naturally there is a lot of care required to get into that domain: we must have a very powerful and accurate correction effort to completely hide the original sins.

This chapter is our attempt to lead the worker newly interested in frequency control of lasers on a guided tour of stabilized lasers, ideally providing enough insight for recruiting yet another colleague into this wonderful arena. As nonlinear optics becomes just part of our everyday tools, the buildup cavities which enhance the nonlinear couplings are taking on a more critical role: this is the reason that we focus on the taming of piezoelectric-based (PZT-based) systems. We then cover locking with other transducers, and present some details about their construction and use. We consider the frequency discriminator, which is a key element for these control systems. The chapter concludes with description of the design and performance of several full practical systems, including subhertz linewidth systems.

*Matthew S. Taubman is now with the Pacific Northwest National Laboratories, Richland WA.

Quantifying Frequency Stability

In thinking about the stability of our lasers, one may first wonder whether time- or frequency-domain pictures will be more powerful and instructive. Experience shows that time-domain perturbations of our lasers are usually associated with unwelcome sounds—door slamming, telephone bells, and loud voices. Eventually these time-localized troubles can be eliminated. But what remains is likely the sum of zillions of smaller perturbations: none too conspicuous, but too many in number to attack individually. This perspective leads to a frequency-domain discussion where we can add the Fourier amplitudes caused by the many little sources. Eventually we are led to idealize our case to a continuum of spectrally-described perturbations. This physical outlook is one reason we will mainly be specifying our performance measures in the frequency domain: We have already removed the few really glaring problems and now begin to see (too) many small ones.

Another important issue concerns the nifty properties of Mr. Fourier's description: in the frequency domain, cascaded elements are represented by the multiplication of their individual transfer functions. If we had chosen instead the time domain, we would need to work with convolutions, nonlocal in time. Today's result in time is the sum of all previous temporal events that have the proper delay to impact us now. So it seems clear that frequency domain is good for analysis. What about describing the results?

Frequency versus Time: Drift—the Allan Variance Method

At the other end of our laser stabilization project, describing the results, it is convenient to measure and record the frequency as a function of time. We can measure the frequency averaged over one-second gating time, for example, and stream 100 points to a file. This would be a good way to see the variations around a mean for the 1-second time intervals. This measurement could be repeated using a succession of gate times, 3 s, 10 s, 30 s, 100 s. . . . Surely it will be attractive to make this measurement just once and numerically combine the data to simulate the longer gate times. Thinking this way brings us a new freedom: we can process this data to recover more than just the mean and the standard deviation. Of course, we can expect to eventually see some drift, particularly over long times. When we look at the drift and slowly varying laser frequency, one wishes for a method to allow us to focus on the random noise effects which are still visible, even with the extended gate times. This is where the resonance physics is, while the drift is mainly due to technical problems. Dave Allan introduced the use of first differences, which has come to be called the Allan Variance method.¹ If we take the difference between adjacent samples of the measured frequency, we focus on the random processes which are averaged down to small, but not insignificant values within each gate time τ . These first differences (normalized by $1/\sqrt{2}$ to account for random noise in each entry) form a new data set which is first-order insensitive to long-term processes such as drift which dominate the directly recorded data.

Essentially the Allan Variance calculation presents us with a display of the laser's fractional frequency variation, σ_y , as a function of the time over which we are interested. At medium times, say τ of a few seconds, most laser stabilization systems will still be affected by the random measurement noise arising from shot noise and perhaps laser technical noise. At longer times the increased signal averaging implies a smaller residual fluctuation due to random processes. It is easy to show that the dependence of σ_y versus τ can be expected to be $1/\tau^{1/2}$, in the domain controlled by random (white) noise. The Allan deviation also has a great utility in compressing our statement of laser stability: we might say, for example, "the (in-)stability is 2×10^{-12} at 1 second, with the $1/\tau^{1/2}$ dependence which shows that only random noise is important out to a time of 300 s."

Allan Deviation Definition

With a counter linked to a computer, it is easy to gather a file of frequency values f_i measured in successive equal gate time intervals, t_g . Usually there is also some dead time, say t_d , while the counter-to-computer data transfers occur via the GPIB connection. This leads to a sample-to-sample time interval

of $t_s = t_g + t_d$. Allan variance is one half of the average squared difference between adjacent samples, and the usually quoted quantity, the Allan Deviation, is the square root of this averaged variance,

$$\sigma_y(\tau) = \left[\frac{1}{2(N-1)} \sum_{n=1}^{N-1} (f_{n+1} - f_n)^2 \right]^{1/2} \quad (1)$$

The dependence of σ_y upon the measuring time τ contains information essential for diagnosis of the system performance. These values for several times can be efficiently calculated from the (large) data set of frequencies observed for a fixed minimum gate time by adding together adjacent measurements to represent what would have been measured over a longer gate time. (This procedure neglects the effects of the small dead-time t_d , which are negligible for the white frequency noise $1/\sqrt{\tau}$ of usual interest but, for systems with drift and increased low-frequency noise, the dead-time effects can seriously impact the apparent results.) In any case, fewer samples will be available when the synthetic gate time becomes very long, so the uncertainty of this noise measurement increases strongly. Usually one insists on three or four examples to reduce wild variations, and so the largest synthetic gate time τ_{\max} will be taken to be the total measurement time/3. For a serious publication we might prefer 5 or 10 such synthetic measurements for the last point on the graph.

The Allan Deviation has one curiosity in the presence of a distinct sinusoidal modulation of the laser's frequency: when the gate time is 1/2 the sinusoid's period, adjacent samples will show the maximum deviation between adjacent measurements, leading to a localized peak in σ_y versus τ . Interestingly, there will be "ghosts" or aliases of this when the gate time/modulation period ratio is 1/4, 1/8, and so on. For longer gate times compared with the modulation, some fractional cycle memories can be expected also. So a clean slope of $-1/2$ for a log-log plot of σ_y versus τ shows that there is no big coherent FM process present.

Historically, Allan Variance has been valuable in locating time scales at which new physical processes must be taken into account. For example, at long times it is usual for a laser or other stable oscillator to reach a level of unchanging σ_y versus τ . We speak of this as a "flicker" floor. It arises from the interplay of two opposing trends: the first is the decreasing random noise with increasing τ (decreasing σ_y versus τ). At longer times one sees an increasing σ_y versus τ , due to drifts in the many system parameters (electronic offsets, temperature . . .), which make our lasers lock at points increasingly offset from the ideal one. If we wait long enough, ever larger changes become likely. So for several octaves of time, the combination of one decreasing and one increasing contribution leads to a flat curve. Eventually significant drift can occur even within one measurement time, and this will be mapped as a domain of rising σ_y increasing as the +1 power of τ .

It is useful to note that the frequency/time connection of the Allan Variance transformation involves very strong data compression and consequently cannot at all be inverted to recover the original data stream in the way we know from the Fourier transform pair. However in the other direction, we can obtain the Allan Deviation from the Phase Spectral Density.²

Spectral Noise Density

As noted earlier, when the number of individual contributions to the noise becomes too large to enumerate, it is convenient to move to a spectral density form of representation. To carry this idea forward, two natural quantities to use would be the frequency deviations occurring at some rate and the narrow bandwidth within which they occur. To work with a quantity that is positive definite and has additive properties, it is convenient to discuss the squared frequency deviations $\langle (f^2_N) \rangle$ which occur in a noise bandwidth B around the Fourier frequency f. This Frequency Noise Power Spectral Density, $S_f \equiv \langle (f^2_N) \rangle / B$, will have dimensions of Hz² (deviation²)/Hz (bandwidth). The summation of these deviations over some finite frequency interval can be done simply by integrating S_f between the limits of interest.

Connecting Allan Deviation and Spectral Density Sometimes one can estimate that the system has a certain spectrum of frequency variations described by $S_f(f)$, and the question arises of what Allan

Deviation this would represent. We prefer to use the Allan presentation only for experimental data. However Ref. 2 indicates the weighted transform from S_f to Allan Variance.

Connecting Linewidth and Spectral Density A small surprise is that an oscillator's linewidth generally will not be given by the summation of these frequency deviations! Why? The answer turns on the interesting properties of Frequency Modulated (FM) signals. What counts in distributing power is the Phase Modulation Index β , which is the peak modulation-induced phase shift or, equivalently, the ratio of the peak frequency excursion compared with the modulation rate. Speaking of pure tone modulation for a moment, we can write the phase-modulated field as

$$E(t) = \sin(\Omega t + \beta \sin(\omega t))$$

$$= J_0(\beta) \exp(i\Omega t) + \sum_{n=1}^{\infty} J_n(\beta) \exp(i(\Omega + n\omega)t) + \sum_{n=1}^{\infty} J_n(\beta) (-1)^n \exp(i(\Omega - n\omega)t) \quad (2)$$

where Ω is the "carrier" frequency, and $\omega = 2\pi f$ and its harmonics are the modulation frequencies. The frequency offset of one of these "sidebands," say the n th one, is n times the actual frequency of the process' frequency f . The strength of the variation at such an n th harmonic decreases rapidly for $n > \beta$ according to the Bessel function $J_n(\beta)$. We can distinguish two limiting cases.

Large excursions, slow frequency rate This is the usual laboratory regime with solid state or HeNe and other gas lasers. The dominant perturbing process is driven by laboratory vibrations that are mainly at low frequencies (5–200 Hz). The extent of the frequency modulation they produce depends on our mechanical design, basically how efficient or inefficient an "antenna" have we constructed to pick up unwanted vibrations. Clearly a very stiff, lightweight structure will have its mechanical resonances at quite high frequencies. In such case, both laser mirrors will track with nearly the same excursion, leading to small differential motion, i.e., low pickup of the vibrations in the laser's frequency. On the other hand, heavy articulated structures, particularly mirror mounts with soft springs, have resonances in the low audio band and lead to big FM noise problems. A typical laser construction might use a stiff plate, say 2 inches thick of Al or honeycomb-connected steel plates. The mirror mounts would be clamped to the plate, and provide a laser beam height of 2 inches above the plate. Neglecting air pressure variations, such a laser will have vibration-induced excursions $(\langle f_N^2 \rangle)^{1/2}$ of $\ll 100$ kHz. An older concept used low expansion rods of say 15 mm diameter Invar, with heavy Invar plates on the ends, and kinematic but heavy mirror mounts. This system may have a vibration-induced linewidth $(\langle f_N^2 \rangle)^{1/2}$ in the megahertz range. Only when the "rods" become several inches in diameter is the axial and transverse stiffness adequate to suppress the acceleration-induced forces. With such massive laser designs we have frequency excursions of tens to thousands of kilohertz, driven by low-frequency laboratory vibrations in a bandwidth $B < 1$ kHz. In this case $(\langle f_N^2 \rangle)^{1/2} \gg B$, and the resulting line shape is Gaussian. The linewidth is given by Ref. 3, $\Delta f_{\text{FWHM}} = [8 \ln(2) (\langle f_N^2 \rangle)]^{1/2} \cong 2.355 (\langle f_N^2 \rangle)^{1/2}$.

The broadband fast, small excursion limit This is the domain in which we can usually end up if we can achieve adequate servo gain to reduce the vibration-induced FM. Since the drive frequency of the perturbation is low, it is often feasible to obtain a gain above 100, particularly if we use a speedy transducer such as an acousto-optic modulator (AOM) or an electro-optic modulator (EOM). In general we will find a noise floor fixed, if by nothing else than the broadband shot noise which forms a minimum noise level in the measurement process. Here we can expect small frequency excursions at a rapid rate, $(\langle f_N^2 \rangle)^{1/2} \ll B$, leading to a small phase modulation index. If we approximate that the Spectral Noise Frequency Density $S_f = (\langle f_N^2 \rangle)/B$ is flat, with the value $S_f \text{ Hz}^2$ (deviation²)/Hz (bandwidth), then the linewidth in this domain is Lorentzian,³ with the $\Delta f_{\text{FWHM}} = \pi S_f = \pi (\langle f_N^2 \rangle)/B$.

This summary of frequency-domain measures is necessarily brief and the interested reader may find additional discussion useful.³⁻⁵ A number of powerful consequences and insights flow from reworking the above discussions in terms of a Phase Noise Power Spectral Density, $S_\phi = S_f/f^2$. The National Institute of Standards and Technology (NIST) Frequency and Time Division publishes collections

of useful tutorial and overview articles from time to time. The currently available volume² covers these topics in more detail. Vendors of rf-domain spectrum analyzers also have useful application notes.⁶

22.2 SERVO PRINCIPLES AND ISSUES^{7,8}

Bode Representation of a Servo System

We will describe our systems by transfer functions, output/input, as a function of Fourier frequency ω . We begin purely in the domain of electronics. The amplifier gain is $G(\omega)$. The electrical feedback is represented as $H(\omega)$. Both will have voltage as their physical domain, but are actually dimensionless in that they are output/input ratios. Considering that we will have to represent phase of these AC signals, both $G(\omega)$ and $H(\omega)$ will generally be complex. It will be fundamental to view these functions with their dependence on frequency, for both the amplitude and phase response.

Imagine a closed loop system with this amplifier as the forward gain $G(\omega)$ between input V_i and output V_o . Some fraction of the output is tapped off and sent back to be compared with the actual input. For more generality we will let $H(\omega)$ represent this feedback transfer ratio. The actual input, minus this sampled output will be our input to our servo amplifier $G(\omega)$. After a line of algebra we find the new gain of the closed loop—in the presence of feedback—is

$$A_{cl} = \frac{V_o}{V_i} = \frac{G(\omega)}{1 + G(\omega)H(\omega)} \quad (3)$$

A particularly instructive plot can be made for the product $G(\omega)H(\omega)$, called the “open loop gain,” which appears in the denominator. In this so-called Bode plot, the gain and phase are separately plotted. Also, from inspection of Eq. (3) we can learn one of the key advantages which feedback brings us: if the feedback factor GH were $\gg 1$, the active gain G would basically cancel out and we would be left with $A_{cl} \sim 1/H$. We imagine this feedback channel will be passive, formed from nearly ideal nondistorting components. The noise, exact value of the gain, and distortion introduced by it are seen to be nearly unimportant, according to the large magnitude of $1 + GH$. Gentle amplifier overload will lead to overtone production, but could alternatively be represented by a decrease of G with signal. Since the output doesn't depend sensitively upon G anyway, we are sure these distortion products and internally generated noise will be suppressed by the feedback. We can identify the denominator $1 + GH$ as the noise and distortion reduction factor.

What is the cost of this reduced dependence on the active components $G(\omega)$ and their defects? Basically it is that the gain is reduced and we must supply a larger input signal to obtain our desired output. For a music system one can then worry about the distortion in the preamplifier system. However, we want to make quiescent lasers, without the slightest hint of noise. So it is nice that the amplification of internal noise is reduced.

To be concrete, the circuit of Fig. 1 represents a common building block in our servo design. It also represents a simple case of feedback. We show it as a current summing input node: the subtraction at the input arises here because the sign of the gain is negative. With the nearly ideal high-gain operational amplifiers now available, $G \gg 1$ and we can closely approximate the closed loop gain by $1/H(\omega)$, yielding a flat gain above and a rising gain below some corner frequency $\omega_0 = 1/\tau_0$, with $\tau_0 = R_f C$. Remember $1/H(\omega)$ is the closed loop gain between V_o and V_i . To find the exact relationship between the signals V_s and V_o , we notice the related voltage-divider effect gives $V_i = (1 - H(\omega))V_s$, which leads to

$$\frac{V_o}{V_s} = -\frac{R_f}{R_i} \frac{(1 + j\omega/\omega_0)}{j\omega/\omega_0} = -\frac{R_f}{R_i} \frac{(1 + j\omega\tau_0)}{j\omega\tau_0} \quad (4)$$

The negative sign arises from the fact that the forward gain is negative. When the corner frequency ω_0 is chosen to be sufficiently high, we may have to consider the bandwidth issue of the OpAmp: $G(\omega)$ could

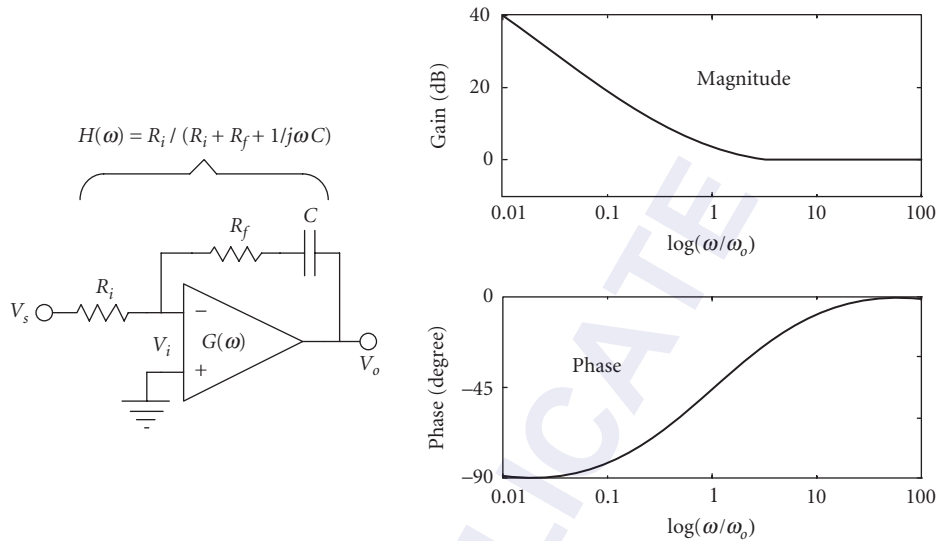


FIGURE 1 Phase and amplitude response of a Proportional-Integral (PI) amplifier circuit. The PI function is implemented using an inverting OpAmp.

start to roll off and no longer satisfy the approximation of $G \gg 1$. A more complex network is needed to compensate for the gain roll-off and that is exactly the topic of feedback we wish to cover below.

Phase and Amplitude Responses versus Frequency

We can plot⁹ the gain magnitude and phase of this elementary feedback example in Fig. 1, where we can see the flat gain at high frequencies and the rising response below ω_0 . Our laser servo designs will need to echo this shape, since the drift of the laser will be greater and greater at low frequencies, or as we wait longer. This will require larger and larger gains at low frequencies (long times) to keep the laser frequency nearby our planned lock point. The phase in Fig. 1 shows the lag approaching 90° at the lowest frequencies. (An overall minus sign is put into the subtractor unit, as our circuit shows an adder.) The time-domain behavior of this feedback system is a prompt inverted output, augmented later by the integration's contribution.

As a first step toward modeling our realistic system, Fig. 2 shows the laser included in our control loop. The servo system's job is to keep the laser output at the value defined by the reference or setpoint input. Some new issues will arise at the high-frequency end with the physical laser, as its piezo-electric transducer (PZT) will have time delay, finite bandwidth, and probably some resonances.

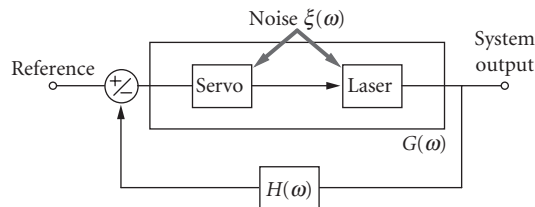


FIGURE 2 Model of laser system, including frequency noise, as part of a servo control loop.

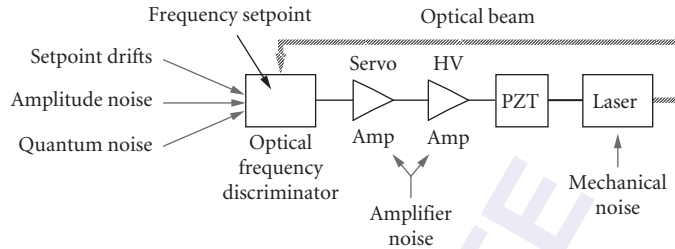


FIGURE 3 Detailed model of a frequency-controlled laser.

One way we should expand the model is to include the laser's operation as a frequency transducer, converting our control voltage to a frequency change. Probably the laser will have some unwanted frequency noises, and in Fig. 3 we can indicate their unwanted contributions arriving in the optical frequency discriminator, which functions like a summing junction. The emitted laser field encounters an optical frequency discriminator and the laser frequency is compared with the objective standard, which we will discuss below. In our diagram we show this laser frequency discriminator's role as an optical frequency-to-voltage converter element. More exactly, laser frequency *differences* from the discriminator's reference setpoint are converted to voltage outputs. Laser amplitude noises (due to the intrinsic property of the laser itself or external beam propagation) and vibration effects on the discriminator will appear as undesired additive noises also.

The first simple idea is that the feedback loop can be closed when the servo information, carried as a voltage signal in our amplifier chain, is converted to a displacement (in meters) by the PZT, then into laser frequency changes by the laser's standing-wave boundary condition. As the length changes, the "accordian" in which the waves are captive is expanded or compressed, and along with it the wavelength and frequency of the laser's light.

A second truth becomes clear as well: there is freedom in designating the division into the forward gain part and the feedback path part. Actually, we probably would like the laser to be tightly locked onto the control cavity/discriminator, and then we will tune the whole system by changing the set point, which is the discriminator's center frequency. This leads us to view the optical frequency discriminator as the summing junction, with the amplifier and PZT transducer as the forward gain part. The output is taken as an optical frequency, which would be directly compared to the setpoint frequency of the discriminator. So the feedback path $H = 1$.

We should consider some magnitudes. Let K_{PZT} represent the tuning action of the PZT transducer, expressed as displacement meter per volt. A typical value for this would be $K_{PZT} = 0.5 \text{ nm/V}$. The laser tunes a frequency interval $c/2L$ for a length change by $\lambda/2$, so the PZT tuning will be $\sim 600 \text{ V/order}$ at 633 nm.

$$K_V = K_{PZT} \frac{2}{\lambda} \frac{c}{2L} \quad (5)$$

So we obtain a tuning sensitivity $K_V \sim 800 \text{ kHz/V}$ tuning for a foot-long laser, assuming a disk-type PZT geometry. See the section below on PZT design.

Measurement Noise as a Performance Limit—It Isn't

Usually our desire for laser stability exceeds the range of the possible by many orders, and we soon wonder about the ultimate limitations. Surely the ultimate limit would be due to measurement noise. However, we rarely encounter the shot-noise-limited case, since the shot noise-limited S/N of a $100 \mu\text{W}$ locking signal is $\sim 6 \times 10^6$ in a 1 Hz bandwidth. (See section on "The Optical Cavity-Based Frequency Discriminator" later.) Rather we are dealing with the laser noise remaining because our

servo gain is inadequate to reduce the laser's intrinsic noise below the shot noise limit, the clear criterion of gain sufficiency. So our design task is to push up the gain as much as possible to reduce the noise, limited by the issue of stability of the thus-formed servo system.

Closed-Loop Performance Expectations When Transducer Resonance Limits the Usable Gain

Servo Stability: Larger Gain at Lower Frequencies, Decreasing to Unity Gain and Below Our need for high gain is most apparent in the low-frequency domain ~ 1 kHz and below. Vibrations abound in the dozens to hundreds of hertz domain. Drifts can increase almost without limit as we wait longer or consider lower Fourier frequencies. Luckily, we are allowed to have more gain at low frequencies without any costs in stability. At high frequencies, it is clear we will not help reduce our noise if our correction is applied too late and so no longer has an appropriate phase. One important way we can characterize the closed-loop behavior of our servo is by a time delay t_{delay} . Here we need to know the delay time before any servo response appears; a different (longer) time characterizes the $1/e$ response. The latter depends on the system gain, while the ultimate high-speed response possible is controlled by the delay until the first action appears. A good criterion is that the useful unity gain frequency can be as high as $f_{\tau} = 1/(2\pi t_{\text{delay}})$, corresponding to 1 rad extra phase-shift due to the delay. Below this ultimate limit we need to increase the gain—increase it a lot—to effectively suppress the laser's increased noise at low frequencies. This brings us to address the closed-loop stability issue.

Closed-Loop Stability Issues

One can usefully trace the damping of a transient input as it repetitively passes the amplifier and transducer, and is reintroduced into the loop by the feedback. Evidently stability demands that the transient is weaker on each pass. The settling dynamics will be more beautiful if the second-pass version of the perturbation is reduced in magnitude and is within say $\pm 90^\circ$ of the original phase. Ringing and long delay times result when the return phasor approaches -1 times the input signal vector, as then we are describing a sampled sinewave oscillation. These time-domain pictures are clear and intuitive, but require treatment in terms of convolutions, so here we will continue our discussion from the frequency-domain perspective that leads to more transparent algebraic forms. We can build up an arbitrary input and response from a summation of sinusoidal inputs. This leads to an output as the sum of corresponding sinusoidal outputs, each including a phase shift.

In our earlier simple laser servo example, no obvious limitation of the available closed-loop gain was visible. The trouble is we left out two fundamental laboratory parasites: time delay, as just noted, and mechanical resonances. We will usually encounter the mechanical resonance problem in any servo based on a PZT transducer. For design details, see the "Practical Issues" section. A reasonable unit could have its first longitudinal resonance at about 25 kHz, with a $Q \sim 10$. In servo terms, the actual mechanical PZT unit gives an added 2-pole roll-off above the resonance frequency and a corresponding asymptotic phase lag of 180° . Including this reality in our model adds another transfer function $R_{\text{PZT}} = \omega_0^2 / (\omega_0^2 + 2\omega\eta\omega_0 + \omega^2)$, where ω_0 is 2π times the resonance frequency, and $\eta = 1/2Q$ is the damping factor of the resonance. This response is shown in Fig. 4.

We now talk of stabilizing this system. The elements are the laser and some means to correct its frequency, a frequency discriminator to measure the difference between the actual and the setpoint frequencies, and a feedback amplifier. Here we propose to do the frequency control by means of a PZT transducer to change the laser frequency. For the present discussion, we assume the frequency discriminator has a flat response. For the feedback amplifier, the first appealing option is to try a pure integrator. The problem then is that we are limited in gain by the peakheight of the resonance which must remain entirely below unity gain to avoid instability. In Fig. 5 case (a) we see that the unity gain frequency is limited to a value of 1.5 kHz. Some margin is left to avoid excessive ringing near the resonant frequency, but it is still visible in the time domain. Techniques that help this case include a roll-off filter between the unity gain and PZT resonance frequencies.

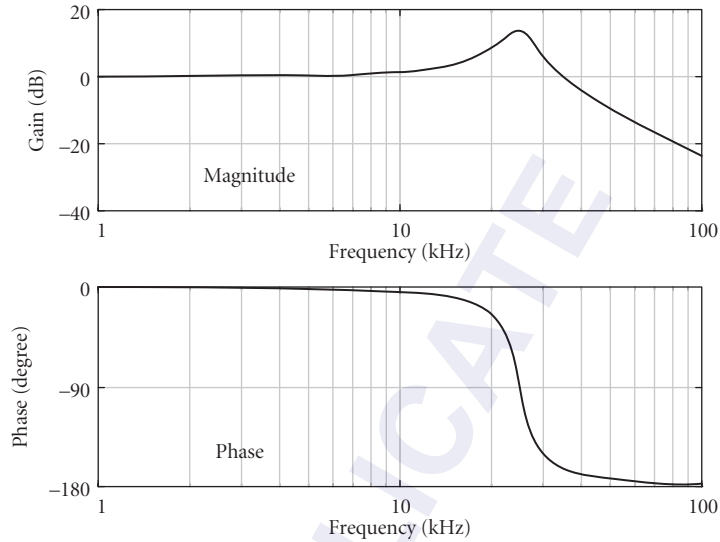


FIGURE 4 The amplitude and phase response of a tubular PZT transducer and an 8-mm-diameter by 5-mm-thick mirror. The resonance is at 25 kHz with a Q of 10.

Figure 5 shows the “open loop” gain function GH of the feedback equation, and the corresponding phase response. We already noted the dangerous response function of -1 where the denominator of Eq. (3) vanishes. In the time-domain iterative picture, the signal changes sign on successive passes and leads to instability/oscillation. We need to deal with care as we approach near this point in order to obtain maximum servo gain: it is useful to consider two stability margins. The *phase stability margin* is the phase of the open-loop function when the gain is unity. It needs to be at least 30° . The *gain*

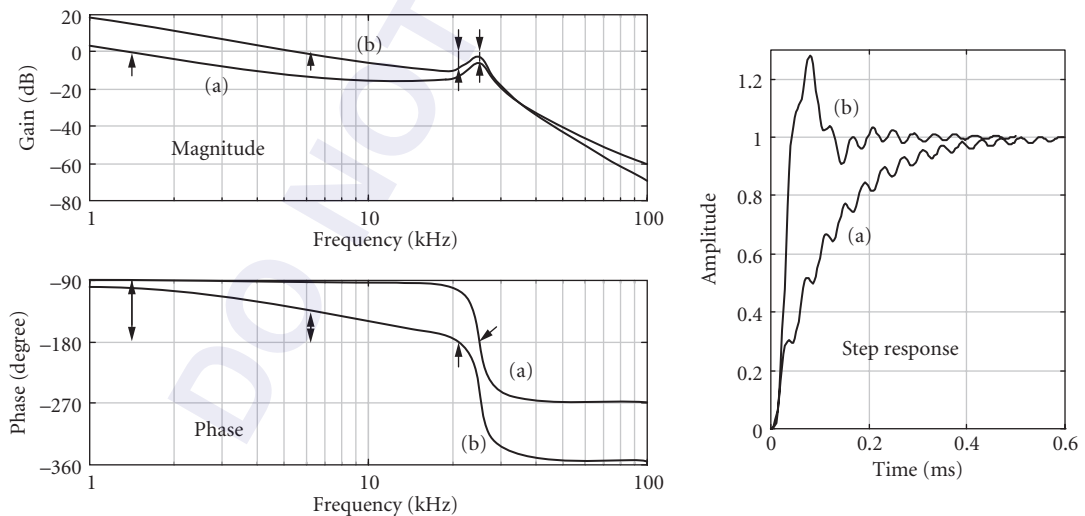


FIGURE 5 (a) Integrator gain function alone. Gain must be limited so that gain is <1 even at the resonance. (b) Single-pole low pass at 6 kHz inserted. Now unity gain can increase to 6 kHz and time response is ~ 3 -fold faster. Small arrows in the graph indicate the phase margin at the unity gain frequency (gain = 0 dB) and gain margin at a phase shift of -180° .

margin is the closed loop gain when the phase is 180° . In Fig. 5 case (a) we see that the phase is not shifted very much until we really “sense” the amplitude increase from the resonance. So this resonance may tend to fix an apparently solid barrier to further servo improvement. But as shown in Fig 5 case (b), just a low-pass to push down the PZT resonance is very helpful.

In fact, there are many ways of improving the low frequency gain of this system. They include: (1) imposing yet another high frequency roll-off (or multi-pole low pass filter) just before the resonance thus pushing its height down and allowing the open loop transfer function to come up, (2) adding lag compensators before the resonance to push the low frequency gain up while keeping the high frequency response relatively unchanged, (3) adding a lead compensator just above the resonance to advance the phase and increase the unity gain point, (4) or placing a notch at the resonant frequency to “cut it out” of the open loop transfer function. The last two options in this list are quite promising and are discussed in more detail below.

Proportional Integral Derivative (PID) Controller versus Notch Filters Like many “absolute” barriers, it is readily possible to shoot ahead and operate with a larger closed loop bandwidth than that represented by the first PZT resonance. The issue is that we must control the lagging phase that the resonance introduces. A good solution is a differentiator stage, or a phase lead compensator, which could also be called a high frequency boost/gain-step circuit. In Fig. 6 case (a) we show the Bode plot of our PZT-implemented laser frequency servo, based on a PID (Proportional Integral Differentiator) controller design. Just a few moments of design pay a huge benefit, as the unity gain frequency has now been pushed to 40 kHz, almost a factor of 2 *above* the PZT’s mechanical resonance. For this PID controller example, unity gain occurs at a 7-fold increased frequency compared with Fig. 5 case (b). Thus at the lower frequencies we would hope to have increased the servo gain by a useful factor of 7x or 17 dB. However, comparison of Fig. 5 (b) and 6 (a) shows that the low frequency gain is hardly changed, even though we greatly increased the servo bandwidth.

So, how *do* we go forward? We could in principle continue to increase the gain and unity gain frequency, but this is not really practical, however, since we will again be limited by additional structure resonances that exist beyond the first resonance. Also, the Derivatives needed to tame these resonances cost low frequency gain, and it is hard to win overall system performance. To make

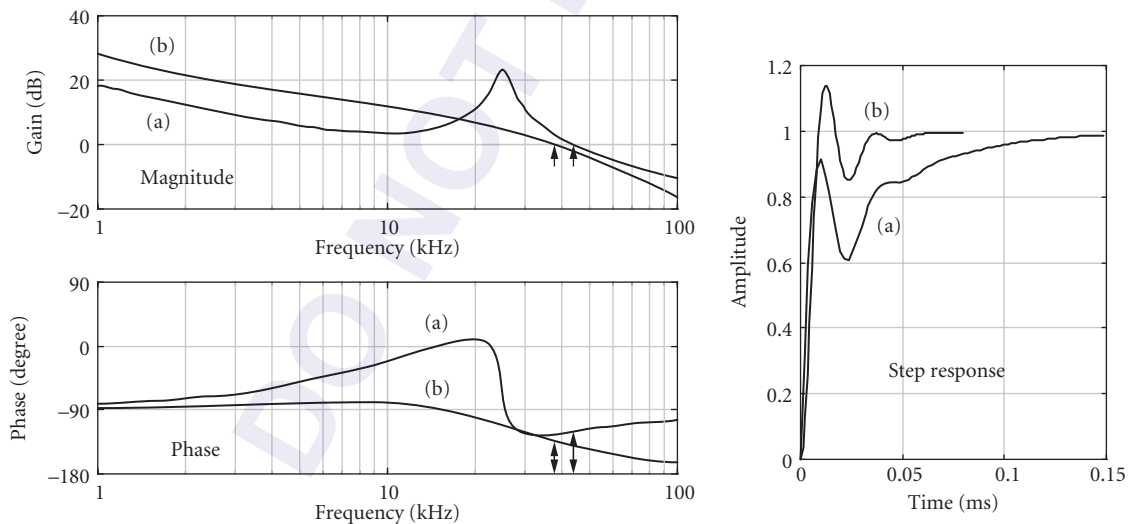


FIGURE 6 Two methods of working through and beyond a resonance. (a) PID controller where the Derivative term advances the phase near the resonance. (b) Adding a notch is a better approach, where the notch function approximates the inverse of the resonance peak. Transient response settles much more quickly. Again, we use the small arrows in the graph to indicate the phase margin at the unity gain frequency.

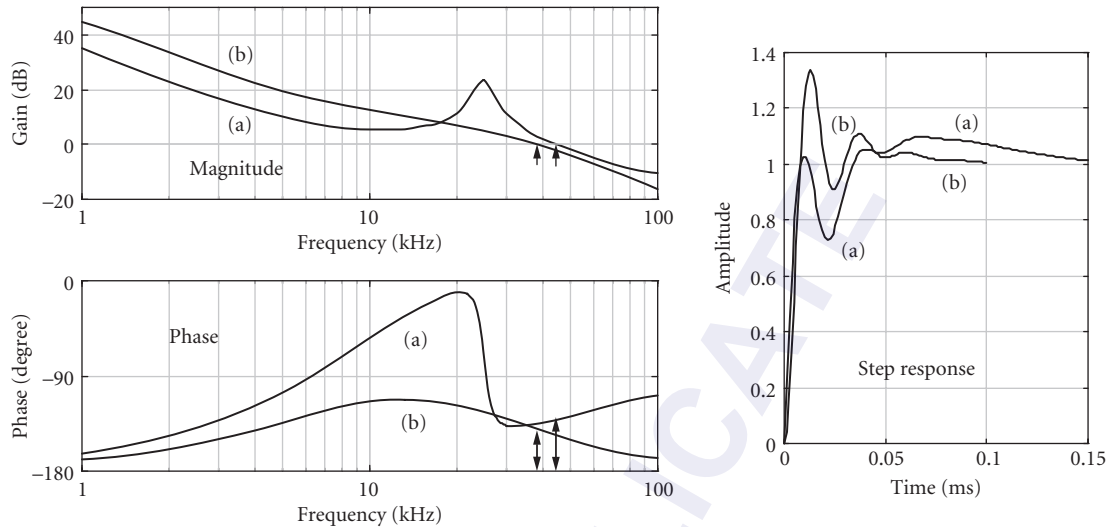


FIGURE 7 Adding an additional PI stage to (a) the PID and (b) the PI-plus-notch stabilizers of Fig. 6. Note that low frequency gain is strongly increased.

progress, we use a notch as an alternative technique to suppress the resonance. Now a D term is not needed, and we can conserve the gain at low frequencies. The notch filter, combined with a PI stage, gives unity gain at higher frequencies, and increases gain for low ones. See Fig. 6 case (b). Then Fig. 7 compares adding another PI stage to the two cases of Fig. 6, adding a PID in (a) and a Notch plus PI in (b). The time-domain approach, shown in Fig. 7, shows case (b) settles rather nicely. And the gain has increased more than 20 dB at frequencies of 1 kHz and below. So this is very encouraging.

While we have come to the cascaded-integrators approach cautiously in this discussion, in fact at least 2 integrators would always be used in practice. Workers with serious gain requirements, for example, the LIGO and VIRGO interferometric gravitational wave detector groups, may use the equivalent of 4 cascaded integrators! Such a design is “conditionally stable” only, meaning that the gain cannot be smoothly reduced or increased. Such aggressive stabilizer designs have their place, but not for a first design!

“Rule of Thumb” PID Design for System with a Transducer Resonance Optimizing servo performance is an elegant art, turned into science by specification of our “cost function” for the system performance shortcomings. In the case that we wish to minimize the time-integrated magnitude of the residuals following a disturbance, one comes to the case studied by Ziegler and Nichols for the PID controller used in a system with a combined roll-off and time delay.⁸ Such a case occurs also in thermal controllers. With only the P term, one first looks for the frequency f_{osc} where the system first oscillates when the gain is increased. The PD corner is then set $1.27\times$ higher than this f_{osc} , the P gain is reset at 0.6 of the oscillation gain, and the PI corner is set at 0.318 times the oscillation frequency. This “rule of thumb” design of the phase compensation produces a transient response which settles reasonably well, so as to minimize the Time Integrated Error. For phase-locking lasers, a cost function with more emphasis on long-lasting errors leads to another kind of “optimum” tuning, but with qualitatively similar results.

When a notch is used to suppress the resonance, there is no longer an anomalous gain at the resonant frequency and one is returned to the same case as in its absence. A reasonable servo approach to using two PI stages is to design with only one, achieving the desired unity gain frequency. The second PI is then added to have its corner frequency at this same point or up to 10-fold lower in frequency, depending on whether we wish the most smooth settling or need the highest feasible low frequency gain. Figures 6b and 7 show the Bode plot of such designs, along with the system’s

closed-loop transient response. An elegant strategy is to use adaptive clamping to softly turn on the extra integrator stage when the error is small enough, thus dynamically increasing the order of the controller when it will not compromise the dynamics of recovery.

22.3 PRACTICAL ISSUES

Here we offer a number of important tidbits that are useful background material for a successful application of the grand schemes discussed above.

Frequency Discriminators for Laser Locking—Overview

So far we have devoted our main effort to addressing the issues of the feedback scheme. Of equal importance is the subject of frequency reference system. After all, a good servo eliminates intrinsic noises of the plant (laser), and replaces them with the measurement noise associated with the reference system. Indeed, development of prudent strategies in high precision spectroscopy and the progress of laser stabilization have been intimately connected to each other through the years,¹⁰ with the vigorous pursuit of resolution and sensitivity resulting in amazing achievements in both fields.

To stabilize a laser, one often employs some kind of resonance information to derive a frequency/phase-dependent discrimination signal. The resonance can be of material origin, such as modes of an optical interferometer; or of natural origin, such as atomic or molecular transitions. If the desired use of a stabilized laser is to be an optical frequency standard, its long-term stability or reproducibility will be key, so the use of a natural resonance is preferred. Reproducibility is a measure of the degree to which a standard repeats itself from unit to unit and upon different occasions of operation. The ultimate reproducibility is limited to the accuracy of our knowledge of the involved transitions of free atoms or molecules. The term “free” means the resonance under study has a minimum dependence on the laboratory conditions, such as the particle moving frame (velocity), electromagnetic fields, collisions, and other perturbations. To realize these goals, modern spectroscopy has entered the realm of quantum-limited measurement sensitivities and exquisite control of internal and external degrees of freedom of atomic motions.

A careful selection of a high-quality resonance can lead to superior system performance and high working efficiency. For example, the combined product of the transition quality factor Q and the potential signal-to-noise ratio (S/N) is a major deciding factor, since this quantity controls the time scale within which a certain measurement precision (fractional frequency) can be obtained. This importance is even more obvious when one considers the waiting time for a systematic study is proportional to the inverse square of ($Q \times S/N$). A narrower transition linewidth of course also helps to reduce the susceptibility to systematic errors. The resonance line shape is another important aspect to explore. By studying the line shape we will find out whether we have come to a complete understanding of the involved transition and whether there are other unresolved small lines nearby ready to spoil our stabilization system.

Sometimes it may not be sufficient to use the natural resonance alone for stabilization work, or may not be necessary. The saturation aspect of the atomic transition limits the attainable S/N . To stabilize a noisy laser we need to use, for example, an optical resonator, which can provide a high-contrast and basically unlimited S/N of the resonance information. Careful study of the design and control of the material properties can bring the stability of material reference to a satisfactory level. See below for a more detailed discussion on this topic.

Ideally, a resonance line shape is even symmetric with respect to the center frequency of the resonance, and deviations from this ideal case will lead to frequency offsets. However, for the purpose of feedback, the resonance information needs to be converted to an odd symmetric discriminator shape: we need to know in which direction the laser is running away from the resonance. A straightforward realization of an error signal using direct absorption technique is to have the laser tuned to the side of resonance.¹¹ The slope of the line is used to convert the laser frequency noise to amplitude information

for the servo loop. This technique is essentially a DC approach and can suffer a huge loss in S/N due to the low-frequency amplitude noise of the laser. A differential measurement technique using dual beams is a requirement if one wishes to establish a somewhat stable operation. With a dual beam approach, the information about the laser noise can be measured twice and therefore it is possible to completely eliminate the technical noise and approach the fundamental limit of shot noise using clever designs of optoelectronic receivers. Conventional dual beam detection systems use delicate optical balancing schemes,¹² which are often limited by the noise and drift of beam intensities, residual interference fringes, drift in amplifiers, and spatial inhomogeneity in the detectors. Electronic auto-cancellation of the photodetector currents has provided near shot noise-limited performance.¹³ Although this process of input normalization helps to increase S/N of the resonance, the limitation on the locking dynamic range remains a problem. The servo loop simply gets lost when the laser is tuned to the tail or over the top of the resonance. Further, it is found that transient response errors basically limit the servo bandwidth to be within the cavity linewidth.¹⁴ Another effective remedy to the DC measurement of resonances is the use of zero-background detection techniques, for example, polarization spectroscopy.^{15,16} In polarization spectroscopy the resonance information is encoded in the differential phase shifts between two orthogonally polarized light beams. Heterodyne detection between the two beams can reveal an extremely small level of absorption-induced polarization changes of light, significantly improving the detection sensitivity. However, any practical polarizer has a finite extinction ratio (ϵ) which limits the attainable sensitivity. Polarization spectroscopy reduces the technical noise level by a factor of $\sqrt{\epsilon}$, with $\epsilon \sim 10^{-7}$ for a good polarizer. Polarization techniques do suffer the problem of long-term drifts associated with polarizing optics.

Modulation techniques are of course often used to extract weak signals from a noisy background. Usually noises of technical origins tend to be more prominent in the low frequency range. Small resonance information can then be encoded into a high-frequency region where both the source and the detector possess relatively small noise amplitudes. Various modulation schemes allow one to compare on-resonant and off-resonant cases in quick succession. Subsequent demodulations (lock-in detection) then simultaneously obtain and subtract these two cases, hence generating a signal channel with no output unless there is a resonance. Lorentzian signal recovery with the frequency modulation method has been well documented.¹⁷ The associated lock-in detection can provide the first, second, and third derivative type of output signals. The accuracy of the modulation waveform can be tested and various electronic filters can be employed to minimize nonlinear mixing among different harmonic channels and excellent accuracy is possible. In fact, the well-established 633-nm HeNe laser system¹⁸ is stabilized on molecular iodine transitions using this frequency dither technique and third harmonic (derivative) signal recovery. Demodulation at the third or higher order harmonics helps to reduce the influence of other broad background features.¹⁹ The shortcoming of the existence of dither on the output beam can be readily cured with an externally implemented "un-dithering" device based on an AOM.²⁰ However, in this type of modulation spectroscopy the modulation frequency is often chosen to be relatively low to avoid distortions on the spectral profile by the auxiliary resonances associated with modulation-induced spectral sidebands. An equivalent statement is that the line is distorted because it cannot reach an equilibrium steady state in the face of the rapidly tuning excitation. This low-frequency operation (either intensity chopping or derivative line shape recovery) usually is still partly contaminated by the technical noise and the achievable signal-to-noise ratio (S/N) is thereby limited. To recover the optimum signal size, large modulation amplitudes (comparable to the resonance width) are also employed, leading to a broadened spectral linewidth. Therefore the intrinsic line shape is modified by this signal recovery process and the direct experimental resolution is compromised.

A different modulation technique was later proposed and developed in the microwave magnetic resonance spectroscopy and similarly in the optical domain.²¹⁻²³ The probing field is phase-modulated at a frequency much larger than the resonance linewidth under study. When received by a square-law photodiode, the pure FM signal will generate no photocurrent at the modulation frequency unless a resonance feature is present to upset the FM balance. Subsequent heterodyne and rf phase-sensitive detection yield the desired signal. The high sensitivity associated with the FM spectroscopy is mainly due to its high modulation frequency, usually chosen to lie in a spectral region where the amplitude noise level of the laser source approaches the quantum (shot noise) limit. The redistribution of some of the carrier

power to its FM sidebands causes only a slight penalty in the recovered signal size. Another advantage of FM spectroscopy is the absence of linewidth broadening associated with low-frequency modulation processes. The wide-spread FM spectra allows each individual component to interact with the spectral features of interest and thereby preserves the ultrahigh resolution capability of contemporary narrow-linewidth lasers.

Since its invention, FM spectroscopy has established itself as one of the most powerful spectroscopic techniques available for high-sensitivity, high-resolution, and high-speed detection. The high bandwidth associated with the radio frequency (rf) modulation enables rapid signal recovery, leading to a high Nyquist sampling rate necessary for a high-bandwidth servo loop. The technique has become very popular in nonlinear laser spectroscopy,²⁴ including optical heterodyne saturation spectroscopy,²³ two-photon spectroscopy,²⁵ Raman spectroscopy,²⁶ and heterodyne four-wave mixing.²⁷ Recent developments with tunable diode lasers have made the FM technique simpler and more accessible. The field of FM-based laser diode detection of trace gas and remote sensing is rapidly growing. In terms of laser frequency stabilization, the rf sideband based Pound-Drever-Hall locking technique²⁸ has become a uniformly adopted fast stabilization scheme in the laser community. The resonance-based error signal in a high-speed operating regime is shown to correspond to the instantaneous phase fluctuations of the laser, with the atom or optical cavity serving the purpose of holding the phase reference. Therefore a properly designed servo loop avoids the response time of the optical phase/frequency storage apparatus and is limited only by the response of frequency-correcting transducers.

In practice some systematic effects exist to limit the ultimate FM sensitivity and the resulting accuracy and stability. Spurious noise sources include residual amplitude modulation (RAM), excess laser noise, and étalon fringes in the optical system.²⁹ A number of techniques have been developed to overcome these problems. In many cases FM sidebands are generated with electro-optic modulators (EOM). A careful design of EOM should minimize the stress on the crystal and the interference between the two end surfaces (using angled incidence or antireflection coatings). Temperature control of the EOM crystal is also important and has been shown to suppress the long-term variation of RAM.³⁰ The RAM can also be reduced in a faster loop using an amplitude stabilizer³¹ or a tuning filter cavity.³² The étalon fringe effect can be minimized by various optical or electronic means.³³ An additional low-frequency modulation (two-tone FM³⁴) can be used to reduce drifts and interference of the demodulated baseline.

In closing this section, we note that a laser is not always stabilized to a resonance but is sometimes referenced to another optical oscillator.³⁵ Of course the working principle does not change: one still compares the frequency/phase of the laser with that of reference. The technique for acquiring the error information is however more straightforward, often with a direct heterodyne detection of the two superposed waveforms on a fast photo detector. The meaning of the fast photo detector can be quite extensive, sometimes referring to a whole table-top system that provides THz-wide frequency gap measurement capabilities.^{36–38} Since it is the phase information that is detected and corrected, an optical phase locked loop usually provides a tight phase coherence between two laser sources. This is attractive in many measurement applications where the relative change of optical phase is monitored to achieve a high degree of precision. Other applications include phase-tracked master-slave laser systems where independent efforts can be made to optimize laser power, tunability and intrinsic noise.

The Optical Cavity-Based Frequency Discriminator

It is difficult to have both sensitive frequency discrimination and short time delay, unless one uses the reflection mode of operation: these issues have been discussed carefully elsewhere.²⁸ With ordinary commercial mirrors, we can have a cavity linewidth of 1 MHz, with a contrast C above 50 percent. We can suppose using 200 μW optical power for the rf sideband optical frequency discriminator, leading to a dc photo current i_0 of $\sim 100 \mu\text{A}$ and a signal current of $\sim 25 \mu\text{A}$. The shot noise of the dc current is $i_n = \sqrt{2ei_0}$ in a 1-Hz bandwidth, leading to an S/N of $\sim 4 \times 10^6$. The frequency noise-equivalent would then be 250 milliHertz/ $\sqrt{\text{Hz}}$. If we manage to design enough useful gain in the controller to suppress

the laser's intrinsic noise below this level, the laser output frequency spectrum would be characterized by this power spectral density. Under these circumstances, according to the earlier discussion in the Introduction and Overview Section, the output spectrum would be Lorentzian, of width $\Delta\nu_{\text{FWHM}} = \pi S_f = \pi (0.25 \text{ Hz})^2/\text{Hz} \sim 0.8 \text{ Hz}$. One comes to impressive predictions in this business! But usually the results are less impressive.

What goes wrong? From measurements of the servo error, we can see that the electronic lock is very tight indeed. However, the main problem is that vibrations affect the optical reference cavity's length and hence its frequency. For example measurements show the JILA Quiet Room floor has a seismic noise spectrum which can be approximated by $4 \times 10^{-9} \text{ m rms}/\sqrt{\text{Hz}}$ from below 1 Hz to about 20 Hz, breaking there to an f^{-2} roll-off. Below 1 Hz the displacement noise climbs as f^{-3} . Horizontal and vertical vibration spectra are similar. Accelerations associated with these motions lead to forces on the reference cavity that will induce mechanical distortion and hence frequency shifts. To estimate the resulting frequency shift, simple approximate analysis leads to a dynamic fractional modulation of the cavity length l by the (colinear) acceleration a , as

$$\left. \frac{\Delta l}{l} \right|_{\text{axial}} = -\frac{\Delta f}{f} = \frac{a \rho l \varepsilon}{2Y} \quad (6)$$

where $Y \sim 70 \text{ GPa}$ is the Young's modulus and $\rho \sim 2.2 \text{ gm/cm}^3$ is the density for the ULE (or Zerodur) spacer. The factor ε ($-1 < \varepsilon < 1$) is a geometrical design factor. For example, suppose the cavity is hanging vertically, suspended from the top. Then the cavity is stretched by its weight, and $\varepsilon = 1$. Using $l = 10 \text{ cm}$ and $a = 1 \text{ g}$, we expect $\Delta l/l = -\Delta f/f \sim 1.5 \times 10^{-8} \rightarrow \sim 8.7 \text{ MHz/g}$, supposing $\lambda = 532 \text{ nm}$. (This is equivalent to 885 kHz/ms^{-2} .) If the cavity were vertical, but supported from below, it would be in compression and $\varepsilon = -1$. Evidently there is an interesting regime in which the cavity is supported near its middle height, where there will be a strong cancelation of the net vertical length change. We return below to this case where $\varepsilon \sim 0$.

First, let us suppose our reference cavity bar is uniformly supported horizontally from a flat horizontal surface. Even in this transverse case, vertically accelerating the interferometer produces length changes through the distortion coupling between the transverse compression and lengthwise extension, the effect of "extrusion of the toothpaste." So the longitudinal displacement of Eq. (6) is reduced by this Poisson ratio $\sigma = 0.17$. Also the vertical weight now comes from the cavity's height, which is now really the spacer's diameter ϕ , typically about 5-fold less than the length. So we have

$$\left. \frac{\Delta l}{l} \right|_{\text{transverse}} = \frac{a \phi \rho \sigma}{2Y} \quad (7)$$

We come to a predicted sensitivity then of $\sim 300 \text{ kHz/g}$ for vertically applied uniform force (equivalent to 30 kHz/ms^{-2}).

Some important things have been so far left out of this discussion. For one, to make a stable reference cavity the details of the mounting and cavity support can be very important, since the expansion coefficient of the metal vacuum envelope is likely three orders of magnitude greater than that of ULE near its critical temperature-stable point. To prevent the vacuum shell's dimensional expansion from causing stresses in the cavity, it makes sense to use a pendulum suspension of the cavity. With two loops around the horizontal bar, forming a dual pendulum suspension, the cavity motion is mainly restricted to the axial direction, and the horizontal acceleration forces at high frequency are filtered down. Now we have the question:

What should be the spacing B between the two suspension loops? Put them close together and the expansion of the metal outer shell has even less impact on the cavity length. But the cavity rod (or bar, or tube) now takes on a stronger static bend, which shortens the cavity and the resulting cross-term in the cosine projection leads to a first-order length response with vertical acceleration noise. Furthermore, the bending-induced misalignment of the cavity mirrors means the intracavity resonant mode will displace laterally across the mirror surface to again have the optimal standing-wave buildup. Certainly the mirrors are rather nicely polished on their surfaces, but at least one is

a curved surface. So with our greedy dream of 10^{-15} frequency stability, wiping the beam vertically across the curvature will introduce disastrous optical length changes.

What about a wider spacing of the supports? Luckily for us the “two-point suspension problem” was addressed by G. B. Airy in the nineteenth century. He established that a support-spacing-to-length ratio of $B/L = 0.577$ was an ideal design for such a suspension, as it restored the parallelism of the two end faces of the measurement bar. A series of JILA experiments explored this domain.³⁹ These showed a vertical acceleration sensitivity of the horizontally suspended bar of 2200 kHz/ms^{-2} at $B/L = 0.11$, reduced to 150 kHz/ms^{-2} at the Airy spacing $B/L = 0.577$. Our “theory” in Eq. (7) doesn’t consider static bending of the bar, but would lead to 90 kHz/ms^{-2} if scaled for the $5.7 \times 7.1 \times 27.7 \text{ cm}$ dimensions of our cavity’s spacer-bar (suspended with the 5.1-cm direction vertical). Regrettably, the sign of the vibration-induced response was not determined: cavity bending shortens the optical path, while vertical squeezing would lengthen it.

Integrating the acceleration produced by the mid-band floor vibration spectrum quoted above leads to a broadband noise of a few dozen hertz in both H and V planes. Left out however is the 1 milli-“g” vibration near 30 Hz due to ac motors in JILA (Pepsi refrigerators!). So we should have a vibration-induced linewidth of something like 1/2 kHz, which correlates adequately well with experience. Passive air-table antivibration measures suppress this vibration (acceleration) spectrum to $\sim 2 \times 10^{-6} \text{ ms}^{-2}/\sqrt{\text{Hz}}$, again roughly flat over 2 to 20 Hz band by filtering the floor’s vibrational noise above $\sim 2 \text{ Hz}$ Fourier frequency. The calculated Fourier frequency at which the phase modulation processes have removed 1/2 the laser carrier power (approximate half-linewidth of the locked laser) is $\sim 1 \text{ Hz}$, but nonmodeled noise led to experimental values more like 5 Hz. Elegant passive vibration-damping suspensions at NIST have led to record-level subhertz cavity-locked laser linewidths.⁴⁰ It has been suggested that much of the remaining noise is associated with thermal mechanical displacement noise in the mirror coatings.⁴¹ Later measurements confirmed that the thermal noise was indeed the dominant source limiting the laser linewidth.⁴²

Returning to the cavity-mounting problem, we introduced the symmetry factor ε in the axial direction, because it is clear that holding the cavity in the midplane seems wise. Then the acceleration-induced net length change would tend be cancelled: one half of the length is under compression, the other half is under tension at a particular moment in the ac vibration cycle. We denote this cancellation by symmetry as ε , with $-1 \leq \varepsilon \leq 1$. Some experiments were made with short vertically mounted cavities.⁴³ The hand-assembly of the central disk limited the observed asymmetry value for our vertical mountings to a ~ 20 -fold reduction of the vibration sensitivity ($\varepsilon \geq 0.05$), to about $\sim 10 \text{ kHz/ms}^{-2}$, measured at the Nd fundamental wavelength. It was directly possible to observe subhertz laser beats! A computer design⁴⁴ for a more optimal cavity is shown as Fig. 8.

Quantum Resonance Absorption⁴⁵

Establishing a long-term stable optical frequency standard requires a natural reference of atomic or molecular origin. Historically, the use of atomic/molecular transitions was limited to those that had accidental overlap with some fixed laser wavelengths. With the advent of tunable lasers, research on quantum absorbers has flourished. A stabilized laser achieves fractional frequency stability

$$\frac{\delta\nu}{\nu} = \frac{1}{Q} \frac{1}{S/N} \frac{1}{\sqrt{\tau}}$$

where Q is the quality factor of the transition involved, S/N is the recovered signal-to-noise ratio of the resonance information, and τ is the averaging time. Clearly one wishes to explore the limits on both resolution and sensitivity of the detected signal. The nonlinear nature of a quantum absorber, while on one hand limiting the attainable S/N , permits sub-Doppler resolutions. With sensitive techniques such as FM-based signal modulation and recovery, one is able to split a MHz scale linewidth by a factor of 10^4 to 10^5 , at an averaging time of 1 s or so. Sub-Hertz long-term stability can be achieved with carefully designed optical systems where residual effects on baseline stability are minimized. However, a pressing question is: How accurate is our knowledge of the center of the resonance? Collisions, electromagnetic fringe fields, probe field wavefront curvature, and probe

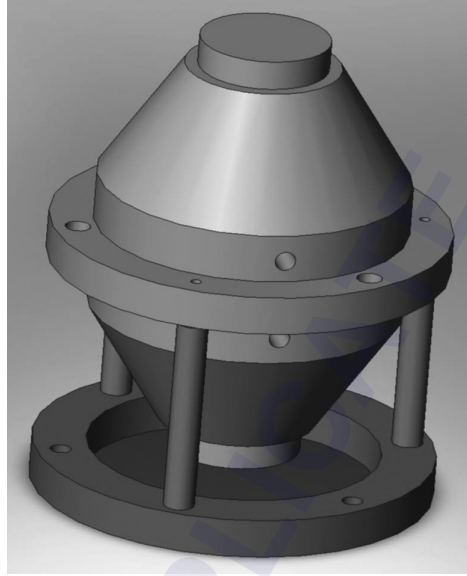


FIGURE 8. Computer model of vibration-resistant optical reference cavity.

power can all bring undesired linewidth broadening and center shifts. Distortion in the modulation process and other physical interactions can produce asymmetry in the recovered signal line shape. These issues will have to be addressed carefully before one can be comfortable talking about accuracy. A more fundamental issue related to time dilation of the reference system (second-order Doppler effect) can be solved in a controlled fashion: one simply knows the sample velocity accurately (e.g., by velocity selective Raman process), or the velocity is brought down to a negligible level using cooling and trapping techniques.

The simultaneous use of quantum absorbers and an optical cavity offers an attractive laser stabilization system. On one hand, a laser prestabilized by a cavity offers a long phase-coherence time, reducing the need of frequent interrogations of the quantum absorber. In other words, information of the atomic transition can be recovered with an enhanced S/N and the long averaging time translates into a finer examination of the true line center. On the other hand, the quantum absorber's resonance basically eliminates inevitable drifts associated with material standards. The frequency offset between the cavity mode and atomic resonance can be bridged by an AOM. In this case the cavity can be made of totally passive elements: mirrors are optically contacted to a spacer made of ultralow expansion material such as ULE or Zerodur. In case that the cavity needs to be made somewhat tunable, an intracavity Brewster plate driven by Galvo or a mirror mounted on PZT is often employed. Of course, these mechanical parts bring additional thermal and vibrational sensitivities to the cavity, along with nonlinearity and hysteresis. Temperature tuning of a resonator is potentially less noisy but slow. Other tuning techniques also exist, for example, through the use of magnetic force or pressure (change of intracavity refractive index or change of cavity dimension by external pressure). An often-used powerful technique called frequency-offset-locking brings the precision of tuning capability to the optical world.¹⁴

Transducers

PZT Transducer Design: Disk versus Tube Designs We will usually encounter the mechanical resonance problem in any servo based on a PZT transducer: Small mirrors clearly are nice as they can have higher resonance frequencies. A mirror, say 7.75 mm $\Phi \times$ 4 mm high, might be waxed onto a

PZT disk 10 mm diameter \times 0.5 mm thick. The PZT, in turn, is epoxied onto a serious backing plate. This needs to be massive and stiff, since the PZT element will produce a differential force between the mirror and the backing plate. At short times there will be a “reduced mass” kind of splitting of the motion between the mirror and the support plate. At lower frequencies, one hates to get a lot of energy coupled into the mirror mount since it will have a wealth of resonances in the sub-kHz range. For this size mirror, the backing plate might be stainless steel, 1 inch diameter by \sim 3/4 inch wedged thickness, and with the PZT deliberately decentered to break down high Q modes. The piston mode will be at \sim 75 kHz.

Tubular PZT Design Often it is convenient to use a tubular form of PZT, with the electric field applied radially across a thin wall of thickness t . This gives length expansion also, transverse to the field using a weaker d_{31} coefficient, but wins a big geometric factor in that the transverse field is generating a length response along the entire tube height h . The PZT tube could be 1/2 inch diameter by 1/2 inch length, with a wall thickness $t=1.25$ mm. This geometry leads to a \sim 7-fold sensitivity win, when $d_{31} \approx 0.7 d_{33}$ is included. Typical dimensions for the mirror might be 12.5 mm diameter \times 7 mm high. The PZT tube also is epoxied onto a serious backing plate. For the high voltage isolation of the PZT electrodes at the tube ends, a thin sheet (say < 0.5 mm) of stiff ceramic, alumina for example, will suffice. An alternative way to provide the electrical isolation of the ends involves removing the silver electrodes for several millimeters at the end. A new technique uses a diamond-charged tubular core drill mounted into a collet in a lathe. The active tool face projects out only 2 mm so that handheld PZT grinding leads to clean electrode removal, inside and out. This end of the PZT tube is attached to the backing mass with strong epoxy. The mirror is attached to the open PZT tube end with melted wax. This is vastly better than epoxy in that it does not warp the optic, and the small energy dissipation occurs at the best place to damp the Q of the PZT assembly. If done well, this unit will have its first longitudinal resonance at about 25 kHz, with a $Q \sim 10$. As noted above, in servo terms, the actual mechanical PZT unit gives an added 2-pole roll-off above the resonance frequency and a corresponding asymptotic phase lag of 180° . So it is useful to design for high resonant frequency and low Q.

Comparing disk and tubular designs, the disk approach can have a three-fold higher resonance frequency, while the tubular design is \sim 7-fold more sensitive. Perhaps more important is the tube’s reduced stiffness, moving the PZT/mirror resonance down into the 20-kHz domain. This brings us to the subject of spectral shaping of the amplifier gain and limitations of servo performance due to electronic issues.

Amplifier Strategies for PZT Driver We enjoy the tubular PZT for its large response per volt and its relatively high resonance frequencies. But it gives a problem in having a large capacitance, for example, of 10 nF in the above design. Even with the high sensitivity of 70 V/order, achieving a tight lock requires high frequency corrections and can lead to a problem in supplying the necessary ac current, supposing that we ask the HV amplifier alone to do the job. An apparent answer is to use a pair of amplifiers, one fast and the other HV, separately driving the two sides of the PZT. This alone doesn’t solve the problem, as the big high-frequency ac current is only returned via the HV amplifier. The answer is to use a crossover network on the HV amplifier side. A capacitor to ground, of perhaps 3- or 5-fold larger value than the PZT will adequately dump the fast currents coming through the PZT’s capacitance. A resistor to this PZT/shunt capacitor junction can go to the HV amplifier. Now this HV amp has indeed more capacitance to drive, but is only needed to be active below a few hundred hertz where the current demand becomes reasonable. An alternative topology sums the two inputs on one side of the PZT.

Other Useful Transducers—Slow but Powerful Commercial multiple wafer designs utilize 100 or more thin PZT sheets mechanically in series and electrically in parallel to produce huge excursions such as $10 \mu\text{m}$ for 100 V. Of course the capacitance is $\sim 0.1 \mu\text{F}$ and the stability leaves something to be desired. These are useful for applications that can tolerate some hysteresis and drift, such as grating angle tuning in a diode laser. When a large dynamic range is needed to accommodate wide tuning range or to correct for extensive laser frequency drifts at low frequencies, a galvo-driven Brewster

plate can be used inside the optical cavity. Typically a Brewster plate inflicts an insertion loss less than 0.1 percent if its angular tuning range is limited within $\pm 4^\circ$. Walk-off of the optical beam by the tuning plate can be compensated with a double-passing arrangement or using dual plates. In the JILA-designed Ti:Sapphire laser, we use the combination of PZT and Brewster plate for the long-term frequency stabilization. The correction signal applied to the laser PZT is integrated and then fed to the Brewster plate to prevent saturation of the PZT channel. At higher frequencies we use much faster transducers, such as AOM and EOM, which are discussed below.

Temperature control of course offers the most universal means to control long-term drifts. Unfortunately the time constant associated with thermal diffusion is usually slow and therefore the loop bandwidth of thermal control is mostly limited to Hertz scale. However, thoughtful designs can sometimes push this limit to a much higher value. For example, a Kapton thin-film heater tape wrapped around the HeNe plasma tube has produced a thermal control unity gain bandwidth in excess of 100 Hz.⁴⁶ The transducer response is reasonably modeled as an integrator above 0.3 Hz and excessive phase shifts associated with the thermal diffusion does not become a serious issue until ~ 200 Hz. This transfer function of the transducer can be easily compensated with an electronic PI filter to produce the desired servo loop response. Radiant heating of a glass tube by incandescent lamps has achieved a time delay < 30 ms and has also been used successfully for frequency control of HeNe lasers.⁴⁷ If a bipolar thermal control is needed, Peltier-based solid state heat pumps (thermoelectric coolers) are available and can achieve temperature differences up to 70°C , or can transfer heat at a rate of 125 W, given a proper configuration of heat sinking. Parallel use of these Peltier devices results in a greater amount of heat transfer while a cascaded configuration achieves a larger temperature difference.

Combining various servo transducers in a single feedback loop requires thorough understanding of each actuator, their gains and phase shifts, and the overall loop filter function one intends to construct. Clearly, to have an attractive servo response in the time domain, the frequency transfer functions of various gain elements need to crossover each other smoothly. A slow actuator may have some resonance features in some low-frequency domain, hence the servo action needs to be relegated to a faster transducer at frequency ranges beyond those resonances. The roll-off of the slow transducer gain at high-frequencies needs to be steep enough, so that the overall loop gain can be raised without exciting the associated resonance. On the other hand, the high-frequency channel typically does not have as large a dynamic range as the slow ones. So one has to pay attention not to overload the fast channel. Again, a steep filter slope is needed to rapidly relinquish the gain of the fast channel toward the low-frequency range. However, we stress here that the phase difference between the two channels at the crossover point needs to be maintained at less than 90° . In the end, predetermined gains and phase shifts will be assigned to each transducer so that the combined filter function resembles a smooth single channel design. Some of these issues will be addressed briefly in the section below on example designs.

Servo Design in the Face of Time Delay: Additional Transducers Are Useful As one wishes for higher servo gain, with stability, it means a higher closed-loop bandwidth must be employed. Eventually the gain is sufficiently large that the intrinsic laser noise, divided by this gain, has become less than the measurement noise involved in obtaining the servo error signal. This should be sufficient gain. However, it may not be usable in a closed-loop scenario, due to excessive time delay. If we have a time delay of t_{delay} around the loop from an injection to the first receipt of correction information, a consideration of the input and response as vectors will make it clear that no real servo noise suppression can occur unless the phase of the response at least approximates that required to subtract from the injected error input to reduce its magnitude on the next cycle through the system. A radian of phase error would correspond to a unity-gain frequency of $1/(2\pi t_{\text{delay}})$, and we find this to be basically the upper useful limit of servo bandwidth. One finds that to correct a diode laser or dye laser to leave residual phase errors of 0.1 rad, it takes about 2 MHz servo bandwidth. This means a loop delay time, at the absolute maximum, of $t_{\text{delay}} = 1/2\pi \cdot 2\text{MHz} = 80$ ns. Since several amplifier stages will be in this rf and servo-domain control amplifier chain, the individual bandwidths need to be substantially beyond the 12 MHz naively implied by the delay spec. In particular rf modulation frequencies need to be unexpectedly large, 20 MHz at least, and octave rf bandwidths need to

be utilized, considering that the modulation content can only be 1/2 the bandwidth. Suppression of even-order signals before detection is done with narrow resonant rf notches.

Of course a PZT transducer will not be rated in the nanosecond regime of time delay. Rather, one can employ an AOM driven by a fast-acting Voltage-Controlled Oscillator to provide a frequency shift. Unfortunately the acoustic time delay from the ultrasonic transducer to the optical interaction seems always to be 400 ns, and more if we are dealing with a very intense laser beam and wish to avoid damage to the delicate AO transducer. The AOM approach works well with diode-pumped solid state lasers, where the bandwidth of major perturbations might be only 20 kHz. By double-passing the AOM the intrinsic angular deflection is suppressed. Usually the AOM prefers linear polarization. To aid separation of the return beam on the input side, a spatial offset can be provided with a collimating lens and roof prism, or with a cat's-eye retroreflector. Amplitude modulation or leveling can also be provided with the AOM's dependence on rf drive, but it is difficult to produce a beam still at the shot noise level after the AOM.

The final solution is an *EOM phase modulator*. In the external beam, this device will produce a phase shift per volt, rather than a frequency shift. So we will need to integrate the control input to generate a rate of change signal to provide to the EOM, in order to have a frequency relationship with the control input.⁴⁸ Evidently this will bring the dual problems of voltage saturation when the output becomes too large, and a related problem, the difficulty of combining fast low-delay response with high-voltage capability. The standard answer to this dilemma was indicated in our PZT section, namely, one applies fast signals and high-voltage signals independently, taking advantage of the fact that the needed control effort at high frequencies tends to cover only a small range. So fast low-voltage amplifier devices are completely adequate, particularly if one multipasses the EOM crystal several times. A full discussion of the crossover issues and driver circuits will be prepared for another publication.

Representative/Example Designs

Diode-Pumped Solid State Laser Diode-pumped solid state lasers are viewed as the most promising coherent light sources in diverse applications, such as communications, remote detection, and high precision spectroscopy. The diode-pumped Nd:YAG laser is probably the most highly developed of the rapidly expanding universe of diode-pumped solid state lasers, and it has enjoyed continuous improvements in its energy efficiency, size, lifetime, and intrinsic noise levels. The laser's free-running linewidth of ~ 10 kHz makes it a straightforward task to stabilize the laser via an optical cavity or an optical phase locked loop. In our initial attempt to stabilize the laser on a high finesse ($F \sim 100,000$, linewidth ~ 3 kHz) cavity, we employ an external AOM along with the laser internal PZT which is bonded directly on the laser crystal. The frequency discrimination signal between the laser and cavity is obtained with 4-MHz FM sidebands detected in cavity reflection. The PZT corrects any slow but potentially large laser frequency noise. Using the PZT alone allows the laser to be locked on the cavity. However, the loop tends to oscillate around 15 kHz and the residual noise level is more than 100 times higher than that obtained with the help of an external AOM. The AOM is able to extend the servo bandwidth to ~ 150 kHz, limited by the propagation time delay of the acoustic wave inside the AOM crystal. The crossover frequency between the PZT and AOM is about 10 kHz. Such a system has allowed us to achieve a residual frequency noise spectral density of 20 mHz/ $\sqrt{\text{Hz}}$. The laser's linewidth relative to cavity is thus a mere 1.3 mHz,⁴⁹ even though the noise spectral density is still 100 times higher than the shot noise. This same strategy of servo loop design has also been used to achieve a microradian level phase locking between two Nd:YAG lasers.⁵⁰

It is also attractive to stabilize the laser directly on atomic/molecular transitions, given the low magnitude of the laser's intrinsic frequency noise. Of course the limited S/N of the recovered resonance information will not allow us to build speedy loops to clean off the laser's fast frequency/phase noise. Rather we will use the laser PZT alone to guide the laser for a long-term stability. An example here is the 1.064 μm radiation from the Nd:YAG, which is easily frequency doubled to 532 nm where strong absorption features of iodine molecules exist.^{51,52} The doubling is furnished with a noncritical phase-matched KNb_3O_7 crystal located inside a buildup cavity. 160 mW of green light

output is obtained from an input power of 250 mW of IR. Only mW levels of the green light are needed to probe the iodine saturated absorption signal. Low vapor pressure (~ 0.5 Pa) of the iodine cell is used to minimize the collision-induced pressure shift and to reduce the influence on baseline by the linear Doppler absorption background. The signal size decreases as the pressure is reduced. However, this effect is partly offset due to the reduced resonance linewidth (less pressure broadening) which helps to increase the slope of the frequency locking error signal. A lower pressure also helps to reduce power-related center frequency shifts since a lower power is needed for saturation. With our 1.2-m long cells, we have achieved an S/N of 120 in a 10-kHz bandwidth, using the modulation transfer spectroscopy.⁵³ (Modulation transfer is similar to FM except that we impose the frequency sideband on the saturating beam and rely on the nonlinear medium to transfer the modulation information to the probe beam which is then detected.) Normalized to 1-s averaging time, this S/N translates to the possibility of a residual frequency noise level of 10 Hz when the laser is locked on the molecular resonance, given the transition linewidth of 300 kHz. We have built two such iodine-stabilized systems and the heterodyne beat between the two lasers permits systematic studies on each system and checks the reproducibility of the locking scheme.⁵⁴ With a 1 second counter gate time, we have recorded the beat frequency between the two lasers. The standard deviation of the beat frequency noise is ~ 20 Hz, corresponding to ~ 14 -Hz rms noise per IR laser, basically a S/N limited performance. The beat record can be used to calculate the Allan standard deviation: starting at 5×10^{-14} at 1 second, decreasing with a slope of $1/\sqrt{\tau}$ up to 100 second. (τ is the averaging time.) After 100-s the deviations reach the flicker noise floor of $\sim 5 \times 10^{-15}$. At present, the accuracy of the system is limited by inadequate optical isolation in the spectrometer and the imperfect frequency modulation process (residual amplitude noise, RAM) used to recover the signal. This subject is under intense active study in our group.⁵⁵

External Cavity Diode Lasers Diode lasers are compact, reliable, and coherent light sources for many different applications.⁵⁶ The linewidth of a free-running diode laser is limited by the fundamental spontaneous emission events, enhanced by the amplitude-phase coupling inside the gain medium. With a low-noise current driver, a typical milliwatt scale AlGaAs diode laser has a linewidth of several MHz. To reduce this fast frequency noise, one typically employs an external cavity formed between one of the diode laser facets and a grating (or an external mirror that retroreflects the first-order grating diffraction).^{57–59} This optical feedback mechanism suppresses the spontaneous emission noise, replaced by much slower fluctuations of mechanical origin. The linewidth of the grating-stabilized external cavity diode laser (ECDL) is usually between 100 kHz and 1 MHz, determined by the quality factor of the optical feedback. The ECDL also offers much better tuning characteristics compared against a solitary diode. To do such tuning, the external grating (or the mirror that feeds the grating-dispersed light back to the laser) is controlled by a PZT for scanning. Synchronous tuning of the grating dispersion and the external cavity mode can be achieved with a careful selection of the grating rotation axis position. Similarly, this PZT-controlled grating can be used to stabilize the frequency of an ECDL. However, owing to the low bandwidth limited by the mechanical resonance of PZT, a tight frequency servo is possible only through fast transducers such as the laser current or intracavity phase modulators.

This hybrid electro-optic feedback system is attractive, and ECDLs have been demonstrated to show hertz level stability under a servo bandwidth of the order of 1 MHz. For a solitary diode, feedback bandwidth of tens of megahertz would have been needed in order to bring the frequency noise down to the same level. However, considering that the optical feedback has a strong impact on the laser frequency noise spectrum, one finds the frequency response of the compound laser system is clearly dependent upon the optical alignment. Therefore, for each particular ECDL system, we need to measure the frequency response function of the laser under the optimally aligned condition. We are dealing with a multichannel feedback system (e.g., PZT plus current), so that designing smooth crossovers between different transducers requires knowledge of the transfer functions of each transducer. Normally the current-induced FM of a solitary diode has a flat response up to ~ 100 kHz, and then starts to roll off in the region between 100 kHz and 1 MHz, initially with a single-pole character. This is due to the time response of the current-induced thermal change of the refractive index inside the diode. (At a faster time scale, the carrier density variation will remain and then dominate

the laser frequency response.) Design of a fast feedback loop needs to take into account this intrinsic diode response. Fortunately the time delay associated with the current response is low, typically below 10 ns.

In our example system, the frequency discrimination signal of the ECDL is obtained from a 100 kHz linewidth cavity with a sampling frequency of 25 MHz. The error signal is divided into three paths: PZT, current modulation through the driver, and direct current feedback to the diode head. The composite loop filter function is shown in Fig. 9. The crossover between the slow current channel and the PZT usually occurs around 1 kHz, in order to avoid the mechanical resonance of the PZT at a few kHz. In our system, the frequency response of the PZT/grating is 10 GHz/V. To furnish this in-loop gain of ~ 1000 at 1 kHz, we need to supply an electronic gain of 0.1, given that the error signal has a slope of 1 V/1 MHz.

Toward the lower frequency range, the PZT gain increases by 40 dB/decade (double integrators) to suppress the catastrophically rising laser frequency noise. It is obvious from Fig. 9 that the intermediate current channel tends to become unstable at a few hundred kHz, due to the excessive phase shift there. The fast current loop, bypassing the current driver to minimize additional time delay and phase shift, has a phase lead compensator to push the unity gain bandwidth to 2 MHz. With this system we can lock the ECDL robustly on the optical cavity, with a residual noise spectral density of 2 Hz/ $\sqrt{\text{Hz}}$, leading to a relative linewidth of 12 Hz. The achieved noise level is about 100 times higher than the fundamental measurement limit set by shot noise. We note in passing that when an ECDL gradually goes out of alignment, the previously adjusted gain of the current loop will tend to make the servo oscillate so we know a new alignment is needed. The laser FM sideband used to generate

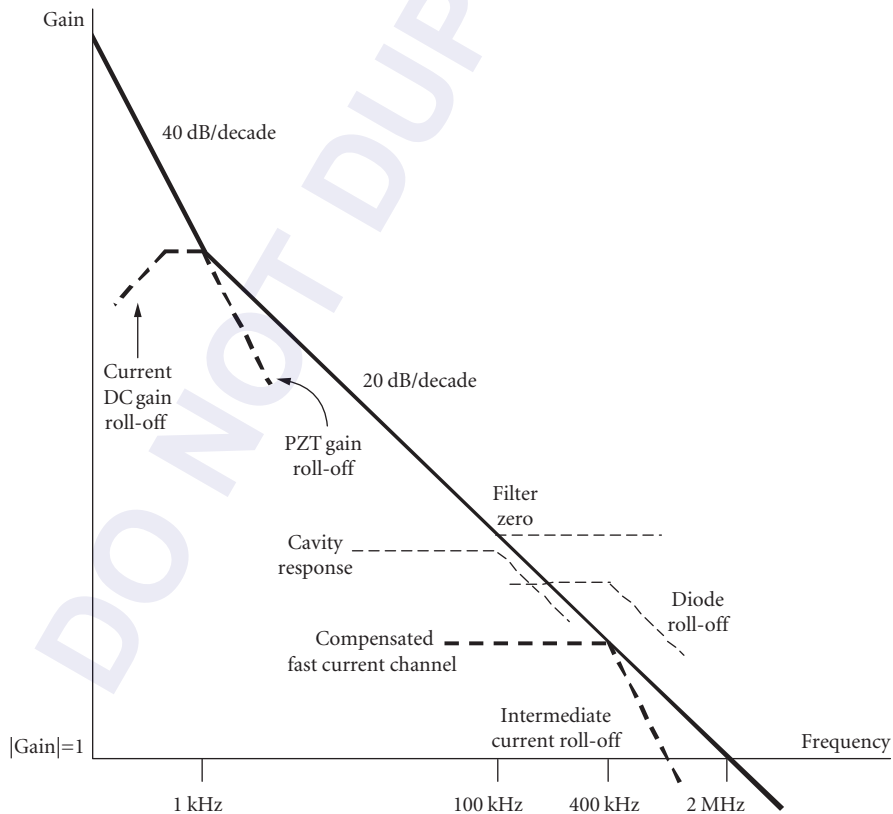


FIGURE 9 The combined loop filter function for ECDL frequency stabilization.

the locking signal is produced directly by current modulation. An electronic filter network is employed to superimpose the slow servo, fast servo, and modulation inputs to the diode. Exercise caution when accessing the diode head, as a few extra mA current increase can lead to drastic output power increase and melted laser facets, all in 1 μ s!

Since 2000 this art of locking diode lasers to cavities has been progressing rapidly and in many labs worldwide. To note just one recent measurement of subhertz linewidths, we may note the results⁶⁰ in JILA, where two diode lasers were locked to two independent vertical cavities, giving 10^{-15} frequency stability at 1-second averaging time. In another Boulder collaboration, subhertz optical beat linewidths have been measured⁶¹ between a JILA diode laser source at 689 nm and another stabilized laser used as the local oscillator for the Hg⁺ clock, at a different wavelength of 1126 nm, and located 4 km away in the NIST labs. Two optical frequency Combs were employed, along with a Nd:YAG laser whose wavelength could be transmitted by the 4-km fiber link. (This Comb technology seems to have gone quickly from being a research topic to a reliable and versatile tool!)

22.4 SUMMARY AND OUTLOOK

The technology of laser frequency stabilization has been refined and simplified over the years and has become an indispensable research tool in any modern laboratories involving optics. Research on laser stabilization has been and still is pushing the limits of measurement science. Indeed, a number of currently active research projects on fundamental physical principles benefit a great deal from stable optical sources and will need a continued progress of the laser stabilization front.⁶² Using extremely stable phase coherent optical sources, we will be entering an exciting era when the LISA interferometer will achieve picometer resolution over a five million kilometer distance in space⁶³ and a few Hertz linewidth of an ultraviolet resonance will be probed with a high S/N .^{64, 65} One has to be optimistic looking at the stabilization results of all different kinds of lasers. To list just a few examples of cw tunable lasers, we notice milliHertz linewidth stabilization (relative to a cavity) for diode-pumped solid state lasers; dozen milliHertz linewidth for Ti:Sapphire lasers; and sub-Hertz linewidths for diode and dye lasers. Long-term stability of lasers referenced to atoms and molecules⁶⁶ has reached mid 10^{-16} levels in an averaging time as short as ~ 300 s. Phase locking between different laser systems is now routinely employed, even for diode lasers that have fast frequency noise.

An important new capability comes from the accurate frequency transfer via noise-compensated optical fiber, which allows groups to collaborate in the testing of their state-of-the-art systems, thus tapping into the fruits of our neighbors' achievements in stabilizing lasers.⁶⁷ The Optical Comb makes it possible to use nearly any stable frequency as the reference. The fiber transport allowed measurement of the inaccuracy issues affecting four potential atomic clocks, based on a trapped Hg⁺ ion,⁶⁸ Al⁺ trapped ions,⁶⁸ on cold and free Ca atoms,⁶⁹ and on cold Sr atoms confined in an optical lattice.⁶⁶ Measurements confirmed that the Sr, Hg⁺, and Al⁺ systems can be expected to yield sub- 1×10^{-16} independent reproducibility (the same quality as *accuracy*, except for the authority of Cs definition as the unit of time.) Thus for the first time there is not one, but even three(!) frequency reference systems which surpass the performance of the present best standard in Physics, namely the Cs Fountain Clock. Likely there will be other optically-based systems (such as Yb in a lattice⁷⁰) which may soon be in this performance category. Clearly there will be great laser and physics fun coming in the next times! For example, an extended series of comparisons between Hg⁺ and Al⁺ was carried out at NIST under strict conditions, and provided a $< 1 \times 10^{-16}$ /year checking for possible drift in the fundamental "constants" of physics.⁶⁸ By the way, do *you* believe that the frequencies of nuclear Mössbauer transitions will stay at the same frequency, as measured by these atomic frequency sources? With the rapid progress in producing VUV Combs by Higher Harmonic Generation,^{71, 72} one can begin to dream about coherent Mössbauer spectroscopy so we can find out. In the optical comparisons, we're looking in the seventeenth decimal digit now!

Another highway toward fundamental physics involves precision laser spectroscopy of simple atoms, such as H and He⁺, but laser cooling is a challenge due to the low available power. The exciting program⁷³ at MPQ in Munich can now generate enough VUV 121.5-nm light to begin this story!

Quantum noise is the usual limit of the measurement process and therefore may be the limit of the stabilization process as well. To circumvent the quantum noise altogether is an active research field itself.⁷⁴ We, however, have not reached this quantum limit just yet. For instance, we have already stated that the Nd:YAG laser should be able to reach microHertz stability if the shot noise is the true limit. What have we done wrong? A part of the deficiency is due to the inadequacy of the measurement process, namely the lack of accuracy. This is because the signal recovery effort—modulation and demodulation process—is contaminated by spurious optical interference effects and RAM associated with the modulation frequency. Every optical surface along the beam path can be a potential time bomb to damage the modulation performance. In cases that some low contrast interference effects are not totally avoidable, we would need to have the whole system controlled in terms of the surrounding pressure and temperature. The degree to which we can exert control of course dictates the ultimate performance. The second fundamental limitation appears to be thermally generated mechanical noise of the optical coatings, which was already noted,^{41,42} but there are already schemes being discussed to buy us another decade or two.

22.5 CONCLUSIONS AND RECOMMENDATIONS

It becomes clear that there are many interlinking considerations involved in the design of laser stabilization systems, and it is difficult to present a full description in a chapter such as this. Still it is hoped that the reader will see some avenues to employ feedback control methods to the laser systems of her current interest. We are optimistic that some of this technology may become commercially available in the future, thus simplifying the user's task.

22.6 ACKNOWLEDGMENTS

The work discussed here has profited from interactions with many colleagues, postdoctoral researchers, and graduate students over many years. In particular we must thank Leo Hollberg and Miao Zhu for their earlier contributions. More recent contributors are Long Sheng Ma, and Mark Notcutt. Especially we thank Mark for his experimental work with the vertically mounted cavities, and for his useful suggestions for improving this text. The now-retired one of us (JLH) declares his joy to see the JILA laboratory prospering in the strong hands of Jun Ye. Jan also is particularly grateful to his wife Lindy for patience “beyond the call of duty” during these many years of laser research. The work at JILA has been supported over many years by the Office of Naval Research, the National Science Foundation, the Air Force Office of Scientific Research, NASA, and the National Institute of Standards and Technology, as part of its frontier research into basic standards and their applications.

22.7 REFERENCES

1. D. W. Allan, “Statistics of Atomic Frequency Standards,” *Proc. IEEE* **54**:221–230 (1966).
2. D. B. Sullivan, D. W. Allan, D. A. Howe, and F. L. Walls, (eds.), *Characterization of Clocks and Oscillators*, NIST Technical Note **1337**, U. S. Government Printing Office, Washington D.C., 1990. See also <http://tf.nist.gov/timefreq/general/generalpubs.htm> (2009).
3. D. S. Elliot, R. Roy and S. J. Smith, “Extracavity Laser Band-Shape and Bandwidth Modification,” *Phys. Rev. A* **26**:12–18 (1982).
4. J. L. Hall and M. Zhu, “An Introduction to Phase-Stable Optical Sources,” in Proc. of the *International School of Phys. “Enrico Fermi,” Course CXVIII, Laser Manipulation of Atoms and Ions*, E. Arimondo, W. D. Phillips, and F. Strumia, (eds.), North Holland, pp 671–702 (1992).
5. M. Zhu and J. L. Hall, “Stabilization of Optical-Phase Frequency of a Laser System—Application to a Commercial Dye-Laser with an External Stabilizer,” *J. Opt. Soc. Am. B* **10**: 802–816 (1993).

6. For similar issues in the microwave/rf field, see the application note *Time Keeping and Frequency Calibration*, Agilent, Palo Alto, Calif., and <http://tf.nist.gov/timefreq/general/generalpubs.htm>.
7. For general references on feedback systems, please refer to *Modern Control Systems*, 3rd ed., Richard C. Dorf, Addison-Wesley Publishing Co., Reading Mass. (1980).
8. *Feedback Control of Dynamic Systems*, 3rd ed., Gene F. Franklin, J. David Powell, and A. Emami-Naeini, Addison-Wesley Publishing Co., Reading, Mass. (1994).
9. The simulations presented here are performed with the following software by The Math Works Inc., Natick, Mass. (1996): Matlab Control Systems Toolbox: User's Guide; Matlab Simulink: User's Guide.
10. See, e.g., J. L. Hall, "Frequency Stabilized Lasers—A Parochial Review," in *Frequency-Stabilized Lasers and Their Applications*, Y. C. Chung, (ed.), *Proc. SPIE* **1837**:2–15 (1993).
11. R. L. Barger, M. S. Sorem, and J. L. Hall, "Frequency Stabilization of a Cw Dye Laser," *Appl. Phys. Lett.* **22**:573–575 (1973).
12. G. D. Houser and E. Garmire, "Balanced Detection Technique to Measure Small Changes in Transmission," *Appl. Opt.* **33**:1059–1062 (1994).
13. K. L. Haller and P. C. D. Hobbs, "Double Beam Laser Absorption Spectroscopy: Shot-Noise Limited Performance at Baseband with a Novel Electronic Noise Canceller," in *Optical Methods for Ultrasensitive Detection and Analysis: Techniques and Applications*, B. L. Fearey, (ed.), *Proc. SPIE* **1435**: 298–309 (1991). Also see <http://www.electrooptical.net> (2009).
14. J. Helmcke, S. A. Lee, and J. L. Hall, "Dye-Laser Spectrometer for Ultrahigh Spectral Resolution—Design and Performance," *Appl. Opt.* **21**:1686–1694 (1982).
15. C. E. Wieman and T. W. Hänsch, "Doppler-Free Laser Polarization Spectroscopy," *Phys. Rev. Lett.* **36**:1170–1173 (1976).
16. T. W. Hänsch and B. Couillaud, "Laser Frequency Stabilization by Polarization Spectroscopy of a Preflecting Reference Cavity," *Opt. Comm.* **35**:441–444 (1980).
17. H. Wahlquist, "Modulation Broadening of Unsaturated Lorentzian Lines," *J. Chem. Phys.* **35**:1708–1710 (1961).
18. T. J. Quinn, "Mise en Pratique of the Definition of the Metre (1992)" *Metrologia* **30**:524–541 (1994).
19. J. Hu, E. Ikonen, and K. Riski, "On the Nth Harmonic Locking of the Iodine Stabilized He-Ne-Laser," *Opt. Comm.* **120**:65–70 (1995); **121**:169 (1995).
20. M. S. Taubman and J. L. Hall, "Cancellation of laser dither modulation from optical frequency standards," *Opt. Lett.* **25**:311–313 (2000).
21. B. Smaller, *Phys. Rev.* **83**:812–820 (1951); R. V. Pound, "Electronic Stabilization of Microwave Oscillators," *Rev. Sci. Instrum.* **17**: 490–505 (1946).
22. G. C. Bjorklund, "Frequency-Modulation Spectroscopy: A New Method for Measuring Weak Absorptions and Dispersions," *Opt. Lett.* **5**:15–17 (1980).
23. J. L. Hall, L. Hollberg, T. Baer, and H. G. Robinson, "Optical Heterodyne Saturation Spectroscopy," *Appl. Phys. Lett.* **39**:680–682 (1981).
24. G. C. Bjorklund and M. D. Levenson, "Sub-Doppler Frequency-Modulation Spectroscopy of I₂," *Phys. Rev. A* **24**:166–169 (1981).
25. W. Zapka, M. D. Levenson, F. M. Schellenberg, A. C. Tam, and G. C. Bjorklund, "Continuous-Wave Doppler-Free Two-Photon Frequency Modulation Spectroscopy in Rb Vapor," *Opt. Lett.* **8**:27–29 (1983).
26. G. J. Rosasco and W. S. Hurst, "Phase-Modulated Stimulated Raman Spectroscopy," *J. Opt. Soc. Am. B* **2**:1485–1496 (1985).
27. J. H. Shirley, "Modulation Transfer Processes in Optical Heterodyne Saturation Spectroscopy," *Opt. Lett.* **7**:537–539 (1982).
28. R. W. P. Drever, J. L. Hall, F. V. Kowalski, J. Hough, G. M. Ford, A. J. Munley, and H. Ward, "Laser Phase and Frequency Stabilization Using an Optical-Resonator," *Appl. Phys. B* **31**:97–105 (1983).
29. P. Werle, "Laser Excess Noise and Interferometric Effects in Frequency-Modulated Diode-Laser Spectrometers," *Appl. Phys. B* **60**:499–506 (1995).
30. J. L. Hall, J. Ye, L.-S. Ma, K. Vogel, and T. Dinneen, "Optical Frequency Standards: Progress and Applications," in *Laser Spectroscopy XIII*, Z.-J. Wang, Z.-M. Zhang, and Y.-Z. Wang, (eds.), World Scientific, Singapore, pp. 75–80 (1998).

31. N. C. Wong and J. L. Hall, "Servo Control of Amplitude-Modulation in Frequency-Modulation Spectroscopy—Demonstration of Shot-Noise-Limited Detection," *J. Opt. Soc. Am. B* **2**:1527–1533 (1985).
32. M. S. Taubman and J. L. Hall, unpublished JILA work (1999).
33. D. R. Hjelm, S. Neegård, and E. Vartdal, "Optical Interference Fringe Reduction in Frequency-Modulation Spectroscopy Experiments," *Opt. Lett.* **20**:1731–1733 (1995).
34. G. R. Janik, C. B. Carlisle, and T. F. Gallagher, "Two-Tone Frequency-Modulation Spectroscopy," *J. Opt. Soc. Am. B* **3**:1070–1074 (1986).
35. J. L. Hall, L.-S. Ma, and G. Kramer, "Principles of Optical Phase-Locking—Application to Internal Mirror He-Ne Lasers Phase-Locked Via Fast Control of the Discharge Current," *IEEE J. Quan. Electron.* **QE-23**:427–437 (1987).
36. M. Kourogi, K. Nakagawa, and M. Ohtsu, "Wide-Span Optical Frequency Comb Generator for Accurate Optical Frequency Difference Measurement," *IEEE J. Quan. Electron.* **QE-29**:2693–2701 (1993).
37. Th. Udem, J. Reichert, R. Holzwarth, and T. W. Hänsch, "Accurate Measurement of Large Optical Frequency Differences with a Mode-Locked Laser," *Opt. Lett.* **24**:881–883 (1999).
38. S. A. Diddams, L.-S. Ma, J. Ye, and J. L. Hall, "Broadband Optical Frequency Comb Generation with a Phase-Modulated Parametric Oscillator," *Opt. Lett.* **24**:1747–1749 (1999).
39. J. L. Hall, M. Notcutt, and J. Ye, "Improving Laser Coherence," in *Laser Spectroscopy XVII*, E. Hinds, A. Ferguson, and E. Riis, (eds.), World Scientific, Singapore, pp. 3–12 (2006).
40. B. C. Young, F. C. Cruz, W. M. Itano, and J. C. Bergquist, "Visible Lasers with Subhertz Linewidths," *Phys. Rev. Lett.* **82**:3799–3802 (1999).
41. K. Numata, A. Kemery, and J. Camp, "Thermal-Noise Limit in the Frequency Stabilization of Lasers with Rigid Cavities," *Phys. Rev. Lett.* **93**:250602 (2004).
42. M. Notcutt, L. S. Ma, A. D. Ludlow, S. M. Foreman, J. Ye, J. L. Hall, "Contribution of Thermal Noise to Frequency Stability of Rigid Optical Cavity via Hertz-Linewidth Lasers," *Phys. Rev. A* **73**:031804 (R) (2006).
43. Notcutt, M. L. S. Ma, J. Ye, J. L. Hall, "Simple and Compact 1-Hz Laser System via an Improved Mounting Configuration of a Reference Cavity," *Opt. Lett.* **30**:1815–1817 (2005).
44. JILA design by Lisheng Chen, 2004. See L. Chen, J. L. Hall, J. Ye, T. Yang, E. Zang, and T. Li, "Vibration-Induced Elastic Deformation of Fabry-Perot Cavities," *Phys. Rev. A* **74**:053801 (2006).
45. Interested readers are referred to the following books for more detailed discussions. *The Quantum Physics of Atomic Frequency Standards*, vol. I and II, Jacques Vanier and Claude Audoin; Adam Hilger, Institute of Physics Publishing Ltd, Bristol, UK (1989).
46. T. M. Niebauer, J. E. Faller, H. M. Godwin, J. L. Hall, and R. L. Barger, "Frequency Stability Measurements on Polarization-Stabilized He-Ne Lasers," *Appl. Opt.* **27**:1285–1289 (1988).
47. Jun Ishikawa, NMIJ, AIST, Tsukuba, Japan, private communications (1996).
48. J. L. Hall and T. W. Hänsch, "External Dye-Laser Frequency Stabilizer," *Opt. Lett.* **9**:502–504 (1984).
49. J. Ye, L.-S. Ma, and J. L. Hall, "Ultrasensitive Detections in Atomic and Molecular Physics: Demonstration in Molecular Overtone Spectroscopy," *J. Opt. Soc. Am. B* **15**:6–15 (1998).
50. J. Ye and J. L. Hall, "Optical Phase Locking in the Microradian Domain: Potential Applications to NASA Spaceborne Optical Measurements," *Opt. Lett.* **24**:1838–1840 (1999).
51. A. Arie and R. L. Byer, "Laser Heterodyne Spectroscopy of $^{127}\text{I}_2$ Hyperfine Structure near 532 nm," *J. Opt. Soc. Am. B* **10**:1990–1997 (1993).
52. M. L. Eickhoff and J. L. Hall, "Optical Frequency Standard at 532 Nm," *IEEE Trans. Instrum. Meas.* **44**:155–158 (1995); and P. Jungner, M. Eickhoff, S. Swartz, J. Ye, J. L. Hall and S. Waltman, *IEEE Trans. Instrum. Meas.* **44**:151–154 (1995).
53. J. Ye, L. Robertsson, S. Picard, L.-S. Ma, and John L. Hall, "Absolute Frequency Atlas of Molecular I-2 Lines at 532 nm," *IEEE Trans. Instrum. Meas.* **48**:544–549 (1999).
54. J. L. Hall, L.-S. Ma, M. Taubman, B. Tiemann, F.-L. Hong, O. Pfister, and J. Ye, "Stabilization and Frequency Measurement of the I-2-Stabilized Nd:YAG laser," *IEEE Trans. Instrum. Meas.* **48**:583–586 (1999).
55. C. Ishibashi, J. Ye, and J. L. Hall, "Issues and Applications in Ultra-Sensitive Molecular Spectroscopy," *Proc. SPIE.* **4634**:58–69 (2002).
56. C. E. Wieman and L. Hollberg, "Using Diode-Lasers for Atomic Physics," *Rev. Sci. Instrum.* **62**:1–20 (1991).

57. M. G. Littman and H. J. Metcalf, "Spectrally Narrow Pulsed Dye Laser without Beam Expander," *Appl. Opt.* **17**:2224–2227 (1978).
58. B. Dahmani, L. Hollberg, and R. Drullinger, "Frequency Stabilization of Semiconductor-Lasers by Resonant Optical Feedback," *Opt. Lett.* **12**:876–878 (1987).
59. K. MacAdam, A. Steinbach, and C. E. Wieman, "A Narrow-Band Tunable Diode-Laser System with Grating Feedback, and a Saturated Absorption Spectrometer for Cs and Rb," *Am. J. Phys.* **60**:1098–1111 (1992).
60. A. D. Ludlow, X. Huang, M. Notcutt, T. Zanon-Willette, S. M. Foreman, M. M. Boyd, S. Blatt, and J. Ye, "Compact, Thermal-Noise-Limited Optical Cavity for Diode Laser Stabilization at 1×10^{-15} ," *Opt. Lett.* **32**:641–643 (2007).
61. S. M. Foreman, A. D. Ludlow, M. H. G. de Miranda, J. Stalnaker, S. A. Diddams, and J. Ye, "Coherent Optical Phase Transfer over a 32-km Fiber with 1 s Instability at 10^{-17} ," *Phys. Rev. Lett.* **99**:153601 (2007).
62. P. Fritschel, G. González, B. Lantz, P. Saha, and M. Zucker, "High Power Interferometric Phase Measurement Limited by Quantum Noise and Application to Detection of Gravitational Waves," *Phys. Rev. Lett.* **80**:3181–3185 (1998).
63. K. Danzmann, "LISA—An ESA Cornerstone Mission for a Gravitational Wave Observatory," *Class. Quantum Grav.* **14**:1399–1404 (1997).
64. J. C. Bergquist, R. J. Rafac, B. C. Young, J. A. Beall, W. M. Itano, and D. J. Wineland, "Sub-Dekahertz Spectroscopy of $^{199}\text{Hg}^+$," in *Laser Frequency Stabilization, Standards, Measurement, and Applications*, J. L. Hall and J. Ye, (eds.), *Proc. SPIE*. **4269**:1–7 (2001).
65. E. E. Eyler, D. E. Chieda, M.C. Stowe, M. J. Thorpe, T. R. Schibli, and J. Ye, "Prospects for Precision Measurements of Atomic Helium using Direct Frequency Comb Spectroscopy," *Eur. Phys. J. D.* **48**:43–55 (2008).
66. D. Ludlow, T. Zelevinsky, G. K. Campbell, S. Blatt, M. M. Boyd, M. H. G. de Miranda, M. J. Martin, J. W. Thomsen, S. M. Foreman, Jun Ye, T. M. Fortier, J. E. Stalnaker, S. A. Diddams, Y. Le Coq, Z. W. Barber, N. Poli, N. D. Lemke, K. M. Beck, and C. W. Oates, "Sr Lattice Clock at 1×10^{-16} Fractional Uncertainty by Remote Optical Evaluation with a Ca clock," *Science* **319**:1805–1808 (2008).
67. S. M. Foreman, K. W. Holman, D. D. Hudson, D. J. Jones, and J. Ye, "Remote Transfer of Ultrastable Frequency References via Fiber Networks," *Rev. Sci. Instrum.* **78**:021101 (2007).
68. T. Rosenband, D. B. Hume, P. O. Schmidt, C. W. Chou, A. Brusch, L. Lorini, W. H. Oskay, R. E. Drullinger, T. M. Fortier, J. E. Stalnaker, S. A. Diddams, W. C. Swann, N. R. Newbury, W. M. Itano, D. J. Wineland, and J. C. Bergquist, "Frequency Ratio of Al^+ and Hg^+ Single-Ion Optical Clocks; Metrology at the 17th Decimal Place," *Science* **319**:1808–1812 (2008).
69. J. E. Stalnaker, Y. Le Coq, T. M. Fortier, S. A. Diddams, C. W. Oates, and L. Hollberg, "Measurement of Excited-State Transitions in Cold Calcium Atoms by Direct Femtosecond Frequency-Comb Spectroscopy," *Phys. Rev. A* **75**:040502 (R) (2007).
70. N. Poli, Z. W. Barber, N. D. Lemke, C. W. Oates, L. S. Ma, J. E. Stalnaker, T. M. Fortier, S. A. Diddams, L. Hollberg, J. C. Bergquist, A. Brusch, S. Jefferts, T. Heavner, and T. Parker, "Frequency Evaluation of the Doubly Forbidden $^1\text{S}_0 \rightarrow ^3\text{P}_0$ Transition in Bosonic Yb-174," *Phys. Rev. A* **77**:050501 (R) (2008).
71. R. J. Jones, K. D. Moll, M. J. Thorpe, and J. Ye, "Phase-Coherent Frequency Combs in the Vacuum Ultraviolet via High-Harmonic Generation inside a Femtosecond Enhancement Cavity," *Phys. Rev. Lett.* **94**:193201 (2005).
72. C. Göhle, T. Udem, M. Herrmann, J. Rauschenberger, R. Holzwarth, H. A. Schuessler, F. Krausz, and T. W. Hänsch, "A Frequency Comb in the Extreme Ultraviolet," *Nature* **436**:234–237 (2005).
73. See <http://www.mpg.mpg.de/~haensch/antihydrogen/index.html> (2009).
74. V. B. Braginsky and F. Ya. Khalili, *Quantum Measurement*, Cambridge University Press, Cambridge, UK (1992).

This page intentionally left blank.

DO NOT DUPLICATE

QUANTUM THEORY OF THE LASER

János A. Bergou^{a,b}
 Berthold-Georg Englert^{a,c,d}
 Melvin Lax^{e,*}
 Marian O. Scully^{a,c}
 Herbert Walther^{c,f,*}
 M. Suhail Zubairy^{a,g}

^a*Institute for Quantum Studies and Department of Physics
 Texas A&M University
 College Station, Texas*

^b*Department of Physics and Astronomy
 Hunter College of the City University of New York
 New York, New York*

^c*Max-Planck-Institut für Quantenoptik
 Garching bei München, Germany*

^d*Abteilung Quantenphysik der Universität Ulm
 Ulm, Germany*

^e*Department of Physics
 City College of the City University of New York
 New York, New York*

^f*Sektion Physik der Universität München
 Garching bei München, Germany*

^g*Department of Electronics
 Quaid-i-Azam University
 Islamabad, Pakistan*

23.1 GLOSSARY

Section 23.3

a_k, a_k^\dagger annihilation, creation operator for photons in the k th mode
 A Einstein coefficient for spontaneous emission
 $\mathbf{A}_k(\mathbf{r})$ k th electric mode function at position \mathbf{r}

*Deceased.

B	Einstein coefficient for absorption and stimulated emission
$\mathbf{B}(\mathbf{r}, t)$	magnetic field at position \mathbf{r} and time t
$\mathbf{B}_k(\mathbf{r})$	k th magnetic mode function at position \mathbf{r}
c	speed of light
\mathbf{e}_k	polarization unit vector of the k th mode
$\mathbf{E}_\perp(\mathbf{r}, t)$	transverse electric field at position \mathbf{r} and time t
h, \hbar	Planck's constant [$\hbar = h/(2\pi)$]
\mathbf{k}, k	wave vector, its length
$(d\mathbf{k})$	three-dimensional volume element in \mathbf{k} space
k_B	Boltzmann's constant
\mathbf{n}_k	propagation unit vector of the k th mode
N	photon number operator
N_e	number of excited-state atoms
N_g	number of ground-state atoms
\mathbf{r}	position vector
$(d\mathbf{r})$	three-dimensional volume element in \mathbf{r} space
t, dt	time, time interval
T	temperature
U_{1ph}	energy density for a one-photon state
dV	volume element
$ \text{vac}\rangle, \langle\text{vac} $	ket and bra of the photon vacuum
$\overline{w^2}$	mean square energy fluctuations
W	spectral-spatial energy density of blackbody radiation
α_k	coherent state amplitudes
$ \{\alpha\}_c\rangle$	ket of a coherent state
δ_{jk}	Kronecker's delta symbol
$\delta_\perp(\mathbf{r})$	transverse delta function at \mathbf{r} (a dyadic)
ϵ_0	dielectric constant [$\epsilon_0 = 1/(\mu_0 c^2) \approx 8.854 \times 10^{-12}$ F/m]
μ_0	permeability constant ($4\pi \times 10^{-7}$ H/m)
$\nu, d\nu$	light frequency, frequency interval
ν_k	frequency of the k th mode
ρ_{ph}	statistical operator of the photon field
ρ	probability amplitude for reflection
τ	probability amplitude for transmission
Ψ_k, Ψ_{jk}	probability amplitudes of the one-photon, two-photon states
$ \{\psi\}_1\rangle, \{\psi\}_2\rangle$	kets for one-photon, two-photon states
∇	gradient vector differential operator

Section 23.4

A	destruction operator for the field
\mathcal{A}	laser gain coefficient
a, a^\dagger	photon ladder operators
\mathcal{B}	laser saturation parameter
b_k, b_k^\dagger	destruction and creation operators for the reservoir modes
\mathcal{D}	largest eigenvalue of the laser equation

$F(t)$	field noise operator
$F_\alpha(t), F_\gamma$	noise operators associated with gain and loss, respectively
$F(1, x, y)$	hypergeometric function
g	radiant frequency stating the strength of atom-photon coupling
g_k	coupling coefficient between reservoir and field
$G^{(1)}(t_0 + t, t_0)$	field correlation function
$\mathcal{H}, \mathcal{H}_0, \mathcal{H}_1$	total, free, and interaction hamiltonians for atom-field interaction
k	wave vector for the field
K	kick operator
\mathcal{L}	superoperator for cavity damping
$M(\tau)$	superoperator describing the effect of a single inverted atom on the field
$\overline{n}, \overline{n^2}$	mean and mean squared number of photons in a laser
n_m	maximum of the photon distribution of a laser
n_{th}	mean number of photons in a thermal reservoir
N_{ex}	number of atoms traversing the cavity during the lifetime of the cavity field
$N(t_i, t, \tau)$	notch function
p	transition dipole moment
$p(n)$	photon distribution function
$P(\tau)$	distribution function for the interaction times
q	a nonnegative integer
Q	Mandel Q function
r	atom injection rate inside a laser and micromaser
$S(\omega)$	spectrum of the laser field
T	temperature
t_m	measurement time
U	time evolution operator
$U_{\text{af}}(\tau)$	atom-field time evolution operator
\mathcal{V}	interaction picture Hamiltonian
$\alpha(t, t')$	gain function of a laser
Γ	atomic decay rate via spontaneous emission
$\gamma, \gamma_a, \gamma_b$	atomic decay rates
$\delta_{nn'}$	Kronecker delta function
ν_k	frequency of the reservoir mode
σ_+, σ_-	ladder operators for a two-level atom
σ_z	atomic inversion operator
ω_0	radiant frequency of an atomic transition
λ	eigenvalue
λ_j	eigenvalues of the laser equation
$\rho(t)$	reduced density operator for the field
ρ_{at}	total density operator for the atom-field system
$\rho_{nn'}$	matrix elements of the field density operator
$\rho_n^{(k)}$	off-diagonal density matrix element
κ	cavity loss rate
χ	square of the ratio of vacuum Rabi frequency and the atomic decay rate
$\theta(t)$	step function
$\phi(t)$	phase of the field

Section 23.5

a, b, c, d	parameters appearing in Table 1
$\frac{E}{E}$	electric field strength mean field
$F(M), F_0$	free energy of a ferromagnet, its value for $M = 0$
g_k	undefined in Eqs. (126) and (127)
$G(E), G_0$	free energy of a laser, its value for $E = 0$
H	external magnetic field
H_{n0}	heating rate
K	one-fourth the spontaneous emission rate
K_{n0}	cooling rate
n_0	number of atoms in the Bose-Einstein condensate
$\langle n_0 \rangle$	mean number of atoms in the condensate
$\langle \dot{n}_0 \rangle$	time derivative of $\langle n_0 \rangle$
$\langle n_k \rangle_{n_0}$	average number of atoms in the k th excited state, given n_0 atoms in the condensate
N	total number of Bose atoms
N', N''	normalization constants in Table 1
M	magnetization of a ferromagnet
$P(M)$	probability density for a ferromagnet
$P(E)$	probability density for a laser
$P(\alpha, \alpha^*)$	P representation for the field
S	injected signal
T_c	critical temperature
W_k	heat bath density of states
x, y	$x = \text{Re } \alpha, y = \text{Im } \alpha$
X	zero-field susceptibility of a ferromagnet
$Z(T, N)$	canonical partition function
α	eigenvalue of the coherent state $ \alpha\rangle$
ε	scaled temperature (inversion) of a ferromagnet (laser) in Fig. 9
$\zeta(3)$	Riemann's zeta function $\zeta(3) = 1.2020569 \dots$
η	scaled thermodynamical variable (in Fig. 9)
$\langle \eta_k \rangle$	average occupation number of the k th heat bath oscillator
$\Theta(\cdot)$	Heaviside's unit step function
\mathbf{k}	undefined in Eq. (15)
ξ	laser analog of X
ρ_{n_0, n_0}	probability for having n_0 atoms in the condensate
$\dot{\rho}_{n_0, n_0}$	time derivative of ρ_{n_0, n_0}
σ	population inversion
σ_t	threshold inversion
$\Phi_e(\eta)$	scaled thermodynamical potential (in Fig. 9)
Ω	trap frequency

Section 23.6

A	phase-shifted destruction operator for a free-electron laser (FEL)
c_b, c_c	probability amplitudes for atom to be in levels $ b\rangle$ and $ c\rangle$, respectively

$\mathcal{D}(\theta)$	phase diffusion function
\mathcal{E}	slowly varying field amplitude
g	coupling constant for the electron-field interaction in an FEL
j	parameter for the gain in an FEL
k	wavevector for the laser field in an FEL
\mathcal{L}	linear gain ($i = j$) and cross-coupling ($i \neq j$) Liouville operators
m	mass of electron
$O(A, A^\dagger)$	arbitrary operator containing A, A^\dagger
p	electron momentum
\bar{p}	eigenvalue of electron momentum operator
P_b, P_c	probability of atom being in states $ b\rangle$ and $ c\rangle$, respectively
P_{emission}	probability of emission of radiation
$S(T)$	time-evolution operator for the electron-photon state
T	electron-photon interaction time
\mathcal{T}	time-ordering operator
z	electron coordinate
α_i	eigenvalue of the coherent state $ \alpha_i\rangle$
α_{ij}	constants depending upon parameters of gain medium in a correlated emission laser (CEL)
β	normalized electron momentum
γ	relativistic factor for an electron
Δ	atom-field detuning in lasing without inversion (LWI)
ν_1, ν_2	frequencies of the two modes in a CEL
$\rho(t_i)$	atomic density operator at initial time t_i
$\rho_{ij}^{(0)}$	initial values of the ij th atomic matrix elements
ρ_i	classical amplitude of the i th mode
$\rho(a_1, a_1^\dagger; a_2, a_2^\dagger)$	reduced density operator for the two-mode field in a CEL
θ_i	phase of the i th field
ω_0	microwave frequency
$\omega_a, \omega_b, \omega_c$	frequencies associated with atomic levels
λ_s	wavelength of the field emitted in an FEL
λ_w	the period of the magnetic wiggler
κ	gain coefficient
$\kappa_{a \rightarrow b}, \kappa_{a \rightarrow c}$	constants depending upon the matrix elements between the relevant levels
ϕ	relative phase between atomic levels
Φ	total phase angle in two-mode CEL schemes
ψ	relative phase $\Phi + \theta_1 - \theta_2$

23.2 INTRODUCTION

Most lasers, and in particular all commercially sold ones, emit electromagnetic radiation whose properties can be accounted for quite well by a *semiclassical* description. In such a treatment, quantum aspects (level spacings, oscillator strengths, etc.) of the matter (atoms, molecules, electron-hole pairs, etc.) that constitute the gain medium are essential, but those of the electromagnetic field are disregarded. Quantum properties of the radiation are, however, of decisive importance for laser systems “at the limit” which reach fundamental bounds for the linewidth, for the regularity of photon statistics, or for other quantities of interest.

Recognition and understanding of these fundamental limitations are furnished by the quantum theory of the laser, whose foundations were laid in the 1960s. The two main approaches, the master-equation formalism and the Langevin method—equivalent in the physical contents and supplementing each other like spouses—can be roughly, and somewhat superficially, associated with the Schrödinger and the Heisenberg pictures of quantum mechanics. The master-equation method corresponds to the former, the Langevin approach to the latter. Both are reviewed in Sec. 23.4, but more room is given to the master-equation treatment. This bias originates in our intention to present a parallel exposition for both the standard laser theory and the theory of the micromaser, which in turn is traditionally and most conveniently treated by master equations.

The micromaser, in which the dynamic is dominated by the strong coupling of a single mode of the radiation field to a single atomic dipole transition, is the prototype of an open, driven quantum system. Accordingly, micromaser experiments are *the* test ground for the quantum theory of the laser; therefore, micromaser theory deserves the special attention that it receives in Sec. 23.4.

As a logical and historical preparation, we recall in Sec. 23.3, the theoretical and experimental facts that are evidence for quantum properties of electromagnetic radiation in general, and the reality of photons in particular. Some special issues are discussed in Secs. 23.5 and 23.6. In Sec. 23.5, we stress the analogy between the threshold behavior of a laser and the phase transition of a ferromagnet, and note the recent lessons about Bose-Einstein condensates taught by this analogy. Section 23.6 summarizes the most important features of some exotic lasers and masers, which exploit atomic coherences or the quantum properties of the atomic center-of-mass motion. Basics of the so-called free-electron laser are reported as well.

The quantum theory of the laser is a central topic in the field of quantum optics. An in-depth understanding of the various facets of quantum optics can be gained by studying the pertinent textbooks.^{1–18}

23.3 SOME HISTORY OF THE PHOTON CONCEPT

Early History: Einstein's Light Quanta

Planck's formula of 1900¹⁹ marks the beginning of quantum mechanics, and in particular of the quantum theory of light. It reads

$$W = \frac{8\pi\nu^2}{c^3} \frac{h\nu}{\exp(h\nu/k_B T) - 1} \quad (1)$$

and relates the spectral-spatial energy density W of blackbody radiation to the frequency ν of the radiation and the temperature T of the blackbody. Boltzmann's constant k_B and Planck's constant h are conversion factors that turn temperature and frequency into energy, and c is the speed of light. A volume dV contains electromagnetic energy of the amount $W d\nu dV$ in the frequency range $\nu \cdots \nu + d\nu$.

The first factor in Eq. (1) is the density of electromagnetic modes. It obtains as a consequence of the classical wave theory of light and owes its simplicity to an implicit short-wavelength approximation. For wavelengths of the order of magnitude set by the size of the cavity that contains the radiation, appropriate corrections have to be made that reflect the shape and size of the cavity. This is of great importance in the context of micromasers, but need not concern us presently.

The second factor in Eq. (1) is the mean energy associated with radiation of frequency ν . It is a consequence of the quantum nature of light. In the limits of very high frequencies or very low ones, it turns into the respective factors of Wien²⁰ and Rayleigh-Jeans:^{21,22}

$$\begin{aligned} & \frac{h\nu}{\exp(h\nu/k_B T) - 1} && \text{(Planck)} \\ \rightarrow & \begin{cases} h\nu \exp\left(-\frac{h\nu}{k_B T}\right) & \text{for } \nu \gg \frac{k_B T}{h} && \text{(Wien)} \\ k_B T & \text{for } \nu \ll \frac{k_B T}{h} && \text{(Rayleigh-Jeans)} \end{cases} && (2) \end{aligned}$$

The relevant frequency scale is set by $k_B T/h$; at a temperature of $T = 288$ K it is about 6×10^{12} Hz, corresponding to a wavelength of $50 \mu\text{m}$.

Ironically, Planck—whose stroke of genius was to interpolate between the two limiting forms of Eq. (2)—was not convinced of the quantum nature of electromagnetic radiation until much later. Legend has it that it was the discovery of Compton scattering in 1923 that did it.²³ We are, however, getting ahead of the story.

The true significance of Planck's formula, Eq. (1), started to emerge only after Einstein²⁴ had drawn the conclusions that led him to his famous *light-quantum hypothesis* of 1905, the *annus mirabilis*. In Pauli's words,

He immediately applied [it] to the photoelectric effect and to Stokes' law for fluorescence, later also to the generation of secondary cathode rays by X-rays and to the prediction of the high frequency limit in the *Bremsstrahlung*.²⁵

Quite a truckload, indeed.

The conflict with the well-established wave theory of light was, of course, recognized immediately, and so the introduction of light quanta also gave birth to the wave-particle duality. Upon its extension to massive objects by de Broglie in 1923 to 1924,^{26–27} it was instrumental in Schrödinger's wave mechanics.²⁸

Taylor's 1909 experiment,²⁹ in which feeble light produced interference fringes, although at most one light quantum was present in the interferometer at any time, addressed the issue of wave-particle duality from a different angle. Its findings are succinctly summarized in Dirac's dictum that “a photon interferes only with itself”^{29,30}—a statement that became the innocent victim of misunderstanding and misquotation in the course of time.³¹

Another important step was taken the same year by Einstein.³² By an ingenious application of thermodynamic ideas to Planck's formula, in particular consequences of Boltzmann's relation between entropy and statistics, he derived an expression for the mean-square energy fluctuations w^2 of the radiation in a frequency interval $\nu \cdots \nu + d\nu$ and a volume dV :

$$\overline{w^2} = \left(\frac{c^3}{8\pi\nu^2} W + h\nu \right) W d\nu dV \quad (3)$$

where W is the spectral-spatial density of Eq. (1), so that $W d\nu dV$ is the mean energy in the frequency interval and volume under consideration. The first term is what one would get if classical electrodynamics accounted for all properties of radiation. There is no room for the second term in a wave theory of light; it is analogous to the fluctuations in the number of gas molecules occupying a given volume. This second term therefore supports Einstein's particle hypothesis of 1905,²⁴ in which electromagnetic energy is envisioned as being concentrated in localized lumps that are somehow distributed over the volume occupied by the electromagnetic wave.

Wave aspects (first term) and particle aspects (second term) enter Eq. (3) on equal footing. Since the thermodynamic considerations have no bias toward either one, one must conclude that Planck's formula, Eq. (1), is unbiased as well. Electromagnetic radiation is as much a particle phenomenon as it is a wave phenomenon.

Einstein left the center stage of quantum theory for some years, returning to it after completing his monumental work on general relativity. In 1913, Bohr's highly speculative postulates³³ had suddenly led to a preliminary understanding of many features of atomic spectra (the anomalous Zeeman effect was one big exception; it remained a bewildering puzzle for another decade). In the course, “quantum theory was liberated from the restriction to such particular systems as Planck's oscillators” (Pauli²⁵).

Here was the challenge to rederive Planck's formula from Bohr's postulates, assuming that they hold for arbitrary atomic systems. Einstein's famous paper of 1917³⁴ accomplished just that, and more.

He considered radiation in thermal equilibrium with a dilute gas of atoms at temperature T . We shall here give a simplified treatment that contains the essential features without accounting for all

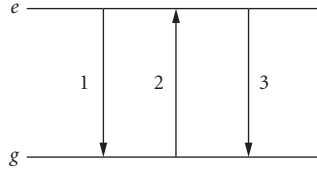


FIGURE 1 Transitions of three kinds happen between the excited level e and the ground level g : (1) spontaneous emission [see Eq. (4)], (2) absorption [see Eq. (5)], and (3) stimulated emission [see Eq. (6)].

details of lesser significance. Suppose that the energy spacing between two atomic levels e and g equals $h\nu$, so that the transition from the more energetic level e to the energetically lower level g (excited to ground state) is accompanied by the emission of a light quantum of frequency ν . According to Einstein, three processes are to be taken into account (see Fig. 1): spontaneous emission, absorption, and stimulated emission. The latter has no analog in Maxwell's electrodynamics.

Each of the three processes leads to a rate of change of the number of gas atoms in states e and g . Denoting these numbers by N_e and N_g we have the following contributions to their time derivatives. The *spontaneous emission* rate is proportional to the number of excited-state atoms:

$$\text{Spontaneous emission} \quad \frac{d}{dt}N_e = -\frac{d}{dt}N_g = -AN_e \quad (4)$$

where A is the first Einstein coefficient of the transition in question. The *absorption* rate is proportional to the number of ground-state atoms and to the energy density W of the radiation:

$$\text{Absorption} \quad \frac{d}{dt}N_e = -\frac{d}{dt}N_g = BWN_g \quad (5)$$

where B is the second Einstein coefficient. The *stimulated emission* rate is proportional to the number of excited state atoms and to the radiation energy density:

$$\text{Stimulated emission} \quad \frac{d}{dt}N_e = -\frac{d}{dt}N_g = BWN_e \quad (6)$$

where the same B coefficient appears as in the absorption rate.

The *detailed balance* between states e and g therefore requires

$$\text{Total} \quad \frac{d}{dt}N_e = -\frac{d}{dt}N_g = -AN_e + BWN_g - BWN_e = 0 \quad (7)$$

under the circumstances of thermal equilibrium. Therefore, W can be expressed in terms of the ratios A/B and N_g/N_e :

$$W = \frac{A/B}{N_g/N_e - 1} = \frac{A/B}{\exp(h\nu/k_b T) - 1} \quad (8)$$

where the second equality recognizes that Boltzmann's factor relates N_e to N_g . [The absence of additional weights here is the main simplification alluded to before; if taken into account, these weights would also require two closely related B coefficients in Eqs. (5) and (6).] In Eq. (8) we encounter the

denominator of Planck's factor from Eq. (2), and Planck's formula, Eq. (1), is recovered in full if the relation

$$A = \frac{8\pi h\nu^3}{c^3} B \quad (9)$$

is imposed on Einstein's coefficients.

The main ingredients in this derivation of Eq. (1) are the postulate of the process of stimulated emission, with a strength proportional to the density of radiation energy, and the relation between the coefficients for spontaneous emission and stimulated emission, Eq. (9).

All of this is well remembered, but there was in fact more to the 1917 paper.³⁴ It also contains a treatment of the momentum exchange between atoms and light quanta, and Einstein succeeded in demonstrating that the Maxwell velocity distribution of the atoms is consistent with the recoil they suffer when absorbing and emitting quanta. Insights gained in his study of Brownian motion (another seminal paper of 1905³⁵) were crucial for this success. When taken together, the considerations about energy balance and those concerning momentum balance are much more convincing than either one could have been alone.

The discovery of the Compton effect in 1923²³ finally convinced Bohr and other skeptics of the reality of the particlelike aspects possessed by light. But Bohr, who until then was decidedly opposed to Einstein's light-quantum hypothesis and the consequent wave-particle duality, did not give in without a last try. Together with Kramers and Slater,³⁶ he hypothesized that perhaps energy-momentum conservation does not hold for each individual scattering event, but only in a statistical sense for a large ensemble. Then one could account for Compton's data without conceding a particle nature to light in general, and X rays in particular. The refined measurements that were immediately carried out by Bothe and Geiger³⁷ showed, however, that this hypothesis is wrong: energy and momentum are conserved in each scattering event, not just statistically.

And then there was, also in 1924, Bose's seminal observation that it is possible to derive Planck's radiation law [Eq. (1)] from purely corpuscular arguments without invoking at all the wave properties of light resulting from Maxwell's field equations.^{38,39} The main ingredient in Bose's argument was the indistinguishability of the particles in question and a new way of counting them—now universally known as *Bose-Einstein statistics*—which pays careful attention to what is implied by their being indistinguishable. In the case of light quanta, an additional feature is that their number is not conserved, because light is easily emitted and absorbed. Massive particles (atoms, molecules, etc.), by contrast, are conserved; therefore, as Einstein emphasized,^{40–42} their indistinguishability has further consequences, of which the phenomenon of Bose-Einstein condensation (or should one rather say “Einstein condensation of a Bose gas”?) is the most striking one.

Quantum Electrodynamics

Theoretical studies of the quantum nature of light had a much more solid basis after Dirac's introduction of quantum electrodynamics (QED) in his seminal paper of 1927.⁴³ The basic ingredients of QED were all present in Dirac's formulation, although it is true that a consistent understanding of QED was not available until renormalized QED was developed 20 years later (see the papers reprinted in Ref. 44). In particular, the photon concept was clarified in the sense described in the following paragraphs.

The infinite number of degrees of freedom of the electromagnetic field—an operator field in QED—become manageable with the aid of a mode expansion. For the transverse part $\mathbf{E}_\perp(\mathbf{r}, t)$ of the electric field, for example, it reads

$$\mathbf{E}_\perp(\mathbf{r}, t) = \sum_k \sqrt{\frac{\hbar\nu_k}{2\epsilon_0}} [a_k(t)\mathbf{A}_k(\mathbf{r}) + a_k^\dagger(t)\mathbf{A}_k^*(\mathbf{r})] \quad (10)$$

The mode functions $\mathbf{A}_k(\mathbf{r})$ are complex vector functions of the position vector \mathbf{r} that are eigenfunctions of the Laplace differential operator,

$$-\nabla^2 \mathbf{A}_k(\mathbf{r}) = \left(\frac{2\pi\nu_k}{c} \right)^2 \mathbf{A}_k(\mathbf{r}) \quad (11)$$

where the eigenvalue is determined by the frequency ν_k (>0) of the mode in question. The boundary conditions that the electric and magnetic field must obey at conducting surfaces imply respective boundary conditions on the $\mathbf{A}_k(\mathbf{r})$ s.

The corresponding mode expansion for the magnetic field is given by

$$\mathbf{B}(\mathbf{r}, t) = \sum_k \sqrt{\frac{\mu_0 \hbar \nu_k}{2}} [-ia_k(t)\mathbf{B}_k(\mathbf{r}) + ia_k^\dagger(t)\mathbf{B}_k^*(\mathbf{r})] \quad (12)$$

where

$$\mathbf{B}_k(\mathbf{r}) = \frac{c}{2\pi\nu_k} \nabla \times \mathbf{A}_k(\mathbf{r}); \quad \mathbf{A}_k(\mathbf{r}) = \frac{c}{2\pi\nu_k} \nabla \times \mathbf{B}_k(\mathbf{r}) \quad (13)$$

relates the two kinds of mode functions to each other.

Among others, the mode functions $\mathbf{A}_k(\mathbf{r})$ have the following important properties:

$$\text{Transverse} \quad \nabla \cdot \mathbf{A}_k(\mathbf{r}) = 0 \quad (14a)$$

$$\text{Orthonormal} \quad \int (d\mathbf{r}) \mathbf{A}_j^*(\mathbf{r}) \cdot \mathbf{A}_k(\mathbf{r}) = \delta_{jk} \quad (14b)$$

$$\text{Complete} \quad \sum_k \mathbf{A}_k(\mathbf{r}) \mathbf{A}_k^*(\mathbf{r}') = \delta_\perp(\mathbf{r} - \mathbf{r}') \quad (14c)$$

The same statements hold for the $\mathbf{B}_k(\mathbf{r})$ s as well. The property in Eq. (14a) states the radiation-gauge condition. The integration in Eq. (14b) covers the entire volume bounded by the conducting surfaces just mentioned; the eigenvalue equation Eq. (11) holds inside this volume, the so-called quantization volume. In the completeness relation in Eq. (14c), both positions \mathbf{r} and \mathbf{r}' are inside the quantization volume, and δ_\perp is the transverse delta function, a dyadic that is explicitly given by

$$\delta_\perp(\mathbf{r}) = \int \frac{(d\mathbf{k})}{(2\pi)^3} \exp(i\mathbf{k} \cdot \mathbf{r}) \left(1 - \frac{\mathbf{k}\mathbf{k}}{k^2} \right) \quad (15)$$

where 1 is the unit dyadic and $k = \sqrt{\mathbf{k} \cdot \mathbf{k}}$ is the length of the wave vector \mathbf{k} integrated over. The transverse character of $\delta_\perp(\mathbf{r})$ ensures the consistency of the properties in Eqs. (14a) and (14c).

The time dependence of $\mathbf{E}_\perp(\mathbf{r}, t)$ and $\mathbf{B}(\mathbf{r}, t)$ stems from the ladder operators $a_k(t)$ and $a_k^\dagger(t)$, which obey the bosonic equal-time commutation relations

$$[a_j, a_k] = 0 \quad [a_j, a_k^\dagger] = \delta_{jk} \quad [a_j^\dagger, a_k^\dagger] = 0 \quad (16)$$

The photon number operator

$$N = \sum_k a_k^\dagger a_k \quad (17)$$

has eigenvalues $N' = 0, 1, 2, \dots$; its eigenstates with $N' = 1$ are the one-photon states, those with $N' = 2$ are the two-photon states, and so on. The unique eigenstate with $N' = 0$ is the photon vacuum.

We denote its ket by $|\mathbf{vac}\rangle$. It is, of course, the joint eigenstate of all *annihilation operators* a_k with eigenvalue zero:

$$a_k |\mathbf{vac}\rangle = 0 \quad \text{for all } k \quad (18)$$

Application of the *creation operator* a_k^\dagger to $|\mathbf{vac}\rangle$ yields a state with one photon in the k th mode:

$$a_k^\dagger |\mathbf{vac}\rangle = \{\text{state with 1 photon of the type } k\} \quad (19)$$

More generally, the ket of a pure one-photon state is of the form

$$|\{\psi\}_1\rangle = \sum_k \psi_k a_k^\dagger |\mathbf{vac}\rangle \quad \text{with} \quad \sum_k |\psi_k|^2 = 1 \quad (20)$$

where $|\psi_k|^2$ is the probability for finding the photon in the k th mode. Similarly, the kets of pure two-photon states have the structure

$$|\{\psi\}_2\rangle = \frac{1}{\sqrt{2}} \sum_{j,k} \psi_{jk} a_j^\dagger a_k^\dagger |\mathbf{vac}\rangle \quad \text{with} \quad \psi_{jk} = \psi_{kj} \quad \text{and} \quad \sum_{j,k} |\psi_{jk}|^2 = 1 \quad (21)$$

and analogous expressions apply to pure states with 3, 4, 5, . . . photons.

Einstein's light quanta are one-photon states of a particular kind. In a manner of speaking, they are localized lumps of electromagnetic energy. In technical terms this means that the energy density

$$\begin{aligned} U_{\text{1ph}}(\mathbf{r}, t) &= \left\langle \{\psi\}_1 \left| : \left[\mathbf{E}^2 / (2\epsilon_0) + \left(\frac{\mu_0}{2} \right) \mathbf{B}^2 \right] : \right| \{\psi\}_1 \right\rangle \\ &= \frac{\hbar}{2} \sum_{j,k} \psi_j^* \psi_k \sqrt{v_j v_k} (\mathbf{A}_k^* \cdot \mathbf{A}_j + \mathbf{B}_k^* \cdot \mathbf{B}_j) \end{aligned} \quad (22)$$

is essentially nonzero in a relatively small spatial region only. The time dependence is carried by the probability amplitudes ψ_k , the spatial dependence by the mode functions \mathbf{A}_k and \mathbf{B}_k . An arbitrarily sharp localization is not possible, but it is also not needed. The pair of colons :: symbolize the injunction to order the operator in between in the *normal* way: all creation operators a_k^\dagger to the left of all annihilation operators a_k . This normal ordering is an elementary feature of renormalized QED.

At high frequencies, or when the quantization volume is unbounded, the eigenvalues of $-\nabla^2$ in Eq. (11) are so dense that the summations in Eqs. (10), (12), and (14b) are effectively integrations, and the Kronecker delta symbol in Eq. (14b) is a Dirac delta function. Under these circumstances, it is often natural to choose plane waves

$$\mathbf{A}_k(\mathbf{r}) \sim \mathbf{e}_k \exp\left(i \frac{2\pi \nu_k}{c} \mathbf{n}_k \cdot \mathbf{r}\right) \quad \mathbf{B}_k(\mathbf{r}) \sim \mathbf{n}_k \times \mathbf{e}_k \exp\left(i \frac{2\pi \nu_k}{c} \mathbf{n}_k \cdot \mathbf{r}\right) \quad (23)$$

for the mode functions. The unit vector \mathbf{e}_k that specifies the polarization is orthogonal to the unit vector \mathbf{n}_k that specifies the direction of propagation.

With

$$\psi_k(t) \sim \exp(-i2\pi \nu_k t) \quad (24)$$

in Eq. (22), one then meets exponential factors of the form

$$\exp\left[i \frac{2\pi \nu_k}{c} (\mathbf{n}_k \cdot \mathbf{r} - ct)\right]$$

As a consequence, an einsteinian light quantum propagates without dispersion, which is the anticipated behavior.

The one-photon energy density in Eq. (22) illustrates the general feature that quantum-mechanical probabilities (the ψ_k s) with their interference properties appear together with the classical interference patterns of superposed mode functions [the $\mathbf{A}_k(\mathbf{r})$ s and $\mathbf{B}_k(\mathbf{r})$ s]. In other words, interference phenomena of two kinds are present in QED: (1) the classical interference of electromagnetic fields in the three-dimensional \mathbf{r} space, and (2) the quantum interference of alternatives in the so-called Fock space; that is, the Hilbert space spanned by the photon vacuum $|\mathbf{vac}\rangle$, the one-photon states $|\{\psi\}_1\rangle$, the two-photon states $|\{\psi\}_2\rangle$, and all multiphoton states.

In the early days of QED, this coexistence of classical interferences and quantum interferences was a research topic, to which Fermi's paper of 1929 "Sulla teoria quantistica delle frange di interferenza" is a timeless contribution.⁴⁵ He demonstrated, at the example of Lippmann fringes, a very general property of single-photon interference patterns: the photon-counting rates, as determined from quantum-mechanical probabilities, are proportional to the corresponding classical intensities.

Electromagnetic radiation is easily emitted and absorbed by antennas, processes that change the number of photons. Accordingly, the number of photons is not a conserved quantity, and therefore states of different photon numbers can be superposed. Particularly important are the *coherent states*

$$|\{\alpha\}_c\rangle = \exp\left(-\frac{1}{2}\sum_k |\alpha_k|^2 + \sum_k \alpha_k a_k^\dagger\right) |\mathbf{vac}\rangle \quad (25)$$

that are characterized by a set $\{\alpha\}_c$ of complex amplitudes α_k . As revealed in Glauber's 1963 papers,^{46–48} they play a central role in the coherence theory of light.

Since the coherent states are common eigenstates of the annihilation operators

$$a_k |\{\alpha\}_c\rangle = |\{\alpha\}_c\rangle \alpha_k \quad (26)$$

the expectation values of the electric and magnetic field operators of Eqs. (10) and (12)

$$\langle\{\alpha(t)\}_c | \mathbf{E}_\perp(\mathbf{r}, t) | \{\alpha(t)\}_c\rangle = \sum_k \sqrt{\frac{h\nu_k}{2\epsilon_0}} [\alpha_k(t) \mathbf{A}_k(\mathbf{r}) + \alpha_k^*(t) \mathbf{A}_k^*(\mathbf{r})] \quad (27)$$

$$\langle\{\alpha(t)\}_c | \mathbf{B}_\perp(\mathbf{r}, t) | \{\alpha(t)\}_c\rangle = \sum_k \sqrt{\frac{\mu_0 h\nu_k}{2}} [-i\alpha_k(t) \mathbf{B}_k(\mathbf{r}) + i\alpha_k^*(t) \mathbf{B}_k^*(\mathbf{r})]$$

have the appearance of classical Maxwell fields. In more general terms, if the statistical operator ρ_{ph} of the photonic degrees of freedom—in other words, the statistical operator of the radiation field—is a mixture of (projectors to) coherent states

$$\rho_{\text{ph}} = \sum_{\{\alpha\}_c} |\{\alpha\}_c\rangle w(\{\alpha\}_c) \langle\{\alpha\}_c| \quad (28)$$

with

$$w(\{\alpha\}_c) \geq 0 \quad \text{and} \quad \sum_{\{\alpha\}_c} w(\{\alpha\}_c) = 1 \quad (29)$$

then the electromagnetic field described by ρ_{ph} is very similar to a classical Maxwell field. Turned around, this says that whenever it is impossible to write a given ρ_{ph} in the form of Eq. (28), then some statistical properties of the radiation are decidedly nonclassical.

During the 20-year period from Dirac's paper of 1927 to the Shelter Island conference in 1947, QED remained in a preliminary state that allowed various studies—the most important ones included the Weisskopf-Wigner treatment of spontaneous emission⁴⁹ and Weisskopf's discovery that the self-energy of the photon is logarithmically divergent⁵⁰—although the not-yet-understood divergences

were very troublesome. The measurement by Lamb and Retherford⁵¹ of what is now universally known as the *Lamb shift*, first reported at the Shelter Island conference, was the crucial experimental fact that triggered the rapid development of renormalized QED by Schwinger, Feynman, and others.

Theoretical calculations of the Lamb shift rely heavily on the quantum properties of the electromagnetic field, and their marvelous agreement with the experimental data proves convincingly that these quantum properties are a physical reality. In other words, photons exist. The same remark applies to the theoretical and experimental values of the anomalous magnetic moment of the electron, one of the early triumphs of Schwinger's renormalized QED,⁵² which finally explained an anomaly in the spectra of hydrogen and deuterium that Pasternack had observed in 1938⁵³ and a discrepancy in the measurements by Millman and Kusch⁵⁴ of nuclear magnetic moments.

Photon-Photon Correlations

Interferometers that exploit not the spatial intensity variations (or, equivalently, the photon-detection probabilities) but correlations between intensities at spatially separated positions became important tools in astronomy and spectroscopy after the discovery of the Hanbury-Brown-Twiss (HB&T) effect in 1954.^{55–57} A textbook account of its classical theory is given in Sec. 4.3 of Ref. 58.

In more recent years, the availability of single-photon detectors made it possible to study the HB&T effect at the two-photon level. The essentials are depicted in Fig. 2. Two light quanta are incident on a half-transparent mirror from different directions, such that they arrive simultaneously. If their frequency contents are the same, it is fundamentally impossible to tell if an outgoing light quantum was reflected or transmitted. This indistinguishability of the two light quanta is of decisive importance in the situation where one is in each output channel. The two cases *both reflected* and *both transmitted* are then indistinguishable and, according to the laws of quantum mechanics, the corresponding probability amplitudes must be added.

Now, denoting the probability amplitudes for single-photon reflection and transmission by ρ and τ , respectively, the probability for one light quantum in each output port is given by

$$|\rho^2 + \tau^2|^2 = \left| \left(\frac{1}{\sqrt{2}} \right)^2 + \left(\frac{i}{\sqrt{2}} \right)^2 \right|^2 = 0 \quad (30)$$

where we make use of $\rho = 1/\sqrt{2}$ and $\tau = i/\sqrt{2}$, which are the values appropriate for a symmetric half-transparent mirror. Thus, the situation of one light quantum in each output port does not occur. Behind the half-transparent mirror, one always finds both light quanta in the same output port.

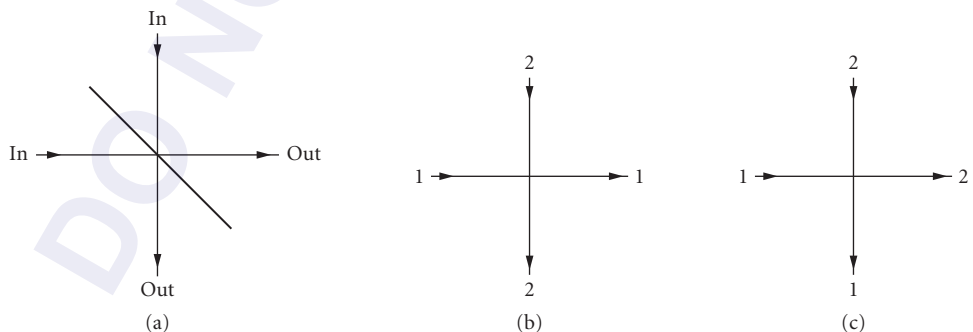


FIGURE 2 Essentials of the Hanbury-Brown-Twiss effect at the two-photon level. (a) Two light quanta are simultaneously incident at a symmetric half-transparent mirror. To obtain one quantum in each output port, the quanta must be either both transmitted (b) or both reflected (c). The probability amplitude for (b) is $(1/\sqrt{2})^2 = 1/2$, that for (c) is $(i/\sqrt{2})^2 = -1/2$. If the two cases are indistinguishable, the total amplitude is $1/2 - 1/2 = 0$.

If the light quanta arrived at very different times, rather than simultaneously, they would be distinguishable and one would have to add probabilities instead of probability amplitudes, so that

$$|\rho^2|^2 + |\tau^2|^2 = \frac{1}{2} \quad (31)$$

would replace Eq. (30). Clearly, there are intermediate stages at which the temporal separation is a fraction of the temporal coherence length and the two quanta are neither fully distinguishable nor utterly indistinguishable. The probability for one light quantum in each output port is then a function of the separation, a function that vanishes when the separation does.

Experiments that test these considerations⁵⁹ employ correlated photon pairs produced by a process known as *parametric downconversion*. Roughly speaking, inside a crystal that has no inversion symmetry a high-frequency photon is absorbed and two lower-frequency photons are emitted, whereby the conservation of energy and momentum imposes geometrical restrictions on the possible propagation directions of the three photons involved. Downconversion sources with a high luminosity are available.^{60,61}

The HB&T effect of Fig. 2 as well as closely related phenomena are crucial in many experiments in which entangled photons are a central ingredient. In particular, the recent realizations^{62,63} of schemes for quantum teleportation⁶⁴ and the experiment⁶⁵ that demonstrated the practical feasibility of quantum-dense coding⁶⁶ are worth mentioning here.

None of these exciting developments could be understood without the quantum properties of radiation. Since the photon concept, in the sense of the discussion of Eqs. (10), (12), (19), (20), and so on is an immediate consequence of these quantum properties, the existence of photons is an established experimental fact beyond reasonable doubt.

23.4 QUANTUM THEORY OF THE LASER

The quantum theory of laser radiation is a problem in nonequilibrium statistical mechanics. There are several alternative, but ultimately equivalent, approaches to the characterization of the field inside the resonator. As is customary in the Scully-Lamb quantum theory, we describe the state of the laser field by a density operator.^{67,68} In this section our main focus is on the review of the equation of motion, the so-called master equation, for this density operator as it emerges from an underlying physical model, with statistical considerations and some simplifying assumptions. The alternative procedure based on the quantum theory of noise sources introduced in Refs. 69 and 70 and summarized in Refs. 71 to 74 will also be briefly reviewed at the end of the section. For a recent, more detailed overview of the quantum Langevin point of view, we refer the reader to Ref. 75.

In general, a laser model should be based on the interaction of multimode fields with multilevel atoms as the active medium, and a detailed consideration of all possible processes among all the levels involved should be given. Decay channels and decay rates, in particular, play a crucial role in determining the threshold inversion and, thus, the necessary pumping rates. Of course, the pumping mechanisms themselves can be quite complicated. It is well established that a closed two-level model cannot exhibit inversion and, hence, lasing. In order to achieve inversion three- and four-level pumping schemes are employed routinely. On the other hand, to illustrate the essential quantum features a single mode field can serve as paradigm. The single-mode laser field inside the resonator interacts with one particular transition of the multilevel system—the lasing transition—and the role of the entire complicated level structure is to establish inversion on this transition—that is, to put more atoms in the upper level than there are in the lower one. If one is not interested in the details of how the inversion builds up, it is possible to adopt a much simpler approach than the consideration of a multilevel-multimode system. In order to understand the quantum features of the single-mode field it is sufficient to focus only on the two levels of the lasing transition and their interaction with the laser field. In this approach, the effect of pumping, decay, and so on in the multilevel structure is simply replaced by an initial condition; it is assumed that the atom is in its upper state immediately before the interaction with the laser mode begins. Since here we are primarily interested

in the quantum signatures of the laser field and not in the largely classical aspects of cavity design, pumping mechanisms, and so on, we shall follow this simpler route from the beginning. The model that accounts for the resonant interaction of a two-level atom with a single quantized mode in a cavity was introduced by Jaynes and Cummings.⁷⁶

We shall make an attempt to present the material in a tutorial way. We first derive an expression for the change of the field density operator due to the interaction with a single two-level atom, initially in its upper state, using the Jaynes-Cummings model. This expression will serve as the seed for both the laser and micromaser theories. We next briefly review how to account for cavity losses by using standard methods for modeling the linear dissipation loss of the cavity field due to mirror transmission. Then we show that with some additional assumptions the single-atom-single-mode approach can be used directly to derive what has become known as the Scully-Lamb master equation for the more traditional case of the laser and the micromaser. The additional assumptions include the Markov approximation or, equivalently, the existence of very different time scales for the atomic and field dynamics so that adiabatic elimination of the atoms and introduction of coarse-grained time evolution for the field become possible. The main difference between the laser and micromaser theories is that the interaction time of the active atoms with the field is governed by the lifetime of the atoms in the laser and by the transit time of the atoms through the cavity in the micromaser. In the laser case, the atoms decay out of the lasing levels into some far-removed other levels, and they are available for the lasing transition during their lifetime on the average. In the micromaser case, the transit time is approximately the same for all atoms in a monoenergetic pumping beam. Thus, the laser involves an extra averaging over the random interaction times. If we model the random interaction times by a Poisson distribution and average the change of the field density operator that is due to a single atom—the kick—over the distribution of the interaction times, we obtain the master equation of a laser from that of the micromaser. Historically, of course, the development was just the opposite: the master equation was derived in the context of the laser much earlier. However, it is instructive to see how the individual Rabi oscillations of single nondecaying atoms with a fixed interaction time, as in the micromaser, give rise to the saturating, nonoscillatory collective behavior of an ensemble of atoms, as in the laser, upon averaging over the interaction times. As applications of this fully quantized treatment we study the photon statistics, the linewidth, and spectral properties. Finally, we briefly discuss other approaches to the quantum theory of the laser.

Time Evolution of the Field in the Jaynes-Cummings Model

We shall consider the interaction of a single two-level atom with a single quantized mode of a resonator using the rotating-wave approximation (for a recent review of the Jaynes-Cummings model see Ref. 77). The arrangement is shown in Fig. 3.

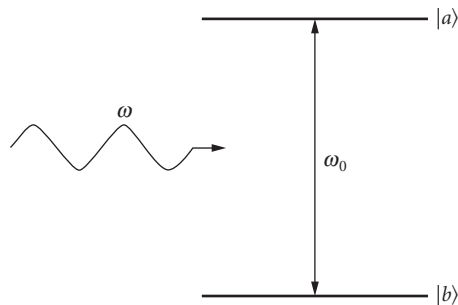


FIGURE 3 Scheme of a two-level atom interacting with a single mode quantized field. The text focuses on the resonant case, $\omega = \omega_0$.

The hamiltonian for this system is given by

$$\mathcal{H} = \mathcal{H}_0 + \mathcal{H}_1 \quad (32)$$

where

$$\mathcal{H}_0 = \frac{1}{2} \hbar \omega_0 \sigma_z + \hbar \omega_0 a^\dagger a \quad (33)$$

and

$$\mathcal{H}_1 = \hbar g (\sigma_+ a + a^\dagger \sigma_-) \quad (34)$$

Here a and a^\dagger are the annihilation and creation operators for the mode. The upper level of the lasing transition is denoted by $|a\rangle$ and the lower level by $|b\rangle$. The atomic lowering and raising operators are expressed in terms of the state vectors as $\sigma_- = (\sigma_+)^{\dagger} = |b\rangle\langle a|$ and the population operator as $\sigma_z = |a\rangle\langle a| - |b\rangle\langle b|$. ω_0 is the frequency of the $|a\rangle - |b\rangle$ transition. For simplicity we assume perfect resonance with the mode. Finally, g is the coupling constant between the atom and the mode. In terms of atomic and field quantities it is given by $g = p\sqrt{\omega_0/2\hbar\epsilon_0 v}$ where p is the transition dipole moment and v the quantization volume (the volume of the active medium, in our case). Again, for simplicity, p is assumed to be real.

In the interaction picture with respect to \mathcal{H}_0 , the interaction hamiltonian becomes

$$\mathcal{V} = \exp\left(\frac{i\mathcal{H}_0 t}{\hbar}\right) \mathcal{H}_1 \exp\left(-\frac{i\mathcal{H}_0 t}{\hbar}\right) = \hbar g (\sigma_+ a + a^\dagger \sigma_-) \equiv \mathcal{H}_1 \quad (35)$$

since \mathcal{H}_0 and \mathcal{H}_1 commute.

In the two-dimensional Hilbert space spanned by the state vectors $|a\rangle$ and $|b\rangle$ the interaction hamiltonian can be written as

$$\mathcal{H}_1 = \hbar g \begin{pmatrix} 0 & a \\ a^\dagger & 0 \end{pmatrix} \quad (36)$$

The time evolution operator for the coupled atom-field system satisfies the equation of motion in this picture

$$\frac{i\hbar dU}{dt} = \mathcal{V} U \quad (37)$$

and since \mathcal{H}_1 is time independent the solution is formally

$$U(\tau) = \exp\left(-\frac{i}{\hbar} \mathcal{V} \tau\right) \quad (38)$$

Using the properties of the σ_- and σ_+ matrices, it is easy to show that $U(\tau)$ can be written in the preceding 2×2 matrix representation as

$$U(\tau) = \begin{pmatrix} \cos(g\tau\sqrt{aa^\dagger}) & -i\frac{\sin(g\tau\sqrt{aa^\dagger})}{\sqrt{aa^\dagger}} a \\ -ia^\dagger \frac{\sin(g\tau\sqrt{aa^\dagger})}{\sqrt{aa^\dagger}} & \cos(g\tau\sqrt{a^\dagger a}) \end{pmatrix} \quad (39)$$

Let us now assume that initially, at t_0 , the atom is in its upper state given by the atomic density operator $\varrho_{\text{at}}(t_0) = |a\rangle\langle a|$ and the field is in an arbitrary state which, in general, can be described by the density operator $\varrho(t_0)$, so that the joint atom-field system is characterized by the initial density operator $\varrho_{\text{af}}(t_0) = \varrho_{\text{at}}(t_0) \otimes \varrho(t_0)$. After the interaction, $\varrho_{\text{af}}(t_0 + \tau) = U(\tau)\varrho_{\text{af}}(t_0)U(\tau)^{-1}$. Our main interest here is in the evolution of the field density operator. This we obtain if we trace the atom-field density operator over the atomic states, yielding

$$\begin{aligned} \varrho(t_0 + \tau) &= \text{Tr}_{\text{atom}}[\varrho_{\text{af}}(t_0 + \tau)] \\ &= \cos(g\tau\sqrt{aa^\dagger})\varrho(t_0)\cos(g\tau\sqrt{aa^\dagger}) + a^\dagger \frac{\sin(g\tau\sqrt{aa^\dagger})}{\sqrt{aa^\dagger}}\varrho(t_0)\frac{\sin(g\tau\sqrt{aa^\dagger})}{\sqrt{aa^\dagger}}a \\ &\equiv M(\tau)\varrho(t_0) \end{aligned} \quad (40)$$

Here in the last step we just introduced the superoperator M , which describes the effect of a single inverted atom on the field and is a key ingredient of laser and micromaser theory. The matrix elements in photon number representation take the form

$$(M(\tau)\varrho)_{nn'} = A_{nn'}(\tau)\varrho_{nn'} + B_{n-1n'-1}(\tau)\varrho_{n-1n'-1} \quad (41)$$

where the coefficients are given by

$$A_{nn'}(\tau) = \cos(g\tau\sqrt{n+1})\cos(g\tau\sqrt{n'+1}) \quad (42)$$

$$B_{nn'}(\tau) = \sin(g\tau\sqrt{n+1})\sin(g\tau\sqrt{n'+1}) \quad (43)$$

For later purposes, we also introduce the change in the state of the field due to the interaction with a single inverted atom as

$$\varrho(t_0 + \tau) - \varrho(t_0) = M(\tau)\varrho(t_0) - \varrho(t_0) = (M - 1)\varrho(t_0) \equiv K\varrho(t_0) \quad (44)$$

The operator K , sometimes called the *kick operator*, contains all the information we will need to build the quantum theory of the single-mode laser and micromaser. In matrix representation, $[K(\tau)\varrho]_{nn'} = [A_{nn'}(\tau) - \delta_{nn'}]\varrho_{nn'} + B_{n-1n'-1}(\tau)\varrho_{n-1n'-1}$.

For more elaborate systems (multimode lasers driven by multilevel atoms, for example) one cannot give M in such a simple analytical form, but the principle remains always the same. One should find the superoperator M or, equivalently, the kick operator $K = M - 1$ which gives the action of a single (possibly multimode) atom on the (possibly multimode) field from the general expression $\varrho(t_0 + \tau) = \text{Tr}_{\text{atom}}[U_{\text{af}}(\tau)\varrho_{\text{at}}(t_0) \otimes \varrho(t_0)U_{\text{af}}(\tau)^{-1}] \equiv M(\tau)\varrho(t_0)$. In order to determine the full time-evolution operator $U_{\text{af}}(\tau)$ of the coupled atom-field system, however, one usually needs to resort to approximation methods, such as perturbation theory, in the more complicated multilevel-multimode cases.

Derivation of the Scully-Lamb Master Equation

In 1954, Gordon, Zeiger, and Townes showed that coherent electromagnetic radiation can be generated in the radio frequency range by the maser (acronym for *microwave amplification by stimulated emission of radiation*).⁷⁸ The first maser action was observed in a beam of ammonia.⁷⁹ The maser principle was extended to the optical domain by Schawlow and Townes,⁸⁰ and also by Prokhorov

and Basov,⁸¹ thus obtaining a laser (acronym for *light amplification by stimulated emission of radiation*). A laser consists of a large ensemble of inverted atoms interacting resonantly with the electromagnetic field inside a cavity. The cavity selects only a specific set of modes corresponding to a discrete set of eigenfrequencies. The active atoms—that is, the ones that are pumped to the upper level of the laser transition—are in resonance with one of the eigenfrequencies of the cavity in the case of the single-mode laser and with a finite set of frequencies in the case of the multimode laser. As discussed in the introduction to this section, for the discussion of the essential quantum features of the radiation field of a laser it is sufficient to confine our treatment to the single-mode case, and that is what we will do for the remainder of this section. A resonant electromagnetic field gives rise to stimulated emission, and the atoms thereby transfer their energy to the radiation field. The emitted radiation is still at resonance. If the upper level is sufficiently populated, this radiation gives rise to further transitions in other atoms. In this way, all the excitation energy of the atoms is transferred to the single mode of the radiation field.

The first pulsed laser operation was demonstrated by Maiman in ruby.⁸² The first continuous wave (CW) laser, a He-Ne gas laser, was built by Javan et al.⁸³ Since then a large variety of systems have been demonstrated to exhibit lasing action. Coherent radiation has been generated this way over a frequency domain ranging from infrared to soft X rays. These include dye lasers, chemical lasers, solid-state lasers, and semiconductor lasers.

Many of the laser properties can be understood on the basis of a semiclassical theory. In such a theory the radiation field is treated classically, but the active medium is given a full quantum-mechanical treatment. Such a theory can readily explain threshold and saturation, transient dynamics, and general dependence on the external parameters (pumping and losses). It is not our aim here to give an account of the semiclassical theory; therefore, we just refer the reader to the ever instructive and wonderfully written seminal paper by Lamb⁸⁴ and a more extended version in Ref. 67. Although quantum effects play only a minor role in usual practical laser applications because of the large mean photon numbers, they are essential for the understanding of the properties of micromasers, in which excited two-level atoms interact one after the other with a single mode of the radiation field.⁸⁵ Nevertheless, the quantum properties of the laser field are of fundamental interest as well. They have been thoroughly investigated theoretically with respect to the photon statistics and the spectrum of the laser. In particular, the quantum limitation of the laser linewidth caused by the inevitably noisy contribution of spontaneous emission has attracted much attention. It gives rise to the so-called Schawlow-Townes linewidth, which is inversely proportional to the laser intensity (see Ref. 80). Because of the importance of stable coherent signals for various high-precision measurements, the problem of the intrinsic quantum-limited linewidth has gained renewed interest recently, and the investigations have been extended to cover bad-cavity lasers and several more exotic systems. In this review, however, we shall restrict ourselves to good-cavity lasers in which the cavity damping time is long compared to all other relevant time scales, and present a fully quantized theory of the most fundamental features.

Cavity Losses To account for the decay of the cavity field through the output mirror of the cavity, we simply borrow the corresponding result from reservoir theory. Its usage has become fairly standard in laser physics and quantum optics (see, for example, Refs. 67 and 68), and here we just quote the general expression without actually deriving it.

$$\left(\frac{d\varrho}{dt}\right)_{\text{loss}} = \mathcal{L}\varrho \equiv -\frac{\kappa}{2}n_{\text{th}}(aa^\dagger\varrho + \varrho aa^\dagger - 2a^\dagger\varrho a) - \frac{\kappa}{2}(n_{\text{th}}+1)(a^\dagger a\varrho + \varrho a^\dagger a - 2a\varrho a^\dagger) \quad (45)$$

This equation refers to a loss reservoir which is in thermal equilibrium at temperature T , with n_{th} being the mean number of thermal photons $n_{\text{th}} = [\exp(\hbar\omega_0/kT) - 1]^{-1}$, and κ is the cavity damping rate. For the laser case, it is sufficient to take the limiting case of a zero temperature reservoir since $\hbar\omega_0 \gg kT$ and n_{th} is exponentially small. We obtain this limit by substituting $n_{\text{th}} = 0$ into Eq. (45). For the description of most micromaser experiments, however, we need the finite temperature version, since even at very low temperatures the thermal photon number is comparable to the total number of photons in the cavity.

The Laser Master Equation After introducing the loss part of the master equation, we now turn our attention to the part that stems from the interaction with the gain reservoir. The gain reservoir is modeled by an ensemble of initially excited two-level atoms allowed to interact with the single-mode cavity field. A central role in our subsequent discussions will be played by the so-called kick operator, $K = M - 1$, describing the change of the field density operator due to the interaction with a single atom. This quantity was introduced in Eq. (44). While in the micromaser case the effect of each of the atoms can be represented by the same kick operator, since in a monoenergetic pumping beam each atom has the same interaction time with the cavity field, this is no longer the case for a laser. In a typical CW gas laser, such as the He-Ne laser, atoms are excited to the upper level of the lasing transition at random times and, more important, they can also interact with the field for a random length of time. The interaction time thus becomes a random variable. Since at any given time the number of atoms is large (about 10^6 to 10^7 active atoms in the lasing volume of a CW He-Ne laser), it is a legitimate approach to describe their effect on the field by an average kick operator. We can arrive at the interaction-time-averaged master equation quickly if we take the average of Eq. (44) with respect to the interaction time τ :

$$(M-1)\varrho(t) = \int_0^\infty d\tau P(\tau)(M(\tau)-1)\varrho(t) \quad (46)$$

where the distribution function for the interaction time $P(\tau)$ is defined as

$$P(\tau) = \gamma e^{-\gamma\tau} \quad (47)$$

This distribution function corresponds to the exponential decay law. Individual atoms can decay from the lasing levels at completely random times, but for an ensemble of atoms the probability of finding an initially excited atom still in the lasing levels in the time interval $(\tau, \tau + d\tau)$ is given by Eq. (47). With increasing τ it is increasingly likely that the atoms have decayed outside the lasing transition. Also note that our model corresponds to an open system: the atoms decay to other non-lasing levels both from the upper state $|a\rangle$ and the lower state $|b\rangle$, and, in addition, we assume that decay rate γ is the same for both levels, as indicated in Fig. 4.

Obviously, these restrictions can be relaxed and, indeed, there are various more general models available. For example, the upper level $|a\rangle$ can have two decay channels. It can decay to the lower level $|b\rangle$ and to levels outside the lasing transition. Or, in some of the most efficient lasing schemes, the lower level decays much faster than the upper one, $\gamma_b \gg \gamma_a$. In these schemes virtually no population builds up in the lower level; hence, saturation of the lasing transition occurs at much higher

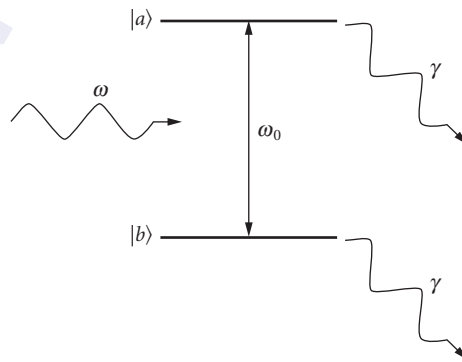


FIGURE 4 Scheme of a two-level atom, with atomic decay permitted, interacting with a single-mode quantized field.

intensities than in lasers with equal decay rates for both levels. These generalizations, however, are easily accounted for (see Ref. 67) and it is not our concern here to provide the most general treatment possible. Instead, we want to focus on the essential quantum features of the laser field and employ the simple two-level model with equal decay rates for both.

The formal averaging in Eq. (46) can be performed most easily if we transform to a particular representation of the density matrix. For our purposes, the photon number representation suffices, although other options are readily available and some of them will be summarized briefly at the end of this section. Taking the n, n' elements of Eq. (46) and using Eqs. (42) and (43), the averaging yields

$$\begin{aligned} [(M-1)\rho]_{nn'} = & -\frac{\chi(n+1+n'+1)+\chi^2(n-n')^2}{1+2\chi(n+1+n'+1)+\chi^2(n-n')^2}\rho_{nn'} \\ & +\frac{2\chi\sqrt{nn'}}{1+2\chi(n+n')+\chi^2(n-n')^2}\rho_{n-1n'-1} \end{aligned} \quad (48)$$

where the notation $\chi = g^2/\gamma^2$ is introduced. Finally, taking into account cavity losses $(\mathcal{L}\rho)_{nn'}$ from Eq. (45), with $n_{\text{th}} = 0$, we obtain the following master equation for our quantum-mechanical laser model:

$$\begin{aligned} \dot{\rho}_{nn'} = & -\frac{\mathcal{N}'_{nn'}\mathcal{A}}{1+\mathcal{N}'_{nn'}\mathcal{B}/\mathcal{A}}\rho_{nn'} + \frac{\sqrt{nn'}\mathcal{A}}{1+\mathcal{N}_{n-1n'-1}\mathcal{B}/\mathcal{A}}\rho_{n-1n'-1} \\ & -\frac{\kappa}{2}(n+n')\rho_{nn'} + \kappa\sqrt{(n+1)(n'+1)}\rho_{n+1n'+1} \end{aligned} \quad (49)$$

Here we introduced the original notations of the Scully-Lamb theory—the linear gain coefficient:

$$\mathcal{A} = 2r\chi \quad (50)$$

the self-saturation coefficient:

$$\mathcal{B} = 4\chi\mathcal{A} \quad (51)$$

and the dimensionless factors:

$$\mathcal{N}' = \frac{1}{2}(n+1+n'+1) + \frac{(n-n')^2\mathcal{B}}{8\mathcal{A}} \quad (52)$$

and

$$\mathcal{N} = \frac{1}{2}(n+1+n'+1) + \frac{(n-n')^2\mathcal{B}}{16\mathcal{A}} \quad (53)$$

Equation (49) is the Scully-Lamb master equation, which is the central equation of the quantum theory of the laser. Along with the notations introduced in Eqs. (50) to (53), it constitutes the main result of this section and serves as the starting point for our treatment of the quantum features of the laser. Among the specific problems we shall consider are the photon statistics, which is the physical information contained in the diagonal elements, and the spectrum, which is the physical information contained in the off-diagonal elements of the field density matrix.

The Micromaser Master Equation The development of the single-atom maser or micro-maser plays a particularly important role in cavity quantum electrodynamics because it realizes one of the most fundamental models, the Jaynes-Cummings hamiltonian. The experimental situation⁸⁶ is very close to the idealized case of a single two-level atom interacting with a single-mode quantized field, as previously discussed, and allows a detailed study of fundamental quantum effects in the atom-field interaction.

In the micromaser, a stream of two-level atoms is injected into a superconducting microwave cavity of very high quality Q . The injection rate r is low enough to ensure that at most only one atom is present inside the cavity at any given time and that most of the time the cavity is empty. The decay time of the high- Q cavity field is very long compared to both the interaction time τ , which is set by the transit time of atoms through the cavity, and the inverse of the single-photon Rabi frequency g^{-1} . In typical experimental situations, however, $g\tau \approx 1$. Therefore, a field is built up in the cavity provided the interval between atomic injections does not significantly exceed the cavity decay time. Sustained oscillation is possible with less than one atom on the average in the cavity.

In addition to the progress in constructing superconducting cavities, advances in the selective preparation of highly excited hydrogenlike atomic states, called *Rydberg states*, have made possible the realization of the micromaser. In Rydberg atoms the probability of induced transitions between adjacent states is very large, and the atoms may undergo several Rabi cycles—that is, several periodic energy exchanges between the atom and the cavity field may take place in the high- Q cavity. The lifetime of Rydberg states for spontaneous emission decay is also very long, and atomic decay can be neglected during the transit time in the cavity.

Here we set out to derive a master equation for the micromaser. For this, we consider a single-mode resonator into which two-level atoms are injected in their upper states. Due to the different time scales, the effect of cavity damping can be neglected during the interaction. Then the effect of a single atom on the field density operator, injected at t_i and interacting with the field for a time τ , is given by Eq. (44) with t_0 replaced by t_i . If several atoms are injected during a time interval Δt which is still short on the time scale governed by the cavity decay time κ^{-1} but long on the time scale of the interaction time τ , then the cumulative effect on the field is simply the sum of changes

$$\Delta \varrho = \sum_{t \leq t_i \leq t + \Delta t} (M(\tau) - 1) \varrho(t_i) = r \int_t^{t + \Delta t} (M(\tau) - 1) \varrho(t_i) dt_i \quad (54)$$

where in the last step we turned the sum into an integral by using the injection rate r as the number of excited-state atoms entering the cavity per unit time.

At this point we introduce some of the most important approximations of laser physics, or reservoir theory in general—the so-called Markov or adiabatic approximation and coarse time graining. These approximations are based on the existence of three very different time scales in the problem. First, the interaction time τ for individual atoms is of the order of the inverse single-photon Rabi frequency g^{-1} , and is the shortest of all. In fact, on the time scale set by the other relevant parameters, it appears as a delta-function-like kick to the state of the field with the kick operator K of Eq. (44). The second time scale is set by r^{-1} , the average time separation between atomic injections. It is supposed to be long compared to τ but short compared to the cavity-damping time κ^{-1} . Thus, we have the following hierarchy of timescales: $\tau \ll r^{-1} \ll \kappa^{-1}$. When we turned the sum in Eq. (54) into an integral we already tacitly assumed that there is a time scale on which the injection appears to be quasicontinuous. We now see that it is the time scale set by κ^{-1} . This is the time scale that governs the time evolution of the cavity field. In the evaluation of the integral in Eq. (54) we assume that Δt is an intermediate time interval such that $r^{-1} < \Delta t < \kappa^{-1}$. Then during this interval the state of the field does not change appreciably, and we can replace $\varrho(t_i)$ on the right-hand side of Eq. (54) by $\varrho(t)$. This is the essential step in transforming the integral equation into a differential equation, and it constitutes what is called the *Markov approximation*. It is also called the *adiabatic approximation* since the field changes very slowly (adiabatically) on the time scale set by the atoms. As a result, $\varrho(t)$ can now be taken out of the integral, and the integration in Eq. (54) yields $\Delta \varrho = r \Delta t (M - 1) \varrho(t)$. Dividing both sides by Δt , we obtain

$$\frac{\Delta \varrho}{\Delta t} = r(M(\tau) - 1) \varrho \quad (55)$$

The left-hand side is not a true derivative; it only appears to be one on the time scale of the cavity decay time. However, if we are interested in the large-scale dynamics of the field, we can still regard it as a good approximation to a time derivative. It is called the *coarse-grained derivative* and the Eq. (55)

now properly describes the time rate of change of the field due to the interaction with an ensemble of active atoms, the gain reservoir.

Equation (55) gives the time rate of change of the field density operator due to the gain reservoir $(d\varrho/dt)_{\text{gain}}$. To this, we add the time rate of change of the density operator due to the cavity losses by hand. For the parameters of the Garching micromaser experiment, $T = 0.5$ K and $\omega_0/2\pi = 21.5$ GHz, yielding $n_{\text{th}} = 0.15$. The thermal background cannot be neglected since, as we shall see, the steady-state field contains but a few photons. The complete master equation for the micromaser, including both gain and loss, is then simply the sum of Eqs. (55) and (45):

$$\frac{d\varrho}{dt} = \left(\frac{d\varrho}{dt}\right)_{\text{gain}} + \left(\frac{d\varrho}{dt}\right)_{\text{loss}} = r(M(\tau)-1)\varrho + \mathcal{L}\varrho \quad (56)$$

For later purposes, we also give the master equation in matrix representation:

$$\begin{aligned} \frac{d\varrho_{nn'}}{dt} = & r[(A_{nn'}(\tau)-1)\varrho_{nn'} + B_{n-1n'-1}(\tau)\varrho_{n-1n'-1}] \\ & - \frac{\kappa}{2}n_{\text{th}}[(n+n'+2)\varrho_{nn'} - 2\sqrt{nn'}\varrho_{n-1n'-1}] \\ & - \frac{\kappa}{2}(n_{\text{th}}+1)[(n+n')\varrho_{nn'} - 2\sqrt{(n+1)(n'+1)}\varrho_{n+1n'+1}] \end{aligned} \quad (57)$$

Here $A_{nn'}(\tau)$ and $B_{nn'}(\tau)$ are given by Eqs. (42) and (43). Equation (57) is identical to the one obtained by more standard methods⁸⁷ and employed widely in the context of micromasers. It forms the basis for most studies (with a few notable exceptions, as discussed at the end of this section) on the quantum statistical properties of the micromaser and, naturally, it will be our starting point as well.

Physics on the Main Diagonal: Photon Statistics

To begin to bring to light the physical consequences of the laser and maser master equations, we shall first focus on the diagonal elements of the field density matrix ϱ_{nn} , which give us the photon-number distribution since $\varrho_{nn} = p(n)$ is the probability of finding n photons in the cavity mode. The case of the laser is sufficiently different from that of the micromaser that we shall deal with them separately.

Laser Photon Statistics Taking the diagonal $n = n'$ elements in Eq. (49) and regrouping the terms, we obtain the following equation for the photon-number probabilities:

$$\dot{p}(n) = -\frac{(n+1)\mathcal{A}}{1+(n+1)\mathcal{B}/\mathcal{A}}p(n) + \kappa(n+1)p(n+1) + \frac{n\mathcal{A}}{1+n\mathcal{B}/\mathcal{A}}p(n-1) - \kappa np(n) \quad (58)$$

Here the overdot stands for the time derivative. Note that diagonal elements are coupled only to diagonal elements. This holds quite generally; Eq. (49) describes coupling along the same diagonal only. For example, elements on the first side diagonal are coupled to other elements on the first side diagonal, and so on, and in general only elements with the same difference $n - n'$ are coupled. The quantity $k = n - n'$ corresponds to elements on the k th side diagonal. Elements on different diagonals are not coupled, which greatly simplifies the solution of laser-related problems.

Before we begin the solution of Eq. (58), we want to give a simple intuitive physical picture of the processes it describes in terms of a probability flow diagram, shown in Fig. 5.

The left-hand side is the rate of change of the probability of finding n photons in the cavity. The right-hand side contains the physical processes that contribute to the change. Each process is

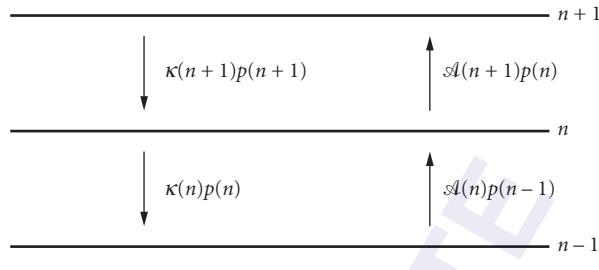


FIGURE 5 Probability flow diagram for the laser.

represented by an arrow in the diagram. The processes are proportional to the probability of the state they are starting from and this will be the starting point of the arrow. The tip of the arrow points to the state the process is leading to. There are two kinds of elementary processes: emission of a photon into the cavity mode (upward arrows) and loss of a photon from the cavity through the output mirror (downward arrows). Furthermore, the processes that start from a given state and end in a different one will decrease the probability of the state they are starting from and increase the probability of the state they are ending at. For example, the first term on the right-hand side describes emission of a photon into the cavity mode provided there are n photons already present before the emission event takes place. Since there will be $n + 1$ photons after the emission, this process decreases the probability of finding n photons, hence the minus sign. The total emission rate $\mathcal{A}(n+1)$ has a contribution from stimulated emission, the $\mathcal{A}(n)$ term, and another one from spontaneous emission, the \mathcal{A} term. The third term on the right-hand side corresponds similarly to emission, conditioned on the presence of $n - 1$ photons in the cavity initially. After the emission, there will be n photons, hence the plus sign. The second term describes the loss of a photon through the cavity mirror, provided there are n photons initially. After the escape of a photon, there will be $n - 1$ photons; therefore, this term decreases $p(n)$. Finally, the last term corresponds similarly to a loss process, with initially $n + 1$ photons in the cavity. After the escape of a photon there will be n photons left, so this term increases the probability $p(n)$ of finding n photons in the cavity.

After this brief discussion of the meaning of the individual terms, we now turn our attention to the solution of the equation. Although it is possible to obtain a rather general time-dependent solution to Eq. (58), our main interest here is in the steady-state properties of the field. To obtain the steady-state photon statistics, we replace the time derivative with zero. Note that the right-hand side of the equation is of the form $F(n + 1) - F(n)$, where

$$F(n) = \kappa n p(n) - \frac{n\mathcal{A}}{1 + n\mathcal{B}/\mathcal{A}} p(n-1) \quad (59)$$

simply meaning that in steady-state $F(n + 1) = F(n)$. In other words, $F(n)$ is independent of n and is, therefore, a constant c . Furthermore, the equation $F(n) = c$ has normalizable solution only for $c = 0$. From Eq. (59) we then immediately obtain

$$p(n) = \frac{\mathcal{A}/\kappa}{1 + n\mathcal{B}/\mathcal{A}} p(n-1) \quad (60)$$

which is a very simple two-term recurrence relation to determine the photon-number distribution. Before we present the solution, a remark is called for here. The fact that $F(n) = 0$ and $F(n + 1) = 0$ hold separately is called the *condition of detailed balance*. As a consequence, we do not need to deal with all four processes affecting $p(n)$. It is sufficient to balance the processes connecting a pair of adjacent levels in Fig. 5, and instead of solving the general three-term recurrence relation resulting from the steady-state version of Eq. (58), it is enough to solve the much simpler two-term recursion, Eq. (60).

It is instructive to investigate the photon statistics in some limiting cases before discussing the general solution. Below threshold the linear approximation holds. Since only very small n states are

populated appreciably, the denominator on the right-hand side of Eq. (60) can be replaced by unity in view of $n\mathcal{B}/\mathcal{A} \ll 1$. Then

$$p(n) = p(0) \left(\frac{\mathcal{A}}{\kappa} \right)^n \quad (61)$$

The normalization condition $\sum_{n=0}^{\infty} p(n) = 1$ determines the constant $p(0)$, yielding $p(0) = (1 - \mathcal{A}/\kappa)$. Finally,

$$p(n) = \left(1 - \frac{\mathcal{A}}{\kappa} \right) \left(\frac{\mathcal{A}}{\kappa} \right)^n \quad (62)$$

Clearly, the condition of existence for this type of solution is $\mathcal{A} < \kappa$. Therefore, $\mathcal{A} = \kappa$ is the threshold condition for the laser. At threshold, the photon statistics change qualitatively and very rapidly in a narrow region of the pumping parameter. It should also be noted that below threshold the distribution function Eq. (62) is essentially of thermal character. If we introduce an effective temperature T defined by

$$\exp\left(-\frac{\hbar\omega_0}{kT}\right) = \frac{\mathcal{A}}{\kappa} \quad (63)$$

we can cast Eq. (62) to the form

$$p(n) = \left[1 - \exp\left(-\frac{\hbar\omega_0}{kT}\right) \right] \exp\left(-\frac{n\hbar\omega_0}{kT}\right) \quad (64)$$

This is just the photon-number distribution of a single mode in thermal equilibrium with a thermal reservoir at temperature T . The inclusion of a finite temperature-loss reservoir to represent cavity losses will not alter this conclusion about the region below threshold.

There is no really good analytical approximation for the region around threshold, although the lowest-order expansion of the denominator in Eq. (60) yields some insight. The solution with this condition is given by

$$p(n) = p(0) \left(\frac{\mathcal{A}}{\kappa} \right)^n \prod_{k=0}^{n-1} \left(1 - \frac{k\mathcal{B}}{\mathcal{A}} \right) \quad (65)$$

This equation clearly breaks down for $n > \mathcal{A}/\mathcal{B} = n_{\max}$, where $p(n)$ becomes negative. The resulting distribution is quite broad, exhibiting a long plateau and a rapid cutoff at n_{\max} . The broad plateau means that many values of n are approximately equally likely; therefore, the intensity fluctuations are large around threshold. The most likely value of $n = n_{\text{opt}}$ can be obtained from the condition $p(n_{\text{opt}} - 1) = p(n_{\text{opt}})$ since $p(n)$ is increasing before $n = n_{\text{opt}}$ and decreasing afterwards. This condition yields $n_{\text{opt}} = (\mathcal{A} - \kappa)/\mathcal{B}$, which is smaller by the factor κ/\mathcal{A} than the value obtained from the full nonlinear equation [cf. Eq. (70) following].

The third region of special interest is the one far above threshold. In this region, $\mathcal{A}/\kappa \gg 1$ and the n values contributing the most to the distribution function are the ones for which $n \gg \mathcal{A}/\mathcal{B}$. We can then neglect 1 in the denominator of Eq. (60), yielding

$$p(n) = \exp\left(-\bar{n} \frac{\bar{n}^n}{n!}\right) \quad (66)$$

with $\bar{n} = \mathcal{A}^2/(\kappa\mathcal{B})$. Thus, the photon statistics far above threshold are poissonian, the same as for a coherent state. This, however, does not mean that far above threshold the laser is in a coherent state. As we shall see later, the off-diagonal elements of the density matrix remain different from those of a coherent state for all regimes of operation.

After developing an intuitive understanding of the three characteristically different regimes of operation, we give the general solution of Eq. (60), valid in all three regimes:

$$p(n) = p(0) \prod_{k=1}^n \frac{(\mathcal{A}/\kappa)}{(1+k\mathcal{B}/\mathcal{A})} \quad (67)$$

The normalization constant $p(0)$ may be expressed in terms of the confluent hypergeometric function

$$\begin{aligned} p(0) &= \left[\sum_{n=0}^{\infty} \frac{\left(\frac{\mathcal{A}}{\mathcal{B}}\right)! \left(\frac{\mathcal{A}^2}{\mathcal{B}\kappa}\right)^n}{\left(n + \frac{\mathcal{A}}{\mathcal{B}}\right)!} \right]^{-1} \\ &= \left[F\left(1; \frac{\mathcal{A}}{\mathcal{B}} + 1; \frac{\mathcal{A}^2}{\mathcal{B}\kappa}\right) \right]^{-1} \end{aligned} \quad (68)$$

In Fig. 6, the photon-number distribution is displayed for various regimes of operation.

It is interesting to note that $p(n)$ is a product of n factors of the form $(\mathcal{A}/\kappa)/(1+k\mathcal{B}/\mathcal{A})$. This is an increasing function of k as long as the factors are larger than 1 and decreasing afterward. The maximum of the distribution function can be found from the condition

$$\frac{(\mathcal{A}/\kappa)}{(1+n_m\mathcal{B}/\mathcal{A})} = 1 \quad (69)$$

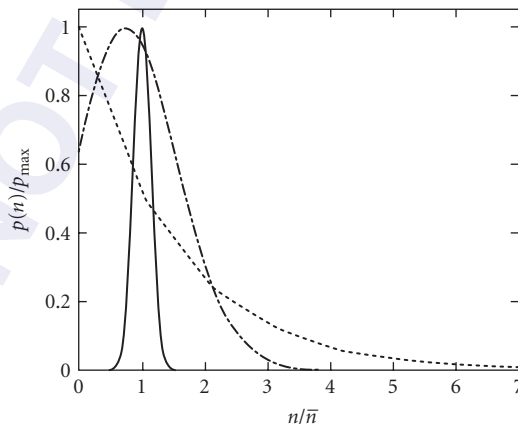


FIGURE 6 Photon statistics of the laser for various regimes of operation. Dotted line: laser 50 percent below threshold ($\mathcal{A}/\kappa = 1/2$), distribution thermal in character, Mandel $Q_M = 1$. Dot-dashed line: laser 10 percent above threshold, ($\mathcal{A}/\kappa = 1.1$), distribution is broad, $Q_M = 5$. Solid line: laser 100 percent above threshold ($\mathcal{A}/\kappa = 2$), distribution super-Poissonian, $Q_M = 1$. Since in the various regimes the actual photon numbers and $p(n)$ values differ by several orders of magnitude, in order to comparatively display the statistics in one figure, we plot $p(n)/p_{\max}$ versus n/\sqrt{n} . This way the maximum of each curve is 1 and unity on the horizontal axis corresponds to the average photon number. $\mathcal{A}/\mathcal{B} = 50$ is used for all plots, for illustrative purposes only. For more realistic values, the above-threshold distributions are much narrower on this scale.

or

$$n_m = \frac{\mathcal{A}}{\mathcal{B}} \frac{\mathcal{A} - \kappa}{\kappa} \quad (70)$$

Clearly, for $\mathcal{A} < \kappa$ the maximum is at $n = 0$ and the distribution is monotonically decreasing, which is characteristic of a thermal distribution. This is in agreement with the previous findings for the below-threshold region. Around $\mathcal{A} = \kappa$ the distribution quickly changes its character. The factor \mathcal{A}/\mathcal{B} governs the magnitude of the photon number while $(\mathcal{A} - \kappa)/\kappa$ is a measure of how far away from threshold the laser is operating. Typical values for CW gas lasers (the He-Ne laser, for example) are $\mathcal{A}/\mathcal{B} \approx 10^6$ and $\kappa \approx 10^6$ Hz. Around threshold, $\mathcal{A} = \kappa$ and the factors appearing in $p(n)$, given by Eq. (69), are effectively unity for a broad range of n . For example, for $\mathcal{A}/\kappa = 1.001$ (i.e., one-tenth of a percent above threshold), the factors are slightly above 1 for $1 < n < 1000$. So, in the threshold region, the distribution very quickly changes from a thermal one, dominated by the vacuum state, to a broad distribution with large intensity fluctuations. Farther above threshold the distribution becomes more and more peaked around n_m and becomes essentially poissonian for $\mathcal{A}/\kappa > 2$.

It is easy to obtain the mean photon number \bar{n} from Eq. (67):

$$\bar{n} = \sum_{n=0}^{\infty} n p(n) \frac{\mathcal{A}}{\mathcal{B}} = \frac{\mathcal{A} - \kappa}{\kappa} + \frac{\mathcal{A}}{\mathcal{B}} p(0) \quad (71)$$

Above threshold, $p(0) \ll 1$ and the last term becomes quickly negligible. Then \bar{n} coincides with n_m , the maximum of the distribution. We can obtain \bar{n}^2 similarly. The result is

$$\bar{n}^2 = \bar{n} + \frac{\mathcal{A}^2}{\mathcal{B}\kappa} \quad (72)$$

Using Eq. (71), the variance can be expressed as

$$\sigma^2 = \bar{n}^2 - \bar{n} = \bar{n} \frac{\mathcal{A}}{\mathcal{A} - \kappa} \quad (73)$$

From here we see that the variance always exceeds that of a poissonian distribution ($\sigma^2 > \bar{n}$), but it approaches one far above threshold. A good characterization of the photon-number distribution is given by the Mandel Q_M parameter:

$$Q_M = \frac{\bar{n}^2 - \bar{n}^2}{\bar{n}} - 1 \quad (74)$$

For our case, it is given by

$$Q_M = \frac{\kappa}{\mathcal{A} - \kappa} \quad (75)$$

Since $Q_M > 0$ above threshold, the field is superpoissonian. Very far above threshold, when $\mathcal{A} \gg \kappa$, Q_M approaches zero, which is characteristic of a poissonian distribution, again in agreement with our discussion of the far-above-threshold region.

Micromaser Photon Statistics As a first application of the micromaser master equation, Eq. (56), we shall study the steady-state photon statistics arising from it. To this end we take the diagonal $n = n'$ elements, and after regrouping the terms we obtain:

$$\begin{aligned} \dot{p}(n) = & -N_{\text{ex}} \sin^2(g\tau\sqrt{n+1})p(n) + (n_{\text{th}} + 1)(n+1)p(n+1) - n_{\text{th}}(n+1)p(n) \\ & + N_{\text{ex}} \sin^2(g\tau\sqrt{n})p(n-1) - (n_{\text{th}} + 1)np(n) + n_{\text{th}}np(n-1) \end{aligned} \quad (76)$$

Here the overdot stands for derivative with respect to the scaled time $t' = \kappa t$. $N_{\text{ex}} = r/\kappa$ is the number of atoms traversing the cavity during the lifetime of the cavity field, and the diagonal matrix elements of the density operator $p(n) = \rho_{nn}$ are the probabilities of finding n photons in the cavity. The various processes in the right-hand side of this equation are again visualized in Fig. 5. They have the structure $F(n+1) - F(n)$ where $F(n+1)$ corresponds to the processes connecting $p(n+1)$ to $p(n)$. In the steady state the left-hand side is zero and $F(n+1) = F(n)$, yielding $F(n) = \text{constant}$. The only normalizable solution to the photon statistics arises when this constant is zero, $F(n) = 0$. Once again, this is the condition of detailed balance because the probability flows between adjacent levels are separately balanced. More explicitly, this leads to the following recurrence relation for the photon-number probabilities:

$$p(n) = \frac{N_{\text{ex}} \sin^2(g\tau\sqrt{n})/n + n_{\text{th}}}{n_{\text{th}} + 1} p(n-1) \quad (77)$$

The solution to this simple recurrence relation is straightforward:

$$p(n) = p(0) \prod_{i=1}^n \frac{N_{\text{ex}} \sin^2(g\tau\sqrt{i})/i + n_{\text{th}}}{n_{\text{th}} + 1} \quad (78)$$

where $p(0)$ is determined from the normalization condition $\sum_{n=0}^{\infty} p(n) = 1$. The photon-number distribution $p(n)$ versus n can be multip peaked in certain parameter regimes. This can be easily understood on the basis of Fig. 5. The gain processes (upward arrows) balance the loss (downward arrows). Since the gain is an oscillatory function of n , several individual peaks [with the property $p(n+i) = p(n)$] will develop for those values of n where the gain perfectly balances the loss. The resulting mean photon number (first moment of the distribution) and photon-number fluctuations (second moment Q_M) versus the scaled interaction parameter $\theta = g\tau\sqrt{N_{\text{ex}}}$ are displayed in Fig. 7.

The mean photon number is an oscillatory function of the scaled interaction time. The oscillations correspond to subsequent Rabi cycles the atoms are undergoing in the cavity as function of the interaction time. The first threshold occurs at $\theta = 1$; the higher ones occur where θ is approximately an integer multiple of 2π . Around the thresholds the micromaser field is superpoissonian; in the parameter region between the thresholds it is subpoissonian, which is a signature of its nonclassicality.

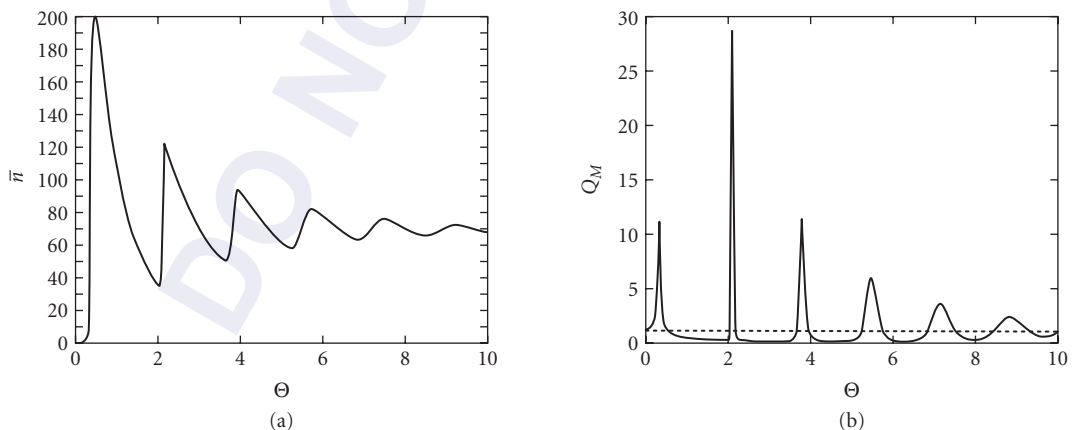


FIGURE 7 (a) Mean photon number and (b) Mandel Q_M parameter versus interaction parameter for the micromaser, for $N_{\text{ex}} = 200$.

Physics off the Main Diagonal: Spectrum and Linewidth

In the subsection on laser photon statistics we have shown that the laser has poissonian photon statistics far above threshold, just as a coherent state would. However, it is erroneous to conclude from this that the laser field is in a coherent state far above threshold. Only the intensity of the laser becomes stabilized due to a delicate balance of the nonlinear gain and loss. Any deviation from the steady-state intensity induces a change that tries to restore the steady-state value; as a result, the intensity is locked to this value. The phase of the field, on the other hand, evolves freely and is not locked to any particular value. In fact, it performs a random walk due to the separate spontaneous emission events. Each such event contributes a small random phase change to the instantaneous phase of the field. The mean of these changes averages to zero, but the mean of the square of these changes remains finite. As a result, the phase undergoes a diffusionlike process and will become uniformly distributed over the 2π interval. Any information contained in the instantaneous phase will be erased in this process. The time scale for the decay of the phase information is given by the rate of the phase diffusion. In the following, we shall determine this characteristic time scale, the so-called phase diffusion constant for the laser and micromaser.

Spectral Properties of the Laser Field The decay of the phase information can be directly read out from the temporal behavior of the two-time correlation function of the field amplitude:

$$g^{(1)}(t_0 + t, t_0) = \frac{\langle a^\dagger(t_0 + t)a(t_0) \rangle}{\langle a^\dagger(t_0)a(t_0) \rangle} \quad (79)$$

With increasing time difference between their amplitudes, the fields become less and less correlated since spontaneous emission randomizes their phases. At steady state, the two-time correlation function [Eq. (79)] depends only on the time difference t and is independent of the choice of the initial time t_0 .

A quantum regression “theorem,”⁸⁸ based on a system-reservoir factorization of the density matrix, was developed to permit the time evolution of the two-time correlation function at steady state to be calculated from the time evolution of the single-time correlation function for a Markov process. The unusual success of this procedure (see Ref. 89 for a comparison between experiment and theory for the phase linewidth, the intensity linewidth, the photo-count distribution, and the spectral moments) requires additional explanation. This was supplied in a proof that regression is valid when the system is markovian.⁹⁰ In the quantum case, the system is only approximately markovian. But this assumption has already been made in all cases for which solutions have been found. Therefore, it is sufficient for us to study the time evolution of the amplitude itself:

$$\langle a^\dagger(t) \rangle = \sum_{n=0}^{\infty} \sqrt{n+1} \varrho_{nn+1} \quad (80)$$

At this point it is useful to define a column vector with the components $\varrho_n^{(k)} \equiv \varrho_{nn+k}$. This way, we simply arrange the elements of the k th diagonal in the form of a vector. For example, elements on the first side diagonal correspond to $k = 1$, and so on. Let us note that the equation of motion for the off-diagonal matrix elements of the density operator can now be written in a simple matrix form

$$\dot{\varrho}_n^{(k)} = A_{nn'} \varrho_{n'}^{(k)} \quad (81)$$

where summation is implied over repeated indexes and the matrix elements $A_{nn'}$ can be read out from Eq. (49). They are given by

$$A_{nn'} = - \left[\frac{\mathcal{N}'_{nn+k} \mathcal{A}}{1 + \mathcal{N}_{nn+k} \mathcal{B}/\mathcal{A}} + \kappa \left(n + \frac{k}{2} \right) \right] \delta_{nn'} + \frac{\sqrt{n(n+k)} \mathcal{A}}{1 + \mathcal{N}_{n-1+n+k-1} \mathcal{B}/\mathcal{A}} \delta_{nn'+1} + \kappa \sqrt{(n+1)(n+k+1)} \delta_{nn'-1} \quad (82)$$

Clearly, A_{nn} is a tridiagonal matrix. Due to its linearity, we can look for the solution of Eq. (81) by the simple exponential Ansatz, $\boldsymbol{\rho}_n^{(k)}(t) = e^{-\lambda t} \boldsymbol{\rho}_n^{(k)}(0)$. With this substitution, Eq. (81) can be written in the form of an eigenvalue equation to determine λ ,

$$\lambda \boldsymbol{\rho}_n^{(k)} = - \sum_{n-1}^{n+1} A_{nn} \boldsymbol{\rho}_n^{(k)} \quad (83)$$

We restrict the following treatment to $k = 1$ because that is what we need for the calculation of the $g^{(1)}$ correlation function. Higher-order correlation functions, $g^{(k)}$ with $k > 1$, are related to $\boldsymbol{\rho}_n^{(k)}$ with $k > 1$. From the structure of the $-A$ matrix, one can show that all eigenvalues are positive. There is a smallest eigenvalue, which we denote by D . This eigenvalue will dominate the longtime behavior of the field amplitude, as can easily be seen from the following considerations. Let us denote the set of eigenvalues by $\{D, \lambda_j\}$ with $j = 1, 2, 3, \dots$. Then $\boldsymbol{\rho}_n^{(1)}(t)$ can be written as

$$\boldsymbol{\rho}_n^{(1)}(t) = \boldsymbol{\rho}_{n0}^{(1)}(0) \exp(Dt) + \sum_{j=1}^{\infty} \boldsymbol{\rho}_{nj}^{(1)}(0) \exp(-\lambda_j t) \quad (84)$$

From this we see that, indeed, the longtime behavior of the off-diagonal elements will be governed by the first term, since the other terms decay faster according to our assumption of D being the smallest positive eigenvalue. Therefore, our task is reduced to the determination of D . In order to obtain an analytical insight, we can proceed as follows. First, let us note that in the longtime limit all elements of the vector $\boldsymbol{\rho}_n^{(1)}(t)$ decay the same way—they are proportional to $\exp(-Dt)$. Therefore, the sum of the elements also decays with the same rate, D , in this limit. It is quite easy to obtain an equation of motion for the sum of the elements. Starting from Eq. (81) and using Eq. (82) for the case $k = 1$, we immediately obtain

$$\dot{\boldsymbol{\rho}}^{(1)} = - \sum_{n=0}^{\infty} \left[\frac{n+3/2 - \sqrt{(n+1)(n+2)}}{1 + (n+3/2)\mathcal{B}/\mathcal{A}} \mathcal{A} + \kappa(n+1/2 - \sqrt{n(n+1)}) \right] \boldsymbol{\rho}_n^{(1)} \quad (85)$$

Here we introduced the notation $\sum_{n=0}^{\infty} \boldsymbol{\rho}_n^{(1)} = \boldsymbol{\rho}^{(1)}$ and used the fact that $\mathcal{N}'_{nn+1} = n+3/2 + \mathcal{B}/(8\mathcal{A}) \approx n+3/2$ and $\mathcal{N}_{nn+1} = n+3/2 + \mathcal{B}/(16\mathcal{A}) \approx n+3/2$ since $\mathcal{B}/\mathcal{A} \approx 10^{-6}$ and can therefore safely be neglected next to $3/2$. In the longtime limit the time derivative on the left-hand side can simply be replaced by $-D$ due to Eq. (84). It is also plausible to assume that in the same limit those values of n will contribute the most that lie in the vicinity of \bar{n} . Then we can expand the coefficients around the steady-state value of the photon number. This is certainly a good approximation in some region above threshold. The key point is that after the expansion the coefficients of $\boldsymbol{\rho}_n^{(1)}$ on the right-hand side become independent of the summation index n and can be factored out from the sum. Then, after the summation, the quantity $\boldsymbol{\rho}^{(1)}$ also appears on the right-hand side of the equation:

$$D \boldsymbol{\rho}^{(1)} = \frac{\mathcal{A} + \kappa}{8\bar{n}} \boldsymbol{\rho}^{(1)} \quad (86)$$

From this we can simply read out the decay rate:

$$D = \frac{\mathcal{A} + \kappa}{8\bar{n}} \quad (87)$$

This quantity, called the *phase diffusion coefficient*, plays a crucial role in determining the transient behavior of the laser as well as its spectral properties. It exhibits the characteristic line narrowing for high intensity, first found by Schawlow and Townes.⁸⁰ The mean amplitude, Eq. (80), can now be written as

$$\langle a^\dagger(t) \rangle = e^{-Dt} \langle a^\dagger(0) \rangle \quad (88)$$

The decay of any initial coherent component of the laser field is governed by the phase diffusion constant, due to the randomization of the initial phase information. The randomization is due to two separate processes, as can be read out from the analytical expression, Eq. (87) of the phase diffusion constant. The part proportional to the spontaneous emission rate \mathcal{A} is due to the random addition of photons to the field via spontaneous emission; the part proportional to the cavity decay rate is due to leaking of vacuum fluctuations into the cavity through the output mirrors. Both processes randomize the phase of the initial field; as a result, the phase performs a random walk with a diffusion rate given by Eq. (87). Ultimately, of course, vacuum fluctuations are also responsible for spontaneous emission.

The phase diffusion constant also determines the linewidth of the spectrum of the laser field. Using the quantum regression theorem, we immediately find that the (nonnormalized) steady-state field correlation function is given by

$$g^{(1)}(t_0+t, t_0) = \langle a^\dagger(t_0+t)a(t_0) \rangle = \bar{n} \exp(i\omega_0 t - Dt) \quad (89)$$

where ω_0 denotes the operating frequency of the laser, as before. The power spectrum is given by the Fourier transform of the field correlation function:

$$S(\omega) = \frac{1}{\pi} \operatorname{Re} \int_0^\infty g^{(1)}(t_0+t, t_0) e^{-i\omega t} dt = \frac{\bar{n}}{\pi} \frac{D}{(\omega - \omega_0)^2 + D^2} \quad (90)$$

This is a lorentzian spectrum centered around the operating frequency, $\omega = \omega_0$. The full width at half-maximum (FWHM) is given by $2D$. Figure 8 depicts the normalized spectrum $S(\omega)/S(\omega_0)$ versus the detuning $\Delta = (\omega - \omega_0)/D$.

It should be emphasized that our method of obtaining the preceding analytical approximations is justified only if the mean photon number is large, the photon-number distribution consists of a single large peak, and cross-coupling between intensity and phase, arising from the nonlinearity of the gain very far above threshold, is negligible. These conditions are met for a laser in some region above threshold. Near the threshold, however, the intensity fluctuations cannot be neglected. From

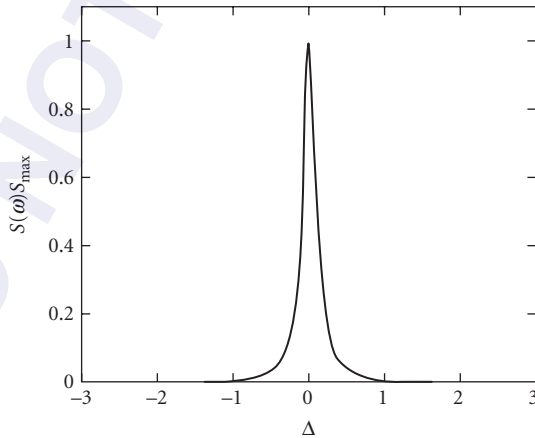


FIGURE 8 Laser spectrum $S(\omega)/S_{\max}$ versus detuning $\Delta = (\omega - \omega_0)/\kappa$, for a laser 100 percent above threshold (note that detuning is in units of bare-cavity linewidth). For our parameters the spectral width D is about 1 percent of the bare-cavity linewidth, an example of the Schawlow-Townes narrowing

numerical studies⁹¹ it was concluded that their contribution to the linewidth is approximately equal to the phase diffusion constant, and the linewidth in this region is about twice what is predicted by Eq. (90); i.e., it is $4D$; and far above threshold (for $\mathcal{A}\kappa > 2$), the linewidth is smaller than the prediction of Eq. (90). These numerical findings are confirmed by a recent analytical approach to the problem.⁹² One of the earliest quantum calculations of the laser linewidth⁶⁹ was based on the quantum theory of noise sources.⁷⁰ The laser linewidth above threshold was established to be due mostly to phase noise,⁶⁹ and it was explained that the factor of 2, too large, was obtained by many previous authors because they had permitted amplitude fluctuations, which were valid near threshold but not far above. These ideas were confirmed by analytic calculations below and above threshold⁹³ and by a numerical solution of the associated Fokker-Planck equation.^{94,95} Moreover, the effects of laser detuning on the linewidth were determined without assuming that the light field decays much slower than the atomic decay rates. It was found that the effective linewidth was a harmonic mean between the field and electronic decay rates, as shown in Eq. (35) of Ref. 69.

The calculations of the phase linewidth done in this section are equivalent to the quasilinear approximation employed in the Langevin noise source procedure in Ref. 93 and shown in Ref. 94 to be valid for dimensionless pump rates $p > 10$. Below threshold, $p < -10$, the components of the electromagnetic field can be treated as uncoupled gaussian variables, leading to the Schawlow-Townes formula. Only the region in the immediate vicinity of threshold requires careful analysis. In that region, it is shown in footnote 10 in Ref. 93 that the coefficients in the Fokker-Planck equation can be expanded in powers of n , retaining only the first nonlinear term. The result then reduces to the rotating-wave Van der Pol oscillator. One advantage of this reduction is that the equation can be scaled in time and in amplitude, leaving the dimensionless pump parameter p as the only remaining parameter. This greatly reduces any subsequent numerical calculations. The ability to retain only one nonlinear term is based only on the requirement that we are dealing with a good oscillator—that is, one whose signal-to-noise ratio is large. Equivalent approximations can be introduced in the density matrix treatment, as well.

The assumption made here that the phase linewidth comes predominantly from the smallest decay eigenvalue was established in Ref. 94, where the actual line shape is shown to be a sum of lorentzians, and the percentage from each lorentzian is calculated. The intensity fluctuations are also expressed as a sum of lorentzians but several modes contribute, not just the lowest. The percentage in each is given in Ref. 94. Part of the reason for this is that the modes approach degeneracy in pairs as one moves above threshold. The lifting of such a degeneracy was shown in Ref. 96 to be associated with a phase transition. Since our system is finite, the phase change occurs gradually and can be observed. The view of lasing as a phase transition will be explored in more detail in Sec. 23.5.

Spectral Properties of the Micromaser Field If we take the $n, n+1$ matrix elements of Eq. (57), then, following the methods of the previous subsection, it is straightforward to derive the diffusion constant. In particular, elements of the first side diagonal can again be arranged in the form of a column vector, and this vector again satisfies an equation of motion similar to Eq. (81), with an appropriate redefinition of the tridiagonal matrix appearing in the equation:

$$\begin{aligned}
 A_{nm'} = & -\{r[1 - \cos(g\tau\sqrt{n+1})\cos(g\tau\sqrt{n+2})] + \kappa(2n_{\text{th}}(n+1) + n+1/2)\}\delta_{nm'} \\
 & + (r\sin(g\tau\sqrt{n})\sin(g\tau\sqrt{N+1}) + \kappa n_{\text{th}}\sqrt{n(n+1)})\delta_{nm'+1} \\
 & + \kappa(n_{\text{th}}+1)\sqrt{(n+1)(n+2)}\delta_{nm'-1}
 \end{aligned} \tag{91}$$

The exponential Ansatz for the decay of the column vector again turns the equation of motion into an eigenvalue equation for the matrix A and the longtime behavior will again be governed by the smallest eigenvalue, which we denote by D . Summing over all elements of the resulting eigenvalue

equation and replacing the coefficients on the right-hand side by their longtime value—that is, expanding the coefficients around $n = \bar{n}$ —will finally yield

$$D = 2r \sin^2 \left(\frac{g\tau}{4\sqrt{\bar{n}}} \right) + \kappa \frac{(2n_{\text{th}} + 1)}{8\bar{n}} \quad (92)$$

for the phase diffusion constant of the micromaser. Here \bar{n} is the mean photon number of the (single-peaked) distribution. This expression was first derived by Scully et al.⁹⁷ It leads to a lorentzian spectrum similar to Eq. (90) but with D appropriately replaced by that of the micromaser. It should be noted, however, that the analytical formula has a more restricted validity than in the case of the laser. Namely, in the case of the micromaser, the photon-number distribution can be multi-peaked, and the simple expansion around the \bar{n} value corresponding to a single dominating peak may not hold. For this more general case several numerical approaches have been developed (see Ref. 98). When Eq. (92) is valid it coincides with the results of the numerical calculations.

For small values of the argument in the sine function, Eq. (92) can be cast to a form which is formally identical to the usual laser phase diffusion constant⁹⁹[cf. Eq. (87)]:

$$\begin{aligned} D &= r \frac{g^2 \tau^2}{8\bar{n}} + \kappa \frac{2n_{\text{th}} + 1}{8\bar{n}} \\ &= \frac{\mathcal{A} + \kappa(2n_{\text{th}} + 1)}{8\bar{n}} \end{aligned} \quad (93)$$

Here we introduced the small signal gain $\mathcal{A} \equiv r g^2 \tau^2$, in analogy to the laser gain. However, for other values of the argument the first term in Eq. (92) may dominate, and the phase diffusion constant can exceed the bare-cavity linewidth, which is a unique quantum feature of the micromaser and makes it distinctly different from the classical Schawlow-Townes-type behavior.

Without going into the specifics, we just mention a few other aspects of the quantum theory of the micromaser. As we have just seen a steady state is reached due to the equilibrium between the gain and loss processes. In some cases, however, a steady state can be reached even in a lossless cavity. This happens if the probability flow in Fig. 5 due to the gain process is interrupted for some value of the interaction parameter. The upward flow is interrupted when

$$g\tau\sqrt{n_q + 1} = q\pi \quad (94)$$

and the downward flow is interrupted when

$$g\tau\sqrt{n_q} = q\pi \quad (95)$$

with $q = 1, 2, \dots$ in both cases. They are called the *upward* and *downward trapping state*, respectively.¹⁰⁰ In such cases the state of the field is a number state. The signature of the number state is a large maximum in the linewidth at the corresponding interaction parameter. It can easily be understood qualitatively: since the state of the field is a number state, it cannot have any phase information. Therefore, the phase randomizes on a very rapid time scale; in other words, the phase correlations decay very rapidly, and a large phase diffusion constant ensues. Further, if the atoms are injected in a coherent superposition of their upper and lower levels in the cavity then, under certain conditions, the so-called tangent and cotangent states of the field may develop.¹⁰¹ Finally, it should be mentioned that the master equations Eqs. (49) and (56) hold only for the case when the time interval between consecutive atomic injections is completely random. Other arrival times statistics, including the case of regular injections, have been investigated by a number of authors.¹⁰² A closely related area of recent theoretical studies pertains to the detection of the statistical properties of the (experimentally inaccessible) intracavity field via monitoring the statistics of the outgoing atoms for poissonian¹⁰³ as well as nonpoissonian¹⁰⁴ pumping. For regular (subpoissonian) pumping, transient oscillations in field correlation function and a corresponding multi-peaked spectrum were predicted.¹⁰⁵

Most of the predictions of this theory have been confirmed by experiments. For example, trapping states have been observed in recent experiments by the Garching group.¹⁰⁶ Without providing an exhaustive list, we just refer the reader to recent progress in the experimental department.¹⁰⁷

Other Approaches

So far we have considered laser theory based on a density-operator approach. An equivalent approach to a laser theory can be formulated using a Heisenberg-Langevin approach. In this approach explicit equations of motion are derived for the field operator.

The quantum-noise operator formalism was presented in essentially its final form by Lax at the 1965 meeting on the *Physics of Quantum Electronics*.⁶⁹ There, it was applied to the laser, calculating, the laser linewidth in its most general form. For the general theory, see the 1966 Brandeis lecture notes of Lax.⁷¹

In the present formulation, our goals are more specific—namely, the Heisenberg-picture quantum theory of the laser. To that end, we will here give a quantum-noise treatment along the lines of that followed in the previous subsection—namely, a single laser mode damped at a rate κ by a (dissipative) reservoir and driven by atoms injected into the laser cavity at random times t_j .

First we discuss the simple example of damping of the field by a reservoir and derive the quantum Langevin equation for the field operator. We then discuss the gain noise in a laser and derive the laser linewidth.

Damping of Field by Reservoir We consider the damping of a single-mode field of frequency ν interacting with a reservoir consisting of simple harmonic oscillators. The system describes, for example, the damping of a single-mode field inside a cavity with lossy mirrors. The reservoir in this case consists of a large number of phononlike modes in the mirror.

The field is described by the creation and destruction operators a^\dagger and a , whereas the harmonic oscillators of frequency $\nu_k = ck$ are described by the operators b_k^\dagger and b_k . The field-reservoir system evolves in time under the influence of the total hamiltonian:

$$\mathcal{H} = \hbar\omega_0\left(a^\dagger a + \frac{1}{2}\right) + \hbar\sum_k \nu_k \left(b_k^\dagger b_k + \frac{1}{2}\right) + \hbar\sum_k g_k (ab_k^\dagger + b_k a^\dagger) \quad (96)$$

Here g_k are the coupling constants and we have made the rotating-wave approximation. We are interested in the evolution of the field operator a . The Heisenberg equations of motion for the field and reservoir are

$$\dot{a}(t) = -i\omega_0 a(t) - i\sum_k g_k b_k(t) \quad (97)$$

$$\dot{b}_k(t) = -i\nu_k b_k(t) - ig_k a(t) \quad (98)$$

The equation for b_k can be formally solved, and the resulting expression is substituted in Eq. (97). In the Weisskopf-Wigner approximation, the annihilation operator in the interaction picture $a = a(t) \exp[i\omega_0(t-t_0)]$ satisfies a Langevin equation

$$\dot{a} = -\frac{\kappa}{2}a + F(t) \quad (99)$$

where

$$F(t) = -i\sum_k g_k b_k(0) \exp[-i(\nu_k - \omega_0)(t-t_0)] \quad (100)$$

is a noise operator. Equation (99) clearly indicates that the damping of the field (represented by the term $-\kappa a/2$) is accompanied by noise.

For the damping of the single-mode field inside a cavity via transmission losses, the damping constant κ is related to the quality factor Q of the cavity via $\kappa = \omega_0/Q$.

Atomic (Gain) Noise and Laser Linewidth As discussed earlier, the natural linewidth of the laser arises due to spontaneous emission by the atoms. In the density-operator approach, a fully nonlinear treatment was followed. Here, we present a simple linear analysis to calculate the laser linewidth in the Heisenberg-Langevin approach. We assume that the atoms are long lived, and that they interact with the cavity field for a time τ . This treatment allows us to include the memory effects inside a laser, and is one of the simplest examples of a nonmarkovian process.^{108,109}

We start with the hamiltonian describing the atom-field interaction:

$$\mathcal{H} = \mathcal{H}_F + \mathcal{H}_{\text{atom}} + \hbar g \sum_i \{ \sigma^i a^\dagger N(t_i, t, \tau) + Hc \} \quad (101)$$

where \mathcal{H}_F and $\mathcal{H}_{\text{atom}}$ describe the field and atoms, respectively; g is the atom-field coupling constant; and σ_i is the lowering operator for the i th atom; Hc is the injunction to add the Hermitian conjugate. The operators a and a^\dagger represent the annihilation and creation operators, and $N(t_i, t, \tau)$ is a notch function which has the value

$$N(t_i, t, \tau) = \begin{cases} 1 & \text{for } t_i \leq t < \tau \\ 0 & \text{otherwise} \end{cases} \quad (102)$$

Using this hamiltonian, we write the equations for the atom-field operators in the interaction picture as

$$\dot{a} = -ig \sum_i \sigma^i N(t_i, t, \tau) - \frac{1}{2} \kappa a(t) + F_\kappa(t) \quad (103)$$

$$\dot{\sigma}^i = ig N(t_i, t, \tau) \sigma_z^i a(t)$$

where the effects of cavity damping are determined by the cavity decay rate κ and the associated Langevin noise source F_κ . Integrating the equation for the atom operator and substituting it into that for the field operator, we obtain

$$\dot{a}(t) = \int_{-\infty}^t dt' \alpha(t, t') a(t') - \frac{1}{2} \gamma a + F_\alpha(t) + F_\kappa(t) \quad (104)$$

where

$$\alpha(t, t') = g^2 \sum_i N(t_i, t, \tau) N(t_i, t', \tau) \sigma_z^i(t') \quad (105)$$

$$F_\alpha(t) = -ig \sum_i N(t_i, t, \tau) \sigma^i(t_i) \quad (106)$$

Here, the noise operator Eq. (106) may be seen to have the moments

$$\langle F_\alpha(t) \rangle = 0 \quad (107)$$

$$\langle F_\alpha^\dagger(t) F_\alpha(t') \rangle = g^2 \sum_{ij} N(t_i, t, \tau) N(t_j, t', \tau) \langle \sigma^{\dagger i}(t_i) \sigma^j(t_j) \rangle \quad (108)$$

Because we are injecting our lasing atoms in the upper state, the atomic average is given by $\langle \sigma^{\dagger i}(t_i) \sigma^j(t_j) \rangle = \delta_{ij}$. After replacing the sum upon i in Eq. (108) by an integration over injection times t_j , we find

$$\langle F_\alpha^\dagger(t) F_\alpha(t') \rangle = r g^2 \{ N(t' - \tau, t, \tau) [t - (t' - \tau)] - N(t', t, \tau) [t - (t' + \tau)] \} \quad (109)$$

where r is the atomic injection rate. The phase variance can then be calculated through the noise operator product:

$$\langle \phi^2(t) \rangle = -\frac{1}{2\pi} \int_0^t dt' \int_0^{t'} dt'' \langle F^\dagger(t') F(t'') \exp\{i[\phi(t') - \phi(t'')]\} \rangle \quad (110)$$

On insertion of Eq. (109) into Eq. (110), the expression for the generalized maser phase diffusion noise $\langle \phi^2(t) \rangle$ is found to be

$$\langle \phi^2(t) \rangle = \left(\frac{\mathcal{A}}{2\bar{n}} \right) \left[\left(\frac{t^2}{\tau} - \frac{t^3}{3\tau^2} \right) \theta(\tau - t) + \left(t - \frac{\tau}{3} \right) \theta(t - \tau) \right] \quad (111)$$

Here $\mathcal{A} = rg^2\tau^2$ is the small-signal gain of the maser [cf. Eq. (93), with $n_{th} = 0$ and using that in steady state $\mathcal{A} = \kappa$]. In the case involving atoms which are injected at random times t_i but which decay via spontaneous emission to far-removed ground states at a rate γ , a similar but more complicated analysis can be carried out. The result in this case is given by

$$\langle \phi^2(t) \rangle = \left(\frac{\mathcal{A}}{2\bar{n}} \right) [t + \gamma^{-1}(e^{-\gamma} - 1)] \quad (112)$$

Here $\mathcal{A} = 2rg^2/\gamma^2$ is the small-signal gain of the laser [cf. Eq. (50)]. In both of the preceding cases, we find that for times $t = t_m$ small compared to the atomic lifetime, the phase diffusion is quadratic in the measurement time t_m ; that is, we now have a phase error which goes as

$$\Delta\phi^2 = \left(\frac{\mathcal{A}t_m}{2\bar{n}} \right) \left(\frac{\gamma t_m}{2} \right) \quad (113)$$

Therefore, we see that the quantum noise due to spontaneous emission is reduced from the Schawlow-Townes linewidth $2D = \mathcal{A}/2\bar{n}$ by the factor $\gamma t_m/2$, which can be a significant reduction for short measurement times. For times long compared to the atomic lifetime, however, the Schawlow-Townes result is obtained from both Eqs. (111) and (112) as expected.

23.5 THE LASER PHASE-TRANSITION ANALOGY

Considerations involving the analogies between phase transitions in ferromagnets, superfluids, and superconductors have emphasized the similarities between these systems near their critical temperatures.¹¹⁰

A natural comparison can be made between second-order phase transitions, such as the order-disorder transitions of ferromagnetic and ferroelectric materials or the vapor-liquid transition of a pure fluid, and the laser threshold. As we have discussed in Sec. 23.4, the state of a laser changes abruptly upon passing through the threshold point. This point is characterized by a threshold population inversion.

The physical basis for this similarity becomes evident when it is recalled that the usual treatments of laser behavior are self-consistent theories. In the laser analysis we assume that each atom evolves in a radiation field due to all the other atoms, and then calculate the field produced by many such evolving atoms in a self-consistent fashion. In this way the laser problem is similar to that of a ferromagnet, in which each spin sees a mean magnetic field due to all the other spins and aligns itself accordingly, thus contributing to the average magnetic field.

Following this point of view, we can discuss the laser theory using the language of second-order phase transitions.

The density matrix of the laser field obeys Eq. (49). The time dependence of the expectation value \bar{E} of the electric field operator $E=(a+a^\dagger)$ is there given by the following equation:

$$\dot{\bar{E}} = \frac{1}{2}(\mathcal{A} - \kappa)\bar{E} - \frac{\mathcal{B}}{2}\bar{E}^3 \quad (114)$$

Here we have assumed that the laser is operating close to threshold ($\mathcal{B}\bar{n}/\mathcal{A} \ll 1$) so that we retain only the terms proportional to \mathcal{B} . In addition, we assume $\bar{E} \gg 1$. We can then replace \bar{E}^3 by \bar{E}^3 and Eq. (114) becomes the well-known result of Lamb's semiclassical theory. The steady-state properties of the laser oscillator are described by the following equation of state:

$$(\mathcal{A} - \kappa)\bar{E} - \mathcal{B}\bar{E}^3 = 0 \quad (115)$$

The threshold condition is given by $\mathcal{A} = \kappa$ as before. Upon putting $\mathcal{A} = a\sigma$, $\mathcal{B} = b\sigma$, and $\kappa = a\sigma_t$ where σ_t is the threshold population inversion, the steady-state solution of Eq. (115) is

$$\begin{aligned} \bar{E} &= 0 && \text{if } \sigma - \sigma_t < 0 \text{ (below threshold)} \\ \bar{E} &= \left[\frac{a}{b} \left(\frac{\sigma - \sigma_t}{\sigma} \right)^{1/2} \right] && \text{if } \sigma - \sigma_t > 0 \text{ (above threshold)} \end{aligned} \quad (116)$$

Equation (116) is formally identical to the equation for a ferromagnet in the Weiss mean-field theory. The electric field E corresponds to the static magnetization M , which is the order parameter in the ferromagnetic transition. The quadratic polarization $P = (\mathcal{A}\bar{E} - \mathcal{B}\bar{E}^3)/2$ in Eq. (115) corresponds to the magnetic field H generated by a magnetization M , and the term $\kappa\bar{E}/2$ corresponds to a local magnetic field which is assumed proportional to M in the mean-field theory. Furthermore, the steady-state points depend $\sigma - \sigma_t$ in the same way that M in the ferromagnetic case depends on $T - T_c$, where T_c is the critical temperature. Therefore, σ and σ_t correspond to T and T_c , respectively. The similarity between these two systems is summarized in Table 1 and illustrated in Fig. 9.

We recall that the probability density $P(M)$ for a ferromagnetic system with magnetization M near a phase transition is given by, in thermal equilibrium,

$$P(M) = N'' \exp\left(-\frac{F(M)}{k_b T}\right) \quad (117)$$

where

$$F(M) = \frac{1}{2}c(T - T_c)M^2 + \frac{1}{4}dT M^4 \quad (118)$$

is the free energy. In the corresponding laser analysis, the probability density for the electromagnetic field $P(E)$ is derived in the form

$$P(E) = N' \exp\left(-\frac{G(E)}{k_b \sigma}\right) \quad (119)$$

For this purpose we transform the laser equation for the density matrix for the field [Eq. (49)] into an equivalent equation in terms of the $P(\alpha, \alpha^*)$ representation defined by

$$\rho = \int d^2\alpha P(\alpha, \alpha^*) |\alpha\rangle\langle\alpha| \quad (120)$$

where $|\alpha\rangle$ is an eigenstate of the annihilation operator a with eigenvalue α . The P representation allows us to evaluate any normally ordered correlation function of the field operators using the

TABLE 1 Summary of Comparison between the Laser and a Ferromagnetic System Treated in a Mean-Field Approximation

Parameter	Ferromagnet	Laser
Order parameter	Magnetization M	Electric field strength E
Reservoir variable	Temperature T	Population inversion σ Threshold inversion σ_t
Coexistence curve*	$M = \Theta(T_c - T) \left[\frac{c}{d} \frac{T - T_c}{T} \right]^{1/2}$	$E = \Theta(\sigma - \sigma_t) \left[\frac{a}{b} \frac{\sigma - \sigma_t}{\sigma} \right]^{1/2}$
Symmetry breaking mechanism	External field H	Injected signal S
Critical isotherm [†]	$M = \left[\frac{H}{dT_c} \right]^{1/2}$	$E = \left[\frac{2S}{b\sigma_t} \right]^{1/2}$
Zero field susceptibility*	$X \equiv \left(\frac{\partial M}{\partial H} \right) \Big _{H=0}$ $= \Theta(T_c - T) [2c(T_c - T)]^{-1}$ $+ \Theta(T - T_c) [c(T - T_c)]^{-1}$	$\xi \equiv \left(\frac{\partial E}{\partial S} \right) \Big _{S=0}$ $= \Theta(\sigma_t - \sigma) \left[\frac{a(\sigma_t - \sigma)}{2} \right]^{-1}$ $+ \Theta(\sigma - \sigma_t) [a(\sigma - \sigma_t)]^{-1}$
Thermo-dynamic potential	$F(M) = \frac{1}{2}c(T - T_c)M^2$ $+ \frac{1}{4}dTM^4$ $- HM + F_0$	$G(E) = -\frac{a}{4}(\sigma - \sigma_t)E^2$ $+ \frac{1}{8}b\sigma E^4$ $- SE + G_0$
Statistical distribution	$P(M) = N'' \exp\left(-\frac{F(M)}{k_B T}\right)$	$P(E) = N' \exp\left(-\frac{G(E)}{k_B \sigma}\right)$

* $\Theta()$ is Heaviside's unit step function.

[†]Value of order parameter at critical point.

methods of classical statistical mechanics. The quantity $P(\alpha, \alpha^*)$ represents the probability density for finding the electric field corresponding to α .

Near threshold, $P(\alpha, \alpha^*)$ obeys the following Fokker-Planck equation:

$$\begin{aligned} \frac{\partial P}{\partial t} = & -\frac{\partial}{\partial \alpha} \left[\frac{1}{2}(\mathcal{A} - \kappa)\alpha P - \frac{1}{2}\mathcal{B}|\alpha|^2 \alpha P \right] \\ & - \frac{\partial}{\partial \alpha^*} \left[\frac{1}{2}(\mathcal{A} - \kappa)\alpha^* P - \frac{1}{2}\mathcal{B}|\alpha|^2 \alpha^* P \right] + \mathcal{A} \frac{\partial^2 P}{\partial \alpha \partial \alpha^*} \end{aligned} \quad (121)$$

The steady-state solution of this equation is given by

$$P(\alpha, \alpha^*) = \mathcal{N} \exp \left[\frac{(\mathcal{A} - k)|\alpha|^2 - \mathcal{B}|\alpha|^4/2}{2\mathcal{A}} \right] \quad (122)$$

where \mathcal{N} is a normalization constant. The P representation can be rewritten in terms of the variables $x = \text{Re } \alpha$ and $y = \text{Im } \alpha$ as

$$P(x, y) = \mathcal{N} \exp \left[-\frac{G(x, y)}{K\sigma} \right] \quad (123)$$

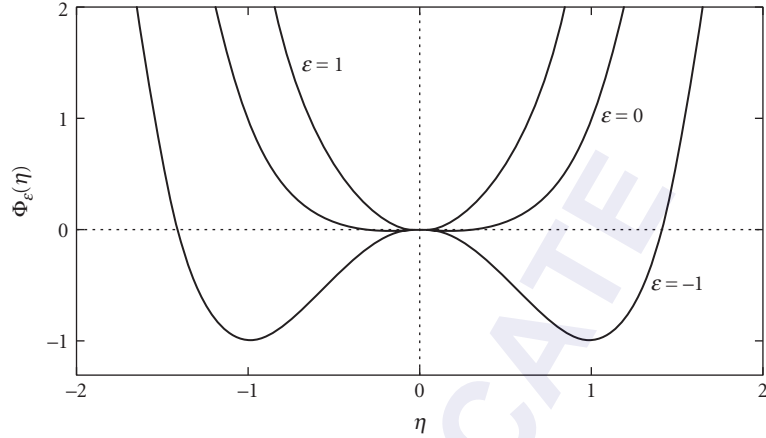


FIGURE 9 Scaled thermodynamical potentials. The $H = 0$ version of $F(M)$ and the $S = 0$ version of $G(E)$ can be expressed in terms of $\Phi_\varepsilon(\eta) = 2\varepsilon\eta^2 + \eta^4$ because

$$\Phi_\varepsilon(\eta) = \begin{cases} \frac{[F(M) - F_0]}{\frac{c^2 T_c}{4d}} & \text{with } M = (c/d)^{1/2} (T/T_c)^{1/4} \eta \\ & \text{and } \varepsilon = (T - T_c)/(TT_c)^{1/2} \\ \frac{[G(E) - G_0]}{\frac{a^2 \sigma_i}{8b}} & \text{with } E = (a/b)^{1/2} (\sigma_i/\sigma)^{1/4} \eta \\ & \text{and } \varepsilon = (\sigma_i - \sigma)/(\sigma\sigma_i)^{1/2} \end{cases}$$

are equivalent to the respective entries in Table 1. The plot shows $\Phi_\varepsilon(\eta)$ for $\varepsilon=1$ (ferromagnet above, T_c , laser below threshold), $\varepsilon=0$ (ferromagnet at T_c , laser at threshold), and $\varepsilon=-1$ (ferromagnet below T_c , laser above threshold).

with

$$G(x, y) = -\frac{1}{4}a(\sigma - \sigma_i)(x^2 + y^2) + \frac{1}{8}b\sigma(x^2 + y^2)^2 \quad (124)$$

Here $K = a/4$ is one-fourth of the gain of one atom, $a(\sigma - \sigma_i) = \mathcal{A} - \kappa$, and $b\sigma = \mathcal{B}$.

We can see that the steady-state situation of the laser corresponds to the minimum value of G , i.e., $\partial G/\partial x = \partial G/\partial y = 0$. These solutions are $x = y = 0$ and $|\alpha|^2 = (x^2 + y^2) = a(\sigma - \sigma_i)/b\sigma$. Thus, for $(\sigma - \sigma_i) < 0$, the only allowed solution is $x = y = 0$. However, for $(\sigma - \sigma_i) > 0$, $x = y = 0$ is an unstable solution as the second derivative of G with respect to x and y is positive. This is seen clearly in Fig. 9, where we have plotted G versus $x = E$ for $y = 0$.

We thus see that G behaves in essentially the same way as the free energy of a thermodynamic system.

It should be emphasized that in the thermodynamic treatment of the ferromagnetic order-disorder transition, there are three variables required: (1) magnetization M , (2) external magnetic field H , and (3) temperature T . In order to have a complete analogy, it is important to realize that in addition to the electric-field-magnetization, population inversion-temperature correspondences, there must exist a further correspondence between the external magnetic field and a corresponding symmetry-breaking mechanism in the laser analysis. As shown in Ref. 111 and illustrated in Fig. 10, this symmetry breaking mechanism in the laser problem corresponds to an injected classical signal S . This leads to a skewed effective free energy.

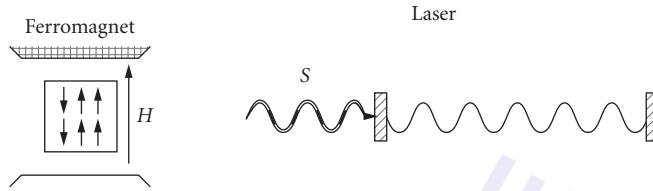


FIGURE 10 Figure depicting the broken symmetry mode of operation for both a ferromagnet and a laser.

An example of how the analogy can provide us with deeper insight is contained in the fact that we are able to *guess* correctly the $P(E)$ for a laser influenced by an injected signal, by analogy with the corresponding magnetic problem in the broken symmetry mode of operation.

More recently we have been turning the tables and using the quantum laser theory to learn about Bose-Einstein condensation (BEC). Recent experiments on BEC in laser-cooled gases,^{112–114} and in He^4 in porous gel,¹¹⁵ have stimulated a wealth of theoretical work^{116–118} on the equilibrium and non-equilibrium properties of confined quantum degenerate gases.^{119–124} Presently the partition function, critical temperature, and other such quantities are of interest for N bosons in a box below T_c . But the canonical ensemble is difficult to use in practice, because the state sums must be restricted to N particles. Indeed, the canonical partition function for a Bose gas of N particles at temperature T has not been so widely studied as one might have thought. To quote Herzog and Olshanii,

To our knowledge there is no simple analytic expression for the canonical partition function in [the case of N bosons in a three-dimensional trap].¹²¹

Furthermore, there are questions of principle concerning the critical temperature and the validity of using phase-transition concepts in a mesoscopic sample having a small number of particles ($N \approx 10^3$). In fact, Uhlenbeck pointed out to Einstein many years ago that BEC rigorously occurs only in the limit of infinite particle number.¹²⁵ Indeed, for a finite number of atoms there is no sharp “critical point” or critical temperature T_c . But the same can be said for the laser threshold. There is a *gradual* transition from disorder to order in both cases. However, as discussed later, even when fluctuations are present, T_c for a Bose gas and the laser threshold inversion are well defined.

Motivated by the preceding, we extend the laser-phase transition analogy to include BEC. We present a new approach to the problem of N bosons in thermal equilibrium below T_c . We emphasize that the present work provides another example¹²⁶ in which steady-state (detailed balance) solutions to nonequilibrium equations of motion provide a supplementary approach to conventional statistical mechanics (e.g., partition-function calculations). The present approach lends itself to different approximations; yielding, among other things, a simple (approximate) analytic expression for the ground-state density matrix for N trapped bosons and the partition function for same.

Thus, we seek a nonequilibrium equation of motion for the ground state of an ideal Bose gas in a three-dimensional harmonic trap coupled to the thermal reservoir, as shown elsewhere.¹²⁷

$$\begin{aligned} \dot{\rho}_{n_0, n_0} = & -K_{n_0} (n_0 + 1) \rho_{n_0, n_0} + K_{n_0-1} n_0 \rho_{n_0-1, n_0-1} \\ & - H_{n_0} n_0 \rho_{n_0, n_0} + H_{n_0+1} (n_0 + 1) \rho_{n_0+1, n_0+1} \end{aligned} \quad (125)$$

The cooling and heating coefficients K_{n_0} and H_{n_0} are given by

$$K_{n_0} = \sum_k 2\pi W_k g_k^2 \langle \eta_k + 1 \rangle \langle n_k \rangle_{n_0} \quad (126)$$

and

$$H_{n_0} = \sum_k 2\pi W_k g_k^2 \langle \eta_k \rangle \langle n_k + 1 \rangle_{n_0} \quad (127)$$

where W_k is the heat-bath density of states, $\langle \eta_k \rangle$ is the average occupation number of the k th heat-bath oscillator, and $\langle \eta_k \rangle_{n_0}$ is the average number of atoms in the k th excited state, given n_0 atoms in the condensate. Here the coefficient K_{n_0} denotes the cooling rate from the excited states to the ground state, and similarly H_{n_0} stands for the heating rate for the ground state.

The heating term is approximately

$$H_{n_0} = \kappa \sum_k \langle \eta(\varepsilon_k) \rangle = \kappa \sum_{\ell, m, n} \left[\exp\left(\frac{\hbar\Omega}{k_B T}\right) (\ell + n + m) - 1 \right]^{-1} \quad (128)$$

In the weak trap limit, this yields

$$H_{n_0} = \kappa \left(\frac{k_B T}{\hbar\Omega} \right)^3 \zeta(3) \quad (129)$$

where $\zeta(3)$ is the Riemann zeta function and Ω is the trap frequency. Likewise, the cooling term in Eq. (125) is governed by the total number of excited state bosons,

$$K_{n_0} = \kappa \sum_k \langle n_k \rangle_{n_0} = \kappa(N - n_0) \quad (130)$$

By writing the equation of motion for $\langle n_0 \rangle$ from Eq. (125), using H_{n_0} in the weak trap limit, and Eq. (130) for K_{n_0} , we find

$$\langle \dot{n}_0 \rangle = \kappa \left[(N+1) \langle n_0 \rangle - \langle (n_0+1)^2 \rangle - \zeta(3) \left(\frac{k_B T}{\hbar\Omega} \right)^3 \langle n_0 \rangle \right] + \kappa(N+1) \quad (131)$$

Noting that near T_c , $\langle n_0 \rangle = N$, we may neglect $\langle (n_0+1)^2 \rangle$ compared to $N \langle n_0 \rangle$, and neglecting the spontaneous emission term $\kappa(N+1)$, Eq. (131) becomes

$$\langle \dot{n}_0 \rangle = \kappa \left[N - \zeta(3) \left(\frac{k_B T}{\hbar\Omega} \right)^3 \right] \langle n_0 \rangle \quad (132)$$

We now define the critical temperature (in analogy with the laser threshold) such that cooling (gain) equals heating (loss) and $\langle \dot{n}_0 \rangle = 0$ at $T = T_c$; this yields

$$T_c = \left(\frac{\hbar\Omega}{k_B} \right) \left[\frac{N}{\zeta(3)} \right]^{1/3} \quad (133)$$

Thus, by defining the critical temperature as that temperature at which the rate of removal of atoms from the ground state equals the rate of addition, we arrive at the usual definition for the critical temperature, even for mesoscopic systems.

23.6 EXOTIC MASERS AND LASERS

Lasing without Inversion

For a long time, it was considered that population inversion was necessary for laser action to take place. Recently, it has been shown both theoretically^{128–130} and experimentally^{131–134} that it is also

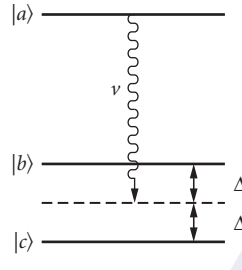


FIGURE 11 Level diagram for lasing without inversion.

possible to achieve lasing without inversion (LWI). In LWI, the essential idea is the cancellation of absorption by atomic coherence and interference.

Consider a system of three-level atoms interacting with a laser field in a cavity. The simple model we will focus on is that of Fig. 11. The atoms have one upper level $|a\rangle$ and two lower levels $|b\rangle$ and $|c\rangle$, with energies $\hbar\omega_a$, $\hbar\omega_b$, and $\hbar\omega_c$, respectively. The cavity field of frequency ν can be detuned from the atomic transition, as shown in the figure. The transitions $|a\rangle \rightarrow |b\rangle$ and $|a\rangle \rightarrow |c\rangle$ are now induced by one classical light field of frequency ν . The transition $|b\rangle \rightarrow |c\rangle$ is dipole forbidden. The atoms are pumped at a rate r_a in a coherent superposition of states

$$\rho(t_i) = \rho_{aa}^{(0)} |a\rangle\langle a| + \rho_{bb}^{(0)} |b\rangle\langle b| + \rho_{cc}^{(0)} |c\rangle\langle c| + \rho_{bc}^{(0)} |b\rangle\langle c| + \rho_{cb}^{(0)} |c\rangle\langle b| \quad (134)$$

Here $\rho_{\alpha\alpha}^{(0)}$ ($\alpha = a, b, c$) are the level populations and $\rho_{\alpha\alpha'}^{(0)}$ ($\alpha \neq \alpha'$) are the atomic coherences. We give a simple argument to show how cancellation of absorption can lead to lasing without inversion in this scheme.

As the levels $|b\rangle$ and $|c\rangle$ are independent, the probability of emission is given by

$$\begin{aligned} P_{\text{emission}} &= P_b + P_c \\ &= (\kappa_{a \rightarrow b} |2\mathcal{E}|^2 + |\kappa_{a \rightarrow c}|^2 |2\mathcal{E}|^2) \rho_{aa}^{(0)} \end{aligned} \quad (135)$$

where $\kappa_{a \rightarrow b}$ and $\kappa_{a \rightarrow c}$ are constants which depend on the matrix element between the relevant levels and the coupling of the atom with the field. On the other hand, the absorption probability is given by

$$\begin{aligned} P_{\text{absorption}} &= \kappa |c_b + c_c|^2 |2\mathcal{E}|^2 \\ &= \kappa (\rho_{bb}^{(0)} + \rho_{cc}^{(0)} + \rho_{bc}^{(0)} + \rho_{cb}^{(0)}) \mathcal{E}^2 \end{aligned} \quad (136)$$

where c_a and c_b are the probability amplitudes for the states $|b\rangle$ and $|c\rangle$. Therefore, the rate of growth of the laser field amplitude, under appropriate conditions, becomes

$$\dot{\mathcal{E}} = \frac{\mathcal{A}}{2} (\rho_{aa}^{(0)} - \rho_{bb}^{(0)} - \rho_{cc}^{(0)} - \rho_{bc}^{(0)} - \rho_{cb}^{(0)}) \mathcal{E} \quad (137)$$

Here \mathcal{A} is a constant. Thus, if the terms $\rho_{bc}^{(0)}$ and $\rho_{cb}^{(0)}$ cancel $\rho_{bb}^{(0)}$ and $\rho_{cc}^{(0)}$, we have

$$\dot{\mathcal{E}} = \frac{\mathcal{A}}{2} \rho_{aa}^{(0)} \mathcal{E} \quad (138)$$

and we can have lasing even if only a small fraction of atoms is in the excited state $|a\rangle$, that is, even if $\rho_{aa} < (\rho_{bb} + \rho_{cc})$.

Physically, the lack of absorption in the three-level system considered here is a manifestation of quantum coherence phenomena. When an atom makes a transition from the upper level to the two lower levels, the total transition probability is the sum of $|a\rangle \rightarrow |b\rangle$ and $|a\rangle \rightarrow |c\rangle$ probabilities. However, the transition probability from the two lower levels to the single upper level is obtained by squaring the sum of the two probability amplitudes. When there is coherence between the two lower levels, this can lead to interference terms yielding a null in the transition probability corresponding to photon absorption.

Correlated (Spontaneous) Emission Laser

As discussed earlier, the fundamental source of noise in a laser is spontaneous emission. A simple pictorial model for the origin of the laser linewidth envisions it as being due to the random phase diffusion process arising from the addition of spontaneously emitted photons with random phases to the laser field. Here we show that the quantum noise leading to the laser linewidth can be suppressed below the standard Schawlow-Townes limit by preparing the atomic systems in a coherent superposition of states as in the Hanle-effect and quantum-beat experiments. In such coherently prepared atoms, the spontaneous emission is said to be *correlated*. Lasers operating via such a phase-coherent atomic ensemble are known as *correlated emission lasers* (CELs).¹³⁵

An interesting aspect of the CEL is that it is possible to eliminate the spontaneous emission quantum noise in the relative linewidths by correlating the two spontaneous emission noise events.

A number of schemes exist in which quantum noise quenching below the standard limit can be achieved. In two-mode schemes a correlation between the spontaneous emission events in two different modes of the radiation field is established via atomic coherence so that the relative phase between them does not diffuse or fluctuate. In a Hanle laser¹³⁶ and a quantum-beat laser¹³⁷ this is achieved by pumping the atoms coherently such that every spontaneously emitting atom contributes equally to the two modes of the radiation, leading to a reduction and even vanishing of the noise in the phase difference. In a two-photon CEL, a cascade transition involving three-level atoms is coupled to only one mode of the radiation field.¹³⁸ A well-defined coherence between the upper and lower levels $|a\rangle$ and $|c\rangle$ leads to a correlation between the light emitted by an $|a\rangle \rightarrow |b\rangle$ and a subsequent $|b\rangle \rightarrow |c\rangle$ transition.

The quantum theory of quantum-beat or Hanle-effect lasers may be conveniently cast in terms of the equation of motion for the density matrix describing the laser radiation field $\rho(a_1, a_1^\dagger; a_2, a_2^\dagger)$; that is,

$$\dot{\rho} = \sum_{ij} \mathcal{L}_{ij} \rho \quad (139)$$

where the linear gain and cross-coupling Liouville operators are given by

$$\mathcal{L}_{ii} \rho = -\frac{1}{2} [\alpha_{ii} \rho a_i a_i^\dagger + \alpha_{ii}^* a_i a_i^\dagger \rho - (\alpha_{ii} + \alpha_{ii}^*) a_i^\dagger \rho a_i] \quad (140)$$

$$\mathcal{L}_{12} \rho = -\frac{1}{2} [\alpha_{12} \rho a_2 a_1^\dagger + \alpha_{21}^* a_2 a_1^\dagger \rho - (\alpha_{12} + \alpha_{21}^*) a_1^\dagger \rho a_2] e^{i\Phi} \quad (141)$$

$$\mathcal{L}_{21} \rho = -\frac{1}{2} [\alpha_{21} \rho a_1 a_2^\dagger + \alpha_i^* a_1 a_2^\dagger \rho - (\alpha_i + \alpha_i^*) a_2^\dagger \rho a_1] e^{i\Phi} \quad (142)$$

Here α_{ij} are constants that depend on the parameters of the gain medium such as detunings, Rabi frequency of the driving field, and so on. When the coherent mixing of levels $|a\rangle$ and $|b\rangle$ is produced via a microwave signal having frequency ω_0 , the phase angle ϕ is given by $\Phi(t) = (\nu_1 - \nu_2 - \omega_0)t - \phi$ where ϕ is the (microwave determined) atomic-phase difference $\phi_a - \phi_b$. In the case of polarization-induced

coherent mixing, the phase angle is $\Phi(t) = (v_2 - v_1)t - \phi$, where ϕ is again the relative phase between levels $|a\rangle$ and $|b\rangle$ but determined this time by the state of elliptic polarization of the pump light used to excite the atoms.

The Liouville equation (31) for the reduced-density operator for the field can be converted into an equivalent Fokker-Planck equation by introducing coherent state representation for a_1 and a_2 and the P representation $P(\alpha, \alpha^*)$ for ρ . If we define the coherent states as

$$a_i |\alpha_i, \alpha_2\rangle = \alpha_i |\alpha_i, \alpha_2\rangle; i = 1, 2 \quad (143)$$

where α is an arbitrary complex number, and we represent α_i as

$$\alpha_i = \rho_i \exp(i\theta_i); i = 1, 2 \quad (144)$$

then the Fokker-Planck equation in terms of ρ_i, θ_i will contain a term which describes diffusion of relative phase $\theta = \theta_1 - \theta_2$ as

$$\dot{P} = \frac{\partial^2}{\partial \theta^2} [\mathcal{D}(0)P] \quad (145)$$

with

$$\mathcal{D} = \frac{1}{16} \left\{ \left(\frac{\alpha_{11}}{\rho_1^2} + \frac{\alpha_{22}}{\rho_2^2} \right) - \frac{(\alpha_{12} + \alpha_{21}^*) e^{-i\psi}}{\rho_1 \rho_2} \right\} \quad (146)$$

and $\psi = \Phi + \theta_1 - \theta_2$ with θ_i being the phase of the i th field. The diffusion constant \mathcal{D} for the relative phase vanish for $\psi = 0, \rho_1 = \rho_2$, and $\alpha_{11} = \alpha_{22} = \alpha_{12} = \alpha_{21}^*$, thus leading to CEL action.

Free-Electron Laser

A coherent emission of radiation in a free-electron laser (FEL) is due to the bunching of a relativistic electron beam propagating along a periodic magnetic structure. The electrons experience a Lorentz force and thus follow oscillating orbits and radiate. This spontaneous emission coupled with the periodic magnetic structure give rise to a periodic ponderomotive potential. The electrons bunch together and radiate coherently.¹³⁹ The spontaneous emission pattern of a relativistic electron of energy $E = \gamma mc^2$ with $\gamma \gg 1$ is mostly in the forward direction. For a magnetic wiggler of period λ_w , the spectrum in the forward direction is symmetric about the wavelength $\lambda_s \cong \lambda_w / 2\gamma^2$. Thus, a change of the periodicity of the wiggler λ_w can be used to tune the coherent light emitted by the FEL over a very wide range.

Many interesting features of FEL can be understood classically. However, the quantum-statistical properties of radiation emitted by FEL exhibit many interesting features such as squeezing and sub-poissonian statistics.¹⁴⁰⁻¹⁴²

Here we describe a free-electron amplifier in the small-signal noncollective regime. Such an FEL can be described by the one-electron nonrelativistic Bambini-Renieri hamiltonian which refers to a moving frame, where the laser and the wiggler frequencies coincide with $\omega = ck/2$.¹⁴³ In this frame, resonance occurs when the electron is at rest; therefore, the electron can be treated nonrelativistically. The hamiltonian is given by

$$\mathcal{H} = \frac{p^2}{2m} + \hbar \omega A^\dagger A + i\hbar g(A - A^\dagger) \quad (147)$$

with $A = a \exp(ikz)$. Here a is the annihilation operator of the laser field, p and z are the electron's momentum and coordinate with $[z, p] = i\hbar$, $[A, A^\dagger] = 1$, $[p, A] = \hbar k A$, m is the effective mass of the electron, and

$$g = \left(\frac{e^2 B}{mk} \right) \left(\frac{2}{V \epsilon_0 \hbar \omega} \right)^{1/2} \quad (148)$$

with V the quantization volume and B the magnetic-field strength of the wiggler field in the moving frame. In Eq. (147) we have already taken the classical limit of the wiggler field. By transforming to the interaction picture we obtain

$$\mathcal{H}_I = ig\hbar \left\{ \exp \left[-\frac{it(\hbar k^2 + 2kp)}{2m} \right] A^\dagger - Hc \right\} \quad (149)$$

We now consider an initial state made up by an electron with momentum p and the field vacuum, i.e., $|\text{in}\rangle = |\bar{p}, 0\rangle$

$$p|\bar{p}, 0\rangle = \bar{p}|\bar{p}, 0\rangle \quad (150)$$

$$A|\bar{p}, 0\rangle = 0 \quad (151)$$

$$A^\dagger|\bar{p}, 0\rangle = |\bar{p} - \hbar k, 1\rangle \quad (152)$$

The final-state expectation value of any operator $O(A, A^\dagger)$ is then

$$\langle \text{out} | O | \text{out} \rangle = \langle \bar{p}, 0 | \bar{s}^\dagger(T) O S(T) | \bar{p}, 0 \rangle \quad (153)$$

where

$$S(T) = T \exp \left[-\frac{i}{\hbar} \int_{-T/2}^{T/2} dt H_I(t) \right] \quad (154)$$

is the time-evolution operator for the electron-photon state.

The evaluation of Eq. (153) is straightforward in the small-signal limit along the lines given in Ref. 141, and we obtain

$$(\Delta A_1)^2 = \frac{1}{4} - \frac{\hbar k^2}{2m} j \frac{\partial j}{\partial \beta} \quad (155a)$$

$$(\Delta A_2)^2 = \frac{1}{4} + \frac{\hbar k^2}{2m} j \frac{\partial j}{\partial \beta} \quad (155b)$$

$$(\Delta A_1)(\Delta A_2) = \frac{1}{4} \quad (155c)$$

$$\Delta n^2 - \langle n \rangle = -\frac{2\hbar k^2}{m} j^3 \frac{\partial j}{\partial \beta} \quad (155d)$$

where

$$j = \left(\frac{2g}{\beta} \right) \sin \left(\frac{\beta T}{2} \right) \quad (156)$$

$$\beta = \frac{k\bar{p}}{m} \quad (157)$$

In our notation, the gain of the free-electron laser is proportional to $-j\partial j/\partial\beta$. Hence, Eqs. (155a) and (155b) show that, depending on the sign of the gain, either A_1 or A_2 is squeezed while, because of Eq. (155c), minimum uncertainty is maintained. Here, we have defined squeezing with respect to the operator A instead of the annihilation operator a of the radiation field. This must be so because we employ electron-photon states, and the annihilation of a photon always comes up to increasing the momentum of the electron by $\hbar k$. Finally, Eq. (155d) shows that we have subpoissonian, poissonian, or superpoissonian statistics if the electron momentum is below resonance ($\beta < 0$), at resonance ($\beta = 0$), or below resonance ($\beta > 0$), respectively.

Exploiting the Quantized Center-of-Mass Motion of Atoms

In the treatment of the interaction of a two-level atom with photons of a single, dynamically privileged mode by the Jaynes-Cummings model, as discussed in Sec. 23.4, the center-of-mass motion of the atom is regarded as classical. This is a well-justified approximation, since the atom's kinetic energy of typically $\sim 10^{-2}$ eV is many orders of magnitude larger than the interaction energy of typically $\sim 10^{-11}$ eV if the atom belongs to a thermal beam. For ultracold atoms, however, matters can be quite different, and the quantum properties of the center-of-mass motion must be taken into account.

Early studies showed that very slow atoms can be reflected at the entry port of a resonator¹⁴⁴ or trapped inside.¹⁴⁵ The reflection probability is considerable even if the photon lifetime is not short as compared with the relatively long interaction time.¹⁴⁶

Whereas Refs. 144 to 146 deal mainly with the mechanical effects on the center-of-mass motion of the atom, the modifications in the maser action are addressed in Refs. 147 to 150. For thermal atoms, the emission probability displays the usual Rabi oscillations (see Sec. 23.4) as a function of the interaction *time*. For very slow atoms, however, the emission probability is a function of the interaction *length* and exhibits resonances such as the ones observed in the intensity transmitted by a Fabry-Perot resonator. The resonances occur when the resonator length is an integer multiple of half the de Broglie wavelength of the atom inside the cavity.

A detailed calculation¹⁴⁷ shows that the emission probability is 50 percent at a resonance, irrespective of the number of photons that are present initially. Owing to this unusual emission probability, a beam of ultracold atoms can produce unusual photon distributions, such as a shifted thermal distribution. In the trilogy (Refs. 148 to 150) this *microwave amplification by z-motion-induced emission of radiation* (mazer) is studied in great detail.

In order to see the mazer resonances for atoms with a certain velocity spread, the interaction length has to be small. Therefore, micromaser cavities of the usual cylindrical shape, for which the smallest cavity length is given by half the wavelength of the microwaves, cannot be used for this purpose. But cavities of the reentrant type (familiar as components of klystrons) allow for an interaction length that is much smaller than the wavelength. With such a device, an experiment with realistic parameters seems possible.¹⁴⁹ As a potential application, we mention that a working mazer could be used as a velocity filter for atoms.¹⁵¹

23.7 ACKNOWLEDGMENTS

The authors gratefully acknowledge the generous support of the Office of Naval Research over the many years spent on completing the works reviewed here. It is also a pleasure to acknowledge the Max-Planck-Institute for Quantum Optics (Garching, Germany) for providing the excellent working atmosphere and the intellectual stimulus for the most interesting works in this field.

23.8 REFERENCES

1. J. R. Klauder and E. C. G. Sudarshan, *Fundamentals of Quantum Optics* (W. A. Benjamin, New York, 1970).
2. R. Loudon, *The Quantum Theory of Light* (Oxford University Press, New York, 1973).
3. W. H. Louisell, *Quantum Statistical Properties of Radiation* (John Wiley, New York, 1973).
4. H. M. Nussenzveig, *Introduction to Quantum Optics* (Gordon and Breach, New York, 1974).
5. M. Sargent III, M. O. Scully, and W. E. Lamb, Jr., *Laser Physics* (Addison-Wesley, Reading, Mass., 1974).
6. L. Allen and J. H. Eberly, *Optical Resonance and Two-Level Atoms* (John Wiley, New York, 1975).
7. H. Haken, *Light*, Vols. I and II (North-Holland, Amsterdam, 1981).
8. P. L. Knight and L. Allen, *Concepts of Quantum Optics* (Pergamon Press, Oxford, 1983).
9. P. Meystre and M. Sargent III, *Elements of Quantum Optics* (Springer-Verlag, Berlin, 1990).
10. C. W. Gardiner, *Quantum Noise* (Springer-Verlag, Berlin, 1991).
11. C. Cohen-Tannoudji, J. Dupont-Roc, and G. Grynberg, *Atom-Photon Interactions* (John Wiley, New York, 1992).
12. H. Carmichael, *An Open Systems Approach to Quantum Optics* (Springer-Verlag, Berlin, 1993).
13. W. Vogel and D.-G. Welsch, *Lectures on Quantum Optics* (Akademie Verlag, Berlin, 1994).
14. J. Peřina, Z. Hradil, and B. Jurčo, *Quantum Optics and Fundamentals of Physics* (Kluwer, Dordrecht, Netherlands, 1994).
15. D. F. Walls and G. J. Milburn, *Quantum Optics* (Springer-Verlag, Berlin, 1994).
16. L. Mandel and E. Wolf, *Optical Coherence and Quantum Optics* (Cambridge University Press, London, 1995).
17. E. R. Pike and S. Sarkar, *Quantum Theory of Radiation* (Cambridge University Press, London, 1995).
18. M. O. Scully and M. S. Zubairy, *Quantum Optics* (Cambridge University Press, London, 1997).
19. M. Planck, *Verh. Phys. Ges.* **2**:202, 237 (1900).
20. W. Wien, *Ann. Physik* **58**:662 (1896).
21. Lord Rayleigh, *Phil. Mag.* **49**:539 (1900).
22. J. H. Jeans, *Phil. Mag.* **10**:91 (1905).
23. A. H. Compton, *Phys. Rev.* **21**:483 (1923).
24. A. Einstein, *Ann. Physik* **17**:132 (1905).
25. W. Pauli, "Einstein's Contributions to Quantum Theory," in *Albert Einstein: Philosopher-Scientist*, P. A. Schilpp (ed.) (Library of Living Philosophers, Evanston, Ill., 1949).
26. L. de Broglie, *C. R. Acad. Sci. Paris* **177**:517 (1923).
27. L. de Broglie, *These* (Masson et Cie., Paris, 1924).
28. E. Schrödinger, *Ann. Physik* **79**:361 (1926).
29. G. I. Taylor, *Proc. Camb. Phil. Soc.* **15**:114 (1909).
30. P. A. M. Dirac, *The Principles of Quantum Mechanics*, 4th ed. (Oxford University Press, Oxford, 1958).
31. R. J. Glauber, *Am. J. Phys.* **63**:12 (1995).
32. A. Einstein, *Phys. Z.* **10**:185 (1909).
33. N. Bohr, *Phil. Mag.* **26**:1,476, 857 (1913).
34. A. Einstein, *Phys. Z.* **18**:121 (1917).
35. A. Einstein, *Ann. Physik* **17**:549 (1905).
36. N. Bohr, H. A. Kramers, and I. C. Slater, *Phil. Mag.* **47**:785 (1924).
37. W. Bothe and H. Geiger, *Z. Phys.* **26**:44 (1924).
38. S. N. Bose, *Z. Phys.* **26**:178 (1924).
39. S. N. Bose, *Z. Phys.* **27**:384 (1924).
40. A. Einstein, "Quantentheorie des einatomigen Gases," in *Sitzungsber. Preuss. Akad. Wiss., Phys.-math. Kl.*, 1924, p. 261.

41. A. Einstein, "Quantentheorie des einatomigen Gases. 2. Abhandlung," in *Sitzungsber. Preuss. Akad. Wiss., Phys.-math. Kl.*, 1925, p. 3.
42. A. Einstein, "Quantentheorie des idealen Gases," in *Sitzungsber. Preuss. Akad. Wiss., Phys.-math. Kl.*, 1925, p. 18.
43. P. A. M. Dirac, *Proc. Roy. Soc. London* **A114**:710 (1927).
44. J. Schwinger (ed), *Quantum Electrodynamics* (Dover, New York, 1958).
45. E. Fermi, *Rend. Lincei* **10**:72 (1929).
46. R. J. Glauber, *Phys. Rev.* **130**:2529 (1963).
47. R. J. Glauber, *Phys. Rev.* **131**:2766 (1963).
48. R. J. Glauber, *Phys. Rev. Lett.* **10**:84 (1963).
49. V. F. Weisskopf and E. P. Wigner, *Z. Phys.* **63**:47 (1930).
50. V. F. Weisskopf, *Phys. Rev.* **56**:72 (1939).
51. W. E. Lamb, Jr., and R. C. Retherford, *Phys. Rev.* **72**:241 (1947).
52. J. Schwinger, *Phys. Rev.* **73**:416 (1948).
53. S. Pasternack, *Phys. Rev.* **54**:1113 (1938).
54. S. Millman and P. Kusch, *Phys. Rev.* **57**:438 (1940).
55. H. Hanbury-Brown and R. Q. Twiss, *Phil. Mag.* **45**:663 (1954).
56. H. Hanbury-Brown and R. Q. Twiss, *Nature* (London) **178**:1046 (1956).
57. H. Hanbury-Brown and R. Q. Twiss, *Proc. Roy. Soc.* **A242**:300 (1957).
58. R. G. Newton, *Scattering Theory of Waves and Particles* (McGraw-Hill, New York, 1966).
59. C. K. Hong, Z. Y. Ou, and L. Mandel, *Phys. Rev. Lett.* **59**:2044 (1987).
60. P. G. Kwiat, K. Mattle, H. Weinfurter, and A. Zeilinger, *Phys. Rev. Lett.* **75**:4337 (1995).
61. P. G. Kwiat, E. Waks, A. G. White, I. Appelbaum, and P. H. Eberhard, *Phys. Rev. A* **60**:773 (1999).
62. D. Bouwmeester, J.-W. Pan, K. Mattle, M. Eibl, H. Weinfurter, and A. Zeilinger, *Nature* (London) **390**:575 (1997).
63. D. Boschi, S. Branca, F. De Martini, L. Hardy, and S. Popescu, *Phys. Rev. Lett.* **80**:1121 (1998).
64. C. H. Bennett, G. Brassard, C. Crepeau, R. Josza, A. Peres, and W. Wootters, *Phys. Rev. Lett.* **70**:1895 (1993).
65. K. Mattle, H. Weinfurter, P. G. Kwiat, and A. Zeilinger, *Phys. Rev. Lett.* **76**:4656 (1996).
66. C. H. Bennett and S. J. Wiesner, *Phys. Rev. Lett.* **69**:2881 (1992).
67. M. Sargent, M. O. Scully, and W. E. Lamb, *Laser Physics* (Addison-Wesley, Reading, Mass., 1974).
68. M. O. Scully and M. S. Zubairy, *Quantum Optics* (Cambridge University Press, Cambridge, 1997).
69. M. Lax, "Phase Noise in a Homogeneously Broadened Maser," in *Physics of Quantum Electronics*, P. L. Kelley, B. Lax, and P. E. Tannenwald (eds.) (McGraw-Hill, New York, 1966).
70. M. Lax, *Phys. Rev.* **145**:110 (1966).
71. M. Lax, "Fluctuations and Coherence Phenomena in Classical and Quantum Physics," in *1966 Brandeis Summer Lecture Series on Statistical Physics*, Vol. 2, M. Chretien, E. P. Gross, and S. Dreser (eds.) (Gordon and Breach, New York, 1968; Mir, Moscow, 1975).
72. W. H. Louisell, *Quantum Statistical Properties of Radiation* (John Wiley, New York, 1973).
73. H. Haken, "Laser Theory," in *Encyclopedia of Physics*, Vol. 35/c, S. Flügge (ed.) (Springer, Berlin, 1970).
74. H. Risken, *The Fokker Planck Equation* (Springer, Heidelberg, 1984).
75. M. Lax, "The Theory of Laser Noise," keynote address, 1990 Conference on Laser Science and Optics. Applications, *Proc. SPIE* **1376**:2 (1991).
76. E. T. Jaynes and F. W. Cummings, *Proc. IEEE* **51**:89 (1963).
77. B. W. Shore and P. L. Knight, *J. Mod. Opt.* **40**:1195 (1993).
78. J. P. Gordon, H. J. Zeiger, and C. H. Townes, *Phys. Rev.* **95**:282L (1955).
79. J. P. Gordon, H. J. Zeiger, and C. H. Townes, *Phys. Rev.* **99**:1264 (1955).
80. A. L. Schawlow and C. H. Townes, *Phys. Rev. A* **112**:1940 (1958).
81. N. G. Basov and A. M. Prokhorov, *Dokl. Ak. Nauk* **101**:47 (1955).

82. T. H. Maiman, *Nature* (London) **187**:493 (1960).
83. A. Javan, W. R. Bennett, and D. R. Herriott, *Phys. Rev. Lett.* **6**:106 (1961).
84. W. E. Lamb, *Phys. Rev.* **134**:A1429 (1964).
85. For recent reviews, see, e.g., B.-G. Englert, M. Löffler, O. Benson, B. Varcoe, M. Weidinger, and H. Walther, *Fortschr. Phys.* **46**:897 (1998); G. Raithel, C. Wagner, H. Walther, L. M. Narducci, and M. O. Scully, in *Advances in Molecular and Optical Physics*, P. Berman (ed.) (Academic, New York, 1994), Suppl. 2.
86. D. Meschede, H. Walther, and G. Müller, *Phys. Rev. Lett.* **54**:551 (1985).
87. P. Filipowicz, J. Javanainen, and P. Meystre, *Phys. Rev. A* **34**:3077 (1986).
88. M. Lax, *Phys. Rev.* **129**:2342 (1963).
89. M. Lax and M. Zwanziger, *Phys. Rev. A* **7**:750 (1973).
90. M. Lax, *Phys. Rev.* **172**:350 (1968).
91. N. Lu, *Phys. Rev. A* **47**:4322 (1993).
92. U. Herzog and J. Bergou, *Phys. Rev. A* **62** (2000). In press.
93. M. Lax, *Phys. Rev.* **160**:290 (1967).
94. R. D. Hempstead and M. Lax, *Phys. Rev.* **161**:350 (1967).
95. H. Risken and H. D. Vollmer, *Z. Physik* **191**:301 (1967).
96. M. Kac, in *1966 Brandeis Summer Lecture Series on Statistical Physics*, Vol. 1, M. Chretien, E. P. Gross, and S. Dreser (eds.) (Gordon and Breach, New York, 1968).
97. M. O. Scully, H. Walther, G. S. Agarwal, T. Quang, and W. Schleich, *Phys. Rev. A* **44**:5992 (1991).
98. N. Lu, *Phys. Rev. Lett.* **70**:912 (1993); N. Lu, *Phys. Rev. A* **47**:1347 (1993); T. Quang, G. S. Agarwal, J. Bergou, M. O. Scully, H. Walther, K. Vogel, and W. P. Schleich, *Phys. Rev. A* **48**:803 (1993); K. Vogel, W. P. Schleich, M. O. Scully, and H. Walther, *Phys. Rev. A* **48**:813 (1993); R. McGowan and W. Schieve, *Phys. Rev. A* **55**:3813 (1997).
99. S. Qamar and M. S. Zubairy, *Phys. Rev. A* **44**:7804 (1991).
100. P. Filipowicz, J. Javanainen, and P. Meystre, *J. Opt. Soc. Am. B* **3**:154 (1986).
101. J. J. Slosser, P. Meystre, and S. L. Braunstein, *Phys. Rev. Lett.* **63**:934 (1989).
102. J. Bergou, L. Davidovich, M. Orszag, C. Benkert, M. Hillery, and M. O. Scully, *Phys. Rev. A* **40**:7121 (1989); J. D. Cresser, *Phys. Rev. A* **46**:5913 (1992); U. Herzog, *Phys. Rev. A* **52**:602 (1995).
103. H. J. Briegel, B.-G. Englert, N. Sterpi, and H. Walther, *Phys. Rev. A* **49**:2962 (1994); U. Herzog, *Phys. Rev. A* **50**:783 (1994); J. D. Cresser and S. M. Pickless, *Phys. Rev. A* **50**:R925 (1994); U. Herzog, *Appl. Phys. B* **60**:S21 (1995).
104. H.-J. Briegel, B.-G. Englert, Ch. Ginzel, and A. Schenzle, *Phys. Rev. A* **49**:5019 (1994).
105. H.-J. Briegel, B.-G. Englert, *Phys. Rev. A* **52**:2361 (1995); J. Bergou, *Quantum and Semiclass. Optics* **7**:327 (1995); U. Herzog and J. Bergou, *Phys. Rev. A* **54**:5334 (1996); *ibid.* **55**:1385 (1997).
106. M. Weidinger, B. T. H. Varcoe, R. Heerlein, and H. Walther, *Phys. Rev. Lett.* **82**:3795 (1999).
107. H. Walther, *Phys. Rep.* **219**:263 (1992).
108. M. O. Scully, G. Süssmann, and C. Benkert, *Phys. Rev. Lett.* **60**:1014 (1988).
109. M. O. Scully, M. S. Zubairy, and K. Wódkiewicz, *Opt. Commun.* **65**:440 (1988).
110. H. E. Stanley, *Introduction to Phase Transitions and Critical Phenomena* (Oxford University Press, Oxford, 1971).
111. V. DeGiorgio and M. O. Scully, *Phys. Rev. A* **2**:1170 (1970).
112. M. Anderson, J. Ensher, M. Matthews, C. Wieman, and E. Cornell, *Science* **269**:198 (1995).
113. C. Bradley, C. Sackett, J. Tollett, and R. Hulet, *Phys. Rev. Lett.* **75**:1687 (1995).
114. K. Davis, M. Mewes, M. Andrews, N. van Druten, D. Durfee, D. Kurn, and W. Ketterle, *Phys. Rev. Lett.* **75**:3969 (1995).
115. M. H. W. Chen, K. I. Blum, S. Q. Murphy, G. K. S. Wong, and J. D. Reppy, *Phys. Rev. Lett.* **61**:1950 (1998).
116. K. Huang, *Statistical Mechanics* (John Wiley, New York, 1987).
117. R. Arnowitt and M. Girardeau, *Phys. Rev.* **113**:745 (1959).

118. G. Baym and C. Pethick, *Phys. Rev. Lett.* **76**:6 (1996).
119. H. Politzer, *Phys. Rev. A* **54**:5048 (1996).
120. S. Grossmann and M. Holthaus, *Phys. Rev. E* **54**:3495 (1996).
121. C. Herzog and M. Olshanii, *Phys. Rev. A* **55**:3254 (1997).
122. P. Navez, D. Bitouk, M. Gajda, Z. Idziaszek, and K. Rzążewski, *Phys. Rev. Lett.* **79**:1789 (1997).
123. M. Wilkens and C. Weiss, *J. Mod. Opt.* **44**:1801 (1997).
124. S. Grossmann and M. Holthaus, *Phys. Rev. Lett.* **79**:3557 (1997).
125. H. Woolf (ed.), *Some Strangeness in the Proportion: A Centennial Symposium to Celebrate the Achievements of Albert Einstein* (Addison-Wesley, Reading, Mass., 1980), p. 524.
126. J. Goldstein, M. O. Scully, and P. Lee, *Phys. Lett.* **A35**:317 (1971).
127. M. O. Scully, *Phys. Rev. Lett.* **82**:3927 (1999).
128. O. Kocharovskaya and Ya I. Khanin, *JETP Lett.* **48**:630 (1988).
129. S. E. Harris, *Phys. Rev. Lett.* **62**:1033 (1989).
130. M. O. Scully, S.-Y. Zhu, and A. Gavrieleides, *Phys. Rev. Lett.* **62**:2813 (1989).
131. A. Nottelmann, C. Peters, and W. Lange, *Phys. Rev. Lett.* **70**:1783 (1993).
132. E. S. Fry, X. Li, D. Nikonov, G. G. Padmabandu, M. O. Scully, A. V. Smith, F. K. Tittel, C. Wang, S. R. Wilkinson, and S.-Y. Zhu, *Phys. Rev. Lett.* **70**:3235 (1993).
133. W. E. van der Veer, R. J. J. van Diest, A. Dönszelmann, and H. B. van Linden van den Heuvell, *Phys. Rev. Lett.* **70**:3243 (1993).
134. A. S. Zibrov, M. D. Lukin, D. E. Nikonov, L. W. Hollberg, M. O. Scully, V. L. Velichansky, and H. G. Robinson, *Phys. Rev. Lett.* **75**:1499 (1995).
135. M. O. Scully, *Phys. Rev. Lett.* **55**:2802 (1985).
136. J. Bergou, M. Orszag, and M. O. Scully, *Phys. Rev. A* **38**:768 (1988).
137. M. O. Scully and M. S. Zubairy, *Phys. Rev. A* **35**:752 (1987).
138. M. O. Scully, K. Wodkiewicz, M. S. Zubairy, J. Bergou, N. Lu, and J. Meyer ter Vehn, *Phys. Rev. Lett.* **60**:1832 (1988).
139. J. M. J. Madey, *J. Appl. Phys.* **42**:1906 (1971).
140. W. Becker, M. O. Scully, and M. S. Zubairy, *Phys. Rev. Lett.* **48**:475 (1985).
141. W. Becker and M. S. Zubairy, *Phys. Rev. A* **25**:2200 (1982).
142. R. Bonifacio, *Opt. Commun.* **32**:440 (1980).
143. A. Bambini and A. Renieri, *Opt. Commun.* **29**:244 (1978).
144. B.-G. Englert, J. Schwinger, A. O. Barut, and M. O. Scully, *Europhys. Lett.* **14**:25 (1991).
145. S. Haroche, M. Brune, and J.-M. Raimond, *Europhys. Lett.* **14**:19 (1991).
146. M. Battocletti and B.-G. Englert, *J. Phys. II France* **4**:1939 (1994).
147. M. O. Scully, G. M. Meyer, and H. Walther, *Phys. Rev. Lett.* **76**:4144 (1996).
148. G. M. Meyer, M. O. Scully, and H. Walther, *Phys. Rev. A* **56**:4142 (1997).
149. M. Löffler, G. M. Meyer, M. Schröder, M. O. Scully, and H. Walther, *Phys. Rev. A* **56**:4153 (1997).
150. M. Schröder, K. Vogel, W. P. Schleich, M. O. Scully, and H. Walther, *Phys. Rev. A* **56**:4164 (1997).
151. M. Löffler, G. M. Meyer, and H. Walther, *Lett.* **41**:593 (1998).

This page intentionally left blank.

DO NOT DUPLICATE

PART

5

DETECTORS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

PHOTODETECTORS

Paul R. Norton

*U.S. Army Night Vision and Electronics Directorate
Fort Belvoir, Virginia*

*Second revision and update from an article by Stephen F. Jacobs**

24.1 SCOPE

The purpose of this chapter is to describe the range of detectors commercially available for sensing optical radiation. Optical radiation over the range from vacuum ultraviolet to the far-infrared or submillimeter wavelength (25 nm to 1000 μm) is considered. We will refer to the following spectral ranges:

25–200 nm	vacuum ultraviolet	VUV
200–400 nm	ultraviolet	UV
400–700 nm	visible	VIS
700–1000 nm	near infrared	NIR
1–3 μm	short-wavelength infrared	SWIR
3–5 μm	medium-wavelength infrared	MWIR
5–14 μm	long-wavelength infrared	LWIR
14–30 μm	very long wavelength infrared	VLWIR
30–100 μm	far-infrared	FIR
100–1000 μm	submillimeter	SubMM

We begin by giving a brief description of the photosensitive mechanism for each type of detector. The usefulness and limitations of each detector type are also briefly described. Definitions of the technical terms associated with the detection process are listed. The concept of sensitivity is defined, and D^* (D^*) is presented as a measure of ideal performance. Examples are then given of the limiting cases for D^* under several conditions. In addition, other detector performance parameters are described which may be of critical interest in a specific application, including spectral response, responsivity, quantum efficiency, noise, uniformity, speed, and stability. Finally, manufacturers' specifications for a range of available detectors are compiled and a list of manufacturers is included for each type of detector.

*In *Handbook of Optics*, first edition, McGraw-Hill, 1978. Section 4, "Nonimaging Detectors," by Stephen F. Jacobs, Optical Sciences Center, University of Arizona, Tucson, Arizona.

The sensitivity of many detectors has reached the limit set by room-temperature background photon fluctuations (radiation shot noise). For these detectors, sensitivity may be enhanced by providing additional cooling, while restricting their spatial field of view and/or spectral bandwidth. At some point, other factors such as amplifier noise may limit the improvement.

Techniques for evaluating detector performance are not covered in this treatment but can be found in Refs. 1–6.

24.2 THERMAL DETECTORS

Thermal detectors sense the change in temperature produced by absorption of incident radiation. Their spectral response can therefore be as flat as the absorption spectrum of their blackened coating and window* will allow. This makes them useful for spectroscopy and radiometry. These detectors are generally operated at room temperature, where their sensitivity is limited by thermodynamic considerations^{7,8} to 3 pW for 1-s measurement time and 1-mm² sensitive area. This limit has been very nearly reached in practice, whereas cooled bolometers have been made to reach the background photon noise limit. Figure 1 illustrates the basic structural elements of a thermal detector.

Construction of the detector seeks to minimize both the thermal mass of the sensitive element and the heat loss from either conductive or convective mechanisms. Heat loss may ideally be dominated by radiation. This allows the incident photon flux to give a maximum temperature rise (maximum signal), but results in a correspondingly slow response time for this class of detectors. The response time τ of thermal detectors is generally slower than 1 ms, depending on thermal capacity C and heat loss per second per degree G , through the relation

$$\tau = C/G$$

A short time constant requires a small C . However, for room-temperature operation, ultimate sensitivity is limited by the mean spontaneous temperature fluctuation:

$$\Delta T = T \sqrt{\frac{k}{C}}$$

where k is Boltzmann's constant. There is thus a trade-off between time constant and ultimate sensitivity.

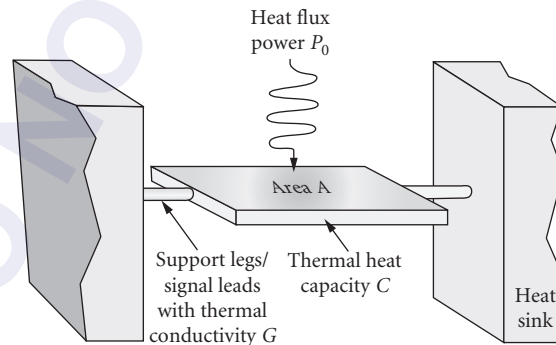


FIGURE 1 Structural elements of a thermal detector. Sensitive area A with a thermal heat capacity of C , supported by leads having thermal conductivity G , and with a heat flux of P_0 incident on the pixel.

*No windows exist without absorption bands anywhere between the visible and millimeter region. Some useful window materials for the far-infrared are diamond, silicon, polyethylene, quartz, and CsI.

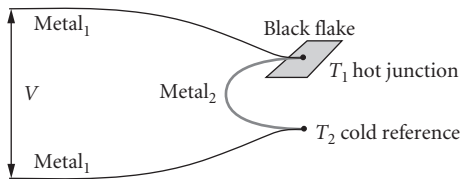


FIGURE 2 Thermocouple detector structure.

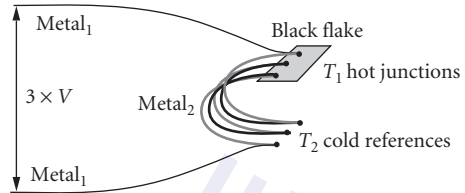


FIGURE 3 Thermopile detector structure.

Thermocouple/Thermopile

The thermocouple receiver, illustrated in Fig. 2, is a thin, blackened flake connected thermally to the junction of two dissimilar metals or semiconductors. Heat absorbed by the flake causes a temperature rise of the junction, and hence a thermoelectric emf is developed which can be measured, for example, with a voltmeter.

Thermocouples are limited in sensitivity by thermal (Johnson-Nyquist) noise but are nevertheless respectably sensitive. Their usefulness lies in the convenience of room temperature operation, wide spectral response, and their rugged construction. Thermocouples are widely used in spectroscopy.

Thermopiles consist of thin-film arrays of thermocouples in series, as illustrated in Fig. 3. This device multiplies the thermocouple signal corresponding to the number of junctions in series. The device may be constructed with half the thermocouples acting as reference detectors attached to a heat sink.

Bolometer/Thermistor

The receiver is a thin, blackened flake or slab, whose impedance is highly temperature dependent—see Fig. 4. The impedance change may be sensed using a bridge circuit with a reference element in the series or parallel arm. Alternatively, a single bolometer element in series with a load and voltage source may be used.

Most bolometers in use today are of the thermistor type made from oxides of manganese, cobalt, or nickel. Their sensitivity closely approaches that of the thermocouple for frequencies higher than 25 Hz. At lower frequencies there may be excess or $1/f$ noise. Construction can be very rugged for systems applications. Some extremely sensitive low-temperature semiconductor and superconductor bolometers are available commercially.

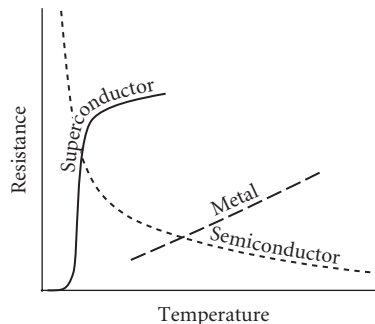


FIGURE 4 Temperature dependence characteristics of three bolometer material types.

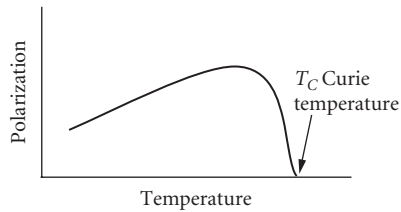


FIGURE 5 Ferroelectric materials exhibit residual polarization with no applied bias.

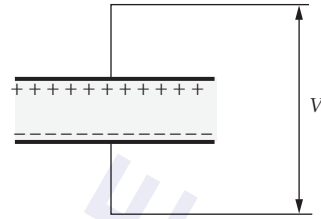


FIGURE 6 The pyroelectric effect produces a surface charge when the temperature changes.

Pyroelectric

Ferroelectric material exhibits a residual polarization in the absence of any electric field, as illustrated in Fig. 5. Dipole moments, initially aligned by applying an external field, result in a surface charge which is normally slowly neutralized by leakage. This polarization is temperature-dependent (pyroelectric effect), and when incident radiation heats an electroded sample, there is a change in surface charge (open-circuit voltage) which is proportional to the incident radiation power—see Fig. 6. Electrically, the device behaves like a capacitor, requiring no bias and therefore exhibiting no current noise. The signal, however, must be chopped or modulated. Sensitivity is limited either by amplifier noise or by loss-tangent noise. Response speed can be engineered, with a proportional decrease in sensitivity, making pyroelectric detectors useful for moderately fast laser pulse detection. Other common applications include power meters. Microphonic noise in applications associated with vibrations can be a problem with some of these devices.

24.3 QUANTUM DETECTORS

Photon detectors respond in proportion to incident photon rates (quanta) rather than to photon energies (heat). Thus, the spectral response of an ideal photon detector is flat on an incident-photon-rate basis but linearly rising with wavelength on an incident-power (per watt) basis. The sensitivity of efficient quantum detectors can approach the limits of photon noise fluctuations provided that the detector temperature is sufficiently low for photon-induced mechanisms to dominate thermally induced mechanisms in the detector. Quantum detectors generally have sub-microsecond time constants. Their main disadvantage is the associated cooling required for optimum sensitivity. (These remarks do not apply to photographic detection, which measures cumulative photon numbers.)

Photoemissive

The radiation is absorbed by a photosensitive surface which usually contains alkali metals (cesium, sodium, or potassium). Incident quanta release photoelectrons (Fig. 7), via the photoelectric effect, which are collected by a positively biased anode. This is called a diode phototube; it can be made the basis for the multiplier phototube (photomultiplier phototube, or photomultiplier) by the addition of a series of biased dynodes which serve as secondary emission multipliers.

In spectral regions where quantum efficiency is high ($\lambda < 550$ nm), the photoemissive detector is very nearly ideal. Sensitivity is high enough to count individual photons. Amplification does not degrade the signal-to-noise ratio. The sensitive area is conveniently large. Photomultiplier signal response time (transit spread time) can be made as short as 0.1 ns. Since the sensitivity in red-sensitive tubes is limited by thermally generated electrons, sensitivity can be improved by cooling.

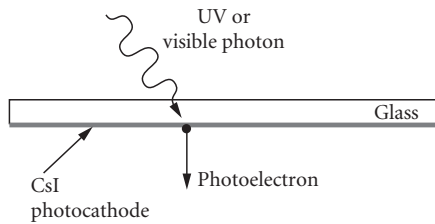


FIGURE 7 Electrons are ejected from a photoemissive surface when excited by photons.

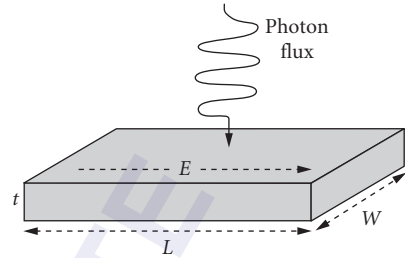


FIGURE 8 Photoconductor device structure.

Photoconductive

The radiation is absorbed by a photoconductive (PC) material, generally a semiconductor, either in thin-film or bulk form, as illustrated in Fig. 8. Each incident quantum may release an electron-hole pair or an impurity-bound charge carrier, thereby increasing the electrical conductivity. The devices are operated in series with a bias voltage and load resistor. Very low impedance photoconductors may be operated with a transformer as the load. Since the impedance of photoconductors varies with device type and operating conditions from less than $50\ \Omega$ to more than $10^{14}\ \Omega$, the load resistor and preamplifier must be chosen appropriately. Figure 9 shows the current-voltage characteristics of a photoconductor.

Photoconductors that utilize excitation of an electron from the valence to conduction band are called *intrinsic* detectors. Those which operate by exciting electrons into the conduction band or holes into the valence band from states within the band—impurity-bound states, quantum wells, or quantum dots—are called *extrinsic* detectors. Figure 10 illustrates these two mechanism types. Intrinsic detectors are most common at the short wavelengths, out to about $20\ \mu\text{m}$. Extrinsic detectors are most common at longer wavelengths. A key difference between intrinsic and extrinsic detectors is that intrinsic detectors do not require as much cooling to achieve high sensitivity at a given spectral response cutoff as extrinsic detectors. Thus, intrinsic photoconductors such as HgCdTe will operate out to 15 to $20\ \mu\text{m}$ at $77\ \text{K}$, while comparable extrinsic detectors with similar cutoff must be cooled below 30 to $40\ \text{K}$.

A further distinction may be made by whether the semiconductor material has a direct or indirect bandgap. This difference shows up near the long-wavelength limit of the spectral response

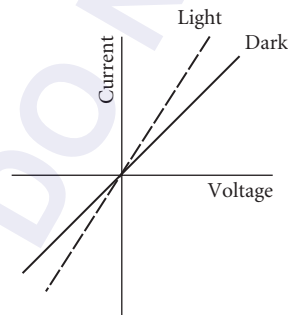


FIGURE 9 Photoconductor current-voltage characteristics in the dark and in the light.

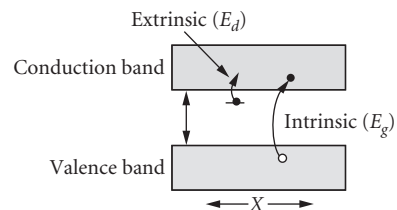


FIGURE 10 Intrinsic detectors excite electrons between the valence and conduction band. Extrinsic detectors excite electrons (or holes) from states within the band to the conduction (valence) band.

where detectors made from direct bandgap materials such as InGaAs, InSb, or HgCdTe have a sharper spectral cutoff than indirect bandgap materials such as silicon and germanium.

Photoconductors can have high quantum efficiency from the visible region out to the far infrared but lack the nearly ideal high amplification of photomultipliers. They are therefore most commonly used in the spectral region beyond $1\ \mu\text{m}$, where efficient photoemitters are unavailable. Photoconductors do, however, provide current gain which is equal to the recombination time divided by the majority-carrier transit time. This current gain leads to higher responsivity than is possible with (nonavalanching) photovoltaic (PV) detectors. For applications where photovoltaic detection would be amplifier noise limited, the larger photoconductive responsivity makes it possible to realize greater sensitivity with the photoconductor. In general, lower-temperature operation is associated with longer-wavelength sensitivity in order to suppress noise due to thermally induced transitions between close-lying energy levels. Ideally, photoconductors are limited by generation-recombination noise in the photon-generated carriers. Response time can be shorter than $1\ \mu\text{s}$ and in some cases response times can be shorter than $1\ \text{ns}$ for small elements. Response across a photoconductive element can be nonuniform due to recombination mechanisms at the electrical contacts, and this effect may vary with electrical bias.

Photovoltaic*

The most widely used photovoltaic detector is the pn junction type (Fig. 11), where a strong internal electric field exists across the junction even in the absence of radiation. Photons incident on the junction of this film or bulk material produce free hole-electron pairs which are separated by the internal electric field across the junction, causing a change in voltage across the open-circuit cell or a current to flow in the short-circuited case.

As with the photoconductor, quantum efficiency can be high from the visible to the very long-wavelength infrared, generally about 20 to $25\ \mu\text{m}$. The limiting noise level can ideally be $\sqrt{2}$ times lower than that of the photoconductor, thanks to the absence of recombination noise. Lower temperatures are associated with longer-wavelength operation. Response times can be less than a nanosecond,

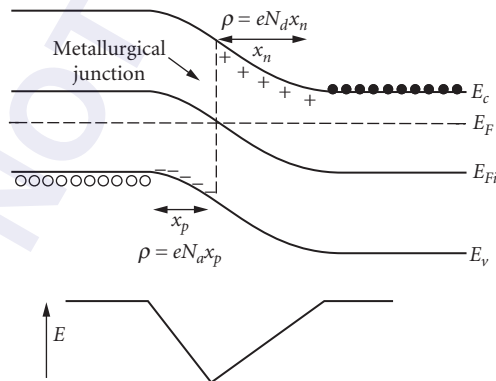


FIGURE 11 pn junction showing how the band bends across the junction. The junction width is given by $W = x_p + x_n$.

*The use of a photovoltaic detector at other than zero bias is often referred to as its photoconductive mode of operation because the circuit then is similar to the standard photoconductor circuit. This terminology is confusing with regard to detection mechanism and will not be used here.

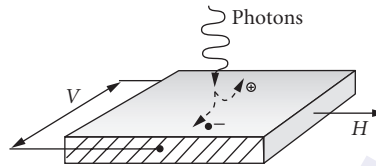


FIGURE 12 Photoelectromagnetic detector configuration—incident photons create electron-hole pairs that diffuse away from the surface. A magnetic field perpendicular to the diffusion separates the charges toward opposite sides, creating a voltage.

and are generally limited by device capacitance and detection-circuit resistance. The *pin* diode has been developed to minimize capacitance for high-bandwidth applications. The advantages of nearly ideal internal amplification have now become available in avalanche photodiodes sensitive out to $1.55\ \mu\text{m}$. This internal gain is most important for high-frequency operation, where external load resistance must be kept small and would otherwise introduce limiting thermal noise, and for situations involving low signal flux where amplifier noise is otherwise dominant.

Photoelectromagnetic

A thin slab of photoconductive material is oriented with radiation incident on a large face and a magnetic field perpendicular to it, as illustrated in Fig. 12. Electron-hole pairs generated by the incident photons diffuse through the material and are separated by the magnetic field, causing a potential difference at opposite ends of the detector.

These detectors require no cooling or biasing electric field but do require a (permanent) magnetic field. Photoelectromagnetic InSb at room temperature has response up to $7.5\ \mu\text{m}$, where it is as sensitive as a thermocouple of equal size, and has a response time less than $1\ \mu\text{s}$. Another competing uncooled detector is InAs, which is far more sensitive out to $3.5\ \mu\text{m}$. HgCdTe is also available in the photoelectromagnetic (PEM) configuration out to the LWIR spectral region. Cooled infrared detectors are one to two orders of magnitude more sensitive.

Photographic

The receiver is an emulsion containing silver halide crystals. Incident photons are absorbed by the halide ion, which subsequently loses its electron. This electron eventually recombines with a silver ion and reduces it to a neutral silver atom. As more photons are absorbed, this process is repeated until a small but stable cluster of reduced silver atoms is formed within the crystal (latent image). Internal amplification is provided by introduction of an electron donor (photographic) developer, which, using the latent image as a catalytic center, reduces all the remaining silver ions within the exposed crystal to neutral silver atoms. The density of reduced crystals is a measure of the total radiation exposure.

The spectral region of sensitivity for photographic detection coincides rather closely with that of the photoemissive detector. For $\lambda > 1.2\ \mu\text{m}$ ($\sim 1\ \text{eV}$) there is too little energy in each photon to form a stable latent image. The basic detection process for both detectors operates well for higher energy, shorter wavelength radiation. The problem in ultraviolet and x-ray operation is one of eliminating nonessential materials, for example, the emulsion which absorbs these wavelengths.

The photographic process is an integrating one in that the output (emulsion density) measures the cumulative effect of all the radiation incident during the exposure time. The efficiency of the photographic process can be very high, but it depends upon photon energy; for example, in the

visible region it takes only 10 to 100 photons to form a stable latent image (developable grain). The photographic process enjoys a large and efficient internal amplification ability (development) wherein the very small energy of the photons' interaction is converted into readily observed macroscopic changes. An extensive discussion of photographic detection is found in Chap. 29, "Photographic Films."

Photoionization

The radiation is absorbed by a gas. If the photon energy exceeds the gas-ionization threshold, ion pairs can be produced with very high efficiency. They are collected by means of an applied voltage. Operating in a dc mode, these detectors are known as ionization and gas-gain chambers. When a pulse mode is used, the detectors are known as proportional and photon (Geiger) counters.

Photoionization detectors have a high sensitivity and a low noise level. They may also be quite selective spectrally since the choice of window and gas independently set upper and lower limits on detectable photon energies. Manufacturers' specifications are not discussed for these detectors as applications are still few enough to be treated as individual cases.⁹

24.4 DEFINITIONS

The following definitions will be used:

Avalanche photodiode (APD) A photodiode designed to operate in strong reverse bias where electron and/or hole impact ionization produces multiplication of photogenerated carriers.

Background temperature The effective temperature of all radiation sources viewed by the detector exclusive of the signal source.

Blackbody D star D_{BB}^* (cm Hz^{1/2}/W also called "Jones") Similar to $D^*(\lambda)$ or $D^*(T_B f)$ except that the source is a blackbody whose temperature must be specified.

Blackbody detectivity D_{BB} (W⁻¹) A measure of detector sensitivity, defined as $DBB (NEP_{BB})^{-1}$.

Blackbody noise-equivalent power NEP_{BB} Same as spectral NEP, except that the source has blackbody spectral character whose temperature must be specified, for example, NEP (500 K, 1, 800) means 500-K blackbody radiation incident, 1-Hz electrical bandwidth, and 800-Hz chopping frequency.

Blackbody responsivity R_{BB} Same as spectral sensitivity except that the incident signal radiation has a blackbody spectrum (whose temperature must be specified).

Blip detector or blip condition Originally meaning background-limited impurity photoconductor, this term has come to mean performance of any detector where the limiting noise is due to fluctuations in the arrival rate of background photons.

Cutoff wavelength λ_c The wavelength at which the detectivity has degraded to one-half its peak value.

Dark current The output current which flows even when input radiation is absent or negligible. (Note that although this current can be subtracted out for the dc measurements, the shot noise on the dark current can become the limiting noise.)

Detective quantum efficiency The square of the ratio of measured detectivity to the theoretical limit of detectivity.

Detective time constant $\tau_d = 1/2\pi f_d$ where f_d is the frequency at which D^* drops to 0.707 ($1/\sqrt{2}$) times its maximum value. A physics convention defines it as $1/e$, or 0.368 of the maximum value.

Dewar A container (cryostat) for holding detector coolant.

Equivalent noise input (ENI) A term meaning nearly the same thing as NEP_{BB} (287 K, 1, f). The difference is that the peak-to-peak value of square-wave chopped input flux is used, rather than the rms value of the sinusoidally chopped input flux. (See recommendation IRE.²)

Excess noise A term usually referring to noises other than generation-recombination, shot, or thermal.

Extrinsic semiconductor transition Incident photons produce a free electron in the conduction band and bound hole at a donor impurity site or a bound electron at an acceptor impurity site and a free hole in the valence band by excitation of an impurity level.

Field of view (FOV) The solid angle from which the detector receives radiation.

Flicker noise See Modulation noise.

Generation noise Noise produced by the statistical fluctuation in the rate of production of photoelectrons.

Generation-recombination (GR) noise Charge carriers are generated both by (optical) photons and (thermal) phonons. Fluctuations in these generation rates cause noise; fluctuations in carrier-recombination times cause recombination noise. The phonon contribution can be removed by cooling. The remaining photon contribution is indistinguishable from radiation shot noise (photon noise). With photovoltaic pn junctions, carriers are swept away before recombination, so that recombination noise is absent.

Guard ring An electrically biased field plate or surrounding diode used in some photodiodes, usually used to control surface recombination effects and thus reduce the leakage current in the detection circuit.

Intrinsic semiconductor transition Incident photons produce a free electron-free hole pair by direct excitation across the forbidden energy gap (valence to conduction band).

Johnson noise Same as Thermal noise.

Jones Unit of measure for D^* $\text{cm Hz}^{1/2}/W$.

Maximized D star, $D^(\lambda_{pk}, f_o)$, $\text{cm Hz}^{1/2}/W$ or Jones* The value of $D^*(\lambda_{pk}, f)$ corresponding to wavelength λ_{pk} and chopping frequency of maximum D^* .

Modulation (or 1/f) noise A consensus regarding the origin(s) of the mechanism has not been established, and although a quantum theory has been proposed, other mechanisms may dominate. As its name implies, it is characterized by a $1/f^n$ noise power spectrum, where $0.8 < n < 2$. This type of noise is prominent in thermal detectors and may dominate the low-frequency noise characteristics of photoconductive and photovoltaic quantum detectors as well as other electronic devices such as transistors and resistors.

Multiplier phototube or multiplier photodiode Phototube with built-in amplification via secondary emission from electrically biased dynodes.

NEI photons/($\text{cm}^2 \text{sec}$) noise equivalent irradiance is the signal flux level at which the signal produces the same output as the noise present in the detector. This unit is useful because it directly gives the photon flux above which the detector will be photon noise limited. See also Spectral noise equivalent power (NEP).

Noise spectrum The electrical power spectral density of the noise.

Nyquist noise Same as Thermal noise.

Photo cell See Photodiode.

Photoconductive gain The ratio of carrier lifetime divided by carrier transit time in a biased photoconductor.

Photodiode The term photodiode has been applied both to vacuum- or gas-filled photoemissive detectors (diode phototubes, or photo cells) and to photovoltaic detectors (semiconductor pn junction devices).

Photomultiplier Same as Multiplier phototube.

Photon counting Digital counting of individual photons from the photoelectrons produced in the detector in contrast to averaging of the photocurrent. This technique leads to very great sensitivity but can be used only for quite low light levels.

Quantum efficiency The ratio of the number of countable output events to the number of incident photons, for example, photoelectrons per photon, usually referred to as a percentage value.

RMS noise $V_{n,rms}$ That component of the electrical output which is not coherent with the radiation signal (generally measured with the signal radiation removed).

RMS signal $V_{s,rms}$ That component of the electrical output which is coherent with the input signal radiation.

Response time τ Same as Time constant.

Responsive quantum efficiency See Quantum efficiency.

Sensitivity Degree to which detector can sense small amounts of radiation.

Shot noise This current fluctuation results from the random arrival of charge carriers, as in a photodiode. Its magnitude is set by the size of the unit charge.

$$i_{n,rms} = (2ei_{dc}\Delta f)^{1/2}$$

Spectral D-double-star $D^{**}(\lambda, f)$ A normalization of D^* to account for detector field of view. It is used only when the detector is background-noise-limited. If the FOV is 2π sr, $D^{**} = D^*$.

Spectral D-star $D^*(\lambda, f)$ ($cm\ Hz^{1/2}/W$ or *Jones*) A normalization of spectral detectivity to take into account the area and electrical bandwidth dependence, for example, $D^*(1\ \mu m, 800\ Hz)$ means D^* at $\lambda = 1\ \mu m$ and chopping frequency 800 Hz; unit area and electrical bandwidth are implied. For background-noise-limited detectors the FOV and the background characteristics must be specified. For many types of detectors this normalization is not valid, so that care should be exercised in using D^* .

Spectral detectivity $D(\lambda)$ (W^{-1}) A measure of detector sensitivity, defined as $D(\lambda) = (NEP_{\lambda})^{-1}$. As with NEP, the chopping frequency electrical bandwidth, sensitive area, and, sometimes, background characteristics should be specified.

Spectral noise equivalent power NEP_{λ} The rms value of sinusoidally modulated monochromatic radiant power incident upon a detector which gives rise to an rms signal voltage equal to the rms noise voltage from the detector in a 1-Hz bandwidth. The chopping frequency, electrical bandwidth, detector area, and, sometimes, the background for characteristics should be specified. $NEP(1\ \mu m, 800\ Hz)$ means noise equivalent power at 1- μm wavelength, 1-Hz electrical bandwidth, and 800-Hz chopping rate. Specification of electrical bandwidth is often simplified by expressing NEP in units of $W/Hz^{1/2}$.

Spectral responsivity $R(\lambda)$ The ratio between rms signal output (voltage or current) and the rms value of the monochromatic incident signal power or photon flux. This is usually determined by taking the ratio between a sample detector and a thermocouple detector. The results are given as relative response/watt or relative response/photon, respectively.

Temperature noise Fluctuations in the temperature of the sensitive element, due either to radiative exchange with the background or conductive exchange with a heat sink, produce a fluctuation in signal voltage. For thermal detectors, if the temperature noise is due to the former, the detector is said to be at its theoretical limit. For thermal detectors:

$$\overline{(\Delta T)^2} = \frac{4kT^2G\Delta f}{K^2 + 4\pi^2 f^2 C^2}$$

where $\overline{(\Delta T)^2}$ = mean square temperature fluctuations

K = thermal conductance

C = heat capacity

Thermal noise (also known as Johnson or Nyquist noise) Noise due to the random motion of charge carriers in a resistive element:

$$V_{n,rms} = (4kTR\Delta f)^{1/2} \quad k = \text{Boltzmann's constant}$$

Thermopile A number of thermocouples mounted in series in such a way that their thermojunctions lie adjacent to one another in the plane of irradiation.

Time constant τ (see also *detective time constant*) A measure of the detector's speed of response. $\tau = 1/(2\pi f_c)$, where f_c is that chopping frequency at which the responsivity has fallen to 0.707 ($1/\sqrt{2}$) times its maximum value. Sometimes a physics convention defines it as $1/e$, or 0.368 of the maximum value:

$$R(f) = \frac{R_0}{(1 + 4\pi^2 f^2 \tau^2)^{1/2}}$$

24.5 DETECTOR PERFORMANCE AND SENSITIVITY

D^*

A figure of merit defined by Jones in 1958 is used to compare the sensitivity of detectors.¹⁰ It is called D^* . Although the units of measure are $\text{cm Hz}^{1/2}/\text{W}$, this unit is now referred to as a *Jones*. D^* is the signal-to-noise (S/N) ratio of a detector measured under specified test conditions, referenced to a 1-Hz bandwidth and normalized by the square root of the detector area (A) in square centimeters. Specified test conditions usually consist of the blackbody signal source temperature, often 500 K for infrared detectors, and the signal chopping frequency. If the background temperature is other than room temperature (295 or 300 K in round numbers), then that should be noted.

By normalizing the measured S/N ratio by the square root of the detector area, the D^* figure of merit recognizes that the statistical fluctuations of the background photon flux incident on the detector (photon noise) are dependent upon the square root of the number of photons and thus increase as the square root of the detector area, while the signal will increase in proportion to the detector area itself. This figure of merit therefore provides a valid comparison of detector types that may have been made and tested in different sizes.

The ultimate limit in S/N ratio for any radiation power detector is set by the statistical fluctuation in photon arrival times. For ideal detectors which are photon-noise-limited, and where only generation noise is present, we shall discuss limiting detectivity for three cases:

1. Photon detector where arrival rate of signal photons far exceeds that of background photons (all other noise being negligible)
2. Photon detector where background photon arrival rate exceeds signal photon rate (all other noise being negligible)
3. Thermal detector, background limited

The rate of signal-carrier generation is

$$n = \eta AN_s \tag{1}$$

where η = detector quantum efficiency and AN_s = average rate of arrival of signal photons.

It can be shown¹¹ that in a bandwidth Δf , the rms fluctuation in carrier-generation rate is

$$\delta n_{rms} = (2P_N \Delta f)^{1/2} \tag{2}$$

where P_N is the frequency dependence of the mean square fluctuations in the rate of carrier generation, that is,

$$P_N = A \int_0^\infty \eta(\nu) (\Delta N)^2 d\nu \quad (3)$$

where $(\Delta N)^2$ is the mean square deviation in the total rate of photon arrivals per unit area and frequency interval including signal and background photons. For thermally produced photons of frequency ν (see Ref. 12).

$$(\Delta N)^2 = \bar{N} \frac{e^{h\nu/kT}}{e^{h\nu/kT} - 1} = \frac{2\pi\nu^2}{c^2} \frac{e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} \quad (4)$$

where \bar{N} is the average rate of photon arrivals per unit area and frequency interval. Then, for the special case of $h\nu \gg kT$

$$\delta n_{\text{rms}} = (2A\eta\bar{N}\Delta f)^{1/2} \quad (5)$$

This is also the case for a laser well above threshold. Here the photon statistics become Poisson, and $(\Delta N)^2 = \bar{N}$ even when $h\nu$ is not greater than kT .

Photon Detector, Strong-Signal Case This is generally a good approximation for visible and higher photon energy detectors since the background radiation is often weak or negligible. When signal photons arrive at a much faster rate than background photons

$$\delta n_{\text{rms}} = (2A\eta\bar{N}_s\Delta f)^{1/2} \quad (6)$$

then

$$\text{NEP} = \frac{N_s A h \nu}{n / \delta n_{\text{rms}} (\Delta f)^{1/2}} = \left(\frac{2N_s A}{\eta} \right)^{1/2} h \nu \quad (7)$$

or the noise-equivalent quantum rate is

$$\text{NEQ} = \left(\frac{2 \times \text{incident photon rate}}{\text{quantum efficiency}} \right)^{1/2} \quad (8)$$

Photon Detector, Background-Limited Case This is usually a good approximation for detecting low signal levels in the infrared where background flux levels exceed signal flux levels in many applications. When the background photon noise rate N_B exceeds the signal photon rate ($N_B \gg N_s$)

$$\delta n_{\text{rms}} \approx (2A\eta\bar{N}_B\Delta f)^{1/2} \quad (9)$$

the noise-equivalent power is

$$\text{NEP} = \frac{N_s A h \nu}{(n / \delta n_{\text{rms}}) (\Delta f)^{1/2}} = \left(\frac{2N_B A}{\eta} \right)^{1/2} h \nu \quad (10)$$

The noise-equivalent quantum rate is

$$\text{NEQ} = \left(\frac{2 \times \text{incident background photon rate}}{\text{quantum efficiency}} \right)^{1/2} \quad (11)$$

or

$$D^* = \frac{A^{1/2}}{\text{NEP}} = \left(\frac{\eta}{2N_B} \right)^{1/2} \frac{1}{h\nu} \quad (12)$$

or

$$\text{Area-normalized quantum detectivity} = \frac{A^{1/2}}{\text{NEQ}} = \left(\frac{\eta}{2N_B} \right)^{1/2} \quad (13)$$

For the general case of a detector with area A seeing 2π sr of blackbody background at temperature T , $(\Delta N)^2$ is that in Eq. (4)

$$\overline{(\delta n)^2} = 2A\Delta f \int_0^\infty \eta(\nu) (\Delta N)^2 d\nu = 4\pi A\Delta f \int_0^\infty \eta(\nu) \nu^2 \frac{e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} d\nu \quad (14)$$

Then for $\Delta f = 1$ Hz,

$$\text{NEP} = \frac{h\nu}{c\eta} \left[4\pi A \int_0^\infty \eta(\nu) \frac{\nu^2 e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} d\nu \right]^{1/2} \quad (15)$$

or

$$D^*(T, \lambda) = \frac{c\eta}{2\pi^{1/2}h\nu} \left[\int_0^\infty \eta(\nu) \nu^2 \frac{e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} d\nu \right]^{-1/2} \quad (16)$$

Assuming $\eta(\nu)$ is independent of frequency but falls back to zero for $\nu < \nu_c$

$$D^*(T, \lambda) = \frac{c\eta^{1/2}}{2\pi^{1/2}h\nu} \left[\int_0^\infty \nu^2 \frac{e^{h\nu/kT}}{(e^{h\nu/kT} - 1)^2} d\nu \right]^{-1/2} \quad (17)$$

Figure 13 shows photon-noise-limited D^* versus cutoff wavelength λ_c for various thermal-background temperatures.¹³ Note that these curves are not independent. $D^*(T, \lambda_c)$ is related to $D^*(T, \lambda'_c)$ by the formula

$$D^*(T, \lambda_c) = \left(\frac{T'}{T} \right)^{5/2} D^*(T, \lambda'_c) \quad \text{where } \lambda'_c = \frac{T}{T'} \lambda_c \quad (18)$$

This relation is useful for determining values of $D^*(T, \lambda_c)$, which do not appear in Fig. 13, in terms of a value of $D^*(T', \lambda'_c)$, which does appear. For example, to find $D^*(1000 \text{ K}, 4 \mu\text{m})$ from the 500-K curve

$$D^*(1000, 4) = \left(\frac{500}{1000} \right)^{5/2} D^* \left(500, 4 \times \frac{1000}{500} \right) = 2.3 \times 10^9 \text{ Jones} \quad (19)$$

If higher accuracy is desired than can be determined from Fig. 13, one can use the preceding formula in combination with Table 1, which gives explicit values of $D^*(\lambda_c)$ versus λ_c for $T = 295 \text{ K}$.

The effect on D^* of using a narrow bandwidth detection system is illustrated in Fig. 14. Such a system may be configured with a cold narrow bandwidth filter, or with a narrow bandwidth amplifier—in order to limit the background flux noise or the electrical bandwidth noise, respectively. Q refers to the factor of the reduction provided.

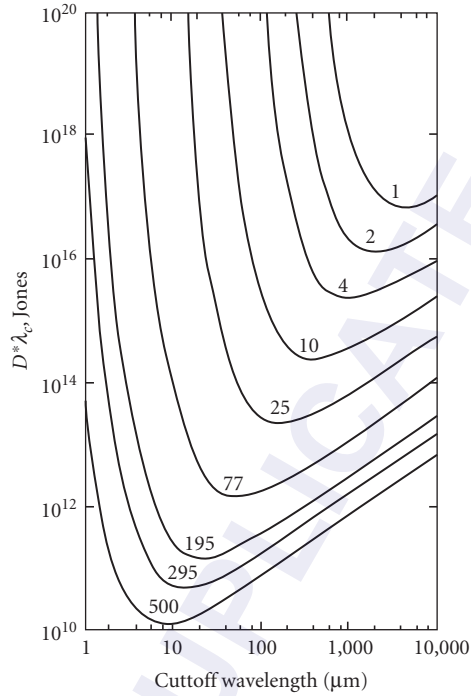


FIGURE 13 Photon-noise-limited D^* at peak wavelength—assumed to be cut-off wavelength—for background temperatures 1, 2, 4, 10, 25, 77, 195, 295, and 500 K (assumes 2π FOV and $\eta = 1$). (Reprinted from Ref. 13.)

TABLE 1 D^* versus λ_c for $T = 295$ K

λ_c μm	$D^*(\lambda_c)$	λ_c μm	$D^*(\lambda_c)$	λ_c μm	$D^*(\lambda_c)$	λ_c μm	$D^*(\lambda_c)$
1	2.19×10^{13}	10	5.35×10^{10}	100	1.67×10^{11}	1000	1.55×10^{12}
2	4.34×10^{13}	20	5.12×10^{10}	200	3.20×10^{11}	2000	3.10×10^{12}
3	1.64×10^{12}	30	6.29×10^{10}	300	4.74×10^{11}	3000	4.64×10^{12}
4	3.75×10^{11}	40	7.68×10^{10}	400	6.28×10^{11}	4000	6.19×10^{12}
5	1.70×10^{11}	50	9.13×10^{10}	500	7.82×10^{11}	5000	7.73×10^{12}
6	1.06×10^{11}	60	1.06×10^{11}	600	9.36×10^{11}	6000	9.28×10^{12}
7	7.93×10^{10}	70	1.21×10^{11}	700	1.09×10^{12}	7000	1.08×10^{13}
8	6.57×10^{10}	80	1.36×10^{11}	800	1.24×10^{12}	8000	1.24×10^{13}
9	5.80×10^{10}	90	1.52×10^{11}	900	1.40×10^{12}	9000	1.39×10^{13}

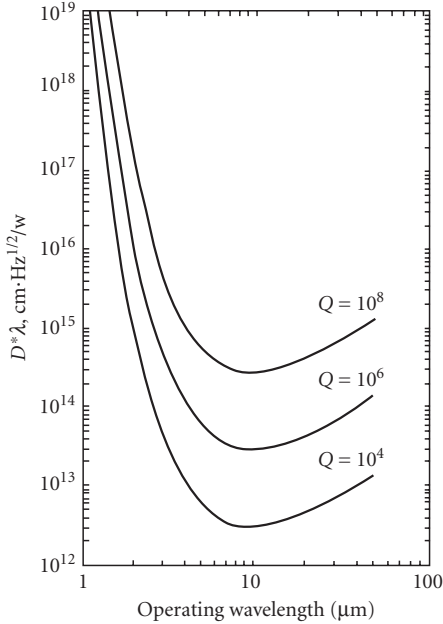


FIGURE 14 Photon noise limit of a narrow-band quantum counter as a function of operating wavelength for a 290-K background, 2π FOV, and $\eta = 1$. (From Ref. 14.)

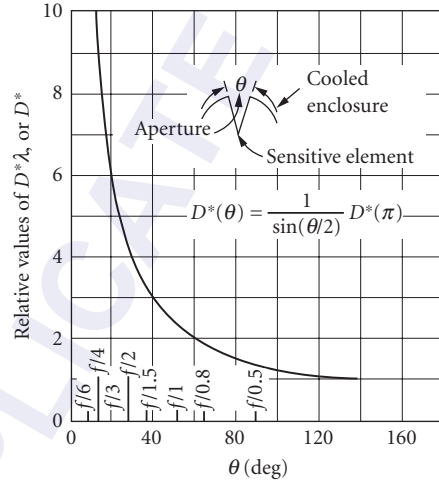


FIGURE 15 Relative increase in photon-noise-limited $D^*(\lambda_{pk})$ or D^* achieved by using a cooled aperture in front of lambertian detector. (From Ref. 14.)

Figure 15 shows the relative increase in photon-noise-limited D^* achievable by limiting the FOV through use of a cooled aperture.

The photon-noise-limited sensitivity shown in Fig. 13 and Table 1 apply to photovoltaic and photoemissive detectors. Figure 14 is for photoconductors. For photoconductors, recombination noise results in a $\sqrt{2}$ reduction in D^* at all wavelengths.

Thermal Detectors Limiting sensitivity of an ideal thermal detector has been discussed previously.^{12,14,15} Assuming no shortwave- or long-wavelength cutoffs exist,

$$D^* = \frac{\varepsilon}{[8\varepsilon\sigma k(T_1^5 + T_2^5)]^{1/2}} = \frac{4 \times 10^{16} \varepsilon^{1/2}}{(T_1^5 + T_2^5)^{1/2}} \text{ Jones} \quad (20)$$

where T_1 = detector temperature
 T_2 = background temperature
 ε = detector emissivity
 σ = Stefan-Boltzmann constant
 k = Boltzmann constant

D^* versus T_2 is plotted for various T_1 in Fig. 16. Figure 17 shows the effect of both short- and long-wavelength cutoffs on bolometer sensitivity,¹⁶ with the ideal photoconductor curve for reference. D^* can be seen to increase rapidly when the cutoff is set to avoid the high flux density from the 300-K background that peaks around 10 μm .

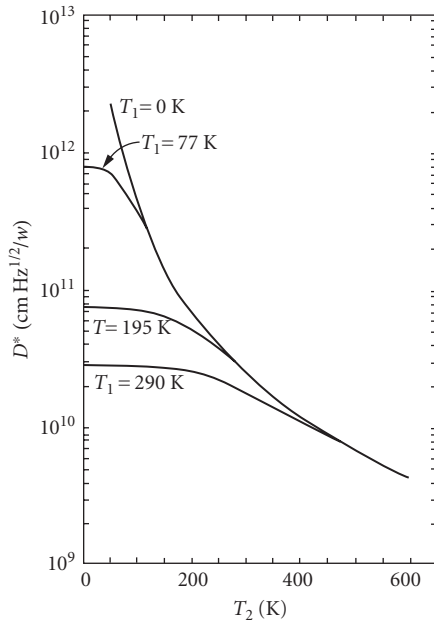


FIGURE 16 Photon-noise-limited D^* for thermal detectors as a function of detector temperature T_1 and background temperature T_2 (2π FOV: $\eta = 1$). (From Ref. 14.)

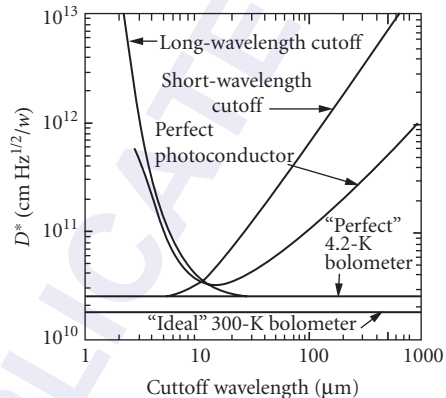


FIGURE 17 The detectivity of a “perfect” bolometer plotted as a function of both short- and long-wavelength cutoffs. Plots for a perfect photoconductor and two other cases are included for comparison. The background temperature is 300 K. (From Ref. 16.)

24.6 OTHER PERFORMANCE PARAMETERS

Spectral Response

Spectral response provides key information regarding how the detector will respond as a function of wavelength or photon energy. Spectral response may be limited by the intrinsic detector material properties, a coating on the detector, or by a window through which the radiation must pass. Relative response is the spectral response ratioed against a detector with a nominally wavelength independent response, such as a thermocouple having a spectrally-broad black coating. Relative response is plotted as a function of wavelength with either a vertical scale of W^{-1} or photon^{-1} . Thermal detectors tend to be spectrally flat in the first case while quantum detectors are generally flat in the second case. The curves are typically shown with the peak value of the spectral response normalized to a value of 1. The spectral response curve can be used together with the blackbody D^* to calculate D^* as a function of wavelength, which is shown in Fig. 18 for selected detectors.

Responsivity and Quantum Efficiency

Responsivity and quantum efficiency are important figures of merit relating to the detector signal output. Responsivity is a measure of the transfer function between the input signal photon power or flux and the detector electrical signal output. Thermal detectors will typically give this responsivity in volts/watt. Photoconductors will usually quote the same units, but will also frequently reference the value to the peak value of relative response per watt from the spectral response curve. This value

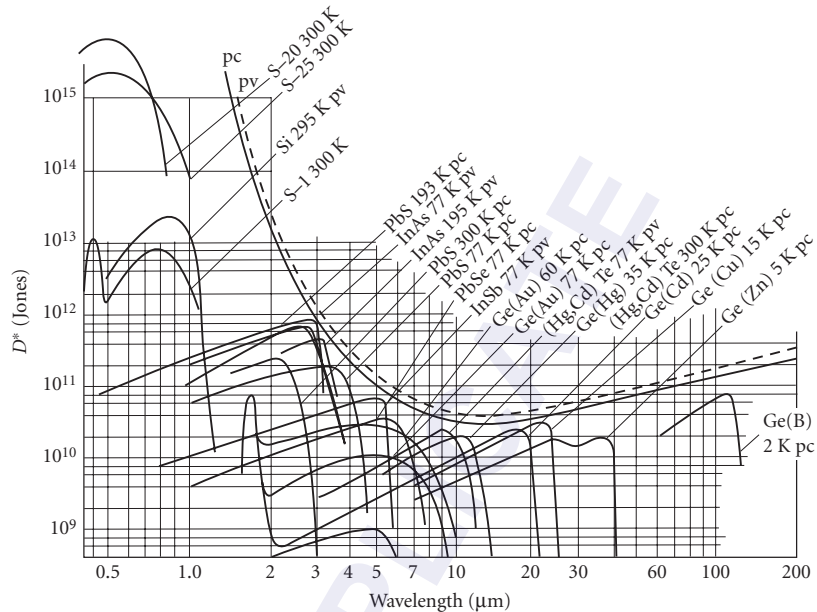


FIGURE 18 D^* versus λ for selected detectors.

is actually realized if the detector bias circuit load resistor is significantly larger than the detector resistance. Photoconductor responsivity is given by:

$$\text{Responsivity}_{\text{peak}} = \frac{\eta q R_d E \tau (\mu_n + \mu_p)}{h \nu \ell} \quad (21)$$

where η is the quantum efficiency, q is the electronic charge, R_d is the detector resistance, E is the electric field, τ is time constant, μ are the mobilities for electrons (n) and holes (p), $h\nu$ is the photon energy, and ℓ is the device length. Photomultiplier tubes and photovoltaic detectors will usually reference the responsivity in amperes per watt, again referenced to peak spectral response.

Detector response performance is also conveyed from the detector quantum efficiency. In the case of photovoltaic detectors which, in the absence of avalanche operation have a gain of unity, quantum efficiency is essentially the current per photon. For a blip photovoltaic detector, the quantum efficiency also determines the D^* . Quantum efficiency is not readily measured for photoconductors and photomultiplier tubes unless the internal gain is carefully calibrated. It is sometimes inferred from the measured D^* for photoconductive devices which are blip—see definition of detective quantum efficiency.

Noise, Impedance, Dark and Leakage Current

Noise has a number of potential origins. Background photon flux-limited detectors have noise dominated by the square root of the number of background photons striking the detector per second [see Eq. (9)]. Other noise sources may contribute or dominate. Among these are

- Johnson, Nyquist, or thermal noise which is defined by the detector temperature and impedance
- Modulation or $1/f$ noise which may dominate at lower frequencies
- Amplifier noise
- Shot noise from dark or leakage current

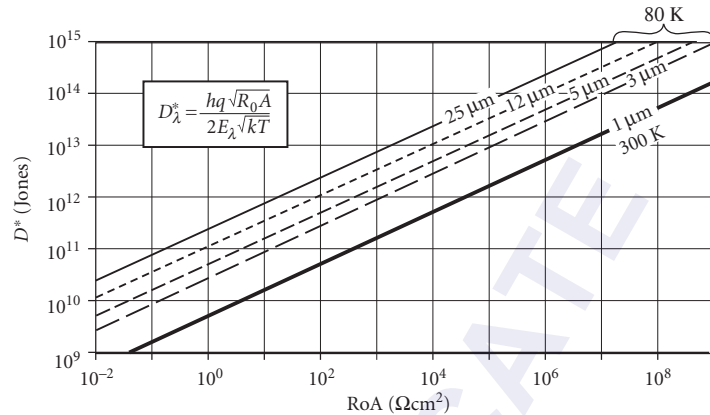


FIGURE 19 Zero-bias impedance-area product (R_0A or shunt resistance per unit area) of a photodiode can limit the D^* as shown. The limiting D^* depends on R_0A , temperature, and photon energy or wavelength. Examples are illustrated for a 1- μm cutoff diode at 300 K, and for 3-, 5-, 12-, and 25- μm cutoff devices at 80 K. Other noise mechanisms, such as photon noise, typically limit D^* to lower values than the highest values shown here.

The impedance of a photodiode may limit performance, depending upon the detector operating conditions. Figure 19 illustrates the diode impedance per unit area (R_0A) limiting value of D^* for silicon detectors at room temperature and longer-wavelength infrared photodiodes at 80 K.

Measurement of the noise as a function of frequency can be valuable for characterizing the relevant noise sources. Selection of an appropriate preamplifier is also critical, particularly for detectors having very low or very high impedance. Integration of preamplifiers together with detectors has significantly improved the overall performance of many detectors. The use of phase-sensitive lock-in amplifiers in combination with a modulated signal can also improve the signal-to-noise ratio.

Uniformity

One cannot assume that the response of a detector will be uniform across its sensitive area. Material inhomogeneity and defects and/or fabrication variables can give rise to nonuniformity. Lateral collection from near the perimeter of a photodiode may give a gradual response decrease away from the edge—this effect will typically be accompanied by a change in response speed as well. Recombination at the electrical contacts to a photoconductor can limit the lifetime, and hence the photoconductive gain, for carriers generated near the contact, a phenomenon called sweep-out. Recombination may be enhanced at surfaces and edges also. Laser spot scanning is useful to check the detector spatial uniformity, although laser sources may not be readily available at all the wavelengths of interest. An alternative method is to move the detector around under a fixed small aperture in conjunction with a light source.

Speed

Detector response speed is often related inversely to detector sensitivity. Thermal detectors often show this characteristic because the signal is proportional to the inverse of response speed, while the noise is amplifier or Johnson limited. Excluding detectors with internal carrier multiplication mechanisms, the best detectors from broad experience seem limited to a D^*f^* product of a few times 10^{17} Jones Hz. D^*f^* may be proportionally higher for devices with gain, since speed can be increased to a greater extent by using a lower value of load resistance without becoming Johnson-noise-limited. The user should be aware that with many detectors it is possible to operate them in a circuit to maximize

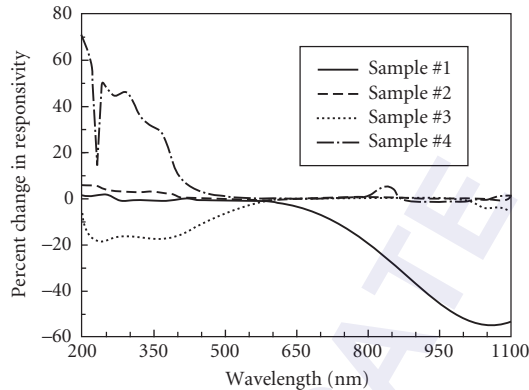


FIGURE 20 After only 3 hours of UV irradiation, these silicon detectors showed great variations in responsivity. (Reprinted from the September 1993 issue of *Photonics Spectra*, © Laurin Publishing Co., Inc.)

sensitivity or speed, but not both at the same time. Speed may vary across the sensitive area of the detector and with temperature, wavelength, and electrical bias.

Stability

Detector performance may change or drift with time. Changes in operating temperature, humidity, and exposure to elevated temperatures as well as to visible, ultraviolet, and high-energy radiation can affect device operation. These effects arise from the temperature dependence of electronic properties in solids, as well as from the critical role played by electrical charge conditions near the surface of many device types. Sensitivity changes in a sample of silicon detectors from four vendors illustrate this point. Wide variations in responsivity change after UV exposure, as shown in Fig. 20. In applications where stability is of significant concern, these effects must be carefully reviewed along with the detector supplier's experience in these matters.

24.7 DETECTOR PERFORMANCE

Manufacturers' Specifications

Table 2 lists the detector materials covered in Sec. 24.7.

TABLE 2 Detector Materials Covered in Sec. 24.7

Thermocouple	GaAsP	Ge:Cu
Thermopile	CdS, CdSe	HgCdTe
Thermistor bolometer	CdTe	PbSnTe
Pyroelectric	GaAs	Ge:Hg
InSb hot electron bolometer	Si	Si:Ga
Ge:Ga bolometer	InGaAs	Si:B
Photoemissive	Ge	Ge:Cu
GaN, InGaN	PbS	Ge:Zn
SiC	InAs	Ge:Ga
TiO ₂	PbSe	Photographic
GaP	InSb	

Detector sensitivity can be the determining factor in the design and performance of a sensor system. Detector performance is subject to the development of improved materials, fabrication techniques, and the ingenuity of device engineers and inventors. The descriptions given here may improve with time, and consultation with manufacturers and users is recommended. Today, the internet can be the quickest and most up-to-date source of currently available manufacturers and specifications for the devices they offer. Many suppliers noted in this section may have gone out of business—a search on the internet is the best choice for finding active vendors. Other than a general Web search, some collections of device suppliers can be found at

<http://www.photonics.com/bgHome.aspx>

<http://laserfocusworld.365media.com/laserfocusworld/search.asp>

<http://www.physicstoday.org/ptbg/search.jsp>

Thermocouple The thermocouple offers broad uniform spectral response, a high degree of stability, and moderate sensitivity. Its slow response and relative fragility have limited its use to laboratory instruments, such as spectrometers.

Compared with thermistors, thermocouples are slower, require no bias, and have higher stability but much lower impedance and responsivity. This increases the amplification required for the thermocouple; however, the only voltage appearing is the signal voltage, so that the serious thermistor problem of bridge-circuit bias fluctuations is avoided. With proper design, performance should not be amplifier-limited but limited instead by the Johnson noise of the thermocouple. Thermocouples perform stably in dc operation, although the instability of dc amplifiers usually favors ac operation.

The inherent dc stability of thermocouples is attractive for applications requiring no moving parts, and recently a relatively rugged solid-backed evaporated thermocouple has been developed whose sensitivity approaches that of the thermistor bolometer.

Sensitivity: $D^* 1 \times 10^9$ Jones for 20-ms response time; spectral response depends on black coating (usually gold black) (see Fig. 21)

Noise: White Johnson noise, falling off with responsivity (see Fig. 22)

Resistance: 5 to 15 Ω typical

Responsivity: 5 V/W (typical), 20 to 25 V/W (selected)

Time constant: 10 to 20 ms (typical)

Operating temperature: Normally ambient

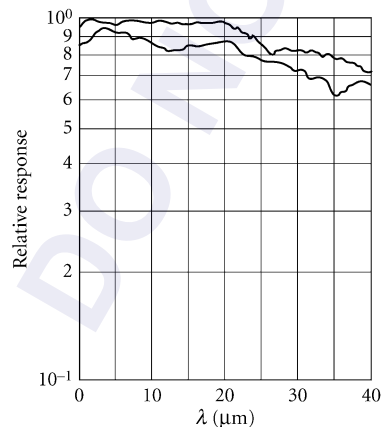


FIGURE 21 Typical thermocouple spectral response curves (CsI window) for two different manufactures.

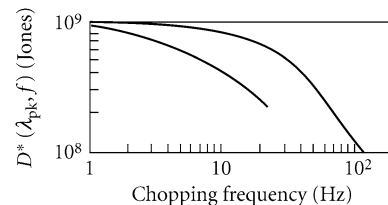


FIGURE 22 Typical thermocouple D^* (noise) frequency response for two manufactures. (From Ref. 1.)

Sensitive area: 0.1×1 to 0.3×3 or 0.6×2 mm (typical)

Linearity: 0.1 percent in region investigated (6×10^{10} to 6×10^8 W incident)

Recommended circuit: Transformer coupled into low-noise (bipolar or JFET) amplifier with good low-frequency noise characteristics

Manufacturers: Perkin-Elmer, Charles Reeder, Beckman Instruments, Farrand, Eppley Laboratory

Thermopile Thermopiles are made by evaporating an array of metal junctions, such as chromel-constantan or manganin-constantan, onto a substrate. The thin-film construction is rugged, but the Coblenz-type may be quite delicate. Wire-wound thermopile arrays are also available which are very robust. Devices with arrays of semiconductor silicon junctions are also available. The array may typically be round, square, or rectangular (for matching a spectrometer slit) and consist of 10 to 100 junctions. Configuration options include matched pairs of junction arrays or compensated arrays to provide an unilluminated reference element. A black coating, such as 3M black or lampblack, is used to provide high absorption over a broad spectral range, as illustrated in Figs. 23 and 24. The housing window may limit the spectral range of sensitivity. Typical applications include power meters and radiometers.

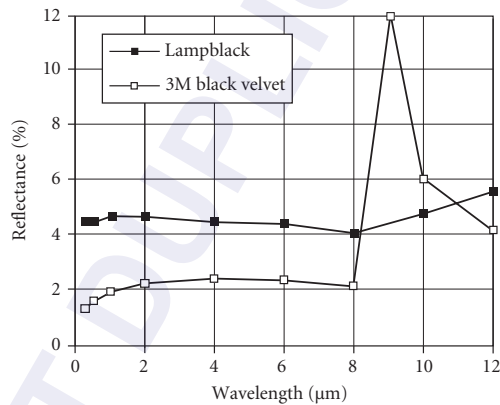


FIGURE 23 Spectral reflectance of two black coatings used in the construction of thermopile detectors. (From Eppley Laboratory studies.)

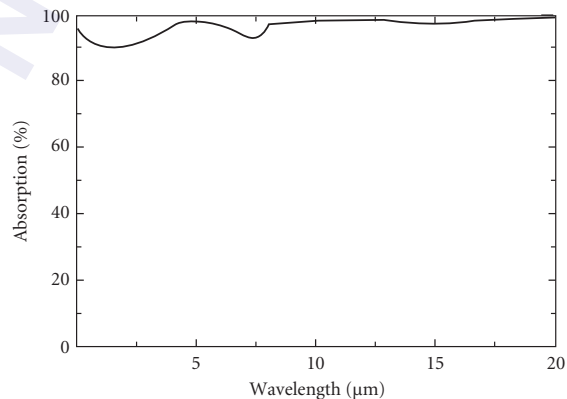


FIGURE 24 Spectral absorption of a thermopile detector coating made from a black metal oxide. (Oriol Corporation.)

Sensitivity: D^* 0.5 to 4×10^8 Jones for 30-ms typical response time. D^* may be dependent upon sensitive area; spectral response depends on black coating and window (see Fig. 23 and Fig. 24)

Noise: White Johnson noise, falling off with responsivity, typical range is 5 to 30 nV/Hz^{1/2}

Resistance: 2Ω to 60 kΩ typical

Responsivity: 4 to 250 V/W (typical) depends on the number of junctions and time constant

Time constant: 10 ms to 2 s (typical)

Operating temperature: Normally ambient

Sensitive area: 0.5 to 6-mm diameter, 0.025×0.025 to 3×3 mm, various rectangular, 0.4×3 , 0.6×2 , 0.6×4 mm (typical)

Recommended circuit: Low noise ($0.5 V_{p-p}$, dc to 1 Hz), low drift with voltage gain of 1000 and input impedance of 1 MΩ

Manufacturers: Armtech, Beckman Instruments, Concept Engineering, Dexter Research Center, Edinburgh Instruments, Eppley Laboratory, Farrand, Gentech, Molectron Detector, Ophir Optronics, Oriel, Scientech, Scitec, Swan Associates

Thermistor Bolometer Thermistors offer reliability, moderate sensitivity, and broad spectral coverage without cooling. Construction is rugged and highly resistant to vibration, shock, and other extreme environments. Response is slower than 1 ms, and trade-off exists between speed and sensitivity.

Thermistor elements are made of polycrystalline Mn, Ni, and Co oxides. In their final form they are semiconductor flakes 10 μm thick, which undergo a temperature resistance change of ~4 percent per Kelvin. Since thermistor resistance changes with ambient temperature enough to alter the biasing significantly, it is usually operated in a bridge circuit, with a nearly identical thermistor shielded from signal radiation and used for a balance resistor.

Sensitivity:

$$\text{NEP} = 8.9 \times 10^{-10} \sqrt{\frac{A(\text{mm}^2)}{\tau_{\text{rms}}}} \text{ W}$$

$$D^* = 1.1 \times 10^9 \sqrt{\tau_{\text{rms}}} \text{ Jones}$$

Spectral response: Depends on coating (usually Zapon lacquer); see Fig. 25.

Quantum efficiency: Depends on blackening coating, typically 80 percent.

Noise: Thermal-noise-limited above 20 Hz ($V_{\text{noise}} = \sqrt{4kTR\Delta f}$); below that, 1/f type noise—see Fig. 26. Used in balanced-bridge circuit (two flakes in parallel); limiting noise due to thermal noise in both flakes.

Resistance: For standard 10-μm-thick flakes, two different resistivities are available: 2.5 MΩ/sq or 250 kΩ/sq. Note that in a bolometer bridge, the resistance between the output connection and ground is half that of single flake.

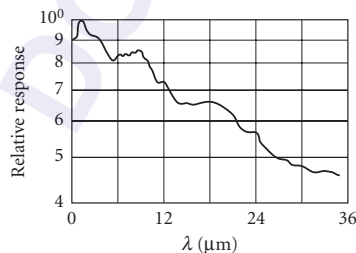


FIGURE 25 Typical thermistor spectral response (no window).

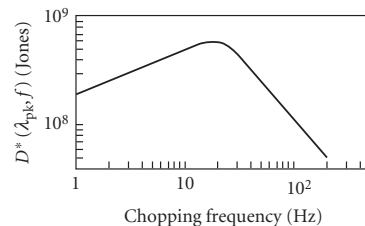


FIGURE 26 Typical thermistor D^* (noise) frequency spectrum. (From Ref. 1.)

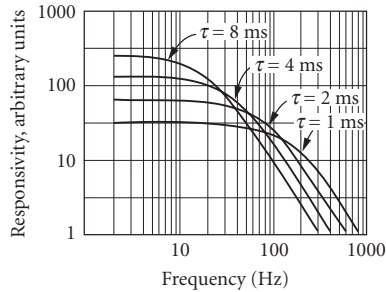


FIGURE 27 Relative response vs. frequency for various time-constant thermistor detectors ($A = 1 \times 1$ mm). (From Barnes Engineering, Bull. 2-100.)

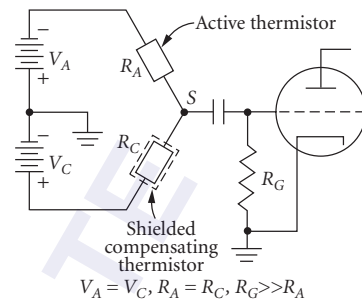


FIGURE 28 Bolometer electrical circuit. (From Barnes Engineering, Bull. 2-100.)

Time constant: τ is 1- to 20-ms standard for nonimmersed detectors and 2- to 10-ms standard for immersed detectors.

Sensitive area: 0.1×0.1 mm, 5×5 mm standard.

Operating temperature: Normally ambient, 285–370 K.

Responsivity: Depends on bias, resistance, area, and time constant $\mathcal{R} \propto \sqrt{R\tau/A} \approx 10^3$ V/W for 0.1×0.1 mm area, 250 k Ω resistance, and $\tau = 4$ ms (see Fig. 27 for frequency-time-constant dependence with given area). Output voltage (responsivity) can be increased to a limited degree by raising bias voltage. Figure 29 shows the deviation from Ohm's law due to heating. Bias should be held below 60 percent of peak voltage. Listed responsivity is that of active flake. In the bridge circuit, responsivity is half this value.

Sensitivity profile: Approximately 10 percent for 10- μ m scan diameter over a 1×1 mm cell.

Linearity: ± 5 percent 10^{-6} to 10^{-1} W/cm 2 .

Recommended circuit: See Figs. 28 and 29.

Manufacturers: Servo Corporation of America, Thermometrics.

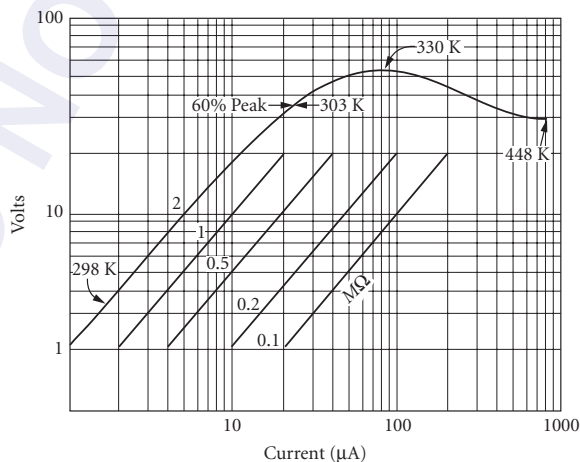


FIGURE 29 Thermistor voltage-current characteristics, showing typical flake temperatures under different conditions. (From Barnes Engineering, Bull. 2-100.)

Pyroelectric Lithium tantalate (LiTaO₃), triglycine sulfate (TGS), and other pyroelectric materials provide an uncooled thermal detector with good sensitivity. The devices are capable of fast response, limited inversely by the preamplifier feedback resistance, but with responsivity and D^* traded for speed. This detector's principle of operation is the pyroelectric effect, which is the change of electric polarization with temperature. Pyroelectric detectors offer rugged construction and the absence of $1/f$ noise because no bias is involved.

Lack of $1/f$ noise, combined with the ability to easily trade off speed and sensitivity, makes pyroelectric detectors useful for scanning applications and energy measurement of pulsed optical sources. In addition, the NEP is independent of area at low frequencies (10 Hz), so that these detectors are useful for large-area applications (preamplifier $1/f$ noise may limit, however). Pyroelectric detectors are useful for calorimetry since the pyroelectric effect is an integrated volume effect and the output signal is unaffected by spatial or temporal distribution of the radiation, up to damage threshold or depolarizing temperature. For higher damage thresholds, lead zirconate titanate ceramic (Clevite PZT-5) exhibits a much smaller pyroelectric effect than TGS, but its high Curie temperature of 638 K makes it more useful than TGS for high-energy applications.

Sensitivity: Sensitive from ultraviolet to millimeter wavelengths. For $\lambda < 2 \mu\text{m}$, TGS must be blackened, which slows response. Normally ($\lambda > 2 \mu\text{m}$) a transparent electrode is used, since TGS absorption is high from 2 to 300 μm . Beyond 300 μm , poor absorption and increased reflectivity reduce sensitivity. Spectral response depends largely on coating. See Fig. 30 for spectral response with modified 3M black. Figure 31 illustrates the relative spectral response for a LiTaO₃ device.

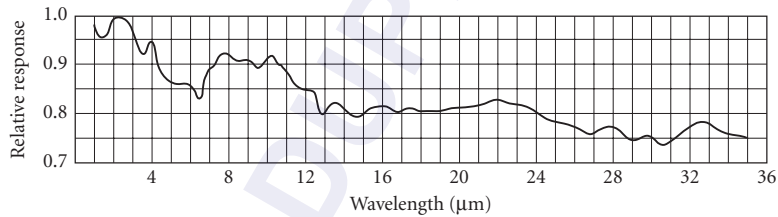


FIGURE 30 Relative spectral response of TGS detectors with modified 3M black. (From Barnes Engineering, Bull. 2-100.)

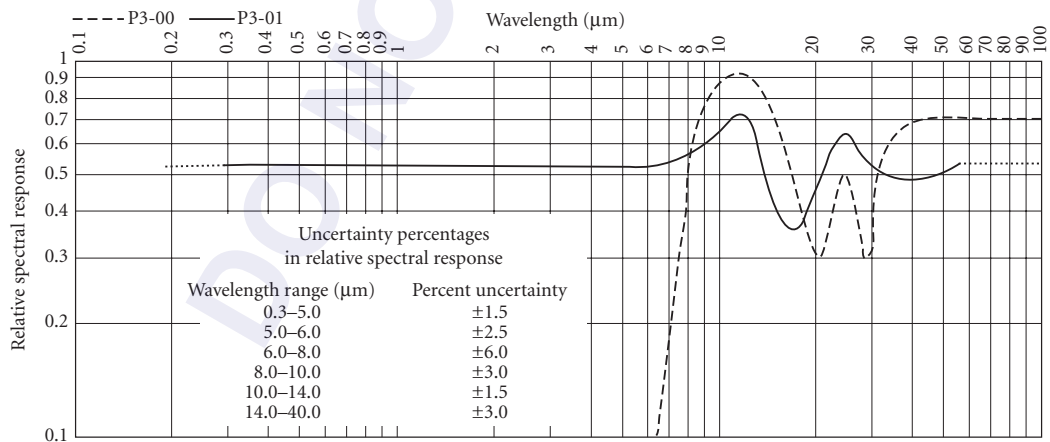


FIGURE 31 Relative spectral response of LiTaO₃ pyroelectric detectors showing both a black spectral coating and an optional coating tuned to the 8 to 14- μm LWIR band. (From Moleclectron Detector, Inc.)

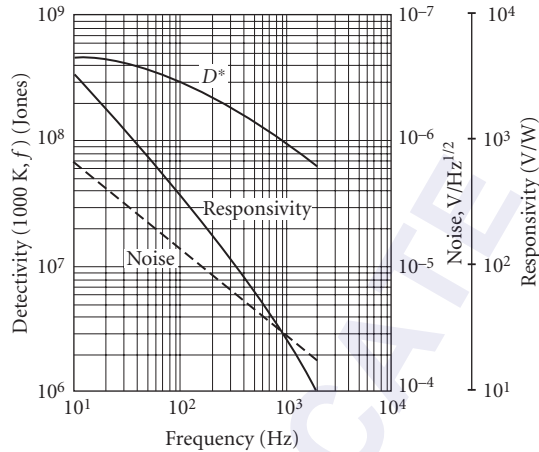


FIGURE 32 Typical D^* , responsivity, and noise versus frequency for TGS ($A = 1 \times 1$ mm; $T = 296$ K). (From Barnes Engineering, Bull. 2-220A.)

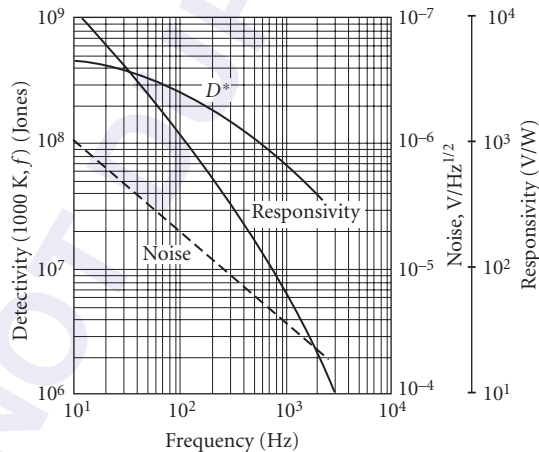


FIGURE 33 Typical D^* , responsivity, and noise versus frequency for TGS ($A = 1 \times 1$ mm; $T = 296$ K). (From Barnes Engineering, Bull. 2-220A.)

D^* is independent of A at low frequencies (10 Hz) (see Figs. 32 and 33). Figure 34 shows NEP versus A for various frequencies.

Quantum efficiency: Depends on coating absorptivity (for 3M black typically $\eta > 75$ percent).

Noise: (See Figs. 32 and 33). Limited by loss-tangent noise up to frequencies that become limited by amplifier short-circuit noise (see Fig. 35).

Operating temperature: Ambient, up to 315 K. Can be repolarized if $T > T_{\text{curie}} = 322$ K for TGS. Irreversible damage at $T = 394$ K (see Fig. 36). Other pyroelectric materials have significantly higher Curie temperatures (398 to 883 K).

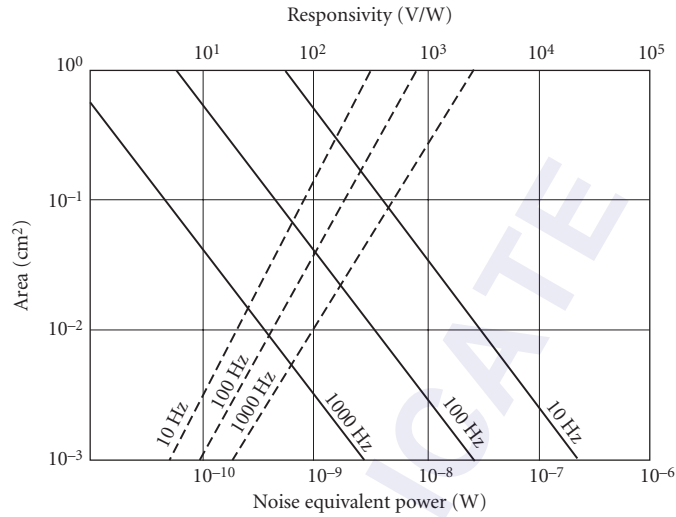


FIGURE 34 Noise equivalent power in watts (broken lines) and responsivity (solid lines) versus TGS detector area for various frequencies. (From Barnes Engineering, Bull. 2-220B.)

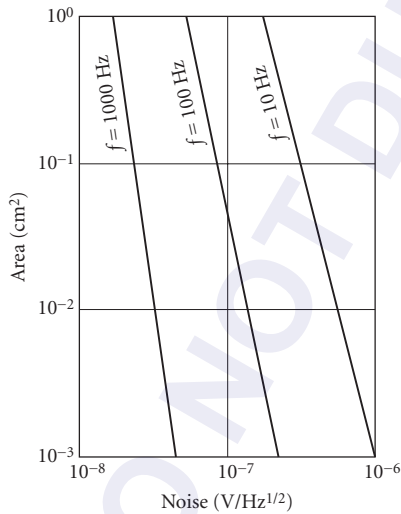


FIGURE 35 TGS noise versus detector area for various operating frequencies. (From Barnes Engineering, Bull. 2-100.)

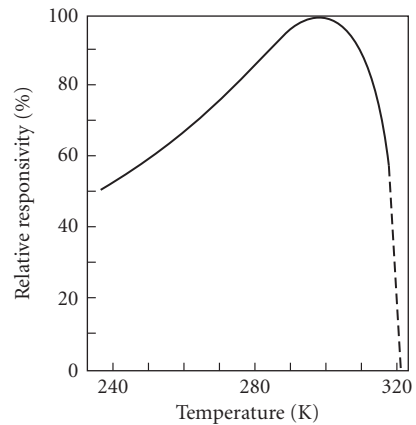


FIGURE 36 Relative responsivity versus temperature for TGS. (From Barnes Engineering, Bull. 2-220B.)

Output impedance: 50 Ω to 10 KΩ, set by built-in amplifier (see Fig. 37).

Responsivity: See Figs. 32 to 34, 36, and 38.

Capacitance: 5 pF for 0.5 × 0.5 mm; 20 pF for 1 × 1 mm; 100 pF for 5 × 5 mm.

Sensitive area: 2 to 50-mm diameter round, 0.5 × 0.5-mm to 10 × 10-mm square, typical.

Time constant: Not pertinent, response speed set by the preamplifier feedback resistor (see Fig. 38).

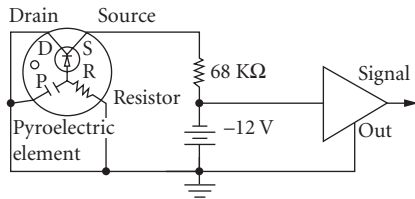


FIGURE 37 Pyroelectric detector amplifier circuit. (From Barnes Engineering, Bull. 2-220A.)

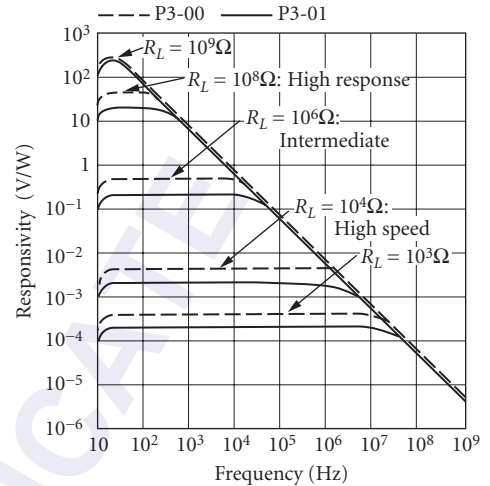


FIGURE 38 Responsivity vs. frequency for two pyroelectric models using various feedback resistors. (From Molectron Detector, Inc.)

Linearity: 5 percent between 10^{-6} and 10^{-1} W/cm².

Sensitivity profile: Depends on coating or transparent electrode; 5 to 7 percent across 12×12 -mm; spot size $< 250 \mu\text{m}$.

Recommended circuit: See Figs. 37 and 39. Field effect transistor (FET) amplification stage usually built in. Since $\mathcal{R} \propto 1/f$, use of an amplifier with $1/f$ noise and gain $\propto f$ is recommended. Then output signal and signal-to-noise ratio are independent of frequency.

Manufacturers: Alrad Instruments, Belov Technology, CSK Optronics, Delta Developments, EG&G Heimann, Electro-Optical Systems, Eltec, Gentec, Graseby, International Light, Laser Precision, Molectron Detector, Oriel, Phillips Infrared Defence Components, Sensor-Physics, Servo Corporation of America, Spiricon, Thermometrics.

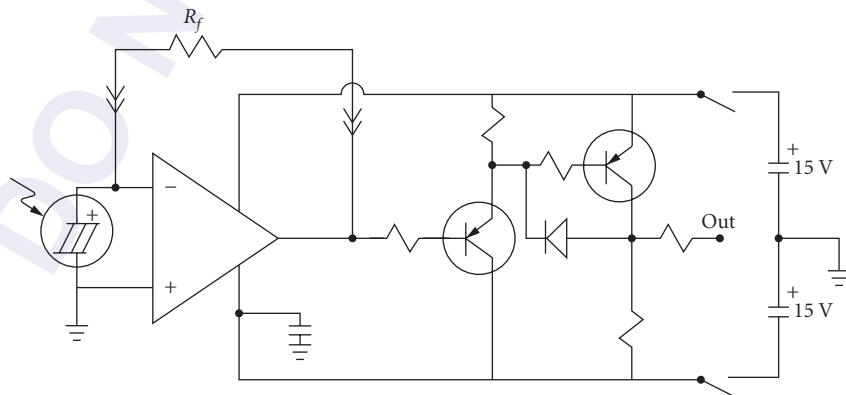


FIGURE 39 Resistive feedback circuit used with LiTaO₃ detectors. (From Molectron Detector, Inc.)

InSb Hot-Electron Bolometer At temperatures of liquid helium and lower, free carriers in indium antimonide (InSb) can absorb radiation in the far-infrared and submillimeter spectral region. Because the mobility of the electrons varies as $T_e^{3/2}$ under these conditions, the conductivity of the material is modulated. This mechanism offers submicrosecond response and broad far-infrared coverage out to millimeter wavelengths but requires liquid-helium cooling and very sophisticated receiver design.

Technically, these devices may be classed as bolometers, since incident radiation power produces a heating effect which causes a change in free-charge mobility. In the normal bolometer, the crystal lattice absorbs energy and transfers it to the free carriers via collisions. However, in InSb bolometers incident radiation power is absorbed directly by free carriers, the crystal lattice temperature remaining essentially constant. Hence the name electron bolometer. Note that this mechanism differs from photoconductivity in that free-electron mobility rather than electron number is altered by incident light (hence there is no photoconductive gain).

Sensitivity: $D^*(2 \text{ mm}, 900) = 4 \times 10^{11}$ Jones (see Fig. 40).

Noise: See Figs. 41 and 42.

Responsivity: 1000 V/W.

Time constant: 250 ns.

Sensitive area: 5×5 mm typical.

Operating temperature: 1.5 to 4.2 K.

Impedance: Without bias, 200Ω ; optimum bias, 150Ω , depends on bias (see Fig. 29).

Recommended circuit: Optimum bias 0.5 mA (see Fig. 43).

Manufacturer: Infrared Laboratories.

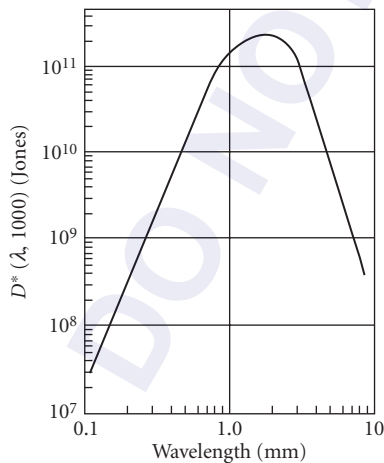


FIGURE 40 D^* versus λ for InSb electron bolometer ($H = 0$). (From Raytheon, *IR Millimeter Wave Detector*, 1967.)

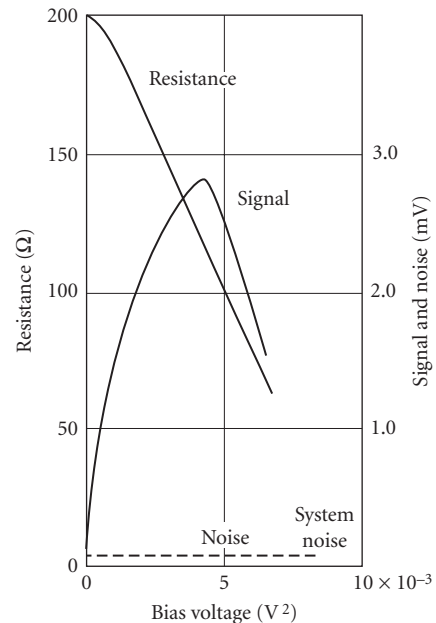


FIGURE 41 InSb electron bolometer, typical resistance, signal, and noise versus bias voltage squared ($T = 5 \text{ K}$; $R_L = 200 \Omega$; gain $= 2.4 \times 10^4$; $F = 1100 \text{ Hz}$). (From Santa Barbara Research Center, *Prelim. Res. Rep.*, 1967.)

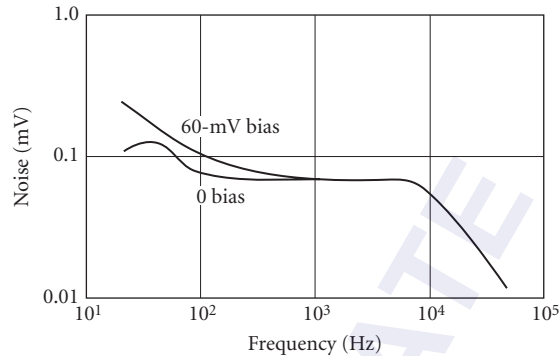


FIGURE 42 InSb electron bolometer, typical noise spectrum ($T = 5$ K; $R_L = 200 \Omega$; gain = 2.4×10^4 ; $\Delta f = 5.6$ Hz). (From Santa Barbara Research Center, Prelim. Res. Rep., 1967.)

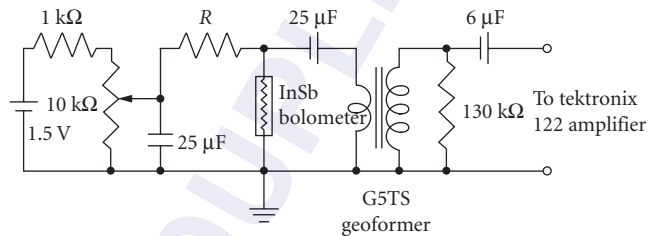


FIGURE 43 Electron bolometer biasing circuit. (From Santa Barbara Research Center, Prelim. Res. Rep., 1967.)

Ge(Ga) Low-temperature Bolometer The Ge(Ga) bolometer offers very high sensitivity and broad spectral coverage in the region 1.7 to 2000 μm . Liquid-helium cooling is required. A trade-off exists between response time (seconds) and sensitivity (10^{-14} -W NEP). Operation at 1000 Hz can be achieved still maintaining 2×10^{-13} W NEP.

Sensitivity: Depends on thermal conductance G (see Figs. 44 and 45). $\text{NEP}(\lambda, 10 \text{ Hz}) = V_n/S = 3 \times 10^{-14}$ W for $G = 1 \mu\text{W/K}$; $A = 1 \text{ mm}^2$, $D^*(\lambda, 10 \text{ Hz}) = 3 \times 10^{13}$ Jones ($Q < 0.2 \mu\text{W}$). For this detector, $\text{NEP} \approx 4T(kG)^{1/2}$ and does not vary with $A^{1/2}$ (T is heat-sink temperature, and k is Boltzmann's constant). Thus D^* cannot be used as a valid means of comparison with other detectors; 300-K background-limited performance is achievable for 2π FOV when the bolometer is operated at 4.25 K with $G = 10^{-3}$ W/K. (For $A = 0.1 \text{ cm}^2$, the time constant is 50 s.)

Responsivity: Typically, responsivity = 2.5×10^5 V/W = $0.7(R/TG)^{1/2}$, where R = resistance. Responsivity, and hence NEP, depends on thermal conductance G , which in turn is set by background power. G ranges from 0.4 to 1000 $\mu\text{W/K}$.

Thermal conductance: Typically $G = 1 \mu\text{W/K}$ for background $Q < 0.2 \mu\text{W}$; note that $Q < 1/2$ (optimum bias power P).

Sensitive area: 0.25 \times 0.25 to 10 \times 10 mm.

Resistance: 0.5 M Ω .

Operating temperature: 2 K (see Fig. 44). In applications where radiation noise can be eliminated there is much to be gained by operating at the lowest possible temperature. Figure 45 shows the theoretical NEP and time constant at 0.5 K, assuming that current noise remains unimportant.

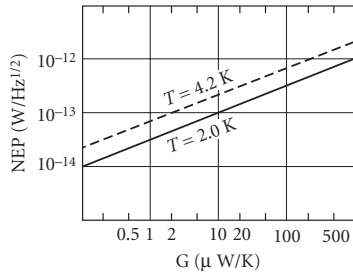


FIGURE 44 Germanium bolometer NEP versus conductance. (*Infrared Laboratories, Inc.*)

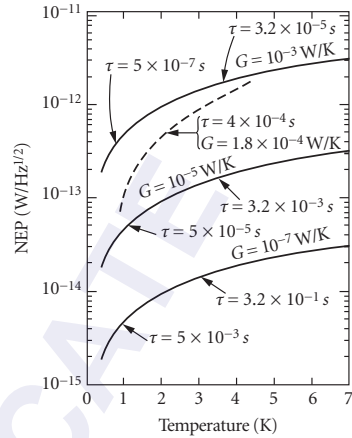


FIGURE 45 Germanium bolometer NEP versus temperature. Solid curves are theoretical: $NEP \approx 4T(kG)^{1/2}$. (*From Ref. 17.*)

Quantum efficiency: Depends on blackened coating and window. For $\lambda < 100 \mu\text{m}$, absorptivity exceeds 95 percent. For $\lambda > 100 \mu\text{m}$, efficiency varies with geometry.

Time constant: Response time constant is proportional to G^{-1} . Therefore, if G must be increased to accommodate larger background, the time constant is decreased proportionally. Responsivity and NEP, however, are degraded as $G^{-1/2}$.

Noise: $V_n = 1 \times 10^{-8} \text{ V/Hz}^{1/2}$; thermal noise is due to R and R_L .

Recommended circuit: Standard photoconductive circuit, with load resistor, grid resistor, and blocking capacitor at low temperature (see Fig. 46). See Fig. 47 for typical electrical characteristics. Bias power $P = 0.1 \text{ TG}$.

Manufacturer: Infrared Laboratories, Inc.

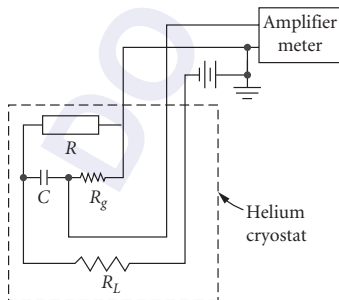


FIGURE 46 Germanium bolometer circuit and cryostat.

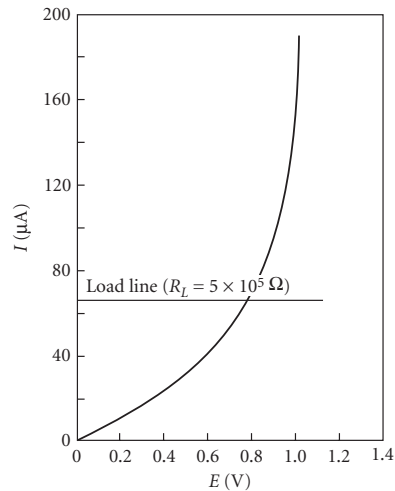


FIGURE 47 Load curve for a typical germanium bolometer ($T = 2.15 \text{ K}$) with load line, showing optimum operating point. (*From Ref. 17.*)

Photoemissive Detectors Photoemissive detector is generally the detector of choice in the UV, visible, and near-IR where high quantum efficiency is available. In the spectral region $\lambda < 600$ nm, the photomultiplier, or multiplier phototube, has close to ideal sensitivity; that is, selected photomultiplier tubes (PMT) are capable of detecting single photon arrivals (but at best only with about 30 percent quantum efficiency) and amplifying the photocurrent (pulse) enormously without seriously degrading the signal-to-noise ratio. Time resolution can be as short as 0.1 ns. Only very specialized limitations have precluded their use for $\lambda < 800$ nm, for example, cost, ruggedness, uniformity of manufacture, or need for still faster response. Recently these limitations have all been met individually but generally not collectively. Where adequate light is available, the simple phototube has advantages over the multiplier phototube in that high voltages are not required, the output level is not sensitive to applied voltage, and dynode fatigue is eliminated.

Microchannel plate tubes (MCPT) are a variant of the photomultiplier tube where the current amplifying dynode structure is replaced by an array of miniature tubes in which the photocathode current is amplified. MCP tubes are more compact than PMTs and are reliable in operating conditions of high environmental stress. The same range of photocathode materials is available in MCPTs as PMTs. MCPTs can provide a wide range of electron gain as available depending upon whether a single MCP or a stack of MCPs is used. The structure of a PMT and MCPT are compared in Fig. 48.

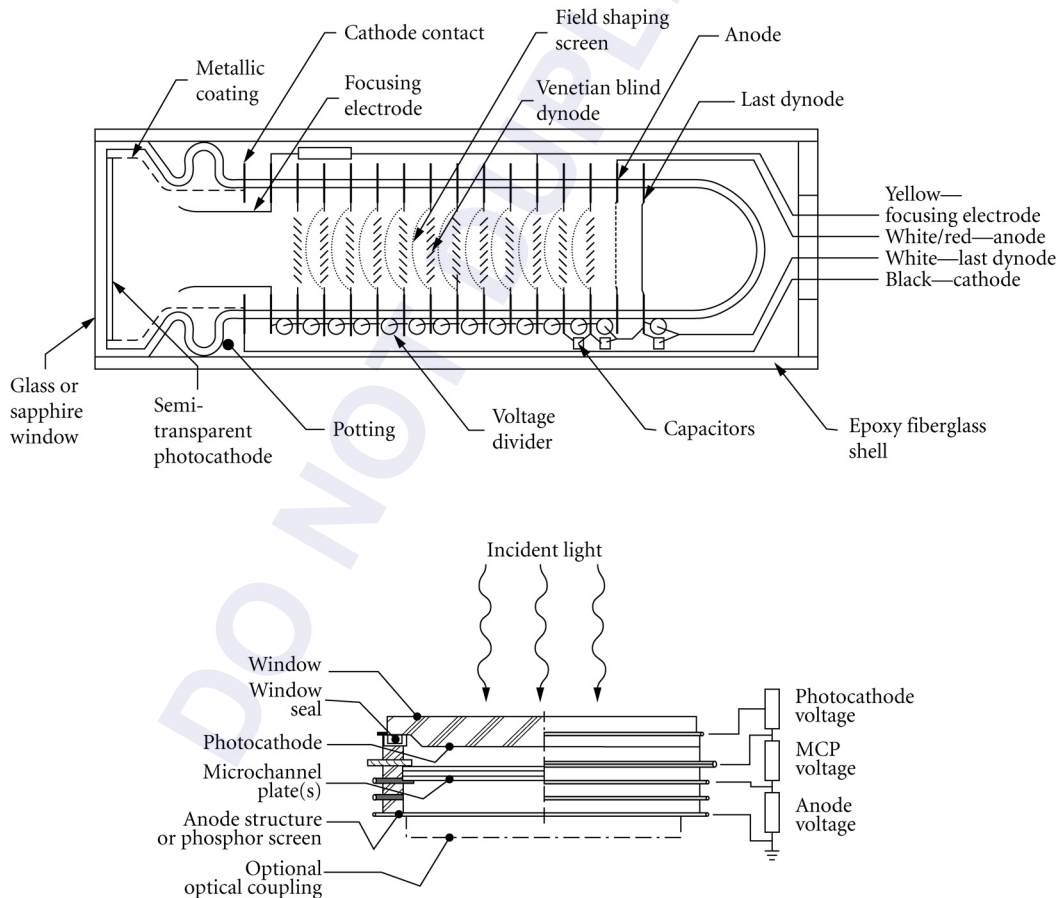


FIGURE 48 Comparison of photomultiplier tube (PMT) and microchannel plate tube (MCPT) construction. (*EMR Photoelectric.*)

Sensitivity In modern phototubes, shot noise due to the cathode dark current is by far the most important noise source. The most common descriptions of phototube sensitivity list both current responsivity (amperes per watt), dark current, and dark noise. Several useful measures of sensitivity are noise equivalent input (NEI) (see Sec. 24.4 “Definitions”), noise equivalent power (NEP), or its reciprocal $D \equiv 1/\text{NEP}$. NEP and NEI in the range 10^{-14} to 10^{-17} W/Hz^{1/2} are not uncommon.

Detectivity is generally limited by dark-current shot noise. Dark current depends on photocathode material, area, and temperature. Thus the best detectivity is obtained with small effective sensitive area. Cooling is especially useful for red-sensitive and near-IR tubes and is generally not worthwhile for others (see “Operating temperature”). Special tube housings which can provide thermoelectric cooling are available.

The spectral response curves shown in Fig. 49 are for the combination of photocathode and window. Historically, this method of description gave rise to the S-response designations, most of which are now obsolete. It is often desirable to separate photocathode response from window transmission. Thus, Fig. 50 shows the quantum efficiency (electrons per incident photon) of a number of photocathodes without window losses. For $\lambda < 400$ nm, each photocathode should maintain its peak quantum efficiency, up to photon energies where multiple photoemissions take place. In Fig. 51, the spectral dependence of quantum efficiency for a variety of modern photocathode/window combinations is illustrated.

D^* is a meaningful figure of merit for phototubes whose sensitivity is limited by dark noise (shot noise on the dark current) and whose emitting photocathode area is clearly defined, but D^* must be used with caution because, although modern phototubes are generally dark-noise-limited devices, they are often limited by noise in signal, that is, the noise content of the signal itself.¹⁸ Serious errors in predicting the detection capability of phototubes will arise if noise in signal is ignored and D^* is presumed to be the important limiting parameter [see Eqs. (6) to (8)]. Very little reliable data are presently available on D^* for photoemitters. However D^* curves for S-1, S-20, and S-25 are shown in Figs. 52 (300 K) and 53 (PMT cooled to 200 K).

Short-wavelength considerations Window considerations are as follows:

For $\lambda > 200$ nm: Windows are essential, as all useful photocathode materials are oxidized and performance would otherwise be destroyed.

$200 \text{ nm} > \lambda > 105$ nm: Photocathode materials are not oxidized by dry air (moisture degrades performance). Windows are optional. LiF windows have shortest known cutoff, 105 nm. For $\lambda < 180$ nm, it is generally advisable to flush with dry nitrogen.

$\lambda < 105$ nm: No windows are available.

Since air absorbs radiation in the region 0.2 to 200 nm (ozone absorbs 200 to 300 nm), it is necessary to include the (windowless) detector in a vacuum enclosure with the source.

A useful technique for avoiding the far-ultraviolet window-absorption problem (provided $\lambda >$ vacuum ultraviolet) is to coat the outside of the window of a conventional PMT with an efficient fluorescent material, for example, sodium salicylate, which absorbs in the ultraviolet and reemits in the blue, and is efficiently detected by most photocathodes.⁹

Solar-blind considerations are as follows—although most photocathodes have high quantum efficiency at short wavelengths, background-noise considerations often preclude their use at very short wavelengths, and very wide bandgap semiconductor photocathodes such as CsI, KBr, Cs₂Te, and Rb₂Te (having peak quantum efficiency a little greater than 10 percent) often give better signal-to-noise ratio. This sacrifice in quantum efficiency to obtain insensitivity to wavelengths greater than those of interest would not be necessary if suitable short-wavelength pass filters were readily available.

For applications where it is desirable that the detector not see much solar radiation, one can use photocathodes whose high work function precludes photoemission for photons of too low an energy. Figure 54 shows quantum efficiency versus λ for three such photocathodes, tungsten, CsI, and Cs₂Te, compared with Cs₃Sb and GaAs(Cs), which are not solar blind.

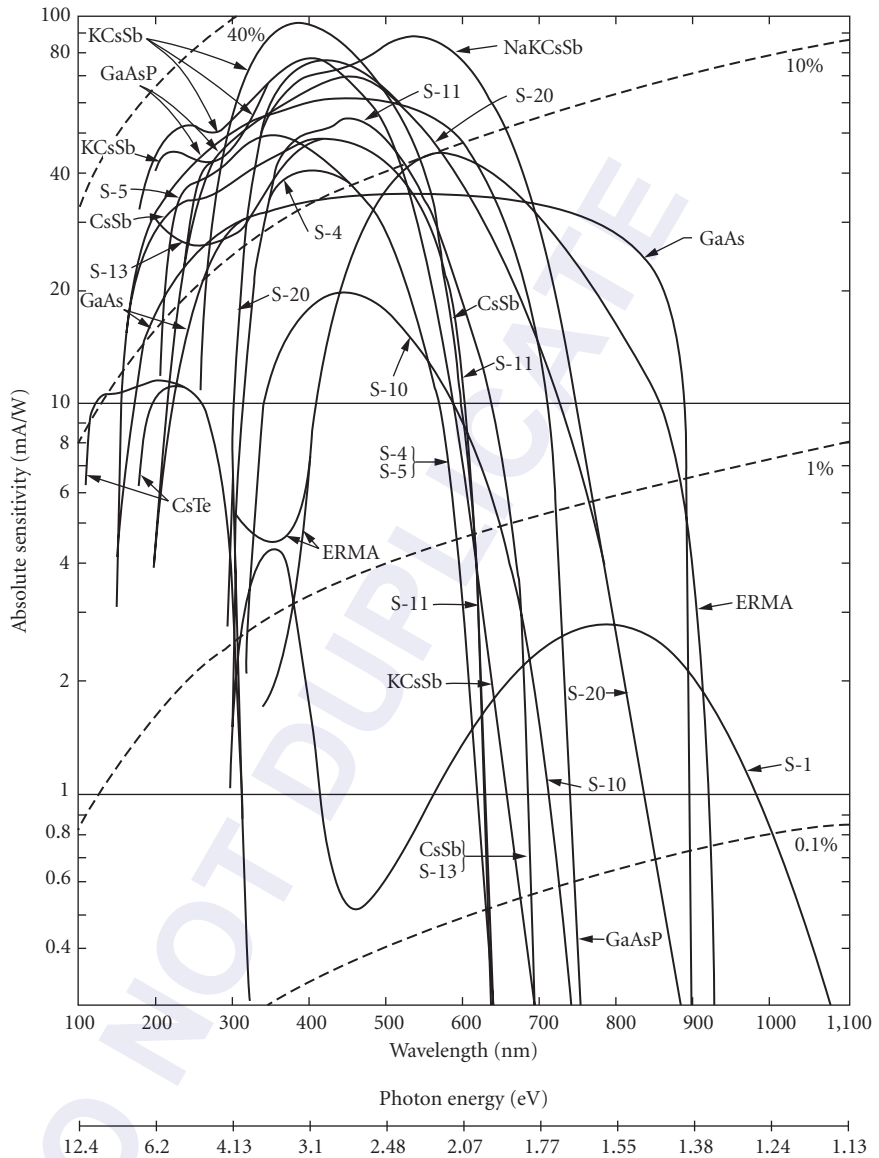


FIGURE 49 Spectral sensitivity of various photoemitters. Dotted lines indicate photocathode quantum efficiency. Chemical formulas are abbreviated to conserve space. S-1 = AgOCs with lime or borosilicate crown-glass window; S-4 = Cs₃Sb with lime or borosilicate crown-glass window (opaque photocathode); S-5 = S₅Sb with ultraviolet-transmitting glass window; S-8 = Cs₃Bi with lime or borosilicate crown-glass window; S-10 = AgBiOCs with lime or borosilicate crown-glass window; S-11 = Cs₃Sb with lime or borosilicate crown-glass window (semitransparent photocathode); S-13 = Cs₃Sb with fused-silica window (semitransparent photocathode); S-19 = Cs₃Sb with fused-silica window (opaque semicathode); S-20 = Na₂ KCsSb with lime or borosilicate glass window. ERMA = extended red multialkali (RCA; IIT uses MA for multialkali). This curve is representative of several manufacturers' products. Many variations of this response are available, for example, trade-offs between short- and long-wavelength response. (From RCA Electronic Components, chart. PIT-701 B.)

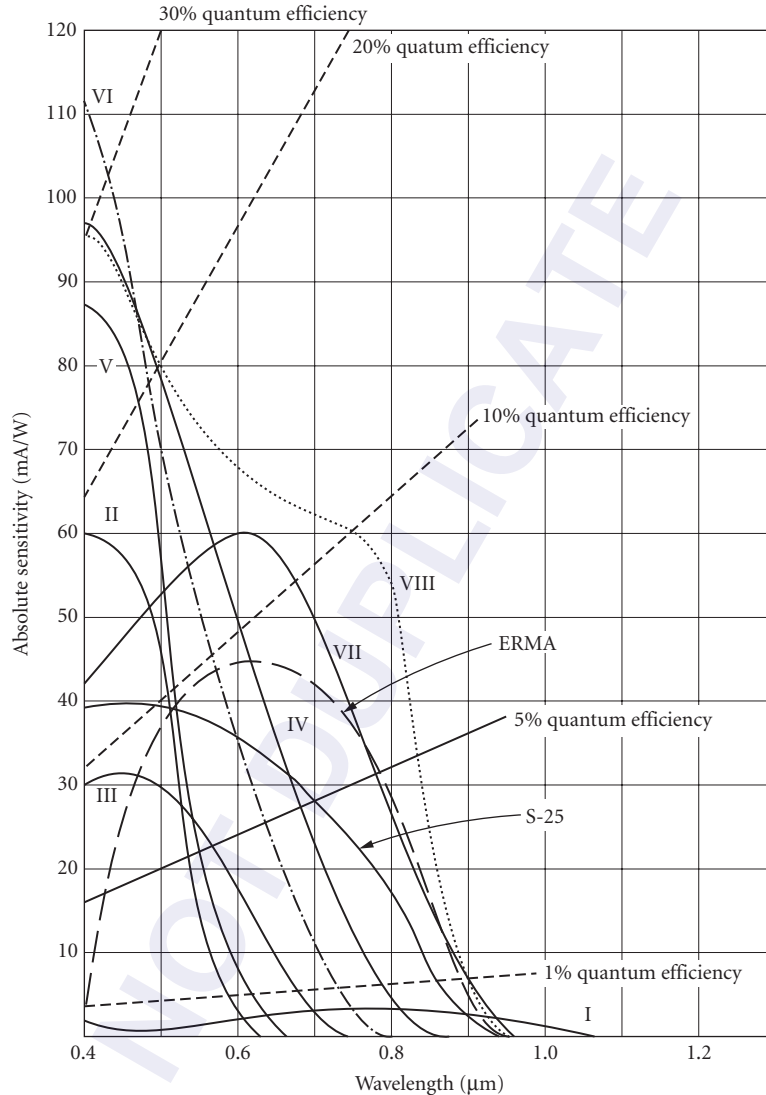
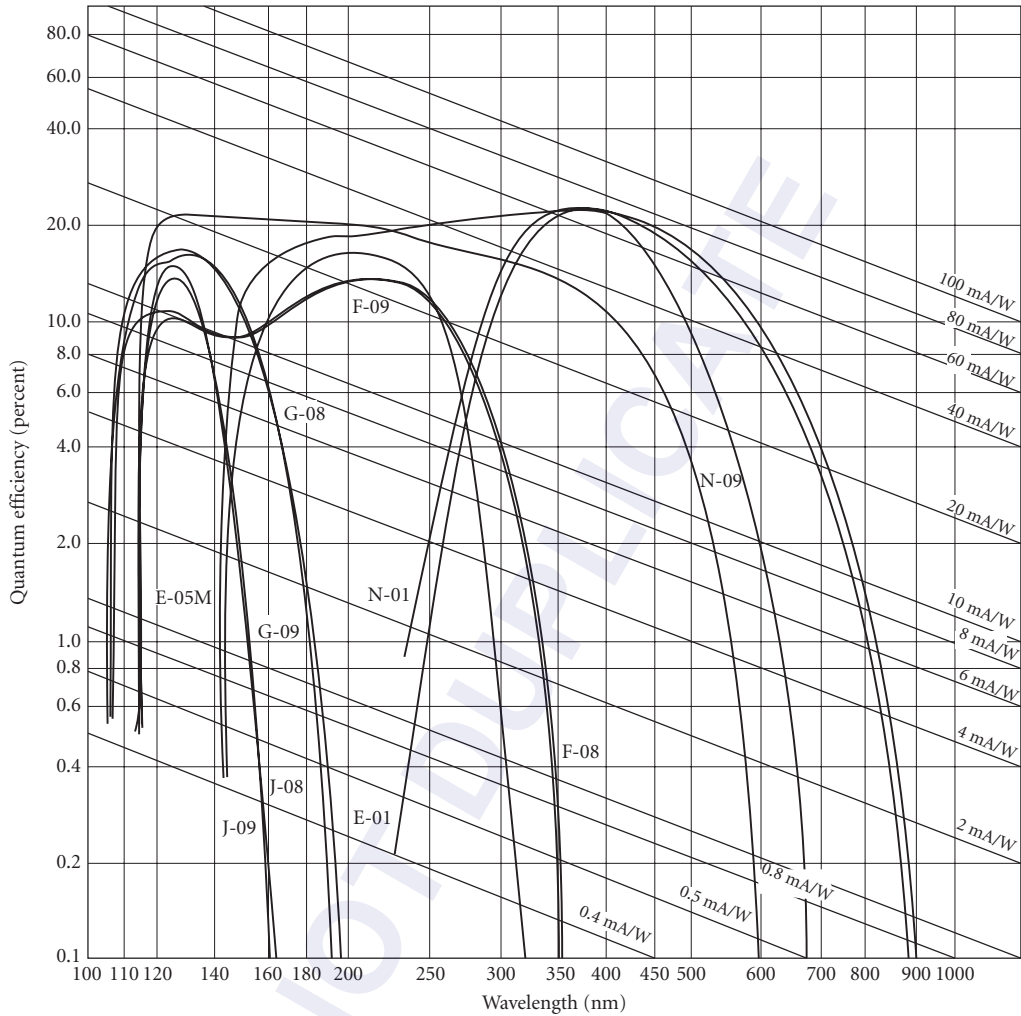


FIGURE 50 Photocathode responsivity and quantum efficiency versus wavelength (no windows). For abbreviations, see Fig. 49. I = S-1, II = S-11, III = S-10, IV = S-20, V = K_2CsSb , VI = $K_2CsSb(O)$, VII = $NaKCsSb_3$, VIII = $GaAs(Cs)$. (Based on material from RCA.)

Quantum efficiency Figures 49, 50, 51, and 54 show photocathode spectral quantum efficiency (probability that one photoelectron is emitted when a single photon is incident). Note that there are fairly few basic photocathode materials and that the window often determines effective quantum efficiency at short wavelengths.

For $\lambda < 40$ nm, a wide variety of photocathode materials are available with high quantum efficiency. Many of these materials, such as tungsten, are not destroyed by being subjected to air, so that open structures can be used, consisting of a photocathode multiplier chain without window. The complete windowless structure is then placed in a vacuum enclosure with the source of radiation.



Photocathode key

Key letter	Description	Long-wavelength cutoff (Note 1)	Long-wavelength sensitivity (Note 2)
E	Tri-alkali (S-20)	850 nm	780 nm
F	Cesium telluride	355 nm	340 nm
G	Cesium iodide	195 nm	185 nm
J	Potassium bromide	165 nm	150 nm
N	High-temperature Bi-alkali	690 nm	640 nm
Q	Rubidium telluride	320 nm	300 nm

Window material key

Key no.	Description	Short-wavelength cutoff*
01	Borosilicate Glass	270 nm
05	UV Grade Sapphire	145 nm
08	UV Grade Lithium Fluoride	105 nm
09	Magnesium Fluoride	115 nm

Note 1—Point at which QE becomes 1% (typical) of peak QE.

Note 2—Point at which QE is 1% (typical).

*10% Energy transmission

FIGURE 51 Quantum efficiency of photocathode/window combinations as a function of wavelength. (EMR Photoelectric)

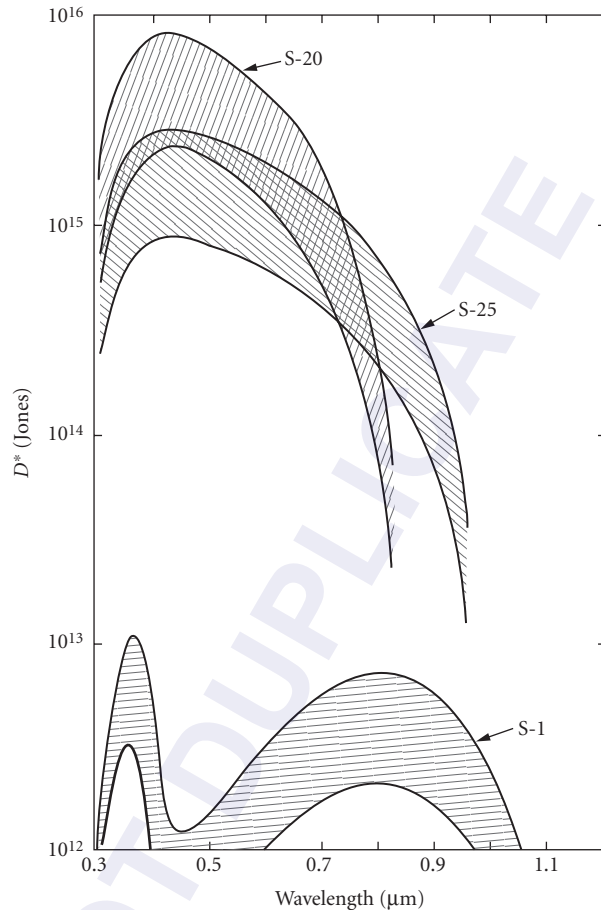


FIGURE 52 Range of D^* for uncooled photomultiplier tubes ($T = 300$ K). For abbreviations, see Fig. 49. S-25 = same as S-20 but different physical processing. (Based on material from RCA)

The quantum efficiency at any wavelength can be calculated from the formula

$$\eta = \frac{\mathcal{R} \times 1239.5}{\lambda} \quad (22)$$

where \mathcal{R} = photocathode response, A/W, and λ = wavelength, nm.

A useful technique for improving quantum efficiency, reported by Livingston¹⁹ and Gunter,²⁰ involves multipassing the photocathode by trapping the light inside the photocathode using a prism.

Responsivity PMT responsivity depends upon photocathode quantum efficiency and subsequent dynode gain. For most purposes, the dynode gain in a well-designed PMT introduces no significant degradation in the photocathode signal-to-noise ratio. Figure 49 shows photocathode response expressed in photocurrent (amperes) per incident radiation power (watts).

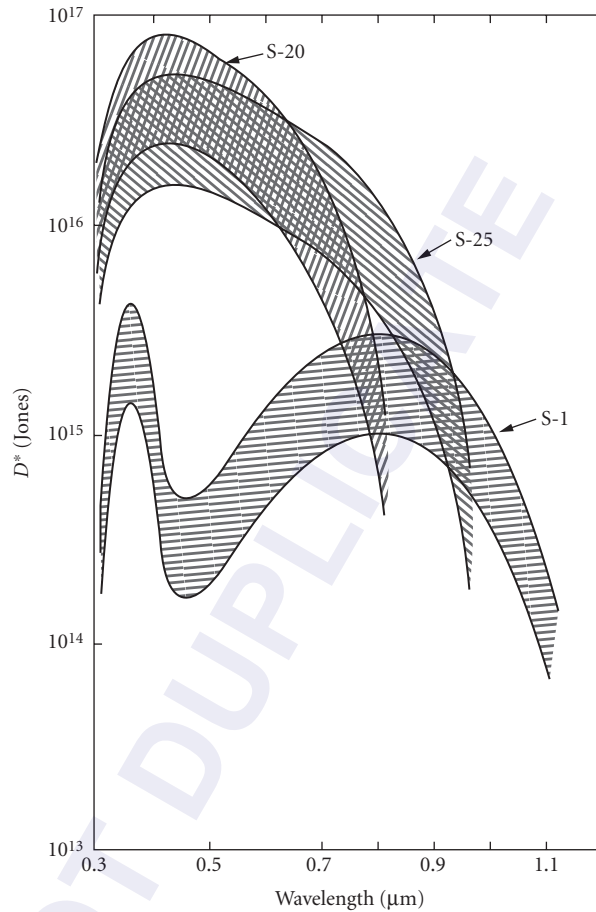


FIGURE 53 Range of D^* for uncooled photomultiplier tubes ($T = 300$ K). For abbreviations, see Fig. 49. S-25 = same as S-20 but different physical processing. (Based on material from RCA) 2 Range of D^* for uncooled photomultiplier tubes ($T = 300$ K). For abbreviations, see Fig. 49. S-25 = same as S-20 but different physical processing. (Based on E.H. Eberhardt, "D* of Photomultiplier Tubes and Image Detectors", ITT Industrial Labs, 1969.)

Noise The limiting noise in a PMT depends on the level of illumination. For low-level detection, limiting noise is the shot noise on the dark current,

$$i_n = (2e j_{\text{dark}} \Delta f)^{1/2} \quad (23)$$

For high illumination levels the shot noise on the signal photocurrent

$$i_n = (2e i_{\text{signal}} \Delta f)^{1/2} \quad (24)$$

far exceeds that on the dark current. Manufacturers usually express noise as photocathode dark current or anode dark current for given gain, which is therefore traceable to photocathode dark current.

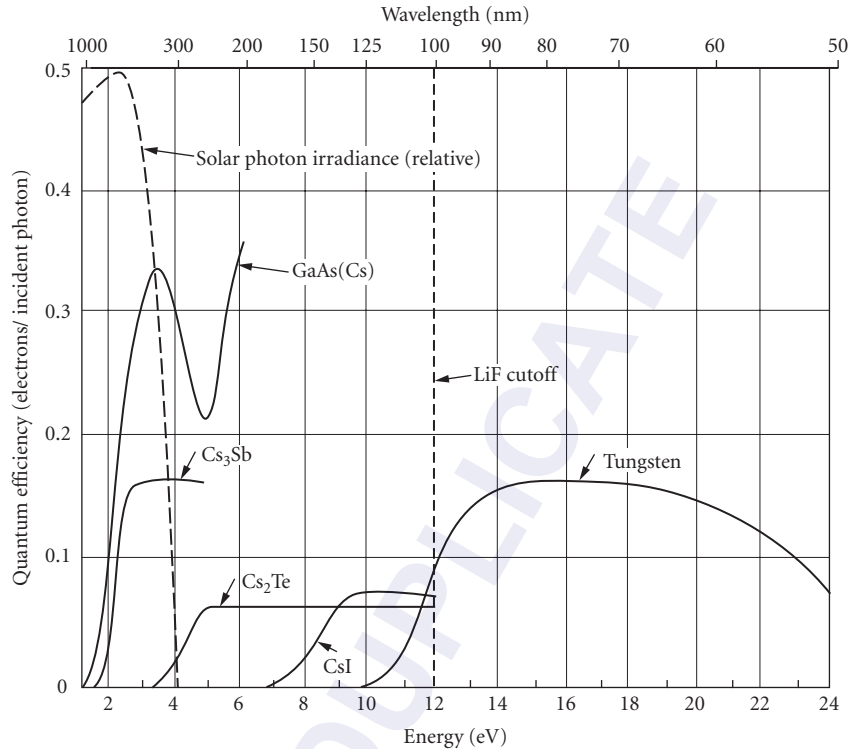


FIGURE 54 Quantum efficiency versus wavelength (photon energy) for several photoemitters.

Photocathode dark current is approximately proportional to photocathode area so that small photocathode effective areas can be expected to have reduced noise. Figure 55 shows how anode dark current and gain increase with applied voltage for a typical PMT.

Minimum detectable power is related to limiting noise through responsivity via

$$\text{NEP} = i_n / \mathcal{R} \quad (25)$$

where \mathcal{R} is in amperes per watt.

Operating temperature Dark current due to thermionic emission, usually greater in red-sensitive tubes, can be reduced by cooling (see Ref. 21). The trialkali (S-20) performance does not benefit from cooling below 255 K. Maximum beneficial cooling (three to four dark counts per second) for AgOCs (S-1), (Cs)Na₂KSb (S-20), and Cs₃Sb (S-11) is 195, 255, and 239 K, respectively. Most photocathodes become noisier as temperature rises above ambient because of increased thermionic emission. Because its thermionic emission starts at a very low value, (Cs)Na₂KSb is a useful photocathode up to temperatures of approximately 373 K.

Response time The rise time of photomultiplier tubes depends chiefly on the spread in transit time during the multiplication process. For photomultiplier tubes, this spread is about 10 ns. Some tubes with specially designed electron optics can give spread as low as 1 ns. The crossed-field PMT makes possible a spread as small as 0.1 ns. Microchannel plate tubes have response times of a nanosecond or less.

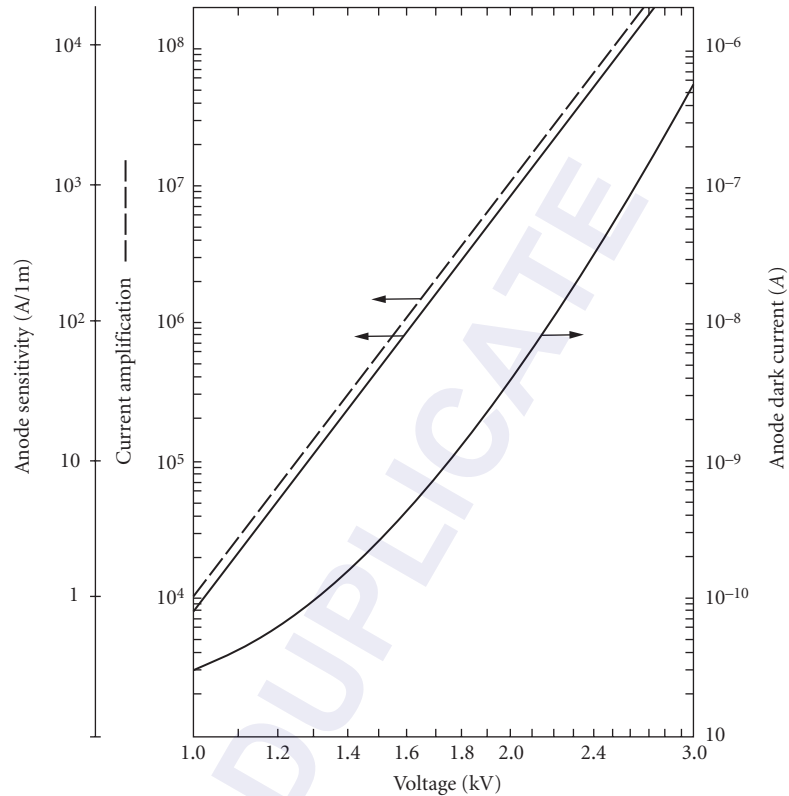


FIGURE 55 Typical current amplification and anode dark current as a function of applied voltage. (Based on E.H. Eberhardt, "D* of Photomultiplier Tubes and Image Detectors," ITT Industrial Labs, 1969.)

For high-speed work (<1 ns rise time), good transmission-line technique must be used to obtain impedance match and to avoid reflection. The bandwidth of the output circuit will depend upon the total capacitance (PMT circuit plus stray capacitances) and the value of the load resistance.

Linearity Photomultiplier tubes are nearly all linear to about 1 percent for cathode currents of 0.1 μA or less. Some tubes may be linear to better than 0.1 percent but must be individually selected.²² Probably most of the nonlinearity results from the dynode structure.

Sensitive area No fundamental limitation. Only recently available with very small effective areas for extremely low dark current. Magnetic focusing has been used so that only a small fraction of the photocathode is used electron-optically.

Sensitivity profile Usually uniform within 20 to 50 percent. Microchannel plate detectors may have uniformity of ± 5 percent.

Stability PMTs are subject to short- and long-term drift which can depend upon anode current, changes in anode current, storage times, and aging or anode life. They are also subject to change if exposed to magnetic fields or changes in temperature. Vibration of the tube may modulate the signal (microphonic effect).

Recommended circuit See Fig. 56.

1. Since a PMT is a current generator, increasing output resistance R_1 increases output voltage. An upper limit to R_1 may be imposed either by the time-constant limitation or by nonlinearity, which results from a space charge produced near the anode when the anode is left nearly floating electrically.
2. The rated photocathode current (referred to anode current through gain) should not be exceeded.
3. Care should be taken not to destroy the photocathode with light (heating).
4. When large currents are drawn, it may affect later dynode interstage voltage and hence gain, causing nonlinearity; for example, in Fig. 44, if the photocurrent from DY 10 to anode becomes comparable to the biasing current, through R_{11} , the gain of the final stage is reduced. This can be avoided by biasing the dynodes with constant-voltage sources.
5. To avoid dynode damage, final dynode current must not exceed the value suggested by the manufacturer.

Photon counting At the photocathode, the shot-noise-limited signal-to-noise ratio (with negligible dark current) is

$$\frac{i_s}{i_n} \frac{i_s}{(2ei_s \Delta f)^{1/2}} \left(\frac{i_s}{2e\Delta f} \right)^{1/2} = \left(\frac{N_s}{2\Delta f} \right)^{1/2} \quad (26)$$

where N_s is the photoelectron rate at the photocathode. Thus, for extremely low levels of illumination, the ideal signal-to-noise ratio becomes very poor. At this point there is much to be gained by abandoning attempts to measure the height of the fluctuating signal (Fig. 57a) in favor of digitally recording the presence or absence of individual pulses (Fig. 57e).

Single photoelectron counting can be achieved by using a pulse amplifier (see Fig. 58), which suppresses spurious dark-noise pulses not identical in amplitude and shape to those produced by photoelectrons.

An upper practical limit for (random) photon counting is set by convenient amplifier bandwidths at about 10^5 s^{-1} . For reasonable (1 percent) statistical accuracy, this implies a 10-MHz bandwidth.

Gallium phosphide dynodes The development of GaP dynodes for increased secondary-electron production^{23,24} makes possible unambiguous discrimination of small numbers of individual photoelectron counts which was not previously possible with lower dynode gains. This is shown in Fig. 59, where the spread in number of secondary electrons ($N \times \text{gain}$) is just $(N \times \text{gain})^{1/2}$.

In addition to the aforementioned fundamental advantage of high dynode gain, the large gain per stage in the first dynodes also helps discriminate against noise introduced by later stages of amplification. Also, fewer stages of amplification are required.

Manufacturers ADIT, EMR Photoelectric, Bicon, Burle, Edinburgh Instruments, Galileo Electro-Optics, Hamamatsu, K and M Electronics, id Quantique, International Light, Optometrics USA, Oriol, Phillips Components, Photek, Photon Technology, Photonis, Penta Laboratories, Thorn EMI, Varo.

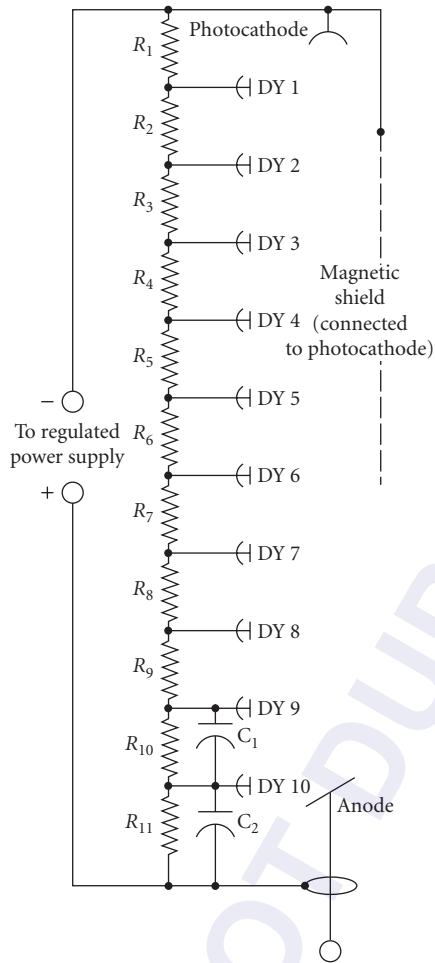


FIGURE 56 $C_1 = 68 \text{ pF} \pm 10 \text{ percent}$, 500 V (dc working); $C_2 = 270 \text{ pF} \pm 10 \text{ percent}$, 500 V (dc working); $R_1 = 220 \text{ k}\Omega \pm 5 \text{ percent}$, 1/4 W; $R_2 = 240 \text{ k}\Omega \pm 5 \text{ percent}$, 1/4 W; $R_3 = 330 \text{ k}\Omega \pm 5 \text{ percent}$, 1/4 W; R_4 to $R_{11} = 220 \text{ k}\Omega \pm 5 \text{ percent}$, 1/4 W. (Based on E.H. Eberhardt, "D* of Photomultiplier Tubes and Image Detectors," ITT Industrial Labs, 1969.)

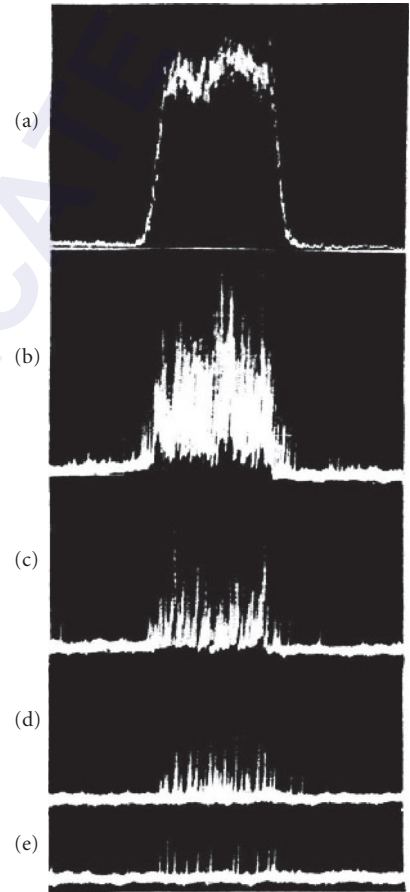


FIGURE 57 Oscilloscope presentation of PMT output when reviewing square-wave chopped light pulse. In (a) to (e) the intensity is reduced and gain is increased commensurately. (Courtesy of E.H. Eberhardt.)

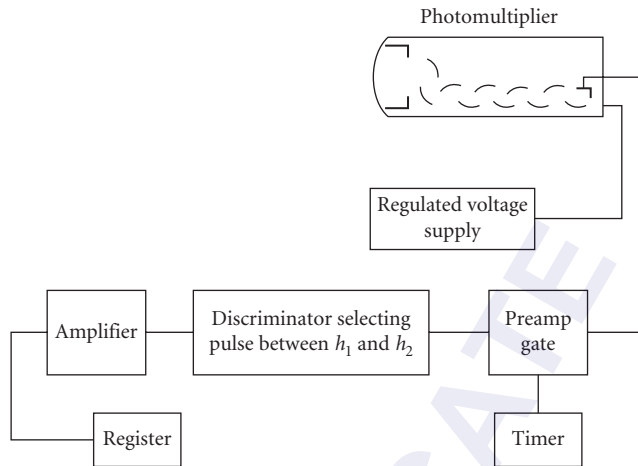


FIGURE 58 Photomultiplier and associated circuits for photon counting. (ITT Report E5.)

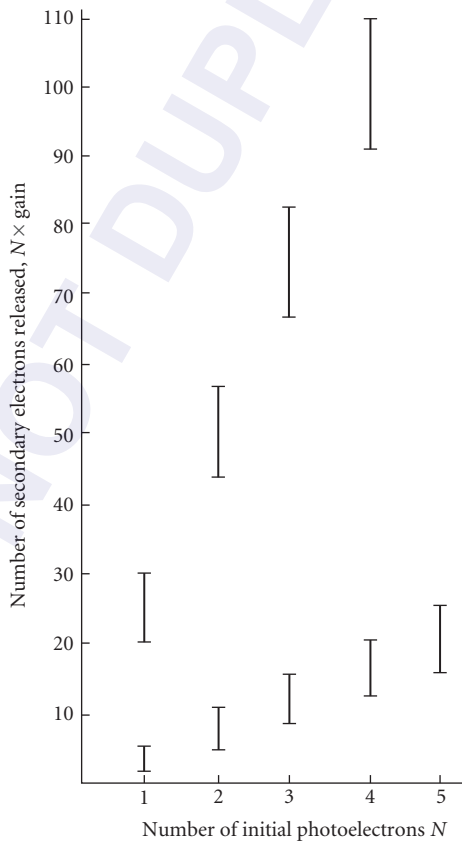


FIGURE 59 Spread in number of secondary electrons for various phototube gains.

GaN and AlGaN Gallium nitride photovoltaic detectors, with a bandgap of 3.39 eV have spectral response in the ultraviolet (UV) from 200 to 365 nm, as illustrated in Fig. 60. By using aluminum-gallium nitride—an alloy mixture of AlN and GaN—the spectral response can be tuned to shorter wavelength cutoffs. Spectral response examples are shown in Figs. 61 to 64 to compare GaN with one particular AlGaN alloy. Some devices may be tailored to custom UV bands, such as UVA (320 to 400 nm), UVB (280 to 320 nm), or UVC (100 to 280).

Response at visible wavelengths is low or absent, so that no special filtering may be required to detect UV in the presence of visible lighting or solar radiation—but see the logarithmic spectral Figs. 62 and 64 to see the degree of longer wavelength response. These solid-state devices are potentially useful for operation at elevated temperatures, in high-vibration environments, and in other environments unsuitable for photomultiplier tubes.

The photoconductive GaN devices use interdigitated contact electrodes because of the very high impedance of the GaN films, but currently there may not be any available commercially.

Response: Photovoltaic 0.1 A/W

Dark current: 0.05-nA photovoltaic

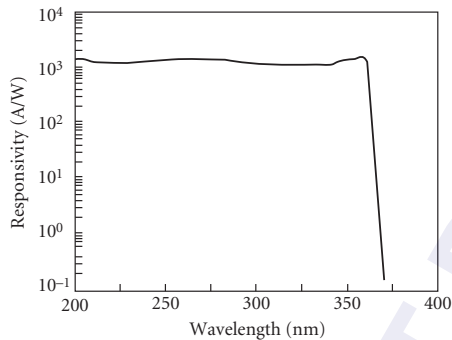


FIGURE 60 Response in amperes per watt for a GaN detector. (Reprinted from *Appl. Phys. Lett.*, vol. 60, no. 23, 1992, p. 2918.)

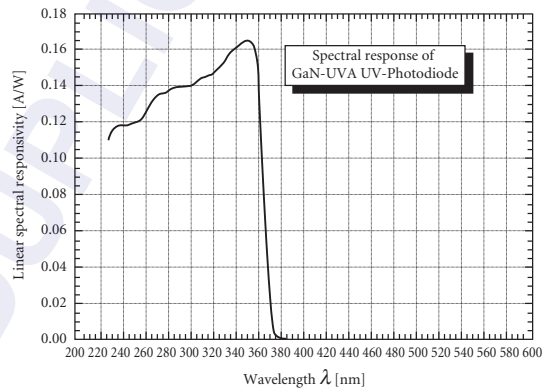


FIGURE 61 Spectral response of a UVA GaN detector shown on a linear vertical scale of amps/watt versus wavelength. (<http://www.boselec.com/products/documents/GaNAlGaNall.pdf>)

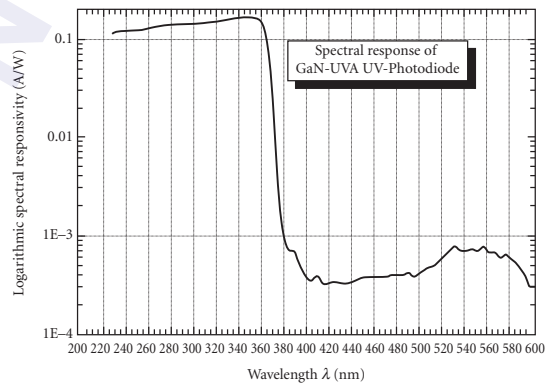


FIGURE 62 Spectral response of a UVA GaN detector shown on a logarithmic vertical scale of amperes/watt versus wavelength. (<http://www.boselec.com/products/documents/GaNAlGaNall.pdf>)

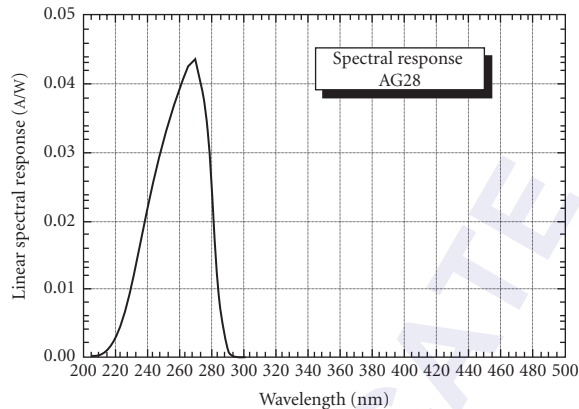


FIGURE 63 Spectral response of a UVC AlGaIn detector shown on a linear vertical scale of amperes/watt versus wavelength. (<http://www.boselec.com/products/documents/GaNAlGaInAll.pdf>.)

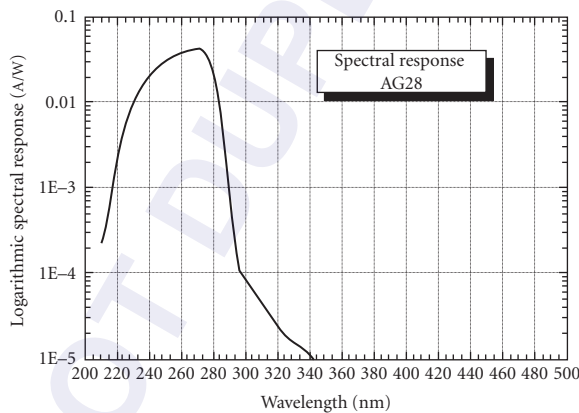


FIGURE 64 Spectral response of a UVC GaN detector shown on a logarithmic vertical scale of amperes/watt versus wavelength. (<http://www.boselec.com/products/documents/GaNAlGaInAll.pdf>.)

Capacitance: 24 pF photovoltaic at 0-V bias

Time constant: Photovoltaic 0.10 ns

Size: 0.076 mm²

Devices with AlGaIn alloys have wider bandgaps and generally lower leakage currents.

Response: Photovoltaic 0.045 A/W.

Dark current: 0.1-pA photovoltaic at 0.1-V reverse bias

Capacitance: 24-pF photovoltaic at 0-V bias

Size: 0.076 mm²

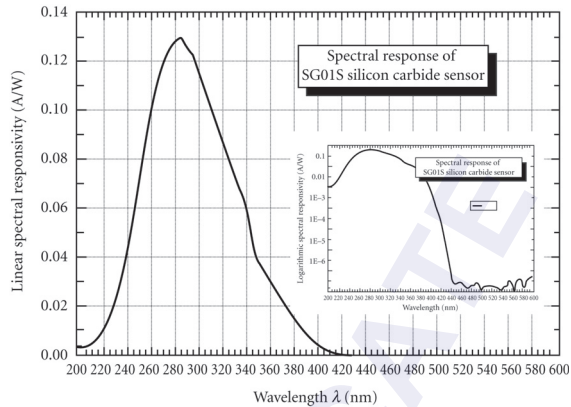


FIGURE 65 Spectral response of an unfiltered, broadband SiC detector shown on a linear vertical scale of amps/watt versus wavelength. The inset shows the same data on a logarithmic scale. (<http://www.boselec.com/products/documents/UVPhotodetectors2-08WWW.pdf>.)

Manufacturers: Advanced Photonix: http://www.advancedphotonix.com/ap_products/standard_GaN.asp?from=leftnav, Boston Electronics: <http://www.boselec.com/products/detuv.html>, Orion Semiconductor: <http://www.orion-semi.com>, SVT Associates: <http://www.svta.com/products/uv/uv.htm>.

SiC Silicon carbide UV detectors are available in photovoltaic structures. The 3-eV bandgap of SiC is slightly narrower than GaN, thereby giving a response that may extend to slightly longer wavelength. However, because the bandgap of SiC is indirect, unlike GaN which is direct, the response cut-on is more gradual in SiC, peaking at a wavelength much shorter than the wavelength corresponding to 3 eV (413 nm)—see Fig. 65. SiC detectors with integrated filters are available.

Response: 0.13 A/W peak

Dark current: 1 fA for a 1 × 1 mm device

Capacitance: 195 pF for a 1 × 1 mm device

Sizes: 0.25 × 0.25 mm, 0.5 × 0.5 mm, 1 × 1 mm

Manufacturers: Boston Electronics: <http://www.boselec.com/products/detuv.html>, Electro Optical Components: http://www.eoc-inc.com/UV_detectors_silicon_carbide_photodiodes.htm

TiO₂ Detectors With a bandgap of 3.2 eV, TiO₂ is another UV photodetector. Photovoltaic devices are made with Schottky diodes. An unfiltered spectral response is shown in Fig. 66. TiO₂ detectors with integrated filters are available.

Response: 0.021 A/W peak

Dark current: 100 pA for a 5.4 × 2.9 mm device

Sizes: 2.2 × 1.9 mm, 5.4 × 2.9 mm

Manufacturer: Boston Electronics: <http://www.boselec.com/products/detuv.html>

GaP Gallium phosphide can provide Schottky photodiodes that cover the UV to mid-visible spectral region as shown in Fig. 67. The bandgap of GaP is 2.26 eV and is indirect, leading to a soft spectral

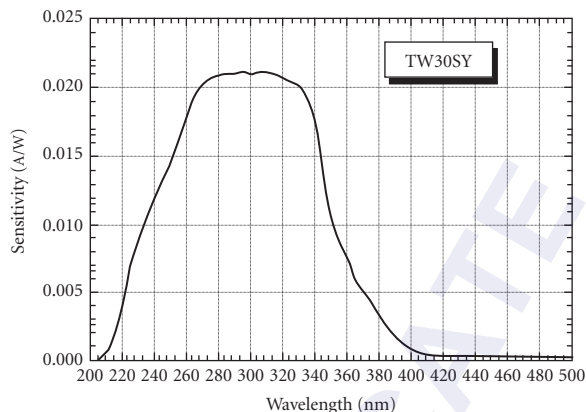


FIGURE 66 Spectral response of a TiO_2 Schottky photodiode detector shown on a linear vertical scale of amperes/watt versus wavelength. (<http://www.boselec.com/products/documents/UVPhotodetectors2-08WWW.pdf>.)

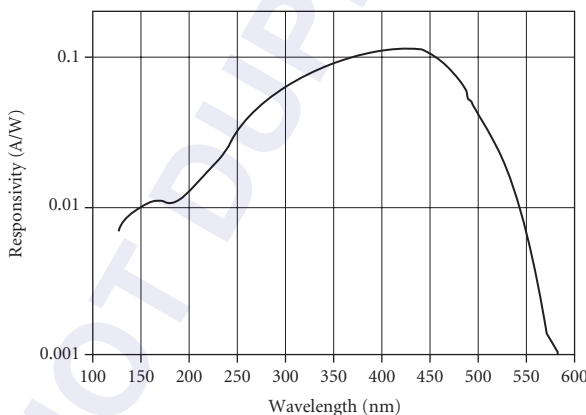


FIGURE 67 Spectral response of a GaP Schottky photodiode detector shown on a logarithmic vertical scale of amps/watt versus wavelength. (<http://www.thorlabs.com/Thorcat/12100/12174-S01.PDF>.)

turn-on and a peak quite far from the wavelength corresponding to the bandgap (549 nm). GaP devices with integrated filters to restrict the response to the UV region are also available.

Response: 0.12 A/W peak

Dark current: 1 nA max for a 2.5×2.5 mm device

NEP @ 440 nm: 1×10^{-14} W/ $\sqrt{\text{Hz}}$ @ 5 V bias

Rise time: 1 nsec @ 5 V bias for a 2.5×2.5 mm device

Fall time: 140 nsec @ 5 V bias for a 2.5×2.5 mm device

Sizes: 1.1×1.1 mm, 2.3×2.3 mm, 2.5×2.5 mm, 4.6×4.6 mm

Manufacturers: Hamamatsu: http://jp.hamamatsu.com/products/sensor-ssd/pd140/pd144/index_en.html, Electro Optical Components: <http://www.eoc-inc.com/ifw/EPD-365-0-2-5.pdf>, Thor Labs: <http://thorlabs.com/thorProduct.cfm?partNumber=FGAP71>

GaAsP Gallium arsenide phosphide alloys can provide photodiodes that cover from the UV to the near-infrared spectral region. The bandgap of GaP is 2.26 eV and is indirect, while that of GaAs is 1.43 eV and is direct. GaAsP alloys from 0 to ~50% GaP are direct bandgap materials while those with higher percentages of GaP are indirect*. A variety of alloys are available, covering the following spectral bands:

Spectral Band (nm)	λ Peak (nm)	Response at Peak (A/W)	Sizes (mm)
400–760	710	0.4	1.3 × 1.3, 2.7 × 2.7, 5.6 × 5.6
300–680	640	0.3	1.3 × 1.3, 2.7 × 2.7, 5.6 × 5.6
300–580	470	0.25	0.8 × 0.8
280–580	470	0.2	0.8 × 0.8
260–400	370	0.06	0.8 × 0.8
190–760	710	0.22	2.3 × 2.3, 4.6 × 4.6
190–680	610	0.18	10.1 × 10.1

Spectral responses of these alloys are shown in Figs. 68 to 73 (ref: http://jp.hamamatsu.com/products/sensor-ssd/pd140/pd143/index_en.html?sort=WAVE_LENGTH4&desc=1&style=F1 for all six figures).

Manufacturer: Hamamatsu: http://jp.hamamatsu.com/products/sensor-ssd/pd140/pd143/index_en.html

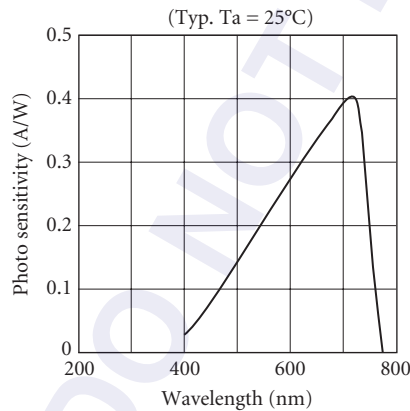


FIGURE 68 Spectral response of a 400 to 760-nm GaAsP photodiode detector with a vertical scale of amps/watt versus wavelength.

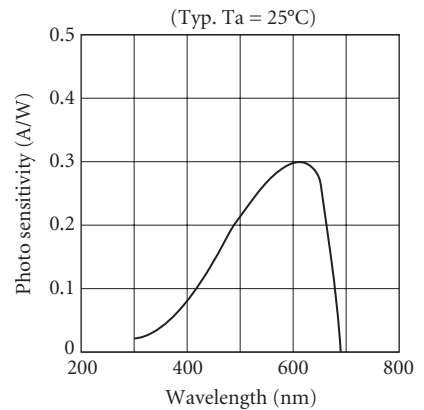


FIGURE 69 Spectral response of a 300 to 680-nm GaAsP photodiode detector with a vertical scale of amps/watt versus wavelength.

*<http://www.iue.tuwien.ac.at/phd/palankovski/node37.html>.

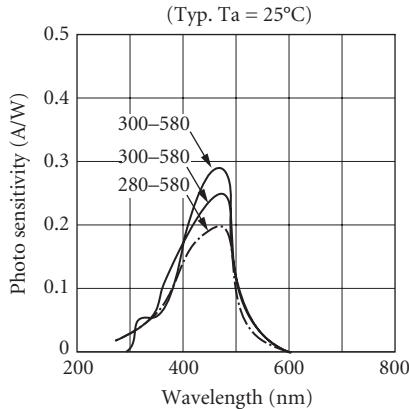


FIGURE 70 Spectral response of 300 to 580 and 280- to 580-nm GaAsP photodiode detectors with a vertical scale of amperes/watt versus wavelength.

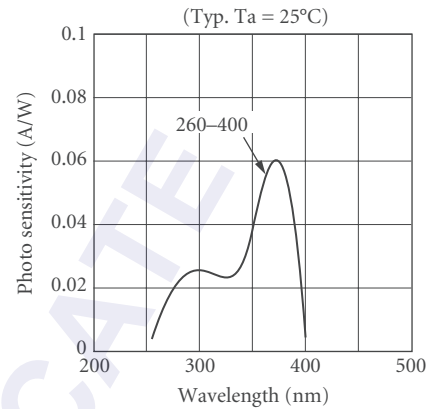


FIGURE 71 Spectral response of a 260- to 400-nm GaAsP photodiode detector with a vertical scale of amperes/watt versus wavelength.

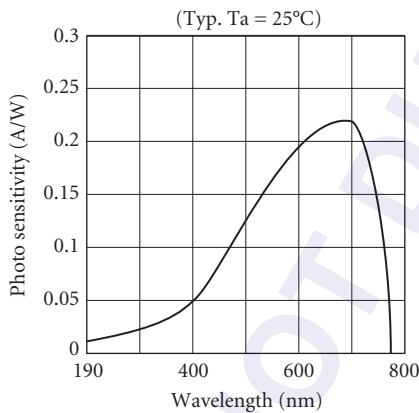


FIGURE 72 Spectral response of a 190- to 760-nm GaAsP photodiode detector vertical scale of amps/watt versus wavelength.

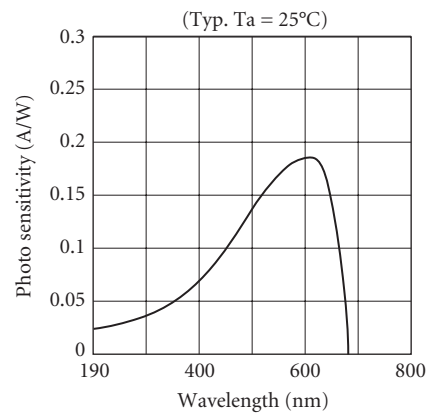


FIGURE 73 Spectral response of a 190- to 680-nm GaAsP photodiode detector vertical scale of amps/watt versus wavelength.

CdS and CdSe Cadmium sulfide and cadmium selenide photoconductors are available for detection of visible light out to 700 to 800 nm. CdS and CdSe films have sheet resistivity in the range of 20 $\text{m}\Omega$ per square at an illumination level of 2 footcandles. The devices are typically made in a linear or serpentine configuration consisting of 2 to 500 squares to maximize the length-to-width ratio. A variety of material “types” are available, offering unique spectral curves for various applications, depending upon the source color. CdS and CdSe are typically slow detectors, with response times of 5 to 100 ms, with speed improving at higher light levels. These devices exhibit “memory” or “history” effects, where the response is dependent upon the storage condition preceding use—the length of storage and time in use, and differences between the storage light level and the light level during use. These history effects may amount to changes in resistance from less than 10 percent to over 500 percent. CdSe has comparably greater memory effect than CdS.

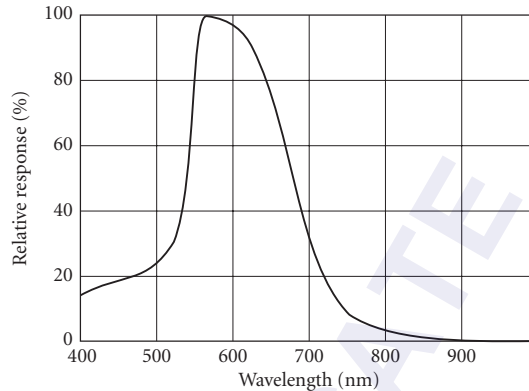


FIGURE 74 Relative spectral response of a “Type 0” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

CdS and CdSe are useful for a variety of commercial applications, both analog and digital, such as camera exposure control, automatic focus and brightness controls, densitometers, night light controls, etc. They are comparatively inexpensive and are available in a wide range of packages and resistance values, including dual cell configurations.

Spectral response: See Figs. 74 to 76.

Resistance and sensitivity: See Figs. 77 to 78.

Temperature coefficient of resistance: See Figs. 79, 80.

Light history effects: See Table 3.

Detector size: 4×4 mm to 12×12 mm approximate, dual elements available.

Manufacturers: In the previous edition, the listed supplier was EG&G VACTEK. Their product line has been acquired by Perkin Elmer. In this transition, all but two of the detector “types” have

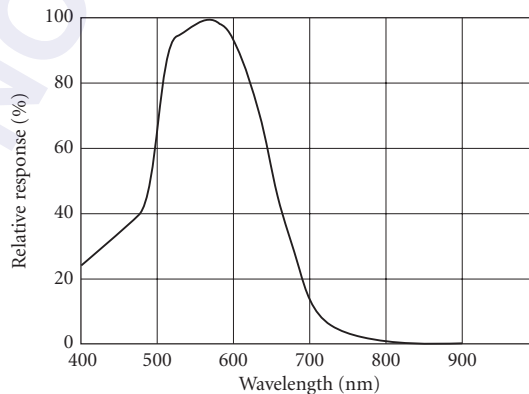


FIGURE 75 Relative spectral response of a “Type 3” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

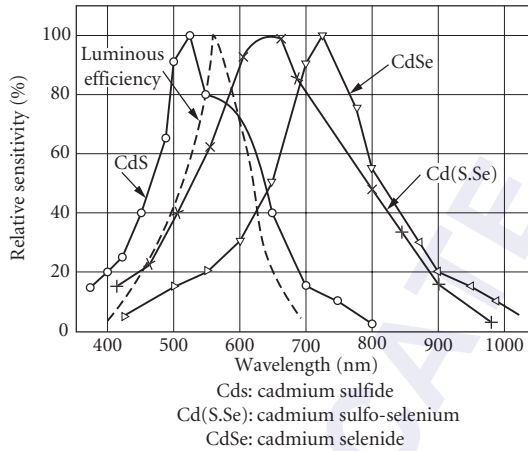


FIGURE 76 Spectral response of CdS, CdSse, and CdSe photocell detectors together with the human eye response or luminous efficiency. (http://www.selcoproducts.com/CFM/photocells/photozell_PDF/Selco_PhotoCells_Construct.pdf.)

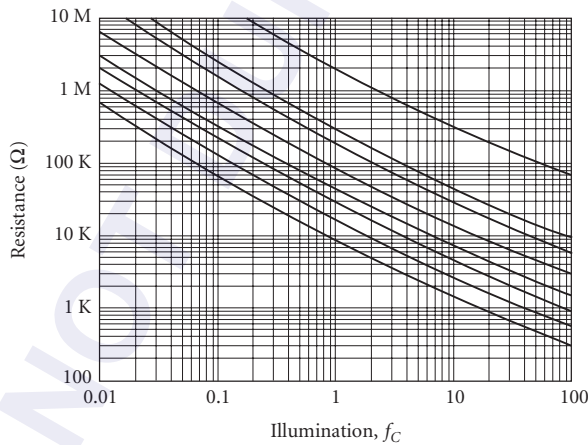


FIGURE 77 Resistance as a function of illumination for a “Type 0” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptois.pdf.)

been dropped—Perkin Elmer now only sells types 0 and 3. Other manufacturers offer comparable varieties and charts from at least one other producer are included.

Jameco Electronics: <http://www.jameco.com/webapp/wcs/stores/servlet/CategoryDisplay?storeId=10001&catalogId=10001&langId=-1&categoryId=151080>, Perkin Elmer: <http://optoelectronics.perkinelmer.com/catalog/Category.aspx?CategoryName=Photocells>, Selco Products: http://www.selcoproducts.com/CFM/photozell_toc.cfm, Silonex: <http://www1.silonex.com/optoelectronics/optophotoc.html>

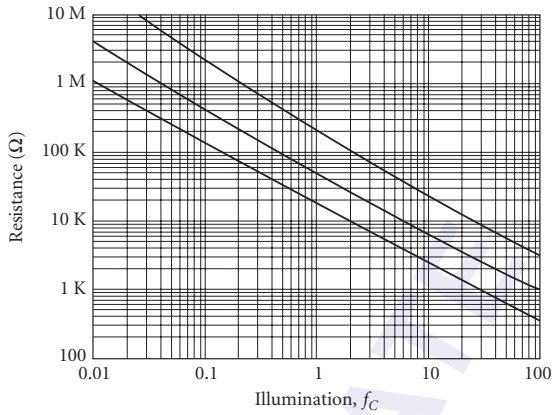


FIGURE 78 Resistance as a function of illumination for a “Type 3” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

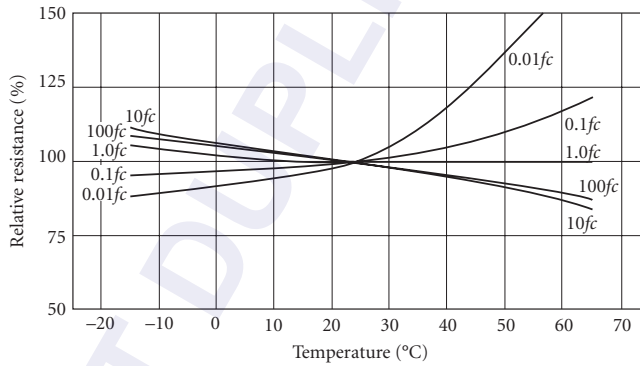


FIGURE 79 Relative resistance as a function of temperature for a “Type 0” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

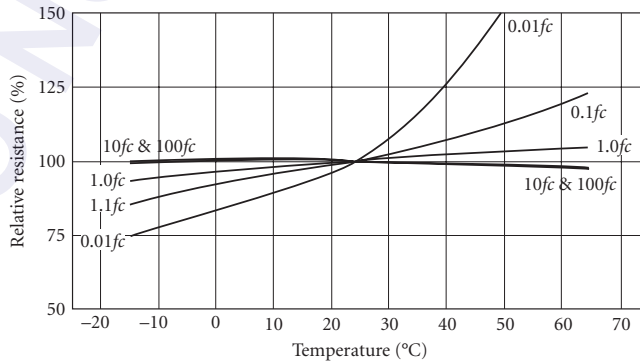


FIGURE 80 Relative resistance as a function of temperature for a “Type 3” CdS photocell. (http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

TABLE 3 Typical Variation of Resistance with Light History Expressed as a Ratio R_{LH}/R_{DH} at Various Test Illumination Levels in Foot Candles.

Illumination (foot candles)	0.01	0.1	1.0	10	100
R_{LH}/R_{DH} ratio	1.55	1.35	1.20	1.10	1.10

(http://optoelectronics.perkinelmer.com/content/RelatedLinks/Brochures/BRO_PhotoconductiveCellsAndAnalogOptoiso.pdf.)

R_{LH} is the resistance after “infinite” exposure to light, while R_{DH} is the resistance after “infinite” exposure to a dark environment. Infinite may be approximated by 24 hours.

CdTe Cadmium telluride and cadmium zinc telluride detectors are chemical group II-VI materials having an energy bandgap of about 1.6 eV, corresponding to a spectral cutoff in the vicinity of 775 nm. These devices, however, are principally used for gamma ray detection because of their high z number which translates into a high absorption coefficient for gamma rays. The principal advantage of CdTe in this application is its ability to operate at room temperature, in comparison with Ge gamma ray detectors which must typically be cooled to 77 K. Figure 81 illustrates the absorption of CdTe as a function of gamma ray energy out to 300 keV.

Sensitivity: See Fig. 81.

Standard sizes: Wafers in 10- and 16-mm diameter; rods $7 \times 2 \times 2$ mm; cubes $2 \times 2 \times 2$ mm.

Standard thickness: 1 and 2 mm.

Bias voltage: 150 to 300 V/cm.

Operating temperature range: -10 to $+55^\circ\text{C}$

Leakage current: 10 to 300 nA

Capacitance: 10 pF

Response time: $< 1 \mu\text{s}$.

Manufacturers: Acrorad: <http://www.acrorad.co.jp/us/index.html>, Aurora, II-VI eV Products: <http://www.evproducts.com/>, Perkin Elmer, Radiation Monitoring Devices: <http://www.rmdinc.com/products/p007.html>

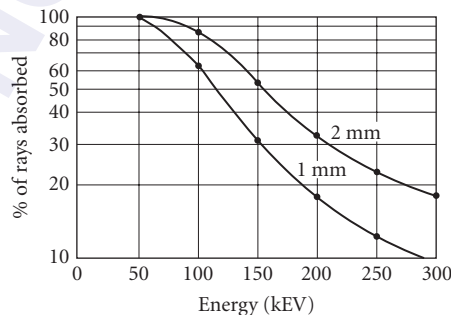


FIGURE 81 The high percentage of rays absorbed by CdTe makes these detectors highly sensitive. At 100 keV, a 2-mm-thick detector absorbs 85 percent of the rays. (*Radiation monitoring devices, Cadmium Telluride brochure.*)

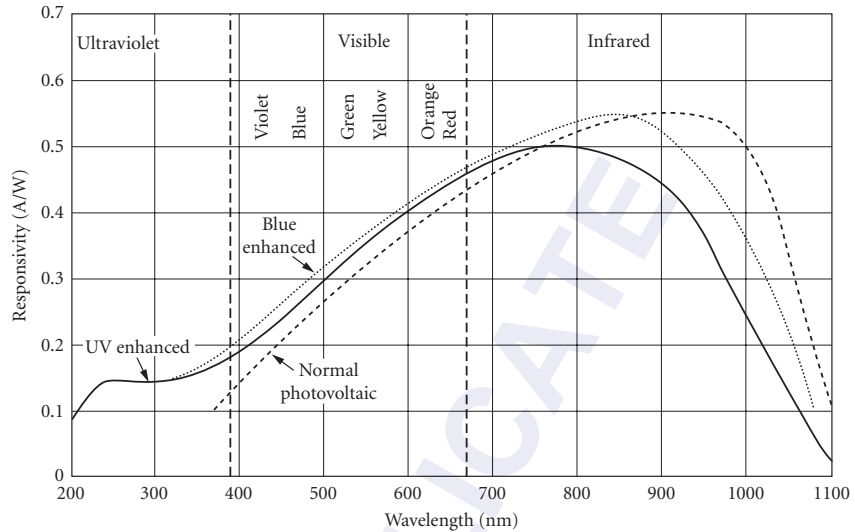


FIGURE 82 Typical spectral response for pn junction, blue-enhanced, and UV-enhanced silicon photodiodes. (UDT Sensors, Inc., *Optoelectronic Components Catalog*.)

Si Silicon photovoltaic detectors are widely available. They are useful at wavelengths shorter than about $1.1\ \mu\text{m}$ and can even be used for x-ray and gamma-ray detection. There are four main silicon detector types:

- pn junction photodiodes, generally formed by diffusion, but ion implantation can also be used.
- pin junctions, which have lower capacitance and hence higher speed, and because of a thicker active region have enhanced near-IR spectral response.
- UV- and blue-enhanced photodiodes
- Avalanche photodiodes with significant internal gain, combining high speed and sensitivity

The main parameters of interest are spectral response (see Fig. 82), time constant, and zero-bias resistance or reverse-bias leakage current. Silicon material has an indirect bandgap and hence the spectral cutoff is not very sharp near its long-wavelength limit as shown in Fig. 82. The effective time constant of pn junction silicon detectors is generally limited by resistance-capacitance (RC) considerations rather than by the inherent speed of the detection mechanism (drift and/or diffusion). High reverse bias may or may not shorten charge collection time, but it generally reduces cell capacitance, and therefore the RC product, therefore, reverse bias usually results in faster response.

On the other hand, increased reverse bias causes increased noise, so that a trade-off exists between speed and sensitivity. For high-frequency applications, load resistance should be made small, although this makes Johnson (thermal) noise comparatively larger, which limits sensitivity (see Fig. 83). In order to keep sensitivity high when using these devices at high frequency, operational (current-mode) amplifiers, which can be built into the detector package, and avalanche photodiodes, which incorporate built-in gain before the load resistor is encountered, have been developed. Very careful regulation of the detector bias is required for stable operation of avalanche photodiodes.

Silicon pn junction photodiodes These are general purpose when high sensitivity is required and time constants on the order of a microsecond are permissible. The device construction is illustrated in Fig. 84. These devices are typically operated in a photovoltaic mode at zero bias, but can be used in a photoconductive mode in which the device is reverse biased.

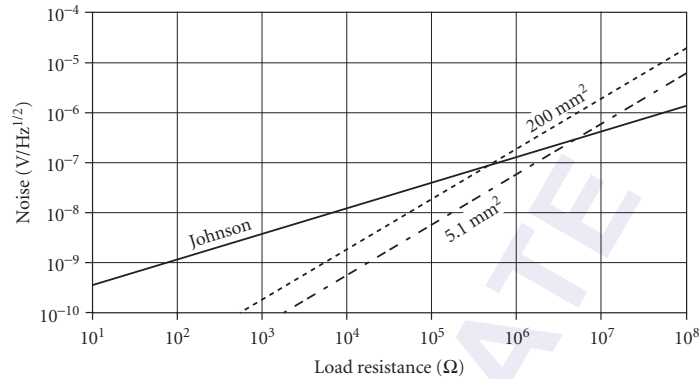


FIGURE 83 Output noise as a function of circuit load resistance for *pin* silicon photodiodes with areas of 5.1 and 200 mm², compared with the Johnson noise of the load resistor. Dark current measured at 10-V reverse bias for the detector with area of 5.1 mm² is 10 nA, and 100 nA for the detector with an area of 200 mm². Note that good preamplifiers have a noise level of about 1 nV/Hz^{1/2}, depending upon the bandwidth. (Detector data from UDT Sensors, *Optoelectronic Components Catalog*.)

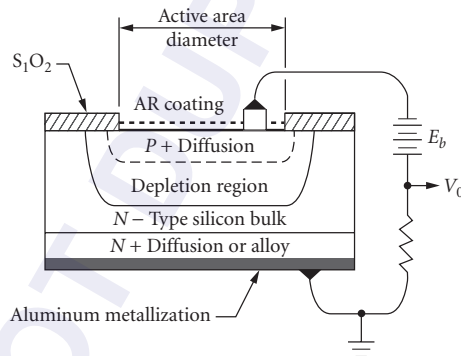


FIGURE 84 Planar diffused silicon photodiode construction. (UDT Sensors, Inc., *Optoelectronic Components Catalog*.)

Sensitivity: $D^*(\lambda_{pk}) \approx \text{mid-}10^{12} \text{ to } 10^{13}$ Jones, $D^*(2800 \text{ K}) \approx 2 \times 10^4$ Jones, becoming amplifier-limited for small-area detectors (see Figs. 85 and 86). D^* can also be estimated from the R_0A product (detector zero-bias resistance or shunt resistance diode area), which is illustrated in Fig. 87, in combination with Fig. 19, which illustrates the dependence of D^* on R_0A product.

Noise: See Figs. 88 (noise vs. bias) and 89 (noise vs. temperature); as T drops, impedance rises, so that decreasing noise current produces increasing noise voltage. However, the signal increases even faster, yielding an improved signal-to-noise ratio with cooling. Figure 90 (noise vs. frequency) shows the dependence on bias.

Capacitance: Capacitance is proportional to area and increases slightly with temperature (see Fig. 91).

Responsivity: See Figs. 82 and 88.

Quantum efficiency: >90 percent quantum efficiency achievable with antireflection coating.

Sensitive area: 0.2 to 600 mm² areas are readily available.

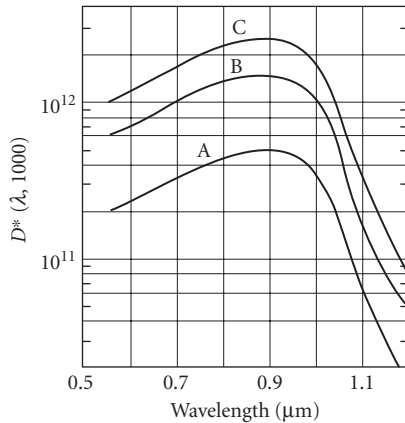


FIGURE 85 D^* versus λ for small-area junction silicon photodiodes: curves A, B, and C correspond to areas of 0.02, 0.2, and 1 cm^2 . The lower D^* for smaller area detector performance is due to amplifier limitations rather than intrinsically poorer D^* , for small-area detectors. (Texas Instruments, *Infrared Devices*, SC-8385-366. Reprinted by permission of Texas Instruments.)

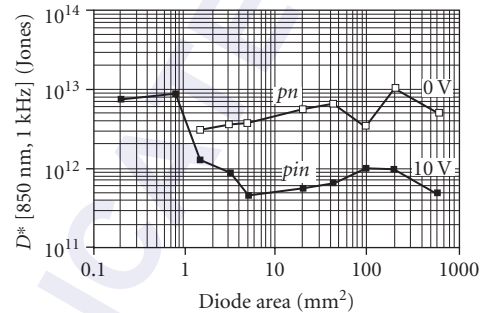


FIGURE 86 D^* as a function of diode area for pn junction silicon photodiodes operated in the photovoltaic mode (0 V) and pin junction diodes operated in the photoconductive or reverse-bias mode (10 V). (UDT Sensors, *Optoelectronic Components Catalog*.)

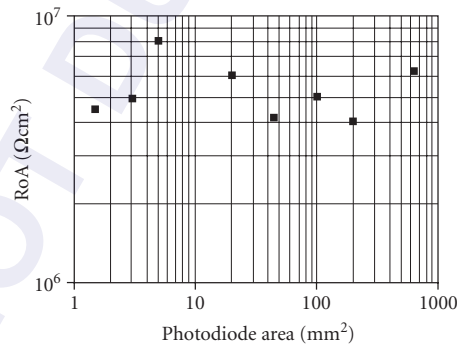


FIGURE 87 The resistance-area product (RoA) at zero bias and 295 K of silicon pn junction photodiodes. The lack of area dependence is evidence that intrinsic properties of the junction, rather than surface effects, are dominant in these devices.

Time constant: Inherently slow for high-sensitivity applications, generally limited by RC (depends directly on device area), but can be limited by carrier diffusion outside the depletion region or by trapping of carriers in deep impurity centers. Typical data for a circuit using a 50- Ω load resistor is illustrated in Fig. 92.

Operating temperature: Ambient, but noise (leakage current) can be reduced by operating at lower temperatures (see Fig. 76 for typical signal and noise vs. temperature).

Uniformity: Typically ± 8 percent across a diode area with a 40- μm focused light spot.

Linearity: 5 percent or better over 10 orders of magnitude flux from 10^{-13} to 10^{-3} W/cm^2 .

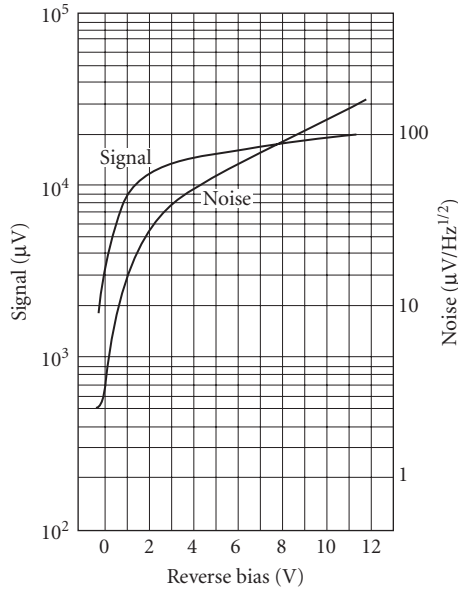


FIGURE 88 Typical *pn* junction signal and noise versus reverse bias ($R_L = 10\text{ M}\Omega$). (Electronuclear Laboratories, Bull. 1053, 1966.)

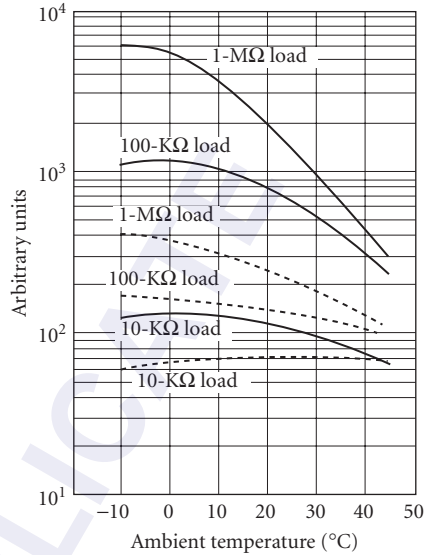


FIGURE 89 Relative signal and noise versus temperature for *pn* junction silicon photodiode at zero bias; — = signal; --- = noise. (Electronuclear Laboratories, Bull. 1052, 1966.)

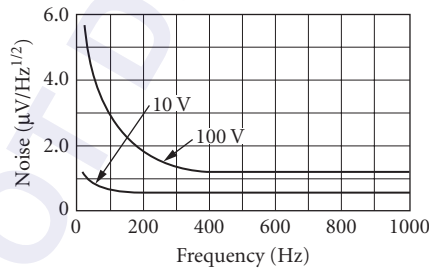


FIGURE 90 Typical *pn* junction and *pin* junction noise-frequency spectrum for different reverse bias ($A = 1 \times 1\text{ mm}$; $R_L = 1\text{ M}\Omega$). (Electronuclear Laboratories, Bull. 1078, 1966.)

Recommended circuit: See Fig. 93. High-impedance FET current-mode amplifier to supply fixed bias voltage, regardless of current.

Stability: See Fig. 20 and section relating to stability. Check with manufacturer.

Manufacturers: Advanced Photonix, EG&G Canada, EG&G Heimann, Edmund Optics, Electro Optical Systems, Electro-Optics Technology, International Radiation Detectors, Janos Technology, Laser Precision Corp, Laser Systems Devices, Melles Griot, Newport/Klinger, Ophir Optronics, Optical Signature, Opto-Electronics, Optometrics, Oriel, Photonic Detectors, RMD, Sapidyne, Scientific Instruments, SEMICOA, Silonex, Spire, UDT Sensors.

Silicon *pin* junction photodiodes The *pin* junction detector is faster but is also somewhat less sensitive than conventional *pn* junction detectors. *PIN* photodiodes have slightly extended red

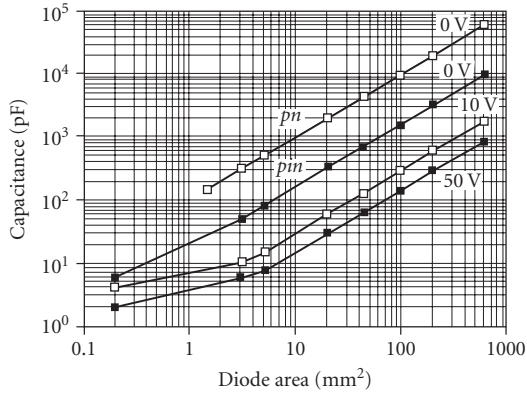


FIGURE 91 Capacitance as a function of detector area for pn junction silicon photodiodes operated in the photovoltaic mode (0 V) and pin junction photodiodes at 0-, 10-, and 50-V reverse bias. The larger depletion width, which is a consequence of the lightly doped “ i ” region in the pin device, gives pin diodes lower capacitance for the same device area. (UDT Sensors, *Optoelectronic Components Catalog*.)

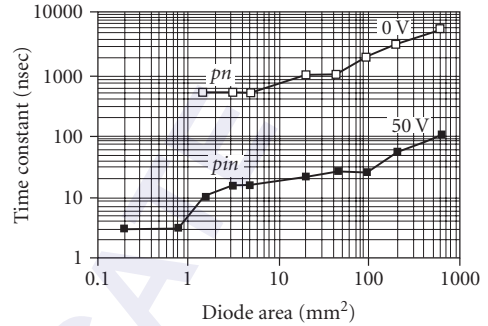
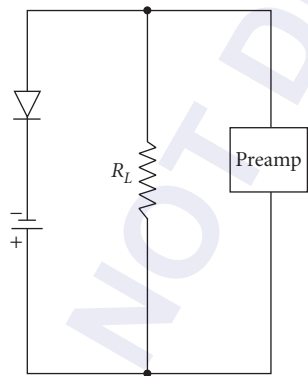
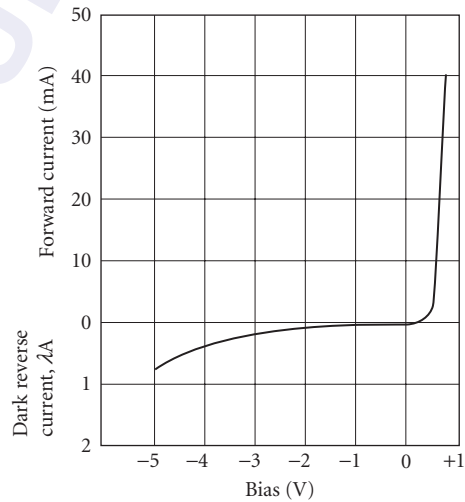


FIGURE 92 Time constant for pn junction silicon photodiodes operated in the photovoltaic mode (0 V) and a pin junction detector in the photoconductive or reverse-bias mode (50 V). A 50- Ω load was used in both cases, which limits sensitivity because of Johnson noise (see also Fig. 83). (UDT Sensors, *Optoelectronic Components Catalog*.)



(a)



(b)

FIGURE 93 pn junction silicon photodiode: (a) recommended circuit; (b) typical electrical characteristics. (Texas Instruments, Bull. SC-8385-366. Reprinted by permission of Texas Instruments.)

response. In the normal pn junction, charge-collection time has a slow and a fast component. The fast component is due to photons absorbed in the depletion layer of the pn junction. Since the electric field in the depletion region is strong, carriers are quickly separated by drift through the electric field across the depletion region. However, photons absorbed deeper in the material, beyond the depletion region, produce carriers which must diffuse to the junction before they are collected, and diffusion times are on the order of a microsecond. This component becomes more significant near the

long-wavelength limit of the spectral response. Application of reverse bias in an ordinary pn junction detector reduces the capacitance, shortening the RC time constant, and increases the width of the depletion layer thereby increasing the fraction of photons absorbed within the high field region and proportionally increasing the fraction of the fast component of the response.

However, the doping level of the ordinary pn junction limits the extent of the depletion layer increase to only 5 to 10 μm at a reverse bias of 50 V (this assumes an abrupt junction with a concentration of $1 \times 10^{15} \text{ cm}^{-3}$). pin detectors incorporate a very lightly doped region between the p - and n -regions that allows a modest reverse bias to form a depletion region the full thickness of the material (500 μm for a typical silicon wafer). Extended red response in a pin device is a consequence of the extended depletion layer width, since longer wavelength photons will be absorbed in the active device region. Unfortunately, the higher dark current collected from generation within the wider depletion layer results in lower sensitivity. Generation of carriers can be minimized by minimizing the concentration of deep-level impurity centers in the detector with careful manufacturing. Operation at lower temperature will also reduce the dark current.

Sensitivity: $D_{pk}^* 1 \times 10^{12}$ Jones for 2-mm² area (depends slightly on bias, see Figs. 86 and 94). For high-speed operation, detectivity is lower (see Fig. 83).

Noise: Depends upon diode area and circuit load resistance. Johnson noise will dominate at low values of load resistance when circuit is optimized for fast response. Preamp noise may also limit. See Fig. 83.

Responsivity: Similar to pn junction. See Fig. 82.

Quantum efficiency: 90 percent quantum efficiency achievable with antireflection coating.

Capacitance: Proportional to detector area. See Fig. 91.

Operating temperature: See Fig. 95.

Time constant: Varies with capacitance (device area); see Fig. 92.

Sensitive area: 0.2 to 600 mm² readily available.

Recommended circuit: Same as for pn junction photodiodes. See Figs. 93 and 96.

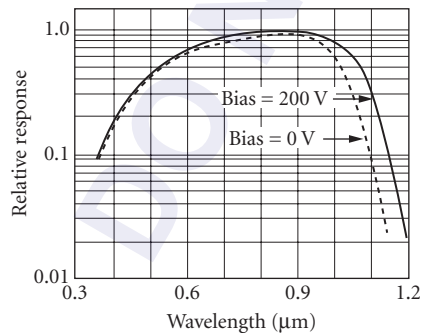


FIGURE 94 Dependence of spectral response on bias for silicon photodiodes. (Electronuclear Laboratories, Bull. 1076, 1968.)

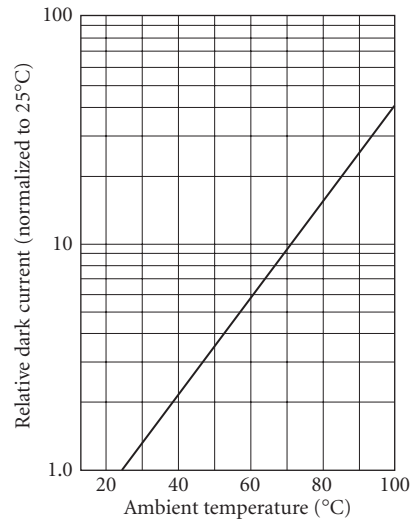


FIGURE 95 Relative dark current versus temperature for pin junction silicon photodiodes—bias 100 V. (Electronuclear Laboratories, Bull. 1076, 1969.)

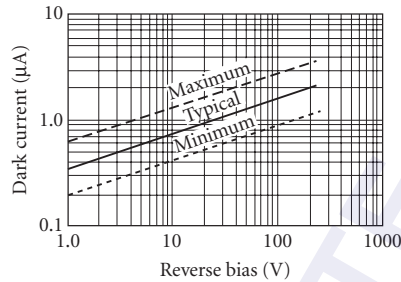


FIGURE 96 Typical dark current versus bias for *pin* silicon photodiode ($A = 1 \times 1 \text{ mm}$). (*Electronuclear Laboratories, Bull. 1078, 1968.*)

Stability: See Fig. 20 and section relating to stability. Check with manufacturer.

Manufacturers: Same as for *pn* junction photodiodes.

UV- and blue-enhanced photodiodes Blue- and UV-enhanced photodiodes may improve the quantum efficiency by 50 to 100 percent over standard photodiodes in the blue and UV spectral region. The quantum efficiency of ordinary *pn* and *pin* junction photodiodes degrades rapidly in the blue and UV spectral regions. This is because the high absorption coefficient of silicon at these wavelengths causes the photocarriers to be generated within the heavily doped *p*- (or *n*-) type contact surface where the lifetime is short due to the high doping and/or surface recombination. Blue- and UV-enhanced photodiodes optimize the response at short wavelengths by minimizing near-surface carrier recombination. This can be achieved by using very thin and highly graded *p* (or *n* or metal Schottky) contacts, by using lateral collection to minimize the percentage of the surface area which is heavily doped, and/or passivating the surface with a fixed surface charge to repel minority carriers from the surface. These devices typically have quartz windows or UV-transmissive glass, compatible with good transmission into the UV spectrum. The user should be aware that UV and higher energy radiation in particular can alter the fixed charge conditions in the surface region of silicon and other detectors (typically in the surface oxide) which can cause the detector performance to drift and/or be unstable (see Fig. 20).

Sensitivity: See Figs. 82 and 97; D^*_{pk} $3\text{--}5 \times 10^{12}$ Jones for diodes with areas of $1\text{--}100 \text{ mm}^2$ at $V_R = 0$, $R_L = 40 \text{ M}\Omega$.

Quantum efficiency: Same as *pn* junction photodiodes, but enhanced in the UV and blue regions by 50 percent or more (see Fig. 82).

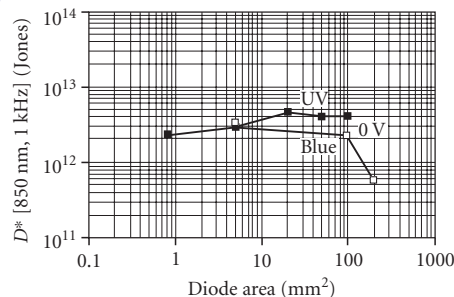


FIGURE 97 D^* as a function of diode area for blue- and UV-enhanced silicon photodiodes operating in the photovoltaic mode (0-V bias). (*UDT Sensors, Optoelectronic Components Catalog.*)

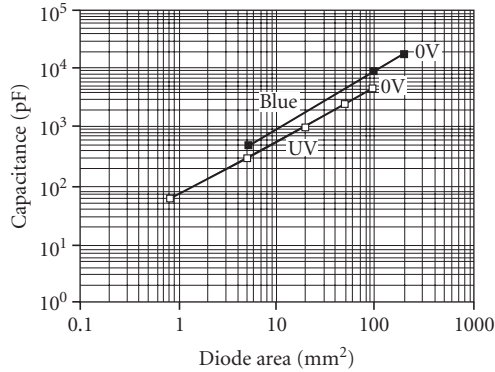


FIGURE 98 Capacitance as a function of detector area for blue- and UV-enhanced silicon photodiodes operated in the photovoltaic mode (0 V). The capacitance per unit area is close to that of pn junction photodiodes shown in Fig. 91. (UDT Sensors, *Opto electronic Components Catalog*.)

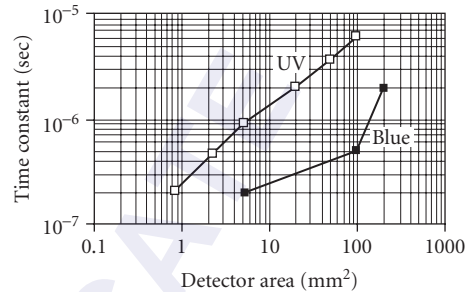


FIGURE 99 Time constant for UV- and blue-enhanced silicon photodiodes as a function of detector area at zero-volts bias. A 50- Ω load was used in both cases. (UDT Sensors, *Optoelectronic Components Catalog*.)

Responsivity: see Fig. 82.

Capacitance: Comparable to pn junction photodiodes (See Fig. 98).

Operating temperature: Ambient.

Time constant: Dependent upon device type. Increases with device area; 200 ns to 6 μ s for areas 1 to 100 mm² (see Fig. 99).

Sensitive area: 1 to 200 mm² readily available.

Recommended circuit: Same as pn junction photodiodes.

Stability: See Fig. 20 and section relating to stability. Check with manufacturer.

Manufacturers: See list for silicon pn junction photodiodes.

Silicon avalanche photodiodes The avalanche photodiode, is especially useful where both fast response and high sensitivity are required. Whereas normal photodiodes become Johnson- or thermal-noise-limited when used with a low-impedance load resistor for fast response, avalanche photodiodes make use of internal multiplication, associated with reverse breakdown in the pn junction in order to keep the detector noise above the Johnson-noise level. [Because the response time is usually RC-limited, small load resistors (often 50 Ω) are used to achieve fast signal response.] However, as the load resistor is decreased, the detector noise voltage decreases in direct proportion, whereas the Johnson noise decreases only as the square root of the load resistor. Thus, the detector noise voltage can become lower than the Johnson noise for load resistance values smaller than a critical value. With an APD device, lower values of load resistors can be used without reaching the critical value because the internal gain boosts the detector noise voltage.

Stable avalanche or multiplication is made possible by a guard-ring construction using n^+pp^+ , Schottky- nn^+ , or $n^+p\pi p^+$ structure; beveled pin structure (see Fig. 100); mesa structures; or other structures which prevent surface breakdown.²⁵ However, very careful bias control is essential for stable performance. An optimum gain exists below which the system is limited by receiver noise and above which shot noise dominates receiver noise and the overall noise increases faster than the signal (Fig. 101).

In addition to fast-response applications, avalanche photodiodes are useful whenever amplifier noise is limiting, for example, small-area devices. Signal-to-noise-ratio improvements of one to two orders of magnitude over a nonavalanche detector can be achieved.

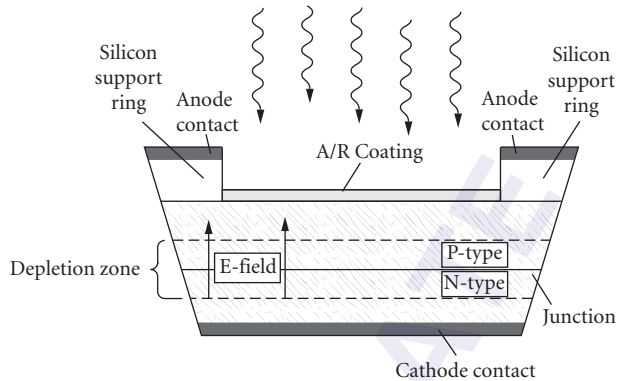


FIGURE 100 Cross section of a beveled-edge silicon avalanche photodiode. The beveled edge prevents early breakdown. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

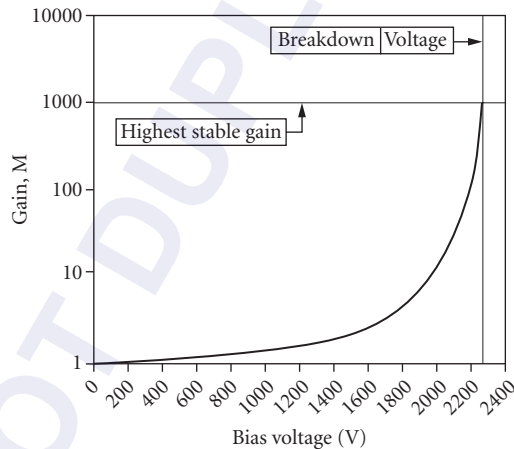


FIGURE 101 Gain as a function of reverse bias can reach 1000. This operating point is very close to breakdown and requires careful bias control. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

APDs are sometimes used in combination with scintillator crystals such as CsI to detect high-energy radiation in the range of 10 to 1000 keV.

Sensitivity: $D^* 3\text{--}5 \times 10^{13}$ Jones (see Fig. 102).

Noise: Function of detector area. As gain increases, noise (dark current) increases (see Fig. 103). Optimum gain is where avalanche noise equals system noise. Thus, optimum gain is a function of system noise.

Responsivity: Photocurrent is the product of the incident optical power in watts, wavelength in micrometers, and quantum efficiency (η) divided by 1.24 and multiplied by the avalanche gain M .

$$I_{\text{photo}} = M(P \eta \lambda / 1.24) \quad (\text{See Figs. 101 and 104})$$

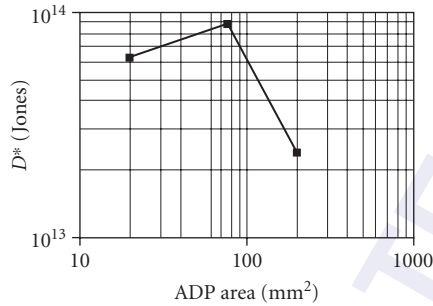


FIGURE 102 D^* for three silicon avalanche photodiodes shown as a function of diode area. (Advanced Photonix. Avalanche Photodiode Catalog.)

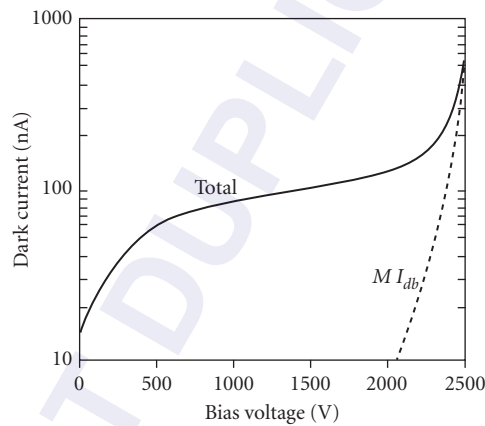


FIGURE 103 Dark current as a function of reverse bias for a 16-mm-diameter APD. At low-bias surface dark current dominates, but avalanche-multiplied bulk dark current increases rapidly as the gain increases. (Advanced Photonix. Avalanche Photodiode Catalog.)

Quantum efficiency: Typically 85-percent peak (see Fig. 104).

Capacitance: Depends on bias and area (see Fig. 105).

Sensitive area: 20 to 200 mm².

Series resistance: Depends on area; typical values are 40 Ω for 5-mm diameter to 5 Ω for 16-mm diameter

Time constant: See Fig. 106.

Recommended circuit: Requires a filtered high-voltage dc supply that itself must have very low noise and a load resistor. The output may be ac- or dc-coupled. (See Fig. 107.)

Operating temperature: 40 to 45°C.

Stability: Exposure to UV or high-energy radiation may affect dark current. See Fig. 20 and section relating to stability. Check with manufacturer.

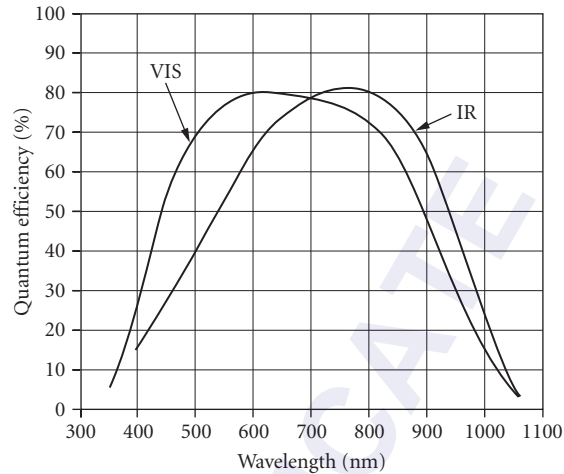


FIGURE 104 APD quantum efficiency at high gain. Adjustment of the oxide deposited on the surface produces two different curves. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

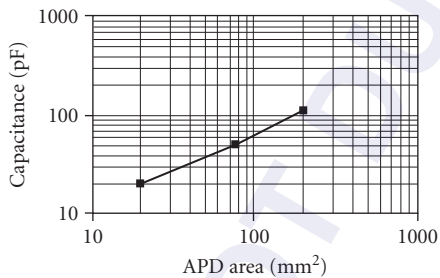


FIGURE 105 Capacitance for three silicon avalanche photodiodes shown as a function of diode area. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

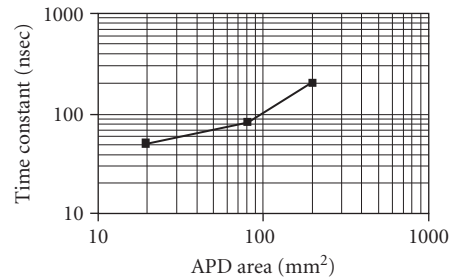


FIGURE 106 Time constant for three silicon avalanche photodiodes shown as a function of diode area. (*Advanced Photonix. Avalanche Photodiode Catalog.*)

Manufacturers: Advanced Photonix, Devar, EG&G Judson, EG&G Vactec, EG&G Canada, Edmund Optics, Electro-Optical Systems, Hamamatsu, Janos Technology, Newport/Klinger, Opto-Electronics (Ontario), Oriel, Photonic Detectors, Photonic Packaging Technologies, RMD, Texas Optoelectronics, Thorn EMI Electron Tubes.

InGaAs Indium gallium arsenide detectors have been developed for optimum performance with fiber-optic communications at 1.3 and 1.55 μm . This detector material has a direct bandgap and represents one of several compound semiconductor alloy systems specially developed for photodetectors. In the case of this alloy of two group III-V chemical compound semiconductors, the ratio of InAs to GaAs controls the spectral cutoff, allowing the detector to be optimized for a particular wavelength. InGaAs detectors have generally been specialized for high-speed applications with

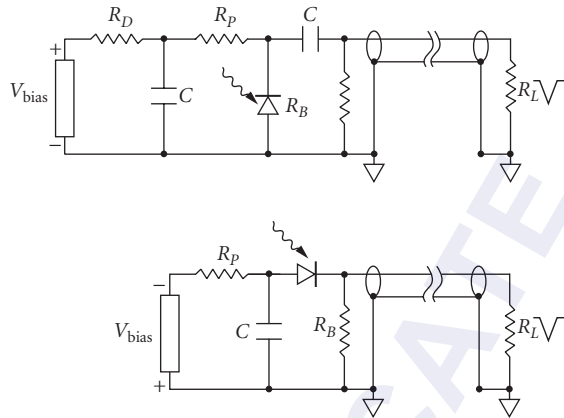


FIGURE 107 AC- and DC-coupled APD circuits with a filter following the power supply and a coaxial cable on the output. (Advanced Photonix, Avalanche Photodiode Catalog.)

optimum sensitivity since these performance factors can drive a fiber-optic system throughput and cost. For this reason, available devices include

- PIN photodiodes
- Avalanche photodiodes

The significance of these devices to fiber-optic applications is reflected in the number of vendors who sell integrated packages of InGaAs photodetectors combined with preamplifiers and fiber-optic pigtailed.

The InGaAs alloy system allows the spectral response to be tailored to longer wavelengths than the quartz fiber-optic bands and devices with cutoffs of 2.2 and 2.6 μm are also available.

InGaAs pin photodiode

Sensitivity: D^* mid- 10^{12} Jones for 1.67- μm cutoff; $D^* \approx 1 \times 10^{12}$ Jones for 2.2- μm cutoff, and $D^* \approx 5 \times 10^{11}$ Jones for 2.6- μm cutoff.

Responsivity: 0.85 to 0.95 A/W in the range of 1.3 to 1.55 μm .

Quantum efficiency: For 1.67- μm cutoff, see Fig. 108.

Dark current: See Fig. 109.

Capacitance: 0.7 to 1.2×10^4 pF/cm for 1.7- μm cutoff; 2.5×10^4 pF/cm for 1.85- μm cutoff; 3×10^4 pF/cm for 2.15- μm cutoff; 5×10^4 pF/cm for 2.65- μm cutoff. See also Fig. 110 for bias dependence.

Time constant: Varies with resistance-capacitance time (see Fig. 111). Since capacitance depends upon reverse bias, the time constant varies proportionally (see Fig. 110 for dependence of capacitance on bias).

Size: 0.05 to 3-mm diameter.

Recommended circuit: Standard photodiode options; zero bias for best sensitivity, reverse bias for maximum speed.

InGaAs avalanche photodiode

Sensitivity: $D^* \approx 5 \times 10^{11}$ Jones for 1.7- μm cutoff. In the fiber-optics industry, the sensitivity is also given in power units of dBm. Figure 112 compares InGaAs *pin*, APD, and Ge APD sensitivities.

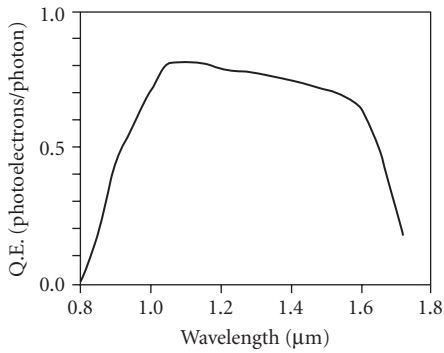


FIGURE 108 Spectral dependence of quantum efficiency for an InGaAs detector having a cutoff of 1.67 μm . (*Sensors Unlimited, data sheet.*)

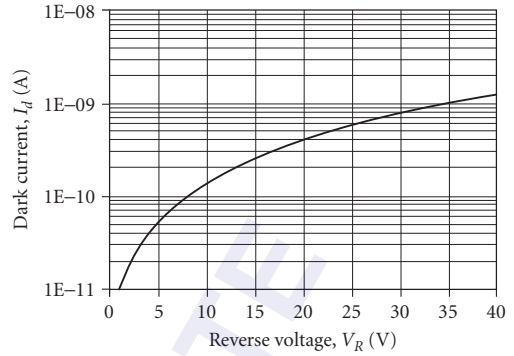


FIGURE 109 Dark current as a function of reverse-bias voltage for a 60- μm -diameter InGaAs detector having a cutoff of 1.67 μm . (*Fermionics, InGaAs Photodiodes.*)

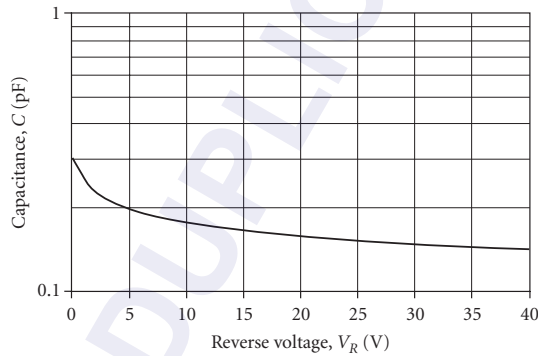


FIGURE 110 Capacitance as a function of reverse-bias voltage for a 60- μm -diameter InGaAs detector having a cutoff of 1.67 μm . (*Fermionics, InGaAs Photodiodes.*)

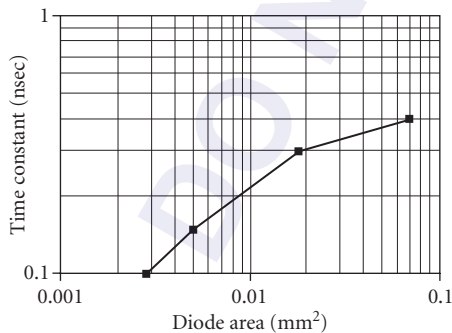


FIGURE 111 Time constant for small InGaAs *pin* photodiodes as a function of diode area measured with a 50- Ω load.

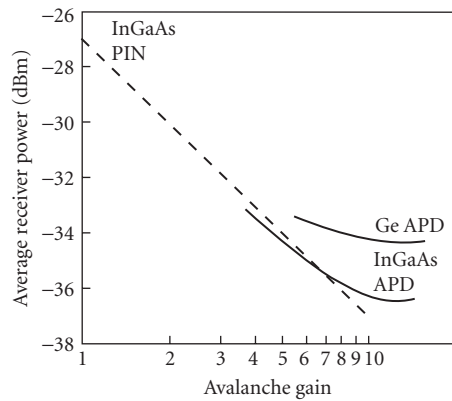
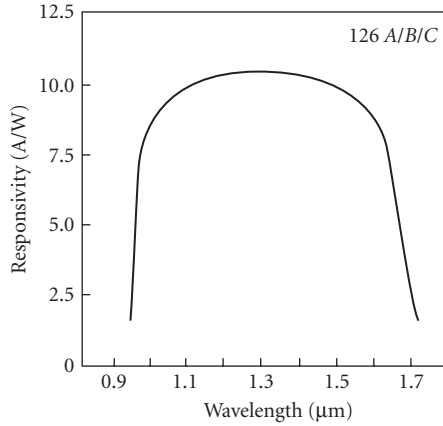


FIGURE 112 APD receiver sensitivity. Typical receiver sensitivity at a receiver rate of 1.7 Gbit/s and $\lambda = 1.3 \mu\text{m}$ for an InGaAs *pin*, Ge APD, and an InGaAs APD. (*AT&T, 126A/B, CASTROTEC InGaAs.*)



Note: Responsivity = (chip quantum efficiency) × gain × λ (μm)/1.24. The minimum chip quantum efficiency is 80%, and the minimum pigtail coupling efficiency is 90%.

FIGURE 113 Responsivity for an InGaAs avalanche photodiode versus wavelength for avalanche gain of 12. (AT&T, 126A/B, C ASTROTEC InGaAs.)

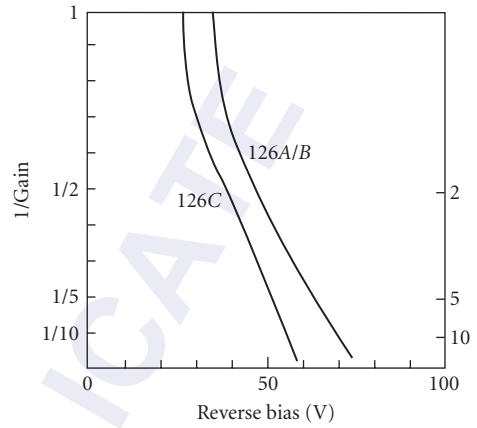


FIGURE 114 Inverse of avalanche gain for an InGaAs avalanche photodiode versus reverse bias. (AT&T, 126A/B, C ASTROTEC InGaAs.)

Spectral response: 1.0 to 1.65 μm; see Fig. 113.

Responsivity: 8 to 10 A/W typical.

Avalanche gain: Critically depends upon reverse bias, see Fig. 114.

Capacitance: At gain of 12, 7.5×10^4 pF/cm².

Bandwidth: Up to 3 GHz; see Figs. 115 and 116.

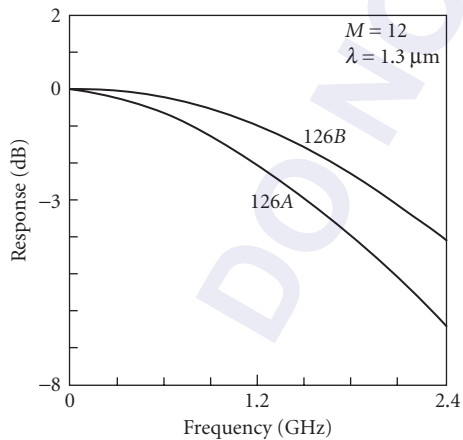


FIGURE 115 Frequency response of an InGaAs APD (126A/B). (AT&T, 126A/B, C ASTROTEC InGaAs.)

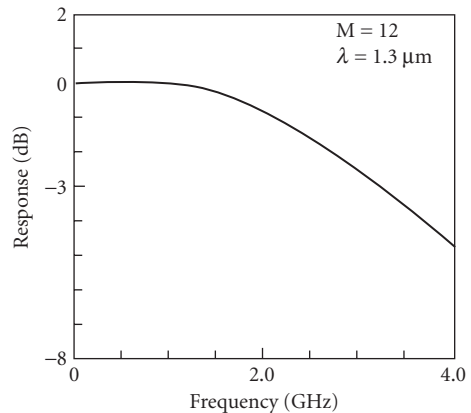


FIGURE 116 Frequency response of InGaAs APD (126C). (AT&T, 126A/B, C ASTROTEC InGaAs.)

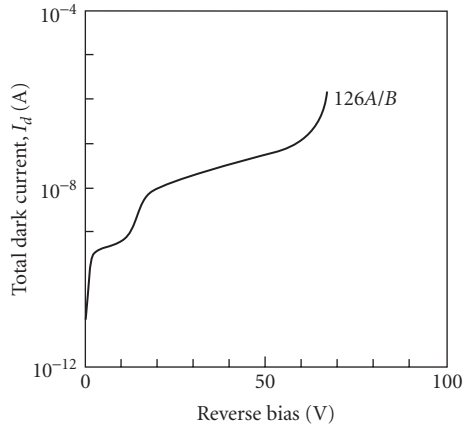


FIGURE 117 Dark current versus reverse bias of InGaAs APD (126A/B). (AT&T, 126A/B, C ASTROTEC InGaAs.)

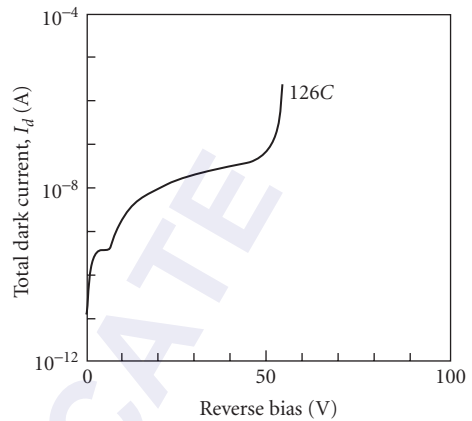


FIGURE 118 Dark current versus reverse bias of InGaAs APD (126C). (AT&T, 126A/B, C ASTROTEC InGaAs.)

Size: 0.04 to 0.5-mm diameter.

Dark current: Dependent upon reverse bias, device structure, and temperature. See Figs. 117, 118, and 119.

Recommended circuit: See Figs. 107 and 120.

Manufacturers: Advanced Photonix, AT&T, EG&G Canada, Edinburgh Instruments, Edmund Optics, Electro-Optical Systems, Electro-Optics Technology, Emcore, Epitaxx, Fermionics, GCA Electronics, Germanium Power Devices, Hamamatsu, New England Photoconductor, New Focus,

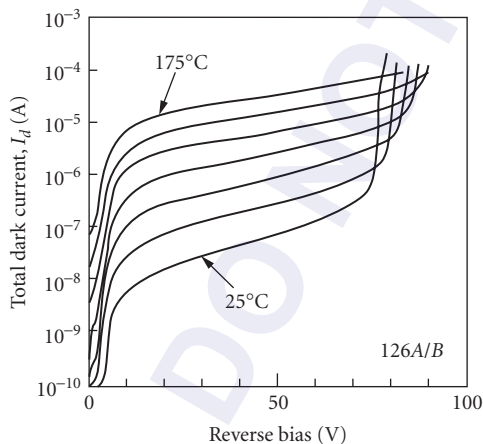


FIGURE 119 Dark current versus voltage of InGaAs APD as a function of temperature at 25°C increments. Note: The temperature dependence of the 126C dark current is the same as the 126A/B. (AT&T, 126A/B, C ASTROTEC InGaAs.)

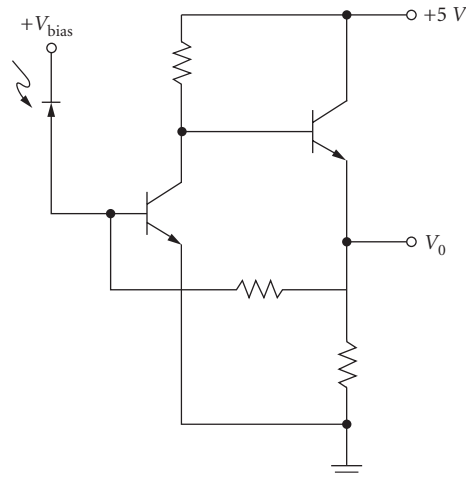


FIGURE 120 Bipolar transimpedance amplifier for InGaAs avalanche photodiode. (AT&T, 126A/B, C ASTROTEC InGaAs; *Optical Fiber Communications*, John M. Senior, © 1985, Prentice-Hall; ISBN-0-13-638248-7.)

Newport, North Coast Scientific, Opto-Electronics, Ortel, Photonic Detectors, Goodrich Sensors Unlimited, Spire, Swan Associates, Telcom Devices, Teledyne Judson Technologies, UDT Sensors.

Ge Germanium intrinsic photodetectors are similar to intrinsic silicon detectors but offer spectral response out to 1.5 to 1.8 μm . *PN* junction photodiodes offer submicrosecond response or high sensitivity from the visible region to 1.8 μm . Zero bias is generally used for high sensitivity and large reverse bias for high speed. As in the case of silicon, germanium has an indirect bandgap and soft spectral cutoff. The previous discussion on silicon detectors applies in general, with the exception that blue- and UV-enhanced devices are not relevant to germanium detectors. Germanium detectors, because of their narrower bandgap, have higher leakage currents at room temperature, compared to silicon detectors. Detector impedance increases about an order of magnitude by cooling 20°C below room temperature. Thus, performance can improve significantly with thermoelectric cooling or cooling to liquid nitrogen temperature.

As with silicon, the device structure and bias configuration can affect spectral response and rise time. Three detector types are available:

- *pn* junction
- *pin* junction
- Avalanche photodiode

Germanium pn and pin

Sensitivity: D^* (peak, 300 Hz, room temperature) $> 2 \times 10^{11}$ Jones, increases significantly with cooling by thermoelectric cooler or liquid nitrogen. (See Figs. 121, 122, and 123.)

Quantum efficiency: > 50 percent with antireflection coating.

Noise: See Figs. 124 and 125.

Responsivity: 0.9 A/W at peak wavelength. See Fig. 121.

Capacitance: Lower for *pin* structure compared with *pn* diode. See Fig. 126.

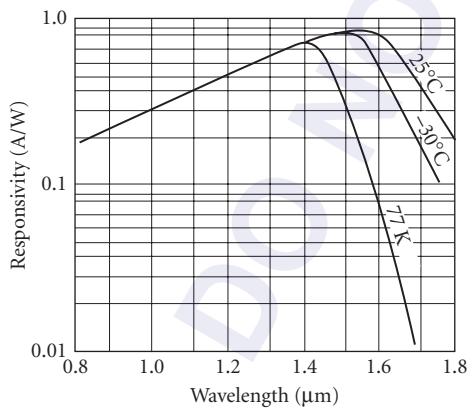


FIGURE 121 Spectral response for a germanium *pn* junction photodiode at three temperatures. (Teledyne Judson Technologies, *J16 germanium photodiodes*, 2008.)

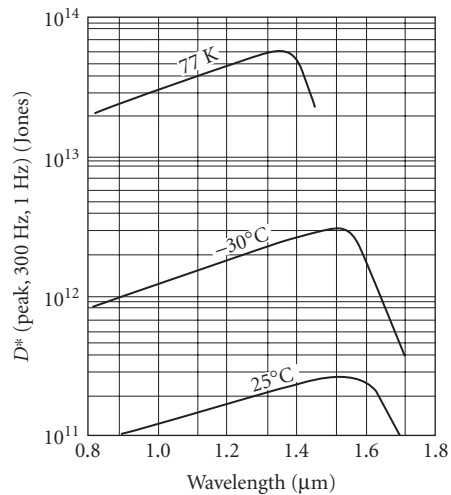


FIGURE 122 D^* as a function of wavelength for a germanium *pn* junction photodiode at three temperatures. (EG&G Judson, *Infrared Detectors*, 1994.)

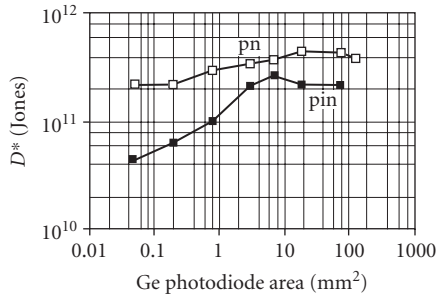


FIGURE 123 D^* for germanium pn and pin photodiodes shown as a function of diode area. (EG&G Judson, *Infrared Detectors*, 1994.)

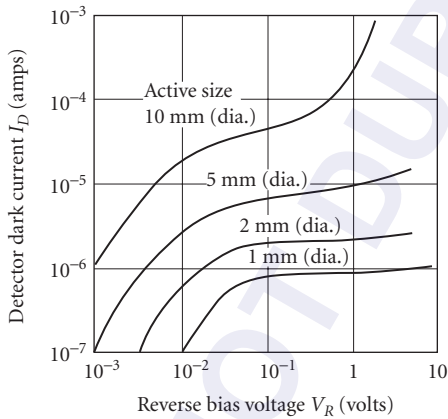


FIGURE 125 Dark current as a function of reverse bias for germanium pn junction photodiodes of different diameters at 25°C. (Teledyne Judson Technologies, *J16 germanium photodiodes*, 2008.)

Time constant: PIN diodes provide faster response. See Fig. 127.

Sensitive area: 0.25 to 13-mm diameter standard.

Operating temperature: Ambient, TE-cooled, or liquid nitrogen.

Profile: ± 2 percent across active area at 1.3 μm .

Linearity: Excellent over 10 orders of magnitude. See Fig. 128.

Recommended circuit: See previous section on silicon photodiodes.

Manufacturers: Edinburgh Instruments, Electro-Optical Systems, Judson, Electro-Optical Systems, Fastpulse Technology, Germanium Power Devices, Infrared Associates, Newport, North Coast Scientific, Opto-Electronics, Oxford Instruments, Scientific Instruments.

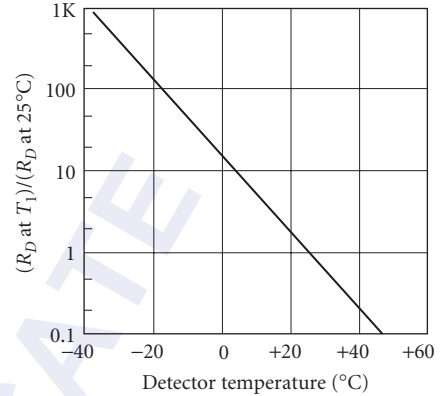


FIGURE 124 Ratio of resistance at temperature T to the resistance at 25°C for a germanium pn junction photodiode. (Teledyne Judson Technologies, *J16 germanium photodiodes*, 2008.)

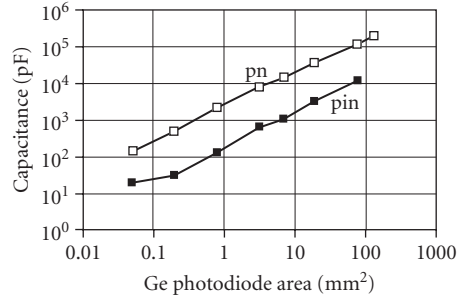


FIGURE 126 Capacitance for germanium pn and pin photodiodes shown as a function of diode area. (EG&G Judson, *Infrared Detectors*, 1994.)

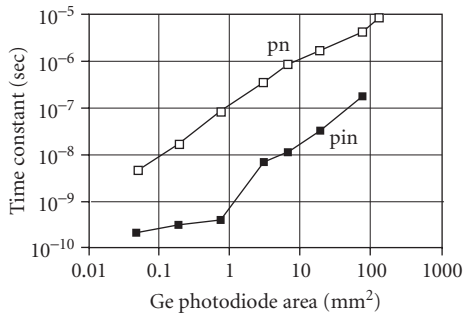


FIGURE 127 Time constant for germanium *pn* and *pin* photodiodes shown as a function of diode area. (EG&G Judson, *Infrared Detectors*, 1994.)

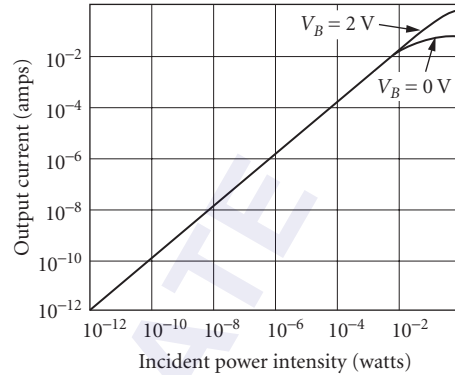


FIGURE 128 Linearity of a germanium *pn* junction photodiode. (EG&G Judson, *Infrared Detectors*, 1994.)

Germanium avalanche photodiode The germanium avalanche photodiode is similar to the silicon APD but has lower optimum gain, longer cut-off wavelength (1.7 μm), and higher leakage current. Germanium APDs combine the sensitivity of a Ge *pn* photodiode and the speed of a *pin* Ge photodiode.

Sensitivity: $D^* 2 \times 10^{11}$ Jones at 30 MHz for a diode with area of 5×10^{-2} mm^2 or about the same as for a *pn* Ge diode with the same area, and about a factor of 4 higher than a *pin* Ge photodiode (compare with Figs. 122 and 123). D^* depends on gain.

Gain: See Fig. 129.

Dark current: See Fig. 130.

Capacitance: 2 pF at 20-V reverse bias for 100- μm diameter, 8 pF at 20-V reverse bias for 300- μm diameter.

Quantum efficiency: 60 to 70 percent at 1.3 μm .

Responsivity: Photocurrent is the product of the incident optical power in watts, wavelength in micrometers, and quantum efficiency (η) divided by 1.24 and multiplied by the avalanche gain M . $I_{\text{photo}} = M(P\lambda\eta/1.24)$. (See Figs. 121 and 129).

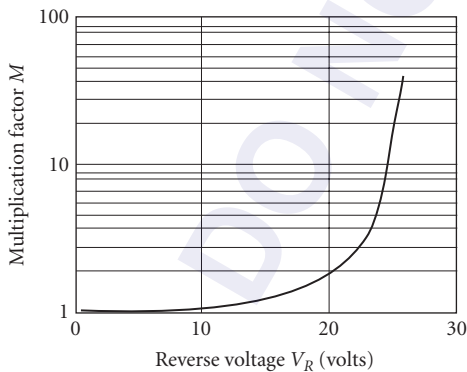


FIGURE 129 Gain as a function of reverse bias for germanium avalanche photodiode. (EG&G Judson, *Infrared Detectors*, 1994.)

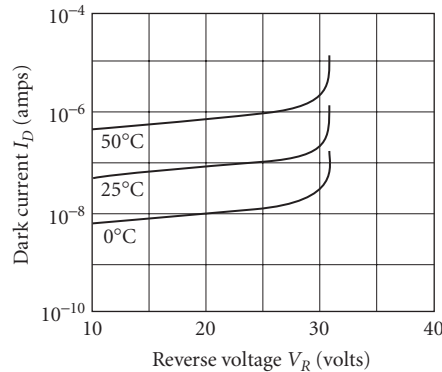


FIGURE 130 Dark current as a function of reverse bias for germanium avalanche photodiode. (EG&G Judson, *Infrared Detectors*, 1994.)

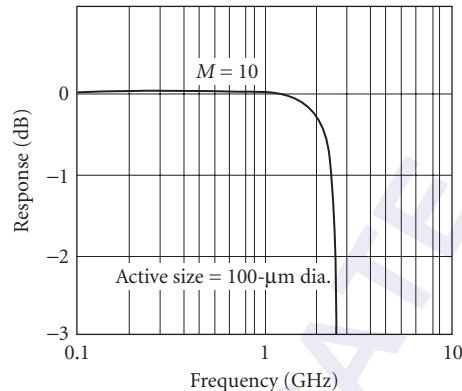


FIGURE 131 Frequency response of a germanium avalanche photodiode at two gain operating points. (EG&G Judson, *Infrared Detectors*, 1994.)

Operating temperature: Ambient or TE-cooled.

Time constant: 0.2 ns for 100- μm diameter; 0.3 ns for 300- μm diameter; both at 1.3 μm , $M = 10$ and with $R_L = 50 \Omega$. See Fig. 131.

Sensitive area: 100- and 300- μm diameter standard.

Recommended circuit: See circuit recommended for Si APD.

Manufacturers: Edmund Optics, Metrotech, North Coast Scientific, Teledyne Judson Technologies.

PbS Photoconductive lead sulfide was one of the earliest and most successful infrared detectors. Even today it is one of the most sensitive uncooled detectors in the 1.3- to 3- μm spectral region. With cooling, PbS sensitivity is competitive with other detectors out to about 4.2 μm ; however, its response time is slow.

Many PbS characteristics can be varied by adjusting the chemistry of the deposition process and/or the post-deposition heat treatment. These characteristics include spectral detectivity, time constant, resistance, and upper limit of operating temperature.²⁶ PbS is generally made by chemical reaction of Pb acetate and thiourea, except for high-temperature (373 K) applications, where evaporation is used. The material is deposited as a thin film (1 to 2 μm thick) on a variety of substrates, such as sapphire. With photolithographic processing, small sensitive areas can be made with comparatively high D^* values.

PbS may be tailored for ambient or room-temperature operation (ATO), intermediate or thermoelectrically cooled operation (ITO), and low-temperature or nitrogen-cooled operation (LTO). They are manufactured differently for particular temperature ranges, as shown in Table 4.

Sensitivity: $D^* 1.5 \times 10^{11}$ Jones at 295 K. See Figs. 132 to 137.

Responsivity: Depends on detector area, bias, resistance, and operating temperature (see Figs. 138 and 139).

Quantum efficiency: Generally limited by incomplete absorption in the thin film to 30 percent as estimated from blip D^* values.

Noise: Dominated by $1/f$ noise at low frequencies. See Figs. 135 and 140.

Time constant: Can be varied in manufacturing. Typical values are 0.2 ms at 295 K, 2 to 5 ms at 193 and 77 K. See Fig. 139.

Sensitive area: Typical sizes are square elements with dimensions of 0.5, 1, 2, and 5 mm on a side.

TABLE 4 PbS Performance Characteristics

	Typical operating temperature, K			
	350	273 (ATO)	193 (ITO)	77 (LTO)
Sensitivity	†	Figs. 132, 138	Figs. 133, 136, 138	Figs. 134, 138
$D^* (\lambda_{\max})/D^* (500 \text{ K})$		105	55	17
Noise, $\text{V}/\text{Hz}^{1/2}$		Fig. 140	Fig. 140	Fig. 140
Dark resistance, $\text{M}\Omega/\text{sq}$	<0.3	<2	<10	<20
Time constant, μs^\ddagger	50	100–500	5000	3000

†At 350 K, cutoff wavelength moves into $\sim 2.4 \mu\text{m}$, with $D^* (\lambda_{\max}) \approx 10^{10}$ Jones.

‡These are typical values; the time constant can be adjusted over two orders of magnitude in fabrication, but D^* is affected.

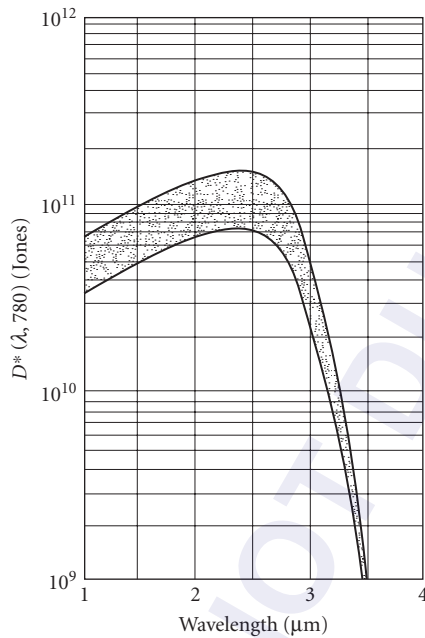


FIGURE 132 Range of spectral detectivities for PbS (ATO) at 295 K; 2π FOV, 295-K background. (Santa Barbara Research Center.)

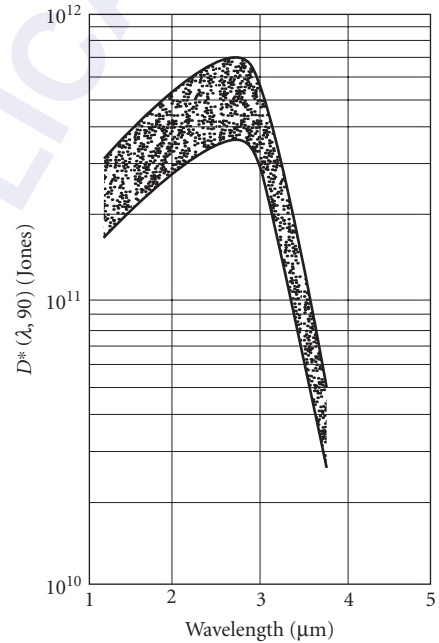


FIGURE 133 Range of spectral detectivities for PbS (ITO) at 193 K; 2π FOV, 295-K background. (Santa Barbara Research Center.)

Capacitance: <1 pF (limited by mounting configuration).

Recommended circuit: Standard photoconductor.

Stability: Exposure to visible and/or UV radiation can induce instability and drift. Stability will recover with storage in the dark, or by baking.

Sensitivity profile: Uniform within 10 percent.

Linearity: Excellent over broad range 10^{-8} to 10^{-3} W.

Manufacturers: Alpha Omega Instruments, Cal-Sensors, Edmund Optics, Electro-Optical Systems, Hamamatsu, New England Photoconductor, Teledyne Judson Technologies, OptoElectronics, Orielt.

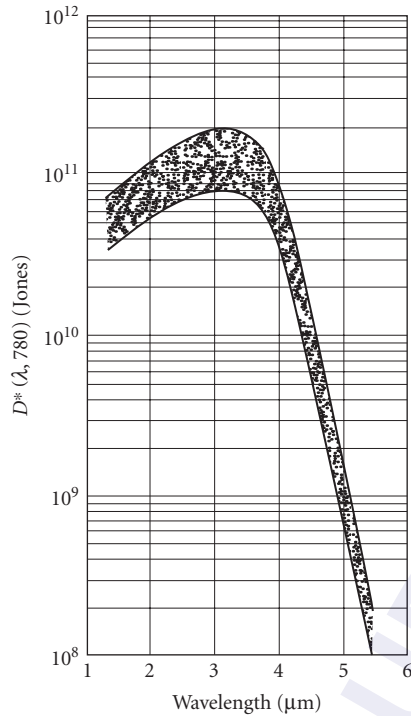


FIGURE 134 Range of spectral detectivities for PbS (LTO) at 77 K; 2π FOV, 295-K background. (Santa Barbara Research Center.)

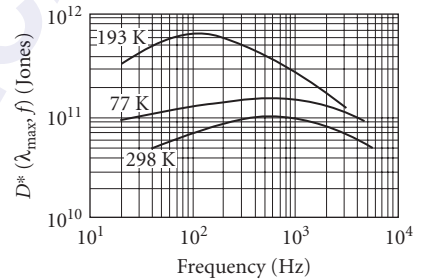


FIGURE 135 Example of detectivity vs. temperature for PbS detectors at various operating temperatures; 2π FOV, 295-K background. (Santa Barbara Research Center.)

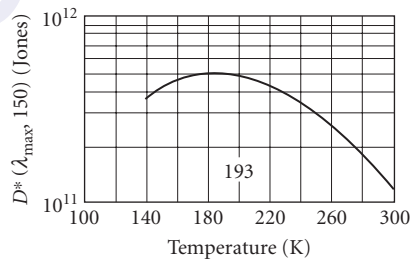


FIGURE 136 Example of detectivity versus temperature for PbS (ITO) detectors; 2π FOV, 295-K background. (Santa Barbara Research Center.)

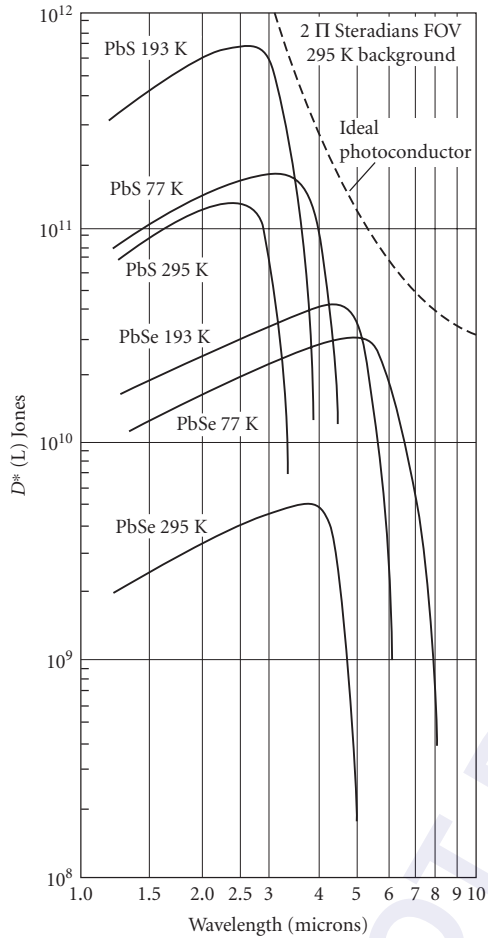


FIGURE 137 D^* versus wavelength for PbS and PbSe detectors operating at temperatures ranging between 77 K and 295 K. (CAL-SENSORS, *Infrared Detectors*.)

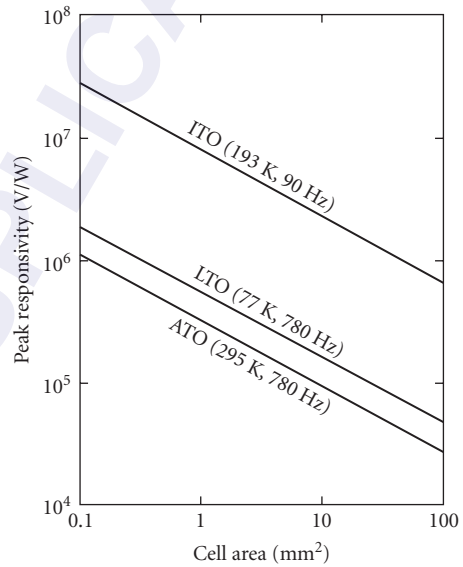


FIGURE 138 PbS typical peak responsivity versus cell area (actual values range within a factor of two of these shown.)

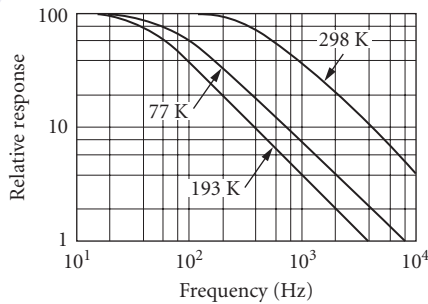


FIGURE 139 Example of signal versus frequency for PbS detectors. (Santa Barbara Research Center.)

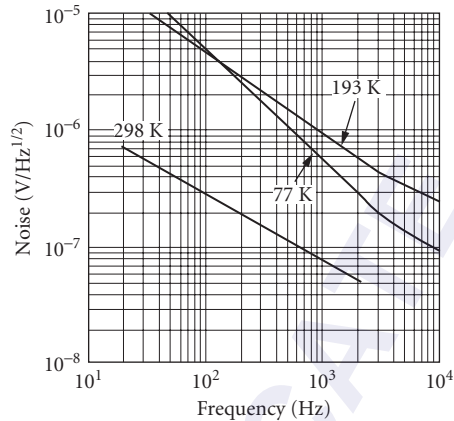


FIGURE 140 Example of noise versus frequency for PbS detectors. (Santa Barbara Research Center.)

InAs (Photovoltaic) InAs detectors are single-crystal, intrinsic, direct-bandgap photovoltaic devices for use in the 1- to 4- μm region (spectral cutoff varies with temperature). At room temperature, InAs provides good sensitivity and submicrosecond response times. At 195 K, InAs performance equals or better the sensitivity of any other detector in the 1 to 3.5- μm region. Devices with sapphire immersion lenses are available to increase signal responsivity for operation at higher temperatures where the detector is thermal-noise-limited. Compared to PbS and PbSe detectors, InAs has very little low frequency ($1/f$) noise if operated in the photovoltaic mode.

Sensitivity: D^* (peak) varies from 1.2×10^9 Jones at 295 K to 6×10^{11} Jones at 77 K. See Fig. 141.

Quantum efficiency: Maximum of about 75 percent without antireflection coating.

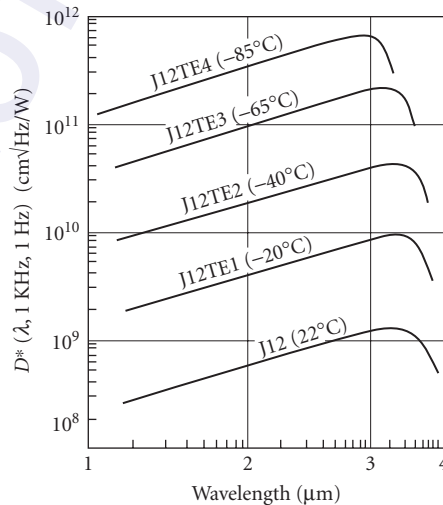


FIGURE 141 D^* versus wavelength for InAs detector operating at temperatures ranging between 77 K and 295 K. (Teledyne Judson Technologies, 2008.)

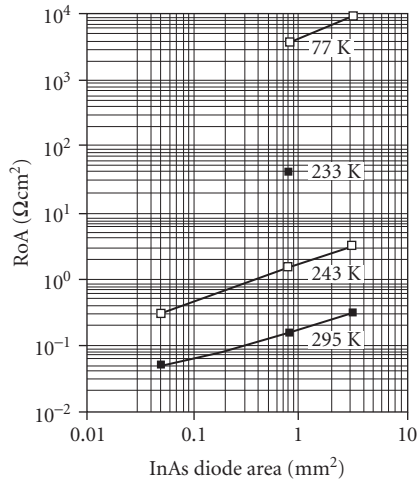


FIGURE 142 R_oA of InAs photodiodes shown as a function of diode area. Lower impedance per unit area for smaller devices indicates that these devices are surface-leakage-limited. (Santa Barbara Research Center.)

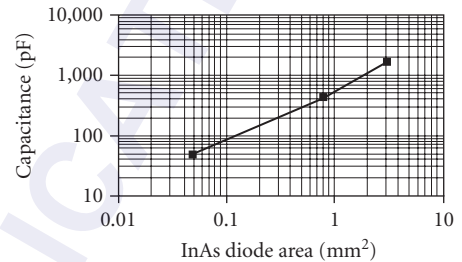


FIGURE 143 Capacitance of InAs photodiodes shown as a function of diode area. Capacitance will not change appreciably with temperature. (Teledyne Judson Technologies, 2008.)

Noise: Low impedance tends to make preamplifier noise dominate at room temperature; background limited (for 300 K background) at operating temperatures below 200 K.

Time constant: Less than 0.5 μ s at all temperatures when low values of load resistor are used to reduce the RC time constant.

Responsivity: 0.5 to 1.25 A/W at peak.

Dynamic resistance: See Fig. 142.

Diode capacitance: See Fig. 143.

Sensitive area: Standard sizes 0.25 to 2-mm diameters.

Operating temperatures: 77 to 300 K.

Linearity: Anticipated to be very good over many decades.

Sensitivity profile: ± 15 percent.

Recommended circuits:

Open circuit: PV InAs detectors with areas less than 2×10^{-2} cm require no bias when operated and can be connected directly into the input stage of amplifier (capacitor ensures elimination of dc bias from amplifier) (Fig. 144a).

Transformer: Useful when using InAs at zero bias, particularly at room temperature where diode impedance is low (Fig. 144b).

Reversed bias: At temperatures greater than 225 K considerable gain in impedance and responsivity is achieved by reverse-biasing (Fig. 144c).

Fast response: To utilize the short intrinsic time constant, it is sometimes necessary to load the detector to lower the RC of the overall circuit (reverse bias will also lower detector capacitance) (Fig. 144d).

Manufacturers: Electro-Optical Systems, Hamamatsu, Teledyne Judson Technologies.

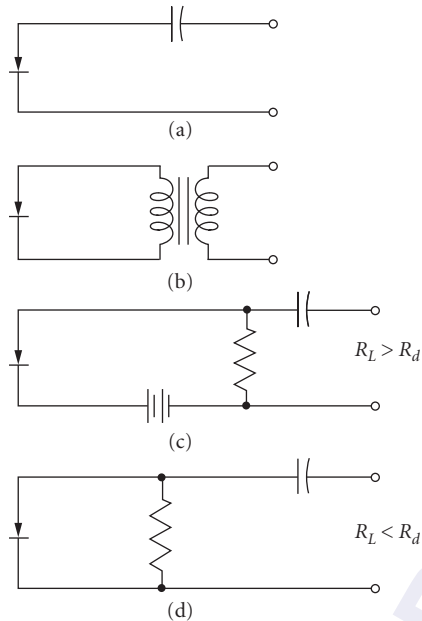


FIGURE 144 Recommended circuits for InAs detectors: (a) open circuit, (b) transformer, (c) reversed bias, (d) fast response.

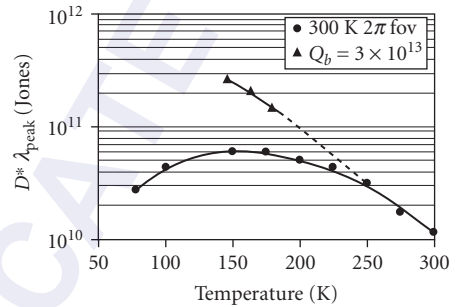


FIGURE 145 D^* of PbSe as a function of a function of temperature for two background flux conditions: high background of 2π field of view and reduced background of 3×10^{13} photons/cm²/s. D^* at the higher background flux reaches a maximum around 160 K because the background noise increases at lower temperatures due to the increase in long-wavelength spectral response of the detector. (Santa Barbara Research Center.)

PbSe Lead selenide is an intrinsic, thin-film photoconductor, whose long-wavelength spectral response and speed of response exceeds that of PbS. At room temperature, PbSe has peak D^* which can exceed 1×10^{10} Jones with a spectral cutoff out to 4.4 μm . At liquid-nitrogen temperature, InSb offers twice the D^* , largely because PbSe offers response out to 7 μm at 77 K, considerably longer than InSb. However, for intermediate temperatures, from 180 K to room temperature, PbSe offers competitive D^* combined with moderately fast response.²⁶ PbSe technology has made a significant advance in the past decade in some vendors being able to reproducibly make high-performance detectors.

Sensitivity: $D^* \approx 1 \times 10^{10}$ Jones at 300 K, increases with cooling (see Figs. 137 and 145). D^* is limited by $1/f$ noise at low frequencies (see Fig. 146).

Response: Figure 147 shows responsivity in amperes per watt for a high-quality detector with a length of 0.016 cm and width of 0.024 cm. Responsivity in volts per watt is obtained by multiplying A/W data by resistance (see Fig. 148). Responsivity will vary inversely with detector length (see Figs. 149 and 150).

Noise: Figure 151 shows the noise as a function of temperature for a detector with a length of 0.016 cm and width of 0.024 cm. Noise as a function of frequency is shown in Fig. 152.

Resistance: Figure 148 shows the resistance as a function of temperature for a detector with a length of 0.016 cm and width of 0.024 cm.

Capacitance: 1 pF (limited by mounting configuration).

Time constant: See Figs. 153 and 154. Time constant will be longer when detector is operated in reduced background flux condition.

Stability: Exposure to visible and/or UV radiation can induce instability and drift. Stability will recover with storage in the dark at room temperature.

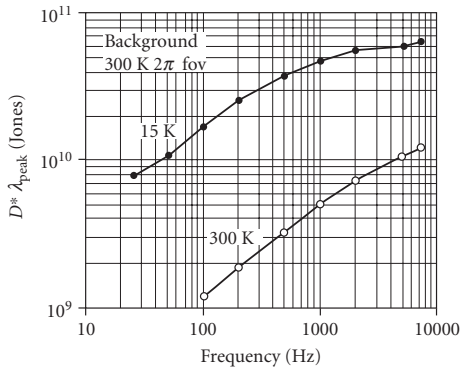


FIGURE 146 D^* of PbSe as a function of frequency for two temperatures. PbSe has considerable $1/f$ noise which reduces D^* at lower frequencies. (Santa Barbara Research Center.)

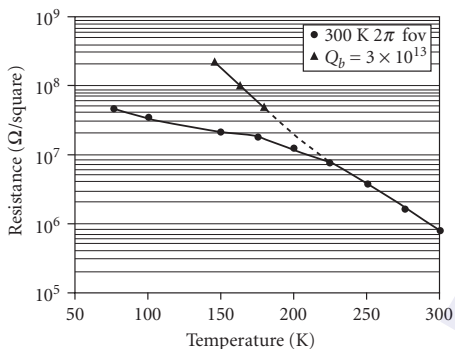


FIGURE 148 Resistance (Ω /square) of PbSe thin-films as a function of temperature for two background flux levels. At any temperature, the absolute value can be varied by altering the manufacturing process in chemical deposition and/or heat treatment. (Santa Barbara Research Center.)

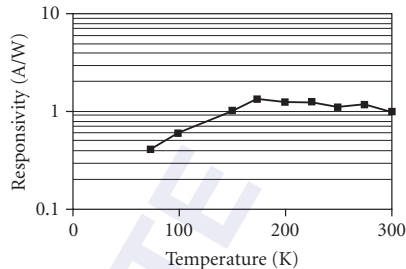


FIGURE 147 Responsivity in A/W of PbSe thin-film photoconductive detectors as a function of temperature for a high background flux level. Multiply by detector resistance to get V/W. (Santa Barbara Research Center.)

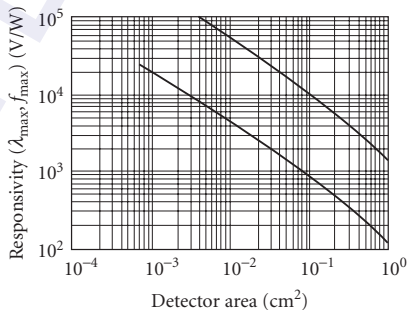


FIGURE 149 Expected range of peak responsivities versus detector size, typical PbSe (ATO) infrared detectors. (Santa Barbara Research Center.)

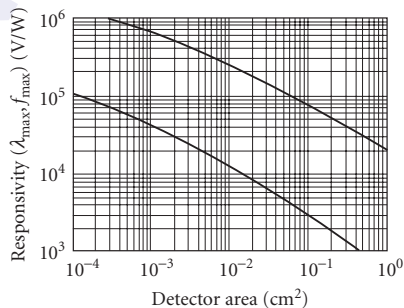


FIGURE 150 Expected range of peak responsivities versus detector size, typical PbSe (ITO and LTO) infrared detectors. (Santa Barbara Research Center.)

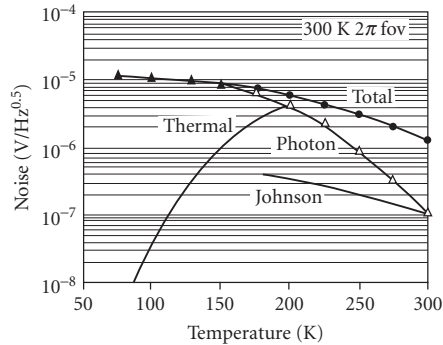


FIGURE 151 Noise voltage (per square root of bandwidth) of PbSe thin-film photoconductive detectors as a function of temperature for a high background flux level. Photon noise is dominant below 200 K. Thermal noise is dominant at higher temperatures. Total noise levels are well above typical preamplifier noise. (Santa Barbara Research Center.)

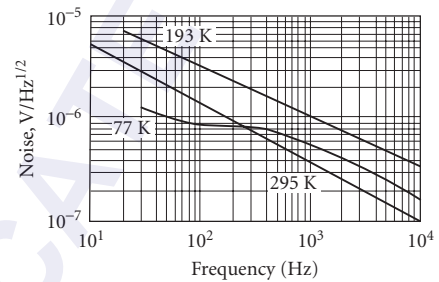


FIGURE 152 Example of noise versus frequency for PbSe detectors (ATO, ITO and LTO types) (1×1 mm). (Santa Barbara Research Center.)

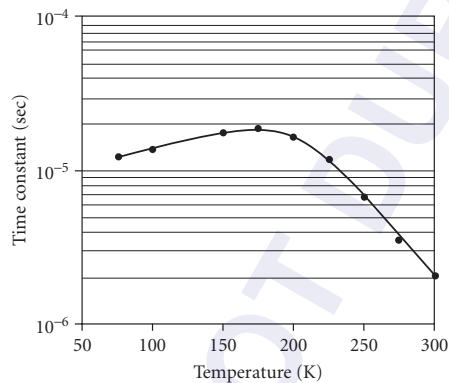


FIGURE 153 Time constant of PbSe thin-film photoconductive detectors as a function of temperature. (Santa Barbara Research Center.)

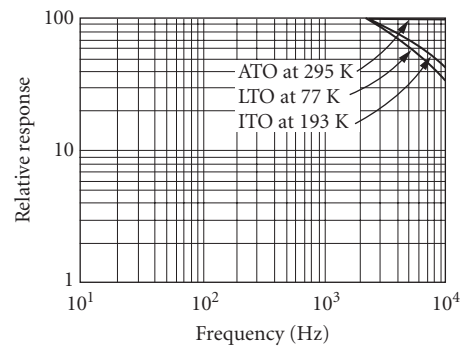


FIGURE 154 Example of signal versus frequency for PbSe detectors (ATO, ITO and LTO). (Santa Barbara Research Center.)

Recommended circuit: Standard photoconductor.

Operating temperature: 77 to 300 K.

Manufacturers: Cal-Sensors, Edmund Optics, Electro-Optical Systems, Hamamatsu

InSb Historically, indium antimonide material has been used for at least four different radiation detector types, two of which, the photoconductive and photoelectromagnetic types, are no longer widely used. We discuss here the intrinsic photovoltaic device. [The very far infrared bolometer (InSb bolometer) was previously discussed in Sec. 24.7.

At 77 K, InSb photodiodes offer background limited sensitivity at medium-to-high background flux conditions in the 1- to 5.5- μm spectral range. At lower temperatures, they provide sensitive detectors at low background flux levels such as in astronomy applications, but with a slightly shortened long-wavelength cutoff. Operation is possible up to as much as 145 K, but because the spectral

response increases with increasing temperature, the detector impedance drops rapidly leading to significant thermal noise.

Sensitivity: Spectral response out to 5.5 μm at 77 K (see Fig. 155). $D^* 1 \times 10^{11}$ Jones, increases with reduced background flux (narrow field of view and/or cold filtering) as illustrated in Fig. 156.

Quantum efficiency: ~ 60–70 percent without antireflection coating. >90 percent with antireflection coating.

Noise: Background current limited over wide range of background flux at 77 K (see Fig. 157).

Time constant: <1 μs .

Responsivity: 3 A/W at 5 μm without antireflection coating.

Noise equivalent power (NEP): Frequency dependence is shown in Fig. 158 for three detector sizes.

Capacitance: Typically 0.05 F/cm².

Impedance: Top-grade detectors have $1\text{--}5 \times 10^6 \Omega\text{cm}^2 R_0A$ product at 77 K, at zero bias, and without background flux.

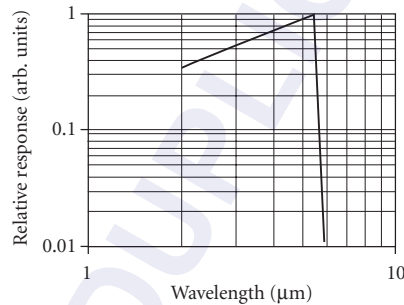


FIGURE 155 Relative spectral response per watt of an InSb photodiode without antireflection coating. The direct bandgap results in a sharp spectral cutoff. (Santa Barbara Research Center.)

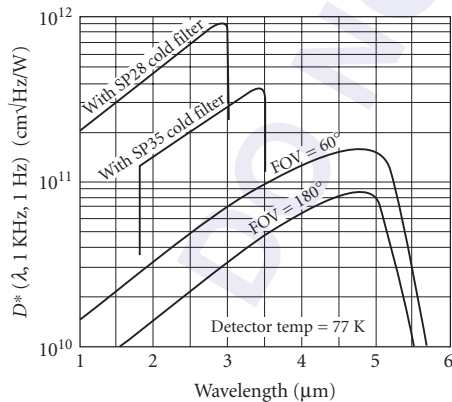


FIGURE 156 D^* as a function of wavelength for an InSb detector operating at 77 K. (Teledyne Judson Technologies, 2008.)

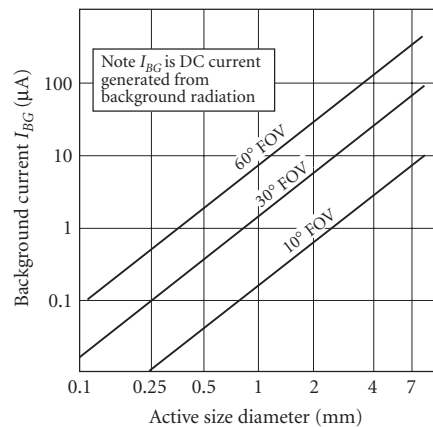


FIGURE 157 Background current as a function of photodiode area for InSb detectors operating at 77 K, shown at three values of the detector field of view. (Teledyne Judson Technologies, 2008.)

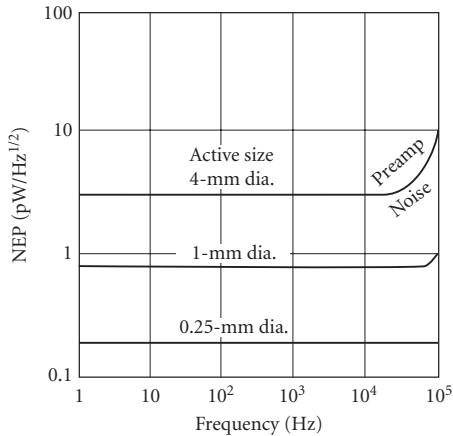


FIGURE 158 NEP as a function of frequency for three sizes of InSb photodiodes. (Teledyne Judson Technologies, 2008.)

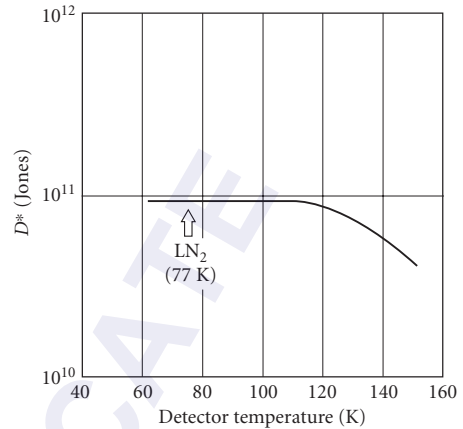


FIGURE 159 InSb photodiode D^* as a function of operating temperature between 77 K and 150 K, for a 2π (180°) FOV. (Teledyne Judson Technologies, 2008.)

Sensitive area: 0.04×0.04 -mm square to 1×5 -mm rectangle; 0.25 to 10-mm diameter.

Operating temperature: Normally 77 K; InSb can be used up to approximately 145 K (see Fig. 159).

Linearity: Linear to ~ 1 -mW/cm² flux.

Sensitivity profile: ± 15 percent or better.

Stability: Devices from some vendors are subject to “flashing,” where exposure to visible or UV flux causes a change in the insulating surface charge thereby causing a change in the diode impedance. The detector typically recovers at room temperature.

Recommended circuit: Same as for Si and Ge photodiodes; zero or reverse bias in combination with a load resistor and low-noise preamplifier. Low-impedance load resistor can be used for obtaining fast response, with consequences of reduced sensitivity.

Manufacturers: L3 Cincinnati Electronics, Edinburgh Instruments, New England Photoconductor, Teledyne Judson Technologies, Electro-Optical Systems, Hamamatsu, Infrared Associates.

Ge:Au Gold-doped germanium detectors are relatively fast single-crystal p-type impurity-doped photoconductors for the 2- to 9- μ m region. Although not the most sensitive detector anywhere in its range of spectral sensitivity, Ge:Au offers respectable sensitivity over a broad spectral region using liquid nitrogen cooling. Sensitivity can be improved by a factor of 2.5 by operating at $T < 65$ K (pumped liquid nitrogen or other cryogen). At these temperatures, Ge:Au becomes background limited.

Sensitivity: See Figs. 160 to 163.

Quantum efficiency: Dependent on wavelength, detector geometry (absorption thickness), antireflection coating, and enclosure (integration chamber can increase absorption). $D_{\lambda_{pk}}/D_{500K}^* = 2.7$ (see Fig. 164).

Noise: See Fig. 165.

Time constant: < 50 ns with full D^* [shorter response times (< 2 ns) can be tailored by heavy concentration of compensating (n -type) dopant and suitable bias circuit (see circuit discussion to follow). Heavy compensation increases resistance, and hence the incoherent signal-to-noise ratio becomes limited by the thermal noise of the load (typically a factor of 2 degradation in the signal-to-noise ratio). Quantum efficiency, however, is not significantly altered, so that a high compensation concentration does not hurt the coherent-detection signal-to-noise ratio].

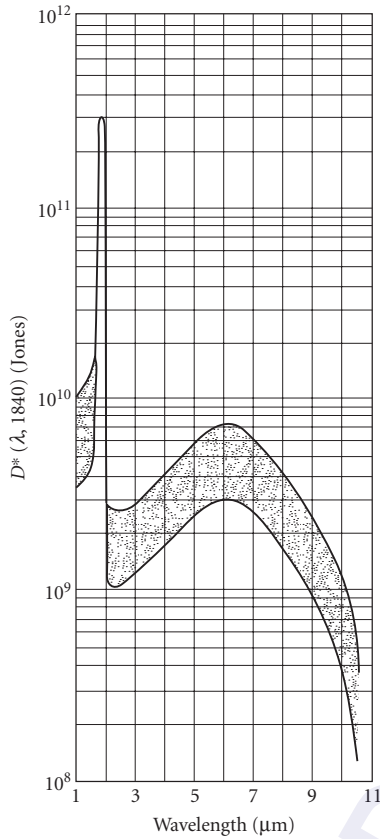


FIGURE 160 D^* versus λ for Ge:Au; $T = 77$ K, 2π FOV; 295-K background. (Santa Barbara Research Center.)

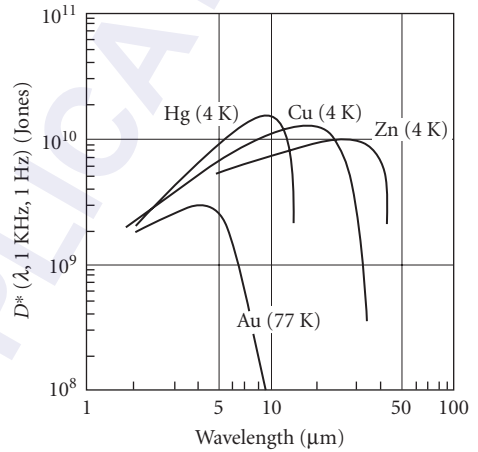


FIGURE 161 D^* as a function of wavelength for extrinsic germanium detectors doped with Au, Hg, Cu, and Zn, for a 300-K 2π (180°) FOV background flux. (EG&G Judson, *Infrared Detectors*, 1994.)

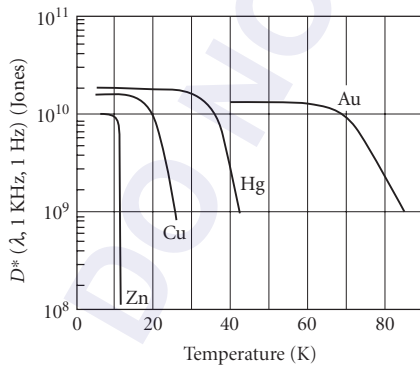


FIGURE 162 D^* as a function of operating temperature for extrinsic germanium detectors doped with Au, Hg, Cu, and Zn for a 300-K 2π (180°) FOV background flux. (EG&G Judson, *Infrared Detectors*, 1994.)

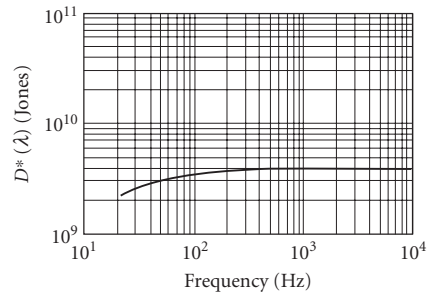


FIGURE 163 Typical D^* versus frequency ($T = 77$ K). (Santa Barbara Research Center.)

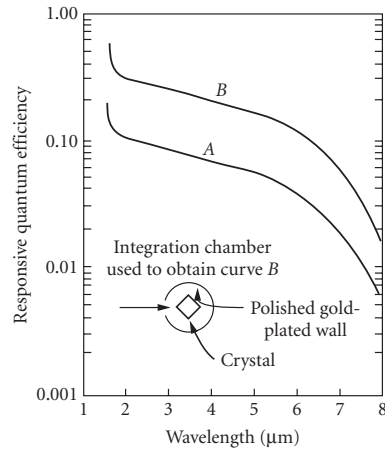


FIGURE 164 Quantum efficiency versus λ for Ge:Au ($T = 78$ K). (Santa Barbara Research Center, internal report.)

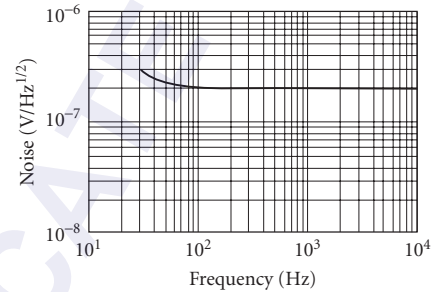


FIGURE 165 Typical noise spectrum for Ge:Au ($T = 77$ K, $A = 1 \times 1$ mm). (Santa Barbara Research Center.)

Responsivity: Dependent upon bias and geometry, typical values are 0.1 to 0.2 A/W at 77 K. Multiply by detector resistance to get V/W.

Dark resistance: Varies with background flux and effective quantum efficiency (see previous quantum efficiency discussion), range may be 20 k Ω to 5 M Ω , or much greater under very low background flux conditions if adequately cooled to limit thermally activated conductivity. (Also see previous time-constant discussion.)

Capacitance: Depends on device geometry and mounting, typically <1 pF.

Sensitive area: 1 to 5-mm diameter standard.

Operating temperature: < 85 K (normally 77 K, but see Fig. 162).

Recommended circuit: Standard photoconductive. See Fig. 166.

Manufacturers: No suppliers are presently known.

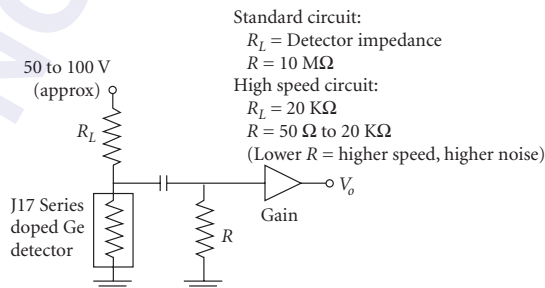


FIGURE 166 Basic operating circuit for extrinsic germanium detectors doped with Au, Hg, Cu, and Zn, for a 300-K 2π (180°) FOV high background flux. If the detector is operated in very low background flux conditions, the detector impedance can become very high. Cooled JFET ($T > \approx 50$ K) or PMOS buffer amplifiers can be helpful in impedance matching under these conditions. (EG&G Judson, *Infrared Detectors*, 1994.)

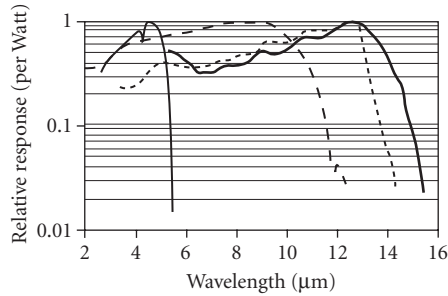


FIGURE 167 Relative spectral response per watt at 80 K for photoconductive HgCdTe detectors with antireflection coating. The curves are normalized to unity at peak value. The spectral cutoff can be adjusted by varying the ratio of HgTe to CdTe in the alloy. (Santa Barbara Research Center.)

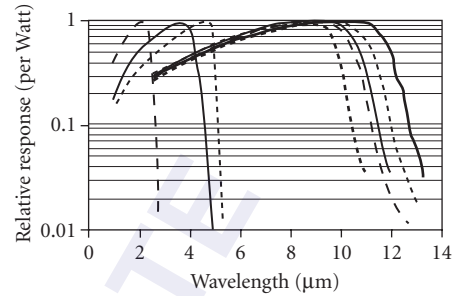


FIGURE 168 Relative spectral response per watt at 80 K for photovoltaic HgCdTe detectors without antireflection coating. The curves are normalized to unity at peak value. The spectral cutoff can be adjusted by varying the ratio of HgTe to CdTe in the alloy. (Santa Barbara Research Center.)

HgCdTe, HgZnTe, HgMnTe, etc. Mercury cadmium telluride is a direct-bandgap compound alloy semiconductor, made of chemical group II and VI elements, whose peak sensitivity at a particular temperature can be adjusted from 1 to 30 μm by varying the ratio of HgTe to CdTe (see Figs. 167 and 168). In addition to HgCdTe, other combinations of chemical groups II and VI elements can be used to produce similar variable spectral cutoff compound alloys, including HgZnTe, HgMnTe, HgCdZnTe, etc. For almost all purposes, HgCdTe will be as good as any other II-VI alloy detector, so we will speak here about it exclusively, but provide some data from the other alloys noted above. Both photoconductive (PC) and photovoltaic (PV) HgCdTe detectors are available for background-limited, high-speed, intrinsic photon detection in the SWIR, MWIR, and LWIR regions. Photoconductive devices with VLWIR response out to 25 μm are also available. HgCdTe detectors can be used at room temperature, with TE cooling, and at 77 K and lower temperatures. Sensitivity generally increases with cooling, depending upon the spectral cutoff and background flux. MWIR, LWIR, and VLWIR spectral range devices are generally operated at 77 K or lower temperature for maximum sensitivity, depending upon the background flux. The PC mode is advantageous when cooling is limited, since the thermal noise of a photoconductor increases less rapidly than a photodiode as the temperature is raised. Photoconductive devices with response out to 25 μm can be usefully operated at liquid nitrogen temperature and are popular for IR spectroscopy for this reason.

HgCdTe photoconductors are fabricated from thin ($\approx 10\text{-}\mu\text{m}$) single-crystal slices or epitaxial layers with metal contacts at each end of the element (see Fig. 8). They are low-impedance devices with 15 to 2000 Ω/square , depending upon the alloy composition, carrier concentration, operating temperature, background flux, and surface treatment. Photoconductor time constants at 77 K may be $\approx 2 \mu\text{s}$ for devices having a 12- μm cutoff, with longer time constants for shorter cutoffs, and shorter time constants for higher operating temperatures. In the case of small detector elements, the time constant may be reduced with increasing bias voltage because photoexcited carriers will be transported to the electrical contacts where they recombine. The spectral noise characteristics of PC HgCdTe typically exhibit $1/f$ noise out to a range of 50 Hz to 1 kHz or more, the value depending upon the detector quality, long-wavelength response, operating temperature, and background flux. White noise levels range from less than 1 $\text{nV}/\sqrt{\text{Hz}}$ (where preamplifier noise may then dominate), up to 20 $\text{nV}/\sqrt{\text{Hz}}$, depending upon detector quality, size, geometry, applied bias, temperature, and alloy composition. Photoconductive HgCdTe detectors are typically antireflection coated with a quarter-wave ZnS film, giving a peak quantum efficiency in the range of 85 to 90 percent, although this figure is only indirectly measured because the PC gain can be much greater than unity. Without antireflection coating, the quantum efficiency is typically 70 percent, limited by the optical index of ≈ 4 .

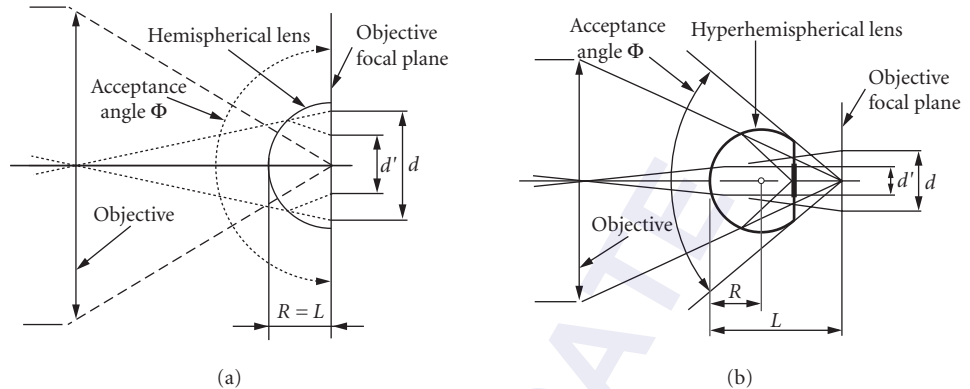


FIGURE 169 Schematic of an optically immersed HgCdZnTe detector with (a) hemispherical lens and (b) hyperhemispherical lens. Dimensions are summarized in Table 5. (Vigo Systems.)

TABLE 5 Summary of the Dimensions Depicted in Fig. 169, and Their Impact on D^* (Vigo Systems)

Parameter	Hemisphere	Hyperhemisphere
Distance, L	$L = R$	$L = R(n + 1)$
d/d'	n	n^2
$D^*_{\text{imm}}/D^*_{\text{non-imm}}$	n	n^2
Acceptance angle, Φ	$\Phi = 180^\circ$	$\Phi = 2 \arcsin(1/n)$
F/#	0.5	1.55

n = index of refraction (approx. 3.3 for GaAs and 2.7 for CdTe)

A class of HgCdTe detectors is offered for detection at TE-cooled and room temperature which are optically “immersed” with a hemispherical or hyperhemispherical lens of Ge, CdTe, GaAs, or other high-index material (see Fig. 169 and Table 5). The lens increases the effective area of the detector without increasing the detector noise, provided the noise is dominated by thermal rather than photon noise as is the case for minimal cooling. The lens must be in intimate contact with the detector surface ($\ll 1 \mu\text{m}$ spacing) to avoid total internal reflection of off-axis rays at the lens-detector interface. Immersed detectors offer up to a factor of n^2 (where n is the optical index) increased detector signal, which can mean an increase in D^* by the same factor for a thermal-noise-limited device. Operation of LWIR PC detectors at TE-cooled and room temperature is generally accompanied by increased $1/f$ noise which dominates out to higher frequencies.

Photovoltaic HgCdTe detectors ideally offer $\sqrt{2}$ higher D^* than detectors operating in the PC mode. Diodes are made in both n^+p and p^+n polarities, depending upon the manufacturer’s capabilities. The R_0A product of HgCdTe photodiodes varies significantly with temperature, spectral cutoff, and device quality. It also varies with the amount of background flux incident on the device. The R_0A product defines the maximum D^* in the limit of reduced background flux (see Fig. 170 and Fig. 19). In addition to theoretically higher D^* , high-quality PV HgCdTe detectors have lower $1/f$ noise than PC HgCdTe devices, with $1/f$ knee frequencies as low as 1 Hz or less. However, the noise of PV detectors increases more rapidly with increasing temperature than for PC detectors, making photodiodes less attractive for applications where cooling is limited. Photodiodes of high quality are more difficult to make than good photoconductors and can be expected to warrant a premium price. Antireflection coating is available from some diode producers, but is not routinely offered.

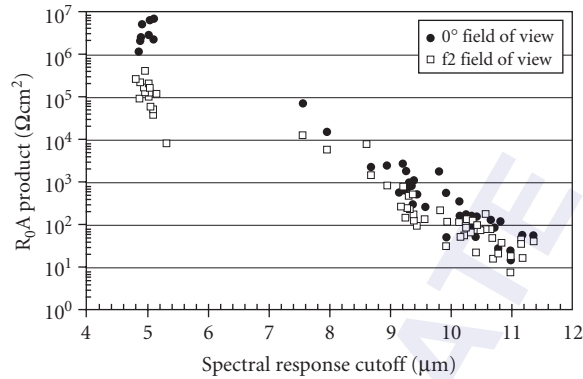


FIGURE 170 R_0A product “trendline” for small (25×25 to $100 \times 100 \mu\text{m}$) HgCdTe photodiodes at 77 K. Data is shown for devices with zero background flux (0° FOV) and with an F/2 field of view (29°) of a 300-K background. For 5- μm spectral cutoff material, the R_0A product is higher at 0° FOV by about an order of magnitude, compared with an F/2 background. At 10 μm , there is less of a difference between the two background conditions. Note that the R_0A product will generally be somewhat lower for larger area diodes.

Both PC and PV HgCdTe detectors are useful for infrared heterodyne detection. When sufficient local oscillator power is available, detector cooling becomes less important, since photon noise can dominate thermal noise at comparatively higher temperatures. Other things being equal, the photovoltaic detector has $\sqrt{2}$ sensitivity (signal-to-noise voltage) advantage over the photoconductor. For 10.6- μm heterodyne detection, 0.1×0.1 -mm HgCdTe *pin* photodiodes with sensitivity near the quantum limit of $\approx 2 \times 10^{-20}$ W/Hz are available with bandwidths up to several gigahertz. Ordinary photodiodes of the same area have bandwidths of several hundred megahertz. Photoconductors make better 10.6- μm heterodyne detectors when cooling is limited to TE-cooled temperatures of 180 K up to room temperature. At 180 K, TE-cooled photoconductors offer bandwidths of 50 to 100 MHz and heterodyne NEPs of 1 to 2×10^{-19} W/Hz. At room temperature, the NEP at 10.6 μm is limited to about 1×10^{-16} W/Hz. Immersion does not improve the performance of minimally cooled heterodyne detectors, since optical gain is already provided by the local oscillator.

Photoconductive HgCdTe

Sensitivity: Adjustable by varying alloy composition (see Figs. 167, 171 to 174).

Dark resistance: 15 to 2000 Ω/sq depending upon temperature, spectral cutoff, and surface passivation.

Responsivity: Varies with spectral cutoff, temperature, detector resistance, element length, and bias voltage or power. See Eq. (21) and Fig. 175 for detector elements with 50×50 - μm dimensions.

Noise: $1/f$ noise is dominant at frequencies below 50 to 1000 Hz for LWIR detectors at 77 K (greater for LWIR at room temperature or TE-cooled). Generation-recombination (thermal or photon) white noise is present beyond the $1/f$ region at a level of less than 10^{-9} V/ $\sqrt{\text{Hz}}$ to 2×10^{-8} V/ $\sqrt{\text{Hz}}$, depending upon spectral cutoff, background flux, responsivity, bias, and operating temperature. Noise and signal rolloff at high frequency is determined by the time constant. See Fig. 176 for an example of the noise spectrum of an LWIR detector at 77 K.

Operating temperature: 77 K and below to 300 K and above for short spectral cutoffs and/or with significant D^* reduction for operation at higher temperatures. Detector immersion can increase D^* at elevated temperatures where thermal noise is dominate.

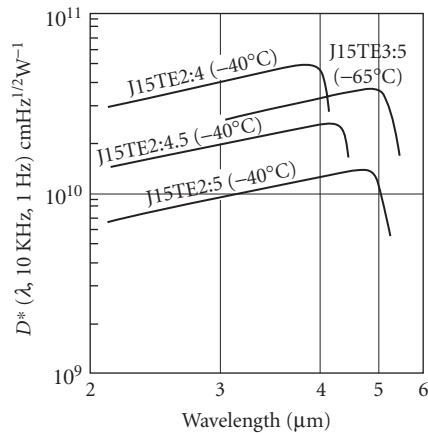


FIGURE 171 Typical D^* as a function of wavelength for a variety of MWIR HgCdTe photoconductors with thermoelectric cooling. (Teledyne Judson Technologies.)

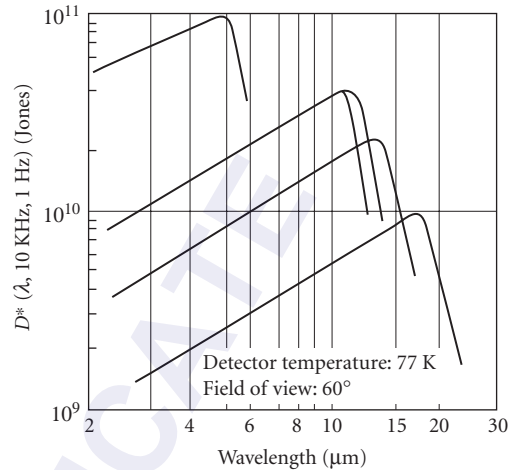


FIGURE 172 Typical D^* as a function of wavelength for a variety of LWIR and VLWIR HgCdTe photoconductors at 77 K. (Teledyne Judson Technologies.)

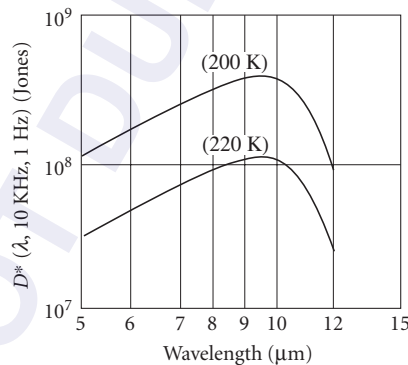


FIGURE 173 Typical D^* as a function of wavelength for LWIR HgCdTe photoconductors at 200 and 220 K. These units are cooled with three- or four-stage thermoelectric coolers. (Teledyne Judson Technologies.)

Linearity: At 77 K linearity begins to degrade at photon flux levels above $\sim 10^{-3}$ W/cm². At 200 K linearity begins to degrade at photon flux levels above ~ 1 W/cm².

Sensitive area: 0.025 to 4-mm linear dimensions.

Quantum efficiency: Typically >70 percent, 85 to 90 percent with antireflection coating.

Capacitance: Low, limited by mounting configuration.

Time constant: 1–2 μ s for LWIR at 77 K (see Fig. 176), depends on spectral cutoff, temperature, doping, and bias.

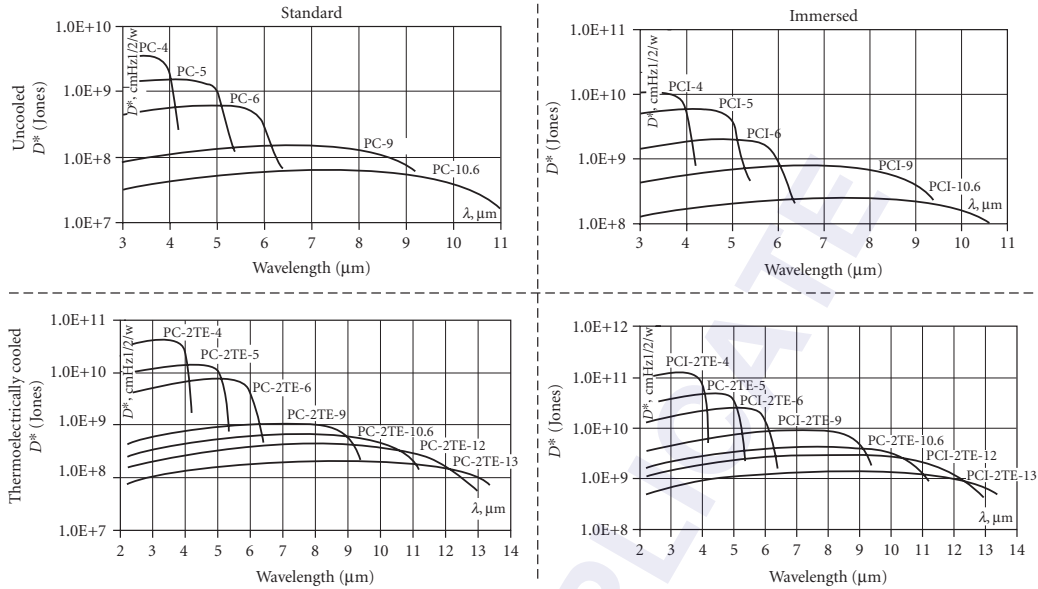


FIGURE 174 D^* for photoconductive HgCdTe detectors as a function of wavelength: *upper left*—ambient temperature operation; *upper right*—ambient temperature operation and immersed; *lower left*—thermoelectrically cooled; *lower right*—thermoelectrically cooled and immersed. (Vigo Systems.)

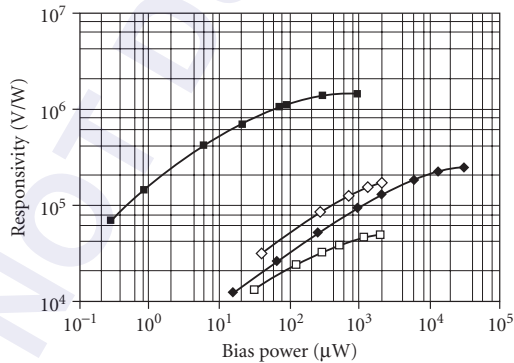


FIGURE 175 Range of peak responsivities for 12- μm cutoff HgCdTe photoconductors at 80 K. These devices have nominal dimensions of $50 \times 50 \mu\text{m}$, and resistance of 50 to $150 \Omega/\text{square}$. (Santa Barbara Research Center.)

Circuit: Standard photoconductive.

Manufacturers: Belov Technology, Boston Electronics, Hamamatsu, Infrared Associates, Kolmar Technologies, Oriel, Teledyne Judson Technologies, Vigo Systems.

Photovoltaic HgCdTe

Sensitivity: Adjustable by varying alloy composition (see Figs. 168, 170, 177, and 178). Also compare Fig. 170 with Fig. 19 for an estimate of the extent to which D^* may increase (up to the R_0A or shunt resistance limit) as the background flux is reduced.

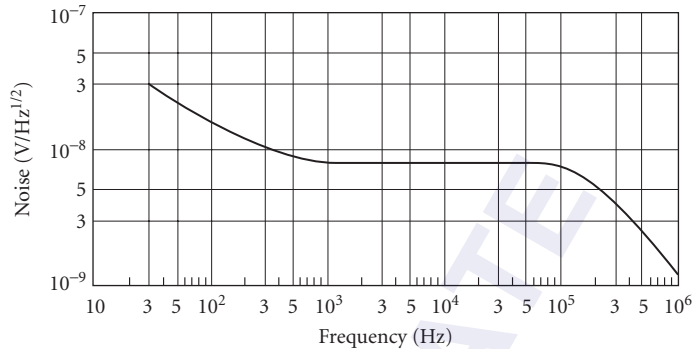


FIGURE 176 Variation of noise with frequency for photoconductive HgCdTe at 77 K. (GEC Marconi Infrared Ltd.)

Time constant: Depends on diode capacitance (area); 10 to 20 ns without bias; 0.5 to 3 ns with reverse bias (some trade-off of sensitivity). Low-capacitance *pin* devices with response out to several gigahertz (0.05 to 0.2-ns time constant) are also available for 10.6- μm CO₂ laser heterodyne detection.

Resistance: Refer to Fig. 170 for the R_0A product at zero bias corresponding to the detector cutoff wavelength (this figure shows very high quality diode impedances) and divide by the diode area. Large-area diodes will have somewhat lower R_0A product than shown in this figure. R_0A varies somewhat with background flux as can be noted from Fig. 170.

Operating temperature: Depends on spectral cutoff; 77 K and lower for LWIR and VLWIR detectors, up to room temperature for SWIR devices. Optical immersion and/or TE cooling will boost the performance for all spectral ranges compared with operation at ambient temperature—see Fig. 178.

Noise: High-quality devices may have flat noise response from 1 Hz out to the high-frequency limit of the time constant. $1/f$ noise may be present in lower quality devices and will increase with reverse bias.

Quantum efficiency: >50 percent (60–75 percent typical) without antireflection coating. Higher with antireflection coating.

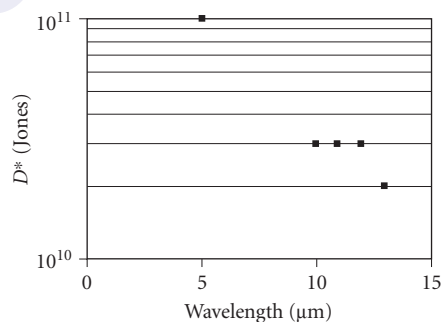


FIGURE 177 D^* specifications for small (50×50 to $250 \times 250 \mu\text{m}$) HgCdTe photodiodes at 77 K as a function of spectral cutoff. Data is shown for devices with 60° FOV background flux. (Fermionics, Mercury Cadmium Telluride MWIR and LWIR Detector Series.)

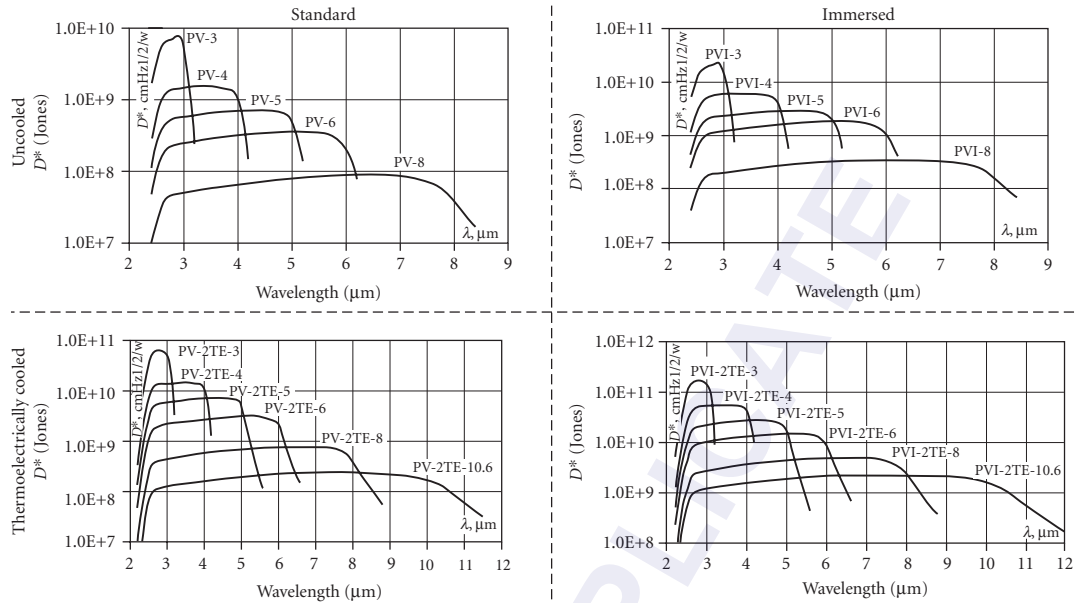


FIGURE 178 D^* for photovoltaic HgCdTe detectors as a function of wavelength: *upper left*—ambient temperature operation; *upper right*—ambient temperature operation and immersed; *lower left*—thermoelectrically cooled; *lower right*—thermoelectrically cooled and immersed. (Vigo Systems.)

Sensitive area: 0.0 to 0.25-mm square, 0.5- and 1-mm diameter.

Capacitance: Depends on junction doping, area, and applied bias (very slightly dependent on spectral cutoff). For standard *pn* junction devices at zero bias and 10^{15} cm^{-3} doping, capacitance is approximately $3 \times 10^4 \text{ pF/cm}$. Significantly lower for *pin* junction devices.

Circuit: Standard photovoltaic circuits, reverse-bias operation to enhance speed and zero bias to maximize D^* .

Manufacturers: Boston Electronics, Kolmar Technologies, Oriel, Raytheon Vision Systems, Vigo Systems.

PbSnTe PbSnTe offers an alternative semiconductor alloy system based upon the IV-VI chemical groups to the II-VI (HgCdMnTeSe) groups previously described for fabricating variable spectral cut-off detectors. Only photovoltaic detectors are available in PbSnTe. This technology has an advantage in the ease of material growth and in the fabrication of good quality photodiode junctions. It has a disadvantage in the very high dielectric constant of the material, combined with relatively high doping concentrations giving high-capacitance (comparatively slow) detectors. For low-frequency applications this is not a disadvantage.

Sensitivity: $D^*_{\text{peak}} > 10^{10}$ Jones (see Fig. 179).

Time constant: $> 50 \text{ ns}$.

Sensitive area: $1 \times 1 \text{ mm}$.

Operating temperature: 77 K.

Circuit: Standard photovoltaic.

Manufacturers: No suppliers are presently known.

Ge:Hg Mercury-doped germanium detectors are fast single-crystal impurity-doped photoconductors, sensitive out to 14 m. Ge:Hg is especially well suited for detection through the 8- to 13- μm

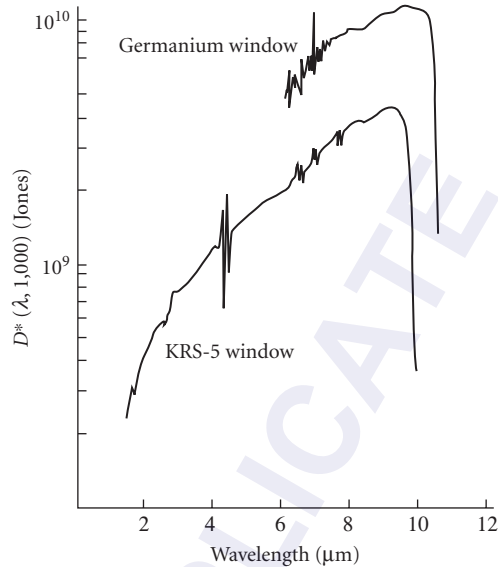


FIGURE 179 Photovoltaic PbSnTe D^* versus wavelength at 77 K and 60° FOV. (Barnes Engineering Co.)

atmospheric window and for detection of near-ambient sources. Unfortunately, its operating temperature must be kept less than 40 K, where it becomes 300-K background limited.

Sensitivity: See Figs. 161, 162, 180, and 181.

Quantum efficiency: 25 to 30 percent.

Noise: See Figs. 182 and 183.

Time constant: 100 ns with 50- Ω load for $T < 28$ K and electric fields < 30 V/cm. (Compensated material is available with ~ 5 -ns time constant with a 50- Ω load. Responsivity then is reduced by 5–10 \times and detectivity is reduced by 2.)

Responsivity: Depends on concentration of compensating impurities, bias, area, and background flux. See Figs. 183 to 185, $\sim 10^5$ V/W.

Dark resistance: Depends on area and FOV: ~ 100 k Ω for 180° FOV (see Fig. 186).

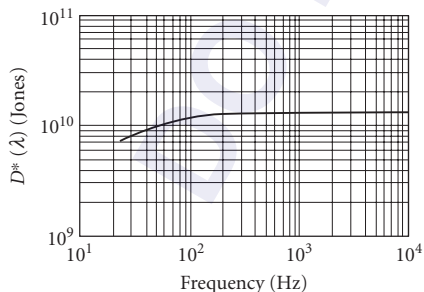


FIGURE 180 Typical D^* versus frequency at 30 K for Ge:Hg; 1×1 -mm area; essentially constant with temperature. (Santa Barbara Research Center.)

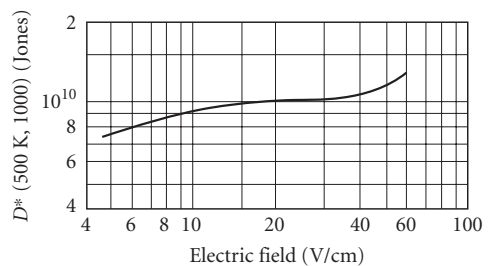


FIGURE 181 D^* versus electric field for Ge:Hg; $T = 5$ K, 90° FOV; 300-K background; 6×10^{-4} cm $^{-2}$ area; Irtran II window. (Reprinted by permission of Texas Instruments.)

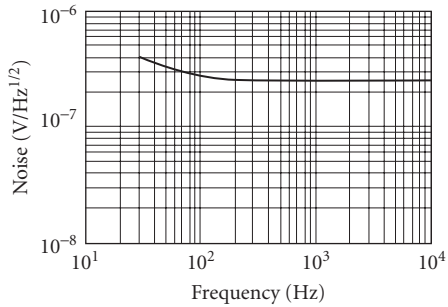


FIGURE 182 Typical noise versus frequency spectrum for Ge:Hg; 1×1 -mm area; essentially constant with temperature. (Santa Barbara Research Center.)

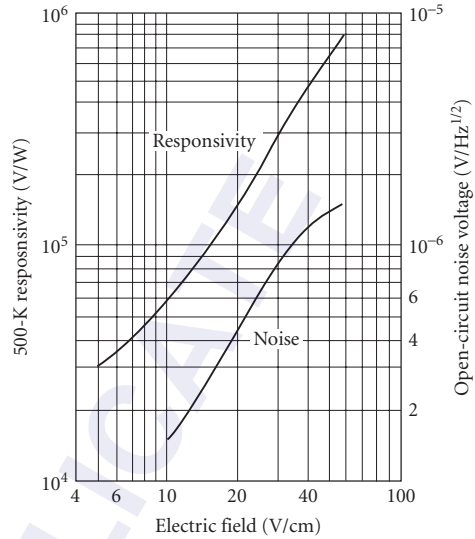


FIGURE 183 Open-circuit responsivity and noise voltage versus electric field for Ge:Hg; $T = 5$ K, 90° FOV; 300-K background; 6×10^{-4} cm² area; Irtran II window. (Reprinted by permission of Texas Instruments.)

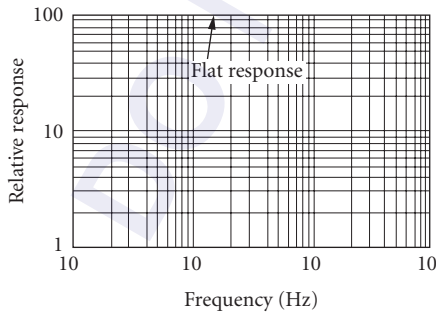


FIGURE 184 Ge:Hg typical relative response versus frequency; $T = 30$ K. (Santa Barbara Research Center.)

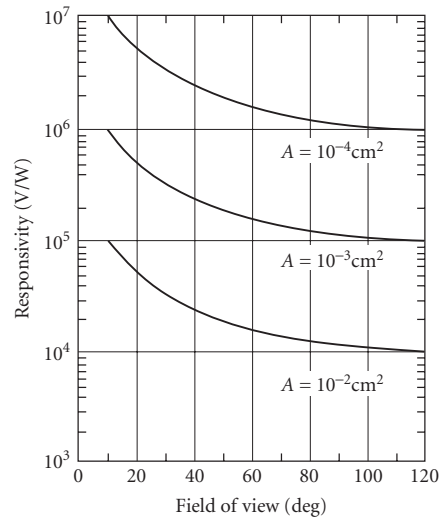


FIGURE 185 Open-circuit responsivity versus FOV for Ge:Hg at 5 K for various detector areas; 300-K background; 500-K blackbody source. (Reprinted by permission of Texas Instruments.)

Capacitance: < 1 pF.

Sensitive area: 1 to 5-mm diameter.

Operating temperature: See Fig. 162.

Linearity: 10^{-3} to 10^{-8} W (size dependent).

Sensitivity profile: ± 15 percent.

Recommended circuit: Standard photoconductive (see Ge: Au). See Fig. 187 for current-voltage characteristics.

Manufacturer: No suppliers are presently known.

Si:Ga Gallium in silicon forms an acceptor level with a binding energy of ~ 72 meV which is the basis of an infrared detector with spectral response out to approximately $17 \mu\text{m}$, as shown in Fig. 188. The exact spectral cutoff and quantum efficiency will vary slightly with the gallium doping concentration. Gallium-doped silicon requires cooling to 20 K or lower for optimum performance. Background-limited performance associated with a quantum efficiency of about 15 percent is achievable over a wide range of background flux levels, provided the operating temperature is low enough to reduce thermal noise below the photon noise level.

Sensitivity: D^* is given by $D^* = 1.1 \times 10^{10} \times \sqrt{(A\lambda\eta)/Q}$ (Jones); where A is detector area, λ is the wavelength in micrometers, η is the quantum efficiency, and Q is the background flux in watts.

Responsivity: 0.9 A/W.

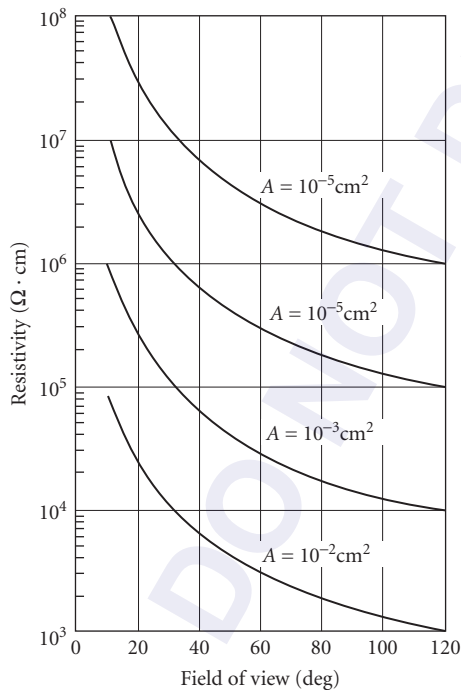


FIGURE 186 Open-circuit resistivity versus FOV for Ge:Hg at 5 K for various detector areas; 300-K background; 500-K blackbody source. (Reprinted by permission of Texas Instruments.)

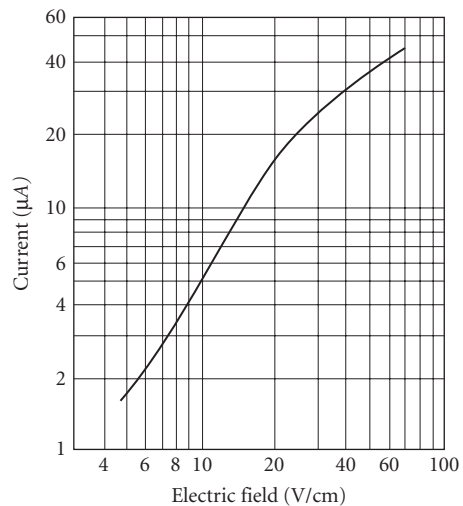


FIGURE 187 Ge:Hg bias current versus electric field; 90° FOV; $T = 5$ K, 300-K background; $6 \times 10^{-4} \text{ cm}^{-2}$ area; Irtran II window. (Reprinted by permission of Texas Instruments.)

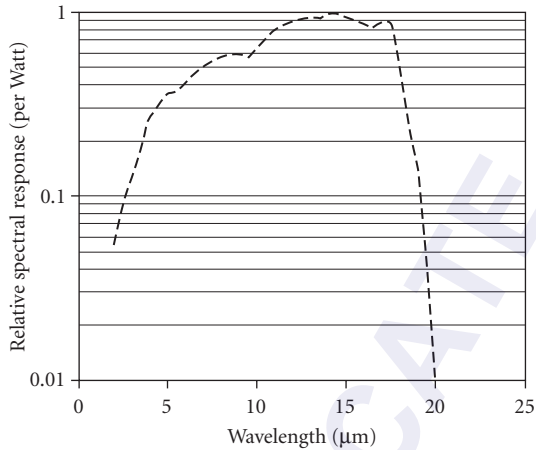


FIGURE 188 Relative spectral response per watt of Si:Ga as a function of wavelength. Data is normalized to unity at the peak response.

Time constant: <1 μs.

Resistance: Depends on background flux and detector bias (similar to Ge:Hg, see Fig. 186).

Capacitance: <1 pF (limited by mounting configuration).

Sensitive area: 0.2 to 2-mm square.

Operating temperature: <20 K.

Recommended circuit: Standard photoconductive.

Manufacturer: Infrared Laboratories.

Si:B Boron in silicon forms an acceptor level with a binding energy of ~45 meV which is the basis of this infrared detector with spectral response out to approximately 30 μm. The exact spectral cut-off and quantum efficiency will vary slightly with the boron-doping concentration. Boron-doped silicon requires cooling to about 15 K or lower for optimum performance. Background-limited performance associated with a quantum efficiency of about 10 percent is achievable over a wide range of background flux levels, provided the operating temperature is low enough to reduce thermal noise below the photon noise level.

Sensitivity: D^* is given by $D^* = 1.1 \times 10^{10} \times \sqrt{(A\lambda\eta)/Q}$ (Jones); where A is detector area, λ is the wavelength in micrometers, η is the quantum efficiency, and Q is the background flux in watts.

Responsivity: 2 A/W.

Time constant: <1 μs.

Resistance: Depends on background flux and detector bias (similar to Ge:Hg, see Fig. 186).

Capacitance: <1 pF (limited by mounting configuration).

Sensitive area: 0.2 to 2-mm square.

Operating temperature: <15 K.

Recommended circuit: Standard photoconductive (see Ge:Au).

Manufacturer: Infrared Laboratories.

Ge:Cu Copper-doped germanium detectors are fast, single-crystal, impurity-doped photoconductors, with high sensitivity in the broad region 2 to 30 μm. Operating temperature must be maintained below 20 K (ideally <14 K). Ge:Cu is then 300-K background-limited, and response time <50 ns.

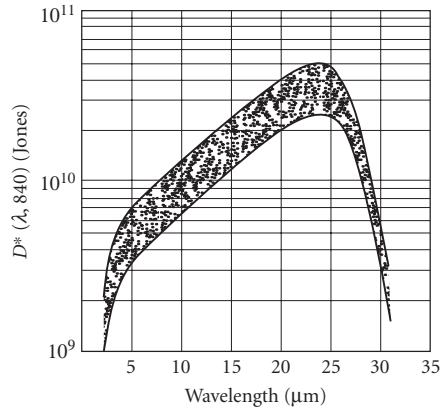


FIGURE 189 Range of spectral detectivities for Ge:Cu; 60° FOV. (Santa Barbara Research Center.)

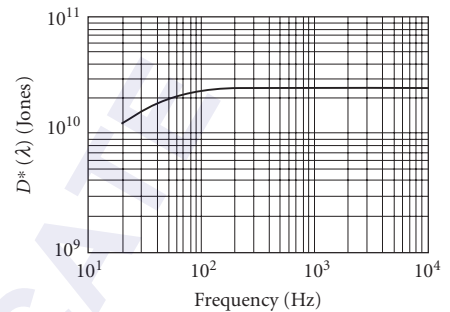


FIGURE 190 Typical D^* vs. frequency for Ge:Cu at 4.2 K. (Santa Barbara Research Center.)

Sensitivity: See Figs 161, 162, 189, and 190.

Noise: See Figs. 191 and 192.

Time constant: ~100 ns (can be doped to be faster, ~5 ns). (See discussion of time constant under Ge:Cu.)

Responsivity: 10^5 V/W (see Figs. 192 to 194).

Dark resistance: Depends on FOV (~100 k Ω for 180° FOV).

Capacitance: <1 pF.

Sensitive area: 1 to 5-mm diameter.

Operating temperature: See Figs. 162 and 194.

Linearity: 10^{-3} – 10^{-8} W/cm² (depends on size).

Sensitivity profile: ± 15 percent.

Stability: Stable in all ambient storage environments tested.

Recommended circuit: See Figs. 166 and 195.

Manufacturer: No suppliers are presently known.

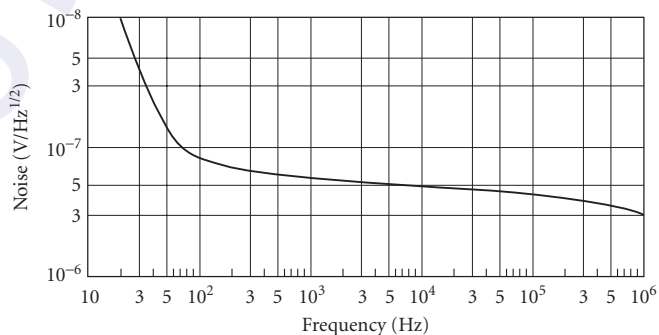


FIGURE 191 Typical noise frequency spectrum for Ge:Cu. (GEC Marconi Infra Red Ltd.)

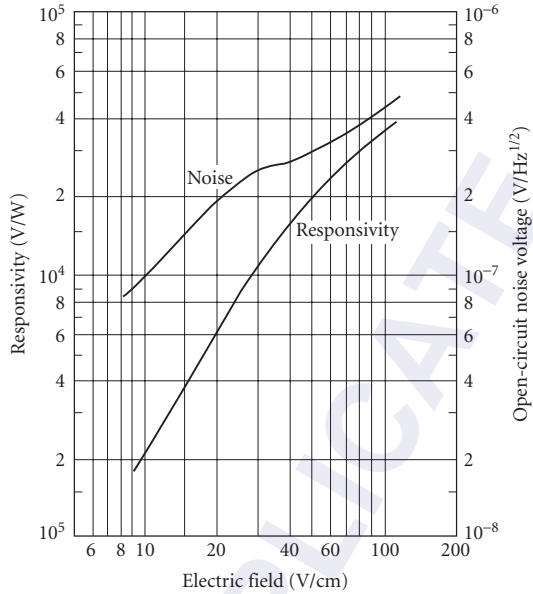


FIGURE 192 Typical noise and responsivity versus biasing field for Ge:Cu; 5 K, 60° FOV; 300-K background, $A = 10^{-2} \text{ cm}^2$, 500-K blackbody. (Reprinted by permission of Texas Instruments.)

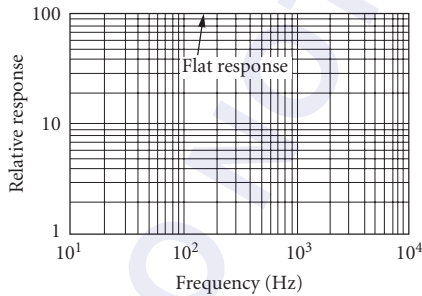


FIGURE 193 Relative response versus frequency for Ge:Cu at 4.2 K. (Santa Barbara Research Center.)

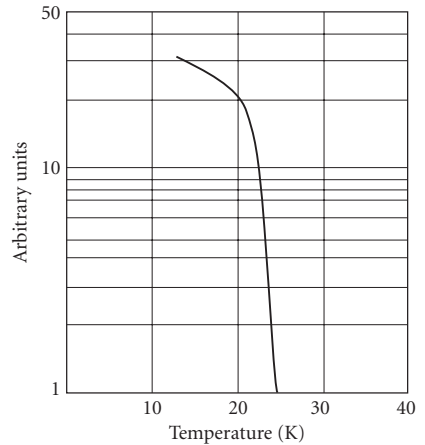


FIGURE 194 Relative responsivity versus temperature for Ge:Cu.

Ge:Zn Very similar to Ge:Cu except that cutoff wavelength moves out to 42 μm and operating temperature should be $<10 \text{ K}$. A relatively low field breakdown limits the responsivity.

Sensitivity: See Figs 161, 162, and 196.

Noise: See Fig. 197 and 198.

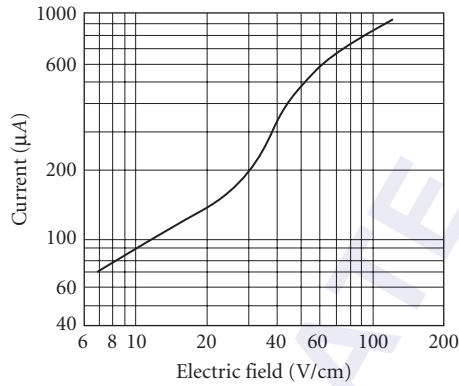


FIGURE 195 Typical current-voltage curve for Ge:Cu, 60° FOV; 300-K background; $A = 10^{-2} \text{ cm}^{-2}$. (Reprinted by permission of Texas Instruments.)

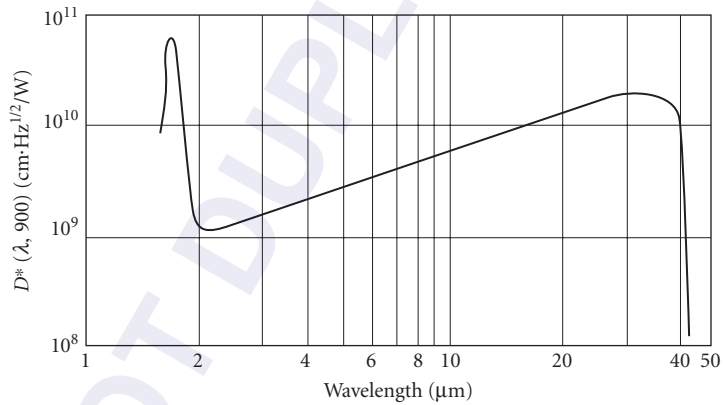


FIGURE 196 D^* versus wavelength for Ge:Zn.

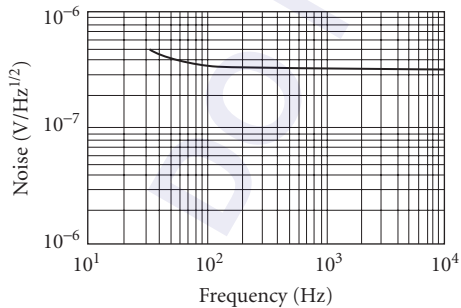


FIGURE 197 Typical noise-frequency spectrum for Ge:Zn at 4.2 K; $A = 1 \times 1 \text{ mm}$. (Santa Barbara Research Center.)

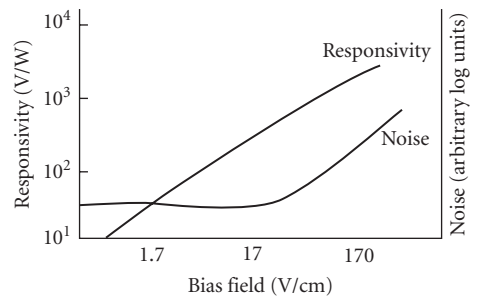


FIGURE 198 Signal and noise for Ge:Zn.

Responsivity: 10^3 V/W (see Fig. 198).

Time constant: <50 ns. (See discussion of time constant under Ge:Au.)

Dark resistance: 0.5 to 5 M Ω / sq (60°-FOV ambient background).

Capacitance: <1 pF (limited by mounting configuration).

Sensitive area: 1-, 2-, 3-, and 5-mm diameters.

Operating temperature: <10 K.

Recommended circuit: Standard photoconductive.

Manufacturer: No suppliers are presently known.

Ge:Ga The elements of B, Al, Ga, In, and Tl from chemical group III form shallow acceptor states (~10 meV) in germanium which are the basis of infrared detectors with spectral response out to approximately 120 μm . Currently, gallium-doped germanium is commercially available, but germanium doped with other group III elements (Ge:B, Ge:Al, Ge:In, Ge:Tl) will give similar detector performance. The small binding energies associated with this detector require cooling to liquid helium temperatures (4.2 K) or lower for optimum performance. Background-limited performance associated with a quantum efficiency of about 7 percent is achievable over a wide range of background flux levels, provided the operating temperature is low enough to reduce thermal noise below the photon noise level.

Sensitivity: D^* is given by $D^* = 1.1 \times 10^{10} \times \sqrt{(A\lambda\eta)/Q}$ (Jones); where A is detector area, λ is the wavelength in μm , η is the quantum efficiency, and Q is the background flux in watts.

Responsivity: 4 A/W.

Time constant: <1 μs .

Resistance: Depends on background flux and detector bias.

Capacitance: <1 pF (limited by mounting configuration).

Sensitive area: 0.5-, 1-, and 2-mm square.

Operating temperature: <4.2 K, best below 3 K.

Recommended circuit: Standard photoconductive.

Manufacturer: Infrared Laboratories.

Photographic In this paragraph we present only the spectral sensitivity of some typical photographic emulsions. See Chap. 29, "Photographic Films," for a more extensive coverage.

The term spectral sensitivity generally has a different meaning when applied to photographic detectors than it does when applied to the other detectors described in this chapter. It comes closer to responsivity than to minimum detectable power or energy.* In Fig. 199 sensitivity is the reciprocal of exposure, expressed in ergs per centimeter, required to produce

$$\text{Density} = \log \frac{1}{\text{transmittance}} = 0.3$$

above gross fog in the emulsion when processed as recommended.

Manufacturers: AGFA, Eastman Kodak, Fuji, Polaroid.

*Work is in progress to evaluate photographic materials in terms of minimum detectable energy, a concept involving the average number of photons necessary to produce a change in density (signal) equal to that of the fog-density fluctuations (noise); see Refs. 10, 26, 27, and 28.

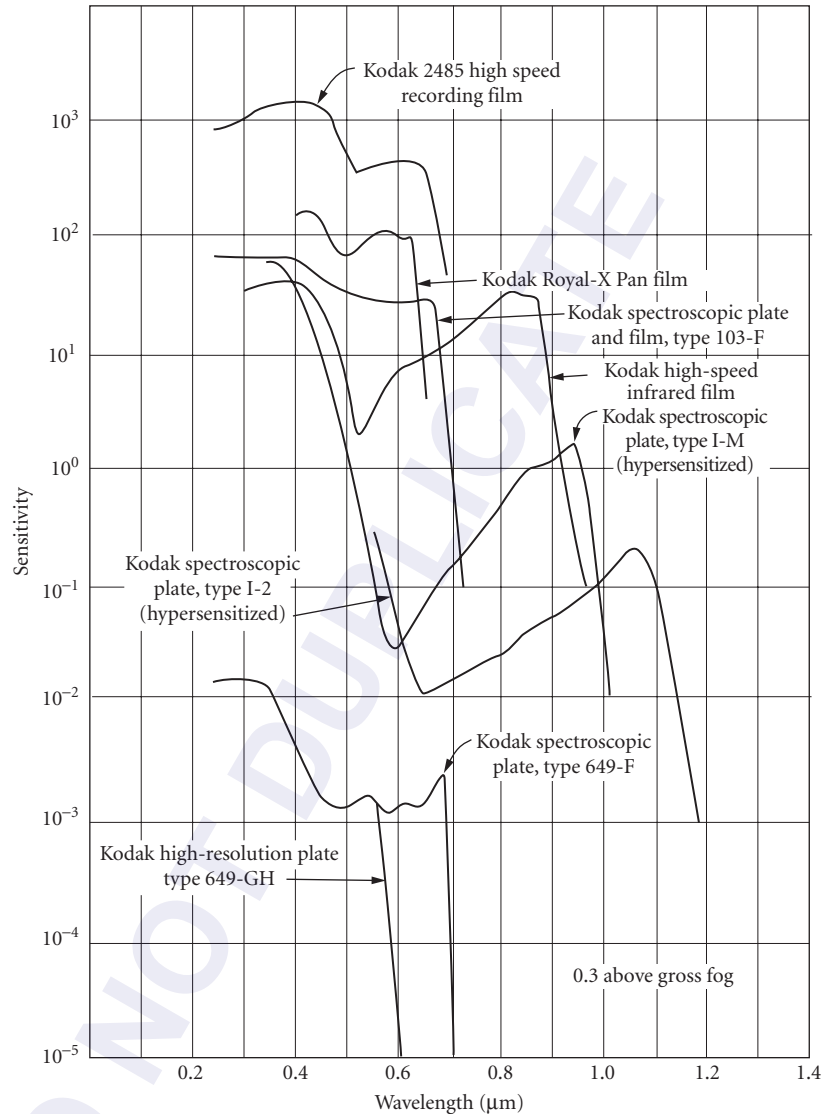


FIGURE 199 Sensitivity versus λ for typical photographic emulsions. (Eastman Kodak.)

24.8 REFERENCES

1. W. L. Wolfe, (ed.), *Handbook of Military and Infrared Technology*, Office of Naval Research, Washington, D.C., 1965.
2. Institute of Radio Engineers, "IRE Standards on Electron Tubes: Methods of Testing," *Proceedings of the IRE* 50(9): 1974-1975 (1962).
3. Radio Corporation of America, "Phototubes and Photocells," *Tech. Man. PT-60* (1963).

4. D. Vincent, *Fundamentals of Infrared Detector Operation and Testing*, Wiley, New York, 1989.
5. E. L. Dereniak and D. Crowe, *Optical Radiation Detectors*, Wiley, New York, 1984.
6. W. Rogatto, ed., *The Infrared and Electro-Optical Systems Handbook*, vol. 3, SPIE Press, Bellingham, Wash., 1993.
7. R. C. Jones, "Factors of Merit for Radiation Detectors," *J. Opt. Soc. Am.* **39**:344 (1949).
8. G. Bauer, "Ein halbleiter Hochohmbolometer mit Tafel, IV," *Phys. Z.* **44**:53 (1943).
9. J. A. R. Samson, *Techniques of Ultraviolet Spectroscopy*, Wiley, New York, 1967.
10. R. Jones, "On the Quantum Efficiency of Photographic Negatives: On the Minimum Energy Detectable by Photographic Materials," *Photogr. Sci. Eng.* **2**:57–65j, 191–204 (1958).
11. D. J. Fink, *Principles of Television Engineering*, McGraw-Hill, New York, 1940.
12. R. A. Smith, F. E. Jones, and R. P. Chasmar, *Detection and Measurement of Infrared Radiation*, Oxford University, London 1957.
13. S. F. Jacobs, and M. Sargent III, "Photon Noise Limited D^* for Low Temperature Backgrounds and Long Wavelengths," *Infrared Phys.* **10**(4):233–235 (1970).
14. P. W. Kruse, L. D. McLaughlin, and R. B. McQuistan, *Elements of Infrared Technology*, Wiley, New York, 1962.
15. P. B. Fellgett, "On the Ultimate Sensitivity and Practical Performance of Radiation Detectors," *J. Opt. Soc. Am.* **39**:970 (1949).
16. F. J. Low, and A. R. Hoffman, "The Detectivity of Cryogenic Bolometers," *Appl. Opt.* **2**:649 (1963).
17. F. J. Low, "Low-Temperature Germanium Bolometer," *J. Opt. Soc. Am.* **51**:1300 (1961).
18. E. H. Eberhardt, "Threshold Sensitivity and Noise Ratings of Multiplier Phototubes," *Appl. Opt.* **6**:251 (1967).
19. W. C. Livingston, "Enhancement of a Photocathode Sensitivity by Total Internal Rejection as Applied to an Image Tube," *Appl. Opt.* **5**:1335 (1966).
20. W. D. Gunter, G. R. Grant, and S. A. Shaw, "Optical Devices to Increase Photocathode Quantum Efficiency," *Appl. Opt.* **9**:251 (1970).
21. M. Cole and D. Ryer, "Cooling PM Tubes for Best Spectral Response," *Electro Optical Systems Design* **4**(6):16–19 (June 1972).
22. H. E. Bennett, "Accurate Method for Determining Photometric Linearity," *Appl. Opt.* **5**:1265 (1966).
23. R. E. Simon, A. H. Sommer, J. J. Tietjen, and B. F. Williams, "New High-Gain Dynode for Photomultipliers," *Appl. Phys. Lett.* **13**:355 (1968).
24. G. A. Morton, H. M. Smith, and H. R. Krall, "Pulse Height Resolution of High Gain First Dynode Photomultipliers," *Appl. Phys. Lett.* **13**:356 (1968).
25. S. M. Sze, *Physics of Semiconductor Devices*, 2d ed., Wiley, New York, p. 773, 1981.
26. T. H. Johnson, "Lead Salt Detectors and Arrays: PbS and PbSe," *Proc. SPIE* **443**: 60–94, (1984).
27. J. C. Marchant, "Exposure Criteria for the Photographic Detection of Threshold Signals," *J. Opt. Soc. Am.* **54**:79 (1964).
28. G. R. Bird, R. C. Jones, and A. E. Ames, "The Efficiency of Radiation Detection by Photographic Films: State-of-the-Art and Methods of Improvement," *Appl. Opt.* **8**:2389 (1969).

24.9 SUGGESTED READINGS

- Sommer, A. H., *Photoemissive Materials*, Wiley, New York, 1968.
- Sommers, H. S., Jr. and E. K. Gritchell, "Demodulation of Low-Level Broadband Optical Signals with Semiconductors," *Proc. IEEE* **54**:1553 (1966).
- Sommer, A. H., and W. B. Teusch, "Demodulation of Low-Level Broadband Optical Signals with Semiconductors, II: Analysis of the Photoconductive Detector," *Proc. IEEE* **52**:144 (1964).
- Sun, C. and T. E. Walsh, "Performance of Broadband Microwave-Biased Extrinsic Photoconductive Detectors at 10.6 μm ," *IEEE J. Quantum Electron.* **6**:450 (1970).

PHOTODETECTION

Abhay M. Joshi

*Discovery Semiconductors, Inc.
Cranbury, New Jersey*

Gregory H. Olsen

*Sensors Unlimited, Inc.
Princeton, New Jersey*

25.1 GLOSSARY

A	photodetector active area
A_0	incident photon flux
B	bandwidth of the photodetector
C	capacitance of the photodetector
D^*	detectivity
E	applied electric field
E_a	activation energy
E_g	bandgap of the semiconductor
E_i	impurity energy state
f	frequency
I_{diff}	diffusion current
I_{g-r}	generation-recombination current
IR_0	unity gain current
IR_1	reverse current generated by avalanche action
I_{tun}	tunneling current
k	Boltzmann's constant
L	distance traveled by a charge carrier
M	photocurrent gain
m	effective mass of a electron
N_A	acceptor impurity concentration on p side
N_D	donor impurity concentration on n side
n	refractive index of the AR coating
q	electron charge
R	sum of the detector series resistance and load resistance

R_o	detector shunt impedance
T	temperature in kelvin
t_n	transit time of electrons
t_p	transit time of holes
t_r	transit time of charge carriers (holes or electrons)
V	applied reverse bias in volts
V_B	breakdown voltage
V_{bi}	built-in potential of a p - n junction
W	depletion width of a p - n junction
α	absorption coefficient of the photodetector's absorption layer
ϵ_s	semiconductor permittivity
η	quantum efficiency of photodetector
θ	tunneling constant
λ	wavelength of incident photons (nm)
λ_{co}	detector cutoff wavelength (10 percent of peak response, nm)
μ	mobility of charge carriers (holes or electrons)
μ_n	mobility of electrons
μ_p	mobility of holes

25.2 INTRODUCTION

The approach of this chapter is descriptive and tutorial rather than encyclopedic. It is assumed that the reader is primarily interested in an overview of how things work. Among the many excellent references to be consulted for further details are Sze's book,¹ and the article by Forrest.² For the latest in photodetector developments, consult recent proceedings of the Society of Photo-optical and Instrumentation Engineers (SPIE) conference or the IEEE Optical Fiber Conference.

A photodetector is a solid-state sensor that converts light energy into electrical energy. According to Isaac Newton, light energy consists of small packets or bundles of particles called *photons*. Albert Einstein, who won a Nobel prize for the discovery of the photoelectric effect, showed that when these photons strike a metal they can excite electrons in it. The minimum photoenergy required to generate (excite) an electron is defined as the *work function* and the number of electrons generated is proportional to the intensity of the light. The semiconductor photodetectors are made from different semiconductor materials such as silicon, germanium, indium gallium arsenide, indium antimonide, and mercury cadmium telluride, to name a few. Each material has a characteristic bandgap energy E_g which determines its light-absorbing capabilities. Light is a form of electromagnetic radiation comprised of different wavelengths (λ). The range of light spectrum is split approximately as ultraviolet (0–400 nm), visible (400–1000 nm), near infrared (1000–3000 nm), medium infrared (3000–6000 nm), far infrared (6000–40,000 nm), extreme infrared (40,000–100,000 nm). The equation between bandgap energy E_g and cutoff wavelength (λ_c) is

$$\lambda_c = \frac{1.24 \times 10^3 \text{ nm}}{E_g \text{ (eV)}} \quad (1)$$

The smaller the bandgap (eV), the farther the photodetector “sees” into the infrared. Table 1 lists some prominent photodetector materials, their bandgaps, and cutoff wavelengths λ_c at room temperature (300 K).

Photodetectors find various applications in fiber-optic communications (800–1600 nm), spectroscopy (400–6000 nm), laser range finding (400–10,600 nm), photon counting (400–1800 nm),

TABLE 1 Important Photodetector Materials

Type	E_g (eV)	λ_c (nm)	Band
Silicon	1.12	1100	Visible
Gallium arsenide	1.42	875	Visible
Germanium	0.66	1800	Near-infrared
Indium gallium arsenide*	0.73–0.47	1700–2600	Near-infrared
Indium arsenide	0.36	3400	Near-infrared
Indium antimonide	0.17	5700	Medium-infrared
Mercury cadmium	0.7–0.1	1700–12500	Near-to-far-infrared

*The alloy composition of indium gallium arsenide and mercury cadmium telluride can be changed to alter the bandgap E_g .

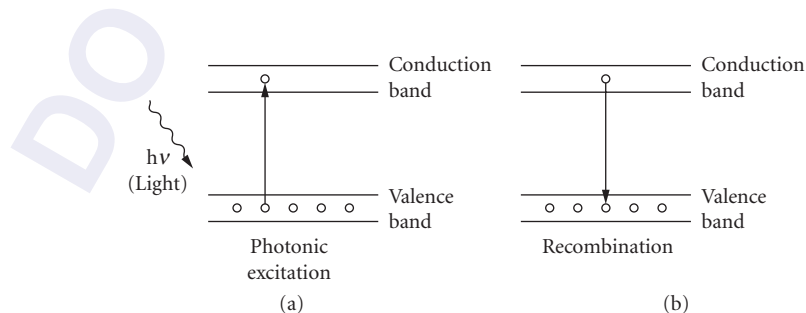
and satellite imaging (200–1200 nm), to name only a few topics. We will discuss three kinds of photodetectors: (1) Photoconductors, (2) *p-i-n* photodetectors (including avalanche diodes), and (3) photogates. Frequently, such detectors need to have high sensitivity, low noise, and high reliability. For fiber-optic applications, the frequency response and the cost can be a critical issue, whereas for infrared applications, many times it is the area of the photodetector. Large area (1 in diameter), high-sensitivity silicon avalanche photodetectors compete with conventional photomultiplier tubes for low light sensing applications in the visible. They offer the advantage of compact size and more rugged construction. We also discuss reliability issues concerning photodetectors and take notice of a few novel photodetector structures.

25.3 PRINCIPLE OF OPERATION

When an electron in the valence band receives external energy in the form of light, the electron may overcome the nuclear attraction and become a “free electron.” When light energy creates this transformation, it is termed *photonic excitation* (see Fig. 1a). The range of energies acquired by these free electrons is termed the *conduction band*. The energy difference between the bottom of this conduction band and the top of the valence band is termed the *energy bandgap* E_g and represents the minimum energy of light that the material can absorb.

However, under the influence of even a small external electric field, the free electrons can “drift” in a specific direction. This is the fundamental principle of a photodetector. Figure 2 shows the three kinds of photodetectors discussed in this chapter. A brief explanation of each kind follows.

All photodetectors can be characterized by their quantum efficiency, detectivity (sensitivity), and response time.¹ Quantum efficiency (QE) is perhaps the most fundamental property, as it determines just how efficiently the device converts incoming photons into conduction electrons.

**FIGURE 1** Photonic excitation and recombination in a semiconductor.

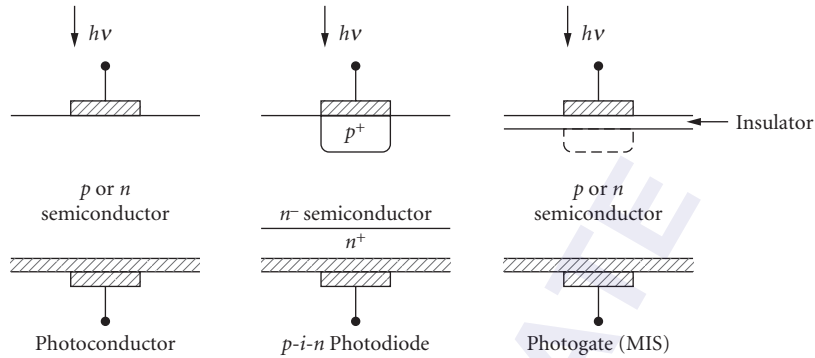


FIGURE 2 Types of photodetectors.

Usually expressed in percentage, quantum efficiency can range from under 1 percent for PtSi Schottky barrier infrared detectors to well over 90 percent for InGaAs *p-i-n* fiber-optic photodetectors. Responsivity (R) is a related term expressed in amps/watt, which determines how much photocurrent is produced by optical power of a given wavelength. Detectivity measures how *sensitive* a detector is; that is, not only its light conversion efficiency, but also its ability to detect low-level light signals. It is limited by the various noise currents (shot, $1/f$, etc.) introduced by the detector. Finally, response time describes how rapidly a detector can respond to a changing light signal. This ranges from milliseconds for certain types of PbS photoconductors to picoseconds for GaAs-like metal-insulator-semiconductor detectors.

These three parameters are frequently traded off. A large-area detector captures more light signal and thus might have greater detectivity. However, its larger capacitance would slow down the device. Similarly, response time in *p-i-n* detectors can be improved by thinning the absorbing region of the detector. However, this in turn cuts down quantum efficiency by reducing the total number of photons absorbed.

Photoconductor

A photoconductor, as the name implies, is a device whose conductivity increases with illumination. It acts as an open switch under dark (or no illumination) and as a closed switch under illumination. The simple equivalent electric circuit is shown in Fig. 3. This basic principle of a photoconductor finds numerous applications in relays and control circuits. An ideal switch should have low resistance in the closed position and, therefore, a pair of ohmic (not-rectifying) contacts are formed to the photoconductor. These ohmic contacts form the electrodes and usually have contact resistance of less than 10 ohms.

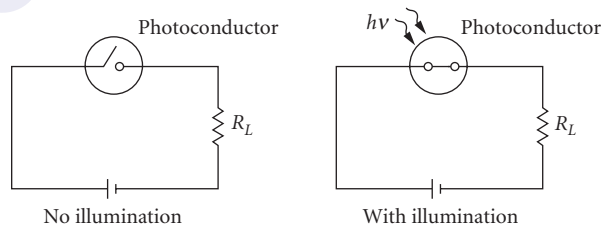


FIGURE 3 Equivalent circuit diagram of a photoconductor.

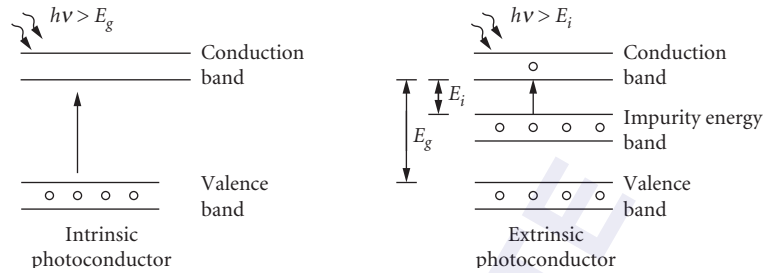


FIGURE 4 Functional diagram of an intrinsic and extrinsic photoconductor.

Types of Photoconductors The two principle photoconductors are (1) intrinsic and (2) extrinsic. In an intrinsic device there is no external impurity atom. However, when an external impurity (dopant) is added to a material, it is termed *extrinsic*. This impurity atom occupies an energy state between the valence band and the conduction band. The functional difference between the intrinsic and extrinsic photoconductor is seen in Fig. 4. For an intrinsic device, the photo excitation ($h\nu$) needs to have energy greater than the bandgap energy E_g and its cutoff wavelength λ_c is given by Eq. (1). But, for an extrinsic one, the photon excitation ($h\nu$) should exceed the impurity energy state E_i and its cutoff wavelength λ_{co} is

$$\lambda_{co} = \frac{1.24 \times 10^3 \text{ nm}}{E_i (\text{eV})} \quad (2)$$

For intrinsic photoconductors, it is extremely difficult to achieve bandgap energies E_g less than 0.1 eV (refer to Table 1). This limits its capability to see in the far infrared and extreme infrared (13,000–100,000 nm) and beyond. This is overcome by extrinsic devices whose E_i value is less than 0.1 eV and is normally done by doping germanium or silicon. However, the extrinsic photoconductor suffers from very low absorption coefficients and, hence, poor quantum efficiencies. Also, since ambient thermal energy can excite carriers, they have to be cooled to liquid nitrogen temperature (77 K) and below, whereas most intrinsic photoconductors can operate at room temperature (300 K).

Photo Gain The sensitivity of a photoconductor is determined by its gain. Photo gain is defined as the ratio of the output signal to the input optical signal. When photons impinge on a photoconductor, they generate electron-hole pairs and, under the influence of external fields, they are attracted toward the anode and cathode. A typical photoconductor is illustrated in Fig. 5 with L being the thickness of the active layer. The transit time (t_r) required for a charge carrier to travel a distance L is given by

$$t_r = \frac{L^2}{\mu V} \quad (3)$$

where V = applied voltage bias and μ = charge mobility.

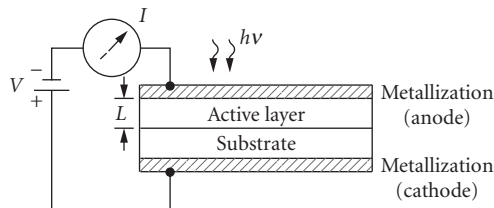


FIGURE 5 A typical photoconductor.

The mobility of electrons μ_n and that of holes μ_p is different, with μ_n being usually far higher than μ_p . This causes a difference in the transit time of electrons and holes. Hence, photon-generated electrons are swept away more quickly than holes which result in a positive charge in the active layer. To maintain the charge neutrality, new electrons are supplied by the external voltage source. Therefore, for one incident photon, more than one electron is circulated in the electric circuit. This results in an “effective gain.” Thus, photoconductor gain can be defined as the ratio of slower transit time t_p to faster transit time t_n .

$$M = \frac{t_p}{t_n} \quad \text{or} \quad \frac{\mu_n}{\mu_p} \quad (4)$$

The slower the transit time, the higher the gain; however, the bandwidth of the device is reduced. Hence, high-gain photoconductors will result in slow devices and vice versa. Such “high-gain, slow devices” can be best utilized for imaging applications.³ For high-speed optical communication applications in the 1000- to 1700-nm spectrum, InGaAs is the material of choice due to its high mobility. Several reports have been published on high-speed InGaAs photoconductors that find practical applications in optical receivers.⁴⁻⁹ (Also see Chap. 26, “High-Speed Photodetectors,” by John Bowers and Yih G. Wey.)

***p-i-n* Photodiode**

Unlike photoconductors, a photodiode has a *p-n* junction, usually formed by diffusion or epitaxy. In a photoconductor, metal contacts are made to either *n*- or *p*-type material. However, a photodiode consists of both *n*- and *p*-type materials across which a natural electric field is generated. This field is known as the *built-in potential* V_{bi} , and its value depends on the bandgap of its material. A silicon *p-n* junction has V_{bi} of 0.7 V whereas in germanium it is 0.3 V. The higher the bandgap E_g , the larger the built-in potential V_{bi} . An important physical phenomenon called *depletion* occurs when a *p*-type semiconductor is merged with an *n*-type semiconductor. After an initial exchange of charge, a potential is built up to prevent further flow of charge. This built-in field creates the depletion width W , which is a region free of any charge carriers and is given by¹⁰

$$W = \sqrt{\frac{2\epsilon_s(N_A + N_D)}{q(N_A N_D)}(V_{bi} - V)} \quad (5)$$

where N_A and N_D are impurity concentrations of *p* side and *n* side, respectively, q is the electron charge, V_{bi} is the built-in potential, and ϵ_s is the semiconductor permittivity, V is the applied bias and is negative for reverse-bias operation. As seen from Eq. (5), under reverse bias the depletion width W increases causing a decrease in the capacitance of the photodiode. A *p-i-n* photodiode is similar to a parallel plate capacitor with the anode-cathode being the two plates and the depletion width W being the separating medium. A typical InGaAs *p-i-n* photodiode is shown in Fig. 6 and its capacitance is given by¹

$$C = \frac{\epsilon_s A}{W} \quad (6)$$

where ϵ_s is the semiconductor permittivity and A is the active area of the photodiode. From first principles, a decrease in capacitance improves the bandwidth B of the photodetector according to

$$B = \frac{0.35}{2.2RC} \quad (7)$$

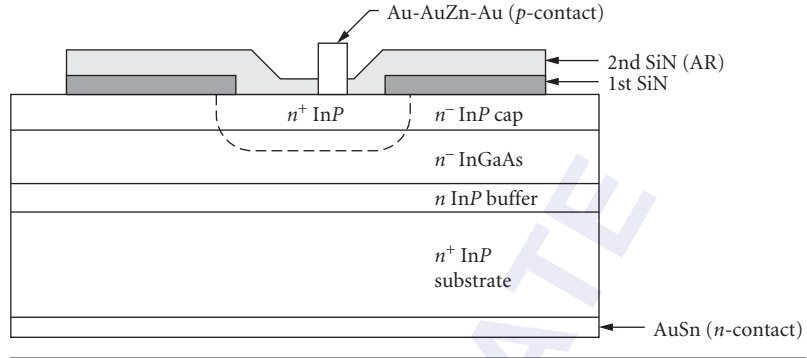


FIGURE 6 A typical InGaAs *p-i-n* photodiode.

where, R is the sum of the detector series resistance and load resistance. For a detailed analysis on high-speed photodetectors, see Chap. 26 by Bowers and Wey.

Dark Current For applications ranging from optical communications (III-V compound semiconductors) to infrared sensing (Si, Ge, III-IV, and II-IV compound semiconductors), a *p-i-n* photodiode must have high sensitivity and low noise. These are largely determined by the dark currents originating in the device. Several authors have published papers on dark currents in InGaAs¹¹⁻¹³ and HgCdTe.^{14,15} The three major components of dark current are (1) diffusion current, (2) generation-recombination current, and (3) tunneling current.

Diffusion current In the nondepleted region of the photodiode, electron-hole pairs are formed by the ambient temperature. These thermally generated carriers diffuse toward the depletion region and produce the diffusion current.

$$I_{\text{diff}} \propto e^{-E_g/kT} \quad (8)$$

where E_g is the bandgap of the photodiode material, k is Boltzmann's constant, and T is the ambient temperature in kelvin. It is clear from Eq. (8) that the diffusion current is higher in a low-bandgap material. Therefore, InSb ($E_g = 0.17$ eV) has far higher diffusion current than silicon ($E_g = 1.12$ eV) and in fact this makes InSb almost useless at room temperature. To overcome this excessive diffusion current, InSb photodiodes are cooled to liquid nitrogen temperature (77 K).

Generation-recombination current The current generated in the depletion region of the photodiode is called the *generation-recombination current*. When impurity trap levels are present within the forbidden gap E_g , trapped carriers can be elevated to the conduction band with less energy than for diffusion current. This "trap-assisted" current is given by

$$I_{g-r} \propto \sqrt{(V_{\text{bi}} - V)} e^{-E_g/2kT} \quad (9)$$

From Eq. (5), the depletion width W is proportional to $V_{\text{bi}} - V$. Hence,

$$I_{g-r} \propto W e^{-E_g/2kT} \quad (10)$$

Generation-recombination current is proportional to the volume of the depletion width and, hence, is reverse-bias-dependent, whereas the diffusion current in Eq. (8) is bias-independent. For high-bandgap semiconductors with bandgaps above 1.0 eV (e.g., silicon), the generation current

usually dominates over the diffusion current at room temperature. However, for low-bandgap material such as indium antimonide, the diffusion current is dominant over generation current at room temperature.

Tunneling current When the electric field in a reverse-biased p - n junction exceeds 10^5 V/cm, a valence band electron can jump to the conduction band due to the quantum mechanical effect¹⁰ called *tunneling* which occurs at high field and with geometrically narrow energy barriers. The tunneling current is given by

$$I_{\text{tun}} \propto EV \exp\left(\frac{-\theta\sqrt{m}}{E} E_g^{3/2}\right) \quad (11)$$

where E is the applied electric field, m is the effective mass of an electron, and θ is a dimensionless constant whose value depends on the tunneling barrier height. Higher doping levels at the p - n junction lead to a narrower depletion width which causes higher electric fields, thus increasing the amount of tunneling current. Low-bandgap photodiodes exhibit much more tunneling than do higher-bandgap diodes. Tunneling shows a weak dependence on temperature, the only minor change being caused by the temperature dependence of the bandgap E_g . This leads to a *decreasing* breakdown voltage with an increasing temperature as opposed to an *increasing* breakdown voltage exhibited by the avalanche effect.

Quantum Efficiency, Responsivity, and Absorption Coefficient Quantum efficiency η is defined as the ratio of electron-hole pairs generated for each incident photon. In a nonavalanche p - i - n photodiode, quantum efficiency is less than unity. Responsivity R is a measured quantity in amps/watt or volts/watt and is related to quantum efficiency by

$$\eta = \frac{(1240)R}{\lambda} \quad (12)$$

where λ is the wavelength in nm of incident photons and R is the responsivity in amps/watt. The value of η is determined by the absorption coefficient α of the semiconductor material and the penetration distance x in the absorbing layer. The light flux A at a distance x with the absorption layer is

$$A = A_0 e^{-\alpha x} \quad (13)$$

where, A_0 is the incident photon flux and α is a strong function of wavelength λ . Figure 7 shows its typical values for a 1- μm -thick undoped $\text{In}_x\text{Ga}_{1-x}\text{As}$, $0 < x < 0.25$.¹⁶ For optimized η , the reflectivity at the semiconductor surface has to be minimized. Hence, an antireflection (AR) coating of proper thickness is deposited on the photodiode surface. For single-layer AR coatings, the proper “quarter-wave” thickness (L) of the AR coating is

$$L = \frac{\lambda}{4n} \quad (14)$$

where n is its refractive index. With good AR coatings, InGaAs photodiodes can achieve quantum efficiencies above 95 percent at 1300 to 1500 nm. For the visible region, silicon photodiodes show high η (90 percent) in the 800-nm range, and the mid-infrared InSb has a typical η of 80 percent at 5000 nm.

Avalanche Photodiodes Avalanche photodiodes (APDs) will be briefly discussed here. For a detailed treatment, see Chap. 26, “High-Speed Photodetectors,” by Bowers and Wey. An avalanche photodiode is a p - i - n diode with a net efficiency or gain greater than unity. This is obtained through the process of “impact ionization” by operating the photodiode at a sufficiently high reverse bias. The typical operating voltage for an InGaAs APD is 75 V, while that for silicon can be as high as 400 V. The

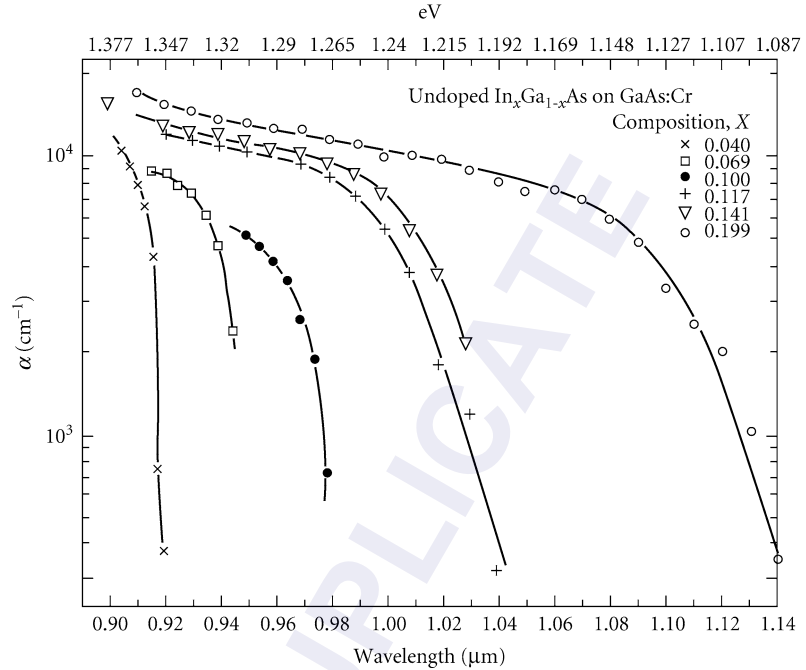


FIGURE 7 Absorption coefficients for 1- μm -thick undoped $\text{In}_x\text{Ga}_{1-x}\text{As}$, $0 < x < 0.25$.¹⁶

impact ionization process is described in Fig. 8. Under the influence of a high electric field ($>5 \times 10^5$ V/cm), electron A gains sufficient kinetic energy to strike atom B with a tremendous force and knock out an electron hole pair $B'-B''$. A is called the “parent” and $B'-B''$ the “child” charge carriers. The child electron B'' moves through a critical distance S and acquires enough kinetic energy to create its own child particles $D'-D''$. The sum effect of the impact ionization of a number of electrons is termed *avalanche multiplication*. Because of this avalanche action, the gain in an APD exceeds unity, reaching useful values above 10 for InGaAs and several hundred for silicon before the multiplied noise begins to exceed the multiplied signal. A solid-state APD is a fast device with gain-bandwidth products that can exceed 20 GHz.^{17,18} In spite of high operating bias, an APD can be designed for low noise operation¹⁹ and used for numerous applications such as photon-counting, laser pulse detection,²⁰ and fiber-optic communication. The gain or avalanche multiplication M of an APD is given by

$$M = \frac{IR_1}{IR_0} \quad (15)$$

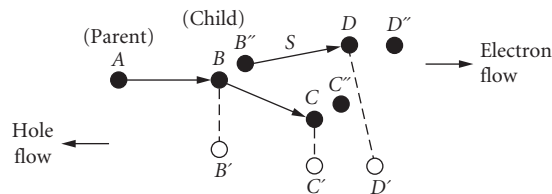


FIGURE 8 Impact ionization process.

where IR_1 is the reverse current generated by avalanche action and IR_0 is the unity gain current. At voltage breakdown of the APD, the multiplication factor M tends to infinity. An empirical relation between the multiplication factor M and reverse bias V is given by^{21,22}

$$M = \frac{1}{1 - (V/V_B)^n} \quad (16)$$

where V is the applied reverse bias and V_B is the breakdown voltage. The factor n varies between 3 and 6, depending on the semiconductor material and its substrate type. Typical gains are on the order of 10 to 20 for germanium and InGaAs APDs, and above 100 for silicon APDs. Due to their lower noise, InGaAs and silicon APDs have better sensitivity than their germanium counterparts.

Extended Wavelength (1000–3000 nm) Photodetectors

Detector materials used for the 1000 to 3000 nm spectrum include InSb, InAs, PbS, HgCdTe, and recently InGaAs. PbS is an inexpensive, reasonably sensitive detector that can operate at relatively high temperatures, even at room temperature. Its major drawback is its slow (typically milliseconds) response time. InAs has higher sensitivity over the 1000 to 3500 nm spectrum and fast response time, but must be cooled thermoelectrically (to 230 K) or cryogenically (to 77 K). InSb has similar properties out to 5500 nm but must definitely be cryogenically cooled. HgCdTe has high sensitivity and speed and it can be operated at room temperature. Indium gallium arsenide was originally developed for fiber-optic applications out to 1.7 μm (using $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$) but it can be used out to 2500 nm by increasing its indium content to $\text{In}_{0.8}\text{Ga}_{0.2}\text{As}$. InGaAs appears to be the *best* detector material for high-temperature operation in the 1000 to 3000 nm spectrum. It has a 10 to 100 times advantage in shunt resistance at room temperature compared to HgCdTe—the previously used material for this wavelength.

Table 2 contains a summary of the data. It is difficult to find data at exactly the same cutoff wavelengths and temperature with the same area device. R_0A was determined (in cases where it was given as such) by simply multiplying two numbers, where, R_0 is the shunt impedance of the detector, and A is the active area of the photodetector.

Photogate (Metal-Insulator-Semiconductor Detector)

The advent of silicon charge-coupled devices (CCDs) has revolutionized the television industry and introduced one of the most popular consumer items to millions of people around the world—the CCD camcorder. From the sandy shores of Hawaii to the ski slopes of Colorado, people have captured

TABLE 2 Comparison of R_0A Values in HgCdTe and InGaAs ($\Omega\text{-cm}^2$)

λ_{co} (nm)	$R_0A(T)$	
	HgCdTe	InGaAs
1400	4×10^4 (292 K)	2.5×10^5 (300 K)
	7×10^6 (230 K)	1.3×10^8 (220 K)
1700	2×10^2 (300 K)	2.5×10^5 (300 K)
	2×10^5 (220 K)	1.3×10^8 (220 K)
2100	7×10^1 (300 K)	2.5×10^3 (300 K)
	7×10^3 (220 K)	6.5×10^5 (220 K)
2500	1×10^1 (300 K)	1.3×10^2 (300 K)
	1×10^3 (210 K)	1.0×10^5 (210 K)

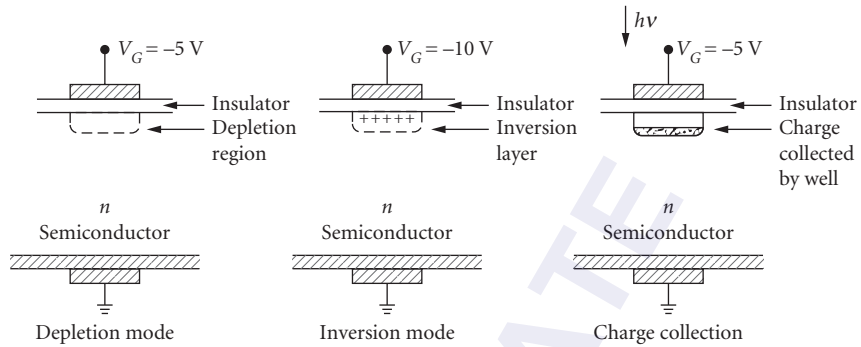


FIGURE 9 Schematic diagram of a photogate.

life's best moments with a CCD camcorder. With its superior imaging quality and noise performance of a few electrons per pixel (<20), a CCD has diverse uses from space imaging to chemical analysis spectroscopy. A CCD is a matrix of metal oxide semiconductor (MOS) devices operated in the "depletion" mode. Each individual MOS device is called a *photogate* and its schematic diagram is shown in Fig. 9. Consider an n -type semiconductor with a negative potential applied to its gate. It repels the negatively charged electrons and create a depletion layer. As the negative potential on the gate is increased, the volume of the depletion region increases further into the bulk. However, the surface potential at the semiconductor-insulator interface also becomes more negative. Finally, with increased gate bias, the surface potential becomes sufficiently high to attract minority carriers (holes). This creates a positive charge at the semiconductor-insulator interface and is termed the *inversion layer*. In a MOS transistor, the inversion layer forms a conducting channel between the source and the drain, and the gate bias needed to achieve inversion is termed the *threshold voltage*. Usually a photogate is operated in depletion at a gate bias lower than the threshold voltage. When incident photons create hole-electron pairs, the minority carriers drift away to the depletion region and the volume of the depletion region shrinks. The total amount of charge that a photogate can collect is defined as its *well capacity*. The total well capacity is decided by the gate bias, the insulator thickness, the area of the electrodes, and the background doping of the semiconductor. Numerous such photogates with proper clocking sequence form a CCD imaging array. For in-depth understanding of CCDs, we refer to Chap. 32 "Visible Array Detectors," by Timothy J. Tredwell.

25.4 APPLICATIONS

The main commercial uses of photodetectors include optical communications and infrared sensing. Although these applications often overlap, optical communication typically involves transmitting data over an optical fiber at higher rates. The format is increasingly digital (telecommunications and data links) at rates from 1 Mbit/s to over 2 Gbit/s.

However, one growing application is cable TV where analog data rates from 1 to over 1000 MHz are most often found. Infrared sensing mostly involves nonfiber applications at sub-MHz analog rates. The property to be detected is usually the amplitude (in watts) and wavelength of the incoming radiation. In digital applications, the wavelength and individual pulse amplitude are relatively fixed, and successful communication occurs simply by distinguishing when the pulse is "on" or "off." Although very weak pulses must sometimes be detected, the actual amplitude of the pulse is irrelevant. The ultimate "resolving power" of the detector is when a weak pulse can no longer be distinguished from background noise, i.e., when the incoming signal strength S equals the background

noise strength N or when the signal-to-noise ratio $S/N = 1$. Thus, the strengths of individual pulses are unimportant as long as the presence of a pulse can be detected. Information is conveyed by the timing sequence of the pulses rather than by the amplitude of the individual pulse. Analog applications, on the other hand, depend critically on the frequency content and amplitude of the transmitted signal. In an AM cable TV transmission system, the detector must be able to linearly reproduce the incoming optical signal as an electrical current of the same frequency content and amplitude, and to minimize intermodulation and harmonic distortion that is invariably produced in the detection of an AM signal.

Infrared applications often involve spectroscopy whereby the detected electrical signal depends on both the optical wavelength and strength of the incoming infrared signal. Thus, the detector must be carefully calibrated in terms of “responsivity” (electrical amps/optical watt) versus wavelength in order to accurately identify the nature of the incoming signal. Identification of gases (e.g., methane, which absorbs light near 1650 nm) depends on these properties. Other “infrared” applications include spectroscopy, remote sensing from satellite, and general laboratory detection. Not *all* infrared applications are analog, however. One notable digital application is LIDAR (Light detection and ranging), which essentially is a form of laser radar. High-intensity light pulses are emitted into the atmosphere (or a gas) which absorbs, scatters, and reemits the laser pulse. The character of the light pulses detected near the source can be used to determine the nature of the gas particles that interact with the light: the absorbing wavelength, the gas density (velocity), and the amount present. Applications include remote pollution monitoring and “windshear detection,” whereby the presence of abrupt changes in wind velocity can be instantly detected at distances of several miles. This application²³ has been demonstrated with laser wavelengths of 2060 and 10,600 nm for use on an aircraft. The 2060-nm system works better in severe storms and does not require a cryogenically cooled detector (as does the 10,600-nm system).

One important noise source in infrared applications is the so-called $1/f$ noise which becomes noticeable at frequencies below 10 MHz. Although poorly understood, this noise is thought to originate at heterointerfaces such as semiconductor-metal contacts and heteroepitaxial interfaces. Photodiode arrays are often used to detect low-light-level signals of a few hundred photons, and they must integrate the signal for 1 second or more. However, with longer integration times, $1/f$ noise may become noticeable and can degrade the S/N ratio and thus, impose an upper limit on the effectiveness of longer integration times. Limiting $1/f$ noise becomes critical for numerous infrared sensing applications and research indicates that surface depletion width at the semiconductor-insulator interface to be a major source of $1/f$ noise in the InGaAs photodiodes.²⁴

One important area for detectors is the array configuration used both for spectroscopy (linear) and imaging (two-dimensional). Linear arrays are used in the so-called multichannel analyzers whereby the detector is placed behind a fixed grating and the instrument functions as “motionless” or “instant” spectrometer with each pixel corresponding to a narrow band of wavelengths. The resolution of the instrument is determined by the number and spacing of pixels, so *narrow* pixel geometries are needed along one direction whereas *tall* pixel geometries are needed along the perpendicular direction to enhance the light collection.

One “mixed” infrared/fiber-optic application is the use of large-area (typically 3 mm diameter) detectors for optical power meters: the optical equivalent of a voltmeter which accurately measures the amount of optical power in watts or dBm (number of decibels above or below 1 mW contained in an incoming beam). The large area ensures large collection efficiency. The most important parameter here is the responsivity and the uniformity of response across the detector.

A “figure-of-merit” for infrared detectors is D^* (D-star), whereby detectors of differing area can be compared. It is related to the noise equivalent power (NEP) in watts, the lowest power a detector can detect at a signal-to-noise ratio of 1 as

$$D^*(\lambda, f, B) = (AB)^{1/2} / \text{NEP} \quad (17)$$

where A is the detector area. The optical bandwidth B (often taken to be 1 Hz), frequency of signal modulation f , and operating wavelength λ must be stated.

25.5 RELIABILITY

In today's global economy of severe competitiveness, new product development and innovation are incomplete without quality assurance and reliability. Reliability is the assurance that a device will perform its stated functions for a certain period of time under stated conditions, and considerable research has been done to improve the reliability of photodetectors.^{25–28} The two major industrial standards for testing semiconductor device reliability are (1) test methods and procedures for microelectronics (MIL. STD. 883C), and (2) Bellcore technical advisory (TA-TSY-00468). The former standard is generic to the semiconductor industry, while the latter is specifically developed for fiber-optic optoelectronic devices.

The tests performed under MIL. STD. 883C comprise the following major groups: (1) *environmental tests*, e.g., moisture resistance, burn-in, seal, dew point, thermal shock, (2) *mechanical tests*, e.g., constant acceleration, mechanical shock, vibration, solderability, and bond strength, and (3) *electrical tests*, e.g., breakdown voltage, transition time measurements, input currents, terminal capacitance, and electrostatic discharge (ESD) sensitivity classification. Under the Bellcore Technical Advisory, each photodetector lot undergoes visual inspection, optical and electrical characterization, and screening. Visual inspection removes any photodiodes with faulty wire bonds or cracks in the glass window or in the insulating films. Table 3 lists the electrical and optical testing performed on every photodiode. After testing, all the devices are sent for screening (burn-in), e.g., some $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$ photodiodes are burned-in at 200°C for 20 hours at –20 V reverse bias to weed out any infant mortality.

Photodetector Life Test (Accelerated Aging)

To predict the lifetime or *mean-time-to-failure* (MTTF), accelerated aging tests are carried out on groups of diodes at several elevated temperatures. For example, the MTTF for 300 μm diameter InGaAs photodiodes was determined on groups of 20 screened devices at elevated temperatures of 200, 230, and 250°C. The failure criterion was a 25 percent increase in the room temperature dark current value.^{28,29} The total lifetest extended over a time period of several years, and every week the samples were cooled to room temperature to check their dark current. Failed devices were removed from the sample population and the remaining good ones put back at the elevated temperature.

From the temperature-dependence of the data, it was observed that the failure mechanism is thermally activated. The Arrhenius relationship calculates the activation energy E_a for thermally activated failure²⁹ as

$$\text{MTTF}(T) = Ce^{(E_a/kT)} \quad (18)$$

where, C is a constant, k is Boltzmann's constant (8.63×10^5 eV/K), and T is the temperature in kelvin. Figure 10 shows the MTTF for three batches of 300 μm diameter $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}/\text{InP}$ photodiodes. Using least-squares fit to the data, the calculated activation energy E_a is 1.31 eV with a correlation coefficient r^2 of 0.99. From Eq. (18) and an experimentally determined activation energy

TABLE 3 Electrical and Optical Testing of Photodetectors

Tests or Measurement	Parameter	Symbol	References
Optical response	Responsivity	R	Bellcore Technology Advisory
	Gain	G	
Electrical performance	Dark current	I_d	TA-TSY-00468 Issue 2, July 1988
	Breakdown voltage	V_{br}	

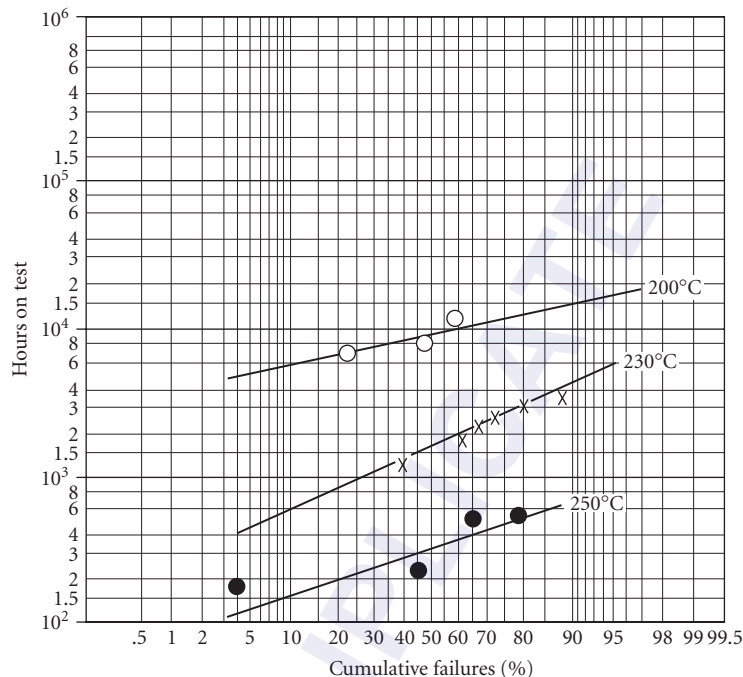


FIGURE 10 Time versus cumulative percent failure of InGaAs diodes at 200, 230, and 250°C lifetest.

of 1.31 eV, the MTTF at 25°C is calculated to be 1.34×10^{14} hours. Such a “geological” lifetime may seem to be an overkill, even for the electronics industry. However, when thousands of these devices are working together in a single system (e.g., telephones), the net MTTF of all these devices chained together may be on the order of only a few years. Thus, continuing improvements in reliability must be an ongoing process. Reliability in most photodetectors is determined by a number of factors including: (1) material quality, (2) processing procedures, (3) planar technology versus mesa technology, and (4) amount of leakage current. Poor material quality can introduce crystal defects such as vacancies and dislocations which can increase the dark current. Higher dark current has been directly linked to lower MTTF.²⁸ Device processing is probably the most crucial item in photodetector reliability. The dielectric (typically silicon nitride) used in planar detector processing serves as a diffusion mask in *p-n*-junction formation and a passivant (termination) for the junction so produced. Any surface states or impurities introduced here can directly increase leakage current and degrade reliability.

An important milestone in detector reliability was the changeover mesa to planar structures.^{2,30} Just as the transistors in the 1950s were first made in mesa form, so were the optical photodetectors of the 1980s, due to their simplicity and ease of fabrication. However, in both cases, reliability issues forced the introduction of the more complex planar structure. A sketch of a mesa and planar photodiode is illustrated in Fig. 11. A mesa photodiode typically is formed by wet chemical etching of an epitaxially grown *p-n* crystal structure, while in a planar process, a *p-n* junction is formed by diffusing a suitable *p* or *n* dopant in an *n*- or *p*-type crystal. It has been shown a planar structure to be more reliable than a mesa one³⁰ because a *p-n* junction is never exposed to ambient conditions in a well-designed planar process. Exposure of the *p-n* junction can cause surface corrosion leading to increased leakage current and, in effect, poorer reliability.³¹

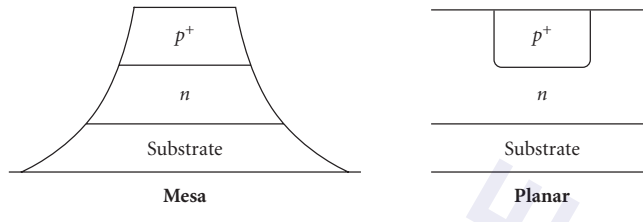


FIGURE 11 Sketch of a mesa and planar photodetector.

25.6 FUTURE PHOTODETECTORS

A lateral p - i - n photodiode and a quantum well infrared photodetector (QWIP) have been developed with better characteristics compared to the photodetector structures. A long wavelength QWIP in the 8000- to 12000-nm band³²⁻³⁵ has posed a severe challenge to the present favorite mercury cadmium telluride photodetectors, while a medium wavelength QWIP in the 3000- to 5000-nm band³⁶ may compete with indium antimonide and platinum silicide photodetectors. A QWIP made from GaAs/AlGaAs heterosystems promises to have higher detectivity (D^*), higher yield due to well-established 3-in wafer GaAs technology, and easier monolithic integration with circuit electronics. A lateral p - i - n diode, as the name implies, has charge carrier flow in a lateral direction compared to the vertical direction in a conventional (vertical) photodiode structure. Because of its process compatibility and simple fabrication, a lateral p - i - n photodiode can be suitably integrated on an optoelectronic integrated circuit (OEIC) chip^{37,38} having numerous field-effect transistors. An OEIC has a lower noise floor due to the reduced stray capacitances and inductances compared to that of hybrid detector-amplifier packages and finds applications in high-speed digital data communication.

Lateral p - i - n Photodetector

The vertical p - i - n structure in Fig. 6 has high sensitivity, low noise, low capacitance, better reliability, and an easy manufacturing process. However, such a vertical structure is nonplanar and therefore harder to integrate on an OEIC. The nonplanarity is also an issue with lasers and LEDs, and optical integration demands surface-emitting LEDs and lasers (SLEDs and SLASERs) over the conventional edge-emitting sources (ELED and ELASER). The cross section of an AlGaAs/GaAs lateral p - i - n photodiode is shown in Fig. 12.

The higher-bandgap AlGaAs layer acts as a surface barrier, reducing the leakage currents. The low-bandgap GaAs layer absorbs the incoming light, and the generated carriers flow to the W-Zn

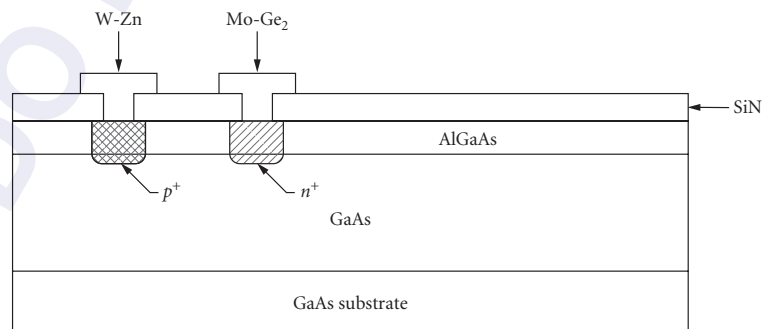


FIGURE 12 Cross section of an AlGaAs/GaAs lateral p - i - n photodiode.³⁸

and Mo-Ge₂ ohmic contacts which act as the p and n regions, respectively. The diffusion of zinc and germanium in the ohmic contacts forms a compositionally graded barrier at the AlGaAs/GaAs interface, rather than an abrupt interface. This smooth barrier helps the lateral p - i - n photodiode to have better speed than a lateral metal-semiconductor-metal photodetector. A comfortable spacing of 3 to 5 μm between the p and n regions gives high quantum efficiency and low capacitance, thus providing all the desirable properties of a vertical p - i - n structure and yet being easier to integrate.

Quantum Well Infrared Photodetector

Quantum well infrared photodetectors (QWIPs) offer *long* wavelength (5000–10,000 nm) infrared detection by using materials whose bandgap normally allows them only to absorb light in the *short* wavelength ($\sim 10,000$ nm) region, e.g., GaAs/AlGaAs. The use of thin (< 500 Å) layers allows the absorbing wavelength to be controlled by material *geometry* rather than its *chemistry*.³⁹

Before discussing the QWIPs, we take the liberty of explaining a few basic terms and concepts of quantum physics. Superlattices or quantum well structures consist of a stack of ultrathin semiconductor layers normally 50 to 500 Å in thickness. Molecular beam epitaxy (MBE) techniques are frequently employed to grow these structures because their characteristically slow growth rate of a few Angstroms per seconds which helps achieve abrupt heterointerfaces. Two semiconductors of different compositions, when stacked together, form a heterointerface. Type III-V compound semiconductors such as AlGaAs/GaAs and InAlAs/InGaAs are the best candidates for growing quantum well structures, as they can be easily doped and their alloy composition readily changed to form semiconductor layers of different bandgaps. Tailoring the bandgap can alter the heterobarriers, creating exciting device results. When quantum well layers have thicknesses less than the electron mean free path (typically 50 to 100 Å), electron and holes cannot have their normal three-dimensional motion. This restricts carriers to move in two dimensions in the plane of the layer.^{2,39} Because of this quantized motion, a new band of discrete energy levels is generated. Carriers no longer obey Boltzmann's statistics¹ and optical absorption becomes more complicated than the conventional band-to-band absorption given by Eq. (1). The absorption of light energy by a quantum well structure can cause an electron to jump from "multiple valence subbands" to "multiple conduction subbands," thereby enabling it to absorb light wavelengths not decided by the *material* properties (bandgap) of the semiconductor layers alone, but by its *geometrical* properties as well.

In QWIPs, the light energy transfers an electron in a bound state to an excited state in the continuum.⁴⁰ Figure 13 shows an AlGaAs/GaAs quantum well structure with L being the width of the

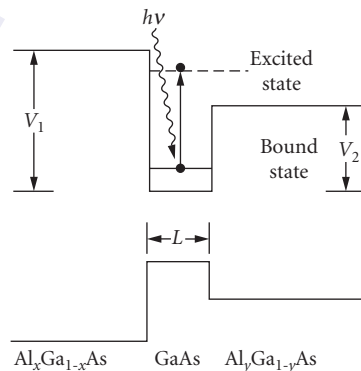


FIGURE 13 Infrared detection with an AlGaAs/GaAs quantum well.⁴⁰

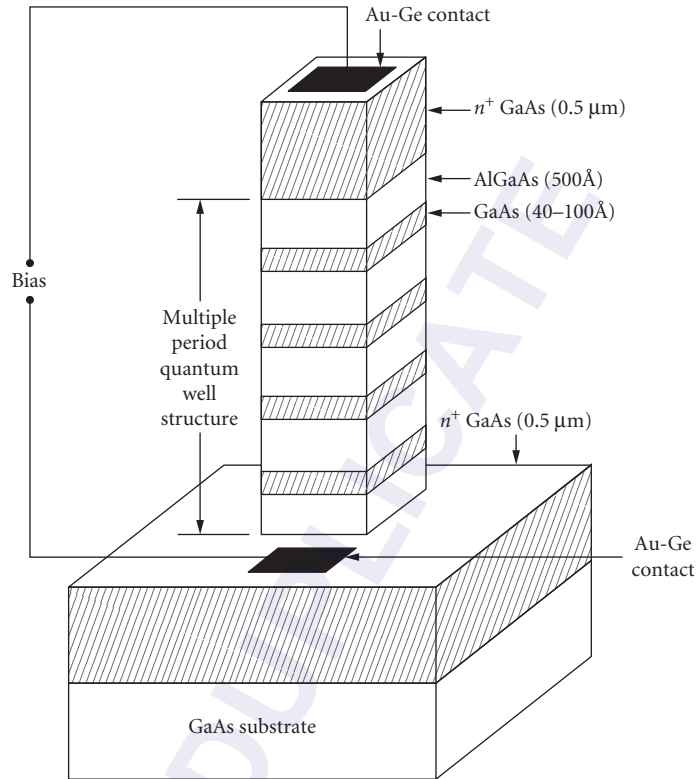


FIGURE 14 Multiple period AlGaAs/GaAs quantum well infrared photodetectors.

well and V_1, V_2 being the barrier heights. The electron excited by the IR radiation is swept out of the doped GaAs well by applying an external electrical field. By controlling the barrier heights V_1, V_2 , and quantum well width L , the spectral response of a QWIP can be changed for the desired IR window of 3000 to 5000 or 8000 to 12,000 nm.⁴⁰ A multiple period quantum well infrared photodetector is illustrated in Fig. 14.³⁴ The n^+ -doped ($2 \times 10^{18} \text{ cm}^{-3}$) GaAs quantum wells are 40 to 100 Å and the undoped AlGaAs barriers are of 500 Å thickness. The multiple period stack is sandwiched between two n^+ GaAs-doped contacts. This photodetector has exhibited a blackbody D^* of $1 \times 10^{10} \text{ cm}/(\text{Hz/W})^{1/2}$ at 68 K for a cutoff wavelength of 10,700 nm. InGaAs/AlInAs superlattices have exhibited blackbody D^* of $2 \times 10^{10} \text{ cm}/(\text{Hz/W})^{1/2}$ at 120 K with peak responsivity at 4000 nm.³⁶

In summary, QWIPs have high detectivity, good uniformity, high yield, multiple spectral windows, and intrinsic radiation hardness for numerous imaging and spectroscopy applications.⁴¹

25.7 ACKNOWLEDGMENTS

We sincerely acknowledge the support of EPITAXX, Inc., Amy Vasger, and Jennifer Romano (Sensors Unlimited) for preparing the manuscript, and Jim Rue for technical discussions.

25.8 REFERENCES

1. S. M. Sze, *Physics of Semiconductor Devices*, Wiley, 1981.
2. S. R. Forrest, "Optical Detectors for Lightwave Communication," *Optical Fiber Telecommunications II*, 1988, pp. 569–599.
3. Z. S. Huang and T. Ando, "A Novel Amplified Image Sensor with a Si:H Photoconductor and MOS Transistors," *IEEE Trans. on Elect. Dev.*, vol. 37, no. 6, 1990, pp. 1432–1438.
4. J. C. Gammel, G. M. Metzger, and J. M. Ballantyne, "A Photoconductor Detector for High Speed Fiber Communication," *IEEE Trans. on Elect. Dev.* vol. ED-28, no. 7, 1981, pp. 841–849.
5. M. V. Rao, P. K. Bhattacharya, and C. Y. Chen, "Low Noise $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$: Fe Photoconductive Detectors for Optical Communications," *IEEE Trans. on Elect. Dev.* vol. ED-33, no. 1, 1986, pp. 67–71.
6. J. C. Gammel, H. Ohno, and J. M. Ballantyne, "High Speed Photoconductive Detectors Using GaInAs," *IEEE J. of Quantum Elect.* vol. QE-17, no. 2, 1981, pp. 269–272.
7. C. Y. Chen, Y. M. Pang, K. Alavi, A. Y. Cho, and P. A. Garbinski, "Interdigitated $\text{Al}_{0.48}\text{In}_{0.52}\text{AsGa}_{0.47}\text{In}_{0.53}\text{As}$ Photoconductive Detectors," *App. Phys. Lett.* 44, 1983, pp. 99–101.
8. M. V. Rao, G. K. Chang, and W. P. Hong, "High Sensitivity, High Speed InGaAs Photoconductive Detector," *Elect. Lett.* vol. 26, no. 11, 1990, pp. 756–757.
9. J. Degani, R. F. Leheny, R. E. Nahory, M. A. Pollack, J. P. Heritage, and J. C. DeWinter, "Fast Photoconductive Detector Using p- $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$ with Response to 1.7 μm ," *Appl. Phys. Lett.* 38, 1981, pp. 27–29.
10. S. M. Sze, *Semiconductor Devices—Physics and Technology*, Wiley, 1985.
11. G. H. Olsen, "Low Leakage, High Efficiency, Reliable VPE InGaAs 1.0–1.7 μm Photodiodes," *IEEE Elect. Dev. Lett.* vol. EDL-2, no. 9, 1981, pp. 217–219.
12. S. R. Forrest, "Performance of $\text{In}_x\text{Ga}_{1-x}\text{As}_y\text{P}_{1-y}$ Photodiodes with Dark Current Limited by Diffusion, Generation Recombination and Tunneling," *IEEE J. of Quantum Elect.* vol. QE-17, no. 2, 1981, pp. 217–226.
13. O. K. Kim, B. V. Dutt, R. J. McCoy, and J. R. Zuber, "A Low Dark-Current, Planar InGaAs *p-i-n* Photodiode with a Quaternary InGaAsP Cap Layer," *IEEE J. of Quantum Elect.* vol. QE-21, no. 2, 1985, pp. 138–143.
14. Y. Yoshida, Y. Hisa, T. Takiguchi, and Y. Komine, "Reduction of Surface Leakage Current in $\text{Cd}_{0.2}\text{Hg}_{0.8}\text{Te}$ Photodiode," *Proc. SPIE* vol. 972, 1988, pp. 39–43.
15. J. C. Flachet, M. Royer, Y. Carpentier, and G. Pichard, "Emission and Detection in the 1 to 3 μm Spectral Range with $\text{Hg}_{1-x}\text{Cd}_x\text{Te}$ Diodes," *Proc. SPIE* vol. 587, 1985, pp. 149–155.
16. R. E. Enstrom, P. J. Zanucchi, and J. R. Apert, "Optical Properties of Vapor-Grown $\text{In}_x\text{Ga}_{1-x}\text{As}$ Epitaxial Films on GaAs and $\text{In}_x\text{Ga}_{1-x}\text{P}$ Substrates," *J. of Appl. Phys.* vol. 45, no. 1, 1974, pp. 300–306.
17. H. W. Ruegg, "An Optimized Avalanche Photodiode," *IEEE Trans. on Elect. Dev.* vol. ED-14, no. 5, 1967, pp. 239–251.
18. K. Taguchi, T. Torikai, Y. Sugimoto, K. Makita, and H. Ishihara, "Planar Structure InP/InGaAsP/InGaAs Avalanche Photodiodes with Preferential Lateral Extended Guard Ring for 1.0–1.6 μm Wavelength Optical Communication Use," *J. of Lightwave Tech.* vol. 6, no. 11, pp. 1643–1655.
19. S. R. Forrest, R. G. Smith, and O. K. Kim, "Performance of $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}/\text{InP}$ Avalanche Photodiodes," *IEEE J. of Quantum Elect.* vol. QE-18, no. 12, 1982, pp. 2040–2047.
20. R. J. McIntyre, "The Distribution of Gains in Uniformly Multiplying Avalanche Photodiodes: Theory," *IEEE Trans. on Elect. Dev.* vol. ED-19, no. 6, 1972, pp. 703–713.
21. A. S. Grove, *Physics and Technology of Semiconductor Devices*, Wiley, New York, 1967.
22. S. L. Miller, "Avalanche Breakdown in Germanium," *Phys. Rev.* vol. 99, 1955, p. 1234.
23. M. E. Storm, "Coherent 2 μm Sources Burst into Windshear Detection," *Laser Focus* vol. 21, 1991, pp. 117–122.
24. A. M. Joshi, G. H. Olsen, V. S. Ban, E. Mykietyn, M. J. Lange, and D. R. Mohr, "Reduction of 1/f Noise in Multiplexed Linear $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$ Detector Arrays via Epitaxial Doping," *IEEE Trans. on Elect. Dev.* vol. 40, no. 2, 1993, pp. 303–308.
25. R. U. Martinelli and R. E. Enstrom, "Reliability of Planar InGaAs/InP Photodiodes Passivated with BoroPhosho-Silicate Glass," *J. of Appl. Phys.* vol. 63, no. 1, 1988, pp. 250–252.
26. A. K. Chin, F. S. Chen, and F. Ermanis, "Failure Mode Analysis of Planar Zinc-Diffused $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}$ *p-i-n* Photodiodes," *J. of Appl. Phys.* vol. 55, no. 6, 1984, pp. 1596–1606.

27. Y. Kuhara, H. Tercuchi, and H. Nishizawa, "Reliability of InGaAs/InP Long Wavelength *p-i-n* Photodiodes Passivated with Polyimide Thin Film," *IEEE J. of Lightwave Tech.* vol. LT-4, no. 7, 1986, pp. 933–937.
28. A. M. Joshi, G. H. Olsen, and S. R. Patil, "Reliability of InGaAs Detectors and Arrays," *Proc. SPIE* vol. 1580, 1991, pp. 34–40.
29. S. R. Forrest, V. S. Ban, G. Gasparian, D. Gay, and G. H. Olsen, "Reliability of Vapor Grown In_{0.53}Ga_{0.47}As/InP *p-i-n* Photodiodes With Very High Failure Activation Energy," *IEEE Elect. Dev. Lett.* vol. 9, no. 5, 1988, pp. 217–219.
30. C. P. Skrimshire, J. R. Farr, D. F. Sloan, M. J. Robertson, P. A. Putland, J. C. D. Stokoe, and R. R. Sutherland, "Reliability of Mesa and Planar InGaAs *p-i-n* Photodiodes," *IEEE Proc.* vol. 137, part J, no. 1, 1990, p. 7478.
31. R. R. Sutherland, J. C. D. Stokoe, C. P. Skrimshire, B. M. Macdonald, and D. F. Sloan, "The Reliability of Planar InGaAs/InP *p-i-n* Photodiodes with Organic Coatings for Use in Low Cost Receiver," *Proc. SPIE* vol. 1174, 1989, pp. 226–232.
32. B. F. Levine, C. G. Bethea, G. Hasnain, V. O. Shen, E. Pelve, R. R. Abbott, and S. J. Hsieh, "High Sensitivity Low Dark Current 10 μm GaAs Quantum Well Infrared Photodetectors," *Appl. Physics Lett.* vol. 56, no. 9, 1990, pp. 851–853.
33. B. F. Levine, C. G. Bethea, V. O. Shen, and R. J. Malik, "Tunable Long-Wavelength Detectors Using Graded Quantum Wells Grown by Electron Beam Source Molecular Beam Epitaxy," *Appl. Phys. Lett.* vol. 57, no. 4, 1990, pp. 383–385.
34. G. Hasnain, B. F. Levine, S. Gunapala, and N. Chand, "Large Photoconductive Gain In Quantum Well Infrared Photodetectors," *Appl. Phys. Lett.* vol. 57, no. 6, 1990, pp. 608–610.
35. E. Pelve, F. Beltram, C. G. Bethea, B. F. Levine, V. O. Shen, S. J. Hsieh, and R. R. Abbott, "Analysis of the Dark Current in Doped Well Multiple Quantum Well AlGaAs Infrared Photodetectors," *J. of Appl. Phys.* vol. 66, no. 11, 1989, pp. 5656–5658.
36. G. Hasnain, B. F. Levine, D. L. Sivco, and A. Y. Cho, "Mid-Infrared Detectors in the 3–5 μm Band Using Bound to Continuum State Absorption in InGaAs/InAlAs Multi-quantum Well Structures," *Appl. Phys. Lett.* vol. 56, no. 8, 1990, pp. 770–772.
37. S. N. Subbarao, D. W. Bechtel, R. J. Menna, J. C. Connolly, R. L. Camisa, and S. Y. Narayan, "2–4 GHz Monolithic Lateral *p-i-n* Photodetector and MESFET Amplifier on GaAs-on-Si," *IEEE Trans. on Microwave Theory and Tech.* vol. 38, no. 9, 1990, pp. 1199–1202.
38. S. Tiwari, J. Burroughs, M. S. Milshtein, M. A. Tischler, and S. L. Wright, "Lateral *p-i-n* Photodetectors with 18 GHz Bandwidth at 1.3 μm Wavelength and Small Bias Voltages," *Tech. Dig. of IEEE Int. Elect. Dev. Mtg.* 1991, pp. 421–425.
39. D. S. Chemla, "Quantum Wells for Photonics," *Physics Today*, May 1995, pp. 56–64.
40. D. D. Coon and R. P. G. Karunasiri, "New Mode of IR Detection Using Quantum Wells," *App Phys. Lett.* vol. 45, no. 6, 1984, pp. 649–651.
41. B. F. Levine, C. G. Bethea, J. W. Stayt, K. G. Glogovski, R. E. Leibenguth, S. D. Gunapala, S. S. Pei, and J. M. Kuo, "Long Wavelength GaAs/Al_xGa_{1-x}As Quantum Well Infrared Photodetectors (QWIPs)," *Proc. SPIE* vol. 1540, 1991, pp. 232–238.

25.9 ADDITIONAL READING

- Dereniak, E. L. and D. G. Crowe, *Optical Radiation Detectors*, Wiley, New York, 1984.
- Olsen, G. H., "Reliable Operation of Lattice Mismatched InGaAs Detectors on Silicon," *Tech. Dig. of IEEE Int. Elect. Dev. Mtg.*, 1990, pp. 145–147.

This page intentionally left blank.

DO NOT DUPLICATE

HIGH-SPEED PHOTODETECTORS

John E. Bowers and Yih G. Wey

*Department of Electrical and Computer Engineering
University of California
Santa Barbara, California*

26.1 GLOSSARY

A	area
A^{**}	modified effective Richardson constant
$A_e(A_h)$	electron (hole) ionization parameters
B	bit rate
C_j	junction capacitance
C_p	pad capacitance
$D_e(D_h)$	diffusion coefficient for electrons (holes)
E	electric field
$e_n(e_p)$	emission functions for electrons (holes)
F	frequency response
f	frequency
f_{3dB}	3-dB bandwidth
G	photoconductor gain
H	transfer function
h	Planck's constant
I_d	dark current
I_{dm}	multiplied dark current
I_{du}	unmultiplied dark current
I_{ph}	photocurrent
i	current
$\langle i_{na}^2 \rangle$	amplifier noise power
J	current density
J_{DIFF}	diffusion component of current density
J_{DRIFT}	drift component of current density

$J_e(J_h)$	electron (hole) component of current density
k	ratio of electron to hole ionization coefficient
k_b	Boltzmann constant
L	absorption layer thickness
$L_e(L_h)$	diffusion length for electrons (holes)
L_s	series inductance
M	multiplication factor
$M_n(M_p)$	electron (hole) initiated multiplication factor
m	electron mass
$n(p)$	electron (hole) density
P	input optical flux
q	electron charge
R	reflectivity
R_L	load resistance
R_s	series resistance
$R_1(R_2)$	reflectivity of the surface (substrate) mirror in a resonant detector
T	temperature
t	time
$t_e(t_h)$	transit time for electrons (holes)
V_B	breakdown voltage
V_j	junction voltage
v_e, v_h	electron and hole velocities
W	thickness of the depleted region
x	position
α	absorption coefficient
α_{FC}	free carrier absorption inside the absorption layer
α_{FCx}	free carrier absorption outside the absorption layer
α_{IB}	interband absorption
α_i	electron ionization rate
α_s	scattering loss
β	propagation constant in a waveguide photodetector
β_i	hole ionization rate
Γ	confinement factor
ϵ	permittivity
η	quantum efficiency
κ	coupling coefficient to the waveguide of a waveguide detector
λ	wavelength
$\mu_e(\mu_h)$	mobility for electrons (holes)
ν	optical frequency
σ	charge density
σ_d	noise current spectral density
τ_e, τ_h	trapping time at a heterojunction for electrons (holes)
τ_{tr}	transit time
$\phi_{bc}(\phi_{bv})$	barrier height for the conduction band (valence band)
ω	angular frequency

26.2 INTRODUCTION

High-speed photodetectors are required for telecommunications systems, for high-capacity local area networks, and for instrumentation. Many different detector structures and materials are required to cover this range of applications. Silicon is one of the most commonly used detector materials for wavelengths from 0.4 to 1.0 μm , while Ge photodetectors are used at longer wavelengths up to 1.8 μm . Silicon and germanium have indirect bandgaps at these wavelengths, which result in relatively small bandwidth-efficiency products. Consequently, for high-speed applications, direct bandgap semiconductors such as III-V materials are more important and are the focus of this chapter. $\text{Ga}_{0.47}\text{In}_{0.53}\text{As}$ with a cutoff wavelength of 1.65 μm is especially useful for telecommunication photodetectors at 1.3 and 1.55 μm . GaAs has a cutoff wavelength around 0.9 μm and is ideal for visible and near-infrared applications.

This chapter will focus on the physics and technology of high-speed photodetectors. The next section discusses the different structures that are possible. Later sections discuss some specific results and motivations for particular structures. The primary limitations to detector speed are discussed, followed by a description of specific photodetector systems. To supplement this chapter, the reader should refer to excellent chapters and articles written specifically about photodetectors,¹ photoconductors,² pin detectors,³ avalanche photodetectors,⁴ phototransistors,⁵ and receivers.^{6,7}

26.3 PHOTODETECTOR STRUCTURES

Many photodetector structures have been demonstrated and many more structures are possible. In this section, we classify the different possible structures and identify a few of the trade-offs. The optimum structure for a given application depends on the required bandwidth, efficiency, saturation power, linearity, ease of integration, and leakage current.

There are four common types of photodetectors: (1) photovoltaic detectors, (2) photoconductive detectors, (3) avalanche photodetectors (APD), and (4) phototransistors. Photovoltaic detectors have blocking contacts and operate under reverse bias. The blocking contact can be a reverse-biased p - n junction or a Schottky contact. The photoconductive detector has identical, nonblocking contacts such as two $n+$ regions in an undoped sample. The avalanche photodetector has a similar configuration to a photovoltaic detector except that it has a high-field region that causes avalanching and results in gain in the detector. Improvements to the basic APD design include separate avalanche and gain regions (SAM APDs), and staircase APDs to increase the ratio of electron to hole (or hole to electron) multiplication rate. Phototransistors are three-terminal devices which have an integrated electronic gain region.

The second criterion is the contact type and configuration. The photogenerated carriers may be collected by means of (1) a vertical current collector, often a p - n or Schottky junction, (2) an interdigitated metal-semiconductor-metal (MSM) structure, or (3) a laterally grown or etched structure. These options are illustrated in Fig. 1a. PIN junctions are usually formed during the growth steps and tend to have low leakage current and high reliability.^{8,9} Schottky junctions are simple to fabricate, but tend to have a large leakage current on narrow-gap semiconductors, such as InGaAs. MSM structures have the advantage of lower capacitance for a given cross-sectional area, but often have longer transit times, limited by the lithography capabilities possible in production. Experimental demonstrations with very fine lines (50 nm) have yielded high-speed devices with good quantum efficiencies. MSM detectors tend to be photoconductive detectors, but one could lower the capacitance for a given area of a PIN detector by using an interdigitated MSM structure with p and n regions under alternating metal fingers.

The third important aspect of photodetector design is the orientation of the light with respect to the wafer and the current collection region (Fig. 1b). Most commercial photodetectors are vertically illuminated and the device area is 10 μm in diameter or larger, which allows simple, high-yield packaging with single-mode optical fibers, or easy alignment to external bulk optics. The problem with this configuration is that the absorbing layer must be thin for a high-speed detector to keep

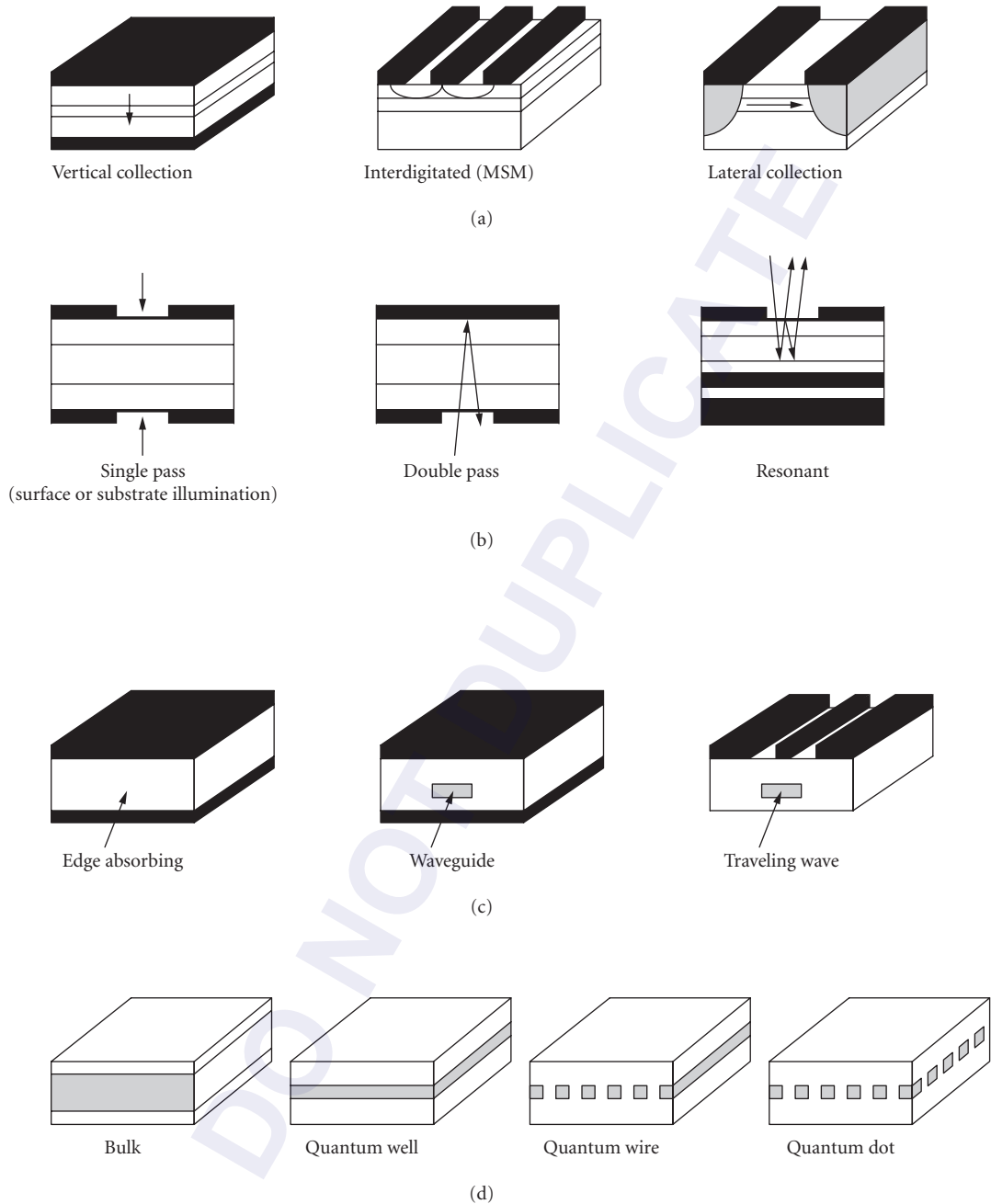


FIGURE 1 Schematic drawings of different types of photodetectors: (a) electrical configuration; (b) optical configuration—vertical illumination; (c) optical configuration—horizontal illumination; and (d) absorbing material.

the transit time of photogenerated carriers short. Consequently, the quantum efficiency is low, and single-pass vertically illuminated photodetectors tend to have bandwidth efficiency products around 30 GHz.³ Bandwidth efficiency products are discussed in greater detail in Sec. 26.5. The bandwidth efficiency product can be increased by allowing two passes by reflecting the light off a metal layer or dielectric mirror.¹⁰ Making a resonant cavity with multiple reflections at particular wavelengths should allow bandwidth efficiency products in excess of 100 GHz.^{11–13} As will be seen below, essentially 100 percent quantum efficiency is possible with bandwidths up to 20 GHz, so there is no need for resonant detectors unless the required bandwidth is above 20 GHz or wavelength selectivity is needed as in a wavelength division multiplexed (WDM) system.

The other class of optical inputs are horizontally illuminated photodetectors (Fig. 1c). The simplest configuration is an edge-illuminated detector. The primary problem with an edge-illuminated detector is that the light is not guided. Diffraction of the incident light causes absorption to occur outside of the high-field region, and slow diffusion tails in the impulse response occur. A solution to this problem is the waveguide detector, where an optical waveguide confines the light to the high-field absorption region.^{14–16} The waveguide efficiency product of this structure can be 100 to 200 GHz. However, it is limited by the capacitance of the structure, particularly if thin intrinsic layers are used for ultrahigh-speed devices. A solution to the capacitance limitation is a traveling wave photodetector where the incoming optical beam is velocity matched with the generated microwave signal.^{17,18} The bandwidth efficiency product is then limited only by loss on the electrical transmission lines, and bandwidth efficiency products of hundreds of GHz are possible. Traveling wave detectors and, to a lesser extent, waveguide detectors have the important advantage that the volume of the light absorption can be quite large and, consequently, these detectors have much higher saturation powers.¹⁹ Velocity matching in these structures requires quite narrow waveguides. Wu and Itoh²⁰ have suggested separating the parts of the optical waveguide with microwave delay lines to achieve velocity matching.

The fourth issue is the type of absorbing material (Fig. 1d) (1) bulk; (2) quantum well; (3) quantum wire; (4) quantum dot; (5 to 7) strained quantum well, wire, or dot; (8) *n-p-i-i* structure. The vast majority of commercial and experimental detectors use bulk material. However, quantum well detectors²¹ are becoming increasingly important in photonic integrated circuits (PICs) because the absorbing quantum well material is also used in other parts of the PIC such as the laser. Quantum wire photodetectors²² have potential advantages in terms of higher-bandwidth-efficiency products, but uniform quantum wires are rather difficult to fabricate. Quantum dot detectors have even higher peak absorption coefficients and more wavelength selectivity, but will probably have problems with slow impulse responses due to trapping of the carriers by the heterojunction. In other quantum-confined detectors, the carriers can be extracted along the quantum wire or well, and this problem can be avoided.²²

The final classification is by means of the lifetime of the material. Conventional detectors have material lifetimes of typically 1 ns and achieve speed by using high field for rapid carrier collection. A second approach is to use low temperature (LT) grown material which has a very short lifetime, perhaps as low as 1 ps. A third approach is to damage the material by means of ion implantation. The final approach is to grow or diffuse in traps into the material such as iron²³ or gold.

If we combine these classifications, we find that 2600 types of photodetectors are possible, and additional subgroups such as superlattice APDs or SAGM APDs increase the total even further. In reality, about 100 types of detectors have been demonstrated. One of the points of this section is that improvements in one type of detector, such as adding a resonant cavity to a PIN detector, can be applied to other types of detectors, such as adding a resonant cavity to an APD. In the following section, we discuss in more detail some of the real limitations to the speed of a detector, and then apply this knowledge to a few important types of detectors.

26.4 SPEED LIMITATIONS

Generally speaking, the bandwidths of most photodetectors are limited by the following factors: (1) carrier transit time, (2) RC time constant, (3) diffusion current, (4) carrier trapping at heterojunctions, and (5) packaging. These limiting factors will be discussed in turn with specific application to *p-i-n* photodiodes.

Carrier Transit Time

In response to light absorbed in a material, the photogenerated carriers in the active region will travel across the high-field region and then be collected by the electrodes. As an example, Fig. 2a shows the p - i - n structure with a photogenerated electron-hole charge sheet of density σ . In response to the electric field, the electron will travel to the right and the hole to the left. This induces a

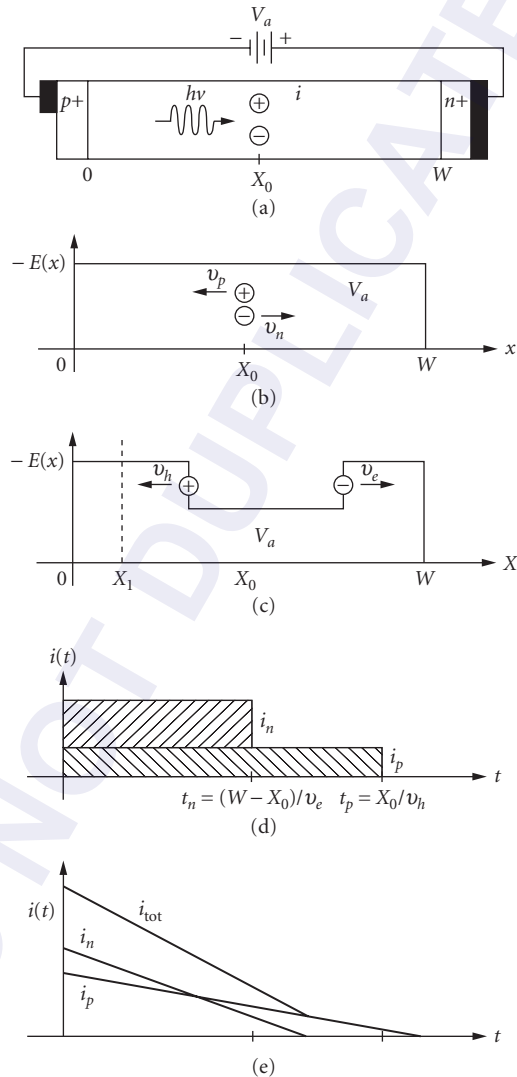


FIGURE 2 (a) Biased p - i - n structure. (b) Electrical field at the time when the electron-hole pairs are generated. (c) The perturbed electrical field due to the separated electron-hole pair. (d) The photocurrent due to single electron-hole pair. (e) The total photocurrent due to uniform illumination across the photodiode.

displacement current and reduces the internal electric field (Fig. 2b and c), which is the cause of saturation in photodetectors. From Gauss's law the difference in the electric field at the position of electron or hole is

$$\Delta E = \frac{q\sigma}{\epsilon} \quad (1)$$

where q is the electron charge, and ϵ is the permittivity. Due to the constant total voltage across the depletion region, the reduced electric field between the electron and hole will be compensated by the increased electric field outside. The rate of change of the electric field at the position $X = X_1$ is

$$\frac{\partial E}{\partial t} = -\frac{(v_e + v_h)\Delta E}{W} \quad (2)$$

where v_e and v_h are the saturation velocities for electrons and holes, respectively. The assumption of saturation velocities is valid at high fields. The displacement current is hence given by

$$i(t) = -\epsilon A \frac{\partial E}{\partial t} = \frac{qAv_e\sigma}{W} + \frac{qv_h\sigma A}{W} \quad (3)$$

The current consists of two components due to the electron and hole currents. The electron current lasts for a time duration of $(W - X_0)/v_e$ and hole current of X_0/v_h . This is shown in Fig. 2d. Here, we note that if the fast carrier (i.e., electron) travels a longer distance, then we have a shorter pulse. The total electron and hole currents are given by

$$i_e(t) = \frac{qv_e A}{W} \int_0^w n(x, t) dx \quad (4)$$

$$i_h(t) = \frac{qv_h A}{W} \int_0^w p(x, t) dx \quad (5)$$

where $n(x, t)$ and $p(x, t)$ represent the electron and hole densities in the depletion region. The total current is the sum of Eqs. (4) and (5).

RC Time Constant

The RC time constant is determined by the equivalent circuit parameters of photodiode. For example, the intrinsic response of the p - i - n diode can be modeled as a current source in parallel with a junction capacitor. The diode series resistance, parasitic capacitance, and load impedance form the external circuit. Figure 3 shows the equivalent circuit of the p - i - n photodiode. The junction capacitance is defined by the edge of the depletion region (or space charge region). The series resistance is due to the ohmic contacts and bulk resistances. In addition, the parasitic capacitance depends on the metallization

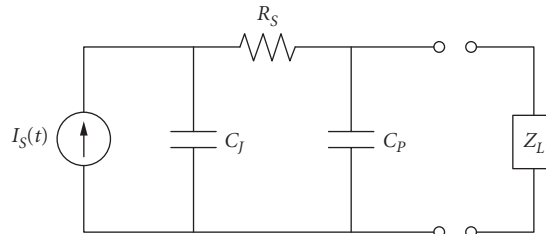


FIGURE 3 Equivalent circuit of a photodiode.

geometry. If the diode series resistance is R_s and a load resistance R_L is used to terminate the device, then the electrical 3-dB bandwidth can be approximated as

$$f_{3\text{dB}} = \frac{1}{2\pi(C_j + C_p)(R_L + R_s)} \quad (6)$$

If the photodiode is bonded by a section of gold wire, additional series inductance will be included in the load impedance. The 3-dB bandwidth due to parasitics in this case is then given in Ref. 3.

Diffusion Current

Diffusion current is important in detectors in which significant absorption occurs in regions outside the high-field region. This effect is reduced to some extent by recombination in these highly doped contact layers. Those carriers within about one diffusion length of the depletion region will have a chance to diffuse into the active region. This diffusion current will contribute a slow tail to the detector impulse response (Fig. 4). The electron diffusion current at the edge of the depletion region is given by

$$J_e = qD_e \frac{\partial n}{\partial x} = qD_e \frac{\Delta n}{L_e} \quad (7a)$$

and

$$J_h = -qD_h \frac{\partial p}{\partial x} = qD_h \frac{\Delta p}{L_h} \quad (7b)$$

where D_e (D_h) and L_e (L_h) are the diffusion coefficient and diffusion length, respectively, for electrons (holes). The diffusion process is a relatively slow process compared with the drift process. Assuming the photocarrier density is n , with the Einstein relation and Eq. (7a), the electron diffusion current can be written as

$$J_{\text{DIFF}} = qn\mu_e \left(\frac{kT/q}{L_e} \right) \quad (8)$$

and the drift current term for electron can be written as

$$J_{\text{DRIF}} = qn\mu_e \mathbf{E} \quad (9)$$

where μ_e is the electron mobility and \mathbf{E} is the electric field. For most devices the electric field inside the depletion region is an order of magnitude larger than $(kT/q)L_e$. For example, the hole diffusion

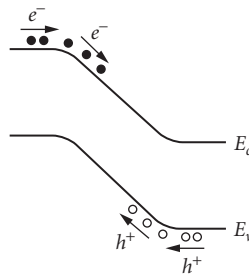


FIGURE 4 The origins of diffusion current.

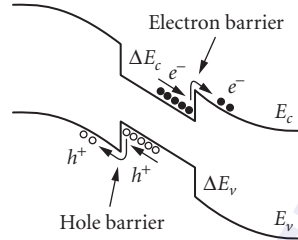


FIGURE 5 Heterostructure carrier trapping effect.

length for GaAs is typically around $10\ \mu\text{m}$ and electric field is very often over $10\ \text{kV/cm}$. However, the diffusion-current terms could last as long as the carrier lifetime and the charge content in the tail can be as large as the drift component due to the slow diffusion times. For high-speed detectors, the diffusion-current problem can be eliminated with a double-heterostructure design that limits the absorbing regions to the high-field intrinsic regions.²⁴

Carrier Trapping

Heterojunctions in photodetectors cause carrier trapping of electrons at conduction band discontinuities and trapping of holes at valence band discontinuities (Fig. 5). Hole trapping is a significant problem in long-wavelength photodetectors because of the large valence band discontinuity at the InP/InGaAs heterojunction. Usually, the emission rate is approximated by thermionic emission. If the interface deep-level recombination rate is significant, the total emission rate will be the sum of the two emission rates. The emission functions for electrons and holes are given by

$$e_n(t) = (1/\tau_e) \exp(-t/\tau_e) u(t) \quad (10a)$$

$$e_p(t) = (1/\tau_h) \exp(-t/\tau_h) u(t) \quad (10b)$$

where τ_e (τ_h) represents the emission time constant for electron (hole) and $u(t)$ is the step function. The rates of thermionic emission of trapped carrier are related to the Schottky barrier height due to the bandgap discontinuity:

$$1/\tau_e = B \exp(-\phi_{bc}/kT) \quad (11)$$

where B is a constant and ϕ_{bc} is the barrier height for the conduction band. The response of the carrier-trap current in time domain is often obtained by convolving an intrinsic current source with the emission function. Since the applied bias will reduce the barrier height, sufficient device bias therefore will increase the emission rate. In order to reduce the barrier height, superlattice or compositional grading is often added at the heterointerface.²⁴

Packaging

The external connections to the photodetector often limit the detector performance. Another problem is that the photodiode is a high impedance load and the device has a reflection coefficient close to unity. One solution to this problem is to integrate a matching resistor with the device.²⁴ This can usually be added using the lower contact layer without adding any additional mask or process steps. Figure 6 shows a Smith chart plot of the impedance of a typical photodetector along with the impedance of a device with an integrated matching resistor. A good match up to 40 GHz

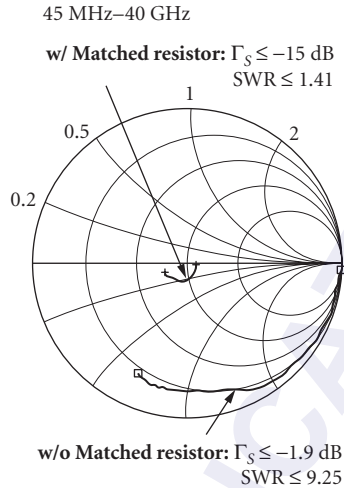


FIGURE 6 Smith chart of the impedance of typical photodiodes and photodiodes with integrated matching resistors.

is achieved. The disadvantage of a load resistor is the reduction in effective quantum efficiency by a factor of two since half of the photocurrent goes through the matching resistor. However, since the load resistance is now one-half, the RC time constant is also cut in half. Bandwidths in excess of 100 GHz have been achieved with quite large devices ($7 \times 7 \mu\text{m}^2$) in this way. A second problem with very high speed devices is the difficulty in building external bias circuits without resonances in the millimeter range. The necessary bias capacitor and load resistor can be integrated with *p-i-n* photodetector without adding any additional mask or process steps by using a large-area *p-i-n* region as the capacitor and using the lower contact layer as the series resistor.²³ A photograph of the device is shown in Fig. 7 along with the device performance. In this case, bandwidths in excess of 100 GHz were achieved.

Optical fiber alignment and packaging are now quite standard. Simplified alignment by means of holes etched in the substrate of back-illuminated photodiodes may allow passive alignment of optical fibers. The photodetectors must be antireflection coated to reduce the reflection to air or optical epoxy. Single dielectric layers are typically used to minimize the reflection at one wavelength. Braun et al.²⁵ have achieved minimum reflectivity at multiple wavelengths with one dielectric layer by using one of the semiconductor layers in a multiple antireflectivity design.

26.5 P-I-N PHOTODETECTORS

Vertically Illuminated *p-i-n* Photodiode

In order to increase the frequency response of the vertically illuminated *p-i-n* photodiode, the efficiency is always sacrificed. As the active layer thickness is reduced, the transit time decreases, and the optical absorption decreases, and there is a trade-off between the efficiency and speed. The external quantum efficiency for a surface-illuminated *p-i-n* diode is given by

$$\eta = (1 - R) \times (1 - e^{-\alpha L}) \quad (12)$$

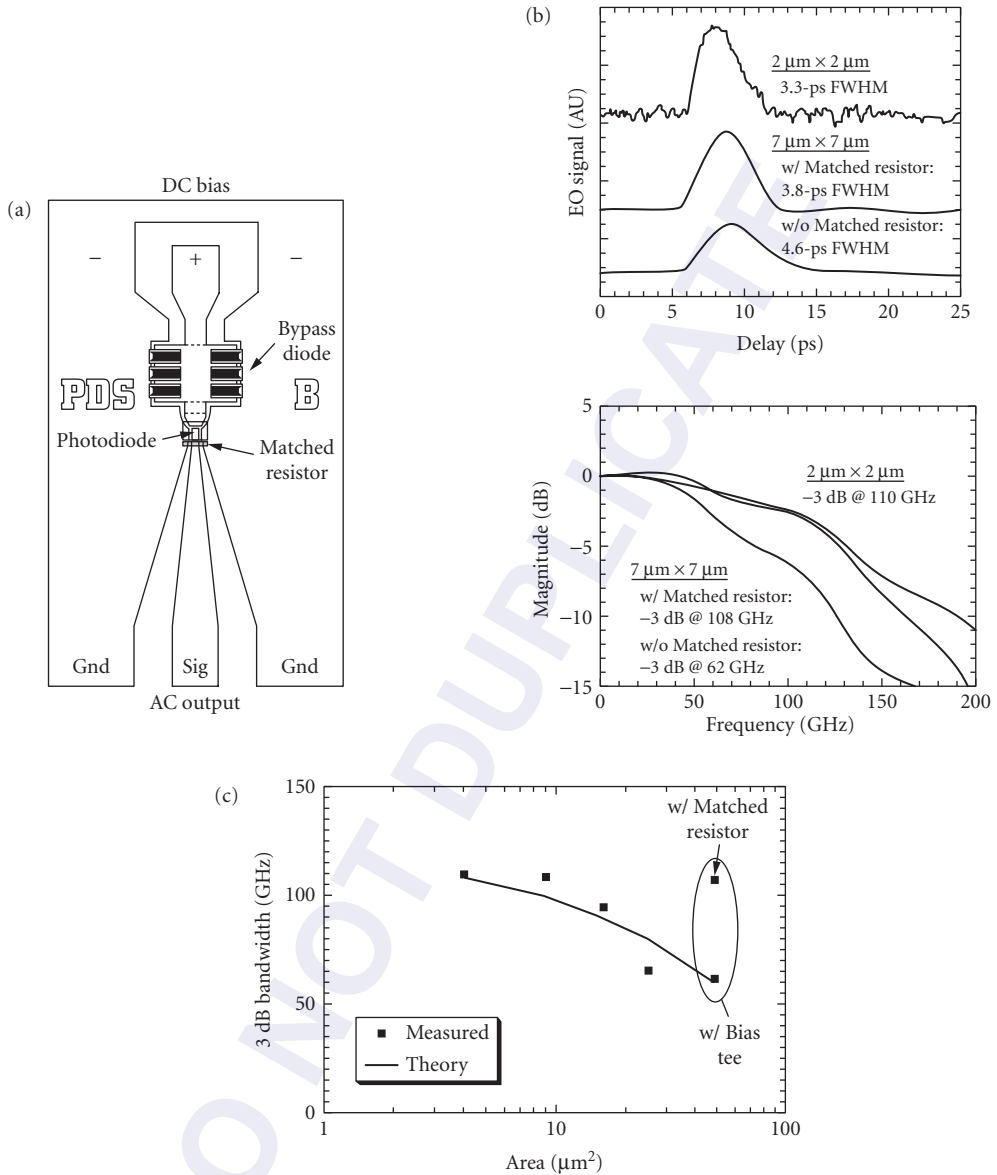


FIGURE 7 (a) Schematic diagram of a $p-i-n$ photodiode with integrated matching resistor and bias circuit. (b) Impulse response of a $2 \times 2 \mu\text{m}$ pin detector compared to $7 \times 7 \mu\text{m}$ detectors with and without matching resistors. (After Ref. 24.) (c) Dependence of measured bandwidth on detector area and comparison to the theory presented in the text.

where R is the surface reflection, α is the absorption coefficient, and L is the active layer thickness. Since the absorption coefficient is a function of wavelength $\alpha = \alpha(\lambda)$, usually α decreases as λ increases. Thus, the diode intrinsic response is wavelength dependent. We can easily see the effect of light absorption on the transit-time-limited bandwidth by comparing the transit-time response^{3,25} for two limiting cases for $\alpha L \rightarrow 0$ and $\alpha L \rightarrow \infty$ when $t_e = t_h = \tau_r$. The transit-time frequency response for a uniformly illuminated detector is

$$|F(\omega)_{\alpha L \rightarrow 0}| = \frac{2}{\omega \tau_r} \left[1 + \frac{\sin^2\left(\frac{\omega \tau_r}{2}\right)}{\left(\frac{\omega \tau_r}{2}\right)^2} - 2 \frac{\sin(\omega \tau_r)}{(\omega \tau_r)} \right]^{1/2} \quad (13)$$

For electron-hole pairs generated near the p side of the intrinsic region, electrons travel across the i region, and the frequency response is given by

$$|F(\omega)_{\alpha L \rightarrow \infty}| = \left| \frac{\sin\left(\frac{\omega \tau_r}{2}\right)}{\left(\frac{\omega \tau_r}{2}\right)} \right| \quad (14)$$

For these two limits, the transit-time-limited bandwidths are $f_{3\text{dB}(\alpha L=0)} = 0.45/\tau_r$ and $f_{3\text{dB}(\alpha L=\infty)} = 0.55/\tau_r$, respectively. For long-wavelength high-speed p - i - n diodes, the absorption layer is often very thin so that $1 - \exp(-\alpha L) \approx \alpha L$. The bandwidth efficiency product for transit-time-limited p - i - n diode is given by³

$$\eta f_{3\text{dB}} = 0.45 \alpha v_s \quad (15)$$

Figure 8 shows the calculated 3-dB bandwidth for GaInAs/InP p - i - n diodes on the device area versus thickness plane for wavelength $\lambda = 1.3 \mu\text{m}$. The horizontal axis is the active layer thickness (which

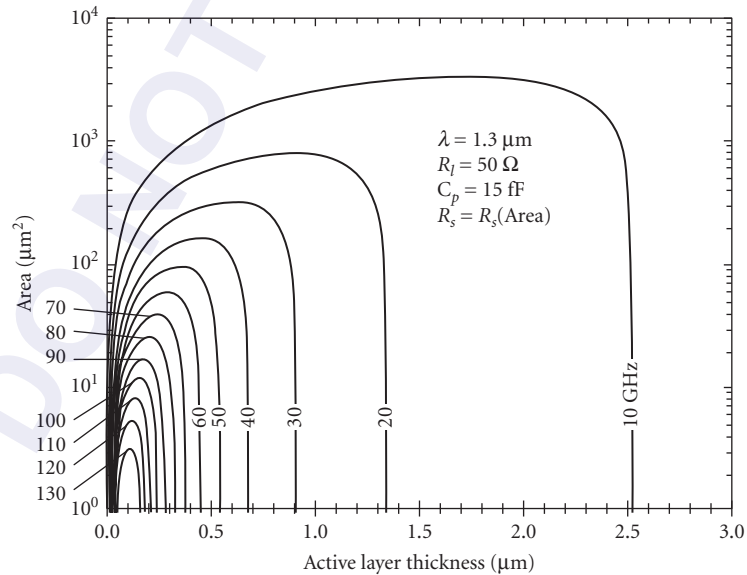


FIGURE 8 Calculated 3-dB bandwidth contours for a GaInAs pin vertically illuminated photodiode.

corresponds to the quantum efficiency) and the vertical axis is the device area. As we can see, when the device active layer thickness decreases, the quantum efficiency of p - i - n diode also decreases due to the insufficient light absorption in the active layer. The capacitance decreases as the device area decreases. Optimization of the device bandwidth is reached when transit-time-limited bandwidth approximately equals the RC limited bandwidth.

To minimize the bonding-pad capacitance, a semi-insulating substrate and thick polyimide layer are often used. Sometimes the series inductance of the bond wire is used to resonate the parasitic capacitance, and this results in a slightly peaked response with an increased 3-dB corner frequency. The electrical transfer function with series inductance is given by

$$H(\omega) = \frac{R_L}{[1 - \omega^2(R_S R_L C_J C_P + L_S(C_J + C_P))] - j[\omega(R_L(C_J + C_P) + R_S C_J) - \omega^3 R_S C_J C_P L_S]} \quad (16)$$

where L_S is the series inductance.

To achieve high detector bandwidth, *double heterostructure* InP/GaInAs/InP p - i - n photodiodes have been fabricated to reduce the diffusion-current problem. However, carrier trapping can limit the impulse response. This effect can be characterized by the emission function $e_{e,h}(t) = (1/\tau_{e,h}) \exp(-t/\tau_{e,h})$ where $\tau_{e,h}$ is emission time for electron (hole). The current-source response due to the electron and hole trapping at the heterointerfaces for p -side illumination is given by

$$\begin{aligned} \frac{i_s(\omega)}{i_s(0)} = & \frac{1}{(1 - e^{-\alpha L})} \left\{ \left(\frac{1 - e^{-j\omega\tau_e}}{j\omega\tau_e} - e^{-\alpha L} \frac{1 - e^{-\alpha L} e^{-j\omega\tau_e}}{j\omega\tau_e - \alpha L} \right) \left(\frac{1}{1 + j\omega\tau_e} \right) \right. \\ & \left. + \left(\frac{1 - e^{-\alpha L} e^{-j\omega\tau_h}}{j\omega\tau_e + \alpha L} - e^{-\alpha L} \frac{1 - e^{-j\omega\tau_h}}{j\omega\tau_h} \right) \left(\frac{1}{1 + j\omega\tau_h} \right) \right\} \quad (17) \end{aligned}$$

and for n -side illumination is given by

$$\begin{aligned} \frac{i_s(\omega)}{i_s(0)} = & \frac{1}{(1 - e^{-\alpha L})} \left\{ \left(\frac{1 - e^{-\alpha L} e^{-j\omega\tau_e}}{j\omega\tau_e + \alpha L} - e^{-\alpha L} \frac{1 - e^{-j\omega\tau_e}}{j\omega\tau_e} \right) \left(\frac{1}{1 + j\omega\tau_e} \right) \right. \\ & \left. + \left(\frac{1 - e^{-j\omega\tau_h}}{j\omega\tau_h} - e^{-\alpha L} \frac{1 - e^{-\alpha L} e^{-j\omega\tau_h}}{j\omega\tau_h - \alpha L} \right) \left(\frac{1}{1 + j\omega\tau_h} \right) \right\} \quad (18) \end{aligned}$$

where $\tau_{e,h}$ is the electron (hole) transit time. Other than the original p - i - n diode response, the extra terms $1/(1 + j\omega\tau_{e,h})$ are due to the trapping effect. For InGaAs/InP heterostructure p - i - n diodes, the valence band offset is larger than the conduction offset and the hole effective mass is much larger than the electron effective mass. Thus, hole trapping is worse than the electron trapping in an InGaAs/InP p - i - n diode.

Waveguide p - i - n Photodiode

The main advantages of waveguide detectors are the very thin depletion region resulting in a very short transit time and the long absorption region resulting in a high bandwidth photodetector with a high saturation power (Fig. 1c). Due to the thin intrinsic layer, it can often operate at zero bias.²⁷ The absorption length of a waveguide detector is usually *designed to be long enough* ($>5 \mu\text{m}$) to ensure full absorption. The waveguide structure design (Fig. 9) is often required to have low coupling loss due to modal mismatch and reasonable effective absorption coefficient. The external quantum efficiency of a waveguide p - i - n detector is^{28,29}

$$\eta = \kappa(1 - R) \frac{\Gamma \alpha_{\text{IB}}}{\alpha} (1 - e^{-\alpha L}) \quad (19)$$

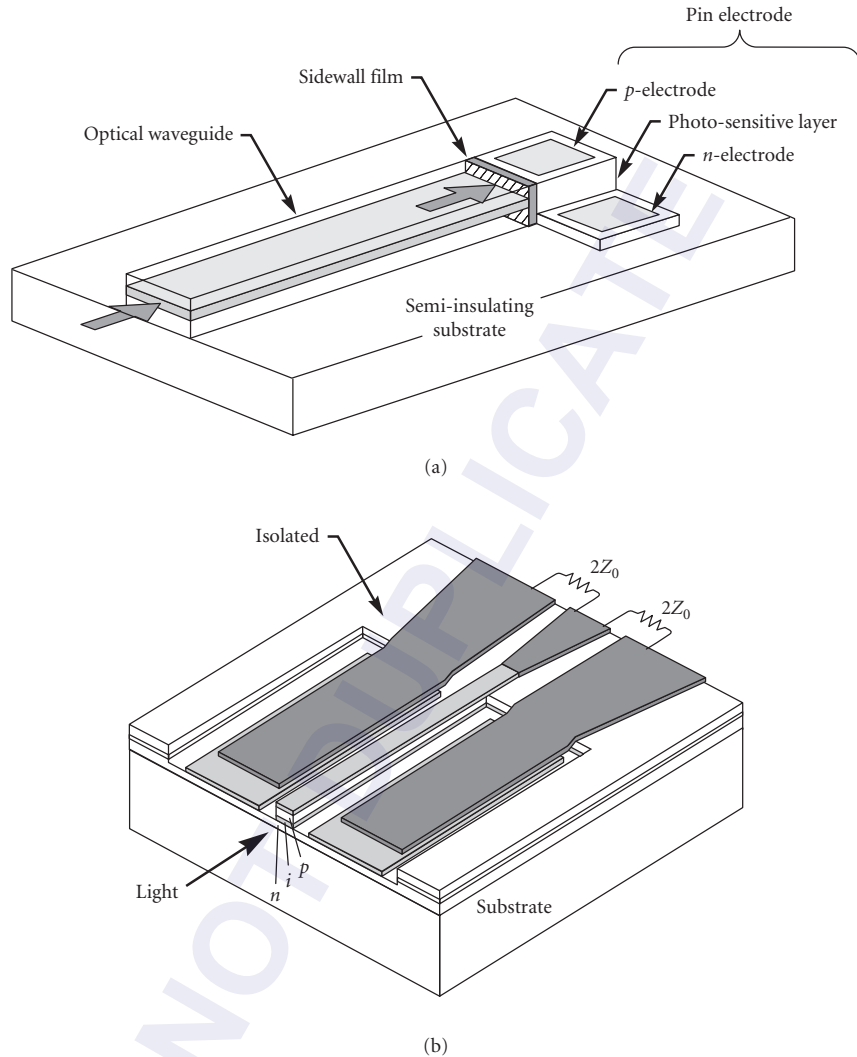


FIGURE 9 Schematic diagram of (a) a waveguide photodetector (after Ref. 14) and (b) a traveling wave photodetector (after Ref. 17).

where κ is the coupling efficiency due to the modal mismatch, Γ is the mode confinement factor, α_{IB} is the interband absorption. The loss coefficient α is given by

$$\alpha = \Gamma \alpha_{IB} + \Gamma \alpha_{FC} + (1 - \Gamma) \alpha_{FCx} + \alpha_s \quad (20)$$

where α_{FC} , α_{FCx} are the free carrier absorption loss inside and outside the absorption layer, α_s is the scattering loss. Kato et al.¹⁴ reported an InGaAs waveguide *p-i-n* diode with bandwidth of 40 GHz. The detector quantum efficiency is 44 percent at 1.55- μm wavelength. The coupling loss estimated by an overlap integral was 2.1 dB. To reduce the coupling loss, it is important to have a good design of the layer structure to reduce modal mismatch and to coat the facet with an antireflecting (AR) film.

Resonant p - i - n Photodiode

A resonant detector utilizes the multiple passes in a Fabry-Perot resonator to achieve high quantum efficiency with thin absorbing layers (Fig. 1b). Since the speed of light is about three orders of magnitude faster than the carrier velocities, the quantum efficiency can be increased without significant pulse broadening due to the effective optical transit time.

The schematic diagram of a resonant cavity-enhanced photodetector is shown in Fig. 10. The efficiency of the resonant detector is given by

$$\eta = \left[\frac{(1 + R_2 e^{-\alpha d})}{1 - 2\sqrt{R_1 R_2} e^{-\alpha d} \cos(2\beta L + \phi_1 + \phi_2) + R_1 R_2 e^{-2\alpha d}} \right] \times (1 - R_1) \times (1 - e^{-\alpha d}) \quad (21)$$

where R_1, R_2 are mirror reflectivities, ϕ_1, ϕ_2 are mirror phase shifts, $\beta = 2\pi/n\lambda$ is the propagation constant, and d is the thickness of active region. The quantum efficiency has its maximum when $2\beta L + \phi_1 + \phi_2 = 2m\pi$ ($m=1, 2, 3, \dots$) and the quantum efficiency is then

$$\eta = \left[\frac{(1 + R_2 e^{-\alpha d})}{(1 - \sqrt{R_1 R_2} e^{-\alpha d})^2} \right] \times (1 - R_1) \times (1 - e^{-\alpha d}) \quad (22)$$

Figure 11 shows the calculated resonant quantum efficiency versus normalized absorption coefficient αd .¹¹ High quantum efficiency is possible even from thin absorption layers.

However, in terms of the fabrication, the material growth, and the structure design, building a high-speed resonant photodetector is not a simple task. The required low resistance and low capacitance with incorporated mirror structure is difficult to achieve due to the significant resistance of the multiple heterojunction mirror stack.

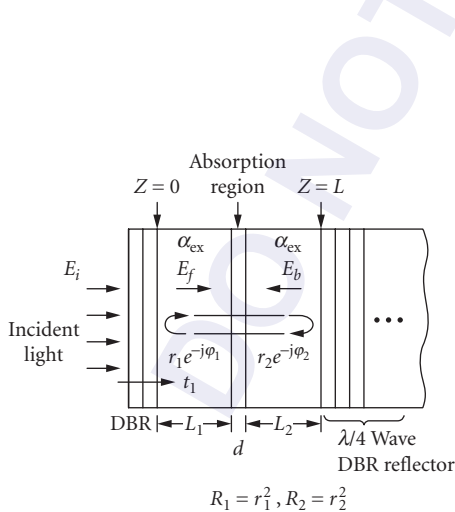


FIGURE 10 Schematic diagram of a resonant photodetector. (After Ref. 13.)

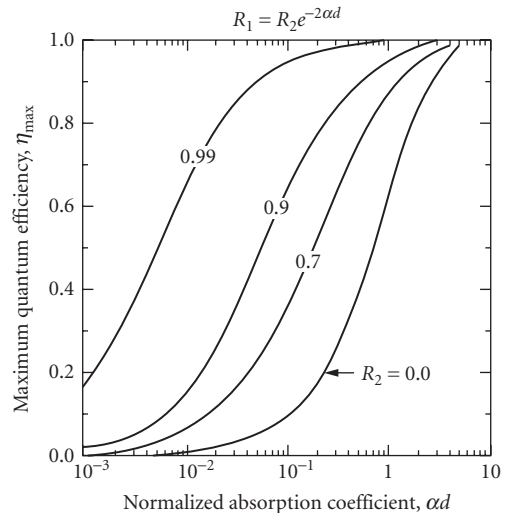


FIGURE 11 Dependence of quantum efficiency on mirror design in a resonant photodetector. (After Ref. 13.)

26.6 SCHOTTKY PHOTODIODE

Schottky photodiodes^{30–32} are especially attractive for integration with FETs and III-V integrated circuits because of their simple material structure and easy fabrication. Figure 12 shows the Schottky barrier structure. For front-illuminated devices, the metal is very thin so that the light can penetrate the metal with very little loss. The J - V characteristic of a Schottky diode is given by³³

$$J = J_0 \left[\exp\left(\frac{qV}{k_B T}\right) - 1 \right] \quad (23)$$

where

$$J_0 = A^{**} T^2 \exp\left(-\frac{q\Phi_b}{k_B T}\right) \quad (24)$$

and ϕ_b is the barrier height and A^{**} is the modified effective Richardson constant.³⁴

The dynamics of photogenerated carriers in a Schottky diode are similar to those of a p - i - n diode (Fig. 2). The dynamics of both electrons and hole must be included in the analysis of a Schottky photodiode, resulting in expressions for the Schottky diode response given in Eqs. (13) to (18) with the exception that there is no diffusion current from the metal layer. The equivalent circuit of a Schottky diode is the same as that of a p - i - n diode.

In high-speed applications, GaAs Schottky diodes in the short-wavelength region with bandwidths over 200 GHz have been reported.³⁶ These devices can be combined with FETs or sampling diodes.^{35,36} Figure 13 shows an integrated Schottky photodiode with a diode sampling circuit. A pair of short voltage pulses are generated by the nonlinear transmission line and a differentiator. The short voltage pulses are used to control the sampling capacitors to measure the photodiode signal. The sampled signal is then passed through a low-pass filter to extract the equivalent time domain waveform. Impulse responses of under 2 ps have been demonstrated in this way.³⁶

In the long-wavelength region, GaInAs Schottky diodes experience high dark current problems due to the relatively low Schottky barrier height at the metal/GaInAs interface.³⁷ An InP or quaternary layer is usually added at the interface in order to increase the Schottky barrier height.³⁸ A graded bandgap layer (e.g., GaInAsP) is then needed at the GaInAs/InP interface to reduce hole trapping.

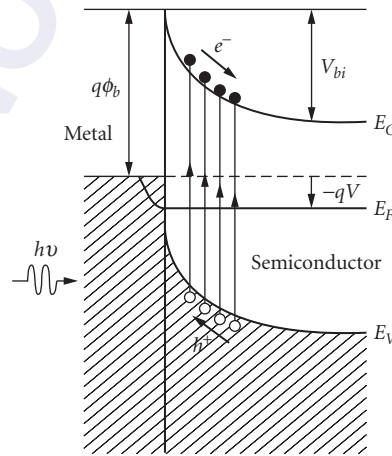


FIGURE 12 Schematic diagram of a Schottky barrier photodiode.

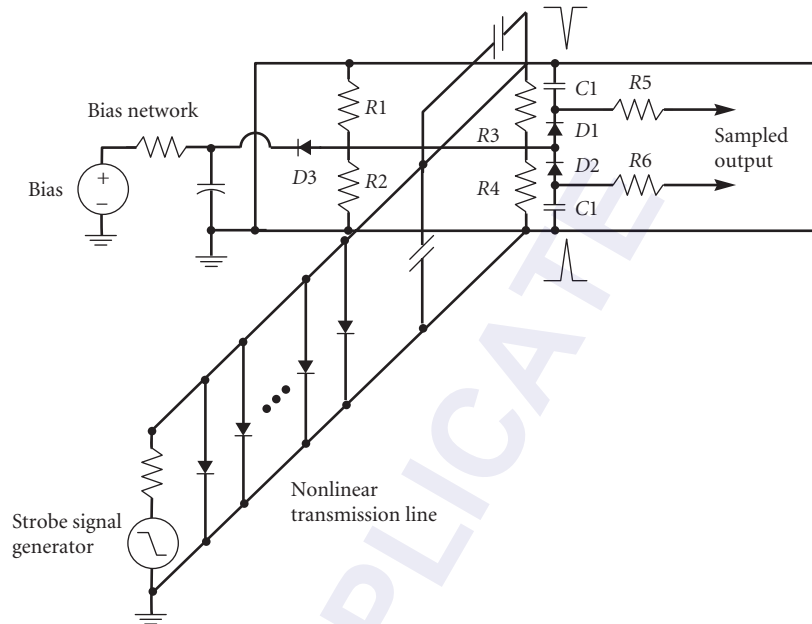


FIGURE 13 Integrated Schottky photodiode and sampling circuit. (After Ref. 35.)

26.7 AVALANCHE PHOTODETECTORS

High-speed avalanche photodiodes (APDs) are widely used in fiber communication. APDs with gain-bandwidth (GB) products in excess of 100 GHz have been reported.^{39–43} In long-wavelength applications, InGaAs/InP APDs are better than Ge APDs due to their lower dark current and lower multiplication noise. The Ge APD also has a limited spectral response at 1.55- μm wavelength. The maximum achievable GB product of InGaAs/InP APDs is predicted to be around 140 GHz,⁴⁴ while the gain-bandwidth product of Si APDs in the near-infrared region can have GB products of over 200 GHz.⁴⁵ InGaAs/InAlAs superlattice avalanche photodiodes have a lower ionization ratio ($k = 0.2$)⁴⁶ than bulk avalanche photodiodes, and lower noise and higher gain-bandwidth product can be achieved.

High-speed GaInAs/InP APDs make use of separated absorption and multiplication layers (SAM APD). Figure 14 shows the simplified one-dimensional APD structure. The narrow bandgap n -GaInAs layer absorbs the incident light. The layer is usually thick ($>1 \mu\text{m}$) to ensure high quantum efficiency. The electric field in the absorption layer is high enough for carriers to travel at saturated velocities, yet is below the field where significant avalanching occurs and the tunneling current is negligible. The wide bandgap InP multiplication layer is thin (a few tenths of a micron) to have shorter multiplication buildup time.^{47,48} The bias is applied to the fully depleted absorption layer in order to obtain effective carrier collection efficiency and, at the same time, electric field in the multiplication region must be high enough to achieve avalanche gain. A guard ring is usually added to prevent premature avalanche breakdown (or microplasma) at the corner of the diffusion edge. To reduce the hole pileup effect, a graded bandgap layer (e.g., superlattice or compositional grading) is often added at the heterointerface between the absorption layer and multiplication layer. This is the so-called separated absorption, grading, multiplication avalanche photodiode (SAGM APD).

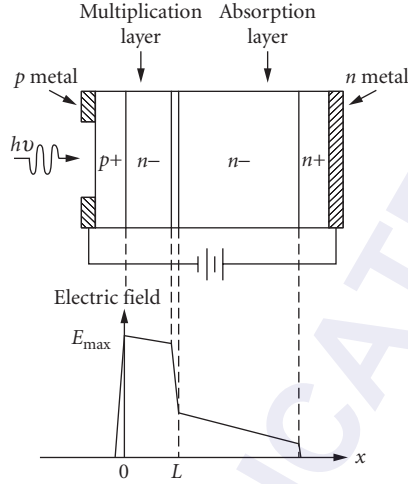


FIGURE 14 Schematic diagram of a SAGM APD.

The multiplication process in APDs can be described by the electron and hole ionization coefficients α_i and β_i . The field dependencies of ionization coefficients for electrons and holes are given by

$$\alpha_i(x) = A_e \exp(-B_e/E(x)) \quad (25a)$$

$$\beta_i(x) = A_h \exp(-B_h/E(x)) \quad (25b)$$

where $A_{e,h}$ and $B_{e,h}$ are constant parameters.⁴⁹ Since the electric field is generally position-dependent, the ionization coefficients are also position-dependent. With Eq. (25a, b) and the electric field distribution, the position-dependence of ionization coefficients can be derived. The multiplied photocurrent in the avalanche region ($0 \leq x \leq W$) including injected electron current density $J_n(0)$, injected hole current density $J_p(0)$, and photo-generation of electron-hole pairs $g(x)$ was derived by Lee et al.⁵⁰ The total photocurrent density is given by

$$J = \frac{J_p(w) \exp\left[-\int_0^w (\alpha_i - \beta_i) dx\right] + J_n(0) + q \int_0^w g(x) \exp\left[-\int_0^x (\alpha_i - \beta) dx'\right] dx}{1 - \int_0^w \alpha_i \exp\left[-\int_0^w (\alpha_i - \beta_i) dx'\right] dx} \quad (26)$$

The electron-initiated and hole-initiated multiplication factors, M_n and M_p , can be obtained by putting $J_p(w) = g(x) = 0$ and $J_n(0) = g(x) = 0$, respectively in Eq. (26):

$$M_n = \frac{J}{J_n(0)} = \frac{1}{1 - \int_0^w \alpha_i \exp\left[-\int_0^w (\alpha_i - \beta_i) dx'\right] dx} \quad (27a)$$

$$M_p = \frac{J}{J_p(w)} = \frac{\exp\left[-\int_0^w (\alpha_i - \beta_i) dx\right]}{1 - \int_0^w \alpha_i \exp\left[-\int_0^w (\alpha_i - \beta_i) dx'\right] dx} \quad (27b)$$

The bandwidth of an APD is limited by the device RC time constant when the multiplication gain M is low (i.e., $M < \alpha_i/\beta_i$). As the multiplication gain increases above the ratio of the electron and hole ionization coefficients (i.e., $M > \alpha_i/\beta_i$), the avalanche buildup time becomes the dominant limitation on 3-dB bandwidth and the product of the multiplication gain and 3-dB bandwidth reaches a constant. The multiplication factor M as a function of frequency was derived by Emmons⁵¹ and is given by

$$M(\omega) \approx \frac{M_o}{\{1 + \omega^2 M_o^2 \tau_1^2\}^{1/2}}, \quad M_o > \alpha_i/\beta_i \quad (28a)$$

$$\tau_1 \approx N(\alpha_i/\beta_i)\tau \quad (28b)$$

where τ_1 is the effective transit time, τ is the multiplication-region transit time, and $N(\beta_i/\alpha_i)$ is a number changing between 1/3 and 2 as β_i/α_i varies from 1 to 10^{-3} . The dc multiplication factor M_o is given by Miller.⁵²

$$M_o = \frac{1}{1 - (V_j/V_B)^n} \quad (29)$$

where V_B is the breakdown voltage, V_j is the junction voltage, and n is an empirical factor ($n < 1$).

The total APD dark current consists of two components. I_{du} is the unmultiplied current which is mainly due to the surface leakage current. I_{dm} is the bulk dark current experiencing the multiplication process. The total dark current is expressed by

$$I_d = I_{du} + MI_{dm} \quad (30)$$

where M is the avalanche gain. The noise current spectral density due to the dark current is given by

$$\sigma_d^2 = 2qI_{du} + 2qI_{dm} M^2 F(M) \quad (31)$$

where $F(M)$ is the avalanche excess noise factor derived by McIntyre.⁵³ Excess noise factors for electron-initiated or hole-initiated multiplication are given by

$$F(M) = F_e(M) = [kM + (1-k)(2-1/M)] \quad (32a)$$

$$F(M) = F_h(M) = \left[\frac{1}{k}M + \left(1 - \frac{1}{k}\right)(2-1/M) \right] \quad (32b)$$

where k is the ratio of the ionization coefficient of holes to electrons ($k = \beta_i/\alpha_i$), and k is assumed to be a constant independent of the position. Figure 15 shows the sensitivity of an APD receiver as a function of k_{eff} which is obtained by weighting the ionization rates over the electric field profile. From Fig. 15 we can see that the smaller the k factor is, the smaller the noise factor is and the better the receiver sensitivity is. Ge APDs have k values close to unity (0.7 to 1.0). GaInAs/InP APDs using an InP multiplication region have $1/k$ values from 0.3 to 0.5. Silicon is an excellent APD material since its k value is 0.02. Therefore, a Si APD has an excellent low dark current noise density and is predominantly used at short wavelengths compared with Ge APDs and GaInAs/InP APDs which are used at longer wavelengths.

In optical receiver applications,^{54,55} the photodetector is used with a low-noise amplifier. The dark current noise power is given by

$$\langle i_{nd}^2 \rangle = 2qI_{du}BI_2 + 2qI_{dm}M^2F(M)BI_2 \quad (33)$$

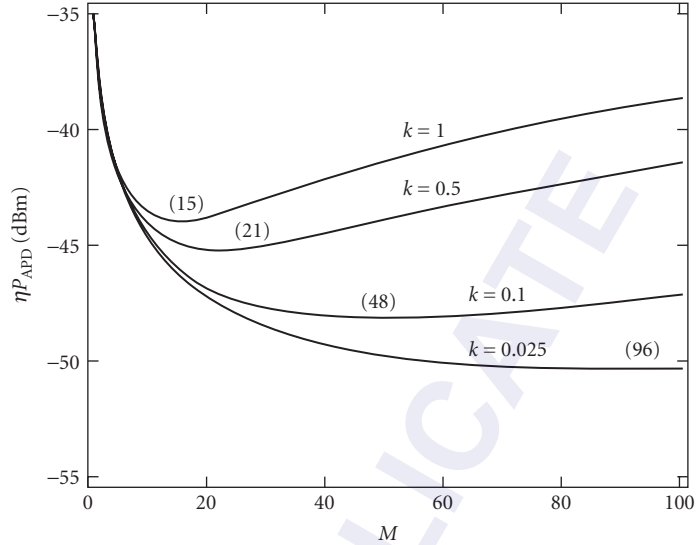


FIGURE 15 Dependence of APD receiver sensitivity on β/α in SAGM APDs. (After Ref. 7.)

where B is the receiver bit rate and I_2 is a parameter depending on the input optical pulse shape. The receiver sensitivity penalty⁵⁶ is given in terms of parameter ε_N ,

$$\bar{\eta P} = (1 + \varepsilon_N) \bar{\eta P}_o \quad (34)$$

where $\bar{\eta P}_o$ is the sensitivity with zero dark current. For example, $\varepsilon_N = 0.023$ for a 0.1-dB penalty. The maximum allowable dark current for a given sensitivity for a p - i - n FET receiver is

$$I_{du} = \frac{\varepsilon_N (2 + \varepsilon_N)}{2qBI_2} \langle i_{na}^2 \rangle \quad (35)$$

where $\langle i_{na}^2 \rangle$ is the amplifier noise power and is proportional to B^3 above 100 MBits/s. Therefore, the maximum allowable dark current is proportional to B^2 . For APD receivers, the maximum allowable dark current I_{dm} as a function of bit rate can be approximated by assuming that the sensitivity penalty is within 1 or 2 dB and optimum gain is constant. In Fig. 16, as we can see, the dark current is proportional to B at lower bit rates and $B^{1.25}$ at higher bit rates. So, as the bit rate increases, the maximum allowable dark current increases.

26.8 PHOTOCONDUCTORS

High-speed photoconductors⁵⁷⁻⁶² have become more important not only because of their simplicity in fabrication and ease of integration with MESFET amplifiers but also because of their useful applications for photodetector and photoconductor sampling gates. Usually the photoconductive film has a high density of defects with the trap energy levels deep within the bandgap to shorten the material lifetime and the detector impulse response. The characteristics of the photoconductive films include (1) high resistivity due to the fact that Fermi level is pinned at the midgap, (2) enhanced optical absorption for photon energy below the bandgap due to the introduction of new bandgap states, and (3) easy fabrication of ohmic contacts possibly due to the enhancement of tunneling through the narrow Schottky barrier with a pinned Fermi level.

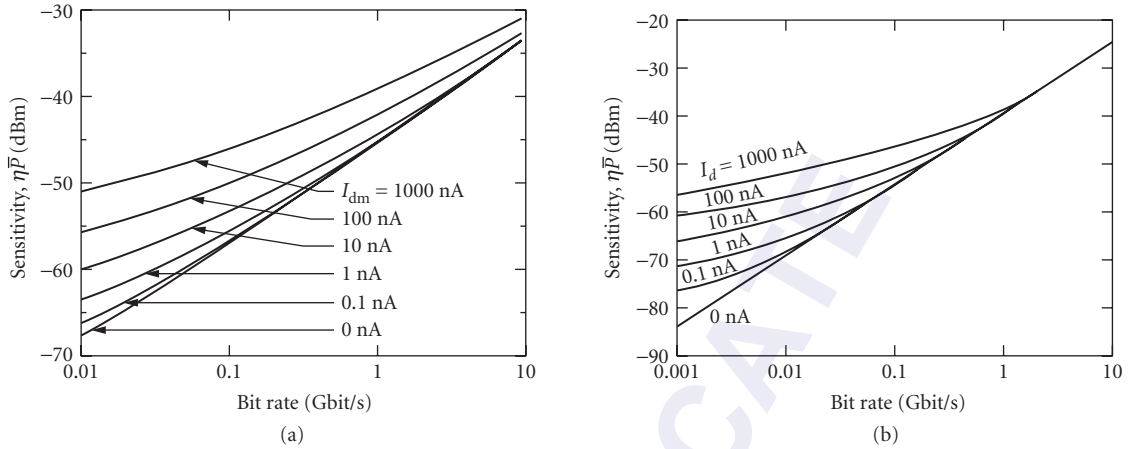


FIGURE 16 Dependence of receiver sensitivity on dark current for an (a) SAGM APD and (b) PIN detector, each with a GaAs FET preamplifier. (After Ref. 6.)

Figure 17 shows a typical photoconductor on a microstrip line structure. The photoconductive film is formed on top of a semi-insulating substrate. A microstrip transmission line consists of microstrip electrodes on top and ground plane on bottom. Under a steady-state illumination, the photogenerated carrier will experience high electrical field and travel to the electrodes. The photo-current is

$$I_{\text{ph}} = \frac{q\eta GP}{h\nu} \quad (36)$$

where η is the external quantum efficiency, G is the photoconductor gain, and P is optical input flux. The photoconductor gain G is given by

$$G = \frac{\tau}{\tau_{\text{tr}}} \quad (37)$$

which is the ratio of carrier lifetime τ to the carrier transit time τ_{tr} . The frequency response of a photoconductive detector is plotted in Fig. 18 for different material lifetimes. In a detector without damage sites, the gain can be quite large at low frequencies. We can see from this figure how the

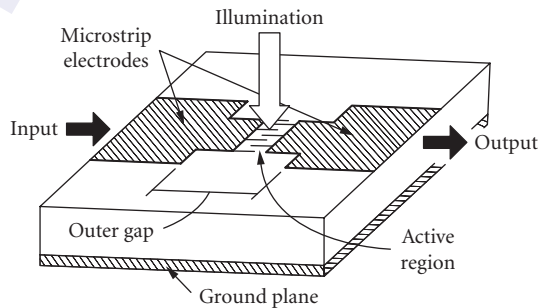


FIGURE 17 Schematic diagram of a high-speed photoconductor.

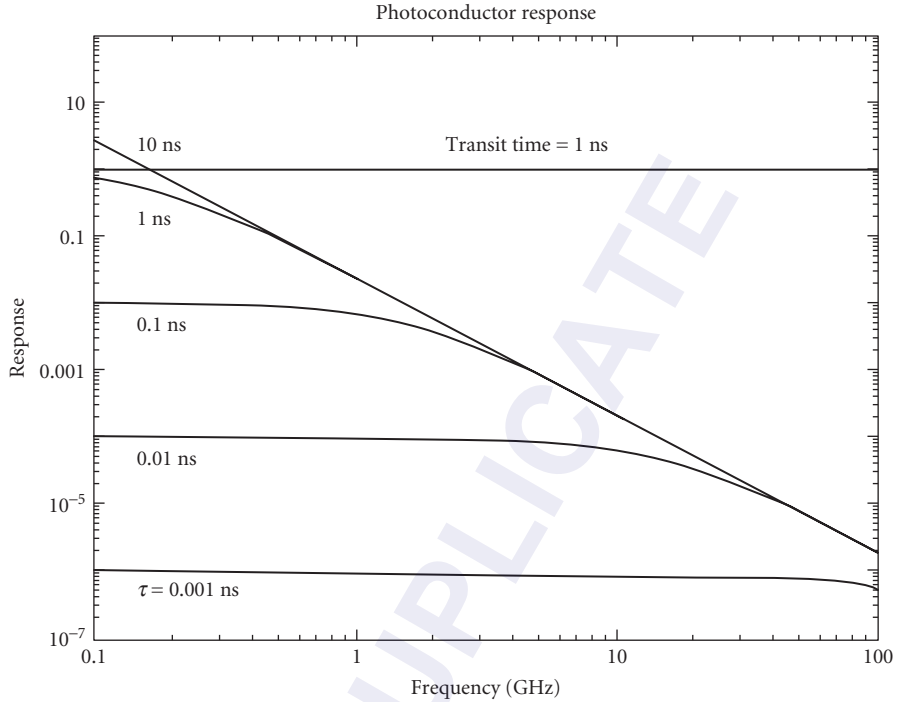


FIGURE 18 Frequency response of a photoconductor.

increased bandwidth is achieved at the expense of quantum efficiency. Using smaller finger separations, higher quantum efficiency can be achieved for a particular bandwidth.

The standard microstrip line configuration has reflection problems in the thickness direction of the substrate and the dispersion characteristics of a microstrip line is worse than that of a coplanar stripline.⁵⁷ Coplanar striplines with “sliding-contact” excitation can have zero capacitance to first order.⁵⁸ The photoconductor using coplanar stripline has been very successful in generating short electrical pulses. To measure the short electric pulse, several techniques can be used such as photoconductor sampling or electro-optic sampling. Both of the above techniques can provide subpicosecond resolution. The coplanar strip line configuration with sliding contact and sampling gate is shown in Fig. 19a. The equivalent circuit is shown in Fig. 19b.⁵⁸ The infinite capacitances represent that the line extends without end in both directions. The generated electrical signal due to the time-varying resistance $R_s(t)$ is

$$V_{\text{out}}(t) = V_b \frac{Z_o}{Z_o + R_s(t) + R_c} \quad (38)$$

where R_c is the contact resistance. If the excitation intensity is sufficiently low to keep $R_s(t) \gg Z_o$, then

$$V_{\text{out}}(t) = V_b \frac{Z_o}{R_s(t) + R_c} \quad (39)$$

The photoconductor resistance $R_s(t)$ can be related to the photoexcited electron-hole pair density $n(t)$:⁵⁹

$$R_s(t) = \frac{L}{qn(t)(\mu_e + \mu_h)wd_e} \quad (40)$$

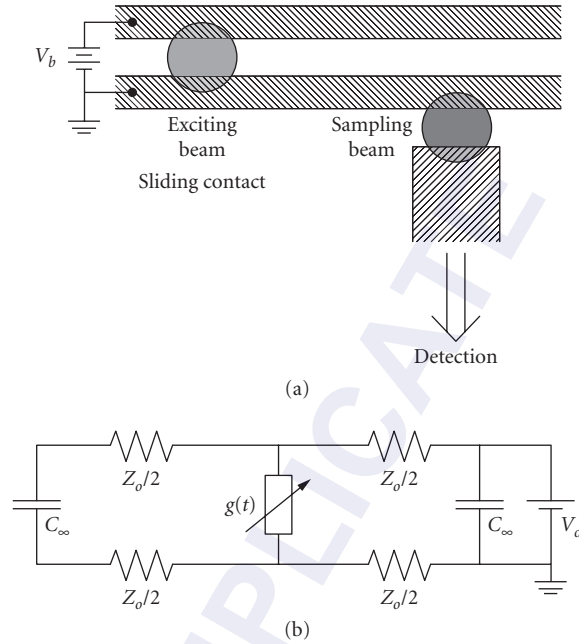


FIGURE 19 (a) Coplanar circuit layout of a photoconductor with sliding contact and (b) Equivalent circuit.

where L is the gap width, w is the width of photoconductive volume and d_e is the effective absorption length. When the pulse width is of the same order of magnitude as carrier lifetime and much shorter than the carrier transit time across the switch gap, the electron-hole pair density is given by

$$n(t) = e^{-t/\tau} \int_0^t e^{t'/\tau} \frac{\eta P_o(t')(1-R)}{h\nu w d_e} dt' \quad (41)$$

Here, we notice that the carrier density is an exponential decay function, so $G(t) = 1/R_s(t)$ is also an exponential decaying function with a time constant τ . This can be explained as a result of the convolution of the laser pulse with an exponential function with carrier life time τ .

The low temperature (LT) grown GaAs⁶⁰ can have both high carrier mobility and subpicosecond carrier lifetime when being compared with that of the ion-implanted photoconductor.⁶¹ The dislocation density in LT GaAs is about the same as that in GaAs epitaxial layer grown at normal substrate temperatures such that the LT GaAs has a mobility as high as that of the bulk material. The resistivity of the LT GaAs is greater than that of semi-insulating GaAs ($>10^{17} \Omega\text{-cm}$) due to its high deep-level concentration. The LT GaAs photoconductive-gap switch in a coplanar strip transmission line configuration has obtained a 1.6-ps (FWHM) response with a 3-dB bandwidth of 220 GHz. Chen et al.⁶² reported a high-speed photodetector utilizing LT GaAs MSM photoconductor. To achieve reasonable quantum efficiency and high-speed response, the optimum design requires carrier transit time approximately equal to carrier lifetime. With this requirement, the carriers not collected fast enough by the electrodes will be consumed by recombination. The response of a 0.2- μm finger and space MSM photodetector was measured by electro-optic sampling system. A 1.2 ps (FWHM) response with a 3-dB bandwidth of 350 GHz is obtained.

26.9 SUMMARY

Photodetector performance has steadily improved over the past decade. High-speed detectors are now available at a variety of wavelengths from 1.65 to 0.4 μm . MSM photoconductors have demonstrated the shortest impulse responses of under a picosecond. For applications that require high speed and high efficiency, the best results have been obtained using two passes through a p - i - n photodetector (30 percent quantum efficiency with 110-GHz bandwidth). Many applications require a high saturation power, and waveguide photodetectors have achieved the best results (20-GHz bandwidth with 0.5-A/W responsivity and 10-mW saturation power). Traveling wave photodetectors appear to offer the ultimate results in high-speed, high-responsivity, high-saturation power detectors. The combination of high-speed photodetectors with optical amplifiers is resulting in superb sensitivity of all bit rates, but requires the fabrication of high-speed photodetectors with at least 10-dBm saturation power.

An increasing amount of attention is being paid to integrating high-speed photodetectors with electronic and photonic circuits. Integration with electronic circuits increases the performance by eliminating the parasitics and limited bandwidth of bonding pads, wires, and connectors. Integration with optical waveguides decreases the optical loss associated with coupling from one device to another and reduces the packaging cost. Integration of photodetectors with optical amplifiers and wavelength tuning elements is a particularly important research direction.

26.10 REFERENCES

1. T. P. Lee and T. Li, "Photodetectors," *Optical Fiber Telecommunications*, S. E. Miller and A. G. Chynoweth (eds.), Academic Press, New York, 1979.
2. D. H. Auston, *Picosecond Optoelectronic Devices*, chap. 4, C. H. Lee (ed.), Academic Press, New York, 1984, pp. 73–117.
3. J. E. Bowers and C. A. Burrus, Jr., "Ultrawide-Band Long-Wavelength p - i - n Photodetectors," *J. Lightwave Tech.*, vol. LT-5, no. 10, October 1987.
4. F. Capasso, "Physics of Avalanche Photodiodes," *Semiconductors and Semimetals*, 22D, W. T. Tsang (ed.), Academic Press, New York, 1985.
5. J. Campbell, *Semiconductors and Semimetals*, 22D, W. T. Tsang (ed.), Academic Press, New York, 1985.
6. B. Kasper, "Receiver Design," *Optical Fiber Telecommunications II*, S. E. Miller and I. P. Kaminow (eds.), Academic Press, Boston, 1988.
7. S. R. Forrest, "Avalanche Photodetector Receiver Sensitivity," *Semiconductors and Semimetals*, vol. 22D, *Lightwave Communications Technology: Photodetectors*, W. Tsang (ed.), Academic Press, New York, 1985.
8. S. R. Sloan, "Processing and Passivation Techniques for Fabrication of High Speed InP/InGaAs/InP Mesa Photodetectors," *Hewlett Packard Journal*, Oct. 1989, p. 69.
9. K. Carey, S. Y. Wang, J. S. C. Chang, and K. Nauka, "Leakage Current in GaInAs/InP Photodiodes Grown by OMVPE," *J. Crystal Growth*, vol. 98, 1989, p. 90.
10. Y. G. Wey, D. L. Crawford, K. Giboney, J. E. Bowers, M. J. Rodwell, P. Silvestre, M. J. Hafich, and G. Y. Robinson, "Ultrafast Graded Double-Heterostructure GaInAs/InP Photodiode," *Appl. Phys. Lett.*, vol. 58, no. 19, 1991, p. 2156.
11. M. S. Unlu, K. Kishino, J. Chyi, L. Aresenault, J. Reed, and S. N. Mohammad, "Resonant Cavity Enhanced AlGaAs/GaAs Heterojunction Phototransistors with an Intermediate InGaAs Layer in the Collector," *Appl. Phys. Lett.*, vol. 57, no. 8, 20 Aug. 1990, p. 750.
12. A. Chin and T. Y. Chang, "Enhancement of Quantum Efficiency in Thin Photodiodes through Absorptive Resonance," *J. Lightwave Tech.*, vol. 9, no. 3, March 1991, p. 321.
13. K. Kishino, M. S. Unlu, J.-I. Chyi, J. Reed, L. Aresenault, and H. Morkoc, "Resonant Cavity-Enhanced (RCE) Photodetectors," *IEEE J. Quantum Electron.*, vol. 27, no. 8, Aug. 1991, p. 2025.
14. K. Kato, S. Hata, A. Kozen, J. Yoshida, and K. Kawano, "High-Efficiency Waveguide InGaAs Pin Photodiode with Bandwidth of over 40 GHz," *IEEE Photon. Tech. Lett.*, vol. 3, no. 5, May 1991, p. 473.
15. D. Wake, S. N. Judge, T. P. Spooner, M. J. Harlow, W. J. Duncan, I. D. Henning, and M. J. O'Mahony, "Monolithic Integration of 1.5 μm Optical Pre-amplifier and PIN Photodetector with a Gain of 20 dB and a Bandwidth of 35 GHz," *Electron. Lett.*, vol. 26, no. 15, July 19, 1990, pp. 1166–1168.

16. R. J. Deri, N. Yasuoka, M. Makiuchi, H. Hamaguchi, O. Wada, A. Kuramata, and R. J. Hawkins, "Integrated Waveguide/Photodiodes with Large Bandwidth and High External Quantum Efficiency," *IEEE Photon. Technol. Lett.*, vol. 2, 1990, pp. 496–498.
17. K. S. Giboney, M. J. W. Rodwell, and J. E. Bowers, "Travelling-Wave Photodetectors," *Photon. Tech. Lett.*, vol. 4, no. 12, Dec. 1992, pp. 1363–1365.
18. H. Taylor, O. Eknoyan, C. S. Park, K. N. Choi, and K. Chang, "Traveling Wave Photodetectors," *SPIE Proc. on Optoelectronic Signal Processing for Phased Array Antennas II*, 1990, p. 59.
19. A. R. Williams, A. L. Kellner, X. S. Jiang, and P. K. L. Yu, "InGaAs/InP Waveguide Photodetector with High Saturation Intensity," *Electron. Lett.*, vol. 28, 1992, p. 2258.
20. M. Wu and T. Itoh, "Ultrafast Photonic to Microwave Transformer (PMT)," *LEOS Topical Meeting on Optical Microwave Interactions*, Paper W1.2, 1993.
21. A. Larsson et al., *J. Quantum Electron.*, vol. 24, 1988, p. 787.
22. D. L. Crawford, R. Nagarajan, and J. E. Bowers, "Comparison of Bulk and Quantum Wire Photodetectors," *Appl. Phys. Lett.*, vol. 58, no. 15, April 1991, pp. 1629–1631.
23. D. Kuhl, F. Hieronymi, E. H. Bottcher, and D. Bimberg, "High-Speed Metal-Semiconductor-Metal Photodetectors on InP: Fe," *IEEE Photon. Tech. Lett.*, vol. 2, no. 8, August 1990, p. 574.
24. Y. G. Wey, K. S. Giboney, J. E. Bowers, M. J. W. Rodwell, P. Silvestre, P. Thiagarajan, and G. Y. Robinson, "110 GHz GaInAs/InP *p-i-n* Photodiodes with Integrated Bias Tees and Matched Resistors," *IEEE Photonic Tech. Lett.*, August 1993.
25. D. M. Braun, "Design of Single Layer Antireflection Coatings for InP/InGaAs/InP Photodetectors for the 1200–1600 nm Wavelength Range," *Appl. Opt.*, vol. 27, 1988, pp. 2006–2011.
26. G. Lucovsky, R. F. Schwarz, and R. B. Emmons, "Transit-Time Considerations in *p-i-n* Diodes," *J. Appl. Phys.*, vol. 35, no. 3, March 1961.
27. J. E. Bowers and C. A. Burrus, "High Speed Zero Bias Waveguide Photodetectors," *Electron. Lett.*, vol. 22, 1986, p. 905.
28. A. Alping, R. Tell, and S. T. Eng, "Photodetection Properties of Semiconductor Laser Diode Detectors," *J. Lightwave Tech.*, vol. LT-4, 1986, pp. 1662–1668.
29. A. Alping, "Waveguide *p-i-n* Photodetectors: Theoretical Analysis and Design Criteria," *IEEE Proceedings*, vol. 136, part J, no. 3, June 1989.
30. S. Y. Wang and D. Bloom, "100 GHz Bandwidth Planar GaAs Schottky Photodiode," *Electron. Lett.*, vol. 19, no. 14, 7 July, 1983, p. 554.
31. D. G. Parker and P. G. Say, "Indium Tin Oxide/GaAs Photodiodes for Millimetric-Wave Applications," *Electron. Lett.*, vol. 22, no. 23, 6 Nov. 1986, p. 1266.
32. H. Kamiyama, Y. Kobayashi, T. Nagatsuma, and T. Kamiya, "Very Short Electrical Pulse Generation by a Composite Planar GaAs Photodetectors," *Jpn. J. Appl. Phys.*, vol. 29, Sept. 1990, p. 1717.
33. S. M. Sze, *Physics of Semiconductor Devices*, 1981, pp. 255–263.
34. M. Missous and E. H. Rhoderick, "On the Richardson Constant for Aluminum/Gallium Arsenide Schottky Diodes," *J. Appl. Phys.*, vol. 69, no. 10, 15 May 1991, p. 7142.
35. M. Kamegawa, K. Giboney, J. Karin, S. Allen, M. Case, R. Yu, M. J. W. Rodwell, and J. E. Bowers, "Picosecond GaAs Monolithic Optoelectronic Sampling Circuit," *Photonics Technology Lett.*, vol. 3, no. 6, June 1991, pp. 567–569.
36. E. Ozbay, K. D. Li, and D. M. Bloom, "2.0 psec GaAs Monolithic Photodetector," *IEEE Photon. Tech. Lett.*, vol. 3, no. 6, June 1991, p. 570.
37. N. Emeis, H. Schumacher, and H. Beneking, "High-Speed GaInAs Schottky Photodetector," *Electron. Lett.*, vol. 21, no. 5, 28 Feb. 1985, p. 181.
38. L. Yang, A. S. Sudbo, R. A. Logan, T. Tanbun-Ek, and W. T. Tsang, "High Performance Fe:InP/GaAs Metal/Semiconductor/Metal Photodetectors Grown by Metalorganic Vapor Phase Epitaxy," *IEEE Photon. Tech. Lett.*, vol. 2, no. 1, January 1990, p. 56.
39. T. Mikawa, H. Kuwatsuka, Y. Kito, T. Kumai, M. Makiuchi, S. Yamazaki, O. Wada, and T. Shirai, "Flip-Chip InGaAs Avalanche Photodiode with Ultra Low Capacitance and Large Gain-Bandwidth Product," *Tech. Digest, ThO₂, OFC 1991*, p. 186.
40. H. Kuwatsuka, T. Mikawa, S. Miura, N. Yasuoka, T. Tanahashi, and O. Wada, "An Al_xGa_{1-x}Sb Avalanche Photodiode with Gain Bandwidth Product of 90 GHz," *Photon. Tech. Lett.*, vol. 2, no. 1, Jan 1990, p. 54.

41. F. Capusso, H. M. Cox, A. L. Hutchinson, N. A. Olsson, and S. G. Hummel, "Pseudo-Quaternary GaInAsP Semiconductors: A New $\text{In}_{0.53}\text{Ga}_{0.47}\text{As}/\text{InP}$ Graded Gap Superlattice and Its Applications to Avalanche Photodiodes," *Appl. Phys. Lett.*, vol. 45, no. 11, 1 December 1984, pp. 1193–1195.
42. L. E. Tarof, "Planar InP/InGaAs Avalanche Photodiodes with a Gain-Bandwidth Product Exceeding 100 GHz," *Tech. Digest, ThO3, OFC* 1991, p. 187.
43. K. Taguchi, T. Torikai, Y. Sugimoto, K. Makito, and H. Ishihara, "Planar-Structure InP/InAsP/InGaAs Avalanche Photodiodes with Preferential Lateral Extended Guard Ring for 1.0–1.6 μm Wavelength Optical Communication Use," *J. Lightwave Tech.*, vol. 6, no. 11, Nov. 1988, p. 1643.
44. T. Shiba, E. Ishimura, K. Takahashi, H. Namizaki, and W. Susaki, "New Approach to the Frequency Response Analysis of an InGaAs Avalanche Photodiode," *J. Lightwave Tech.*, vol. 6, no. 10, Oct. 1988, p. 1502.
45. K. Berchtold, O. Krumpolz, and J. Suri, "Avalanche Photodiodes with a Gain-Bandwidth Product of More Than 200 GHz," *Appl. Phys. Lett.*, vol. 26, no. 10, 15 May 1975, p. 585.
46. T. Kagawa, H. Asai, and Y. Kawamura, "An InGaAs/InAlAs Superlattice Avalanche Photodiode with a Gain Bandwidth Product of 90 GHz," *IEEE Photon. Tech. Lett.* vol. 3, no. 9, September 91, pp. 815–817.
47. H. Imai and T. Kaneda, "High-Speed Distributed Feedback Lasers and InGaAs Avalanche Photodiodes," *J. Lightwave Tech.*, vol. 6, no. 11, Nov. 1988, p. 1643.
48. H. C. Hsieh and W. Sargeant, "Avalanche Buildup Time of an InP/InGaAsP/InGaAs APD at High Gain," *J. Quantum Electron.*, vol. 25, no. 9, Sept. 1989, p. 2027.
49. F. Osaka, T. Mikawa, and T. Kaneda, "Impact Ionization of Electrons and Holes in (100)-Oriented Ga $_{1-x}$ In $_x$ As $_y$ P $_{1-y}$," *IEEE J. Quantum Electron.*, vol. QE-21, no. 9, September 1985, pp. 1326–1338.
50. C. A. Lee, R. A. Logan, R. L. Batdorf, J. J. Kleimack, and W. Wiegmann, "Ionization Rates of Holes and Electrons in Silicon," *Phys. Rev.*, vol. 134, 1964, pp. A761–A773.
51. R. B. Emmons, "Avalanche-Photodiode Frequency Response," *J. Appl. Phys.*, vol. 38, no. 9, August 1967, p. 3705.
52. S. L. Miller, "Avalanche Breakdown in Germanium," *Phys. Rev.*, vol. 99, Aug. 1955, pp. 1234–1241.
53. R. J. McIntyre, "Multiplication Noise in Uniform Avalanche Junctions," *IEEE Trans. Electron. Devices*, vol. ED-13, Jan. 1966, pp. 164–168.
54. B. L. Kasper and J. C. Campell, "Multigigabit-per-Second Avalanche Photodiode Lightwave Receivers," *J. Lightwave Tech.*, vol. LT-5, no. 10, October 1987, p. 1351.
55. M. Brain and T. P. Lee, "Optical Receiver for Lightwave Communication Systems," *J. Lightwave Tech.*, vol. LT-3, no. 6, December 1985, p. 1281.
56. T. V. Muoi, "Receiver Design for High-Speed Optical-Fiber Systems," *J. Lightwave Tech.*, vol. LT-2, no. 3, June 1984, p. 243.
57. J. A. Caldmain and G. Mourou, "Subpicosecond Electrooptic Sampling: Principles and Applications," *IEEE J. Quantum Electron.*, vol. QE-22, 1986, pp. 69–78.
58. D. R. Grischkowsky, M. B. Ketchen, C.-C. Chi, I. N. Duling, III, N. J. Halas, J.-M. Halbout, and P. G. May, "Capacitance Free Generation and Detection of Subpicosecond Electrical Pulses on Coplanar Transmission Lines," *IEEE J. Quantum Electron.*, vol. 24, no. 2, February 1988, pp. 221–225.
59. W. C. Nunnally and R. B. Hammond, "Optoelectronic Switch for Pulsed Powers," *Picosecond Optoelectronic Devices*, C. H. Lee (ed.), Academic Press, Orlando, Fla., 1984, pp. 373–398.
60. F. W. Smith, H. Q. Le, V. Diadiuk, M. A. Hollis, A. R. Calawa, S. Gupta, M. Frankel, D. R. Dykaar, G. A. Mourou, and T. Y. Hsiang, "Picosecond GaAs-Based Photoconductive Optoelectronic Detectors," *Appl. Phys. Lett.*, vol. 54, no. 10, 6 March 1989, p. 890.
61. N. G. Paulter, A. J. Gibbs, and D. N. Sinha, "Fabrication of High-Speed GaAs Photoconductive Pulse Generators and Sampling Gates by Ion Implantation," *IEEE Trans. Electron Device*, vol. 35, no. 12, December 1988, pp. 2343–2348.
62. Y. Chen, S. Williamson, and T. Brock, "1.2 ps High Sensitivity Photodetector/Switch Based on Low-Temperature-Grown GaAs," Postdeadline papers, CPDP 10/591, CLEO 1991.

SIGNAL DETECTION AND ANALYSIS

John R. Willison

*Stanford Research Systems, Inc.
Sunnyvale, California*

27.1 GLOSSARY

A	dimensionless material constant for $1/f$ noise
C	capacitance (farads)
I	current (amps)
$I_{\text{shot noise}}$	shot noise current (amps)
k	Boltzmann's constant
q	electron charge (Coulombs)
R	resistance (ohms)
S/N	signal-to-noise ratio
T	temperature (Kelvin)
$V_{\text{Johnson, rms}}$	RMS Johnson noise voltage (V)
Δf	bandwidth (Hz)

27.2 INTRODUCTION

Many optical systems require a quantitative measurement of light. Applications range from the very simple, such as a light meter using a photocell and a d'Arsenval movement, to the complex, such as the measurement of a fluorescence lifetime using time-resolved photon counting.

Often, the signal of interest is obscured by noise. The noise may be fundamental to the process: photons are discrete quanta governed by Poisson statistics which gives rise to shot noise. Or, the noise may be from more mundane sources, such as microphonics, thermal emf's, or inductive pickup.

This chapter describes methods for making useful measurements of weak optical signals, even in the presence of large interfering sources. The chapter will emphasize the electronic aspects of the problem. Important details of optical systems and detectors used in signal recovery are covered in Chap. 24, "Photodetectors," by Paul R. Norton; Chap. 25, "Photodetection," by Abhay M.

Joshi and Gregory H. Olsen; and Chap. 26, “High-Speed Photodetectors,” by John E. Bowers and Yih G. Wey in this volume.

27.3 PROTOTYPE EXPERIMENT

Figure 1 details the elements of a typical measurement situation. We wish to measure light from the source of interest. This light may be obscured by light from background sources. The intensity of the source of interest, and the relative intensity of the interfering background, will determine whether some or all of the techniques shown in Fig. 1 should be used.

Optics

An optical system is designed to pass photons from the source of interest and reject photons from background sources. The optical system may use spatial focusing, wavelength, or polarization selection to preferentially deliver photons from the source of interest to the detector.

There are many trade-offs to consider when designing an optical system. For example, if the source is nearly monochromatic and the background is broadband, then a monochromator may be used to improve the signal-to-background ratio of the light reaching the detector. However, if the source of interest is an extended isotropic emitter, then a monochromator with narrow slits and high f -number will dramatically reduce the number of signal photons from the source which can be passed to the detector. In this case, the noise of the detector and amplifiers which follow the optical system may dominate the overall signal-to-noise ratio (S/N).

Photodetectors

There are many types of nonimaging photodetectors. Key criteria to select a photodetector for a particular application include sensitivity for the wavelength of interest, gain, noise, and speed. Important details of many detector types are given in other chapters in the *Handbook*. Operational details (such as bias circuits) of photomultipliers which are specific to boxcar integration and photon counting are discussed in Sec. 27.5.

Amplifiers

In many applications, the output of the detector must be amplified or converted from a current to a voltage before the signal may be analyzed. Selection criteria for amplifiers include type (voltage or transconductance), gain, bandwidth, and noise.

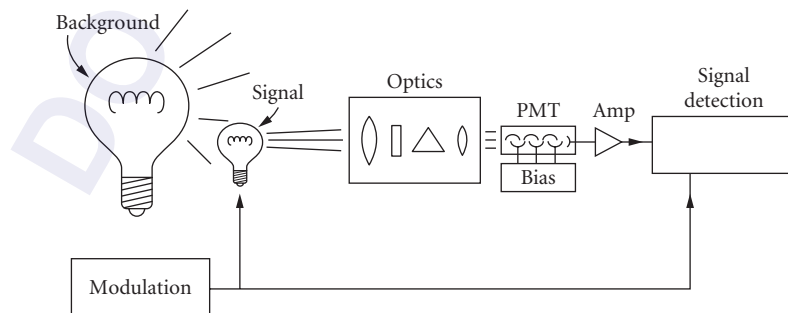


FIGURE 1 Prototypical optical measurement.

Signal Analysis

There are two broad categories of signal analysis, depending on whether or not the source is modulated. Modulating the source allows the signal to be distinguished from the background. Often, source modulation is inherent to the measurement. For example, when a pulsed laser is used to induce a fluorescence, the signal of interest is present only after the laser fires. Other times, the modulation is “arranged,” as when a cw source is chopped. Sometimes the source cannot be modulated or the source is so dominant over the background that modulation is unnecessary.

27.4 NOISE SOURCES

An understanding of noise sources in a measurement is critical to achieving signal-to-noise performance near theoretical limits. The quality of a measurement may be substantially degraded by a trivial error. For example, a poor choice of termination resistance for a photodetector may increase current noise by several orders of magnitude.¹

Shot Noise

Light and electrical charge are quantized, and so the number of photons or electrons which pass a point during a period of time are subject to statistical fluctuations. If the signal mean is M photons, its standard deviation (noise) will be \sqrt{M} , hence the $S/N = M/\sqrt{M} = \sqrt{M}$. The mean M may be increased if the rate is higher or the integration time is longer. Short integration times or small signal levels will yield poor S/N values. Figure 2 shows the S/N , which may be expected as a function of current level and integration time for a shot-noise limited signal.

“Integration time” is a convenient parameter when using time-domain signal recovery techniques. “Bandwidth” is a better choice when using frequency-domain techniques. The rms noise current in the bandwidth Δf Hz due to a “constant” current, I amps, is given by

$$I_{\text{shot noise}} = \sqrt{(2qI\Delta f)} \quad (1)$$

where $q = 1.6 \times 10^{-19}$ C

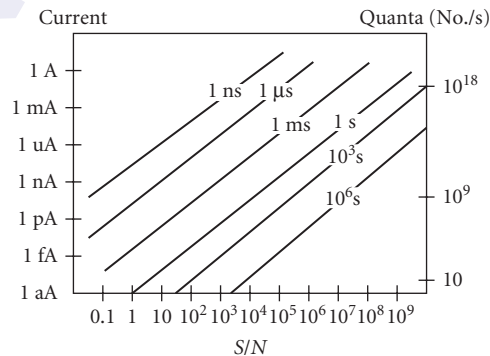


FIGURE 2 Signal-to-noise versus flux and measurement time.

Johnson Noise

The electrons which allow current conduction in a resistor are subject to random motion which increases with temperature. This fluctuation of electron density will generate a noise voltage at the terminals of the resistor. The rms value of this noise voltage for a resistor of R ohms, at a temperature of T Kelvin, in a bandwidth of Δf Hz is given by

$$V_{\text{Johnson,rms}} = \sqrt{(4kTR\Delta f)} \quad (2)$$

where k is Boltzmann's constant. The noise voltage in a 1-Hz bandwidth is given by

$$V_{\text{Johnson,rms}}(\text{per } \sqrt{\text{Hz}}) = 0.13 \text{ nV} \times \sqrt{(R(\text{ohms}))} \quad (3)$$

Since the Johnson noise voltage increases with resistance, large-value series resistors should be avoided in voltage amplifiers. For example, a 1-k Ω resistor has a Johnson voltage of about 4.1 nV/ $\sqrt{\text{Hz}}$. If detected with a 100-MHz bandwidth, the resistor will show a noise of 41- μV rms, which has a peak-to-peak value of about 200 μV .

When a resistor is used to terminate a current source, or as a feedback element in a current-to-voltage converter, it will contribute a noise current equal to the Johnson noise voltage divided by the resistance. Here, the noise current in a 1-Hz bandwidth is given by

$$I_{\text{Johnson,rms}}(\text{per } \sqrt{\text{Hz}}) = 130 \text{ pA} / \sqrt{R(\text{ohms})} \quad (4)$$

As the Johnson noise current increases as R decreases, small-value resistors should be avoided when terminating current sources. Unfortunately, small terminating resistors are required to maintain a wide frequency response. If a 1-k Ω resistor is used to terminate a current source, the resistor will contribute a noise current of about 4.1 pA/ $\sqrt{\text{Hz}}$, which is about 1000 \times worse than the noise current of an ordinary FET input operational amplifier.

1/f Noise

The voltage across a resistor carrying a constant current will fluctuate because the resistance of the material used in the resistor varies. The magnitude of the resistance fluctuation depends on the material used: carbon composition resistors are the worst, metal film resistors are better, and wire wound resistors provide the lowest 1/f noise. The rms value of this noise source for a resistance of R ohms, at a frequency of f Hz, in a bandwidth of Δf Hz is given by

$$V_{1f,\text{rms}} = IR \times \sqrt{(A\Delta f/f)} \quad (5)$$

where the dimensionless constant A has a value of about 10^{-11} for carbon. In a measurement in which the signal is the voltage across the resistor (IR), the $S/N = 3 \times 10^5 \sqrt{(f/\Delta f)}$. Often, this noise is a troublesome source of low-frequency noise in voltage amplifiers.

Nonessential Noise Sources

There are many discrete noise sources which must be avoided in order to make reliable low-level light measurements. Figure 3 shows a simplified noise spectrum on log-log scales. The key features in this noise spectra are frequencies worth avoiding: diurnal drifts (often seen via input offset drifts with temperature), low frequency (1/f) noise, power line frequencies and their harmonics, switching power supply and crt display frequencies, commercial broadcast stations

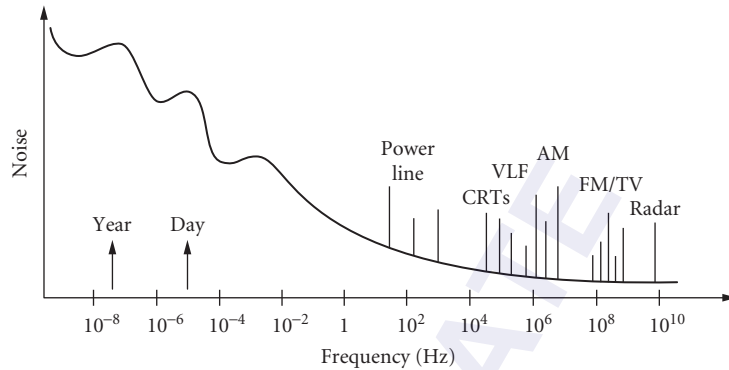


FIGURE 3 Simplified noise spectrum.

(AM, FM, VHF, and UHF TV), special services (cellular telephones, pagers, etc.), microwave ovens and communications, to RADAR and beyond.

Your best alternatives for avoiding these noise sources are

1. Shield to reduce pickup.
2. Use differential inputs to reject common mode noise.
3. Bandwidth limit the amplifier to match expected signal.
4. Choose a quiet frequency for signal modulation when using a frequency-domain detection technique.
5. Trigger synchronously with interfering source when using a time-domain detection technique.

Common ways for extraneous signals to interfere with a measurement are illustrated in Fig. 4a to f.

Noise may be injected via a stray capacitance as in Fig. 4a. The stray capacitance has an impedance of $1/j\omega C$. Substantial currents may be injected into low-impedance systems (such as transconductance inputs), or large voltages may appear at the input to high-impedance systems.

Inductive pickup is illustrated in Fig. 4b. The current circulating in the loop on the left will produce a magnetic field which in turn induces an emf in the loop on the right. Inductive noise pickup may be reduced by reducing the areas of the two loops (by using twisted pairs, for example), by increasing the distance between the two loops, or by shielding. Small skin depths at high frequencies allow nonmagnetic metals to be effective shields; however, high- μ materials must be used to shield from low frequency magnetic fields.

Resistive coupling, or a “ground loop,” is shown in Fig. 4c. Here, the detector senses the output of the experiment plus the IR voltage drop from another circuit which passes current through the same ground plane. Cures for ground-loop pickup include grounding everything to the same point, using a heavier ground plane, providing separate ground return paths for large interfering currents, and using a differential connection between the signal source and amplifier.

Mechanical vibrations can create electrical signals (microphonics) as shown in Fig. 4d. Here, a coaxial cable is charged by a battery through a large resistance. The voltage on the cable is $V = Q/C$. Any deformation of the cable will modulate the cable’s capacitance. If the period of vibration which causes the deformation is short compared to the RC time constant then the stored charge on the cable, Q , will remain constant. In this case, a 1-ppm modulation of the cable capacitance will generate an ac signal with an amplitude of 1 ppm of the dc bias on the cable, which may be larger than the signal of interest.

The case of magnetic microphonics is illustrated in Fig. 4e. Here, a dc magnetic field (the earth’s field or the field from a permanent magnet in a latching relay, for example) induces an emf in the signal path when the magnetic flux through the detection loop is modulated by mechanical motion.

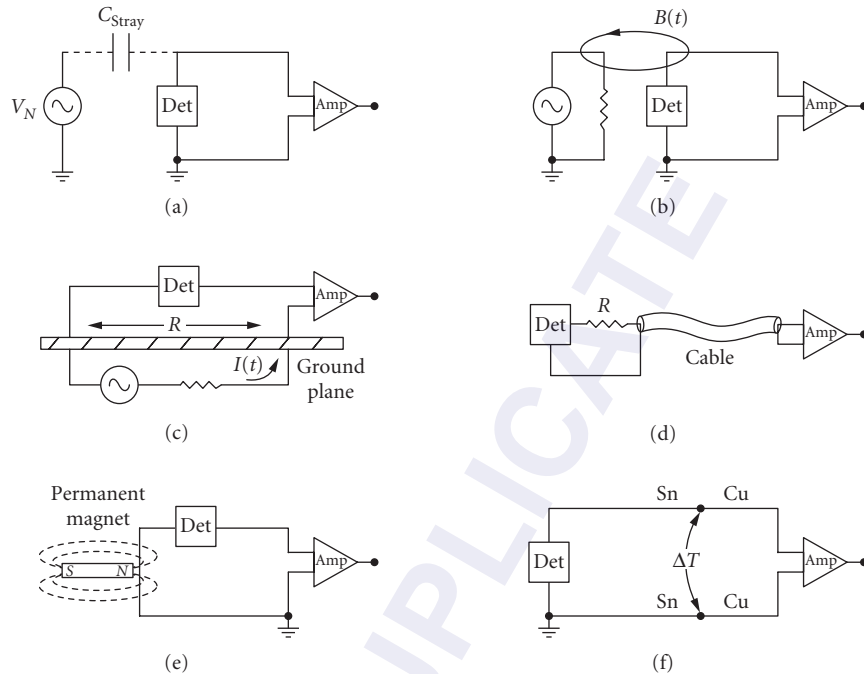


FIGURE 4 Coupling of noise sources.

Unwanted thermocouple junctions are an important source of offset and drift. As shown in Fig. 4f, two thermocouple junctions are formed when a signal is connected to an amplifier. For typical interconnect materials (copper, tin) one sees about $10 \mu\text{V}/^\circ\text{C}$ of offset. These extraneous junctions occur throughout instruments and systems: their impact may be eliminated by making ac measurements.

27.5 APPLICATIONS USING PHOTOMULTIPLIERS

Photomultiplier tubes (PMTs) are used for detection of light from about 200 to 900 nm. Windowless PMTs can be used from the near UV through the x-ray region, and may also be used as particle detectors. Their low noise, high gain, wide bandwidth, and large dynamic range have placed them in many applications. They are the only detectors which may be recommended for low-noise photon counting applications.^{2,3}

In this chapter, we are primarily concerned with the electrical characteristics of PMTs. Understanding these characteristics is important if we are to realize the many desirable features of these devices.

A schematic representation of a PMT, together with a typical bias circuit, is shown in Fig. 5. While the concepts depicted here are common to all PMTs, the particulars of biasing and termination will change between PMT types and applications. PMTs have a photocathode, several dynodes (6 to 14), and an anode. They are usually operated from a negative high voltage, with the cathode at the most negative potential, each successive dynode at a less negative potential, and the anode near ground. An incident photon may eject a single photoelectron from the photocathode which will strike the first dynode with an energy of a few hundred volts. A few (2– to 5) electrons will be ejected from the first dynode by the impact of the photoelectron: these electrons will in turn strike

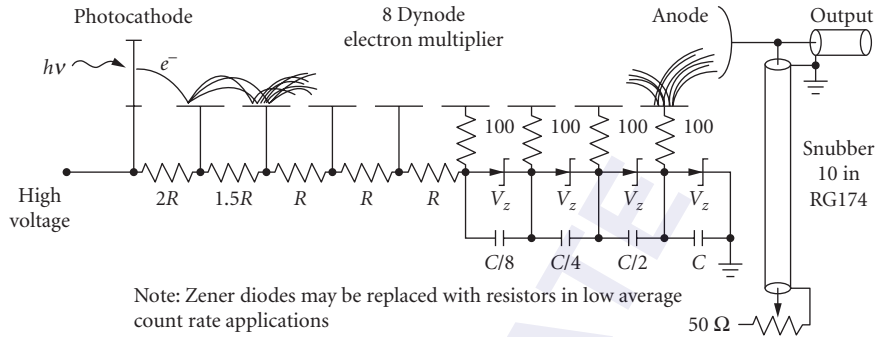


FIGURE 5 PMT base for photon counting or fast integration.

the second dynode, ejecting more electrons. The process continues at each dynode until all of the electrons are collected by the anode.

Quantum Efficiency

The quantum efficiency (QE) of a PMT is a measure of the probability that a photon will eject a photoelectron at the photocathode. The QE depends on the type of material used in the cathode and the wavelength of light. QEs may be as high as 10 to 30 percent at their peak wavelength. The cathode material will also affect the dark count rate from the PMT: a cathode with good red sensitivity may have a high dark count rate.

Gain

A PMT's gain depends on the number of dynodes, the dynode material, and voltage between the dynodes. PMT gains range from 10^3 to 10^7 . The anode output from the PMT will typically go to an electronic amplifier. To avoid having the system noise be dominated by the amplifier's noise, the PMT should be operated with enough gain so that the dark current times the gain is larger than the amplifier's input current noise.

Bandwidth

The frequency response, speed, rise time, and pulse-pair resolution of PMTs depend on the structure of the dynode multiplier chain. The leading edges of the anode output have transition times from 2 to 20 ns. Trailing edges are usually about three times slower. Much faster PMTs, with rise times on order 100 ps, use microchannel plate multipliers.

When using gated integrators to measure PMT outputs, the pulse width of the anode signal should be less than the gate width so that timing information is not lost. For photon counting, the pulse width should be smaller than the pulse-pair resolution of the counter/discriminator to avoid saturation effects. When using lock-in amplifiers, pulse width is usually not important, since the slowest PMTs will have bandwidths well above the modulation frequency.

Pulse Height

In pulsed experiments, the criterion for a detectable signal often depends on the electrical noise environment of the laboratory and the noise of the preamplifier. In laboratories with Q-switched lasers or pulsed discharges, it is difficult to reduce the noise on any coaxial cable below a few millivolts.

A good, wide bandwidth preamplifier will have about $1.5 \text{ nV}/\sqrt{\text{Hz}}$, or about $25\text{-}\mu\text{V}$ rms over a 300-MHz bandwidth. Peak noise will be about 2.5 times the rms noise, so it is important that the PMT provide pulses of greater than 1-mV amplitude.

Use manufacturer's specifications for the current gain and rise time to estimate the pulse amplitude from the PMT:

$$\text{Amplitude (mV)} = 4 \times \text{gain (millions)/rise time (ns)} \quad (6)$$

This formula assumes that the electrons will enter a $50\text{-}\Omega$ load in a square pulse whose duration is twice the rise time. (Since the rise time will be limited by the bandwidth of the preamplifier, use the larger of the amplifier or PMT rise times in this formula.)

If the PMT anode is connected via a $50\text{-}\Omega$ cable to a large load resistance, then the pulse shape may be modeled by the lumped parameters of the cable capacitance (about 100 pF/meter for RG-58) and the termination resistance. All of the charge in the pulse is deposited on the cable capacitance in a few nanoseconds. The voltage on the load will be $V = Q/C$ where C = cable capacitance. This voltage will decay exponentially with a time constant of RC where R is the load resistance in ohms. In this case, the pulse height will be

$$\text{Amplitude (mV)} = 160 \times \text{gain (millions)/cable } C \text{ (pF)} \quad (7)$$

The current gain of a PMT is a strong function of the high voltage applied to the PMT. Very often, PMTs will be operated well above the high voltage recommended by the manufacturer, and thus substantially higher current gains ($10\times$ to $100\times$ above specs). There are usually no detrimental effects to the PMT as long as the anode current is kept well below the rated value.

Dark Counts

PMTs are the quietest detectors available. The primary noise source is thermionic emission of electrons from the photocathode and from the first few dynodes of the electron multiplier. PMT housings which cool the PMT to about -20°C can dramatically reduce the dark counts (from a few kHz to a few Hz). The residual counts arise from radioactive decays of materials inside the PMT and from cosmic rays.

PMTs which are specifically designed for photon counting will specify their noise in terms of the rate of output pulses whose amplitudes exceed some fraction of a pulse from a single photon. More often, the noise is specified as an anode dark current. Assuming the primary source of dark current is thermionic emission from the photocathode, the dark count rate is given by

$$\text{Dark count (kHz)} = 6 \times \text{dark current (nA)/gain (millions)} \quad (8)$$

PMT Base Design

PMT bases which are designed for general-purpose applications are not appropriate for photon counting or fast-gated integrator applications (gates < 10 to 20 ns). General-purpose bases will not allow high count rates, and often cause problems such as double counting and poor plateau characteristics. A PMT base with the proper high-voltage taper, bypassing, snubbing, and shielding is required for good time resolution and best photon counting performance.

Dynode Biasing A PMT base provides bias voltages to the PMTs photocathode and dynodes from a single, negative, high-voltage power supply. The simplest design consists of a resistive voltage divider. In this configuration the voltage between each dynode, and thus the current gain at each dynode, is the same. Typical current gains are three to five, so there will typically be four electrons leaving the first dynode, with a variance of about two electrons. This large relative variance (due to the small

number of ejected electrons) gives rise to large variations in the pulse height of the detected signal. Since statistical fluctuations in pulse height are dominated by the low gain of the first few stages of the multiplier chain, increasing the gain of these stages will reduce pulse-height variations and so improve the pulse-height distribution. This is important for both photon counting and analog detection. To increase the gain of the first few stages, the resistor values in the bias chain are increased to increase the voltage in the front end of the multiplier chain. The resistor values are tapered slowly so that the electrostatic focusing of electrons in the multiplier chain is not adversely affected.⁴

Current for the electron multiplier is provided by the bias network. Current drawn from the bias network will cause the dynode potentials to change, thus changing the tube gain. This problem is of special concern in lifetime measurements. The shape of exponential decay curves will be changed if the tube gain varies with count rate. To be certain that this is not a problem, lifetime measurements should be repeated at reduced intensity. The problem of gain variation with count rate is avoided if the current in the bias network is about 20 times the output current from the PMT's anode.

There are a few other methods to avoid this problem which do not require high bias currents. These methods depend on the fact that the majority of the output current is drawn from the last few dynodes of the multiplier:

1. Replace the last few resistors in the bias chain with Zener diodes. As long as there is some reverse current through a Zener, the voltage across the diodes is nearly constant. This will prevent the voltage on these stages from dropping as the output current is increased.
2. Use external power supplies for the last few dynodes in the multiplier chain. This approach dissipates the least amount of electrical power since the majority of the output current comes from lower-voltage power supplies. However, it is the most difficult to implement.
3. If the average count rate is low, but the peak count rate is high, then bypass capacitors on the last few stages may be used to prevent the dynode voltage from dropping (use $20\times$ the average output current for the chain current). For a voltage drop of less than 1 percent, the stored charge on the last bypass capacitor should be $100\times$ the charge output during the peak count rate. For example, the charge output during a 1-ms burst of a 100-MHz count rate, each with an amplitude of 10 mV into $50\ \Omega$ and a pulse width of 5 ns, is $0.1\ \mu\text{C}$. If the voltage on the last dynode is 200 Vdc, then the bypass capacitor for the last dynode should have a value given by

$$C = 100Q/V = 100 \times 0.15\text{C}/200\text{V} = 0.05\ \mu\text{F} \quad (9)$$

The current from higher dynodes is smaller so the capacitors bypassing these stages may be smaller. Only the final four or five dynodes need to be bypassed, usually with a capacitor which has half the capacitance of the following stage. To reduce the voltage requirement for these capacitors, they are usually connected in series.

Bypassing the dynodes of a PMT may cause high-frequency ringing of the anode output signal. This can cause multiple counts for a single photon or poor time resolution in a gated integrator. The problem is significantly reduced by using small resistors between the dynodes and the bypass capacitors.

Snubbing Snubbing refers to the practice of adding a network to the anode of the PMT to improve the shape of the output pulse for photon counting or fast-gated integrator applications. This "network" is usually a short piece of $50\text{-}\Omega$ coax cable which is terminated into a resistor of less than $50\ \Omega$. The snubber will delay, invert, and sum a small portion of the anode signal to itself. Snubbing should not be used when using a lock-in amplifier since the current conversion gain of a $50\text{-}\Omega$ resistor is very small.

There are four important reasons for using a snubber network:

1. Without some dc resistive path between the anode and ground, anode dark current will charge the signal cable to a few hundred volts (last dynode potential). When the signal cable is connected to an amplifier, the stored charge on the cable may damage the front end of the instrument. PMT bases without a snubber network should include a $100\text{-M}\Omega$ resistor between the anode and ground to protect the instruments.

2. The leading edge of the output current pulse is often much faster than the trailing edge. A snubber network may be used to sharply increase the speed of the trailing edge, greatly improving the pulse pair resolution of the PMT. This is especially important in photon counting applications.
3. Ringing (with a few nanoseconds period) is very common on PMT outputs. A snubber network may be used to cancel these rings which can cause multiple counts from a single photon.
4. The snubber network will help to reverse terminate reflections from the input to the preamplifier.

The round-trip time in the snubber cable may be adjusted so that the reflected signal cancels anode signal ringing. This is done by using a cable length with a round-trip time equal to the period of the anode ringing.

Cathode Shielding Head-on PMTs have a semitransparent photocathode which is operated at negative high voltage. Use care so that no objects near ground potential contact the PMT near the photocathode.

Magnetic Shielding Electron trajectories inside the PMT will be affected by magnetic fields. A field strength of a few gauss can dramatically reduce the gain of a PMT. A magnetic shield made of a high permeability material should be used to shield the PMT.

PMT Base Summary

1. Taper voltage divider for higher gain in the first stages.
2. Bypass last few dynodes in pulsed applications.
3. Use a snubber circuit to shape the output pulse for photon counting or fast-gated integration.
4. Shield the tube from electrostatic and magnetic fields.

27.6 AMPLIFIERS

Several considerations are involved in choosing the correct amplifier for a particular application. Often, these considerations are not independent, and compromises will be necessary. The best choice for an amplifier depends on the electrical characteristics of the detector, and on the desired gain, bandwidth, and noise performance of the system.

Voltage Amplifiers

High Bandwidth Photon counting and fast-gated integration require amplifiers with wide bandwidth. A 350-MHz bandwidth is required to preserve a 1-ns rise time. The input impedance to these amplifiers is usually 50 Ω in order to terminate coaxial cables into their characteristic impedance. When PMTs (which are current sources) are connected to those amplifiers, the 50- Ω input impedance serves as the current-to-voltage converter for the PMT anode signal. Unfortunately, the small termination resistance and wide bandwidth yield a lot of current noise.⁵

High Input Impedance It is important to choose an amplifier with a very high input impedance and low-input bias current when amplifying a signal from a source with a large equivalent resistance. Commercial amplifiers designed for such applications typically have a 100-M Ω input impedance. This large input impedance will minimize attenuation of the input signal and reduce the Johnson noise current drawn through the source resistance, which can be an important noise source. Field effect transistors (FETs) are used in these amplifiers to reduce the input bias current to the amplifiers. Shot noise on the input bias current can be an important noise component, and temperature drift of the input bias current is a source of drift in dc measurements.⁶

The bandwidth of a high-input impedance amplifier is often determined by the RC time constant of the source, cable, and termination resistance. For example, a PMT with 1 meter of RG-58 coax (about 100 pF) terminated into a 1-M Ω resistor will have a bandwidth of about 1600 Hz. A smaller resistance would improve the bandwidth, but increase the Johnson noise current.

Moderate Input Impedance Bipolar transistors offer an input noise voltage which may be several times smaller than the FET inputs of high-input impedance amplifiers, as low as $1 \text{ nV}/\sqrt{\text{Hz}}$. Bipolar transistors have larger input bias currents, hence larger shot noise current, and so should be used only with low-impedance ($<1 \text{ k}\Omega$) sources.

Transformer Inputs When ac signals from very low source impedances are to be measured, transformer coupling offers very quiet inputs. The transformer is used to step up the input voltage by its turns-ratio. The transformer's secondary is connected to the input of a bipolar transistor amplifier.

Low Offset Drift Conventional bipolar and FET input amplifiers exhibit input offset drifts on the order of $5 \mu\text{V}/\text{C}$. In the case where the detector signal is a small dc voltage, such as from a bolometer, this offset drift may be the dominant noise source. A different amplifier configuration, chopper-stabilized amplifiers, essentially measure their input offsets and subtract the measured offset from the signal. A similar approach is used to "autozero" the offset on the input to sensitive voltmeters. Chopper-stabilized amplifiers exhibit very low input offsets with virtually no input offset drift.

Differential The use of "true-differential" or "instrumentation" amplifiers is advised to provide common mode rejection to interfering noise, or to overcome the difference in grounds between the voltage source and the amplifier. This amplifier configuration amplifies the difference between two inputs, unlike a single-ended amplifier, which amplifies the difference between the signal input and the amplifier ground. In high-frequency applications, where good differential amplifiers are not available or are difficult to use, common mode choke may be used to isolate disparate grounds.

Transconductance Amplifiers

When the detector is a current source (or has a large equivalent resistance) then a transconductance amplifier should be considered. These amplifiers (current-to-voltage converters) offer the potential of lower noise and wider bandwidth than a termination resistor and a voltage amplifier; however, some care is required in their application.⁷

A typical transconductance amplifier configuration is shown in Fig. 6. An FET input op amp would be used for its low-input bias current. (Op amps with input bias currents as low as 50 fA are readily available.) The detector is a current source, I_o . Assuming an ideal op amp, the transconductance gain is $A = V_{\text{out}}/I_{\text{in}} = R_f$, and the input impedance of the circuit is R_{in} to the op amp's virtual null. (R_{in} allows negative feedback, which would have been phase shifted and attenuated by the

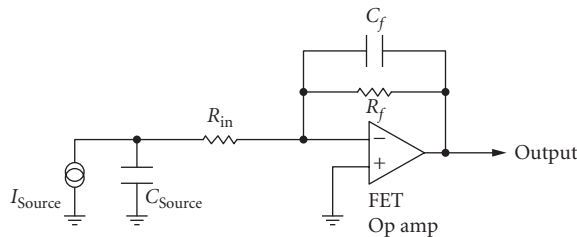


FIGURE 6 Typical transconductance amplifier.

source capacitance at high frequencies, to assure stability.) Commercial transconductance amplifiers use R_s 's as large as $10\text{ M}\Omega$, with R_{in} 's, which are typically $R_f/1000$. A low-input impedance will ensure that current from the source will not accumulate on the input capacitance.

This widely used configuration has several important limitations which will degrade its gain, bandwidth, and noise performance. The overall performance of the circuit depends critically on the source capacitance, including that of the cable connecting the source to the amplifier input. Limitations include:

1. The "virtual null" at the inverting input to the op amp is approximately R_f/A_v , where A_v is the op amp's open loop gain at the frequency of interest. While op amps have very high gain at frequencies below 10 Hz (typically a few million), these devices have gains of only a few hundred at 1 kHz. With an R_f of $1\text{ G}\Omega$, the virtual null has an impedance of $5\text{ M}\Omega$ at 1 kHz, hardly a virtual null. If the impedance of the source capacitance is less than the input impedance, then most of the ac input current will go to charging this capacitance, thereby reducing the gain.
2. The configuration provides high gain for the voltage noise at the noninverting input of the op amp. At high frequencies, where the impedance of the source capacitance is small compared to R_{in} , the voltage gain for noise at the noninverting input is R_f/R_{in} , typically about 1000. As FET input op amps with very low bias currents tend to have high-input-voltage noise, this term can dominate the noise performance of the design.
3. Large R_f 's are desired to reduce the Johnson noise current; however, large R_f 's degrade the bandwidth. If low values of R_f are used, the Johnson noise current can dominate the noise performance of the design.
4. To maintain a flat frequency response, the size of the feedback capacitance must be adjusted to compensate for different source capacitances.

As many undesirable characteristics of the transconductance amplifier can be traced to the source capacitance, a system may benefit from integrating the amplifier into the detector, thereby eliminating interconnect capacitance. This approach is followed in many applications, from microphones to CCD imagers.

27.7 SIGNAL ANALYSIS

Unmodulated Sources

For unmodulated sources, a strip-chart recorder, voltmeter, A/D converter, or oscilloscope may be used to measure the output of the amplifier or detector. In the case of low-light-level measurement, continuous photon counting would be the method of choice.

A variety of problems are avoided by modulating the signal source. When making dc measurements, the signal must compete with large low-frequency noise sources. However, when the source is modulated, the signal may be measured at the modulation frequency, away from these large noise sources.

Modulated Sources

When the source is modulated, one may choose from gated integration, boxcar averaging, transient digitizers, lock-in amplifiers, spectrum analyzers, gated photon counters, or multichannel scalars.

Gated Integration A measurement of the integral of a signal during a period of time can be made with a gated integrator. Commercial devices allow gates from about 100 ps to several milliseconds. A gated integrator is typically used in a pulsed laser measurement. The device can provide shot-by-shot data which is often recorded by a computer via an A/D converter. The gated integrator is

recommended in situations where the signal has a very low duty cycle, low pulse repetition rate, and high instantaneous count rates.⁸

The noise bandwidth of the gated integrator depends on the gate width: short gates will have wide bandwidths, and so will be noisy. This would suggest that longer gates would be preferred; however, the signal of interest may be very short-lived, and using a gate which is much wider than the signal will not improve the S/N .

The gated integrator also behaves as a filter: the output of the gated integrator is proportional to the average of the input signal during the gate, so frequency components of the input signal which have an integral number of cycles during the gate will average to zero. This characteristic may be used to “notch out” specific interfering signals.

It is often desirable to make gated integration measurements synchronously with an interfering source. (This is the case with time-domain signal detection techniques, and not the case with frequency-domain techniques such as lock-in detection.) For example, by locking the pulse repetition rate to the power-line frequency (or to any submultiple of this frequency) the integral of the line interference during the short gate will be the same from shot to shot, which will appear as a fixed offset at the output of the gated integrator.

Boxcar Averaging Shot-by-shot data from a gated integrator may be averaged to improve the S/N . Commercial boxcar averagers provide linear or exponential averaging. The averaged output from the boxcar may be recorded by a computer or used to drive a strip-chart recorder. Figure 7 shows a gated integrator with an exponential averaging circuit.

Lock-In Amplifiers Phase-sensitive synchronous detection is a powerful technique for the recovery of small signals which may be obscured by interference that is much larger than the signal of interest. In a typical application, a cw laser which induces the signal of interest will be modulated by an optical chopper. The lock-in amplifier is used to measure the amplitude and phase of the signal of interest relative to a reference output from the chopper.⁹

Figure 8 shows a simplified block diagram for a lock-in amplifier. The input signal is ac-coupled to an amplifier whose output is mixed (multiplied by) the output of a phase-locked loop which is locked to the reference input. The operation of the mixer may be understood through the trigonometric identity

$$\cos(\omega_1 t + \Phi) * \cos(\omega_2 t) = \frac{1}{2} \{ \cos[(\omega_1 + \omega_2)t + \Phi] + \cos[(\omega_1 - \omega_2)t + \Phi] \} \quad (10)$$

When $\omega_1 = \omega_2$ there is a dc component of the mixer output, $\cos \Phi$. The output of the mixer is passed through a low-pass filter to remove the sum frequency component. The time constant of the filter is selected to reduce the equivalent noise bandwidth: selecting longer time constants will improve the S/N at the expense of longer response times.

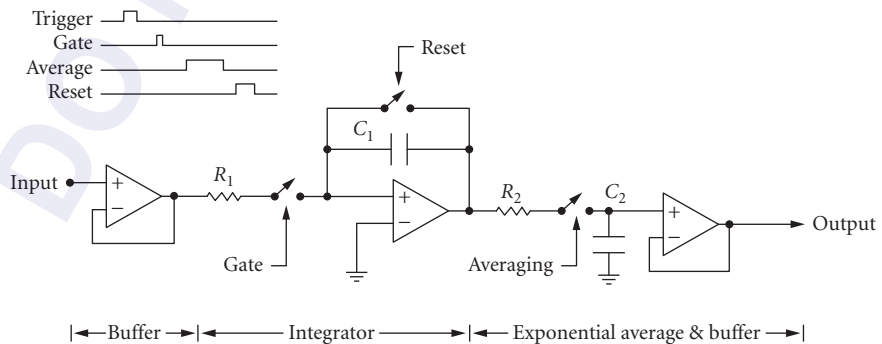


FIGURE 7 Gated integrator and exponential averager.

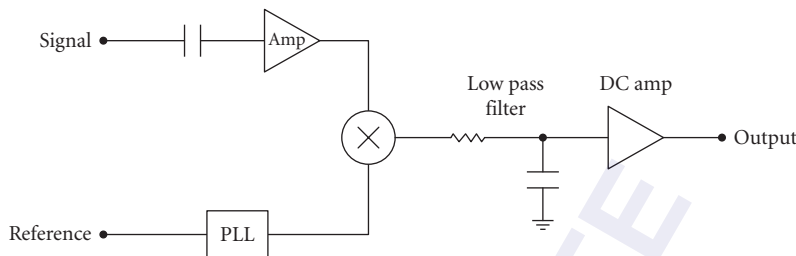


FIGURE 8 Lock-in amplifier block diagram.

The simplified block diagram shown in Fig. 8 is for a “single-phase” lock-in amplifier, which measures the component of the signal at one set phase with respect to the reference. A dual-phase lock-in has another channel which measures the component of the signal at 90° relative to the first channel, which allows simultaneous measurement of the amplitude and phase of the signal.

Digital signal processing (DSP) techniques are rapidly replacing the older analog techniques for the synchronous detection of the signal. In these instruments, the input signal is digitized by a fast, high-resolution A/D converter, and the signal’s amplitude and phase are determined by high-speed computations in a digital signal processor. To maintain the 100-kHz bandwidth of the analog designs, the DSP designs must complete a quarter million 16-bit A/D conversions and 20 million multiply-and-accumulate operations each second. Many artifacts of the analog designs are eliminated by the DSP approach; for example, the output drift and dynamic range of the instruments are dramatically improved.¹⁰

Photon Counting Photon counting techniques offer several advantages in the measurement of light: very high sensitivity (count rates as low as 1 per minute can be a usable signal level), large dynamic range (signal levels as high as 100 MHz can be counted, allowing a 195-dB dynamic range), discrimination against low-level noise (analog noise below the discriminator thresholds will not be counted), and ability to operate over widely varying duty cycles.¹¹

Key elements of a photon counting system include a high-gain PMT operated with sufficiently high voltage so that a single photoelectron will generate an anode pulse of several millivolts into a $50\text{-}\Omega$ load, a fast discriminator to generate logic pulses from anode signals which exceed a set threshold, and fast-gated counters to integrate the counts.

Transient Photon Counting In situations where the time evolution of a light signal must be measured (LIDAR, lifetime measurements, chemical kinetics, etc.) transient photon counters allow the entire signal to be recorded for each event. In these instruments, the discriminated photon pulses are summed into different bins depending on their timing with respect to a trigger pulse. Commercial instruments offer 5-ns resolution with zero dead-time between bins. The time records from many events may be summed together in order to improve the S/N .¹²

Choosing the “Best” Technique Which instrument is best suited for detecting signals from a photomultiplier tube? The answer is based on many factors, including the signal intensity, the signal’s time and frequency distribution, the various noise sources and their time-dependence and frequency distribution.

In general, the choice between boxcar averaging (gated integration) and lock-in detection (phase-sensitive detection) is based on the time behavior of the signal. If the signal is fixed in frequency and has a 50 percent duty cycle, lock-in detection is best suited. This type of experiment commonly uses an optical chopper to modulate the signal at some low frequency. Signal photons occur at random times during the “open” phase of the chopper. The lock-in detects the average difference between the signal during the “open” phase and the background during the “closed” phase.

To use a boxcar averager in the same experiment would require the use of very long, 50 percent duty cycle gates since the photons can arrive anywhere during the “open” phase. Since the gated integrator is collecting noise during this entire gate, the signal is easily swamped by the noise. To correct for this, baseline subtraction can be used where an equal gate is used to measure the background during the “closed” phase of the chopper and subtracted from the “open” signal. This is then identical to lock-in detection. However, lock-in amplifiers are much better suited to this, especially at low frequencies (long gates) and low signal intensities.

If the signal is confined to a very short amount of time, then gated integration is usually the best choice for signal recovery. A typical experiment might be a pulsed laser excitation where the signal lasts for only a short time (100 ps to 1 μ s) at a repetition rate of up to 10 kHz. The duty cycle of the signal is much less than 50 percent. By using a narrow gate to detect signal only when it is present, noise which occurs at all other times is rejected. If a longer gate is used, no more signal is measured but the detected noise will increase. Thus, a 50 percent duty cycle gate would not recover the signal well and lock-in detection is not suitable.

Photon counting can be used in either the lock-in or the gated mode. Using a photon counter is usually required at very low signal intensities or when the use of a pulse height discriminator to reject noise results in an improved S/N . If the evolution of a weak light signal is to be measured, a transient photon counter or multichannel scaler can greatly reduce the time required to make a measurement.

27.8 REFERENCES

1. P. Horowitz and W. Hill, *The Art of Electronics*, Cambridge, New York, 1989, p. 428–447.
2. Photomultiplier Tubes, Hamamatsu Company catalog, 1988.
3. Photomultipliers, Thorn EMI Company catalog, 1990.
4. G. A. Morton and H. M. Smith, “Pulse Height Resolution of High Gain First Dynode Photomultipliers,” *Appl. Phys. Lett.* vol. 13, 1968, p. 356.
5. Model SR445 Fast Preamplifier, *Operation and Service Manual*, Stanford Research Systems, 1990.
6. Model SR560 Low Noise Preamplifier, *Operation and Service Manual*, Stanford Research Systems, 1990.
7. Model SR570 Low Noise Current Amplifier, *Operation and Service Manual*, Stanford Research Systems, 1992.
8. Fast Gated Integrators and Boxcar Averagers, *Operation and Service Manual*, Stanford Research Systems, 1990.
9. Model SR510 Lock-in Amplifier, *Operation and Service Manual*, Stanford Research Systems, 1987.
10. Model SR850 DSP Lock-in Amplifier, *Operation and Service Manual*, Stanford Research Systems, 1992.
11. Model SR400 Gated Photon Counter, *Operation and Service Manual*, Stanford Research Systems, 1988.
12. Model SR430 Multichannel Scaler/Averager, *Operation and Service Manual*, Stanford Research Systems, 1989.

This page intentionally left blank.

DO NOT DUPLICATE

THERMAL DETECTORS

William L. Wolfe

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

Paul W. Kruse

*Consultant
Edina, Minnesota*

28.1 GLOSSARY

DTGS	deuterated triglycine sulfate
p	pyroelectric coefficient
R_e	electrical resistance
R_{th}	thermal resistance
\mathfrak{R}	responsivity
S	Seebeck coefficient
TGS	triglycine sulfate
Z	figure of merit
τ_e	electrical time constant
τ_{th}	thermal time constant

28.2 THERMAL DETECTOR ELEMENTS¹

Introduction

Thermal detectors (transducers) of optical radiation are generally considered to be those devices that absorb the radiation, increase their own temperature, and provide a resultant electrical signal. There are several types, divided according to the physical mechanism that converts the temperature change to a resultant electrical one. The oldest are bolometers and thermocouples. The bolometer changes its electrical resistance as a result of the temperature increase; the thermocouple changes its contact potential difference. There are several different types of bolometers, including thermistor, semiconducting, superconducting, carbon, and metallic. They may also be subdivided according to whether they operate at room or cryogenic temperature. Thermocouples vary according to the materials that are joined, and are sometimes connected in series to generate thermopiles.

Pyroelectric detectors make use of the property of a change in the internal polarization as a function of the change in temperature, the pyroelectric effect. Golay cells and certain variations make use of the expansion of a gas with temperature. All of these detectors are governed by the fundamental equation of heat absorption in the material. Many reviews and two books of collected reprints² provide additional information.

Thermal Circuit Theory

In the absence of joulean heating of the detector element, the spectrum of the temperature difference $d\Delta T$ is given in terms of the spectrum of the absorbed power \tilde{P} (the power is P)

$$d\Delta\tilde{T} = \frac{\epsilon\tilde{P}}{G(1+i\omega\tau)} \quad (1)$$

where G is the thermal conductance, given by the product of the thermal conductivity times the cross-sectional area of the path to the heat sink and divided by the length of the path to that heat sink. The time constant τ is the product of the thermal resistance and the heat capacitance. The thermal resistance is the reciprocal of the thermal conductance, while the thermal capacitance is the thermal capacity times the mass of the detector. In the absence of joulean heating, this is a simple, single time constant thermal circuit, for which the change in temperature is given by

$$d\tilde{T} = \frac{\epsilon\tilde{P}}{G(1+i\omega\tau)} \quad (2)$$

The absorbed power is equal to the incident power times the absorptance α of the material:

$$\tilde{P} = \epsilon\tilde{P}_i \quad (3)$$

The absorptance α is usually written as ϵ (which is legitimate according to Kirchhoff's law) since α is also used for the relative temperature coefficient of resistance (some writers use η).

$$\alpha = \frac{1}{R} \frac{dR}{dT} \quad (4)$$

As radiation is absorbed, part of the heat is conducted to the sink. Some of it gives rise to an increase in temperature. Some is reradiated, but this is usually quite small and is ignored here. The dc responsivity of a thermal detector is proportional to the emissivity and to the thermal resistance. The greater proportion of radiation that is absorbed, the greater will be the responsivity. The less heat that is conducted to the sink, the greater will be the temperature rise. The time constant is a true thermal time constant, the product of thermal resistance and capacitance. The greater the heat capacitance, the more heat necessary for a given temperature increase, and the less heat conducted to the sink, the more available for temperature increase. A high absorptance is accomplished by the use of a black coating, and a sufficient amount of it. Thus, there is a direct conflict between high speed and high responsivity.

The Ideal Thermal Detector³⁻⁶

The ideal thermal detector has a noise that is associated only with the thermal fluctuations of the heat loss to the heat sink, and this coupling is purely radiative. Then the noise equivalent power (NEP) is given by

$$\text{NEP} = \sqrt{16A\sigma kT^5/\epsilon} \quad (5)$$

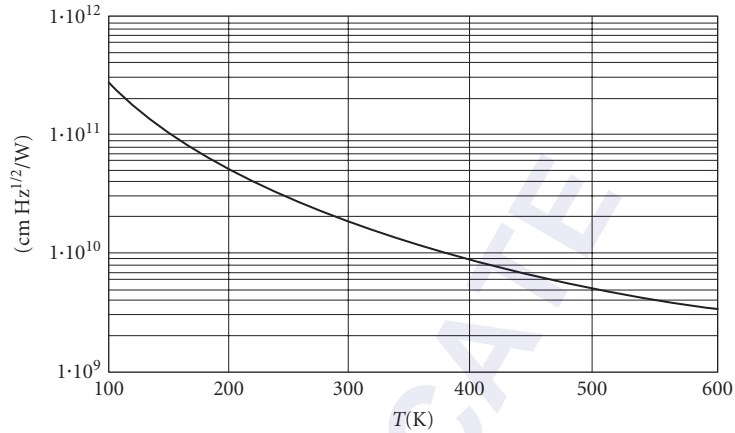


FIGURE 1 Theoretical specific detectivity for ideal thermal detectors.

where it is assumed that the detector is irradiated by a hemisphere of blackbody radiation at the same temperature T as the detector. The corresponding specific detectivity, assuming that the signal varies as the area and the noise as its square root, is

$$D^* = \frac{\epsilon^{1/2}}{4\sqrt{\sigma k T^5}} \quad (6)$$

where the detector and background are at the same temperature.

For circumstances in which the detector is in a cooled chamber, the total radiation from the sources at various temperatures must be calculated. Figure 1 shows the specific detectivity of a background limited ideal thermal detector as a function of the temperature of the surround.

No detector is ideal, and every one will be limited by the signal loss due to incomplete absorption at the surface and any transmission losses by the optical system that puts the radiation on the detector. The detector will also have noise that arises from its conductive coupling to the heat sink, and probably Johnson noise as well. The conductive mean square power fluctuation is given by

$$\langle P^2 \rangle = 4kT^2G \quad (7)$$

The Johnson noise power density is $4kT$. Therefore, the total mean square power fluctuation is given by

$$\langle P^2 \rangle = 4kT[GT + 4\epsilon A \sigma T^4 + 1] \quad (8)$$

Bolometers

Most single-element bolometers are connected in a voltage divider network, as shown in Fig. 2. A stable voltage supply is used to develop a current and consequent voltage drop across the two resistors. One is the detector, while the other should be a matching element to eliminate signals arising from a change in the ambient temperature. It should match the detector in both resistance and in the

temperature coefficient of resistance. Usually another, but blinded, detector is used. The expression for power conservation is

$$\begin{aligned}
 C \frac{d\Delta T}{dt} + G\Delta T &= \frac{d(i^2 R)}{dt} \Delta T + P \\
 C \frac{d\Delta T}{dt} + G\Delta T &= \frac{V^2(R_1 - R)}{(R_1 + R)^3} \frac{dR}{dT} \Delta T + P \\
 C \frac{d\Delta T}{dt} + \left[G - \frac{V^2 R \alpha}{(R_1 + R)^2} \frac{(R_1 - R)}{(R_1 + R)} \right] \Delta T &= P
 \end{aligned} \tag{9}$$

The solution to this is a transient that has an RC time constant, where R is the reciprocal of the bracketed term, and C is the thermal capacitance, and the same steady-state term given above. The transient decays as long as G is greater than the rest of the bracket, but the detector burns up if not. This is still another reason for matching the resistances. The dc responsivity is a function of the construction parameters, including the path to the sink, the bias voltage, and the relative change of resistance with temperature. The different bolometers are divided according to how their resistances change with temperature. (R_1 and R represent slightly different values of R_D .)

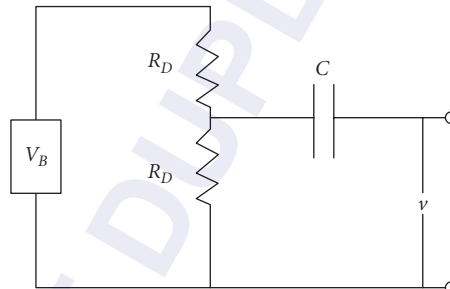


FIGURE 2 Balanced voltage divider circuit for a thermal detector.

Metal Bolometers These have a linear change in resistance with temperature that may be expressed as

$$R = R_0 [1 + \gamma(T - T_0)] \tag{10}$$

Therefore the thermal coefficient is

$$\alpha = \frac{\gamma}{1 + \gamma(T - T_0)} \tag{11}$$

This coefficient always decreases with temperature, and burnout does not occur. The coefficient is approximately equal to the inverse of the temperature, and is therefore never very high.

Semiconductor Bolometers These have an exponential change of resistance with temperature, given by

$$R = R_0 e^{\beta/T} \tag{12}$$

so that

$$\alpha = -\beta/T^2 \tag{13}$$

The value of β depends upon the particular material. These detectors can burn out. Two basic types exist: (1) those that are used at low temperatures and (2) those that are used at about room temperature.

The most used low-temperature bolometer⁷ is germanium in a bath of liquid helium. Pure germanium is transparent in the infrared, but with enough compensated doping it becomes a good conductor with a high-temperature coefficient of resistance.⁸ Typical concentrations are about 10^{16}cm^{-3} of gallium and 10^{15} of indium. Even these are not sufficient at wavelengths shorter than $10\ \mu\text{m}$ since the free-carrier absorption is proportional to wavelength. In such a case a black coating is sometimes used. Improvements have been made since Low's first work.⁹⁻¹¹

Superconducting Bolometers These make use of the extremely large thermal coefficient of resistance at the transition temperature.¹²⁻¹⁴ Originally they needed to be controlled very carefully, or a small change in ambient conditions (on the order of $0.01\ \text{K}$) could cause an apparent signal of appreciable magnitude. A more recent version¹⁵ incorporates an evaporated thin film on an anodized aluminum block that is coupled to a helium bath by a brass rod. The detector has a time constant of about $3\ \mu\text{s}$ due to this high thermal conductance and a good NEP of about $10^{-13}\ \text{WHz}^{-1/2}$. It still must be controlled to about $10^{-5}\ \text{K}$, and this is accomplished with a heater current and control circuit.

Recently developed materials not only have high-temperature transition points but also have more gradual transitions, and provide a better compromise between good responsivity and the requirement for exquisite control.¹⁶

Carbon Bolometers These are a form of semiconductor bolometers that have been largely superseded by germanium bolometers. They are made of small slabs of carbon resistor material, connected to a metal heat sink by way of a thin mylar film. Although their responsivities are comparable to germanium bolometers, their noise is several orders of magnitude higher.¹⁷

Thermocouples and Thermopiles^{18,19}

A *thermocouple* is made by simply joining two dissimilar conductors. A good pair has a large relative Seebeck coefficient and gives rise to a potential difference. The materials also have large electrical conductivities and small thermal ones, so there is little voltage drop across the length and a small thermal gradient. Although there are many different couples (many are not even used for radiation detection), those most often used for this application are bismuth telluride, copper, and constantan. The expression for the responsivity is given in terms of the relative Seebeck coefficient S_{12} (the difference in the voltage change with temperature between the two materials) and the expression derived above for the thermal circuit

$$\mathfrak{R} = \frac{S\epsilon}{G(1+i\omega\tau)} \quad (14)$$

Good materials are those that have a large Seebeck coefficient, a high electrical conductivity, and a small thermal conductivity, and the figure of merit is often defined as

$$Z_{12} = \frac{S_{12}^2}{\left[\sqrt{G_1/\sigma_1} + \sqrt{G_2/\sigma_2}\right]^2} \quad (15)$$

Thermopiles are arrays of thermocouples connected in series. They are manufactured in two ways. Some are carefully wound wires with junctions aligned in the desired pattern, while others are evaporated with the pattern determined by masking operations. Most of the "bulk" thermopiles are wrapped on appropriate mandrels to obtain rigidity. Both kinds are obtainable in a variety of sizes

and patterns that correspond to such things as spectrometer slits, centering annuli, and staggered arrays for moderate-sensitivity thermal imaging.

The Golay Cell²⁰

This detector is used mostly for laboratory operations, as it is slow and fragile, although it has high sensitivity. It is a gas-filled chamber that has a thin membrane at one end and a blackened detector area at the other. Light on the blackened surface causes the increase in temperature; this is transferred to the gas which therefore expands. The membrane bulges, and the amount of the bulge is sensed by some sort of optical lever²¹ or even change in capacity of an electrical element.²² Other versions do not use a blackened surface, but allow the radiation to interact with the gas directly, in which case they are spectral detectors that are “tuned” to the absorption spectrum of the gas.²³

Pyroelectric Detectors²⁴

Some crystals which do not have a center of symmetry experience an electric field along a crystal axis. This internal electric field results from an alignment of electric dipoles (known as polarization), and is related to the crystal temperature. In these ferroelectric crystals, this results in a charge being generated and stored on plates connected to the crystal. Polarization disappears above the so-called Curie temperature that is characteristic of each material. Thus, below the Curie temperature, a change in temperature results in a current. The equation for the response of a pyroelectric detector is

$$\mathfrak{R} = \frac{\omega p A_d \epsilon R_e R_{th}}{(1 + i\omega\tau_{th})(1 + i\omega\tau_e)} \quad (16)$$

where ω is the radian frequency, p is the pyroelectric coefficient, A_d is the detector area, R_e is the electrical resistance, R_{th} is the thermal resistance, τ_{th} is the thermal time constant, and τ_e is the electrical time constant. The relation is shown in Fig. 3, where the responsivity is plotted as a function of frequency. In the low-frequency region the responsivity rises directly as the frequency. This is a result of the ac operation of a pyroelectric. At the (radian) frequency that is the reciprocal of the slower (usually the thermal) time constant, the response levels off. This is the product of the ac rise and the thermal rolloff. Then, when the frequency corresponding to the shorter time constant is reached, the response rolls off.

Type II pyroelectric detectors work on a slightly different mechanism, which is still not fully understood. The electrodes are on the sensing surface of the detector and parallel to the polar axis. In these crystals, the temperature change is not uniform at the onset of radiation and the primary

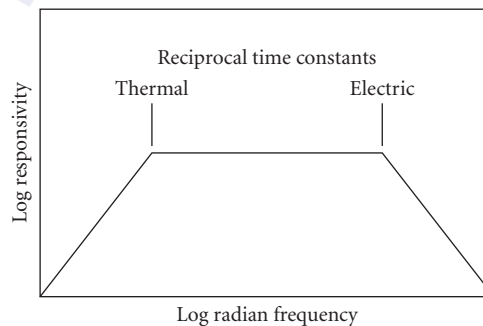


FIGURE 3 Responsivity asymptotes versus frequency.

TABLE 1 General Properties of Thermal Detectors

Type	Operating Temperature (K)	$D^* \times 10^8$ (cmHz ^{1/2} W ⁻¹)	NEP $\times 10^{-10}$ (WHz ^{-1/2})	Time Constant (m)	Size (mm ²)
Silicon bolometer	1.6		3×10^{-5}	8	0.25–0.70
Metal bolometer	2–4	1		10	
Thermistor bolometer	300	1–6		1–8	0.01–10
Germanium bolometer	2–4		0.005	0.4	1.5
Carbon bolometer	2–4		0.03	10	20
Superconducting bolometer (NbN)	15		0.2	0.5	5×0.25
Thermocouples	300		2–10	10–40	0.1 \times 1 to 0.3 \times 3
Thermopiles	300			3.3–10	1–100
Pyroelectrics	300	2–5		10–100 [†]	2 \times 2
Golay cell	300	10	0.6	10–30	10

[†]Shorter values can be obtained at the expense of NEP (for laser detection).

and secondary pyroelectric effects take place, thereby generating a body electric charge distribution.²⁵ Materials most often used for these detectors are TGS (triglycine sulfate), DTGS (deuterated TGS), Li₂SO₄, LiNbO₃, LiTaO₃, and PLZT (lead lanthanum zirconate titanate). TGS is the most used for specialized sensor systems, but has a relatively low Curie point. For higher-temperature operation, usually LiTaO₃ or PLZT is used in the general laboratory environment.

The two advantages of the pyroelectric detector over the other thermal detectors, bolometers, and thermopiles, are its responsivity and its capability of rapid response. The response time and responsivity are traded by choice of the load resistor in the circuit. For instance, with a 100-M Ω load the time constant can be 1 ms and the responsivity 100 V/W, but with a 1-M Ω load the values would be 10 μ s and 1 V/W.

Summary of Elemental Thermal Detector Properties

Although the user should contact suppliers for detailed information, this section provides overall property information about thermal detectors. There are several cautions about summary data. Most detectors can be tailored to have somewhat different properties. Improvements have often been made since the publication of these results. Not all parameters are available in all combinations. Table 1 does, however, give the general flavor of the performance of different thermal detectors.

28.3 ARRAYS

Introduction

As pointed out earlier, thermal detector response is governed by the thermal response time, which is the ratio of the pixel heat capacity C to the thermal conductance G of the heat leakage mechanism. High pixel responsivity is associated with high thermal isolation, i.e., low thermal conductance. Thermal detector design is driven by the thermal isolation structure. It is the structure which determines the extent to which the pixel performance can approach the temperature fluctuation noise limit and, ultimately, the background fluctuation noise limit. Given the value of G associated

with the heat loss mechanism, the pixel heat capacity must be designed appropriately to attain the required thermal response time. Response times in the millisecond range are compatible with high thermal isolation; response times in the microsecond range are not. Thus, two-dimensional arrays of thermal detectors which operate at TV frame rates (30 Hz in the United States) are under development for applications in thermal imagers.

Noise Equivalent Temperature Difference

Whereas elemental detectors are usually described by such figures of merit as NEP and D^* , arrays have been described by a noise equivalent temperature difference (NETD) associated with their use in a camera under certain specific conditions. It is defined as the change in temperature of a blackbody which fills the field of view of a pixel of an infrared imaging system that gives rise to a change of unity in the signal-to-noise ratio at the output of the system. The measurement of the NETD should, however, be with the flooding of several pixels to avoid fringing effects and with an SNR (signal-to-noise ratio) well above 1 to obtain good accuracy. The pixel is defined as the subtense of a single element of the array. The NETD can be written in several different forms. Perhaps the simplest is

$$\text{NETD} = \frac{\sqrt{A_d B}}{D^* (dP_d / dT)} \quad (17)$$

where D^* is the specific detectivity, A_d is the area of a single pixel, B is the system bandwidth and (dP_d / dT) represents the change in power on the detector element per unit change in temperature in the spectral band under consideration. This form does not include the system noise, which is often included by the manufacturers in their calculations. In Eq. (17) the change in power with respect to temperature is

$$\frac{dP_d}{dT} = \frac{A_d \tau_o}{4FN^2} \int_{\lambda_1}^{\lambda_2} \frac{dM}{dT} d\lambda \quad (18)$$

where τ_o is the optics transmission, FN is the focal ratio (defined as the effective focal length divided by the entrance pupil diameter), and M is the radiant emittance of the source. This is almost the definition of the specific detectivity. The NETD can also be written in terms of the responsivity \mathfrak{R} , since the detectivity and responsivity are related in the following way:

$$D^* = \frac{\sqrt{A_d B}}{P} \frac{V_s}{V_N} = \sqrt{A_d B} \frac{\mathfrak{R}}{V_N} \quad (19)$$

where V_s is the signal voltage at the sensor and V_N is the rms noise voltage of a pixel in the bandwidth B . Therefore

$$\text{NETD} = \frac{V_N}{\mathfrak{R} (\partial P_d / \partial T)_{\lambda_1 - \lambda_2}} \quad (20)$$

The power on the detector is related to the power on the aperture by the optical transmission τ_o . The expression can also be formulated in terms of the source radiance, L

$$\text{NETD} = \frac{4FN^2 \sqrt{B}}{D^* \tau_a \tau_o \pi D \Delta \theta (\partial L / \partial T)_{\lambda_2 - \lambda_1}} \quad (21)$$

where D is the diameter of the aperture, $\Delta\theta$ is the angular subtense of a pixel, L is the source radiance, and τ_a is the atmospheric transmission. One last form can be generated by recognizing that, for an isotropic radiator, the radiance is the radiant emittance divided by π :

$$\text{NETD} = \frac{4FN^2V_N}{A_D\tau_a\tau_o\mathfrak{R}(T_s)(\partial M/\partial T)} \quad (22)$$

In this form of the expression for NETD, it is not necessary that the noise be white, nor is it necessary that the noise not include system noise. Whether or not system noise is included should be clearly stated.

Theoretical Limits

Figure 4 illustrates the theoretical limits of thermal arrays having the parameters shown and operating at 300 and 85 K as a function of thermal conductance. The performance of real thermal arrays with those parameters lies on or above the sloping line. As the conductance G is reduced (better thermal isolation), the noise equivalent temperature difference NETD is reduced (improves) until the background limit is reached, when radiant power exchange between the array and the background becomes the dominant heat transfer mechanism. Reducing the detector temperature to 85 K appropriate to a bolometer operating at the transition edge of the high-temperature superconductor $\text{YBa}_2\text{Cu}_3\text{O}_{7-x}$ (YBCO) reduces the NETD by $\sqrt{2}$ and allows the limit to be reached with less thermal isolation (higher G value).

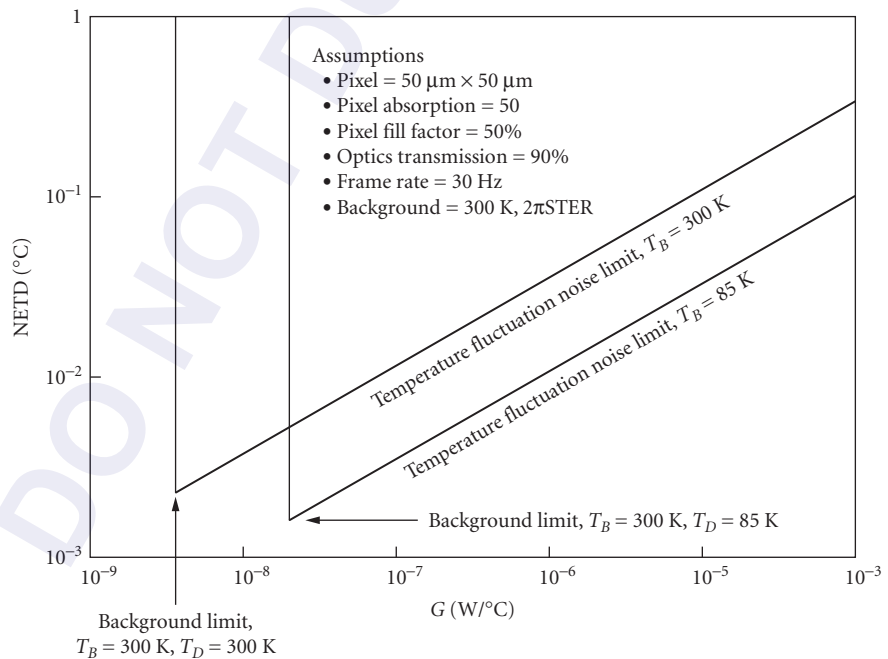


FIGURE 4 Temperature fluctuation noise limit and background fluctuation noise limit of uncooled and cryogenic thermal detector arrays.

Arrays fall into two categories: monolithic and hybrid. Monolithic arrays are prepared on a single substrate, e.g., silicon, upon which the detecting material is deposited in the form of a thin film which is subsequently processed into an array. Hybrid arrays are prepared in two parts: (1) the read-out electronics arrays, usually in silicon, and (2) the detecting material array, usually in wafer form which is thinned by lapping, etching, and polishing. These two arrays are mated by a technique such as flip-chip bonding. Here the interconnection at each pixel must have a sufficiently high electrical conductivity, yet a sufficiently low thermal conductivity—a difficult requirement. If array cost considerations are important, then the monolithic approach, especially in silicon, is the more desirable.

Resistive Bolometer Arrays

The development of resistive bolometric arrays has proceeded along two paths: uncooled arrays and cryogenic arrays. Large, uncooled bolometric arrays have been developed at Honeywell by a team lead by R. A. Wood.^{26,27} Silicon microstructure technology is employed to produce the arrays, a process resembling the fabrication of integrated circuits. Twelve arrays are prepared on a 4-in-diameter silicon wafer. Each monolithic array consists of 240×336 pixels; each pixel is $50 \times 50 \mu\text{m}$. The detecting material is a thin film of vanadium oxide. A Si_3N_4 membrane having a thermal conductance of $1 \times 10^{-7} \text{ WC}^{-1}$ supports the vanadium oxide at each pixel, as shown in Fig. 5. Bipolar transistors implanted in the silicon substrate act as pixel switches for the matrix-addressed array. The response is optimized for the 8- to $14\text{-}\mu\text{m}$ spectral interval. The thermal response time is adjusted for a 30-Hz frame rate. Each pixel is addressed once per frame by a $5\text{-}\mu\text{s}$ pulse. A thermoelectric stabilizer maintains the array at ambient temperature. Other than a one-shot shutter, the camera has no moving parts.

The measured NETD of the camera with F/1 optics at 300 K is 0.04 K. Given the G value of $1 \times 10^{-7} \text{ WK}^{-1}$ it can be seen from Fig. 4 that the array is within a factor of 4 of the temperature fluctuation noise limit. Furthermore, the pixel thermal isolation is so complete that there is no measurable thermal spreading among the pixels.

Linear resistive bolometric arrays of the high-temperature superconductor YBCO on silicon microstructures have been prepared by Johnson et al.,²⁸ also of Honeywell. A two-dimensional array is under development.²⁹ The monolithic arrays operate at the transition edge from 70 to 90 K. As was true for the uncooled arrays, the superconducting ones employed a silicon nitride membrane to support the thin film and provide thermal isolation. Excess noise at the contacts limited the performance

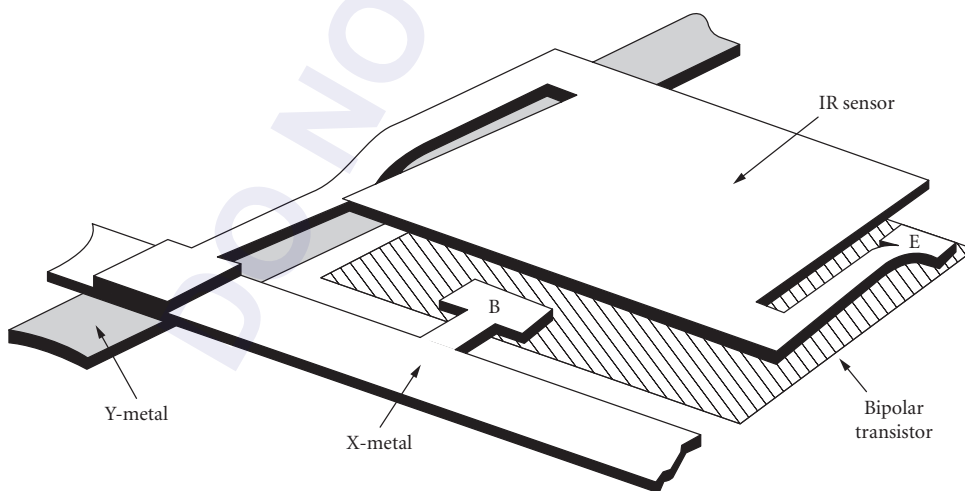


FIGURE 5 Monolithic silicon microbolometer.^{26,27} (© 1992 IEEE.)

of the 12-element linear array. With no excess noise, the calculated NETD²⁹ of a 240×336 array with $50 \mu\text{m}$ pixels and F/1 optics would be 0.002 K , which is near the 300 K background limit, as shown in Fig. 4.

Pyroelectric Hybrid Arrays

Linear and two-dimensional pyroelectric uncooled arrays have been under development for more than two decades.^{20–32} The arrays employ hybrid construction, in which a bulk pyroelectric ceramic material such as lithium tantalate or lead zirconate is mechanically thinned, etched, and polished, then bump-bonded to a silicon substrate containing readout electronics,³³ as shown in Fig. 6. Reticulation is usually employed to prevent lateral heat conduction through the pyroelectric material. The theoretical system NETD of a two-dimensional uncooled array with F/1 optics is estimated to be 0.1 K .³³ Two-dimensional uncooled arrays operating in the $8\text{--}14\text{-}\mu\text{m}$ region having 100×100 pixels, each $100 \times 100 \mu\text{m}$, are available commercially.³⁴ Their NETD (with F/1 speed) is 0.35 K . A two-dimensional pyroelectric monolithic array employing a thin film of lead titanate on a silicon microstructure is under development.³⁵

Ferroelectric bolometer arrays, also known as field-enhanced pyroelectric arrays, have been developed by Texas Instruments.^{36,37} Operation depends upon the temperature dependence of the spontaneous polarization and dielectric permittivity in a ferroelectric ceramic near the Curie temperature. Barium strontium titanate (BST), the selected material, has its composition (barium-to-strontium ratio) adjusted during preparation so the Curie point is 22°C . A thermoelectric stabilizer is employed to hold the BST near 22°C such that the absorbed infrared radiation changes the temperature and thus the dielectric properties. The effect is similar to the pyroelectric effect; however, a voltage is applied to enhance the signal. Construction of this array is naturally similar to that of the pyroelectric array, described above, as shown in Fig. 7. Reticulation of the ceramic is frequently applied to these arrays as well. A radiation chopper is required as both the pyroelectric and ferroelectric effects depend upon the change in temperature. The Texas Instruments BST array,

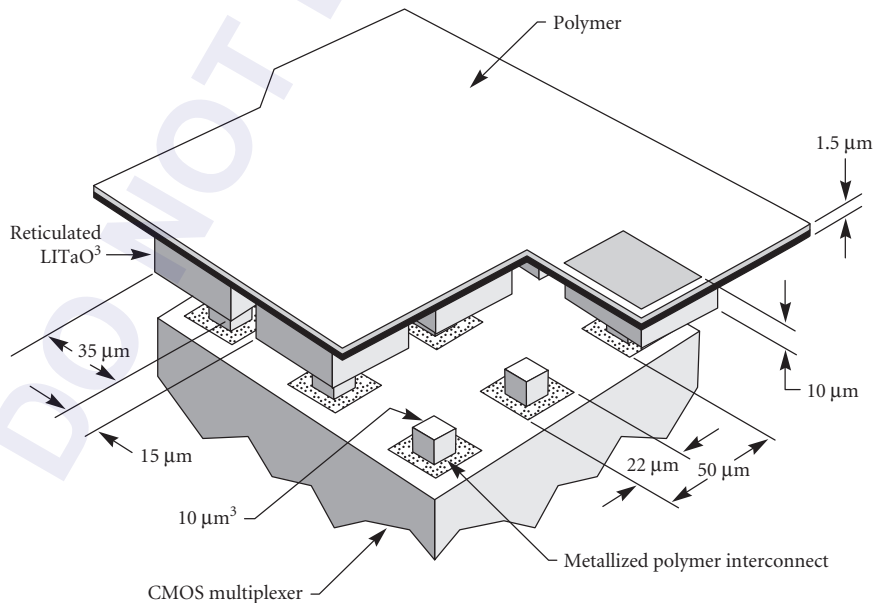


FIGURE 6 Hybrid pyroelectric array structure.³³

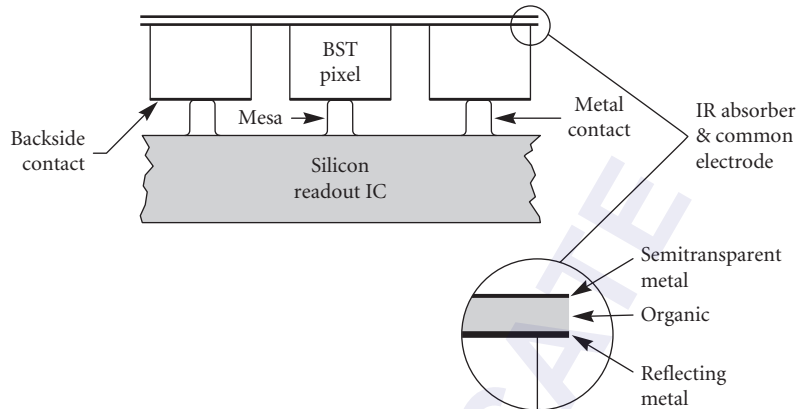


FIGURE 7 Ferroelectric bolometer hybrid array.^{36,37}

incorporating 80,000 pixels, each about $50 \times 50 \mu\text{m}$, which are matrix addressed, has an NETD of less than 0.1°C (with F/1 optics).

Thermoelectric Arrays

Thermoelectric arrays prepared by silicon micromachining have been described by Choi and Wise.³⁸ Series-connected, thin-film thermocouples, i.e., a thermopile, are prepared on a silicon microstructure, the “hot” junctions (receiving the thermal radiation) on a silicon nitride/silicon dioxide membrane and the “cold” shield junctions on the surrounding silicon substrate. Both 64- and 96-pixel microthermopile linear arrays in silicon microstructures have been prepared by Honeywell,³⁹ each microthermopile consisting of several nickel iron/chromium thermocouples connected in series, as shown in Fig. 8. The “hot” junctions are deposited on silicon nitride membranes, whereas the “cold” junctions are on the silicon substrates. A camera incorporating the linear array has been employed to image moving targets such as automobiles. With an F/0.73 lens, the measured NETD is 0.10°C .

Since the first publication of this *Handbook*, many advances have been made in these arrays. The suppliers have improved sensitivity somewhat but have increased the number of pixels and decreased their size. The reader should check with the manufacturers for the latest information.

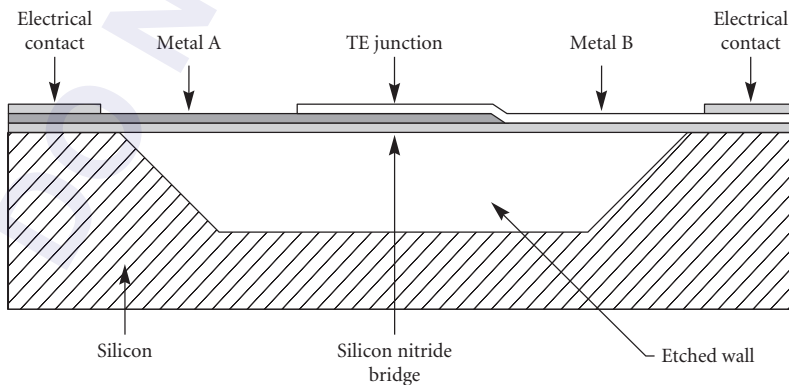


FIGURE 8 Monolithic thermoelectric array.³⁹ (© 1991 Instrument Society of America. Reprinted with permission from the Symposium for Innovation in Measurement Science.)

28.4 REFERENCES

1. P. W. Kruse, L. D. McGlauchlin, and R. B. McQuistan, *Elements of Infrared Technology*, John Wiley and Sons, New York, 1962.
2. R. D. Hudson and J. W. Hudson, *Infrared Detectors*, Halsted Press, New York, 1975; and F. R. Arams, *Infrared to Millimeter Wave Detectors*, Artech House, Dedham, Mass., 1973.
3. E. H. Putley, *Optical and Infrared Detectors*, R. J. Keyes (ed.), Springer-Verlag, Berlin, 1980, chap. 3.
4. R. A. Smith, F. E. Jones, and R. P. Chasmar, *The Detection and Measurement of Infrared Radiation*, Oxford University Press, 1968.
5. R. C. Jones, "The Ultimate Sensitivity of Radiation Detectors," *J. Opt. Soc. Am.* **37**:879 (1974).
6. S. Nudelman, "The Detectivity of Infrared Detectors," *Appl. Opt.* **1**:627 (1962).
7. F. J. Low, "Low Temperature Germanium Bolometer," *J. Opt. Soc. Am.* **51**:1300 (1961).
8. S. Zwerdling, R. A. Smith, and J. P. Thierault, "A Fast High Responsivity Bolometer for the Very Far Infrared," *Infrared Physics* **8**:271 (1968).
9. N. Coron, "Infrared Helium Cooled Bolometers in the Presence of Background Radiation: Optimal Parameters and Ultimate Performance," *Infrared Physics* **16**:411 (1976).
10. N. Coron, G. Dambier, and J. Le Blanc, *Infrared Detector Techniques for Space Research*, V. Manno and J. Ring (eds.), Reidel, Dordrecht, 1972.
11. N. Coron, G. Dambier, J. Le Blanc, and J. P. Moliac, "High Performance, Far Infrared Bolometer Working Directly in a Helium Bath," *Rev. Sci. Instr.* **46**:492 (1975).
12. W. H. Andrews, R. M. Milton, and W. De Sorbo, "A Fast Superconducting Bolometer," *J. Opt. Soc. Am.* **36**:518 (1946).
13. R. M. Milton, "A Superconducting Bolometer for Infrared Measurements," *Chem. Rev.* **39**:419 (1946).
14. N. J. Fuson, "The Infrared Sensitivity of Superconducting Bolometers," *J. Opt. Soc. Am.* **38**:845 (1948).
15. G. Gallinaro and R. Varone, "Construction and Calibration of a Fast Superconducting Bolometer," *Cryogenics* **15**:292 (1975).
16. K. B. Bhasin and V. O. Heinen (eds.), "Superconductivity Applications for Infrared and Microwave Devices," *Proc. SPIE* **1292** (1990). (Includes many other references.)
17. W. S. Boyle and K. F. Rodgers, *J. Opt. Soc. Am.* **49**:66 (1959).
18. D. F. Hornig and B. J. O'Keefe, "Design of Fast Thermopiles and the Ultimate Sensitivity of Thermal Detectors," *Rev. Sci. Instr.* **18**:7 (1947).
19. P. B. Felgett, "Dynamic Impedance and the Sensitivity of Radiation Thermocouples," *Proc. Phys. Soc.* **B62**:351 (1949).
20. M. J. E. Golay, "Theoretical Considerations in Heat and Infrared Detection with Particular Reference to the Pneumatic Detector," *Rev. Sci.* **18**:347 (1947); "Pneumatic Infrared Detector," *ibid.*, **18**:357 (1947); "Theoretical and Practical Sensitivity of the Pneumatic Detector," *ibid.*, **20**:816 (1949).
21. J. R. Hickey and D. B. Daniels, "Modified Optical System for the Golay Detector," *Rev. Sci. Instr.* **40**:732 (1969).
22. M. Chatanier and G. Gauffre, *IEEE Transactions Instr. and Meas.* **IMEE** **179** (1973).
23. A detector once made by Patterson Moos and cited by R. DeWaard and E. Wormser, "Description and Properties of Various Thermal Detectors," *Proc. IRE* **47**:1508 (1959).
24. E. H. Putley, *Semiconductors and Semimetals*, vol. 5, R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1970, chap. 6, "The Pyroelectric Detector;" vol. 12, 1977, chap. 7.
25. Zu-Sheng Wang and Jian-Qi Zhang, "The Mechanism of Type II Pyroelectric Detectors," *Infrared Phys.* **33**(6):481–486 II (1993).
26. R. A. Wood, C. J. Han, and P. W. Kruse, "Integrated Uncooled Infrared Detector Imaging Array," *Proc. of the 1992 IEEE Solid State Sensor and Actuator Workshop*, Hilton Head Island, S.C., pp. 132–135.
27. R. A. Wood, "Uncooled Thermal Imaging with Monolithic Silicon Focal Plane Arrays," *Proc. SPIE* **2020**: Infrared Tech. XIX (1993).
28. B. R. Johnson, T. Ohnstein, C. J. Han, R. Higashi, P. W. Kruse, R. A. Wood, H. Marsh, and S. B. Dunham, "High- T_c Superconductor Microbolometer Arrays Fabricated by Silicon Micromachining," *IEEE Trans. Appl. Superconductivity* **3**:2856 (1993).

29. B. R. Johnson and P. W. Kruse, "Silicon Microstructure Superconducting Microbolometer Infrared Arrays," *Proc SPIE* **2020**:Infrared Technology XIX (1993).
30. E. H. Putley, "The Pyroelectric Detector," *Semiconductors and Semimetals*, vol. 5, *Infrared Detectors*, R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1970.
31. P. A. Manning, D. E. Burgess, and R. Watton, "A Linear Pyroelectric Array IR Sensor," *Proc SPIE* **590**:2 (1985).
32. R. Watton and M. V. Mansi, "Performance of a Thermal Imager Employing a Hybrid Pyroelectric Detector Array with MOSFET Readout," *Proc. SPIE* **865**:79 (1987).
33. N. Butler and S. Iwasa, "Solid State Pyroelectric Imager," *Proc SPIE* **1685**:146 (1992).
34. GEC-Marconi Materials Technology Ltd., 9360 Ridgehaven Court, San Diego, CA 92123.
35. B. E. Cole, R. D. Horning, and P. W. Kruse, "PbTiO₃ Deposited by an Alternating Dual-Target Ion-Beam Sputtering Technique," *Ferroelectric Thin Films II*, A. I. Kingon, E. R. Myers, and B. Tuttle (eds.), *Materials Research Society Symposium Proc.*, **243**:185 (1992).
36. C. Hanson, H. Beratan, R. Owen, M. Corbin, and S. McKenney "Uncooled Thermal Imaging at Texas Instruments," *Infrared Detectors: State of the Art, Proc. of SPIE* **1735**:17 (1992).
37. C. M. Hanson, "Uncooled Ferroelectric Thermal Imaging," *Proc. SPIE* **2020**: Infrared Technology XIX (1993).
38. I. H. Choi and K. D. Wise, "A Silicon-Thermopile-Based Infrared Sensing Array for Use in Automated Manufacturing," *IEEE Trans. on Electron Devices* **ED-33**:72 (1986).
39. M. Listvan, M. Rhodes, and M. L. Wilson, "On-Line Thermal Profiling for Industrial Process Control," *Proc. of the Instrument Society of America, Symposium for Innovation in Measurement Science*, Geneva, NY, August 1991.

PART

6

**IMAGING
DETECTORS**

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

Joseph H. Altman

*Institute of Optics
University of Rochester
Rochester, New York*

29.1 GLOSSARY

A	area of microdensitometer sampling aperture
a	projective grain area
D	optical transmission density
D_R	reflection density
DQE	detective quantum efficiency
$d(\mu)$	diameter of microdensitometer sampling aperture stated in micrometers
E	irradiance/illuminance (depending on context)
\mathcal{G}	Selwyn granularity coefficient
g	absorbance
H	exposure
IC	information capacity
M	modulation
M_e	angular magnification
m	lateral magnification
NEQ	noise equivalent quanta
$P(\lambda)$	spectral power in densitometer beam
Q'	effective Callier coefficient
q	exposure stated in quanta/unit area
R	reflectance
S	photographic speed
$S(\lambda)$	spectral sensitivity
S/N	signal-to-noise ratio of the image
T	transmittance
$T(\nu)$	modulation transfer factor at spatial frequency ν

t	duration of exposure
$WS(\nu)$	value of Wiener (or power) spectrum for spatial frequency ν
γ	slope of D-log H curve
ν	spatial frequency
$\rho(\lambda)$	spectral response of densitometer
$\sigma(D)$	standard deviation of density values observed when density is measured with a suitable sampling aperture at many places on the surface
$\sigma(T)$	standard deviation of transmittance
$\phi(\tau)$	Autocorrelation function of granular structure

29.2 STRUCTURE OF SILVER HALIDE PHOTOGRAPHIC LAYERS

The purpose of this chapter is to review the operating characteristics of silver halide photographic layers. Descriptions of the properties of other light-sensitive materials, such as photoresists, can be found in Ref. 4.

Silver-halide-based photographic layers consist of a suspension of individual crystals of silver halide, called *grains*, dispersed in gelatin and coated on a suitable “support” or “base.” The suspension is termed an *emulsion* in the field. The grains are composed of AgCl, AgClBr, AgBr, or AgBrI, the listing being in order of increasing sensitivity. Grain size ranges from less than 0.1 μm (“Lippmann” emulsions) to 2 to 3 μm , depending on the intended use of the coating. The number of grains per square centimeter of coating surface is usually very large, of the order of 10^6 to 10^8 grains/ cm^2 . The weights of silver and gelatin coated per unit area of support vary depending on intended use; usually both fall in the range 1 to 10 g/m^2 . The silver-to-gel ratio may also vary depending on intended use. Typically, the emulsion may be about 30 to 40 percent silver by weight, but some special-purpose materials, such as films to record Schumann-wavelength-region radiation, contain very little gelatin.

For modern materials, both the emulsion and the coating structure can be very complex. The emulsion layer is much more than silver halide in gelatin, containing additional agents such as hardeners, antifoggants, fungicides, surfactants, static control agents, etc. Likewise, the coating structure may be very complex. Even some black-and-white materials consist of layers of two different emulsions coated one over the other and a thin, clear layer of gelatin is often coated over the emulsion(s) to provide some mechanical protection. In the case of color films, as many as 15 layers may be superimposed, some of them of the order of 1 μm thick. The thickness of the complete coating varies from about 3 μm to about 25 μm in normal films.

Commercially available emulsions are coated on a variety of glass, plastic (film), and paper supports (or “bases”). Glass is used for mechanical rigidity, spatial stability, or surface flatness.

Two different types of plastic are available commercially as film supports: cellulose acetate and polyethylene terephthalate (trade names “Cronar,” “Estar,” and “Mylar”). Of the two types, Mylar is superior in strength, flexibility, and spatial stability. However, the material is birefringent and its physical properties may be different in orthogonal directions. Also, these directions may not be aligned with the length or width of the sample. The anisotropic properties arise from the method of manufacture. Although not as tough as Mylar, cellulose acetate is, of course, fully adequate for most purposes. Also, this material is isotropic and easier to slit and perforate. Typical supports for roll films are around 4 mils (102 μm) thick, and for sheet films, 7 mils (178 μm), and other thicknesses are available. Most films are also coated on the back side of the support. The “backing” may be a layer of clear gelatin applied for anticurl protection, or of gelatin dyed with a dye that bleaches during processing, and provides both anticurl and antihalation protection. Lubricants and antistatic agents may also be coated, either on the front or back of the film. Properties of supports are discussed in Ref. 1.

29.3 GRAINS

The grain is the radiation-sensing element of the film or plate. It is a face-centered cubic crystal, with imperfections in the structure. For the most part, the grains act as independent receptors. In general, the larger grains are faster. The properties of the individual grains are controlled by the precipitation conditions and the after-precipitation treatment. Details of these matters are proprietary, but some discussion is given in Refs. 2 and 3. From the user's standpoint, the important fact is that when the grain is exposed to sufficient radiation it forms a "latent-image speck" and becomes *developable* by a solid-state process called the "Gurney-Mott mechanism." An excellent review of grains and their properties is given by Sturmer and Marchetti in chap. 3 of Ref. 4.

29.4 PROCESSING

The exposed halide layer is converted to a usable image by the chemical processes of development and fixation.

Development consists of reducing exposed crystals from silver halide to metallic silver, and a developing agent is an alkaline solution of mild reducer that reduces the grains having latent image specks, while not attacking the unexposed grains. Generally, once development starts the entire grain is reduced if the material is allowed to remain in the developer solution. Also, in most cases adjacent grains will not be affected, although developers can be formulated that will cause adjacent grains to be reduced ("infectious developers"). The number of quanta that must be absorbed by a grain to become developable is relatively small, of the order 4 to 40, while the developed grain contains on the order of 10^6 atoms. The quantum yield of the process is thus very high, accounting for the speed of silver-based materials.

The remainder of the process consists essentially of removing the undeveloped halide crystals which are still light-sensitive. The "fixer" is usually an acid solution of sodium thiosulfate $\text{Na}_2\text{S}_2\text{O}_3$, called "hypo" by photographers. The fixing bath usually serves as a gelatin hardener also. The thiosulfate reacts with the halide of the undeveloped grains to form soluble silver complexes, which can then be washed out of the emulsion layer. It is worth noting that proper washing is essential for permanent images. Additional treatments to promote permanence are available. Processing is discussed in detail in Refs. 2, 3, and 5, and image permanence in Ref. 6.

The exposed and processed silver halide layer thus consists of an array of grains of metallic silver, dispersed in a gelatin matrix. In color films, the silver is removed, and the "grains" are tiny spheres of dyed gelatin (color materials will be discussed below). Either type of grain acts as an absorber; in addition, the metallic silver grains act as scatterers. The transmittance or reflectance of the layer is thus reduced and, from the user's standpoint, this change constitutes the response of the layer.

29.5 EXPOSURE

From fundamental considerations it is apparent that the dimensions of exposure must be energy per unit area. Exposure H is defined by

$$H = Et \quad (1)$$

where t is the time for which the radiation is allowed to act on the photosensitive layer, and therefore E must be the irradiance on the layer. The symbol H is used here for exposure in accordance with international standards, but it should be noted that in many publications, especially older ones, E is used for exposure and I for irradiance, so that the defining expression for exposure becomes $E = It$.

Strictly speaking, H and E in Eq. (1) should be in radiometric units. However, photographic exposures are customarily stated not in radiometric but in photometric units. This is done mostly

for historical reasons; the English scientists Hurter and Drifffield, who pioneered photographic sensitometry in 1891, measured the incident flux in their experiments in lumen per square meter, or lux. Their unit of exposure was thus the lux-second (old term, meter-candle-second). Strictly speaking, of course, weighting the incident flux by the relative visibility function is wrong or at least unnecessary, but in practice it works well enough because in most cases the photographer wishes to record what he or she sees, i.e., the visible spectrum. Conversion between radiometric and photometric units is discussed by Altman, Grum, and Nelson.⁷

Also, it should be noted that equal values of the exposure product (Et) may produce different outputs on the developed film because of a number of *exposure effects* which are described in the literature.⁸ A complete review of radiometry and photometry is given in Chap. 34 in this volume of the *Handbook*.

29.6 OPTICAL DENSITY

As noted above, the result of exposure and processing is a change in the transmittance or reflectance of the layer. However, in photography, the response is usually measured in terms of the *optical density*, hereafter called the “density” in this chapter. For films (transmitting samples), density is defined by

$$D = -\log T = \log 1/T \quad (2)$$

where T is the transmittance. (Note: throughout this chapter, “log” indicates the base-10 logarithm.) For either silver grains or color grains, the *random dot* model of density predicts that

$$D = 0.434nag \quad (3)$$

where n = the number of grains per unit area of surface
 a = the average projective grain area
 g = the absorbance of the grain

Absorbance in turn is defined as $g = 1 - (T + R)$ where T and R are the transmittance and reflectance of the grain. For silver grains, the absorbance is taken as unity. The above expression, sometimes known as “Nutting’s law,” is based on a geometric approach, and does not take into account any scattering by the grains. However, of course, opaque silver particles on a clear background will act as scatterers, and in fact multiple scattering usually occurs in developed silver layers. This produces an increase in the density of such layers by a factor of 2 to 3 times from that predicted by Nutting’s law. For color films, the refractive index of the gelatin in the dyed spheres is only negligibly different from that of the surround so that such layers are not scatterers. Even in the case of silver films, however, Nutting’s law provides a useful model. Since for a given population of grains a and g will remain effectively constant, the law states that density should vary as the number of grains per unit area of surface. This fact is easily verified with a microscope.

Transmission Density

Transmission density is measured in a densitometer. It is worth noting that the device actually measures the transmittance of the sample and then displays the negative log of the result. In a normal or *macro* densitometer the sampling aperture area A is typically 1 mm² or more in size. When A is small, say, 0.1 mm² or less, the device becomes a *micro* densitometer. Microdensitometers present special problems and will be discussed below.

Because the scattered light may not reach the sensor of the densitometer when silver layers are measured, it is necessary to specify the angular subtenses of both the incident (influx) and emergent (efflux) beams at the sample. Clearly, if scattered light is lost to the sensor, the indicated density of the sample will *increase*. Four types of transmission density are described in an ISO standard,⁹ of

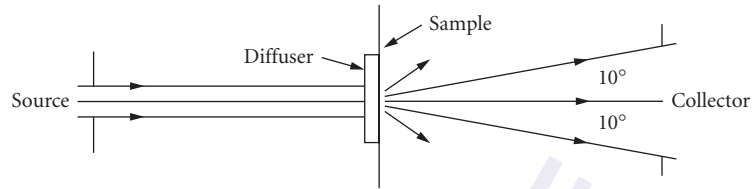


FIGURE 1 Optical system for measuring ISO/ANSI diffuse density with a 20° collection angle. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

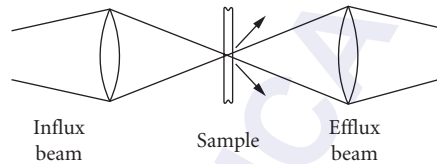


FIGURE 2 Optical conditions for projection density measurement.

which two are principally important to the user. The first of these is *diffuse density*, which is diagrammed in Fig. 1. As can be seen, a collimated incident beam illuminates an opal glass diffuser. The emulsion side of the sample is placed in contact with this diffuser, and the flux contained within a cone angle of $\pm 10^\circ$ is collected and evaluated by the sensor. The reverse of this arrangement yields the same density values and is also permitted by the standard. This is the type of density normally measured in practice. Physically, it corresponds to the conditions of contact printing.

The other case that is important in practice is projection density, which is diagrammed in Fig. 2. This case simulates the use of the layer in an optical system. As the figure shows, light is lost to the efflux system in projection density, the exact amount depending on the numerical aperture of the optics involved and the scattering characteristics of the sample. Thus, the projection density of a silver film is usually greater than the diffuse density. The effective Callier coefficient Q' may be defined by

$$Q' = \frac{\text{project density}}{\text{diffuse density}} \quad (4)$$

This factor can be measured experimentally. For silver films and $f/2$ optics, $Q' \approx 1.3$; for color films $Q' \approx 1.0$.

Nonneutral (Color) Density

In many cases, silver densities can be treated as neutrals. For dye densities, i.e., color films, measured density depends on the spectral characteristics of both the dye and the densitometer. The spectral response of the instrument is given by

$$\rho_\lambda = P_\lambda S_\lambda F_\lambda \quad (5)$$

where ρ_λ = the response at wavelength λ

P_λ = the power in the densitometer beam at λ

S_λ = the spectral sensitivity of the sensor at λ

F_λ = the transmittance of the densitometer optics at λ , specifically including any filters placed in the densitometer beam

The measured density of a nonneutral layer is then

$$D = \log \left[\frac{\int_{\lambda_1}^{\lambda_2} \rho_\lambda d\lambda}{\int_{\lambda_1}^{\lambda_2} \rho_\lambda T_\lambda d\lambda} \right] \quad (6)$$

where T_λ is the transmittance of the layer at wavelength λ , and the wavelength limits are set by the distributions. The response ρ_λ of the system is adjusted to be equal to that of the readout device with which the film is to be used.

Thus, for example, if the sample is to be viewed by an observer, ρ_λ is made equal to the visibility function, and the resulting measurement is called visual density, etc. Instrument responses have been standardized for sensitometry of color films.¹⁰

Reflection Density

When the emulsion is coated on paper the density is measured by reflection. Reflection density is then defined by

$$D_R = -\log R \quad (7)$$

where R is reflectance, measured under suitable geometric conditions. The measurement of reflection density is also described in the standards literature.¹¹

29.7 THE D-LOG H CURVE

In routine sensitometry, samples receive a series of exposures varying by some constant factor, such as $\times 2$ or $\times \sqrt{2}$. After processing, the measured densities are plotted against the common logarithm of the exposures that produced them. The resulting curve is known as the “D-log H curve,” or the “H & D” curve (after Hurter and Driffled, the previously mentioned pioneers in the field). A typical D-log H curve is shown in Fig. 3.

As shown, the curve is divided into three regions, known as the “toe,” “straight-line portion,” and “shoulder,” respectively. The fact that an appreciable straight-line portion is found in many cases is

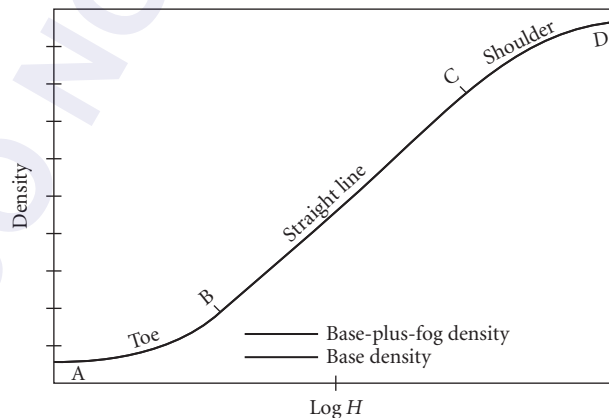


FIGURE 3 Typical D-log H curve for a negative photographic material. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

not an indication that the film is a linear responder in this region because, of course, both axes of the plot are logarithmic. It is worth noting here that the equation of the straight-line portion of this curve can be written

$$D = \gamma(\log H - \log C) \quad (8)$$

where γ is the slope of the straight-line portion and C is the exposure at the point where the extrapolated straight line cuts the exposure axis. Taking antilogarithms, Eq. (8) becomes

$$T = \left(\frac{H}{C}\right)^{-\gamma} \quad (8a)$$

If $\gamma = -1$, $T = (1/C)H$, and for this special case, the system becomes linear over the exposure range corresponding to the straight-line portion. A negative value of γ indicates a *positive* image.

A number of useful performance parameters for films are taken from their D-log H curves, as follows.

1. *Fog*: For most films, a certain number of grains will be reduced even though they have received no exposure at all, or insufficient exposure to form a latent image speck. The resulting density is called *fog*. Since it is not exposure-related, and since it tends to veil any information recorded in the toe, excessive fog is very undesirable. For many purposes, the fog density plus the density of the support are subtracted from the gross density to give the value of the net density resulting from the exposure. More complicated formulas for correcting the film's response for fog grains have been proposed, but are rarely used.
2. *Gamma*: Traditionally, the slope of the straight-line portion of the D-log H curve is called the "gamma." Gamma is a crude measure of the contrast with which the original object is reproduced; it would be a good measure of this contrast if the object were in fact recorded entirely on the straight-line portion of the response curve. However, for many purposes, notably pictorial photography, an appreciable part of the toe is used. This fact led Niederpruem, Nelson, and Yule to propose the use of an average gradient that included part of the toe as a measure of the contrast of the reproduction.¹² This quantity is called the *contrast index*. Since it includes part of the toe, contrast index is less than gamma. Since information is often recorded in the toe, it is convenient to generalize the meaning of γ to refer to the gradient anywhere along the D-log H curve, and this is done in this chapter. Note that in this case, the traditional gamma is the maximum value the gradient attains.
3. *Latitude*: Latitude can be defined as the log exposure range between the point in the toe and the point in the shoulder between which the gradient is equal to or greater than the minimum value required for acceptable recording. Clearly, the latitude of the film must be at least equal to the log exposure range of the object for proper recording. In many practical cases the film's latitude easily exceeds the required range. Note that the latitude is determined in part by the maximum density that the film can produce.
4. *Speed*: Speed is defined by the general expression

$$S = \frac{K}{H_{\text{ref}}} \quad (9)$$

where K is a proportionality factor and H_{ref} is the exposure required to produce some desired effect. Since the desired effect varies depending on the type of film and the application, H_{ref} also varies. Also, H_{ref} can be stated in either radiometric or photometric units. If H_{ref} is given in radiometric units, the proportionality factor K is set to unity, and the resulting values are termed "radiometric speeds." Although radiometric speeds are the fundamental speed values, they are rarely used in practice because, as previously noted, exposures are usually given in photometric units. In this case, the factor K takes on different values that depend not only on H_{ref} , but also on the characteristics of the exposure meter, which is standardized.¹³ Varying the factor K allows a single meter to be used with

all kinds of films and applications, which is a practical necessity. Thus, for example, the photometric speed (usually simply the “speed”) of black-and-white pictorial films is evaluated from

$$S = \frac{0.8}{H_{0.1}} \quad (9a)$$

where $H_{0.1}$ is the exposure in lux-seconds required to produce a density of 0.1 above the densities of the base plus fog in a specified process. This density level has been shown empirically to be predictive of excellent tone reproduction quality in the print. Similarly, for color slide films

$$S = \frac{10}{H_m} \quad (10)$$

where H_m is the exposure to reach a specified position on the film’s D-log H curve. Again, H_m was established empirically. The above two examples show how the two factors involved in determining a photometric or practical speed can vary. A number of other speeds have been defined and are described in the literature.¹⁴

Variation of Sensitometry with Processing

The rate of reduction of the exposed photographic grain depends on the characteristics of the grain itself, the formulation of the developer, and its temperature. In general, the reaction is allowed to continue until substantially all exposed grains have been reduced, and ended before fog becomes excessive. Many modern films are hardened and able to withstand processing temperatures up to, say, 40°C. Development times are often chosen on the basis of convenience and usually run in the order of 5 to 10 minutes in nonmachine processing. As development time and/or temperature are increased, the amount of density generated naturally increases. Thus for a given film a whole family of response curves can be produced, as shown in Fig. 4. As development time is lengthened, gamma and contrast index also increase. Typical behavior of these parameters is shown in the inset of Fig. 4.

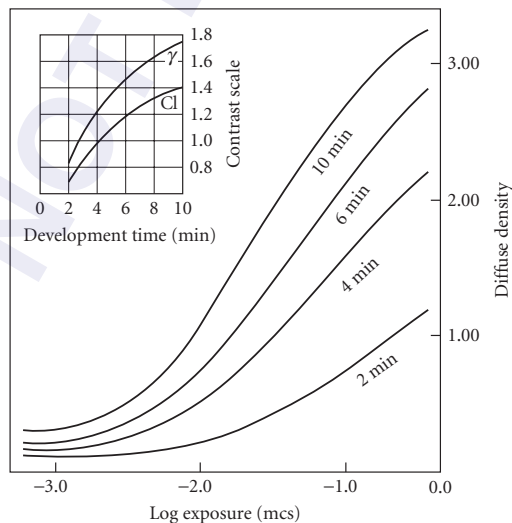


FIGURE 4 A family of characteristic curves for development times as shown, with corresponding plots of contrast index and gamma, mcs or the meter-candle second, as stated on p. 29.6, is the old term for lux-second or lxs. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

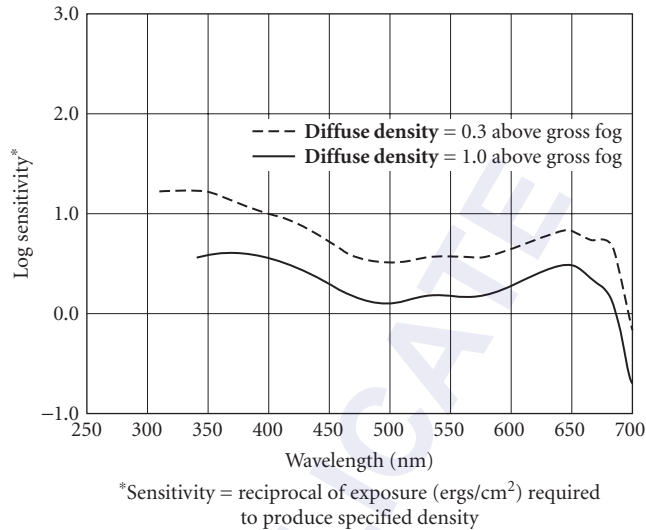


FIGURE 5 Spectral sensitivity curves for a modern negative material sensitized to about 690 nm. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

29.8 SPECTRAL SENSITIVITY

The spectral absorption of the silver halide grain extends only to about 500 nm, and thus the inherent sensitivity of the grain is limited to regions shorter than that limit. However, it was discovered by Vogel in 1873 that sensitivity could be extended to longer wavelengths by dyeing the grains, and over the years the effective sensitivity range was extended into the infrared. Presently, materials are available usefully sensitized out to about 1.2 μm , but it should be noted that IR materials tend to have poor shelf life, and may require special handling. For our purposes, we may say that the sensitizing dye absorbs the incident energy in the required spectral region and transfers the energy to the grain in a manner that produces the required latent-image speck. The mechanisms are discussed in Ref. 3, chap. 10.

The spectral sensitivity of a photographic layer is usually specified by a family of curves showing $\log(1/H_D)$ vs λ , where H_D is the exposure in ergs/cm² of wavelength λ required to produce some stated density. Spectral sensitivity values are thus radiometric speeds. Typical curves are shown in Fig. 5. In practical work, three broad classes of sensitization are recognized, which are called “color-blind” (or blue-sensitive), orthochromatic (additionally sensitized to green), and panchromatic (additionally sensitized to green and red). Most modern general-purpose materials are panchromatic.

In general, the shape of the spectral sensitivity curve follows that of the spectral absorption of the layer. It should also be noted that the gradient of the D-log H curve will be affected by the absorption of the layer. Gamma may therefore vary as a function of wavelength, and is generally somewhat lower in the blue and UV regions of the spectrum. This means that if the material is being used as a radiometer, it must be calibrated at the wavelength(s) of interest.

29.9 RECIPROCITY FAILURE

It was noted above that the exact response of a photographic layer may change due to exposure effects. Of these, the phenomenon of *reciprocity failure* is the most important in practical photography.

By definition [Eq. (1)] exposure is simply the product of the irradiance and time, and nothing in this definition specifies the magnitude of either factor. However, the developed *density* resulting from a given calculated exposure is often found to depend on the *rate* at which the radiation is supplied, all other factors being held constant. Broadly speaking, exposure times of about 0.01 to 1.0 second are most efficient in producing density, the exact values depending on the film involved. Times much outside the above range tend to produce lower density for the same calculated exposure. The emulsion-maker has some ways of minimizing the effect, and usually attempts to optimize the response for the exposure times expected to be used with the material. Reciprocity failure is discussed in detail in Ref. 3, chap. 4, sec. II.

The loss of efficiency for short-time and correspondingly high-irradiance exposures is termed “high-intensity reciprocity failure,” and that for long exposure times (low irradiances) is termed “low-intensity reciprocity failure.” The Gurney-Mott mechanism explains both types of failure well. Note also that the names are misnomers; the terms should, of course, be high and low *irradiance*.

Only limited data are available, but the gradient of the D-log H curve tends to decrease as exposure times are shortened. An example is shown in data published by Hercher and Ruff.¹⁵ Limited data also indicate that the speed loss due to high-intensity failure stabilizes for times shorter than about 10^{-5} second.¹⁶ Essentially, the amount of reciprocity failure is independent of exposing wavelength.¹⁷ This is to be expected.

Reciprocity failure may be a considerable problem for photographers working in specialized time domains, such as oscilloscope photography or astronomy. Astronomers have been able to devise user treatments for minimizing low-intensity failure.¹⁸ For the practical photographer, reciprocity failure sometimes appears as a problem in color photography. If the RF of the three sensor layers differs, the resulting picture may be “out of balance,” i.e., grays may reproduce as slightly tinted. This is very undesirable, and correction filter recommendations are published for some films for various exposure times.

29.10 DEVELOPMENT EFFECTS

Besides exposure effects the final density distribution in the developed image may be affected by “development effects” arising from chemical phenomena during development. Various names such as “border effect,” “fringe effect,” “Eberhard effect,” etc., are applied to these phenomena; what we shall here term “edge effect” may be important in practice.

Consider a sheet of black-and-white film developing in a tray, and assume for purposes of discussion that there is no motion of the developer. Since developing agent must be oxidized as halide is reduced, and since the by-products of this reaction may themselves be development inhibitors, it can be seen that local variations of developer activity will be produced in the tray, with the activity decreasing as density increases. Agitation of the solution in the tray reduces the local variations, but usually does not eliminate them entirely, because it is the developer that has diffused into the gelatin matrix that is actually reacting. Now an “edge” is a boundary between high- and low-density areas, as shown in Fig. 6. Because of the local exhaustion and the diffusion phenomena, the variation of developer activity within the layer will be as shown by the dotted line in the figure. The result is that the developed density near the edge on the low-density side tends to decrease, and on the high-density side tends to increase, as also shown in the figure. In other words, the density distribution at the edge is changed; this actually occurs to some degree in much practical work and has interesting consequences, as will be discussed below.

The local exhaustion of the developer may also be important in color films where development in one layer (see below) may affect the response of an adjacent layer. In color photography, the phenomenon is called “interimage effect.”

29.11 COLOR PHOTOGRAPHY

Color photography has been extensively reviewed by Kapecki and Rodgers.¹⁹ With one exception at the time of writing, all commercially available color films employ subtractive color reproduction. The exception is an instant film for color slides marketed by Polaroid Corporation, which employs

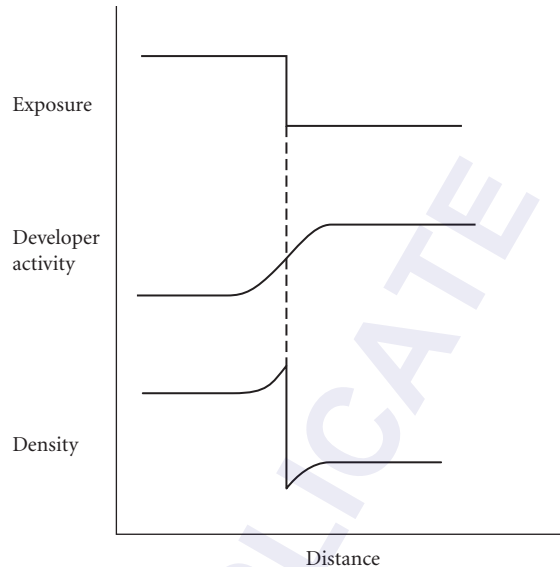


FIGURE 6 Distributions of developer activity and density resulting from a step-function input in the presence of edge effect.

additive color. The mechanism involves a reseau of very fine red, green, and blue stripes, which provide the color separation during the taking of the photograph, and also the color when the (reversed) image is projected.²⁰

The basic structure of most other color films is similar to that shown in Fig. 7. The incoming light first encounters a nonspectrally-sensitized emulsion layer, which records the blue-light elements of the scene. The next layer is a yellow, or minus-blue filter, the purpose of which is to prevent any blue light from reaching the other two emulsion layers. This yellow filter is usually composed of "Carey Lea," or colloidal silver dispersed in gelatin. Such sols are yellow. The reason for using Carey Lea silver is that all metallic silver is removed from the film during processing anyway, and the necessary removal of the filter layer is thus accomplished automatically.

Moving downward in the stack, the next layer is an ortho-sensitized emulsion. Since any blue light has been blocked by the yellow filter, this layer records the green-light elements of the scene. The final layer in the stack is sensitized to red light but not to green light and this layer serves to record the red elements of the scene. Modern films usually contain many more than the four layers indicated here, but the operating principles of the film can be discussed in terms of such a "tripack." In accordance with the principles of subtractive color reproduction, the images in the three separation

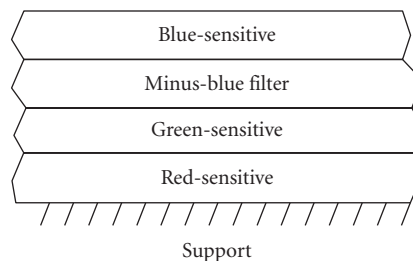


FIGURE 7 Schematic tripack structure of a color film.

records are then converted to yellow, magenta, and cyan dyes, respectively, and the final image is composed of these dyes, without the developed silver, which is either removed chemically or left behind, depending on the exact material.

Within this broad framework, films can be separated into two classes: chromogenic and nonchromogenic. In the former class, the dyes are not coated in the film, but are formed during the processing by a reaction called “coupling.” In coupling, the by-products of the halide reduction reaction serve as components for a second reaction in which dye is formed; the amount of dye thus increases as density increases. The components of the dye-forming reaction (i.e., other than the development by-products) may be present in the developing solution (Eastman Kodak Co., “Kodachrome”) or, more generally, coated within the various layers (Kodacolor, Polaroid “One Film,” Agfachrome, Fujicolor, etc.). After the required dyes are formed, the developed silver is removed by a chemical process termed “bleaching.”

The advantage of incorporating the couplers in the various layers of the tripack is simpler processing, but because of the additional material in the layers, such films tend to be not as sharp as the nonincorporated-coupler types. Chromogenic color films are available both as slide materials, in which the film undergoes a reversal process,²¹ and as color negative—color print materials in which a dye negative is formed and then printed onto a color paper whose structure is fundamentally similar to that of the films.

The principal example of the nonchromogenic film is the Polaroid Instant Color Film. In this film the yellow, magenta, and cyan dyes are actually coated in the structure, along with the blue, green, and red-sensitive emulsions. When development takes place in a given layer, the corresponding dye is immobilized. The dye that has *not* been immobilized in the three layers migrates to a “receiver” layer, where it is mordanted. Since the amount of dye that migrates *decreases* as the original density *increases*, the result is a positive color image formed in the receiver. The material has been described in more detail in a paper by Land,²² and also in chap. 6 of Ref. 4.

The image in a color film thus essentially consists of three superimposed dye images. Typical spectrophotometric curves for dyes formed by coupling reactions are shown in Fig. 8. The density of any one of these dye layers taken by itself is known as an *analytical density*. Note, however, that all the dyes show some “unwanted” absorption—that is, absorption in spectral regions other than the specific region that the dye is supposed to control. Thus the total density of the layer at any wavelength is the sum of the contributions of all three dyes; this type of density is known as “integral density.” The integral density curve is also shown in Fig. 8.

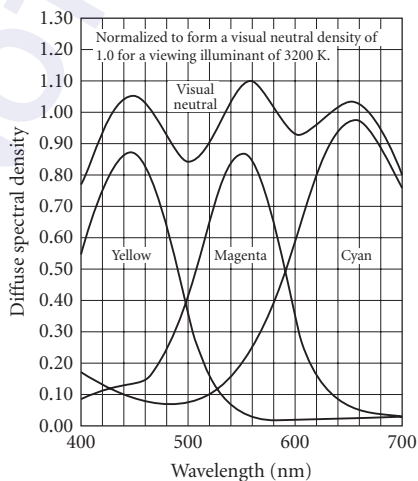


FIGURE 8 Spectral dye density curves for a color film. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

The reproduction of color by photographic systems has been discussed by Hunt²³ and many others.²⁴ In general, exact colorimetric reproduction is not achieved, but for most purposes the reproduction of the hues and luminances in the original scene is satisfactory. Evans,²⁵ in fact, observes that under the right conditions the “magnitude of the reproduction errors that can be tolerated is astonishing.” One aspect that is critical for many workers, especially expert photographers, is the ability of the system to produce good “balance,” i.e., to reproduce a neutral as a neutral. This requirement is so important that one of the types of density that is measured for dye layers is the *equivalent neutral density* (END) which is defined as the visual neutral density that a dye layer would produce if combined with the correct amounts of the other two dyes (whatever those correct amounts may be). When the ENDs of the three layers are equal, the system is in balance, and color-film sensitometry is therefore often done in terms of ENDs. Further discussion of color sensitometry and densitometry may be found in Ref. 3, chap. 18.

29.12 MICRODENSITOMETERS

As indicated above, a microdensitometer is a densitometer designed to measure the density of a small area. The sampling apertures are typically slits which may be as narrow as 1 to 2 μm in nominal width. The sample is scanned over the aperture, creating a record of density as a function of position on the sample surface, i.e., distance.

In practical instruments, the small sampling aperture dimensions are achieved by projecting an enlarged image of the film onto a physical aperture. The optical system produces some effects not encountered in macrodensitometers, as follows:

1. In general, microdensitometers measure projection, or semispecular, density. As already noted, projection density is higher than diffuse density for silver layers, and the exact value of the effective Callier coefficient Q' depends in part on the numerical aperture of the optical system. Thus two microdensitometers fitted with optics of different NAs may give different density values for the same sample. Furthermore, macrodensity data are usually in terms of diffuse density, so that data from the microdensitometer should be corrected if intercomparisons are to be made. The effective Callier coefficient for the specific optics-sample combination is easily determined by measuring suitable areas of sample both in the microdensitometer and a macrodensitometer, and taking the ratios of the values.
2. The presence of stray light in the system tends to lower the measured density. This problem is especially troublesome in microdensitometry because of the types of images that are often encountered. Thus, for example, when an interface between clear and dense areas—that is, an edge—is scanned, stray light will distort the record in the manner shown in Fig. 9. If the image is that of a star or spectroscopic line, this behavior results in an artificially low density reading. It is very important to control stray light as completely as possible in microdensitometry.

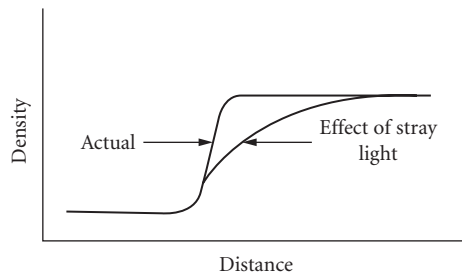


FIGURE 9 Effect of stray light in the microdensitometer on the apparent density distribution at an edge.

3. One feature of the usual microdensitometer optical system, installed for control of stray light, is the “preaperture,” a field stop that limits the area of the sample that is illuminated. Because of this preaperture, and the optical system, a normal microdensitometer is a partially coherent system. This means that the instrument may respond to path-length differences in the sample as well as to density differences. This is undesirable since in practice it is the density differences that carry the information. Partial coherence in microdensitometers has been studied by Thompson and by Swing among others, and the results are summarized by Dainty and Shaw.²⁶ It has been shown that the coherence effects can be minimized by satisfying two conditions. In the first condition, the width W of the preslit

$$W \geq \frac{4\lambda}{NA_{in}} \quad (11a)$$

The second condition is

$$\frac{NA_{in}}{NA_{eff}} = 1 + \frac{\nu_s}{\nu_o} \quad (11b)$$

where λ = the wavelength of the light
 NA_{in} = numerical aperture of the influx optics
 NA_{eff} = numerical aperture of the efflux optics
 ν_s = maximum spatial frequency in the sample
 ν_o = spatial frequency cutoff of the scanning objective

If, for example $\nu_o = 3\nu_s$, then $NA_{in}/NA_{eff} \geq 1.3$. A microdensitometer arranged to minimize coherence problems is called a *linear* microdensitometer. Note that the two conditions above conflict with conditions commonly adopted to control stray light.

29.13 PERFORMANCE OF PHOTOGRAPHIC SYSTEMS

The following discussion of performance is limited to those aspects which are properties of the system, and excludes such aspects as the skill of the photographer, etc. Of those aspects, which we may term “technical” quality parameters, the most important is *tone reproduction*. The subject is divided into two areas: subjective tone reproduction and objective tone reproduction. Subjective tone reproduction is concerned with the relation between the brightness sensations produced in the observer’s mind when the scene is viewed and when the picture is viewed. Since the sensation of brightness depends markedly on the viewing conditions and the observer’s state of adaptation, the subjective tone reproduction of a given picture is not constant, and this aspect of the general subject is not often measured in the photographic laboratory. It is discussed by Kowalski.²⁷ Objective tone reproduction is concerned with the reproduction of the luminance and luminance differences of the scene as luminances in the final output. Tone reproduction studies apply equally well to projected images, prints, transparencies, and video images, but a negative-print system is usually assumed for discussion. It is easy to show that the log luminance of a print area, $\log L_p = C - D_p$, where C is a constant determined by the illuminance on the print and D_p is the density of the print area. Thus tone reproduction curves are usually plotted in the form of the print density versus the log luminance of the corresponding scene element (Fig. 10). Although both the scene luminances and the print densities are fixed quantities, the viewer’s reaction still depends on the illumination level at which the picture is seen.

It has been shown empirically that for paper prints viewed under typical “room lighting conditions,” the preferred tone reproduction curve is the solid line in Fig. 10. Perfect objective tone reproduction, defined as the case where $\Delta D_p = -\Delta \log L_{sc}$ for all scene luminance levels, would be

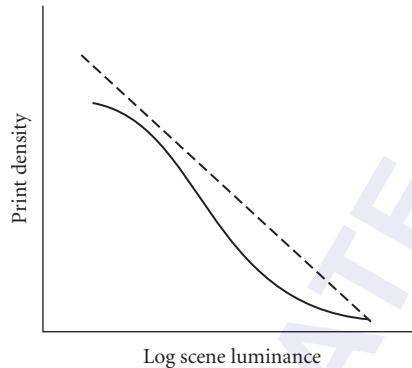


FIGURE 10 Preferred tone reproduction for viewing under typical room lighting conditions.

the dotted line in the figure. Thus under typical room lighting conditions viewers prefer a reproduction that has somewhat more contrast and less density than the “perfect” result. This preference will, however, change with illumination and stray light levels.

The exact shape of the tone reproduction curve obtained with a given system depends on the shape of the negative and positive D-log H curves, and also on the stray-light characteristic of the camera system. (Excessive stray light in the camera can be very deleterious to tone reproduction quality.) Using a graphic method devised by Jones and described by Nelson and others,²⁸ the effect of the D-log H curves and the stray light on the tone reproduction can be studied.

29.14 IMAGE STRUCTURE

The other two technical quality parameters of a photographic system are its sharpness and graininess, to use the most familiar terms for these properties. These properties are often lumped together under the general term “image structure.” Actually, in photoscience the term *sharpness* and *graininess* are reserved for the subjective aspects of the phenomena, and measurement of these properties requires psychometric testing. Since such testing is expensive and time-consuming, objective correlates of both properties have been defined, and methods for measurement have been established. The objective correlate of sharpness is termed *acutance*, and of graininess is termed *granularity*. Image structure data for various kinds of materials are given in Table 1, along with speed and contrast values.

29.15 ACUTANCE

The original proposal for measuring acutance was made by Higgins and Jones²⁹ in 1952. It involved calculating a value from a microdensitometer trace of a test edge. However, the visual processes that occur when an observer views an edge are complex, and the straightforward calculation proposed by Higgins and Jones fails to predict the sensation of sharpness produced by some edge distributions. At about the same time, the optical transfer function (OTF) and related concepts began to be widely used in optics, and these were soon applied to photographic materials also.

The concepts of the point and line spread functions are essentially identical in optics and photoscience; in the photographic layer the smearing of the point (or line) input is caused by diffraction around the grains, refraction through them, and reflection from them. These phenomena are usually

TABLE 1 Performance Data for Types of Materials^a

Type	SPD ^b	γ	Gran. ^c	ν_{50} ^d
B&W microfilm	80	3.0	6	>200
B&W very slow camera neg.	25	0.5–3.5	5–7	80–145
B&W slow camera neg.	100	0.5–1.1	8–9	65–120
B&W fast camera neg.	400	0.5–1.0	10–14	50–100
B&W very fast camera neg.	1600–3200	0.5–1.0	18	70
Color neg. very slow	25–50	0.65	4–5	40–60
Color neg. slow	100–160	0.60–0.80	4–6	30–70
Color neg. fast	400	0.65–0.80	5–7	25–40
Color neg. very fast	1000–1600	0.80	8–11	25–35
Color rev. slide very slow	25–50	1.8–2.3	9–10	30–40
Color rev. slide slow	100	2.0–3.0	10–13	25–30
Color rev. slide fast	400	2.0–2.4	15–20	20
Color rev. slide very fast	800–1600	2.2–2.8	22–28	16–20
Instant print films ^e		–1.7 ^f	NA	3–4
—Black and white	3000	–1.6	NA	2.5
—Color	100			

^aData are as of early 1993 and were obtained from publications of the manufacturers listed in Sec. 29.21. They are presented as published. Note that products are frequently changed or improved.

^bSpeeds are calculated in different ways for various classes of product. The values given are suitable for use with standard exposure meters.

^cValues represent 1000 \times the standard deviation of the diffuse density, measured at an average density of 1.0 using a 48- μ m circular aperture. The exact granularity of a print depends on the characteristics of the print material and the printer as well as the granularity of the negative.

^dValues show the spatial frequency at which the modulation transfer function is 50 percent.

^eMTF values apply to the final print.

^fNegative sign arises from the definition of gamma for the case of a positive image.

lumped together and termed “scattering,” and have been treated by Gasper and dePalma.³⁰ Likewise, the concept of the optical transfer function, or the Fourier transform of the LSF, is basically the same in optics and photography. However, three important differences should be noted for the photographic case. (1) The emulsion is isotropic, so that the PSF and LSF are always symmetrical, and the complex OTF reduces to the modulation transfer function (MTF) only. (2) Unlike lenses, photographic layers are stationary, but are generally nonlinear. Therefore, all data and calculations must be in terms of *exposure* or allied quantities. When the calculations are complete, the results are converted to density via the D-log H curve. (3) The presence of edge effects tends to raise the MTF curve, so that for low frequencies the measured response values are often found to be greater than 100 percent. The subject is treated in detail by Dainty and Shaw.³¹ Typical MTF curves for a film, showing the overshoot, are given in Fig. 11. Data for MTF curves of various types of films are also given in Table 1. The value given shows the spatial frequency for which the transfer factor drops to 50 percent.

The chain relating the MTF curve to image sharpness is the same as in optics: a high MTF curve transforms to a narrow spread function, and this in turn indicates an abrupt transition of exposure—and therefore density—across the edge. Thus MTF response values greater than unity, although mathematically anomalous, indicate improved image sharpness, and this is found to be the case in practice. As a matter of fact, edge effects are often introduced deliberately to improve sharpness. This is done either by adding suitable compounds to the coating itself, or by adjusting the developer formulation. The process is similar to the electronic “crispensing” often used in television.

An index of sharpness can be computed from the MTF data by a procedure first suggested by Crane and later modified by Gendron.³² These workers were interested mainly in films, but they recognized that the film is one component of a system; e.g., a color slide system involves a camera lens, the film (and process), a projector lens, the screen, and the observer’s eye. The MTF of the system is then the cascaded MTFs of these components. Gendron suggested that the area under the cascaded MTF be taken as the stimulus that produced the sensation of sharpness. A formula was proposed

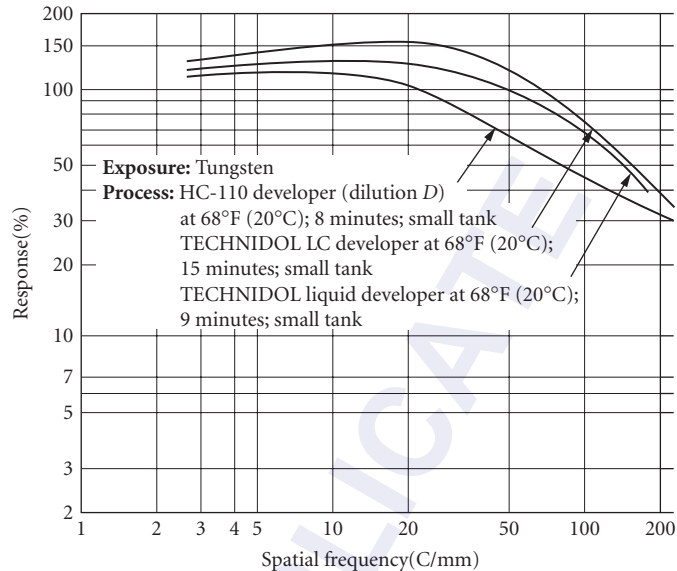


FIGURE 11 Modulation transfer functions of Kodak Technical Pan Film for three conditions of development. Note that the response at low frequencies exceeds 100 percent. (Reprinted courtesy Eastman Kodak Co. © Eastman Kodak Co.)

that produced a sharpness index scaled to 100. This index, now called CMT-Acutance (to distinguish it from the original Higgins-Jones acutance) has been found to correlate well with subjective data, and is used in the industry. Note that the treatment above has been simplified for the sake of brevity; details are available in Gendron's paper.

29.16 GRAININESS

The granular nature of the photographic image is one of its most significant characteristics. It may appear to the observer as an unpleasant roughness in what should be uniform areas, but it may also interfere with extracting information from the image. In the former case it is an aesthetic problem; in the latter case the structure is equivalent to noise in a communications channel. In either case it is desirable to measure the phenomenon objectively and in engineering terms.

The procedure now used for making these objective measurements was proposed by Selwyn in 1935.³³ He postulated that if the density of a uniformly exposed and processed layer were measured at many places using a suitable sampling aperture, the population of density values so obtained would be approximately gaussianly distributed around the mean. This being so, the variability for a given mean density is completely described by the standard deviation of the values. This quantity is termed the *rms-granularity*, $\sigma(D)$, and has indeed been shown to correlate with the subjective graininess.³⁴ It might be noted that calculating the standard deviation of the density is mathematically improper, since the underlying transmittance values are being multiplied. To avoid this problem the rms-granularity is formally defined by

$$\sigma(D) = \frac{0.434\sigma(T)}{\bar{T}} \quad (12)$$

where T is transmittance. However, it can be shown that when $\sigma(T)$ is small compared to \bar{T} , the error involved in calculating directly in density is small, and this is often done in practice.

Selwyn also showed that the measured value of $\sigma(D)$ depended upon A , the area of the sampling aperture. The product $\sigma(D)A^{1/2}$ may be termed the *Selwyn coefficient* \mathcal{G} ; for black-and-white films exposed to light, Selwyn showed that it should be constant, and this relation is called "Selwyn's law." Thus $\sigma(D)A^{1/2}$ is a measure of sample graininess no matter what the size of the sampling aperture. Unfortunately, the Selwyn coefficient does not remain constant with changes of aperture size for very important classes of samples. Selwyn's law may fail for prints and enlargements, black-and-white or color, for many color materials even if not enlarged, and also for radiographs, especially screened radiographs. For such materials $\sigma(D)$ still increases as A decreases, but not at a rate sufficient to keep the Selwyn coefficient constant, and it (the coefficient) is therefore not useful as an objective measure of graininess.

Stultz and Zweig³⁵ found that they obtained good correlation between $\sigma(D)$ and the sensation of graininess when the sampling aperture was selected in accordance with the rule

$$d(\mu) \cdot M_\theta \approx 515 \quad (13)$$

where $d(\mu)$ is the diameter of the sampling aperture in μm , and M_θ is the angular magnification³⁶ at which the photograph is seen by the viewer. M_θ is readily calculated from the relation

$$M_\theta = \frac{m}{4V} \quad (13a)$$

where m is the ordinary lateral magnification between the film image and the image presented to the viewer, and V is the viewing distance in meters.

An American standard³⁷ exists for the measurement of rms-granularity. This standard specifies that samples be scanned with a 48- μm -diameter aperture; for such an aperture, rms-granularity values for commercial films range from about 0.003 to 0.050 at an average density of 1.0. In practice, these values are often multiplied $\times 1000$ to eliminate the decimals. It is worth noting that, experimentally, the measurement of rms-granularity is subject to many sources of error, such as sample artifacts. The standard discusses sources of error and procedures for minimizing them, and is recommended reading for those who must measure granularity.

While in practice rms-granularity serves well as an objective correlate of graininess, the situation is complicated by the fact that there are two broadly different types of granular pattern. Silver grains are small, opaque, and in nearly all cases are situated randomly and independently in the coating. The granular structure in an enlargement, however, is composed of clusters of print-stock grains that reproduce the exposure pattern coming from the enlarged negative grains. This type of granular pattern tends to be large and soft-edged compared to the pattern arising from the primary grains. The patterns found in such samples as color films and screened radiographs are generally similar to those found in enlargements. Microdensitometer traces of these two structures are illustrated in Fig. 12. A little thought will show that the two patterns shown in the figure might have the same mean and standard deviation, and yet the two patterns look entirely different. When different types of patterns are involved, the rms-granularity above is not a sufficient descriptor. The work by Bartleson which showed the correlation between graininess and rms-granularity was done with color negative films having similar granular structures.

As discussed by Dainty and Shaw,³⁸ further objective analysis of granular patterns may be carried out in terms of the *autocorrelation function*:

$$\phi(\tau) = \lim_{x \rightarrow \infty} \frac{1}{2x} \int_{-x}^{+x} \delta(x) \delta(x + \tau) dx \quad (14)$$

where $\delta(x) = D_x - \bar{D}$
 D_x = the density reading at point x on the sample
 \bar{D} = mean density
 $\delta(x + \tau) = D_{x+\tau} - \bar{D}$
 τ = a small increment of distance

Note carefully that for the sake of simplicity it has been assumed that the sample is scanned by a very long, narrow slit, so that the autocorrelation function reduces to a one-dimensional function.

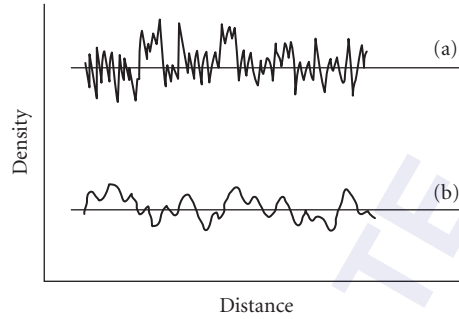


FIGURE 12 Microdensitometer traces of (a) primary silver grains and (b) granular structure of a print.

If the sample is scanned by a small circular aperture it will be two-dimensional. The use of a slit is common in practice. For the case where $\tau = 0$, we have

$$\phi(0) = \lim_{x \rightarrow \infty} \frac{1}{2x} \int_{-x}^{+x} [\delta(x)]^2 dx = \sigma^2 \quad (15)$$

Let $\phi(\tau)$ be calculated for several different values of τ , including values smaller than the slit width. Since when $\tau < w$ the slit will contain some of the same grains at points x and $x + \tau$, correlation is observed. If there is no spatial correlation in the actual sample, $\phi(\tau) \rightarrow 0$ for $\tau > w$. Positive values of $\phi(\tau)$ for $\tau > w$ are an indication of spatial correlation in the sample; that is, large-scale grain.

In practice, it is convenient to carry out the analysis of large-scale patterns in the spatial frequency domain. The Wiener-Khinchin theorem states that what is called the “Wiener spectrum” or *power spectrum of the granularity* distribution is the Fourier transform of the autocorrelation function; that is,

$$\text{WS}(\nu) = \int_{-\infty}^{\infty} \phi(\tau) e^{i2\pi\nu\tau} d\tau \quad (16)$$

and also

$$\phi(\tau) = \int_{-\infty}^{\infty} \text{WS}(\nu) e^{-i2\pi\nu\tau} d\nu \quad (16a)$$

where ν = spatial frequency.

In practice, an approximation of the Wiener spectrum is usually obtained by a direct Fourier transform of the granularity trace itself according to the expression

$$\text{WS}(\nu) = \lim_{x \rightarrow \infty} \frac{1}{2x} \overline{\left| \int_{-x}^x \delta(x) e^{i2\pi\nu x} dx \right|^2} \quad (17)$$

where $\delta(x)$ has the same meaning as before, and the horizontal bar indicates that the average value of several different runs should be taken in order to provide a reasonable value for the approximation.

Returning to the definition of the autocorrelation function [Eq. (14)] it can be seen that the measured values of $\delta(x)$ and $\delta(x + \tau)$ are each actually the true values convolved with the slit response. By the convolution theorem of Fourier transforms, this means that in frequency space

$$\text{WS}(\nu)_{\text{measured}} = \text{WS}(\nu)_{\text{true}} xT(\nu)^2 \quad (18)$$

where $T(\nu)$ is the modulation transfer factor of the measuring system at spatial frequency ν . If $T(\nu)^2$ is divided out of $WS(\nu)_{\text{measured}}$ the underlying net value $WS(\nu)$ is obtained. When this is done for black-and-white film samples exposed to light, the underlying spectrum is found to be flat. For the kinds of samples for which Selwyn's law fails, on the other hand, the net Wiener spectrum is found to contain excess low-frequency power and lack high frequencies. Finally, by the properties of Fourier transforms, the area under the Wiener spectrum curve is the zero-value of the autocorrelation function. But by Eq. (15) this zero-value equals σ^2 for the sample, so that the area under the WS curve gives the rms-granularity of the sample. This can also be seen intuitively, since the value of the WS at special frequency ν is simply the noise power of σ^2 in a spatial frequency band $\nu \pm \Delta\nu/2$.

Doerner³⁹ has shown that the rms-granularity of a negative can be tracked through a printer system in terms of the Wiener spectrum of the negative and the modulation transfer function of the printer system (which includes the MTF of the print stock itself). Doerner's expression is

$$WS(\nu)_p = WS(\nu)_n \gamma_p^2 A(\nu)_{\text{sys}}^2 + WS(\nu)_{\text{ps}} \quad (19)$$

where n and p indicate negative and print respectively, "ps" the print stock itself, and the other symbols have their previous meanings. The granularity/graininess of a print thus depends on the contrast of the print stock and the spatial frequency response of the printer system as well as on the graininess of the negative.

29.17 SHARPNESS AND GRAININESS CONSIDERED TOGETHER

In the foregoing discussion, we have considered the sharpness and graininess aspects of the picture separately. But real photographs frequently suffer from less-than-optimum sharpness and graininess *both*. Bartleson⁴⁰ has studied the subjective quality of such photographs. He concluded that quality was not a linear combination of sharpness and graininess. Instead, "... quality tends to be determined primarily by the poorer of the two attributes. . . . If graininess is high, the print will likely be low in quality regardless of how sharp it may be or, conversely, if sharpness is low, so also will quality be low regardless of how grainy the print appears. . . ."

Bartleson's results are of interest in assessing the quality of electronic images, since, at least at the time of writing, such images exhibit low graininess, but also low sharpness. On the basis of Bartleson's work, such images would be judged to be of low subjective quality. In a comparison of electronic and film imagery published in 1990, Ikenoue and Tabei⁴¹ rated the quality of the former as poor because of the sharpness level.

29.18 SIGNAL-TO-NOISE RATIO AND DETECTIVE QUANTUM EFFICIENCY

Since the information in a photograph is normally carried by the density variation, we may usefully define the output signal-to-noise ratio of the photography by

$$S/N_{\text{out}} = \frac{\Delta\bar{D}}{\sigma} = \frac{\gamma}{\sigma} \cdot \Delta\log H \quad (20)$$

where $\Delta\bar{D}$ is the mean density difference between an element to be detected and its surround, and σ is the rms-granularity of the surround. By Selwyn's law, $\sigma(D)$ varies as the area of the sampling

aperture changes; it is convenient here to take the sampling aperture area A as equal to the area of the image element. Furthermore, for ΔH sufficiently small we may write

$$S/N_{\text{out}} = 0.434 \frac{\gamma}{\sigma} \cdot \frac{\Delta H}{H} \quad (21)$$

Note that γ/σ is a property of the film; it is termed the *detectivity*. $\Delta H/H$, on the other hand, is a property of the object; in fact, it is the object contrast. Practical tests have shown that S/N_{out} should be 4 to 5 if the element is to be detected against its surround, and 8 to 10 if it is to be recognized.⁴²

In 1946, Albert Rose of RCA published a paper⁴³ in which he discussed the performances of the TV pickup tube, the photographic layer, and the human eye on a unified basis. His approach was to compare their performances with that of an "ideal device," that is, a radiation detector whose performance was limited only by the quantum nature of the incoming radiation. Such a perfect detector would report the arrival of every incoming quantum, and add no noise to the signal.

In 1958, R. Clark Jones of Polaroid expanded Rose's work with specific application to photographic layers.⁴⁴ Jones proposed the term *detective quantum efficiency* (DQE) for Rose's performance indicator, and defined it by the expression

$$\text{DQE} = \left[\frac{S/N_{\text{out}}}{S/N_{\text{ideal}}} \right]^2 = \left[\frac{S/N_{\text{out}}}{S/N_{\text{max}}} \right]^2 = \left[\frac{S/N_{\text{out}}}{S/N_{\text{in}}} \right]^2 \quad (22)$$

where S/N_{out} is the signal-to-noise ratio produced by the actual device, and S/N_{ideal} is the ratio that would be produced by the ideal device, given the same input. By the definition of the ideal device, $S/N_{\text{ideal}} = S/N_{\text{max}} = S/N_{\text{in}}$, where S/N_{in} is the signal-to-noise ratio in the input, and is due to the quantum nature of the input. The ratio is squared to make the DQE compatible with various other concepts.

It is easy to derive an expression for the S/N ratio of the input; when this is combined with Eq. (21) the result is

$$\text{DQE} = \left[\frac{0.434\gamma}{\mathcal{G}} \right]^2 \cdot \frac{1}{q} \quad (23)$$

where \mathcal{G} is the Selwyn coefficient and q is the average exposure received by the image in *quanta per unit area*. Since $(1/q)$ is the radiometric speed of the film, it can be seen that the DQE is a performance parameter that combines the gain (gamma), noise, and speed of the layer. As written, the DQE does not involve the sharpness aspect, but this can be included also.⁴⁵

Since Eq. (23) involves standard photographic parameters, it is readily evaluated for a given material. When this is done, it is found that the DQE of typical materials is on the order of 1 to 3 percent, peaking sharply at low densities in the case of black-and-white films. Note that DQE can be calculated for any sensor for which S/N_{out} can be derived. It is interesting to compare the 1 to 3 percent values given above with those of other sensors. Thus, for example, Jones gives a value of 1 percent for the human eye, and 6 percent for an image orthicon tube. On the other hand, a suitable photographic layer recording electrons may approach 100 percent DQE, and a value of 30 percent has been reported for a screened x-ray film.⁴⁶ DQE has also served as a useful approach to considering silver halide mechanisms; see, for example, a paper by Bird, Jones, and Ames that appeared in *Applied Optics* in 1969.⁴⁷

Another figure of merit allied to DQE is the *noise equivalent quanta* (NEQ) which is defined as the number of quanta that a perfect detector would need to produce a record having the same S/N ratio as the system under consideration. It can be shown that

$$q' = \text{DQE} \times q \quad (24)$$

where q is again the number of quanta/unit area used by the real system, and q' the NEQ of a unit area of image.

29.19 RESOLVING POWER

The basic procedure used to measure photographic resolving power follows that used in optics. An American Standard exists.⁴⁸ The standard provides for a suitable test object to be reduced optically onto the material to be tested. (Strictly speaking, the test thus determines the resolution of the lens-film combination, but the resolution capability of the lenses specified is high compared to that of the film.) The developed image is then studied in a microscope to determine the highest spatial frequency in which the observer is “reasonably confident” that the structure of the test pattern can still be detected. Thus, as in optics, the “last resolved” image is a threshold image. However, unlike the optical case, the limit is set not only by the progressive decrease in image modulation as spatial frequency increases, but also by the granular nature of the image.

If an exposure series of the test pattern is made, it is found that the spatial frequency of the limiting pattern goes through a maximum. It is customary to report the spatial frequency limit for the optimum exposure as the resolving power of the film.

Photographic resolving power is now not much measured, but it retains some interest as an example of a signal-to-noise ratio phenomenon. Consider the modulation in the various triplet images in the pattern. By definition,

$$M = \frac{H_{\max} - H_{\min}}{H_{\max} + H_{\min}} = \frac{\Delta H}{2\bar{H}} \quad (25)$$

and from Eq. (21),

$$\Delta D = 0.868 \gamma M \quad (26)$$

so that ΔD within the pattern varies as the modulation, which in turn decreases as the spatial frequency increases. Furthermore, the area of each of the triplets in the pattern (assuming the ISO pattern configuration) = $2.5^2 \lambda^2 = 2.5^2 / \nu^2$. Assuming that Selwyn's law holds, it follows that $\sigma = \mathcal{G}/A^{1/2} = 0.4 \mathcal{G}\nu = C\nu$, where we have lumped the constants. As the spatial frequency increases, the S/N of the triplet decreases because ΔD decreases and the effective rms-granularity increases. The resolving power limit comes at the spatial frequency where the S/N ratio drops to the limit required for resolution. Such a system has been analyzed by Schade.⁴⁹

29.20 INFORMATION CAPACITY

The information capacity, or number of bits per unit area that can be stored on a photographic layer, depends on the size of the point spread function, which determines the smallest element that can be recorded, and on the granularity, which determines the number of gray levels that can be reliably distinguished. The exact capacity level for a given material depends on the acceptable error rate; for one set of fairly stringent conditions, Altman and Zweig⁵⁰ reported levels up to 160×10^6 bits/cm². Jones⁵¹ has published an expression giving the information capacity of films as

$$IC = \frac{1}{2} \iint_0^\infty \log_2 \left(1 + \frac{S(\nu_x, \nu_y)}{N(\nu_x, \nu_y)} \right) d\nu_x d\nu_y \quad (27)$$

where S and N are the spectral distributions of power in the system's signal and noise. For a given spatial frequency, the signal power is given by

$$S = WS_i(\nu) \text{MTF}^2(\nu) \quad (28)$$

where $WS_i(\nu)$ is the value of the Wiener spectrum of the input at ν and $\text{MTF}(\nu)$ is the value of the film's MTF at that frequency. The information capacity thus depends on the frequency response and noise of the system, as would be expected. The matter is discussed by Dainty and Shaw.⁵²

29.21 LIST OF PHOTOGRAPHIC MANUFACTURERS

Agfa Photo Division
Agfa Corporation
100 Challenger Road
Ridgefield, NJ 07660

Ilford Photo Corporation
70 West Century Boulevard
Paramus, NJ 07653

E. I. duPont de Nemours and Co.
Imaging Systems Department
666 Driving Park Avenue
Rochester, NY 14613

3M Company
Photo Color Systems Division
3M Center
St. Paul, MN 55144-1000

Eastman Kodak Co.
343 State Street
Rochester, NY 14650
Tel. 1-800-242-2424 for product info.

Polaroid Corporation
784 Memorial Drive
Cambridge, MA 02139
Tel. 1-800-255-1618 for product info.

Fuji Photo Film USA
555 Taxter Road
Elmsford, NY 10523

29.22 REFERENCES

1. W. Thomas (ed.), *SPSE Handbook of Photographic Science and Engineering*, Wiley, New York, 1973, sec. 8.
2. B. H. Carroll, G. C. Higgins, and T. H. James, *Introduction to Photographic Theory*, Wiley, New York, 1980, chaps. 6-9.
3. T. H. James (ed.), *The Theory of the Photographic Process*, 4th ed., Macmillan, New York, 1977.
4. J. Sturge, Vivian Walworth, and Alan Shepp (eds.), *Imaging Processes and Materials*, Van Nostrand Reinhold, New York, 1989, chap. 3. (See also Ref. 3, chap. 4.)
5. G. Haist, *Modern Photographic Processing*, vols. I and II, Wiley, New York, 1979.
6. Klaus Hendricks, chap. 20 of Ref. 4.
7. J. H. Altman, F. Grum, and C. N. Nelson, *Phot. Sci. Eng.* **17**:513 (1973).
8. Reference 3, chaps. IV-VII.
9. ANSI/ISO, May 2, 1991. (See also ANSI/ISO, May 1, 1984, which sets forth standardized terminology.)
10. R. W. G. Hunt, *The Reproduction of Colour*, 4th ed., Fountain Press, Tolworth, England, 1987, pp. 247-257.
11. *American National Standard*, ANSI PH2.2, 1984, R1989, PH2.17 1985.
12. C. J. Niederpruem, C. N. Nelson, and J. A. C. Yule, *Phot. Sci. Eng.* **10**:35 (1966).
13. *American National Standard*, ANSI PH3.49, 1971, R1987.
14. H. N. Todd and R. D. Zakia, *Phot. Sci. Eng.* **8**:249 (1964). (See also publications of the American National Standards Institute.)
15. M. Hercher and B. J. Ruff, *J. Opt. Soc. Am.* **57**:103 (1967).
16. W. F. Berg, *Proc. Roy. Soc. of London*, ser. A, **174**:5599 (1940).
17. Reference 2, p. 141.
18. *Scientific Imaging with Kodak Films and Plates*, Publication P315, Eastman Kodak Co., Rochester, N.Y., 1987, p. 63.
19. J. Kapecki and J. Rodgers, *Kirk-Othmer Encyclopedia of Chemical Technology*, 4th ed., vol. 6, Wiley, New York, 1993, p. 965.
20. S. H. Liggero, K. J. McCarthy, and J. A. Stella, *J. Imaging Technology* **10**:1 (1984).

21. Reference 5, chap. 7.
22. E. H. Land, *Phot. Sci. Eng.* **16**:247 (1972).
23. Reference 10, chap. 11 et seq.
24. See, for example, F. R. Clapper in Ref. 3, sec. II, chap. 19.
25. R. M. Evans, *Eye, Film, and Camera in Color Photography*, Wiley, New York, 1959, p. 180.
26. B. J. Thompson, *Progress in Optics VII*, E. Wolf (ed.), North-Holland Publishing Co., Amsterdam, 1969; R. E. Swing, *J. Opt. Soc. Am.* **62**:199 (1972); Dainty and Shaw, Ref. 27, chap. 9.
27. P. Kowalski, *Applied Photographic Theory*, Wiley, New York, p. 77 et seq.
28. C. N. Nelson, Ref. 3, chap. 19.
29. G. C. Higgins and L. A. Jones, *J. Soc. of Mot. Pict. & TV Engrs.* **58**:277 (1952).
30. J. Gasper and J. J. dePalma, Ref. 3, chap. 20.
31. J. C. Dainty and R. Shaw, *Image Science*, Academic Press, London, 1974, chaps. 6 and 7. (See also J. C. Dainty, *Optica Ada* **18**:795, 1971.)
32. R. M. Gendron, *J. Soc. of Mot. Pict. & TV Engrs.* **82**:1009 (1973).
33. E. W. H. Selwyn, *Photography Journal* **75**:571 (1935).
34. C. J. Bartleson, *Photography Journal* **33**:117 (1985).
35. K. F. Stultz and H. J. Zweig, *J. Opt. Soc. Am.* **49**:693 (1959).
36. F. W. Sears, *Optics*, Addison-Wesley, Reading, Mass., 1949, p. 159.
37. *American National Standard*, ANSI PH2.40, 1985, R1991.
38. Reference 31, chap. 8.
39. E. C. Doerner, *J. Opt. Soc. Am.* **52**:669 (1962).
40. C. J. Bartleson, *J. Phot. Sci.* **30**:33 (1982).
41. S. Ikenoue and M. Tabei, *J. Imaging Sci.* **34**:187 (1990).
42. Reference 2, p. 335.
43. A. Rose, *J. Mot. Pict. & TV Engrs.* **47**:273 (1946).
44. R. C. Jones, *Phot. Sci. Eng.* **2**:57 (1958).
45. Reference 31, chap. 8, p. 311 et seq.
46. Reference 4, table 3.2, p. 73.
47. G. R. Bird, R. C. Jones, and A. E. Ames, *Appl. Opt.* **8**:2389 (1969).
48. *American National Standard*, ANSI PH2.33, 1983, R1990.
49. O. H. Schade, *J. Soc. Mot. Pict. & TV Engrs.* **73**:81 (1964).
50. J. H. Altman and H. J. Zweig, *Phot. Sci. Eng.* **7**:173 (1963).
51. R. Clark Jones, *J. Opt. Soc. Am.* **51**:1159 (1961).
52. Reference 31, chap. 10.

John D. Baloga

Imaging Materials and Media
Eastman Kodak Company
Rochester, New York

30.1 INTRODUCTION

Photographic materials are optical devices. Their fabrication and technical understanding encompass the field of optics in the broadest sense. Light propagation, absorption, scattering, and reflection must be controlled and used efficiently within thin multilayer coatings containing tiny light-detecting silver halide crystals and chemistry. Many subdisciplines within classical optics, solid-state theory, and photochemistry describe these processes.

In Chap. 20 of Ref. 1, Altman sweeps broadly through the basic concepts and properties of photographic materials. His brief, high-level descriptions are still generally valid for today's photographic products, and the reader will profit by reviewing that summary. This chapter focuses more sharply on four fundamental photographic technologies intimately related to the broad field of optics, then gives an overview of photographic materials available today.

Section 30.2 discusses the optical properties of multilayer color photographic films and papers. This section outlines the basic structure of the device and describes the principal function of each layer. Light absorption, reflection, scattering, and diffraction properties are discussed in the context of providing minimum optical distortion to the final image.

Silver halide light detection crystals are solid-state devices described by quantum solid-state photophysics and chemistry in addition to crystallography. These devices absorb light to create an electron-hole pair. Delocalization of the electron in the conduction band of the crystal, trapping, reduction of interstitial silver ions, nucleation and growth of clusters of silver atoms, and reversible regression of these phenomena must be controlled and optimized to produce the most sensitive light detectors that also possess long-term stability.

Section 30.3 describes the basic photophysics of silver halides according to our best understanding today. It outlines the solid-state properties most important to image detection and amplification. Surface chemical treatment and internal doping are discussed in the context of providing silver halide emulsions having the highest sensitivity to light.

A color photographic image is formed by high-extinction and high-saturation dyes structurally designed by chemists to optimize their color, permanence, and other useful characteristics. Their properties in both the ground state and photoexcited states are important to color and stability.

Section 30.4 briefly outlines the photochemistry of photographic image dyes. Excited-state properties are described that are important to the photostability of these dyes for image permanence.

Color science guides us to optimize all factors in photographic materials conspiring to render an accurate and pleasing image. These include spectral discrimination during the light detection phase and color correction to compensate for imperfect spectral detection and imperfect image dyes.

Section 30.5 sketches the photophysics and color science of photographic spectral sensitizers with an aim toward describing how modern photographic films sense the world in color nearly as the human eye sees it.

Today there exists a large diversity of photographic films. In addition to different manufacturers' brands, films differ by speed, type, color attributes, format, and a multitude of other factors. This can be confusing. And what are the differences between a consumer film bought in a drugstore and a more expensive professional film?

Section 30.6 gives an overview of the different types of films available today. Differences between high- and low-speed films are described with an understanding of the origins of these differences. Consumer films are compared to professional films and some of the special needs of each type of customer are described. Some general guidelines are offered for film choices among color reversal films, black-and-white films, and color negative films.

In addition to Chap. 20 in Ref. 1, several other texts contain useful information about photographic materials. Kapecki and Rodgers² offer a highly lucid and digestible sketch of basic color photographic film technology. Besides Chap. 20 in Ref. 1, this is a good place for the technically astute but nonpractitioner to gain a high level understanding of basic color photography. The technically detailed treatise by James³ is perhaps the best single comprehensive source containing the science and technology of the photographic system outside photographic manufacturers' internal proprietary libraries. Hunt⁴ provides a good comprehensive treatise on all aspects of color science and additionally gives a good technical review of color photographic materials. Chapter 6 in Ref. 1 contains a useful overview of the theory of light-scattering by particles, Chap. 9 contains a good overview of optical properties of solids, and Chap. 26 on colorimetry provides a good introduction to the basic concepts in that field.

30.2 THE OPTICS OF PHOTOGRAPHIC FILMS AND PAPERS

Introduction

Color photographic materials incorporate design factors that optimize their performance across all features deemed important to customers who use the product. These materials are complex in composition and structure. Color films typically contain over 12 distinct optically and chemically interacting layers. Some of the newest professional color reversal films contain up to 20 distinct layers.* Each layer contributes a unique and important function to the film's final performance. Careful design and arrangement of the film's various layers ultimately determine how well the film satisfies the user's needs.

One very important customer performance feature is film *acutance*, a measure of the film's ability to clearly record small detail and render sharp edges. Because images recorded on color film are frequently enlarged, image structure attributes such as acutance are very important. Magnifications such as those used to place the image on a printed page challenge the film's ability to clearly record fine detail.

The ability of a photographic material to record fine detail and sharp edges is controlled by two factors. The first includes light scatter, diffraction, and reflections during the exposure step. These are collectively called *optical* effects. This factor dominates the film's ability to record fine detail and sharp edges and is therefore critical to the film's design.

*Fujichrome Velvia 50 RVP professional film.

Chemical adjacency effects collectively known as *Eberhard effects* (see Chap. 21 in Ref. 3) are the second factor. In today's color and black-and-white film products these are present during film development and are caused by accumulating development by-products such as iodide ions or inhibitors released by development inhibitor-releasing (DIR) chemicals that restrain further silver development in exposed regions of the image, leading to a chemical unsharp mask effect (pp. 274 and 365 of Ref. 4). These chemical signals also give rise to the film's interlayer interimage effects (HE; see p. 278 of Ref. 4) used for color correction. Film sharpness and HE are strongly related.

This section describes multilayer factors that contribute to light scatter, reflection, diffraction, and absorption in photographic materials. Factors that influence the nonoptical part of film acutance, such as Eberhard effects, were briefly reviewed by Altman.¹ More information about these processes can be found in references cited therein.

Structure of Color Films

Color film structures contain image-recording layers, interlayers, and protective overcoat layers. These layers suspend emulsions and chemistry in a hardened gelatin binder coated over a polyacrylate or polyester support. Figure 1 shows the principal layers in a color multilayer film structure.

The overcoat and ultraviolet (UV) protective layers contain lubricants; plastic beads 1 to 5 μm in size called *matte* to impart surface roughness that prevents sticking when the film is stored in roll form; UV-light-absorbing materials; and other ingredients that improve the film's handling characteristics and protect the underlying light-sensitive layers from damage during use and from exposure to invisible UV light. Ultraviolet light is harmful for two reasons: it will expose silver halide emulsions, thereby rendering an image from light invisible to humans, and it promotes photodecomposition of image dyes, leading to dye fade over time (p. 977 of Ref. 2). Visible light exposure must first pass through these overcoats, but they typically do little harm to acutance as they contain nothing that seriously scatters visible light.

The blue-light-sensitive yellow dye imaging layers appear next. Silver halides, with the exception of AgCl—used primarily in color papers—have an intrinsic blue sensitivity even when spectrally sensitized to green or red light. Putting the blue-light-sensitive layers on top avoids incorrect

Overcoat
UV protective layer
Fast yellow layer
Slow yellow layer
Yellow filter layer
Barrier layer
Fast magenta layer
Mid magenta layer
Slow magenta layer
Optional magenta filter layer
Barrier layer
Fast cyan layer
Mid cyan Layer
Slow cyan layer
AHU layer
Plastic support
Optional pelloid AHU layer

FIGURE 1 Simplified diagram showing the key layers in a color photographic film (not drawn to scale).

exposure leading to color contamination because a blue-light-absorbent yellow filter layer, located beneath the blue-sensitive imaging layers, absorbs all blue light that has passed completely through the blue-sensitive layers.

A collection of layers adjacent to one another that contain silver halide sensitized to a common spectral band is called a *color record*. Color records in films are always split into two or three separate layers. This arrangement puts the highest-sensitivity emulsion in the upper (fast) layer, where it gets maximum light exposure for photographic speed, and places low-sensitivity emulsions into the lower layer (or layers), where they provide exposure latitude and contrast control.

One or more interlayers separate the yellow record from the green-light-sensitive magenta record. These upper interlayers filter blue light to avoid blue light exposure color contaminating the red- and green-sensitive emulsion layers below. This behavior is called *punch through* in the trade. These interlayers also contain oxidized developer scavengers to eliminate image dye contamination between color records caused by oxidized color developer formed in one color record diffusing into a different color record to form dye.

The traditional yellow filter material is Carey Lea silver (CLS), a finely dispersed colloidal form of metallic silver, which removes most of the blue light at wavelengths less than 500 nm. The light absorption characteristics of this finely dispersed metallic silver are accounted for by Mie theory.⁵ Unfortunately this material also filters out some green and red light, thus requiring additional sensitivity from the emulsions below. Bleaching and fixing steps in the film's normal processing remove the yellow CLS from the developed image.

Some films incorporate a yellow filter dye in place of CLS. Early examples contained yellow dyes attached to a polymer mordant in the interlayer. A *mordant* is a polymer that contains charged sites, usually cationic, that binds an ionized anionic dye by electrostatic forces. These dyes are removed during subsequent processing steps. Although these materials are free from red and green light absorption, it is a challenge to make them high in blue light extinction to avoid excessive chemical loads in the interlayer with consequent thickening.

Most recently, solid-particle yellow filter dyes⁶ are seeing more use in modern films. These solid dyes are sized to minimize light scatter in the visible region of the spectrum and their absorption of blue light is very strong. They can be made with very sharp spectral cuts and are exceptionally well suited as photographic filter dyes. In some films solid-particle magenta filter dyes⁷ are also used in the interlayers below the magenta imaging layers. Solid-particle filter dyes are solubilized and removed or chemically bleached colorless during the film's normal development process.

Because the human visual system responds most sensitively to magenta dye density, the green-light-sensitive image recording layers are positioned just below the upper interlayers, as close as practical to the source light exposure side of the film. This minimizes green light spread caused when green light passes through the turbid yellow record emulsions. It gives the maximum practical film acutance to the magenta dye record.

A lower set of interlayers separates the magenta record from the red-light-sensitive cyan dye record. These lower interlayers give the same type of protection against light and chemical contamination as do the upper interlayers. The magenta filter dye is absent in some films because red-light-sensitive emulsions are not as sensitive to green light exposure as to blue light exposure.

Located beneath the cyan record, just above the plastic film support, are antihalation undercoat (AHU) layers. The black absorber contained in this bottom layer absorbs all light that has passed completely through all imaging layers, thereby preventing its reflection off the gel-plastic and plastic-air interfaces back into the imaging layers. These harmful reflections cause a serious type of light spread called *halation*, which is most noticeable as a halo around bright objects in the photographic image.

The opacity required in the AHU is usually obtained by using predeveloped black filamentary silver, called *gray gel*, which is bleached and removed during the normal photographic processing steps. Alternatively, in many motion picture films a layer of finely divided carbon suspended in gelatin (*rem jet*) is coated on the reverse side of the film. When placed in this position the layer is called an *AHU pelloid*. It is physically removed by scrubbing just before the film is developed. The newest motion picture films incorporate a black solid-particle AHU filter dye that is solubilized and removed during normal development of the film and does not require a separate scrubbing step.

The overall thickness of a film plays an important role in minimizing harmful effects from light scatter. Because color films are structured with the yellow record closest to the light exposure source, it is especially important to minimize thickness of all film layers in the yellow record and below it, because light scattered in the yellow record progressively spreads as it passes to lower layers. The cyan record shows the strongest dependence on film thickness because red light passes through both yellow and magenta record emulsions en route to the cyan record. Because both emulsions often contribute to red light scatter, the cyan record suffers a stronger loss in acutance with film thickness.

Structure of Color Papers

The optical properties of photographic paper merit special consideration because these materials are coated in very simple structures on a highly reflective white Baryta-coated paper support. *Baryta* is an efficient diffuse reflector consisting of barium sulfate powder suspended in gelatin that produces isotropically distributed reflected light with little absorption.

Photographic papers generally contain about seven layers. On top are the overcoat and UV protective layers, which serve the same functions as previously described for film.

The imaging layers in color papers contain silver chloride emulsions for fast-acting development, which is important to the industry for rapid throughput and productivity. Generally only one layer is coated per color record, in contrast to films, which typically contain two or three layers per color record. The order of the color records in photographic papers differs from that in films due mainly to properties of the white Baryta reflective layer.

Because color paper is photographically slow, exposure times on the order of seconds are common. Most light from these exposures reflects turbidly off the Baryta layer. The imaging layers getting the sharpest and least turbid image are those closest to the reflective Baryta where light spread is least, not those closest to the top of the multilayer as is the case with film.

In addition, the Baryta layer as coated is not smooth. Its roughness translates into the adjacent layer, causing nonuniformities in that layer. Because the human visual system is most forgiving of physical imperfections in yellow dye, the yellow color record must be placed adjacent to the Baryta to take the brunt of these imperfections.

The magenta color record is placed in the middle of the color paper multilayer, as close to the Baryta layer as possible, since sharpness in the final image is most clearly rendered by magenta dye. This leaves the cyan record nearest to the top of the structure, just below the protective overcoats.

The magenta color record in the middle of the multilayer structure is spaced apart from the other two color records by interlayers containing oxidized developer scavengers to prevent cross-record color contamination, just as in films. However, no filter dyes are needed in color paper because silver chloride emulsions have no native sensitivity to light in the visible spectrum.

Light Scatter by Silver Halide Crystals

Because they consist of tiny particles, silver halide emulsions scatter light. Scattering by cubic and cubo-octahedral emulsions at visible wavelengths is most intense when the emulsion dimension ranges from 0.3 to 0.8 μm , roughly comparable to the wavelengths of visible light. This type of scattering is well characterized by Mie⁸ theory and has been applied to silver halides.⁹ We are often compelled to use emulsions having these dimensions in order to achieve photographic speed.

For photographic emulsions of normal grain size and concentration in gelatin layers of normal thickness, multiple scattering predominates. This problem falls within the realm of radiative transfer theory. Pitts¹⁰ has given a rigorous development of radiative transfer theory to the problem of light scattering and absorption in photographic emulsions. A second approach that has received serious attention is the Monte Carlo technique.¹¹ DePalma and Gasper¹² were able to obtain good agreement between their modulation transfer functions (MTFs) determined by a Monte Carlo calculation and experimentally measured MTFs for a silver halide emulsion layer coated at various thicknesses.

Scattering by yellow record emulsions is especially harmful because red and green light must first pass through this record en route to the magenta and cyan records below. All other things being equal, the amount of light scattered by an emulsion layer increases in proportion to the total amount of silver halide coated in the layer.

Color films are constructed with low-scattering yellow record emulsions whenever possible. The amount of silver halide coated in the yellow record is held to a minimum consistent with the need to achieve the film's target sensitometric scale in yellow dye.

High-dye-yield yellow couplers¹³ having high coupling efficiency are very useful in achieving silver halide mass reductions in the yellow record of color reversal films. These provide high yellow dye densities per unit silver halide mass, thereby reducing the amount of silver halide needed to achieve target sensitometry. Some upper-scale (high-density) increase in granularity often results from using these couplers in a film's fast yellow layer, but the benefits of reduced red and green light scatter overcome this penalty, especially because the human visual system is insensitive to yellow dye granularity.

Tabular emulsion grains offer a way to achieve typical photographic speeds using large-dimension emulsions, typically 1.0 to 3.0 μm in diameter, although smaller- and larger-diameter crystals are sometimes used. These do not scatter light as strongly at high angles from normal incidence with respect to the plane of the film as do cubic or octahedral emulsions having comparable photographic speeds. However, diffraction at the edges and reflections off the crystal faces can become a detrimental factor with these emulsions.

Silver halide tabular crystals orient themselves to lie flat in the plane of the gelatin layer. This happens because shear stress during coating of the liquid layer stretches the layer along a direction parallel to the support and also because the water-swollen thick layer compresses flat against the support after it dries by a ratio of roughly 20:1.

Reflections can become especially harmful with tabular emulsion morphologies since light reflecting from the upper and lower faces of the crystal interferes, leading to resonances in reflected light. The most strongly reflected wavelengths depend on the thickness of the tabular crystal. The thickness at which there is a maximum reflectance occurs at fractional multiples of the wavelength¹⁴ given by:

$$t = \frac{\left(m + \frac{1}{2}\right)\lambda}{2n}$$

where t is the thickness of the tabular crystal, λ is the wavelength of light, n is the refractive index of the crystal, and m is an integer.

In an extreme case, tabular emulsions act like partial mirrors reflecting light from grain to grain over a substantial distance from its point of origin. This is called *light piping* by analogy with light traveling through an optical fiber (see Fig. 2). It is especially serious in the cyan record, where red light often enters at angles with respect to perpendicular incidence caused by scattering in upper layers.

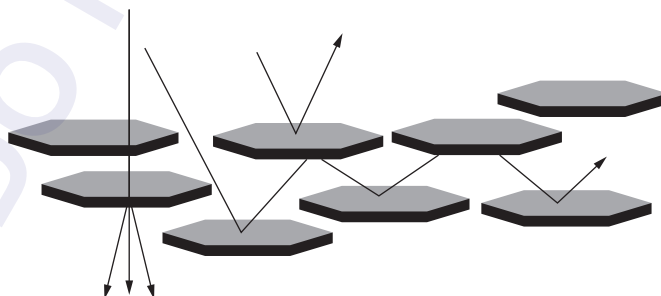


FIGURE 2 Light undergoes reflection, diffraction, and multiple reflections that may lead to light piping in film layers containing tabular silver halide crystals.

Another cause of light piping is the presence of small-grain emulsion contaminants within a tabular grain population. A highly turbid contaminant emulsion, even if present in small amounts, can scatter light at sufficient angles to induce severe light piping through the tabular grain matrix. Modern emulsion science takes great pains to produce tabular emulsions free from highly scattering contaminants.

Gasper¹⁵ treated the problem of scattering from thin tabular silver halide crystals in a uniform gelatin layer. For very thin and large-diameter tabular silver bromide crystals (thickness $<0.06\ \mu\text{m}$, diameter $>1.0\ \mu\text{m}$) light at normal incidence is scattered almost equally in the forward and backward hemispheres. As the grain thickness increases there appears increasing bias for forward scatter in a narrow cone. The efficiencies for scatter and absorption were found to be independent of grain diameter for crystal diameters larger than $\sim 1.0\ \mu\text{m}$. For such crystals, the backscatter and absorption efficiencies are approximately equal to the specular reflectance and absorption of a semiinfinite slab of the same material and thickness.

But the forward scattering efficiency does not approximate the specular transmittance of the semiinfinite slab as a result of interference between the directly scattered forward field and the mutually coherent unscattered field, a process analogous to diffraction that does not occur for the semiinfinite slab with no edge. The specific turbidity depends on grain thickness, but not diameter for diameter $>1.0\ \mu\text{m}$.

A light-absorbing dye is sometimes added to a film to reduce the mean free path of scattered and reflected light leading to an acutance improvement. This improvement happens at the expense of photographic speed, since light absorbed by the dye is not available to expose the silver halide emulsions. However, the trade-off is sometimes favorable if serious light piping is encountered.

Light-absorbing interlayers between color records are sometimes used in color films to help eliminate the harmful effects of reflected light. For example, some films have a magenta filter dye interlayer between the magenta and cyan records. In addition to its usefulness in reducing light punchthrough (green light passing through the magenta record to expose the red-sensitized layers), this filter layer helps eliminate harmful reflections of green light off cyan record tabular emulsions, which bounce back into the magenta record. These reflections, although potentially useful for improving photographic green speed, can be harmful to magenta record acutance.

30.3 THE PHOTOPHYSICS OF SILVER HALIDE LIGHT DETECTORS

Chapter 20 in Ref. 1 gave a very brief sketch of the general characteristics of silver halide crystals used for light detection and amplification in photographic materials. These crystals, commonly termed *emulsions*, are randomly dispersed in gelatin binder layers in photographic films and papers. For photographic applications the most commonly used halide types are AgCl, AgBr, and mixed halide solid solutions of AgCl_xI_y and AgBr_xI_y . In addition, pure phases of AgCl and AgI are sometimes grown epitaxially on AgBr crystals to impart special sensitivity and development properties.

Upon exposure to light, silver halide crystals form a latent image (LI) composed of clusters of three to hundreds of photoreduced silver atoms either within the interior or most usefully on the surface of the crystal where access by aqueous developing agents is easiest.* Higher light exposures result in larger numbers of silver atoms per latent image cluster on average in addition to exposing a larger number of crystals on average.

The detection of light by a silver halide crystal, subsequent conversion to an electron hole pair, and ultimate reduction of silver ions to metallic silver atoms is in essence a solid-state photophysical process. The application of solid-state theory to the photographic process began with Mott and

*Some specialized developing agents can etch into the silver halide crystal to develop the internal latent image (LI), but the most commonly used color negative and color paper developers have little capability to do this. They are primarily surface-developing agents. The color reversal black-and-white developer can develop slightly subsurface LI in color reversal emulsions.

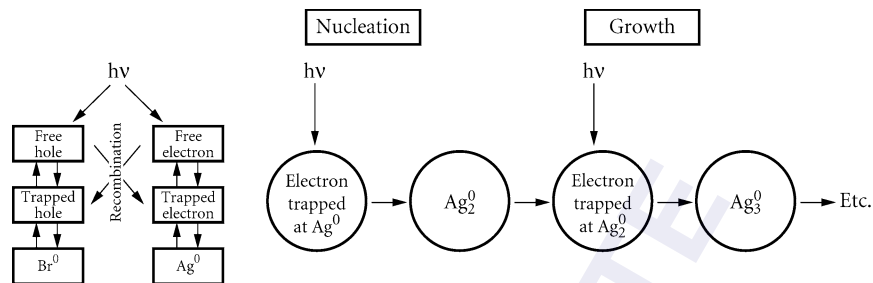


FIGURE 3 Diagram showing the basic features of the nucleation and growth mechanism for formation of a developable latent image site by light exposure of a silver halide crystal.

Gurney¹⁶ and was extended by Seitz.¹⁷ Additional reviews were published^{18–20} as our understanding of these processes grew.

According to the nucleation and growth mechanism,¹⁸ the photographic LI forms through a series of reversible electronic and ionic events (see Fig. 3). An electron is generated in the conduction band either by direct band gap photoexcitation of the crystal in its intrinsic spectral absorption region or by injection from a photoexcited spectral sensitizing dye on the crystal's surface. As this mobile electron travels through the crystal it becomes shallowly and transiently trapped at lattice imperfections carrying a positive charge. The ionized donor may be an interstitial silver ion or a surface defect. The kinked edge of a surface terrace, having a formal charge of $+1/2$ is one such site.

The electron may sample many such traps before becoming more permanently trapped at a preferred site. Trapping of the electron occurs in less than 0.1 ns at room temperature.²¹ The shallowly trapped electron orbits about the trapping site in a large radius. At this stage the formal charge of the trap site becomes $-1/2$ and adjacent surface lattice ions relax, thereby increasing the depth of the trap.²² In octahedral AgBr emulsions this shallow-deep transition occurs within 10 ns at room temperature.²¹

The negative charge formally associated with the trapped electron attracts a mobile interstitial Ag^+ ion, which the trapped electron reduces to an Ag^0 atom. The site's charge reverts to $+1/2$. A second photoelectron is then trapped at the site and reacts with a second interstitial Ag^+ ion to form Ag_2^0 and so on as additional photons strike the crystal. Frenkel defect concentrations—interior plus surface—in the range of $10^{14}/\text{cm}^3$ have been reported for AgBr crystals.¹⁹ Latent image formation depends critically upon the existence of partially charged defects on the surface of the crystal.

Because the single-atom Ag^0 state is metastable and the Ag_2^0 state is not usefully detectable from gross fog by developing solutions (Chap. 4 of Ref. 3), a minimum of three photons must strike the crystal to form a detectable LI site, the second photon within the lifetime of the metastable trapped Ag^0 atom.* The transient existence of one or two silver atom sites has not been directly observed but is strongly inferred by indirect evidence.²⁰

The highest photoefficiencies are achieved after silver halide emulsions are heated with aqueous sulfur and gold salts to produce Ag_2S and AgSAu species on the surface of the crystal.²² Most of the evidence for sulfur sensitization suggests these sites improve electron trapping, perhaps by making trapping lifetimes longer at preexisting kink sites. Gold sensitization accomplishes at least two things: it reduces the size of the latent image center needed to catalyze development and it stabilizes silver atoms formed during the exposure process (Chap. 5 of Ref. 3).

Ideally only a single latent image site forms per exposed silver halide crystal. If more sites nucleate on a crystal, these can compete for subsequent photoelectrons, leading to losses in efficiency related to LI dispersity. One remarkable aspect of this process is that despite the large number of

*A metastable silver atom may dissociate back into an electron-hole pair, and the electron can migrate to another trapping site to form a new silver atom. This may happen many times in principle. The second electron must encounter the silver atom before electron-hole recombination or other irreversible electron loss takes place.

kink sites and chemically sensitized sites on a given crystal, in the most efficient silver halide emulsions only a single LI center forms per grain. The shallow-deep transition that occurs after the initial electron trapping selects one preferred defect out of the plethora of other possibilities for subsequent electron trapping.

The two predominant crystal faces that occur on photographically useful silver halide grains are [100] and [111]. Both surfaces have a negative charge²³ ranging from -0.1 to -0.3 V. A subsurface space charge layer rich in interstitial Ag⁺ ions compensates for the charged surface.²⁴ The [100] surfaces are flat with steps containing jogs that are either positively or negatively charged kink sites. The positive kinks on the surface form shallow electron traps.

In contrast, the [111] surface is not well characterized and is believed to be sensitive to surface treatment. Calculations of surface energies suggest that several arrangements may coexist, including microfaceted [100] planes and half layers of all silver or all halide ions in arrangements with hexagonal symmetry.^{24,25}

All silver halides absorb light in the near-UV region of the spectrum. The absorption edge extends to longer wavelengths from AgCl to AgBr to AgI. AgCl and AgBr have indirect band gap energies of 3.29 and 2.70 eV, respectively.¹⁹ Absorption of light produces electrons in a conduction band whose minimum lies at the zone center and holes in the valence band whose maximum is at an L point. Carriers produced at excess energy rapidly thermalize at room temperature since electron mobility is limited by phonon scattering.

Electron mobility in AgBr is about 60 cm²/Vs at room temperature.¹⁹ In AgCl the hole mobility is low, but in AgBr it is substantial—only about a factor of 30 lower than electron mobility.²⁰ Carriers generated by exposure will explore the crystal in a random walk as they scatter off phonons and transiently trap at charged defects. It is expected that electrons can fully sample a 1- μ m crystal within less than 1 μ s.

In most commercial photographic materials today, only about 25 percent to 50 percent of the electrons injected into the crystal's conduction band contribute to a developable LI site. The rest are wasted through electron-hole recombination, formation of internal LI that is inaccessible to the developing agent, or formation of multiple nucleation sites termed *dispersity*. The deliberate addition of dopants that act as shallow electron traps reduces the time electrons spend in the free carrier state and thereby limits their propensity to recombine with trapped holes.

Polyvalent transition metal ions are frequently doped into silver halide crystals to limit reciprocity effects, control contrast, and reduce electron hole recombination inefficiencies.²⁶ They act as electron or hole traps and are introduced into the crystal from aqueous hexa-coordinated complexes during precipitation. They generally substitute for (AgX₆)⁵⁻ lattice fragments. Complexes of Ru and Ir have been especially useful. Because the dopant's carrier-trapping properties depend on the metal ion's valence state and the stereochemistry of its ligand shell, there are many opportunities to design dopants with specific characteristics. Dopants incorporating Os form deep electron traps in AgCl emulsions with an excess site charge of +3. An effective electron residence lifetime of 550 seconds has been measured for these dopants at room temperature.²⁰ They are used to control contrast at high-exposure regions in photographic papers.

Quantum sensitivity is a measure of the average number of photons absorbed per grain to produce developability in 50 percent of an emulsion population's grains (Chap. 4 in Ref. 27). This microscopic parameter provides a measure of the photoefficiency of a given emulsion; the lower the quantum sensitivity value, the more efficient the emulsion. The quantum sensitivity measurement gives a cumulative measure of latent image formation, detection, and amplification stages of the imaging chain.

Quantum sensitivity has been measured for many emulsions. The most efficient emulsions specially treated by hydrogen hypersensitization yield a quantum sensitivity of about three photons per grain.²⁸ Although hydrogen hypersensitization is not useful for general commercial films because it produces emulsions very unstable toward gross fog, this sensitivity represents an ambitious goal for practical emulsions used in commercial photographic products. The most efficient practical photographic emulsions reported to date have a quantum sensitivity of about five to six photons per grain.²⁹ The better commercial photographic products in today's market contain emulsions having quantum sensitivities in the range of 10 to 20 photons per grain.

30.4 THE STABILITY OF PHOTOGRAPHIC IMAGE DYES TOWARD LIGHT FADE

Introduction

Azomethine dyes are formed in most photographic products by reaction between colorless couplers and oxidized *p*-phenylenediamine developing agents to form high-saturation yellow, magenta, and cyan dyes (Chap. 12 in Ref. 3). Typical examples are shown in Fig. 4. The diamine structural element common to all three dyes comes from the developing agent, where R_1 and R_2 represent alkyl groups. The R group extending from the coupler side of the chromophore represents a lipophilic ballast, which keeps the dye localized in an oil phase droplet. In some dyes it also has hue-shifting properties.

Heat, humidity, and light influence the stability of these dyes in photographic films and papers.³⁰ Heat and humidity promote thermal chemical reactions that lead to dye density loss. Photochemical processes cause dyes to fade if the image is displayed for extended periods of time. Ultraviolet radiation is especially harmful to a dye's stability, which is partly why UV absorbers are coated in the protective overcoat layers of a color film or paper.

Stability toward light is especially important for color papers, where the image may be displayed for viewing over many years. It is less important in color negative film, which is generally stored in the dark and where small changes in dye density can often be compensated for when a print is made. Light stability is somewhat important for color reversal (slide and transparency) film, since a slide is projected through a bright illuminant, but projection and other display times are short compared to values for color papers. A similar situation exists for movie projection films, where each frame gets a short burst of high-intensity light but the cumulative exposure is low.

Evaluation of the stability of dyes toward light is difficult because the time scale of the photochemical reactions, by design, is very slow or inefficient. The quantum yields of photochemical fading of photographic dyes, defined as the fraction of photoexcited dyes that fade, are on the order of 10^{-7} or smaller.^{2,31} Accelerated testing using high-intensity illumination can sometimes give misleading results if the underlying photochemical reactions change with light intensity (p. 266 in Ref. 4).

Given that dyes fade slowly over time, it is best if all three photographic dyes fade at a common rate, thereby preserving the color balance of the image (p. 267 in Ref. 4). This rarely happens. The perceived light stability of color photographic materials is limited by the least stable dye. Historically, this has often been the magenta dye, whose gradual fade casts a highly objectionable green tint to a picture containing this dye. This is most noticeable in images containing memory colors, such as neutral grays and skin tones.

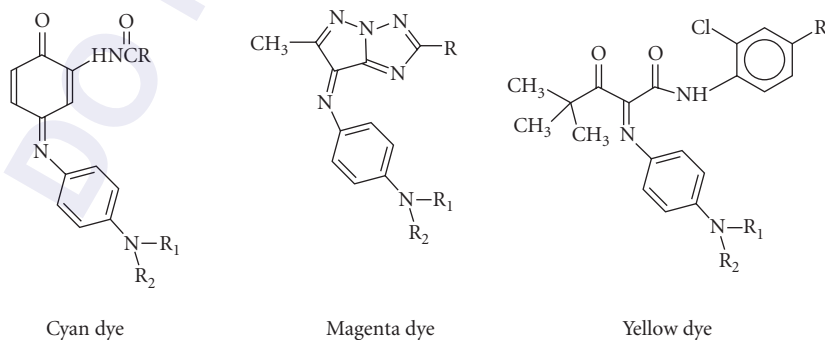


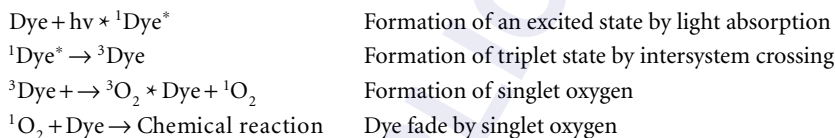
FIGURE 4 Typical examples of azomethine dyes used in photographic materials.

Photochemistry of Azomethine Dyes

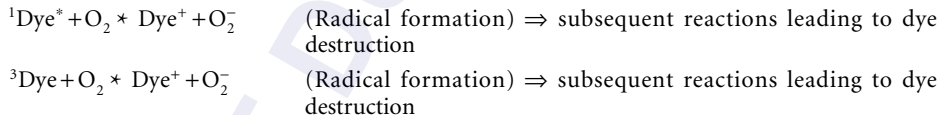
Upon absorption of a photon, a dye becomes reversibly excited to the singlet state. For azomethine dyes, the lifetime of the excited state is estimated to be on the order of picoseconds.^{32–39} Similarly, the lifetime of the triplet state of azomethine dyes, which is expected to form rapidly by intersystem crossing from the singlet excited state, has been estimated to be in the nanosecond range.³⁶ The short lifetime of these species seems to be at the basis of the observed low quantum yields of photochemical fading of azomethine dyes.

The nature of the initial elementary reactions involving excited or triplet states of dyes is not well understood. For some magenta dyes, the fading mechanism was reported to be photooxidative.³⁷ Cyan dye light fade appears to involve both reductive and oxidative steps³⁸ following photoexcitation of the dye.

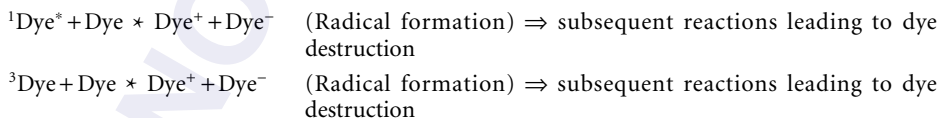
Singlet oxygen has been traditionally postulated to be involved in the oxidative pathway.³⁹ This high-energy species can, in principle, be formed by energy transfer between a triplet dye and molecular oxygen, and can subsequently induce dye destruction, although reaction with singlet oxygen does not always lead to decomposition.⁴⁰



Studies of the photophysical properties of photographic dyes, however, have shown that in general, azomethine dyes are good quenchers of singlet oxygen.⁴¹ This implies that the formation of singlet oxygen by the triplet state of the dyes should be inefficient. An alternative feasible role for molecular oxygen in dye fade involves electron transfer mechanisms.



The chemistry of dye fade may also be started by electron transfer between excited or triplet dye and a dye molecule in its ground state.



Excited State Properties

The few papers that have been published about excited singlet-state properties of azomethine dyes have mainly focused on pyrazolotriazole magenta dyes.^{32–35} These dyes have fluorescence quantum yields on the order of 10^{-4} at room temperature, but increase to ~ 1 in rigid organic glasses at 77 K.⁴⁰ The fluorescence quantum yield is the fraction of photoexcited molecules that emit a quantum of light by fluorescence from an excited state to the ground state having the same multiplicity (usually singlet to singlet). The results at room temperature imply singlet-state lifetimes of only a few picoseconds.⁴⁰

Room-temperature flash photolysis has identified a very short-lived fluorescent state plus a longer-lived nonfluorescent transient believed to be an excited singlet state whose structure is twisted compared to the ground state. This is consistent with the temperature-dependent results that allow rapid conformational change to the nonfluorescent singlet species at room temperature,

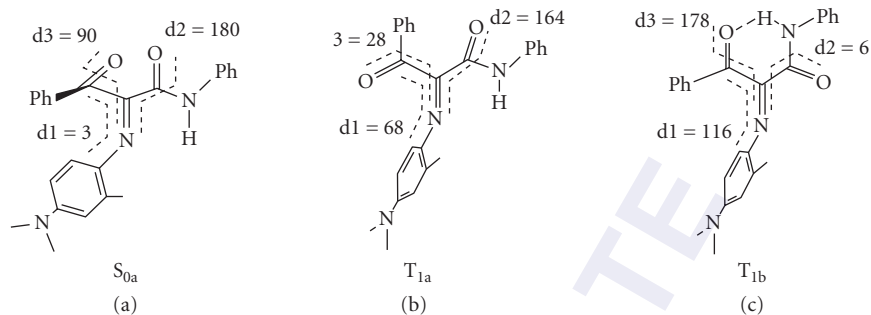


FIGURE 5 (a) Ground state conformation and (b and c) two conformations of the lowest triplet state.

in addition to intersystem crossing to the triplet excited state. But this conformational change is presumably restrained in a cold organic glass matrix, thereby allowing the excited singlet state time to fluoresce back to the ground state.

Investigation of the triplet states of azomethine dyes has proven difficult, and there are no definitive reports of direct observation of the triplet states. There are no reliable reports of phosphorescence from the triplet state of these dyes, although they have been studied by flash photolysis.^{32–35,41–45} Using indirect energy transfer methods, triplet lifetimes have been estimated to be less than 10 ns.³⁶ Similar studies⁴³ gave triplet energies of 166 to 208 kJ mol⁻¹ for yellow dyes, 90 to 120 kJ mol⁻¹ for magenta dyes, and about 90 kJ mol⁻¹ for cyan dyes.

Computational studies have been carried out recently on the ground and lowest triplet states of yellow azomethine dyes.⁴⁵ Consistent with crystallographic studies,^{46,47*} the calculated lowest-energy ground state conformation S_0 is calculated to have the azomethine C=N bond coplanar with the anilide C=O carbonyl, but perpendicular to the C=O carbonyl of the ketone. The energy of the vertical transition triplet state having this conformation is quite high, calculated to be 230 kJ mol⁻¹ above the ground state. Energy minimization on the lowest triplet energy surface starting with the geometry of S_0 results in a relaxation energy of 100 kJ mol⁻¹ and leads to T_{1a} , characterized by a substantial twist around the C=N bond and by increased planarity of the keto carbonyl with the azomethine plane. These conformational changes offer an efficient mechanism for stabilizing the lowest triplet state in these dyes. A second triplet conformer T_{1b} with an energy 25 kJ mol⁻¹ below T_{1a} was also calculated (see Fig. 5).

The low energy values calculated for T_{1a} and T_{1b} relative to S_{0a} (130 and 105 kJ mol⁻¹, respectively) have led to a new interpretation of the reported results⁴³ of energy transfer to S_{0a} . Using the nonvertical energy transfer model of Balzani,⁴⁸ it is shown that observed rates of energy transfer to S_{0a} can be consistent with formation of two distinct triplet states: one corresponding to a higher-energy excited triplet state (T_2) 167 kJ mol⁻¹ above S_{0a} , and the second corresponding to one conformation of the lowest triplet state (T_1) with an energy of 96 kJ mol⁻¹ above S_{0a} , in good agreement with the calculated T_{1b} configuration. The large structural differences between S_{0a} and the conformers of T_1 can explain the lack of phosphorescence in this system.

Light Stabilization Methods

Stabilization methods focus on elimination of key reactants (such as oxygen by various barrier methods), scavenging of reactive intermediates such as free radicals transiently formed during photodecomposition, elimination of ultraviolet light by UV absorbers, or quenching of photoexcited species in the reaction sequence.

*Other unpublished structures exhibiting the same conformation of S_{0a} have been solved at Kodak.

Stabilizer molecules⁴⁹ have been added to photographic couplers to improve the stability of their dyes toward light fade. Although the exact mode by which these stabilizers operate has not been disclosed, it is probable that some quench the photoexcited state of the dyes, thereby preventing further reaction leading to fading; some scavenge the radicals formed during decomposition steps; and some may scavenge singlet oxygen.

Polymeric materials added to coupler dispersions also stabilize dyes. These materials are reported to improve thermal dye decomposition⁵⁰ by physically separating reactive components in a semi-rigid matrix or at least decreasing the mobility of reactive components within that matrix. They may provide benefits for light stability of some dyes by a similar mechanism.

Photographic dyes sometimes form associated complexes or microcrystalline aggregates (p. 322 of Ref. 4). These have been postulated to either improve light stability by electronic-to-thermal (phonon) quenching of excited states, or to degrade light stability by concentrating the dye to provide more available dye molecules for reaction with light-induced radicals. Both postulates may be correct depending upon the particular dyes.

Modern photographic dyes have vastly improved light stabilities over dyes of the past. For example, the magenta dye light stability of color paper dyes has improved by about two orders of magnitude between 1942 and the present time.⁵¹ New classes of dyes^{52,53} have proved especially resistant to light fade and have made memories captured by photographic materials truly archival.

30.5 PHOTOGRAPHIC SPECTRAL SENSITIZERS

Introduction

Spectral sensitizing dyes are essential to the practice of color photography. In 1873, Vogel⁵⁴ discovered that certain organic dyes, when adsorbed to silver halide crystals, extend their light sensitivity to wavelengths beyond their intrinsic ultraviolet-blue sensitivity. Since then, many thousands of dyes have been identified as silver halide spectral sensitizers having various degrees of utility, and dyes exist for selective sensitization in all regions of the visible and near-infrared spectrum. These dyes provide the red, green, and blue color discrimination needed for color photography.

The most widely used spectral sensitizing dyes in photography are known as *cyanine dyes*. One example is shown in Fig. 6. These structures are generally planar molecules as shown.

A good photographic spectral sensitizing dye must adsorb strongly to the surface of the silver halide crystal. The best photographic dyes usually self-assemble into aggregates on the crystal's surface. Aggregated dyes may be considered partially ordered two-dimensional dye crystals.⁵⁵ Blue-shifted aggregates are termed *H-aggregates*, while red-shifted ones, which generally show spectral narrowing and excitonic behavior, are termed *J-aggregates*. Dyes that form J-aggregates are generally most useful as silver halide spectral sensitizers.

A good sensitizing dye absorbs light of the desired spectral range with high efficiency and converts that light into a latent image site on the silver halide surface. The relative quantum efficiency⁵⁶ of sensitization is defined as the number of quanta absorbed at 400 nm in the AgX intrinsic absorption region to produce a specified developed density, divided by the number of like quanta absorbed only by dye within its absorption band. For the best photographic sensitizing dyes this number is only slightly less than unity.

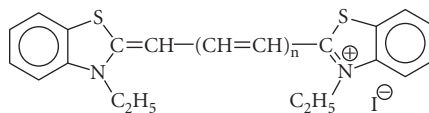


FIGURE 6 Typical cyanine spectral sensitizing dye used in photographic materials.

Adsorption isotherms for J-aggregated cyanine dyes on AgX follow Langmuir behavior with a saturation coverage that indicates closely packed structures.⁵⁷ The site area per molecule derived from the Langmuir isotherms implies dye adsorption along its long axis with heterocyclic rings perpendicular to the surface. Recent scanning tunneling microscope work^{58,59} has convincingly detailed the molecular arrangement within these aggregates, confirming general expectations.

The Photophysics of Spectral Sensitizers on Silver Halide Surfaces

Sensitizing dyes transfer an electron from the dye's excited singlet state to the conduction band of AgX. Evidence for this charge injection comes from extensive correlations of sensitization with the electrochemical redox properties of the dyes.^{22,60,61} In the adsorbed state a good sensitizing dye has an ionization potential smaller than that of the silver halide being sensitized. This allows the energy of the dye excited state to lie at or above the silver halide conduction band even though the energy of the photon absorbed by the dye is less than the AgX band gap. Dye ionization potentials are generally well correlated with their electrochemical oxidation potentials.

The triplet energy level of a typical cyanine dye lies about 0.4 to 0.5 eV (about 35 to 50 kJ mol⁻¹) below the radiationally excited singlet level.⁶² In most cases, this triplet state seems to have little effect on spectral sensitization at room temperature, although it may influence sensitization in special cases.⁶³

Electron transfer takes place by tunneling from the excited-state dye to the silver halide conduction band. Simple calculations verify that penetration of the potential barrier will compete favorably with de-excitation of the dye by fluorescence emission (p. 253 of Ref. 3). The interaction time between an excited spectral sensitizer and silver bromide appears to occur between 10⁻¹³ and 10⁻¹⁰ (p. 237 of Ref. 3). Factors favoring irreversible flow of electrons from dye to silver halide are delocalization within the conduction band, the negative space charge layer on the surface of the crystal, trapping by remote sites on the silver halide surface, and irreversible trapping to form a latent image.

Free electrons from dye sensitization appear in the silver halide conduction band having the same mobility and lifetime as those formed by intrinsic absorption. These electrons participate in latent image formation by the usual mechanism.

After electron transfer, the oxidized dye radical cation becomes the hole left behind. Because there is evidence that a single dye molecule may function repeatedly,⁶⁴ the dye cation "hole" must be reduced again. This can occur by electron tunneling from an occupied site on the crystal's surface, whose energy level is favorable for that process. A bromide ion at a kink is one such site. This leaves a trapped hole on the crystal surface that may be thermally liberated to migrate through the crystal. This hole can lead to conduction band photoelectron loss by recombination. A supersensitizer molecule (discussed later) may also trap the hole.

The formation of the J-aggregate exciton band has significant effects on the excited-state dynamics of the dye. There is the possibility of coherent delocalization of the exciton over several molecules in the aggregate.⁶⁵ The exciton is mobile and can sample over 100 molecules within its lifetime.⁶⁶ This mobility means that the exciton is sensitive to quenching by traps within the aggregate structure. The overall yield of sensitization generally increases to maximum, then decreases somewhat as the aggregate size increases. Optimum sensitizations usually occur below surface monolayer coverage.

Sometimes the spectral sensitivity of a dye is increased by addition of a second substance. If the added sensitivity exceeds the sum of both sensitizers individually, the increase is super-additive and the second substance is called a *supersensitizer*. Maximum efficiency of a supersensitizer often occurs in the ratio of about 1:20 where the supersensitizer is the dilute component.

The supersensitizer may increase the spectral absorption of the sensitizer by inducing or intensifying formation of a J-aggregate. In some cases these changes are caused by a mutual increase in adsorption to the grain surface, as when the sensitizer and supersensitizer are ions of opposite charge.⁶⁷ However, the most important supersensitizers appear to increase the fundamental efficiency of spectral sensitization as measured by the relative quantum yield.

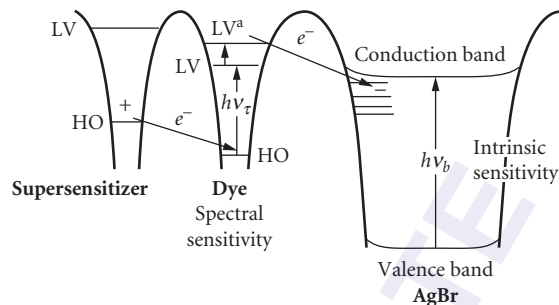


FIGURE 7 Simplified energy level diagram for a spectral sensitizer dye plus a supersensitizer dye adsorbed to a silver bromide surface.

The low concentration of supersensitizer relative to sensitizer suggests trapping of an exciton moving through the quasicrystalline array of J-aggregated sensitizer molecules. One hypothesis is that the supersensitizer molecule traps the exciton, providing a site for facile electron transfer into the silver halide.⁶⁸

Gilman and coworkers^{69–72} proposed that the supersensitization of J-aggregating dyes takes place via hole trapping by the supersensitizer molecules. The exciton moving through the aggregate may be self-trapped at a site adjacent to a supersensitizer molecule. By hypothesis, the highest occupied (HO) electron energy level of the supersensitizer is higher than that of the partially empty HO level of the excited sensitizer molecule. An electron transfers from the supersensitizer to the sensitizer molecule. The sensitizer is reduced to an electron-rich, anion-free radical while the supersensitizer is oxidized to an electron-deficient cation radical of lower energy than the original hole, thereby localizing the hole on the supersensitizer (see Fig. 7).

Following electron transfer from the supersensitizer, the electron-rich free radical anion left on the sensitizer will possess an occupied electron level of higher energy than the excited level of the parent sensitizer. The reduction process thereby raises the level of the excited electron with respect to the conduction band of AgX, with a consequent increase in probability of electron tunneling into the silver halide conduction band.

The various proposed mechanisms of supersensitization are not mutually exclusive. Exciton trapping and hole trapping may operate together. Underlying all mechanisms are the roles of exciton migration throughout the aggregate, its interruption at or near the supersensitizer site, and a more facile transfer of the electron into the AgX conduction band at that localized state.

Spectral sensitizing dyes, especially at high surface coverage, desensitize silver halides in competition with their sensitization ability. Red spectral sensitizers generally desensitize more than green or blue spectral sensitizers do. Desensitization can be thought of as reversible sensitization. For example, a mobile conduction band electron can be trapped by a dye molecule to form a dye radical anion⁷³ or by reduction of a hole trapped on the dye.^{63,74–76} Under normal conditions of film use (not in a vacuum), a dye radical anion may transfer the excess electron irreversibly to an O₂ molecule with formation of an O₂⁻ anion. Either postulate leads to irreversible loss of the photoelectron and consequent inefficiency in latent image formation.

Color Science of Photographic Spectral Sensitizers

Human vision responds to light in the range of 400 to 700 nm. The human eye is known to contain two types of light-sensitive cells termed *rods* and *cones*, so named because of their approximate shapes. Cones are the sensors for color vision. There is strong evidence for three types of cones sometimes termed long (L), middle (M), and short (S) wavelength receptors (Chap. 1 in Ref. 1). The normalized spectral responses of the human eye receptors (Chap. 26, Table 5 in Ref. 1) are

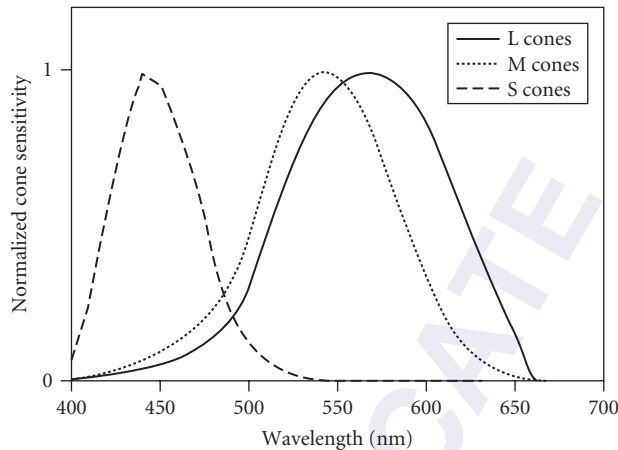


FIGURE 8 Normalized spectral sensitivity of the human eye.

shown in Fig. 8. The ability of humans to distinguish differences in spectral wavelength is believed to rely on difference signals created in the retina and optic nerve and interpreted by the brain from responses to light by the three cone detectors (Chap. 9 in Ref. 4).

The human eye's red sensitivity function peaks at around 570 nm and possesses little sensitivity at wavelengths longer than 650 nm. Also, there exists considerable overlap between red and green sensitivity functions, and to a lesser extent overlap between blue and green and blue and red sensitivity functions. This overlap in sensitivity functions, combined with signal processing in the retina, optic nerve, and brain, allows us to distinguish subtle color differences while simultaneously appreciating rich color saturation.

Photographic films cannot do this. The complex signal processing needed to interpret closely overlapping spectral sensitivity functions is not possible in a photographic material (Chap. 9 in Ref. 4). Some relief from this constraint might be anticipated by considering Maxwell's principle.⁷⁷

Maxwell's principle states that any visual color can be reproduced by an appropriate combination of three independent color stimuli called *primaries*. The amount of each primary per wavelength needed to reproduce all spectral colors defines three color-matching functions (Chap. 6 in Ref. 1) for those primaries. If there exists a set of color-matching functions that minimize overlap between spectral regions, a photographic film with spectral sensitivity like that set of color-matching functions and imaging with the corresponding primary dyes would faithfully reproduce colors.

The cyan, magenta, and yellow dye primaries used in photographic materials lead to theoretical color-matching functions (Chap. 19-II in Ref. 3; Ref. 78)* having more separation than human visual sensitivity, but possessing negative lobes as shown in Fig. 9. This is impossible to achieve in practical photographic films (Chap. 9 in Ref. 4), although the negative feedback from HE effects provides an imperfect approximation. Approximate color-matching functions have been proposed⁸⁰ that contain no negative lobes. But these possess a very short red sensitivity having a high degree of overlap with the green sensitivity, similar to that of the human eye.

It is therefore not possible to build a perfect photographic material possessing the simultaneous ability to distinguish subtle color differences, especially in the blue-green and orange spectral regions, and to render high-saturation colors. Compromises are necessary.

*The theoretical color-matching functions were calculated using the analysis described in Ref. 78, Eq. (9a-c) for block dyes having trichromatic coefficients in Table II, and using the Judd-Vos modified XYZ color matching functions in Chap. 26, Table 2 of Ref. 1.

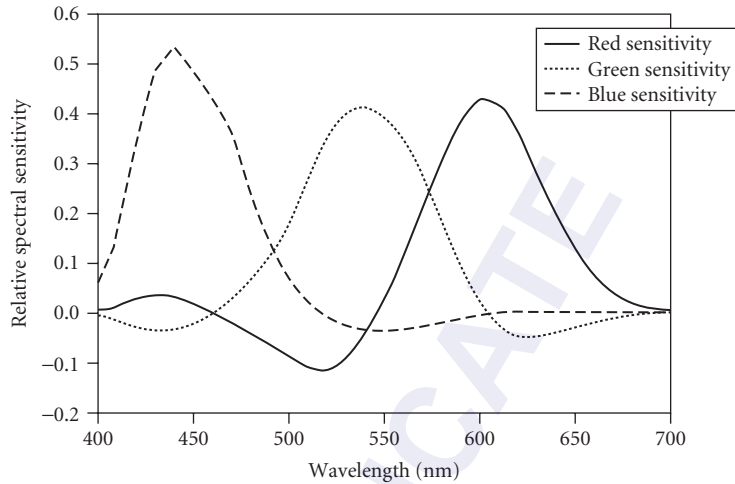


FIGURE 9 Spectral sensitivity for idealized photographic dyes. Note the large negative lobe around 525 nm.

Most photographic materials intended to image colors at high saturation possess spectral sensitivities having small overlap, and in particular possess a red spectral sensitivity centered at considerably longer wavelengths than called for by theoretical color-matching functions. This allows clean spectral detection and rendering among saturated colors but sacrifices color accuracy.

For example, certain *anomalous reflection colors* such as the blue of the lobelia flower possess a sharp increase in spectral reflectance in the long red region not seen by the human eye but detected strongly by a photographic film's spectral sensitivity (see Fig. 10). This produces a serious color

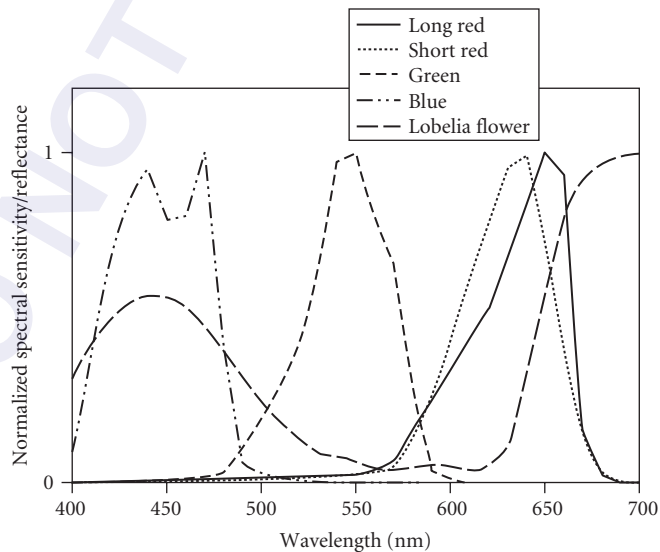


FIGURE 10 Spectral sensitivity for a real photographic film plus reflectance of the lobelia flower.

accuracy error (the flower photographs pink, not blue).⁸⁰ Among other such colors are certain green fabric dyes that image brown or gray. To correct for these hue errors, a special color reversal film was made whose general characteristic involves a short red spectral sensitivity.* This film does a considerably better job of eliminating the most serious hue errors, but unfortunately does not possess high color saturation. It is important in catalogue and other industries where accurate color reproduction is of paramount importance, but most consumers prefer saturated colors.

More recently, color negative films possessing short red spectral sensitivities combined with powerful development inhibitor anchimeric releasing (DIAR) color correction chemistry and masking coupler technology have produced films[†] simultaneously having high color saturation plus accurate color reproduction, at least toward anomalous reflection colors. In addition, these films render improved color accuracy under mixed daylight plus (red-rich) incandescent lighting conditions. These films maintain the high color saturation capability important to consumers but also provide more accurate color hues—although they are not perfect.

A new type of color negative film has been described^{§1} that possesses a fourth color record whose function approximates the large negative lobe in the theoretical red color-matching function of photographic dyes. Several embodiments have been described and are employed in modern films. The most practical utilizes a fourth color record containing silver halide emulsions sensitive to short green light plus a magenta DIR color correction coupler.[‡]

In simple terms, this special color record responds to short green light in the negative lobe region by releasing a development inhibitor that travels into the red-sensitive cyan color record, thereby reducing cyan dye in response to short green light, an effect opposite to that of normal red light exposure. This film's advantages reside in its ability to more accurately record colors in the blue-green spectral region, plus improved color accuracy under certain fluorescent lights that have a pronounced blue-green emission component.

Modern films are moving toward shorter-wavelength red spectral sensitivity. This allows improved color accuracy by shifting the red spectral sensitivity closer to that of the human eye. In addition, films have been designed to crudely approximate the principal negative lobe in the theoretical red color-matching function for photographic dye primaries, thereby producing a more accurate rendition of certain blue-green colors. Although not perfect, these techniques allow modern photographic films to record the world of color in all its richness and subtleties nearly as we ourselves see it.

30.6 GENERAL CHARACTERISTICS OF PHOTOGRAPHIC FILMS

Today there is available a large variety of films including different film speeds, types, and sizes. Many films are offered in both professional and consumer categories. This section describes some general characteristics of films as they relate to film speed and type. Some unique characteristics of professional films are contrasted to consumer films. Finally, color reversal, black-and-white, and color negative films are compared in more detail.

Low-Speed Films Compared to High-Speed Films

Higher granularity, higher sensitivity to ambient background radiation, and higher sensitivity to x rays are fundamental penalties that come with higher photographic speed. High-speed films often suffer additional penalties compared to their lower-speed counterparts, but these additional penalties

*Kodak Ektachrome 100 Film, EPN—a color reversal professional film.

†Kodak Gold 100, Kodak Gold 200, and Kodak Max 400 color negative films.

‡Fuji Realia 100 Film; Fuji Superia 200 and 400 film; Fuji Nexia 200 and 400 film; and Fuji 160 NPL, NPS, and NC professional films. These are all color negative films.

are not intrinsic. They represent conscious choices made by the manufacturer nearly always traceable to the desire to minimize granularity.

As one example, 400-speed color slide film generally has lower color saturation than its 100-speed counterpart. The reduced color saturation comes as a result of lower interlayer interimage effects (IIE) caused by higher iodide content in the mixed silver bromo-iodide crystals, which, in turn, is needed to achieve higher speeds from a smaller sized crystal and is ultimately done to reduce the photographic granularity penalty.

Photographic Speed and Photographic Granularity To understand the intrinsic connection between photographic speed and photographic granularity, recall that light quanta striking a flat surface follow Poisson statistics (Chap. 1 of Ref. 27). Therefore, a uniform light exposure given to a photographic film delivers photon hits Poisson-distributed across the plane of the film. Silver halide crystals are on the order of $1\ \mu\text{m}^2$ in projected area. On average in a film exposure (the product of the intensity of the radiation and the duration) appropriate for an EI100-speed film, the most sensitive silver halide crystals receive:*

E_q (red light; 650 nm)	~48 photons per square micrometer
E_q (green light; 550 nm)	~40 photons per square micrometer
E_q (blue light; 450 nm)	~33 photons per square micrometer

Because the most efficient practical photographic emulsions to date have a quantum sensitivity of about five photons per grain,²⁹ these approximate numbers suggest there is enough light (on average) striking a film during an EI100-speed film exposure to form latent image on efficient $1\text{-}\mu\text{m}^2$ silver halide grains. But an EI800-speed exposure at three stops down begins to push the limits of modern emulsion efficiency at five green photons per square micron.

Silver halide crystals within the film are light collectors. Just as large-area mirrors in telescopes collect more photons and are more sensitive detectors of starlight than smaller mirrors, all other things being equal, larger silver halide crystals likewise collect more photons to become more sensitive light detectors.

Strong correlates exist between silver halide surface area and photographic speed. Photographic speed often linearly follows the average surface area of the silver halide crystal population (p. 969 of Ref. 2): speed $\sim \bar{a}$.

According to the well-known Siedentopf formula (p. 60 in Ref. 27), the root-mean-square photographic granularity at density D is given by $G \sim D\hat{a}$. In this formula, the symbol \hat{a} represents the average cross-sectional area of the image particle projected onto the film plane. To a good approximation, this is proportional (not necessarily linear) to the mass of developed silver regardless of whether the image particle is a dye cloud or silver deposit.

These considerations demonstrate the intrinsic and well-known correlation between photographic speed and photographic granularity. For a given crystal morphology, as the size of the silver halide crystal increases, both surface area and mass increase, leading to a speed plus granularity increase.

At a fundamental level, to lower the granularity of a higher-speed photographic film the sensitivity to light of the silver halide crystals must improve. Then smaller silver halide crystals can be used to attain speed. Factors that contribute to this problem were discussed in Sec. 30.2.

In general, larger silver halide crystals suffer more inefficiencies in the latent-image-forming process because latent-image dispersity, electron-hole recombination, and other irreversible photon waste processes become more problematic as the grain size increases (p. 993 in Ref. 2). Overcoming these inefficiencies at all speeds, but especially at high photographic speeds, is a major challenge of modern photographic science.

Photographic Speed and Sensitivity to High-Energy Radiation High-energy subatomic particles induce developable latent image formation in photographic films. The process by which high-energy

*Analysis based on a typical EI100-speed color reversal film where the exposure falls near the photographically fastest, most sensitive part of the sensitometric curve. This analysis is described on p. 47 of Ref. 27.

charged particles produce latent image usually involves energy transfer to valence band electrons by inelastic scattering. Some of these scattered electrons are promoted into the silver halide conduction band and go on to produce latent image by the usual processes. Higher-speed films are more sensitive to this radiation because their larger silver halide crystals are hit (on average) more frequently and more collisions occur by passage through a larger crystal than through a smaller crystal.

Photographic films are slowly exposed by ambient background radiation (Chap. 23 in Ref. 3) caused by cosmic rays; ubiquitous stone and masonry building materials, which contain trace amounts of radioactive elements; and other low-level radiation sources. These radiation sources emit various types of charged and uncharged high-energy particles and photons including alpha particles, neutrons, γ rays, β particles, muons, and others. This results in a shorter ambient shelf life for very-high-speed films unless special storage precautions are taken.

Photographic emulsions are also sensitive to x rays. The formation of latent image silver by x rays (and γ rays) is due to the action of high-energy electrons emitted during the absorption or inelastic scattering of the electromagnetic radiation by matter (Chap. 23 in Ref. 3). These electrons may enter the conduction band of silver halide and proceed to latent image formation in the usual way. If the primary emitted electron is of sufficiently high energy, it may create secondary electrons of lower energy by inelastic scattering with other electrons. These secondary electrons enter the conduction band.

High-speed films have more sensitivity to x rays than do low-speed films. Some airline travelers ask airport security agents to hand-check their film rather than allowing it to pass through the baggage x-ray scanner. This precaution is prudent for higher-speed film because airport x-ray machines have become stronger emitters recently, especially at airports with special security concerns.

Despite the higher price, inherent penalties, and willful compromises that characterize high-speed films, they offer access to more picture space under low-light conditions. Among the advantages a high-speed film offers are the following:

1. Pictures under low light where flash is not permitted, such as theaters and museums.
2. Pictures under low light using slower lenses such as zoom and telephoto lenses.
3. Pictures under low light when the subject lies beyond the reach of flash.
4. Pictures under low light taken in rapid sequence without waiting for a flash to recharge.
5. Pictures under low light when flash disturbs the subject, as in some nature photography.
6. Pictures of people without flash. Flash is responsible for a major defect known as *red eye*, caused when light from a flash reflects off the retina of the eye.
7. Pictures using a stopped-down camera aperture for improved depth of field, so objects over a wider range of distance from the camera are imaged in focus.
8. Reduced dependence on flash for longer battery life.
9. Low-cost one-time-use cameras (OTUCs) when flash units are not needed.
10. Ability to freeze action using a rapid shutter speed.

For these and other reasons, a high-speed film may be the best choice for many film applications despite its shortcomings compared to lower-speed alternatives.

Professional Film Compared to Amateur Film

Most film manufacturers offer film in a professional category and an amateur (sometimes called *consumer*) category. These differ in many respects.

Consumer film is targeted for the needs of the amateur snap shooter. Drugstores, supermarkets, department stores, and other outlets frequented by nonprofessional photographers sell this film. Because most amateurs are price conscious, discount brands of film are available having some performance compromises that may not be important to some consumers.

To keep the price low, manufacturing tolerances for consumer film performance attributes are wider than those for professional film. In a well-manufactured film these departures from performance

aims are small enough to go unnoticed by most average consumers but large enough to be objectionable to highly discerning professionals.

People want pleasing color in their pictures that renders the main subject attractive to the eye. For most amateurs this means saturated colors enhanced to some extent beyond their true appearance in the original scene. However, this preference is highly selective since over-saturated skin tones are objectionable. Other *memory colors*—those that we recognize and know what they should be, such as green grass and blue sky—must also be rendered carefully to produce the most pleasing color appearance (Chap. 5 in Ref. 4). The best consumer films accommodate these color preferences.

Exposures on a roll of consumer film are sometimes stored for long times. It is not unusual for amateurs to take pictures spanning two successive Christmas seasons on the same roll of film. Thus, latent image stability under a wide range of temperature and humidity conditions is very important to a good consumer film. On the other hand professional photographers develop their film promptly, making this requirement far less demanding in a professional film.

Most amateurs are not skilled photographers. It is advantageous for consumer film to be easy to use and provide good-looking pictures under a wide variety of conditions. Consumer film incorporates more tolerance toward over- and underexposures.

Film for amateurs is offered in three standard formats. The most common format is the 35-mm cassette. Consumer cameras that accept this film have become compact, lightweight, and fully automatic. The camera senses the film speed through the DX coding on the film cassette and adjusts its shutter, aperture, and flash accordingly. Autofocus, autorewind, and automatic film advance are standard on these cameras, making the picture-taking experience essentially point and shoot.

The new Advanced Photo System (APS) takes ease of use several steps further by incorporating simple drop-in film loading and a choice of three popular print aspect ratios from panoramic to standard size prints. A magnetic coating on the film records digital data whose use and advantages are just beginning to emerge.

The third format is 110 film. Although not as popular as 35-mm and APS, it is often found in inexpensive cameras for gift packages. Film for this format comes in a cartridge for easy drop-in loading into the camera.

Professional photographers demand more performance from their film than do typical amateurs. These people shoot pictures for a living. Their competitive advantage lies in their ability to combine technical skill and detailed knowledge of their film medium with an artistic eye for light, composition, and the color of their subject.

Film consistency is very important to professionals, and they pay a premium for very tight manufacturing tolerances on color balance, tone scale, exposure reciprocity, and other properties. They buy their film from professional dealers, catalogues, and in some cases directly from the manufacturer. Professionals may buy the same emulsion number in large quantity, then pretest the film and store it in a freezer to lock in its performance qualities. Their detailed knowledge about how the film responds to subtleties of light, color, filtration, and other characteristics contributes to their competitive advantage.

Professional film use ranges from controlled studio lighting conditions to harsh and variable field conditions. This leads to a variety of professional films designed for particular needs. Some photographers image scenes under incandescent lighting conditions in certain studio applications, movie sets, and news events where the lights are balanced for television cameras. These photographers need a tungsten-balanced film designed for proper color rendition when the light is rich in red component, but deficient in blue compared to standard daylight.

Catalogue photographers demand films that image colors as they appear to the eye in the original scene. Their images are destined for catalogues and advertisement media that must display the true colors of the fabric. This differs from the wants of most amateurs, who prefer enhanced color saturation.

Some professional photographers use high color saturation for artistic effect and visual impact. These photographers choose film with oversaturated color qualities. They may supplement this film by using controlled light, makeup on models, lens gels, or other filtration techniques to control the color in the image.

Professional portrait photographers need film that gives highly pleasing skin tones. Their films typically image at moderate to low color saturation to soften facial skin blotches and marks that most people possess but don't want to see in their portraits. In addition, portrait film features highly

controlled upper and lower tone scale to capture the subtle texture and detail in dark suits and white wedding dresses and the sparkle in jewelry.

Professional photographers image scenes on a variety of film sizes ranging from standard 35-mm size to 120-size, and larger 4×5 to 8×10 sheet-film formats that require special cameras. Unlike cameras made for amateurs, most professional cameras are fully controllable and adjustable by the user. They include motor drives to advance the film rapidly to capture rapid sequence events. Professional photographers often take many pictures to ensure that the best shot has been captured.

Many specialty films are made for professional use. For example, some films are sensitized to record electromagnetic radiation in the infrared, which is invisible to the human eye. Note that these films are not thermal detectors. Infrared films are generally *false colored*, with the infrared-sensitized record producing red color in the final image, the red-light-sensitized record producing green color, and the green/blue-sensitized record producing blue color. Other professional films are designed primarily for lab use, such as duplicating film used to copy images from an original capture film.

Some film applications are important to a professional but are of no consequence to most amateurs. For example, professionals may knowingly underexpose their film and request *push processing*, whereby the film is held in the developing solution longer than normal. This brings up an underexposed image by rendering its midtone densities in the normal range, although some loss in shadow detail is inevitable. This is often done with color transparency film exposed under low-ambient-light conditions in the field. Professional films for this application are designed to provide invariant color balance not only under normal process times but also under a variety of longer (and shorter) processing times.

Multipop is another professional technique whereby film is exposed to successive flashes or *pops*. Sometimes the photographer moves the camera between pops for artistic effect. The best professional films are designed to record a latent image that remains nearly invariant toward this type of exposure.

Consumer film is made for amateur photographers and suits their needs well. Professional film comes in a wide variety of types to fit the varied needs of highly discerning professional photographers who demand the highest quality, specialized features, and tightest manufacturing tolerances—and are willing to pay a premium price for it.

Color Reversal Film

Color reversal films, sometimes called *color transparency* or *color slide* films, are sold in professional and amateur categories. This film is popular among amateurs in Europe and a few other parts of the world but has largely been replaced by color negative film among amateurs in North America. However, it is still very popular among advanced amateurs, who favor it for the artistic control this one-stage photographic system offers them, the ease of proofing and editing a large number of images, and the ability for large-audience presentations of their photographic work. Because the image is viewed by transmitted light, high amounts of image dye in these films can project intense colors on the screen in a dark auditorium for a very striking color impact.

Color reversal film is the medium of choice for most professional commercial photographers. It offers them the ability to quickly proof and select their most favored images from a large array of pictures viewed simultaneously on a light table. The colors are more intense compared to a reflection print, providing maximum color impact to the client. And professional digital scanners have been optimized for color reversal images, making it convenient to digitally manipulate color reversal images and transcribe them to the printed page. Today's professional imaging chain infrastructure has been built around color reversal film, although digital camera image capture is making inroads, especially for low-quality and low-cost applications.

Large-format color reversal film, from 120 size to large sheet formats, is often shot by professionals seeking the highest-quality images destined for the printed page. A larger image needs little or no magnification in final use and avoids the deterioration in sharpness and grain that comes from enlarging smaller images. Additionally, art directors and other clients find it easier to proof larger images for content and artistic appeal.

EI400-speed color reversal films generally possess low color saturation, low acutance, high grain, and high contrast compared to lower-speed color reversal films. The high contrast can lead to some loss in highlight and shadow detail in the final image. Except for very low-light situations, EI100 is

the most popular speed for a color reversal film. It offers the ultrahigh image quality professionals and advanced amateurs demand.

Modern professional color reversal films maintain good color balance under push processing conditions. They also perform very well under a variety of exposure conditions including multipop exposures and very long exposure times.

When professional photographers believe nonstandard push or pull process times are needed, they may request a “clip test” (sometimes also called a “snip test”) whereby some images are cut from the roll of film and processed. Adjustments to the remainder of the roll are based on these preprocessed images. Some photographers anticipate a clip test by deliberately shooting sacrificial images at the beginning of a roll.

Professional color reversal films can be segregated into newer modern films and older *legacy films*. Modern films contain the most advanced technology and are popular choices among all photographers. Legacy films are popular among many professional photographers who have come to know these films’ performance characteristics in detail after many years of experience using them.

Legacy films continue to enjoy significant sales because the huge body of personal technical and artistic knowledge about a film accumulated by a professional photographer over the years contributes to his or her competitive advantage. This inertia is a formidable barrier to change to a new, improved, but different film whose detailed characteristics must be learned.

In some cases, the improved features found in modern films are not important to a particular professional need. For example, a large 4 × 5 sheet film probably needs little or no enlargement, so the poorer grain and sharpness of a legacy film may be quite satisfactory for the intended application. The detailed knowledge the professional has about how the legacy film behaves in regard to light and color may outweigh the image structure benefits in a modern film. Also, all films possess a unique tone scale and color palette. A professional may favor some subtle and unique feature in the legacy film’s attributes.

Kodachrome is a special class of color reversal legacy films. Unlike all other types of color reversal films, this film contains no dye-forming incorporated couplers. The dyes are imbibed into the film during processing so the film’s layers are coated very thin. This results in outstanding image sharpness. Moreover, many photographers prefer Kodachrome’s color palette for some applications.

Kodachrome dyes are noted for their image permanence. Images on Kodachrome film have retained their colors over many decades of storage, making this film attractive to amateurs and professionals alike who want to preserve their pictures. However, modern incorporated coupler color reversal films also have considerably improved dark storage image stability compared to their predecessors.

Table 1 lists some color reversal films available today. This is not a comprehensive list. Most film manufacturers improve films over time, and portfolios change frequently.

TABLE 1 Color Reversal Capture Films

Modern Films	Legacy Films	Kodachrome Films	Tungsten Films
Kodak Ektachrome 100VS (very saturated color)	Fuji Velvia 50 RVP (very saturated color)	Kodachrome 25 PKM	Kodak Ektachrome 64T EPY
Kodak Ektachrome E100S (standard color, good skin tones)	Kodak Ektachrome 64 EPR (standard color)	Kodachrome 64 PKR	Fujichrome 64T RTP
Fuji Astia 100 RAP (standard color, good skin tones)	Kodak Ektachrome 100 Plus EPP (standard color)	Kodachrome 200 PKL	Agfachrome 64T
Fuji Provia 100 RDPII (standard color)	Kodak Ektachrome 100 EPN (accurate color)		Kodak Ektachrome 160T EPT
Agfachrome 100 (standard color)	Kodak Ektachrome 200 EPD (lower color saturation)		Kodak Ektachrome 320T EPJ
Kodak Ektachrome E200 (can push 3 stops)	Agfachrome 200		
	Kodak Ektachrome 400 EPL		
	Fuji Provia 400		
	Agfachrome 400		

Unless a photographer has some particular reason for choosing a legacy color reversal film, the modern films are generally a better overall choice. More detailed information about these films' uses and characteristics can be found in the film manufacturers' Web pages on the Internet.

Black-and-White Film

Black-and-white (B&W) films and papers are sold through professional product supply chains. Color films and papers have largely displaced B&W for amateur use, with the exception of advanced amateurs who occasionally shoot it for artistic expression.

The image on a B&W film or print is composed of black metallic silver. This image has archival stability under proper storage conditions and is quite stable even under uncontrolled storage. Many remarkable and historically important B&W images exist that date from 50 to over 100 years ago. This is remarkable considering the primitive state of the technology and haphazard storage conditions.

Pioneer nature photographer Ansel Adams reprinted many of his stored negatives long after he stopped capturing images in the field. This artistic genius recreated new expressions in prints from scenes captured many years earlier in his career because his negatives were well preserved.

Ansel Adams worked with B&W film and paper. Many photographic artists work in B&W when color would distract from the artistic expression they wish to convey. For example, many portraits are imaged in B&W.

B&W films pleasingly image an extended range of tones from bright light to deep shadow. They are panchromatically sensitized to render colors into grayscale tones. A B&W film's tone scale is among its most important characteristics.

Photographers manipulate the contrast of a B&W image by push- or pull-processing the negative (varying the length of time in the developer solution) or by printing the negative onto a select contrast paper. Paper contrast grades are indexed from 1 through 5, with the higher index giving higher contrast to the print. Multigrade paper is also available whereby the contrast in the paper is changed by light filtration at the enlarger.

It is generally best to capture the most desired contrast on the negative because highlight or shadow information lost on the negative cannot be recovered at any later printing stage. Paper grades are designed for pleasing artistic effect, not for highlight or shadow recovery.

Legacy films are very important products in B&W photography. Although typically a bit grainier and less sharp than modern B&W films, their tone scale characteristics, process robustness, and overall image rendition have stood the test of time and are especially favored by many photographers. Besides, the B&W photographic market is small so improvements to this line of films occur far less frequently than do improvements to color films.

Unlike the case with color films, which are designed for a single process, a photographer may choose among several developers for B&W films. Kodak Professional T-Max developer is a good choice for the relatively modern Kodak T-Max B&W film line. Kodak developer D-76 is also very popular and will render an image with slightly less grain and slightly less speed. It is a popular choice for legacy films. Kodak Microdol-X developer is formulated to give very low grain at the expense of noticeable speed loss. Kodak developer HC110 is popular for home darkrooms because of its low cost and ease of use, but it renders a grainier image compared to other developers. Other manufacturers, notably Ilford Ltd., carry similar types of B&W developers.

A new type of B&W film has emerged that is developed in a standard color negative process. This incorporated-coupler 400-speed film forms a black-and-white image from chemical dyes instead of metallic silver. Among its advantages are very fine grain at normal exposures in the commonly available color negative process. Its shortcomings are the inability to control contrast on the negative by push or pull processing and its objectionable high grain when underexposed compared to a comparably underexposed 400-speed legacy film. And many photographers prefer the tone scale characteristics found in standard B&W films.

Table 2 lists some B&W films available today. This is not a comprehensive list. Modern films generally feature improved grain and sharpness compared to legacy films. However, the legacy films are favorite choices among many photographers because of their forgiving process insensitivity, the

TABLE 2 Black and White Films

Modern Films	Legacy Films	Specialty Films
Kodak T-Max 100 Professional (fine grain, high sharpness)	Ilford Pan F Plus 50	Kodak Technical Pan 25 (fine grain and very high sharpness)
Ilford Delta 100 Pro (fine grain, high sharpness)	Kodak Plus-X 125	Kodak IR (infrared sensitive)
Kodak T-Max 400 Professional (finer grain compared to legacy films)	Ilford FP4 Plus 125	Ilford SFX 200 (extended long red sensitivity but not IR)
Ilford Delta 400 Pro (finer grain compared to legacy films)	Kodak Tri-X 400	Ilford Chromogenic 400 (incorporated couplers, color negative process)
Fuji Neopan 400 (finer grain compared to legacy films)	Ilford HP5 Plus 400	Kodak Professional T400 CN (incorporated couplers, color negative process)
Kodak T-Max P3200 Professional (push process to high speed, high grain)		
Ilford Delta P3200 Pro (push process to high speed, high grain)		

experience factor of knowing their detailed behavior under many conditions, and the characteristic look of their images. More detailed information about these films' uses and characteristics can be found in the film manufacturers' Web pages on the Internet.

Color Negative Film

Today's color negative film market is highly segmented. It can be most broadly divided into consumer films and professional films. The consumer line is further segmented into 35-mm film, 110 format film (a minor segment whose image size is 17×13 mm), single-use cameras, and the new Advanced Photo System (APS). Each segment serves the needs of amateur photographers in different ways. The basic films are common across segments, with the main difference being camera type.

Consumer color negative film is a highly competitive market with many film and camera manufacturers offering products. In general, films for this marketplace offer bright color saturation that is most pleasing to amateurs. The higher-speed films show progressively more grain than do the lower-speed films. The highest-speed films are less often used when enlargement becomes significant, such as in APS and 110 format, but are prevalent in single-use cameras.

Consumer 35-mm Films Films in 35-mm format form the core of all consumer color negative products. These same basic films are also found in single-use cameras, 110 format cameras, and APS cameras. Film speeds include 100, 200, 400, and 800. All manufacturers offer films at 100, 200, and 400 speed, but only a few offer films at speeds of 800 and above.

The 400-speed film is most popular for indoor shots under low light and flash conditions. Lower-speed films are most often used outdoors when light is plentiful.

The best consumer films feature technology for optimized color rendition including the high color saturation pleasing to most consumers, plus accurate spectral sensitization for realistic color rendition under mixed lighting conditions of daylight plus fluorescent and incandescent light, which is often present inside homes. Technology used for accurate spectral sensitization was briefly described for Fujicolor Superia and Kodak Gold films in Sec. 30.5.

Kodak consumer films segment into Kodak Gold and Kodak Royal Gold films. Kodak Gold consumer films emphasize colorfulness, while Kodak Royal Gold films emphasize fine grain and high sharpness.

Films made in 110 format generally fall into the 100- to 200-speed range because the enlargement factors needed to make standard-size prints place a premium on fine grain and high sharpness.

Single-Use Cameras The single-use camera marketplace is extremely competitive because of the high popularity of this film and camera system among amateurs. These low-cost units can be bought at all consumer outlets and are especially popular at amusement parks, theme parks, zoos, and similar family attractions. After the film is exposed, the photographer returns the entire film plus camera unit for processing. Prints and negatives are returned to the customer while the camera body is returned to the manufacturer, repaired as needed, reloaded with film, and repackaged for sale again. These camera units are recycled numerous times from all over the world.

Single-use cameras come in a variety of styles including a waterproof underwater system, low-cost models with no flash, regular flash models, an APS style offering different picture sizes, and a panoramic style. Single-use cameras typically contain 800-speed film, except for APS and panoramic models, which usually contain 400-speed film due to enlargement demands for finer grain.

Advanced Photo System (APS) APS is the newest entry into consumer picture taking. The size of the image on the negative is 29×17 mm, about 60 percent of the image size of standard 35-mm film (image size 36×24 mm). This puts a premium on fine grain and high sharpness for any film used in this system. Not all film manufacturers offer it. Speeds offered are 100, 200, and 400, with the higher speeds being the most popular.

Some manufacturers offer a B&W film in APS format. The film used is the incorporated-coupler 400-speed film that forms a black-and-white image from chemical dyes and is developed in a standard color negative process.

The APS camera system features simple drop-in cassette loading. Unlike the 35-mm cassette, the APS cassette contains no film leader to thread into the camera. When loaded, the camera advances the film out of the cassette and automatically spools it into the camera, making this operation mistake-proof. Three popular print sizes are available; panoramic, classic, and high-definition television (HDTV) formats. These print sizes differ in their width dimension.

A thin, nearly transparent magnetic layer is coated on the APS film's plastic support. Digital information recorded on this layer, including exposure format for proper print size plus exposure date and time, can be printed on the back of each print. Higher-end cameras offer prerecorded titles—for example, "Happy Birthday"—that can be selected from a menu and placed on the back of each print. Additional capabilities are beginning to emerge that take more advantage of this magnetic layer and the digital information it may contain.

The APS system offers easy mid-roll change on higher-end cameras. With this feature, the last exposed frame is magnetically indexed so the camera "remembers" its position on the film roll. A user may rewind an unfinished cassette to replace it with a different cassette (different film speed, for example). The original cassette may be reloaded into the camera at a later time. The camera will automatically advance the film to the next available unexposed frame.

Markings on the cassette indicate whether the roll is unexposed, partially exposed, fully exposed, or developed. There is no uncertainty about whether a cassette has been exposed or not. Unlike the 35-mm system, where the negatives are returned as cut film, the negatives in the APS system are rewound into the cassette and returned to the customer for compact storage. An index print is given to the customer with each processed roll so that each numbered frame in the cassette can be seen at a glance when reprints are needed.

Compact film scanners are available to digitize pictures directly from an APS cassette or 35-mm cut negatives. These digitized images can be uploaded into a computer for the full range of digital manipulations offered by modern photo software. Pictures can be sent over the Internet, or prints can be made on inkjet printers. For the best photo-quality prints, special photographic ink cartridge assemblies are available with most inkjet printers to print onto special high-gloss photographic-quality paper.

Table 3 summarizes in alphabetical order many brands of 35-mm consumer film available today worldwide. Some of these same basic films appear in single-use cameras, APS format, and 110 format, although not all manufacturers listed offer these formats. This is not a comprehensive list. Because this market is highly competitive, new films emerge quickly to replace existing films. It is not unusual for a manufacturer's entire line of consumer films to change within three years. More detailed information about these films' uses and characteristics can be found in the film manufacturers' Web pages on the Internet.

TABLE 3 Consumer 35-mm Color Negative Films

Manufacturer and Brand Name					
Agfacolor HDC Plus	100	200	400		
Fujicolor Superia	100	200	400	800	
Ilford Colors	100	200	400		
Imation HP	100	200	400		
Kodak Gold	100	200	Max400	Max800	
Kodak Royal Gold	100	200	400		1000
Konica Color Centuria	100	200	400	800	
Polaroid One Film	100	200	400		

Professional Color Negative Film Professional color negative film divides roughly into portrait and wedding film and commercial and photojournalism film. Portrait and wedding films feature excellent skin tones plus good neutral whites, blacks, and grays. These films also incorporate accurate spectral sensitization technology for color accuracy and excellent results under mixed lighting conditions. The contrast in these films is about 10 percent lower than in most consumer color negative films for softer, more pleasing skin tones, and its tone scale captures highlight and shadow detail very well.

The most popular professional format is 120 size, although 35-mm and sheet sizes are also available. Kodak offers natural color (NC) and vivid color (VC) versions of Professional Portra film, with the vivid color having a 10 percent contrast increase for colorfulness and a sharper look. The natural color version is most pleasing for pictures having large areas of skin tones, as in head and shoulder portraits.

Commercial films emphasize low grain and high sharpness. These films offer color and contrast similar to those of consumer films. Film speeds of 400 and above are especially popular for photojournalist applications, and the most popular format is 35 mm. These films are often push processed.

Table 4 summarizes in alphabetical order many brands of professional color negative film available today worldwide. This is not a comprehensive list. Because this market is highly competitive, new films emerge quickly to replace existing films. More detailed information about these films' uses and characteristics can be found in the film manufacturers' Web pages on the Internet.

Silver halide photographic products have enjoyed steady progress during the past century. This chapter has described most of the technology that led to modern films. Most noteworthy among these advances are

1. Efficient light management in multilayer photographic materials has reduced harmful optical effects that caused sharpness loss and has optimized beneficial optical effects that lead to efficient light absorption.
2. Proprietary design of silver halide crystal morphology, internally structured halide types, transition metal dopants to manage the electron-hole pair, and improved chemical surface treatments has optimized the efficiency of latent image formation.

TABLE 4 Professional Color Negative Films

Manufacturer and Brand Name (Commercial Film)					
Agfacolor Optima II Prestige	100	200	400		
Fujicolor Press			400	800	
Kodak Professional Ektapress	PJ100		PJ400	PJ800	
Konica Impresa	100	200			3200 SRG
Manufacturer and Brand Name (Portrait Film)					
Agfacolor Portrait		160 XPS			
Fujicolor		160 NPS	400 NPH	800 NHGII	
Kodak Professional Portra		160 NC	400 NC		Pro 1000
		160 VC	400 VC		
Konica Color Professional		160			

3. Transition metal dopants and development-modifying chemistry have improved process robustness, push processing, and exposure time latitude.
4. New spectral sensitizers combined with powerful chemical color correction methods have led to accurate and pleasing color reproduction.
5. New image dyes have provided rich color saturation and vastly improved image permanence.
6. New film and camera systems for consumers have made the picture-taking experience easier and more reliable than ever before.

Photographic manufacturers continue to invest research and product development resources in their silver halide photographic products. Although digital electronic imaging options continue to emerge, consumers and professionals can expect a constant stream of improved silver halide photographic products for many years to come.

30.7 REFERENCES

1. M. Bass, E. Van Stryland, D. Williams, and W. Wolfe (eds.), *Handbook of Optics*, vol. 1, 2d ed., McGraw-Hill, New York, 1995.
2. J. Kapecki and J. Rodgers, *Kirk-Othmer Encyclopedia of Chemical Technology*, vol. 6, 4th ed., John Wiley & Sons, New York, 1993.
3. T. H. James (ed.), *The Theory of the Photographic Process*, 4th ed., Macmillan, New York, 1977.
4. R. W. G. Hunt, *The Reproduction of Colour*, Fountain Press, Tolworth, England, 1987.
5. E. Klein and H. J. Metz, *Photogr. Sci. Eng.* **5**:5 (1961).
6. R. E. Factor and D. R. Diehl, U. S. Patent 4,940,654, 1990.
7. R. E. Factor and D. R. Diehl, U. S. Patent 4,855,221, 1989.
8. G. Mie, *Ann. Physik.* **25**:337 (1908); M. Kerker, *The Scattering of Light and Other Electromagnetic Radiation*, Academic Press, New York, 1969.
9. D. H. Napper and R. H. Ottewill, *J. Photogr. Sci.* **11**:84 (1963); *J. Colloid Sci.* **18**:262 (1963); *Trans. Faraday Soc.* **60**:1466 (1964).
10. E. Pitts, *Proc. Phys. Soc. Lond.* **67B**:105 (1954).
11. J. M. Hammersley and D. C. Handscomb, *Monte Carlo Methods*, John Wiley & Sons, New York, 1964.
12. J. J. DePalma and J. Gasper, *Photogr. Sci. Eng.* **16**:181 (1972).
13. M. Motoki, S. Ichijima, N. Saito, T. Kamio, and K. Mihayashi, U. S. Patent 5,213,958, 1993.
14. Miles V. Klein, *Optics*, John Wiley & Sons, New York, 1970, pp. 582–585.
15. J. Gasper, Eastman Kodak Company, unpublished data.
16. N. F. Mott and R. W. Gurney, *Electronic Processes in Ionic Crystals*, Oxford University Press, London, 1940.
17. F. Seitz, *Rev. Mod. Phys.* **23**:328 (1951).
18. J. F. Hamilton, *Adv. Phys.* **37**:359 (1988).
19. A. P. Marchetti and R. S. Eachus, *Advances in Photochemistry*, vol. 17, John Wiley & Sons, New York, 1992, p. 145.
20. R. S. Eachus, A. P. Marchetti, and A. A. Muentner, *Ann. Rev. Phys. Chem.* **50**:117 (1999).
21. R. J. Deri, J. P. Spoonhower, and J. F. Hamilton, *J. Appl. Phys.* **57**:1968 (1985).
22. T. Tani, *Photographic Sensitivity*, Oxford University Press, Oxford, UK, 1995, p. 235.
23. L. M. Slifkin and S. K. Wonnell, *Solid State Ionics* **75**:101 (1995).
24. J. F. Hamilton and L. E. Brady, *Surf. Sci.* **23**:389 (1970).
25. R. C. Baetzold, Y. T. Tan, and P. W. Tasker, *Surf. Sci.* **195**:579 (1988).
26. R. S. Eachus and M. T. Olm, *Annu. Rep. Prog. Chem., Sect. C Phys. Chem.* **83**:3 (1986).
27. J. C. Dainty and R. Shaw, *Image Science*, Academic Press, London, 1974.

28. R. K. Hailstone, N. B. Liebert, M. Levy, R. T. McCleary, S. R. Girolmo, D. L. Jeanmaire, and C. R. Boda, *J. Imaging Sci.* **3**(3) (1988).
29. J. D. Baloga, "Factors in Modern Color Reversal Films," *IS&T 1998 PICS Conference*, p. 299 (1998).
30. R. J. Tuite, *J. Appl. Photogr. Eng.* **5**:200 (1979).
31. W. F. Smith, W. G. Herkstroeter, and K. I. Eddy, *Photogr. Sci. Eng.* **20**:140 (1976).
32. P. Douglas, *J. Photogr. Sci.* **36**:83 (1988).
33. P. Douglas, S. M. Townsend, P. J. Booth, B. Crystall, J. R. Durrant, and D. R. Klug, *J. Chem. Soc. Faraday Trans.* **87**:3479 (1991).
34. F. Wilkinson, D. R. Worrall, and R. S. Chittock, *Chem. Phys. Lett.* **174**:416 (1990).
35. F. Wilkinson, D. R. Worrall, D. McGarvy, A. Goodwin, and A. Langley, *J. Chem. Soc. Faraday Trans.* **89**:2385 (1993).
36. P. Douglas, S. M. Townsend, and R. Ratcliffe, *J. Imaging Sci.* **35**:211 (1991).
37. P. Egerton, J. Goddard, G. Hawkins, and T. Wear, *Royal Photographic Society Color Imaging Symposium*, Cambridge, UK, September 1986, p. 128.
38. K. Onodera, T. Nishijima, and M. Sasaki, *Proceedings of the International Symposium on the Stability and Conservation of Photographic Images*, Bangkok, Thailand, 1986.
39. Y. Kaneko, H. Kita, and H. Sato, *Proceedings of IS & T's 46th Annual Conference*, 1993, p. 299.
40. R. J. Berry, P. Douglas, M. S. Garley, D. Clarke, and C. J. Winscom, "Photophysics and Photochemistry of Azomethine Dyes," *IS & T 1998 PICS Conference*, p. 282 (1998).
41. W. F. Smith, W. G. Herkstroeter, and K. I. Eddy, *J. Am. Chem. Soc.* **97**:2164 (1975).
42. W. G. Herkstroeter, *J. Am. Chem. Soc.* **95**:8686 (1973).
43. W. G. Herkstroeter, *J. Am. Chem. Soc.* **97**:3090 (1975).
44. P. Douglas and D. Clarke, *J. Chem. Soc. Perkin Trans.* **2**:1363 (1991).
45. F. Abu-Hasanayn, *Book of Abstracts, 218th ACS National Meeting*, New Orleans, LA, August 22–26, 1999.
46. K. Hidetoshi et al. *International Congress of Photographic Science meeting*, Belgium, 1998.
47. N. Saito and S. Ichijima, *International Symposium on Silver Halide Imaging*, 1997.
48. V. Balzani, F. Bolletta, and F. Scandola, *J. Am. Chem. Soc.* **102**:2152 (1980).
49. R. Jain and W. R. Schleigh, U. S. Patent 5,561,037, 1996.
50. O. Takahashi, H. Yoneyama, K. Aoki, and K. Furuya, "The Effect of Polymeric Addenda on Dark Fading Stability of Cyan Indoaniline Dye," *IS & T 1998 PICS Conference*, p. 329 (1998).
51. R. L. Heidke, L. H. Feldman, and C. C. Bard, *J. Imag. Tech.* **11**(3):93 (1985).
52. S. Cowan and S. Krishnamurthy, U. S. Patent 5,378,587 (1995).
53. T. Kawagishi, M. Motoki, and T. Nakamine, U. S. Patent 5,605,788 (1997).
54. H. W. Vogel, *Berichte* **6**:1302 (1873).
55. J. E. Maskasky, *Langmuir* **7**:407 (1991).
56. J. Spence and B. H. Carroll, *J. Phys. Colloid Chem.* **52**:1090 (1948).
57. A. H. Hertz, R. Danner, and G. Janusonis, *Adv. Colloid Interface Sci.* **8**:237 (1977).
58. M. Kawasaki and H. Ishii, *J. Imaging Sci. Technol.* **39**:210 (1995).
59. G. Janssens, J. Gerritsen, H. van Kempen, P. Callant, G. Deroover, and D. Vandenbroucke, *The Structure of H-, J-, and Herringbone Aggregates of Cyanine Dyes on AgBr(111) Surfaces*. Presented at ICPS 98 International Conference on Imaging Science, Antwerp, Belgium (1998).
60. P. B. Gilman, *Photogr. Sci. Eng.* **18**:475 (1974).
61. J. Lenhard, *J. Imaging Sci.* **30**:27 (1986).
62. W. West, "Scientific Photography," in *Proceedings of the International Conference at Liege, 1959*, H. Sauvenier (ed), Pergamon Press, New York, 1962, p. 557.
63. T. Tani, *Photogr. Sci. Eng.* **14**:237 (1970).
64. J. Eggert, W. Meidinger, and H. Arens, *Helv. Chim. Acta.* **31**:1163 (1948).
65. J. M. Lanzafame, A. A. Muentner, and D. V. Brumbaugh, *Chem. Phys.* **210**:79 (1996).

66. A. A. Muentner and W. Cooper, *Photogr. Sci. Eng.* **20**:121 (1976).
67. W. West, B. H. Carroll, and D. H. Whitcomb, *J. Phys. Chem.* **56**:1054 (1952).
68. R. Brunner, A. E. Oberth, G. Pick, and G. Scheibe, *Z. Elektrochem.* **62**:146 (1958).
69. P. B. Gilman, *Photogr. Sci. Eng.* **11**:222 (1967).
70. P. B. Gilman, *Photogr. Sci. Eng.* **12**:230 (1968).
71. J. E. Jones and P. B. Gilman, *Photogr. Sci. Eng.* **17**:367 (1973).
72. P. B. Gilman and T. D. Koszelak, *J. Photogr. Sci.* **21**:53 (1973).
73. H. Sakai and S. Baba, *Bull. Soc. Sci. Photogr. Jpn.* **17**:12 (1967).
74. B. H. Carroll, *Photogr. Sci. Eng.* **5**:65 (1961).
75. T. Tani, *Photogr. Sci. Eng.* **15**:384 (1971).
76. T. A. Babcock, P. M. Ferguson, W. C. Lewis, and T. H. James, *Photogr. Sci. Eng.* **19**:49 (1975).
77. E. J. Wall, *History of Three Color Photography*, American Photographic Publishing Company, Boston, MA, 1925.
78. A. C. Hardy and F. L. Wurzburg Jr., "The Theory of Three Color Reproduction," *J. Opt. Soc. Am.* **27**:227 (1937).
79. M. L. Pearson and J. A. C. Yule, *J. Color Appearance* **2**:30 (1973).
80. S. G. Link, "Short Red Spectral Sensitizations for Color Negative Films," *IS & T 1998 PICS Conference*, p. 308 (1998).
81. Y. Nozawa and N. Sasaki, U. S. Patent 4,663,271 (1987).

IMAGE TUBE INTENSIFIED ELECTRONIC IMAGING

C. Bruce Johnson

*Johnson Scientific Group Inc.
Phoenix, Arizona*

Larry D. Owen

*NuOptics International
Phoenix, Arizona*

31.1 GLOSSARY

B_s	phosphor screen brightness, photometric units
CCDs	charge-coupled devices
CIDs	charge-injection devices
E_i	image plane illuminance, lux
E_s	scene illuminance, lux
e	electronic charge, coulombs
FO	fiberoptic
FOV	field-of-view, degrees
fc	illuminance, photometric, foot candles = lm/ft^2
f_N	spatial Nyquist frequency, cycle/mm
f_{fto}	limiting resolution at fiberoptic taper output
F_{si}	input window signal flux
ftL	luminance, photometric (brightness), foot Lamberts = lm/ft^2
G_m	VMCP electron gain, e/e
HVPS	high-voltage power supply
II	image intensifier
LLL	low-light-level
lx	illuminance, photometric, lux = lm/m^2
M_{fot}	magnification of fiberoptic taper
MCP	microchannel plate
MTF	modulation transfer function, 0 to 1.0
N_{essa}	number of stored SSA electrons per input photoelectron, e/photon
N_f	total number of frames, #
N_p	number of photoelectrons, #

$N_{ps}(\lambda)$	number of photons per second, photon/s
PDA	photodiode arrays
P	phosphor screen efficiency, photon/eV
$P_p(\lambda)$	radiometric power spectral distribution, W
QLI	quantum limited imaging
Q_{ssa}	stored SSA charge per input photoelectron from the photocathode, C
R_s	scene reflectance, ratio
R_{sn}	signal-to-noise ratio, ratio
$S(\lambda)$	absolute spectral sensitivity, mA/W
$S(f)$	squarewave response versus frequency, cycles/mm
SIT	silicon-intensifier-target vidicon
SNR	signal-to-noise ratio
sb	luminance, photometric (brightness), stilbs = cd/cm ²
SSA	silicon self-scanned array
T_f	filter transmission, 0 to 1.0
T_{fot}	transmission of fiberoptic taper, 0 to 1.0
T_n	lens T-number = $FN/\sqrt{\tau_0}$
T_{ssa}	transmission of fiberoptic window on the SSA, 0 to 1.0
V_a	phosphor screen, actual applied voltage, V
V_d	phosphor screen, "dead-voltage," V
V_m	VMCP applied potential, V
V_s	MCP-to-screen applied potential, V
$Y(\lambda)$	quantum yield (electrons/photon), percent
Y_k	quantum yield, photoelectrons/photon
Y_{ssa}	SSA quantum yield, e/photon
τ_e	the exposure period, s
τ_i	CCD charge integration period, s
τ_o	lens transmission, 0 to 1.0
Φ_p	photon flux density, photon/m ² /s

31.2 INTRODUCTION

It is appropriate to begin our discussion of image tube intensified (II) electronic imaging with a brief review of natural illumination levels. Figure 1 illustrates several features of natural illumination in the range from full sunlight to overcast night sky conditions. Various radiometric and photometric illuminance scales are shown in this figure. Present silicon self-scanned array (SSA) TV cameras, having frame rates of 1/30 to 1/25 s, operate down to about 0.5 lx minimum illumination.

The generic term *self-scanned array* is used here to denote any one of several types of silicon solid-state sensors available today which are designed for optical input. Among these are charge-coupled devices (CCDs), charge-injection devices (CIDs), and photodiode arrays (PDAs). Vol. II, Chaps. 32, "Visible Array Detectors," and 33, "Infrared Detector Arrays," contain detailed information on these types of optical imaging detectors. Specially designed low-light-level (LLL) TV cameras making use of some type of image intensifier must be used for lower exposures, i.e., lower illumination and/or shorter exposures.

The fundamental reason for using an II SSA camera instead of a conventional SSA camera is that low-exposure applications require the low-noise optical image amplification provided by an II

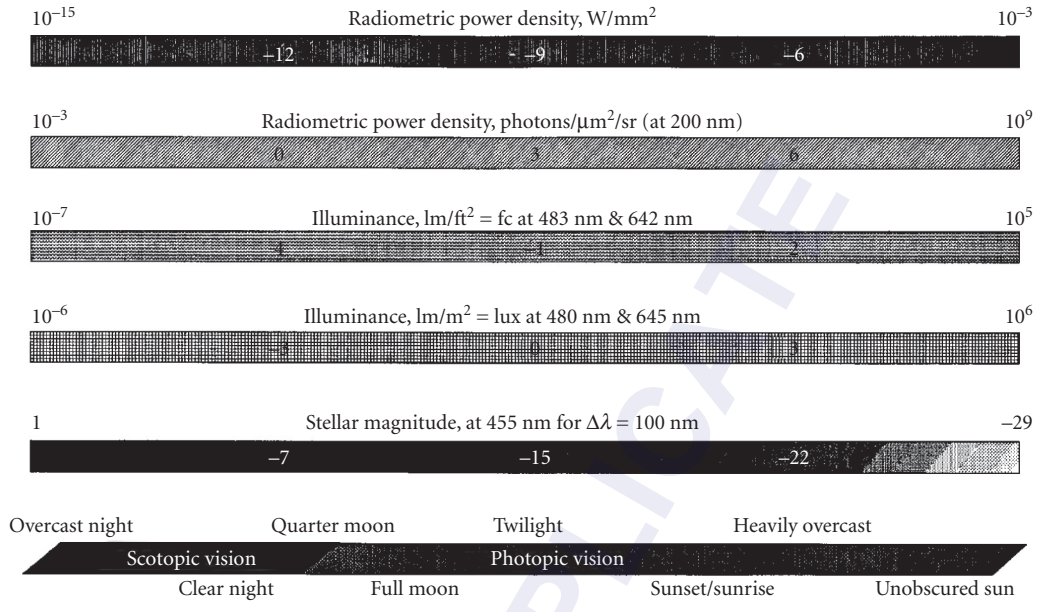


FIGURE 1 Various optical illumination ranges.

to produce a good signal-to-noise ratio from the SSA camera. Other important applications arise because of the ability to electronically shutter IIs as fast as 1 ns or less and the higher sensitivity of IIs in certain spectral regions. The following sections deal with the optical interface between the object and the II SSA, microchannel plate proximity-focused IIs, and II SSA detector assemblies. By using auto-iris lenses and controlling both the electronic gain and gating conditions of the II, II SSA cameras can provide an interscene dynamic range covering the full range of twelve orders of magnitude shown in Fig. 1. Several applications for II SSAs are discussed later in the chapter under “Applications.”

31.3 THE OPTICAL INTERFACE

It is necessary to begin our analysis of II SSA cameras with a brief discussion of the various ways to quantify optical input and exposure. Two fundamental systems are used to specify input illumination: radiometric and photometric. These systems are briefly described, and the fundamentals of optical image transfer are discussed. Detailed aspects of radiometry, photometry, and optical image transfer are discussed in Vol. II, Chap. 34, “Radiometry and Photometry” and Vol. I, Chap. 4, “Transfer Function Techniques.” However, enough information is presented in this chapter to allow the reader to properly design, analyze, and apply II SSA imaging technology for a wide variety of practical applications.

Quantum Limited Imaging Conditions

Quantum limited imaging (QLI) conditions exist in a wide variety of applications. An obvious one is that of LLL TV imaging at standard frame rates, i.e., 33-ms exposure periods, under nighttime illumination conditions. For example, under full moonlight input faceplate illumination conditions,

only ~ 1000 photons enter a $10 \times 10 \mu\text{m}^2$ image pixel in a 33-ms frame period. Assuming a quantum yield of 10 percent, an average of only 100 electrons is generated, and the maximum SNR achievable in each pixel and each frame is only $\sqrt{100} = 10$. Alternatively, under full unobscured sunlight input faceplate illumination conditions, an electronically gated camera with gatewidth limited exposure period of 10 ns produces a total of $(1\text{E}9 \text{ photons}/\mu\text{m}^2/\text{s}) (10 \lambda \times 10 \mu\text{m}^2)(10 \text{ ns}) = 1000$ photons, or the same SNR as for the LLL operating conditions noted above. These are both clearly QLI operating conditions. II SSA camera technology is used to obtain useful performance in both of these types of applications. Without the use of an II, a bare SSA does not meet the requirements for useful SNR under these conditions.

Radiometry

The unit of light flux in the radiometric system is the watt. The watt can be used anywhere in the optical spectrum to give the number of photons per second (N_{ps}) as a function of wavelength (λ). Since the photon energy $E_p(\lambda)$ is

$$E_p = \frac{hc}{\lambda} \quad (1)$$

where h is Planck's constant and c is the velocity of light in vacuum, the radiometric power $P_p(\lambda)$, in watts, is given by

$$P_p(\lambda) = \left(\frac{hc}{\lambda} \right) \cdot N_{\text{ps}}(\lambda) \quad (2)$$

or

$$P_p(\lambda) = (2 \cdot 10^{-25}) \cdot \frac{N_{\text{ps}}(\lambda)}{\lambda} \quad (3)$$

where N_{ps} is the number of photons per second. Alternatively, the photon rate is given by

$$N_{\text{ps}}(\lambda) = (5 \times 10^{24}) \lambda P_p(\lambda) \quad \text{photons/s} \quad (4)$$

For example, one milliwatt of 633-nm radiation from an He-Ne laser is equivalent to $(5\text{E}24)(633\text{E}-9)(1\text{E}-3) = 3.2\text{E}15$ photons/s.

Radiometric flux density, in W/m^2 , represents a photon rate per unit area, and radiometric exposure per unit area is the product of the flux density times the exposure period. The active surface of a photoelectronic detector produces a current density in response to an optical flux density input, while a total signal charge is produced per unit area in the same detector during a given exposure period.

Rose¹ has shown that all types of optical detectors, e.g., photographic, electronic, or the eye, are subject to the same fundamental limits in terms of signal-to-noise ratio (R_{sn}), optical input, and exposure period. In summary, the noise in a measured signal of N_p photoelectrons during a fixed exposure period is $\sqrt{N_p}$, so that

$$R_{\text{sn}} = \sqrt{N_p} \quad (5)$$

The brightness (B_s) of a scene that produces this signal in a square pixel of dimensions ($y \cdot y$), as a result of the optical transfer and conversion from the source to the detector, possibly through a medium that absorbs, scatters, and focuses photons, is

$$B_s = \frac{C \cdot N_p}{y^2} \quad (6)$$

where C is a constant. In terms of signal-to-noise ratio,

$$B_s = \frac{C \cdot R_{sn}^2}{y^2} \quad (7)$$

Thus, for twice the signal-to-noise ratio, the scene brightness must be increased four times, or the throughput of the optical system must be quadrupled, etc. Also, if the pixel size is reduced by a factor of two, the same changes in scene brightness or optical throughput must be made in order to maintain the same signal-to-noise ratio. Under QLI conditions, higher resolution necessarily requires more input flux density for equal signal-to-noise ratio, and higher resolution inherently implies less sensitivity. The Rose limit should be used often as a proof check on design and performance estimates of LLL and other QLI imaging systems.

As an example, assume a simple imaging situation such as a single pixel, e.g., a star in the night-time sky, and an II SSA camera having an objective lens of diameter D_o . Also assume that the starlight is filtered, to observe only a narrow wavelength band, and that the photon flux density from the star is Φ_p (photon/m²/s). The number of photoelectrons produced at the photocathode of the II SSA detector (N_p) is given by

$$N_p = \Phi_p \cdot T_f \cdot \left(\frac{\pi D_o^2}{4} \right) \cdot \tau_o \cdot Y_k \cdot \tau_c \quad (8)$$

where T_f is the filter transmission, τ_o is the lens transmission, Y_k is the quantum yield of the window/photocathode assembly in the II SSA camera, and τ_c is the exposure period. Note that the II SSA camera parameters which determine the rate of production of signal photoelectrons are filter transmission, lens diameter, quantum yield, and exposure period. The key one is of course the lens diameter, and not lens f -number, for this kind of imaging; it is important, however, for extended sources such as terrestrial scenes.

Photometry and the Camera Lens

A lens on the II SSA camera is used to image a scene onto the input window/photocathode assembly of the II SSA. The relationship between the scene (E_s) and II SSA image plane (E_i) illuminances in lux (lx) is

$$E_i = \frac{\pi \cdot E_s \cdot R_s \cdot \tau_o}{(4 \cdot FN^2 \cdot (m+1)^2)} \quad (9)$$

where R_s is the scene reflectance, τ_o is the optical transmission of the lens, FN is the lens f -number, and m is the scene-to-image magnification. If E_s is in foot-lamberts, then the π is dropped and E_i is in footcandles.

Alternatively, Eq. (9) becomes

$$E_i = \frac{E_s \cdot R_s}{(4 \cdot T_n^2 \cdot (m+1)^2)} \quad (10)$$

using the T-number of the lens, where

$$T_n = \frac{FN}{\sqrt{\tau_o}} \quad (11)$$

The sensitivity of an II is usually given in two forms, i.e., “white-light” luminous sensitivity, in units of $\mu\text{A}/\text{lm}$, and absolute spectral sensitivity, in units of A/W as a function of wavelength, as discussed later in the section “Input Window/Photocathode Assemblies” in Sec. 31.4.

Example: A scene having an average reflectance of 50 percent receives LLL “full-moon” illumination of $1.0\text{E} - 2$ fc. If a lens having a T-number of 3.0 is used, and the scene is at a distance of 100 m

from a lens with a focal length of 30 mm, what is the input illumination at the II SSA? Since the distance to the scene is much longer than the focal length of the lens, the magnification is much smaller than unity and m can be neglected. Thus,

$$E_i = \frac{E_s R_s}{(4 \cdot T_n^2)} \quad (12)$$

For the given values, the input illumination at the II SSA is bound to be $E_i = (1.0E - 2 \text{ fc}) / (0.50) / (4(3.0)^2) = 1.4E - 4 \text{ fc}$.

General Considerations

It is of prime importance in any optoelectronic system to couple the maximum amount of signal input light into the primary detector surface, e.g., the window/photocathode assembly of an II SSA. In order to achieve the maximum signal-to-noise ratio, the modulation transfer function of the input optic and the spectral sensitivity of the II SSA must be carefully chosen. As shown in Fig. 2, the spectral sensitivity of a silicon SSA is much different than that of a Gen-3 image intensifier tube.

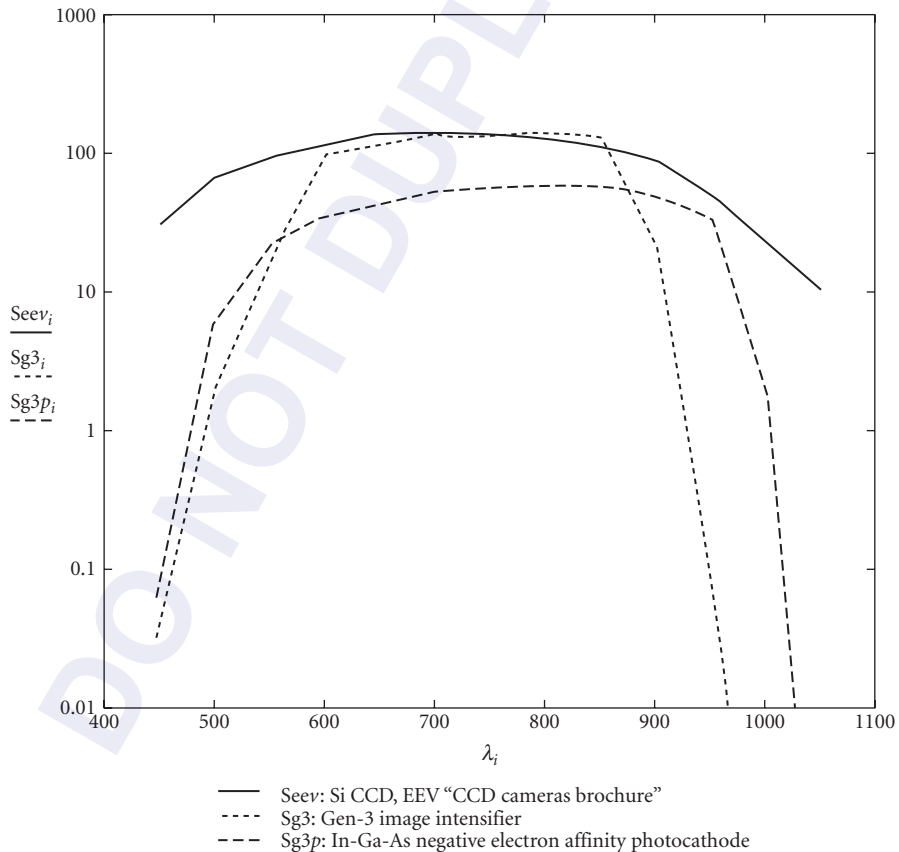


FIGURE 2 Absolute spectral sensitivity S (mA/W) versus wavelength λ (nm) of a frame-transfer type of CCD, a gen-III image intensifier, and an II having an In-Ga-As negative electron affinity photocathode.

Thus, an optimized objective lens design for a CCD will be much different than that for an II SSA. The dynamic range characteristics are also very different, since IIs will handle seven orders-of-magnitude interscene dynamic range, using a combination of II gain control and electronic duty-cycle gating, while SSAs will only provide about two orders of magnitude.²

Several factors must be considered if the overall system resolution and sensitivity are to be optimized. For example, the spectral responses of many optical input SSAs and/or lenses used in commercial cameras have been modified by using filters to reduce the red and near-ir responses to give more natural flesh tones. In an II SSA the filter may have little effect if the filter is on the SSA. The filter should not be used in the objective lens for the II SSA since a major portion of the signal will be filtered out. If a color SSA is to be used in an intensified system using relay lens coupling, sacrifice of both sensitivity and resolution will result. This is due to the matrix color filter used in these SSA chip designs. Most of the signal will go into green bandpass filter elements, and very little will go into the blue and red elements. The color matrix filter is usually bonded to the surface of the SSA chip; thus these SSA types are not used for fiberoptically coupled II SSAs.

The ideal objective lens design for an II SSA needs to be optically corrected over the spectral range of sensitivity of the II and the spectral range of interest. For special-purpose photosensitivity covering portions of the uv, blue, or near-ir spectral regions, appropriate adjustments must be made in the lens design. Although they may be adequate for many applications, it is very seldom that a commercial CCTV lens is optimized for nighttime illumination, or other LLL or QLI conditions.

Another very important part of an optimized II SSA camera design is to make the proper choice of II and SSA formats. This subject is discussed in detail later under "Fiberoptic-Coupled II/SSAs," under Sec. 31.5. The input of the II SSA system is the II, and the most likely choice will be one with an 18-mm active diameter, since the widest choice of II features is available in this size. Image intensifiers are also available having 25- and 12-mm active diameters, but these are generally more expensive. Regarding the SSA standard format sizes, the standard commercial TV formats are named by a longtime carryover from the days when vidicons were used extensively. Thus 2/3-, 1/2-, and 1/3-in format sizes originally referred to the diameters of the vidicon envelope and not the actual image format.

31.4 IMAGE INTENSIFIERS

An image intensifier (II) module, when properly coupled to an SSA camera, produces a low-light-level electronic imaging capability that is extremely useful across a broad range of application areas, including spectral analysis, medical imaging, military cameras, nighttime surveillance, high-speed optical framing cameras, and astronomy. An immediate advantage of using an II is that its absolute spectral sensitivity can be chosen from a wide variety of window/photocathode combinations to yield higher sensitivity than that of a silicon SSA. Since recently developed IIs are very small, owing to the use of microchannel plate (MCP) electron multipliers, the small size of a solid-state SSA camera is not severely compromised. In summary, advantages of using MCP IIs are

- Long life
- Low power consumption
- Small size and mass
- Rugged
- Very low image distortion
- Linear operation
- Wide dynamic range
- High-speed electronic gating, e.g., a few nanoseconds or less

An image intensifier can be thought of as an active optical element which transforms an optical image from one intensity level to another, amplifying the entire image at one time, i.e., all pixels are

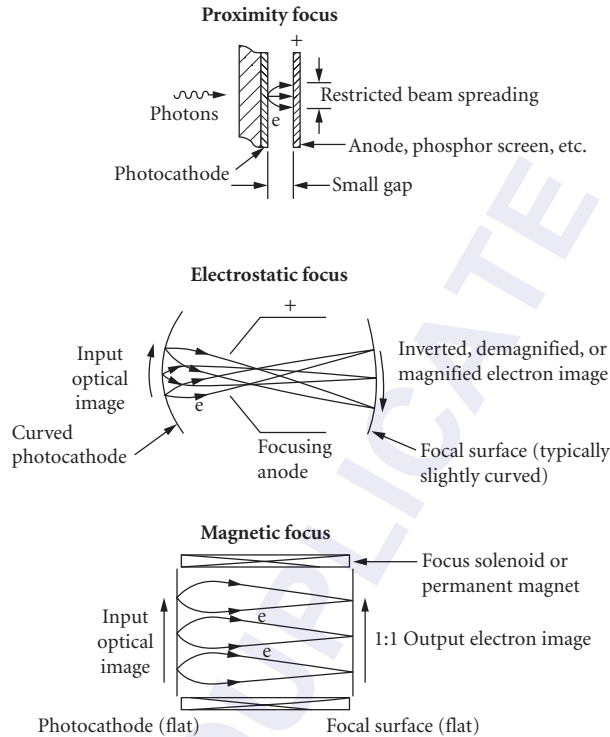


FIGURE 3 Electron lenses.

amplified in parallel and relatively independent of each other. In most cases, the resultant output image is more intense than that of the input image. The level of image amplification depends on the composite efficiency of all the conversion steps of the process involved in the image intensification operation and the basic definition of amplification. The term *image intensifier* is generally used to refer to a device that transforms visible and near-visible light into brighter visible images. Devices which convert nonvisible radiation, e.g., uv or ir, into visible images are generally referred to as *image converters*. For simplicity we refer to both types of image amplifiers/converters as “IIs” in this chapter.

Three general families of IIs exist, shown schematically in Fig. 3, that are based upon the three kinds of electron lenses used to extract the signal electrons from the photocathode, namely,

- Proximity focus IIs
- Electrostatic focus IIs
- Magnetic focus IIs

The first image tubes used a “proximity-focus” electron lens.³ Having inherently low gain and resolution, the proximity-focus lens was dropped in favor of electrostatic focus and magnetic focus IIs. The so-called Generation-0 and Generation-1 image tubes made for the U.S. Army used electrostatically focused IIs. The input end of the silicon-intensifier-target (SIT) vidicon also made use of electrostatic focusing. Magnetic focusing was used extensively in the old TV camera tubes, e.g., image orthicons, image isocons, and vidicons, and also for large-active-area and high-resolution IIs for specialized military and scientific markets.

With the development of the MCP, which was achieved for the U.S. Army’s Generation-2 types of night-vision devices, it became practical to use a proximity-focused electron lens again to meet the

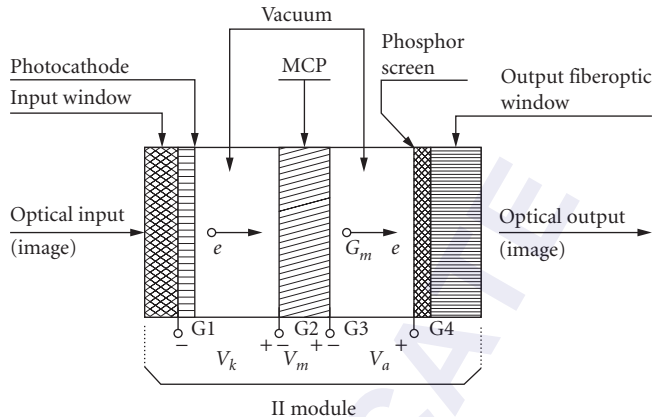


FIGURE 4 Schematic design of a proximity-focused MCP image intensifier tube module.

needs for extremely small and low-mass IIs. These “Gen-2” tubes are being used extensively for military night-vision applications, e.g., night-vision goggles for helicopter pilots, individual soldier helmet mounted night-vision goggles, etc. The most recently developed “Gen-3” IIs have higher sensitivity and limiting resolution characteristics than Gen-2 IIs, and they are used in similar night-vision systems.

Both the Gen-2 and Gen-3 types of IIs are available for use as low-noise, low-light-level amplifiers in II SSA cameras. In addition, by choosing special input window/photocathode combinations outside the military needs for Gen-2 and Gen-3 devices, a very wide range of II SSA spectral sensitivities can be achieved, well beyond silicon’s range. For II SSA camera applications, we will focus our attention exclusively on the use of proximity-focused MCP IIs because of their relative advantages over other types of IIs.

The basic components of a proximity-focused MCP II are shown schematically in Fig. 4. This type of II contains an input window, a photocathode, a microchannel plate, a phosphor screen, and an output window. The *photocathode* on the vacuum side of the *input window* converts the input optical image into an electronic image at the vacuum surface of the photocathode in the II. The *microchannel plate* (MCP) is used to amplify the electron image pixel-by-pixel. The amplified electron image at the output surface of the MCP is reconverted to a visible image using the *phosphor screen* on the vacuum side of the *output window*. This complete process results in an output image which can be as much as 20,000 to 50,000 times brighter than what the unaided eye can perceive. The input window can be either plain transparent glass, e.g., Corning type 7056, fiberoptic, sapphire, fused-silica, or virtually any optical window material that is compatible with the high-vacuum requirements of the II. The output window can be glass, but it is usually fiberoptic, with the fibers going straight through or twisted 180° for image inversion in a short distance.

A block diagram of a generalized high-voltage power supply (HVPS) used to operate the II is given in Fig. 5. For dc operation, the basic HVPS provides the following typical voltages:

$$\begin{aligned} V_k &= 200 \text{ V} \\ V_m &= 800 \text{ V for an MCP} \\ &\quad (V_m = 1600 \text{ V for a VMCP}) \\ &\quad (V_m = 2400 \text{ V for a ZMCP}) \\ V_a &= 6000 \text{ V} \end{aligned}$$

For high-speed electronic gating of the II, the photocathode is normally gated off by holding the G1 electrode a few volts positive with respect to the G2 electrode. Then, to gate the tube on and off

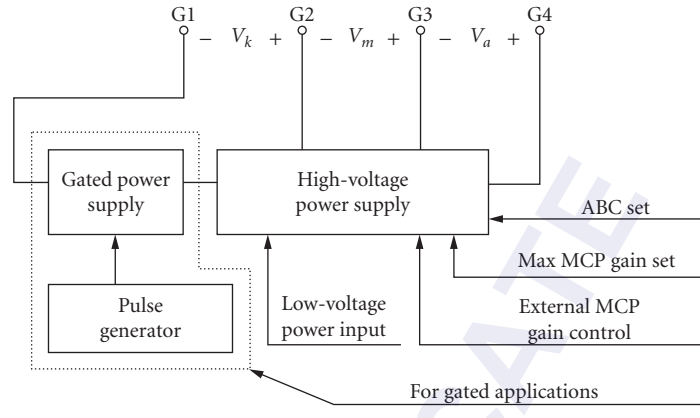


FIGURE 5 MCP image intensifier high-voltage supply.

for a short period, a pulse generator is used to control the output of the gated power supply to the normal gated on condition, i.e., $V_k = 200$ V with the polarity as shown in Fig. 5.

The dc HVPSs for IIs draw very little power, and they can be operated continuously using two AA cells, e.g., 3-V input voltage, for about 2 days. These dc HVPSs are available in small flat-packs or wraparound versions. Gated HVPSs, excluding the pulse generator, are generally at least two times larger than their dc counterparts.

In operation, an input image is focused onto the input window/photocathode assembly, producing a free-electron image pattern which is accelerated across the cathode-to-MCP gap by an applied bias voltage V_k . Electrons arriving at the MCP are swept into the channels, causing secondary electron emission gain due to the potential V_m applied across the MCP input and output electrodes. Finally, the amplified electron image emerging from the output end of the MCP is accelerated by the voltage V_a applied across the MCP-to-phosphor screen gap so that they strike an aluminized phosphor screen on a glass or FO output window with an energy of about 6 keV. This energy is sufficient to produce an output image which is many times brighter than the input image. The brightness gain of the MCP II is proportional to the product of the window/photocathode sensitivity to the input light, the gain of the MCP, and the conversion efficiency of the phosphor-screen/output-window assembly. Each of these key components and/or assemblies is discussed in more detail in the following sections of this section.

Input Window/Photocathode Assemblies

The optical spectral range of sensitivity of an II, or the II SSA that it is used in, is determined by the combination of the optical transmission properties of the window and the spectral sensitivity of the photocathode. In practice, a photocathode is formed on the input window in a high-vacuum system to produce the window/photocathode assembly as shown in Fig. 6. This assembly is then vacuum-sealed onto the II body assembly, and the finished II is then removed from the vacuum system. This type of photocathode processing is called *remote processing* (RP), because the alkali metal generators, antimony sources, and/or other materials used to form the photocathode are located outside of the vacuum II tube. Since there is no room for these photocathode material generators, remote processing must be used for MCP IIs. Also, IIs made using remote processing are found to have significantly less spurious dark current emission than the older Gen-0 and Gen-1 types of IIs having internally processed photocathodes.

The short wavelength cutoff of a window/photocathode assembly is determined by the optical transmission characteristic of the window, i.e., its thickness and material composition. The absolute

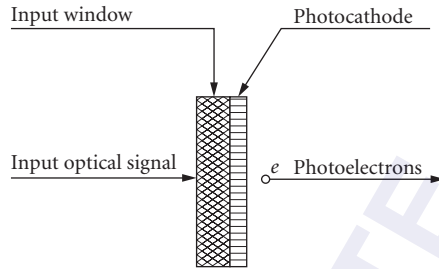


FIGURE 6 Input window/photocathode assembly.

spectral sensitivity of the photocathode determines the midrange and long wavelength cutoff characteristics of the assembly. Photocathode materials having longer wavelength cutoffs also have lower bandgap energies and generally higher thermionic emission than photocathodes with shorter wavelength cutoffs.

The spectral quantum efficiencies of various window/photocathode combinations are shown in Fig. 7 for comparison. Useful spectral bands range from the uv to the near-ir, depending upon the particular combination chosen. This figure shows the spectral sensitivity advantages that can be achieved with II SSAs. Other advantages are discussed throughout this chapter.

Note that the window/photocathode spectral quantum efficiency [$Y(\lambda)$] curves given in Fig. 7 represent the ratio of the average number of photoelectrons produced per input photon as a function of wavelength λ . Alternatively, window/photocathode response can be specified in terms of

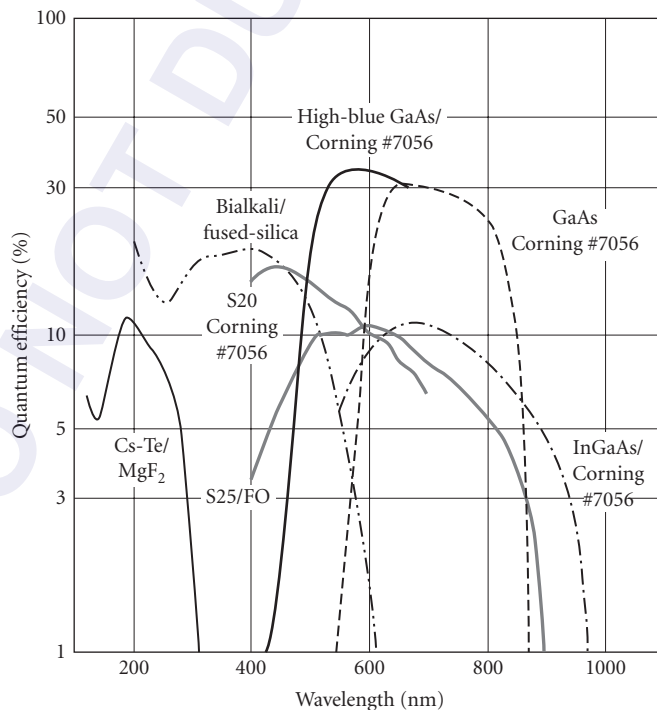


FIGURE 7 Window/cathode spectral quantum efficiencies.

absolute spectral sensitivity [$S(\lambda)$], or defined as the ratio of photocathode current per watt incident as a function of wavelength. These two parameters are related by the convenient equation

$$Y(\lambda) = \frac{124 \cdot S(\lambda)}{\lambda} \tag{13}$$

where Y is the quantum yield in percent, S is the absolute sensitivity in mA/W, and λ is the wavelength in nm.

Microchannel Plates

The development of the microchannel plate (MCP) was a revolutionary step in the art of making IIs. Although developed for and used in modern military passive night-vision systems, MCP IIs are being used today in nearly all II SSA cameras.

An MCP is shown schematically in Fig. 8. Microchannel plates are close-packed-hexagonal arrays of channel electron multipliers. With a voltage V_m applied across its input and output electrodes, the MCP produces a low-noise gain G_m , e.g., a small electron current (I_{in}) from a photocathode produces an output current $G_m I_{in}$. In addition to its function as a low-noise current amplifier, the MCP retains the current density pattern or “electron image” from its input to output electrodes. It is also possible to operate two MCPs (VCMP) or three MCPs (ZMCP) in face-to-face contact to achieve electron gains as high as about $1E7$ e/e in an II tube, as shown in Fig. 9. Other general characteristics of these types of MCP assemblies are also given in Fig. 9.

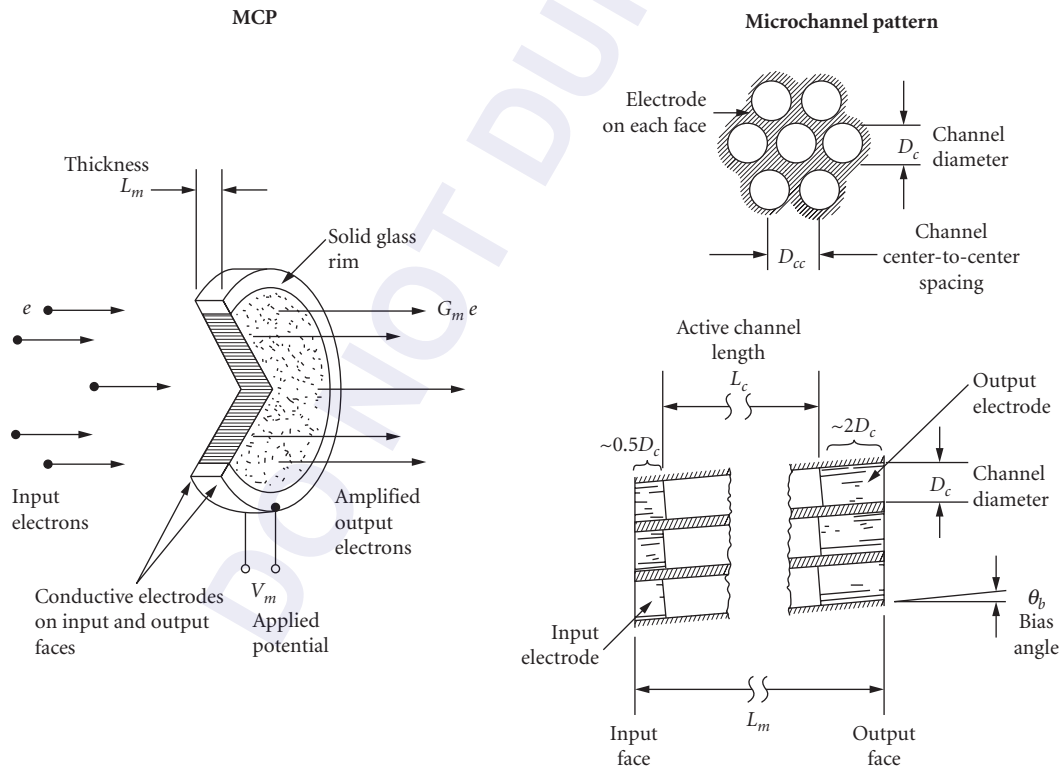


FIGURE 8 MCP parameters.




	V_m max (kV)	G_m max (e/e)	Pulse height distribution (% FWHM)	Relative limiting resolution (units)
MCP: 	1.0	1E3	Negative exponential	1.00
VMCP: 	2.0	1E5	120	0.71
ZMCP: 	3.0	1E7	80	0.50

FIGURE 9 General characteristics of MCPs, VMCPs, and ZMCPs.

The approximate limiting spatial resolutions of MCPs depend upon the channel center-to-center spacings, as follows:

Channel Diameter (μm)	Channel Center-to-Center Spacing (μm)	Approximate Limiting Resolution (lp/mm)
4	6	83
6	8	63
8	10	50
10	12	42
12	15	33

As shown in Fig. 8, MCPs are made to have channel axes that make a “bias angle” (θ_b) with respect to the normal to its input and output faces. This bias angle improves electron gain and reduces noise factor by reducing “boresighting” of electrons into the channels. The MCP bias current or “strip current” (I_s) that results from the voltage applied to the MCP sets an upper limit to the maximum linear dynamic range of the MCP. Generally, when the output current density of the MCP is in excess of about 10 percent of the strip current density, the MCP ceases to remain a linear amplifier. Conventional MCPs have strip current densities of about $1 \mu\text{A}/\text{cm}^2$, and recent high-output-technology MCPs (HOT MCPs)⁴ have become available that have strip current densities as high as about $40 \mu\text{A}/\text{cm}^2$. Electron-gain characteristics of MCP assemblies are given approximately by the equation and associated parameters shown in Table 1.

TABLE 1 MCP Gain Equation and Gain Parameters

$G_m(V_m) = \left(\frac{V_m}{V_c} \right)^g$			
Type	V_c (V)	g (units)	(L_m/D_c) (units)
MCP	350	8.5	40
MCP	530	13	60
VMCP	700	17	80
ZMCP	1050	25	120

Power noise factors for conventional MCPs, used in Gen-2 IIs, and “filmed-MCPs,” used in Gen-3 IIs, are approximately 2.0 and 3.5, respectively. Detailed information on MCP gain, noise factors, and other parameters are given by Eberhardt.⁵ Note that MCP gain is a strong function of the channel length-to-diameter ratio. The parameter V_c in the gain equation is the “crossover” voltage for the channel, i.e., it is the MCP applied voltage at which the gain is exactly unity.

Phosphor Screens

Output spectral and temporal characteristics of a wide variety of screens are given in an Electronic Industries Association publication.⁶ The phosphor materials covered in this publication are listed in Table 2. Both the old “P-type” and the new two-letter phosphor designations are given in this table. Any of these phosphor screen materials can be used in proximity-focused MCP IIs. However, one very commonly used phosphor is the type KA (P20) because it has a high conversion efficiency, its output spectral distribution matches the sensitivity of a silicon SSA reasonably well, it is fast enough for conventional 1/30-s frame times, it has high resolution, and it is typically used in direct-view night-vision IIs.

The three main components of an aluminized phosphor-screen/output-window assembly, of the type used in a proximity focused MCP II, are shown schematically in Fig. 10. An aluminum film electrode is deposited on the electron input side of the phosphor to accelerate the MCP output to high energy, e.g., about 6 keV, and to increase the conversion efficiency of the assembly by reflecting light toward the output window. The phosphor itself is deposited on the glass or fiberoptic output window.

Decay times, or persistence, and relative output spectral distributions for a variety of phosphor types are given in Fig. 11. Key phosphor assembly parameters that should be accounted for in the design of MCP II SSAs are MCP-to-phosphor applied potential (V_a), effective “dead-voltage” resulting from electron transmission losses in the aluminum film, phosphor screen energy input-to-output conversion efficiency, optical transmission of the glass or fiberoptic window, sine-wave MTF of the assembly, phosphor persistence, and output spectral distribution.

Before specifying the use of a particular phosphor, the operational requirements of the II SSA camera should be reviewed. The phosphor persistence should be short compared to the SSA frame

TABLE 2 Worldwide Phosphor-Type Designation System*

P1	GJ	P20	KA	P38	LK
P2	GL	P21	RD	P39	GR
P3	YB	P22	X(XX)	P40	GA
P4	WW	P23	WG	P41	YD
P5	BJ	P24	GE	P42	GW
P6	WW	P25	LJ	P43	GY
P7	GM	P26	LC	P44	GX
P10	ZA	P27	RE	P45	WB
P11	BE	P28	KE	P46	KG
P12	LB	P29	SA	P47	BH
P13	RC	P31	GH	P48	KH
P14	YC	P32	GB	P49	VA
P15	GG	P33	LD	P51	VC
P16	AA	P34	ZB	P52	BL
P17	WF	P35	BG	P53	KJ
P18	WW	P36	KF	P55	BM
P19	LF	P37	BK	P56	RF
				P57	LL

*Cross reference: old-to-new designations.

Source: Adapted from Electronic Industries Association Publication, no. 116-A, 1985.

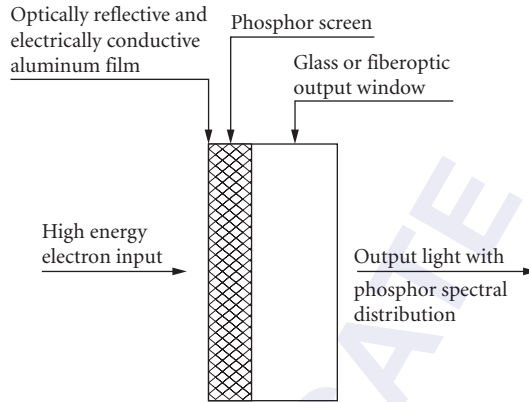


FIGURE 10 Aluminized phosphor screen and window assembly.

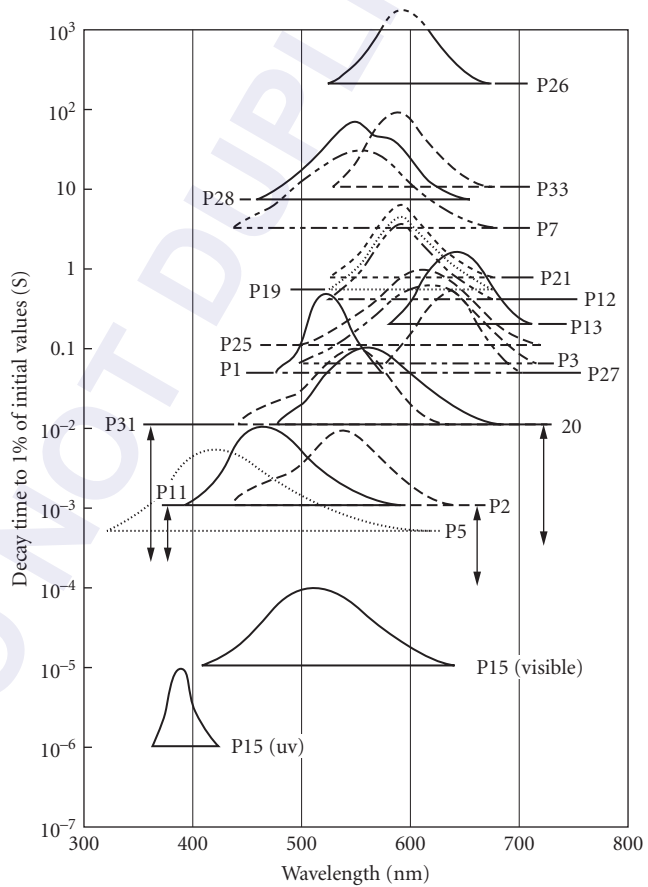


FIGURE 11 Phosphor screen decay times and spectral outputs. (Reprinted with permission from United Mineral and Chemical Co.)

time to minimize image smear due to rapidly moving objects. Also, the absolute conversion efficiency of the phosphor assembly and its relative output spectral distribution should be spectrally matched⁷ to the sensitivity of the SSA for maximum coupling efficiency.

Typical absolute spectral response characteristics, i.e., the phosphor spectral efficiency (radiated watts per nanometer per watt excitation) as a function of wavelength, of aluminized phosphor screens are given in Ref. 7. The associated phosphor screen efficiencies are also given in this reference in three different ways:

- Typical quantum yield factor: photons out per eV input
- Typical absolute efficiency: radiated watts per watt excitation
- Typical luminous equivalent: radiated lumens per radiated watt

For example, a type KA(P20) aluminized phosphor-screen/glass window assembly is found to have its peak output at 560 nm and a typical quantum yield factor of 0.063 photons/eV. Thus, an electron which leaves the MCP and strikes the assembly with 6 keV of energy, and for a “dead-voltage” of 3 kV, approximately $(6 - 3) \text{ keV} \times 0.063 \text{ photons/eV} = 190$ photons will be produced at the output.

Proximity-Focused MCP IIs

By combining the image transfer and conversion properties of the three major proximity-focused MCP II assemblies discussed earlier, i.e.,

- Input window/photocathode
- Microchannel plate
- Phosphor screen/output window

the operational characteristics of the II itself, as shown in Fig. 4, can be determined.

For example, consider an II CCD application for a space-based astronomical telescope that requires more than 10 percent quantum yield at 200 nm, but minimum sensitivity beyond 300 nm. It is desired that the top end of the dynamic range be at an input window signal flux (F_{si}) of 1000 photon/pixel/s at 250 nm. Let the CCD have a 1-in vidicon format, i.e., an active area of $11.9 \times 8.9 \text{ mm}^2$, with 325 vertical columns and 244 horizontal rows of pixels. The limiting resolution of even a dual-MCP (VMCP) image tube has a limiting resolution that is significantly higher than the horizontal pixel spatial Nyquist frequency (f_N) in the CCD, so that the pixel size at the input to the II will be essentially the same as that of the CCD. Let us rough-in an II design by making the following additional assumptions:

MCP-to-phosphor applied potential (V_a)	6000 V
Phosphor screen type	KA (P20)
Phosphor screen/window-quantum yield (P_q)	0.06 photon/eV
Phosphor screen dead voltage (V_d)	3000 V
CCD charge integration period (τ_i)	33 ms
CCD pixel full-well charge	1 pC = $6.3E6$ e

An II with an 18-mm active diameter can be used, since the diagonal of the CCD active area is 14.9 mm. From Fig. 7, the $\text{MgF}_2/\text{Cs-Te}$ window/photocathode assembly will be chosen, having a quantum yield (Y_k) of 0.12 at 200 nm, to meet the spectral sensitivity requirements.

Let's now proceed to estimate the required gain of the MCP structure, decide what kind of an MCP structure to use, and determine its operating point. A first-order estimate of the stored pixel charge (Q_{ccd}) for the given input signal flux density is

$$Q_{ccd} = F_{si} \cdot Y_k \cdot G_m \cdot (V_a - V_d) \cdot P_q \cdot Y_{ccd} \cdot \tau_i \quad (14)$$

Since

$$F_{si} = 1000 \text{ photon/pixel/s}$$

$$Y_k = 0.10 \text{ e/photon}$$

$$Y_{ccd} = 0.3 \text{ e/photon}$$

it is found that $Q_{ccd} = G_m (178 \text{ e/pixel})$. Setting this charge equal to the full-well pixel charge gives $G_m = 6.3E6 \text{ e/pixel}/(178 \text{ e/pixel}) = 3.5E4 \text{ e/e}$. This MCP assembly gain is easily satisfied by using a VMCP. From Table 1, it is found that the gain of a VMCP is given approximately by $G_m = (V_m/700)^{17} = 3.5E4 \text{ e/e}$. Solving for V_m gives $V_m = 1300 \text{ V}$.

Thus, a first-order estimate for the general requirements to be placed in the II to do the job is as follows:

Active diameter	18 mm
Quality area	(11.9 × 8.9 mm)
Input window/photocathode	Fused-Silica/Cs-Te
MCP assembly	VMCP
Aluminized phosphor screen assembly	KA/FO window

Coupling this II to the specified FO input window CCD, e.g., by using a suitable optical cement, will meet the specified objective. Other parameters like the dark count rate per pixel as a function of temperature, the DQE of the II CCD, cosmetic, uniformity of sensitivity, and other specifications will have to be considered as well before completing the design.

Recent “Generations” of MCP IIs The most impressive improvement in direct-view night-vision devices has come with the advent of Gen-3 technology. The improvement, which is most apparent at very low light levels, is mainly due to the use of GaAs as the photocathode material. At higher light levels, e.g., half-moon to full-moon conditions, the Gen-2+ gives somewhat better performance. Key to the detection of objects under LLL conditions is the efficiency of the photocathode; the Gen-3 sensitivity is typically a factor of 3 higher. Also, the spectral response of Gen-3 matches better to the night sky spectral illumination. This equates to being able to see at almost one decade lower scene illumination with Gen-3. A summary of proximity-focused MCP image intensifier general characteristics is given in Table 3.

TABLE 3 Summary of Proximity-Focused MCP Image Intensifier General Characteristics

Minimum Active Diameter (mm)	Input Window Material*	Spectral Sensitivity Range (nm)	MCP Assembly Type	Temperature Rating		Output Window Material	Minimum Limiting Resolution (lp/mm)	Technology Type
				Storage (°C)	Operating (°C)			
11.3	FS	160–850	MCP	–55, +65	–20, +40	FO	25	Gen-2
12.0	FS, G, FO	160–900	MCP, VMCP, ZMCP	–57, +65	–51, +45	FO, G	45, 29, 20	Gen-2
17.5	FS, G, FO	600–900	MCP, VMCP, ZMCP	–57, +65	–51, +45	FO, G	45, 25, 20	Gen-2
17.5	G, FO	500–1100	MCP, VMCP, ZMCP	–57, +95	–51, +52	FO, G	45, 25, 20	Gen-3
25.0	FS, G, FO	160–900	MCP	–57, +65	–51, +45	FO, G	40	Gen-2
25.0	G, FO	500–1100	MCP	–57, +95	–51, +52	FO, G	40	Gen-3

*FS = fused silica; G = Corning #7056 glass; FO = fiberoptic.

Technology Type	Options	
	Photocathode	Phosphor
Gen-2	All but GaAs, InGaAs	Wide selection
Gen-3	GaAs, InGaAs	Wide selection

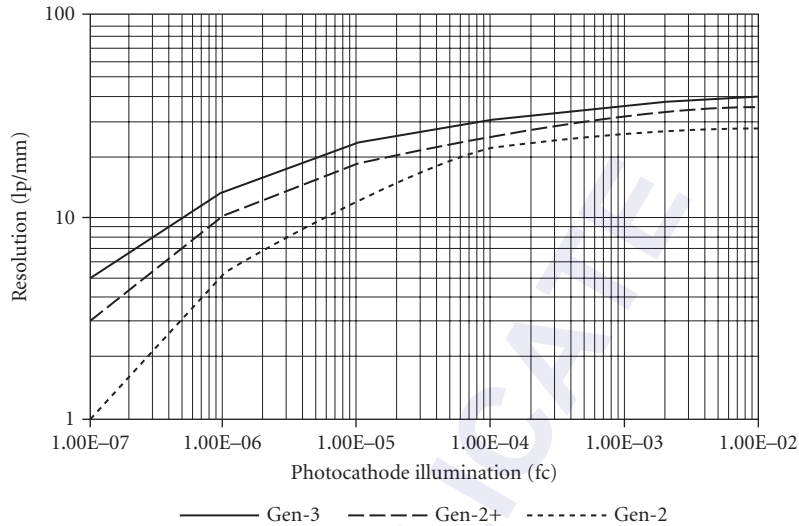


FIGURE 12 Image intensifier tube resolution curves.

For systems design work, it is useful to know the approximate characteristics of the three most recent generations in terms of II resolution versus photocathode illumination. The resolution transfer curves shown in Fig. 12 give the II resolution, observable by the eye, as a function of input illumination for Gen-2, Gen-2+, and Gen-3 IIs. These curves do not include system optics degradations, except in the sense that a human observer made the resolution measurements using a 10-power eyepiece in viewing the output image of the II.

Improved Performance Gen-2 IIs Recent enhancements in the dynamic range performance of Gen-2 IIs for direct-view applications have been made which also benefit II SSA camera performance. Improvement goals were to increase both the usable output brightness and the LLL gain of Gen-2 IIs. Night-vision devices are normally used at light levels ranging from full moon to just below quarter-moon, or in dark city environments with ample scattered light. It is important to have good contrast over as wide a light-level range as possible. To get this extended dynamic range, the gain should be held nearly constant to as high a level as possible, for improved contrast at the high-light levels. Any gain improvement should be attained with little or no increase in noise, to ensure good performance at the minimum light levels. Reducing the objective lens f -number as low as possible also improves system performance and gain. However, f -number reduction by itself may create problems in the system dynamic range if the II and its power supply assembly is not appropriately adjusted to match the optical throughput.

Figure 13 shows the extended dynamic range of a Gen-2+ II and power supply assembly, as compared to the typical MIL-SPEC Gen-2 assembly. Increasing the gain in a standard Gen-2 assembly by increasing the gain control voltage, i.e., the MCP voltage, will not give the same benefits as the Gen-2+. Ideally, a change of one unit in input brightness should result in a proportional output brightness change. The increased near-linear gain range up to higher-output light levels in the Gen-2+ improves the contrast at the higher levels. Brightness limiting begins reducing the gain to hold the output brightness constant after the automatic brightness control (ABC) limit of the power supply is reached. The increased gain of the Gen-2+ improves the performance at the lowest-light levels as well.

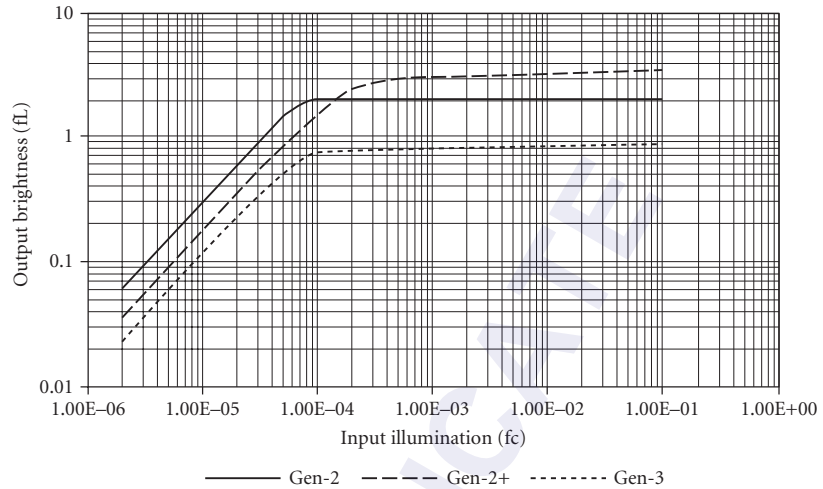


FIGURE 13 Output versus input transfer characteristics of Gen-2, Gen-2+, and Gen-3 II/power supply assemblies.

31.5 IMAGE INTENSIFIED SELF-SCANNED ARRAYS

There are several reasons to consider using an II SSA instead of an SSA alone. One obvious reason is to achieve LLL sensitivity. Figure 14 shows the limiting resolution versus faceplate illumination characteristic of CID camera operating in the unintensified and intensified modes.⁸ It is seen that

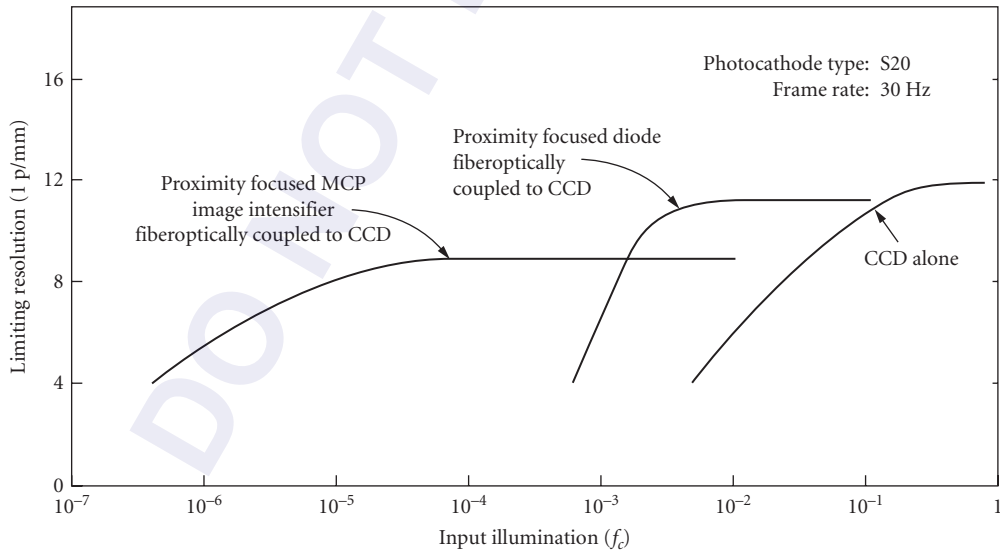


FIGURE 14 Resolution versus input illumination characteristics of a conventional optical input CCD camera and the same camera fiberoptically coupled to an MCP image intensifier tube. (From Ref. 8.)

LLL sensitivity is achieved by coupling the CID to an image intensifier tube, albeit at the expense of reduced high-light resolution. Other reasons for using an II SSA are

- High-speed electronic gating, down to a few nanoseconds, for framing cameras, LADAR, smoke and fog penetration
- Improved spectral sensitivity
- Use in a TV camera system that operates automatically under lighting conditions ranging from nighttime to full daylight conditions.
- High-sensitivity and high-speed-gated optical multichannel analyzers (OMAs)

Fiberoptic-Coupled II/SSAs

Figure 15 shows a schematic design of a fiberoptically (FO) coupled II SSA assembly. These designs are modular, since an II module is optically coupled to an SSA module. Virtually any type of image tube can be optically coupled to an SSA. The fiberoptically coupled design shown in Fig. 15 requires the use of an II having a fiberoptic output window and an SSA having an SSA input window. A fiberoptic taper, instead of a simple unity magnification FO window, is also generally required to efficiently couple the output of the II into the SSA, and this is shown in Fig. 15 as a separate module. The various fiberoptic modules are joined at interfaces 1, 2, and 3, using optical cement, optical grease, immersion oil, or "air." For the highest-resolution image transfer across these interfaces, it is necessary that the gap length at each interface be kept short, and the numerical aperture of the fiberoptic windows should be kept as low as possible, consistent with the SNR and gain requirements. It has been shown⁹ that the first interface can be eliminated by making the fiberoptic taper part of the II and depositing the phosphor screen directly onto it, and interface 3 can also be eliminated by coupling the fiberoptic taper directly to the SSA. The properties of the image transfer and conversion components shown in Fig. 15 can be used to estimate the overall performance characteristics of the fiberoptically coupled II SSA camera.

The terminology used to define SSA image format sizes derives from the earlier vidicon camera tube technology. The mass, volume, and power requirements of vidicon cameras are much larger than SSA cameras. Vidicons also have image distortion and gamma characteristics which must be accounted for, whereas SSAs and II SSSAs using proximity-focused IIs are nearly distortion-free with

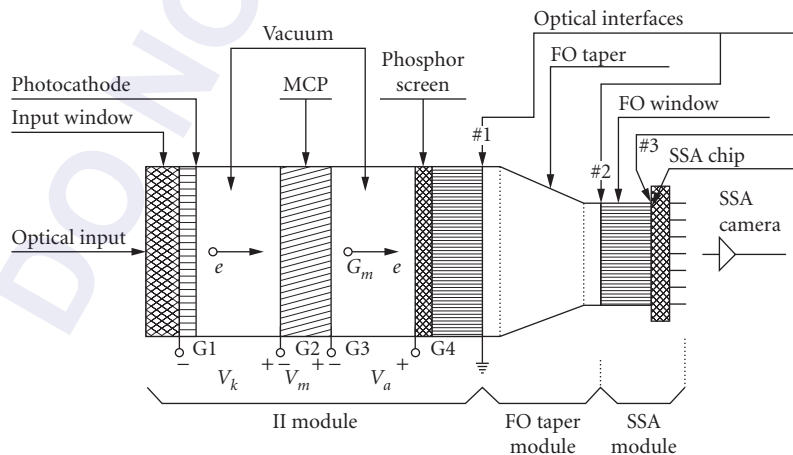


FIGURE 15 Schematic design of fiberoptically coupled II SSA assembly.

TABLE 4 Comparison of Basic Image Intensifier Diameters, SSA Format Sizes, Matching Fiberoptic Taper Magnifications, and Limiting Resolutions at the Fiberoptic Taper Output Surface (for 45 lp/mm Intensifier)

Image Intensifier Active Dia. (mm)	SSA Format		Diagonal (mm)	(M_{fot}) FOT Magnification	(f_{ito}) Limiting Resolution at FOT Output (lp/mm)
	Vidicon (in)	(mm)			
25	1	11.9 × 8.9	14.9	0.596	76
25	2/3	8.8 × 6.6	11.0	0.440	102
18	1	11.9 × 8.9	14.9	0.828	54
18	2/3	8.8 × 6.6	11.0	0.611	74
18	1/2	6.5 × 4.85	8.1	0.451	100
12	2/3	8.8 × 6.6	11.0	0.917	49
12	1/2	6.5 × 4.85	8.1	0.676	67
12	1/3	4.8 × 3.6	6.0	0.500	90

linear, i.e., unity gamma, input/output transfer characteristics over wide intrascene dynamic ranges. Table 4 gives the basic II active diameters, SSA format sizes, SSA active-area diagonal lengths, fiberoptic taper magnifications (M_{fot}) required to couple II outputs to the SSAs, and limiting resolutions (f_{ito}) at the fiberoptic taper output. Figure 16 shows schematically the relative sizes of the standard active diameters of IIs and the standard SSA formats.

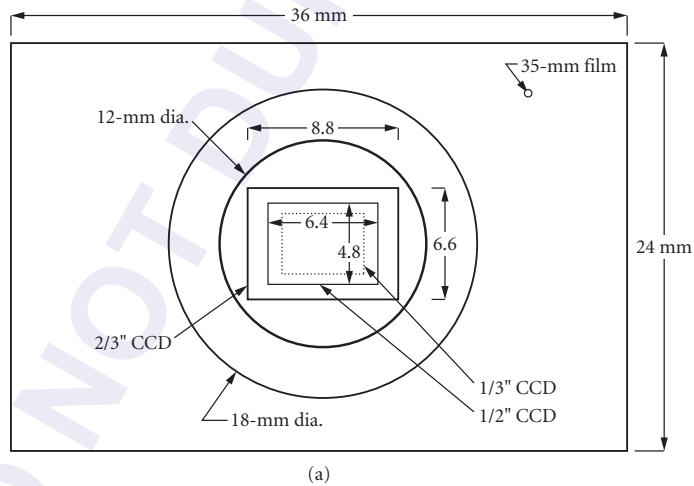


FIGURE 16a Typical 35-mm film, image intensifier and SSA formats.

	18-mm Φ image tube	12-mm Φ image tube	1" CCD	2/3" CCD	1/2" CCD	1/3" CCD	35-mm film	Units
Diagonal	18.0	12.0	14.9	11.0	8.1	6.0	43.3	mm
Vertical	10.8	7.2	8.9	6.6	4.9	3.6	24.0	mm
Horizontal	14.4	9.6	11.9	8.8	6.5	4.8	36.0	mm
Area	155.5	69.1	105.9	58.1	31.5	17.3	864.0	mm ²

(b)

FIGURE 16b Typical dimensions for image intensifiers and SSAs using 3:4 format.

The present limiting resolution range of MCP IIs is 36 to 51 lp/mm. In an II SSA, the resolution of the II should be matched, in some sense, to that of the SSA. For example, it is unwise to use a low-resolution II and fiberoptic lens combination with a much higher resolution CCD.

Lens-Coupled II SSAs

Figure 17 is a schematic design for a lens-coupled II SSA assembly. The differences between this design and the fiberoptic-coupled II SSA design described earlier are that the output window of the II can be either fiberoptic or glass, and a lens is used instead of an FO taper to couple the output optical image from the II directly into a conventional optical input SSA, i.e., no FO window is required at the SSA. Although the lens-coupling efficiency is lower, its image distortion and resolution performance is superior to the FO-coupled design. Also, the chance for possible adverse rf interference at the sensitive input to the SSA camera from the II high-voltage power supply is less than for the lens-coupled design.

Parameters to Specify Typical parameters to specify for an MCP II SSA detector assembly, using either fiberoptic or lens-coupling, are as follows:

- Sensitivity
 - White-light (2856 K) ($\mu\text{A}/\text{lm}$)
 - Spectral sensitivity (mA/W versus nm)
 - Sensitivity (mA/W at specified wavelength)
- EBI (lm/cm^2 at 23°C)
- MCP applied potential for 10 K fl/fc luminous gain (V)
- Horizontal resolution at specified input illumination (TVL)
- Shades-of-gray (units)
- Cosmetic properties
 - Uniformity (percent)
 - Bright spots (number allowable in format zone) Dark spots (number allowable in format zone)
- Burn-in (procedure)

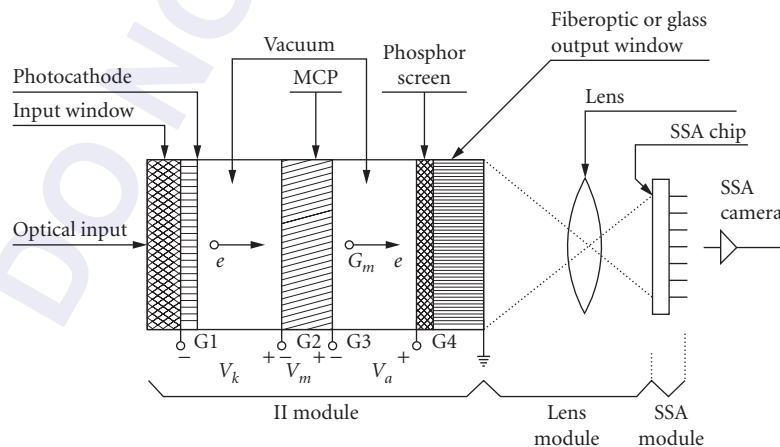


FIGURE 17 Schematic design of lens-coupled II SSA assembly.

- Mechanical specifications
- Dimensions (interface drawing)
- Mass (g)
- Environmental (specified)

Electron-Bombarded SSA

Since the early work by Abraham et al.¹⁰ which showed the feasibility of achieving useful electron gain by electron bombardment (EB) of a silicon diode in a photomultiplier tube, several attempts have been made to achieve similar operation using an SSA specially designed for EB input, instead of optical input. The charge gain (G_{eb}) resulting from the electron bombardment is given by

$$G_{eb} = \frac{(V_a - V_d)}{3.6} \quad (15)$$

where V_a is the acceleration voltage and V_d is the “dead-voltage” of the EBSSA. It was quickly found that successful CCD operation could not be obtained by simply bombarding the normal optical input side of the chip with electrons, because interface states soon form which prevent readout of the chip and other problems. By thinning a CCD chip to 10 to 15 μm from the “backside” and operating in a backside EB-mode, useful performance is achieved. In this way, 100 percent of the silicon chip is sensitive to incident photoelectrons, and it becomes technically feasible to make EBSSA cameras.

Proximity Focused EBSSAs A proximity-focused EBSSA is shown schematically in Fig. 18. In this design, the input light enters the window/photocathode assembly to generate the signal photoelectrons which are accelerated to about 10-keV energy and bombard the thinned backside of the EBSSA. Note that no MCP, no MCP-to-screen gap, no phosphor screen/output window assembly, and no fiberoptic or lens coupling is used to transfer the electronic image to the SSA for readout. Thus, higher limiting resolution is attainable. Also, the power noise factor associated with the EBSSA gain process is lower than that of MCP devices, and image lag is eliminated because no phosphor is used. Early work on proximity-focused EBDDs was done by Barton et al.,¹¹ Williams,¹² and Cuny et al.¹³ By 1979, a 100×160 pixel TI CCD was used in this type of detector and put into a miniature TV camera. With an acceleration voltage of $V_a = 15$ kV, an electron gain of 2000 was achieved, along with a Nyquist limited resolution of 20 lp/mm. Recent advances have brought this technology closer

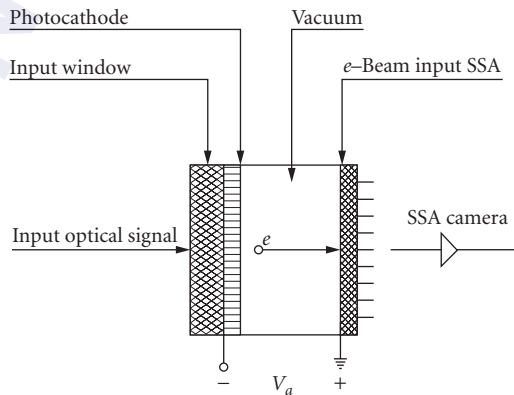


FIGURE 18 Electron bombarded SSA (EBSSA).

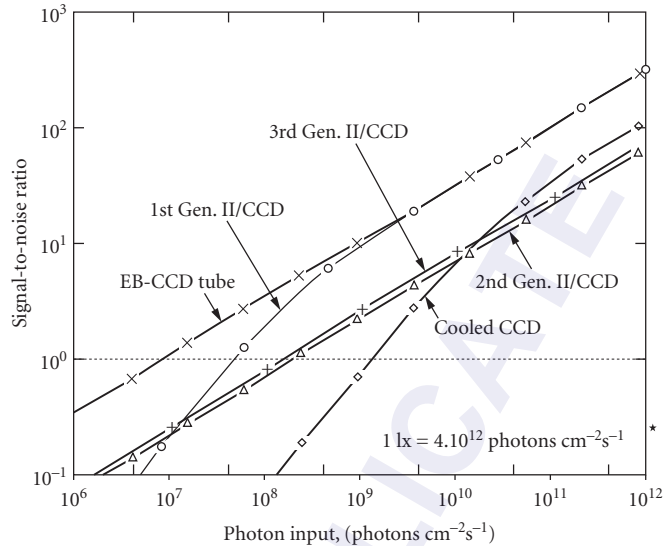


FIGURE 19 Comparison of the signal-to-noise ratio of various optoelectronic imagers versus the photon input. (From Ref. 14.)

to extensive usage possibilities. Richard et al.¹⁴ have compared the SNR characteristics of an EB CCD tube, various other types of II CCDs, and bare CCDs. Their results are shown in Fig. 19.

In order to achieve its full performance capabilities, the energy of the bombarding electrons must be absorbed by the active silicon SSA material, photoelectrons must not be lost, the exposure of the EBSSA to high-energy electrons should not cause a life problem, and it must be possible to read out the stored charge pattern in the SSA. It is found that recombination phenomena at the EB-input face can be reduced with a p^+ passivation layer, e.g., by using $3E17 \text{ cm}^{-3}$ boron doping, which reduces back-diffusion of signal electrons, front-diffusion of “dark” charges from the rear face, reduced diffusion length, separation of holes and electrons by the built-in electric field, and higher surface conductivity, thus better voltage stability, at the rear face.

Internally processed (IP) and remotely processed (RP) or “transfer” photocathodes have been used in EBSSAs. It is generally found that the internal processing produces consistently higher-background and spurious noise problems due to field emission from tube body parts and the photocathode. Both types of photocathode processes have yielded long-life EBCCD detectors.

Proven applications to date for EBSSA detectors:

- Photon-counting wavefront sensor (adaptive optics), European Space Organization 3.6-m telescope at La Silla, Chile
- NASA, Goddard Space Flight Center, Oblique Imaging EB CCD uv sensitive camera

Advantages of EBSSA cameras over MCP II-based II-SSAs:

- No image lag
- Higher resolution
- Single photoelectron detection per frame per pixel
- Higher DQE

Digital II SSA Cameras Consider a photon-counting imaging detector consisting of an MCP image intensifier tube (II) that is fiberoptically coupled to a silicon solid-state self-scanned array

(SSA) chip in a TV camera. Incoming photons at wavelength λ pass through the input window of the II and produce an average quantum yield of Y_k photoelectrons per photon at the photocathode. The resulting photoelectrons (e) are accelerated into the MCP electron multiplier assembly. Amplified output electrons from this low-noise electron multiplier are accelerated into an aluminized phosphor screen on the output window of the II. The number of output photons from the II per photoelectron is proportional to the electron gain in the MCP (G_m), the effective electron bombardment energy at the phosphor screen (ϵV_s), and finally the electron-input to photon-output conversion efficiency (P) at the phosphor screen. As discussed earlier the optical transmission of the input window and the actual quantum yield of the photocathode are usually factored together in the average quantum yield parameter Y_k , and the optical transmission of the output window is also normally factored together with the actual conversion efficiency of the phosphor screen in the screen efficiency parameter P .

The output photon pulse from the II, resulting from the single detected input photon, is coupled into the SSA via the fiberoptic taper, which matches the output size of the II to the size of the SSA, and a fiberoptic window on the SSA. This photon pulse is then converted to an electron signal charge packet (Q_{ssa}) at the SSA. The number of electrons stored per pixel in the SSA depends upon the area of the photon pulse at the SSA, the spatial distribution of photons in this pulse, and the area per pixel in the SSA. Thus, in addition to the above II factors, the stored charge in the SSA per photoelectron is also proportional to the optical transmissions of the FO taper (T_{fot}) and SSA window (T_{ssa}), and the quantum yield of the SSA (Y_{ssa}).

By using two or three conventional MCPs in cascade, i.e., VMCPs or ZMCPs, the gain can be made so large that it completely overrides any normal room-temperature thermal dark current in an SSA at a conventional RS-170 rate. In this photon-counting mode of operation, a charge signal above a preset threshold value is looked for. When it is found in a given pixel, a "1" is stored in memory for that pixel's address, "0s" are stored in pixel addresses where this condition is not met, and the entire frame is read out. By reading out a total of N_f frames, the dynamic range can be made as high as N_f if the dark count rate is negligible. Thus, photon-counting imaging can achieve a very large dynamic range.

Another advantage of photon-counting imaging is that the image resolution can also be made very high by centroiding the detected charge packets in the SSA. Since the performance of a centroiding camera depends upon the signal-processing algorithm, this will not be analyzed here. Instead, the reader is referred to several references in which centroiding is discussed.¹⁵

Let us next calculate the stored charge and number of stored electrons in a photon-counting II SSA per photoelectron. Assume that a proximity-focused VMCP II is coupled to the SSA with a fiberoptic taper. For our analysis, some typical values will be used for the operating voltage and gain of a VMCP: the acceleration voltage between the VMCP and the phosphor screen, the efficiency of an aluminized type KA (P20) phosphor screen, the optical transmissions of a fiberoptic taper and an SSA fiberoptic window, and the quantum yield of an SSA.

Definitions for the parameters that will be used are summarized as follows:

Q_{ssa}	stored SSA charge per input photoelectron from the photocathode
Y_k	photocathode quantum yield, e/photon
G_m	VMCP electron gain, e/e
V_m	VMCP applied potential, V
V_s	MCP-to-screen applied potential, V
V_d	phosphor screen "dead-voltage"
P	phosphor screen efficiency, photon/eV
T_{fot}	transmission of fiberoptic taper
T_{ssa}	transmission of fiberoptic window on the SSA
Y_{ssa}	quantum yield of SSA, e/photon
e	electron charge, 1.6E-19 C
N_{essa}	number of stored SSA electrons per input photoelectron

Using these definitions, the general equation for the charge stored in the SSA per input photoelectron is given by Eq. (16).

$$Q_{ssa} = e \cdot G_m \cdot (V_s - V_d) \cdot P \cdot T_{\text{fot}} \cdot T_{ssa} \cdot Y_{ssa} \quad (16)$$

Thus, Q_{ssa} is given by the product of the VMCP gain, the effective electron bombardment energy at the aluminized phosphor screen, the conversion efficiency of the phosphor screen assembly, the transmissions of the FO taper and the SSA's FO window, and finally the quantum yield of the SSA.

For

$$V_m = 1380 \text{ V}$$

the VMCP electron gain is

$$G_m(V_m) = \left(\frac{V_m}{700 \text{ V}} \right)^{17}$$

$$G_m(V_m) = 1 \times 10^5$$

By using the following values for the additional parameters

$$V_s = 5500 \text{ V}$$

$$V_d = 2500 \text{ V}$$

$$P = 0.06 \text{ photon/eV}$$

$$T_{\text{fot}} = 0.6$$

$$T_{ssa} = 0.8$$

$$Y_{ssa} = 0.5 \text{ e/photon}$$

it is found that the stored charge per photoelectron is

$$Q_{ssa} = 7 \times 10^{-13} \text{ C}$$

and that the number of electrons stored in the SSA per input photoelectron is

$$N_{\text{essa}} = Q_{ssa} / e$$

$$N_{\text{essa}} = 4 \times 10^6 \text{ electrons}$$

Since the full-well or saturation charge for an SSA pixel is on the order of 1 pC, this VMCP II/SSA assembly is seen to qualify as a photon-counting imaging detector.

Modulation Transfer Function and Limiting Resolution The modulation transfer function (MTF) of an II SSA camera is determined by a convolution of the individual MTFs of the camera lens, II, II-to-SSA-coupling fiberoptic or lens, fiberoptic-to-fiberoptic interfaces, fiberoptic-to-SSA interface, SSA, etc. There are several ways to determine the MTF of an existing II SSA camera.

For example, the II SSA camera can be focused on a spatial frequency burst pattern, i.e., a periodic pattern of black and white bars in which the spatial frequency of the bars increases in one direction. Alignment of the pattern's bars with pixel columns and readout of the modulation of the spatial pattern in the pixel row direction gives the squarewave MTF of the camera $S(f)$, where f is the spatial frequency in cycles/mm. Conversion of this square-wave MTF to a sine-wave MTF is accomplished by using the Fourier transform at a given frequency:

$$T(f) = \left(\frac{\pi}{4} \right) \cdot \left(S(f) + \frac{S(3f)}{3} - \frac{S(5f)}{5} + \frac{S(7f)}{7} + \dots \right) \quad (17)$$

By calculating several values of $T(f)$ from the known square-wave function, the sine-wave MTF can be determined and used for camera system optical image transfer analysis. The limiting resolution is often taken to be the spatial frequency value for this sine-wave MTF at which the modulation drops to a few percent.

Also, to use the above example as an illustration, the MTF of an II SSA is a function of the direction in which the spatial frequency burst pattern is aligned. Self-scanned array pixels are not generally square, and the distances between centers of pixels in the horizontal and vertical directions are not generally the same. Thus, the corresponding MTFs in the horizontal and vertical directions are different, and the MTF is a function of the angle that the burst pattern makes with the rows and columns of pixels.

A convenient specification of the spatial frequency response of an SSA or an II SSA camera is the Nyquist frequency (f_N), defined to be the reciprocal of twice the distance between the pixels. For example, if an SSA has rows of pixels spaced on 20- μm center-to-center, then the horizontal Nyquist spatial frequency is

$$f_{N,h} = \frac{1}{(2 \cdot (0.02 \text{ mm}))} = 25 \text{ cycles/mm}$$

Assuming an II limiting resolution (f_{II}) of 32 cycle/mm, and assuming that the MTFs are gaussian, then an estimated value for the limiting resolution of the II SSA camera is

$$\left(\frac{1}{f_{\text{cam}}^2} \right) = \left(\frac{1}{f_{II}^2} \right) + \left(\frac{1}{f_{\text{ssa}}^2} \right) \quad (18)$$

or

$$f_{\text{cam}} = \frac{f_{II} \cdot f_{\text{ssa}}}{\sqrt{(f_{II}^2) + (f_{\text{ssa}}^2)}} \quad (19)$$

For the values used in this example, $f_{\text{cam}} = 20$ cycle/mm. Although this gaussian estimate is convenient to use, a more exact estimate can be made by multiplying the various component sine-wave MTFs to find the II SSA camera's MTF, and from this its limiting resolution can also be found.

For example, actual MTF measurements⁹ on an II SSA camera, having a proximity-focused 18-mm active diameter MCP II fiberoptically coupled to a CCD with an $m = 8 \text{ mm}/18 \text{ mm} = 0.44$ magnification taper, showed that the MTF of the CCD, referred to the II input, is given by $T_{\text{CCD}}(f) = \exp - (f/9.0)^{1.4}$, where f is the spatial frequency in cycles/mm. The Nyquist frequency of the CCD was $f_N = 20$ cycle/mm, and the MTF of the complete camera was found to be $T_{\text{II SSA}}(f) = \exp - (f/6.3)^{1.1}$, both referred to the II input.

31.6 APPLICATIONS

In time, it is expected that most of the quantum-limited and LLL TV applications will use some form of II SSA camera, instead of an intensified vidicon-based camera. A few of the major application areas for II SSA cameras are highlighted in this section.

Optical Multichannel Analyzers

Optical multichannel analyzers (OMAs) are instruments used to measure optical radiation in linear patterns, e.g., spectra or two-dimensional images. Photographic film, single-channel photomultiplier tubes, and TV camera tubes, e.g., vidicons, have been replaced by SSAs, e.g., CCDs, and II SSAs, e.g., image tube intensified CCDs or photodiode arrays (PDAs). Four distinct application areas exist for OMAs using II SSA detectors:¹⁶

Application	Detector*	Time
	Type	Resolution
Spectroscopy	IILPDA	50 ms
Time-resolved pulsed laser spectroscopy	GIILPDA	<5 ns
Time-resolved pulsed laser imaging spectroscopy	GISSCCD	<5 ns/spectrum
Time-resolved imaging	GISSCCD	<5 ns/spectrum

*IILPDA = image intensified photodiode array; GIILPDA = gated image intensified photodiode array; GISSCCD = gated image intensified charge-coupled device.

In each of these types of systems, the incoming radiation is converted to charge packets that are stored in each pixel of the SSA. Each line and/or field of pixels is read out by a suitable camera electronics, and the pixel charge values are stored in a computer for subsequent processing and analysis. The sensitivity of the II is chosen for best performance over the range of wavelengths being investigated. The dynamic range of OMAs can be as high as 18 bits. In comparison with the older single-element scanning system, modern II SSA-based OMAs acquire spectra up to 1000 times faster and/or with higher SNR during a given measuring period. Three common OMA applications are Raman spectroscopy, multiple input spectroscopy, and small-angle light scattering.

Range Gating and LADAR

Range gating is becoming increasingly important because the required technology now exists at an affordable cost and the signal-to-noise ratio improvement is much higher than nongated conventional TV imaging. A gated laser sends out a laser pulse of only a few nanoseconds duration while the II SSA camera is gated off. No reflection from scattering or reflection in the medium between the camera and the object is allowed to be registered in the II SSA camera.¹⁷ The II SSA camera is gated on only at the moment when the light packet from the object returns to the camera and, after exposing for the duration of the outgoing pulse, it is returned to a gated-off condition. By repeating this process and controlling the image-storing conditions, high-contrast images having high signal-to-noise ratios can be achieved. Another obvious advantage of the range-gated system is that the time-of-flight between pulse output and receipt gives the range to the object being viewed, thus leading to the realization of a *laser detection and ranging* (LADAR) system.

Microchannel plate image tubes offer high-speed gating and spectral response advantages to LADAR systems. They can be electronically gated to a few nanoseconds, i.e., providing distance resolutions of a few feet for LADAR systems. Present MCP IIs offer the user a broad range of spectral sensitivity, including near-ir imaging at 1060 nm, so that powerful and efficient lasers may be used for optimum LADAR system performance. Also, the ruggedness, extremely small size, and low power drain characteristics of II SSA cameras make them very attractive for LADAR applications for spacecraft and unmanned autonomous vehicles.

Day/Night Cameras

Full day-night interscene dynamic range capability, while maintaining a high signal-to-noise ratio, is achievable using an II SSA camera in conjunction with an auto-iris camera lens.¹⁸ The principle of operation of this type of camera can be described as follows. Assume that operation begins at the lowest light level to be encountered. The MCP voltage in the II is set to operate with high SNR for the camera lens, an auto-iris lens, set to its lowest f -number. As the light level is increased, the system operates over two orders-of-magnitude dynamic range. For four orders-of-magnitude higher light level inputs, the f -setting of the auto-iris lens is increased by a feedback circuit, driven by the peak-to-peak video output signal from the SSA camera. The effective exposure of the SSA is next automatically reduced as the light level increases by four-and-one-half orders-of-magnitude, again to

maintain a high SNR, by duty-cycle gating the MCP II. Finally, another two-and-one-half orders-of-magnitude in input light level are accommodated by reducing the gain of the MCP, i.e., by reducing its applied operating potential. The total interscene dynamic range achievable with this type of automatically controlled day/night camera is 13 orders-of-magnitude. In addition to its wide dynamic range capability, this type of camera is also able to make rapid narrow-band spectral samples across a wide spectral range, e.g., from the uv to the near-ir.

Mosaic II SSA Cameras

For very high amounts of image information throughput, multiple SSAs are used to read out large area IIs. For example, a 75-mm active diameter MCP II can be coupled fiberoptically to four individual SSA cameras. The fiberoptic couplers are made to butt against each other at the output of the II, and their output ends are optically coupled to the SSAs. Parallel readout of the SSAs is then accomplished, giving the advantage of high-resolution readout without the disadvantage of having to use a wide bandwidth video electronic system or lower frame rates.

One such II SSA camera is designed for x-ray radiology image input.¹⁹ The x-ray input image is converted to a visible light image at a scintillator screen that is in optical contact with the 6-in-diameter fiberoptic window. This scintillator/window assembly is, in turn, coupled to the input of a 6-in-diameter proximity-focused diode II, i.e., without an MCP, for modest light gain and good image quality. Six fiberoptic tapers in a 2×3 matrix couple the images from the six adjacent output sections of the II to six CCD cameras which operate in parallel to continuously read out the converted x-ray image.

Other Applications

Other applications for II SSA cameras are the following:

- Semiconductor circuit inspection
- Astronomical observations
- X-ray imaging
- Coronary angiography
- Mammography
- Nondestructive testing
- Multispectral video systems

Active Imaging Active imaging is becoming increasingly important because the required technology now exists at an affordable cost and the signal-to-noise ratio improvement is much higher than nonactive conventional TV imaging. Two types of active imaging presently exist, i.e., “line-scanned” and “range-gated.”

In a line-scanned imaging system, a narrow beam from a cw laser is raster-scanned across the object to be viewed, and the resulting reflected light is collected by a lens and detector assembly which receives and measures light from the illuminated field-of-view (FOV). Large FOV scenes can be scanned in a short period of time, which is a major advantage of this system. The signal from the detector is finally processed by a video electronics system and displayed, for direct viewing, or image processed as required. Limiting-resolution is set by the beam diameter achievable at the object. Thus systems operating in space; atmospheric and underwater environments have significantly different limiting-resolution characteristics.

In a range-gated type of system, a gated laser sends out a laser pulse of only a few nanoseconds duration while the TV camera is gated off. No reflection from scattering or reflection in the medium between the camera and the object is allowed to be registered in the TV camera. The TV camera is

gated on only at the moment when the light packet from the object returns to the camera and, after exposing for the duration of the outgoing pulse, the TV camera is returned to a gated-off condition. By repeating this process and controlling the image-storing conditions, high-contrast images having high signal-to-noise ratios can be achieved.

Another obvious advantage of the range-gated system is that the time of flight between pulse output and receipt gives the range to the object viewed, thus leading to the realization of a laser detection and ranging (LADAR) system. MicroChannel plate image tubes offer high-speed gating and spectral response advantages to LADAR systems. They can be electronically gated to a few nanoseconds, i.e., distance resolutions of a few feet, and in some parts of the optical spectrum they offer high sensitivity.

31.7 REFERENCES

1. A. Rose, "A Unified Approach to the Performance of Photographic Film, Television Pickup Tubes, and the Human Eye," *J. Soc./Motion Picture Engrs.* **47**:273–294 (1946).
2. D. E. Caudle, "Dynamic Range Enhancement Techniques for Gated, Solid-State Intensified Cameras," *SPIE* **1155**:104–109 (1990).
3. G. Hoist, J. H. de Boer, and C. F. Veenemans, *Physica* **1**:297 (1934).
4. HOT MCP™ is a trademark of Galileo Electro-Optics Corp.
5. E. H. Eberhardt, "An Operational Model for MicroChannel Plate Devices," *IEEE Trans. Nucl. Sci.* **NS-28**:712–717 (1981).
6. "Optical Characteristics of Cathode Ray Tubes," EIA Publication, no. 116-A, 1985.
7. H. P. Westman, ed., *Reference Data for Radio Engineers*, 5th ed., Howard W. Sams & Co, Inc., Indianapolis, 1969, pp. 16–34–16–37.
8. J. J. Cuny, T. F. Lynch, and C. B. Johnson, "Proximity Focused Image Tube Intensified Charge Injection Device (CID) Camera for Low Light Level Television," *SPIE* **203**:75–79 (1979).
9. G. M. Williams, Jr., "A High Performance LLLTV CCD Camera for Nighttime Pilotage," *SPIE* **1655**:14–32 (1992).
10. J. M. Abraham, L. G. Wolfgang, and C. N. Inskeep, "Application of Solid-State Elements to Photoemission Devices," *Adv. EEP* **22B**:671 (1966).
11. J. B. Barton, J. J. Curry, and D. R. Collins, "Performance Analysis of EBS-CCD Imaging Tubes/Status of ICCD Development," *Proc. Int. Conf. Apps. of CCDs*, San Diego, Calif., 1975.
12. J. T. Williams, "Test Results on Intensified Charge Coupled Devices," *SPIE* **78**:78–82 (1976).
13. J. J. Cuny, T. F. Lynch, and C. B. Johnson, "Small Intensified Charge Injection Device Cameras for Low Light Television," *MEDE '79 Conf. Proc.*, Interavia SA, Publishers, 1979, pp. 836–845.
14. J. C. Richard, M. Vittot, and J. C. Rebuffie, "Recent Developments and Applications of Electron-Bombarded CCD in Imaging," *Proc. Conf. on Photoelectronic Image Devices*, SEP91, London, IOP Pub. Ltd., 1992.
15. A. Boksenberg and D. E. Burgess, "An Image Photon Counting System for Optical Astronomy," *Adv. EEP* **33B**:835–849 (1972).
16. H. W. Messinger, "Modern Optical Multichannel Analyzers Capture Images Without Film," *Laser Focus World*, March 1992, pp. 91–94.
17. R. A. Sturz and D. E. Caudle, "Capabilities of New Cost-Effective Near Infrared Imaging," *Adv. Imaging*, April 1993, pp. 60–62.
18. D. E. Caudle, "Low Light Level Imaging Systems Application Considerations and Calculations," *SPIE* **1346**:54–63 (1990).
19. H. Roehrig et al., "High Resolution X-Ray Imaging Device," *SPIE* **1072**:88–99 (1989).

Timothy J. Tredwell

*Sensor Systems Division
Imager Systems Development Laboratory
Eastman Kodak Company
Rochester, New York*

32.1 GLOSSARY

A	area of the pixel
C_{FD}	total capacitance of the floating diffusion in a CCD output
C_g	gate capacitance per unit area
C_r	readout line capacitance
J_D	total dark current per unit area
J_s	surface generation current
L_e	diffusion length of electrons in silicon
L_p	diffusion length of holes in silicon
N_A	p -type dopant concentration in silicon
N_D	n -type dopant concentration in silicon
N_e	total number of electrons collected in a pixel
$N(0)$	number of photons entering the silicon
$N(x)$	number of photons remaining a distance x below the surface
n_i	intrinsic carrier concentration in silicon
$P(x)$	probability that an electron-hole pair generated a distance x from the surface will be collected before recombination
q	electron charge
S_o	surface recombination velocity
T_{int}	integration time of light in an image sensor
$T(\lambda)$	transmission of light
V_{bi}	built-in voltage for a silicon pn junction
$W(V)$	width of the depletion layer at a given bias voltage in the MOS capacitor or junction-photodiode
$\alpha(\lambda)$	absorption coefficient of light

ϵ_s	silicon dielectric constant
μ_e	electron mobility in silicon
μ_p	hole mobility in silicon
τ_o	depletion-layer lifetime
τ_p	minority carrier hole lifetime in silicon
Φ_s	electrostatic potential at the silicon-silicon dioxide, also called surface potential

32.2 INTRODUCTION

Since the invention of the image sensor in 1964, solid state image sensors have advanced in resolution, sensitivity, and image quality to the point where they have replaced other methods of converting visible light to electronic signals in nearly all imaging applications. There are two types of image sensors: *area image sensors*, which are used in cameras, and *linear sensors*, which are used in scanning applications. Cameras using area image sensors dominate the camcorder, video and broadcast, machine vision, scientific, and medical fields. Area image sensors for camcorder applications are typically 400,000 picture elements in resolution, 60 dB in dynamic range, and have noise levels of a few tens of electrons. Area image sensors for scientific applications may have resolutions of over six million elements, dynamic ranges exceeding 80 dB, and noise levels approaching a single electron. Scanners employing linear solid-state image sensors dominate facsimile, document scanner, digital copier, and film scanner applications. Linear sensors range from 2000-element monochrome arrays with 40 dB of dynamic range used in facsimile applications to 8000 or more element trilinear arrays with 80 dB of dynamic range for high-performance color scanning applications.

The steps involved in image sensing consist of (1) converting the incoming photons to charge at picture element (pixel), and (2) transferring that charge to an output amplifier and converting the charge to a voltage or current signal which can be sensed by circuits external to the sensor. The image-sensing elements are described first, followed by readout elements. Sensor architectures for area and linear sensors are described next.

32.3 IMAGE SENSING ELEMENTS

There are four basic types of structures which are used for image sensing: the junction photodiode, the photocapacitor, the pinned (p^+np) photodiode, and the photoconductor.* The first three are generally fabricated in single-crystal silicon as part of the image sensor; the photoconductor is usually fabricated from amorphous silicon deposited over the image sensor. The photoconversion process begins with the absorption of a photon in silicon resulting in the generation of a single electron-hole pair. The absorption of light at a particular wavelength is given by

$$N(x, \lambda) = N(0)e^{-\alpha(\lambda)x} \quad (1)$$

where $N(x)$ is the number of photons remaining at a distance x below the surface, $N(0)$ is the number of photons entering, and $\alpha(\lambda)$ is the absorption coefficient.^{1,2} The absorption coefficient is shown in Fig. 1 as a function of wavelength λ for single-crystal silicon, doped polycrystalline silicon, and hydrogenated amorphous silicon. The absorption depth is defined as the inverse of the absorption coefficient [$d(\lambda) = 1/\alpha(\lambda)$]. In single-crystal silicon, the absorption depth is 0.4 μm in the blue (450 nm), 1.5 μm

*For some scientific applications in which high quantum efficiency and fill factor are essential, the silicon wafer is thinned to 10 μm or less in thickness and illuminated from the backside. The frontside contains an area charge coupled device, which is used to collect the photogenerated carriers.

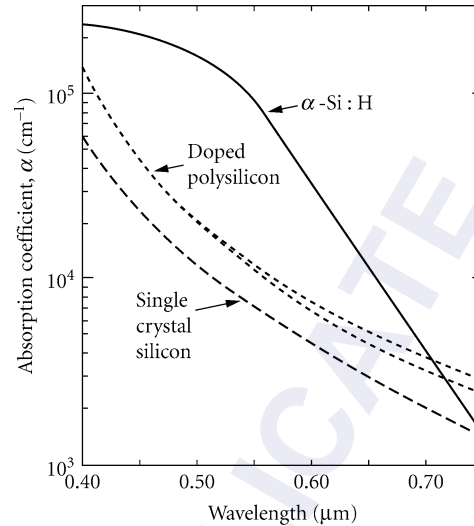


FIGURE 1 Absorption coefficient for light in single-crystal silicon, heavily doped polycrystalline silicon, and hydrogenated amorphous silicon as a function of wavelength.

in the green (550 nm), and 3.0 μm in the red (640 nm). In the infrared, the absorption depth increases to 10.5 μm at 800 nm. Beyond 1100 nm, the absorption is virtually zero because the photon energy is less than the 1.1-eV silicon bandgap.

Junction Photodiode

The junction photodiode is one of the most common image-sensing elements. The physical structure and band diagram of the junction photodiode are shown in Fig. 2a for a p -type substrate. The n -type region is formed by ion implantation or diffusion of phosphorous or arsenic to a depth of 2000 to 10,000 \AA into the p -type silicon. The n -type dopant region is usually graded, with the highest concentration at the surface. The gradient in n -type dopant concentration results in a gradient in electrostatic potential which accelerates photogenerated carriers (holes in the n -type region) away from the surface. This reduces loss of photogenerated carriers to surface recombination. The photodiode is typically operated with a reverse bias V of 1 to 5 V. A depletion layer is formed between the n - and p -type regions.* The width of the depletion layer $W(V)$ for an abrupt $n + p$ junction is given by

$$W(V) = \sqrt{\frac{2\epsilon_s(V + V_{bi})}{qN_A}} \quad (2)$$

where q is the electronic charge, ϵ_s is the silicon dielectric constant, N_A is the p -type dopant concentration, and V_{bi} is the built-in voltage given by

$$V_{bi} = \frac{kT}{q} \ln \frac{N_A}{n_i} \quad (3)$$

*For a detailed review of the device physics of junction diodes and MOS capacitors, see Sze, *Physics of Semiconductor Devices*, Wiley, New York, 1969.

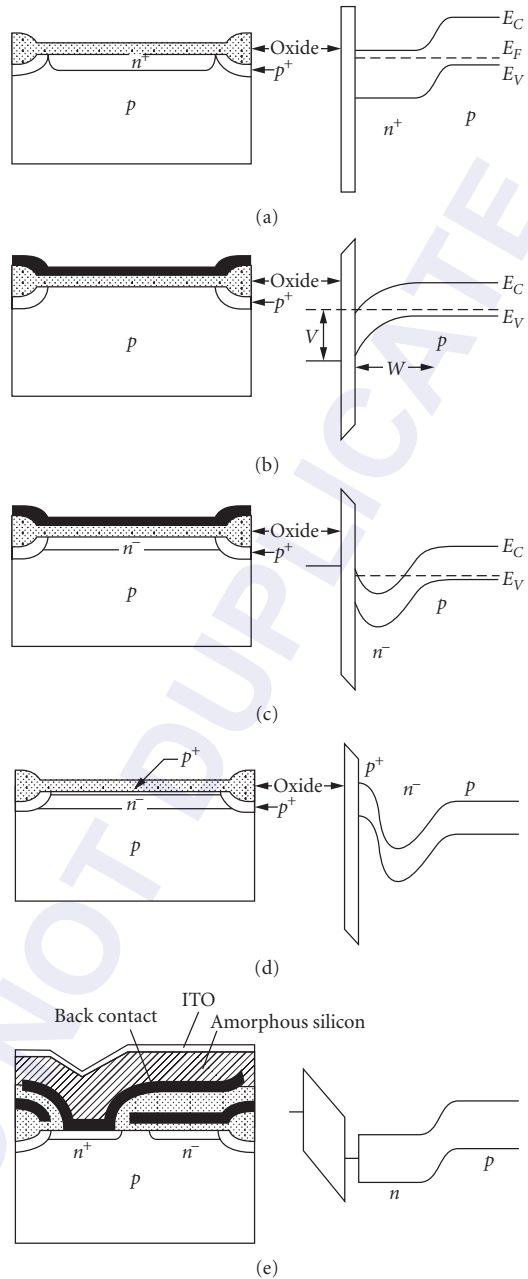


FIGURE 2 Cross-sectional diagrams and band diagrams for (a) junction photodiode; (b) surface-channel MOS capacitor; (c) buried-channel MOS capacitor; (d) pinned or hole-accumulated diode; and (e) amorphous silicon photoconductor.

where n_i is the intrinsic carrier concentration. For silicon doped at $1 \times 10^{16} \text{ cm}^{-3}$, the depletion width at 5-V reverse bias is 0.8 μm .

If the diode is illuminated, some of the photons will be absorbed in the n -region, some in the depletion layer, and the remainder in the p -type substrate. The quantum efficiency $\eta(\lambda)$ is the ratio of the charge collected to the number of photons incident on the diode (i.e., 100-percent quantum efficiency refers to one electron-hole pair collected for every incident photon). The quantum efficiency depends on three factors: transmission $T(\lambda)$ of light through the overlying layers into silicon, absorption of light in silicon, and the probability $P(x)$ that an electron-hole pair generated a distance x from the surface will be collected before recombination:

$$\eta(\lambda) = T(\lambda) \int_{x=0}^{x=\infty} (1 - e^{-\alpha(\lambda)x}) P(x) dx \quad (4)$$

The transmission of light through the overlying layers into the silicon can be calculated using standard multilayer interference models.

The collection of the photogenerated charge takes place by two processes: drift and diffusion. Drift is the movement of electrons and holes due to an electric field. Even for small electric fields, the transport of carriers by drift will dominate diffusion. This is the case in the depletion region. Outside the depletion region, such as in the p -type substrate, there is no electric field, and carrier transport occurs by diffusion. For example, a photon absorbed in the p -type region will excite an electron from the valence band to the conduction band. The electron will move in a three-dimensional random walk until it recombines or encounters the edge of the depletion region, where it is swept across the junction by the electric field. The probability that an electron at a distance x' from the junction can diffuse to the junction before recombining is given by

$$P(x') = e^{-x'/L_e} \quad \text{where} \quad L_e = \sqrt{\frac{KT}{q}} \mu_e \tau_e \quad (5)$$

in which L_e is the diffusion length, μ_e is the electron mobility, and τ_e is the electron lifetime (typically $\sim 1 \mu\text{s}$). For a 1- μs lifetime, the electron diffusion length in p -type silicon is 50 μm . Electrons generated from photons absorbed at less than the diffusion length from the junction have a high probability of collection.*† Similarly, photons absorbed in the n -type layer create electron-hole pairs. The holes must travel by diffusion to the junction in order to be collected. The probability that a hole a distance x' from the junction can diffuse to the junction is given by

$$P(x') = e^{-x'/L_p} \quad \text{where} \quad L_p = \sqrt{\frac{KT}{q}} \mu_p \tau_p \quad (6)$$

in which L_p is the diffusion length, μ_p is the hole mobility, and τ_p is the hole lifetime. For a 1- μs lifetime, the hole diffusion length in n -type silicon is 30 μm . Since the n -type region in an $n + p$ junction is less than 1 μm thick, the holes have no difficulty diffusing through the n -type region to the junction unless the n -type region is so heavily damaged or so heavily doped that the hole lifetime is

*In image sensors, the collection of photogenerated carriers can be complicated by a variety of factors. The doping concentration may not be uniform on either the n - or p -sides. On the n -side, the dopant concentration is designed to decrease from the surface to the junction, building in a potential gradient for holes away from the surface and toward the junction, preventing surface recombination. On the p -side, the dopant concentration may not be uniform owing to the use of epitaxial layers or wells diffused into silicon. Additionally, the carrier lifetime may not be uniform. Impurity gettering, a process used in many image sensors to remove metallic contaminants will leave a region of crystalline defects in the silicon starting 20 to 50 μm beneath the silicon surface. The defects result in a very short electron lifetime in this region. Finally, diffusion takes place in three dimensions; a carrier absorbed beneath a given pixel in an array will diffuse laterally by the same amount it diffuses vertically. This can cause it to be collected in adjacent pixels.

†See, for example, Lavine et al., "Steady State Photocurrent Collection in Silicon Imaging Devices," *IEEE Transactions on Electron Devices* ED-30:1123–1133 (Sept. 1983).

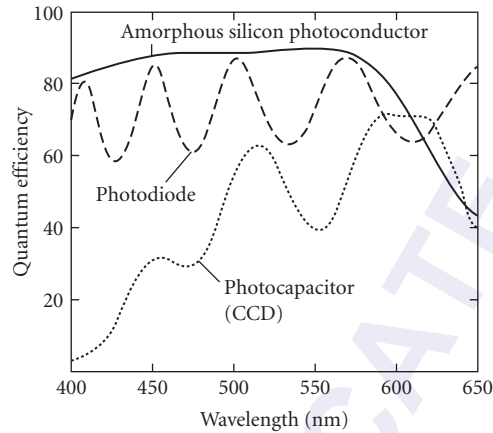


FIGURE 3 Quantum efficiency as a function of wavelength for photocapacitor, photodiode, and amorphous silicon photoconductor.

very short. More typically, loss of quantum efficiency on the n -side results from recombination at the surface.

The quantum efficiency of a junction photodiode is illustrated as a function of wavelength in Fig. 3. In the ultraviolet, the light is absorbed very near the surface and some of the photo-generated charge can be lost to surface recombination. In the 420- to 700-nm region, most of the light is absorbed close to the junction and is easily collected. The structure in the quantum efficiency illustrated in Fig. 3 results from the multilayer reflections of light in the oxide layer which overlies the photodiodes used in image sensors. The positions of the minima and maxima depend on the thicknesses and indices of refraction of the layers overlying the silicon. Beyond 800 nm, some of the photons are absorbed sufficiently deep in silicon that the carriers recombine before they can diffuse to the junction; this results in a decrease in quantum efficiency at longer wavelengths.

In most image sensing applications, the junction diode at each picture element is used not only to collect the photogenerated carriers but also to store the carriers until they can be read out. In an imaging array, each photodiode would be reset to a reverse bias V_r by a MOS gate. The capacitance of the diode of area A for an abrupt $n + p$ junction is given by $C(V) = \epsilon_s A / [W(V)]$ where $W(V)$ is the depletion width.

When a photogenerated carrier is collected by the junction, it is stored on the junction until it is read out. Storage of a carrier will cause the voltage on the junction to decrease by $q/C(V)$. When the photogenerated charge is removed from the junction during readout, the junction voltage is restored to its original value.

If so much photogenerated charge is stored on the photodiode that the voltage drops to zero before the charge can be read out, additional charge cannot be collected and will diffuse into the p -type substrate. This condition is referred to as *saturation*. The diffusion of excess charge into neighboring photosites is called *blooming*.

One of the difficulties encountered in using the junction photodiode in image sensor applications is *image lag*. The combination of the capacitance of the photodiode and the channel resistance of the MOS transfer gate used to read out the diode give rise to a time constant for transferring the photogenerated charge from the diode to the readout structure. As a result, not all the charge can be completely drained from the diode during the short reset times typically used in imaging applications. The remaining charge is drained in successive readouts, causing an afterimage. This effect is called *image lag*.³

MOS Capacitor

The MOS capacitor consists of the silicon substrate (taken to be p -type in this section), a thin layer of silicon dioxide (typically 200 to 1000 Å thick), and an electrode (typically polycrystalline silicon doped heavily n -type with phosphorous). The physical structure and the band diagrams of the surface channel MOS capacitor are shown in Fig. 2*b* for a p -type substrate. If the gate of the MOS capacitor is biased positive, a depletion layer is created in the p -type silicon substrate. The depth of the depletion layer depends on the substrate doping, gate voltage, and oxide thickness. The calculation of the depth of the depletion layer is somewhat more complex than the photodiode* and depends on the electrostatic potential at the surface, called the surface potential Φ_s . For values typical of image sensors ($N_A = 1 \times 10^{15}$, 5-V gate bias, 500-Å oxide thickness) the depletion layer is 2.4 μm deep. On the edges of the photocapacitor (see Fig. 2*b*) is a heavily doped p -type region overlaid by a thick (2000- to 5000-Å) oxide layer. This is called the field or *channel stop* region. Because of the heavier p -type doping, the field is not depleted by the voltage on the gate. The channel stops confine the electrons to the channel region.

If the MOS capacitor is illuminated, a fraction of the light will be reflected, a fraction will be absorbed in the polysilicon, and the rest will be transmitted into the silicon substrate. The absorption coefficient of heavily doped polysilicon is shown as a function of wavelength in Fig. 1. The absorption coefficient is $4 \times 10^4 \text{ cm}^{-1}$ at 450 nm and $1.2 \times 10^4 \text{ cm}^{-1}$ at 550 nm. For a 3000-Å-thick polysilicon electrode, less than 30 percent of the blue light and less than 70 percent of the green light is transmitted through the polysilicon. The photons entering the silicon are absorbed, either in the depletion layer of the MOS capacitor or in the undepleted p -type silicon beneath the depletion layer. Those photogenerated electrons created in the undepleted p -type region move by diffusion until they are captured by the depletion layer or they recombine. The electrons are held at the silicon-SiO₂ interface, where they remain until they are read out. The quantum efficiency as a function of wavelength is illustrated in Fig. 3. The efficiency is low in the blue owing to absorption in the polysilicon. The structure in the quantum efficiency as a function of wavelength is due to multilayer interference in the polysilicon-oxide-silicon stack.

Charge is stored in the MOS capacitor at the silicon-silicon dioxide interface as a layer of sheet-charge only a few hundred angstroms thick. As more photogenerated charge is added, the surface potential decreases. If sufficient photogenerated charge is added, the surface potential becomes zero and no additional charge can be stored. This condition is saturation.

The capacitance per unit area on which photogenerated charge is stored is the parallel capacitance of the oxide and the depletion layer:

$$C^{-1} = \left(\frac{t_{\text{ox}}}{\epsilon_{\text{ox}}} + \frac{W(\Phi_s)}{\epsilon_{\text{si}}} \right) \quad (7)$$

where the depletion width $W(\Phi_s)$ is give by

$$W(\Phi_s) = \sqrt{\frac{2\epsilon_{\text{si}}\Phi_s}{qN_A}} \quad (8)$$

in nearly all cases, the oxide capacitance is the dominant term. As a result, the storage capacity of the MOS capacitor is significantly larger than the junction diode in which the charge is stored only on the depletion capacitance.

A variant of the surface-channel photocapacitor is the buried-channel photocapacitor. The structure and band diagram are shown in Fig. 2*c*. In this device, a lightly-doped n -type region is diffused or implanted into the silicon surface early in the fabrication process. This n -type region is sufficiently lightly doped that it is fully depleted. The n -type dopant in the buried channel results in

*To determine the depletion depth, the surface potential Φ_s must be calculated which depends on the voltage on the gate of the MOS capacitor, the oxide thickness, substrate doping, and weakly on temperature. See reference on p. 32.3.

a band diagram with a potential minimum, or well, for electrons just below the surface.* This well is separated from the surface by a few thousand angstroms in distance and about 1 V in potential. When the buried-channel photocapacitor is illuminated, the electrons collect in the buried channel and do not contact the surface. The primary purpose of the buried channel is to prevent electrons from being trapped by interface states at the silicon-silicon dioxide interface.

The photocapacitor has several advantages over the junction photodiode. These include higher storage capacity per unit area, zero lag readout, and generally lower dark current. The principal disadvantage is the low quantum efficiency in the blue. In some applications, a transparent electrode, such as indium-tin-oxide (ITO) may be substituted to improve the blue response.⁴ ITO has very low absorption over the visible (420 to 750 nm) and can be deposited sufficiently conductive for use as a gate electrode in an image sensor.

Another approach to achieving high quantum efficiency is the thinned backside-illuminated charge-coupled device (CCD). In this approach, the CCD (which is an array of MOS capacitors) is fabricated on the frontside of a silicon wafer. The wafer is then thinned from the backside to 10 μm or less in thickness. The backside is passivated to prevent surface recombination. Photons entering the backside are absorbed in silicon beneath the MOS capacitors (usually buried channel). The photogenerated carriers diffuse to the capacitors, where they are held until they are read out. This device has quantum efficiency similar to the photodiode in the visible. However, because the silicon is thin, some of the photons at wavelengths beyond 700 nm will not be absorbed and so the quantum efficiency falls off beyond 700 nm. Owing to their extreme complexity in process and their extremely fragile design, backside-illuminated image sensors are limited to special scientific applications, especially astronomy.

Pinned Photodiode

The third type of photosensitive element is the pinned (p^+np) photodiode.⁵ This is sometimes called the *hole accumulation diode*, or HAD. This element combines the best features of the photodiode and photocapacitor, offering the high blue response of the photodiode with the high charge capacity, zero lag, and low dark current of the buried-channel photocapacitor. The pinned photodiode consists of a very shallow (<2000 Å) P^+ layer overlying an n -type buried-channel region. The structure and band diagrams are shown in Fig. 2d. The p^+ surface layer, which contacts the p^+ channel stop region on the sides, holds the electrostatic potential at the surface at 0 V. When the photodiode is illuminated, the photogenerated electrons are held in the n -type buried-channel region just below the surface.

The quantum efficiency of the pinned photodiode is nearly identical to that of the photodiode shown in Fig. 2. Because there is no overlying polysilicon electrode, the blue response is very high, similar to that of the photodiode. Because the buried-channel region can be completely emptied, the pinned photodiode does not have the lag of a normal junction (n^+p or p^+n) photodiode. The pinned photodiode has been the most widely used image-sensing element in interline area CCDs used for camcorders and for industrial and medical cameras. The pinned photodiode is also used in some linear image sensors, particularly where low image lag is critical.

Photoconductor

The last type of photosensitive device is the photoconductor. The most common material for the photoconductor is hydrogenated amorphous silicon, although other material systems have been explored. Amorphous silicon photoconductors have been used for two types of image sensors: area image

*For a review of the calculation of electrostatic potential and charge capacity of buried-channel MOS capacitors and charge-coupled devices, see B. C. Burke, G. Lubberts, E. A. Trabka, and T. J. Tredwell, *IEEE Transactions on Electron Devices* ED-31(4):423 (April 1984).

sensors, where it is deposited on top of an area array to improve fill factor (i.e., the proportion of the picture element which is photosensitive), and contact linear sensors, where it is deposited on large ceramic or glass substrates to fabricate very long line or linear arrays.

The structure of an amorphous silicon photoconductor on a CCD image sensor and the corresponding band diagrams are shown in Fig. 2e. The hydrogenated amorphous silicon photoconductor consists of a back electrode, an undoped amorphous silicon layer approximately 1 μm thick, and a transparent top electrode.⁶ Additional doped amorphous silicon or silicon nitride layers may be added to the amorphous silicon front or back surfaces to improve performance. Photons absorbed in the amorphous silicon generate electron-hole pairs. Photogenerated electrons and holes are quickly swept to the back and front electrodes, respectively, because of the high electric field in the photoconductor. When the amorphous silicon is used as part of an area image sensor, the electrons on the back electrode can be transferred into the readout element; when used as part of a contact linear array, the voltage change can be amplified and read out through a multiplexer.

The quantum efficiency of an amorphous silicon photoconductor is shown in Fig. 3. Owing to the wider bandgap of the amorphous silicon, photons of wavelength greater than about 650 nm are not absorbed. The advantage of the amorphous silicon is the high quantum efficiency across the visible wavelengths and the ability to fabricate devices either on top of area CCDs for higher fill factor or to process on glass or ceramic for very large linear sensors. The disadvantages include charge trapping at defects in the amorphous silicon and low carrier mobility, both of which lead to field-dependent nonlinear response and image lag when used in an image sensor. Improvements in material and device technology have mitigated many of these disadvantages at the expense of process and device complexity.⁷

Antiblooming in Charge-Sensing Elements

Blooming in image sensors occurs when the charge generated in an image-sensing element exceeds its capacity. If no method is provided to remove this excess charge, it will be injected either onto the readout element (CCD or MOS readout line) or into the substrate. If the excess charge is injected onto the readout element, it will usually appear as a bright line in the image. If it is injected into the substrate, the charge can diffuse in a circular pattern and be collected by neighboring elements.

There are two basic types of antiblooming circuits: lateral and vertical.^{8,*} In lateral antiblooming, illustrated in Fig. 4a and b, a MOS antiblooming gate and an antiblooming drain are provided adjacent to each image-sensing element. Excess charge on the sensing element overflows the antiblooming gate onto the antiblooming drain. The antiblooming drains of all elements on the array are connected and the current sunk in a bias supply.

In a vertical antiblooming structure, illustrated in Fig. 4c and d for a pinned photodiode sensing element, the image-sensing element is fabricated in a shallow, lightly doped *p*-well. A large (10- to 30-V) bias is applied to the *n*-type silicon substrate, causing the *p*-well underneath the photodiode to become completely depleted. Once the charge on the diode exceeds its capacity, the excess charge flows over the saddle-point in the *p*-well, under the diode, and into the substrate. The substrate is connected to a bias supply which sinks the blooming current. The vertical antiblooming structure has the advantage of requiring no additional silicon area, so that antiblooming is achieved with no reduction in fill factor. Lateral antiblooming, on the other hand, requires additional silicon area in each pixel, thereby reducing fill factor. The vertical antiblooming has the disadvantage of lower quantum efficiency at wavelengths longer than about 500 nm. Because any light absorbed below the photodiode or photocapacitor is drained into the substrate, the quantum efficiency of photodiodes or photocapacitors with vertical antiblooming is reduced.

*Other types of antiblooming are occasionally used in image-sensing arrays. One of these is charge pumping, in which a MOS gate is repeatedly clocked in order to cause excess charge held underneath it to recombine at interface states at the silicon-silicon dioxide interface. In charge pumping, the MOS gate is pulsed sufficiently negative to cause accumulation, resulting in the interface states filling with holes. When the gate is pulsed out of accumulation, excess electrons can recombine with the holes.

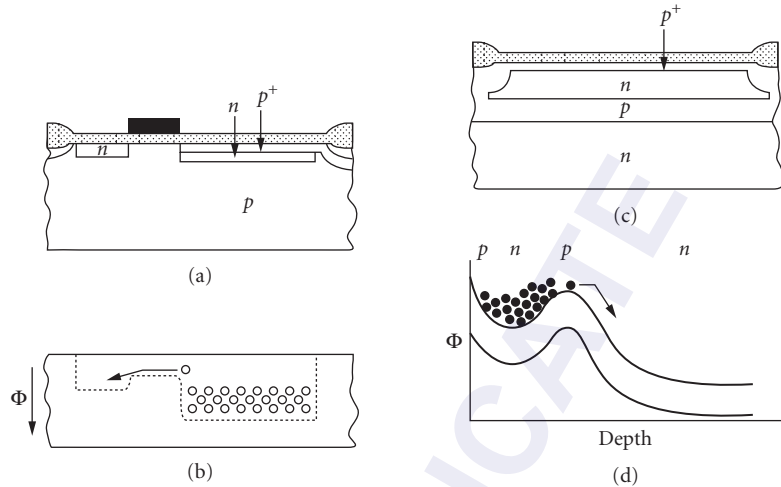


FIGURE 4 Antiblooming methods for image sensors: (a) cross section of lateral antiblooming structure; (b) illustration of electrostatic potential and charge overflow in lateral antiblooming; (c) cross section of vertical antiblooming; and (d) band diagram and illustration of charge overflow in vertical antiblooming.

Dark Current in Photosensing Elements

Signal in photosensing elements is the result of collection of electron-hole pairs generated by the absorption of light. However, charge is generated at each photosensing element even in the absence of light. This generation is called dark current and is a result of thermal generation of electron-hole pairs. The thermal generation occurs at defects, such as impurities or crystalline defects, in the bulk of silicon and at surface states at the silicon-silicon dioxide interface.

There are four sources of dark current: (1) diffusion current, which is the thermal generation of carriers in the undepleted n - and p -type regions; (2) depletion layer generation current, which occurs in the depletion layer of a diode or MOS capacitor; (3) surface generation current, which is the generation of electron-hole pairs at interface states at the Si-SiO₂ interface; and (4) leakage, which refers to generation at extended defects such as impurity clusters or stacking faults, particularly in the presence of a large electric field. These sources of dark current are illustrated for a photodiode and a photocapacitor in Fig. 5. The generation of the electron-hole pairs in both diffusion current and depletion-layer generation-recombination current occurs almost exclusively at impurity sites. Impurities with energy levels near midgap, such as gold, copper, and iron, are particularly effective in the thermal generation of charge. The depletion-layer generation current is given by*

$$J_g(V) = \frac{qn_i W(V)}{\tau_o} \quad (9)$$

where $W(V)$ is the width of the depletion layer at a given bias voltage in the MOS capacitor or junction-photodiode, n_i is the intrinsic carrier concentration, and τ_o is the depletion-layer lifetime, for carriers.

*The expressions here are for the generation current. The total current is the sum of the generation and recombination current given by

$$J_{gr}(V) = (qn_i W(V)/\tau_o)(e^{qV/2kT} - 1)$$

However, for diodes under 0.2 V or more of reverse bias ($qV/kT > 8$), such as would be the case in an image sensor, the recombination current is negligible.

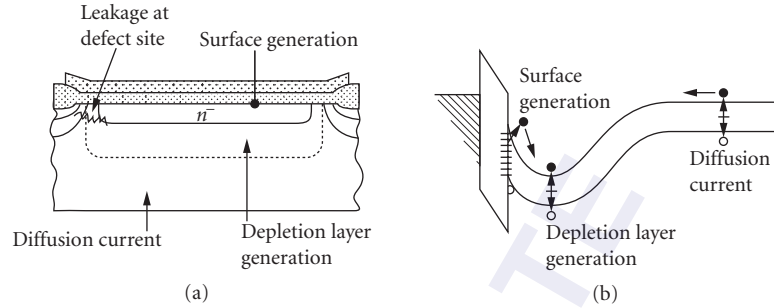


FIGURE 5 (a) Cross section of buried-channel MOS capacitor and (b) band diagram of buried-channel MOS capacitor illustrating mechanisms for dark current generation.

Typical values of τ_0 in clean MOS processes would be 1 to 10 ms and values of J_{gr} would be 30 to 300 pA/cm².

The diffusion generation current for a *p*-type region is given by*

$$J_{\text{diff}} = \frac{qn_i^2 L_e}{N_a \tau_e} \quad (10)$$

Since the intrinsic carrier concentration n_i depends on temperature as $e^{-E_g/2kT}$, depletion-layer generation current and diffusion current will have different temperature dependences. Depletion-layer generation current will increase as $e^{-E_g/2kT}$, which corresponds to a doubling in dark current for every 9 to 11°C increase in temperature near room temperature. Diffusion current will increase as $e^{-E_s/2kT}$, which corresponds to a doubling every 4.5 to 5.5°C increase in temperature.

The surface generation current J_s occurs almost exclusively at regions where the depletion layer intersects the Si-SiO₂ interface, such as the surface region between the *n*⁺- and *p*-regions around a photodiode or in the depleted surface under a MOS capacitor. It is given by

$$J_s = \frac{qn_i s_o}{2} \quad (11)$$

where s_o is the surface recombination velocity. A typical value of J_s for a MOS surface would be 100 pA/cm².^{9†} The surface generation depends on temperature in the same manner as the depletion-layer current. In very clean MOS processes (i.e., low concentration of metallic impurities in silicon), the surface current is often the largest component of the overall dark current.

Leakage current occurs at extended defects in silicon, such as impurity clusters and stacking faults, particularly when these defects are in a depletion layer and so are subject to a high electric field. While there is no single analytical expression for the current generated by an extended defect,

*The full expression for the diffusion current from the undepleted *p*-region in a diode is

$$J_{\text{diff}} = \frac{qn_i^2 L_e}{N_a \tau_e} (e^{qV/kt} - 1)$$

However, for reverse biases more than 0.2 V ($qV/kt > 8$), as would be the case in a photodiode in an image sensor, only the generation term is important.

[†]In calculating the total surface current, the current density J_s is multiplied by the area of the depleted surface. For a *pn* junction diode, this would be the area of depleted surface region around the diode between the *n*- and *p*-regions (i.e., approximately the junction perimeter times the depletion-layer width); for a MOS capacitor it would be the entire area under the MOS capacitor unless sufficient sheet charge of electrons had formed at the surface to invert the surface. The surface recombination velocity is given by $S_o = N_{st} v_{th} \sqrt{\sigma_n \sigma_p}$ where N_{st} is the density of surface states near midgap, v_{th} is the thermal velocity (2.0×10^7 cm/s) and σ_n and σ_p are the electron and hole cross sections for midgap surface states. See reference on p. 32.3 for a complete explanation.

they are characterized by very high values of dark current ($\gg 1$ nA/cm²), a very strong dependence on the electric field, and, in regions of high electric field, only a very slight dependence on temperature. In an image sensor, they are visible as “bright spots” in a few isolated pixels against the otherwise low level of background thermal generation.¹⁰

Values for the dark current vary widely because of variation in the amount of impurities in the silicon; the dark current can range from 0.01 nA/cm² in very high quality image sensors to >10 nA/cm² in sensors with significant metallic contamination. The total number of electrons N_e collected in a pixel is

$$N_e = JAT_{\text{int}}/q \quad (12)$$

where J is the total dark current per unit area, A is the area of the pixel, and T_{int} is the integration time. For a 1/2-in format CCD such as would be used in a camcorder, typical values would be a dark current of 0.5 nA/cm², a pixel area of 100 μm^2 , and an integration time of 1/30 second; the number of thermally generated electrons would be 105 electrons per pixel. Image sensors developed for scientific purposes might have a dark current 5 to 10 times lower at room temperature. These same devices might also be operated below room temperature to reduce the thermal generation to levels below one electron per pixel.

There are two types of noise associated with the charge generated by the dark current: shot noise and pattern noise. The shot noise due to the dark current is the square root of the number of dark electrons in a pixel. Pattern noise is due to pixel-to-pixel variations in the dark current and is often highly correlated between neighboring pixels. A numerical value for the dark pattern noise is typically obtained by using the standard deviation of the pixel values in the dark from a large region of an imager, where the values are obtained by averaging over many frames to eliminate shot noise and other temporal noise sources.*

32.4 READOUT ELEMENTS

The readout element transfers the charge from the image-sensing element (photodiode, photocapacitor, or photoconductor) to the output of the image sensor. In a linear sensor, the readout is only in one direction. In an area sensor, both x and y readout is required. There are two basic types of readout elements: charge-coupled devices, or CCDs, and x - y addressed photodiode arrays (typically called MOS arrays owing to the MOS transistors used in the addressing of the pixels). CCDs are by far the most widely used readout elements owing to their very low noise. Nearly all consumer camcorders, facsimile machines, scanners, copiers, and professional and scientific cameras utilize CCDs. Applications of MOS arrays are largely restricted to those where addressing of an individual pixel or subarray is required.

Charge-Coupled Device (CCD)

CCD Operation The CCD works by moving packets of charge physically at or near surface of silicon from the image-sensing element to an output, where the charge packet is converted into a voltage. The CCDs is formed by an array of overlapping MOS capacitors.^{11–13} There are a number of different types of CCDs, including four-phase, three-phase, two-phase, virtual (single-) phase, and ripple-clocked CCDs. The number of phases refers to the number of separately clocked elements

*A histogram of the dark current values of the pixels is rarely gaussian. In most cases, it exhibits an extended tail at high dark current values due to pixels with crystalline defects. The histogram may also exhibit quantization due to pixels with integral number of impurities (i.e., 1, 2, or 3 gold atoms); the quantization may even be used as a signature of the impurity present. See, for example, McColgin et al., “Effects of Deliberate Metal Contamination on CCD Image Sensors,” *Materials Research Society Symposium Proceedings*, 262: 769 (1992) and “Dark Current Quantization in CCD Image Sensors,” *Proc. 1992 International Electron Device Meeting*, Washington, D.C., 113–116 (1992).

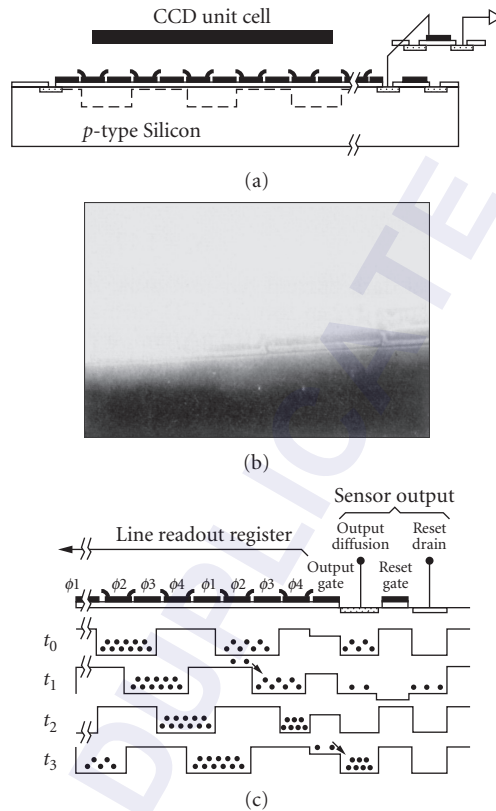


FIGURE 6 (a) Cross-sectional diagram of unit cell of a four-phase CCD; (b) scanning electron micrograph of a unit cell of a four-phase CCD; and (c) illustration of charge transfer along a four-phase CCD, showing both transfer along the register and at the sensor output.

within a single stage of the CCD. These are described later. For understanding the principle of charge transfer in a CCD, consider the four-phase CCD illustrated in Fig. 6a.

The four-phase CCD physically comprises a silicon substrate (assumed to be *p*-type for purposes of illustration), a gate oxide 300 to 1000 Å thick, and overlapping polysilicon electrodes 1000 to 5000 Å thick which have been heavily doped with phosphorus to lower their resistance. For a four-phase CCD, two levels of electrode are required. The first level is deposited, then defined photolithographically to form two phases (ϕ_1 and ϕ_3). A thin (500-Å) oxide is grown over the first level of polysilicon to insulate it from the second level. The second level of polysilicon is then deposited, doped with phosphorus, and defined to form the other two phases (ϕ_2 and ϕ_4). An electron micrograph of a four-phase CCD with six-micron long gates is shown in Fig. 6b.

The process of charge transfer in a four-phase CCD is illustrated in Fig. 6c. In order to hold a packet of electrons, two adjacent gates (ϕ_2 and ϕ_3 , for example) would be held at a high positive potential ($\sim +5$ V) while the other two phases would be held at a low potential (~ 0 V). A depletion layer, or well, is formed under ϕ_2 and ϕ_3 , allowing electrons to be held at or below the surface. The other two phases, ϕ_1 and ϕ_4 , serve as potential barriers, keeping the charge packet under ϕ_2 and ϕ_3 . To transfer the electrons through silicon, the electrode ahead of the charge packet (ϕ_4) is clocked

positive and the electrode behind (ϕ_2) is clocked negative. The electrons move along the silicon surface following the positive potential. This is repeated through all four phases, during which the charge packet is moved forward one pixel.

The CCD may be either surface-channel or buried-channel. In a surface-channel CCD (Figs. 2b and 6a) the electrons are held at the silicon surface. Surface-channel CCDs are rarely used owing to the trapping of electrons at interface states at the silicon surface. At the silicon-SiO₂ interface there is a density of states of about $1 \times 10^{10} \text{cm}^{-2} \text{eV}^{-1}$. These states can trap electrons from one charge packet and reemit the electrons into a later charge packet. This results in transfer inefficiency. In the buried-channel CCD (Figs. 2c and 6d), a lightly doped *n*-type layer is formed in silicon. This *n*-type layer results in a potential well for electrons just below the surface rather than at the surface (Fig. 2b). As a result, the electrons remain separated from the interface and are not trapped by interface states. This results in the ability to transfer charge from one stage to another with very high efficiency. Virtually all CCDs for image-sensing applications utilize buried-channel CCDs. The clock voltages used to drive CCD gates typically swing by 5 to 8 V between high and low levels.

For buried-channel CCDs, transfer rates of up to 20 MPix/s can usually be achieved without special design considerations. Above this rate, special attention must be paid to optimizing the electric field along the direction of transfer to speed up charge transfer. Transfer rates exceeding 50 MPix/s have been achieved in optimized CCD designs.¹⁴

CCD Output At the output of the CCD is a circuit to convert the charge packets into a voltage signal. By far the most common type of output circuit is the floating diffusion with source follower amplifier. The floating diffusion output is shown schematically in Fig. 7a and a photograph of the end of a CCD shift register with the first stage of the amplifier is shown in Fig. 7b. The floating diffusion output consists of an output gate (OG), a floating diffusion, a reset gate (RG), and a reset drain. The floating diffusion is an *n*⁺-type region between the output gate and the reset gate. The floating diffusion is connected to the gate of a source-follower amplifier. The output gate is held at a fixed dc voltage, creating a barrier potential over which the packet of electrons can be transferred onto the floating diffusion when the last phase of the CCD register is clocked to its low (~ 0 V) voltage.

The sequence of events in the CCD output is shown in Fig. 6c. As the charge packet is transferred along the CCD, it arrives at the last phase (ϕ_4) before the output gate (time t_2 in Fig. 6c). When ϕ_4 is clocked low (time t_3), the packet of electrons is transferred over the output gate onto the floating diffusion (time t_0). The voltage of the floating diffusion changes by an amount

$$V = Nq/C \quad (13)$$

where N is the number of electrons and C is the total capacitance of the floating diffusion itself, the interconnect to the source-follower amplifier and the input capacitance of the amplifier. In typical CCDs, this capacitance would be in the 10- to 50-fF range. Because this capacitance is so small, there is a large change in voltage for a small change in charge. The charge-to-voltage conversion is an important parameter in CCD design; for a 10-fF capacitance the charge-to-voltage conversion is 16 $\mu\text{V}/\text{electron}$. After the change in voltage has been sensed by the amplifier, the charge packet must be removed from the floating diffusion before the next packet arrives. This is achieved by clocking the reset gate positive (time t_1 in Fig. 6c), allowing the charge packet to flow from the floating diffusion to the reset drain. The reset drain is held at a constant positive voltage, typically ~ 10 V. The reset drain is then turned off and the floating diffusion is prepared to accept the next packet of electrons.

The change in voltage from the floating diffusion is typically buffered on-chip by a source-follower (Fig. 7c). A two-stage source follower is most often used. The first stage utilizes very small-dimensioned FET transistors in order to minimize the input capacitance. The second uses much larger FET transistors in order to achieve sufficient drive current to overcome off-chip capacitances such as package and lead capacitances on the circuit board. The source-follower amplifier is typically designed to have a bandwidth on the order of ten times the CCD pixel rate. For high-pixel-rate applications (> 10 MPix/s), three-stage source-follower amplifiers are employed.

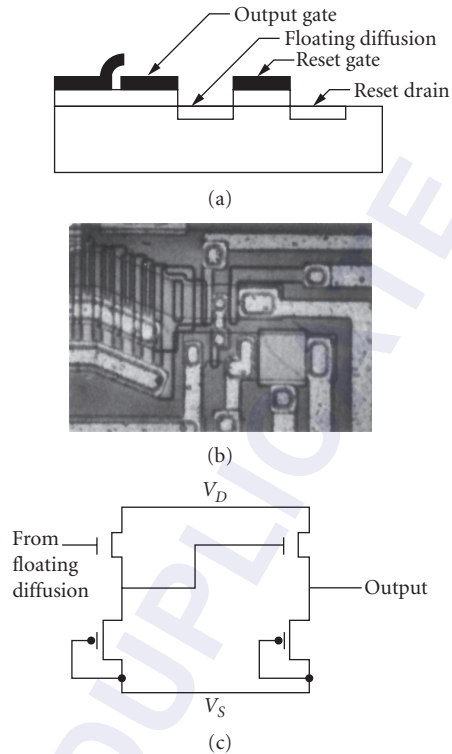


FIGURE 7 (a) Cross section of the floating diffusion output for a CCD; (b) photomicrograph of a floating diffusion output including the first stage of the source-follower amplifier; and (c) two-stage source-follower amplifier with on-chip loads.

For special purpose applications, such as CCD signal processing, nondestructive readout is required. Floating gate output amplifiers are used in these applications. These amplifiers are similar to the floating diffusion except that the floating diffusion is replaced by a MOS gate which is connected to the MOS amplifier. Other outputs, such as buried-channel JFET structures, have seen limited use in very low noise applications.¹⁵

Types of CCDs In addition to the four-phase CCD described here, there are a number of other types of CCD shift registers, including three-phase, two-phase, and virtual phase. The different types are illustrated in Fig. 8. The four-phase CCD (Fig. 8a) was described previously. The three-phase (Fig. 8b) consists of three different layers of polysilicon electrodes. The charge is normally held under one of the three; during transfer, the gate ahead of the charge packet is clocked positive and the gate holding the charge is clocked negative in order to transfer the charge one gate ahead. The three-phase CCD has the advantage of a shorter unit cell than the four-phase but at the expense of additional processing complexity (i.e., a third polysilicon layer).

In the two-phase CCD, each phase has a barrier and a well region. The barrier is formed by doping the barrier region slightly less n -type than the well region, making the electrostatic potential in the barrier region a few volts lower than the well region for the same gate voltage. Electrons will flow over the barrier region and be held in the well region. There are two methods of fabricating a two-phase CCD. In the first method (Fig. 8c), two separate gates are used for each phase; one gate receives

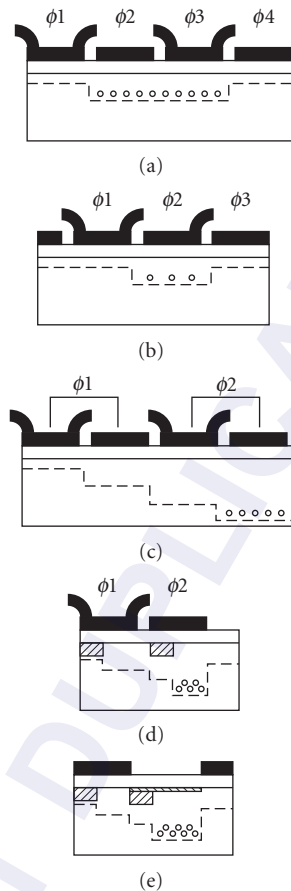


FIGURE 8 Types of CCD registers: (a) four-phase; (b) three-phase; (c) pseudo-two-phase; (d) true two-phase; and (e) virtual phase. In the pseudo-two-phase CCD, each phase consists of two polysilicon gates, one of which is offset in potential from the other by an implanted threshold adjustment. In the true two-phase CCD, each phase consists of a single polysilicon gate in which the implanted threshold-adjust region is formed underneath a portion of the gate.

a threshold-adjust ion implantation in order to create the barrier region. The two gates, however, are connected and thus driven at the same voltage. In the other method (Fig. 8d), the barrier and well regions are created under a single polysilicon gate. The two-phase CCD is very commonly used for three reasons. First, it requires only two clocks (ϕ_1 and ϕ_2), simplifying drive circuitry. Second, the clocks are complementary. This reduces clock feedthrough. Third, the two-phase has a higher horizontal density, especially the two-phase with barrier and well regions under the same gate.

The virtual phase CCD^{16,17} (Fig. 8e) consists of one gate with both barrier and well regions within it and a second region, the virtual phase, in which a shallow high-dose ion implantation is used to create a heavily doped surface region which pins the surface potential at 0 V. The virtual phase has both barrier and well regions within it and acts the same as a polysilicon gate held at a constant voltage. Charge is first transferred from the clocked phase into the virtual phase, then the charge is transferred from the virtual phase into the next clocked phase. The virtual phase CCD has a number of advantages, including higher quantum efficiency than two-polysilicon-level area CCD sensors and simpler clocking. The disadvantages of the virtual phase include the need for larger voltage

swings on the single-clock and high-clock feedthrough into the output due to the lack of complementary clocks.

CCD Characteristics The four major performance parameters of a CCD shift register and output amplifier are charge-handling capacity, charge-transfer efficiency, charge-to-voltage conversion ratio, and noise. The charge capacity is the number of electrons which can be held and transferred in the CCD shift register. As charge is added to a shift register, a point is reached at which excess charge cannot be held; the excess charge either overflows into adjacent pixels or overflows into the bulk beneath the CCD or, in the case of a buried-channel CCD, overflows the barrier to the Si-SiO₂ surface. The charge capacity is a function of the device design, device layout, and CCD process. Figure 9a shows the electrostatic potential and charge distribution in a buried-channel CCD at three levels of charge: approximately one-quarter of saturation, at saturation, and beyond saturation. Below saturation, the electrons fill the center of the buried channel in the region of largest potential, separated from the surface by approximately 0.2 μm in distance and about 500 mV in potential. As more charge is added, the electron distribution spreads toward the surface and the potential barrier to the surface drops. Beyond saturation, the electrons contact the surface directly, resulting in charge-transfer inefficiency and blooming to neighboring pixels. The charge capacity of most CCD shift registers is of the order 1×10^{12} electrons/cm² of area in which the charge is held. In most CCD cells, the charge is held in only a fraction of the total cell area both along and across the CCD register. Typically, area CCD image sensors are designed with CCD charge capacities in the range of 50,000 to 200,000 electrons. Linear CCD image sensors, because of the larger amount of silicon area available for the CCD shift register, often are designed for 100,000 to 1,000,000 electrons.

The second major performance parameter for CCD shift registers is charge transfer efficiency. The transfer of charge from one stage to the next is neither instantaneous nor complete, limiting both transfer rate and the total number of stages in the CCD. There are two intrinsic mechanisms governing charge transfer: drift and diffusion. Drift is the movement of charge in the presence of an electric field. There are two origins of the electric field seen by an electron during charge transfer:

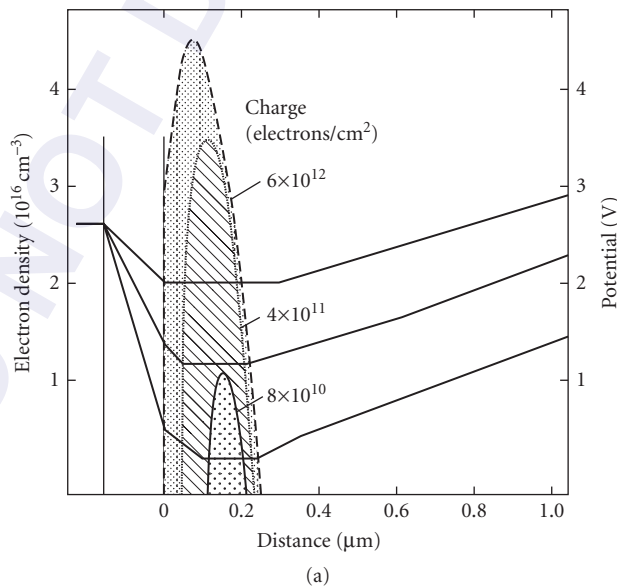
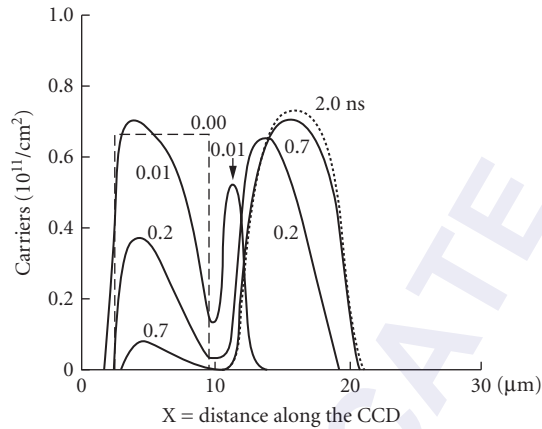
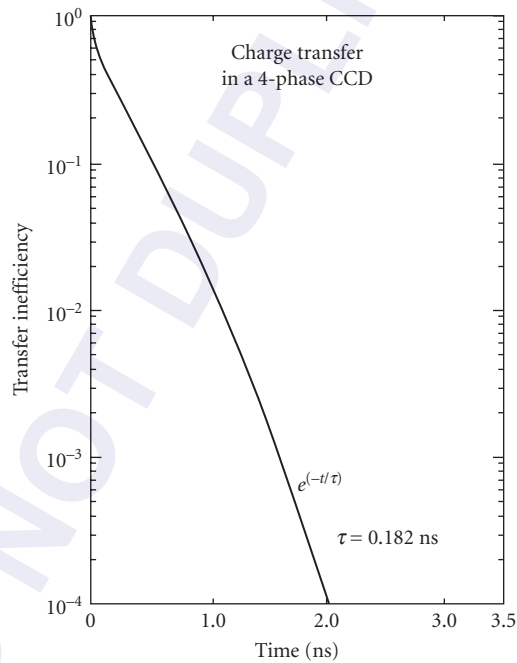


FIGURE 9 (a) Charge distribution and electrostatic potential as a function of distance from the surface for a buried-channel CCD for three levels of charge density.



(b)



(c)

FIGURE 9 (b) Charge density as a function of distance along a CCD for various times following the beginning of charge transfer from one 8- μm stage to another and (c) charge transfer inefficiency as a function of time from the start of charge transfer for the same example. (Continued)

the self-induced field resulting from the other electrons under the gate and the externally induced or fringing field. During the early stages of charge transfer, the self-induced field is large and is the dominant factor. After the charge concentration under the gate has decreased to a low value, the remainder of the charge transfer will be governed either by diffusion or by drift due to externally induced or fringing fields.

For a surface-channel CCD, the self-induced fields can be estimated from the formula:

$$E = -\frac{q}{C} \frac{dN}{dx} \quad (14)$$

where C is the gate capacitance per unit area and $N(x)$ is the density of electrons as a function of distance x from the edge of the electrode. In the early stages of charge transfer, both N and dN/dx are large. As the transfer proceeds, both N and dN/dx become small and the charge transfer is governed by fringing fields and diffusion.

Fringing fields are due to the two-dimensional nature of the electrostatic potential. If the charge is being transferred from gate 1 to 2, the potential will change smoothly between the two gates. The effect of the potential from one gate will typically extend 1 to 3 μm into a neighboring gate. The charge within the range of the fringing field will move by drift to the neighboring gate. Charge out of the range of the fringing fields will move by diffusion.

Figure 9b shows an example of charge transfer calculated for charge transfer from one 8- μm -long CCD gate into an adjacent gate. The charge density is shown as a function of distance along the CCD at several times after the start of charge transfer. At the start of the transfer, all the charge is under the first gate. At 0.01 ns into the transfer, the charge has moved from the edge of the first gate into the second gate, as a result of self-induced drift. By 0.2 ns, approximately half the charge has been transferred. By 0.7 ns, over 90 percent of the charge has moved into the second gate, leaving a residual in the first. At this point, the self-induced drift is sufficiently small that drift due to fringing fields and diffusion are the dominant mechanisms. Figure 9c shows the charge transfer inefficiency as a function of time for this example. Two slopes are evident in the transfer inefficiency. In the first 0.5 ns, the charge is transferred rapidly owing to the self-induced drift. For times longer than 0.7 ns, the transfer is due to fringing fields in this example.

For CCDs with longer gates or lower fringing fields than the example above, the final ~10 percent of the charge must transfer by diffusion. Charge transfer by diffusion follows an exponential time dependence.

$$N(t) = N(0)e^{-t/\tau_{\text{diff}}} \quad (15)$$

where, the time constant τ_{diff} for diffusion is

$$\tau_{\text{diff}} = \frac{4L_g^2}{\pi^2 D} \quad (16)$$

Here, L_g is one gate length and $D = KT\mu_e/q$.

For electrons at room temperature, $D = 25.8 \text{ cm}^2/\text{s}$. For an 8- μm long gate, the diffusion time constant is ~10 ns. To achieve a transfer inefficiency below 2×10^{-5} per transfer, 11 time-constants, or 110 ns in this example, are required. For this reason CCDs are typically designed with gate lengths shorter than ~8 μm and to build-in fringing fields.

In the simplest case for very low levels of transfer inefficiency, the total transfer inefficiency in a CCD register is the product of the number of stages N in the register and the transfer inefficiency per stage. Each stage will require two or more transfers. Virtual phase and two-phase require two transfers per stage, three-phase requires three transfers, and four-phase requires four transfers. The inefficiency per stage, however, will likely depend in a complicated manner on the amount of charge in the charge packet, the amount of charge in preceding charge packets, the voltages, and frequency of operation. The charge in the preceding packets will affect the filling of both bulk and interface traps.

In addition to the intrinsic sources of transfer inefficiency, there are a variety of extrinsic sources. These include surface and bulk traps and potential wells and barriers. The traps and the potential obstacles hold back an amount of charge from a charge packet. The charge is reemitted to later charge packets. The inefficiency due to traps depends on whether the traps have been filled by preceding charge packets and the emission time constants of the traps, and so is not modeled in a simple manner.

The intrinsic sources of noise in CCD shift registers include dark current and output amplifier noise. In addition, other sources of noise not intrinsic to the CCD itself include noise due to clock feedthrough from the CCD clocks to the output signal and noise in external electronics. The generation of dark current in CCD shift registers is the same as described at the end of Sec. 32.3 for photosensing elements. Associated with the dark current are both shot noise and pattern noise. The magnitude of the pattern noise in a CCD shift register is reduced over that of a single element since the charge packet averages the dark current over many pixels as it is transferred to the output.

The noise associated with the CCD output consists of the Johnson or thermal noise, the $1/f$ noise of the output amplifier, and the kTC noise associated with resetting the floating diffusion. The kTC noise is due to uncertainty in the amount of charge remaining on the floating diffusion following reset owing to thermal fluctuations in the reset gate. The rms noise σ_n in the number of electrons caused by kTC noise is given by

$$\sigma_n = \sqrt{\frac{kT}{q}} C_{\text{FD}} \quad (17)$$

where C_{FD} is one total floating diffusion capacitance. For a 10-fF floating diffusion capacitance, the rms noise is ~ 40 electrons.

The kTC noise may be eliminated entirely by use of a signal-processing technique called correlated double sampling (CDS). In correlated double sampling, the output level is sampled before the charge packet is transferred onto the floating diffusion and again after transfer of the charge packet (times t_2 and t_0 respectively in Fig. 6c). The two values are subtracted either by an analog circuit or by digital subtraction. Any uncertainty in the voltage level of the floating diffusion following reset is subtracted and thus the kTC noise eliminated. The output amplifier low-frequency noise, called $1/f$ noise because of the inverse frequency ($1/f$) shape of the noise power spectrum, is also reduced but not eliminated by correlated double sampling. The total noise of a CCD output amplifier is in the range of 7 to 40 rms electrons per pixel depending on the amplifier design and the operating speed. Values of a single rms electron or less have been obtained for very slow pixel rates under cooled conditions.¹⁸

MOS Readout

The other major category of readout structures is the MOS readout. The individual light-sensing elements (photodiodes, photocapacitors, or photoconductors) at each pixel are connected to a readout line by means of a transfer gate. Each pixel along the readout line is addressed separately by addressing circuitry. When a particular pixel is addressed, the transfer gate is turned on and the charge transferred from the pixel to the readout line. An amplifier at the end of the readout line senses the change in voltage or current resulting from the charge transfer.

Typically, the pixels would be addressed serially along the line. The first pixel would be addressed, causing the charge from the image-sensing element to be transferred onto the readout line. The voltage change or current would be sensed, the readout line reset to its original voltage if necessary, and the next pixel addressed. This is different from a CCD. In the CCD, charge from all pixels is transferred into the CCD register simultaneously. Individual pixels or groups of pixels cannot easily be addressed in a random fashion by the CCD, but this random addressing can be accomplished readily by the MOS readout.

There are several types of MOS readout devices. These include the CID,¹⁹ the AMI,²⁰ and the CMD²¹ in addition to the normal MOS array. The CID has no readout line. Each pixel consists of two overlapping gates, one controlled by a row address and the other by a column address.¹¹ When neither a row nor a column of a particular pixel is being addressed, the photogenerated charge is held under both gates and can be transferred between them. When a row of a pixel is addressed, the charge transfers onto the column gate. Then both row and column are addressed, the charge is injected into the substrate, and the current sensed. The CID is not widely used in visible imaging applications because the charge conversion sensitivity is very poor and noise very high compared to the CCD or the other MOS architectures.

In the amplified MOS imager (AMI), the image-sensing element at each pixel consists of a phototransistor rather than a simple photodiode. The photogenerated charge is stored on the gate of the MOS transistor.¹⁹ When a particular pixel is addressed, the photogenerated charge modulates the transistor current. This current amplification at each pixel helps to overcome many of the noise and speed limitations of conventional MOS arrays.

MOS readout differs in an important way from CCD readout. In MOS readout, the charge is transferred from a single pixel onto a readout line and the change in voltage or current in the readout line is sensed. In a CCD, the charge packets are kept intact while being transferred physically to a low-capacitance output. The lower sensitivity of a simple MOS array can be illustrated as follows. The change in voltage on the readout line is given by $V = Nq/C$ where N is the number of electrons, and C is the readout line capacitance. Because the readout line covers the full length of the array, its capacitance is in the picofarad range (typically 2 to 10 pF depending on design and process). This compares to the 10-fF capacitance for the CCD output. As a result, the voltage swings on the readout line are very small (16 nV/electron for a 10-pF readout line capacitance). This leads to a high sensitivity to clock noise due to capacitive feedthrough of the row and column address clocks onto the readout line. The feedthrough may be many times larger than the signal in most MOS sensors. Once the charge has been transferred onto the readout line, it is sensed either by a current-sensitive amplifier or by a voltage-sensitive amplifier, followed by a reset of the readout line to its original voltage.

CCD readout has the advantages of very high sensitivity and low noise. However, CCD readouts are limited in charge-handling capacity, while MOS readouts are capable of carrying very large amounts of charge and so are not as limited on the high end of the dynamic range. However, because the MOS readout line has much higher capacitance than the CCD, the sensitivity is lower and the noise is higher. Another difference is in the readout architecture. The CCD readout is essentially serial and not suited to random readout or partial-array readout. The MOS array, however, can be addressed in a manner similar to a memory, making it well-suited to pixel or partial-array addressing.

32.5 SENSOR ARCHITECTURES

Solid-state image sensors are classified into two basic groupings: linear and area. Linear sensors include single-line arrays, multilinear arrays for color scanning, and time delay and integrate (TDI) arrays for low-light-level scanning. Area sensor architectures include the frame transfer CCD, the interline transfer CCD, and various forms of MOS x - y addressed arrays.

Linear Image Sensor Arrays

Linear sensors are used almost exclusively in scanning systems for scanning documents, film, and three-dimensional still objects. There are two basic classes of scanning systems: contact scanners and reduction scanners. These are illustrated in Figs. 10*a* and *b*. In reduction scanners (Fig. 10*a*), the sensor is smaller than the document to be scanned; lenses are used to image the document onto the sensor. In contact scanners (Fig. 10*b*), the sensor is the same width as the item to be scanned, usually a document. Relay optics is used between the sensor and the document. Selfoc lenses (Fig. 10*c* and *d*) and roof-mirror-lens arrays are the two types of relay optics used most frequently.

There are three basic architectures for linear sensing arrays: MOS line arrays, CCD linear and multilinear sensors, and time-delay and integrate (TDI) sensors. These architectures are illustrated in Fig. 11. The MOS array is used most often in contact scanning applications where material or processing problems make CCD arrays impractical. These applications include arrays fabricated from polysilicon or amorphous silicon on nonsilicon substrates, arrays covering large distances, or arrays requiring special processing (such as logarithmic amplification) at each pixel. The CCD linear and multilinear arrays are used most often in reduction scanning where wide dynamic range

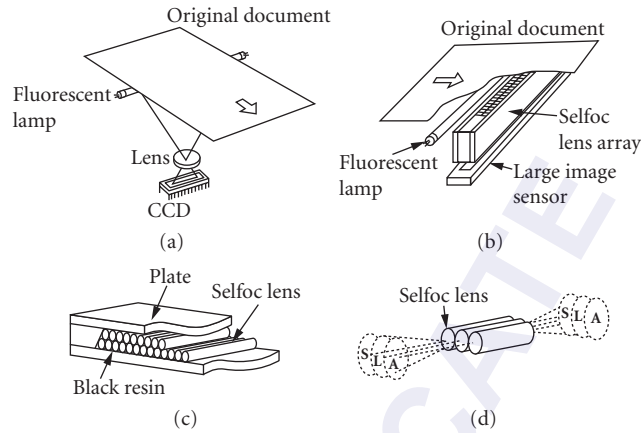


FIGURE 10 (a) Reduction document scanner using linear image sensor, in which the image on the page is reduced by a lens onto the line array; (b) contact document scanner in which the length of the image sensor is the same as the document width and one-to-one relay optics is used to transfer the image from the document to the array; (c) diagram of a selffoc lens array used as transfer optics between document and array in a contact scanner; and (d) operation of a selffoc lens array.

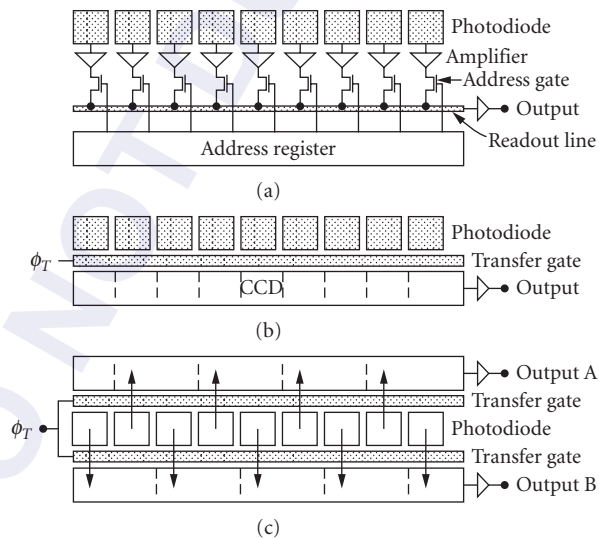


FIGURE 11 Architectures for linear image sensors: (a) MOS line array consisting of photodiodes, preamplifier, MOS switches addressed by an address register, and a readout line with amplifier; (b) linear CCD image sensor consisting of photo-diodes, transfer gate, and CCD readout; and (c) linear CCD image sensor with two CCD output registers, one for the odd diodes and the other for the even diodes, for higher horizontal pitch.

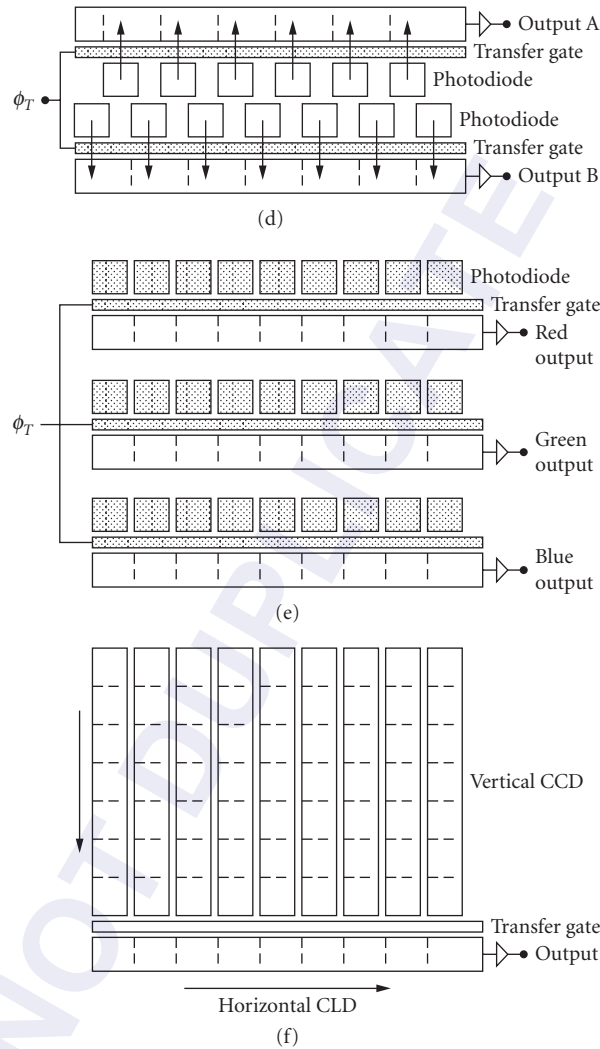


FIGURE 11 (d) Staggered linear CCD image sensor with two rows of photodiodes offset by one-half pixel to increase horizontal sampling; (e) trilinear CCD in which three CCDs are fabricated on the same silicon die, each with its own color filter; and (f) time-delay and integrate (TDI) array, in which the charge in the vertical registers is clocked in phase with the motion of the scene or document being imaged in order to increase signal-to-noise. (*Continued*)

and small pixel size are required. However, contact arrays are also often realized by butting multiple CCD arrays end-to-end. TDI arrays are used in very low-light-level scanning applications where integration over many lines is required to achieve adequate signal-to-noise ratio.

The MOS linear array (Fig. 11a) consists of individual photosensing elements, an amplifier, an address switch, and a readout line and amplifier. The photosensing element is usually a photodiode,

although photoconductors are used in contact scanning arrays fabricated from amorphous silicon. Since the charge generated in the diode is generally very small (in the hundreds to thousands of electrons), a simple amplifier is usually placed at each pixel to drive the high-capacitance readout line. The use of amplification at each pixel can allow some signal-processing functions, such as logarithmic amplification, clipping, triggering and latching, etc., to be performed at the pixel. A MOS switch placed after the amplifier allows each pixel to be addressed in sequence; the switch is driven by an address register. At the end of the readout line is an amplifier which may buffer and/or amplify the signal. The MOS array has the advantage of process simplicity and the ability to perform signal processing at each pixel; it has the disadvantage of low signal level (because of the large readout line capacitance) and pattern noise introduced by feedthrough from the switching transistor.

The linear CCD image sensor is the most often used architecture for scanning applications owing to its low noise, high sensitivity, wide dynamic range, and small pixel pitch. Figure 11*b* shows the simplest type of linear CCD, in which a single row of photodiodes is connected to a single CCD register via a transfer gate. In operation, the signal is integrated on the image-sensing element (generally a photodiode) for a line time. The horizontal CCD is then stopped, the transfer gate opened, and the charge transferred from all the photodiodes simultaneously to the CCD. The transfer time is typically a few microseconds. The transfer gate is then closed, integration resumed for the next line, and the CCD clocked to read out the charge packets. Many arrays also feature antiblooming for situations where the light level may not be controlled, as well as electronic shuttering, which allows an integration time on the photodiodes to be less than the readout time of the CCD.

For linear-sensing applications requiring a higher pixel density, a double-sided readout is often used (Fig. 11*c*). In this architecture, a CCD array is placed on either side of the line of photodiodes and charge transfer from the diodes alternates between the top and bottom CCDs. This architecture uses lower horizontal clock rates and a higher pixel pitch, since the diode pitch can usually be made smaller than the CCD pitch. The charge packets from the two arrays may be multiplexed into one output if desired. The disadvantage of this architecture is the slight differences between even and odd pixels, due to slight differences in the two outputs (or slight differences due to multiplexing the two registers).

Another architecture which is used to further decrease the sampling pitch is the staggered linear array (Fig. 11*d*). In the staggered array, two rows of photodiodes are offset by a half pixel. The two rows are read out by CCD arrays. The first array is delayed (usually in a digital line store) and then combined with the second to form a double-density scan.

Multilinear arrays (Fig. 11*d*) have been developed for color scanning applications. In this architecture, several (usually three) linear arrays are combined on the same silicon die separated by a distance equivalent to an integral number of scan lines. Color filters (either integral or in close proximity) are aligned over the arrays. External digital line delays are used to realign the three arrays. Separate electronic shuttering may be provided for each array in order to adjust for differences in intensity in each of the bands.

The third major class of line arrays is time-delay and integrate, or TDI, arrays.^{22,23} The TDI architecture is shown in Fig. 11*f*. TDI arrays are used when inadequate signal-to-noise ratio from a single-line array requires averaging over multiple lines. Applications for TDI arrays include high-speed document scanners and space-based imaging systems. Instead of a single row of photodiodes, the TDI array utilizes CCD stages in the vertical dimension which are clocked synchronously with the movement of the document to be scanned. The signal level in a TDI array increases linearly with the number of stages. The noise level, however, increases most as the square root of the number of stages.

Area Image Sensor Arrays

There are three major classes of area image sensor architectures: MOS diode arrays, frame-transfer CCDs, and interline-transfer CCDs. Within each there are a number of variations. CCDs have come to dominate the majority of applications owing to their higher sensitivity. However, MOS arrays are still used for specialized applications where addressability or high readout rate is important.

TABLE 1 Image Area Dimensions and Pixel Dimensions for Various Format Image Sensors

The format name is given in inches based on historic image tube formats. The pixel dimensions are based on a 484×768 pixel array.

Optical format	1 in	2/3 in	1/2 in	1/3 in	1/4 in
Active area (4:3 aspect ratio)					
Height (mm)	9.6	6.6	4.8	3.3	2.4
Width (mm)	12.8	8.8	6.4	4.4	3.2
Diagonal (mm)	16	11	8	5.5	4
Pixel dimensions ($484 \text{ lines} \times 768 \text{ pixels}$)					
Height (μm)	19.8	13.6	9.9	6.8	4.9
Width (μm)	16.7	11.5	8.3	5.7	4.2

Historically, the physical dimensions of the active imaging areas of CCD arrays for consumer and commercial applications are specified by the size of the vidicon tube which it replaces. The common format sizes include 1/4, 1/3, 1/2, and 1 in. In most cases, the aspect ratio is 4:3, reflecting the television standard. Table 1 lists the formats and the corresponding dimensions of the imaging area of the array.

The standards for the number of vertical lines in arrays for consumer and professional video are usually based on the corresponding television standards, including NTSC, PAL, and the various HDTV standards. NTSC has 484 active lines, PAL has 575, the Japanese HDTV standard has 1035, and the European HDTV standard has 1150. In all cases the lines are interlaced; i.e., odd lines are read out in one field and even lines in the next. Historically, for sensors for NTSC television, the number of pixels horizontally in a line has been associated with multiples of the color subcarrier frequency; common horizontal pixel counts include 384, 576, and 768. Sensors for the Japanese HDTV standard typically have 1920 horizontal pixels. The pixels are rectangular rather than square. Table 1 lists approximate pixel dimensions in micrometers for a few common formats for a $484(\text{V}) \times 768(\text{H})$ pixel image sensor.

For industrial, scientific, graphics electronic photography, digital television, and multimedia applications, however, the nonsquare pixels and interlaced readout of sensors based on television standards are significant disadvantages. Image sensors designed for these industrial and commercial applications typically have square pixels and progressive scan readout. In interlaced NTSC readout, for example, the first field consists of lines 1, 3, 5 \dots 483 and is read out in the first 1/60 second of a frame. The second field consists of lines 2, 4, 6 \dots 484 and is read out in the second 1/60 second of a frame. The resulting temporal and spatial displacement between the two fields is undesirable for these applications. In addition, digital compression of interlaced scan moving images, especially with motion estimation, is difficult and also introduces artifacts.

In progressive scan readout each line is read out sequentially. There is no even or odd field, only a single frame. As a result there are no temporal and spatial sampling displacement differences. However, for a given resolution and frame rate, the readout rate of a progressive scan image sensor is double that of an interlaced scan. In addition, for these applications, the number of pixels is often based on powers of 2: such as 512×512 or 1024×1024 —facilitating memory mapping and image processing.

MOS Area Array Image Sensors The architecture of MOS area arrays is illustrated in Fig. 12.²⁴ It consists of the imaging array, vertical and horizontal address registers, and output amplifiers. The pixel of a MOS array consists of an image-sensing element (photodiode, photocopacitor, phototransistor, or photoconductor), a row-address gate, and a vertical readout line. The row-address gate is bussed horizontally across the array and is driven from a row-address register on the side(s) of the array. At the start of a line, a single row is addressed, causing the charge from all the photodiodes in

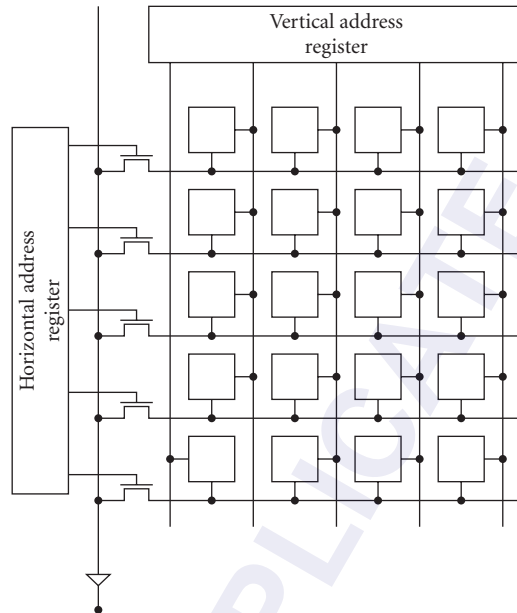


FIGURE 12 MOS photodiode array, consisting of vertical and horizontal address registers and readout line.

a row to be transferred onto the vertical readout line. Horizontal address gates are placed at the bottom of the vertical readout line. Buffer amplifiers to drive the horizontal readout line may also be placed at the end of the vertical readout line. The horizontal address register then serially addresses each vertical readout line, sequentially turning on the horizontal address gates. After the horizontal addressing is completed, the readout lines may be reset and precharged and the next row addressed.

There is a wide range of variations on this basic architecture. An example is the charge modulation device, or CMD,^{25,26} array in which a phototransistor is placed at each pixel site in order to achieve amplification and to achieve high currents to drive the capacitance of the vertical and horizontal readout lines. Other examples include arrays with sophisticated charge collection circuits at the end of each vertical readout line.

Frame Transfer CCD Image Sensors CCD area arrays fall into two categories: frame transfer and interline transfer. The simplest form of frame transfer CCD is the full-frame type, shown in Fig. 13a. A photograph of a few pixels of a frame transfer CCD is shown in Fig. 14a. The array consists of a single image area composed of vertical CCDs and a single horizontal register with an output amplifier at its end. In this architecture, the pixel consists of a single stage of a vertical CCD. This type of device requires an external shutter. When the shutter is opened, the entire surface of the sensor is exposed and the charge is collected in the CCD potential wells at each pixel. After the shutter is closed, the sensor is read out a row at a time by clocking a row of the vertical register into the horizontal register, then clocking the horizontal register to read out the row through the output amplifier. For higher readout rates dual horizontal CCDs are used in parallel. The full-frame CCD has the advantage of progressive scan readout high fill factor, very low noise, and wide dynamic range. However, it requires an external shutter. It is most often used in still electronic photography, scientific, industrial, and graphics applications.

For motion imaging applications, a shutter is not practical. In order to overcome the need for a shutter, frame transfer CCDs incorporate a storage area in addition to imaging area. For interlaced video applications, this storage area is sufficiently large to hold a field (242 lines in NTSC television).

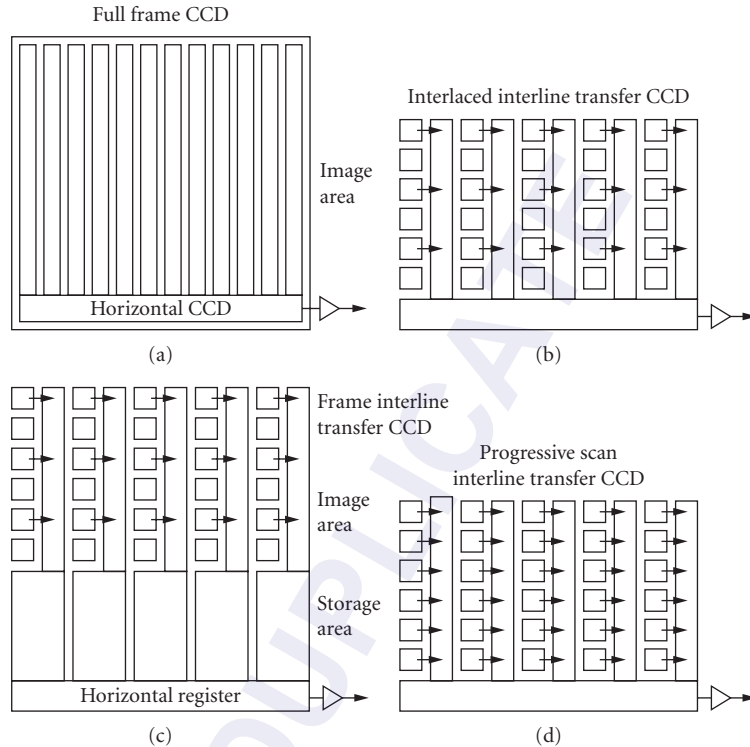


FIGURE 13 (a) Architecture of full-frame CCD; (b) architecture of interlaced interline transfer CCD; (c) architecture of frame interline transfer CCD, in which a storage area is provided to reduce smear during readout; and (d) progressive scan interline transfer CCD, in which every photodiode is read out into the vertical CCD simultaneously.

The image area consists of vertical CCDs. For interlaced NTSC video, there are 242 pixels vertically in the image area. Interlace is achieved by changing the gates under which integration is performed in the even and odd fields. For example, for a four-phase CCD, integration would be performed under phases 1 and 2 in one field and 3 and 4 in another, effectively shifting the sampling area by half a pixel in each field. The storage area consists also of 242 pixels vertically. The device operation is as follows. The image area integrates for a field time and the photogenerated charge held in the vertical CCDs. The vertical CCDs are then clocked in order to rapidly transfer the charge from the image area into the storage area. Because the sensor is still under illumination, this transfer time must be much shorter than the integration time. This transfer typically requires 0.2 to 0.5 ms. The storage area is then read out a row at a time by transferring a row into the horizontal register and clocking this register. While this readout is occurring, the image area is integrating the next field. The most significant disadvantage of frame transfer CCDs is the image smear caused by illumination during the transfer from the image to the storage area. This smear can be on the order of 3 percent.

In both frame transfer and full-frame devices, the light must pass through the polysilicon electrodes before being absorbed in silicon. Owing to the high absorption of short wavelengths in the polysilicon, the quantum efficiency in the blue is only about 20 percent and in the green is about 50 percent. Figure 14b shows the quantum efficiency for a full-frame image sensor with polysilicon gates. Three approaches have been used to improve the efficiency: the virtual phase CCD, transparent electrodes, and backside illumination. In the virtual phase CCD (see “Types of

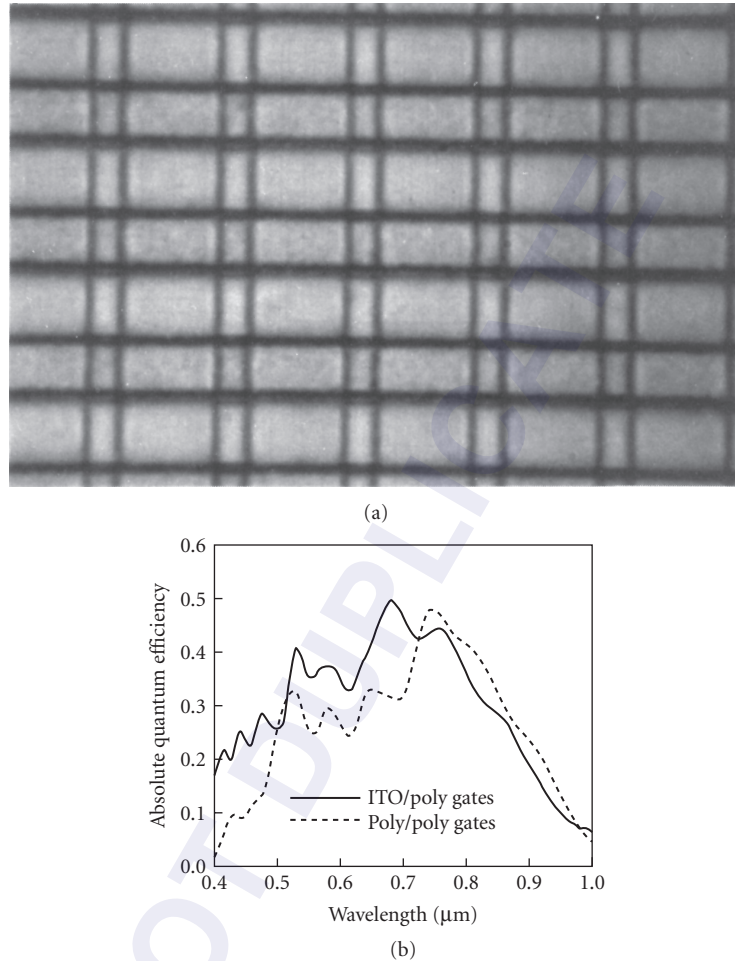


FIGURE 14 (a) Photograph of a few pixels in the image area of a full-frame CCD and (b) quantum efficiency as a function of wavelength for a full-frame CCD with polysilicon and with transparent indium-tin-oxide electrodes.

CCDs” in Sec. 32.4), the second polysilicon electrode is replaced with a very shallow highly doped p^+ layer, very similar to the pinned photodiode (Fig. 2d). Since there is no electrode over this phase to absorb the light, higher quantum efficiency is achieved, particularly at wavelengths less than 500 nm. Backside illumination provides nearly unity quantum efficiency but requires that the sensor be thinned to less than 10 μm . Owing to its cost, backside thinning is employed only in sensors for very specialized scientific or aerospace applications. Indium-tin-oxide, or ITO, is the most commonly used transparent electrode.²⁷ Usually it is substituted for the second level of polysilicon. Figure 14b also shows the quantum efficiency of a full-frame device in which ITO has been substituted for one of the polysilicon gates.

Interline Transfer CCD Image Sensors The interline transfer CCD is fundamentally different from the frame-transfer CCD in that, in addition to the vertical CCD, the pixel also contains a separate image-sensing element (photodiode, pinned photodiode, photocapacitor, or photoconductor) and a

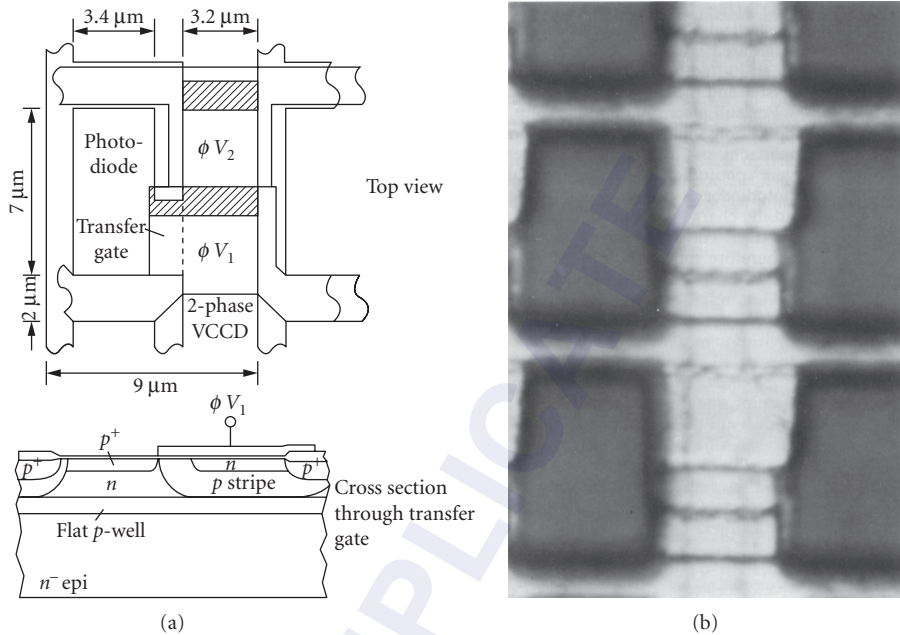


FIGURE 15 (a) Diagram of a pixel of an interline transfer CCD and (b) photograph of a pixel from an interline transfer CCD.

transfer region between the photodiode and the vertical CCD. Interline CCDs with photodiodes for sensing elements are considered first. Figure 13*b* illustrates the architecture of an interline transfer CCD and Fig. 15 illustrates a pixel of the CCD. The CCD and the transfer region between the diode and CCD are covered with a light shield (typically aluminum, although metal silicides are also used). The light shield prevents any light from entering the vertical CCD registers, allowing them to be read out while the sensor is illuminated. When the sensor is illuminated, the photogenerated charge is held on the photodiode. During the vertical retrace interval at the end of a field, the photogenerated charge is transferred into the vertical CCD by clocking the CCD gate over the transfer region. Once the charge has been transferred from the diodes into the vertical CCDs, the diodes resume integrating and the vertical CCDs are clocked in order to transfer a row at a time from the image area into the horizontal CCD.

For consumer and most commercial applications, interlaced interline CCDs are used to be consistent with television standards. In interlaced interline CCDs there is one vertical CCD stage for every two photodiodes. During the retrace time before the first field the charge from the odd rows of photodiodes is transferred into the vertical CCDs; the charge is then transferred out a row at a time into the horizontal CCD. During the retrace time before the second field the charge from the even rows of photodiodes is transferred into the vertical CCDs; once again, the charge is transferred out a row at a time into the horizontal CCD. In NTSC television the fields are approximately 1/60 second long; in PAL the field time is 1/50 second.

Because the vertical CCDs in an interline CCD are covered by a light shield, very little stray light is absorbed in the CCD. However, due to light scattering under the lightshield and lateral diffusion of photogenerated electrons, some charge can reach the vertical registers as they are read out during a field. This results in smear. For consumer applications the level of smear is not noticeable. However, for especially demanding applications such as television broadcast cameras, a field storage area is added to the bottom of the image area. This architecture is called the frame interline transfer, or FIT CCD²⁸ (Fig. 13*c*). Following transfer of the photogenerated charge from the diodes

into the vertical CCDs, the vertical CCDs are clocked to rapidly shift the charge from the image area into the storage area. This transfer typically takes less than 0.5 ms, reducing smear 30-fold. One line at a time is transferred from the storage area to the horizontal register and the readout.

For still electronic photography, scientific, computer-related, graphics and professional applications, interlaced video is not desirable. For these applications a progressive scan architecture is utilized. In a progressive scan interline CCD there is a full CCD stage for every photodiode.²⁸ Following integration, the photogenerated charge from all the photodiodes is transferred into the vertical CCDs. The photodiodes resume integration and the charge packets from the vertical registers are transferred into the horizontal a row at a time. The progressive scan interline CCD requires twice the vertical CCD density and is therefore more complicated to fabricate. However, progressive scan readout provides many advantages in image quality for both motion and still imaging.

Nearly all interline CCDs used in camcorder, broadcast camera, or commercial applications utilize vertical antiblooming in order to prevent blooming when the sensor is illuminated beyond saturation. A cross-section diagram of an interline CCD with vertical antiblooming is shown in Fig. 15*b*. The device illustrated uses a pinned photodiode photosensing element. The CCD is built on an *n*-type silicon substrate. A *p*-well is formed about 2 μm deep in the *n*-type silicon. An *n*-type buried-channel is then formed followed by a *p*+ surface layer. The *n*-type substrate is reverse biased with respect to the *p*-well. When the photodiode is illuminated above saturation, the excess electrons spill out of the *n*-type buried channel and into the substrate. Owing to the vertical overflow, however, photogenerated carriers from photons absorbed below the *p*-well are drawn into the substrate and are not collected by the diode. Thus the quantum efficiency of these devices falls rapidly at wavelengths beyond 550 nm. Figure 16 shows the quantum efficiency of an interline CCD with vertical antiblooming as a function of wavelength. The internal quantum efficiency of the photodiode is nearly 100 percent to about 550 nm, after which it decreases. However, since the photodiode only occupies about 20 percent of the pixel, and this aperture is typically reduced further by the light shield in order to eliminate optical scattering into the vertical CCD, the actual efficiency of the device is only about 15 percent. The photoresponse is linear at low signal levels but becomes nonlinear at charge levels near saturation.²⁹

Two major approaches are used to improve the fill factor (and therefore the quantum efficiency) of interline CCDs. One uses microlens arrays to focus the light incident on a pixel onto the photodiode and the other a vertical integration of image sensors with amorphous silicon photoconductors to achieve a high fill factor. Figure 17 shows a microlens array on top of an interline CCD.³⁰ The

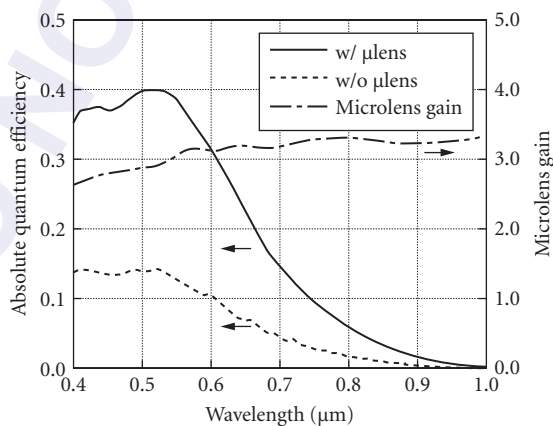


FIGURE 16 Quantum efficiency as a function of wavelength for an interline transfer CCD with and without microlens array.

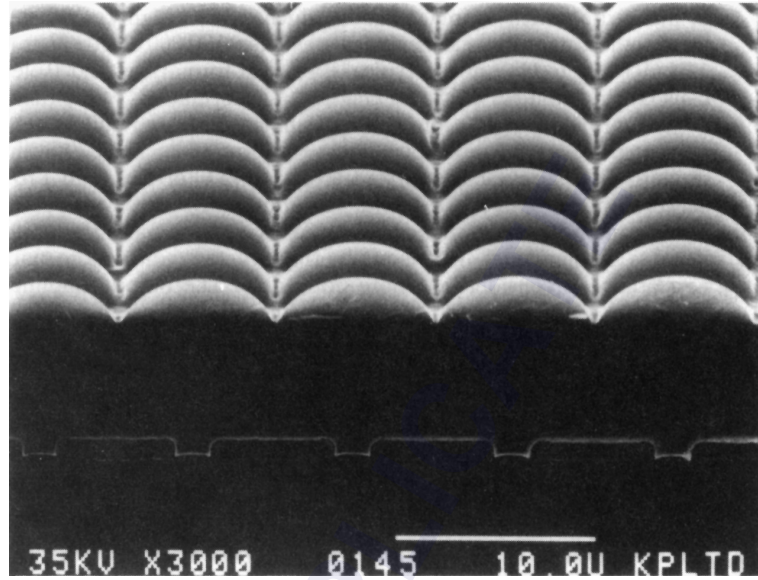


FIGURE 17 Scanning electron micrograph of microlens array on top of an interline transfer CCD.

lenses are formed by coating a spacer layer on the wafer of CCD devices followed by a lens-forming layer. The lens-forming layer is patterned and then reflowed to form the lens arrays. The quantum efficiency of an interline CCD with and without a microlens array is shown in Fig. 18. A 3-fold improvement in quantum efficiency is achieved because light from nearly the entire pixel area is focused onto the photodiode.

The structure of the photoconductor and its band diagram in the second approach⁷ are illustrated in Fig. 2e. The amorphous silicon is about 1 μm thick; owing to its high absorption coefficient in the visible, it can achieve nearly 100 percent internal quantum efficiency. The back contact of the

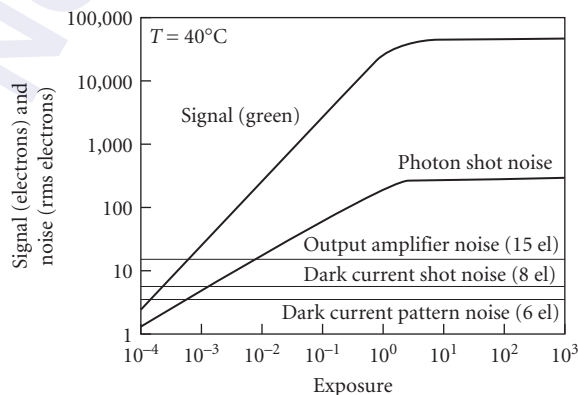


FIGURE 18 Quantum efficiency of interline CCD with and without microlens array.

photoconductor contacts the diode in the interline CCD. Light absorbed in the amorphous silicon generates electron-hole pairs. The amorphous silicon is biased such as to create a high field which sweeps out the electrons to the back contact, where they can be stored on the diode until they are transferred into the CCD. Because the photoconductor is fabricated on top of the CCD pixel, it can have nearly 100 percent fill factor. In addition, because of the wide bandgap of the amorphous silicon, its dark current (due to thermal generation of carriers) is often lower than the single-crystal silicon. However, the amorphous silicon photoconductors suffer difficulties related to charge trapping. Owing to impurities and dangling or strained bonds, there is a high density of traps in amorphous silicon. These can trap charge carriers and reemit them at a later time, causing image lag. Because the trapping and detrapping is field-dependent, it can also result in nonlinear response.

CCD Performance In all CCDs, both interline and frame transfer, the signal output is linear with the illumination except at levels approaching saturation. This linear response is in marked contrast to image tubes, which exhibit highly nonlinear response. For interline CCDs, the total charge capacity ranges from 100,000 electrons for larger cells (such as the $13.6(\text{V}) \times 11.6(\text{H})\text{-}\mu\text{m}$ cell typical for a 2/3-in format) to 20,000 electrons or less for smaller cells (such as the $6.8 \times 5.8\text{-}\mu\text{m}$ cell in a 1/3-in format 484×768 -pixel CCD).

The principal noise sources in both interline and full-frame image sensors are dark current pattern and shot noises and output amplifier noise. To illustrate the contributions, CCDs have dark current levels at 40°C of less than 1 nA/cm^2 . For a 2/3-in format sensor, this would correspond to about 320 electrons per pixel dark level. The corresponding shot noise would be 18-rms electrons and the pattern noise would typically be at a similar level. However, because of the nonrandom nature of pattern noise, its appearance is considerably more noticeable. The output amplifier noise would also be at about the 15-rms electron level. Thus, the overall noise for this example would be in the 30-electron range and the dynamic range would be over 2000 for a charge capacity of 90,000 electrons. For scientific applications where the sensor can be cooled and the readout performed at a lower frequency, noise levels less than 5 electrons can be achieved and even subelectron noise has been reported.¹⁸

Color Imaging

Silicon based CCDs are monochrome in nature. That is, they have no natural ability to determine the varying amounts of red, green, and blue (RGB) illumination presented to the photodetectors. There are three techniques to extract color information.

1. *Color Sequential* (Fig. 19)—A color image can be created using a CCD by taking three successive exposures while switching in optical filters having the desired RGB transmittances. This approach is normally used only to provide still images of stationary scenes. The resulting image is then reconstructed off-chip. The advantage to this technique is that resolution of each color can remain that of the CCD itself. The disadvantage is that three exposures are required, reducing frame times by more than a factor of three. Color misregistration can also occur due to subject or camera motion. The filter switching assembly also adds to the mechanical complexity of the system.
2. *Three-Chip Color* (Fig. 19)—Three-chip color systems use an optical system to split the scene into three separate color images. A dichroic prism beam splitter is normally used to provide RGB images. Color images can then be detected by synchronizing the outputs of the three CCDs. The disadvantage to such a system is that the optical complexity is very high and registration between sensors is difficult.
3. *Integral Color Filter Arrays (CFA)* (Fig. 19)—Instead of performing the color filtering off-chip, filters of the appropriate characteristics can be fabricated above individual photosites.^{31,32} This approach can be performed during device fabrication using dyed (e.g., cyan, magenta, yellow) photoresists in various patterns. The major problem with this approach is that each pixel is sensitive to only one color. Off-chip processing is required to “fill in” the missing color information between pixels.^{33,34}

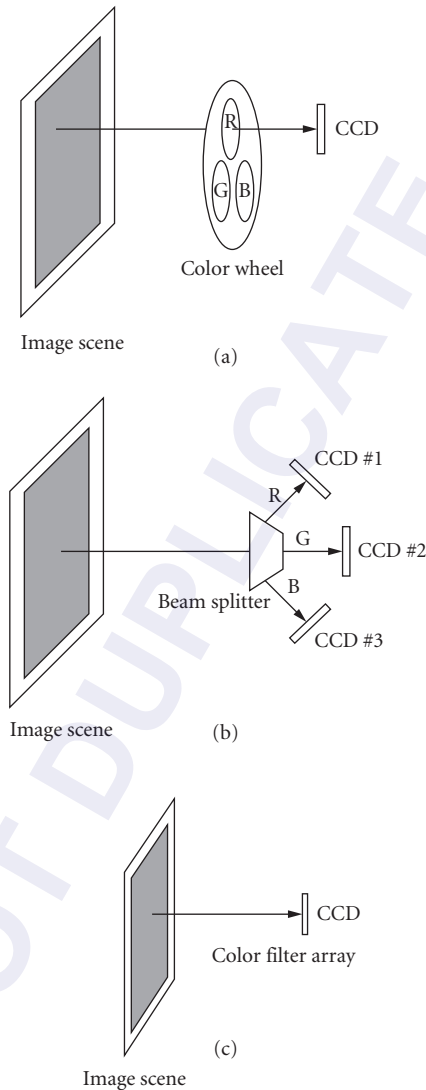


FIGURE 19 Methods of color separation in cameras with area image sensors: (a) color sequential using a color wheel; (b) a prism with dichroic beam splitters and three image sensors; and (c) single-chip image sensor with integral color filter array.

In order to minimize size, weight, and cost, most consumer color camcorders use a CCD sensor with an integral CFA. The photosites are covered with individual color filters—for example, a red, green, and blue striped filter, or a green, magenta, cyan, and yellow mosaic filter. Some popular CFA patterns are shown in Fig. 20. Because each photosite can sense only one color, the color sampling is not coincident. For example, a blue pixel might be seeing a white line, while nearby red and green pixels are seeing a dark line in the scene. As a result, high-frequency luminance edges can be aliased into bright color bands. These color bands depend not only on the color filter pattern used, but also

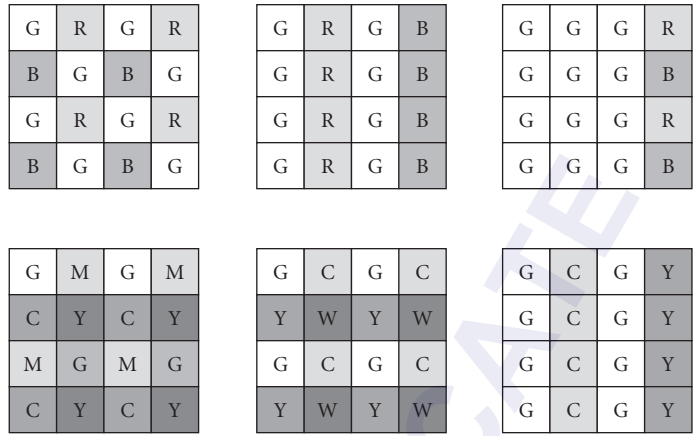


FIGURE 20 Common color filter array patterns where, R = Red, G = Green, B = Blue, Y = Yellow, M = Magenta, C = Cyan, and W = White.

on the optical prefilter and CFA interpolation algorithm. The color bands are caused by aliasing, which is a property of any sampled system. Aliasing occurs when the frequency of the input signal is greater than the Nyquist limit of one-half the sampling frequency. If the input frequency is well below the Nyquist limit, there are many samples per cycle. This allows the input to be reconstructed perfectly, if a proper reconstruction filter is used. When the input frequency is greater than the Nyquist limit, there are less than two samples per cycle. The sample values now define a new curve, which has a frequency lower than the input frequency. In effect, the high frequency takes on the alias of a lower frequency. Aliasing is a particular problem with color sensors, since the sampling phase is different for the different color photosites. Therefore, the aliased signal has different phases for different colors. This creates the color bands.

The color aliasing is reduced by using an optical anti-aliasing or “blur” filter, positioned in front of the color CCD sensor.³⁵ Blur filters are typically made of birefringent quartz, with the crystal axis oriented at a 45° angle, as shown in Fig. 21. In this orientation, the birefringent quartz exhibits the double refraction effect. An unpolarized input ray emerges as two polarized output rays, labelled *o*- and *e*-rays. The output ray separation is proportional to the filter’s thickness, *T*. A 1.5-mm-thick plate will give a separation of about 9 μm. Figure 21 shows a simple “two-spot” filter. More complex filters use three or more pieces of quartz cemented in a stack.

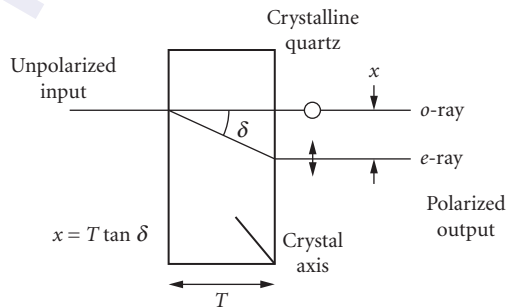


FIGURE 21 Birefringent blur filter used to reduce aliasing in single-chip color image sensors.

32.6 REFERENCES

1. G. Lubberts and B. C. Burkey, "Optical and Electrical Properties of Heavily Phosphorous-Doped Epitaxial Silicon Layers," *J. Appl. Phys.*, **55**(3):760–763 (Feb. 1, 1984).
2. O. S. Heavens, *Optical Properties of Thin Solid Films*, Dover, New York, 1965.
3. N. Teranishi et al., "An Interline CCD Image Sensor with Reduced Image Lag," *IEEE Trans. Electron Devices* **ED-31**(12):1829–1833 (Dec. 1984).
4. S. Kosman et al., "A Large Area 1.3 Megapixel Full-Frame CCD Image Sensor with a Lateral Overflow Drain and a Transparent Gate Electrode," *Proc. IEEE International Electron Device Meeting*, Washington, D.C. (Dec. 1990).
5. B. Burkey et al., "The Pinned Photodiode for an Interline Transfer CCD Image Sensor," *Proc. 1984 International Electron Device Meeting*, Washington, D.C. (Dec. 1986), p. 28.
6. M. Sasaki, H. Ihara, and Y. Matsunaga, "A 2/3-in 400k-Pixel Sticking-Free Stact CCD Image Sensor," *IEEE Trans. Electron Devices*, **ED-28**(11) (Nov. 1993).
7. H. Yamashita et al., "A 2/3-Inch 2-Megapixel Stack CCD Imager," *Proceedings of the 1994 International Solid State Circuits Conference*, p. 224.
8. E. G. Stevens et al., "A 1-Megapixel, Progressive-Scan Image Sensor with Antiblooming Control and Lag-Free Operation," *IEEE Trans. Electron Devices*, **ED-38**(5) (May 1991).
9. D. K. Schroder, "The Concept of Generation and Recombination Lifetimes in Semiconductors," *IEEE Trans. Electron Devices* **ED-29**(8) (Aug. 1982), p. 17.
10. J. Van der Spiegel and G. J. Declerck, *Solid State Electronics* **27**:147 (1984).
11. Carlo, H. Sequin and Michael F. Tompsett, *Charge Transfer Devices Advances in Electronics and Electron Physics*, suppl. 8, Academic Press, Inc., New York, 1975.
12. M. J. Howes and D. V. Morgan, (eds.), *Charge Coupled Devices and Systems*, Wiley, Chicago, 1979.
13. D. F. Barbe et al., *Charge-Coupled Devices*, Springer-Verlag, New York, 1980.
14. M. Azuma et al., "A Fixed Pattern Noise Free 2/3", 1.3 M-Pixel CCD Image Sensor for HDTV Camera System," *Proc. 1991 IEEE International Solid-State Circuits Conference*, San Francisco, 212–213 (Feb. 1991).
15. Y. Matsunaga and S. Ohsawa, "1/3 Inch Interline Transfer CCD Image Sensor with Negative Feedback 94 dB Dynamic Range Charge Detector," *Proc. 1991 IEEE International Solid-State Circuits Conference*, San Francisco, 210–211 (Feb. 1991).
16. J. Janesick, "Open Pinned-Phase CCD Technology," *Proc. SPIE*, vol. 1159, San Diego (Nov. 1989).
17. J. Hyncek, "Virtual Phase Technology: A New Approach to Fabrication of Large-Area CCDs," *IEEE Trans. Electron Devices* **ED-28**(5):483–489 (May 1981).
18. J. Janesick et al., "Fano-Noise-Limited CCDs," *Optical and Optoelectric Applied Science and Engineering Symposium: X-Ray Instrumentation in Astronomy*, San Diego (Aug. 14–19, 1988).
19. J. Carbone et al., "New Low Noise Random Access, Radiation Resistant and Large Format Charge Injection Device (CID) Imagers," *Proc. of SPIE Conference 1900* (Jan. 31–Feb. 4, 1993).
20. M. Sugawara et al., "An Amplified MOS Imager Suited for Image Processing," *Proc. 1994 IEEE International Solid-State Circuits Conference*, San Francisco, 228–229 (Feb. 1994).
21. M. Ogata et al., "A Small Pixel CMD Image Sensor," *IEEE Trans. Electron Devices*, **ED-38**(5):1005–1010 (May 1991).
22. D. H. McCann, M. H. White, A. P. Turley, and R. A. Frosch, "Time Delay and Integration Detectors Using Charge Transfer Devices," U.S. Patent Number 4,280,141, July 1981.
23. Thompson et al., "Time-Delay-and-Integration Charge-Coupled Devices Using Tin Oxide Gate Technology," *IEEE Trans. Electron Devices* **ED-25**(2):132–134 (Feb. 1978).
24. S. Ohba et al., "MOS Area Sensor: Part 2: Low-Noise MOS Area Sensor with Antiblooming Photodiodes," *IEEE Trans. Electron Devices* **ED-27**(8):1682–1687 (Aug. 1980).
25. T. Nakamura et al., "A New MOS Image Sensor Operating in a Non-Destructive Readout Mode," *Proc. 1986 IEEE International Electron Device Meeting*, Washington, D.C., 353–356 (Dec. 1986).
26. Nomoto et al., "A 2/3 Inch 2M-Pixel CMD Image Sensor with Multi-Scanning Functions," *Proc. 1993 IEEE International Solid-State Circuits Conference*, San Francisco, 196–197 (Feb. 1993).

27. D. H. McCann et al., "Buried Channel CCD Imaging Arrays with Tin Oxide Transparent Gates," *IEEE International Solid State Circuits Conference*, pp. 30, 31, 261, 262 (Feb. 1978).
28. K. Harada, "A 2/3-Inch, 2M-Pixel Frame Interline Transfer CCD HDTV Image Sensor," *Proc. 1992 IEEE International Solid-State Circuits Conference*, San Francisco, 170–171 (Feb. 1992).
29. E. Stevens, "Photoresponse Nonlinearity of Solid State Image Sensors with Antiblooming Protection," *IEEE Trans. Electron Devices* **ED-38**(2):229–302 (Feb. 1991).
30. A. Weiss et al., *J. Electrochem. Soc.* 133:110C (1986).
31. P. Dillion et al., "Fabrication and Performance of Color Filter Arrays for Solid-State Imagers," *IEEE Trans. Electron Devices* **ED-25**:97–101 (1978).
32. P. Dillion, D. Lewis, and F. Kaspar, "Color Imaging System Using a Single CCD Array," *IEEE Trans. Electron Devices*, **ED-25**:102–107 (Feb. 1978).
33. K. A. Parulski, et al., "A Digital Color CCD Imaging System Using Custom VLSI Circuits," *IEEE Trans. Consumer Electronics* **35**(3):382–389 (Aug. 1989).
34. K. A. Parulski, "Color Filters and Processing Alternatives for One-Chip Cameras," *IEEE Trans. Electron Devices* **ED-32**(8):1381–1389 (Aug. 1985).
35. J. E. Greivenkamp, "Color Dependent Optical Prefilter for the Suppression of Aliasing Artifacts," *Appl. Opt.* **29**(5):676–684, (Feb. 1990).

INFRARED DETECTOR ARRAYS

Lester J. Kozłowski

*Altasens, Inc.
Westlake Village, California*

Walter F. Kosonocky*

*New Jersey Institute of Technology
University Heights
Newark, New Jersey*

33.1 GLOSSARY

A_{det}	detector area
A_I	gate modulation current gain (ratio of integration capacitor current to load current)
A_V	amplifier voltage gain
C_{amp}	amplifier capacitance
C_{det}	detector capacitance
C_{FIS}	fill-and-spill gate capacitance for a Tompsett type CCD input
C_{fb}	CTIA feedback or Miller capacitance
C_{gd}	FET gate-drain overlap capacitance
C_{gs}	FET gate-source capacitance
C_L	CTIA band-limiting load capacitance
C_{out}	sense node capacitance at the CCD output
C_T	effective feedback (transcapacitance) or integration capacitance for a capacitive transimpedance amplifier
$C_{T\lambda}$	spectral photon contrast
cte	charge transfer efficiency
$D_{\lambda\text{pk}}^*$	peak detectivity ($\text{cm}\cdot\text{Hz}^{1/2}/\text{W}$ or Jones)
D_{bb}^*	blackbody detectivity ($\text{cm}\cdot\text{Hz}^{1/2}/\text{W}$ or Jones)
D_{th}^*	thermal detectivity ($\text{cm}\cdot\text{Hz}^{1/2}/\text{W}$ or Jones)
e^-	electron
E_g	detector energy gap
$f/\#$	conventional shorthand for the ratio of the focal length of a lens to its diameter
f_{chop}	chopper frequency

*Deceased.

f_{frame}	display frame rate
f_{knee}	frequency at which the $1/f$ noise intersects the broadband noise
f_s	spatial frequency (cycles/radian)
$g_{m, \text{LOAD}}$	gate transconductance of the load FET in the gate-modulated input circuit
g_m	gate transconductance of a Field Effect Transistor
h	Planck's constant
I_D	FET drain current
I_{det}	detector current
I_{photo}	detector photocurrent
k	Boltzmann constant
K_{amp}	amplifier FET noise spectral density at 1 Hz
K_{det}	detector noise spectral density at 1 Hz
K_{FET}	FET noise spectral density at 1 Hz
L	length-to-width ratio of a bar chart (always set to 7)
MRT	minimum resolvable temperature (K)
MTF	modulation transfer function for the optics, detector, readout, the integration process, or the composite sensor
n	detector junction ideality or diffusion constant
$N_{\text{amp}, 1/f}$	number of noise carriers for one integration time due to amplifier FET $1/f$ noise
$N_{\text{amp}, \text{white}}$	number of noise carriers for one integration time due to amplifier FET white noise
N_c	number of photo-generated carriers integrated for one integration time
n_{det}	detector junction ideality or diffusion constant
NE ΔT	noise equivalent temperature difference (K)
n_{FET}	subthreshold FET ideality
N_{FPA}	composite (total) FPA noise in carriers
$N_{\text{KTC, channel}}$	CTIA broadband channel noise in carriers
$N_{\text{load, white}}$	number of noise carriers for one integration time due to CTIA load FET white noise
N_{os}	display overscan ratio
N_{PHOTON}	shot noise of photon background in carriers
NSD	noise spectral density of a detector or field effect transistor; the $1/f$ noise is often specified by the NSD at a frequency of 1 Hz
N_{sf}	source follower noise
N_{ss}	serial scan ratio
q	electron charge in coulombs
Q_B	photon flux density (photons/cm ² -s) incident on a focal plane array
Q_D	charge detected in a focal plane array for one integration time
Q_{max}	maximum charge signal at saturation
R_{det}	detector resistance
R_{LOAD}	gate modulation load resistance
R_0	detector resistance at zero-bias resistance
R_0A	detector resistance-area product at zero-bias voltage
R_r	detector resistance in reverse-bias resistance
S/N	signal-to-noise ratio
SNR _T	target signal-to-noise ratio
S_V	readout conversion factor describing the ratio of output voltage to detected signal carriers
T	operating temperature

tce	thermal coefficient of expansion
TCR	thermal coefficient of resistance for bolometer detectors
T_D	time constant for correlated double sampling process normally set by Nyquist rate
t_{int}	integration time
U	residual nonuniformity
V_{br}	detector reverse-bias breakdown voltage, sometimes defined as the voltage where $R_r = R_0$
V_D	FET drain voltage
V_{det}	detector bias voltage
V_{DS}	FET drain-to-source voltage
V_G	FET gate voltage
v_n	measured rms noise voltage
ΔA_I	gate modulation current gain nonuniformity
Δf	noise bandwidth (Hz)
ΔI_{photo}	differential photocurrent
ΔT	scene temperature difference creating differential photocurrent ΔI_{photo}
ΔV_S	signal voltage for differential photocurrent ΔI_{photo}
Δx	horizontal detector subtense (mradian)
Δy	vertical detector subtense (mradian)
η	detector quantum efficiency
η_{BLIP}	percentage of BLIP
$\eta_{\text{inj, DI}}$	injection efficiency of detector current into the source-modulated FET of the direct injection input circuit
η_{inj}	injection efficiency of detector current
η_{noise}	injection efficiency of DI circuit noise into integration capacitor
η_{pc}	quantum efficiency of photoconductive detector
η_{pv}	quantum efficiency of photovoltaic detector
λ_c	detector cutoff wavelength (50 percent of peak response, μm)
σ_{det}	noise spectral density of total detector noise including photon noise
$\sigma_{\text{input, ir}}$	noise spectral density of input-referred input circuit noise
σ_{LOAD}	noise spectral density of input-referred load noise
$\sigma_{\text{mux, ir}}$	noise spectral density of input-referred multiplexer noise
σ_{VT}	rms threshold voltage nonuniformity across an FPA
τ_{amp}	amplifier time constant (s)
τ_{eye}	eye integration time (s)
τ_o	optical transmission
ω	angular frequency (radians)
$\langle e_{\text{amp}} \rangle$	buffer amplifier noise for buffered direct injection circuit

33.2 INTRODUCTION

Infrared sensors have been available since the 1940s to detect, measure, and image the thermal radiation emitted by all objects. Due to advanced detector materials and microelectronics, large scanning and staring focal plane arrays (FPA) with few defects are now readily available in the short wavelength infrared (SWIR; 1 to 3 μm), medium wavelength infrared (MWIR; ≈ 3 to 5 μm), and long wavelength infrared (LWIR; ≈ 8 to 14 μm) spectral bands. We discuss in this chapter the disparate FPA technologies, including photon and thermal detectors, with emphasis on the emerging types.

IR sensor development has been driven largely by the military. Detector requirements for missile seekers and forward looking infrared (FLIR) sensors led to high-volume production of photoconductive (PC) HgCdTe arrays starting in the 1970s. Though each detector requires direct connection to external electronics for purposes of biasing, signal-to-noise ratio (SNR) enhancement via time delay integration (TDI), and signal output, the first-generation FPAs displaced the incumbent Pb-salt (PbS, PbSe) and Hg-doped germanium devices, and are currently being refined using custom analog signal processing,¹ laser-trimmed solid-state preamplifiers, etc.

Size and performance limitations of first-generation FLIRs necessitated development of self-multiplexed FPAs with on-chip signal processing. Second-generation thermal imaging systems use high-density FPAs with relatively few external connections. Having many detectors that integrate longer, low-noise multiplexing and on-chip TDI (in some scanning arrays), second-generation FPAs offer higher performance and design flexibility. Video artifacts are suppressed due to the departure from ac-coupling and interlaced raster scan, and external connections are minimized. Fabricated in monolithic and hybrid methodologies, many detector and readout types are used in two basic architectures (staring and scanning). In a monolithic FPA, the detector array and the multiplexing signal processor are integrated in a single substrate. The constituents are fabricated on separate substrates and interconnected in a hybrid FPA.

FPAs use either photon or thermal detectors. Photon detection is accomplished using intrinsic or extrinsic semiconductors and either photovoltaic (PV), photoconductive (PC), or metal insulator semiconductor (MIS) technologies. Thermal detection relies on capacitive (ferro- and pyroelectric) or resistive bolometers. In all cases, the detector signal is coupled into a multiplexer and read out in a video format.

Infrared Applications

Infrared FPAs are now being applied to a rapidly growing number of civilian, military, and scientific applications such as industrial robotics and thermography (e.g., electrical and mechanical fault detection), medical diagnosis, environmental and chemical process monitoring, Fourier transform IR spectroscopy and spectroradiometry, forensic drug analysis, microscopy, and astronomy. The combination of high sensitivity and passive operation is also leading to many commercial uses. The passive monitoring provided by the addition of infrared detection to gas chromatography-mass spectroscopy (GC-MS), for example, yields positive chemical compound and isomer detection without sample alteration. Fusing IR data with standard GC-MS aids in the rapid discrimination of the closely related compounds stemming from drug synthesis. Near-IR (0.7 to 0.1 μm) and SWIR spectroscopy and fluorescence are very interesting near-term commercial applications since they pave the way for high-performance FPAs in the photochemical, pharmaceutical, pulp and paper, biomedical, reference quantum counter, and materials research fields. Sensitive atomic and molecular spectroscopies (luminescence, absorption, emission, and Raman) require FPAs having high quantum efficiency, low dark current, linear transimpedance, and low read noise.

Spectral Bands

The primary spectral bands for infrared imaging are 3 to 5 and 8 to 12 μm because atmospheric transmission is highest in these bands. These two bands, however, differ dramatically with respect to contrast, background signal, scene characteristics, atmospheric transmission under diverse weather conditions, and optical aperture constraints. System performance is a complex combination of these and the ideal system requires dual band operation. Factors favoring the MWIR include its higher contrast, superior clear-weather performance, higher transmissivity in high humidity, and higher resolution due to $\sim 3 \times$ smaller optical diffraction. Factors favoring the LWIR include much-reduced background clutter (solar glint and high-temperature countermeasures including fires and flares have much-reduced emission), better performance in fog, dust, and winter haze, and higher immunity to atmospheric turbulence. A final factor favoring the LWIR, higher

S/N ratio due to the greater radiance levels, is currently moot because of technology limitations. Due to space constraints and the breadth of sensor applicability, we focus on target/background metrics in this section.

The signal collected by a visible detector has higher daytime contrast than either IR band because it is mainly radiation from high-temperature sources that is subsequently reflected off earth-based (ambient temperature; ≈ 290 K) objects. The high-temperature sources are both solar (including the sun, moon, and stars) and synthetic. Since the photon flux from high- and low-temperature sources differs greatly at visible wavelengths from day to night, scene contrasts of up to 100 percent ensue.

Reflected solar radiation has less influence as the wavelength increases to a few microns since the background radiation increases rapidly and the contrast decreases. In the SWIR band, for example, the photon flux density from the earth is comparable to visible room light (10^{13} photons/cm²-s). The MWIR band ($\sim 10^{15}$ photon flux density) has lower, yet still dynamic, daytime contrast, and can still be photon-starved in cold weather or at night.

The net contribution from reflected solar radiation is even lower at longer wavelengths. In the LWIR band, the background flux is equivalent to bright sunlight ($\approx 10^{17}$ photons/cm²-s). This band thus has even lower contrast and much less background clutter, but the “scene” and target/background metrics are similar day and night. Clear-weather performance is relatively constant.

Depending on environmental conditions, however, IR sensors operating in either band must discern direct emission from objects having temperatures very near the average background temperature (290 K) in the presence of the large background and degraded atmospheric transmissivity. Under conditions of uniform thermal soak, such as at diurnal equilibrium, the target signal stems from minute emissivity differences.

The spectral photon incidence for a full hemispheric surround is

$$Q = \tau_{cf} \int_{\lambda_1}^{\lambda_2} Q_{\lambda}(\lambda) d\lambda \quad (1)$$

if a zero-emissivity bandpass filter having in-band transmission τ_{cf} , cut-on wavelength λ_1 , and cutoff wavelength λ_2 is used (zero emissivity obtained practically by cooling the spectral filter to a temperature where its self-radiation is negligible). The photon flux density, Q_B (photons/cm²-s), incident on a focal plane array is

$$Q_B = \frac{1}{4(f/\#)^2 + 1} Q \quad (2)$$

where $f/\#$ is the conventional shorthand for the ratio of the focal length to the diameter (assumed circular) of the limiting aperture or lens. The cold shield $f/\#$ limits the background radiation to a field-of-view consistent with the warm optics to eliminate extraneous background flux and concomitant noise. The background flux in the LWIR band is approximately two orders of magnitude higher than in the MWIR.

The spectral photon contrast, $C_{T\lambda}$, is the ratio of the derivative of spectral photon incidence to the spectral photon incidence, has units K^{-1} , and is defined

$$C_{T\lambda} = \left(\frac{\partial Q}{\partial T} \right) / \left(\frac{Q}{Q} \right) \quad (3)$$

Figure 1 is a plot of $C_{T\lambda}$ for several MWIR subbands (including 3.5 to 5, 3.5 to 4.1, and 4.5 to 5 μm) and the 8.0 to 12 μm LWIR spectral band. The contrast in the MWIR bands at 300 K is 3.5 to 4 percent compared to 1.6 percent for the LWIR band. While daytime MWIR contrast is even higher due to reflected sunlight, an LWIR FPA offers higher sensitivity if it has the larger capacity needed for storing the larger amounts of photogenerated (due to the higher background flux) and detector-generated carriers (due to the narrow bandgap). The photon contrast and the background flux are key parameters that determine thermal resolution as will be described later under “Performance Figures of Merit.”

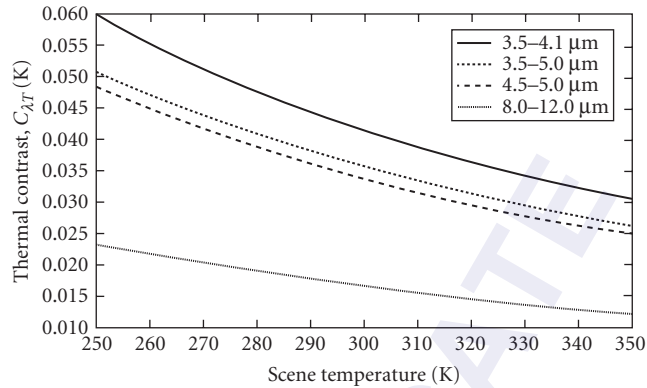


FIGURE 1 Spectral photon contrast in the MWIR and LWIR.

Scanning and Staring Arrays

The two basic types of FPA are scanning and staring. The simplest scanning device consists of a linear array as shown in Fig. 2a. An image is generated by scanning the scene across the strip. Since each detector scans the complete horizontal field-of-view (one video raster line) at standard video frame rates, each resolution element or pixel has a short integration time and the total detected charge can usually be accommodated.

A staring array (Fig. 2b) is the two-dimensional extension of a scanning array. It is self-scanned electronically, can provide enhanced sensitivity, and is suitable for lightweight cameras. Each pixel

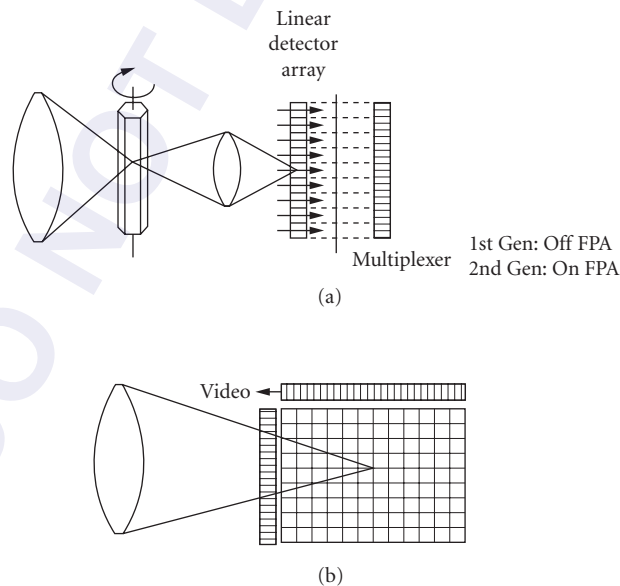


FIGURE 2 Scanning (a) and staring (b) focal plane arrays.

is a dedicated resolution element, but synchronized dithering of sparsely populated arrays is sometimes used to enhance the effective resolution, minimize spatial aliasing, and increase the effective number of pixels. Although theoretically charge can be integrated for the full frame time, the charge-handling capacity is inadequate at terrestrial LWIR backgrounds.

Detectors

Infrared detectors convert IR photons and energy to electrical signals. Many types are used in FPAs (as shown in Fig. 3²) including photon and thermal detectors that address diverse requirements spanning operating temperatures from 4 K to room temperature. Figure 4 compares the quantum efficiencies of several detector materials.

Intrinsic detectors³ usually operate at higher temperatures than extrinsic devices, have higher quantum efficiencies, and dissipate less power. Backside-illuminated devices, consisting of an absorbing epitaxial layer on a transparent substrate, are used in hybrid FPAs and offer the advantages of nearly 100 percent active detector area, good mechanical support, and high quantum efficiency. The most popular intrinsic photovoltaics are HgCdTe and InSb. These detectors are characterized by their quantum efficiency (η), zero-bias resistance (R_0), reverse-bias resistance (R_r), junction ideality or diffusion constant n , excess noise (if any) versus bias, and reverse-bias breakdown voltage (V_{br}), which is sometimes defined as the voltage where $R_r = R_0$.

PtSi is a photovoltaic Schottky barrier detector (SBD) which is the most mature for large monolithic FPAs. IR detection is via internal photoemission over a Schottky barrier (0.21 to 0.23 eV). Characteristics include low (≈ 0.5 percent for broadband 3.5 to 5.0 μm) but very uniform quantum efficiency, high producibility that is limited only by the Si readout circuits, full compatibility with VLSI technology, and soft spectral response with peak below 2 μm and zero response just beyond 6 μm . Internal photoemission dark current requires cooling below 77 K.

HgCdTe is the most popular intrinsic photoconductor, and various linear arrays in several scanning formats are used worldwide in first-generation FLIRs. For reasons of producibility and cost, HgCdTe photoconductors have historically enjoyed a greater utilization than PV detectors despite the latter's higher quantum efficiency, higher D^* by a factor of $(2\eta_{pv}/\eta_{pc})^{1/2}$, and superior modulation transfer function (MTF). Nevertheless, not all photoconductors are good candidates for FPAs due to their low detector impedance. This includes the intrinsic materials InSb and HgCdTe.

The most popular photoconductive material system for area arrays is doped extrinsic silicon (Si: x ; where x is In, As, Ga, Sb, etc.), which is made in either conventional or impurity band conduction [IBC or blocked impurity band (BIB)] technologies. Early monolithic arrays were doped-Si devices, due primarily to compatibility with the silicon readout. Extrinsic photoconductors must be made relatively thick (up to 30 mils; doping density of IBCs, however, minimizes this thickness requirement but does not eliminate it) because they have much lower photon capture cross section than intrinsic detectors. This factor adversely affects their MTF in systems having fast optics.

Historically, Si:Ga and Si:In were the first mosaic focal plane array PC detector materials because early monolithic approaches were compatible with these dopants. Nevertheless, problems in fabricating the detector contacts, early breakdown between the epitaxial layer and the detector material (double injection), and the need for elevated operating temperatures helped force the general move to monolithic PtSi and intrinsic hybrids.

The most advanced extrinsic photoconductors are IBC detectors using Si:As and Si:Ga.⁴ These have reduced recombination noise (negating the $\sqrt{2}$ superiority in S/N that PV devices normally have) and longer spectral response than standard extrinsic devices due to the higher dopant levels. IBC detectors have a unique combination of PC and PV characteristics, including extremely high impedance, PV-like noise (reduced recombination noise since IBC detectors collect carriers both from the continuum and the "hopping" impurity band), linear photoconductive gain, high uniformity, and superb stability. The photo-sensitive layer in IBCs is heavily doped to achieve hopping-type conduction. A thin, lightly doped ($10^{10}/\text{cm}^2$) silicon layer blocks the hopping current before

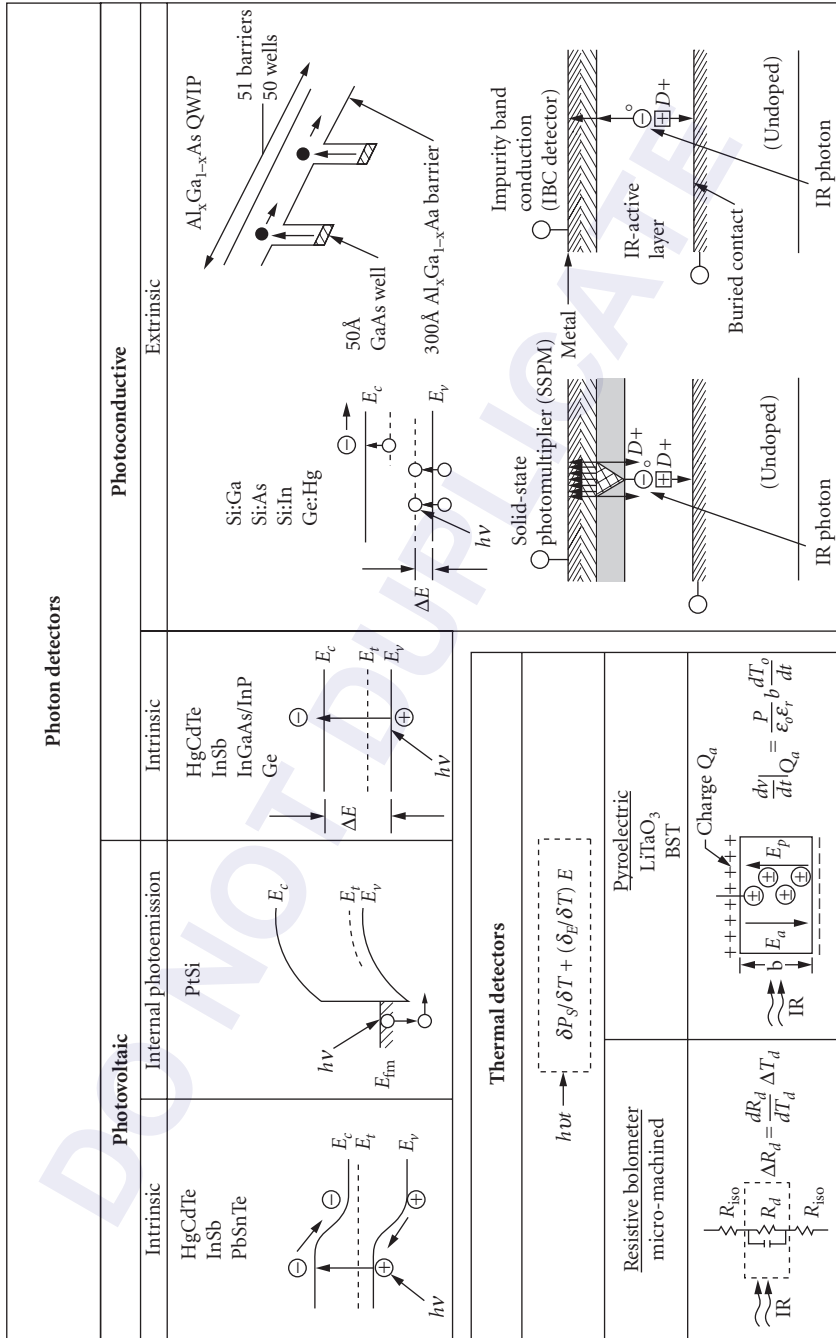


FIGURE 3 Photon and thermal detectors.

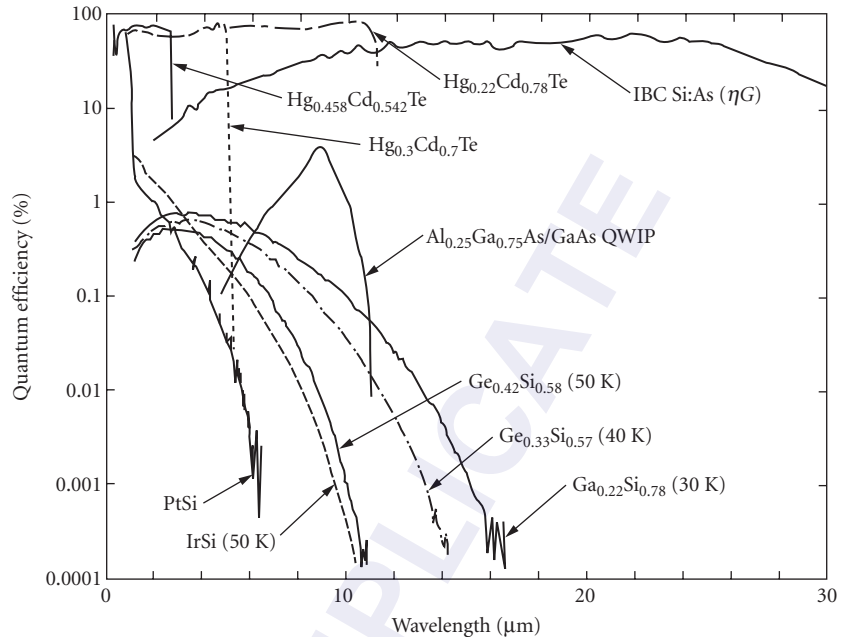


FIGURE 4 Quantum efficiency versus wavelength for several detector materials.

it reaches the device electrode to reduce noise. Specially doped IBCs (see cross-sectional views in Fig. 3) operate as solid-state photomultipliers (SSPM) and visible light photon counters (VLPC) in which photoexcited carriers are amplified by impact ionization of impurity-bound carriers.⁵ The amplification allows counting of individual photons at low flux levels. Standard SSPMs respond from 0.4 to 28 μm .

An alternative custom tunable detector is the GaAs/AlGaAs quantum well infrared photodetector (QWIP). Various QWIP photoconductive⁶ and photovoltaic⁷ structures are being investigated as low-cost alternatives to II-VI LWIR detectors like HgCdTe. Infrared detection in the typical PC QWIP is via intersubband or bound-to-extended-state transitions within the multiple quantum well superlattice structure. Due to the polarization selection rules for transitions between the first and second quantum wells, the photon electric field must have a component parallel to the superlattice direction. Light absorption in *n*-type material is thus anisotropic with zero absorption at normal incidence. The QWIP detector's spectral response is narrowband, peaked about the absorption energy. The wavelength of peak response can be adjusted via quantum well parameters and can be made bias-dependent.

Various bolometers, both resistive and capacitive (pyroelectric), are also available. Bolometers sense incident radiation via energy absorption and concomitant change in device temperature in both cooled moderate-performance and uncooled lower-performance schemes. Much recent research, which was previously highly classified, has focused on both hybrid and monolithic uncooled arrays and has yielded significant improvements in the detectivity of both resistive and capacitive bolometer arrays. The resistive bolometers currently in development consist of a thin film of a temperature-sensitive resistive material film which is suspended above a silicon readout. The pixel support struts provide electrical interconnect and high thermal resistance to maximize pixel sensitivity. Recent work has focused on the micromachining necessary to fabricate mosaics with low thermal conductance using monolithic methodologies compatible with silicon.

Capacitive bolometers sense a change in elemental capacitance and require mechanical chopping to detect incident radiation. The most common are pyroelectric detectors. J. Cooper⁸ suggested the use of pyroelectric detectors in 1962 as a possible solution for applications needing a low-cost IR FPA with acceptable performance. These devices have temperature-dependent spontaneous polarization. Ferroelectric detectors are pyroelectric detectors having reversible polarization. There are over a thousand pyroelectric crystals, including several popularly used in hybrid FPAs; e.g. lithium tantalate (LiTaO_3), triglycine sulfate (TGS), and barium strontium titanate⁹ (BaSrTiO_3).

33.3 MONOLITHIC FPAs

A monolithic FPA consists of a detector array and the readout multiplexer integrated on the same substrate. The progress in the development of the monolithic FPAs in the last two decades has been strongly influenced by the rapid advances in the silicon VLSI technology. Therefore, the present monolithic FPAs can be divided into three categories reflecting their relationship to the silicon VLSI technology. The first category includes the “complete” monolithic FPAs in which the detector array and the readout multiplexer are integrated on the same silicon substrate using processing steps compatible with the silicon VLSI technology. They include the extrinsic Si FPAs reported initially in the 1970s,¹⁰ FPAs with Schottky barrier,¹¹ heterojunction detector FPAs, and microbolometer FPAs.¹²

The second category will be referred to here as the “partial monolithic” FPAs. This group includes narrowband detector arrays of HgCdTe ¹³ and InSb ¹⁴ integrated on the same substrate only with the first level of multiplexing, such as the row and column readout from a two-dimensional detector array. In this case the multiplexing of the detected signal is completed by additional silicon IC chips usually packaged on the imager focal plane.

The third category represents “vertically integrated” photodiode (VIP) FPAs. These FPAs are functionally similar to hybrids in the sense that a silicon readout multiplexer is used with the narrow-bandgap HgCdTe detectors. However, while in the hybrid FPA the completed HgCdTe detector array is typically connected by pressure contacts via indium bumps to the silicon multiplex pads; in the case of the vertically integrated FPAs, HgCdTe chips are attached to a silicon multiplexer wafer and then the fabrication of the HgCdTe photodiodes is completed including the deposition and the definition of the metal connections to the silicon readout multiplexer.

In the following sections we will review the detector readout structures, and the main monolithic FPA technologies. It should also be noted that most of the detector readout techniques and the architectures for the monolithic FPAs were originally introduced for visible silicon imagers. This heritage is reflected in the terminology used in the section.

Architectures

The most common structures for the photon detector readout and architectures of monolithic FPAs are illustrated schematically in Figs. 5 and 6.

MIS Photogate FPAs: CCD, CID, and CIM Most of the present monolithic FPAs use either MIS photogates or photodiodes as the photon detectors. Figure 5 illustrates a direct integration of the detected charge in the potential well of a MIS (photogate) detector for a charge coupled device (CCD) readout in (a), a charge injection device (CID) readout in (b), and a charge-integration matrix (CIM) readout in (c). The unique characteristic of the CCD readout is the complete transfer of charge from the integration well without readout noise. Also in a CCD FPA the detected charge, Q_D , can be transferred via potential wells along the surface of the semiconductor that are induced by clock voltages but isolated from the electrical pickup until it is detected by a low-capacitance (low kTC noise) on-chip amplifier. However, because of the relatively large charge transfer losses ($\sim 10^{-3}$ per transfer) and a limited charge-handling capacity, the use of nonsilicon CCD readout has

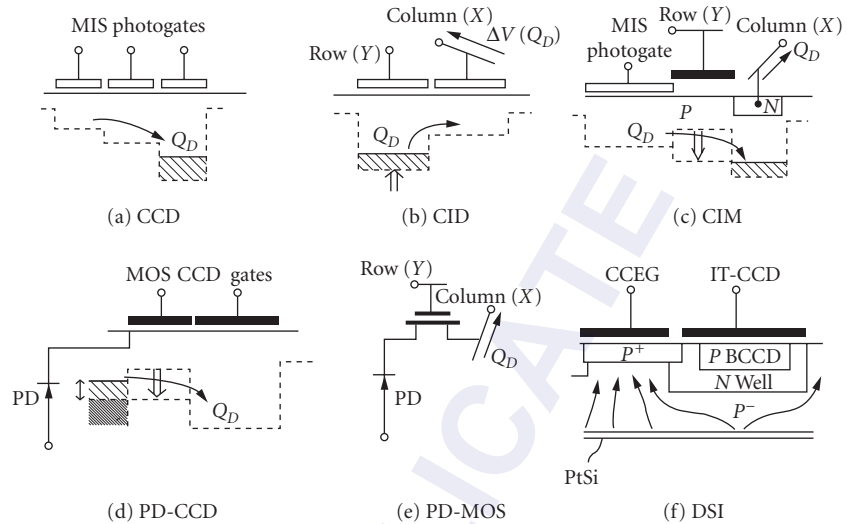


FIGURE 5 Photodetector readout structures.

been limited mainly to HgCdTe TDI FPAs. Such TDI imager is shown in Fig. 6a as a frame transfer (FT) type CCD area imager performing a function of a line sensor with the effective optical integration time increased by the number of TDI elements (CCD stages) in the column CCD registers. In this imager the transfer of the detected charge signal between CCD wells of the vertical register is adjusted to coincide with the mechanical motion of the image.

In the FT-CCD TDI FPA, the vertical registers perform the functions of charge detection and integration as well as transfer. The detected image is transferred one line at a time from the parallel vertical registers to the serial output registers. From there it is transferred at high clock rate to produce the output video. Similar TDI operation can also be produced by the interline-transfer (IT) CCD architecture shown in Fig. 6e. However, in the case of IT-CCD readout, the conversion of infrared radiation into charge signal photodetection is performed by photodiodes.

In the CID readout, see Fig. 5b, the detected charge signal is transferred back and forth between the potential wells of the MIS photogates for nondestructive X-Y addressable readout, $\Delta V(Q_D)$, that is available at a column (or a row) electrode due to the displacement current induced by the transfer of the detected charge signal, Q_D . At the end of the optical integration time, the detected charge is injected into the substrate by driving both MIS capacitors into accumulation.

CID FPAs with column readout for single-output-port and parallel-row readout are illustrated schematically in Fig. 6b and c, respectively. Another example of a parallel readout is the CIM FPA shown in Fig. 6d. The parallel readout of CID and CIM FPAs is used to overcome the inherent limitation on charge-handling capacity of these monolithic FPAs by allowing a short optical integration time with fast frame readout and off-chip charge integration by supporting silicon ICs.

Silicon FPAs: IT-CCD, CSD, and MOS FPAs The monolithic FPAs fabricated on silicon substrate take advantage of well-developed silicon VLSI process technology. Therefore, silicon ($E_g = 1.1$ eV), which is transparent to infrared radiation having wavelength longer than $1.0 \mu\text{m}$, is often used to produce monolithic CCD and MOS FPAs with infrared detectors that can be formed on silicon substrate. In the 1970s there was great interest in the development of monolithic silicon FPAs with extrinsic Si:In and Si:Ga photoconductors. However, since the early 1980s most progress was reported on monolithic FPAs with Schottky-barrier photodiodes, GeSi/Si heterojunction photodiodes, vertically integrated photodiodes, and resistive microbolometers. With the exception of the

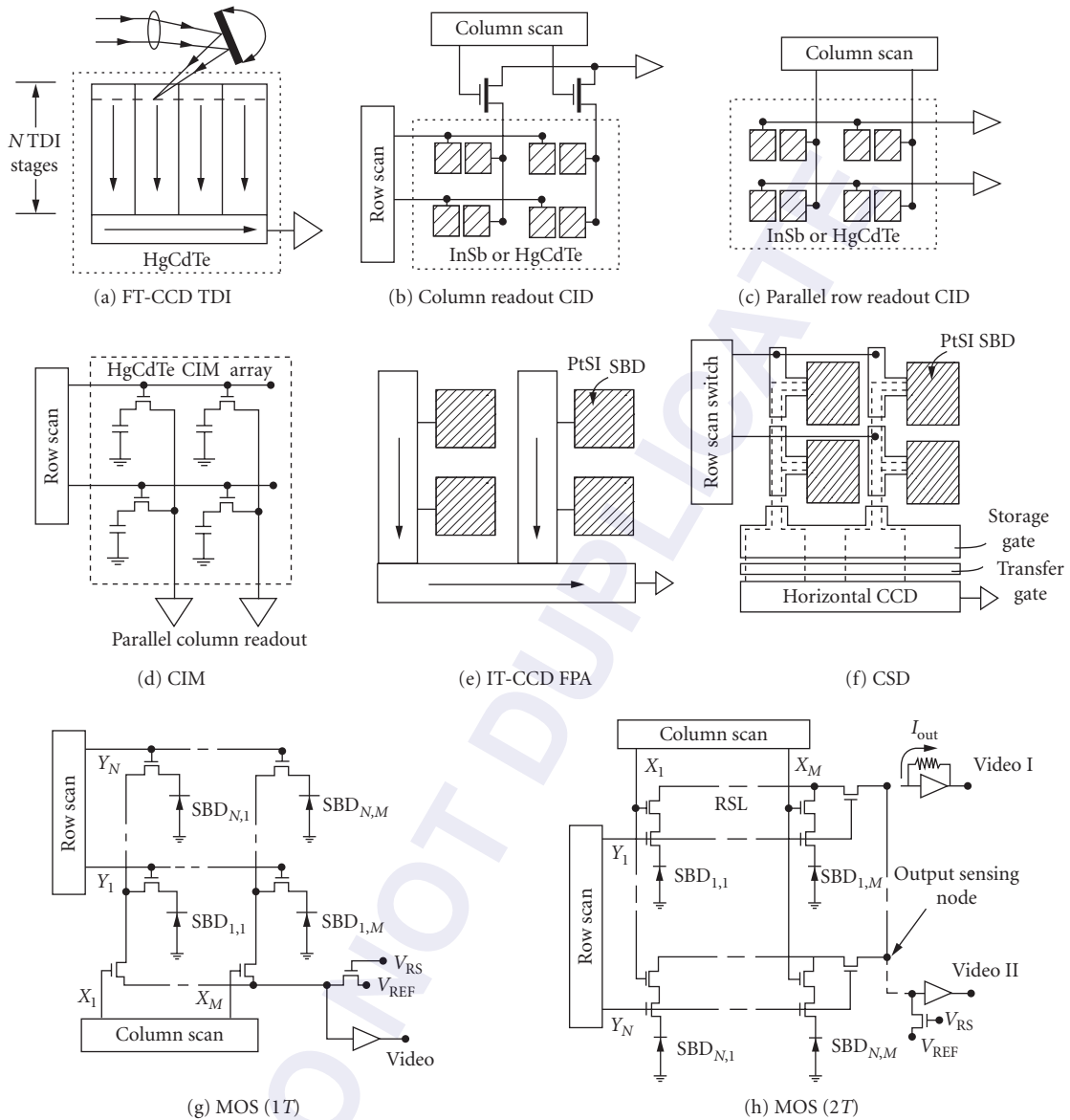


FIGURE 6 Monolithic FPA architectures.

resistive microbolometers, in the form of thin-film semiconductor photoresistors formed on micro-machined silicon structures, all of the above infrared detectors can be considered to be photodiodes and can be read out either by IT-CCD or MOS monolithic multiplexer.

The photodiode (PD) CCD readout, see Fig. 5d, is normally organized as the interline-transfer (IT) CCD staring FPA, shown in Fig. 6e. The IT-CCD readout has been used mostly for PtSi Schottky-barrier detectors (SBDs). The operation of this FPA consists of direct integration of the detected charge signal on the capacitance of the photodiode. At the end of the optical integration time, a frame readout is initiated by a parallel transfer of the detected charge from the photodiodes to the

parallel vertical CCD registers. From there the detected image moves by parallel transfers one line at a time into the horizontal output register for a high-clock-rate serial readout.

The design of the Schottky-barrier IT-CCD FPA involves a trade-off between fill factor (representing the ratio of the active detector area to the pixel area) and maximum saturation charge signal (Q_{mix}). This trade-off can be improved by the charge sweep device (CSD) architecture, shown in Fig. 6f, that has also been used as a monolithic readout multiplexer for PtSi and IrSi SBDs.

In CSD FPA the maximum charge signal is limited only by the SBD capacitance since its operation is based on transferring the detected charge signal from one horizontal line corresponding to one or two rows of SBDs (depending on the type of the interlacing used) into minimum geometry vertical CCD registers. During the serial readout of the previous horizontal line, the charge signal is swept into a potential well under the storage gate by low-voltage parallel clocking of the vertical registers. Then, during the horizontal blanking time, the line charge signal is transferred in parallel to the horizontal CCD register for serial readout during the next horizontal line time.

The main advantage of the silicon CCD multiplexer is relatively low readout noise, from a few electrons to the order of several tens of electrons (depending on video rate, sense capacitance, and CCD technology), so that a shot-noise-limited operation can be achieved at relatively low signal levels. But as the operating temperature is lowered below 60 K, the charge transfer losses of buried-channel CCDs (BCCDs) become excessive due to the freeze-out of the BCCD implant. Therefore, detectors requiring operation at 40 K or lower are more compatible with MOS readout device technology.

Photodiode (PD) MOS readout (see Fig. 5e) represents another approach to construction of an X-Y addressable silicon multiplexer. These types of monolithic MOS multiplexers used for readout of PtSi SBDs are illustrated in Fig. 6g and h.

A single-output-port FPA with one MOSFET switch per detector, MOS (1T), is shown in Fig. 6g. During FPA readout, the vertical scan switch transfers the detected charge signal from one row of detectors to the column lines. Then the column lines are sequentially connected by MOSFET switches to the output sense line under the control of the horizontal scan switch. The main limitation of the MOS (1T) FPA is a relatively high readout noise (on the order of 100 electrons/pixel) due to sensing of small charge signals on large-capacitance column lines. This readout noise can be decreased with a row readout MOS (2T) FPA having two MOSFET switches per detector. In this case, low readout noise can be achieved using current sensing, it is limited by the noise of the amplifier, and for voltage sensing it can be reduced by correlated double sampling (CDS).¹⁵ A readout noise of 300 rms electrons/pixel was achieved at Sarnoff for a 640×480 low-noise PtSi MOS (2T) FPA designed with row buffers and 8:1 multiplexing of the output lines;¹⁶ 2T MOS FPA read noise of $60e^-$ was later achieved by reading via capacitive transimpedance amplifier column buffers in $0.5 \mu\text{m}$ CMOS technology.¹⁷

An alternative form of the MOS (1T) FPA architecture is an MOS FPA with parallel column readout for fast frame operation. This silicon VIP FPA, resembling CIM architecture in Fig. 6d, can be used with HgCdTe vertically integrated PV detectors.

Direct-Charge-Injection Silicon FPAs All of the silicon monolithic FPAs thus far described use separately defined detectors. A direct-charge-injection type monolithic silicon FPA with a single detector surface is a PtSi direct Schottky injection (DSI) imager that is made on thinned silicon substrate having a CCD or MOS readout on one side and PtSi SBD charge-detecting surface on the other side.¹⁸ A cross-sectional area of one pixel of this FPA for IT-CCD readout is shown in Fig. 5f. In the operation of this imager, the *p*-type buried-channel CCD formed in an *n*-well removes charge from a P^+ charge-collecting electrode that in turn depletes a high-resistivity *p*-type substrate. Holes injected from the PtSi SBD surface into the *p*-type substrate drift through the depleted *p*-type substrate to the P^+ charge-collecting electrode. The advantages of the DSI FPA include 100 percent fill factor, a large maximum charge due to the large capacitance between the charge collecting electrode and the overlapping gate, and that the detecting surface does not have to be defined. A 128×128 IT-CCD PtSi DSI FPA was demonstrated;¹⁹ however, the same basic structure could also be used with other internal photoemission surfaces such as IrSb or Ge:Si.

Microbolometer FPAs A microbolometer FPA for uncooled applications consists of thin-film semiconductor photoresistors micromachined on a silicon substrate. The uncooled IR FPA is fabricated as an array of microbridges with a thermoresistive element in each microbridge. The

resistive microbolometers have high thermal coefficient of resistance (TCR) and low thermal conductance between the absorbing area and the readout circuit which multiplexes the IR signal. As each pixel absorbs IR radiation, the microbridge elemental resistance changes accordingly with its temperature.

Metal films have traditionally been used to make the best bolometer detectors because of their low $1/f$ noise. These latest devices use semiconductor films of 500 Å thickness having TCR of 2 percent per °C. The spacing between the microbridge and the substrate is selected to maximize the pixel absorption in the 8- to 14- μm wavelength range. Standard photolithographic techniques pattern the thin film to form detectors for individual pixels. The thin film TCR varies over an array by ± 1 percent, and produces responsivity of 70,000 V/W in response to 300 K radiation. This has been sufficient to yield $0.1^\circ\text{C NE } \Delta T$ with an $f/1$ lens. Potential low-cost arrays at prices similar to that of present large IC memories are possible with this technology.

Scanning and Staring Monolithic FPAs

In an earlier section we have reviewed the available architectures for the construction of two-dimensional scanning TDI, scanning, and staring FPAs. The same basic readout techniques, however, are also used for line-sensing imagers with photodiodes and MIS photogate detectors. For example, a line-scanning FPA corresponds to a vertical column CCD with the associated photodiodes of an IT-CCD, a column readout CID, a MOS (IT) FPA with only one row of detectors. However, since the design of a staring FPA is constrained by the size of the pixels, there is more space available for the readout multiplexer of a line of detectors. Therefore, design of the monolithic silicon multiplexer for a line detector array may also resemble the complexity of silicon multiplexer for hybrid FPAs.

33.4 HYBRID FPAs

Hybrid FPAs are made by interconnecting, via either direct or indirect means, a detector array to a multiplexing readout. Several approaches are pursued in both two- and three-dimensional configurations. Hybrids are typically made by either epoxying detector material to a processed silicon wafer (or readout) and subsequently forming the detectors and electrical interconnects by, for example, ion-milling; by mating a fully processed detector array to a readout to form a “two-dimensional” hybrid;²⁰ or by mating a fully processed detector array to a stack of signal processors to form a three-dimensional stack (Z-hybrid or “3D-IC”). The detector is usually mounted on top of the multiplexer and infrared radiation impinges on the backside of the detector array. Indium columns typically provide electrical and, often in conjunction with various epoxies, mechanical interconnect.

Hybrid methodology allows independent optimization of the detector array and the readout. Silicon is the preferred readout material due to performance and the leveraging of the continuous improvements funded by commercial markets. Diverse state-of-the-art processes and lithography are hence available at a fraction of their original development cost.

Thermal expansion match In a hybrid FPA, the detector array is attached to a multiplexer which can be of a different material. In cooling the device from room temperature to operating temperature, mechanical strain builds up in the hybrid due to the differing coefficients of thermal expansion. Hybrid integrity requires detector material that has minimum thermal expansion mismatch with silicon. Based on this criterion, the III-V and II-VI detectors are favored over Pb-salts. Silicon-based detectors are matched perfectly to the readout; these include doped-Si and PtSi. The issue of hybrid reliability has prompted the fabrication of II-VI detectors on alternative substrates to mitigate the mismatch. HgCdTe, for example, is being grown on sapphire (PACE-I),²¹ GaAs (liftoff techniques are available for substrate thinning or removal), and silicon in addition to the lattice-matched Cd(Zn)Te substrates. Detector growth techniques²² include liquid phase epitaxy and vapor phase epitaxy (VPE). The latter includes metal organic chemical vapor deposition (MOCVD) and molecular beam epitaxy (MBE).

Hybrid Readout

Hybrid readouts perform the functions of detector interface, signal processing, and video multiplexing.²³ The hybrid FPA readout technologies include

- Surface channel charge coupled device (SCCD)
- Buried channel charge coupled device (BCCD)
- x - y addressed switch-FET (SWIFET) or direct readout (DRO) FET arrays
- Combination of MOSFET and CCD (MOS/CCD)
- Charge-injection device (CID)

Early hybrid readouts were either CIDs²⁴ or CCDs, and the latter are still popular for silicon monolithics. However, x - y arrays of addressed MOSFET switches are superior for most hybrids for reasons of yield, design flexibility, simplified interface, and direct leveraging of Moore's Law for ongoing improvements and cost reduction. The move to the FET-based, direct readouts is key to the dramatic improvements in staring array producibility and is a consequence of the spin-off benefits from the silicon memory markets. DROs are fabricated with high yield and are fully compatible with advanced processes that are available at captive and commercial foundries. We will thus focus our discussion on these families. Though not extensively, CCDs are still sometimes used in hybrid FPAs.²⁵

Nonsilicon readouts Readouts have been developed in Ge, GaAs, InSb, and HgCdTe. The readout technologies include monolithic CCD, charge-injection device (CID), charge-injection matrix (CIM), enhancement/depletion (E/D) MESFET (GaAs), complementary heterostructure FET (C-HFET; GaAs),²⁶ and JFET (GaAs and Ge). The CCD, CID, and CIM readout technologies generally use MIS detectors for monolithic photon detection and signal processing.

The CID and CIM devices rely on accumulation of photogenerated charge within the depletion layer of a MIS capacitor that is formed using a variety of passivants (including CVD and photo-SiO₂, anodic SiO₂, and ZnS). A single charge transfer operation then senses the accumulated charge. Device clocking and signal readout in the CIDs and CIMs relies on support chips adjacent to the monolithic IR FPA. Thus, while the FPA is monolithic, the FPA assembly is actually a multichip hybrid.

CCDs have been demonstrated in HgCdTe and GaAs.²⁷ The n -channel technology is preferred in both materials for reasons of carrier mobility and device topology. In HgCdTe, for example, n -MOSFETs with CVD SiO₂ gate dielectric have parameters that are in good agreement with basic silicon MOSFET models. Fairly elaborate circuits have been demonstrated on CCD readouts, e.g., an on-chip output amplifier containing a correlated double sampler (CDS).

GaAs has emerged as a material that is very competitive for niche applications including IR FPAs. Since GaAs has very small thermal expansion mismatch with many IR detector materials including HgCdTe and InSb, large hybrids are possible, and VPE detector growth capability suggests future development of composite monolithic FPAs. The heterostructure (H-)MESFET and C-HFET technologies are particularly interesting for IR FPAs because low $1/f$ noise has been demonstrated; noise spectral densities at 1 Hz of as low as 0.5 $\mu\text{V}/\sqrt{\text{Hz}}$ for p -HIGFET and 2 $\mu\text{V}/\sqrt{\text{Hz}}$ for the enhancement H-MESFET²⁸ at 77 K have been achieved. The H-MESFET has the advantage of greater fabrication maturity (16 K SRAM and 64 \times 2 readout demonstrated), but the C-HFET offers lower power dissipation.

Direct Readout Architectures The DRO multiplexer consists of an array of FET switches. The basic multiplexer has several source follower stages that are separated at the cell, row, and column levels by MOSFET switches which are enabled and disabled to perform pixel access, reset, and multiplexing. The signal voltage from each pixel is thus direct-coupled through the cascaded source follower architecture as shown, for example, in Fig. 7. Shift registers generate the various clock signals; a minimum of externally supplied clocks is required. Since CMOS logic circuitry is used, the clock levels do not require precise adjustment for optimum performance. The simple architecture also gives high functional yield even in readout materials less mature than silicon.

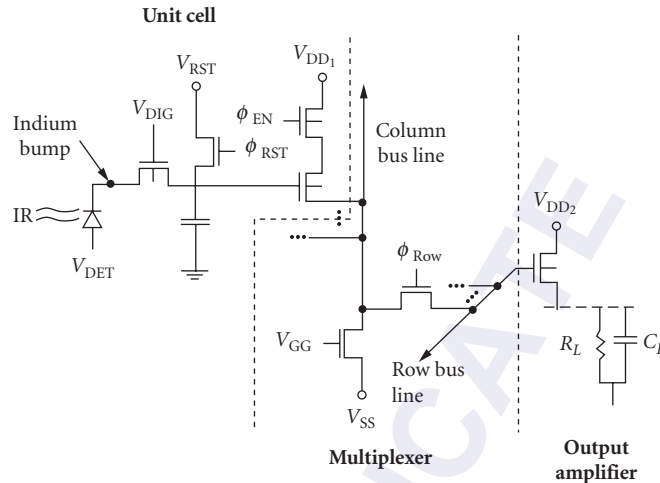


FIGURE 7 Direct readout schematic (shown with direct injection input).

Owing to the relatively low internal impedances beyond the input circuit, multiplexer noise is usually negligible. The inherent dynamic range is often >100 dB and the FPA dynamic range is limited only by the output-referred noise of the input circuit and the maximum signal excursion. The minimum read noise for DROs in imaging IR FPAs is typically capacitor reset noise. Correlated double samplers are thus used to suppress the reset noise for highest possible SNR at low integrated signal level.

In addition to excellent electrical characteristics, the DRO has excellent electro-optical properties including negligible MTF degradation and no blooming. Crosstalk in DRO-based FPAs is usually detector-limited since the readouts typically have low (<0.005 percent) electrical crosstalk. DROs also have higher immunity to clock feedthrough noise due to their smaller clock capacitances. Substrate charge pumping, which causes significant FET backgating²⁹ and transconductance degradation in SCCDs, is low in DROs.

X-Y addressing and clock generation Both static and dynamic shift registers are used to generate the clock signals needed for cell access, reset, and pixel multiplexing. Static registers offer robust operation and increased hardness to ionizing radiation in trade for increased FET count and preference for CMOS processes. Dynamic registers use fewer transistors in NMOS or PMOS processes, but require higher voltages, have lower maximum clocking rate, and must be carefully designed to avoid being affected by incident radiation.

Dynamic shift registers use internal bootstrapping to regenerate the voltage at each tap. The circuit techniques limit both the lowest and highest clock rates and require fine-tuning of the MOSFET design parameters for the specific operating frequency. More importantly, the high internal voltages stress conventional CMOS processes.

Electronically scanned staring FPAs The inability to integrate photogenerated charge for full staring frame times is often handled by integration time management. Since the photon background for the full 8- to 12- μm spectral band is over two orders of magnitude larger than the typical MWIR passband, and since the LWIR detector dark current is several orders of magnitude larger than similarly sized MWIR devices, LWIR FPA integration duty cycle can be quite poor. It is sometimes prudent to concede the limited duty cycle by electronically scanning the staring readout. Electronic scanning refers to a modified staring FPA architecture wherein the FPA is operated like a scanning FPA but without optomechanical means. Sensitivity is enhanced beyond that of a true scanning FPA by, for example, using multiple readout bus lines to allow integration times longer than one row time. The sharing reduces circuit multiplicity and frees unit cell real estate to share circuitry and

larger integration capacitance, and to use the otherwise parasitic bus capacitance to further increase capacity. More charge can thus be integrated even though the duty cycle is preset to $1/N$, where N is the number of elements on each common bus.

Time delay integration scanning FPAs While no longer extensively used for staring readouts, SCCDs are used in scanning readouts to incorporate on-focal plane TDI since they have higher dynamic range than FET bucket brigades. Dynamic range >72 dB and as high as 90 dB are typically achieved with high TDI efficacy. Two architectures dominate. In one, the CCD is integrated adjacent to the input circuit in a contiguous unit cell. In the second, the input circuit is segregated from the CCD in a sidecar configuration. The latter offers superior cell-packing density and on-chip signal processing in trade for circuit complexity. Figure 8 shows the schematic circuit for a channel of a scanning readout having capacitive transimpedance amplifier input circuit (discussed earlier), common TDI channel bus, and fill-and-spill³⁰ input to a sidecar SCCD TDI. This scheme integrates CMOS and CCD processes for much on-chip signal processing in very fine orthoscan pitch.³¹

The readout conversion factor, i.e., volts out per electrons in, for the sidecar CTIA scheme is

$$S_v = \frac{\Delta V}{e^-} = \frac{C_{F/S}}{C_T} A_{V_1} A_{V_2} \frac{q}{C_{out}} \quad (4)$$

where $C_{F/S}$ is the fill-and-spill gate capacitance, C_T is the integration/feedback capacitance, A_{V_x} characterizes the various source follower gains, and C_{out} is the sense node capacitance at the CCD output. The ratio of $C_{F/S}$ to C_T sets a charge gain that allows design-tailoring for managing dynamic range or lowering input-referred noise. High charge-gain yields read noise that is limited by the input circuit and not by the transfer noise³² of the high-carrier-capacity SCCD.

MOSFET bucket brigades are also used as TDI registers since simpler, all-MOS designs and processes can be used. Advantages include compatibility with standard MOS and CMOS, and capability for external clocking using specific CMOS-compatible clock levels. The latter potential advantage is mitigated in the sidecar TDI scheme by appropriately sizing the SCCD registers and the charge gain to yield the desired CCD clock levels, for example. Disadvantages include higher TDI register noise due to kTC noise being added at each transfer and limited signal excursion.

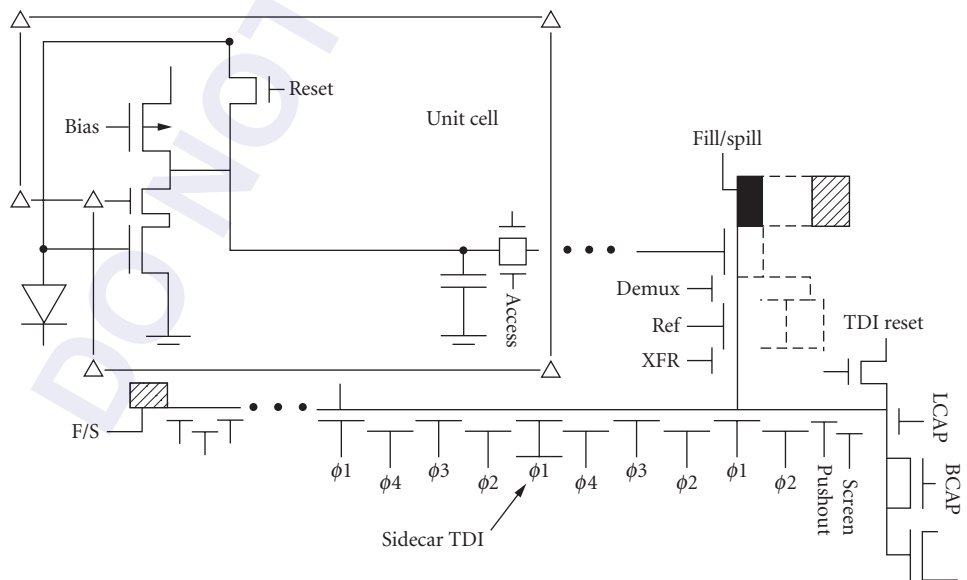


FIGURE 8 Sidecar TDI with capacitive transimpedance amplifier input circuit.

Output circuits Output circuitry is usually kept to a minimum to minimize power dissipation. Circuit design thus tends to focus on the trades between voltage-mode and current-mode output amplifiers, although on-chip signal processing is now at system-on-chip (SoC) level of sophistication, including low-speed A/D conversion, switched-capacitor filtering,³³ and on-chip nonuniformity correction. Voltage-mode outputs offer better S/N performance across a wider range in backgrounds for a given readout transimpedance. Current-mode outputs offer wider bandwidth and better drive capability at higher clock frequencies.

Detector Interface: Input Circuit After the incoming photon flux is converted into a signal by the detector, it is coupled into the readout via a detector interface circuit.³⁴ Signal input is optical in a monolithic FPA, so signal conditioning is limited. In hybrid FPAs and some composite material monolithics, the signal is injected electrically into the readout. The simplest input schemes offering the highest mosaic densities include direct detector integration (DDI) and direct injection (DI). More complex schemes trade simplicity for input impedance reduction [buffered direct injection (BDI) and capacitive transimpedance amplification (CTIA)], background suppression (e.g., gate modulation), or ultralow read noise with high speed (CTIA). We briefly describe the more popular schemes and their performance. Listed in Table 1 are approximate performance-describing equations for comparing the circuits schematically shown in Fig. 9.

Direct detector integration Direct detector integration (Fig. 9a), also referred to as source follower per detector (SFD), is used at low backgrounds and long frame times (frame rates typically ≤ 15 Hz in large staring arrays). Photocurrent is stored directly on the detector capacitance, thus requiring the detector to be heavily reverse-biased to maximize dynamic range. The changing detector voltage modulates the gate of a source follower whose drive FET is in the cell and whose current source is common to all the detectors in a column or row. The limited cell area constrains the source followers' drive capability and thus the bandwidth.

The DDI unit cell typically consists of the drive FET, cell enable transistor(s), and reset transistor(s). A detector site is read out by strobing the appropriate row clock, thus enabling the output source follower. The DDI circuit is capable of read noise as low as a few electrons per pixel.

Direct injection Direct injection (Fig. 9b) is perhaps the most widely used input circuit due to its simplicity and high performance. The detector directly modulates the source of a MOSFET. The direct coupling requires that detectors with p -on- n polarity, as is the case with InSb and most photovoltaic LWIR detectors, interface p -type FETs (and vice versa) for carrier collection in the integration capacitor. In surface channel CCDs, the input transistor's drain is virtual, as formed by a fully enhanced well, and often doubles as the integration capacitor.

Practical considerations, including limited charge-handling capacity, constrain the DI input to operation with high-impedance MWIR or limited cutoff ($\lambda_c \leq 9.5 \mu\text{m}$) LWIR detectors. The associated background photocurrent for the applications where direct injection can be used mandates that the DI FET operate subthreshold.³⁵ The subthreshold gate transconductance, g_m , is independent of FET geometry:³⁶

$$g_m = \left(\frac{\partial I_D}{\partial V_G} \right) \Big|_{V_{DS}=\text{constant}} = \frac{q \left(\eta_{\text{inj}} \left\{ I_{\text{photo}} + \frac{V_{\text{det}}}{R_{\text{det}}} - I_{\text{det}_0} (e^{(qV_{\text{det}}/nkT)} - 1) \right\} \right)}{nkT} \approx \frac{qI_D}{nkT} \quad (5)$$

The injection efficiency, η_{inj} , of detector current into the DI FET is

$$\eta_{\text{inj, DI}} = \frac{g_m R_{\text{det}}}{1 + g_m R_{\text{det}}} \left[\frac{1}{1 + \frac{j\omega C_{\text{det}} R_{\text{det}}}{1 + g_m R_{\text{det}}}} \right] \quad (6)$$

where R_{det} and C_{det} are the detectors' dynamic resistance and capacitance, respectively. Poor DI circuit bandwidth occurs at low-photon backgrounds due to low g_m .

TABLE 1 Focal Plane Array Performance

Input Circuit	Percentage of BLIP	Detector Noise	Input-Referred Circuit Noise	Input-Referred MUX Noise	Transimpedance
Direct Detector Integration	$\left[\frac{(\eta A_{\text{det}} Q_{\beta} \tau_{\text{int}})^{1/2}}{\eta A_{\text{det}} Q_{\beta} \tau_{\text{int}} + q R_{\text{det}}^{\text{int}}} + N_{\text{f}}^2 + \frac{2kT \tau_{\text{int}}}{q R_{\text{det}}^{\text{int}}} \right]^{1/2}$	$I_{\text{det}}^2 = \left[\frac{4kT}{R_{\text{det}}} + 2qI_{\text{det}} \right] \Delta f + \int \left(\frac{K_{\text{det}}}{f} \right) df$	$N_{\text{sf}} \approx \sqrt{2} \int_n^M \sqrt{V^2(f) \frac{(1 - \cos 2\pi f T_D)}{[1 + (2\pi f T_D)^2]} df}$ $S_V = \left(\frac{C_{\text{det}} A_V}{q} \right)^{-1}$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{A_V^2} kTC_{\text{mux}} \Delta f$	$\left(\frac{t_{\text{int}}}{C_{\text{det}}} \right) A_V$
Direct Injection	$\frac{(2qI_{\text{photo}} \Delta f)^{1/2}}{(\sigma_{\text{det}}^2 + \sigma_{\text{input,ir}}^2 + \sigma_{\text{mux,ir}}^2)^{1/2}}$		$\sigma_{\text{input,ir}}^2 = \int \left[\frac{\Delta f}{1 + \omega^2 C_{\text{det}}^2 R_{\text{det}}^2} \left(\frac{8}{3} kT g_m + \frac{K_{\text{FET}}}{f \alpha} \right) df \right]$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{\eta_{\text{lin}}^2} kTC_{\text{input}} \Delta f$	$\left(\frac{t_{\text{int}}}{C_{\text{int}}} \right) A_V$
Buffered Direct Injection	$\frac{(2qI_{\text{photo}} \Delta f)^{1/2}}{(\sigma_{\text{det}}^2 + \sigma_{\text{input,ir}}^2 + \sigma_{\text{mux,ir}}^2)^{1/2}}$		$\sigma_{\text{input,ir}}^2 = \int \left[\eta_{\text{noise}}^2 \left(\frac{8}{3} kT g_m + \frac{K_{\text{FET}}}{f \alpha} \right) + A_{\text{amp}}^2 \left(\epsilon_{\text{amp}} \right) df \right]$ $\eta_{\text{noise}} = \frac{1 + j\omega R_{\text{det}} C_{\text{int}}}{1 + (1 + A_V) g_m R_{\text{det}} + j\omega R_{\text{det}} (1 + A_V) C_{\text{int}}}$ $A_{\text{amp}} = \left(\frac{g_m}{R_{\text{det}}} \right) \frac{(1 + j\omega R_{\text{det}} C_{\text{det}})}{1 + (1 + A_V) g_m R_{\text{det}} + j\omega C_{\text{det}} (1 + A_V) C_{\text{int}} R_{\text{det}} }$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{\eta_{\text{lin}}^2} kTC_{\text{input}} \Delta f$	$\left(\frac{t_{\text{int}}}{C_{\text{int}}} \right) A_V$
Chopper-Stabilized BDI	$\frac{(2qI_{\text{photo}} \Delta f)^{1/2}}{(\sigma_{\text{det}}^2 + \sigma_{\text{input,ir}}^2 + \sigma_{\text{mux,ir}}^2)^{1/2}}$		$\sigma_{\text{input,ir}}^2 = \int \left[\eta_{\text{noise}}^2 \left(\frac{8}{3} kT g_m \right) + A_{\text{amp}}^2 \left(\epsilon_{\text{amp}} \right) df \right]$ $\eta_{\text{noise}} = \left(\frac{g_m}{1 + (1 + A_V) g_m R_{\text{det}} + j\omega R_{\text{det}} C_{\text{int}}} \right) \frac{1 + j\omega R_{\text{det}} C_{\text{int}}}{1 + j\omega R_{\text{det}} C_{\text{det}}}$ $A_{\text{amp}} = \left(\frac{g_m}{R_{\text{det}}} \right) \frac{(1 + j\omega R_{\text{det}} C_{\text{det}})}{1 + (1 + A_V) g_m R_{\text{det}} + j\omega R_{\text{det}} C_{\text{det}} (1 + A_V) C_{\text{int}} R_{\text{det}} }$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{\eta_{\text{lin}}^2} kTC_{\text{input}} \Delta f$	$\left(\frac{t_{\text{int}}}{C_{\text{int}}} \right) A_V$
Gate Modulation (FET Load)	$\frac{(2qI_{\text{photo}} \Delta f)^{1/2}}{(\sigma_{\text{det}}^2 + \sigma_{\text{load}}^2 + \sigma_{\text{input,ir}}^2 + \sigma_{\text{mux,ir}}^2)^{1/2}}$		$\sigma_{\text{input,ir}}^2 = \int \left[\frac{\Delta f}{A_V^2} \left(2qI_{\text{input}} + \frac{K_{\text{FET,input}}}{f \alpha} \right) df \right]$	$\sigma_{\text{mux,ir}}^2 = \frac{1}{A_{\text{lin}}^2} kTC_{\text{input}} \Delta f$	$\left(\frac{A_V t_{\text{int}}}{C_{\text{int}}} \right) A_V$
Capacitive Transimpedance Amplifier	$\left[\frac{\eta A_{\text{det}} Q_{\beta} \tau_{\text{int}}}{\eta A_{\text{det}} Q_{\beta} \tau_{\text{int}} + \frac{2kT \tau_{\text{int}}}{q R_{\text{det}}^{\text{int}}} + N_{\text{amp,lif}}^2 + N_{\text{amp,white}}^2 + N_{\text{load,white}}^2 \right]^{1/2}$	$N_{\text{amp,lif}} \approx \frac{C_{\text{det}} K_{\text{amp}} \sqrt{2}}{q} \left[\frac{5\tau_{\text{int}}}{\tau_{\text{amp}}} \right]$ $N_{\text{amp,white}} \approx \frac{C_{\text{det}}}{q} \sqrt{\frac{8}{3} \frac{kT}{g_m \tau_{\text{amp}}}}$ $N_{\text{load,white}} \approx \frac{C_{\text{det}}}{q A_V \text{amp} A_V^2} \sqrt{\frac{2kT}{C_L}}$		$I_{\text{mux,ir}}^2 = kTC_{\text{input}} \Delta f$	$Z_T = \frac{\tau_{\text{int}}}{C_T} A_V$ $C_T = \frac{(C_{\text{gd}} + C_{\text{gs}} + C_{\text{det}}) + A_V (C_{\text{fb}} + C_{\text{gd}})}{A_V}$

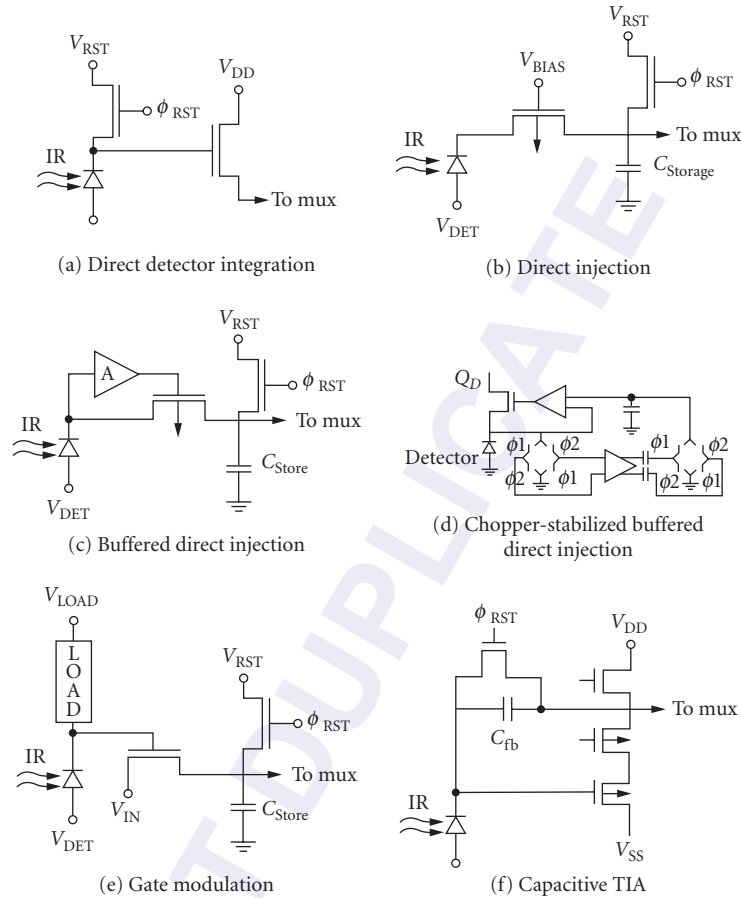


FIGURE 9 Hybrid FPA detector interface circuits.

The injection efficiency varies across an FPA due to FET threshold, detector bias, and detector resistance nonuniformity. Changes in detector current create detector bias shifts since the input impedance is relatively high. In extreme cases a large offset in threshold gives rise to excess detector leakage current and $1/f$ noise, in addition to fixed pattern noise. The peak-to-peak threshold voltage nonuniformity spans the range from ≈ 1 mV for silicon p -MOSFETs to over 125 mV for some GaAs-based readouts.

Depending upon the interface to the multiplexing bus, the noise-limiting capacitance, C_{input} , is approximately the integration capacitance or the combined integration and bus capacitance. Some direct-injection cells are thus buffered with a source follower (see Fig. 7). Omitting the source follower reduces the readout transimpedance (due to charge splitting between the integration capacitor and thus bus line capacitance) in trade for larger integration capacitance since more unit-cell real estate is available.

Increasing the pixel density has required a continuing reduction in cell pitch. Figure 10 plots the charge-handling capacity as functions of cell pitch and minimum gate length for representative DI designs using the various minimum feature lengths. Also plotted is the maximum capacity assuming the cell is composed entirely of integration capacitor ($225 \text{ \AA} \text{ iO}_2$). A $27\text{-}\mu\text{m}$ DI cell, fabricated in $1.25\text{-}\mu\text{m}$ CMOS, thus had similar cell capacity as an earlier $60\text{-}\mu\text{m}$ DI cell in $3\text{-}\mu\text{m}$ CMOS. Limitations on cell real estate, operating voltage, and the available capacitor dielectrics nevertheless

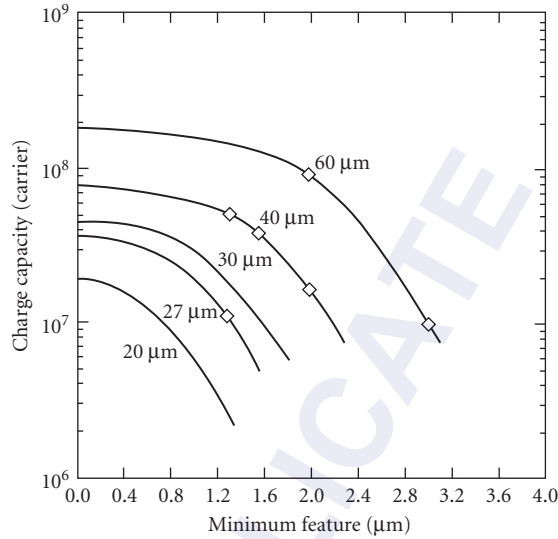


FIGURE 10 Direct injection charge-handling capacity versus cell pitch for standard CMOS processes.

dictate maximum integration times that are often shorter than the total frame time. This duty cycling equates to degradation in detective quantum efficiency.

Buffered direct injection A significant advantage of the source-coupled input is MOSFET noise suppression. This suppression is implied in the η_{BLIP} expression shown in Table 1. When injection efficiency is poor, however, MOSFET noise becomes a serious problem along with bandwidth. These deficiencies are ameliorated via buffered direct injection (BDI),³⁷ wherein a feedback amplifier with open-loop gain $-A_v$ (Fig. 9c) is added to the DI circuit. The buffering increases injection efficiency to near-unity, increases bandwidth by orders of magnitude, and suppresses the DI FET noise.

BDI has injection efficiency

$$\eta_{\text{inj}} \equiv \frac{g_m R_{\text{det}} (1 + A_v)}{1 + g_m R_{\text{det}} (1 + A_v) + j\omega R_{\text{det}} [(1 + A_v) C_{\text{amp}} + C_{\text{det}}]} \quad (7)$$

where C_{amp} is the Miller capacitance of the amplifier. Circuit bandwidth is maximized by lowering the amplifiers' Miller capacitance to provide detector-limited frequency response that is lower than that possible with DI by the factor $(1 + A_v)$.

The noise margin of the BDI circuit is superior to DI, even though two additional noise sources associated with the feedback amplifier are added. The dominant circuit noise stems from the drive FET in the amplifier. The noise power for frequencies less than $1/(2\pi R_{\text{det}} C_{\text{det}})$ is directly proportional to the detector impedance. Amplifier noise is hence a critical issue with low impedance ($<1 \text{ M}\Omega$) detectors including long wavelength photovoltaics operating at temperatures above 80 K.

Chopper-stabilized buffered direct injection The buffered direct injection circuit has high $1/f$ noise with low-impedance detectors. While the MOSFET $1/f$ noise can be suppressed somewhat by reverse-biasing the detectors to the point of highest resistance, detector $1/f$ noise may then dominate. Other approaches include enlarging MOSFET gate area, using MOS input transistors in the lateral bipolar mode, or using elaborate circuit techniques such as autozeroing and chopper stabilization. Chopper stabilization is useful if circuit real estate is available, as in a scanning readout. Figure 9d shows a block diagram schematic circuit of chopper-stabilized BDI.

Chopper stabilization refers to the process of commutating the integrating detector node between the inverting and noninverting inputs of an operational amplifier having open-loop gain, A_v . This chopping process shifts the amplifier's operating frequency to higher frequencies where the amplifier's noise is governed by white noise, not $1/f$ noise. At chopping frequencies $f_{\text{chop}} \gg f_c$, the equivalent low-frequency input noise of the chopper amplifier is equal to the original amplifier white-noise component.³⁸ The amplifier's output signal is subsequently demodulated and filtered to remove the chopping frequency and harmonics. This scheme also reduces the input offset nonuniformity by the reciprocal of the open-loop gain, thereby generating uniform detector bias. Disadvantages include high circuit complexity and the possibility of generating excess detector noise via clock feedthrough-induced excitation of traps, particularly with narrow bandgap photovoltaic detectors.

Gate modulation Signal processing can be incorporated in a small unit cell by using a gate-modulated input structure (c.f. MOSFET load gate modulation in Fig. 9e).³⁹ The use of an MOSFET as an active load device, for example, provides dynamic range management via automatic gain control and user-adjustable background pedestal offset since the detector current passes through a load device with resistance R_{LOAD} . The differential gate voltage applied to the input FET varies for a change in photocurrent, ΔI_{photo} , as

$$\Delta V_G = R_{\text{LOAD}} \eta_{\text{inj, DI}} \Delta I_{\text{photo}} \quad (8)$$

The current injected into the integration capacitor is

$$I_{\text{input}} = g_m R_{\text{LOAD}} \eta_{\text{inj, DI}} \Delta I_{\text{photo}} \quad (9)$$

The ratio of I_{input} to I_{photo} is the current gain, A_I , which is

$$A_I = \frac{g_m}{g_{m, \text{LOAD}}} \eta_{\text{inj, DI}} \quad (10)$$

The current gain can self-adjust by orders of magnitude depending on the total detector current. Input-referred read noise of tens of electrons has thus been achieved with high-impedance SWIR detectors at low-photon backgrounds. The same circuit has also been used at LWIR backgrounds with LWIR detectors having adequate impedance for good injection efficiency.

The current gain expression suggests a potential shortcoming for imaging applications since the transfer characteristic is nonlinear, particularly when the currents in the load and input FETs differ drastically. In conjunction with tight specifications for threshold uniformity, pixel functionality can be decreased, dynamic range degraded, and imagery dominated by spatial noise. The rms fractional gain nonuniformity (when operating subthreshold) of the circuit is approximately

$$\frac{\Delta A_I}{A_I} = \frac{q \sigma_{VT}}{n_{\text{FET}} kT} \quad (11)$$

where σ_{VT} is the rms threshold nonuniformity. At 80 K, state-of-the-art σ_{VT} of 0.5 mV for a 128 × 128 FPA, and $n = 1$, the minimum rms nonuniformity is ≈ 7 percent.

Capacitive transimpedance amplifier (CTIA) Many CTIAs have been successfully demonstrated. The most popular approach uses a simple CMOS inverter⁴⁰ for feedback amplification (Fig. 9f). Others use a more elaborate differential amplifier. The two schemes differ considerably with respect to open-loop voltage gain, bandwidth, power dissipation, and cell real estate. The CMOS inverter-based CTIA is more attractive for high-density arrays. The latter is sometimes preferred for scanning or Z-hybrid applications where real estate is available primarily to minimize power dissipation.⁴¹

In either case, photocurrent is integrated directly onto the feedback capacitor of the transimpedance amplifier. The minimum feedback capacitance is set by the amplifier's Miller capacitance and defines the maximum circuit transimpedance. Since the Miller capacitance can be made very small (<5 fF), the resulting high transimpedance yields excellent margin with respect to

downstream system noise. The transimpedance degrades when the circuit is coupled to large detector capacitances, so reducing pixel size serves to minimize read noise and the circuits' attractiveness will continue to increase in the future.

The CTIA allows extremely small currents to be integrated with high efficiency and tightly regulated detector bias. The amplifier open-loop gain, $A_{V,amp}$, ranges from as low as on the order of ten to higher than several thousand for noncascoded inverters, and many thousands for cascoded inverter and differential amplifier designs. The basic operating principle is to apply the detector output to the inverting input of a high-gain CMOS differential amplifier operated with capacitive feedback. The feedback capacitor is reset at the detector sampling rate. The noise equivalent input voltage of the amplifier is referred to the detector impedance, just as in other circuits.

The CTIA's broadband channel noise sets a lower limit on the minimum achievable read noise, is the total amplifier white noise, and can be approximated for condition of large open-loop gain by

$$N_{kTC,channel}^2 = \frac{kTC_{fb}}{q^2} \left[\frac{C_{det} + C_{fb}}{C_L + \frac{C_{fb}C_{det}}{C_{fb} + C_{det}}} \right] \quad (12)$$

where C_{fb} is the feedback capacitance including the Miller capacitance and integration capacitance and C_L is the output load capacitance. This expression provides an intuitive formula for minimizing noise: detector capacitance must be low (i.e., minimize detector shunting capacitance which reduces closed-loop gain) and output capacitance high (i.e., limit bandwidth). Of the three amplifier noise sources listed in Table 1, amplifier $1/f$ noise is often largest. For this reason, the CMOS inverter-based CTIA has best performance with p -on- n detectors and p -MOSFET amplifier FET due to the lower $1/f$ noise.

33.5 PERFORMANCE: FIGURES OF MERIT

In the early days of infrared technology, detectors were characterized by the noise equivalent power (NEP) in a 1-Hz bandwidth. This was a good specification for single detectors, since their performance is usually amplifier limited. The need to compare detector technologies for application to different geometries and the introduction of FPAs having high-performance on-board amplifiers and small parasitics necessitated normalization to the square root of the detector area for comparing S/N. R. C. Jones⁴² thus introduced detectivity (D^*), which is simply the reciprocal of the normalized NEP and has units $\text{cm} \cdot \sqrt{\text{Hz}}/\text{W}$ (or Jones).

While D^* is well-suited for specifying infrared detector performance, it can be misleading to the uninitiated since the D^* is highest at low background. An LWIR FPA operating at high background with background-limited performance (BLIP) S/N can have a D^* that is numerically lower than for a SWIR detector having poor S/N relative to the theoretical limit. Several figures of merit that have hence proliferated include other ways of specifying detectivity: e.g., thermal D^* (D_{th}^*), blackbody D^* (D_{bb}^*), peak D^* (D_{pk}^*), percentage of BLIP (%BLIP or η_{BLIP}), and noise equivalent temperature difference ($NE\Delta T$). Since the final output is an image, however, the ultimate figure of merit is how well objects of varying size are detected and resolved in the displayed image. The minimum resolvable temperature (MRT) is thus a key benchmark. These are briefly discussed in this section.

Detectivity (D^*)

D^* is the S/N ratio normalized to the electrical bandwidth and detector area. In conjunction with the optics area and the electrical bandwidth, it facilitates system sensitivity estimation. However, D^* can be meaningless unless the test conditions, including magnitude and spectral distribution of the flux source (e.g., blackbody temperature), detector field-of-view, chopping frequency (lock-in amplifier),

background temperature, and wavelength at which the measurement applies. D^* is thus often quoted as “blackbody,” since the spectral responsivity is the integral of the signal and background characteristics convolved with the spectral response of the detector. D_{bb}^* specifications are often quantified for a given sensor having predefined scene temperature, filter bandpass, and cold shield $f/\#$ using a generalized expression.

Peak detectivity is sometimes preferred by detector engineers specializing in photon detectors. The background-limited peak detectivity for a photovoltaic detector is

$$D_{\lambda_{\text{pk}}}^* = \sqrt{\frac{\eta}{2Q_B}} \frac{\lambda_{\text{pk}}}{hc} \quad (13)$$

and refers to measurement at the wavelength of maximum spectral responsivity. For detector-limited scenarios, such as at higher operating temperatures or longer wavelengths (e.g., $\lambda_c > 12 \mu\text{m}$ at operating temperature less than 78 K or $\lambda_c > 4.4 \mu\text{m}$ at >195 K), the peak detectivity is limited by the detector and not the photon shot noise. In these cases the maximum detector-limited peak D^* in the absence of excess bias-induced noise (both $1/f$ and shot noise) is

$$D_{\lambda_{\text{pk}}}^* = \frac{\eta q}{2} \sqrt{\frac{R_0 A}{kT}} \frac{\lambda}{hc} \quad (14)$$

The $R_0 A$ product of a detector thus describes detector quality even though other parameters may actually be more relevant for FPA operation.

Percentage of BLIP

Whereas D^* compares the performance of dissimilar detectors, FPA designers often need to quantify an FPA's performance relative to the theoretical limit at a specific operating background. Percentage of BLIP, η_{BLIP} , is one such parameter and is simply the ratio of photon noise to composite FPA noise

$$\eta_{\text{BLIP}} = \left(\frac{N_{\text{PHOTON}}^2}{N_{\text{PHOTON}}^2 + N_{\text{FPA}}^2} \right)^{1/2} \quad (15)$$

NE ΔT

The NE ΔT of a detector represents the temperature change, for incident radiation, that gives an output signal equal to the rms noise level. While normally thought of as a system parameter, detector NE ΔT and system NE ΔT are the same except for system losses (conservation of radiance). NE ΔT is defined:

$$\text{NE } \Delta T = v_n \left(\frac{\partial T}{\partial Q} \right) \bigg/ \left(\frac{\partial V_s}{\partial Q} \right) = v_n \frac{\Delta T}{\Delta V_s} \quad (16)$$

where v_n is the rms noise and ΔV_s is the signal measured for the temperature difference ΔT . It can be shown that

$$\text{NE } \Delta T = \left(\tau_o C_{T\lambda} \eta_{\text{BLIP}} \sqrt{N_c} \right)^{-1} \quad (16a)$$

where τ_o is the optics transmission, $C_{T\lambda}$ is the thermal contrast from Fig. 1, and N_c is the number of photogenerated carriers integrated for one integration time, t_{int} :

$$N_c = \eta A_{\text{det}} t_{\text{int}} Q_B \quad (16b)$$

The distinction between an integration time and the FPA's frame time must be noted. It is often impossible at high backgrounds to handle the large amount of carriers generated over frame times compatible with standard video frame rates. The impact on system D^* is often not included in the FPA specifications provided by FPA manufacturers. This practice is appropriate for the user to assess relative detector quality, but must be coupled with usable FPA duty cycle, read noise, and excess noise to give a clear picture of FPA utility. Off-FPA frame integration can be used to attain a level of sensor sensitivity that is commensurate with the detector-limited D^* and not the charge-handling-limited D^* .

The inability to handle a large amount of charge nevertheless is a reason why the debate as to whether LWIR or MWIR operation is superior is still heated. While the LWIR band should offer order-of-magnitude higher sensitivity, staring readout limitations often reduce LWIR imager sensitivity to below that of competing MWIR cameras. However, submicron photolithography, alternative dielectrics, and refinements in readout architectures are ameliorating this shortfall and LWIR FPAs having sensitivity superior to MWIR counterparts are available. Figure 11 shows the effect on high quantum efficiency FPA performance and compares the results to PtSi at TV-type frame rate. The figure illustrates BLIP and measured NE ΔT s versus background temperature for several spectral bands assuming a 640×480 DI readout multiplexer (27- μm pixel pitch and ≈ 1 - μm -minimum

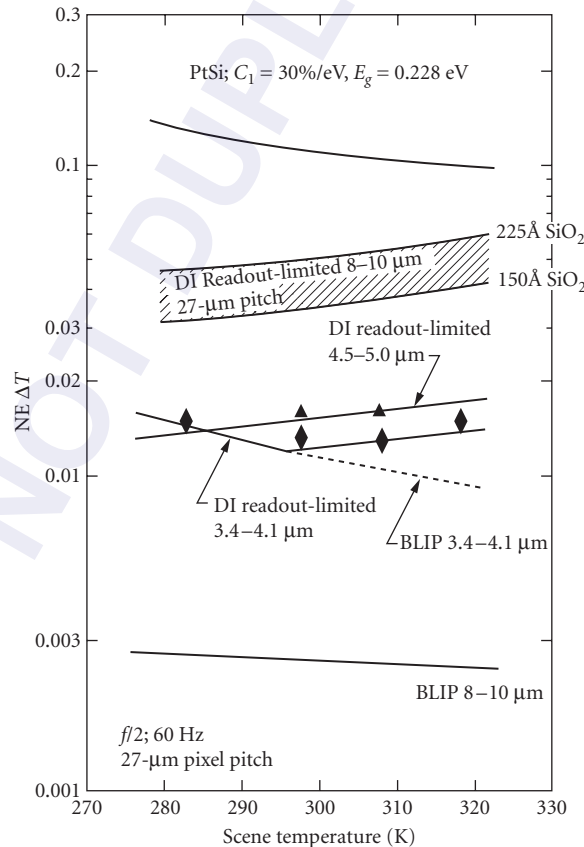


FIGURE 11 NE ΔT versus background temperature for several prominent spectral bands.

feature) is hybridized to high-quantum-efficiency MWIR and LWIR detectors. Though the device has large charge-handling capacity, there is large shortfall in the predicted LWIR FPA performance relative to the BLIP limit. The measured MWIR FPA performance values, as shown by the data points, are in good agreement with the predicted trends.

Spatial noise Estimation of IR sensor performance must include a treatment of spatial noise that occurs when FPA nonuniformities cannot be compensated correctly. This requires consideration of cell-to-cell response variations. Mooney et al.⁴³ comprehensively discussed the origin of spatial noise. The total noise determining the sensitivity of a staring array is the composite of the temporal noise and the spatial noise. The spatial noise is the residual nonuniformity U after application of nonuniformity compensation, multiplied by the signal electrons N . Photon noise, equal to \sqrt{N} , is the dominant temporal noise source for the high infrared background signals for which spatial noise is significant (except for TE-cooled or uncooled sensors). The total noise equivalent temperature difference is

$$\text{Total NE } \Delta T = \frac{\sqrt{N+U^2N^2}}{\frac{\partial N}{\partial T}} = \frac{\sqrt{1/N+U^2}}{\frac{1}{N} \frac{\partial N}{\partial T}} \quad (17)$$

where $\partial N/\partial T$ is the signal change for a 1 K source temperature change. The denominator, $(\partial N/\partial T)/N$, is the fractional signal change for a 1 K source temperature change. This is the relative scene contrast due to $C_{T\lambda}$ and the FPA's transimpedance.

The dependence of the total NE ΔT on residual nonuniformity is plotted in Fig. 12 for 300 K scene temperature, two sets of operating conditions, and three representative detectors: LWIR HgCdTe, MWIR HgCdTe, and PtSi. Operating case A maximizes the detected signal with $f/1.4$ optics, 30-Hz frame rate, and 3.4 to 5.0- μm passband. Operating case B minimizes the solar influence by shifting the passband to 4.2 to 5.0 μm and trades off signal for the advantages of lighter, less expensive optics ($f/2.0$) at 60-Hz frame rate. Implicit in the calculations are charge-handling capacities of 30 million e^- for MWIR HgCdTe, 100 million e^- for LWIR HgCdTe, and 1 million for PtSi. The sensitivity at the lowest nonuniformities is independent of nonuniformity and limited by the

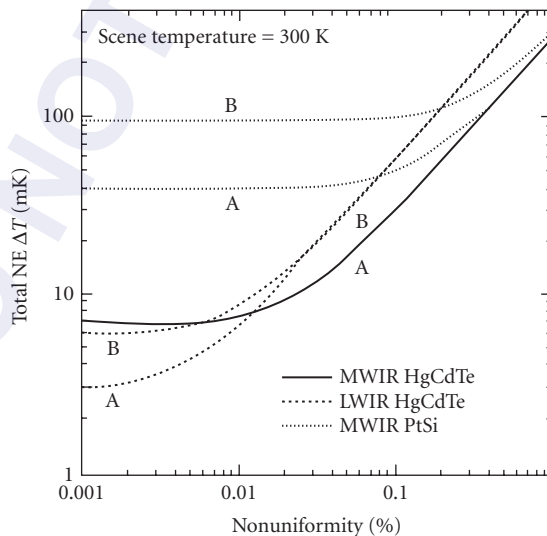


FIGURE 12 Total NE ΔT versus nonuniformity at 300 K scene temperature.

shot noise of the detected signal. The LWIR sensitivity advantage is achieved only at nonuniformities less than 0.01 percent, which is comparable to that achieved with buffered input circuits. At the reported direct-injection MWIR HgCdTe residual nonuniformity of 0.01 to 0.02 percent, the total $NE\Delta T$ is about 0.007 K, which is comparable to the MWIR BLIP limit. At the reported PtSi residual nonuniformity of 0.05 percent with direct detector integration, total $NE\Delta T$ is higher at 0.04K, but exceeds the BLIP limit for the lower quantum efficiency detectors.

Minimum Resolvable Temperature

The minimum resolvable temperature (MRT) is often the preferred figure of merit for imaging infrared sensors. MRT is a function of spatial resolution and is defined as the signal-to-noise ratio required for an observer to resolve a series of standard four-bar targets. While many models exist due to the influence of human psycho-optic response, a representative formula⁴⁴ is

$$MRT(f_s, T_{SCENE}) = \frac{2SNR_t NE \Delta T(T_{SCENE})}{MRT(f_s)} \left[\frac{f_s^2 \Delta x \Delta y}{L \tau_{eye} f_{frame} N_{OS} N_{SS}} \right]^{1/2} \quad (18)$$

where f_s is the spatial frequency in cycles/radian, a target signal-to-noise ratio (SNR_t) of five is usually assumed, the MTF describes the overall modulation transfer function including the optics, detector, readout, and the integration process, Δx and Δy are the respective detector subpixels in mRad, τ_{eye} is the eye integration time, f_{frame} is the display frame rate, N_{OS} is the overscan ratio, N_{SS} is the serial scan ratio, and L is the length-to-width ratio of a bar chart (always set to 7). While the MRT of systems with temporal noise-limited sensitivity can be adequately modeled using the temporal $NE\Delta T$, scan noise in scanned system and fixed pattern noise in staring cameras requires that the MRT formulation be appropriately modified.

Shown in Fig. 13 are BLIP (for 70 percent quantum efficiency) MRT curves at 300 K background temperature for narrow-field-of-view (high-resolution) sensors in the MWIR and LWIR spectral bands. Two LWIR curves are included to show the impact of matching the diffraction-limited blur

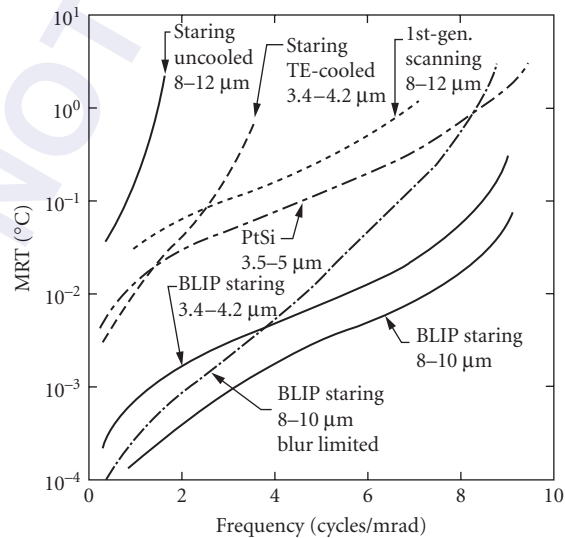


FIGURE 13 BLIP MRT for representative staring FPA configurations in the various bands.

to the pixel pitch versus $2\times$ oversampling of the blur. The latter case commonly arises, for example, when the LWIR FPA is miniaturized to minimize die size for enhancing yield and, for hybrid FPAs, alleviating thermal expansion mismatch. Also included for comparison are representative curves for first-generation scanning, staring uncooled, staring TE-cooled and staring PtSi sensors assuming 0.1, 0.1, 0.05, and 0.1 K NE Δ Ts, respectively. Theoretically, the staring MWIR sensors have order of magnitude better sensitivity while the staring LWIR bands have two orders of magnitude better sensitivity than the first-generation sensor. In practice, due to charge-handling limitations, an LWIR sensor has only slightly better MRT than the MWIR sensor. The uncooled sensor is useful for short-range applications such as a driver's aid in modern automobiles; the TE-cooled sensor provides longer range than the uncooled, but less than the high-density PtSi-based cameras.

33.6 CURRENT STATUS AND FUTURE TRENDS

Status

After over three decades of ongoing development, today's second-generation infrared focal plane arrays have typically 1000 times more pixels and up to 10 times higher sensitivity than first-generation devices. FPA format is consequently now set by application need, rather than a technological barrier. Nevertheless, there is ongoing motivation for fully achieving theoretically limited performance, especially at elevated operating temperatures. Development is hence continuing in the form of third-generation FPA technology. While specifications for third-generation FPAs encompass a broad range of needs, common objectives require overcoming recurring practical limits with respect to sensitivity, frame rate, power dissipation, multispectral capability, and cost. The common methodology going forward, regardless of specific mission requirements, is to dramatically increase on-FPA functionality via system-on-chip integration. While second-generation technology was largely driven by defense-oriented R&D funding, it is likely that third-generation technology will more directly leverage emerging commercial foundry capability for 3D-IC assembly; this emerging commercial interest in 3D-IC integration should dramatically lower infrared FPA cost while further improving FPA performance.

Figure 14 summarizes the typical performance of the most prominent detector technologies. The figure, a plot of the D_{th}^* (300 K, 0° field-of-view) versus operating temperature, clearly shows the performance advantage that the intrinsic photovoltaics have over the other technologies. Thermal detectivity is used here to compare the various technologies for equivalent NE Δ T irrespective of wavelength. While the extrinsic silicon detectors offer very high sensitivity, high producibility, and very long cutoff wavelengths, the very low operating temperature is often prohibitive. Also shown is the relatively low and slightly misleading detectivity of PtSi, which is offset by its

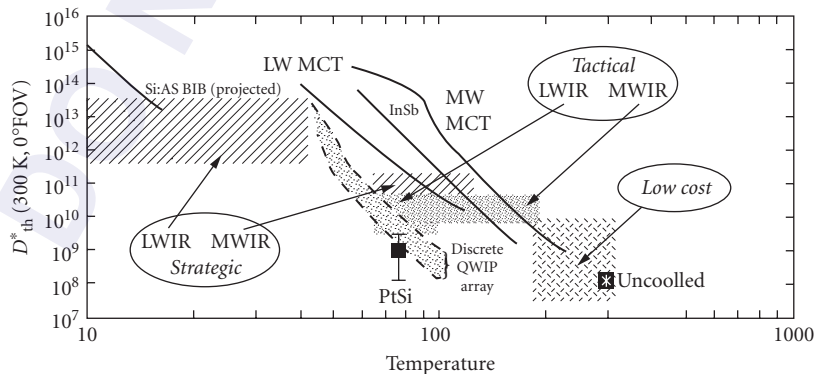


FIGURE 14 Comparison of photon detector D_{th}^* (300 K, 0°) for various IR technologies for equivalent NE Δ T.

Similar in performance at cryogenic temperatures, InSb and HgCdTe have comparable array size and pixel yield at MWIR cutoff wavelengths. Wavelength tunability and high radiative efficiency, however, have often made HgCdTe the preferred material because the widest possible bandgap semiconductor can be configured and thus the highest possible operating temperature achieved for a given set of operating conditions. The associated cooling and system power requirements can thus be optimally distributed.

FPA costs are currently very similar for all the second-generation FPA technologies. Though it is often argued that FPA cost for IR cameras will be irrelevant once it reaches a certain minimum production volume, nevertheless the key determinants as to which FPA technology becomes ubiquitous in the coming decade are availability and cost.

Future Trends and Technology Directions

The 1980s saw the maturation of PtSi and the emergence of HgCdTe and InSb as producible MWIR detector materials. Many indicators pointed to the 1990s being the decade that IR FPA technology would enter the consumer marketplace, but penetration did not actually occur until a decade later. Figure 16 shows the chronological development of IR FPAs including monolithic and hybrid technologies. Specifically compared is the development of various hybrids (primarily MWIR) to PtSi, the pace-setting technology with respect to array size. The hybrid FPA data includes Pb-salt, HgCdTe, InSb, and PtSi devices. This database suggests that monolithic PtSi led all other technologies with respect to array size by about 2 years and clearly shows the thermal mismatch barrier confronted by hybrid FPA developers in the mid-1980s.

In addition to further increases in pixel density to $>10^{16}$ pixels, several trends are clear. Future arrays will have much more on-chip signal processing, need less cooling, have higher sensitivity (particularly intrinsic LWIR FPAs), and offer multispectral capability. If the full performance potential of the uncooled technologies is realized, either the microbolometer arrays or, less likely, the pyroelectric arrays will capture the low-cost markets. It is not unreasonable that the uncooled arrays may obsolete “low-cost” PtSi, QWIP, and HIP FPAs, and render the intrinsic TE-cooled developments inconsequential. To improve hybrid reliability, alternative detector substrate materials including silicon along with alternative readout materials will become sufficiently mature to begin monolithic integration of the intrinsic materials with highest radiative performance. True optoelectronic FPAs consisting of IR sensors with optical output capability may be developed for reducing

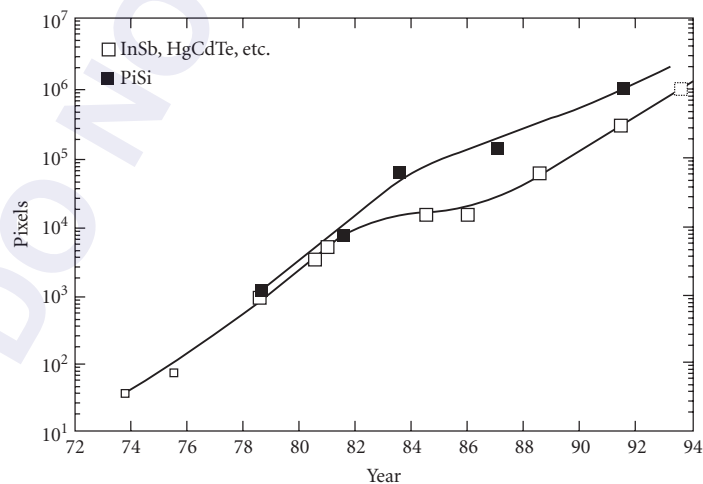


FIGURE 16 Chronological development of IR FPAs.

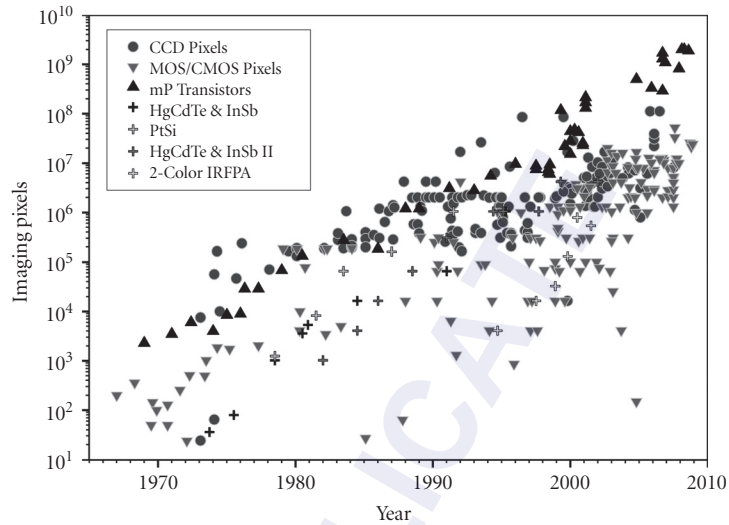


FIGURE 17 Integrated Circuit Chronology.

thermal loading and improving immunity to noise pickup. Availability of inexpensive, commercial devices is imminent along with the development of IR neural networks for additional signal processing capability.

In conclusion, it is likely instructive to compare the development of infrared sensors with both visible imaging sensors and other commercial integrated circuits, including semiconductor memory and microprocessors. All these products commonly share the benefits derived from Moore's law.⁵⁶ Fig. 17 is such a chronology showing notable devices of all types. Early infrared devices tended to develop at a pace only a few years behind the overlying trend known as Moore's law. Multicolor sensors initially developed at a rapid pace since the infrastructure was already in place. All infrared devices lag visible sensor development, which is driven by consumer demand. The largest visible sensors now boast about 100 million pixels; these foreshadow the size of upcoming infrared sensors once reliability issues and packaging costs are fully contained.

33.7 REFERENCES

1. W. P. McCracken, "CCDing in the Dark," *IEEE Spectrum*. (1992).
2. P. R. Norton, "Infrared Image Sensors," *Opt. Eng.* **30**:11 (1991).
3. A. Rogalski and J. Piotrowski, "Intrinsic Infrared Detectors," *Progress in Quantum Electronics* 12, Pergamon Press, 1988, pp. 2–3.
4. S. B. Stetson, D. B. Reynolds, M. G. Stapelbroek, and R. L. Stermer, "Design and Performance of Blocked-Impurity-Band Detector Focal Plane Arrays," *SPIE Proceedings*, vol. 686 (1986).
5. M. D. Petroff, G. Stapelbroek, and W. A. Kleinhans, "Solid State Photomultiplier," U.S. Patent No. 4,586,068 (1983).
6. B. F. Levine, K. K. Choi, C. G. Bethea, J. Walker, and R. J. Malik, *Appl. Phys. Lett.* **50**:1092 (1987).
7. C. S. Wu et al., "Novel GaAs/AlGaAs Multi-quantum-Well Schottky Junction Device and Its Photovoltaic LWIR Detection," *IEEE Trans. Electron Devices*, **ED-39**:2 (1992).
8. J. Cooper, *Rev. Sci. Instrum.* **33**:92 (1962).

9. C. Hanson, H. Beratan, R. Owen, M. Corbin, and S. McKenney, "Uncooled Thermal Imaging at Texas Instruments," *SPIE Proceedings*, vol. 1735, 1992.
10. D. H. Pommerrenig, "Extrinsic Silicon Focal Plane Array," *SPIE Proceedings*, vol. 443, 1983.
11. (a) W. F. Kosonocky, "Review of Infrared Image Sensors with Schottky-Barrier Detectors," *Optoelectronics—Devices and Technologies*, vol. 6, no. 2, December 1991, pp. 173–203. (b) W. F. Kosonocky, "State-of-the-Art in Schottky-Barrier IR Image Sensors," *SPIE Proceedings*, vol. 1681, Orlando, Fla., 1992.
12. Private communications with N. A. Foss from Sensors and Systems Center, Honeywell, Inc., Bloomington, Minn.
13. C. G. Roberts, "HgCdTe Charge Transfer Focal Planes," *SPIE Proceedings*, vol. 443, 1983.
14. M. D. Gibbons et al., "Status of IrSb Charge Injection Device (CID) Detection Technology," *SPIE Proceedings*, vol. 443, 1983.
15. M. H. White et al., "Characterization of Surface CCD Image Arrays at Low Light Levels," *IEEE J. Solid-State Circuits*, vol. SC-9, 1974, pp. 1–12.
16. D. J. Sauer, F. V. Shallcross, F. L. Hsueh, G. M. Meray, P. A. Levine, H. R. Gilmartin, T. S. Villani, B. J. Esposito, and J. R. Tower, "640 × 480 MOS PtSI IR Sensor," *SPIE Proceedings*, vol. 1540, 1991, pp. 285–296.
17. L. J. Kozlowski et al., "Comparison of Passive and Active Pixel Schemes for CMOS Visible Imagers," *SPIE* 3360, 1998.
18. W. F. Kosonocky, T. S. Villani, F. V. Shallcross, G. M. Meray, and J. J. O'Neill, III, "A Schottky-Barrier Image Sensor with 100% Fill Factor," *SPIE Proceedings*, vol. 1308, 1990, pp. 70–80.
19. M. Denda, M. Kimata, S. Iwade, N. Yutani, T. Kondo, and N. Tsubouchi, "4 × 4096-Element SWIR Multispectral Focal Plane Array," *SPIE Proceedings*, vol. 819–824, 1987, pp. 279–286.
20. K. Chow, J. P. Rode, D. H. Seib, and J. D. Blackwell, "Hybrid Infrared Focal Plane Arrays," *IEEE Trans. Electron Devices* **ED-29**(1) (January, 1982).
21. E. R. Gertner, W. E. Tennant, J. D. Blackwell, and J. P. Rode, "HgCdTe on Sapphire—A New Approach to Infrared Detector Arrays," *J. Cryst. Growth* **72**:465 (1985).
22. E. R. Gertner, *Ann. Rev. Mater. Sci.* **15**:303–328 (1985).
23. E. R. Fossum, "Infrared Readout Electronics," *SPIE Proceedings*, vol. 1684 (1992).
24. D. F. Barbe, "Imaging Devices Using the Charge-Coupled Concept," *Proc. IEEE* **63**(1) (1975).
25. P. Felix, M. Moulin, B. Munier, J. Portmann, and J.-P. Reboul, "CCD Readout of Infrared Hybrid Focal Plane Arrays," *IEEE Trans. Electron Devices* **ED-27**(1) (1980).
26. S. T. Baier, "Complementary Heterostructure (CHFET) Readout Technology for Infrared Focal Plane Arrays," *SPIE Proceedings*, vol. 1684 (1992).
27. R. Sahai, R. L. Pierson, R. J. Anderson, E. H. Martin, E. A. Sover, and J. Higgins, "GaAs CCD's with Transparent (ITO) Gates for Imaging and Optical Signal Processing," *IEEE Electron Device Lett.* **EDL-4** (1983).
28. L. J. Kozlowski and R. E. Kezer, "2 × 64 GaAs Readout for IR FPA Application," *SPIE Proceedings*, vol. 1684 (1992).
29. J. S. Bugler and P. G. A. Jespers, "Charge Pumping in MOS Devices," *IEEE Trans. on Electron Devices* **ED-16**:3 (1969).
30. M. F. Tompsett, "Surface Potential Equilibration Method of Setting Charge in Charge-Coupled Devices," *IEEE Trans. Electron Devices* **ED-22**:6 (1975).
31. L. J. Kozlowski, K. Vural, W. E. Tennant, R. E. Kezer, and W. E. Kleinhans, "10 × 132 CMOS/CCD Readout with 25 μm Pitch and On-Chip Signal Processing Including CDS and TDI," *SPIE Proceedings*, vol. 1684 (1992).
32. M. F. Tompsett, "The Quantitative Effects of Interface States on the Performance of Charge-Coupled Devices," *IEEE Trans. Electron Devices* **ED-20**:1 (1973).
33. P. W. Bosshart, "A Multiplexed Switched Capacitor Filter Bank," *IEEE Journal Solid-State Circuits*, **SC-15**:6 (1980).
34. W. S. Chan, "Detector-Charge-Coupled Device (CCD) Interface Methods," *SPIE Proceedings*, vol. 244 (1980).
35. R. M. Swanson and J. D. Meindl, "Ion-Implanted Complementary MOS Transistors in Low-Voltage Circuits," *IEEE Journal of Solid-State Circuits* **SC-7**:2 (1972).
36. R. Troutman and S. N. Chakravarti, "Subthreshold Characteristics of Insulated-Gate Field-Effect Transistors," *IEEE Trans. Circuit Theory* **CT-20**:6 (1973).

37. N. Bluzer and R. Stehlik, "Buffered Direct Injection of Photocurrents into Charge-Coupled Devices," *IEEE Journal of Solid-State Circuits* **SC-13**:1 (1978).
38. C. C. Enz, E. A. Vittoz, and F. Kruppenacher, "A CMOS Chopper Amplifier," *IEEE Journal of Solid-State Circuits* **SC-22**:3 (1987).
39. S. G. Chamberlain and J. P. Y. Lee, "A Novel Wide Dynamic Range Silicon Photodetector and Linear Imaging Array," *IEEE Trans. Electron Devices* **ED-31**:2 (1984).
40. F. Kruppenacher, E. Vittoz, and M. DeGrauwe, "Class AB CMOS Amplifier for Micropower SC Filters," *Electron. Lett.* **17**:13 (1981).
41. E. Vittoz and J. Fellrath, "CMOS Analog Integrated Circuits Based on Weak Inversion Operation," *IEEE Journal Solid State Circuits* **SC-12**:3 (1977).
42. R. D. Hudson, *Infrared System Engineering*, John Wiley and Sons, 1969.
43. J. M. Mooney et al., "Responsivity Nonuniformity Limited Performance of Infrared Staring Cameras," *Opt. Eng.* **28**:1151 (1989).
44. D. L. Shumaker, J. T. Wood, and C. R. Thacker, *FLIR Performance Handbook*, DCS Corporation, Alexandria, Va. (1988).
45. N. Yutani, M. Kimata, H. Yagi, J. Nakanishi, S. Nagayoshi, and N. Tsubouchi, "1040 × 1040 Element PtSi Schottky-Barrier IR Image Sensor," IEDM, Washington, D.C., 1991.
46. T. S. Villani, W. F. Kosonocky, F. V. Shallcross, J. V. Groppe, G. M. Meray, J. J. O'Neill, III, and B. J. Esposito, "Construction and Performance of a 320 × 244-Element IR-CCD Imager with PtSi SBDs," *SPIE Proceedings*, vol. 1107-01, 1989, pp. 9-21.
47. D. J. Sauer, F. V. Shallcross, F. L. Hsueh, G. M. Meray, P. A. Levine, H. R. Gilmartin, T. S. Villani, B. J. Esposito, and J. R. Tower, "640 × 480 MOS PtSi IR Sensor," *SPIE Proceedings*, vol. 1540, 1991, pp. 285-296.
48. K. Konuma, N. Teranishi, S. Tohyama, K. Masubuchi, S. Yamagata, T. Tanaka, E. Oda, Y. Moriyama, N. Takada, and N. Yoshioka, "324 × 487 Schottky-Barrier Infrared Imager," *IEEE Trans. Electron Devices*, **37**(3):629-635 (1990).
49. K. Konuma, S. Tohyama, A. Tanabe, K. Masubuchi, N. Teranishi, T. Saito, and T. Muramatsu, "A 648 × 487 Pixel Schottky-Barrier Infrared CCD Image Sensor," *1991 ISSCC Digest of Tech. Papers*, 1991, pp. 156-157.
50. H. Elabd, Y. Abedini, J. Kim, M. Shih, J. Chin, K. Shah, J. Chen, F. Nicol, W. Petro, J. Lehan, M. Duron, M. Manderson, H. Balopole, P. Coyle, P. Cheng, and W. Shieh, "488 × 512 and 244 × 256-Element Monolithic PtSi Schottky IR Focal Plane Arrays," *SPIE Proceedings*, vol. 1107-29, presented at *SPIE Aerospace Sensor Symposium*, Orlando, Fla., March 1989.
51. E. T. Nelson, K. Y. Wong, S. Yoshizumi, D. Rockafellow, W. DesJardin, M. Elzinga, J. P. Lavine, T. J. Tredwell, R. P. Khosla, P. Sorlie, B. Howe, S. Brickman, and S. Reformat, "Wide Field of View PtSi Infrared Focal Plane Array," *SPIE Proceedings*, vol. 1308, 1990, pp. 36-44.
52. M. Kimata, M. Denda, N. Yutani, S. Iwade, and N. Tsubouchi, "A 512 × 512-Element PtSi Schottky-Barrier Infrared Image Sensor," *IEEE J. Solid-State Circuits*, vol. SC-22, no. 6, 1987, pp. 1124-1129.
53. H. Yagi, N. Yutani, S. Nagayoshi, J. Nakanishi, M. Kimata, and N. Tsubouchi, "Improved 512 × 512 IR CSD with Large Fill Factor and Large Saturation Level," *SPIE Proceedings*, vol. 1685-04, 1992.
54. J. Edwards, J. Gates, H. Altin-Mees, W. Connelly, and A. Thompson, "244 × 400 Element Hybrid Platinum Silicide Schottky Focal Plane Array," *SPIE Proceedings*, vol. 1308, 1990, pp. 99-100.
55. J. L. Gates, W. G. Connelly, T. D. Franklin, R. E. Mills, F. W. Price, and T. Y. Wittwer, "488 × 640-Element Hybrid Platinum Silicide Schottky Focal Plane Array," *SPIE Proceedings*, vol. 1540, *Infrared Technology XVII*, 1991, pp. 262-273.
56. G. E. Moore, "Cramming More Components onto Integrated Circuits," *Electronics Magazine* **38**:8, 1965.

This page intentionally left blank.

DO NOT DUPLICATE

PART

7

RADIOMETRY
AND
PHOTOMETRY

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

RADIOMETRY AND PHOTOMETRY

Edward F. Zalewski

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

34.1 GLOSSARY

A	area
A_1, A_2, A_s, A_d	area of surface 1, surface 2, a source, a detector, respectively
A_r	area of an image on the retina of a human eye
A_p	area of the pupil of a human eye
A_{in}, A_{out}, A_{sph}	area of an input port, output port, and sphere surface, respectively
b	distance from optic axis
c	the speed of light in a vacuum
C_e	photon-to-electron conversion efficiency, i.e., quantum efficiency of a photodetector
D	diameter
dA	infinitesimal element of area
dA_1, dA_2, dA_s, dA_d	infinitesimal element of area of surface 1, surface 2, a source, a detector, respectively
dL_λ	infinitesimal change in radiance per wavelength interval
dT	infinitesimal change in temperature
$d\Phi_{12}$	infinitesimal amount of radiant power transferred from point 1 to point 2
$d\lambda$	infinitesimal wavelength interval
$d\nu$	infinitesimal frequency interval
$d\Omega$	infinitesimal change in solid angle
E	irradiance, the incident radiant power per the projected area of a surface
E_v	illuminance, the photometric equivalent of irradiance
E_r	average illuminance in an image on the retina of a human eye
E_T	retinal illuminance in units of trolands
$E_T(\lambda)$	photopic retinal illuminance from a monochromatic source in trolands

$E'_T(\lambda)$	scotopic retinal illuminance from a monochromatic source in trolands
$E_{r\lambda}$	retinal spectral irradiance in absolute units: $\text{W nm}^{-1}\text{m}^{-2}$
f	focal length
$f\#$	F -number
g	fraction of light lost through the input and output ports of an averaging sphere
h	Planck's constant
h_s, h_d	object (source) height, image (detector) height
I	radiant intensity, the emitted or reflected radiant power per solid angle
i	photoinduced current from a radiation detector
I_v	luminous intensity, the photometric equivalent of radiant intensity
k	Boltzmann's constant
K_m	luminous efficacy (i.e., lumen-to-watt conversion factor) for photopic vision
K'_m	luminous efficacy for scotopic vision
K_{ab}	nonlinearity correction factor for a photodetector
L	radiance, the radiant power per projected area and solid angle
L_{12}	radiance from point 1 into the direction of point 2
L_a, L_b	radiance in medium a , in medium b
L_e	radiance within the human eye
L_λ	radiance per wavelength interval
L_ν	radiance per frequency interval
L_v	luminance, the photometric equivalent of radiance
$L_v(\lambda)$	luminance of a monochromatic light source
M	exitance, the emitted or reflected radiant power per the projected area of a source
m	mean value
N	photon flux, the number of photons per second
n	index of refraction
n_a, n_b	index of refraction in medium a , in medium b
n_e	index of refraction of the ocular medium of the human eye
n_s, n_d	index of refraction in the object (i.e., source) region, in the image (i.e., detector) region
N_λ	photon flux per wavelength interval
N_ν	photon flux per frequency interval
$N_{E\lambda}$	photon flux irradiance on the retina of a human eye
Q	radiant energy
Q_λ	radiant energy per wavelength interval
Q_ν	radiant energy per frequency interval
R	responsivity of a photodetector, i.e., electrical signal out per radiant signal in
r	radius
r_s, r_d, r_{sph}	radius of a source, detector, sphere, respectively
$R(\lambda)$	spectral (i.e., per wavelength interval) responsivity of a photodetector
s	distance
s_{12}	length of the light ray between points 1 and 2
s_{sd}	length of the light ray between points on the source and detector
s_{pr}	distance from the pupil to the retina in a human eye

T	absolute temperature
t	time
U	photon dose, the total number of photons
$V(\lambda)$	spectral luminous efficiency function (i.e., peak normalized human visual spectral responsivity) for photopic vision
$V'(\lambda)$	spectral luminous efficiency function for scotopic vision
w	width
x_i	the i th sample in a set of measurements
α	absorptance, fraction of light absorbed
β_a, β_b	angle of incidence or refraction
γ	absorption coefficient of a solute
δ	angle of rotation between crossed polarizers
ε	emittance of a blackbody simulator
E	étendue
η	total number of sample measurements
θ_s, θ_d	angle between the light ray and the normal to a point on the surface of a source, of a detector
θ_1, θ_2	angle between the light ray and the normal to a surface at point 1, at point 2
κ	concentration of a solute
λ	wavelength
ν	frequency
ρ	fraction of light scattered or reflected
σ	standard deviation
σ_m	standard deviation of the mean
τ	transmittance, radiant signal out per radiant signal into a material
$\tau(\lambda)$	spectral (i.e., per wavelength interval) transmittance
$\tau_c(\lambda)$	spectral transmittance of the ocular medium of the human eye
Φ	radiant power or equivalently radiant flux
ϕ	half angle subtended by a cone
Φ_{in}, Φ_{out}	incoming radiant power, outgoing radiant power
Φ_r	luminous flux at the retina of the human eye
Φ_λ	radiant power per wavelength interval
Φ_ν	radiant power per frequency interval
Φ_v	photopic luminous flux, radiant power by photopic detectable human vision
Φ'_v	scotopic luminous flux, radiant power detectable by scotopic human vision
Ω	solid angle, a portion of the area on the surface a sphere per of the square of the sphere radius
Ω_a, Ω_b	solid angle in medium a , in medium b

34.2 INTRODUCTION

Radiometry is the measurement of the energy content of electromagnetic radiation fields and the determination of how this energy is transferred from a source, through a medium, and to a detector. The results of a radiometric measurement are usually obtained in units of power, i.e., in watts. However, the result may also be expressed as photon flux (photons per second) or in units of energy

(joules) or dose (photons). The measurement of the effect of the medium on the transfer of radiation, i.e., the absorption, reflection, or scatter, is usually called *spectrophotometry* and will not be covered here. Rather, the assumption is made here that the radiant power is transferred through a lossless medium.

Traditional radiometry assumes that the propagation of the radiation field can be treated using the laws of geometrical optics. That is, the radiant energy is assumed to be transported along the direction of a ray and interference or diffraction effects can be ignored. In those situations where interference or diffraction effects are significant, the flow of energy will be in directions other than along those of the geometrical rays. In such cases, the effect of interference or diffraction can often be treated as a correction to the result obtained using geometrical optics. This assumption is equivalent to assuming that the energy flow is via an incoherent radiation field. This assumption is widely applicable since most radiation sources are to a large degree incoherent. For a completely rigorous treatment of radiant energy flow, the degree of coherence of the radiation must be considered via a formalism based on the theory of electromagnetism as derived from Maxwell's equations.^{1,2} This complexity is not necessary for most of the problems encountered in radiometry.

In common practice, radiometry is divided according to regions of the spectrum in which different measurement techniques are used. Thus, vacuum ultraviolet radiometry, intermediate-infrared radiometry, far-infrared radiometry, and microwave radiometry are considered separate fields, and all are distinguished from radiometry in the visible and near-visible optical spectral region.

The reader should note that there is considerable confusion regarding the nomenclatures of the various radiometries. The terminology for radiometry that we have inherited is dictated not only by its historical origin,³ but also by that of related fields of study. By the late 1700s, techniques were developed to measure light using the human eye as a null detector in comparisons of sources. At about the same time, radiant heating effects were studied with liquid-in-glass thermometers and actinic (i.e., chemical) effects of solar radiation were studied by the photoinduced decomposition of silver compounds into metallic silver. The discovery of infrared radiation in 1800 and ultraviolet radiation in 1801 stimulated a great deal of effort to study the properties of these radiations. However, the only practical detectors of ultraviolet radiation at that time were the actinic effects—for infrared radiation it was thermometers and for visible radiation it was human vision. Thus actinometry, radiometry, and photometry became synonymous with studies in the ultraviolet, infrared, and visible spectral regions. Seemingly independent fields of study evolved and even today there is confusion because the experimental methods and terminology developed for one field are often inappropriately applied to another. Vestiges of the confusion over what constitutes photometry and radiometry are to be found in many places. The problems encountered are not simply semantic, since the confusion can often lead to substantial measurement error.

As science progressed, radiometry was in the mainstream of physics for a short time at the end of the nineteenth century, contributing the absolute measurement base that led to Planck's radiation law and the discovery of the quantum nature of radiation. During this period, actinic effects, which were difficult to quantify, became part of the emerging field of photochemistry. In spite of the impossibility of performing an absolute physical measurement using the human eye, it was photometry, however, that grew to dominate the terminology and technology of radiant energy measurement practice in this period. At the beginning of the nineteenth century, the reason photometry was dominant was that the most precise (not absolute) studies of radiation transfer relied on the human eye. By the end of the nineteenth century, the growth of industries such as electric lighting and photography became the economic stimulus for technological developments in radiation transfer metrology and supported the dominance of photometry. Precise photometric measurements using instrumentation in which the human eye was the detector continued into the last half of the twentieth century. The fact that among the seven internationally accepted base units of physical measurement there remains one unit related to human physiology—the candela—is an indication of the continuing economic importance of photometry.

Presently, the recommended practice is to limit the term photometry to the measurement of the ability of electromagnetic radiation to produce a visual sensation in a physically realizable manner, that is, via a defined simulation of human vision.⁴⁻⁵ Radiometry, on the other hand, is used to describe the measurement of radiant energy independent of its effect on a particular detector.

Actinometry is used to denote measurement of photon flux (photons per second) or dose (total number of photons) independent of the subsequent photophysical, photochemical, or photobiological process. Actinometry is a term that is not extensively used, but there are current examples where measurement of the “actinic effect of radiation” is an occasion to produce a new terminology for a specific photoprocess, such as for the Caucasian human skin reddening effect commonly known as sunburn. We do not attempt here to catalog the many different terminologies used in photometry and radiometry, instead the most generally useful definitions are introduced where appropriate.

This chapter begins with a discussion of the basic concepts of the geometry of radiation transfer and photon flux measurement. This is followed by several approximate methods for solving simple radiation transfer problems. Next is a discussion of radiometric calibrations and the methods whereby an absolute radiant power or photon flux measurement is obtained. The discussion of photometry that follows is restricted to measurements employing physical detectors rather than those involving a human observer. Because many esoteric terms are still in use to describe photometric measurements, the ones most likely to be encountered are listed and defined in the section on photometry.

It is not the intention that this chapter be a comprehensive listing or a review of the extensive literature on radiometry and photometry; only selected literature citations are made where appropriate. Rather, it is hoped that the reader will be sufficiently introduced to the conceptual basis of these fields to enable an understanding of other available material. There are many texts on general radiometry. Some of the recent books on radiometry are listed in the reference section.⁶⁻¹⁰ In addition, the subject of radiometry or photometry is often presented as a subset of another field of study and can therefore be found in a variety of texts. Several of these texts are also listed in the reference section.¹¹⁻¹⁴ Finally, the reader will also find material related to radiometry, photometry, colorimetry, and spectrophotometry in Chaps. 34 to 40 in this volume and Chap. 10, “Colorimetry” in Vol. III.

34.3 RADIOMETRIC DEFINITIONS AND BASIC CONCEPTS

Radiant Power and Energy

For a steadily emitting source, that is a radiation source with a continuous and stable output, radiometric measurement usually implies measurement of the power of the source. For a flashing or single-pulse source, radiometric measurement implies a measurement of the energy of the source.

Radiometric measurements are traditionally measurements of thermal power or energy. However, because of the quantum nature of most photophysical, photochemical, and photobiological effects, in many applications it is not the measurement of the thermal power in the radiation beam but measurement of the number of photons that would provide the most physically meaningful result. The fact that most radiometric measurements are in terms of watts and joules is due to the history of the field. The reader should examine the particular application to determine if a measurement in terms of photon dose or photon flux would not be more meaningful and provide insight for the interpretation of the experiment. (See section on “Actinometry” later in this chapter.)

Radiant Energy Radiant energy is the energy emitted, transferred, or received in the form of electromagnetic radiation.

Symbol: Q *Unit:* joule (J)

Radiant Power Radiant power or radiant flux is the power (energy per unit time t) emitted, transferred, or received in the form of electromagnetic radiation.

Symbol: Φ *Unit:* watt (W)

$$\Phi = \frac{dQ}{dt} \quad (1)$$

Geometrical Concepts

The generally accepted terminology and basic definitions for describing the geometry of radiation transfer are presented below. More extensive discussions of each of these definitions and concepts can be found in the references.⁴⁻¹⁴

The concepts of irradiance, intensity, and radiance involve the density of the radiant power (or energy) over area, solid angle, and area times solid angle, respectively.

In situations where the density or distribution of the radiation on a surface is the required quantity, then it is the irradiance that must be measured. An example of where an irradiance measurement would be required is the exposure of a photosensitive surface such as the photoresists used in integrated circuit manufacture. The irradiance distribution over the surface determines the local degree of exposure of the photoresist. A nonuniform irradiance distribution will result in overexposure and/or underexposure of regions across the piece and results in a defect in manufacture.

In an optical system where the amount of radiation transfer through the system is important, then it is the radiance that must be measured. The amount of radiation passing through the optical system is determined by the area of the source from which the radiation was emitted and the field of view of the optic, also known as the solid angle or collection angle. Radiance is often thought of as a property of a source, but the radiance at a detector is also a useful concept.

Both irradiance and radiance are defined for infinitesimal areas and solid angles. However, in practice, measurements are performed with finite area detectors and optics with finite fields of view. Therefore all measurements are in fact measurement of average irradiance and average radiance.

Irradiance and radiance must be defined over a projected area in order to account for the effect of area change with angle of incidence. This is easily seen from the observation that the amount of a viewed area diminishes as it is tilted with respect to the viewer. Specifically, the view of the area falls off as the cosine of the angle between the normal to the surface and the line of sight. This effect is sometimes called the *cosine law of emission* or the *cosine law of irradiation*.

Intensity is a term that is part of our common language and often a point of confusion in radiometry. Strictly speaking, intensity is definable only for a source that is a point. An average intensity is not a measurable quantity since the source must by definition be an infinitesimal point. All intensity measurements are an approximation, since a true point source is physically impossible to produce. It is an extrapolation of a series of measurements that is the approximation of the intensity. An accurate intensity measurement is one that is made at a very large distance and, consequently, with a very small signal at the detector and an unfavorable signal-to-noise ratio. Historically speaking, however, intensity is an important concept in photometry and, to a much lesser extent, it has some application in radiometry. Intensity is a property of a source, not a detector.

Irradiance Irradiance is the ratio of the radiant power incident on an infinitesimal element of a surface to the projected area of that element, dA_d , whose normal is at an angle θ_d to the direction of the radiation.

Symbol: E Unit: watt/meter² (W m⁻²)

$$E = \frac{d\Phi}{\cos\theta_d dA_d} \quad (2)$$

Exitance The accepted convention makes a distinction between the irradiance, the surface density of the radiation incident on a radiation detector (denoted by the subscript d), and the exitance, the surface density of the radiation leaving the surface of a radiation source (denoted by the subscript s).

Exitance is the ratio of the radiant power leaving an infinitesimal element of a source to the projected area of that element of area dA_s , whose normal is at an angle θ_s to the direction of the radiation.

Symbol: M Unit: watt/meter² (W m⁻²)

$$M = \frac{d\Phi}{\cos\theta_s dA_s} \quad (3)$$

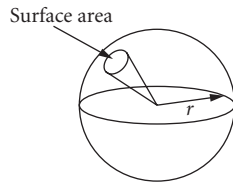


FIGURE 1 The solid angle at the center of the sphere is the surface area enclosed in the base of the cone divided by the square of the sphere radius.

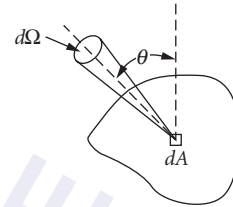


FIGURE 2 The radiance at the infinitesimal area dA is the radiant flux divided by the solid angle times the projection of the area dA onto the direction of the flux.

Intensity Radiant intensity (often simply “intensity”) is the ratio of the radiant power leaving a source to an element of solid angle $d\Omega$ propagated in the given direction.
Symbol: I Unit: watt/steradian (W sr^{-1})

$$I = \frac{d\Phi}{d\Omega} \quad (4)$$

Note that in the field of physical optics, the word *intensity* refers to the magnitude of the Poynting vector and thus more closely corresponds to irradiance in radiometric nomenclature.

Solid Angle The solid angle is the ratio of the area on the surface of a sphere to the square of the radius r of the sphere. This is illustrated in Fig. 1.
Symbol: Ω Unit: steradian (sr)

$$d\Omega = \frac{dA}{r^2} \quad (5)$$

It follows from the definition that the solid angle subtended by a cone of half angle ϕ , the apex of which is at the center of the sphere, is given by

$$\Omega = 2\pi(1 - \cos\phi) = 4\pi \sin^2 \frac{\phi}{2} \quad (6)$$

Radiance Radiance, shown in Fig. 2, is the ratio of the radiant power, at an angle θ_s to the normal of the surface element, to the infinitesimal elements of both projected area and solid angle. Radiance can be defined either at a point on the surface of either a source or a detector, or at any point on the path of a ray of radiation.

Symbol: L Unit: watt/steradian meter² ($\text{W sr}^{-1}\text{m}^{-2}$)

$$L = \frac{d\Phi}{\cos\theta_s dA_s d\Omega} \quad (7)$$

Radiance plays a special role in radiometry because it is the propagation of the radiance that is conserved in a lossless optical system; see “Radiance Conservation Theorem, Homogeneous Medium.” Radiance was often referred to as the brightness or the specific intensity, but this terminology is no longer recommended.

Spectral Dependence of Radiometric Quantities

Polychromatic Radiation Definitions For polychromatic radiation, the spectral distribution of radiant power (or radiant energy) is denoted as either radiant power (energy) per wavelength interval or radiant power (energy) per frequency interval.

Symbol: $\Phi_\lambda(Q_\lambda)$ Unit: watt/nanometer (W nm^{-1}); joule/nanometer (J nm^{-1});
or

Symbol: $\Phi_\nu(Q_\nu)$ Unit: watt/hertz (W Hz^{-1}); joule/hertz (J Hz^{-1})

It follows that $\Phi_\lambda d\lambda$ is the radiant power in the wavelength interval λ to $\lambda + d\lambda$, and $\Phi_\nu d\nu$ is the radiant power in the frequency interval ν to $\nu + d\nu$. The total radiant power over the entire spectrum is therefore

$$\Phi = \int_0^\infty \Phi_\lambda d\lambda \quad (8a)$$

or

$$\Phi = \int_0^\infty \Phi_\nu d\nu \quad (8b)$$

If λ is the wavelength in the medium corresponding to the frequency ν , and since $\nu = c/n\lambda$, where c is the speed of light in a vacuum and n is the index of refraction of the medium, then

$$d\nu = -\frac{c}{n\lambda^2} d\lambda \quad (9)$$

and

$$\lambda\Phi_\lambda = \nu\Phi_\nu \quad (10)$$

Since the wavelength changes with the index of refraction of the medium, it is becoming more common to use the vacuum wavelength, $\lambda = c/\nu$. It is particularly important in high-accuracy applications to state explicitly whether or not the vacuum wavelength is being used.

Spectral versions of the other radiometric quantities, i.e., radiant energy, radiance, etc., are defined similarly.

Polychromatic Radiation Calculations As an example of the application of the concept of the spectral dependence of a radiometric quantity, consider the calculation of the response of a radiometer consisting of a detector and a spectral filter. The spectral responsivity of a detector $R(\lambda)$ is the ratio of the output signal to the radiant input at each wavelength λ . The output is usually an electrical signal, such as a photocurrent i , and the input is a radiometric quantity, such as radiant power. The spectral transmittance of a filter $\tau(\lambda)$ is the ratio of the output radiant quantity to the input radiant quantity at each wavelength λ . For a spectral radiant power Φ_λ , the photocurrent i of the radiometer is

$$i = \int_0^\infty R(\lambda)\tau(\lambda)\Phi_\lambda d\lambda \quad (11)$$

In practice, either the responsivity of the detector or the transmittance of the filter are nonzero only within a limited spectral range. The integral need be evaluated only within the wavelength limits where the integrand is nonzero.

Photometry

The radiation transfer concepts, i.e., geometrical principles, of photometry are the same as those for radiometry. The exception is that the spectral responsivity of the detector, the human eye, is specifically defined. Photometric quantities are related to radiometric quantities via the spectral efficiency functions defined for the photopic and scotopic CIE Standard Observer. The generally accepted values of the photopic and scotopic human eye response function are represented in the "Photometry" section in Table 2.

Luminous Flux The photometric equivalent of radiant power is luminous flux, and the unit that is equivalent to the watt is the lumen. Luminous flux is spectral radiant flux weighted by the appropriate eye response function. The definition of luminous flux for the photopic CIE Standard Observer is

Symbol: Φ_v Unit: lumen (lm)

$$\Phi_v = K_m \int \Phi_\lambda V(\lambda) d\lambda \quad (12)$$

where $V(\lambda)$ is the spectral luminous efficiency function and K_m is the luminous efficacy for photopic vision. The spectral luminous efficacy is defined near the maximum, $\lambda_m = 555$ nm, of the photopic efficiency function to be approximately 683 lm W^{-1} .

Definitions of the Density of Luminous Flux

Illuminance Illuminance is the photometric equivalent of irradiance; that is, illuminance is the luminous flux per unit area.

Symbol: E_v Unit: lumen/meter² (1 m m^{-2})

$$E_v = \frac{d\Phi_v}{\cos\theta_d dA_d} = \frac{d[K_m \int \Phi_\lambda V(\lambda) d\lambda]}{\cos\theta_d dA_d} \quad (13)$$

Luminous intensity Luminous intensity is the photometric equivalent of radiant intensity. Luminous intensity is the luminous flux per solid angle. For historical reasons, the unit of luminous intensity, the candela—not the lumen—is defined as the base unit for photometry. However, the units for luminous intensity can either be presented as candelas or lumens/steradian.

Symbol: I_v Unit: candela or lumen/steradian (cd or lm sr^{-1})

$$I_v = \frac{d\Phi_v}{d\Omega} = \frac{d[K_m \int \Phi_\lambda V(\lambda) d\lambda]}{d\Omega} \quad (14)$$

Luminance Luminance is the photometric equivalent of radiance. Luminance is the luminous flux per unit area per unit solid angle.

Symbol: L_v Unit: candela/meter² (cdm^{-2})

$$L_v = \frac{d\Phi_v}{\cos\theta_s dA_s d\Omega} = \frac{d[K_m \int \Phi_\lambda V(\lambda) d\lambda]}{\cos\theta_s dA_s d\Omega} \quad (15)$$

Actinometry

Radiant Flux to Photon Flux Conversion Actinometric measurement practice closely follows that of general radiometry except that the quantum nature of light rather than its thermal effect is emphasized. In actinometry, the amount of electromagnetic radiation being transferred is measured in units of photons per second (photon flux). The energy of a single photon is

$$Q = h\nu \quad (16)$$

where ν is the frequency of the radiation and h is Planck's constant, $6.6261 \times 10^{-34} \text{ J s}$. For monochromatic radiant power Φ_λ , measured as watts and wavelength λ , measured as nanometers, the number of photons per second N_λ in the monochromatic radiant beam is

$$N_\lambda = 5.0341 \times 10^{15} n \lambda \Phi_\lambda \quad (17)$$

Photon Dose and the Einstein Dose is the total number of photons impinging on a sample. For a monochromatic beam of radiant power Φ_λ that irradiates a sample for a time t seconds, the dose U measured as Einsteins is

$$U = 8.3593 \times 10^{-9} n \lambda \Phi_\lambda t \quad (18)$$

The Einstein is a unit of energy used in photochemistry. An Einstein is the amount of energy in one mole (Avogadro's number, 6.0221×10^{23}) of photons.

TABLE 1 Radiation Transfer Terminology, Spectral Relationships

	Radiometric	Photometric	Actinometric
Base quantity:	Radiant power (also radiant flux)	Luminous flux	Photon flux
Units:	Watts/nanometer	Lumens	Photons/second
Conversion:	—	[W/nm] $K_m V(\lambda)$	[W/nm] $\lambda (hc)^{-1}$
Surface density:	Irradiance	Illuminance	Photon flux irradiance
Solid angle density:	Radiant intensity	Luminous intensity	Photon flux intensity
Solid angle and surface density:	Radiance	Luminance	Photon flux radiance

Conversions between Radiometry, Photometry, and Actinometry

Conversions between radiometric, photometric, and actinometric units is not simply one of determining the correct multiplicative constant to apply. As seen previously, the conversion between radiant power and photon flux requires that the spectral character of the radiation be known. It was also shown that, for radiometric to photometric conversions, the spectral distribution of the radiation must be known. Furthermore, there is an added complication for photometry where one must also specify the radiant power level in order to determine which CIE Standard Observer function is appropriate. Table 1, which summarizes the spectral radiation transfer terminology, may be helpful to guide the reader in determining the relationship between radiometric, photometric, and actinometric concepts. In Table 1, the power level is assumed to be high enough to restrict the photometric measurements to the range of the photopic eye response function.

Basic Concepts of Radiant Power Transfer

Radiance Conservation Theorem, Homogeneous Medium In a lossless, homogeneous isotropic medium, for a perfect optical system (i.e., having no aberrations) and ignoring interference and diffraction effects, the radiance is conserved along a ray through the optical system. In other words, the spectral radiance at the image always equals the spectral radiance at the source.

It follows from Eq. (7), the definition of radiance, that for a surface A_1 with radiance L_{12} in the direction of a second surface A_2 with radiance L_{21} in the direction to a first surface, and joined by a light ray of length s_{12} , the net radiant power exchange between elemental areas on each surface is given by

$$\Delta\Phi = d\Phi_{12} - d\Phi_{21} = \frac{(L_{12} - L_{21}) \cos\theta_1 \cos\theta_2 dA_1 dA_2}{s_{12}^2} \quad (19)$$

where θ_1 and θ_2 are the angles between the ray s_{12} and the normals to the surfaces A_1 and A_2 , respectively. The transfer of radiant power and the terminology used in this discussion is depicted in Fig. 3.

The total amount of radiation transferred between the two surfaces is given by the integral over both areas as follows:

$$\Phi = \iint \frac{(L_{12} - L_{21}) \cos\theta_1 \cos\theta_2}{s_{12}^2} dA_1 dA_2 \quad (20)$$

This is the generalized radiant power transfer equation for net exchange between two sources. In the specialized case of a source and receiver, the radiant power emitted by a receiver is zero by definition. In this case, the term L_{21} in Eq. (20) is zero.

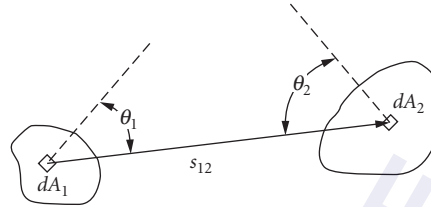


FIGURE 3 The radiant flux transferred between the infinitesimal areas dA_1 to dA_2 .

Refractive Index Changes In the case of a boundary between two homogeneous isotropic media having indices of refraction n_a and n_b , the angles of incidence and refraction at the interface β_a and β_b are related by Snell's law. If the direction of the light is oblique to the boundary between n_a and n_b , the solid angle change at the boundary will be

$$d\Omega_a = \frac{n_b^2 \cos \beta_b}{n_a^2 \cos \beta_a} d\Omega_b \quad (21)$$

Therefore the radiance change across the boundary will be

$$\frac{L_a}{n_a^2} = \frac{L_b}{n_b^2} \quad (22)$$

This result is obtained directly by substituting the optical path for the distance in Eq. (19) and considering that the radiance transferred across the boundary between the two media is unchanged. Optical path is the distance within the medium times the index of refraction of the medium.

In the case of an optical system having two or more indices of refraction, the radiance conservation theorem is more precisely stated as: In a lossless, homogeneous isotropic medium, for a perfect optical system (i.e., having no aberrations) and ignoring interference and diffraction effects, at a boundary between two media having different indices of refraction the radiance divided by the square of the refractive index is conserved along a ray through the optical system.

Radiative Transfer through Absorbing Media For radiation transmitted through an absorbing and/or scattering medium, the radiance is not conserved. This is not only because of the loss due to the absorption and/or scattering but the medium could also emit radiation. The emitted light will be due to thermal emission (see the discussion on blackbody radiation later in this chapter). In some cases, the medium may also be fluorescent. Fluorescence is the absorption of radiant energy at one wavelength with subsequent emission at a different wavelength.

Historically, the study of radiative transfer through absorbing and/or scattering media dealt with the properties of stellar atmospheres. Presently, there is considerable interest in radiative transfer measurements of the earth and its atmosphere using instruments on board satellites or aircraft. An accurate measure of the amount of reflected sunlight (approximately 400 to 2500 nm) or the thermally emitted infrared (wavelengths >2500 nm) requires correction for the absorption, scattering, and, in the infrared, the emission of radiation by the atmosphere. This specialized topic will not be considered here. Detailed discussion is available in the references.¹⁵⁻¹⁷

34.4 RADIANT TRANSFER APPROXIMATIONS

The solution to the generalized radiant power transfer equation is typically quite complex. However, there are several useful approximations that in some instances can be employed to obtain an estimate of the solution of Eq. (20). We shall consider the simpler case of a source and a detector rather

then the net radiant power exchange between two sources, since this is the situation commonly encountered in an optical system. In this case, Eq. (20) becomes

$$\Phi = \iint \frac{L \cos \theta_s \cos \theta_d}{s_{sd}^2} dA_s dA_d \quad (23)$$

where the subscripts s and d denote the source and detector, respectively. Here it is assumed the detector behaves as if it were a simple aperture. That is, it responds equally to radiation at any point across its surface and from any direction. Such a detector is often referred to as a cosine corrected detector. Of course, deviations from ideal detection behavior within the spatial and angular range of the calculation reduces the accuracy of the calculation.

Point-to-point Approximation: Inverse Square Law

The simplest approximations are obtained by assuming radiant flux transfer between a point source emitting uniformly in all directions and a point detector. The inverse square law is an approximation that follows directly from the definitions of intensity, solid angle, and irradiance, Eqs. (2), (4), and (5), respectively. The irradiance (at an infinitesimal area whose normal is along the direction of the light ray) times the square of the distance from a point source equals the intensity of the source

$$I = \frac{\Phi}{A} s^2 = E s^2 \quad (24)$$

The relationship between the uniformly emitted radiance and the intensity of a point source is obtained similarly from Eqs. (4) and (7):

$$L = \frac{I}{A_s} \quad (25)$$

These point-to-point relationships are perhaps most important as a test of the accuracy of a radiation transfer calculation at the limit as the areas approach zero.

Lambertian Approximation: Uniformly Radiant Areas

Lambertian Sources A very useful concept for the approximation of radiant power transfer is that of a source having a radiance that is uniform across its surface and uniformly emits in all directions from its surface. Such a uniform source is commonly referred to as a lambertian source.

For the case of a lambertian source, Eq. (23) becomes

$$\Phi = L \iint \frac{\cos \theta_s \cos \theta_d}{s_{sd}^2} dA_s dA_d \quad (26)$$

Configuration factor The double integral in Eq. (26) has been given a number of different names: configuration factor, radiation interaction factor, and projected solid angle. There is no generally accepted terminology for this concept, although configuration factor appears most frequently. Analytical solutions to the double integral have been found for a variety of different shapes of source and receiver. Tabulations of these exact solutions to the integral in Eq. (26) are usually found in texts on thermal engineering,^{18,19} under the heading of radiant heat transfer or configuration factor.

Radiation transfer between complex shapes can often be determined by using various combinations of configuration factors. This technique is often referred to as configuration factor algebra.¹⁸ The surfaces are treated as pieces, each with a calculable configuration factor, and the separate configuration factors are combined to obtain the effective configuration factor for the complete surface.

Étendue The double integral in Eq. (26) is often used as a means to characterize the flux-transmitting capability of an optical system in a way that is taken to be independent of the radiant properties of the source. Here the double integral is written as being over area and solid angle:

$$\Phi = L \iint \cos \theta_d dA_s d\Omega \quad (27)$$

In this case, the surface of the lambertian source is assumed perpendicular to the optic axis and to lie in the entrance window of the optical system. The solid angle is measured from a point on the source to the entrance pupil. The étendue E of an optical system of refractive index n is defined as

$$E = n^2 \iint \cos \theta_d dA_s d\Omega \quad (28)$$

Equation (28) is sometimes referred to as the throughput of an optical system.

Total flux into a hemisphere The total amount of radiation emitted from a lambertian source of area dA_s into the hemisphere centered at dA_s (or received by a hemispherical, uniform detector centered at dA_s) is obtained from integrating Eq. (26) over the area A_d . Note that the ray s_{sd} is everywhere normal to the surface of the hemisphere; i.e., $\cos \theta_d = 1$.

$$\Phi = L\pi \int dA_s \quad (29)$$

Using Eq. (3), the definition of the exitance, the radiance at each point on the surface of the source is

$$L = \frac{M}{\pi} \quad (30)$$

Because of the relationship expressed in Eq. (30), Eq. (26) is often written in terms of the exitance.

$$\Phi = M\pi \iint \frac{\cos \theta_s \cos \theta_d}{s_{sd}^2} dA_s dA_d \quad (31)$$

In this case, the factor π is considered to be part of the configuration factor. Note again that there is no generally accepted definition of the configuration factor.

Radiation transfer between a circular source and detector The particular case of radiation transfer between circular apertures, the centers of which are located along the same optical axis as shown in Fig. 4, is a configuration common to many optical systems and is therefore illustrated here. The

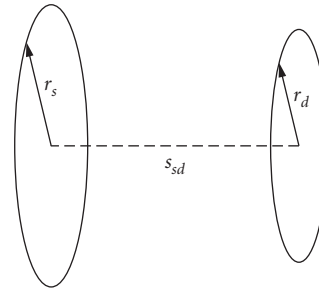


FIGURE 4 Radiant flux transfer between two circular apertures normal and concentric to the axis joining them.

radius of the source (or first aperture) is r_s , the detector (second aperture) radius is r_d , and the distance between the centers is s_{sd} . The exact solution of the integral in Eq. (26) yields

$$\Phi = \frac{2L(\pi r_s r_d)^2}{r_s^2 + r_d^2 + s_{sd}^2 + [(r_s^2 + r_d^2 + s_{sd}^2)^2 - 4r_s^2 r_d^2]^{1/2}} \quad (32)$$

This result can be approximated for the case where the sum of the squares of the distance and radii is large compared to the product of the radii, that is, $(r_s^2 + r_d^2 + s_{sd}^2) \gg 2r_s r_d$, so that Eq. (32) reduces to

$$\Phi \cong \frac{L(\pi r_s r_d)^2}{r_s^2 + r_d^2 + s_{sd}^2} \quad (33)$$

From this expression the irradiance at the detector can be obtained

$$E = \frac{\Phi}{A_d} \cong \frac{LA_s}{r_s^2 + r_d^2 + s_{sd}^2} \cong \frac{LA_s}{s_{sd}^2} \quad (34)$$

where A_s is the area of the lambertian disk and A_d is the detector area. The approximation at the extreme right is obtained by assuming that the radii are completely negligible with respect to the distance. This is the same result that would be obtained from a point-to-point approximation.

Off-axis irradiance: cosine-to-the-fourth approximation Equation (34) describes the irradiance from a small lambertian disk to a detector on the ray axis and where both surfaces are perpendicular to the ray. If the detector is moved off-axis by a distance b as depicted in Fig. 5, the ray from A_s to A_d will then be at an angle with respect to the normal at both surfaces as follows

$$\theta_s = \theta_d = \theta = \tan^{-1} \left(\frac{b}{s_{sd}} \right) \quad (35)$$

The projected areas are then $(A_s \cos \theta)$ and $(A_d \cos \theta)$. In addition, the distance from the source to the detector increases by the factor $(1/\cos \theta)$. The radiant power at a distance b away from the axis therefore decreases by the fourth power of the cosine of the angle formed between the normal to the surface and the ray.

$$\Phi \cong \frac{LA_s A_d}{s_{sd}^2} \cos^4 \theta \quad (36)$$

Since the radiance is conserved for propagation in a lossless optical system, Eq. (36) also approximates the radiant power from an off-axis region of a large lambertian source received at a small detector. The approximate total radiant power received at the detector would then be the sum of the radiant power contributed by each region of the source.

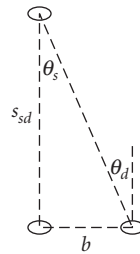


FIGURE 5 illustration of the cosine-fourth effect on irradiance, displacement of the receiving surface by a distance b .

Spherical lambertian source In order to compute the radiant power at a point at a distance s_{sd} from the center of a spherical lambertian source of radius r_{sph} , it is not necessary to explicitly solve the integrals in the radiation transfer equation. The solution is readily obtained from the symmetry of the lambertian sphere. Using the relationship between the exitance and radiance of a lambertian source [Eq. (30)], the total radiation power emitted by the source is obtained from the product of the surface area of the source times the exitance.

$$\Phi = 4\pi^2 r_{sph}^2 L \quad (37)$$

The radiant power is isotropically emitted. Therefore, the irradiance at any point on an enclosing sphere of radius s_{sd} is the total radiant power divided by the area of the enclosing sphere.

$$E = \frac{\pi r_{sph}^2 L}{s_{sd}^2} \quad (38)$$

Note that the irradiance from a spherical lambertian source follows the inverse square law at all distances from the surface of the sphere. The intensity of a spherical lambertian source is

$$I = \pi r_{sph}^2 L \quad (39)$$

Radiant Flux Transfer through a Lambertian Reflecting Sphere A lambertian reflector is a surface that uniformly scatters a fraction ρ of the radiation incident upon it.

$$L = \frac{\rho E}{\pi} \quad (40)$$

where E is the irradiance.

A spherical enclosure whose interior is coated with a material that approximates a lambertian reflector is a widely used tool in radiometry and photometry.²⁰ Such spheres are used either for averaging a nonuniform radiant power distribution (averaging sphere) or for measuring the total amount of radiant power emitted from a source (integrating sphere).

The sphere has the useful property whereby the solid angle subtended by any one section of the wall times the projected area is constant over all other points on the inside surface of the sphere. Therefore, if radiation falling on any point within the sphere is uniformly reflected, the reflected radiation will be uniformly distributed, i.e., produce uniform irradiance, throughout the interior. This result follows directly from the symmetry of the sphere.

Consider a sphere of radius r_{sph} and the radiant power transfer between two points on the inner surface. The normals to the two points are radii of the sphere and form an isosceles triangle when taken with the ray joining the points. Therefore, the angles between the ray and the normals to each point are equal. From Eq. (26)

$$\Phi = L \iint \frac{\cos^2 \theta}{s_{sd}^2} dA_s dA_d \quad (41)$$

The length of the ray joining the points is $2r_{sph} \cos \theta$. The irradiance is therefore

$$E = \frac{\Phi}{A_d} = \frac{L A_s}{4r_{sph}^2} \quad (42)$$

which is independent of the angle θ . If Φ_{in} is the radiant power entering the sphere, the irradiance at any point on the sphere after a single reflection will be

$$E = \frac{\rho \Phi_{in}}{4\pi r_{sph}^2} \quad (43)$$

A fraction ρ of the flux will be reflected and again uniformly distributed over the sphere. After multiple reflections the irradiance at any point on the wall of the sphere is

$$E = \frac{(\rho + \rho^2 + \rho^3 + \dots)\Phi_{\text{in}}}{4\pi r_{\text{sph}}^2} = \frac{\rho\Phi_{\text{in}}}{(1-\rho)A_{\text{sph}}} \quad (44)$$

where A_{sph} is the surface area of the sphere. The flux Φ_{out} exiting the sphere through a port of area A_{out} is

$$\Phi_{\text{out}} = \frac{\rho\Phi_{\text{in}}A_{\text{out}}}{(1-\rho)A_{\text{sph}}} \quad (45)$$

In Eq. (45) it is assumed that the loss of radiation at the entrance and exit ports is negligible and does not affect the symmetry of the radiation distribution.

The effect of the radiation lost through the entrance and exit ports is approximated as follows. After the first reflection, the fraction of radiation lost in each subsequent reflection is equal to the combined areas of the ports divided by the sphere area. Therefore the fraction reflected within the sphere is

$$g = 1 - \frac{A_{\text{in}} + A_{\text{out}}}{A_{\text{sph}}} \quad (46)$$

Using this in Eq. (44) yields

$$\Phi_{\text{out}} = \frac{\rho\Phi_{\text{in}}A_{\text{out}}}{(1-\rho g)A_{\text{sph}}} \quad (47)$$

Since the sphere is approximately a lambertian source, the radiance at the exit port is

$$L = \frac{\rho\Phi_{\text{in}}}{(1-\rho g)\pi A_{\text{sph}}} \quad (48)$$

Radiometric Effect of Stops and Vignetting

Refer to Fig. 6 for an illustration of these definitions. The *aperture stop* of an optical system is an aperture near the entrance to the optical system that determines the size of the bundle of rays leaving the source that can enter the optical system.

The *field stop* is an aperture within the optical system that determines the maximum angle of the rays that pass through the aperture stop that can reach the detector. The position and area of the field stop determines the field of view of the optical system. The field stop limits the extent of the source that is represented in its image at the detector.

The image of the aperture stop in object space, i.e., in the region of the source, is the *entrance pupil*. The image of the aperture stop in image space, i.e., in the region of the detector, is the *exit pupil*. Light rays that pass through the center of the aperture stop also pass through the centers of the images of the aperture stop at the entrance and exit pupils. Since all of the light entering the optical system must pass through the aperture stop, all of the light reaching the detector appears to pass through the exit pupil.

The field stop defines the solid angle within the optical system, the system field of view. When viewed from the image, the field stop of an optical system takes on the radiance of the object being imaged. This is a useful radiometric concept since a complex optical system can often be approximated as an exit pupil having the same radiance as the object being imaged (modified by the system transmission losses). The direction in which the radiation in the image appears to be emitted is,

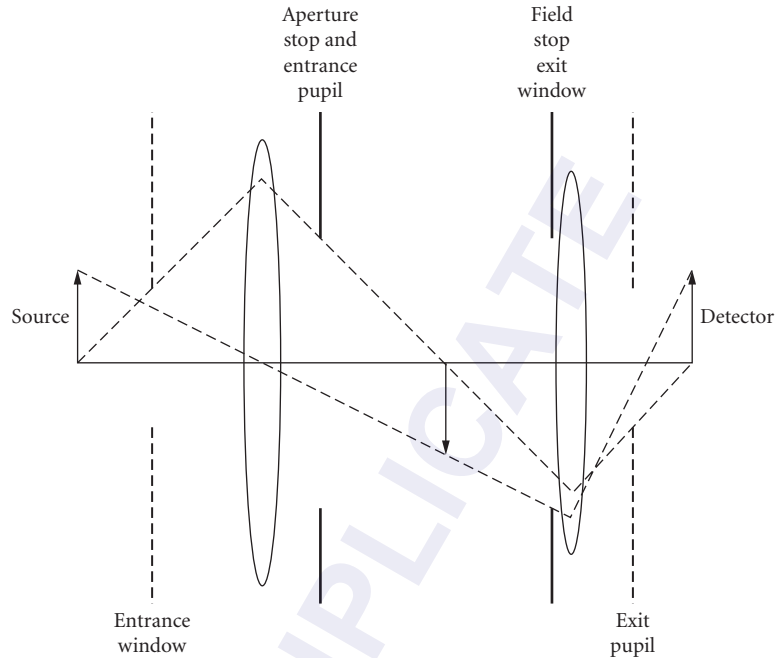


FIGURE 6 Schematic showing the relative positions of the stops, pupils, and windows in a simple optical system.

of course, limited by the aperture at the field stop. A word of caution: if the object is small, its image will be limited by diffraction effects and its radiance will depart even further from the extended source (large area) approximations used here.

The *entrance window* is the image of the field stop at the source and the *exit window* is the image of the field stop at the detector. If the field stop coincides with the detector, i.e., the detector is in the image plane of the optical system, then the entrance window will correspond with the object plane on the source. If the field stop does not coincide with the image plane at the detector, then because of parallax, different portions of the source will be visible from different points within the exit pupil. This condition, known as *vignetting*, causes a decrease in the irradiance at the off-axis points on the detector or image plane.

Approximate Radiance at an Image

Aplanatic Optical Systems Except for rays that lie on the optic axis, the radiance of an image must be based on a knowledge of the image quality since any aberrations introduced by the optical system divert some of the off-axis rays away from the image.

Consider a well-corrected optical system that is assumed to be aplanatic for the source and image points. That is, the optical system obeys Abbe's sine condition which is

$$n_s h_s \sin \theta_s = n_d h_d \sin \theta_d \quad (49)$$

where n_s and n_d are the refractive indices of the object (source) and image (detector) spaces, h_s and h_d are the object and image heights, and θ_s and θ_d are the angles between the off-axis rays and the

optic axis in object and image space. From Eq. (27) the flux radiated by a small lambertian source of area A_s into the solid angle of the optical system is

$$\Phi \cong 2\pi L A_s \int_0^{\theta_s} \cos\theta \sin\theta d\theta = \pi L A_s \sin^2 \theta_s \quad (50)$$

The differential of the solid angle is obtained from Eq. (6). Since Φ is the radiant flux at the image and A_d is the area of the image, the irradiance at the image is

$$E = \frac{\Phi}{A_d} = \frac{\pi L A_s}{A_d} \sin^2 \theta_s \quad (51)$$

If h_s and h_d are the radii of circular elements A_s and A_d , then according to Abbe's sine condition

$$\frac{A_s}{A_d} = \frac{n_d^2 \sin^2 \theta_d}{n_s^2 \sin^2 \theta_s} \quad (52)$$

The irradiance at the image is

$$E = \frac{\pi L n_d^2}{n_s^2} \sin^2 \theta_d \quad (53)$$

Numerical Aperture and F-number The quantity $n_d \sin \theta_d$ in Eq. (53) is called the *numerical aperture* of the imaging system. The irradiance of the image is proportional to the square of the numerical aperture. Geometrically speaking, the image irradiance increases with the angle of the cone of light converging on the image.

Another approximate measure of the image irradiance of an optical system is the *F-number*, $f\#$ (sometimes called the focal ratio) defined by the ratio of focal length (in the image space) f to the diameter D of the entrance pupil. For a source at a very large distance

$$f\# = \frac{f}{D} = \frac{1}{2 \tan \theta_d} \cong \frac{1}{2 \sin \theta_d} \quad (54)$$

The approximate image irradiance expressed in terms of the *F-number* is

$$E \cong \pi L \left(\frac{n_d}{2n_s f\#} \right)^2 \quad (55)$$

34.5 ABSOLUTE MEASUREMENTS

An absolute measurement, often referred to as an absolute calibration, is a measurement that is based upon, i.e., derived from, one of the internationally recognized units of physical measurement. These units are known as the SI units (Système International d'Unités²¹). The absolute SI base units are the meter, second, kilogram, kelvin, ampere, candela, and mole. The definitions of the SI units, the methods for their realization, or their physical embodiment are a matter of international agreement under the terms of the 1875 Treaty of the Meter. A convenient method (but often not a sufficient condition) for achieving absolute accuracy is to obtain traceability to one of the SI units via a calibration transfer standard issued by one of the national standards laboratories. The United States standards laboratory is the National Institute of Standards and Technology (NIST, formerly the National Bureau of Standards).

A relative measurement is one that is not required to be traceable to one of the SI units. Relative measurements are usually obtained as the ratio of two measurements. An example of a relative

measurement is the determination of the transmittance of an optical material wherein the ratio of the output radiant power to the input radiant power is measured; the measurement result is independent of SI units.

Absolute Accuracy and Traceability

Establishment of legal traceability to an SI unit requires that one obtain legally correct documentation, i.e., certification, of the device that serves as the calibration transfer standard and sometimes of the particular measurement process in which the device is to be used. Certification of legal traceability within each nation is obtained from the national standards laboratory of that nation. Often another nation's standards laboratory can be used to establish legal traceability, provided that there exists the legal framework for mutual recognition of the legality of each other's standards.

In order to establish accurate traceability to an SI unit, one needs to determine the total accumulated error arising from: (1) the realization of the base SI unit; (2) if applicable, the derivation of an associated measurement quantity; (3) if applicable, scaling to a higher or lower value; and (4,5...) transfer of the calibration from one device to another. The last entries must include the instability of the calibration transfer devices; the others may or may not involve a transfer device.

Legal traceability to SI units does not guarantee accurate traceability and vice versa. In order to obtain accurate traceability, it is not necessary to prove traceability to a national standards laboratory. Instead, the measurement must trace back to one of the SI units. However, it is usually convenient to establish accurate traceability via one of the national standards laboratories. The degree of convenience and accuracy will depend upon the accuracy of the measurement method and type of calibration transfer device available from the particular national standards laboratory.

Although relative measurements do not require traceability to one of the SI base units often to satisfy legal requirements traceability to a national standards laboratory may be necessary.

Types of Errors, Uncertainty Estimates, and Error Propagation

It is almost pointless to state a value for an absolute or relative measurement without an estimate of the uncertainty and the degree of confidence to be placed in the uncertainty estimate. Verification of the accuracy and the confidence limits is not only desirable but is often a legal requirement.

The accuracy and the uncertainty of a measurement are synonymous. The usual terminology is that a measurement is "accurate to within $\pm x$ " or "uncertain to within $\pm x$ ", where x is either a fraction (percent) of the measured value or an interval within which the true value is known to within some degree of confidence. The degree of confidence in the uncertainty estimate is the confidence interval or σ -level.

Errors are classified as type A errors, also known as random errors, and type B, or systematic errors. Type A errors are the variations due to the effects of uncontrolled variables. The magnitude of these effects is usually small and successive measurements form a random sequence. Type B errors are not detectable as variations since they do not change for successive measurements with a given apparatus and measurement method. Type B errors arise because of differences between the ideal behavior embodied in fundamental laws of physics and real behavior embodied in an experimental simulation of the ideal. A type B error could also be a function of the quantity being measured; for example, in a blackbody radiance standard using the freezing point of a metal and its defined temperature instead of the true absolute temperature.

Type A errors are estimated using standard statistical methods. If the distribution of the measurements is known (e.g., either Gaussian, which is often called a normal distribution, or Poisson), then one uses the formalism appropriate to the distribution. Unless enough data is obtained to establish that the distribution is not Gaussian, it is usual to assume a gaussian distribution. A brief discussion of Gaussian statistical concepts and terminology is given here to guide the reader in interpreting or determining the uncertainty in a radiometric or photometric calibration. A thorough discussion of these topics is available via the web from NIST.²²

The mean value m , the standard deviation σ , and the standard deviation of the mean σ_m , of a set of measurement x_i , are estimated for a small sample from a gaussian distribution of measurements as follows:

$$m = \sum_i \frac{x_i}{\eta} \quad (56)$$

$$\sigma^2 = \frac{\sum (x_i - m)^2}{\eta - 1} \quad (57)$$

$$\sigma_m = \frac{\sigma}{\sqrt{\eta}} \quad (58)$$

where $i = 1$ to η , and η is the total number of measurements.

The standard deviation is an estimate of the spread of the individual measurements within a sample, and it approaches a constant value as η is increased.

The standard deviation of the mean is an estimate of the spread of the values of the mean that would be obtained from several different sets of sample measurements. The standard deviation of the mean decreases as the number of samples in a set increases, since the estimate of the mean approaches the true mean for an infinite data set. The standard deviation of the mean is used in the estimate of the confidence interval assigned to the reported value of the mean.

The degree of confidence to which a reported value of the mean is valid is known as the *confidence interval* (CI). If it is assumed that a very large set of measurements has been sampled, then the CI is often given in terms of the number of standard deviations of the mean (one- σ level, two- σ level, etc.) within which the type A error of a reported value is known.

The CI is the probability that the mean from a normal distribution will be within the estimated uncertainty. That is, for a z -percent confidence interval, z -percent of the measurements will fall inside and $(100 - z)$ percent will fall outside of the uncertainty estimate. For small measurement samples from a gaussian distribution, Student's t -distribution is used to estimate the CI. Tables of Student's t -distribution along with discussions concerning its use are presented in most statistics textbooks. For large sets of measurements, a one- σ level corresponds approximately to a 68-percent CI, a two- σ level to a 95-percent CI, and a three- σ level to a 99.7-percent CI.

The reader is cautioned about using the σ level designation to describe the CI for a small sample of measurements. As an example of the small versus large sample difference, consider two data sets, one consisting of three samples and the other ten. Using a Student's t -distribution to estimate the CI for the three sample set, the one- σ , two- σ , and three- σ levels correspond to CIs of 61 percent, 86 percent, and 94 percent, respectively. For the ten-sample set, the respective CIs are 66 percent, 93 percent, and 99 percent. It can be quite misleading to state only the σ level of the uncertainty estimate without an indication of the size of the measurement set from which it was drawn. In order to avoid misleading accuracy statements, it is recommended that, instead of simply reporting the σ level, either the estimated CI be reported or the standard deviation of the mean be reported along with the number of measurement samples obtained.

Type B error estimates are either educated guesses of the magnitude of the difference between the real and the ideal or they are the result of an auxiliary measurement. If an appropriate auxiliary experiment can be devised to measure a systematic or type B error, then it need no longer be considered an error. The result obtained from the auxiliary measurement can usually be used as a correction factor. If a correction factor is applied, then the uncertainty is reduced to the uncertainty associated with the auxiliary experiment.

Most of the effort in high-accuracy radiometry and photometry is devoted to reducing type B errors. The first rule for reducing type B errors is to ensure that the experiment closely simulates the ideal. The second rule is that the differences between real and ideal should be investigated and that a correction be applied. Unlike type A errors for which an objective theory exists, the educated guess for a type B error is often subjective. For type B errors, neither a confidence interval nor a σ level is objectively quantifiable.

Error propagation, error accumulation, or a combined uncertainty analysis is the summation of all the type A and type B uncertainties that contribute to the final measurement in the chain. Because type A errors are truly random, they are uncorrelated and the accumulated type A error is obtained from the square-root of the sum of the squares (also known as root-sum-square, RSS) of the several type A error estimates. Type B uncertainties, however, may be either correlated or uncorrelated. If they are uncorrelated, the total uncertainty is the RSS of the several estimates. Type B uncertainties that are correlated must be arithmetically summed in a way that accounts for their correlation. Therefore, it is usually desirable to partition type B uncertainties so that they are uncorrelated.

Absolute Sources

Planckian or Blackbody Radiator A blackbody, or planckian, radiator is a thermal radiation source with a predictable absolute radiance output. An ideal blackbody is a uniform, i.e., lambertian, source of radiant power having a predictable distribution over area, solid angle, and wavelength. It is used as a standard radiance source from which the other radiometric quantities, e.g., irradiance, intensity, etc., can be derived.

Blackbody simulators are in widespread use not only at national standards laboratories but also in many other industrial, academic, and government laboratories. Blackbody simulators are commercially available from a number of manufacturers and cover a wide range of temperatures and levels of accuracy. Because they are in such widespread use as absolute standard sources for a variety of radiometric applications, particularly in the infrared, they are discussed here in some detail. Furthermore, since many practical sources of radiation can be approximated as a thermal radiation source, a blackbody function is often used in developing the radiometric model of an optical system.

An ideal blackbody is a completely enclosed volume containing a radiation field which is in thermal equilibrium with the isothermal walls of the enclosure that is at a known absolute temperature. The radiation in equilibrium with the walls does not depend upon the shape or constitution of the walls provided that the cavity dimensions are much larger than the wavelengths involved in the spectrum of the radiation.

Since the radiometric properties of a blackbody source are completely determined by its temperature, the SI base unit traceability for blackbody-based radiometry is to the kelvin. Because of recent improvements in the accuracy of absolute detector-based measurements thermodynamic temperatures are obtained by radiometric detector methods.²³

Since the radiation field and the walls are in equilibrium, the energy in the radiation field is determined by the temperature of the walls. The relationship between the absolute temperature T and the spectral radiance L_λ is given by Planck's law:

$$L_\lambda = \frac{2hc^2}{n^2\lambda^5} [e^{(hc/n\lambda kT)} - 1]^{-1} \quad (59)$$

Here h is Planck's constant, c is the speed of light in a vacuum, k is Boltzmann's constant, λ is the wavelength, and n is the index of refraction of the medium. Incorporating the values of the constants in this equation yields,

Spectral radiance units: $W\ m^{-2}\ sr^{-1}\ \mu m^{-1}$

$$L_\lambda = \frac{1.1910 \times 10^8}{n^2\lambda^5} [e^{(1.4388 \times 10^4/n\lambda T)} - 1]^{-1} \quad (60)$$

It follows that the peak of the spectrum of a blackbody is determined by its temperature (Wein displacement law).

$$n\lambda_{\max} T = 2898\ \mu mK \quad (61)$$

It is often useful to measure blackbody spectral radiance in units of photons per second N_λ . The form of Planck's law in this case is

Spectral radiance units: photons $s^{-1} m^{-2} sr^{-1} \mu m^{-1}$

$$N_{\lambda} = \frac{2c}{n\lambda^4} [e^{(hc/n\lambda kT)} - 1]^{-1} \quad (62)$$

The peak of this curve is not at the same wavelength as in the case of radiance measured in units of power. Wein's displacement law for blackbody radiance measured in photons per second is

$$n\lambda_{\max} T = 3670 \mu m K \quad (63)$$

In other applications, the spectral distribution of the blackbody radiation may be required in units of photons per second per frequency interval (symbol: N_{ν}). This form of Planck's law is
Spectral radiance units: photons $s^{-1} m^{-2} sr^{-1} Hz^{-1}$

$$N_{\nu} = \frac{2n^2\nu^2}{c^2} [e^{(h\nu/kT)} - 1]^{-1} \quad (64)$$

and that of Wein's displacement law is

$$\frac{T}{\nu_{\max}} = 1.701 \times 10^{-11} K Hz^{-1} \quad (65)$$

Planck's law integrated over all wavelengths (or frequencies) leads to the Stefan-Boltzmann law which describes the temperature dependence of the total radiance of a blackbody. For blackbody radiance measured as radiant power, the Stefan-Boltzmann law is

Radiance units: $W m^{-2} sr^{-1}$

$$L = 1.8047 \times 10^{-8} n^2 T^4 \quad (66)$$

Equation (66) is the usual form of the Stefan-Boltzmann law; however, it can also be derived for blackbody radiance measured as photon flux.

Radiance units: photons $s^{-1} m^{-2} sr^{-1}$

$$N = 4.8390 \times 10^{14} n^2 T^3 \quad (67)$$

The preceding expressions are valid provided that the cavity dimensions are much larger than the wavelengths involved in the spectrum of the radiation. The restriction imposed by the cavity dimension may lead to significant errors in very high accuracy radiometry or very long wavelength radiometry. For example, in a cube 1 mm on a side and at a wavelength of 1 μm , the approximate correction to Planck's equation is only 3×10^{-7} ; however, if the measurement is made within a 1-nm bandwidth or less, the root mean square fluctuation of the signal is about 2×10^{-3} which may not be negligible. Recent work describes how well the Planck and Stefan-Boltzmann equations describe the radiation in small cavities and at long wavelengths.²⁴⁻²⁶

Blackbody Simulators An ideal blackbody, being completely enclosed, does not radiate into its surrounds and therefore cannot serve as an absolute radiometric source. A blackbody simulator is a device that does emit radiation but only approximates the conditions under which Planck's law is valid. In general, a blackbody simulator is an enclosure at some fixed temperature with a hole in it through which some of the radiation is emitted. Some low-accuracy blackbody simulators are fabricated as a flat surface held at a fixed temperature.

A blackbody simulator can be used as an absolute source provided that the type B errors introduced by the deviations from the ideal Planck's-law conditions are evaluated and the appropriate corrections are applied. In a blackbody simulator there are three sources of type B error: inaccurate surface temperature, nonequilibrium between the radiant surface and the radiation field due to openings in the enclosure, and nonuniformity in the temperature of the radiant surface.

Calculation of the effect of a temperature error on the spectral radiance is obtained from the derivative of Planck's law with respect to temperature.

$$\frac{dL_\lambda}{L_\lambda} = \frac{hc}{n\lambda kT} [1 - e^{-(hc/n\lambda kT)}]^{-1} \frac{dT}{T} \quad (68)$$

Since the radiation field is in equilibrium with the surface of the cavity, it is the absolute temperature of the surface that must be measured. It is usually impractical to have the thermometer located on the emitting surface and it is the temperature within the wall that is measured. The difference between the temperature within the wall and the surface must therefore be measured, or calculated from a thermal model, and the correction applied.

The error due to nonequilibrium occurs because a practical radiation source cannot be a completely closed cavity. The correction factor for the effect on the radiance due to the escaped radiation is obtained from application of Kirchhoff's law. Simply stated, Kirchhoff's law states that the absorptive power of a material is equal to its emissive power. According to the principle of detailed balancing, for a body to be in equilibrium in a radiation field, the absorption of radiation by a given element of the surface for a particular wavelength, state of polarization, and in a particular direction and solid angle must equal the emission of that same radiation. If this were not true, the body would either emit more than it absorbs or vice versa, it would not be in equilibrium with the radiation, and it would either heat up or cool off.

Radiation impinging upon a body is either reflected, transmitted, or absorbed. The fraction of the incident radiation that is reflected ρ (reflectance), plus the fraction absorbed α (absorptance), plus the fraction transmitted τ (transmittance), is equal to one.

$$1 = \rho + \alpha + \tau \quad (69)$$

From Kirchhoff's law for a surface in radiative equilibrium, the fraction of absorbed radiation equals the fraction emitted ε (emittance or emissivity). Therefore, the sum of the reflectance, transmittance, and emittance must also be equal to one. If the body is opaque, the transmittance is zero and the emittance is just equal to one minus the reflectance.

$$\varepsilon = 1 - \rho \quad (70)$$

For a body not in an enclosed volume to be in equilibrium with a radiation field, it must absorb all the radiation impinging upon it, because any radiation lost through reflection will upset the equilibrium. An emittance less than one is the measure of the departure from a perfect absorber and, therefore, it is a measure of the radiance change due to the departure from closed-cavity equilibrium. In general, cavities with an emittance nearly equal to unity are those for which the size of the hole is very small in comparison to the size of the cavity.

Temperature nonuniformity modifies the radiant flux over the whole cavity in much the same way as the presence of a hole in that it is a departure from equilibrium. Radiation loss from the region of the cavity near the hole is typically larger than from other regions and this loss produces a temperature change near the hole and a nonuniformity along the cavity wall. In addition, the temperature nonuniformity is another source of uncertainty in the absolute temperature. In practice, the limiting factor in the accuracy of a high-emittance blackbody simulator is typically the nonuniformity of the temperature.

Accurate calculation of the emittance of a cavity radiator requires a detailed knowledge of the geometry of the cavity and the viewing system. This is a radiance transfer calculation and, in order to perform it accurately, one must know the angular emitting or reflecting properties of the cavity surface. The regions that contribute most to the accuracy of the calculation are those that radiate directly out the hole into the direction of the solid angle of the optical detection system.

There are many methods of calculating the emittance. The most popular are based upon the assumption of uniform emission that is independent of direction, i.e., lambertian emission. One can calculate the spectral emittance and temperature of each element along the cavity wall and sum the contribution from each element to the cavity radiance. Extensive discussion of the diffuse emittance and temperature nonuniformity calculation methods can be found elsewhere.²⁷⁻³¹

Instead of calculating the emittance of a cavity directly, the problem may be transformed into one of calculating the absorptance for a ray incident from the direction in which the emittance is required.³²⁻³⁴ The quantity to be calculated in this case is that fraction of the radiation entering the hole from a particular direction which is subsequently reflected out of the hole into a hemisphere.

Real surfaces are not perfectly diffuse reflectors and often have a higher reflectance in the specular direction. A perfect specularly reflecting surface is at the other extreme for calculating the emittance of a blackbody simulator. In some applications, a specular black surface might perform better than a diffusely reflecting one, particularly if the viewing geometry is highly directional and well known. The calculation of the emittance of a cavity made from a perfectly specular reflector is obtained in terms of the number of reflections undergone by an incident ray before it leaves the cavity.³⁵

One can reduce the error due to temperature nonuniformity by reducing the emittance of those regions along the cavity wall that do not contribute radiation directly to that emitted from the cavity.³⁶ That is, by fabricating the “hidden” portions of the cavity wall from a specular, highly reflecting material and by proper orientation of these surfaces, the highly reflecting surfaces absorb almost none of the radiation but reflect it back to the highly absorbing surfaces. Since the highly reflective surfaces absorb and emit very little radiation, their temperature will have a minimal effect on the equilibrium within the cavity.

In high-accuracy applications, it is preferable to measure rather than calculate the emittance of the blackbody cavity. This can be done either by comparison of the radiance of the device under test to that of a higher-quality blackbody simulator (emittance closer to unity) or by a direct measurement of the reflectance of the cavity.³⁷ Accurate measurement of thermal nonuniformity by measurement of the variations in the radiance from different regions within the cavity is made difficult by the fact that radiance variations depend not only on the local temperature but also upon the emittance of the region.

Synchrotron Radiation A synchrotron is an electronic radiation source that if well-characterized has a predictable absolute radiance output. A synchrotron source is a very nonuniform, i.e., highly directional and highly polarized, radiance standard in contrast to a blackbody which is uniform and unpolarized. However, like a blackbody, a synchrotron has a predictable spectral output and it is useful as a standard radiance source from which the other radiometric quantities, e.g., irradiance, intensity, etc., can be derived.

Classical electrodynamic theory predicts that an accelerated charged particle will emit radiation. A synchrotron is a type of electron accelerator where the electron beam is accelerated in a closed loop and synchrotron radiation is the radiation emitted by the electrons undergoing acceleration. The development of these and other charged particle accelerators led to closer experimental and theoretical scrutiny of the radiation emitted by an accelerated charged particle. These studies culminated in Schwinger’s complete theoretical prediction, including relativistic effects, of the spectral and angular distribution of the radiation emitted by a beam in a particle accelerator.³⁸ The accuracy of Schwinger’s predictions have been verified in numerous experimental studies.³⁹⁻⁴¹

Schwinger’s theoretical model of the absolute amount of radiation emitted by an accelerated charged particle is analogous to the Planck equation for blackbody sources in that both predict the behavior of an idealized radiation source. Particle accelerators, when compared to even the most elaborate blackbodies, are, however, far more expensive. Furthermore, in order to accurately predict the spectral radiance of the beam in a particle accelerator, much detailed information is required of the type not found for most accelerators. Accurately predictable radiometric synchrotron sources are consequently found only in a few laboratories throughout the world.

The magnitude of the radiant power output from a synchrotron source is proportional to the number of electrons in the beam and their velocity, i.e., the number of electrons per second or current and their energy. Therefore, synchrotron radiometry is traceable to the SI unit of electricity, the ampere.

Because absolute synchrotron sources are so rare, a detailed discussion of Schwinger’s model of synchrotron radiation and the various sources of uncertainty will not be presented here. It is generally useful, however, to know some of the characteristics of synchrotron radiation. For example, the radiance from a synchrotron beam is highly polarized and very nonuniform: radiant power is almost entirely in the direction of the electron velocity vector and tangent to the electron beam. The peak of the synchrotron radiation spectrum varies from the vacuum ultraviolet to the soft x-ray region depending upon the energy in the beam. Higher-energy beams have a shorter wavelength peak: 1-GeV

peaks near 10 nm, 6-GeV peaks near 0.1 nm. Radiant power decreases to longer wavelengths by very roughly two decades for every decade increase of wavelength, so that for the typical radiometric-quality synchrotron source, there is usually sufficient energy to perform accurate radiometric measurements in the visible for intercomparison to other radiometric standards.⁴¹

Absolute Detectors

Electrical Substitution Radiometers An electrical substitution radiometer, often called an electrically calibrated detector, is a device for measuring absolute radiant power by comparison to electrical power.⁸ As a radiant power standard, an electrical substitution radiometer can be used as the basis for the derivation of the other radiometric quantities (irradiance, radiance, or intensity) by determining the geometrical distribution (either area and/or solid angle) of the radiation.^{42–48} Since an electrical substitution radiometer measures the spectrally total radiant power, it is used primarily for the measurement of monochromatic sources or those with a known relative spectral distribution.

An electrical substitution radiometer consists of a thermal detector (i.e., a thermometer) that has a radiation-absorbing surface and an electrical heater within the surface, or the heater is in good thermal contact with the surface. When the device is irradiated, the thermometer senses the temperature of the radiantly heated surface. The radiation source is then blocked and the power to the electrical heater adjusted to reproduce the temperature of the radiantly heated surface. The electrical power to the heater is measured and equated to the radiant power on the surface. The absolute base for this measurement is the electrical power measurement which is traceable to the SI ampere. In order for the measurement to be accurate, differences between the radiant and electrical heating modes must be evaluated and the appropriate corrections applied.

Electrical substitution radiometers predate the planckian radiator as an absolute radiometric standard.^{49–50} They were the devices used to quantify the radiant power output of the experimental blackbody simulators studied at the end of the nineteenth century. Electrical substitution radiometers are in widespread use today and are commercially available in a variety of forms that can be classified either as to the type of thermometer, the type of radiant power absorber, or the temperature at which the device operates.

Early electrical substitution radiometers operated at ambient temperature and used either a thermocouple, a thermopile, or a bolometer as the detector. Thermopile- and bolometer-based radiometers are presently used in a variety of applications. They have been refined over the years to produce devices of either greater accuracy, sensitivity, and/or faster response time. Thermopile-based, ambient temperature electrical substitution radiometers used for radiant power (and laser power, see later discussion) measurements at about the 1-mW level at several national standards laboratories have estimated uncertainties reported to be within ± 0.1 percent.^{51–53} Electrical substitution radiometers have also been used for very high-accuracy absolute radiant power measurements of the total solar irradiance both at the surface of the earth,^{51,53} and above its atmosphere.⁵² The type of high-accuracy radiometer used at various national standards laboratories is a custom-built device and is not commercially available in general. On the other hand, electrical substitution radiometers for solar and laser power measurements at a variety of accuracy levels are available commercially.

An ambient temperature electrical substitution radiometer based on a pyroelectric as the thermal detector was developed in the 1970s.^{54,55} A pyroelectric detector is a capacitor containing a dielectric with a temperature-sensitive spontaneous electrical polarization; a change in temperature results in a change in polarization. Small and rapid changes of polarization are readily detectable, making the pyroelectric a sensitive and fast thermal detector. It is most useful as a detector of a pulsed or chopped radiant power signal. During the period when the radiant power signal is blocked, electrical power can be introduced to a heater in the absorptive surface of the radiometer. As in the method for a thermopile- or bolometer-based electrical substitution radiometer, the electrical power is adjusted to equal the heating produced by the radiant power signal. Chopping can be done at a reasonable frequency, hence the electrical heating can be adjusted to achieve a balance in a comparatively short time. Because the radiant-to-electrical heating balance is more rapidly obtained, a pyroelectric radiometer is often more convenient to use than a thermopile radiometer. Pyroelectric

electrical substitution radiometers are generally more sensitive but are usually less accurate than the room temperature thermopile or bolometer electrical substitution radiometers.

Electrical substitution radiometers are further distinguished by two types of radiant power absorber configurations: a flat surface coated with a highly absorbing material or a cavity-shaped, light-trapping detector. Cavity-shaped radiometers are usually more accurate over a greater spectral range than flat-surface radiometers. However, a flat-surface receiver can usually be fabricated with less thermal mass than a cavity-shaped receiver and therefore may have greater sensitivity and/or a faster response time.

Electrical substitution radiometers are further distinguished by the temperature at which the electrical-to-radiant-power comparison is performed. In the last two decades there have been significant advancements^{56,57} made in instruments that perform the radiant-to-electrical comparison at a temperature near to that of liquid helium (4.2 K). Such devices are known as cryogenic electrical substitution radiometers or electrically calibrated cryogenic detectors and they are commercially available. Cryogenic electrical substitution radiometers are presently the most accurate absolute radiometric devices; the uncertainty of some measurements of radiant power has been estimated to be within ± 0.005 percent.⁴⁷

Sources of Error in Electrical Substitution Radiometers The relative significance of each of the possible sources of error and the derived correction factor depends upon the type of radiometer being used and the particular measurement application. It is possible to determine the total error occurring in equating radiant to electrical power and thence the accuracy of the traceability to the absolute electrical SI unit of measurement. Most manufacturers provide extensive characterization of their instruments. In such cases, the traceability to SI units is independent of radiometric standards such as blackbodies, hence electrical substitution radiometers are sometimes called absolute detectors. A commercially produced electrical substitution radiometer is capable of far greater accuracy (within ± 0.01 percent for the cryogenic instruments) than any of the typical radiometric transfer devices available from a national standards laboratory. Hence, establishing traceability through a radiometric standard from a national standards laboratory is almost pointless for a cryogenic electrical substitution radiometer.

The sources of error in an electrical substitution radiometer can be divided into three categories: errors in traceability to the absolute base unit, errors due to differences in the radiant-versus-electrical heating modes of operation, and errors arising in a particular application. The major error sources common to all electrical substitution radiometers as well as some of the less common are briefly described here. An extensive listing and description of all of these errors is given in Ref. 8, Chap. 1.

Electrical power measurement accuracy is first determined by the accuracy of the voltage and resistance standards (or voltmeter and resistance meter) used to measure voltage and current. Electrical power measurement accuracy within ± 0.01 percent is readily achievable and if needed it can be improved by an order of magnitude or better. Additional error is possible due to improper electrical measurement procedures such as those giving rise to ground-loops (improper connection to earth).

Differences between electrical-versus-radiant heating appear as differences in radiative, conductive, or convective losses. Most of these differences can be measured and a correction factor applied to optimize accuracy. The most obvious example is probably that of the radiative loss due to reflection from the receiver surface. Less obvious perhaps is the effect due to extraneous heating in the portion of the electrical conductors outside the region defined by the voltage sensing leads.

Differences between electrical heating and radiant heating may also arise due to spatial nonuniformity of the thermal sensor and/or differences in the heat conduction paths in the electrical-versus-radiant heating modes. These effects are specific to the materials and design of each radiometer. The electrical heater is typically buried within the device, whereas radiant heating occurs at the surface, so that the thermal conductivity paths to the sensor may be very different. Also, the distribution of the radiant power across the receiver is usually quite different compared to the distribution of the electrical heating. A detailed thermal analysis is required to create a design which minimizes these effects, but for optimum accuracy, the measurement of the magnitude of the nonuniformity effects is required to test the thermal model. Nonuniformity can be measured either by placing small auxiliary electrical heaters in various locations or by radiative heating of the receiver in several regions by moving a small spot of light across the device.

It should be noted that the thermal conduction path differences may also be dependent upon the environment in which the radiometer is to be operated. For example, atmospheric-pressure-dependent differences between the electrical-to-radiant power correction factor have been detected for many radiometers. These differences are, of course, greatest for a device for which the correction factors have been characterized in a normal atmosphere and which is then used in a vacuum.

Application-dependent errors arise from a variety of sources. Some examples are window transmission losses if a window is used, the accuracy of the aperture area and diffraction corrections are critical for measurements of irradiance; and, if a very intense source such as the sun is measured, heating of the instrument case and the body of the aperture could be an important correction factor. The last effect might also be very sensitive to atmospheric pressure changes.

Photoionization Devices Another type of absolute detector is a photoionization detector which can be used for absolute photon flux, i.e., radiant power, measurements of high-energy photon beams. Since a photoionization detector is a radiant power standard like the electrical substitution radiometer, it can in principle be used as the basis for the derivation of the other radiometric quantities (irradiance, radiance, or intensity) by determining the geometrical distribution (either area and/or solid angle) of the radiation.

A photoionization detector is a low-pressure gas-filled chamber through which a beam of high-energy (vacuum ultraviolet) photons is passed between electrically charged plates, the electrodes. The photons absorbed by the gas, if of sufficient energy, ionize the gas and enable a current to pass between the electrodes. The ion current is proportional to the number of photons absorbed times the photoionization yield of the gas and is, therefore, proportional to the photon flux.

The photoionization yield is the number of electrons produced per photon absorbed. If the photon is of sufficiently high energy, the photoionization yield is 100 percent for an atomic gas. The permanent atomic gases are the rare gases: helium, neon, argon, krypton, and xenon. Their photoionization yields have been measured relative to each other and shown to be 100 percent over specific wavelength ranges.^{58,59} If an ionization chamber is constructed properly and filled with the appropriate gas so that all of the radiation is absorbed, then the number of photons per second incident on the gas is simply equal to the ion current produced. If instead of measuring the ion current one were to measure each pulse produced by a photon absorption, then one would have a photon counter.

Carefully constructed ion current measurement devices have been used as absolute detectors from 25 to 102.2 nm and photon counters from 0.2 to 30 nm. Careful construction implies that all possible systematic error sources have either been eliminated or can be estimated, with an appropriate correction applied. Because of the difficulty in producing accurate and well-characterized devices, ion chambers and high-energy photon counters are not claimed to be high-accuracy radiometric devices. Furthermore, they are limited to applications in vacuum ultraviolet radiometry and are consequently of restricted interest.

Predictable Quantum Efficiency Devices A useful and quite economical type of absolute detector is a predictable quantum efficiency (PQE) device using high-quality silicon photodiodes. Quantum efficiency is the photon flux-to-photocurrent conversion efficiency. Because there have been many technological advancements made in the production of solid-state electronics, it is now possible to obtain very high quality silicon photodiodes whose performance is extremely close to that of the theoretical model.^{60,61} The technique for predicting the quantum efficiency of a silicon photodiode is also known as the self-calibration of a silicon photodiode.^{62,63} It is a relatively new absolute radiometric technique, quite simple to implement and of very high accuracy.^{64,65}

Conversion of a detector calibration from spectral responsivity $R(\lambda)$, in units of A/W, to quantum efficiency, i.e., photon-to-electron conversion efficiency, is as follows:

$$C_e = 1239.85 \frac{R(\lambda)}{\lambda} \quad (71)$$

where λ is the in-vacuum wavelength in nm and C_e is in units of electrons per photon.

As in the case of the other absolute detectors discussed previously, a PQE device is used for absolute photon flux, i.e., radiant power, measurements. It can also be used as the basis for the derivation

of the other radiometric quantities such as irradiance, radiance, or intensity. The extension to other radiometric measurements is by the determination of the geometrical distribution (area and/or solid angle) of the radiation. Also, like other absolute detectors, it measures spectrally total flux (within its spectral response range) and is therefore used primarily for the measurement of monochromatic sources or those with a known relative spectral distribution.

In a solid-state photodiode, the process for the conversion of a photon to an electronic charge is as follows. Photons not lost through reflection or by absorption in a coating at the front surface are absorbed in the semiconductor—if the photon is of high enough energy. To be absorbed, the photon energy must be greater than the band gap; the band gap for silicon is 1.11 eV (equivalent wavelength, 1.12 μm). In silicon, the absorption of a photon causes a promotion of a charge carried to the conduction band. Absorption of very high energy photons will create charge carriers with sufficient energy to promote a second, third, or possibly more charge carriers into the conduction band by collision processes. However, for silicon, the photon energy throughout the visible spectral range is insufficient for such impact ionization processes to occur. Therefore, in the visible to near-ir spectral region (about 400 to 950 nm), one absorbed photon produces one electron in the conduction band of silicon.

In a photodiode, impurity atoms diffused into a portion of the semiconductor material create an electric field. The internal electric field causes the newly created charge carriers to separate, eventually promoting the flow of an electron in an external measurement circuit. The efficiency with which the charge carriers are collected depends upon the region of the photodiode in which they are created. In the electric field region of a high-quality silicon photodiode, this collection efficiency has been demonstrated to approach 100 percent to within about 0.01 percent. Outside the field region, the collection efficiency can be determined by simple electrical bias measurements.

For the spectral regions in which the collection efficiency is 100 percent, the only loss in the photon-to-electron conversion process is due to reflection from the front surface of the detector. Several silicon photodiodes can be positioned to more effectively collect the radiation, acting as a light trap.^{66,67} If the radiation reflected from the first photodiode is directed to a second photodiode, then onto a third photodiode, etc., almost all the radiation will be collected in a small number of reflections. The photocurrents from all of the photodiodes are then summed and the total current (electrons per second) will be nearly equal, within 0.1 percent or less, to the photon flux (photons per second).

The more common type of silicon photodiode is the pn-type (positive charge impurity diffused into negative charge impurity starting material). High-quality pn-type detectors have their high collection efficiency in the long wavelength visible to near-ir spectral region. On the other hand, np-type silicon photodiodes have high collection efficiency in the short wavelength spectral region. At this time, the silicon photodiodes with the highest quantum efficiency (closest to ideal behavior) in the blue spectral region are the np-type devices, while nearly ideal red region performance is obtained with pn-type devices. The predictable quantum efficiency technique for silicon photodiodes has been demonstrated⁶⁴⁻⁶⁷ to be absolutely accurate to within ± 0.1 percent from about 400 nm to 900 nm.

A disadvantage of the light-trap geometry is the limited collection angle (field of view) of the device. Light-trap silicon photodiode devices are now commercially available using large area devices and a compact light-trap configuration that maximizes the field of view.

An np-type silicon photodiode trap detector optimized for short-wavelength performance and a pn-type silicon photodiode trap detector optimized for long-wavelength performance can be used as an almost ideal radiometric standard. The pair covers the 400- to 900-nm spectral range, has direct absolute SI base unit traceability via convenient electrical standards, and they are sufficiently independent to be meaningfully cross-checked to verify absolute accuracy and long-term stability. These detectors are not only useful radiometric standards by themselves but can be used with various source standards to either verify the absolute accuracy or to correct for the instabilities in the source standards.

The concept of a PQE light-trapping device is extendable to other high-quality photodiodes. Very recently, InGaAs devices with nearly 100-percent collection efficiency in the 1000- to 1600-nm spectral range have been developed. A light-trapping device employing these new detectors is now commercially available.

Calibration Transfer Devices

The discussion to this point focused on absolute radiometric measurements using methods that in themselves can be made traceable to absolute SI units. It is often more convenient (and sometimes required by contractual agreements) to obtain a device that has been calibrated in radiometric units at one of the national standards laboratories. Specific information as to the type and availability of various calibration transfer devices and calibration services may be obtained by directly contacting any of the national standards laboratories in the world. The products and services offered by the various standards laboratories cover a range of applications and accuracies, and differ from country to country.

Radiometric calibration transfer devices are either sources or detectors. The calibration transfer sources are either incandescent, tungsten filament lamps, deuterium lamps, or argon arc discharge sources.⁶⁸⁻⁷⁰ Generally, calibration transfer detectors are photodiodes of silicon, germanium, or indium gallium arsenide. The most prevalent calibration transfer sources are incandescent lamps and the typical calibration transfer detector is a silicon photodiode.⁷¹

The commonly available spectral radiance calibration transfer devices that span the 250- to 5000-nm region are typically tungsten strip filament lamps. Lamps calibrated in the 250- to 2500-nm region by a national standards laboratory are available. Lamps calibrated in the 2.5- to 5- μm region by comparison to a blackbody are commercially available. These devices are calibrated within specific geometrical constraints: the area on the filament, and the direction and solid angle of observation. The calibration is reported at discrete wavelengths, for a specified setting of the current through the filament and the ambient laboratory temperature. The optimum stability of spectral radiance is obtained with vacuum rather than gas-filled lamps, and with temperature controlled, i.e., water-cooled electrodes. Vacuum lamps cannot be operated at high filament temperatures and consequently do not have sufficient uv output. Gas-filled lamps cover a broader spectral and dynamic range and are the more commonly available calibration transfer device.

The commonly available spectral irradiance calibration transfer devices that span the 250- to 2500-nm region are tungsten coiled filament lamps. These are usually gas-filled lamps that have a halogen additive to prolong filament life and enable higher-temperature operation. Lamps calibrated in the 250- to 2500-nm region by a national standards laboratory are available. These devices are calibrated within the specific geometrical constraints of the distance and the direction with respect to a location on the lamp base or the filament. The calibration is reported at discrete wavelengths, for a specified setting of the current through the filament and the ambient laboratory temperature. Because the filament is operated at a higher temperature, the spectral irradiance lamps are usually less stable than the radiance lamps.

The drift of an incandescent lamp's radiance or irradiance output is not reliably predictable. It is for this reason that the calibration is most reliably maintained not by an individual lamp but by a group of lamps. The lamps are periodically intercompared and the average radiance (irradiance) of the group is considered to be the calibration value. The calculated differences between the group average and the individual lamps is used as a measure of the performance of the individual lamp. Lamps that have drifted too far from the mean are either recalibrated or replaced.

Spectral radiance and irradiance calibration transfer devices for the vacuum to near-uv (from about 160 to 400 nm) are typically available as deuterium lamps.

The commonly available calibrated transfer detectors for the 250- to 1100-nm spectral region are silicon photodiodes and for the 1000- to 1700-nm region, they are either germanium or indium gallium arsenide photodiodes. The calibration is reported at discrete wavelengths in absolute responsivity units (A/W) or irradiance response units ($\text{A cm}^2/\text{W}$). In the first case, the calibration of the detector is performed with its active area underfilled, while in the second case, it is overfilled. If the detector is fitted with a precision aperture and if its spatial response is acceptably uniform, then the area of the aperture can be used to calculate the calibration in either units. The conditions under which the device was calibrated should be reported. The critical parameters are the location and size of the region within the active area in which it was calibrated, the radiant power in the calibration beam (alternately the photocurrent), and the temperature at which the calibration was performed. The direction in which the device was calibrated is usually assumed to be normal to its surface and

the irradiation geometry is usually that from a nearly collimated beam. Significant departures from normal incidence or near collimation should be noted.

Lasers

Power and Energy Measurement Lasers are highly coherent sources and the previous discussion of radiometry has been limited to the radiometry of incoherent sources. Nevertheless, the absolute power (or energy) in a laser beam can be determined to a very high degree of accuracy (within ± 0.01 percent in some cases) using some of the detector standards discussed here. The most accurate laser power measurements are made with cryogenic and room temperature electrical substitution radiometers and with predictable quantum efficiency devices. In order to measure the laser power (energy) it is necessary to ensure that all the radiation is impinging on the sensitive area of the detector and, if the absolute detector characterization was obtained at a different power (energy) level, that the detector is operating in a linear fashion. For pulsed lasers, the peak power may substantially exceed the dynamic range of the detector's linear performance. (A discussion of detector linearity is presented later in this chapter.) Furthermore, caution should be exercised to ensure that the detector not be damaged by the high photon flux levels achieved with many lasers.

In addition to ensuring that the detector intercept all of the laser beam, it is necessary to determine that all coherence effects have been eliminated (or minimized and corrected).^{72,73} The predominant effects of coherence are, first, interference effects at windows or beam splitters in the system optics and, second, diffraction effects at aperture edges. The use of wedged windows will minimize interference effects, and proper placement of apertures or the use of specially designed apertures⁷⁴ will minimize diffraction effects.

Lasers as a Radiometric Characterization Tool It should be noted that lasers, particularly the cw (continuous wave) variety, are particularly useful as characterization tools in a radiometric laboratory. Some of their applications are instrument response uniformity mapping, detector-to-detector spectral calibration transfer, polarization sensitivity, linearity verifications, and both diffuse and specular reflectance measurements.

Lasers are highly polarized and collimated sources of radiation. It is usually simple to construct an optical system as required for each measurement using mostly plane and spherical mirrors and to control scattered light with baffles and apertures. Lasers are high-power sources so that the signal-to-noise levels obtained are very good. If the power level is excessive it can usually be easily attenuated. Also, care must be taken to avoid local saturation of a detector at the peak of a laser's typical gaussian beam profile. They are highly monochromatic so that spectral purity, i.e., out-of-band radiation, is not usually a problem. However, in very high accuracy, within <0.1 percent, measurements, lasing from weaker lines may be significant and additional spectral blocking filters could be required.

Lasers are not particularly stable radiation sources. This problem is overcome by putting a beam splitter and stable detector into the optical system near the location of the measurement. The detector either serves to monitor the laser beam power and thereby supplies a correction factor to compensate for the instability, or its output is used to actively stabilize the laser.⁴³ In the latter case, an electronically controllable attenuator, such as an electro-optical, acousto-optical, or a liquid crystal system, is used to continuously adjust the power in the laser beam at the beam splitter. Feedback stabilization systems for cw lasers, both the electro-optical and liquid crystal type, are commercially available. For the highest-accuracy measurements, i.e., optimum signal-to-noise ratios, it is necessary both to actively stabilize the laser source and also to monitor the beam power close to the measurement in order to correct for the residual system drifts.

Various Type B Error Sources

Offset Subtraction One common error source, which is often simply an oversight, is the incorrect (or sometimes neglected) adjustment of an instrument reading for electronic and radiometric offsets.

This is often called the dark signal or dark current correction since it is obtained by shutting off the radiation source and reading the resulting signal. The shuttered condition needs to be close to radiant zero, at least within less than the expected accuracy of the measurement.

A dark signal reading is usually easy to achieve in the visible and near-visible spectral regions. However, in the long-wavelength infrared a zero radiance source is one that is at a temperature of ideally 0 K. Often an acceptably cold shutter is not easily obtained so that the radiance, i.e., temperature, of the “zero” reference source must be known in order to determine the true instrument offset.

Scattered Radiation and Size of Source Effect An error associated with the offset correction is that of scattered radiation from regions outside the intended optical path of the measurement system. Often, by judicious placement of the shutter, the principal light path can be blocked while the scattered light is not. In this case, the dark signal measurement includes the scattered light which is then subtracted from the measurement of the unshuttered signal. It is not possible to formulate a general scattered light elimination method so that each radiometric measurement system needs to be evaluated on an individual basis. The effects of scattered radiation can often be significantly reduced by using an optical chopper, properly placed, and lock-in amplifier system to read the output of the photodetector.

No optical element will produce a perfect image and there will be an error due to geometrically introduced stray light. Sharp edges between bright and dark regions will be blurred by aberrations, instrument fabrication errors, scattering due to roughness and contamination of the optical surfaces, and scattered light from baffles and stops within the instrument enclosure. Diffraction effects will also introduce stray light. Light originating from the source will be scattered out of the region of the image and light from the area surrounding the source will be scattered into the image. The error resulting from scattering at the objective lens or mirror is related to the size of the source since the scattering is proportional to the irradiance of the objective element. Thus, the error introduced by the lack of image quality is commonly referred to as the size-of-source effect.

The effect of the aberrations on the radiant power both into and out of an image can, in principle, be calculated. The diffraction-related error can also be calculated in some situations.⁷³⁻⁷⁵ However, the effect due to scattering is very difficult to model accurately and usually will have to be measured. In addition, the amount of scattering can be expected to change in time due to contamination of the optics, baffles, and stops. It is often more practical to measure the size of source effect and determine a correction factor for the elimination of this systematic error.

There are two different methods for measuring the size-of-source effect. The first method measures the response of the instrument as the size of the source is increased from the area imaged to the total area of the source. In the second method, a dark target of the same size as the image is placed at the imaged region on the source, and the surrounding area is illuminated. The second method has the advantage in that the effect being measured is the error signal above zero, whereas in the first method a small change in a large signal is being sought. In either case, the total error signal is measured; it includes aberrations, diffraction, and scattering effects.

Polarization Effects These are often significant perturbations of radiant power transfer due to properties of the radiation field other than its geometry. One such possible error is that due to the polarization state of the radiation field. The signal from a photodetector that is polarization-sensitive will be dependent upon the relative orientation of the polarization state of the radiation with respect to the detector orientation. Examples of polarization-dependent systems are grating monochromators and radiation transfer through a scattering medium or at a reflecting surface. In principle, the polarization state of the radiation field may be included in the geometrical transfer equation as a discrete transformation that occurs at each boundary or as a continuous transformation occurring as a function of position in the medium. Often it is sufficient to perform a calibration at two orthogonal rotational positions of the instrument or its polarization-sensitive components. However, it is recommended that other measurements at rotations intermediate between the two orthogonal measurements be included to test if the maximum and minimum polarization sensitivities have been sampled. The average of the maximum and minimum polarization measurements is then the calibration factor of the instrument for a nonpolarized radiation source.

Detector Nonlinearity

Nonlinearity measurement by superposition of sources Another possibly significant error source is photodetector and/or the electronic signal processing system nonlinearity. If the calibration and subsequent measurements are performed at the same radiant power level, then nonlinearity errors are avoided. Often conditions require that the measurements be performed over a range of power levels. In general, a separate measurement is required either to verify the linearity of the photodetector (and/or the electronics) or to deduce the form of the nonlinearity function in order to apply the appropriate correction.^{72,76–78}

The typical form of a nonlinearity appears as a saturation of the photoelectronic process at high irradiance levels. At low radiant flux levels what often appears to be a nonlinearity may be the result of failing to apply a dark signal or offset correction. There are, of course, other effects that will appear as a nonlinearity of the photodetector and/or electronics.⁷⁹

Either the linearity of the detector and electronics can be directly verified by experiment or it can be determined by comparison to a photodetector/electronics system of verified linear performance. It is useful to note that several types of silicon photodiodes using a transimpedance or current amplifier have been demonstrated to be linear within ± 0.1 percent over up to eight decades for most of its principal spectral range.⁷⁸

The fundamental experimental method for determining the dynamic range behavior of a photodetector is the superposition-of-sources method.^{76–78} The principle of the method is as follows. If a photodetector/electronics system is linear, then the arithmetic sum of the individual signals obtained from different radiant power sources should equal the signal obtained when all the sources irradiate the photodetector at the same time. There are many variations of the multiple source linearity measurement method using combinations of apertures or beam splitters. A note of caution: Interference effects must be avoided when combining beams split from the same source or when combining highly coherent sources such as lasers.

The difference between the arithmetic sum and the measured signal from the combined sources is used as the nonlinearity correction factor. Consider the superposition of two sources having approximately equal radiant powers ϕ_a and ϕ_b , which when combined have a radiant power of $\phi_{(a+b)}$. The signals from the photodetector when irradiated by the individual and combined sources is i_a , i_b , and $i_{(a+b)}$. The following equation would be equal to unity for a linear detector:

$$K_{ab} = \frac{i_{(a+b)}}{i_a + i_b} \quad (72)$$

For a calibration performed at the radiant power level ϕ_a (or ϕ_b), the detector responsivity is R and

$$i_a = R\phi_a \quad (73)$$

For a measurements at the higher radiant power level $\phi_{(a+b)}$,

$$i_{(a+b)} = K_{ab} R\phi_{(a+b)} \quad (74)$$

Scaling up to much higher radiant power levels (or down to lower levels) requires repeated application of the superposition-of-sources method. For example, in order to scale up to the next higher radiant power level, the source outputs from the first level are increased to match the second level (e.g., by using larger apertures). The increased source outputs are then combined to reach a third level and a new correction factor calculated. The process is repeated to cover the entire dynamic range of a photodetector/electronics system in factor-of-two steps.

Note that when type B errors, such as the interference effects noted above, are eliminated, the accumulated uncertainty in the source superposition method is the accumulated imprecision of the individual measurements.

Various nonlinearity measurement methods Other techniques for determining the dynamic range behavior of a photodetector are derivable from predictable attenuation techniques.⁷² One such method is based upon chopping the radiation signal using apertures of known area in a rotating disk. This is often referred to as Talbot's law: the average radiant power from a source viewed through the

apertures of a rotating disk is given by the product of the radiant power of the source and the transmittance of the disk. The transmittance of the disk is the ratio of the open area to the blocked area of the disk. The accuracy of this technique depends upon the accuracy with which the areas are known and may also be limited by the time dependence of the photodetector and/or electronics.

Another predictable attenuation technique is based upon the transmittance obtained when rotating, i.e., crossing two polarizers.

$$\tau = \tau_0 \cos^4 \delta \quad (75)$$

Here δ is the angle of rotation between the linear polarization directions of the two polarizers and τ_0 is the transmittance at $\delta = 0^\circ$. This technique, of course, assumes ideal polarizers that completely extinguish the transmitted beam at $\delta = 90^\circ$, and its accuracy is limited by polarization efficiency of the polarizers.

A third predictable attenuation technique is the application of Beer's law which states that the transmittance of a solution is proportional to the concentration κ of the solute

$$\tau = e^{-\gamma \kappa} \quad (76)$$

Here γ is the absorption coefficient of the solute. The accuracy of this technique depends upon the solubility of the solute and the absence of chemical interference, i.e., concentration-dependent chemical reactions.

Time-dependent Error For measurements of pulsed or repeatedly chopped sources of radiation, the temporal response of the detector could introduce a time-dependent error. A photodetector that has a response that is slow compared to the source's pulse width or the chopping frequency will not have reached its peak signal during the short time interval. Time-dependent error is avoided by determining if the detector's frequency response is suitable before undertaking the calibration and measurement of pulsed or chopped radiation sources.

Nonuniformity The nonuniformity of the distribution of radiation over an image or within the area sampled in an irradiance or radiance measurement may lead to an error if the response of the instrument is nonuniform over this area. The calibration factor for a nonuniform instrument will be different for differing distributions of radiation. The size of the error will depend upon the relative magnitudes of the source and instrument nonuniformities and it is a very difficult error to correct. This type of error is usually minimized either by measuring only sources that are uniform or by ensuring that the instrument response is uniform. It is usually easier to ensure that the instrument response is uniform.

Nonideal Aperture For very high-accuracy radiometric calibrations, the error due to the effect of the land on an aperture must be correctly taken into account. An ideal aperture is one that has an infinitesimally thin edge that intercepts the radiation beam. In practice an aperture will have a surface of finite thickness at its edge. This surface is referred to as the land; see Fig. 7. The effective

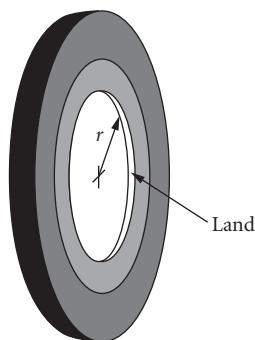


FIGURE 7 The nonideal aperture showing the location of the land.

radius of the aperture will be slightly reduced by vignetting caused by the land on its edge, assuming that the vignetting is small compared to the aperture radius and that all the radiation reflected by the land eventually falls on the detector (i.e., the land has a highly reflective surface). The effective radius of the aperture r' is

$$r' = r \left[1 - \frac{(1-\rho)w}{s} \right] \quad (77)$$

Here ρ is the reflectance, s is the distance between this aperture and the mirror or lens (or other aperture), and w is the width of the land.

Spectral Errors

Wavelength error In those cases where a spectrally selective element, such as a monochromator or a filter, is included in the radiometric instrument, spectral errors must be taken into account.⁴⁴ The first type of error is called the wavelength error and is due to misassignment of the wavelength of the spectrum of the filter or the monochromator in the instrument. That is, either the monochromator used to calibrate the filter transmission or the monochromator within the radiometric instrument has an error in its wavelength setting. This error is eliminated by calibration of the monochromator wavelength setting using one or more atomic emission lines from either a discharge lamp (usually a mercury and/or a rare gas lamp) or one of the many spectra of the elements available as a hollow cathode lamp. The wavelengths of most atomic emission lines are known with an accuracy that exceeds the requirements of radiometric calibrations.

A special note regarding the wavelength error and the use of interference filters in a radiometric instrument. The typical angular sensitivity of an interference filter is 0.1 nm per angular degree of rotation. If the transmission of the interference filter is measured in a collimated beam and then used in a convergent beam there will be an error due to the angularly dependent shift in the spectral shape of the transmission. For an accurate radiometric instrument it is important to measure interference filter transmission in nearly the same geometry as it will be used. Furthermore, the temperature coefficient of the transmission of an interference filter is about 0.2 nm K⁻¹. Therefore, it is also important in subsequent measurements to ensure that the filter remains at the same temperature as that used during the calibration.

It should also be noted that in order to accurately determine the spectral transmission of a monochromator it is necessary to completely fill the aperture of the dispersing element in the monochromator. This usually means that the field of view of the monochromator must be filled by the beam from the spectral calibration instrument.

Out-of-band radiation error The second type of spectral error is called the out-of-band or spectrally stray radiation error. This error is due to the radiation transmitted at both longer and shorter wavelengths that are beyond the edges of the principle transmission band of the filter or monochromator. This radiation is not taken into account if the limits in the integral in Eq. (11) are restricted to the edges of the principle transmission band. Although the relative amount of radiation transmitted at any wavelength beyond the edges of the principle transmission band may appear to be negligible with respect to the amount of radiation within the band, it is the spectrally total radiation that “leaks through” that is the error signal. It is therefore necessary to determine the transmission of the filter or monochromator over the entire spectrum of either the detector’s response and/or the light source’s output, whichever is greater. If the out-of-band radiation effect is small, it is possible to determine the correction factor from nominal values or limited accuracy measurements of the out-of-band spectra of the source, detector, and filter (monochromator).

Temperature Dependence The effect of temperature on the various elements in a radiometric instrument must not be overlooked. Unless the temperature of the instrument at the time of its calibration is maintained during subsequent applications, there may be substantial changes introduced in the instrument calibration factor that could well be above the uncertainty of its traceability to absolute SI units. The simple solution is to control the temperature of the system from the calibration

to the subsequent measurements. A more practical solution is often to measure the relative change in the instrument calibration factor as a function of temperature and then apply this as a correction factor to account for the temperature difference in the subsequent measurements.

34.6 PHOTOMETRY

Definition and Scope

Photometry is the measurement of radiation in a way that characterizes its effectiveness in stimulating the normal human visual system.^{4–5,80–82} Since visual sensation is a subjective experience, it is not directly quantifiable in absolute physical units. Attempts to quantify human visual sensation, therefore, were by comparison to various specified sources of light. The first sources used as standards were candles and later flames of prescribed construction. About the turn of the century, groups of incandescent lamps were selected as photometric standards and eventually a planckian radiator at a specific temperature was adopted by international agreement as the standard source. At present, the SI base unit for photometry, the candela, is no longer defined in terms of a given light source but is related to the radiant intensity by a multiplicative constant. Therefore, either an absolute source or detector can be used to establish an internationally recognized photometric calibration. Furthermore, there is no need for a human observer to effect a quantitative photometric measurement.

Photometry, as discussed here, is more precisely referred to as physical photometry to distinguish it from psychophysical photometry. Early photometric calibrations relied on human observers to compare an unknown light source to a standard. Presently, photometric calibrations are based on measurements using physical instruments. The instrument simulates human visual response either by having a detector with a spectral response that approximates that of the CIE standard observer or by using the CIE standard observer spectral response function in the data analysis.

Psychophysical photometry is the measurement of the effectiveness of light in individual observers and is more generally referred to as visual science. An individual's visual system may differ from that of the CIE standard observer defined for physical photometry, and these differences are sometimes important in experiments in visual science.

Photometry is restricted to the measurement of the magnitude of the visual sensation without regard to color, although it is well known that the perception of brightness is highly dependent on color in many circumstances. Measurement of the human response to color in terms of color matching is known as colorimetry. See Vol. III, Chap. 10, "Colorimetry."

Under reasonable light levels, the human eye can detect a difference of as little as 0.5 percent between two adjacent fields of illumination. For fields of illumination which are not adjacent, or are viewed at substantially different times, the eye can only detect differences of 10 to 20 percent. A discussion of the performance of the human visual system can be found in Vol. III, Chap. 2, "Visual Performance." Extensive treatments of photometry can be found in Walsh,⁸⁰ and Wyszecki and Stiles.⁸¹

Photopic, Scotopic, and Mesopic Vision

Electromagnetic radiation of sufficient power and in the wavelength range from about 360 to 830 nm, will stimulate the human visual system and elicit a response from an observer. The spectral range given here is the range over which measurements in physical photometry are defined. The range of reasonably perceptible radiation is usually given as about 400 to 700 nm. After light enters through the optical system of the eye—the cornea, iris or pupil, lens, and vitreous humor—the next stage of the visual response occurs in the retina. The retina contains two types of receptor cells: cones, which are the dominant sensors when the eye is adapted to higher radiance levels of irradiation (*photopic* vision), and rods, the dominant sensors at lower radiance levels (*scotopic* vision). Between the higher and lower levels of light adaptation is the region of *mesopic* vision, the range of radiance levels where both cones and rods contribute in varying degrees to the visual process.

Three types of cones having different spectral sensitivity functions exist in the normal human eye. The brain is able to distinguish colors by comparison of the signals from the three cone types. Of the three cone types, only the middle- and long-wavelength-sensitive cones contribute to the photopic sensation of the radiation entering the eye. The relative spectral sensitivity functions of the photopically and scotopically adapted human eye have been measured for a number of observers. From averages of these measurements, a set of values has been adopted by international agreement as the spectral efficiency for the CIE standard observer for photopic vision and another set for the CIE standard observer for scotopic vision (CIE, Commission Internationale de l'Eclairage). Because of the complexity of the spectral sensitivity of the eye at intermediate irradiation levels, there is no standard of spectral efficiency for mesopic vision. Values of the photopic and scotopic spectral efficiency functions are listed in Table 2.

TABLE 2 Photopic and Scotopic Spectral Luminous Efficiency Functions

Wavelength	Photopic	Scotopic	Wavelength	Photopic	Scotopic
375	0.00002	—	575	0.91540	0.1602
380	0.00004	0.00059	580	0.87000	0.1212
385	0.00006	0.00111	585	0.81630	0.0899
390	0.00012	0.00221	590	0.75700	0.0655
395	0.00022	0.00453	595	0.69490	0.0469
400	0.00040	0.00929	600	0.63100	0.03315
405	0.00064	0.01852	605	0.56680	0.02312
410	0.00121	0.03484	610	0.50300	0.01593
415	0.00218	0.0604	615	0.44120	0.01088
420	0.00400	0.0966	620	0.38100	0.00737
425	0.00730	0.1436	625	0.32100	0.00497
430	0.01160	0.1998	630	0.26500	0.00334
435	0.01684	0.2625	635	0.21700	0.00224
440	0.02300	0.3281	640	0.17500	0.00150
445	0.02980	0.3931	645	0.13820	0.00101
450	0.03800	0.455	650	0.10700	0.00068
455	0.04800	0.513	655	0.08160	0.00046
460	0.06000	0.567	660	0.06100	0.00031
465	0.07390	0.620	665	0.04458	0.00021
470	0.09098	0.676	670	0.03200	0.00015
475	0.11260	0.734	675	0.02320	0.00010
480	0.13902	0.793	680	0.01700	0.00007
485	0.16930	0.851	685	0.01192	0.00005
490	0.20802	0.904	690	0.00821	0.00004
495	0.25860	0.949	695	0.00572	0.00003
500	0.32300	0.982	700	0.00410	0.00002
505	0.40730	0.998	705	0.00293	0.00001
510	0.50300	0.997	710	0.00209	0.00001
515	0.60820	0.975	715	0.00148	0.00001
520	0.71000	0.935	720	0.00105	0.00000
525	0.79320	0.880	725	0.00074	0.00000
530	0.86200	0.811	730	0.00052	0.00000
535	0.91485	0.733	735	0.00036	0.00000
540	0.95400	0.650	740	0.00025	0.00000
545	0.98030	0.564	745	0.00017	0.00000
550	0.99495	0.481	750	0.00012	0.00000
555	1.00000	0.402	755	0.00008	0.00000
560	0.99500	0.3288	760	0.00006	0.00000
565	0.97860	0.2639	765	0.00004	0.00000
570	0.95200	0.2076	770	0.00003	0.00000

Photometric quantities can be calculated or measured as either photopic or scotopic quantities. Adaptation to luminance levels of $\geq 3 \text{ cd m}^{-2}$ (see further discussion) in the visual field usually leads to photopic vision, whereas adaptation to luminance levels of $\leq 3 \times 10^{-5} \text{ cd m}^{-2}$ usually leads to scotopic vision. Photopic vision is normally assumed in photometric measurements and photometric calculations unless explicitly stated to be otherwise.

Basic Concepts and Terminology

As noted in the earlier section on “Photometry,” the principles of photometry are the same as those for radiometry with the exception that the spectral responsivity of the detector is defined by general agreement to be specific approximations of the relative spectral response functions of the human eye. Photometric quantities are related to radiometric quantities via the spectral efficiency functions defined for the photopic and scotopic CIE standard observers.

Luminous Flux If physical photometry were to have been invented after the beginning of the twentieth century, then the physical basis of measurement might well have been the relationship between visual sensations and the energy of the photons and their flux density. It would follow naturally because vision is a photobiological process that is more closely related to the quantum nature of the radiation rather than its thermal heating effects. However, because of the weight of historical precedent, the basis of physical photometry is defined as the relationship between visual sensation and radiant power and its wavelength. The photometric equivalent of radiant power is luminous flux, and the unit that is equivalent to the watt is the lumen.

Luminous flux, Φ_v , is the quantity derived from spectral radiant power by evaluating the radiation according to its action upon the CIE standard observer.

$$\Phi_v = K_m \int \Phi_\lambda V(\lambda) d\lambda \quad (78)$$

where $V(\lambda)$ is the spectral efficiency function for photopic vision listed in Table 2, and K_m is the luminous efficacy for photopic vision. The spectral luminous efficacy is defined near the maximum, $\lambda_m = 555 \text{ nm}$, of the photopic efficiency function to be

$$K_m = 683 \frac{V(\lambda_m)}{V(555.016 \text{ nm})} \cong 683 \text{ lmW}^{-1} \quad (79)$$

The definitions are similar for scotopic vision

$$\Phi'_v = K'_m \int \Phi_\lambda V'(\lambda) d\lambda \quad (80)$$

where $V'(\lambda)$ is the spectral luminous efficiency function for scotopic vision listed in Table 2, and K'_m is the luminous efficacy for scotopic vision. The scotopic luminous efficiency function maximum occurs at $\lambda_m = 507 \text{ nm}$. The defining equation for K'_m is

$$K'_m = 683 \frac{V'(\lambda_m)}{V'(555.016 \text{ nm})} \cong 1700 \text{ lmW}^{-1} \quad (81)$$

The spectral shifts indicated in Eqs. (79) and (81) are required in order to obtain the precise values for the photopic and scotopic luminous efficacies. The magnitudes of the shifts follow from the specification of an integral value of frequency instead of wavelength in the definition of the SI base unit for photometry, the candela.

Luminous Intensity, Illuminance, and Luminance The candela, abbreviated cd, is defined by international agreement to be the luminous intensity in a given direction of a source that emits monochromatic radiation of frequency $540 \times 10^{12} \text{ Hz}$ (equal to 555.016 nm) and that has a radiant

intensity of $1/683 \text{ W sr}^{-1}$ in that direction. The spectral luminous efficacy of radiation at $540 \times 10^{12} \text{ Hz}$ equals 683 lm W^{-1} for all states of visual adaptation.

Because of the long history of using a unit of intensity as the basis for photometry, the candela was chosen as the SI base unit instead of the lumen, notwithstanding the fact that intensity is, strictly speaking, measurable only for point sources.

The functional form of the definitions of illuminance, luminous intensity, and luminance were presented in Eqs. (13), (14), and (15). The concepts are briefly reviewed here for the sake of convenience.

Luminous intensity is the photometric equivalent of radiant intensity, that is, luminous intensity is the luminous flux per solid angle. The symbol for luminous intensity is I_v . The unit for luminous intensity is the candela.

Illuminance is the photometric equivalent of irradiance, that is, illuminance is the luminous flux per unit area. The symbol for illuminance is E_v . The typical units for illuminance are lumens/meter².

Luminance is the photometric equivalent of radiance. Luminance is the luminous flux per unit area per unit solid angle. The symbol for luminance is L_v . The units for luminance are typically candelas/meter². In many older treatises on photometry, the term brightness is often taken to be equivalent to luminance, however, this is no longer the accepted usage.

In present usage, luminance and brightness have different meanings. In visual science (psychophysical photometry), two spectral distributions that have the same luminance typically do not have the same brightness. Operationally, spectral distributions of equal luminance are established with a psychophysical technique called heterochromatic flicker photometry. The observer views two spectral distributions that are rapidly alternated in time at the same spatial location, and the radiance of one is adjusted relative to the other to minimize the appearance of flicker. Spectral distributions of equal brightness are established with heterochromatic brightness matching, in which the two spectral distributions are viewed side-by-side and the radiance of one is adjusted relative to the other so that the fields appear equally bright. Though repeatable matches can easily be set with each technique, flicker photometric matches and brightness matches differ for many pairs of spectral distributions.

Photometric radiation transfer calculations and measurements are performed using the same methods and approximations that apply to the radiometric calculations discussed earlier. The exception, of course, is that the spectral sensitivity of the detector is specified.

Retinal Illuminance

In vision research it is frequently required to determine the effectiveness of a uniform, extended field of light (i.e., a large lambertian source that overfills the field of view of the eye) by estimating the illuminance on the retina. If it is assumed that the cornea, lens, and vitreous humor are lossless, then the luminous flux Φ_v in the image on the retina can be approximated from the conservation of the source luminance L_v as follows [see Eq. (22)],

$$L_v = \frac{L_e}{n_e^2} = \frac{\Phi_r S_{pr}^2}{n_e^2 A_r A_p} = E_r \frac{S_{pr}^2}{n_e^2 A_p} \quad (82)$$

where L_e is the radiance within the eye, n_e is the index of refraction of the ocular medium (the index of refraction of air is 1), Φ_r is the luminous flux at the retina, A_r is the area of the image of the retina, A_p is the area of the pupil, s_{pr} is the distance from the pupil to the retina, and E_r is the average illuminance in the image. Therefore, the average illuminance on the retina is

$$E_r = L_v \frac{n_e^2 A_p}{S_{pr}^2} \quad (83)$$

The luminance can, of course, be in units of either photopic or scotopic cd m^{-2} . The area of the pupil is measurable, but the distance between the pupil and retina is typically not available.

Therefore, a unit of retinal illuminance that avoids the necessity of determining this distance has been defined in terms of just the source luminance and pupil area. This unit is the troland, abbreviated td, and is defined as the retinal illumination for a pupil area of 1 mm^2 produced by a radiating surface having a luminance of 1 cd m^{-2} .

$$E_T = L_v A_p \quad (84)$$

Although it may be construed as an equivalent unit, one troland is *not equal* to one microcandela. The source is not a point but is infinite in extent. The troland is useful for relating several vision experiments where sources of differing luminance levels and pupil areas have been used.

The troland is, furthermore, not a measure of the actual illuminance level on the retina since the distance, index of refraction, and transmittance of the ocular medium are not included. For a schematic eye, which is designed to include many of the optical properties of the typical human eye, the effective distance between the pupil and the retina including the effect of the index of refraction is 16.7 mm^{83} (see also Vol. III, Chap. 10, "Colorimetry"). For the schematic eye with a 1-mm^2 pupil area, the effective solid angle at the retina is approximately 0.0036 sr . The retinal illuminance equivalent to one troland is therefore 0.0036 lm m^{-2} times the ocular transmittance.

Recall from the section on "Radiometric Effects of Stops and Vignetting" the effect of the aperture stop on the light entering an optical system; that is, all of the light entering the optical system appears to pass through the exit pupil, and the image of the aperture stop on the retina is the exit pupil. If the source is uniform and very large so that it overfills the field of view of the eye, the illuminance on the retina is independent of the distance between the source and the eye. If there is no intervening optic between the eye and the source, then the pupil is the aperture stop. However, if one uses an optical system to image the source into the eye, then the aperture stop need not be the pupil. An external optical system enables both the use of a more uniform, smaller source and the precise control of the retinal illumination by adjustment of an external aperture. The first configuration is called the newtonian view and the second is the maxwellian view of a source⁸⁰ (see also Vol. III, Chap. 5, "Optical Generation of the Visual Stimulus").

Though the troland is a very useful and a commonly used photometric unit among vision researchers, it should be interpreted with some caution in situations where one wishes to draw quantitative inferences about the effect of light falling on the retina. The troland is not a precise predictor because, besides not including transmission losses, no angular information is conveyed. The photoreceptors exhibit directional sensitivity where light entering through the center of the pupil is more effective than light entering through the pupil margin (the Stiles-Crawford effect, see Vol. III, Chap. 1, "Optics of the Eye"). Finally, specifying retinal illuminance in photometric units of trolands does not completely define the experimental conditions because the spectral distribution of the light on the retina is unspecified. Rather, the relative spectral responsivity of the eye (including the spectral dependence of the transmittance) is assumed by the inclusion of the $V(\lambda)$ or $V'(\lambda)$ functions. Experiments performed under mesopic conditions will be particularly prone to error.

If one measures the absolute spectral radiance of the light source, then Eq. (83) may be used in the radiometric form; that is, one substitutes $E_{r\lambda}$ and L_λ for E_r and L_v . The retinal spectral irradiance will then be in absolute units: $\text{W nm}^{-1}\text{m}^{-2}$.

Because the process of vision is a photobiological effect determined by the number and energy of the incident photons, the photon flux irradiance may be a more meaningful measure of the effect of the light on the retina. Using the radiometric form of Eq. (83) and the conversion to photon flux in Eq. (17), the photon flux irradiance $N_{E\lambda}$ on the retina is as follows.

$$N_{E\lambda} = 5.03 \times 10^{15} \lambda \tau_e(\lambda) L_\lambda \frac{n_e^2 A_p}{s_{pr}^2} = 1.80 \times 10^{13} \lambda \tau_e(\lambda) L_\lambda A_p \quad (85)$$

The ocular transmittance is included in this expression as $\tau_e(\lambda)$. The wavelength is in nm and, for radiance in units of $\text{W m}^{-2} \text{sr}^{-1} \text{nm}^{-1}$, the photon flux irradiance is in units of $\text{photons s}^{-1} \text{m}^{-2} \text{nm}^{-1}$.

If one uses a monochromatic light source, then a relationship between a monochromatic troland $E_T(\lambda)$ and the photon flux irradiance may be derived.

$$N_{E\lambda} = 1.80 \times 10^{13} \lambda A_p \tau_e(555) \frac{L_v(\lambda)}{K_m V(\lambda)} = 1.53 \times 10^{10} \lambda \frac{E_T(\lambda)}{V(\lambda)} \quad (86)$$

Only the transmittance at the peak of the $V(\lambda)$ curve, $\tau_e(555) = 0.58^{83}$ needs to be included since the spectral dependence of the ocular transmittance is already included in the $V(\lambda)$ function. The term $L_v(\lambda)$ is the luminance of a monochromatic light source.

An equivalent expression can be derived for the scotopic form of the monochromatic troland $E'_T(\lambda)$. Here, an ocular transmittance at 505 nm of 0.55^{83} has been used.

$$N_{E\lambda} = 5.82 \times 10^9 \lambda \frac{E'_T(\lambda)}{V'(\lambda)} \quad (87)$$

The reader is reminded that Eqs. (86) and (87) are valid only for a monochromatic source.

Absolute Photometric Calibrations

Photometric calibrations are in principle derived from the SI base unit for photometry, the candela. However, as one can see from the definitions of the candela and the other photometric quantities, photometric calibrations are in fact derived from absolute radiometric measurements using either a planckian radiator or an absolute detector. Typically, the relationship between illuminance and irradiance, Eq. (13), is used as the defining equation in deriving a photometric calibration.

The photometric calibration transfer devices available from national standards laboratories are usually incandescent lamps of various designs.⁸⁴ The photometric quantities commonly offered as calibrations are luminous intensity and total luminous flux.

The luminous intensity of a lamp, at a specified minimum distance and in a specified direction, is derived from a calibration of the spectral irradiance of the lamp, in the specified direction and at measured distance(s). The radiometric-to-photometric conversion [see Eq. (13)] is used to convert from spectral irradiance to illuminance. The inverse-square-law approximation, Eq. (24), is then used to derive luminous intensity.

Total luminous flux is a measure of all the flux emitted in every direction from a lamp. Total luminous flux is derived from illuminance (or luminous intensity) by measuring the flux emitted in all directions around the lamp. This procedure is known as goniophotometry. For an illuminance-based derivation, the total flux is the average of all the illuminance measurements times the surface area of the sphere described by the locus of the points at which the average illuminance was sampled. In the case of an intensity-based derivation, the total flux is the average of all the intensity measurements times 4π steradians. These are, in principle, the calculation methods for goniophotometry. In practice, the average illuminance (or intensity) is measured in a number of zones of fixed area (or solid angle) around the lamp. The product of the illuminance times the area of the zone (or the intensity times the solid angle of the zone) is the flux. The flux from each of the zones is then summed to obtain the total flux from the lamp.

A number of national standards laboratories provide luminance calibration transfer devices. These are typically in the form of a translucent glass plate that is placed at a specified distance and direction from a luminous intensity standard. One method of deriving the luminance calibration of the lamp/glass unit is to restrict the area of the glass plate with an aperture of known area. The intensity of the lamp/glass combination is then calibrated by comparison to an intensity standard lamp and the average luminance calculated by dividing the measured intensity by the area of the aperture.

Some national standards laboratories also offer calibrations of photometers,⁸⁵ also known as illuminance meters. A photometer is a photodetector that has been fitted with a filter to tailor its relative (peak normalized) spectral responsivity to match that of the CIE standard photopic observer.

Calibration of a photometer is usually obtained by reference to a luminous intensity standard positioned at a measured distance from the detector aperture. The inverse-square-law approximation is invoked to obtain the value of the illuminance at the measurement distance.

Other Photometric Terminology

Foot-candles, Foot-lamberts, Nits, etc The following units of illuminance are often used in photometry, particularly in older texts:

lux (abbreviation: lx) = lumen per square meter
 phot (abbreviation: ph) = lumen per square centimeter
 meter candle = lumen per square meter
 footcandle (abbreviation: fc) = lumen per square foot

One foot-candle = 0.0929 lux.

The following units of luminance are often used:

nit (abbreviation: nt) = candela per square meter
 stilb (abbreviation: sb) = candela per square centimeter

It is sometimes the practice, particularly in illuminating engineering, to express the luminance of an actual surface in any given direction in relation to the luminance of a uniform, diffuse, i.e., lambertian, source that emits one lumen per unit area into a solid angle of π steradians [see Eq. (30)]. This concept is one of relative luminance and its units (given following) are not equatable to the units of luminance. Furthermore, in spite of what may appear to be a similarity, this concept is not, strictly speaking, the photometric equivalent of the exitance of a source because, in its definition, the integral of the flux over the entire hemisphere is referenced. In other words, the equivalence to exitance is true only for a perfectly uniform radiance source. For all other sources it is the luminance in a particular direction divided by π . The units for luminance normalized to a lambertian source are:

1 apostilb (abbreviation: asb) = $(1/\pi)$ candela per square meter
 1 lambert (abbreviation: L) = $(1/\pi)$ candela per square centimeter
 1 foot-lambert (abbreviation: fL) = $(1/\pi)$ candela per square foot

Sometimes the total flux of a source is referred to as its *candlepower* or *spherical candlepower*. This term refers to a point source that uniformly emits in all directions, that is, into a solid angle of 4π steradians. Such a source does not exist, of course, so that the terminology is more precisely stated as the *mean spherical candle power*, which is the mean value of the intensity of the source averaged over the total solid angle subtended by a sphere surrounding the source.

Distribution Temperature, Color Temperature Distribution temperature is an approximate characterization of the spectral distribution of the visible radiation of a light source. Its use is restricted to sources having relative spectral outputs similar to that of a blackbody such as an incandescent lamp. The mathematical expression for evaluating distribution temperature is

$$\int_{\lambda_1}^{\lambda_2} \left[1 - \frac{\phi_x(\lambda)}{a\phi_b(\lambda, T)} \right]^2 d\lambda \Rightarrow \text{minimum} \quad (88)$$

where $\phi_x(\lambda)$ is the relative spectral radiant power distribution function of the test source, $\phi_b(\lambda, T)$ is the relative spectral radiant power distribution function of a blackbody at the temperature T , and a is an arbitrary constant. The limits of integration are the limits of visible radiation. Since distribution temperature is only an approximation, the exact values of the integration limits are arbitrary; typical limits are 400 and 750 nm. Values of a and T are adjusted simultaneously until the value of

the integral is minimized. The temperature of the best-fit blackbody function is the distribution temperature.

Color temperature and correlated color temperature are defined in terms of the perceived color of a source and are obtained by determining the chromaticity of the radiation rather than its relative spectral distribution. Because they are not related to physical photometry they are not defined in this section of the *Handbook*. These quantities do not provide information about the spectral distribution of the source except when the source has an output that closely approximates a blackbody. Although widely, and mistakenly, used in general applications to characterize the relative spectral radiant power distribution of light sources, color temperature and correlated color temperature relate only to the three types of cone cell receptors for photopic human vision and the approximate manner in which the human brain processes these three signals. As examples of incompatibility, there are the obvious differences between the spectral sensitivity of the human eye and physical receptors of visible optical radiation, e.g., photographic film, TV cameras. In addition, ambiguities occur in the ability of humans to distinguish the perceived color from different spectral distributions (metameric pairs). These ambiguities and the spectral sensitivity functions of the eye are not replicated by physical measurement systems. Caution must be exercised when using color temperature or correlated color temperature to predict the performance of a physical measurement system.

34.7 REFERENCES

1. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon, Oxford, 1980.
2. E. Wolf, "Coherence and Radiometry," *J. Opt. Soc. Am.* **68**:6 (1978).
3. J. Geist, *McGraw-Hill Encyclopedia of Science and Technology*, McGraw-Hill, New York, 1987, p. 156.
4. W. Blevin et al., "Principles Governing Photometry," *Metrologia* **19**:97 (1983).
5. "The Basis of Physical Photometry," 2d ed., *Commission Internationale de L'Eclairage Publ. No. 18.2*, Central Bureau of the CIE, Vienna, 1983.
6. R. W. Boyd, *Radiometry and the Detection of Optical Radiation*, Wiley, New York, 1983.
7. F. Grum and R. J. Becherer, *Optical Radiation Measurements*, Academic Press, New York, 1979.
8. F. Hengstberger, ed. *Absolute Radiometry*, Academic Press, New York, 1989.
9. C. L. Wyatt, *Radiometric Calibration: Theory and Methods*, Academic Press, New York, 1978.
10. A. C. Parr, R. U. Datla, and J. L. Gardner, eds., *Optical Radiometry*, Elsevier, Amsterdam, 2005.
11. E. L. Dereniak and D. G. Crowe, *Optical Radiation Detectors*, Wiley, New York, 1984, pp. 1–14, and Apps. A and C.
12. M. V. Klein and T. E. Furtak, *Optics*, 2d ed., Wiley, New York, 1986, pp. 203–222.
13. T. J. Quinn, *Temperature*, Academic Press, New York, 1983, pp. 284–363.
14. W. S. Smith, *Modern Optical Engineering*, 2d ed., McGraw-Hill, New York, 1990, pp. 135–136, 142–145, and 205–231.
15. M. E. Chahine, D. J. McCleese, P. W. Rosenkranz, and D. H. Staelin, in *Manual of Remote Sensing*, 2d ed., American Society of Photogrammetry, Falls Church, VA, 1983, pp. 172–179.
16. E. Hansen and L. D. Travis, "Light Scattering in Planetary Atmospheres," *Space Science Reviews* **16**:527 (1974).
17. Y. J. Kaufman, in Ghassem Asrar (ed.), *Theory and Applications of Optical Remote Sensing*, Wiley, New York, 1989, pp. 350–378.
18. H. Y. Wong, *Handbook of Essential Formulae and Data on Heat Transfer for Engineers*, Longman, London, 1977, pp. 89–128.
19. E. M. Sparrow and R. D. Cess, *Radiation Heat Transfer*, Brooks-Cole, Belmont, CA, 1966.
20. D. G. Goebel, "Generalized Integrating Sphere Theory," *Appl. Opt.* **6**:125 (1967).
21. "Le Système International d'Unités," 3d ed., Bureau International des Poids et Mesures, Sèvres, France, 1977.

22. B. N. Taylor and C. E. Kuyatt, *Guidelines for Evaluating and Expressing the Uncertainty of NIST Measurement Results*, NIST Technical Note 1297, U.S. Government Printing Office, Washington, D.C., 1994.
23. N. P. Fox, J. E. Martin, and D. H. Nettleton, "Absolute Spectral Radiometric Determination of The Thermodynamic Temperatures of The Melting/Freezing Points of Gold, Silver and Aluminium," *Metrologia* **28**:357 (1991).
24. H. P. Baltes and E. R. Hilf, *Spectra of Finite Systems*, Bibliographisches Institut Mannheim, Vienna, Zurich, 1976.
25. H. P. Baltes, "Deviations from the Stefan-Boltzmann Law at Low Temperatures," *Appl. Phys.* **1**:39 (1973).
26. H. P. Baltes, "Planck's Law for Finite Cavities and Related Problems," *Infrared Phys.* **16**:1 (1976).
27. E. M. Sparrow, L. U. Ulbers, and E. R. Eckert, "Thermal Radiation Characteristics of Cylindrical Enclosures," *J. Heat Transfer* **C84**:188 (1962).
28. B. A. Peavy, "A Note on the Numerical Evaluation of Thermal Radiation Characteristics of Diffuse Cylindrical and Conical Cavities," *J. Res. Natl. Bur. Stds.* **70C**:139 (1966).
29. R. E. Bedford and C. K. Ma, "Emissivities of Diffuse Cavities: Isothermal and Non-isothermal Cones and Cylinders," *J. Opt. Soc. Amer.* **64**:339 (1974).
30. R. E. Bedford and C. K. Ma, "Emissivities of Diffuse Cavities, II: Isothermal and Non-isothermal Cylindrocones," *J. Opt. Soc. Amer.* **65**:565 (1975).
31. R. E. Bedford and C. K. Ma, "Emissivities of Diffuse Cavities, III: Isothermal and Non-isothermal Double Cones," *J. Opt. Soc. Amer.* **66**:724 (1976).
32. J. C. de Vos, "Evaluation of the Quality of a Blackbody," *Physica* **20**:669 (1954).
33. T. J. Quinn, "The Calculation of the Emissivity of Cylindrical Cavities Giving near Blackbody Radiation," *Brit. J. Appl. Phys.* **18**:1105 (1967).
34. E. M. Sparrow and V. K. Johnson, "Absorption and Emission Characteristics of Diffuse Spherical Enclosures," *J. Heat Transfer* **C84**:188 (1962).
35. T. J. Quinn, "The Absorptivity of a Specularly Reflecting Cone for Oblique Angles of View," *Infrared Phys.* **21**:123 (1981).
36. T. J. Quinn and J. E. Martin, "Blackbody Source in the -50 to +200°C Range for the Calibration of Radiometers and Radiation Thermometers," *Appl. Opt.* **30**:4486 (1991).
37. E. F. Zalewski, J. Geist, and R. C. Willson, "Cavity Radiometer Reflectance," *Proc. of the SPIE* **196**:152 (1979).
38. J. Schwinger, "On the Classical Radiation of Accelerated Electrons," *Phys. Rev.* **75**:1912 (1949).
39. D. H. Tombouljian and P. L. Hartman, "Spectral and Angular Distribution of Ultraviolet Radiation from the 300-Mev Cornell Synchrotron," *Phys. Rev.* **102**:1423 (1956).
40. D. Lemke and D. Labs, "The Synchrotron Radiation of the 6-GeV DESY Machine as a Fundamental Radiometric Standard," *Appl. Opt.* **6**:1043 (1967).
41. N. P. Fox, P. J. Key, P. J. Riehle, and B. Wende, "Intercomparison between Two Independent Primary Radiometric Standards in the Visible and near Infrared: A Cryogenic Radiometer and the Electron Storage Ring BESSY," *Appl. Opt.* **25**:2409 (1986).
42. L. P. Boivin, "Calibration of Incandescent Lamps for Spectral Irradiance by Means of Absolute Radiometers," *Appl. Opt.* **19**:2771 (1980).
43. E. F. Zalewski and W. K. Gladden, "Absolute Spectral Irradiance Measurements Based on the Predicted Quantum Efficiency of a Silicon Photodiode," *Opt. Pura y Aplicada* **17**:133 (1984).
44. L. P. Boivin and A. A. Gaertner, "Realization of a Spectral Irradiance Scale in the Near Infrared at the National Research Council of Canada," *Appl. Opt.* **28**:6082 (1992).
45. T. J. Quinn and J. E. Martin, "A Radiometric Determination of the Stefan-Boltzmann Constant and Thermodynamic Temperatures between -40°C and +100°C," *Phil. Trans. Roy. Soc. London* **316**:85 (1985).
46. V. E. Anderson and N. P. Fox, "A New Detector-based Spectral Emission Scale," *Metrologia* **28**:135 (1991).
47. N. P. Fox, J. E. Martin, and D. H. Nettleton, "Absolute Spectral Radiometric Determination of the Melting/freezing Points of Gold, Silver and Aluminum," *Metrologia* **28**:357 (1991).
48. L. Jauniskis, P. Foukal, and H. Kochling, "Absolute Calibration of an Ultraviolet Spectrometer Using a Stabilized Laser and a Cryogenic Radiometer," *Appl. Opt.* **31**:5838 (1992).
49. F. Kurlbaum, "Über eine Methode zur Bestimmung der Strahlung in Absolutem Maass un die Strahlung des schwarzen Körpers zwischen 0 und 100 Grad," *Wied. Ann.* **65**:746 (1898).

50. K. Ångström, "The Absolute Determination of the Radiation of Heat with the Electrical Compensation Pyrheliometer, with Examples of the Application of this Instrument," *Astrophys. J.* **9**:332 (1899).
51. W. R. Blevin and W. J. Brown, "Development of a Scale of Optical Radiation," *Austr. J. Phys.* **20**:567 (1967).
52. R. C. Willson, "Active Cavity Radiometer Type V," *Appl. Opt.* **19**:3256 (1980).
53. L. P. Boivin and F. T. McNeely, "Electrically Calibrated Absolute Radiometer Suitable for Measurement Automation," *Appl. Opt.* **25**:554 (1986).
54. J. Geist and W. R. Blevin, "Chopper-Stabilized Radiometer Based on an Electrically Calibrated Pyroelectric Detector," *Appl. Opt.* **12**:2532 (1973).
55. R. J. Phelan and A. R. Cook, "Electrically Calibrated Pyroelectric Optical-radiation Detector," *Appl. Opt.* **12**:2494 (1973).
56. J. E. Martin, N. P. Fox, and P. J. Key, "A Cryogenic Radiometer for Absolute Radiometric Measurements," *Metrologia* **21**:147 (1985).
57. C. C. Hoyt and P. V. Foukal, "Cryogenic Radiometers and Their Application to Metrology," *Metrologia* **28**:163 (1991).
58. J. A. R. Samson, "Absolute Intensity Measurements in the Vacuum Ultraviolet," *J. Opt. Soc. Amer.* **54**:6 (1964).
59. F. M. Matsunaga, R. S. Jackson, and K. Watanabe, "Photoionization Yield and Absorption Coefficient of Xenon in the Region of 860–1022Å," *J. Quant. Spectrosc. Radiat. Transfer* **5**:329 (1965).
60. J. Geist, "Quantum Efficiency of the p–n Junction in Silicon as an Absolute Radiometric Standard," *Appl. Opt.* **18**:760 (1979).
61. J. Geist, W. K. Gladden, and E. F. Zalewski, "The Physics of Photon Flux Measurements with Silicon Photodiodes," *J. Opt. Soc. Amer.* **72**:1068 (1982).
62. E. F. Zalewski and J. Geist, "Silicon Photodiode Absolute Spectral Response Self-calibration," *Appl. Opt.* **19**:1214 (1980).
63. J. Geist, E. F. Zalewski, and A. R. Schaefer, "Spectral Response Self-calibration and Interpolation of Silicon Photodiodes," *Appl. Opt.* **19**:3795 (1980).
64. J. L. Gardner and W. J. Brown, "Silicon Radiometry Compared to the Australian Radiometric Scale," *Appl. Opt.* **26**:2341 (1987).
65. E. F. Zalewski and C. C. Hoyt, "Comparison Between Cryogenic Radiometry and the Predicted Quantum Efficiency of Silicon Photodiode Light Traps," *Metrologia* **28**:203 (1991).
66. E. F. Zalewski and C. R. Duda, "Silicon Photodiode Device with 100 Percent External Quantum Efficiency," *Appl. Opt.* **22**:2867 (1983).
67. N. P. Fox, "Trap Detectors and Their Properties," *Metrologia* **28**:197 (1991).
68. J. H. Walker, R. D. Saunders, J. K. Jackson, and D. A. McSparron, *Spectral Irradiance Calibrations*, National Bureau of Standards Special Publication No. 250-20, U. S. Government Printing Office, Washington, D.C., 1987.
69. J. H. Walker, R. D. Saunders, and A. T. Hattenburg, *Spectral Radiance Calibrations*, National Bureau of Standards Special Publication No. 250-1, U.S. Government Printing Office, Washington, D.C., 1987.
70. J. Z. Klose, J. M. Bridges, and W. R. Ott, *Radiometric Standards in the Vacuum Ultraviolet*, National Bureau of Standards Special Publication No. 250-3, U.S. Government Printing Office, Washington, D.C., 1987.
71. E. F. Zalewski, *The NBS Photodetector Spectral Response Calibration Transfer Program*, National Bureau of Standards Special Publication No. 250-17, U.S. Government Printing Office, Washington, D.C., 1987.
72. W. Budde, *Physical Detectors of Optical Radiation*, Academic Press, New York, 1983.
73. L. P. Boivin, "Some Aspects of Radiometric Measurements Involving Gaussian Laser Beams," *Metrologia* **17**:19 (1981).
74. L. P. Boivin, "Reduction of Diffraction Errors in Radiometry by Means of Toothed Apertures," *Appl. Opt.* **17**:3323 (1978).
75. W. R. Blevin, "Diffraction Losses in Photometry and Radiometry," *Metrologia* **6**:31 (1970).
76. C. L. Sanders, "A Photocell Linearity Tester," *Appl. Opt.* **1**:207 (1962).
77. C. L. Sanders, "Accurate Measurements of and Corrections for Non-linearities in Radiometers," *J. Res. Natl. Bur. Stand.* **A76**:437 (1972).
78. W. Budde, "Multidecade Linearity Measurements on Silicon Photodiodes," *Appl. Opt.* **18**:1555 (1979).

79. A. R. Schaefer, E. F. Zalewski, and J. Geist, "Silicon Detector Non-linearity and Related Effects," *Appl. Opt.* **22**:1232 (1983).
80. J. W. T. Walsh, *Photometry*, Dover, New York, 1965.
81. G. Wyszecki and W. S. Stiles, *Colour Science: Concepts and Methods*, Wiley, New York, 1967.
82. "Light as a True Visual Quantity: Principles of Measurement," *Commission Internationale de L'Eclairage Publ. No. 41*, Central Bureau of the CIE, Vienna, 1978.
83. E. N. Pugh, "Vision: Physics and Retinal Physiology," in R. C. Atkinson, R. J. Herrnstein, G. Lindsey, and R. D. Luce, (eds.), *Steven's Handbook of Experimental Psychology*, 2d ed., Wiley, New York, 1988, pp. 75–163.
84. R. L. Booker and D. A. McSparron, *Photometric Calibrations*, National Bureau of Standards Special Publication No. 250–15, U.S. Government Printing Office, Washington, D.C., 1987.
85. "Methods of Characterizing the Performance of Radiometers and Photometers," *Commission Internationale de L'Eclairage Publ. No. 53*, Central Bureau of the CIE, Vienna, 1982.

This page intentionally left blank.

DO NOT DUPLICATE

MEASUREMENT OF TRANSMISSION, ABSORPTION, EMISSION, AND REFLECTION

James M. Palmer*

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

35.1 GLOSSARY

A	Area
a	absorption
bb	blackbody
c_2	second radiation constant
E	irradiance
f	bidirectional scattering distribution function
i	internal
L	radiance
P	electrical power
R	reflectance factor
r	reflection
T	temperature
t	transmission
α	absorptance
α'	absorption coefficient
ε	emittance (emissivity)
θ, ϕ	angles
ρ	reflectance
σ	Stefan Boltzmann constant
τ	transmittance
Φ	power (flux)
Ω	projected solid angle

*Deceased.

35.2 INTRODUCTION AND TERMINOLOGY

When radiant flux is incident upon a surface or medium, three processes occur: transmission, absorption, and reflection. Figure 1 shows the ideal case, where the transmitted and reflected components are either specular or perfectly diffuse. Figure 2 shows the transmission and reflection for actual surfaces.

The symbols, units, and nomenclature employed in this chapter follow the established usage as defined in *ISO Standards Handbook 2*,¹ Cohen and Giacomo,² and Taylor.³ Additional general terminology applicable to this chapter is from ASTM,⁴ IES,⁵ IES,⁶ Drazil,⁷ and CIE.⁸ The prefix *spectral* is used

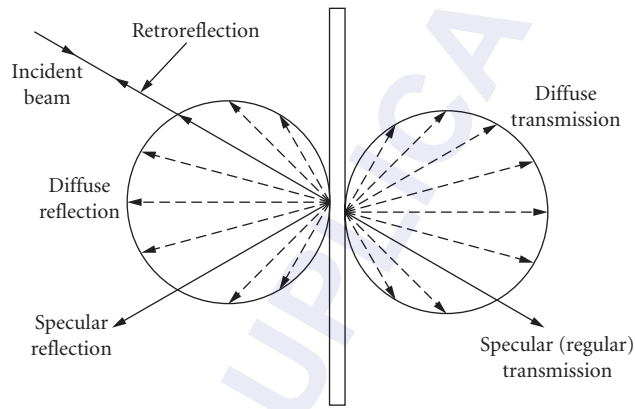


FIGURE 1 Idealized reflection and transmission.

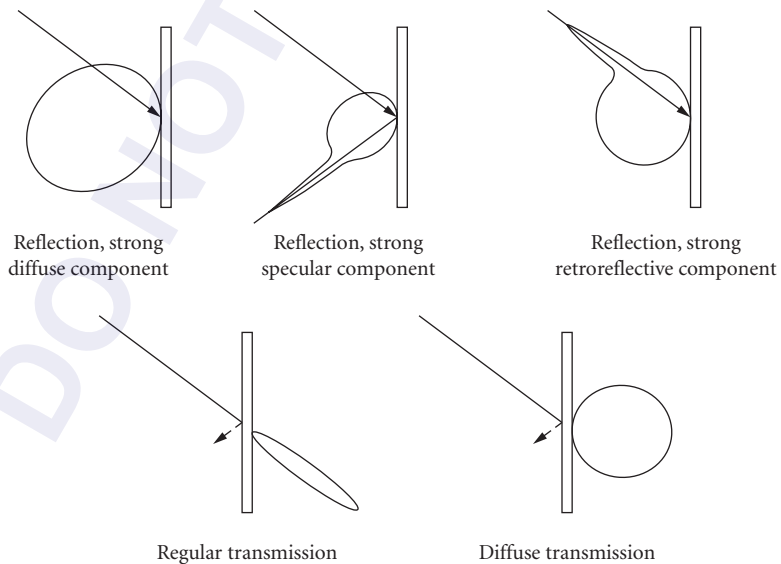


FIGURE 2 Actual reflection and transmission.

to denote a characteristic at a particular wavelength and is indicated by the symbol λ . The absence of the *spectral* prefix implies integration over all wavelengths with a source function included (omission of the source function is meaningless except where the characteristic is constant with wavelength).

There has been a continuing dialog over terminology, particularly between the suffixes *-ance* and *-ivity*.⁹⁻¹³ The suggested usage reserves terms ending with *-ivity* (such as transmissivity, absorptivity, and reflectivity) for properties of a pure material and employs the suffix *-ance* for the characteristics of a specimen or sample. For example, one can distinguish between the *reflectivity* of pure aluminum and the *reflectance* of a particular sample of 6061-T6 aluminum with a natural oxide layer. This distinction can be extended to differentiate between emissivity (of a pure substance) and emittance (of a sample). This usage of *emittance* should not be confused with the older term *radiant emittance*, now properly called *radiant exitance*. In this chapter, the suffix *-ance* will be used exclusively inasmuch as the measurement of radiometric properties of materials is under discussion.

35.3 TRANSMITTANCE

Transmission is the term used to describe the process by which incident radiant flux leaves a surface or medium from a side other than the incident side, usually the opposite side. The spectral transmittance $\tau(\lambda)$ of a medium is the ratio of the transmitted spectral flux $\Phi_{\lambda t}$ to the incident spectral flux $\Phi_{\lambda i}$, or

$$\tau(\lambda) = \frac{\Phi_{\lambda t}}{\Phi_{\lambda i}} \quad (1)$$

The transmittance τ is the ratio of the transmitted flux Φ_t to the incident flux Φ_i , or

$$\tau = \frac{\int_0^\infty \tau(\lambda) \Phi_{\lambda i} d\lambda}{\int_0^\infty \Phi_{\lambda i} d\lambda} \neq \int \tau(\lambda) d\lambda \quad (2)$$

Note that the integrated transmittance is *not* the integral over wavelength of the spectral transmittance, but must be weighted by a source function Φ_λ as shown.

The transmittance may also be described in terms of radiance as follows:

$$\tau = \frac{\int_0^\infty \int_{\Omega} L_{\lambda i} d\Omega_i d\lambda}{\int_0^\infty \int_{\Omega} L_{\lambda t} d\Omega_t d\lambda} \quad (3)$$

where $L_{\lambda i}$ represents the spectral radiance $L_{\lambda i}(\lambda; \theta_p, \phi_p)$ incident from direction (θ_p, ϕ_p) , $L_{\lambda t}$ represents the spectral radiance $L_{\lambda t}(\lambda; \theta_t, \phi_t)$ transmitted in direction (θ_t, ϕ_t) , and $d\Omega$ is the elemental projected solid angle $\sin \theta \cos \theta d\theta d\phi$.

The bidirectional transmittance distribution function (BTDF, symbol f_t) relates the transmitted radiance to the radiant incidence as

$$f_t(\lambda; \theta_i, \phi_i) \equiv \frac{dL_{\lambda t}}{dL_{\lambda i} d\Omega_i} = \frac{dL_{\lambda t}}{dE_{\lambda i}} (\text{sr}^{-1}) \quad (4)$$

Geometrically, transmittance can be classified as specular, diffuse, or total, depending upon whether the specular (regular) direction, all directions other than the specular, or all directions are considered.

35.4 ABSORPTANCE

Absorption is the process by which incident radiant flux is converted to another form of energy, usually heat. Absorptance is the fraction of incident flux that is absorbed. The absorptance α of an element is defined by $\alpha = \Phi_a / \Phi_i$. Similarly, the spectral absorptance $\alpha(\lambda)$ is the ratio of spectral power absorbed $\Phi_{\lambda a}$ to the incident spectral power $\Phi_{\lambda i}$,

$$\alpha = \frac{\int_0^{\infty} \alpha(\lambda) \Phi_{\lambda i} d\lambda}{\int_0^{\infty} \Phi_{\lambda i} d\lambda} \neq \int_{\lambda} \alpha(\lambda) d\lambda \quad (5)$$

An absorption coefficient α' (cm^{-1} or km^{-1}) is often used in the expression $\tau_i = e^{-\alpha' t}$, where τ_i is internal transmittance and t is pathlength (cm or km).

35.5 REFLECTANCE

Reflection is the process where a fraction of the radiant flux incident on a surface is returned into the same hemisphere whose base is the surface and which contains the incident radiation. The reflection can be specular (in the mirror direction), diffuse (scattered into the entire hemisphere), or a combination of both. Table 1¹⁴ shows a wide range of materials that have different goniometric (directional) reflectance characteristics.

TABLE 1 Goniometric Classification of Materials¹⁴

Material Classification	Scatter*	σ^\dagger	γ^\ddagger	Structure [§]	Example	
Exclusively reflecting materials	None	0	$\cong 0$	None	Mirror	
	Weak	≤ 0.4	$\leq 27^\circ$	Micro	Matte aluminum	
				Macro	Retroreflectors	
$\tau = 0$	Strong	> 0.4	$> 27^\circ$	None	Laquer & enamel coatings	
				Micro	Paint films, BaSO ₄ , Halon	
				Macro	Rough tapestries, road surfaces	
Weakly transmitting, strongly reflecting Materials	None	0	$\cong 0$	None	Sunglasses, color filters cold mirrors	
	Weak	≤ 0.4	$\leq 27^\circ$	Micro	Matte-surface color filters	
				Macro	Glossy textiles	
	Strong	> 0.4	$> 27^\circ$	None	Highly turbid glass	
$\tau \leq 0.35$				Micro	Paper	
				Macro	Textiles	
	Strongly transmitting materials	None	0	$\cong 0$	None	Window glass
		Weak	≤ 0.4	$\leq 27^\circ$	None	Plastic film
					Micro	Ground glass
	$\tau > 0.35$				Macro	Ornamental glass
					prismatic glass	
Strong		> 0.4	$> 27^\circ$	None	Opal glass	
			Micro	Ground opal glass		
			Macro	Translucent acrylic plastic with patterned surface		

*It is suggested that the diffusion factor is appropriate for strongly diffusing materials and that the half-angle is better for weakly diffusing materials.

[†] γ is a half-value angle, the angle from the normal where the radiance has dropped to one-half the value at normal.

[‡] σ is a diffusion factor, the ratio of the mean of radiance measured at 20° and 70° to the radiance measured at 5° from the normal, when the incoming radiation is normal. $\sigma = [L(20) + L(70)]/[2L(5)]$. It gives an indication of the spatial distribution of the radiance, and is unity for a perfect (Lambertian) diffuser.

[§]Structure refers to the nature of the surface. In a microscattering structure, the scatterers cannot be resolved with the unaided eye. The macrostructure scatterers can be readily seen.

The most general definition for reflectance ρ is the ratio of the radiant flux reflected Φ_r to the incident radiant flux Φ_i , or

$$\rho = \frac{\Phi_r}{\Phi_i} \quad (6)$$

Spectral reflectance is similarly defined at a specified wavelength λ as

$$\rho(\lambda) = \frac{\Phi_{\lambda r}}{\Phi_{\lambda i}} \quad (7)$$

(Spectral) reflectance factor (symbol R) is the ratio of (spectral) flux reflected from a sample to the (spectral) flux which would be reflected by a perfect diffuse (Lambertian) reflector.

No single descriptor of reflectance will suffice for the wide range of possible geometries. The fundamental geometric descriptor of reflectance is the bidirectional reflectance distribution function (BRDF, symbol f_r). It is defined as the differential element of reflected radiance dL_r in a specified direction per unit differential element of radiant incidence dE_i , also in a specified direction,¹⁵ and carries unit of sr^{-1} :

$$f_r(\theta_i, \phi_i, \theta_r, \phi_r) = \frac{dL_r(\theta_i, \phi_i; \theta_r, \phi_r; E_i)}{dE_i(\theta_i, \phi_i)} \quad [\text{sr}^{-1}] \quad (8)$$

The polar angle θ is measured from the surface normal and the azimuth angle ϕ is measured from an arbitrary reference in the surface plane, most often the plane containing the incident beam. The subscripts i and r refer to the incident and reflected beams, respectively.

By integrating over varying solid angles, Nicodemus et al.,¹⁵ based upon earlier work by Judd,¹⁶ defined nine goniometric reflectances, and by extension, nine goniometric reflectance factors. These are shown in Tables 2 and 3 and Fig. 3. In these tables, the term *directional* refers to a differential solid angle $d\omega$ in the direction specified by (θ, ϕ) . *Conical* refers to a cone of finite extent centered in direction (θ, ϕ) ; the solid angle ω of the cone must also be specified.

Details on these definitions and further discussion can be found in ASTM STP475,⁴ ASTM E808,¹⁷ Judd,¹⁶ Nicodemus,¹⁸ Nicodemus,¹⁹ and Nicodemus et al.¹⁵

TABLE 2 Nomenclature for Nine Types of Reflectance¹⁵

1. Bidirectional reflectance	$d\rho(\theta_i, \phi_i; \theta_r, \phi_r) = f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
2. Directional-conical reflectance	$\rho(\theta_i, \phi_i; \omega_r) = \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
3. Directional-hemispherical reflectance	$\rho(\theta_i, \phi_i; 2\pi) = \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
4. Conical-directional reflectance	$d\rho(\omega_i; \theta_r, \phi_r) = (d\Omega_r / \Omega_i) \int_{\omega_i} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
5. Biconical reflectance	$\rho(\omega_i; \omega_r) = (1/\Omega_i) \int_{\omega_i} \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
6. Conical-hemispherical reflectance	$\rho(\omega_i; 2\pi) = (1/\Omega_i) \int_{\omega_i} \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
7. Hemispherical-directional reflectance	$d\rho(2\pi; \theta_r, \phi_r) = (d\Omega_r / \pi) \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i$
8. Hemispherical-conical reflectance	$\rho(2\pi; \omega_r) = (1/\pi) \int_{2\pi} \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$
9. Bihemispherical reflectance	$\rho(2\pi; 2\pi) = (1/\pi) \int_{2\pi} \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r d\Omega_i$

TABLE 3 Nomenclature for Nine Types of Reflectance Factor¹⁵

1. Bidirectional reflectance factor	$R(\theta_i, \phi_i; \theta_r, \phi_r) = \pi f_r(\theta_i, \phi_i; \theta_r, \phi_r)$
2. Directional-conical reflectance factor	$R(\theta_i, \phi_i; \omega_r) = (\pi/\Omega_r) \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
3. Directional-hemispherical reflectance factor	$R(\theta_i, \phi_i; 2\pi) = \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_r$
4. Conical-directional reflectance factor	$R(\omega_i; \theta_r, \phi_r) = (\pi/\Omega_i) \int_{\omega_i} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i$
5. Biconical reflectance factor	$R(\omega_i; \omega_r) = [\pi/(\Omega_i \Omega_r)] \int_{\omega_i} \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i d\Omega_r$
6. Conical-hemispherical reflectance factor	$R(\omega_i; 2\pi) = (1/\Omega_i) \int_{\omega_i} \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i d\Omega_r$
7. Hemispherical-directional reflectance factor	$R(2\pi; \theta_r, \phi_r) = \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i$
8. Hemispherical-conical reflectance factor	$R(2\pi; \omega_r) = (1/\Omega_r) \int_{2\pi} \int_{\omega_r} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i d\Omega_r$
9. Bihemispherical reflectance factor	$R(2\pi; 2\pi) = (1/\pi) \int_{2\pi} \int_{2\pi} f_r(\theta_i, \phi_i; \theta_r, \phi_r) d\Omega_i d\Omega_r$

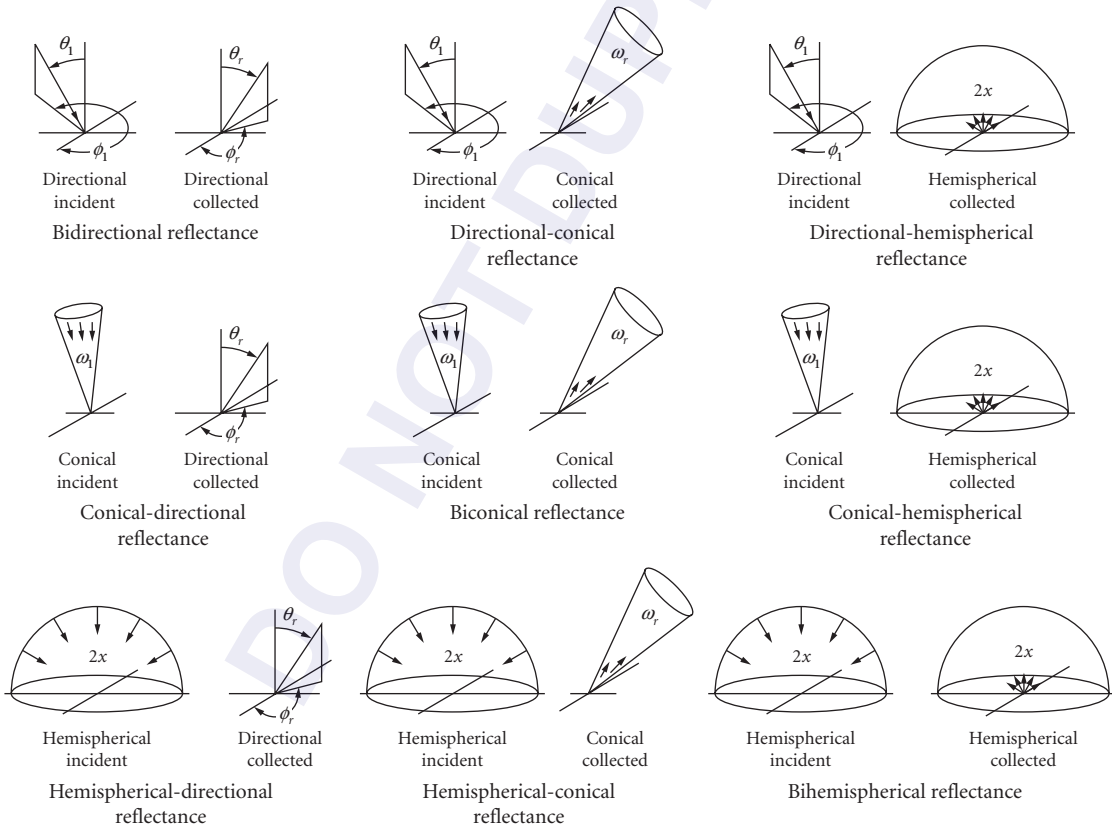


FIGURE 3 Nine geometrical definitions of reflectance.

35.6 EMITTANCE

Emittance (ε) is the ratio of the radiance of an object or surface to the radiance of a blackbody (planckian radiator) at the same temperature. It is therefore dimensionless and can assume values between 0 and 1 for thermal radiators at equilibrium. Spectral emittance $\varepsilon(\lambda)$ is the emittance at a given wavelength. If a radiator is neutral with respect to wavelength, with a constant spectral emittance less than unity, it is called a graybody.

$$\varepsilon = \frac{L}{L^{bb}}, \quad \varepsilon(\lambda) = \frac{L_\lambda}{L_\lambda^{bb}} \quad (9)$$

Directional emittance $\varepsilon(\theta, \phi)$ is defined by

$$\varepsilon(\theta, \phi) = \frac{L(\theta, \phi)}{L^{bb}} \quad (10)$$

Note that if the body is nongray, its emittance is dependent upon temperature inasmuch as the integral must be weighted by the source (Planck) function.

$$\varepsilon = \frac{\int_0^\infty \varepsilon(\lambda) L_\lambda^{bb} d\lambda}{\int_0^\infty L_\lambda^{bb} d\lambda} = \frac{1}{\pi} \frac{\int_0^\infty \varepsilon(\lambda) L_\lambda^{bb} d\lambda}{\sigma T^4} \quad (11)$$

35.7 KIRCHHOFF'S LAW

In a closed system at thermal equilibrium, conservation of energy necessitates that emitted and absorbed fluxes be equal. Since the radiation field in such a system is isotropic (the same in all directions), the directional spectral emittance and the directional spectral absorptance must be equal, i.e.,

$$\varepsilon(\lambda; \theta, \phi) = \alpha(\lambda; \theta, \phi) \quad (12)$$

This statement was first made by Kirchhoff.²⁰ Strictly, this equation holds for each orthogonal polarization component, and for it to be valid as written, the total radiation must have equal orthogonal polarization components. Kirchhoff's law is often simplified to the declaration $\alpha = \varepsilon$; however, this is not a universal truth; it may only be applied under a limited set of conditions. The geometrical and spectral averaging (integration) is governed by a specific set of rules as demonstrated by Siegel and Howell.²¹ Table 4, adapted from Siegel and Howell²¹ and Grum and Becherer,²² shows the various geometrical and spectral conditions under which the absorptance may be related to the emittance.

35.8 RELATIONSHIP BETWEEN TRANSMITTANCE, REFLECTANCE, AND ABSORPTANCE

Radiant flux incident upon a surface or medium undergoes transmission, reflection, and absorption. Application of conservation of energy leads to the statement that the sum of the transmission, reflection, and absorption of the incident flux is equal to unity, or

$$\alpha + \tau + \rho = 1 \quad (13)$$

TABLE 4 Summary of Absorptance-Emittance Relations²¹

Quantity	Equality	Required Conditions
Directional spectral	$\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$	None other than thermal equilibrium
Directional total	$\alpha(\theta, \phi, T_a) = \varepsilon(\theta, \phi, T_a)$	(1) Spectral distribution of incident energy proportional to blackbody at T_a , or (2) $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of wavelength
Hemispherical spectral	$\alpha(\lambda, T_a) = \varepsilon(\lambda, T_a)$	(1) Incident radiation independent of angle, or (2) $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of angle
Hemispherical total	$\alpha(T_a) = \varepsilon(T_a)$	(1) Incident energy independent of angle and spectral distribution proportional to blackbody at T_a , or (2) Incident energy independent of angle and $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of wavelength, or (3) Incident energy at each angle has spectral distribution proportional to blackbody at T_a and $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of angle, or (4) $\alpha(\lambda; \theta, \phi, T_a) = \varepsilon(\lambda; \theta, \phi, T_a)$ independent of angle and wavelength

In the absence of nonlinear effects (i.e., the Raman effect, etc.),

$$\alpha(\lambda) + \tau(\lambda) + \rho(\lambda) = 1 \quad (14)$$

If the situation is such that one of the above Kirchhoff-type relations is applicable, then emittance ε may be substituted for absorptance α in the previous equations, or

$$\varepsilon = 1 - \tau - \rho \quad \varepsilon(\lambda) = 1 - \tau(\lambda) - \rho(\lambda) \quad (15)$$

35.9 MEASUREMENT OF TRANSMITTANCE

A knowledge of the transmission of optical materials and elements, gaseous atmospheres, and various liquids is necessary throughout the realm of optics. Most of these measurements are made with commercial spectrophotometers. It is beyond the scope of this chapter to discuss the design and operation of spectrophotometric equipment except sample-handling practices. For further discussion, see Gram and Becherer,²² ASTM E275,²³ and ASTM E409.²⁴

Conventional spectrophotometers are of the double-beam configuration, where the output is the ratio of the signal in the sample beam to the signal in the reference beam plotted as a function of wavelength. It is incumbent upon the experimenter to ensure that the only difference between the two beams is the unknown. Therefore, if liquid or gas cells are employed, one should be placed in each beam. For gas cells, an equal amount of carrier gas should be injected into each cell, with the unknown to be sample placed in only one cell, destined for the sample beam. For liquids, an equal amount of solute should be placed in each cell. A critical issue with liquid and solid samples is the beam geometry. Most spectrophotometers feature converging beams in the sample space. If the

optical path length (the product of index of refraction and actual distance) for each beam is not identical, a systematic difference is presented to either the entrance slit or the detector. In addition, some specimens (e.g., interference filters) are susceptible to errors when measured in a converging beam. Most instruments also have a single monochromator which is susceptible to stray radiation, the limiting factor when trying to make measurements of samples that are highly absorbing in one spectral region and transmitting in another. Some recent instruments feature linear detector arrays along with single monochromators to allow the acquisition of the entire spectrum in several milliseconds; these are particularly applicable to reaction rate studies.

Conventional double-beam instruments are limited by these factors to uncertainties on the order of 0.1 percent. For lower uncertainties, the performance deficiencies found in double-beam instruments can largely be overcome by the use of a single-beam architecture. The mode of operation is sample-in–sample-out. If the source is sufficiently stable with time, the desired spectral range can be scanned without the sample, then rescanned with the sample in place. Otherwise, the spectrometer can be set at a fixed wavelength and alternate readings with and without the sample in place must be made. Care should be taken to ensure that the beam geometry is not altered between sequential readings.

To achieve the ultimate in performance from conventional spectrophotometry, several design characteristics should be included. A double monochromator is essential to minimize stray light. The beam geometry in the sample compartment should be highly collimated to avoid focus shifts with optically thick samples. Some form of beam integration, such as an integrating sphere or other diffuser, should be employed to negate the effects of nonuniform detectors and beam shifts. An exemplary instrument is the high-accuracy spectrophotometer developed by the National Institute for Standards and Technology (NIST), described in Mielenz and Eckerle,²⁵ Mielenz et al.,²⁶ Venable et al.,²⁷ Eckerle,²⁸ and Eckerle et al.²⁹ A particularly useful review is Eckerle et al.³⁰ Similar laboratory instruments have also been built elsewhere by Clarke,³¹ Freeman,³² and Zwinkles and Gignac.³³

Numerous other instruments have been described in the literature; some have been designed for singular or limited purposes while others have a more universal appeal. Use of integrating spheres is common, both for the averaging effects and for the isolation of the specular and diffuse components, as shown in Fig. 4.^{14,34} Several useful instruments are described by Karras,³⁵ Taylor,³⁶ Zerlaut and Anderson,³⁷ Clarke and Larkin,³⁸ and Kessell.³⁹

Conventional instruments lack a wide dynamic range because there are simply not enough photons available in a narrow bandpass in a reasonable time. Solutions include Fourier transform spectrometers with a large multiplex advantage, the use of tunable lasers, and heterodyne spectrometry.⁴⁰

Simple instruments can be purchased or constructed for specific purposes. For example, solar transmittance can be determined using either the natural sun (if available) or simulated solar radiation as the source. A limited degree of spectral isolation can be achieved with an abridged spectrophotometer using narrow bandpass interference or glass absorption filters.

Several publications have suggested methods for making accurate and repeatable measurements including Hughes,⁴¹ Mielenz,⁴² Venable and Hsia,⁴³ Burke and Mavrodineanu,⁴⁴ ASTM F768,⁴⁵ ASTM E971,⁴⁶ ASTM E903,⁴⁷ and ASTM E179.⁴⁸ Calibration and performance assessment of spectrophotometers includes photometric accuracy, linearity, stray light analysis, and wavelength calibration.

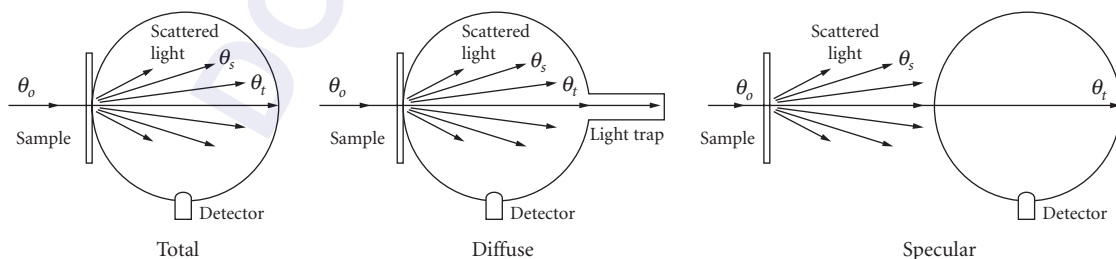


FIGURE 4 Measurement of total, diffuse, and specular transmittance using an integrating sphere.

Particular attention should be paid to luminescent samples that absorb radiant energy in one spectral region and re-emit it at longer wavelengths. Pertinent references include Hawes,⁴⁹ Bennett and Ashley,⁵⁰ ASTM E387,⁵¹ ASTM E275,²³ and ASTM, E409.²⁴

Standards of spectral transmittance are available as Standard Reference Materials from NIST. These take the form of metal-on-glass, metal-on-quartz, and solid glass filters. Most are used for verifying the photometric scale or for checking the wavelength calibration of a recording spectrophotometer. Descriptions of their development and use are given in Mavrodineanu and Baldwin,⁵² Mavrodineanu and Baldwin,⁵³ Eckerle et al.,³⁰ and Hsia.⁵⁴ The standardization laboratories of several countries occasionally conduct international intercomparisons of traveling standards. Recent intercomparisons have been reported in Eckerle et al.⁵⁵ and Fillinger and Andor.⁵⁶

35.10 MEASUREMENT OF ABSORPTANCE

In most cases, absorptance is not directly measured, but is inferred from transmission measurements, with appropriate corrections for reflection losses. These corrections can be calculated from the Fresnel equations if the surfaces are polished and the index of refraction is known. For materials where the absorption is extremely small, this method is unsatisfactory, as the uncertainties are dominated by the reflection contribution. In this case, direct measurements (such as laser calorimetry) must be made as discussed by Lipson et al.⁵⁷ and Hordvik.⁵⁸

35.11 MEASUREMENT OF REFLECTANCE

Instrumentation for the measurement of reflectance takes many forms. Only a few of the definitions for reflectance (Table 2) and reflectance factor (Table 3) have been adopted as standard configurations. The biconical configuration with small solid angles is most suited to a measurement of specular (regular, in the mirror direction) reflection. A simple reflectometer for the absolute measurement of specular reflectance was devised by Strong^{59,60} and is shown schematically in Fig. 5. Numerous detail improvements have been made on this fundamental design, including the use of averaging spheres. Designs range from simple⁶¹⁻⁶⁵ to complex.⁶⁶ Some reflectometers have been built specifically to measure at normal incidence.⁶⁷⁻⁶⁹ Measurement methods and data interpretation are also given in ASTM F768,⁴⁵ ASTM D523,⁷⁰ ASTM F1252,⁷¹ Hernicz and DeWitt,⁷² and Snail et al.⁷³

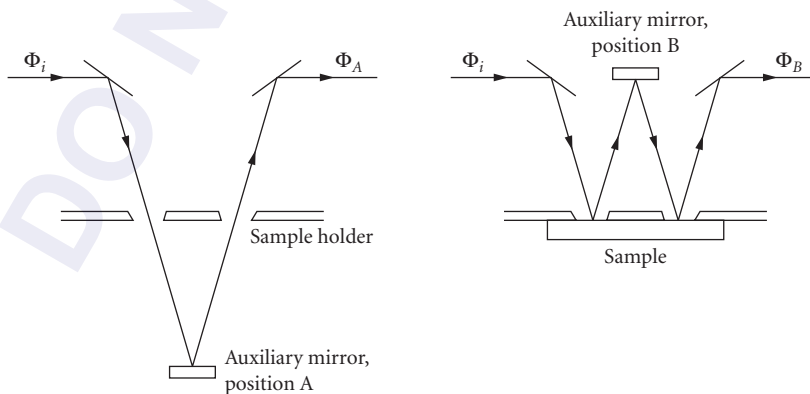


FIGURE 5 The Strong "VW" reflectometer.

The characterization of appearance of materials involves measurements of reflectance, both diffuse and specular. Numerous procedures and instruments have been devised for goniophotometry, the measurement of specular gloss with biconical geometry. Measurements are made at several angles from normal (20° , 30° , 45° , 60° , 75° , and 85°) depending upon the material under scrutiny. Further details can be found in ASTM C347,⁷⁴ ASTM E167,⁷⁵ ASTM D523,⁷⁰ ASTM E1349,⁷⁶ ASTM E179,⁷⁷ ASTM E430,⁷⁸ Erb,⁷⁹ and Hunter.⁸⁰

The measurement of diffuse reflectance can be accomplished using any one of the nine definitions from Table 2 and integrating where necessary. One could, for example, choose to measure the bidirectional reflectance distribution function (BRDF) as a function of incident beam parameters and to integrate over the hemisphere, but this would be a tedious process, and the large amount of data generated would be useful only to those involved with detailed materials properties research. Most practical measurements of diffuse reflectance involve the use of an integrating sphere. Several papers have discussed the general theory of the integrating sphere.^{81,82}

In the visible and near-IR spectral regions, the integrating sphere is the instrument of choice for both specular and diffuse specimens. Many papers have been written detailing instruments, methods, and procedures, some of which are shown in Fig. 6. The specular component of the reflected flux can be included to determine the total reflectance (Fig. 6a) or excluded to measure just the diffuse component (Fig. 6b). The angle of incidence can be varied by placing the sample at the center of the sphere (Fig. 6c), Edwards et al.⁸³ Others making contributions include Clarke and Compton,⁸⁴ Clarke and Larkin,⁸⁵ Dunkle,⁸⁶ Egan and Hilgeman,⁸⁷ Goebel et al.,⁸⁸ Hisdal,^{89,90} Karras,³⁵ McNicholas,⁹¹ Richter and Erb,⁹² Sheffer et al.,^{93,94} Taylor,⁹⁵ and Venable et al.²⁷ Some of these methods have been incorporated into standard methods and practices, such as ASTM C523,⁹⁶

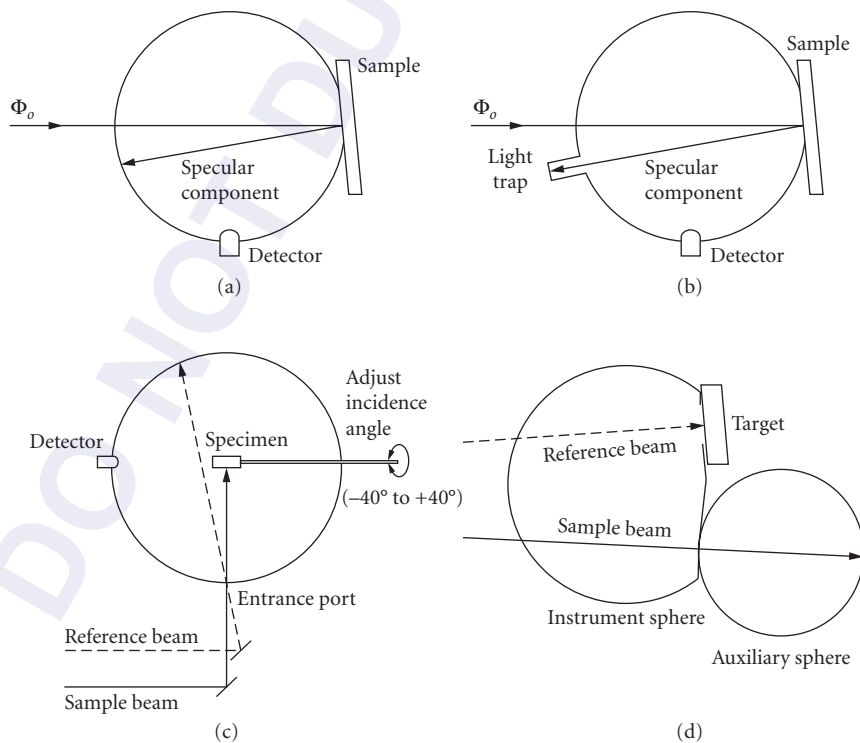


FIGURE 6 Measurement of diffuse reflectance using an integrating sphere.

ASTM E429,⁹⁷ ASTM E903,⁴⁷ CIE,⁹⁸ and IES.⁹⁹ Most integrating sphere measurements require reference to some form of an artifact standard, but the double-sphere method (Fig. 6d) produces absolute diffuse reflectance.^{37,38,100,101} Lindberg¹⁰² has demonstrated a method to scale relative measurements to absolute.

Alternative forms of hemispherical irradiation and/or collection have been described, several of which are shown in Fig. 7. Specular hemispherical (Fig. 7a), paraboloidal (Fig. 7b), and ellipsoidal (Fig. 7c) collectors have been used, particularly in those spectral regions where integrating sphere coatings are difficult to obtain.^{86,103–109} The Helmholtz reciprocity principle has been invoked to demonstrate the reversibility of the source and the collector.¹¹⁰ Hemispherical irradiation has also been employed by placing a cooled sample coplanar with the wall of a furnace as shown in Fig. 7d.^{86,111–113}

The procedures and instrumentation for the measurement of reflectance factor are identical with those for diffuse reflectance using the $0^\circ/45^\circ$ or $45^\circ/0^\circ$ geometry with annular, circumferential, or uniplanar illumination or viewing. Reference is made to a white reflectance standard characterized for reflectance factor, which must be compared with a perfect diffuse reflector. Pertinent references are ASTM E1349,⁷⁶ ASTM E97,¹¹⁴ ASTM E1348,¹¹⁵ Hsia and Weidner,¹¹⁶ and Taylor.^{36,117}

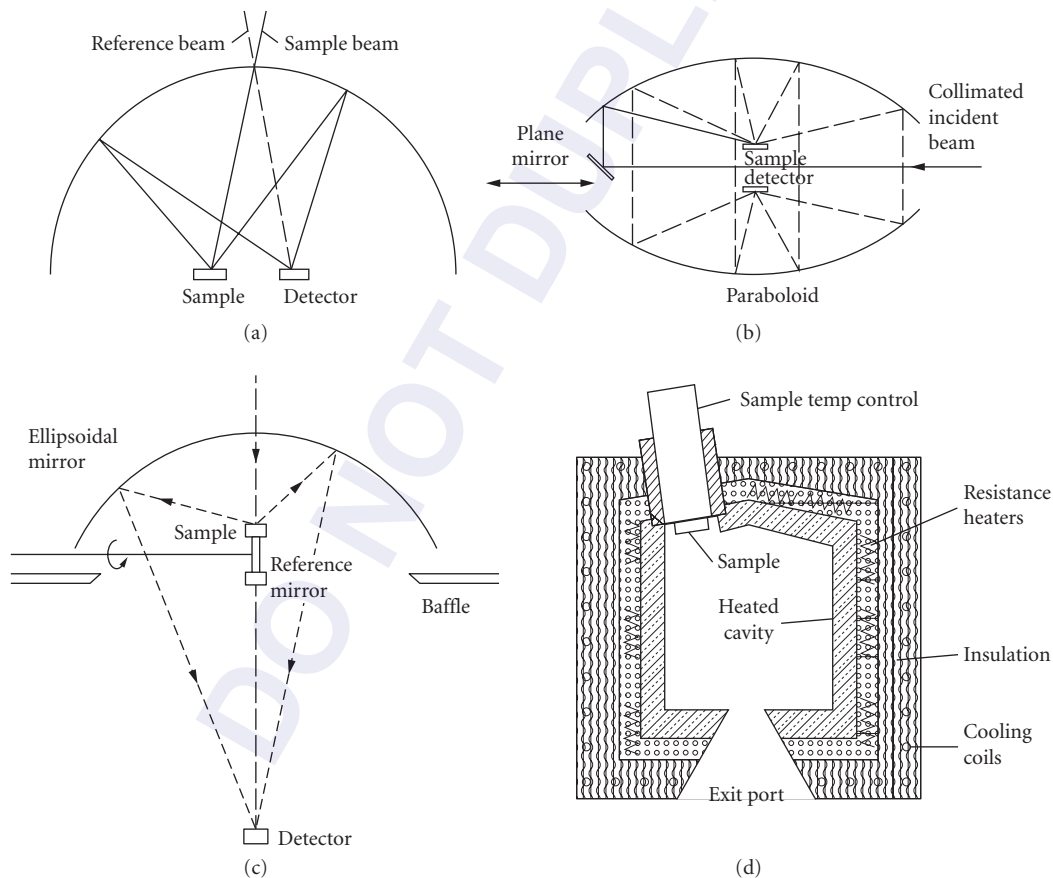


FIGURE 7 Measurement of diffuse reflectance using alternate methods.

Orbiting sensors measure the radiance of the earth-atmosphere system with some known geometry (generally nadir) and in well-defined wavelength bands. The quantity of interest is reflectance, as it is related to factors such as crop assessment, mineralization, etc. Corrections must be made for the atmospheric absorption, emission, and scattering and for the BRDF of the target. BRDF has been characterized in the field using the sun as the source, as described by Duggin.^{118,119}

Laboratory measurements of BRDF are made using goniometers where the sample-source angle and sample-receiver angle are independently adjustable. Coherent sources (lasers) are employed for the characterization of smooth specimens where a large source power is necessary for adequate SNR for off-specular angles and where speckle is not a concern. Incoherent sources (xenon arcs, black-body simulators, or tungsten-halogen) sources are often employed with spectral filters for more diffuse specimens. Similar measurements and techniques are employed to characterize bidirectional transmittance distribution function (BTDF) and bidirectional scattering distribution function (BSDF). For further information, see Asmail,¹²⁰ ASTM E1392,¹²¹ and Bartell et al.¹²²

The measurement of retroreflection poses the situation that the return beam coincides with the incident beam. The usual solution is to employ a beam splitter in the system, allowing the incident beam to pass and the return beam to be reflected. This immediately imposes both a significant loss in flux and the situation where the beam reaching the sample is partially polarized, unless a non-polarizing beamsplitter is employed. In addition, it is imperative that the reflected component of the incident beam be well-trapped, as the radiometer is looking in the same direction. The special vocabulary for retroreflection is given in CIE¹²³ and ASTM E808.¹⁷ Test methods are given in ASTM E810,¹²⁴ ASTM E809,¹²⁵ and Venable and Johnson.¹²⁶ An instrument specifically designed to measure retroreflection is detailed in Eckerle et al.¹²⁷

Most measurements of reflectance are relative and require artifact standards. Specular reflectance measurements are occasionally made using the absolute technique shown in Fig. 5, but are more commonly done using a simple reflectance attachment for a commercial spectrophotometer, requiring a calibrated reference. Freshly deposited metallic films have been used, with the assumption that an individual coating is the same as accepted data as shown in Table 5 for Al and Au.¹²⁸⁻¹³³ Standards for specular reflectance are also available.¹³⁵⁻¹³⁷

Diffuse reflectance standards exist in several forms. Ideally, a perfect ($\rho = 1$, lambertian) diffuser would be used, particularly for measuring reflectance factor.¹³⁸⁻¹⁴⁰ Certain materials approach the ideal over a limited angular and wavelength range. Historically, MgO was used in the visible spectrum.^{141,111} It was replaced first by BaSO₄¹⁴² and more recently by PTFE^{143,144} and ASTM E259.¹⁴⁵ Table 5 also shows a typical 6°/hemispherical reflectance factor for PTFE.¹⁴⁶ This fine, white powder, when pressed to a density of 1 g/cm³, is close to ideal over a wide spectral range. It is not quite lambertian, showing a falloff of BRDF at angles far removed from the specular direction¹⁴⁷ and exhibiting a slight amount of retroreflection. It may also be slightly luminescent when excited by far-ultraviolet.¹⁴⁸

In the infrared, two materials have proven useful. Flowers of sulfur¹⁴⁹⁻¹⁵¹ is suitable over the spectral range 1 to 15 μm . Gold is highly reflective and very stable. To be useful as a diffuse reflectance standard, it must be placed on top of a lambertian surface. Several substrates for gold have been suggested, including sandpaper¹⁵² and flame-sprayed aluminum.

PTFE is a satisfactory laboratory standard but is not well-suited for field use as it is not particularly rugged and is highly adsorbant and therefore subject to contamination. Several solutions have been proposed for working standards, including Eastman integrating sphere paint (BaSO₄), Vitriolite tile, and the Russian MS20 and MS14 opal glasses. These materials are, in general, more rugged, stable, and washable than PTFE. Further details can be found in CIE¹⁵³ and Clarke et al.¹⁵⁴

Discussion on the fabrication, calibration, and properties of various diffuse reflectance standards can be found in ASTM E259,¹⁵⁵ Budde,^{156,157} Egan and Hilgeman,¹⁵⁸ Fairchild and Daoust,¹⁴⁷ Morren et al.,¹⁵⁹ TAPPI,¹⁶⁰ and Weidner.^{161,162} International intercomparisons of laboratory standards of diffuse reflectance have been reported in Budde et al.,¹⁶³ and Weidner and Hsia,¹⁴⁶ and IES.⁹⁹

There are no standards available for retroreflection. However, NIST offers a Measurements Assurance Program (MAP) to enable laboratories to make measurements consistent with other national standards.

TABLE 5 Reflectance Standards

Wavelength (nm)	Aluminum	Gold	PTFE (6°/hemi)
250		0.295	0.973
300	0.921	0.346	0.984
350	0.921	0.330	0.990
400	0.919	0.360	0.993
450	0.918	0.358	0.993
500	0.916	0.453	0.994
550	0.916	0.800	0.994
600	0.912	0.906	0.994
650	0.906	0.947	0.994
700	0.898	0.963	0.994
750	0.886	0.970	0.994
800	0.868	0.973	0.994
850	0.868	0.973	0.994
900	0.891	0.974	0.994
950	0.924	0.974	0.994
1000	0.940	0.974	0.994
1100		0.975	0.994
1200	0.964	0.975	0.993
1300		0.975	0.992
1400		0.975	0.991
1500	0.974	0.975	0.992
1600		0.975	0.992
1700		0.976	0.990
1800		0.976	0.990
1900		0.976	0.985
2000	0.978	0.976	0.981
2100		0.976	0.968
2200		0.976	0.977
2300		0.976	0.972
2400		0.976	0.962
2500	0.979	0.977	0.960

35.12 MEASUREMENT OF EMITTANCE

Measurements of emittance can be done in several ways. The most direct method involves forming a material into the shape of a cavity in such a way that near-blackbody radiation is emitted. A measurement then compares the radiation from a location within the formed cavity to radiation from a flat, outside surface of the material, presumably at the same temperature.¹⁶⁴ The cavity can take the form of a cylinder, cone, or sphere. Similarly, a small-diameter, deep hole can be drilled into a specimen and radiation from the surface compared to radiation from the hole. Care must be taken that the specimen is isothermal and that the reflected radiation is considered. The definitive measurements of several materials, such as tungsten,¹⁶⁵ were determined in this fashion. The significant advantage in this direct method is that it is relative, depending on neither absolute radiometry or thermometry, but only requiring that the radiometer or spectroradiometer be linear over the dynamic range of the measurement. This linearity is also determinable by relative measurements.

If a variable-temperature blackbody simulator and a suitable thermometer are available, the specimen can be heated to the desired temperature T_s and the blackbody simulator temperature T_{bb} can be adjusted such that its (spectral) radiance matches that of the specimen. Then the (spectral)

emittance is calculable using the following equations for spectral emittance $\varepsilon(\lambda)$ and emittance ε (for a graybody only).

$$\varepsilon(\lambda) = \frac{e^{c_2/\lambda T_s} - 1}{e^{c_2/\lambda T_{bb}} - 1}, \quad \varepsilon = \frac{T_{bb}^4}{T_s^4} \quad (16)$$

If an absolutely calibrated radiometer and a satisfactory thermometer are available, a direct measurement can be made, as L_b is calculable if the temperature is known. Again, the reflected radiation must be considered.

Simple “inspection meter” techniques have been developed, and instrumentation is commercially available to determine the hemispherical emittance over a limited range of temperatures surrounding ambient. These instruments provide a single number, as they integrate both spatially and spectrally. A description of the technique can be found in ASTM E408.¹⁶⁶

Measurements of spectral emittance are most often made using spectral reflectance techniques, invoking Kirchhoff’s law along with the assumption that the transmittance is zero. A review of early work is found in Dunn et al.,¹⁶⁷ and Millard and Streed.¹⁶⁸ The usual geometry of interest is the directional-hemispherical. This can be achieved by either hemispherical irradiation-directional collection, or, using Helmholtz reciprocity,¹¹⁰ directional irradiation-hemispherical collection. Any standard reflectometry technique is satisfactory.

A direct method for the measurement of total (integrated over all wavelengths) hemispherical emittance is to use a calorimeter as shown in Fig. 8. A heated specimen is suspended in the center of a large, cold, evacuated chamber. The vacuum minimizes gaseous conduction and convection. If the sample suspension is properly designed, the predominant means of heat transfer is radiation. The chamber must be large to minimize the configuration factor between the chamber and the specimen. The chamber is cooled to T_c to reduce radiation from the chamber to the specimen. The equation used to determine emittance ε is

$$\varepsilon = \frac{P}{\sigma A(T_s^4 - T_c^4)} \quad (17)$$

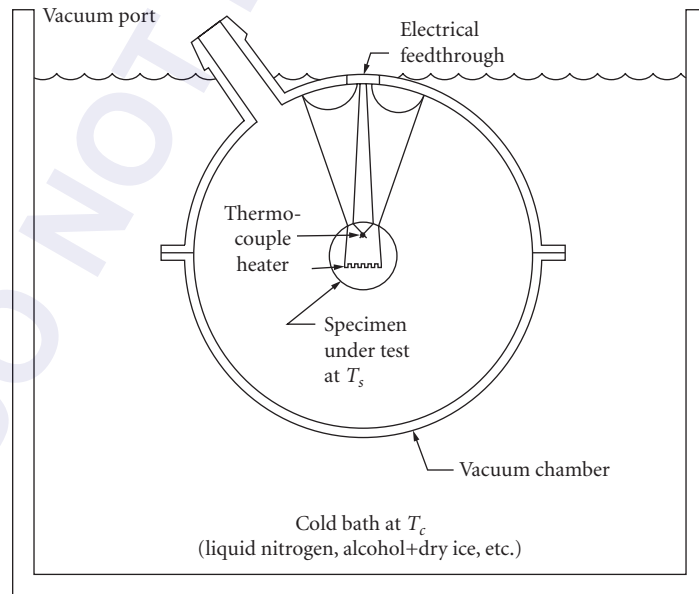


FIGURE 8 Calorimetric measurement of total hemispherical emittance.

where P is the power input to the specimen heater necessary to maintain an equilibrium specimen temperature T_s and A is the specimen area. The equation has been simplified with the aid of the following assumptions: (1) no thermal conduction from the specimen to the chamber, (2) no convective losses, (3) equilibrium has been achieved, and (4) the specimen area is much less than the chamber area. The power can be supplied electrically by means of a known heater or optically via a window in the chamber. In the latter case, a direct measurement of the ratio of solar absorptance α_s to thermal emittance ϵ_T can be directly obtained if the optical source simulates solar radiation. By varying the input power, the emittance can be determined as a function of temperature. There are numerous small corrections to account for geometry, lead conduction, etc. Details can be found in ASTM C835,¹⁶⁹ ASTM E434,¹⁷⁰ Edwards,¹⁷¹ and Richmond and Harrison.¹⁷²

Several attempts have been made to define and characterize artifact standards of spectral emittance for direct measurements.^{10,11} These were specimens of a thermally stable metal (i.e., Inconel) which were calibrated for emittance as a function of wavelength at several temperatures. No such standards are currently available. Interlaboratory comparisons have been made and reported.¹⁷³

Special problems include measurements at cryogenic temperatures¹⁷⁴ and effects of partially transparent materials.¹⁷⁵ Some additional references relating to emittance and its measurement are ASTM E307,¹⁷⁶ ASTM E423,¹⁷⁷ Clarke and Larkin,⁸⁵ DeWitt,¹⁷⁸ DeWitt and Richmond,¹⁷⁹ Hornbeck,¹⁸⁰ Millard and Streed,¹⁶⁸ Redgrove,¹⁸¹ Sparrow et al.,¹⁶⁴ Stierwalt,¹⁸² and Wittenberg.¹⁸³

35.13 REFERENCES

1. ISO, *Units of Measurement*, ISO Standards Handbook 2, International Organization for Standardization, Geneva, 1982.
2. E. R. Cohen and P. Giacomo, *Symbols, Units, Nomenclature and Physical Constants in Physics*, Document IUPAP-25, International Union of Pure and Applied Physics, 1987.
3. B. N. Taylor, "The International System of Units (SI)," *NIST Special Publication 330*, National Institute of Standards and Technology, Washington, D.C., 1991.
4. ASTM, "Nomenclature and Definitions Applicable to Radiometric and Photometric Characteristics of Matter," *ASTM Special Technical Publication 475*, ASTM, Philadelphia (1971).
5. IES Nomenclature Committee, "Proposed American National Standard Nomenclature and Definitions for Illuminating Engineering (proposed revision of Z7.1-R-1973)," *J. Illum. Eng. Soc.* **8**:2 (1979).
6. IES Nomenclature Committee, *American National Standard Nomenclature and Definitions for Illuminating Engineering*, ANSI/IES RP-16-1986. Illuminating Engineering Society of North America, N.Y., 1986.
7. J. V. Drazil, *Quantities and Units of Measurement: A Dictionary and Handbook*, Mansell, London, (1983).
8. CIE, "International Lighting Vocabulary," *CIE Publ. 17.4*, CIE, Paris (1987).
9. A. G. Worthing, "Temperature Radiation Emissivities and Emittances," in *Temperature, Its Measurement and Control in Science and Industry*, Reinhold, N.Y., p. 1164 (1941).
10. J. C. Richmond, "Physical Standards of Emittance and Reflectance," in H. H. Blau and H. Fischer (eds), *Radiative Transfer from Solid Material*, Macmillan, N.Y., 1962.
11. J. C. Richmond, W. N. Harrison, and F. J. Shorten, "An Approach to Thermal Emittance Standards, in J. C. Richmond (ed.), *Measurement of Thermal Radiation Properties of Solids*, NASA SP-31, NASA, Washington, D.C., 1963.
12. W. L. Wolfe, "Proclivity for Emissivity," *Appl. Opt.* **21**:1 (1982).
13. J. C. Richmond, J. J. Hsia, V. R. Weidner, and D. B. Wilmering, *Second-Surface Mirror Standards of Spectral Specular Reflectance*, NBS Special Publication SP260-79, U.S. National Bureau of Standards, Washington, D.C., 1982.
14. CIE, "Radiometric and Photometric Characteristics of Materials and their Measurement," *CIE Publication 38*, CIE, Paris (1977).
15. F. E. Nicodemus, J. C. Richmond, and J. J. Hsia, *Geometrical Considerations and Nomenclature for Reflectance*, NBS Monograph 160, U.S. National Bureau of Standards, Washington, D.C., 1977.
16. D. B. Judd, "Terms, Definitions and Symbols in Reflectometry," *J. Opt. Soc. Am.* **57**:445 (1967).

17. ASTM, "Terminology for Retroreflection and Retroreflectors," *ASTM E808*, ASTM, Philadelphia (1981).
18. F. E. Nicodemus, "Directional Reflectance and Emissivity of an Opaque Surface," *Appl. Opt.* **4**:767 (1965).
19. F. E. Nicodemus, "Reflectance Nomenclature and Directional Reflectance and Emissivity," *Appl. Opt.* **9**:1474 (1970).
20. G. Kirchhoff, "On the Relation between the Radiating and Absorbing Powers of Different Bodies for Light and Heat," *Phil. Mag.* **20**:1 (1860).
21. R. Siegel, and J. R. Howell, *Thermal Radiation Heat Transfer*, 2d ed., Hemisphere, N.Y., p. 63 (1981).
22. F. Grum, and R. J. Becherer, "Radiometry," in *Optical Radiation Measurements*, vol. 1, Academic, N.Y., p. 115 (1979).
23. ASTM, "Practice for Describing and Measuring Performance of UV/VIS/Near-IR Spectrophotometers," *ASTM E275*, ASTM, Philadelphia (1989).
24. ASTM, "Procedure for Description and Performance in the Spectrophotometer," *ASTM E409*, ASTM, Philadelphia (1990).
25. K. D. Mielenz, and K. L. Eckerle, *Design, Construction, and Testing of a New High Accuracy Spectrophotometer*, NBS Tech Note 729, U.S. National Bureau of Standards, Washington, D.C., 1972.
26. K. D. Mielenz, K. L. Eckerle, R. P. Madden, and J. Reader, "New Reference Spectrophotometer," *Appl. Opt.* **12**:1630 (1973).
27. W. H. Venable, J. J. Hsia, and V. R. Weidner, *Development of an NBS Reference Spectrophotometer for Diffuse Transmittance and Reflectance*, NBS Tech Note TN594-11, U.S. National Bureau of Standards, Washington, D.C., 1976.
28. K. L. Eckerle, *Modification of an NBS Reference Spectrophotometer*, NBS Technical Note TN913, U.S. National Bureau of Standards, Washington, D.C., 1976.
29. K. L. Eckerle, V. R. Weidner, J. J. Hsia, and Z. W. Chao, *Extension of a Reference Spectrophotometer into the Near Infrared*, NBS Technical Note TN1175, U.S. National Bureau of Standards, Washington, D.C., 1983.
30. K. L. Eckerle, J. J. Hsia, K. D. Mielenz, and V. R. Weidner, *Regular Spectral Transmittance*, NBS Special Publication SP250-6, U.S. National Bureau of Standards, Washington, D.C., 1987.
31. F. J. J. Clarke, "High-Accuracy Spectrophotometry at the National Physical Laboratory," *J. Res. Natl. Bur. Stand.* **A76**:375 (1972).
32. G. H. C. Freeman, "The New Automated Reference Spectrophotometer at NPL," in C. Burgess and K. D. Mielenz (eds.), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1986, p. 69.
33. J. C. Zwinkles and D. S. Gignac, "Design and Testing of a New High-Accuracy Ultraviolet-Visible- Near-Infrared Spectrophotometer," *Appl. Opt.* **31**:1557 (1992).
34. A. Roos, "Interpretation of Integrating Sphere Signal Output for Nonideal Transmitting Samples," *Appl. Opt.* **30**:468 (1991).
35. E. Karras, "The Use of the Ulbricht Sphere in Measuring Reflection and Transmission Factors," *J Opt. Soc. Am.* **11**:96 (1921).
36. A. H. Taylor, "A Simple Portable Instrument for the Absolute Measurement of Reflection and Transmission Factors," *Sci. Papers Bur. Standards* **17**:1 (1922).
37. G. A. Zerlaut and T. E. Anderson, "Multiple-Integrating Sphere Spectrophotometer for Measuring Absolute Spectral Reflectance and Transmittance," *Appl. Opt.* **20**:3797 (1981).
38. F. J. J. Clarke and J. A. Larkin, "Measurement of Total Reflectance, Transmittance and Emissivity over the Thermal IR Spectrum," *Infrared Phys.* **25**:359 (1985).
39. J. Kessell, "Transmittance Measurements in the Integrating Sphere," *Appl. Opt.* **25**:2752 (1986).
40. A. L. Migdall, B. Roop, Y. C. Zheng, J. E. Hardis, and G. J. Xia, "Use of Heterodyne Detection to Measure Optical Transmittance over a Wide Range," *Appl. Opt.* **29**:5136 (1990).
41. H. K. Hughes, "Beer's Law and the Optimum Transmittance in Absorption Measurements," *Appl. Opt.* **2**:937 (1963).
42. K. D. Mielenz, "Physical Parameters in High-Accuracy Spectrophotometry," in R. Mavrodineanu, J. I. Schultz, and O. Menis (eds), *Accuracy in Spectrophotometry and Luminescence Measurements*, NBS SP378, U.S. National Bureau of Standards, Washington, D.C., 1973.
43. W. H. Venable and J. J. Hsia, *Describing Spectrophotometric Measurements*, NBS Technical Note TN594-9, U.S. National Bureau of Standards, Washington, D.C., 1974.

44. R. W. Burke and R. Mavrodineanu, *Standard Reference Material: Accuracy in Analytical Spectrophotometry*, NBS Special Publication SP260-81, U.S. National Bureau of Standards, Washington, D.C., 1983.
45. ASTM, "Specular Reflectance/Transmittance of Optically Flat Coated/Non-coated Specimens," *ASTM F768*, ASTM, Philadelphia (1987).
46. ASTM, "Test for Photometric Transmittance/Reflectance of Materials to Solar Radiation," *ASTM E971*, ASTM, Philadelphia (1988).
47. ASTM, "Standard Test Method for Solar Absorptance, Reflectance and Transmittance of Materials Using Spectrophotometers with Integrating Spheres," *ASTM E903*, ASTM, Philadelphia (1988).
48. ASTM, "Standard Guide for Selection of Geometric Conditions for Measurement of Reflectance and Transmittance Properties of Materials," *ASTM E179*, ASTM, Philadelphia (1990).
49. R. C. Hawes, "Technique for Measuring Photometric Accuracy," *Appl. Opt.* **10**:1246 (1971).
50. J. M. Bennett and E. J. Ashley, "Calibration of Instruments Measuring Reflectance and Transmittance," *Appl. Opt.* **11**:1749 (1972).
51. ASTM, "Estimating Stray Radiant Energy of Spectrophotometers," *ASTM E387*, ASTM, Philadelphia (1984).
52. R. Mavrodineanu and J. R. Baldwin, *Standard Reference Materials: Glass Filters as a Standard Reference Material for Spectrophotometry—Selection, Preparation, Certification, Use*, NBS Special Publication SP260-51, U.S. National Bureau of Standards, Washington, D.C., 1975.
53. R. Mavrodineanu and J. R. Baldwin, *Standard Reference Materials: Metal-on-Quartz Filters as a Standard Reference Material for Spectrophotometry*, NBS Special Publication SP260-68, U.S. National Bureau of Standards, Washington, D.C., 1980.
54. J. J. Hsia, "National Scales of Spectrometry in the U.S.," in C. Burgess and K. D. Mielenz (eds.), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1987.
55. K. E. Eckerle, E. Sutter, G. H. Freeman, G. Andor, and L. Fillinger, "International Intercomparison for Transmittance," *Metrologia* **27**:33 (1990).
56. L. Fillinger and G. Andor, "International Intercomparison of Transmittance Measurement," *CIE Journal* **7**:21 (1990).
57. H. G. Lipson, L. H. Skolnik, and D. L. Stierwalt, "Small Absorption Coefficient Measurement by Calorimetric and Spectral Emittance Techniques," *Appl. Opt.* **13**:1741 (1974).
58. A. Hordvik, "Measurement Techniques for Small Absorption Coefficients: Recent Advances," *Appl. Opt.* **16**:2827 (1977).
59. J. Strong, *Procedures in Experimental Physics*, Prentice-Hall, N.Y., 1938, p. 376.
60. J. Strong, *Procedures in Applied Optics*, Dekker, N.Y., 1989, p. 162.
61. C. Castellini, G. Emiliani, E. Masetti, P. Poggi, and P. P. Polato, "Characterization and Calibration of a Variable-Angle Absolute Reflectometer," *Appl. Opt.* **29**:538 (1990).
62. R. S. Ram, O. Prakash, J. Singh, and S. P. Varma, "Simple Design for a Reflectometer," *Opt. Eng.* **30**:467 (1991).
63. R. F. Weeks, "Simple Wide Range Specular Reflectometer," *J. Opt. Soc. Am.* **48**:775 (1958).
64. V. R. Weidner and J. J. Hsia, "NBS Specular Reflectometer-Spectrophotometer," *Appl. Opt.* **19**:1268 (1980).
65. D. K. Zhuang and T. L. Yang, "Spectral Reflectance Measurements Using a Precision Multiple Reflectometer in the UV and VUV Range," *Appl. Opt.* **28**:5024 (1989).
66. H. E. Bennett and W. F. Koehler, "Precision Measurements of Absolute Specular Reflectance with Minimized Systematic Errors," *J. Opt. Soc. Am.* **50**:1 (1960).
67. A. Bittar and J. D. Hamlin, "High-Accuracy True Normal-Incidence Absolute Reflectometer," *Appl. Opt.* **23**:4054 (1984).
68. G. Boivin and J. M. Theriault, "Reflectometer for Precise Measurement of Absolute Specular Reflectance at Normal Incidence," *Rev. Sci. Instrum.* **52**:1001 (1981).
69. J. E. Shaw and W. R. Blevin, "Instrument for the Absolute Measurement of Direct Spectral Reflectances at Normal Incidence," *J. Opt. Soc. Am.* **54**:334 (1964).
70. ASTM, "Test Method for Specular Gloss," *ASTM D523*, ASTM, Philadelphia (1988).

71. ASTM, "Test Method for Measuring Optical Reflectivity of Transparent Materials," *ASTM F1252*, ASTM, Philadelphia (1990).
72. R. S. Hernicz and D. P. DeWitt, "Evaluation of a High Accuracy Reflectometer for Specular Materials," *Appl. Opt.* **12**:2454 (1973).
73. K. A. Snail, A. A. Morrish, and L. M. Hanssen, "Absolute Specular Reflectance Measurements in the Infrared," *Proc. SPIE* **692**:143 (1986).
74. ASTM, "Test for Reflectivity and Coefficient of Scatter of White Porcelain Enamels," *ASTM C347*, ASTM, Philadelphia (1983).
75. ASTM, "Recommended Practice for Goniophotometry of Objects and Materials," *ASTM E167*, ASTM, Philadelphia (1987).
76. ASTM, "Standard Test Method for Reflectance Factor and Color by Spectrophotometry Using Bidirectional Geometry," *ASTM E1349*, ASTM, Philadelphia (1990).
77. ASTM, "Measurement of Gloss of High-Gloss Surfaces by Goniophotometry," *ASTM E430*, ASTM, Philadelphia (1983).
78. W. Erb, "Computer-Controlled Gonioreflectometer for the Measurement of Spectral Reflection Characteristics," *Appl. Opt.* **19**:3789 (1980).
79. W. Erb, "High Accuracy Gonioreflectance Spectrometry," Chap. 2.2 in C. Burgess and K. D. Mielenz (eds), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1987.
80. R. S. Hunter, *The Measurement of Appearance*, Wiley, N.Y., 1975.
81. J. A. J. Jacquez and H. F. Kuppenheim, "Theory of the Integrating Sphere," *J. Opt. Soc. Am.* **45**:460 (1955).
82. D. G. Goebel, "Generalized Integrating Sphere Theory," *Appl. Opt.* **6**:125 (1967).
83. D. K. Edwards, J. T. Gier, K. E. Nelson, and R. D. Roddick, "Integrating Sphere for Imperfectly Diffuse Samples," *J. Opt. Soc. Am.* **51**:1279 (1961).
84. F. J. J. Clarke and J. A. Compton, "Correction Methods for Integrating Sphere Measurement of Hemispherical Reflectance," *Color. Res. Appl.* **11**:253 (1986).
85. F. J. J. Clarke and J. A. Larkin, "Improved Techniques for the NPL Hemispherical Reflectometer," *Proc. SPIE* **917**:7 (1988).
86. R. V. Dunkle, "Spectral Reflectance Measurements," in F. J. Clauss (ed.), *Surface Effects on Spacecraft Materials*, Wiley, N.Y., 1960, p. 117.
87. W. G. Egan and T. Hilgeman, "Integrating Spheres for Measurements between 0.185 micrometers and 12 micrometers," *Appl. Opt.* **14**:1137 (1975).
88. D. G. Goebel, B. P. Caldwell, and H. K. Hammond III, "Use of an Auxiliary Sphere with a Spectroreflectometer to Obtain Absolute Reflectance," *J. Opt. Soc. Am.* **56**:783 (1966).
89. B. J. Hisdal, "Reflectance of Perfect Diffuse and Specular Samples in the Integrating Sphere," *J. Opt. Soc. Am.* **55**:1122 (1965a).
90. B. J. Hisdal, "Reflectance of Nonperfect Surfaces in the Integrating Sphere," *J. Opt. Soc. Am.* **55**:1255 (1965b).
91. H. J. McNicholas, "Absolute Methods in Reflectometry," *J. Res. Natl. Bur. Stand.* **1**:29 (1928).
92. W. Richter and W. Erb, "Accurate Diffuse Reflectance Measurements in the IR Spectral Range," *Appl. Opt.* **26**:4620 (1987).
93. D. Sheffer, U. P. Oppenheim, D. Clement, and A. D. Devir, "Absolute Reflectometer for the 0.8–2.5 micrometer Region," *Appl. Opt.* **26**:583 (1987).
94. D. Sheffer, U. P. Oppenheim, and A. D. Devir, "Absolute Measurements of Diffuse Reflectances in the x°/d Configuration," *Appl. Opt.* **30**:3181 (1991).
95. A. H. Taylor, "Errors in Reflectometry," *J. Opt. Soc. Am.* **25**:51 (1935).
96. ASTM, "Test for Light Reflectance of Acoustical Materials by the Integrating Sphere Method," *ASTM C523*, ASTM, Philadelphia (1984).
97. ASTM, "Measurement and Calculation of Reflecting Characteristics of Metallic Surfaces Using Integrating Sphere Instruments," *ASTM E429*, ASTM, Philadelphia (1987).
98. CIE, "Absolute Methods for Reflection Measurements," *CIE Publication 44*, CIE, Paris (1979).

99. IES Testing Procedures Committee, "IES Approved Method for Total and Diffuse Reflectometry," IES LM-44-1985, *J. Illum. Eng. Soc.* **19**:195 (1985).
100. J. A. Van den Akker, L. R. Dearth, and W. M. Shilcox, "Evaluation of Absolute Reflectance for Standardization Purposes," *J. Opt. Soc. Am.* **56**:250 (1966).
101. W. H. Venable, J. J. Hsia, and V. R. Weidner, "Establishing a Scale of Directional-Hemispherical Reflectance Factor I: The Van den Akker Method," *J. Res. Natl. Bur. Stand.* **82**:29 (1977).
102. J. D. Lindberg, "Absolute Diffuse Reflectance from Relative Reflectance Measurements," *Appl. Opt.* **26**:2900 (1987).
103. P. Y. Barnes and J. J. Hsia, "45°/0° Bidirectional Reflectance Distribution Function Standard Development," *Proc. SPIE* **1165**:165 (1989).
104. W. R. Blevin and W. J. Brown, "An Infrared Reflectometer with a Spheroidal Mirror," *J. Sci. Instrum.* **42**:1 (1965).
105. S. T. Dunn, J. C. Richmond, and J. C. Weibel, "Ellipsoidal Mirror Reflectometer," *J. Res. Nat. Bur. Stand.* **70C**:75 (1966b).
106. L. M. Hanssen and K. A. Snail, "Infrared Diffuse Reflectometer for Spectral, Angular and Temperature Resolved Measurements," *Proc. SPIE* **807**:148 (1987).
107. P. L. Hartman and E. Logothetis, "An Absolute Reflectometer for Use at Low Temperatures," *Appl. Opt.* **3**:255 (1964).
108. J. T. Neu, R. S. Dummer, and O. E. Myers, "Hemispherical Directional Ellipsoidal Infrared Spectroreflectometer," *Proc. SPIE* **807**:165 (1987).
109. B. E. Wood, J. G. Pipes, A. M. Smith, and J. A. Roux, "Hemi-Ellipsoidal Mirror Infrared Reflectometer," *Appl. Opt.* **15**:940 (1976).
110. F. J. J. Clarke and D. J. Parry, "Helmholtz Reciprocity: Its Validity and Application to Reflectometry," *Light. Res. Technol.* **17**:1 (1985).
111. J. T. Agnew and R. B. McQuistan, "Experiments Concerning Infrared Diffuse-Reflectance Standards in the Range 0.8 to 20.0 Micrometers," *J. Opt. Soc. Am.* **43**:999 (1953).
112. J. T. Gier, R. V. Dunkle, and J. T. Bevans, "Measurement of Absolute Spectral Reflectivity from 1.0 to 15 microns," *J. Opt. Soc. Am.* **44**:558 (1954).
113. D. C. Reid and E. D. McAlister, "Measurement of Spectral Emissivity from 3 μ to 15 μ ," *J. Opt. Soc. Am.* **49**:78 (1959).
114. ASTM, "Test Method for (45-0) Directional Reflectance Factor of Opaque Specimens by Broad-Band Filter Reflectometry," *ASTM E97*, ASTM, Philadelphia (1987).
115. ASTM, "Standard Test Method for Reflectance Factor and Color by Spectrophotometry Using Hemispherical Geometry," *ASTM E1348*, ASTM, Philadelphia (1990).
116. J. J. Hsia and V. R. Weidner, "NBS 45-degree/Normal Reflectometer for Absolute Reflectance Factors," *Metrologia* **17**:97 (1981).
117. A. H. Taylor, "The Measurement of Diffuse Reflection Factors and a New Absolute Reflectometer," *J. Opt. Soc. Am.* **4**:9 (1920).
118. M. J. Duggin, "The Field Measurement of Reflectance Factors," *Photogram. Eng. Rem. Sens.* **46**:643 (1980).
119. M. J. Duggin and W. R. Philipson, "Field Measurement of Reflectance: Some Major Considerations," *Appl. Opt.* **21**:2833 (1982).
120. C. Asmail, "Bidirectional Scattering Distribution Function (BSDF): A Systematized Bibliography," *J. Res. Natl. Inst. Stand. Technol.* **96**:215 (1991).
121. ASTM, "Standard Practice for Angle Resolved Optical Scatter Measurements on Specular or Diffuse Surfaces," *ASTM E1392*, ASTM, Philadelphia (1990).
122. F. O. Bartell, E. L. Dereniak, and W. L. Wolfe, "Theory and Measurement of Bidirectional Reflectance Distribution Function (BRDF) and Bidirectional Transmittance Distribution Function (BTDF)," *Proc. SPIE* **257**:154 (1980).
123. CIE, "Retroreflection: Definition and Measurement," *CIE Publ. 54*, CIE, Paris (1982).
124. ASTM, "Test Method for Coefficient of Retroreflection on Retroreflective Sheeting," *ASTM E810*, ASTM, Philadelphia (1981).
125. ASTM, "Standard Practice for Measuring Photometric Characteristics of Retroreflectors," *ASTM E809*, ASTM, Philadelphia (1991).

126. W. H. Venable and N. L. Johnson, "Unified Coordinate System for Retroreflectance Measurements," *Appl. Opt.* **19**:1236 (1980).
127. K. L. Eckerle, J. J. Hsia, V. R. Weidner, and W. H. Venable, "NBS Reference Retroreflectometer," *Appl. Opt.* **19**:1253 (1980).
128. G. Hass and J. E. Waylonis, "Optical Constants and Reflectance and Transmittance of Evaporated Aluminum in the Visible and Ultraviolet," *J. Opt. Soc. Am.* **51**:719 (1961).
129. H. E. Bennett, J. M. Bennett, and E. J. Ashley, "Infrared Reflectance of Evaporated Aluminum Films," *J. Opt. Soc. Am.* **52**:1245 (1962).
130. H. E. Bennett, M. Silver, and E. J. Ashley, "Infrared Reflectance of Aluminum Evaporated in Ultra-High Vacuum," *J. Opt. Soc. Am.* **53**:1089 (1963).
131. G. Hass and R. E. Thun, *Physics of Thin Films*, vol. 2, Academic, N.Y., 1964, p. 337.
132. J. M. Bennett and E. J. Ashley, "Infrared Reflectance and Emittance of Silver and Gold Evaporated in Ultra-High Vacuum," *Appl. Opt.* **4**:221 (1965).
133. G. Hass, "Reflectance and Preparation of Front-Surface Mirrors for Use at Various Angles of Incidence from the Ultraviolet to the Far Infrared," *J. Opt. Soc. Am.* **72**:27 (1982).
134. J. C. Richmond and J. J. Hsia, *Preparation and Calibration of Standards of Spectral Specular Reflectance*, NBS Special Publication SP260-38, U.S. National Bureau of Standards, Washington, D.C., 1972.
135. J. C. Richmond, "Rationale for Emittance and Reflectivity," *Appl. Opt.* **21**:1 (1982).
136. J. F. Verrill, "Physical Standards in Absorption and Reflection Spectrometry," Chap. 3.1 in C. Burgess and K. D. Mielenz (eds), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1987.
137. V. R. Weidner and J. J. Hsia, *Standard Reference Materials: Preparation and Calibration of First Surface Aluminum Mirror Spectral Reflectance Standard*, NBS Special Publication SP260-75, U.S. National Bureau of Standards, Washington, D.C., 1982.
138. W. Erb, "Requirements for Reflection Standards and the Measurement of their Reflection Value," *Appl. Opt.* **14**:493 (1975).
139. W. Erb and W. Budde, "Properties of Standard Materials for Reflection," *Color Res. Appl.* **4**:113 (1979).
140. D. Scheffer, U. P. Oppenheim, and A. D. Devir, "Absolute Reflectometer for the Mid-Infrared Region," *Appl. Opt.* **29**:129 (1990).
141. W. E. K. Middleton and C. L. Sanders, "The Absolute Spectral Diffuse Reflectance of Magnesium Oxide," *J. Opt. Soc. Am.* **41**:419 (1951).
142. F. Grum and G. W. Luckey, "Optical Sphere Paint and a Working Standard of Reflectance," *Appl. Opt.* **7**:2289 (1968).
143. F. Grum and M. Saltzman, "New White Standard of Reflectance," *CIE Publication 36*, CIE, Paris (1976).
144. V. R. Weidner and J. J. Hsia, "Reflection Properties of Pressed Polytetrafluorethylene Powder," *J. Opt. Soc. Am.* **71**:856 (1981).
145. ASTM, "Preparation of Pressed Power White Reflectance Factor Transfer Standards for Hemispherical Geometry," *ASTM E259*, ASTM, Philadelphia (1992).
146. V. R. Weidner and J. J. Hsia, *Spectral Reflectance*, NBS Special Publication SP250-8, U.S. National Bureau of Standards, Washington, D.C., 1987.
147. M. D. Fairchild and D. J. O. Daoust, "Goniospectrophotometric Analysis of Pressed PTFE Powder for Use as a Primary Transfer Standard," *Appl. Opt.* **27**:3392 (1988).
148. R. D. Saunders and W. R. Ott, "Spectral Irradiance Measurements: Effect of UV-Produced Luminescence in Integrating Spheres," *Appl. Opt.* **15**:827 (1976).
149. M. Kronstein, R. J. Kraushaar, and R. E. Deacle, "Sulfur as a Standard of Reflectance in the Infrared," *J. Opt. Soc. Am.* **53**:458 (1963).
150. S. T. Dunn, "Application of Sulfur Coatings to Integrating Spheres," *Appl. Opt.* **4**:877 (1965).
151. R. Tkachuk and F. D. Kuzina, "Sulfur as a Proposed Near Infrared Reflectance Standard," *Appl. Opt.* **17**:2817 (1978).
152. T. W. Stuhlinger, E. L. Dereniak, and F. O. Bartell, "Bidirectional Distribution Function of Gold-Plated Sandpaper," *Appl. Opt.* **20**:2648 (1981).
153. CIE, "A Review of Publications on Properties and Reflection Values of Material Reflection Standards," *CIE Publication 46*, CIE, Paris (1979b).

154. F. J. J. Clarke, F. A. Garforth, and D. J. Parr, "Goniophotometric and Polarization Properties of White Reflection Standard Materials," *Light. Res. Technol.* **15**:133 (1983).
155. ASTM, "Practice for Preparation of Reference White Reflectance Standards," *ASTM E259*, ASTM, Philadelphia (1987).
156. W. Budde, "Standards of Reflectance," *J. Opt. Soc. Am.* **50**:217 (1960).
157. W. Budde, "Calibration of Reflectance Standards," *J. Res. Natl. Bur. Stand.* **A80**:585 (1976).
158. W. G. Egan and T. Hilgeman, "Retroreflectance Measurements of Photometric Standards and Coatings," *Appl. Opt.* **15**:1845 (1976).
159. L. Morren, G. Vandermeersch, and P. Antoine, "A Study of the Reflection Factor of Usual Photometric Standards in the Near Infrared," *Light. Res. Technol.* **4**:243 (1972).
160. TAPPI, "Calibration of Reflectance Standards for Hemispherical Geometry," TAPPI Standard TIS 0804-07 in *1990 TAPPI Test Methods*, TAPPI, Atlanta, 1990.
161. V. R. Weidner, *Standard Reference Materials: White, White Opal Glass Diffuse Spectral Reflectance Standards for the Visible Spectrum*, NBS Special Publication SP260-82, U.S. National Bureau of Standards, Washington, D.C., 1983.
162. V. R. Weidner, "Gray Scale of Diffuse Reflectance for the 250–2500 nm Wavelength Range," *Appl. Opt.* **25**:1265 (1986).
163. W. Budde, W. Erb, and J. J. Hsia, "International Intercomparison of Absolute Reflectance Scales," *Color Res. Appl.* **7**:24 (1982).
164. E. M. Sparrow, P. D. Kruger, and R. P. Heinisch, "Cavity Methods for Determining the Emittance of Solids," *Appl. Opt.* **12**:2466 (1973).
165. J. C. DeVos, "A New Determination of the Emissivity of Tungsten Ribbon," *Physica* **20**:690 (1954).
166. ASTM, "Test for Total Normal Emittances of Surfaces Using Inspection Meter Techniques," *ASTM E408*, ASTM, Philadelphia (1990).
167. S. T. Dunn, J. C. Richmond, and J. F. Parmer, "Survey of Infrared Measurement Techniques and Computational Methods in Radiant Heat Transfer," *J. Spacecraft Rockets* **3**:961 (1966a).
168. J. P. Millard and E. R. Streed, "A Comparison of Infrared Emittance Measurements and Measurement Techniques," *Appl. Opt.* **8**:1485 (1969).
169. ASTM, "Total Hemispherical Emittance of Surfaces from 20 to 1400C," *ASTM C835*, ASTM, Philadelphia (1988).
170. ASTM, "Test for Calorimetric Determination of Hemispherical Emittance and the Ratio of Solar Absorbance to Hemispherical Emittance Using Solar Simulation," *ASTM E434*, ASTM, Philadelphia (1990).
171. D. K. Edwards, "Thermal Radiation Measurements," Chap. 9 in E. R. G. Eckert and R. J. Goldstein (eds), *Measurement Techniques in Heat Transfer*, AGARD 130, Technivision, Slough, England, 1970, p. 353.
172. J. C. Richmond and W. N. Harrison, "Equipment and Procedures for Evaluation of Total Hemispherical Emittance," *Am. Ceram. Soc. Bull.* **39**:668 (1960).
173. R. R. Willey, "Results of a Round-Robin Measurement of Spectral Emittance in the Mid-Infrared," *Proc. SPIE* **807**:140 (1987).
174. D. Weber, "Spectral Emissivity of Solids in the Infrared at Low Temperatures," *J. Opt. Soc. Am.* **49**:815 (1959).
175. R. Gardon, "The Emissivity of Transparent Materials," *J. Am. Ceram. Soc.* **39**:278 (1956).
176. ASTM, "Test for Normal Spectral Emittance at Elevated Temperatures," *ASTM E307*, ASTM, Philadelphia (1990).
177. ASTM, "Test for Normal Spectral Emittance at Elevated Temperatures of Non-Conducting Specimens," *ASTM E423*, ASTM, Philadelphia (1990).
178. D. P. DeWitt, "Inferring Temperature from Optical Radiation Measurements," *Opt. Eng.* **25**:596 (1986).
179. D. P. DeWitt and J. C. Richmond, "Theory and Measurement of the Thermal Radiative Properties of Metals," in *Techniques of Metals Research*, vol. 6, Wiley, N.Y., 1970.
180. G. A. Hornbeck, "Optical Methods of Temperature Measurement," *Appl. Opt.* **5**:179 (1966).
181. J. S. Redgrove, "Measurement of the Spectral Emissivity of Solid Materials," *Measurement (UK)* **8**:90 (1990).
182. D. L. Stierwalt, "Infrared Spectral Emittance Measurements on Optical Materials," *Appl. Opt.* **5**:1911 (1966).
183. A. M., Wittenberg, "Determination of Total Emittance of a Nongray Surface," *J. Appl. Phys.* **39**:1936 (1968).

35.14 FURTHER READING

- ASTM, *ASTM Standards on Color and Appearance Measurement*, 3d ed., ASTM, Philadelphia, 1991. Blau, Jr, H. H. and H. Fischer (eds.), *Radiative Transfer from Solid Materials*, Macmillan, N.Y., 1962. Burgess, C. and K. D. Mielenz (eds.), *Advances in Standards and Methodology in Spectrophotometry*, Elsevier, Amsterdam, 1987.
- Clauss, F. J. (ed.), *First Symposium, Surface Effects on Spacecraft Materials*, Wiley, N.Y., 1960.
- Frei, R. W. and J. D. MacNeil, *Diffuse Reflectance Spectroscopy in Environmental Problem-Solving*, CRC Press, Cleveland, 1973.
- Grum, F. and R. J. Becherer, "Radiometry," in *Optical Radiation Measurements*, vol. 1, Academic, N.Y., 1979.
- Hammond III, H. K., and H. L. Mason (eds.), "Selected NBS Papers on Radiometry and Photometry," NBS Special Publication SP300-7, *Precision Measurement and Calibration*, U.S. National Bureau of Standards, Washington, D.C., 1971.
- Hunter, R. S., *The Measurement of Appearance*, Wiley, N.Y., 1975.
- Kortum, G., *Reflectance Spectroscopy*, Springer-Verlag, N.Y., 1969.
- Nimeroff, L. (ed.), "Selected NBS Papers on Colorimetry," NBS Special Publication SP300-9, *Precision Measurement and Calibration*, U.S. National Bureau of Standards, Washington, D.C., 1972.
- Richmond, J. C. (ed.), "Measurement of Thermal Radiation Properties of Solids," *NASA Special Publication SP-31*, National Aeronautics and Space Administration, Washington, D.C., 1963.
- Walsh, J. W. T., *Photometry*, 3d ed., Dover, N.Y., 1958.
- Wendlandt, W. W. and H. G. Hecht, *Reflectance Spectroscopy*, Interscience, N.Y., 1966.

This page intentionally left blank.

DO NOT DUPLICATE

RADIOMETRY AND PHOTOMETRY: UNITS AND CONVERSIONS

James M. Palmer*

*College of Optical Sciences
University of Arizona
Tucson, Arizona*

36.1 GLOSSARY†

A	area, m^2
A_{proj}	projected area, m^2
E, E_e	radiant incidence (irradiance), W m^{-2}
E_p, E_q	photon incidence, $\text{s}^{-1} \text{m}^{-2}$
E_v	illuminance, lm m^{-2} [lux (lx)]
I, I_e	radiant intensity, W sr^{-1}
I_p, I_q	photon intensity, $\text{s}^{-1} \text{sr}^{-1}$
I_v	luminous intensity, candela (cd)
K_m	absolute luminous efficiency at λ_p for photopic vision, 683 lm/W
K'_m	absolute luminous efficiency at λ_p for scotopic vision, 1700 lm/W
$K(\lambda)$	absolute spectral luminous efficiency, photopic vision, lm/W
$K'(\lambda)$	absolute spectral luminous efficiency, scotopic vision, lm/W
L, L_e	radiance, $\text{W m}^{-2} \text{sr}^{-1}$
L_p, L_q	photon radiance (photonance), $\text{s}^{-1} \text{m}^{-2} \text{sr}^{-1}$
L_v	luminance, cd sr^{-1}
M, M_e	radiant exitance, W m^{-2}
M_p, M_q	photon exitance, $\text{s}^{-1} \text{m}^{-2}$
Q, Q_e	radiant energy, joule (J)
Q_p, Q_q	photon energy, J or eV
Q_v	luminous energy, lm s^{-1}

*Deceased.

†Note. The subscripts are used as follows: e (energy) for radiometric, v (visual) for photometric, and q (or p) for photonic. The subscript e is usually omitted; the other subscripts may also be omitted if the context is unambiguous.

$\mathfrak{R}(\lambda)$	spectral responsivity, A/W or V/W
$V(\lambda)$	relative spectral luminous efficiency, photopic vision
$V'(\lambda)$	relative spectral luminous efficiency, scotopic vision
$V_q(\lambda)$	relative spectral luminous efficiency for photons
λ	wavelength, nm or μm
λ_p	wavelength at peak of function, nm or μm
Φ, Φ_e	radiant power (radiant flux), watt (W)
Φ_p, Φ_q	photon flux, s^{-1}
Φ_v	luminous power (luminous flux), lumen (lm)
ω	solid angle, steradian (sr)
Ω	projected solid angle, sr

36.2 INTRODUCTION AND BACKGROUND

After more than a century of turmoil, the symbols, units, and nomenclature (SUN) for radiometry and photometry seem somewhat stable at present. There are still small, isolated pockets of debate and even resistance; by and large, though, the community has settled on the International System of Units (SI) and the recommendations of the International Organization for Standardization (ISO)¹ and the International Union of Pure and Applied Physics (IUPAP).²

The seed of the SI was planted in 1799 with the deposition of two prototype standards, the meter and the kilogram, into the Archives de la République in Paris. Noted physicists Gauss, Weber, Maxwell, and Thompson made significant contributions to measurement science over the next 75 years. Their efforts culminated in the Convention of the Meter, a diplomatic treaty originally signed by representatives of 17 nations in 1875 (currently there are 48 member nations). The Convention grants authority to the General Conference on Weights and Measures (CGPM), the International Committee for Weights and Measures (CIPM), and the International Bureau of Weights and Measures (BIPM). The CIPM, along with a number of subcommittees, suggests modifications to the CGPM. In our arena, the subcommittee is the Consultative Committee on Photometry and Radiometry (CCPR). The BIPM, the international metrology institute, is the physical facility that is responsible for realization, maintenance, and dissemination of standards.

The SI was adopted by the CGPM in 1960 and is the official system in the 48 member states. It currently consists of seven base units and a much larger number of derived units. The base units are a choice of seven well-defined units that, by convention, are regarded as independent. The seven base units are as follows:

1. Meter
2. Kilogram
3. Second
4. Ampere
5. Kelvin
6. Mole
7. Candela

The derived units are those that are formed by various combinations of the base units.

International organizations involved in the promulgation of SUN include the International Commission on Illumination (CIE), the IUPAP, and the ISO. In the United States, the American National Standards Institute (ANSI) is the primary documentary (protocol) standards organization. Many other scientific and technical organizations publish recommendations concerning the use of SUN for their scholarly journals. Several examples are the Illuminating Engineering Society

TABLE 1 Projected Areas of Common Shapes

Shape	Area	Projected area
Flat rectangle	$A = L \times W$	$A_{\text{proj}} = L \times W \cos \beta$
Circular disc	$A = \pi r^2 = \pi d^2/4$	$A_{\text{proj}} = \pi r^2 \cos \beta = (\pi d^2 \cos \beta)/4$
Sphere	$A = 4\pi r^2 = \pi d^2$	$A_{\text{proj}} = A/4 = \pi r^2$

(IESNA), the International Astronomical Union (IAU), the Institute for Electrical and Electronic Engineering (IEEE), and the American Institute of Physics (AIP).

The terminology employed in radiometry and photometry consists principally of two parts: (1) an adjective distinguishing between a radiometric, photonic, or photometric entity, and (2) a noun describing the underlying geometric or spatial concept. In some instances, the two parts are replaced by a single term (e.g., radiance).

There are some background concepts and terminology that are needed before proceeding further.

Projected area is defined as the rectilinear projection of a surface of any shape onto a plane normal to the unit vector. The differential form is $dA_{\text{proj}} = \cos(\beta)dA$, where β is the angle between the local surface normal and the line of sight. Integrate over the (observable) surface area to get

$$A_{\text{proj}} = \int_A \cos \beta dA \quad (1)$$

Some common examples are shown in Table 1.

Plane angle and solid angle are both derived units in the SI system. The following definitions are from National Institute of Standards and Technology (NIST) SP811.³

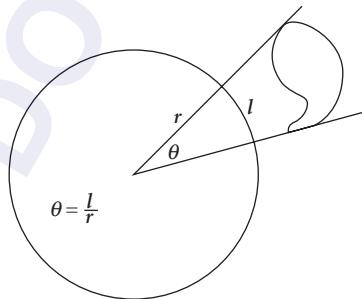
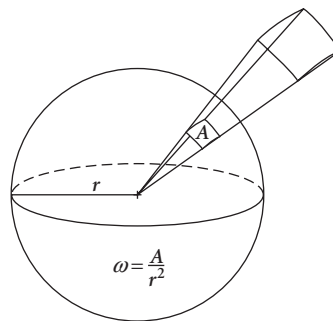
The radian is the plane angle between two radii of a circle that cuts off on the circumference an arc equal in length to the radius.

The abbreviation for the radian is *rad*. Since there are 2π rad in a circle, the conversion between degrees and radians is $1 \text{ rad} = (180/\pi)$ degrees. (See Fig. 1.)

A solid angle is the same concept that is extended to three dimensions.

One steradian (sr) is the solid angle that, having its vertex in the center of a sphere, cuts off an area on the surface of the sphere equal to that of a square with sides of length equal to the radius of the sphere.

The solid angle is the ratio of the spherical area (A_{proj}) to the square of the radius, r . The spherical area is a projection of the object of interest onto a unit sphere, and the solid angle is the surface area of that projection, as shown in Fig. 2. Divide the surface area of a sphere by the square of its radius to find that there are 4π sr of solid angle in a sphere. One hemisphere has 2π sr. The accepted symbols


FIGURE 1 Plane angle.

FIGURE 2 Solid angle.

for solid angle are either lowercase Greek omega (ω) or uppercase Greek omega (Ω). I recommend using ω exclusively for solid angle, and reserving Ω for the advanced concept of projected solid angle ($\omega \cos \theta$). The equation for solid angle is $d\omega = dA_{\text{proj}}/r^2$. For a right circular cone, $\omega = 2\pi(1 - \cos \theta)$, where θ is the half-angle of the cone.

Both plane angles and solid angles are dimensionless quantities, and they can lead to confusion when attempting dimensional analysis. For example, the simple inverse square law, $E = I/d^2$ appears dimensionally inconsistent. The left side has units W m^{-2} while the right side has $\text{W sr}^{-1} \text{m}^{-2}$. It has been suggested that this equation be written $E = I\Omega_0/d^2$, where Ω_0 is the unit solid angle, 1 sr. Inclusion of the term Ω_0 will render the equation dimensionally correct, but will far too often be considered a free variable rather than a constant equal to 1, leading to erroneous results.

Isotropic and lambertian both carry the meaning “the same in all directions” and, regrettably, are used interchangeably.

Isotropic implies a spherical source that radiates the same amount in all directions [i.e., the intensity (W/sr or cd) is independent of direction]. The term *isotropic point source* is often heard. No such entity can exist, for the energy density would necessarily be infinite. A small, uniform sphere comes very close. Small in this context means that the size of the sphere is much less than the distance from the sphere to the plane of observation, such that the inverse square law is applicable. An example is a globular tungsten lamp with a milky white diffuse envelope some 10 cm in diameter, as viewed from a distance greater than 1 m. From our vantage point, a distant star is considered an isotropic point source.

Lambertian implies a flat radiating surface. It can be an active emitter or a passive, reflective surface. The intensity decreases with the cosine of the observation angle with respect to the surface normal (Lambert’s law). The radiance ($\text{W m}^{-2} \text{sr}^{-1}$) or luminance (cd m^{-2}) is independent of direction. A good approximation is a surface painted with a quality matte, or flat, white paint. The intensity is the product of the radiance, L , or luminance, L_v , and the projected area A_{proj} . If the surface is uniformly illuminated, it appears equally bright from whatever direction it is viewed. Note that the flat radiating surface can be used as an elemental area of a curved surface.

The ratio of the radiant exitance (power per unit area, W m^{-2}) to the radiance (power per unit projected area per unit solid angle, $\text{W m}^{-2} \text{sr}^{-1}$) of a lambertian surface is a factor of π and not 2π . This result is not intuitive, as, by definition, there are 2π sr in a hemisphere. The factor of π comes from the influence of the $\cos\theta$ term while integrating over a hemisphere.

A sphere with a lambertian surface illuminated by a distant point source will display a radiance that is maximum at the point where the local normal coincides with the incoming beam. The radiance will fall off with a cosine dependence to zero at the terminator. If the intensity (integrated radiance over area) is unity when viewing from the direction of the source, then the intensity when viewing from the side is $1/\pi$. Think about this and ponder whether our Moon has a lambertian surface.

36.3 SYMBOLS, UNITS, AND NOMENCLATURE IN RADIOMETRY

Radiometry is the measurement of optical radiation, defined as electromagnetic radiation within the frequency range from 3×10^{11} to 3×10^{16} hertz (Hz). This range corresponds to wavelengths between 0.01 and 1000 micrometers (μm) and includes the regions commonly called *ultraviolet*, *visible*, and *infrared*.

Radiometric units can be divided into two conceptual areas: (1) those having to do with power or energy, and (2) those that are geometric in nature. In the first category are:

Energy is an SI-derived unit, measured in joules (J). The recommended symbol for energy is Q . An acceptable alternate is W .

Power (also called *radiant flux*) is another SI-derived unit. It is the derivative of energy with respect to time, dQ/dt , and the unit is the watt (W). The recommended symbol for power is uppercase Greek phi (Φ). An accepted alternate is P .

Energy is the integral of power over time, and it is commonly used with integrating detectors and pulsed sources. Power is used for continuous sources and nonintegrating detectors. The radiometric quantity power can now be combined with the geometric spatial quantities area and solid angle.

Irradiance (also referred to as *flux density* or *radiant incidence*) is an SI-derived unit and is measured in W m^{-2} . Irradiance is power per unit area *incident* from all directions within a hemisphere onto a surface that coincides with the base of that hemisphere. A related quantity is *radiant exitance*, which is power per unit area *leaving* a surface into a hemisphere whose base is that surface. The symbol for irradiance is E , and the symbol for radiant exitance is M . Irradiance (or radiant exitance) is the derivative of power with respect to area, $d\Phi/dA$. The integral of irradiance or radiant exitance over area is power. There is no compelling reason to have two quantities carrying the same units, but it is convenient.

Radiant intensity is an SI-derived unit and is measured in W sr^{-1} . Intensity is power per unit of solid angle. The symbol is I . *Intensity* is the derivative of power with respect to solid angle, $d\Phi/d\Omega$. The integral of radiant intensity over solid angle is power.

A great deal of confusion surrounds the use and misuse of the term *intensity*. Some use it for W sr^{-1} ; some use it for W m^{-2} ; others use it for $\text{W m}^{-2} \text{sr}^{-1}$. It is quite clearly defined in the SI system, in the definition of the base unit of luminous intensity—the candela. Attempts are often made to justify these different uses of intensity by adding adjectives like *optical* or *field* (used for W m^{-2}) or *specific* (used for $\text{W m}^{-2} \text{sr}^{-1}$). In the SI system, the underlying geometric concept for intensity is quantity per unit solid angle. For more discussion, see Palmer.⁴

Radiance is an SI-derived unit and is measured in $\text{W m}^{-2} \text{sr}^{-1}$. Radiance is a directional quantity, power per unit projected area per unit solid angle. The symbol is L . Radiance is the derivative of power with respect to solid angle and projected area, $d\Phi/d\omega dA \cos\theta$, where θ is the angle between the surface normal and the specified direction. The integral of radiance over area and solid angle is power.

Photon quantities are also common. They are related to the radiometric quantities by the relationship $Q_p = hc/\lambda$, where Q_p is the energy of a photon at wavelength λ , h is Planck's constant, and c is the velocity of light. At a wavelength of $1 \mu\text{m}$, there are approximately 5×10^{18} photons per second in a watt. Conversely, one photon has an energy of about $2 \times 10^{-19} \text{ J (W/s)}$ at $1 \mu\text{m}$.

36.4 SYMBOLS, UNITS, AND NOMENCLATURE IN PHOTOMETRY

Photometry is the measurement of light, electromagnetic radiation detectable by the human eye. It is thus restricted to the wavelength range from about 360 to 830 nanometers (nm; $1000 \text{ nm} = 1 \mu\text{m}$). Photometry is identical to radiometry *except* that everything is weighted by the spectral response of the nominal human eye. *Visual photometry* uses the eye as a comparison detector, while *physical photometry* uses either optical radiation detectors constructed to mimic the spectral response of the nominal eye, or spectroradiometry coupled with appropriate calculations to do the eye response weighting.

Photometric units are basically the same as the radiometric units, except that they are weighted for the spectral response of the human eye and have strange names. A few additional units have been introduced to deal with the amount of light that is reflected from diffuse (matte) surfaces. The symbols used are identical to the geometrically equivalent radiometric symbols, except that a subscript v is added to denote *visual*. Table 2 compares radiometric and photometric units.

The SI unit for light is the *candela* (unit of luminous intensity). It is one of the seven base units of the SI system. The candela is defined as follows:⁵

The candela is the luminous intensity, in a given direction, of a source that emits monochromatic radiation of frequency 540×10^{12} hertz and that has a radiant intensity in that direction of $1/683$ watt per steradian.

The candela is abbreviated as *cd*, and its symbol is I_v . This definition was adopted by the 16th CGPM in 1979. The candela was formerly defined as the luminous intensity, in the perpendicular direction,

TABLE 2 Comparison of Radiometric and Photometric Units

Quantity	Radiometric	Photometric
Power	Φ : watt (W)	Φ_v : lumen (lm)
Power per area	E, M : $W\ m^{-2}$	E_v : $lm\ m^{-2} = lux$ (lx)
Power per solid angle	I : $W\ sr^{-1}$	I_v : $lm\ sr^{-1} = candela$ (cd)
Power per area per solid angle	L : $W\ m^{-2}\ sr^{-1}$	L_v : $lm\ m^{-2}\ sr^{-1} = cd\ m^{-2} = nit$

of a surface of $1/600,000\ m^2$ of a blackbody at the temperature of freezing platinum under a pressure of 101,325 newtons per square meter ($N\ m^{-2}$). This earlier definition was initially adopted in 1948 and later modified by the 13th CGPM in 1968. It was abrogated in 1979 and replaced by the current definition.

The 1979 definition was adopted for several reasons.^{6–10} First, the realization of the candela using a platinum blackbody was extraordinarily difficult—only several were ever built, and there were large variations between the units realized by different national laboratories based upon the state of platinum at its freezing point. The difficulty in fabricating and operating the platinum point blackbody created an unacceptable uncertainty in the value of the candela. For example, if the platinum blackbody temperature is slightly off, possibly because of temperature gradients in the ceramic crucible or contamination of the platinum, the freezing point may change or the temperature of the cavity may differ. The sensitivity of the candela to a slight change in temperature is significant. At a wavelength of 555 nm, a change in temperature of only 1 K results in a luminance change approaching 1 percent. Second, the unit of the candela was realized on the specific broadband radiation, whose spectral power distribution was not known with satisfactory accuracy (because the platinum fix point temperature was not precisely known), thus there were large uncertainties in determining photometric quantities of various other practical light sources from their spectral power distributions. Third, recent advances in radiometry based on absolute radiometers offered new possibilities for realization of the candela using a much simpler device with much lower uncertainties if the candela is defined in relation to watt. In 1977, through an international comparison among several national laboratories, the Comité International des Poids et Mesures (CIPM) determined the numerical relationship (683 lm/W at 555 nm) to be recommended for the new standard for candela so that the magnitude of the unit was kept consistent with the previous unit of the candela.

The value of 683 lm/W was selected based upon the best measurements with existing platinum freezing point blackbodies at several national standards laboratories. It has varied over time from 620 to nearly 700 lm/W, depending largely upon the assigned value of the freezing point of platinum. The value of $1/600,000\ m^2$ was chosen to maintain consistency with prior standards. Note that neither the old nor the new definition of the candela say anything about the spectral responsivity of the human eye. There are additional definitions that include the characteristics of the eye, but the base unit (candela) and those SI units derived from it are “eyeless.”

Note also that in the definition of the candela, there is no specification for the spatial distribution of intensity. Luminous intensity, while often associated with an isotropic point (i.e., small) source, is a valid specification for characterizing any highly directional light source, such as a spotlight or an LED.

One other issue: since the candela is no longer independent but is now defined in terms of other SI-derived quantities, there is really no need to retain it as an SI base quantity. It remains so for reasons of history and continuity and perhaps some politics.

The *lumen* is an SI-derived unit for luminous flux (power). The abbreviation is *lm*, and the symbol is Φ_v . The lumen is derived from the candela and is the luminous flux that is emitted into unit solid angle (1 sr) by an isotropic point source having a luminous intensity of 1 cd. The lumen is the product of luminous intensity and solid angle (cd sr). It is analogous to the unit of radiant flux (watt), differing only in the eye response weighting. If a light source is isotropic, the relationship between lumens and candelas is $1\ cd = 4\pi\ lm$. In other words, an isotropic source that has a luminous intensity of 1 cd emits $4\pi\ lm$ into space, which is $4\pi\ sr$. Also, $1\ cd = 1\ lm\ sr^{-1}$, which is analogous to the equivalent radiometric definition.

If a source is not isotropic, the relationship between candelas and lumens is empirical. A fundamental method used to determine the total flux (lumens) is to measure the luminous intensity (candelas) in many directions using a goniophotometer, and then numerically integrate over the

entire sphere. Later on, this “calibrated” lamp can be used as a reference in an integrating sphere for routine measurements of luminous flux.

The SI-derived unit of luminous flux density, or illuminance, has a special name: *lux*. It is lumens per square meter (lm m^{-2}), and the symbol is E_v . Most light meters measure this quantity, as it is of great importance in illuminating engineering. The IESNA’s *Lighting Handbook*¹¹ has some 16 pages of recommended illuminances for various activities and locales, ranging from morgues to museums. Typical values range from 100,000 lx for direct sunlight to between 20 and 50 lx for hospital corridors at night.

Luminance should probably be included on the list of SI-derived units, but it is not. Luminance is analogous to radiance, giving the spatial and directional dependences. It also has a special name, *nit*, and is candelas per square meter (cd m^{-2}) or lumens per square meter per steradian ($\text{lm m}^{-2} \text{sr}^{-1}$). The symbol is L_v . Luminance is most often used to characterize the “brightness” of flat-emitting or -reflecting surfaces. A common use is the luminance of a laptop computer screen. They typically have between 100 and 250 nits, and the sunlight-readable ones have more than 1000 nits. Typical CRT monitors have luminances between 50 and 125 nits.

Other Photometric Units

There are other photometric units, largely historical. The literature is filled with now obsolete terminology, and it is important to be able to properly interpret these terms. Here are several terms for illuminance that have been used in the past.

$$\begin{aligned} 1 \text{ meter-candle} &= 1 \text{ lx} \\ 1 \text{ phot (ph)} &= 1 \text{ lm cm}^{-2} = 10^4 \text{ lx} \\ 1 \text{ footcandle (fc)} &= 1 \text{ lm ft}^{-2} = 10.76 \text{ lx} \\ 1 \text{ milliphot} &= 10 \text{ lx} \end{aligned}$$

Table 3 is useful to convert from one unit to another. Start with the unit in the leftmost column and multiply it by the factor in the table to arrive at the unit in the top row.

There are two classes of units that are used for luminance. The first is conventional, directly related to the SI unit, the cd m^{-2} (nit).

$$\begin{aligned} 1 \text{ stilb} &= 1 \text{ cd cm}^{-2} = 10^4 \text{ cd m}^{-2} = 10^4 \text{ nit} \\ 1 \text{ cd ft}^{-2} &= 10.76 \text{ cd m}^{-2} = 10.76 \text{ nit} \end{aligned}$$

The second class was designed to “simplify” characterization of light that is reflected from diffuse surfaces by incorporating within the definition the concept of a perfect diffuse reflector (lambertian, reflectance $\rho = 1$). If 1 unit of illuminance falls upon this ideal reflector, then 1 unit of luminance is reflected. The perfect diffuse reflector emits $1/\pi$ units of luminance per unit of illuminance. If the reflectance is ρ , then the luminance is ρ/π times the illuminance. Consequently, these units all incorporate a factor of $1/\pi$.

$$\begin{aligned} 1 \text{ lambert (L)} &= (1/\pi) \text{ cd cm}^{-2} = (10^4/\pi) \text{ cd m}^{-2} = (10^4/\pi) \text{ nit} \\ 1 \text{ apostilb} &= (1/\pi) \text{ cd m}^{-2} = (1/\pi) \text{ nit} \\ 1 \text{ foot-lambert (ft-lambert)} &= (1/\pi) \text{ cd ft}^{-2} = 3.426 \text{ cd m}^{-2} = 3.426 \text{ nit} \\ 1 \text{ millilambert} &= (10/\pi) \text{ cd m}^{-2} = (10/\pi) \text{ nit} \\ 1 \text{ skot} &= 1 \text{ milliblondel} = (10^{-3}/\pi) \text{ cd m}^{-2} = 10^{-3}/\pi \text{ nit} \end{aligned}$$

TABLE 3 Illuminance Unit Conversions

	fc	lx	phot	milliphot
1 fc (lm/ft^2) =	1	10.764	0.0010764	1.0764
1 lx (lm/m^2) =	0.0929	1	0.0001	0.1
1 phot (lm/cm^2) =	929	10,000	1	0.001
1 milliphot =	0.929	10	0.1	1

TABLE 4 Illuminance Unit Conversions*

	nit	stilb	cd/ft ²	apostilb	lambert	ft-lambert
1 nit(cd/m ²) =	1	10 ⁻⁴	0.0929	π	$\pi/10000$	0.0929 π
1 stilb (cd/cm ²) =	10,000	1	929	10 ⁴ π	π	929 π
1 cd/ft ² =	10.764	1.0764 × 10 ⁻³	1	10.764 π	$\pi/929$	π
1 apostilb =	1/ π	10 ⁴ / π	0.0929/ π	1	10 ⁻⁴	0.0929
1 lambert =	10 ⁴ / π	1/ π	929/ π	10 ⁴	1	929
1 ft · lambert =	10.76/ π	1/(929 π)	1/ π	10.764	1.076 × 10 ⁴	1

*Note: Photometric quantities are the result of an integration over wavelength. It therefore makes no sense to speak of spectral luminance or the like.

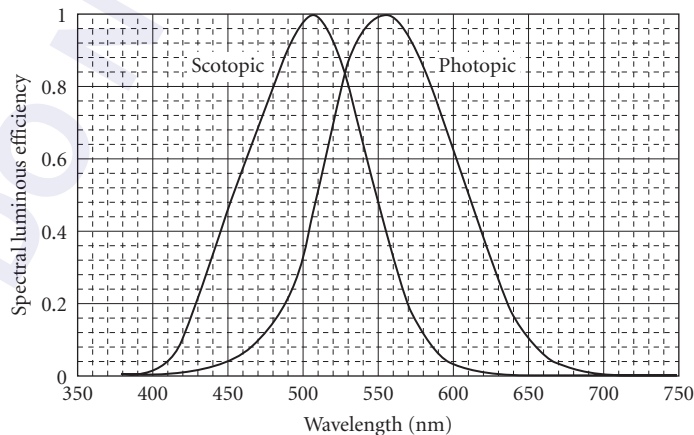
Table 4 is useful to convert from one unit to another. Start with the unit in the leftmost column and multiply it by the factor in the table to arrive at the unit in the top row.

Human Eye The SI base unit and units derived therefrom have a strictly physical basis; they have been defined monochromatically at a wavelength of 555 nm. But the eye does not see all wavelengths equally. For other wavelengths or for band or continuous-source spectral distributions, the spectral properties of the human eye must be considered. The eye has two general classes of photosensors: *cones* and *rods*.

Cones The cones are responsible for light-adapted vision; they respond to color and have high resolution in the central foveal region. The light-adapted relative spectral response of the eye is called the *spectral luminous efficiency function for photopic vision*, $V(\lambda)$, and is published in tabular form.¹² This empirical curve, shown in Fig. 3, was first adopted by the CIE in 1924. It has a peak that is normalized to unity at 555 nm, and it decreases to levels below 10⁻⁵ at about 370 and 785 nm. The 50 percent points are near 510 and 610 nm, indicating that the curve is slightly skewed. A logarithmic representation is shown in Fig. 4.

More recent measurements have shown that the 1924 curve may not best represent typical human vision. It appears to underestimate the response at wavelengths shorter than 460 nm. Judd,¹³ Vos,¹⁴ and Stockman and Sharpe¹⁵ have made incremental advances in our knowledge of the photopic response.

Rods The rods are responsible for dark-adapted vision, with no color information and poor resolution when compared with the foveal cones. The dark-adapted relative spectral response of the

**FIGURE 3** Spectral luminous efficiency for photopic and scotopic vision.

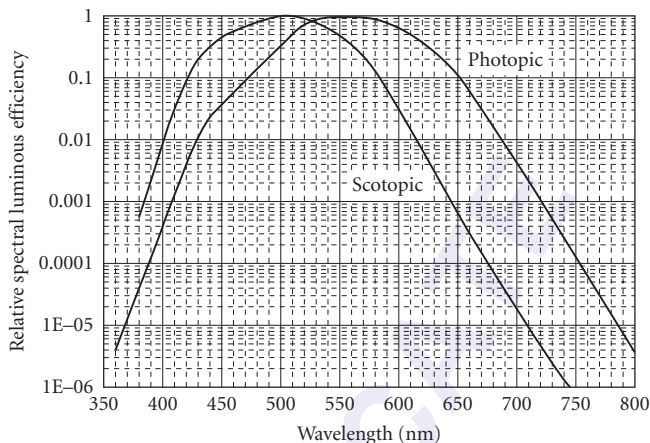


FIGURE 4 Spectral luminous efficiency for photopic and scotopic vision (log scale).

eye is called the *spectral luminous efficiency function for scotopic vision*, $V(\lambda)$, also published in tabular form.¹² Figures 3 and 4 also show this empirical curve, which was adopted by the CIE in 1951. It is defined between 380 and 780 nm. The $V(\lambda)$ curve has a peak of unity at 507 nm, and it decreases to levels below 10^{-3} at about 380 and 645 nm. The 50 percent points are near 455 and 550 nm.

Photopic (light-adapted cone) vision is active for luminances that are greater than 3 cd m^{-2} . Scotopic (dark-adapted rod) vision is active for luminances that are lower than 0.01 cd m^{-2} . In between, both rods and cones contribute in varying amounts, and in this range the vision is called *mesopic*. There have been efforts to characterize the composite spectral response in the mesopic range for vision research at intermediate luminance levels. Definitive values at 1-nm intervals for both photopic and scotopic spectral luminous efficiency functions may be found in CIE.¹² Values at 5-nm intervals are given by Zalewski.¹⁶

The relative spectral luminous efficiency curves can be converted for use with photon flux (s^{-1}) by multiplying by the spectrally dependent conversion from watts to photons per second. The results are shown in Fig. 5. The curves are similar to the spectral luminous efficiency curves, with the peaks shifted to slightly shorter wavelengths, and the skewness of the curves is different. This function

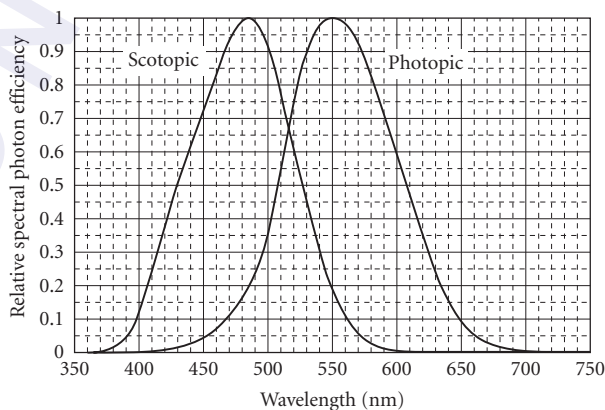


FIGURE 5 Spectral photon efficiency for scotopic and photopic vision.

can be called $V_q(\lambda)$ for photopic response or $V'_q(\lambda)$ for scotopic response. The conversion to an absolute curve is made by multiplying by the response at the peak wavelength. For photopic vision ($\lambda = 550$ nm), $K_{mp} = 2.45 \times 10^{-16}$ lm/photon s^{-1} . There are, therefore, 4.082×10^{15} photon s^{-1} lm^{-1} at 550 nm, and more at all other wavelengths. For scotopic vision ($\lambda_p = 504$ nm), $K'_{mp} = 6.68 \times 10^{-16}$ lm/photon s^{-1} . There are 1.497×10^{15} photon s^{-1} lm^{-1} at 504 nm, and more at all other wavelengths.

Approximations The $V(\lambda)$ curve appears similar to a Gaussian (normal) function. A nonlinear regression technique was used to fit the Gaussian shown in Eq. (2) to the $V(\lambda)$ data

$$V(\lambda) \cong 1.019e^{-285.4(\lambda - 0.559)^2} \quad (2)$$

The scotopic curve can also be fit with a Gaussian, although the fit is not quite as good as the photopic curve. My best fit is

$$V'(\lambda) \cong 0.992e^{-321.0(\lambda - 0.503)^2} \quad (3)$$

The results of the curve fitting are shown in Figs. 6 and 7. These approximations are satisfactory for application with continuous spectral distributions, such as sunlight, daylight, and incandescent sources. Calculations have demonstrated errors of less than 1 percent with blackbody sources from 1500 K to more than 20,000 K. The equations must be used with caution for narrow-band or line sources, particularly in those spectral regions where the response is low and the fit is poor.

Usage The SI definition of the candela was chosen in strictly physical terms at a single wavelength. The intent of photometry, however, is to correlate a photometric observation to the visual perception of a human observer. The CIE introduced the two standard spectral luminous efficiency functions $V(\lambda)$ (photopic) and $V'(\lambda)$ (scotopic) as spectral weighting functions, and they have been approved by the CIPM for use with light sources at other wavelengths. Another useful function is the CIE $V_M(\lambda)$ Judd-Vos modified $V(\lambda)$ function,¹⁴ which has increased response at wavelengths that are shorter than 460 nm. It is identical to the $V(\lambda)$ function for wavelengths that are longer than 460 nm. This function, while not approved by CIPM, represents more realistically the spectral responsivity of the eye. More recently, studies on cone responses have led to the proposal of a new, improved luminous spectral efficiency curve, with the suggested designation $V_2^*(\lambda)$.¹⁵

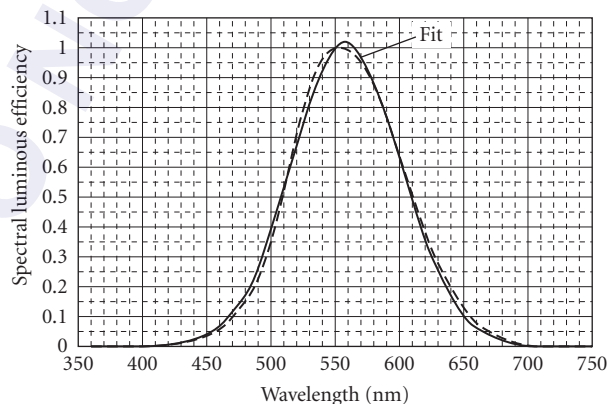


FIGURE 6 Gaussian fit to photopic relative spectral efficiency curve.

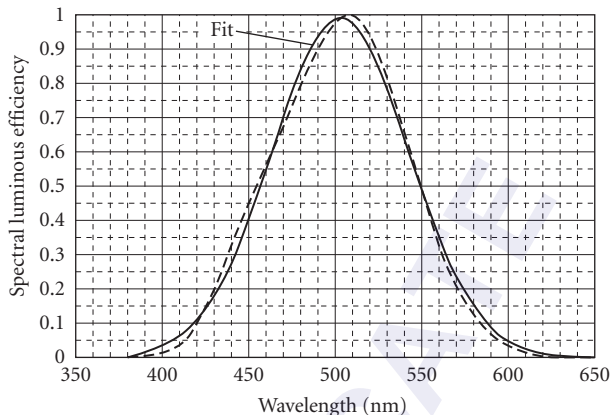


FIGURE 7 Gaussian fit to scotopic relative spectral efficiency curve.

36.5 CONVERSION OF RADIOMETRIC QUANTITIES TO PHOTOMETRIC QUANTITIES

The definition of the candela states that there are 683 lm W^{-1} at a frequency of 540 terahertz (THz), which is very nearly 555 nm (vacuum or air), the wavelength that corresponds to the maximum spectral responsivity of the photopic (light-adapted) human eye. The value 683 lm W^{-1} is K_m , the absolute luminous efficiency at λ_p for photopic vision. The conversion from watts to lumens at any other wavelength involves the product of the power (watts), K_m , and the $V(\lambda)$ value at the wavelength of interest. For example, a 5-mW laser pointer has $0.005 \text{ W} \times 0.032 \times 683 \text{ lm W}^{-1} = 0.11 \text{ lm}$. $V(\lambda)$ is 0.032 at 670 nm. At 635 nm, $V(\lambda)$ is 0.217, and a 5-mW laser pointer has $0.005 \text{ W} \times 0.217 \times 683 \text{ lm W}^{-1} = 0.74 \text{ lm}$. The shorter-wavelength laser pointer will create a spot that has nearly seven times the luminous power as the longer-wavelength laser.

Similar calculations can be done in terms of photon flux at a single wavelength. As was shown previously, there are $2.45 \times 10^{16} \text{ lm}$ in 1 photon s^{-1} at 550 nm, the wavelength that corresponds to the maximum spectral responsivity of the light-adapted human eye to photon flux. The conversion from lumens to photons per second at any other wavelength involves the product of the photon flux (s^{-1}) and the $V_p(\lambda)$ value at the wavelength of interest. For example, again compare laser pointers at 670 and 635 nm. As shown before, a 5-mW laser at 670 nm [$V_p(\lambda) = 0.0264$] has a luminous power of 0.11 lm. The conversion is $0.11 \times 4.082 \times 10^{15} / 0.0264 = 1.68 \times 10^{16} \text{ photon s}^{-1}$. At 635 nm [$V_p(\lambda) = 0.189$], the 5-mW laser has a luminous power of 0.74 lm. The conversion is $0.74 \times 4.082 \times 10^{15} / 0.189 = 1.6 \times 10^{16} \text{ photon s}^{-1}$. The 635-nm laser delivers just 5 percent more photons per second.

In order to convert a source with nonmonochromatic spectral distribution to a luminous quantity, the situation is decidedly more complex. The spectral nature of the source must be known, as it is used in an equation of the form

$$X_v = K_m \int_0^{\infty} X_\lambda V(\lambda) d\lambda \quad (4)$$

where X_v is a luminous term, X_λ is the corresponding spectral radiant term, and $V(\lambda)$ is the photopic spectral luminous efficiency function. For X_v , luminous flux (lm) may be paired with spectral power (W nm^{-1}), luminous intensity (cd) with spectral radiant intensity ($\text{W sr}^{-1} \text{ nm}^{-1}$), illuminance (lx) with spectral irradiance ($\text{W m}^{-2} \text{ nm}^{-1}$), or luminance (cd m^{-2}) with spectral radiance ($\text{W m}^{-2} \text{ sr}^{-1} \text{ nm}^{-1}$). This equation represents a weighting, wavelength by wavelength, of the radiant spectral term

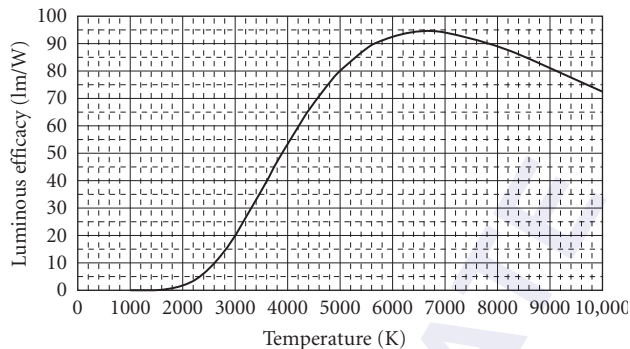


FIGURE 8 Luminous efficacy of blackbody radiation versus temperature (K).

by the visual response at that wavelength. The constant K_m is the maximum spectral luminous efficiency for photopic vision, 683 lm W^{-1} . The wavelength limits can be set to restrict the integration to only those wavelengths where the product of the spectral term X_λ and $V(\lambda)$ is nonzero. Practically, the limits of integration need only extend from 360 to 830 nm, limits specified by the CIE $V(\lambda)$ function. Since this $V(\lambda)$ function is defined by a table of empirical values,¹² it is best to do the integration numerically. Use of the Gaussian equation [Eq. (2)] is only an approximation.

For source spectral distributions that are blackbody-like (thermal source, spectral emissivity constant between 360 and 830 nm) and of known source temperature, it is straightforward to convert from power to luminous flux and vice versa. Equation (4) is used to determine a scale factor for the source term X_λ . Figure 8 shows the relationship between total power and luminous flux for blackbody (and graybody) radiation as a function of blackbody temperature. The most efficient temperature for the production of luminous flux is near 6630 K.

There is nothing in the SI definitions of the base or derived units concerning the eye response, so there is some flexibility in the choice of the weighting function. The choice can be made to use a different spectral luminous efficiency curve, perhaps one of the newer ones. The equivalent curve for scotopic (dark-adapted) vision can also be used for work at lower light levels. The $V'(\lambda)$ curve has its own constant, K'_m , the maximum spectral luminous efficiency for scotopic vision. K'_m is defined as 1700 lm/W at the peak wavelength for scotopic vision (507 nm). This value was deliberately chosen such that the absolute value of the scotopic curve at 555 nm coincides with the photopic curve, 683 lm/W at 555 nm. Some researchers are referring to “scotopic lumens,” a term that should be discouraged because of the potential for misunderstanding. In the future, expect to see spectral weighting to represent the mesopic region as well.

The CGPM has approved the use of the CIE $V(\lambda)$ and $V'(\lambda)$ curves for determination of the value of photometric quantities of luminous sources.

36.6 CONVERSION OF PHOTOMETRIC QUANTITIES TO RADIOMETRIC QUANTITIES

The conversion from watts to lumens in the previous section required only that the spectral function, X_λ , of the radiation be known over the spectral range from 360 to 830 nm, where $V(\lambda)$ is nonzero. Attempts to go in the other direction, from lumens to watts, are far more difficult. Since the desired quantity was inside of an integral, weighted by a sensor spectral responsivity function, the spectral function, X_λ , of the radiation must be known over the entire spectral range where the source emits, not just the visible.

For a monochromatic source in the visible spectrum (between the wavelengths of 380 and 860 nm), if the photometric quantity (e.g., lux) is known, apply the conversion $K_m \times V(\lambda)$ and determine the

radiometric quantity (e.g., $W\ m^{-2}$). In practice, the results that one obtains are governed by the quality of the $V(\lambda)$ correction of the photometer and the knowledge of the wavelength of the source. Both of these factors are of extreme importance at wavelengths where the $V(\lambda)$ curve is steep (i.e., other than very close to the peak of the $V(\lambda)$ curve).

Narrowband sources, such as LEDs, cause major problems. Typical LEDs have spectral bandwidths ranging from 10- to 40-nm full width at half-maximum (FWHM). It is intuitive that in those spectral regions where the $V(\lambda)$ curve is steep, the luminous output will be greater than that predicted using the $V(\lambda)$ curve at the peak LED wavelength. This expected result increases with wider-bandwidth LEDs. Similarly, it is also intuitive that the luminous output is less than that predicted by using the $V(\lambda)$ curve when the peak LED wavelength is in the vicinity of the peak of the $V(\lambda)$ curve. Therefore, there must be two wavelengths where the conversion ratio (lm/W) is largely independent of LED bandwidth. An analysis of this conversion ratio was done using a Gaussian equation to represent the spectral power distribution of an LED and applying Eq. (4). Indeed, two null wavelengths were identified (513 and 604 nm) where the conversion between radiometric and photometric quantities is constant (independent of LED bandwidth) to within 0.2 percent up to an LED bandwidth of 40 nm. These wavelengths correspond (approximately) to the wavelengths where the two maxima of the first derivative of the $V(\lambda)$ curve are located.

At wavelengths between these two null wavelengths (513 and 604 nm), the conversion ratio (lm/W) decreases slightly with increasing bandwidth. The worst case occurs when the peak wavelength of the LED corresponds with the peak of $V(\lambda)$. It is about 5 percent lower for bandwidths up to 30 nm, increasing to nearly 10 percent for 40-nm bandwidth. At wavelengths outside of the null wavelengths, the conversion (lm/W) increases with increasing bandwidth, and the increase is greater when the wavelength approaches the limits of the $V(\lambda)$ curve. Figure 9 shows that factor by which a conversion ratio (lm/W) should be multiplied as a function of LED bandwidth, with the peak LED wavelength as the parameter. Note that the peak wavelength of the LED is specified as the radiometric peak and not as the dominant wavelength (a color specification). The dominant wavelength shifts with respect to the radiometric peak, the difference increasing with bandwidth.

Most often, LEDs are specified in luminous intensity [cd or $millicandela$ (mcd)]. The corresponding radiometric unit is watts per steradian (W/sr). In order to determine the radiometric power (watts), the details of the spatial distribution of the radiant intensity must be known prior to integration.

For broadband sources, a photometric quantity cannot in general be converted to a radiometric quantity unless the radiometric source function, Φ_λ , is known over all wavelengths. However, if the source spectral distribution is blackbody-like (thermal source, spectral emissivity constant between 360 and 830 nm), and the source temperature is also known, then an illuminance can be converted to a spectral irradiance curve over that wavelength range. Again, use Eq. (4) to determine a scale factor

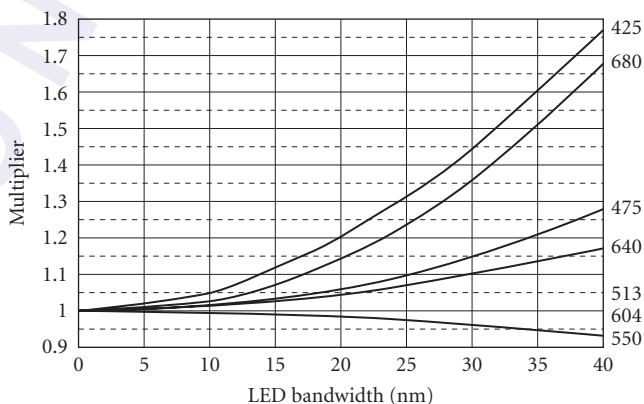


FIGURE 9 Multiplier for converting LED luminous intensity to radiometric intensity.

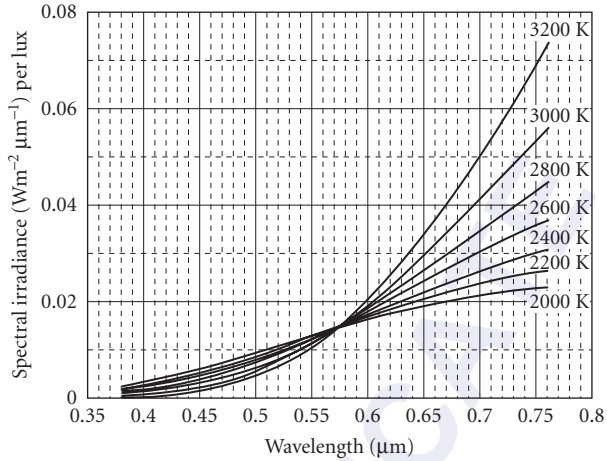


FIGURE 10 Spectral irradiance versus wavelength of blackbody radiation versus temperature.

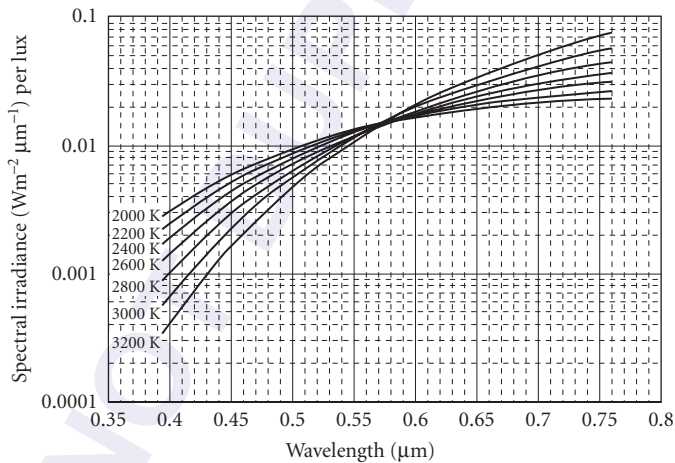


FIGURE 11 Spectral irradiance versus wavelength of blackbody radiation versus temperature.

for the source term, X_λ . Figures 10 and 11 show the calculated spectral irradiance versus wavelength for blackbody radiation with source temperature as the parameter.

36.7 RADIOMETRIC/PHOTOMETRIC NORMALIZATION

In radiometry and photometry, there are two mathematical activities that fall into the category called *normalization*. The first is bandwidth normalization. The second doesn't have an official title but involves a conversion between photometric and radiometric quantities, used to convert detector responsivities from amperes per lumen to amperes per watt.

The measurement equation relates the output signal from a sensor to its spectral responsivity and the spectral power that is incident upon it:

$$S = \int_0^{\infty} \Phi_{\lambda} \mathfrak{R}(\lambda) d\lambda \quad (5)$$

An extended discourse on bandwidth normalization¹⁷ showed that the spectral responsivity of a sensor (or detector) can be manipulated to yield more source information than is immediately apparent from the measurement equation. The sensor weights the input spectral power according to its spectral responsivity, $\mathfrak{R}(\lambda)$, such that only a limited amount of information about the source can be deduced. If either the source or the sensor has a sufficiently small bandwidth such that the spectral function of the other term does not change significantly over the passband, the equation simplifies to

$$S = \Phi_{\lambda} \cdot \mathfrak{R}(\lambda) \cdot \Delta\lambda \quad (6)$$

where $\Delta\lambda$ is the passband. Spectroradiometry and multifilter radiometry, using narrow-bandpass filters, take advantage of this simplified equation. For those cases where the passband is larger, the techniques of bandwidth normalization can be used. The idea is to substitute for $\mathfrak{R}(\lambda)$ an equivalent response that has a uniform spectral responsivity, \mathfrak{R}_n , between wavelength limits λ_1 and λ_2 and zero response elsewhere. Then, the signal is given by

$$S = \mathfrak{R}_n \int_{\lambda_1}^{\lambda_2} \Phi_{\lambda} d\lambda \quad (7)$$

and now the integrated power between wavelengths λ_1 and λ_2 is determined. There are many ways of assigning values for λ_1 , λ_2 , and \mathfrak{R}_n for a sensor. Some of the more popular methods were described by Nicodemus¹⁷ and Palmer.¹⁸ An effective choice is known as the *moments method*,¹⁹ an analysis of the zeroth, first, and second moments of the sensor spectral responsivity curve. The derivation of this normalization scheme involves the assumption that the source function is exactly represented by a second-degree polynomial. If this condition is met, the moments method of determining sensor parameters yields exact results for the source integral. In addition, the results are completely independent of the source function. The errors encountered are related to deviation of the source function from the said second-degree polynomial.

Moments normalization has been applied to the photopic spectral luminous efficiency function, $V(\lambda)$, and the results are given in Table 5 and shown in Fig. 12. These values indicate the skewed nature of the photopic and scotopic curves as the deviation from the centroid and the peak wavelengths. The results can be applied to most continuous sources, like blackbody and tungsten radiation, which are both continuous across the visible spectrum. To demonstrate the effectiveness of moments normalization, the blackbody curve was multiplied by the $V(\lambda)$ curve for temperatures ranging from 1000 to 20,000 K to determine a photometric function [e.g., lumens per square meter (or lux)]. Then, the blackbody curve was integrated over the wavelength interval between λ_1 and λ_2 to determine the equivalent (integrated between λ_1 and λ_2) radiometric

TABLE 5 Bandwidth Normalization on Spectral Luminous Efficiency

	Photopic	Scotopic
Peak wavelength (λ_p)	555 nm	507 nm
Centroid wavelength (λ_c)	560.19 nm	502.40 nm
Short wavelength (λ_1)	487.57 nm	436.88 nm
Long wavelength (λ_2)	632.81 nm	567.93 nm
Moments bandwidth	145.24 nm	131.05 nm
Normalized \mathfrak{R}_n	0.7357	0.7407
Absolute \mathfrak{R}_n	502.4 lm/W	1260 lm/W

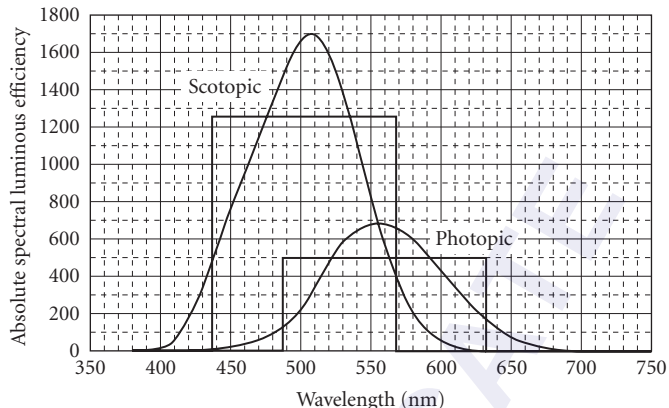


FIGURE 12 Absolute spectral luminous efficiency functions and equivalent normalized bandwidths

function (e.g., in-band watts per square meter). The ratio of lux to watt per square meter is 502.4 ± 1.0 (3σ) over the temperature range from 1600 to more than 20,000 K. This means that the in-band (487.6 to 632.8 nm) irradiance for a continuous blackbody-like source can be determined using a photometer that is properly calibrated in lux and that is well corrected for $V(\lambda)$. Simply divide the reading in lux by 502.4 to get the in-band irradiance in watts per square meter between 487.6 and 632.8 nm.

If a photometer is available with a $V'(\lambda)$ correction, calibrated for lux (scotopic) with K'_m of 1700 lm/W, a similar procedure is effective. The integration takes place over the wavelength range from 436.9 to 567.9 nm. The ratio of lux (scotopic) to watts per square meter is 1260 ± 2 (3σ) over the temperature range from 1800 K to more than 20,000 K. This means that the in-band (436.9 to 567.9 nm) irradiance for a continuous blackbody-like source can be determined using a $V'(\lambda)$ -corrected photometer that is properly calibrated in lux (scotopic). Simply divide the reading in lux (scotopic) by 1260; the result is the in-band irradiance in watts per square meter between 436.9 and 567.9 nm.

A common problem is the interpretation of specifications for photodetectors, which are given in photometric units. An example is a photomultiplier with an S-20 photocathode, which has a typical responsivity of 200 $\mu\text{A}/\text{lm}$. Given this specification and a curve of the relative spectral responsivity, the problem is to determine the output when exposed to a known power from an arbitrary source.

Photosensitive devices, in particular vacuum photodiodes and photomultiplier tubes, are characterized using CIE Illuminant A, a tungsten source at 2854 K color temperature. The illuminance is measured using a photooptically corrected photometer, and this illuminance is applied to the device under scrutiny. This technique is satisfactory only if the source being used is spectrally comparable to Illuminant A. If a source with a different spectral distribution is used, a photometric normalization must be done. Eberhart²⁰ generated a series of conversion factors for various sources and standardized detector spectral responsivities (S-1, S-11, S-20, etc.).

The luminous flux from any source is given by

$$\Phi_v = K_m \int_{360}^{830} \Phi_\lambda V(\lambda) d\lambda \quad (8)$$

and the output of a detector when exposed to the said luminous flux is

$$S = \int_0^\infty \Phi_\lambda \mathfrak{R}(\lambda) d\lambda \quad (9)$$

where Φ_λ is spectral radiant flux, $V(\lambda)$ is the spectral luminous efficiency of the photopic eye, $\mathfrak{R}(\lambda)$ is the absolute spectral responsivity of the photodetector, and S is the photodetector signal. The luminous responsivity of the detector when exposed to this source is

$$\mathfrak{R} = \frac{\int_0^\infty \Phi_\lambda \mathfrak{R}(\lambda) d\lambda}{K_m \int_{360}^{830} \Phi_\lambda V(\lambda) d\lambda} \text{ A/lm} \quad (10)$$

This luminous responsivity is specific to the source that is used to make the measurement, and it cannot be applied to other sources with differing spectral distributions.

36.8 OTHER WEIGHTING FUNCTIONS AND CONVERSIONS

The general principles that are outlined here can be applied to action spectra other than those already defined for the human eye (photopic and scotopic). Some action spectra take the form of defined spectral regions, such as UVA (315–400 nm), UVB (280–315 nm), UVC (100–280 nm), IR-A (770–1400 nm), IR-B (1400–3000 nm), and IR-C (3000–10⁶ nm). Others are more specific. $A(\lambda)$ is for aphakic hazard, $B(\lambda)$ is for photochemical blue-light hazard, $R(\lambda)$ is for retinal thermal hazard, and $S(\lambda)$ is an actinic ultraviolet action spectrum.²¹ PPF (a.k.a. PhAR) is a general action spectrum for plant growth. Many others have been defined, including those for erythema (sunburn), skin cancer, psoriasis treatment, mammalian and insect vision, and other generalized chemical and biological photoeffects. Various conversion factors from one action spectrum to another are scattered throughout the popular and archival literature. Ideally, they have been derived via integration over the appropriate spectral regions.

36.9 REFERENCES

1. ISO, "Quantities and Units—Part 6. Light and Related Electromagnetic Radiations," *ISO Standards Handbook, Quantities and Units*, 389.15 (1993).
2. IUPAP, *Symbols, Units, Nomenclature and Fundamental Constants in Physics*, prepared by E. Richard Cohan and Pierre Giacomo, Document IUPAP-25 (SUNAMCO 87-1), International Union of Pure and Applied Physics, 1987.
3. NIST, *Guide for the Use of the International System of Units (SI)*, prepared by Barry N. Taylor, NIST Special Publication SP811 (1995). (Available in PDF format from <http://physics.nist.gov/cuu/>.)
4. J. M. Palmer, "Getting Intense on Intensity," *Metrologia* **30**:371 (1993).
5. BIPM (Bureau International des Poids et Mesures), *The International System of Units (SI)*, 7th ed., 1998.
6. <http://www.bipm.org/en/CGPM/db/16/3/>.
7. W. R. Blevin and B. Steiner, "The Redefinition of the Candela and the Lumen," *Metrologia* **11**:97–104 (1975).
8. O. C. Jones, "Proposed Changes to the SI System of Photometric Units," *Lighting Research and Technology* **10**:37–40 (1978).
9. W. R. Blevin, "The Candela and the Watt," *CIE Proc.* P-79-02 (1979).
10. CGPM, *Comptes Rendus des Séances de la 16e Conférence Générale des Poids et Mesures*, Paris 1979, BIPM, Sèvres, France (1979).
11. IESNA, *Lighting Handbook: Reference and Application*, M. S. Rea, ed., Illuminating Engineering Society of North America, New York, 1993.
12. CIE, *The Basis of Physical Photometry*, CIE 18.2, Vienna, 1983.

13. D. B. Judd, "Report of U.S. Secretariat Committee on Colorimetry and Artificial Daylight," *Proceedings of the Twelfth Session of the CIE (Stockholm)*, CIE, Paris, 1951. (The tabular data is given in G. Wyszecki and W. S. Stiles, *Color Science*, 2nd ed., Wiley, New York, 1982.)
14. J. J. Vos, "Colorimetric and Photometric Properties of a 2-Degree Fundamental Observer," *Color Research and Application* 3:125 (1978).
15. A. Stockman and L. T. Sharpe, "Cone Spectral Sensitivities and Color Matching," *Color Vision: From Genes to Perception*, K. Gegenfurtner and L. T. Sharpe, eds., Cambridge, 1999. (This information is summarized at the Color Vision Lab at UCSD Web site located at cvision.uscd.edu/.)
16. E. F. Zalewski, "Radiometry and Photometry," chap. 24, *Handbook of Optics*, vol. II, McGraw-Hill, New York, 1995.
17. F. E. Nicodemus, "Normalization in Radiometry," *Appl. Opt.* 12:2960 (1973).
18. J. M. Palmer, "Radiometric Bandwidth Normalization Using r.m.s. Methods," *Proc. SPIE* 256:99 (1980).
19. J. M. Palmer and M. G. Tomasko, "Broadband Radiometry with Spectrally Selective Detectors," *Opt. Lett.* 5:208 (1980).
20. E. H. Eberhart, "Source-Detector Spectral Matching Factors," *Appl. Opt.* 7:2037 (1968).
21. ANSI, *Photobiological Safety of Lamps*, American National Standards Institute RP27.3-96 (1996).

36.10 FURTHER READING

Books, Documentary Standards, Significant Journal Articles

- American National Standard Nomenclature and Definitions for Illuminating Engineering*, ANSI Standard ANSI/IESNA RP-16 96.
- W R. Blevin and B. Steiner, "Redefinition of the Candela and the Lumen," *Metrologia* 11:97 (1975).
- C. DeCusatis, *Handbook of Applied Photometry*, AIP Press, 1997. (Authoritative, with pertinent chapters written by technical experts at BIPM, CIE, and NIST. Skip chapter 4!)
- J. W. T. Walsh, *Photometry*, Constable, London, 1958. (The classic!)

Publications Available on the World Wide Web

- All you ever wanted to know about the SI is contained at BIPM and at NIST. Available publications (highly recommended) include the following:
- "The International System of Units (SI)," 7th ed. (1998), direct from BIPM. This is the English translation of the official document, which is in French. Available in PDF format at www.bipm.fr/.
- NIST Special Publication SP330, "The International System of Units (SI)." The U.S. edition (meter rather than metre) of the above BIPM publication. Available in PDF format from <http://physics.nist.gov/cuu/>.
- NIST Special Publication SP811, "Guide for the Use of the International System of Units (SI)," Available in PDF format from <http://physics.nist.gov/cuu/>.
- Papers published in recent issues of the NIST Journal of Research are also available on the Web in PDF format from mvl.nist.gov/pub/nistpubs/jres/jres.htm. Of particular relevance is "The NIST Detector-Based Luminous Intensity Scale," vol. 101, p. 109 (1996).

Useful Web Sites

- AIP (American Institute of Physics): www.aip.org
- ANSI (American National Standards Institute): www.ansi.org/
- BIPM (International Bureau of Weights and Measures): www.bipm.fr/

CIE (International Commission on Illumination): www.de.co.at/cie/

Color Vision Lab at UCSD: cvision.uscd.edu/

CORM (Council for Optical Radiation Measurements): www.corm.org

IESNA (Illuminating Engineering Society of North America): www.iesna.org/

ISO (International Standards Organization): www.iso.ch/

IUPAP (International Union of Pure and Applied Physics): www.physics.umanitoba.ca/iupap/

NIST (National Institute of Standards and Technology): physics.nist.gov/

OSA (Optical Society of America): www.osa.org

SPIE (International Society for Optical Engineering): www.spie.org

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

RADIOMETRY AND PHOTOMETRY FOR VISION OPTICS*

Yoshi Ohno

*Optical Technology Division
National Institute of Standards and Technology
Gaithersburg, Maryland*

37.1 INTRODUCTION

Radiometry is the measurement of optical radiation, which is electromagnetic radiation in the frequency range between 3×10^{11} Hz and 3×10^{16} Hz. This range corresponds to wavelengths between 10 nm and 1000 μm , and includes the regions commonly called the ultraviolet, the visible, and the infrared. Typical radiometric units include watt (radiant flux), watt per steradian (radiant intensity), watt per square meter (irradiance), and watt per square meter per steradian (radiance).

Photometry is the measurement of light, which is defined as electromagnetic radiation detectable by the human eye. It is thus restricted to the visible region of the spectrum (wavelength range from 360 nm to 830 nm), and all the quantities are weighted by the spectral response of the eye. Photometry uses either optical radiation detectors constructed to mimic the spectral response of the eye, or spectroradiometry coupled with appropriate calculations for weighting by the spectral response of the eye. Typical photometric units include lumen (luminous flux), candela (luminous intensity), lux (illuminance), and candela per square meter (luminance).

The difference between radiometry and photometry is that radiometry includes the entire optical radiation spectrum (and often involves spectrally resolved measurements), while photometry deals with the visible spectrum weighted by the response of the eye. This chapter provides some guidance in photometry and radiometry (Refs. 1 through 6 are available for further details). The terminology used in this chapter follows international standards and recommendations.⁷⁻⁹

37.2 BASIS OF PHYSICAL PHOTOMETRY

The primary aim of photometry is to measure visible optical radiation, light, in such a way that the results correlate with the visual sensation of a normal human observer exposed to that radiation. Until about 1940, visual comparison measurement techniques were predominant in photometry.

*Chapter 1 in Vol. III gives further treatment of issues in radiometry, and Chap. 35 in this volume describes measurement of optical properties of materials.

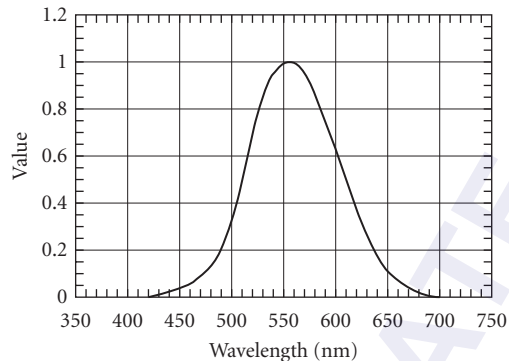


FIGURE 1 CIE $V(\lambda)$ function.

In modern photometric practice, measurements are made with photodetectors. This is referred to as *physical photometry*. In order to achieve the aim of photometry, one must take into account the characteristics of human vision. The relative spectral responsivity of the human eye was first defined by the CIE (Commission Internationale de l'Éclairage) in 1924,¹⁰ and redefined as part of colorimetric standard observers in 1931.¹¹ Called the *spectral luminous efficiency for photopic vision*, or $V(\lambda)$, it is defined in the domain from 360 nm to 830 nm, and is normalized at its peak, 555 nm (Fig. 1). This model has gained wide acceptance. The values were republished by CIE in 1983,¹² and published by CIPM (Comité International des Poids et Mesures) in 1982¹³ to supplement the 1979 definition of the candela. (The tabulated values of the function at 1-nm increments are available in Refs. 12 through 15.) In most cases, the region from 380 nm to 780 nm suffices for calculation with negligible errors because the value of the $V(\lambda)$ function falls below 10^{-4} outside this region. Thus, a photodetector having a spectral responsivity matched to the $V(\lambda)$ function replaced the role of human eyes in photometry.

Radiometry concerns physical measurement of optical radiation in terms of optical power, and in many cases, as a function of its wavelength. As specified in the definition of the candela by CGPM (Conférence Générale des Poids et Mesures) in 1979¹⁶ and by CIPM in 1982,¹³ a photometric quantity X_v is defined in relation to the corresponding radiometric quantity $X_{e,\lambda}$ by the equation:

$$X_v = K_m \int_{360\text{nm}}^{830\text{nm}} X_{e,\lambda} V(\lambda) d\lambda \quad (1)$$

The constant, K_m , relating the photometric quantities and radiometric quantities, is called the *maximum spectral luminous efficacy (of radiation) for photopic vision*. The value of K_m is given by the 1979 definition of candela that defines the spectral luminous efficacy of light at the frequency 540×10^{12} Hz (at the wavelength 555.016 nm in standard air) to be 683 lm/W. The value of K_m is calculated as $683 \times V(555.000 \text{ nm})/V(555.016 \text{ nm}) = 683.002 \text{ lm/W}$.¹² K_m is normally rounded to 683 lm/W with negligible errors.

It should be noted that $V(\lambda)$ is defined for the *CIE standard photometric observer for photopic vision*, which assumes additivity of sensation and a 2° field of view at relatively high luminance levels (higher than approximately 1 cd/m^2). The human vision in this level is called photopic vision. The spectral responsivity of human vision deviates significantly at very low levels of luminance (less than approximately 10^{-3} cd/m^2). This type of vision is called scotopic vision. Its spectral responsivity, peaking at 507 nm, is designated by $V'(\lambda)$, which was defined by CIE in 1951,¹⁷ recognized by CIPM in 1976,¹⁸ and republished by CIPM in 1982.¹³ Human vision in the region between photopic vision and scotopic vision is called mesopic vision. While there have been active researches in this area,¹⁹ there is no internationally accepted spectral luminous efficiency function for the mesopic region yet. In current practice, almost all photometric quantities are given in terms of photopic vision, even at low light levels. Quantities in scotopic vision are seldom used except for special calculations for research purposes. (Further details of the contents in this section are given in Ref. 12.)

37.3 PHOTOMETRIC BASE UNIT—THE CANDELA

The history of photometric standards dates back to the early nineteenth century, when the intensity of light sources was measured in comparison with a standard candle using visual bar photometers. At that time, the flame of a candle was used as a unit of luminous intensity that was called the *candle*. The old name for luminous intensity *candle power* came from this origin. Standard candles were gradually superseded by flame standards of oil lamps, and in the early twentieth century investigations on platinum point blackbodies began at some national laboratories. An agreement was first established in 1909 among several national laboratories to use such a blackbody to define the unit of luminous intensity, and the unit was recognized as the *international candle*. This standard was adopted by the CIE in 1921. In 1948, it was adopted by the CGPM¹⁶ with a new Latin name *candela* with the following definition:

The candela is the luminous intensity, in the perpendicular direction, of a surface of 1/600000 square meter of a blackbody (full radiator) at the temperature of freezing platinum under a pressure of 101325 newton per square meter.

Although the 1948 definition served to establish the uniformity of photometric measurements in the world, difficulties in fabricating the blackbodies and in improving accuracy were addressed. Beginning in the mid-1950s, suggestions were made to define the candela in relation to the unit of optical power, watt, so that complicated source standards would not be necessary. Finally, in 1979, a new definition of the candela was adopted by the CGPM¹⁶ as follows:

The candela is the luminous intensity, in a given direction, of a source that emits monochromatic radiation of frequency 540×10^{12} hertz and that has a radiant intensity in that direction of (1/683) watt per steradian.

The value of K_m (683 lm/W) was determined in such a way that the consistency from the prior unit was maintained, and was determined based on the measurements by several national laboratories. (Technical details on this redefinition of the candela are reported in Refs. 20 and 21.) This 1979 redefinition of the candela has enabled photometric units to be derived from radiometric units using a variety of techniques. (The early history of photometric standards is described in greater detail in Ref. 22.)

37.4 QUANTITIES AND UNITS IN PHOTOMETRY AND RADIOMETRY

In 1960, the SI (Système International) was established, and the candela became one of the seven SI base units.²³ (For further details on the SI, Refs. 23 through 26 may be consulted.) Several quantities and units, defined in different geometries, are used in photometry and radiometry. Table 1 lists the photometric quantities and units, along with corresponding quantities and units for radiometry.

While the candela is the SI base unit, the luminous flux (lumen) is perhaps the most fundamental photometric quantity, as the other photometric quantities are defined in terms of lumen with an appropriate geometric factor. The definitions of these photometric quantities are given in the following sections. (The descriptions given here are somewhat simplified from the rigorous definitions for ease of understanding. Refer to Refs. 7 through 9 for official rigorous definitions.)

Radiant Flux and Luminous Flux

Radiant flux (also called *optical power* or *radiant power*) is the energy Q (in joules) radiated by a source per unit of time, expressed as

$$\Phi = \frac{dQ}{dt} \quad (2)$$

The unit of radiant flux is the *watt* ($W = J/s$).

TABLE 1 Quantities and Units Used in Photometry and Radiometry

Photometric Quantity	Unit	Relationship with Lumen	Radiometric Quantity	Unit
Luminous flux	lm (lumen)		Radiant flux	W (watt)
Luminous intensity	cd (candela)	lm sr ⁻¹	Radiant intensity	W sr ⁻¹
Illuminance	lx (lux)	lm m ⁻²	Irradiance	W m ⁻²
Luminance	cd m ⁻²	lm sr ⁻¹ m ⁻²	Radiance	W sr ⁻¹ m ⁻²
Luminous exitance	lm m ⁻²		Radiant exitance	W m ⁻²
Luminous exposure	lx · s		Radiant exposure	W m ⁻² · s
Luminous energy	lm · s		Radiant energy	J (joule)
Total luminous flux	lm (lumen)		Total radiant flux	W (watt)
Color temperature	K (kelvin)		Radiance temperature	K (kelvin)

Luminous flux (Φ_v) is the time rate of flow of light as weighted by $V(\lambda)$. The unit of luminous flux is the *lumen* (lm). It is defined as

$$\Phi_v = K_m \int_{\lambda} \Phi_{e,\lambda} V(\lambda) d\lambda \tag{3}$$

where ($\Phi_{e,\lambda}$) is the spectral concentration of radiant flux as a function of wavelength λ . The term luminous flux is often used in the meaning of total luminous flux in photometry (see the following subsection entitled “Total Radiant Flux and Total Luminous Flux”).

Radiant Intensity and Luminous Intensity

Radiant intensity (I_e) or luminous intensity (I_v) is the radiant flux (luminous flux) from a point source emitted per unit solid angle in a given direction, as defined by

$$I = \frac{d\Phi}{d\Omega} \tag{4}$$

where $d\Phi$ is the radiant flux (luminous flux) leaving the source and propagating in an element of solid angle $d\Omega$ containing the given direction. The unit of radiant intensity is W/sr, and that of luminous intensity is the *candela* (cd = lm/sr). (See Fig. 2.)

Solid Angle The solid angle (Ω) of a cone is defined as the ratio of the area (A) cut out on a spherical surface (with its center at the apex of that cone) to the square of the radius (r) of the sphere, as given by

$$\Omega = \frac{A}{r^2} \tag{5}$$

The unit of solid angle is *steradian* (sr), which is a dimensionless unit. (See Fig. 3.)

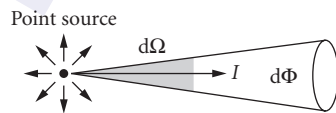


FIGURE 2 Radiant intensity and luminous intensity.

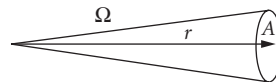


FIGURE 3 Solid angle.

Irradiance and Illuminance

Irradiance (E_e) or *illuminance* (E_v) is the density of incident radiant flux or luminous flux at a point on a surface, and is defined as radiant flux (luminous flux) per unit area, as given by

$$E = \frac{d\Phi}{dA} \quad (6)$$

where $d\Phi$ is the radiant flux (luminous flux) incident on an element dA of the surface containing the point. The unit of irradiance is W/m^2 , and that of illuminance is *lux* ($\text{lx} = \text{lm}/\text{m}^2$). (See Fig. 4.)

Radiance and Luminance

Radiance (L_e) or *luminance* (L_v) is the radiant flux (luminous flux) per unit solid angle emitted from a surface element in a given direction, per unit projected area of the surface element perpendicular to the direction. The unit of radiance is $\text{W sr}^{-1} \text{m}^{-2}$, and that of luminance is cd/m^2 . These quantities are defined by

$$L = \frac{d^2\Phi}{d\Omega \cdot A \cdot \cos\theta} \quad (7)$$

where $d\Phi$ is the radiant flux (luminous flux) emitted (reflected or transmitted) from the surface element and propagating in the solid angle $d\Omega$ containing the given direction. dA is the area of the surface element, and θ is the angle between the normal to the surface element and the direction of the beam. The term $dA \cos \theta$ gives the projected area of the surface element perpendicular to the direction of measurement. (See Fig. 5.)

Radiant Exitance and Luminous Exitance

Radiant exitance (M_e) or *luminous exitance* (M_v) is defined to be the density of radiant flux (luminous flux) leaving a surface at a point. The unit of radiant exitance is W/m^2 and that of luminous exitance is lm/m^2 (but it is not lux). These quantities are defined by

$$E = \frac{d\Phi}{dA} \quad (8)$$

where $d\Phi$ is the radiant flux (luminous flux) leaving the surface element. Luminous exitance is rarely used in the general practice of photometry. (See Fig. 6.)

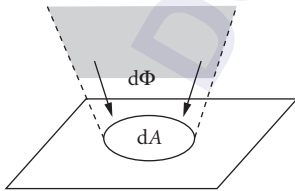


FIGURE 4 Irradiance and illuminance.

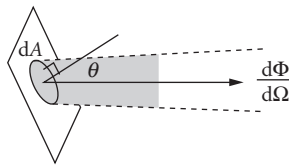


FIGURE 5 Radiance and luminance.

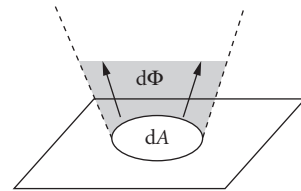


FIGURE 6 Radiant exitance and luminous exitance.

Radiant Exposure and Luminous Exposure

Radiant exposure (H_c) or *luminous exposure* (H_v) is the time integral of irradiance $E_c(t)$ or illuminance $E_v(t)$ over a given duration Δt , as defined by

$$H = \int_{\Delta t} E(t) dt \quad (9)$$

The unit of radiant exposure is J m^{-2} , and that of luminous exposure is *lux · second* ($\text{lx} \cdot \text{s}$).

Radiant Energy and Luminous Energy

Radiant energy (Q_c) or *luminous energy* (Q_v) is the time integral of the radiant flux or luminous flux (Φ) over a given duration Δt , as defined by

$$Q = \int_{\Delta t} \Phi(t) dt \quad (10)$$

The unit of radiant energy is *joule* (J), and that of luminous energy is *lumen · second* ($\text{lm} \cdot \text{s}$).

Total Radiant Flux and Total Luminous Flux

Total radiant flux or *total luminous flux* (Φ_v) is the geometrically total radiant (luminous) flux of a light source. It is defined as

$$\Phi = \int_{\Omega} I d\Omega \quad (11)$$

or

$$\Phi = \int_A E dA \quad (12)$$

where I is the radiant (luminous) intensity distribution of the light source and E is the irradiance (illuminance) distribution over a given closed surface surrounding the light source. If the radiant (luminous) intensity distribution or the irradiance (illuminance) distribution is given in polar coordinates (θ, ϕ) , the total radiant (luminous) flux of the source Φ is given by

$$\Phi = \int_{\phi=0}^{2\pi} \int_{\theta=0}^{\pi} I(\theta, \phi) \sin\theta d\theta d\phi \quad (13)$$

or

$$\Phi = r^2 \int_{\phi=0}^{2\pi} \int_{\theta=0}^{\pi} E(\theta, \phi) \sin\theta d\theta d\phi \quad (14)$$

For example, the total luminous flux of an isotropic point source having luminous intensity of 1 cd would be 4π lm.

Radiance Temperature and Color Temperature

Radiance temperature (unit: kelvin) is the temperature of the Planckian radiator for which the radiance at the specified wavelength has the same spectral concentration as for the thermal radiator considered.

Color temperature (unit: kelvin) is the temperature of a Planckian radiator with radiation of the same chromaticity as that of the light source in question. This term is commonly used to specify the

colors of incandescent lamps whose chromaticity coordinates are practically on the blackbody locus. The next two terms are also important in photometry.

Distribution temperature (unit: kelvin) is the temperature of a blackbody with a spectral power distribution closest to that of the light source in question, and is used for quasi-Planckian sources such as incandescent lamps (refer to Ref. 27 for details).

Correlated color temperature (unit: kelvin) is the temperature of the Planckian radiator whose perceived color most closely resembles that of the light source in question. Correlated color temperature is used for sources with a spectral power distribution significantly different from that of Planckian radiation (e.g., discharge lamps; refer to Ref. 28 for details).

Relationship Between SI Units and English Units

The SI units as described previously should be used in all radiometric and photometric measurements according to international standards and recommendations on SI units. However, some English units are still rather widely used in some countries, including the United States. The use of these non-SI units is discouraged. The definitions of these English units are given in Table 2 for conversion purposes only.

The definition of footlambert is such that the luminance of a perfect diffuser is 1 fL when illuminated at 1 fc. In the SI unit, the luminance of a perfect diffuser would be $1/\pi$ (cd/m²) when illuminated at 1 lx. For convenience of changing from English units to SI units, the conversion factors are listed in Table 3. For example, 1000 lx is the same illuminance as 92.9 fc, and 1000 cd/m² is the same luminance as 291.9 fL. (Conversion factors to and from some more units are given in Ref. 5.)

Troland

This unit is not an SI unit, not used in metrology, and is totally different from all other photometric units mentioned previously. It is introduced here because this unit is commonly used by vision scientists. Troland is defined as the retinal illuminance when a surface of luminance one candela per square meter is viewed through a pupil at the eye (natural or artificial) of area one square millimeter. Thus, the troland value, T , for the luminance, L (cd/m²), of an external field and the pupil size, p (mm²), is given by

$$T = L \cdot P \tag{15}$$

TABLE 2 English Units and Definition

Unit	Quantity	Definition
Footcandle (fc)	Illuminance	Lumen per square foot (lm ft ⁻²)
Footlambert (fL)	Luminance	$1/\pi$ candela per square foot (π^{-1} cd ft ⁻²)

TABLE 3 Conversion between English Units and SI Units

To Obtain the Value in	Multiply the Value in	By
lx from fc	fc	10.764
fc from lx	lx	0.09290
cd/m ² from fL	fL	3.4263
fL from cd/m ²	cd/m ²	0.29186
m (meter) from feet	feet	0.30480
mm (millimeter) from inch	inch	25.400

or, for pupil size p (m²),

$$T = L \cdot 10^6 \times p \quad (16)$$

The troland value is not the real illuminance in lux on the retina, but is a quantity proportional to it. Since the natural pupil size changes with luminance level, luminance changes do not have a proportional visual effect. Thus, troland value rather than luminance is often useful in visual experiments. There is no simple or defined conversion between troland value (for natural pupil) and luminance (cd/m²), without knowing the actual pupil size. (Further details of this unit can be found in Ref. 29.)

37.5 PRINCIPLES IN PHOTOMETRY AND RADIOMETRY

Several important theories in practical photometry and radiometry are introduced in this section.

Inverse Square Law

Illuminance E (lx) at a distance d (m) from a point source having luminous intensity I (cd) is given by

$$E = \frac{I}{d^2} \quad (17)$$

For example, if the luminous intensity of a lamp in a given direction is 1000 cd, the illuminance at 2 m from the lamp in this direction is 250 lx. Note that the inverse square law is valid only when the light source is regarded as a point source. Sufficient distances relative to the size of the source are needed to assume this relationship.

Lambert's Cosine Law

The luminous intensity of a Lambertian surface element is given by

$$I(\theta) = I_n \cos \theta \quad (18)$$

(See Fig. 7.)

Lambertian Surface A surface whose luminance is the same in all directions of the hemisphere above the surface.

Perfect (Reflecting/Transmitting) Diffuser A Lambertian diffuser with a reflectance (transmittance) equal to 1.

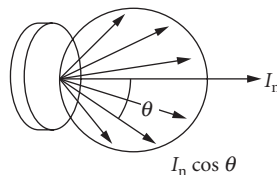


FIGURE 7 Lambert's cosine law.

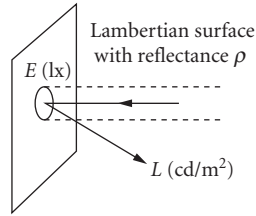


FIGURE 8 Relationship between illuminance and luminance.

Relationship between Illuminance and Luminance

The luminance L (cd/m^2) of a Lambertian surface of reflectance ρ , illuminated by E (lx) is given by

$$L = \frac{\rho \cdot E}{\pi} \quad (19)$$

(See Fig. 8.)

Reflectance (ρ) The ratio of the reflected flux to the incident flux in a given condition. The value of ρ can be between 0 and 1.

In the real world, there is no existing perfect diffuser nor perfectly Lambertian surfaces, and Eq. 19 does not apply. For real object surfaces, the following terms apply.

Luminance Factor (β) Ratio of the luminance of a surface element in a given direction to that of a perfect reflecting or transmitting diffuser, under specified conditions of illumination. The value of β can be larger than 1. For a Lambertian surface, reflectance is equal to the luminance factor. Equation (19) for real object is restated using β as

$$L = \frac{\beta \cdot E}{\pi} \quad (20)$$

Luminance Coefficient (q) Quotient of the luminance of a surface element in a given direction by the illuminance on the surface element, under specified conditions of illumination,

$$q = \frac{L}{E} \quad (21)$$

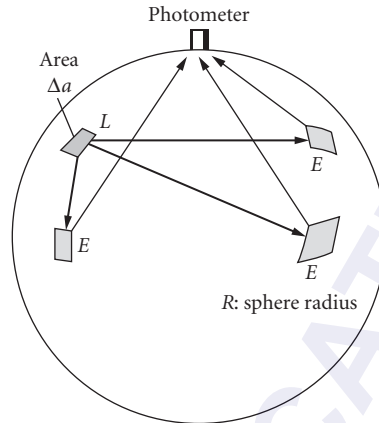
Using q , the relationship between luminance and illuminance is thus given by

$$L = q \cdot E \quad (22)$$

Luminance factor corresponds to radiance factor, and luminance coefficient corresponds to radiance coefficient in radiometry. BRDF (bidirectional reflectance distribution function) is also used for the same concept as radiance coefficient.

Integrating Sphere

An integrating sphere is a device to make a spatial integration of luminous flux (or radiant flux) generated (or introduced) in the sphere and to detect it with a single photodetector. In the case of measurement of light sources, the spatial integration is made over the entire solid angle (4π).


FIGURE 9 Flux transfer in a sphere.

In Fig. 1, assuming that the integrating sphere wall surfaces are perfectly Lambertian, the illuminance E on any part of the sphere wall created by luminance L of an element Δa is given by

$$E = \frac{L\Delta a}{4R^2} \quad (23)$$

where R is the radius of the sphere. This equation holds no matter where the two surface elements are. In other words, the same amount of flux incident anywhere on the sphere wall will create an equal illuminance on the detector port. In the case of actual integrating spheres, the surface is not perfectly Lambertian, but due to interreflections of light in the sphere, the distribution of reflected light will be uniform enough to assume the relationship of Eq. (23). (See Fig. 9.)

The direct light from an actual light source is normally not uniform; thus, it must be shielded from the detector. When a light source with luminous flux Φ is operated in a sphere having reflectance ρ , the flux created by interreflections is given by

$$\Phi(\rho + \rho^2 + \rho^3 + \dots) = \Phi \cdot \frac{\rho}{1 - \rho} \quad (24)$$

Then, the illuminance E_d created by all the interreflections is given by

$$E_d = \frac{\Phi \cdot \rho}{1 - \rho} \cdot \frac{1}{4\pi \cdot R^2} \quad (25)$$

The sphere efficiency (E_d/Φ) is strongly dependent on reflectance ρ due to the term $1 - \rho$ in the denominator. For example, the detector signal at $\rho = 0.98$ is 10 times larger than at $\rho = 0.8$.

Planck's Law

The spectral radiance of a blackbody at a temperature T (K) is given by

$$I_c(\lambda, T) = c_1 n^{-2} \pi^{-1} \lambda^{-5} \left[\exp\left(\frac{c_2}{n\lambda T}\right) - 1 \right]^{-1} \quad (26)$$

where $c_1 = 2\pi hc^2 = 3.7417749 \times 10^{-16} \text{ W} \cdot \text{m}^2$, $c_2 = hc/k = 1.438769 \times 10^{-2} \text{ m} \cdot \text{K}^2$ (1986 CODATA from Ref. 9), h is Planck's constant, c is the speed of light in vacuum, k is the Boltzmann constant, n ($= 1.00028$) is the refractive index of standard air (12),³⁰ and λ is the wavelength.

Wien's Displacement Law

Taking the partial derivative of the Planck's equation with respect to temperature T , and setting the result equal to zero, the solution yields the relationship between the peak wavelength λ_m for Planck's radiation and temperature T (K), as given by

$$\lambda_m T = 2897.8 \mu\text{m} \cdot \text{K} \quad (27)$$

This shows that the peak wavelength of blackbody radiation shifts to shorter wavelengths as the temperature of the blackbody increases.

Stefan-Boltzmann's Law

The (spectrally total) radiant exitance M_e from a blackbody in a vacuum is expressed in relation to the temperature T (K) of the blackbody, in the form

$$M_e(T) = \int_0^\infty M_e(\lambda, T) d\lambda = \sigma T^4 \quad (28)$$

where $M_e(\lambda, T)$ is the spectral radiant exitance of the blackbody, and σ is the Stefan-Boltzmann constant, equal to $5.67051 \times 10^{-8} \text{ W} \cdot \text{m}^{-2} \cdot \text{K}^{-4}$ (1986 CODATA from Ref. 9). Using this value, the unit for M_e is $\text{W} \cdot \text{m}^{-2}$.

37.6 PRACTICE IN PHOTOMETRY AND RADIOMETRY

Photometry and radiometry are practiced in many different areas and applications, dealing with various light sources and detectors, and cannot be covered in this chapter. Various references are available on practical measurements in photometry and radiometry.

Further references in practical radiometry include books on absolute radiometry,³¹ optical detectors,³² spectroradiometry,³³ photoluminescence,³⁴ radiometric calibration,³⁵ etc. There are a number of publications from CIE that are regarded as international recommendations or standards. CIE publications in radiometry include reports on absolute radiometry,³⁶ reflectance,³⁷ spectroradiometry,³⁸ detector spectral response,³⁹ photobiology and photochemistry,⁴⁰ etc. There are also a number of publications from the National Institute of Standards and Technology (NIST) in radiometry, on spectral radiance,⁴¹ spectral irradiance,⁴² spectral reflectance,⁴³ spectral responsivity,⁴⁴ and so on (Ref. 45 provides greater depths of knowledge in radiometry).

For practical photometry, Ref. 4 provides the latest information on standards and practical measurements of photometry in many aspects. A recent publication from NIST⁴⁶ is also available. CIE publications are also available on many subjects in photometry, including characterization of illuminance meters and luminance meters,⁴⁷ luminous flux measurement,⁴⁸ measurements of LEDs,⁴⁹ characteristics of displays,⁵⁰ and many others. A series of measurement guide documents are published from the Illuminating Engineering Society of North America (IESNA) for operation and measurement of particular types of lamps⁵¹⁻⁵³ and luminaires. The American Society for Testing and Materials (ASTM) provides many useful standards and recommendations on optical properties of materials and color measurements.⁵⁴ Colorimetry is a branch of radiometry and is becoming increasingly important among color imaging industry and multimedia applications. The basis of colorimetry is provided by CIE publications^{28,55,56} and many other authoritative references are available.^{29,57}

37.7 REFERENCES

1. F. Grum and R. J. Becherer, *Optical Radiation Measurements, Vol. 1 Radiometry*, Academic Press, San Diego, CA, 1979.
2. R. McCluney, *Introduction to Radiometry and Photometry*, Artech House, Norwood, MA, 1994.
3. W. L. Wolfe, *Introduction to Radiometry*, SPIE—The International Society for Optical Engineering, P.O. Box 10, Bellingham, WA 98227-0010, 1998.
4. Casimer DeCusatis (ed.), *OSA/AIP Handbook of Applied Photometry*, AIP Press, Woodbury, NY, 1997.
5. *IES Lighting Handbook, 8th edition, Reference and Application*, Illuminating Engineering Society of North America, New York, 1993.
6. J. M. Palmer, Radiometry and Photometry FAQ, <http://www.optics.arizona.edu/Palmer/rpfaq/rpfaq.htm>.
7. *International Lighting Vocabulary*, CIE Publication 17.4 (1987).
8. *International Vocabulary of Basic and General Terms in Metrology*, BIPM, IEC, IFCC, ISO, IUPAC, IUPAP, and OIML, 1994.
9. *Quantities and Units*, ISO Standards Handbook, 3rd edition, 1993.
10. CIE Compte Rendu, p. 67 (1924).
11. CIE Compte Rendu, Table II, pp. 25–26 (1931).
12. *The Basis of Physical Photometry*, CIE Publication 18.2 (1983).
13. CIPM, *Comité Consultatif de Photométrie et Radiométrie 10e Session—1982*, BIPM, Pavillon de Breteuil, F-92310 Sèvres, France (1982).
14. *Principles Governing Photometry*, Bureau International Des Poids et Mesures (BIPM) Monograph, BIPM, F-92310 Sèvres, France (1983).
15. CIE Disk D001, Photometric and Colorimetric Tables (1988).
16. CGPM, *Comptes Rendus des Séances de la 16e Conférence Générale des Poids et Mesures*, Paris 1979, BIPM, F-92310 Sèvres, France (1979).
17. CIE Compte Rendu, Vol. 3, Table II, pp. 37–39 (1951).
18. CIPM Procès-Verbaux 44, 4 (1976).
19. *Mesopic Photometry: History, Special Problems and Practical Solutions*, CIE Publication 81 (1989).
20. W. R. Blevin and B. Steiner, *Metrologia* 11:97 (1975).
21. W. R. Blevin, “The Candela and the Watt,” *CIE Proc.* P-79-02 (1979).
22. J. W. T. Walsh, *Photometry*, Constable, London, 1953.
23. *Le Système International d’Unité (SI), The International System of Units (SI)*, 6th edition, Bur. Intl. Poids et Mesures, Sèvres, France, 1991.
24. B. N. Taylor, *Guide for the Use of the International System of Units (SI)*, Natl. Inst. Stand. Technol. Spec. Publ. 811, 1995.
25. B. N. Taylor (ed.), *Interpretation of the SI for the United States and Metric Conversion Policy for Federal Agencies*, Natl. Inst. Stand. Technol. Spec. Publ. 814, 1991.
26. *SI Units and Recommendations for the Use of Their Multiples and of Certain Other Units*, ISO 1000: 1992, International Organization for Standardization, Geneva, Switzerland, 1992.
27. *CIE Collection in Photometry and Radiometry*, Publication No. 114/4, 1994.
28. *Colorimetry* 3rd edition, CIE Publication 15:2004 (2004).
29. G. Wyszecki and W. S. Stiles, *Color Science*, John Wiley and Sons, Inc., New York, 1982.
30. W. R. Blevin, “Corrections in Optical Pyrometry and Photometry for the Refractive Index of Air,” *Metrologia* 8:146 (1972).
31. F. Hengstberger (ed.), *Absolute Radiometry*, Academic Press, San Diego, CA, 1989.
32. W. Budde, *Optical Radiation Measurements, Vol. 4—Physical Detectors of Optical Radiation*, Academic Press, Orlando, FL, 1983.
33. H. Kostkowski, *Reliable Spectroradiometry*, Spectroradiometry Consulting, P.O. Box 2747, La Plata, MD 20646-2747, USA.

34. K. D. Mielenz (ed.), *Optical Radiation Measurements, Vol. 3, Measurement of Photoluminescence*, Academic Press, New York, 1982.
35. C. L. Wyatt, *Radiometrie Calibration: Theory and Models*, Academic Press, New York, 1978.
36. *Electrically Calibrated Thermal Detectors of Optical Radiation (Absolute Radiometers)*, CIE Publication 65 (1985).
37. *Absolute Methods for Reflection Measurements*, CIE Publication 44 (1979).
38. *The Spectroradiometric Measurement of Light Sources*, CIE Publication 63 (1984).
39. *Determination of the Spectral Responsivity of Optical Radiation Detectors*, CIE Publication 64 (1984).
40. *CIE Collection in Photobiology and Photochemistry*, CIE Publication 106 (1993).
41. J. H. Walker, R. D. Saunders, and A. T. Hattenburg, *Spectral Radiance Calibrations*, NBS Special Publication 250-1, 1987.
42. J. H. Walker, R. D. Saunders, J. K. Jackson, and D. A. McSparron, *Spectral Irradiance Calibrations*, NBS Special Publication 250-20, 1987.
43. P. Y. Barnes, E. A. Early, and A. C. Parr, *Spectral Reflectance*, NIST Special Publication 250-48, 1998.
44. T. C. Larason, S. S. Bruce, and A. C. Parr, *Spectroradiometric Detector Measurements*, NIST Special Publication 250-41, 2008.
45. F. Nicodemos (ed.), *Self-Study Manual on Optical Radiation Measurements*, NBS Technical Note 910 Series, Parts 1-12, 1978-1985.
46. Y. Ohno, *Photometric Calibrations*, NIST Special Publication 250-37 (1997).
47. *Methods of Characterizing Illuminance Meters and Luminance Meters*, CIE Publication 69 (1987).
48. *Measurement of Luminous Flux*, CIE Publication 84 (1989).
49. *Measurement of LEDs*, 2nd edition, CIE Publication 127: 2007 (2007).
50. *The Relationship between Digital and Colorimetric Data for Computer-Controlled CRT Displays*, CIE Publication 122, 1996.
51. *IES Approved Method for the Electric and Photometric Measurement of Fluorescent Lamps*, IESNA LM-9.
52. *Electrical and Photometric Measurements of General Service Incandescent Filament Lamps*, IESNA LM-45.
53. *Electrical and Photometric Measurements of Compact and Fluorescent Lamps*, IESNA LM-66.
54. *ASTM Standards on Color and Appearance Measurement*, 5th edition, 1996.
55. ISO 11664-2:2008(E)/CIE S 014-2/E:2006, CIE Colorimetry Part 2: Standard Illuminants for Colorimetry.
56. ISO 11664-1:2008(E)/CIE S 014-1/E:2006, CIE Colorimetry Part 1: Standard Colorimetric Observers.
57. F. Grum and C. J. Bartleson (eds.), *Optical Radiation Measurements, Vol. 2, Color Measurement*, Academic Press, New York, 1980.

This page intentionally left blank.

DO NOT DUPLICATE

Carolyn J. Sher DeCusatis

*Pace University
White Plains, New York*

38.1 INTRODUCTION

Spectroradiometry is the measurement of the spectral content of optical radiation. This has many important applications. The measure of terrestrial, direct, solar spectral irradiance between 295 and 305 nm can be used to calculate atmospheric ozone thickness.¹ More close to home, irradiance measurements are used to characterize light fixtures, and solar UV spectroradiometry methods also apply to the measurement of artificial sources that mimic the sun for applications like phototherapy to treat seasonal affective disorder (SAD) and tanning booths.² Transmission spectra are used to analyze the chemical composition of samples, such as the concentration of chlorophyll in a solution. Spectral reflectance quantifies the color of surfaces, with many practical applications to building, lighting, and design. Spectral responsivity is a necessary part of calibrating photodetectors.

38.2 DEFINITIONS, CALCULATIONS, AND FIGURES OF MERIT

Defining Quantities

There is a relationship between radiometric, photometric, and spectroradiometric quantities. Radiometric and photometric quantities, such as irradiance and luminous flux have been defined in other chapters of this book. Photometric quantities that are similar to radiometric quantities, such as radiant energy versus luminous energy, have the same symbol with a subscript of γ . In general, the spectroradiometric quantity that is defined by the similar radiometric quantity is preceded by the term “spectral,” and designated by the symbol λ , either in parenthesis or with a subscript.

Spectral Irradiance is the quantity most frequently measured in spectroradiometry.¹ Irradiance E is the total radiant flux incident on an element of surface divided by the surface area of that element

$d\Phi/dA$, in watts per meter squared. The average spectral irradiance \bar{E}_λ is the irradiance for a wavelength interval, or

$$\bar{E}_\lambda = \frac{\Delta\Phi}{\Delta\lambda \cdot \Delta A} \quad (1)$$

where $\Delta\Phi$ is the radiant flux within a wavelength interval $\Delta\lambda$ incident on a surface area ΔA . As the area and wavelength are made smaller, \bar{E}_λ becomes the spectral irradiance for that wavelength,

$$E_\lambda(\lambda, P) = \frac{d^2\Phi}{d\lambda \cdot dA} \quad (2)$$

where P is the position and λ is the wavelength. The SI unit for spectral irradiance is the watt per meter cube. However, more intuitive units, such as microwatt per square centimeter of area and per nanometer of wavelength ($\mu\text{W} \cdot \text{cm}^{-2} \cdot \text{nm}^{-1}$) are also used.

The spectral irradiance is dependant on the position, the size of the solid angle subtended, and the orientation of the surface. This leads to the cosine dependence, where the spectral irradiance of a point source is proportional to cosine θ , where θ is the angle between the normal to the surface and the direction of the source, and the inverse square law, relating to the distance between the point source and the detector.

When the radiation from a point source is emitted into a solid angle $d\Omega_s$, the *spectral radiance* L_λ is the flux per unit solid angle and per unit projected area perpendicular to the specified direction, per wavelength.

$$L_\lambda = \frac{d^3\Phi}{dA \cdot \cos\theta \cdot d\Omega \cdot d\lambda} \quad (3)$$

The unit for spectral radiance is watt \cdot centimeter⁻² \cdot steradian⁻¹ \cdot nanometer⁻¹.

Spectral radiance represents the flux density at a point for a particular direction through that point. While it is not usually the desired quantity an experiment is designed to measure, it is required for quantitatively analyzing data, and is also very useful for flux transfer calculations.

If you know the wavelength dependence of the radiant flux, you can also calculate *radiant flux* Φ , and *luminous flux* Φ_γ .

$$\phi = \int_\lambda \Phi(\lambda) d\lambda \quad \text{in watts} \quad (4)$$

$$\phi_\gamma = 683 \int_\lambda \Phi(\lambda) V(\lambda) d\lambda \quad \text{in lumens} \quad (5)$$

where $V(\lambda)$ is the relative photopic luminous efficiency curve (normalized at 555 nm). The absolute luminous efficacy at 555 nm is 683 lumens per watt.

Spectral transmittance $\tau(\lambda)$, which is widely measured during spectrometry, is

$$\tau(\lambda) = \phi^t(\lambda) / \phi^i(\lambda) \quad (6)$$

where t and i refer to transmitted and incident flux. The transmittance has two parts: the regular and diffuse transmittance. The regular transmittance follows Snell's law, while the diffuse transmittance is scattered by the roughness of the surface.

Similarly, *spectral reflectance* $\rho(\lambda)$ is also a flux ratio, but in this case, it is the ratio of the reflected radiant flux to the incident flux.

$$\rho(\lambda) = \phi^r(\lambda) / \phi^i(\lambda) \quad (7)$$

where r refers to reflected flux. The total reflectance is also composed of specular and diffuse components.

The *spectral responsivity* $R(\lambda)$ is the current from a detector divided by the incident flux for a specific wavelength.

$$R(\lambda) = r(\lambda) / \phi(\lambda) \quad (8)$$

where $r(\lambda)$ is the electric signal generated by the photodetector. The units of responsivity are amps per watt. Spectral responsivity is an important part of calibrating photodetectors. The spectroradiometer itself will be calibrated using a standard lamp.

Calculations

Since spectroradiometers provide the spectral content of light, it is natural to use their data to calculate color values. The *tristimulus values* for sources can be calculated,³

$$X = \sum_{\lambda=380}^{780} \bar{E}(\lambda) \bar{x}(\lambda) \Delta\lambda \quad (9)$$

$$Y = \sum_{\lambda=380}^{780} \bar{E}(\lambda) \bar{y}(\lambda) \Delta\lambda \quad (10)$$

$$Z = \sum_{\lambda=380}^{780} \bar{E}(\lambda) \bar{z}(\lambda) \Delta\lambda \quad (11)$$

where $E(\lambda)$ is the irradiance in watt · meters⁻², and \bar{x} , \bar{y} , and \bar{z} are the CIE 1931 spectral tristimulus values.

When analyzing the reflection off an object, the values assume the color as seen under a standard light source.

$$X = k \sum_{\lambda=380}^{780} \bar{E}(\lambda) \Gamma(\lambda) \bar{x}(\lambda) \Delta\lambda \quad (12)$$

$$Y = k \sum_{\lambda=380}^{780} \bar{E}(\lambda) \Gamma(\lambda) \bar{y}(\lambda) \Delta\lambda \quad (13)$$

$$Z = k \sum_{\lambda=380}^{780} \bar{E}(\lambda) \Gamma(\lambda) \bar{z}(\lambda) \Delta\lambda \quad (14)$$

where $\Gamma(\lambda)$ is the *spectral reflectance* or *transmittance data*, and

$$k = \frac{100}{\sum_{\lambda=380}^{780} \bar{E}(\lambda) \bar{y}(\lambda) \Delta\lambda} \quad (15)$$

Since $\bar{y}(\lambda)$ is the photopic curve, k is a constant that can be used to couple the colorimetric (photometric) quantities with the radiometric ones. This can be expressed in equation form as

$$E_v [\text{lm cm}^{-2}] = 683 Y [\text{W cm}^{-2}] \quad (16)$$

because absolute luminous efficacy of the photopic curve at 555 nm is 683 lumens per watt.

The CIE 1931 chromaticity x , y , z coordinates are

$$x = \frac{X}{X + Y + Z} \quad (17)$$

$$y = \frac{Y}{X + Y + Z} \quad (18)$$

$$z = \frac{Z}{X + Y + Z} \quad (19)$$

Similarly, the R, G, B tristimulus values are⁴

$$R = k \sum_{\lambda=380}^{780} P(\lambda) \bar{r}(\lambda) \Delta\lambda \quad (20)$$

$$G = k \sum_{\lambda=380}^{780} P(\lambda) \bar{g}(\lambda) \Delta\lambda \quad (21)$$

$$B = k \sum_{\lambda=380}^{780} P(\lambda) \bar{b}(\lambda) \Delta\lambda \quad (22)$$

where $P(\lambda)$ is the spectral power distribution in watts, and \bar{r} , \bar{g} , and \bar{b} are the color-matching functions of the CIE 1931 Colorimetric Observer.

The UCS 1960 u and v coordinates, and the UCS 1976 u' and v' coordinates are³

$$v = \frac{6y}{12y - 2x + 3} = \frac{2}{3}v' \quad (23)$$

$$u = \frac{4x}{12y - 2x + 3} = u' \quad (24)$$

The CIE LAB/LUV color space calculations (1976) are calculated using the tristimulus values normalized equally to $Y = 100$. X_n , Y_n , and Z_n are the tristimulus values of the reference white. The coordinates can be defined in either the $L^*a^*b^*$ color space, or the $L^*u^*v^*$ color space. When X/X_n , Y/Y_n , and Z/Z_n are all greater than 0.01,

$$L^* = 116 \left(\frac{Y}{Y_n} \right)^{1/3} - 16 \quad (25)$$

$$a^* = 500 \left[\left(\frac{X}{X_n} \right)^{1/3} - \left(\frac{Y}{Y_n} \right)^{1/3} \right] \quad (26)$$

$$b^* = 200 \left[\left(\frac{Y}{Y_n} \right)^{1/3} - \left(\frac{Z}{Z_n} \right)^{1/3} \right] \quad (27)$$

Otherwise

$$L^* = 116 \left[f \left(\frac{Y}{Y_n} \right) - \left(\frac{16}{116} \right) \right] \quad (28)$$

$$a^* = 500 \left[f \left(\frac{X}{X_n} \right)^{1/3} - f \left(\frac{Y}{Y_n} \right)^{1/3} \right] \quad (29)$$

$$b^* = 200 \left[f \left(\frac{Y}{Y_n} \right)^{1/3} - f \left(\frac{Z}{Z_n} \right)^{1/3} \right] \quad (30)$$

where

$$f(Y/Y_n) = \left(\frac{Y}{Y_n} \right)^{1/3} \quad \text{for } Y/Y_n > 0.008856 \quad (31)$$

$$f(Y/Y_n) = 7.787 \left(\frac{Y}{Y_n} \right) + 16/116 \quad \text{for } Y/Y_n \leq 0.008856$$

and $f(X/X_n)$ and $f(Z/Z_n)$ are defined in the same way, and

$$u^* = 13L^*(u' - u'_n) \quad (32)$$

$$v^* = 13L^*(v' - v'_n) \quad (33)$$

These values represent a comparison to a standard illuminant for sources and an ideal white object illuminated by a standard illuminant for objects.³

The *correlated color temperature* is calculated using interpolation from a table of 30 isothermperature lines. Robertson's method, which uses successive approximation, should be accurate to within 0.1 μ rad, with a maximum error from 1600 to 3000 K of less than 0.2 K plus the measurement uncertainty. This technique should only be used for sources with chromaticities farther than 0.01 from the Planckian locus.³

When using a spectroradiometer or spectrophotometer to measure transmission through a sample, the *Beer-Lambert law*, also known as *Beer's law* is used to calculate the *concentration* of a sample solution. The absorbance A is defined as

$$A = -\log_{10}(\tau) \quad (34)$$

where τ is the transmittance. The Beer-Lambert Law states:

$$A = cl\alpha \quad (35)$$

where c is the concentration, l is the path length, and α is the absorption coefficient.

The absorption coefficient is related to the wavelength by

$$\alpha = \frac{4\pi k}{\lambda} \quad (36)$$

where k is the extinction coefficient.

Figures of Merit

Spectroradiometry measurements have large errors compared to other physical measurements. During an intercomparison of solar ultraviolet monitoring between 14 instruments the measured solar irradiances agreed to within 3 percent when the instruments remained outdoors, but the spectral irradiance responsivities changed upon moving the instruments.⁵ In a 2002 intercomparison study by the project Quality Assurance of Ultraviolet Measurements in Europe (QASUME), the spread of absolute irradiance between spectroradiometers was 12 percent (± 6 percent).⁶ There are two major reasons for the large uncertainties:

- The measurement has many dimensions—it is dependent on the magnitude of the flux, its position on the entrance aperture, its direction, its wavelength distribution, and its polarization.
- The instability of measuring instruments and standards, which are very dependent on room conditions such as temperature, and are frequently off by 1 percent or more.

Potential errors in spectroradiometer measurements include measurement noise, detector instability, wavelength instability, nonlinearity, directional and positional effects, spectral scattering, spectral distortion, polarization effects, and size of source effect.¹

These errors can be characterized as³

- Random noise from the detector, electronics, and light source
- Systematic errors from
 - The measurement of the geometry
 - The calibration, including uncertainty from the calibration standard

Noncosine collection of light
 Stray light
 Nonlinearity of the detector and its electronics
 Dark noise subtraction errors

- Periodic errors from drifts due to temperature, humidity, air movement, electronics, beating of AC sources, and changes in stray light

The way of calculating uncertainty was standardized in 1992 by the International Committee for Weights and Measures (CIPM), and the *Guide to the Expression of Uncertainty in Measurement* was published in 1993.⁷ There are two ways to determine the uncertainty of a component: A. statistically or B. “usually based on scientific judgment using all the relevant information available.”⁸

In spectroradiometry, Type B evaluation finds the upper and lower limits of the value (or correction), and then assumes a probability distribution between these values to obtain the standard uncertainty. If you have nothing to base the probability distribution on, you are instructed to assume that it is rectangular (uniform).

For example, for value a where $a = (a_+ - a_-)/2$, where a_+ is the upper limit and a_- is the lower limit, the standard uncertainty is $a/\sqrt{3}$, but assuming a triangular distribution makes it $a/\sqrt{6}$, and a Gaussian distribution makes it $a/3$.

The collected uncertainties for each identified potential error combines as the square root of the sum of the squares, called the *suggested* or *overall uncertainty*. Often this value is multiplied by a constant, under current international practice of value 2, to form the expanded uncertainty.

CIPM requires that all the standard uncertainties and their derivation are included in the uncertainty report.^{7,8}

If your value for spectral irradiance can be expressed of the form

$$E_{\lambda}^{\text{report}} = E_{\lambda}^{\text{observed}} + c_1 + c_2 + c_3 + \dots + c_n \quad (37)$$

where $E_{\lambda}^{\text{report}}$ is the reported spectral irradiance, $E_{\lambda}^{\text{observed}}$ is the measured value of the spectral irradiance, and the c_i 's are the corrections mentioned earlier in this section, then by the CIPM method, the uncertainty U is calculated from the uncertainties u of the parts as follows

$$U = 2\sqrt{u^2(E_{\lambda}^{\text{observed}}) + u^2(c_1) + u^2(c_2) + u^2(c_3) + \dots + u^2(c_n)} \quad (38)$$

Spectroradiometers are calibrated by use of a standard, which has a known value with a reported error. Assuming your system is linear,

$$E_{\lambda}^{\text{obs}} = \frac{S}{S^s} E_{\lambda}^s \quad (39)$$

where S is the measurement and the superscript s refers to the standard.

Then, the uncertainty of the observed spectral irradiance can be calculated using

$$u(E_{\lambda}^{\text{obs}}) = E_{\lambda}^{\text{obs}} \sqrt{\left(\frac{u(S)}{S}\right)^2 + \left(\frac{u(S^s)}{S^s}\right)^2 + \left(\frac{u(E_{\lambda}^s)}{E_{\lambda}^s}\right)^2} \quad (40)$$

The values for $u(S)$ and $u(S^s)$ are typically calculated by a Type B evaluation, while the standard lamp's uncertainty $u(E_{\lambda}^s)$ is calculated from the uncertainty U reported by the standard lamp's supplier.¹

38.3 GENERAL FEATURES OF SPECTRORADIOMETRY SYSTEMS

There are four parts in every spectroradiometer system

- Input or fore-optics
- A monochromator
- A detector
- Electronics and software to analyze data

There is a fifth aspect to every spectroradiometer system, although it is not usually included on these types of lists, because it is not a “part” of the spectroradiometer. However, I believe it is important to consider it while considering other fundamental parts of the system, because it is essential for the accurate measurement of optical radiation:

- Calibration, usually using standard lamps, reflectance standards, or a standard detector

The Input or Fore-Optics

The input optics gathers light from a specified field of view. The layout determines the quantity which is being measured. For example, when measuring spectral irradiance, the light must be diffused, so an integrating sphere or diffusing plate is used, but when measuring spectral radiance, imaging optics control the solid angle and source area, so a focusing mirror is typically used. Transmittance can be measured by placing a light source at the entrance of the system, and measuring the signal twice: with and without the object to be measured. However, an instrument dedicated to measuring transmittance may have a double beam optical design where the light source is split and recombines at the photodetector.

For measurements in the ultraviolet, or light below 190 nm, the radiation is absorbed by the oxygen in the air, so the whole system will be designed to be enclosed and under vacuum.

Telescoping input optics are used when the sources are large distances from the measurement system, turning a spectroradiometer into a telespectroradiometer. Mounting a microscope to the entrance port of the monochromator can make it possible to measure small radiating sources. Fiber-optic probes can be coupled directly to the monochromator, or in combination with any of the previously mentioned input optical devices.³

The Monochromator

The monochromator is the heart of the spectroradiometric system, because separating the radiation into its component wavelengths is the fundamental aspect of the system. While monochromators used to be made with prisms, they are now always made with diffraction gratings. Monochromators come in different sizes; a large monochromator will be more accurate, but a smaller monochromator can be easier to evacuate to measure the ultraviolet, or place in a dry carbon dioxide-free enclosure to measure the infrared.

The monochromator is designed to collimate and focus light. After the entrance slit, light hits a collimating element. Since light is often diverging when it reaches the slit, a concave mirror can form it into a collimated beam directed at the grating.

Generally, the grating will rotate so this beam hits it at different angles. There are also monochromators with curved gratings, but they are limited in wavelength range.

The grating equation, which defines the wavelength of the diffracted flux to the angle of diffraction, is

$$m\lambda = d(\sin\theta \pm \sin\beta) \quad (41)$$

where m is an integer known as the *order of diffraction*, λ is the wavelength, d is the distance between grooves, θ is the angle of incidence, and β is the angle of diffraction. To remove higher orders ($m > 1$), *blocking filters* that absorb short wavelengths while transmitting long wavelengths are used.

The maximum theoretical groove density is $2/\lambda$, although a practical limit is usually 0.85 of the maximum.

The efficiency of a monochromator is directly proportional to large grating area, short focal length, high groove density, long slit length, and high transmittance.¹ However, there is a trade-off between this efficiency and accuracy, which requires large focal lengths. The f -number is the focal distance divided by the entrance slit. Large f -numbers (>3 or 4) are required because the mirrors are spherical, not parabolic, and introduce errors.³

The *dispersion* is the width of the band of wavelengths per unit of slit width, in nm/mm. The *band pass* is the spectral interval that may be isolated. If the dispersion at a groove density k is known, then the band pass can be calculated

$$B = (n_k D(n_k) S) / n \quad (42)$$

where B is the bandpass (in nm), n_k (in grooves/mm) is the known groove density where $D(n_k)$ is its dispersion (in nm/mm), S is the slit width (in mm), and n (in grooves/mm) is the groove density of the grating used.

Bandpass should be small for the best precision. However, there is a trade-off between bandpass and *geometrical etendue* G the light gathering power of an optical system.

$$G = \frac{hnmG_A B}{F 10^6} \quad (43)$$

where h is the slit height (mm), n groove density (groove/mm), m is order of the grating, G_A area of the grating (mm^2), B is the bandpass (nm), and F is the focal length (mm). The ratio h/F implies that the etendue may be increased by making the height of the entrance slit larger. However, this does not work as well in practice; increasing the height of the slit will increase stray light and may also increase the system aberrations, reducing the bandpass.

Geometrical etendue is a limiting function of system *throughput*.⁹

High signal, for which high throughput is necessary, is limited by bandpass. Sometimes the slit size is determined by other factors, like the field of view. But the slit size might be chosen to optimize other factors. In this case, monochromatic sources behave differently from broadband sources, and mixed sources are a combination of the two. For a fluorescent lamp, which is a mixed source, as you decrease the bandpass, the peaks due to the monochromatic lines become much higher in proportion to the broad emission spectra of the phosphors.³

In night vision systems, stray light becomes an important factor which affects system design. To limit stray light, a double monochromator system is used, where the output of the first monochromator is the input of the second one. This can reduce typical stray light levels from 10^{-4} to 10^{-8} .³

Some spectroradiometers are designed with multichannel detectors within the monochromator, or so multichannel detectors can be easily installed. Multichannel detectors can also eliminate the need to scan, reducing moving parts and allowing for longer integration times or quick measurements of unstable or short-lived sources. However, they are not suitable for all types of spectral irradiance and radiance measurements, and spectral transmission, reflection and responsivity require the monochromatic light to exit the monochromator and interact with samples.

The Detector

The desired wavelength range will strongly influence the type of detector used. In Table 1 the approximate wavelength ranges of different spectroradiometry detectors are shown. Other important factors in choosing a detector include the dynamic range, sensitivity, and response time required for the data, as well as environmental factors determining how rugged a detector is needed.

TABLE 1 Approximate Wavelength Ranges of Different Spectroradiometry Detector Types

Detector	Wavelength Range
PMT	200–850 nm
Si photodiode	200–1100 nm
Ge photodiode	1100–1800 nm
InGaAs photodiode	850–1700 nm
PbS photoconductor	1–4.5 μm
PbSe photodiode	1–4.5 μm
InSb photodiode	1.5–5 μm
Pyroelectric	500 nm–50 μm
CCD	200–1100 nm
InGas PDA	800–1700 nm

The two types of detectors most commonly used in spectroradiometry are photomultiplier tubes and semiconductor devices, although thermal detectors have some very limited applications. Detector sensitivity is measured in noise equivalent power (NEP) or equivalent noise input (ENI), which basically mean the same thing, the minimum detectable signal, whose units can be taken in watts.¹ The detectivity is the inverse of the NEP.

Photomultiplier tubes are the most sensitive detectors when used in their wavelength region, with ENIs ranging from 10^{-15} W at 1100 nm to 10^{-17} W from 850 nm to 200 nm and 5×10^{-16} W at 110 nm. They are usually used from 200–850 nm, their range of greatest sensitivity. Silicon photodiodes have NEPs reported of 2×10^{-14} W at 1100 nm, 10^{-15} W at 850 nm, 5×10^{-15} W at 350 nm, and 10^{-14} W between 300 and 200 nm. Silicon photodiodes are used in these wavelength ranges when the signal is sufficiently large, because they are more temperature stable and more rugged.¹

From 1100 to 1800 nm, germanium photodiodes are the most sensitive, with NEPs in the 10^{-13} W range, and lead sulfide photoconductors the most sensitive from 1800–3800 nm, with NEPs ranging from 4×10^{-13} W to 4×10^{-12} W, although InAs may also be used in this range.¹

Photomultipliers, germanium photodiodes, and PbS photoconductors would be cooled, but silicon photodiodes can operate at room temperature (25°C).

Thermal detectors are less sensitive, having NEPs of about 6×10^{-10} W, and have a flat response.¹

When comparing detectors, it is also common to compare the normalized detectivity. There is some confusion in the nomenclature, in that sometimes this is D^* , and sometimes D^* refers to the specific detectivity, which has a different definition. Therefore, I will refer to the normalized detectivity as D_N .

$$D_N = D\sqrt{AB_w} \quad (44)$$

where D is the detectivity, A is the detector area, and B_w is the detector bandwidth.

For example, the D_N of an InAs photodiode is roughly the same as that of a PbS photoconductor at 3000 nm, and the rise time for InAs is one thousandth of PbS. At that specific wavelength, the two detector types have equivalent sensitivity, but InAs has a faster response. However, the D_N drops off more rapidly for InAs at longer and shorter wavelengths than for PbS.¹

Multichannel detectors are sometimes used for spectroradiometry. They have the advantage in spectral irradiance and radiance measurements that they can eliminate the need to scan, reducing moving parts and allowing for longer integration times, as well as making measurements of nonstable sources or short-lived, such as flashbulbs, explosions, and solar measurements during changing weather conditions, possible. Silicon photodiode linear arrays (PDAs) and charge-coupled detectors (CCDs) both singly and with microchannel plates (intensified arrays) are used, as well as a combination of microchannel plates with resistive film called resistive anode (MCP-RA's).¹ These detectors are temperature dependent and therefore require cooling. Noise is also a factor.

The NEP of array detectors is measured with respect to integration time. This is because array detectors operate in the capacitive-discharge mode, which collects charge for a period of time before

TABLE 2 A Comparison of the NEP in Watts of Five Different Multichannel Detector Types with 5-Second Integration Times.

Detector Type	250 nm	550 nm	850 nm	1100 nm
Si PDA at -40°C	1.4×10^{-15}	2.7×10^{-16}	2.2×10^{-16}	1.3×10^{-15}
CCD at -110°C	6.0×10^{-18}	1.6×10^{-18}	1.0×10^{-18}	2.5×10^{-17}
Intensified Si PDA at -40°C	1.3×10^{-18}	8.2×10^{-19}	1.1×10^{-17}	
Intensified CCD at -110°C	1.3×10^{-18}	8.3×10^{-19}	1.1×10^{-17}	
MCP-RA at -30°C	2.0×10^{-18}	1.3×10^{-18}	7.3×10^{-18}	

This data is compiled from Tables 15.4, 15.6, 15.7, 15.9, and 15.11 of Ref. 1.

discharging, which will be discussed in more detail in the next section on electronics and software. However, in Table 2 you can see a comparison of the NEP of the five different types of detectors at integration times of 5 seconds. It can be noted that CCD detectors have NEPs comparable with photomultiplier tubes, even without the advantage of multichannel detection.

While the detectors compared above were all cooled, there exist many commercial spectroradiometers with multichannel detectors which are not cooled for applications where sensitivity is less important, such as projector calibration and film and video post production.

Signal to noise ratio is also a major concern when comparing detector types, but that is also dependent on the integration time. As a general rule, the signal to noise ratio is much lower for PDAs than other multichannel detectors, but at 0.5 second integration times, CCDs are superior.¹

Spectral scattering is a more important concern with multichannel detectors than with single detectors. Interference filters are used to block the short wavelength scattered flux.

The Electronics and Software

In “The Detectors” section we have discussed photomultiplier tubes, semiconductor detectors (photodiodes and photoconductors), PDA-based detectors, and CCD-based detectors.

Photomultiplier tubes use a 250- to 2500-V power supply, and the gain of the detector is adjusted by changing the voltage of the power supply. The signal is the anode current, and the simplest way of measuring it is measuring the voltage drop over a load resistor. However, many problems can be eliminated by using an operational amplifier with a feedback resistor instead, creating a transimpedance current to voltage converter. For even lower signal to noise ratios, integrated current is used as the signal, and the amplifier is chosen to have a short enough time constant that the entire anode current is integrated during the entire sampling time. It is important to make sure that the integrating time is properly matched with the scanning time.

Photodiodes are usually used in the photovoltaic or unbiased mode to minimize noise. In this mode a voltage amplifier with a high-input impedance measures the voltage generated across the photodiode.

Photoconductors decrease in resistance when absorbing radiant flux. The detector is placed in series with a load resistor and a bias voltage. As the resistance of the detector decreases there is an increase in the voltage over the load resistor.

Lock-in amplifiers are often used with semiconductor detectors to decrease $1/f$ noise.

PDA-based detectors operate in the capacitive-discharge mode where the photodiodes store charge as well as sense light. The diodes are reverse biased for 5 V and then electrically isolated while being exposed to light. After an integration time, the charge required to bring the diode back to 5 V is measured.

CCD-based detectors have arrays of pixels that convert radiation to charge. Parallel gates between the pixels can release the charge at will. The charge is then moved along the channels of the arrays by changing the voltage at the gates, until it reaches an output shift register at the end of the channel. The charge is then amplified and measured. Sometimes neighboring pixel's values are intentionally combined, which is known as *binning* or *pixel summation*.

While the readout of a single detector's signal can be as simple as a digital voltmeter, in most cases spectroradiometric data will be recorded by a computer. Often it will be part of an automated system that also controls the drive mechanism. After the information is recorded by the computer, it is analyzed by the computer software. Examples of typical calculations performed on spectroradiometric data are in Sec. 38.2 "Definitions, Calculations, and Figures of Merit" under Calculations.

Calibration

To calibrate a spectroradiometer for irradiance and radiance, standard lamps are usually used, although standard detectors are available. While spectroradiometric standards available for use over the visible are based on the spectral radiance of a blackbody as defined by Planck's radiation law, most commercially available blackbodies are used primarily in the infrared at wavelengths above 1000 nm. Blackbodies suitable for the visible are very expensive because they must operate at temperatures of 2500 K or higher, and are not practical for normal laboratory calibrations.³

Standard lamps introduce a certain degree of error; in the ultraviolet intercomparison in 2002 where the error was 12 percent (± 6 percent), 6 percent of that error was attributed to the calibration lamps.⁶ In earlier intercomparisons, it was more like 1.4 percent from 250 to 2400 nm, but 3 to 4 percent in the infrared.¹

Different lamps are used to calibrate over different wavelength regions. From wavelengths of 250 to 2400 nm, tungsten lamps are used; below 250 nm, a deuterium lamp, argon miniarc, or synchrotron radiation are used.¹

In Table 3 the wavelength range of frequently used spectral irradiance and spectral radiance source standards is listed.

For a discussion of the history of calibration standards, and some additional standards see reference.³

The FEL 1000-W tungsten lamp is the most widely used source standard for spectral irradiance. It is a commercially available 1000-W clear quartz envelope tungsten-halogen, coiled coil filament lamp that is modified to a medium bipost base with 1/4-inch diameter stainless steel posts. It is 5 inches high with a filament about 1 inch long and 1/4 inch in diameter. It is operated at 8-A DC and about 120 V. It is mounted base down with the steel post vertical and the optic axis of the spectroradiometer is horizontal. For more details of the calibration procedure, see reference.¹

It is important to recognize that any individual standard lamp may develop problems. It is recommended to have three standard lamps, so if one changes significantly, you still will have two that agree. It is also standard procedure to transfer the calibration of your newly acquired standards to working standards, lamps that you calibrate using your detector and the NIST traceable standards that have been shipped to you. These lamps can be commercially acquired, but then you must age them for 40 hours at 120-V DV and check them for stability; they must have a drift of less than 0.5 percent at 650 nm in 24 hours. The working standard should be compared to your three purchased standards after 50 to 100 hours of use.

The responsivity of the spectroradiometer can be modeled from a measurement equation, although it is necessary to actually calibrate any system.

TABLE 3 The Wavelength Range of Frequently Used Spectral Irradiance and Spectral Radiance Source Standards

Irradiance	
Tungsten FEL lamp	250–2400 nm
Deuterium lamp	165–350 nm
Argon mini-arc	90–350 nm
Radiance	
Tungsten strip lamp	225–2400 nm
Blackbody below 1000 K	1000–4000 nm
Argon miniarc	90–350 nm

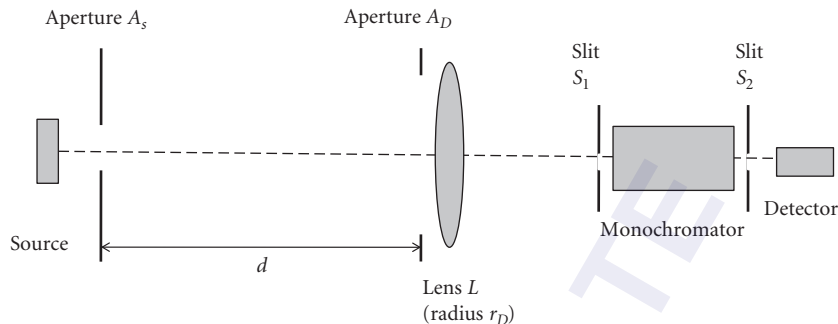


FIGURE 1 A schematic of the experimental setup to measure the spectroradiometer's responsivity.

If the monochromator is set to a wavelength λ_0 , the spectral responsivity $R(\lambda_0, \lambda)$ can be separated into two parts: the slit scattering function $\rho(\lambda_0, \lambda)$ which is an envelope around wavelength λ_0 and the overall responsivity as a function of wavelength $R^f(\lambda)$

$$R(\lambda_0, \lambda) = \rho(\lambda_0, \lambda) R^f(\lambda) \quad (45)$$

When put in an experimental setup shown in Fig. 1, the signal response for a monochromator, ignoring small corrections, is

$$r(\lambda_0) = \frac{A_s A_D}{(r_s^2 + r_D^2 + d^2)} \int_{\lambda} L_{\lambda}(\lambda) \rho(\lambda_0, \lambda) R^f(\lambda) d\lambda \quad (46)$$

where A_s is a circular aperture in front of a standard radiometric source of radius r_s , A_D is an aperture a distance d away from the source in front of a focusing lens L of radius r_D , and $L_{\lambda}(\lambda)$ is of the Lambertian calibration source.¹⁰

If the constants are collected in a term C , and a suitably averaged luminance is used, the signal response can be calculated by

$$r(\lambda_0) = CL_{\lambda}(\lambda_0) \int_{\lambda} \rho(\lambda_0, \lambda) R^f(\lambda) d\lambda \quad (47)$$

Detector standards are also available. In this case, the assumption is that any drift is not very wavelength specific. It is possible to interpolate between calibration laser wavelengths using a blackbody,¹⁰ but that is not usually done, because it is very time consuming and probably an insignificant difference within the larger sources of error.¹ Detectors that presently are suitable for use in transfer standards include silicon, germanium, and InGas photodiodes, and certain types of thermal detectors. The basic approach is to have the two detectors measure the same sources. A typical source is a continuous wave (CW) laser directed at the entrance aperture of an integrating sphere. Types of lasers which can be used are helium-neon, argon, krypton, helium-cadmium, Nd:Yag, and Ti:sapphire.¹¹ When using a detector standard to calibrate a spectroradiometer's spectral irradiance, the area of the entrance aperture and spectral slit width become important and must be taken into account, so it is not a simple measurement.¹

Spectroradiometers are used for reflectance measurements as well as irradiance and radiance measurements. For these purposes, reflectance standards are used for calibration. A number of reflectance standards are available which have been developed for spectroscopy applications, as well as to calibrate colorimeters and spectrophotometers. Specular reflectance standards are calibrated mirrors, and diffuse reflectance standards are made of material similar to the inside of integrating spheres.

Labsphere makes Spectralon into a diffuse white standard and a selection of diffuse gray and color standards, which are calibrated and NIST traceable. Halon, a trade name for polytetrafluoroethylene

(PTFE) powder (which is also used to coat integrating spheres) is also used to make NIST traceable reflectance standards. In a round-robin intercomparison of bidirection diffuse reflectance (BRDR) four types of diffuse reflectors (spectralon, halon, sintered halon, and vacuum deposited aluminum on a ground aluminum surface) were measured at five laboratories. These four types of standards were chosen because of their different scattering mechanisms; Spectralon and pressed PTFE scatter from the bulk, aluminum scatters from the surface, and sintered PTFE scatters from both the bulk and the surface. The purpose of this experiment was to test the laboratories, not the standards, but there was general agreement with the NIST specifications to 2 percent.¹²

38.4 TYPICAL SPECTRORADIOMETRY SYSTEM DESIGNS

Spectral Irradiance and Radiance

The input optics for spectral irradiance measurements require a diffuser in order to eliminate or reduce directional, positional, and polarization effects. This diffuser can be an integrating sphere coated with Halon, which is the best choice for sunlight and large or irregularly shaped sources, a plane reflector diffuser coated with BaSO_4 , or a transmitting diffuser made of teflon.

In Fig. 2 we see a block diagram of a spectroradiometer designed to measure spectral irradiance. The basic steps that have to be taken are the signal must be diffused (cosine corrected), then wavelength selected by the monochromator, after which it is detected by the detector, amplified and analyzed.

In certain special cases, such as when measuring the spectral irradiance of point sources or columnated sources, and if the spectroradiometer responds uniformly over the angular field viewed, no input optics are necessary. However, most of the time input optics are necessary. In Fig. 3 we see a typical setup of the input optics for measuring the irradiance of a large source using a small integrating sphere with a 3-mm-thick PTFE coating (or BaSO_4 in the 310 to 350 nm wavelength region, where PTFE fluoresces weakly), and a circular entrance port of 1 cm² and a rectangular exit port with dimensions

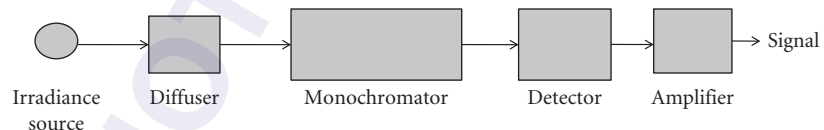


FIGURE 2 Block Diagram to measure spectral irradiance.

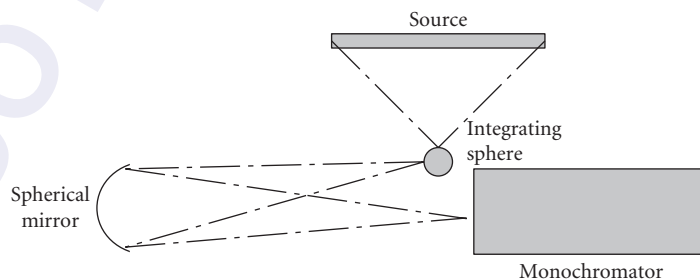


FIGURE 3 Typical setup of the input optics for measuring the irradiance of a large source using a small integrating sphere and a spherical mirror.

3 mm × 12 mm. The radiation which has one reflection does not get into the monochromator because it is reflected away from the spherical mirror. The radiation that does reach the monochromator is independent of the direction and position of the incident flux. However, it is worth noting that the attenuation of a system of this type is large.

In Fig. 4 we see more compact input optics to measure irradiance. The disadvantages are that the sphere is harder to reach, position, and orient, and more stray flux reaches the monochromator.

In Fig. 5 we see typical input optics using a plane diffuser instead of an integrating sphere. In Fig. 6 we see typical input optics using a transmitting diffuser.

In Fig. 7 we see the schematic of a single monochromator and Fig. 8 we see the schematic of a double monochromator. Both of these designs assume an external detector outside the exit slit. In Fig. 9 we see the layout of a single grating monochromator with a built-in multichannel detector.

The input optics for measuring spectral radiance form an image of the source on the entrance port of the monochromator. Possible geometries include a plane and spherical mirror which focus

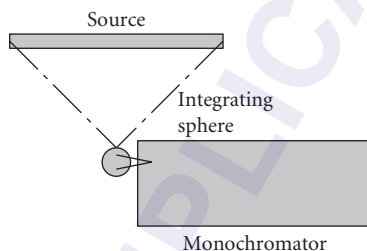


FIGURE 4 Measuring irradiance directly using an integrating sphere.

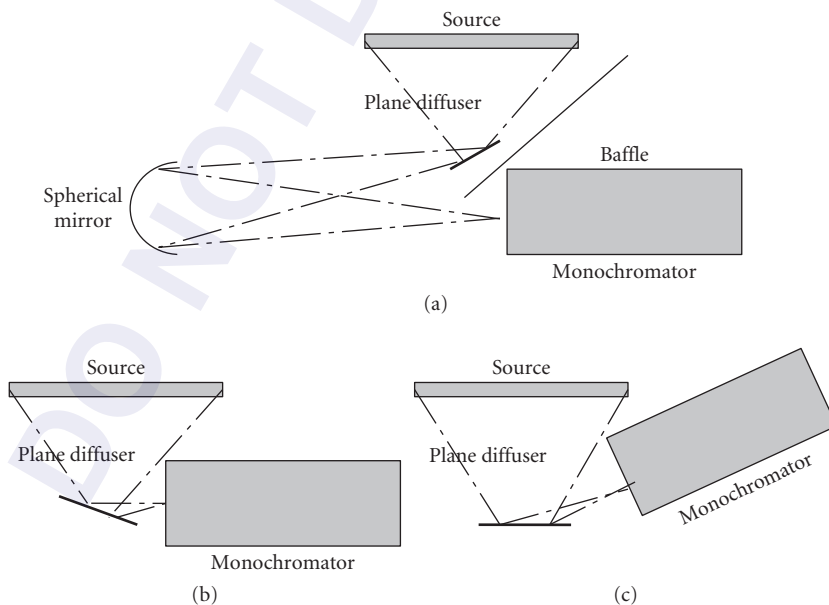


FIGURE 5 Typical setups using a plane diffuser to measure irradiance: (a) with a spherical mirror, (b) and (c) directly reflected off the diffuser into the sphere.

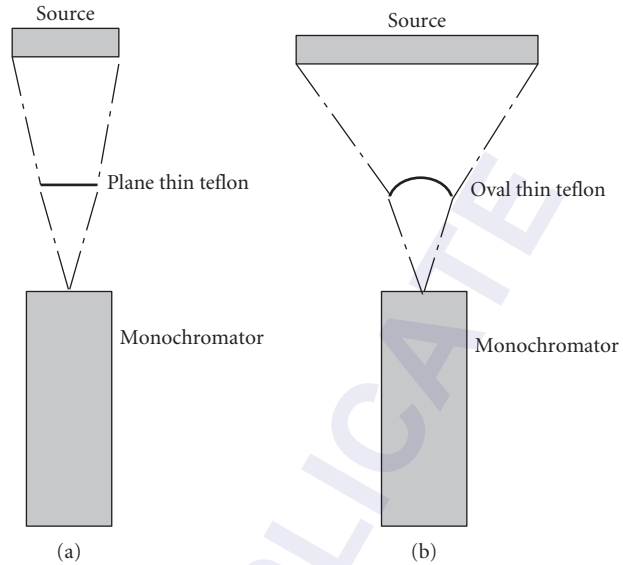


FIGURE 6 Typical input optics for irradiance measurements using a (a) plane thin teflon diffuser and (b) oval shaped thin teflon diffuser.

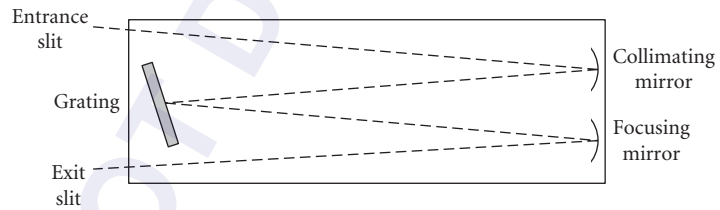


FIGURE 7 Schematic of a single monochromator.

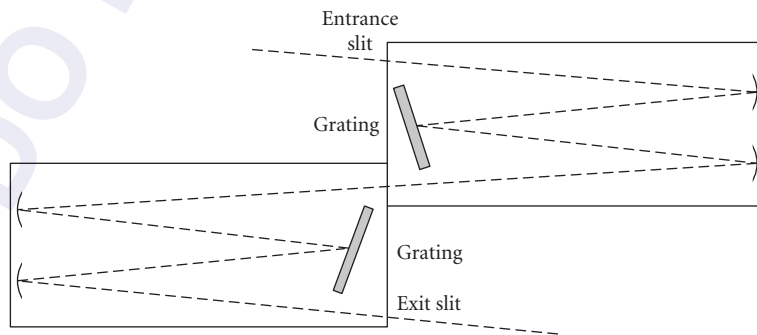


FIGURE 8 Schematic of a double monochromator.

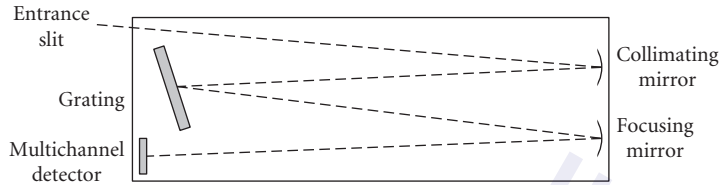


FIGURE 9 Schematic of a single monochromator with a built-in multichannel detector.

the radiation on the entrance slit, and a field of view baffle attachment to limit the acceptance angle of the monochromator.

In Fig. 10 we see the simplest input optics for measuring spectral radiance. However, it is much more likely you are using a system that measures both irradiance and radiance. In that case, your input optics are more likely to look like Fig. 11, where a mirror will reflect the light away from the integrating sphere, and toward the spherical mirror.

The same monochromator and detection electronics would be used for irradiance and radiance measurements.

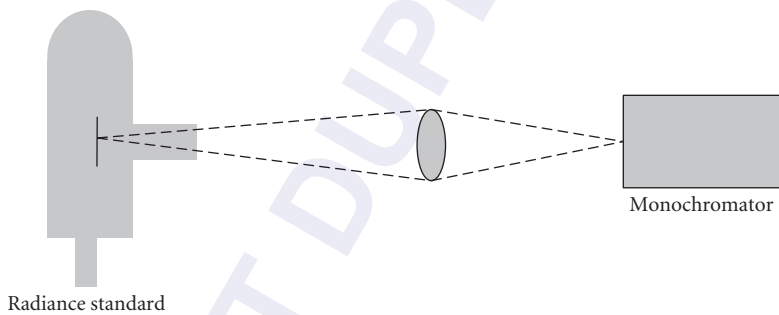


FIGURE 10 Measuring spectral radiance.

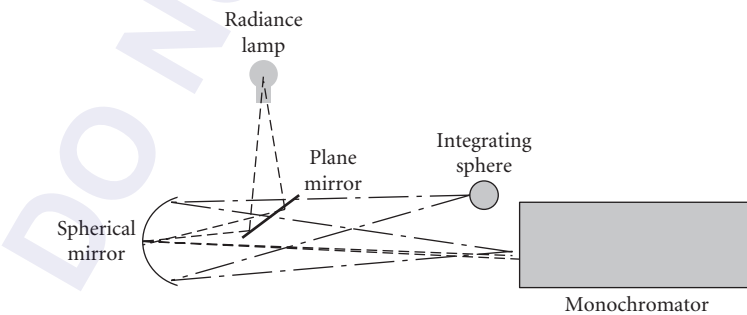


FIGURE 11 System for measuring both irradiance and radiance. The radiance setup adds a plane mirror and measures the radiance lamp. The irradiance setup removes the plane mirror and measures light after it passes the integrating sphere.

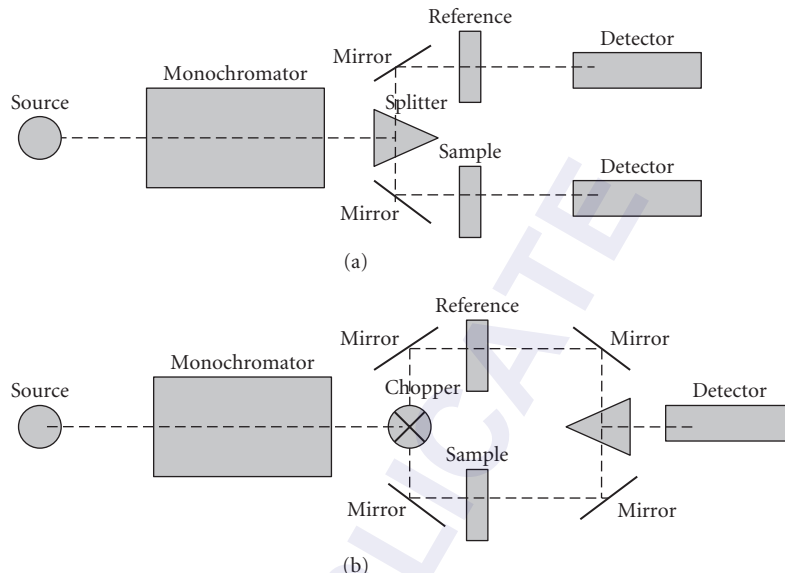


FIGURE 12 Spectrophotometer measuring spectral transmission: (a) dual-beam design and (b) double beam design.

Spectral Transmittance

Spectral transmittance measurements use the source and monochromator to create monochromatic radiation, rather than to measure it. The goal is to pass the radiation through a sample, to measure the sample's properties. This application is of vast importance to biology and chemistry, and there are many commercially available spectroradiometers dedicated to measuring transmittance and reflectance, known as *spectrophotometers*.

Chemists use absorbance spectroscopy to obtain qualitative and quantitative information about samples, using the Beer-Lambert law, also known as Beer's law, as was discussed in Sec. 38.2. Most spectrophotometers use a dual or double beam configuration, as is shown in Fig. 12. This experimental design measures regular transmittance, the signal that passes directly through the sample without being scattered and follows Snell's law. The output is the ratio of the signal in the sample beam to the signal in the reference beam with respect to wavelength. It is necessary to ensure that the only difference between the sample beam and the reference beam is the quantity to be measured, which implies that liquid cells with equal amounts of solute, or gas cells with equal amount of carrier gas, should be placed in the reference beam.¹³

Because of fluorescence, broadband illumination may have different results than monochromatic illumination. This should be considered when measuring transmittance and the setup should approximate the same manner in which the material will be used. Total spectral transmission, which is a combination of regular and diffuse transmission requires the addition of an integrating sphere after the beam transmits through the sample.³

Spectral Reflectance

Manufacturers use spectral reflectance information to provide color information about inks and textiles. There exist several spectral libraries (USGS,¹⁴ Johns Hopkins,¹⁵ JPL¹⁶) which contain almost 2000 spectra of powdered materials for use in spectroscopy, measured by using spectrophotometers and spectrometers in the diffuse reflectance mode.

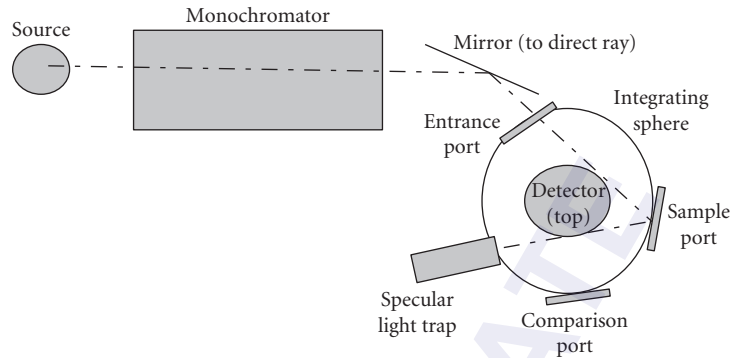


FIGURE 13 Schematic of a spectroradiometer measuring diffuse spectral reflectance.

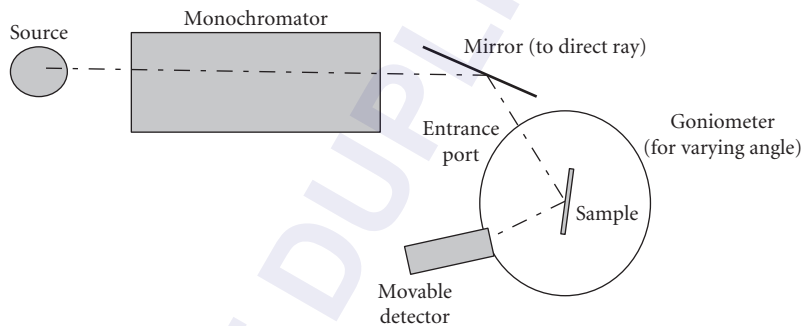


FIGURE 14 A schematic of a spectroradiometer measuring specular spectral reflectance.

Figure 13 is a spectroradiometer configured for measuring diffuse spectral reflectance. There is a double beam design, allowing for the comparison method of using a diffuse reflectance standard. An integrating sphere collects the diffuse radiation, with a removable light trap allows the user to block specular reflectance. If the light trap is not in place, total reflectance, rather than diffuse reflectance, is measured. The sample to be measured is placed in the sample port, and the standard is placed in the comparison port. The detector is perpendicular to the samples and the incident radiation.

Figure 14 is a spectroradiometer configured for measuring specular spectral reflectance. Specular reflectance can be measured at various angles of incidence, including 0° for a 100 percent reading. A calibrated mirror (specular reflectance standard) is not necessary in this design.

Spectral Responsivity

In 38.3 General Features of Spectroradiometer Systems under “Calibration” we discussed how to calculate the spectral responsivity (see Fig. 1) of the spectroradiometer. Once that is known, spectroradiometers can be used to find the spectral responsivity of other detectors.

A spectroradiometer can be configured to measure detector spectral responsivity. The first step uses the spectroradiometer’s standard detector to measure the monochromatic flux or irradiance of the source, generating a function $r^s(\lambda)$ for the standard detector, which has a known responsivity of $R^s(\lambda)$.

These values should not be confused with the responsivity of the total spectroradiometer system. This responsivity of the spectroradiometer system includes both the scatter from the slits and the detector responsivity. In this case we are singling out the detector responsivity. A NIST traceable standard silicon detector is usually used for measurements over the visible spectrum.³

Then the detector is replaced with the detector to be tested, and a signal response of $r^t(\lambda)$ is generated. The responsivity of the test detector $R^t(\lambda)$ will be

$$R^t(\lambda) = r^t(\lambda) / \Phi(\lambda) = r^t(\lambda) R^s(\lambda) / r^s(\lambda) \quad (48)$$

38.5 REFERENCES

1. H. J. Kostkowski (1997) *Reliable Spectroradiometry*, La Plata, MD: Spectroradiometry Consulting.
2. *A Guide to Spectroradiometry: Instruments & Applications for the Ultraviolet*, Reading: Bentham Instruments (1997) (p. 24—tanning booth reference).
3. W. E. Schneider and R. Young (1997) “Spectroradiometry Methods,” *Handbook of Applied Photometry*, Casimer DeCusatis (ed.), Chap. 8, pp. 239–287. New York: AIP Press.
4. J. D. Schanda (1997) “Colorimetry,” *Handbook of Applied Photometry*, C. DeCusatis (ed.), Chap. 10, p. 347. New York: AIP Press.
5. E. Early, A. Thompson, C. Johnson, J. DeLuisi, P. Disterhoft, D. Wardle, and E. Wu, et al., “The 1995 North American Interagency Intercomparison of Ultraviolet Monitoring Spectroradiometers,” *Journal of Research of the National Institute of Standards and Technology* **103**: 15 (1997).
6. A. R. Webb and D. Cotton (2002) “Report of Ispra Intercomparison, May 2002” *Assurance of Ultraviolet Measurements in Europe (QASUME)*. lap.physics.auth.gr/qasume/Files/EvalPdfs/QASUMEREPORT.pdf. Accessed May 19, 2009.
7. ISO (1993) *Guide to the Expression of Uncertainty in Measurement*, International Organization for Standardization, 1, rue de Varembé, Case postale 56, CH-1211 Genève 20, Switzerland.
8. B. N. Taylor and C. E. Kuyatt (1993) *Guidelines for Evaluating and Expressing the Uncertainty of NIST Measurement Results*, NIST Technical note 1297, Washington, D.C.: U.S. Government Printing Office.
9. J. M. Lerner and A. Thevenon (1988) *The Optics of Spectroscopy: A Tutorial*, vol. 2.0, pp. 1–56. Jobin-Yvon: Instruments SA Inc.
10. R. U. Datla and A. C. Parr (2005) “Introduction to Optical Radiometry,” *Optical Radiometry*, A. C. Parr and R. U. Datla (eds.), Chap. 3, pp. 97–154. New York: Elsevier/AIP Press.
11. L. P. Boivin (2005) “Realization of Spectral Responsivity Scales,” *Optical Radiometry*, A. C. Parr and R. U. Datla (eds.), Chap. 3, pp. 97–154. New York: Elsevier/AIP Press.
12. E. A. Early, P. Y. Barnes, B. C. Johnson, J. J. Butler, C. J. Bruegge, S. F. Biggar, P. R. Spyak, and M. M. Pavlov, “Bidirectional Reflectance Round-Robin in Support of the Earth Observing System Program,” *Journal of Atmospheric and Oceanic Technology* **17**: 1077 (2000).
13. J. M. Palmer (2001) “The Measurement of Transmission, Absorption, Emission and Reflection,” *Handbook of Optics*, M. Bass (ed.), I.25.5. New York: McGraw Hill.
14. R. N. Clark, G. A. Swayze, A. J. Gallagher, T. V. V. King, and W. M. Calvin (1993) “The U. S. Geological Survey, Digital Spectral Library: Version 1: 0.2 to 3.0 microns,” *U.S. Geological Survey Open File Report 93–592*, 1340 pages, <http://speclab.cr.usgs.gov>. Accessed May 19, 2009.
15. J. W. Salisbury, L. S. Walter, N. Vergo, and D. M. D’Aria (1992) *Infrared (2.1–25 μm) Spectra of Minerals*, Baltimore, MD: The Johns Hopkins University Press.
16. C. I. Grove, S. J. Hook, and E. D. Paylor (1992) *Laboratory Reflectance Spectra of 160 Minerals, 0.4 to 2.5 Micrometers*. JPL-Publication 92-2. Pilot Land Data System. Pasadena, California: Jet Propulsion Laboratory.

This page intentionally left blank.

DO NOT DUPLICATE

NONIMAGING OPTICS: CONCENTRATION AND ILLUMINATION

William Cassarly

*Optical Research Associates
Pasadena, California*

39.1 INTRODUCTION

Nonimaging optics is primarily concerned with efficient and controlled transfer of radiation. *Nonimaging* refers to the fact that image formation is not a fundamental requirement for efficient radiation transfer; however, image formation is not excluded and is useful in many cases. The two main aspects of radiation transfer that nonimaging optics attempts to solve are maximizing radiation transfer and creating a controlled illuminance distribution. These two problem areas are often described as *concentration* and *illumination*. Many systems require a blending of the concentration and illumination aspects of nonimaging optics, with an important example being the creation of a uniform illuminance distribution with maximum collection efficiency.

Solar collection is an area in which one of the dominant requirements is to maximize the concentration of flux collected by a receiver of a given size (e.g., the irradiance). Any small patch of area on the receiver surface can collect radiation from a hemispherical solid angle. A related problem exists in the detection of low levels of radiation (e.g., Cherenkov counters,¹ infrared detection,² and scintillators³). Nonimaging optics provides techniques to simultaneously maximize the concentration and collection efficiency. The compound parabolic concentrator⁴ is the most well known nonimaging device and has been extensively investigated over the past 40 years.

In addition, there are situations in which collected solid angle must be less than a hemisphere, but the flux density must still be maximized. For a given solid angle and receiver size, the problem is then to maximize the average radiance rather than the irradiance. For example, many fibers only propagate flux that enters the lightpipe within the fiber's numerical aperture (NA), and minimizing the size of the fiber is a key aspect to making practical fiber optic lighting systems.⁵ Projector systems often require maximum flux density at the "film gate," but only the flux reimaged by the projection lens is of importance.^{6,7}

The collected flux must often satisfy certain uniformity criteria. This can be especially important when the luminance of the source of radiation varies or if the flux is collected in a way that introduces nonuniformities. Nonimaging optics provides techniques to control the uniformity. Uniformity control can sometimes be built into the optics used to collect the radiation. In other cases, a separate component is added to the optical system to improve the uniformity.

Uniformity control tends to be created using two main approaches: tailoring and superposition. Tailoring transfers the flux in a controlled manner and uses knowledge of the source's radiance distribution.⁸ Superposition takes subregions of a starting distribution and superimposes them to provide an averaging effect that is often less dependent on the details of the source than tailoring. In many systems, a combination of tailoring and superposition is used. Examples of nonimaging approaches to uniformity control include mixing rods, lens arrays, integrating spheres, faceted reflectors, and tailored surfaces.

The emphasis of this chapter is on the transfer of incoherent radiation. Many of the techniques can be applied to coherent radiation, but interference and diffraction are only mentioned in passing.

39.2 BASIC CALCULATIONS

Some of the terms used in nonimaging optics have yet to be standardized. This section is an effort to collate some of the most common terms.

Photometric and Radiometric Terminology

Understanding nonimaging optics can sometimes be confusing if one is not familiar with photometric and radiometric terminology. This is especially true when reading papers that cross disciplines because terminology usage has not been consistent over time. Reference 9 discusses some of these terminology issues.

This chapter is written using photometric terminology. Photometry is similar to radiometry except that photometry weights the power by the response of the human eye. For readers more familiar with radiometric terminology, a chart showing the relationship between some of the most common photometric terms and radiometric terms is shown in Table 1. Reference 10 describes radiometry and photometry terminology in more detail.

Exitance and *emittance* are similar terms to *irradiance*; however, they denote the case of flux at the surface of a source, whereas irradiance applies to any surface. There are numerous textbooks and handbooks on radiometry and photometry.^{11–16}

Etendue

Etendue describes the integral of the area and the angular extents over which a radiation transfer problem is defined. Etendue is used to determine the trade-off between the required area and angular extents in nonimaging optic designs. Reference 17 provides a brief review of various etendue descriptions with copious references to other work.

TABLE 1 Photometric and Radiometric Terminology

Quantity	Radiometric	Photometric
Power or flux	Watt (W)	Lumen (lm)
Power per unit area	Irradiance, W/m ²	Illuminance, lm/m ² = lux (lx)
Power per unit solid angle	Radiant intensity, W/sr	Luminous intensity, lm/sr = candela (cd)
Power per unit solid angle per unit projected area or Power per unit projected solid angle per unit area	Radiance, W/m ² -sr	Luminance, cd/m ²

One definition of etendue is

$$\text{etendue} = n^2 \iint \cos(\theta) dA d\Omega \quad (1)$$

where n is the index of refraction and θ is the angle between the normal to the differential area dA and the centroid of the differential solid angle $d\Omega$.

In phase space nomenclature (e.g., Ref. 4, Appendix A), etendue is described by

$$\text{etendue} = \iint dx dy n dL n dM = \iint dx dy dp dq \quad (2)$$

where dL and dM are differential changes in the direction cosines (L, M, N), $dx dy$ is the differential area, and $dp dq$ is the differential projected solid angle within the material of index n . The term *phase space* is often used to describe the area and solid angle over which the etendue integral is performed.

Luminance

Luminance divided by the index of refraction squared is the ratio of the differential flux $d\Phi$ to the differential etendue:

$$L/n^2 = d\Phi / d \text{ etendue} = d\Phi / [n^2 \cos(\theta) dA d\Omega] \quad (3)$$

The L/n^2 over a small area and a small angular extent is constant for a blackbody source (e.g., Ref. 18, p. 189). A consequence of constant L/n^2 is that if optical elements are added to modify the apparent area of the small region, then the angular extent of this small area must also change.

If the source of radiation is not a blackbody source, then flux that passes back to the source can either pass through the source or be reflected by the source. In this case, L/n^2 can increase. One example occurs with discharge sources where a spherical mirror is used to reflect flux back through the emitting region. Another example is observed by evaluating the luminance of tungsten in a coiled filament. The luminance is higher at the filament interior surface than the exterior surface because the interior filament surface emits radiation and also reflects radiation that is emitted from other coils.

Lambertian

In many situations the luminance of a radiation source does not vary as a function of angle or position. Such a source is often called a *Lambertian radiator*. In some nonimaging literature, the term *isotropic* is used. The context in which *isotropic* is used should be evaluated because *isotropic* is sometimes used to describe isotropic intensity instead of isotropic luminance.

If constant L/n^2 can be assumed for a given system, then etendue is an important tool for understanding the trade-offs between angular and spatial distributions. Etendue has been used to understand the limits to concentration,⁴ projection display illumination,¹⁹ and backlit display illumination.²⁰ Reference 21 investigates the situation where there is a spectral frequency shift.

In the imaging community, etendue conservation arises in many situations and is often described using the Lagrange invariant. Because imaging systems can often make paraxial assumptions, the approximation $\tan(\theta) = \sin(\theta)$ is often used; however, $\sin(\theta)$ should be used when the collection angles become large.

Clipped Lambertian

An aperture with a clipped Lambertian distribution is one where the source of flux appears to be Lambertian, but only over a finite range of angles. Outside of that range of angles, there is no flux and the range of angles is assumed to be constant across the aperture. The most common example is

when the source is at infinity and the flux at a planar aperture is considered. A clipped Lambertian distribution is also often found at an intermediate surface within an optical system. The terms *limited Lambertian* and *restricted Lambertian* are also used.

Generally, a clipped Lambertian distribution is defined across an aperture and the illuminance across the aperture is constant (e.g., a spatially uniform clipped Lambertian distribution). Another common distribution is a spatially uniform apodized Lambertian. In this case, the angular distribution is spatially uniform but the luminance is not Lambertian. A clipped Lambertian is a special case of an apodized Lambertian. The output of a fiber-optic cable is often assumed to have a spatially uniform apodized Lambertian distribution.

The etendue for clipped Lambertian situations is straightforward to compute. Consider the case of an infinite strip of width $2R$ with a clipped Lambertian distribution defined between $\pm\theta_{\max}$ relative to the surface normal. The etendue per unit length for this 2D clipped Lambertian case is

$$\text{etendue}_{2D} = n(2R)(2\sin\theta_{\max}) \tag{4}$$

2D refers to a trough or extruded geometry and typically assumes that the distribution is infinite in the third dimension. Mirrors can often be placed at the ends of a finite length source so that the source appears to be infinitely long.

For a Lambertian disk with a half cone angle of θ_{\max} , the etendue is

$$\text{etendue}_{3D} = n^2 \text{Area} \pi \sin^2 \theta_{\max} = n^2 \pi R^2 \pi \sin^2 \theta_{\max} \tag{5}$$

In a system where the etendue is preserved, these etendue relationships highlight that increasing either θ_{\max} or R requires a reduction in the other, as depicted in Fig. 1.

Hotell Strings

The etendue relationships described by Eqs. (4) and (5) are primarily defined for the case of a clipped Lambertian source. Such situations arise when the source is located at infinity (e.g., solar collection) or when considering the output of a fiber-optic cable. When the angular distributions vary spatially, Ref. 22 provides a method to compute the etendue that is very simple to use with 2D systems. The method is straightforward for the case of a symmetric system or an off-axis system, and even if there is an absorber positioned between the two apertures. The method is depicted in Fig. 2, where the etendue of the radiation that can be transferred between AB and CD is computed.

For a rotationally symmetric system, the etendue between two apertures (e.g., left side of Fig. 2) has been shown by Ref. 23 to be $(\pi^2/4)(AD-AC)^2$. Reference 24 has provided a generalized treatment of the 3D case.

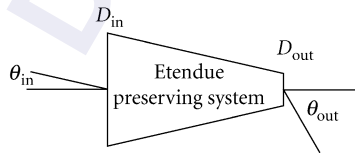


FIGURE 1 Graphical representation of etendue preservation.

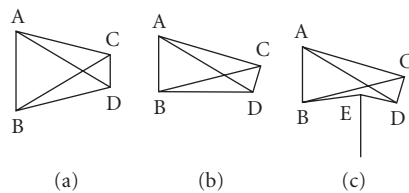


FIGURE 2 Hotell's crossed string relationship for computing etendue. This figure shows a symmetric case (a) and an asymmetric case (b) where the 2D etendue between aperture AB and CD is $(AD-AC + BC-BD)$. In (c), radiation is blocked below point E and the etendue is $(AD-AC + BC-BED)$.

Solid Angle and Projected Solid Angle

Solid angle is the ratio of a portion of the area of a sphere to the square of the sphere radius (see Chap. 3, and Refs. 24 and 25) and is especially important when working with sources that radiate into more than a hemisphere. Solid angle is directly applicable to point sources; however, when the size of the sphere is large enough, the solid angle is still used to characterize extended sources. For solid angle to be applied, the general rule of thumb is that the sphere radius should be greater than 10 times the largest dimension of the source, although factors of 30 or more are required in precision photometry (e.g., Ref. 26, pp. 4.29–4.31).

The etendue is sometimes called the *area-solid angle product*; however, this term can often cause confusion because etendue is actually the projected-area-solid-angle product [$\cos \theta dA d\Omega$] or area-projected-solid-angle product [$dA \cos \theta d\Omega$]. Reference 27 discusses some distinctions between projected solid angle (PSA) and solid angle. An important implication of the cosine factor is that the solid angle for a cone with half angle θ is

$$\text{Solid angle}_{\text{cone}} = 2\pi[1 - \cos\theta] = 4\pi \sin^2(\theta/2) \quad (6)$$

but the projected solid angle for the cone is

$$\text{Projected solid angle}_{\text{cone}} = \pi \sin^2 \theta \quad (7)$$

For a hemisphere, the solid angle is 2π and the PSA is π . This factor of 2 difference is often a source of confusion. PSA and solid angle are pictured in Fig. 3.

If the luminance at the receiver is constant, then PSA times luminance provides illuminance. Nonuniform luminance can be handled using weighted averages.²⁸

In the 2D case, the projected solid angle analog is simply projected angle, and is $2 \sin \theta$ for angles between $\pm\theta$ or $|\sin \theta_1 - \sin \theta_2|$ for angles between θ_1 and θ_2 .

Concentration

Concentrators can be characterized by the ratio of the output area to the input area.²⁹ For an *ideal* system, the etendue at the input aperture and the output aperture are the same, which leads to the ideal concentration relationships

$$\text{Concentration}_{2D} = n_{\text{out}} \sin \theta_{\text{out}} / (n_{\text{in}} \sin \theta_{\text{in}}) \quad (8)$$

$$\text{Concentration}_{3D} = n_{\text{out}}^2 \sin^2 \theta_{\text{out}} / (n_{\text{in}}^2 \sin^2 \theta_{\text{in}}) \quad (9)$$

where the input and output distributions are clipped Lambertians, θ_{in} is the maximum input angle, and θ_{out} is the maximum output angle. Maximum concentration occurs when $\sin \theta_{\text{out}} = 1$, so that the

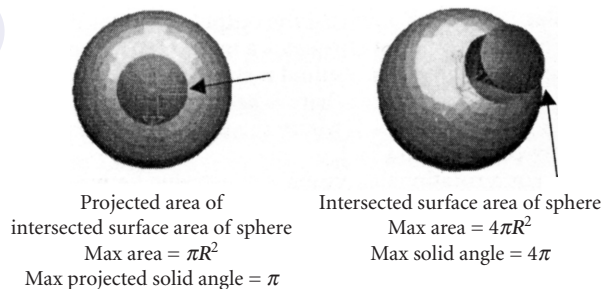


FIGURE 3 PSA (left) versus solid angle (right).

maximum concentration ratios for 2D and 3D symmetric systems in air are $1/\sin \theta_{in}$ and $1/\sin^2 \theta_{in}$, respectively. Concentrating the flux within a dielectric medium further increases the concentration ratio by a factor of n_{out}^2 in the 3D case and n_{out} in the 2D case.

In general, maximum concentration occurs when the phase space of the receiver is completely filled (i.e., all portions of the receiver are illuminated from all angles). Maximum concentration can be obtained if the etendue of the input is greater than the etendue of the output. In this case, however, the efficiency of the system is less than 1. One example of this is provided by a two-stage system where the second stage is an ideal maximum concentrator and the etendue of the first stage is greater than the etendue of the second stage.^{30,31} Such a system is not conventionally considered ideal because not all of the flux from the first stage is coupled into the entrance aperture of the second stage.

Dilution

Dilution is a term that is used to describe the situation in which the phase space of the receiver is unfilled. Consider the situation where a Lambertian disk of radius 1 is coupled directly to a receiver of radius 2. In this case, the area of the receiver is underfilled and represents *spatial dilution*. An example of *angular dilution* occurs when flux is transferred to a receiver using multiple discrete reflector elements. If there are gaps between the reflector elements, then the angular distribution of the flux incident on the receiver possesses gaps. The phrase *etendue loss* is also used to describe situations in which dilution is present.

39.3 SOFTWARE MODELING OF NONIMAGING SYSTEMS

The computer simulation of nonimaging systems is quite different from the computer simulation of imaging systems. Modeling of nonimaging systems typically requires three major items that imaging designers either do not need or tend to use infrequently. These are nonsequential ray tracing, spline surfaces, and modeling extended sources. Extended sources are often modeled using Monte Carlo techniques. Other differences include the need to model more complex surface properties including scattering surfaces; new methods to present simulation results including luminance, illuminance, and intensity distributions; and improved visualization because of the more complex geometries typically involved in nonimaging systems.

Nonsequential Ray Tracing

Nonsequential ray tracing means that the order of the surfaces with which the rays interact is not predetermined. In fact, rays can hit a single surface multiple times, which is especially common when skew rays are investigated. Nonsequential surfaces are especially important in the software modeling of lightpipes^{32,33} and prisms.

Spline Surfaces

Spline surfaces have been successfully applied to the design of nonimaging optical systems. Design work in this area includes genetic algorithm-based approaches,^{34,35} neural networks,³⁶ and variational principles.³⁷ Automotive headlamp designs also use spline surfaces routinely. Spline surfaces are also extremely important in computer-aided drafting (CAD) data exchange formats (e.g., IGES, STEP, and SAT) where Non-Uniform Rational B Splines (NURBS) have found tremendous utility.

Monte Carlo Techniques

Monte Carlo simulations are used to determine the intensity and/or illuminance distribution for optical systems. Typically, Monte Carlo techniques are used to model the source and/or surface properties of nonspecular surfaces. Reference 38 discusses some of the issues regarding Monte Carlo simulations. Some nonimaging systems that have been evaluated using Monte Carlo ray traces include an end-to-end simulation of a liquid crystal display (LCD) projector system,³⁹ light-emitting diodes,⁴⁰ integrating spheres,^{41,42} scintillating fibers,^{43,44} stray light,^{45,46} and automotive headlamps.⁴⁷

Source Modeling

Monte Carlo modeling of sources typically falls into three categories: point sources, geometric building blocks, and multiple camera images. Point sources are the simplest approach and are often used in the early stages of a design because of simplicity. Measured data, sometimes called *apodization data*, can be applied to a point source so that the intensity distribution of the source matches reality. The projected area of the source as seen from a given view angle can often be used to estimate the apodization data^{8,48} when measurements are unavailable. Since the luminance of a point source is infinite, point source models are of limited utility when etendue limitations must be considered. Geometric models create the source through the superposition of simple emitters (e.g., disks, spheres, cylinders, cubes, ellipsoids, toroids). The simple emitters are typically either Lambertian surface emitters or isotropic volume emitters, and the emitters can have surface properties to model effects that occur when flux propagates back through the emitter. The simple emitters are then combined with models of the nonemissive geometry (e.g., bulb walls, electrodes, scattering surfaces) to create an accurate source model. Apodization of the spatial and/or angular distributions of sources can help to improve the accuracy of the source model. In the extreme, independent apodization files can be used for a large number of view angles. This often results in the use of multiple source images,^{49,50} which has seen renewed interest now that charge-coupled device (CCD) cameras have matured.^{51–55}

The most typical Monte Carlo ray trace is one in which the rays traced from the source are independent of the optical system under investigation. When information about the portions of the source's spatial and angular luminance distribution that are contributing to a region of the intensity/illuminance distribution is available, importance sampling can be used to improve the efficiency of the ray trace. An example of such an approach occurs when an $f/1$ condenser is used to collect the flux from a source. Use of importance sampling means that rays outside of the $f/1$ range of useful angles are not traced.

Backward Trace

When a simulation traces from the source to the receiver, rays will often be traced that land outside of the region of the receiver that is under investigation. A backward ray trace can eliminate those rays by only tracing from the observation point of interest. The efficiency of the backward ray trace becomes dependent upon knowing which ray directions will hit the receiver. Depending upon the details of the system and prior knowledge of the system, backward tracing can often provide far more efficient use of traced rays when only a few observation points are to be investigated. Use of the backward ray trace and examples for a number of cases have been presented.²⁸ Backward ray tracing has also been called the *aperture flash mode*.⁵⁶

Field Patch Trace

Another trace approach is the field patch, where rays are traced backward from a point within the optical system and those rays are used to determine which rays to trace forward. Reference 56 describes the prediction of a distribution by summing results from multiple field patches. This

type of approach is especially useful if the effect of small changes in the system must be quantified because shifting the distributions and resuming can approximate small changes.

Software Products

There are many optical software vendors that have features that are primarily aimed toward the illumination market. A detailed comparison of the available products is difficult because of the speed at which software technology progresses. A survey of optical software products has been performed by Illuminating Engineering Society of North America (IESNA) and published in *Lighting Design and Application*.^{57,58}

39.4 BASIC BUILDING BLOCKS

This section highlights a number of the building blocks for designing systems to collect flux when the angles involved are large. Many of these building blocks are also important for imaging optics, which reflects the fact that image formation is not a fundamental requirement for nonimaging optics but that image formation is often useful.

Spherical Lenses

Spherical lenses are used in nonimaging systems although the aberrations can limit their use for systems with high collection angles. The aplanatic points of a spherical surface⁵⁹ are sometimes used to convert wide angles to lower angles and remain free of all orders of spherical aberration and low orders of coma and astigmatism. There are three cases⁶⁰ in which a spherical surface can be aplanatic:

1. The image is formed at the surface (e.g., a field lens).
2. Both object and image are located at the center of curvature.
3. The object and image are located on radii that are rn/n' and rn'/n away from the center of curvature of the spherical surface.

Case 3 is used to construct hyperhemispheric lenses that are often used in immersed high-power microscope objectives. Such lenses can suffer from curvature of field and chromatism (see Ref. 60, pp. 258–262, and Ref. 61). Some example aplanats are shown in Fig. 4.

Hyperhemispheric lenses can also be used in LED packages⁶² and have been used in photographic-type objectives and immersed IR detectors. Aplanatic surfaces for use as concentrators have been investigated⁶³ and can be nearly ideal if the exit surface is nonplanar (Ref. 4, pp. 37–38).

Aspheric Lenses

Spherical lenses have been the workhorse of the imaging industry because of the simplicity and accuracy with which they can be manufactured. Design requirements, especially the speed of the system, often necessitate the use of aspheric surfaces in nonimaging optics. Fortunately, the accuracy with which the surface figure must be maintained is often less severe than that of an imaging system. Aspheric lenses are often used as condensers in projection systems including LCD projectors, automotive headlamps, and overhead projectors.

The classic piano aspheric examples are the conic lenses where the eccentricity is equal to $1/n$ with the convex surface facing the collimated space and eccentricity = n with the piano surface facing the collimated space (see, for example, Ref. 64, pp. 112–113, and Ref. 65, pp. 100–103).

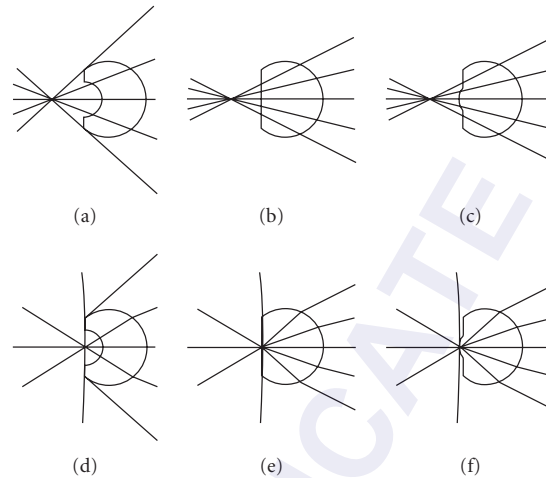


FIGURE 4 Several aplanats used to concentrate a virtual source. (a–c) The rays superimposed on top of the aplanats. (d–f) the ray trace. (d) A meniscus lens with hyperhemispheric outer surface and a spherical surface centered about the high concentration point. (e) A hyperhemisphere. (f) A hyperhemisphere with curved output surface.

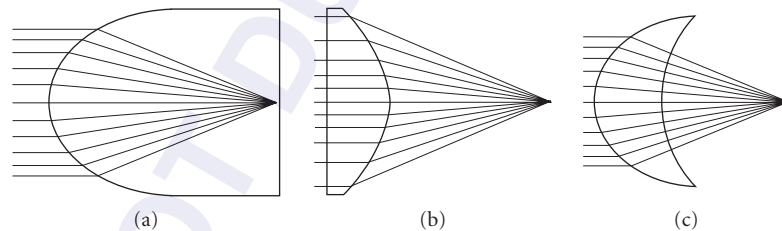


FIGURE 5 Conic collimators. (a) $e = 1/n$ and $f = Rn/(n - 1)$. (b) $e = n$ and $f = R/(n - 1)$. (c) $e = 1/n$ with other surface being spherical. $e =$ eccentricity, $n =$ index of refraction, $k = -e^2$.

Examples are shown in Fig. 5. The paraxial focal lengths are $Rn/(n - 1)$ and $R/(n - 1)$, respectively, where R is the radius of curvature. These examples provide correction for spherical aberration but still suffer from coma.

Refractive aspheric surfaces are often used in the classic Schmidt system.^{60,66} Reference 67 has described a procedure for determining the aspheric profile for systems with two aspheric surfaces that has been adapted to concentrators by Minano.⁶⁸

Fresnel Lenses

The creation of a lens using a symmetric arrangement of curved prisms is typically attributed to A. J. Fresnel and commonly called a *Fresnel lens*. General descriptions of Fresnel lenses are available.^{69,70}

Fresnel lenses provide a means to collect wide angular distributions with a device that can be easily molded. Standard aspheric lenses can be molded, but the center thickness compared to the edge thickness adds bulk and fabrication difficulties.

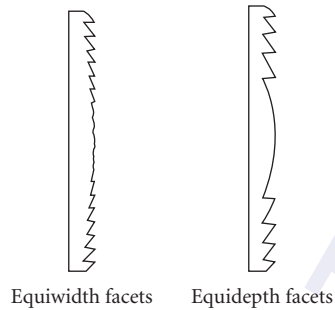


FIGURE 6 Two main types of Fresnel lenses: equiwidth and equidepth.

In imaging applications, Fresnel lenses are used, but the system designer must be careful to consider the impact of the facets on the image quality. The facet structure is also important in nonimaging applications where the use of facets introduces losses either by blocking flux or by underfilling the aperture, which typically results in an etendue loss.

Fresnel lenses have two main variations: constant-depth facet structures and constant-width facet structures (see Fig. 6). As the width of the facets becomes smaller, the effects of diffraction should be considered.⁷¹ Although Fresnel lenses are typically created on planar substrates, they can be created on curved surfaces.^{72,73}

The design of a Fresnel lens requires consideration of the riser angle (sometimes called the *draft angle*) between facets.^{74,75} A typical facet is depicted in Fig. 7. The riser angle should be designed to minimize its impact on the rays that are controlled by the facets that the riser connects.

Total Internal Reflection Fresnel Lenses As the Fresnel lens collection angle increases, losses at the outer facets increase.⁷⁶ One solution is the use of total internal reflection (TIR) facets. In this case, the flux enters the substrate material, hits the TIR facet, and then hits the output surface. The entering, exiting, and TIR facets can also have power in more sophisticated designs. The basic TIR Fresnel lens idea has been used in beacon lights since at least the 1960s.⁷⁷ TIR Fresnel lenses have received more recent design attention for applications including fiber illumination, LEDs, condensers, and concentrators.^{76,78–82} Vanderwerf⁸³ investigates achromatic issues with catadioptric Fresnel lenses.

In some applications, the TIR Fresnel lens degenerates into one refractive facet and two or more TIR facets. Combined TIR and refractive designs have been proposed for small sources with significant attention to LEDs.^{84–86}

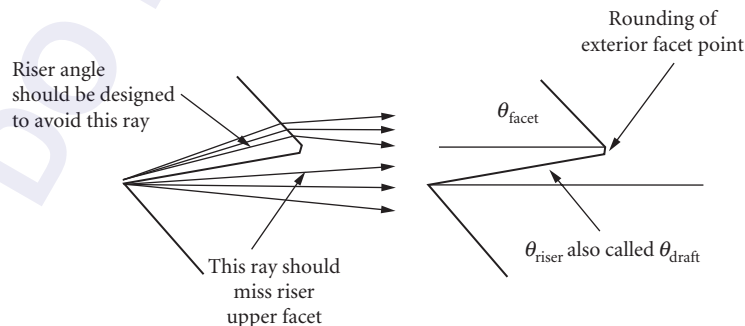


FIGURE 7 Facet terminology and design issues.

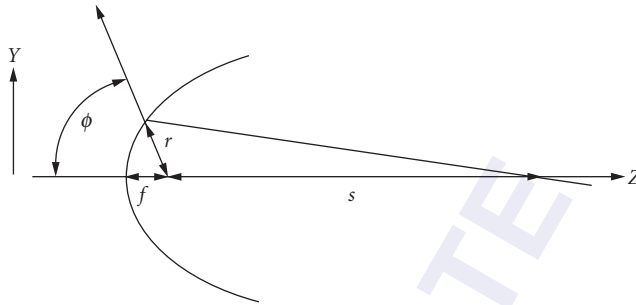


FIGURE 8 Conic reflector showing r , f , s , and ϕ .

Conic Reflectors

Reflectors made from conic sections (parabolas, ellipses, and hyperbolas) are commonly used to transfer radiation from a source to a receiver. They provide spherical-aberration-free transfer of radiation from one focal point to the other. The problem is that the distance from a focal point to a point on the surface of the reflector is not constant except for the case of a sphere. This introduces a nonuniform magnification (coma) that can result in a loss of concentration.

Many textbooks describe conic surfaces using equations that have an origin at the vertex of the surface or at center of symmetry. Reference 17 (page 1.34) shows a number of different forms for a conic surface. In nonimaging systems, it is often convenient to define conic reflectors using an origin shifted away from the vertex. Reference 87 describes these surfaces in Cartesian coordinates as *apo-vertex surfaces*. They can also be described in polar coordinates where the focal point is typically the origin. The polar definition of a conic surface⁸ is

$$r = \frac{f(1+e)}{1+e \cos \phi}$$

where r is the distance from the foci to the surface, f is the distance from one focus to its nearest vertex, ϕ is the angle from the foci to the vertex, e is $s/(s+2f)$, and s is the distance between the two foci. In Cartesian coordinates, $z = r \sin \phi$ and $y = r \cos \phi$. An example is shown in Fig. 8.

Macrofocal Reflectors

One can consider the standard conic reflector to map the center of a sphere to a single point. What can sometimes be more valuable for nonimaging optics is to map the edge of the sphere to a single point. Reflectors to provide this mapping have been called macrofocal reflectors.⁸⁸ They have also been called extinction reflectors⁸ because of the sharp edge in the illuminance distribution that they can produce. Reference 89 describes an illumination system that provides a sharp cutoff using tilted and offset parabolic curves.

Involute

A reflector that is used in many nonimaging systems is an *involute*. An involute reflector sends tangential rays from the source back onto themselves. One way to produce an involute for a circle is to unwind a string that has been wrapped about the circle. The locus of points formed by the string equals the involute. In Cartesian coordinates, the equation for the involute of a circle is⁹⁰

$$x = r(\sin \theta - \theta \cos \theta) \quad (10)$$

$$y = -r(\cos \theta + \theta \sin \theta) \quad (11)$$

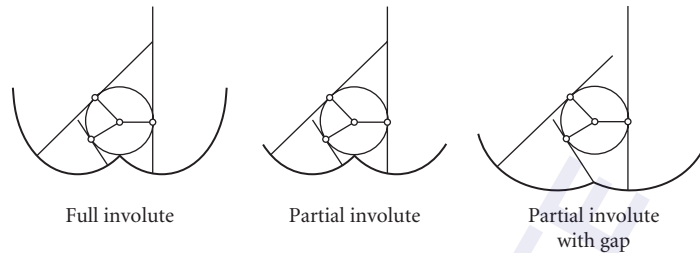


FIGURE 9 Full involute, partial involute, and involute with gap. Rays tangent to the source reflect back toward the same location on the source.

where r is the radius of the circle, $\theta = 0$ at the cusp point, and the origin is at the center of the circle. An example is shown in Fig. 9, where rays are drawn to highlight the fact that rays tangent to the circle reflect off of the involute and retrace their path.

The term *full involute* has been used to describe an involute that is continued until the output aperture is tangent to the surface of the source. A *partial involute* is an involute that is less than a full involute. An involute with gap is also shown in Fig. 9, and is sometimes called a *modified involute*.⁹¹ Involute segments can also be applied to noncircular convex shapes and to disjoint shapes.⁹²

39.5 CONCENTRATION

The transfer of radiation in an efficient manner is limited by the fact that the luminance of the radiation from a source in a lossless system cannot increase as the radiation propagates through the system. A consequence is that the etendue cannot be reduced and maximizing the luminance means avoiding dilution. A less understood constraint that can impact the transfer of radiation is that rotationally symmetric optical systems conserve the skewness of the radiation about the optical axis (see for example Ref. 4, Chap. 2.8, and more recently Ref. 93). This section describes some of the basic elements of systems used to transfer radiation with an emphasis on those that preserve etendue.

Many papers on nonimaging optics can be found in the Society of Photooptical Instrumentation Engineers (SPIE) proceedings. SPIE conferences specifically focused on nonimaging optics include Refs. 94–99. Reference 100 also provides a large number of selected papers. Textbooks on nonimaging optics include Refs. 4 and 101. Topical articles include Refs. 102 and 103.

Discussions of biological nonimaging optic systems have also been written.^{104–108}

Tapered Lightpipes

Lightpipes can be used to transport flux from one location to another. This can be an efficient method to transport flux, especially if total internal reflection (TIR) at the lightpipe surface is utilized to provide lossless reflections. Lightpipes can provide annular averaging of the flux. In addition, if the lightpipe changes in size from input to output over a given length (e.g., the lightpipe is tapered), then the change in area produces a corresponding change in angles. If the taper occurs over a long enough distance, then in most cases the etendue will be preserved (see, for example, Ref. 109). If the taper occurs over too short a length, the etendue will not be preserved, which can result in an increase in angles or possibly rays lost, such as when they are reflected back toward the input.

An interesting geometric approach to determining if meridional rays in a conical lightpipe can propagate from input to output end was shown by Williamson,¹¹⁰ and is informally called the Williamson construction or a tunnel diagram. An example is shown in Fig. 10, where the duplicate copies of the primary lightpipe are all centered about a single point. Where the rays cross the copies

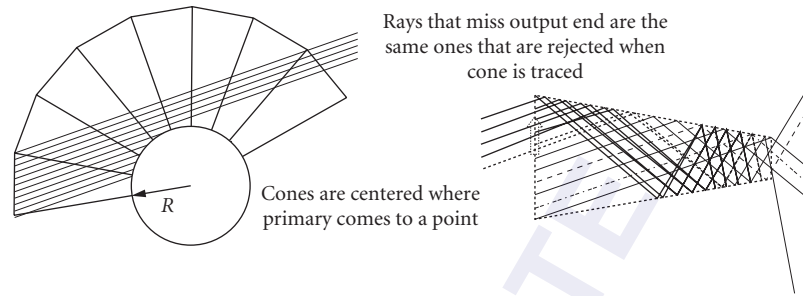


FIGURE 10 Williamson construction. Multiple copies of the primary cone are drawn. Where rays cross, these copies are the same as where they hit the actual cone. Rays that do not cross the output end are reflected back to the input end.

shows the position and at what angle the rays will hit the primary lightpipe. As seen in Fig. 10, the rays that do not cross the output end of any of the copies of the cone are also the rays that are reflected back out the input end of the cone. Other analyses of tapered geometries have been performed by Welford⁴ (pp. 74–76), and Meyer¹¹¹ provides numerous references.

Witte¹¹² has performed skew ray analysis using a variation of the Williamson construction. Witte also shows how a reference sphere centered at the point where the cone tapers to a point can be used to define a pupil when cones are used in conjunction with imaging systems analysis. Burton¹¹³ simplified the expressions provided by Witte.

Vector flux investigations have shown that a cone is an ideal concentrator for a sphere (Ref. 114, p. 539). An important implication of this result, which is similar to the information presented by Witte,¹¹² is that a cone concentrator can be analyzed by tracing from the cone input aperture to the surface of a sphere. If the rays hit the sphere surface, then they will also hit the sphere after propagating through the lightpipe. This is true for both meridional and skew rays.

CPC

A compound parabolic collector (CPC) can be used to concentrate the radiation from a clipped Lambertian source in a nearly etendue-preserving manner. Welford⁴ provides a thorough description of CPCs. An example CPC is shown in Fig. 11, where the upper and lower reflective surfaces are tilted parabolic surfaces. The optical axis of the parabolic curves is not the same for the upper and lower curves, hence the use of the term *compound*.

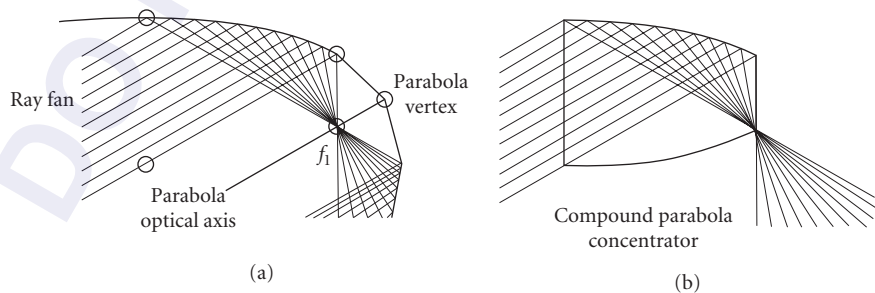


FIGURE 11 Compound parabolic concentrator (CPC). (a) Upper parabolic curve of CPC. Circles are drawn to identify the parabola vertex and the edges of the CPC input and output apertures. (b) Compound parabola concentrator.

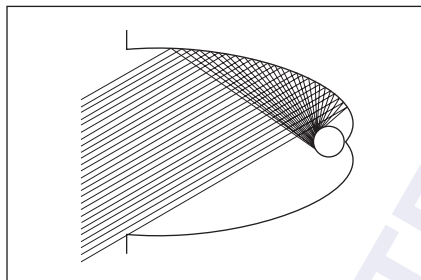


FIGURE 12 A CPC-type reflector for a tubular source.

In 2D, the CPC has been shown to be an ideal concentrator.⁴ In 3D, a small fraction of skew rays are rejected and bound the performance. For a CPC, this skew ray loss increases as the concentration ratio increases. Optimization procedures have produced designs that provide slightly better performance than the CPC in 3D.³⁷ Molledo¹¹⁵ has investigated a crossed cylindrical CPC configuration.

The length of a CPC concentrator is

$$\text{Length} = (R_{\text{in}} + R_{\text{out}}) / \tan(\theta_{\text{in}}) \quad (12)$$

A particularly important aspect of the CPC is that the length is minimized for the case where the input port can see the exit port. If the distance between the input and output ports is made smaller than the CPC length, then rays higher than the maximum design angle could propagate directly from input to output aperture. The reduction in concentration is often small, so many practical systems use truncated CPC-like designs.^{90,116–118} A cone that is the same length as an untruncated CPC will provide lower concentration, but as the cone is lengthened, the performance can exceed the performance of the CPC.

In the literature, CPC is sometimes used to denote any ideal concentrator. For nonplanar receivers, the reflector shape required to create an ideal concentrator is not parabolic. The convention is to use CPC-type.

An example of a CPC-type concentrator for a tubular receiver is shown in Fig. 12. Part of the reflector is an involute and the remainder reflects rays from the extreme input angle to the tangent of the receiver. The design of CPC-type concentrators for nonplanar sources has been described.^{4,116,119,120} Designs with gaps between reflector and receiver^{121–123} and with prisms attached to the receiver^{124,125} have also been described.

CPC-type geometries are also used in laser pump cavities.¹²⁶

Arrays of CPCs or CPC-like structures have also been applied to the liquid crystal displays,^{127,128} illumination,⁸⁹ and solar collection.¹²⁹

CEC

When the source of radiation is located at a finite distance away from the input port of the concentrator, a construction similar to the CPC can be used; however, the reflector surface is now elliptical.¹³⁰ In a similar manner, if the source of radiation is virtual, then the reflector curvature becomes hyperbolic. These two constructions are called compound elliptical concentrators (CEC) and compound hyperbolic concentrators (CHC). A CEC is shown in Fig. 13, where the edge of the finite size source is reimaged onto the edge of the CEC output aperture. Hottel strings can be used to compute the etendue of the collected radiation (see Sec. 39.2). The CPC is similar to the CEC, except the edge of the source is located at infinity for the CPC.

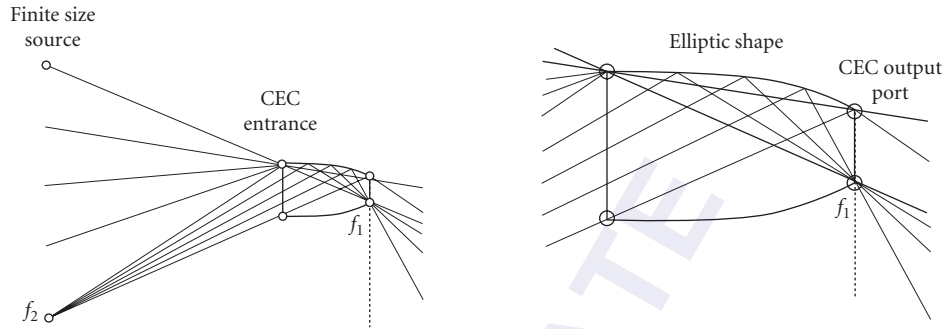


FIGURE 13 Compound elliptical concentrator (CEC). Rays originating at finite size source and collected at the CEC entrance are concentrated at the CEC output port. Edge of source is imaged onto edge of output aperture.

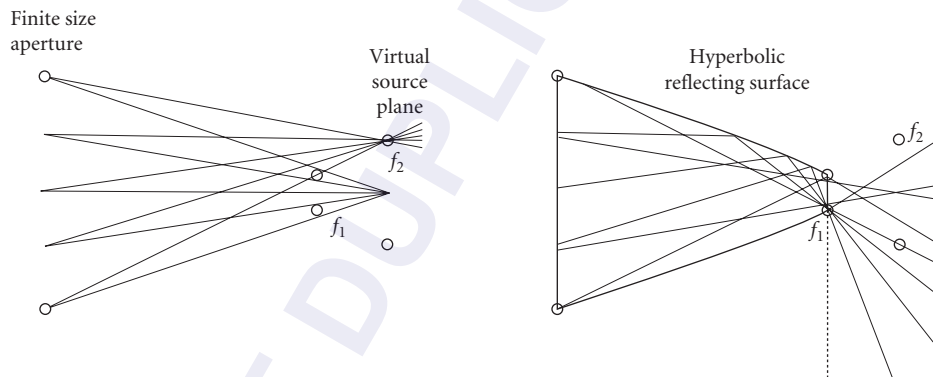


FIGURE 14 Compound hyperbolic concentrator for use with virtual source. Edge of virtual source is imaged onto the edge of the output port.

CHC

A CHC is shown in Fig. 14 where the rays that pass through the finite-size aperture create a virtual source. The CHC concentrates the flux. Arrangements where the CHC is just a cone are possible.¹³¹ Gush¹³² also describes a hyperbolic cone-channel device.

If the virtual source is located at the output aperture of the concentrator, then an ideal concentrator in both 2D and 3D can be created (see Fig. 15). Generically, this is a hyperboloid of revolution and is commonly called a *trumpet*.¹³³ The trumpet maps the edge of the virtual source to $\pm 90^\circ$. The trumpet can require an infinite number of bounces for rays that exit the trumpet at nearly 90° , which can limit the performance when finite-reflectivity mirrors are used.

DCPC

The CPC construction can also be applied when the exit port is immersed in a material with a high index of refraction. Such a device is called a *dielectric compound parabolic concentrator* (DCPC).¹³⁴ This allows concentrations higher than those found when the exit port is immersed in air. In this immersed exit port case, the standard CPC construction can still be used, but the exit port area is now n times smaller in 2D and n^2 smaller in 3D. Another motivation for a DCPC is that the reflectivity

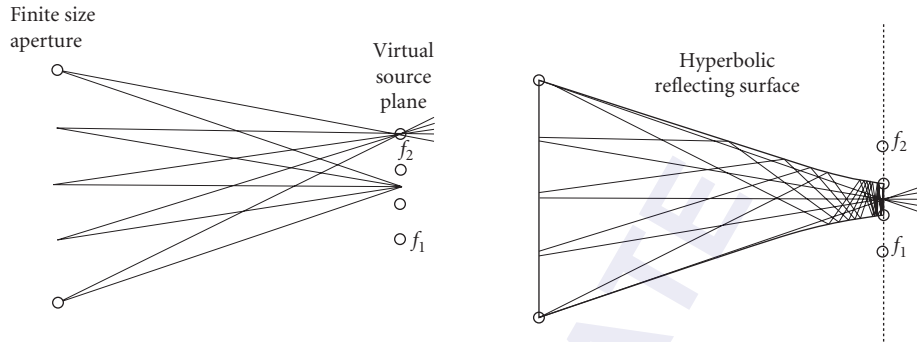


FIGURE 15 Trumpet with exit aperture located at virtual source plane. Rays “focused” at edge of virtual source exit the trumpet at $\pm 90^\circ$.

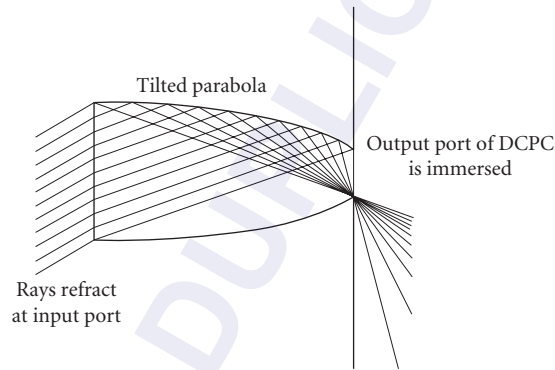


FIGURE 16 Dielectric CPC: same as CPC, except output port is immersed.

losses for a hollow CPC can be minimized because TIR can provide lossless reflections. A DCPC is shown in Fig. 16. A DCPC that uses frustrated TIR at the DCPC to receiver interface has been investigated.¹³⁵

A hollow CPC can also be attached to a DCPC to minimize the size of the DCPC.¹³⁶ Extraction of flux from a high-index medium into a lower-index medium has been investigated^{137,138} using faceted structures.

Multiple Surface Concentrators

CPC-type concentrators can produce a long system when θ_{in} is small. To avoid excessively long systems, other optical surfaces are typically added. Often, a CPC-type concentrator is combined with a primary element such as a condensing lens or a parabolic reflector. There are also designs where the optical surface of the primary and the optical surface of the concentrator are designed together, which often blurs the distinction between primary and concentrator.

Lens + Concentrator There are numerous examples of lenses combined with a nonimaging device in the literature. One approach is to place a lens with the finite-size aperture shown in Figs. 13, 14, or 15 and illuminate the lens with a relatively small angular distribution. The lens produces higher angles that are then concentrated further using the conventional CPC, CEC, CHC, or trumpet. Such a combination typically produces a shorter package length than a CPC-type concentrator.

A lens/CHC can be attractive because the CHC is placed between the lens and lens focal plane, whereas the CEC is placed after the lens focal plane. Thus, the lens/CHC provides a shorter package than a lens/CEC. Both the lens curvature and the CHC reflector can be optimized to minimize the effects of aberrations, or the combination can be adjusted so that the CHC turns into a cone. Investigations of lens/cone combinations include Williamson,¹⁰ Collares-Pereira,⁷⁵ Hildebrand,¹³⁹ Welford,⁴ and Keene.¹⁴⁰ Use of a cone simplifies fabrication complexity.

After the theoretical implications of the trumpet were recognized, Winston¹⁴¹ and O’Gallagher¹³³ investigated the case of a lens/trumpet.

When TIR is used to implement the mirrors, the mirror reflection losses can be eliminated, with the drawback that the package size grows slightly. A design procedure for a lens mirror combination with maximum concentration was presented by Ning,¹⁴² who coined the term *dielectric total internal reflecting concentrator* (DTIRC) and showed significant package improvements over the DCPC.

Eichhorn¹⁴³ describes the use of CPC, CEC, and CHC devices in conventional optical systems, and uses explicit forms for the six coefficients of a generalized quadric surface.¹⁴⁴ A related “imaging” lens and reflector configuration is the Schmidt corrector system.¹⁴⁵

Arrays of lenses with concentrators have also been investigated for use in illumination.¹⁴⁶

Mirror + Concentrator There have been numerous investigations of systems that have a parabolic primary and a nonimaging secondary.^{137,147–150} Other two-stage systems have been investigated.^{29–31,151–153} Kritchman¹⁵⁴ has analyzed losses due to aberrations in a two-stage system.

Asymmetric concentrators have also been investigated. Winston¹⁵² showed that an off-axis parabolic system can improve collection by tilting the input port of a secondary concentrator relative to the parabolic axis. Other asymmetric designs have been investigated.^{155,156}

Related two-mirror “imaging” systems (e.g., Ref. 60, Chap. 16) such as the Ritchey-Chretien can be used as the primary. Ries¹⁵⁷ also investigated a complementary Cassegrain used for concentration.

Simultaneous Multiple Surfaces Minano and coworkers have investigated a number of concentrator geometries where the edge of the source does not touch the edge of a reflector. The procedure has been called simultaneous multiple surfaces (SMS) and builds on the multisurface aspheric lens procedure described by Schultz.⁶⁷ Minano uses the nomenclature *R* for refractive, *X* for reflective (e.g., *refleXive*), and *I* when the main mode of reflection is TIR. The *I* surface is typically used for refraction the first time the ray intersects the surface and TIR for the second ray intersection. The *I* surface is sometimes mirrored over portions of the surface if reflection is desired but the angles do not satisfy the TIR condition. Illustrative citations include RR,¹⁵⁸ RX,^{159,160} and RXI.^{161,162} Example RX and RXI devices are shown in Fig. 17. SMS can be used to design all reflecting configurations.

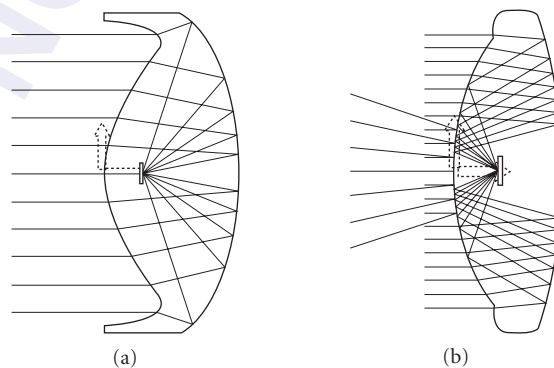


FIGURE 17 RX (a) and RXI (b) concentrators. The receiver is immersed in both cases. In the RXI case, the central portion of the refractive surface is mirrored.

Restricted Exit Angle Concentrators with Lenses

If the angular distribution of the flux is constant across an aperture and the centroid of the distribution is normal to the aperture, then the distribution is called *telecentric*. Many optical systems require the transfer of flux from one location to another. If the desired distribution of flux at the second location is also telecentric, then the system is called doubly telecentric. An afocal imaging system can be modified to provide a doubly telecentric imaging system when used with finite conjugates. An example is shown in Fig. 18.

In nonimaging systems, flux must often be transferred from one aperture to another and the angular distribution at the second aperture must be constant across the aperture. However, the point-by-point mapping is not required in nonimaging systems. A single lens can provide this transfer of radiation, as shown in Fig. 19. This type of lens is a collimator and is similar to a Fourier transform lens; it has been called a *beam transformer*.¹⁶³

The aberrations introduced by a lens can limit the etendue preservation of the aggregate flux collected by the lens. The increase in etendue tends to become more severe as range of angles collected by the lens increases. Somewhere around $f/1$, the change can become quite significant.

θ_1/θ_2

θ_1/θ_2 concentrator^{29,163,164} maps a limited-angle Lambertian distribution with maximum angle θ_1 into another limited-angle Lambertian with maximum angle θ_2 . One version of the θ_1/θ_2 is a compound parabolic construction with a cone replacing part of the small end of the CPC (Ref. 4, pp. 72–74, sec. 5.3). A picture is shown in Fig. 20. The length of a θ_1/θ_2 is given by the same equation as the CPC [e.g., Eq. (12)]. If the θ_1/θ_2 is hollow and $\theta_2 = 90^\circ$, then the cone disappears and the θ_1/θ_2 becomes a CPC.

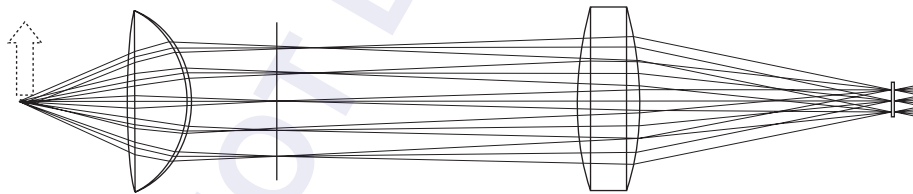


FIGURE 18 Doubly telecentric system showing a twofold increase in size and a twofold decrease in angles.

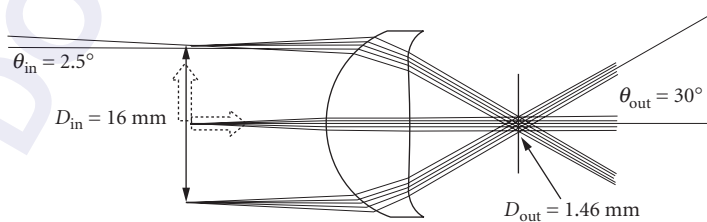


FIGURE 19 Condenser to provide an angle to area transformation and produce a distribution that is symmetric about the optical axis at both input and output planes. $D_{in} \sin \theta_{in} = D_{out} \sin \theta_{out}$.

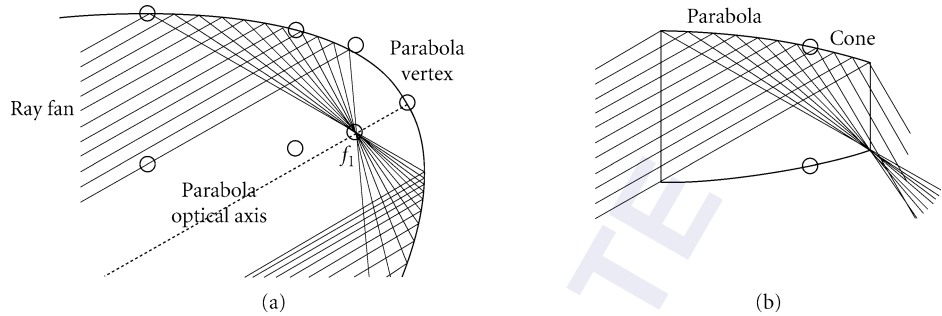


FIGURE 20 θ_1/θ_2 implemented using cone with compound parabolic. (a) Construction information. (b) θ_1/θ_2 , (b) compound parabola + cone.

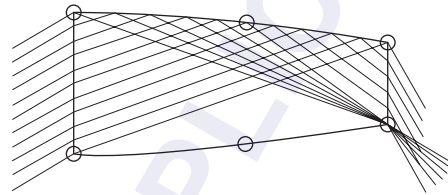


FIGURE 21 Dielectric θ_1/θ_2 implemented using cone with compound parabolic. The rays refract at the input port, total internally reflect (TIR) at the sidewall, and then refract at the output port.

To reduce losses introduced by nonunity mirror reflectivities, a θ_1/θ_2 converter can also be created using a solid piece of transparent material. The design is similar to the hollow θ_1/θ_2 and must take into account the refraction of rays at the input and output port. A dielectric θ_1/θ_2 converter is shown in Fig. 21.

If the θ_1/θ_2 is an all-dielectric device, nearly maximum concentration for an air exit port can be obtained without the mirror losses of a hollow CPC; however, Fresnel surface reflections at the input and output ports should be considered.

Similar to the lens-mirror combination for maximal concentration, the θ_1/θ_2 can be implemented with a lens at the input port. Going one step further, the output port of the θ_1/θ_2 can also incorporate a lens. Similar to a Galilean telescope, the θ_1/θ_2 with lenses at both the input and output ports can provide a short package length compared to embodiments with a flat surface at either port.¹⁶⁴

In many cases, a simple tapered cone can be used to implement a θ_1/θ_2 converter. One reason is that when the difference between θ_1 and θ_2 is small, the compound parabolic combined with a cone is nearly identical to a simple cone. By using an approximation to the ideal shape, some of the flux within θ_1 maps to angles outside of θ_2 . In many practical systems, the input distribution does not have a sharp transition at θ_1 , so a small mapping outside of θ_2 has minor impact. In general, the transition can be made sharper by increasing the cone length; however, there may be multiple local minima for a specific situation.¹⁶⁵

Uniformity at the output port of a θ_1/θ_2 compared to a CPC has been investigated.¹⁶⁶ Tabor¹¹⁸ describes hot spots that can occur with CPCs used in solar energy systems. Edmonds¹²⁵ uses a prism where hot spots that may occur are defocused after propagating through the prism. Emmons¹⁶⁷ discusses polarization uniformity with a dielectric 90/30 converter.

Two CPCs can be placed “throat-to-throat” to provide a θ_1/θ_2 converter that also removes angles higher than the desired input angle.¹⁶⁸ The removal of higher angles is a result of the CPCs ability to transfer rays with angles less than θ_1 and reject rays with angles greater than θ_1 . Two nonimaging concentrators have also been investigated for use in fiber-optic coupling.¹⁶⁹

Tight lightpipe bends are also made possible by inserting a bent lightpipe between the throats of the two concentrators.³³ In the special case of $\theta_1 = \theta_2 = 90^\circ$, a bent lightpipe has also been used as an ideal means to transport the flux for improved packaging.¹⁷⁰

2D versus 3D Geometries

Skew Ray Limits One of the standard design methods for concentrators has been to design a system using a meridional slice and then simply rotate the design about the optical axis. In some cases, such as the CPC used with a disk receiver, this may only introduce a small loss; however, there are cases where the losses are large.

The standard CPC-type concentrator for use with a circular receiver^{90,171} is an example where this loss is large. To assess the loss, compare the 2D and 3D cases. In 2D a concentrator can be designed where the input port is $2R_{in}$ with $a \pm \theta_{in}$ angular distribution; the size of the tube-shaped receiver is $2\pi R_{tube} = 2R_{in} \sin \theta_{in}$. The 2D case provides maximum concentration. In the 3D case where the 2D design is simply spun about the optical axis, the ratio of the input etendue to the output etendue is $\pi R_{in}^2 \pi \sin^2 \theta_{in} / (4\pi R_{tube}^2) = \pi^2 / 4 \sim 2.5$. The dilution is quite large. Feuermann¹⁷² provides a more detailed investigation of this type of 2D-to-3D etendue mismatch.

Ries⁹³ shows how skew preservation in a rotationally symmetric system can limit system performance even though the etendue of the source and receiver may be the same. The skewness distributions (*detendue/dskewness* as a function of skewness) for an example disk, sphere, and cylinder aligned parallel to the optical axis are shown in Fig. 22. The etendues of the disk, sphere, and cylinder are all the same. A second example where a cylinder is coupled to a disk is shown in Fig. 22b. This figure highlights losses and dilution. For skew values where the skewness distribution for the receiver is greater than the source, dilution occurs. Where the source is greater than the receiver, losses occur. If the size of the disk is reduced, then the dilution will also decrease, but the losses increase. If the disk is made bigger, then the losses decrease, but the dilution increases. Skew ray analysis of inhomogeneous sources and targets has also been considered.¹⁷³

Star Concentrator One way of avoiding the skew ray limit in a system with a rotationally symmetric source and a rotationally symmetric receiver is to avoid the use of rotationally symmetric optics to couple the flux from source to target. A device called a *star concentrator* has been devised¹⁷⁴ and shown to improve the concentration. The star concentrator derives its name from the fact that the cross section of the concentrator looks like a star with numerous lobes. The star concentrator is designed using a global optimization procedure. Performance of the star concentrator, assuming a reflectivity of 1, provides performance that exceeds the performance possible using a rotationally

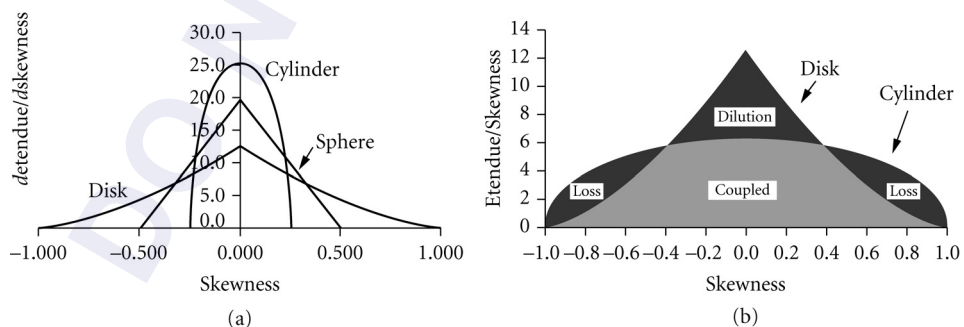


FIGURE 22 Skewness distributions (Ries⁹³). (a) Skewness distribution for a disk with unit radius, a cylinder with length = 2 and radius = 0.25, and a sphere of radius 0.5. The etendue for all three cases is constant, (b) Comparison of the coupling of a disk to a cylinder where the coupled flux is distinguished from the losses and the dilution.

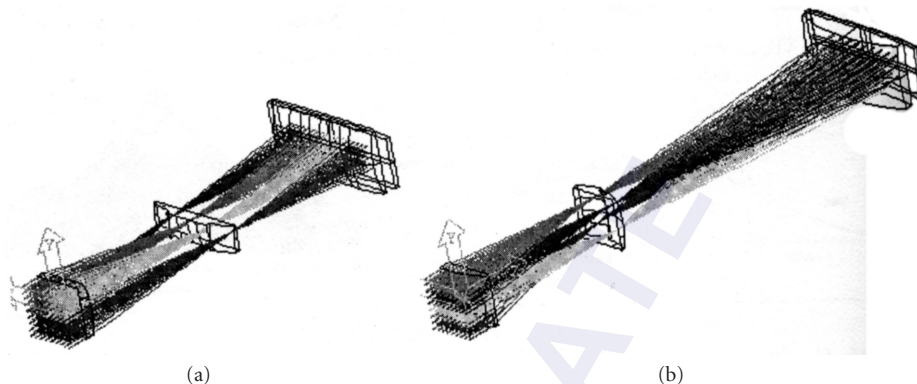


FIGURE 23 Square-to-rectangular mapping with an image dissector (a) and a lens array (b) where the output angular distribution is symmetric. The image dissector uses an array of discrete field lenses. The lens array uses a single field lens.

symmetric system. The star concentrator approach has also been investigated for the case of coupling from a cylindrical source to a rectangular aperture.^{7,175}

Image Dissectors The problem of coupling flux from a rotationally symmetric source to a spectrometer slit motivated some earlier efforts to overcome the skewness limit. Benesch¹⁷⁶ shows a system with two sets of facets that “dissect” an X by Y aperture into N regions of size X/N by Y and relay those images to create an X/N by NY region. The dissection can be done while preserving the NA of the input (e.g., the etendue is preserved but the skew is changed). Benesch calls this system an *optical image transformer* and it is similar in function to Bowen’s image slicer.¹⁷⁷ Refractive versions are also possible (see for example Ref. 178). One issue that must be addressed for some systems, which use these dissectors, is that the adjoining regions of the X/N by Y areas can produce a nonuniform illumination distribution. Laser arrays can also use dissectors to improve performance.^{179,180}

Tandem-lens-array approaches can be used that provide the image dissection and superimpose all the images rather than placing them side by side.^{181,182}

An example of a refractive image dissector compared to a lens array is shown in Fig. 23. The image dissector breaks an input square beam into four smaller square regions that are placed side by side at the output. The lens array breaks the input beam into four rectangles that are superimposed at the output. Both systems use field lenses at the output. The lens array does not have a gap between mapped regions at the output plane compared to the dissector.

Fiber bundles can also be used as image dissectors. For example, the fibers can be grouped in a round bundle at one end and a long skinny aperture at the other end. Feuermann¹⁸³ has provided a recent investigation of this idea. If the fibers in the bundle have a round cross section, tapering the ends of each of the fibers can help improve the uniformity at the output.¹⁸⁴ Special systems that use slowly twisted rectangular light guides have also been investigated for scintillators,^{185–187} document scanning,^{188,189} and microlithography.¹⁹⁰

Bassett¹⁹¹ also shows that arrays of ideal fibers can be used to create ideal concentrators, and Forbes¹⁹² uses the same approach to create a θ_1/θ_2 transformer.

Geometrical Vector Flux

The edge ray design method does not show how the jumble of rays from multiple reflections yields a Lambertian output. The geometric vector formalism was investigated in an attempt to understand the process. Summaries can be found in Bassett¹⁰¹ and Welford.⁴ The proof that a trumpet provides ideal performance has been performed using the flow line concept.¹¹⁴

The vector flux formalism provides the lines of flow and loci of constant geometric vector flux J . The vector J is the PSA for three orthogonal patches where the PSA on the front side of a patch can subtract from the PSA on the back side. The magnitude of J describes what the PSA would be for a surface aligned perpendicular to J . The flow line concept uses the idea that an ideal mirror can be placed along lines of flow without modifying the flow.

In the context of vector flux, the CPC, θ_1/θ_2 , and light cone are discussed in Winston.¹⁴¹ The light cone is demonstrated using an orange in Winston.⁹⁵ The 2D CEC and θ_1/θ_2 are discussed by Barnett.^{193,194} Gutierrez¹⁹⁵ uses Lorentz geometry to show that other rotationally symmetric ideal concentrators can be obtained using surfaces that are paraboloid, hyperboloid, and ellipsoid. Greenman¹⁹⁶ also describes vector flux for a number of ideal concentrators.

Edge Rays

The design of nonimaging optics has been facilitated by the use of the edge ray method^{4,68,101,197} and is sometimes called the string method.^{4,8} In the design of an ideal concentrator, the main principle is that extreme rays at the input aperture must also be extreme rays at the output aperture. Minano¹⁹⁸ describes the edge rays as those that bound the others in phase space. In many cases, the result is that the edge of the source is imaged to the edge of the receiver, and the reflector surface to perform the mapping is a compound conic. When the location of the source edge or the receiver edge is not constant (e.g., a circular source), the required reflectors are not simple conies (e.g., see the preceding section). Gordon¹⁹⁹ describes a complementary construction that highlights the fact that what is conventionally the “outside” of a concentrator can also be used in nonimaging optics.

Davies²⁰⁰ further explores the continuity conditions required by the edge-ray principle. Ries²⁰¹ refines the definition of the edge ray principle using topology arguments. Rabl²⁰² shows examples where the analysis of specular reflectors and Lambertian sources can be performed using only the edge rays.

Tailoring a reflector using edge rays is also described in the next section.

Inhomogeneous Media

A medium with a spatially varying index of refraction can provide additional degrees of freedom to the design of nonimaging optics. In imaging optics, examples of inhomogeneous media include Maxwell’s fisheye and the Luneburg lens.^{4,18} In concentrator design, 2D and 3D geometries have been investigated.^{198,203}

39.6 UNIFORMITY AND ILLUMINATION

In many optical systems, an object is illuminated and a lens is used to project or relay an image of the illuminated object to another location. Common examples include microscopes, slide projectors, overhead projectors, lithography, and machine vision systems. In a projection system the uniformity of the projected image depends upon the object being illuminated, the distribution of flux that illuminates the object, and how the projection optics transport the object modulated flux from object to image. The portion of the system that illuminates the object is often called the *illumination subsystem*.

Beam forming is another broad class of illumination systems. In beam forming, there is often no imaging system to relay an image of the illuminated object. Since the human eye is often used to view the illuminated object, the eye can be considered an imaging subsystem. Application examples include commercial display cases, museum lighting, room lighting, and automotive headlamps.

Approaches to provide uniform object illumination are discussed in this section. The discussed approaches include Kohler/Abbe illumination, integrating cavities, mixing rods, lens array, tailored

optics, and faceted structures. Many of the approaches discussed in this section can be designed to preserve etendue, but the design emphasis is on uniformity so they have been placed in this section rather than the preceding one.

Classic Projection System Uniformity

In classical projection systems, there is typically a source, a transparency (or some other item to be illuminated), and a projection lens. There are two classic approaches to obtaining uniformity in those systems: Abbe and Kohler illumination (see Fig. 24). In both cases, a source illuminates a transparency that is then often imaged onto a screen. *Transparency* is a general term that includes items such as spatial light modulators, film, slides, gobos, and microscope slides.

Some general references regarding Abbe/Kohler illumination include Jones,²⁰⁴ Kingslake,²⁰⁵ O'Shea,²⁰⁶ Wallin,²⁰⁷ Weiss,²⁰⁸ Ray (pp. 455–461),⁶⁶ Bradbury,²⁰⁹ Inoue,²¹⁰ and Born.¹⁸ A description of Kohler with partial coherence is included in Lacombat.²¹¹

Abbe Abbe (also called Nelsonian and Critical illumination) images the source onto the transparency. When the source is spatially uniform, the Abbe approach works fine. A frosted incandescent bulb is a good example. Sometimes a diffuser is inserted in the system and the diffuser is imaged onto the transparency. Fluorescent tubes are another case where a light source provides excellent spatial uniformity. The luminance of sources that are spatially uniform is typically lower than that of sources that are not uniform.

Kohler Even though the source's illuminance distribution may be nonuniform, there are often uniform regions of the source's intensity distribution. This phenomenon is used in the design of Kohler-type systems where the source is imaged into the projection lens aperture stop and a mapping of the source intensity distribution illuminates the transparency. One reason this approach works well is that far from the source, the illuminance distribution over a small area is independent of where the flux originated. This tends to provide uniformity over portions of the source's intensity distribution.

A nearly spherical reflector is sometimes added behind the source to improve the collection efficiency of a Kohler illumination system. Typically, the spherical reflector produces an image of the source that is aligned beside the primary source. The direct radiation from the source and the flux

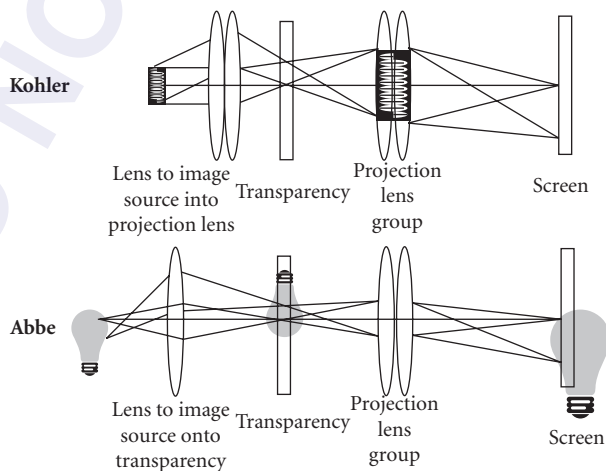


FIGURE 24 Abbe versus Kohler layout comparison.

from the image of the source then become the effective source for the Kohler illuminator. If the source is relatively transparent to its own radiation, then the image of the source can be aligned with the source itself. Blockage by electrodes and leads within the bulb geometry often limit the flux that can be added using the spherical reflector. Scattering and Fresnel surface reflections also limit the possible improvement. Gordon²¹² provides a recent description of a back reflector used with filament sources.

In Kohler illuminators, tailoring the illuminance distribution at the transparency to provide a desired center-to-edge relative illumination can be performed using the edge ray approach described by Zochling.²¹³ Medvedev⁸² also applied the edge ray principle to the design of a projection system. The other tailoring approaches described later in this section may also be relevant.

Integrating Cavities

Integrating cavities are an important example of nonimaging optics that are used in precision photometry,¹⁶ backlighting systems, collection systems (Basset,¹⁰¹ Sec. 7, and Steinfeld²¹⁴), scanners,⁴¹ and reflectometers.²¹⁵ They provide a robust means for creating a distribution at the output port that is independent of the angular and spatial luminance distributions at the input port, thereby removing “hot spots” that can be obtained using other methods. Integrating cavities are limited by a trade-off between maximizing the transmission efficiency and maximizing the luminance at the output port. The angular distribution exiting the sphere is nominally an unclipped Lambertian, but there are a number of techniques to modify the Lambertian distribution.²¹⁶

Nonuniformities with Integrating Spheres Integrating cavities, especially spherically shaped cavities, are often used in precision photometry where the flux entering the sphere is measured using the flux exiting the output port of the sphere. For precision photometry, the angular distribution of the exiting flux must be independent of the entering distribution. In general, the flux is scattered multiple times before exiting the sphere, which provides the required independence between input and output distributions. The two main situations where this independence is not obtained are when the output port views the input port and when the output port views the portion of the sphere wall in which the direct illumination from the input port illuminates the sphere wall.

Figure 25 shows an integrating sphere with five ports. A source is inserted into the upper port and the other four ports are labeled. Figure 26 shows the intensity distribution that exits output port 1 of the sphere shown in Fig. 25. The structure in the exiting intensity distribution is roughly Lambertian, but there are nonuniformities in the distribution that occur at angles that can “see” the source or the other output ports. To highlight this effect, rotated views of the sphere are included in Fig. 26. The bright peak in the intensity distribution occurs at the view angle where the output port

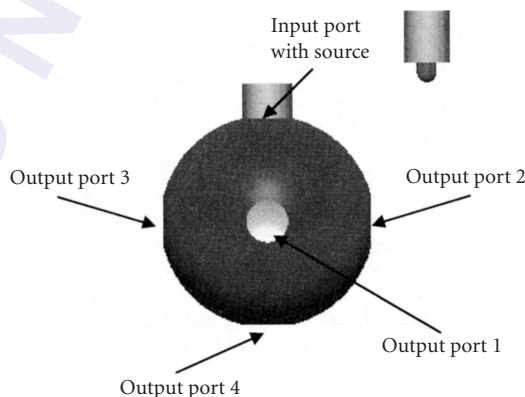


FIGURE 25 Spherical integrating cavity with one input port and four output ports.

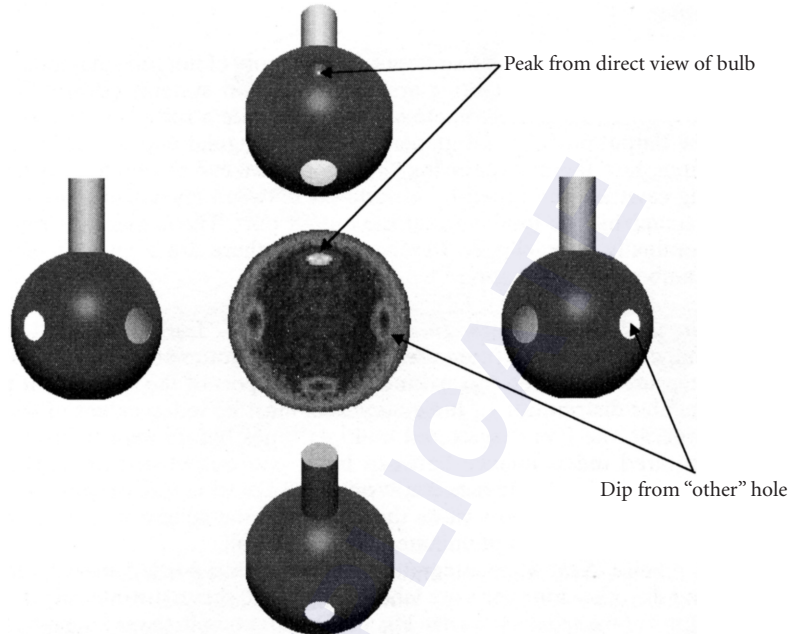


FIGURE 26 Relationship between the intensity distribution from the flux that exits an integrating sphere to what you see when looking back into the sphere.

can “see” direct radiation from the source. The dips in the intensity distribution occur at view angles where the output port can “see” the unused ports in the sphere.

If the source is replaced by a relatively collimated source that illuminates only a small region of the cavity interior (e.g., a collimated laser), then the output port intensity also shows a peak as shown in Fig. 27.

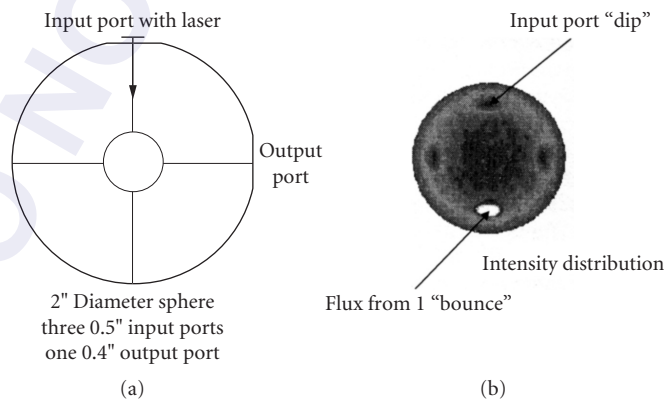


FIGURE 27 Intensity distribution (a) for an integrating sphere (b) that is illuminated with a highly collimated source. The flux that exits the sphere with only one bounce produces a higher intensity than the flux that rattles around and exits in a more uniform manner. The intensity distribution is shown normal to the orientation of the output port.

Three methods that are often used in precision photometry to create an output intensity distribution that is independent of the input involve the use of baffles, satellite spheres, and diffusers. Baffles block radiation from exiting the sphere directly, thereby forcing the exiting flux to scatter multiple times before exiting the sphere. A *satellite sphere* is a secondary sphere that is added to the primary sphere. Diffusers superimpose flux from multiple angles to provide an averaging effect that removes observable structure.

Finite-thickness walls at the exit port will also alter the view angle of external ports and can improve the on-axis intensity because light that reflects off of the walls can go back into sphere and have a chance to add to the output.

Efficiency versus Luminance An important aspect of an integrating cavity is that structure in the angular and spatial distribution of the source's light flux are removed by multiple reflections from the diffuse reflector. Under the assumption that the flux inside of an integrating sphere is scattered in a Lambertian manner, Goebel²¹⁷ has provided a description of the relationship between the fraction of flux that enters the sphere and the fraction of flux that exits the sphere (see also the Technical Information chapter in the LabSphere catalog). Goebel's results can be used to show that the ratio of the flux that exits the sphere via one port to the flux that exits out another port is

$$\text{Efficiency} = \eta = \frac{f_{\text{output}} R_{\text{sphere}}}{1 - f_{\text{input}} R_{\text{input}} - f_{\text{output}} R_{\text{output}} - f_{\text{sphere}} R_{\text{sphere}}} \quad (13)$$

where $f_{\text{sphere}} = (\text{total sphere surface area} - \text{port areas})/\text{total sphere surface area}$

$f_{\text{output}} = \text{output port area}/\text{total sphere surface area}$

$f_{\text{input}} = \text{input port area}/\text{total sphere surface area}$

$R_x = \text{reflectivity of } x$

When the input and output have zero reflectivity, the result is

$$\text{Efficiency} = \eta = \frac{f_{\text{output}} R_{\text{sphere}}}{1 - f_{\text{sphere}} R_{\text{sphere}}} \quad (14)$$

Averaging the input flux over a hemisphere gives a ratio of the luminance at the output port to the luminance at the input port of

$$\frac{\text{Output luminance}}{\text{Input luminance}} = \frac{\eta / \text{area}_{\text{output}}}{1 / \text{area}_{\text{input}}} = \frac{\eta f_{\text{input}}}{f_{\text{output}}} = \frac{R_{\text{sphere}} f_{\text{input}}}{1 - f_{\text{input}} R_{\text{input}} - f_{\text{output}} R_{\text{output}} - f_{\text{sphere}} R_{\text{sphere}}} \quad (15)$$

When the input and output ports have zero reflectivity, then the result is

$$\frac{\text{Output luminance}}{\text{Input luminance}} = \frac{R_{\text{sphere}} f_{\text{input}}}{1 - f_{\text{sphere}} R_{\text{sphere}}} \quad (16)$$

The efficiency and the ratio of input to output luminances are shown in Fig. 28 for the case of $R_{\text{sphere}} = 98.5$ percent.

When the fractional input port size equals the fractional output port size, the sphere transmission is less than 50 percent and the output luminance is less than 50 percent of the input luminance. Increasing the output port size provides higher transmission, but the luminance of the output port is now less than 50 percent of that of the input port. Reducing the output port size can increase the output luminance to values greater than the input luminance; however, the transmission is now less than 50 percent. Thus, an integrating cavity offers the possibility of improving uniformity at the expense of either luminance or efficiency.

Integrating Cavity with Nonimaging Concentrator/Collectors The flux exiting an integrating cavity tends to cover a hemisphere. Adding a nonimaging collector, such as a CPC, to the output port can

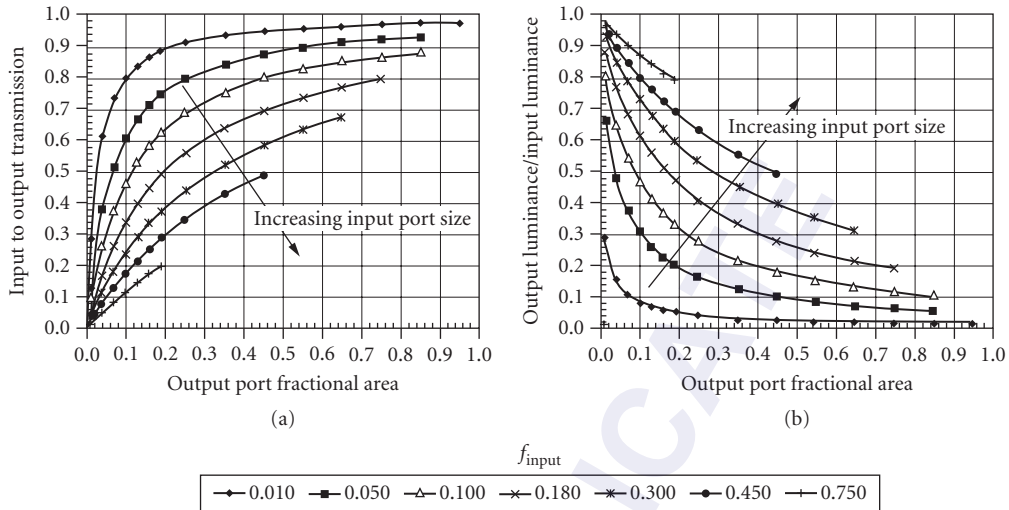


FIGURE 28 Input to output transmission (a) and relative output luminance (b) for the case of sphere reflectivity = 0.985. The different curves are for various input port fractional areas.

reduce the angular spread. The integrating cavity can also include a concentrator at the input port so as to minimize the size of the input port if the angular extent of the radiation at the input port does not fill a hemisphere.

A CEC can be added to the output port of an integrating cavity so that the portion of the sphere surface that contributes to the output distribution is restricted to a region that is known to have uniform luminance.²¹⁵ Flux outside of the CEC field of view (FOV) is rejected by the CEC. CPCs and lenses with CHCs have also been investigated.^{215,218,219}

Modifying the Cavity Output Distribution Prismatic films can also be placed over the exit port of an integrating cavity, which is an approach that has been used by the liquid crystal display community (e.g., the BEFII films commercialized by 3M). The prismatic films increase the intensity at the angles that exit the sphere by reflecting the undesired angles back into the sphere so that they can be rescattered and have another chance to exit the sphere. For systems that require polarized light, nonabsorbing polarizing films can also be used to return the undesired polarization state back into the sphere.

Another approach to selectively controlling the exit angles is the addition of a hemispherical or hemiellipsoid reflector with a hole in the reflector. Webb²²⁰ described this approach for use with lasers.

In a similar manner, reflecting flux back through the source has been used with discharge sources, and is especially effective if the source is transparent to its own radiation.

Mixing Rods (Lightpipes)

When flux enters one end of a long lightpipe, the illuminance at the other end can be quite uniform. The uniformity depends upon the spatial distribution at the lightpipe input, the intensity distribution of the input flux, and the shape of the lightpipe.³⁰⁶ Quite interestingly, straight round lightpipes often do not provide good illuminance uniformity, whereas some shapes like squares and hexagons do provide good illuminance uniformity. This is shown in Fig. 29.

Shapes (Round versus Polygon) Square lightpipes have been used very successfully to provide extremely uniform illuminance distributions. To understand how a lightpipe works, consider the

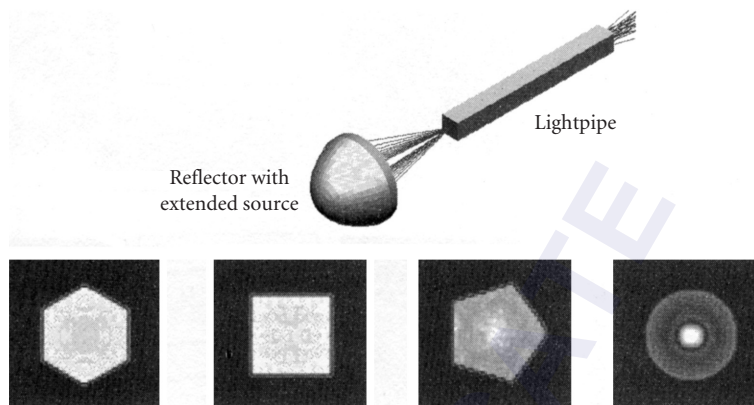


FIGURE 29 When flux from a source is coupled into a mixing rod, the illuminance at the output end of the mixing rod is significantly affected by the shape of the lightpipe. Square and hexagonal lightpipes work much better than round or many other regular polygons.

illuminance distribution that would result at the output face of the lightpipe if the lightpipe's sidewalls were removed. Including the sidewalls means that various regions of the illuminance distribution with no lightpipe are superimposed. In the case of a square, the multiple regions add in a controlled manner such that there is improved uniformity by virtue of the superposition of multiple distributions. If the source distribution is symmetric about the optical axis, then distributions with a positive slope are superimposed onto distributions with a negative slope. This means that the uniformity from N distributions can be better than the addition of N uncorrelated distributions that would provide uniformity with a $1/\sqrt{N}$ standard deviation.

One way to view the operation of a rectangular mixing rod is to consider the illuminance distribution that would occur at the output face of the mixing rod if no sidewall reflections were present. The sidewalls force subregions of this unreflected illuminance distribution to superimpose in a well-controlled manner.²²¹ This is shown in Fig. 30, where a Williamson-type construction is shown on the left. The illuminance distribution at the lightpipe output with no sidewalls is also shown, as is the superposition of the subregions when sidewall reflection is included. Every other subregion is inverted because of the reflection at the sidewall. This inversion provides a difference between a lightpipe used for homogenization and a lens array.

Figure 31 shows the results of a Monte Carlo simulation where a source with a Gaussian illuminance distribution and a Gaussian intensity distribution are propagated through a square lightpipe. The illuminance at the lightpipe input and output faces is shown in the top of the figure. The illuminance distribution in the lower right shows the illuminance distribution if no sidewall reflections occurred. Lines are drawn through the no-sidewall reflection illuminance distribution to highlight the subregions that are superimposed to create the illuminance distribution when the lightpipe is present. The lower left illuminance distribution results from propagating the flux from the output end of the lightpipe back to the input end. Multiple images of the source are present, similar to the multiple images observed in a kaleidoscope. This kaleidoscope effect is one reason why lightpipe mixing is sometimes considered to be the superposition of multiple virtual sources.

In general, if the cross-sectional area is constant along the lightpipe length, and the area is covered completely by multiple reflections with respect to straight sidewalls of the lightpipe, then excellent uniformity can be obtained. The shapes that provide this *mirrored tiling* include squares, rectangles, hexagons, and equilateral triangles.²²² A square sliced along its diagonal and an equilateral triangle sliced from the center of its base to the apex will also work. Pictures of these shapes

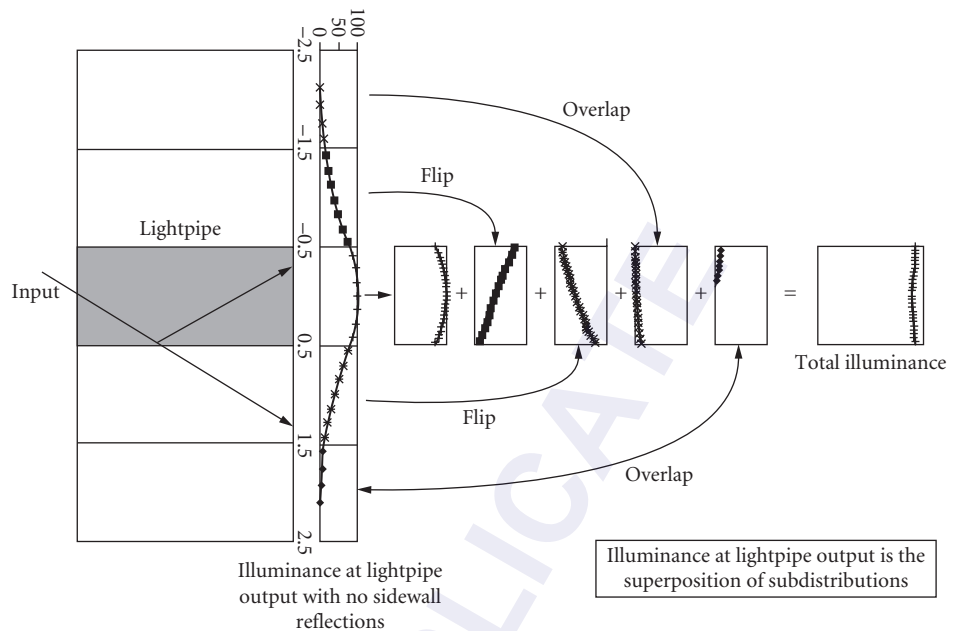


FIGURE 30 Flip-and-fold approach²²¹ for use with rectangular mixing rods. Subregions of the illuminance at the lightpipe output with no sidewall reflections are superimposed to create a uniform distribution. The lightpipe is shown shaded and a Williamson construction is included. The path of one ray is shown with and without sidewall reflection.

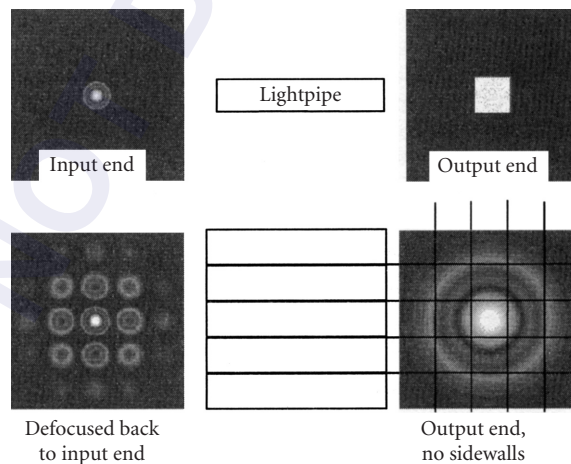


FIGURE 31 Example with a source that has a Gaussian intensity distribution. If the source is propagated to the end of the lightpipe without sidewall reflections, the illuminance distribution at the output can be sliced into distinct regions. These regions are superimposed to create the illuminance distribution that actually occurs with the sidewall reflections. If the flux at the lightpipe output is propagated virtually back to the input end, the multiple images can be observed.

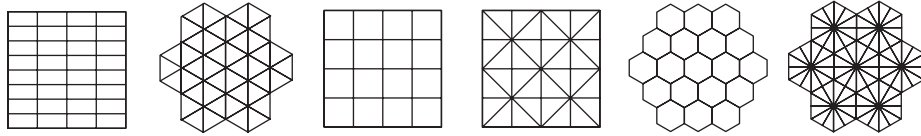


FIGURE 32 Mirror-tiled shapes that can ensure uniformity with adequate length.

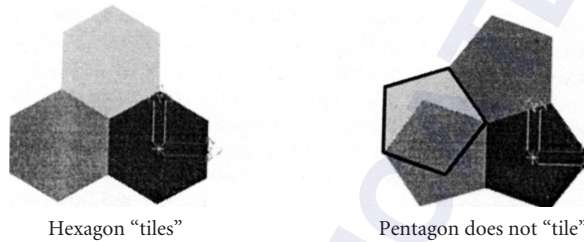


FIGURE 33 Diagram showing why pentagon does not tile like the hexagon. Simulation results are shown in Fig. 29.

arranged in a mirror-tiled format are shown in Fig. 32. The lightpipe length must be selected so as to provide an adequate number of mirrored regions. In particular, the regions near the edge that do not fill a complete mirror-tiled subregion must be controlled so that they do not add a bias to the output distribution.

Shapes that do not mirror-tilde may not provide ideal uniformity as the lightpipe length is increased, at least for a small clipped Lambertian source placed at the input to the lightpipe. One explanation is shown in Fig. 33.

Round lightpipes can provide annular averaging of the angular distribution of the flux that enters the input end. Because of this, the illuminance distribution along the length of the lightpipe tends to be rotationally symmetric; however, the illuminance distribution tends to provide a peak in the illuminance distribution at the center of the lightpipe. The magnitude of this peak relative to the average varies along the length of the lightpipe. If the distribution entering the lightpipe is uniform across the whole input face, and the angular distribution is a clipped Lambertian, then the round lightpipe maintains the uniformity that was present at the input end. However, if either the input spatial distribution or the input angular distributions are not uniform, then nonuniform illuminance distributions can occur along the lightpipe length. This presents the interesting effect that a source can be coupled into a square lightpipe of adequate length to provide a uniform illuminance distribution at the lightpipe output face; however, if the angular distribution is nonuniform and coupled into a round lightpipe, then the illuminance at the output face of the round lightpipe may not be uniform.

Periodic Distributions If the unmirrored illuminance distribution has a sinusoidal distribution with certain periods that match the width of the lightpipe, then the peaks of the sinusoid can overlap. A Fourier analysis of the lightpipe illuminance distribution can help to uncover potential issues.²²¹

Length The length of the lightpipe that is required to obtain good uniformity will depend upon the details of the source intensity distribution. In general, overlapping a 9×9 array of mirrored regions provides adequate uniformity for most rotationally symmetric distributions. For an $f/1$ distribution with a hollow lightpipe, this means that the lightpipe aspect ratio (length/width) should be greater than 6:1. For a lightpipe made from acrylic, the aspect ratio needs to be more like 10:1.

Solid versus Hollow In addition to a solid lightpipe requiring a longer length to obtain a given uniformity, the design of the solid case should consider the following: Fresnel surface losses at the input

and output ends, material absorption, cleanliness of the sidewalls (heat shrink Teflon can help here), and chipping of the corners. For the hollow case, some considerations include dust on the inside, coatings with reflectivities of less than 100 and angular/color dependence, and chips in mirrors where the sidewalls join together.²²¹

With high-power lasers and high-power xenon lamps, a solid coupler may be able to handle the average power density, but the peak that occurs at the input end may introduce damage or simply cause the solid lightpipe to shatter.

Hollow lightpipes can also be created using TIR structures.^{223,224}

Angular Uniformity Although the illuminance at the output end of a lightpipe can be extremely uniform, the angular distribution may not be, especially the “fine” structure in the distribution. One method to smooth this fine structure issue is to add a low-angle diffuser at the lightpipe output end.²²⁵ A diffuser angle that is about the angular difference between two virtual images is often sufficient. Since multiple images of the source are usually required to obtain illuminance uniformity at the lightpipe output, the diffuser does not increase the etendue much. Making the lightpipe longer reduces the angular difference between neighboring virtual images and the required diffuser angle is reduced accordingly. Combinations of lightpipes and lens arrays have also been used.²²⁶

Tapered Lightpipes If the lightpipe is *tapered*, meaning that the input and output ends have different sizes, excellent uniformity can still be obtained. Since the lightpipe tends to need a reasonable length to ensure uniformity, the tapered lightpipe can provide reasonable preservation of etendue. This etendue preservation agrees with Garwin’s³ adiabatic theorem. In general, for a given light pipe length, better mixing is obtained by starting with high angles at the light pipe input and tapering to lower angles than if the input distribution starts with the lower angles.²²⁷ A Williamson construction can be used to explain why the tapered pipe has more virtual images than a straight tapered case. Tapered lightpipes have been used with multimode lasers²²⁸ and solar furnaces.²²²

An example with a small Lambertian patch placed at the entrance to a tapered lightpipe and a straight lightpipe is shown in Fig. 34. The input NA for the untapered case is 0.5/3 for the 100-mm-long

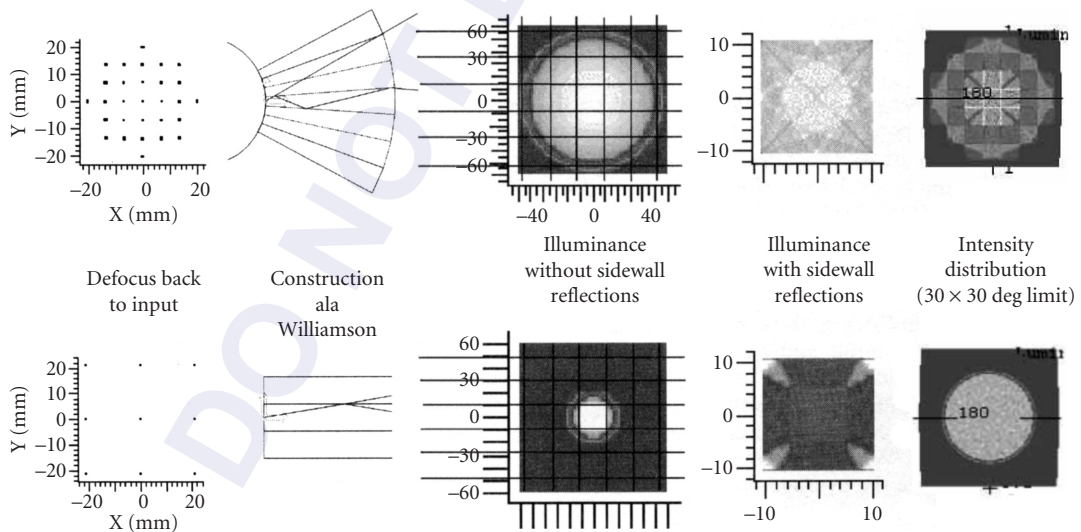


FIGURE 34 Tapered lightpipe example. The upper and lower rows are for the tapered and untapered cases, respectively. Moving left to right, the drawings show the source when defocused from the lightpipe output back to the lightpipe input, a Williamson construction, the no-sidewall illuminance at the lightpipe output, the illuminance at the lightpipe output, and the angular distribution of the flux exiting the lightpipe.

straight lightpipe and 0.5 for the 100-mm-long tapered lightpipe. The input size of the straight lightpipe is a square with a 30-mm diagonal. The input size of the tapered lightpipe is a 10-mm diagonal and the output size is a 30-mm diagonal. The Williamson construction shows that the extreme rays for the 0.5-NA meridian have almost three bounces for the tapered case, but only slightly more than one bounce for the extreme ray in the 0.5/3-NA straight lightpipe case. The illuminance distribution that would occur at the lightpipe output if the lightpipe sidewalls were removed is shown in Fig. 34. Gridlines are superimposed on the picture to highlight that the tapered case has numerous regions that will be superimposed, but the superposition will be small for the straight lightpipe case. The tapered case provides superior uniformity, although the output NA in the tapered case is higher than 0.5/3 because the angular distribution has “tails.” These results are obtained using a source that is a small Lambertian patch. The comparison is typically far less drastic when the input face of the lightpipe is better filled.

Applications of Mixing Rods Mixing rods have found use in a number of applications. Examples include projection systems,²²⁹ providing uniformity with a bundle of fibers using discharge lamps,²²⁷ fiber-to-fiber couplers,²³⁰ and solar simulators.²²²

Mixing rods have also been used with coherent sources. Some applications include semiconductor processing;²³¹ hyperpigmented skin lesions;^{232,233} lithography, laser etching, and laser scanning;²²⁸ and high-power lasers.²³⁴ If the coherence length of the laser is not small compared to the path length for overlapping distributions, then speckle effects must be considered.^{231,234} Speckle effects can be mitigated if the lightpipe width is made large compared to the size of the focused laser beam used to illuminate the input end of the lightpipe (e.g., Refs. 231 and 236). A negative lens for a laser illuminated lightpipe is described by Grojean;²³⁵ however, this means that etendue of the source is less than the etendue of the lightpipe output.

Speckle can be controlled in lightpipes by adding a time-varying component to the input distribution. Some methods include placing a rotating diffuser at the lightpipe input end and moving the centroid of the distribution at the input end. In principle, the diffuser angle can be small so that the change in etendue produced by the diffuser can also be small.

Coherent interference can also be used advantageously, such as in the use of mixing rods of specific lengths that provide sub-Talbot plane reimaging. This was successfully applied to an external cavity laser.²³⁷

In addition to uniformity and the classic kaleidoscope, mixing rods can also be used in optical computing systems by reimaging back to the lightpipe input.^{238,239}

Lens Arrays

Lens arrays are used in applications including Hartmann sensors,²⁴⁰ diode laser array collimators,²⁴¹ beam deflection,²⁴² motion detectors,²⁴³ fiber-optic couplers,²⁴⁴ commercial lighting, and automotive lighting. In commercial and automotive lighting, lens arrays are often used to improve uniformity by breaking the flux into numerous subregions and then superimposing those subregions together. This superposition effect is illustrated in Fig. 35, where fans of rays illuminate three lenslets and are superimposed at the focal plane of a condensing lens. The condensing lens is not used if the area to be illuminated is located far from the lens array.

If the incident flux is collimated, then the superposition of wavefronts provides improved uniformity by averaging the subregions. However, when the wavefront is not collimated, the off-axis beams do not completely overlap the on-axis beams (see Fig. 35*b*). Tandem lens arrays add additional degrees of design freedom and can be used to eliminate the nonoverlap problem (see Fig. 35*c*).

Literature Summary Some journal citations that describe the use of lens arrays with incoherent sources include Zhidkova's²⁴⁵ description of some issues regarding the use of lens arrays in microscope illumination. Ohuchi²⁴⁶ describes a liquid crystal projector system with an incoherent source and describes the use of two tandem lens arrays with nonuniform sizes and shapes for the second lens array. This adjustment to the second lens array provides some adjustment to the coma introduced by the conic reflector.

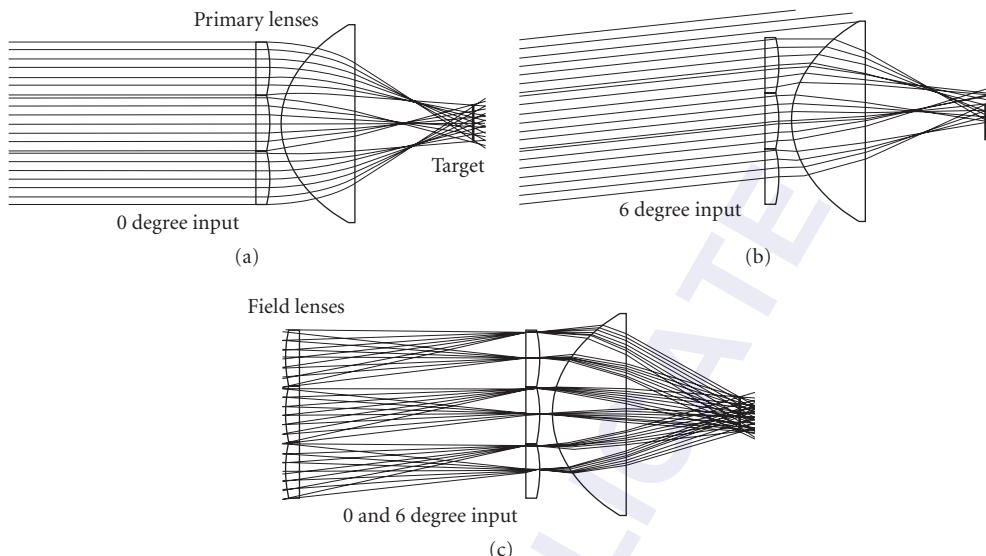


FIGURE 35 Lens arrays used for illumination. (a) Superposition of three subregions of the input wavefront. (b) Off-axis angles are displaced relative to the on-axis angles. (c) Tandem lens array case, where the addition of an array of field lenses allows the superposition to work with a wide range of angles.

Rantsch²⁴⁷ shows one of the earliest uses of tandem lens arrays. Miles²⁴⁸ also uses a tandem lens array system and is concerned about the distribution in the pupil. Ohtu²⁴⁹ uses two sets of tandem lens arrays in a lithography system where one array controls the spatial distribution of the image and the other controls the angular (pupil) distribution of the image. Kudo²²⁶ combines a lens array and lightpipe to provide the same type of control over the illuminance and intensity distributions. Matsumoto²⁵⁰ adds a second set of tandem lens arrays with a sparse fill to correct the fact that an odd spatial frequency term in the illuminance distribution can result in a final nonuniformity, independent of the number of lenslets. Van den Brandt²⁵¹ shows numerous configurations using tandem lens arrays, including the use of nonuniform sizes and shapes in the second lens array. Watanabe²⁵² provides design equations for tandem lens arrays in a converging wavefront.

Journal publications that include lens arrays often involve the use of lens arrays with coherent sources rather than incoherent sources. Deng²⁵³ shows a single lens array with a coherent source and shows that the distribution at the target is a result of diffraction from the finite size of each lenslet and interference between light that propagates through the different lenslets. Nishi²⁵⁴ modifies the Deng system to include structure near the lenslet edges so as to reduce the diffraction effects. Ozaki²⁵⁵ describes the use of a linear array of lenses with an excimer laser, including the use of lenslets where the power of the lenslets varies linearly with distance from the center of the lens array. Glockner²⁵⁶ investigates the performance of lens arrays under coherent illumination as a function of statistical variation in the lenslets. To reduce interference effects with coherent sources, Bett²⁵⁷ describes the use of a hexagonal array of Fresnel zone plates where each zone incorporates a random phase plate. Kato²⁵⁸ describes random phase plates.

Single-Lens Arrays Single-lens-array designs for uniformity tend to have two extremes. One occurs when adding the lens array significantly changes the beam width. This is the beam-forming category where the shape of the beam distribution is essentially determined by the curvature and shape of the lenslets. The other extreme occurs when adding the lens array does not significantly change the beam width. This is the beam-smearing category and is useful for removing substructure in the beam, such as structure from filament coils or minor imperfections in the optical elements. There are many designs that lie between these two extremes.

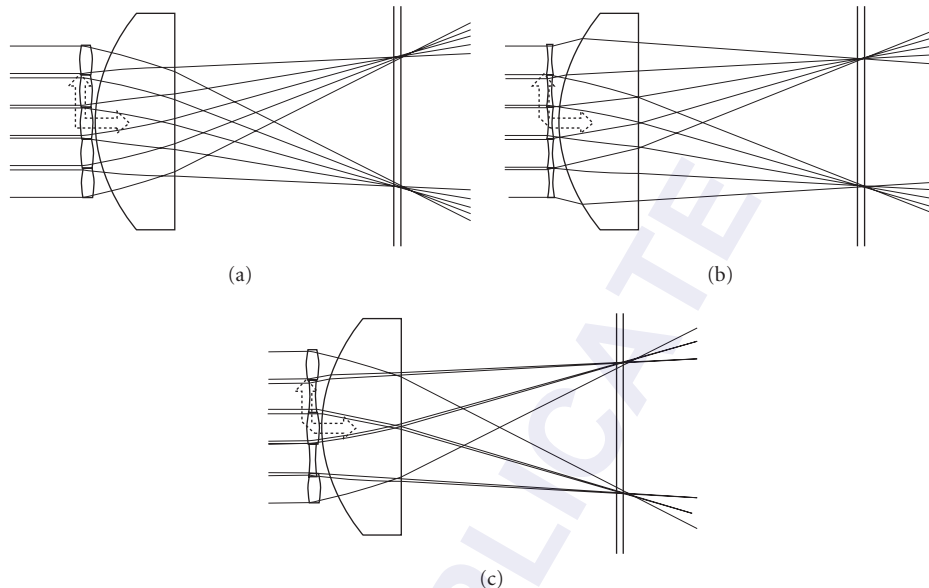


FIGURE 36 Superposition using an array of positive lenses (a), negative lenses (b), and also a hybrid system (c) that uses negative and positive lenses. A collimated wavefront passing through each lenslet is mapped over the target. The hybrid system offers the possibility of a continuous surface with no slope discontinuities.

Figure 36 depicts the beam-forming case where the source divergence is assumed to be small. Positive and/or negative lenses can be used to create the distribution. Hybrid systems where mixtures of positive and negative lenses are used provide the added benefit that the slope discontinuity between neighboring lenses can be eliminated. Slope discontinuities can degrade performance in many systems.

Tandem-Lens Arrays The use of a single-lens array has two issues that sometimes limit its use: First, the lens array increases the etendue because the angles increase but the area stays the same. Second, the finite size of the source introduces a smearing near the edges of the “uniform” distribution. Adding a second lens array in tandem with the first can minimize these limitations.

Figure 37a shows a 1:1 imager in an Abbe (critical) illumination system. Figure 37b shows the same imager with two-lens arrays added. Adding the lens arrays has changed the system from Abbe to Kohler (see the preceding text). In addition, as shown in Fig. 38, the lenslets provide channels that can be designed to be independent of one another. There are three main reasons why two tandem-lens arrays provide uniformity:

1. Each channel is a Kohler illuminator, which is known to provide good uniformity assuming a slowly varying source-intensity distribution.
2. Adding all the channels provides improved uniformity as long as the nonuniformities of one channel are independent of the nonuniformities of the other channels (e.g., a factor of $1/\sqrt{N}$ improvement by adding independent random variables).
3. For symmetric systems, the illuminance from channels on one side of the optical axis tends to have a slope that is equal but opposite to the slope on the other side of the optical axis. The positive and negative slopes add to provide a uniform distribution.

Reflector/Lens-Array Combinations Using a reflector to collect the flux from the source can provide a significant increase in collection efficiency compared to the collection efficiency of a single lens. The

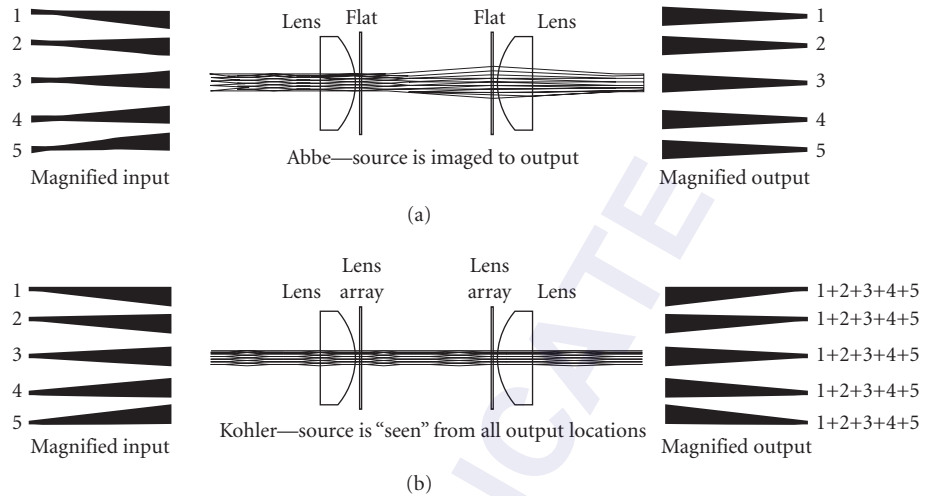


FIGURE 37 Critical illumination (Abbe, *a*) compared to Kohler (*b*). In Abbe, the source is imaged to the output. In Kohler, most of the source is “seen” from all output points.

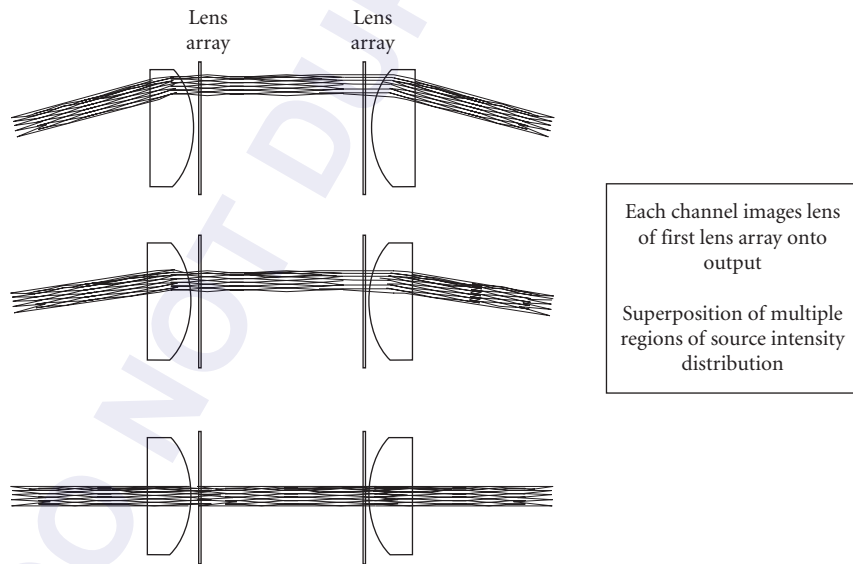


FIGURE 38 Multiple channels in a tandem-lens array. Each channel images a subregion of the flux at the first lens array. The flux from the different channels is superimposed.

drawback is that a single conic reflector does not generally provide uniform illuminance. The use of a lens array in this situation allows the illuminance distribution to be broken into multiple regions, which are then superimposed on top of each other at the output plane.

Examples showing the illuminance distribution at the first lens array, the second lens array, and the output plane are shown in Fig. 39 for the case of a lens collector, Fig. 40 for the case of a parabola

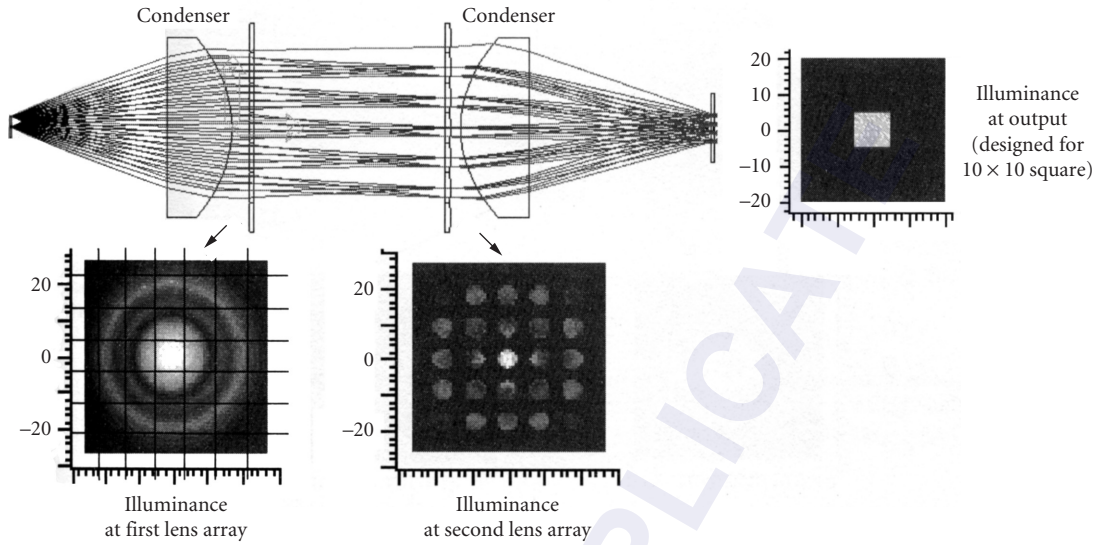


FIGURE 39 Source/condenser/tandem lens arrays/condenser configuration. Illuminance at first lens array is broken into subregions that are superimposed at the target by the second lens array/second condenser. All units are millimeters.

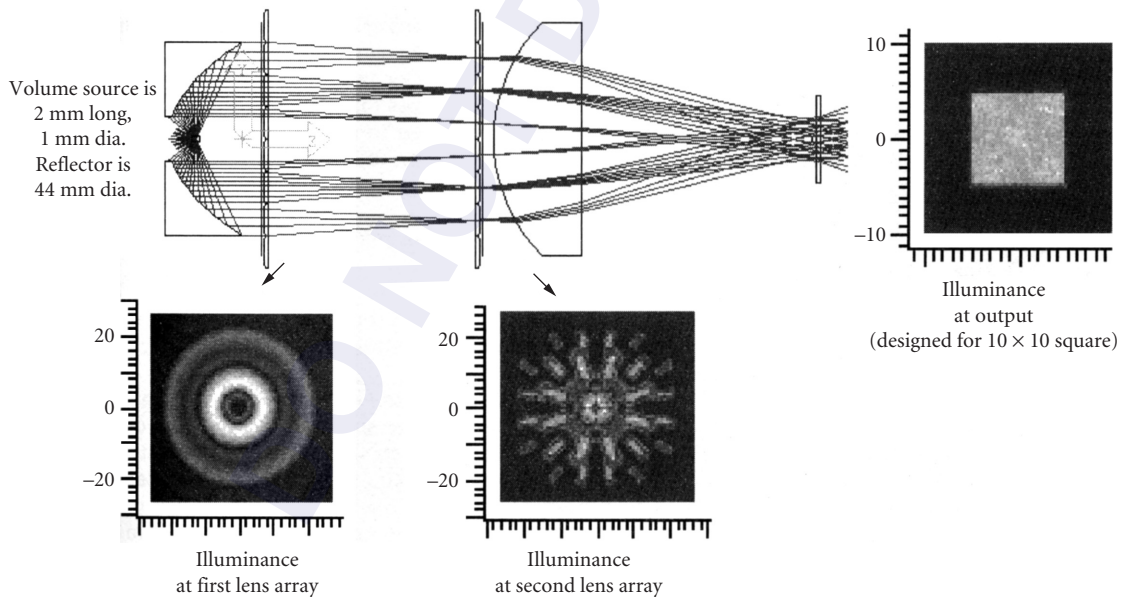


FIGURE 40 Parabola/tandem lens arrays/condenser configuration. Nonuniform magnification from the reflector (coma) is seen in the array of images at the second lens array. The illuminance at the target is the superposition of subregions at the first lens array. All units are millimeters.

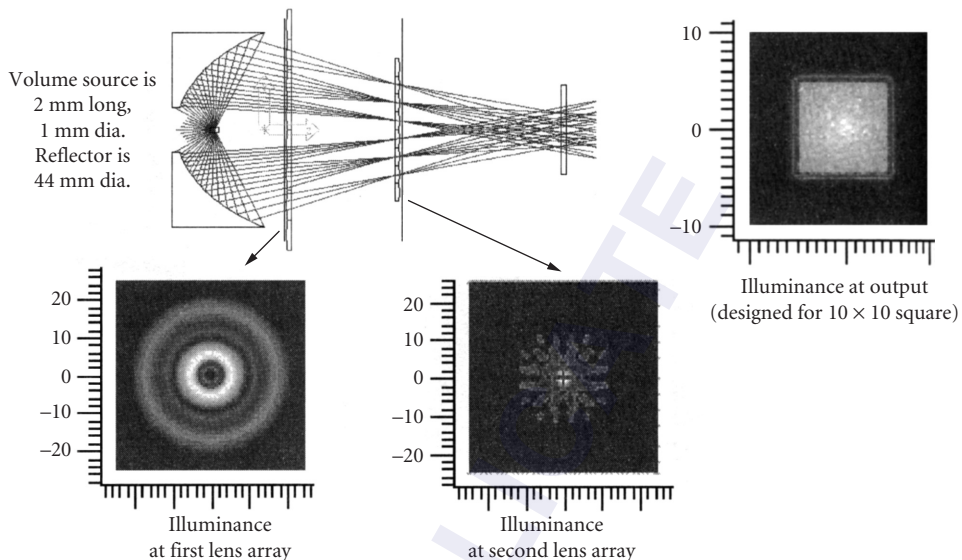


FIGURE 41 Ellipse/tandem lens arrays/condenser configuration. Nonuniform magnification from the reflector (coma) is seen in the array of images at the second lens array. The illuminance at the target is the superposition of subregions at the first lens array. The spacing between lenses is different for the two lens arrays. All units are millimeters.

collector, and Fig. 41 for the case of an elliptical collector. These examples show that nonuniform illuminance distribution at the first lens array is broken into regions. Each region is focused onto the second set of lenses, which creates an array of images of the source. The lenslets in the second array image the illuminance distribution at the first lens array onto the output plane. Enhanced uniformity is obtained because the distributions from each of the channels are superimposed at the target. The discrete images of the source that are seen at the second lens array indicate that the pupil of the output distribution is discontinuous.

Improvements to this basic idea include the use of lens arrays where the lenslets in the second lens array do not all have the same size/shape.²⁵¹ The curvature of each lenslet in the first array is decentered to “aim” the flux toward the proper lenslet in the second array. When nonuniform lenslet size/shapes are used, the etendue loss that results from the use of a conic reflector can be minimized. Such a system can be used to create asymmetric output distributions without a correspondingly misshapen pupil distribution (see earlier section on image dissectors).

Tailored Optics

The tailoring of a reflector to provide a desired distribution with a point or line source has been explored.^{8,259–265} Some simple cases are available in closed form.^{88,266,267} Refractive equivalents have also been investigated.^{268,269}

In a manner similar to the lens configurations shown in Fig. 36, tailored reflectors can produce converging or diverging wavefronts. Examples are shown in Fig. 42. The converging reflector creates a smeared image of the source between the reflector and the target. The diverging reflector creates a smeared image of the source outside of the region between the reflector and target. This terminology was developed assuming that the target is far from the reflector. The terms *compound hyperbolic* and *compound elliptic* are also used to describe these two types of reflector configurations.²⁷⁰

If the flux from the two sides of the reflector illuminates distinct sides of the target, then there are four main reflector classifications (e.g., Ref. 26, pp. 6.20–6.21; Refs. 91 and 260). Examples of the

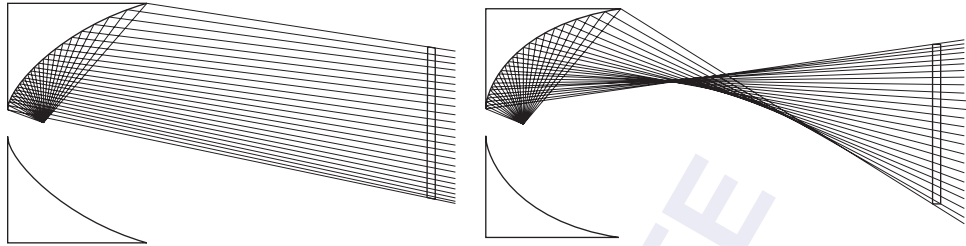


FIGURE 42 Two main classes of reflector types—diverging (*left*) and converging (*right*). The target is shown on the right of each of the two figures.

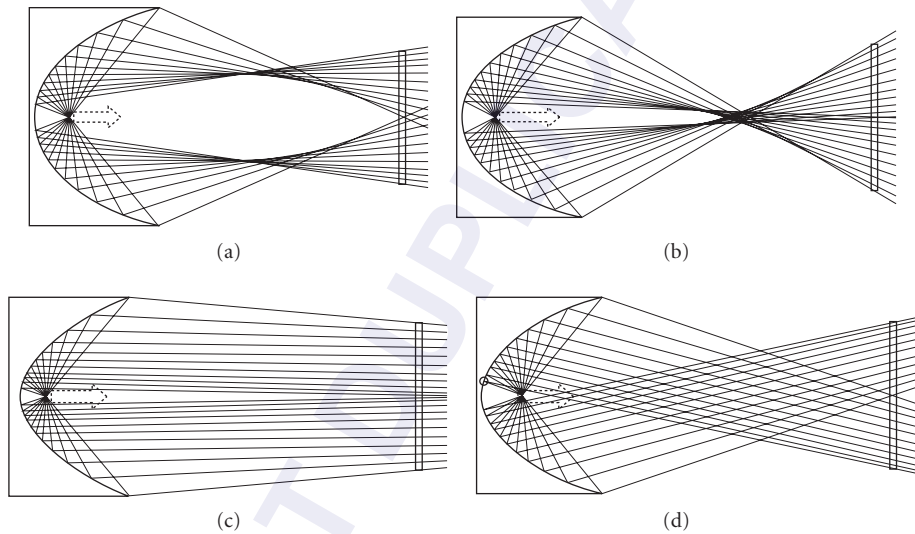


FIGURE 43 Four types of one-to-one mapping reflectors where the flux from each side of the reflector does not overlap at the target. (a) Uncrossed converging. (b) Crossed converging. (c) Uncrossed diverging. (d) Crossed diverging.

four cases are pictured in Fig. 43. They are the crossed and uncrossed versions of the converging and diverging reflectors, where *crossed* is relative to the center of the optical axis. Elmer²⁷¹ uses the terms *single* and *divided* for uncrossed and crossed, respectively. The uncrossed converging case tends to minimize issues with rays passing back through the source. The four cases shown provide one-to-one mappings because each point in the target receives flux from only one portion of the reflector.

Tailored Reflectors with Extended Sources Elmer⁸ describes some aspects of the analysis/design of illumination systems with extended sources using what are now called *edge rays* (see also Ref. 271). Winston²⁶⁶ sparked renewed interest in the subject. The following is a brief summary of some of the edge-ray papers that have been published in recent years. Some concentration related citations are included. Edge rays were also discussed earlier.

Gordon²⁷² adds a gap between the source and reflector to offset the cosine cubed effect and obtain a sharp cutoff. Gordon²⁷³ explores the case where the extreme direction is a linear function of polar angle. Friedman¹⁴⁷ uses a variable angle to design a secondary for a parabolic primary. Gordon²⁷⁴ describes a tailored edge ray secondary for a Fresnel primary, in particular for a heliostat field. Winston²⁷⁰ describes the tailoring of edge rays for a desired functional rather than maximal

concentration. Ries²⁷⁵ describes the tailoring based on a family of edge rays. Rabl²⁷⁶ tailors the reflector by establishing a one-to-one correspondence between target points and edge rays for planar sources. Ong^{277,278} explores tailored designs using full and partial involutes. Jenkins²⁷⁹ also describes a partial involute solution and generalizes the procedure as an integral design method.²⁸⁰ Ong⁹¹ explores the case of partial involutes with gaps. Gordon²⁸¹ relates the string construction to tailored edge ray designs for concentrators.

Faceted Structures

Nonuniform distributions can be made uniform through the use of reflectors where the reflector aims multiple portions of the flux from the source toward common locations at the target. In many cases, the resulting reflector has a faceted appearance. If the source etendue is significantly smaller than the etendue of the illumination distribution, then smearing near the edges of the uniform region can be small. However, if the source etendue is not negligible compared to the target, then faceted reflectors tend to experience smearing similar to the case of one lens array.

The term *faceted reflector* identifies reflectors composed of numerous distinct reflector regions. Elmer²⁷⁰ uses the term *multiphase* and Ong⁹¹ uses the term *hybrid* to describe these types of reflectors. If there are more than two regions, then the regions are often called *facets*. As with lens arrays, there are two extremes for faceted reflector designs. One is the beam-smearing category, where the facets remove substructure in the beam. The other extreme is the beam-forming category, where each facet creates the same distribution. If all facets create the same distribution, Elmer²⁷¹ calls the reflector *homogeneous*. If the distributions are different, then the term *inhomogeneous* is used.

For a given meridional slice, each facet can create an image of the source either in front of or behind the reflector. Using the convergent/divergent terminology to distinguish where the blurred image of the source occurs (see Fig. 42), a meridional slice of a faceted reflector can be composed of an array of convergent or divergent facets. These two cases are shown in Fig. 44 *a* and *b*. Flat facets are

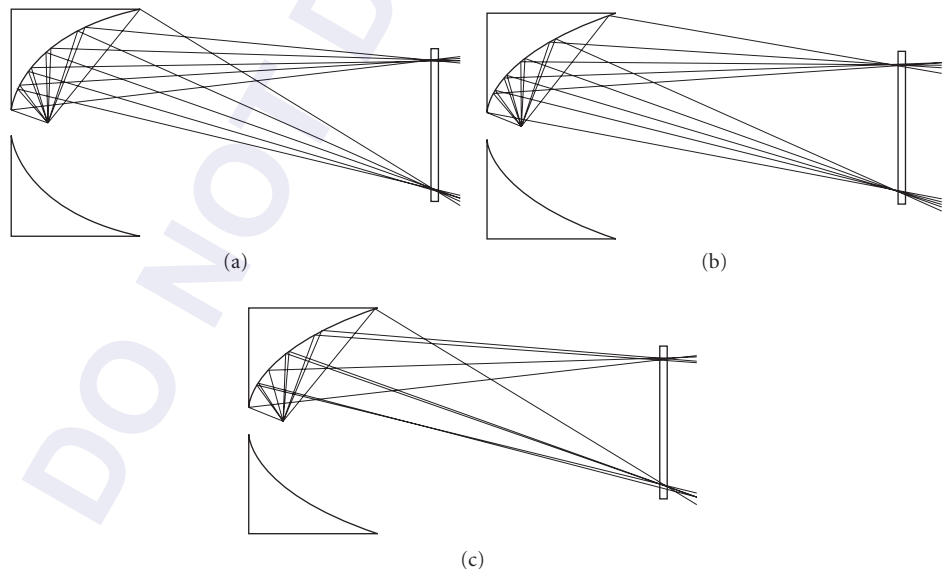


FIGURE 44 Some faceted reflector profiles. Convergent (*a*), divergent (*b*), and mixed convergent/divergent (*c*) profiles are all shown with five facets for the upper half of the reflector. Rays from the source that hit near the edges of the facets are shown. The rays for the upper facets are highlighted for all three cases. The rays for the second facet are highlighted in the mixed case.

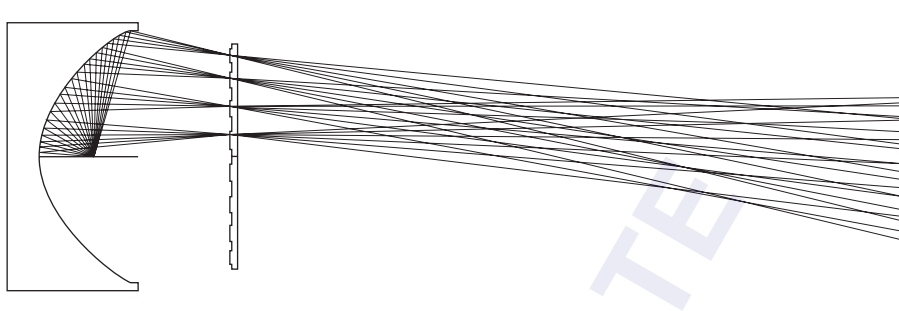


FIGURE 45 Reflector with array of concave facets and an array of lenslets. The lenslets image the facets.

often desired because of the simplicity of fabrication and fall under the divergent category because the image of the source is behind the reflector.

The intersection between facets of the same type introduces a discontinuity in the slope of the reflector. Mixed convergent/divergent versions are possible, as shown in Fig. 44c. Mixed versions offer the possibility of minimizing discontinuities in the slope of the reflector curve.

If all the distributions are superimposed, then improved uniformity can arise through superposition if sufficient averaging exists. Consider a manufacturing defect where a subregion of the reflector that is about the size of a facet is distorted. If a superposition approach is used, then the effect of the defect is small because each facet provides only a small contribution to target distribution. However, if a one-to-one mapping approach is used, then the same defect distorts an entire subregion of the target distribution because there is no built-in redundancy.

The effects of dilution should be considered when faceted optics are used. A single-lens array provides insufficient degrees of freedom to provide homogeneous superposition without introducing dilution, which is why a second lens array is often added (e.g., the tandem-lens-array configuration). Schering²⁸² shows lenslets placed on a reflector combined with a second array of lenslets. A straightforward extension of lenslets on a reflector is the use of concave facets, as pictured in Fig. 45.

Some Literature on Faceted Structures If the source is highly collimated (e.g., a laser) and has a Gaussian distribution, then only a couple of regions of the wavefront need to be superimposed to provide excellent uniformity. Descriptions of optical implementations that emphasize nearly Gaussian wavefronts include a four-surface prism by Kawamura,²⁸³ a wedge prism by Lacombat,²¹¹ and a multipass wedge prism by Sameshima.²⁸⁴ Austin²⁸⁵ has described a prism system for an excimer laser. Brunsting²⁸⁶ has described a system with nonplanar wedges. Bruno²⁸⁷ and Latta²²⁸ have discussed examples of mirrored systems with very few superimposed wavefronts.

For the more general case, uniformity can be obtained by superimposing a number of different regions of the beam. Dourte²⁸⁸ describes a 2D array of facets for laser homogenization. Ream²⁸⁹ discusses a system with a convex array of flat facets for a laser. Doherty²⁹⁰ describes a segmented conical mirror that can be used to irradiate circular, annular, and radial surfaces with a laser source. Dagenais²⁹¹ describes the details of a system with nonplanar facets for use with a laser in paint stripping. Dickey²⁹² shows a laser-based system composed of two orthogonal linear arrays of facets and includes curvature in the facet structures to minimize the effects of diffraction. Experimental data for two orthogonal linear arrays of facets [called a *strip mirror integrator* (SMI)] are provided by Geary.²⁹³ Lu²⁹⁴ has described a refractive structure using a multiwedge prism structure. Henning²⁹⁵ and Unnebrink²⁹⁶ show curved facets for use with UV laser systems.

Laser-based systems provide the degree of freedom that the wavefront divergence can be small, although the effects of diffraction and coherence must be considered in those designs. The divergence of nonlaser sources is often an important issue.

Jolley²⁵⁹ describes a faceted mirror design with disc and spherical sources. David²⁹⁷ describes faceted reflectors like the ones used in commercial lighting products and highlights the fact that facets on a

conic reflector base can improve uniformity with very little increase in “spot size” for long linear sources. Coaton²⁹⁸ provides a brief description of a faceted reflector design with a point source. Donohue²⁹⁹ describes faceted structures in headlamps.

There are many patents that describe the use of faceted structures for incoherent sources. Examples include Wiley,³⁰⁰ who uses the facets to remove filament image structure; Laudenschlager,³⁰¹ who describes design equations for a faceted reflector; Heimer,³⁰² who uses a faceted reflector in a lithography system; and Hamada,³⁰³ who uses a faceted structure with a linear lamp.

Ronneld³⁰⁴ uses a corrugated reflector to improve the angular acceptance of a 2D CPC in solar applications. Receiver uniformity can also be improved in CPCs using irregular surfaces (Ref. 4, pp. 153–161).

Faceted structures have often been used in complex lightpipes such as those used in the illumination of backlit liquid-crystal displays and instrument panels. For instrument-panel application, more recent efforts have been directed toward using aspheric surfaces.³⁰⁵

39.7 ACKNOWLEDGMENTS

Special acknowledgement is given to Doug Goodman, whose list of uniformity references and OSA course notes provided an important starting point for portions of this chapter. The support of Optical Research Associates (ORA) during the manuscript preparation is also appreciated, as are the editorial and terminology discussions with many of the technical staff at ORA. Simulation and ray trace results were all generated using either LightTools® or CodeV®.

39.8 REFERENCES

1. H. Hinterberger, L. Lavoie, B. Nelson, R. L. Sumner, J. M. Watson, R. Winston, and D. M. Wolfe, “The Design and Performance of a Gas Cerenkov Counter with Large Phase-Space Acceptance,” *Rev. Sci. Instrum.* **41**(3):413–418 (1970).
2. D. A. Harper, R. H. Hildebrand, R. Stiening, and R. Winston, “Heat Trap: An Optimized Far Infrared Field Optics System,” *Appl. Opt.* **15**:53–60 (1976) (Note: References are on page 144 of vol. 15).
3. R. L. Garwin, “The Design of Liquid Scintillation Cells,” *Rev. Sci. Instrum.* **23**:755–757 (1952).
4. R. Winston, J. C. Miñano, and P. Benítez, *Nonimaging Optics*, Elsevier Academic Press, San Diego, CA, 2005.
5. W. J. Cassarly and J. M. Davenport, “Fiber Optic Lighting: The Transition from Specialty Applications to Mainstream Lighting,” *SAE* 1999-01-0304 (1999).
6. A. E. Rosenbluth and R. N. Singh, “Projection Optics for Reflective Light Valves,” *Proc. SPIE* **3634**:87–111 (1999).
7. J. Bortz, N. Shatz, and R. Winston, “Advanced Nonrotationally Symmetric Reflector for Uniform Illumination of Rectangular Apertures,” *Proc. SPIE* **3781**:110–119 (1999).
8. W. Elmer, *The Optical Design of Reflectors*, 3rd ed., Academic Press, New York, 1989.
9. J. Palmer, Frequently Asked Questions, www.optics.arizona.edu/Palmer/rpfaq/rpfaq.htm, 1999.
10. J. M. Palmer, “Radiometry and Photometry: Units and Conversions,” in *OSA Handbook of Optics*, 3rd ed., vol. III, chap. 3, McGraw-Hill, New York, 2009.
11. A. Stimson, *Photometry and Radiometry for Engineers*, John Wiley & Sons, 1974.
12. F. Grum and R. J. Becherer, *Optical Radiation Measurements*, vol. 1, *Radiometry*, Academic Press, 1979.
13. W. Budde, *Optical Radiation Measurements*, vol. 4, *Physical Detectors of Optical Radiation*, Academic Press, 1983.
14. A. D. Ryer, *Light Measurement Handbook*, International Light, Inc., 1997.
15. W. L. Wolfe, *Introduction to Radiometry*, *Proc. SPIE Press*, 1998.
16. W. R. McCluney, *Introduction to Radiometry and Photometry*, Artech, 1994.
17. D. Goodman, “Geometrical Optics,” in *OSA Handbook of Optics*, 3rd ed., vol. I, chap. 1, McGraw-Hill, New York, 2009.

18. M. Born and E. Wolf, *Principles of Optics*, Cambridge Press, pp. 522–525, 1980.
19. M. S. Brennessoltz, “Light Collection Efficiency for Light Valve Projection Systems,” *Proc. SPIE* **2650**:71–79 (1996).
20. B. A. Jacobson, R. Winston, and P. L. Gleckman, “Flat-Panel Fluorescent Backlights with Reduced Illumination Angle: Backlighting Optics at the Thermodynamic Limit,” *SID 92 Digest*, 423–426 (1992).
21. H. Ries, “Thermodynamic Limitations of the Concentration of Electromagnetic Radiation,” *J. Opt. Soc. Am.* **72**(3):380–385 (1982).
22. H. Hottel, “Radiant Heat Transmission,” in *Heat Transmission*, W. H. McAdams, ed., 3rd ed., McGraw Hill, New York, 1954.
23. R. Winston, “Cone-Collectors for Finite Sources,” *Appl. Opt.* **17**(5):688–689 (1978).
24. G. H. Derrick, “A Three-Dimensional Analogue of the Hottel String Construction for Radiation Transfer,” *Opt. Acta* **32**:39–60 (1985).
25. C. Weiner, *Lehrbuch der darstellenden Geometrie*, vol. 1, Leipzig, 1884 (see F. E. Nicodemus, *Self Study Manual on Optical Radiation Measurements*, part 1, *Concepts*, NBS, Washington, March 1976).
26. M. S. Rea (ed.), *Lighting Handbook*, Illuminating Engineering Society of North America, 1993.
27. F. O. Bartell, “Projected Solid Angle and Black Body Simulators,” *Appl. Opt.* **28**(6):1055–1057 (March 1989).
28. D. G. Koch, “Simplified Irradiance/Illuminance Calculations in Optical Systems,” *Proc. SPIE* **1780**:226–242, 1992.
29. A. Rabl and R. Winston, “Ideal Concentrators for Finite Sources and Restricted Exit Angles,” *Appl. Optics* **15**:2880–2883 (1976).
30. E. Harting, D. R. Mills, and J. E. Giutronich, “Practical Concentrators Attaining Maximal Concentration,” *Opt. Lett.* **5**(1):32–34 (1980).
31. A. Luque, “Quasi-Optimum Pseudo-Lambertian Reflecting Concentrators: An Analysis,” *Appl. Opt.* **19**(14): 2398–2402 (1980).
32. R. M. Saraiji, R. G. Mistrick, and M. F. Modest, “Modeling Light Transfer through Optical Fibers for Illumination Applications,” *IES* 128–139 (Summer 1996).
33. T. L. Davenport, W. J. Cassarly, R. L. Hansler, T. E. Stenger, G. R. Allen, and R. F. Buelow, “Changes in Angular and Spatial Distribution Introduced into Fiber Optic Headlamp Systems by the Fiber Optic Cables,” *SAE*, Paper No. 981197, 1998.
34. S. Doyle and D. Corcoran, “Automated Mirror Design Using an Evolution Strategy,” *Opt. Eng.* **38**(2):323–333 (1999).
35. I. Ashdown, “Non-imaging Optics Design Using Genetic Algorithms,” *J. Illum. Eng. Soc.* **3**(1):12–21 (1994).
36. Y. Nakata, “Multi B-Spline Surface Reflector Optimized with Neural Network,” *SAE* **940638**:81–92 (1994).
37. N. E. Shatz and J. C. Bortz, “Inverse Engineering Perspective on Nonimaging Optical Design,” *Proc. SPIE* **2538**:136–156, 1995.
38. C. Gilray and I. Lewin, “Monte Carlo Techniques for the Design of Illumination Optics,” *IESNA Annual Conference Technical Papers*, Paper #85, pp. 65–80, July 1996.
39. J. C. Schweyen, K. Garcia, and P. L. Gleckman, “Geometrical Optical Modeling Considerations for LCD Projector Display Systems,” *Proc. SPIE* **3013**:126–140 (1997).
40. D. Z. Ting and T. C. McGill, “Monte Carlo Simulation of Light-Emitting Diode Light-Extraction Characteristics,” *Opt. Eng.* **34**(12):3545–3553 (1995).
41. M. Kaplan, “Monte Carlo Calculation of Light Distribution in an Integrating Cavity Illuminator,” *Proc. SPIE* **1448**:206–217 (1991).
42. B. G. Crowther, “Computer Modeling of Integrating Spheres,” *Appl. Opt.* **35**(30):5880–5886 (1996).
43. R. C. Chaney, “Monte Carlo Simulation of Gamma Ray Detectors using Scintillation Fibers,” *EUV, X-ray, and Gamma-Ray Instrumentation for Astronomy and Atomic Physics*; Proceedings of the Meeting, San Diego, California, Aug. 7–11, 1989 (A90-50251 23–35), SPIE 1989.
44. J. A. Bamberg, “Scintillation Detector Optimization Using GUERAP-3 Radiation Scattering in Optical Systems, Proceedings of the Seminar, Huntsville, Alabama, September 30–October 1, 1980 (A81-36878 16-74) p. 86–93, SPIE 1981.
45. B. K. Likeness, “Stray Light Simulation with Advanced Monte Carlo Techniques,” *Stray-Light Problems in Optical Systems*; Proceedings of the Seminar, Reston, VA, April 18–21, 1977 (A78-40270 17-35) pp. 80–88, SPIE 1977.

46. E. R. Freniere, "Simulation of Stray Light in Optical Systems with the GUERAP III," *Radiation Scattering in Optical Systems*; Proceedings of the Seminar, Huntsville, AL, September 30–October 1, 1980 (A81-36878 16-74) pp. 78–85, SPIE 1981.
47. N. Shatz, J. Bortz, and M. Dassanayake, "Design Optimization of a Smooth Headlamp Reflector to SAE/DOT Beam-Shape Requirements," SAE 1999.
48. T. Hough, J. F. Van Derlofske, and L. W. Hillman, "Measuring and Modeling Intensity Distributions of Light Sources in Waveguide Illumination Systems," *Opt. Eng.* **34**(3):819–823 (1995).
49. R. E. Levin, "Photometric Characteristics of Light-Controlling Apparatus," *Illum. Eng.* **66**(4):202–215 (1971).
50. R. J. Donohue and B. W. Joseph, "Computer Synthesized Filament Images from Reflectors and through Lens Elements for Lamp Design and Evaluation," *Appl. Opt.* **14**(10):2384–2390 (1975).
51. I. Ashdown, "Near Field Photometry: A New Approach," *J. IES* 163–180 (Winter 1993).
52. R. D. Stock and M. W. Siegel, "Orientation Invariant Light Source Parameters," *Opt. Eng.* **35**(9):2651–2660 (1996).
53. M. W. Siegel and R. D. Stock, "Generalized Near-Zone Light Source Model and Its Application to Computer Automated Reflector Design," *Opt. Eng.* **35**(9):2661–2679 (1996).
54. P. V. Shmelev and B. M. Khana, "Near-Field Modeling versus Ray Modeling of Extended Halogen Light Source in Computer Design of a Reflector," *Proc. SPIE* (July 27–28, 1997).
55. R. Rykowski and C. B. Wooley, "Source Modeling for Illumination Design," 3130B-27, *Proc. SPIE* (July 27–28, 1997).
56. T. P. Vogl, L. C. Lintner, R. J. Pegis, W. M. Waldbauer, and H. A. Unvala, "Semiautomatic Design of Illuminating Systems," *Appl. Opt.* **11**(5):1087–1090 (1972).
57. "1998 IESNA Software Survey," *Lighting Design and Application* **28**(10):53–62 (1998).
58. "1999 IESNA Software Survey," *Lighting Design and Application* **29**(12):39–48 (1999).
59. W. T. Welford, *Aberrations of Optical Systems*, Adam Hilger, Bristol, 1986, pp. 158–161.
60. W. J. Smith, *Modern Lens Design*, McGraw Hill, New York, 1992.
61. Y. Shimizu and H. Takenaka, "Microscope Objective Design," in *Advances in Optical and Electron Microscopy*, vol. 14, Academic Press, San Diego, CA, 1994.
62. G. P. Smestad, "Nonimaging Optics of Light-Emitting Diodes: Theory and Practice," *Proc. SPIE* **1727**: 264–268 (1992).
63. W. T. Welford and R. Winston, "Two-Dimensional Nonimaging Concentrators with Refracting Optics," *J. Opt. Soc. Am.* **69**(6):917–919 (1979).
64. R. Kingslake, *Lens Design Fundamentals*, Academic Press, San Diego, CA, 1978.
65. E. Hecht and A. Zajac, *Optics*, Addison-Wesley, 1979.
66. S. F. Ray, *Applied Photographic Optics*, Focal Press, 1988.
67. G. Schultz, "Achromatic and Sharp Real Image of a Point by a Single Aspheric Lens," *Appl. Opt.* **22**(20): 3242–3248 (1983).
68. J. C. Minano and J. C. Gonzalez, "New Method of Design of Nonimaging Concentrators," *Appl. Opt.* **31**(16): 3051–3060, 1992.
69. E. A. Boettner and N. E. Barnett, "Design and Construction of Fresnel Optics for Photoelectric Receivers" *J. Opt. Soc. Am.* **41**(11):849–857 (1951).
70. Fresnel Technologies, Fort Worth, TX. See <http://www.fresneltech.com/pdf/FresnelLenses.pdf>
71. S. Sinzinger and M. Testorf, "Transition Between Diffractive and Refractive Microoptical Components," *Appl. Opt.* **34**(26):5670–5676 (1995).
72. F. Erismann, "Design of Plastic Aspheric Fresnel Lens with a Spherical Shape," *Opt. Eng.* **36**(4):988–991 (1997).
73. D. J. Lamb and L. W. Hillman, "Computer Modeling and Analysis of Veiling Glare and Stray Light in Fresnel Lens Optical Systems," *Proc. SPIE* **3779**:344–352 (1999).
74. J. F. Goldenberg and T. S. McKechnie, "Optimum Riser Angle for Fresnel Lenses in Projection Screens," U.S. Patent 4,824,227, 1989.
75. M. Collares-Pereira, "High Temperature Solar Collector with Optimal Concentration: Nonfocusing Fresnel Lens with Secondary Concentrator," *Solar Energy* **23**:409–419 (1979).

76. W. A. Parkyn and D. G. Pelka, "Compact Nonimaging Lens with Totally Internally Reflecting Facets," *Proc. SPIE* **1528**:70–81 (1991).
77. R. I. Nagel, "Signal Lantern Lens," U.S. Patent 3,253,276, 1966.
78. W. A. Parkyn, P. L. Gleckman, and D. G. Pelka, "Converging TIR Lens for Nonimaging Concentration of Light from Compact Incoherent Sources," *Proc. SPIE* **2016**:78–86 (1993).
79. W. A. Parkyn and D. G. Pelka, "TIR Lenses for Fluorescent Lamps," *Proc. SPIE* **2538**:93–103 (1995).
80. W. A. Parkyn and D. G. Pelka, "New TIR Lens Applications for Light-Emitting Diodes," *Proc. SPIE* **3139**:135–140 (1997).
81. V. Medvedev, W. A. Parkyn, and D. G. Pelka, "Uniform High-Efficiency Condenser for Projection Systems," *Proc. SPIE* **3139**:122–134 (1997).
82. V. Medvedev, D. G. Pelka, and W. A. Parkyn, "Uniform LED Illuminator for Miniature Displays," *Proc. SPIE* **3428**:142–153 (1998).
83. D. F. Vanderwerf, "Achromatic Catadioptric Fresnel Lens," *Proc. SPIE* **2000**:174–183 (1993).
84. J. Spigulis, "Compact Illuminators, Collimators and Focusers with Half Spherical Input Aperture," *Proc. SPIE* **2065**:54–59 (1994).
85. D. Silvergate, "Collimating Compound Catadioptric Immersion Lens," U.S. Patent 4,770,514, 1988.
86. McDermott, "Angled Elliptical Axial Lighting Device," U.S. Patent 5,894,196, 1999.
87. D. Korsch, *Reflective Optics*, Academic Press, San Diego, CA, 1991.
88. D. E. Spencer, L. L. Montgomery, and J. F. Fitzgerald, "Macrofocal Conies as Reflector Contours," *J. Opt. Soc. Am.* **55**(1):5–11 (1965).
89. D. R. Philips, "Low Brightness Louver," U.S. Patent 2,971,083, 1961.
90. H. P. Baum and J. M. Gordon, "Geometric Characteristics of Ideal Nonimaging (CPC) Solar Collectors with Cylindrical Absorber," *Solar Energy* **33**(5):455–158 (1984).
91. P. T. Ong, J. M. Gordon, and A. Rabl, "Tailored Edge-Ray Designs for Illumination with Tubular Sources," *Appl. Opt.* **35**(22):4361–4371 (1996).
92. I. M. Bassett and G. H. Derrick, "The Collection of Diffuse Light onto an Extended Absorber," *Optical Quantum Electronics* **10**:61–82 (1978).
93. H. Ries, N. Shatz, J. Bortz, and W. Spirkl, "Performance Limitations of Rotationally Symmetric Nonimaging Devices," *J. Opt. Soc. Am. A* **14**(10):2855–2862 (1997).
94. M. Ruda (ed.), "International Conference on Nonimaging Concentrators," *Proc. SPIE* **441** (1983).
95. R. Winston and R. L. Holman (eds.), "Nonimaging Optics: Maximum Efficiency Light Transfer," *Proc. SPIE* **1528** (1991).
96. R. Winston and R. L. Holman (eds.), "Nonimaging Optics: Maximum Efficiency Light Transfer II," *Proc. SPIE* **2016** (1993).
97. R. Winston (ed.), "Nonimaging Optics: Maximum Efficiency Light Transfer III," *Proc. SPIE* **2538** (1995).
98. R. Winston (ed.), "Nonimaging Optics: Maximum Efficiency Light Transfer IV," *Proc. SPIE* **3139** (1997).
99. R. Winston (ed.), "Nonimaging Optics: Maximum Efficiency Light Transfer V," *Proc. SPIE* **3781** (1999).
100. R. Winston (ed.), "Selected Papers on Nonimaging Optics," *Proc. SPIE Milestone Series* **106** (1995).
101. I. M. Bassett, W. T. Welford, and R. Winston, "Nonimaging Optics for Flux Concentration," in *Progress in Optics*, E. Wolf (ed.), 1989, pp. 161–226.
102. R. Winston, "Nonimaging Optics," *Sci. Am.* 76–81 (March 1991).
103. P. Gleckman, J. O'Gallagher, and R. Winston, "Approaching the Irradiance of the Sun Through Nonimaging Optics," *Optics News*: 33–36 (May 1989).
104. M. F. Land, "The Optical Mechanism of the Eye of Limulus," *Nature* **280**:396–397 (1979).
105. M. F. Land, "Compound Eyes: Old and New Optical Mechanisms," *Nature* **287**:681–685 (1980).
106. R. Levi-Seti, D. A. Park, and R. Winston, "The Corneal Cones of Limulus as Optimized Light Concentrators," *Nature* **253**:115–116 (1975).
107. D. A. Baylor and R. Fettiplace, "Light Path and Photon Capture in Turtle Photoreceptors," *J. Physiol.* **248**(2): 433–464 (1975).
108. R. Winston and J. Enoch, "Retinal Cone Receptor as Ideal Light Collector," *J. Opt. Soc. Am.* **61**(8):1120–1121 (1971).

109. R. L. Garwin, "The Design of Liquid Scintillation Cells," *Rev. Sci. Instr.* **23**:755–757 (1952).
110. D. E. Williamson, "Cone Channel Condensor," *J. Opt. Soc. Am.* **42**(10):712–715 (1952).
111. J. H. Myer, "Collimated Radiation in Conical Light Guides," *Appl. Opt.* **19**(18):3121–3123 (1980).
112. W. Witte, "Cone Channel Optics," *Infrared Phys.* **5**:179–185 (1965).
113. C. H. Burton, "Cone Channel Optics," *Infrared Phys.* **15**:157–159 (1975).
114. R. Winston and W. T. Welford, "Ideal Flux Concentrators as Shapes That Do Not Disturb the Geometrical Vector Flux Field: A New Derivation of the Compound Parabolic Concentrator," *J. Opt. Soc. Am.* **69**(4):536–539 (1979).
115. A. G. Molledo and A. Luque, "Analysis of Static and Quasi-Static Cross Compound Parabolic Concentrators," *Appl. Opt.* **23**:2007–2020 (1984).
116. R. Winston and H. Hinterberger, "Principles of Cylindrical Concentrators for Solar Energy," *Solar Energy* **17**:255–258 (1975).
117. W. R. McIntire, "Truncation of Nonimaging Cusp Concentrators," *Solar Energy* **23**:351–355 (1979).
118. H. Tabor, "Comment—The CPC Concept—Theory and Practice," *Solar Energy* **33**(6):629–630 (1984).
119. A. Rabl and R. Winston, "Ideal Concentrators for Finite Sources and Restricted Exit Angles," *Appl. Opt.* **15**:2880–2883 (1976).
120. A. Rabl, N. B. Goodman, and R. Winston, "Practical Design Considerations for CPC Solar Collectors," *Solar Energy* **22**:373–381 (1979).
121. W. R. McIntire, "New Reflector Design Which Avoids Losses through Gaps between Tubular Absorber and Reflectors," *Solar Energy* **25**:215–220 (1980).
122. R. Winston, "Ideal Flux Concentrators with Reflector Gaps," *Appl. Opt.* **17**(11):1668–1669 (1978).
123. R. Winston, "Cavity Enhancement by Controlled Directional Scattering," *Appl. Opt.* **19**:195–197 (1980).
124. F. Bloisi, P. Cavaliere, S. De Nicola, S. Martellucci, J. Quartieri, and L. Vicari, "Ideal Nonfocusing Concentrator with Fin Absorbers in Dielectric Rhombuses," *Opt. Lett.* **12**(7):453–155 (1987).
125. I. R. Edmonds, "Prism-Coupled Compound Parabola: A New Look and Optimal Solar Concentrator," *Opt. Lett.* **11**(8):490–492 (1986).
126. J. D. Kuppenheimer, "Design of Multilamp Nonimaging Laser Pump Cavities," *Opt. Eng.* **27**(12):1067–1071 (1988).
127. D. Lowe, T. Chin, T. L. Credelle, O. Tezucar, N. Hariston, J. Wilson, and K. Bingaman, "SpectraVue: A New System to Enhance Viewing Angle of LCDs," *SID 96 Applications Digest* 39–42 (1996).
128. J. L. Henkes, "Light Source for Liquid Crystal Display Panels Utilizing Internally Reflecting Light Pipes and Integrating Sphere," U.S. Patent 4,735,495, 1988.
129. R. Winston, "Principles of Solar Concentrators of a Novel Design," *Solar Energy* **16**:89–94 (1974).
130. R. Winston, "Cone-Collectors for Finite Sources," *Appl. Opt.* **17**(5):688–689 (1978).
131. M. Collares-Pereira, A. Rabl, and R. Winston, "Lens-Mirror Combinations with Maximal Concentration," *Appl. Opt.* **16**(10):2677–2683 (October 1977).
132. H. P. Gush, "Hyberbolic Cone-Channel Condensor," *Opt. Lett.* **2**:22–24 (1978).
133. J. O'Gallagher, R. Winston, and W. T. Welford, "Axially Symmetric Nonimaging Flux Concentrators with the Maximum Theoretical Concentration Ratio," *J. Opt. Soc. Am. A* **4**(1):66–68 (1987).
134. R. Winston, "Dielectric Compound Parabolic Concentrators," *Appl. Opt.* **15**(2):291–292 (1976).
135. J. R. Hull, "Dielectric Compound Parabolic Concentrating Solar Collector with a Frustrated Total Internal Reflection Absorber," *Appl. Opt.* **28**(1):157–162 (1989).
136. R. Winston, "Light Collection within the Framework of Geometrical Optics," *J. Opt. Soc. Am.* **60**(2):245–247 (1970).
137. D. Jenkins, R. Winston, R. Bliss, J. O'Gallagher, A. Lewandowski, and C. Bingham, "Solar Concentration of 50,000 Achieved with Output Power Approaching 1kW," *J. Sol. Eng.* **118**:141–144 (1996).
138. H. Ries, A. Segal, and J. Karni, "Extracting Concentrated Guided Light," *Appl. Opt.* **36**(13):2869–2874 (1997).
139. R. H. Hildebrand, "Focal Plane Optics in Far-Infrared and Submillimeter Astronomy," *Opt. Eng.* **25**(2):323–330 (1986).
140. J. Keene, R. H. Hildebrand, S. E. Whitcomb, and R. Winston, "Compact Infrared Heat Trap Field Optics," *Appl. Opt.* **17**(7):1107–1109 (1978).

141. R. Winston and W. T. Welford, "Geometrical Vector Flux and Some New Nonimaging Concentrators," *J. Opt. Soc. Am.* **69**(4):532–536 (1979).
142. X. Ning, R. Winston, and J. O'Gallagher, "Dielectric Totally Internally Reflecting Concentrators," *Appl. Opt.* **26**(2):300–305 (1987).
143. W. L. Eichhorn, "Designing Generalized Conic Concentrators for Conventional Optical Systems," *Appl. Opt.* **24**(8):1204–1205 (1985).
144. W. L. Eichhorn, "Generalized Conic Concentrators," *Appl. Opt.* **21**(21):3887–3890 (1982).
145. G. H. Smith, *Practical Computer-Aided Lens Design*, Willmann-Bell, VA, 1998, pp. 379–383.
146. K. W. Beeson, I. B. Steiner, and S. M. Zimmerman, "Illumination System Employing an Array of Microprisms," U.S. Patent 5,521,725, 1996.
147. R. P. Friedman, J. M. Gordon, and H. Ries, "New High-Flux Two-Stage Optical Designs for Parabolic Solar Concentrators," *Solar Energy* **51**:317–325 (1993).
148. P. Gleckman, J. O'Gallagher, and R. Winston, "Concentration of Sunlight to Solar-Surface Levels Using Non-imaging Optics," *Nature* **339**:198–200 (May 1989).
149. H. Ries and W. Sprickl, "Nonimaging Secondary Concentrators for Large Rim Angle Parabolic Trough with Tubular Absorbers," *Appl. Opt.* **35**:2242–2245 (1996).
150. J. O'Gallagher and R. Winston, "Test of a 'Trumpet' Secondary Concentrator with a Paraboloidal Dish Primary," *Solar Energy* **36**(1):37–14 (1986).
151. A. Rabl, "Comparison of Solar Concentrators," *Solar Energy* **18**:93–111 (1976).
152. R. Winston and W. T. Welford, "Design of Nonimaging Concentrators as Second Stages in Tandem with Image-Forming First-Stage Concentrators," *Appl. Opt.* **19**(3):347–351 (1980).
153. X. Ning, R. Winston, and J. O'Gallagher, "Optics of Two-Stage Photovoltaic Concentrators with Dielectric Second Stages," *Appl. Opt.* **26**(7):1207–1212 (1987).
154. E. M. Kritchman, "Second-Stage Concentrators—A New Formalism," *J. Opt. Soc. Am. Lett.* **73**(4):508–511 (1983).
155. D. R. Mills and J. E. Giutronich, "New Ideal Concentrators for Distant Radiation Sources," *Solar Energy* **23**:85–87 (1979).
156. D. R. Mills and J. E. Giutronich, "Asymmetrical Non-imaging Cylindrical Solar Concentrators," *Solar Energy* **20**:45–55 (1978).
157. H. Ries and J. M. Gordon, "Double-Tailored Imaging Concentrators," *Proc. SPIE* **3781**:129–134 (1999).
158. J. C. Minano and J. C. Gonzalez, "New Method of Design of Nonimaging Concentrators," *Appl. Opt.* **31**(16):3051–3060 (1992).
159. J. C. Minano, P. Benitez, and J. C. Gonzalez, "RX: A Nonimaging Concentrator," *Appl. Opt.* **34**(13):2226–2235 (1995).
160. P. Benitez and J. C. Minano, "Ultrahigh-Numerical-Aperture Imaging Concentrator," *J. Opt. Soc. Am. A* **14**(8):1988–1997 (1997).
161. J. C. Minano, J. C. Gonzalez, and P. Benitez, "A High-Gain, Compact, Nonimaging Concentrator: RXI," *Appl. Opt.* **34**(34):7850–7856 (1995).
162. J. L. Alvarez, M. Hernandez, P. Benitez, and J. C. Minano, "RXI Concentrator for 1000× Photovoltaic Conversion," *Proc. SPIE* **3781**:30–37 (1999).
163. J. F. Forkner, "Aberration Effects in Illumination Beams Focused by Lens Systems," *Proc. SPIE* **3428**:73–89 (1998).
164. X. Ning, "Three-Dimensional Ideal Θ_1/Θ_2 Angular Transformer and Its Uses in Fiber Optics," *Appl. Opt.* **27**(19):4126–4130 (1988).
165. A. Timinger, A. Kribus, P. Doron, and H. Ries, "Optimized CPC-type Concentrators Built of Plane Facets," *Proc. SPIE* **3781**:60–67 (1999).
166. J. P. Rice, Y. Zong, and D. J. Dummer, "Spatial Uniformity of Two Nonimaging Concentrators," *Opt. Eng.* **36**(11):2943–2947 (1997).
167. R. M. Emmons, B. A. Jacobson, R. D. Gengelbach, and R. Winston, "Nonimaging Optics in Direct View Applications," *Proc. SPIE* **2538**:42–50 (1995).
168. D. B. Leviton and J. W. Leitch, "Experimental and Raytrace Results for Throat-to-Throat Compound Parabolic Concentrators," *Appl. Opt.* **25**(16):2821–2825 (1986).

169. B. Moslehi, J. Ng, I. Kasimoff, and T. Jansson, "Fiber-Optic Coupling Based on Nonimaging Expanded-Beam Optics," *Opt. Lett.* **14**(23):1327–1329 (1989).
170. M. Collares-Pereira, J. F. Mendes, A. Rabl, and H. Ries, "Redirecting Concentrated Radiation," *Proc. SPIE* **2538**: 131–135 (1995).
171. A. Rabl, "Solar Concentrators with Maximal Concentration for Cylindrical Absorbers," *Appl. Opt.* **15**(7): 1871–1873 (1976). See also an erratum, *Appl. Opt.* **16**(1):15 (1977).
172. D. Feuermann and J. M. Gordon, "Optical Performance of Axisymmetric Concentrators and Illuminators," *Appl. Opt.* **37**(10):1905–1912 (1998).
173. J. Bortz, N. Shatz, and H. Ries, "Consequences of Etendue and Skewness Conservation for Nonimaging Devices with Inhomogeneous Targets," *Proc. SPIE* **3139**:28 (1997).
174. N. E. Shatz, J. C. Bortz, H. Ries, and R. Winston, "Nonrotationally Symmetric Nonimaging Systems that Overcome the Flux-Transfer Performance Limit Imposed by Skewness Conservation," *Proc. SPIE* **3139**:76–85 (1997).
175. N. E. Shatz, J. C. Bortz, and R. Winston, "Nonrotationally Symmetric Reflectors for Efficient and Uniform Illumination of Rectangular Apertures," *Proc. SPIE* **3428**:176–183 (1998).
176. W. Benesch and J. Strong, "The Optical Image Transformer," *J. Opt. Soc. Am.* **41**(4):252–254 (1951).
177. I. S. Bowen, "The Image-Slicer, a Device for Reducing Loss of Light at Slit of Stellar Spectrograph," *Astrophysical J.* **88**(2):113–124 (1938).
178. W. D. Westwood, "Multiple Lens System for an Optical Imaging Device," U.S. Patent 4,114,037, 1978.
179. J. Bernges, L. Unnebrink, T. Henning, E. W. Kreutz, and R. Poprawe, "Novel Design Concepts for UV Laser Beam Shaping," *Proc. SPIE* **3779**:118–125 (1999).
180. J. Endriz, "Brightness Conserving Optical System for Modifying Beam Symmetry," U.S. Patent 5,168,401, 1992.
181. T. Mori and H. Komatsuda, "Optical Integrator and Projection Exposure Apparatus Using the Same," U.S. Patent 5,594,526, 1997.
182. J. A. Shimizu and P. J. Janssen, "Integrating Lens Array and Image Forming Method for Improved Optical Efficiency," U.S. Patent 5,662,401, 1997.
183. D. Feuermann, J. M. Gordon, and H. Ries, "Nonimaging Optical Designs for Maximum-Power-Density Remote Irradiation," *Appl. Opt.* **37**(10):1835–1844 (1998).
184. F. C. Genovese, "Fiber Optic Line Illuminator with Deformed End Fibers and Method of Making Same," U.S. Patent 4,952,022, 1990.
185. P. Gorenstein and D. Luckey, "Light Pipe for Large Area Scintillator," *Rev. Sci. Instrum.* **34**(2):196–197 (1963).
186. W. Gibson, "Curled Light Pipes for Thin Organic Scintillators," *Rev. Sci. Instrum.* **35**(8):1021–1023 (1964).
187. H. Hinterberger and R. Winston, "Efficient Design of Lucite Light Pipes Coupled to Photomultipliers," *Rev. Sci. Instrum.* **39**(3):419–420 (1968).
188. H. D. Wolpert, "A Light Pipe for Flux Collection: An Efficient Device for Document Scanning," *Lasers Appl.* 73–74 (April 1983).
189. H. Karasawa, "Light Guide Apparatus Formed from Strip Light Guides," U.S. Patent 4,824,194, 1989.
190. D. A. Markle, "Optical Transformer Using Curved Strip Waveguides to Achieve a Nearly Unchanged F/Number," U.S. Patent 4,530,565, 1985.
191. I. M. Bassett and G. W. Forbes, "A New Class of Ideal Non-imaging Transformers," *Opt. Acta* **29**(9):1271–1282 (1982).
192. G. W. Forbes and I. M. Bassett, "An Axially Symmetric Variable-Angle Nonimaging Transformer," *Opt. Acta* **29**(9):1283–1297 (1982).
193. M. E. Barnett, "The Geometric Vector Flux Field within a Compound Elliptical Concentrator," *Optik* **54**(5):429–432 (1979).
194. M. E. Barnett, "Optical Flow in an Ideal Light Collector: The Θ_1/Θ_2 Concentrator," *Optik* **57**(3):391–400 (1980).
195. Gutierrez, J. C. Minano, C. Vega, and P. Benitez, "Application of Lorentz Geometry to Nonimaging Optics: New 3D Ideal Concentrators," *J. Opt. Soc. Am. A* **13**(3):532–540 (1996).
196. P. Greenman, "Geometrical Vector Flux Sinks and Ideal Flux Concentrators," *J. Opt. Soc. Am. Lett.* **71**(6):777–779 (1981).

197. W. T. Welford and R. Winston, "On the Problem of Ideal Flux Concentrators," *J. Opt. Soc. Am.* **68**(4):531–534 (1978). See also addendum, *J. Opt. Soc. Am.* **69**(2):367 (1979).
198. J. C. Minano, "Design of Three-Dimensional Nonimaging Concentrators with Inhomogeneous Media," *J. Opt. Soc. Am. A* **3**(9):1345–1353 (1986).
199. J. M. Gordon, "Complementary Construction of Ideal Nonimaging Concentrators and Its Applications," *Appl. Opt.* **35**(28):5677–5682 (1996).
200. P. A. Davies, "Edge-Ray Principle of Nonimaging Optics," *J. Opt. Soc. Am. A* **11**:1256–1259 (1994).
201. H. Ries and A. Rabl, "Edge-Ray Principle of Nonimaging Optics," *J. Opt. Soc. Am.* **10**(10):2627–2632 (1994).
202. A. Rabl, "Edge-Ray Method for Analysis of Radiation Transfer among Specular Reflectors," *Appl. Opt.* **33**(7):1248–1259 (1994).
203. J. C. Minano, "Two-Dimensional Nonimaging Concentrators with Inhomogeneous Media: A New Look," *J. Opt. Soc. Am. A* **2**(11):1826–1831 (1985).
204. M. T. Jones, "Motion Picture Screen Light as a Function of Carbon-Arc-Crater Brightness Distribution," *J. Soc. Mot. Pic. Eng.* **49**:218–240 (1947).
205. R. Kingslake, *Applied Optics and Optical Engineering*, vol. II, Academic Press, New York, 1965, pp. 225–226.
206. D. O'Shea, *Elements of Modern Optical Design*, Wiley, New York, 1985, pp. 111–114 and 384–390.
207. W. Wallin, "Design of Special Projector Illuminating Systems," *J-SMPTE* **71**:769–771 (1962).
208. H. Weiss, "Wide-Angle Slide Projection," *Inf. Disp.* 8–15 (September/October 1964).
209. S. Bradbury, *An Introduction to the Optical Microscope*, Oxford University Press, 1989, pp. 23–27.
210. R. Oldenbourg and M. Shribak, "Microscopes," in *OSA Handbook of Optics*, 3rd ed., vol. I, chap. 28, McGraw-Hill, New York, 2009.
211. M. Lacombat, G. M. Dubroeuq, J. Massin, and M. Brevignon, "Laser Projection Printing," *Solid State Technology* **23**(115): (1980).
212. J. M. Gordon and P. T. Ong, "Compact High-Efficiency Nonimaging Back Reflector for Filament Light Sources," *Opt. Eng.* **35**(6):1775–1778 (1996).
213. G. Zochling, "Design and Analysis of Illumination Systems," *Proc. SPIE* **1354**:617–626 (1990).
214. A. Steinfeld, "Apparent Absorptance for Diffusely and Specularly Reflecting Spherical Cavities," *Int. J. Heat Mass Trans.* **34**(7):1895–1897 (1991).
215. K. A. Snail and L. M. Hanssen, "Integrating Sphere Designs with Isotropic Throughput," *Appl. Opt.* **28**(10):1793–1799 (1989).
216. D. P. Ramer and J. C. Rains, "Lambertian Surface to Tailor the Distribution of Light in a Nonimaging Optical System," *Proc SPIE* **3781** (1999).
217. D. G. Goebel, "Generalized Integrating Sphere Theory," *Appl. Opt.* **6**(1):125–128 (1967).
218. L. Hanssen and K. Snail, "Nonimaging Optics and the Measurement of Diffuse Reflectance," *Proc. SPIE* **1528**:142–150 (1991).
219. D. B. Chenault, K. A. Snail, and L. M. Hanssen, "Improved Integrating-Sphere Throughput with a Lens and Nonimaging Concentrator," *Appl. Opt.* **34**(34):7959–7964 (1995).
220. R. H. Webb, "Concentrator for Laser Light," *Appl. Opt.* **31**(28):5917–5918 (1992).
221. D. S. Goodman, "Producing Uniform Illumination," OSA Annual Meeting, Toronto, October 5, 1993.
222. M. M. Chen, J. B. Berkowitz-Mattuck, and P. E. Glaser, "The Use of a Kaleidoscope to Obtain Uniform Flux over a Large Area in a Solar or Arc Imaging Furnace," *Appl. Opt.* **2**:265–271 (1963).
223. L. A. Whitehead, R. A. Nodwell, and F. L. Curzon, "New Efficiency Light Guide for Interior Illumination," *Appl. Opt.* **21**(15):2755–2757 (1982).
224. S. G. Saxe, L. A. Whitehead, and S. Cobb, "Progress in the Development of Prism Light Guides," *Proc. SPIE* **692**:235–240 (1986).
225. K. Jain, "Illumination System to Produce Self-Luminous Light Beam of Selected Cross-Section, Uniform Intensity and Selected Numerical Aperture," U.S. Patent 5,059,013, 1991.
226. Y. Kudo and K. Matsumoto, "Illuminating Optical Device," U.S. Patent 4,918,583, 1990.
227. W. J. Cassarly, J. M. Davenport, and R. L. Hansler, "Uniform Lighting Systems: Uniform Light Delivery," *SAE 950904*, **SP-1081**:1–5 (1995).

228. M. R. Latta and K. Jain, "Beam Intensity Uniformization by Mirror Folding," *Opt. Comm.* **49**:27 (1984).
229. C. M. Chang, K. W. Lin, K. V. Chen, S. M. Chen, and H. D. Shieh, "A Uniform Rectangular Illuminating Optical System for Liquid Crystal Light Valve Projectors," *Euro Display* '96:257–260 (1996).
230. L. J. Coyne, "Distributive Fiber Optic Couplers Using Rectangular Lightguides as Mixing Elements," *Proc. FOC '79*:160–164 (1979).
231. M. Wagner, H. D. Geiler, and D. Wolff, "High-Performance Laser Beam Shaping and Homogenization System for Semiconductor Processing," *Meas. Sci. Technol. UK* **1**:1193–1201 (1990).
232. K. Iwasaki, Y. Ohyama, and Y. Nanaumi, "Flattening Laserbeam Intensity Distribution," *Lasers Appl.* **76** (April 1983).
233. K. Iwasaki, T. Hayashi, T. Goto, and S. Shimizu, "Square and Uniform Laser Output Device: Theory and Applications," *Appl. Opt.* **29**:1736–1744 (1990).
234. J. M. Geary, "Channel Integrator for Laser Beam Uniformity on Target," *Opt. Eng.* **27**:972–977 (1988).
235. R. E. Grojean, D. Feldman, and J. F. Roach, "Production of Flat Top Beam Profiles for High Energy Lasers," *Rev. Sci. Instrum.* **51**:375–376 (1980).
236. B. Fan, R. E. Tibbetts, J. S. Wilczynski, and D. F. Witman, "Laser Beam Homogenizer," U.S. Patent 4,744,615, 1988.
237. R. G. Waarts, D. W. Nam, D. E. Welch, D. R. Scifres, J. C. Ehlert, W. Cassarly, J. M. Finlan, and K. Flood, "Phased 2-D Semiconductor Laser Array for High Output Power," *Proc. SPIE* **1850**:270–280 (1993).
238. J. R. Jenness Jr., "Computing Uses of the Optical Tunnel," *Appl. Opt.* **29**(20):2989–2991 (1990).
239. L. J. Krolak and D. J. Parker, "The Optical Tunnel—A Versatile Electrooptic Tool," *J. Soc. Motion Pict. Telev. Eng.* **72**:177–180 (1963).
240. T. D. Milster and T. S. Tkaczyk, "Miniature and Micro-Optics," in *OSA Handbook of Optics*, 3rd ed., vol. I, chap. 22, McGraw-Hill, New York, 2009.
241. R. A. Sprague and D. R. Scifres, "Multi-Beam Optical System Using Lens Array," U.S. Patent 4,428,647, 1984.
242. K. Flood, B. Cassarly, C. Sigg, and J. Finlan, "Continuous Wide Angle Beam Steering Using Translation of Binary Microlens Arrays and a Liquid Crystal Phased Array," *Proc. SPIE* **1211**:296–304 (1990).
243. W. Kuster and H. J. Keller, "Domed Segmented Lens System," U.S. Patent 4,930,864, 1990.
244. W. J. Cassarly, J. M. Davenport, and R. L. Hansler, "Uniform Light Delivery Systems," SAE, Paper No. 960490 (1996).
245. N. A. Zhidkova, O. D. Kalinina, A. A. Kuchin, S. N. Natarovskii, O. N. Nemkova, and N. B. Skobeleva, "Use of Lens Arrays in Illuminators for Reflected-Light Microscopes," *Opt. Mekh. Promst.* **55**:23–24 (September 1988); *Sov. J. Opt. Technol.* **55**:539–541 (1988).
246. S. Ohuchi, T. Kakuda, M. Yatsu, N. Ozawa, M. Deguchi, and T. Maruyama, "Compact LC Projector with High-Brightness Optical System," *IEEE Trans. Consumer Electronics* **43**(3):801–806 (1997).
247. K. Rantsch, L. Bertele, H. Sauer, and A. Merz, "Illuminating System," U.S. Patent 2,186,123, 1940.
248. J. R. Miles, "Lenticulated Collimating Condensing System," U.S. Patent 3,296,923, 1967.
249. M. Ohtu, "Illuminating Apparatus," U.S. Patent 4,497,013, 1985.
250. K. Matsumoto, M. Uehara, and T. Kikuchi, "Illumination Optical System," U.S. Patent 4,769,750, 1988.
251. A. H. J. Van den Brandt and W. Timmers, "Optical Illumination System and Projection Apparatus Comprising Such a System," U.S. Patent 5,098,184, 1992.
252. F. Watanabe, "Illumination Optical System," U.S. Patent 5,786,939, 1998.
253. X. Deng, X. Liang, Z. Chen, W. Yu, and R. Ma, "Uniform Illumination of Large Targets Using a Lens Array," *Appl. Opt.* **25**:377–381 (1986).
254. N. Nishi, T. Jitsuno, K. Tsubakimoto, M. Murakami, M. Nakatsuma, K. Nishihara, and S. Nakai, "Aspherical Multi Lens Array for Uniform Target Illumination," *Proc. SPIE* **1870**:105–111 (1993).
255. Y. Ozaki and K. Takamoto, "Cylindrical Fly's Eye Lens for Intensity Redistribution of an Excimer Laser Beam," *Appl. Opt.* **28**:106–110 (1989).
256. S. Glocker and R. Goring, "Investigation of Statistical Variations between Lenslets of Microlens Arrays," *Appl. Opt.* **36**(19):4438–4445 (1997).
257. T. H. Bett, C. N. Danson, P. Jinks, D. A. Pepler, I. N. Ross, and R. M. Stevenson, "Binary Phase Zone-Plate Arrays for Laser-Beam Spatial-Intensity Distribution Conversion," *Appl. Opt.* **34**(20):4025–4036 (1995).

258. Y. Kato, K. Mima, N. Miyanaga, S. Arinaga, Y. Kitagawa, M. Nakatsuka, and C. Yamanaka, "Random Phasing of High-Power Lasers for Uniform Target Acceleration and Plasma-Instability Suppression," *Phys. Rev. Lett.* **53**(11):1057–1060 (1984).
259. L. B. W. Jolley, J. M. Waldram, and G. H. Wilson, *The Theory and Design of Illuminating Engineering Equipment*, John Wiley & Sons, New York, 1931.
260. D. G. Burkhard and D. L. Shealy, "Design of Reflectors Which Will Distribute Sunlight in a Specified Manner," *Solar Energy* **17**:221–227 (1975).
261. D. G. Burkhard and D. L. Shealy, "Specular Aspheric Surface to Obtain a Specified Irradiance from Discrete or Continuous Line Source Radiation: Design," *Appl. Opt.* **14**(6):1279–1284 (1975).
262. O. K. Kusch, *Computer-Aided Optical Design of Illuminating and Irradiating Devices*, ASLAN Publishing House, Moscow, 1993.
263. T. E. Horton and J. H. McDermit, "Design of Specular Aspheric Surface to Uniformly Radiate a Flat Surface Using a Nonuniform Collimated Radiation Source," *J. Heat Transfer*, 453–458 (1972).
264. J. S. Schruben, "Analysis of Rotationally Symmetric Reflectors for Illuminating Systems," *J. Opt. Soc. Am.* **64**(1):55–58 (1974).
265. J. Murdoch, *Illumination Engineering, from Edison's Lamp to the Laser*, MacMillan Publishing, 1985.
266. R. Winston, "Nonimaging Optics: Optical Design at the Thermodynamic Limit," *Proc. SPIE* **1528**:2–6 (1991).
267. D. Thackeray, "Reflectors for Light Sources," *J. Photog. Sci.* **22**:303–304 (1974).
268. P. W. Rhodes and D. L. Shealy, "Refractive Optical Systems for Irradiance Redistribution of Collimated Radiation: Their Design and Analysis," *Appl. Opt.* **19**:3545–3553 (1980).
269. W. A. Parkyn, "Design of Illumination Lenses via Extrinsic Differential Geometry," *Proc. SPIE* **3428**: 154–162 (1998).
270. R. Winston and H. Ries, "Nonimaging Reflectors as Functionals of the Desired Irradiance," *J. Opt. Soc. Am. A* **10**(9):1902–1908 (1993).
271. W. B. Elmer, "The Optics of Reflectors for Illumination," *IEEE Trans. Industry Appl.* **IA-19**(5):776–788 (September/October 1983).
272. J. M. Gordon, P. Kashin, and A. Rabl, "Nonimaging Reflectors for Efficient Uniform Illumination," *Appl. Opt.* **31**(28):6027–6035 (1992).
273. J. M. Gordon and A. Rabl, "Nonimaging Compound Parabolic Concentrator-Type Reflectors with Variable Extreme Direction," *Appl. Opt.* **31**(34):7332–7338 (1992).
274. J. Gordon and H. Ries, "Tailored Edge-Ray Concentrators as Ideal Second Stages for Fresnel Reflectors," *Appl. Opt.* **32**(13):2243–2251 (1993).
275. H. R. Ries and R. Winston, "Tailored Edge-Ray Reflectors for Illumination," *J. Opt. Soc. Am. A* **11**(4): 1260–1264 (1994).
276. A. Rabl and J. M. Gordon, "Reflector Design for Illumination with Extended Sources: The Basic Solutions," *Appl. Opt.* **33**(25):6012–6021 (1994).
277. P. T. Ong, J. M. Gordon, A. Rabl, and W. Cai, "Tailored Edge-Ray Designs for Uniform Illumination of Distant Targets," *Opt. Eng.* **34**(6):1726–1737 (1995).
278. P. T. Ong, J. M. Gordon, and A. Rabl, "Tailoring Lighting Reflectors to Prescribed Illuminance Distributions: Compact Partial-Involute Designs," *Appl. Opt.* **34**(34):7877–7887 (1995).
279. D. Jenkins and R. Winston, "Tailored Reflectors for Illumination," *Appl. Opt.* **35**(10):1669–1672 (1996).
280. D. Jenkins and R. Winston, "Integral Design Method for Nonimaging Concentrators," *J. Opt. Soc. Am. A* **13**(10):2106–2116 (1996).
281. J. M. Gordon, "Simple String Construction Method for Tailored Edge-Ray Concentrators in Maximum-Flux Solar Energy Collectors," *Solar Energy* **56**:279–284 (1996).
282. H. Schering and A. Merz, "Illuminating Device for Projectors," U.S. Patent 2,183,249, 1939.
283. Y. Kawamura, Y. Itagaki, K. Toyoda, and S. Namba, "A Simple Optical Device for Generating Square Flat-Top Intensity Irradiation from a Gaussian Laser Beam," *Opt. Comm.* **48**:44–46 (1983).
284. T. Sameshima and S. Usui, "Laser Beam Shaping System for Semiconductor Processing," *Opt. Comm.* **88**:59–62 (1992).
285. L. Austin, M. Scaggs, U. Sowada, and H.-J. Kahlert, "A UV Beam-Delivery System Designed for Excimers," *Photonics Spectra* 89–98 (May 1989).

286. A. Brunsting, "Redirecting Surface for Desired Intensity Profile," U.S. Patent 4327972, 1982.
287. R. J. Bruno and K. C. Liu, "Laserbeam Shaping for Maximum Uniformity and Minimum Loss," *Laser Appl.* 91–94 (April 1987).
288. D. D. Dourte, C. Mesa, R. L. Pierce, and W. J. Spawr, "Optical Integration with Screw Supports," U.S. Patent 4,195,913, 1980.
289. S. L. Ream, "A Convex Beam Integrator," *Laser Focus*: 68–71 (November 79).
290. V. J. Doherty, "Design of Mirrors with Segmented Conical Surfaces Tangent to a Discontinuous Aspheric Base," *Proc. SPIE* 399:263–271 (1983).
291. D. M. Dagenais, J. A. Woodroffe, and I. Itzkan, "Optical Beam Shaping of a High Power Laser for Uniform Target Illumination," *Appl. Opt.* 24:671–675 (1985).
292. F. M. Dickey and B. D. O'Neil, "Multifaceted Laser Beam Integrators: General Formulation and Design Concepts," *Opt. Eng.* 27:999–1007 (1988).
293. J. Geary, "Strip Mirror Integrator for Laser Beam Uniformity on a Target," *Opt. Eng.* 28(8):859–864 (August 1989).
294. B. Lu, J. Zheng, B. Cai, and B. Zhang, "Two-Dimensional Focusing of Laser Beams to Provide Uniform Irradiation," *Opt. Comm.* 149:19–26 (1998).
295. T. Henning, L. Unnebrink, and M. Scholl, "UV Laser Beam Shaping by Multifaceted Beam Integrators: Fundamentals and Principles and Advanced Design Concepts," *Proc. SPIE* 2703:62–73 (1996).
296. L. Unnebrink, T. Henning, E. W. Kreutz, and R. Poprawe, "Optical System Design for Excimer Laser Materials Processing," *Proc. SPIE* 3779:413–422 (1999).
297. S. R. David, C. T. Walker, and W. J. Cassarly, "Faceted Reflector Design for Uniform Illumination," *Proc. SPIE* 3482:437–446 (June 12, 1998).
298. J. R. Coaton and A. M. Marsden, *Lamps and Lighting*, 4th ed., John Wiley & Sons, New York, 1997.
299. R. J. Donohue and B. W. Joseph, "Computer Design of Automotive Lamps with Faceted Reflectors," *IES* 36–42 (October 1972).
300. E. H. Wiley, "Projector Lamp Reflector," U.S. Patent 4,021,659, 1976.
301. W. P. Laudenschlager, R. K. Jobe, and R. P. Jobe, "Light Assembly," U.S. Patent 4,153,929, 1979.
302. R. J. Heimer, "Multiple Apparent Source Optical Imaging System," U.S. Patent 4,241,389, 1980.
303. H. Hamada, K. Nakazawa, H. Take, N. Kimura, and F. Funada, "Lighting Apparatus," U.S. Patent 4,706,173, 1987.
304. M. Ronnelid, B. Perers, and B. Karlsson, "Optical Properties of Nonimaging Concentrators with Corrugated Reflectors," *Proc. SPIE* 2255:595–602 (1994).
305. D. J. Lamb, J. F. Van Derlofske, and L. W. Hillman, "The Use of Aspheric Surfaces in Waveguide Illumination Systems for Automotive Displays," SAE Technical Paper 980874, 1998.
306. W. J. Cassarly and T. L. Davenport, "Non-Rotationally Symmetric Mixing Rods," International Optical Design Conference, *Proc. SPIE* 6342: (June 2006).

This page intentionally left blank.

DO NOT DUPLICATE

Anurag Gupta

*Optical Research Associates
Tucson, Arizona*

R. John Koshel

*Photon Engineering LLC and
College of Optical Sciences
University of Arizona
Tucson, Arizona*

40.1 GLOSSARY

Illuminance. Luminous flux incident on a surface per unit projected area in the direction of emission relative to the surface normal. $1 \text{ lux} = 1 \text{ lumen/m}^2$.

Intensity. Luminous flux emitted by a source per unit solid angle in a given direction. $1 \text{ candela} = 1 \text{ lumen/steradian}$.

Luminance. Luminous flux emitted in a given direction per unit solid angle per unit projected area in the direction of emission relative to the surface normal. $1 \text{ nit} = 1 \text{ lumen}/(\text{m}^2 \times \text{steradian})$.

Luminaire. Lamp or lighting fixture that includes optical components, baffles, housing, and electronics.

Display. It refers to many things: a computer monitor, a projected image, and a piece of art or a decorative item.

Backlighting. It refers to illumination of an object from behind. The object can be opaque, translucent, or transparent. The light source is usually large in size and diffuse.

40.2 INTRODUCTION

Lighting is an area of science that includes the interaction between light and people in their daily lives. The primary goals of lighting are to provide the illumination to perform tasks, direct people to desired locations and provide a sense of security. Additionally, lighting has a profound effect on mood and sleep-wake cycles of all living beings. As mankind has advanced technologically, the needs fulfilled by lighting have also increased to additionally provide relaxation, alter moods, attract people, provide entertainment, create virtual environments and improve human productivity.

Unlike most of other fields in optics, lighting is a more subjective field than objective—it is based on the emission aspects and how it is perceived within its surroundings. The subjective nature is based upon our interactions with lighting shaped by vision biology and brain perception. The capabilities of the human eye largely determine the detectable range and variations (in both time and space) in colors and brightness. Perception depends upon the brain's interpretation of the input from the eye and is influenced by past experiences. For example, although brightness and the color

of lit objects are mathematically represented by luminance metrics as a function of distance, direction and wavelength, the perception of these quantities is context dependent based upon the observer's experiences and environment. Thus, the same illumination levels in two different environments can be perceived as two drastically different lighting outcomes.

The importance of understanding lighting goes beyond the basic need to illuminate objects or surroundings. The importance of representing lighting accurately has been exemplified in art over the centuries.¹ Understanding lighting in the context of human perception and the ability to represent and simulate lighting on computers has attracted attention in fields such as virtual reality for training especially in the fields of medicine, aviation and computer-aided design, and entertainment such as video games and movies. As our understanding and computational capabilities have increased over time, so has the sophistication and quality of virtual environments and the entertainment media.

An effective lighting design is contextual, cultural and is well integrated into its surroundings. Light interacts intimately with everything it impinges upon. Every object, including humans, plants and architectural elements, has distinct scatter, reflective, transmissive, and absorptive properties that are dependent upon wavelength and direction of light. As a result, the lit objects contribute to the appearance of the scene. Lighting therefore must complement in concert with the architecture and its surroundings in form, composition, and style to meet our expectations for the lit environment. For example, our expectations for the lighting environment of a casino are quite different to that of a sports stadium, a retail store, an office or a hospital. In each case, there is a distinct purpose that directs the layout of the light sources with respect to the objects. Therefore, the lighting must change in each case as dramatically as the differences between the contexts of lighting. If it does not, then responses to the lit environment are often not positive, typically making, for example, a poor work environment for an office, a lackluster sales venue for retail, or an uncomfortable room in someone's home. This subjectivity is the most difficult aspect of lighting design. After all, the success of any lighting scheme depends upon being able to meet the expectations of its users. Simply said, we know a good lighting design as soon as we see it, but it has limited quantifiable metrics.

The history of modern lighting traces its origin back to the advent of artificial light sources: the incandescent light bulb, more than a century ago. Currently there is a wide variety of sources to choose from: incandescent (includes halogen), discharge (fluorescent, high-intensity discharge, sodium vapor), lasers, electroluminescent (includes LEDs and OLEDs), and daylight. The choice of light source depends on operating characteristics, cost, efficiency, and safety. The luminaire optics play a critical role in shaping the light distribution from the source. Thus the choice of the source and luminaire optics is an integral part of the lighting design process.

Any lighting design must comply with the relevant government regulations for the specific purpose of lighting. For example, safety is critical in transportation lighting, so there are stringent regulations on the relative intensity distribution from street lights and automobile headlamps to minimize glare and achieve the desired visibility. There are increasing numbers of government mandates on using efficient sources in many countries, such as the banning of inefficient incandescent sources in favor of fluorescent and LED light sources. In addition to regulations, there are guidelines for best lighting practices in various situations. These guidelines are published by national and international committees such as CIE (Commission Internationale de l'Éclairage or The International Commission on Illumination), IESNA (Illumination Engineering Society of North America), SAE International (Society of Automotive Engineers), CIBSE (Chartered Institution of Building Services Engineers), and many others. These groups publish guidelines on many aspects of lighting in the interiors and exteriors of homes, offices, educational institutions, hospitals, entertainment facilities, malls, industrial complexes, sports stadiums, theatres, museum, streets, parks, transportation, and even underwater. The goals of these guidelines are to help, at a minimum, to create designs that are functional, efficient and provide safe and comfortable lighting. It is in the hands of the lighting designer to add to this mandated objective illumination to achieve the desired subjective lighting—in other words, the aesthetics.

In this chapter, we provide an introduction to many facets of lighting by touching upon the perceptual and biological factors that guide lighting design (Sec. 40.3), design elements and methods

to create functional and aesthetically pleasing designs (Sec. 40.4), technology of sources and design of relevant optics (Sec. 40.5), and measurement of lighting conditions (Sec. 40.6). We end the chapter with application examples on interior and exterior lighting in Sec. 40.7. These sections provide comprehensive data on source selection and guidelines for best practices in general and in specific application areas such as lighting for offices, homes, healthcare, retail, and transportation. Although very important aspects of lighting engineering, we do not discuss electronic control mechanisms, maintenance, and commissioning of a lighting design due to limited space.

We make extensive use of terms used in Radiometry and Photometry. The reader is encouraged to refer to the chapters on various aspects of radiometry and photometry in Chaps. 34 to 39 in this volume.

40.3 VISION BIOLOGY AND PERCEPTION

The lighting design is guided by our understanding of perception and vision biology. In this section, we touch upon various aspects of biology of vision and perception.

Vision Biology

Biological aspects of human visual system response² that are parameterized and used in lighting design are visual acuity (resolution, vernier, recognition, and stereoscopic), color sensitivity, accommodation, field of view, and adaptability to color and brightness changes.

Once adapted, the human visual system response does not change appreciably with time. The human visual system responds to over eleven orders of magnitude of luminance. However, at any given instance, only 2 to 3 orders of magnitude are adapted to by the eye. It takes up to 60 minutes to fully adapt to lighting conditions. Light adaptation takes place via change in the eye pupil size that controls the amount of light entering the eye (2 mm to 8 mm), photochemical processes in the retinal cells and neural processes that respond to change in the luminance below 600 nits when cone cells have not fully bleached. Neural adaptation occurs quickly, within the first 200 ms, and allows us to adapt to 2 to 3 orders of magnitude of luminance fluctuation, such as in lit spaces of building interiors. Pupil adaptation via change in the pupil size takes up to few minutes and allows us to adapt up to 1 to 2 orders of magnitude of luminance fluctuation. For tasks that require good color discrimination, several minutes to an hour are needed for this color adaptation. The higher the luminance, the shorter is the adaptation time. Based upon luminance, three vision conditions exist:

1. Photopic vision occurs when bright illumination conditions in the visible (luminance >3 nits) exist. Both the rod and the cone cells in the retina are excited, but the rods are saturated and thus effectively ignored except for peripheral vision. Full color vision with highest resolution is possible. Most indoor lighting conditions ensure photopic enabling lighting conditions.
2. Scotopic vision occurs when low illumination level in the visible (luminance <0.001 nits) exists. Only the rod cells in the retina are excited. No color vision occurs and only low-resolution peripheral vision is possible. Lighting design is usually not done for the scotopic domain but it must take into account our peripheral vision capability wherever possible to help guide the individual toward relevant objects.
3. Mesopic vision is described by the transition between photopic and scotopic. The rod cells are excited and the cone cells are partially excited. The eye has limited color discrimination and resolution capabilities. Most outdoor artificial lighting conditions operate in this region.

The luminance and color distributions should be such that it includes individuals with impaired vision that cannot be corrected such as reduced illumination at the retina (50 percent for a 50-year old as compared to a 20-year old) due to smaller pupil size and reduced transmission efficiency of the eye, reduced contrast, increased glare sensitivity, loss of accommodation and reduced field of view caused by ageing, macular degeneration, cataract, or glaucoma. For most such cases

increased illuminance, slow transients from the light to dark regions and an environment with reduced glare help considerably.³

Perception

It is the perception of lit environment that eventually determines the acceptability and adequateness of lighting conditions. Lighting perception is ultimately determined by the human visual system response and brain processes that are individual dependent. It is the latter that brings considerable subjectivity to lighting perception. Lighting design criteria for specific environments are often guided by studies that determine average perceptual responses. These guidelines are continually evolving, in response to geopolitical factors, research on vision biology and technology.

Depending upon an object's optical and physical characteristics and how it is viewed in relation to its surroundings, lighting can alter its perceived visual attributes. Visual attributes of physical objects are defined in terms of brightness, lightness, hue, saturation, transparency, and glossiness. We explain each of these attributes below and also show that not every object has each of these attributes.

Brightness is the perceptual correlate of luminance. Although, in the absence of a background, the perception of brightness of an object is proportional to its luminance in a logarithmic manner (log of brightness is proportional to log of luminance), the perception of brightness is dependent upon adaptation to the surroundings and the relative hue and saturation of the object. Figure 1 shows an example of how the brightness perception is altered by surroundings.⁵⁰ Different portions of a uniform luminance bar appear to have different brightness depending upon the local background. Similarly, car headlights appear much brighter in the dark than in daylight. Both vision biology and perception play a role here. Due to low ambient lighting, the pupil dilates and admits much more light (up to 16 times in night than in the day). Therefore, any bright source leads to sudden pupil contraction and ensuing discomfort. The sources appear even brighter due to a dark background, an effect similar to Fig. 1. Color perception is similarly affected. Color-saturated objects tend to appear brighter and vice versa.

Lightness is the perceptual correlate of diffuse reflectance. Our perception of an object being dark or light depends upon our estimate of its reflectance. For example, a piece of white paper appears white irrespective of the illumination falling upon it as we know from prior experience that it has a high reflectance. We tend to perceive it whiter than other objects of lesser reflectance even when the flux reflected by white paper in relatively low illumination is lower than a grey object that is placed under higher illumination. This confusion would not occur if we view a small region of the object through an aperture without knowing anything about the object. The aperture masks the contextual information and forces us to make objective judgments.



FIGURE 1 Perception of brightness. The rectangular bar in the center has constant luminance yet it seems brighter against a darker background at the right hand side.

Hue is the perception of dominant wavelength of the color spectrum transmitted or reflected by the object. It helps us judge if an object appears close to a known color such as red, blue, yellow, purple, green, colorless, or a combination of multiple colors such as bluish-green. Saturation is the perception of the extent of color purity. Highly saturated colors have a narrow band of wavelengths, centered on the hue. Transparency is the perception of the degree of light penetration into the object. Glossiness is the perception of smoothness of a surface relative to a matte finish.

Depending upon the properties of an object, the lighting conditions and how it is viewed, one or more visual attributes can be determined. For example, if an object is viewed through an aperture that hides the information of its surroundings or is illuminated such as nothing other than the object is visible then only attributes like brightness, hue, and saturation can be observed. A self-illuminated object or an object that gives the appearance of being so (such as a lit computer monitor) does not display lightness or glossiness, but if it stops emitting light, these attributes can be observed under external illumination.

Our brain attempts to maintain a perceptual constancy of shape, color, size, and lightness of lit objects based on past experiences and the context in which an object exists. Perceptual constancy allows us to maintain a certain perception of an object or a scene under changing viewing conditions. For example, we can recognize an apple as such when viewed at different angles, surroundings, or lit conditions. We can perceive a skyscraper as a tall building even if seen from far away and having a small image on the retina. Similarly, a tilted book appears rectangular although the image in the retina is trapezoidal. We can perceive the colors of common objects around us fairly correctly while wearing mildly colored glasses. Lighting can help maintain or alter perceptual constancies. To maintain perceptual constancies the lighting conditions must be such as there is adequate ambient lighting with high-color-rendering sources and without any disability glare. We discuss the terms such as color rendering, glare, and the impact of ambient lighting in the next section. The position of the light sources must be obvious to the observer even if not directly visible to establish the directionality of illumination. Direction of illumination is important as we tend to expect light to come from above and our perception of lit objects takes that into account at a subconscious level.

The relative distribution of light can impact the perception of the space itself and as such can impact human behavior.⁴ For example, there is a general tendency among humans to be attracted to brighter regions of space. It can be demonstrated by viewing people that there is a clear trend toward walking on brighter areas of pathways or facing brightly illuminated areas of a restaurant. That is why brightly lit shopping malls are quite effective in attracting traffic as compared to a poorly lit shopping complex. The knowledge of human behavior toward lit space is also used in making more effective lit object displays and navigational aids. There are four distinct categories of light distribution⁵ that affect the perception of space: privacy, relaxation, visual clarity, and spaciousness. An effect of privacy can be created by utilizing nonuniform high-brightness illumination across the vertical surfaces in the room with dark spaces in the occupant domains and low ambient luminance. An effect of relaxation can be created by using nonuniform warm (correlated color temperature (CCT) <3500 K) ambient light across the room. Visual clarity in the environment is emphasized with cool light (CCT >4000 K), high and uniform brightness near the center of the room and at all task planes. In addition, higher emphasis is given to ceiling and horizontal surfaces. A sense of spaciousness can be implemented with uniform room illumination and relatively higher levels of brightness on the walls and ceilings.

Summary

Understanding perception and biology of vision helps us develop models to describe desired lighting conditions. An example is the development of the relative visual performance (RVP) model. This model has been extensively developed by Mark Rea and his colleagues over the last few decades⁶ to obtain the relative visual performance of a given task under different lighting conditions and establish lighting guidelines. RVP is provided as a probability of performing a visual task successfully under given lighting conditions. Task performance depends upon both visual and nonvisual aspects

of the task. To obtain the true impact of lighting conditions on task performance, it is necessary to isolate those tasks for evaluation that are dominated by the visual component. The impact of the nonvisual components is minimized by quantifying their effect to the fullest possible extent and subtracting it from the overall task performance. A key finding of the RVP model is that the visual performance improves rapidly as the luminance contrast between the task and the background increases up to 40 percent. Beyond this value, the improvement in visual performance is negligible with increase in contrast. Luminance contrast in this context is defined as

$$\text{Luminance contrast} = 100(L_{\text{Task}} - L_{\text{background}}) / L_{\text{background}} \quad (1)$$

Visual performance curves (performance metric versus luminance contrast) can be evaluated for various task sizes and background luminance. A major limitation of the existing RVP model is its limited validity to only those tasks that are quantifiable by task size, luminance contrast, and background luminance and only under the conditions of foveal vision. We need more sophisticated task performance models to cover a larger range of tasks.

In the sections to follow, design guidelines make use of the understanding of lighting perception and vision biology to justify their development.

40.4 THE SCIENCE OF LIGHTING DESIGN

The lighting design process begins with identifying the needs to be addressed. An understanding of the functions of lighting and knowledge of basic building blocks helps us in making a preliminary design. The quality of lighting is determined by its ability to fulfill human needs in an economical and environmental friendly manner while at the same time complementing the architecture in form, composition, style, codes, and standards. We discuss the lighting design process in the following subsections:

Design considerations. We discuss the factors involved in creating a lighting environment for an application. The categories discussed are: goals, context, illuminance, color, visual discomfort, trespass and light induced damage of objects.

Functions of lighting. Here we discuss the four primary functions of lighting: ambient, task, decorative and accent.

Lighting geometries to achieve specific functions of lighting. Here we discuss the building blocks for lighting design.

Properties of objects and their impact on lit scene are discussed.

Modeling. Here we discuss the techniques used to simulate a lit environment in order to achieve the best design.

Design Considerations

Goals Lighting for any application must take into account the needs (both human and nonhuman such as plants or animals), lighting economics, environment impact, and architectural aspects of the application. Human needs include the desired degree of visibility, comfort, ability to perform the needed tasks, social communication, ambience, and aesthetics.

Context Lighting helps create a perceptual environment or ambience to suit a specific application such as office, home, lobby, restaurant, casino, or a sports stadium. Appropriate selection of lighting schemes and luminaires that complement the architecture and interior design helps in achieving the desired ambience.

TABLE 1 IESNA Guidelines on Illumination Categories and Average Illuminance Levels

Category	Average Illuminance (lx)
Public spaces	30
Simple orientation or short visits in a new environment	50
Working spaces where simple visual tasks are performed	100
Performance of visual tasks of high contrast and large size	300
Performance of visual tasks of high contrast and small size, or visual tasks of low contrast and large size	500
Performance of visual tasks of low contrast and small size	1,000
Performance of visual tasks near vision threshold	3,000–10,000

Large size: Object's projected solid angle subtense at the eye $>4.0 \times 10^{-6}$ sr.

Small size: Object's projected solid angle subtense at the eye $\leq 4.0 \times 10^{-6}$ sr but not near the visual acuity limit.

Low contrast: ≤ 0.3 but greater than visual threshold.

High contrast: >0.3 .

Illuminance—Horizontal and Vertical Horizontal and vertical illuminances refer to illuminance distribution on horizontal and vertical planes respectively. Table 1 describes IESNA guidelines on illumination categories and average illuminance levels needed for each.⁷ These guidelines do not apply to special situations that involve setting up a particular ambience or focusing on an object for emphasis.

The uniformity of luminance/illuminance is generally defined by the ratio of maximum to minimum luminance/illuminance. The need for uniformity across the field of view and across the entire space depends upon the application. The human eye is a brightness detector and is thus responsive to changes in luminance. To calculate luminance, illuminance, and the reflectivity of surfaces must be taken into account. For Lambertian surfaces (surfaces of constant luminance, independent of the viewing direction),

$$\text{Luminance} = (\text{illuminance} \times \text{surface reflectivity}) / \pi \quad (2)$$

Generally, a luminance uniformity of 0.7 within the field of view across the task is considered adequate as the eye is unable to detect these variations. Not all tasks require high luminance uniformity. For example, in tasks involving inspection of 3D objects, nonuniform illumination is able to highlight the geometrical features, especially surface textures much better due to being able to provide better depth perception. In lit environments, a nonuniform light distribution is used to provide perceptions of privacy or exclusivity, for example in retail lighting or in restaurants. A luminance ratio (maximum luminance: minimum luminance) of greater than 15:1 is generally considered undesirable. In the applications section, we discuss the recommended luminance ratios for several scenarios in a variety of applications.

Color Lighting influences the color appearance of lit objects. For most lighting applications, white light is the standard form of illumination: exteriors (landscape, roadways, buildings, city, and stadiums) and interiors (homes, offices, restaurants, museums, industrial complexes, and shopping malls). Saturated colors in lighting are used only in special applications like indicators or signals, displays, color-specific industrial applications such as those involving color discrimination, special effects in casinos, hotels, malls, or discotheques, and so forth. The chapter on colorimetry in this *Handbook* (Vol. III, Chap. 10) provides an excellent introduction to the subject of color.

Depending upon needs such as aesthetics, task performance, and color-dependent reflective properties of objects, an appropriate light spectrum must be selected. The desired light spectrum can be achieved by using light sources that emit in that spectrum, by using static filters on the sources to tailor the spectrum or by spatial and time-averaged color mixing. Spatial color mixing involves using multiple light sources that emit in different portions of the desired spectrum but are laid out in such a manner that the lit environment or object appear to be illuminated by light

having the combined spectrum of individual light sources. Time-averaged color mixing involves high-frequency (>60 Hz) mixing of different portions of the light spectrum in different proportions. If the frequency of color mixing is high enough, the brain perceives a specific color based on the time average of the varying spectra used in their respective strengths. For example, in modern digital projectors, a color wheel is used that has various segments of different spectral transmission. When it rotates through those segments, one can create the appearance of any color within the color gamut of the light source. In the section on LEDs, we discuss color mixing to create white light or any desired color. Spatial and temporal color mixing can be used together to create a variety of effects. Although color mixing can achieve the visual perception of any desired color in emission, its ability to color render the object it illuminates depends upon the product of incident light spectrum and wavelength-dependent reflectance of the object. In this section we discuss how white light is specified and how its ability to render objects is estimated.

It is tedious to choose light sources if we have to analyze the spectrum and its impact on common objects. For white light, the color rendering index (CRI) and the correlated color temperature (CCT) help provide a quick estimate of the appearance of light and its color rendering of lit objects. For example, at a CCT of 2700 to 6500 K, a CRI ≥ 70 is adequate for common situations such as in offices and homes, a CRI ≥ 50 is sufficient for most industrial tasks and a CRI ≥ 90 is needed for stringent color discrimination tasks such as hospital surgery, paint mixing, or color matching. In the future, it is likely that different metrics based on color-appearance models would be in use. This is especially true with the increasing use and availability of a variety of light sources, especially LEDs, which have significantly different emission spectra from incandescent sources.

The CCT is the temperature of the planckian (perfect blackbody) radiator in Kelvin whose perceived color most closely resembles that of a given stimulus at the same brightness and under specified viewing conditions.⁸ To find the CCT, the nearest point on the planckian locus is considered but only in a perceptually uniform color space. The isotherms across the planckian locus in a uniform color space are represented as normal to locus curve but in a nonuniform color space such as XYZ, these are no longer perpendicular to the locus.

Color rendering is defined as the effect of an illuminant on the color appearance of objects by conscious or subconscious comparison with their color appearance under a reference illuminant.⁸ Two functional reference illuminants are currently used for calculating CRI: (1) for a source CCT up to 5000 K, a blackbody at the same color temperature is used and (2) for a source CCT above 5000 K, one of the phases of daylight is used. The phase of daylight selected is such that its chromaticity is within $1/(15E6 \text{ K})$ to that of the test source. It is defined by a mathematical formula based on the CCT of the test source.⁹ For all cases, a CRI value of 100 is considered to be a perfect match between the test source and the illuminant. The precise steps and calculations needed to calculate the CRI are recommended by the CIE.¹⁰ Here we summarize the results:

$$\text{CRI}(R_a) = 100 - 4.6\bar{E}_{UVW} \quad (3)$$

where R_a refers to the general CRI and \bar{E}_{UVW} is the average of the Euclidean distances between the color coordinates of the reflected reference and test light sources from the first 8 out of the 14 CIE-prescribed test samples (see Ref. 10). The color coordinates of the test source are chromatically adapted to the reference source and are expressed in the CIE 1964 color space. The CRI calculation is valid only when the color difference between the test source and the reference source is not large.

Although CRI is a useful metric and is widely used, it has several shortcomings. The existing method of calculating CRI is not perceptually well correlated. The high CRI predicted for sources with extreme CCT do not have good color-rendering properties. CRI is not valid for sources such as discrete spectral LEDs where CCT cannot be defined. CRI cannot be used to evaluate white-light LEDs with nonuniform spectra; the CRI predicted is quite low although the quality of white light appears to be better. The practice of using only eight test samples, none of which are saturated creates situations where high CRI sources do not color render color-saturated objects correctly. The CRI formulation can be modified to include the impact of different reference illuminants for different source types, color spaces, test sample set, different chromatic adaptation formulas or even a

reduced focus on absolute color fidelity. CIE reviews various propositions in this regard and updates its recommendations.

Visual Discomfort Visual discomfort can be caused by a variety of reasons.¹¹ Many of these reasons are context dependent where the expectation of the nature of lit environment determines the suitability of lighting. It also depends on cultural differences between various groups of people. For example, lighting flicker in a dance club may be desirable as opposed to almost all other situations. Similarly, the preference of color among different cultural groups varies considerably. Visual discomfort occurs when lighting creates perceptual confusion. For example, if the pattern of illumination is such that surfaces of higher reflectance reflect less light than the surfaces of lower reflectance, perceptual confusion may result. The causes of visual discomfort that are more specific to lighting are summarized as follows:

Insufficient lighting Insufficient lighting, especially for task performance, results in eye strain besides reduced task performance. For different tasks, there are different levels of horizontal and vertical illuminance necessary to be considered adequate. Table 1 describes suggested horizontal illuminance levels needed for certain situations. Lighting communities across the world have established guidelines for illuminance levels for a wide variety of specific tasks and environments.

Uniformity Visual discomfort occurs when the uniformity across the field of view is not as expected. A high uniformity can be as undesirable as a high degree of nonuniformity especially when considered across the entire visual field of view. Both can cause severe eye strain. In the section on applications, we provide examples of preferred luminance ratios between task and its vicinity.

Glare Glare results when there are regions of unexpected very high levels of luminance in the field of view. Glare can be direct or indirect. Direct glare occurs when a light source or a portion of it is visible such as in overhead lamps, automobile headlights, or direct sunlight. Indirect glare occurs when the light is reflected or scattered directly into the eye, such as the sky reflecting off a lake surface and obscuring the view beneath the water surface. Glare comes in many forms:¹²

1. Flash blindness: This is caused by a sudden onset of bright light leading to temporary bleaching of retinal pigment.
2. Paralyzing glare: This is caused by sudden illumination with bright light that can temporarily “freeze” the movements of the observer.
3. Distracting glare: This is caused by flashing bright sources of light in the peripheral vision field.
4. Retinal damage: When the light is bright enough to cause retinal damage.
5. Saturation or dazzle: When a large portion of the vision field is dominated by bright source(s), which can be alleviated by wearing low transmittance eye glasses.
6. Adaptation: When one enters from a low ambient illumination region to a bright ambient illumination region without a transition region to help in vision adaptation.
7. Disability: This is caused by intraocular light scattering which reduces the luminance contrast [See Eq. (1)] of the task image at the retina. The impact is loss in task performance due to reduced visibility.
8. Discomfort: When the glare causes discomfort or distraction but does not affect the task visibility to the extent of limiting its performance

Note that several forms of glare can exist concurrently, especially discomfort with the other types. Discomfort and disability glare are the most commonly experienced glare forms. There are several mechanisms available to estimate the impact of these forms of glare. We discuss briefly the CIE-recommended models for disability and discomfort glare.

Disability glare Disability glare occurs when the luminance contrast [Eq. (1)] of the task image at the retina falls due to the superposition of intraocular scattered light on the retinal image.

This background noise from scattered light can be thought of as viewing through a veil. Therefore, it makes sense to define an equivalent glare (or veiling) luminance (EVL) that mimics the impact of disability glare. The luminance contrast C is described as

$$C = \left| \frac{(L_b + L_{\text{EVL}}) - (L_{\text{object}} + L_{\text{EVL}})}{(L_b + L_{\text{EVL}})} \right| = \left| \frac{L_b - L_{\text{object}}}{L_b + L_{\text{EVL}}} \right| \quad (4)$$

where L is object luminance and b is background.

As the equivalent veiling luminance L_{EVL} increases, contrast at the retina C reduces. The contrast reduction is especially severe during difficult viewing conditions such as fog or nighttime when the object luminance is low. That is why car headlights are a much stronger glare source in the night than in the day.

Equivalent veiling luminance L_{EVL} is defined as:¹³

$$L_{\text{EVL}} = \sum_i \left[\frac{10}{\theta_i^3} + \left\{ \frac{5}{\theta_i^2} + \frac{0.1p}{\theta_i} \right\} \cdot \left\{ 1 + \left(\frac{A}{62.5} \right)^4 \right\} + 0.0025p \right] \cdot E_i \quad (5)$$

where E_i = illuminance at the eye due to i th glare source

θ_i = angle of the glare source (in degrees) from the line of sight, $0.1^\circ < \theta_i < 100^\circ$

p = eye pigmentation factor (0 for black eyes, 0.5 for brown eyes, 1.0 for light eyes, and 1.2 for very light-blue eyes)

A = age of the viewer in years

For young adults (<35 years of age) and for glare source angle, $1^\circ < \theta_i < 30^\circ$, Eq. (5) approximates to

$$L_{\text{EVL}} = \sum_i 10(E_i/\theta_i^2) \quad (6)$$

The EVL as described above is strictly due to intraocular scatter. In practice, it is necessary to add to this the luminance from external scatterers, such as fog or dust in the atmosphere, to yield the correct retinal contrast.

Equation (5) is currently the most sophisticated recommended treatment for disability glare. It can be applied toward a wide variety of circumstances, including indoor lighting, street lighting, and bright sky at a tunnel's exit.

For road lighting, the CIE recommendation on disability glare¹⁴ is given by a percentage threshold increment (TI) described by Eq. (7). TI is limited between 10 and 15 percent.

$$TI = 65(L_{\text{EVL}}/L^{0.8}) \quad (7)$$

where L_{EVL} = equivalent veiling luminance as described by Eq. (5)

L = average road surface luminance

Discomfort glare There are many formulations that describe discomfort glare. Each model is prescribed for well defined geometries and sources. Discomfort glare has been described by the visual comfort probability (VCP) model¹⁵ in North America, British Glare Index system¹⁶ (CIBSE) and the European glare limiting system.¹⁷⁻¹⁹ Each of these systems has validity under specific constraints. Many luminaire manufacturers in North America and Europe provide VCP or glare index tables for worst case scenarios. CIE has proposed a Unified Glare Rating (UGR) model²⁰ to replace these systems. We describe here the formulation recommended by CIE.

The UGR formula is described by Eq. (8). UGR values range from 5 to 30, the higher values signifying a greater level of discomfort. For home and offices, UGR is specified at <20 and for industrial application it is >20. This formula is valid for source areas between 0.005 and 1.5 m². For smaller

TABLE 2 Glare Specification for Large Sources Such as an Illuminated Ceiling

Maximum Average Illuminance (lx)	UGR
300	13
600	16
1000	19
1600	22

sources, UGR overestimates the glare and for larger sources, UGR underestimates the glare. Therefore, for ranges outside the validity of UGR, CIE has provided detailed prescriptions to tackle small, large, and complex luminaires.²¹ For source areas smaller than 0.005 m², Eq. (9a) is recommended. In practice, any bare incandescent lamp, frosted or clear qualifies as small. For source areas larger than 1.5 m², but not as large as an illuminated ceiling or uniform indirect lighting (see Section, "Lighting Geometries," for a definition of indirect lighting), UGR is modified into a large room glare rating (GGR) and is described by Eq. (9b). The same UGR and GGR values represent an identical level of discomfort. For very large sources, only the maximum average illuminance values correspond to a specific UGR rating, as shown in Table 2. For nonuniform indirect source, CIE has provided guidelines.²¹ Each of the glare formulations discussed in this section are independent of the light spectrum.

Equations (9a) and (9b) are expressed for a single small and a single large source, respectively. Equation (9) is valid for viewing angles greater than 5° from the line of sight. For a combination of sources of different sizes, Eq. (8) must be modified to include glare from sources specified by equations (9a) and (9b).

$$\text{UGR} = 8 \log_{10} \left(\frac{0.25}{L_b} \right) \sum_i \frac{L_i^2 \omega_i}{P_i^2} \quad (8)$$

$$\text{UGR}_{\text{SingleSmallSource}} = 8 \log_{10} \left(\frac{0.25}{L_b} \right) \frac{200(I^2/R^2)}{P_i^2} \quad (9a)$$

$$\text{GGR}_{\text{SingleLargeSource}} = \text{UGR} + \{1.18 - (0.18/\text{CC})\} 8 \log [2.55\{1 + (E_d/220)\} / \{1 + (E_d/E_i)\}] \quad (9b)$$

where L_b = average luminance of the field of view without the luminaire or glare source

L_i = luminance of the i th luminaire in the observer's direction

ω_i = solid angle of the i th luminaire subtended at the observer's eye

P_i = Guth Position index²² of the i th luminaire. [It is a function of angular deviations (vertical and horizontal) from the line of sight and valid up to 53° of deviation from the line of sight. See Eq. (10).]

I = luminous intensity of the small source expressed in lumens per steradians. The source must be >5° away from the line of sight

E_d = direct illuminance at the eye due to the source

E_i = indirect illuminance at the eye = πL_b

CC = ceiling coverage = (area projected by the source at nadir)/(area lit by the source)

$$P_i = \exp\{(35.2 - 0.31889\alpha - 1.22e^{-2\alpha/9})10^{-3}\beta + (21 + 0.26667\alpha - 0.002963\alpha^2)10^{-5}\beta^2\} \quad (10)$$

where α = angle of the plane containing the observer's line of sight and a line from the observer to the source from the vertical direction. [Vertical direction is the height above (orthogonal to) the floor on which the observer is positioned.]

β = angle between the observer's line of sight and the line from the observer to the source

For roadway lighting, the glare rating is affected by driver fatigue, vehicular speed, and whether the person in the car is driving or not. All these effects must be dealt with comprehensively to formulate a single glare rating model that is largely driver independent. There is ongoing research in this field to develop better models for evaluating glare for road and vehicular lighting.^{23,24} Current CIE recommendation for limiting the discomfort glare for road lighting is identical to disability glare as described by Eq. (7).

Veiling reflections Veiling reflections are reflections from the task surface that result in the luminance contrast [Eq. (1)] reduction of the task itself. Veiling reflections can be evaluated from the source-task-eye geometry. It has been found that 90 percent of test subjects find a luminance contrast reduction of 25 percent as acceptable.²⁵ Veiling reflections are sometimes used as highlights to reveal the specular nature of a display (a lit object).

Miscellaneous Design Issues Other issues include designing against unacceptable light pollution or trespass, light-induced degradation or damage of objects and flicker from light sources. UV and IR filters are used with luminaires when there is a potential of damage to artwork or object displays. Flicker is more noticeable with high levels of the percentage modulation, area of visual field that is impacted by it, or the adaptation luminance. The impact of flicker can be reduced by using high frequency electronic ballasts or multiphase power supplies for different sources.

Functions of Lighting

There are four functions of lighting: ambient, task, decorative, and accent.²⁶ For each function, there are several implementation geometries. For most applications, more than one of these functions is necessary. We first discuss each of these lighting functions and then the lighting geometries used to accomplish these functions.

Ambient lighting fills up the space and is integral to almost any lighting scheme and yet is commonly ignored. It reduces the difference between the magnitudes of vertical and horizontal illuminance. As a result, it reduces glare, softens shadows, and provides a well-lit appearance. A common mistake is to consider any light that is illuminating the space as ambient light. For example, ceilings with recessed down lights with a narrow angular spread cause harsh shadows of objects on the ground. Even facial features show unflattering shadows. This is a result of insufficient level of vertical illuminance. So although there seems to be enough light to illuminate the space, it does not achieve a proper balance between vertical and horizontal illuminances resulting in cast shadows. To achieve proper ambient illumination, it is necessary to use those lighting geometries that spread light into large angles and from many directions. Ambient lighting is provided by large overhead luminaires with diffusers, torchieres, wall sconces, cove lighting, cornice lighting, valence and wall slot lighting, illuminated ceilings and wall washings. Figure 2 illustrates some of these schemes. Figure 4 illustrates ambient lighting with wall sconces.

Task lighting is used to provide sufficient illumination at the task plane such as a desk or work plane. Task lighting should be free of glare and shadows caused by the illuminated objects such as shadows from the hand and body or shadows from machine parts. Several lamp types and lighting geometries are available to reduce the impact of shadows. Lamps such as Banker or Bouillotte (Fig. 24) or large overhead luminaires with diffusers are good choices. The light emanating from these lamps is spread over a wide range of angles from an effectively large source. In a Banker lamp, a significant portion of light from the source undergoes multiple reflections from the inside of the luminaire before exiting. Large fluorescent light sources in a vertical configuration in the Bouillotte lamp provide excellent vertical and horizontal illuminance. Lamps with batwing lenses achieve cross illumination or lighting from two different directions overlapping in the task region. A batwing lens has a linear prism array (Fig. 23c) at the front of the source that leads to spreading of the light in predominantly two directions from each source point. Side lights installed in the vicinity of the task region increase vertical illuminance, and when used with overhead lighting provide excellent task lighting. Sometimes backlighting is necessary for certain tasks that involve transparent objects. Task lights must have

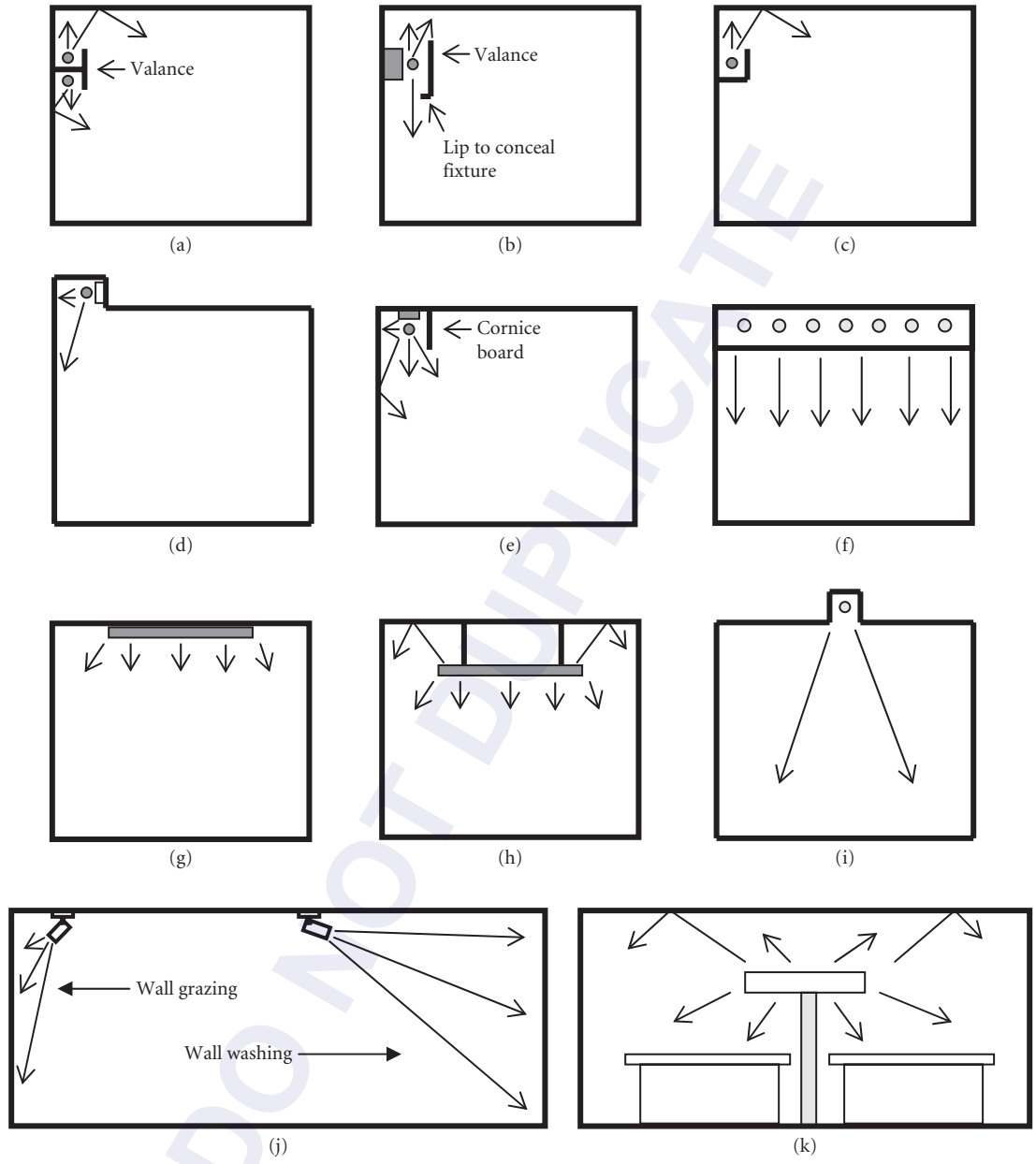


FIGURE 2 Various lighting geometry components. (a) and (b) Valance lighting. (c) Cove lighting. (d) and (e) Wall slot and Cornice lighting for wall illumination. It can be used to create an effect of a floating ceiling. (f) Illuminated ceiling. There is a diffuser below the row of line sources. This scheme can also be used to create illuminated wall panels. Backlighting is another technique. (g) Overhead luminaire. (h) Suspended luminaire. (i) Recessed downlight. (j) Wall-grazing and wall-washing illumination emphasizes and flattens the wall textures respectively. (k) Furniture-integrated lighting system. In cases (a) to (i), the sources are line sources such as fluorescent tubes along the length of the wall. In each case, the neighboring surfaces whether wall or shielding surfaces are coated with high reflectance diffuse paints. The distance of light source from the ceiling/wall determines the uniformity of illumination across the respective regions.



FIGURE 3 Accent lighting. (See also color insert.) (Courtesy of Pegasus Associates Lighting, www.pegasusassociates.com.)

the light sources well shielded by the luminaire and baffles to avoid direct glare. For large overhead sources, louvers (Fig. 24e) are used to limit the direct source view. Veiling reflections are avoided by ensuring that the source-task-eye geometry does not direct the reflections off the task into the eye.

Decorative lighting, as the term implies, is used to add sparkle to the lit environment. Decorative lighting is most effective when it appears to provide most of the lighting in a scene. Using decorative lighting to provide other functions of lighting by increasing the lamp brightness is not a suitable solution, although it increases the illumination in the region. It draws too much attention toward the lamp itself and creates undesirable levels of source brightness against the background leading to glare and unsightly shadows; therefore, decorative lighting must be immersed in an environment with good ambient light. Decorative lighting is provided by low-wattage table or floor lamps, gas lights, chandeliers, sconces, bare light sources, backlighting, light art, and torchieres. See the use of chandelier in Fig. 5.

Accent lighting primarily provides highlighting of objects such as artwork (Fig. 3), plants and decorative objects within a lit environment. Track lights, adjustable recessed lights, uplights, and backlights are commonly used to provide accent lighting.

Lighting Geometries

Lighting geometries can essentially be broken down into four broad classifications. These classifications are not mutually exclusive.

1. **Direct Lighting:** Most of the light (>90 percent) from a luminaire is targeted toward a certain region such as in downlighting in an office by overhead light fixtures. Applications of direct lighting include all the lighting functions such as ambient, task, decorative, and accent.
2. **Indirect Lighting:** An object or a region is illuminated by light that is not directly coming from the source. For example, the light from a lamp is directed toward the ceiling, wall, or even a



FIGURE 4 Wall sconces for providing ambient lighting and the much needed vertical illumination in various situations. (See also color insert.) (Courtesy of Lightcrafters Inc., www.lightcrafters.com.)

diffuse reflecting region of the luminaire. The reflections illuminate the space around the luminaire. Indirect lighting tends to give a more spacious appearance and eliminates shadows. Its primary application is to provide ambient lighting. See Figs. 4, 5, and 24c.

3. **Diffuse Lighting:** When lighting does not appear to come from any specific direction. Examples of diffuse lighting include indirect lighting and certain direct lighting geometries such as overcast skies, large area lighting fixtures with diffusing or prismatic optics or large spatial extent fluorescent lamps. It is mostly used for ambient illumination but can also be used to create local areas of high and uniform brightness for task lighting or accent lighting.
4. **Direct-Indirect/Semi-Direct/Semi-Indirect Lighting:** These terms typically apply to lamps that distribute some portion of their light toward the target and the remaining portion toward a surface that reflects (specular and/or diffuse) light toward the target. Semi-direct typically refers to the case where 60 to 90 percent of the light is directed toward the target while semi-indirect refers to the case where 60 to 90 percent of the light is directed away from the target. Applications include all the lighting functions: ambient, task, accent, and decorative. See Fig. 24d for an example of direct-indirect lighting fixture.

Figure 2 illustrates several implementations of these lighting geometries⁵ within the confines of indoor lighting. A practical lighting layout is likely to include one or more of these concepts.



FIGURE 5 Indirect lighting with cove lighting in a restaurant using light strips. The chandelier provides the decorative lighting without significantly contributing to any other lighting function. (See also color insert.) (Courtesy of Pegasus Associates Lighting, www.pegasusassociates.com.)

Properties of Objects and Their Impact on Lit-Scene

Objects in a lit-scene have geometrical (shape, location, and orientation) and optical properties (wavelength dependent absorption, reflection, transmission, and diffraction). Reflection and transmission include both diffuse and specular components. Lit-scene characteristics are then in-part determined by geometrical and optical properties of the objects and also in part determined by the interaction of these characteristics with the light sources in the environment. The object-light source interaction gives rise to such phenomena as glare, nonuniform illumination and color change. Major indoor features that impact a lighting environment are walls, floors (including large area carpeting) room partitions and ceilings. Minor features consist of temporary objects such as wall coverings, furniture, partitions, art, displays and plants. Major outdoor features that impact lighting environment are ground, buildings, vegetation, distant and close landscape features. Minor outdoor features are temporary objects.

System Layout and Simulation

Based upon the lighting design criteria, appropriate lamps with desired lighting schemes are selected. However, to obtain the desired illuminance distribution over time, appropriate calculations and simulations are needed. In this section, we discuss the tools for system simulation.

For any given system layout, a certain light level is needed. Equation (11) provides an approximation of the number and type of luminaires needed to obtain a certain horizontal illuminance over a work plane. This equation is useful when expressed as a summation of individual luminaires with their specific constants.

$$E_{\text{maintained}} = FnLLFCU/A \quad (11)$$

where $E_{\text{maintained}}$ = average illuminance maintained
 F = total rated luminaire lumen
 n = number of luminaires

CU = coefficient of utilization (It defines the percentage of light from the lamp reaching the work plane. CU depends upon the relative placement of lamp and work plane and illuminance distribution at the work plane corresponding to the geometry.)

A = work plane area

LLF = total light loss factor

LLF²⁷ accounts for the lamp output reduction over time. LLF has both recoverable and nonrecoverable components. Nonrecoverable factors are due to permanent degradation of the luminaire surface, reduction of output due to deviation from ideal operating temperature and environment (convection and ambient temperature), inefficiency of the electronic drive components and deviation in the operating position (tilt) from the ideal position. Recoverable factors are those whose effects can be mitigated by regular cleaning of the luminaire, operation in a clean environment, and regular replacement of bulbs or sources upon their natural degradation with time. LLF is a product of all these factors. Many of these factors also affect the intensity profile of the lamp output and can therefore affect the CU. The *Lighting Handbook* by IESNA²⁸ lists a detailed procedure for calculating the room surface dirt depreciation factor and luminaire dirt depreciation factor. Other factors can be obtained either by lamp manufacturer specification data or by measurement in an as-used configuration of each lamp before use.

System simulation can be done in two ways: manually or using specialized software. To do it manually, the lumen method [see Eq. (11)] and the zonal cavity method are used. The zonal cavity method²⁹ involves modifying the CU by calculating the effects of room geometry, wall reflectance, luminaire intensity, luminaire suspension distance and workplane height. The impact of various parameters is available in the form of look-up tables that can be consulted to provide estimates. These methods are quite powerful but outside the allowable space in this chapter, please consult the mentioned references.

System simulation with specialized software allows unparalleled accuracy and flexibility. Modern software tools allow easy modeling of 3D geometries, material and source properties. Sophisticated analysis tools allow for calculation of luminance, intensity, illuminance and chromaticity at any location, and they also provide photorealistic rendering of lit models that account for specular as well as diffuse reflections. Accurate system modeling involves the following steps:

1. Model the 3D geometry of the lit region. Windows, skylights, and light shelves must be modeled with their coverings: blinds or glazings with their optical properties.
2. Model the surfaces and paints. The reflectance is a combination of specular and diffuse components. For those surfaces and paints that cannot be simply described by specular or Lambertian properties, the bidirectional reflectance distribution function (BRDF) is used.³⁰ The dependency of BRDF is composed of the incident direction and the direction of observation. This concept has been explained in much more detail in Chap. 7, "Control of Stray Light," in this volume. BRDF measurements are mostly obtained experimentally. These measurements are available sometimes by paint vendors or makers of specific surfaces. It is also possible to simulate BRDF approximately by assigning texture and reflective properties to the surface and performing a ray trace. Once the reflectance properties of each surface are available, they are assigned to the surfaces in the optical model and allow accurate photorealistic rendering of the model for all cases of source and viewer locations.
3. Model the lamps: luminaire geometry and source model. Source models are increasingly being made available by the lamp vendors. If the model is not available, source measurements may be needed. Source models consist of a collection of rays that represent the output from the source. Each ray is described by its position and direction cosines in 3D space. Daylight can be simulated by creating two sources outside the model: sun and sky and assigning diffuse reflective properties (10 to 20 percent) to the ground outside the model. There are a variety of ways to model the sun. One way is to represent it as a Lambertian disk of its angular extent (0.52°) and a luminance value that provides the insolation on the earth's surface at the geographic location ($\sim 1000 \text{ W/m}^2$). The final step is to locate the Sun's position relative to the model. The sky is modeled as a large Lambertian disk of a prescribed luminance. For example, a clear sky is represented by 8000 nits and an overcast sky is represented by 2000 nits.
4. Decide upon the location of the measurement surfaces. These could be real or virtual. A photorealistic scene rendering helps in realizing the most promising layouts.

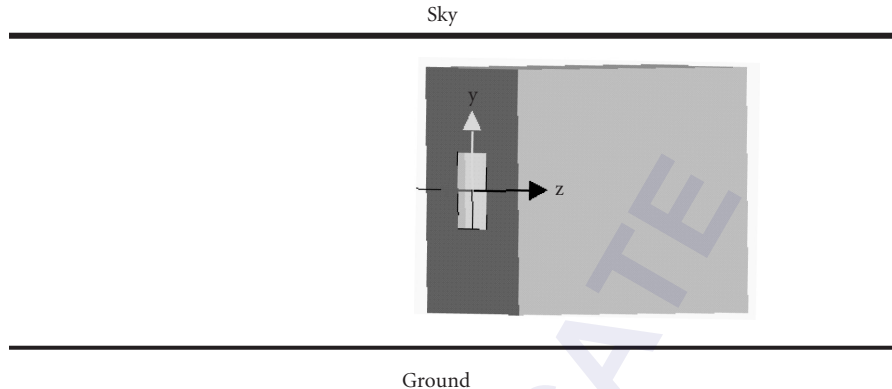


FIGURE 6 System layout. (Source: Clear sky 8000 cd/m^2 . Room Size: $3 \text{ m} \times 3 \text{ m}$. Window size: $1 \text{ m} \times 1 \text{ m}$. Window location: center of the wall. Wall reflectance: 35 percent. Roof Reflectance: 60 percent. Floor reflectance: 40 percent. Ground reflectance exterior to room: 20 percent.)

Figures 6 and 7 show an example that illustrates the above steps.³¹ We simulate a room with one window. All the light received in the room is daylight from a clear sky. In this example, we calculate the light distribution across the floor for two cases: with and without assuming that the room surfaces are diffuse reflective. Figure 7 shows the impact of including reflectances from various objects which results in a wider spread of illumination across the floor.

It is easy to make the model more complex by adding internal light sources, objects with a variety of surface properties, skylights, influence of sun and so forth. Next, we discuss the software tools available to perform such simulations.

Software Tools There are extensive software tools to assist the lighting designer in simulating illumination systems. The software includes computer-aided design (CAD), source modeling, optical analysis and design, and computer graphics. Each of these plays a crucial role in the design process, so they are described in the following subsections. Some of the software areas have applicability in a number of aspects of the design process. Finally, we do not mention the explicit names of the software packages since they are continuously evolving, and we make use of certain ones in our daily lives, so we do not want to bias the discussion presented here.

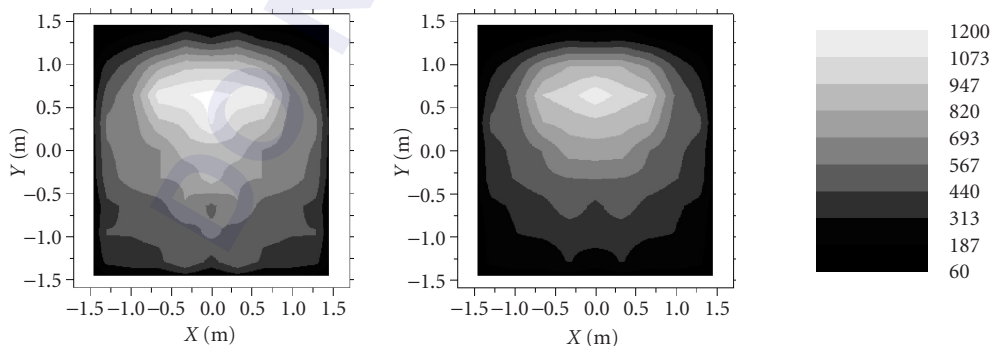


FIGURE 7 Horizontal illuminance (lux) on the floor (a) with and (b) without taking into account the room reflectance.

CAD software Computer-aided design software is used to build the geometry of any system and is then interfaced with machine tools to fabricate the components. Lighting design makes great use of CAD software to ease the process of integration of complex optical components and the mechanics that hold them and the electronics. Not only are there self-standing, mechanical CAD software packages, but native CAD geometry generation capabilities in the optical analysis software packages discussed later. The latter provide a wealth of tools for the generation of complex illumination systems, but they do not have the range of tools that mechanical CAD software provides. The software can be broken down into two subsets: surface based and solid based. The former implies that each surface is defined separately (e.g., a cube is made up of six separate surfaces), while the latter implies that each object is defined (e.g., a cube is made from one function call). Of course tools are provided in each code such that the design process is simplified (e.g., in a surface-based code a macro could generate all of the six surfaces of the cube with one function call). Solid-based codes tend to provide a more efficient process to enter the geometry of the system, but surface-based codes have a longer market history. None of the mechanical CAD packages provide optical analysis and design tools, but they do provide hooks to integrate into them. In fact, some of the optical analysis and design software companies have developed plug-ins that allow the user to specify optical characteristics such as materials and surfaces. These tools assist with the design process by requiring only one iteration of assigning such properties and simplify the transfer of the geometry to the optical analysis software. The transfer of the geometry is typically accomplished by two formats: International Graphics Exchange Specification (IGES) and Standard for the Exchange of Product model data (STEP). Other protocols including proprietary ones, DXF, DWG, STL, and SAT are understood by certain optical design and analysis tools.

IGES IGES is a standard first published in 1980 by the National Bureau of Standards (now the National Institute of Standards and Technology, NIST) as NBSIR 80-1978, and then approved by an ANSI committee as Version 1.0.³² It was first known as Digital Representation for Communication of Product Definition Data. It is essentially a surface/curve-based method to represent the geometry that comprises a component or system. The IGES standard was updated through the years, with Version 5.3 in 1996 being the last published version.³³ STEP (see the next section) was to replace IGES, but IGES remains the prevalent neutral-based method to transfer geometry data. Most optical analysis codes can read and interpret various versions of IGES, but the surface-based optics codes tend to have better performance (i.e., fewer import errors).

STEP STEP is a standard first published in 1994 by the International Standards Organization (ISO) as ISO 10303 with the goal of replacing IGES. It is essentially a solid-based method to represent a system, but it also has 2D and database aspects.³⁴ It has been updated through the years and is now comprised of many part and application protocols. STEP is developed and maintained by the ISO technical committee TC 184, Technical Industrial automation systems and integration, subcommittee SC4 Industrial data.³⁵ Most optical analysis codes can read and interpret various versions of STEP, but the solid-based optics codes tend to have better performance (i.e., fewer import errors).

Source modeling software In order to perform accurate simulation of a lighting system it is imperative that an accurate source model is used. There are essentially three methods to accomplish this:

- Generation of ray data based on manufacturer data sheets
- Experimental measurements of the emitted radiation to create ray sets
- Modeling of the geometry of the source and the physics of emission to create ray sets

The first method creates either ray sets or Illuminating Engineering Society (IES) intensity distributions that can be used to model the performance of a lighting system that incorporates the prescribed source. The second, experimental measurement uses a goniometer or similar device to measure the output of the source both as a function of position and angle; therefore, it measures the luminance distribution of the source. The last is based on modeling of the source components, such as filament, base, and glass envelope in an incandescent bulb; electrodes, glass envelope, and

base for an HID lamp; and die, epoxy dome, and reflector cup for a LED. Based on the physics of the source, rays are assigned to the emission areas. A number of optical analysis software companies are providing geometrical and ray source model libraries to their customers. All three methods are based on the data, measurements, or model of a single source, called the nominal source. Thus, there is the opportunity for error between the nominal source and what is used in your fabricated system.

Source manufacturers and architectural lighting software tend to use the IES source files to specify their sources. These files are not as precise as the other two methods, but they provide a good enough and fast method to implement the source emission characteristics into the design process. Companies that make these accurate experimental measurements keep libraries of the data, so that their customers can include them in their lighting system models. The CAD software allows a designer to make a complex model of the source, and then the optical analysis software allows the rays to be generated. This method also provides accurate source models, but with the only time expenditure to develop such models. It has been found that methods that employ both the geometry and experimental emission measurements provide the highest level of accuracy while also giving an avenue to model the tolerances of the emission.^{36,37}

Optical modeling software As previously stated, optical analysis and design software not only allows for the modeling of optical systems, but it also provides tools for inclusion of geometry, source modeling, and rendering (discussed later). There are both solid-based and surface-based codes. Most optical phenomena occur at the surface interfaces, such as reflection, refraction, scattering, and diffraction, but there are volume effects, such as scattering, absorption, and emission. Therefore, though a single software package is based around solids or surfaces, it must be able to effectively model the other type of phenomena. A number of the codes are generic in nature, such that they can handle virtually any type of system or application, from backlights to luminaires to biomedical applications. There are application specific codes in a number of areas, especially external automotive lighting and architectural illumination. The software is increasingly adding tools such as source modeling macros, optimization, tolerancing, and rendering.

There are essentially three types of software packages that can be used to model a lighting system: optically based ray tracing, lighting-based radiosity, and computer graphics rendering. Ray tracing is simply the tracing of a multitude of rays from sources through the optical system. It is quite accurate, only limited both by the characterization of the geometry and the assigned optical properties within the model and the number of rays that are traced. Radiosity algorithms are essentially scatter-based methods that propagate approximate wavefronts from one object to another. Initially, Lambertian scatter properties of all objects were assumed, but more recently radiosity implementations propagation based on the bidirectional surface distribution function (BSDF), generally, or the respective reflective (BRDF) or transmissive (BTDF) forms, specifically, have been developed.³⁸ Ray tracing can handle both specular and diffuse reflections, but radiosity is limited to diffuse reflections. Ray tracing has a number of benefits including accuracy and utility from the near field to far field. Radiosity is for the most part limited to far field calculations, where the approximations of the propagation model are minimized. As the distance between the source and target is reduced, the limitations of the scatter-based propagation inhibit accuracy. The primary limitation of ray tracing is the calculation time, which is several orders of magnitude more than radiosity. Thus, hybrid methods that employ both ray tracing and radiosity are in use. The goal of hybrid methods is to obtain the benefits of both ray tracing and radiosity in a single algorithm.

In the next three subsections the three types of software packages are discussed in more detail. Notably, the lines between these three types of software packages is disappearing, especially between the lighting design and computer graphics sectors. In each of these sections the applicability to lighting design and modeling is provided. These software packages are seeing rapid growth, so consultation of the literature on active research on future development is suggested.

Optical design and analysis software There are two types of optical design and analysis software: imaging system software and general analysis software. The first is typically called lens design software, and has limited utility to the design of lighting systems. The second uses nonsequential ray tracing from the source to the target. This process allows the illumination distribution at the

target to be accurately determined. This type of software often includes, at the discretion of the user, such features as spectrum, coherence, polarization, and so forth. Thus, the accuracy is only limited by the user input. The design of the actual luminaire is best done with this type of software. It allows the designer to design efficient systems that effectively couple and broadcast the emission from the light source. These codes also provide tools for optimization and tolerancing of the luminaire. However, due to lengthy computation time, optical design and analysis codes have limited (but increasing, due to the advances in computer speeds) utility in determination of a lit scene (i.e., rendering).

Lighting design software Unlike optical design and analysis software, lighting design software typically makes use of radiosity algorithms.³⁹ Radiosity codes are quite fast and have acceptable accuracy in the far field. The use of Lambertian or BSDF scatter properties allows the diffuse reflection from objects to be quickly ascertained and propagated further into the system. The diffuse emission is both collected at the observation location and is cascaded to other objects in the scene. In order to obtain a higher convergence speed objects are typically parameterized with polygons, while optical phenomena such as refraction or specular reflection are approximated or even ignored. This type of software thus provides at worst a first-order approximation of the lit appearance such that architects and lighting designers can view the results of their design work. Lit-scene rendering or illumination in the far field of a luminaire is the best use of lighting design software. More advanced software packages can then be employed, such as those that employ some semblance of ray tracing—see the previous and next sections.

Computer graphics software In the past two decades the computer graphics community has grown rapidly. The tools they develop and employ are useful in the illumination community of optics. Foremost they have tools to model the simulated look of an unlit or lit lighting system, called unlit-and lit-appearance modeling respectively. These tools are important in illumination since acceptance of a system is often based on subjective criteria such as appearance. Thus, these tools provide such before potential costly and time-consuming manufacture. Additionally, the computer graphics community uses both ray-based and scatter-based radiosity methods, while also employing hybrid methods. These methods are especially geared to the rendering of scenes in video games, movies, and other types of visual media. Thus, they tend to have the least amount of accuracy since the completion of numerous images in a timely manner is demanded.

Computer graphics software makes direct use of forward ray tracing (from the source to observer) and reverse ray tracing (from the observer to the source). The latter is especially useful for the rendering of scenes where there is a discrete viewpoint, like that of a virtual observer. These codes and algorithms are increasingly being used to model the lit appearance of illumination systems such as luminaires and lightpipes. This process involves some form of ray tracing and/or radiosity calculations and then employs vision biology (Sec. 40.3). As an example, consider a star-shaped taillight as shown in Fig. 8.^{40,41} The taillight is oriented at several angles such that the effects of changing ones aspect to the lit taillight is taken into account. The combination of these different angular views of the taillight provides the luminance distribution, or essentially what an observer would see by walking around the lamp. A ray tracing method employing a pupil collection of 15° is used for each of the plots within Fig. 8. Note that saturation of the retinal cones is included, which is evidenced by the whitish appearance of the filament at the center of the lit patterns.

Furthermore, the resulting intensity pattern for the lit-appearance model can be projected into a scene to provide a rendering of what the illumination from the lamp will look like. For example consider Fig. 9, which shows three view aspects of an automobile headlight designed to meet standards (see section on vehicular lighting): (a) the driver's perspective, (b) 20 m above and behind the drive, and (c) the bird's eye view.⁴² These renderings are quite accurate since the goal is to completely mimic the lit-scene appearance prior to fabrication. These types of renderings can be extended to any scene that involves sources and objects with accurate optical characteristics applied to them. Figure 10a shows the rendering of a south-facing office room⁴³ in Tucson, Arizona. The illumination was modeled from average direct and diffuse insolation data for November 15 at noon for this location. The CAD model was generated from architectural blueprints of the facility. Surface reflectances were determined from first principles. The scene outside the window was created from a

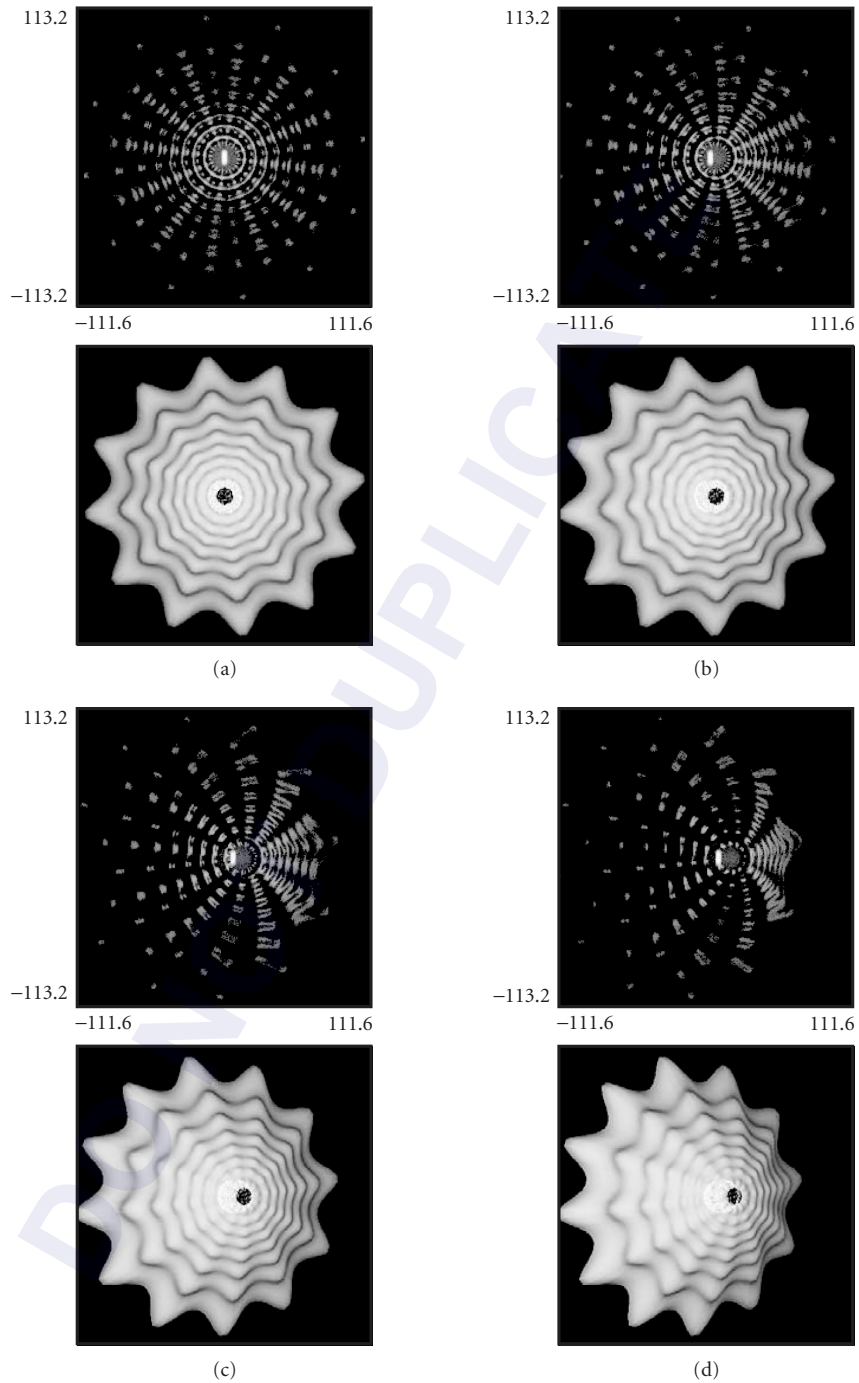


FIGURE 8 Views of the lit appearance (upper) of a star-shaped taillight (lower) at four horizontal angles of (a) 0° ; (b) 10° ; (c) 20° ; and (d) 30° . (See also color insert.) (Used with permission from SPIE;⁴⁰ Developed with Advanced Systems Analysis Program from Breault Research Organization.)

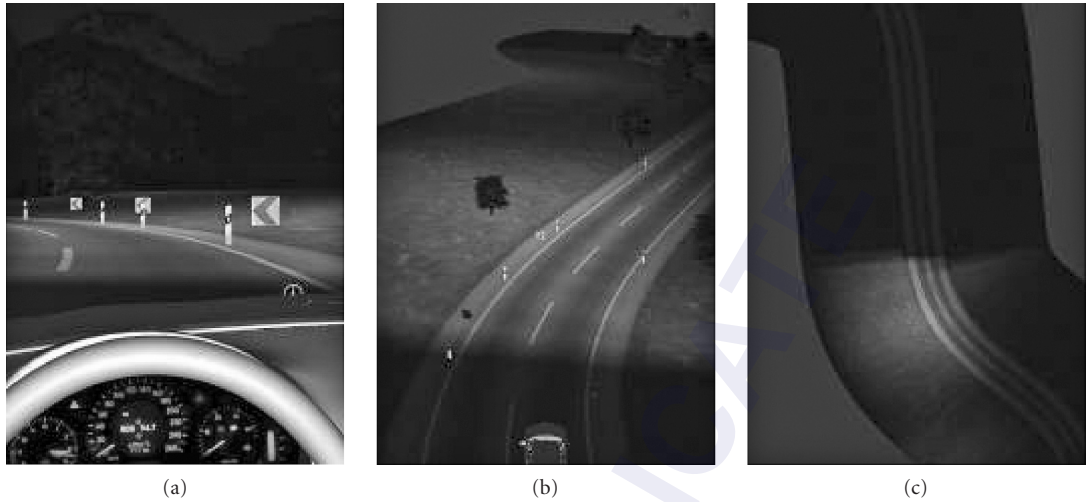


FIGURE 9 Three perspectives of lit-scene renderings from a low-beam headlamp: (a) driver's view; (b) 20 m above and behind automobile; and (c) bird's eye view. (See also color insert.) (Developed with *LucidShape* and *LucidDrive* from Brandenburg, GMBH.)

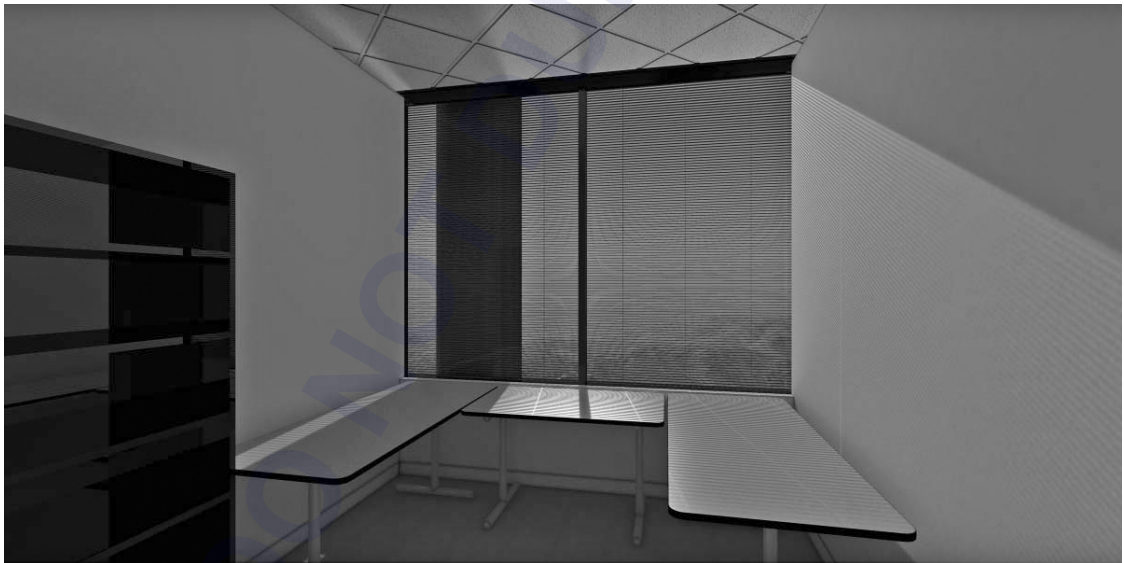


FIGURE 10a Rendering of a lit office room. (See also color insert.) (Developed with *LightTools* from Optical Research Associates.)

digital photograph taken in the mountains outside of Tucson. 100,000,000 rays were traced from the source throughout the model. Similarly, Fig. 10b shows the rendering of a desk surface lit by interior, incandescent lighting.⁴⁴ All surfaces, except the three objects on the desk (shown in wireframe to ease view through the objects), are diffuse Lambertian reflectors. The three objects: wine glass, ice cube, and crystal ball, display the effects of specular refraction and total internal reflection.

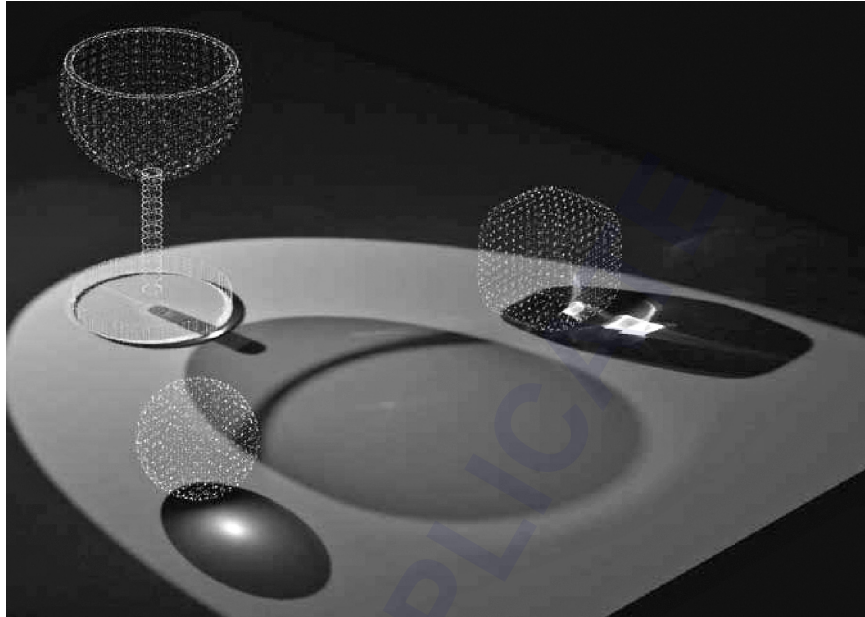


FIGURE 10b Rendering of a lit desk with three objects located on it (wine glass, ice cube, and crystal ball) to show both diffuse and specular effects. (See also color insert.) (Developed with *FRED* from Photon Engineering.)

40.5 LUMINAIRES

A luminaire is a packaged light source consisting of emitter, optics to baffle or redirect light, fixture and electrical components. Luminaires form an integral part of the lighting design by the virtue of the illumination they provide and also by their appearances. In this section, we discuss the optical components of luminaires: types of light sources, both artificial and natural (daylight), and the design of luminaires with optical components such as reflectors, lenses, lightguides, fibers, windows, skylights and baffles to achieve the desired light distribution.

Types of Light Sources

Light sources are optically characterized by their luminous spectrum (lumens as a function of wavelength), intensity (in candela), and efficiency (lumens per watt, lpw). For white light, CCT and CRI are derived from the spectrum of the source and are typically reported on the data sheets to give an estimate of its appearance and its ability to color render lit objects. Other source characteristics are cost, safety, government regulations, package size/type, lamp sockets, electrical driver requirements and constraints related to environment or operating conditions. Cost is defined in terms of luminous flux per dollar per hour of use, replacement, and initial installation costs.

We discuss daylighting and the following artificial light sources: incandescent, fluorescent, high-intensity discharge (HID), light-emitting diode (LED), electrodeless HID, electroluminescent, nuclear, laser, bare discharge, low-pressure sodium (LPS), and short arc sources. In the next few subsections we discuss various light sources in terms of operating principles, construction, and packaging. Refer to Table 3, which provides the performance comparison of various lamps; and Table 4, which

TABLE 3 General Lamp Characteristics for Most Lighting Applications

	Watts	Efficacy lpw	CCT K	CRI	Lifetime K hours	Notes
Standard Incandescent	40–100 ^a	10–17 ^T	2700	>95	0.75	Undesired IR radiation, fire hazard. Naturally low flicker, instant-on and dimming to 0.
Tungsten-Halogen	300 ^a	20	2850–3200	>95	<6	Same as standard incandescent.
Fluorescent/CFL	5–55	15–100 60–70 ^T	3000–6500 4100 ^T	50–98 70–85 ^T	5–20 ^b	Complex ballasts for good dimming range, short start-up, low hum, low flicker. Hg disposal issue, EMI.
HID-Hg	50–1000	30–65	3900–5700 ^c	15–20 higher with phosphor coatings	16–24	Complex ballasts for start-up and flicker control, poor dimming, high flicker, explosion, fire and UV hazard, lamp disposal issue, high start-up and re-strike intervals, poor color stability with time, operating position affects performance.
HID-MH	30–18 K 50–1000 ^T	75–125	2500–6000	60–70	7.5–20	Same as HID-Hg. Lifetime <1 K hours for >10 kW lamps.
HID-HPS	175–1000	45–150	1900–2700	22–85	7.5–24	Same as HID-Hg, CRI is inversely proportional to efficacy. 110 lpw corresponds to CRI 20, CCT 1900–2100 K.
HID-CMH	20–400	70–90	3000–4200	80–96	7.5–20	Same as HID-Hg except for good color stability with time and performance independent of operating position.
Electrodeless	4	60+	2700–6500	50–98	15–30 ^d , <100 ^e	Complex ballast, EMI.
LPS LED^f	20–100	80–150	1800	0	14–18	LEDs (large chip or arrays) are likely to replace most of the existing lamp sources. Efficacy can reach up to and beyond 200 ^g lpw with the theoretical limit being the CIE standard observer luminous efficacy curve, lifetimes up to 100 K hours. Key advantages are: fast turn-on times, color tunability, high dimming range, low voltage operation, no Hg or Lead, no UV or IR, no catastrophic failures and scalable packages. Phosphor coated LEDs can provide white light at desired CRI and CCT.

HID—high-intensity discharge, Hg—mercury, MH—metal halide, HPS—high-pressure sodium, CMH—ceramic mercury halide, Xe—xenon, LED—light-emitting diode, CFL—compact fluorescent lamps, ^T—typical.

^aAlso available in kW. Lifetimes shrinks to a few hundred hours, ^bRegular Fluorescent have typical lifetimes >10 K hours, Compact Fluorescent lamps have lifetimes >5 K hours, ^cCCT 5700 K at CRI 15, CCT 3900 K at CRI 50, ^ddeveloping technology, ^ewith integrated ballasts, ^fseparate ballast, ^gonly light generation efficiency inside the LED die.

lists various lamp types, currently available packages and their applications. Figure 11 shows various lamp packages. The alphabetic designation is explained in Fig. 11. The numeric designation with the letters is the diameter specified in 1/8th of an inch. For example, MR16 refers to a multifaceted reflector, 2 in. in diameter.

Incandescent Sources These are thermal sources that emit electromagnetic radiation from a heated filament, which is a phenomenon known as *incandescence*. Incandescent light sources are being steadily replaced with fluorescent lamps, LEDs, and potentially with electrodeless lamps. Other sources are able to provide similar or improved performance at higher efficiency and lifetime leading to reduced operational costs.

Modern day incandescent sources use a tungsten filament.⁴⁵ Tungsten has a high melting point (3382°C), high ductility, high conductivity, and low thermal expansion that make it a preferred material for use as a lamp filament. Tungsten is alloyed with tiny amounts of potassium (60 ppm), aluminum oxide (10 ppm), and silicon (1 ppm) to give it high strength near its melting point. This allows lamp operation close to the melting point, thus improving its efficacy. Sometimes tungsten

TABLE 4 Common Lamp Packages and Applications

	Available Packaging	Current Applications
Tungsten	A-line, elliptical, decorative (B, C, CA, F, G, M), PAR, reflector (R, BR, ER), appliance and indicators (S), tubular (T)	General purpose, 3-way, reader, decorative (chandelier, Globe, ceiling fan), Track and Recessed (Indoor floodlight and spot light), outdoor (post & lantern, pathway, garden & deck, motion-sensing & security, bug light, yard stake), appliances, colored lamp, display, exit sign, freshwater and saltwater aquarium, heat lamp, marine, nightlight, party, recreation vehicle, plant, rough service, sewing machine, shatter resistant, terrarium, vacuum cleaner, airport, emergency, city lighting, projection, photoflood, filmstrip, retail display, restaurants.
Tungsten-Halogen	A-line, decorative (B, G, F, T10), PAR, AR, MR, single ended, double ended	General purpose, 3-way, reader, decorative (chandelier, Globe, ceiling fan), Track and Recessed (Indoor floodlight), outdoor (post & lantern, pathway, garden & deck, motion-sensing & security, yard stake), camera light, microfilm, curio cabinet, display, enlarger & printer, equipment, fiber optics, landscape lighting, projection, stage & studio, torchiere, airport, emergency, city lighting, special service, monuments, museums, heat lamp.
Fluorescent (Linear)	Straight linear (T5, T6, T8, T10, T12), circular (T9), U-shaped (T8, T12), grooved (PG17)	Kitchen, bath, shop, work light, Appliances, blacklight, blacklight blue, cold temperature, colored lamp, freshwater and saltwater aquarium, plant, shatter resistant, terrarium, stage & studio, projection, diazo reprographic, germicidal, gold UV blocking, superstores, warehouses.
Compact Fluorescent	Plug-in (2-pin, 4-pin, proprietary), self-ballasted (decorative, reflectors, proprietary)	General purpose, reader, decorative (chandelier, Globe, ceiling fan), Track and Recessed (Indoor floodlight and spot light), outdoor (post & lantern, pathway, garden & deck, bug light), appliances, blacklight blue, facilities, hospitality, office, plant, restaurant, retail display, saltwater aquarium, terrarium, torchiere, warehouse.
HID - Hg	A-line, elliptical, reflector	Street lighting.
HID - MH	Elliptical, PAR, single ended, double ended, tubular	Street lighting, sports lighting, decorative lighting of architectural wonders.
HID - HPS	Elliptical, double ended, tubular	Street lighting, horticulture.
HID - CMH	Elliptical, par, single ended, double ended, tubular	Track and Recessed (Indoor floodlight and spot light), retail display.
Miniature	B, G, R, RP, S, T, TL, discharge	Outdoor (post & lantern, pathway, garden & deck, motion-sensing & security, yard stake), automotive (headlamp, fog, daytime running, parking, directional front & rear, tail, stop, high mount stop, side-marker front & rear, backup/cornering, instrument, license plate, glove compartment, map, dome, step/convenience, truck/cargo and under hood), flashlight, landscape lighting, low voltage, marine, telephone, traffic signal, emergency.
Sealed Beam	PAR, rectangular	Outdoor (motion-sensing & security, yard stake), automotive headlamp, railway, shatter resistant, stage & studio, directional lighting, aircrafts, tractors, airport, emergency, city lighting.
LPS Electrodeless	Tubular T, P	Street lighting, parking lots. Any application where long lifetimes and high efficacy are needed due to high replacement costs and/or difficult access. Road signs, warehouses.
LED	MR16, miniature (2, 3, 5 mm)	LEDs can potentially replace most of other light sources being used for various applications. Currently being used in building interiors (homes, offices, commercial places such as health clubs, hospitals rooms) stairways and pathways, flashlights, traffic lights, signs, digital projection, instrument indicator panels, backlighting applications, decorative, display.

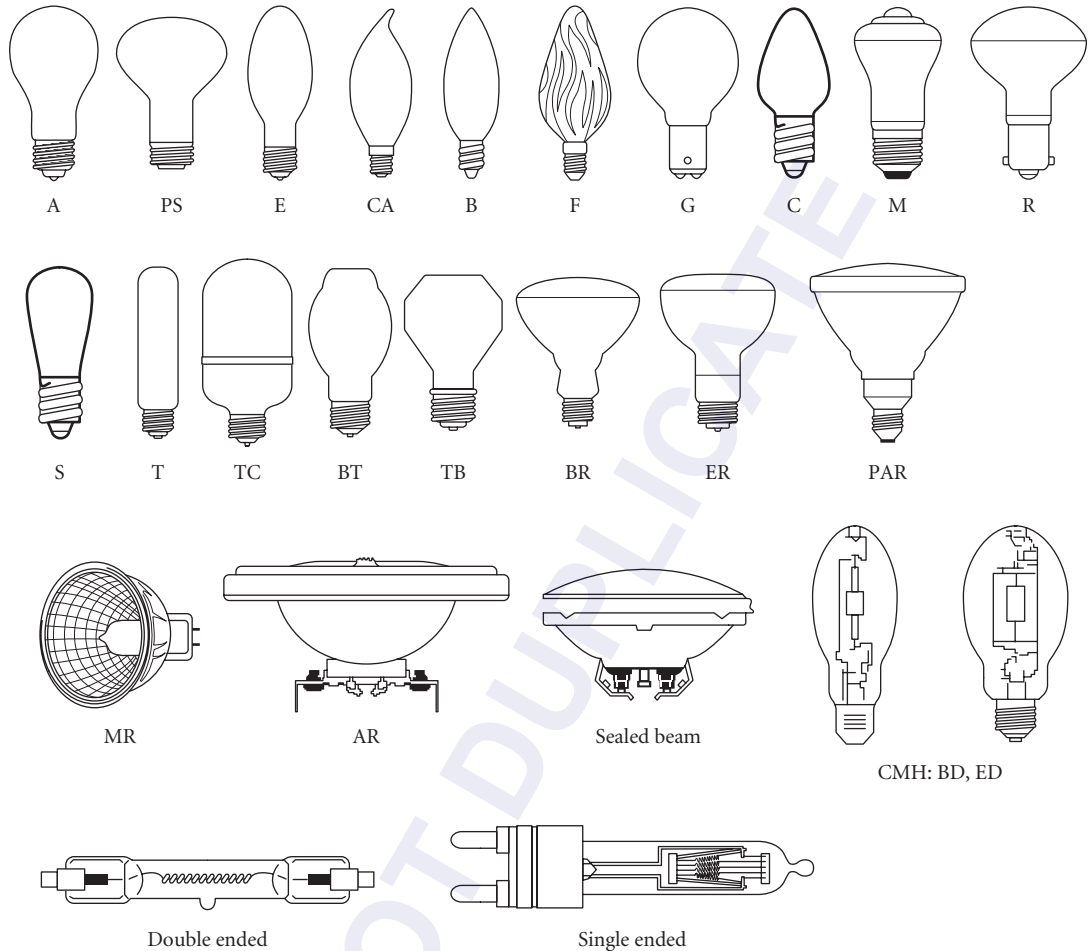
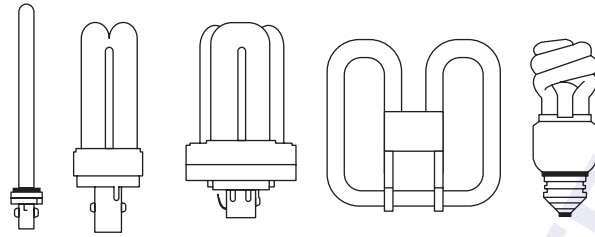


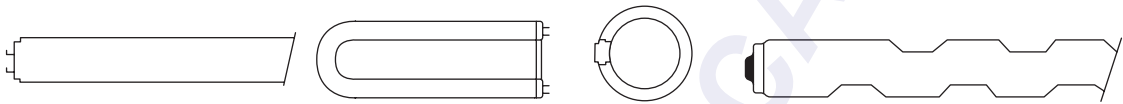
FIGURE 11a Various lamp package representations. Shape and sizes are not to scale. Various sizes are available for each package type. For each bulb shape, a variety of bases are available. A—arbitrary spherical shape tapered to narrow neck, B—bulged or bullet shape, BT—bulged tubular, C—conical, CA—conical with blunt tip, E—elliptical, blunt tip, ED—elliptical with dimple in the crown, F—flame shaped, decorative, G—globe, M—mushroom-shaped with rounded transitions, MR—multifaceted reflector, PS—pear shaped with straight neck, PAR—parabolic aluminized reflector, BD—bulged with dimple in crown, S—straight, T—tubular, TB—Teflon bulb, TL—tubular with lens in crown. (Illustration courtesy of General Electric Company.)

is alloyed with rhenium (3 to 25 percent) to make it more ductile at low temperatures and achieve higher recrystallization temperatures thereby giving the lamp a longer life. Alloying tungsten with thorium provides increased strength, better machinability and high recrystallization temperatures. Such filaments are used for very high voltage applications.

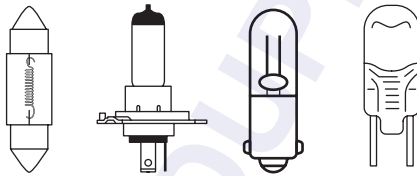
Figure 12 shows the radiating characteristic of tungsten at 3000 K and its comparison with a blackbody emitter. It differs from the blackbody due to wavelength-dependent low emissivity. A perfect blackbody has an emissivity of one for all wavelengths. The hotter the filament, the higher are the luminous flux radiated per watt, the percentage of luminous flux of the total radiation, and the CCT. However, the lifetime is inversely proportional to the filament temperature as filament evaporation is



CFLs: Biax, double biax, triple biax, 2D, spiral



Fluorescent T, U-line, circline and grooved PG



Miniature: neon, festoon, automotive, TL

FIGURE 11b Various lamp package representations. Shape and sizes are not to scale. (Illustration courtesy of General Electric Company.)

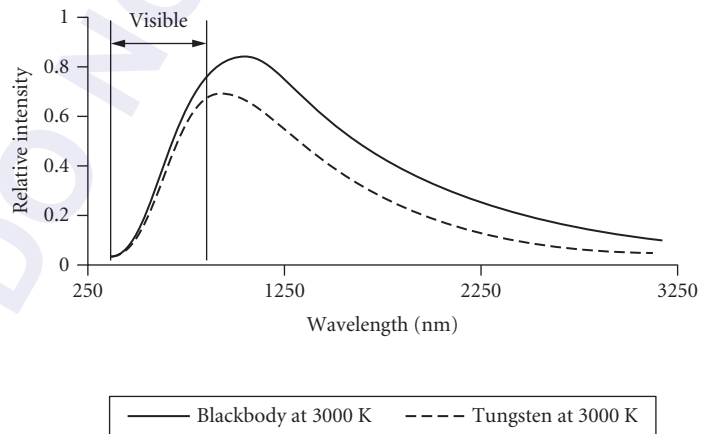


FIGURE 12 Blackbody at 3000 K versus tungsten filament at 3000 K.

the primary mode of lamp failure. The best engineered solution consists of the hottest possible lamp filament at an acceptable lifetime, voltage rating, and packaging. Near its melting point, an uncoiled tungsten wire has a luminous efficacy of 53 lm/W. To achieve an acceptable lifetime, the tungsten filament is operated at much lower temperatures. The luminous efficacy of typical incandescent lamps with tungsten filaments ranges from about 5 to 20 lm/W for lamp wattages 5 to 300, respectively.

A typical incandescent lamp consists of an evacuated glass envelope; filament, with or without fill gases; filament leads; and base. Figure 13 shows the key characteristics.

The bulb material is application dependent. For most applications soda lime glass is used. For applications requiring heat resistant glass, borosilicate, quartz, or aluminosilicate is used. When the bulb envelope is frosted from the inside or coated with powdered silica, it provides diffuse illumination from the bulb surface and masks the bright filaments from direct view. The bulb envelope material can be used to filter the radiation to alter the CCT. Daylight application bulbs filter out the longer wavelengths to provide a higher CCT.

Lead-in wires are made of borax coated dumet (alloys of nickel, copper, and iron). Dumet is able to form a glass-metal seal. It is important to match the thermal expansion of the lead-in wires with the envelope material. When high bulb-envelope temperatures are involved, molybdenum strips bonded to lead-in wires are used for glass-metal seals. The bulb base is cemented to the bulb envelope and is designed to withstand the operating temperatures. The filament itself comes in various configurations depending upon the application. Various filament configurations are a straight wire (designated as S),

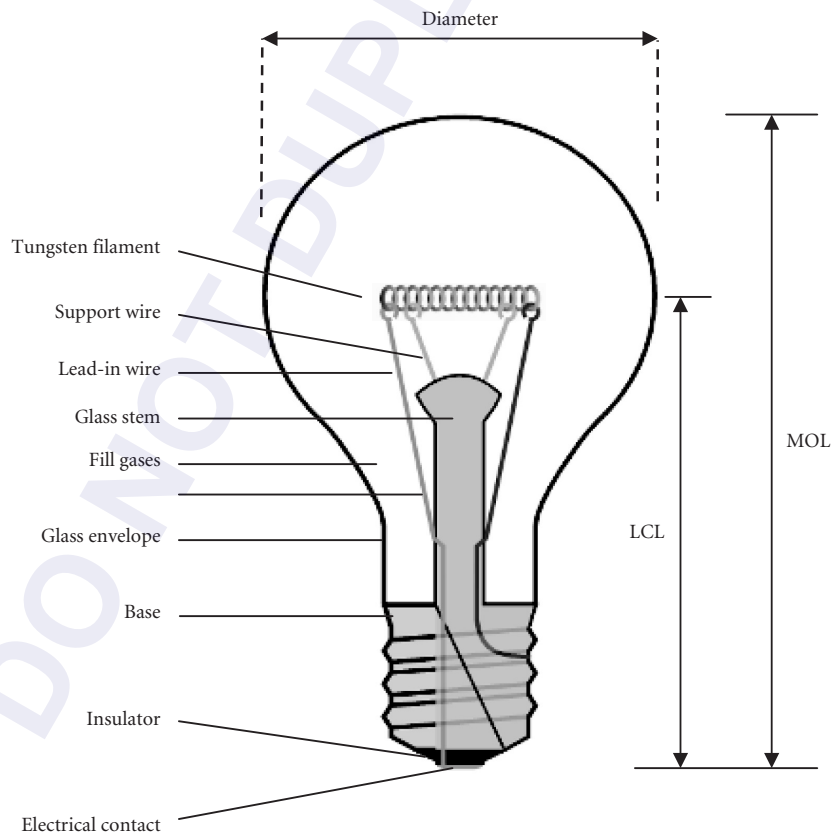


FIGURE 13 An incandescent bulb. (MOL—maximum overall length, LCL—light center length.) (Adapted from Wikimedia Commons.)

a coiled wire (designated as C) or a coiled coil (the coiled wire is further coiled onto itself, designated as CC). See Fig. 16 in Chap. 15, "Artificial Sources," in this volume. Coiling increases the surface area per unit length of the filament as well as volume packing density. This allows higher operation temperature at higher efficiency due to relative reduction in heat lost to convection. Multiple supports are used to reduce filament vibration. The number and type of supports depend upon the bulb operating characteristics. A higher number of supports is needed for rough/vibration service lamps. These lamps have low wattages and efficiency and operate in environments that involve shock and vibration for the lamp. Heavier filaments withstand vibrations much better and need fewer filament supports.

Fill gases at low pressure are used to prolong the lamp life at high operating temperatures (high efficacy and power) by reducing the evaporation rate of tungsten. For most applications, mixtures of argon and nitrogen are used as fill gases. More expensive gases such as krypton and xenon are added to the mixture when the higher efficacy justifies the cost increase. Bromine is added to create a new class of lamps called tungsten-halogen lamps which we discuss next.

Tungsten-halogen lamps operate at substantially higher temperatures leading to higher CCT and efficacy. In addition, they have longer lifetimes. At high temperatures, evaporated tungsten combines with the halogen gas and forms a gaseous compound. This gaseous compound circulates throughout the bulb via convection. When this gaseous compound comes in contact with the hottest parts of the lamp such as electrodes the halogen compound breaks down. Tungsten is re-deposited on the electrodes and the halogen is freed up to take part in the halogen regenerative cycle. For the halogen regenerative cycle to work, the bulb envelope must reach a high temperature ($>250^{\circ}\text{C}$). At lower temperatures, the tungsten will deposit on the bulb envelope instead and lead to lamp blackening and filament thinning, a common failure mode of incandescent lamps. Operating a tungsten-halogen lamp at less than the rated voltage inhibits the halogen regenerative cycle and takes away the longevity advantage. Tungsten halogen lamps have relatively compact bulb envelopes to allow for high bulb-envelope temperatures. The bulb envelope may also have an IR reflective coating to enhance the heat density inside the bulb. Compact bulb sizes make them good candidates for use with reflectors (such as a PAR) to provide directional properties for the lamp emission. The bulb envelopes for such lamps are made of heat resistant material to withstand high temperatures. The high filament temperature in halogen lamps generates UV. The UV rays must be blocked by a UV absorbing lamp cover or UV absorbing but heat resistant bulb envelope such as high-silica or aluminosilicate.

Flicker in incandescent lamps is naturally very low due to slow response of filament temperature to voltage fluctuations caused by power supply frequency or noise.

Incandescent lamps fail when bulb blackening or filament notching (thinning of the filament by evaporation) reduces the output substantially or the filament breaks either by vibration or by complete evaporation of some filament portion. With the exception of halogen lamps, operating the incandescent lamps at less than rated voltage dramatically extends the lifetime (by reducing the filament evaporation rate) at the cost of reduced efficacy, CCT, and luminous flux. In situations, where long lifetimes are needed such as when the lamps are located in a hard to replace areas and/or under tough environmental conditions, heavy filament lamps that are operated at less than the rated voltage are used. But using lower operating voltage to extend the lifetime is not always economical once the increased cost of electricity due to reduced efficacy and compensation of the reduced flux by using more bulbs is taken into account.

Fluorescent Lamps These are luminous sources based on light emission by excited states of phosphors, a phenomenon known as *fluorescence*. These phosphors are typically excited by UV emission due to spectral line transitions across gases such as mercury vapor and/or rare gases such as Xe and Ar. Most commonly available fluorescent lamps are mercury-vapor-based although mercury-free fluorescent sources are available. Phosphors are now available that are excitable by visible light. Such phosphors are also being used for LED-based light sources to produce white light. Fluorescent lamps are actively replacing incandescent light sources and in many cases are themselves being replaced by LEDs.

As shown in Fig. 14, a fluorescent lamp consists of a closed tubular fluorescent material coated glass envelope filled with low-pressure mercury vapor, electrodes at each end of the tube and a ballast to provide high-strike voltage across the electrodes and limit the current during operation. A high-strike voltage across the electrodes initiates a gas discharge in the mercury vapor with light

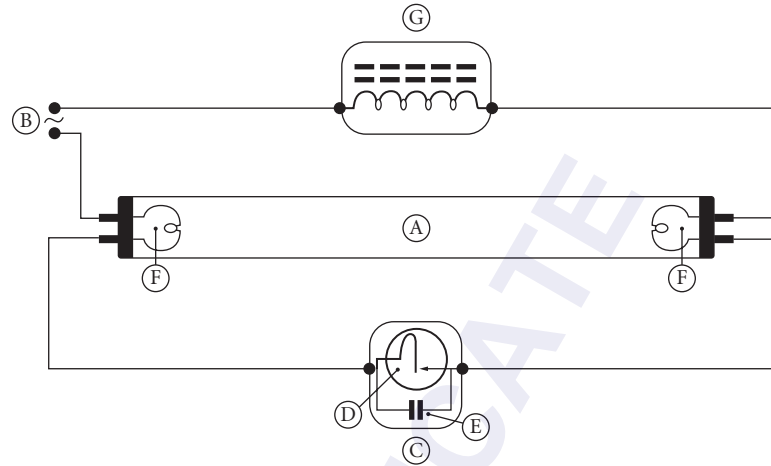


FIGURE 14 A preheat fluorescent lamp circuit using an automatic starting switch. A—fluorescent tube with fill gases, B—power, C—starter, D—switch (bimetallic thermostat), E—capacitor, F—filaments, and G—ballast. (Courtesy of Wikipedia Commons.)

emissions across various wavelengths. The UV emission lines of mercury, mostly at 254 nm, excite the phosphors coated on the inside of the tube, cause them to fluoresce and lead to emission in the visible. The fill gases consist of mercury vapor at low pressure (10^{-5} atm) for UV emission and mixtures of inert gasses such as argon, krypton, neon, or xenon at a relatively higher pressure (10^{-3} atm). Inert gases help in lowering the strike voltage across the electrodes as discussed later.

Mercury-free fluorescent lamps mostly use excimers (excited dimers of rare gases and/or their halides with Xe being popular) to produce UV to excite the phosphors. High-wattage operation and high lifetime is achievable but the efficacy is low as compared to the mercury-based lamps. Xenon-filled fluorescent lamps are also available but far less efficient than excimer based. Both excimer and Xe-based fluorescent lamps have additional advantages: instant-on, instant restrike, and color stability.

Fluorescent lamp envelope material is made of soft soda lime glass. This glass material blocks all UV. The length and diameter of the fluorescent tubes affect the efficacy and operating characteristics (voltage, current, and temperature). Longer tubes need higher voltage and power but provide higher efficacy. The efficacy is optimized for a certain tube diameter. To reduce the angular extent of emission, fluorescent tubes are coated across a prescribed region for high reflection in the visible. In such cases, only the uncoated regions of the tube emit light.

The fluorescent coating consists of different mixtures of phosphor salts.^{46,47} Phosphor salts consist of oxides, sulfides, selenides, halides, or silicates of zinc, cadmium, manganese, aluminum, silicon, or rare earth materials. Inclusions to the base phosphor material help tailor the emission characteristics. There is a large variety of phosphors available. Each phosphor emits light in one or more narrow wavelength bands. A fluorescent coating consists of one or more phosphor materials to produce a desired CRI and CCT for white light emission or specific spectral characteristics. Halophosphates (wide band) and triphosphors (blends of three narrow band red, green, and blue rare-earth phosphors) are commonly used for general lighting applications. Figure 15 shows an example of the spectrum from a halophosphate phosphor. Figure 17 shows the spectrum of a standard fluorescent lamp.

Due to the variety of fluorescent lamp spectra available, it is possible to achieve the desired CCT, CRI, and color requirements of an application. For example, an application may require a specific blend of white light such as “warm white” (CCT 3000 K, CRI 53), “cool white” (CCT 4100 K, CRI 62), “daylight” (CCT 6500 K, CRI 79), or special requirements such as aquarium lighting for providing optimal plant and coral growth and color enhancing of the display. CFLs use triphosphor coating to produce a high CRI (>90) in order to effectively replace the incandescent lamps. Special application

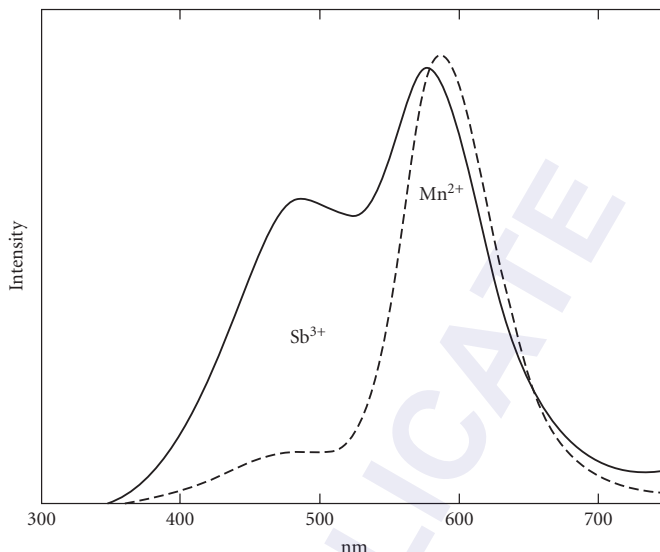


FIGURE 15 Emission spectrum of halophosphate phosphor:⁴⁷ Sb^{3+} , Mn^{2+} activated $\text{Ca}^5(\text{PO}_4)^3(\text{Cl}, \text{F})$. The ratio of Sb^{3+} and Mn^{2+} species can be adjusted to adjust the spectral distribution.

fluorescent lamps include backlight and tanning lamps. These lamps have phosphors that convert short-wave UV into long-wave UV for applications in tanning (UVA and UVB), detecting materials that fluoresce at long UV (urine, paints, or dyes), or attracting insects.

The lamp electrodes are coated with materials such as oxides of barium, calcium, and strontium to provide low-temperature thermionic emission (electron emission from a heated electrode). Under a potential, these electrons accelerate and ionize the inert gas atoms by impact ionization. Each ionization event generates more electrons that are available to accelerate and further ionize leading to an avalanche that rapidly lowers the gas conductivity. Eventually mercury atoms are ionized. The operating voltage drops and a steady current is established. UV is emitted by transitions across the excited states of mercury atoms. The electrodes can be operated in either of the two modes: cold cathode (arc mode) or hot cathode (glow mode). In the hot cathode mode, the electrodes are preheated to improve the thermionic emission and lower the strike voltage. To enable preheating, electrodes at each end are in an incandescent bulb-like configuration with a tungsten filament (straight wire or coiled). The ballast circuitry allows for preheating (up to 1100°C) before the voltage strike. Hot cathode operation allows operation at relatively higher power and over larger tube sizes. In contrast, cold cathode can be a single cylindrical pin electrode at each end. The strike voltage across the electrodes is relatively higher (~10×). The high voltage strips the electrons from the cathode at ambient temperature and initiates the breakdown process of the gas. The coatings on the electrode surface amplify production of secondary electrons, which are produced when high-energy ions and electrons collide against the cathode. These electrons further increase the gas conductivity. Cold cathode fluorescent lamps (CCFL) are compact (~3-mm diameter) and are used in applications such as thin monitors (i.e., LCDs), backlights, timers, photocells, dimmers, closets, and bathrooms. CCFLs are less efficient (<50 lm/W) but provide instant-on and have long lifetimes (50,000 hours). CCFLs operate at high surface temperatures and require complex power supplies which are fairly compact.

The primary functions of the lamp ballast are to provide a high-strike voltage to start the lamp, give regulated current supply during lamp operation and sometimes provide for cathode preheating for rapid restart applications. Often starter circuitry is deployed to preheat the electrodes. It is critical to use the correct lamp-ballast combination for proper operation. Lamp ballast can be a current limiting resistor, magnetic ballast or electronic high-frequency ballast (most modern ballasts). Due to high-frequency operation of electronic ballasts, flicker in fluorescent lamps is reduced to

almost unnoticeable. Flicker is caused by fast response of the fluorescent lamps to the voltage fluctuations in the lamp power supply caused by noise or operating power supply frequency.

Lamp failure mode consists of thermionic emission electrode coating degradation (a function of strike voltage and the number of strikes), phosphor degradation, mercury loss (diffusion or absorption by lamp materials) and ballast malfunction. Fluorescent lamps fail to operate far outside the ambient temperature range for which they are designed. Most fluorescent lamps are designed to operate in an ambient temperature of $\sim 20^{\circ}\text{C}$. For operation at low temperatures, special cold start circuitry and mercury amalgams are needed.

High-Intensity Discharge (HID) and Low-Pressure Sodium (LPS) Lamps These sources emit light across the spectral line transition of enclosed gases by electrical discharge. The source spectrum also includes background thermal radiation due to the heated electrodes and plasma. These sources are similar to fluorescent sources in basic physics involving discharge and emission of light in the UV and visible due to transitions between the excited states of gas atoms. Unlike fluorescent lamps, the electrodes are separated by less than 1 mm up to a few inches. The enclosed gases are at a pressure that is three orders of magnitude higher than in fluorescent sources. As a result, HID sources are far brighter than fluorescent sources with much higher lumen output. The following types of HID lamps are commonly available: mercury vapor (Hg), metal halides (MH), high-pressure sodium (HPS or nicknamed as White SON), and ceramic metal halides (CMH). CMH combines the advantages of MH and HPS technologies. Each HID lamp technology has very different performance and operating characteristics. An LPS lamp is similar to HID lamps in construction and operation with some differences that are identified as we describe the HID lamps below. Figure 16 describes the construction of various HID and LPS lamps. Figure 17 shows the relative spectrum of various HID lamps and comparison with the fluorescent lamp spectrum.

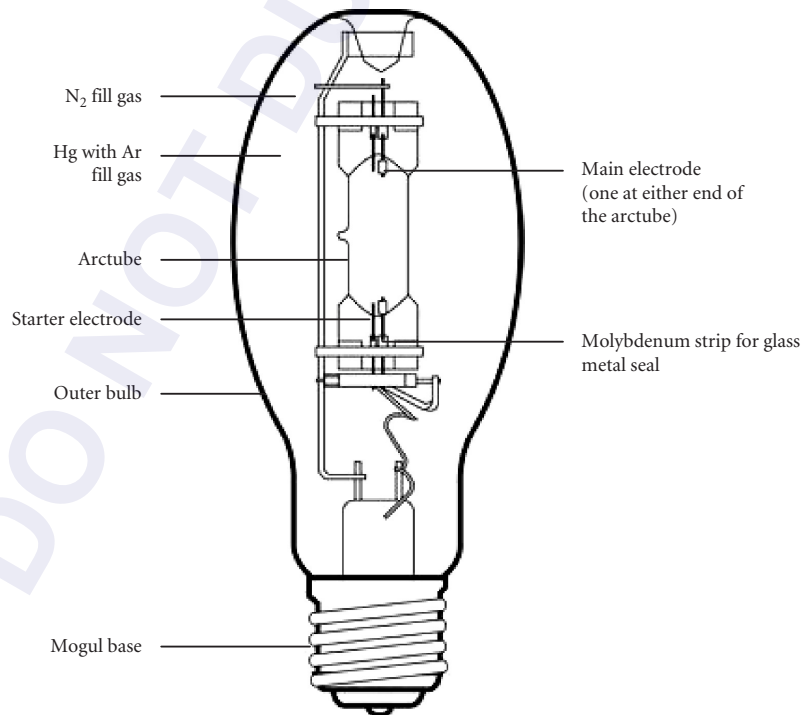


FIGURE 16a Mercury lamp construction. (Illustration courtesy of General Electric Company.)

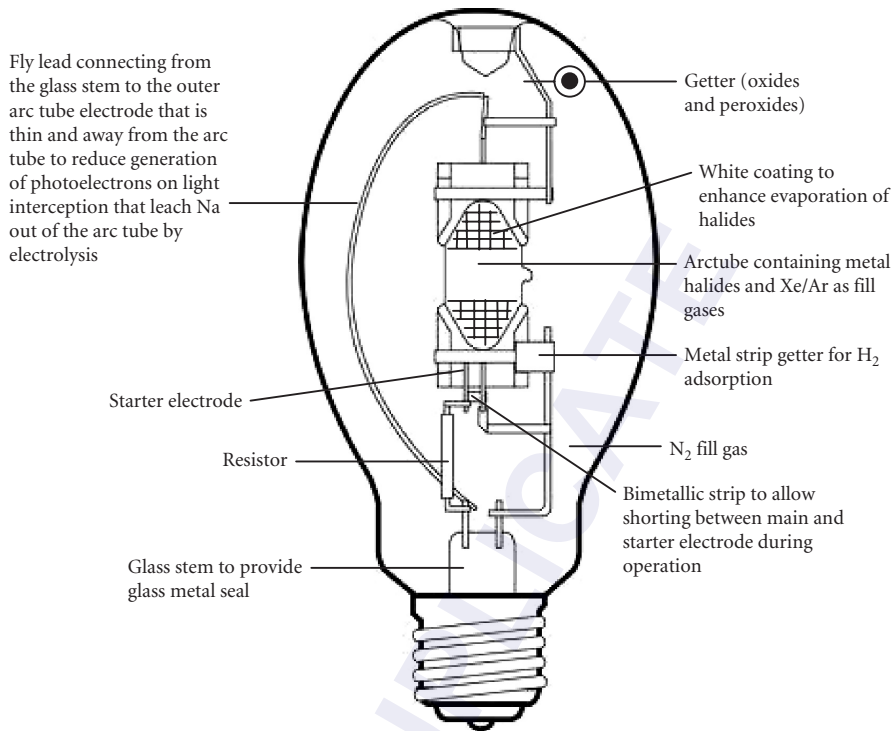


FIGURE 16b Mercury halide lamp construction. (Illustration courtesy of General Electric Company.)

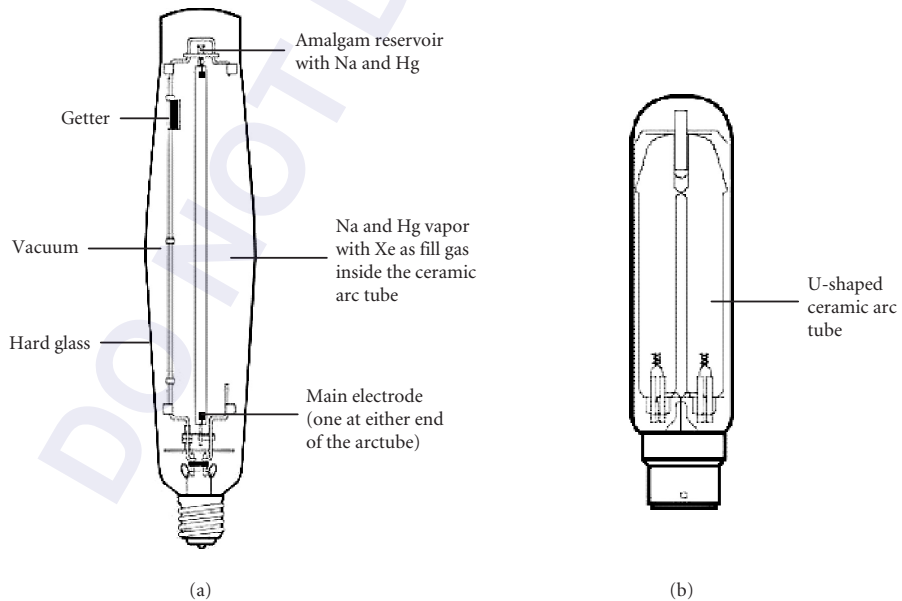


FIGURE 16c (a) High-pressure sodium lamp construction and (b) low-pressure sodium lamp construction. (Illustration courtesy of General Electric Company.)

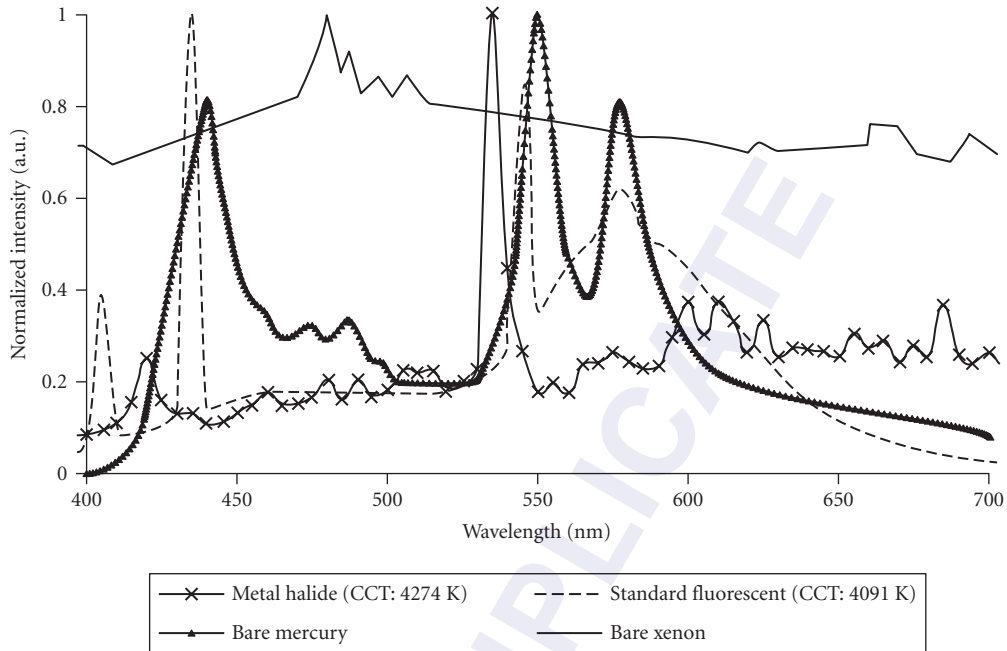


FIGURE 17 Representative lamp spectra of standard fluorescent, bare mercury, bare xenon, and metal halide lamps.

An HID lamp consists of two glass envelopes: inner and outer. The inner glass envelope or the arc tube is made of quartz for MH and Hg and alumina ceramic for HPS, LPS, and CMH lamps. The arc tube houses the discharge gases at high pressure (several atmospheres) and the tungsten electrodes. The outer glass envelope is made of borosilicate and sometimes with soft glass in Hg lamps. It absorbs UV and insulates the arc tube from outer convection currents and from large ambient temperature ranges. It houses the lead-in wires, circuitry to help initiate the high-voltage discharge, getters in case of MH to absorb impurities and has a vacuum (HPS) or low-pressure nitrogen (MH and Hg) to prevent shorting of the lead-in wires. The outer-envelope Hg lamp is sometimes coated with phosphors to provide white light at high CRI and CCT like in regular fluorescent lamps.

Tungsten electrodes are coated with various oxides in a tungsten matrix (except in MH lamps where gases can react with such electrodes) to slow down evaporation and assist in thermionic emission when heated. Starter electrodes, when used in MH and Hg lamps, assist in arc initiation via an electric field between the starter electrode and the adjacent main electrode. During operation the starter electrode is removed from the active circuitry with a bimetallic strip, otherwise premature lamp failure results.

An LPS lamp is similar in construction and operation to HID lamps. Key differences are lower arc tube pressure (0.7 atm), long arc length with a U-shaped arc tube. The arc tube gases include sodium vapor and small amounts of Ne, Ar, or Xe as startup gases.

Different HID lamp types have different gas mixtures. An easily ionizable gas such as Argon (Ar), Neon (Ne), or Xenon (Xe) is used in the arc tube to help in arc initiation. Similarly, mercury is used in most HID lamps to achieve high pressure and improved color rendering. An HPS lamp used sodium (Na)-Hg amalgam. Mercury-free HPS lamps are also available. An LPS lamp uses sodium vapor. The MH lamps use halides of metals in addition to mercury, argon, and xenon.

A low-pressure sodium lamp (LPS) has a CRI of zero due to spectral emission only at 589.0 nm and 589.6 nm (sodium-D line). An HPS lamp has broader spectrum and thus better CRI due to

pressure broadening of the sodium-D line. At very high pressures (27 atm), the sodium-D line is self absorbed by the cooler outer layers of the arc leading to a narrow spectral hole around 589 nm. Mercury atoms help in further broadening of the red end of the pressure broadened sodium-D line due to Van Der Waals forces. A CRI from 22 to 85 is achievable depending upon the pressure which can be greater than 90 atm. MH lamps achieve a rich spectrum due to the line spectra of metals like sodium, tin, dysprosium, holmium, thulium, scandium, iron, or cesium. An MH lamp may have mixtures of halides of one or more metals to achieve the desired efficacy, CRI, and CCT.

An MH lamp using a single metal halide compound can also be used to generate discrete spectral output: orange (sodium), green (thallium), blue (indium), and UV (iron). The metal-halide compound is stable at low temperatures and does not react with the arc tube material unlike some metals. At high temperatures, near the arc, the metal-halide compound breaks down and provides the spectral line emission from the metals. The operating temperature is lower than would be the case where the metals are evaporated to see the spectral emission lines. CRI is usually traded for lifetime and efficacy. CMH lamps combine the advantages of MH and HPS by using a poly crystalline alumina as the bulb envelope material. This material does not allow diffusion of metals, especially sodium or reaction of metals with the bulb material. The bulb is operated at a much higher temperature and pressure than the MH lamps. These advantages lead to high color stability with high CRI, uniformity, and efficacy over the lifetime in spite of the bulb material allowing only ~90 percent transmission.

HID lamps require several minutes of start-up time (time to reach stable output) and restrike. A long start-up is due to the time taken to reach a stable operating temperature and pressure within the arc tube. The restrike time interval results from the need to have low pressure inside the arc tube for arc initiation. Complex ballasts are needed to provide startup, restrike, and stable operation with constant current. To improve startup, restrike, and operations at low voltage, a high voltage-low current pulsed start is used. Sometimes, multiple arc tubes within the outer bulb envelope are used to provide faster restrike. Only one arc tube operates at a time in such lamps. Xe-based HID lamps are capable of instant-on and restrike. Automotive HID lamps use Xe with metal halides to improve the start and restrike times dramatically.

The physical orientation of HID lamps such as Hg and MH during operation is far more important than with the other lamps. Due to convection, within and outside the lamp, different portions of any lamp, not just HID, are heated to different temperatures. The lamp engineering must take this into account and ensure that the hottest regions do not constitute a failure mode either by design or by providing instructions to the user for best operating configuration. In MH and Hg HID, the high convection roll within the long arc tube has an overwhelming effect on the arc shape and position under gravity. It can make the arc shape curved and lead to nonuniform degradation of the electrode tips and impact the light output, lifetime, and light distribution patterns. HPS lamps, however, can be operated in any position primarily due to a compact arc tube at high gas pressure (5 to 27 atmospheres).

Lamp degradation and failure occur due to electrode degradation by evaporation, arc tube blackening due to electrode material deposition, loss of gas pressure, and selective diffusion of gases leading to change in the lamp color. Arc tube blackening also leads to the rise in the arc tube temperature leading to an analogous rise in pressure and operating voltage. The effect is especially pronounced in HPS where lamp cycling can occur: as the lamp cools down, it is able to restrike but after some time temperature rises to the point that it shuts down.

Electrodeless Lamps As the name suggests, electrodeless lamps do not have any electrodes internal to the bulb envelope. As a result lifetime is not limited by electrode degradation. The concept behind these sources is over a century old.⁴⁸ These sources are being sought to replace conventional light sources where high flux is needed at low operational costs (long lifetimes and high efficiency). There are two kinds of electrodeless lamps: induction lamps (IL) also known as *electrodeless fluorescent lamps* and microwave powered lamps also known as *electrodeless sulfur lamps* (ESL). Extraordinary high bulb life times (>25,000 hours) are possible due to lack of electrodes that degrade under operation. The causes of lamp failure are due to electrical components rather than the bulb itself, which implies a greater lifetime.

In each case, the goal is to excite a discharge with an EM field without the need of electrodes inside the bulb. Alternating magnetic fields in IL or microwaves in ESL initiate the discharge by

accelerating the free electrons of a gas with low ionization potential, such as argon or krypton. Free electrons are created in the gas by a spark from a high-voltage pulse across two electrodes in the vicinity of the bulb. These free electrons ionize the gas atoms by impact ionization. Ionization yields more free electrons and ions and the process resumes, eventually resulting in plasma formation. Excited states of gas atoms produce light via spectral transitions across various wavelengths. In case of IL, mercury is present in addition to argon/krypton to produce UV from excited mercury atoms. The UV excites the phosphor coating on the inner surface of the bulb envelope and emits white light just like a regular fluorescent lamp.

In ESL, microwaves are used to produce an intense plasma inside a rotating quartz ball containing argon/krypton and sulfur. The rotation of the quartz ball helps in stabilizing the fill for uniform emission as well as convective cooling with a fan to prevent its meltdown. Initially, the microwaves create a high-pressure (several atmospheres) noble gas plasma. This heats sulfur to a high temperature resulting in brightly emitting plasma. Light emission by sulfur plasma is due to the molecular emission spectra of the sulfur dimer molecules (S_2). The spectrum is continuous across the visible and has >70 percent of its emission in the visible. It peaks at 510 nm, giving a greenish hue. The resultant CRI is 79 at a CCT 6000 K. The lamp spectrum can be modified with additives such as calcium bromide, lithium iodine, or sodium iodide or by using an external color filter.

Electroluminescent Sources Electroluminescent sources are materials that emit light in response to an electric field. Examples of such materials include powdered ZnS doped with copper or silver, thin film ZnS doped with manganese, natural blue diamond (pure diamond with boron as a dopant), III-V semiconductors or inorganic LED materials such as AlGaAs, phosphors coated on a capacitor plane and powered by pulsating current, organic LED (OLED) also known as light-emitting polymer or organic electroluminescent. The sources can operate at low electrical power with simple circuitry.

Electroluminescent sources are commonly used for providing illumination across small regions such as indicator panels. LEDs are already a major lighting source that is rapidly replacing incandescent and fluorescent light sources. The remainder of this section discusses LEDs and OLEDs.

LEDs emit light by electron-hole pair recombination across the P-N junction of a diode. The wavelength of the emitted light corresponds to the band gap (energy gap between valence and conduction bands) across which the electron hole pair is created. The degeneracy in the valence and conduction bands leads to a closely spaced band of wavelengths that constitute light from the LED. Narrow spectral bandwidth enables applications that require saturated colors. Although LEDs are not available for every desired color, it is possible to combine LEDs of different colors to create any color within the color gamut defined by these LEDs. As such, color mixing has become an important field. LEDs emitting in specific bandwidths can be combined with sources with continuous spectra to either enhance a certain spectral region or to provide an easy dynamic color control. One easy method of combining multiple LEDs is using lightguides. Lightguides with rippled surface texture along the cross section are particularly efficient in combining multiple colors with excellent uniformity in a short path length.^{49,50}

LED lamps may have an array of small LED chips or a single large chip to achieve the desired power levels. The directionality is controlled by appropriately mounted optics. DC operation of LEDs makes them flicker free sources. Figure 18 shows the structure of a simple packaged LED. Chapter 17 in this volume is dedicated to the subject of LEDs.

LED packages come in various sizes and shapes. Surface mount LEDs (SMD) have minimum packaging and are almost a bare die. LED packages are also offered in multicolored die formats.

Over the years, a variety of materials for LEDs have been used with the goal of obtaining higher efficiencies and different colors across the visible spectrum. LED technology is fast evolving with ever increasing brightness, lifetime, colors, and materials and decreasing costs. The available materials, at the time this chapter was written, are listed in Table 5.

Most lighting applications require white light. Some of the processes for making white-light LEDs are listed below:

- Arrays of small red, green, and blue dies placed in close proximity in a single LED package. Good color mixing takes place in angular space.
- Color-mixing of red, green, and blue colors using lightguide or other optical means.^{50,52}

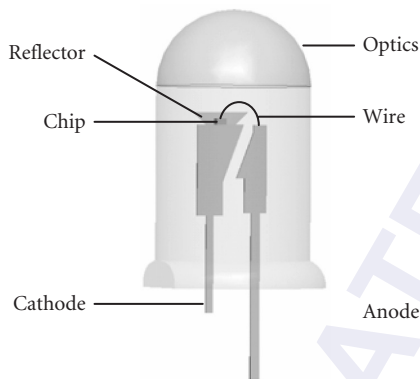


FIGURE 18 Structure of a simple packaged LED.

TABLE 5 LED Materials and Emitted Colors

Material	Color
AlGaAs	red, infra
AlGaP	green
AlGaInP	higher brightness orange, orange-red, yellow, and green
GaAsP	orange, orange-red, orange, and yellow
GaN	green and blue
InGaN	blue (450–470 nm), near UV, bluish-green, and blue
SiC (as substrate)	blue
Si (as substrate)	blue
Sapphire (as substrate)	blue
ZnSe	blue
Diamond	UV
AlN, AlGaN, AlGaInN	UV(<210 nm) ⁵¹
Organic light-emitting diodes (OLED)	Red, green, and blue

- Phosphor excitation by blue or UVLEDs.
- Novel techniques like quantum dot blue LEDs or homoepitaxially grown ZnSe blue LEDs on a ZnSe substrate. The active region emits blue light while the substrate emits yellow light.^{53,54}

There is a continuous push to improve the LED efficiency and brightness while keeping the lifetimes high. High brightness LEDs became possible due to large area chips, efficient heat extraction and better light extraction from the chip. Internal quantum efficiency of LEDs can be increased by placing emitters inside a cavity⁵⁵ to increase the radiative recombination rate. Due to the high internal Fresnel reflections and lateral waveguiding, a lot of light fails to exit the chip. Techniques such as texturing the surface with photonic crystals assist in increasing the light extraction from large dies.^{56–58} Figure 19 shows the internal structure of a photonic crystal LED.

Organic LEDs (OLEDs)^{59,60} in contrast to inorganic LEDs are size-scalable light sources with richer color spectra. OLEDs can be used to create flexible transparent lighting solutions as they can be printed on a malleable substrate with transparent electrodes. Currently OLEDs are being used for displays and are competing with LCD flat panels. Conceptually, OLEDs are no different from inorganic

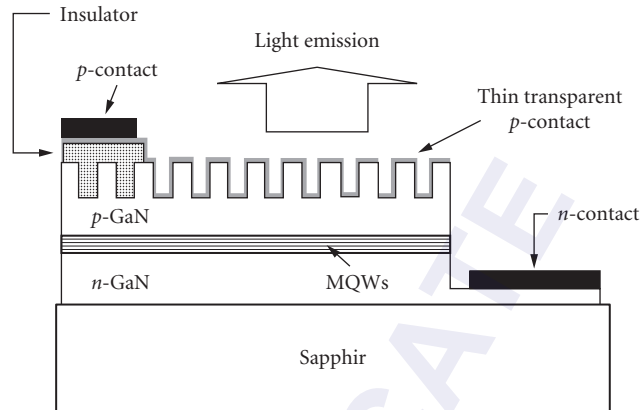


FIGURE 19 Internal Structure of a Photonic Crystal LED. MQW refers to “multiple quantum wells.” (Courtesy: Seoul National University, Korea, http://optics.org/cws/article/research/23635/1/sem_image.)

semiconductor based LEDs. An OLED deploys layers of organic materials on polymer substrates to form conductive and emissive layers connected to a cathode and anode respectively. Much of the OLED research is aimed at making them brighter and longer lasting.

LED failure causes include damage due to degradation of the active layers with time (spontaneously or in operation); plastic package degradation due to ambient UV; electrostatic discharge; current crowding or inhomogeneous current distribution across the junction leading to hot spots; and thermal stresses causing rupture of the LED package, diffusing of the metal contact material into the die material at high currents, high output leading to facet melting and phosphor degradation in white LEDs, and degradation of organic layers in OLEDs.

Miscellaneous Artificial Light Sources *Neon signs* are essentially cold cathode-like operation of a fluorescent tube without phosphors. A low-pressure mixture of noble gases such as neon, argon, helium, xenon and a small amount of mercury is used in the discharge tube. Neon emits a reddish-orange color; argon emits blue; and krypton, helium, and xenon emit over a wider spectrum. Colored filter glass can be used for making different colors.

Short arc sources function almost identical to HID lamps with the only difference being that the electrodes are much closer (less than 1 mm to 12 mm). The gases are mercury (with argon), mercury-xenon, pure xenon, or metal halides (with mercury and argon). These lamps are primarily used in illuminating high loss systems where the source étendue needs to be as small as possible. Projectors, medical optical instruments, metrology instruments, and daylight or solar simulators use such lamps. These lamps have lifetimes from a few hundred hours up to 10,000 hours. The light sources are typically used with a reflector (parabolic or ellipsoidal). Sometimes the reflector is an integral part of the source/lamp package.

Pure Xe arc lamps have an instant-on capability and provide high CRI (>80) at high CCT (>6000 K). Digital cinema projectors and flash tubes (for warning signs, entertainment applications, camera flash lights, and warning or emergency signs and indicators) often deploy pure Xe-arc sources.

Lasers are used for visual displays for entertainment. The subject of Lasers is discussed in Chap. 16 in this volume.

Nuclear sources are self luminous light sources that function by phosphor excitation caused by beta radiation from radioactive materials such as tritium. These light sources are used to illuminate tiny spaces such as watches or displays of instrument panels in very low ambient light.

Glow lamps are low wattage arc sources with gases such as argon emitting in the UV to excite UV-excitable materials or neon to emit orange light to be used as indicator lights.

Carbon arc sources are now obsolete but still find applications in illuminating small areas with bright light under demanding environmental conditions (such as outer space). An arc is struck across a pair of carbon rods and the incandescence from the heated carbon rods provides the light.

Gas lights that operate by the burning of gases like methane, ethylene or hydrogen are used with appropriate lanterns primarily for decorative applications.

Natural Sources: Daylight Daylight can be utilized in the lighting design of buildings to provide a pleasing environment that enhances physiological well-being and productivity and also energy savings during the day by reducing the need for artificial lighting and solar influx contribution to building heating. Daylight is primarily used for ambient lighting. It can be used for task lighting when integrated with electrical lighting. Daylight constitutes direct sunlight, scattered sunlight from the atmosphere, reflected sunlight from the clouds, and reflected light from the surroundings such as ground (especially snow) and objects such as buildings. The solar spectrum changes with atmospheric conditions and so does scattered light from the sky or reflected light from the ground. The CCT of Sun is 1000 to 5500 K, clear blue sky is 10000 to 100000 K, overcast sky is 4500 to 7000 K and clear sky with sunlight is 5000 to 7000 K. Figure 20 shows an example of the solar spectrum at the ground, at noon, at Golden, Colorado. Note the IR content in the spectrum and the shift in spectrum from noon to evening.

Designing for daylight requires close attention to many factors:

- Goals of providing daylight: physiological well being of occupants and/or energy savings.
- Intended distribution of light inside the building during the day and in the night. Penetration of daylight into the interiors and the impact of reflectivity from various surfaces.
- Impact of daylight on materials such as wall paints, artwork, plastic materials, furniture, and plants. The UV component of daylight is generally harmful to most materials via solarization of plastics or fading of paints and stains. If the impact of UV is not known and acceptable, then the UV should be rejected by the optical system components through coatings or materials.
- Integration of daylighting-based building design with other controls such as for electrical lighting, cooling and any automated systems.

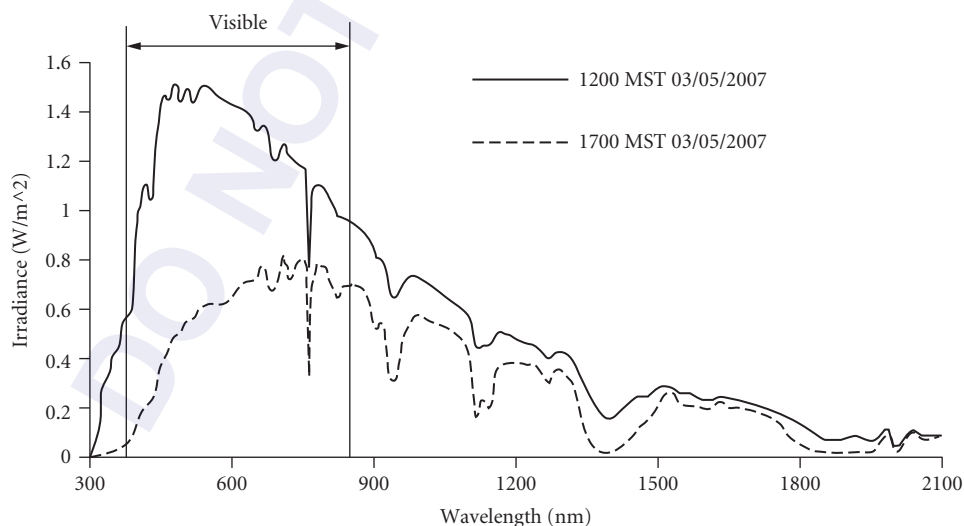


FIGURE 20 Solar spectrum recorded at ground at Golden, Colorado at two different times during the day. (Courtesy of National Renewable Energy Laboratory, USA.)

- Location (latitude and longitude) and orientation of the building relative to the surroundings. Simulating the daylight entering the building during various times of the day and the year.
- Design for reduction or elimination of glare from daylight sources. Internal layout of the building during intended use will play a key role.
- Outside view requirements via windows. A minimum window size^{61–63} is needed for a given geographic location, and a certain window location is dictated by aesthetic reasons including the desired view from the window. The minimum window size can also depend upon the prevalent building regulations based on wall size or floor size. Under such circumstances, the daylight entering must be controlled or balanced with interior lighting and cooling for uniformity, glare, and heat management. Daylight can be controlled by providing appropriate window shades, glazings, or even additional windows or skylights at appropriate places.

It is not always possible to obtain energy savings with daylight due to increased heat load, especially during summers and the need to supplement daylight with electrical lighting during the night. Infrared rejection is achieved by using special glazings. To limit glare, glazing such as high reflectance, low transmittance, electrically controlled transmission, tilted glazings with angularly dependent transmittance or IR reflecting films are used. Adequate consideration must be given to the window view impact due to any of the daylight control mechanisms. A very low-transmittance glazing can make the outside view appear gloomy on a bright day.

Luminaire Design

The source is the starting point of the luminaire design. The optics performs at a minimum two functions: first to capture the light from source and second to transfer light efficiently to the desired distribution at the target (i.e., illuminance, intensity, and/or luminance). The choice of the optics is also a challenging process. First, the designer must select if refractive, reflective, or both types of optics (i.e., hybrid) are to be used in the design. Reflective optics have been a standard for most sources, with flat refractive optics used at the output aperture in order to protect the other optics. Until recently reflectors were standard conics such as ellipses, hyperbolas, parabolas, or spheres. Standard refractors in use include cover plates (still called a lens), pillow lens arrays, Fresnel lenses, and other faceted designs. However, with advances in LEDs, refractive optics are increasingly being used for advanced function. Solid-state lighting optics are either plastic or glass that surround individual or a limited number of LEDs. They are hybrid optics that use at least total internal reflection (TIR) and refraction at the input and output facets, but they can also use reflection, scattering, and even diffraction. With the advent of better technology, especially manufacturing capabilities and software for modeling, faceted and continuously smooth surfaces parameterized with nonuniform rational B-splines (NURBS) are being developed. For high-performance, injection molding of plastic or glass is being increasingly used. For refractive optics, the surfaces are left bare, but for reflective optics there is vacuum metallization of the surfaces. Reflectors are also made with stamping methods, but these optics tend to have lower performance.

Baffles (louvers) are often considered by many to be optics of the system; however, many, rather than shaping a distribution through redirection of the light, block or even absorb incident radiation. Thus, while some baffles are reflective and provide some shaping of the illumination distribution, they most often achieve shaping through subtraction rather than addition. They are primarily added to alleviate glare, trespass, and pollution, but they are also included for aesthetic reasons, correcting errors in the design process, and to hide structure within the luminaire.

The source coupling aspect is the focus of the next section, following that is a discussion of the design of the optics and subsections on baffling in the form of luminaire cutoff classification and types of luminaires. Following this process from the source to the optics to the baffle to the target, while always considering the perception aspects of the illumination (see Sec. 40.3), means that aesthetically pleasing while technically sound lighting can be developed.

Étendue and Source Coupling For reflective optics, the source-coupling components also act as the transfer optics, often in conjunction with a front, protective lens. With the advances in source

technology that use hybrid, dielectric components, coupling of the sources, is increasingly important in order to improve upon system efficiency. Typically, individual LEDs are placed in recesses in the dielectric optic, and by obeying the conditions for TIR, all of the emitted light can be captured by the coupler and then transferred to following optics that shape the emitted distribution.

In all cases, the term étendue describes the flux transfer characteristics of the optics, starting with the coupling optics, of an optical system, such as a luminaire. Étendue is a geometrical quantity that is the integrated product of an area and solid angle. In paraxial form it is the Lagrange Invariant, but in nonparaxial form it is given by

$$\mathcal{E} = n^2 \iint_{\text{pupil}} \cos \theta dA d\Omega \quad (12)$$

where \mathcal{E} is the étendue, dA is the differential area, $d\Omega$ is the differential solid angle, θ is the angle with respect to the surface of interest (i.e., normal), and n is the index of refraction of the space. The limits of integration in area are over some aperture (e.g., a lens clear aperture or a reflector exit aperture), while the solid angle integration is over the limits that are passed by the aperture. For example, consider a source of area A_s that emits into a half angle of θ_0 from every point on the surface. The étendue for this source is

$$\mathcal{E} = n^2 A_s \int_0^{2\pi} \int_0^{\theta_0} \cos \theta \sin \theta d\theta d\phi = \pi n^2 A_s \sin^2 \theta_0 \quad (13)$$

In lossless optical systems, étendue is conserved. Thus, in order to design the most efficient luminaire, one must continue to match the étendue as one progresses through the optical components of the system. For example, if the source of Eq. (13) is used for a luminaire, one must keep this étendue quantity consistent. If one desires to reduce the angular spread of the output from an optic ($\theta_0 > \theta_{\text{optic}}$), then the area of the optic must be increased ($A_s < A_{\text{optic}}$). The counter also holds true: to reduce the real extent ($A_s > A_{\text{optic}}$), one must increase the angle ($\theta_0 < \theta_{\text{optic}}$). An expression for conservation of étendue in a generalized form is

$$dx dy dp dq = dx' dy' dp' dq' \quad (14)$$

where the dx and dy terms are the differential position terms and the dp and dq terms are the differential optical direction cosine terms, which are equivalent to ndL and ndM respectively. More information about étendue can be found in Chap. 39, “Nonimaging Optics: Concentration and Illumination,” in this volume. Another factor related to étendue through a differential is skewness, which denotes the twist on individual rays of light in an optical system. Skewness is also invariant and implies that transfer from one source geometry (e.g., a square) cannot be transformed to another source geometry (e.g., a circle) without loss except if some rotational asymmetry is added to the optical system. Further information about étendue and associated terms like skewness can be found in the literature.⁶⁴

Luminaire Design Methods There are a multitude of design principles for the design of the optics of a luminaire. Fundamentally, most design methods are based on the basic conic shapes as listed in Table 6. Each of these shapes provides a basic intensity distribution at its output aperture. However, increasing demands of tailored light distributions and also increased efficiency require perturbations to these basic design forms. Furthermore, the topics of light trespass, light pollution, and glare are receiving a wealth of attention from ordinance and regulatory agencies. To reduce glare issues it is best to use diffuse optics with a well-defined cutoff. In the field of nonimaging optics (see Chap. 39), the edge-ray theorem provides a means to have a well delineated cutoff. The edge ray is defined by the maximum extent of the source, thus providing a maximum cone of light from the reflector designed around the source shape. However, the edge-ray principle is passive with respect to the luminance distribution of the source—it contends for the maximum extents but not the physical distribution of light in the radiation pattern.

Thus, tailoring methods have been developed. The tailoring methods specify the shape of the optics based upon the luminance distribution of the source and the desired illumination pattern

TABLE 6 Basic Conic Shapes, Their Conic Provide at Their Output Aperture Constant, and the Basic Intensity Distribution That They Provide at Their Output Aperture

Shape	Conic Constant (k)	Base Intensity Distribution
Hyperbolic	$k < -1$	Diverging; far-field applications
Parabolic	$k = -1$	Collimating; pseudocollimation applications
Elliptic	$-1 < k < 0$	Converging; near-field applications
Spheric	$k = 0$	Converging; self-imaging applications

(i.e., luminance, illuminance, and/or intensity) at the target. These methods are extensive and beyond the confines of this chapter. The reader is encouraged to consult the *Handbook* chapter on nonimaging optics (Chap. 39) for a brief introduction, the theoretical book on nonimaging optics by Winston, Miñano, and Benítez,⁶⁴ and the applied book on nonimaging optics by Chaves.⁶⁵

Luminaire Cutoff Classification A cutoff ensures that light from the luminaire is restricted above the horizon with respect to the lamp geometry. Cutoff is designed into the luminaire through the optics (i.e., edge-ray designs) and/or the integration of baffles. While most applications do not require cutoff classification, except for those used on the exterior, such as automotive, roadway, and landscape lighting, most designers include such to make effective lighting systems by alleviating potential glare, trespass, and pollution concerns. Often strict cutoff guidelines are mandated by governmental standards, such as for automotive, traffic signal, and roadway lighting. The goals are to provide the required lighting level to its users, while also alleviating light pollution and light trespass. Light pollution is light that is directed up into the atmosphere, causing sky glow, which is especially present in urban settings. The reduction of light pollution is a growing trend being addressed by the astronomy community. When light is incident on surfaces outside the intended illumination region, it is called *light trespass*. The impact of light trespass from roadway lighting is a major concern in residential areas.

For both trespass and pollution the luminaire cutoffs provide a protocol to reduce both. Automotive lighting does not use the criteria presented here, but rather use a set of governmental standards. Roadway and external lighting make the most use cutoff criteria as shown in Table 7. See Fig. 22 for a depiction of the angles listed in Table 7.⁶⁶

Luminaire Classification System In 2007 the IESNA published research results for refinement of the cutoff classification system of the previous section.⁶⁷ This study focused on light distribution in front of the luminaire (forward light), behind the luminaire (back light), and above the luminaire (uplight) as shown in Fig. 21. They found the photometric luminaire efficiency ($\eta_{luminaire}$) to be

$$\eta_{luminaire} = 100 \frac{\Phi_{forward} + \Phi_{back} + \Phi_{uplight}}{\Phi_{source}} \tag{15}$$

where $\Phi_{forward}$, Φ_{back} , $\Phi_{uplight}$, and Φ_{source} are the integrated fluxes in lumens over the solid angles shown in Fig. 21 for forward light, back light, uplight, and the bare source, respectively.

TABLE 7 Amount of Emitted Light Criteria for the Luminaire Cutoff Classification⁶⁶

Type	Horizon and above (90° and greater)	10° below Horizon (80° to 90°)	Remainder (0° to 80°)
Full Cutoff	0%	≤10%	≥90%
Cutoff	<2.5%	≤10%	≥87.5%
Semicutoff	<5%	≤20%	≥75%
Noncutoff	No restrictions over entire angular space		

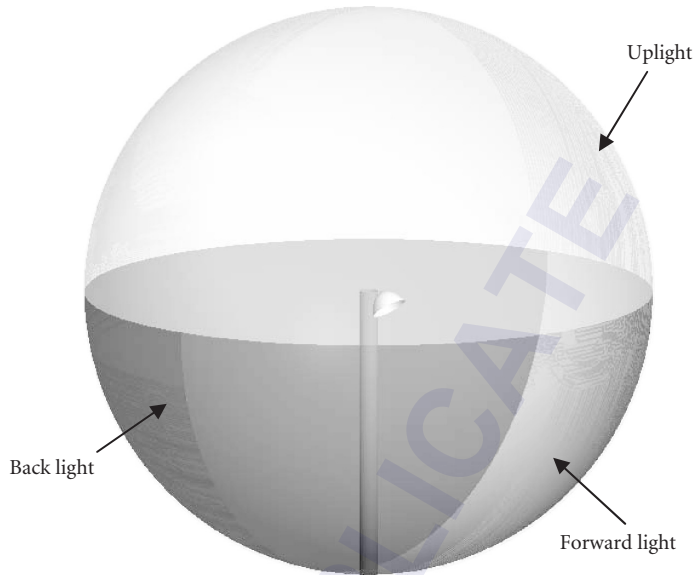


FIGURE 21 Lighting classification system zones for forward light, back light, and uplight based upon the exit aperture of the luminaire.

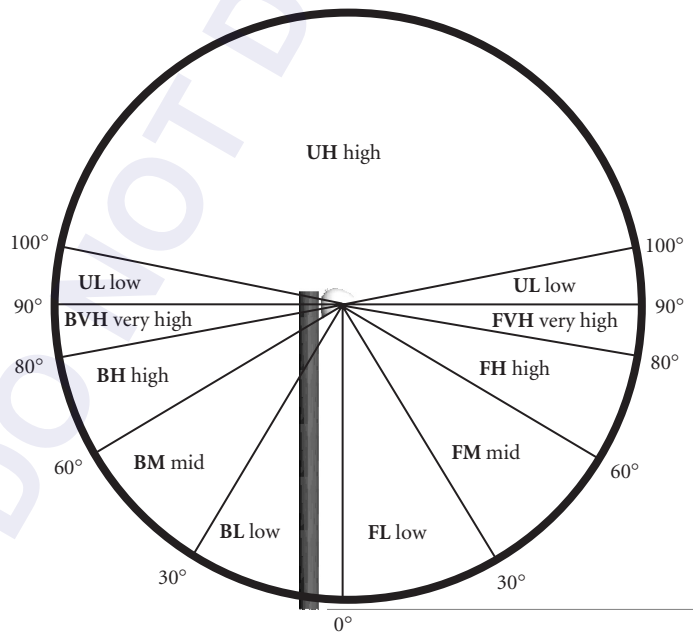


FIGURE 22 Layout of the light classification system subzones. (See also color insert.)

These three zones are then broken down into a total of 10 subzones as shown in Fig. 22. The lumens in each subzone is measured with respect to the bare source flux, as in Eq. (15), in each of these subzones and then reported as an evaluation of the luminaire. These subzones indicate the distribution of light over several regions rather than restricted to the 80° and greater as per the previous section. Standard goals include reduction of light above the horizon (i.e., all upright subzones and the BVH and FVH subzones) and the desired uniformity over the other zones. This new classification system provides for better control of the illumination such that both vertical and horizontal surface illuminances can be addressed in the design process. Horizontal surface illuminance criteria are met by increasing the flux in the BM to FM subzones. Vertical surface illuminance criteria are met by increasing the flux in the BM, BH, FM, and FH subzones.

Luminaire Optics There are innumerable schemes to the design and availability of luminaire optics, both for artificial and natural sources. Methods employing reflectors, refractors, TIR optics, and combinations thereof have been developed. Additionally, the optics are both specular and diffuse or a combination of these two are used. In the next few subsections, we provide examples of commonly available luminaire optics. Finally, as per Table 6, the shapes of the optics are typically based on conic shapes.

Luminaire optics for artificial sources While the design of the optics of the luminaire is to provide a desired illumination distribution, we learned from the previous sections that light cutoff is important in the design of the luminaire. In fact, the ability to hide or shield the source from direct view in order to reduce glare concerns is as equally important as obtaining the desired illumination distribution. A reflector as shown in Fig. 23a has a fairly wide direct view of the source, denoted as the shielding angle or similarly the cutoff angle. Room lamps are perfect examples of this dual requirement since it is typical to use a diffuse shade around the source. The source and shade (which also acts as a diffuse reflector) provide the desired general illumination, while the shade provides the requisite cutoff. Besides using the body of the reflector to hide the direct view of the source, louvers (see Fig. 23b), a prismatic or Fresnel lens (see Fig. 23c), or a bulb shield (see Fig. 23d) are

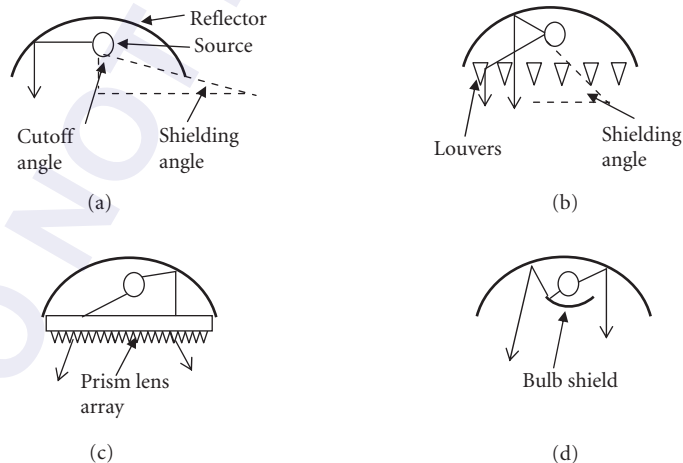


FIGURE 23 Glare issues of four luminaire geometries: (a) the reflector providing the cutoff; (b) louvers, flat or parabolic, increasing the shielding angle; (c) a lens array, prismatic or Fresnel, which directs more of the light emission downward thus reducing the direct intensity level of the source; and (d) the inclusion of a bulb shield to ensure that all emission strikes the primary reflector at least once. (Adapted from Ref. 5, p. 71.)

added to the luminaire. For the baffle case two standard options are typically used: vertical, strongly absorbing louvers, or specular reflecting parabolic louvers. The former inhibits direct view of the source in the luminaire while also absorbing the glare-inducing light. The latter also reduces the direct view angle of the source, but it uses the parabolic, specular louvers to direct the glare-producing light into directed radiation typically outside the direct view of an observer. The lens, prismatic or Fresnel, coupled to a reflector causes most of the emission to be directed downward, thus frustrating the direct imaging of the source. This obscuration of the source means glare effects are reduced. The bulb shield, which is typically spherical, ensures that all light from the source is incident on the primary reflector of the luminaire, thus completely hiding the source from direct view and in turn greatly reducing the glare potential.

The specifics of the lamp design depend on the application. Figure 24 shows some representative luminaire designs. Figure 24a shows a banker's desk lamp, which has multiple bounces within the glass envelope optic. It creates a wide illumination distribution, but the colored glass and the multiple bounces softens the appearance of the bulb located within. The envelope optic can be positioned by the user to alleviate source glare while also providing the desired illumination over a desk surface. Figure 24b shows a Bouillotte table lamp, which uses two vertically oriented fluorescent lamps. The shade hides part of the bulbs from direct view, and the coating on the fluorescent tubes causes near-Lambertian emission; therefore, the illumination for this lamp is wide and glare issues are minimal. Figure 24c shows an indirect RLM lamp fixture that is suspended from a ceiling. The indirect nature

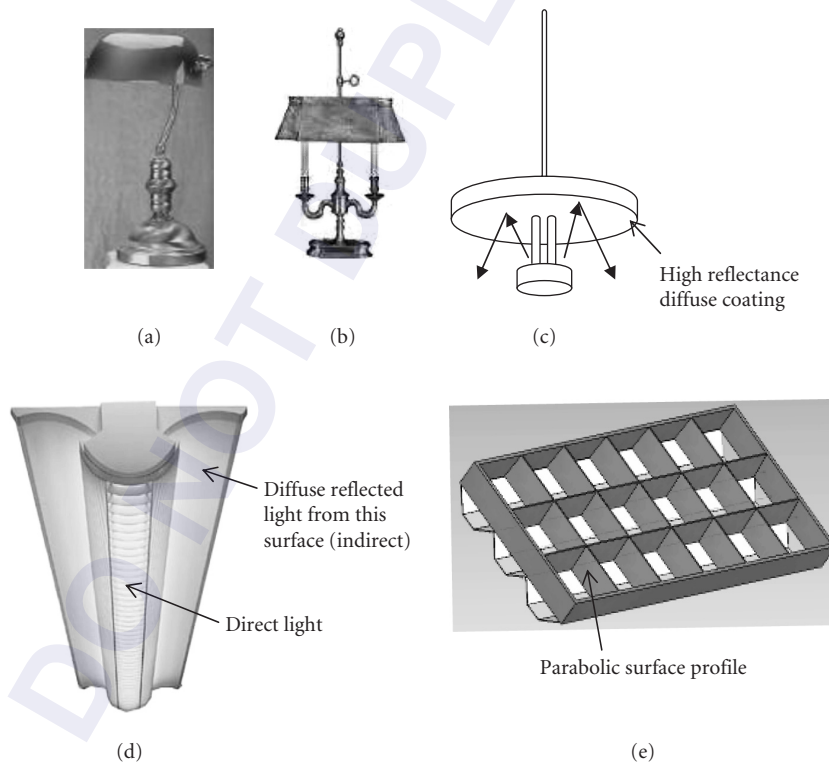


FIGURE 24 Depictions of luminaires: (a) Banker's lamp: multiple bounces inside the reflector create a wide angled uniform illumination; (b) Bouillotte lamp: vertical fluorescent tubes provide diffuse illumination; (c) indirect lighting with RLM fixture where the top surface reflects light into a wide angular range; (d) overhead direct-indirect lighting fixture using fluorescent tubular bulbs; and (e) parabolic louvered trough reflector for fluorescent tubes. (See also color insert.)

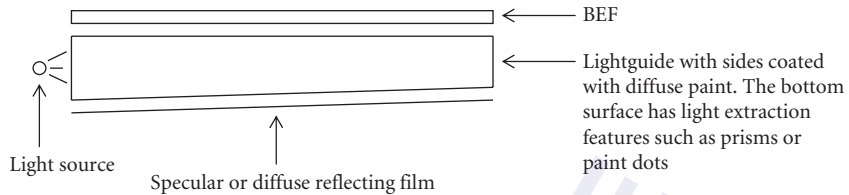


FIGURE 25 An edge lit backlight. (BEF: Brightness enhancing films.)

is due to the inclusion of a bulb shield and a highly reflective diffuse primary reflector thus alleviating glare. Figure 24*d* depicts a direct-indirect overhead, fluorescent lamp fixture. The white, diffuse wing structures provide indirect lighting, while the central section is louvered direct lighting. The direct lighting provides task illumination needs, while the indirect lighting provides a general light level for the room. Finally, Fig. 24*e* shows a standard overhead office luminaire: a series of parabolic troughs (called a troffer) in which long, tubular fluorescent lamps are located. The parabolic louvers reduce glare concerns by increasing the shielding angle, while also increasing the task illumination due to the specular reflectivity of the vanes. Instead of using the louvers, a pillow lens array or prismatic lens can be used; however, computer monitor glare issues can arise with these lenses. Other luminaire geometries not presented here include track lighting, recessed lights, chandeliers, and spot lamps. Please see Refs. 71 and 5 for more information.

Backlighting Backlighting is used extensively in photography to separate the background from the subject and create 3D effects. It can be used for the same purpose in interior lighting to illuminate displays from behind. Backlighting is used extensively in signage, to illuminate instrument panels and device display panels such as laptop and cellphone screens. A typical edge lit backlight operation is illustrated in Fig. 25. Light enters the planar lightguide⁶⁸ (for example, 50 mm × 50 mm × 5 mm) from the thin edges and bounces around. Carefully designed and positioned light extraction features such as prism or spheres deflect the light out of the lightguide. As a result the entire planar surface is lit and appears as a planar light source. Modern backlights use sources such as CCFLs and LEDs. With LEDs, dynamic multicolored backlit displays are possible as lightguides allow for efficient color mixing. Edge lit backlighting can be used in direct lighting of large spaces such as living rooms by creating illuminated ceilings, walls, or artificial windows. Figure 26 shows one such application where backlit ceiling tiles are used to create an illuminated ceiling or an artificial skylight. A picture of sky and vegetation is superimposed on the tiles to create an effect of natural sky with vegetation on the roof. Backlights could be replaced by OLEDs which provide not only illumination but also information content.

Luminaire optics for daylight sources Daylighting schemes^{69,70} involve careful layout of windows, skylights, skytubes, and controls such as shades, window overhangs, window depths, light shelves, and hybrids with electrical lighting. For effective daylight illumination, it is necessary to determine the access to sunlight by taking into account the sun path across the sky during the day and across different months and the impact of neighboring buildings and ground features. Daylight can be used for city lighting too with careful planning. Heliostats or large plane mirrors have been used atop buildings or even mountains to direct sunlight into the city interior.

Layout of windows and skylights can utilize schemes such as side lighting, top-level lighting, and clerestory lighting. Placement of nonview providing windows for providing daylight must be designed carefully for maximum daylight penetration and controlled glare. Side lighting allows daylight in from the walls usually at eye level, top lighting allows daylight in from skylights, and clerestory lighting allows daylight in from the side windows near the roof above the eye level. Clerestory windows provide more uniform ambient illumination over a larger region. However, proper attention must be given to glare from the sky or direct sun by using baffles. The depth of windows near the ceilings or window overhangs also helps in limiting such glare. Figure 27 shows several different



FIGURE 26 A conference room with artificial skylight made up of backlit ceiling image tiles. (See also color insert.) (Courtesy of The Sky Factory, LLC.)

layouts of windows and skylights that bring daylight into the interiors.^{71,72} The location and number of openings for daylight determine the penetration and uniformity of the illumination achieved. The height and the slope of the ceilings determine the penetration and the illumination gradient toward the room. For example, for daylight from windows located near the ground, having a tall ceiling allows light to penetrate to the farthest ends of the room. Similarly a sloping ceiling with windows located near the ground allows a gentler illumination gradient from the window to interior.

The reflective properties of the walls and the ceiling must be taken into consideration when simulating the impact of daylight. High reflectance paints on the ceiling can be used to spread the daylight entering from the window portion near the ceiling into the building interiors. Window blinds, shades, and mechanical louvers are commonly used to control daylight. Sometimes partition walls around certain window sections are used to control glare or even limit the extent of illuminated region such as artwork in museums.

Light shelves are used interior or exterior to the windows to allow daylight from the sky without glare. Light shelves consist of large horizontal sections hanging below the top edge of the window. A layout of mirrors or even just a plane aluminum sheet directs the daylight from the sky toward the ceiling or deep into the room without hitting the ground. A reflective ceiling will scatter daylight into the room. Combinations and variations of these schemes can be used to suit a particular situation. Figure 28 illustrates the concept of light shelves. Suncatchers are similar in concept as shown in Fig. 29.

Skytubes are lightguides that carry daylight into the building interiors. The inside walls are specular with high reflectance. These tubes may be straight or curved, as needed, to transport daylight to different regions. Curved lightguides have a symmetric cross section to maximize light transfer efficacy across the bends.

Hybrid daylighting systems consist of daylight integrated with electrical light as shown in Fig. 30.

Solar lighting systems include the use of tracked or untracked mirrors, lenses or apertures for collecting sunlight and channeling it into the building interior via skytubes, lightguides or fiber-optical cables. Heliostats or large plane mirrors are used at locations with access to the Sun and they direct that light into building interiors or even city interiors. Optimally, these mirrors track the Sun.

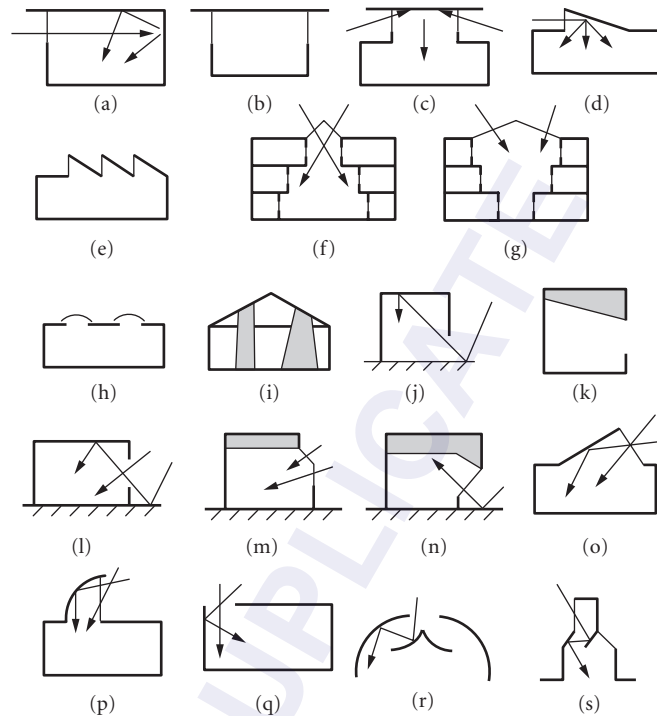


FIGURE 27 Windows for gathering daylight: (a) Unilateral, side lighting. (b) Bilateral, side lighting (c) Roof monitoring. (d) Clerestory. (e) Sawtooth skylighting. (f) Atrium. (g) Litrium. As opposed to atrium it provides best light distribution to adjacent spaces. (h) Top skylighting. (i) Skywells: straight and splayed sections. Shaded region shows the lit region. Splayed section distributes light farther, more evenly and with reduced peak brightness. (j) Window near the ground. Deep ceiling allow deeper light penetration. (k) Tilted ceiling reduced the illumination gradient from the window to interior. (l) Window in the mid-section allows direct skylight and ground reflected skylight. (m) “Greenhouse” opening for overcast sky. (n) “Overbite” opening for ground reflected skylight. (o) Tilted glazing with clerestory opening to allow sunlight from a wider range of elevations. (p) Capturing daylight via top lighting. (q) Using high reflectance wall adjacent to top lighting to provide glare free light. (r) and (s) Top lighting via lightguiding.

Figure 31 shows use of a solar light pipe (SLP) to bring sunlight deep into the building interior. The atrium is 140 ft deep, 60 ft long, and 9 ft across. Without the SLP, the view in the building interior was a dark concrete wall. A rooftop heliostat captures sunlight and directs in down into a prismatic glass cone. The glass refracts the incoming sunlight horizontally onto an outboard cylinder of open weave fabric; creating a glowing, translucent, 120-ft-long tube of diffuse sunlight. This display is visible from the 14 floors of atrium offices and also from the ground floor lobby, elevator lobbies, and the library. The SLP projects a 10-in diameter sunburst onto the lobby floor (Fig. 31*b*). Besides injecting daylight into dark spaces, the SLP’s unique design provides a compelling and a very dynamic visual focus for the atrium occupants. It constantly updates their understanding of the Sun, the sky, and the weather patterns and reconnects them to the otherwise solar rhythms of the day and seasons. At night, powerful searchlights use the “at rest” heliostat and SLP to inject a shifting palette of colored light into atrium.

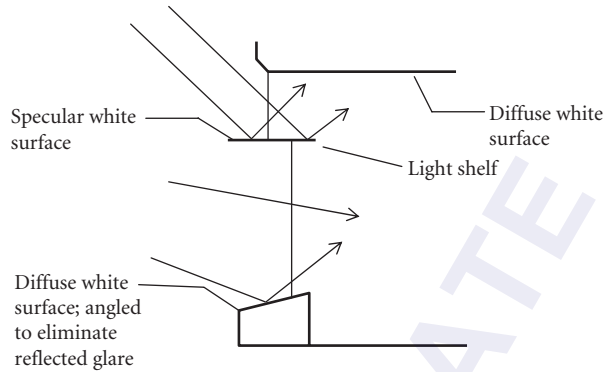


FIGURE 28 Light shelf to limit glare from direct skylight and redirect light to the ceiling.⁷² The light shelf can be curved in shape and moveable angularly. The position of the light shelf relative to the window can be interior, exterior, or both. Exterior light shelf provides shading while interior light shelf limits glare. Blinds or moveable shades can further help in glare control.

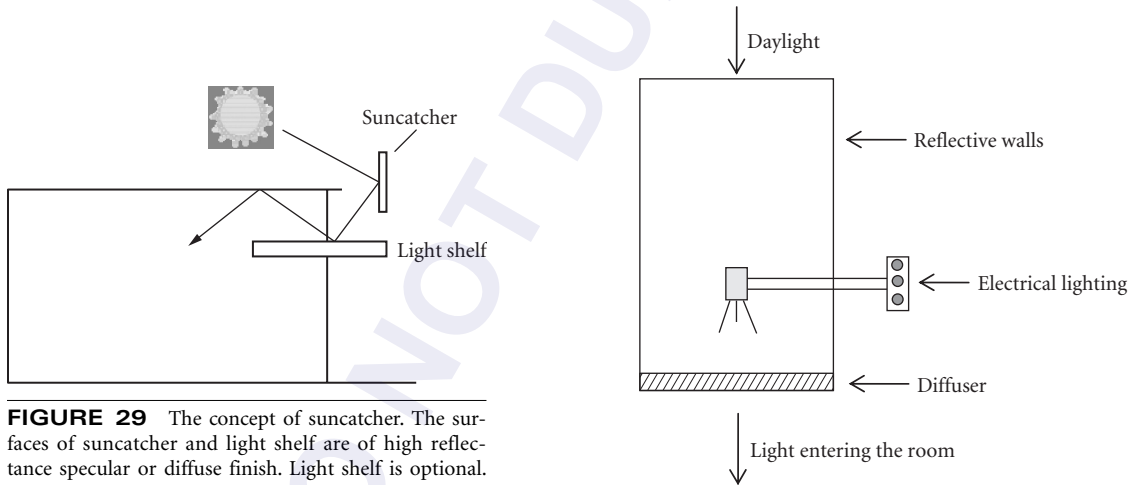


FIGURE 29 The concept of suncatcher. The surfaces of suncatcher and light shelf are of high reflectance specular or diffuse finish. Light shelf is optional. Suncatcher reduces the view, eliminates direct glare and increased illumination when the sun is at a particular location.

FIGURE 30 A skytube or skywell with integrated Daylight and electrical light.

The use of modern CAD software, especially radiosity based, makes it easy simulate the daylight (that includes sun-tracking, scattered and reflected light from the sky, ground, and neighboring buildings throughout the day and the year), its interaction with optical components such as optical fibers, lightguides, lenses, and diffusers; its penetration through the windows into the interior; impact of reflection from walls, ceilings or furniture, and interaction with interior lighting.

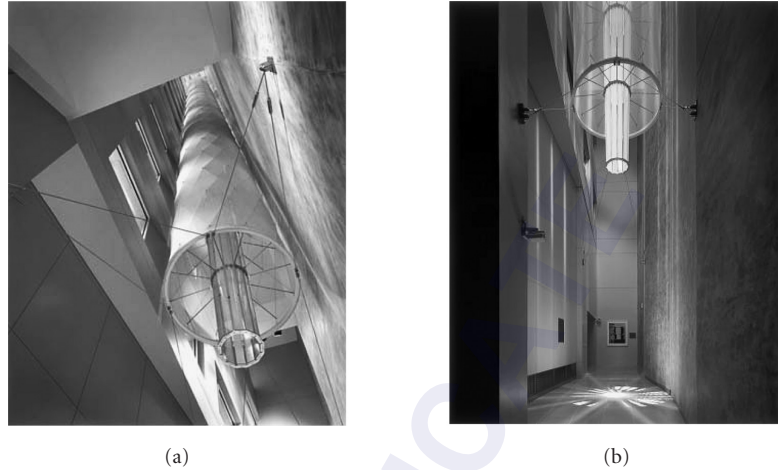


FIGURE 31 A Solar light pipe. (a) A 140-ft-tall light gathering and distributing device that presents daylight down into the core of a building that has no other access to daylight. (b) Light projected (10-in diameter) on the floor. (See also color insert.) (Photograph by Paul Warchol; Courtesy of Carpenter Norris Consulting.)

40.6 LIGHTING MEASUREMENTS

In this section, we discuss briefly the tools and techniques used for measuring light, 3D object profiles and optical properties of objects, in the context of lighting applications. The quantities most relevant for light measurement are horizontal/vertical illuminance and luminance of lamps and lit objects as a function of wavelength. To measure the optical properties of objects, luminance measurements in transmission and reflection as a function of wavelength, magnitude, angle of viewing, and direction and angle of incident light are needed. To model the objects in CAD software, measurement of 3D object profiles is needed. These measurements help in picking the lamps and their placement geometry to achieve desired lighting goals.

For measuring light from the source, we discuss the following instruments: illuminance meters and goniometers. For measuring the optical properties of objects, including luminaire optics, in reflection or in transmission, we discuss instruments such as reflectometers and luminance meters. For the measurement of luminaires conformance to the design, we present two methods: CMM and laser scanning.

Illuminance Meters

Illuminance meters or luxmeters measure illuminance in lumens/area. These are typically handheld devices, as shown in Fig. 32 that consist of a photodiode with a photopic correction filter and a cosine corrector on top of it. The photopic correction filter multiplies the incident light spectrum by the eye's photopic response to convert the incident energy in watts into lumens. The cosine corrector consists of a diffuser such as a plastic disk or flashed opal disk or a small integrating sphere with a knife edge entrance port. The goal of the diffuser is to provide a constant relative distribution of angles to the detector regardless of the incident angle of light on the diffuser. High Fresnel losses at high angles of incidence at the surface of the diffusers such as plastic disk or opal glass can still introduce errors. These errors can be minimized by allowing some light to leak through the edges of the diffuser and then using a screening ring to prevent additional error due to leaked light.⁷³



FIGURE 32 Handheld lighting measurement instruments. (a) Simple Luxmeter; (b) Luxmeter (from Labsphere) with an integrating sphere as a diffuser; and (c) luminance meter (from Konica Minolta).

Luminance Meters

Luminance meters measure luminance in lumen/area/solid angle. These are handheld instruments (Fig. 32) and are equivalent to using a lens and placing a luminance meter at the image location of the image formed by the lens of the target object. Thus it consists of a detector, a photopic correction filter or a color filter, a cosine corrector and an imaging lens. The lens is used to image the region of the object to be measured on the entrance to the detector (or the cosine corrector surface). Since the NA of the lens is known, the solid angle is known. This allows us to obtain the luminance by dividing the detected illuminance by the solid angle. Care must be taken to focus the lens only on the region whose average luminance has to be evaluated. Apertures with different field of view are provided to limit the region over which luminance is evaluated.

These instruments can be made more sophisticated by using a high-quality CCD detector array to image an entire scene and provide luminance distribution across it. The imaging lens has a flat field and is telecentric in image (detector) space.

The use of neutral density filter helps in expanding the dynamic range of the instrument. The measurement on color coordinates in various colors spaces is provided by using multiple detectors, each with a characteristic color filter. The incoming spectrum is multiplied by the transmission spectrum of the color filter for each detector. The relative signals on the detectors help in determining the color coordinates. Similarly the CCT can be evaluated.

Reflectometers

Reflectometers measure the reflectance from samples for cases such as reflectance as a function of wavelength, angle of incidence, angle of viewing and polarization. Depending upon the requirements, not all of these parameters are needed. The sample surfaces may be diffuse, specular, or a mix of two. When samples must be characterized by BRDF, scatterometers are used. Reflectometers are discussed in Chap. 35 in this volume and Scatterometers are discussed in Chap. 1 in Vol. V of this *Handbook*.

Goniometers

Goniometers are instruments that measure the irradiance or illuminance distribution of a source or luminaire. They accomplish this by measuring the illumination distribution at a number of points. By locating the detector in the far field with respect to the source or luminaire, one is actually measuring the intensity distribution. The far field is when the irradiance/illuminance distribution closely matches that of the respective intensity distribution, which means the inverse

TABLE 8 Three Types of Standard Goniophotometer with Their Respective Standard Coordinate Systems⁷⁴

	Type A	Type B	Type C
Polar axis	Vertical	Horizontal	Vertical
Vertical angle designation	Y	V	V
Horizontal angle designation	X	H	L
Range of vertical angles*	$Y \in [-90^\circ, 90^\circ]$	$V \in [-180^\circ, 180^\circ]$	$V \in [0^\circ, 180^\circ]$
Range of horizontal angles	$X \in [-180^\circ, 180^\circ]^\dagger$	$H \in [-90^\circ, 90^\circ]^*$	$L \in [0^\circ, 360^\circ]^\ddagger$
Straight ahead/down	Ahead: $Y = 0^\circ, X = 0^\circ$	Ahead: $V = 0^\circ, H = 0^\circ$	Down: $V = 0^\circ, L = 0^\circ$
Primary applications	Optical systems Automotive lighting	Floodlights	Indoor lighting Roadway lighting

*The lower angle is in the nadir direction while the upper angle is in the zenith direction.

†The lower angle is measured to the left from the perspective of the luminaire.

‡Measured from the primary axis of the luminaire.

square law applies. A good rule of thumb is ten times or greater the greatest extent of the emitter. For example, a luminaire with a 100-mm-diameter aperture indicates that measurement should be made at no closer than 1 m. Of course, better results are obtained as the distance between the source and detector is increased. By rotating either the source or detector with respect to the other, the full intensity distribution can be measured. The inclusion of a rotation device on either the luminaire or detector while the other remains fixed defines a goniometer. For the purposes of the lighting community, the luminous intensity is measured, so goniometers are better known as goniophotometers in the community.

There are three standard types of goniophotometers (Type A, Type B, and Type C), which are listed with their design criteria in Table 8. There are also three coordinate systems used for the purposes of the measurement reporting. These spherical coordinate systems are also labeled Types A, B, and C, and they typically conform to the type of goniophotometers.⁷⁴ However, individual goniophotometers may be of one type and report measurements in another coordinates system. Table 8 assumes that the goniophotometers types and coordinates systems agree.

For Types A and B goniophotometers the detector is held fixed while the luminaire is located in the rotating device. Typically guidelines or standards define the distance that the detector is located away from the luminaire. For example, for U.S. and European automotive headlamps this distance is 75 ft and 25 m, respectively. Type C goniophotometers have the detector rotate around the horizontal axis of the luminaire while the luminaire is rotated around its vertical axis. This setup is important for sources that have limitations to orientation (e.g., metal halide arc lamps). Figure 33 shows a Type C goniometer that is used to measure the luminance distribution from sources.

There are many other types of goniophotometers, especially those labeled as snapshot systems. Snapshot systems capture an “image” of the intensity distribution through one measurement. Examples include systems with several detectors; rapid scanning systems for smaller sources such as LEDs; camera-capture systems that incorporate an intermediate diffuse, reflective screen; and tapered-fiber bundles integrated to detector arrays.⁷⁴

Surface Measurement Systems

Not only are the illumination distributions, spectra, scatter distributions, and optical characteristics measured, but also the geometry of the optics. These measurements are done to characterize the geometry of the fabricated optic in regard to the design. This step is done, in conjunction with tests from the previous testing sections, when there is a disagreement between the laboratory and modeling results. There are essentially two methods in use: coordinate measurement method (CMM) and laser scanning.^{75,76} CMM uses a probe that is drawn across the surface of the optic, which provides the (x, y, z) data at a series of points on the surface. This method is analogous to using a spherometer

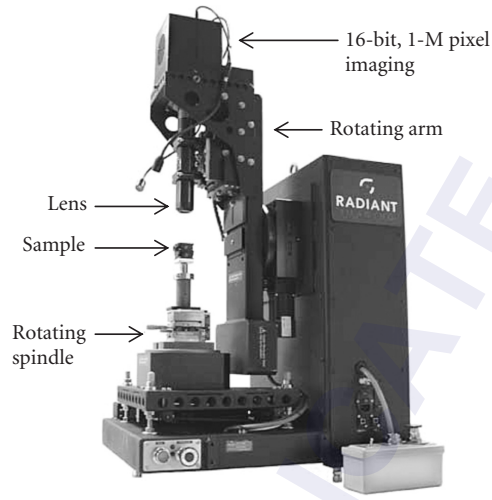


FIGURE 33 Photograph of a source measurement goniometer that is used to ascertain the luminance distribution of the source. The system wobble (electro-mechanical-software runout) is $15\ \mu\text{m}$ to allow for measuring small light sources like LED die. (See also color insert.) (Courtesy of Radiant Imaging, Inc., model SIG 400.)

to measure the sag and thus curvature of a lens surface. CMM is a contact method so the optic under test can be detrimentally affected (e.g., scratched) during measurement. Noncontact methods such as laser scanning essentially replace the mechanical probe with a laser beam. The position is measured through location of the reflected spot, triangulation, and/or time of flight as the laser is scanned across the part. Software can then determine the slopes of the test points, and thus rebuild the shape through a point cloud. If the component under test reflects well, then the surface must be coated or a fine, white powder can be placed across the surfaces. Laser scanning is a rapid method to ascertain the part shape. It provides a multitude of points in a short time; however, laser scanning accuracy is currently limited in comparison to CMM.

40.7 LIGHTING APPLICATION AREAS

In this section, we discuss several lighting application areas. We broadly divide the applications into two areas: interior lighting (office, residential, retail and healthcare) and exterior lighting (industrial and transportation). This set of applications is limited; we provide design concepts, metrics and data reference on major application areas. These ideas could be readily applied to other areas of lighting such as entertainment, sports, theatre, and so forth. Each field has its own region-specific and time-dependent guidelines and standards to implement.

The data provided here is primarily from IESNA suggested guidelines, thus exceptions are expected. For example, an office lobby is likely to have far dimmer ambient lighting than dictated by general office lighting guidelines. In essence high-end retail guidelines provide the design goals in order to provide a desired ambience. Although we have provided the guideline data to indicate the relevant specification parameters, by no means is this complete and the guidelines and standards are evolving with time, place, and technology advances. The reader must check the prevalent guidelines for the exact numbers to use. Previous sections of this chapter describe in detail the general aspects

of the lighting design process, principles, and techniques for the applications described here. Each of the following subsections provides insights into the specifics of the lighting design process for the given application.

Interior Lighting

Interior lighting includes the office, retail, residential, and health-care facility subfields. In each of these areas there are established guidelines to provide a pleasing environment to the users of the space. Of special importance are required illuminance levels, luminance ratios, and the reduction of glare. The following four subsections describe each of these application areas in more detail.

Office Lighting Modern office lighting has assumed greater importance than in the past as more people work in them. The goals of office lighting include efficient task performance, energy-efficient lighting, nonmonotonic ambient lighting that provides a balance between horizontal and vertical illuminance and minimum glare.

The task surface in offices has historically been the horizontal desk surface. However, computer monitors are ubiquitous these days and present a self-luminous vertical task surface. In addition it is specular and reflects light which can cause glare. Lighting goals are to provide adequate task illumination and adequate light in the task vicinity to eliminate eye strain due to varying brightness and disability glare. This is ensured by limiting the maximum luminance ratio of task to nontask regions. Veiling reflections (see section on glare) from computer monitor screens and paper surfaces must be avoided by taking into account source task and eye geometry. The glare from the computer monitor arises typical from ceiling source, and it is avoided by putting a limit upon the maximum ceiling luminance. Glare across the horizontal surface can be eliminated by using low luminance, wide-area lighting from overhead luminaires or special desk lamps (Fig. 24). Daylight from windows at desk level or near the ceiling can be a source of glare and steps such as daylight control via blinds or changing the task-eye geometry with respect to the glare source must be taken into account.

Balancing of daylight with electrical lighting, preferably with automatic controls, is not only energy efficient but also preferred by the workers. Therefore electrical lighting with high CRI (>70) is desirable as it mimics daylight illumination quality better.

Table 9 summarizes the common specifications and major guidelines on office lighting. The reflectivity of room parameters such as walls and ceiling is provided to aid in room modeling to aid in determining the layout of luminaires.

Retail Lighting The goal of lighting for the retail environment is to attract the customer, facilitate evaluation of the merchandise by the customer, and provide light for completing the transaction.⁷⁷ There are varying lighting methods dependent on the type of retail store, what goods are being illuminated, and the background. Table 10 provides various guidelines on lighting levels dependent on the type of store and what is being illuminated. Note that this table is not complete by any standard, so the reader is encouraged to consult the literature (see, for example, Ref. 77) for more specific information for a given retail lighting environment. Note that there are particular design issues that have varying importance levels for each type of retail outlet, such as, glare reduction in jewelry and china stores.⁷⁸

The circulation areas are those not typically used for the display of merchandise such as walkways, aisles, and foyers. The general areas are those for the generic display of merchandise. Perimeter areas are the walls where merchandise is placed for sale. Feature areas are where important displays are positioned. Horizontal illuminance values are listed for all of the columns in Table 10 except the perimeter areas, which provide vertical illuminance values. The feature areas have peak illuminance values of 5:1 to 10:1 compare to the respective general area peak illuminance. Of particular note in Table 10, the trends indicate

- Bulk stores, or those described as “big box,” tend to have much higher illuminance levels with higher uniformity across areas. This type of illumination provides the user with an abundance of light to fully inspect the merchandise, associated labeling, and compare to similar products.

TABLE 9 Summary of Specifications and Guidelines for Office Lighting

Specifications	Goals
Maximum luminance ratios:	
Task to neighboring areas	3:1 to 1:3
Task to distant regions	10:1 to 1:10
Max ambient illuminance (lx)	500
Max ceiling luminance (nits)	<1000 (in the presence of computer monitors) <425 nits does not cause glare
Reflectivity of room objects* (%)	
Ceiling	≥80
Walls	50–70
Partition	40–70
Furniture	25–45
Floors	20–40
Corridor Floors	>20 (of the illuminance of the adjacent areas)
CRI	>70
Lighting schemes	Direct, indirect and direct-indirect
Key considerations	Eliminate shadows on faces or tasks, glare control, provide a spacious ambience
Ceiling uniformity (max: min)	<8:1 (for indirect lighting only)
Common light sources	Daylighting, LED, CMH, fluorescent, and CFL
Glare sources to be eliminated or reduced	Veiling reflections from direct sources or ceilings on computer monitor screen, reflections from any specular surface including walls and desks

*Valid for matte or diffuse finish.

TABLE 10 Suggested Illuminance Levels for Circulation, General, Perimeter, and Feature Areas for Various Types of Retail Stores⁷⁷

Type of Retail Store	Circulation Area Illuminance (lx)	General Area Illuminance (lx)	Perimeter Area Illuminance (lx)	Feature Area Illuminance (lx)
Warehouse	250–300	750–850	750–850	3750–8500
Supermarket	250–300	750–1000	750–1000	3750–5000
Discount	250–300	750–850	750–850	3750–8500
Mass merchant	200–250	500–600	750–850	2000–5000
Department	200–250	400–500	500–750	2000–3500
Upscale	150–200	300–400	400–800	1500–4000
Specialty	200–250	400–500	500–750	2000–3500
Upscale	150–200	300–400	400–800	1500–4000
Boutique	80–120	200–300	200–600	1000–3000
Jewelry	80–120	200–600	200–600	1000–6000
Upscale China	80–120	200–600	200–600	1000–6000
Drugstore	250–300	750–850	750–850	3750–8500
Home	200–250	400–500	500–750	2000–3500
Furniture	80–120	200–300	200–600	1000–3000

- More specialized or exclusive, the lighting levels tend to be reduced. This type of lighting provides a more intimate environment between the customer and the salesperson with illumination to highlight the item under evaluation.

Essentially, Table 10 can be broken down into three categories: mass merchandising, department, and exclusive stores. At the lower end, mass merchandising, one typically specifies ambient illuminance

TABLE 11 Suggested Illuminance Values for Areas within a Department Store⁷⁷

Department Store Area	Horizontal Illuminance (lx)	Vertical Illuminance (lx)	Very Important Design Issues
Alteration room	500	300	Color appearance and source geometry
Dressing area	300	50	Color appearance and object modeling
Fitting area	500	300	Color appearance and object modeling
Stock room	300	50	
Sales area	300		Direct glare
General merchandise area	500		Color appearance
Feature display	2,500		Appearance of area, color appearance, direct glare, and object modeling
Display window	2,000 (day)		Appearance of area, color appearance, daylighting, object modeling, and reflected glare
	500 (night)		
Feature	10,000 (day)		
	5,000 (night)		

between 750 and 1000 lx and a CCT of 3500 to 4100 K. Department stores are in the range of 400 to 600 lx with a color temperature around 3500 K. High-end stores have ambient light levels of 150 to 300 lx and color temperatures of 2700 to 3000 K.

The values provided in Table 10 are guidelines, and values are expected to vary dependent on the background, the items being illuminated, and any external lighting. The reasoning behind this is that observers see luminance rather than illuminance. Thus, the lighting designer must take into account the reflectance, both diffuse and specular, from the merchandise and background. Thus, Table 10 assumes that there is constant reflectance between the features and the background, such that illuminance ratios of 5:1 to 10:1 are specified, when in actuality luminance ratios in this range are prescribed.

Within any type of store there are different lighting levels dependent on the application. Consider, for example, a department store, which is made up of several different environments, from the entrance areas to fitting areas, to general displays areas. Table 11 provides illuminance level guidelines for typical areas in a department store. It also includes the very important design issues for each of the retail areas.

The methods of lighting a given area are dependent on the application of the space. Ambient lighting provides the baseline lighting level, while the addition of secondary lighting units provides the increased illuminance values as listed in Tables 10 and 11. Table 12 provides a synopsis of these other application space lighting demands, including the typical luminaire used and design issues.

As previously noted the lighting designer must remain cognizant of external lighting conditions when specifying the artificial retail lighting environment. A large facet of external lighting is better known as daylighting, and of particular concern are the varying light level that is provided through the day, direct glare through windows and doors, and the strong background to incorrectly situated merchandise. Other aspects of the retail lighting environment include

- The CRI should be 70 or greater for most environments.
- Transitions between spaces in stores should have luminance ratios of no greater than 3:1 for similar neighboring spaces, greater than 3:1 when there is a distinct transition between the neighboring spaces, and 5:1 to 10:1 for abrupt transitions.
- The perimeter area illuminance should be greater than the overhead area in order to draw the attention of the shopper to the merchandise.
- Calculating baseline lighting levels for retail spaces one should use diffuse reflection values of office spaces.

Residential Lighting The goals of residential lighting are to provide ambient illumination to create a pleasing ambience due to a well-lit environment; sufficient task lighting in workspaces such

TABLE 12 Specific Application Space Lighting for Illumination of Merchandise Locations

Application Space Lighting	Typical Luminaires Used	Design Issues
Ambient Light		
Mass merchandising department	Fluorescent and halogen fluorescent, recessed fluorescent and halogen	Uniformity
High end	Recessed fluorescent and halogen, track lighting	Flexibility for change
Perimeter	Fluorescent, incandescent, or HID wall wash, track or recessed spot lighting	Uniform vertical illuminance, hidden luminaires
Rack	Recessed or track spot lighting	Direct glare reduction
Shelf	Ambient lighting with recessed surface lighting	Direct glare reduction, hidden luminaires
Counter	Focused downlighting	Direct glare reduction, 3 to 5 times ambient lighting
Mirror	Downlighting	Glare reduction, color appearance, object modeling, and consistent lighting with use of product
Showcase	Fluorescent and fiber-optic lighting	Hidden luminaires
Accent	Small (point-like) sources	Provide luminance ratios of 5:1 to 15:1
Decorative	Sconces, chandelier, table, and torchiere lamps	For high-end stores to provide a desired look and feel

as office, garage, kitchen, and workshop; and decorative lighting and accent lighting for lit object displays. A variety of lighting techniques are used to create layers of lighting that perform different lighting functions (see book by Whitehead on “Residential Lighting” in Ref. 26). Like retail lighting, residential lighting can be quite complex and creative. Although we summarize many common guidelines for residential lighting requirements in Table 13, other specialized guidelines are more relevant in some areas of the house, such as retail lighting guidelines for the kitchen, office lighting guidelines for the home office, industrial lighting guidelines for task lighting in garage or workshop and exterior lighting guidelines for landscaping and the house facade. Thus depending upon the purpose of a specific region of the house the lighting scheme must be tailored. However, the different lighting schemes used across the house must be gently blended into one another. For example, exterior lighting for residences includes the lighting of entry, walkways, and landscape. Although the primary goals are to provide direction, safety, identification and aesthetic appearance, achieving a balance between interior and exterior lighting makes the interior and exterior spaces extensions of one another. Residential lighting should be customized to maximize the comfort of the inhabitants especially when the inhabitants are disabled or elderly. For example, elderly people require much higher levels of illumination especially for task performance.

Other than those areas that require excellent task lighting such as offices or workshops, the limits on luminance ratios are relatively relaxed when compared to other lighting applications, such as the maximum residential luminance ratio of 5:1 between Task and neighboring areas is higher than that of office lighting of 3:1. The idea of limiting the maximum luminance ratio is to minimize visual discomfort caused by disability glare and adapting to variations in brightness. The integration of daylight with artificial lighting is highly desirable.

Health-Care Facility Lighting The lighting in health-care facilities, which includes hospitals, out-patient clinics, chronic and extended care centers, and other facilities, requires careful understanding of the lighting requirements for not only the patients but also the individuals working therein. The lighting must be pleasing and comforting to the patients in order to put them at ease and assist in their healing. The patients and visitors have a wide range of ages, with the majority being elderly;

TABLE 13 Summary of Specifications and Guidelines for Residential Lighting

Specifications	Goals
Maximum luminance ratios	
Task to neighboring areas	5:1 to 1:5 (in general) 3:1 for demanding tasks such as sewing
Task to distant regions	10:1 to 1:10
Luminaire to ceiling	20:1
Hallway or stair to adjacent area	1:5
Reflectivity of room objects* (%)	
Ceiling	60–90
Walls, curtains, draperies†	35–60
Floors†	40–70
Average luminance for luminaire (nits)	1700
Maximum luminance for luminaire (nits)	2700 (in utility areas)
CRI	>80 (kitchen and clothing closets)
Lighting schemes	Direct, indirect, and direct-indirect
Key considerations	Eliminate shadows on faces or tasks, glare control, providing a spacious and cozy ambience as well cozy as desired
Common light sources	Daylighting, LED, fluorescent, and CFL.
Glare sources to be eliminated or reduced	Direct glare from light sources, veiling reflections, indirect glare from shiny objects

*Valid for matte or diffuse finish.

†Reflectance of walls and floor can be increased by 40% and 25%, respectively, to improve visual task lighting where needed.

therefore, there is a large variance in the response to lighting, especially increased glare sensitivity, loss of contrast sensitivity, the need for higher lighting levels, and slow adaptation to changes in brightness as one ages.⁷⁹ The lighting levels also need to provide for the medical staff such that they can effectively carry out their job—from meticulous and demanding surgery to patient interviews and diagnoses to manning the check-in desk. Finally, since in a number of these facilities patients may be there for extended periods, circadian system illumination levels that conform to the human biological clock are the norm. As can be seen, there is an extensive range of tasks, observers, and daily requirements for illumination in health-care facilities, thus, making the lighting design a challenging assignment.

Foremost is the need to specify the light requirements for a given location based on the tasks to be performed there. Table 14 provides a synopsis of the illuminance level guidelines and important design issues for a number of locations in health-care facilities; however, since these are only guidelines, controls for the area illumination are often available to the medical staff and/or patients. This flexibility in the control of lighting levels based on the specific function of the environment and even mood of the occupants is important in health-care facilities. Note that there are many other areas and functions than can be included in Table 14, so the reader is encouraged to see Ref. 79 for this additional information.

The operating room environment is especially challenging since the lighting is preferentially directed to the task area; however, this has the potential of creating large luminance ratio variation within the room. Practices suggest three regions within the operating room with the following luminance ratios: task area to the surgical table of 3:1 or less and task area to the background (i.e., walls) of 10:1 or less. Additionally, there are specific guidelines for the finishes for the surfaces within the surgery theatre as provided in Table 15. In most cases all surfaces are white or pastels with a matte finish to reduce reflected glare. The lighting in an operating room is provided by a multitude of sources including ambient and even daylighting, directed spotlights, surgeon headlamps, and increasingly fiber-optic lighting.

Concurrent to the design guidelines of Table 14, the lighting designer must remain cognizant of the specifications for reflectances of the walls, floor, ceiling, and other objects that occupy the design area.

TABLE 14 Suggested Illuminance Values and Other Criteria in Health-Care Facilities⁷⁹

Health Care Facility Area and Current Function	Horizontal Illuminance (lx)	Vertical Illuminance (lx)	Important Design Issues
Patient room			
Normal	200+, as home	30+, as home	Patient control, CRI >80, daylighting
Examination	1000	300	Glare reduction, CRI >80, CCT >3000 K, doctor control
Observation (night)	30	30	Red/amber CCT, no higher than 18" off floor
Nursing station	300	50	Glare reduction, CRI >80
Critical care unit			
Normal	300	50	Patient control, CRI >85, daylighting
Examination	1000	300	Glare reduction, CRI >85, CCT >3000 K, doctor control
Nursery	100	30	Glare reduction, CRI >80, lighting control, daylighting
Mental health center	As per other functions	As per other functions	Daylighting, CRI >80, CCT of 4100–5000 K with fluorescent, 3500 K otherwise, glare reduction
Operating room			
Normal	1000	500	CRI >85, matched light
Table	25000+ in 20-cm spot (adjustable)	1000	Sources for illumination, shadow, and glare reduction, doctor control, CCT of 3500–6700 K, high uniformity
Dental unit			
Normal	300	50	
Examination	24000+ in central spot	500	As operating room
Radiography unit	50, but depends on test	30, but depends on test	CRI >80, glare reduction, doctor control
Pharmacy	300	100	Glare reduction, high uniformity, CRI >80

TABLE 15 Suggested Reflectances for the Various Objects in a Health-Care Facility Room for a General Application and the Operating Room Environ⁷⁹

Surface	General Reflectance	Operating Room Reflectance
Ceiling	70–80%	90–100%
Walls	40–60%	60%
Furniture/fabrics	25–45%	0–30%
Equipment	25–45%	25–45%
Floors	20–40%	20–30%

Table 15 provides this data for two locations, general and operating room. Using these values in conjunction with the optical design process one can find suitable illumination configurations to provide the Table 14 guidelines.

Industrial Lighting The goals of industrial lighting are to provide energy efficient lighting with adequate task and ambient lighting, direction, safety, and visual comfort. Many common requirements for industrial lighting requirements are summarized in Table 16.

TABLE 16 Summary of Commonly Suggested Specifications and Guidelines for Industrial Lighting

Specifications	Goals
Maximum luminance ratios	
Task to neighboring areas	3:1 to 1:3 (in general)
Task to distant regions	10:1 to 1:10
Luminaire (including daylight sources) to adjacent surfaces	20:1
Anywhere within the FOV	40:1
Reflectivity of objects* (%)	
Ceiling	80–90
Walls	40–60
Floors	≥20
Desk/bench tops, machines, equipment	25–45
CRI	>65
Lighting schemes	Direct, indirect, and direct-indirect
Key considerations	Eliminate shadows on faces or tasks, provide high illuminance uniformity, sufficient illuminance, and glare control
Common light sources	Daylighting, LED, fluorescent, HID, and CFL
Glare sources to be eliminated or reduced	Direct glare from light sources, veiling glare, indirect glare from shiny objects
Direct view of the luminaire (deg)	>25 (preferably >45)

*Valid for matte or diffuse finish.

The ambient lighting is provided by daylighting and/or large area, overhead, wide angle luminaires (direct and semidirect). Task lighting is provided by fixed and portable direct or diffuse light sources. Backlights are used for translucent task surfaces. Direction of illumination with respect to the view angle(s) is important in tasks where surface features of the object cause shadows to enhance depth perception. Illumination at grazing incidence is used to highlight specific feature that can scatter light and render themselves visible. To emphasize specular but uneven surfaces, specific patterns of light (like bright and dark lines) can be projected onto the task surface to be reflected into the viewer's eyes. For lighting in regions where there is a high density of workstations or areas where there are several similar process, a high level of uniformity in horizontal illuminance is recommended. In such areas, variation of horizontal illuminance should be less than 1/6th the average horizontal illuminance. Otherwise, the variations in luminance across the work space must be within the luminance ratios prescribed in Table 16. The level of horizontal illuminance needed varies with task. In industrial lighting, the quality of horizontal illuminance is of special importance as efficient and safe task performance is needed. The reader is referred to IESNA published guidelines for horizontal illuminance values for a wide variety of tasks in different industries.⁸⁰

In addition to providing the task lighting and ambient lighting, it is often necessary to provide additional lighting dedicated for emergency, safety, and security within the industrial complex and its exteriors.

Visual discomfort must be avoided especially during task performance and in situation where safety could be compromised. Special attention should be given to situations causing veiling reflections and glare. The various limits placed on the luminance ratios in Table 16 between task and non-task regions help eliminate the impact of glare.

Exterior Lighting

Of course the primary aspect of exterior lighting is to provide illumination during hours of darkness. The lighting not only provides illumination for general use, it also provides safety and security; indication of direction of travel for paths, roadways, and so forth; and architectural enhancement.

External lighting is attempting to replicate the illumination of the sun; however, it can in no fashion accomplish this feat. The illumination provided by any lighting source cannot match that of the sun, which means that the sky will appear dark rather than blue; numerous sources are required to provide the necessary illumination level so glare from the many sources is a major factor; distance from the light sources, mesopic vision and even scotopic levels may be demanded; and the many sources can confuse the viewer such that objects can be difficult to discern from one another.⁸¹ Additionally, this section is rather broad in scope, including design aspects such as roadway lighting, path lighting, outdoor event lighting, and façade illumination. Thus, the reader may need to consult specific literature for a certain type of lighting. Please consult the transportation lighting section for further information on vehicular and roadway. In the next two subsections details are presented about the issues present for external lighting design and design examples, excluding the transportation applications discussed later.

External Lighting Design Issues There are several topics that a lighting designer for external spaces must consider. First and foremost is the specific application guidelines (for example, see the section on illuminating roadways), glare issues, light pollution and trespass, and perception issues. All of these factors are interrelated, and the subsections herein provide insight into them.

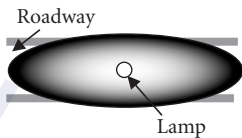
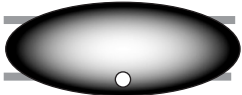

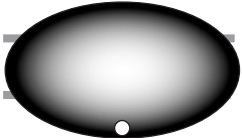
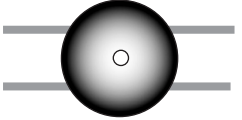
Glare A primary issue of external lighting is glare. Since users of a space in the hours of darkness are relying on nonphotopic vision, glare can blind the viewer as one approaches a bright source thus not only causing the scotopic vision receptors (e.g., rods) to naturally saturate but also the photopic vision receptors (e.g., cones). This “blinding” is due to disability glare, which will hide objects or reduce their contrast. Other levels of glare can also inhibit observation: discomfort glare, which is attributed to a large variance on and around the object under view; and annoyance glare, which is light that in the opinion of the observer should not be present, such as light trespass.

The major source of glare is directly from the source in comparison to the object under observation. Comparison between the luminaire luminance to the object’s illuminance is the factor that often defines the level of glare, from disabling to discomfort. Depending on the orientation of the object one either compares the horizontal or vertical surface luminance to the luminaire luminance. The horizontal luminance is used for horizontal surfaces such as walkways and roads, while the vertical luminance is used for vertical surfaces such as building façades, people, and structures.

Light Pollution and Light Trespass Associated to glare is light pollution and light trespass. Light pollution, also called sky glow, is light that is directed upward to the atmosphere.⁸² This sky glow hides stars from observation, and on cloudy days provides a glow, typically red in hue, which can distract from one’s appreciation of the view of their surroundings. Astronomers have been particularly vocal in the reduction of light pollution, and there is increasingly consideration of energy efficiency demands. In 2008 it was hypothesized that over \$10B U.S. was wasted through light pollution,⁸³ which means that this amount has only increased since then. Like light pollution, light trespass is unwanted light, but in this case it is light that illuminates objects outside of the intended illumination region.⁸⁴ Such stray light enters through windows, illuminates someone’s property, or can impair the vision of drivers. Light pollution and trespass arise from improperly designed luminaires. There are increasing regulations for the design of outdoor lighting to alleviate such concerns. A primary correcting factor is the use of luminaire cutoffs as presented in the section on luminaire design.

External Lighting Example Pole lighting is the most common form of outdoor illumination providing a significant amount of light, safety and security, and indication of the path of travel for walkways, roadways, and so forth. In the United States, there are five primary forms of classifications of pole-mounted luminaires, but there are numerous variations based upon the lighting requirements, the shape of the illumination region, and illumination level demands.⁸⁵ Table 17 provides descriptions of the five types, the typical application(s) of the type, and an iso-illuminance plot for each of the types. Per the iso-illuminance plots of Table 17, the larger the illumination distribution the more ground is illuminated, which means not only the roadway is illuminated but also sidewalks, shoulders, and other neighboring areas. Thus, Type IV and V luminaires tend to have more light trespass.

TABLE 17 Description, Applications, and Iso-Illuminance Plots for the Five Primary Luminaire Types Used for Pole-Mounted Outdoor Lighting⁸⁵

Type	Description	Applications	Iso-Illuminance Plot
I	Little setback to the road or mounted over the roadway; narrow oval illumination distribution	Roadways	
II	Moderate setback to pedestrian area; moderate oval illumination distribution	Pedestrian areas	
III	Large setback to the road; Moderate elliptical illumination distribution	Pedestrian areas roadways, parking lots	
IV	Great setback to the roadway; large elliptical illumination distribution	Roadways	
V	Mounted over pedestrian area rotationally symmetric illumination	Pedestrian areas parking lots	

Transportation Lighting

Transportation lighting includes the subfields of vehicular and roadway lighting. In both cases rather than just guidelines, there are many stringent governmental standards that must be met by the illumination systems. For example, there are governmental standards for traffic lights and vehicular taillights and headlights. If these standards are not met, then the lighting systems are not legal for use on our roadways.

Vehicular Lighting Vehicular lighting is the external illumination aspects of vehicles, including automobiles, motorcycles, snowmobiles, emergency vehicles, airplanes, ships, and heavy machinery including construction, industrial and agricultural. Ground vehicles are especially important due to their pervasiveness in society. For the remainder of this section, we highlight ground vehicle standards, but there are similar standards for other types of vehicles. In the United States, the U.S. Federal Government standardizes the lighting requirements through the Federal Motor Vehicle Safety Standards (FMVSS) from the offices of the National Highway Traffic Safety Administration (NHTSA) within the Department of Transportation (DOT). In Europe and a number of other countries spanning the globe, the United Nations Economic Commission for Europe (UNECE or better known as ECE) sets the standards. The FMVSS calls upon the standards delineated by the Society of Automotive Engineers (SAE) to provide the explicit lighting requirements for distinct lighting systems. While FMVSS 108 provides the framework for lighting systems on ground vehicles within the United States,⁸⁶ the SAE provides the accepted standards to design into the lighting systems.

For example, SAE standard J581 provides the upper-beam requirements to be an accepted high beam on U.S. roads,⁸⁷ while R113 is the accepted ECE standard.⁸⁸ The associated low-beam standards are SAE J582 in the U.S.⁸⁹ and ECE R112 in Europe,⁹⁰ but there is currently an active dialogue to harmonize the standards between the U.S., European, and Japanese markets. The SAE J1735 standard is currently addressing the low-beam harmonization, and in time it will also address the high-beam requirements. The end results will be, excepting the inherent difference of left-hand (e.g., United Kingdom) and right-hand (e.g., United States) traffic, harmonization has the goal of making the lighting standards the same at as many places possible across the globe. This increased standardization means that the design and fabrication costs will be reduced for manufacturers.

There are essentially two types of lighting on a ground vehicle: forward lighting for illuminating the road surface and nearby surroundings and indicator and warning lighting including turn signals, brake lights, side markers, and tail lights. Each vehicular light system is defined by its own set of regulations including the distribution of light, the color and characteristics of the protective lens, mounting requirements, and a number of enviromechanical tests including vibration, dust, moisture, corrosion, and warpage. In the United States, the test procedures and protocols are typically governed by SAE J575, while each individual standard provides the photometric requirements. In the United States, the photometric requirements are the luminous intensity distribution at a distance of 18.3 m from the lamp, while the ECE requirements are for the illuminance distribution at a distance of 25 m. The typical instrument to make this measurement is a goniometer that holds the lamp and allows it to be rotated so that a different angle is incident on a fixed photodetector at the end of an absorbing light tunnel. By rotating the lamp within the goniometer and making a series of measurements the full distribution of light can be determined; however, rather than measuring across the whole lit region, the standards dictate a series of test points and test areas that must be measured. In the remainder of this section, we highlight two applications: low-beam headlamp and stop lamp. In each subsection, the requirements for both U.S. and ECE standards are provided, a depiction of a typical distribution of light, discussion of design strategies, and additional optical requirements of the standards. Note that while we highlight these two applications and their respective standards, there are numerous other optical standards governing the external lighting systems on vehicles (see Refs. 86 to 92) for further information about these additional standards).

Design and Standards for a Low-Beam Headlamp A headlamp, as shown in Fig. 34a, must illuminate the road surface for the driver, illuminate to the sides for both the driver and any other observer, provide a low level of illumination to oncoming drivers, and in all cases remove the formation of hot spots above the horizon such that glare is not a concern to oncoming traffic. Of note in Fig. 34a is

- High-beam luminaire:
 - Rightmost recess of the lamp.
 - Faceted reflector, but other options include smooth, tailored reflector (e.g., NURBS surface as designed in CAD); faceted lens in conjunction with smooth reflector (e.g., parabolic); and projection lens in conjunction with reflector.

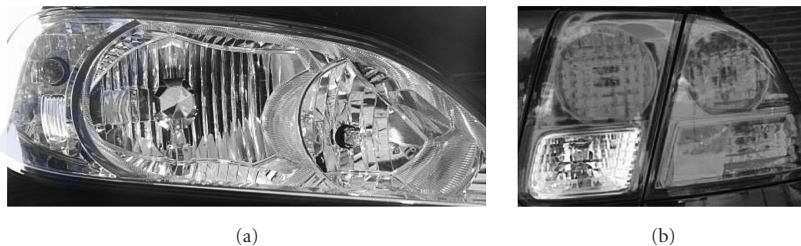


FIGURE 34 (a) A faceted headlamp including high-beam (right), low-beam (middle), and turn signal (left) luminaire. Note the yellowish tinge of the turn signal, which is due to the coating placed on the bulb used therein. (b) A faceted taillight including the following functions: tail (upper left), stop (upper right), turn signal (lower right), reflex reflector (lower middle), and backup (lower left). (See also color insert.)

- Filament source, but other options include high-intensity discharge lamp or array of white-light LEDs.
- Low-beam luminaire:
 - Middle recess of the lamp.
 - Faceted reflector, but other options are as the high-beam luminaire above.
 - Filament source, but other options are as the high-beam luminaire above.
 - Bulb shield: reflective structure in the middle of the lamp and covers direct view of the filament source. The bulb shield greatly alleviates direct light above the horizon.
- Turn signal luminaire:
 - Leftmost recess of the lamp.
 - Faceted reflector, but other options are as the high-beam luminaire above.
 - Filament source, but other options are as the high-beam luminaire above. Note that the bulb has a yellow coating placed on its glass envelope, which then provides the yellowish appearance of such lamps. This coating provides the color as specified by its respective standard.

Headlamp reflectors have a parabolic base shape in order to provide a high degree of collimation in the forward direction. The reflector, both faceted and smooth, is then deformed to provide the distribution that meets standards. Designing such deformations can be difficult, but there are software codes to assist in the process. Design guidelines include a goal of having the emitted radiation only incident once on the reflector; avoid secondary interactions with the bulb shield, bulb, or reflector shelves (i.e., the reflector sides, as shown in Fig. 34a); and angling intersegment fillers (i.e., spaces between the facets) such that they cannot be directly seen from the source emission region. In order to avoid secondary interactions with the bulb and its respective shield, one typically angles the reflected light across the luminaire, except that striking near the reflector near the source. The latter is directed away from the source. In systems with projection and faceted lenses, the reflector is there to capture bulb emission, while the lens provides the required distribution that meets standards. Due to demanding illumination requirements, it is best if most of the bulb emission is incident on the reflector; therefore, in almost all cases the filament or arc is oriented along the axis of the reflector (i.e., orthogonal to the filament as shown in Fig. 13). This axial filament orientation also ensures less light goes above the horizon in the low-beam luminaire.

Table 18 provides the SAE J582 luminous intensity requirements, while Table 19 provides the ECE R1 12 requirements. Figure 35 shows a typical SAE luminous intensity distribution that meets the requirements of Table 18. Figure 36 shows a typical ECE illuminance distribution that meets the requirements of Table 19. The SAE standards of Table 18 and Fig. 35 are in the units of candelas

TABLE 18 SAE J582 Standard Photometric Requirements for Auxiliary Low-Beam Lamps⁸⁹

Test Point (Deg, $\pm 0.25^\circ$)	Max Luminous Intensity (cd)	Min Luminous Intensity (cd)
10U to 90U	75	—
1.5U–1L to L	300	—
1.5U–1R to R	300	—
0.5U–1L to L	400	—
0.5U to 1R to 3R	400	—
0.5D–1R to 3R	25,000	2,000
0.5D–1L to L	10,000	—
5D–4R	—	3,000
5D–4L	—	3,000
1D–1R	—	10,000
3D–3R	5,000	—
4D–V	3,000	—
2.5D–15L	—	1,500
2.5D–15R	—	1,500

TABLE 19 ECE R112 Standard Photometric Requirements for Class A Passing (i.e., Low-Beam) Lamp for Right-Hand Traffic⁹⁰

Test Point Label	Horizontal Location (mm)	Vertical Location (mm)	Max. Illuminance (lx)	Min. Illuminance (lx)
B50L	1500L	250U	0.4	—
75R	500R	250D	—	6
75L	1500L	250D	12	—
50L	1500L	375D	15	—
50R	750R	375D	—	6
50V	V	375D	—	—
25L	3960L	750D	—	1.5
25R	3960R	750D	—	1.5
Zone III	See Fig. 36	See Fig. 36	0.7	—
Zone IV	2250L to 2250R	375D to 750D	—	2
Zone I	L to R	750D to D	20	—

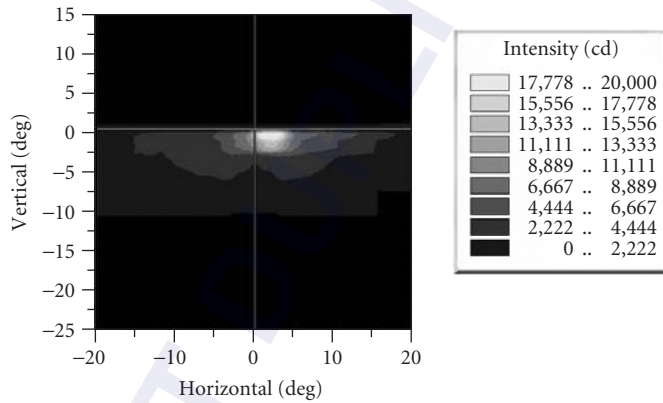


FIGURE 35 Luminous intensity (cd) distribution for the SAE low-beam requirements of Table 18. (See also color insert.)

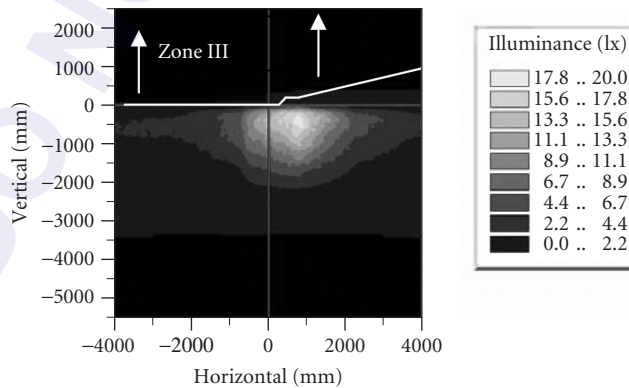


FIGURE 36 Illuminance (lx) distribution for the ECE passing/low-beam requirements of Table 19. (See also color insert.)

(cd, lumens/steradian). The test points are given in degrees and a letter designation, which mean U = up direction, D = down direction, L = left direction, and R = right direction as measured from the point (H, V), which is the center point where H = horizontal and defines the horizon, and V = vertical and defines the lane of traffic. The ECE standards of Table 19 and Fig. 36 are in the units of lux (lx, lumens/m²), and, while the letter designations still hold, the test point locations and zones are given as position coordinates in millimeters (mm).

Design and Standards for a Stop Lamp A taillight, as shown in Fig. 34*b*, is typically comprised of a number of different functions, such as a stop light, turn signal, backup lamp, and so forth. The taillight of Fig. 34*b* has the following optics in its functions.

- Tail lamp (upper left, red): for night-driving conditions or when headlamps are on. Indicates to following traffic the presence of the vehicle during reduced lighting conditions. The governmental standards are SAE J585 and ECE R7.
- Stop lamp (upper right, red): indicates when the driver has applied the brakes. This luminaire is also lit for night driving conditions but a lower output level. The governmental standards are SAE J586 and ECE R7.
- Turn Signal Lamp (lower right, red): indicates the driver is to make a turn in the designated direction. The governmental standards are SAE J588 and ECE R6.
- Reflex reflectors (middle right, red): this area, directly neighboring the turn signal is comprised of prism structures that provide retroreflection to the driver of following vehicles. It is important for dark-driving conditions to highlight the presence of this vehicle to following drivers or indicate its presence when the vehicle is not in operation. Note that this function of the taillight is passive in the sense that no internal light source is used. The governmental standards are SAE J594 and ECER3.
- Backup lamp (lower left, white): this luminaire indicates when the vehicle is in reverse. The governmental standards are SAE J593 and ECE R6.

In Fig. 34*b*, faceted reflectors are used for all the active luminaires. Typically, a transverse filament is used, as per Fig. 13, since the illumination standards for taillight functions are quite broad angularly. Thus, direct radiation from the filament is useful to filling in the required light distribution, while the reflected component of the illumination fills in the required hot spots. LEDs are replacing filament lamps to become the norm for most taillight functions. For LED sources, refractive optics that also employ total internal reflection provide a better means to meet the illumination requirements.

In the remainder of this section, we focus our attention on the design and standards of the stop lamp function, but the other taillight functions have similar requirements. The stop lamp is to inform following vehicles and pedestrians that the vehicle is slowing. In Fig. 34*b*, the stop lamp is located on the upper right of the taillight. Table 20 provides the SAE J586 luminous intensity requirements,⁹¹ while Table 21 provides the ECE R7.⁹²

Figure 37 shows a typical SAE luminous intensity distribution that meets the requirements of Table 20. Figure 38 shows a typical ECE illuminance distribution that meets the requirements of Table 21. The SAE standards of Table 20 and Fig. 37 are in the units of candelas (cd, lumens/steradian). The ECE standards of Table 21 and Fig. 38 are in the units of candelas (cd, lumens/steradian). In all cases, the letter designations as per the previous section (headlamps), but in this case the angles are with respect to rear direction of the vehicle.

Roadway Lighting There are multiple subfields that make up roadway lighting: street lighting, roadway signage, tunnel lighting, and integration of collocated pedestrian, bike, and analogous areas. The goal of all forms of roadway lighting is to reduce the potential of accidents, aid in the flow of traffic, provide a higher level of safety and security, and assist in commerce during hours of darkness. There are a multitude of light sources that affect roadway lighting, including the external lighting from vehicles (see section on vehicular lighting); direct roadway lighting, which is the subject of this section; traffic lights, as governed in the United States by the International Transportation Engineers

TABLE 20 SAE J586 Standard Photometric Requirements for a Stop Lamp*

Zone	Test Point (Deg.)	1 Lit Section	2 Lit Sections	3 Lit Sections
		Min. Luminous Intensity (cd)	Min. Luminous Intensity (cd)	Min. Luminous Intensity (cd)
I	10U-5L	9.6	11.4	13.2
	5U-20L	6	7.2	9
	5D-20L	6	7.2	9
	10D-5L	9.6	11.4	13.2
	Zone Total	52	62	74
II	5U-V	18	21	24
	H-10L	24	28.2	33
	5D-V	18	21	24
	Zone Total	100	117	135
III	5U-V	42	49.2	57
	H-5L	48	57	66
	H-V	48	57	66
	H-5R	48	57	66
	5D-V	42	49.2	57
Zone Total	380	449	520	
IV	5U-V	18	21	13.2
	H-10R	24	28.2	9
	5D-V	18	21	9
Zone Total	100	117	13.2	
V	10U-5R	9.6	11.4	74
	5U-20R	6	7.2	24
	5D-20R	6	7.2	33
	10D-5R	9.6	11.4	24
	Zone Total	52	62	135
Maximum	Any point above	300	360	420

*The stop lamp is comprised of up to three distinct lit sections over the extent of the taillight for the stop lamp function. Each zone has a number of test point minima that must be realized and the summed total of all test points within a zone. The last line of the table indicates the maximum luminous intensity that can be measured at any test point.⁹¹

TABLE 21 ECE R7 Minimum and Maximum Photometric Requirements for a Stop Lamp*

Test Point (Deg.)	1 Lamp Illumination Level		2 Lamp Illumination Levels	
	Min. Luminous Intensity (cd)	Max. Luminous Intensity (cd)	Min. Luminous Intensity (cd)	Max. Luminous Intensity (cd)
10U-5L	12	37	6	16
10U-5R	12	37	6	16
5U-20L	6	18.5	3	8
5U-10L	12	37	6	16
5U-H	42	129.5	21	56
5U-10R	12	37	6	16
5U-20R	6	18.5	3	8
V-10L	21	64.75	10.5	28
V-5L	54	166.5	27	72
V-H	60	185	30	80
V-5R	54	166.5	27	72
V-10R	21	64.75	10.5	28
5D-20L	6	18.5	3	8
5D-10L	12	37	6	16
5D-H	42	129.5	21	56
5D-10R	12	37	6	16
5D-20R	6	18.5	3	8
10D-5L	12	37	6	16
10D-5R	12	37	6	16

*The two categories are for a lamp that is either lit or not lit (e.g., 1 lamp illumination level) and for a lamp that has an unlit, partially lit, and fully lit state (e.g., 2 lamp illumination levels).⁹²

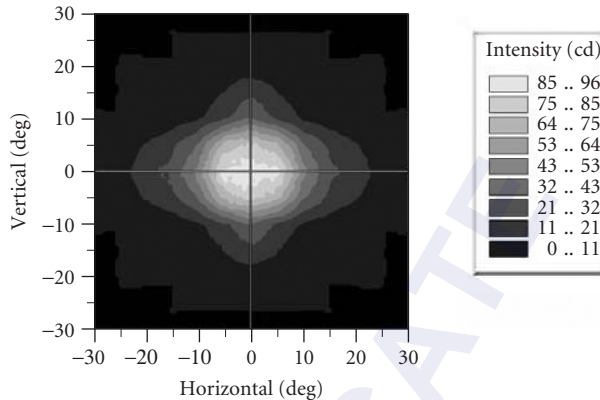


FIGURE 37 Luminous intensity (cd) distribution for the SAE stop lamp requirements of Table 20 (1 lit section). (See also color insert.)

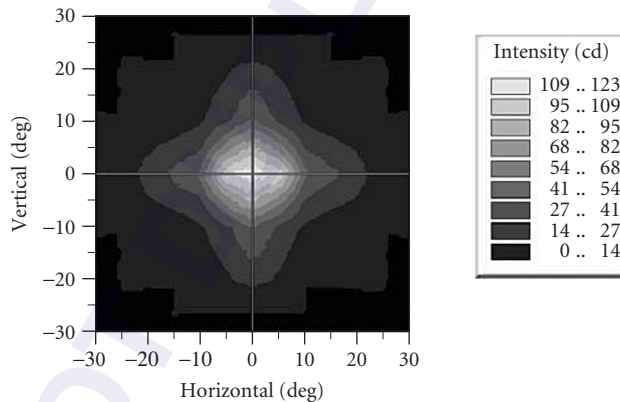


FIGURE 38 Luminous intensity (cd) distribution for the R7 stop lamp requirements of Table 21 (1 lamp illumination level). (See also color insert.)

society; and lighting from other sources such as residential, industrial, and retail. An added difficulty of roadway lighting is the harsh conditions in which they reside. There are stringent maintenance demands that include replacement of sources, cleaning of the optics, and trimming of foliage around the luminaire.⁹³

Street Lighting There are essentially three metrics for defining the lighting of a roadway: illuminance (lux, lumens/m²), luminance (nit, lumens/sr/m²), and small target visibility (STV).⁹⁴ Illuminance modeling provides the level of illumination across a surface. Luminance modeling predicts the level of light (i.e., glare), both direct and reflected, that is directed to a driver. STV is the visibility of a target array (18 × 18 cm² and 50 percent reflective) on the road. Table 22 provides the illuminance, luminance, and STV guidelines for a number of road types, four road surfaces in dry conditions, and the interaction level with pedestrians. The illuminance and luminance ratios of Table 22 provide limits such that disabling glare does not blind drivers or pedestrians. Table 23 qualifies the road types, the four road surfaces, and the pedestrian interaction levels of Table 22.

TABLE 22 Illuminance, Luminance, and STV Guidelines for the Type of Road, the Road Surface Classification, and the Interaction Level with Pedestrians⁹⁴

Road	Ped. Level	Road Surface				Illum. Uni. Ratio E_{avg}/E_{min}	Avg. Lum. L_{avg} (cd/m ²)	Avg. Lum. Uni. Ratio L_{avg}/L_{min}	Max. Lum. Uni. Ratio L_{max}/L_{min}	Veiling Lum. Ratio L_{Vmax}/L_{avg}	STV Weight Avg. VL	STV Avg. Lum. L_{avg} (cd/m ²)		STV Lum. Uni. Ratio L_{max}/L_{min}
		R1 (lx)	R2/R3 (lx)	R4 (lx)								<7.3 m	≥7.3 m	
Class A Freeway	NA	6	9	8	3	0.6	3.5	6	0.3	3.2	0.5	0.4	6	
Class B Freeway	NA	4	6	5	3	0.4	3.5	6	0.3	2.6	0.4	0.3	6	
Expressway	High	10	14	13	3	1	3	5	0.3	3.8	0.5	0.4	6	
	Medium	8	12	10	3	0.8	3	5	0.3	3.8	0.5	0.4	6	
	Low	6	9	8	3	0.6	3.5	6	0.3	3.8	0.5	0.4	6	
Major	High	12	17	15	3	1.2	3	5	0.3	4.9	1	0.8	6	
	Medium	9	13	11	3	0.9	3	5	0.3	4	0.8	0.7	6	
	Low	6	9	8	3	0.6	3.5	6	0.3	3.2	0.6	0.6	6	
Collector	High	8	12	10	4	0.8	3	5	0.4	3.8	0.6	0.5	6	
	Medium	6	9	8	4	0.6	3.5	6	0.4	3.2	0.5	0.4	6	
	Low	4	6	5	4	0.4	4	8	0.4	2.7	0.4	0.4	6	
Local	High	6	9	8	6	0.6	6	10	0.4	2.7	0.5	0.4	10	
	Medium	5	7	6	6	0.5	6	10	0.4	2.2	0.4	0.3	10	
	Low	3	4	4	6	0.3	6	10	0.4	1.6	0.3	0.3	10	

TABLE 23 Road Types and Surfaces, and Pedestrian Interaction Levels as per Table 22⁹⁴

Road Type	
Class A Freeway	Divided high traffic highways with full access control
Class B Freeway	All other divided highways with full access control
Expressway	Divided highways with limited access control
Major	Primary roadways within and leaving metropolitan areas
Collector	Service roads connecting major and local roadways
Local	Provide direct access to residential, retail, and industrial areas
Road Surface Class	
R1	Portland cement concrete and asphalt with at least 12% artificial brightener aggregates; treat as diffuse
R2	Asphalt road surface with a minimum 60% gravel aggregate; treat as both diffuse and specular
R3	Asphalt with dark aggregates; slightly specular
R4	Smooth asphalt surface; specular
Pedestrian Interaction	
High	Large amount of pedestrian traffic during hours of darkness (100 or more pedestrians in one block during an hour): retail and entertainment areas
Medium	Moderate amount of pedestrian traffic during hours of darkness (11 to 100 pedestrians in one block during an hour): office area, apartment area, older city neighborhoods
Low	Little pedestrian traffic during hours of darkness (10 or less pedestrians in one block during an hour): rural and suburban areas

The luminaire cutoff classification, as described earlier, is used extensively in the design of roadway lighting. The luminaire cutoffs suppress glare to drivers and pedestrians while also alleviating light pollution and light trespass.

There are a number of other road types (e.g., sidewalks, bike paths, and intersections), environmental conditions (e.g., fog, rain, and wet roads), and special considerations with interaction with pedestrians. The reader is encouraged to consult Ref. 94 for these additional circumstances.

Sign Lighting A number of governmental standards have been developed for the lighting of road signs, but a set of guidelines have been developed by the IESNA.⁹⁵ Light for signs can be from external sources (e.g., car headlamps) or an internal source for a transmissive sign. Externally lit signs either make use of associated lights that illuminate it or retroreflectors that reflect back to the observer. In the United States the Federal Government provides the standards that must be met for lit roadway signage.⁹⁶

Tunnel Lighting Tunnel lighting is an important factor to increase drive safety while also allowing drivers to maintain their speed. In the United States a structure is considered a tunnel when it is greater than 25 m in length.⁹⁷ Distances greater than this require additional lighting to supplement any available daylight. Tunnel lighting is broken up into a number of regions including the approach, adaptation, threshold, transition, interior, and exit zones. Each of these zones has different guidelines to ensure driver comfort; however, these guidelines are dependent on the time of the year, the average speed of traffic, and the presence of oncoming traffic in undivided tunnels. For more information please consult Ref. 97.

40.8 ACKNOWLEDGMENTS

We are grateful to Jim McGuire and Kevin Thompson for providing valuable feedback on this chapter. We would like to acknowledge the financial support on literature purchase for this chapter by Optical Research Associates and our past and present employers for allowing us the use of their software: Optical Research Associates for LightTools, Photon Engineering for FRED, Lambda Research for TracePro and Breault Research Organization for ASAP.

40.9 REFERENCES

1. M. Johnson, D. G. Stork, S. Biswas, and Y. Furuichi, "Inferring Illumination Direction Estimated from Disparate Sources in Painting: An Investigation into Jan Vermeer's Girl with a Pearl Earring," *SPIE Symposium on Electronic Imaging*, San Jose, CA, 2008.
2. M. S. Rea, *The IESNA Lighting Handbook*, 9th edition, IESNA, New York, NY, 2000, Chapter 3.
3. IESNA 1988 Lighting for the Aged and Partially Sighted Committee, *Recommended Practice for Lighting and the Visual Environment for Senior Living*, RP-28-98, 1998.
4. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, New York, NY, 2003, p. 87.
5. M. D. Egan and V. Olgay, *Architectural Lighting*, 2nd edition, McGraw Hill, New York, 2002, p. 219.
6. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, 2003, p. 141.
7. M. S. Rea, *The IESNA Lighting Handbook*, 9th Edition, IESNA, New York, NY, 2000, pp. 10–13.
8. Commission Internationale de l'Éclairage (CIE), *International Lighting Vocabulary*, 4th edition, Publication 17.4, 1987.
9. Commission Internationale de l'Éclairage (CIE), *Colorimetry*, 3rd edition, Publication 15, 2004.
10. Commission Internationale de l'Éclairage (CIE), *Method of Measuring and Specifying Color Rendering Properties of Light Sources*, Publication 13.3, 1995.
11. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, New York, NY, 2003, p. 162.
12. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, New York, NY, 2003, p. 171.
13. Commission Internationale de l'Éclairage (CIE), *CIE Equations for Disability Glare*, Publication, 146, 2002.
14. Commission Internationale de l'Éclairage (CIE), 1995c, *Recommendations for the Lighting of Roads for Motor and Pedestrian Traffic*, CIE Technical Report 115, Vienna, 1995.
15. S. K. Guth, "A Method for the Evaluation of Discomfort Glare," *Illumination Engineering*, 57:351–64 (1963).
16. Chartered Institution of Building Services Engineers (CIBSE), *Technical Memorandum 10: The Calculation of Glare Indices*, CIBSE, London, 1985.
17. D. Fischer, "The European Glare Limiting Method," *Lighting Research and Technology*, 4:97–100, 1972.
18. G. Sollner, "Glare from Luminous Ceilings," *Lichttechnik*, 24:557–560, 1972.
19. G. Sollner, "Subjective Appraisal of Discomfort Glare in High Halls, Illuminated by Luminaires for High Intensity Discharge Lamps," *Lichttechnik*, 26:169–172, 1974.
20. Commission Internationale de l'Éclairage (CIE), 1995b, *Discomfort Glare in Interior Lighting*, Publication 117, 1995.
21. Commission Internationale de l'Éclairage (CIE), *Glare from Small, Large and Complex Sources*, Publication 147, 2002.
22. R. Levin, "Position Index in VCP Calculations," *Journal of the Illuminating Engineering Society*, pp. 99–105, January 1975. (Equation adapted from M. S. Rea, *The IESNA Lighting Handbook*, 9th Edition, IESNA, New York, NY, 2000, p. 9–26.)
23. P. R. Boyce, *Human Factors in Lighting*, 2nd edition, Taylor and Francis, New York, NY, 2003, Chapter 10.
24. R. B. Gibbons and C. J., Edwards, "A Review of Disability and Discomfort Glare Research and Future Direction," *18th Biennial TRB Visibility Symposium*, College Station, TX, 2007.
25. H. H. Bjorset and E. A. Frederiksen, "A Proposal for Recommendations for the Limitation of the Contrast Reduction in Office Lighting," *Proc. CIE 19th session*, Kyoto, Japan, 1979.
26. R. Whitehead, *Residential Lighting—A Practical Guide*, John Wiley and Sons, Newark, NJ, 2004, p. 5.
27. M. S. Rea, *The IESNA Lighting Handbook*, 9th edition, IESNA, New York, NY, 2000, pp. 9–17.
28. M.S. Rea, *The IESNA Lighting Handbook*, 9th edition, IESNA, New York, NY, 2000, pp. 9–21.
29. M. S. Rea, *The IESNA Lighting Handbook*, 9th edition, IESNA, New York, NY, 2000, pp. 9–29.
30. J. C. Stover, *Optical Scattering: Measurement and Analysis*, 2nd edition, SPIE, Bellingham, WA, 1995.
31. A. Gupta, *Simulation in Light Tools*, an illumination design software by Optical Research Associates, 2008.
32. Wikipedia, en.wikipedia.org/wiki/IGES (as of April 17, 2008).

33. U. S. Product Data Association, *Initial Graphics Exchange Specification IGES 5.3*, Charlestown, SC (1996), downloadable PDF version available at www.uspro.org/documents/IGES5-3 for Download. pdf (1997, as of April 17, 2008).
34. ISO TC 184/SC4, *STEP Application Handbook ISO 10303*, Version 3, Prepared by SCRA (North Charleston, SC, 2006); downloadable PDF version available at www.tc184_sc4.org/SC4_Open/SC4_Standards_Developers_Info/Files/STEP_application_handbook_63006.pdf (2006, as of April 17, 2008).
35. Wikipedia, en.wikipedia.org/wiki/ISO_10303 (as of April 17, 2008).
36. M. S. Kaminski, K. J. Garcia, M. A. Stevenson, M. Frate, and R. J. Koschel, "Advanced Topics in Source Modeling," *SPIE Proc. of Source Modeling I* **4775**, 46 (Society of Photo-Optical Instrumentation Engineers, Bellingham, WA, 2002).
37. M. A. Stevenson, M. S. Kaminski, M. Frate, and R. J. Koschel, "Modeling Filament-Based Sources for System Tolerancing," *SPIE Proc. of Modeling and Characterization of Light Sources* **4775**, 67, Bellingham, WA, 2000.
38. T. McReynolds and D. Blythe, *Advanced Graphics Programming Using OpenGL*, Elsevier Morgan Kaufmann, San Francisco, CA (2005).
39. F. X. Sillion and C. Puech, *Radiosity and Global Illumination*, Morgan Kaufmann, San Francisco, CA (1994).
40. R. J. Koschel, "Lit Appearance Modeling of Illumination Systems," *SPIE Proc. of Novel Optical Systems Design and Optimization V* **4768**:65 (Society of Photo-Optical Instrumentation Engineers, Bellingham, WA, 2002).
41. D. Scott and R. J. Koschel, Rendering in "ASAP," optical illumination software by Breault Research Organization.
42. R. J. Koschel, Rendering in "LucidShape" and "LucidDrive," vehicular illumination software by Brandenburg, GMBH.
43. M. Zollers, Photorealistic rendering in "LightTools," optical illumination software by Optical Research Associates.
44. D. Scott, Photorealistic rendering in "FRED," optical illumination software by Photon Engineering.
45. R. E. Smallwood, "Refractory Metals and Their Industrial Applications: Asymposium," ASTM committee B-10 on Reactive and Refractory Metals and Alloys, 1982.
46. S. Shionoya and W. M. Yen, *Phosphor Handbook*, CRC Press LLC, Boca Raton, FL, 1999.
47. A. M., Srivastava and C. R. Ronda, "Phosphors," *The Electrochemical Society Interface*, summer 2003.
48. N. Tesla, "Experiments with Alternate Currents of Very High Frequency and Their Application to Methods of Artificial Illumination," delivered before the American Institute of Electrical Engineers, Columbia College, NY, May 20, 1891.
49. W.J. Cassarly and T.L.R. Davenport, "Non-Rotationally Symmetric Mixing Rods," *Proceedings of SPIE*, 6342, July 2006.
50. W.J. Cassarly, "High-Brightness LEDs," *Optics and Photonics News*, 19–23, January 2008.
51. "LEDs Move into the Ultraviolet," www.physicsworld.com, May 17, 2006
52. W.J. Cassarly and T.L.R. Davenport, "Non-rotationally Symmetric Mixing Rods," *Proceedings of SPIE*, 6342, July 2006.
53. H. Chen, C. Hsu, H. Hong, "InGaN-CdSe-ZnSe Quantum Dots White LEDs," *Photonics Technology Letters*, IEEE **18**(1):193–195 (Jan. 1, 2006).
54. "Joint Venture to Make ZnSe White LEDs," <http://optics.org/cws/article/research/16534>, accessed on Oct. 13, 2008.
55. S. J. Smith, E. M. Purcell, "Visible Light from Localized Charges Moving across a Grating," *Physical Review* **92**(4):1069 (1953).
56. E. Yablanovich, "Inhibited Spontaneous Emission in Solid-State Physics and Electronics," *Physical Review Letters*, **55**:20, 2059, 1987.
57. J. J. Wierer, M. R. Krames, J. E. Epler, N. F. Gardner, J. R. Wendt, M.I. M. Sigalas, S. R. J. Brueck, D. Li, and M. Shagam, "III-nitride LEDs with Photonic Crystal Structures," *Proceedings of SPIE*, **5739**:102–107, 2005.
58. D. L. Barton and A. J. Fischer, "Photonic Crystals Improve LED Efficiency," *SPIE Newsroom*, 10.1117/2.1200603.0160, 2006.
59. C. W. Tang and S. A. VanSlyke, "Organic Electroluminescent Diodes," *Applied Physics Letters*, **51**:913, 1987.
60. K. Müllen and U. Scherf, *Organic Light Emitting Devices: Synthesis, Properties and Applications*, 1st edition, Wiley-VCH, Weinheim, March 2006.

61. E. Ne'eman, and R. G. Hopkinson, "Critical Minimum Acceptable Window Size: A Study of Window Design and Provision of a View," *Lighting Research and Technology*, **2:1**, 17–27, 1970.
62. E. C. Keighley, "Visual Requirements and Reduced Fenestration in Offices: A Study of Multiple Apertures and Window Area," *Building Science*, **8:4**, 321–331, 1973.
63. A. M. Ludlow, "The Functions of Windows in Buildings," *Lighting Research and Technology*, **8:2**, 57–68, 1976.
64. R. Winston, J. C. Miñano, and P. Benítez, *Nonimaging Optics*, Elsevier Academic Press, Burlington, MA, 2005.
65. J. Chaves, *Introduction to Nonimaging Optics*, CRC Press, Taylor and Francis Group, Boca Raton, FL, 2008.
66. *American National Standard Practice for Roadway Lighting*, The Standard Practice Subcommittee of the IESNA Roadway Lighting Committee, ANSI/IESNA RP-8-001, p. 6 IESNA, New York, NY, 1999, reaffirmed 2005.
67. *Luminaire Classification System for Outdoor Luminaires*, Luminaire Classification Task Group of IESNA, IESNA TM-15-07 (revised), IESNA, New York, NY, 2007.
68. W. J. Cassarly, D. Jenkins, A. Gupta, R. J. Koschel, "Hidden Devices That Light Our World Lightpipes," *Optics and Photonics News*, 34–39, August 2001.
69. J. K. Holton, "Daylighting of Buildings," *US National Bureau of Standards*, NBSIR 76-1098, October 1976.
70. W. M. C. Lam, *Sunlight as Formgiver of Architecture*, Von Nostrand Reinhold, New York, 1986.
71. M. S. Rea, *The IESNA Lighting Handbook*, 9th Edition, IESNA, New York, NY, 2000.
72. M. D. Egan and V. Olgyay, *Architectural Lighting*, 2nd edition, McGraw Hill, New York, NY, 2002.
73. Commission Internationale de l'Éclairage, *Methods of Characterizing Illuminance Meters and Luminance Meters: Performance, Characteristics and Specifications*, Publication 69, Bureau Central de la CIE, Vienna, 1987.
74. A. V. Arecchi, T. Messadi, and R. J. Koschel, *Field Guide to Illumination*, SPIE Press, Bellingham, WA (2007), pp. 90–92.
75. M. S. Kaminski and R. J. Koschel, "Methods of Tolerancing Injection-Molded Parts for Illumination Systems," *Proceedings of SPIE, Design of Efficient Illumination Systems*, 5186, 61, Bellingham, WA, 2003.
76. R. J. Koschel, "Illumination System Tolerancing," *Proceedings of SPIE, Optical System Alignment and Tolerancing*, 6676, 667604, Bellingham, WA, 2007.
77. *Recommended Practice for Lighting Merchandise Areas (A Store Lighting Guide)*, The IESNA Merchandise Lighting Committee, IESNA RP-2-01, IESNA, New York, NY, 2001. p. 2.
78. *Ibid*, p. 5.
79. *Lighting for Hospitals and Health Care Facilities*, The IESNA Committee for Health Care Facilities, ANSI/IESNA RP-29-06, IESNA, New York, NY, 2006, p. 4.
80. *Recommended Practice for Lighting Industrial Facilities*, Appendix A2, by the IESNA Industrial Lighting Committee, New York, NY, 2001.
81. *Lighting for Exterior Environments and IESNA Recommended Practice*, The IESNA Outdoor Environment Lighting Committee, IESNA RP-33-99, IESNA, New York, NY, 1998.
82. *Addressing Obtrusive Light (Urban Sky Glow and Light Trespass) in Conjunction with Roadway Lighting*, the Obtrusive Light Subcommittee of the IESNA Roadway Lighting Committee, IESNA TM-10-00, IESNA, New York, NY, 2000.
83. International Dark Sky Association, U.S. House of Representatives Briefing by Lee Cooper, June 20, 2008.
84. *Light Trespass: Research Results and Recommendations*, The Obtrusive Light Subcommittee of the IESNA Roadway Lighting Committee, IESNA TM-11-00, IESNA, New York, NY, 2000.
85. *A Discussion of Appendix E—"Classification of Luminance Light Distributions"*, "Roadway Standard Practice Subcommittee of the IESNA Roadway Lighting Committee, IESNA TM-3-95, IESNA, New York, NY, 1995.
86. *Code of Federal Regulations*, Title 49 Transportation, Volume 6, Chapter 5, Part 571, Federal Motor Vehicle Safety Standards, USA Federal Government (Washington, D.C., 2007), pp. 279–353; online at http://edocket.access.gpo.gov/cfr_2007/octqtr/pdf/49cfr571.108.pdf. Accessed 13 October 2008.
87. *SAE Ground Vehicle Lighting Standards Manuals*, HS-34, 2001 edition, SAE International, Warrendale, PA, 2001, p. 8.
88. *UNECE Standards*, ECE Transaction 505, Revision 2, Addendum 112, Regulation No. R1 13, online at <http://www.unece.org/trans/main/wp29/wp29regs/r113e.pdf>, 2001. Accessed 13 October 2008.

89. *SAE Ground Vehicle Lighting Standards Manuals*, HS-34, 2001 edition, SAE International, Warrendale, PA, 2001, p. 9.
90. *UNECE Standards*, ECE Transaction 505, Revision 2, Addendum 111, Regulation No. R112, online at <http://www.unece.org/trans/main/wp29/wp29regs/r112e.pdf>, 2001. Accessed 13 October 2008.
91. *SAE Ground Vehicle Lighting Standards Manuals*, HS-34, 2001 edition, SAE International, Warrendale, PA, 2001, p. 134.
92. *UNECE Standards*, ECE Transaction 505, Revision 4, Addendum 6, Regulation No. R7, online at <http://www.unece.org/trans/main/wp29/wp29regs/r007r4e.pdf>, 2006. Accessed 13 October 2008.
93. *Design Guide for Roadway Lighting Maintenance*, The Subcommittee on Lighting Maintenance and Light Sources of the IESNA Roadway Lighting Committee, IESNA DG-4-03, IESNA, New York, NY, 2003.
94. *American National Standard Practice for Roadway Lighting*, The Standard Practice Subcommittee of the IESNA Roadway Lighting Committee, ANSI/IESNA RP-8-001, IESNA (New York, NY, 1999, reaffirmed 2005), p. 2.
95. *IESNA Recommended Practice for Roadway Sign Lighting*, The Sign Lighting Subcommittee of the IESNA Roadway Lighting Committee, IESNA RP-19-01, IESNA, New York, NY, 2001.
96. *Manual on Uniform Traffic Control Devices*, Sections 2A-11 to 2A-15, Federal Highway Administration, U.S. Department of Transportation, Washington, D.C., 1988.
97. *American National Standard Practice for Tunnel Lighting*, The Tunnel Lighting Subcommittee of the IESNA Roadway Lighting Committee, IESNA RP-22-05, IESNA, New York, NY, 2005.

This page intentionally left blank.

DO NOT DUPLICATE

INDEX

Index note: The *f* after a page number refers to a figure, the *n* to a note, and the *t* to a table.

- Abbe illuminated eyepieces, **12.12**, **12.12f**
Abbe illumination system, **39.23**, **39.23f**,
39.34, 39.35f
Abbe's sine condition, **34.19**
Aberration coefficients, **3.10–3.11**
Aberration curves (in lens design), **2.1–2.6**
 considerations for, **2.5–2.6**
 field plots of, **2.4–2.5**
 transverse ray plots of, **2.2–2.4**
Aberrations:
 balancing of, **11.30**, **11.35–11.36**, **11.36t**
 evaluation of, **3.9–3.11**
 (See also *specific aberrations*, e.g.: Axial chromatic aberration)
Absolute detectors, **34.27–34.30**
 electrical substitution radiometers, **34.27–34.29**
 photoionization devices, **34.29**
 predictable quantum efficiency devices,
 34.29–34.30
Absolute measurements, **34.20–34.37**
 accuracy and traceability of, **34.21**
 error propagation in, **34.22**
 error types in, **34.21–34.23**
 relative vs., **34.20–34.21**
 and uncertainty estimates, **34.21–34.23**
Absolute responsivity units (A/W), **34.31**
Absolute sources (of radiation), **34.23–34.27**
 blackbody radiator, **34.23–34.24**
 blackbody simulators, **34.24–34.26**
 synchrotron radiation, **34.26–34.27**
Absorbing media:
 in photodetectors, **26.4f**, **26.5**
 radiant power transfer through, **34.13**
Absorbing substrate (AS) chips, **17.7**, **17.7t**
Absorptance:
 defined, **35.4**
 measurement of, **35.10**
 in thermal detectors, **28.2**
 and transmittance/reflectance, **35.7**, **35.8**, **35.8t**
Absorption:
 defined, **35.4**
 quantum resonance, **22.16**, **22.17**
 stimulated, **16.7–16.8**, **16.8f**
Absorption coefficient, **32.2–32.4**, **32.3f**
 of *pin* photodiodes, **25.8**, **25.9f**
 of visible array detectors, **32.2–32.3**, **32.3f**
Absorption rate, **23.8**
Ac lamps, **15.32f**
Accent lighting, **40.14**, **40.14f**
Accuracy:
 of absolute measurements, **34.21**
 of CGHs, **14.6–14.7**, **14.6f**, **14.7f**
 as measure of systemic errors, **12.2**
Achromatic lenses, athermalized, **1.16**, **1.16f**
Achromatism, **1.14–1.15**, **1.15f**
Actinic effects (of radiation), **34.6**, **34.7**
Actinic ultraviolet action spectrum, **36.17**
Actinometry:
 conversions between radiometry/photometry
 and, **34.12**, **34.12t**
 defined, **34.7**, **34.11**
Activated-phosphor sources (of radiation), **15.49**
Active athermalization, **6.24**, **6.24f**
Active imaging, **31.29–31.30**
Active mechanical athermalization, **8.11**, **8.11f**
Acutance:
 defined, **30.2**
 of photographic systems, **29.17–29.19**,
 29.18t, **29.19f**
Adaptation, in vision, **40.9**
Additive damping, **3.18**
Adiabatic approximation, **23.21** (See also
 Markov approximation)
Advanced Photo System (APS), **30.21**,
 30.26, **30.27t**
Afocal systems:
 as attachments, **1.8**, **1.9f**
 first-order layout for, **1.7**, **1.7f**

- Agfachrome, 29.14
 AHU pelloid, 30.4
 Air-spaced doublet lens, 6.7
 Air-spaced triplet lens, 6.21, 6.22*f*
 Airway beacon lamps, 15.11
 Allan Deviation, 22.2–22.4
 Allan Variance method, 22.2, 22.3
 Alloy disordering, 19.24
 Alphanumeric displays, LED, 17.31–17.32
 Aluminized phosphor-screen/window assembly, 31.14, 31.15*f*
 Aluminum, diamond turning and, 10.4
 Aluminum gallium arsenide (AlGaAs) emitters, 17.32
 Aluminum gallium arsenide (AlGaAs) LEDs, 17.17, 17.17*t*, 17.28*f*
 Aluminum gallium arsenide (AlGaAs) quantum well photodetectors, 25.16–25.17, 25.16*f*, 25.17*f*
 Aluminum gallium arsenide (AlGaAs) substrate, 17.22
 Aluminum gallium nitride (AlGaN) alloy photovoltaic detectors, 24.46
 Aluminum gallium nitride (AlGaN) substrate, 17.22
 Aluminum indium gallium nitride (AlInGaN) material systems, 18.1–18.2, 18.2*f*, 18.4
 Aluminum indium gallium phosphide (AlInGaP) LEDs, 17.18, 17.19*f*
 Aluminum indium gallium phosphide (AlInGaP) material systems, 18.1, 18.5
 Aluminum indium gallium phosphide (AlInGaP) substrate, 17.22
 Aluminum mirrors, mounting of, 6.20*f*
 Ambient lighting, 40.12, 40.13*f*, 40.15*f*
 Ambient temperature electrical substitution radiometers, 34.27
 American Institute of Physics (AIP), 36.3
 American National Standards Institute (ANSI), 4.11, 36.2
 American Society for Testing and Materials (ASTM), 37.11
 AMI (amplified MOS imager) MOS readout, 32.21
 Amorphous silicon photoconductors, 32.4*f*, 32.31, 32.32
 Amplification, 16.9
 Amplifier strategies, for PZTs, 22.18
 Amplifiers, 27.2
 properties of, 16.3
 selection of, 27.10–27.12
 transconductance, 27.11–27.12, 27.11*f*
 voltage, 27.10–27.11
 Amplitude gating, 21.7
 Amplitude modulation (AM), 19.36
 Amplitude response, frequency vs., 22.6–22.7
 Analytical density, 29.14
 -ance (suffix), 35.3
 Angle measurement, 12.10–12.17
 autocollimeters for, 12.11–12.12, 12.11*f*, 12.12*f*
 interferometric methods of, 12.14
 levels (tools) for, 12.13–12.14, 12.13*f*, 12.14*f*
 mechanical methods of, 12.10–12.11, 12.11*f*
 in prisms, 12.14–12.16, 12.15*f*–12.17*f*
 theodolites for, 12.13
 Angle solves, 3.6
 Angular dilution, 39.6
 Angular uniformity, 39.31
 Annular flanges, 6.3–6.4, 6.4*f*
 Annular polynomials:
 for defocus, 11.38*f*, 11.39
 Zernike, 11.13–11.21, 11.14*f*, 11.16*f*, 11.17*t*–11.21*t*
 Annunciator assemblies, 17.30
 Anodes, in photomultipliers, 27.6, 27.7, 27.7*f*
 Anomalous reflection colors, 30.17
 Antiblooming, 32.9, 32.10*f*
 Antihalation undercoat (AHU) layers, 30.4
 AOM transducers, 22.20
 APART (stray light analysis program), 7.11
 Aperture(s):
 data about, 3.4
 nonideal, 34.35–34.36, 34.35*f*
 numerical, 34.20, 39.1
 in optical design software, 3.6
 Aperture flash mode, 39.7
 Aperture placement (in stray light suppression), 7.5–7.10
 aperture stops, 7.6–7.7, 7.7*f*, 7.8*f*
 field stops, 7.7, 7.8*f*, 7.9*f*, 34.18–34.19, 34.19*f*
 Lyot stops, 7.8–7.10, 7.8*f*–7.11*f*
 Aperture stops, 1.4, 3.4, 7.6–7.7, 7.7*f*, 7.8*f*, 34.18, 34.19*f*
 Aphakic hazard, 36.17
 Aplanatic optical systems, 34.19–34.20
 Aplanats, 39.8, 39.9*f*
 Apodization, 39.7

- Apostilb (unit), 34.43, 36.7, 36.8*t*
 Apovortex surfaces, 39.11
 Approximate transfer factor (ATF), 4.1, 4.4
 Arc radiation sources (*see specific arcs, e.g.:*
 Argon arcs)
 Area image sensor arrays, 32.24–32.32
 about, 32.2
 CCD
 frame transfer, 32.26–32.28, 32.27*f*, 32.28*f*
 interline transfer, 32.28–32.32,
 32.29*f*–32.31*f*
 performance of, 32.32
 image area dimensions for, 32.25*t*
 MOS, 32.25–32.26, 32.26*f*
 Area-solid angle product, 39.5
 Argon arcs, 15.12, 15.13
 Argon ion lasers, 16.14, 16.15*f*, 16.30
 Array-mode selection, semiconductor, 19.27
 Array-mode stability, semiconductor, 19.27
 Artificial sources (of radiation), 15.3–15.53
 about, 15.3–15.4
 commercial, 15.13–15.53
 activated-phosphor sources, 15.49
 blackbody simulators, 15.14, 15.15*f*, 15.16*f*
 carbon arc sources, 15.21–15.24, 15.23*f*,
 15.24*f*, 15.25*t*–15.27*t*, 15.28*f*
 concentrated arc lamps, 15.47–15.48,
 15.48*f*, 15.49*f*
 glow modulator tubes, 15.49, 15.50*f*,
 15.51*f*, 15.52*t*
 high-energy sources, 15.40
 high-pressure enclosed arc, 15.24,
 15.28–15.34, 15.29*f*–15.35*f*
 hydrogen and deuterium arc lamps,
 15.49, 15.53*f*
 incandescent nongaseous sources,
 15.15–15.21, 15.17*f*–15.22*f*
 low-pressure enclosed arc, 15.35–15.47,
 15.36*f*, 15.36*t*–15.43*t*, 15.44*f*–15.47*f*,
 15.46*t*, 15.47*t*
 special-purpose sources, 15.53
 luminaire optics for, 40.45–40.47, 40.45*f*,
 40.46*f*
 and radiation law, 15.4–15.7, 15.5*f*, 15.5*t*, 15.6*f*
 standardized laboratory sources, 15.7–15.13
 baseline standard of radiation, 15.9, 15.9*f*,
 15.10*f*, 15.12*f*
 blackbody cavity theory, 15.7–15.9, 15.8*f*
 working standards of radiation,
 15.9–15.13, 15.10*f*, 15.12*f*, 15.13*f*
 ASAP (optical software), 7.25
 Aspheric lenses, 39.8, 39.9, 39.9*f*
 Aspheric measuring system, 10.12*f*
 Aspheric surfaces, 3.5
 Aspherical optics fabrication, 9.7–9.8, 9.7*f*
 Aspherical wavefront measurement,
 13.23–13.27
 holographic compensators, 13.25, 13.25*f*,
 13.26*f*
 infrared interferometry, 13.25
 Moiré tests, 13.26–13.27
 refractive or reflective compensators, 13.24,
 13.24*f*, 13.25
 sub-Nyquist interferometry, 13.27
 two-wavelength interferometry, 13.25, 13.26
 wavefront stitching, 13.27, 13.27*f*
 Assembly tolerances, 5.8
 Astigmatism, 2.3, 2.3*f*
 Astronomical telescopes, 1.7*f*
 Athermal laser beam expanders, 8.13–8.14
 Athermalization, 1.15–1.16, 1.16*f*, 6.22–6.24
 active, 6.24, 6.24*f*
 intrinsic, 8.7–8.8, 8.7*f*
 mechanical, 8.8–8.12
 active, 8.11, 8.11*f*
 by image processing, 8.12
 part active, part passive, 8.11–8.12, 8.12*f*
 passive, 8.8–8.10, 8.8*f*–8.10*f*
 optical, 8.12–8.15, 8.13*t*
 about, 8.12–8.13
 athermal laser beam expanders, 8.13–8.14
 diffractive optics usage, 8.15
 of separated components, 8.14, 8.15
 three-material solutions, 8.14, 8.14*t*, 8.15*t*
 passive, 6.22, 6.23*f*, 6.24
 single material design, 6.22, 6.23*f*
 Athermalized achromatic lenses, 1.16, 1.16*f*
 Atomic (gain) noise, 23.34–23.35
 Attosecond optics, 21.1–21.9
 about, 21.2
 driving lasers in, 21.4–21.6
 carrier-envelope offset frequency, 21.5
 carrier-envelope phase, 21.5*f*, 21.6
 carrier-envelope phasemeter, 21.6
 chirped pulse amplification, 21.5
 chirped pulse amplifiers, 21.6
 single-shot *f*-to-2*f* interferometer, 21.6
 high-harmonic generation, 21.2, 21.2*f*
 phase-matching in, 21.4
 ponderomotive potential in, 21.3

- Attosecond optics (*Cont.*):
- pulse characterization, 21.8–21.9, 21.8f
 - attosecond streak camera, 21.9
 - FROG-CRAB, 21.9
 - RABITT, 21.9
 - second-order autocorrelator, 21.9
 - pulse generation, 21.6–21.8, 21.7f
 - amplitude gating, 21.7
 - attosecond pulse train, 21.6, 21.7
 - double optical gating, 21.8
 - polarization gating, 21.7–21.8
 - two-color gating, 21.7
 - quantum trajectories in, 21.3–21.4
 - semiclassical model of, 21.3
 - single isolated pulses in, 21.4
 - strong field approximation in, 21.3
- Attosecond pulse, 21.8–21.9, 21.8f
- attosecond streak camera, 21.9
 - FROG, 21.9
 - FROG-CRAB, 21.9
 - generation of, 21.6–21.8, 21.7f
 - amplitude gating, 21.7
 - attosecond pulse train, 21.6, 21.7
 - double optical gating, 21.8
 - polarization gating, 21.7–21.8
 - two-color gating, 21.7
 - RABITT, 21.9
- Attosecond pulse train, 21.6, 21.7
- Attosecond streak camera, 21.9
- Augur recombination, 19.17, 19.17f
- Autocollimeters:
- angle measurement with, 12.11–12.12, 12.11f, 12.12f
 - curvature measurement with, 12.19–12.20, 12.20f
 - defined, 12.11–12.12
- Automatic brightness control (ABC), 31.18
- Automatic spherometers, 12.19
- Autoset levels (tools), 12.14, 12.14f
- Avalanche multiplication, 25.9
- Avalanche photodetectors (APDs):
- high-speed, 26.17–26.20, 26.18f, 26.20f, 26.21f
 - improvements in, 26.3
- Avalanche photodiodes, 24.62–24.70, 24.63f–24.70f, 24.72–24.73, 24.72f, 24.73f, 25.8–25.10, 25.9f
- defined, 24.10
 - germanium, 24.70f, 24.72–24.73, 24.72f, 24.73f
 - InGaAs, 24.66–24.70, 24.66f–24.69f
 - silicon, 24.62–24.65, 24.63f–24.66f
- Avogadro's number, 34.11
- Axial chromatic aberration, 1.14, 2.2, 2.3f
- Axial gap prevention, 6.21, 6.22f
- Axial rays, 1.4, 1.11f, 1.12
- Azimuth angle, 35.5
- Azomethine dyes, 30.10f
- excited state properties of, 30.11–30.12, 30.12f
 - formation of, 30.10
 - photochemistry of, 30.11
- Back light, 40.43, 40.44f
- Background temperature, 24.10
- Background-limited performance (BLIP), of infrared detector arrays, 33.24
- Backing, film, 29.4
- Backlighting, 40.1, 40.12, 40.47, 40.47f, 40.48f
- Backscattering, 20.13–20.15, 20.15f
- Backward trace, 39.7
- Baffles:
- cone-shaped secondary, 7.3–7.4, 7.3f, 7.4f
 - with integrating cavities, 39.26
 - in lighting design, 40.41, 40.45f, 40.46
 - shields for, 7.9–7.10, 7.9f, 7.10f
 - in stray light suppression, 7.10, 7.11
 - two-stage, 7.10
- Balanced spherical aberrations, 11.30
- Ballasts:
- in fluorescent lamps, 40.32–40.33
 - in HID lamps, 40.36
- Band pass, 38.8
- Bandwidth:
- of amplifiers, 27.10
 - gain-bandwidth, 26.17
 - normalization of, 36.14–36.16, 36.15t, 36.16f
 - of photomultipliers, 27.7
- “Bang-bang” zoom, 1.12
- Banker lamps, 40.12, 40.46, 40.46f
- Bar spherometers, 12.19, 12.19f
- Bar-code reading, 17.34
- Bare source light, 40.43
- Barium strontium titanate (BST), 28.11, 28.12
- Baryta layer, 30.5
- Batwing lenses, 40.12
- Baud rate, 17.33
- Beacon lamps, airway, 15.11
- Beam splitters, 13.7, 34.32
- Beam transformers, 39.18, 39.18f
- Beam-forming illumination systems, 39.22, 39.39

- Beam-smearing faceted reflectors, 39.39
- Beer-Lambert law (Beer's law), 16.9, 34.35, 38.5
- Bell-clamping (edging fabrication step), 9.6
- Bellcore, 19.39, 19.41
- Bevel gauges, 12.10, 12.11, 12.11f
- Bevel placement (on vanes), 7.13, 7.14f
- Bias angle, 31.13
- Biased *pin* photodetectors, 26.6f
- Biconical reflectance, 35.5t, 35.6f, 35.6t
- Bidirectional reflectance, 35.5t, 35.6f, 35.6t
- Bidirectional reflectance distribution function (BRDF), 7.1, 7.18–7.19, 7.22, 35.5, 35.13, 37.9
- Bidirectional scattering distribution function (BSDF), 7.2, 7.23, 7.24f, 7.25f, 35.13
- Bidirectional transmittance distribution function (BTDF), 35.3, 35.13
- Bihemispherical reflectance, 35.5t, 35.6f, 35.6t
- Binning, 38.10
- Bipolar transistors, 27.11
- Black-and-white (B&W) film, 29.4, 30.24–30.25, 30.25t
- Blackbody cavity theory, 15.7–15.9
- Blackbody D star, 24.10
- Blackbody detectivity, 24.10
- Blackbody noise-equivalent power, 24.10
- Blackbody radiation, 15.4–15.6, 15.5t
emittance of, 34.25–34.26
sources of, 15.14, 15.15f, 15.16f, 34.23–34.24
temperature vs., 36.12, 36.12f, 36.14, 36.14f
working standards for, 15.14, 15.16f
- Blackbody responsivity, 24.10
- Blackbody simulators, 15.14, 15.15f, 15.16f, 34.24–34.26
- Black-light fluorescent lamps, 15.35, 15.36t
- Bleaching, in film development, 29.14
- Blindness, flash, 40.9
- Blip detector (blip condition), 24.10
- Blocked impurity band (BIB), 33.7
- Blocking contacts, 26.3
- Blocking filters, 38.8
- Blocking filters, 38.8
- Blood gas analysis, 17.34
- Blooming:
antiblooming, 32.9, 32.10f
in image sensors, 32.6, 32.9
- Blue emitters, in LED technology, 17.18, 17.19
- Blue light, color film and, 29.13, 29.13f, 30.3–30.4
- Blue semiconductor lasers, 19.7
- Blue-enhanced photodiodes, 24.55f, 24.61–24.62, 24.61f, 24.62f
- Blur filters, 32.34, 32.34f
- “Boat grown” technique, 17.21
- Bode representation, of servo system, 22.5–22.6, 22.6f
- Bolometers, 24.5, 28.3–28.5, 28.4f
about, 28.1
carbon, 28.5
detectivity of perfect, 24.17, 24.18f
germanium low-temperature, 24.31–24.32, 24.32f, 24.33f
indium antimonide hot-electron, 24.29, 24.30, 24.30f, 24.31f
as infrared detectors, 33.9–33.10
metal, 28.4, 28.7t
properties of, 28.7t
resistive arrays of, 28.10–28.11, 28.10f
semiconductor, 28.4–28.5
superconducting, 28.5
thermistor, 24.24–24.25, 24.24f, 24.25f, 28.7t
- Bonded mountings, 6.13–6.15, 6.15f, 6.16f
- Boresight tolerances, 5.8
- Boron-doped silicon (Si:B) detectors, 24.95f, 24.96
- Bose-Einstein condensation (BEC), 23.39
- Bose-Einstein statistics, 23.9
- Bouillotte lamps, 40.12, 40.46, 40.46f
- Boules, glass, 9.3
- Boundary conditions (of optics):
defined, 3.17
methods for handling, 3.18–3.19
specification of, 4.12
- Boxcar averaging, of modulated signal sources, 27.13, 27.13f, 27.15
- Bragg reflectors, 19.41
- Bridgeman technique, 17.21
- Brightness:
of carbon arcs, 15.23f
luminance vs., 34.40
perception of, 40.4, 40.4f
of scene, 31.4–31.5
- British Glare Index (CIBSE), 40.10
- Broad bandwidth solid-state lasers, 16.34–16.35, 16.34f
- Bromine, in light bulbs, 40.30
- Brunning distance-measuring interferometers, 12.8–12.9, 12.8f
- Buffered direct injection (BDI), 33.19t, 33.20f, 33.21–33.22

- Build-and-test evaluation (for stray light suppression), 7.28
- Built-in potential, 25.6
- Bulb blackening, 40.30
- Bulb shield, 40.45*f*, 40.46
- Bulk material photodetectors, 26.4*f*, 26.5
- Bulk-grown materials, 17.8
- Buried crescent lasers, 19.24, 19.25*f*
- Buried heterostructure (BH) lasers, 19.8, 19.9*f*, 19.20*t*, 19.24, 19.36*f*
- Buried TRS (BTRS) lasers, 19.19, 19.20*t*, 19.21*f*
- Buried V-groove-substrate inner stripe (BVSIS), 19.19, 19.20*t*
- Buried-channel CCDs (BCCDs), 32.14, 33.13
- Buried-channel MOS capacitors, 32.4*f*, 32.7–32.8
- Burnished mounting, 6.3, 6.3*f*
- Cable TV (CATV), 25.11–25.12
- Cadmium selenide (CdSe) photoconductors, 24.49–24.52, 24.52*f*
- Cadmium sulfide (CdS) photoconductors, 24.49–24.52, 24.51*f*–24.53*f*
- Cadmium telluride (CdTe) detectors, 24.52, 24.54, 24.54*f*
- Cadmium zinc telluride (CdZnTe) detectors, 24.52
- Calibration:
 - artificial sources of radiation for (*see* Artificial sources (of radiation))
 - legal traceability of, 34.21
 - photometric, 34.42–34.43
 - radiometric, 34.31–34.32
 - self-calibration, 34.29
 - spectroradiometric, 38.11–38.13, 38.11*t*, 38.12*f*
- Calibration transfer devices, 34.31–34.32
- Callier coefficient, 29.7
- Candela (unit), 34.37, 34.39, 36.4, 36.5, 37.3, 37.4
- Candle power, 37.3
- Capacitive bolometers, 33.10
- Capacitive transimpedance amplifier (CTIA), 33.19*t*, 33.20*f*, 33.22–33.23
- Capacitors, MOS, 32.4*f*, 32.7–32.8
- Capillary mercury-arc lamps, 15.30–15.31, 15.31*f*
- Carbon arc light sources, 40.40
- Carbon arc sources (of radiation), 15.21–15.24, 15.23*f*, 15.24*f*, 15.25*t*–15.27*t*, 15.28*f*
- Carbon bolometers, 28.5, 28.7*t*
- Carbon-dioxide lasers, 16.16, 16.16*f*, 16.30
- Carey Lea silver (CLS), 29.13, 30.4
- Carrier confinement, 17.12, 17.12*f*–17.14*f*, 17.13, 17.17
- Carrier density, 19.30–19.33
- Carrier transit time, 26.6–26.7, 26.6*f*
- Carrier trapping, 26.9, 26.9*f*
- Carrier-envelope offset, 20.4
- Carrier-envelope offset frequency, 21.5
- Carrier-envelope (CE) phase:
 - of chirped pulse amplifiers, 21.6
 - of lasers, 21.5, 21.5*f*
- Carrier-envelope phasemeters, 21.6
- Cassegrain design, 7.3*f*, 7.11, 7.14, 7.14*f*, 7.16, 7.16*f*, 7.19, 7.20*f*
- Cathodes:
 - photo-, 27.6, 27.7*f*
 - shielding of, 27.10
- Cavity(-ies):
 - distributed feedback lasers, 16.29
 - integrating (*see* Integrating cavities, of nonimaging optics)
 - mode-locking, 16.27–16.29, 16.28*f*
 - modifying output distribution of, 39.27
 - properties of, 16.3
 - Q-switching, 16.26–16.27, 16.27*f*
 - ring lasers, 16.29
 - stability of, 16.23–16.25, 16.24*f*, 16.25*f*
 - unstable resonators, 16.25–16.26, 16.26*f*
- Cavity dumping, 16.27
- Cavity losses, 23.18
- Cavity-shaped radiometers, 34.28
- Ceilings, illuminated, 40.13*f*
- Cellulose acetate film, 29.4
- Cenco Company, 15.47
- Center for Optics Manufacturing, 9.4
- Center-of-mass motion of atoms, 23.45
- Centrally obscured system (*see* Cassegrain design)
- Centration, of spherical lenses, 9.8
- Channel stop region, 32.7
- Channeled substrate planar (CSP) lasers, 19.20*t*, 19.21*f*, 19.22, 19.36*f*
- Charge integration matrix (CIM), 33.10–33.11, 33.11*f*, 33.12*f*
- Charge pumping, 32.9*n*
- Charge sweep devices (CSDs), 33.12*f*, 33.13

- Charge-coupled detector area image sensor arrays:
 frame transfer, 32.26–32.28, 32.27*f*, 32.28*f*
 interline transfer, 32.28–32.32, 32.29*f*–32.31*f*
 performance of, 32.32
- Charge-coupled detectors (CCDs), 25.10, 25.11, 25.11*f*, 31.1, 32.12–32.20, 38.9–38.10, 38.10*t*
 characteristics of, 32.17–32.20, 32.17*f*, 32.18*f*
 electronics of, 38.10
 image sensing with, 32.8
 linear arrays of, 32.21–32.24, 32.22*f*, 32.23*f*
 MIS photogate FPAs for, 33.10–33.11, 33.11*f*, 33.12*f*
 multilinear arrays of, 32.21, 32.23*f*, 32.24
 operation of, 32.12–32.14, 32.13*f*
 output of, 32.14–32.15, 32.15*f*
 performance of, 32.32
 readout from, 32.12–32.21, 32.13*f*
 types of, 32.15–32.17, 32.16*f*
- Charge-injection devices (CIDs), 31.1, 32.20, 33.10–33.11, 33.11*f*, 33.12*f*
- Chartered Institution of Building Services Engineers (CIBSE), 40.2
- Chemical beam epitaxy (CBE), 19.7
- Chemical-assisted ion beam etching (CAIBE), 19.39
- Chirped pulse amplification (CPA), 21.5, 21.5*f*, 21.6
- Chopper-stabilized amplifiers, 27.11
- Chopper-stabilized BDI, 33.19*t*, 33.20*f*, 33.21–33.22
- Chromatic aberrations, 2.2–2.4, 2.3*f*, 2.4*f*
- Chromium lasers, 16.34, 16.35
- Chromogenic film, 29.14
- Circle polynomials, 11.36*t*, 11.39
 isometric plots/interferograms/PSFs for defocus, 11.38*f*
 and noncircular pupils, 11.37, 11.39
 radial, 11.7, 11.9*f*–11.10*f*
 Zernike, 11.4, 11.6–11.12, 11.8*t*–11.9*t*, 11.9*f*–11.11*f*, 11.12*t*
- Circular discs, projected area of, 36.3*t*
- Cladding layers, 19.4
- Clarity, perception of visual, 40.5
- Clip test (of photographic film), 30.23
- Clipped Lambertian distribution, 39.3–39.4
- Clock generation, 33.16
- Closed-loop performance (in servo systems), 22.8
- Closed-loop stability issues (in servo systems), 22.8–22.12, 22.9*f*
 PID controller vs. notch filters, 22.10–22.11, 22.10*f*, 22.11*f*
 rule-of-thumb PID design for system with transducer resonance, 22.11–22.12
- Coarse-grained derivative, 23.21
- Coatings:
 lens specifications for, 4.10
 of photographic film, 29.4
- Coblentz-type thermopiles, 24.23
- Coherent states, 23.12
- Coiling, of light bulb filament, 40.30
- Cold cathode fluorescent lamps (CCFLs), 40.32
- Collares-Pereira, M., 39.17
- Collector power (in stray light suppression), 7.2
- Collectors (*see* Concentrators, nonimaging)
- Collimators:
 autocollimators, 12.12
 conic, 39.8, 39.9*f*
- Collisional broadening, emission-line, 16.5
- Colloidal silver, 29.13
- Color(s):
 anomalous reflection, 30.17
 in LEDs, 40.37
 and lighting design, 40.7–40.9
 mixing of, 40.8
 science of, 30.15–30.18, 30.16*f*, 30.17*f*
- Color aliasing, 32.34
- Color density, 29.7–29.8
- Color filter arrays (CFAs), 32.32–32.34, 32.33*f*, 32.34*f*
- Color imaging architectures, 32.32–32.34
 integral filter arrays, 32.32–32.34, 32.33*f*, 32.34*f*
 sequential, 32.32, 32.33*f*
 three-chip, 32.32, 32.33*f*
- Color negative films, 30.25–30.28, 30.27*t*
- Color photographic films:
 about, 30.2
 coating of, 29.4
 negative, 30.25–30.28, 30.27*t*
 reversal, 30.22–30.24, 30.23*t*
 structure of, 29.12–29.15, 29.13*f*, 29.14*f*, 30.3–30.5, 30.3*f*
- Color photographic paper, 30.5
- Color records, 30.4
- Color rendering, 40.8–40.9
- Color rendering index (CRI), 40.8

- Color reversal films, **30.2**, **30.22–30.24**, **30.23t**
 Color sequential systems, **32.32**, **32.33f**
 Color slide films (*see* Color reversal films)
 Color space calculations, **38.4–38.5**
 Color temperature, **34.44**, **37.4t**, **37.6–37.7**,
38.5, **40.8**
 Color transparency films (*see* Color reversal
 films)
 Color-center lasers, **16.35**
 Colorimetry, **37.11**
 Coma, with spherical aberration, **2.4**, **2.4f**
 Combined recombination, **17.3**
 Combined servo transducers, **22.19**
 Commercial sources (of radiation), **15.13–15.53**
 activated phosphor, **15.49**
 blackbody simulators, **15.14**, **15.15f**
 carbon arcs, **15.21–15.24**, **15.23f**, **15.24f**,
15.25t–15.27t, **15.28f**
 concentrated arcs, **15.47–15.49**, **15.48f**, **15.49f**
 glow modulator tubes, **15.49**, **15.50f**, **15.51f**,
15.52t
 high-energy, **15.40**
 high-pressure enclosed arc, **15.24**, **15.28–15.34**
 compact-source arcs, **15.31–15.34**,
15.32f–15.35f
 Lucalox lamps, **15.30**, **15.31f**
 mercury arcs, **15.29–15.31**, **15.30f**, **15.31f**
 multivapor arcs, **15.29**, **15.31f**
 Uviarc, **15.28–15.29**, **15.29f**, **15.30f**
 hydrogen and deuterium arcs, **15.49**, **15.53f**
 incandescent nongaseous, **15.15–15.21**
 comparisons, **15.19**, **15.19f**
 gas mantle, **15.17**, **15.18**, **15.19f**
 globar, **15.17**, **15.18f**
 Nernst glower, **15.14**, **15.15**, **15.17**, **15.17f**
 quartz-envelope lamps, **15.20**, **15.21**
 tungsten-filament lamps, **15.19**, **15.20**, **15.**
20f–15.22f
 low-pressure enclosed arc, **15.35–15.47**
 black-light fluorescent lamps, **15.35**, **15.36t**
 electrodeless discharge lamps, **15.36**, **15.44**
 germicidal lamps, **15.35**
 hollow cathode lamps, **15.35**,
15.37t–15.43t, **15.44f**
 Pluecker spectrum tubes, **15.47**, **15.47f**,
15.47t
 spectral lamps, **15.44**, **15.45**, **15.45f**,
15.46f, **15.46t**
 Sterilamps, **15.35**, **15.36f**
 special-purpose, **15.53**
- Commission Internationale de l'Eclairage
 (CIE), **40.2**
 publications from, **37.11**
 standard photometric observer, **37.2**
 Common path interferometers, **13.9**, **13.11f**
 Compact fluorescent lights (CFLs), **40.25t**,
40.26t, **40.28f**, **40.31**
 Compact-source arcs, **15.31–15.34**,
15.32f–15.35f
 Compensators:
 Dall, **13.24**, **13.24f**
 holographic, **13.25**
 Offner, **13.24**, **13.24f**
 reflective, **13.24**, **13.24f**, **13.25**
 refractive, **13.24**, **13.24f**, **13.25**
 and tolerances, **3.21**, **5.7**
 Complete monolithic FPAs, **33.10**
 Compound elliptical collectors (CECs), **39.14**,
39.15f, **39.27**, **39.37**
 Compound hyperbolic collectors (CHCs),
39.15, **39.15f**, **39.16f**, **39.37**
 Compound lens, thermal defocus of, **8.4**, **8.5f**
 Compound parabolic collectors (CPCs),
39.13–39.14, **39.13f**, **39.14f**, **39.18**, **39.19**,
39.19f
 Compressively strained QW lasers, **19.16f**, **19.17**
 Compton effect, **23.9**
 Computer graphics software, **40.21–40.23**,
40.22f–40.24f
 Computer numeric control (CNC) systems, **9.4**
 Computer-aided design (CAD) software, **40.19**
 Computer-generated holograms (CGHs),
14.1–14.9
 about, **14.1–14.3**
 accuracy limitations of, **14.6–14.7**, **14.6f**, **14.7f**
 discussion of, **14.9**
 experimental results from, **14.7–14.9**, **14.7f**,
14.8f
 interferometers using, **14.4–14.5**, **14.4f**, **14.5f**
 plotting of, **14.3–14.4**, **14.3f**
 sample, **14.2f**
 Concave facets, **39.40**, **39.40f**
 Concentrated arc lamps, **15.47–15.49**, **15.48f**,
15.49f
 Concentration:
 of radiation, **39.1**, **39.5**, **39.6**
 of solution, **38.5**
 Concentrators, nonimaging, **39.12–39.22**
 calculation of, **39.5**, **39.6**
 compound elliptical collectors, **39.14**, **39.15f**

- Concentrators, nonimaging (*Cont.*):
 compound hyperbolic collectors, 39.15, 39.15f, 39.16f
 compound parabolic collectors, 39.13–39.14, 39.13f, 39.14f
 dielectric compound parabolic collectors, 39.15, 39.16, 39.16f
 edge rays, 39.22
 geometrical vector flux, 39.21–39.22
 inhomogeneous media, 39.22
 integrating cavities with, 39.26, 39.27
 and lenses, 39.16–39.17
 and mirrors, 39.17
 multiple surface concentrators, 39.16–39.17, 39.17f
 restricted exit angle concentrators with lenses, 39.18, 39.18f
 RX, 39.17, 39.17f
 RXI, 39.17, 39.17f
 star, 39.20, 39.21
 tapered lightpipes, 39.12–39.13, 39.13f
 θ_1/θ_2 concentrators, 39.18–39.20, 39.19f
 2D vs. 3D, 39.20–39.21, 39.20f, 39.21f
- Condensers, first-order layout for, 1.10–1.11, 1.11f
- Condition of detailed balance (term), 23.23
- Conduction band, 17.4, 17.4f, 17.5f
- Conduction bandgap, 25.3, 25.3f
- Cones (eye receptors), 30.15, 30.16f, 34.37, 34.38, 36.8, 36.8f, 36.9f
- Cone-shaped secondary baffle, 7.3–7.4, 7.3f, 7.4f
- Confidence interval (CI), 34.22
- Configuration factor algebra, 34.14
- Confocal cavity technique, 12.20, 12.20t
- Confocal parameter, 16.23
- Conic collimators, 39.8, 39.9f
- Conic reflectors, 39.11, 39.11f
- Conic surfaces, 3.5
- Conical (term), 35.5
- Conical-directional reflectance, 35.5t, 35.6f, 35.6t
- Conical-hemispherical reflectance, 35.5t, 35.6f, 35.6t
- Conservation, of radiant power transfer, 34.13f
- “Conservation of complexity,” 3.7
- Constraints:
 defined, 3.17
 methods for handling, 3.18–3.19
- Constricted double-heterostructure large optical cavity (CDH-LOC), 19.19, 19.20t, 19.21f
- Consultative Committee on Photometry and Radiometry (CCPR), 36.2
- Contact scanners, 32.21, 32.22f
- Contact stresses, 6.21
- Contacting, in wafer processing, 17.24
- Contamination levels (in stray light suppression), 7.18–7.19, 7.18t, 7.19f–7.21f
- Continuous polishers (CPs), 9.7
- Continuous ring flanges, 6.4f, 6.11
- Continuous wave (cw) lasers, 23.18, 34.32
- Continuous wave (cw) power, 19.19, 19.22, 19.22f
- Convergent reflectors, 39.38–39.40, 39.38f, 39.39f
- Conversion factors:
 for English and SI units, 37.7t
 for photometric and radiometric quantities, 36.11–36.14, 36.12f–36.14f
- Convex surfaces, testing of, 14.5, 14.5f
- Coordinate measurement machines (CMMs), 9.6
- Coordinate measurement method (CMM), 40.53, 40.54
- Copper, 17.28
- Copper vapor lasers (CVLs), 16.12, 16.13f, 16.30
- Copper-doped germanium (Ge:Cu) detectors, 24.84f, 24.85f, 24.96, 24.97, 24.97f–24.99f
- Corner cube prisms, 12.16
- Cornice lighting, 40.13f
- Corning ULE, 6.18
- Correlated color temperature (CCT), 34.44, 37.7, 38.5, 40.8
- Correlated double sampling (CDS), 33.13
- Correlated emission lasers (CELs), 23.42–23.43
- Cosine law, 37.8, 37.8f
- Cosine-to-the-fourth approximation, 34.16
- Coupled cavity lasers, 19.37f, 19.38
- Coupling:
 of circulating pulses, 20.12–20.15, 20.12f, 20.15f
 étendue and source, 40.41–40.42
 in film development, 29.14
 gain, 19.29
 interface, 20.14
 output, 16.13
 phase-conjugated, 20.14
 repetition-rate, 20.14–20.15, 20.15f

- Coupling noise, resistive, 27.5, 27.6f
- Cove lighting, 40.13f, 40.16f
- Critical illumination (*see* Abbe illumination system)
- Critical objects (in stray light suppression), 7.2
 imaged, 7.4, 7.5f
 real-space, 7.2–7.4, 7.3f, 7.4f
- Crossed reflectors, 39.38f
- Crossed string relationship, 39.4, 39.4f
- Cryogenic electrical substitution radiometers, 34.28
- Crystalline optics, 9.8
- Current density, 19.12–19.13, 19.12f, 19.13f
- Current-confined constricted double-heterostructure large optical cavity (CC-CDH-LOC), 19.19, 19.20t, 19.21f
- Curvature measurement, 12.17–12.25
 mechanical methods of, 12.17–12.19, 12.18f, 12.19f, 12.19t
 optical methods of, 12.19–12.21, 12.20f, 12.20t, 12.21f
- Cusp surface, of diamond-turned optics, 10.10, 10.10f
- Cutoff wavelength, 24.10
- Cyanine dyes, 30.13, 30.13f
- Dall compensators, 13.24, 13.24f
- Damped least-squares (DLS) method, 3.17–3.19
- Damping:
 additive, 3.18
 of field by reservoir, 23.33–23.34
- Damping factor, 3.18
- Dark counts (of photomultipliers), 27.8
- Dark current:
 absorption coefficient, 25.8, 25.9f
 in CCDs, 32.20
 correction of, 34.33
 defined, 24.10
 diffusion current, 25.7
 generation-recombination current, 25.7–25.8
 histogram of, 32.12n
 in photosensing elements, 24.19–24.20, 32.10–32.12, 32.11f
 in *pin* photodiodes, 25.7–25.8
 quantum efficiency, 25.8
 responsivity, 25.8
 tunneling current, 25.8
- Dark regions, 17.28
- Dark signal correction, 34.33
- Dark-line defects, 17.28
- Dashed rays, 1.12, 1.12f
- Daylight:
 as natural light source, 40.40–40.41
 simulation of, 40.17
 spectrum of, 40.40f
- Daylighting schemes, luminaires for, 40.47–40.50, 40.49f–40.51f
- Day/night cameras, 31.28–31.29
- Dazzle, 40.9
- dc carbon arcs, 15.25t
- dc lamps, 15.32f
- Decay time, 16.4
- Decorative lighting, 40.14, 40.16f
- Deep-diffused stripe (DDS) lasers, 19.23
- Defocus:
 aberrations of, 11.30
 annular polynomials for, 11.38f, 11.39
 thermal, 8.4, 8.5f
- Density (of photographic films), 29.6–29.8, 29.7f
- Density of states, 19.9–19.10, 19.10f
- Density-operator approach, to quantum theory of lasers, 23.14–23.33
 derivation of Scully-Lamb master equation, 23.17–23.22
 cavity losses, 23.18
 laser master equation, 23.19–23.20, 23.19f
 micromaser master equation, 23.20–23.22
 photon statistics, 23.22–23.27
 laser, 23.22–23.26, 23.23f, 23.25f
 micromaser, 23.26–23.27, 23.27f
 spectrum, 23.28–23.33
 laser field, 23.28–23.31, 23.30f
 micromaser field, 23.31–23.33
 time evolution of the field in Jaynes-Cummings model, 23.15–23.17, 23.15f
- Depletion, of charge, 25.6
- Depletion layer generation current, 32.10–32.11, 32.11f
- Depth of focus, 4.7–4.8, 4.8f
- Design software, optical (*see* Optical design software)
- Detailed balance, condition of, 23.23
- Detective quantum efficiency (DQE), 24.10, 29.1, 29.23
- Detective time constant, 24.10
- Detectivity:
 of film, 29.23
 of infrared detector arrays, 33.23–33.24
 normalized, 38.9

- Directional total absorptance, **35.8t**
 Direct-reading autocollimators, **12.12**
 Disability glare, **40.9–40.10**
 Discomfort, visual, **40.9–40.12**
 Discomfort glare, **40.9–40.12, 40.11t**
 Discrete energy levels, **16.4**
 Disk PZT transducers, **22.17–22.18**
 Dislocation reduction, **18.2, 18.2f**
 Dispersion (of light), **38.8**
 Dispersivity, **30.9**
 Displacement current, **26.7**
 Displays (term), **40.1** (*See also specific displays, e.g.: Monolithic LED displays*)
 Distance measurement (*see* Length measurements)
 Distortion plot, **2.4, 2.4f**
 Distortion tolerances, **5.8**
 Distracting glare, **40.9**
 Distributed backscattering, **20.14**
 Distributed Bragg reflector (DBR) lasers, **19.38, 19.40**
 Distributed feedback (DFB) lasers, **16.29, 19.36, 19.38**
 Distributed grating surface-emitting lasers, **19.40–19.41, 19.40f**
 Distribution temperature, **34.43–34.44, 37.7**
 Divergent reflectors, **39.38–39.40, 39.38f, 39.39f**
 D-log H curve (for photographic films), **29.8–29.10, 29.8f**
 Domes, mounting of, **6.11, 6.12f**
 Doped extrinsic silicon, **33.7, 33.8f**
 Doping, substrate, **17.20**
 Doppler broadening, **16.5, 16.6, 16.6f, 16.9**
 Doppler linewidth (*see* Full width at half maximum)
 Double heterojunction (DH) LEDs, **17.13, 17.13f–17.15f**
 Double heterostructure (DH) lasers, **19.4, 19.5f, 19.7, 19.12–19.15, 19.18–19.19, 19.19f**
 Double heterostructure *pin* photodiodes, **26.13**
 Double monochromators, **35.9, 38.15f**
 Double optical gating, **21.8**
 Double sampling, correlated, **33.13**
 Double-beam spectrophotometers, **35.8–35.9**
 Double-channel planar buried heterostructure (DC-PBH) lasers, **19.24, 19.25f, 19.34f**
 Double-pass photodetectors, **26.4f**
 Doublet lens, air-spaced, **6.7**
 Downconversion, parametric, **23.14**
 Downhill optimizer, **3.17**
 Draft angle, **39.10**
 Drag-wiping (cleaning), **10.9**
 Drift:
 in CCDs, **32.17**
 frequency vs. time, **22.2**
 low offset, **27.11**
 photogenerated charge collection by, **32.5**
 thermocouple junctions as source of, **27.6, 27.6f**
 Driving lasers, in attosecond optics, **21.4–21.6, 21.5f**
 Drop-in assembly, **6.6, 6.6f**
 Dual beam detection, **22.13**
 Dual in-line octocouplers, **17.32, 17.32f**
 Dumet (alloy), **40.29**
 Dye lasers, **16.31–16.32, 16.32f, 20.15–20.16**
 Dye-forming reaction, in film development, **29.14**
 Dyes:
 azomethine, **30.10f, 30.11–30.12, 30.12f**
 cyanine, **30.13, 30.13f**
 light-absorbing, **30.7**
 photographic, **30.10–30.13, 30.10f, 30.12f**
 yellow filter, **30.4**
 Dynodes, in photomultipliers, **27.6–27.9, 27.7f**
 Eberhard effects, **30.3**
 Eccentric pupil design (*see* Z-system)
 ECE (United Nations Economic Commission for Europe), **40.63–40.64**
 Edge lit backlight, **40.47, 40.47f**
 Edge rays, **39.22, 39.38**
 Edge-absorbing photodetectors, **26.4f**
 Edge-emitting lasers (ELASERs), **25.15**
 Edge-emitting LEDs (ELEDs), **25.15**
 Edge-illuminated photodetectors, **26.4f, 26.5**
 Edging step (of optics fabrication), **9.6**
 Einstein (unit), **34.11**
 Einstein's light quanta, **23.6–23.9, 23.8f**
 Einstein's particle hypothesis, **23.7**
 Elastomeric mountings, **6.4, 6.4f, 6.5, 6.12**
 Electrical contact (light bulb), **40.29f**
 Electrical parasitics (laser), **19.34–19.35, 19.34f**
 Electrical substitution radiometers, **34.27–34.29**
 Electrical transfer function, with series inductance, **26.13**
 Electrodeless discharge lamps, **15.36, 15.44**
 Electrodeless fluorescent lamps, **40.36–40.37**
 Electrodeless lamps, **40.25t, 40.26t, 40.36–40.37**

- Electrodeless sulfur lamps (ESLs), **40.36–40.37**
- Electroluminescent light sources, **40.37–40.39**,
40.38f, **40.38t**, **40.39f**
- Electron bombardment (EB), **31.23**
- Electron current, **26.7**
- Electron lenses, **31.8**, **31.8f**
- Electron-bombarded SSAs (EBSSAs),
31.23–31.27
- digital cameras, **31.24–31.26**
- modulation transfer function and limiting
resolution of, **31.26–31.27**
- proximity-focused, **31.23–31.24**, **31.23f**,
31.24f
- Electronically scanned staring FPAs,
33.16–33.17
- Electro-optic modulators (EOMs), **22.14**, **22.20**
- Electro-Optical Industries, Inc. (EOI), **15.14**,
15.15f, **15.16f**
- Elliptical polynomials, **11.21**, **11.25–11.27**,
11.26t–11.27t, **11.36t**, **11.38f**
- Emission, stimulated (*see* Stimulated emission)
- Emission lasers, correlated, **23.42–23.43**
- Emission linewidth (of radiation), **16.4–16.7**,
16.6f, **16.7f**
- Emission-line broadening, **16.4–16.7**, **16.6f**, **16.7f**
- Emissivity:
of blackbody cavity, **15.7**, **15.8f**
tungsten, **40.28f**
- Emittance, **39.2**
- and absorptance, **35.8**
- calculating, **34.25–34.26**
- defined, **35.7**
- measurement of, **35.14–35.16**, **35.15f**
- Emitted photon wavelength, **17.4–17.5**
- Emitters:
AlGaAs, **17.32**
blue, **17.18**, **17.19**
GaAsP, **17.32**
- Emulsions, photographic, **24.100**, **24.101f**, **29.4**,
30.7
- Enclosed arcs, **15.24**, **15.28–15.47**
- high-pressure, **15.24**, **15.28–15.34**
- capillary mercury-arc lamps, **15.30–15.31**,
15.31f
- compact-source arcs, **15.31–15.34**,
15.32f–15.35f
- Lucalox lamps, **15.30**, **15.31f**
- mercury arcs, **15.29**, **15.30f**
- multivapor arcs, **15.29**, **15.31f**
- Uviarc, **15.28–15.29**, **15.29f**, **15.30f**
- Enclosed arcs (*Cont.*):
low-pressure, **15.35–15.47**
- black-light fluorescent lamps, **15.35**, **15.36t**
- electrodeless discharge lamps, **15.36**, **15.44**
- germicidal lamps, **15.35**
- hollow cathode lamps, **15.35**, **15.37–15.43t**,
15.44f
- Pluecker spectrum tubes, **15.47**, **15.47f**,
15.47t
- spectral lamps, **15.44**, **15.45**, **15.45f**,
15.46f, **15.46t**
- Sterilamps, **15.35**, **15.36f**
- End loss, **19.6**
- Energy:
levels of, **16.4**, **16.7**
- luminous, **37.4t**, **37.6**
- measurement of, **34.32**
- nomenclature for, **36.4**, **36.5**
- radiant, **34.7**, **37.4t**, **37.6**
- units of, **34.5–34.6**
- Energy band structure, **17.3–17.6**, **17.3f–17.5f**
- Energy bandgap, **25.3**, **25.3f**
- English units, and SI units, **37.7**, **37.7t**
- Entrance pupil, **34.18**, **34.19f**
- Entrance window, **34.19**, **34.19f**
- Environmental specifications, optical, **4.10**
- Epitaxial growth, **17.8**, **17.21**
- Epitaxial technology (for LEDs), **17.21–17.23**
- Epoxy, in indicator lamps, **17.29**
- Equivalent neutral density (END), **29.15**
- Equivalent noise input (ENI), **24.11**
- Equivalent veiling luminance (EVL), **40.10**
- Error functions, in optical design software,
3.17, **3.19–3.20**
- Error types, in absolute measurements,
34.21–34.23
- Étendue, **34.15**
- defined, **40.42**
- geometrical, **38.8**
- in nonimaging optics, **39.2**, **39.3**,
39.4f, **39.5**
- and source coupling, **40.41–40.42**
- Étendue loss, **39.6**
- Evaluation function (of optical software),
3.8–3.16
- of aberrations, **3.9–3.11**
- paraxial ray-trace, **3.8–3.9**, **3.9f**
- ray-trace, **3.11–3.13**, **3.12f**
- by spot-diagram analysis, **3.13–3.16**
- Event-driven programs (optical software), **3.7**

- Exact rays (term), 3.3, 3.11–3.12
- Excimer lasers, 16.30–16.31, 16.31*f*
- Excimers, in fluorescent lamps, 40.31
- Excited state, of azomethine dyes, 30.11–30.12, 30.12*f*
- Exciton recombination, 17.6
- Exit pupil, 34.18, 34.19*f*
- Exitance, 39.2
- defined, 34.8
 - luminous, 37.4*t*, 37.5, 37.5*f*
 - radiant, 15.4–15.6, 15.5*t*, 15.6*f*; 37.4*t*, 37.5, 37.5*f*
- Exposure:
- luminous, 37.4*t*, 37.6
 - of photographic films, 29.5–29.6
 - radiant, 37.4*t*, 37.6
- Extended baffle shields, 7.9–7.10, 7.9*f*, 7.10*f*
- Extended wavelength photodetectors, 25.10, 25.10*t*
- Exterior lighting, 40.61–40.62, 40.63*t*
- External cavity diode lasers (ECDLs), 22.21–22.23, 22.22*f*
- Extreme infrared (IR) light, 25.2
- Extrinsic photoconductors, 25.5, 25.5*f*
- Extrinsic photodetectors, 24.7, 24.7*f*
- Extrinsic semiconductor transition, 24.11
- Eye, human (*see* Human eye)
- Fabrication, optical, 9.3–9.9
- about, 9.3
 - aspherical, 9.7–9.8, 9.7*f*
 - crystalline, 9.8
 - by diamond turning (*see* Diamond turning)
 - diamond turning vs. traditional, 10.6
 - guide to methods of, 10.3*t*
 - material formation for, 9.3–9.4
 - methods of, 10.3*t*
 - plano, 9.7
 - spherical, 9.4–9.6
- Fabry-Perot interferometers, 16.19*f*
- Faceted reflectors, 39.10*f*, 39.39–39.41, 39.39*f*, 39.40*f*
- Failures per 10⁹ hours (FITS), 17.25
- False-colored infrared film, 30.22
- Far infrared (FIR) radiation, 24.3, 25.2
- Far ultraviolet radiation, 15.12, 15.13
- Fatigue, thermal, 17.25
- Federal Motor Vehicle Safety Standards (FMVSS), 40.63
- Feedback:
- optical, 16.2
 - resonant optical, 19.38, 19.38*f*
 - stabilization of, 34.32
- Femtoseconds, 20.1
- Fermi occupation functions, 19.11
- Ferroelectric bolometer arrays, 28.11, 28.12, 28.12*f*
- Ferroelectric detectors, 33.10
- Fiber lasers, 16.34
- Fiber optics:
- fiber alignment for, 17.33
 - focal-length measurement with, 12.25
 - LED considerations with, 17.33–17.34
 - in nonimaging optics, 39.21
- Fiber-optic octocouplers, 17.33–17.34
- Fiberoptic-coupled (FO) II SSAs, 31.20–31.22, 31.20*f*, 31.21*f*, 31.21*t*
- Field curvature plot, 2.4, 2.4*f*, 2.5
- Field effect transistors (FETs), 27.10
- Field lenses, first-order layout for, 1.8, 1.10, 1.10*f*
- Field of view (FOV), 3.4, 7.11, 24.11, 31.1
- Field patch trace, 39.7–39.8
- Field plots, of aberration curves, 2.4–2.5
- Field stops, 7.7, 7.8*f*, 7.9*f*, 34.18–34.19, 34.19*f*
- Field-enhanced pyroelectric arrays (*see* Ferroelectric bolometer arrays)
- Figures of merit (FOM), 24.13, 33.23–33.28
- for infrared photodetectors, 25.12
 - minimum resolvable temperature (MRT), 33.27–33.28, 33.27*f*
 - NE ΔT , 33.25*f*, 33.26*f*
 - in spectroradiometry, 38.5–38.6
- Filament notching, 40.30
- Filaments:
- lamp, 15.20*f*
 - light bulb, 40.25, 40.27, 40.29–40.30, 40.29*f*
- Fill gases, light bulb, 40.29*f*, 40.30
- Film, photographic (*see* Photographic films)
- Filters:
- blocking, 38.8
 - blur, 32.34, 32.34*f*
 - in II SSA cameras, 31.7
 - infrared, 40.12
 - interference, 34.36
 - neutral density, 40.52
 - notch, 22.10–22.11, 22.10*f*, 22.11*f*
 - UV, 40.12
- Finitely distant objects, systems with, 1.6, 1.6*f*

- First-order layout techniques, 1.3–1.16
- achromatism, 1.14–1.15, 1.15*f*
 - afocal attachments, 1.8, 1.9*f*
 - afocal systems, 1.7, 1.7*f*
 - athermalization, 1.15–1.16, 1.16*f*
 - axial/principal rays, 1.12
 - component power minimization, 1.13, 1.13*f*
 - condensers, 1.10–1.11, 1.11*f*
 - defined, 1.4
 - field lenses, 1.8, 1.10, 1.10*f*
 - magnifiers and microscopes, 1.8
 - ray-tracing, 1.4–1.5
 - reasonableness of layout, 1.13–1.14
 - two-component systems, 1.5–1.7
 - zoom or varifocal systems, 1.11–1.12
- “Fitness for use,” 17.25
- Five-axis machining, 10.7*f*
- 5 × 7 matrix LED displays, 17.31–17.32
- Fixed interferogram evaluation, 13.14–13.15
- Fixed-orientation mirrors, 6.17
- Fixer, film, 29.5
- Fizeau interferometers, 12.14, 13.8–13.9, 13.9*f*, 13.10*f*, 13.18, 14.4, 14.5*f*
- Flaming arcs, 15.23, 15.24*f*, 15.26*t*–15.27*t*
- Flanges:
- annular, 6.3–6.4, 6.4*f*
 - continuous ring, 6.4*f*, 6.11
- Flash blindness, 40.9
- Flash lamps, 16.16–16.17, 16.17*f*
- Flex-Pivots (flexures), 6.19
- Flexure mountings, 6.5, 6.5*f*, 6.15–6.17, 6.16*f*
- Flicker:
- in fluorescent lamps, 40.32–40.33
 - impact of, 40.12
 - in incandescent lamps, 40.30
- Flicker floor, 22.3
- Flicker noise, 24.11
- Flight control, 19.3
- Flip chip packaging, 18.6, 18.6*f*
- Flip-and-fold approach, 39.29*f*
- Floating diffusion, 32.14, 32.15*f*
- Floating gate output amplifiers, 32.15
- Fluorescence, 34.13, 40.30
- Fluorescent lamps, 40.30–40.33
- applications for, 40.26*t*
 - characteristics of, 40.25*t*
 - construction of, 40.33*f*, 40.34*f*
 - elements of, 40.31*f*
 - emission spectrum of, 40.32*f*, 40.35*f*
 - types of, 40.28*f*
- Flux:
- luminous, 37.4, 37.4*t*, 37.6
 - radiant, 37.3, 37.4*t*
 - total luminous, 37.4*t*, 37.6
 - total radiant, 37.4*t*, 37.6
- Flux budget, for fiber optics, 17.33
- Flux density (*see* Irradiance)
- Fly-by-light (FBL), 19.3
- FM spectroscopy, 22.13–22.14
- F-number, 34.20
- Focal depth, 4.7–4.8, 4.8*f*
- Focal length, 1.5–1.7, 12.21–12.25
- fiber optics, 12.25
 - focimeters, 12.22–12.23, 12.22*f*, 12.23*f*
 - Fourier transforms, 12.24
 - microlenses, 12.24
 - Moiré deflectometry, 12.23, 12.24*f*
 - nodal slide bench, 12.22, 12.22*f*
 - Talbot autoimages, 12.23, 12.24
- Focal plane arrays (FPAs), 33.3
- hybrid, 33.14–33.23
 - microbolometer, 33.13–33.14
 - MIS photogate, 33.10–33.11, 33.11*f*, 33.12*f*
 - monolithic, 33.10–33.14
- Focal ratio, 34.20
- Focimeters, 12.22–12.23, 12.22*f*, 12.23*f*
- Focus athermalization techniques, 6.22–6.24
- active athermalization, 6.24, 6.24*f*
 - passive athermalization, 6.22, 6.23*f*, 6.24
 - single material designs, 6.22, 6.23*f*
- Focus distance, 1.5–1.7
- Focus shift, thermal, 8.2–8.4, 8.3*t*, 8.4*t*
- Fog, film, 29.9
- Fokker-Planck equation, 23.37
- Foot-candle (unit), 34.43, 36.7, 36.7*t*, 37.7*t*
- Foot-lambert (unit), 34.43, 36.7, 36.8*t*, 37.7, 37.7*t*
- Forbes method, 3.20
- Fore-optics, 38.7
- Forward light, 40.43, 40.44*f*
- Forward looking infrared (FLIR), 33.4
- Foucault test, 12.19, 13.2–13.3, 13.2*f*, 13.3*f*
- Fourier analysis, of interferograms, 13.16–13.17, 13.17*f*
- Fourier transform spectrophotometers, 35.9
- Fourier transforms, for focal-length determination, 12.24
- Four-phase CCDs, 32.13–32.14, 32.13*f*, 32.16*f*
- Frame interline transfer (FIT) CCDs, 32.27*f*, 32.29

- Frame transfer (FT) CCD image sensors, 32.26–32.28, 32.27f, 32.28f, 32.32
- Frame transfer (FT) CCD TDI FPAs, 33.11, 33.12f
- Free-electron lasers (FELs), 16.36–16.37, 16.37f, 23.43–23.45
- Free-electron-laser (FEL) lamps, 15.11, 15.12, 15.13f
- Frequency:
- and drift, 22.2
 - phase and amplitude responses vs., 22.6–22.7, 22.6f, 22.7f
 - stability of, 22.2
- Frequency comb, 20.1, 20.2, 20.7–20.9
- Frequency discriminators:
- for laser locking, 22.12–22.14
 - optical cavity-based, 22.14–22.16, 22.17f
- Frequency modulation (FM), 19.36, 19.36f, 22.4
- Frequency shift keying (FSK), 19.36
- Frequency-controlled lasers, 22.7, 22.7f
- Frequency-resolved optical gating (FROG), 21.8, 21.9
- Frequency-resolved optical gating for complete reconstruction of attosecond bursts (FROG-CRAB), 21.9
- Frequency-selective-feedback lasers, 19.37f, 19.38
- Fresnel lenses, 39.9–39.10, 39.10f, 40.45f, 40.46
- Full width at half maximum (FWHM), 16.5, 16.6, 20.3, 21.2
- Functional specifications (optical design), 4.2
- Fundamental array mode, 19.27
- Furniture-integrated lighting system, 40.13f
- Gain:
- defined, 16.9
 - in photomultipliers, 27.7
- Gain coefficient, 16.9–16.10
- Gain coupling, 19.29
- Gain medium, 16.3
- Gain saturation, 16.10
- Gain stability margin, 22.9–22.10
- Gain-bandwidth (GB), 26.17
- Gain-coupled arrays, 19.27
- Gain-guided phased array, 19.28f
- Galilean telescopes, 1.7f
- Gallium aluminum arsenide (GaAlAs) LEDs, 17.12, 17.12f, 17.13f
- Gallium arsenide (GaAs) lasers, 19.7
- Gallium arsenide (GaAs) LEDs, 17.8, 17.9, 17.9f
- Gallium arsenide (GaAs) semiconductor diode lasers, 16.18–16.19, 16.18f
- Gallium arsenide phosphide (GaAsP) emitters, 17.32
- Gallium arsenide phosphide (GaAsP) LEDs, 17.9–17.10, 17.10f, 17.15–17.17
- energy band diagram for, 17.5f
 - homojunction in, 17.10, 17.11f
 - light degradation in, 17.27f
 - performance summary of chips in, 17.16t
- Gallium arsenide phosphide (GaAsP) photodiodes, 24.49, 24.49f, 24.50f
- Gallium arsenide phosphide (GaAsP) substrate, 17.22
- Gallium arsenide (GaAs) quantum well photodetectors, 25.16–25.17, 25.16f, 25.17f
- Gallium nitride (GaN) photovoltaic detectors, 24.42, 24.43, 24.45f, 24.46, 24.46f, 24.47
- Gallium nitride (GaN) substrate, 17.22
- Gallium phosphide (GaP), 17.16, 17.21–17.22
- Gallium phosphide (GaP) dynodes, 24.42, 24.44f
- Gallium phosphide (GaP) photodiodes, 24.47–24.49, 24.48f
- Gallium phosphide (GaP) substrate, 17.20–17.22
- Gallium-doped germanium (Ge:Ga) infrared detectors, 24.100
- Gallium-doped silicon (Si:Ga) infrared detectors, 24.95, 24.95f, 24.96, 24.96f
- Galvo-driven Brewster plates, 22.18–22.19
- Gas chromatography-mass spectroscopy (GC-MS), 33.4
- Gas lights, 40.40
- Gas mantle, 15.17–15.19, 15.19f
- Gaseous laser gain media, 16.30–16.31, 16.31f
- Gas-filled lamps, 34.31
- Gate modulation, 33.19t, 33.20f, 33.22
- Gated integration, 27.12–27.13, 27.13f, 27.15
- Gating:
- amplitude, 21.7
 - double optical, 21.8
 - frequency resolved, 21.8, 21.9
 - FROG, 21.9
 - FROG-CRAB, 21.9
 - polarization, 21.7–21.8
 - two-color, 21.7
- Gauss illuminated eyepieces, 12.12

- Gaussian intensity distribution, **39.28, 39.29f**
 Gaussian line shape, **16.6f**
 Gaussian mode, **16.21**
 Gaussian parameters (optical design), **4.5–4.6, 4.6t**
 Gaussian-shaped beam, **16.22**
 General Conference on Weights and Measures (CGPM), **36.2**
 General Electric, **15.29, 15.30, 15.48, 19.29t**
 General system data (optics), **3.3, 3.4**
 Generating step (of optics fabrication), **9.4**
 Generation noise, **24.11**
 Generation-recombination (GR) current, **25.7–25.8**
 Generation-recombination (GR) noise, **24.11**
 Geometrical configuration factor (GCF), **7.1, 7.2, 7.22**
 Geometrical etendue, **38.8**
 Geometrical optical transfer function (GOTF), **3.16**
 Geometrical optics, **3.16**
 Geometrical vector flux, **39.21–39.22**
 Geometry-controlled lasers, **19.37f, 19.38**
 Germanium (Ge) avalanche photodiodes, **24.70f, 24.72–24.73, 24.72f, 24.73f**
 Germanium (Ge) bolometers, **28.5, 28.7t**
 Germanium (Ge) detectors:
 copper-doped, **24.84f, 24.85f, 24.96, 24.97, 24.97f–24.99f**
 gallium-doped infrared, **24.100**
 gold-doped, **24.83–24.85, 24.84f–24.86f**
 intrinsic photodetectors, **24.70–24.73, 24.70f–24.73f**
 mercury-doped, **24.84f, 24.92–24.95, 24.93f–24.95f**
 pn and *pin*, **24.70–24.71, 24.70f–24.72f**
 zinc-doped, **24.84f, 24.98–24.100, 24.99f**
 Germanium gallium arsenide (GeGaAs) photodiodes, **34.31**
 Germanium (Ge) intrinsic photodetectors, **24.70–24.73, 24.70f–24.73f**
 Germanium (Ge) low-temperature bolometers, **24.31–24.32, 24.32f, 24.33f**
 Germanium photodiodes, **38.9, 38.9t**
 Germicidal lamps, **15.35**
 Glare:
 and exterior lighting, **40.62**
 limiting of, **40.10, 40.41**
 and visual discomfort, **40.9–40.12, 40.11t**
 and windows, **40.41**
 Glare stops (*see* Lyot stops)
 Glass:
 formation of optical, **9.3–9.4**
 optical, **5.9**
 as photographic film emulsion, **29.4**
 tolerances for, **5.9**
 Glass envelope, light bulb, **40.29, 40.29f**
 Glazing, window, **40.41**
 Global, **15.17, 15.18f, 15.19, 15.19f**
 Glossiness, **40.5**
 Glow lamps, **40.39**
 Glow modulator tubes, **15.49, 15.50f, 15.51f, 15.52t**
 Gobs, glass, **9.4**
 Golay cell detectors, **28.2, 28.6, 28.7t**
 Gold, diamond turning and, **10.5**
 Gold-doped germanium (Ge:Au) detectors, **24.83–24.85, 24.84f–24.86f**
 Gold-germanium (Au-Ge) alloys, **17.24**
 Goldpoint blackbody, **15.9**
 Gold-zinc (Au-Zn) alloys, **17.24**
 Goniometers (goniophotometers), **12.10, 40.52–40.53, 40.53f, 40.54f**
 Gouffé method, **15.7–15.9, 15.8f**
 Graded-index separate-confinement heterostructure (GRIN SCH) quantum lasers, **19.14, 19.14f**
 Gradient-freeze technique, **17.21**
 Grains and graininess:
 of photographic films, **29.5**
 of photographic images, **29.18t, 29.19–29.22, 29.21f**
 of silver halide crystals, **29.4**
 Granularity, photographic film speed and, **30.19**
 Grating equation, **38.7–38.8**
 Grating surface-emitting laser array, **19.40–19.41, 19.40f**
 Gray gel, **30.4**
 Graybody, **35.7**
 Green light, color film and, **29.13, 29.13f**
 Green-emitting AlInGaP devices, **17.18**
 Grinding, aspheric, **9.8**
 Ground loop noise, **27.5, 27.6f**
 Ground state, **16.4**
 Grown homojunctions, LED, **17.8, 17.9, 17.9f**
 Growth techniques:
 for epitaxial layers, **17.21–17.23**
 substrate, **17.20, 17.21**
 Guard ring, **24.11**

- Guide to the Expression of Uncertainty in Measurement* (IS), 38.6
- Gurney-Mott mechanism, 29.5
- H-aggregates, 30.13
- Haidinger interferometers, 12.14
- Halation, 30.4
- Halogen lamps, 15.11, 15.12, 15.13*f*, 40.25*t*, 40.26*t*, 40.30
- Halon, 38.12–38.13
- Halophosphates, 40.31
- Hamiltonian rays, 3.12
- Hanbury-Brown-Twiss (HB&T) effect, 23.13, 23.13*f*, 23.14
- Hard mounting, of optics, 6.1–6.4, 6.3*f*, 6.4*f*
- Hartmann test, 13.4–13.6, 13.5*f*
- Hartmann-Shack test, 13.6–13.7, 13.6*f*
- H&D curve (*see* D-log H curve)
- Headlamps:
 - design of, 40.21, 40.23, 40.23*f*
 - low-beam, 40.64–40.67, 40.64*f*, 40.65*t*, 40.66*f*, 40.66*t*
- Health-care facility lighting, 40.58–40.60, 40.60*t*
- Heat-pipe blackbody furnace, 15.9*f*
- Height solves, 3.6
- Heisenberg-Langevin approach, to quantum theory of lasers, 23.33–23.35
- Helium-cadmium (He-Cd) lasers, 16.6, 16.15, 16.15*f*, 16.30
- Helium-neon (He-Ne) lasers, 16.15, 16.15*f*, 16.30
- Hemispherical emittance, 35.15
- Hemispherical total absorptance, 35.8*t*
- Hemispherical-conical reflectance, 35.5*t*, 35.6*f*, 35.6*t*
- Hemispherical-directional reflectance, 35.5*t*, 35.6*f*, 35.6*t*
- Hemispherical-spectral absorptance, 35.8*t*
- Heterodyne interferometers, 13.22
- Heterojunction lasers, 19.4
- Heterojunctions, 17.12, 17.12*f*–17.15*f*, 17.13, 17.17, 26.9
- Hewlett-Packard double-frequency distance-measuring interferometer, 12.9–12.10, 12.9*f*
- Hexagonal polynomials, 11.21, 11.22*t*–11.25*t*, 11.36*t*, 11.38*f*, 11.39
- High-accuracy spectrophotometers, 35.9
- High-brightness visible LEDs (HB-LEDs), 18.1–18.6
 - about, 18.1
 - epitaxial growth of, 18.3
 - packaging of, 18.5–18.6, 18.5*f*, 18.6*f*
 - processing of, 18.3–18.4, 18.3*f*
 - semiconductor material systems for, 18.1–18.2
 - solid-state lighting with, 18.4–18.5, 18.4*f*
 - structure for modern InGaN, 18.2*f*
 - substrates for, 18.2–18.3, 18.2*f*
- High-dye-yield yellow couplers, 30.6
- High-energy radiation, 15.40, 30.19–30.20
- High-gain oscillators, 20.10–20.12, 20.11*f*
- High-intensity carbon arc lamps, 15.21–15.23, 15.24*f*
- High-intensity discharge (HID) lamps, 40.33–40.36
 - applications for, 40.26*t*
 - characteristics of, 40.25*t*
 - CMH, 40.25*t*, 40.26*t*, 40.33, 40.36
 - construction of, 40.34*f*
 - emission spectrum of, 40.35*f*
 - Hg, 40.25*t*, 40.26*t*, 40.33, 40.36
 - HPS, 40.25*t*, 40.26*t*, 40.33
 - MH, 40.25*t*, 40.26*t*, 40.33, 40.35, 40.35*f*, 40.36
- High-intensity reciprocity failure, of photographic films, 29.12
- High-order harmonic generation, 21.2, 21.2*f*
- High-power diode lasers, 19.19–19.23, 19.20*t*, 19.21*f*, 19.22*f*
- High-power laser arrays, 19.26–19.30, 19.28*f*, 19.28*t*, 19.29*t*, 19.30*f*
- High-power lasers, 19.24, 19.25*f*, 19.25*t*
- High-power semiconductor lasers, 19.18–19.30
 - arrays in, 19.26–19.29, 19.28*f*, 19.28*t*, 19.29–19.30, 19.29*t*, 19.30*f*
 - commercial diode, 19.19–19.23, 19.20*t*, 19.21*f*, 19.22*f*
 - future directions for, 19.23–19.26, 19.25*f*, 19.25*t*, 19.26*t*, 19.27*f*
 - mode-stabilized lasers with reduced facet intensity, 19.18–19.19, 19.19*f*
- High-power strained QW lasers, 19.26, 19.26*t*
- High-pressure, short-arc xenon lamps, 15.35*f*
- High-pressure enclosed arcs, 15.24, 15.28–15.34
 - capillary mercury-arc lamps, 15.30–15.31, 15.31*f*

- High-pressure enclosed arcs (*Cont.*):
 compact-source arcs, 15.31–15.34,
 15.32f–15.35f
 Lucalox lamps, 15.30, 15.31f
 mercury arcs, 15.29, 15.30f
 multivapor arcs, 15.29, 15.31f
 Uviarc, 15.28–15.29, 15.29f, 15.30f
- High-pressure mercury-arc lamps,
 15.29f, 15.30f
- High-speed modulation, of semiconductor
 lasers, 19.30–19.36, 19.31f–19.36f
- High-speed optical recording systems, 19.3
- High-speed photoconductors, 26.20–26.23,
 26.21f–26.23f
- High-speed photodetectors, 26.1–26.24
 about, 26.3
 avalanche photodetectors, 26.17–26.20,
 26.18f, 26.20f, 26.21f
 photoconductors, 26.20–26.23, 26.21f–26.23f
pin photodiodes, 26.10, 26.12–26.15
 resonant, 26.15, 26.15f
 vertically illuminated, 26.3, 26.4f, 26.5,
 26.10, 26.12–26.13, 26.12f
 waveguide, 26.13–26.14, 26.14f
- Schottky photodiodes, 26.16, 26.16f, 26.17f
- speed limitations on, 26.5–26.10
 carrier transit time, 26.6–26.7, 26.6f
 carrier trapping, 26.9, 26.9f
 diffusion current, 26.8, 26.8f, 26.9
 packaging, 26.9–26.10, 26.10f, 26.11f
 RC time constant, 26.7–26.8, 26.7f
 structures of, 26.3–26.5, 26.4f
- High-speed photographic films, 30.18–30.20
- High-voltage power supply (HVPS), 31.1,
 31.9, 31.10f
- Hindle mounts, 6.19, 6.19f
- Hole current, 26.7
- Hole-accumulated photodiodes (HADs), 32.4f,
 32.8
- Hollow cathode lamps, 15.35, 15.37t–15.43t,
 15.44f
- Hollow lightpipes, 39.30–39.31
- Holograms, computer-generated (*see* Computer-
 generated holograms)
- Holographic compensators, 13.25
- Homogeneous broadening, emission-line, 16.5,
 16.6, 16.6f, 16.9
- Homogeneous reflectors, 39.39
- Homogeneous temperature change, 8.2–8.6,
 8.3t, 8.4t, 8.5f, 8.6f
- Homojunction lasers, 19.4
- Homojunctions, LED, 17.8–17.10, 17.9f–17.11f
- Horizontal illuminance, 40.7, 40.18f
- Horizontally illuminated photodetectors,
 26.4f, 26.5
- Hot cathode fluorescent lamps, 40.32
- Hot-electron bolometers, 24.29, 24.30, 24.30f,
 24.31f
- Hottel strings, 39.4, 39.14
- Hub mounting, 6.17, 6.18f
- Hubble telescope, 11.4, 13.24
- Hue, 40.5
- Human eye, 30.15–30.16, 30.16f, 34.6
 cones in, 36.8, 36.8f, 36.9f
 rods in, 36.8–36.10, 36.8f, 36.9f
 wavelengths detectable by the, 36.8–36.10,
 36.8f, 36.9f
- Humidity specifications, for lenses, 4.10
- Hybrid arrays, pyroelectric, 28.11–28.12,
 28.11f, 28.12f
- Hybrid FPAs:
 direct readout architectures of,
 33.15–33.18
 electronically scanned staring FPAs,
 33.16–33.17
 output circuits, 33.18
 TDI scanning FPAs, 33.17, 33.17f
 X-Y addressing and clock generation,
 33.16
- input circuits of, 33.18–33.23, 33.19t, 33.20f
 buffered direct injection, 33.19t, 33.20f,
 33.21
 capacitive transimpedance amplifier,
 33.19t, 33.20f, 33.22–33.23
 chopper-stabilized BDI, 33.19t, 33.20f,
 33.21–33.22
 direct detector integration, 33.18, 33.19t,
 33.20f
 direct injection, 33.18–33.21, 33.19t,
 33.20f, 33.21f
 gate modulation, 33.19t, 33.20f, 33.22
 thermal expansion match in, 33.14
- Hybrid reflectors (*see* Faceted reflectors)
- Hyde maxim, 3.22
- Hydrogen arc lamps, 15.49, 15.53f
- Hypo (film fixer), 29.5
- Ideal mode-locked lasers, 20.7
- Ideal thermal detectors, 28.2–28.3, 28.3f
- Ilford Photo Corporation, 29.25

- Illuminance, *37.4t*, *37.5*, *37.5f*, *39.2t*
 defined, *34.11*, *34.40*, *40.1*
 guidelines on levels of, *40.7t*
 and lighting design, *40.7*
 and luminance, *37.9*, *37.9f*
 retinal, *34.40–34.42*
 uniformity of, *40.7*
 unit conversions for, *36.7t*, *36.8t*
 units of, *34.43*
- Illuminance meters, *34.42*, *40.51*, *40.52f*
- Illuminated ceilings, *40.13f*
- Illuminated eyepieces, *12.12*, *12.12f*
- Illuminated objects (in stray light suppression),
 7.5, *7.5f*, *7.6f*
- Illuminating Engineering Society (IES),
 40.19
- Illumination:
 guidelines on, *40.7t*
 in nonimaging objects, *39.1*
 (See also Uniform illumination, of
 nonimaging optics)
- Illumination Engineering Society of North
 America (IESNA), *36.2*, *36.3*, *37.11*, *40.2*,
 40.7t
- Illumination subsystem, of nonimaging optics,
 39.22
- Image dissectors, *39.21*, *39.21f*
- Image height, *1.4*
- Image intensifiers (IIs), *31.7–31.18*,
 31.8f–31.10f
 defined, *31.8*
 input window/photocathode assemblies for,
 31.10–31.12, *31.11f*, *31.12f*
 MCP IIs, *31.7*, *31.9*, *31.9f*, *31.10f*
 and microchannel plates, *31.12–31.14*,
 31.12f, *31.13f*, *31.13t*
 phosphor screen assemblies for, *31.14–31.16*,
 31.14t, *31.15f*
 proximity-focused MCP IIs, *31.16–31.18*,
 31.17t, *31.18f*, *31.19f*
- Image irradiance, *4.7*
- Image lag, *32.6*
- Image processing, *8.12*
- Image quality, *4.6–4.7*
- Image sensors, *32.2–32.12*, *32.3f*, *32.21–32.34*
 antiblooming in, *32.9*, *32.10f*
 area arrays of, *32.24–32.32*, *32.25t*
 CCD performance, *32.32*
 frame transfer CCDs, *32.26–32.28*,
 32.27f, *32.28f*
- Image sensors, area arrays of (*Cont.*):
 interline transfer CCDs, *32.28–32.32*,
 32.29f–32.31f
 MOS, *32.25–32.26*, *32.26f*
 color imaging with, *32.32–32.34*, *32.33f*, *32.34f*
 dark current in, *32.10–32.12*, *32.11f*
 junction photodiodes, *32.3–32.6*, *32.4f*, *32.6f*
 linear arrays of, *32.21–32.24*, *32.22f*, *32.23f*
 MOS capacitors, *32.7–32.8*
 photoconductors, *32.8–32.9*
 pinned photodiodes, *32.8*
- Image size, *1.4*
- Image specifications, for lenses, *4.3*, *4.6–4.8*, *4.8f*
- Image structure, of photographic systems, *29.17*
- Imaged critical objects, *7.4*, *7.5f*
- Image-intensified (II) electronic imaging,
 31.1–31.30
 about, *31.2–31.3*, *31.3f*
 applications for, *31.27–31.30*
 active imaging, *31.29–31.30*
 day/night cameras, *31.28–31.29*
 mosaic II SSA cameras, *31.29*
 optical multichannel analyzers, *31.27–31.28*
 range gating and LADAR, *31.28*
 image-intensifier modules of, *31.7–31.18*,
 31.8f–31.10f
 input window/photocathode assemblies,
 31.10–31.12, *31.11f*, *31.12f*
 microchannel plates, *31.12–31.14*, *31.12f*,
 31.13f, *31.13t*
 phosphor screens, *31.14–31.16*, *31.14t*,
 31.15f
 proximity-focused MCP IIs, *31.16–31.18*,
 31.17t, *31.18f*, *31.19f*
 optical interface of, *31.3–31.7*
 considerations, *31.6–31.7*, *31.6f*
 photometry and camera lens, *31.5–31.6*
 quantum limited imaging conditions,
 31.3–31.4
 radiometry, *31.4–31.5*
 self-scanned arrays, *31.19–31.27*, *31.19f*
 electron-bombarded, *31.23–31.27*, *31.23f*,
 31.24f
 fiberoptic-coupled, *31.20–31.22*, *31.20f*,
 31.21f, *31.21t*
 lens-coupled, *31.22–31.23*, *31.22f*
 Image-intensified self-scanned arrays (II SSAs),
 31.19–31.27, *31.19f*
 for active imaging, *31.29–31.30*
 camera for, *31.5–31.6*

- Image-intensified self-scanned arrays (II SSAs)
(*Cont.*):
 electron-bombarded, 31.23–31.27, 31.23f, 31.24f
 fiberoptic-coupled, 31.20–31.22, 31.20f, 31.21f, 31.21t
 lens-coupled, 31.22–31.23, 31.22f
- Impedance:
 in amplifiers, 27.10–27.11
 in photodetectors, 24.19–24.20
- Impurity band conduction (IBC), 33.7
- Incandescence, 40.25
- Incandescent sources (of radiation), 40.25, 40.27–40.30
 calibration of, 34.31
 characteristics of, 40.25t
 elements of, 40.29f
 nongaseous, 15.15–15.21
 comparisons of, 15.19, 15.19f
 gas mantle, 15.17, 15.18, 15.19f
 globar, 15.17, 15.18f
 Nernst glower, 15.14, 15.15, 15.17, 15.17f
 quartz-envelope lamps, 15.20, 15.21
 tungsten-filament lamps, 15.19, 15.20, 15.20f–15.22f
 tungsten emissivity in, 40.28f
- Index-guided lasers, 19.8, 19.27, 19.28f
- Indicator lamps, LED, 17.27f, 17.29–17.30, 17.29f
- Indirect glare, 40.9
- Indirect lighting, 40.14, 40.15, 40.15f, 40.16f, 40.46f
- Indirect semiconductors, 17.4, 17.4f, 17.5f, 17.6
- Indium antimonide (InSb) hot-electron bolometers, 24.29, 24.30, 24.30f, 24.31f
- Indium antimonide (InSb) intrinsic photovoltaic detectors, 24.80–24.83, 24.82f, 24.83f
- Indium arsenide (InAs) photovoltaic detectors, 24.75, 24.77–24.78, 24.77f–24.79f
- Indium gallium arsenic phosphide (InGaAsP) laser material system, 19.7
- Indium gallium arsenide (InGaAs) detectors, 24.65–24.70, 24.66f–24.69f
- Indium gallium arsenide (InGaAs) photodetectors, 25.10, 25.10t
- Indium gallium arsenide (InGaAs) photodiodes, 24.66–24.70, 24.66f–24.69f, 34.31
- Indium gallium nitride (InGaN) HB-LEDs, 18.2f
- Indium phosphide (InP) laser material system, 19.7
- Induction lamps (ILs), 40.36–40.37
- Inductive pickup noise, 27.5, 27.6f
- Inductively coupled plasma (ICP), 18.3
- Industrial lighting, 40.60–40.61, 40.61f
- “Infant mortality period,” 17.26, 17.26f
- Infectious film developers, 29.5
- Infinitely distant objects, systems with, 1.5–1.6, 1.6f
- Information capacity, of photographic systems, 29.24
- Infrared detector arrays, 33.1–33.31
 about, 33.3–33.4
 applications for, 33.4
 current status of, 33.28–33.30, 33.28f, 33.29f, 33.29t
 future trends and technology directions of, 33.30–33.31, 33.30f, 33.31f
 hybrid FPAs, 33.14–33.23
 detector interface input circuit, 33.18–33.23, 33.19t, 33.20f, 33.21f
 hybrid readout, 33.15–33.23
 readout, 33.17f
 thermal expansion match in, 33.14
 monolithic FPAs, 33.10–33.14
 direct-charge-injection silicon FPAs, 33.11f, 33.13
 microbolometer FPAs, 33.13–33.14
 MIS photogate FPAs, 33.10–33.11, 33.11f, 33.12f
 scanning and staring, 33.14
 silicon FPAs, 33.11–33.13, 33.11f, 33.12f
 operating principles of, 33.7–33.10, 33.8f, 33.9f
 performance of, 33.23–33.28
 detectivity, 33.23–33.24
 minimum resolvable temperature, 33.27–33.28, 33.27f
 NE ΔT , 33.24–33.27, 33.25f, 33.26f
 percentage of BLIP, 33.24
 scanning and staring, 33.6–33.7, 33.6f
 spectral bands for, 33.4–33.5, 33.6f
- Infrared detectors:
 gallium-doped germanium, 24.100
 gallium-doped silicon, 24.95, 24.95f, 24.96, 24.96f
- Infrared film, 30.22
- Infrared filters, 40.12
- Infrared interferometry, 13.25
- Infrared LED chips, 17.8, 17.9, 17.9f
- Infrared photodetectors, 25.12, 25.15

- Infrared (IR) radiation, **34.6, 40.41**
 extreme, **25.2**
 far, **24.3, 25.2**
 forward looking, **33.4**
 long-wavelength, **24.3, 33.3–33.5, 33.6f**
 medium-wavelength, **24.3, 25.2, 33.3, 33.5, 33.6f**
 near, **24.3, 25.2**
 short-wavelength, **24.3, 33.3, 33.5**
 very long-wavelength, **24.3**
- Infrared radiometry, standards for, **15.11–15.12, 15.12f**
- Inhomogeneous broadening, emission-line, **16.5, 16.6f**
- Inhomogeneous media, **39.22**
- Inhomogeneous reflectors, **39.39**
- Injection-locked lasers, **19.37f, 19.38**
- Input circuits, of hybrid FPAs, **33.18–33.23, 33.19t, 33.20f**
 buffered direct injection, **33.19t, 33.20f, 33.21**
 capacitive transimpedance amplifier, **33.19t, 33.20f, 33.22–33.23**
 chopper-stabilized BDI, **33.19t, 33.20f, 33.21–33.22**
 direct detector integration, **33.18, 33.19t, 33.20f**
 direct injection, **33.18–33.21, 33.19t, 33.20f, 33.21f**
 gate modulation, **33.19t, 33.20f, 33.22**
- Input optics, **38.7**
- Input windows, of image intensifiers, **31.9–31.12, 31.9f, 31.11f, 31.12f**
- Institute for Electrical and Electronic Engineering (IEEE), **36.3**
- Insulators, light bulb, **40.29f**
- Integral color filter arrays (CFAs), **32.32–32.34, 32.33f, 32.34f**
- Integral density, **29.14**
- Integrated lasers, with 45° mirror, **19.39–19.40, 19.39f**
- Integrated transmittance, **35.3**
- Integrating cavities, of nonimaging optics, **39.24–39.27**
 efficiency vs. luminance, **39.26, 39.27f**
 modifying cavity output distribution, **39.27**
 with nonimaging concentrator/collectors, **39.26, 39.27**
 nonuniformities with spherical, **39.24–39.26, 39.24f, 39.25f**
- Integrating spheres (devices), **35.9, 35.11–35.13, 35.11f, 37.9–37.10, 37.10f**
- Integrating-bucket phase shifting, **13.21, 13.21f**
- Intensity:
 defined, **34.9, 40.1**
 luminous, **39.2t**
 nomenclature for, **36.4**
 radiant, **39.2t**
- Interface coupling, **20.14**
- Interference filters, **34.36**
- Interferograms, **13.14–13.18**
 from direct interferometry, **13.17–13.18**
 fixed, **13.14–13.15**
 Fourier analysis of, **13.16–13.17, 13.17f**
 interpolation of, **13.15–13.16**
- Interferometers, **13.7–13.12**
 Brunning distance-measuring, **12.8–12.9, 12.8f**
 common-path, **13.9, 13.11f**
 computer-generated holograms for, **14.4–14.5, 14.4f, 14.5f**
 distance-measuring, **12.7–12.10, 12.7f–12.9f**
 Fabry-Perot, **16.19f**
 Fizeau, **12.14, 13.8–13.9, 13.9f, 13.10f, 13.18, 14.4, 14.5f**
 Haidinger, **12.14**
 heterodyne, **13.22**
 lateral-shearing, **12.14, 13.9–13.12, 13.11f, 13.12f**
 Michelson, **12.5, 12.6, 12.14**
 microinterferometers, **10.13, 10.13f**
 multiple-pass, **13.13**
 multiple-reflection, **13.13**
 nonreacting, **12.7**
 point diffraction, **13.11f**
 radial-shearing, **13.12, 13.13f**
 reversing-shearing, **13.12, 13.13f**
 rotational-shearing, **13.12, 13.13f**
 sensitivity of, **13.13–13.14, 13.14f**
 single-shot f -to- $2f$, **21.6**
 Twyman-Green, **13.7–13.8, 13.7f, 13.8f, 13.18**
 Zernike phase-contrast method applied to, **13.13–13.14, 13.14f**
 (*See also specific interferometers, e.g.: Lateral-shearing interferometers*)
- Interferometric plots, for orthonormal aberrations, **11.36–11.37, 11.37f, 11.38f**
- Interferometry:
 of angles, **12.14**
 direct, **13.17–13.18**

- Interferometry (*Cont.*):
 infrared, 13.25
 of medium distances, 12.6–12.10, 12.7f–12.9f
 phase-shifting, 13.18–13.23, 13.18f–13.20f
 heterodyne interferometer, 13.22
 integrating bucket method, 13.21, 13.21f
 phase errors, 13.22
 phase stepping method of, 13.20, 13.20f
 phase-lock method, 13.23, 13.23f
 simultaneous measurement, 13.22
 two steps plus one method, 13.21, 13.22
 pulse-train, 20.12, 20.12f
 of small distances, 12.5, 12.6
 sub-Nyquist, 13.27
 two-wavelength, 13.25, 13.26
- Interior lighting, 40.55–40.61
 for health-care facilities, 40.58–40.60, 40.60t
 for industry, 40.60–40.61, 40.61f
 for offices, 40.55, 40.56t
 for residences, 40.57, 40.58, 40.59t
 for retail, 40.55–40.57, 40.56t–40.58t
- Interlayer interimage effects (IIEs), 30.19
- Interline transfer (IT) CCD image sensors, 32.28–32.32, 32.29f–32.31f
- Interline-transfer (IT) CCD FPAs, 33.11–33.13, 33.12f
- Internally processed (IP) photocathodes, 31.24
- International Astronomical Union (IAU), 36.3
- International Bureau of Weights and Measures (BIPM), 36.2
- International candle (unit), 37.3
- International Commission on Illumination (CIE), 36.2
- International Committee for Weights and Measures (CIPM), 36.2, 38.6
- International Graphics Exchange Specification (IGES), 40.19
- International Standards Organization (ISO), 4.10, 4.11, 36.2, 40.19
- International System of Units (*see* SI units)
- International Union of Pure and Applied Physics (IUPAP), 36.2
- Intervalence band absorption (IVBA), 19.17
- Interwoven pulse trains, 20.13
- Intrinsic athermalization, 8.7–8.8, 8.7f
- Intrinsic infrared detectors, 33.7, 33.8f
- Intrinsic photoconductors, 25.5, 25.5f
- Intrinsic photodetectors, 24.7, 24.7f
 germanium, 24.70–24.73, 24.70f–24.73f
 indium antimonide photovoltaic, 24.80–24.83, 24.82f, 24.83f
- Intrinsic semiconductor transition, 24.11
- Invariants, 1.11
- Inverse square law, 34.14, 37.8
- Inversion layer, in MOS transistors, 25.11, 25.11f
- Inverted channel substrate planar (ICSP) lasers, 19.20t, 19.23
- Involute reflectors, 39.11–39.12, 39.12f
- Ion bombardment strip lasers, 19.8, 19.9f
- Ion current measurement devices, 34.29
- Irradiance:
 defined, 34.8, 36.5, 37.4t, 37.5, 39.2t
 excitance and emittance vs., 39.2
 image specifications for, 4.7
 spectral, 36.14, 38.1–38.2, 38.11t, 38.13–38.16, 38.13f–38.16f
- Irradiance response units, 34.31
- Isoelectronic dopants, 17.16
- Isoelectronic trap, 17.6, 17.6f
- Isolated pulses, in attosecond optics, 21.4
- Isometric plots, for orthonormal aberrations, 11.36–11.37, 11.37f, 11.38f
- Isotope broadening, 16.6
- Isotropic (term), 36.4, 39.3
- Isotropic point source, 36.4
- Iterated rays, 3.12
- J-aggregates, 30.13, 30.14
- Jaynes-Cummings model, 23.15–23.17, 23.15f
- Jet polishing, 9.6
- Johnson noise (*see* Thermal noise)
- Johnson noise power density, 28.3
- Jones (unit), 24.11, 24.13
- Joule (unit), 34.5–34.6, 37.6
- Judd-Vos modified function, 36.10
- Junction photodiodes, 32.3–32.6, 32.4f, 32.6f
- Kaleidoscope effect, 39.28
- Keck telescope, 11.4
- Kick operator, 23.17
- Kirchhoff's law, 34.25, 35.7, 35.8t
- Knife-edge test (*see* Foucault test)
- Kodachrome film, 29.14, 30.23
- Kodacolor, 29.14
- Kodak Gold film, 30.25
- Kodak Royal Gold film, 30.25

- Kodak Technical Pan Film, **29.19t**
- Kohler illumination, **1.11, 1.11f, 39.23–39.24, 39.23f, 39.34, 39.35f**
- Laboratory sources (of radiation), **15.7–15.13**
 baseline standard for, **15.9, 15.9f, 15.10f, 15.12f**
 blackbody cavity theory, **15.7–15.9, 15.8f**
 working standards for, **15.9–15.13, 15.10f, 15.12f, 15.13f**
- Labsphere, **38.12**
- Lagrangian rays, **3.12**
- Lamb shift, **23.13**
- Lambda Research Corporation, **7.27**
- Lambert (unit), **34.43, 36.7, 36.8t**
- Lambertian approximation (of radiant flux transfer), **34.14–34.18**
 and lambertian sources, **34.14–34.17, 34.15f, 34.16f**
 radiant flux transfer through lambertian reflecting sphere, **34.17–34.18**
- Lambertian surface, **37.8**
- Lambert's cosine law, **37.8, 37.8f**
- Lamps:
 configurations of, **15.20f**
 modeling of, **40.17**
 standards for, **15.11**
 (See also *specific types of lamps, e.g.: Airway beacon lamps*)
- Lamps for Scientific Purposes* (G. M. B. H. Osram), **15.20**
- Land (term), **34.35, 34.35f**
- Lapping step (of optics fabrication), **9.5**
- Large-area detectors, **25.12**
- Laser(s), **12.7, 16.1–16.37**
 about, **16.2–16.3**
 diagram of, **16.2f**
 electromagnetic spectrum involving, **16.2, 16.3, 16.3f**
 and laser gain medium, **16.4–16.19**
 emission linewidth and line broadening of radiating species, **16.4–16.7, 16.6f, 16.7f**
 energy levels and radiation, **16.4**
 gain saturation, **16.10**
 optimization of output coupling from laser cavity, **16.13, 16.14**
 population inversions, **16.8–16.10, 16.12–16.13, 16.13f, 16.14f**
 pumping techniques to produce inversions, **16.14–16.19**
- Laser(s), and laser gain medium (*Cont.*):
 stimulated absorption and emission, **16.7–16.8, 16.8f**
 threshold conditions, **16.10–16.12, 16.11f**
 as light sources, **40.39**
 in measurement, **12.2, 12.6**
 and optical cavities or resonators, **16.19–16.25, 16.19f**
 configurations and cavity stability, **16.23–16.25, 16.24f, 16.25f**
 longitudinal laser modes, **16.20, 16.20f**
 transverse laser modes, **16.21–16.23, 16.21f–16.23f**
 probability flow diagram for, **23.23f**
 quantum theory of (*see* Quantum theory of lasers)
 as radiometric characterization tool, **34.32**
 semiconductor arrays of, **19.26–19.29, 19.28f, 19.28t**
 special laser cavities, **16.25–16.29**
 distributed feedback lasers, **16.29**
 mode-locking, **16.27–16.29, 16.28f**
 Q-switching, **16.26–16.27, 16.27f**
 ring lasers, **16.29**
 unstable resonators, **16.25–16.26, 16.26f**
 two-dimensional high-power arrays of, **19.29–19.30, 19.29t, 19.30f**
 as two-level system, **20.23–20.24, 20.25t**
 types of, **16.29–16.37, 16.31f, 16.32f, 16.34f, 16.35f, 16.37f**
 (See also *related topics, e.g.: Laser stabilization*)
- Laser beam expanders, athermal, **8.13–8.14**
- Laser cavities, **16.25–16.29**
 distributed feedback lasers, **16.29**
 mode-locking, **16.27–16.29, 16.28f**
 Q-switching, **16.26–16.27, 16.27f**
 ring lasers, **16.29**
 unstable resonators, **16.25–16.26, 16.26f**
- Laser detection and ranging (LADAR), **31.28, 31.30**
- Laser field, spectral properties of the, **23.28–23.31**
- Laser gain media, **16.4–16.19**
 dielectric solid-state, **16.32–16.34**
 gaseous, **16.30–16.31, 16.31f**
 liquid, **16.31–16.32, 16.32f**
 properties associated with, **16.4–16.19**
 emission linewidth and line broadening of radiating species, **16.4–16.7, 16.6f, 16.7f**
 energy levels and radiation, **16.4**

- Laser gain media, properties associated with (*Cont.*):
 gain saturation, 16.10
 optimization of output coupling from laser cavity, 16.13, 16.14
 population inversions, 16.8–16.10, 16.12–16.13, 16.13f, 16.14f
 pumping techniques to produce inversions, 16.14–16.19
 stimulated absorption and emission, 16.7–16.8, 16.8f
 threshold conditions with mirrors, 16.10–16.12, 16.11f
 in vacuum, 16.36–16.37, 16.37f
- Laser linewidth, 23.34–23.35
- Laser locking, frequency discriminators for, 22.12–22.14
- Laser master equation, 23.19–23.20, 23.19f
- Laser phase-transition analogy, 23.35–23.40, 23.37t, 23.38f, 23.39f
- Laser photon statistics, 23.22–23.26, 23.23f, 23.25f
- Laser power measurement, 34.32
- Laser resonators, 16.20f, 16.23–16.26, 16.23f–16.26f
- Laser scanning, 40.54
- Laser scribing, 17.24–17.25
- Laser stabilization, 22.1–22.24
 about, 22.1
 Allan Deviation, 22.2–22.3
 and frequency discriminators for laser locking, 22.12–22.14
 frequency vs. time drift, 22.2
 future directions for, 22.23–22.24
 and optical cavity-based frequency discriminators, 22.14–22.16, 22.17f
 quantifying frequency stability, 22.2
 and quantum resonance absorption, 22.16, 22.17
 representative/example designs of, 22.20–22.23, 22.22f
 and servos, 22.5–22.12
 Bode representation of servos, 22.5–22.6, 22.6f
 closed-loop performance, 22.8
 closed-loop stability issues, 22.8–22.12, 22.9f–22.1f, 22.10–22.12
 measurement noise, 22.7–22.8
 phase and amplitude responses, 22.6–22.7, 22.6f, 22.7f
- Laser stabilization (*Cont.*):
 spectral noise density, 22.3–22.5
 and transducers, 22.17–22.20
- Laser stripe structures, 19.8, 19.9f
- Lasers and Optronics Buying Guide*, 15.14
- Lasing, without inversion, 23.40–23.42, 23.41f
- Latent image (LI), 30.7
- Latent-image speck, 29.5
- Lateral aberration curves (*see* Transverse ray plots)
- Lateral antiblooming, 32.9, 32.10f
- Lateral color, 1.14, 2.5, 2.5f
- Lateral *pin* photodetectors, 25.15–25.16, 25.15f
- Lateral-collection photodetectors, 26.4f
- Lateral-shearing interferometers, 12.14, 13.9–13.12, 13.11f, 13.12f
- Lathe assembly technique, 6.7–6.8, 6.8f
- Lead salt lasers, 19.7–19.8
- Lead selenide (PbSe) detectors, 24.76f, 24.78, 24.79, 24.79f–24.82f
- Lead sulfide (PbS) photoconductors, 24.73–24.74, 24.74f–24.77f, 24.74t
- Lead tin telluride (PbSnTe) photovoltaic detectors, 24.92, 24.93f
- Lead-in wires, light bulb, 40.29, 40.29f
- Leakage current, 24.19–24.20, 32.10–32.12, 32.11f
- Leaky-mode arrays, 19.29
- Least-squares method, 3.17–3.19
- Legacy films, 30.23–30.25
- Legal traceability (of calibration), 34.21
- Legendre polynomials, 11.5, 11.30
- Length measurements, 12.2–12.10
 interferometric measurement, 12.5–12.10, 12.7f–12.9f
 stadia and range finders, 12.2–12.4, 12.3f, 12.4f
 standards for, 12.2
 time-based and optical radar, 12.4, 12.5, 12.6f
- Lens(es), 39.32–39.37, 39.33f
 assembly adjustment of, 6.8–6.11, 6.10f
 coating specifications for, 4.10
 component specifications for, 4.2, 4.3
 and concentrators, 39.16–39.18, 39.18f
 data entry for, 3.2–3.8
 humidity specifications for, 4.10
 for II electronic imaging, 31.5–31.6
 image specifications for, 4.3, 4.6–4.8, 4.8f
 literature on nonimaging, 39.32, 39.33

- Lens(es) (*Cont.*):
- multicomponent assemblies of (*see* Multicomponent lens assemblies)
 - optical parameters for, 4.9
 - perfect, 4.4
 - reflector/lens-array combinations, 39.34–39.37, 39.36f, 39.37f
 - single-lens arrays, 39.33–39.34, 39.34f
 - tandem-lens arrays, 39.34, 39.35f–39.37f
 - thermal defocus of, 8.4, 8.5f
 - thermal focus shift of, 8.2–8.4, 8.3t, 8.4t (*See also specific lenses, e.g.: Air-spaced doublet lens*)
- Lens design:
- aberration curves in, 2.1–2.6
 - software for, 40.20
 - and tolerancing calculations, 5.10
- Lens setup routine (in optical software), 3.7
- Lens-coupled II SSAs, 31.22–31.23, 31.22f
- Levels (tools), 12.13–12.14, 12.13f, 12.14f
- Lever mechanism mountings, 6.19, 6.19f
- Lifetime classification, of photodetector materials, 26.5
- Light:
- spectrum of, 25.2
 - speed of, 12.2
- Light amplification by stimulated emission of radiation (*see* Laser(s))
- Light detection and ranging (LIDAR), 25.12
- Light distribution, 40.5
- Light extraction, 17.6–17.8, 17.7t
- Light loss factor (LLF), 40.17
- Light piping, 30.6–30.7, 30.6f
- Light pollution, 40.43, 40.62
- Light quanta, 23.6–23.9, 23.8f
- Light scattering, 30.5–30.7, 30.6f
- Light shelves, 40.48, 40.50f
- Light sources, 40.24–40.41, 40.28f
- applications for, 40.26t
 - carbon arc sources, 40.40
 - characteristics of, 40.25t
 - daylight, 40.40–40.41, 40.40f
 - electrodeless lamps, 40.36–40.37
 - electroluminescent sources, 40.37–40.39, 40.38f, 40.38t, 40.39f
 - fluorescent lamps, 40.30–40.33, 40.31f, 40.32f
 - glow lamps, 40.39
 - high-intensity discharge lamps, 40.33–40.36, 40.33f–40.35f
- Light sources (*Cont.*):
- incandescent sources, 40.25, 40.27–40.30, 40.28f, 40.29f
 - low-pressure sodium lamps, 40.33–40.36, 40.33f–40.35f
 - nuclear sources, 40.39
 - pure Xe arc lamps, 40.39
 - short arc sources, 40.39
 - types of, 40.27f
- Light stability, 30.10
- Light stabilization, 30.12–30.13
- Light trespass, 40.43, 40.62
- Light-absorbing dye, 30.7
- Light bulbs, 40.27–40.30
- base of, 40.29, 40.29f
 - CFL/fluorescent/minature, 40.28f
 - elements of, 40.29f
 - shapes of, 15.20f, 15.30f
 - sizes of, 15.30f
 - types of, 40.27t
- Light-emitting diodes (LEDs), 17.1–17.34
- conversion of, luminous intensity to radiant intensity, 36.13, 36.13f
 - device structures of, 17.8–17.15
 - diffused homojunctions, 17.9–17.10, 17.10f, 17.11f
 - double heterojunctions, 17.13, 17.13f–17.15f
 - grown homojunctions, 17.8, 17.9, 17.9f
 - single heterojunctions, 17.12, 17.12f
- epitaxial technology for, 17.21–17.23
- lamps with, 40.37–40.39
- applications for, 40.26t
 - characteristics of, 40.25t
 - materials/emitted colors of, 40.38t
 - photonic crystal, 40.39f
 - structure of, 40.38f
- and LED-based products, 17.29–17.32, 17.29f–17.31f
- and light extraction, 17.6–17.8, 17.7t
- and light-generation processes, 17.2–17.6, 17.3f–17.6f
- material systems for, 17.15–17.19
- AlInGaP system, 17.18, 17.19f
 - $\text{Al}_x\text{Ga}_{1-x}\text{As}$ system, 17.17, 17.17t
 - blue LED technology, 17.18, 17.19
 - $\text{GaAs}_{1-x}\text{P}_x$ system, 17.15–17.17, 17.16t
- octocouplers in, 17.32–17.34, 17.32f
- production levels for, 17.2
- quality/reliability of, 17.25–17.28

- Light-emitting diodes (LEDs) (*Cont.*):
 substrate technology for, 17.20–17.21, 17.20*t*
 wafer processing for, 17.23–17.25
 (*See also specific light-emitting diodes, e.g.:*
 High-brightness visible LEDs)
- Lighting, 40.1–40.71
 about, 40.1–40.3
 exterior, 40.61–40.62, 40.63*t*
 functions of, 40.12–40.14, 40.13*f*–40.16*f*
 insufficient, 40.9
 interior, 40.55–40.61
 health-care facility lighting, 40.58–40.60,
 40.60*t*
 industrial lighting, 40.60–40.61, 40.61*f*
 office lighting, 40.55, 40.56*t*
 residential lighting, 40.57, 40.58, 40.59*t*
 retail lighting, 40.55–40.57, 40.56*t*–40.58*t*
 perception of, 40.4–40.6, 40.4*f*
 for transportation, 40.63–40.71
 roadway lighting, 40.67, 40.69–40.71,
 40.70*t*, 40.71*t*
 vehicular lighting, 40.63–40.67, 40.64*f*,
 40.65*t*, 40.66*f*, 40.66*t*, 40.68*t*, 40.69*f*
 vision biology, 40.3–40.6
 (*See also Luminaires*)
- Lighting design, 40.6–40.23
 and color, 40.7–40.9
 and context, 40.6
 and functions of lighting, 40.12–40.14,
 40.13*f*–40.16*f*
 geometries in, 40.13*f*, 40.14–40.15, 40.15*f*,
 40.16*f*
 goals of, 40.6
 and illuminance, 40.7, 40.7*t*
 and properties of objects and impact, 40.16
 system layout and simulation in, 40.16–
 40.23, 40.18*f*, 40.22*f*–40.24*f*
 and visual discomfort, 40.9–40.12, 40.11*t*
- Lighting Design and Application (IESNA),
 39.8
- Lighting design software, 40.21
- Lighting geometries, 40.13*f*, 40.14–40.15,
 40.15*f*, 40.16*f*
- Lighting Handbook (IESNA), 36.7, 40.17
- Lighting measurement, 40.51–40.54
 goniometers, 40.52–40.53, 40.53*t*, 40.54*f*
 illuminance meters, 40.51, 40.52*f*
 luminance meters, 40.52, 40.52*f*
 reflectometers, 40.52
 surface measurement systems, 40.53, 40.54
- Lighting system layout and simulation,
 40.16–40.23, 40.18*f*
 computer graphics software for, 40.21–40.23,
 40.22*f*, 40.23*f*
 IGES standard for, 40.19
 optical analysis and design software for,
 40.20
 optical design and analysis software for,
 40.20–40.21
 software tools for, 40.18–40.23, 40.22*f*–40.24*f*
 source modeling software for, 40.19–40.20
 STEP standard for, 40.19
- Lightness (term), 40.4
- Light-output degradation, 17.26–17.28, 17.28*f*
- Lightpipes, 39.27–39.32
 angular uniformity of, 39.31
 applications for, 39.32
 length of, 39.30
 periodic distributions of, 39.30
 shapes of, 39.27–39.30, 39.28*f*–39.30*f*
 solid vs. hollow, 39.30–39.31
 tapered, 39.12–39.13, 39.13*f*, 39.31–39.32,
 39.31*f*
- LightTools (optical software), 7.26
- Light-trap silicon photodiodes, 34.30
- Limiting resolution, of EBSSAs, 31.27
- Linear image sensor arrays, 32.2, 32.21–32.24,
 32.22*f*, 32.23*f*
- Linear lasers, 20.18–20.19, 20.19*f*
- Linear variable differential transformers
 (LVDTs), 10.12
- Linearity, of photoemissive detectors, 24.41
- Line-scanned imaging systems, 31.29
- Linewidth, spectral density and, 22.4–22.5
- Lippmann emulsions, 29.4
- Liquid Encapsulated Czochralski (LEC)
 technique, 17.20, 17.21
- Liquid laser gain media, 16.31–16.32, 16.32*f*
- Liquid phase epitaxy (LPE), 17.21–17.22, 19.6,
 19.19
- Lit-appearance modeling, 40.21–40.23, 40.23*f*,
 40.24*f*
- Lithographic etching, 18.3–18.4, 18.3*f*
- Lithographic projection lens, 6.10*f*
- LLL sensitivity, 31.19, 31.20
- Local area networks, 19.3
- Localized avalanche breakdown, 17.28
- Lock-in amplifiers, 27.14, 27.14*f*, 27.15, 38.10
- Long wavelength lasers, 19.8

- Long wavelength QW lasers, 19.17–19.18, 19.17f
- Longitudinal laser modes, 16.19f, 16.20, 16.20f
- Long-wavelength infrared (LWIR), 24.3, 33.3–33.5, 33.6f
- Lorentzian line shape, 16.6f
- Louvers, 40.41, 40.45f, 40.46
- Low-beam headlamps, 40.64–40.67, 40.64f, 40.65t, 40.66f, 40.66t
- Low-intensity carbon arc lamps, 15.21–15.24, 15.24f, 15.28f
- Low-intensity reciprocity failure, of photographic films, 29.12
- Low-level-light (LLL) TV imaging, 31.1–31.4
- Low-pressure enclosed arcs, 15.35–15.47
 - black-light fluorescent lamps, 15.35, 15.36t
 - electrodeless discharge lamps, 15.36, 15.44
 - germicidal lamps, 15.35
 - hollow cathode lamps, 15.35, 15.37t–15.43t, 15.44f
 - Pluecker spectrum tubes, 15.47, 15.47f, 15.47t
 - spectral lamps, 15.44, 15.45, 15.45f, 15.46f, 15.46t
 - Sterilamps, 15.35, 15.36f
- Low-pressure lamps, 15.36f
- Low-pressure sodium (LPS) lamps, 40.33–40.36
 - applications for, 40.26t
 - characteristics of, 40.25t
 - construction of, 40.34f
 - emission spectrum of, 40.35f
- Low-speed photographic films, 30.18–30.20
- Low-temperature bolometers, 24.31–24.32, 24.32f, 24.33f, 28.5
- Low-temperature (LT) grown photoconductors, 26.23
- Lucalox lamps, 15.30, 15.31f
- Lumen (unit), 36.6, 37.6, 39.2t
- Lumen lighting simulation, 40.17
- Luminaires, 40.24–40.50
 - applications for, 40.26t
 - calculation of needed, 40.16–40.17
 - characteristics of, 40.25t
 - classification system for, 40.43–40.45, 40.43t, 40.44f
 - defined, 40.1
 - design of, 40.41–40.50
 - conics shapes and intensity distribution, 40.43t
 - etendue and source coupling, 40.41–40.42
 - design of (*Cont.*):
 - luminaire classification system, 40.43–40.45, 40.43t, 40.44f
 - methods, 40.42–40.43, 40.43t
 - light sources for, 40.24–40.41, 40.27f, 40.28f
 - carbon arc sources, 40.40
 - daylight, 40.40–40.41, 40.40f
 - electrodeless lamps, 40.36–40.37
 - electroluminescent sources, 40.37–40.39, 40.38f, 40.38t, 40.39f
 - fluorescent lamps, 40.30–40.33, 40.31f, 40.32f
 - glow lamps, 40.39
 - high intensity discharge lamps, 40.33–40.36, 40.33f–40.35f
 - incandescent sources, 40.25, 40.27–40.30, 40.28f, 40.29f
 - low-pressure sodium lamps, 40.33–40.36, 40.33f–40.35f
 - nuclear sources, 40.39
 - pure Xe arc lamps, 40.39
 - short arc sources, 40.39
 - optics of, 40.45–40.50
 - for artificial sources, 40.45–40.47, 40.45f, 40.46f
 - backlighting, 40.47, 40.47f, 40.48f
 - for daylight sources, 40.47–40.50, 40.49f–40.51f
- Luminance, 37.4t, 37.5, 37.5f, 39.2t
 - calibration of, 34.42
 - defined, 34.11, 34.40, 36.7, 40.1
 - and illuminance, 37.9
 - of integrating cavities, 39.26
 - in nonimaging optics, 39.2t, 39.3
 - uniformity of, 40.7
 - units of, 34.43
- Luminance contrast, 40.6, 40.10
- Luminance meters, 40.52, 40.52f
- Luminance ratio, 40.7
- Luminous efficacy, 18.5
- Luminous efficiency, 17.15
- Luminous energy, 37.4t, 37.6
- Luminous exitance, 37.4t, 37.5, 37.5f
- Luminous exposure, 37.4t, 37.6
- Luminous flux, 15.11, 34.10–34.11, 34.39, 34.42, 37.4, 37.4t, 37.6, 38.2
- Luminous flux density, 34.11
- Luminous intensity, 15.11, 34.11, 34.40, 37.4, 37.4t, 39.2t
- Lux (unit), 34.43, 36.7, 36.7t

- Luxmeters, 40.52*f*
 Lux-second, 29.6
 Lyot stops, 7.8–7.10, 7.8*f*–7.11*f*
- Macrofocal reflectors, 39.11
 Magnesium, as *p*-type impurity, 17.23
 Magnetic shielding, 27.10
 Magnetorheological finishing (MRF), 9.5
 Magnification, 1.4
 Magnifiers, first-order layout for, 1.8
 Main event loop, 3.7
 Maksutov sphere, 14.7, 14.8*f*
 Maksutov test, 14.7–14.9, 14.8*f*
 Mandel Q_M parameter, 23.26
 Markov approximation, 23.21
 Martin Black coating, 7.14–7.17, 7.14*f*–7.16*f*,
 7.23, 7.25*f*
 Masers, micro- (*see* Micromasers)
 Master-oscillator power amplifier (MOPA),
 19.41
 Materials:
 formation of, for optics, 9.3–9.4
 specifications for, 4.9
 tolerancing and properties of, 5.9
 Matte, 30.3
 Maximized D star, 24.11
 Maximum spectral luminous efficiency
 (of radiation), 37.2
 Maxwell's principle, 30.16
 Mazers (microwave amplification by z-motion-
 induced emission of radiation), 23.45
 Mean-square-spot size (MSS), 3.15
 Mean-time-to-failure (MTTF), 25.13, 25.14
 Measurement(s), 35.8–35.16
 of absorptance, 35.10
 of emittance, 35.14–35.16, 35.15*f*
 of lighting, 40.51–40.54, 40.52*f*, 40.53*t*, 40.54*f*
 of reflectance, 35.10–35.13, 35.10*f*–35.12*f*,
 35.14*t*
 terminology, 35.2–35.3, 35.2*f*
 of transmittance, 35.8–35.10, 35.9*f*
 (*See also specific types of measurement, e.g.:*
 Nonlinearity measurement)
 Measurement noise, 22.7–22.8
 Measurements Assurance Program (MAP), 35.13
 Mechanical athermalization, 8.8–8.12
 active, 8.11, 8.11*f*
 by image processing, 8.12
 part active, part passive, 8.11–8.12, 8.12*f*
 passive, 8.8–8.10, 8.8*f*–8.10*f*
- Mechanical scribing, 17.24
 Mechanical specifications:
 for lenses, 4.9
 optical vs., 4.2
 Mechanical tolerances, 5.2
 Mechanical vibrations, 27.5, 27.6*f*
 Mechanically clamped mountings, 6.12, 6.13*f*
 Medium-wavelength infrared (MWIR) radia-
 tion, 24.3, 25.2, 33.3, 33.5, 33.6*f*
 Memory colors, 30.21
 Mercury arc lamps, 15.29, 15.29*f*–15.31*f*, 15.34*f*
 Mercury cadmium telluride (HgCdTe)
 detectors, 24.86–24.92, 24.87*f*
 infrared, 33.7, 33.8*f*
 photoconductors, 24.86*f*, 24.88–24.90,
 24.89*f*–24.91*f*
 photodetectors, 25.10, 25.10*t*
 photovoltaic, 24.86*f*, 24.88*f*, 24.90–24.92,
 24.91*f*, 24.92*f*
 Mercury-doped germanium detectors, 24.84*f*,
 24.92–24.95, 24.93*f*–24.95*f*
 Mercury-free fluorescent lamps, 40.31
 Mercury-halide fluorescent lamps, 40.33*f*
 Mercury-vapor fluorescent lamps, 40.30,
 40.31, 40.33*f*
 Mercury-xenon lamps, 15.34*f*
 Meridional rays, 3.3
 Merit function, in optical design software, 3.17
 Mesa etching, 18.3–18.4, 18.3*f*
 Mesa photodiodes, 25.14, 25.15*f*
 Mesopic vision, 34.37, 36.9, 37.2
 Metal insulator semiconductor (MIS) photo-
 gate FPAs, 33.10–33.11, 33.11*f*, 33.12*f*
 Metal insulator semiconductors (MISs), 33.4
 Metallic mirrors, mounting of, 6.19–6.20, 6.20*f*
 Metalorganic chemical vapor deposition
 (MOCVD), 19.6–19.7, 19.20*t*, 19.23
 Metalorganic vapor phase epitaxy, 17.21, 17.22
 Metal-oxide semiconductors (MOSs):
 linear arrays of, 32.21–32.24, 32.22*f*, 32.23*f*
 readouts from, 32.20–32.21
 Metal-oxide-semiconductor (MOS) area array
 image sensors, 32.25–32.26, 32.26*f*
 Metal-oxide-semiconductor (MOS) capacitors,
 32.4*f*, 32.7–32.8
 Metal-oxide-semiconductor (MOS) detectors,
 25.11, 25.11*f*
 Metal-semiconductor-metal (MSM)
 photodetectors, 26.3, 26.4*f*
 Meter (unit), 12.2

- Meter, 1875 Treaty of the, **34.20**, **36.2**
 Meter candle (unit), **34.43**, **36.7**
 Metrology, optical, **12.1–12.25**
 angle measurements, **12.10–12.17**
 autocollimeters, **12.11–12.12**, **12.11f**, **12.12f**
 interferometric methods, **12.14**
 levels (tools), **12.13–12.14**, **12.13f**, **12.14f**
 mechanical methods, **12.10–12.11**, **12.11f**
 in prisms, **12.14–12.16**, **12.15f–12.17f**
 theodolites, **12.13**
 curvature measurements, **12.17–12.25**
 mechanical methods, **12.17–12.19**, **12.18f**,
 12.19f, **12.19t**
 optical methods, **12.19–12.21**, **12.20f**,
 12.20t, **12.21f**
 of diamond-turned optics, **10.12–10.13**,
 10.12f, **10.13f**
 focal length measurements, **12.21–12.25**,
 12.22f–12.24f
 length measurements, **12.2–12.10**
 interferometers, **12.5–12.10**, **12.7f–12.9f**
 stadia and range finders, **12.2–12.4**, **12.3f**,
 12.4f
 time-based and optical radar, **12.4**, **12.5**,
 12.6f
 straightness measurements, **12.10**
 terminology, **12.2**
 Michelson interferometers, **12.5**, **12.6**, **12.14**
 Microbolometer FPAs, **33.13–33.14**
 Microchannel plate tubes (MCPTs), **24.32**,
 24.33f, **24.40**
 Microchannel plates (MCPs), **31.1**, **31.9**, **31.9f**,
 31.12–31.14, **31.12f**, **31.13f**, **31.13t**
 Microchannel-plate image intensifiers
 (MCP IIs), **31.7**
 high-voltage power supply for, **31.9**, **31.10f**
 proximity-focused, **31.9**, **31.9f**, **31.16–31.18**,
 31.17t, **31.18f**, **31.19f**
 Microdensitometers, **29.6**, **29.15–29.16**, **29.15f**
 Microinterferometers, **10.13**, **10.13f**
 Microlens arrays, **32.30**, **32.31**, **32.31f**
 Microlenses, **12.24**
 Micromaser master equation, **23.20–23.22**
 Micromasers, **23.26–23.27**, **23.27f**, **23.45**
 Microminiature lamps, **15.53**
 Microplasma formation, **17.28**
 Microscopes:
 first-order layout for, **1.8**
 Nomarski, **10.11**, **10.11f**
 traveling, **12.20**, **12.21**, **12.21f**
 Microwave powered lamps (*see* Electrodeless
 sulfur lamps)
 Military Sensing Information Analysis Center
 (SENSIAC), **15.6**
 Millilambert (unit), **36.7**
 Milliphot (unit), **36.7**, **36.7t**
 Miniature lamps, **15.53**, **40.26t**, **40.28f**
 Minimum resolvable temperature (MRT)
 (of infrared detector arrays), **33.2**,
 33.27–33.28, **33.27f**
 Minkowitz, S., **12.7**
 Minkowitz distance-measuring interferometers,
 12.7, **12.8**, **12.8f**
 Minority-carrier recombination, **17.2**
 Mirror scatter relationships (stray light), **7.18**
 Mirrored tiling, **39.28**, **39.30**, **39.30f**
 Mirrors:
 on amplifiers, **16.3**
 and concentrators, **39.17**
 mounting of [*see* Mounting (of optical
 components)]
 nonabsorbing, **19.23–19.24**
 threshold conditions with, **16.10–16.12**,
 16.11f
 Mixing rods (*see* Lightpipes)
 Modal dispersion, **17.34**
 Mode-locked lasers, **16.27–16.29**, **16.28f**, **20.7**
 Mode-stabilized lasers, **19.18–19.19**, **19.19f**
 Modulation, as laser stabilization technique,
 22.13–22.14
 Modulation noise, **24.11**
 Modulation transfer function (MTF), **4.1**, **4.2**,
 4.4, **4.5f**, **4.8f**, **29.18**, **29.18t**, **29.19f**, **30.5**,
 31.26–31.27, **33.2**, **33.7**
 Modulators, electro-optic, **22.14**, **22.20**
 MOFSET bucket brigades, **33.17**
 Moiré deflectometry, **12.23**, **12.24f**
 Moiré tests, of spherical aberrations, **1**
 3.26–13.27
 Molecular beam epitaxy (MBE), **17.21–17.23**,
 19.6, **19.7**, **25.16**
 Moments normalization (technique),
 36.15–36.16, **36.15t**, **36.16f**
 Monochromators, **35.9**, **38.7–38.8**, **38.14–38.16**,
 38.14f–38.16f
 Monolithic built-up mirror substrate, **6.18**, **6.18f**
 Monolithic FPAs, **33.10–33.14**, **33.11f**, **33.12f**
 Monolithic LED displays, **17.30**
 Monolithic silicon bolometers, **28.10–28.11**,
 28.10f

- Monolithic two-dimensional (2D) laser arrays, **19.39**
- Monte Carlo simulations, **39.7**
- Mordants, **30.4**
- Mosaic II SSA cameras, **31.29**
- Mounting (of optical components), **6.1–6.24**, **6.17f**, **6.18f**
 and contact stresses, **6.21**
 of domes, **6.11**, **6.12f**
 hard, **6.2–6.4**
 of individual rotationally symmetric optics, **6.2–6.5**
 lever-mechanism, **6.19**, **6.19f**
 of moderate-sized mirrors, **6.17–6.20**, **6.20f**
 in multicomponent lens assemblies, **6.5–6.11**
 drop-in assembly, **6.6**, **6.6f**
 lathe assembly, **6.7–6.8**, **6.8f**
 lens adjustments at assembly, **6.8–6.11**, **6.10f**
 “poker chip” assembly, **6.8**, **6.9f**
 tightly toleranced assembly, **6.7**, **6.7f**
 of small mirrors/prisms, **6.11–6.17**
 bonded mountings, **6.13–6.15**, **6.15f**
 elastomeric mountings for mirrors, **6.12**
 flexure mountings for small mirrors/prisms, **6.15–6.17**, **6.16f**
 mechanically clamped mountings, **6.12**, **6.13f**
 spring-loaded mountings, **6.13**, **6.14f**
 soft, **6.4–6.5**
 and temperature effects, **6.21–6.24**, **6.22f**
 of windows, **6.11**, **6.11f**
- Multichannel detectors, **38.9**
- Multicomponent lens assemblies, **6.5–6.11**
 drop-in, **6.6**, **6.6f**
 lathe, **6.7–6.8**, **6.8f**
 and lens adjustments, **6.8–6.11**, **6.10f**
 “poker chip,” **6.8**, **6.9f**
 tightly toleranced, **6.7**, **6.7f**
- Multiphase reflectors, **39.39**
- Multiple quantum well (MQW) LEDs, **18.1**, **18.2f**
- Multiple surface concentrators, **39.16–39.17**, **39.17f**
- Multiple wafer transducers, **22.18–22.19**
- Multiple-pass interferometers, **13.13**
- Multiple-reflection interferometers, **13.13**
- Multiple-segment LED displays, **17.31**
- Multiplier photodiodes, **24.11**
- Multiplier phototubes (*see* Photomultiplier tubes (PMTs))
- Multiplier tubes, **24.6**
- Multipop (film exposure technique), **30.22**
- Multiquantum well (MQW) lasers, **19.14**, **19.15**, **19.15f**, **19.24**, **19.25t**
- Multiquantum-well buried heterostructure (MQW BH) lasers, **19.25t**
- Multispeed choppers, **15.14**
- Multivapor arcs, **15.29**, **15.30f**, **15.31f**
- Murty’s lateral shear interferometer, **12.14**, **13.12f**
- Mylar film, **29.4**
- National Physical Laboratory (NPL), **15.4**
- National Search Engine for Standards, **4.11**
- Natural broadening, emission-line, **16.5**, **16.6**
- Natural color (NC) film, **30.27**
- Natural linewidth (of transition), **16.5**
- Near infrared (NIR) radiation, **24.3**, **25.2**
- Negative color photographic films, **30.25–30.28**, **30.27t**
- Neodymium (Nd) glass lasers, **16.32–16.33**
- Neodymium-yttrium vanadate (Nd:VO₄) lasers, **16.33**
- Neodymium-yttrium-aluminum-garnet (Nd:YAG) lasers, **16.32–16.33**
- Neodymium-yttrium-lithide-fluoride (Nd:YLF) lasers, **16.33**
- Neon signs, **40.39**
- Nernst glower, **15.14**, **15.15**, **15.17**, **15.17f**, **15.19**, **15.19f**
- Nesonian illumination (*see* Abbe illumination system)
- Neutral density filters, **40.52**
- Nit (unit), **34.43**, **36.7**, **36.8t**
- Nitride LEDs, **17.19**, **18.3f**
- Nitrogen doping, **17.16**
- Nitrogen-doped GaAsP, **17.16**, **17.21–17.22**
- Nodal slide bench, **12.22**, **12.22f**
- Noise, **27.3–27.6**
 1/*f*, **27.4**
 atomic, **23.34–23.35**
 in CCDs, **32.20**, **32.32**
 excess, **24.11**
 generation, **24.11**
 generation-recombination, **24.11**
 ground loop, **27.5**, **27.6f**
 inductive pickup, **27.5**, **27.6f**
 measurement, **22.7–22.8**
 modulation, **24.11**
 nonessential, **27.4–27.6**, **27.5f**, **27.6f**

- Noise (*Cont.*):
 pattern, 32.12
 in photodetectors, 24.19–24.20, 24.20f
 in photoemissive detectors, 24.39, 24.40, 24.41f
 resistive coupling, 27.5, 27.6f
 RMS, 24.12
 shot, 24.12, 27.3, 27.3f, 32.12
 and signal detection, 27.1
 spatial, 33.26–33.27, 33.26f
 stray capacitance, 27.5, 27.6f
 temperature, 24.12
 thermal, 24.13, 27.4, 32.20
- Noise equivalent irradiance (NEI), 24.11
- Noise equivalent power (NEP), 24.10, 24.12, 24.14, 25.12, 28.2, 38.9, 38.10, 38.10t
- Noise equivalent quanta (NEQ), 29.1, 29.23
- Noise equivalent temperature difference (NETD), 28.8–28.9, 33.2, 33.24–33.27, 33.25f, 33.26f
- Noise spectral density (NSD), 33.2
- Noise spectrum, 24.11
- Nomarski microscope, 10.11, 10.11f
- Nonabsorbing mirrors (NAMs), 19.23–19.24
- Nonblackbody radiation source, 15.10
- Nonchromogenic film, 29.14
- Noncircular pupils, 11.4, 11.37, 11.39
- Non-diffraction-limited optics, 8.6
- Nonequilibrium errors, 34.25
- Nonessential noise, 27.4–27.6, 27.5f, 27.6f
- Nonessential ray tracing, 40.20–40.21
- Nonideal aperture, 34.35–34.36, 34.35f
- Nonimaging concentrators (*see* Concentrators, nonimaging)
- Nonimaging optics, 39.1–39.41
 about, 39.1–39.2
 aspheric lenses in, 39.8, 39.9, 39.9f
 calculations for, 39.2–39.6
 clipped Lambertian distribution, 39.3–39.4
 concentration, 39.5, 39.6
 dilution, 39.6
 etendue, 39.2, 39.3, 39.4f
 Hottel strings, 39.4, 39.4f
 Lambertian, 39.3
 luminance, 39.3
 projected solid angle, 39.5, 39.5f
 solid angle, 39.5, 39.5f
 concentration of, 39.12–39.22
 2D vs. 3D geometries, 39.20–39.21, 39.20f, 39.21f
 calculation, 39.5, 39.6
- Nonimaging optics, concentration of (*Cont.*):
 compound elliptical collectors, 39.14, 39.15f
 compound hyperbolic collectors, 39.15, 39.15f, 39.16f
 compound parabolic collectors, 39.13–39.14, 39.13f, 39.14f
 dielectric compound parabolic collectors, 39.15, 39.16, 39.16f
 edge rays, 39.22
 geometrical vector flux, 39.21–39.22
 inhomogeneous media, 39.22
 multiple surface concentrators, 39.16–39.17, 39.17f
 restricted exit angle concentrators with lenses, 39.18, 39.18f
 tapered lightpipes, 39.12–39.13, 39.13f
 θ_1/θ_2 concentrators, 39.18–39.20, 39.19f
 conic reflectors in, 39.11, 39.11f
 Fresnel lenses in, 39.9–39.10, 39.10f
 involute reflectors in, 39.11–39.12, 39.12f
 macrofocal reflectors in, 39.11
 software modeling of, 39.6–39.8
 spherical lenses in, 39.8, 39.9f
 terminology, 39.2, 39.2t
 uniform illumination of, 39.22–39.41
 classic projection system uniformity, 39.23–39.24, 39.23f
 faceted structures, 39.39–39.41, 39.39f, 39.40f
 integrating cavities, 39.24–39.27, 39.24f, 39.25f, 39.27f
 lens arrays, 39.32–39.37, 39.33f–39.37f
 lightpipes, 39.27–39.32, 39.28f–39.31f
 tailored reflectors, 39.37–39.39, 39.38f
- Nonimaging software modeling, 39.6–39.8
- Noninterferometric optical testing, 13.1–13.7
 Foucault test, 13.2–13.3, 13.2f, 13.3f
 Hartmann test, 13.4–13.6, 13.5f
 Hartmann-Shack test, 13.6–13.7, 13.6f
 Ronchi test, 13.3–13.4, 13.3f, 13.4f
- Nonlinearity correction factor, 34.33
- Nonlinearity measurement, 34.34–34.35
- Nonneutral (color) density, of photographic films, 29.7–29.8
- Nonradiative recombination, 17.2
- Nonreacting interferometers, 12.7
- Nonsequential ray tracing, 39.6
- Nonsequential surfaces data, 3.6–3.7
- Nonuniform rational B-splines (NURBS), 39.6, 40.41

- Nonuniformity, of radiation distribution, 34.25, 34.28, 34.35
- Normal equations, 3.18
- Normalized detectivity, 38.9
- Normalized detector irradiance (NDI), 7.22
- Notch filters, PID controller vs., 22.10–22.11, 22.10f, 22.11f
- np* silicon photodiodes, 34.30
- Nuclear light sources, 40.39
- Numeric displays, LED, 17.30–17.31, 17.30f, 17.31f
- Numerical aperture (NA), 34.20, 39.1
- Numerically controlled machines (CNCs), 12.11
- Nutting's law, 29.6
- Nyquist condition, 13.27
- Nyquist noise (*see* Thermal noise)
- Object counting, reflexive sensors for, 17.34
- Objective tone reproduction, 29.16–29.17, 29.17f
- Octocouplers, 17.32–17.34, 17.32f
- Off-axis angles, 7.23
- Off-axis chromatic aberrations, 2.2–2.4, 2.3f, 2.4f
- Off-axis irradiance, 34.16, 34.16f
- Off-axis rejection (OAR), 7.23
- Office lighting, 40.55, 40.56t
- Offner compensators, 13.24, 13.24f
- Offset, thermocouple junctions as source of, 27.6, 27.6f
- Offset drift, 27.11
- Offset subtraction error, 34.32–34.33
- Ohmic contact, 17.13
- 1/*f* noise, 25.12, 27.4
- 110 photographic film, 30.21, 30.25
- Open arcs, 15.22
- Open-loop gain function, 22.9–22.10, 22.9f
- Open-tube diffusion, 17.24
- Optical analysis software, 40.20
- Optical athermalization, 8.12–8.15, 8.13t–8.15t
- Optical cavities, 19.18, 19.19f
- Optical cavity technique, 12.20, 12.20f
- Optical cavity-based frequency discriminators, 22.14–22.16, 22.17f
- Optical choppers, 27.14
- Optical communication systems, 19.3
- Optical components:
 purchasing of, 9.9
 specifications for systems vs., 4.3
- Optical confinement factor, 19.11
- Optical density, of photographic films, 29.6–29.8, 29.7f
- Optical design software, 3.1–3.24, 40.20–40.21
 about, 3.2
 and computing environment, 3.21
 data entry for, 3.2–3.8
 design process flowchart, 3.3f
 evaluation function of, 3.8–3.16
 aberrations, 3.9–3.11
 paraxial analysis, 3.8–3.9, 3.9f
 ray tracing, 3.11–3.13, 3.12f
 spot-diagram analysis, 3.13–3.16
 global optimization with, 3.21
 optimization function of, 3.16–3.21
 programming considerations for, 3.7–3.8
 purchasing of, 3.22–3.24
 setup routine in, 3.7
 simulation with, 3.21
- Optical image transformers, 39.21
- Optical multichannel analyzers (OMAs), 31.27–31.28
- Optical parametric oscillators (OPOs), 20.20–20.22, 20.20f, 20.21f
- Optical path difference (OPD), 2.1, 2.6, 3.12–3.13, 3.12f, 8.1, 8.7, 13.14–13.15
- Optical power (*see* Radiant flux (power))
- Optical pulse(s), 20.2–20.15
 coupling of circulating, 20.12–20.15, 20.12f, 20.15f
 in high gain oscillators, 20.10–20.12, 20.11f
 in ideal cavity, 20.6–20.7
 and pulse train, 20.2–20.9
 single, 20.2–20.3, 20.3f
 toward steady-state, 20.9–20.12
- Optical pumping, 16.16–16.19, 16.16f–16.18f
- Optical radar, 12.4, 12.5
- Optical software (for stray light suppression), 7.24–7.27
 advantages/disadvantages of, 7.29, 7.29f
 ASAP, 7.25
 CODE V, 7.26
 FRED, 7.25–7.26
 LightTools, 7.26
 SPEOS, 7.27
 TracePro, 7.27
 ZEMAX, 7.26–7.27
- Optical specifications, 4.1–4.12
 about, 4.1–4.2
 element description, 4.8–4.10
 environmental, 4.10

- Optical specifications (*Cont.*):
 image, 4.3, 4.6–4.8, 4.8f
 mechanical vs., 4.2
 preparing, 4.5–4.6, 4.6t
 presentation of, 4.10–4.11
 problems with writing, 4.11–4.12
 for systems vs. components, 4.3
 wavefront, 4.3–4.5, 4.5t
- Optical tolerances (*see* Tolerances)
- Optical transfer function (OTF), 29.17
- Optimization function (of optical software),
 3.16–3.21
 by damped least-squares method, 3.17–3.19
 and error functions, 3.19–3.20
 global, 3.21
 multiconfiguration, 3.20
 by orthonormalization, 3.19
 by simulated annealing, 3.19
 and tolerancing, 3.20–3.21
- OPTIS (simulation software), 7.27
- Optoelectronic integrated circuit (OEIC) chip,
 25.15
- Optronics Laboratories, Incorporated, 15.49
- Ordinary rays, 3.12
- Organic dye lasers, 16.31–16.32, 16.32f
- Organic LEDs (OLEDs), 40.37–40.39
- Organometallic vapor phase deposition
 (OMVPE), 19.6–19.7, 19.20t, 19.23
- Orthonormal polynomials, 11.3–11.40
 and aberration balancing, 11.30, 11.35–
 11.36, 11.36t
 about, 11.4–11.5
 circle, for noncircular pupils, 11.37, 11.39
 defined, 11.5–11.6
 discussion of, 11.39–11.40
 elliptical, 11.21, 11.25–11.27, 11.26t–11.27t
 hexagonal, 11.21, 11.22t–11.25t
 isometric, interferometric, and PSF plots for
 orthonormal aberrations, 11.36–11.37,
 11.37f, 11.38f
 rectangular, 11.27–11.28, 11.28t, 11.29t
 slit, 11.30, 11.35t
 square, 11.30, 11.31t–11.34t
 Zernike annular, 11.13–11.21, 11.14f, 11.16f,
 11.17t–11.21t
 Zernike circle, 11.6–11.12, 11.8t–11.9t,
 11.9f–11.11f, 11.12t
- Orthonormalization, 3.19
- Oscillation, relaxation, 16.12, 19.31–19.34,
 19.31f
- Oscillators:
 high-gain, 20.10–20.12, 20.11f
 optical-parametric, 20.20–20.22,
 20.20f, 20.21f
- Out-of-band radiation errors, 34.36
- Output amplifier noise, in CCDs, 32.20
- Output circuits, direct readout architectures,
 33.18
- Output coupling mirror, 16.11
- Output gate (OG), 32.14
- Output windows, in proximity-focused MCP
 IIs, 31.9, 31.9f
- Overhead lighting, 40.12, 40.13f, 40.14, 40.46f
- Overloosening and overtightening (tolerancing
 problems), 5.11
- Oxide stripe lasers, 19.8, 19.9f
- Packages and packaging:
 HB-LED, 18.5–18.6, 18.5f, 18.6f
 of photodetectors, 26.9–26.10, 26.10f, 26.11f
 reliability of LED, 17.25–17.26
- Paint modeling, 40.17
- Paralyzing glare, 40.9
- Parametric downconversion, 23.14
- Paraxial ray tracing, 3.5, 3.8–3.9, 3.9f
- Paraxial rays, 3.3
- PART (stray light analysis program), 7.11
- Part active, part passive athermalization,
 8.11–8.12, 8.12f
- Particle hypothesis, Einstein's, 23.7
- Particle pumping, 16.14–16.16, 16.15f, 16.16f
- Parts per million (ppm), 17.25
- Passivation, in wafer processing, 17.23
- Passive athermalization, 6.22, 6.23f, 6.24,
 8.8–8.10, 8.8f–8.10f
- Pattern effect, of light pulses, 19.32, 19.32f
- Pattern noise, 32.12
- Penalty function method, 3.18, 3.19
- Pen-Ray, 15.36f
- Pentaprisms, 6.14f, 6.15f, 12.12, 12.12f
- Perception, of lit environment, 40.1–40.2,
 40.4–40.6, 40.4f
- Perceptual constancy, 40.5
- Perfect lens, 4.4
- Perfect reflecting diffusers, 37.8
- Perfect transmitting diffusers, 37.8
- Periscopes, 1.10, 1.10f
- Petzval surface, 2.4, 2.5
- Phase conjugated coupling, 20.14
- Phase diffusion coefficient, 23.29

- Phase errors, phase-shifting interferometry and, 13.22
- Phase modulation index, 22.4
- Phase response, frequency vs., 22.6–22.7, 22.6f, 22.7f
- Phase space, 39.3
- Phase stability margin, 22.9
- Phase-interruption broadening, emission-line, 16.5
- Phase-lock phase shifting, 13.23
- Phase-locked laser arrays, 19.26, 19.27, 19.28f, 19.28t
- Phase-matching, in attosecond optics, 21.4
- Phase-shifting interferometry, 13.18–13.23, 13.18f–13.20f
- heterodyne interferometer, 13.22
 - integrating bucket method, 13.21, 13.21f
 - phase errors, 13.22
 - phase-lock method, 13.23, 13.23f
 - phase-stepping method, 13.20, 13.20f
 - simultaneous measurement, 13.22
 - two-steps-plus-one method, 13.21, 13.22
- Phase-stepping phase shifting, 13.20, 13.20f
- Phonon broadening, emission-line, 16.5
- Phonon coupling, 17.16
- Phosphor salts, 40.31
- Phosphor screens:
 - of image intensifiers, 31.14–31.16, 31.14t, 31.15f
 - in proximity-focused MCP IIs, 31.9, 31.9f
- Phosphor-type designation system, 31.14t
- Phot (unit), 34.43, 36.7, 36.7t
- Photo cell (*see* Photodiodes)
- Photo excitation, 25.5
- Photo gain, in photoconductors, 25.5–25.6, 25.5f
- Photocapacitors, MOS (*see* Metal-oxide-semiconductor capacitors)
- Photocathodes, 27.6, 27.7f, 31.9, 31.9f
- assemblies of, 31.10–31.12, 31.11f, 31.12f
 - internally and remotely processed, 31.10, 31.24
- Photochemical blue-light hazard, 36.17
- Photoconductive (PC) arrays, 33.4
- Photoconductive gain, 24.11
- Photoconductors, 24.7–24.8, 24.7f
- electronics of, 38.10
 - fabrication of, 32.2
 - image sensing with, 32.8–32.9
 - operating principles of, 25.4–25.6, 25.4f, 25.5f
- Photoconductors (*Cont.*):
 - photo gain in, 25.5–25.6, 25.5f
 - types of, 25.5, 25.5f
 - (*See also specific photoconductors, e.g.:*
 - Amorphous silicon photoconductors)
- Photodetection, 25.1–25.17
- about, 25.2–25.3
 - applications for, 25.11–25.12
 - future directions in, 25.15–25.17, 25.15f–25.17f
 - materials for, 25.3t
 - operating principles of, 25.3–25.11, 25.3f, 25.4f
 - extended wavelength photodetectors, 25.10, 25.10t
 - photoconductors, 25.4–25.6, 25.4f, 25.5f
 - photogate, 25.10, 25.11, 25.11f
 - pin* photodiodes, 25.6–25.10, 25.7f, 25.9f
 - reliability of, 25.13–25.14, 25.13t
- Photodetectors, 24.3–24.101
- AlGaN alloy photovoltaic detectors, 24.46
 - CdS photoconductors, 24.49–24.52, 24.51f–24.53f
 - CdSe photoconductors, 24.49–24.52, 24.52f
 - CdTe detectors, 24.52, 24.54, 24.54f
 - CdZnTe detectors, 24.52
 - GaAsP photodiodes, 24.49, 24.49f, 24.50f
 - GaN photovoltaic detectors, 24.42, 24.43, 24.45f, 24.46, 24.46f, 24.47
 - GaP photodiodes, 24.47–24.49, 24.48f
 - Ge intrinsic photodetectors, 24.70–24.73, 24.70f–24.73f
 - Ge low-temperature bolometers, 24.31–24.32, 24.32f, 24.33f
 - Ge:Au detectors, 24.83–24.85, 24.84f–24.86f
 - Ge:Cu detectors, 24.84f, 24.85f, 24.96, 24.97, 24.97f–24.99f
 - Ge:Ga infrared detectors, 24.100
 - Ge:Hg detectors, 24.84f, 24.92–24.95, 24.93f–24.95f
 - Ge:Zn detectors, 24.84f, 24.98–24.100, 24.99f
 - HgCdTe detectors, 24.86–24.92, 24.86f–24.92f
 - InAs photovoltaic detectors, 24.75, 24.77–24.78, 24.77f–24.79f
 - InGaAs detectors, 24.65–24.70, 24.66f–24.69f
 - InGaAs photodiodes, 34.31
 - InSb hot-electron bolometers, 24.29, 24.30, 24.30f, 24.31f
 - InSb intrinsic photovoltaic detectors, 24.80–24.83, 24.82f, 24.83f
 - life test for, 25.13, 25.14, 25.14f, 25.15f

- Photodetectors (*Cont.*):
- manufacturers' specifications for, 24.21–24.101, 24.21*t*
 - noise, impedance, dark and leakage current in, 24.19–24.20, 24.20*f*
 - PbS photoconductors, 24.73–24.74, 24.74*f*–24.77*f*, 24.74*t*
 - PbSe detectors, 24.76*f*, 24.78, 24.79, 24.79*f*–24.82*f*
 - PbSnTe photovoltaic detectors, 24.92, 24.93*f*
 - performance/sensitivity of, 24.13–24.18
 - background-limited case, 24.14–24.17, 24.16*f*, 24.16*t*, 24.17*f*
 - strong-signal case, 24.14
 - thermal detectors, 24.17, 24.18*f*
 - photoemissive detectors, 24.32–24.42, 24.33*f*, 24.35*f*–24.41*f*, 24.43*f*, 24.44*f*
 - photographic emulsions, 24.100, 24.101*f*
 - planar, 25.14, 25.15*f*
 - pyroelectric detectors, 24.26–24.29, 24.26*f*–24.29*f*
 - quantum, 24.6–24.10, 24.7*f*–24.9*f*
 - quantum well, 25.16–25.17, 25.16*f*, 25.17*f*
 - quantum well infrared, 25.15
 - responsivity and quantum efficiency of, 24.18, 24.19
 - Si photovoltaic detectors, 24.52*f*, 24.54–24.65, 24.55*f*–24.66*f*
 - Si:B detectors, 24.95*f*, 24.96
 - SiC UV detectors, 24.47, 24.47*f*
 - Si:Ga infrared detectors, 24.95, 24.95*f*, 24.96, 24.96*f*
 - and signal detection, 27.2
 - spectral response of, 24.18, 24.19*f*
 - speed of, 24.20, 24.21
 - stability of, 24.21, 24.21*f*
 - terminology, 24.10–24.13
 - thermal detectors, 24.4–24.6, 24.4*f*–24.6*f*
 - thermistor bolometers, 24.24–24.25, 24.24*f*, 24.25*f*, 28.7*t*
 - thermocouples, 24.22–24.23, 24.22*f*
 - thermopiles, 24.23–24.24, 24.23*f*
 - TiO₂ UV detectors, 24.47, 24.48*f*
 - types of, 25.3, 25.4*f*
 - uniformity of, 24.20
- (*See also specific photodetectors, e.g.: Avalanche photodetectors*)
- Photodiode CCDs (PD-CCDs), 33.11*f*, 33.12
- Photodiode linear arrays (PDAs), 38.9, 38.10, 38.10*t*
- Photodiode MOSs (PD-MOSs), 33.11*f*, 33.13
- Photodiodes (PDs):
- CCD, 33.11*f*, 33.12
 - defined, 24.11
 - electronics of, 38.10
 - GaAsP, 24.49, 24.49*f*, 24.50*f*
 - GaP, 24.47–24.49, 24.48*f*
 - Ge avalanche, 24.70*f*, 24.72–24.73, 24.72*f*, 24.73*f*
 - GeGaAs, 34.31
 - InGaAs, 34.31
 - InGaAs avalanche, 24.66–24.70, 24.66*f*–24.69*f*
 - junction, 32.3–32.6, 32.4*f*, 32.6*f*
 - MOS, 33.11*f*, 33.13
 - pin* (*see pin photodiodes*)
 - p⁺np*, 32.8
 - silicon, 34.29, 34.30
 - silicon avalanche, 24.62–24.65, 24.63*f*–24.66*f*
 - silicon *pn*, 24.52*f*, 24.55–24.58, 24.55*f*–24.59*f*
 - UV-enhanced, 24.55*f*, 24.61–24.62, 24.61*f*, 24.62*f*
- (*See also specific photodiodes, e.g.: Avalanche photodiodes*)
- Photoelectromagnetic (PEM) detectors, 24.9, 24.9*f*
- Photoemissive detectors, 24.6, 24.7*f*
- gallium phosphide dynodes, 24.42, 24.44*f*
 - linearity of, 24.41
 - manufacturers of, 24.42
 - noise from, 24.39, 24.40, 24.41*f*
 - operating temperature of, 24.40
 - photon counting for, 24.42, 24.43*f*, 24.44*f*
 - quantum efficiency of, 24.35*f*–24.38*f*, 24.36–24.38
 - recommended circuit for, 24.42, 24.43*f*
 - response time of, 24.40
 - responsivity of, 24.35*f*, 24.38
 - sensitive area of, 24.41
 - sensitivity of, 24.34, 24.35*f*–24.39*f*
 - sensitivity profile of, 24.41
 - short-wavelength considerations for, 24.34, 24.40*f*
 - specifications for, 24.32–24.42, 24.33*f*
 - stability of, 24.41
- Photogates, 25.10, 25.11, 25.11*f*
- Photographic detectors, 24.9–24.10
- Photographic dyes, 30.10–30.13, 30.10*f*, 30.12*f*
- Photographic emulsions, 24.100, 24.101*f*, 29.4, 30.7

- Photographic film speed, 29.9–29.10
 in Advanced Photo System, 30.26
 in color negative film, 30.25
 and granularity, 30.19
 high vs. low, 30.18–30.20
 and sensitivity to high-energy radiation,
 30.19–30.20
- Photographic films, 29.3–29.16, 30.18–30.28
 about, 30.2–30.3
 black-and-white (B&W) film, 30.24–30.25,
 30.25*t*
 color, 29.12–29.15, 29.13*f*, 29.14*f*
 color negative film, 30.25–30.28, 30.27*t*
 color reversal film, 30.22–30.24, 30.23*t*
 development effects on, 29.12, 29.13*f*
 D-log H curve for, 29.8–29.10, 29.8*f*, 29.10*f*
 exposure of, 29.5–29.6
 grain element of, 29.5
 granularity of, 30.19
 high-speed vs. low-speed, 30.18–30.20
 and light scattering by silver halide crystals,
 30.5–30.7, 30.6*f*
 microdensitometers, 29.15–29.16, 29.15*f*
 optical density of, 29.6–29.8, 29.7*f*
 processing of, 29.5
 professional vs. amateur film, 30.20–30.22
 reciprocity failure of, 29.11–29.12
 spectral sensitivity of, 29.11, 29.11*f*
 speed of, 30.18–30.20
 structure of color, 30.3–30.5, 30.3*f*
 structure of silver halide photographic layers
 in, 29.4
- Photographic materials, 30.1–30.28
 about, 30.1–30.2
 dyes, 30.10–30.13
 about, 30.10, 30.10*f*
 excited state properties, 30.11–30.12, 30.12*f*
 light stabilization methods, 30.12–30.13
 photochemistry of azomethine dyes, 30.11
 films, 30.18–30.28
 black-and-white film, 30.24–30.25, 30.25*t*
 color negative film, 30.25–30.28, 30.27*t*
 color reversal film, 30.22–30.24, 30.23*t*
 professional vs. amateur film,
 30.20–30.22
 speed, 30.18–30.20
 optics of, 30.2–30.7
 about, 30.2–30.3
 light scatter by silver halide crystals,
 30.5–30.7, 30.6*f*
- Photographic materials, optics of (*Cont.*):
 structure of color films, 30.3–30.5, 30.3*f*
 structure of color papers, 30.5
 and photographic spectral sensitizers,
 30.13–30.18
 about, 30.13–30.14, 30.14*f*
 color science, 30.15–30.18, 30.16*f*, 30.17*f*
 photophysics of spectral sensitizers on silver
 halide surfaces, 30.14–30.15, 30.15*f*
 silver halide light detectors, 30.7–30.9, 30.8*f*
- Photographic papers:
 about, 30.2–30.3
 and light scattering by silver halide crystals,
 30.5–30.7, 30.6*f*
 structure of color, 30.5
- Photographic spectral sensitizers (*see* Spectral
 sensitizers, photographic)
- Photographic systems, 29.16–29.25
 acutance of, 29.17–29.19, 29.18*t*, 29.19*f*
 detective quantum efficiency of, 29.23
 graininess in, 29.19–29.22, 29.21*f*
 image structure of, 29.17
 information capacity of, 29.24
 manufacturers of, 29.25
 performance of, 29.16–29.17, 29.17*f*, 29.18*t*
 resolving power of, 29.24
 sharpness in, 29.22
 signal-to-noise ratio of, 29.22–29.23
- Photoionization detectors, 24.10
- Photoionization devices, 34.29
- Photoionization yield, 34.29
- Photometry, 34.37–34.44, 36.1–36.17, 36.3*f*,
 36.3*t*
 about, 36.2–36.4
 approximations, 36.10, 36.10*f*, 36.11*f*
 basis of physical, 37.1–37.2, 37.2*f*
 calibrations in, 34.42–34.43
 concepts/terminology of, 34.10–34.11,
 34.38*t*, 34.39–34.40, 34.43–34.44, 39.2,
 39.2*t*
 conversion between radiometric and
 photometric quantities, 34.12*t*,
 36.11–36.14, 36.12*f*–36.14*f*
 defined, 34.37, 37.1
 and human eye, 36.8–36.10, 36.8*f*, 36.9*f*
 illuminance-luminance relationship, 37.9,
 37.9*f*
 integrating sphere device, 37.9–37.10, 37.10*f*
 inverse square law, 37.8
 Lambert's cosine law, 37.8, 37.8*f*

- Photometry (*Cont.*):
 normalization, 36.14–36.17, 36.15*t*, 36.16*f*
 and photopic/scotopic/mesopic vision,
 34.37–34.39, 34.38*t*
 practice in, 37.11
 quantities and units in, 37.4–37.8, 37.4*t*,
 37.5*f*, 37.7*t*
 radiometry vs., 34.6
 retinal illuminance, 34.40–34.42
 symbols/nomenclature of, 36.5–36.10,
 36.6*t*–36.8*t*
 weighting functions, 36.17
- Photomicrographic lamps, 15.47–15.49,
 15.48*f*, 15.49*f*
- Photomultiplier tubes (PMTs), 24.11,
 24.32–24.34, 24.33*f*, 24.38–24.42, 24.38*f*,
 24.39*f*, 24.43*f*, 38.9, 38.9*t*
 applications for, 27.6–27.10, 27.7*f*
 base design of, 27.8–27.10
 electronics of, 38.10
- Photon(s), 23.6–23.14, 25.2, 34.30, 36.4
 Einstein's light quanta, 23.6–23.9, 23.8*f*
 photon-photon correlations, 23.13–23.14,
 23.13*f*
 quantum electrodynamics, 23.9–23.13
- Photon counting, 27.15
 defined, 24.12
 of modulated signal sources, 27.14
 in photoemissive detectors, 24.42, 24.43*f*,
 24.44*f*
- Photon density, 19.30–19.35
- Photon detectors:
 background-limited case of, 24.14–24.17,
 24.16*f*, 24.16*t*, 24.17*f*
 strong-signal case of, 24.14
- Photon dose, 34.6, 34.11
- Photon Engineering, 7.25
- Photon flux, 34.5, 34.11
- Photon infrared detectors, 33.7, 33.8*f*
- Photonic excitation, 25.3, 25.3*f*
- Photonics Directory of Optical Industries*, 15.14
- Photopic vision, 34.37–34.39, 34.38*t*,
 36.8*f*–36.10*f*, 37.2, 40.3
- Photovoltaic (PV) arrays, 33.4
- Photovoltaic detectors, 24.8, 24.8*f*, 24.9
 aluminum gallium nitride, 24.46
 gallium nitride, 24.42, 24.43, 24.45*f*, 24.46,
 24.46*f*, 24.47
 indium antimonide, 24.80–24.83, 24.82*f*,
 24.83*f*
- Photovoltaic detectors (*Cont.*):
 indium arsenide, 24.75, 24.77–24.78,
 24.77*f*–24.79*f*
 lead tin telluride, 24.92, 24.93*f*
 mercury cadmium telluride, 24.86*f*, 24.88*f*,
 24.90–24.92, 24.91*f*, 24.92*f*
 silicon, 24.52*f*, 24.54–24.65, 24.55*f*–24.66*f*
- Photovoltaic Schottky barrier detectors (SBDs),
 33.7, 33.8*f*
- Physical optics, 3.16
- Physical photometry, 34.37, 36.4, 37.2
- Pickups, 3.6
- Piezoelectric transducers (PZTs):
 amplifier strategies for, 22.18
 disk vs. tube, 22.17–22.18
- Piezoelectric-based (PZT-based) systems, 22.1
- pin* junctions, 26.3
- pin* photodetectors:
 biased, 26.6*f*
 high-speed, 26.10, 26.12–26.15, 26.12*f*,
 26.14*f*, 26.15*f*
 lateral, 25.15–25.16, 25.15*f*
- pin* photodiodes, 24.54, 24.55*f*–24.60*f*,
 24.58–24.61, 25.4*f*, 25.6–25.10, 25.7*f*, 32.4*f*
 absorption coefficient of, 25.8, 25.9*f*
 avalanche photodiodes, 25.8–25.10, 25.9*f*
 dark current in, 25.7–25.8
 diffusion current of, 25.7
 equivalent circuit of, 26.7, 26.7*f*
 generation-recombination current of,
 25.7–25.8
 germanium, 24.70–24.71, 24.70*f*–24.72*f*
 InGaAs, 24.66, 24.67*f*
 operating principles of, 25.6–25.10,
 25.7*f*, 25.9*f*
 quantum efficiency of, 25.8
 resonant, 26.15, 26.15*f*
 responsivity of, 25.8
 silicon, 24.55*f*–24.57*f*, 24.58–24.61, 24.59*f*,
 24.60*f*
 tunneling current of, 25.8
 vertically illuminated, 26.3, 26.4*f*, 26.5,
 26.10, 26.12–26.13, 26.12*f*
 waveguide, 26.13–26.14, 26.14*f*
- Pitch-based edging (fabrication step), 9.6
- Pixel summation, 38.10
- Planar buried heterostructure (PBH) lasers,
 double-channel, 19.24, 19.25*f*, 19.34*f*
- Planar diffused silicon photodiodes, 24.56*f*
- Planar photodetectors, 25.14, 25.15*f*

- Planckian radiation (*see* Blackbody radiation)
- Planck's formula, 23.6, 23.7
- Planck's law, 34.23, 34.24, 37.10–37.11
- Plane diffusers, 38.14, 38.14*f*, 38.15*f*
- Plano optics fabrication, 9.7
- Plastic, as photographic film emulsion, 29.4
- Plastic-packaged LEDs, 17.25–17.26
- Platings, diamond turning of, 10.5
- Platinum silicon (PtSi) infrared detectors, 33.7, 33.8*f*, 33.29, 33.29*t*
- Pluecker spectrum tubes, 15.47, 15.47*f*, 15.47*t*
- pn* junctions, 24.8, 24.8*f*
- pn* photodetectors, 24.70–24.71, 24.70*f*–24.72*f*
- pn* photodiodes, 24.52*f*, 24.54–24.58, 24.55*f*–24.59*f*, 34.30
- p⁺np* photodiodes, 32.8
- Point diffraction interferometers, 13.11*f*
- Point source irradiance transmittance (PSIT), 7.23
- Point source normalized irradiance transmittance (PSNIT), 7.22–7.23
- Point source power transmittance (PSPT), 7.23
- Point source transmittance (PST), 7.5, 7.6*f*, 7.22–7.23
- Point spread function (PSF), 3.16, 11.36–11.37, 11.37*f*, 11.38*f*
- Point-to-point approximation (of radiant flux transfer), 34.14
- “Poker chip” assembly, 6.8, 6.9*f*
- Polar angles, 35.5
- Polarization gating, 21.7–21.8
- Polarization-dependent systems, 34.33
- Polaroid Corporation, 29.12, 29.25
- Polaroid Instant Color Film, 29.14
- Polaroid “One Film,” 29.14
- Pole-mounted luminaires, 40.62, 40.63*t*
- Polishers, continuous, 9.7
- Polishing step (of optics fabrication), 9.5–9.6, 9.8
- Polychromatic radiation, 34.9–34.10
- Polyethylene terephthalate film, 29.4
- Polynomial numbering, 11.39
- Polynomial-ordering number, 11.7
- Polynomials, orthonormal (*see* Orthonormal polynomials)
- Polytetrafluoroethylene (PTFE), 35.13, 38.12–38.13
- Ponderomotive potential, 21.3
- Population inversions, 16.8–16.10, 16.12–16.13
described, 16.8–16.10
mechanism for achieving, 16.12–16.13, 16.13*f*, 16.14*f*
optical pumping for, 16.16–16.19, 16.16*f*–16.18*f*
particle pumping for, 16.14–16.16, 16.15*f*, 16.16*f*
semiconductor diode laser pumping for, 16.19
- Positive image (in photography), 29.9
- Potassium dihydrogen phosphate (KDP), 10.2
- Power measurement, for lasers, 34.32
- Power spectrum of granularity, 29.21
- Power supply, high-voltage, 31.9, 31.10*f*
- Predictable quantum efficiency (PQE) devices, 34.29–34.30
- Preloads, 6.2
- Primaries (colors stimuli), 30.16
- Principal rays, 1.4, 1.11*f*, 1.12
- Prisms, 40.45*f*, 40.46
angle measurement in, 12.14–12.16, 12.15*f*–12.17*f*
mounting of, 6.11–6.17
bonded mountings, 6.13–6.15, 6.15*f*, 6.16*f*
flexure mountings, 6.15–6.17, 6.16*f*
mechanically clamped mountings, 6.12, 6.13*f*
spring-loaded mountings, 6.13, 6.14*f*
penta, 6.14*f*, 6.15*f*, 12.12, 12.12*f*
right-angle, 12.15, 12.16, 12.16*f*, 12.17*f*
Risley, 12.4, 12.4*f*
rotating glass block, 12.4, 12.4*f*
sliding, 12.4, 12.4*f*
Zerodur, 6.16, 6.16*f*
- Procedural programs (optical software), 3.7
- Projected area, in photometry/radiometry, 36.3–36.4, 36.3*f*, 36.3*t*
- Projected solid angle (PSA), 39.5, 39.5*f*
- Projection density, 29.7, 29.7*f*
- Projection lenses, 6.6*f*
- Projection systems, 39.23–39.24, 39.23*f*
- Proportional integral derivative (PID) controllers, 22.10–22.12, 22.10*f*, 22.11*f*
- Proportional-integral (PI) amplifier circuit, 22.5*f*
- Proton stripe lasers, 19.35*f*
- Proton-bombardment-defined lasers, 19.41
- Proximity-focus electronic lens, 31.8, 31.23–31.24, 31.23*f*, 31.24*f*

- Proximity-focused MCP IIs, **31.9, 31.9f, 31.16–31.18, 31.17t, 31.18f, 31.19f**
- Psychophysical photometry, **34.37**
- P*-type impurities, **17.23**
- Pulse height, of photomultipliers, **27.7–27.8**
- Pulse train interferometry, **20.12, 20.12f**
- Pulse trains:
 about, **20.3–20.5, 20.4f, 20.5f**
 attosecond, **21.6, 21.7**
 and backscattering, **20.13–20.15, 20.15f**
 soliton solution and steady-state, **20.5–20.9**
- Pulsed lasers, **23.18**
- Pumping (for population inversions), **16.14–16.19**
 optical, **16.16–16.19, 16.16f–16.18f**
 particle, **16.14–16.16, 16.15f, 16.16f**
 semiconductor diode laser, **16.19**
- Punch through (in color films), **30.4**
- Pupils, noncircular, **11.4**
- Push processing (of film), **30.22**
- Pyramidal error, **12.14, 12.15, 12.15f**
- Pyroelectric detectors, **24.6, 24.6f, 24.26–24.29, 24.26f–24.29f, 28.2, 28.6, 28.6f, 28.7, 28.7t, 33.10**
- Pyroelectric electrical substitution radiometers, **34.27–34.28**
- Pyroelectric hybrid arrays, **28.11–28.12, 28.11f, 28.12f**
- Q-switched lasers, **16.26–16.27, 16.27f**
- Quality, image, **4.6–4.7**
- Quality Assurance of Ultraviolet Measurements in Europe (QASUME), **38.5**
- Quantized center-of-mass motion, of atoms, **23.45**
- Quantum box, **19.18**
- Quantum cascade lasers, **16.36**
- Quantum dot, **16.7, 19.18, 26.4f, 26.5**
- Quantum efficiency (QE), **25.3, 25.4, 34.29**
 defined, **24.12**
 detective, **29.23**
 of photodetectors, **24.18, 24.19**
 of photoemissive detectors, **24.35f–24.38f, 24.36–24.38**
 of photomultipliers, **27.7**
 of *pin* photodiodes, **25.8**
- Quantum electrodynamics (QED), **9.6, 23.9–23.13**
- Quantum limited imaging (QLI), **31.3–31.4**
- Quantum photodetectors, **24.6–24.10, 24.7f–24.9f**
- Quantum resonance absorption, **22.16, 22.17**
- Quantum sensitivity, **30.9**
- Quantum theory of lasers, **23.14–23.35**
 about, **23.5–23.6**
 density-operator approach to, **23.14–23.33**
 derivation of Scully-Lamb master equation, **23.17–23.22, 23.19f**
 photon statistics, **23.22–23.27, 23.23f, 23.25f, 23.27f**
 spectral properties, **23.28–23.33, 23.30f**
 spectrum, **23.28–23.33**
 time evolution of the field in Jaynes-Cummings model, **23.15–23.17, 23.15f**
 Heisenberg-Langevin approach to, **23.33–23.35**
- Quantum trajectories, in attosecond optics, **21.3–21.4**
- Quantum well (QW) detectors, **26.4f, 26.5**
- Quantum well infrared photodetectors (QWIPs), **25.15–25.17, 25.16f, 25.17f, 33.9**
- Quantum well (QW) lasers, **16.7, 19.9–19.18, 19.20t**
 GRIN SCH single, **19.14, 19.14f**
 long wavelength, **19.17–19.18, 19.17f**
 schematic of, **19.10f**
 strained, **19.15–19.17, 19.16f**
 threshold modal gain, **19.12–19.15, 19.12f, 19.13f, 19.15f**
- Quantum well (QW) photodetectors, **25.16–25.17, 25.16f, 25.17f**
- Quantum wire, **16.7, 19.18, 26.4f, 26.5**
- Quartz-envelope lamps, **15.20, 15.21**
- QW ridge (QWR) waveguide lasers, **19.19, 19.20t**
- Rabi cycles and Rabi cycling, **23.21**
 defined, **20.24**
 off resonance, **20.27**
 on resonance, **20.26–20.27, 20.26f, 20.27f**
- Rack-stack laser arrays, **19.29, 19.30f**
- Radial circle polynomial, **11.7, 11.9f–11.10f**
- Radial shearing interferometers, **13.12, 13.13f**
- Radian (rad), **36.3**
- Radiance, **34.9, 34.9f, 37.4t, 37.5, 38.2, 38.11t, 38.13–38.16, 38.13f–38.16f, 39.2t**
- Radiance conservation theorem, **34.12–34.13**
- Radiance temperature (unit), **37.4t, 37.6**
- Radiance units, **34.24**

- Radiant energy, 34.7, 37.4t, 37.6
- Radiant exitance (emittance), 15.4–15.6, 15.5t, 15.6f, 35.3, 37.4t, 37.5, 37.5f
- Radiant exposure, 37.4t, 37.6
- Radiant flux (power), 34.7, 34.11, 34.17–34.18, 36.4, 36.6t, 37.3, 37.4t, 37.6, 38.2
- Radiant incidence (*see* Irradiance)
- Radiant intensity, 36.4, 37.4, 37.4t, 39.2t
- Radiant power transfer, 34.12–34.13, 34.13f
- Radiant transfer approximations, 34.13–34.20
- approximate radiance at an image, 34.19–34.20
 - lambertian, 34.14–34.18, 34.15f, 34.16f
 - point-to-point, 34.14
 - radiometric effect of stops and vignetting, 34.18–34.19, 34.19f
- Radiation, 34.23–34.27
- actinic effects of, 34.6, 34.7
 - artificial sources of (*see* Artificial sources (of radiation))
 - baseline standard of, 15.9, 15.9f, 15.10f, 15.12f
 - from blackbodies, 34.23–34.24
 - from blackbody simulators, 34.24–34.26
 - between circular source and detector, 34.15–34.16, 34.15f
 - commercial sources of (*see* Commercial sources (of radiation))
 - incandescent sources of (*see* Incandescent sources (of radiation))
 - infrared (*see* Infrared (IR) radiation)
 - and lasers, 16.4
 - photographic film speed and sensitivity to high-energy, 30.19–30.20
 - polychromatic, 34.9–34.10
 - from synchrotrons, 34.26–34.27
 - through absorbing media, 34.13
 - transfer of, 7.21–7.22
 - ultraviolet (*see* Ultraviolet (UV) radiation)
 - working standards of, 15.9–15.13, 15.10f, 15.12f, 15.13f
- Radiation law, 15.4–15.7, 15.5f, 15.5t, 15.6f
- Radiative lifetimes, 16.4, 17.4
- Radiative recombination, 17.2
- Radiators, blackbody, 34.23–34.24
- Radio frequency (rf) modulation, 22.14
- Radiometers and radiometry, 34.3–34.37, 36.1–36.17, 36.3f, 36.3t
- about, 34.5–34.7, 36.2–36.4
- Radiometers and radiometry (*Cont.*):
- approximate (*see* Radiant transfer approximations)
 - cavity-shaped, 34.28
 - concepts/terminology of, 34.7, 39.2, 39.2t
 - conversion between radiometric and photometric quantities, 34.12t, 36.11–36.14, 36.12f–36.14f
 - defined, 37.1
 - electrical substitution, 34.27–34.29
 - geometrical concepts of, 34.8–34.9, 34.9f
 - of II electronic imaging, 31.4–31.5
 - illuminance-luminance relationship, 37.9f
 - integrating sphere device, 37.9–37.10, 37.10f
 - laser as characterization tool for, 34.32
 - normalization, 36.14–36.17, 36.15t, 36.16f
 - photometry vs., 34.6
 - Planck's law, 37.10–37.11
 - quantities and units in, 37.3–37.7, 37.4t, 37.5f, 37.7t
 - spectral dependence of, 34.9–34.10
 - Stefan-Boltzmann's law, 37.11
 - symbols/units/nomenclature of, 36.4–36.5
 - thermopile-based, 34.27
 - weighting functions, 36.17
 - Wien's displacement law, 37.11
- Range finders, 12.3–12.4, 12.3f, 12.4f
- Range gating, 31.28–31.30
- Ray displacement, 3.12
- Ray intercept curves, 2.2–2.4, 3.13
- Ray sets, 3.20
- Ray tracing, 1.4–1.5, 1.4f, 3.11–3.13, 3.12f
- in lighting simulation, 40.20
 - nonsequential, 39.6
 - in optical design software, 3.11–3.13
 - paraxial, 3.5, 3.8–3.9, 3.9f
- Rays:
- axial rays, 1.4, 1.11f, 1.12
 - dashed rays, 1.12, 1.12f
 - edge rays, 39.22, 39.38
 - exact, 3.3, 3.11–3.12
 - hamiltonian rays, 3.12
 - iterated rays, 3.12
 - lagrangian rays, 3.12
 - meridional rays, 3.3
 - ordinary, 3.12
 - paraxial, 3.3
 - principal rays, 1.4, 1.11f, 1.12
- RC time constant, 26.7–26.8, 26.7f
- Reactive ion etching (RIE), 18.3, 19.39

- Readouts (of visible array detectors),
 32.12–32.21
 CCD, 32.12–32.20, 32.13f, 32.15f–32.18f
 MOS, 32.20–32.21
- Real-space critical objects, 7.2–7.4, 7.3f, 7.4f
- Reasonableness, of layout, 1.13–1.14
- Recessed lighting, 40.13f
- Reciprocity failure, of photographic films,
 29.11–29.12
- Recombination:
 combined, 17.3
 exciton, 17.6
 in GaAs, 17.8, 17.9, 17.9f
 minority-carrier, 17.2
 nonradiative, 17.2
 radiative, 17.2
- Reconstruction of attosecond beating by
 interference of two-photon transition
 (RABITT), 21.9
- Rectangular polynomials, 11.27–11.28, 11.28t,
 11.29t, 11.36t
- Red light, and color film, 29.13, 29.13f
- Reduction scanners, in linear sensors,
 32.21, 32.22f
- Reflectance:
 classification of materials by, 35.4t
 defined, 35.4–35.5
 geometrical definitions of, 35.6f
 and illuminance/luminance, 37.9
 measurement of, 35.10–35.13, 35.10f–35.12f
 nomenclature for, 35.5t, 35.6t
 spectral, 38.2, 38.17–38.18, 38.18f
 standards of, 35.14t
 and transmittance/absorptance, 35.7, 35.8,
 35.8t
- Reflecting telescopes, 11.4
- Reflection(s):
 actual/idealized, 35.2f
 defined, 35.4
 veiling, 40.12
- Reflection density, of photographic films,
 29.8
- Reflective compensators, for spherical
 aberrations, 13.24, 13.24f, 13.25
- Reflective-refractive (RX) concentrators, 39.17,
 39.17f
- Reflectometers, 35.10, 40.52
- Reflectors:
 conic, 39.11, 39.11f
 convergent, 39.38–39.40, 39.38f, 39.39f
- Reflectors (*Cont.*):
 CPC-type (*see* Compound parabolic
 collectors)
 divergent, 39.8f, 39.9f, 39.38–39.40
 faceted, 39.39–39.41, 39.39f, 39.40f
 headlamp, 40.64
 homogeneous/inhomogenous, 39.39
 involute, 39.11–39.12, 39.12f
 and lens array combinations, 39.34–39.37,
 39.36f, 39.37f
 luminaire, 40.45, 40.45f
 macrofocal, 39.11
 tailored, 39.37–39.39, 39.38f
- Reflexive sensors, 17.34
- Refraction index, 17.34
- Refractive compensators, for spherical
 aberrations, 13.24, 13.24f, 13.25
- Refractive index, 3.6, 34.13
- Relative measurements, absolute vs.,
 34.20–34.21
- Relative visual performance (RVP) model,
 40.5–40.6
- Relaxation oscillation, 16.12, 19.31–19.34,
 19.31f
- Rem jet, 30.4
- Remotely processed (RP) photocathodes,
 31.10, 31.24
- Repetition rate coupling, 20.14–20.15, 20.15f
- Rescattering model, semiclassical, 21.3
- Reset gate (RG), 32.14
- Residential lighting, 40.57, 40.58, 40.59t
- Residual amplitude modulation (RAM),
 22.14
- Resistive bolometers, 28.10–28.11, 28.10f, 33.9
- Resistive coupling noise, 27.5, 27.6f
- Resolution, of optical system, 4.6
- Resolving power, photographic, 29.24
- Resonant optical feedback, 19.38, 19.38f
- Resonant photodetectors, 26.4f
- Resonant *pin* photodiodes, 26.15, 26.15f
- Response time:
 defined, 24.12
 of photodetectors, 25.4
 of photoemissive detectors, 24.40
- Responsive quantum efficiency (*see* Quantum
 efficiency)
- Responsivity:
 blackbody, 24.10
 of photodetectors, 24.18, 24.19, 25.4
 of photoemissive detectors, 24.35f, 24.38

- Responsivity (*Cont.*):
 of *pin* photodiodes, 25.8
 spectral, 24.12, 38.3, 38.18–38.19
 of spectroradiometers, 38.11–38.12, 38.12*f*
- Restricted exit angle concentrators,
 39.18, 39.18*f*
- Retail lighting, 40.55–40.57, 40.56*t*–40.58*t*
- Reticles and reticulation, 12.13, 28.11, 28.12
- Retina, 34.37
- Retinal damage, 40.9
- Retinal illuminance, 34.40–34.42
- Retinal thermal hazard, 36.17
- Retroreflection, measurement of, 35.13
- Reverse bias, 26.3
- Reversing shear interferometers, 13.12, 13.13*f*
- Rhenium, 40.27
- Ribbon-type tungsten filaments, 15.20*f*
- Ridge waveguide (RWG) lasers, 19.8, 19.9*f*
- Right-angle prisms, 12.15, 12.16, 12.16*f*, 12.17*f*
- Ring flanges, continuous, 6.4*f*, 6.11
- Ring lasers, 16.29
 with additional Kerr crystal, 20.17–20.18,
 20.17*f*
 dye, 20.15–20.16
 Ti:sapphire, with saturable absorber,
 20.16–20.17, 20.16*f*, 20.17*f*
 in two-level system analogy, 20.24, 20.25*f*
- Rings, aperture, 3.20
- Risley prisms, 12.4, 12.4*f*
- Ritchey-Chretien two-mirror imaging system,
 39.17
- RLM lamp, 40.46, 40.46*f*, 40.47
- RMS noise, 24.12
- RMS signal, 24.12
- rms-granularity, 29.19–29.21
- Roadway lighting, 40.67, 40.69–40.71
 and disability glare, 40.10
 and discomfort glare, 40.12
 sign lighting, 40.71
 street lighting, 40.69–40.71, 40.70*t*, 40.71*t*
 tunnel lighting, 40.71
- Robertson's correlated color temperature
 calculation, 38.5
- Rods (eye receptors), 30.15, 30.16*f*, 34.37,
 36.8–36.10, 36.8*f*, 36.9*f*
- Rome Air Development Center, 7.19
- Ronchi test, 13.3–13.4, 13.3*f*, 13.4*f*
- Roof-mirror-lens arrays, 32.21, 32.22*f*
- Room temperature vulcanizing (RTV) sealing
 compound, 6.4
- Root-mean-square (rms) wavefront error, 4.1,
 4.3, 4.7, 4.8
- Rotating glass block prisms, 12.4, 12.4*f*
- Rotational shear interferometers, 13.12, 13.13*f*
- Rotationally symmetric aspheric lenses, 9.7, 9.7*f*
- Rotationally symmetric optics:
 hard mounting of, 6.2–6.4, 6.3*f*, 6.4*f*
 soft mounting of, 6.4–6.5, 6.4*f*, 6.5*f*
- Ruby lasers, 16.12, 16.13*f*, 16.32
- Rule-of-thumb PID design, 22.11–22.12
- “Rule-of-thumb” tolerance, 6.2
- Rydberg states, 23.21
- Sapphire (Ti:Al₂O₃) lasers, titanium-doped,
 16.34, 16.34*f*
- Sapphire (Ti:Al₂O₃) ring lasers, titanium-
 doped, 20.16–20.17, 20.16*f*, 20.17*f*
- Sapphire substrate (for HB-LEDs), 18.2, 18.3,
 18.5, 18.6
- Satellite spheres, 39.26
- Saturated colors, 40.7
- Saturation, 40.5, 40.9
- Sawing (of LEDs), 17.24, 17.25
- Scanning arrays, 33.6, 33.6*f*, 33.14
- Scanning FPAs, 33.17, 33.17*f*
- Scanning white light interferometry (SWLI),
 10.13
- Scattered radiation effect, 34.33
- Scattering:
 and photographic film, 29.18
 rescattering, 21.3
 by silver halide crystals, 30.5–30.7, 30.6*f*
 spectral, 38.10
 surface, 7.23
- Scattering sensors, 17.34
- Schawlow-Townes linewidth, 23.18
- Schmidt-Cassegrain design, 7.20
- Schottky barrier detectors (SBDs), 33.7, 33.8*f*,
 33.12–33.13
- Schottky contact, 26.3
- Schottky junctions, 26.3
- Schottky photodiodes, 26.16, 26.16*f*, 26.17*f*
- Scotopic vision, 34.37–34.39, 34.38*t*, 36.8*f*,
 36.9, 36.9*f*, 36.10, 36.11*f*, 37.2, 40.3
- Scully-Lamb master equation, 23.15,
 23.17–23.22
 cavity losses, 23.18
 laser master equation, 23.19–23.20, 23.19*f*
 micromaser master equation, 23.20–23.22
- Sealed beam lights, 40.26*t*

- Sealed-ampoule diffusion, 17.23
- Secondary spectrum, of radiation, 1.14, 1.15
- Second-order autocorrelator, 21.9
- Seidel (third-order monochromatic) aberrations, 3.9–3.10
- Seidel astigmatism, 11.27, 11.30, 11.35, 11.39, 11.40
- Self-athermalized, 8.7–8.8, 8.7f
- Self-calibration, of silicon photodiodes, 34.29
- Selfoc lenses, 32.21, 32.22f
- Self-phase modulation, 20.6
- Self-scanned array, 31.2
- Selwyn coefficient, 29.20
- Selwyn's law, 29.20
- Semiconductor bolometers, 28.4–28.5
- Semiconductor laser pumping, 16.19
- Semiconductor lasers, 16.35–16.36, 16.35f, 19.1–19.43
- applications for, 19.3–19.4
 - arrays of, 19.26–19.29, 19.28f, 19.28t
 - fabrication and configurations of, 19.6–19.8, 19.9f
 - gain mechanism of, 19.4
 - high-power semiconductor lasers, 19.18–19.30
 - arrays, 19.26–19.29, 19.28f, 19.28t
 - commercial, 19.19–19.23, 19.20t, 19.21f, 19.22f
 - future directions for, 19.23–19.26, 19.25f, 19.25t, 19.26t, 19.27f
 - mode-stabilized lasers, 19.18–19.19, 19.19f
 - two-dimensional, 19.29–19.30, 19.29t, 19.30f
 - high-speed modulation of, 19.30–19.36, 19.31f–19.36f
 - operation of, 19.4–19.6, 19.4f–19.6f
 - quantum cascade lasers, 16.36
 - quantum well lasers, 19.9–19.18
 - GRIN SCH single, 19.14, 19.14f
 - long wavelength, 19.17–19.18, 19.17f
 - schematic of, 19.10f
 - strained, 19.15–19.17, 19.16f
 - threshold modal gain, 19.12–19.15, 19.12f, 19.13f, 19.15f
 - spectral properties of, 19.36–19.39, 19.37f, 19.38f
 - surface-emitting lasers, 19.39–19.41
 - distributed grating, 19.40–19.41, 19.40f
 - integrated laser with 45° mirror, 19.39–19.40, 19.39f
 - vertical cavity, 16.36, 19.41, 19.42f, 19.43t
- Semiconductor photodetectors, 25.2, 38.9, 38.9t
- Semiconductors:
- arrays of, 16.36
 - direct, 17.4, 17.5f
 - indirect, 17.4, 17.4f, 17.5f, 17.6
 - material systems for, 18.1–18.2
 - properties of substrates for, 17.20t
 - waveband structure of, 17.3–17.6, 17.3f–17.5f
- Semisealed-ampoule diffusion, 17.24
- Sensitive area, of photoemissive detectors, 24.41
- Sensitivity:
- defined, 24.12
 - film spectral, 29.11, 29.11f
 - film speed and, to high-energy radiation, 30.19–30.20
 - of interferometers, 13.13–13.14, 13.14f
 - of photoemissive detectors, 24.34, 24.35f–24.39f, 24.41
 - quantum, 30.9
- Sensitometry variation, with film processing, 29.10, 29.10f
- Sensors:
- area arrays of, 32.24–32.32, 32.25t
 - about, 32.2
 - frame transfer CCD, 32.26–32.28, 32.27f, 32.28f
 - image area dimensions for, 32.25t
 - interline transfer CCD, 32.28–32.32, 32.29f–32.31f
 - linear image, 32.2, 32.21–32.24, 32.22f, 32.23f
 - metal-oxide-semiconductor, 32.25–32.26, 32.26f
 - image, 32.2–32.12, 32.3f, 32.21–32.34
 - antiblooming in, 32.9, 32.10f
 - color imaging with, 32.32–32.34, 32.33f, 32.34f
 - dark current in, 32.10–32.12, 32.11f
 - junction photodiodes, 32.3–32.6, 32.4f, 32.6f
 - linear arrays of, 32.21–32.24, 32.22f, 32.23f
 - MOS capacitors, 32.7–32.8
 - photoconductors, 32.8–32.9
 - pinned photodiodes, 32.8
 - LED detectors in, 17.34
 - staggered linear CCD, 32.23f, 32.24
 - time-delay-and-integrate linear, 32.23f, 32.24

- Separate confinement heterostructure waveguide, 19.24
- Separated absorption, grading, and multiplication layer APDs (SAGM APDs), 26.17, 26.18*f*, 26.20*f*
- Separated absorption and multiplication layer APDs (SAM APDs), 26.3, 26.17
- Servo stability, 22.8
- Servos, 22.5–22.12
 - Bode representation of, 22.5–22.6, 22.6*f*
 - closed-loop performance, 22.8
 - closed-loop stability issues, 22.8–22.12, 22.9*f*–22.11*f*
 - design with time delay, 22.19–22.20
 - measurement noise not a performance limit, 22.7–22.8
 - phase and amplitude responses vs. frequency, 22.6–22.7, 22.6*f*, 22.7*f*
- Seven-segment LED displays, 17.10, 17.11*f*
- Shapes, projected areas of common, 36.3–36.4, 36.3*f*, 36.3*t*
- Sharpness, of photographic images, 29.18, 29.19, 29.22
- Shells, mounting of, 6.11, 6.12*f*
- Shock specifications, for lenses, 4.10
- Short arc light sources, 15.34, 15.35*f*, 40.39
- Short-wavelength infrared (SWIR), 24.3, 33.3, 33.5
- Shot noise, 24.12, 27.3, 27.3*f*, 32.12
- SI units, 34.20, 37.7, 37.7*t*
- Side lighting, 40.12
- Sidecar TDI, 33.17, 33.17*f*
- Sign lighting, 40.71
- Signal analysis, 27.12–27.15
 - boxcar averaging, 27.13, 27.13*f*
 - categories of, 27.3
 - gated integration, 27.12–27.13, 27.13*af*
 - lock-in amplifiers, 27.13, 27.14, 27.14*f*
 - photon counting, 27.14
 - selection of technique, 27.14–27.15
 - transient photon counting, 27.14
 - of unmodulated sources, 27.12
- Signal detection, 27.1–27.12
 - and amplifiers, 27.10–27.12, 27.11*f*
 - and noise sources, 27.3–27.6, 27.3*f*, 27.5*f*, 27.6*f*
 - photomultiplier applications in, 27.6–27.10, 27.7*f*
 - technique selection for, 27.2–27.3, 27.2*f*
- Signal-to-noise ratio (S/N), 22.12, 27.1, 27.3, 29.1, 29.22–29.23, 33.2, 38.10
- Silicon (Si):
 - and diamond turning, 10.5
 - doped extrinsic, 33.7, 33.8*f*
 - Si:Ga infrared detectors, 24.95, 24.95*f*, 24.96, 24.96*f*
- Silicon avalanche photodiodes (APDs), 24.62–24.65, 24.63*f*–24.66*f*
- Silicon bolometers, 28.7*t*
- Silicon carbide (SiC) LED devices, 17.18
- Silicon carbide (SiC) substrate (for HB-LEDs), 18.2, 18.3
- Silicon carbide (SiC) UV detectors, 24.47, 24.47*f*
- Silicon CCDs (SCCDs), 33.11–33.13, 33.11*f*, 33.12*f*, 33.17
- Silicon nitride layer, 17.23
- Silicon oxide (SiO₂) passivation, 18.4
- Silicon oxynitride layer, 17.23
- Silicon (Si) photoconductors, 32.4*f*, 32.31, 32.32
- Silicon (Si) photodiodes, 38.9, 38.9*t*
 - avalanche, 24.62–24.65, 24.63*f*–24.66*f*
 - high-quality, 34.30
 - light-trap, 34.30
 - np*, 34.30
 - pin*, 24.55*f*–24.57*f*, 24.58–24.61, 24.59*f*, 24.60*f*
 - pn*, 24.52*f*, 24.55–24.58, 24.55*f*–24.59*f*, 34.30
 - self-calibration of, 34.29
 - UV- and blue-enhanced, 24.55*f*, 24.61–24.62, 24.61*f*, 24.62*f*
- Silicon (Si) photovoltaic detectors, 24.54–24.65, 24.55*f*, 24.56*f*
 - avalanche photodiodes, 24.62–24.65, 24.63*f*–24.66*f*
 - pin* photodiodes, 24.55*f*–24.57*f*, 24.58–24.61, 24.59*f*, 24.60*f*
 - pn* photodiodes, 24.52*f*, 24.55–24.58, 24.55*f*–24.59*f*
 - UV- and blue-enhanced photodiodes, 24.55*f*, 24.61–24.62, 24.61*f*, 24.62*f*
- Silicon-intensifier-target (SIT) vidicons, 31.8
- Silver, colloidal, 29.13
- Silver halide crystals, 30.1, 30.5–30.7, 30.6*f*
- Silver halide light detectors, 30.7–30.9, 30.8*f*
- Silver halide surfaces, 29.4, 30.14–30.15, 30.15*f*
- Simple lens, thermal focus shift of, 8.2–8.4, 8.3*t*, 8.4*t*

- Simulated annealing, 3.19
 Simultaneous measurement, in phase-shifting interferometry, 13.22
 Simultaneous multiple surfaces (SMSs), 39.17, 39.17f
 Sine plate, 12.10, 12.11f
 Single heterojunction LEDs, 17.12, 17.12f
 Single isolated pulses, 21.4
 Single material designs, 6.22, 6.23f
 Single monochromators, 38.15f, 38.16f
 Single optical pulse, 20.2–20.3, 20.3f, 20.6–20.7
 Single point diamond turning (SPDT), 6.1, 6.20
 Single quantum well (SQW) LEDs, 18.1
 Single-frequency lasers, 19.37f
 Single-lens arrays, 39.33–39.34, 39.34f
 Single-longitudinal-mode lasers, 19.38
 Single-pass photodetectors, 26.4f
 Single-shot f -to- $2f$ interferometers, 21.6
 Single-use cameras, 30.26
 Size-of-source effect, 34.33
 Skew ray limits, 39.20, 39.20f
 Skewness, 40.42
 Skot (unit), 36.7
 Skytubes, 40.50f
 Skywells, 40.50f
 Sliding prisms, 12.4, 12.4f
 Slit polynomials, 11.30, 11.35t, 11.36t
 Slot interrupters, 17.34
 Small-signal gain coefficient, 16.10
 Smoke detectors, 17.34
 Snubbing, 27.9–27.10
 Society of Automotive Engineers (SAE), 40.2, 40.63–40.64
 Society of Photooptical and Instrumentation Engineers (SPIE), 25.2, 39.12
 Soft mounting, 6.1, 6.4–6.5, 6.4f
 Software:
 for lighting simulation, 40.18–40.23
 for nonimaging modeling, 39.6–39.8
 for optical design (*see* Optical design software)
 for stray light suppression, 7.24–7.27
 Solar collection, 39.1
 Solar light pipes (SLPs), 40.49, 40.51f
 Solid angles, 34.9, 34.9f, 37.4, 39.5, 39.5f
 Solid lightpipes, 39.30–39.31
 Solid-state lasers, 16.12, 16.13, 16.17–16.18, 16.18f, 22.20–22.21
 Solid-state lighting, 18.4–18.5, 18.4f
 Solid-state photomultipliers (SSPM), 33.9
 Soliton solution, 20.5–20.9
 Solves (term), 3.5
 Source coupling, 40.41–40.42
 Source diameter, for fiber optics, 17.33
 Source modeling, 39.7
 Source modeling software, 40.19–40.20
 Source modulation, 27.3
 Spaciousness, perception of, 40.5
 Spatial dilution, 39.6
 Spatial noise, 33.26–33.27, 33.26f
 Special-purpose sources (of radiation), 15.53
 Specifications, optical (*see* Optical specifications)
 Speckle effects, 39.32
 Spectra Diode Labs, 19.29, 19.29t, 19.41
 Spectral (term), 35.2, 35.3
 Spectral D-double star, 24.12
 Spectral density, 22.3–22.4
 Spectral dependence (of radiometric quantities), 34.9–34.10
 Spectral detectivity, 24.12
 Spectral D-star, 24.12
 Spectral emittance, 35.7, 35.15
 Spectral errors, 34.36
 Spectral irradiance, 38.1–38.2, 38.11t, 38.13–38.16, 38.13f–38.16f
 Spectral irradiance calibration transfer devices, 34.31
 Spectral irradiance lamps, 15.11, 15.12, 15.13f
 Spectral lambertian source, 34.17
 Spectral lamps, 15.44, 15.45, 15.45f, 15.46f, 15.46t
 Spectral luminous efficiency, for photopic vision, 36.8, 36.8f, 36.9f, 36.16f, 37.2
 Spectral noise density, 22.3–22.5
 Spectral noise equivalent power, 24.12
 Spectral properties:
 of laser field, 23.28–23.31, 23.30f
 of micromaser field, 23.31–23.33
 of semiconductor lasers, 19.36–19.39, 19.37f, 19.38f
 Spectral radiance, 38.2, 38.11t, 38.13–38.16, 38.13f–38.16f
 Spectral radiance calibration transfer devices, 34.31
 Spectral radiance ribbon filament lamps, 15.11, 15.12f
 Spectral radiance units, 34.23–34.24

- Spectral reflectance, 35.4, 35.5, 38.2, 38.17–38.18, 38.18f
- Spectral response, of photodetectors, 24.18, 24.19f
- Spectral responsivity, 24.12, 38.3, 38.18–38.19
- Spectral scattering, 38.10
- Spectral sensitivity, of photographic films, 29.11, 29.11f
- Spectral sensitizers, photographic, 30.13–30.18
 about, 30.13–30.14, 30.14f
 color science of, 30.15–30.18, 30.16f, 30.17f
 photophysics of, on silver halide surfaces, 30.14–30.15, 30.15f
- Spectral transmittance, 35.10, 38.2, 38.17, 38.17f
- Spectrally stray radiation errors, 34.36
- Spectralon, 38.12, 38.13
- Spectrophotometers and spectrophotometry, 34.6, 35.8–35.9, 38.17, 38.17f
- Spectroradiometry, 38.1–38.19
 about, 38.1
 calculations for, 38.3–38.5
 calibration of, 38.11–38.13, 38.11t, 38.12f
 computer software for, 38.11
 detectors in, 38.8–38.10, 38.9t, 38.10t
 electronics of, 38.10
 errors in, 38.5–38.6
 figures of merit in, 38.5–38.6
 input (fore-) optics in, 38.7
 monochromators in, 38.7–38.8
 quantities used in, 38.1–38.2
 spectroradiometers, 38.18, 38.18f
 system designs in, 38.13–38.19
 spectral irradiance/radiance, 38.13–38.16, 38.13f–38.16f
 spectral reflectance, 38.17–38.18, 38.18f
 spectral responsivity, 38.18–38.19
 spectral transmittance, 38.17, 38.17f
- Specular reflectance, 35.10, 35.13
- Specular transmittance, 35.3, 35.9f
- Specular vanes, 7.17, 7.17f
- Speed:
 of LEDs, 17.33
 of photodetectors, 24.20, 24.21
 (See also Photographic film speed)
- SPEOS (optical software), 7.27
- Sphere(s):
 aberrations in, 11.30
 integrating (see Integrating spheres)
- Sphere(s) (*Cont.*):
 nonuniformities with integrating, 39.24–39.26, 39.24f, 39.25f
 projected area of, 36.3–36.4, 36.3f, 36.3t
- Spherical lambertian source, 34.17
- Spherical lenses, 39.8, 39.9f
- Spherical optics fabrication, 9.4–9.6
- Spherochromatism, 2.2
- Spherometers, 12.18–12.19, 12.18f, 12.19f, 12.19t
- Spline surfaces, 39.6
- Spokes, aperture, 3.20
- Spontaneous emission lasers (see Correlated emission lasers)
- Spontaneous emission rate, 23.8
- Spot-diagram analysis, 3.13–3.16
- Spring-loaded mountings, 6.13, 6.14f
- Square polynomials, 11.30, 11.31t–11.34t, 11.36t
- Stability:
 light, 30.10
 of photodetectors, 24.21, 24.21f
 of photoemissive detectors, 24.41
- Stabilization, light, 30.12–30.13 (See also Laser stabilization)
- Stable resonators, 16.23, 16.23f
- Stadia, 12.2–12.3, 12.3f
- Staggered linear CCD image sensor, 32.23f, 32.24
- Stagnation, 3.17
- Staircase APDs, 26.3
- Standard for the Exchange of Product model data (STEP), 40.19
- Standards:
 baseline, of radiation sources, 15.9f, 15.10f, 15.12f
 for detectors, 38.12–38.13
 for infrared radiometry, 15.11–15.12, 15.12f
 international, 4.11
 for length measurements, 12.2
 for lighting, 40.19
 for lighting system layout and simulation, 40.19
 for optical image quality, 4.6
 published, 4.10
 of reflectance, 35.14t
 search engine for, 4.11
 of spectral transmittance, 35.10
 for vehicular lighting, 40.63–40.64, 40.66f, 40.66t
 working, of radiation sources, 15.9–15.13, 15.10f, 15.12f, 15.13f
- Star concentrators, 39.20, 39.21

- Staring arrays, 33.6–33.7, 33.6f, 33.14
 Staring FPAs, 33.16–33.17, 33.29, 33.29t
 Steady-state pulse train, 20.5–20.9
 Stefan-Boltzmann law, 34.24, 37.11
 Steradian (sr), 36.3, 37.4, 37.4f
 Sterilamps, 15.35, 15.36f
 Stilb (unit), 34.43, 36.7, 36.8t
 Stimulated absorption, 16.7–16.8, 16.8f
 Stimulated emission, 16.2, 16.7–16.9, 16.8f, 23.8
 Stop lamps, 40.64f, 40.67, 40.68t, 40.69f
 Stop shifting, 2.5, 2.6f
 Stops:
 aperture, 34.18, 34.19f
 field, 34.18–34.19, 34.19f
 Straddling springs, 6.13, 6.14f
 Straightness measurement, 12.10
 Strained QW lasers, 19.15–19.17, 19.16f
 Stray capacitance noise, 27.5, 27.6f
 Stray light, 29.15–29.16, 29.15f
 Stray light suppression, 7.1–7.32
 about, 7.1–7.2
 aperture placement in, 7.5–7.10
 aperture stops, 7.6–7.7, 7.7f, 7.8f
 field stops, 7.7, 7.8f, 7.9f
 Lyot stops, 7.8–7.10, 7.8f–7.11f
 baffles in, 7.10, 7.11
 and BRDF characteristics, 7.23, 7.24f
 Cassegrain design with aperture stop at
 primary (example), 7.3f
 contamination levels in, 7.18–7.19, 7.18t,
 7.19f–7.21f
 evaluation methods for, 7.27–7.29, 7.29f
 illuminated objects in, 7.5, 7.5f, 7.6f
 imaged critical objects in, 7.4, 7.5f
 information sources on, 7.31–7.32
 issues with, 7.30–7.31
 and point source transmittance definitions,
 7.22–7.23
 radiation transfer equation for, 7.21–7.22
 real-space critical objects in, 7.2–7.4, 7.3f, 7.4f
 software for, 7.24–7.27
 and stray radiation paths, 7.22
 strut design in, 7.20, 7.21, 7.21f
 and surface scattering characteristics, 7.23
 vane spacing and depth in, 7.13–7.17
 angle considerations, 7.13–7.16, 7.14f, 7.15f
 bevel placement, 7.13, 7.14f
 depth considerations, 7.16, 7.16f, 7.17f
 specular vanes, 7.17, 7.17f
 vanes in, 7.11–7.12, 7.12f, 7.13f
 Stray radiation paths, 7.9, 7.22
 Street lighting, 40.69–40.71, 40.70t, 40.71t
 Stress tolerance, 6.3
 Stretched segment displays, 17.30–17.31,
 17.30f, 17.31f
 Strip mirror integrator (SMI), 39.40
 Strong field approximation, 21.3
 Strong VW reflectometer, 35.10f
 Strut design (in stray light suppression), 7.20,
 7.21, 7.21f
 Subjective tone reproduction, 29.16
 Submillimeter (SubMM) radiation, 24.3
 Subminiature lamps, 15.53
 Sub-Nyquist interferometry, 13.27
 Substrate(s):
 absorbing, 17.7, 17.7t
 for HB-LEDs, 18.2–18.3, 18.2f
 LED, 17.20–17.21, 17.20t
 mirror, 6.17–6.18, 6.17f, 6.18f
 transparent, 17.7, 17.7t
 Suncatchers, 40.48, 40.50f
 Superconducting bolometers, 28.5, 28.7t
 Superposition (of uniformity), 39.2, 39.32,
 39.33f
 Superposition-of-sources nonlinearity mea-
 surement, 34.33
 Supersensitizers, 30.14, 30.15, 30.15f
 Support wires, light bulb, 40.29f, 40.30
 Surface emitting lasers (SELs) (SLASERs),
 19.39–19.41, 19.39f, 19.40f, 19.42f, 19.43t,
 25.15
 Surface emitting LEDs (SLEDs), 25.15
 Surface finishing, of diamond-turned optics,
 10.9–10.11, 10.9f–10.11f
 Surface generation current, 32.10, 32.11,
 32.11f
 Surface measurement systems, 40.53, 40.54
 Surface mount device (SMD) package, 18.5f
 Surface mount LEDs (SMDs), 40.37
 Surface profilometers, 9.6
 Surface scattering, 7.23
 Surface-channel CCDs, 32.14
 Surface-channel MOS capacitors, 32.4f, 32.7
 Surfaces, modeling of, 40.17
 Suspended luminaires, 40.13f
 Synchronotron radiation, 34.26–34.27
 Synchronous pumping, 16.29
 System specifications, for lenses, 4.3
 Système International (SI), 12.2, 36.2, 37.3
 (See also SI units)

- Taillights, 40.21, 40.22*f*, 40.64*f*, 40.67, 40.68*t*, 40.69*f*
- Tailored (T) reflectors, 39.37–39.39, 39.38*f*
- Tailoring (of uniformity), 39.2
- Talbot autoimages, 12.23, 12.24
- Talbot's law, 34.33–34.34
- Tandem-lens arrays, 39.34, 39.35*f*–39.37*f*
- Tapered lightpipes, 39.12–39.13, 39.13*f*, 39.31–39.32, 39.31*f*
- Task lighting, 40.12, 40.14
- Taylor-Hobson Form TalySurf, 9.6
- Technical specifications, 4.2
- Tehis method, 40.53, 40.54
- Telecentric distribution, 39.18, 39.18*f*
- Telescope(s):
- astronomical, 1.7*f*
 - Galilean, 1.7*f*
 - Hubble, 11.4, 13.24
 - Keck, 11.4
 - reflecting, 11.4
 - unit magnification Galilean, 12.4, 12.4*f*
- Telescoping input optics, 38.7
- Temperature:
- color, 37.4*t*, 37.6–37.7
 - correlated color, 37.7, 38.5
 - distribution, 37.7
 - and mounted optics, 6.21–6.24, 6.22*f*–6.24*f*
 - radiance, 37.4*t*, 37.6
- Temperature control, of PZT transducers, 22.19
- Temperature noise, 24.12
- Temperature specifications, for lenses, 4.10
- Temperature-dependence effects, 34.36–34.37
- Templates (for curvature measurement), 12.17
- Tensile-strained QW lasers, 19.16, 19.16*f*, 19.17
- Test plates (for curvature measurement), 12.17
- Testing, 13.1–13.27
- aspherical wavefront measurement, 13.23–13.27
 - holographic compensators, 13.25, 13.25*f*, 13.26*f*
 - infrared interferometry, 13.25
 - Moiré tests, 13.26–13.27
 - refractive or reflective compensators, 13.24, 13.24*f*, 13.25
 - sub-Nyquist interferometry, 13.27
 - two-wavelength interferometry, 13.25, 13.26
 - wavefront stitching, 13.27, 13.27*f*
 - computer-generated holograms in (see Computer-generated holograms)
 - Testing (*Cont.*):
 - of convex surfaces, 14.5
 - interferogram evaluation, 13.14–13.18
 - direct interferometry, 13.17–13.18
 - fixed interferograms, 13.14–13.15
 - Fourier analysis of interferograms, 13.16–13.17, 13.17*f*
 - global and local interpolation of interferograms, 13.15–13.16
 - interferometric, 13.7–13.12
 - common path interferometer, 13.9, 13.11*f*
 - Fizeau interferometer, 13.8–13.9, 13.9*f*, 13.10*f*
 - lateral shearing interferometers, 13.9–13.12, 13.11*f*, 13.12*f*
 - multiple-pass interferometers, 13.13
 - multiple-reflection interferometers, 13.13
 - radial, rotational, and reversal shearing interferometers, 13.12, 13.13*f*
 - sensitivity of interferometers, 13.13–13.14, 13.14*f*
 - Twyman-Green interferometer, 13.7–13.8, 13.7*f*, 13.8*f*
 - Zernike phase-contrast method applied to interferometers, 13.13–13.14, 13.14*f*
 - noninterferometric, 13.1–13.7
 - Foucault test, 13.2–13.3, 13.2*f*, 13.3*f*
 - Hartmann test, 13.4–13.6, 13.5*f*
 - Hartmann-Shack test, 13.6–13.7, 13.6*f*
 - Ronchi test, 13.3–13.4, 13.3*f*, 13.4*f*
 - phase-shifting interferometry, 13.18–13.23, 13.18*f*–13.20*f*
 - heterodyne interferometer, 13.22
 - integrating bucket method, 13.21, 13.21*f*
 - phase errors, 13.22
 - phase stepping, 13.20, 13.20*f*
 - phase-lock method, 13.23, 13.23*f*
 - simultaneous measurement, 13.22
 - two steps plus one method, 13.21, 13.22
 - in wafer processing, 17.24
 - Thef-number, 38.8
 - Theodolites, 12.13
 - Thermal arrays, 28.7–28.12
 - about, 28.7–28.8
 - noise equivalent temperature difference in, 28.8–28.9
 - pyroelectric hybrid, 28.11–28.12, 28.11*f*, 28.12*f*
 - resistive bolometer, 28.10–28.11, 28.10*f*
 - theoretical limits of, 28.9–28.10, 28.9*f*
 - thermoelectric, 28.12, 28.12*f*

- Thermal circuit theory, 28.2
- Thermal coefficient of resistance (TCR), 33.2, 33.14
- Thermal compensation, 8.1–8.15
 about, 8.2
 and effect of thermal gradients, 8.6–8.7
 and homogeneous thermal effects, 8.2–8.5, 8.3*t*, 8.4*t*, 8.5*f*
 intrinsic athermalization, 8.7–8.8, 8.7*f*
 mechanical athermalization, 8.8–8.12, 8.8*f*–8.12*f*
 optical athermalization, 8.12–8.15, 8.13*t*–8.15*t*
 tolerable homogeneous temperature change, 8.5–8.6, 8.6*f*
- Thermal defocus, of compound lens, 8.4, 8.5*f*
- Thermal detector(s), 24.4–24.6, 24.4*f*, 28.1–28.12, 38.9, 38.9*t*
 arrays of, 28.7–28.12
 about, 28.7–28.8
 noise equivalent temperature difference, 28.8–28.9
 pyroelectric hybrid arrays, 28.11–28.12, 28.11*f*, 28.12*f*
 resistive bolometer arrays, 28.10–28.11, 28.10*f*
 theoretical limits, 28.9–28.10, 28.9*f*
 thermoelectric arrays, 28.12, 28.12*f*
- bolometer, 24.5*f*, 28.3–28.5, 28.4*f*
- Golay cell, 28.6
- ideal, 28.2–28.3, 28.3*f*
- performance/sensitivity of, 24.17, 24.18*f*
- properties of, 28.7, 28.7*t*
- pyroelectric, 24.6, 24.6*f*, 28.7
- and thermal circuit theory, 28.2
- thermistor, 24.5
- thermocouple, 28.4
- thermopile, 24.5*f*, 28.4–28.5
- Thermal expansion, 33.14
- Thermal fatigue, 17.25
- Thermal focus shift, 8.2–8.4, 8.3*t*, 8.4*t*
- Thermal gradients, effect of, 8.6–8.7
- Thermal infrared detectors, 33.7, 33.8*f*
- Thermal noise, 24.13, 27.4, 32.20
- Thermal properties, of high-power lasers, 19.26, 19.27*f*
- Thermal stability, of plastic packaging materials, 17.26
- Thermistor bolometers, 24.24–24.25, 24.24*f*, 24.25*f*, 28.7*t*
- Thermistors, 24.5
- Thermocouple junctions, noise from, 27.6, 27.6*f*
- Thermocouples, 24.5, 28.7*t*
 about, 28.1
 manufacturers' specifications for, 24.22–24.23, 24.22*f*
 as thermal detectors, 28.4
- Thermoelectric arrays, 28.12, 28.12*f*
- Thermopiles, 24.5, 28.7*t*
 defined, 24.13
 manufacturers' specifications for, 24.23–24.24, 24.23*f*
 as thermal detectors, 28.4–28.5
- θ_1/θ_2 concentrators, 39.18–39.20, 39.19*f*
- Thick window chips, 17.7, 17.7*t*
- Thin doublet, 1.15–1.16
- Thin lenses, 1.5
- Thin teflon diffusers, 38.15*f*
- Thin-disk lasers, 16.18
- 35-mm photographic films, 30.21, 30.25
- Thoria (in incandescent lights), 40.27
- Threaded retaining rings, 6.3, 6.3*f*
- 3D concentrators, 2D vs., 39.20–39.21, 39.20*f*, 39.21*f*
- Three-chip color systems, 32.32, 32.33*f*
- Three-material athermal solutions, 8.14, 8.14*t*, 8.15*t*
- Three-phase CCDs, 32.15, 32.16*f*
- Three-step rescattering model, 21.3
- Threshold carrier density, 19.12, 19.12*f*, 19.13, 19.13*f*
- Threshold current, 19.6, 19.6*f*
- Threshold modal gain, 19.12, 19.12*f*, 19.13, 19.13*f*
- Threshold voltage, 25.11
- Tightly toleranced assembly, 6.7, 6.7*f*
- Time delay integration (TDI), 33.4
- Time delay integration (TDI) linear sensors, 32.23*f*, 32.24
- Time delay integration (TDI) scanning FPAs, 33.17, 33.17*f*
- Time evolution of the field, 23.15–23.17, 23.15*f*
- Time-averaged color mixing, 40.8
- Time-based measurement, 12.2, 12.4, 12.5, 12.6*f*
- Time-dependent error, 34.35
- Time-of-flight distance measurement, 12.4, 12.5

- Titanium oxide (TiO_2) UV detectors, 24.47, 24.48*f*
- Titanium-doped sapphire ($\text{Ti:Al}_2\text{O}_3$) lasers, 16.34, 16.34*f*
- Titanium-doped sapphire ($\text{Ti:Al}_2\text{O}_3$) ring lasers, 20.16–20.17, 20.16*f*, 20.17*f*
- Tolerance budgeting, 5.3
- Tolerance verification, 5.3
- Tolerances, 5.2–5.8
- assembly, 5.8
 - basis for, 5.2–5.3
 - boresight, 5.8
 - budgeting of, 5.3
 - distortion, 5.8
 - optical vs. mechanical, 5.2
 - verification of, 5.3
 - wavefront, 5.3–5.7, 5.4*f*, 5.5*f*, 5.5*t*, 5.6*t*, 5.7*f*
- Tolerancing, 5.8–5.11
- and aberration balancing, 11.35, 11.36
 - about, 5.1–5.2
 - and material properties, 5.9
 - measurement practices for, 5.8–5.9
 - and optimization, 3.20–3.21
 - problems in, 5.11
 - procedures for, 5.9–5.10
 - shop practices for, 5.8
- Tone reproduction, 29.16–29.17, 29.17*f*
- Total flux into a hemisphere, 34.15
- Total hemispherical emittance, 35.15, 35.15*f*
- Total internal reflection (TIR), 39.12, 39.17, 40.41
- Total internal reflection (TIR) Fresnel lenses, 39.10
- Total luminous flux, 37.4*t*, 37.6
- Total radiant flux, 37.4*t*, 37.6
- Total transmittance, 35.3, 35.9*f*
- Traceability:
- of absolute measurements, 34.21
 - errors in, 34.28
- TracePro (optical software), 7.27
- Transconductance amplifiers, 27.11–27.12, 27.11*f*
- Transducer resonance, 22.8, 22.11–22.12
- Transducers, 22.17–22.20
- Transformers, in voltage amplifiers, 27.11
- Transient photon counting, 27.14
- Transmission, 4.7
- actual/idealized, 35.2*f*
 - defined, 35.3
- Transmission density, of photographic films, 29.6–29.7, 29.7*f*
- Transmissive sensors, 17.34
- Transmittance, 35.3
- measurement of, 35.8–35.10, 35.9*f*
 - and reflectance/absorptance, 35.7, 35.8, 35.8*t*
 - spectral, 38.2, 38.17, 38.17*f*
- Transmitter speed, for fiber optics, 17.33
- Transparency, 39.23, 40.5
- Transparency point, 19.5
- Transparent substrate (TS) chips, 17.7, 17.7*t*
- Transportation lighting, 40.63–40.71
- roadway lighting, 40.67, 40.69–40.71, 40.70*t*, 40.71*t*
 - vehicular lighting, 40.63–40.67, 40.64*f*, 40.65*t*, 40.66*f*, 40.66*t*, 40.68*t*, 40.69*f*
- Transverse electromagnetic mode (TEM), 16.21–16.23, 16.22*f*
- Transverse junction stripe (TJS) lasers, 19.8, 19.9*f*, 19.23–19.24, 19.36*f*
- Transverse laser modes, 16.21–16.23, 16.21*f*–16.23*f*
- Transverse ray plots, 2.2–2.4, 3.13
- Traveling microscopes, 12.20, 12.21, 12.21*f*
- Traveling wave photodetectors, 26.4*f*, 26.5, 26.14*f*
- Treaty of the Meter of 1875, 34.20, 36.2
- Triphosphors, 40.31, 40.32*f*
- Triplet lens, air-spaced, 6.21, 6.22*f*
- Tristimulus values, 38.3–38.4
- Troffers, fluorescent luminaire, 40.47
- Troland (unit), 34.41–34.42, 37.7, 37.8
- Trough reflectors, 40.46*f*, 40.47
- Trumpet (term), 39.15, 39.16*f*, 39.17
- Tubular PZT transducers, 22.17–22.18
- Tungsten:
- in HID lamps, 40.35
 - in incandescent lights, 40.25, 40.27, 40.29
- Tungsten lamps, 15.13, 40.26*t*, 40.28*f*
- Tungsten-arc lamps, 15.47–15.48, 15.48*f*, 15.49*f*
- Tungsten-filament lamps, 15.11, 15.12, 15.13*f*, 15.19, 15.20, 15.20*f*–15.22*f*, 34.31
- Tungsten-halogen lamps, 15.11, 15.12, 15.13*f*, 40.25*t*, 40.26*t*, 40.30
- Tunnel diagram (*see* Williamson construction)
- Tunnel lighting, 40.71
- Tunneling current, 25.8
- Twin-channel lasers (TCLs), 19.27
- Twin-channel substrate mesa (TCSM) lasers, 19.20*t*, 19.21*f*, 19.23

- Twin-ridge structure (TRS) lasers, **19.19**,
19.20t, **19.21f**, **19.22–19.23**
- 2D (term), **39.4**
- 2D concentrators, 3D vs., **39.20–39.21**, **39.20f**,
39.21f
- 2D high-power laser arrays, **19.29–19.30**,
19.29t, **19.30f**
- Two-color gating, **21.7**
- Two-component systems, first-order layout for,
1.5–1.7
- Two-interference pattern distance-measuring
interferometer, **12.7**, **12.7f**
- Two-mirror imaging system, **39.17**
- Two-phase CCDs, **32.15–32.16**, **32.16f**
- Two-stage baffle, **7.10**
- Two-step rescattering model, **21.3**
- Two-steps-plus-one phase shifting,
13.21, **13.22**
- Two-wavelength interferometry, **13.25**, **13.26**
- Twyman-Green interferograms, **13.10f**, **13.18f**
- Twyman-Green interferometers, **13.7–13.8**,
13.7f, **13.8f**
- Type A errors (in absolute measurement),
34.21–34.23
- Type B errors and error sources (in absolute
measurement), **34.32–34.37**
defined, **34.21–34.23**
nonideal aperture, **34.35–34.36**, **34.35f**
nonlinearity of detector, **34.34–34.35**
nonuniformity, **34.35**
offset subtraction, **34.32–34.33**
polarization effects, **34.33**
scattered radiation effect, **34.33**
size-of-source effect, **34.33**
spectral errors, **34.36**
temperature-dependence effects, **34.36–34.37**
time-dependent error, **34.35**
- Ultrashort cavity microlasers, **19.39**
- Ultrashort optics, **20.1–20.28**
about, **20.1–20.2**
cavities with two circulating pulses,
20.15–20.22
linear lasers, **20.18–20.19**, **20.19f**
optical parametric oscillators, **20.20–20.22**,
20.20f, **20.21f**
ring dye lasers, **20.15–20.16**
ring lasers, **20.17–20.18**, **20.17f**
Ti:sapphire ring lasers, **20.16–20.17**,
20.16f, **20.17f**
- Ultrashort optics (*Cont.*):
coupling of circulating pulses, **20.12–20.15**,
20.12f, **20.15f**
optical pulses and pulse trains, **20.2–20.9**
single optical pulse, **20.2–20.3**, **20.3f**
soliton solution and steady-state pulse
train, **20.5–20.9**
train of pulses, **20.3–20.5**, **20.4f**, **20.5f**
and quantum mechanical two-level system,
20.22–20.28
coherent interaction, **20.22–20.23**
experimental demonstration, **20.24–20.27**,
20.25f–20.27f
impact of analogy, **20.27–20.28**
laser as two-level system, **20.23–20.24**,
20.25t
Rabi cycling, **20.26–20.27**, **20.26f**, **20.27f**
steady-state pulse, **20.9–20.12**, **20.11f**
- Ultrasonic-assisted machining, **10.5**
- Ultraviolet (UV) detectors:
silicon carbide, **24.47**, **24.47f**
TiO₂, **24.47**, **24.48f**
- Ultraviolet (UV) enhanced photodiodes,
24.55f, **24.61–24.62**, **24.61f**, **24.62f**
- Ultraviolet (UV) filters, **40.12**
- Ultraviolet (UV) radiation, **34.6**
and color film, **30.3**
far, **15.12**, **15.13**
spectrum of, **25.2**
vacuum, **24.3**
- Uncrossed reflectors, **39.38**, **39.38f**
- Unified Glare Rating (UGR), **40.10–40.11**, **40.11t**
- Uniform illumination, of nonimaging optics,
39.22–39.41
with classic projection systems, **39.23–39.24**,
39.23f
faceted structures in, **39.39–39.41**, **39.39f**,
39.40f
integrating cavities in, **39.24–39.27**, **39.24f**,
39.25f, **39.27f**
lens arrays in, **39.32–39.37**, **39.33f–39.37f**
lightpipes in, **39.13f**, **39.27–39.32**,
39.28f–39.30f
tailored reflectors, **39.37–39.39**, **39.38f**
- Uniformity:
angular, **39.31**
control of, **39.1–39.2**
of luminance/illuminance, **40.7**, **40.13f**
of photodetectors, **24.20**
and visual discomfort, **40.9**

- Unit conversions:
 - for English and SI units, 37.7, 37.7t
 - for illuminance, 36.7t, 36.8t
 - for photometric and radiometric quantities, 36.11–36.14, 36.12f–36.14f
- Unit magnification Galilean telescope, 12.4, 12.4f
- Unlit-appearance modeling, 40.21
- Unmodulated signal sources, 27.12
- Unstable resonators, 16.25–16.26, 16.26f
- Unstrained QW lasers, 19.15–19.16, 19.16f
- Uplight, 40.43, 40.44f, 40.45
- U.S. Air Force three-bar target, 4.6
- Useful life period, of LEDs, 17.26, 17.26f
- Uviarc, 15.28–15.29, 15.29f, 15.30f

- Vacuum, laser gain media in, 16.36–16.37, 16.37f
- Vacuum lamps, 34.31
- Vacuum ultraviolet (VUV) radiation, 24.3
- Valence band, 17.3, 17.4, 17.4f
- Valence lighting, 40.13f
- Vanes (in stray light suppression), 7.11–7.17
 - defined, 7.11–7.12, 7.12f, 7.13f
 - placement design for, 7.12f
 - and scatter path, 7.13f
 - spacing and depth of, 7.13–7.17, 7.14f–7.17f
- Vapor exposure, in LED packaging, 17.26
- Vapor phase epitaxy (VPE), 17.21, 17.22
- Variable temperature blackbody, 15.10f
- Variable-orientation mirrors, 6.17
- Varifocal systems, first-order layout for, 1.11–1.12
- Vector flux, 39.21–39.22
- Vehicular lighting, 40.63–40.67, 40.64f, 40.65t, 40.66f, 40.66t, 40.68t, 40.69f
- Veiling reflections, 40.12, 40.14
- Verification (of tolerance), 5.3
- Vertical antiblooming, 32.9, 32.10f
- Vertical Bridgeman technique, 17.21
- Vertical cavity lasers, 19.41, 19.42f, 19.43t
- Vertical cavity semiconductor lasers, 16.36
- Vertical cavity surface-emitting lasers (VCSELs), 16.36
- Vertical illuminance, 40.7
- Vertically integrated photodiode (VIP) FPAs, 33.10
- Vertically illuminated *pin* photodiodes, 26.3, 26.4f, 26.5, 26.10, 26.12–26.13, 26.12f

- Very-long-wavelength infrared (VLWIR) radiation, 24.3
- Very-long-wavelength semiconductor lasers, 19.7–19.8
- Vibration specifications, for lenses, 4.10
- Vibration-resistant optical reference cavity, 22.16, 22.17f
- Vignetting, 3.4, 34.19
- Virtual phase CCDs, 32.16–32.17, 32.16f
- Visible array detectors, 32.1–32.34
 - about, 32.2
 - image sensing elements of, 32.2–32.12, 32.3f
 - antiblooming, 32.9, 32.10f
 - dark current, 32.10–32.12, 32.11f
 - junction photodiode, 32.3–32.6, 32.4f, 32.6f
 - MOS capacitor, 32.7–32.8
 - photoconductor, 32.8–32.9
 - pinned photodiode, 32.8
 - readout elements of, 32.12–32.21
 - CCD, 32.12–32.20, 32.13f, 32.15f–32.18f
 - MOS, 32.20–32.21
 - sensor architectures of, 32.21–32.34
 - area image sensor arrays, 32.24–32.32, 32.25t, 32.26f–32.31f
 - color imaging, 32.32–32.34, 32.33f, 32.34f
 - linear image sensor arrays, 32.21–32.24, 32.22f, 32.23f
- Visible light photon counters (VLPCs), 33.9
- Visible (VIS) radiation, 24.3, 25.2
- Vision, 40.3–40.6, 40.9
 - biology of, 40.3–40.4
 - and perception, 40.4–40.5
 - photopic/scotopic/mesopic, 34.37–34.39, 37.2
 - (See also Human eye)
- Visual clarity, perception of, 40.5
- Visual discomfort, 40.9–40.12, 40.11t
- Visual discomfort probability (VCP), 40.10
- Visual photometry, 36.4
- Visual science, 34.37
- Vivid color (VC) film, 30.27
- Voltage amplifiers, 27.10–27.11

- Wafer processing, 17.23–17.25
- Wall slot lighting, 40.13f
- Wall-grazing illumination, 40.13f
- Wall-washing illumination, 40.13f
- Watanabe, F., 39.33
- Watt (unit), 39.2t

- Wave modulation distance meter, 12.5, 12.6f
- Waveband materials, 8.3t, 8.4t
- Waveband structure of semiconductors, 17.3–17.6, 17.3f–17.5f
- Wavefront error (*W*), 4.1, 4.3, 4.7, 4.8
- Wavefront measurement, aspherical (*see* Aspherical wavefront measurement)
- Wavefront stitching, 13.27, 13.27f
- Wavefront tolerancing, 5.3–5.7, 5.4f, 5.5f, 5.5t, 5.6t, 5.7f
- Wavefronts, from lenses, 4.3–4.5, 4.5t
- Waveguide photodetectors, 26.4f, 26.5
- Waveguide *pin* photodiodes, 26.13–26.14, 26.14f
- Wavelength, in fiber optics, 17.33–17.34
- Wavelength errors, 34.36
- Wearout period, 17.26, 17.26f
- Weighting functions, 36.17
- Well capacity, 25.11
- Welsbach mantle, 15.17, 15.18
- Whiffletrees (lever mechanisms), 6.19
- White light, 18.4–18.5, 18.4f, 40.7, 40.8, 40.24
- White surfaces, reflectivity of, 17.31
- White-light LEDs, 40.37, 40.38
- WI 9 lamps, 15.21f
- WI 14 lamps, 15.21f
- WI 16/G lamps, 15.21f, 15.22f
- WI 17/G lamps, 15.22f
- WI 40/G lamps, 15.22f
- WI 41/G lamps, 15.22f
- Wiener spectrum, 29.21
- Wien's displacement law, 15.7, 34.23, 34.24, 37.11
- Williamson construction, 39.12–39.13, 39.13f, 39.28, 39.29f, 39.31, 39.32
- Window/photocathode assemblies, of image intensifiers, 31.10–31.12, 31.11f, 31.12f
- Windows:
and daylight sources, 40.41, 40.47, 40.48, 40.49f, 40.50f
mounting of optical, 6.11, 6.11f, 6.12f
- Wire-wound thermopile arrays, 24.23
- Work function (of photons), 25.2
- Xenon lamps, 15.34f, 15.35f, 40.31, 40.35f
- X-ray lasers, 16.31
- X-Y addressing, 33.16
- Y-coupled junctions, 19.27, 19.29
- Yellow filter dyes, 30.4
- Yellow light, 29.13, 29.13f
- Yttrium aluminum garnet (YAG) phosphor, 18.4
- ZEMAX (optical software), 7.26–7.27
- Zernike phase-contrast test, 13.13–13.14, 13.14f
- Zernike polynomials, 5.9
annular, 11.13–11.21, 11.14f, 11.17t–11.21t
circle, 11.4, 11.6–11.12, 11.8t–11.9t, 11.9f–11.11f, 11.12t, 11.39
- Zerodur prisms, 6.16, 6.16f
- Zinc, 17.23
- Zinc diffusion, 17.9–17.10, 17.10f
- Zinc doping, 17.20
- Zinc oxide (ZnO) doped GaP, 17.16, 17.21–17.22
- Zinc selenide (ZnSe) LED devices, 17.19
- Zinc-doped germanium (Ge:Zn) detectors, 24.84f, 24.98–24.100, 24.99f
- Zirconium arc lamps, 15.47, 15.48f
- Zonal cavity lighting simulation, 40.17
- Zoom systems, 1.11–1.12, 3.20
- Z-system (eccentric pupil design), 7.11, 7.12f, 7.15–7.17, 7.15f, 7.17f, 7.19, 7.21f

COLOR PLATES

DO NOT DUPLICATE

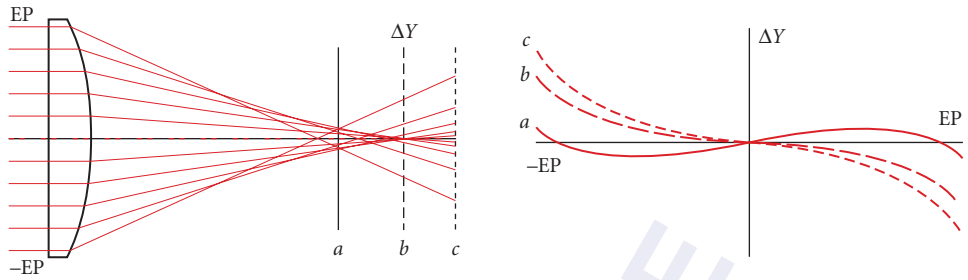


FIGURE 2.1 (Left) Rays exiting a lens are intercepted at three evaluation planes. (Right) Ray intercept curves plotted for the evaluation planes: (a) at the point of minimum ray error (circle of least confusion); (b) at the paraxial image plane; and (c) outside the paraxial image plane.

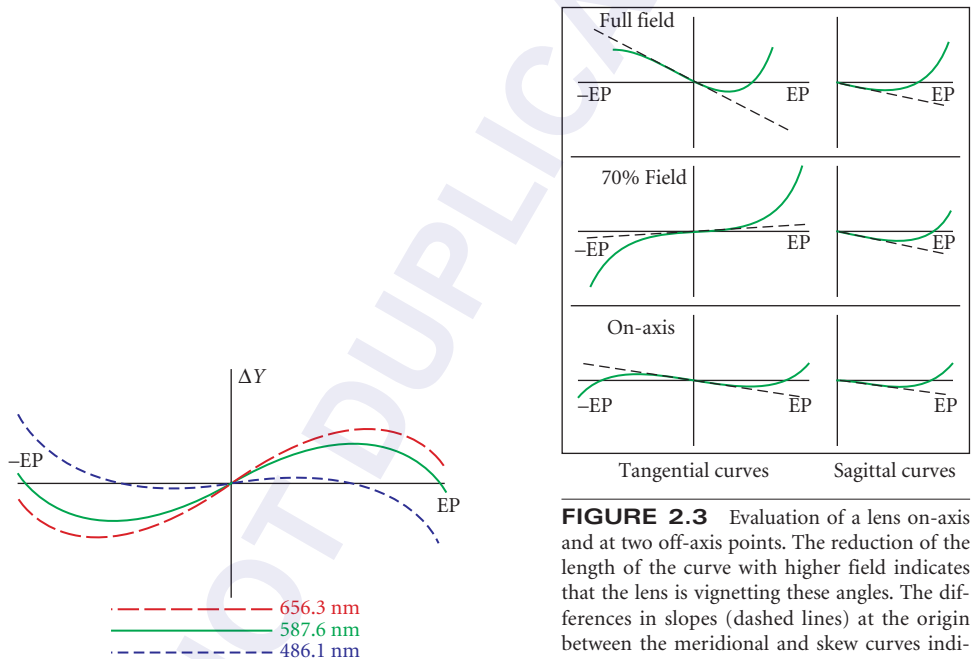


FIGURE 2.2 Meridional ray intercept curves of a lens with spherical aberration plotted for three colors.

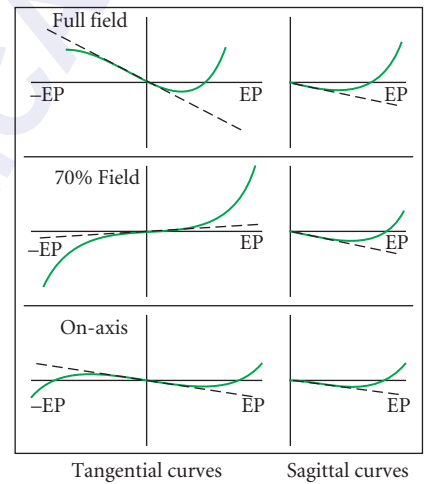


FIGURE 2.3 Evaluation of a lens on-axis and at two off-axis points. The reduction of the length of the curve with higher field indicates that the lens is vignetting these angles. The differences in slopes (dashed lines) at the origin between the meridional and skew curves indicate that the lens has astigmatism at these field angles. The variation in the slopes with field indicates the presence of field curvature.

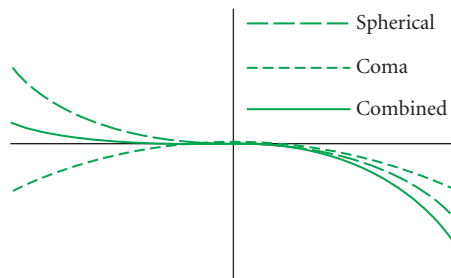


FIGURE 2.4 Ray intercept curve showing coma combined with spherical aberration.

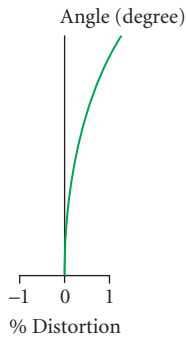


FIGURE 2.5 Field curve: distortion plot. The percentage distortion is plotted as a function of field angle. Note that the axis of the dependent variable is the horizontal axis.

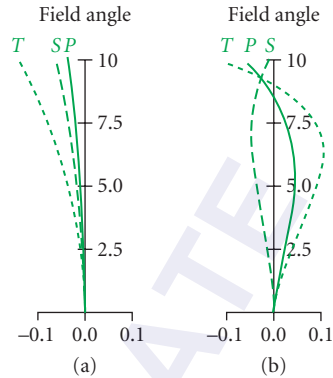


FIGURE 2.6 Field curve: field curvature plot. The locations of the tangential T and sagittal S foci are plotted for a full range of field angles. The Petzval surface P is also plotted. The tangential surface is always three times farther from the Petzval surface than from the sagittal surface: (a) an uncorrected system and (b) a corrected system.

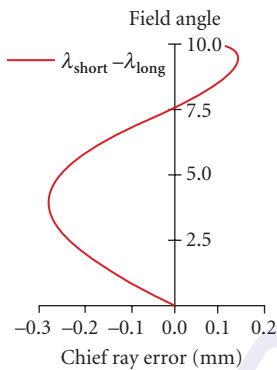


FIGURE 2.7 Field curve: lateral color plot. A plot of the transverse ray error between red and blue chief ray heights in the image plane for a full range of field angles. Here the distance along the horizontal axis is the color error in the image plane.

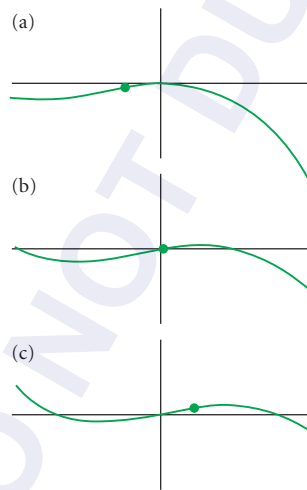


FIGURE 2.8 The effect of stop shifting on the meridional ray intercept curves of a double Gauss lens. (a) Stop located in front of the normal centrally located stop. (b) Stop at the normal stop position. (c) Stop behind the normal stop position. The dot locates the point on the curve where the origin is located for case (b).



FIGURE 40.3 Accent lighting.



FIGURE 40.4 Wall sconces for providing ambient lighting and the much needed vertical illumination in various situations.

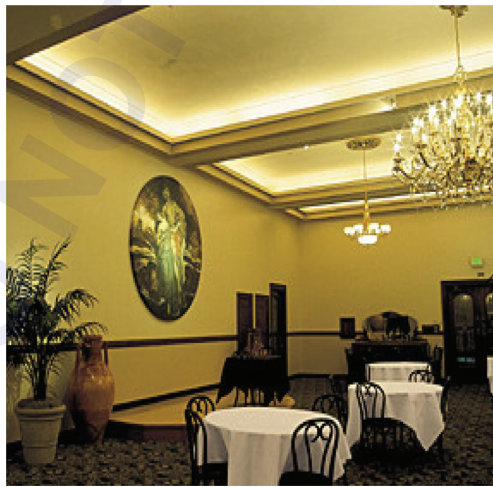


FIGURE 40.5 Indirect lighting with cove lighting in a restaurant using light strips. The chandelier provides the decorative lighting without significantly contributing to any other lighting function.

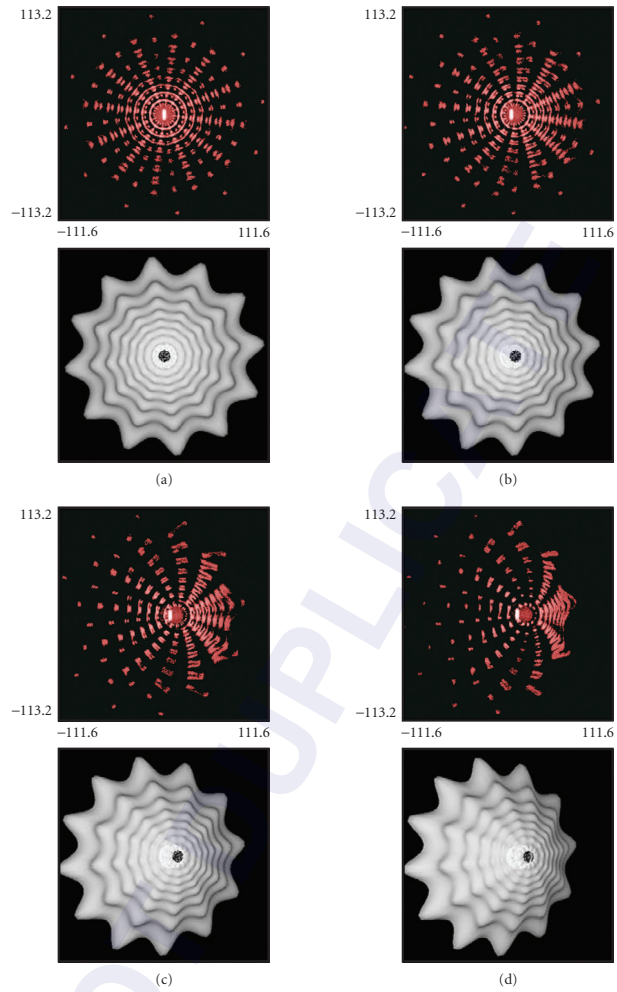


FIGURE 40.8 Views of the lit appearance (upper) of a star-shaped taillight (lower) at four horizontal angles of (a) 0°; (b) 10°; (c) 20°; and (d) 30°.

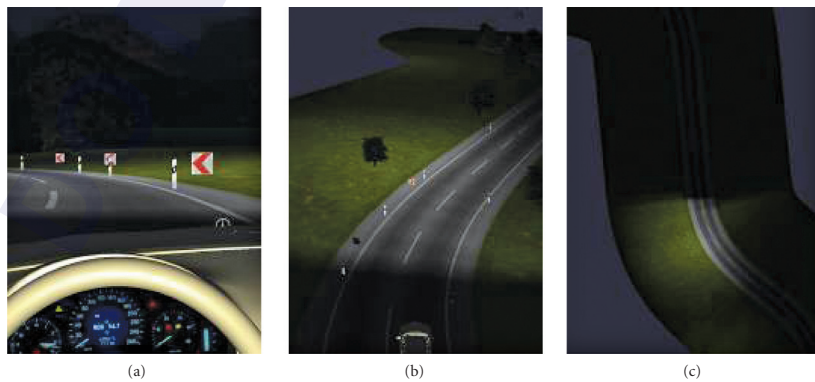


FIGURE 40.9 Three perspectives of lit-scene renderings from a low-beam headlamp: (a) driver's view; (b) 20 m above and behind automobile; and (c) bird's eye view.

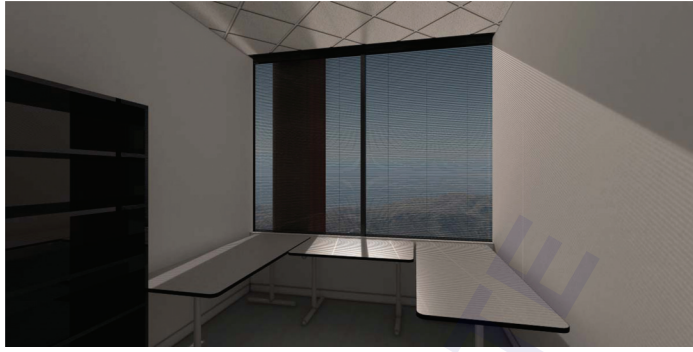


FIGURE 40.10a Rendering of a lit office room.

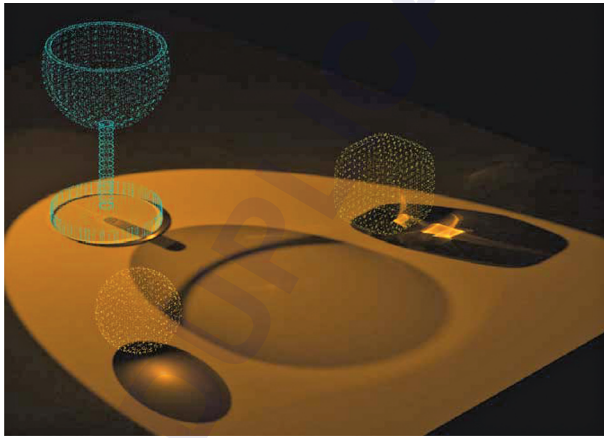


FIGURE 40.10b Rendering of a lit desk with three objects located on it (wine glass, ice cube, and crystal ball) to show both diffuse and specular effects.

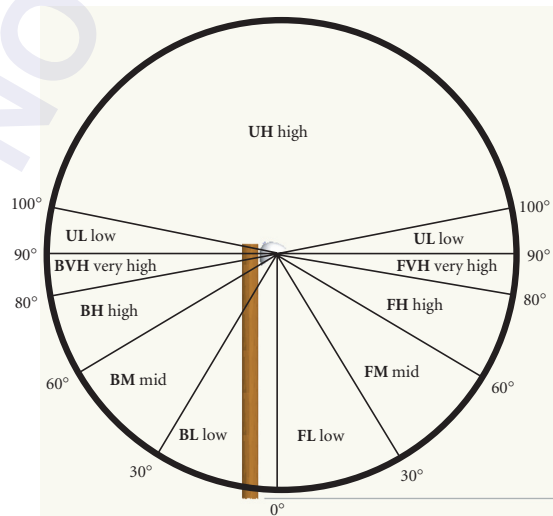


FIGURE 40.22 Layout of the light classification system subzones.



FIGURE 40.24 Depictions of luminaires: (a) Bankers lamp: multiple bounces inside the reflector create a wide angled uniform illumination; (b) Bouillotte lamp: vertical fluorescent tubes provide diffuse illumination; (c) indirect lighting with RLM fixture where the top surface reflects light into a wide angular range; (d) overhead direct-indirect lighting fixture using fluorescent tubular bulbs; and (e) parabolic louvered trough reflector for fluorescent tubes.



FIGURE 40.26 A conference room with artificial skylight made up of backlit ceiling image tiles.

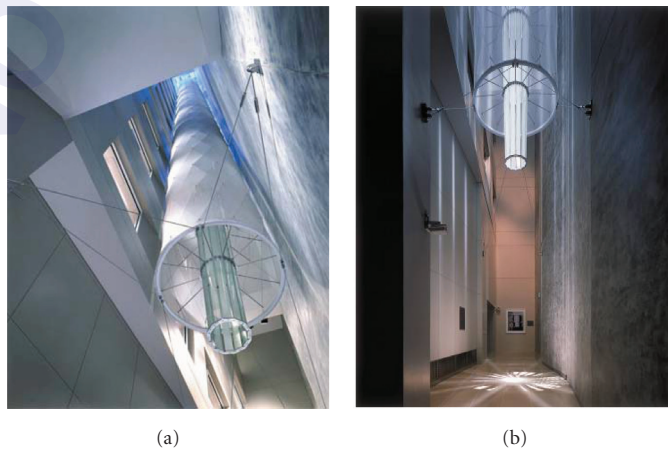


FIGURE 40.31 A Solar light pipe. (a) A 140-ft-tall light gathering and distributing device that presents daylight down into the core of a building that has no other access to daylight. (b) Light projected (10-in diameter) on the floor.

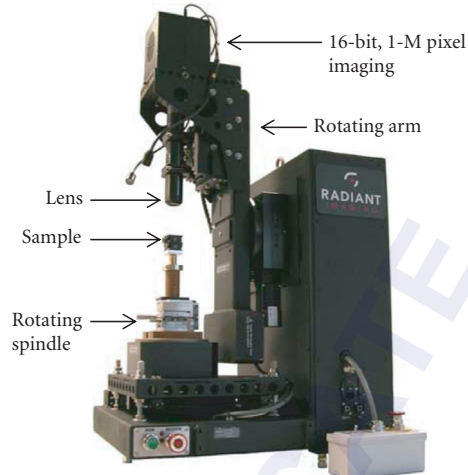


FIGURE 40.33 Photograph of a source measurement goniometer that is used to ascertain the luminance distribution of the source. The system wobble (electro-mechanical-software runout) is $15\ \mu\text{m}$ to allow for measuring small light sources like LED die.



FIGURE 40.34 (a) A faceted headlamp including high-beam (right), low-beam (middle), and turn signal (left) luminaire. Note the yellowish tinge of the turn signal, which is due to the coating placed on the bulb used therein. (b) A faceted taillight including the following functions: tail (upper left), stop (upper right), turn signal (lower right), reflex reflector (lower middle), and backup (lower left).

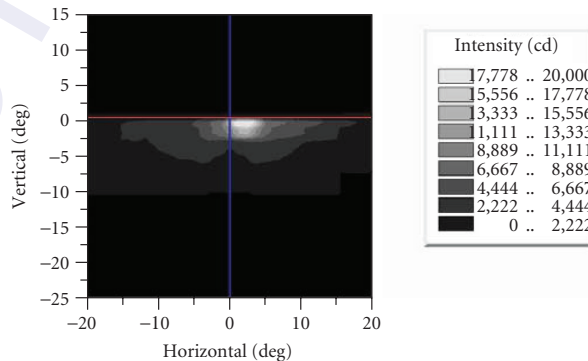


FIGURE 40.35 Luminous intensity (cd) distribution for the SAE low-beam requirements of Table 18.

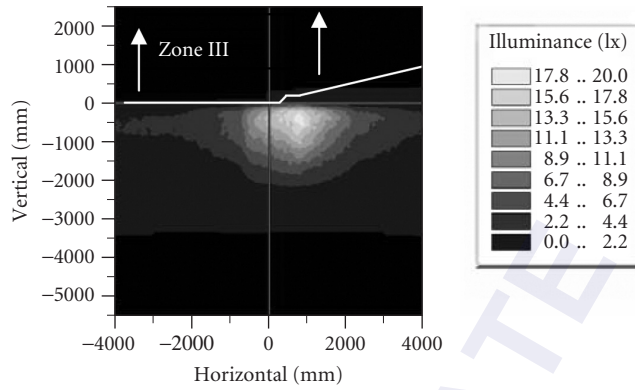


FIGURE 40.36 Illuminance (lx) distribution for the ECE passing/low-beam requirements of Table 19.

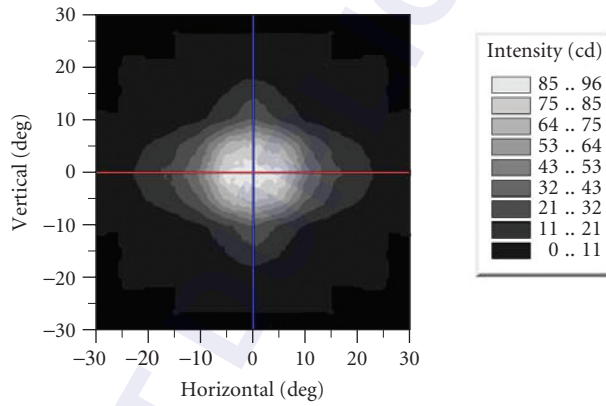


FIGURE 40.37 Luminous intensity (cd) distribution for the SAE stop lamp requirements of Table 20 (1 lit section).

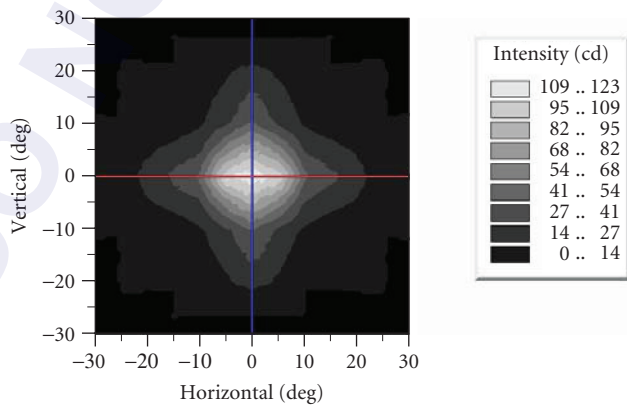


FIGURE 40.38 Luminous intensity (cd) distribution for the R7 stop lamp requirements of Table 21 (1 lamp illumination level).

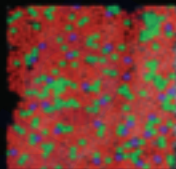
Third Edition

Sponsored by the Optical Society of America

HANDBOOK OF OPTICS

Volume III

Vision and Vision Optics



Editor-in-Chief:
Michael Bass

Associate Editors:
Casimer M. DeCusatis
Jay M. Enoch
Vasudevan Lakshminarayanan
Guifang Li
Carolyn MacDonald
Virendra N. Mahajan
Eric Van Stryland

OSA[®]

HANDBOOK OF OPTICS

DO NOT DUPLICATE

ABOUT THE EDITORS

Editor-in-Chief: Dr. Michael Bass is professor emeritus at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Associate Editors:

Dr. Casimer M. DeCusatis is a distinguished engineer and technical executive with IBM Corporation.

Dr. Jay M. Enoch is dean emeritus and professor at the School of Optometry at the University of California, Berkeley.

Dr. Vasudevan Lakshminarayanan is professor of Optometry, Physics, and Electrical Engineering at the University of Waterloo, Ontario, Canada.

Dr. Guifang Li is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Dr. Carolyn MacDonald is a professor at the University at Albany, and director of the Center for X-Ray Optics.

Dr. Virendra N. Mahajan is a distinguished scientist at The Aerospace Corporation.

Dr. Eric Van Stryland is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

HANDBOOK OF OPTICS

Volume III
Vision and Vision Optics

THIRD EDITION

Sponsored by the
OPTICAL SOCIETY OF AMERICA

Michael Bass Editor-in-Chief

*CREOL, The College of Optics and Photonics
University of Central Florida,
Orlando, Florida*

Jay M. Enoch Associate Editor

*School of Optometry, University of California at Berkeley
Berkeley, California
and
Department of Ophthalmology
University of California at San Francisco
San Francisco, California*

Vasudevan Lakshminarayanan Associate Editor

*School of Optometry and Departments of Physics and Electrical Engineering
University of Waterloo
Waterloo, Ontario, Canada*



New York Chicago San Francisco Lisbon London Madrid
Mexico City Milan New Delhi San Juan Seoul
Singapore Sydney Toronto

This page intentionally left blank.

DO NOT DUPLICATE

COVER ILLUSTRATIONS

A Photograph Taken of a Lady Viewing Her Face Using One of the World's Oldest Ground and Polished Mirrors.

The oldest known manufactured mirrors (ground and polished), made of obsidian (volcanic glass) have been found in ancient Anatolia in the ruins of the City of Çatal Hüyük = "mound at a road-fork." The locations where the mirrors were discovered were dated 6000 to 5900 B.C.E. by Mellaart and his coworkers. That city is located in the South Konya Plane of Modern Turkey. Thus, these mirrors are about 8000 years old (B.P.). The obsidian was transported over a distance of more than one hundred miles to the city for processing. These mirrors can be found at the Museum of Anatolian Civilizations in Ankara. One cannot fail to be impressed by the quality of this image seen by reflectance from this ancient mirror! These mirrors had been buried twice. There is an extended history of processing of obsidian at that site for scrapers, spear, and arrow points and other tools. This very early city contained an estimated 10,000 individuals at that time(!); it was a center for development of modern agriculture, Indo-European languages, various crafts, etc., and had established road connections and trade relations [Enoch, J., *Optom. Vision Sci.* **83**(10):775–781, 2006]. (*This figure is published with permission of Prof. Mellaart, the Director of the Museum of Anatolian Civilizations, the author, and the editor of the Journal.*)

Waveguide Modal Patterns in Vertebrate Eyes (Including Human).

This illustration demonstrates the variety of waveguide modal patterns observed in freshly removed retinas obtained from normal human, monkey, and rat retinas [Enoch, J., *J. Opt. Soc. Am.* **53**(1):71–85, 1963]. These modal patterns have been recorded in paracentral retinal receptors. Reverse path illumination was employed. These modes were photographed in near monochromatic light. This figure provides representative modal patterns observed and recorded near terminations of these photoreceptor waveguides. With variation of wavelength, at cutoff (please refer to the "V" parameter), it is possible to witness sharp modal pattern alterations. In this figure, the intent was to show the classes of modal patterns observed in these retinal receptors. (*This figure is reproduced with permission of JOSA and the author.*)

Photoreceptors in the Human Eye.

This figure shows the first map ever made of the spatial arrangement of the three cone classes in the human retina. The three colors (red, green, and blue) indicate cones that are sensitive to the long, middle, and short wavelength ranges of the visible spectrum and are classified as L, M, and S cones. The image was recorded from a living human eye using the adaptive optics ophthalmoscope, which was developed by David Williams' lab at the University of Rochester [Liang, J., Williams, D. R., and Miller, D. (1997). Supernormal vision and high-resolution retinal imaging through adaptive optics, *J. Opt. Soc. Am. A* **14**:2884–2892]. This image was first published in the journal *Nature* [Roorda, A., and Williams, D. R. (1999). The arrangement of the three cone classes in the living human eye, *Nature* **397**:520–522]. (*Courtesy of Austin Roorda and David Williams.*)

This page intentionally left blank.

DO NOT DUPLICATE

CONTENTS

Contributors	xiii
Brief Contents of All Volumes	xv
Editors' Preface	xxi
Preface to Volume III	xxiii
Glossary and Fundamental Constants	xxvii

Chapter 1. Optics of the Eye *Neil Charman* 1.1

1.1	Glossary / 1.1
1.2	Introduction / 1.3
1.3	Ocular Parameters and Ametropia / 1.4
1.4	Ocular Transmittance and Retinal Illuminance / 1.8
1.5	Factors Affecting In-Focus Retinal Image Quality / 1.12
1.6	Final Retinal Image Quality / 1.21
1.7	Depth-of-Focus and Accommodation / 1.28
1.8	Eye Models / 1.36
1.9	Two Eyes and Stereopsis / 1.38
1.10	Movements of the Eyes / 1.42
1.11	Conclusion / 1.45
1.12	References / 1.45

Chapter 2. Visual Performance *Wilson S. Geisler and Martin S. Banks* 2.1

2.1	Glossary / 2.1
2.2	Introduction / 2.2
2.3	Optics, Anatomy, Physiology of the Visual System / 2.2
2.4	Visual Performance / 2.14
2.5	Acknowledgments / 2.41
2.6	References / 2.42

Chapter 3. Psychophysical Methods *Denis G. Pelli and Bart Farell* 3.1

3.1	Introduction / 3.1
3.2	Definitions / 3.2
3.3	Visual Stimuli / 3.3
3.4	Adjustments / 3.4
3.5	Judgments / 3.6
	Magnitude Estimation / 3.8
3.6	Stimulus Sequencing / 3.9
3.7	Conclusion / 3.9
3.8	Tips from the Pros / 3.10
3.9	Acknowledgments / 3.10
3.10	References / 3.10

Chapter 4. Visual Acuity and Hyperacuity *Gerald Westheimer* 4.1

4.1	Glossary / 4.1
4.2	Introduction / 4.2
4.3	Stimulus Specifications / 4.2

- 4.4 Optics of the Eye's Resolving Capacity / 4.4
- 4.5 Retinal Limitations—Receptor Mosaic and Tiling of Neuronal Receptive Fields / 4.5
- 4.6 Determination of Visual Resolution Thresholds / 4.6
- 4.7 Kinds of Visual Acuity Tests / 4.7
- 4.8 Factors Affecting Visual Acuity / 4.9
- 4.9 Hyperacuity / 4.14
- 4.10 Resolution, Superresolution, and Information Theory / 4.15
- 4.11 Summary / 4.16
- 4.12 References / 4.16

Chapter 5. Optical Generation of the Visual Stimulus *Stephen A. Burns and Robert H. Webb* 5.1

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.1
- 5.3 The Size of the Visual Stimulus / 5.2
- 5.4 Free or Newtonian Viewing / 5.2
- 5.5 Maxwellian Viewing / 5.4
- 5.6 Building an Optical System / 5.8
- 5.7 Light Exposure and Ocular Safety / 5.18
- 5.8 Light Sources / 5.19
- 5.9 Coherent Radiation / 5.19
- 5.10 Detectors / 5.21
- 5.11 Putting It Together / 5.21
- 5.12 Conclusions / 5.24
- 5.13 Acknowledgments / 5.24
- 5.14 General References / 5.25
- 5.15 References / 5.26

Chapter 6. The Maxwellian View: with an Addendum on Apodization *Gerald Westheimer* 6.1

- 6.1 Glossary / 6.1
- 6.2 Introduction / 6.2
- 6.3 Postscript (2008) / 6.13

Chapter 7. Ocular Radiation Hazards *David H. Sliney* 7.1

- 7.1 Glossary / 7.1
- 7.2 Introduction / 7.2
- 7.3 Injury Mechanisms / 7.2
- 7.4 Types of Injury / 7.3
- 7.5 Retinal Irradiance Calculations / 7.7
- 7.6 Examples / 7.8
- 7.7 Exposure Limits / 7.9
- 7.8 Discussion / 7.11
- 7.9 References / 7.15

Chapter 8. Biological Waveguides *Vasudevan Lakshminarayanan and Jay M. Enoch* 8.1

- 8.1 Glossary / 8.1
- 8.2 Introduction / 8.2
- 8.3 Waveguiding in Retinal Photoreceptors and the Stiles-Crawford Effect / 8.3
- 8.4 Waveguides and Photoreceptors / 8.3
- 8.5 Photoreceptor Orientation and Alignment / 8.5
- 8.6 Introduction to the Models and Theoretical Implications / 8.8
- 8.7 Quantitative Observations of Single Receptors / 8.15
- 8.8 Waveguide Modal Patterns Found in Monkey/Human Retinal Receptors / 8.19
- 8.9 Light Guide Effect in Cochlear Hair Cells and Human Hair / 8.24

- 8.10 Fiber-Optic Plant Tissues / 8.26
- 8.11 Sponges / 8.28
- 8.12 Summary / 8.29
- 8.13 References / 8.29

Chapter 9. The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution *Jay M. Enoch and Vasudevan Lakshminarayanan* 9.1

- 9.1 Glossary / 9.1
- 9.2 Introduction / 9.2
- 9.3 The Problem and an Approach to Its Solution / 9.3
- 9.4 Sample Point-by-Point Estimates of SCE-1 and Integrated SCE-1 Data / 9.6
- 9.5 Discussion / 9.13
- 9.6 Teleological and Developmental Factors / 9.14
- 9.7 Conclusions / 9.14
- 9.8 References / 9.15

Chapter 10. Colorimetry *David H. Brainard and Andrew Stockman* 10.1

- 10.1 Glossary / 10.1
- 10.2 Introduction / 10.2
- 10.3 Fundamentals of Colorimetry / 10.3
- 10.4 Color Coordinate Systems / 10.11
- 10.5 Matrix Representations and Calculations / 10.24
- 10.6 Topics / 10.32
- 10.7 Appendix—Matrix Algebra / 10.45
- 10.8 References / 10.49

Chapter 11. Color Vision Mechanisms *Andrew Stockman and David H. Brainard* 11.1

- 11.1 Glossary / 11.1
- 11.2 Introduction / 11.3
- 11.3 Basics of Color-Discrimination Mechanisms / 11.9
- 11.4 Basics of Color-Appearance Mechanisms / 11.26
- 11.5 Details and Limits of the Basic Model / 11.31
- 11.6 Conclusions / 11.79
- 11.7 Acknowledgments / 11.85
- 11.8 References / 11.86

Chapter 12. Assessment of Refraction and Refractive Errors and Their Influence on Optical Design *B. Ralph Chou* 12.1

- 12.1 Glossary / 12.1
- 12.2 Introduction / 12.3
- 12.3 Refractive Errors / 12.3
- 12.4 Assessment of Refractive Error / 12.5
- 12.5 Correction of Refractive Error / 12.8
- 12.6 Binocular Factors / 12.15
- 12.7 Consequences for Optical Design / 12.17
- 12.8 References / 12.17

Chapter 13. Binocular Vision Factors That Influence Optical Design *Clifton Schor* 13.1

- 13.1 Glossary / 13.1
- 13.2 Combining the Images in the Two Eyes into One Perception of the Visual Field / 13.3
- 13.3 Distortion of Space by Monocular Magnification / 13.13

- 13.4 Distortion of Space Perception from Interocular Aniso-Magnification (Unequal Binocular Magnification) / 13.16
- 13.5 Distortions of Space from Convergence Responses to Prism / 13.19
- 13.6 Eye Movements / 13.19
- 13.7 Coordination and Alignment of the Two Eyes / 13.20
- 13.8 Effects of Lenses and Prism on Vergence and Phoria / 13.25
- 13.9 Prism-Induced Errors of Eye Alignment / 13.27
- 13.10 Head and Eye Responses to Direction (Gaze Control) / 13.29
- 13.11 Focus and Responses to Distance / 13.30
- 13.12 Video Head Sets, Head's Up Displays and Virtual Reality: Impact on Binocular Vision / 13.31
- 13.13 References / 13.35

Chapter 14. Optics and Vision of the Aging Eye *John S. Werner, Brooke E. Scheffrin, and Arthur Bradley* 14.1

- 14.1 Glossary / 14.1
- 14.2 Introduction / 14.2
- 14.3 The Graying of the Planet / 14.2
- 14.4 Senescence of the Eye's Optics / 14.4
- 14.5 Senescent Changes in Vision / 14.14
- 14.6 Age-Related Ocular Diseases Affecting Visual Function / 14.22
- 14.7 The Aging World from the Optical Point of View: Presbyopic Corrections / 14.27
- 14.8 Conclusions / 14.30
- 14.9 Acknowledgments / 14.30
- 14.10 References / 14.30

Chapter 15. Adaptive Optics in Retinal Microscopy and Vision *Donald T. Miller and Austin Roorda* 15.1

- 15.1 Glossary / 15.1
- 15.2 Introduction / 15.2
- 15.3 Properties of Ocular Aberrations / 15.4
- 15.4 Implementation of AO / 15.7
- 15.5 Application of AO to the Eye / 15.15
- 15.6 Acknowledgments / 15.24
- 15.7 References / 15.24

Chapter 16. Refractive Surgery, Correction of Vision, PRK and LASIK *L. Diaz-Santana and Harilaos Ginis* 16.1

- 16.1 Glossary / 16.1
- 16.2 Introduction / 16.2
- 16.3 Refractive Surgery Modalities / 16.9
- 16.4 Laser Ablation / 16.15
- 16.5 Acknowledgments / 16.19
- 16.6 References / 16.19

Chapter 17. Three-Dimensional Confocal Microscopy of the Living Human Cornea *Barry R. Masters* 17.1

- 17.1 Glossary / 17.1
- 17.2 Introduction / 17.3
- 17.3 Theory of Confocal Microscopy / 17.3
- 17.4 The Development of Confocal Instruments / 17.3
- 17.5 The Scanning Slit and Laser Scanning Clinical Confocal Microscopes / 17.6
- 17.6 Clinical Applications of Confocal Microscopy / 17.8
- 17.7 Perspectives / 17.9
- 17.8 Summary / 17.10
- 17.9 Acknowledgments / 17.10
- 17.10 References / 17.10

Chapter 18. Diagnostic Use of Optical Coherence Tomography in the Eye *Johannes F. de Boer* 18.1

- 18.1 Glossary / 18.1
- 18.2 Introduction / 18.2
- 18.3 Principle of OCT: Time Domain OCT / 18.3
- 18.4 Principle of OCT: Spectral Domain OCT / 18.5
- 18.5 Principle of OCT: Optical Frequency Domain Imaging / 18.7
- 18.6 SD-OCT Versus OFDI / 18.9
- 18.7 Sensitivity Advantage of SD-OCT Over TD-OCT / 18.9
- 18.8 Noise Analysis of SD-OCT Using Charge Coupled Devices (CCDs) / 18.9
- 18.9 Signal to Noise Ratio and Autocorrelation Noise / 18.11
- 18.10 Shot-Noise-Limited Detection / 18.12
- 18.11 Depth Dependent Sensitivity / 18.13
- 18.12 Motion Artifacts and Fringe Washout / 18.15
- 18.13 OFDI at 1050 NM / 18.15
- 18.14 Functional Extensions: Doppler OCT and Polarization Sensitive OCT / 18.18
- 18.15 Doppler OCT and Phase Stability / 18.18
- 18.16 Polarization Sensitive OCT (PS-OCT) / 18.20
- 18.17 PS-OCT in Ophthalmology / 18.24
- 18.18 Retinal Imaging with SD-OCT / 18.27
- 18.19 Conclusion / 18.29
- 18.20 Acknowledgment / 18.30
- 18.21 References / 18.30

Chapter 19. Gradient Index Optics in the Eye *Barbara K. Pierscionek* 19.1

- 19.1 Glossary / 19.1
- 19.2 Introduction / 19.2
- 19.3 The Nature of an Index Gradient / 19.2
- 19.4 Spherical Gradients / 19.2
- 19.5 Radial Gradients / 19.3
- 19.6 Axial Gradients / 19.5
- 19.7 The Eye Lens / 19.5
- 19.8 Fish / 19.6
- 19.9 Octopus / 19.7
- 19.10 Rat / 19.7
- 19.11 Guinea Pig / 19.8
- 19.12 Rabbit / 19.8
- 19.13 Cat / 19.9
- 19.14 Bovine / 19.9
- 19.15 Pig / 19.11
- 19.16 Human/primate / 19.12
- 19.17 Functional Considerations / 19.14
- 19.18 Summary / 19.15
- 19.19 References / 19.15

Chapter 20. Optics of Contact Lenses *Edward S. Bennett* 20.1

- 20.1 Glossary / 20.1
- 20.2 Introduction / 20.2
- 20.3 Contact Lens Material, Composition, and Design Parameters / 20.3
- 20.4 Contact Lens Power / 20.6
- 20.5 Other Design Considerations / 20.20
- 20.6 Convergence and Accommodation Effects / 20.25
- 20.7 Prismatic Effects / 20.30
- 20.8 Magnification / 20.31
- 20.9 Summary / 20.34
- 20.10 Acknowledgments / 20.34
- 20.11 References / 20.34

Chapter 21.	Intraocular Lenses <i>Jim Schwiegerling</i>	21.1
21.1	Glossary / 21.1	
21.2	Introduction / 21.2	
21.3	Cataract Surgery / 21.4	
21.4	Intraocular Lens Design / 21.5	
21.5	Intraocular Lens Side Effects / 21.20	
21.6	Summary / 21.22	
21.7	References / 21.22	
Chapter 22.	Displays for Vision Research <i>William Cowan</i>	22.1
22.1	Glossary / 22.1	
22.2	Introduction / 22.2	
22.3	Operational Characteristics of Color Monitors / 22.3	
22.4	Colorimetric Calibration of Video Monitors / 22.20	
22.5	An Introduction to Liquid Crystal Displays / 22.34	
22.6	Acknowledgments / 22.40	
22.7	References / 22.40	
Chapter 23.	Vision Problems at Computers <i>Jeffrey Anshel and James E. Sheedy</i>	23.1
23.1	Glossary / 23.1	
23.2	Introduction / 23.4	
23.3	Work Environment / 23.4	
23.4	Vision and Eye Conditions / 23.9	
23.5	References / 23.12	
Chapter 24.	Human Vision and Electronic Imaging <i>Bernice E. Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allebach</i>	24.1
24.1	Introduction / 24.1	
24.2	Early Vision Approaches: The Perception of Imaging Artifacts / 24.2	
24.3	Higher-Level Approaches: The Analysis of Image Features / 24.6	
24.4	Very High-Level Approaches: The Representation of Aesthetic and Emotional Characteristics / 24.9	
24.5	Conclusions / 24.10	
24.6	Additional Information on Human Vision and Electronic Imaging / 24.11	
24.7	References / 24.11	
Chapter 25.	Visual Factors Associated with Head-mounted Displays <i>Brian H. Tsou and Martin Shenker</i>	25.1
25.1	Glossary / 25.1	
25.2	Introduction / 25.1	
25.3	Common Design Considerations among All HMDs / 25.2	
25.4	Characterizing HMD / 25.7	
25.5	Summary / 25.10	
25.6	Appendix / 25.10	
25.7	Acknowledgments / 25.12	
25.8	References / 25.12	

CONTRIBUTORS

- Jan P. Allebach** *Electronic Imaging Systems Laboratory, School of Electrical and Computer Engineering, Purdue University, West Lafayette, Indiana (CHAP. 24)*
- Jeffrey Anshel** *Corporate Vision Consulting, Encinitas, California (CHAP. 23)*
- Martin S. Banks** *School of Optometry, University of California, Berkeley, California (CHAP. 2)*
- Edward S. Bennett** *College of Optometry, University of Missouri, St. Louis, Missouri (CHAP. 8)*
- Arthur Bradley** *School of Optometry, Indiana University, Bloomington, Indiana (CHAP. 14)*
- David H. Brainard** *Department of Psychology, University of Pennsylvania, Philadelphia, Pennsylvania (CHAPS. 10, 11)*
- Stephen A. Burns** *School of Optometry, Indiana University, Bloomington, Indiana (CHAP. 5)*
- Neil Charman** *Department of Optometry and Vision Sciences, University of Manchester, Manchester, United Kingdom (CHAP. 1)*
- B. Ralph Chou** *School of Optometry, University of Waterloo, Waterloo, Ontario, Canada (CHAP. 12)*
- William Cowan** *Department of Computer Science, University of Waterloo, Waterloo, Ontario, Canada (CHAP. 22)*
- Johannes F. de Boer** *Department of Physics, VU University, Amsterdam, and Rotterdam Ophthalmic Institute, Rotterdam, The Netherlands (CHAP. 18)*
- Jay M. Enoch** *School of Optometry, University of California at Berkeley, Berkeley, California (CHAPS. 8, 9)*
- Bart Farell** *Institute for Sensory Research, Syracuse University, Syracuse, New York (CHAP. 3)*
- Wilson S. Geisler** *Department of Psychology, University of Texas, Austin, Texas (CHAP. 2)*
- Harilaos Ginis** *Institute of Vision and Optics, University of Crete, Greece (CHAP. 16)*
- Vasudevan Lakshminarayanan** *School of Optometry and Departments of Physics and Electrical Engineering, University of Waterloo, Waterloo, Ontario, Canada (CHAPS. 8, 9)*
- Barry R. Masters** *Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts (CHAP. 17)*
- Donald T. Miller** *School of Optometry, Indiana University, Bloomington, Indiana (CHAP. 15)*
- Thrasyvoulos N. Pappas** *Department of Electrical and Computer Engineering, Northwestern University Evanston, Illinois (CHAP. 24)*
- Denis G. Pelli** *Psychology Department and Center for Neural Science, New York University, New York (CHAP. 3)*
- Barbara K. Pierscionek** *Department of Biomedical Sciences, University of Ulster, Coleraine, United Kingdom (CHAP. 19)*
- Bernice E. Rogowitz** *IBM T. J. Watson Research Center, Hawthorne, New York (CHAP. 24)*
- Austin Roorda** *School of Optometry, University of California, Berkeley, California (CHAP. 15)*
- L. Diaz-Santana** *Department of Optometry and Visual Science, City University, London, United Kingdom (CHAP. 16)*
- Brooke E. Scheffrin** *Department of Psychology, University of Colorado, Boulder, Colorado (CHAP. 14)*
- Clifton Schor** *School of Optometry, University of California, Berkeley, California (CHAP. 13)*
- Jim Schwiegerling** *Department of Ophthalmology, University of Arizona, Tucson, Arizona (CHAP. 21)*

- James E. Sheedy** *College of Optometry, Pacific University, Forest Grove, Oregon (CHAP. 23)*
- Martin Shenker** *Martin Shenker Optical Design, Inc., White Plains, New York (CHAP. 25)*
- David H. Sliney** *Consulting Medical Physicist, Fallston, Maryland, and Retired, U.S. Army Center for Health Promotion and Preventive Medicine, Laser/Optical Radiation Program, Aberdeen Proving Ground, Maryland (CHAP. 7)*
- Andrew Stockman** *Department of Visual Neuroscience, UCL Institute of Ophthalmology, London, United Kingdom (CHAPS. 10, 11)*
- Brian H. Tsou** *Air Force Research Laboratory, Wright Patterson AFB, Ohio (CHAP. 25)*
- Robert H. Webb** *The Schepens Eye Research Institute, Boston, Massachusetts (CHAP. 5)*
- John S. Werner** *Department of Ophthalmology & Vision Science, University of California, Davis, Sacramento, California (CHAP. 14)*
- Gerald Westheimer** *Division of Neurobiology, University of California, Berkeley, California (CHAPS. 4, 6)*

BRIEF CONTENTS OF ALL VOLUMES

VOLUME I. GEOMETRICAL AND PHYSICAL OPTICS, POLARIZED LIGHT, COMPONENTS AND INSTRUMENTS

PART 1. GEOMETRICAL OPTICS

Chapter 1. General Principles of Geometrical Optics *Douglas S. Goodman*

PART 2. PHYSICAL OPTICS

Chapter 2. Interference *John E. Greivenkamp*

Chapter 3. Diffraction *Arvind S. Marathay and John F. McCalmont*

Chapter 4. Transfer Function Techniques *Glenn D. Boreman*

Chapter 5. Coherence Theory *William H. Carter*

Chapter 6. Coherence Theory: Tools and Applications *Gisele Bennett, William T. Rhodes, and J. Christopher James*

Chapter 7. Scattering by Particles *Craig F. Bohren*

Chapter 8. Surface Scattering *Eugene L. Church and Peter Z. Takacs*

Chapter 9. Volume Scattering in Random Media *Aristide Dogariu and Jeremy Ellis*

Chapter 10. Optical Spectroscopy and Spectroscopic Lineshapes *Brian Henderson*

Chapter 11. Analog Optical Signal and Image Processing *Joseph W. Goodman*

PART 3. POLARIZED LIGHT

Chapter 12. Polarization *Jean M. Bennett*

Chapter 13. Polarizers *Jean M. Bennett*

Chapter 14. Mueller Matrices *Russell A. Chipman*

Chapter 15. Polarimetry *Russell A. Chipman*

Chapter 16. Ellipsometry *Rasheed M. A. Azzam*

PART 4. COMPONENTS

Chapter 17. Lenses *R. Barry Johnson*

Chapter 18. Afocal Systems *William B. Wetherell*

Chapter 19. Nondispersive Prisms *William L. Wolfe*

Chapter 20. Dispersive Prisms and Gratings *George J. Zisis*

Chapter 21. Integrated Optics *Thomas L. Koch, Frederick J. Leonberger, and Paul G. Suchoski*

Chapter 22. Miniature and Micro-Optics *Tom D. Milster and Tomasz S. Tkaczyk*

Chapter 23. Binary Optics *Michael W. Farn and Wilfrid B. Veldkamp*

Chapter 24. Gradient Index Optics *Duncan T. Moore*

PART 5. INSTRUMENTS

Chapter 25. Cameras *Norman Goldberg*

Chapter 26. Solid-State Cameras *Gerald C. Holst*

Chapter 27. Camera Lenses *Ellis Betensky, Melvin H. Kreitzer, and Jacob Moskovich*

Chapter 28. Microscopes *Rudolf Oldenbourg and Michael Shribak*

Chapter 29. Reflective and Catadioptric Objectives *Lloyd Jones*

- Chapter 30. Scanners *Leo Beiser and R. Barry Johnson*
- Chapter 31. Optical Spectrometers *Brian Henderson*
- Chapter 32. Interferometers *Parameswaran Hariharan*
- Chapter 33. Holography and Holographic Instruments *Lloyd Huff*
- Chapter 34. Xerographic Systems *Howard Stark*
- Chapter 35. Principles of Optical Disk Data Storage *Masud Mansuripur*

VOLUME II. DESIGN, FABRICATION, AND TESTING; SOURCES AND DETECTORS; RADIOMETRY AND PHOTOMETRY

PART 1. DESIGN

- Chapter 1. Techniques of First-Order Layout *Warren J. Smith*
- Chapter 2. Aberration Curves in Lens Design *Donald C. O'Shea and Michael E. Harrigan*
- Chapter 3. Optical Design Software *Douglas C. Sinclair*
- Chapter 4. Optical Specifications *Robert R. Shannon*
- Chapter 5. Tolerancing Techniques *Robert R. Shannon*
- Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.*
- Chapter 7. Control of Stray Light *Robert P. Breault*
- Chapter 8. Thermal Compensation Techniques *Philip J. Rogers and Michael Roberts*

PART 2. FABRICATION

- Chapter 9. Optical Fabrication *Michael P. Mandina*
- Chapter 10. Fabrication of Optics by Diamond Turning *Richard L. Rhorer and Chris J. Evans*

PART 3. TESTING

- Chapter 11. Orthonormal Polynomials in Wavefront Analysis *Virendra N. Mahajan*
- Chapter 12. Optical Metrology *Zacarias Malacara and Daniel Malacara-Hernández*
- Chapter 13. Optical Testing *Daniel Malacara-Hernández*
- Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant*

PART 4. SOURCES

- Chapter 15. Artificial Sources *Anthony LaRocca*
- Chapter 16. Lasers *William T. Silfvast*
- Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman*
- Chapter 18. High-Brightness Visible LEDs *Winston V. Schoenfeld*
- Chapter 19. Semiconductor Lasers *Pamela L. Derry, Luis Figueroa, and Chi-shain Hong*
- Chapter 20. Ultrashort Optical Sources and Applications *Jean-Claude Diels and Ladan Arissian*
- Chapter 21. Attosecond Optics *Zenghu Chang*
- Chapter 22. Laser Stabilization *John L. Hall, Matthew S. Taubman, and Jun Ye*
- Chapter 23. Quantum Theory of the Laser *János A. Bergou, Berthold-Georg Englert, Melvin Lax, Marian O. Scully, Herbert Walther, and M. Suhail Zubairy*

PART 5. DETECTORS

- Chapter 24. Photodetectors *Paul R. Norton*
- Chapter 25. Photodetection *Abhay M. Joshi and Gregory H. Olsen*
- Chapter 26. High-Speed Photodetectors *John E. Bowers and Yih G. Wey*
- Chapter 27. Signal Detection and Analysis *John R. Willison*
- Chapter 28. Thermal Detectors *William L. Wolfe and Paul W. Kruse*

PART 6. IMAGING DETECTORS

- Chapter 29. Photographic Films *Joseph H. Altman*
- Chapter 30. Photographic Materials *John D. Baloga*

- Chapter 31. Image Tube Intensified Electronic Imaging *C. Bruce Johnson and Larry D. Owen*
 Chapter 32. Visible Array Detectors *Timothy J. Tredwell*
 Chapter 33. Infrared Detector Arrays *Lester J. Kozlowski and Walter F. Kosonocky*

PART 7. RADIOMETRY AND PHOTOMETRY

- Chapter 34. Radiometry and Photometry *Edward F. Zalewski*
 Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection *James M. Palmer*
 Chapter 36. Radiometry and Photometry: Units and Conversions *James M. Palmer*
 Chapter 37. Radiometry and Photometry for Vision Optics *Yoshi Ohno*
 Chapter 38. Spectroradiometry *Carolyn J. Sher DeCusatis*
 Chapter 39. Nonimaging Optics: Concentration and Illumination *William Cassarly*
 Chapter 40. Lighting and Applications *Anurag Gupta and R. John Koschel*

VOLUME III. VISION AND VISION OPTICS

- Chapter 1. Optics of the Eye *Neil Charman*
 Chapter 2. Visual Performance *Wilson S. Geisler and Martin S. Banks*
 Chapter 3. Psychophysical Methods *Denis G. Pelli and Bart Farell*
 Chapter 4. Visual Acuity and Hyperacuity *Gerald Westheimer*
 Chapter 5. Optical Generation of the Visual Stimulus *Stephen A. Burns and Robert H. Webb*
 Chapter 6. The Maxwellian View with an Addendum on Apodization *Gerald Westheimer*
 Chapter 7. Ocular Radiation Hazards *David H. Sliney*
 Chapter 8. Biological Waveguides *Vasudevan Lakshminarayanan and Jay M. Enoch*
 Chapter 9. The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution *Jay M. Enoch and Vasudevan Lakshminarayanan*
 Chapter 10. Colorimetry *David H. Brainard and Andrew Stockman*
 Chapter 11. Color Vision Mechanisms *Andrew Stockman and David H. Brainard*
 Chapter 12. Assessment of Refraction and Refractive Errors and Their Influence on Optical Design
B. Ralph Chou
 Chapter 13. Binocular Vision Factors That Influence Optical Design *Clifton Schor*
 Chapter 14. Optics and Vision of the Aging Eye *John S. Werner, Brooke E. Scheffrin, and Arthur Bradley*
 Chapter 15. Adaptive Optics in Retinal Microscopy and Vision *Donald T. Miller and Austin Roorda*
 Chapter 16. Refractive Surgery, Correction of Vision, PRK, and LASIK *L. Diaz-Santana and Harilaos Ginis*
 Chapter 17. Three-Dimensional Confocal Microscopy of the Living Human Cornea *Barry R. Masters*
 Chapter 18. Diagnostic Use of Optical Coherence Tomography in the Eye *Johannes F. de Boer*
 Chapter 19. Gradient Index Optics in the Eye *Barbara K. Pierscionek*
 Chapter 20. Optics of Contact Lenses *Edward S. Bennett*
 Chapter 21. Intraocular Lenses *Jim Schwiegerling*
 Chapter 22. Displays for Vision Research *William Cowan*
 Chapter 23. Vision Problems at Computers *Jeffrey Anshel and James E. Sheedy*
 Chapter 24. Human Vision and Electronic Imaging *Bernice E. Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allebach*
 Chapter 25. Visual Factors Associated with Head-Mounted Displays *Brian H. Tsou and Martin Shenker*

VOLUME IV. OPTICAL PROPERTIES OF MATERIALS, NONLINEAR OPTICS, QUANTUM OPTICS

PART 1. PROPERTIES

- Chapter 1. Optical Properties of Water *Curtis D. Mobley*
 Chapter 2. Properties of Crystals and Glasses *William J. Tropf, Michael E. Thomas, and Eric W. Rogala*
 Chapter 3. Polymeric Optics *John D. Lytle*

- Chapter 4. Properties of Metals *Roger A. Paquin*
Chapter 5. Optical Properties of Semiconductors *David G. Seiler, Stefan Zollner, Alain C. Diebold, and Paul M. Amirtharaj*
Chapter 6. Characterization and Use of Black Surfaces for Optical Systems *Stephen M. Pompea and Robert P. Breault*
Chapter 7. Optical Properties of Films and Coatings *Jerzy A. Dobrowolski*
Chapter 8. Fundamental Optical Properties of Solids *Alan Miller*
Chapter 9. Photonic Bandgap Materials *Pierre R. Villeneuve*

PART 2. NONLINEAR OPTICS

- Chapter 10. Nonlinear Optics *Chung L. Tang*
Chapter 11. Coherent Optical Transients *Paul R. Berman and D. G. Steel*
Chapter 12. Photorefractive Materials and Devices *Mark Cronin-Golomb and Marvin Klein*
Chapter 13. Optical Limiting *David J. Hagan*
Chapter 14. Electromagnetically Induced Transparency *Jonathan P. Marangos and Thomas Halfmann*
Chapter 15. Stimulated Raman and Brillouin Scattering *John Reintjes and M. Bashkansky*
Chapter 16. Third-Order Optical Nonlinearities *Mansoor Sheik-Bahae and Michael P. Hasselbeck*
Chapter 17. Continuous-Wave Optical Parametric Oscillators *M. Ebrahim-Zadeh*
Chapter 18. Nonlinear Optical Processes for Ultrashort Pulse Generation *Uwe Siegner and Ursula Keller*
Chapter 19. Laser-Induced Damage to Optical Materials *Marion J. Soileau*

PART 3. QUANTUM AND MOLECULAR OPTICS

- Chapter 20. Laser Cooling and Trapping of Atoms *Harold J. Metcalf and Peter van der Straten*
Chapter 21. Strong Field Physics *Todd Ditmire*
Chapter 22. Slow Light Propagation in Atomic and Photonic Media *Jacob B. Khurgin*
Chapter 23. Quantum Entanglement in Optical Interferometry *Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver, and Jonathan P. Dowling*

VOLUME V. ATMOSPHERIC OPTICS, MODULATORS, FIBER OPTICS, X-RAY AND NEUTRON OPTICS

PART 1. MEASUREMENTS

- Chapter 1. Scatterometers *John C. Stover*
Chapter 2. Spectroscopic Measurements *Brian Henderson*

PART 2. ATMOSPHERIC OPTICS

- Chapter 3. Atmospheric Optics *Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman*
Chapter 4. Imaging through Atmospheric Turbulence *Virendra N. Mahajan and Guang-ming Dai*
Chapter 5. Adaptive Optics *Robert Q. Fugate*

PART 3. MODULATORS

- Chapter 6. Acousto-Optic Devices *I-Cheng Chang*
Chapter 7. Electro-Optic Modulators *Georgianne M. Purvinis and Theresa A. Maldonado*
Chapter 8. Liquid Crystals *Sebastian Gauza and Shin-Tson Wu*

PART 4. FIBER OPTICS

- Chapter 9. Optical Fiber Communication Technology and System Overview *Ira Jacobs*
Chapter 10. Nonlinear Effects in Optical Fibers *John A. Buck*
Chapter 11. Photonic Crystal Fibers *Philip St. J. Russell and G. J. Pearce*
Chapter 12. Infrared Fibers *James A. Harrington*

- Chapter 13. Sources, Modulators, and Detectors for Fiber Optic Communication Systems *Elsa Garmire*
Chapter 14. Optical Fiber Amplifiers *John A. Buck*
Chapter 15. Fiber Optic Communication Links (Telecom, Datacom, and Analog) *Casimer DeCusatis and Guifang Li*
Chapter 16. Fiber-Based Couplers *Daniel Nolan*
Chapter 17. Fiber Bragg Gratings *Kenneth O. Hill*
Chapter 18. Micro-Optics-Based Components for Networking *Joseph C. Palais*
Chapter 19. Semiconductor Optical Amplifiers *Jay M. Wiesenfeld and Leo H. Spiekman*
Chapter 20. Optical Time-Division Multiplexed Communication Networks *Peter J. Delfyett*
Chapter 21. WDM Fiber-Optic Communication Networks *Alan E. Willner, Changyuan Yu, Zhongqi Pan, and Yong Xie*
Chapter 22. Solitons in Optical Fiber Communication Systems *Pavel V. Mamyshev*
Chapter 23. Fiber-Optic Communication Standards *Casimer DeCusatis*
Chapter 24. Optical Fiber Sensors *Richard O. Claus, Ignacio Matias, and Francisco Arregui*
Chapter 25. High-Power Fiber Lasers and Amplifiers *Timothy S. McComb, Martin C. Richardson, and Michael Bass*

PART 5. X-RAY AND NEUTRON OPTICS

Subpart 5.1. Introduction and Applications

- Chapter 26. An Introduction to X-Ray and Neutron Optics *Carolyn A. MacDonald*
Chapter 27. Coherent X-Ray Optics and Microscopy *Qun Shen*
Chapter 28. Requirements for X-Ray diffraction *Scott T. Misture*
Chapter 29. Requirements for X-Ray Fluorescence *George J. Havrilla*
Chapter 30. Requirements for X-Ray Spectroscopy *Dirk Lützenkirchen-Hecht and Ronald Frahm*
Chapter 31. Requirements for Medical Imaging and X-Ray Inspection *Douglas Pfeiffer*
Chapter 32. Requirements for Nuclear Medicine *Lars R. Furenlid*
Chapter 33. Requirements for X-Ray Astronomy *Scott O. Rohrbach*
Chapter 34. Extreme Ultraviolet Lithography *Franco Cerrina and Fan Jiang*
Chapter 35. Ray Tracing of X-Ray Optical Systems *Franco Cerrina and M. Sanchez del Rio*
Chapter 36. X-Ray Properties of Materials *Eric M. Gullikson*

Subpart 5.2. Refractive and Interference Optics

- Chapter 37. Refractive X-Ray Lenses *Bruno Lengeler and Christian G. Schroer*
Chapter 38. Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region *Malcolm R. Howells*
Chapter 39. Crystal Monochromators and Bent Crystals *Peter Siddons*
Chapter 40. Zone Plates *Alan Michette*
Chapter 41. Multilayers *Eberhard Spiller*
Chapter 42. Nanofocusing of Hard X-Rays with Multilayer Laue Lenses *Albert T. Macrander, Hanfei Yan, Hyon Chol Kang, Jörg Maser, Chian Liu, Ray Conley, and G. Brian Stephenson*
Chapter 43. Polarizing Crystal Optics *Qun Shen*

Subpart 5.3. Reflective Optics

- Chapter 44. Reflective Optics *James Harvey*
Chapter 45. Aberrations for Grazing Incidence Optics *Timo T. Saha*
Chapter 46. X-Ray Mirror Metrology *Peter Z. Takacs*
Chapter 47. Astronomical X-Ray Optics *Marshall K. Joy and Brian D. Ramsey*
Chapter 48. Multifoil X-Ray Optics *Ladislav Pina*
Chapter 49. Pore Optics *Marco Beijersbergen*
Chapter 50. Adaptive X-Ray Optics *Ali Khounsary*
Chapter 51. The Schwarzschild Objective *Franco Cerrina*
Chapter 52. Single Capillaries *Donald H. Bilderback and Sterling W. Cornaby*
Chapter 53. Polycapillary X-Ray Optics *Carolyn MacDonald and Walter M. Gibson*

Subpart 5.4. X-Ray Sources

- Chapter 54. X-Ray Tube Sources *Susanne M. Lee and Carolyn MacDonald*
Chapter 55. Synchrotron Sources *Steven L. Hulbert and Gwyn P. Williams*
Chapter 56. Laser Generated Plasmas *Alan Michette*
Chapter 57. Pinch Plasma Sources *Victor Kantsyrev*
Chapter 58. X-Ray Lasers *Greg Tallents*
Chapter 59. Inverse Compton X-Ray Sources *Frank Carroll*

Subpart 5.5. X-Ray Detectors

- Chapter 60. Introduction to X-Ray Detectors *Walter M. Gibson and Peter Siddons*
Chapter 61. Advances in Imaging Detectors *Aaron Couture*
Chapter 62. X-Ray Spectral Detection and Imaging *Eric Lifshin*

Subpart 5.6. Neutron Optics and Applications

- Chapter 63. Neutron Optics *David Mildner*
Chapter 64. Grazing-Incidence Neutron Optics *Mikhail Gubarev and Brian Ramsey*

DO NOT DUPLICATE

EDITORS' PREFACE

The third edition of the *Handbook of Optics* is designed to pull together the dramatic developments in both the basic and applied aspects of the field while retaining the archival, reference book value of a handbook. This means that it is much more extensive than either the first edition, published in 1978, or the second edition, with Volumes I and II appearing in 1995 and Volumes III and IV in 2001. To cover the greatly expanded field of optics, the *Handbook* now appears in five volumes. Over 100 authors or author teams have contributed to this work.

Volume I is devoted to the fundamentals, components, and instruments that make optics possible. Volume II contains chapters on design, fabrication, testing, sources of light, detection, and a new section devoted to radiometry and photometry. Volume III concerns vision optics only and is printed entirely in color. In Volume IV there are chapters on the optical properties of materials, non-linear, quantum and molecular optics. Volume V has extensive sections on fiber optics and x ray and neutron optics, along with shorter sections on measurements, modulators, and atmospheric optical properties and turbulence. Several pages of color inserts are provided where appropriate to aid the reader. A purchaser of the print version of any volume of the *Handbook* will be able to download a digital version containing all of the material in that volume in PDF format to one computer (see download instructions on bound-in card). The combined index for all five volumes can be downloaded from www.HandbookofOpticsOnline.com.

It is possible by careful selection of what and how to present that the third edition of the *Handbook* could serve as a text for a comprehensive course in optics. In addition, students who take such a course would have the *Handbook* as a career-long reference.

Topics were selected by the editors so that the *Handbook* could be a desktop (bookshelf) general reference for the parts of optics that had matured enough to warrant archival presentation. New chapters were included on topics that had reached this stage since the second edition, and existing chapters from the second edition were updated where necessary to provide this compendium. In selecting subjects to include, we also had to select which subjects to leave out. The criteria we applied were: (1) was it a specific application of optics rather than a core science or technology and (2) was it a subject in which the role of optics was peripheral to the central issue addressed. Thus, such topics as medical optics, laser surgery, and laser materials processing were not included. While applications of optics are mentioned in the chapters there is no space in the *Handbook* to include separate chapters devoted to all of the myriad uses of optics in today's world. If we had, the third edition would be much longer than it is and much of it would soon be outdated. We designed the third edition of the *Handbook of Optics* so that it concentrates on the principles of optics that make applications possible.

Authors were asked to try to achieve the dual purpose of preparing a chapter that was a worthwhile reference for someone working in the field and that could be used as a starting point to become acquainted with that aspect of optics. They did that and we thank them for the outstanding results seen throughout the *Handbook*. We also thank Mr. Taisuke Soda of McGraw-Hill for his help in putting this complex project together and Mr. Alan Tourtlotte and Ms. Susannah Lehman of the Optical Society of America for logistical help that made this effort possible.

We dedicate the third edition of the *Handbook of Optics* to all of the OSA volunteers who, since OSA's founding in 1916, give their time and energy to promoting the generation, application, archiving, and worldwide dissemination of knowledge in optics and photonics.

Michael Bass, Editor-in-Chief

Associate Editors:

Casimer M. DeCusatis

Jay M. Enoch

Vasudevan Lakshminarayanan

Guifang Li

Carolyn MacDonald

Virendra N. Mahajan

Eric Van Stryland

This page intentionally left blank.

DO NOT DUPLICATE

PREFACE TO VOLUME III

Volume III of the *Handbook of Optics*, Third Edition, addresses topics relating to vision and the eye which are applicable to, or relate to the study of optics. For reasons we do not understand fully, in recent years, there seems to have been a tendency for the optics and the vision science communities (the latter group was known in earlier times as “physiological optics”) to drift somewhat apart. Physiological optics had become a meaningful component within optics during the latter part of the nineteenth century. As but one example, we urge interested readers to read H. von Helmholtz’s masterful three-volume *Handbook of Physiological Optics* (third edition) which was translated by J. P. C. Southall, of Columbia University, into English by the Optical Society of America in the 1920s.¹ It should also be noted that Allvar Gullstrand received the Nobel Prize in Physiology/Medicine in 1911 for his work on model eyes, which was a direct application of thick lens theory. Gullstrand was not only a professor of ophthalmology at the University of Uppsala, but was also a professor of physiological and physical optics at that institution. He also added five new chapters to the first volume of Helmholtz’s treatise published in 1909. Not only is this a remarkable scientific work, but much of it remains applicable today! The simple fact is that the two groups, optical science and vision science need each other, or, alternatively, are effectively “joined at the hip.” Thus, here, we seek to provide a broad view of vision, vision processes, and discussions of areas where vision science interacts with the ever broadening field of optics.

Obviously, no treatment such as this one can be complete, but we have tried here to present applicable topics in an orderly manner. In the current edition, we have taken a wide-ranging view of vision and its relationship with optics. In particular, in recent years, we have seen a rapid increase of interest in new technologies and applications in the areas of adaptive optics (AO), scanning laser ophthalmoscopy (SLO), and optical coherence tomography (OCT), amongst others. Separately, there has been rapid growth of refractive surgery (LASIK, etc.), use of intraocular lenses (IOLs), and other forms of visual corrections. And, we do not overlook the incredible expansion of information technology, the broad utilization of computer, video, and other forms of displays which have been employed in myriad applications (with associated implications in vision).

We want to call the reader’s attention to the three cover illustrations. Here, of course, our choices were many! We have chosen one for its historical value, it is a photograph taken of a modern young lady viewing herself in an obsidian mirror in bright sunlight. That obsidian mirror, buried twice for extended time periods in its history, is ca. 8000 years old (!), and it is one of a number of the oldest known mirrors. These items are displayed in the Museum of Anatolian Civilizations, located in Ankara, Turkey² and/or at the Konya Museum in Konya which is in the south-central valley of Turkey and is located near to the dig site. This photograph falls into the evolving field of archaeological optics (not treated in this edition).^{3,4} Please consider the quality of the image in that “stone-age” mirror which was manufactured during the mesolithic or epipaleolithic period! A second figure displays a wave-guide modal pattern obtained radiating from a single human or primate photoreceptor in the early 1960s. There is further discussion of this topic in Chap. 8 on biological waveguides. The third figure is of human parafoveal cone photoreceptors taken from the living human eye by Austin Roorda (see Chap. 15). It was obtained using adaptive optics technology. The long (seen in red), middle (seen in green), and short (seen in blue) wavelength absorbing pigments contained in these individual cone photoreceptors are readily defined.

Please note, with the formation of the section on radiometry and photometry in this edition, the chapter addressing such measurements (as they pertain to visual optics), written by Dr. Yoshi Ohno, was relocated to Volume II, Chap. 37. A new, relatively brief, chapter on radiometry and photometry associated with the Stiles-Crawford effect (of the first kind) (Chap. 9) has been added. It was added after the chapter on biological waveguides (Chap. 8). This chapter fitted more logically there (where

the Stiles-Crawford effects are discussed) than in the new section on radiometry and photometry. The new chapter raises issues suggesting that revision is needed for specification of the visual stimulus (retinal illuminance) in certain test situations, particularly if the entrance pupil of the eye is larger than about 3 mm in diameter.

The outline of the “Vision” section of the this edition offers the reader a logical progression of topics addressed. Volume III leads off with an extensive chapter on optics of the eye written by Neil Charman (Chap. 1). This material has been considerably expanded from the earlier version in the second edition. William S. Geisler and Martin S. Banks reproduced their earlier chapter on visual performance (Chap. 2); and Denis G. Pelli and Bart Farrell similarly repeated their chapter on psychophysical methods used to test vision (Chap. 3). We are pleased that Prof. Gerald Westheimer wrote a new chapter on visual acuity and hyperacuity for this edition (Chap. 4). Professors Stephen A. Burns and Robert H. Webb repeated their chapter on optical generation of the visual stimulus (Chap. 5). Gerald Westheimer also kindly allowed the editors to reproduce a “classic” article he wrote some years ago in *Vision Research* on the topic of the Maxwellian view for inclusion in this edition. He also added a valuable addendum updating material in that discussion (Chap. 6). These chapters as a group provide a valuable introduction to this volume. As in other sections of the *Handbook of Optics*, all material written is intended to be at the first-year graduate student level. So saying, it is intended to be readable and readily appreciated by all parties.

The next set of topics is intended to broaden the discussion in several specific areas of interest. Chap. 7 by David H. Sliney addresses radiation hazards associated with vision and vision testing. Vasudevan Lakshminarayanan and Jay M. Enoch address biological waveguides in Chap. 8. This rapidly broadening subject grew out of work on the Stiles-Crawford effects, that is, “the directional sensitivity of the retina,” first reported in 1933. The 75th anniversary of this discovery was celebrated recently at the meeting of the Optical Society of America held in Rochester, New York, in 2008. In Chap. 9, Enoch and Lakshminarayanan speak of issues associated with the specification of the visual stimulus and the integration of the Stiles-Crawford effect of the first kind (SCE-1). These are meaningful matters associated with photometric and radiometric characterization of visual stimuli.

In Chaps. 10 and 11 David H. Brainard and Andrew Stockman address issues associated with color vision and colorimetry. The associate editors felt strongly that it was necessary to expand coverage of these topics in the third edition of the handbook. The chapter on refraction and refractive techniques (Chap. 12) was prepared by a different author, B. Ralph Chou, in this edition. He also broadened the topics covered from those treated in the second edition. Clifton Schor updated his chapter on binocular vision (Chap. 13), and John S. Werner, Brooke E. Scheffrin, and Arthur Bradley updated the discussion on the optics of vision and the aging eye (Chap. 14), a very important topic due to the rapid increase in aging of populations occurring worldwide!

The next portion of this volume addresses new/emerging technology. Donald T. Miller and Austin Roorda teamed-up to write about the rapidly evolving field of adaptive optics (AO) (Chap. 15). The reader will find that AO techniques are being combined with other emerging techniques in order to enhance utility of instruments, those both in development and also now appearing on the market. Included are scanning laser ophthalmoscopy (SLO), optical coherence tomography (OCT), and flood illumination. That is, these emerging technologies offer additional unique advantages. Interestingly, SLO technology originated some years ago in the laboratory of Prof. Robert H. Webb (Chap. 5). Optical coherence tomography (OCT), a powerful new tool useful in ophthalmic examinations (both for study of the anterior and posterior segments of the eye), is discussed by the researcher, Prof. Dr. Johannes F. deBoer (Chap. 18). New techniques for refractive surgery have been addressed by Harilaos Ginis and L. Diaz-Santana in Chap. 16. Dr. Barry R. Masters has considered confocal imaging of the cornea in Chap. 17; and Prof. Barbara K. Pierscionek has addressed the current state of graded index of refraction in the eye lens (GRIN profiles) in Chap. 19. Perhaps the Pierscionek chapter might have been better placed earlier in the order of things. Edward S. Bennett addresses the always very lively field of contact lens optics in Chap. 20 and Dr. Jim Schweigerling considers the optics of intraocular lenses in Chap. 21.

Clearly, we cannot overlook imaging and display problems associated with modern optical science. Thus, from an information processing point of view as applied to optical problems, William Cowan considers displays for vision research in Chap. 22. And Jeffery Anshel has added to, modified

and updated, the chapter which James E. Sheedy had written in the second edition of the *Handbook of Optics* (Chap. 23). That chapter addressed visual problems and needs of the observer when using computers—often for long periods of time. Bernice Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allerbach discussed human vision and electronic imaging, which is a major issue in the modern environment (Chap. 24). Finally, Brian H. Tsou and Martin Shenker address visual problems associated with heads-up displays (Chap. 25). The latter problems have been a particular challenge in the aviation industry.

Thus, we have tried to cover reasonably “the waterfront” of interactions between man/eye and instruments (human engineering/ergonomics, if one prefers). And we have sought to address directly issues related to optics per se. These changes have resulted in a longer volume than in the past. So saying, we wish to emphasize that the material is not encyclopedic; that is, we wish we had more material on eye movements, as well as subjects such as aniseikonia or problems encountered by individuals with unequal image sizes in their two eyes. Relating man to the optical instruments or images presented is no small thing from a number of points of view. So saying, we sought to achieve a reasonable balance in topics presented, and in the lengths of these discussions. Obviously, we were constrained by time and availability of authors. We thank all of our diligent group of authors spread out around the world and their helpers, as well as the editorial team at McGraw-Hill and those representing the Optical Society of America, particularly Editor-in-Chief, Michael Bass!

REFERENCES

1. H. von Helmholtz, *Handbuch der Physiologischen Optik*, 3d ed., 3 volumes, 1911; with commentary by A. Gullstrand, J. von Kries, W. Nagel., and a chapter by Christine Ladd-Franklin. *Handbook/Treatise on Physiological Optics*, translated to English by H. P. C. Southall. Published by the Optical Society of America, 1924. Reprinted by Dover Press, New York City, 1962. (The three original volumes were combined into two print volumes in the Dover Edition.)
2. Jay M. Enoch, “History of Mirrors Dating Back 8000 Years,” *Optometry and Vision Science* **83**:775–781, 2006.
3. Jay M. Enoch, “Archaeological Optics,” chap. 27 in Arthur H. Guenther (ed.): *International Trends in Applied Optics*, International Commission on Optics, vol. 5, Bellingham, Wash., SPIE Press, Monograph: PM 119. 2002, pp. 629–666. (ISBN: 0-8194-4510-X).
4. Jay M. Enoch, “Archeological Optics,” *Journal of Modern Optics* **54**:1221–1239, 2007.

Jay M. Enoch and Vasudevan Lakshminarayanan
Associate Editors

This page intentionally left blank.

DO NOT DUPLICATE

GLOSSARY AND FUNDAMENTAL CONSTANTS

Introduction

This glossary of the terms used in the Handbook represents to a large extent the language of optics. The symbols are representations of numbers, variables, and concepts. Although the basic list was compiled by the author of this section, all the editors have contributed and agreed to this set of symbols and definitions. Every attempt has been made to use the same symbols for the same concepts throughout the entire Handbook, although there are exceptions. Some symbols seem to be used for many concepts. The symbol α is a prime example, as it is used for absorptivity, absorption coefficient, coefficient of linear thermal expansion, and more. Although we have tried to limit this kind of redundancy, we have also bowed deeply to custom.

Units

The abbreviations for the most common units are given first. They are consistent with most of the established lists of symbols, such as given by the International Standards Organization ISO¹ and the International Union of Pure and Applied Physics, IUPAP.²

Prefixes

Similarly, a list of the numerical prefixes¹ that are most frequently used is given, along with both the common names (where they exist) and the multiples of ten that they represent.

Fundamental Constants

The values of the fundamental constants³ are listed following the sections on SI units.

Symbols

The most commonly used symbols are then given. Most chapters of the Handbook also have a glossary of the terms and symbols specific to them for the convenience of the reader. In the following list, the symbol is given, its meaning is next, and the most customary unit of measure for the quantity is presented in brackets. A bracket with a dash in it indicates that the quantity is unitless. Note that there is a difference between units and dimensions. An angle has units of degrees or radians and a solid angle square degrees or steradians, but both are pure ratios and are dimensionless. The unit symbols as recommended in the SI system are used, but decimal multiples of some of the dimensions are sometimes given. The symbols chosen, with some cited exceptions, are also those of the first two references.

RATIONALE FOR SOME DISPUTED SYMBOLS

The choice of symbols is a personal decision, but commonality improves communication. This section explains why the editors have chosen the preferred symbols for the *Handbook*. We hope that this will encourage more agreement.

Fundamental Constants

It is encouraging that there is almost universal agreement for the symbols for the fundamental constants. We have taken one small exception by adding a subscript B to the k for Boltzmann's constant.

Mathematics

We have chosen i as the imaginary unit arbitrarily. IUPAP lists both i and j , while ISO does not report on these.

Spectral Variables

These include expressions for the wavelength λ , frequency ν , wave number σ , ω for circular or radian frequency, k for circular or radian wave number and dimensionless frequency x . Although some use f for frequency, it can be easily confused with electronic or spatial frequency. Some use $\tilde{\nu}$ for wave number, but, because of typography problems and agreement with ISO and IUPAP, we have chosen σ ; it should not be confused with the Stefan-Boltzmann constant. For spatial frequencies we have chosen ξ and η , although f_x and f_y are sometimes used. ISO and IUPAP do not report on these.

Radiometry

Radiometric terms are contentious. The most recent set of recommendations by ISO and IUPAP are L for radiance [$\text{Wcm}^{-2}\text{sr}^{-1}$], M for radiant emittance or exitance [Wcm^{-2}], E for irradiance or incidence [Wcm^{-2}], and I for intensity [Wsr^{-2}]. The previous terms, W , H , N , and J , respectively, are still in many texts, notably Smith⁴ and Lloyd⁵ but we have used the revised set, although there are still shortcomings. We have tried to deal with the vexatious term *intensity* by using *specific intensity* when the units are $\text{Wcm}^{-2}\text{sr}^{-1}$, *field intensity* when they are Wcm^{-2} , and *radiometric intensity* when they are Wsr^{-1} .

There are two sets of terms for these radiometric quantities, which arise in part from the terms for different types of reflection, transmission, absorption, and emission. It has been proposed that the *ion* ending indicate a process, that the *ance* ending indicate a value associated with a particular sample, and that the *ivity* ending indicate a generic value for a "pure" substance. Then one also has reflectance, transmittance, absorptance, and emittance as well as reflectivity, transmissivity, absorptivity, and emissivity. There are now two different uses of the word emissivity. Thus the words *exitance*, *incidence*, and *sterance* were coined to be used in place of emittance, irradiance, and radiance. It is interesting that ISO uses radiance, exitance, and irradiance whereas IUPAP uses radiance, exitance [*sic*], and irradiance. We have chosen to use them both, i.e., emittance, irradiance, and radiance will be followed in square brackets by exitance, incidence, and sterance (or vice versa). Individual authors will use the different endings for transmission, reflection, absorption, and emission as they see fit.

We are still troubled by the use of the symbol E for irradiance, as it is so close in meaning to electric field, but we have maintained that accepted use. The spectral concentrations of these quantities, indicated by a wavelength, wave number, or frequency subscript (e.g., L_λ) represent partial differentiations; a subscript q represents a photon quantity; and a subscript ν indicates a quantity normalized to the response of the eye. Thereby, L_ν is luminance, E_ν illuminance, and M_ν and I_ν luminous emittance and luminous intensity. The symbols we have chosen are consistent with ISO and IUPAP.

The refractive index may be considered a radiometric quantity. It is generally complex and is indicated by $\tilde{n} = n - ik$. The real part is the relative refractive index and k is the extinction coefficient. These are consistent with ISO and IUPAP, but they do not address the complex index or extinction coefficient.

Optical Design

For the most part ISO and IUPAP do not address the symbols that are important in this area.

There were at least 20 different ways to indicate focal ratio; we have chosen FN as symmetrical with NA; we chose f and efl to indicate the effective focal length. Object and image distance, although given many different symbols, were finally called s_o and s_i since s is an almost universal symbol for distance. Field angles are θ and ϕ ; angles that measure the slope of a ray to the optical axis are u ; u can also be $\sin u$. Wave aberrations are indicated by W_{ijk} , while third-order ray aberrations are indicated by σ_i and more mnemonic symbols.

Electromagnetic Fields

There is no argument about \mathbf{E} and \mathbf{H} for the electric and magnetic field strengths, Q for quantity of charge, ρ for volume charge density, σ for surface charge density, etc. There is no guidance from Refs. 1 and 2 on polarization indication. We chose \perp and \parallel rather than p and s , partly because s is sometimes also used to indicate scattered light.

There are several sets of symbols used for reflection transmission, and (sometimes) absorption, each with good logic. The versions of these quantities dealing with field amplitudes are usually specified with lower case symbols: r , t , and a . The versions dealing with power are alternately given by the uppercase symbols or the corresponding Greek symbols: R and T versus ρ and τ . We have chosen to use the Greek, mainly because these quantities are also closely associated with Kirchhoff's law that is usually stated symbolically as $\alpha = \epsilon$. The law of conservation of energy for light on a surface is also usually written as $\alpha + \rho + \tau = 1$.

Base SI Quantities

length	m	meter
time	s	second
mass	kg	kilogram
electric current	A	ampere
temperature	K	kelvin
amount of substance	mol	mole
luminous intensity	cd	candela

Derived SI Quantities

energy	J	joule
electric charge	C	coulomb
electric potential	V	volt
electric capacitance	F	farad
electric resistance	Ω	ohm
electric conductance	S	siemens
magnetic flux	Wb	weber
inductance	H	henry
pressure	Pa	pascal
magnetic flux density	T	tesla
frequency	Hz	hertz
power	W	watt
force	N	newton
angle	rad	radian
angle	sr	steradian

Prefixes

Symbol	Name	Common name	Exponent of ten
F	exa		18
P	peta		15
T	tera	trillion	12
G	giga	billion	9
M	mega	million	6
k	kilo	thousand	3
h	hecto	hundred	2
da	deca	ten	1
d	deci	tenth	-1
c	centi	hundredth	-2
m	milli	thousandth	-3
μ	micro	millionth	-6
n	nano	billionth	-9
p	pico	trillionth	-12
f	femto		-15
a	atto		-18

Constants

c	speed of light vacuo [299792458 ms ⁻¹]
c_1	first radiation constant = $2\pi^2 h = 3.7417749 \times 10^{-16}$ [Wm ²]
c_2	second radiation constant = $hc/k = 0.014838769$ [mK]
e	elementary charge [$1.60217733 \times 10^{-19}$ C]
g_n	free fall constant [9.80665 ms ⁻²]
h	Planck's constant [$6.6260755 \times 10^{-34}$ Ws]
k_B	Boltzmann constant [1.380658×10^{-23} JK ⁻¹]
m_e	mass of the electron [$9.1093897 \times 10^{-31}$ kg]
N_A	Avogadro constant [6.0221367×10^{23} mol ⁻¹]
R_∞	Rydberg constant [10973731.534 m ⁻¹]
ϵ_0	vacuum permittivity [$\mu_0^{-1}c^{-2}$]
σ	Stefan-Boltzmann constant [5.67051×10^{-8} Wm ⁻¹ K ⁻⁴]
μ_0	vacuum permeability [$4\pi \times 10^{-7}$ NA ⁻²]
μ_B	Bohr magneton [$9.2740154 \times 10^{-24}$ JT ⁻¹]

General

B	magnetic induction [Wbm ⁻² , kgs ⁻¹ C ⁻¹]
C	capacitance [f, C ² s ² m ⁻² kg ⁻¹]
C	curvature [m ⁻¹]
c	speed of light in vacuo [ms ⁻¹]
c_1	first radiation constant [Wm ²]
c_2	second radiation constant [mK]
D	electric displacement [Cm ⁻²]
E	incidence [irradiance] [Wm ⁻²]
e	electronic charge [coulomb]
E_v	illuminance [lux, lmm ⁻²]
E	electrical field strength [Vm ⁻¹]
E	transition energy [J]
E_g	band-gap energy [eV]
f^g	focal length [m]
f_f	Fermi occupation function, conduction band
f_v	Fermi occupation function, valence band

FN	focal ratio (f /number) [—]
g	gain per unit length [m^{-1}]
g_{th}	gain threshold per unit length [m^{-1}]
H	magnetic field strength [Am^{-1} , $\text{Cs}^{-1} \text{m}^{-1}$]
h	height [m]
I	irradiance (see also E) [Wm^{-2}]
I	radiant intensity [Wsr^{-1}]
I	nuclear spin quantum number [—]
I	current [A]
i	$\sqrt{-1}$
$\text{Im}()$	imaginary part of
J	current density [Am^{-2}]
j	total angular momentum [$\text{kg m}^2 \text{sec}^{-1}$]
$J_1()$	Bessel function of the first kind [—]
k	radian wave number $=2\pi/\lambda$ [rad cm^{-1}]
k	wave vector [rad cm^{-1}]
k	extinction coefficient [—]
L	sterance [radiance] [$\text{Wm}^{-2} \text{sr}^{-1}$]
L_v	luminance [cdm^{-2}]
L	inductance [h, $\text{m}^2 \text{kg C}^2$]
L	laser cavity length
L, M, N	direction cosines [—]
M	angular magnification [—]
M	radiant exitance [radiant emittance] [Wm^{-2}]
m	linear magnification [—]
m	effective mass [kg]
MTF	modulation transfer function [—]
N	photon flux [s^{-1}]
N	carrier (number) density [m^{-3}]
n	real part of the relative refractive index [—]
\tilde{n}	complex index of refraction [—]
NA	numerical aperture [—]
OPD	optical path difference [m]
P	macroscopic polarization [C m^{-2}]
$\text{Re}()$	real part of [—]
R	resistance [Ω]
r	position vector [m]
S	Seebeck coefficient [VK^{-1}]
s	spin quantum number [—]
s	path length [m]
S_o	object distance [m]
S_i	image distance [m]
T	temperature [K, C]
t	time [s]
t	thickness [m]
u	slope of ray with the optical axis [rad]
V	Abbe reciprocal dispersion [—]
V	voltage [V , $\text{m}^2 \text{kg s}^{-2} \text{C}^{-1}$]
x, y, z	rectangular coordinates [m]
Z	atomic number [—]

Greek Symbols

α	absorption coefficient [cm^{-1}]
α	(power) absorptance (absorptivity)

ϵ	dielectric coefficient (constant) [—]
ϵ	emittance (emissivity) [—]
ϵ	eccentricity [—]
ϵ_1	Re (ϵ)
ϵ_2	Im (ϵ)
τ	(power) transmittance (transmissivity) [—]
ν	radiation frequency [Hz]
ω	circular frequency = $2\pi\nu$ [rads ⁻¹]
ω	plasma frequency [Hz]
λ	wavelength [μm , nm]
σ	wave number = $1/\lambda$ [cm ⁻¹]
σ	Stefan Boltzmann constant [Wm ⁻² K ⁻¹]
ρ	reflectance (reflectivity) [—]
θ, ϕ	angular coordinates [rad, °]
ξ, η	rectangular spatial frequencies [m ⁻¹ , r ⁻¹]
ϕ	phase [rad, °]
ϕ	lens power [m ⁻²]
Φ	flux [W]
χ	electric susceptibility tensor [—]
Ω	solid angle [sr]

Other

\Re	responsivity
$\exp(x)$	e^x
$\log_a(x)$	log to the base a of x
$\ln(x)$	natural log of x
$\log(x)$	standard log of x : $\log_{10}(x)$
Σ	summation
Π	product
Δ	finite difference
δx	variation in x
dx	total differential
∂x	partial derivative of x
$\delta(x)$	Dirac delta function of x
δ_{ij}	Kronecker delta

REFERENCES

1. Anonymous, *ISO Standards Handbook 2: Units of Measurement*, 2nd ed., International Organization for Standardization, 1982.
2. Anonymous, *Symbols, Units and Nomenclature in Physics*, Document U.I.P. 20, International Union of Pure and Applied Physics, 1978.
3. E. Cohen and B. Taylor, "The Fundamental Physical Constants," *Physics Today*, 9 August 1990.
4. W. J. Smith, *Modern Optical Engineering*, 2nd ed., McGraw-Hill, 1990.
5. J. M. Lloyd, *Thermal Imaging Systems*, Plenum Press, 1972.

William L. Wolfe
 College of Optical Sciences
 University of Arizona
 Tucson, Arizona

OPTICS OF THE EYE

Neil Charman

*Department of Optometry and Vision Sciences
University of Manchester
Manchester, United Kingdom*

1.1 GLOSSARY

F, F' focal points
 N, N' nodal points
 P, P' principal points

Equation (1)

r distance from axis
 R_0 radius of curvature at corneal pole
 p corneal asphericity parameter

Equation (2)

s distance from Stiles-Crawford peak
 η/η_{\max} relative luminous efficiency
 ρ coefficient in S-C equation

Equation (3)

d pupil diameter
 P_A ratio of effective to true pupil area

Transmittance and reflectance

$T_E(\lambda)$ total transmittance of the eye media
 $R_R(\lambda)$ reflectance of the retina
 λ wavelength

Equation (4)

(θ, ϕ) angular direction coordinates in visual field
 $\delta(\lambda)$ wavelength interval
 $L_{e\lambda}(\theta, \phi)$ spectral radiance per unit wavelength interval per unit solid angle in direction (θ, ϕ)
 $p(\theta, \phi)$ area of pupil as seen from direction (θ, ϕ)

$t(\theta, \phi, \lambda)$ fraction of incident radiation flux which is transmitted by the eye
 $m(\theta, \phi, \lambda)$ areal magnification factor

Equation (5)

I normalized illuminance
 z dimensionless diffraction unit

Equation (6)

γ angular distance from center of Airy diffraction pattern
 d pupil diameter

Equation (7)

θ_{\min} angular resolution by Rayleigh criterion

Equation (8)

R spatial frequency
 R_R reduced spatial frequency

Equation (9)

ΔF dioptric error of focus
 g number of Rayleigh units of defocus

Equation (10)

β angular diameter of retinal blur circle

Equation (11)

$T(R)$ modulation transfer function

Equation (13)

$R_x(\lambda)$ chromatic difference in refraction with respect to 590 nm

Equation (14)

L_{eq} equivalent veiling luminance
 E illuminance produced by glare source at eye
 ω angle between direction of glare source and visual axis

Equations (15) and (16)

$M_{IT}(R)$ threshold modulation on the retina
 $M_{OT}(R)$ external threshold modulation

Equation (17)

DOF_{go} total depth-of-focus for an aberration-free eye according to geometrical optics
 ΔF_{tol} tolerable error of focus
 β_{tol} tolerable angular diameter of retinal blur circle

Equation (18)

DOF_{po} total depth-of-focus for an aberration-free eye according to physical optics

Equation (19)

OA objective amplitude of accommodation

Equation (20)

l object distance
 p interpupillary distance
 δl minimum detectable difference in distance
 $\delta \theta$ stereo acuity

Equation (21)

- M transverse magnification
 N factor by which effective interpupillary distance is increased

1.2 INTRODUCTION

The human eye (Fig. 1) contains only a few optical components. However, in good lighting conditions when the pupil is small (2 to 3 mm), it is capable of near diffraction-limited performance close to its axis. Each individual eye also has a very wide field of view (about 65, 75, 60, and 95 deg in the superior, inferior, nasal, and temporal semimeridians, respectively, for a fixed frontal direction of gaze, the exact values being dependent upon the individual's facial geometry). The binocular field, where the two monocular fields overlap, has a lateral extent of about 120 deg. Optical image quality, while somewhat degraded in the peripheral field is, in general, adequate to meet the needs of the neural network which it serves, since the spatial resolution of the neural retina falls rapidly away from the visual axis (the latter joins the point of regard, the nodal points and the fovea). The orientation of the visual axis typically differs by a few degrees from that of the optical axis, as the fovea, where neural resolution is optimal, is usually slightly displaced from the intersection of the optical axis with the retina.^{1,2} Control of ocular aberrations is helped by aspheric optical surfaces and by the gradients of refractive index in the lens, the lens index progressively reducing from the lens center toward its outer layers. Off-axis aberrations are further reduced by the eye's approximation to a homocentric system, in which the optical and detector surfaces are concentric with a common center of curvature at the aperture stop.³ Although aberration levels increase and optical image quality falls as the pupil dilates at lower light levels (to reach a maximum diameter of about 8 mm, corresponding to a numerical aperture of about 0.25 mm), neural performance also declines, so that optical and neural performances remain reasonably well matched. When the eye is in its basic "relaxed" state it is nominally in focus for distant objects. In the younger eye (<50 years) the power of the crystalline lens can be increased to allow near objects to be clearly focused, a process known as *accommodation*. These general characteristics will now be discussed in more detail.

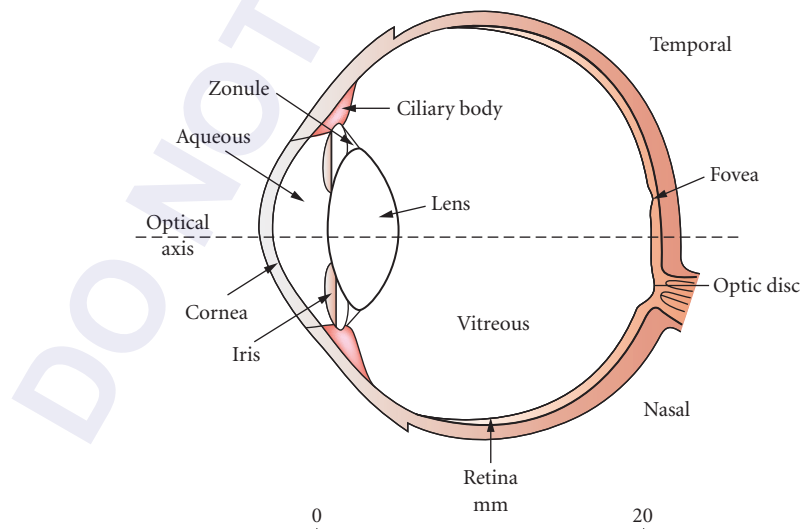


FIGURE 1 Schematic horizontal section of the eye. The bar gives the approximate scale.

1.3 OCULAR PARAMETERS AND AMETROPIA

Variation in Ocular Parameters across the Population

In the first few years of life there is substantial overall ocular growth, with axial length increasing from about 17 mm at birth to about 23 mm by the age of 3 years.⁴⁻⁶ Some growth continues thereafter, the axial length in adults being about 24 mm.⁷⁻¹² Although most dimensions remain almost stable during adulthood, the lens continues to grow in thickness and volume, but not diameter, throughout life. The thickness increases approximately linearly by about 50 percent between birth and the age of 70,^{10,13-15} there being accompanying changes in its index distribution.¹⁶⁻¹⁹ Surface curvatures, component separations, and axial lengths show considerable variation (~10 percent) between individuals of the same age, although the refractive indices of the cornea, vitreous, and aqueous humours are essentially constant.^{7-9,20,21} Figure 2 shows some typical measured adult distributions of the values of several parameters: most of the distributions shown are approximately normal (dashed curves), but this is not true for the axial length.

As noted earlier, the optical surfaces may be aspheric. The form of the anterior cornea is particularly significant. Owing to the large refractive index change at its anterior surface it contributes about three-quarters of the total refractive power of the eye. Its topography also has obvious relevance to contact lens design and to the monochromatic aberrations of the eye. It is often modeled as a conicoid, of the form

$$r^2 + pz^2 - 2R_0z = 0 \quad (1)$$

where the axis of symmetry lies in the z direction, r is the distance perpendicular to the axis, and R_0 is the radius of curvature at the pole of the surface. Figure 3 shows experimental measurements of the distribution of the parameter p .^{22,23} The distribution is fairly wide and peaks at a value of

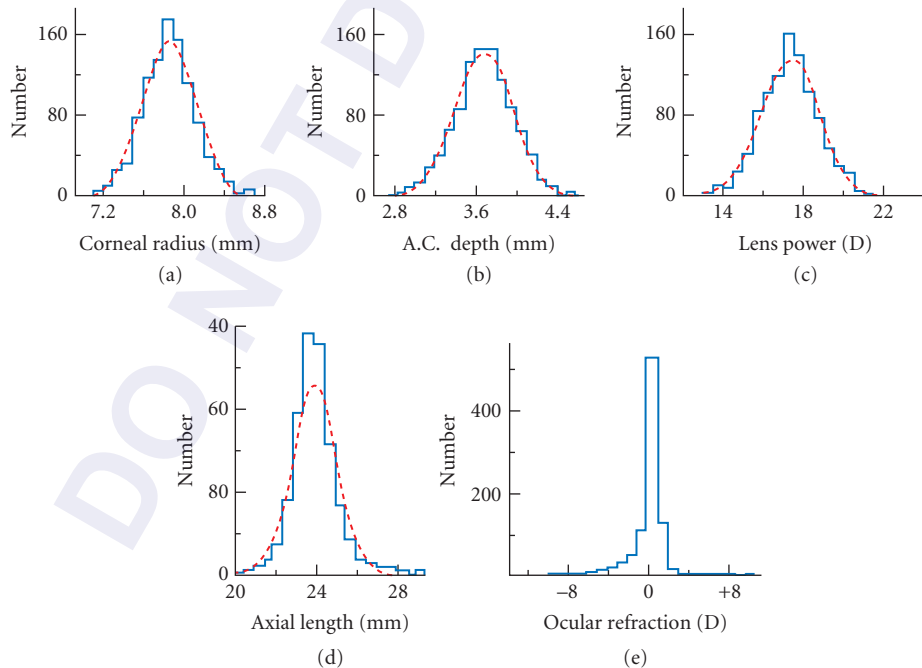


FIGURE 2 Distributions of some dimensional parameters for the adult human eye. (After Stenström.²¹)

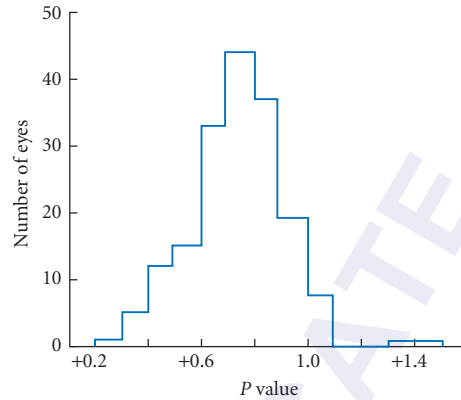


FIGURE 3 Histogram showing the distribution of the corneal asphericity parameter p . A value of $p = 1$ corresponds to a spherical surface, values between 0 and 1.0 represent prolate (flattening) ellipsoids and $p > 1.0$ oblate (steepening) ellipsoids. (Based on Refs. 22 and 23.)

about +0.8, corresponding to a flattening ellipsoid in which the radius of curvature is smallest at the center of the cornea and increases toward the corneal periphery. At the corneal vertex the radius of curvature is about 7.8 ± 0.25 mm.^{22,23}

The least understood optical feature of the eye is the distribution of refractive index within the lens. As noted earlier, the lens grows throughout life,^{10,13,14} with new material being added to the surface layers (the cortex). The oldest part of the lens is its central region (the nucleus). While there is general agreement that the refractive index is highest at the lens center and falls toward its outer layers, the exact form of the gradients involved has proved difficult to measure.^{16–19} Description is complicated by the fact that the shape of the lens and its gradients change when the eye accommodates to view near objects and with age. To illustrate the general trend of the changes with age, Fig. 4 shows some

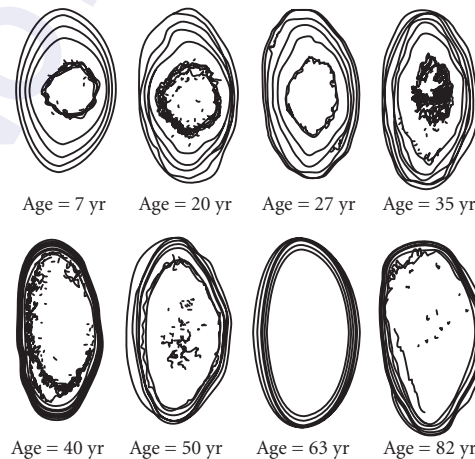


FIGURE 4 Contours of refractive index in lenses of different ages (7 to 82 years). The contour interval is 0.01. (After Ref. 18.)

recent *in vitro* iso index contours for isolated lenses, obtained using magnetic resonance imaging: when measured *in vitro* the lenses take up their fully accommodated form. It can be seen that, with age, the region of relatively constant index at the lens center increases in volume and that the gradient becomes almost entirely confined to the surface layers of the lens. The central index remains constant at 1.420 ± 0.075 and the surface index at 1.371 ± 0.004 .¹⁸ These index distributions have been modeled as a function of age and accommodation by several authors. (see, e.g., Refs. 24–30)

In addition to these general variations, each eye may have its own idiosyncratic peculiarities, such as small tilts or lateral displacements of surfaces, lack of rotational symmetry about the axis, or irregularities in the shape or centration of the pupil.^{1,2,31} These affect both the refractive error (*ametropia*) and the higher-order aberrations.

Ocular Ametropia

If the combination of ocular parameters is such that, with accommodation relaxed, a distant object of regard is focused on the retinal fovea, the region where the density of the cone receptors is highest and photopic neural performance is optimal (Fig. 1), the eye is *emmetropic*. This condition is often not achieved, in which case the eye is *ametropic*. If the power of the optical elements is too great for the axial length, so that the image of the distant object lies anterior to the retina, the eye is *myopic*. If, however, the power is insufficient, the eye is *hypermetropic* (or *hyperopic*). These defects can be corrected by the use of appropriately powered diverging (myopia) or converging (hypermetropia) spectacle or contact lenses to respectively reduce or increase the power of the lens-eye combination. (see Chap. 20 by Edward S. Bennett and William J. Benjamin for reviews.) Spherical ametropia tends to be associated with axial length differences, that is, myopic eyes tend to be longer than emmetropic eyes while hyperopic eyes are shorter.³²

In some individuals the ocular dioptrics lack rotational symmetry, one or more optical surfaces being toroidal, tilted, or displaced from the axis. This leads to the condition of ocular astigmatism, in which on the visual axis two longitudinally separated, mutually perpendicular, line images of a point object are formed. In the vast majority of cases (*regular astigmatism*) the meridians of maximal and minimal power are perpendicular: there is a strong tendency for the principal meridians to be approximately horizontal and vertical but this is not always the case. Eyes in which the more powerful meridian is vertical are often described as having *with-the-rule astigmatism* and those in which it is horizontal as having *against-the-rule astigmatism*. The former is more common. Correction of astigmatism can be achieved by including an appropriately oriented cylindrical component in any correcting lens. It is sometimes convenient to talk of the *best-, mean-, or equivalent-sphere* correction. This is the power of spherical lens which brings the circle of least confusion onto the retina: its value is $S + C/2$, where S and C are respectively, the spherical and cylindrical dioptric components of the correction.

In addition to spectacle and contact lens corrections, surgical methods of correction for both spherical and astigmatic errors are now widely used. Most common are those using excimer lasers which essentially reshape the anterior surface of the cornea by selectively ablating its tissue across the chosen area to appropriately modify its spherocylindrical power. A popular current method is laser-assisted keratomileusis (LASIK) which involves cutting a thin, uniform “flap” of material from the anterior cornea and then ablating the underlying corneal stroma to change its curvature. The flap is then replaced. (see Chap. 16 by L. Diaz-Santana and Harilaos Giniş for details.) Intraocular lenses can be used to replace the crystalline lens when the latter has lost transparency due to cataract: single-vision, bifocal and multifocal designs are available and efforts are being made to develop lenses of dynamically varying power to simulate the accommodative abilities of the younger eye. (see Chap. 21 by Jim Schwiegerling.)

Figure 5 shows representative data for the frequency of occurrence of different spherical²¹ and astigmatic³⁴ errors in western adults. Myopia is more common in many eastern populations. Note particularly that the spherical errors are not normally distributed and that a state of near-emmetropia is most common. It is believed that the correlation of component values required to achieve near-emmetropia is achieved partly as a result of genetic factors and partly as a result of

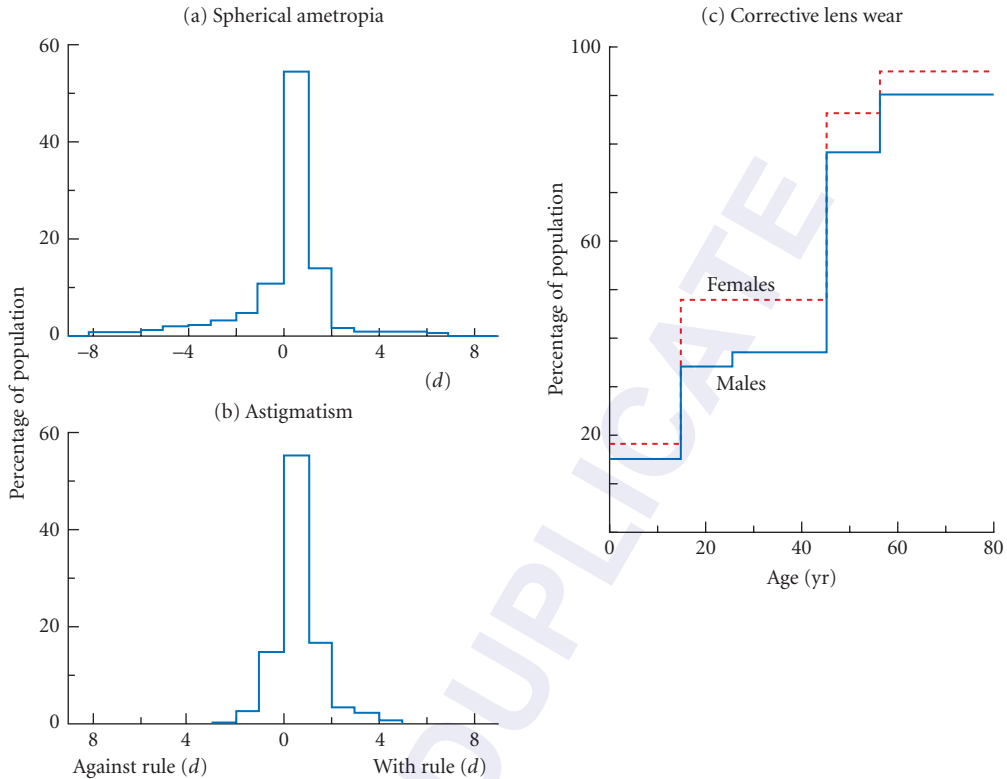


FIGURE 5 Typical adult data for the frequency of occurrence of spherical and cylindrical (astigmatic) refractive errors, together with the fraction of the population wearing corrective spectacle or contact lenses. (a) Spherical errors (diopters) in young adult males. (After Stenstrom.²¹) (b) Cylindrical errors. (Based on Lyle.³⁴) Cases in which the meridian of greatest power is within 30 deg of the horizontal (against-the-rule) and within 30 deg of the vertical (with-the-rule) are shown: the remaining 4 percent of the population have axes in oblique meridians. (c) Percentage of the population wearing lenses, as a function of age. (After Farrell and Booth.³⁵)

environmentally influenced growth processes which drive the development of the young eye toward emmetropia (*emmetropization*).³³

Not all individuals with ametropia actually wear a correction; the fraction of the population that typically does so is shown in Fig. 5. The increase in lens wear beyond the age of 40 is due to the need for a near correction for close work, a condition known as *presbyopia*. This arises as a result of the natural, progressive failure with age of the eye's own accommodation system (see "Age-dependent Changes in Accommodation").

The widespread existence of ametropia among users of visual instruments such as telescopes and microscopes, means that it is desirable to make provision for focusing the eyepiece to compensate for any spherical refractive error of the observer. This is particularly the case where the eyepiece contains a graticule. Since the refractive errors of the two eyes of an individual may not be identical (*anisometropia*), differential focusing should be provided for the eyepieces of binocular instruments. As correction for cylindrical errors is inconvenient to incorporate into eyepieces, astigmatic users of instruments must usually wear their normal refractive correction. For spectacle wearers, where the distance of the lenses in front of the eyes is usually 10 to 18 mm, this implies that the exit pupil of the instrument must have an adequate eye clearance or eye relief (at least 20 mm and preferably 25 mm) to avoid contact between the spectacle lens and the eyepiece and allow the instrument's full field to be seen.

1.4 OCULAR TRANSMITTANCE AND RETINAL ILLUMINANCE

The amount, spectral distribution, and polarization properties of the light reaching the retina are modified with respect to the original stimulus in a way that depends upon the pupil diameter and the transmittance characteristics of the eye.

Pupil Diameter

The circular opening in the iris, located approximately tangential to the anterior surface of the lens, plays the important role of aperture stop of the eye. It therefore controls the amount of light flux reaching the retina, as well as influencing retinal image quality through its effects on diffraction, aberration, and depth-of-focus (see Sec. 1.7). It may also affect the amount of scattered light reaching the retina, particularly in older eyes where cataract is present.

What is normally measured and observed is the image of the true pupil as viewed through the cornea, that is, the entrance pupil of the eye. This is some 13 percent larger in diameter than the true pupil. Although ambient lighting and its spatial distribution have the most important influence on entrance pupil diameter (Fig. 6) the latter is also affected by many other factors including age, accommodation, emotion, and drugs.^{36,37} For any scene luminance, the pupils are slightly smaller under binocular conditions of observation.³⁸ The gradual constriction with age³⁶ helps to account for the poorer visual performance of older individuals under dim lighting conditions in comparison with younger individuals.³⁹

The pupil can respond to changes in light level at frequencies up to about 4 Hz.³⁷ Shifts in pupil center of up to 0.6 mm may occur when the pupil dilates^{40,41} and these may be of some significance in relation to the pupil-dependence of ocular aberration and retinal image quality.

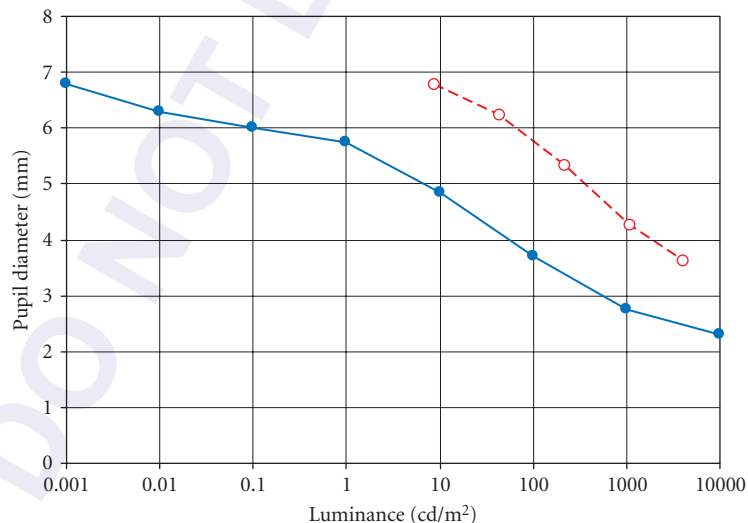


FIGURE 6 Entrance pupil diameter as a function of scene luminance for an extended visual field and young observers: the filled symbols and full curve show the weighted average of 6 studies. (After Farrell and Booth.³⁵) Larger pupils are observed when the illuminated field is of smaller area: the dashed curve and open symbols show data for a 10 deg illuminated field. (After Winn et al.³⁶)

It has been suggested^{42,43} that the major value of light-induced pupillary constriction is that it reduces retinal illuminance and hence prepares the eye for a return to darkness: following a change to a dark environment the dilation of the mobile pupil allows substantially better stimulus detection during the first few minutes of dark adaptation than would be found with a fixed pupil.

Transmittance

Light may be lost by spectrally varying reflection, scattering, and absorption in any of the media anterior to the retina.⁴⁴ Fresnel reflection losses are in general small, the maximum being 3 to 4 percent at the anterior cornea. Wavelength-dependent absorption and scattering are much more important: both tend to increase with age.

The measured transmittance⁴⁵⁻⁴⁹ depends to some extent on the measuring technique, in particular on the extent to which scattered light is included, but representative data are shown in Fig. 7. The transmittance rises rapidly above about 400 nm, to remain high in the longer wavelength visible and near infrared. It then falls through several absorption bands, due mainly to water, to reach zero at about 1400 nm.

Although most of the absorption at ultraviolet wavelengths below 300 nm occurs at the cornea (where it may disrupt the surface cells, leading to photokeratitis, e.g., snow blindness or welder's flash, the lowest damage thresholds of about $0.4 \text{ J} \cdot \text{cm}^{-2}$ being at 270 nm⁵⁰), there is also substantial absorption in the lens at the short wavelength (roughly 300–400 nm) end of the visible spectrum. This lenticular absorption increases markedly with age,⁵¹⁻⁵⁴ the lens becoming progressively yellower in appearance,⁵⁵ and can adversely affect color vision.^{56,57} Most of the absorption occurs in the lens nucleus.⁵⁸ There is evidence that UV absorption in the lens may be a causative factor in some types of cataract.⁵⁹ Excessive visible light at the violet-blue end of the spectrum is thought to cause accelerated aging and resultant visual loss at the retinal level⁶⁰ so that lenticular absorption may have a protective function. (See also Chap. 7 by David H. Sliney.)

In the foveal region a thin layer of macular pigment, extending over the central few degrees of the retina^{61,62} and lying anterior to the receptor outer segments^{63,64} absorbs heavily at shorter wavelengths (Fig. 7). It has been argued that this absorption is helpful in reducing the blurring effects of longitudinal chromatic aberration⁶⁵ and in protecting the foveal receptors, which are responsible for detailed pattern vision, against blue-light damage.⁶⁶ It is, however, notable that the amount of macular pigment varies widely between individuals.⁶⁴

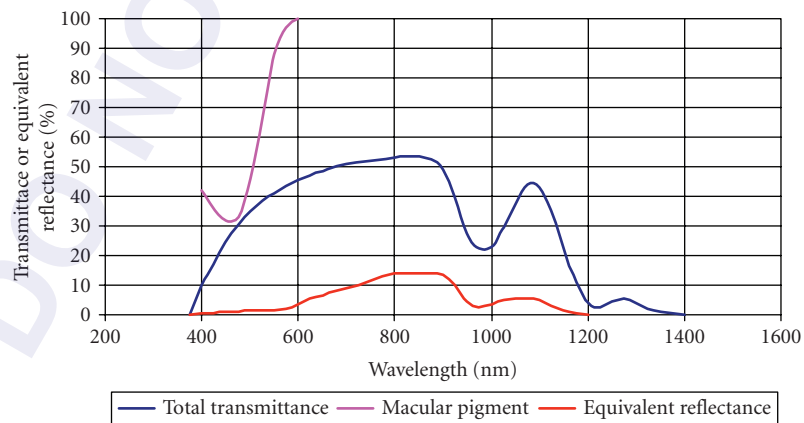


FIGURE 7 Spectral dependence of the overall transmittance of the ocular media⁴⁶ and the equivalent reflectance of the retina.⁷⁵⁻⁷⁷ Also shown is the transmittance at the fovea of the macular pigment.⁶⁵

Since the cornea, lens, and retinal nerve fibre layer all show birefringence, the polarization characteristics of the light entering the eye are modified before it reaches the outer segments of the retinal receptors. In general these effects are of little practical significance, although they can be demonstrated and measured by suitable methods.⁶⁷

The Stiles-Crawford Effect

One complicating factor when considering the effectiveness of the light flux which enters the eye as a stimulus to vision is the Stiles-Crawford effect of the first kind SCEI.⁶⁸ (See also Chaps. 8 and 9 by Jay M. Enoch and Vasudevan Lakshminarayanan in this volume.) This results in light which enters the periphery of the entrance pupil to reach a given retinal location being less effective at stimulating the retina than light which passes through the pupil center (Fig. 8). The effect varies slightly with the individual and is not always symmetric about the center of the pupil. Under photopic conditions, giving cone vision, there is typically a factor of about 8 between the central effectiveness and that at the edge of a fully dilated 8-mm pupil; the effect is much weaker under rod-dominated, scotopic conditions (see Fig. 8).⁶⁹⁻⁷¹ In practice, unless pupil-dilating drugs are used, the natural pupil will only be large under scotopic conditions and will normally be constricted at photopic levels (see Fig. 6): the influence of SCE I is still significant, however.

Many equations have been proposed to fit photopic data of the type illustrated in Fig. 8. The simplest, due to Stiles,⁷² can be written:

$$\log_{10}(\eta/\eta_{\max}) = -\rho s^2 \quad (2)$$

where η/η_{\max} is the relative luminous efficiency and s is the distance within the entrance pupil from the function peak (in mm). Values of the constant ρ equal to about 0.07 are typical, the value varying somewhat with wavelength.⁷²

It is evident that the Stiles-Crawford effect of the first kind results in the photopic retinal stimulus being somewhat weaker than that predicted on the basis of visible pupil area. This can be accounted for by using an effective pupil area instead of the actual entrance pupil area. Moon and Spencer⁷³ suggested that the ratio P_A of the effective to the true pupil areas could be approximated by:

$$P_A = 1 - 0.0106d^2 + 0.0000417d^4 \quad (3)$$

where d is the pupil diameter in mm.

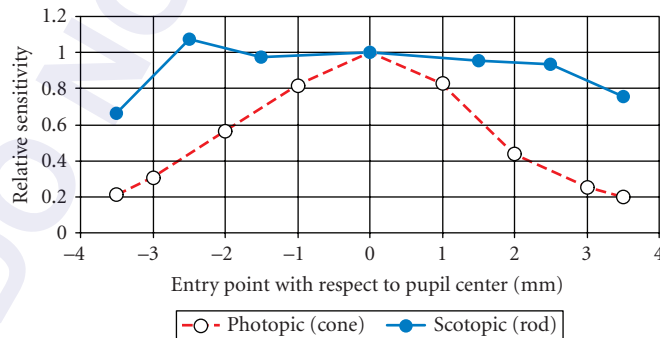


FIGURE 8 The Stiles-Crawford effect (SCE I) under photopic (open circles) and scotopic (filled circles) conditions for the same observer, measured at a position 6 deg from the fovea. The relative luminous efficiency is plotted as a function of the horizontal pupillary position of the beam. (After van Loo and Enoch.⁶⁹)

The Stiles–Crawford effect of the second kind (SCE II) involves a shift in the hue of monochromatic light as it enters different areas of the pupil.^{70,72,74}

Although Fresnel reflection and lenticular absorption variations with pupil position may play minor roles, there seems little doubt that the Stiles–Crawford effects mainly involve the waveguide properties of the outer segments of the small-diameter receptors. (See Refs. 63 and 64, and Chap. 8 by Vasudevan Lakshminarayanan and Jay M. Enoch for reviews.) Effectively, the receptor can only trap light efficiently if the latter is incident in directions within a few degrees of the receptor axis. For this reason, the receptor outer segments must always be aligned toward the exit pupil of the eye, rather than perpendicular to the local surface of the eyeball. An obvious advantage of such directional sensitivity is that it helps to suppress the degrading effects of intraocular stray light. Such light will be ineffective at stimulating the receptors if it is incident at oblique angles on the retina. As will be discussed in Chap. 8, it appears that SCE I acts as an amplitude apodizing filter in its effects upon retinal image quality.

Retinal Reflectance

The function of the optical system of the eye is to deliver an image to the retina. Nevertheless, a significant amount of this light is reflected when it reaches the retina and underlying structures. Although such light does not contribute usefully to the process of vision, its existence allows the development of a variety of clinical instruments for examination of the retina, measurement of refraction, and other purposes.

Due to its double passage through the eye media, the emergent flux at any wavelength is proportional to the equivalent reflectance $T_E(\lambda)^2 R_R(\lambda)$, where $T_E(\lambda)$ is the total transmittance of the eye media and $R_R(\lambda)$ is the true retinal reflectance. Equivalent reflectance rises with wavelength across the visible spectrum, to become quite high in the near infrared.^{62,75–78} Representative values are given Fig. 7. Absolute levels of equivalent reflectance are affected by the pigmentation of an individual eye. At the violet end of the visible spectrum the equivalent reflectance falls markedly with age,⁷⁹ due to the decreased transmittance of the lens. In the same spectral region, equivalent reflectance is usually lower within the immediate area of the fovea, due to the low transmittance of the macular pigment.

The high equivalent reflectance in the infrared is particularly useful in allowing measurements to be made of, for example, refraction or aberration, at wavelengths which are essentially invisible to the patient or subject. In practice, depending upon the wavelength, light may be reflected from structures anywhere between the anterior surface of the retina and the choroid/sclera interface. Although details of the nature of the reflections are still imperfectly understood, shorter visible wavelengths appear to penetrate less deeply before reflection occurs, while the infrared is reflected from the anterior sclera. At the shorter wavelengths the reflection is almost specular but becomes more diffuse as longer visible and infrared wavelengths.^{79–81} Waveguiding effects within the receptors may play some role in determining the nature of the reflection and the angular distribution of the reflected light.^{82–85}

Ocular Radiometry and Retinal Illuminance

If we confine ourselves to uniform object fields subtending at least 1 deg at the eye, so that blurring due to diffraction, aberration or defocus has little effect, the retinal image at moderate field angles would also be expected to be uniform across its area, except at the edges. Wyszecki and Stiles³⁷ show that the retinal irradiance in a wavelength interval $\delta\lambda$ corresponding to an external stimulus of spectral radiance $L_{e\lambda}(\theta, \phi)$ per unit wavelength interval per cm^2 per unit solid angle of emission, in a direction with respect to the eye given by the angular coordinates (θ, ϕ) is:

$$\frac{L_{e\lambda}(\theta, \phi) \cdot \delta\lambda \cdot p(\theta, \phi) \cdot t(\theta, \phi, \lambda)}{m(\theta, \phi, \lambda)} \quad (4)$$

where $p(\theta, \phi, \lambda)$ cm^2 is the apparent area of the pupil as seen from the direction (θ, ϕ) ; $t(\theta, \phi, \lambda)$ is the fraction of the incident radiant flux transmitted through the eye; and $m(\theta, \phi, \lambda)$ is an areal

magnification factor (cm^2) relating the area of the retinal image to the angular subtense of the stimulus at the eye, which will vary somewhat with the parameters of the individual eye. If required, the pupil area $p(\theta, \phi, \lambda)$ can be modified using Eq. (3) to take account of the Stiles-Crawford effect.

Use of Eq. (4) near to the visual axis is straightforward. In the peripheral field, however, complications arise. With increasing field angle the entrance pupil appears as an ellipse of increasing eccentricity and reduced area: the ratio of the minor diameter to the major diameter falls off somewhat more slowly than the cosine of the field angle.^{86,87} Also, due to the retina lying on the curved surface of the quasi-spherical eyeball, both the distance between the exit pupil of the eye and the retina, and the retinal area corresponding to the image of an object of constant angular subtense diminish with field angle. Remarkably, theoretical calculations⁸⁸⁻⁹¹ show that these pupil and retinal effects tend to compensate one another, so that an extended field of constant luminance (i.e. a *Ganzfeld*) results in a retinal illuminance which is almost constant with peripheral field angle. This theoretical result is broadly confirmed by practical measurements,⁹¹ showing that, from the photometric point of view, the design of the eye as a wide-angle system is remarkably effective.

A useful discussion of the photometric aspects of point and extended sources in relation to the eye is given by Wright.⁹²

1.5 FACTORS AFFECTING IN-FOCUS RETINAL IMAGE QUALITY

The optical quality of the retinal image is degraded by the effects of diffraction, monochromatic and chromatic aberration, and scattering. The image is often further blurred by defocus, due to errors in refraction and accommodation. Under many conditions the latter may be the dominant cause of image degradation. It will, however, be convenient to consider the question of focus in a separate section.

The Aberration-Free (Diffraction-Limited) Eye

In the absence of aberration or variation in transmittance across the pupil, the only factors influencing the retinal image quality in monochromatic light at optimal focus would be the diffraction effects associated with the finite wavelength, λ , of the light and the pupil diameter, d . For such an eye, the point-spread function (PSF) is the well-known Airy diffraction pattern⁹³ whose normalized illuminance distribution, $I(z)$, takes the form:

$$I(z) = \left[\frac{2J_1(z)}{z} \right]^2 \quad (5)$$

where $J_1(z)$ is the Bessel function of the first kind of order 1 of the variable z . In the case of the eye, the dimensionless distance z has the value:

$$z = \frac{d\pi \cdot \sin\gamma}{\lambda} \quad (6)$$

γ is the angular distance from the center of the pattern, measured at the second nodal point, this being equal to the corresponding angular distance in the object space, measured at the first nodal point.^{94,95} The angular resolution θ_{\min} for two neighboring equally luminous incoherent object points, as given by the Rayleigh criterion is then:

$$\theta_{\min} = \frac{1.22 \lambda}{d} \text{ rad} \quad (7)$$

θ_{\min} is about 1 minute of arc when d is 2.3 mm and λ is 555 nm. Evidently the Rayleigh criterion is somewhat arbitrary, since it assumes that the visual system can just detect the 26 percent drop in irradiance between the adjacent image peaks; for small pupils actual visual performance is usually somewhat better than this limit.⁹⁶ The size of the PSF increases on either side of optimal focus.⁹⁷⁻⁹⁹

Images of more complex objects can be obtained by regarding the image as the summation of an appropriate array of PSFs, that is, by convolving the PSF with the object radiance distribution.⁹⁴ The in-focus line-spread function (LSF), the edge image,¹⁰⁰ and the modulation transfer function (MTF) also take standard forms. The phase transfer function (PTF) is always zero because the diffractive blur is rotationally symmetrical: this also makes the LSF, edge image, and MTF independent of orientation. Figure 9 shows the form of the MTF as a function of focus.¹⁰¹ Extensive tables of numerical values are given by Levi.¹⁰² Relative spatial frequencies, R_R , in Fig. 9 have been normalized in terms of the cutoff value beyond which the modulation transfer is always zero. Errors of focus have been expressed in what Levi calls “Rayleigh units,” i.e., the number of quarter wavelengths of wavefront aberration at the edge of the pupil. Note from Fig. 9 that the modulation transfer is most sensitive to defocus at intermediate normalized spatial frequencies ($R_R \approx 0.5$).^{103,104}

To convert the units of relative spatial frequency and Rayleighs to true spatial frequencies R c/deg and dioptric errors of focus ΔF respectively, the following relations may be used:

$$\begin{aligned} R &= \frac{10^6 \times d \times R_R}{\lambda} \text{ c/rad} \\ &= \frac{1.746 \times 10^4 \times d \times R_R}{\lambda} \text{ c/deg} \end{aligned} \quad (8)$$

$$\Delta F = \frac{2 \times 10^{-3} \lambda g}{d^2} \text{ diopters} \quad (9)$$

where the entrance pupil diameter d is in mm, the wavelength λ is in nm, and g is the number of Rayleighs of defocus.

Figure 10 illustrates the variation in these parameters as a function of the ocular pupil diameter, d , for the case where $\lambda = 555$ nm. In such green light and with the typical photopic pupil diameter of 3 mm, the cutoff frequency imposed by diffraction is about 100 c/deg and one Rayleigh of defocus corresponds to about 0.12 D.

Sets of diffraction-limited ocular MTF curves for various specific combinations of pupil diameter, wavelength, and defocus have been illustrated by several authors (e.g., Refs. 105–111).

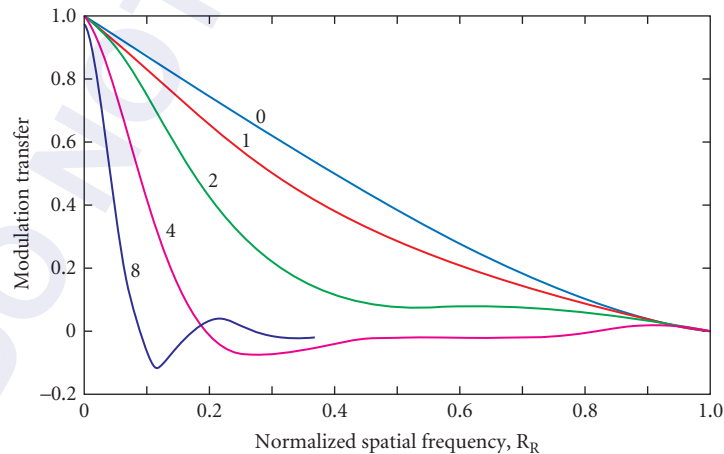


FIGURE 9 Modulation transfer functions for a diffraction-limited optical system with a circular pupil working in monochromatic light. It suffers from the errors of focus indicated. Defocus is expressed in Rayleighs, that is, the number of quarter-wavelengths of defocus wavelength aberration. (Based on Levi.¹⁰²)

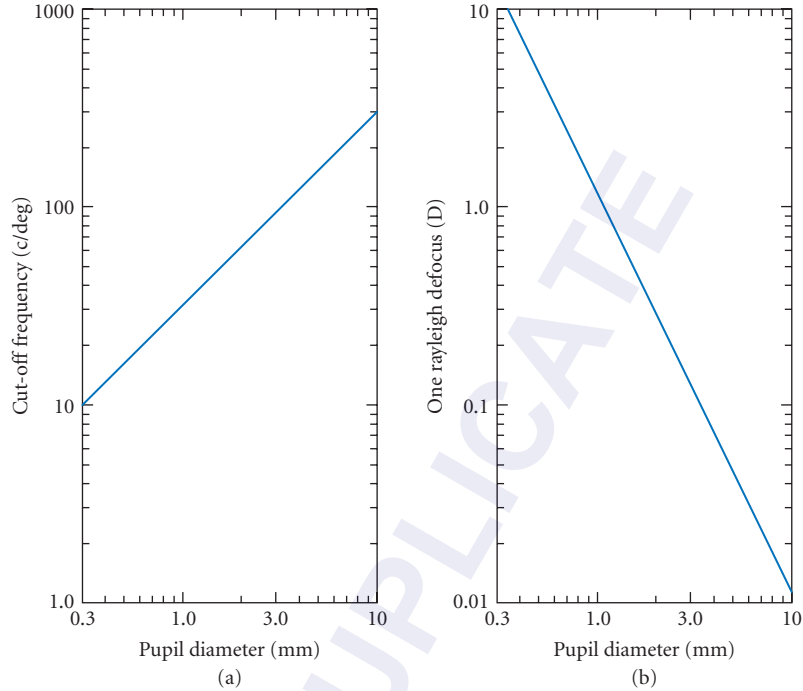


FIGURE 10 Values of (a) cutoff frequency $R_R = 1.0$ and (b) dioptic defocus equivalent to one Rayleigh, for a diffraction-limited eye working in monochromatic light of wavelength 555 nm.

When errors of focus become large, the geometric approximation in which the defocus PSF is a uniform blur circle becomes increasingly valid.^{101,110,112–115} The angular diameter, β of the retinal blur circle for a pupil diameter d mm and error of focus ΔF diopters is

$$\beta = \frac{0.18 \cdot d \cdot \Delta F}{\pi} \text{ deg} \quad (10)$$

The corresponding geometrical optical MTF is

$$T(R) = \frac{2J_1(\pi\beta R)}{\pi\beta R} \quad (11)$$

where R is the spatial frequency (c/deg) and $J_1(\pi\beta R)$ is the Bessel function of the first kind of order 1 of ($\pi\beta R$). Smith¹¹⁴ gives detailed consideration to the range of ocular parameters under which the geometrical optical approximation may reasonably be applied and Chan et al.¹¹⁵ have demonstrated experimentally that Eq. (10) usefully predicts blur circle diameters for pupils between 2 and 6 mm in diameter and defocus between 1 and 12 D.

Monochromatic Ocular Aberrations

In many early studies of ocular aberration, it was common to assume that the eye was symmetrical about a unique optical axis which differed only slightly from the visual axis. Thus it was expected that spherical aberration would dominate on axis, with oblique astigmatism, field curvature, and

coma becoming progressively more important as the field angle increased, although the curved surface of the retina would tend to reduce the impact of field curvature. The assumption was that the brain would adapt to any distortion, so that this aberration was not important. This simple picture has been progressively modified with the realisation that exact symmetry about a unique optical axis rarely occurs and that the fovea does not normally lie where the nominal optical axis intercepts the retina (the difference is usually a few degrees—see e.g., Refs. 116–118 for discussion of possible axes for eye). As a result, the patterns of both on- and off-axis aberrations are more complex than was expected on the basis of early, simple models of the eye.

Recent years have, in fact, seen an enormous expansion in the literature of the aberrations of the eye, fueled by the development of commercial aberrometers capable of measuring the characteristics of the individual eye within a few seconds^{119,120} and the demands of refractive surgery, where it was recognized that although earlier techniques nominally corrected refractive error, they often resulted in poor visual outcomes due to higher than normal levels of residual aberration.¹¹⁹ Currently, the optical defects of the eye are usually described in terms of its wavefront aberration under specified conditions, the overall level of aberration being expressed as the root-mean-square (RMS) wavefront error. Different types of aberration are quantified in terms of the coefficients of the corresponding Zernike polynomials of polar coordinates in the pupil (e.g., Refs. 119, 121–124, see also Chap. 11 by Virendra N. Mahajan in Vol. II as well as Chap. 4 by Virendra N. Mahajan and Chap. 5 by Robert Q. Fugate in Vol. V), defined according to OSA recommendations,^{123,124} although this approach has been criticized as being inappropriate for eyes in which the wavefront aberration shows locally abrupt variation, as in, for example, some postsurgical cases.^{125,126} In the recommended formulation, each Zernike coefficient gives the root-mean-square (RMS) wavefront error (in microns) contributed by the particular Zernike mode: the overall RMS error is given by the square root of the sum of the squares of the individual coefficients. The set of Zernike coefficients thus gives detailed information on the relative and absolute importance of the different aberrational defects of any particular eye for the specified conditions of measurement. In the Zernike description, first-order polynomials simply describe wavefront tilt (i.e. prismatic effects) and have no effect on image quality. Second-order polynomials describe the spherocylindrical errors of focus which can normally be negated by optical corrections, such as spectacles or contact lenses. It is the higher-order (third and greater) polynomials which represent the aberrations. The third-order modes include vertical and horizontal primary coma, and the fourth-order primary spherical aberration. Iskander et al.¹²⁷ have illustrated the effects of some of the individual Zernike aberrations on the retinal images of a selection of objects. The values of the Zernike coefficients for any particular eye will, of course, vary with pupil diameter, accommodation, and field angle.

Aberrations on the Visual Axis Several large-scale studies have addressed the question of the variation of aberrations between individuals.^{128–133} Others have considered changes of aberration with such specific factors as pupil diameter,¹³⁴ age,^{134–138} accommodation,^{139–144} refractive error,^{146,147} and time.^{148–152}

Figure 11 shows recent mean data for the variation in total higher-order, axial, RMS wavefront error with pupil diameter for different age groups.¹³⁴ As would be expected, aberration levels tend to increase with pupil diameter: they also increase with age. The Maréchal criterion¹⁵³ suggests that near diffraction-limited performance will be given if the RMS error is less than $\lambda/14$, corresponding to about 0.04 microns in the green region of the spectrum. It can be seen that this level of aberration is typically present when the pupil diameter is about 3 mm in younger eyes. As illustrated in Fig. 9, such a pupil diameter is generally found under luminance conditions of a few hundred cd/m^2 , corresponding to that occurring on cloudy days. Thus, in most eyes, wavefront aberration is likely to have only a minor impact on vision under daylight conditions.

To give some insight into the image degradation caused by any level of RMS wavefront aberration, we can roughly evaluate its blurring effects by equating them with those of an “equivalent defocus,” that is the spherical error in focus which produces the same magnitude of RMS aberration for the same pupil size. The equivalent defocus is given by:

$$\text{Equivalent defocus (diopters)} = \frac{16.3^{1/2} \cdot \text{RMS error}}{d^2} \quad (12)$$

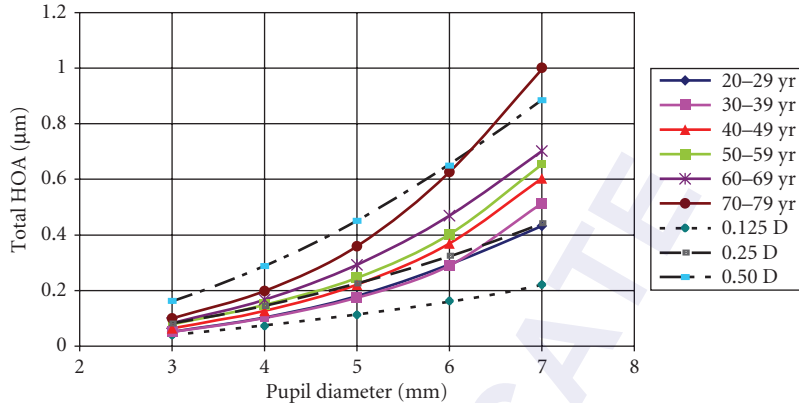


FIGURE 11 Plots of total mean RMS wave aberration as a function of pupil size for different age groups. (After Applegate et al.¹³⁴) The dashed curves show levels of equivalent defocus of 0.125, 0.25, and 0.50 D.

where the RMS aberration is measured in microns and the pupil diameter, d , in mm. As examples, the dashed curves in Fig. 11 indicate equivalent defocus levels of 0.125, 0.25, and 0.50 D. For younger subjects (20–39 years), the mean HOA is always lower than an equivalent defocus level of 0.25 D, except at the largest, 7 mm, pupil diameter. For comparison, the reliability of clinical refractive techniques is around ± 0.3 D.^{154,155} Although the assumption that equal RMS errors produce equal degradation of vision is not completely justified,^{156,157} it is evident that, in most younger, normal eyes, the impact on vision of optical blur due to axial monochromatic aberrations is likely to be modest under most conditions, although this may not be true for a minority of individuals.

When the coefficients of the individual Zernike polynomials are considered for eyes in which the accommodation is relaxed for distance vision, several large-scale studies involving many hundred normal individuals give very similar results.^{128–134} As an example, the study by Applegate and his colleagues¹³⁴ generated mean values for the magnitudes of different types of third- and fourth-order aberration for different pupil sizes and ages (coefficients for still higher-order Zernike modes are usually much smaller). Table 1 gives examples of their values for different age groups. Note that,

TABLE 1 Mean Values of the Coefficients (μm) and Their Standard Deviations for Individual Third- and Fourth-Order Zernike Modes for 3- and 6-mm Pupils and Different Subject Age Groups*

Age (yrs)	Pupil Diameter (mm)	RMS WFE (μm) Trefoil	RMS WFE (μm) Coma	RMS WFE (μm) Tetrafoil	RMS WFE (μm) 2nd Astig.	RMS WFE (μm) Sph. ab.
20–29	3	0.029 ± 0.018	0.028 ± 0.019	0.011 ± 0.010	0.011 ± 0.007	0.013 ± 0.013
30–39	3	0.027 ± 0.017	0.031 ± 0.022	0.010 ± 0.004	0.015 ± 0.008	0.014 ± 0.010
40–49	3	0.038 ± 0.023	0.036 ± 0.020	0.014 ± 0.008	0.014 ± 0.009	0.016 ± 0.011
50–59	3	0.043 ± 0.027	0.048 ± 0.028	0.019 ± 0.016	0.018 ± 0.011	0.014 ± 0.011
60–69	3	0.041 ± 0.021	0.047 ± 0.026	0.023 ± 0.019	0.017 ± 0.011	0.027 ± 0.013
70–79	3	0.059 ± 0.031	0.055 ± 0.026	0.024 ± 0.014	0.020 ± 0.010	0.030 ± 0.022
20–29	6	0.141 ± 0.089	0.137 ± 0.076	0.051 ± 0.025	0.063 ± 0.035	0.132 ± 0.108
30–39	6	0.139 ± 0.089	0.136 ± 0.087	0.056 ± 0.030	0.055 ± 0.027	0.130 ± 0.090
40–49	6	0.187 ± 0.083	0.169 ± 0.089	0.073 ± 0.048	0.071 ± 0.037	0.193 ± 0.110
50–59	6	0.189 ± 0.097	0.198 ± 0.145	0.072 ± 0.051	0.073 ± 0.039	0.197 ± 0.115
60–69	6	0.196 ± 0.115	0.238 ± 0.134	0.088 ± 0.068	0.097 ± 0.070	0.235 ± 0.141
70–79	6	0.292 ± 0.175	0.339 ± 0.170	0.113 ± 0.064	0.093 ± 0.060	0.311 ± 0.153

*The third-order modes are third-order trefoil and coma, the fourth-order are tetrafoil, secondary astigmatism (2nd astig.) and spherical aberration (sph. ab.). The eyes are accommodated for distance vision.

Source: Applegate et al.¹³⁴

where appropriate, the coefficients for similar, but differently oriented, polynomials have been combined. Evidently for smaller, 3-mm pupils, third-order coma and trefoil aberrations tend to dominate over fourth-order aberrations but spherical aberration becomes comparable to coma for the larger 6-mm pupil.

The results of another study¹²⁸ are shown in Fig. 12a, where in this case the second-order coefficients are included. Note that the second-order coefficients are much larger than those of the higher orders implying, not surprisingly, that the optical defects of many eyes are dominated by simple sphero cylindrical refractive errors.

A somewhat different picture emerges if we average the signed coefficients of the higher-order Zernike modes, rather than their absolute values (Fig. 12b). It is striking that the coefficients of most modes now have means close to zero, although individual eyes may have substantial aberration, as is shown by the relatively large standard deviations. A notable exception is the $j = 12$, Z_4^0 spherical aberration

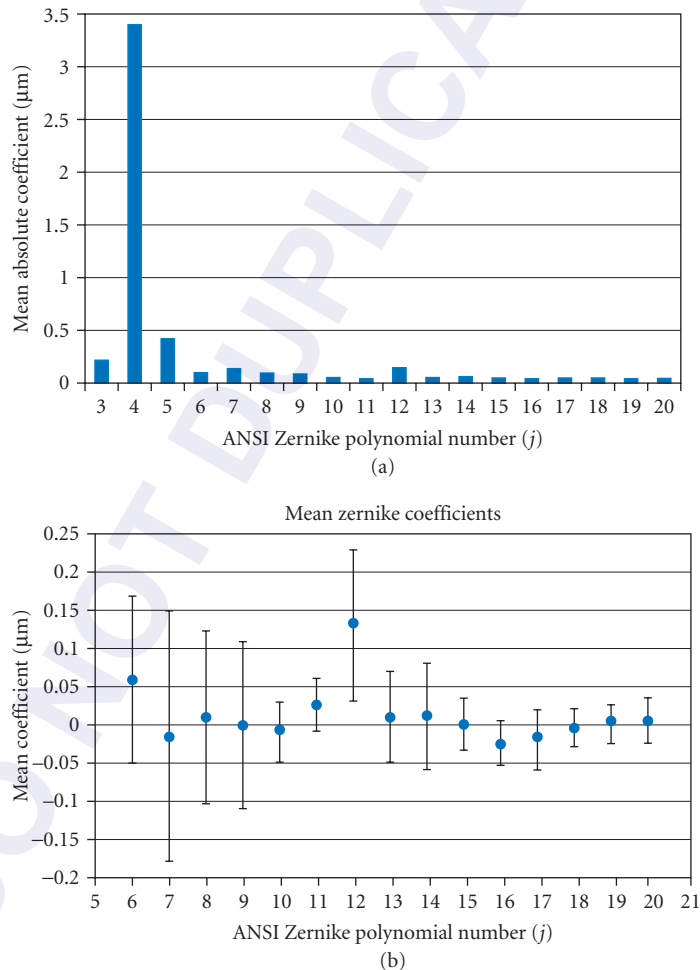


FIGURE 12 Typical data for the wavefront aberration of normal eyes with relaxed accommodation: 109 subjects, 5.7-mm pupil diameter. (a) Means of the absolute values of the coefficients of the Zernike modes from the second to the fifth orders. (b) Means of the signed values of each of the coefficients; among the higher-order coefficients, only that for $j = 12$ (C_4^0 , spherical aberration) has a value which differs significantly from zero. (Based on Porter et al.¹²⁸)

mode, where the mean is positive and differs significantly from zero. Thus the picture that emerges is that most eyes have a central tendency to be free of all higher-order aberration, except for spherical aberration, which shows a significant bias toward slightly positive (under corrected) values. The Zernike coefficients of individual eyes vary randomly about these mean values in a way that presumably depends upon the idiosyncratic surface tilts, decentrations, and other asymmetries of the eye.

When the contributions made by the different optical components of the eye are considered, it appears that, with accommodation relaxed, the overall level of ocular aberration in the young adult is reduced by there being a balance between the contributions of the cornea and the lens. This is particularly the case for spherical aberration and horizontal coma,^{158–162} although the compensation may not be found in young eyes with high levels of total aberration¹⁶³ or in older eyes.¹⁶⁴ The mechanism by which compensation might be achieved has been discussed by Artal and his colleagues.^{2,165} It is interesting to note that there is at least some evidence that there may be additional neural compensation for the aberrations, this being specific to the individual.¹⁶⁶

Since the shape and gradient index characteristics of the lens change with both accommodation and age, this affects the balance between the corneal and internal aberrations. As accommodation increases, spherical aberration tends to change from positive to negative.^{139–145} With age, aberrations with relaxed accommodation at fixed pupil diameter also increase.^{134–138} However, under normal conditions, the pupil diameter at constant light level decreases with age,³⁷ reducing the ocular aberration: image quality therefore remains almost constant with age, although retinal illuminance is lower.¹³⁵ Higher-order aberrations generally show, at most, only a very weak dependence on refractive error,^{147,167} although the balance between the horizontal coma of the cornea and internal optics may be affected.¹⁶⁸

Finally we note that measured higher-order aberrations of any individual eye show small fluctuations over time^{148–152} with frequencies up to at least 20 Hz. Although the causes of these fluctuations remain to be fully elucidated, the lower-frequency components undoubtedly involve such factors as tear film changes and the cardiopulmonary system.^{152,169,170} Lid pressures during such activities as reading may also produce longer-term changes.^{171–173}

Off-Axis Aberrations Off-axis, on average increasing amounts of second-order aberration (defocus and astigmatism) are encountered (Fig. 13). These may show substantial variations with the individual and with the meridian under study, as may also the relationship between the tangential and sagittal image shells and the retinal surface.^{174–187} While there is little systematic change in the mean oblique astigmatism with the axial refraction of the eye, myopes tend to have a relatively hyperopic peripheral mean-sphere refraction, while that in hyperopes tends to be relatively myopic.^{177,178,183,185} There may be small changes in peripheral refraction with accommodation¹⁸¹ and age.^{186,187}

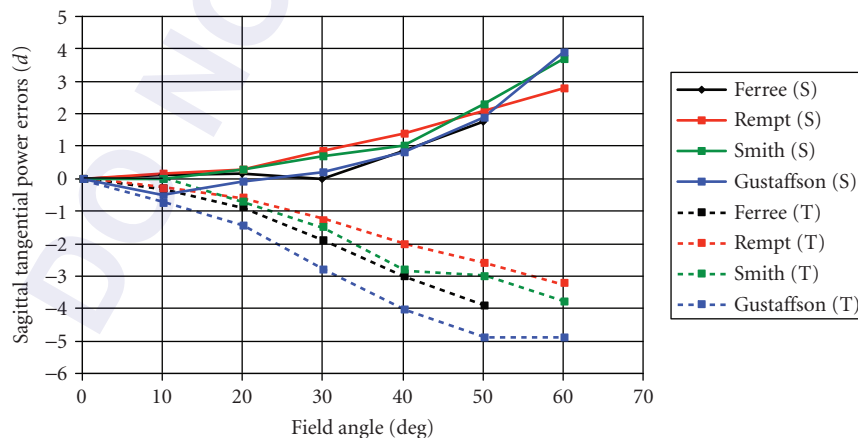


FIGURE 13 Oblique astigmatism in human eyes. T and S refer to the tangential and sagittal image shells, respectively. (After Ferree et al.,¹²⁴ Jenkins,¹⁷⁵ Rempt et al.,¹⁷⁶ Smith et al.,¹⁸¹ Gustafsson et al.¹⁸⁴)

Higher-order wave aberrations have been studied as a function of field angle by several groups.^{188–192} They are generally much less important than the second-order defocus terms but third-order, coma-like terms rise with field angle to be considerably higher than those in axial vision.^{188–192} As on axis, there appears to be a degree of balancing between the aberrations associated with the anterior cornea and those of the lens.¹⁹² One problem in using Zernike polynomials to describe off-axis aberrations is that the associated entrance pupils are elliptical rather than circular as required if the Zernike approach is to be used: scaling methods to overcome this difficulty have been devised.^{193–195}

Chromatic Aberration

Chromatic aberration arises from the dispersive nature of the ocular media, the refractive index, and hence the ocular power, being higher at shorter wavelengths. Constringence values for the ocular media are generally quoted¹⁹⁶ as being around 50, although there is evidence that this may need modification.^{197,198} Atchison and Smith¹⁹⁹ have recently discussed the available data for the different media and recommend the use of Cauchy's equation to fit experimental values in the visible and allow extrapolation into the near infrared. Both longitudinal or axial chromatic aberration (LCA) and transverse or lateral chromatic aberration (TCA) occur (see, e.g., Refs. 200 and 201 for reviews). For LCA, what is normally measured experimentally is not the change in power of the eye across the spectrum but rather the change in its refractive error, or the chromatic difference of refraction. There are only minor differences in the results of different studies of LCA (e.g., Refs. 175, 202–205) and the basic variation in ocular refraction with wavelength, equivalent to about 2 D of LCA across the visible spectrum is well established (Fig. 14). Atchison and Smith¹⁹⁹ suggest that when the chromatic difference data are set to be zero at 578 nm they can be well fitted by the Cauchy equation

$$R_x(\lambda) = 1.60911 - 6.70941 \times \frac{10^5}{\lambda^2} + 5.55334 \times \frac{10^{10}}{\lambda^4} - 5.5998 \times \frac{10^{15}}{\lambda^6} \text{ diopters} \quad (13)$$

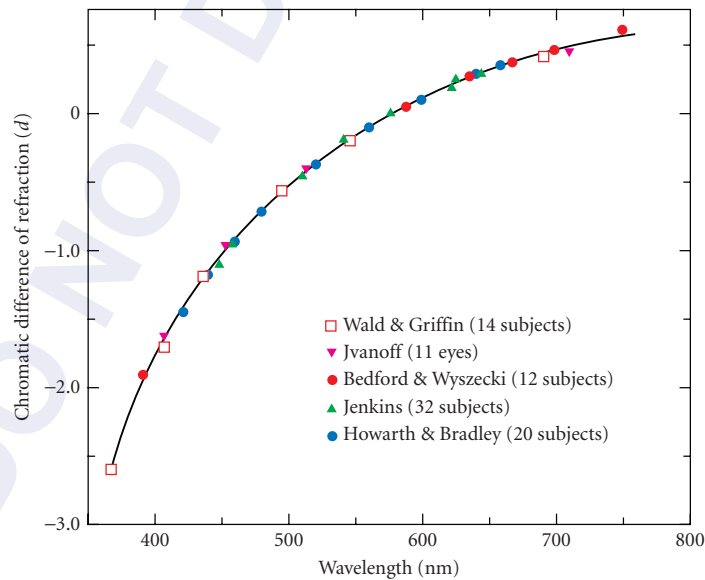


FIGURE 14 Representative sets of average data for the longitudinal chromatic aberration of the eye. The data for all subjects and studies have been displaced in power so that the chromatic aberration is always zero at 578 nm. (Based on Jenkins,¹⁷⁵ Wald and Griffin,²⁰² Bedford and Wyszecki,²⁰⁴ Howarth and Bradley,²⁰⁵)

where the wavelength is in nanometers and the difference in refraction, $R_x(\lambda)$, is in diopters. LCA has little effect on visual acuity for high-contrast objects in white light.²⁰⁶ This appears to be because the spectral weighting introduced by the photopic luminosity curve, which is heavily biased toward the central region of the spectrum, results in the primary effect of the LCA being to degrade modulation transfer at intermediate spatial frequencies rather than at the high spatial frequencies involved in the high-contrast acuity limit.¹⁰⁷ Alternatively, McLellan et al.²⁰⁷ have argued that it is the interaction between the monochromatic and chromatic aberrations of the eye that helps to minimize any change in image quality with wavelength across the spectrum.

TCA results in a wavelength-dependent change in magnification in the retinal image, the red image being larger than the blue. Thus in white light the image of an off-axis point is drawn out into a short radial spectrum whose length in a paraxial model increases with the field angle.²⁰⁸ TCA therefore affects modulation transfer for tangentially oriented grating components of images.^{201,209,210}

For a centered system, TCA would be zero on the optical axis. Since in the eye there is a ~ 5 deg difference (called angle α) in orientation between the visual axis joining the fixation point, nodal points and fovea of the eye, and the approximate optical axis, some foveal TCA might be expected. Remarkably, however, the center of the pupil in most eyes lies almost exactly on the visual axis^{1,200,211} so that actual values of foveal TCA are typically only of the order of 0.7 min arc.^{1,211} Although this value is less than predicted by simple eye models, it is still large enough to cause some orientation-dependent image degradation. This may be substantially increased if an artificial pupil is used which, for any reason, becomes decentered: such a situation may arise when using visual instrumentation having exit pupils which are much smaller than the entrance pupil of the eye.

With binocular viewing, the TCA associated with small decentrations of the natural pupils and of the foveas from the optical axis leads to the phenomenon of *chromostereopsis*, whereby objects of different colors placed at the same physical distance may appear to the observer to be at different distances.^{212–218} The exact effect varies with the individual and, when artificial pupils are used, with the separation of the pupils. It is thus of some practical significance in relation to the design of instruments, such as binocular microscopes, in which interpupillary distance settings may not always be optimal for the observer.²¹⁶

In the periphery, TCA may play a significant role in limiting the detection of tangential as opposed to radial gratings.^{210,219} There is as yet no consensus as to its magnitude, although some measurements have been made.²²⁰

Intraocular Scattered Light and Lenticular Fluorescence

A variety of regular and irregular small-scale inhomogeneities exist within the optical media of the eye. These may serve to scatter light during its passage between the anterior cornea and the retinal receptors. Further stray light may arise for reflections at the various optical surfaces and the retina itself, and some light may also penetrate the nominally opaque iris and sclera to reach the internal eye (*diaphany*), particularly in eyes with low pigmentation, as in albinos. The main effect of such light is to reduce the contrast in the retinal image.

Quantitative studies of the effects of stray light on vision were pioneered by Holladay²²¹ and Stiles,²²² who expressed its impact in terms of an equivalent veiling luminance, L_{eq} $\text{cd} \cdot \text{m}^{-2}$, that would produce the same masking effect as a glare source giving an illuminance E lux at the eye, as a function of the angular distance, ω deg, between the glare source and the fixation point. Vos et al.²²³ have summarized more recent work by the approximate relationship:

$$L_{\text{eq}} = \frac{29E}{(\omega + 0.13)^{2.8}} \quad \text{where} \quad 0.15 \text{ deg} < \omega < 8 \text{ deg} \quad (14)$$

This expression relates to young adult eyes. Scattering increases throughout life by a factor of at least 2 to 3 times^{224–227} and glare formulas can be modified to take account of this (e.g., Ref. 228). Roughly a quarter of the stray light comes from the cornea^{229,230} and a further quarter from the retinal reflections.^{231,232} The rest comes almost entirely from the lens,²³³ there being little contribution from the aqueous or vitreous humors in normal healthy eyes.

Ohzu and Enoch²³⁴ attempted to measure a retinal MTF which included the effects of forward scatter, by focusing grating images onto the anterior surface of an excised retina *in vitro* and measuring the image transfer to the far side. They argued that the receptor outer segments effectively act as a fiber optics bundle and that transmission through this bundle produces image degradation which supplements to the main optical elements of the eye. More recent psychophysical measurements^{235,236} suggest, however, that forward scatter in the inner retina is negligible, implying that postmortem changes in the retina may have degraded Ohzu and Enoch's MTFs.

In addition to the general effects described above, a variety of more regularly organized, wavelength-dependent, scattering effects in the form of annular haloes or "star" patterns may occur (see e.g., Refs. 237 and 238 for reviews). These are most easily observed with point sources in an otherwise dark field, particularly if the pupil is large. Some of these result from diffraction due to ocular structures with quasi-regular spacing, for example, the corneal epithelial cells or lens fibers,^{239,240} others relate to small-scale refractive irregularities or to higher-order aberrations.²⁴¹

Forward scattering of light affects the precision of aberrometry,^{137,241} while back scatter is of importance in allowing anterior ocular structures to be viewed by clinical examination techniques such as slit-lamp biomicroscopy and Scheimpflug photography.²⁴²

Stray light may also arise as a result of fluorescence in the crystalline lens. This increases with age and with the occurrence of cataract,²⁴³ largely through the progressive accumulation of fluorogens. Under some circumstances, the emitted fluorescent light can cause a slight reduction in low-contrast acuity in older individuals²⁴⁴ but in younger adults the effects are probably of little practical significance,²⁴⁵ since the cornea absorbs most of the potentially activating short wavelength light.

1.6 FINAL RETINAL IMAGE QUALITY

Experimental estimates of the final quality of the retinal image can be made in three main ways: by calculation from wavefront aberration or similar data, by psychophysical methods, and by direct measurement of the light distribution on the retina using a double-pass ophthalmoscopic technique. Although each method has its limitations, the various methods yield compatible results in the same subjects and collectively produce a reasonably consistent picture of the changes in retinal image quality with pupil diameter, retinal location, and age.

Image Quality on the Visual Axis

Calculation from Aberration Data The optical transfer function (OTF) can be calculated by auto-correlation of the complex pupil function with its complex conjugate, using methods originally devised by Hopkins.²⁴⁶ The pupil function gives the variation in amplitude and phase across the exit pupil of the system. The phase at each point can be deduced from the corresponding value of the wavefront aberration (each wavelength of aberration corresponds to 2π radians of phase). It is often assumed that the amplitude across the pupil is uniform but if imagery under photopic conditions is being considered, it may be more correct to take account of the Stiles-Crawford effect (SCE I) by including appropriate amplitude apodization, ideally on an individual basis.^{209,247,248} This suggestion is supported by some experimental evidence.^{249,250} The point- and line-spread functions can also be directly calculated from the wavefront aberration (see Chap. 4 by Glenn D. Boreman in Vol. I and Chap. 4 by Virendra N. Mahajan and Chap. 5 by Robert Q. Fugate in Vol. V) and appropriate software is often included with current commercial aberrometers.

The attractive feature of this approach is that it can allow the OTF (i.e., both the modulation and phase transfer functions) to be calculated for any orientation. On the other hand, it fails to include the effects of any scattered light and hence may give too optimistic a view of the final retinal image quality, particularly in older eyes in which scattering is high: high levels of intraocular scatter may have the additional effect of reducing the reliability and validity of the aberrometer estimates of wavefront aberration, the exact effects depending upon the design of the particular aberrometer.²⁵¹

Van Meeteren²⁰⁹ argued that it was appropriate to multiply the aberration-derived MTFs by the MTF derived by Ohzu and Enoch²³⁴ for image transfer through the retina. When this is done, the results agree quite well with those found by the double-pass ophthalmoscope technique;²⁵² although, as noted earlier, the Ohzu and Enoch MTF may overestimate the effects of retinal degradation. A further problem with MTFs derived from wavefront measurements is that most aberrosopes estimate the form of the wavefront from measurements made at a limited number of points across the pupil.^{253,254} Variations in aberration on a small spatial scale may therefore remain undetected, with consequent uncertainties in the derived OTFs; this problem was probably greatest with the early designs of aberroscope, such as the crossed cylinder device.²⁵²

Psychophysical Comparison Method This method depends upon the comparison of modulation (contrast) thresholds for a normally viewed series of sinusoidal gratings of differing spatial frequencies with those for similar gratings which are produced directly on the retina by interference techniques.^{256,257}

Suppose an observer directly views a sinusoidal grating of spatial frequency R . Then if the grating has modulation $M_0(R)$, the modulation of the retinal image will be $M_0(R) \cdot T(R)$, where $T(R)$ is the modulation transfer of the eye at this spatial frequency, under the wavelength and pupil diameter conditions in use. If now the modulation of the grating is steadily reduced until it appears to be just at threshold, the threshold modulation $M_{IT}(R)$ on the retina will be given by

$$M_{IT}(R) = M_{OT}(R) \cdot T(R) \quad (15)$$

where $M_{OT}(R)$ is the measured modulation of the external grating at threshold. The reciprocal of $M_{OT}(R)$ is the corresponding conventional contrast sensitivity and its measurement as a function of R corresponds to the procedure used to establish the contrast sensitivity function.

It is clear that $M_{IT}(R)$ corresponds to the threshold for the retina/brain portion of the visual system. If its value can be independently established, it will be possible to determine $T(R)$. $M_{IT}(R)$, in fact, can be measured by bypassing the dioptics of the eye and their aberrations and forming a system of interference fringes directly on the retina. This procedure was originally suggested by Le Grand^{256,257} and has since been progressively improved^{258–264} (see Ref. 257 for review). Two mutually coherent point sources are produced close to the nodal points of the eye and the two resultant divergent beams overlap on the retina to generate a system of Young's fringes, whose angular separation, γ radians, is given by $\gamma = \lambda/a$, where λ is the wavelength and a is the source separation, both measured in air. If the sources have equal intensity, the fringes will nominally be of unit modulation. Fringe modulation can be controlled by varying the relative intensities of the two sources, by adding a uniform background, or by modulating the two sources with a temporal square-wave and introducing a phase difference between the two modulations. The contrast threshold $M_{IT}(R)$ for the retina brain can then be measured as a function of R , allowing the modulation transfer for the ocular dioptics, $T(R)$, to be deduced from the relationship:

$$T(R) = \frac{M_{IT}(R)}{M_{OT}(R)} \quad (16)$$

There are some problems with this approach. Both sets of thresholds are affected by stray light, but in different ways. The determination of the external modulation threshold involves light entering the full pupil whereas for the interferometric measurements only two small regions of the pupil are used. There may also be problems in maintaining the same threshold criterion for the two types of grating, particularly when they may differ in color, field size, and possibly speckle characteristics. Lastly, although the method can in principle give ocular MTFs for any orientation, it yields no phase information and hence the PTF cannot be determined.

Some other psychophysical methods have been suggested²⁶⁵ but as yet they have not been widely employed.

Ophthalmoscopic (Double-Pass) Methods When the image of an object is thrown on the retina, some of the light will be reflected back out of the eye and can be collected by an appropriate observing system

to form an external image. If, for example, the object is a narrow line, the external image will be the LSF for the double passage of the eye. It is usual to make the assumption that the retina acts as a diffuse reflector,²⁶⁶ coherence of the image light being lost in the reflection. In this case, earlier workers assumed that the image suffered two identical stages of image degradation, so that the MTF deduced from the Fourier transform of the external LSF was the square of the single-pass MTF. Flamant's pioneering study²⁶⁷ with this method used photographic recording, but later workers have all used electronic imaging methods, initially slit-scanning arrangements with photomultipliers to record LSFs, and latterly low-noise CCD cameras which allow PSFs to be recorded.^{266,268–276} An important advance was the realization that in the simple form of the method as employed in earlier studies, when the same pupil acted as aperture stop for both the entering and exiting light paths, information on odd-order aberrations and on transverse chromatic aberration was lost.²⁷⁴ While the estimates of MTF were potentially correct, the PTF could not be measured. This problem can be overcome for monochromatic aberrations by arranging the entering and exiting beams so that the entrance pupil is smaller than the exit pupil.^{275,276} If the entrance pupil is small enough for the initial image to be effectively diffraction limited, the true single-pass OTF (i.e., both the MTF and the PTF) can be deduced from the double-pass PTF, at least up to the cutoff frequency imposed by the small entrance pupil. Some theoretical aspects of this problem have been discussed by Diaz-Santana and Dainty.²⁷⁷

The double-pass method has been used to explore the extent to which poorer retinal image quality contributes to the deterioration in visual performance that is observed in older eyes,²⁷⁸ and to demonstrate the changes in retinal image quality with accommodation that are caused by aberrational change in the crystalline lens.²⁷⁹ An adaptation allows the basic method to be used to determine an "index of diffusion" designed to characterize the optical deficit in eyes with age- and disease-related abnormalities of the anterior segment, using encircled energy measurements of the double-pass PSF.²⁸⁰

In all variations of the double-pass method, one problem is that light levels in the outer parts of any spread function are low, leading to possible truncation errors and to overestimation of the MTF.²⁸¹ Vos et al.²²³ attempted to overcome the truncation problem by combining the ophthalmoscopic estimates of the PSF with measurements of wider-angle entoptic stray light, to produce a realistic estimate of the full light profiles in the foveal white-light PSF. A vexing question which has yet to be fully answered is the identity of the layer or layers at which the retinal reflection occurs: it seems likely that this is wavelength dependent. If more than one layer is involved, the estimated MTF will be somewhat too low. However there is evidence that any effect of retinal thickness on the estimated MTF is small²⁸² and that scattered light from the choroid and deeper retina is guided through the receptors on its return through the pupil.²⁸³

Comparison between Methods Only a few direct comparisons of MTF measurements have been made by different techniques on the same eyes. Campbell and Gubisch²⁶⁶ found that their double-pass MTFs were lower than those determined by the interferometric psychophysical method.²⁶⁰ A similar result was found by Williams et al.,²⁸⁴ who noted that agreement between the two techniques was better if green rather than red light was used for the double-pass measurements, presumably as a result of reduced retinal and choroidal scatter.²⁸⁵ Although MTFs derived from early aberrometers, which only sampled the pupil at a small number of points, tended to be markedly higher than those derived by other methods, if green light is used with young eyes the three basic methods appear to yield very similar MTFs: increased entoptic scatter in older eyes may cause larger differences. Liang and Williams²⁸⁶ give comparative MTF results obtained by the three techniques with three subjects with 3-mm pupils. The greatest discrepancies appear to be at intermediate spatial frequencies, where values of modulation transfer increase in the order double-pass, interferometric, and wave aberration derived.

Summary of Observed Optical Performance When the eye is corrected for any spherocylindrical refractive error, all investigators agree that, near the visual axis, the eye's performance in monochromatic light is reasonably close to the limit set by diffraction for pupil diameters up to 2 mm. As pupil size is increased further, aberration starts to play a more important role. The increasing impact of aberration is illustrated by a consideration of the changes in the Strehl intensity ratio,²⁸⁷ the ratio of the maximum irradiance in the PSF to that which would be found in a truly diffraction-limited

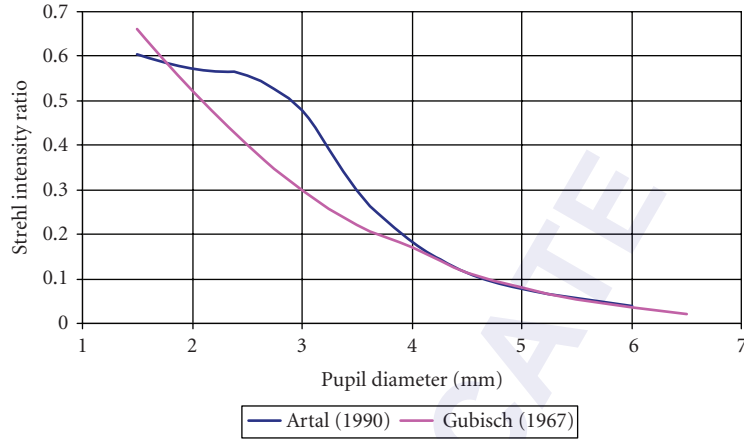


FIGURE 15 Changes in Strehl ration with pupil diameter. (Based on data from Artal²⁸⁸ and Gubisch.²⁸⁸)

system. Typical published values^{89,288} are shown in Fig. 15. Although differences between individuals may be expected, the ratio falls steadily as the pupil diameter is increased, indicating the increasing impact of aberrations. For the smallest natural pupil diameters, the value approaches, but does not quite reach, the figure of 0.8 which is usually accepted as the minimum required for an optical system to closely approximate to being diffraction limited:²⁸⁷ direct measures of the MTF for a 1.5-mm pupil in comparison with the diffraction-limited case support this finding.²⁷⁶

The changing balance between diffractive and aberrational effects results in optimal overall performance usually being achieved with pupil diameters of about 2.5 to 3 mm,^{266,286,288,289} corresponding to the diameters of natural pupils under bright, photopic conditions. For still larger pupils, the degrading effects of aberration dominate and modulation transfer falls. Examples²⁸⁶ of typical estimates of MTF, in this case based on wavefront measurements for different pupil diameters in young eyes, are shown in Fig. 16. The MTF at constant pupil diameter tends to deteriorate with

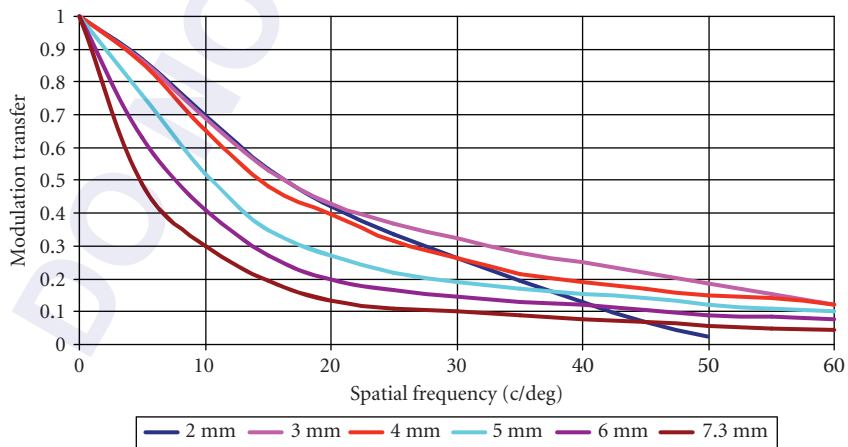


FIGURE 16 Examples of typical foveal MTFs in young, adult eyes for the pupil diameters indicated. (After Liang and Williams.²⁸⁶)

age, due to both increased aberration and scatter.²⁷⁸ It should be noted that, although transmission through the eye undoubtedly changes the polarization state of the incident light, the retinal image quality is apparently nearly independent of the initial state of polarization. However, quality as recorded by double-pass measurements may vary if polarizing elements are included in both the entering and exiting light paths.²⁹⁰

Effects of Aberration Correction on Visual Performance Early attempts²⁹¹ to improve axial visual performance by correction of monochromatic aberrations were failures, largely because it was assumed that spherical aberration was always dominant and the same in all eyes. Correction of longitudinal chromatic aberration also brought no improvement in acuity²⁹² although modest improvements at intermediate spatial frequencies could be demonstrated.⁹⁹ The advent of aberrometers capable of rapidly measuring the aberrations of an individual eye has refocused attention on the possibility that customized aberration correction might yield significantly enhanced vision (e.g., Refs. 293–298).

Several potential methods for correcting the monochromatic aberrations have been suggested. The most effective of these is adaptive optics.^{299,300} This is impractical for everyday use but is being vigorously explored as a way of obtaining improved images of the retina for clinical diagnostic or research purposes (e.g., Ref. 301), since near diffraction-limited performance can be achieved through the full, 7 to 8 mm, drug-dilated eye pupil. For normal purposes, three methods have been suggested. In the first, spatially modulated excimer laser ablation of the cornea is used to correct the measured aberrations of the individual eye, as well as any second-order refractive error.^{293–298} While such “wavefront-guided” ablation has been effective in reducing ocular aberrations consequent upon laser refractive surgery, it has yet to succeed in producing eyes which are completely aberration-free, largely because of factors such as subject-dependent variations in healing which influence the final level of aberration. The second possible method is the wearing of customized contact lenses, where the local thicknesses of the lens are manipulated to yield variations in path length which compensate for the ocular wavefront aberration.^{302,303} The problem with this approach is that the lens must be stable against both rotation and decentration on the eye if aberration correction is to be maintained.^{304–308} The limited control of lens movement occurring in practice means that adequate correction is difficult to achieve in normal eyes. Improved performance may, however, be given in clinically abnormal eyes with high levels of aberration, such as those of keratoconics.³⁰⁷ The last suggested method, which has yet to be successfully demonstrated, is to incorporate the aberration correction in an intraocular lens: this has the theoretical advantage that the position of such a lens and its correction should be stable in the optical path.

Some of the potential benefits of aberration correction have been demonstrated by Yoon and Williams.³⁰⁰ Figure 17 shows some of their results for eyes under cycloplegia. The contrast sensitivity function for either a 3- or a 6-mm pupil was measured under four conditions: white light with only spherocylindrical (second-order) refractive error corrected; with chromatic (but not monochromatic) aberration additionally removed by viewing through a narrow-band green filter; with monochromatic (but not chromatic) aberrations corrected using adaptive optics; with both monochromatic and chromatic aberrations corrected by using adaptive optics with a green filter. The retinal illuminances under the various conditions were kept constant with neutral density filters at 14.3 trolands for the 3-mm pupil and 57 trolands for the 6-mm pupil: these correspond to the lower end of the photopic range (a natural 6-mm pupil diameter is reasonably typical for this level of retinal illuminance). Yoon and Williams³⁰⁰ express their results in terms of the “visual benefit” at each spatial frequency, that is, the ratio of the contrast sensitivity under a particular condition to that achieved with white-light gratings and just the spherocylindrical refractive error corrected.

It can be seen that useful performance gains are given if both monochromatic and chromatic aberration can be corrected, particularly for the larger pupil. The benefits are less impressive (<2 at all spatial frequencies) if only monochromatic aberration is corrected, as would in practice be the case with corneal ablation, contact lens, or other corrections. It has been argued that these gains may be still smaller under real-life conditions, due to such factors as inaccurate accommodation, aberrational changes, and, at the larger natural pupil diameters occurring under mesopic and scotopic conditions, the limits to performance set by the neural parts of the visual system.^{309,310} Although in principle

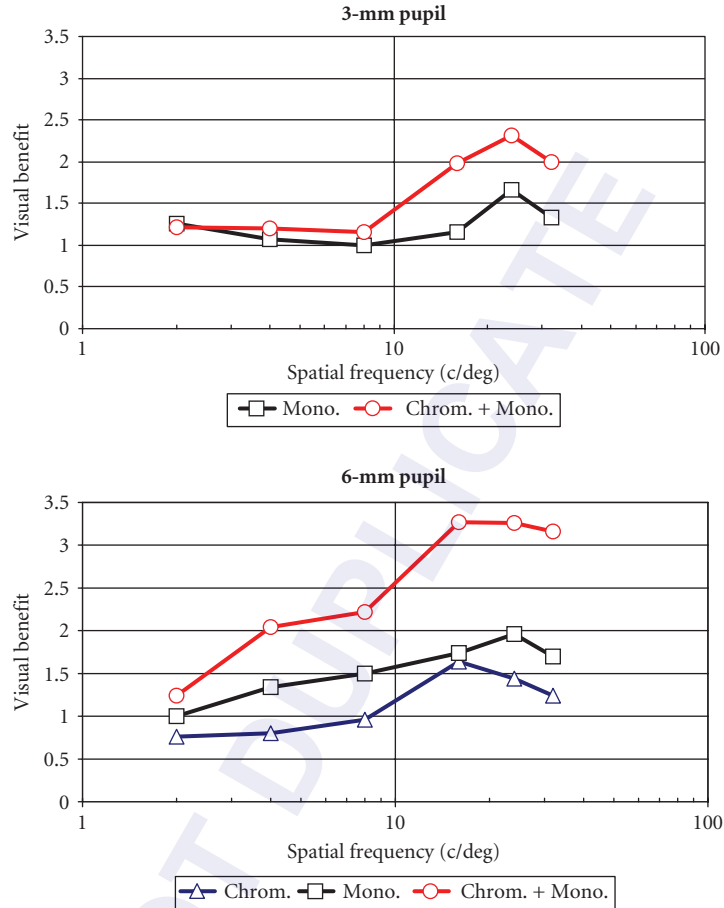


FIGURE 17 Visual benefit in contrast sensitivity for a subject with (a) a 3-mm pupil and (b) a 6-mm pupil. In all cases the eye has an optimal spherocylindrical refractive correction. The benefit is shown for (i) correction of chromatic aberration only (green light), (ii) correction of monochromatic aberration only (adaptive optics, white light), and (iii) correction of both chromatic and monochromatic aberration (adaptive optics, green light). (Based on Yoon and Williams.³⁰⁰)

chromatic aberration of the eye can be corrected,^{204,291,311–315} the required multielement lens systems are relatively bulky and decentration effects may mean that new problems arise in controlling transverse chromatic aberration.³¹⁶

Off-Axis Image Quality

Fewer measurements have been made of off-axis image quality. Psychophysical measurements are difficult to carry out and aberration data, although available, do not yet appear to have been used to calculate MTFs. Most of the available measurements have therefore been made by the double-pass technique.^{272, 317–319} As noted earlier, astigmatism usually increases with field angle, so that the PSFs

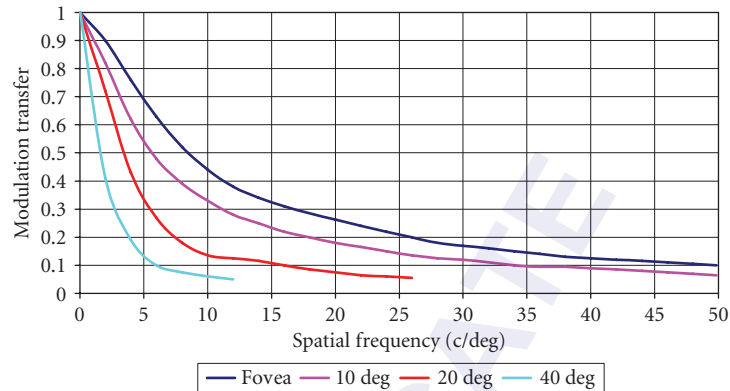


FIGURE 18 MTFs at different field angles for a 3-mm pupil. The position of focus is that corresponding to the circle of least confusion, with no correction for oblique astigmatism. (Based on Williams et al.³¹⁹)

normally lack rotational symmetry and different results for MTF and PTF are obtained at different orientations and levels of focus. Some representative data are shown in Fig. 18. The results shown are for the case when the circle of least confusion lies on the retina. Note that the MTF falls with field angle and that, for objects with strongly oriented structure, image quality will be improved if the appropriate focal line can be brought onto the retina, or if the peripheral astigmatism is corrected. Williams et al.³¹⁹ have suggested that the falloff in optical performance with field angle is advantageous in helping to prevent aliasing in peripheral vision.

Off-axis MTFs for different pupil diameters and field angles have been approximated as the sum of two exponential functions^{318,319} or as a single exponential.³²⁰

Retinal Image Quality with Visual Instruments

When an object (e.g., a display screen) is observed with the naked eye, with no intervening optics, the spatial frequency spectrum of the retinal image is simply the spectrum of the object multiplied by the OTF of the eye under the observing conditions in use. The situation is, however, more complex when instruments such as microscopes, telescopes, or binoculars are used. In general, it is incorrect to assume that the OTF of the instrument-eye combination is simply the product of the individual OTFs. Instead, for any field position the aberrated wavefront from the instrument is modified by the wavefront aberrations of the eye, the two wavefronts being added together point by point across the common pupil. Under favorable circumstances where the component aberrations are of opposite sign, this *coherent coupling*^{321–324} can obviously result in the summed aberration being smaller than the individual wavefront aberrations, and hence in performance being superior to that which would be expected on the basis of the product of the individual OTFs.

Two factors are of particular importance in this coupling: first, the size and positional relationship between the instrument's exit pupil and the ocular entrance pupil and, second, the state of focus of the eye. It is clear that if the entrance pupil of the eye is smaller than the exit pupil of the instrument (the two pupils coinciding) the eye pupil will constitute the aperture stop of the combined system. Movement of the smaller eye pupil within the instrument's exit pupil will sample different areas of the instrumental wavefront aberration and may therefore influence retinal image quality. Conversely, if the exit pupil is smaller than the eye pupil, the instrument's pupil provides the aperture stop and controls the aberrations. For example, with an exit pupil of 2-mm diameter or less centered on the eye pupil, the combined performance would be governed almost entirely by the instrumental aberrations,

since the eye aberrations would normally be small enough to make a negligible contribution to the combined wavefront aberration. Any focus error in the instrument or accommodation error in the observer will add an additional defocus term to the wavefront aberration. This may again be beneficial under some circumstances: for example, instrumental field curvature may be compensated for by changes in ocular accommodation (Refs. 325 and 326, and see later under “Vergence Input”).

Considerations of this type mean that the MTF of any instrument alone may not be a very good guide to the visual performance that can be achieved with it.^{327,328} One possible solution is to test the instrument in combination with a model eye which simulates as accurately as possible the aberrations and focus of the real eye.³²⁹ Burton and Haig³³⁰ explored the tolerance of the visual system to different levels and types of wavefront aberration (see also Ref. 331) and concluded that instrumental criteria based on the Strehl ratio may ultimately be more useful than MTF tests.³²⁸ This suggestion has, however, been disputed by Mouroulis and Zhang³³² who feel that other criteria are more informative.

Chromatic aberrations are also of importance. The role of longitudinal chromatic aberration has been considered by Mouroulis and Woo,^{332,323} who find that under some circumstances quite substantial amounts of instrumental LCA can be tolerated. Further studies support a tolerance of the order of 2.5 min arc for transverse chromatic aberration, with TCA having greater effects on contrast sensitivity than high-contrast resolution.³³⁴ The general question of the design and testing of visual instruments is well reviewed in Mouroulis.³³⁵

1.7 DEPTH-OF-FOCUS AND ACCOMMODATION

In the foregoing it has been tacitly assumed that the eye is always optimally focused. In practice, since clear vision is required over a range of distances, exact focus may not always be achieved. It is of interest that, even if objects are at a fixed distance, small changes in focus may be required to optimize modulation transfer at different spatial frequencies in the presence of some types of aberration (e.g., spherical aberration^{334–339}). With the increasing availability of aberrometers, considerable efforts are being devoted to establishing robust criteria for “best focus” (i.e., ocular refraction) from wavefront data.^{340–342}

Ocular Depth-of-Focus

As in any optical system, very small changes in ocular focus have little effect on the retinal image but image quality deteriorates progressively as the error of focus increases (e.g., Fig. 9). For many practical purposes we would like to know how large the dioptric error of focus can be before image quality becomes unsatisfactory. However, this immediately raises the question of what we mean by “unsatisfactory.” Atchison et al. have, for example, defined the concepts of “noticeable,” “troublesome,” and “objectionable” blur.³⁴³ Noticeable blur is the defocus level at which blur of a set of letters first becomes detectable, troublesome blur is where the lack of clarity in the letters starts to be irritating, although the letters may still be readable, and objectionable blur is the level of blur which is unacceptable: the three dioptric limits of blur were found to be in the ratio of about 1.0:1.7:2.3, respectively. Ciuffreda et al. find very similar results when using minor variants of the same criteria.³⁴⁴ In fact, as will be discussed below, values of depth-of-focus are strongly dependent upon the methodology and conditions used (see Ref. 345 for a recent review).

In the geometrical optical approximation, for an aberration-free eye the angular diameter of the retinal blur circle β degrees increases linearly with the error of focus ΔF diopters and pupil diameter d mm [Eq. (10)]. Thus the limits of the depth-of-focus correspond to the blur circle diameter reaching some assumed tolerable limit β_{tol} degrees: the corresponding value of ΔF_{tol} is then determined from [Eq. (10)] to yield for the total geometrical depth-of-focus DOF_{go}

$$\text{DOF}_{\text{go}} = 2\Delta F_{\text{tol}} = \frac{(34.9\beta_{\text{tol}})}{d} \text{ diopters} \quad (17)$$

The exact value obtained is, then, dependent on the assumed value of β_{tol} . If, say, a value of 2 min arc is taken, and the pupil diameter is 3 mm, the geometric DOF is about 0.4 D. Note that the geometric DOF is inversely dependent on the pupil diameter.

When the effects of physical optics are considered, for a diffraction-limited system it is conventional to use the Rayleigh criterion and to say that the limits of depth-of-focus are set by the requirement that the optical path difference between light from the center and edge of the pupil should not exceed a quarter-wavelength [one Rayleigh unit of defocus, see Eq. (9)]. For the eye, this implies that the total physical optical depth-of-focus (DOF_{po}) should correspond to 2 Rayleigh units, that is,

$$\text{DOF}_{\text{po}} = \frac{4 \cdot 10^{-3} \lambda}{d^2} \text{ diopters} \quad (18)$$

where λ is in nm and d is in mm. Optimal focus will lie midway through this total depth-of-focus. Note that, unlike the geometrical approximation, [Eq. (18)] predicts that depth-of-focus will be inversely proportional to the square of the pupil diameter.

In reality, the ocular depth-of-focus depends on a variety of additional factors. From the purely optical point of view, even if the eye is diffraction-limited the rate of loss in modulation transfer with defocus is spatial-frequency-dependent (Fig. 9), so that the detectable error of focus is likely to depend upon the spatial-frequency content of the object under observation, as well as the pupil diameter. Low spatial frequencies are relatively insensitive to focus change.^{94-96,346} Both monochromatic and chromatic aberrations will further modify the through-focus characteristics of the retinal image focus and, in general, will tend to increase the depth-of-focus.

Equally importantly, the perceptible focus changes will depend upon the neural characteristics of the visual system.³⁴⁷ Under many conditions, the limited capabilities of the retina/brain system will mean that defocus tolerance may become larger than that expected on purely optical grounds. With larger pupils and photopic conditions, the Stiles-Crawford effect may play a role in reducing the effective pupil diameter and increasing the depth-of-focus (see “The Stiles-Crawford effect”). At low luminances, since only low spatial frequencies can be perceived³⁴⁸ there is an increased tolerance to defocus blur.^{349,350} a similarly increased tolerance is found in low-vision patients at photopic levels.⁹⁶

Figure 19 shows a selection of experimental depth-of-focus data obtained by various techniques.^{349,351-354} Although the exact results depend upon the methods and criteria used, it is clear that for larger pupils the visual depth-of-focus substantially exceeds the purely optical Rayleigh-limit predictions. Of more importance from the practical viewpoint, errors of focus in excess of about 0.2 to 0.5 D are likely to lead to perceptible image blur under photopic conditions with pupils of diameter 3 to 4 mm. Thus the eye must either change its focus (accommodate) for the range of distances that are of interest in everyday life (e.g., from infinity down to 0.2 m or less, corresponding to vergences from 0 to 5 D), or use some form of optical aid such as reading glasses.

It is of interest that Goss and Grosvenor³⁵⁵ concluded from their review of the available clinical literature that conventional refraction is repeatable to within 0.25 D in approximately 75 percent of cases and to within 0.50 D in 95 percent of cases: errors in prescription of 0.25 D have been shown to produce dissatisfaction in many patients and a significant loss in both acuity and contrast sensitivity.³⁵⁶

The Accommodation Response

As with any focusing system, several aspects of accommodation are of interest: its range or amplitude, its speed, its stability, and its time-averaged steady-state characteristics. All of these are age-dependent and in most of what follows the behavior of young, adult (age around 15 to 35 years) eyes is described. Discussion of the response is complicated by the fact that, under normal binocular conditions of observation, accommodation (i.e., the focusing system) is intimately linked with the vergence system which ensures that the eyes converge appropriately to bring the images of any object of regard onto the foveas of the two eyes (see “Vergence Input” and “Movements of the Eyes”). Due to this linkage, accommodation can drive convergence and vice versa. The pupil usually contracts

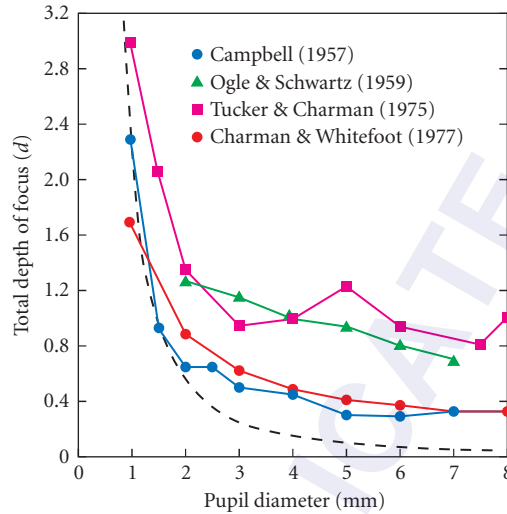


FIGURE 19 Examples of experimental measurements of photopic, total monocular depth-of-focus as a function of pupil diameter (optimal focus lies midway through the total depth-of-focus). The criteria for determining the depths-of-focus are Campbell³⁴⁹—just-perceptible blur for a small disc, one subject, white light; Ogle and Schwartz³⁵¹—50-percent probability of resolving a 20/25 checkerboard, mean of three subjects, white light; Tucker and Charman³⁵²—80-percent probability of achieving 90 percent of the optimal Snellen acuity, mean of two subjects, white light; Charman and Whitefoot³⁵³—detectable movement of laser speckles, mean of six subjects, 633 nm. The dashed line gives the depth-of-focus based on Rayleigh’s quarter-wavelength criterion for an aberration-free eye in 555-nm light.

during near vision (often called accommodative miosis), although it appears that such contraction is not directly driven by accommodation but is a comovement.^{357–362} The three functions (accommodation, vergence, pupil) are sometimes known as the *near triad*. One important result of this linkage is that both the dynamic and static accommodation responses in the two eyes are always essentially the same, even when one eye is occluded.^{363,364} Accommodation is driven by the responses of the cones of the retina and becomes progressively less effective as lighting levels are lowered from photopic levels through the mesopic region: at scotopic levels it is inoperative (see “Accuracy of Response”).^{365,366}

Although details of the neural and physiological mechanisms responsible for accommodation are beyond the scope of this section (see, e.g., Refs. 367–370 for reviews), it will be helpful to remember that, in simple terms, the crystalline lens is supported around its equator by a system (the zonule) of approximately radially oriented zonular fibres. The far ends of these fibres are anchored in the ciliary body or muscle and neighboring structures. The ciliary body effectively forms a ring surrounding the lens (Fig. 1). The lens and its thin enclosing capsule are elastic and, if the lens is free of the forces applied to it by the zonular fibres, its structure is such that its surfaces will naturally assume the relatively steep curvatures required for near vision.

Under conditions of distant viewing, the inner diameter of the ciliary ring is relatively large. This leads to a correspondingly high tension in the zonular fibers. These tensional forces when applied to

the periphery of the elastic lens and capsule cause the lens surfaces to flatten slightly, reducing the optical power to the value required for a clear distance focus. During active accommodation, the ciliary ring reduces in diameter. This relaxes the tension in the zonular fibres and allows the surfaces of the elastic lens to assume a steeper curvature and the power of the lens to increase. Attempts have been made to model this process (e.g., Refs. 371–374).

Dynamics of Response Typical records of a young subject's responses to abrupt changes in the position of an accommodation target are shown in Fig. 20. It can be seen that the change in position is followed by a short reaction time, or latency (about 0.4 s) during which the response remains unchanged. The response then progresses to its new level, the minimum response time typically being around 0.6 s.^{360,375–380} Under conditions where there are plentiful cues to target distance this response usually occurs in a single sweep, but if binocular and other cues to target position apart from defocus blur are eliminated, the response may become uncertain and may initially be in the wrong direction, leading to longer response times.^{381–383} The response times become larger for larger dioptric changes.^{377–380}

Another way of characterizing the response dynamics is in terms of their frequency response characteristics.³⁸⁴ These may be assessed by determining the gain and phase of the response obtained when a target is observed whose vergence is changing sinusoidally with time, as a function of the temporal frequency. It appears that when gain and phase are plotted as a function of the frequency both vary in an essentially linear fashion, with the cutoff frequency at which response no longer tracks the stimulus being about 2 Hz.³⁸⁵ It should be stressed that these characteristics are not the output of a simple reflex system but depend upon higher-order involvement. They are strongly influenced by training and motivation,^{386–389} and, with repetitive stimuli, by the knowledge that the required response has a predictable periodic form.³⁹⁰ When the response to an abrupt unexpected step change in target vergence is analyzed in terms of the corresponding frequency response, much larger phase lags are given.^{385,390} Thus the system does not behave linearly.³⁸⁵

The importance of perceptual factors in relation to accommodation is exemplified by studies in which the distance of the target (a Maltese cross) is kept constant but its lateral scale is varied sinusoidally.

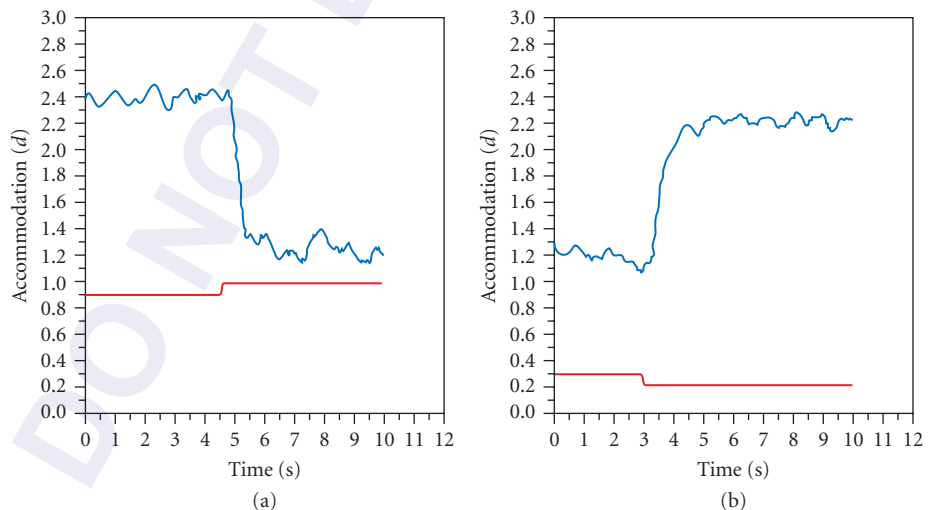


FIGURE 20 Examples of monocular accommodation responses to abrupt step changes in the vergence of an accommodation stimulus. Note the time interval (latency or reaction time) before each response starts to change and the time taken to complete the response (response time). Fluctuations in accommodation can be seen. The lower traces merely show the times at which the stimulus changes (between vergence levels of 2.38 and 1.33 D) occur.

Observers usually interpret this size variation at constant distance as a variation in distance of a target of constant linear size, and change their accommodation accordingly, even though this blurs the retinal image.³⁸¹ A similar effect is observed when binocular accommodation is stimulated by apparent depth elicited by pairs of stereograms at fixed distance.³⁹²

Stability of Response When a target at fixed distance is viewed with steady fixation, the accommodation response typically shows small fluctuations (~ 0.25 D) occurring at frequencies up to about 5 Hz (see Fig. 20): the fluctuations are correlated in the two eyes (see Refs. 393 and 394 for reviews). It appears that the lower frequency (0–0.6 Hz) components of the fluctuations may be under neural control but that higher frequencies (1.0–2.3 Hz) are related to physiological rhythms and other sources of noise.³⁹⁵ The power in the low-frequency components increases as the target luminance is reduced, while that in the high-frequency components remains stable.^{395,396} Similarly, the low-frequency components increase when the pupil is small and the depth-of-focus is large.³⁹⁷ A peak in the frequency spectrum at frequencies of 1 to 2 Hz is often observed^{398,399} and is correlated with the arterial pulse.^{400,401} The fluctuations appear to have their origin in the crystalline lens and tend to increase in magnitude as the mean accommodation level increases,^{397,402–404} although there is some suggestion that the higher-frequency components diminish at very high stimulus levels.³⁹⁵ There is some disagreement on the exact nature of the changes with age.^{395,405}

The possible role of these fluctuations in relation to accommodation control remains contentious. Some have argued that they simply represent plant noise^{384,406} and are of no utility. Others suggest that they could both guide the initial direction of the response and help to maintain accurate focus,^{101,399,407–411} the basic hypothesis being that if a fluctuation in one direction improves the clarity of the image the control system responds by moving the mean level of accommodation in that direction. Current opinion generally favors the concept that the high-frequency fluctuations represent plant noise and that any role in accommodation control is played by the lower-frequency components (below about 0.6 Hz), and that these would be primarily involved in the maintenance of the response at a steady level, rather than in rapid response changes.⁴¹⁰ In any case, under most circumstances the fluctuations appear to produce little variation in visual acuity,⁴¹² although their effects can just be detected.^{413,414}

Accuracy of Response Following pioneering work by Morgan,⁴¹⁵ numerous studies have shown that, rather than there being one-to-one matching between the dioptric stimulus level and the corresponding accommodation response, steady-state errors are an intrinsic part of the accommodation control system. This is true under both monocular and binocular conditions, although the magnitude of the errors differs in the two states.^{415,416} These focus errors are often the major cause of foveal retinal image degradation, rather than the higher-order aberrations discussed earlier (see “Factors Affecting Retinal Image Quality”). Figure 21a illustrates a schematic response/stimulus curve. This emphasizes that the steady-state is usually characterized by overaccommodation (accommodative “lead”) for distant objects and underaccommodation (“lag”) for near objects. The range of stimulus vergence over which there is no noticeable image blur is termed the *subjective amplitude of accommodation*; it obviously includes depth-of-focus effects. The corresponding, somewhat smaller, range of actual accommodation response is the *objective amplitude of accommodation*. Note that a clinically “emmetropic” or “corrected” subject is usually left with a small myopic refractive error for objects at infinity, ocular depth-of-focus being relied upon to give clear vision under these circumstances.

The slope of the quasi-linear region of the curve depends upon the observing conditions, including target form and contrast,^{100,101,417–426} ocular pupil diameter,^{427–429} luminance level,^{350,365,366,417} and the acuity of the observer.^{418,430–433} The common feature is that, as the quality of the stimulus available to the visual system degrades, the slope of the central region of the response/stimulus curve diminishes (see, e.g., Fig. 21b, where the slope falls as the luminance is reduced and cone vision and visual acuity are gradually lost). To a close approximation, for any individual, as the stimulus degrades the curve pivots about the point for which stimulus and response are equal, that is, where the curve crosses the ideal one-to-one stimulus/response line. It is of interest that in many studies it appears that there is

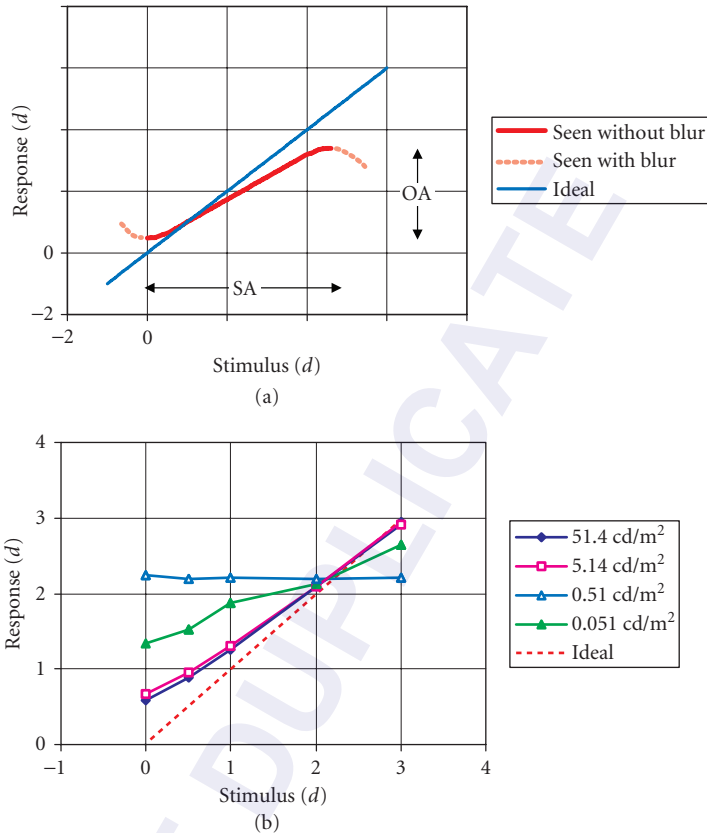


FIGURE 21 (a) Schematic appearance of a typical accommodation response/stimulus curve under photopic conditions. The black portion of the curve represents the range over which the target is seen without noticeable blur, the grey portions where the lags or leads exceed the ocular depth-of-focus and the target appears blurred. OA is the objective amplitude of accommodation and SA the subjective amplitude. The “ideal” line corresponds to equal response and stimulus values. (b) Typical changes in the lower part of the response/stimulus curve with target luminance. (Based on Johnson³⁶⁶) Target luminances are as indicated and the independently measured dark focus for the same subject is 2.0 D, corresponding closely to the cross-over point of the two curves.

a linear relationship between the slope of the central region of the response/stimulus curve and the eye’s minimum angle of resolution under the target and observing conditions in use.⁴³⁴

The extreme case in this process is where the stimulus is either a uniform photopic field completely lacking in detail (a Ganzfeld), or the field is completely dark. In both cases no spatial information is available to the accommodation control system and the response/stimulus curve becomes completely flat. The response remains constant at the somewhat myopic value at which, under normal photopic conditions with a structured target, response equals stimulus (see, e.g., Fig. 21b). The refractive states under these stimulus-free, light and dark conditions are known as *empty field myopia* and *dark focus* respectively, and for any individual these states are essentially the same.^{435–438}

These observations have led to the concept that this intermediate myopic level of accommodation is, in fact, the *tonic* level of accommodation (sometimes called the resting state or equilibrium level)

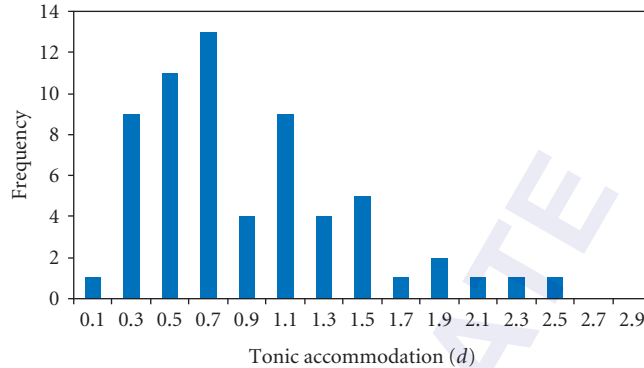


FIGURE 22 Frequency distribution of tonic accommodation as estimated by the dark focus. The values shown represent the difference in autorefractor measurements of refraction in photopic conditions and in darkness. (After McBrien and Millodot.⁴³⁹) The mean value of tonic accommodation and its standard deviation are 0.91 ± 0.53 D (62 young adult subjects).

to which the system returns in the absence of an adequate stimulus. This implies that accommodation must actively change from this level to view both distant and near targets. The tonic level varies somewhat with the individual (Fig. 22) but appears to have a mean value of around 1 D.^{437,439}

Current theories suggest that these characteristics are dictated by the balance between the sympathetic and parasympathetic innervations to the system.^{368,440–442}

Vergence Input As noted earlier, the convergence required to maintain single vision of a near object during binocular observation provides an input to the control system which supplements that due to pure defocus blur. This may be particularly useful in stimulating accommodation under conditions when the accommodation stimulus alone is relatively ineffective, for example, under mesopic night-driving conditions where the onset of *night (or twilight) myopia* is delayed when viewing is binocular.⁴⁴³

Control system models for the combined accommodation-convergence system, involving both fast feedback loops and slow adaptive elements, have been proposed by various authors (see, e.g., Refs. 368, 440, 444–446).

Application to Instrumentation From the practical point of view, it is evident that errors in accommodation and the resultant degradation in retinal image quality will be minimized if targets such as displays which are viewed by the naked eye are placed at vergences which approximately match the typical tonic levels of observers, that is, distances of about 1 m. Acuity will then be optimal³⁶⁶ and the accommodative effort, and hence potential fatigue, will be minimized. With color displays, the longitudinal chromatic aberration of the eye means that more accommodative effort is required for red/black symbol combinations than for blue/black combinations. Only minor variations in response have, however, been reported when two-color isoluminant symbol/background combinations are used:^{447–450} there is no evidence that viewing multichromatic targets results in a less stable response than that for black-and-white targets.⁴⁵¹

It is of interest that, when focusing visual instruments, most observers set the focus so that the light leaving the eyepiece is slightly divergent, that is, they prefer to accommodate slightly when viewing the instrumental image.⁴⁵² This *instrument myopia*⁴⁵³ correlates closely with the empty field myopia and dark focus of the individual observer,^{435–438} suggesting that these myopic states (the anomalous myopias) have a common origin and that the instrument focus is selected to minimize

accommodative effort. This implies that fixed-focus visual instruments should not be designed to have parallel light leaving their eyepieces but that focus should be set so that the image appears at a vergence of about -1 D.⁴⁵² In any binocular instrument where the imagery is at infinity and the optical axes are arranged to be parallel, proximal (psychic) convergence and accommodation often occur, that is most young users tend to converge and accommodate slightly, leading to a loss in visual performance.^{454,455} This is presumed to be due to the perception that the targets lie within the “black box” constituted by the instrument. Smith et al.⁴⁵⁶ have discussed the various problems that arise when the focus of binocular instruments is not compatible with the angle between the eyepiece tubes.

Many visual instruments display field curvature and, provided that the vergences involved are negative, this can be at least partly compensated for by accommodation by the observer as the field is scanned.^{325,326}

Age-Dependent Changes in Accommodation As the lens ages and thickens, its elastic constants change. This, in combination with other ocular changes, causes the efficiency of the accommodation system to diminish (see, e.g., Refs. 367–370 for reviews).

The most obvious change is in the amplitude of accommodation (Fig. 23). The mean subjective amplitude declines steadily from the age of about 10, to reach a small, constant level of about 1 D by the age of 50–55.⁴⁵⁷ This residual subjective amplitude represents depth-of-focus rather than true accommodation and the corresponding objective amplitude is zero.⁴⁵⁸ Longitudinal measurements of objective amplitudes suggest that the decline for any individual is linear with age.^{459–461} Although the exact values of the constants depend upon the individual subject and the method of measurement, the amplitude changes might typically be described by an equation of the form

$$OA = 12.7 - 0.27 (\text{age}) \quad (19)$$

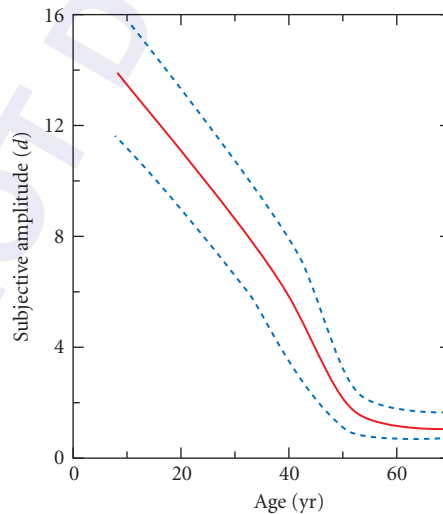


FIGURE 23 Transverse data for the changes in subjective amplitude of accommodation with age. (After Duane.⁴⁵⁶) The full curve represents the mean result, the dashed curves the limit of the range found at each age. The data refer to 4200 eyes, amplitudes being measured at the spectacle point 14 mm anterior to the cornea.

where OA is the objective amplitude in diopters and the age is in years.⁴⁶¹ Note that the conventional use of a value of 250 mm (i.e., an amplitude of accommodation of 4 D) for the “least distance of distinct vision” in the calculation of the nominal magnification of eyepieces and other devices is unlikely to reflect the effective magnification obtained by the individual user.

Apart from the change in amplitude, the general characteristics of the response/stimulus curve remain robust against age change up to about 35 years, when the slope starts to decline.^{362,462,463} The frequency response of the accommodation system maintains its cutoff value at around 2 Hz up to the age of about 40 but at any lower frequency the gain tends to decrease and the phase lag to increase with age. There are also changes in the time constants of the responses to step changes in stimulus.^{360,465–467}

The loss in amplitude becomes of practical significance when it declines sufficiently to begin to make it difficult to read at normal reading distances (around 1/3 m). It is tiring to exercise the full available amplitude of accommodation, so that problems usually start to arise at around the age of 40 when the subjective amplitude has fallen to about 5 D. These difficulties steadily increase as the amplitude declines further. Initially these problems may be eased by increases in reading distance, or the use of high light levels to increase depth-of-focus by constricting the pupil, but by the age of about 45 all emmetropes or corrected ametropes inevitably require some form of optical correction for near work. (See Chap. 14 by John S. Werner, Brooke E. Scheffrin, and Arthur Bradley.)

1.8 EYE MODELS

Optical models of the eye have long been used in attempts to better understand its imaging characteristics, the development of refraction, and the optical effects of spectacle and contact lenses (see, e.g., Refs. 468–470 for reviews). More recently such models have assumed additional importance in relation to the effects of a variety of surgical procedures for modifying the optics of the eye, such as corneal refractive surgery and intraocular lens implantation.⁴⁷¹ They have also been used in the evaluation of retinal radiation hazards.^{472,473} The earlier models were usually limited to describing behavior in the paraxial region, where rays make small angles with the axis: surfaces were assumed to be spherical. However recent years have seen the development of a succession of increasingly sophisticated paraxial and wide-angle models, incorporating such features as aspheric surfaces, gradient index optics, and accommodation.

Paraxial Models

Figure 24 shows examples of three typical types of paraxial model with spherical refracting surfaces and media characterized by single values of refractive index, the relevant parameters being given in Table 2. These paraxial models are useful for predicting the approximate dimensions of paraxial images (1 deg in the visual field corresponds to about 0.29 mm on the retina) and their changes with accommodation and correcting lenses. However, they are of more doubtful value in indicating retinal image quality on axis or in describing images in the peripheral field, and cannot predict the aberrations of real eyes.

In constructing such models, the values of the parameters are normally selected to be reasonably representative of those found in real eyes but are then adjusted slightly to make the eye emmetropic. In *schematic eyes*, the cornea and lens are each represented by a pair of surfaces (although sometimes the lens is divided into central nuclear and outer cortical regions, these being assigned different refractive indices). In the *simplified schematic eye* a single surface is used for the cornea. It can be seen in Fig. 24 that in both these eyes the two nodal points are very close together, as also are the two principal points. This has encouraged the development of *reduced eye* models, containing a single refractive surface, where each pair of nodal and principal points collapses to a single point. The cardinal points of reduced eyes are very close to those of the more complex paraxial models. Accommodation has been incorporated into several schematic eyes,^{469,476–478} some of which incorporate aspheric surfaces and lenticular index gradients.

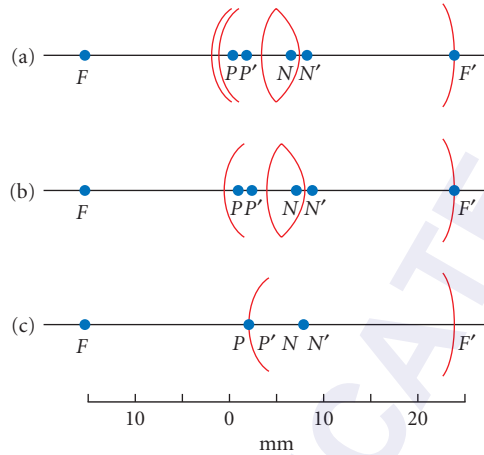


FIGURE 24 Examples of paraxial eye models of the human eye. In each case F, F' ; P, P' ; N, N' represent the first and second focal, principal, and nodal points respectively. (a) Unaccommodated schematic eye with four refracting surfaces. (*Le Grand and El Hage*.⁴⁷⁴) (b) Simplified schematic eye with three refracting surfaces. (*Emsley*.⁴⁷⁵) (c) Reduced eye with a single refracting surface. (*Emsley*.⁴⁷⁵) Note that progressive reduction in the number of surfaces used in the model produces only minor changes in the positions of the cardinal points.

TABLE 2 Parameters of Some Paraxial Models of the Human Eye

		Schematic Eye ⁴⁷⁴	Simplified Schematic Eye ⁴⁷⁵	Reduced Eye ⁴⁷⁵
Radii of surfaces (mm)	Anterior cornea	7.80	7.80	5.55
	Posterior cornea	6.50	-	-
	Anterior lens	10.20	10.00	-
	Posterior lens	-6.00	-6.00	-
Distances from anterior cornea (mm)	Posterior cornea	0.55	-	-
	Anterior lens	3.60	3.60	-
	Posterior lens	7.60	7.20	-
	Retina	24.20	23.90	-
	1st principal point P	1.59	1.55	0
	2nd principal point P'	1.91	1.85	0
	1st nodal point N	7.20	7.06	5.55
	2nd nodal point N'	7.51	7.36	5.55
1st focal point F		-15.09	-14.99	-16.67
	2nd focal point F'	24.20	23.90	22.22
Refractive indices	Cornea	1.3771	-	4/3
	Aqueous humour	1.3374	1.333	4/3
	Lens	1.4200	1.416	4/3
	Vitreous humour	1.3360	1.333	4/3

Wide-Angle Models

Following gradual advances in our knowledge of the form of the aspheric surfaces of the cornea and lens, together with the lenticular gradients of refractive index, several authors have produced sophisticated wide-angle eye models and have attempted to validate them by comparing the model predictions with the off-axis aberrations of real eyes.^{478–487} It is of interest that a very simple reduced eye model with a single aspheric refracting surface and a suitably placed pupil (the “Indiana Eye”) can simulate the chromatic, spherical, and oblique astigmatic aberrations typically found in real eyes.^{488–491}

As yet, none of these models is completely successful, but progressive refinement in the experimental data should lead to corresponding improvements in the modeling of the eye’s overall optical performance. Since, as discussed earlier, there are wide variations in both the on- and off-axis performance of individual eyes, it may be that “personalized” eye models, incorporating parameters as measured for the particular eye, will eventually prove to be of greatest value for use in predicting the outcomes of refractive surgery.

1.9 TWO EYES AND STEREOPSIS

Although binocular vision confers a variety of advantages, ranging from an extension of the field of view to a lowering of contrast thresholds,⁴⁹² attention here will be largely confined to its relevance for stereopsis and stereoscopic instruments. The relationship between the typical monocular and the binocular fields within which stereopsis can occur is shown in Fig. 25.

Just as the second-order wavefront errors (ocular refractions) of the two eyes are usually broadly similar, recent evidence suggests that the higher-order aberrations tend to show mirror symmetry between the two eyes:^{128–132,286,493} cone directionality (SCE I) also appears to be mirror symmetric.⁴⁹⁴ In cases where marked interocular differences in image quality occur, it is possible that the brain can make selective use of the better of the two retinal images under any particular set of observing conditions due to some form of probability summation. Certainly the apparently drastic technique of monovision contact or intraocular lens correction, in which one eye of a presbyopic observer receives a distance correction and the other a near correction, appears to work well for many patients and to yield acceptable vision over a range of distances.^{495–497} It may be the studies of monocular performance that can sometimes give an unduly pessimistic view of the optical information available under binocular conditions.

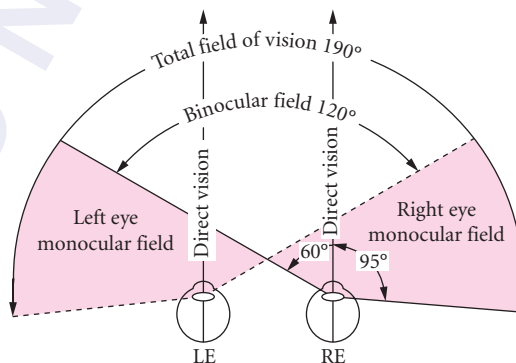


FIGURE 25 Approximate horizontal angular extents of the monocular and binocular fields of vision.

Basics of Stereoscopic Acuity

Due to the lateral separation of the eyes in the head, the apparent angular separations of objects at differing distances are slightly different for the two eyes. The resultant disparities between the two retinal images can be used by the visual system to estimate the relative distances of the objects, although absolute judgment is usually much more dependent upon such monocular cues as perspective and size constancy.

In Fig. 26a, suppose that the observer can just detect that the two object points A and B lie at different distances, l and $l + \delta l$, respectively. Then the corresponding stereo acuity $\delta\theta$ is given by:

$$\delta\theta = \theta_R - \theta_L = \alpha_1 - \alpha_2$$

Approximating all angles as being small [i.e., $l \gg \delta l$, p , where p is the lateral separation of the nodal points of the two eyes, or interpupillary distance (IPD)], and using the binomial expansion with omission of higher-order terms in δl yields:

$$\delta\theta = \frac{p \cdot \delta l}{l^2} \quad \text{or} \quad \delta l = \frac{l^2 \cdot \delta\theta}{p} \quad (20)$$

where $\delta\theta$ is in radians. Thus the minimum detectable difference in object distance is directly proportional to the square of the viewing distance and inversely proportional to the separation between the eyes (see, e.g., Schor and Flom⁴⁹⁸ for a more detailed analysis).

Figure 26b plots this approximate predicted value of the just-detectable difference in distance as a function of the distance l , on the assumption that $p = 65$ mm and $\delta\theta = 10$ sec arc. Note that discrimination of depth becomes very poor at distances in excess of about 500 m.

The interpupillary distance (IPD) varies somewhat between individuals and population groups.⁴⁹⁹ Typical distributions for men and women are illustrated in Fig. 27. Values range between about 50 and 76 mm. In the real world, of course, binocular cues to distance are supplemented by

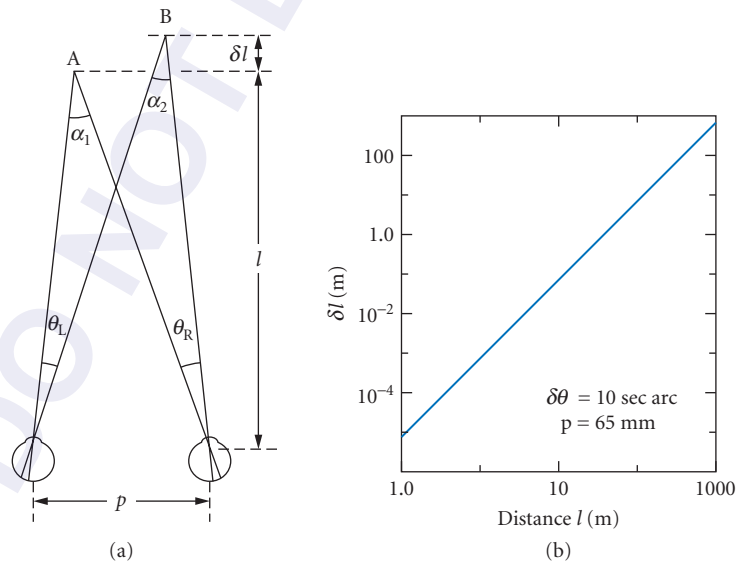


FIGURE 26 (a) Geometry of stereopsis. It is assumed that points A and B can just be discriminated in depth. (b) Theoretical just discriminable distance δl as a function of the object distance l for the assumed values of p and $\delta\theta$ indicated.

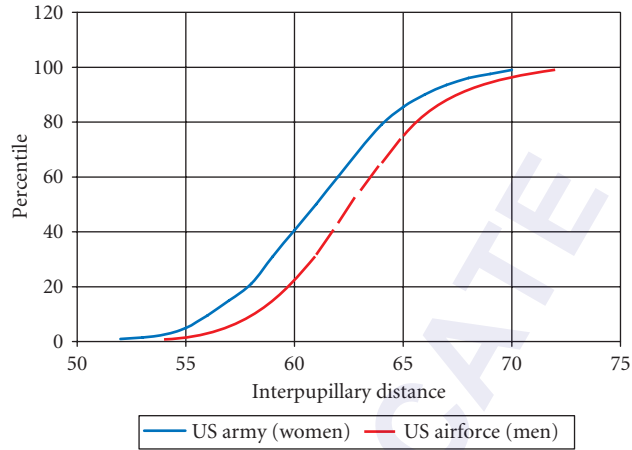


FIGURE 27 Typical cumulative frequency distributions of interpupillary distance for men and women.⁴⁹⁹

a variety of monocular cues such as perspective, overlay (interposition), size constancy, texture, motion parallax, etc.⁵⁰⁰

For any observer, $\delta\theta$ varies with such parameters as the target luminance and angular distance from fixation, and the observation time allowed (e.g., Refs. 501–504), being optimal close to fixation, with high luminance and extended observation times.⁵⁰⁵ Values also vary with the nature of the task involved. Clinical tests of stereoacuity (e.g., Refs. 506 and 507) which are usually carried out at a distance of 40 cm and are calibrated for an assumed IPD, p , of 65 mm, typically yield normal stereoacuties of about 20 to 40 sec arc (see Fig. 28).⁵⁰⁸ Two- or three-needle⁵⁰⁹ or similar tests carried out at longer distances usually give rather smaller values, of around 5 to 10 sec arc.

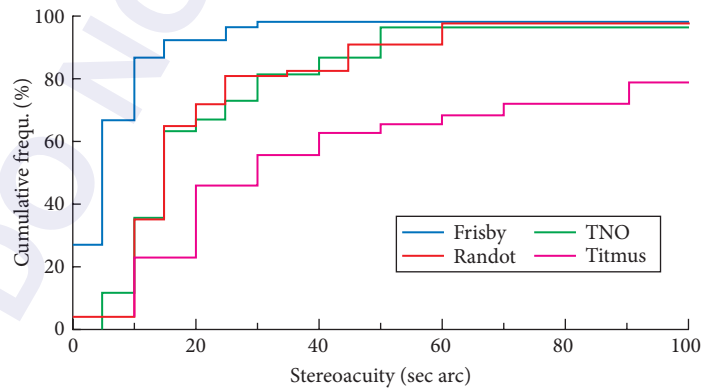


FIGURE 28 Cumulative frequency distribution for stereoscopic acuity, as measured by various clinical tests at a viewing distance of 400 mm, based on a sample of 51 adult subjects with normal binocular vision.⁵⁰⁸

Stereoscopic and Related Instruments

The stereoscopic acuity of natural viewing can be enhanced both by extending the effective IPD from p to Np with the aid of mirrors or prisms and by introducing transverse magnification M into the optical paths before each eye. This nominally has the effect of changing the just-detectable distance to

$$\delta l = \frac{l^2 \cdot \delta\theta}{MNp} \quad (21)$$

although in practice this improvement in performance is not always fully realized. Such changes will in general vary the spatial relationships in the perceived image, so that the object appears as having either enhanced or reduced depth in proportion to its lateral scale. If, for example, the magnification M is greater than one but N is unity, the object appears nearer but foreshortened. Simple geometrical predictions of such effects are, however, complicated by a variety of factors such as the reduced depth-of-field of magnifying system (see, e.g., Refs. 501, 510–512).

As a result of the range of IPD values encountered among different potential users (Fig. 27), it is important that adequate adjustment of IPD be provided (preferably covering 46 to 78 mm), with a scale so that users can set their own IPD. Convergence should be appropriate to the distance at which the image is viewed, that is, the angle between the visual axes should be approximately 3.7 deg for each diopter of accommodation exercised.⁵¹³

Tolerances in Binocular Instrumentation and the Problem of Aniseikonia

In the foregoing, it has been assumed that the images available to the two eyes can be successfully fused to yield full binocular vision. Although this does not demand exact alignment of the two images, since, for example, horizontal misalignment can be compensated for by appropriate convergence or divergence between the visual axes of the user's eyes, such compensation is only possible over a limited range. Several authors have suggested instrumental tolerances appropriate to different circumstances (particularly duration of use) for vertical misalignment, convergence error, divergence error, magnification difference, and rotation in the images presented to the eyes by binocular instruments (see, e.g., Refs. 513–515 for reviews). Recommended tolerances on divergence and vertical misalignment (usually around 3 to 10 min arc) are smaller than those on convergence (around 10 to 20 min arc). These tolerances relate to the eyes: evidently if, for example, an afocal binocular instrument provides lateral magnification, the tolerances on the alignment of the tubes of the right- and left-eye optical systems will be reduced accordingly. Cyclofusion to compensate for relative image rotation is difficult to maintain on a sustained basis: any tolerance is independent of magnification.

Particular interest attaches to the possibility of a magnification difference between the two images (*aniseikonia*).⁵¹⁶ In practice, such magnification differences are probably most likely to occur as a result of refractive spectacle corrections in patients with markedly different refractions in the two eyes (*anisometropia*).⁵¹⁷ Strictly speaking we are not concerned with the retinal images as such but with the corresponding perceived images, since it is possible that these may differ in size as a result of neural as well as purely optical factors. It is conventional to express the relevant magnification differences in percentage terms. As a rough guide, magnification differences less than 1 percent usually present no problems and those greater than 5 percent often make it difficult to maintain fusion, leading to diplopia or suppression. It is the range between about 1 and 5 percent of perceived size difference where disturbing problems in spatial perception may arise but binocular vision is still possible. These may in turn lead to symptoms such as eyestrain, nausea, and headaches.

In eyes demanding a correction for both spherical and astigmatic errors, size differences may vary in different meridians, depending upon the axes of the correcting cylinders. The spatial distortions resulting from horizontal and vertical size differences are different in nature but, in general, objects may appear tilted and unpleasant instabilities in the perception of space may arise as the head or eyes turn. To give a rough estimate of the conditions under which significant aniseikonia may occur, we make use of the crude approximation that the spectacle magnification, expressed in percentage terms

is about $100aF$, where a mm is the distance of the spectacle lens from the eye and F diopter is its power. Thus the spectacle magnification is around 1.5 percent per diopter (positive for converging lenses, negative for diverging lenses). Thus, with a spectacle correction, problems may begin to arise even with quite modest degrees of anisometropia: these are reduced with contact lenses,⁵¹⁸ which give much lower spectacle magnification since $a \approx 0$. Even for ametropic patients without significant anisometropia, size differences will arise if one eye is corrected with, for example, corneal refractive surgery and the other with a spectacle lens. Similarly, correction of one eye with an intraocular lens following cataract surgery while the other eye remains spectacle corrected may also cause difficulties.⁵¹⁶ Methods for measuring aniseikonia and for its control are discussed in Refs. 516 and 517; fortunately many patients adapt to modest image size differences.

1.10 MOVEMENTS OF THE EYES

Movements of the eyes are of significance in relation to visual optics from two major points of view. First, since under photopic conditions both optical and neural performance are optimal on the visual axis, the eye movement system must be capable of rapidly directing the eyes so that the images of the detail of interest fall on the central foveas of both eyes, where visual acuity is highest (*gaze shifting* leading to *fixation*). A scene is explored through a series of such fixational movements (*saccades*) between different points within the field.^{519,520} Second, the system must be capable of maintaining the images on the two foveas both when the object is fixed in space (*gaze holding*) and, ideally, when it is moving. Any lateral movement of the images with respect to the retina is likely to result in degraded visual performance, due to the limited temporal resolution of the visual system and the falloff in acuity with distance from the central fovea.

These challenges to the eye movement control system are further complicated by the fact that the eyes are mounted in what is, in general, a moving rather than a stationary head. Movements of the eyes therefore need to be linked to information derived from the vestibular system or labyrinth of the inner ear, which signals rotational and translational accelerations of the head. The compensatory *vestibulo-ocular responses* take place automatically (i.e., they are reflex movements) whereas the fixational changes required to foveate a new object are voluntary responses. Details of the subtle physiological mechanisms which have evolved to meet these requirements will be found in Refs. 521–524.

Basics

Each eye is moved in its orbit by the action of three pairs of extraocular muscles attached to the outside of the globe. Their action rotates the eye about an approximate center of rotation lying some 13.5 mm behind the cornea, although there is in fact no unique fixed point within the eye or orbit around which the eye can rotate to all possible positions that it can assume.^{521,523} In particular, the “center of rotation” for vertical eye movements lies about 2 mm nearer the cornea than that for horizontal eye movements (mean values 12.3 mm and 14.9 mm, respectively).⁵²⁵

Although the two eyes can scan a field extending about 45 deg in all directions from the straight-ahead or *primary* position, in practice eye movements rarely exceed about 20 deg, fixation on more peripheral objects usually being achieved by a combination of head and eye movements.

If the angle between the two visual axes does not change during the movement, the latter is described as a *version* (or conjugate) movement. However, the lateral separation of the eyes in the head implies the need for an additional class of movements to cope with the differing convergence requirements of objects at different distances. These movements, which involve a change in the angle between the visual axes, are called *vergence* (or disjunctive) movements. Fixational changes may in general involve both types of movement, which appear to be under independent neurological control.⁵²⁶

Characteristics of the Movements

The version movements involved in bringing a new part of the visual field onto the fovea (saccades) are very rapid, being completed in around 100 ms with angular velocities often reaching more than 700 deg/s, depending upon the amplitude of the movement^{527,528} (see Fig. 29a): the saccadic latency is about 200 ms.⁵²⁹ Note that the larger saccades tend initially to be inaccurate. Interestingly, it appears that during the saccade, when the image is moving very rapidly across the retina, vision is largely, although not completely, suppressed. This saccadic suppression (or perhaps, more appropriately, saccadic attenuation) results in, for example, the retinal thresholds for brief flashes of light being elevated, the elevation commencing some 30 to 40 ms before the actual saccadic movement starts. The subject is normally unaware of this temporary impairment of vision and the exact mechanisms responsible for it remain controversial,^{522,529} although an explanation may lie in the masking effect of the clear images available before and after the saccade.⁵³⁰

Smooth voluntary pursuit movements of small targets which are moving sinusoidally in a horizontal direction are accurate at temporal frequencies up to a few hertz. Their peak velocities range up to about 25 deg/s. In practice, when a small moving target is tracked, the following movements usually consist of a mixture of smooth movements and additional corrective saccades (e.g., Ref. 531; see Fig. 29b). With repetitive or other predictable stimuli, tracking accuracy tends to improve markedly with experience, largely through the reduction of the phase lag between target and eye.

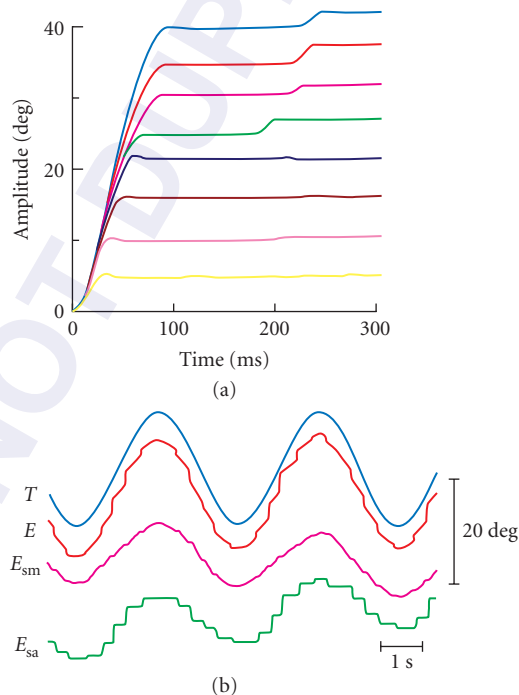


FIGURE 29 (a) Time course of saccades of different sizes. (After Robinson.⁵²⁷) The traces have been superimposed so that the beginning of each saccade occurs at time zero. (b) Separation of smooth (E_{sm}) and saccadic (E_{sa}) components of foveal tracking. (After Collewijn and Tamminga.⁵³¹) T is the target position and E is the eye position.

When the image of a substantial part of the visual field moves uniformly across the retina, the eyes tend to rotate in a following movement to approximately stabilize the retinal image, until eventually the gaze is carried too far from the primary position, when a rapid anticompany flick-back occurs. Thus a continuously moving visual scene, such as the view from a train, causes a series of slow following movements and fast recoveries, so that a record of the eye movements takes a quasi-regular sawtooth form. At slow object speeds, the angular velocity of the slow phase is close to that of the field, but as the field velocity rises, greater lags in eye velocity occur, until at about 100 deg/s the following movements break down. Although this *optokinetic nystagmus* is basically a reflex, it can be influenced by any instructions given to the subject.

The angular speeds of the vergence eye movements (about 5 to 10 deg/s per degree of vergence⁵³²) are usually described as being much lower than those for version movements. However, studies with more natural 3-D targets suggest that it is possible for speeds to be much higher (up to 200 deg/s for 35-deg vergence movements^{533,534}). Vergence movements typically follow a quasi-exponential course with a time constant of about 300 ms. Their latency is about 200 ms. Although the primary stimulus for vergence movements is disparity, that is, the difference in position of the images in the two eyes with respect to their foveas, vergence movements can also be driven by accommodation (see “The Accommodation Response”) and by perceptual factors such as perspective in line drawings.⁵³⁵

The vergence system is primarily designed to cope with the differences in the directions of the visual axes of the two eyes which arise from their horizontal separation in the head. However, some vertical vergence movements and relative torsional movements about the approximate direction of the visual axis can also occur. Such fusional movements can compensate for corresponding small relative angular misalignments in the eyes or in the optical systems of binocular or biocular instruments. The maximum amplitude of these movements is normally small (about 1 deg for vertical movements and a few degrees for cyclofusional movements, see, e.g., Ref. 336). They also take longer to complete than horizontal vergence movements (some 8 to 10 s as compared to about 1 s). Due to the limited effectiveness of vertical and torsional vergence movements, it is typically recommended that in instruments where the two eyes view an image through separate optical systems a 10 min arc tolerance be set for vertical misalignment and a 30 min arc limit be set for rotational differences between the two images (see “Tolerances in Binocular Instrumentation”).

Stability of Fixation

When an observer attempts to maintain fixation on a stationary target, it is found that the eyes are not steady, but that a variety of small-amplitude movements occur. These *miniature eye movements* can be broken down into three basic components: tremor, drift, and microsaccades. The frequency spectrum of the tremor falls essentially linearly with the logarithm of the frequency above 10 Hz, extending to about 200 Hz.⁵³⁷ The amplitude is small (probably less than the angle subtended at the nodal point of the eye by the diameter of the smallest foveal cones, i.e., about 24 sec arc). Drift movements are much larger and slower, with amplitudes of about 2 to 5 min arc at velocities around 4 min/s.⁵³⁸ The errors in fixation brought about by the slow drifts (which are usually dissociated in the two eyes) are corrected microsaccades: these are correlated in the two eyes. There are large intersubject differences in both mean microsaccade amplitude (from 1 to 23 min arc) and intersaccadic interval (from about 300 ms to 5 s),⁵³⁹ which probably reflect voluntary control under experimental conditions.

The overall stability of fixation can be illustrated by considering the statistical variation in the point of fixation (Fig. 30).^{540,541} For most of the time the point of regard lies within a few minutes of arc of the target. Although it has been suggested that these small eye movements could have some role in reducing potential aliasing problems,⁵⁴² experiments measuring contrast sensitivity for briefly presented interference fringes on the retina suggest that this is unlikely.⁵⁴³ Interestingly, when a suitable optical arrangement is used to counteract these small changes in fixation and stabilize the image on the retina, the visual image may fragment or even disappear completely,⁵³⁸ so that these small movements are important for normal visual perception. Fincham’s suggestion³⁸¹ that small eye movements are of importance to the ability of the eye to respond correctly to a change in accommodation stimulus has never been properly explored.

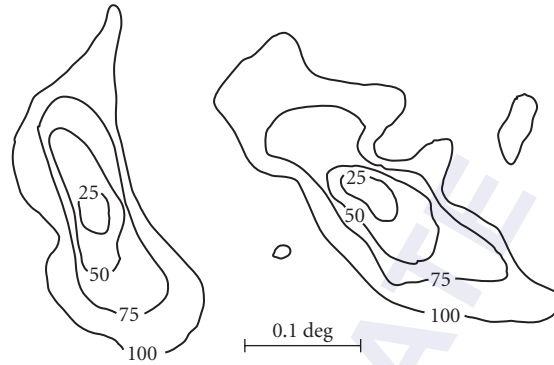


FIGURE 30 Stability of fixation for two subjects. The contours define areas within which the point of fixation was to be found for 25, 50, 75, and 100 percent of the time. (After Bennet-Clark.⁵⁴¹)

1.11 CONCLUSION

Recent years have seen considerable advances in our understanding of the aberrations of the eye and their dependence on factors such as pupil size, age, and field angle. The refractive index distribution of the lens remains to be fully elucidated, although substantial progress has been made. In general, it appears that the optical characteristics of the eye are such that the natural retinal image quality under different conditions of lighting, field, and age is well matched to the corresponding needs of the neural parts of the visual system. Those aberrations which are present may well be useful in such roles as guidance of the accommodation response, expansion of the depth-of-focus, and control of aliasing.

1.12 REFERENCES

1. M. Rynders, B. Lidkea, W. Chisholm, and L. N. Thibos, "Statistical Distribution of Foveal Transverse Aberration, Pupil Centration and Angle ψ in a Population of Young Adult Eyes," *J. Opt. Soc. Am. A* **12**:2348–2357 (1995).
2. J. Tabernero, A. Benito, E. Alc3n, and P. Artal, "Mechanism of Compensation of Aberrations in the Human Eye," *J. Opt. Soc. Am. A* **24**:3274–3283 (2007).
3. S. Guidarelli, "Off-axis Imaging in the Human Eye," *Atti d. Fond. G. Ronchi* **27**:449–460 (1972).
4. J. S. Larsen, "The Sagittal Growth of the Eye IV: Ultrasonic Measurements of the Axial Growth of the Eye from Birth to Puberty," *Acta Ophthalmologica* **49**:873–886 (1971).
5. S. J. Isenberg, D. Neumann, P. Y. Cheong, L. Ling, L. C. McCall and A. J. Ziffer, "Growth of the Internal and External Eye in Term and Preterm Infants," *Ophthalmology* **102**:827–830 (1995).
6. F. C. Pennie, I. C. J. Wood, C. Olsen, S. White, and W. N. Charman, "A Longitudinal Study of the Biometric and Refractive Changes in Full-Term Infants During the First Year of Life," *Vision Res.* **41**:2799–2810 (2001).
7. A. Sorsby, B. Benjamin, M. Sheridan, and J. M. Tanner, "Emmetropia and Its Aberrations," *Spec. Rep. Ser. Med. Res. Coun.* no. 293, HMSO, London, 1957.
8. A. Sorsby, B. Benjamin, and M. Sheridan, "Refraction and Its Components during Growth of the Eye from the Age of Three," *Spec. Rep. Ser. Med. Res. Coun.* no. 301, HMSO, London, 1961.

9. A. Sorsby and G. A. Leary, "A Longitudinal Study of Refraction and Its Components During Growth," *Spec. Rep. Ser. Med. Res. Coun.* no. 309, HMSO, London, 1970.
10. R. A. Weale, *A Biography of the Eye*, H. K. Lewis, London, 1982, pp. 94–120.
11. T. Grosvenor, "Reduction in Axial Length with Age: an Emmetropizing Mechanism for the Adult Eye?" *Am. J. Optom. Physiol. Opt.* **64**:657–663 (1987).
12. F. A. Young and G. A. Leary, "Refractive Error in Relation to the Development of the Eye," *Vision and Dysfunction, Vol 1, Visual Optics and Instrumentation*, W. N. Charman (ed.), Macmillan, London, 1991, pp. 29–44.
13. J. F. Koretz, P. C. Kaufman, M. W. Neider, and P. A. Geckner, "Accommodation and Presbyopia in the Human Eye—Aging of the Anterior Segment," *Vision Res.* **29**:1685–1692 (1989).
14. A. Glasser, M. A. Croft, and P. L. Kaufman, "Aging of the Human Crystalline Lens," *Int. Ophthalmol. Clin.* **41**(2):1–15 (2001).
15. J. E. Koretz, C. A. Cook, and P. L. Kaufman, "Aging of the Human Lens: Changes in Lens Shape upon Accommodation and with Accommodative Loss," *J. Opt. Soc. Am. A* **19**:144–151 (2002).
16. B. K. Pierscionek, D. Y. C. Chan, J. P. Ennis, G. Smith, and R. C. Augusteyn, "Non-destructive Method of Constructing Three-dimensional Gradient Index Models for Crystalline Lenses: I Theory and Experiment," *Am. J. Optom. Physiol. Opt.* **65**:481–491 (1988).
17. B. K. Pierscionek, "Variations in Refractive Index and Absorbance of 670 nm Light with Age and Cataract Formation in the Human Lens," *Exp. Eye Res.* **60**:407–414 (1995).
18. C. E. Jones, D. A. Atchison, R. Meder, and J. M. Pope, "Refractive Index Distribution and Optical Properties of the Isolated Human Lens Measured Using Magnetic Resonance Imaging," *Vision Res.* **45**:2352–2366 (2005).
19. C. E. Jones, D. A. Atchison, and J. M. Pope, "Changes in Lens Dimensions and Refractive Index with Age and Accommodation," *Optom. Vis. Sci.* **84**:990–995 (2007).
20. A. Steiger, *Die Entstehung der sphaerischen Refraktionen des menschlichen Auges*, Karger, Berlin, 1913.
21. S. Stenström, "Untersuchungen über die Variation und Kovariation der optischen Elements des menschlichen Auges," *Acta Ophthalmologica* suppl. 26 (1946) [Translated by D. Woolf as: "Investigation of the Variation and Covariation of the Optical Elements of Human Eyes," *Am. J. Optom. Arch. Am. Acad. Optom.* **25**:218–232; 286–299; 340–350; 388–397; 438–449; 496–504 (1948)].
22. P. M. Kiely, G. Smith, and L. G. Carney, "The Mean Shape of the Human Cornea," *Optica Acta* **29**:1027–1040 (1982).
23. M. Guillon, D. P. M. Lydon, and C. Wilson, "Corneal Topography: a Clinical Model," *Ophthal. Physiol. Opt.* **6**:47–56 (1986).
24. D. A. Atchison and G. Smith, "Continuous Gradient Index and Shell Models of the Human Lens," *Vision Res.* **35**:2529–2538 (1995).
25. L. F. Garner and G. Smith, "Changes in Equivalent and Gradient Refractive Index of the Crystalline Lens with Accommodation," *Optom. Vis. Sci.* **74**:114–119 (1997).
26. G. Smith and B. K. Pierscionek, "The Optical Structure of the Lens and its Contribution to the Refractive Status of the Eye," *Ophthal. Physiol. Opt.* **18**:21–29 (1998).
27. H. T. Kasprzak, "New Approximation for the Whole Profile of the Human Lens," *Ophthal. Physiol. Opt.* **20**:31–43 (2000).
28. R. Navarro, F. Palos, and L. M. González, "Adaptive Model of the Gradient Index of the Human Lens. I. Formulation and Model of Aging *ex vivo* Lenses," *J. Opt. Soc. Am. A* **24**:2175–2185 (2007).
29. R. Navarro, F. Palos, and L. M. González, "Adaptive Model of the Gradient Index of the Human Lens. II. Optics of the Accommodating Aging Lens," *J. Opt. Soc. Am. A* **24**:2911–2920 (2007).
30. J. A. Diaz, C. Pizarro, and J. Arasa, "Single Dispersive Gradient-index Profile for the Aging Human Lens," *J. Opt. Soc. Am. A* **25**:250–261 (2008).
31. P. Rosales and S. Marcos, "Phakometry and Lens Tilt and Decentration using a Custom-developed Purkinje Imaging Apparatus: Validation and Measurements," *J. Opt. Soc. Am. A* **23**:509–520 (2006).
32. C. F. Wildsoet, "Structural Correlates of Myopia," *Myopia and Nearwork*, M. Rosenfield and B. Gilmartin (eds.), Oxford, Butterworth-Heinemann, 1998, pp. 31–56.
33. C. F. Wildsoet, "Active Emmetropization—Evidence for its Existence and Ramifications for Clinical Practice," *Ophthal. Physiol. Opt.* **17**:279–290 (1997).

34. W. M. Lyle, "Changes in Corneal Astigmatism with Age," *Am. J. Optom. Arch. Am. Acad. Optom.* **48**:467–478 (1971).
35. R. J. Farrell and J. M. Booth, *Design Handbook for Imagery Interpretation Equipment*, Boeing Aerospace Co., Seattle, Sec. 3. 2, p. 8 (1984).
36. B. Winn, D. Whitaker, D. Elliott, and N. J. Phillips, "Factors Affecting Light-adapted Pupil Size in Normal Human Subjects," *Invest. Ophthalmol. Vis. Sci.* **35**:1132–1137 (1994).
37. I. E. Loewenfeld, *The Pupil: Anatomy, Physiology and Clinical Applications*, vols. I and II, Butterworth-Heinemann, Oxford, 1999.
38. P. Reeves, "Rate of Pupillary Dilation and Contraction," *Psychol. Rev.* **25**:330–340 (1918).
39. R. A. Weale, *Focus on Vision*, Hodder and Stoughton, London, 1982, p. 133.
40. G. Walsh, "The Effect of Mydriasis on the Pupillary Centration of the Human Eye," *Ophthalm. Physiol. Opt.* **8**:178–182 (1988).
41. M. A. Wilson, M. C. W. Campbell, and P. Simonet, "Change of Pupil Centration with Change of Illumination and Pupil Size," *Optom. Vis. Sci.* **69**:129–136 (1992).
42. J. M. Woodhouse, "The Effect of Pupil Size on Grating Detection at Various Contrast Levels," *Vision Res.* **15**:645–648 (1975).
43. J. M. Woodhouse and F. W. Campbell "The Role of the Pupil Light Reflex in Aiding Adaptation to the Dark" *Vision Res.* **15**:649–65 (1975).
44. G. Wyszecki and W. S. Stiles, *Colour Science*, Wiley, New York, 1967, pp. 214–219
45. E. Ludvigh and E. F. McCarthy, "Absorption of Visible Light by the Refractive Media of the Human Eye," *Arch. Ophthalmol.* **20**:37–51 (1938).
46. E. A. Boerttner and J. R. Wolter, "Transmission of Ocular Media," *Invest. Ophthalmol.* **1**:776–783 (1962).
47. W. J. Geeraets and E. R. Berry, "Ocular Spectral Characteristics as Related to Hazards for Lasers and Other Sources," *Am. J. Ophthalmol.* **66**:15–20 (1968).
48. T. J. T. P. van den Berg and H. Spekrijse, "Near Infrared Light Absorption in the Human Eye Media," *Vision Res.* **37**:249–253 (1997).
49. G. L. Savage, C. A. Johnson, and D. I. Howard, "A Comparison of Noninvasive Objective and Subjective Measurements of the Optical density of Human Ocular Media," *Optom. Vis. Sci.* **78**:386–395 (2001).
50. D. G. Pitts, "The Ocular Effects of Ultraviolet Radiation," *Am. J. Optom. Physiol. Opt.* **55**:19–53 (1978).
51. J. Pokorny, V. C. Smith, and M. Lutze, "Aging of the Human Lens," *Applied Opt.* **26**:1437–1440 (1987).
52. P. A. Sample, F. D. Esterson, R. N. Weinreb, and R. M. Boynton, "The Aging Lens: *In vivo* Assessment of Light Absorption in 84 Human Eyes," *Invest. Ophthalmol. Vis. Sci.* **29**:1306–1311 (1988).
53. R. A. Weale, "Age and the Transmittance of the Human Crystalline Lens," *J. Physiol. (London)* **395**:577–587 (1988).
54. N. P. A. Zagers and D. van Norren, "Absorption of the Eye Lens and Macular Pigment Derived from the Reflectance of Cone Receptors," *J. Opt. Soc. Am. A* **21**:2257–2268 (2004).
55. S. Lerman, *Radiant Energy and the Eye*, Balliere Tindall, London, 1980.
56. K. Knoblauch, F. Saunders, M. Kasuda, R. Hynes, M. Podgor, K. E. Higgins, and F. M. de Monasterio, "Age and Illuminance Effects in the Farnsworth-Munsell 100-hue Test," *Applied Opt.* **26**:1441–1448 (1987).
57. K. Sagawa and Y. Takahashi, "Spectral Luminous Efficiency as a Function of Age," *J. Opt. Soc. Am. A* **18**:2659–2667 (2001).
58. J. Mellerio, "Yellowing of the Human Lens: Nuclear and Cortical Contributions," *Vision Res.* **27**:1581–1587 (1987).
59. C. Schmidtt, J. Schmidtt, A Wegener, and O. Hockwin, "Ultraviolet Radiation as a Risk Factor in Cataractogenesis," *Risk Factors in Cataract Development, Dev. Ophthalmol.* **17**:Karger, Basel, 1989, pp. 169–172.
60. J. Marshall, "Radiation and the Ageing eye," *Ophthalm. Physiol. Opt.* **5**:241–263 (1985).
61. A. Stanworth and E. J. Naylor, "The Measurement and Clinical Significance of the Haidinger Effect," *Trans. Ophthalmol. Soc. UK* **75**:67–79 (1955).
62. P. E. Kilbride, K. B. Alexander, M. Fishman, and G. A. Fishman. "Human Macular Pigment Assessed by Imaging Fundus Reflectometry," *Vision Res.* **26**:663–674 (1989).

63. D. M. Snodderly, J. D. Auran, and F. C. Delori, "The Macular Pigment. II Spatial Distribution in Primate Retinas," *Invest. Ophthalmol. Vis. Sci.* **5**:674–685 (1984).
64. B. R. Hammond, B. R. Wooten, and D. M. Snodderly, "Individual Variations in the Spatial Profile of Human Macular Pigment," *J. Opt. Soc. Am. A* **14**:1187–1196 (1997).
65. V. M. Reading and R. A. Weale, "Macular Pigment and Chromatic Aberration," *J. Opt. Soc. Am.* **64**:231–234 (1974).
66. G. Haegerstrom-Portnoy, "Short-Wavelength Cone Sensitivity Loss with Aging: A Protective Role for Macular Pigment?," *J. Opt. Soc. Am. A* **5**:2140–2145 (1988).
67. L. J. Bour, "Polarized Light and the Eye," *Vision and Visual Dysfunction*, vol. 1, *Visual Optics and Instrumentation*, W. N. Charman (ed.), Macmillan, London, 1991, pp. 310–325.
68. W. S. Stiles and B. H. Crawford, "The Luminous Efficiency of Rays Entering the Eye Pupil at Different Points," *Proc. Roy. Soc. London B* **112**:428–450 (1933).
69. J. A. van Loo and J. M. Enoch, "The Scotopic Stiles-Crawford Effect," *Vision Res.* **15**:1005–1009 (1975).
70. J. M. Enoch and H. E. Bedell, "The Stiles-Crawford Effects," *Vertebrate Photoreceptor Optics*, J. Enoch and F. L. Tobey (eds.), Springer-Verlag, Berlin, 1981, pp. 83–126.
71. J. M. Enoch and V. Lakshminarayanan, "Retinal Fibre Optics," *Vision and Visual Dysfunction*, vol. 1, *Visual Optics and Instrumentation*, W. N. Charman (ed.), Macmillan, London, 1991, pp. 280–309.
72. W. S. Stiles, "The Luminous Efficiency of Monochromatic Rays Entering the Eye Pupil at Different Points and a New Colour Effect," *Proc. R. Soc. London B* **123**:90–118 (1937).
73. P. Moon and D. E. Spencer, "On the Stiles-Crawford Effect," *J. Opt. Soc. Am.* **34**:319–329 (1944).
74. M. Alpern, "The Stiles-Crawford Effect of the Second Kind (SCII): A Review," *Perception* **15**:785–799 (1986).
75. D. van Norren and L. F. Tiemeijer, "Spectral Reflectance of the Human Eye," *Vision Res.* **26**:313–330 (1986).
76. F. C. Delori and K. P. Pflibsen, "Spectral Reflectance of the Human Fundus," *Applied Opt.* **28**:1061–1077 (1989).
77. A. Elsner, S. A. Burns, J. J. Weiter, and F. C. Delori, "Infrared Imaging of Sub-retinal Structures in the Human Ocular Fundus," *Vision Res.* **36**:191–205 (1996).
78. F. C. Delori and S. A. Burns, "Fundus Reflectance and the Measurement of Crystalline Lens Density," *J. Opt. Soc. Am. A* **13**:215–226 (1996).
79. F. W. Campbell and R. W. Gubisch, "Optical Quality of the Human Eye," *J. Physiol (London)* **186**:558–578 (1966).
80. W. N. Charman and J. A. M. Jennings, "Objective Measurements of the Longitudinal Chromatic Aberration of the Human Eye," *Vision Res.* **16**:999–1005 (1976).
81. J. van de Kraats, T. T. J. M. Berendschot, and D. van Norren, "The Pathways of Light Measured in Fundus Reflectometry," *Vision Res.* **36**:2229–2247 (1996).
82. P. J. Delint, T. T. J. M. Berendschot, and D. van Norren, "Local Photoreceptor Alignment Measured with a Scanning Laser Ophthalmoscope," *Vision Res.* **37**:243–248 (1997).
83. J. C. He, S. Marcos, and S. A. Burns, "Comparison of Cone Directionality Determined by Psychophysical and Reflectometric Techniques," *J. Opt. Soc. Am. A* **16**:2363–2369 (1999).
84. N. P. A. Zagers, J. van de Kraats, T. T. J. M. Berendschot, and D. van Norren, "Simultaneous Measurement of Foveal Spectral Reflectance and Cone-photoreceptor Directionality," *Applied Opt.* **41**:4686–4696 (2002).
85. S. S. Choi, N. Doble, J. Lin, J. Christou, and D. R. Williams, "Effect of Wavelength on *in vivo* Images of the Human Cone Mosaic," *J. Opt. Soc. Am. A* **22**:2598–2605 (2005).
86. B. S. Jay, "The Effective Pupillary Area at Varying Perimetric Angles," *Vision Res.* **1**:418–428. (1962).
87. D. A. Atchison and G. Smith, *Optics of the Human Eye*, Butterworth-Heinemann, Oxford, 2000, pp. 25–27.
88. H. E. Bedell and L. M. Katz, "On the Necessity of Correcting Peripheral Target Luminance for Pupillary Area," *Am. J. Optom. Physiol. Opt.* **59**:767–769 (1982).
89. W. N. Charman, "Light on the Peripheral Retina," *Ophthal. Physiol. Opt.* **9**:91–92 (1989).
90. K. P. Pflibsen, O. Pomarentzeff, and R. N. Ross, "Retinal Illuminance Using a Wide-angle Model of the Eye," *J. Opt. Soc. Am. A* **5**:146–150 (1988).
91. A. C. Koosman and F. K. Witmer, "Ganzfeld Light Distribution on the Retina of Humans and Rabbit Eyes: Calculations and *in vitro* Measurements," *J. Opt. Soc. Am. A* **3**:2116–2120 (1986).

92. W. D. Wright, *Photometry and the Eye*, Hatton Press, London, 1949.
93. G. Airy, "On the Diffraction of an Object Glass with Circular Aperture," *Trans. Camb. Philos. Soc.* **5**:283–291 (1835).
94. G. M. Byram, "The Physical and Photochemical Basis of Resolving Power. I. The Distribution of Illumination in Retinal Images," *J. Opt. Soc. Am.* **34**:571–591 (1944).
95. U. Hallden, "Diffraction and Visual Resolution. I. The Resolution of Two Point Sources of Light," *Acta Ophthalmol.* **51**:72–79 (1973).
96. L. A. Riggs, "Visual Acuity," *Vision and Visual Perception*, C. H. Graham (ed.), Wiley, New York, 1966, pp. 321–349.
97. E. Lommel, "Die Beugungerscheinungen einer Kreisrunden Oeffnung und eines runden Schirmchens," *abh der K. Bayer Akad. D. Wissenschaft* **15**:229–328 (1884).
98. E. H. Linfoot and E. Wolf, "Phase Distribution Near Focus in an Aberration-Free Diffraction Image," *Proc. Phys. Soc. B* **69**:823–832 (1956).
99. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon, Oxford, 1993, pp. 435–449.
100. H. Struve, "Beitrag zur der Theorie der Diffraction an Fernrohren," *Ann. Physik. Chem.* **17**:1008–1016 (1882).
101. H. H. Hopkins, "The Frequency Response of a Defocused Optical System," *Proc. Roy. Soc. London A* **231**:91–103 (1955).
102. L. Levi, *Handbook of Tables for Applied Optics*, CRC Press, Cleveland, 1974.
103. W. N. Charman, "Effect of Refractive Error in Visual Tests with Sinusoidal Gratings," *Brit. J. Physiol. Opt.* **33**(2):10–20 (1979).
104. G. E. Legge, K. T. Mullen, G. C. Woo, and F. W. Campbell, "Tolerance to Visual Defocus," *J. Opt. Soc. Am A* **4**:851–863 (1987).
105. G. Westheimer, "Pupil size and Visual Resolution," *Vision Res.* **4**:39–45 (1964).
106. F. W. Campbell and D. G. Green, "Optical and Retinal Factors Affecting Visual Resolution," *J. Physiol. (London)* **181**:576–593 (1965).
107. F. W. Campbell and R. W. Gubisch, "The effect of Chromatic Aberration on Visual Acuity," *J. Physiol. (London)* **192**:345–358 (1967).
108. W. N. Charman and J. Tucker, "Dependence of Accommodation Response on the Spatial Frequency Spectrum of the Observed Object," *Vision Res.* **17**:129–139 (1977).
109. W. N. Charman and J. Tucker, "Accommodation as a Function of Object Form," *Am. J. Optom. Physiol. Opt.* **55**:84–92 (1978).
110. W. N. Charman and J. A. M. Jennings, "The Optical Quality of the Retinal Image as a Function of Focus" *Br. J. Physiol. Opt.* **31**:119–134 (1976).
111. W. N. Charman and G. Heron, "Spatial Frequency and the Dynamics of the Accommodation Response," *Optica Acta* **26**:217–228 (1979).
112. H. H. Hopkins, "Geometrical Optical treatment of Frequency Response," *Proc. Phys. Soc. B* **70**:1162–1172 (1957).
113. G. A. Fry, "Blur of the Retinal Image" *Progress in Optics*, vol. 8, E. Wolf (ed.), North Holland, Amsterdam, 1970, pp. 53–131.
114. G. Smith, "Ocular defocus, Spurious Resolution and Contrast Reversal," *Ophthal. Physiol. Opt.* **2**:5–23 (1982).
115. C. Chan, G. Smith, and R. J. Jacobs, "Simulating Refractive Errors: Source and Observer Methods," *Am. J. Optom. Physiol. Opt.* **62**:207–216 (1985).
116. S. Duke-Elder and D. Abrams, *System of Ophthalmology*, vol. V, *Ophthalmic Optics and Refraction*, Kimpton, London, 1970, pp. 134–139.
117. R. B. Rabbetts, *Clinical Visual Optics*, 3rd ed., Butterworth-Heinemann, Oxford, 1998, pp. 220–221.
118. N. Thibos and A. Bradley, "Modelling the Refractive and Neurosensor Systems of the Eye," P. Mouroullis, (ed.), *Visual Instrumentation: Optical design and Engineering Principles*, McGraw Hill, New York, 1999, pp. 101–159.
119. R. R. Krueger, R. A. Applegate, and S. M. Macrae, (eds.), *Wavefront Customized Visual Correction: the Quest for Supervision II*, Slack Inc., Thorofare, (2004).

120. D. A. Atchison, "Recent Advances in the Measurement of Monochromatic Aberrations in Human Eyes," *Clin. Exper. Optom.* **88**:5–26 (2005).
121. D. A. Atchison, "Recent Advances in the Representation of Monochromatic Aberrations of Human Eyes," *Clin. Exper. Optom.* **87**:138–148 (2004).
122. W. N. Charman, "Wavefront Technology: Past, Present and Future," *Contact Lens Ant. Eye* **28**:75–92 (2005).
123. L. N. Thibos, R. A. Applegate, J. T. Schwiegerling, and R. Webb, "Standards for Reporting the Optical Aberrations of Eyes," *J. Refract. Surg.* **18**:S652–S660 (2002).
124. ANSI, American Standards Institute, *American National Standards for Ophthalmics—Methods for Reporting Optical Aberrations of Eyes*, ANSI Z. 80. 28–2004 (2004).
125. M. K. Smolek and S. D. Klyce, "Zernike Polynomial Fitting Fails to Represent All Visually Significant Corneal Aberrations," *Invest. Ophthalmol. Vis. Sci.* **44**:4676–4681 (2003).
126. S. D. Klyce, M. D. Karon, and M. K. Smolek, "Advantages and Disadvantages of the Zernike Expansion for Representing Wave Aberration of the Normal and Aberrated Eye," *J. Refract. Surg.* **20**:S537–541 (2004).
127. R. Iskander, M. J. Collins, B. Davis, and L. G. Carney, "Monochromatic Aberrations and Characteristics of Retinal Image Quality," *Clin. Exp. Optom.* **83**:315–322 (2000).
128. J. Porter, A. Guirao, I. G. Cox, and D. R. Williams, "The Human Eye's Monochromatic Aberrations in a Large Population," *J. Opt. Soc. Am. A* **18**:1793–1803 (2001).
129. L. N. Thibos, X. Hong, A. Bradley, and X. Cheng, "Statistical Variation of Aberration Structure and Image Quality in a Normal Population of Healthy Eyes," *J. Opt. Soc. Am. A* **19**:2329–2348 (2002).
130. J. F. Castéjon-Mochón, N. Lopez-Gil, A. Benito, and P. Artal, "Ocular Wavefront Statistics in a Normal Young Population," *Vision Res.* **42**:1611–1617 (2002).
131. L. Wang and D. D. Koch, "Ocular Higher-Order Aberrations in Individuals Screened for Refractive Surgery," *J. Cataract Refract. Surg.* **29**:1896–1903 (2003).
132. N. V. Netto, R. Abrosio, T. T. Shen, and S. E. Wilson, "Wavefront Analysis in Normal Refractive Surgery Candidates," *J. Refract. Surg.* **21**:332–338 (2005).
133. T. O. Salmon and C. van de Pol, "Normal-eye Zernike Coefficients and Root-Mean-Square Wavefront Errors," *J. Cataract Refract. Surg.* **32**:2064–2074 (2006).
134. R. A. Applegate, W. J. Donnelly, J. D. Marsack, and D. E. Koenig, "Three-Dimensional Relationship between Higher-Order Root-Mean-Square Wavefront Error, Pupil Diameter and Aging," *J. Opt. Soc. Am. A* **24**:578–587 (2007).
135. R. I. Calver, M. J. Cox, and D. B. Elliott, "Effects of Aging on the Monochromatic Aberrations of the Human Eye," *J. Opt. Soc. Am. A* **16**:2069–2078 (1999).
136. J. S. McLellan, S. Marcos, and S. A. Burns, "Age-Related Changes in Monochromatic Wave Aberrations of the Human Eye," *Invest. Ophthalmol. Vis. Sci.* **42**:1390–1395 (2001).
137. T. Kuroda, T. Fujikado, S. Ninomiya, N. Maeda, Y. Hirohara, and T. Mihashi, "Effect of Aging on Ocular Light Scatter and Higher order Aberration," *J. Refract. Surg.* **18**:S598–S602 (2002).
138. I. Brunette, J. M. Bueno, M. Parent, H. Hamam, and P. Simonet, "Monochromatic Aberrations as a Function of Age, from Childhood to Advanced Age," *Invest. Ophthalmol. Vis. Sci.* **44**:5438–5446 (2003).
139. D. Atchison, M. J. Collins, C. F. Wildsoet, J. Christensen, and M. D. Waterworth, "Measurement of Monochromatic Ocular Aberrations of Human Eyes as a Function of Accommodation by the Howland Aberroscope Technique," *Vision Res.* **35**:313–323 (1995).
140. S. Ninomiya, T. Fujikado, T. Kuroda, N. Maeda, Y. Tano, T. Oshika, Y. Hirohara, and T. Mihashi, "Changes in Ocular Aberration with Accommodation," *Am. J. Ophthalmol.* **134**:924–926 (2002).
141. J. C. He, S. A. Burns, and S. Marcos, "Monochromatic Aberrations in the Accommodated Human Eye," *Vision Res.* **40**:41–48 (2000).
142. C. A. Hazel, M. J. Cox, and N. C. Strang, "Wavefront Aberration and Its Relationship to the Accommodative Stimulus-Response Function in Myopic Subjects," *Optom. Vis. Sci.* **80**:151–158 (2003).
143. S. Plainis, H. S. Ginis, and A. Pallikaris, "The Effect of Ocular Aberrations on Steady-State Errors of Accommodative Response," *J. Vision* **5**(5):466–477 (2005).
144. H. Cheng, J. K. Barnett, A. S. Vilupuru, J. D. Marsack, S. Kasthurirangan, R. A. Applegate, and A. Roorda, "A Population Study on Changes in Wavefront Aberration with Accommodation," *J. Vision* **4**(4):272–280 (2004).

145. H. Radhakrishnan and W. N. Charman, "Age-Related Changes in Ocular Aberrations with Accommodation," *J. Vision* **7**(7):1–21 (2007).
146. M. P. Paquin, H. Hamam, P. Simonet, "Objective Measurement of Optical Aberrations in Myopic Eyes," *Optom. Vis. Sci.* **79**:285–291 (2002).
147. X. Cheng, A. Bradley, X. Hong, and L. N. Thibos, "Relationship between Refractive Error and Monochromatic Aberrations of the Eye," *Optom. Vis. Sci.* **80**:43–49 (2003).
148. H. Hofer, P. Artal, B. Singer, J. Laragon, and D. R. Williams, "Dynamics of the Eye's Aberration," *J. Opt. Soc. Am. A* **18**:497–506 (2001).
149. L. Diaz-Santana, C. Torti, I. Munro, P. Gasson, and C. Dainty, "Benefit of Higher Closed-Loop Bandwidth in Ocular Adaptive Optics," *Opt. Express* **11**:2597–2605 (2003)
150. T. Nirmaier, G. Pudasaini, and J. Bille, "Very Fast Wavefront Measurements at the Human Eye with a Custom CMOS-based Hartmann-Shack Sensor," *Opt. Express* **11**:2704–2716. (2003).
151. M. Zhu, M. J. Collins, and R. Iskander, "Microfluctuations of Wavefront Aberration of the Eye," *Ophthal. Physiol. Opt.* **24**:562–571 (2004).
152. K. M. Hampson, I. Munro, C. Paterson, and C. Dainty, "Weak Correlation between the Aberration Dynamics of the Human Eye and the Cardiopulmonary System," *J. Opt. Soc. Am. A* **22**:1241–1250 (2005).
153. A. Maréchal, "Etude des Effets Combinés de la Diffraction et des Aberrations Géométriques sur l'image d'un Point Lumineux," *Rev. d'Optique* **26**:257–277 (1947).
154. J. Perrigin, D. Perrigin, and T. Grosvenor, "A Comparison of Clinical Refractive Data Obtained by Three Examiners," *Am. J. Optom. Physiol. Opt.* **59**:515–519 (1982).
155. M. A. Bullimore, R. E. Fusaro, and C. W. Adams, "The Repeatability of Automated and Clinical Refraction," *Optom. Vis. Sci.* **75**:617–622 (1998).
156. R. A. Applegate, C. Ballentine, H. Gross, E. J. Sarver, and C. A. Sarver, "Visual Acuity as a Function of Zernike Mode and Level of Root Mean Square Error," *Optom. Vis. Sci.* **80**:97–105 (2003).
157. R. A. Applegate, J. D. Marsack, R. Ramos, and E. J. Sarver, "Interaction between Aberrations to Improve or Reduce Visual Performance," *J. Cataract Refract. Surg.* **29**:1487–1495 (2003).
158. S. G. El Hage and F. Berny, "Contribution of the Crystalline Lens to the Spherical Aberration of the Eye," *J. Opt. Soc. Am.* **63**:205–211 (1973).
159. P. Artal and A. Guirao, "Contribution of the Cornea and Lens to the Aberrations of the Human Eye," *Opt. Lett.* **23**:1713–1715 (1998).
160. P. Artal, A. Guirao, E. Berrio, and D. R. Williams, "Compensation of Corneal Aberrations by the Internal Optics in the Human Eye," *J. Vision* **1**:1–8 (2001).
161. M. Mrochen, M. Jankov, M. Bueeler, and T. Seiler, "Correlation between Corneal and Total Wavefront Aberrations in Myopic Eyes," *J. Refract. Surg.* **19**:104–112 (2003).
162. J. E. Kelly, T. Mihashi, and H. C. Howland, "Compensation of Cornea Horizontal/Vertical Astigmatism, Lateral Coma and Spherical Aberration by the Internal Optics of the Eye," *J. Vision* **4**:262–271 (2004).
163. J. C. He, J. Gwiazda, F. Thorn, and R. Held, "Wave-front Aberrations in the Anterior Corneal Surface and the Whole Eye," *J. Opt. Soc. Am. A* **20**:1155–1163 (2003).
164. P. Artal, E. Berrio, and A. Guirao, "Contribution of the Cornea and Internal Surfaces to the Change of Ocular Aberrations with Age," *J. Opt. Soc. Am. A* **19**:137–143 (2002).
165. P. Artal, A. Benito, J. Tabernero, "The Human Eye is an Example of Robust Optical Design," *J. Vision* **6**:1–7 (2006).
166. P. Artal, L. Chen, E. J. Fernández, B. Singer, S. Manzanera, and D. R. Williams, "Neural Compensation for the Eye's Optical Aberrations," *J. Vision* **4**:281–287 (2004).
167. W. N. Charman, "Aberrations and Myopia," *Ophthal. Physiol. Opt.* **25**:285–301 (2005).
168. L. Llorente, S. Barbero, D. Cano, C. Dorronsoro, and S. Marcos, "Myopic versus Hyperopic Eyes: Axial Length, Corneal Shape and Optical Aberrations," *J. Vision* **4**(4):288–298 (2004).
169. R. Montés-Mico, J. L. Alió, G. Muñoz, J. J. Perez-Santonja, and W. N. Charman, "Postblink Changes in Total and Corneal Optical Aberrations," *Ophthalmology* **111**:758–767, (2004).
170. K. Y. Li and G. Yoon, "Changes in Aberrations and Retinal Image Quality due to Tear Film Dynamics," *Optics Express* **14**:12552–12559 (2006).

171. T. Buehren, M. J. Collins, and L. Carney, "Corneal Aberrations and Reading," *Optom. Vis. Sci.* **80**:159–166 (2003).
172. T. Buehren, M. J. Collins, and L. Carney, "Near Work Induced Wavefront Aberrations in Myopia," *Vision Res.* **45**:1297–1312 (2005).
173. W. Han, W. Kwan, J. Wang, S. P. Yip, and M. Yap, "Influence of Eyelid Position on Wavefront Aberration," *Ophthalm. Physiol. Opt.* **27**:66–75 (2007).
174. C. E. Ferree, G. Rand, and C. Hardy, "Refraction for the Peripheral Field of Vision," *Arch. Ophthalmol. (Chicago)* **5**:717–731 (1931).
175. T. C. A. Jenkins, "Aberrations of the Eye and Their Effects on Vision, Parts I and II" *Br. J. Physiol. Opt.* **20**:50–91 and 161–201 (1963).
176. F. Rempt, J. F. Hoogerheide, and W. P. H. Hoogenboom, "Peripheral Retinoscopy and the Skiagram," *Ophthalmologica* **162**:1–10 (1971).
177. J. F. Hoogerheide, F. Rempt, and W. P. H. Hoogenboom, "Acquired Myopia in Young Pilots." *Ophthalmologica* **163**:209–215 (1971).
178. M. Millodot, "Effect of Ametropia on Peripheral Refraction," *Am. J. Physiol. Opt.* **58**:691–695 (1981).
179. M. Millodot, "Peripheral Refraction in Aphakic Eyes," *Am. J. Optom. Physiol. Opt.* **61**:586–589 (1984).
180. D. R. Williams, P. Artal, R. Navarro, M. J. McMahon, and D. H. Brainard, "Off-Axis Optical Quality and Retinal Sampling in the Human Eye," *Vision Res.* **36**:1103–1114 (1996).
181. G. Smith, M. Millodot, and N. McBrien, "The Effect of Accommodation on Oblique Astigmatism and Field Curvature of the Human Eye," *Clin. Exp. Optom.* **71**:119–125 (1988).
182. A. Guirao and P. Artal, "Off-Axis Monochromatic Aberrations Estimated from Double-Pass Measurements in the Human Eye," *Vision Res.* **39**:207–217 (1999).
183. D. O. Mutti, R. I. Scholtz, N. E. Friedman, and K. Zadnik, "Peripheral Refraction and Ocular Shape in Children," *Invest. Ophthalmol. Vis. Sci.* **41**:1022–1030 (2000).
184. J. Gustafsson, E. Terenius, J. Buckheister, and P. Unsbo, "Peripheral Astigmatism in Emmetropic Eyes," *Ophthalm. Physiol. Opt.* **21**:393–400 and 491 (2001).
185. A. Seidemann, F. Schaefel, A. Guirao, N. Lopez-Gil, and P. Artal, "Peripheral Refractive Errors in Myopic, Emmetropic, and Hyperopic Young Subjects," *J. Opt. Soc. Am. A* **19**:2363–2373 (2002).
186. D. A. Atchison, N. Pritchard, S. D. White, and A. M. Griffiths, "Influence of Age on Peripheral Refraction," *Vision Res.* **45**:715–720 (2005).
187. W. N. Charman and J. A. M. Jennings, "Longitudinal Changes in Peripheral Refraction with Age," *Ophthalm. Physiol. Opt.* **26**:447–455 (2006).
188. R. Navarro, E. Moreno, and C. Dorransoro, "Monochromatic Aberrations and Point-Spread Functions of the Human Eye across the Visual Field," *J. Opt. Soc. Am. A* **15**:2522–2529 (1998).
189. D. A. Atchison and D. H. Scott, "Monochromatic Aberrations of Human Eyes in the Horizontal Visual Field," *J. Opt. Soc. Am. A* **19**:2180–2184 (2002).
190. D. A. Atchison, S. D. Lucas, R. Ashman, and M. A. Huynh, "Refraction and Aberration across the Horizontal Central 10° of the Visual Field," *Optom. Vis. Sci.* **83**:213–221 (2006).
191. D. A. Atchison, "Higher Order Aberrations Across the Horizontal Visual Field," *J. Biomed Optics* **11**(3): article 034026 (2006).
192. D. A. Atchison, "Anterior Corneal and Internal Contributions to Peripheral Aberrations of Human Eyes," *J. Opt. Soc. Am. A* **21**:355–359 (2004).
193. D. A. Atchison, D. H. Scott, and W. N. Charman, "Hartmann-Shack Technique and Refraction Across the Horizontal Visual Field," *J. Opt. Soc. Am. A* **20**:965–973 (2003).
194. L. Lundström and P. Unsbo, "Transformation of Zernike Coefficients:Scaled, Translated, and Rotated Wavefronts with Circular and Elliptical Pupils," *J. Opt. Soc. Am. A* **24**:569–577 (2007).
195. D. A. Atchison, D. H. Scott, and W. N. Charman, "Measuring Ocular Aberrations in the Peripheral Field Using Hartmann-Shack Aberrometry," *J. Opt. Soc. Am. A* **24**:2963–2973 (2007).
196. Y. Le Grand, *Form and Space Vision*, M. Millodot and G. C. Heath (trans.), Indiana UP, Bloomington, 1967, pp. 5–23.
197. J. G. Sivak and T. Mandelman, "Chromatic Dispersion of the Ocular Media," *Vision Res.* **22**:997–1003 (1982).

198. T. Mandelman and J. G. Sivak, "Longitudinal Chromatic Aberration of the Vertebrate Eye," *Vision Res.* **23**:1555–1559 (1983).
199. D. A. Atchison and G. Smith, "Chromatic Dispersions of the Ocular Media of Human Eyes," *J. Opt. Soc. Am. A* **22**:29–37 (2005).
200. L. N. Thibos, A. Bradley, F. D. L. Still, X. Zhang, and P. A. Howarth, "Theory and Measurement of Ocular Chromatic Aberration," *Vision Res.* **30**:33–49 (1990).
201. L. N. Thibos, A. Bradley, and X. Zhang, "Effect of Ocular Chromatic Aberration on Monocular Visual Performance," *Optom. Vis. Sci.* **68**:599–607 (1991).
202. G. Wald and D. R. Griffin, "The Change of Refractive Power of the Human Eye in Dim and Bright Light," *J. Opt. Soc. Am.* **37**:321–336 (1947).
203. A. Ivanoff, *Les Aberrations de l'Oeil*, Revue d'Optique, Paris, 1953.
204. R. E. Bedford and G. Wyszecki, "Axial Chromatic Aberration of the Human Eye," *J. Opt. Soc. Am.* **47**:564–565 (1957).
205. A. Howarth and A. Bradley, "The Longitudinal Chromatic Aberration of the Human Eye, and Its Correction," *Vision Res.* **26**:361–366 (1986).
206. H. Hartridge, "The Visual Perception of Fine Detail," *Phil. Trans. R. Soc. A* **232**:519–671 (1947).
207. J. S. McLellan, S. Marcos, P. M. Prieto, and S. A. Burns, "Imperfect Optics may be the Eye's Defence Against Chromatic Blur," *Nature* **417**:174–176 (2002).
208. P. A. Howarth, "The Lateral Chromatic Aberration of the Human Eye," *Ophthal. Physiol. Opt.* **4**:223–236 (1984).
209. A. van Meeteren, "Calculations on the Optical Modulation Transfer Function of the Human Eye for White Light," *Optica Acta* **21**:395–412 (1974).
210. L. N. Thibos, "Calculation of the Influence of Lateral Chromatic Aberration on Image Quality across the Visual Field," *J. Opt. Soc. Am. A* **4**:1673–1680 (1987).
211. P. Simonet and M. C. W. Campbell, "The Optical Transverse Chromatic Aberration on the Fovea of the Human Eye," *Vision Res.* **30**:187–206 (1990).
212. J. J. Vos, "Some New Aspects of Color Stereoscopy," *J. Opt. Soc. Am.* **50**:785–790 (1960).
213. B. N. Kishto, "The Colour Stereoscopic Effect," *Vision Res.* **5**:3131–329 (1965).
214. J. M. Sundet, "The Effect of Pupil Size Variation on the Colour Stereoscopic Phenomenon," *Vision Res.* **12**:1027–1032 (1972).
215. R. C. Allen and M. L. Rubin, "Chromostereopsis," *Survey Ophthalmol.* **26**:22–27 (1981).
216. D. A. Owens and H. W. Leibowitz, "Chromostereopsis and Small Pupils," *J. Opt. Soc. Am.* **65**:358–359 (1975).
217. M. Ye, A. Bradley, L. N. Thibos, and X. X. Zhang, "Interocular Differences in Transverse Chromatic Aberration Determine Chromostereopsis for Small Pupils," *Vision Res.* **31**:1787–1796 (1991).
218. J. Faubert, "Seeing Depth in Colour: More Than Just What Meets the Eyes," *Vision Res.* **34**:1165–1186 (1994).
219. L. N. Thibos, F. E. Cheney, and D. J. Walsh, "Retinal Limits to the Detection and Resolution of Gratings," *J. Opt. Soc. Am. A* **4**:1524–1529 (1987).
220. Y. U. Ogbo and H. E. Bedell, "Magnitude of Lateral Chromatic Aberration Across the Retina of the Human Eye," *J. Opt. Soc. Am. A* **4**:1666–1672 (1987).
221. L. L. Holladay, "The Fundamentals of Glare and Visibility," *J. Opt. Soc. Am.* **12**:271–319 (1926).
222. W. S. Stiles, "The Effect of Glare on the Brightness Threshold," *Proc. R. Soc. Lond. B* **104**:322–351 (1929).
223. J. J. Vos, J. Walraven, and A. van Meeteren, "Light Profiles of the Foveal Image of a Point Source," *Vision Res.* **16**:215–219 (1976).
224. M. J. Allen and J. J. Vos, "Ocular Scattered Light and Visual Performance as a Function of Age," *Am. J. Optom. Arch. Am. Acad. Optom.* **44**:717–727 (1967).
225. A. Spector, S. Li, and J. Sigelman, "Age-Dependent Changes in the Molecular Size of Human Lens Proteins and Their Relationship to Light Scatter," *Invest. Ophthalmol.* **13**:795–798 (1974).
226. R. P. Hemenger, "Intraocular Light Scatter in Normal Visual Loss with Age," *Applied Opt.* **23**:1972–1974 (1984).

227. R. A. Weale, "Effects of Senescence," *Vision and Visual Dysfunction*, vol. 5, *Limits to Vision*, J. J. Kulikowski, V. Walsh, and J. Murray (eds.), Macmillan, London, 1991, pp. 277–285.
228. W. Adrian and A. Bhanji, "A Formula to Describe Straylight in the Eye as a Function of the Glare Angle and Age," *1st Int. Symp. on Glare*, Orlando Fla. , Oct 24 and 25, 1991, The Lighting Research Institute, New York, 1991, pp. 185–191.
229. J. J. Vos and J. Boogaard, "Contribution of the Cornea to Entoptic Scatter," *J. Opt. Soc. Am.* **53**:869–873 (1963).
230. R. M. Boynton and F. J. J. Clarke, "Sources of Entoptic Scatter in the Human Eye," *J. Opt. Soc. Am.* **54**:110–119, 717–719 (1964).
231. J. J. Vos, "Contribution of the Fundus Oculi to Entoptic Scatter," *J. Opt. Soc. Am.* **53**:1449–1451 (1963).
232. J. J. Vos and M. A. Bouman, "Contribution of the Retina to Entoptic Scatter," *J. Opt. Soc. Am.* **54**:95–100 (1964).
233. R. P. Hemenger, "Small-angle Intraocular Scattered Light: A Hypothesis Concerning its Source," *J. Opt. Soc. Am. A* **5**:577–582 (1988).
234. J. H. Ohzu and J. M. Enoch, "Optical Modulation by the Isolated Human Fovea," *Vision Res.* **12**:245–251 (1972).
235. D. A. Williams, "Visibility of Interference Fringes Near the Resolution Limit," *J. Opt. Soc. Am. A* **2**:1087–1093.
236. D. I. A. MacLeod, D. R. Williams, and W. Makou, "A Visual Non-linearity Fed by Single Cones," *Vision Res.* **32**:347–363 (1992).
237. D. A. Palmer, "Entopic Phenomena," *Vision and Visual Dysfunction*, vol. 1, *Visual Optics and Instrumentation*, W. N. Charman (ed.), Macmillan, London, 1991, pp. 345–370.
238. R. B. Rabbetts, "Entoptic Phenomena," *Clinical Visual Optics*, 3rd ed., Butterworth-Heinemann, Oxford, 1998, pp. 421–429.
239. A. Caldecott and W. N. Charman, "Diffraction Haloes Resulting from Corneal Oedema and Epithelial Cell Size," *Ophthalm. Physiol. Opt.* **22**:209–213 (2002).
240. T. J. T. P. van den Berg, M. P. J. Hagenouw, and J. E. Coppens, "The Ciliary Corona: Physical Model and Simulation of the Fine Needles Radiating from Points Light Sources," *Invest. Ophthalmol. Vis. Sci.* **46**:2627–2632 (2005).
241. R. Navarro and M. A. Losada, "Shape of Stars and Optical Quality of the Human Eye," *J. Opt. Soc. Am. A* **14**:353–359 (1997).
242. D. B. Henson, *Optometric Instrumentation*, 2nd ed., Butterworth-Heinemann, Oxford, 1996.
243. S. Siik, P. J. Airaksinen, A. Tuulonen, and H. Nieminen, "Autofluorescence in Cataractous Human Lens and its Relationship to Light Scatter," *Acta Ophthalmol.* **71**:388–392 (1993).
244. D. B. Elliott, K. C. H. Yang, K. Dumbleton, and A. P. Cullen, "Ultra-Violet Induced Lenticular Fluorescence: Intraocular Straylight Affecting Visual Function," *Vision Res.* **33**:1827–1833 (1993).
245. T. C. D. Whiteside, *Problems of Vision in Flight at High Altitudes*, Butterworths, London, 1957.
246. H. H. Hopkins, "The Application of Frequency Response Techniques in Optics," *Proc. Phys. Soc.* **79**:889–919 (1962).
247. H. Metcalf, "Stiles-Crawford Apodization," *J. Opt. Soc. Am.* **55**:72–74 (1965).
248. D. A. Atchison, A. Joblin, and G. Smith, "Influence of Stiles-Crawford Apodization on Spatial Visual Performance," *J. Opt. Soc. Am. A* **15**:2545–2550 (1998).
249. D. A. Atchison, D. H. Scott, N. C. Strang, and P. Artal, "Influence of Stiles-Crawford Apodization on Visual Acuity," *J. Opt. Soc. Am. A* **19**:1073–1083 (2002).
250. X. Zhang, A. Bradley, and L. N. Thibos, "Apodization by the Stiles-Crawford Effect Moderates the Visual Impact of Retinal Image Defocus," *J. Opt. Soc. Am. A* **16**:812–820 (1999).
251. M. J. Cox, D. A. Atchison, and D. H. Scott, "Scatter and Its Implications for the Measurement of Optical Image Quality in Human Eyes," *Optom. Vis. Sci.* **80**:56–68 (2003).
252. G. Walsh, W. N. Charman, and H. C. Howland, "Objective Technique for the Determination of Monochromatic Aberrations of the Human Eye," *J. Opt. Soc. Am. A* **1**:987–992 (1984).
253. L. Diaz-Santana, G. Walker, and S. X. Bara, "Sampling Geometries for Ocular Aberrometry: a Model for Evaluation of Performance," *Opt. Express* **13**:8801–8818 (2005).

254. L. Llorente, S. Marcos, C. Dorronsoro, and S. A. Burns, "Effect of Sampling on Real Ocular Aberration Measurements," *J. Opt. Soc. Am. A* **24**:2783–2796 (2007).
255. X. Hong, L. N. Thibos, A. Bradley, R. L. Woods, and R. A. Applegate, "Comparison of Monochromatic Ocular Aberrations Measured with an Objective Crossed-Cylinder Aberroscope and a Shack-Hartmann Aberrometer," *Optom. Vis. Sci.* **80**:15–25 (2003).
256. Y. Le Grand, "Sur une Mode de Vision Eliminant les Défauts Optiques de l'Oeil," *Rev. d'Optique Theor. Instrum.* **15**:6–11 (1936).
257. W. N. Charman and P. Simonet, "Yves Le Grand and the Assessment of Retinal Acuity Using Interference Fringes," *Ophthalm. Physiol. Opt.* **17**:164–168 (1997).
258. G. Westheimer, "Modulation Thresholds for Sinusoidal Light Distributions on the Retina," *J. Physiol. (London)* **152**:67–74 (1960).
259. A. Arnulf and O. Dupuy, "La Transmission des Contrastes par la Système Optique de l'Oeil et les Seuls de Modulation Réiniens," *C. r. hebdom. Seanc. Acad. Paris* **250**:2757–2759 (1960).
260. F. W. Campbell and D. G. Green, "Optical and Retinal Factors Affecting Visual Resolution," *J. Physiol. (London)* **181**:576–593 (1965).
261. S. Berger-Lheureux, "Mesure de la Fonction de Transfert de Modulation du Système Optique de l'Oeil et des Seuls de Modulation Réiniens," *Rev. Opt. Theor. Instrum.* **44**:294–323 (1965).
262. L. J. Bour, "MTF of the Defocused Optical System of the Human Eye for Incoherent Monochromatic Light," *J. Opt. Soc. Am.* **70**:321–328 (1980).
263. N. Sekiguchi, D. R. Williams, and D. H. Brainard, "Aberration-Free Measurement of the Visibility of Isoluminant Gratings," *J. Opt. Soc. Am. A* **10**:2105–2116 (1993).
264. N. Sekiguchi, D. R. Williams, and D. H. Brainard, "Efficiency in Detection of Isoluminant and Isochromatic Interference Fringes," *J. Opt. Soc. Am. A* **10**:2118–2133 (1993).
265. J. Rovamo, J. Mustonen, and R. Näsänen, "Two Simple Methods for Determining the Optical Transfer Function of the Human Eye," *Vision Res.* **34**:2493–2502 (1994).
266. F. W. Campbell and R. W. Gubisch, "Optical Quality of the Human Eye," *J. Physiol. (London)* **186**:5589–578 (1966).
267. F. Flamant, "Etude de la Repartition de la Lumière dans l'Image Réinienne d'une Fente," *Rev. Opt. Theor. Instrum.* **34**:433–459 (1955).
268. J. Krauskopf, "Light distribution in Human Retinal Images," *J. Opt. Soc. Am.* **52**:1046–1050 (1962).
269. J. Krauskopf, "Further Measurements in Human Retinal Images," *J. Opt. Soc. Am.* **54**:715–716 (1964).
270. G. Westheimer and F. W. Campbell, "Light Distribution in the Image Formed by the Living Human Eye," *J. Opt. Soc. Am.* **52**:1040–1045 (1962).
271. R. Rohler, U. Miller, and M. Aberl, "Zür Messung der Modulationsübertragungsfunktion des lebenden menschlichen Augen in reflektieren Licht," *Vision Res.* **9**:407–428 (1969).
272. J. A. M. Jennings and W. N. Charman, "Off-Axis Image Quality in the Human Eye," *Vision Res.* **21**:445–455 (1981).
273. J. Santamaria, P. Artal, and J. Bescos, "Determination of the Point Spread Function of Human Eyes Using a Hybrid Optical-Digital Method," *J. Opt. Soc. Am. A* **4**:109–114 (1987).
274. P. Artal, S. Marcos, R. Navarro, and D. Williams, "Odd Aberrations and Double Pass Measurements of Retinal Image Quality," *J. Opt. Soc. Am. A* **12**:195–201 (1995).
275. R. Navarro and M. A. Losada, "Phase Transfer and Point-Spread Function of the Human Eye Determined by a New Asymmetric Double-pass Method," *J. Opt. Soc. Am. A* **12**:2385–2392 (1995).
276. P. Artal, I. Iglesias, N. López-Gil, and D. G. Green, "Double-Pass Measurements of the Retinal Image Quality with Unequal Entrance and Exit Pupil Size and the Reversibility of the Eye's Optical System," *J. Opt. Soc. Am. A* **12**:2358–2366 (1995).
277. L. Diaz-Santana and J. C. Dainty, "Effects of Retinal Scattering in the Ocular Double-Pass Procedure," *J. Opt. Soc. Am. A* **18**:1437–1444 (2001).
278. P. Artal, M. Ferro, I. Miranda, and R. Navarro, "Effects of Aging in Retinal Image Quality," *J. Opt. Soc. Am. A* **10**:1656–1662 (1993).
279. N. López-Gil, I. Iglesias, and P. Artal, "Retinal Image Quality in the Human Eye as Function of the Accommodation," *Vision Res.* **38**:2897–2907 (1998).
280. G. Westheimer and J. Liang, "Evaluating Diffusion of Light in the Eye by Objective Means," *Invest. Ophthalmol. Vis. Sci.* **35**:2652–2657 (1994).

281. J. F. Simon and P. M. Denieul, "Influence of the Size of Test Field Employed in Measurements of Modulation Transfer Function of the Eye," *J. Opt. Soc. Am.* **63**:894–896 (1973).
282. P. Artal and R. Navarro, "Simultaneous Measurement of Two-point-spread Functions at Different Locations across the Human Fovea," *Applied Opt.* **31**:3646–3656 (1992).
283. S. C. Choi, N. Doble, J. Lion, J. Christou, and D. R. Williams, "Effect of Wavelength on *in vivo* Images of the Human Cone Mosaic," *J. Opt. Soc. Am. A* **22**:2598–2605 (2005).
284. D. R. Williams, D. H. Brainard, M. J. McMahon, and R. Navarro, "Double-Pass and Interferometric Measures of the Optical Quality of the Eye," *J. Opt. Soc. Am. A* **11**:3123–3135 (1994).
285. N. López-Gil and P. Artal, "Comparison of Double-pass Estimates of the Retinal-image Quality Obtained with Green and Near-Infrared Light," *J. Opt. Soc. Am. A* **14**:961–971 (1997).
286. R. W. Gubisch, "Optical Performance of the Human Eye," *J. Opt. Soc. Am.* **57**:407–415 (1967).
287. P. Artal, "Calculations of Two-dimensional Foveal Retinal Images in Real Eyes," *J. Opt. Soc. Am. A* **7**:1374–1381 (1990).
288. W. T. Welford, *Aberrations of Optical Systems*, Adam Hilger, Bristol, 1986, pp. 243–244.
289. J. Liang and D. R. Williams, "Aberrations and Retinal Image Quality in the Normal Human Eye," *J. Opt. Soc. Am. A* **14**:2873–2883 (1997).
290. J. M. Bueno and P. Artal, "Polarization and Retinal Image Quality Estimates in the Human Eye," *J. Opt. Soc. Am. A* **18**:489–496 (2001).
291. A. C. S. van Heel, "Correcting the Spherical and Chromatic Aberrations of the Eye," *J. Opt. Soc. Am.* **36**:237–239 (1946).
292. H. Hartridge, "Visual Acuity and Resolving Power of the Eye," *J. Physiol. (London)* **57**:57–62 (1962).
293. S. A. Klein, "Optimal Corneal Ablation for Eyes with Arbitrary Hartmann-Shack Aberrations," *J. Opt. Soc. Am. A* **15**:2580–2588 (1998).
294. J. Schweigerling and R. W. Snyder, "Custom Photorefractive Keratectomy Ablations for the Correction of Spherical and Cylindrical Refractive Error and Higher-Order Aberrations," *J. Opt. Soc. Am. A* **15**:2572–2579 (1998).
295. J. Schweigerling, "Theoretical Limits to Visual Performance," *Surv. Ophthalmol.* **45**:139–146 (2000).
296. T. Seiler, M. Mrochen, and M. Kaemmerer, "Operative Correction of Ocular Aberrations to Improve Visual Acuity," *J. Refract. Surg.* **16**:S619–S622 (2000).
297. S. M. Macrae, R. R. Krueger, and R. A. Applegate (eds.), *Customized Corneal Ablation: the Quest for Supervision*, Slack Inc., Thorofare, N.J., 2001.
298. R. R. Krueger, R. A. Applegate, S. M. MacRae (eds.), *Wavefront Customized Visual Correction: the Quest for Supervision II*, Slack Inc., Thorofare, N.J., 2004.
299. J. Liang, D. R. Williams, and D. T. Miller, "Supernormal Vision and High-Resolution Retinal Imaging through Adaptive Optics," *J. Opt. Soc. Am. A* **14**:2884–2892 (1997).
300. G-Y. Yoon and D. R. Williams, "Visual Performance after Correcting Monochromatic and Chromatic Aberrations of the Eye," *J. Opt. Soc. Am. A* **19**:286–275, (2002).
301. D. C. Chen, S. M. Jones, D. A. Silva, and S. S. Olivier, "High-Resolution Adaptive Optics Scanning Laser Ophthalmoscope with Dual Deformable Mirrors," *J. Opt. Soc. Am. A* **24**:1305–1312 (2007).
302. N. Lopez-Gil, J. F. Castejon-Mochon, A. Benito, J. M. Marin, G. Loe-a-Foe, G. Marin, B. Fermiger, D. Renard, D. Joyeux, N. Chateau, and P. Artal, "Aberration Generation by Contact Lenses with Aspheric and Asymmetric Surfaces," *J. Refract. Surg.* **18**:603–609 (2002).
303. D. A. Chernyak, and C. E. Campbell, "System for the Design, Manufacture, and Testing of Custom Lenses with Known Amounts of High-Order Aberrations," *J. Opt. Soc. Am. A* **20**:2016–2021 (2002).
304. S. Bara, T. Mancebo, and E. Moreno-Barriuso, "Positioning Tolerances for Phase Plates Compensating Aberrations of the Human Eye," *Applied Opt.* **39**:3413–3420 (2000).
305. A. Guirao, D. R. Williams, and I. G. Cox, "Effect of Rotation and Translation on the Expected Benefit of an Ideal Method to Correct the Eye's Higher-order Aberrations," *J. Opt. Soc. Am. A* **18**:1003–1015 (2001).
306. G. Yoon and T. M. Jeong, "Effect of the Movement of Customized Contact Lens on Benefit in Abnormal Eyes," *J. Vision* **3**(12):38a (2003).
307. J. De Brabander, N. Chateau, G. Marin, N. Lopez-Gil, E. van der Worp, and A. Benito, "Simulated Optical Performance of Custom Wavefront Soft Contact Lenses for Keratoconus," *Optom. Vis. Sci.* **80**:637–643 (2003).

308. N. Lopez-Gil, N. Chateau, J. Castejon-Monchon, and A. Benito, "Correcting Ocular Aberrations by Soft Contact Lenses," *S. Afr. Optom.* **62**:173–177 (2003).
309. W. N. Charman and N. Chateau, "The Prospects for Super-Acuity: Limits to Visual Performance after Correction of Monochromatic Ocular Aberration," *Ophthalm. Physiol. Opt.* **23**:479–493 (2003).
310. W. N. Charman, "Ablation Design in Relation to Spatial Frequency, Depth-of-focus and Age," *J. Refract. Surg.* **20**:S572–S579 (2004).
311. L. C. Thomson and W. D. Wright, "The Colour Sensitivity of the Retina Within the Central Fovea of Man," *J. Physiol. (London)* **105**:316–331 (1947).
312. I. Powell, "Lenses for Correcting Chromatic Aberration of the Eye," *Applied Opt.* **20**:4152–4155 (1981).
313. A. L. Lewis, M. Katz, and C. Oehrlein, "A Modified Achromatizing Lens," *Am. J. Optom. Physiol. Opt.* **59**:909–911 (1982).
314. J. A. Díaz, M. Irlbauer, and J. A. Martínez, "Diffractive-Refractive Hybrid Doublet to Achromatize the Human Eye," *J. Mod. Opt.* **51**:2223–2234 (2004).
315. Y. Benny, S. Manzanera, P. M. Prieto, E. N. Ribak, and P. Artal, "Wide-Angle Chromatic Aberration Corrector for the Human Eye," *J. Opt. Soc. Am. A* **24**:1538–1544 (2007).
316. X. Zhang, A. Bradley, and L. N. Thibos, "Achromatizing the Human Eye: the Problem of Chromatic Parallax," *J. Opt. Soc. Am. A* **8**:686–691 (1991).
317. J. A. M. Jennings and W. N. Charman, "Optical Image Quality in the Peripheral Retina," *Am. J. Optom. Physiol. Opt.* **55**:582–590 (1978).
318. R. Navarro, P. Artal, and D. R. Williams, "Modulation Transfer of the Human Eye as a Function of Retinal Eccentricity," *J. Opt. Soc. Am. A* **10**:201–212 (1993).
319. D. R. Williams, P. Artal, R. Navarro, M. J. McMahon, and D. H. Brainard, "Off-Axis Optical Quality and Retinal Sampling in the Human Eye," *Vision Res.* **36**:1103–1114 (1996).
320. J. A. M. Jennings and W. N. Charman, "Analytic Approximation of the Off-Axis Modulation Transfer Function of the Eye," *Vision Res.* **37**:697–704 (1997).
321. P. B. DeVelis and G. B. Parrent, "Transfer Function for Cascaded Optical Systems," *J. Opt. Soc. Am.* **57**:1486–1490 (1967).
322. I. Overington, "Interaction of Vision with Optical Aids," *J. Opt. Soc. Am.* **63**:1043–1049 (1973).
323. I. Overington, "The Importance of Coherence of Coupling When Viewing Through Visual Aids," *Opt. Laser Technol.* **52**:216–220 (1973).
324. I. Overington, "Some Considerations of the Role of the Eye as a Component of an Imaging System," *Optica Acta* **22**:365–374 (1975).
325. W. N. Charman and H. Whitefoot, "Astigmatism, Accommodation and Visual Instrumentation," *Applied Opt.* **17**:3903–3910 (1978).
326. P. Mouroulis, "On the Correction of Astigmatism and Field Curvature in Telescopic Systems," *Optica Acta* **29**:1133–1159 (1982).
327. G. J. Burton and N. D. Haig, "Criteria for Testing of Afocal Instruments," *Proc. Soc. Photo-Opt. Instrum. Eng.* **274**:191–201 (1981).
328. N. D. Haig and G. J. Burton, "Effects of Wavefront Aberration on Visual Instrument Performance, and a Consequential Test Technique," *Applied Opt.* **26**:492–500.
329. J. Eggert and K. J. Rosenbruch, "Vergleich der visuell und der photolektrisch gemessenen Abbildungs-guete von Fernroehren," *Optik* **48**:439–450 (1977).
330. G. J. Burton and N. D. Haig, "Effects of the Seidel Aberrations on Visual Target Discrimination," *J. Opt. Soc. Am. A* **1**:373–385 (1984).
331. R. Legras, N. Chateau, and W. N. Charman, "Assessment of Just Noticeable Differences for Refractive Errors and Spherical Aberration Using Visual Simulation," *Optom. Vis. Sci.* **81**:718–728 (2004).
332. P. Mouroulis and H. Zhang, "Visual Instrument Image Quality Metrics and the Effects of Coma and Astigmatism," *J. Opt. Soc. Am. A* **9**:34–42 (1992).
333. P. Mouroulis and G. C. Woo, "Chromatic Aberration and Accommodation in Visual Instruments," *Optik* **80**:161–166 (1989).
334. P. Mouroulis, T. G. Kim, and G. Zhao, "Transverse Color Tolerances for Visual Optical Systems," *Applied Opt.* **32**:7089–7094 (1993).

335. P. Mouroulis (ed.), *Visual Optical Systems*, McGraw-Hill, New York, 1999.
336. M. Koomen, R. Skolnik, and R. Tousey, "A Study of Night Myopia," *J. Opt. Soc. Am.* **41**:80–90 (1951).
337. D. G. Green and F. W. Campbell, "Effect of Focus on the Visual Response to a Sinusoidally Modulated Spatial Stimulus," *J. Opt. Soc. Am.* **55**:1154–1157 (1965).
338. W. N. Charman and J. A. M. Jennings, "The Optical Quality of the Retinal Image as a Function of Focus," *Br. J. Physiol. Opt.* **31**:119–134 (1976).
339. W. N. Charman, J. A. M. Jennings, and H. Whitefoot, "The Refraction of the Eye in Relation to Spherical Aberration and Pupil Size," *Br. J. Physiol. Opt.* **32**:78–93 (1978).
340. L. N. Thibos, "Unresolved Issues in the Prediction of Subjective Refraction from Wavefront Aberration Maps," *J. Refract. Surg.* **20**:S533–S536 (2004).
341. X. Cheng, A. Bradley, and L. N. Thibos, "Predicting Subjective Judgement of Best Focus with Objective Image Quality Metrics," *J. Vision* **4**:310–321 (2004).
342. L. N. Thibos, X. Hong, A. Bradley, and R. A. Applegate, "Accuracy and Precision of Objective Refraction from Wavefront Aberrations," *J. Vision* **4**:329–351 (2004).
343. D. A. Atchison, S. W. Fisher, C. A. Pedersen, and P. G. Ridall, "Noticeable, Troublesome and Objectionable Limits of Blur," *Vision Res.* **45**:1967–1974 (2005).
344. K. J. Ciuffreda, A. Selenow, B. Wang, B. Vasudevan, G. Zikos, and S. R. Ali, "Bothersome Blur: A Functional Unit of Blur Perception," *Vision Res.* **46**:895–901 (2006).
345. B. Wang and K. J. Ciuffreda, "Depth-of-Focus of the Human Eye: Theory and Clinical Implications," *Surv. Ophthalmol.* **51**:75–85 (2006).
346. D. G. Green and F. W. Campbell, "Effect of Focus on the Visual Response to a Sinusoidally Modulated Spatial Stimulus," *J. Opt. Soc. Am.* **55**:1154–1157 (1965).
347. D. G. Green, M. K. Powers, and M. S. Banks, "Depth of Focus, Eye Size and Visual Acuity," *Vision Res.* **20**:827–835 (1980).
348. F. L. van Nes and M. A. Bouman, "Spatial Modulation Transfer in the Human Eye," *J. Opt. Soc. Am.* **57**:401–407 (1967).
349. F. W. Campbell, "The Depth-of-Field of the Human Eye," *Optica Acta* **4**:157–164 (1957).
350. J. Tucker and W. N. Charman, "Depth-of-Focus and Accommodation for Sinusoidal Gratings as a Function of Luminance," *Am. J. Optom. Physiol. Opt.* **63**:58–70 (1986).
351. K. N. Ogle and J. T. Schwartz, "Depth-of-focus for the Human Eye," *J. Opt. Soc. Am.* **49**:273–280 (1959).
352. J. Tucker and W. N. Charman, "The Depth-of-Focus of the Human Eye for Snellen Letters," *Am. J. Optom. Physiol. Opt.* **52**:3–21 (1975).
353. W. N. Charman and H. Whitefoot, "Pupil Diameter and Depth-of-Field of the Human Eye as Measured by Laser Speckle," *Optica Acta* **24**:1211–1216 (1977).
354. D. A. Atchison, W. N. Charman, and R. L. Woods, "Subjective Depth-of-Focus of the Eye," *Optom. Vis. Sci.* **74**:511–520 (1997).
355. D. A. Goss and T. Grosvenor, "Reliability of Refraction—A Literature Review," *J. Am. Optom. Assoc.* **67**:619–630 (1996).
356. A. D. Miller, M. J. Kris, and A. C. Griffiths, "Effect of Small Focal Errors on Vision," *Optom. Vis. Sci.* **74**:521–526 (1997).
357. I. E. Loewenfeld, *The Pupil: Anatomy, Physiology and Clinical Applications*, Butterworth-Heinemann, London, 1999, pp. 295–317.
358. M. Stakenburg, "Accommodation without Pupillary Constriction," *Vision Res.* **31**:267–273 (1991).
359. N. J. Phillips, B. Winn, and B. Gilmartin, "Absence of Pupil Response to Blur-driven Accommodation," *Vision Res.* **32**:1775–1779 (1992).
360. F. Schaeffel, H. Wilhelm, and E. Zrenner, "Inter-Individual Variability in the Dynamics of Natural Accommodation in Humans: Relation to Age and Refractive Errors," *J. Physiol. (London)* **462**:301–320 (1993).
361. S. Kasthurirangan and A. Glasser, "Age Related Changes in the Characteristics of the Near Pupil Response," *Vision Res.* **46**:1393–1403 (2006).
362. H. Radhakrishnan and W. N. Charman, "Age-related Changes in Static accommodation and Accommodative Miosis," *Ophthalm. Physiol. Opt.* **27**:342–352 (2007).

363. F. W. Campbell, "Correlation of Accommodation between the Two Eyes," *J. Opt. Soc. Am.* **50**:738 (1960).
364. G. Heron, B. Winn, J. R. Pugh, and A. S. Eadie, "Twin Channel Infrared Optometer for Recording Binocular Accommodation," *Optom. Vis. Sci.* **66**:123–129 (1989).
365. F. W. Campbell, "The Minimum Quantity of Light Required to Elicit the Accommodation Reflex in Man," *J. Physiol. (London)* **123**:357–366 (1954).
366. C. A. Johnson, "Effects of Luminance and Stimulus Distance on Accommodation and Visual Resolution," *J. Opt. Soc. Am.* **66**:138–142 (1976).
367. D. A. Atchison, "Accommodation and Presbyopia," *Ophthal. Physiol. Opt.* **15**:255–2272 (1995)
368. K. J. Ciuffreda, "Accommodation, the Pupil and Presbyopia," *Borish's Clinical Refraction*, W. J. Benjamin (ed.), Saunders, Philadelphia, 1998, pp. 77–120.
369. A. Glasser and P. I. Kaufman, "The Mechanism of Accommodation in Primates," *Ophthalmology* **106**:863–872 (1999).
370. W. N. Charman, "The Eye in Focus: Accommodation and Presbyopia," *Clin. Exp. Optom.* **91**:207–225 (2008).
371. H. J. Wyatt, "Application of a Simple Mechanical Model of Accommodation to the Aging Eye," *Vision Res.* **33**:731–738 (1993).
372. A. P. A. Beers and G. L. van der Heijde, "In vivo Determination of the Biomechanical Properties of the Component Elements of the Accommodation Mechanism," *Vision Res.* **34**:2897–2905 (1994).
373. S. J. Judge and H. J. Burd, "Modelling the Mechanics of Accommodation and Presbyopia," *Ophthal. Physiol. Opt.* **22**:397–400 (2002).
374. H. Martin, R. Guthoff, T. Terwee, and K-P Schmitz, "Comparison of the Accommodation Theories of Coleman and Helmholtz by Finite Element Simulations," *Vision Res.* **45**:2910–2915 (2005).
375. F. W. Campbell and G. Westheimer, "Dynamics of the Focussing Response of the Human Eye," *J. Physiol. (London)* **151**:285–295 (1960).
376. S. D. Phillips, D. Shirachi, and L. Stark, "Analysis of Accommodation Times Using Histogram Information," *Am. J. Optom. Arch. Am. Acad. Optom.* **49**:389–401 (1972).
377. D. Shirachi, J. Liu, M. Lee, J. Jang, J. Wong, and L. Stark, "Accommodation Dynamics. I Range Non-Linearity," *Am. J. Optom. Physiol. Opt.* **55**:531–541 (1978).
378. J. Tucker and W. N. Charman, "Reaction and Response Times for Accommodation," *Am. J. Optom. Physiol. Opt.* **56**:490–503 (1979).
379. S. Kasthurirangan, A. S. Vilupuru, and A. Glasser, "Amplitude Dependent Accommodative Dynamics in Humans," *Vision Res.* **43**:2945–2956 (2003).
380. S. Kasthurirangan and A. Glasser, "Influence of Amplitude and Starting Point on Accommodative Dynamics in Humans," *Invest. Ophthalmol. Vis. Sci.* **46**:3463–3472 (2005).
381. E. F. Fincham, "The Accommodation Reflex and its Stimulus," *Br. J. Ophthalmol.* **35**:381–393 (1951).
382. F. W. Campbell and G. Westheimer, "Factors Influencing Accommodation Responses of the Human Eye," *J. Opt. Soc. Am.* **49**:568–571 (1959).
383. A. Troelstra, B. L. Zuber, D. Miller, and L. Stark, "Accommodative Tracking: a Trial and Error Function," *Vision Res.* **4**:585–594 (1964).
384. L. Stark, *Neurological Control Systems: Studies in Bioengineering*, Plenum Press, New York, 1968, Sect. III.
385. W. N. Charman and G. Heron, "On the linearity of Accommodation Dynamics," *Vision Res.* **40**:2057–2066 (2000).
386. E. Marg E. "An Investigation of Voluntary as Distinguished from Reflex Accommodation," *Am. J. Optom. Arch. Am. Acad. Optom.* **28**:347–356 (1951).
387. T. N. Cornsweet and H. D. Crane, "Training the Visual Accommodation System," *Vision Res.* **13**:713–715 (1973).
388. R. R. Provine and J. M. Enoch, "On Voluntary Ocular Accommodation," *Percept. Psychophys.* **17**:209–212 (1975).
389. R. J. Randle and M. R. Murphy, "The Dynamic Response of Visual Accommodation over a Seven-Day Period," *Am. J. Optom. Physiol. Opt.* **51**:530–540 (1974).
390. G. J. van der Wildt, M. A. Bouman, and J. van der Kraats, "The Effect of Anticipation on the Transfer Function of the Human Lens System," *Optica Acta* **21**:843–860 (1974).

391. P. B. Kruger and J. Pola, "Changing Target Size is a Stimulus for Accommodation," *J. Opt. Soc. Am. A* **2**:1832–1835 (1985).
392. T. Takeda, K. Hashimoto, N. Hiruma, and Y. Fukui, "Characteristics of Accommodation toward Apparent Depth," *Vision Res.* **39**:2087–2097 (1999).
393. W. N. Charman and G. Heron, "Fluctuations in Accommodation: A Review," *Ophthalm. Physiol. Opt.* **8**:153–164 (1988).
394. B. Winn and B. Gilmartin, "Current Perspective on Microfluctuations of Accommodation," *Ophthalm. Physiol. Opt.* **12**:252–256 (1992).
395. K. Toshida, F. Okuyama, and T. Tokoro, "Influences of the Accommodative Stimulus and Aging on the Accommodative Microfluctuations," *Optom. Vis. Sci.* **75**:221–226 (1998).
396. L. S. Gray, B. Winn, and B. Gilmartin, "Effect of Target Luminance on Microfluctuations of Accommodation," *Ophthalm. Physiol. Opt.* **13**:258–265 (1993).
397. L. R. Stark and D. A. Atchison, "Pupil Size, Mean Accommodation Response and the Fluctuations of Accommodation," *Ophthalm. Physiol. Opt.* **17**:316–323 (1997).
398. F. W. Campbell, J. G. Robson and G. Westheimer, "Fluctuations in Accommodation under Steady Viewing Conditions," *J. Physiol. (London)* **145**:579–594.
399. J. C. Kotulak and C. M. Schor, "Temporal Variations in Accommodation during Steady-State Conditions," *J. Opt. Soc. Am. A* **3**:223–227 (1986).
400. B. Winn, J. R. Pugh, B. Gilmartin, and H. Owens, "Arterial Pulse Modulates Steady-State Accommodation," *Curr. Eye Res.* **9**:971–975 (1990).
401. M. J. Collins, B. Davis, and J. Wood, "Microfluctuations of Steady-State Accommodation and the Cardiopulmonary System," *Vision Res.* **35**:2491–2502 (1995).
402. G. L. van der Heijde, A. P. A. Beers, and M. Dubbelman, "Microfluctuations of Steady-State Accommodation Measured with Ultrasonography," *Ophthalm. Physiol. Opt.* **16**:216–221 (1996).
403. P. Denieul, "Effects of Stimulus Vergence on Mean Accommodation Response, Microfluctuations of Accommodation and Optical Quality of the Human Eye," *Vision Res.* **22**:561–569 (1982).
404. M. Zhu, M. J. Collins, and D. R. Iskander, "The Contribution of Accommodation and the Ocular Surface to the Microfluctuations of Wavefront Aberration of the Eye," *Ophthalm. Physiol. Opt.* **26**:439–446 (2006).
405. G. Heron and C. Schor, "The Fluctuations of Accommodation and Ageing," *Ophthalm. Physiol. Opt.* **15**:445–449 (1995).
406. L. Stark and Y. Takahashi, "Absence of an Odd-Error Signal Mechanism in Human Accommodation," *IEEE Trans. Biomed. Eng.* **BME-12**:138–146 (1965).
407. M. Alpern, "Variability of Accommodation during Steady Fixation at Various Levels of Illuminance," *J. Opt. Soc. Am.* **48**:193–197 (1958).
408. H. D. Crane, *A Theoretical Analysis of the Visual Accommodation System in Humans*, Report NASA CR-606, NASA, Washington, D.C., 1966.
409. G. K. Hung, J. L. Semmlow, and K. J. Ciuffreda, "Accommodative Oscillation Can Enhance Average Accommodation Response: A Simulation Study," *IEEE Trans. Syst. Man. Cybern.* **SMC-12**:594–598 (1982).
410. W. N. Charman, "Accommodation and the Through-Focus Changes of the Retinal Image," *Accommodation and Vergence Mechanisms in the Visual System*, O. Franzén, H. Richter, and L. Stark (eds.), Birkhäuser Verlag, Basel, 2000, pp. 115–127.
411. B. Winn, "Accommodative Microfluctuations: a Mechanism for Steady-State Control of Accommodation," *Accommodation and Vergence Mechanisms in the Visual System*, O. Franzén, H. Richter, and L. Stark (eds.), Birkhäuser Verlag, Basel, 2000, pp. 129–140.
412. M. Millodot, "Effet des Microfluctuations de l'Accommodation sur l'Acuité Visuelle," *Vision Res.* **8**:73–80 (1968).
413. G. Walsh and W. N. Charman, "Visual sensitivity to Temporal Change in Focus and its Relevance to the Accommodation Response," *Vision Res.* **28**:1207–1221 (1988).
414. B. Winn, W. N. Charman, J. R. Pugh, G. Heron, and A. S. Eadie, "Perceptual Detectability of Ocular Accommodation Microfluctuations," *J. Opt. Soc. Am. A* **6**:459–462 (1989).

415. M. W. Morgan, "Accommodation and Its Relationship to Convergence," *Am. J. Optom. Arch. Am. Acad. Optom.* **21**:183–185 (1944).
416. H. Krueger, "Schwankungen der Akkommodation des menschlichen Auges bei mon- und binokularer Beobachtung," *Albrecht v. Graefes Arch. Ophthalmol.* **205**:129–133 (1978).
417. M. C. Nadell and H. A. Knoll, "The Effect of Luminance, Target Configuration and Lenses upon the Refractive State of the Eye. Parts I and II," *Am. J. Optom. Arch. Am. Acad. Optom.* **33**:24–42 and 86–95 (1956).
418. G. C. Heath, "Influence of Visual Acuity on Accommodative Responses of the Eye," *Am. J. Optom. Physiol. Opt.* **33**:513–534 (1956).
419. J. Tucker, W. N. Charman and P. A. Ward, "Modulation Dependence of the Accommodation Response to Sinusoidal Gratings," *Vision Res.* **26**:1693–1707 (1986).
420. D. A. Owens, "A Comparison of Accommodative Responsiveness and Contrast Sensitivity for Sinusoidal Gratings," *Vision Res.* **20**:159–167 (1980).
421. L. J. Bour, "The Influence of the Spatial Distribution of a Target on the Dynamic Response and Fluctuations of the Accommodation of the Human Eye," *Vision Res.* **21**:1287–1296 (1981).
422. J. Tucker and W. N. Charman, "Effect of Target Content at Higher Spatial Frequencies on the Accuracy of the Accommodation Response," *Ophthalm. Physiol. Opt.* **7**:137–142 (1987).
423. K. J. Ciuffreda, M. Dul, and S. K. Fisher, "Higher-Order Spatial Frequency Contribution to Accommodative Accuracy in Normal and Amblyopic Observers," *Clin. Vis. Sci.* **1**:219–229 (1987).
424. J. C. Kotulak and C. M. Schor, "The Effects of Optical Vergence, Contrast, and Luminance on the Accommodative Response to Spatially Bandpass Filtered Targets," *Vision Res.* **27**:1797–1806 (1987).
425. K. J. Ciuffreda, M. Rosenfield, J. Rosen, A. Azimi, and E. Ong, "Accommodative Responses to Naturalistic Stimuli," *Ophthalm. Physiol. Opt.* **10**:168–174 (1990).
426. K. J. Ciuffreda, "Accommodation to Gratings and More Naturalistic Stimuli," *Optom. Vis. Sci.* **68**:243–260 (1991).
427. H. Ripps, N. B. Chin, I. M. Siegel, and G. M. Breinin, "The Effect of Pupil Size on Accommodation, Convergence and the AC/A ratio," *Invest Ophthalmol.* **1**:127–135 (1962).
428. R. T. Hennessy, R. Iida, K. Shiina, and H. W. Leibowitz, "The Effect of Pupil Size on Accommodation," *Vision Res.* **16**:587–589 (1976).
429. P. A. Ward and W. N. Charman, "Effect of Pupil Size on Steady-State Accommodation," *Vision Res.* **25**:1317–1326 (1985).
430. G. C. Heath, "Accommodative Responses of Totally Color Blind Observers," *Am. J. Optom. Arch. Am. Acad. Optom.* **33**:457–465 (1956).
431. J. Otto and D. Safra, "Ergebnisse objectiver Akkommodationsmessungen an Augen mit organisch bedingtem Zentralskotomn," *Albrecht v. Graefes Arch. Klin. Ophthalmol.* **192**:49–56 (1974).
432. I. C. J. Wood and A. T. Tomlinson, "The Accommodative Response in Amblyopia," *Am. J. Optom. Physiol. Opt.* **52**:243–247 (1975).
433. K. J. Ciuffreda and D. Rumpf, "Contrast and Accommodation in Amblyopia," *Vision Res.* **25**:1445–1447 (1985).
434. W. N. Charman "Static Accommodation and the Minimum Angle of Resolution," *Am. J. Optom. Physiol. Opt.* **63**:915–921 (1986).
435. H. W. Leibowitz and D. A. Owens, "Anomalous Myopias and the Intermediate Dark Focus of Accommodation," *Science* **189**:1121–11128 (1975).
436. H. W. Leibowitz and D. A. Owens, "Night Myopia and the Intermediate Dark Focus of Accommodation," *J. Opt. Soc. Am.* **65**:133–147 (1975).
437. H. W. Leibowitz and D. A. Owens, "New Evidence for the Intermediate Position of Relaxed Accommodation," *Doc. Ophthalmol.* **46**:1121–1128 (1978).
438. G. Smith, "The Accommodative Resting States, Instrument Accommodation and Their Measurement," *Optica Acta* **30**:347–359 (1983).
439. N. A. McBrien and M. Millodot, "The Relationship Between Tonic Accommodation and Refractive Error," *Invest. Ophthalmol. Vis. Sci.* **28**:997–1004 (1987).
440. F. M. Toates, "Accommodation Function of the Human Eye," *Physiol. Rev.* **52**:828–863 (1972).

441. M. Rosenfield, K. J. Ciuffreda, G. K. Hung, and B. Gilmartin, "Tonic Accommodation—a Review. 1. Basic Aspects," *Ophthalm. Physiol. Opt.* **13**:266–284 (1993).
442. M. Rosenfield, K. J. Ciuffreda, G. K. Hung, and B. Gilmartin, "Tonic Accommodation—a Review. 2. Accommodative Adaptation and Clinical Aspects," *Ophthalm. Physiol. Opt.* **14**:265–277 (1994).
443. H. W. Leibowitz, K. W. Gish, and J. B. Sheehy, "Role of Vergence Accommodation in Correcting for Night Myopia," *Am. J. Optom. Physiol. Opt.* **65**:383–386 (1988).
444. J. L. Semmlow and G. K. Hung, "The Near Response: Theories of Control," *Vergence Eye Movements: Basic and Clinical Concepts*, C. M. Schor and K. J. Ciuffreda (eds.), Butterworths, Boston, 1983, pp. 175–195.
445. G. K. Hung, K. J. Ciuffreda, and M. Rosenfield, "Proximal Contribution to a Linear Static Model of Accommodation and Vergence," *Ophthalm. Physiol. Opt.* **16**:31–41 (1996).
446. C. M. Schor and S. R. Bharadwaj, "Pulse-Step Models of Control Strategies for Dynamic Ocular Accommodation and Disaccommodation," *Vision Res.* **46**:242–258 (2006).
447. W. N. Charman and J. Tucker, "Accommodation and Color," *J. Opt. Soc. Am.* **68**:459–471 (1978).
448. J. V. Lovasik and H. Kergoat, "Accommodative Performance for Chromatic Displays," *Ophthalm. Physiol. Opt.* **8**:443–440 (1988).
449. W. N. Charman, "Accommodation Performance for Chromatic Displays," *Ophthalm. Physiol. Opt.* **9**:459–463 (1989).
450. W. R. Bobier, M. C. W. Campbell, and M. Hinch, "The Influence of Chromatic Aberration on the Static Accommodative Response," *Vision Res.* **32**:823–832 (1992).
451. D. A. Atchison, N. C. Strang, and L. R. Stark, "Dynamic Accommodation Responses to Stationary Colored Targets," *Optom. Vis. Sci.* **81**:699–711 (2004).
452. R. Home and J. Poole, "Measurement of the Preferred Binocular Dioptric Settings at High and Low Light Level," *Optica Acta* **24**:97 (1977).
453. M. F. Wesner and R. J. Miller, "Instrument Myopia Conceptions, Misconceptions, and Influencing Factors," *Doc. Ophthalmol.* **62**:281–308 (1986).
454. G. G. Heath, "Components of Accommodation," *Am. J. Optom. Arch. Am. Acad. Optom.* **33**:569–579 (1956).
455. S. C. Hokoda and K. J. Ciuffreda, "Theoretical and Clinical Importance of Proximal Vergence and Accommodation," *Vergence Eye Movements: Basic and Clinical Concepts*, C. M. Schor and K. K. Ciuffreda (eds.), Butterworths, Boston, 1983, pp. 75–97.
456. G. Smith, K. C. Tan, and M. Letts, "Binocular Optical Instruments and Binocular Vision," *Clin. Exper. Optom.* **69**:137–144 (1986).
457. A. Duane, "Studies in Monocular and Binocular Accommodation with Their Clinical Applications," *Am. J. Ophthalmol. Ser. 3* **5**:865–877 (1922).
458. D. Hamasaki, J. Ong, and E. Marg, "The Amplitude of Accommodation in Presbyopia," *Am. J. Optom. Arch. Am. Acad. Optom.* **33**:3–14 (1956).
459. H. W. Hofstetter, "A Longitudinal Study of Amplitude Changes in Presbyopia," *Am. J. Optom. Arch. Am. Acad. Optom.* **42**:3–8 (1965).
460. C. Ramsdale and W. N. Charman, "A Longitudinal Study of the Changes in Accommodation Response," *Ophthalm. Physiol. Opt.* **9**:255–263 (1989).
461. W. N. Charman, "The Path to Presbyopia: Straight or Crooked?," *Ophthalm. Physiol. Opt.* **9**:424–430 (1989).
462. J. A. Mordi and K. J. Ciuffreda, "Static Aspects of Accommodation: Age and Presbyopia," *Vision Res.* **38**:1643–1653 (1998).
463. M. Kalsi, G. Heron, and W. N. Charman, "Changes in the Static Accommodation Response with Age," *Ophthalm. Physiol. Opt.* **21**:77–84 (2001).
464. G. Heron and W. N. Charman, "Accommodation as a function of age and the linearity of the Response Dynamics," *Vision Res.* **44**:3119–3130 (2004).
465. G. Heron, W. N. Charman, and C. Schor, "Dynamics of the Accommodation Response to Abrupt Changes in Target Vergence as a Function of Age," *Vision Res.* **41**:507–519 (2001).
466. J. A. Mordi and K. J. Ciuffreda, "Dynamic Aspects of Accommodation: Age and Presbyopia," *Vision Res.* **44**:591–601 (2004).
467. S. Kasthurirangan and A. Glasser, "Age Related Changes in Accommodation Dynamics in Humans," *Vision Res.* **46**:1507–1519 (2006).

468. G. Smith, "Schematic Eyes: History, Description and Applications," *Clin. Exp. Optom.* **78**:176–189 (1995).
469. D. A. Atchison and G. Smith, *Optics of the Human Eye*, Butterworth-Heinemann, Oxford, 2000, pp. 39–47 and 160–179.
470. R. B. Rabbetts, *Clinical Visual Optics*, 3rd ed., Butterworth-Heinemann, Oxford, 1998, pp. 207–229.
471. A. E. A. Ridgway, "Intraocular Lens Implants," *Vision and Visual Dysfunction*, vol. 1, *Visual Optics and Instrumentation*, W. N. Charman (ed.), Macmillan, London, 1991, pp. 120–137.
472. D. Sliney and M. Wolbarsht, *Safety with Lasers and Other Optical Sources*, Plenum, New York, 1980.
473. D. H. Sliney, "Measurement of Light and the Geometry of Exposure of the Human Eye," *Vision and Visual Dysfunction*, vol. 16, *The Susceptible Visual Apparatus*, J. Marshall (ed.), Macmillan, London 1991, pp. 23–29.
474. Y. Le Grand and S. G. El Hage, *Physiological Optics*, Springer-Verlag, Berlin, 1980, pp. 64–66.
475. H. H. Emsley, *Visual Optics*, vol. 1, 5th ed., Hatton Press, London, 1953.
476. J. W. Blaker, "Toward an Adaptive Model of the Human Eye," *J. Opt. Soc. Am.* **70**:220–223 (1980).
477. J. W. Blaker, "A Comprehensive Model of the Aging, Accommodative, Adult Eye," *Technical Digest on Ophthalmic and Visual Optics*, vol 2, Optical Society of America, Washington, D.C., 1991, pp. 28–31.
478. A. Popielek-Masajada and H. T. Kasprzak, "A New Schematic Eye Model Incorporating Accommodation," *Optom. Vis. Sci.* **76**:720–727 (1999).
479. W. Lotmar, "Theoretical Eye Model with Aspherics," *J. Opt. Soc. Am.* **61**:1522–1529 (1971).
480. A. C. Kooijman (1983) "Light Distribution of the Retina of a Wide-Angle Theoretical Eye," *J. Opt. Soc. Am.* **73**:1544–1550 (1983).
481. R. Navarro, J. Santamaria, and J. Bescós, "Accommodation-Dependent Model of the Human Eye with Aspherics," *J. Opt. Soc. Am. A* **2**:1273–1281 (1985).
482. M. C. M. Dunne and D. A. Barnes, "Schematic Modelling of Peripheral Astigmatism in Real Eyes," *Ophthalm. Physiol. Opt.* **7**:235–239 (1987).
483. S. Patel, J. Marshall, and F. W. Fitzke, "Model for Predicting the Optical Performance of the Eye in Refractive Surgery," *Refract. Corneal Surg.* **9**:366–375 (1993).
484. H-L Liou and N. A. Brennan, "Anatomically Accurate, Finite Model Eye for Optical Modeling," *J. Opt. Soc. Am. A* **14**:1684–1695 (1997).
485. I. Escudero-Sanz and R. Navarro, "Off-Axis Aberrations of a Wide-angle Schematic Eye Model," *J. Opt. Soc. Am. A* **16**:1881–1891 (1999).
486. Y-J Liu, Z. Q. Wang, L. -P. Song, and G. -G. Mu, "An Anatomically Accurate Eye Model with a Shell-Structure Lens," *Optik* **116**:241–246 (2005).
487. A. V. Goncharov and C. Dainty, "Wide-Field Schematic Eye Models with Gradient-Index Lens," *J. Opt. Soc. Am. A* **24**:2157–2174 (2007).
488. L. N. Thibos, M. Ye, X. Zhang, and A. Bradley, "The Chromatic Eye: a New Reduced-eye Model of Ocular Chromatic Aberration in Humans," *Applied Opt.* **32**:3594–3600 (1992).
489. L. N. Thibos, M. Ye, X. Zhang, and A. Bradley, "Spherical Aberration of the Reduced Schematic Eye with Elliptical Refracting Surface," *Optom. Vis. Sci.* **74**:548–556 (1997).
490. Y. -Z. Wang and L. N. Thibos, "Oblique (Off-axis) Astigmatism of the Reduced Schematic Eye with Elliptical Refracting Surface," *Optom. Vis. Sci.* **74**:557–562 (1997).
491. D. A. Atchison, "Oblique Astigmatism of the Indiana Eye," *Optom. Vis. Sci.* **75**:247–248 (1998).
492. F. W. Campbell and D. G. Green, "Monocular versus Binocular Visual Acuity," *Nature* **208**:191–192 (1965).
493. M. Lombardo, G. Lombardo, and S. Serrao, "Interocular High-Order Corneal Wavefront Aberration Symmetry," *J. Opt. Soc. Am. A* **23**:777–787 (2006).
494. S. Marcos and S. A. Burns, "On the Symmetry between Eyes of Wavefront Aberration and Cone Directionality," *Vision Res.* **40**:2437–2447 (2000).
495. M. J. Collins and A. S. Bruce, "Factors Influencing Performance with Monovision," *J. Brit. Contact Lens Assoc.* **17**:83–89 (1994).
496. J. Meyler, "Presbyopia," *Contact Lens Practice*, N. Efron (ed.), Butterworth-Heinemann, Oxford, 2002, pp. 261–274.

497. B. J. W. Evans, "Monovision: A Review," *Ophthal. Physiol. Opt.* **27**:417–439 (2007).
498. C. M. Schor and M. C. Flom, "The Relative Values of Stereopsis as a Function of Viewing Distance," *Am. J. Optom. Arch. Am. Acad. Optom.* **46**:805–809 (1969).
499. R. S. Harvey, "Some Statistics of Interpupillary Distance," *Optician* **184**(4766):29 (1982).
500. R. Sekuler and R. Blake, *Perception*, 5th ed., McGraw-Hill, New York, 2005.
501. N. A. Valyus, *Stereoscopy*, Focal Press, London, 1966.
502. A. Lit, "Depth Discrimination Thresholds as a Function of Binocular Differences of Retinal Illuminance at Scotopic and Photopic Levels," *J. Opt. Soc. Am.* **49**:746–752 (1959).
503. S. C. Rawlings and T. Shipley, "Stereoscopic Acuity and Horizontal Angular Distance from Fixation," *J. Opt. Soc. Am.* **59**:991–993 (1969).
504. W. P. Dwyer and A. Lit, "Effect of Luminance-matched Wavelength on Depth Discrimination at Scotopic and Photopic Levels of Target Illuminance," *J. Opt. Soc. Am.* **60**:127–131 (1970).
505. A. Arditi, "Binocular Vision," *Handbook of Perception and Human Performance*, vol. 1, K. R. Boff, L. Kaufman, and J. P. Thomas (eds.), Wiley, New York, 1986, Chapter 23.
506. R. R. Fagin and J. R. Griffin, "Stereoacuity Tests: Comparison of Mathematical Equivalents," *Am. J. Optom. Physiol. Opt.* **59**:427–438 (1982).
507. R. B. Rabbetts, *Clinical Visual Optics*, 4th ed., Elsevier, Oxford, 2007, Chapter 11.
508. G. Heron, S. Dholakia, D. E. Collins, and H. McLaughlan, "Stereoscopic Thresholds in Children and Adults," *Am. J. Optom. Physiol. Opt.* **62**:505–515 (1985).
509. H. J. Howard, "A Test for Judgement of Distance," *Am. J. Ophthalmol.* **2**:656–675 (1919).
510. G. Westheimer, "Effect of binocular Magnification Devices on Stereoscopic Depth Resolution," *J. Opt. Soc. Am.* **45**:278–280 (1956).
511. W. N. Charman and J. A. M. Jennings, "Binocular Vision in Relation to Stereoscopic Instruments and Three-Dimensional Displays," *Vision and Visual Dysfunction*, vol. 1, *Visual Optics and Instrumentation*, W. N. Charman (ed.), Macmillan, London, 1991, pp. 326–344.
512. D. B. Diner and D. H. Fender, *Human Engineering in Stereoscopic Viewing Devices*, Plenum Press, New York, 1993, pp. 49–65.
513. G. Smith and D. A. Atchison, *The Eye and Visual Optical Instruments*, Cambridge UP, Cambridge, 1997, pp. 727–746.
514. E. Peli, "Optometric and Perceptual Issues with Head-Mounted Displays," *Visual Instrumentation: Optical Design and Engineering Principles*, P. Mouroulis (ed.), McGraw-Hill, New York, 1999, pp. 205–276.
515. T. L. Williams, "Testing of Visual Instrumentation," *Visual Instrumentation: Optical Design and Engineering Principles* P. Mouroulis (ed.), McGraw-Hill, New York, 1999, pp. 353–421.
516. J. A. M. Jennings, "Binocular Vision through Correcting Lenses: Aniseikonia," *Vision and Visual Dysfunction*, vol. 1, *Visual Optics and Instrumentation*, W. N. Charman (ed.), Macmillans, London, 1991, pp. 163–182.
517. M. A. Taylor Kulp, T. W. Raasch, and M. Polasky, "Patients with Anisometropia and Aniseikonia," *Clinical Refraction*, W. J. Benjamin (ed.), Saunders, Philadelphia, 1998, pp. 1134–1159.
518. B. Winn, R. G. Ackerley, C. A. Brown, F. K. Murray, J. Prais, and M. F. St John, "Reduced Aniseikonia in Axial Ametropia with Contact Lens Correction," *Ophthal. Physiol. Opt.* **8**:341–344 (1988).
519. A. L. Yarbus, *Eye Movements and Vision*, Plenum, New York, 1967.
520. M. Land, N. Mennie, and J. Rusted, "The Roles of Vision and Eye Movements in the Control of Activities of Daily Living," *Perception* **28**:1311–1328 (1999).
521. M. A. Alpern, "Movements of the Eyes," *The Eye*, vol. 3, *Muscular Mechanisms*, H. Davson (ed.), Academic Press, New York, 1969, pp. 5–214.
522. R. H. S. Carpenter, *Movements of the Eyes*, 2nd ed., Pion, London, 1988.
523. R. H. S. Carpenter (ed.), *Vision and Visual Dysfunction*, vol. 8, *Eye Movements*, Macmillan, London, 1991.
524. R. J. Leigh and D. S. Zee, *The Neurology of Eye Movements*, 4th ed., Oxford UP, Oxford, 2006.
525. G. Fry and W. W. Hill, "The Mechanics of Elevating the Eye," *Am. J. Optom. Arch. Am. Acad. Optom.* **40**:707–716 (1963).

526. C. Rashbass and G. Westheimer, "Independence of Conjunctive and Disjunctive Eye Movements," *J. Physiol. (London)* **159**:361–364 (1961).
527. D. A. Robinson, "The Mechanics of Human Saccadic Eye Movements," *J. Physiol. (London)* **174**:245–264 (1964).
528. A. T. Bahill, A. Brockenbrough, and B. T. Troost, "Variability and Development of a Normative Data Base for Saccadic Eye Movements," *Invest. Ophthalmol.* **21**:116–125 (1981).
529. R. H. S. Carpenter, "The Neural Control of Looking," *Current Biology* **10**(8):R291–R293 (2000).
530. F. W. Campbell and R. H. Wurtz, "Saccadic Omission: Why We Do Not See a Grey Out during a Saccadic Eye Movement," *Vision Res.* **18**:1297–1303 (1978).
531. H. Collewijn and E. P. Tamminga, "Human Smooth Pursuit and Saccadic Eye Movements during Voluntary Pursuit of Different Target Motions on Different Backgrounds," *J. Physiol. (London)* **351**:217–250 (1984).
532. C. Rashbass and G. Westheimer, "Disjunctive Eye Movements," *J. Physiol. (London)* **159**:339–360 (1961).
533. C. J. Erkelers, J. van der Steen, R. M. Steinman, and H. Collewijn, "Ocular Vergence under Natural Conditions. I Continuous Changes of Target Distance Along the Median Plane," *Proc. R. Soc. Lond.* **B236**:417–440 (1989).
534. C. J. Erkelers, R. M. Steinman, and H. Collewijn, "Ocular Vergence Under Natural Conditions. II Gaze Shifts between Real Targets Differing in Distance and Direction," *Proc. R. Soc. Lond.* **B236**:441–465 (1989).
535. J. T. Enright, "Perspective Vergence: Oculomotor Responses to Line Drawings," *Vision Res.* **27**:1513–1526 (1987).
536. A. E. Kertesz, "Vertical and Cyclofusional Disparity Vergence," *Vergence Eye Movements: Clinical and Applied*, C. M. Schor and K. J. Ciuffreda (eds.), Butterworths, Boston, 1983, pp. 317–348.
537. J. M. Findlay, "Frequency Analysis of Human Involuntary Eye Movement," *Kybernetik* **8**:207–214 (1971).
538. R. W. Ditchburn, *Eye Movements and Visual Perception*, Clarendon Press, Oxford, 1973.
539. J. Nachmias, "Determinants of the Drift of the Eye During Monocular Fixation," *J. Opt. Soc. Am.* **51**:761–766 (1961).
540. J. Nachmias, "Two-dimensional motion of the retinal image during monocular fixation," *J. Opt. Soc. Am.* **49**:901–908 (1959).
541. H. C. Bennet-Clark, "The Oculomotor Response to Small Target Displacements," *Optica Acta* **11**:301–314 (1964).
542. W. H. Marshall and S. A. Talbot, "Recent Evidence for Neural Mechanisms in Vision Leading to a General Theory of Sensory Acuity," *Biol. Symp.* **7**:117–164 (1942).
543. O. Packer and D. R. Williams, "Blurring by Fixational Eye Movements," *Vision Res.* **32**:1931–1939 (1992).

This page intentionally left blank.

DO NOT DUPLICATE

VISUAL PERFORMANCE

Wilson S. Geisler

*Department of Psychology
University of Texas
Austin, Texas*

Martin S. Banks

*School of Optometry
University of California
Berkeley, California*

2.1 GLOSSARY

A	amplitude
a	interpupillary distance
A_p	effective area of the entrance pupil
C	contrast
c_o	maximum concentration of photopigment
d	horizontal disparity
d_e	distance from the image plane to the exit pupil of an optical system
d_θ	average disparity between two points and the convergence point
$E(\lambda)$	photopic spectral illuminance distribution
$E_e(\lambda)$	spectral irradiance distribution
f	spatial frequency of a sinusoid
I_o	half-bleaching constant
J_o	maximum photocurrent
l	length of outer segment
$L(\lambda)$	photopic spectral luminance distribution
$L_e(\lambda)$	spectral radiance distribution
m	magnification of the exit pupil relative to the actual pupil
N	total effective photons absorbed per second
n_r	index of refraction of the media where the image plane is located
$n(\lambda)$	spectral photon-flux irradiance distribution
p	proportion of unbleached photopigment
t_o	time constant of photopigment regeneration
$t(\lambda)$	transmittance function of the ocular media
$V(\lambda)$	standard photopic spectral sensitivity function of the human visual system

$\alpha(\lambda)$	absorptance spectrum
ΔC	contrast increment or decrement
Δf	frequency increment or decrement
Δz	distance between any pair of points in the depth dimension
ε	retinal eccentricity
$\varepsilon(\lambda)$	extinction spectrum
θ	convergence angle of the eyes
θ	orientation of a sinusoid
κ	collection area, or aperture, of a photoreceptor
λ	wavelength of the light in a vacuum
ξ	isomerization efficiency
Σ	covariance matrix for a gaussian noise process
σ_o	half-saturation constant
τ	time interval
τ_{opt}	optimum time interval
$\Phi(\cdot)$	cumulative standard normal probability function
ϕ	phase of a sinusoid

2.2 INTRODUCTION

Physiological optics concerns the study of (1) how images are formed in biological eyes, (2) how those images are processed in the visual parts of the nervous system, and (3) how the properties of image formation and neural processing manifest themselves in the perceptual performance of the organism. The previous chapter reviewed image formation; this chapter briefly describes the neural processing of visual information in the early levels of the human visual system, and summarizes, somewhat more extensively, what is known about human visual performance.

An enormous amount of information about the physical environment is contained in the light reaching the cornea of the eye. This information is critical for many of the tasks the human observer must perform, including identification of objects and materials, determination of the three-dimensional structure of the environment, navigation through the environment, prediction of object trajectories, manipulation of objects, and communication with other individuals. The performance of a human observer in a given visual task is limited by the amount of information available in the light at the cornea and by the amount (and type) of information encoded and transmitted by the successive stages of visual processing. This chapter presents a concise description of visual performance in a number of fundamental visual tasks. It also presents a concise description of the physiological and psychological factors believed to underlie performance in those tasks. Two major criteria governed the selection of the material to be presented. First, we attempted to focus on quantitative data and theories that should prove useful for developing rigorous models or characterizations of visual performance. Second, we attempted to focus on data and theories that have a firm empirical basis, including at least some knowledge of the underlying biological mechanisms.

2.3 OPTICS, ANATOMY, PHYSIOLOGY OF THE VISUAL SYSTEM

Image Formation

The processing of visual information begins with the optics of the eye, which consists of three major components: the cornea, pupil, and lens. The optical components are designed to form a sharp image at the layer of the photoreceptors in the retina. Neil Charman (Chap. 1) discusses many of

the details concerning how these optical components affect image quality. We briefly describe a few of the most useful formulas and methods for computing the approximate size, location, quality, and intensity of images formed at the photoreceptors. Our aim is to provide descriptions of image formation that might prove useful for developing models or characterizations of performance in perceptual tasks.

The sizes and locations of retinal images can be found by projecting points on the objects along straight lines through the image (posterior) nodal point until they intersect the retinal surface. The intersections with the retinal surface give the image locations corresponding to the points on the objects. The angles between pairs of projection lines are *visual angles*. Image locations are usually described in terms of the visual angles that projection lines make with respect to a reference projection line (the *visual axis*), which passes through the nodal point and the center of the fovea. The radial visual angle between a projection line from an object and the visual axis is the object's *eccentricity*. Image sizes are often described by the visual angles between key points on the object. For purposes of computing the size and location of images, it is usually sufficient to use a simplified model of the eye's optics, such as the *reduced eye*, which consists of a single spherical refracting surface (radius of curvature = 5.5 mm) and a retinal surface located 16.7 mm behind the nodal point (see Fig. 2c, Chap. 1).

A general method for computing the quality of retinal images is by convolution with a point-spread function $h(x, y)$. Specifically, if $o(x, y)$ is the luminance (or radiance) of the object, and $i(x, y)$ is the image illuminance (or irradiance), then

$$i(x, y) = o(x, y) ** h(x, y) \quad (1)$$

where $**$ represents the two-dimensional convolution operator. The shape of the point-spread function varies with wavelength and with retinal location. The precise way to deal with wavelength is to perform a separate convolution for each wavelength in the spectral luminance distribution of the object. In practice, it often suffices to convolve with a single point-spread function, which is the weighted average, across wavelength, of the monochromatic point-spread functions, where the weights are given by the shape of the spectral luminance distribution. To deal with retinal location, one can make use of the fact that the human point-spread function changes only gradually with retinal eccentricity out to about 20 deg;¹ thus, a large proportion of the visual field can be divided into a few annular regions, each with a different point-spread function.

Calculation of the point-spread function is normally accomplished by finding the transfer function, $H(u, v)$.^{2,3} (The point-spread function can be obtained, if desired, by an inverse Fourier transform.) The transfer function is given by the autocorrelation of the generalized pupil function followed by normalization to a peak value of 1.0:

$$T(u, v) = p(x, y) e^{iW(x, y, \lambda)} \otimes p(x, y) e^{-iW(x, y, \lambda)} \Big|_{\mathbf{x} = (\lambda d_e / mn_r) \mathbf{u}, \mathbf{y} = (\lambda d_e / mn_r) \mathbf{v}} \quad (2)$$

$$H(u, v) = \frac{T(u, v)}{T(0, 0)} \quad (3)$$

where λ is the wavelength of the light in a vacuum, n_r is the index of refraction of the media where the image plane is located, d_e is the distance from the image plane to the exit pupil of the optical system, and m is the magnification of the exit pupil relative to the actual pupil. The generalized pupil function is the product of the simple pupil function, $p(x, y)$ (transmittance as a function of position within the actual pupil), and the aberration function, $e^{iW(x, y, \lambda)}$. The exit pupil is the apparent pupil, when viewed from the image plane. The size of and distance to the exit pupil can be found by ray-tracing a schematic eye. For the Le Grand eye, the relevant parameters are approximately as follows: $m = 1.03$, $d_e = 20.5$ mm, and $n = 1.336$. Average values of the monochromatic and chromatic aberrations are available (see Chap. 1) and can be used as estimates of $W(x, y, \lambda)$.

Equation (3) can be used to compute approximate point-spread functions (and hence image quality) for many stimulus conditions. However, for some conditions, direct measurements (or psychophysical measurements) of the point-spread function are also available (see Chap. 1), and are easier to deal with. For broadband (white) light and a well-accommodated eye, the axial point-spread functions directly measured by Campbell and Gubisch⁴ are representative. A useful set

of monochromatic point-spread functions was measured at various eccentricities by Navarro et al.¹ These point-spread functions can be used directly in Eq. (1) to compute approximate image quality.

The approximate retinal irradiance for extended objects is given by the following formula:

$$E_e(\lambda) = \frac{A_p}{278.3} L_e(\lambda) t(\lambda) \quad (4)$$

where $E_e(\lambda)$ is the retinal spectral irradiance distribution (watts \cdot m⁻² \cdot nm⁻¹), $L_e(\lambda)$ is the spectral radiance distribution of the object (watts \cdot m⁻² \cdot sr⁻¹ \cdot nm⁻¹), $t(\lambda)$ is the transmittance of the ocular media (see Chap. 1), and A_p is the effective area of the entrance pupil (mm²). [Note, $A_p = \iint p(x/m', y/m') dx dy$, where m' is magnification of the entrance pupil relative to the actual pupil; the entrance pupil is the apparent size of the pupil when viewed from outside the eye.]

Photopic retinal illuminance, $E(\lambda)$ (candelas \cdot nm⁻¹), is computed by an equivalent formula where the spectral radiance distribution, $L_e(\lambda)$, is replaced by the spectral luminance distribution, $L(\lambda)$ (candelas \cdot m⁻² \cdot nm⁻¹), defined by

$$L(\lambda) = 683 V(\lambda) L_e(\lambda) \quad (5)$$

where $V(\lambda)$ is the standard photopic spectral sensitivity function of the human visual system.⁵

In theoretical calculations, it is often useful to express light levels in terms of photon flux rather than in terms of watts. The photon-flux irradiance on the retina, $n(\lambda)$ (quanta \cdot sec⁻¹ \cdot deg⁻² \cdot nm⁻¹) is computed by multiplying the retinal irradiance, $E_e(\lambda)$, by 8.4801×10^{-8} which converts m² to deg² (based upon the *reduced eye*), and by λ/ch which converts watts into quanta/sec (where c is the speed of light in a vacuum, and h is Planck's constant). Thus,

$$n(\lambda) = 1.53 \times 10^6 A_p L_e(\lambda) t(\lambda) \lambda \quad (6)$$

and, by substitution of Eq. (5) into Eq. (6),

$$n(\lambda) = 2.24 \times 10^3 A_p \frac{L(\lambda)}{V(\lambda)} t(\lambda) \lambda \quad (7)$$

Most light-measuring devices report radiance, $L_e(\lambda)$, or luminance, $L(\lambda)$; Eqs. (6) and (7) allow conversion to retinal photon-flux irradiance, $n(\lambda)$. For more details on the calculation of retinal intensity, see Wyszecki and Stiles.⁵

Image Sampling by the Photoreceptors

The image formed at the receptor layer is described by a four-dimensional function $n(x, y, t, \lambda)$, which gives the mean photon-flux irradiance (quanta \cdot sec⁻¹ \cdot deg² \cdot nm⁻¹) as a function of space (x, y), time (t), and wavelength (λ). This four-dimensional function also describes the photon noise in the image. Specifically, photon noise is adequately described as an inhomogeneous Poisson process; thus, the variance in the number of photons incident in a given interval of space, time, and wavelength is equal to the mean number of photons incident in that same interval.

The photoreceptors encode the (noisy) retinal image into a discrete representation in space and wavelength, and a more continuous representation in time. The image sampling process is a crucial step in vision that can, and often does, result in significant information loss. The losses occur because physical and physiological constraints make it impossible to sample all four dimensions with sufficiently high resolution.

As shown in the schematic diagram in Fig. 1, there are two major types of photoreceptors: rods and cones. They play very different functional roles in vision; rods subserve vision at low light levels and cones at high light levels. There are three types of cones, each with a different spectral sensitivity (which is the result of having different photopigments in the outer segment). The "long" (L), "middle" (M), and "short" (S) wavelength cones have peak spectral sensitivities at wavelengths of approximately 570, 540, and 440 nm, respectively. Information about the spectral wavelength distribution of

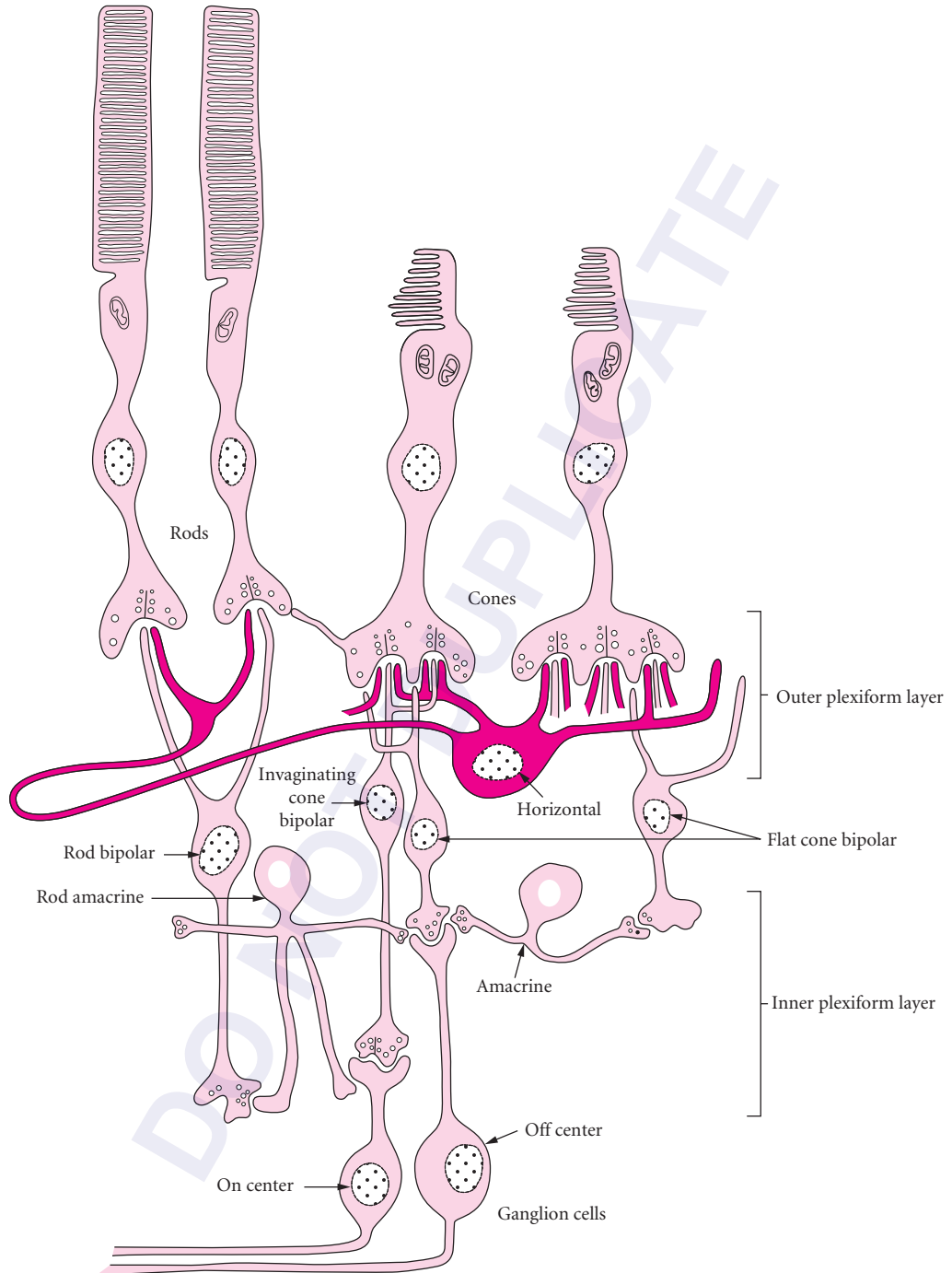


FIGURE 1 Schematic diagram of the retinal neurons and their major synaptic connections. Note that rods, rod bipolar cells, and rod amacrine cells are absent in the fovea. (From Ref. 263.)

the light falling on the retina is encoded by the relative activities of the L, M, and S cones. All rods have the same spectral sensitivity (and the same photopigment), peaking at about 500 nm.

The quality of spatial, temporal, and wavelength information encoded by the photoreceptors depends upon: (1) the spatial distribution of the photoreceptors across the retina, (2) the efficiency with which individual photoreceptors absorb light at different wavelengths (the absorptance spectrum), (3) the area over which the individual photoreceptors collect light (the receptor aperture), and (4) the length of time over which the individual photoreceptors integrate light.

The spatial distribution of cones and rods is highly nonuniform. Figure 2a shows the typical density distribution of rod and cone photoreceptors across the retina (although there are individual differences; e.g., see Ref. 6). Cone density decreases precipitously with eccentricity; rods are absent in the fovea and reach a peak density at 20 deg. If the receptor lattice were perfectly regular, the highest unambiguously resolvable spatial frequency (the Nyquist limit), would be half the linear density (in $\text{cells} \cdot \text{deg}^{-1}$). Under normal viewing conditions, the Nyquist limit does not affect vision in the fovea because the eye's optics eliminate spatial frequencies at and above the limit. However, the

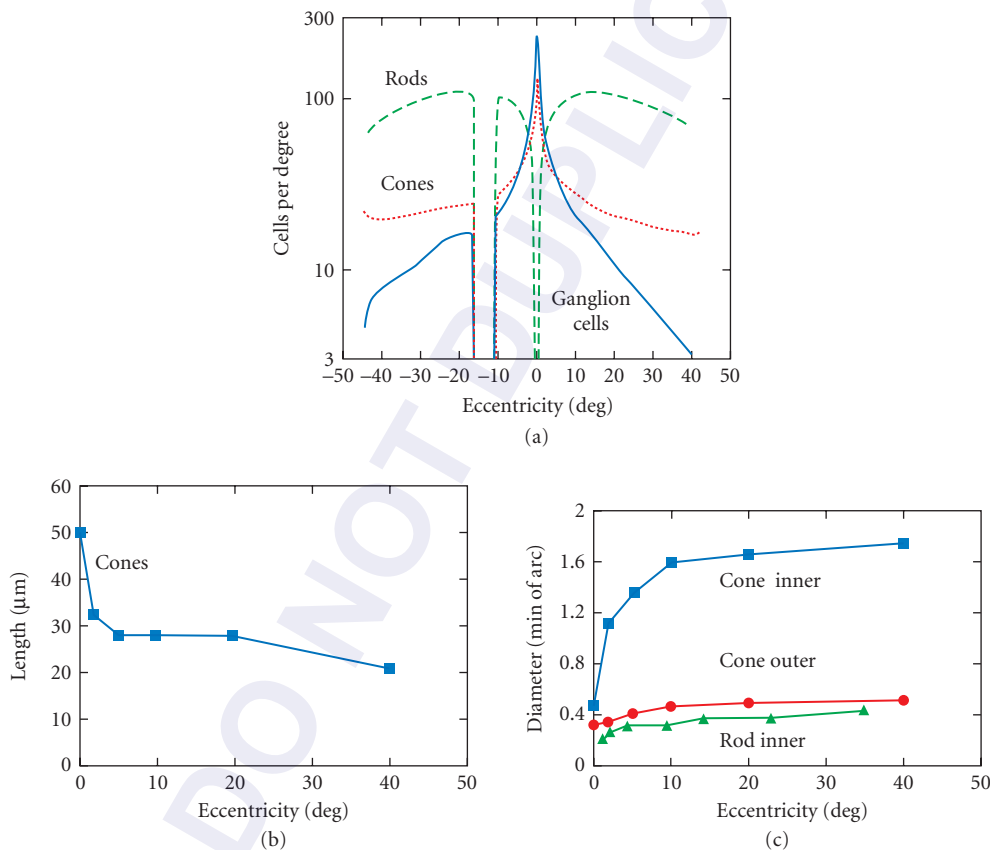


FIGURE 2 (a) Linear density of cones, rods, and ganglion cells as a function of eccentricity in the human retina. (The data were modified from Refs. 6 and 32.) Conversion from cells/mm^2 to $\text{cells}/\text{deg}^2$ was computed assuming a posterior nodal point 16.68 mm from the retina, and a retinal radius of curvature of 12.1 mm. Conversion to cells/deg was obtained by taking the square root of areal density. Ganglion cell density in the central 10 deg was derived assuming a 3:1 ratio of ganglion cells to cones in the fovea.³² (b) Human cone outer segment length. (Modified from Ref. 125.) (c) Human cone inner segment, cone outer segment, and rod diameter as a function of eccentricity. (Modified from Ref. 125.)

presentation of interference fringes (which avoid degradation by the eye's optics) yields visible spatial aliasing for spatial frequencies above the Nyquist limit.⁷⁻⁹

The densities and retinal distributions of the three cone types are quite different. The S cones form a rather regular lattice comprising less than 2 percent of cones in the central fovea and somewhat less than 10 percent of the cones elsewhere;¹⁰⁻¹² they may be absent in the central 20' to 25' of the fovea.¹³ It is much more difficult to distinguish individual L and M cones anatomically, so their densities and distributions are less certain. Psychophysical evidence indicates that the ratio of L to M cones is approximately 2:1,^{14,15} but the available physiological data in monkey suggest a ratio closer to 1:1.^{16,17}

From the Beer-Lambert law, the absorptance spectrum of a receptor depends upon the concentration of the photopigment, $c_o p$, in the receptor outer segment, the length of the outer segment, l , and the extinction spectrum, $\epsilon(\lambda)$, of the photopigment,

$$\alpha(\lambda) = 1 - 10^{-l c_o p \epsilon(\lambda)} \quad (8)$$

where c_o is the concentration of photopigment in the dark-adapted eye and p is proportion of unbleached photopigment. The specific absorptance spectra for the different classes of receptor are described in Chap. 11 by Andrew Stockman and David H. Brainard. At the peak wavelength, the photoreceptors absorb approximately 50 percent of the incident photons (although there are some variations with eccentricity).

Because of photopigment bleaching and regeneration processes within the photoreceptor, the absorption spectra of the photopigments change depending upon the history of photon absorptions. For example, prolonged exposure to high intensities depletes a significant fraction of the available photopigment, reducing the overall optical density. Reflection densitometry measurements¹⁸⁻²⁰ have shown, to first approximation, that the proportion p of available photopigment at a given point in time is described by a first-order differential equation:

$$\frac{dp}{dt} = \frac{n(\lambda)\xi(1 - 10^{-l c_o p \epsilon(\lambda)})}{l c_o} + \frac{1 - p}{t_o} \quad (9)$$

where $n(\lambda)$ is the photon-flux irradiance distribution ($\text{quanta} \cdot \text{sec}^{-1} \cdot \text{deg}^{-2} \cdot \text{nm}^{-1}$), t_o is the exponential time constant of regeneration, and ξ is the isomerization efficiency (which is believed to be near 1.0). For broadband light, Eq. (9) simplifies and can be expressed in terms of retinal illumination:¹⁹

$$\frac{dp}{dt} = -\frac{Ip}{Q_e} + \frac{(1 - p)}{t_o} \quad (10)$$

where I is the retinal illumination (in trolands*) and Q_e is the energy of a flash (in troland \cdot sec) required to reduce the proportion of unbleached photopigment to $1/e$. For the cone photopigments, the time constant of regeneration is approximately 120 sec, and Q_e (for broadband light) is approximately 2.4×10^6 troland \cdot sec. For rods, the time constant of regeneration is approximately 360 sec, and Q_e is approximately 1.0×10^7 scotopic troland \cdot sec. Equation (10) implies that the steady-state proportion of available photopigment is approximately

$$p(I) = \frac{I_o}{I_o + I} \quad (11)$$

where $I_o = Q_e/t_o$. (I_o is known as the *half-bleaching constant*.) Equation (11) is useful when computing photon absorptions at high ambient light levels; however, photopigment bleaching is not significant at low to moderate light levels. Equation (7) also implies that bleaching and regeneration can produce changes in the shapes of the absorptance spectra (see Chap. 10).

The lengths of the photoreceptor outer segments (and hence their absorptances) change with eccentricity. Figure 2b shows that the cone outer segments are longer in the fovea. Rod outer segments are approximately constant in length. If the dark-adapted photopigment concentration, c_o , is constant

*The troland is defined to be the retinal illumination produced by viewing a surface with a luminance of 1 cd/m^2 through a pupil with an area of 1 mm^2 .

in a given class of photoreceptor (which is likely), then increasing outer segment length increases the amount of light collected by the receptor (which increases signal-to-noise ratio).

The light collection area of the cones is believed to be approximately 70 to 90 percent of the cross-sectional area of the inner segment at its widest point (see Fig. 1), although direct measurements are not available (see Refs. 21 and 22). The collection area of the rod receptors is equal to the cross-sectional area of the inner and outer segments (which is the same). Figure 2c plots inner-segment diameter as a function of eccentricity for cones and rods. As can be seen, the cone aperture increases with eccentricity, while the rod aperture is fairly constant. Increasing the cone aperture increases the light collected (which increases signal-to-noise ratio), but slightly reduces contrast sensitivity at high spatial frequencies (see later).

The data and formulas above [Eqs. (6) to (10)] can be used to compute (approximately) the total effective photons absorbed per second, N , in any receptor, at any eccentricity:

$$N = \int \kappa \xi \alpha(\lambda) n(\lambda) d\lambda \quad (12)$$

where κ is the collection area, or aperture, of the receptor (in deg^2).

Photons entering through the pupillary center have a greater chance of being absorbed in the cones than photons entering through the pupillary margins. This phenomenon, known as the Stiles-Crawford effect, can be included in the above calculations by modifying the pupil function [Eqs. (3) and (4); also see Chaps. 1 and 8]. We also note that the cones are oriented toward the exit pupil.

The temporal response properties of photoreceptors are more difficult to measure than their spatial integration properties. Figure 3 shows some physiological measurements of the photocurrent

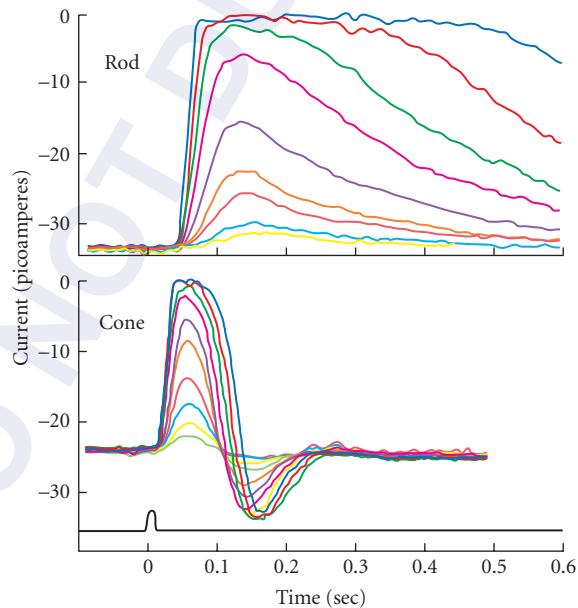


FIGURE 3 Photocurrent responses of macaque rod and cone photoreceptors to flashes of light. Each trace represents the response to a different intensity level. The flash intensities were varied by factors of two. (From Ref. 23. Copyright © 1989 by Scientific American, Inc. All rights reserved.)

responses of primate (macaque) rods and cones at low and moderate light levels.²³ Rods integrate photons over a substantially longer duration than cones do. In addition, rods produce reliable responses to single photon absorptions, but cones do not.

Photoreceptors, like all neurons, have a limited dynamic range. Direct recordings from receptors in macaque have shown that the peak responses of photoreceptors to brief flashes of light are adequately described by either a modified Michaelis-Menton function

$$f(z) = \frac{J_o z^n}{z^n + \sigma_o^n} \quad (13)$$

or an exponential function

$$f(z) = J_o (1 - 2^{-z^n / \sigma_o^n}) \quad (14)$$

where J_o is the maximum photocurrent, σ_o is the half-saturation constant (the value of z that produces exactly half the maximum response), and n is an exponent that has a value near 1.0.^{17,24,25} As the above equations imply, the photoreceptors respond approximately linearly to transient intensity changes at low to moderate light levels, but respond nonlinearly (ultimately reaching saturation) at higher light levels. Recent electroretinogram (ERG) measurements suggest that the flash response of human photoreceptors is similar to those of the macaque.²⁶

The nonlinear response saturation suggests that the receptors should not be able to encode high image intensities accurately (because any two intensities that saturate receptor response will be indistinguishable). However, the cones avoid saturation effects under many circumstances by adjusting their gain (σ_o) depending upon the ambient light level; the rods do not adjust their gain nearly as much. The gain adjustment is accomplished in a few ways. Photopigment depletion (see earlier) allows multiplicative gain changes, but only operates at very high ambient light levels. Faster-acting mechanisms in the phototransduction sequence are effective at moderate to high light levels.²⁷ At this time, there remains some uncertainty about how much gain adjustment (adaptation) occurs within the cones.^{17,24,28}

Retinal Processing

The information encoded by the photoreceptors is processed by several layers of neurons in the retina. The major classes of retinal neuron in the primate (human) visual system and their typical interconnections are illustrated schematically in Fig. 1. Although the retina has been studied more extensively than any other part of the nervous system, its structure and function are complicated and not yet fully understood. The available evidence suggests that the primate retina is divided into three partially overlapping neural pathways: a *rod pathway*, a *cone parvo pathway*, and a *cone magno pathway* (for reviews see Refs. 29–31). These pathways, illustrated in Fig. 4, carry most of the information utilized in high-level visual processing tasks. There are other, smaller pathways (which will not be described here) that are involved in functions such as control of the pupil reflex, accommodation, and eye movements. Although Fig. 4 is based upon the available evidence, it should be kept in mind that there remains considerable uncertainty about some of the connections.

Retinal Anatomy As indicated in Fig. 4, the photoreceptors (R, C) form electrical synapses (gap-junctions) with each other, and chemical synapses with *horizontal cells* (H) and *bipolar cells* (RB, MB+, MB–, DB+, DB–). The electrical connections between receptors are noninverting* and hence produce simple spatial and temporal summation. The dendritic processes and the axon-terminal processes of the horizontal cells provide spatially extended, negative-feedback connections from cones onto cones and from rods onto rods, respectively. It is likely that the horizontal cells play an important role in creating the surround response of ganglion cells. Rod responses probably do not influence cone responses (or vice versa) via the horizontal cells.

*By *noninverting* we mean that changes in the response of the presynaptic neuron, in a given direction, produce changes in the response of the postsynaptic neuron in the same direction.

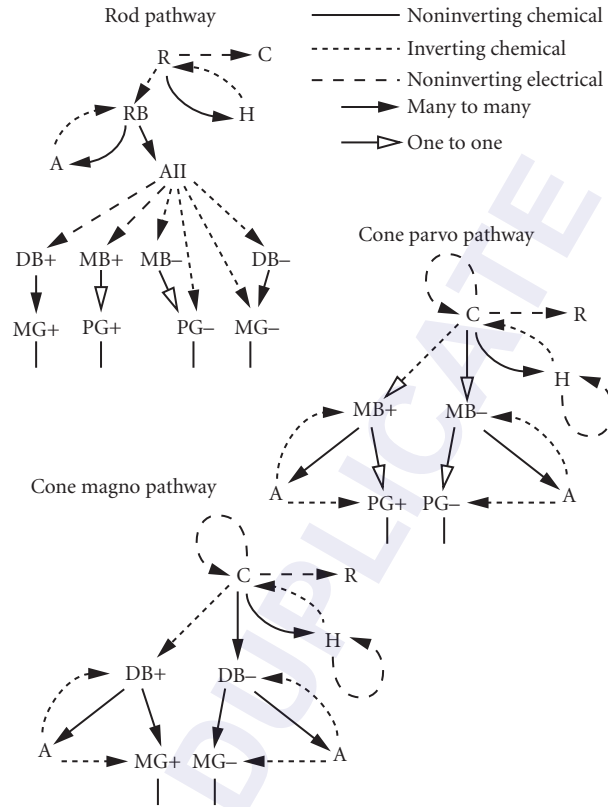


FIGURE 4 Schematic diagram of neural connections in the primate retina. R (rod), C (cone), H (horizontal cell), RB (rod bipolar), MB+ (midget on-bipolar), MB- (midget off-bipolar), DB+ (diffuse on-bipolar), DB- (diffuse off-bipolar), A (amacrine), AII (rod amacrine), PG+ (P or midget on-center ganglion cell), PG- (P or midget off-center ganglion cell), MG+ (M or parasol on-center ganglion cell), MG- (M or parasol off-center ganglion cell).

In the rod pathway, rod bipolar cells (RB) form noninverting chemical synapses with *amacrine* cells (AII, A). The rod amacrine cells (AII) form electrical (noninverting) synapses with the on-center bipolar cells (MB+, DB+), and form chemical (inverting) synapses with off-center bipolar cells (MB-, DB-) and with off-center ganglion cells (PG-, MG-). Other amacrine cells (A) form reciprocal negative-feedback connections to the rod bipolar cells.

In the cone pathways, on-bipolar cells (MB+, DB+) form chemical synapses with amacrine cells (A) and with on-center *ganglion* cells (PG+, MG+).* Off-bipolar cells (FMB, FDB) form chemical synapses with amacrine cells and off-center ganglion cells (PG-, MG-).

The ganglion cells are the output neurons of the retina; their myelinated axons form the optic nerve. In the fovea and parafovea, each P ganglion cell (PG+, PG-) forms synapses with only one midget bipolar cell (MB+, MB-), and each midget bipolar cell forms synapses with only one cone;

*The parvo ganglion cells and magno ganglion cells are also referred to as *midget ganglion cells* and *parasol ganglion cells*, respectively.

however, each cone makes contact with a midget on-bipolar and a midget-off bipolar. Thus, the P pathway (in the fovea) is able to carry fine spatial information. Current evidence also suggests that the P pathway carries most of the chromatic information (see Chap. 11). The M ganglion cells (MG+, MG−) form synapses with diffuse bipolar cells (DB+, DB−), and each diffuse bipolar cell forms synapses with many cones. Thus, the M pathway is less able to carry fine spatial information. However, the larger spatial summation (and other boosts in gain) give the M pathway greater sensitivity to changes in contrast. Furthermore, the larger sizes of neurons in the M pathway provide somewhat faster responses and hence somewhat faster transmission of information.

Figure 2a shows the density of ganglion cells as a function of eccentricity in the human retina;³² for macaque retina, see Ref. 33. At each eccentricity, the P cells (PG+, PG−) comprise approximately 80 percent of the ganglion cells and the M cells (MG+, MG−) approximately 10 percent. Half the P and M cells are on-center (PG+ and MG+) and half are off-center (PG− and MG−). Thus, the approximate densities of the four major classes of ganglion cells can be obtained by scaling the curve in Fig. 2a. As can be seen, ganglion cell density decreases more quickly with eccentricity than does cone density, but in the fovea there are at least 2 to 3 ganglion cells per cone (a reasonable assumption is that there are two P cells for every cone). The Nyquist limit of the retina would appear to be set by the cone density in the fovea and by the ganglion cell density in the periphery.

Retinal Physiology Ganglion cells transmit information to the brain via action potentials propagating along axons in the optic nerve. The other retinal neurons transmit information as graded potentials (although amacrine cells generate occasional action potentials). Much remains to be learned about exactly how the computations evident in the responses of ganglion cells are implanted within the retinal circuitry.*

The receptive fields[†] of ganglion cells are approximately circular and, based upon their responses to spots of light, can be divided into a center region and an antagonistic, annular, surround region.³⁴ In on-center ganglion cells, light to the center increases the response, whereas light to the surround suppresses the response; the opposite occurs in off-center ganglion cells.

For small to moderate contrast modulations around a steady mean luminance, most P and M cells respond approximately linearly, and hence their spatial and temporal response properties can be usefully characterized by a spatiotemporal transfer function. The spatial components of the transfer function are adequately described as a difference of gaussian functions with separate space constants representing the center and surround;^{35–37} see Fig. 5a. The available data suggest that the surround diameter is typically 3 to 6 times the center diameter.³⁷ The temporal components of the transfer function have been described as a cascade of simple feedback and feed-forward linear filters (e.g., Ref. 38; see Fig. 5b). There also appears to be a small subset of M cells that is highly nonlinear, similar to the Y cells in cat.³⁹

The response amplitude of M and P cells as a function of sinewave contrast is reasonably well described by a Michaelis-Menton function [i.e., by Eq. (14)], but where z and σ_0 are contrasts rather than intensities. As indicated by Fig. 5, the M ganglion cells are (over most spatial and temporal frequencies) about 4 to 10 times more sensitive than the P ganglion cells; in other words, σ_0 is 4 to 10 times smaller in M cells than P cells.

When mean (ambient) light level increases, the high-frequency falloffs of the temporal transfer functions of M and P cells shift toward higher frequencies, corresponding to a decrease in the time constants of the linear filters.^{38,40} The effect is reduction in gain (an increase in σ_0), and an increase in temporal resolution.

Relatively little is known about how the spatial transfer functions of primate ganglion cells change with mean light level. In cat, the relative strength of the surround grows with mean luminance, but the space constants (sizes) of the center and surround components appear to change relatively little.^{41,42}

*In this subsection we describe the electrophysiology of M and P cells. This description is based on a composite of data obtained from ganglion cells and geniculate cells. The response properties of M and P ganglion cells and M and P geniculate cells are very similar.

[†]The receptive field of a neuron is defined to be the region of the visual field (or, equivalently, of the receptor array) where light stimulation has an effect on the response of the neuron.

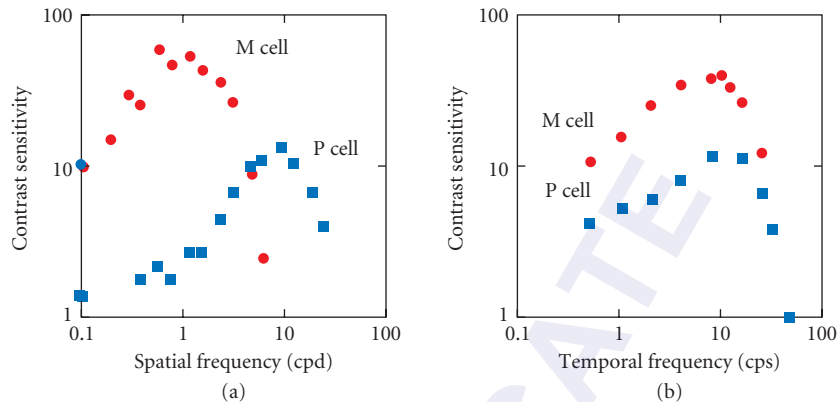


FIGURE 5 Typical spatial and temporal contrast sensitivity functions of P and M ganglion cells. (a) Spatial contrast sensitivity functions for sinusoidal gratings with a mean intensity of 1400 td, drifting at 5.2 cps. The symbols are modulation thresholds for a uniform field flickering sinusoidally at 5.2 cps. The threshold criterion was approximately 10 extra impulses per second. (Adapted from Ref. 37.) (b) Temporal contrast sensitivity functions for drifting sinusoidal gratings with mean intensities of 4600 td (P cell) and 3600 td (M cell), and spatial frequencies of 3 cpd (P cell) and 1.6 cpd (M cell). The threshold criterion was approximately 10 extra impulses per second. (Adapted from Ref. 38.)

The major effect of increasing adaptation luminance is an increase in the magnitude of the low-frequency falloff of the spatial transfer function.

The center diameters of both P and M cells increase with eccentricity,^{42,43} roughly in inverse proportion to the square root of ganglion-cell density (cf., Fig. 2a). However, the precise relationship has not been firmly established because of the difficulty in measuring center diameters of ganglion cells near the center of the fovea, where the optical point-spread function is likely to be a major component of center diameter.³⁷

Central Visual Processing

Most of the information encoded in the retina is transmitted via the optic nerve to the lateral geniculate nucleus (LGN) in the thalamus. The neurons in the LGN then relay the retinal information to the primary visual cortex (V1). Neurons in V1 project to a variety of other visual areas. Less is known about the structure and function of V1 than of the retina or LGN, and still less is known about the subsequent visual areas.

Central Anatomy The LGN is divided into six layers (see Fig. 6); the upper four layers (the parvocellular laminae) receive synaptic input from the P ganglion cells, the lower two layers (the magnocellular laminae) receive input from the M ganglion cells. Each layer of the LGN receives input from one eye only, three layers from the left eye and three from the right eye. The LGN also receives input from other brain areas including projections from the reticular formation and a massive projection (of uncertain function) from layer 6 of V1. The total number of LGN neurons projecting to the visual cortex is slightly larger than the number of ganglion cells projecting to the LGN.

The segregation of M and P ganglion cell afferents into separate processing streams at the LGN is preserved to some extent in subsequent cortical processing. The magnocellular neurons of the LGN project primarily to neurons in layer 4c α in V1, which project primarily to layer 4b. The neurons in layer 4b project to several other cortical areas including the middle-temporal area (MT) which appears to play an important role in high-level motion processing (for reviews see Refs. 44–46).

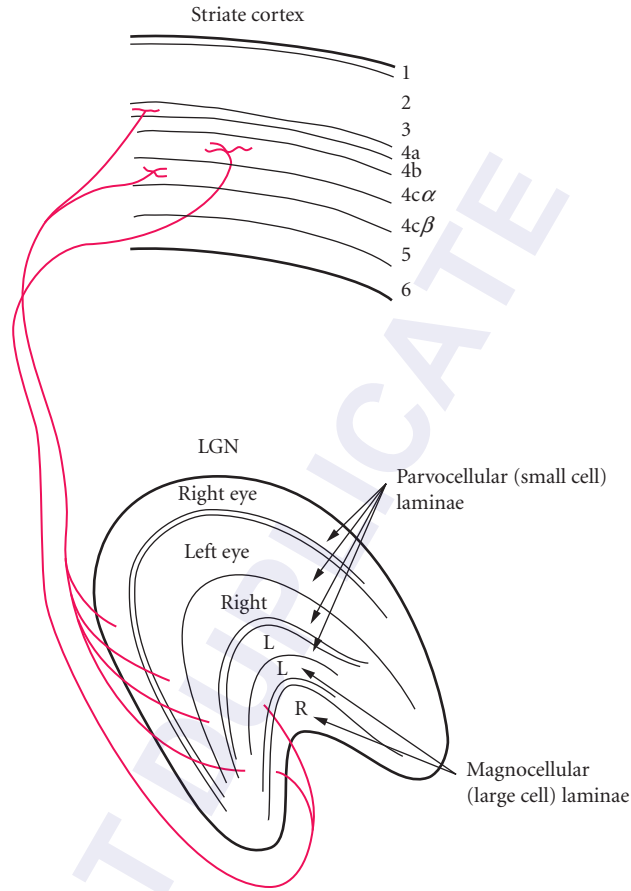


FIGURE 6 Schematic diagram of a vertical section through the lateral geniculate nucleus and through striate cortex (V1) of the primate visual system. (Reproduced from Ref. 49.)

The parvocellular neurons of the LGN project to layers 4a and 4c β in V1. The neurons in layers 4a and 4c β project to the superficial layers of V1 (layers 1, 2, and 3), which project to areas V2 and V3. Areas V2 and V3 send major projections to area V4 which sends major projections to IT (infero-temporal cortex). The superficial layers of V1 also send projections to layers 5 and 6.

It has been hypothesized that the magno stream subserves crucial aspects of motion and depth perception and that the parvo stream subserves crucial aspects of form and color perception. However, current understanding of the areas beyond V1 is very limited; it is likely that our views of their functional roles in perception will change substantially in the near future. For reviews of cortical functional anatomy see, for example, Refs. 46–49.

Central Physiology The receptive field properties of LGN neurons are quite similar to those of their ganglion-cell inputs.^{50–52} LGN neurons have, on average, a lower spontaneous response rate than ganglion cells, and are more affected by changes in alertness and anesthetic level, but otherwise display similar center/surround organization, and similar spatial, temporal, and chromatic response properties (however, see Ref. 53).

The receptive-field properties of neurons in the primary visual cortex are substantially different from those in the retina or LGN; V1 neurons have elongated receptive fields that display a substantial degree of selectivity to the size (spatial frequency), orientation, direction of motion of retinal stimulation, and binocular disparity.^{54–56} Thus, the probability that a cortical neuron will be activated by an arbitrary retinal image is much lower than for retinal and LGN neurons.

Each region of the visual field is sampled by cortical neurons selective to the full range of sizes, orientations, directions of motion, and binocular disparities. When tested with sinewave gratings, cortical neurons have a spatial-frequency bandwidth (full width at half height) of 1 to 1.5 octaves and an orientation bandwidth of 25 to 35 degrees.^{57,58} The number of neurons in area V1 is more than two orders of magnitude greater than the number of LGN neurons;⁵⁹ thus, there are sufficient numbers of cortical neurons in each region of the visual field for their receptive fields to tile the whole spatial-frequency plane transmitted from retina to cortex several times. However, the number of cortical cells encoding each region of the spatial-frequency plane is still uncertain. The temporal frequency tuning of V1 neurons is considerably broader than the spatial-frequency tuning with peak sensitivities mostly in the range of 5 to 10 cps.⁶⁰ The direction selectivity of V1 neurons varies from 0 to 100 percent with an average of 50 to 60 percent.^{57*} Early measurements of disparity tuning in primate suggested that cortical neurons fall into three major categories: one selective to crossed disparities, one selective to uncrossed disparities, and one tuned to disparities near zero.⁵⁶ Recent evidence⁶¹ in agreement with evidence in cat cortex^{62,63} suggests a more continuous distribution of disparity tuning.

The chromatic response properties of V1 neurons are not yet fully understood.⁴⁹ There is some evidence that discrete regions in the superficial layers of V1, called cytochrome oxidase “blobs,” contain a large proportion of neurons responsive to chromatic stimuli;⁶⁴ however, see Ref. 65.

The spatial-frequency tuning functions of cortical cells have been described by a variety of simple functions including Gabor filters, derivatives of gaussian filters, and log Gabor filters (e.g., Ref. 66). The temporal-frequency tuning functions have been described by difference of gamma filters. The full spatiotemporal tuning (including the direction-selective properties of cortical cells) has been described by quadrature (or near quadrature) pairs of separable spatiotemporal filters with spatial and temporal components drawn from the types listed above (e.g., Ref. 67).

Essentially all V1 neurons display nonlinear response characteristics, but can be divided into two classes based upon their nonlinear behavior. *Simple cells* produce approximately half-wave rectified responses to drifting or counterphase sinewave gratings; *complex cells* produce a large unmodulated (DC) response component with a superimposed half-wave or full-wave rectified response component.⁶⁸ Simple cells are quite sensitive to the spatial phase or position within the receptor field; complex cells are relatively insensitive to position within the receptive field.⁵⁴

Most simple and complex cells display an accelerating response nonlinearity at low contrasts and response saturation at high contrasts, which can be described by a Michaelis-Menton function [Eq. (11)] with an exponent greater than 1.0.⁶⁹ These nonlinearities impart the important property that response reaches saturation at approximately the same physical contrast independent of the spatial frequency, orientation, or direction of motion of the stimulus. The nonlinearities sharpen the spatiotemporal tuning of cortical neurons while maintaining that tuning largely independent of contrast, even though the neurons often reach response saturation of 10 to 20 percent contrast.^{70,71}

Unlike retinal ganglion cells and LGN neurons, cortical neurons have little or no spontaneous response activity. When cortical cells respond, however, they display noise characteristics similar to retinal and LGN neurons; specifically, the variance of the response is approximately proportional to the mean response;^{72,73} unlike Poisson noise, the proportionality constant is usually greater than 1.0.

2.4 VISUAL PERFORMANCE

A major goal in the study of human vision is to relate performance—for example, the ability to see fine detail—to the underlying anatomy and physiology. In the sections that follow, we will discuss some of the data and theories concerning the detection, discrimination, and estimation of contrast,

*Direction selectivity is defined as $100 \times (R_p - R_n)/R_p$, where R_p is the magnitude of response in the preferred direction and R_n is the magnitude of response in the nonpreferred direction.

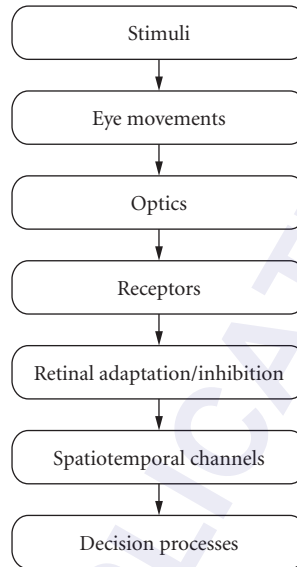


FIGURE 7 Simple information-processing model of the visual system representing the major factors affecting contrast detection and discrimination.

position, shape, motion, and depth by human observers. In this discussion, it will be useful to use as a theoretical framework the information-processing model depicted in Fig. 7. The visual system is represented as a cascade of processes, starting with eye movements, that direct the visual axis toward a point of interest, proceeding through a number of processing stages discussed in the previous sections and in Chaps. 1 and 11, and ending with a decision. When a human observer makes a decision about a set of visual stimuli, the accuracy of this decision depends on all the processes shown in Fig. 7. Psychophysical experiments measure the performance of the system as a whole, so if one manipulates a particular variable and finds it affects observers' performance in a visual task, one cannot readily pinpoint the responsible stage(s). The effect could reflect a change in the amount of information the stimulus provides, in the fidelity of the retinal image, in representation of the stimulus among spatiotemporal channels, or in the observer's decision strategy. Obviously, it is important to isolate effects due to particular processing stages as well as possible.

We begin by discussing the decisions a human observer might be asked to make concerning visual stimuli. Humans are required to perform a wide variety of visual tasks, but they can be divided into two general categories: *identification tasks and estimation tasks*.^{*} In an identification task, the observer is required to identify a visual image, sequence of images, or a part of an image, as belonging to one of a small number of discrete categories. An important special case is the *discrimination task* in which there are only two categories (e.g., Is the letter on the TV screen an *E* or an *F*?). When one of the two categories is physically uniform, the task is often referred to as a *detection task* (e.g., Is there a letter *E* on the screen or is the screen blank?). In the estimation task, the observer is required to estimate the value of some property of the image (e.g., What is the height of the letter?). The distinction between estimation and identification is quantitative, not qualitative; strictly speaking, an estimation task can

^{*}Visual tasks can also be categorized as *objective* or *subjective*. Objective tasks are those for which there is, in principle, a precise physical standard against which performance can be evaluated (e.g., a best or correct performance). The focus here is on objective tasks because performance in those tasks is more easily related to the underlying physiology.

be regarded as an identification task with a large number of categories. There are two fundamental measures of performance in visual tasks: the *accuracy* with which the task is performed, and the *speed* at which the task is performed. For further discussion of methods for measuring visual performance, see Chap. 3 by Denis G. Pelli and Bart Farell.

Referring again to the information-processing model of Fig. 7, we note that it is difficult to isolate the decision strategy; the best one can do is to train observers to adhere as closely as possible to the same decision strategy across experimental conditions. So, whether the experiment involves simple identification tasks, such as discrimination or detection procedures, or estimation tasks, observers are generally trained to the point of using a consistent decision strategy.

The stimulus for vision is distributed over space, time, and wavelength. As mentioned earlier, we are primarily concerned here with spatial and temporal variations in intensity; variations in wavelength are taken up in Chap. 11.

Before considering human visual performance, we briefly describe *ideal-observer theory* (e.g., Refs 74 and 75), which has proven to be a useful tool for evaluating and interpreting visual performance.

Ideal-Observer Theory

In attempting to understand human visual performance, it is critically important to quantify performance limitations imposed by the information contained in the stimulus. There is a well-accepted technique, based on the theory of ideal observers,⁷⁵ for so doing. As applied to vision (see Refs. 76–79), this approach is based on comparing human observers' performance to that of an ideal observer for the same stimuli. Ideal observers have three main elements: (1) a precise description of the stimuli, including any random variation or noise; (2) descriptions of how the stimuli are modified by the visual stages of interest (e.g., optics, receptors); and (3) an optimal decision rule. Such observers specify the best possible level of performance (e.g., the smallest detectable amount of light) given the variability in the stimulus and information losses in the incorporated visual stages. Because ideal observer performance is the best possible, the thresholds obtained are a measure of the information available in the stimulus (to perform the task) once processed by the incorporated stages. Poorer performance by the human observer must be due to information losses in the unincorporated stages. Thus, ideal-observer theory provides the appropriate physical benchmark against which to evaluate human performance.^{75,80} One can, of course, use models of anatomical and physiological mechanisms to incorporate various processing stages into an ideal observer (as shown in Fig. 7); comparisons of human and ideal observer performance in such cases can allow one to compute precisely the information transmitted and lost by mechanisms of interest. This allows assessment of the contributions of anatomical and physiological mechanisms to overall visual performance.⁷⁸

It is important to recognize that ideal-observer theory is not a theory of human performance (humans generally perform considerably below ideal); thus, the theory is not a substitute for psychophysical and physiological modeling. However, the theory is crucial for understanding the stimuli and the task. In many ways, measuring and reporting the information content of stimuli with an ideal observer is as fundamental as measuring and reporting the basic physical dimensions of stimuli.

To illustrate the concepts of ideal-observer theory, consider an identification task with response alternatives (categories) α_1 through α_m , where the probability of a stimulus from category α_i equals q_i . Suppose, further, that the stimuli are static images, with an onset at some known point in time. (By a static image, we mean there is no temporal variation within a stimulus presentation except for photon noise.) Let \mathbf{Z} be the list of random values (e.g., photon counts) at each sample location (e.g., photoreceptor) in some time interval τ for a single presentation of a stimulus:

$$\mathbf{Z} = (Z_1, \dots, Z_n) \quad (15)$$

where Z_p is the sample value at the i th location.

Overall accuracy in an identification task is always optimized (on average) by picking the most probable stimulus category given the sample data and the prior knowledge. Thus, if the goal is to

maximize overall accuracy in a given time interval τ , then the maximum average percent correct (PC_{opt}) is given by the following sum:

$$PC_{\text{opt}}(\tau) = \sum_z q_* p_\tau(\mathbf{z}|\alpha_*) \quad (16)$$

where $p_\tau(\mathbf{z}|\alpha_*)$ is the probability density function associated with \mathbf{Z} for stimulus category α_* , and the subscript $*$ represents the category j for which the product $q_* p_\tau(\mathbf{z}|\alpha_*)$ is maximum. (Note, \mathbf{Z} is a random vector and \mathbf{z} is a simple vector.)

Suppose, instead, that the goal of the task is to optimize speed. Optimizing speed at a given accuracy level is equivalent to finding the minimum stimulus duration, $\tau_{\text{opt}}(\epsilon)$, at a criterion error rate ϵ . The value of $\tau_{\text{opt}}(\epsilon)$ is obtained by setting the left side of Eq. (16) to the criterion accuracy level $(1 - \epsilon)$ and then solving for τ .

For the detection and discrimination tasks, where there are two response categories, Eq. (16) becomes

$$PC_{\text{opt}}(\tau) = \frac{1}{2} + \frac{1}{2} \sum_z |q_1 p_\tau(\mathbf{z}|\alpha_1) - (1 - q_1) p_\tau(\mathbf{z}|\alpha_2)| \quad (17)$$

[The summation signs in Eqs. (16) and (17) are replaced by an integral sign if the probability density functions are continuous.]

Because the sums in Eqs. (16) and (17) are over all possible vectors $\mathbf{z} = (z_1, \dots, z_n)$, they are often not of practical value in computing optimal performance. Indeed, in many cases there is no practical analytical solution for ideal-observer performance, and one must resort to Monte Carlo simulation.

Two special cases of ideal-observer theory have been widely applied in the analysis of psychophysical discrimination and detection tasks. One is the ideal observer for discrimination or detection tasks where the only source of stimulus variability is photon noise.^{77,81,82} In this case, optimal performance is given, to close approximation, by the following formulas:^{83,84}

$$PC_{\text{opt}}(\tau) = q + 1(1 - q)\Phi\left(\frac{c}{d'} + \frac{d'}{2}\right) - q\Phi\left(\frac{c}{d'} - \frac{d'}{2}\right) \quad (18)$$

where $c = \ln[(1 - q)/q]$, $\Phi(\cdot)$ is the cumulative standard normal probability distribution, and

$$d' = \frac{\tau^{1/2} \sum_{i=1}^n (b_i - a_i) \ln(b_i/a_i)}{\left[\sum_{i=1}^n (b_i + a_i) \ln^2(b_i/a_i) \right]^{1/2}} \quad (19)$$

In Eq. (19), a_i and b_i are the average numbers of photons per unit time at the i th sample location, for the two alternative stimuli (α_1 and α_2). For equal presentation probabilities ($q = 0.5$), these two equations are easily solved to obtain $\tau_{\text{opt}}(\epsilon)$.^{85,86} In signal-detection theory⁷⁵ the quantity d' provides a criterion-independent measure of signal detectability or discriminability.

Figure 8 provides an example of the use of ideal-observer theory. It compares human and ideal performance for detection of square targets of different areas presented briefly on a uniform background. The symbols show the energy required for a human observer to detect the target as a function of target area. The lower curve shows the absolute performance for an ideal observer operating at the level of photon absorption in the receptor photopigments; the upper curve shows the same curve shifted vertically to match the human data. The differences between real and ideal performance represent losses of information among neural processes. The data show that neural efficiency, $[d'_{\text{real}}]^2/[d'_{\text{ideal}}]^2$, for the detection of square targets is approximately 1/2 percent for all but the largest target size.

Another frequently applied special case is the ideal observer for detection and discrimination tasks where the signal is known exactly (SKE) and white or filtered gaussian noise (i.e., image or pixel noise) has been added. Tasks employing these stimuli have been used to isolate and measure central mechanisms that limit discrimination performance^{80,87,88} to evaluate internal (neural) noise levels in

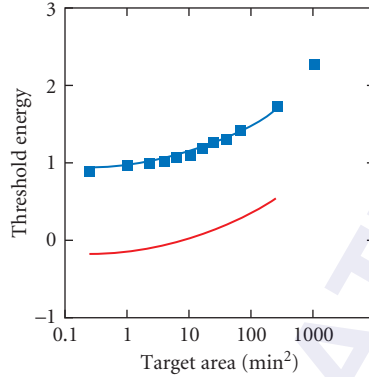


FIGURE 8 Comparison of real and ideal performance for detection of square targets on a uniform background of 10 cd/m^2 (viewed through a 3-mm artificial pupil). The symbols represent human thresholds as a function of target area. Threshold energy is in units of $\text{cd/m}^2 \times \text{min}^2 \times \text{sec}$. The lower curve is ideal performance at the level of the photoreceptors; the upper curve is the ideal performance shifted vertically to match the human data. (Adapted from Ref. 264.)

the visual system,^{79,89} and to develop a domain of applied psychophysics relevant to the perception of the noisy images, such as those created by radiological devices and image enhancers.⁹⁰⁻⁹² For the gaussian-noise-limited ideal discriminator, Eq. (18) still applies, but d' is given by the following:

$$d' = \frac{E(L|\alpha_2) - E(L|\alpha_1)}{\sqrt{\text{VAR}(L)}} \quad (20)$$

where

$$L = [\mathbf{m}_2 - \mathbf{m}_1]' \Sigma^{-1} \mathbf{Z} \quad (21)$$

In Eq. (21), \mathbf{Z} is the column vector of random values from the sample locations (e.g., pixels), $[\mathbf{m}_2 - \mathbf{m}_1]'$ is a row vector of the differences in the mean values at each sample point (i.e., $b_1 - a_1, \dots, b_n - a_n$), and Σ is the covariance matrix resulting from the gaussian noise (i.e., the element σ_{ij} of Σ is the covariance between the i th and j th sample location). In the case of white noise, Σ is a diagonal matrix with diagonal elements equal to σ^2 , and thus,*

$$d' = \frac{1}{\sigma} \sqrt{\sum (b_i - a_i)^2} \quad (22)$$

As an example, Fig. 9 compares human and ideal performance for amplitude discrimination of small targets in white noise as a function of noise spectral power density (which is proportional to σ^2). The solid black line of slope 1.0 shows the absolute performance of an ideal observer operating at the level of the cornea (or display screen); thus, the difference between real and ideal performance represents all losses of information within the eye, retina, and central visual pathways. For these conditions, efficiency for amplitude discrimination of the targets ranges from about 20 to 70 percent, much

*It should be noted that the performance of the photon-noise-limited observer [Eq. (19)] becomes equivalent to an SKE white noise observer [Eq. (22)] for detection of targets against intense uniform backgrounds.

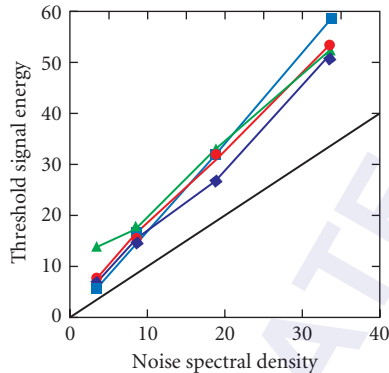


FIGURE 9 Comparison of real (symbols) and ideal (solid black line) performance for contrast discrimination of spatially localized targets in white noise backgrounds with a mean luminance of 154 cd/m^2 . Noise spectral density and target signal energy are in units of 10^{-7} deg^2 (the arrow indicates a noise standard deviation of 26 percent contrast per pixel). (■) Gaussian target (standard deviation = 0.054 deg); (●) Gaussian damped sinewave target (9.2 cpd , standard deviation = 0.109 deg); (▲) Gaussian damped sinewave target (4.6 cpd , standard deviation = 0.217 deg); (◆) Sinewave target (4.6 cpd , $0.43 \times 0.43 \text{ deg}$). (Adapted from Ref. 87.)

higher than for detecting targets in uniform backgrounds (cf., Fig. 8). Efficiencies are lower in Fig. 8 than in Fig. 9 primarily because, with uniform backgrounds, internal (neural) noise limits human, but not ideal, performance; when sufficient image noise is added to stimulus, it begins to limit both real and ideal performance.

The following sections are a selective review of human visual performance and underlying anatomical and physiological mechanisms. The review includes spatial and temporal contrast perception, light adaptation, visual resolution, motion perception, and stereo depth perception, but due to space limitations, excludes many other active areas of research such as color vision (partially reviewed Chaps. 10 and 11), eye movements,⁹³ binocular vision,^{94,95} spatial orientation and the perception of layout,^{96,97} pattern and object recognition,⁹⁶⁻⁹⁹ visual attention,¹⁰⁰ visual development,^{101,102} and abnormal vision.¹⁰³

Contrast Detection

The study of contrast detection in the human visual system has been dominated for the last 20 years by methods derived from the concepts of linear systems analysis. The rationale for applying linear systems analysis in the study of visual sensitivity usually begins with the observation that a complete characterization of a visual system would describe the output resulting from any arbitrary input. Given that the number of possible inputs is infinite, an exhaustive search for all input-output relationships would never end. However, if the system under study is linear, then linear systems analysis provides a means for characterizing the system from the measurement of a manageably small set of input-output relationships. Of course, the visual system has important nonlinearities, so strictly speaking, the assumptions required by linear systems analysis are violated in general. Nonetheless, there have been several successful applications of linear systems analysis in vision.

In the case of spatial vision, linear systems analysis capitalizes on the fact that any spatial retinal illumination distribution, $I(x, y)$ can be described exactly by the sum of a set of basis functions, such as the set of sinusoids. Sinusoids are eigenfunctions of a linear system, which implies that the system response to a sinusoidal input can be completely characterized by just two numbers; an amplitude change and a phase change. Indeed, it is commonly assumed (usually with justification) that spatial phase is unaltered in processing, so only one number is actually required to describe the system response to a sinusoid. Linear systems analysis provides a means for predicting, from such characterizations of system responses to sinusoids, the response to any arbitrary input. For these reasons, the measurement of system responses to spatial sinusoids has played a central role in the examination of human spatial vision. We will begin our discussion with the spatial sinusoid and the system response to it.

A spatial sinusoid can be described by

$$I(x, y) = A \sin[2\pi f(x \cos(\theta) + y \sin(\theta)) + \phi] + \bar{I} \quad (23)$$

where \bar{I} is the space-average intensity of the stimulus field (often expressed in trolands), A is the amplitude of the sinusoid, f is the spatial frequency (usually expressed in cycles \cdot deg $^{-1}$ or cpd), θ is the orientation of the pattern relative to axes x and y , and ϕ is the phase of the sinusoid with respect to the origin of the coordinate system. The contrast C of a spatial sinewave is defined as $(I_{\max} - I_{\min}) / (I_{\min} + I_{\min})$, where I_{\max} is the maximum intensity of the sinusoid and I_{\min} is the minimum; thus $C = A / \bar{I}$. When the sinusoid is composed of vertical stripes, θ is zero and Eq. (23) reduces to

$$I(x, y) = A \sin[2\pi f x + \phi] + \bar{I} \quad (24)$$

When the contrast C of a sinusoidal grating (at a particular frequency) is increased from zero (while holding \bar{I} fixed), there is a contrast at which it first becomes reliably detectable, and this value defines the contrast detection threshold. It is now well-established that contrast threshold varies in a characteristic fashion with spatial frequency. A plot of the reciprocal of contrast at threshold as a function of spatial frequency constitutes the *contrast sensitivity function* (CSF). The CSF for a young observer with good vision under typical indoor lighting conditions is shown in Fig. 10 (blue squares).

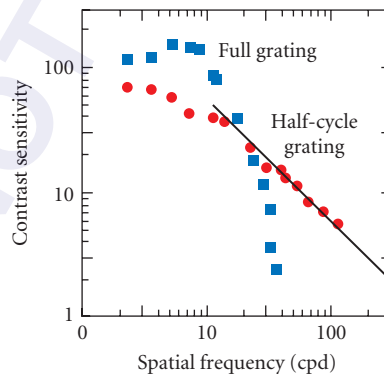


FIGURE 10 Contrast sensitivity as a function of spatial frequency for a young observer with good vision under typical indoor lighting conditions. The blue squares represent the reciprocals of contrast at detection threshold for sinusoidal grating targets. The red circles represent the reciprocals of contrast thresholds for half-cycle sinusoids. The solid line represents the predicted threshold function for half-cycle gratings for a linear system whose CSF is given by the blue squares. (Adapted from Ref. 193.)

The function is bandpass with a peak sensitivity at 3 to 5 cpd. At those spatial frequencies, a contrast of roughly 1/2 percent can be detected reliably (once the eye has adapted to the mean luminance). At progressively higher spatial frequencies, sensitivity falls monotonically to the so-called high-frequency cutoff at about 50 cpd; this is the finest grating an observer can detect when the contrast C is at its maximum of 1.0. At low spatial frequencies, sensitivity falls as well, although (as we will see) the steepness of this low-frequency rolloff is quite dependent upon the conditions under which the measurements are made. It is interesting that most manipulations affect the two sides of the function differently, a point we will expand later.

The CSF is not an invariant function; rather the position and shape of the function are subject to significant variations depending on the optical quality of the viewing situation, the average luminance of the stimulus, the time-varying and chromatic qualities of the stimulus, and the part of the retina on which the stimulus falls. We will examine all of these effects because they reveal important aspects of the visual processes that affect contrast detection.

Before turning to the visual processes that affect contrast detection, we briefly describe how CSF measurements and linear systems analysis can be used to make general statements about spatial contrast sensitivity. Figure 10 plots contrast sensitivity to targets composed of extended sinusoids and to targets composed of one half-cycle of a cosinusoid. Notice that sensitivity to high spatial frequencies is much higher with the half-cycle cosinusoids than it is with extended sinusoids.

The half-cycle wave forms can be described as cosinusoids multiplied by rectangular functions of widths equal to one-half period of the cosinusoid. The truncating rectangular function causes the frequencies in the pattern to “splatter” to higher and lower values than the nominal target frequency. Multiplication of the Fourier transforms of the half-cycle targets by the CSF obtained with extended sinusoids yields an estimate of the visual system’s output response to the half-cycle targets. From there, use of a simple decision rule allows one to derive predicted half-cycle contrast sensitivities for a linear system. One expects to find greater visibility to high spatial frequencies with half-cycle gratings. The quantitative predictions are represented by the solid line in Fig. 10 and they match the observations rather well. There are numerous cases in the literature in which the CSF yields excellent predictions of the visibility of other sorts of patterns, but there are equally many cases in which it does not. We will examine some of the differences between those two situations later.

All of the information-processing stages depicted in Fig. 7 affect the CSF, and most of those effects have now been quantified. We begin with eye movements.

Eye Movements Even during steady fixation, the eyes are in constant motion, causing the retina to move with respect to the retinal images.¹⁰⁴ Such motion reduces retinal image contrast by smearing the spatial distribution of the target, but it also introduces temporal variation in the image. Thus, it is not obvious whether eye position jitter should degrade or improve contrast sensitivity. It turns out that the effect depends on the spatial frequency of the target. At low frequencies, eye movements are beneficial. When the image moves with respect to the retina, contrast sensitivity improves for spatial frequencies less than about 5 cpd.^{105,106} The effect of eye movements at higher spatial frequencies is less clear, but one can measure it by presenting sinusoidal interference fringes (which bypass optical degradations due to the eye’s aberrations) at durations too short to allow retinal smearing and at longer durations at which smearing might affect sensitivity.¹⁰⁷ There is surprisingly little attenuation due to eye position jitter; for 100-ms target presentations, sensitivity at 50 cpd decreases by only 0.2 to 0.3 log units relative to sensitivity at very brief durations. Interestingly, sensitivity improves at long durations of 500 to 2000 ms. This may be explained by noting that the eye is occasionally stationary enough to avoid smearing due to the retina’s motion.

Thus, the eye movements that occur during steady fixation apparently improve contrast sensitivity at low spatial frequencies and have little effect at high frequencies. This means that with steady fixation there is little need for experimenters to monitor or eliminate the effects of eye movements in measurements of spatial contrast sensitivity at high spatial frequencies.

Optics The optical transfer function (OTF) of an optical system describes the attenuation in contrast (and the phase shift) that sinusoidal gratings undergo when they pass through the optical system. As shown in Chap. 1 the human OTF (when the eye is well accommodated) is a low-pass function

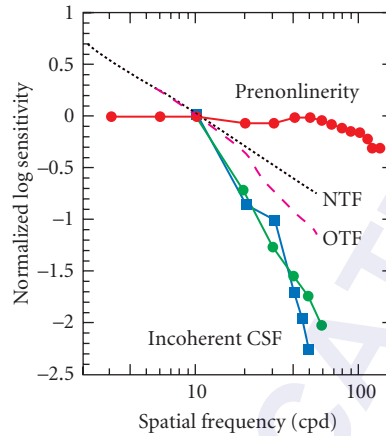


FIGURE 11 The contrast sensitivity function (CSF) and the transfer functions for various visual processing stages. All functions have been normalized to a log sensitivity of 0.0 at 10 cpd. The blue squares represent the CSF measured in the conventional way using incoherent white light.²² The red circles represent the transfer function of the prenonlinearity filter measured.²² The dotted line represents the neural transfer function that results from a bank of photon-noise-limited spatial filters of constant bandwidth.¹²² The dashed line represents the optical transfer function for a 3.8-mm pupil.⁴ The green circles are the product of the three transfer functions.

whose shape depends on pupil size and eccentricity. The dashed curve and the blue squares in Fig. 11 show the foveal OTF and the foveal CSF, respectively, for a pupil diameter of 3.8 mm. At high spatial frequencies, the foveal CSF declines at a faster rate than the foveal OTF. Thus, the OTF can only account for part of the high-frequency falloff of the CSF. Also, it is obvious that the OTF (which is a low-pass function) cannot account for the falloff in the CSF at low spatial frequencies.

Receptors Several receptor properties affect spatial contrast sensitivity and acuity. As noted earlier, the receptors in the fovea are all cones, and they are very tightly packed and arranged in a nearly uniform hexagonal lattice. With increasing retinal eccentricity, the ratio of cones to rods decreases steadily (see Fig. 2a) and the regularity of the lattice diminishes.

Several investigators have noted the close correspondence between the Nyquist frequency and the highest-detectable spatial frequency^{108–111} and have argued that the geometry of the foveal cone lattice sets a fundamental limit to grating acuity. However, more importance has been placed on this relationship than is warranted. This can be shown in two ways. First, grating acuity varies with several stimulus parameters including space-average luminance, contrast, and temporal duration. For example, grating acuity climbs from about 5 cpd at a luminance of 0.001 cd/m² to 50 to 60 cpd at luminances of 100 cd/m² or higher.¹¹² Obviously, the cone lattice does not change its properties in conjunction with changes in luminance, so this effect has nothing to do with the Nyquist limit. The only conditions under which the Nyquist limit might impose a fundamental limit to visibility would be at high luminances and contrasts, and relatively long target durations. Second, and more importantly, several investigators have shown that observers can actually detect targets whose frequency components are well above the Nyquist limit. For example, Williams¹¹³ presented high-frequency sinusoidal gratings using laser inter-ferometry to bypass the contrast reductions that normally occur in passage through the eye's optics. He found that observers could detect gratings at frequencies as high as 150 cpd, which

is 2.5 times the Nyquist limit. Observers did so by detecting the gratings' "aliases," which are Moiré-like distortion products created by undersampling the stimulus.¹¹⁴ The CSF under these conditions is smooth with no hint of a loss of visibility at the Nyquist limit. The observer simply needs to switch strategies above and below the Nyquist limit; when the targets are below the limit, the observer can detect the undistorted targets, and when the targets are above the limit, the observer must resort to detecting the targets' lower frequency aliases.⁷ On the other hand, when the task is to resolve the stripes of a grating (as opposed to detecting the grating), performance is determined by the Nyquist limit over a fairly wide range of space-average luminances.¹¹⁵

In addition to receptor spacing, one must consider receptor size in the transmission of signals through the early stages of vision. As noted earlier, the prevailing model of cones holds that the inner segment offers an aperture to incoming photons. In the fovea, the diameter of the inner segment is roughly $0.5 \text{ min}^{6,21}$ (see Fig. 2c). Modeling the cone aperture by a cylinder function of this diameter, one can estimate the low-pass filtering due to such an aperture from the cylinder's Fourier transform, which is a first-order Bessel function whose first zero occurs at 146 cpd. However, the entrance to a structure only a few wavelengths of light wide cannot be described accurately by geometric optics,¹¹⁶ so the receptor aperture is generally modeled by a gaussian function whose full width at half height is somewhat smaller than the anatomical estimates of the aperture diameter.²²

The modulation transfer of the neural stages prior to a compressive nonlinearity (thought to exist in the bipolar cells of the retina or earlier) has been measured in the following way.²² Laser interferometry was used to bypass the contrast reductions that normally occur as light passes through the eye's optics. The stimuli were two high-contrast sinusoids of the same spatial frequency but slightly different orientations; these created a spatial "beat" of a different orientation and frequency than the component gratings. Passing two sinusoids through a compressive nonlinearity, like the ones known to exist early in the visual system, creates a number of distortion products at various frequencies and orientations. MacLeod et al.²² varied the orientations and frequencies of the components so as to create a "beat" at 10 cpd, and then measured contrast sensitivity to the "beat." Because the "beat" was always at the same spatial frequency, the filtering properties of postnonlinearity stages could not affect the measurements and, therefore, any effect of component frequency on contrast sensitivity had to be due to the filtering properties of prenonlinearity stages. Observers were able to detect the "beat" produced by component gratings of frequencies as high as 140 cpd. This implies that the prenonlinearity filter has exceptionally large bandwidth. By inverse Fourier transformation of these data, one can estimate the spatial extent of the filter; it appears to have a full width at half height of about 16 arcsec, which is a bit smaller than current estimates of the cone aperture. Thus, the early neural stages are able to pass information up to extremely high spatial frequencies once the optics are bypassed.

The red circles in Fig. 11 display the estimate of the transfer function of the prenonlinearity neural filter for foveal viewing. Notice that it is virtually flat up to 140 cpd. Also shown (blue squares) is the CSF obtained in the conventional manner with incoherent white light from a CRT display. This function represents the product of the optics and the prenonlinearity and postnonlinearity spatial filters. As can be seen, the bandwidth of the conventional CSF is much smaller than the bandwidth of the prenonlinearity filter. Thus, optics and postnonlinearity filtering are the major constraints to contrast sensitivity at higher spatial frequencies.

As mentioned earlier, the transfer function of the optics is shown as the dashed line. The product of the optical transfer function and the prenonlinearity transfer function would be very similar to the dashed line because the prenonlinearity function is virtually flat. Thus, the product of optics and prenonlinearity filter does not match the observed CSF under conventional viewing conditions at all; the postnonlinearity filter must also contribute significantly to the shape of the CSF. The mismatch between the dashed line and the CSF is obvious at low and high spatial frequencies. We consider the high-frequency part first; the low-frequency part is considered later under "Adaptation and Inhibition."

Spatial Channels Spatial filters with narrow tuning to spatial frequency have been demonstrated in a variety of ways.^{117,118} For example, the visibility of a sinusoidal target is reduced by the presence of a narrowband noise masker whose center frequency corresponds to the spatial frequency of the sinusoid, but is virtually unaffected by the presence of a masker whose center frequency is more than

1.5 octaves from the frequency of the sinusoid (Ref. 119; see “Contrast Discrimination and Contrast Masking” later in chapter). These spatial mechanisms, which have been described mathematically by Gabor functions and other functions,¹²⁰ correspond to a first approximation with the properties of single neurons in the visual cortex.⁵⁸ They seem to have roughly constant bandwidths, expressed in octaves or log units, regardless of preferred frequency.^{58,117} As a consequence, mechanisms responding to higher spatial frequencies are smaller in vertical and horizontal spatial extent than mechanisms responding to lower frequencies.^{121,122} Thus, for a given stimulus luminance, a higher-frequency mechanism receives fewer photons than does a lower-frequency mechanism. The mean number of photons delivered to such a mechanism divided by the standard deviation of the number of photons—the signal-to-noise ratio—follows a square root relation,¹²³ so the signal-to-noise ratio is inversely proportional to the preferred frequency of the mechanism. The transfer function one expects for a set of such mechanisms should be proportional to $1/f$ where f is spatial frequency.¹²² This function is represented by the dotted line in Fig. 11. If we use that relation for describing the postnonlinearity filter, and incorporate the measures of the optical and prenonlinearity filters described above, the resultant is represented by the green circles. The fit between the observed CSF and the resultant is good, so we can conclude that the filters described here are sufficient to account for the shape of the high-frequency limb of the human CSF under foveal viewing conditions.

By use of the theory of ideal observers (see earlier), one can calculate the highest possible contrast sensitivity a system could have if limited by the optics, photoreceptor properties, and constant bandwidth spatial channels described above. As can be seen (green circles), the high-frequency limb of the CSF of such an ideal observer is similar in shape to that of human observers,^{122,124} although the rate of falloff at the highest spatial frequencies is not quite as steep as that of human observers.¹²⁴ More importantly, the absolute sensitivity of the ideal observer is significantly higher. Some of the low performance of the human compared to the best possible appears to be the consequence of processes that behave like noise internal to the visual system and some appears to be the consequence of employing detecting mechanisms that are not optimally suited for the target being presented.^{79,89} The high detection efficiencies observed for sinewave grating patches in the presence of static white noise (e.g., Ref. 87; Fig. 9) suggest that poor spatial pooling of information is not the primary factor responsible for the low human performance (at least when the number of cycles in the grating patches are kept low).

Optical/Retinal Inhomogeneity The fovea is undoubtedly the most critical part of the retina for numerous visual tasks such as object identification and manipulation, reading, and more. But, the fovea occupies less than 1 percent of the total retinal surface area, so it is not surprising to find that nonfoveal regions are important for many visual skills such as selecting relevant objects to fixate, maintaining posture, and determining one’s direction of self-motion with respect to environmental landmarks. Thus, it is important to characterize the visual capacity of the eye for nonfoveal loci as well.

The quality of the eye’s optics diminishes slightly from 0 to 20 deg of retinal eccentricity and substantially from 20 to 60 deg;¹ see Chap. 1. As described earlier and displayed in Fig. 2, the dimensions and constituents of the photoreceptor lattice and the other retinal neural elements vary significantly with retinal eccentricity. With increasing eccentricity, cone and retinal ganglion cell densities fall steadily and individual cones become broader and shorter. Given these striking variations in optical and receptor properties across the retina, it is not surprising that different aspects of spatial vision vary greatly with retinal eccentricity. Figure 12 shows CSFs at eccentricities of 0 to 40 deg. With increasing eccentricity, contrast sensitivity falls off at progressively lower spatial frequencies. The high-frequency cutoffs range from approximately 2 cpd at 40 deg to 50 cpd in the fovea. Low-frequency sensitivity is similar from one eccentricity to the next. There have been many attempts to relate the properties of the eye’s optics, the receptors, and postreceptor mechanisms to contrast sensitivity and acuity. For example, models that incorporate eccentricity-dependent variations in optics and cone lattice properties, plus the assumption of fixed bandwidth filters (as in Ref. 122), have been constructed and found inadequate for explaining the observed variations in contrast sensitivity; specifically, high-frequency sensitivity declines more with eccentricity than can be explained by information losses due to receptor lattice properties and fixed bandwidth filters alone.¹²⁵ Adding a low-pass filter to the model representing the convergence of cones onto retinal ganglion cells³² (i.e., by incorporating variation in receptive field center diameter) yields a reasonably accurate account of eccentricity-dependent variations in contrast sensitivity.¹²⁵

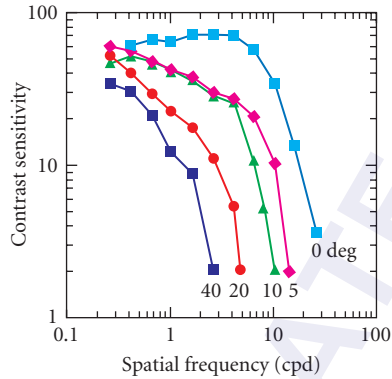


FIGURE 12 CSFs at different retinal eccentricities. (Adapted from Ref. 125.)

Adaptation and Inhibition Contrast sensitivity is also dependent upon the space-average intensity (i.e., the mean or background luminance) and the prior level of light adaptation. The blue symbols in Fig. 13 show how threshold amplitudes for spatial sinewaves vary as a function of background luminance for targets of 1 and 12 cpd.¹²⁶ [Recall that sinewave amplitude is contrast multiplied by background intensity; see discussion preceding Eq. (24).] The backgrounds were presented continuously and the observer was allowed to become fully light-adapted before thresholds were measured. For the low-spatial-frequency target (blue circles), threshold increases linearly with a slope of 1.0 above the low background intensities (1 log td). A slope of 1.0 in log-log coordinates implies that

$$A_T = k\bar{I} \quad (25)$$

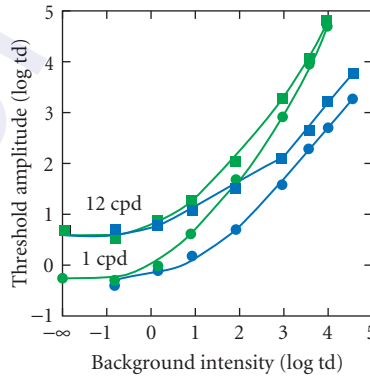


FIGURE 13 Amplitude threshold for gaussian-damped sinusoidal targets with 0.5-octave bandwidths as a function of background intensity. The targets were presented for 50 ms. The green symbols are thresholds measured at the onset of backgrounds flashed for 500 ms in the dark-adapted eye. The blue symbols are thresholds measured against steady backgrounds after complete adaptation to the background intensity. The difference between the green and blue symbols shows the effect of adaptation on contrast sensitivity. (Adapted from Ref. 126.)

where A_T is the amplitude threshold (or equivalently, it implies that contrast threshold, C_T , is constant). This proportional relationship between threshold and background intensity is known as *Weber's law*. For the higher spatial-frequency target (blue squares), threshold increases linearly with a slope of 0.5 at intermediate background intensities (0 to 3 log td) and converges toward a slope of 1.0 at the highest background intensities. A slope of 0.5 implies that

$$A_T = k\sqrt{I} \quad (26)$$

(or equivalently, that contrast threshold, C_T , decreases in inverse proportion to the square root of background intensity). This relationship between amplitude threshold and the square root of background intensity is known as the *square root law* or the *DeVries-Rose law*. As suggested by the data in Fig. 13, the size of the square root region gradually grows as spatial frequency is increased. The effect on the CSF is an elevation of the peak and of the high-frequency limb relative to the low-frequency limb as mean luminance is increased, thereby producing a shift of the peak toward higher spatial frequencies.^{127,128} This result explains why humans prefer bright environments when performing tasks that require high spatial precision.

The visual system is impressive in its ability to maintain high sensitivity to small contrasts (as produced, for example, by small changes in surface reflectance) over the enormous range of ambient light levels that occur in the normal environment. In order to have high sensitivity to small contrasts, the response gain of visual neurons must be kept high, but high gain leaves the visual system vulnerable to the effects of response saturation (because neurons are noisy and have a limited response range). The visual system's solution to this "dynamic-range problem" is threefold: (1) use separate populations of receptors (rods and cones) to detect contrasts in the lower and upper ranges of ambient intensity, (2) adjust the retinal illumination via the pupil reflex, and (3) directly adjust the sensitivity of individual receptors and neurons (e.g., Refs. 129, 130). All three components are important, but the third is the most significant and is accomplished within the receptors and other neurons through photochemical and neural adaptation. The combined effect of the photochemical and neural adaptation mechanisms on foveal contrast detection is illustrated by the difference between the open and solid data points in Fig. 13 (all the data were obtained with a fixed pupil size). The green symbols show cone detection threshold for the 1- and 12-cpd targets on a background flashed in the dark-adapted eye (i.e., there was no chance for light adaptation) as a function of background intensity. The difference between the green and blue symbols shows the improvement in detection sensitivity due to photochemical and neural adaptation within the cone system.

The pupil reflex and photopigment depletion (see "Image Sampling by the Photoreceptors") are *multiplicative adaptation* mechanisms; they adjust the sensitivity of the eye by effectively scaling the input intensity by a multiplicative gain factor that decreases with increasing average intensity. The pupil reflex operates over most of the intensity range (see Chap. 1); photopigment depletion is only effective in the cone system above about 4 log td (see Fig. 13) and is ineffective in the rod system, at least over the range of relevant intensities (see earlier discussion). There is considerable psychophysical evidence that a substantial component of the remaining improvement in sensitivity illustrated in Fig. 12 is due to multiplicative neural adaptation mechanisms.^{129,131-134} However, multiplicative adaptation alone cannot reduce threshold below a line of slope 1.0, tangent to the threshold curve in dark-adapted eye.^{129,134} The remaining improvements in sensitivity (the difference between the tangent lines and the blue symbols) can be explained by neural *subtractive adaptation* mechanisms.^{134,135} There is evidence that the multiplicative and subtractive adaptation mechanisms each have fast and slow components.¹³⁴⁻¹³⁷ Threshold functions, such as those in Fig. 13, can be predicted by simple models consisting of a compressive (saturating) nonlinearity [Eq. (13)], a multiplicative adaptation mechanism, a subtractive adaptation mechanism, and either constant additive noise or multiplicative (Poisson) neural noise.

A common explanation of the low-frequency falloff of the CSF (e.g., Fig. 12) is based on center/surround mechanisms evident in the responses of retinal and LGN neurons. For example, compare the CSFs measured in individual LGN neurons in Fig. 5 with the human CSF in Fig. 12. The center and surround mechanisms are both low-pass spatial filters, but the cutoff frequency of the surround is much lower than that of the center. The center and surround responses are subtracted neurally, so

at low spatial frequencies, to which center and surround are equally responsive, the neuron's response is small. With increasing spatial frequency, the higher resolution of the center mechanism yields a greater response relative to the surround, so the neuron's response increases.

While the center/surround mechanisms are undoubtedly part of the explanation, they are unlikely to be the whole story. For one thing, the time course of the development of the low-frequency falloff is slower than one expects from physiological measurements.¹⁰⁵ Second, the surround strength in retinal and LGN neurons is on average not much more than half the center strength;³⁷ this level of surround strength is consistent with the modest low-frequency falloffs observed under transient presentations, but inconsistent with the steep falloff observed for long-duration, steady-fixation conditions. Third, there is considerable evidence for slow subtractive and multiplicative adaptation mechanisms, and these ought to contribute strongly to the shape of the low-frequency limb under steady fixation conditions.^{137–139} Thus, the strong low-frequency falloff of the CSF under steady viewing conditions may reflect, to some degree, slow subtractive and multiplicative adaptation mechanisms whose purpose is to prevent response saturation while maintaining high contrast sensitivity over a wide range of ambient light levels. The weaker low-frequency falloff seen for brief presentations may be the result of the fast-acting surround mechanisms typically measured in physiological studies; however, fast-acting local multiplicative and subtractive mechanisms could also play a role.

Temporal Contrast Detection Space-time plots of most contrast-detection stimuli (plots of intensity as a function of x, y, t) show intensity variations in both space and time; in other words, most contrast detection stimuli contain both spatial and temporal contrast.

The discussion of contrast detection has, so far, emphasized spatial dimensions (e.g., spatial frequency), but there is, not surprisingly, a parallel literature emphasizing temporal dimensions (for a recent review see Ref. 140). Figure 14 shows de Lange's¹⁴¹ measurements of temporal contrast sensitivity

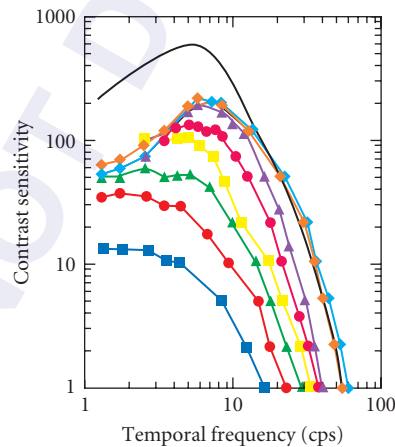


FIGURE 14 Temporal contrast sensitivity functions measured in the fovea at several mean luminances for a 2° circular test field. The test field was modulated sinusoidally and was imbedded in a large, uniform background field of the same mean luminance. Mean luminance: (■) 0.375 td, (●) 1 td, (▲) 3.75 td, (■) 10 td, (●) 37.5 td, (▲) 100 td, (◆) 1000 td, (◆) 10,000 td. (Data from Ref. 141.) The black curve is the MTF derived from the impulse response function of a single macaque cone reported by Schnapf et al.¹⁷ The cone MTF was shifted vertically (in log-log coordinates) for comparison with the shape of the high-frequency falloff measured psychophysically.

functions at several mean luminances for uniform discs of light modulated in intensity sinusoidally over time. The major features of these results have been confirmed by other researchers.^{142,143} Like spatial CSFs, temporal CSFs typically have a bandpass shape. At high mean luminances, the peak contrast sensitivity occurs at approximately 8 cycles per second (cps) and the cutoff frequency (also known as the critical flicker frequency) at approximately 60 cps. Contrast sensitivity increases with increasing mean luminance, larger increases (in log-log coordinates) being observed at middle and high temporal frequencies. For low temporal frequencies (below 6 cps), contrast sensitivity is nearly constant (i.e., Weber's law is observed) for mean intensities greater than 10 td or so.

The entire temporal CSF can be fit by a weighted difference of linear low-pass filters, where one of the filters has a relatively lower cutoff frequency, and/or where one of the filters introduces a relatively greater time delay (see Ref. 140). However, to fit the temporal CSFs obtained under different stimulus conditions, the time constants and relative weights of the filters must be allowed to vary. The low-frequency falloff may be partly due to the biphasic response of the photoreceptors (see Fig. 3), which appears to involve a fast subtractive process.¹⁷ Surround mechanisms, local subtractive mechanisms, and local multiplicative mechanisms also probably contribute to the low-frequency falloff.

The factors responsible for the high-frequency falloff are not fully understood, but there is some evidence that much of the high-frequency falloff is the result of temporal integration within the photoreceptors.^{144,145} The black curve in Fig. 14 is the MTF of the macaque cone photoreceptor response derived from measured impulse-response functions¹⁷ (see Fig. 3). The impulse responses were obtained in the linear response range by presenting brief, dim flashes on a dark background. Although the background was dark, there appears to be little effect of light adaptation on cone responses until background intensity exceeds at least a few hundred trolands;^{17,24} thus, the cone impulse-response functions measured in the dark ought to be valid up to moderate background intensities. Following Boynton and Baron,¹⁴⁴ we have shifted the MTF vertically in order to compare the shape of the cone high-frequency falloff to that of the temporal CSF. The shapes are similar, lending support to the hypothesis that cone temporal integration is responsible for the shape of the high-frequency limb of the temporal CSF.

It is well established that a pulse of light briefer than some critical duration will be just visible when the product of its duration and intensity is constant. This is *Bloch's law*, which can be written:

$$I_T T = k \quad \text{for} \quad T < T_c \quad (27)$$

where k is constant, T_c is the critical duration and I_T is the threshold intensity at duration T . Naturally, Bloch's law can also be stated in terms of contrasts by dividing the intensities in the above equation by the background intensity. Below the critical duration, log-log plots of threshold intensity versus duration have a slope of -1 . Above the critical duration, threshold intensity falls more slowly with increasing duration than predicted by Bloch's law (slope between -1 and 0) and, in some cases, it even becomes independent of duration (slope of 0).^{76,146} It is sometimes difficult to determine the critical duration because the slopes of threshold-versus-duration plots can change gradually, but under most conditions Bloch's critical duration declines monotonically with increasing background intensity from about 100 ms at 0 log td to about 25 ms at 4 log td.

Bloch's law is an inevitable consequence of a linear filter that only passes temporal frequencies below some cutoff value. This can be seen by computing the Fourier transform of pulses for which the product $I_T T$ is constant. The amplitude spectrum is given by $I_T \sin(T\pi f)/\pi f$ where f is temporal frequency; this function has a value $I_T T$ at $f=0$ and is quite similar over a range of temporal frequencies for small values of T . Indeed, to a first approximation, one can predict the critical duration T_c from the temporal CSFs shown in Fig. 14, as well as the increase in T_c with decreasing background intensity.¹⁴⁰

The shape of the temporal sensitivity function depends upon the spatial frequency of the stimulus. Figure 15 shows temporal CSFs measured for spatial sinewave gratings whose contrast is modulated sinusoidally in time.¹⁴⁷ With increasing spatial frequency, the temporal CSF evolves from a bandpass function with high sensitivity (as in Fig. 14) to a low-pass function of lower sensitivity.^{128,147} The solid curves passing through the data for the high spatial-frequency targets and the dashed curves passing through the data for the low spatial-frequency targets are identical except for vertical shifting. Similar behavior occurs for spatial CSFs at high spatial frequencies: changes in temporal frequency cause a vertical sensitivity shift (in log coordinates). Such behavior shows that the spatiotemporal CSF is separable at high temporal and spatial frequencies; that is, sensitivity can be predicted from the product

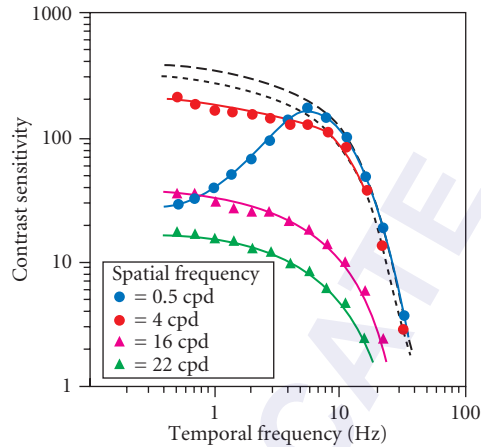


FIGURE 15 Temporal contrast sensitivity functions measured for several different spatial frequencies. The contrast of spatial sinewave gratings were modulated sinusoidally in time. The solid curves passing through the data for the high spatial-frequency targets and the dashed curves passing through the data for the low-frequency targets are identical except for vertical shifting. (Reproduced from Ref. 147.)

of the spatial and temporal CSFs at the appropriate frequencies. This finding suggests that the same anatomical and physiological mechanisms that determine the high-frequency slope of spatial and temporal CSFs determine sensitivity for the high-frequency regions of the spatiotemporal CSF.

The spatiotemporal CSF is, however, clearly not separable at low spatial and temporal frequencies, so an explanation of the underlying anatomical and physiological constraints is more complicated. The lack of spatiotemporal separability at low spatiotemporal frequencies has been interpreted by some as evidence for separate populations of visual neurons tuned to different temporal frequencies^{143,148} and by others as evidence for different spatiotemporal properties of center and surround mechanisms within the same populations of neurons.^{147,149} It is likely that a combination of the two explanations is the correct one. For a review of the evidence concerning this issue prior to 1984, see Ref. 140; for more recent evidence, see Refs. 150, 151, and 152.

The available evidence indicates that the temporal CSF varies little with retinal eccentricity once differences in spatial resolution are taken into account.^{153,154}

Chromatic Contrast Detection The discussion has, so far, concerned luminance contrast detection (contrast detection based upon changes in luminance over space or time); however, there is also a large literature concerned with chromatic contrast detection (contrast detection based upon changes in wavelength distribution over space or time). Although the focus of this chapter is on achromatic vision, it is appropriate to say a few words about chromatic contrast sensitivity functions. A chromatic spatial CSF is obtained by measuring thresholds for the detection of sinewave gratings that modulate spatially in wavelength composition, without modulating in luminance.* This is typically accomplished by adding two sinewave gratings of different wavelength distributions but the

*The standard definition of luminance is based upon “the $V(\lambda)$ function” [see Eq. (5)] which is an average of spectral sensitivity measurements on young human observers. Precision work in color vision often requires taking individual differences into account. To do this a psychophysical procedure, such as flicker photometry, is used to define luminance separately for each observer in the study. For more details, see Ref. 5 and Chap. 11.

same luminance amplitude. The two gratings are added in opposite spatial phase so that the mean luminance is constant across the pattern; such a pattern is said to be *isoluminant*. Formally, an isoluminant sinewave grating can be defined as follows:

$$I(x, y) = (A \sin(2\pi f x) + \bar{I}_1) + (-A \sin(2\pi f x) + \bar{I}_2) \quad (28)$$

where the first term on the right represents the grating for one of the wavelength distributions, and the second term on the right for the other distribution. The terms \bar{I}_1 and \bar{I}_2 are equal in units of luminance. Contrast sensitivity at spatial frequency f is measured by varying the common amplitude, A . (Note that $I(x, y)$ is constant independent of A , and that when $A = 0.0$ the image is a uniform field of constant wavelength composition.) Chromatic contrast is typically defined as $C = 2A / (\bar{I}_1 + \bar{I}_2)$.

The green circles in Fig. 16 are an isoluminant CSF for a combination of red and green gratings.¹⁵⁵ The squares are a luminance (isochromatic) CSF measured with the same procedure, on the same subjects. As can be seen, there are three major differences between luminance and chromatic CSFs:¹⁵⁶ (1) the high-frequency falloff of the chromatic CSFs occurs at lower spatial frequencies (acuity is worse); (2) the chromatic CSF has a much weaker (or absent) low-frequency falloff, even under steady viewing conditions; and (3) chromatic contrast sensitivity is greater than luminance contrast sensitivity at low spatial frequencies.

Two of the factors contributing to the lower high-frequency limb of the chromatic (red/green) CSF are the large overlap in the absorption spectra of the L (red) and M (green) cones, and the overlap (when present) of the wavelength distributions of the component gratings. Both of these factors reduce the effective contrast of the grating at the photoreceptors. A precise way to quantify these effects (as well as those of the other preneural factors) is with an ideal-observer analysis (e.g., Refs. 78, 124, 157). For example, one can physically quantify the equivalent luminance contrast of a chromatic grating by (1) computing the detection accuracy of an ideal observer, operating at the level of the photoreceptor photopigments, for the chromatic grating at the given chromatic contrast, and then (2) computing the equivalent luminance contrast required to bring the ideal observer to the same accuracy level.^{124, 158} The magenta circles in Fig. 16 are the chromatic CSF data replotted in terms

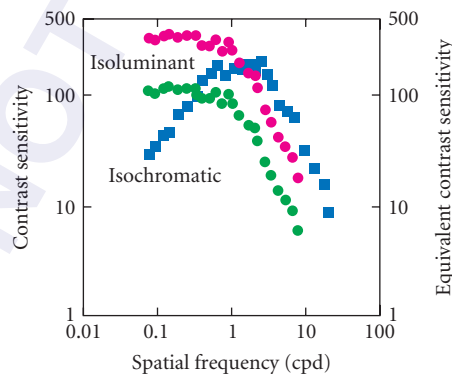


FIGURE 16 Isoluminant and isochromatic contrast sensitivity functions. The blue squares are isochromatic (luminance) CSFs for a green grating (526 nm). The green circles are isoluminant (chromatic) CSFs for gratings modulating between red (602) and green (526 nm). The magenta circles are the isoluminant data replotted as equivalent luminance contrast sensitivity (right axis). The isochromatic CSF plots the same on the left and right axes. (Data from Ref. 155.)

of equivalent luminance contrast. (The luminance CSF data are, of course, unaffected and hence are the same on both scales.) If the differences between the high-frequency falloffs of the luminance and chromatic CSFs were due to preneural factors alone, then the high-frequency limbs (magenta circles and blue squares) would superimpose. The analysis suggests that there are also neural factors contributing to the difference between high-frequency limbs.

One difficulty in interpreting the high-frequency limb of chromatic CSFs is the potential for luminance artifacts due to chromatic aberration.¹⁵⁶ Recently Sekiguchi et al.¹²⁴ eliminated chromatic aberration artifacts by producing isoluminant gratings with a laser interferometer which effectively bypasses the optics of the eye. Their results are similar to those in Fig. 16.

Comparison of the magenta circles and blue squares in Fig. 16 shows that the neural mechanisms are considerably more efficient at detecting chromatic contrast than luminance contrast at spatial frequencies below 2 cpd. One possible explanation is based upon the receptive-field properties of neurons in the retina and LGN (e.g., see Refs. 48 and 159). Many retinal ganglion cells (and presumably bipolar cells) are chromatically opponent in the sense that different classes of photoreceptors dominate the center and surround responses. The predicted CSFs of linear receptive fields with chromatically opponent centers and surrounds are qualitatively similar to those in Fig. 16.

Chromatic temporal CSFs have been measured in much the same fashion as chromatic spatial CSFs except that the chromatic sinewaves were modulated temporally rather than spatially [i.e., x represents time rather than space in Eq. (28)]. The pattern of results is also quite similar:^{94,156,160} (1) the high-frequency falloff of chromatic CSFs occurs at lower temporal frequencies, (2) the chromatic CSF has a much weaker low-frequency falloff, and (3) chromatic contrast sensitivity is greater than luminance contrast sensitivity at low temporal frequencies. Again, the general pattern of results is qualitatively predicted by a combination of preneural factors and the opponent receptive field properties of retinal ganglion cells (e.g., see Ref. 48).

Contrast Discrimination and Contrast Masking

The ability to discriminate on the basis of contrast is typically measured by presenting two sinusoids of the same spatial frequency, orientation, and phase, but differing contrasts. The common contrast C is referred to as the *pedestal contrast* and the additional contrast ΔC in one of the sinusoids as the *increment contrast*. Figure 17 shows typical *contrast discrimination functions* obtained from such measurements. The increment contrast varies as a function of the pedestal contrast. At low

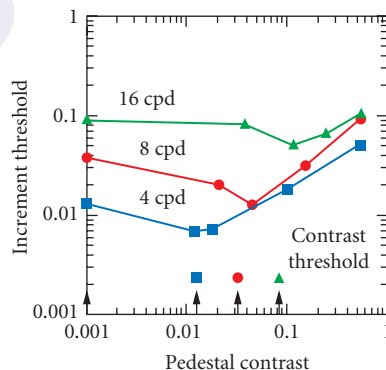


FIGURE 17 Contrast discrimination functions at different spatial frequencies. The just-discriminable contrast increment is plotted as a function of the contrast common to the two targets. The arrows indicate contrast thresholds at the indicated spatial frequencies. (Adapted from Ref. 78.)

pedestal contrasts, ΔC falls initially and then at higher contrasts becomes roughly proportional to $C^{0.6}$. Because the exponent is generally less than 1.0, Weber's law does not hold for contrast discrimination. Ideal observer calculations show that pedestal contrast has no effect on discrimination information; therefore, the variations in threshold with pedestal contrast must be due entirely to neural mechanisms. The dip at low pedestal contrasts has been modeled as a consequence of an accelerating nonlinearity in the visual pathways (e.g., Ref. 161) or the consequence of observer uncertainty (e.g., Ref. 162). The evidence currently favors the former.⁸⁹ Contrast discrimination functions at different spatial frequencies^{78,163} and at different eccentricities¹⁶⁴ are, to a first approximation, similar in shape; that is to say, the functions can be superimposed by shifting the data along the log pedestal and log increment axes by the ratio of detection thresholds.

A generalization of the contrast-discrimination experiment is the contrast-masking experiment in which the increment and pedestal differ from one another in spatial frequency, orientation and/or phase. In contrast-masking experiments, the increment and pedestal are usually referred to as the *target* and *masker*, respectively. Figure 18a shows typical contrast-masking functions obtained for sinewave maskers of different orientations. The orientation of the sinewave target is vertical (0 deg); thus, when the orientation of the masker is also vertical (squares), the result is just a contrast discrimination function, as in Fig. 17. For suprathreshold maskers (in this case, maskers above 2 percent contrast), the threshold elevation produced by the masker decreases as the difference between target and masker orientation increases; indeed, when the target and masker are perpendicular (diamonds), there is almost no threshold elevation even at high masker contrasts. Also, the slope of the contrast-masking function (in log-log coordinates) decreases as the difference in orientation increases. Finally, notice that threshold enhancement at low masking contrasts occurs only when the masker has approximately the same orientation as the target.

Figure 18b shows a similar set of masking functions obtained for sinewave maskers of different spatial frequencies. The effect of varying masker spatial frequency is similar to that of varying masker orientation: with greater differences in spatial frequency between masker and target, masking is reduced, the slope of the contrast-masking function is reduced, and the threshold enhancement at low masker contrasts becomes less evident.

The substantial variations in threshold seen as a function of masker spatial frequency and orientation provide one of the major lines of evidence for the existence of multiple orientation and spatial-frequency-tuned channels in the human visual system;^{119,165} for a review, see Ref. 166. Masking data (such as those in Fig. 18) have been analyzed within the framework of specific multiple-channel models in order to estimate the tuning functions of the spatial channels (e.g., Ref. 167). In one popular model, each channel consists of an initial linear filter followed by a static (zero memory) nonlinearity

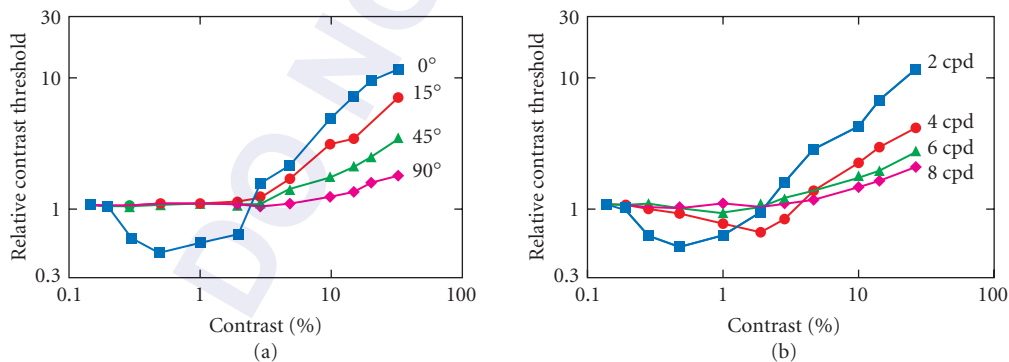


FIGURE 18 Contrast masking functions for sinewave-grating maskers of various orientations and spatial frequencies. The target was a 2-cpd sinewave grating. (a) Contrast-masking functions for 2-cpd maskers of various orientations. (b) Contrast-masking functions for 0° (vertical) maskers of various spatial frequencies. (Adapted from Ref. 265.)

that is accelerating at low contrasts and compressive at high contrasts (e.g., Refs. 161, 167, 168). Wilson and colleagues have shown that the spatial-channel tuning functions estimated within the framework of this model can predict discrimination performance in a number of different tasks (for a summary, see Ref. 169). However, recent physiological evidence^{70,71} and psychophysical evidence^{169,170} suggests that the compressive component of the nonlinearity is a broadly tuned, multiplicative gain mechanism, which is fundamentally different from the accelerating component. The full implications of these findings for the psychophysical estimation of channel-tuning functions and for the prediction of discrimination performance are not yet known.

Masking paradigms have also been used to estimate the tuning of temporal channels.^{151,171}

Contrast Estimation

We have argued that the differential visibility of coarse and fine patterns can be understood from an analysis of information losses in the early stages of vision: optics, receptor sampling, the inverse relation between the preferred spatial frequency and size of tuned spatial mechanisms, and retinal adaptation mechanisms. It is important to note, however, that the differential visibility described by the CSF in Fig. 11 does not relate directly to perceptual experience. For example, if you hold this page up close and then at arm's length, the apparent contrast of the text does not change appreciably even though the spatial frequency content (in cycles per degree) does. This observation has been demonstrated experimentally by asking observers to adjust the contrast of a sinusoidal grating of one spatial frequency until it appeared to match the contrast of a sinusoid of a different frequency.¹⁷² For example, consider two gratings, one at 5 cpd and another at 20 cpd; the former is near the peak of the CSF and the latter is well above the peak so it requires nearly a log unit more contrast to reach threshold. The 5-cpd target was set to a fixed contrast and the observer was asked to adjust the contrast of the 20-cpd target to achieve an apparent contrast match. When the contrast of the 5-cpd target was set near threshold, observers required about a log unit more contrast in the higher-frequency target before they had the same apparent contrast. The most interesting result occurred when the contrast of the 5-cpd grating was set to a value well above threshold. Observers then adjusted the contrast of the 20-cpd grating to the same physical value as the contrast of the lower-frequency target. This is surprising because, as described above, two gratings of equal contrast but different spatial frequencies produce different retinal image contrasts (see Chap. 1). In other words, when observers set 5- and 20-cpd gratings to equal physical contrasts, they were accepting as equal in apparent contrast two gratings whose retinal image contrast differed substantially. This implies that the visual system compensates at suprathreshold contrasts for the defocusing effects of the eye's optics and perhaps for the low-pass filtering effects of early stages of processing. This phenomenon has been called *contrast constancy*; the reader might recognize the similarity to deblurring techniques used in aerial and satellite photography (e.g., Ref. 173).

Visual Acuity

The ability to perceive high-contrast spatial detail is termed *visual acuity*. Measurements of visual acuity are far and away the most common means of assessing ocular health¹⁷⁴ and suitability for operating motor vehicles.¹⁷⁵ The universal use of visual acuity is well justified for clinical assessment,¹⁷⁶ but there is evidence that it is unjustified for automobile licensing (e.g., Ref. 177).

As implied by the discussion of contrast sensitivity, eye movements, optics, receptor properties, and postreceptor neural mechanisms all conspire to limit acuity; one factor may dominate in a given situation, but they all contribute. To assess acuity, high-contrast patterns of various sizes are presented at a fixed distance. The smallest pattern or smallest critical pattern element that can be reliably detected or identified is taken as the threshold value and is usually expressed in minutes of arc. Countless types of stimuli have been used to measure visual acuity, but the four illustrated in Fig. 19 are most common.

Gratings are used to measure minimum angle of resolution (MAR); the spatial frequency of a high-contrast sinusoidal grating is increased until its modulation can no longer be perceived. The

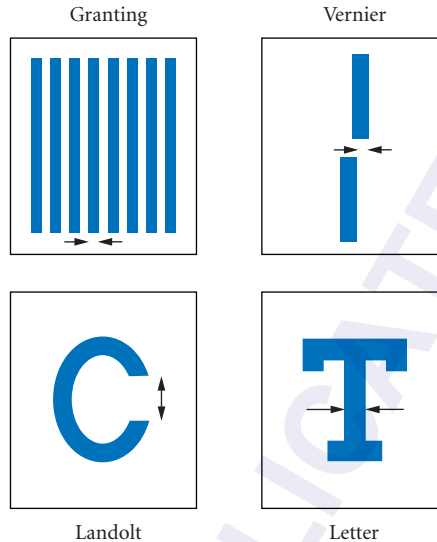


FIGURE 19 Targets commonly used to assess visual acuity. The dimension of the target that is taken as the acuity threshold is indicated by the arrows.

reader should note that this measurement is equivalent to locating the high-frequency cutoff of the spatial CSFs shown in previous sections. Under optimal conditions, the finest detectable grating bars are about 30 arcsec wide.¹¹²

The Landolt ring target is a high-contrast ring with a gap appearing in one of four positions. Threshold is defined as the smallest gap that can be correctly located. Under optimal conditions, threshold is again about 30 arcsec or slightly smaller.¹¹²

The most common test in clinical and licensing situations is the letter acuity task. The stimuli are a series of solid, high-contrast letters, varying in size. Threshold is defined by the smallest identifiable stroke width which, under optimal conditions, is nearly 30 arcsec.

Finally, the vernier acuity task involves the presentation of two nearly aligned line segments. Threshold is defined by the smallest visible offset from colinearity. Under optimal conditions, vernier acuity is about 5 arcsec.^{178,179} The special label, *hyperacuity*, has been provided for spatial thresholds, like vernier acuity, that are smaller than the intercone distance.¹⁸⁰

In practice, testing environments can vary widely along dimensions that affect visual acuity. These dimensions include illumination, the care chosen to correct refractive errors, and more. The National Academy of Sciences (1980) has provided a set of recommended standards for use in acuity measurements including use of the Landolt ring.

Although these various acuity tasks are clearly related, performance on the four tasks is affected differently by a variety of factors and, for this reason, one suspects that the underlying mechanisms are not the same.^{110,181} Figure 20 shows how grating and vernier acuities vary with retinal eccentricity. It is obvious from this figure that no single scaling factor can equate grating and vernier acuity across retinal eccentricities; a scaling factor that would equate grating acuity, for example, would not be large enough to offset the dependence of vernier acuity on retinal eccentricity.

Indeed, there have been several demonstrations that the two types of acuity depend differently on experimental conditions. For example, grating acuity appears more susceptible to optical defects and contrast reductions, whereas vernier, Landolt, and letter acuities are more susceptible to a visual condition called *amblyopia*.¹⁸² There are numerous models attempting to explain why different types

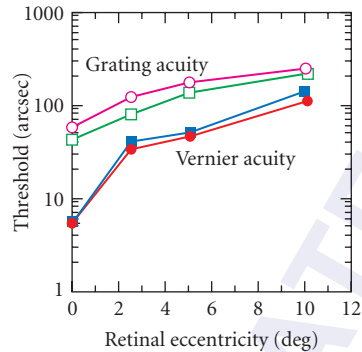


FIGURE 20 Grating and vernier acuity as a function of retinal eccentricity. The open symbols represent the smallest detectable bar widths for a high-contrast grating. The filled symbols represent the smallest detectable offset of one line segment with respect to another. Circles and squares are the data from two different observers. (Adapted from Ref. 110.)

of acuity are affected dissimilarly by experimental conditions, particularly eccentric viewing (e.g., Refs. 183 and 184).

The explanation of spatial discrimination thresholds smaller than the distance between photoreceptors, such as vernier acuity, has become a central issue in contemporary models of spatial vision. One view is that the fine sensitivity exhibited by the hyperacuities¹⁷⁹ reveals a set of local spatial primitives or features that are encoded particularly efficiently in the early stages of vision.^{185–187} The other view is that hyperacuity can be understood from an analysis of information processing among the spatially tuned mechanisms described in the earlier sections.^{183,188,189} Evidence favoring the second view is manifold. It has been shown, for example, that offset thresholds smaller than the grain of the photoreceptor lattice are predicted from measurements of the linear filtering and signal-to-noise properties of retinal ganglion cells.¹⁸⁸ Likewise, human vernier thresholds can be predicted in many situations from measurements of contrast discrimination thresholds.¹⁹⁰ Also, ideal-observer analysis has demonstrated that the information available at the level of the photoreceptors predicts better performance in hyperacuity tasks than in acuity (resolution) tasks,¹⁹¹ suggesting that any relatively complete set of spatially tuned mechanisms should produce better performance in the hyperacuity tasks. The remaining puzzle for proponents of the second view, however, is why vernier acuity and the other hyperacuities diminish so strikingly with increasing retinal eccentricity.

Pattern Discrimination

An important aspect of spatial vision is the ability to distinguish one pattern from another. In this section, we consider the ability to distinguish simple, suprathreshold patterns that vary along a single dimension. The dimensions considered are orientation, size or spatial frequency, and position or phase.

The detectability of a pattern varies to some degree with orientation. For most observers, for example, targets are most visible when oriented vertically or horizontally and least when oriented obliquely.¹⁹² This *oblique effect* is demonstrable for periodic and aperiodic targets and is largest for fine-detail targets. The cause of the effect appears to be differential sensitivity among orientation-tuned mechanisms, presumably in the visual cortex.¹⁹³ The ability to discriminate gratings differing in orientation depends on several stimulus parameters including target length, contrast, and

spatial frequency, but observers can generally discriminate extended, high-contrast targets that differ by 1 deg or less.

The ability to discriminate gratings differing in spatial frequency or aperiodic stimuli differing in size depends critically on the reference frequency or size. That is to say, the ratio of the frequency discrimination threshold, Δf , to the reference frequency, f , is roughly constant at about 0.03 for a wide range of reference frequencies;¹⁹⁴ similarly, the ratio of size discrimination threshold to reference size is roughly constant for a variety of reference sizes, except when the reference is quite small.¹⁷⁹

The encoding of spatial phase is crucial to the identification of spatial patterns. Its importance is demonstrated quite compellingly by swapping the amplitude and phase spectra of two images; the appearance of such hybrid images corresponds to a much greater degree with their phase spectra than with their amplitude spectra.¹⁹⁵ The conventional explanation is that the phase spectrum determines the spatial structure of an image, but this is perhaps too simplistic because the amplitude spectra of most natural images are much more similar to one another than are the phase spectra (e.g., Refs. 196 and 197). Also, given that the result was obtained using global Fourier transforms, it does not have direct implications about the relative importance of phase coding within spatially localized channels.¹⁹⁸ Nonetheless, such demonstrations illustrate the important relationship between the encoding of spatial phase and pattern identification.

In phase discrimination tasks, the observer distinguishes between patterns—usually periodic patterns—that differ only in the phase relationships of their spatial frequency components. The representation of spatial phase in the fovea is very precise: observers can, for instance, discriminate relative phase shifts as small as 2 to 3 deg (of phase angle) in compound gratings composed of a fundamental and third harmonic of 5 cpd;¹⁹⁹ this is equivalent to distinguishing a positional shift of about 5 arcsec. Phase discrimination is much less precise in extrafoveal vision.^{200–202} There has been speculation that the discrimination anomalies observed in the periphery underlie the diminished ability to distinguish spatial displacements,¹¹⁰ to segment textures on the basis of higher-order statistics,^{203,204} and to identify complex patterns such as letters,²⁰⁵ but detailed hypotheses linking phase discrimination to these abilities have not been developed.

One model of phase discrimination²⁰⁶ holds that local energy computations are performed on visual inputs and that local energy peaks are singled out for further analysis. The energy computation is performed by cross-correlating the waveform with even- and odd-symmetric spatial mechanisms of various preferred spatial frequencies. The relative activities of the even- and odd-symmetric mechanisms are used to represent the sort of image feature producing the local energy peak. This account is supported by the observation that a two-channel model, composed of even- and odd-symmetric mechanisms, predicts many phase discrimination capabilities well (e.g., Refs. 207–209). The model also offers an explanation for the appearance of illusory Mach bands at some types of edges and not others.²⁰⁶ For this account to work, however, one must assume that odd-symmetric, and not even-symmetric, mechanisms are much less sensitive in the periphery than in the fovea because some relative phase discriminations are extremely difficult in extrafoveal vision and others are quite simple^{202,209,210} (but see Ref. 211).

Motion Detection and Discrimination

The ability to discriminate the form and magnitude of retinal image motion is critical for many fundamental visual tasks, including navigation through the environment, shape estimation, and distance estimation. There is a vast literature on motion perception (some representative reviews from somewhat different perspectives can be found in Refs. 45, 212, 213). Here, we briefly discuss simple motion discrimination in the frontoparallel plane, and (as an example of more complicated motion processing) observer heading discrimination.

Detection and/or discrimination of motion in the frontoparallel plane is influenced by a number of factors, including contrast, wavelength composition, spatial-frequency content, initial speed, and eccentricity. For simplicity, and for consistency with the other sections in this chapter, we consider here only experiments employing sinewave-grating or random-dot stimuli.

Contrast sensitivity functions for detection of moving (drifting) sinewave-grating targets have the interesting property that they are nearly shape invariant on a log spatial-frequency axis; furthermore,

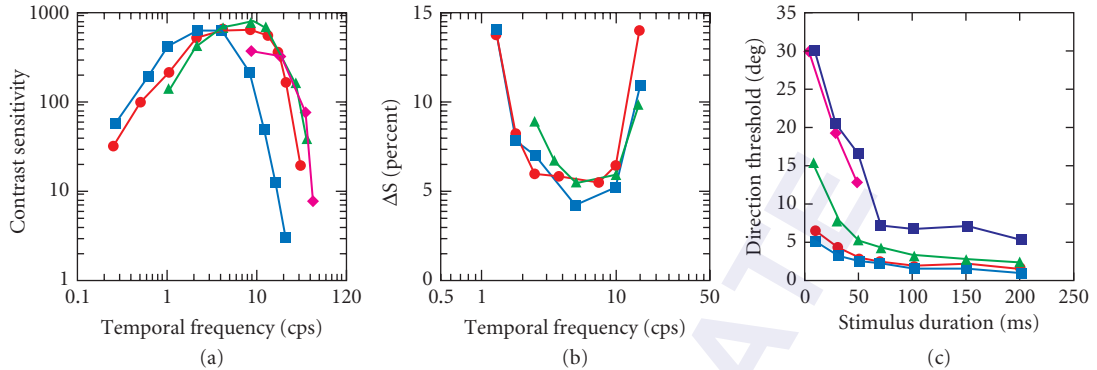


FIGURE 21 Motion detection and discrimination thresholds. (a) Contrast sensitivity for drifting sinewave gratings plotted as a function of temporal frequency. Each curve is for a different drift speed; (■) 1 deg/s, (●) 10 deg/s, (▲) 100 deg/s; (◆) 800 deg/s. (Adapted from Ref. 214.) (b) Speed discrimination thresholds plotted as a function of temporal frequency. Each curve is for a different drift speed; (●) 5 deg/s, (■) 10 deg/s, (▲) 20 deg/s. (Adapted from Ref. 215.) (c) Direction threshold for random dot patterns as a function of stimulus duration. Each curve is for a different drift speed; (■) 1 deg/s, (▲) 4 deg/s, (●) 16 deg/s, (■) 64 deg/s, (◆) 256 deg/s. (Adapted from Ref. 216.)

as shown in Fig. 21a, they are nearly superimposed at high speeds and low spatial frequencies, when plotted as a function of temporal frequency.^{106,214*} This latter result corresponds to the fact (mentioned earlier) that spatial CSFs are relatively flat at low spatial frequencies and high temporal frequencies.^{128,147} In interpreting Fig. 21a, it is useful to note that the velocity (V) of a drifting sine-wave grating is equal to the temporal frequency (f_t) divided by the spatial frequency (f_s):

$$V = \frac{f_t}{f_s} \quad (29)$$

Measurements of CSFs for moving stimuli provide relatively little insight into how the visual system extracts or represents motion information. For example, in one common paradigm¹⁰⁶ it is possible for subjects to perform accurately without “seeing” any motion at all. Greater insight is obtained by requiring the observer to discriminate between different aspects of motion such as speed or direction.

Representative measurements of speed discrimination for drifting sinewave grating targets are shown in Fig. 21b. The figure shows the just-detectable change in speed as a function of temporal frequency; each curve is for a differential initial or base speed. Similar to the motion CSFs, speed discrimination is seen to be largely dependent upon the temporal frequency and relatively independent of the speed.²¹⁵ As can be seen, the smallest detectable changes in speed are approximately 5 percent, and they occur in the temporal frequency range of 5 to 10 cps, which is similar to the temporal frequency where contrast sensitivity for drifting gratings is greatest (cf., Fig. 21a). Interestingly, 5 to 10 cps is the temporal frequency range of strongest response for most neurons in the macaque’s primary visual cortex.⁶⁰ For random-dot patterns, the smallest detectable change in speed is also approximately 5 percent.²¹⁶

Representative measurements of direction discrimination for drifting random-dot patterns are shown in Fig. 21c, which plots direction threshold in degrees as a function of stimulus duration for a wide range of speeds. Direction discrimination improves with duration and is a U-shaped function of dot speed. Under optimal conditions, direction discrimination thresholds are in the neighborhood of 1 to 2°.^{216,217}

The available evidence suggests that direction and speed discrimination improve quickly with contrast at low contrasts but show little improvement when contrast exceeds a few percent.^{215,218,219} Motion discrimination is relatively poor at isoluminance.^{220,221} The variations in motion discrimination with eccentricity can be largely accounted for by changes in spatial resolution.²²²

*The temporal frequency of a drifting grating is the number of stripes passing a fixed reference point per unit time (seconds).

Psychophysical evidence suggests that there are at least two mechanisms that can mediate motion discrimination—a *short-range mechanism* which appears to be closely associated with motion-selective neurons in the early levels of cortical processing, and a *long-range mechanism* which appears to be associated with more “inferential” calculations occurring in later levels of cortical processing. Clear evidence for this view was first observed in an apparent motion paradigm in which observers were required to judge the shape of a region created by laterally displacing a subset of random dots in a larger random-dot pattern.²²³ (These patterns are examples of *random-dot kinematograms*, the motion analog to random-dot stereograms.) Observers could accurately judge the shape of the region only if the displacements (1) were less than approximately 15 minarc ($D_{\max} = 15 \text{ min}$); (2) occurred with a delay of less than approximately 100 ms; and (3) occurred within the same eye (i.e., dichoptic presentations failed to produce reliable discrimination). These data suggest the existence of motion-sensitive mechanisms that can only detect local, monocular correlations in space-time.

The fact that motion-selective neurons in the early levels of the cortex can also only detect local spatiotemporal correlations suggests that they may be the primary source of neural information used by the subsequent mechanisms that extract shape and distance information in random-dot kinematograms. The receptive fields of motion-selective neurons decrease in size as a function of optimal spatial frequency, and they increase in size as a function of retinal eccentricity. Thus, the hypothesis that motion-selective neurons are the substrate for shape extraction in random-dot kinematograms is supported by the findings that D_{\max} increases as spatial-frequency content is decreased via low-pass filtering,²²⁴ and that D_{\max} increases as total size of the kinematogram is increased.²²⁵

Accurate shape and motion judgments can sometimes be made in apparent-motion displays where the stimulation is dichoptic and/or the spatiotemporal displacements are large. However, in these cases the relevant shape must be clearly visible in the static frames of the kinematogram. The implication is that higher-level, inferential analyses are used to make judgments under these circumstances. A classic example is that we can infer the motion of the hour hand on a clock from its long-term changes in position even though the motion is much too slow to be directly encoded by motion-selective cortical neurons.

The evidence for short-term and long-term motion mechanisms raises the question of what mechanisms actually mediate performance in a given motion discrimination task. The studies described above (e.g., Fig. 21) either randomized stimulus parameters (such as duration and spatial frequency) and/or used random-dot patterns; hence, they most likely measured properties of the short-range mechanisms. However, even with random-dot patterns, great care must be exercised when interpreting discrimination experiments. For example, static pattern-processing mechanisms may sometimes contribute to motion discrimination performance even though performance is at chance when static views are presented at great separations of space or time. This might occur because persisting neural responses from multiple stimulus frames merge to produce “virtual” structured patterns (e.g., Glass patterns²²⁶). Similar caution should be applied when interpreting discrimination performance with random-dot stereograms (see below).

The discussion of motion perception to this point has focused on the estimation of velocities of various stimuli. Such estimation is importantly involved in guiding locomotion through and estimating the three-dimensional layout of cluttered environments. Motion of an observer through a rigid visual scene produces characteristic patterns of motion on the retina called the *optic flow field*.^{227,228} Observers are able to judge the direction of their own motion through the environment based upon this optical flow information alone.^{229,230} In fact, accurate motion perception is possible even for random-dot flow fields, suggesting that the motion-selective neurons in the early levels of the visual cortex are a main source of information in perceiving self-motion. But, how is the local motion information provided by the motion-selective neurons used by subsequent mechanisms to compute observer motion?

Gibson^{227,228} proposed that people identify their direction of self-motion with respect to obstacles by locating the source of flow: the focus of expansion. Figure 22a depicts the flow field on the retina as the observer translates while fixating ahead. Flow is directed away from the focus of expansion and this point corresponds to the direction of translation. Not surprisingly, observers can determine their heading to within ± 1 deg from such a field of motion.²³⁰

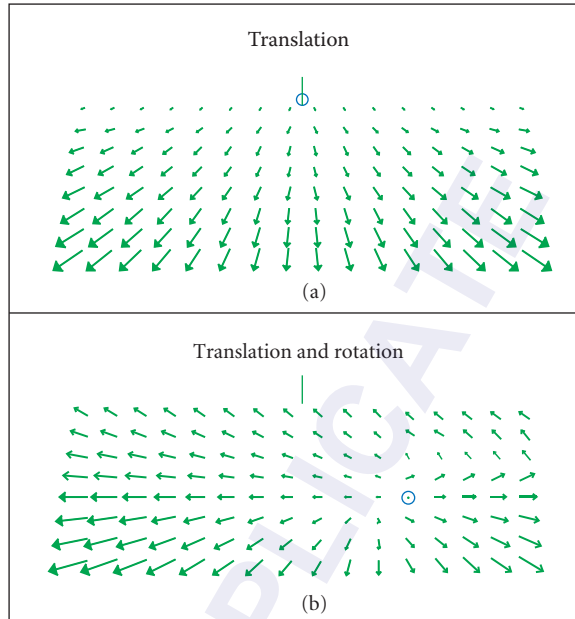


FIGURE 22 Optic flow fields resulting from forward translation across a rigid ground plane. (a) Flow field in the retinal image when the observer translates straight ahead while maintaining constant eye and head position; the heading is indicated by the small vertical line. (b) Retinal flow field when the observer translates straight ahead while making an eye movement to maintain fixation on the circle; again the heading is indicated by the small vertical line. (Adapted from Ref. 231.)

The situation becomes much more complicated once eye/head movements are considered. Figure 22*b* illustrates this by portraying the flow resulting from forward translation while the observer rotates the eyes to maintain fixation on a point off to the right. This motion does not produce a focus of expansion in the image corresponding to the heading; the only focus corresponds to the point of fixation. Consequently, heading cannot be determined by locating a focus (or singularity) in the retinal flow field. Nonetheless, human observers are able to locate their heading to within ± 1.5 deg.²³¹

Recent theoretical efforts have concentrated on this problem of computing heading in the presence of rotations due to eye/head movements. There are two types of models. One holds that observers measure the rotational flow components due to eye/head movements by means of an extraretinal signal; that is, the velocity of rotations is signaled by proprioceptive feedback from the extraocular and/or neck muscles or by efferent information.^{232–236} The rotational flow components are then subtracted from the flow field as a whole to estimate the heading. The other type of model holds that people determine heading in the presence of rotations from retinal image information alone (e.g., Refs. 233–236). These models hypothesize that the visual system decomposes the flow field into rotational and translational components by capitalizing on the fact that flows due to translation and rotation depend differently on the scene geometry. Once the decomposition is performed, the heading can be estimated from the translational components (see Fig. 22*a*).

Current experimental evidence favors the extraretinal signal model,^{237,238} but there is also evidence supporting the retinal-image model when the eye-head rotations are slow.²³¹

The optic flow field also contains information specifying the relative depths of objects in the visual scene: The velocities of retinal image motions are, to a first approximation, proportional to the inverse depths of the corresponding objects.²³⁹ Human observers are quite sensitive to this depth

cue as evidenced by the fact that they perceive depth variation for differential motions of as small as $\frac{1}{2} - \frac{1}{3}$ arcsec per s.²⁴⁰

Binocular Stereoscopic Discrimination

Light from the three-dimensional environment is imaged onto the two-dimensional array of photoreceptors, so the dimension associated with distance from the eye is collapsed in the two-dimensional retinal image. The images in the two eyes generally differ, however, because the eyes are separated by 6 to 7 cm. The differences between the images, which are called *binocular disparities*, are used by the visual system to recover information about the distance dimension. The depth percept that results from binocular disparity is called *stereopsis*. Binocular disparity (stereo) information is formally equivalent to the monocular information that results from a step translation of one eye by 6 to 7 cm. This observation implies a close theoretical connection between the computation of distance from motion and distance from binocular disparity. As indicated in an earlier section, the extraction of stereo and motion information presumably begins in the early stages of cortical processing (i.e., in V1 of the macaque), where many neurons are selective for disparity and direction of motion.

The geometry of stereopsis is described in Chap. 1, Fig. 25. Here, we briefly consider the nature of the information available for computing distance from disparity. The distance between any pair of points in the depth dimension (Δz) is related to the horizontal disparity (d) and the horizontal convergence angle of the eyes (θ) by the following formula:^{*}

$$\Delta z = \frac{ad}{(\theta + d_\theta + d/2)(\theta + d_\theta - d/2)} \quad (30)$$

where d_θ is the average disparity between the two points and the convergence point, and a is the interpupillary distance (the distance between the image nodal points of the eyes). If d_θ is zero, then the points are centered in depth about the convergence point. The locus of all points where the disparity d equals zero is known as the (Vieth-Müller) *horopter*. As Eq. (30) indicates, the computation of absolute distance between points from horizontal disparity requires knowledge of the “eye parameters” (the interpupillary distance and the angle of convergence). However, relative distance information is obtained even if the eye parameters are unknown (in other words, Δz is monotonic with d). Furthermore, absolute distance can be recovered in principle without knowing the angle of convergence from use of horizontal and vertical binocular disparities (e.g., Refs. 241 and 242).

Under optimal conditions, the human visual system is capable of detecting very small binocular disparities and hence rather small changes in distance. In the fovea, disparities of roughly 10 arcsec can be detected if the test objects contain sharp, high-contrast edges and if they are centered in depth about a point on or near the horopter ($d_\theta = 0.0$).^{178,243} At a viewing distance of 50 cm, 10 s corresponds to a distance change of about 0.02 cm, but at 50 m, it corresponds to a distance change of about 2 m.

In measuring stereoscopic discrimination performance, it is essential that the task cannot be performed on the basis of monocular information alone. This possibility can be assessed by measuring monocular and binocular performance with the same stimuli; if monocular performance is substantially worse than binocular performance, then binocular mechanisms per se are used in performing the stereoscopic task.^{243,244} Random-dot stereograms^{244,245} are particularly effective at isolating binocular mechanisms because discriminations cannot be performed reliably when such stereograms are viewed monocularly.^{244,246}

Stereopsis is affected by a number of factors. For example, the contrast required to produce a stereoscopic percept varies with the spatial frequency content of the stimulus. This result is quantified by the contrast sensitivity function for stereoscopic discrimination of spatially filtered random-dot stereograms; the function is similar in shape to CSFs measured for simple contrast detection,²⁴⁷ but

^{*}This formula is an approximation based on the assumption of relatively small angles and object location near the midsagittal plane; it is most accurate for distances under a few meters.

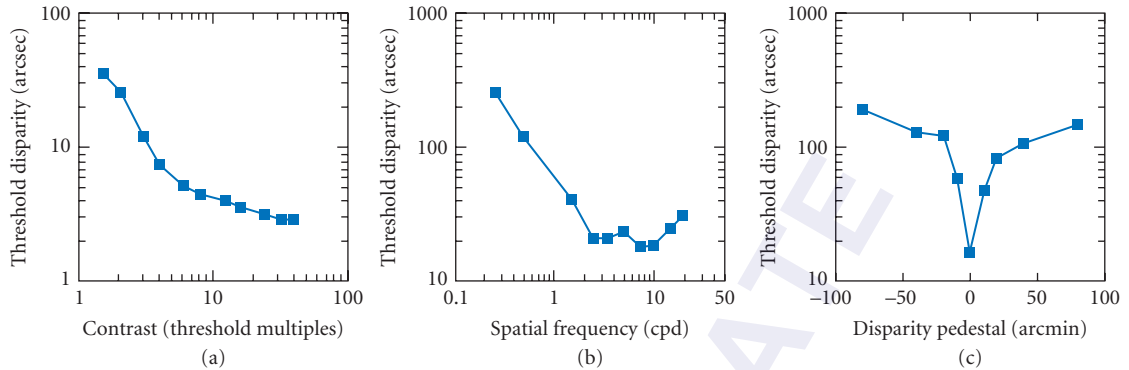


FIGURE 23 Disparity discrimination thresholds as a function of contrast, spatial frequency, and pedestal disparity. (a) Disparity threshold as a function of luminance contrast (in contrast threshold units) for dynamic random-dot stereograms. (Adapted from Ref. 256.) (b) Disparity threshold as a function of spatial frequency for sinewave-grating stereograms. (Adapted from Ref. 255.) (c) Disparity threshold as a function of the disparity pedestal (distance from the convergence plane) for difference-of-gaussian (DOG) stereograms. (Adapted from Ref. 258.).

detection occurs at lower contrasts than stereopsis occurs. In other words, there is a range of contrasts that are clearly detectable, but insufficient to yield a stereoscopic percept.²⁴⁷

The fact that the CSFs for contrast detection and stereoscopic discrimination have similar shapes suggests that common spatial-frequency mechanisms are involved in the two tasks. This hypothesis receives some support from masking and adaptation experiments demonstrating spatial-frequency tuning^{248–251} and orientation tuning²⁵² for stereo discrimination, and from electrophysiological studies demonstrating that cortical cells selective for binocular disparity are also usually selective for orientation^{55,56,62} and spatial frequency.²⁵³

The smallest detectable disparity, which is called *stereoacuity*, improves as luminance contrast is increased.^{254–256} As shown in Fig. 23a, stereoacuity improves approximately in inverse proportion to the square of contrast at low contrasts, and in inverse proportion to the cube root of contrast at high contrasts. Stereoacuity is also dependent upon spatial frequency; it improves in inverse proportion to target spatial frequency over the low spatial-frequency range, reaching optimum near 3 cpd (Fig. 23b).²⁵⁵ Both the inverse-square law at low contrasts and the linear law at low spatial frequencies are predicted by signal-detection models that assume independent, additive noise and simple detection rules.^{255,256}

Stereoacuity declines precipitously as a function of the distance of the test objects from the convergence plane (Fig. 23c).^{257,258} For example, adding a disparity pedestal of 40 minarc to test objects reduces acuity by about 1 log unit. This loss of acuity is not the result of losses of information due to the geometry of stereopsis; it must reflect the properties of the underlying neural mechanisms. [Note in Eq. (30) that adding a disparity to both objects is equivalent to changing the value of d_θ .]

Much like the spatial channels described earlier under “Spatial Channels,” there are multiple channels tuned to different disparities, but it is unclear whether there are a small number of such channels—“near,” “far,” and “tuned”²⁵⁹—or a continuum of channels.^{56,63,260,261} Models that assume a continuum of channels with finer tuning near the horopter predict both the sharp decline in stereoacuity away from the horopter^{261,262} and the shapes of tuning functions that have been estimated from adaptation and probability summation experiments.^{260,261}

2.5 ACKNOWLEDGMENTS

This project was supported in part by NIH grant EY02688 and AFOSR grant F49620-93-1-0307 to WSG and by NIH grant HD 19927 to MSB.

2.6 REFERENCES

1. R. Navarro, P. Artal, and D. R. Williams, "Modulation Transfer Function of the Human Eye as a Function of Retinal Eccentricity," *Journal of the Optical Society of America A* **10**:201–212 (1993).
2. J. D. Gaskill, *Linear Systems, Fourier Transforms, and Optics*, Wiley, New York, 1978.
3. J. W. Goodman, *Introduction to Fourier Optics*, Physical and Quantum Electronics Series, H. Heffner and A. E. Siegman (eds.), McGraw-Hill, New York, 1968.
4. F. W. Campbell and R. W. Gubisch, "Optical Quality of the Human Eye," *Journal of Physiology* **186**:558–578 (London, 1966).
5. G. Wyszecki and W. S. Stiles, *Color Science: Concepts and Methods, Quantitative Data and Formulae*, 2d ed., Wiley, New York, 1982.
6. C. A. Curcio, et al., "Human Photoreceptor Topography," *Journal of Comparative Neurology* **292**:497–523 (1990).
7. D. R. Williams, "Aliasing in Human Foveal Vision," *Vision Research* **25**:195–205 (1985).
8. L. N. Thibos, F. E. Cheney, and D. J. Walsh, "Retinal Limits to the Detection and Recognition of Gratings," *Journal of the Optical Society of America* **4**:1524–1529 (1987).
9. D. R. Williams, "Topography of the Foveal Cone Mosaic in the Living Human Eye," *Vision Research* **28**:433–454 (1988).
10. R. E. Marc and H. G. Sperling, "Chromatic Organization of Primate Cones," *Science* **196**:454–456 (1977).
11. F. M. de Monasterio, et al., "Density Profile of Blue-Sensitive Cones Along the Horizontal Meridian of Macaque Retina," *Investigative Ophthalmology and Visual Science* **26**:283–288 (1985).
12. P. Ahnelt, C. Keri, and H. Kolb, "Identification of Pedicles of Putative Blue-Sensitive Cones in the Human Retina," *Journal of Comparative Neurology* **293**:39–53 (1990).
13. D. R. Williams, D. I. A. MacLeod, and M. M. Hayhoe, "Foveal Tritanopia," *Vision Research* **21**:1341–1356 (1981).
14. J. L. Nerger and C. M. Cicerone, "The Ratio of L Cones to M Cones in the Human Parafoveal Retina," *Vision Research* **32**:879–888 (1992).
15. C. M. Cicerone, "Color Appearance and the Cone Mosaic in Trichromacy and Dichromacy," *Color Vision Deficiencies*, Y. Ohta (ed.) Kugler & Ghedini, Amsterdam, 1990, pp. 1–12.
16. J. K. Bowmaker, "Visual Pigments and Colour in Primates," *From Pigments to Perception: Advances in Understanding Visual Processes*, A. Valberg and B. L. Lee (eds.) Plenum Press, New York, 1991.
17. J. L. Schnapf et al., "Visual Transduction in Cones of the Monkey Macaca Fascicularis," *Journal of Physiology* **427**:681–713 (London, 1990).
18. W. A. H. Rushton, "The Difference Spectrum and the Photosensitivity of Rhodopsin in the Living Human Retina," *Journal of Physiology* **134**:11–297 (London, 1956).
19. W. A. H. Rushton and G. H. Henry, "Bleaching and Regeneration of Cone Pigments in Man," *Vision Research* **8**:617–631 (1968).
20. M. Alpern and E. N. Pugh, "The Density and Photosensitivity of Human Rhodopsin in the Living Retina," *Journal of Physiology* **237**:341–370 (London, 1974).
21. W. H. Miller and G. D. Bernard, "Averaging over the Foveal Receptor Aperture Curtails Aliasing," *Vision Research* **23**:1365–1369 (1983).
22. D. I. A. MacLeod, D. R. Williams, and W. Makous, "A Visual Nonlinearity Fed by Single Cones," *Vision Research* **32**:347–363 (1992).
23. J. L. Schnapf and D. A. Baylor, "How Photoreceptor Cells Respond to Light," *Scientific American* **256**(4):40–47 (1987).
24. J. M. Valetton and D. van Norren, "Light-Adaptation of Primate Cones: An Analysis Based on Extracellular Data," *Vision Research* **23**:1539–1547 (1982).
25. D. A. Baylor, B. J. Nunn, and J. L. Schnapf, "The Photocurrent, Noise and Spectral Sensitivity of Rods of the Monkey Macaca Fascicularis," *Journal of Physiology* **357**:576–607 (London, 1984).
26. D. C. Hood and D. G. Birch, "A Quantitative Measure of the Electrical Activity of Human Rod Photoreceptors Using Electroretinography," *Visual Neuroscience* **5**:379–387 (1990).

27. E. N. Pugh and T. D. Lamb, "Cyclic GMP and Calcium: The Internal Messengers of Excitation and Adaptation in Vertebrate Photoreceptors," *Vision Research* **30**(12):1923–1948 (1990).
28. D. C. Hood and D. G. Birch, "Human Cone Receptor Activity: The Leading Edge of the A-Wave and Models of Receptor Activity," *Visual Neuroscience* **10**:857–871 (1993).
29. R. W. Rodieck, "The Primate Retina," *Comparative Primate Biology, Neurosciences*, H. D. Steklis and J. Erwin (eds.), Liss, New York, 1988, pp. 203–278.
30. P. Sterling, "Retina," *The Synaptic Organization of the Brain*, G. M. Sheperd (ed.), Oxford University Press, New York, 1990, pp. 170–213.
31. H. Wässle and B. B. Boycott, "Functional Architecture of the Mammalian Retina," *Physiological Reviews*, **71**(2):447–480 (1991).
32. C. A. Curcio and K. A. Allen, "Topography of Ganglion Cells in the Human Retina," *Journal of Comparative Neurology* **300**:5–25 (1990).
33. H. Wässle et al., "Retinal Ganglion Cell Density and Cortical Magnification Factor in the Primate," *Vision Research* **30**(11):1897–1911 (1990).
34. S. W. Kuffler, "Discharge Patterns and Functional Organization of the Mammalian Retina," *Journal of Neurophysiology* **16**:37–68 (1953).
35. R. W. Rodieck, "Quantitative Analysis of Cat Retinal Ganglion Cell Response to Visual Stimuli," *Vision Research* **5**:583–601 (1965).
36. C. Enroth-Cugell and J. G. Robson, "The Contrast Sensitivity of Retinal Ganglion Cells of the Cat," *Journal of Physiology* **187**:517–552 (London, 1966).
37. A. M. Derrington and P. Lennie, "Spatial and Temporal Contrast Sensitivities of Neurones in Lateral Geniculate Nucleus of Macaque," *Journal of Physiology* **357**:219–240 (London, 1984).
38. K. Purpura et al., "Light Adaptation in the Primate Retina: Analysis of Changes in Gain and Dynamics of Monkey Retinal Ganglion Cells," *Visual Neuroscience* **4**:75–93 (1990).
39. R. M. Shapley and H. V. Perry, "Cat and Monkey Retinal Ganglion Cells and Their Visual Functional Roles," *Trends in Neuroscience* **9**:229–235 (1986).
40. B. B. Lee et al., "Luminance and Chromatic Modulation Sensitivity of Macaque Ganglion Cells and Human Observers," *Journal of the Optical Society of America* **7**(12):2223–2236 (1990).
41. B. G. Cleland and C. Enroth-Cugell, "Quantitative Aspects of Sensitivity and Summation in the Cat Retina," *Journal of Physiology* **198**:17–38 (London, 1968).
42. A. M. Derrington and P. Lennie, "The Influence of Temporal Frequency and Adaptation Level on Receptive Field Organization of Retinal Ganglion Cells in Cat," *Journal of Physiology* **333**:343–366 (London, 1982).
43. J. M. Cook et al., "Visual Resolution of Macaque Retinal Ganglion Cells," *Journal of Physiology* **396**:205–224 (London, 1988).
44. J. H. R. Maunsell and W. T. Newsome, "Visual Processing in Monkey Extrastriate Cortex," *Annual Review of Neuroscience* **10**:363–401 (1987).
45. K. Nakayama, "Biological Image Motion Processing: A Review," *Vision Research* **25**(5):625–660 (1985).
46. W. H. Merigan and J. H. R. Maunsell, "How Parallel Are the Primate Visual Pathways," *Annual Review of Neuroscience* **16**:369–402 (1993).
47. D. C. Van Essen, "Functional Organization of Primate Visual Cortex," *Cerebral Cortex*, A. A. Peters and E. G. Jones (eds.), Plenum, New York, 1985, pp. 259–329.
48. R. L. DeValois and K. K. DeValois, *Spatial Vision*, Oxford University Press, New York, 1988.
49. P. Lennie et al., "Parallel Processing of Visual Information," *Visual Perception: The Neurophysiological Foundations*, L. Spillman and J. S. Werner (eds.), Academic Press, San Diego, 1990.
50. R. L. DeValois, I. Abramov, and G. H. Jacobs, "Analysis of Response Patterns of LGN Cells," *Journal of the Optical Society of America* **56**:966–977 (1966).
51. T. N. Wiesel and D. H. Hubel, "Spatial and Chromatic Interactions in the Lateral Geniculate Body of the Rhesus Monkey," *Journal of Neurophysiology* **29**:1115–1156 (1966).
52. E. Kaplan and R. Shapley, "The Primate Retina Contains Two Types of Ganglion Cells, with High and Low Contrast Sensitivity," *Proceedings of the National Academy of Sciences* **83**:125–143 (U.S.A., 1986).
53. E. Kaplan, K. Purpura, and R. M. Shapley, "Contrast Affects the Transmission of Visual Information through the Mammalian Lateral Geniculate Nucleus," *Journal of Physiology* **391**:267–288 (London, 1987).

54. D. H. Hubel and T. N. Wiesel, "Receptive Fields and Functional Architecture of Monkey Striate Cortex," *Journal of Physiology* **195**:215–243 (London, 1968).
55. D. Hubel and T. Wiesel, "Cells Sensitive to Binocular Depth in Area 18 of the Macaque Monkey Cortex," *Nature* **225**:41–42 (1970).
56. G. F. Poggio and B. Fischer, "Binocular Interaction and Depth Sensitivity in Striate and Prestriate Cortex of Behaving Rhesus Monkey," *Journal of Neurophysiology* **40**:1392–1405 (1977).
57. R. L. DeValois, E. W. Yund, and N. Hepler, "The Orientation and Direction Selectivity of Cells in Macaque Visual Cortex," *Vision Research* **22**:531–544 (1982).
58. R. L. DeValois, D. G. Albrecht, and L. G. Thorell, "Spatial Frequency Selectivity of Cells in Macaque Visual Cortex," *Vision Research* **22**:545–559 (1982).
59. H. B. Barlow, "Critical Limiting Factors in the Design of the Eye and Visual Cortex," *Proceedings of the Royal Society of London*, series B **212**:1–34 (1981).
60. K. H. Foster et al., "Spatial and Temporal Frequency Selectivity of Neurons in Visual Cortical Areas V1 and V2 of the Macaque Monkey," *Journal of Physiology* **365**:331–363 (London, 1985).
61. G. F. Poggio, F. Gonzalez, and F. Krause, "Stereoscopic Mechanisms in Monkey Visual Cortex: Binocular Correlation and Disparity Selectivity," *Journal of Neuroscience* **8**:4531–4550 (1988).
62. H. B. Barlow, C. Blakemore, and J. D. Pettigrew, "The Neural Mechanism of Binocular Depth Discrimination," *Journal of Physiology* **193**:327–342 (London, 1967).
63. S. LeVay and T. Voigt, "Ocular Dominance and Disparity Coding in Cat Visual Cortex," *Visual Neuroscience* **1**:395–414 (1988).
64. M. S. Livingstone and D. H. Hubel, "Segregation of Form, Color, Movement and Depth: Anatomy, Physiology and Perception," *Science* **240**:740–750 (1988).
65. P. Lennie, J. Krauskopf, and G. Sclar, "Chromatic Mechanisms in Striate Cortex of Macaque," *Journal of Neuroscience* **10**:649–669 (1990).
66. M. J. Hawken and A. J. Parker, "Spatial Properties of Neurons in the Monkey Striate Cortex," *Proceedings of the Royal Society of London*, series B **231**:251–288 (1987).
67. D. B. Hamilton, D. G. Albrecht, and W. S. Geisler, "Visual Cortical Receptive Fields in Monkey and Cat: Spatial and Temporal Phase Transfer Function," *Vision Research* **29**:1285–1308 (1989).
68. B. C. Skottun et al., "Classifying Simple and Complex Cells on the Basis of Response Modulation," *Vision Research* **31**(7/8):1079–1086 (1991).
69. D. G. Albrecht and D. H. Hamilton, "Striate Cortex of Monkey and Cat: Contrast Response Function," *Journal of Neurophysiology* **48**(1):217–237 (1982).
70. D. G. Albrecht and W. S. Geisler, "Motion Selectivity and the Contrast-Response Function of Simple Cells in the Visual Cortex," *Visual Neuroscience* **7**:531–546 (1991).
71. D. J. Heeger, "Nonlinear Model of Neural Responses in Cat Visual Cortex," *Computational Models of Visual Processing*, M. S. Landy and A. Movshon (eds.), MIT Press; Cambridge, Mass., 1991, pp. 119–133.
72. D. J. Tolhurst, J. A. Movshon, and A. F. Dean, "The Statistical Reliability of Signals in Single Neurons in the Cat and Monkey Visual Cortex," *Vision Research* **23**:775–785 (1983).
73. W. S. Geisler et al., "Discrimination Performance of Single Neurons: Rate and Temporal-Pattern Information," *Journal of Neurophysiology* **66**:334–361 (1991).
74. H. L. Van Trees, *Detection, Estimation, and Modulation Theory*, Wiley, New York, 1968.
75. D. M. Green and J. A. Swets, *Signal Detection Theory and Psychophysics*, Krieger, New York, 1974.
76. H. B. Barlow, "Temporal and Spatial Summation in Human Vision at Different Background Intensities," *Journal of Physiology* **141**:337–350 (London, 1958).
77. A. Rose, "The Sensitivity Performance of the Human Eye on an Absolute Scale," *Journal of the Optical Society of America* **38**:196–208 (1948).
78. W. S. Geisler, "Sequential Ideal-Observer Analysis of Visual Discrimination," *Psychological Review* **96**:267–314 (1989).
79. D. G. Pelli, "The Quantum Efficiency of Vision," *Vision: Coding and Efficiency*, C. Blakemore (ed.), Cambridge University Press, Cambridge, 1990, pp. 3–24.
80. H. B. Barlow, "The Efficiency of Detecting Changes of Density in Random Dot Patterns," *Vision Research* **18**:637–650 (1978).

81. H. de Vries, "The Quantum Character of Light and Its Bearing upon Threshold of Vision, the Differential Sensitivity and Visual Acuity of the Eye," *Physica* **10**:553–564 (1943).
82. H. B. Barlow, "Increment Thresholds at Low Intensities Considered as Signal/Noise Discriminations," *Journal of Physiology* **136**:469–488 (London, 1957).
83. H. Helstrom, "The Detection and Resolution of Optical Signals," *IEEE Transactions on Information Theory* **IT-10**:275–287 (1964).
84. W. S. Geisler, "Physical Limits of Acuity and Hyperacuity," *Journal of the Optical Society of America A* **1**:775–782 (1984).
85. W. S. Geisler and K. Chou, "Separation of Low-Level and High-Level Factors in Complex Tasks: Visual Search," *Psychological Review* 1994 (in press).
86. M. E. Rudd, "Quantal Fluctuation Limitations on Reaction Time to Sinusoidal Gratings," *Vision Research* **28**:179–186 (1988).
87. A. E. Burgess et al., "Efficiency of Human Visual Signal Discrimination," *Science* **214**:93–94 (1981).
88. D. Kersten, "Statistical Efficiency for the Detection of Visual Noise," *Vision Research* **27**:1029–1040 (1987).
89. G. E. Legge, D. Kersten, and A.E. Burgess, "Contrast Discrimination in Noise," *Journal of the Optical Society of America A* **4**:391–404 (1987).
90. A. E. Burgess, R. F. Wagner, and R. J. Jennings, "Human Signal Detection Performance for Noisy Medical Images," *Proceedings of IEEE Computer Society International Workshop on Medical Imaging*, 1982.
91. K. J. Myers et al., "Effect of Noise Correlation on Detectability of Disk Signals in Medical Imaging," *Journal of the Optical Society of America A* **2**:1752–1759 (1985).
92. H. H. Barrett et al., "Linear Discriminants and Image Quality," *Image and Vision Computing* **10**:451–460 (1992).
93. R. H. S. Carpenter, "Eye Movements," *Vision and Visual Dysfunction*, J. Cronley-Dillon (ed.), Macmillan Press, London, 1991.
94. D. Regan and C. W. Tyler, "Some Dynamic Features of Colour Vision," *Vision Research* **11**:1307–1324 (1971).
95. G. F. Poggio and T. Poggio, "The Analysis of Stereopsis," *Annual Reviews of Neuroscience* **7**:379–412 (1984).
96. I. P. Howard, *Human Visual Orientation*, Wiley, New York, 1982.
97. I. P. Howard, "The Perception of Posture, Self Motion, and the Visual Vertical," *Handbook of Perception and Human Performance*, K. R. Boff, L. Kaufman, and J. P. Thomas (eds.), Wiley, New York, 1986.
98. I. Rock, "The Description and Analysis of Object and Event Perception," *Handbook of Perception and Human Performance*, K. R. Boff, L. Kaufman, and J. P. Thomas (eds.), Wiley, New York, 1986.
99. R. J. Watt, "Pattern Recognition by Man and Machine," *Vision and Visual Dysfunction*, J. Cronly-Dillon (ed.), Macmillan Press, London, 1991.
100. A. Treisman, "Properties, Parts, and Objects," *Handbook of Perception and Human Performance*, K. R. Boff, L. Kaufman, and J. P. Thomas (eds.), Wiley, New York, 1986.
101. M. S. Banks and P. Salapatek, "Infant Visual Perception," *Handbook of Child Psychology*, M. M. Haith and J. J. Campos (eds.), Wiley, New York, 1983.
102. K. Simons, *Normal and Abnormal Visual Development*, Springer, New York, 1993.
103. J. Marshall, "The Susceptible Visual Apparatus," *Vision and Visual Dysfunction*, J. Cronly-Dillon (ed.), Macmillan Press, London, 1991.
104. F. Ratliff and L. A. Riggs, "Involuntary Motions of the Eye During Monocular Fixation," *Journal of Experimental Psychology* **46**:687–701 (1950).
105. L. E. Arend, "Response of the Human Eye to Spatially Sinusoidal Gratings at Various Exposure Durations," *Vision Research* **16**:1311–1315 (1976).
106. D. H. Kelly, "Motion and Vision II. Stabilized Spatio-Temporal Threshold Surface," *Journal of the Optical Society of America* **69**:1340–1349 (1979).
107. O. Packer and D. R. Williams, "Blurring by Fixational Eye Movements," *Vision Research* **32**:1931–1939 (1992).
108. D. G. Green, "Regional Variations in the Visual Acuity for Interference Fringes on the Retina," *Journal of Physiology* **207**:351–356 (London, 1970).

109. J. Hirsch and R. Hylton, "Quality of the Primate Photoreceptor Lattice and Limits of Spatial Vision," *Vision Research* **24**:347–355 (1984).
110. D. M. Levi, S. A. Klein, and A. P. Aitsebaomo, "Vernier Acuity, Crowding and Cortical Magnification," *Vision Research* **25**:963–977 (1985).
111. D. S. Jacobs and C. Blakemore, "Factors Limiting the Postnatal Development of Visual Acuity in Monkeys," *Vision Research* **28**:947–958 (1987).
112. S. Shlaer, "The Relation Between Visual Acuity and Illumination," *Journal of General Physiology* **21**:165–188 (1937).
113. D. R. Williams, "Visibility of Interference Fringes Near the Resolution Limit," *Journal of the Optical Society of America A* **2**:1087–1093 (1985).
114. R. N. Bracewell, *The Fourier Transform and Its Applications: Networks and Systems*, S. W. Director (ed.), McGraw-Hill, New York, 1978.
115. N. J. Coletta and H. K. Clark, "Change in Foveal Acuity with Light Level: Optical Factors," *Ophthalmic and Visual Optics*, Optical Society of America, Monterey, Calif., 1993.
116. J. M. Enoch and V. Lakshminarayanan, "Retinal Fibre Optics," *Vision and Visual Dysfunctions*, J. Cronly-Dillon (ed.), Macmillan Press, London, 1991.
117. C. B. Blakemore and F. W. Campbell, "On the Existence of Neurones in the Human Visual System Selectively Sensitive to the Orientation and Size of Retinal Images," *Journal of Physiology* **203**:237–260 (London, 1969).
118. H. R. Wilson and J. R. Bergen, "A Four Mechanism Model for Threshold Spatial Vision," *Vision Research* **19**:19–32 (1979).
119. C. F. Stromeyer and B. Julesz, "Spatial Frequency Masking in Vision: Critical Bands and the Spread of Masking," *Journal of the Optical Society of America* **62**:1221–1232 (1972).
120. J. G. Daugman, "Uncertainty Relation for Resolution in Space, Spatial Frequency, and Orientation Optimized by Two-Dimensional Visual Cortical Filters," *Journal of the Optical Society of America* **2**:1160–1169 (1985).
121. E. R. Howell and R. F. Hess, "The Functional Area for Summation to Threshold for Sinusoidal Gratings," *Vision Research* **18**:369–374 (1978).
122. M. S. Banks, W. S. Geisler, and P. J. Bennett, "The Physical Limits of Grating Visibility," *Vision Research* **27**:1915–1924 (1987).
123. A. Rose, "The Relative Sensitivities of Television Pickup Tubes, Photographic Film, and the Human Eye," *Proceedings of the Institute of Radio Engineers* **30**:293–300 (1942).
124. N. Sekiguchi, D. R. Williams, and D. H. Brainard, "Efficiency for Detecting Isoluminant and Isochromatic Interference Fringes," *Journal of the Optical Society of America* **10**:2118–2133 (1993).
125. M. S. Banks, A. B. Sekuler, and S. J. Anderson, "Peripheral Spatial Vision: Limits Imposed by Optics, Photoreceptors, and Receptor Pooling," *Journal of the Optical Society of America* **8**:1775–1787 (1991).
126. P. T. Kortum and W. S. Geisler, "Contrast Sensitivity Functions Measured on Flashed Backgrounds in the Dark-Adapted Eye," *Annual Meeting of the Optical Society of America*, Optical Society of America, Albuquerque, NM, 1992.
127. F. L. Van Nes and M. A. Bouman, "Spatial Modulation Transfer in the Human Eye," *Journal of the Optical Society of America* **57**:401–406 (1967).
128. D. H. Kelly, "Adaptation Effects on Spatio-Temporal Sine-Wave Thresholds," *Vision Research* **12**:89–101 (1972).
129. D. C. Hood and M. A. Finkelstein, "Sensitivity to Light," *Handbook of Perception and Human Performance*, K. R. Boff, L. Kaufman, and J. P. Thomas (eds.), John Wiley and Sons, New York, 1986.
130. J. Walraven et al., "The Control of Visual Sensitivity: Receptor and Postreceptor Processes," *Visual Perception: The Neurophysiological Foundations*, L. Spillman and J. S. Werner (eds.), Academic Press, San Diego, 1990.
131. K. J. W. Craik, "The Effect of Adaptation on Subjective Brightness," *Proceedings of the Royal Society of London*, series B **128**:232–247, 1940.
132. D. C. Hood et al., "Human Cone Saturation as a Function of Ambient Intensity: A Test of Models of Shifts in the Dynamic Range," *Vision Research* **19**:983–993 (1978).
133. W. S. Geisler, "Adaptation, Afterimages and Cone Saturation," *Vision Research* **18**:279–289 (1978).

134. W. S. Geisler, "Effects of Bleaching and Backgrounds on the Flash Response of the Visual System," *Journal of Physiology* **312**:413–434 (London, 1981).
135. M. M. Hayhoe, N. E. Benimoff, and D. C. Hood, "The Time-Course of Multiplicative and Subtractive Adaptation Process," *Vision Research* **27**:1981–1996 (1987).
136. W. S. Geisler, "Mechanisms of Visual Sensitivity: Backgrounds and Early Dark Adaptation," *Vision Research* **23**:1423–1432 (1983).
137. M. M. Hayhoe, M. E. Levin, and R. J. Koshel, "Subtractive Processes in Light Adaptation," *Vision Research* **32**:323–333 (1992).
138. C. A. Burbeck and D. H. Kelly, "Role of Local Adaptation in the Fading of Stabilized Images," *Journal of the Optical Society of America A* **1**:216–220 (1984).
139. U. Tulunay-Keesey et al., "Apparent Phase Reversal During Stabilized Image Fading," *Journal of the Optical Society of America A* **4**:2166–2175 (1987).
140. A. B. Watson, "Temporal Sensitivity," *Handbook of Perception and Human Performance*, K. R. Boff, L. Kaufman, and J. P. Thomas (eds.), John Wiley and Sons, New York, 1986.
141. H. de Lange, "Research into the Dynamic Nature of the Human Fovea-Cortex Systems with Intermittent and Modulated Light. I. Attenuation Characteristics with White and Colored Light," *Journal of the Optical Society of America* **48**:777–784 (1958).
142. D. H. Kelly, "Visual Responses to Time-Dependent Stimuli. I. Amplitude Sensitivity Measurements," *Journal of the Optical Society of America* **51**:422–429 (1961).
143. J. A. J. Roufs, "Dynamic Properties of Vision-I. Experimental Relationships Between Flicker and Flash Thresholds," *Vision Research* **12**:261–278 (1972).
144. R. M. Boynton and W. S. Baron, "Sinusoidal Flicker Characteristics of Primate Cones in Response to Hererchromatic Stimuli," *Journal of the Optical Society of America* **65**:1091–1100 (1975).
145. D. H. Kelly, R. M. Boynton, and W. S. Baron, "Primate Flicker Sensitivity: Psychophysics and Electrophysiology," *Science* **194**:177–179 (1976).
146. G. E. Legge, "Sustained and Transient Mechanisms in Human Vision: Temporal and Spatial Properties," *Vision Research* **18**:69–81 (1978).
147. J. G. Robson, "Spatial and Temporal Contrast Sensitivity Functions of the Visual System," *Journal of the Optical Society of America* **56**:1141–1142 (1966).
148. J. J. Kulikowski and D. J. Tolhurst, "Psychophysical Evidence for Sustained and Transient Mechanisms in Human Vision," *Journal of Physiology*, **232**:149–163 (London, 1973).
149. C. Burbeck and D. H. Kelly, "Spatiotemporal Characteristics of Visual Mechanisms: Excitatory-Inhibitory Model," *Journal of the Optical Society of America* **70**:1121–1126 (1980).
150. M. B. Mandler and W. Makous, "A Three Channel Model of Temporal Frequency Perception," *Perception* **24**:1881–1887 (1984).
151. S. R. Lehky, "Temporal Properties of Visual Channels Measured by Masking," *Journal of the Optical Society of America A* **2**:1260–1272 (1985).
152. R. J. Snowden and R. F. Hess, "Temporal Properties of Human Visual Filters: Number, Shapes, and Spatial Covariation," *Vision Research* **32**:47–59 (1992).
153. J. J. Koenderink et al., "Perimetry of Contrast Detection Thresholds of Moving Spatial Sine Wave Patterns I. The Near Peripheral Field (0°–8°)," *Journal of the Optical Society of America* **68**:845–849 (1978).
154. V. Virsu et al., "Temporal Contrast Sensitivity and Cortical Magnification," *Vision Research* **22**:1211–1263 (1982).
155. K. T. Mullen, "The Contrast Sensitivity of Human Colour Vision to Red-Green and Blue-Yellow Chromatic Gratings," *Journal of Physiology* **359**:381–400 (London, 1985).
156. G. J. C. van der Horst and M. A. Bouman, "Spatiotemporal Chromaticity Discrimination," *Journal of the Optical Society of America* **59**:1482–1488 (1969).
157. M. S. Banks and P. J. Bennett, "Optical and Photoreceptor Immaturities Limit the Spatial and Chromatic Vision of Human Neonates," *Journal of the Optical Society of America A: Optics and Image Science* **5**:2059–2079 (1988).
158. J. R. Jordan, W. S. Geisler, and A. C. Bovik, "Color as a Source of Information in the Stereo Correspondence Process," *Vision Research* **30**:1955–1970 (1990).

159. P. Lennie and M. D’Zmura, “Mechanisms of Color Vision,” *CRC Critical Reviews in Neurobiology* **3**:333–400 (1988).
160. D. H. Kelly, “Luminous and Chromatic Flickering Patterns Have Opposite Effects,” *Science* **188**:371–372 (1975).
161. G. E. Legge and J. M. Foley, “Contrast Masking in Human Vision,” *Journal of the Optical Society of America* **70**:1458–1470 (1980).
162. D. G. Pelli, “Uncertainty Explains Many Aspects of Visual Contrast Detection and Discrimination,” *Journal of the Optical Society of America A* **2**:1508–1532 (1985).
163. A. Bradley and I. Ohzawa, “A Comparison of Contrast Detection and Discrimination,” *Vision Research* **26**:991–997 (1986).
164. G. E. Legge and D. Kersten, “Contrast Discrimination in Peripheral Vision,” *Journal of the Optical Society of America A* **4**:1594–1598 (1987).
165. F. W. Campbell and J. J. Kulikowski, “Orientation Selectivity of the Human Visual System,” *Journal of Physiology* **187**:437–445 (London, 1966).
166. N. Graham, *Visual Pattern Analyzers*, Oxford, New York, 1989.
167. H. R. Wilson, D. K. McFarlane, and G. C. Philips, “Spatial Frequency Tuning of Orientation Selective Units Estimated by Oblique Masking,” *Vision Research* **23**:873–882 (1983).
168. H. R. Wilson, “A Transducer Function for Threshold and Suprathreshold Human Vision,” *Biological Cybernetics* **38**:171–178 (1980).
169. H. R. Wilson, “Psychophysics of Contrast Gain,” *Annual Meeting of the Association for Research in Vision and Ophthalmology*, Investigative Ophthalmology and Visual Science, Sarasota, Fla., 1990.
170. J. M. Foley, “Human Luminance Pattern Vision Mechanisms: Masking Experiments Require a New Theory,” *Journal of the Optical Society of America A* (1994), in press.
171. R. F. Hess and R. J. Snowden, “Temporal Properties of Human Visual Filters: Number, Shapes, and Spatial Covariation,” *Vision Research* **32**:47–59 (1992).
172. M. A. Georgeson and G. D. Sullivan, “Contrast Constancy: Deblurring in Human Vision by Spatial Frequency Channels,” *Journal of Physiology* **252**:627–656 (London, 1975).
173. D. B. Gennery, “Determination of Optical Transfer Function by Inspection of Frequency-Domain Plot,” *Journal of the Optical Society of America* **63**:1571–1577 (1973).
174. D. D. Michaels, *Visual Optics and Refraction: A Clinical Approach*, Mosby, St. Louis, 1980.
175. W. N. Charman, “Visual Standards for Driving,” *Ophthalmic and Physiological Optics* **5**:211–220 (1985).
176. G. Westheimer, “Scaling of Visual Acuity Measurements,” *Archives of Ophthalmology* **97**:327–330 (1979).
177. C. Owsley et al., “Visual/Cognitive Correlates of Vehicle Accidents in Older Drivers,” *Psychology and Aging* **6**:403–415 (1991).
178. R. N. Berry, “Quantitative Relations among Vernier, Real Depth and Stereoscopic Depth Acuties,” *Journal of Experimental Psychology* **38**:708–721 (1948).
179. G. Westheimer and S. P. McKee, “Spatial Configurations for Visual Hyperacuity,” *Vision Research* **17**:941–947 (1977).
180. G. Westheimer, “Visual Acuity and Hyperacuity,” *Investigative Ophthalmology* **14**:570–572 (1975).
181. G. Westheimer, “The Spatial Grain of the Perifoveal Visual Field,” *Vision Research* **22**:157–162 (1982).
182. D. M. Levi and S. A. Klein, “Vernier Acuity, Crowding and Amblyopia,” *Vision Research* **25**:979–991 (1985).
183. S. A. Klein and D. M. Levi, “Hyperacuity Thresholds of 1 sec: Theoretical Predictions and Empirical Validation,” *Journal of the Optical Society of America A* **2**:1170–1190 (1985).
184. H. R. Wilson, “Model of Peripheral and Amblyopic Hyperacuity,” *Vision Research* **31**:967–982 (1991).
185. H. B. Barlow, “Reconstructing the Visual Image in Space and Time,” *Nature* **279**:189–190 (1979).
186. F. H. C. Crick, D. C. Marr, and T. Poggio, “An Information-Processing Approach to Understanding the Visual Cortex,” *The Organization of the Cerebral Cortex*, F. O. Smith (ed.), MIT Press, Cambridge, 1981.
187. R. J. Watt and M. J. Morgan, “A Theory of the Primitive Spatial Code in Human Vision,” *Vision Research* **25**:1661–1674 (1985).
188. R. Shapley and J. D. Victor, “Hyperacuity in Cat Retinal Ganglion Cells,” *Science* **231**:999–1002 (1986).

189. H. R. Wilson, "Responses of Spatial Mechanisms Can Explain Hyperacuity," *Vision Research* **26**:453–469 (1986).
190. Q. M. Hu, S. A. Klein, and T. Carney, "Can Sinusoidal Vernier Acuity Be Predicted by Contrast Discrimination?" *Vision Research* **33**:1241–1258 (1993).
191. W. S. Geisler and K. D. Davila, "Ideal Discriminators in Spatial Vision: Two-Point Stimuli," *Journal of the Optical Society of America A* **2**:1483–1497 (1985).
192. S. Appelle, "Perception and Discrimination As a Function of Stimulus Orientation: The 'Oblique Effect' in Man and Animals," *Psychological Bulletin* **78**:266–278 (1972).
193. F. W. Campbell, R. H. S. Carpenter, and J. Z. Levinson, "Visibility of Aperiodic Patterns Compared with That of Sinusoidal Gratings," *Journal of Physiology* **204**:283–298 (London, 1969).
194. F. W. Campbell, J. Nachmias, and J. Jukes, "Spatial Frequency Discrimination in Human Vision," *Journal of the Optical Society of America* **60**:555–559 (1970).
195. A. V. Oppenheim and J. S. Lim, "The Importance of Phase in Signals," *Proceedings of the IEEE* **69**:529–541 (1981).
196. D. J. Field, "Relations Between the Statistics of Natural Images and the Response Properties of Cortical Cells," *Journal of the Optical Society of America A* **4**:2379–2394 (1987).
197. J. Huang and D. L. Turcotte, "Fractal Image Analysis: Application to the Topography of Oregon and Synthetic Images," *Journal of the Optical Society of America A* **7**:1124–1129 (1990).
198. M. J. Morgan, J. Ross, and A. Hayes, "The Relative Importance of Local Phase and Local Amplitude in Patchwise Image Reconstruction," *Biological Cybernetics* **65**:113–119 (1991).
199. D. R. Badcock, "Spatial Phase or Luminance Profile Discrimination," *Vision Research* **24**:613–623 (1984).
200. I. Rentschler and B. Treutwein, "Loss of Spatial Phase Relationships in Extrafoveal Vision," *Nature* **313**:308–310 (1985).
201. R. F. Hess and J. S. Pointer, "Evidence for Spatially Local Computations Underlying Discrimination of Periodic Patterns in Fovea and Periphery," *Vision Research* **27**:1343–1360 (1987).
202. P. J. Bennett and M. S. Banks, "Sensitivity Loss Among Odd-Symmetric Mechanisms and Phase Anomalies in Peripheral Vision," *Nature* **326**:873–876 (1987).
203. B. Julesz, E. N. Gilbert, and J. D. Victor, "Visual Discrimination of Textures with Identical Third-Order Statistics," *Biological Cybernetics* **31**:137–147 (1978).
204. I. Rentschler, M. Hubner, and T. Caelli, "On the Discrimination of Compound Gabor Signals and Textures," *Vision Research* **28**:279–291 (1988).
205. G. E. Legge et al., "Psychophysics of Reading I. Normal Vision," *Vision Research* **25**:239–252 (1985).
206. M. C. Morrone and D. C. Burr, "Feature Detection in Human Vision: A Phase Dependent Energy Model," *Proceedings of the Royal Society of London, series B* **235**:221–245 (1988).
207. D. J. Field and J. Nachmias, "Phase Reversal Discrimination," *Vision Research* **24**:333–340 (1984).
208. D. C. Burr, M. C. Morrone, and D. Spinelli, "Evidence for Edge and Bar Detectors in Human Vision," *Vision Research* **29**:419–431 (1989).
209. P. J. Bennett and M. S. Banks, "The Effects of Contrast, Spatial Scale, and Orientation on Foveal and Peripheral Phase Discrimination," *Vision Research* **31**:1759–1786, 1991.
210. A. Toet and D. M. Levi, "The Two-Dimensional Shape of Spatial Interaction Zones in the Parafovea," *Vision Research* **32**:1349–1357 (1992).
211. M. C. Morrone, D. C. Burr, and D. Spinelli, "Discrimination of Spatial Phase in Central and Peripheral Vision," *Vision Research* **29**:433–445 (1989).
212. S. Anstis, "Motion Perception in the Frontal Plane: Sensory Aspects," *Handbook of Perception and Human Performance*, K. R. Boff, L. Kaufman, and J. P. Thomas (eds.), John Wiley and Sons, New York, 1986.
213. C. C. Hildreth and C. Koch, "The Analysis of Visual Motion: From Computational Theory to Neuronal Mechanisms," *Annual Review of Neuroscience* **10**:477–533 (1987).
214. D. C. Burr and J. Ross, "Contrast Sensitivity at High Velocities," *Vision Research* **22**:479–484 (1982).
215. S. P. McKee, G. H. Silverman, and K. Nakayama, "Precise Velocity Discrimination despite Random Variations in Temporal Frequency and Contrast," *Vision Research* **26**(4):609–619 (1986).
216. B. De Bruyn and G. A. Orban, "Human Velocity and Direction Discrimination Measured with Random Dot Patterns," *Vision Research* **28**(12):1323–1335 (1988).

217. S. N. J. Watamaniuk, R. Sekuler, and D. W. Williams, "Direction Perception in Complex Dynamic Displays: The Integration of Direct Information," *Vision Research* **29**:47–59 (1989).
218. K. Nakayama and G. H. Silverman, "Detection and Discrimination of Sinusoidal Grating Displacements," *Optical Society of America* **2**(2):267–274 (1985).
219. A. Pantle, "Temporal Frequency Response Characteristics of Motion Channels Measured with Three Different Psychophysical Techniques," *Perception and Psychophysics* **24**:285–294 (1978).
220. P. Cavanagh and P. Anstis, "The Contribution of Color to Motion in Normal and Color-Deficient Observers," *Vision Research* **31**:2109–2148 (1991).
221. D. T. Lindsey and D. Y. Teller, "Motion at Isoluminance: Discrimination/Detection Ratios for Moving Isoluminant Gratings," *Vision Research* **30**(11):1751–1761 (1990).
222. S. P. McKee and K. Nakayama, "The Detection of Motion in the Peripheral Visual Field," *Vision Research* **24**:25–32 (1984).
223. O. J. Braddick, "Low-Level and High-Level Processes in Apparent Motion," *Philosophical Transactions of the Royal Society of London B* **290**:137–151 (1980).
224. J. J. Chang and B. Julesz, "Displacement Limits for Spatial Frequency Filtered Random-Dot Cinematograms in Apparent Motion," *Vision Research* **23**(12):1379–1385 (1983).
225. C. L. Baker and O. J. Braddick, "Does Segregation of Differently Moving Areas Depend on Relative or Absolute Displacement?" *Vision Research* **22**:851–856 (1982).
226. L. Glass, "Moire Effect from Random Dots," *Nature* **243**:578–580 (1969).
227. J. J. Gibson, *The Senses Considered as Perceptual Systems*, Houghton Mifflin, Boston, 1966.
228. J. J. Gibson, *The Perception of the Visual World*, Houghton Mifflin, Boston, 1950.
229. J. E. Cutting, *Perception with an Eye to Motion*, MIT Press, Cambridge, 1986.
230. W. H. Warren, M. W. Morris, and M. Kalish, "Perception of Translational Heading from Optical Flow," *Journal of Experimental Psychology: Human Perception and Performance* **14**:646–660 (1988).
231. W. H. Warren and D. J. Hannon, "Eye Movements and Optical Flow," *Journal of the Optical Society of America A* **7**:160–169 (1990).
232. E. von Hoist, "Relations between the Central Nervous System and the Peripheral Organs," *Animal Behavior* **2**:89–94 (1954).
233. D. J. Heeger and A. D. Jepson, "Subspace Methods for Recovering Rigid Motion. I: Algorithm and Implementation," *University of Toronto Technical Reports on Research in Biological and Computational Vision*, RBCV-TR-90-35.
234. H. C. Longuet-Higgins and K. Prazdny, "The Interpretation of a Moving Retinal Image," *Proceedings of the Royal Society of London B* **208**:385–397 (1980).
235. J. A. Perrone, "Model for Computation of Self-Motion in Biological Systems," *Journal of the Optical Society of America A* **9**:177–194 (1992).
236. J. H. Rieger and D. T. Lawton, "Processing Differential Image Motion," *Journal of the Optical Society of America A* **2**:354–360 (1985).
237. C. S. Royden, M. S. Banks, and J. A. Crowell, "The Perception of Heading during Eye Movements," *Nature* **360**:583–585 (1992).
238. A. V. van den Berg, "Robustness of Perception of Heading from Optic Flow," *Vision Research* **32**:1285–1296 (1992).
239. J. J. Gibson, P. Olum, and F. Rosenblatt, "Parallax and Perspective During Aircraft Landings," *American Journal of Psychology* **68**:373–385 (1955).
240. B. J. Rogers and M. Graham, "Similarities Between Motion Parallax and Stereopsis in Human Depth Perception," *Vision Research* **22**:261–270 (1982).
241. B. J. Rogers and M. F. Bradshaw, "Vertical Disparities, Differential Perspective and Binocular Stereopsis," *Nature* **361**:253–255 (1993).
242. J. E. W. Mayhew and H. C. Longuet-Higgins, "A Computational Model of Binocular Depth Perception," *Nature* **297**:376–378 (1982).
243. G. Westheimer and S. P. McKee, "What Prior Unocular Processing Is Necessary for Stereopsis?" *Investigative Ophthalmology and Visual Science* **18**:614–621 (1979).

244. B. Julesz, "Binocular Depth Perception of Computer-Generated Patterns," *Bell System Technical Journal* **39**:1125–1162 (1960).
245. C. M. Aschenbrenner, "Problems in Getting Information Into and Out of Air Photographs," *Photogrammetric Engineering*, **20**:398–401 (1954).
246. B. Julesz, "Stereoscopic Vision," *Vision Research* **26**:1601–1612 (1986).
247. J. P. Frisby and J. E. W. Mayhew, "Contrast Sensitivity Function for Stereopsis," *Perception* **7**:423–429 (1978).
248. C. Blakemore and B. Hague, "Evidence for Disparity Detecting Neurones in the Human Visual System," *Journal of Physiology* **225**:437–445 (London, 1972).
249. T. B. Felton, W. Richards, and R. A. Smith, "Disparity Processing of Spatial Frequencies in Man," *Journal of Physiology* **225**:349–362 (London, 1972).
250. B. Julesz and J. E. Miller, "Independent Spatial Frequency Tuned Channels in Binocular Fusion and Rivalry," *Perception* **4**:315–322 (1975).
251. Y. Yang and R. Blake, "Spatial Frequency Tuning of Human Stereopsis," *Vision Research* **31**:1177–1189 (1991).
252. J. S. Mansfield and A. J. Parker, "An Orientation-Tuned Component in the Contrast Masking of Stereopsis," *Vision Research* **33**:1535–1544 (1993).
253. R. D. Freeman and I. Ohzawa, "On the Neurophysiological Organization of Binocular Vision," *Vision Research* **30**:1661–1676 (1990).
254. D. L. Halpern and R. R. Blake, "How Contrast Affects Stereoacuity," *Perception* **17**:483–495 (1988).
255. G. E. Legge and Y. Gu, "Stereopsis and Contrast," *Vision Research* **29**:989–1004 (1989).
256. L. K. Cormack, S. B. Stevenson, and C. M. Schor, "Interocular Correlation, Luminance Contrast, and Cyclopean Processing," *Vision Research* **31**:2195–2207 (1991).
257. C. Blakemore, "The Range and Scope of Binocular Depth Discrimination in Man," *Journal of Physiology* **211**:599–622 (London, 1970).
258. D. R. Badcock and C. M. Schor, "Depth-Increment Detection Function for Individual Spatial Channels," *Journal of the Optical Society of America A* **2**:1211–1216 (1985).
259. W. Richards, "Anomalous Stereoscopic Depth Perception," *Journal of the Optical Society of America* **61**:410–414 (1971).
260. L. K. Cormack, S. B. Stevenson, and C. M. Schor, "Disparity-Tuned Channels of the Human Visual System," *Visual Neuroscience* **10**:585–596 (1993).
261. S. B. Stevenson, L. K. Cormack, and C. M. Schor, "Disparity Tuning in Mechanisms of Human Stereopsis," *Vision Research* **32**:1685–1694 (1992).
262. S. R. Lehky and T. J. Sejnowski, "Neural Model of Stereoacuity and Depth Interpolation Based on a Distributed Representation of Stereo Disparity," *Journal of Neuroscience* **10**:2281–2299 (1990).
263. C. H. Bailey and P. Gouras, "The Retina and Phototransduction," *Principles of Neural Science*, E. R. Kandel and J. H. Schwartz (eds.), Elsevier Science Publishing Co., Inc., New York, 1985, pp. 344–355.
264. K. D. Davila and W. S. Geisler, "The Relative Contributions of Pre-Neural and Neural Factors to Areal Summation in the Fovea," *Vision Research* **31**:1369–1380 (1991).
265. J. Ross and H. D. Speed, "Contrast Adaptation and Contrast Masking in Human Vision," *Proceedings of the Royal Society of London, series B* **246**:61–69 (1991).

This page intentionally left blank.

DO NOT DUPLICATE

PSYCHOPHYSICAL METHODS

Denis G. Pelli

*Psychology Department and Center for Neural Science
New York University
New York*

Bart Farell

*Institute for Sensory Research
Syracuse University
Syracuse, New York*

3.1 INTRODUCTION

Psychophysical methods are the tools for measuring perception and performance. These tools are used to reveal basic perceptual processes, to assess observer performance, and to specify the required characteristics of a display. We are going to ignore this field's long and interesting history,¹ and much theory as well.^{2,3} Here we present a formal treatment, emphasizing the theoretical concepts of psychophysical measurement. For practical advice in setting up an experiment, please turn to our user's guide.⁴ Use the supplied references for further reading.

Consider the psychophysical evaluation of the suitability of a visual display for a particular purpose. A home television to be used for entertainment is most reasonably assessed in a "beauty contest" of subjective preference,⁵ whereas a medical imaging display must lead to accurate diagnoses^{6,7} and military aerial reconnaissance must lead to accurate vehicle identifications.⁸ In our experience, the first step toward defining a psychophysically answerable question is to formulate the problem as a task that the observer must perform. One can then assess the contribution of various display parameters toward that performance. Where precise parametric assessment is desired it is often useful to substitute a simple laboratory task for the complex real-life activity, provided one can either demonstrate, or at least reasonably argue, that the laboratory results are predictive.

Psychophysical measurement is usually understood to mean measurement of behavior to reveal internal processes. The experimenter is typically not interested in the behavior itself, such as pressing a button, which merely communicates a decision by the observer about the stimulus.* This chapter reviews the various decision tasks that may be used to measure perception and performance and evaluates their strengths and weaknesses. We begin with definitions and a brief review of visual stimuli. We then explain and evaluate the various psychophysical tasks, and end with some practical tips.

*Psychophysical measurement can also be understood to include noncommunicative physiological responses such as pupil size, eye position, electrical potentials measured on the scalp and face, and even BOLD fMRI responses in the brain, which might be called "unintended" responses. (These examples are merely suggestive, not definitive. Observers can decide to move their eyes and, with feedback, can learn to control many other physiological responses. Responses are "unintended" only when they are not used for overt communication by the observer.) Whether these unintended responses are called psychophysical or physiological is a matter of taste. In any case, decisions are usually easier to measure and interpret, but unintended responses may be preferred in certain cases, as when assessing noncommunicative infants and animals.

3.2 DEFINITIONS

At the highest level, an *experiment* answers a question about how certain “experimental conditions” affect observer performance. *Experimental conditions* include stimulus parameters, observer instruction, and anything else that may affect the observer’s state. Experiments are usually made up of many individual measurements, called “trials,” under each experimental condition. Each trial presents a stimulus and collects a response—a decision—from the observer.

There are two kinds of decision tasks: judgments and adjustments. It is useful to think of one as the inverse of the other. In one case the experimenter gives the observer a stimulus and asks for a classification of the stimulus or percept; in the other case the experimenter, in effect, gives the observer a classification and asks for an appropriate stimulus back. Either the experimenter controls the stimulus and the observer makes a *judgment* based on the resulting percept, or the observer *adjusts* the stimulus to satisfy a perceptual criterion specified by the experimenter (e.g., match a sample). Both techniques are powerful. Adjustments are intrinsically subjective (because they depend on the observers’ understanding of the perceptual criterion), yet they can often provide good data quickly and are to be preferred when applicable. But not all questions can be formulated as adjustment tasks. Besides being more generally applicable, judgments are often easier to analyze, because the stimulus is under the experimenter’s control and the task may be objectively defined. Observers typically like doing adjustments and find judgments tedious, partly because judgment experiments usually take much longer.

An obvious advantage of adjustment experiments is that they measure physical stimulus parameters, which may span an enormous dynamic range and typically have a straightforward physical interpretation. Judgment tasks measure human performance (e.g., frequency of seeing) as a function of experimental parameters (e.g., contrast). This is appropriate if the problem at hand concerns human performance per se. For other purposes, however, raw measures of judgment performance typically have a very limited useful range, and a scale that is hard to interpret. Having noted that adjustment and judgment tasks may be thought of as inverses of one another, we hasten to add that in practice they are often used in similar ways. Judgment experiments often vary a stimulus parameter on successive trials in order to find the value that yields a criterion judgment. These “sequential estimation methods” are discussed in Sec. 3.6. The functional inversion offered by sequential estimation allows judgment experiments to measure a physical parameter as a function of experimental condition, like adjustment tasks, while retaining the judgment task’s more rigorous control and interpretation.

Distinguishing between judgment and adjustment tasks emphasizes the kind of response that the observer makes. It is also possible to subdivide tasks in a way that emphasizes the stimuli and the question posed. In a *detection* task there may be any number of alternative stimuli, but one is a blank, and the observer is asked only to distinguish between the blank and the other stimuli. Slightly more general, a *discrimination* task may also have any number of alternative stimuli, but one of the stimuli, which need not be blank, is designated as the reference, and the observer is asked only to distinguish between the reference and other stimuli. A decision that distinguishes among more than two categories is usually called an *identification* or *classification*.⁹ All decision tasks allow for alternative responses, but two alternatives is an important special case.¹⁰

As normally used, the choice of term, *detection* or *discrimination*, says more about the experimenter’s way of thinking than it does about the actual task faced by the observer. This is because theoretical treatments of detection and discrimination usually allow for manipulation of the experimental condition by introduction of an extraneous element, often called a “mask” or “pedestal,” that is added to every stimulus. Thus, one is always free to consider a discrimination task as detection in the presence of a mask. This shift in perspective can yield new insights (e.g., Refs. 11–15). Since there is no fundamental difference between detection and discrimination,¹⁶ we have simplified the presentation below by letting detection stand in for both. The reader may freely substitute “reference” for “blank” (or suppose the presence of an extraneous mask) in order to consider the discrimination paradigm.

The idea of “threshold” plays a large role in psychophysics. Originally deterministic, *threshold* once referred to the stimulus intensity above which the stimulus was always distinguishable from blank, and below which it was indistinguishable from blank. In a discrimination task one might refer to a “discrimination threshold” or a “just-noticeable difference.” Nowadays the idea is statistical; we know that the observer’s probability of correct classification rises as a continuous function of stimulus

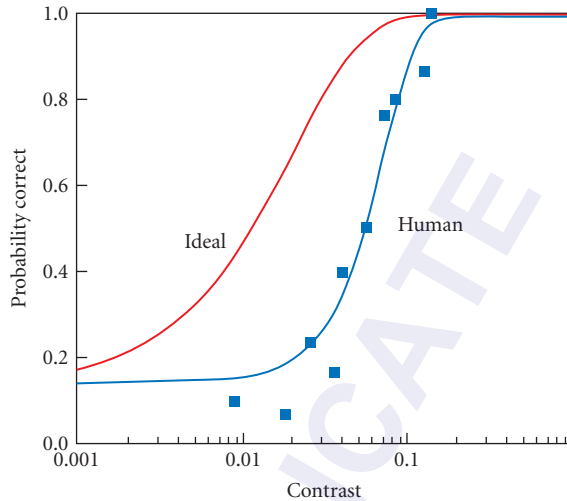


FIGURE 1 Probability of correctly identifying a letter in noise, as a function of letter contrast. The letters are bandpass filtered. Gaussian noise was added independently to each pixel. Each symbol represents the proportion correct in 30 trials. The solid curve through the points is a maximum likelihood fit of a Weibull function. The other curve represents a similar maximum likelihood fit to the performance of a computer program that implements the ideal letter classifier.⁴⁰ Efficiency, the squared ratio of threshold contrasts, is 9 percent. (Courtesy of Joshua A. Solomon.)

intensity (see Fig. 1). Threshold is defined as the stimulus intensity (e.g., contrast) corresponding to an arbitrary level of performance (e.g., 82 percent correct). However, the old intuition, now called a “high threshold,” still retains a strong hold on everyone’s thinking for the good reason that the transition from invisible to visible, though continuous, is quite abrupt, less than a factor of two in contrast.

Most psychophysical research has concentrated on measuring thresholds. This has been motivated by a desire to isolate low-level sensory mechanisms by using operationally defined tasks that are intended to minimize the roles of perception and cognition. This program is generally regarded as successful—visual detection is well understood (e.g., Ref. 17)—but leaves most of our visual experience and ability unexplained. This has stimulated a great deal of experimentation with suprathreshold stimuli and nondetection tasks in recent years.

3.3 VISUAL STIMULI

Before presenting the tasks, which are general to all sense modalities (not just vision), it may be helpful to briefly review the most commonly used visual stimuli. Until the 1960s most vision research used a spot as the visual stimulus (e.g., Ref. 18). Then cathode ray tube displays made it easy to generate more complex stimuli, especially sinusoidal gratings, which provided the first evidence for multiple “spatial frequency channels” in vision.¹⁹ Sinusoidal grating patches have two virtues. A sinusoid at the display always produces a sinusoidal image on the retina.* And most visual mechanisms are selective in space and in spatial frequency, so it is useful to have a stimulus that is restricted in both domains.

*This is strictly true only within an isoplanatic patch, i.e., a retinal area over which the eye’s optical point spread function is unchanged.

Snellen,²⁰ in describing his classic eye chart, noted the virtue of letters as visual stimuli—they offer a large number of stimulus alternatives that are readily identifiable.^{21,22} Other commonly used stimuli include annuli, lines, arrays of such elements, and actual photographs of faces, nature, and military vehicles. There are several useful texts on image quality, emphasizing signal-to-noise ratio.^{23–26} Finally, there has been some psychophysical investigation of practical tasks such as reading,²⁷ flying an airplane,²⁸ or shopping in a supermarket.²⁹

The stimulus alternatives used in vision experiments are usually parametric variations along a single dimension, most commonly contrast, but frequently size and position in the visual field. *Contrast* is a dimensionless ratio: the amplitude of the luminance variation within the stimulus, normalized by the background luminance. Michelson contrast (used for gratings) is the maximum minus the minimum luminance divided by the maximum plus the minimum. Weber contrast (used for spots and letters) is the maximum deviation from the uniform background divided by the background luminance. RMS contrast is the root-mean-square deviation of the stimulus luminance from the mean luminance, divided by the mean luminance.

3.4 ADJUSTMENTS

Adjustment tasks require that the experimenter specify a perceptual criterion to the observer, who adjusts the stimulus to satisfy the criterion. Doubts about the observer's interpretation of the criterion may confound interpretation of the results. The adjustment technique is only as useful as the criterion is clear.

Threshold

Figure 2 shows contrast sensitivity (the reciprocal of the threshold contrast) for a sinusoidal grating as a function of spatial and temporal frequency.³⁰ These thresholds were measured by what is probably the most common form of the adjustment task, which asks the observer to adjust the stimulus contrast up and down to the point where it is “just barely detectable.” While some important studies have collected their data in this way, one should bear in mind that this is a vaguely specified criterion. What should the observer understand by “barely” detectable? Seen half the time? In order to adjust to threshold, the observer must form a subjective interpretation and apply it to the changing percept. It is well known that observers can be induced (e.g., by coaching) to raise or lower their criterion, and when comparing among different observers it is important to bear in mind that social and personality factors may lead to systematically different interpretations of the same vague instructions. Nevertheless, these subjective effects are relatively small (about a factor of two in contrast) and many questions can usefully be addressed, in at least a preliminary way, by quick method-of-adjustment threshold settings. Alternatively, one might ignore the mean of the settings and instead use the standard deviation to estimate the observer's discrimination threshold.³¹

Nulling

Of all the many kinds of adjustments, nulling is the most powerful. Typically, there is a simple basic stimulus that is distorted by some experimental manipulation, and the observer is given control over the stimulus and asked to adjust it so as to cancel the distortion (e.g., Ref. 32). The absence of a specific kind of distortion is usually unambiguous and easy for the observer to understand, and the observer's null setting is typically very reliable.

Matching

Two stimuli are presented, and the observer is asked to adjust one to match the other. Sometimes the experiment can be designed so that the observer can achieve a perfect match in which the stimuli are utterly indistinguishable, which Brindley³³ calls a “Class A” match. Usually, however, the stimuli

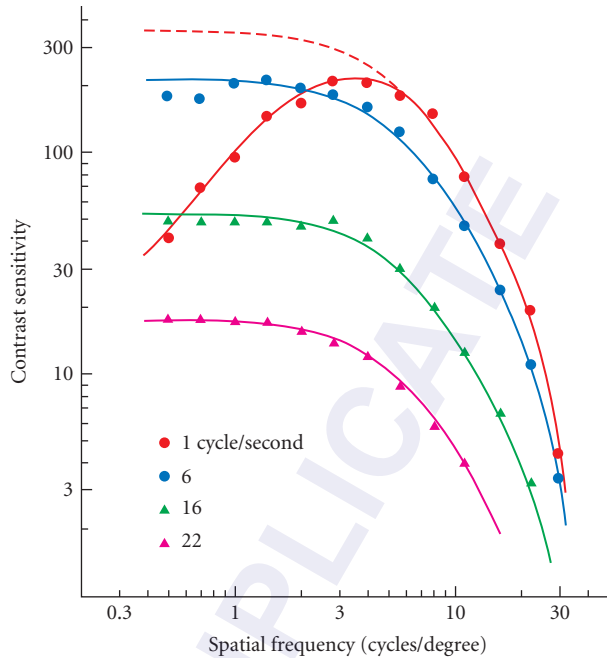


FIGURE 2 Spatial contrast sensitivity (reciprocal of threshold contrast) functions for sinusoidal gratings temporally modulated (flickered) at several temporal frequencies. The points are the means of four method-of-adjustment measurements and the curves (one with a dashed low-frequency section) differ only in their positions along the contrast-sensitivity scale. (From Robson,³⁰)

are obviously different and the observer is asked to match only a particular aspect of the stimuli, which is called a “Class B” match. For example, the observer might be shown two grating patches, one fine and one coarse, and asked to adjust the contrast of one to match the contrast of the other.³⁴ Or the observer might see two uniform patches of different colors and be asked to match their brightnesses.³⁵ Observers (and reviewers for publication) are usually comfortable with matching tasks, but, as Brindley points out, it is amazing that observers can seemingly abstract and compare a particular parameter of the multidimensional stimuli in order to make a Class B match. Matching tasks are extremely useful, but conclusions based on Class B matches may be less secure than those based on Class A matches because our understanding of how the observer does the task is less certain.

Magnitude Production

The observer is asked to adjust a stimulus to match a numerically specified perceptual criterion, e.g., “as bright as a 60-watt light bulb.” The number may have a scale (watts in this case) or be a pure number.² The use of pure numbers, without any scale, to specify a perceptual criterion is obviously formally ambiguous, but in practice many experimenters report that observers seem comfortable with such instructions and produce stable results that are even reasonably consistent among different observers. Magnitude production, however, is rarely used in visual psychophysics research.

3.5 JUDGMENTS

Judgment tasks ask the observer to classify the stimulus or percept. They differ primarily in the number of alternative stimuli that may be presented on a given trial and the number of alternative responses that the observer is allowed.

The Ideal Observer

When the observer is asked to classify the stimulus (not the percept) it may be useful to consider the mathematically defined ideal classifier that would yield the most accurate performance using only the information (the stimuli and their probabilities) available to the observer.^{36–42} Obviously this would be an empty exercise unless there is some known factor that makes the stimuli hard to distinguish. Usually this will be visual noise: random variations in the stimulus, random statistics of photon absorptions in the observer's eyes, or random variations in neural processes in the observer's visual system. If the stimuli plus noise can be defined statistically at some site—at the display, as an image at the observer's retinae, as a pattern of photon absorptions, or as a spatiotemporal pattern of neural activity—then one can solve the problem mathematically and compute the highest attainable level of performance. This ideal often provides a useful point of comparison in thinking about the actual human observer's results. A popular way of expressing such a comparison is to compute the human observer's efficiency, which will be a number between 0 and 1. For example, in Fig. 1 at threshold the observer's efficiency for letter identification is 9 percent. As a general rule, the exercise of working out the ideal and computing the human observer's efficiency is usually instructive but, obviously, low human efficiencies should be interpreted as a negative result, suggesting that the ideal is not particularly relevant to understanding how the human observer does the task.

Yes-No

The best-known judgment task is yes-no. It is usually used for detection, although it is occasionally used for discrimination. The observer is either asked to classify the stimulus, "Was a nonblank stimulus present?" or classify the percept, "Did you see it?" The observer is allowed only two response alternatives: yes or no. There may be any number of alternative stimuli. If the results are to be compared with those of an ideal observer, then the kind of stimulus, blank or nonblank, must be unpredictable.

As with the method-of-adjustment thresholds discussed above, the question posed in a yes-no experiment is fundamentally ambiguous. Where is the dividing line between yes and no on the continuum of internal states between the typical percepts generated by the blank and nonblank stimuli? Theoretical considerations and available evidence suggest that observers act as if they reduced the percept to a "decision variable," a pure magnitude—a number if you like—and compared that magnitude with an internal criterion that is under their conscious control.^{40,43} Normally we are not interested in the criterion, yet it is troublesome to remove its influence on the results, especially since the criterion may vary between experimental conditions and observers. For this reason, most investigators no longer use yes-no tasks.

As discussed next, this pesky problem of the observer's subjective criterion can be dealt with explicitly, by using "rating scale" tasks, or banished, by using unbiased "two-alternative forced choice" (2afc) tasks. Rating scale is much more work, and unless the ratings themselves are of interest, the end result of using either rating scale or 2afc is essentially the same.

Rating Scale

In a rating scale task the observer is asked to rate the likelihood that a nonblank stimulus was presented. There must be blank and nonblank stimulus alternatives, and there may be any number of alternative ratings—five is popular—but even a continuous scale may be allowed.⁴⁴ The endpoints of

the rating scale are “The stimulus was definitely blank” and “The stimulus was definitely nonblank,” with intermediate degrees of confidence in between. The results are graphed as a receiver operating characteristic, or ROC, that plots one conditional probability against another. The observer’s ratings are transformed into yes-no judgments by comparing them with an external criterion. Ratings above the criterion become “yes” and those below the criterion become “no.” This transformation is repeated for all possible values of the external criterion. Finally, the experimenter plots—for each value of the criterion—the probability of a yes when a nonblank stimulus was present (a “hit”) against the probability of a yes when a blank stimulus was present (a “false alarm”). In medical contexts the hit rate is called “sensitivity” and one minus the false alarm rate is called “specificity.” Figure 3 shows an ROC curve for a medical diagnosis;⁴⁵ radiologists examined mammograms and rated the likelihood that a lesion was benign or malignant.

In real-life applications the main value of ROC curves is that they can be used to optimize yes-no decisions based on ratings, e.g., whether to refer a patient for further diagnosis or treatment. However, this requires knowledge of the prior stimulus probabilities (e.g., in Fig. 3, the incidence of disease in the patient population), the benefit of a hit, and the cost of a false alarm.^{6,7} These conditions are rarely met. One usually can estimate prior probability and assess the cost of the wasted effort caused by the false alarms, but it is hard to assign a commensurate value to the hits, which may save lives through timely treatment.

The shape of the ROC curve has received a great deal of attention in the theoretical detection literature, and there are various mathematical models of the observer’s detection process that can account for the shape.^{46–48} However, unless the actual situation demands rating-based decisions, the ROC shape has little or no practical significance, and the general practice is to summarize the ROC curve by the area under the curve. The area is 0.5 when the observers’ ratings are independent of the stimuli (i.e., useless guessing). The area can be at most 1—when the observer makes no mistakes. The

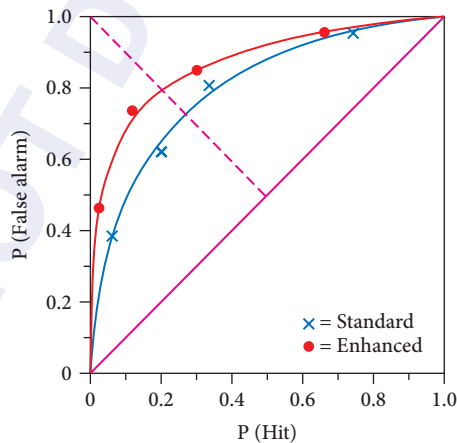


FIGURE 3 Example of an empirical ROC. Six radiologists attempted to distinguish between malignant and benign lesions in a set of 118 mammograms, 58 malignant and 60 benign, first when the mammograms were viewed in the usual manner (“standard”), and then—“enhanced”—when they were viewed with two aids, including a checklist of diagnostic features. The ratings were “very likely malignant,” “probably malignant,” “possibly malignant,” “probably benign,” and “very likely benign.” The areas under the curves are 0.81 and 0.87. (From Swets.⁴⁵)

area can descend below 0.5 when the observer reverses the categories of blank and nonblank. We'll see in a moment that a result equivalent to ROC area can usually be obtained with much less effort by doing a two-alternative forced choice experiment instead.

Two-Alternative Forced Choice

This task is traditionally characterized by two separate stimulus presentations, one blank and one nonblank, in random order. The two stimuli may be presented successively or side by side. The observer is asked whether the nonblank stimulus was first or second (or on the left or right). We noted above that in yes-no tasks observers seem to reduce the stimulus to a decision variable, the magnitude upon which they base their decisions. The 2afc task is said to be “unbiased” because the observer presumably chooses the presentation that generated the higher magnitude, without referring to any subjective internal criterion. At the beginning of this section we said that all judgment tasks consist of the presentation of a stimulus followed by a judgment. In this view, we might consider the two presentations in the 2afc task to be a single stimulus. The two possible composite stimuli to be discriminated are reflections of one another, either in space or time. The symmetry of the two alternatives suggests that the observer's choice between them may be unbiased.

Other related tasks are often called “two-alternative forced choice” and are similarly claimed to be unbiased. There is some confusion in the literature over which tasks should be called “2afc.” In our view, the “2afc” label is of little consequence. What matters is whether the task is unbiased, i.e., are the alternative stimuli symmetric for the observer? Thus a yes-no discrimination of blank and nonblank stimuli may be biased even though there are two response alternatives and the choice is forced, whereas it may be reasonable to say that the judgment of the orientation of a grating that is either horizontal or vertical is unbiased even though there is only a single presentation. We suggest that authors wishing to claim that their task is unbiased say so explicitly and state why. This claim might be based on a priori considerations of the symmetry between the stimuli to be discriminated, or on a post hoc analysis of relative frequencies of the observer's responses.

In theory, if we accept the assumptions that each stimulus presentation produces in the observer a unidimensional magnitude (one number, the decision variable), that the observer's ratings and 2afc decisions are based, in the proper way, on this magnitude, and that these magnitudes are stochastically independent between presentations, then the probability of a correct response on a 2afc trial must equal the area under the ROC curve.⁴⁷ Nachmias⁴³ compared 2afc proportion correct and ROC area empirically, finding that ROC area is slightly smaller, which might be explained by stimulus-induced variations in the observer's rating criteria.

MAGNITUDE ESTIMATION

In the inverse of magnitude production, a stimulus is presented and the observer is asked to rate it numerically.² Some practitioners provide a reference (e.g., a stimulus that rates 100), and some don't, allowing observers to use their own scale. Magnitude estimation and rating scale are fundamentally the same. Magnitude estimation experiments typically test many different stimulus intensities a few times to plot mean magnitude versus intensity, and rating-scale experiments typically test few intensities many times to plot an ROC curve at each intensity.

Response Time

In practical situations the time taken by the observer to produce a judgment usually matters, and it will be worthwhile recording it during the course of the experiment. Some psychophysical research has emphasized response time as a primary measure of performance in an effort to reveal mental processes.⁴⁹

3.6 STIMULUS SEQUENCING

So far we have discussed a single trial yielding a single response from the observer. Most judgments are stochastic, so judgment experiments usually require many trials. An uninterrupted sequence of trials is called a *run* (or a *block*). There are two useful methods of sequencing trials within a run.

Method of Constant Stimuli

Experimenters have to worry about small, hard-to-measure variations in the observer's sensitivity that might contaminate comparisons of data collected at different times. It is therefore desirable to run the trials for the various conditions as nearly simultaneously as possible. One technique is to interleave trials for the various conditions. This is the classic "method of constant stimuli." Unpredictability of the experimental condition and equal numbers of trials for each condition are typically both desirable. These are achieved by using a randomly shuffled list of all desired trials to determine the sequence.

Sequential Estimation Methods

One can use the method of constant stimuli to measure performance as a function of a signal parameter—let us arbitrarily call it intensity—and determine, by interpolation, the threshold intensity that corresponds to a criterion level of performance.* This approach requires hundreds of trials to produce a precise threshold estimate. Various methods have been devised that obtain precise threshold estimates in fewer trials, by using the observer's previous responses to choose the stimulus intensity for the current trial. The first methods were simple enough for the experimenter to implement manually, but as computers appeared and then became faster, the algorithms have become more and more sophisticated. Even so, the requisite computer programs are very short.

In general, there are three stages to threshold estimation. First, all methods, implicitly or explicitly, require that the experimenter provide a confidence interval around a guess as to where threshold may lie. (This bounds the search. Lacking prior knowledge, we would have an infinite range of possible intensities. Without a guess, where would we place the first trial? Without a confidence interval, where would we place the second trial?) Second, one must select a test intensity for each trial based on the experimenter's guess and the responses to previous trials. Third, one must use the collected responses to estimate threshold. At the moment, the best algorithm is called ZEST,⁵⁰ which is an improvement over the popular QUEST.⁵¹ The principal virtues of QUEST are that it formalizes the three distinct stages, and implements the first two stages efficiently. The principal improvement in ZEST is an optimally efficient third stage.

3.7 CONCLUSION

This chapter has reviewed the practical considerations that should guide the choice of psychophysical methods to quickly and definitely answer practical questions related to perception and performance. Theoretical issues, such as the nature of the observer's internal decision process, have been de-emphasized. The question of how well we see is answerable only after we reduce the question to measurable performance of a specific task. The task will be either an adjustment—for a quick answer when the perceptual criterion is unambiguous—or a judgment—typically to find threshold by sequential estimation.

*The best way to interpolate frequency-of-seeing data is to make a maximum likelihood fit by an S-shaped function.⁵² Almost any S-shaped function will do, provided it has adjustable position and slope.⁵³

The success of psychophysical measurements often depends on subtle details: the seemingly incidental properties of the visual display, whether the observers receive feedback about their responses, and the range of stimulus values encountered during a run. Decisions about these matters have to be taken on a case-by-case basis.

3.8 TIPS FROM THE PROS

We asked a number of colleagues for their favorite tips.

- Experiments often measure something quite different from what the experimenter intended. Talk to the observers. Be an observer yourself.
- Viewing distance is an often-neglected but powerful parameter, trivially easy to manipulate over a 100:1 range. Don't be limited by the length of your keyboard cable.
- Printed vision charts are readily available, offering objective measurement of visibility, e.g., to characterize the performance of a night-vision system.^{20,21,54,55}
- When generating images on a cathode ray tube, avoid generating very high video frequencies (e.g., alternating black and white pixels along a horizontal raster line) and very low video frequencies (hundreds of raster lines per cycle) since they are typically at the edges of the video amplifier's passband.⁵⁶
- Liquid crystal displays (LCD) have largely replaced cathode ray tube (CRT) displays in the market place. LCDs are fine for static images, but have complicated temporal properties that are hard to characterize. Thus, CRTs are still preferable for presentation of dynamic images, as they allow you to know exactly what you are getting.⁵⁷
- Consider the possibility of aftereffects, whereby past stimuli (e.g., at high contrast or different luminance) might affect the visibility of the current stimulus.^{58–61}
- Drift of sensitivity typically is greatest at the beginning of a run. Do a few warm-up trials at the beginning of each run. Give the observer a break between runs.
- Allow the observer to see the stimulus once in a while. Sequential estimation methods tend to make all trials just barely detectable, and the observer may forget what to look for. Consider throwing in a few high-contrast trials, or defining threshold at a high level of performance.
- Calibrate your display before doing the experiment, rather than afterward when it may be too late.

3.9 ACKNOWLEDGMENTS

Josh Solomon provided Fig. 1. Tips were contributed by Al Ahumada, Mary Hayhoe, Mary Kaiser, Gordon Legge, Walt Makous, Suzanne McKee, Eugenio Martinez-Uriega, Beau Watson, David Williams, and Hugh Wilson. Al Ahumada, Katey Burns, and Manoj Raghavan provided helpful comments on the manuscript. Supported by National Eye Institute grants EY04432 and EY06270.

3.10 REFERENCES

1. E. G. Boring, *Sensation and Perception in the History of Experimental Psychology*, Irvington Publishers, New York, 1942.
2. G. A. Gescheider, *Psychophysics: Methods, Theory, and Application*, 2d ed., Lawrence Erlbaum and Associates, Hillsdale, N.J., 1985, pp. 174–191.
3. N. A. Macmillan and C. D. Creelman, *New Developments in Detection Theory*, Cambridge University Press, Cambridge, U.K., 1991.

4. B. Farell, and D. G. Pelli, "Psychophysical Methods, or How to Measure a Threshold and Why," R. H. S. Carpenter and J. G. Robson (eds.), *Vision Research: A Practical Guide to Laboratory Methods*, Oxford University Press, New York, 1999.
5. P. Mertz, A. D. Fowler, and H. N. Christopher, "Quality Rating of Television Images," *Proc. IRE* **38**:1269–1283 (1950).
6. J. A. Swets, and R. M. Pickett, *Evaluation of Diagnostic Systems: Methods from Signal Detection Theory*, Academic Press, New York, 1982.
7. C. E. Metz, "ROC Methodology in Radiologic Imaging," *Invest. Radiol.* **21**:720–733 (1986).
8. F. Scott, "The Search for a Summary Measure of Image Quality—A Progress Report," *Photographic Sci. and Eng.* **12**:154–164 (1968).
9. F. G. Ashby, "Multidimensional Models of Categorization," *Multidimensional Models of Perception and Cognition*, F. G. Ashby (ed.), Lawrence Erlbaum Associates, Hillsdale, N.J., 1992.
10. D. G. Pelli, N. J. Majaj, N. Raizman, C. J. Christian, E. Kim, and M. C. Palomares, "Grouping in Object Recognition: The Role of a Gestalt Law in Letter Identification," *Cognitive Neuropsychology* (2009) In press.
11. F. W. Campbell, E. R. Howell, and J. G. Robson, "The Appearance of Gratings with and without the Fundamental Fourier Component," *J. Physiol.* **217**:17–18 (1971).
12. B. A. Wandell, "Color Measurement and Discrimination," *J. Opt. Soc. Am. A* **2**:62–71 (1985).
13. A. B. Watson, A. Ahumada, Jr., and J. E. Farrell, "The Window of Visibility: A Psychophysical Theory of Fidelity in Time-Sampled Visual Motion Displays," *NASA Technical Paper, 2211*, National Technical Information Service, Springfield, Va. 1983.
14. E. H. Adelson, and J. R. Bergen, "Spatiotemporal Energy Models for the Perception of Motion," *J. Opt. Soc. Am. A* **2**:284–299 (1985).
15. S. A. Klein, E. Casson, and T. Carney, "Vernier Acuity as Line and Dipole Detection," *Vision Res.* **30**:1703–1719 (1990).
16. B. Farell, and D. G. Pelli, "Psychophysical Methods," *A Practical Guide to Vision Research*, J. G. Robson and R. H. S. Carpenter (eds.), Oxford University Press, New York, 1999.
17. N. V. S. Graham, *Visual Pattern Analyzers*, Oxford University Press, Oxford, 1989.
18. H. B. Barlow, "Temporal and Spatial Summation in Human Vision at Different Background Intensities," *J. Physiol.* **141**:337–350 (1958).
19. F. W. Campbell, and J. G. Robson, "Application of Fourier Analysis to the Visibility of Gratings," *J. Physiol.* **197**:551–566 (1968).
20. H. Snellen, *Test-Types for the Determination of the Acuteness of Vision*, London: Norgate and Williams, 1866.
21. D. G. Pelli, J. G. Robson, and A. J. Wilkins, "The Design of a New Letter Chart for Measuring Contrast Sensitivity," *Clin. Vis. Sci.* **2**:187–199 (1988).
22. D. G. Pelli, and J. G. Robson, "Are Letters Better than Gratings?," *Clin. Vis. Sci.* **6**:409–411 (1991).
23. J. C. Dainty, and R. Shaw, *Image Science*, Academic Press, New York, 1974.
24. E. H. Linfoot, *Fourier Methods in Optical Image Evaluation*, Focal Press, New York, 1964.
25. D. E. Pearson, *Transmission and Display of Pictorial Information*, John Wiley & Sons, New York, 1975.
26. O. H. Schade, Sr., *Image Quality: A Comparison of Photographic and Television Systems*, RCA Laboratories, Princeton, N.J., 1975.
27. G. E. Legge, D. G. Pelli, G. S. Rubin, and M. M. Schleske, "Psychophysics of Reading—I. Normal Vision," *Vision Res.* **25**:239–252 (1985).
28. J. M. Rolf and K. J. Staples, *Flight Simulation*, Cambridge University Press, Cambridge, U.K., 1986.
29. D. G. Pelli, "The Visual Requirements of Mobility," *Low Vision: Principles and Application*, G. C. Woo (ed.), Springer-Verlag, New York, 1987, pp. 134–146.
30. J. G. Robson, "Spatial and Temporal Contrast-Sensitivity Functions of the Visual System," *J. Acoust. Soc. Am.* **56**:1141–1142 (1966).
31. R. S. Woodworth and H. Schlosberg, *Experimental Psychology*, Holt, Rinehart, and Winston, New York, 1963, pp. 199–200.
32. P. Cavanagh and S. Anstis, "The Contribution of Color to Motion in Normal and Color-Deficient Observers," *Vision Res.* **31**:2109–2148 (1991).

33. G. A. Brindley, *Physiology of the Retina and the Visual Pathways*, Edward Arnold Ltd., London, 1960.
34. M. A. Georgeson and G. D. Sullivan, "Contrast Constancy: Deblurring in Human Vision by Spatial Frequency Channels," *J. Physiol.* **252**:627–656 (1975).
35. R. M. Boynton, *Human Color Vision*, Holt Rinehart and Winston, New York, 1979, pp. 299–301.
36. W. W. Peterson, T. G. Birdsall, and W. C. Fox, "Theory of Signal Detectability," *Trans. IRE PGIT* **4**:171–212 (1954).
37. W. P. Tanner, Jr. and T. G. Birdsall, "Definitions of d' and η as Psychophysical Measures," *J. Acoust. Soc. Am.* **30**:922–928 (1958).
38. H. L. Van Trees, *Detection, Estimation, and Modulation Theory*, Wiley, New York, 1968.
39. W. S. Geisler, "Sequential Ideal-Observer Analysis of Visual Discriminations," *Psychol. Rev.* **96**:267–314 (1989).
40. D. G. Pelli, "Uncertainty Explains Many Aspects of Visual Contrast Detection and Discrimination," *J. Opt. Soc. Am. A* **2**:1508–1532 (1985).
41. D. G. Pelli, "The Quantum Efficiency of Vision," *Vision: Coding and Efficiency*, C. Blakemore (ed.), Cambridge University Press, Cambridge, U.K., 1990, pp. 3–24.
42. D. G. Pelli, C. W. Burns, B. Farell, and D. C. Moore-Page, "Feature Detection and Letter Identification," *Vision Res.* **46**(28):4646–4674 (2006). See Appendix A.
43. J. Nachmias, "On the Psychometric Function for Contrast Detection," *Vision Res.* **21**:215–223 (1981).
44. H. E. Rockette, D. Gur and C. E. Metz, "The Use of Continuous and Discrete Confidence Judgments in Receiver Operating Characteristic Studies of Diagnostic-Imaging Techniques," *Invest. Radiol.* **27**:169–172 (1992).
45. J. A. Swets, "Measuring the Accuracy of Diagnostic Systems," *Science* **240**:1285–1293 (1988).
46. J. Nachmias and R. M. Steinman, "Brightness and Discriminability of Light Flashes," *Vision Res.* **5**:545–557 (1965).
47. D. M. Green and J. A. Swets, *Signal Detection Theory and Psychophysics*, Krieger Press, Huntington, N.Y., 1974.
48. L. W. Nolte and D. Jaarsma, "More on the Detection of One of M Orthogonal Signals," *J. Acoust. Soc. Am.* **41**:497–505 (1967).
49. R. D. Luce, *Response Times: Their Role in Inferring Elementary Mental Organization*, Oxford University Press, New York, 1986.
50. P. E. King-Smith, S. S. Grigsby, A. J. Vingrys, S. C. Benes and A. Supowit, "Efficient and Unbiased Modifications of the QUEST Threshold Method: Theory, Simulations, Experimental Evaluation and Practical Implementation," *Vision Res.* **34**:885–912 (1994).
51. A. B. Watson and D. G. Pelli, "QUEST: A Bayesian Adaptive Psychometric Method," *Percept Psychophys.* **33**:113–120 (1983).
52. A. B. Watson, "Probability Summation Over Time," *Vision Res.* **19**:515–522 (1979).
53. D. G. Pelli, "On the Relation Between Summation and Facilitation," *Vision Res.* **27**:119–123 (1987).
54. S. Ishihara, *Tests for Color Blindness*, 11th ed., Kanehara Shuppan, Tokyo, 1954.
55. D. Regan and D. Neima, "Low-Contrast Letter Charts as a Test of Visual Function," *Ophthalmology* **90**:1192–1200 (1983).
56. D. G. Pelli and L. Zhang, "Accurate Control of Contrast on Microcomputer Displays," *Vision Res.* **31**:1337–1350 (1991).
57. D. H. Brainard, D. G. Pelli and T. Robson, "Display Characterization," In J. Hornak (ed.), *Encyclopedia of Imaging Science and Technology*, Wiley, 2002, pp. 172–188.
58. C. Blakemore and F. W. Campbell, "Adaptation to Spatial Stimuli," *J. Physiol.* **200**(1):11–13 (1969).
59. C. Blakemore and F. W. Campbell, "On the Existence of Neurones in the Human Visual System Selectively Sensitive to the Orientation and Size of Retinal Images," *J. Physiol.* **203**:237–260 (1969).
60. T. N. Cornsweet, *Visual Perception*, Academic Press, New York, 1970.
61. F. S. Frome, D. I. A. MacLeod, S. L. Buck and D. R. Williams, "Large Loss of Visual Sensitivity to Flashed Peripheral Targets," *Vision Res.* **21**:1323–1328 (1981).

4

VISUAL ACUITY AND HYPERACUITY

Gerald Westheimer

*Division of Neurobiology
University of California
Berkeley, California*

4.1 GLOSSARY

Airy disk. Point-spread function in the image of a diffraction-limited optical instrument with a circular pupil.

Diffraction limit. Minimum dissipation of spatial information in the imaging of an optical system, due to the aperture restriction in the propagation of electromagnetic energy.

Fovea. Region in the center of the retina where receptor elements are most closely packed and resolution highest.

Hyperacuity. Performance in task where thresholds are substantially lower than the grain of the receiving layer.

Light. Visually evaluated radiant energy. In this chapter radiant and luminous energy terms are used interchangeably.

Optical-transfer function. Modulation in the transmitted images of spatial sinusoids, as a function of their spatial frequency; it is complex, that is, has amplitude and phase terms.

Psychophysics. Procedure for studying an observer's performance by relating the variables of physical stimuli to measurements of associated responses.

Point-spread function. Spatial distribution of energy in the image of a point object.

Snellen letters. Alphanumeric characters of defined size and shape used in standard clinical testing of visual acuity.

Spatial frequency. Number of cycles of a sinusoidal grating target per unit distance. Commonly cycles/degree visual angle.

Superresolution. Ability to garner knowledge of spatial details in an optical image based on previous available information, either by extrapolation or averaging.

Vernier acuity. Performance limit in the alignment of two abutting line segments; it is the prime example of hyperacuity.

Visual acuity. Performance limit in distinguishing spatial details in visual object.

Visual angle. Angle subtended by an object at the center of the eye's entrance pupil; it is a measure of distance in the retinal image.

Equation (1). Point-spread function of a purely diffraction-limited optical imaging system with a round pupil.

θ angular subtense of radius of Airy's disk

λ wavelength of radiation

a diameter of aperture

Equation (2). Specification of contrast in the spatial distribution.

L_{\max} , L_{\min} Luminance of maximum, minimum, respectively

4.2 INTRODUCTION

Visual acuity—literally sharpness—refers to the limit of the ability to discriminate spatial partitioning in the eye's object space. As a psychophysical measure, its analysis encompasses

- The physics—in this case optics—of the stimulus situation
- The anatomical and physiological apparatus within the organism that processes the external stimulus
- The operations leading to the generation of a response

Measurement of visual acuity involves the organism as a whole, even though it is possible to identify the performance limits of only a segment of the full operation, for example, the eyeball as purely an optical instrument, or the grain of the receptor layer of the retina, or neural activity in the visual cortex. But the term acuity is reserved for the behavioral essay and therefore necessarily includes the function of all components of the arc reaching from physical object space to some indicator of the response of the whole organism. It is a psychophysical operation. The fact that it includes a component that usually involves an observer's awareness does not preclude it from being studied with any desired degree of rigor.

4.3 STIMULUS SPECIFICATIONS

Specification of the stimulus is an indispensable preliminary. For this purpose, an Euclidean object space containing visual targets is best defined by a coordinate system with its origin at the center of the eye's entrance pupil and its three axes coinciding with those of the eye. The observer is usually placed to make the vertical (y) axis that of gravity, and the horizontal (x) axis orthogonal and passing through equivalent points in the two eyes; distances from the observer are measured along the z axis. When an optical device is associated with the eye, its principal axes should accord, either by positioning the device or the observer. On occasions when both eyes of an observer are involved, the origin of the coordinate system is located at the midpoint of the line joining the two eyes. More details of possible coordinate systems have been described elsewhere.¹

The ocular structure most relevant to the spatial dissection of an observer's object world is the retina, and the inquiry begins with the quest for the most instructive way of relating the *distal* and *proximal* stimuli, that is, the actual objects and their retinal images. Here it is achieved by using the center of the entrance pupil of the eye as the origin of the coordinate system which has at least two advantages. First, while associated with the eye's imagery, it belongs to object space and hence is objectively determinable by noninvasive means. In a typical human eye, the entrance pupil is located about 3 mm behind the corneal vertex and is not far from round. Its center is therefore an operationally definable point. The second reason for choosing it as a reference point involves the nature of the eye's image-forming apparatus. As seen in Fig. 1, the bundle of rays from a point source converging toward the retina is centered on the ray emerging from the center of the eye's exit pupil, which is the optical conjugate of the center of the entrance pupil. Regardless of the position

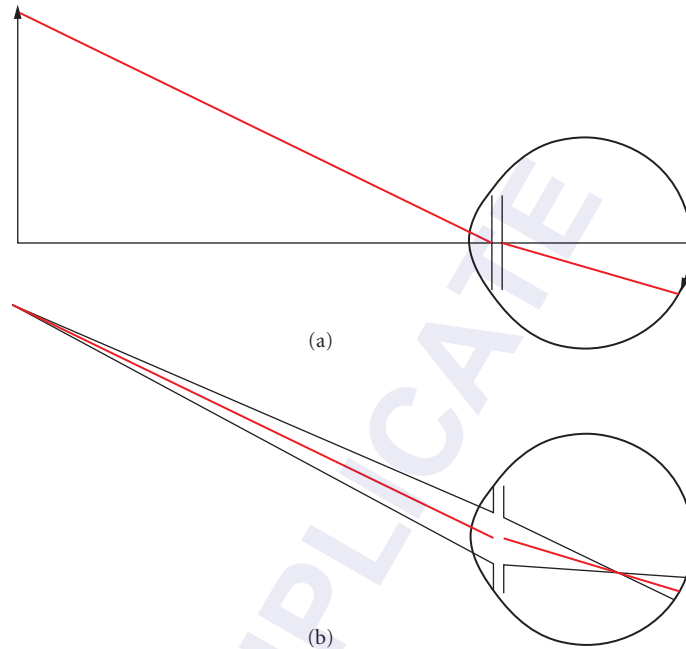


FIGURE 1 Schematic diagram of imaging in human eye. (a) The retinal image size of a target is represented by the angle it subtends at the eye's entrance pupil. (b) The position of the image of a point is most effectively demarcated by the intercept of the image-sided chief ray, which is the center of the light distribution even if the eye is out of focus.

of the geometrical image with respect to the retina, that is, regardless of the state of defocus, the center of the retinal image patch, blurred or sharp, will be defined by the intersection of that ray, called *chief ray*, with the retina. When two object points are presented to the eye, the retinal distance corresponding to their separation is given by the intercept of the chief rays from these objects. In this manner, the three-dimensional object space of the eye has been collapsed into the two-dimensional one of the retinal surface. What has been lost, and needs to be specified separately, is the object's distance from the eye along the chief ray. But all objects, in or out of focus, anywhere along a given chief ray share a single retinal location, or, to rephrase it, the coordinates on the retinal surface are homologous to angular coordinates within the object-sided sheaf of rays converging on the center of the eye's entrance pupil.

Hence in the specification of retinal distances it suffices to identify corresponding angles in the eye's object space, objectively determinable measures. Units of measurement are the radian, or, degrees or minutes of arc. At a distance of 57 cm, a 1-cm object subtends 1 deg, a 0.16-mm object 1 arcmin. At the standard eye-chart distance of 6 m (20 ft) the limb of a 20/20 letter is just under 2 mm wide.

The next specification to be considered is that of the luminous intensity impinging on the eye. One starts with the most elemental stimulus: the luminous intensity of a point source is given in the internationally agreed-on unit of candela ($\text{lumens} \cdot \text{steradian}^{-1}$). This again is object-sided and objectively determinable. Extended sources are measured in terms of luminance ($\text{lumens} \cdot \text{steradian}^{-1} \cdot \text{unit area}^{-1}$, in practice $\text{cd} \cdot \text{m}^{-2}$). Hence the specification of visual acuity targets in the eye's object space requires, apart from their observation distance, their spatial extent in angular measure at the eye's entrance pupil and the luminance of the background from which or against which they are formed. The luminous energy reaching the retina differs in that it depends on the pupil area, the absorption in the eye's

media, and retinal factors considered elsewhere. Attenuation in the passage through the ocular media will differ from one eye to another and is prominently age dependent (Chap. 1). (Polarization and coherence properties of the incoming luminous energy are generally not relevant, but see special conditions analyzed in Chap. 14). How this luminous energy is distributed on the retinal surface depends on the transfer characteristics of the eye's optics which will now be considered.

4.4 OPTICS OF THE EYE'S RESOLVING CAPACITY

The spatial distribution of energy reaching the retina will differ from that incident on the cornea by being subject to spread produced by the eye's imaging.

Light Spread in the Retinal Image and Limits of Resolution

The two modes of proceeding in this discussion, via the point-spread or the contrast-transfer functions, are equivalent (Fig. 2). So long as one remains in the realm of optics and does not enter that of neural and psychophysical processing where linearity is not guaranteed, it is permissible to transfer back and forth between the two.

Point-Spread Function When the object space is restricted to a single point, the spatial distribution of energy in the image is called the point-spread function and describes the spread introduced by passage through the eye's optics. Even in ideal focus, the spread depends on the eye's aperture and the wavelength of the electromagnetic energy; the object-sided distribution then has θ , that is, angle subtended by the distance from its center to the first zero, given by

$$\theta = \frac{1.22\lambda}{a} \quad (1)$$

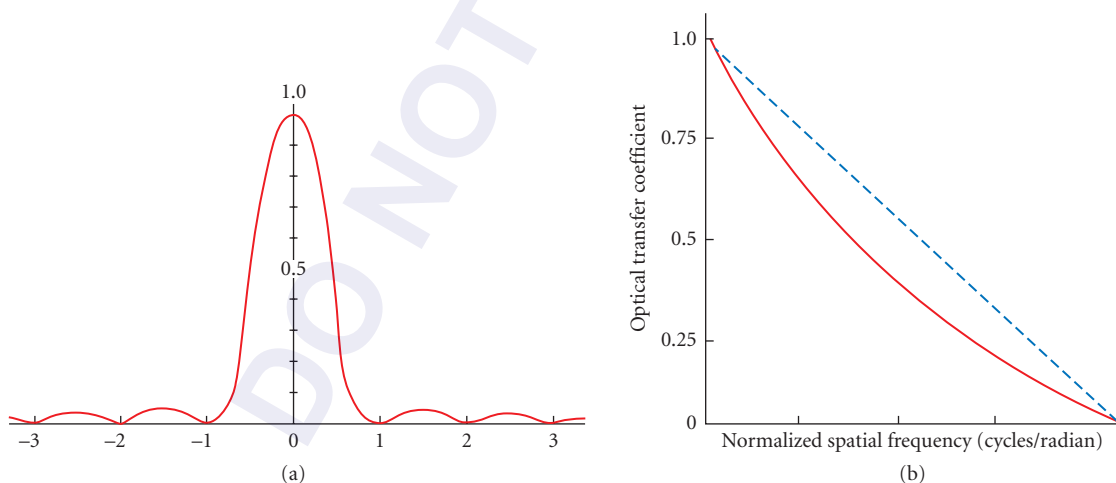


FIGURE 2 Retinal light distribution in an idealized optical system like the eye's. (a) Diffraction-limited point-spread function (Airy's disk). Unit distance is given by $1.22 \lambda/a$ in radians, where λ is the wavelength of light and a the diameter of the round pupil, both in the same units of length. (b) The optical contrast-transfer function for a square pupil (blue dashed line) and a round pupil (red solid line). It descends to zero at a spatial frequency equal to a/λ cycles/radian.

A point can never be imaged smaller than a patch of such a size; the light distribution for any object is the convolution of that of the points constituting it. In practice there are additional factors due to the aberrations and scattering in each particular eye.

Spatial-Frequency Coordinates and Contrast-Transfer Function A fundamental property of optical imagery in its application to the eye is its linearity. Hence a permissible description of light distributions is in terms of their spatial Fourier spectra, that is, the amplitudes and phases of those spatial sinusoidal intensity distribution in terms of their spatial frequency that, when superimposed, will exactly reconstruct the original distribution. The limit here is the cutoff spatial frequency at which the optical transfer coefficient of the eye reaches zero (Fig. 2*b*).

Either of the two descriptors of the optical transfer between the eye's object and image spaces, the point-spread function, that is, light spread in the image of a point, and the contrast-transfer function, that is, the change in amplitude and phase that the component sinusoids (as a function of spatial frequency in two angular dimensions) experience as they are transferred from the eye's object space to the retinal image, is complete and the transposition between the two descriptors is uncomplicated.

Because resolution relates to the finest detail that can be captured in an image, the interest is in the narrowness of the point-spread function (e.g., width at half-height) or, equivalently, the high-frequency end of the contrast-transfer function. The absolute limit imposed by diffraction gives a bound to these two functions but this is actually achieved only in fully corrected eyes with pupil diameters below about 3 mm, when other dioptric deficits are minimal. Concentrating on light of wavelength 555 nm, for which the visual system during daylight is most sensitive, Fig. 2, taken over directly from diffraction theory, illustrates the best possible performance that may be expected from a normal human eye under ordinary circumstances. The cutoff spatial frequency then is near 90 cycles/degree⁻¹ and the diameter of Airy's disk about 1.5 arcmin.

A great deal of effort has gone into the determination of the actual point-spread functions of eyes which include factors such as aberrations. Only under very exceptional circumstances—and nowadays also with the aid of adaptive optics—does one get better imagery than what is shown in Fig. 2. When the pupil diameter is increased, theoretical improvement due to narrowing of the Airy disk is counteracted by aberrations which become more prominent as the outer zones of the pupil are uncovered. The effect of refractive errors on imaging can also be described in theory, but then phase changes in the contrast-transfer function enter because its complex nature (in the mathematical sense) can no longer be ignored (see under "Defocus").

4.5 RETINAL LIMITATIONS—RECEPTOR MOSAIC AND TILING OF NEURONAL RECEPTIVE FIELDS

Also amenable to physical analysis are the limitations imposed on the resolving capacity of the eye by the structure of the retina. All information that is handed on to the neural stages of vision is in the first place partitioned by the elements of the receptor layer, each of which has an indivisible spatial signature. The spacing of the receptors is not uniform across the retina, nor is the individuality of their local sign necessarily retained in further processing. Ultimately, the information transfer to the brain is confined by the number of individual nerve fibers emerging from the retina. Nevertheless it is instructive to inquire into the grain of the retina in the foveal region where there is certainly at least one optic nerve fiber for each receptor.

Figure 3 is a cross section of the layer of a primate retina and shows an approximately hexagonal array of cones, whose average spacing in the center is of the order of 0.6 arcmin. No matter what else is in play, human visual resolution cannot be better than is allowed by this structure. The center of the fovea is only about 0.5 deg in diameter; the further one proceeds into the retinal periphery the coarser the mosaic and the lower the ratio of receptors to optic nerve fibers. This explains the reason for our highly developed oculomotor system with its quick ability to transfer foveal gaze to different eccentric and even moving targets.

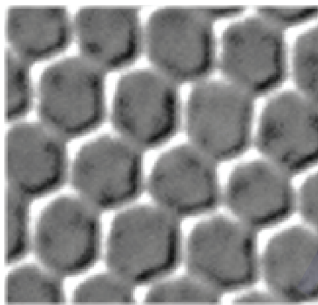


FIGURE 3 Histological cross section of the retinal mosaic in the primate fovea. Each receptor represents an object-sided angle of about 0.6 arcmin.

The neural elements of the retina are not passive transducers but actively rearrange the optical signals that reach the receptors. After transmission to the brain, processing of these neural signals involves interaction from other regions and modification as a result of such factors as attention and memory. As yet the neural circuitry interposed between the optical image on the retina and the individual's acuity response has not reached a level of understanding equivalent to that of the optical and receptor stages.

4.6 DETERMINATION OF VISUAL RESOLUTION THRESHOLDS

Awareness of the limitations imposed by the optics and anatomy of the eye is, of course, of value, but visual acuity is, in the end, a function of the operation of the whole organism: what are the finest spatial differences that can be distinguished? In answering this question, attention has to be paid to the manner of obtaining the measurements.

Since there is scatter in individual determinations, the number of trials will be controlled by the needed precision. The armamentarium of psychophysical procedures allows determination of a threshold with arbitrary precision. The standard optometric visual acuity chart, in use for 150 years, is a textbook case of effective employment of this approach. The ensemble of test symbols, the information content per symbol, the accepted answers, and the scaling of the steps and number of trials (letters in each row) have all been optimized for quick and reliable acuity identification.

A psychophysical threshold is a number along a scale of a variable (e.g., distance between double stars) at which a correct response is made in a predetermined proportion of trials. In a good experiment it is accompanied by a standard error. Thus if in a particular situation the two-star resolution threshold is found to be 0.92 ± 0.09 arcmin and the 50-percent criterion was employed, this means that on 50 percent of occasions when the separation was 0.92 inch the observer would say "yes" (the percentage increasing with increasing separation) and that the scatter of the data and the number of observations were such that if the whole experiment were repeated many times, the value would be expected to be between 0.83 inch and 1.01 inch in 19 out of 20 runs of data.

The distinction is often made between detection and discrimination thresholds. In sophisticated detection procedures, stimulus presentations are alternated randomly with blanks. A count is kept of the number of times the observer gives a "yes" answer when there was no stimulus, the so-called false positives. An elaborate analytical procedure can then be deployed to examine the internal

“noise” against which the incoming sensory signal has to compete.² This methodology is appropriate when there is a blank as one of the alternatives in the test, for example, in the measurement of the high spatial-frequency cutoff for grating resolution. In the bulk of acuity determinations the observer has to discriminate between at least two alternative configurations each presented well above detection threshold, and what has to be safeguarded against are bias errors.

It is the current practice in vision research to observe as many of these niceties of psychophysical methodology as possible (Chap. 3). However, in clinical and screening situations less time and observer engagement is available, with, as a consequence, diminished reliability and repeatability of findings.

4.7 KINDS OF VISUAL ACUITY TESTS

The common denominator of acuity tests is the determination of the finest detectable spatial partitioning, and this can be done in many ways. The closest to the optical concept of resolving power are the two-point or two-line patterns, whose minimum separation is measured at which they are seen as double (Fig. 4*a*). More popular are the two-bar experiments, most effectively implemented in a matrix of 3×3 elements, in which either the top and bottom rows, or the right and left columns

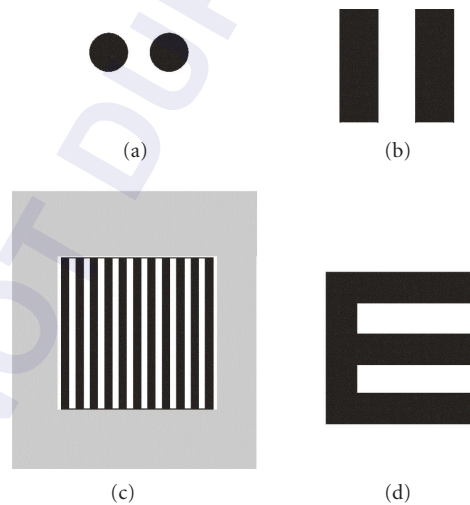


FIGURE 4 Patterns used in visual acuity tests and the associated response criteria: (a) Two point resolution. (b) Koenig bars in a 3×3 matrix. (c) Grating resolution. (d) Letters, as used in the clinical Snellen chart. Observers respond to these questions: You will be shown either a single or a double star. Was it “one” or “two?” (a). You will be shown a two-line pattern. Were the lines vertical or horizontal? (b). You will be shown a field that is either blank or contains vertical stripes? Was it “blank” or “striped?” (c). You will be shown an alphanumeric character. What letter or digit was it? (d).

have a different contrast than the middle row or column (Fig. 4*b*), the observer's responses being limited to "horizontal" and "vertical." The size of the matrix elements is increased to determine the observer's threshold. Overall size is no clue and the response is based on the detection of the internal image structure.

Some brightness variables are available to the experimenter. In the old days, the lines would be black on a white background whose luminance would be specified. With the advent of oscilloscopic displays this can now be white on black or, more generally, brighter and darker than a uniform background. Contrast is then a further variable, usually defined by the Michelson formula

$$\frac{(L_{\max} - L_{\min})}{(L_{\max} + L_{\min})} \quad (2)$$

With the advent of the Fourier approach to optics, grating targets have become popular. For the purposes of acuity, the highest spatial-frequency grating is determined at which, with 100-percent modulation, the field is just seen as striped rather than uniform (Fig. 4*c*). The phenomenon of spurious resolution, described below, makes this test inadvisable when focus errors may be at play. The role of grating targets in acuity measurements differs from that in modulation sensitivity tests, where gratings with a range of spatial periods are demodulated till they are no longer detectable. This process yields the modulation sensitivity curve, discussed below.

Since they were first introduced in the second half of the 19th century, the standard for clinical visual acuity are the Snellen letters—alphanumeric characters, each drawn within a 5×5 matrix, with limb thickness as the parameter (Fig. 4*d*). From the beginning it was accepted that the resolution limit of the human eye is 1 arcmin, and hence the overall size of the Snellen letter for normal acuity is 5 arcmin, or 9.5 mm at a distance of 6 m or 20 ft (optical infinity for practical purposes). When such letters can be read at 20 ft, visual acuity is said to be 20/20. Letters twice this size can normally be read at 40 ft; an observer who can only read such double-sized letters at 20 ft has 20/40 acuity. The charts are usually assembled in lines of about 8 letters for 20/20 and progressively fewer for the lower ratings, with just a single letter for 20/200. Because in acuity determinations the error is proportional to the size of the letters³ and the sequence of letter sizes in charts is usually logarithmic,⁴ Snellen acuity is often converted to a fraction, 20/20 becoming 1.0, 20/40 becoming 0.5, and so on. When some of the letters in a line are missed, say 2 in the 20/20 line, a score of $20/20 - 2$ is recorded. For example, if there are 7 letters in the 20/25 (0.8) line of which 3 are missed, and the next line is 20/30 (0.67) a numerical value of $0.74 [0.8 - (3/7) (0.8 - 0.67)]$ can be entered for statistical purposes.

Snellen charts have been made available in many alphabets, but letters may not have equal legibility, even in English. Hence a stripped-down version of letter acuity is often used. A single letter E can be shown in four or even eight orientations and the observer asked to respond, perhaps by pointing a hand with outstretched fingers. The detection of the location of a gap in an annulus, that is, the distinction between the letter O and an oriented C, called the Landolt C test after its inventor is particularly useful. Both the E and Landolt C tests can be fitted into the tradition of Snellen letters by, for the 20/20 targets, generating them with 1' line within a $5 \times 5'$ matrix, with progressive size increases to arrive at a solid numerical acuity value. These two tests make no demands on the subjects' literacy and have the virtue that the effect of guessing is a known quantity. Even here, a minor problem arises because of the "oblique effect," a small performance deficit in oblique orientations over the horizontal and vertical.

The development of the visual system in infants and the early detection of visual anomalies has sparked the design of infant visual acuity tests, usually depending on the observation of eye movements to targets of interest whose size can be progressively diminished.⁵ Apart from this "preferential looking" technique, optokinetic nystagmus, that is, the involuntary eye tracking of large moving fields, can be effectively utilized to measure acuity by progressively diminishing the size of details until tracking fails.

As outlined so far, all the tests have in common the need for the subject to be aware and cooperative, though not necessary literate. When these conditions are absent, other procedures have to be adopted, for example, recording the signals from the eye into the central nervous system from the scalp. These are not further described here.

4.8 FACTORS AFFECTING VISUAL ACUITY

Visual acuity performance will be diminished whenever any of the contributing functions have not been optimized. Initial analysis concentrates on optical and retinal factors; not enough is known about the subsequent central neural stages to differentiate all the possible ways in which their operation can be encumbered or rendered inefficient. Such factors as attention, training, and task familiarity are clearly relevant. A treatment of the subject from the clinical point of view is available in textbooks.^{6,7}

Pupil When the optical line-spread function has been widened for whatever reason, resolution will obviously suffer. This is the case when the pupil is too small—2 mm or less—when diffraction widens it or, in the presence of aberrated wavefronts, when it is very large—usually 6 mm or larger.

Defocus Focus errors are of particular interest, because one of the most ubiquitous applications of visual acuity testing is to ascertain the best refractive state of eyes for purposes of spectacle, contact lens, or surgical correction. Ever since the establishment of the current optometric routines in the late 19th century, rules of thumb have existed for the relationship between refractive error and unaided acuity. One of these is shown in Fig. 5. But this does not take into account a patient's pupil size which governs depth of focus, the possible presence of astigmatism, and higher-order aberrations, nor some complications arising from the nature of out-of-focus imagery. When a spherical wavefront entering the eye's image space does not have its center at the retina, the imagery can be described as having a phase error that increases as the square of the distance from the center of the aperture.⁹ The contrast-transfer function, which is the Fourier transform of the complex (i.e., amplitude and phase) pupil aperture function, then does not descend to the cutoff spatial frequency monotonically, but shows oscillatory behavior (Fig. 6); in some regions in the spatial-frequency spectrum it dips below zero and the grating images have their black and white stripes reversed. This so-called spurious resolution means that if one views grating pattern under these conditions and gradually increases spatial frequency, stripes will first be visible, then disappear at the zero-crossing of the transfer function, then reappear with inverted contrast and so on. The seen image of objects like Snellen letters will undergo even more complex changes, with the possibility that “spurious” recognition is achieved with specific states of pupil size and defocus.

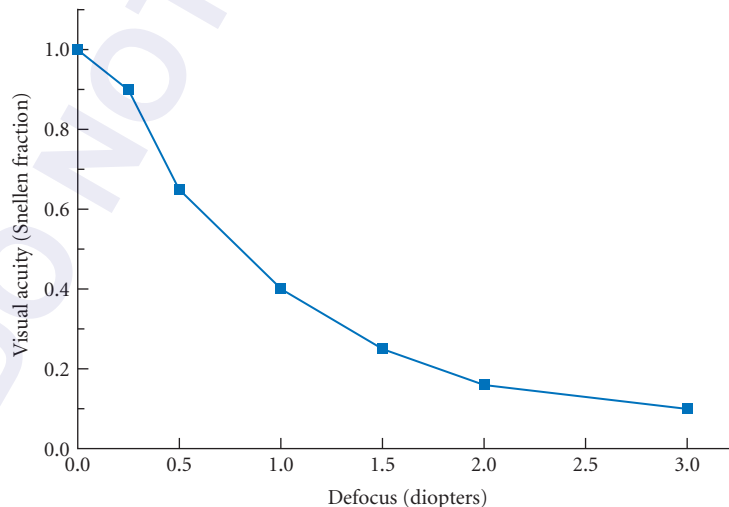


FIGURE 5 Visual acuity in a typical eye as a function of uncorrected spherical refractive error in diopters. (Adapted from Laurance.⁸)

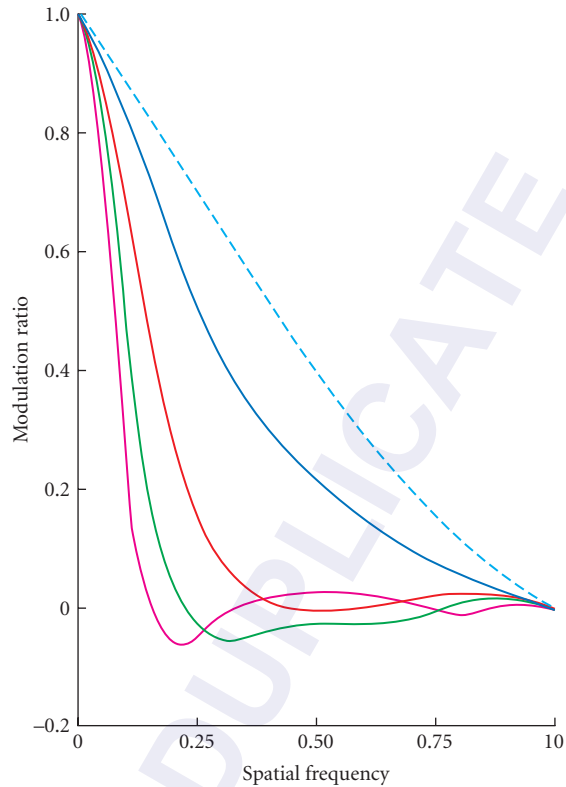


FIGURE 6 Normalized optical-transfer function for various degrees of defocus, showing regions of the spatial-frequency spectrum in which the coefficients are negative and the contrast of grating targets is reversed in the image compared to that in the object. (*Adapted from Hopkins.*¹⁰) For an eye with a 3-mm round pupil and wavelength 560 nm, the cutoff spatial frequency denoted by the normalized value 1.0 on the axis of abscissas is 1.53 cycles/arcmin and the five curves show the theoretical response for 0, 0.23, 0.31, 0.62, and 1 diopters defocus.

Color The diffraction equations have wavelength as an explicit variable. The width of the point-spread function varies inversely with wavelength. This is, however, only a minor factor where the effect of color on visual acuity is concerned. More immediately involved is the eye's chromatic aberration, giving defocus of about 1 diopter at the extremes of the visual spectrum, where also more energy is needed to compensate for the fact that the luminous efficiency of the eye peaks as 555 nm for the photopic and 500 nm for the scotopic (rod) system. In practice, therefore, each situation will have to be handled individually depending on the wavelength distribution of the particular stimulus.

Retinal Eccentricity Due to the coarsening of the grain of the anatomical connections in the peripheral retina, visual acuity falls off with increasing eccentricity (Fig. 7) and this is more pronounced in the extreme nasal than temporal field of view of each eye. In binocular vision, when

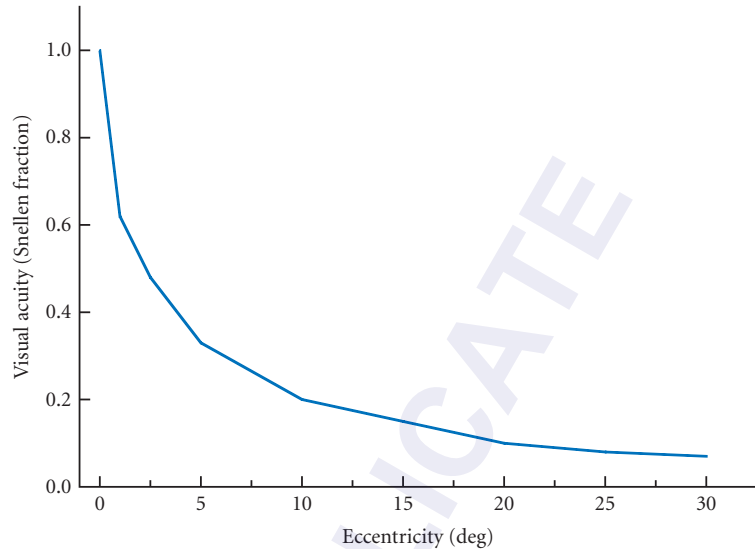


FIGURE 7 Expected visual acuity in various locations in the peripheral visual field. (Adapted from Wertheim.¹¹)

both eyes are in play, acuity is usually better than in monocular vision, not only because the better eye covers any possible deficiency of the other, but also because of probability summation due to the two retinas acting as independent detectors.

Luminance Acuity measured for black symbols against a white background, that is, for 100 percent contrast (Fig. 8) remains constant from about $10 \text{ cd} \cdot \text{m}^{-2}$ up. Below about $1 \text{ cd} \cdot \text{m}^{-2}$ the photopic system drops out and rods, which are absent in the fovea, take over. Their luminosity curve peaks at

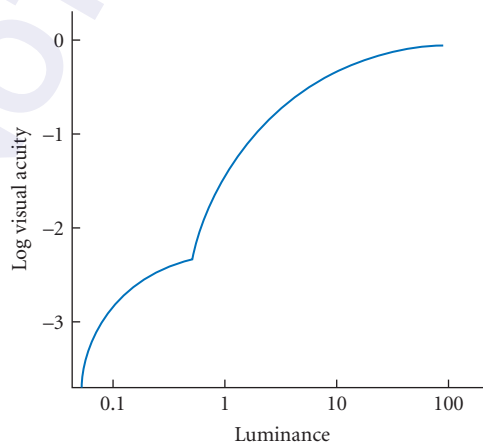


FIGURE 8 Visual acuity for targets as a function of their luminance. (Adapted from Shlaer.¹²) The rod-cone break occurs at a light level equivalent to that of a scene lit by full-moon light.

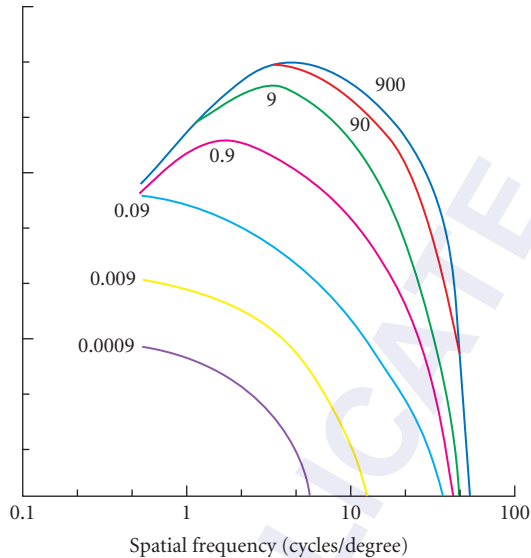


FIGURE 9 Modulation sensitivity (contrast detection) curve of the human visual apparatus as a function of spatial frequency at different light levels. (From van Nes and Bouman.¹³) Visual acuity is equivalent to the highest spatial frequency at which a response can still be obtained (intersection with the x axis) and the data here map well on those in Fig. 8.

about 500 nm. Also, they are color blind, and subject to considerable spatial summation and adaptation, that is, become more sensitive with increased time in the dark, up to as much as 30 to 45 min.

Contrast Most clearly shown with the use of grating stimuli, there is reduction in performance at both the low and high spatial frequency ends of the spectrum (Fig. 9). Because contrast sensitivity is lowered in various ocular abnormalities, particularly scatter or absorption in the media, it has been found valuable to perform acuity measurements with low-contrast charts.^{14,15} Reversing the contrast polarity, that is, presenting bright letters against a dark background, has virtue for eyes with a wide point-spread function and light scatter¹⁶ and indeed improves acuity performance in some older eyes.¹⁷

Time Time is decidedly a factor in visual acuity. For very short presentations, about 20 ms or less, the eye integrates all the light flux, what matters then is the product of the intensity and the duration, in any combination. But acuity improves with duration in the several hundred millisecond range where light detection no longer depends on duration but only intensity (Fig. 10).

Surround A prominent effect in visual acuity is that of crowding, where the presence of any contour close to the resolution target interferes with performance (Fig. 11).

Practice Effects Surprisingly, in the fovea of normal observers, practice in the task does not confer any further advantage—it seems that optimum performance has been entrained through continuous exercise of the facility in everyday situation. Peripheral acuity can, however, be improved by training.^{20,21}

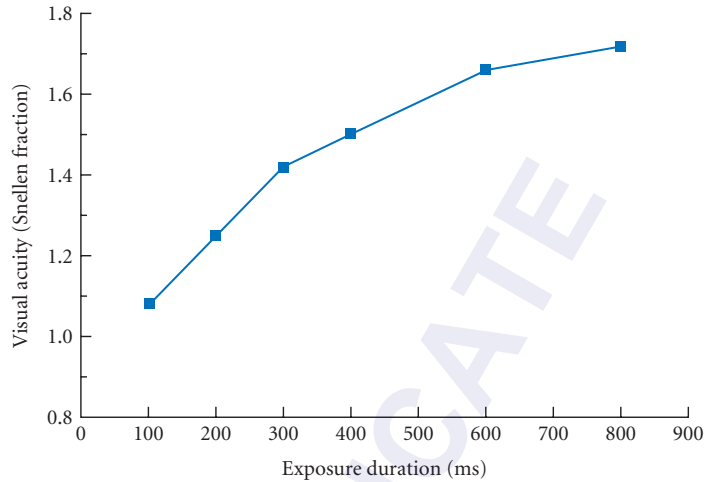


FIGURE 10 Visual acuity as function of exposure duration. (From Baron and Westheimer.¹⁸)

Stage of Development and Aging Visual acuity shows a steep increase in the first few months of life and, if a secure measurement can be obtained, is not far from normal at least by the third year.⁵ The aging eye is subject to a large number of conditions that impair acuity (Chap. 14); it is their presence or absence that determines any individual patient's status. Consequently age decrements of acuity can range from none to severe.

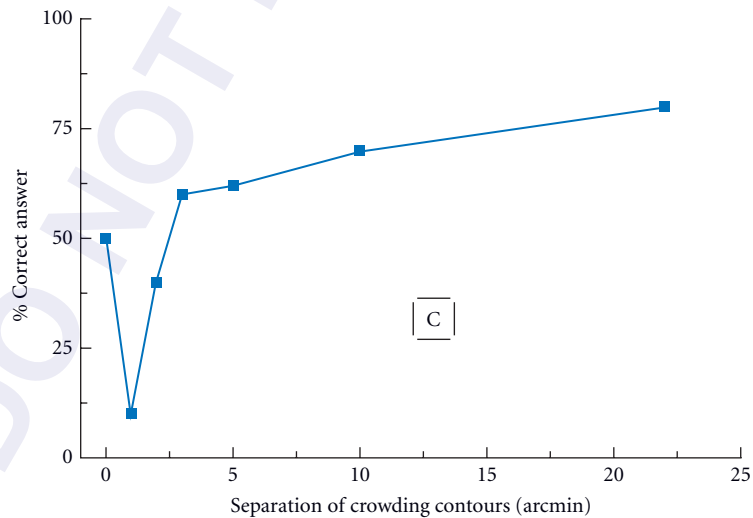


FIGURE 11 Data showing the “crowding” effect in visual acuity. A standard letter is surrounded by bars on all four sides and the performance drops when the bars are separated from the edge of the letter by the thickness of the letter's line-width. (Adapted from Flom, Weymouth, and Kahneman.¹⁹)

4.9 HYPERACUITY

For a long time human spatial discriminations have been known where thresholds are markedly lower than the resolution limit. For vernier acuity, the foveal alignment threshold for two abutting lines is just a few arcsecs, as compared with the 1 arcmin or so of ordinary resolution limit.

Such high precision of localization is shared by many kinds of pattern elements, both in the direction joining them, in their alignment, and in deviations from rectilinearity (Fig. 12). The word hyperacuity is applied to this discrimination of relative position, in recognition that it surpasses by at least an order of magnitude the traditional acuity. Whereas the limitation of the latter are mainly in the resolving capacity of the eye's optics and retinal mosaic, hyperacuity depends on the neural visual system's ability to extract subtle differences within the spatial patterns of the optical image on the retina.

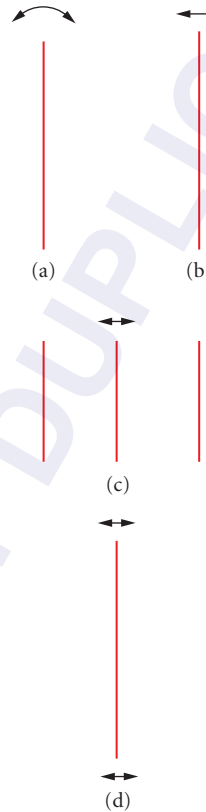


FIGURE 12 Configuration in which location differences can be detected with a precision higher than the resolution limit by up to an order of magnitude. These hyperacuity tasks do not contradict any laws of optics and are the result of sophisticated neural circuits identifying the centroids of light distributions and comparing their location. Thresholds are just a few arcsec in the foveal discrimination of these patterns from their null standards: (a) orientation deviation from the vertical of a short line; (b) alignment or vernier acuity; (c) bisection of a spatial interval; and (d) deviation from straightness of a short line.

Localization discriminations in the hyperacuity range are performed by identification of the centroid of the retinal light distributions²² of the involved pattern components, for example, abutting lines in the vernier task. The information is available in the image and can be described equally well in the domains of light distribution and its Fourier spectrum. It is testimony of sophisticated neural processing to arrive at the desired decision. Resolution and localization acuity share many attributes, though usually not in the same numerical measure.

Like ordinary visual acuity, hyperacuity is susceptible to crowding and to the oblique effect, but the two classes of spatial discrimination do not share all the other attributes.²³ Specifically, hyperacuity is more robust to reduction in exposure duration and diminishes more steeply with retinal eccentricity. It therefore follows that the neural processing apparatus is different and involves subtly recognition of differences in the excitation state of a neural population. In this respect it is similar to a whole host of fine discriminations, for example, those in the visual domains of color and stereoscopic depth.

4.10 RESOLUTION, SUPERRESOLUTION, AND INFORMATION THEORY

The threshold difference between ordinary visual acuity and hyperacuity raises the question whether any fundamental physical principles are being disobeyed.²⁴

Resolution and Superresolution

In the first instance, all of vision must satisfy the laws of physical optics according to which no knowledge can be acquired about an object that is contained in the region beyond the cutoff spatial frequency decreed by diffraction theory. Here the concepts associated with the term *superresolution* (Chaps. 3 and 4 in Vol. I) are relevant; they have been expanded since first formulated as involving an extrapolation: If the predominant features in which the spatial-frequency spectrum of two objects differ are located beyond the cutoff, but are always accompanied by a characteristic signature within it, then in principle it is possible to make the distinction between the two objects from detailed study of the transmitted spectrum and extrapolate from that to arrive at a correct decision. More recently the word has also been used to describe a procedure by which several samples of noisy transmitted spatial-frequency spectra from what is known to be the same object are superimposed and averaged. In both uses of the word, no diffraction theory limit has been breached; rather, knowledge is secured from detailed analyses of observed energy distributions on the basis of prior information—in the first case that the difference in the spatial-frequency spectrum inside the cutoff limit is always associated with those beyond it, in the second that the several spectra that are averaged arose from the same target.

Resolution and Information

Information theory has been of help in understanding some potentially knotty problems in defining resolution.²⁵ Traditionally, the resolving power of an ideal optical instrument is regarded to have been reached when two stars are separated by the width of the Airy disk, when the images of the two stars partially overlap and the dip between the two peaks is about 19 percent. This so-called Rayleigh criterion usually allows the receptive apparatus to signal the two peaks and the trough as separable features—provided its spatial grain is fine enough and the intensity difference detectable. It is then possible to decide, without prior knowledge, that there are two stars or that a spectral line is double. But when there is prior knowledge that the source can only be either single or double, and if double, what their relative intensity is, then the decision can be made for smaller separations, because the light distribution resulting from the overlapping images of two adjoining sources differs in predictable fashion from that of a single source, even when the dip between the peaks is less than at the

Rayleigh limit and even when there is no dip at all. This specific example highlights how information theory, as first formulated by Shannon, enters the discussion: in quantifying the transmitted information, prior knowledge needs to be factored in.

The diffraction image of a point is an expression of the uncertainty principle: it describes the probability distribution of the location of the source of an absorbed photon. The larger the number of absorbed photons, the more secure the knowledge of the source's location *provided always that it is known that it has remained the same source as the photons are being accumulated*. With sufficient number of absorbed photons, and the assurance that they arose from the same source, it is possible to assign location with arbitrary precision. Hence, from the standpoint of optics, there is nothing mysterious about the precision of locating visual target that is now called hyperacuity. Rather it draws attention to the physiological and perceptual apparatus that enables this precision to be attained.

4.11 SUMMARY

The human visual system's capacity to discriminate spatial details is governed by the eye's optical imagery, by the grain of the retinal receiving layer, by the physiological apparatus of the retina and the central nervous system, and by methodological considerations of recording observers' responses. Overall, under the best conditions of ordinary viewing, the resolution limit is close to that governed by the diffraction theory for the involved optical parameters. Performance is impaired whenever any of the optical, anatomical, physiological, or perceptual components does not operate optimally. There is an additional class of visual threshold based not on resolution of spatial detail but on locating relative object position. Because thresholds then are at least one order of magnitude better, they are called hyperacuity. They do not contravene any laws of optics but are testimony to sophisticated neural processing that can identify the location of the centroid of retinal light distributions with high precision.

4.12 REFERENCES

1. G. Westheimer, "The Visual System and Its Stimuli," *The Senses. A Comprehensive Reference*, R. H. Masland, (ed.), Academic Press, Oxford, 2008.
2. N. A. Macmillan and C. D. Creelman, *Detection Theory*, 2nd ed., Erlbaum, New York, 2005.
3. G. Westheimer, "Scaling of Visual Acuity Measurements," *Arch. Ophthalmol.* **97**(2):327–330 (1979).
4. I. L. Bailey and J. E. Lovie, "New Design Principles for Visual Acuity Letter Charts," *Am. J. Optom. Physiol. Opt.* **53**:740–745 (1976).
5. D. Y. Teller, "First Glances: The Vision of Infants. The Friedenwald Lecture," *Invest. Ophthalmol. Vis. Sci.* **38**:2183–2203 (1997).
6. I. Borish, *Clinical Refraction*, Saunders, Philadelphia, 1998.
7. G. Westheimer, "Visual Acuity," *Adler's Physiology of the Eye*, 10th ed., P. L. Kaufman and A. Alm, (eds.), Mosby, St. Louis, 2003, pp. 453–469.
8. L. Laurence, *Visual Optics and Sight Testing*, 3rd ed., School of Optics, London, 1926.
9. G. Westheimer, "Optical Properties of Vertebrate Eyes," *Handbook of Sensory Physiology*, M. G. F. Fuortes, (ed.), Springer-Verlag, Berlin, 1972, pp. 449–482.
10. H. H. Hopkins, "The Application of Frequency Response Techniques in Optics," *Proc. Physical Soc.* **79**:889–918 (1962).
11. T. Wertheim, "Ueber die indirekte Sehschärfe," *Z. Psychol.* **7**:172–189 (1894).
12. S. Shlaer, "The Relation between Visual Acuity and Illumination," *J. gen. Physiol.* **21**:167–188 (1937).
13. F. L. van Nes and M. A. Bouman, "Spatial Modulation Transfer in the Human Eye," *J. Opt. Soc. Am.* **57**:401–406 (1967).

14. D. G. Pelli, J. G. Robson, and A. J. Wilkins, "The Design of a New Letter Chart for Measuring Contrast Sensitivity," *Clinical Vision Sciences* **2**:187–199 (1988).
15. D. Regan, "Low Contrast Letter Charts and Sine Wave Grating Tests in Ophthalmological and Neurological Disorders," *Clinical Vision Sciences* **2**:235–250 (1988).
16. G. Westheimer, "Visual Acuity with Reversed-Contrast Charts: I. Theoretical and Psychophysical Investigations," *Optom. Vis. Sci.* **80**:745–748 (2003).
17. G. Westheimer, P. Chu, W. Huang, T. Tran, and R. Dister, "Visual Acuity with Reversed-Contrast Charts: II. Clinical Investigation," *Optom. Vis. Sci.* **80**:749–750 (2003).
18. W. S. Baron and G. Westheimer, "Visual Acuity as a Function of Exposure Duration," *J. Opt. Soc. Am.* **63**(2):212–219 (1973).
19. M. C. Flom, F. W. Weymouth, and D. Kahneman, "Visual Resolution and Contour Interaction," *J. Opt. Soc. Am.* **53**:1026–1032 (1963).
20. B. L. Beard, D. M. Levi, and L. N. Reich, "Perceptual Learning in Parafoveal Vision," *Vision Res.* **35**:1679–1690 (1995).
21. G. Westheimer, "Is Peripheral Visual Acuity Susceptible to Perceptual Learning in the Adult?," *Vision Res.* **41**(1):47–52 (2001).
22. G. Westheimer and S. P. McKee, "Integration Regions for Visual Hyperacuity," *Vision Res.* **17**(1):89–93 (1977).
23. G. Westheimer, "Hyperacuity," *Encyclopedia of Neuroscience*, L. A. Squire, (ed.), Academic Press, Oxford, 2008.
24. G. Westheimer, "Visual Acuity and Hyperacuity: Resolution, Localization, Form," *Am. J. Optom. Physiol. Opt.* **64**(8):567–574 (1987).
25. G. Toraldo di Francia, "Resolving Power and Information," *J. Opt. Soc. Am.* **45**:497–501 (1955).

This page intentionally left blank.

DO NOT DUPLICATE

OPTICAL GENERATION OF THE VISUAL STIMULUS

Stephen A. Burns

*School of Optometry
Indiana University
Bloomington, Indiana*

Robert H. Webb

*The Schepens Eye Research Institute
Boston, Massachusetts*

5.1 GLOSSARY

A	area
D	distance
E_r	illuminance at the retina (retinal illuminance)
f	focal length
f_e	focal length of the eye
L_s	intrinsic luminance of the source
td	troland (the unit of retinal illuminance)
x	position
Δ	change in position
Φ	flux (general)
Φ_p	flux at the pupil
τ	transmittance

We have also consistently used subscripted and/or subscripted versions of A and S for area, D for distance, f for focal lengths, x for positions, and Δ for changes in position.

5.2 INTRODUCTION

This chapter presents basic techniques for generating, controlling, and calibrating the spatial and temporal pattern of light on the retina (the visual stimulus). It deals with the optics of stimulus generation and the control of light sources used in the vision laboratory. Generation of stimuli by computer video displays is covered in detail in Chap. 22. Units for measuring radiation are discussed in Chap. 34, "Radiometry and Photometry," by Edward Zalewski and Chap. 37, "Radiometry and Photometry for Vision Optics," by Yoshi Ohno in Vol. II.

5.3 THE SIZE OF THE VISUAL STIMULUS

The size of a visual stimulus can be specified either in terms of the angle which the stimulus subtends at the pupil, or in terms of the physical size of the image of the stimulus formed at the retina, for most purposes vision scientists specify stimuli in terms of the angular subtense at the pupil. If an object of size h is viewed at distance D then we express its angular extent as

$$\text{Degrees visual angle} = 2 \left(\tan^{-1} \left(\frac{h}{2D} \right) \right) \quad (1)$$

or, for angles less than about 10 deg

$$\text{Degrees visual angle} \cong \frac{360h}{2\pi D} = \frac{57.3h}{D} \quad (2)$$

In this chapter we specify stimuli in terms of the actual retinal area, as well as the angular extent. Angular extent has the advantage that it is independent of the eye, that is, it can be specified totally in terms of externally measurable parameters. However, an understanding of the physical dimensions of the retinal image is crucial for understanding the interrelation of eye size, focal length, and light intensity, all of which are part of the design of optical systems for vision research.

5.4 FREE OR NEWTONIAN VIEWING

There are two broad classes of optical systems that are used in vision research. Free viewing, or *newtonian* viewing, forms an image of a target on the retina with minimal accessory optics (Fig. 1a).

Retinal Illuminance

Photometric units incorporate the overall spectral sensitivity of the eye, but that is the only special property of this set of units. That is, there are no terms specific to color (for a discussion of colorimetry, see Chap. 10) and no allowances for the details of visual perception. The eye is treated as a linear detector which integrates across wavelengths. The retinal illuminance of an object in newtonian view is determined by the luminous intensity of the target and the size of the pupil of the eye. The luminous power at the pupil [dimensions are luminous power or energy per unit time, the SI units are lumens (lm)]

$$\Phi_p = L_s A_s \frac{A_p}{D^2} \quad (3)$$

where L_s is the luminance of the source [the SI units are lumens per meter squared per steradian ($\text{lm}/\text{m}^2/\text{sr}$) or candelas per meter squared (cd/m^2)], A_s is the source area, A_p is the area of the pupil, and D is the distance from the pupil to the source [so (A_p/D^2) is the solid angle the pupil subtends].

The area of the image element on the retina is

$$A'_R = A_s m^2 = A_s \left(\frac{f_e}{D} \right)^2 \quad (4)$$

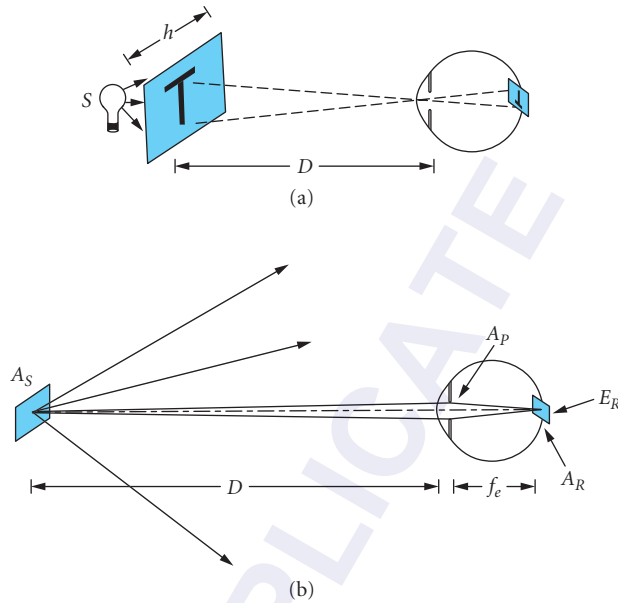


FIGURE 1 (a) In free viewing or Newtonian viewing the eye's optics are used to image a target onto the retina. (b) Computation of the retinal illuminance in free-viewing. A_s the area of the source, D the distance between the source and the eye, A_p the area of the pupil, f_e the optical focal length of the eye, A_r the area of the image of the source on the retina, and E_r the retinal illuminance.

where m is the magnification of the eye and f_e is the effective focal length of the eye (Fig. 1b). From this we compute the illuminance at the retina as

$$E_R = \frac{\Phi_p}{A'_R} = \frac{A_p L_s}{(f_e)^2} \quad (5)$$

(The SI units of illuminance are lm/m^2). A typical value for f_e is 16.67 mm.^{1,2} Note that the retinal illuminance does not depend on the distance. That is, the retinal illuminance when viewing an extended source such as a video screen is independent of the viewing distance and is dependent only on the size of the eye's pupil, the luminance of the source, and the focal length of the eye. In most cases, the focal length of the viewer's eye is not known, but it is possible to measure the size of the pupil. For this reason a standard unit for specifying retinal illuminance was developed, the *troland*.

The Troland (td) The troland is a unit of illuminance (luminous power per unit area). The troland quantifies the luminous power per unit area at the retina (the retinal illuminance). One troland was defined as the illuminance at the retina when the eye observes a surface with luminance = 1 cd/m^2 through a pupil having an area of 1 mm^2 . Using Eq. (5) and a standard f_e of 16.67 mm we find that the troland is defined by

$$1 \text{ td} = 0.0035 \text{ lumens}/\text{m}^2 \quad (6)$$

This definition ties the troland to the illuminance on the retina of a standard eye assuming no transmission loss in the ocular media at 555 nm. Thus, two observers with different-size eyes, different-size pupils, or different relative losses in the ocular media, viewing the same surface, will have different retinal illuminances. Wyszecki and Stiles² have recommended that the term *troland value* be used to distinguish the trolands computed for a standard eye from the actual retinal illuminance. General usage is that the retinal illuminance is determined simply by measuring the luminance of a surface in cd/m^2 and multiplying this value by the area of the pupil in mm^2 .

Limitations of Free Viewing: An Example

There are two major limitations to newtonian view systems. The first is that the retinal illuminance is limited. For instance, a 60-W frosted incandescent bulb can produce a 120,000-td field, but to obtain a uniform 20 deg field it must be placed 17 inches from the observer's eye. This requires an accommodative effort that not all observers can make. A comparable illuminance at more realistic distances, or with variable focus, requires larger light sources or more elaborate optical systems. The second limitation of free viewing is that variations in pupil size are not readily controlled. This means that the experimenter cannot specify the retinal illuminance for different stimulus conditions or for different individuals. Maxwellian view optical systems solve these problems.

5.5 MAXWELLIAN VIEWING

Figure 2a shows a simple *Maxwellian* view system. The key factor that distinguishes the maxwellian view system is that the illumination source is made optically conjugate to the pupil of the eye. As a result, the target which forms the stimulus is not placed at the source, but rather at a separate plane optically conjugate to the retina. In the system of Fig. 2a the plane conjugate to the retina, where a target should be placed, lies at the focal point of lens L_2 , between the source and the lens (see Chap. 6). Light from the source is diverging at this point, and slight changes in target position will cause changes in both the plane of focus and the magnification of the target. For this reason most real maxwellian view systems use multiple lenses. Figure 2b shows such a maxwellian view system where the single lens has been replaced by two lenses. The first lens collimates light from the source and the second forms an image of the source at the pupil. This places the retinal conjugate plane at the focal plane of L_1 . We label the conjugate planes starting at the eye, so R_1 is the first plane conjugate to the retina and P_1 is the first plane conjugate to the pupil.

Control of Focus and the Retinal Conjugate Plane

The focus of the maxwellian view system is controlled by varying the location of the target.³⁻⁵ Moving the target away from R_1 allows the experimenter either to adjust for ametropia or to require the observer to accommodate. To see this we compute where lens L_1 (Fig. 2b) places the image of the target relative to the eye. If the target is at the focal point of L_1 , the image is at infinity, which means there is a real image in focus at the retina for an emmetropic eye. If we move the target toward lens L_1 , the lens of an emmetrope must shorten the focus to bring the image back into focus on the retina. We quantify this required change in focus as the change in the dioptric power of the eye (or the optical system if spectacle correction is used) where the dioptric power is simply the inverse of the focal length measured in meters. If the target is displaced by Δ from the focal point of L_1 then using the newtonian form of the lens formula ($x'x = f^2$) we find that

$$\text{Change in dioptric power} = \frac{\Delta}{(f_1)^2} \quad (7)$$

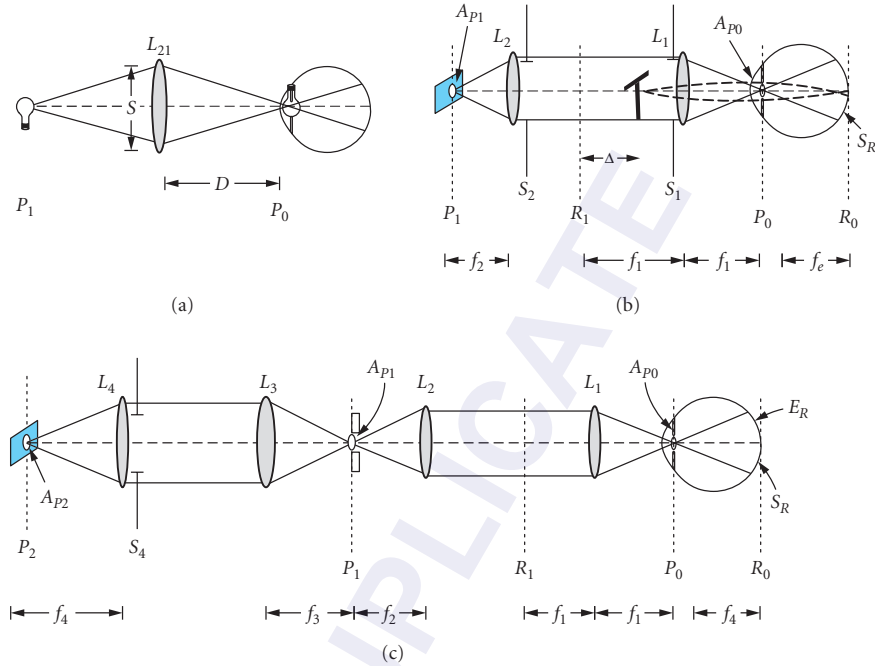


FIGURE 2 (a) In a minimal maxwellian view system a lens L_{21} images a source in the plane of the eye's pupil. The lens is shown at a distance D from the eye equal to twice its focal length. The field stop (S) in the aperture of L_{21} . (b) A maxwellian view optical system with an accessible retinal conjugate plane (R_1). Lens L_2 collects light from the source and collimates it. Lens L_1 images the source in the plane of the eye's pupil (P_0). Pupil conjugate planes (P_i) and retinal conjugate planes (R_i) are labeled starting at the eye. For an emmetrope lens L_1 images the retina (R_0) at R_1 . This is the plane where a visual target will appear in best focus for an emmetropic observer. Moving the target away from R_1 a distance Δ changes the plane of focus of the retinal image, requiring the emmetrope to accommodate. The maximum size of the retinal image (S_R) is limited for this system to the aperture of lens L_2 (S_2). The pupil is imaged at position P_1 , where the light source is placed. (c) Shows a more complex system where a second set of pupil conjugate and retinal conjugate planes have been added. An artificial pupil (A_{p1}) is added at P_1 . It is the image of this pupil at the source (A_{p2}) that limits the retinal illuminance produced in this optical configuration. Other symbols as above.

If the target move toward the lens L_1 , this compensating power must be positive (the focal length of the eye or eye plus spectacle shortens). If the target moves from f_1 away from lens L_1 , the required compensation is negative. If the eye accommodates to keep the target in focus on the retina, then the size of the retinal image is unchanged by the change in position.^{3,4}

Size

For this system the area of the illuminated retinal field is determined by the limiting aperture in a retinal conjugate plane between the light source and the retina (the field stop S_2 in Fig. 2b) and is

$$S_R = \left(\frac{f_e}{f_1} \right)^2 S_2 \quad (8)$$

or for the schematic eye we use

$$S_R = \left(\frac{16.67}{f_1} \right)^2 S_2 \quad (9)$$

where f_e is the focal length of the eye (in mm) and f_1 is the focal length of lens L_1 . This formula is used to compute the physical size of an image on the retina. Note that the linear magnification is (f_e/f_1) and the areal magnification is $(f_e/f_1)^2$. The angular subtense of the field is the same as in Eq. (1), substituting the diameter of S_2 for h and f_1 for the distance D . Also note that the angular extent is independent of details of the optics of the eye. This convenience is the main reason that angular extent is the most widely used unit for specifying the stimulus in vision research.

Retinal Illuminance

One of the *principal advantages* of a maxwellian view system is that it provides a large, uniformly bright field of view. This mode of illumination is called Kohler illumination in microscopy. The light available in a maxwellian view optical system is determined by two main factors, the luminance of the light source, and the effective pupillary aperture of the optical system being used. In Fig. 2b we see that lens L_2 collects light from the source. The amount of light collected is the maximum amount of light that can be presented to the eye and is^{1,3}

$$\Phi = A_{p1} L_s \frac{S_2}{(f_2)^2} \quad (10)$$

where A_{p1} is the area of source being considered (or, equivalently, a unit area of the source), L_s is the luminance of the source (in $\text{lm}/\text{m}^2/\text{sr}$), S_2 is the aperture of lens L_2 , and f_2 is the focal length of lens L_2 . We have used the area of lens L_2 , for this example, although in actual practice this area might be reduced by later field stops. This quantity will cancel out in the subsequent calculations. Finally, if all of the light collected by lens L_2 is distributed across the retina, then the retinal illuminance is

$$E_R = \frac{\Phi}{S_R} = A_s L_s \left(\frac{f_1}{f_2 f_e} \right)^2 \quad (11)$$

where S_R is obtained from Eq. (8). Note that only the luminance of the source, the source area, and the focal lengths of the optical elements are important in setting the retinal illuminance. Using different diameters for lens L_2 will change the amount of light collected and the size of the retinal area illuminated. Such a change will be reflected by changing the area parameter in both Eqs. (8) and (10), which cancel.

In Eq. (11) the area of the source image (A_{p1}) at the pupil is

$$A_{p0} = A_{p1} \left(\frac{f_1}{f_2} \right)^2 \quad (12)$$

If the eye's pupil (A_{p0}) is smaller than this source image then the eye's pupil is limiting and

$$E_R = \frac{A_{p0} L_s}{(f_e)^2} \quad (13)$$

Equation (13) for the maxwellian view system is identical to Eq. (5) which we obtained for the newtonian view system, but now the entire field is at the retinal illuminance set by the source luminance. Thus, in maxwellian view a large, high retinal illuminance field can be readily obtained.

To control the size of the entry pupil rather than allowing it to fluctuate with the natural pupil, most maxwellian view systems use a pupillary stop. Figure 2c shows a system where a stop has been introduced at an intermediate pupil conjugate plane A_{p1} . This has the advantage of placing the pupillary stop of the system conjugate to the eye's pupil. The projection of A_{p1} at the source is A_{p2} and it is A_{p2} that limits the available luminous area of the source. Lenses L_3 and L_4 image the source onto the artificial pupil, and lenses L_1 and L_2 image the artificial pupil in the plane of the eye's pupil (A_{p0}). The retinal illuminance of this more complex system can be computed as follows: the field stop over which light is collected by lens L_4 is S_4 and can be computed by projecting the retinal area illuminated (S_R) back to lens L_4 .

$$S_4 = S_R \left(\frac{f_1 f_3}{f_2 f_e} \right)^2 \quad (14)$$

A_{p2} is the area of the source which passes through the artificial pupil:

$$A_{p2} = \left(\frac{f_4}{f_3} \right)^2 A_{p1} \quad (15)$$

Therefore the total amount of usable light collected is

$$\Phi = L_s A_{p2} \frac{S_4}{(f_4)^2} = S_R L_s A_{p1} \left(\frac{f_1}{f_2 f_e} \right)^2 \quad (16)$$

and the retinal illuminance is

$$E_R = \frac{\Phi}{S_R} = L_s A_p \left(\frac{f_1}{f_2 f_3} \right)^2 \quad (17)$$

Note that, as in Eq. (13), $[A_{p1}(f_1/f_2)^2]$ is the size of the image of the aperture A_{p1} when measured at the exit pupil (A_{p0}). Thus, even in this more complex case we find that the retinal illuminance is dependent only on the source luminance, the pupillary area measured at P_0 , and the focal length of the eye.

Advantages of Maxwellian Viewing: Example Revisited

The strength of a maxwellian view system is that properties of the target (focus, size, shape) can be controlled independently from retinal illuminance. If the source is larger than the limiting pupil, then the maximum retinal illuminance is the luminance of the source, scaled only by the exit pupil and eye size [Eqs. (13) and (17)]. The retinal illuminance is controlled by pupillary stops, and the target size, shape, and focus are controlled by retinal stops. Additionally, focus can be controlled independently of retinal illuminance and image size.

If we use the 60-W frosted bulb that we used in the example of newtonian viewing as a source in a maxwellian view optical system we find that we can produce the same maximum retinal illuminance for a given pupil size [Eqs. (5) and (13)]. However, with the maxwellian view system relatively inexpensive, achromat lenses will allow us to generate a field size greater than 20 deg. In addition, we can dispense with the frosted bulb (which was needed in free viewing to give a homogeneously illuminated target) and use an unfrosted tungsten halogen bulb to produce in excess of 1 million td.

Controlling the Spatial Content of a Stimulus

The spatial frequency content of a retinal image is usually measured in cycles per degree (of visual angle). The frequencies present in the retinal image are those present in the target, cut off by the limiting aperture of the pupillary stop—either the eye's pupil or an artificial pupil in a pupillary plane *between the target and the eye*.²⁵ Apertures placed on the source side of the target do nothing to the spatial frequency content of the retinal image, leaving the pupil of the eye controlling the image, not the artificial pupil. (For a more detailed treatment of the effects of pupils on the spatial content of images see Refs. 6–10. For a discussion of the spatial resolving power of the eye's optics see Chaps. 1 and 4 in this volume of the *Handbook*).

Positioning the Subject

One of the practical disadvantages in using maxwellian view systems is the need to position the eye's pupil at the focal point of the maxwellian lens (L_1 in Fig. 2c) and to maintain its position during an experimental session. One technique for stabilizing the position of the eye is to make a wax impression of the subject's teeth (a bite bar). By attaching the bite bar to a mechanism that can be accurately positioned in three dimensions (such as a milling machine stage) the eye can be aligned to the optical system and, once aligned, the eye position is maintained by having the subject bite loosely on the bite bar. Alignment of the eye can be achieved either by observing the eye, as described at the end of this chapter, or by using the subject to guide the alignment process as described in the following paragraphs.

The first step in alignment is to position the eye at the proper distance from lens L_1 . This can be accomplished by noting that, when the eye's pupil is at the plane of the source image, slight movements of the head from side to side cause the target to dim uniformly. If the eye's pupil is not at the proper distance from lens L_1 then moving the head from side to side causes the target to be occluded first on one side, then the other. By systematically varying the distance of the eye from lens L_1 , the proper distance of the eye can be determined. It is then necessary to precisely center the eye's pupil on the optical axis of the apparatus. There are several approaches to centering the eye.

1. The physical center of the pupil can be located by moving the translation stage such that the pupil occludes the target first on one side, then on the other. The center of these two positions is then the center (in one dimension) of the pupil. The process is then repeated for vertical translations.
2. The entry position that produces the optimum image quality of the target can be determined. One variant is to use a target that generates strong chromatic aberration (for instance, a red target on a blue background). The head can then be moved to minimize and center the chromatic fringes.¹¹ This process defines the achromatic axis of the eye.
3. The eye can be positioned to maximize the brightness of the target, which centers the Stiles-Crawford maximum in the pupil with respect to the exit pupil of the instrument.
4. A set of stimuli can be generated that enter the eye from different pupil positions. If the eye is centered, then all stimuli will be seen. If the eye is not centered, then the pupil occludes one of the pupil entry positions and part of the stimulus array disappears. In this case, the subject merely has to keep his or her head positioned such that all stimuli are visible.

5.6 BUILDING AN OPTICAL SYSTEM

Alternating Source Planes and Retinal Planes in a Controlled Manner

The separation of retinal and pupil conjugate planes in a maxwellian view system allows precise control over the spatial and temporal properties of the visual stimulus. By placing the light source conjugate to the pupil of the eye, every point on the source projects to every point in the retinal image

and vice versa. Thus, to control the whole retinal image, such as turning a light on and off, manipulation of a pupil conjugate plane is optimal. To control the shape of the retinal image without altering the entry pupil characteristics, variation at the retinal conjugate planes is required. However, there is an exception to this rule. Light from the edges of the image traverse the pupil conjugate plane at a higher angle than light from the center of the image. For small fields (<15-deg diameter) the angular dependence is minimal, but for larger fields it can be significant and filters should be placed in a collimated light path.

Combining Lights in an Optical System

To control aspects of the visual stimulus independently, different light sources and targets, each with its own set of filters and shutters, can be combined. We will call these separate “channels.” Three different techniques allow for combining optical channels, beamsplitters, beam separators or reflective apertures, and diffusers. All three methods are demonstrated in Fig. 3.

Beamsplitters A beamsplitter both transmits and reflects a fraction of the incident light. Any surface with an index of refraction change can be a beamsplitter. In vision we generally use either a cube or a plate beamsplitter. By locating the beamsplitter in a collimated portion of the optical system (Fig. 3), but away from a retinal conjugate plane, two channels can be combined. For channels with different spectral compositions the beamsplitter can be dichroic, reflecting some wavelengths and transmitting others, usually by interference effects. Plate beamsplitters have the disadvantage

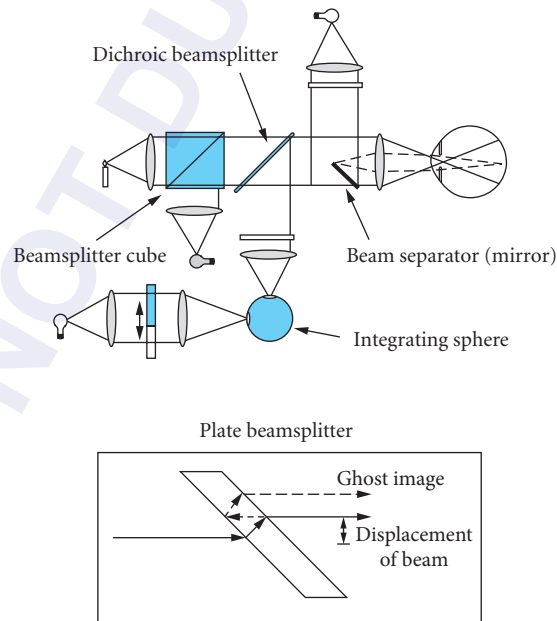


FIGURE 3 Techniques for combining light from different light sources or optical paths. The inset illustrates potential problems that need to be considered when using a plate style of beamsplitter. These problems are eliminated by the use of a cube beamsplitter or wedge beamsplitter.

that there are secondary reflections, slightly displaced (inset, Fig. 3) from the main beam. In an image-carrying beam this produces “ghost” images, slightly displaced from the main image. The displacement decreases with decreasing thickness. A pellicle is simply a plate beamsplitter with negligible thickness. A wedge beamsplitter displaces a ghost image a long way, and Liebman has used this to design a tapered beamsplitter with no secondary reflections.¹² Unlike plate beamsplitters, the beamsplitter cube has the advantage that there is no lateral displacement of an image.

Beam Separators (or Reflective Apertures) A beam separator combines two optical channels into one, while keeping the light spatially distinct. For instance, a mirror placed halfway across the aperture at 45 deg combines two channels at 90 deg to each other. Such a separator in a retinal conjugate plane is imaged at the retina to produce a bipartite field (Fig. 3). The mirror must be oriented to hide the bevel from the subject and must have a high-quality edge. A transparent plate evenly silvered in a pattern and optically sandwiched between right angle prisms is a convenient beam separator with good edges.

Diffusers, Integrating Spheres, and Optical Fibers Diffusers also mix light. Diffusers are used when either the spatial uniformity of the final beam or thorough spatial mixing of light sources is important. The direction of light remitted from the diffuser is independent of the incident light, which simplifies the combination of sources. Both integrating spheres and integrating bars have been used to combine lights (Fig. 4a and b, respectively). Diffusers have been widely used in the construction of colorimeters.² In the La Jolla colorimeter¹³ light is passed through a filter assembly with three monochromatic filters assembled edge to edge. Moving the filter assembly across the light path changes the relative proportions of the light impinging on each filter, which changes the average color of the light. An integrating sphere then completely mixes the colors producing a variable chromaticity light source without the spatial inhomogeneity of the filters. The major disadvantages of using diffusers are light loss, the need for careful optical baffling (since any light impinging on the diffuser gets mixed into the system) and the gradual deterioration of the diffuser due to dirt and dust. While the light loss from an

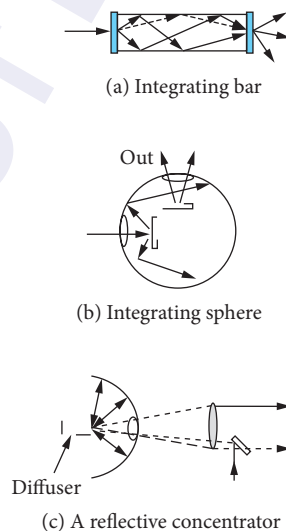


FIGURE 4 An example of different types of beam combiners based on diffusers. In general it is critical that high-efficiency diffusing materials be used in all of these techniques.

integrating sphere is ideally very low, slight decreases in the efficiency of the diffusing surfaces cause large losses in the output. Collection of light from a simple diffuser can be increased by the use of a spherical mirror¹⁴ (Fig. 4c). Similarly, a noncoherent fiber-optic bundle can be used to mix lights.

Mixing coherent light sources with diffusers requires movement to minimize interference (speckle) effects. Either liquid diffusers¹⁵ which work by Brownian motion or moving diffusers¹⁶ can be used.

Lenses

Lens orientation affects any optical system. For instance, the most curved surface of a lens should be toward the more collimated light. This makes both sides of the lens share the work of refraction and decreases aberrations. Achromats reduce both spherical and chromatic aberrations (they have three or more refracting surfaces) and are typically designed to form a collimated beam.

Typically, the goal in aligning an optical system is to place all lenses upon a common optical axis, oriented such that a ray of light traversing that axis is not deviated. A simple prescription is

1. Define the system axis using pinholes, pointers, or a laser beam, with all lenses removed. In a multichannel optical system it is often easiest to start from the position of the eye's pupil (P_0).
2. Introduce mirrors, beamsplitters, and other nonrefractive optics in their proper positions, being careful to keep the beam at the same height. A technique for ensuring a constant height, such as a movable pinhole, also helps.
3. Starting at the alignment source, set each lens, one at a time, such that its Boys points lie on the system axis. The Boys points are the focused reflections from the various curved surfaces of lenses (see Ref. 17 for a detailed discussion of using Boys points for alignment). When using a laser they are the "bull's eye" interference patterns. Lateral translation of the lens moves the first (brightest) Boys point and rotation moves the second, so convergence is rapid.
4. To determine the position of a lens along the optical axis, a photographic loop (a magnifying lens) and a target grid (graph paper) can be used. By placing the graph paper at a known location such as at the exit pupil, and then looking with the photographic loop into the system (turn off the laser first!) the exact position of the next pupil conjugate point can be determined. To position a lens to collimate light from a source, place a mirror in the collimated beam, reflecting the light back into the lens. When the lens is placed at one focal length from the source, the mirror causes the source to be imaged back onto itself (autocollimation).

Field Quality

The uniformity of the illumination of the retinal image is controlled by the size and emission characteristics of the source and by the uniformity of illumination of the target. Ideally, a source emits uniformly into all angles subtended by the collection lens (lens L_2 in Fig. 2). LEDs do not meet this requirement and produce inhomogeneous illumination unless the light is further diffused. Tungsten-halogen sources with coiled-coil filaments uniformly illuminate the collecting lens and can produce uniform retinal images. Problems arise if a retinal conjugate is placed too close to a pupil conjugate plane. This occurs most often when the target is placed at the end of a long collimated portion of the optical path. In this case, the structure of the source becomes visible at the retinal plane. In general, collimation lengths should be kept to less than roughly twice the focal length of the collimating lens. It is also important to control the pupillary conjugate points; in some cases multiple pupil conjugate apertures can help.¹⁸

Controlling Wavelength

Spectral composition of a light source can be varied either by absorbing the unwanted light (in a filter), or by redirecting it (by either interference or refraction). Table 1 presents techniques for controlling wavelength. Interference filters can have blocking filters that absorb wavelengths outside the specified transmission band. However, a blocking of 10^{-3} may not be sufficient for the human visual

TABLE 1 Controlling Wavelength

Technique	Properties	Advantages	Disadvantages/Comments
Absorption filters: aka dye, film, gel, or pigment filters	Available in both gels and glass forms. Generally broadband.	Can be extremely stable over time. Wide range of colors available. Reasonably priced.	Broadband, can be damaged by too much light, as encountered with lasers and focused beams. Most fluoresce.
Narrowband interference	Half-width about 10 nm, usually bonded to an absorption filter for blocking distant wavelengths.	Convenient, readily available, and moderately priced.	If not fully sealed, humidity can cause long-term degradation. Wavelength depends on angle of incidence. Need good blockers.
Broadband interference (see also dichroic beam splitters)	Half-width 40–100 nm.	Same as narrowband interference filters, less need for blocking filters.	
Grating monochromators	Tunable, bandwidth can be set by slits.	Can have low stray light.	Must match aperture of system to the interior grating to minimize stray light. Expensive.
Prism monochromators	Produces a spectrum by refraction.	High throughput.	Bulky.
Interference wedges (spatially varying interference filters)	Allow continuous change of wavelength by changing the position of light incidence.	Small, easy to use.	Some have more leakage than a monochromator, may not be adequately blocked for higher-order transmission. Must be placed at a pupil conjugate point to avoid chromatic wedging on the retina.
Special sources (Na, Cd, etc.)	Discharge sources that produce a few exactly specified spectral lines.	Many lines are available. Very pure wavelengths. Larger aperture sources are available.	Mostly replaced by lasers which are stable, cheap, and easier to use but are point sources (see Table 5).

system. For instance, a narrow 670-nm filter may transmit only 10^{-3} of the intensity at 555 nm, but the human visual system is about 30 times as sensitive to 555 nm as to 670 nm and can integrate over the entire visible spectrum. Thus, the visual effectiveness of the “blocked” light can be considerable. Monochromatic and narrowband light sources and lasers can also be used both to provide a specific wavelength and to calibrate the wavelength scale of monochromators.

Wavelength selection filters are best placed in a collimated beam to avoid changing the passband to shorter wavelengths. This shift is asymmetric, and the maximum transmission changes as the angle of incidence decreases from 90 deg.² Differently colored lights can be combined using broadband interference filters (dichroic filters) placed at an angle, but using narrowband filters designed for normal incidence at 45 deg can introduce unwanted polarization properties.² Filters designed for use at 45 deg are readily available.

Turning the Field On and Off

Light from an optical channel can be clocked entirely by a shutter. Shutters should generally be located at a pupil conjugate plane. At this position the source is focused to a small area and thus a smaller, faster shutter can be used. In addition, manipulation of a pupil plane changes the retinal image uniformly (see above). Parameters of interest are the speed, repetition rate (and duty cycle), and the transmission of the shutter in both its open and closed states. Table 2 summarizes details on common types of shutters. Shutters that work by means of polarization generally need to be placed in a collimated beam, and their extinction (transmission when off) and maximum transmission are wavelength-dependent.

Controlling Intensity

The visual system operates over about 11 log units of luminance (11 orders of magnitude, or 10^{11}), and within this range can make discriminations of 0.1 percent or better. This level of visual capability requires a two-stage control of retinal illuminance. The overall retinal illuminance is set by

TABLE 2 Types of Shutters

Shutter	Speed	Size of Aperture (mm)	T off/T on	Comments
Mechanical shutter	ms	1 mm and up	0	Larger apertures are slower. Most designs cannot run in a continuous mode, but these are ideal for low duty-cycle use.
Galvanometers	ms	A few mm for high speeds	0	Operating at high rates requires fairly careful driver design but they can be used to excellent effect and are commercially available with drivers. Can be run continuously.
Choppers	kHz	Variable	0	For continuous on-off flicker choppers are ideal. They have no attenuation when open and no transmission when closed. Small, feedback-controlled devices are available off the shelf. Stepper motors with vanes mounted on the shaft can also be used.
Acousto-optic modulator (AOM)	μ s	<2 mm	$\sim 10^{-3}$	AOMs have small aperture and high f -number and thus work best for lasers. They are inherently chromatic and orientation-sensitive. Polarization effects can be important.
LCD shutters	ms	A few cm	$\sim 10^{-2}$	Speeds are increasing rapidly, works by polarization. Maximum transmission is less than 50 percent. Use in collimated beam.
LCD displays	μ s	Work as video screens; individual pixels can be μ m.	$\sim 10^{-2}$	Work by polarization. Speed is increasing with development of displays for computer use.
Kerr cells	ps	cm	$\sim 10^{-2}$	Require high voltages, work by polarization.

discrete neutral density filters, while smaller steps are set by either a variable density filter or an electronic system. Due to limited precision in most measuring instruments, it is hard to calibrate an individual filter with a density greater than 3.0.

Varying the Intensity of the Field The intensity of a channel can be varied by using fixed or variable filters, or by controlling the radiance of the source. For simple, relatively slow (1/min) changes in retinal illuminance, neutral density filters are appropriate. For faster changes, modulators or direct variation of the source radiance are typically required.

Filters Neutral density filters either have a uniform attenuation across their spatial extent (fixed filters) or vary in their attenuation characteristics depending on the spatial location (variable filters or “wedges”). Absorbing filters should be used in a collimated beam to avoid different path lengths through the filter. Wedges are used at a pupil conjugate plane to avoid imaging the density gradient on the retina. Fixed filters are typically used to set fixed parameters (such as the average luminance), variable filters to control the brightness continuously (for instance, in an increment threshold test). Table 3 presents some types of neutral density filters.

Modulators Passing linearly polarized light through a rotating polarizer produces sinusoidal modulation of the transmitted light. Careful optical design allows generation of an excellent variable modulation, flicker stimulus.^{19,20} Acousto-optic modulators (AOM) can be used to vary a light rapidly, but have only a small aperture. Mechanical methods have also been used for temporally modulating a target.^{21–23} For the most part, mechanical modulators and moving polarizers have been replaced by direct control of light sources by high-speed shutters and by video systems (see Chap. 22 in this volume of the *Handbook*).

TABLE 3 Neutral Density Filters

Method	Advantages	Disadvantages	Other Comments
Metal film neutral density filters	Readily available, stable, spectral neutrality is good.	Must be cleaned with care, pinholes can cause a problem. Interreflections between filters can cause deviations from density obtained when calibrated alone.	Silver and inconel are most common. Tilt to keep reflections out of the system.
Multilayer dielectric filters	Easier to clean.	Expensive, chromatic.	
Metal film “wedges”	Same as metal film filters.	Same as metal film filters.	Tilt to aim reflections out of the system. Use in a pupil conjugate plane to avoid spatial variation of image.
Absorbing gel filters	Less expensive, easily available at good photography stores.	They become brittle with age. Somewhat chromatic. Damaged by intense light.	Keep dry and away from a hot lamp. Use in a collimated light path.
Sector disk and mechanical devices (cams, variable slits, etc.)	Spectrally flat.	Must be used in a pupil plane to get desired results.	Use in a pupil conjugate (or in a uniform beam before an integrating sphere). Variable slits can be a simple, inexpensive way to control intensity.
Photographic film or film wedges	Cheap, easily made by flashing a film and developing.	Spectral neutrality will be poor in relation to commercial filters. Must be calibrated. Stability over time will depend on the film base used.	Best used with monochromatic light.

Varying the Source Varying the radiance of the light source is fast, and is particularly straightforward with electro-optical sources such as LEDs and diode lasers. However, control of the source is not limited to these devices and can be used with thermal sources (incandescent bulbs), xenon arcs, and fluorescent lights as well. There are four major problems to be overcome when directly varying a light source: (1) nonlinear current versus radiance characteristics, (2) changes in the current-radiance relation over time, (3) temporal hysteresis in the current-radiance relation, and (4) changes in the wavelength distribution of the source with changes in current.

For instance, LEDs have been considered ideal light sources by some investigators, yet they show the same deviations from ideal behavior common to other sources^{24–29} (see Chap. 15, “Artificial Sources,” by Anthony LaRocca, Chap. 17, “Light-Emitting Diodes,” by Roland H. Haitz, M. George Craford, and Robert M. Weissman, and Chap. 18, “High-Brightness Visible LEDs,” by Winston V. Schoenfeld in Vol. II). The current/radiance relation of LEDs depends on the internal temperature of the device. Thus, the relation changes with both the measurable *temperature* (due to both the environment and to the time average current) and with the immediate current history (for instance, it was just driven at a high current). Finally, the temperature also affects the probability of electronic transitions in the semiconductor, which can change the *wavelength* emitted. For LEDs used in visual work this wavelength dependence has been measured to be on the order of a 1.6-nm change in dominant wavelength with variation of the duty cycle for an ultrabright red LED.²⁶ Thermal sources, such as a tungsten-halogen bulb, undergo especially large changes in spectral output with changes in current. With these sources, only a narrow spectral band should be used if the source is controlled. In addition, the thermal inertia of an incandescent source precludes rapid modulation, although slow modulation (~1 Hz) can be achieved. Even fluorescent sources which have been widely used in some areas of research show significant changes in spectral output with time.³⁰

While heat sinks can help to stabilize the response characteristics of many devices, linear control still requires careful driver design. Drivers can use either an analog or binary design. Analog drivers control the source radiance by varying the current, while digital drivers turn the source either on or off, with the average radiance set by the proportion of time that the source is on (the duty cycle). An advantage of the binary scheme is that transistors dissipate less power in the on and off states than in intermediate states.

There are four major approaches to linearization:

1. Calibrate the current-luminance relation and build a driver with the inverse nonlinearity. For simple applications this technique is adequate. However, for demanding applications the driver becomes complex and the demands for stability in the driver raise new problems.
2. Calibrate the current-luminance relation and use an inverse nonlinearity in computer memory (a lookup table). This technique is quick and easy to implement. A static lookup table will not compensate for dynamic nonlinearities (hysteresis) but a variation of this approach, known as *delta modulation*, will. With delta modulation a linear detector is used to precalibrate the output of a binary driver to produce the desired waveform. By later playing back the binary sequence, the waveform can be re-created exactly. Thus, delta modulation can be used to compensate for all source nonlinearities except for changes in wavelength. The disadvantage is that the waveform must be precalibrated.
3. Vary the ratio of the on and off time periods (the duty cycle) of the source using a binary driver. With a fixed cycle time, the on time can be varied (pulse-width modulation or PWM). PWM works fairly well, but is sensitive to capacitance and nonlinear switching effects that can alter the waveform for either short on or off periods. A similar approach is to use fixed (short) pulses and vary the frequency of the pulses (pulse-frequency modulation or PFM). This approach has the advantage that every pulse is identical, thus capacitance and switching effects are minimized. PFM has been used to control the luminance of LEDs²⁷ and AOMs¹⁶ linearly over a 1000:1 luminance range.
4. Detect the light with a linear photodetector and use feedback within the driver to linearize the output. By using a PIN photodiode in photovoltaic mode (see Chap. 24, “Photodetectors,” by Paul R. Norton in Vol. II), it is possible to construct a very linear circuit.^{28,31} Light feedback can be used with either analog or binary drivers.

Generating Complex Temporal and Spatial Patterns

Any type of light source or visual stimulus, from video monitors to street scenes, can be integrated into an optical system, giving improved control of luminance and pupil position. Thus, almost anything can be used as a target. Stimuli can be moved without varying the position of light entry in the eye's pupil by (1) translating a target in the retinal plane, (2) using a cathode ray tube (CRT) or a liquid crystal display (LCD) in a retinal conjugate plane, or (3) rotating a mirror in a pupil conjugate plane.²²

Rotating a mirror in a pupil conjugate plane changes the angle at which light enters the eye, but not the pupil entry position. This technique has been used for generating motion and counterphase gratings²¹ and for decreasing the apparent inhomogeneity of a stimulus by moving the image at a rate above flicker fusion. Galvanometers are commonly used to generate the motion (Fig. 5).

Video Monitors and CRTs The combination of a video display with a maxwellian view optical system allows the experimenter the advantages of precise control of the pupil, use of a broad spectral range of backgrounds (provided by the traditional optics), and the fine spatial control available in a video system. The video monitor or CRT is placed at a retinal conjugate plane and the other optical channels are added to it. By passing the combined light path through an aperture at a pupil conjugate plane it is possible to control the pupil size for the video system, insert achromatizing lenses, and the like (see Fig. 8 later in the chapter).

Liquid Crystal Displays One type of video display that is of increasing utility is the liquid crystal display (LCD) (Chap. 22). LCDs, unlike video monitors, work by transmission or reflection; they are not self-luminous. Thus, while the wavelength composition of a monitor is determined by the characteristics of its phosphors, the LCD can be used with *any* light source. LCD displays work by changing the polarization state of local regions of the display (pixels). The pattern of pixels is controlled by a video signal or a computer. Like most polarization devices, the extinction ratio is wavelength-dependent and the proportion of light transmitted in the off state may be relatively high. With color versions of the newer active matrix LCDs it is possible to pass three wavelengths of monochromatic light through the color LCD and have a spatially controlled, high-intensity display with the maximum possible gamut. In this case the LCD is placed at a retinal conjugate plane.

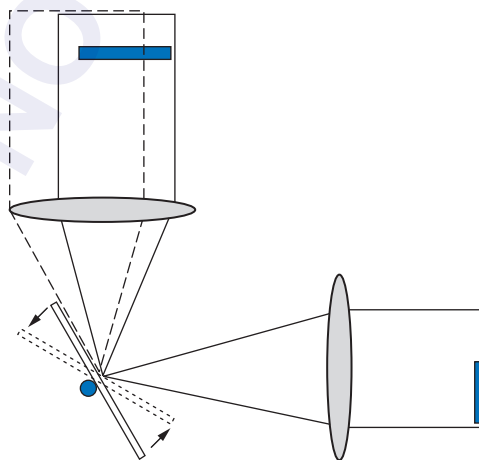


FIGURE 5 A movable mirror located in a pupil conjugate plane allows motion of the retinal image without moving the pupil conjugate images.

Calibration

There are several approaches to calibrating the retinal illuminance of a maxwellian view system. All depend on either a standard source of known luminance (a reference source) or a calibrated, linear^{31,32} photodetector. In the past, calibration typically depended on referring the illuminance of an optical system to a carefully maintained standard light or by measuring the luminance of a diffuser. Well-calibrated photometers and radiometers are now readily available, and we will cover techniques based on using these (the reader is referred to Westheimer³ for other approaches). We assume that all of the light arriving at the exit pupil of the instrument enters the eye and that the exit pupil can be adequately approximated as circular with a radius of r .

Measurement of Power at Exit Pupil With the first technique (Fig. 6a) a calibrated detector is placed at the exit pupil of the maxwellian view device. First, the total luminous power at the exit pupil is measured. Next, the retinal area illuminated is computed from the geometry of the stimulus [Eq. (12)]. From these two quantities the retinal illuminance can be computed. This calculation specifically assumes that all of the power is uniformly distributed across the retinal image. Conversion of retinal illuminance from lm/m^2 to trolands is achieved using Eq. (6).

Measurement of the Illuminance on a Detector In this technique the illuminance falling on a detector at a fixed distance from the exit pupil is measured (Fig. 6b; also see Ref. 33). A circular source of luminance L and radius r produces an illuminance E on a detector at distance d^2 . If we assume that the source dimensions affect the calibration negligibly ($r < d/10$; $r < f_e/10$), then

$$E = \frac{L\pi r^2}{(d^2)} \quad (18)$$

We want to relate E as measured by the detector to E_R , the retinal illuminance, which is

$$E_R = \frac{L\pi r^2}{(f_e)^2} \quad (19)$$

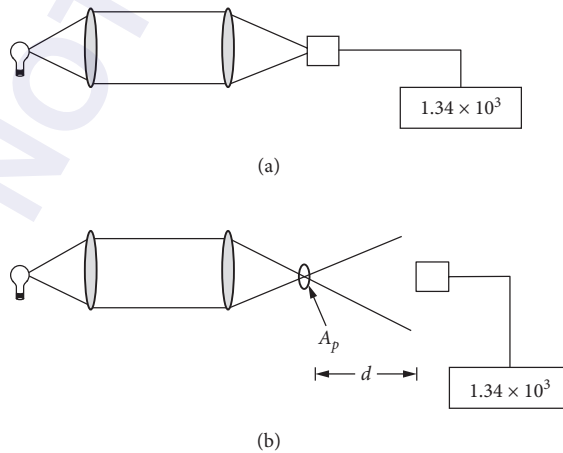


FIGURE 6 Techniques for calibrating an optical system. (a) Measure the luminous power of the exit pupil and assume that it is uniformly distributed in the retinal image. (b) Measure the illuminance produced by an exit pupil of radius r at a detector located at distance d from the exit pupil.

Combining Eqs. (18) and (19) we have

$$E_R = \frac{(d)^2}{(f_e)^2} E \quad (20)$$

If the radius of the exit pupil (r_p) cannot be ignored then we need to account for its properties as a light source. If the exit pupil is circular and can be approximated as a lambertian emitter, then the illumination on the detector [Eq. (18)] is

$$E = \frac{L\pi r^2}{(r_p^2 + d^2)} \quad (21)$$

(Ref. 2, Chap. 56, Table I(4.4)). Likewise, the retinal illuminance [Eq. (20)] is

$$E_R = \frac{(r_p^2 + d^2)}{(r_p^2 + (f_e)^2)} E \quad (22)$$

The troland value of E_R can be computed from Eq. (16).

5.7 LIGHT EXPOSURE AND OCULAR SAFETY

There are two main mechanisms by which light can damage the eye. The first is simply by heating: too much radiation burns the retina. The second mechanism is photochemical. Light, especially short-wavelength light, causes photochemical oxidation. The by-products of this oxidation are toxic to the retina (see Chap. 7).

Figure 7 shows the relation of the danger thresholds as defined by the ANSI 136 standard,³⁴ expressed in tds and lumens/m², for a 440-, a 550-, and a 670-nm light. These thresholds are for field sizes greater than 1 deg. The damage threshold for visible light is quite high, and thus lights capable of producing damage are intensely unpleasant. However, care should be taken to exclude IR and UV light. Most

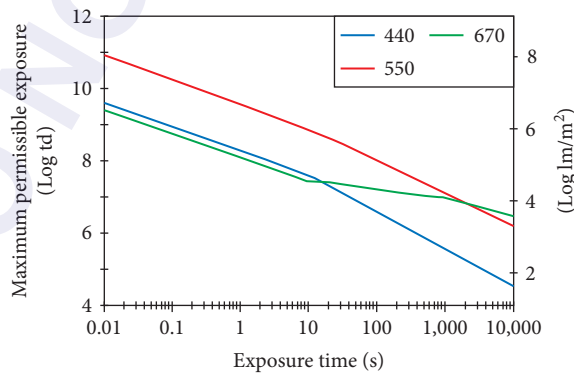


FIGURE 7 The maximum permissible retinal illuminance at different exposure times for 440 nm (blue line), 550 nm (red line), and 670 nm (green line) lights and extended fields according to the ANSI standards.

light sources emit considerable radiant energy in the infrared, and sources such as xenon arcs can emit ultraviolet radiation. Failures of blocking filters in these spectral regions may be visually undetectable. For safety it is advisable to have separate blocking filters to exclude unwanted radiation. It should also be noted that when using coherent light sources speckle can cause large focal variations in the retinal illuminance and the safety limits should be lowered by a factor of between 10 and 100 times.³⁵ Point sources and lasers have different safety standards that take into account the motion of the eye.

5.8 LIGHT SOURCES

Common light sources are presented in Table 4 and lasers are presented in Table 5 and in Chap. 15. Wyszecki and Stiles² present detailed discussions of light sources. The considerations for deciding on a light source are stability and uniformity of the luminous area, intrinsic luminosity of the source, and the spectral distribution of light from the source. Light sources of various types are discussed in Chaps. 15 to 23 in Vol. II.

5.9 COHERENT RADIATION

Coherent light sources can be used for generation of interference fringes,³⁶ a common use of lasers in vision research.⁸⁻¹⁰ Light from the laser is directed into two similar channels. The light from the two channels arrives at the pupil of the eye separated by some distance d . Light spreading from the two beams then overlaps on the retina and creates an interference pattern in the intersection zone. The spacing of the interference pattern is controlled by the angle at which the beams intersect at the retina. The orientation of the pattern is controlled by the relative orientation of the two pupils. However, HeNe and easily controlled solid-state visible lasers are also excellent general-purpose light sources. The major drawback to their use is speckle.

Generation of Speckle Fields

Speckle arises when coherent light is reflected from an optically rough surface. The surface roughness causes variations in the path length between the retina and adjacent areas on the surface. This variation causes phase differences from point to point in the retinal image. The spatial frequency of

TABLE 4 Common Light Sources

Type	Properties	Comments
Tungsten	Broadband, thermal emitter.	Available in a number of filament and power choices, hotter bulbs are short-lived.
Tungsten-halogen	Broadband, thermal emitter.	Inclusion of halogen regeneration cycle allows tungsten-halogen to run hotter (more blue light and a higher luminance) at a longer rated life than tungsten. If running at below the rated wattage the bulbs can be inverted to extend the lifetime.
Xenon arcs	Nearly flat spectrum, high luminance, some emission lines are present.	UV light can create ozone which must be vented. Tendency for the arc to move relative to the electrodes, causing movement of the exit pupil in maxwellian view systems. Luminous area is fairly small.
Fluorescent tube	Broad spectral output with superimposed emission lines. Easily obtainable, efficient.	Can be modulated rapidly. But there may be problems with compound phosphors (see Ref. 28).
Light-emitting diodes (LEDs)	Small, cheap, efficient light sources. Easily controlled.	Can be modulated rapidly, come in a variety of wavelength and power ranges (see Chap. 13).

TABLE 5 Types of Lasers

Type	Most Common Visible Wavelengths	Comments	Typical Power	Typical Noise %
Argon	514, 488, (other lines are available, depending upon the design). 444 nm	Power of each wavelength varies with particular laser design and cost. Most common wavelengths are 514 nm and 488 nm. Fairly expensive, have some high-frequency noise (>300 kHz) that may be important in scanning or short exposure applications.	5 mW–5 W	1
HeCd			0.5–50 mW	3
HeNe	543, 594, 632.8, also have orange lines.	Most common and least expensive gas laser. The 632.8-nm designs are easily available, inexpensive, and long-lived. Other wavelengths have many of the same advantages but power and beam options are limited and cost is higher.	0.1–20 mW	0.5
Krypton Dye lasers	588, 677. Variable.	Expensive, tend to be unstable, large. Expensive, need to be pumped by another laser. Require costly and toxic dyes to operate, not particularly stable.	50 mW–5 W 50 mW–5 W	3
Slab diode lasers	~630 nm, ~670 nm, ~690 nm, some in blue.	Development is very rapid with new wavelengths and power options. The beam geometry is not ideal but they have great potential as light sources. Can be easily controlled using standard electronic techniques as are common with LEDs.	0.2–100 mW	1 percent or better if temperature stabilized.
VCSELS vertical cavity lasers (microlasers)	Infrared.	Have high impedance and are easily damaged by static electricity. Rapid developments at shorter wavelengths. These are under rapid development. They are solid-state lasers but have superior efficiency and beam properties as compared to slab lasers. Low impedance.	1 mW	
Frequency doubled YAG	532 nm.	Readily available, solid-state designs. Expensive.	5–50 mW	<1 percent.
Tunable solid-state and solid dye lasers	Variable, pumped by lasers.	An area of rapid development. By using materials that fluoresce, and tuning the laser cavity, the output wavelength can be varied while maintaining excellent beam qualities.	mW–W	Depends on pump laser.

speckle depends on pupil size, since the size of those adjacent areas is smaller for bigger pupils. Thus, for a very small pupil, speckles will be large. Speckle has been used to develop optometers,³⁷⁻³⁹ and to generate a pattern on the retina in the presence of poor optic media.^{40,41} However, when using lasers as general-purpose light sources, speckle needs to be minimized or eliminated. There are three ways to despeckle a source: *spatial filtering*, *temporal despeckling*, and *raster scanning*.

In spatial filtering¹⁰ the light remains spatially coherent. A light source is imaged onto a pinhole which diffracts the light so that a spherical wave emerges from the pinhole. A lens is then used to collimate the light, resulting in a plane wave. Filtering should be restricted to those experiments that require the use of coherent light. The introduction of dust or even structure in the anterior segment of the eye can introduce undesirable diffraction effects.

Temporal despeckling^{14,15} uses the temporal integration of the visual system to blur the speckle field. Since speckle arises as a result of surface irregularities at a very small scale, a small amount of motion can decrease speckle contrast considerably.

Scanning moves a diffraction-limited spot of light across the retina to create a visual pattern. Typically, the pattern is a raster pattern, like that of a television, and the stimulus is generated by temporal variation in the intensity of the beam.⁴²⁻⁴⁴ Since at any one time only a single, diffraction-limited, retinal region is illuminated, then there is no opportunity for speckle.

5.10 DETECTORS

Detectors are used for measuring the light produced by an optical system. They can generally be characterized by their quantum sensitivity, their spectral response curve, and the temporal amplitude spectrum of their noise (see Ref. 45 and Chap. 24, "Photodetectors," by Paul R. Norton and Chap. 27, "Signal Detection and Analysis," by John R. Willison in Vol. II). For most purposes detectors can be treated as detecting all light power incident upon their active area, but there are limitations that make it advisable to use them for normally incident light whenever possible.⁴⁶ Table 6 presents the most common detectors (also see Chaps. 24 to 28 in Vol. II). At low light levels photomultiplier tubes (PMT)⁴⁷ or an avalanche photodiodes (APD)⁴⁸ can be used in a photon-counting mode (Chap. 26, "High-Speed Photodetectors," by John E. Bowers and Yih G. Wey in Vol. II).

5.11 PUTTING IT TOGETHER

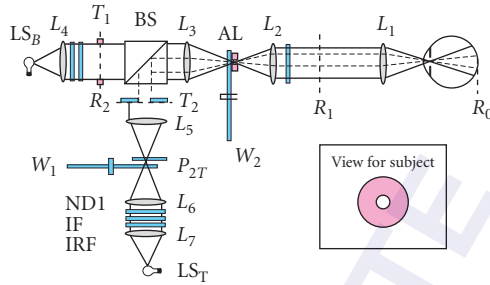
This section briefly describes construction of more complex maxwellian view systems. The goal is simply to help the researcher get started building an optical device by presenting real examples.

A Two-Channel Maxwellian View Apparatus

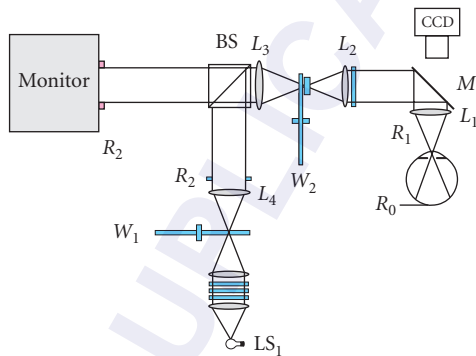
We first consider a simple two-channel maxwellian view device (Fig. 8a) for measuring the detectability of a small circular target on a large circular background. The background is provided via the straight optical channel. It has a light source LS_B followed by a lens L_4 that collimates the light from the source. A field stop (target T_1) controls the size of the background field. Adjacent to the target is a beamsplitter, followed by another lens L_3 . L_3 creates a source image where we place a 2.0-log-unit circular neutral density wedge, a 2-mm (diameter) stop (artificial pupil) and an achromatizing lens AL. An achromatizing lens minimizes the effects of the chromatic aberrations of the eye.⁴⁹⁻⁵¹ The artificial pupil will be imaged at 1:1 in the pupil of the eye, thus providing a limiting pupil. Lenses L_2 and L_1 relay the artificial pupil to the pupil of the eye. The test channel is derived from the second light source LS_T . Light is collimated by L_7 , passed through a heat-rejecting filter IRF, interference filter IF₁, and a neutral density filter ND₁. An image of LS_T is then formed at P_{2T} by L_6 . In this plane we set an electromechanical shutter and a neutral density wedge. We then collimate

TABLE 6 Detectors

Detector	Sensitivity	Speed	Problems/Limitations
Photomultipliers	Best in the blue, extended multi-alkali cathodes extend sensitivity into the red and near IR. Many designs have very low noise.	Fast, MHz, capable of photon counting.	Low sensitivity in the red and infrared, can be damaged by high light levels, though recovery can be aided by leaving in the dark for a long time (months). Fragile. Require high voltage (kV).
Avalanche photodiodes	Have the silicon sensitivity curve, peaking in the near IR. Higher quantum efficiency than PMTs.	Fast, MHz, capable of photon counting, better quantum efficiency than a PMT, but higher noise level.	Noise increases rapidly at avalanche voltage, light can affect the breakdown point. Higher noise level than PMTs, but the higher quantum sensitivity makes them better for video rate imaging. Small sensitive area.
CCDs, CIDs	Very linear, easy to use, sensitivity peaks in the near IR. These have a high quantum efficiency.	Variable, depending upon implementation, integrate between readings (kHz). With care they can be operated at high rates (MHz); however, in most photometers the circuitry is specifically designed for low-noise, lower-frequency operation.	Blooming and charge spread can affect spatial properties at high radiances. ⁵² CIDs are resistant to blooming.
PIN silicon photodiodes	Easy to use, sensitivity peaks in the near IR.	Very slow <10 Hz.	
Thermal (bolometers)	Measures energy, thus flat spectral responsivity.		The only real use for these devices in vision research is for calibration of the spectral sensitivity of another detector. In general, it is better to have the detector calibrated against an NIST standard by a manufacturer.
See also photoresistors, vacuum photodiodes, and phototransistors	Seldom used.		



(a) Two channel increment threshold apparatus



(b) Using a video display in a maxwellian view and monitoring eye position

FIGURE 8 (a) An example of a two-channel increment threshold maxwellian view device. Two light sources are used (LS_B and LS_T) to form a background and a test channel respectively. The two channels are combined using a beamsplitter (BS). The intensity of the test field is controlled by a circular neutral density wedge W_1 , and the intensity of both channels is controlled using wedge W_2 . Additional details are provided in the text. (b) A similar apparatus, except now a video monitor is superimposed on a background derived from source LS_1 . An infrared sensitive CCD camera is used to monitor the pupil through the back of a dichroic beamsplitter plate (M_1) which reflects visible light and transmits infrared light (cold mirror). Infrared illumination of the eye is provided by infrared LEDs (not shown).

the light with lens L_5 , pass it through an aperture T_2 that creates the target, and combine it with the background light at the beamsplitter. In some cases it may be desirable to obtain both the test and background light from a single light source. This is readily achieved by the addition of mirrors and avoids the possibility that slight fluctuations in the sources might produce variability in the ratio of the target and background illuminances. When appropriate, use a single light source for multiple optical channels.

A similar optical design allows full spatial control of the stimulus by incorporating a video monitor (Fig. 8b) or an LCD (not shown). In general, the retinal image of the monitor will need to be minified. As shown, the size of the retinal image of the monitor is set by the ratio of the focal length of lenses L_1 , L_2 , and L_3 [Eq. (14)]. The main advantages of the hybrid system over free-viewing the monitor are the ability to control the pupil precisely, the ability to add spectrally pure backgrounds (and, if only low contrasts are needed, the background can be quite bright), and the ability to align an achromatizing lens precisely and to monitor the eye position. We show a simple monitoring system, where the eye is monitored by a CCD camera, through the back side of a “cold mirror” (M1). A cold mirror is a dichroic beamsplitter that reflects visible light and transmits infrared radiation. IR-emitting LEDs are used to illuminate the eye diffusely. Using the full aperture of the lens L_1 means that the depth of focus will be small, and the eye can be precisely positioned (or monitored) in three dimensions. The resolution of the camera can be relatively low.

To use an LCD in transmission mode rather than a video monitor we can simply introduce the display at the retinal conjugate plane (R_2) in Fig. 8a. Again, it will be desirable to decrease the magnification of the LCD on the retina by choosing L_3 to have a longer focal length than lenses L_1 and L_2 .

An Interferometer

We next describe a research interferometer that was built using modern electro-optical components.¹⁰ To study the spatial sampling of the photoreceptor matrix of the eye, this instrument needed to have the following properties: (1) generation of spatial patterns on the retina with a frequency content above that imaged by the optics of the eye, (2) rapid and precise control of the modulation (contrast) and frequency (spacing) of the spatial patterns, (3) minimization of speckle, and (4) control of the two entrance pupils of the interferometer to displace them symmetrically about the center of the eye’s pupil. We outline below how these objectives were achieved. The researcher should refer to the original paper for a more detailed account.

Interference can be used to generate diffraction patterns on the surface of the retina. Two coherently related sources are imaged in the plane of the eye’s pupil (Fig. 9a). The spacing of the resulting diffraction pattern depends on the separation of the two sources in the eye’s pupil, and the orientation of the pattern depends on the relative position of the two sources. Williams¹⁰ (see also Refs. 8 and 9 for other examples) constructed a modified Mach-Zender interferometer⁵² (Fig. 9b) that satisfies the requirements outlined above.

Modulation of the interference pattern is controlled using AOMs. Each of the beams is either on or off, with a 50 percent duty cycle. If both beams are on at the same time (in phase), the interference pattern is seen at full contrast. If one is off when the other is on (counterphase), then there is no interference pattern. Intermediate phases of the two beams will cause intermediate contrasts of the interference pattern (Fig. 9c). Thus, both beams are always operated in the same state (avoiding luminance artifacts due to temporal nonlinearities), and all control of modulation is electronic. The use of an all-electronic design allows precise control of contrast without audible cues to the subject. Spatial masking from speckle is minimized by placing spatial filters in a pupil conjugate plane. The spacing of the interference pattern is controlled by a rotatable cube. Each beam traverses the cube in opposite directions (Fig. 9d). As the cube is rotated the beams, and thus the images in the eye’s pupil, are displaced symmetrically.

5.12 CONCLUSIONS

We have covered only some of the basics of stimulus generation. We urge the interested reader to look elsewhere in this *Handbook*. In addition, the following references are good general sources for this topic.

5.13 ACKNOWLEDGMENTS

Supported by NIH EYO4395 (Burns) and DOE DE-FG02-91ER61229 (Webb). We thank Francois Delori for aid in computing the triland values for the ANSI standards and David Williams for valuable editorial help.

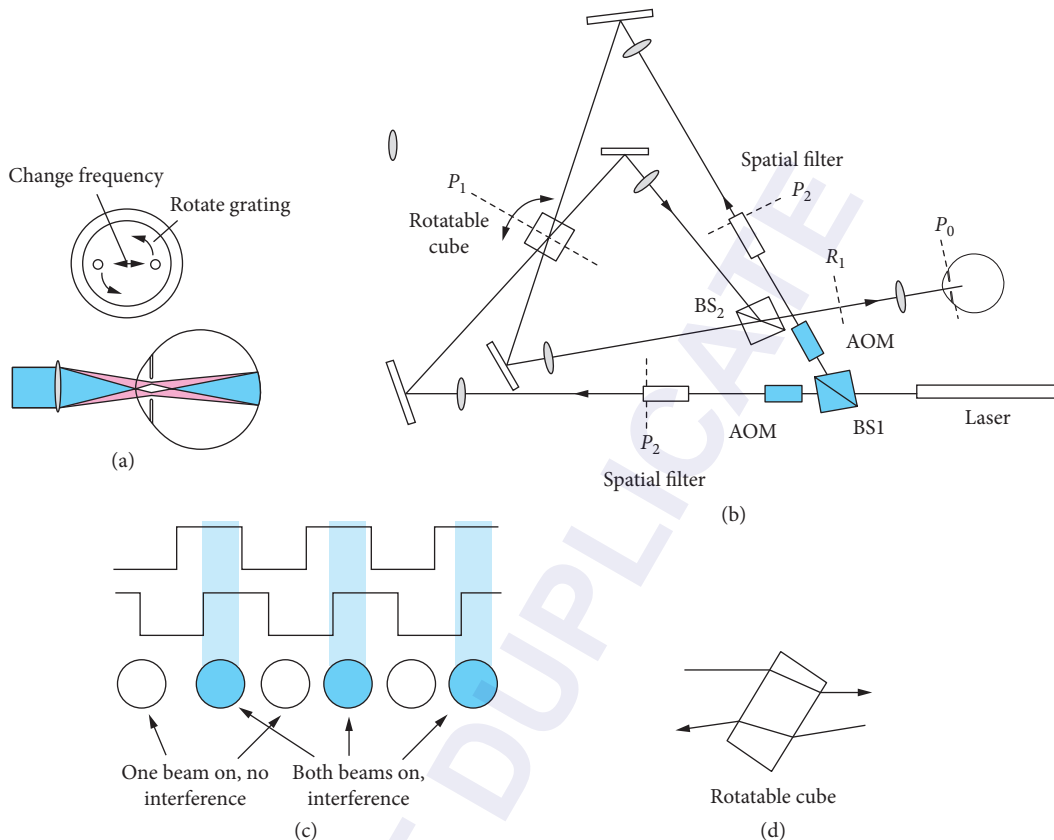


FIGURE 9 An example of a modern interferometer for vision research. (Modified from Ref. 10.) (a) Factors controlling the spacing and orientation of the interference pattern on the retina. Increasing the distance between the entry pupils decreases the spacing of the pattern (increases the spatial frequency in cycles per deg). It is also possible to change the orientation of the interference pattern by rotating the two entry pupil in the plane of the eye's pupil. (b) A schematic of William's interferometer. This is a modified Mach-Zender interferometer and includes acousto-optic modulators (AOM) to control the modulation of the interference pattern, spatial filters to reduce speckle, a rotatable cube to control spacing and beam splitters to separate the two beams (BS1) and to recombine them (BS2). (For additional details see Ref. 10.) (c) Contrast of the interference pattern is controlled by electronically varying the phase of the signals from the two AOMs. Each AOM is square wave modulated. The interference pattern is produced only when both beams are on at the same time. Thus, varying the relative phase of the two beams varies the proportion of the time when both are on, and thus the contrast. (d) The spacing of the beams in the plane of the eye's pupil is varied by rotating a deviation plate. The two beams propagate through the glass plate in opposite directions. Thus, as the plate is rotated, the beams are displaced equal but opposite amounts, resulting in symmetric deviations of the beams about the center of the eye's pupil. The actual apparatus uses an anti-reflection coated cube.

5.14 GENERAL REFERENCES

- J. R. Meyer-Arendt, "Radiometry and Photometry: Units and Conversion Factors," *Applied Optics* 7:2081–2084 (1968).
 G. Westheimer, "The Maxwellian View," *Vision Res.* 6:669–682 (1966).
 J. W. T. Walsh, *Photometry*, Constable, London, 1953.

- G. Wyszecki and W. S. Stiles, *Color Science*, Wiley, New York, 1982.
- R. M. Boynton, "Vision," J. B. Sidowski (ed.), *Experimental Methods and Instrumentation in Psychology*, McGraw-Hill, New York, 1966.
- W. J. Smith, *Modern Optical Engineering*, 2d ed., McGraw-Hill, New York, 1990.

5.15 REFERENCES

1. Y. LeGrand, *Light Color and Vision*, Wiley, New York, 1957.
2. G. Wyszecki and W. S. Stiles, *Color Science*, Wiley, New York, 1982.
3. G. Westheimer, "The Maxwellian View," *Vis. Res.* **6**:669–682 (1966).
4. F. W. Fitzke, A. L. Holden, and F. H. Sheen, "A Maxwellian-View Optometer Suitable for Electrophysiological and Psychophysical Research," *Vis. Res.* **25**:871–874 (1985).
5. R. J. Jacobs, I. L. Bailey, and M. A. Bullimore, "Artificial Pupils and Maxwellian View," *Appl. Opt.* **31**:3668–3677 (1992).
6. G. M. Byram, "The Physical and Photochemical Basis of Visual Resolving Power Part I: The Distribution of Illumination in Retinal Images," *J. Opt. Soc. Amer.* **34**:571–591 (1944).
7. O. Dupuy, "La perception visuelle," *Vis. Res.* **8**:1507–1520 (1968).
8. F. W. Campbell and D. G. Green, "Optical and Retinal Factors Affecting Visual Resolution," *J. Physiol.* **181**:576–593 (1965).
9. F. W. Campbell and R. W. Gubisch, "Optical Quality of the Human Eye," *J. Physiol.* **186**:558–578 (1966).
10. D. R. Williams, "Aliasing in Human Foveal Vision," *Vis. Res.* **25**:195–205 (1985).
11. W. D. Wright, *Researches into Normal and Defective Colour Vision*, Kimpton, London, 1946.
12. T. W. Liebman, "Wedge Plate Beam Splitter without Ghosts Reflections," *Appl. Opt.* **31**:5905–5906 (1992).
13. R. M. Boynton and A. L. Nagy, "La Jolla Analytic Colorimeter," *J. Opt. Soc. Amer.* **72**:666–667 (1982).
14. R. H. Webb, "Concentrator for Laser Light," *Appl. Opt.* **31**:5917–5918 (1992).
15. J. R. Krauskopf, D. R. Williams, and D. H. Heeley, "Computer Controlled Color Mixer with Laser Primaries," *Vis. Res.* **21**:951–953 (1981).
16. S. A. Burns, M. R. Kreitz, and A. E. Eisner, "Apparatus Note: A Computer Controlled, Two Color, Laser-Based Optical Stimulator for Vision Research," *Appl. Opt.* **30**:2063–2065 (1991).
17. C. A. Taylor and B. J. Thompson, "Some Improvements in the Operation of the Optical Diffractometer," *J. Sci. Instr.* **34**:439–447 (1957).
18. R. M. Boynton, M. M. Hayhoe, and D. I. A. MacLeod, "The Gap Effect: Chromatic and Achromatic Visual Discrimination as Affected by Field Separation," *Optica Acta* **24**:159–177 (1977).
19. M. C. Rynders and L. N. Thibos, "Single Channel Sinusoidally Modulated Visual Signal Generator, with Variable Temporal Contrast," *J. Opt. Soc. Amer.* **A10**:1642–1650 (1993).
20. V. V. Toi and P. A. Gronauer, "Visual Stimulator," *Rev. Sci. Instr.* **49**:1403–1406 (1978).
21. G. A. Fry, "Square Wave Grating Convolved with a Gaussian Spread Function," *J. Opt. Soc. Amer.* **58**:1415–1416 (1968).
22. T. W. Butler and L. A. Riggs, "Color Differences Scaled by Chromatic Modulation Sensitivity Functions," *Vis. Res.* **18**:1407–1416 (1978).
23. D. Vincent, "Amplitude Modulation with a Mechanical Chopper," *Appl. Opt.* **25**:1035–1036 (1986).
24. T. E. Cohn, "A Fast, Accurate Monochromatic Light Stimulus Generator," *Am. J. Optom. and Arch. Am. Acad. Optom.* **49**:1028–1030 (1972).
25. R. W. Nygaard and T. E. Frumkes, "LEDs: Convenient, Inexpensive Sources for Visual Experimentation," *Vis. Res.* **22**:435–440 (1982).
26. M. Yamashita and S. Takeuchi, "Temperature-Compensated Pulsed Reference Light Source Using a LED," *Review of Sci. Instr.* **54**:1795–1796 (1983).

27. W. H. Swanson, T. Ueno, V. C. Smith, and J. Pokorny, "Temporal Modulation Sensitivity and Pulse-Detection Thresholds for Chromatic and Luminance Perturbations," *J. Opt. Soc. Amer.* **A4**:1992–2005 (1987).
28. G. R. Cole, C. F. Stromeyer, III, and R. E. Kronauer, "Visual Interaction with Luminance and Chromatic Stimuli," *J. Opt. Soc. Amer.* **7**:128–140 (1990).
29. T. U. Watanabe, N. Mori, and F. Nakamura, "Technical Note: A New Suberbright LED Stimulator: Photodiode Feedback Design for Linearizing and Stabilizing Emitted Light," *Vis. Res.*, 1992.
30. J. D. Mollen and P. G. Polden, "On the Time Constant of Tachistoscopes," *Quarterly J. Exp. Psychol.* **30**:555–568 (1978).
31. C. L. Sanders, "Accurate Measurements of and Corrections for Nonlinearities in Radiometers," *J. Res. Natl. Bur. Stand., sec. A*, **76**:437–453 (1972).
32. R. G. Frehlich, "Estimation of the Nonlinearity of a Photodetector," *Appl. Opt.* **31**:5927–5929 (1992).
33. R. W. Nygaard and T. E. Frumkes, "Calibration of the Retinal Illuminance Provided by Maxwellian Views," *Vision Res.* **22**:433–434 (1982).
34. American National Standards Institute, "Safe Use of Lasers," Z-136.1 ANSI, N.Y., 1976.
35. D. L. Fried, "Laser Eye Safety: The Implication of Ordinary Speckle Statistics of Speckle-Speckle Statistics," *J. Opt. Soc. Amer.* **71**:914–916 (1981).
36. M. Francon, *Optical Interferometry*, Academic Press, N.Y., 1966.
37. R. T. Hennessy and H. W. Leibowitz, "Laser Optometer Incorporating the Badal Principle," *Behav. Res. Meth. Instr.* **4**:237–239 (1972).
38. N. W. Charman, "On the Position of the Plane of Stationarity in Laser Refraction," *Am. J. Optom. Physiol. Opt.* **51**:832–838 (1974).
39. A. Morrell and W. N. Charman, "A Bichromatic Laser Optometer," *Am. J. Optom. Physiol. Opt.* **64**:790–795 (1987).
40. W. W. Dawson and M. C. Barris, "Cortical Responses Evoked by Laser Speckle," *Invest. Ophthalmol. Vis. Sci.* **17**:1207–1212 (1978).
41. J. Fukuhara, H. Uozotot, S. Nojima, M. Saishin, and S. Nakao, "Visual-Evoked Potentials Elicited by Laser Speckle Patterns," *Invest. Ophthalmol. Vis. Sci.* **24**:1400–1407 (1983).
42. R. H. Webb, G. W. Hughes, and F. C. Delori, "Confocal Scanning Laser Ophthalmoscope," *Appl. Opt.* **26**:1492–1499 (1987).
43. R. H. Webb, "Optics for Laser Rasters," *Appl. Opt.* **23**:3680–3683 (1984).
44. A. E. Eisner, S. A. Burns, R. W. Webb, and G. H. Hughes, "Reflectometry with a Scanning Laser Ophthalmoscope," *Applied Optics* **31**:3697–3710 (1992).
45. R. H. Webb and G. H. Hughes, "Detectors for Video Rate Scanning Imagers," *Appl. Optics* **32**:6227–6235 (1993).
46. J. Durnin, C. Reece, and L. Mandel, "Does a Photodetector Always Measure the Rate of Arrival of Photons?" *J. Opt. Soc. Amer.* **71**:115–117 (1981).
47. D. van Norren and J. van der Kraats, "A Continuously Recording Retinal Densitometer," *Vision Res.* **21**:897–905 (1981).
48. P. P. Webb, R. J. McIntyre, and J. Conradi, "Properties of Avalanche Photodiodes," *RCA Review* **35**:234–278 (1974).
49. R. E. Bedford and G. Wyszecki, "Axial Chromatic Aberration of the Eye," *J. Opt. Soc. Amer.* **47**:564–565 (1957).
50. A. L. Lewis, M. Katz, and C. Oehrlein, "A Modified Achromatizing Lens," *Am. J. Optom. Physiol. Opt.* **59**:909–911 (1982).
51. I. Powell, "Lenses for Correcting Chromatic Aberration of the Eye," *Appl. Opt.* **20**:4152–4155 (1981).
52. D. Malacara, *Optical Shop Testing*, chap. 4.5, Wiley Interscience, New York, 1991.
53. M. Marchywka and D. G. Socker, "Modulation Transfer Function Measurement Technique for Small-Pixel Detectors," *Appl. Opt.* **31**:7198–7213 (1992).

This page intentionally left blank.

DO NOT DUPLICATE

THE MAXWELLIAN VIEW*† WITH AN ADDENDUM ON APODIZATION

Gerald Westheimer

*Division of Neurobiology
University of California
Berkeley, California*

Abstract—Optical questions arising in the so-called Maxwellian View method of illuminating the retina have been analyzed theoretically. Problems discussed in detail include those of photometry, of magnification, of focus, and finally of pupil size insofar as it relates to the transmission of spatial frequencies in coherent and incoherent illumination.

6.1 GLOSSARY

Diffraction theory. Formulation for calculating the distribution of electromagnetic energy when the propagation of radiation is restricted by apertures.

Emmetropia. Refractive state of eye in which the retina is conjugate to infinity. Its absence is called ametropia.

Fourier theory of optics. Object-image relationship in which the basic object structure is a sinusoidal grating and all other targets are described in their Fourier terms.

In the Fourier Theory of Optics, the complex amplitude of the electromagnetic vibrations in the plane of focus is the Fourier transform of that in the pupil plane.

Incoherent light sources. Sources for which light originating from different elements, when overlapping in forming the image, sums in intensity rather than by interference of the amplitudes of the electromagnetic vibrations.

Pupil of eye. Aperture stop, limiting the width of bundle of light entering the eye. The entrance pupil is the virtual aperture stop situated in the eye's object space; it can be directly measured and related to objects without further consideration of the eye's refracting apparatus.

Transfer function. Ratio of amplitudes of images to their generating objects for targets whose intensity changes sinusoidally, as a function of spatial frequency.

*This chapter is the full text of an article on the Maxwellian view published in 1966 in the journal *Vision Research* (vol. 6, p. 669–682, 1966). This article is republished with permission from Pergamon Press. Since this is a classic article, it is reprinted in full here. Prof. Westheimer has added a post-script to this article.

†Aided in part by a Grant NB-03154 from the National Institute of Neurological Diseases and Blindness, U.S. P.H.S. and a contract between the Office of Naval Research and the University of California.

6.2 INTRODUCTION

In his experiments on color mixing, James Clerk Maxwell (1968) devised a way by which he could increase the quantity of monochromatic light reaching the retina of his eye. Using the sky as his source, he produced a spectrum by means of a prism and then isolated narrow regions of this spectrum with slits. He thus had monochromatic sources, but the special virtue of his approach is that he then imaged each slit on the pupil of his eye by a lens and saw the lens “uniformly illuminated with light.” The procedure of imaging a light source in the eye’s pupil, instead of looking at it directly, has since been applied widely and is now known as Maxwellian viewing. Its simplest form is illustrated schematically in Fig. 1, and compared there with ordinary viewing.

In the following discussion certain simplifying assumptions are used. (a) The eye is taken to be perfectly transparent; if the transmissivity of the ocular media is known, this can be multiplied into the final results. (b) The Stiles-Crawford effect is not taken into account. (c) No distinction is made between the plane of the eye’s entrance pupil and its anterior principal plane—the error introduced is relatively unimportant here. (d) The eye is assumed to be accommodated for the object plane which is not necessarily at infinity. (e) Angles are regarded as small enough that the approximation $\cos \theta = 1$ is adequate.

In Fig. 1(a) and Fig. 1(b), the primary source of light is a small self-luminous patch of size $S(\text{m}^2)$ and luminance $B(\text{cd}/\text{m}^2)$; it is imaged directly on the retina in the case of Fig. 1(b), ordinary viewing, and it is imaged by a lens of aperture $L(\text{m}^2)$ in the pupil of the eye, in the case of Maxwellian viewing, Fig. 1(a).

Concise discussions of the Maxwellian view can be found in Legrand (Vol. II) and in Walsh (1958). The present paper deals with several aspects of using the Maxwellian view in vision experiments.

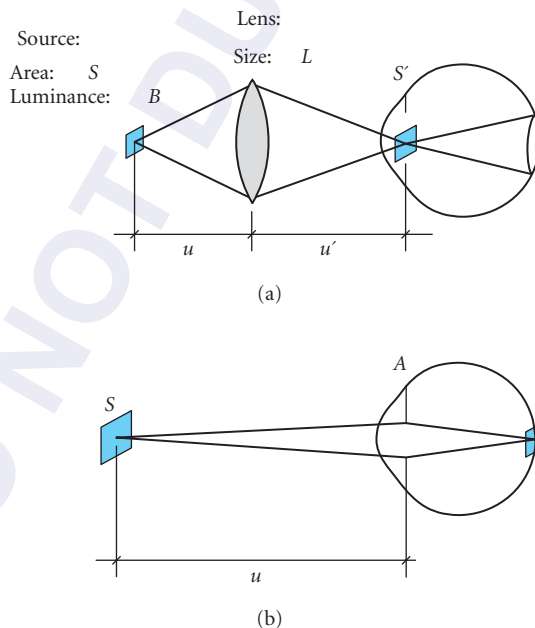


FIGURE 1 (a) Schematic diagram of the principles of Maxwellian viewing. A source of area S is imaged by a lens in the pupil of the eye. The image of the source has area $S' = Su^2/u'^2$. (b) Ordinary viewing of an extended source of area S placed at a distance from an eye whose entrance pupil has area A .

Retinal Illuminance

An immediate concern in any arrangement in which visual stimuli are presented, is the quantity of light reaching the subject's retina. The retinal illuminance should really be measured in terms of the luminous flux per unit area of retina. But there is the difficulty that, while we may be able to specify the light reaching the eye, additional assumptions are necessary to specify the retinal illuminance. The assumptions concern such things as the transmissivity of the ocular media and the size relationship between objects in front of the eye and their images on the retina. To obviate this, an approach is here adopted which retains, for as long as possible, all quantities in realistic units measurable in the eye's object space. For example, the size of the stimulated retinal area is expressed in terms of the object-sided visual angle, rather than square millimeters of retina. Knowledge of the dimension of the structures of the eye permits conversion of one set of units to others, but the advantage claimed for the present approach is that its validity is independent of the constants of a particular eye and it may be applied readily, therefore, to non-standard conditions, and to the eyes of other species.

Take the case of ordinary viewing first (Fig. 1*b*). A luminous surface [area $S(\text{m}^2)$, luminance $B(\text{cd}/\text{m}^2)$] is observed by an eye with an entrance pupil [area $A(\text{m}^2)$], directly facing the surface and at a distance $u(\text{m})$ from it. The total luminous flux entering the eye is equal to the product of the flux emitted by the source per unit solid angle (BS) and the solid angle subtended by the eye's entrance pupil at the source (A/u^2). This flux ($B.S.A./u^2$) covers a retinal region represented in the eye's object space by the solid angle subtended by the source at the center of the entrance pupil (S/u^2). The measure of the retinal illuminance which retains all quantities in the eye's object space is luminous flux per unit solid visual angle and it is given by

$$E (\text{lm/sterad}) = B \cdot A \quad (1a)$$

where B is the source's luminance in $\text{cd}/\text{unit area}$ and A is the eye's pupillary area. It is to be noted that the retinal illuminance is independent of the size of the source and the observation distance.

The most widely used unit of retinal illuminance, the troland, is defined by the condition in which a surface of luminance $1 \text{ cd}/\text{m}^2$ is observed by an eye with an entrance pupil area of 1 mm^2 . A troland is thus equal to $10^{-6} \text{ lm/sterad visual angle}$. A steradian is not a convenient unit of solid visual angle. We are usually interested in much smaller areas, often of the order of magnitude of single retinal receptors. Thus, it may sometimes be better to give the luminous flux not in each steradian of visual angle, but in each square minute of arc of which there are $(60 \times 57.3)^2 = 11.8 \times 10^6$ per steradian. A useful unit in which to express the rate at which light is impinging on the eye then is luminous flux per square minute of arc visual angle for each square millimeter of entrance pupil area. When this pupillary area is $A \text{ mm}^2$ and when the source has luminance $B \text{ cd}/\text{m}^2$, we have

$$E (\text{trolands}) = B \cdot A \quad (1b)$$

or

$$E (\text{lm/sterad}) = 10^{-6} B \cdot A \quad (1c)$$

or

$$E (\text{lm/sq. min of arc}) = \frac{B \cdot A 10^{-6}}{11.8 \times 10^{-6}} = 0.85 \times 10^{-13} B \cdot A \quad (1d)$$

In the human eye, cones in the center of the fovea, and rods in an optimal region for scotopic vision (about 10° peripheral) are packed about 4 to a sq. min of arc, and thus the next kind of derived unit would be

$$E (\text{lm/receptor}) = 0.21 \times 10^{-13} B \cdot A \quad (1e)$$

Each scotopic lumen represents a flux of 1.46×10^{15} quanta (507 m μ) per second (Aguilar and Stiles, 1954) and hence we have yet another derived unit

$$E \text{ (quanta (507)/receptor} \times \text{sec)} = 30 B \cdot A \quad (1f)$$

Note that in equations (1 b-f), B is the luminance of the extended source measured in cd/m², and A is the pupil area in square millimeters.

All the above units are still in the eye's object space, and before the number of quanta that are absorbed in each receptor per second can be computed, it is necessary to know the proportion of externally incident quanta that is transmitted by the ocular media, and the proportion of the latter that is absorbed by the visual pigment. Rushton (1956b) has given an estimate of 0.1 for these two factors combined for human rod vision.

Consider the case of Maxwellian viewing, i.e. in which the source is not imaged on the retina but in the eye's entrance pupil. We may distinguish two important categories, one in which the image of the source (area $S'' = S \cdot u^2/u^2$) is larger than the pupil area (A) and the other in which it is smaller. In the former, it is A that limits the light entering the eye. The retinal area illuminated in each case is given, in steradians of visual angle, by (L/u^2) where L is the area (m²) of the clear aperture of the field lens, and u is the distance between this lens and the eye. The total luminous flux reaching the image S'' is $B \cdot S \cdot L/u^2$. If the pupil area A is smaller than S'' , the proportion of this total flux that can enter the eye is A/S'' .

It is understood that when the pupil size remains constant, the retinal illuminance when viewing a Lambert surface is independent of observation distance. This is so because the angle subtended by an element of area increases when it is brought closer, but so does the angle the pupil subtends at the surface, and hence the fraction of the total radiation emitted by the source. The two effects exactly counterbalance each other.

A similar kind of invariance obtains when using an extended source in a Maxwellian view. So long as the image of the source is larger than the pupil, it is immaterial with what magnification the source is imaged in the pupil plane. The closer the source comes to the lens (i.e. the smaller u) the larger the flux per unit area there. Yet at the same time, for a portion of lens surface of given size to subtend a fixed angle at the eye, the distance between the lens and the eye has to be controlled and this in turn controls the lens power and the magnification of the image of the source in the pupil. The result is that the flux per unit angular size of lens surface per unit pupil area always remains constant, so long as all the available pupil area is occupied by the image of the source.

The actual flux entering the eye is either $S \cdot B \cdot L \cdot u^2$ for $S' < A$ and $S \cdot B \cdot L \cdot A/S' u^2$ for $S' > A$. In the latter case, i.e. when the source image is larger than the pupil, the flux is equal to $B \cdot L \cdot A/(u')^2$ which is spread over an angular area $L (u'^2)$: the flux incident per unit solid visual angle is, therefore,

$$E \text{ (lm/sterad)} = B \cdot A$$

which is the same as that in the case of ordinary viewing.

The principal advantage of the Maxwellian view technique is here brought out; the retinal illuminance attainable is at most equal to that of ordinary viewing of the light source, but a larger area of the retina may be illuminated at that level. This is apparent in Fig. 1. When the pupil area is smaller than the image of the source at the pupil, the retinal illuminance, i.e. flux per unit retinal area, is the same in the two cases, only in the case of Maxwellian viewing the visual angle is equal to L/u^2 , while otherwise it is S/u^2 .

When the image of the source is smaller than the pupil, the retinal illuminance will be smaller in proportion to the fraction of the pupil covered by the image of the source, i.e.

$$E \text{ (lm/sterad)} = B \cdot A \quad S'/A = B \cdot S' \quad (2)$$

The determination of the retinal illuminance in an experimental arrangement employing Maxwellian viewing follows from the above consideration. When the effective entrance pupil of the eye (real or artificial) is completely covered by the image of the source, the retina is illuminated at the same flux/unit area as if the eye with that particular pupil were looking directly at the source—allowance having to be made for transmission losses in the optical system between the source and the eye.

When the image of the source is smaller than the pupil, an indirect method of determining the retinal illuminance may be employed, predicated by the fact that all the light coming towards the eye out of the Maxwellian view optical system enters the eye's pupil. The retinal illuminance in this case, by Eq. (2), is $B \cdot S'$ lm/sterad target area. Figure 1(a) shows how this retinal illuminance is in fact achieved by having the equivalent of a source emitting $B \cdot S'$ lm/sterad, i.e. a source of $B \cdot S'$ candelas, imaged in the pupil. A good method* of measuring the resulting retinal illuminance is to find the number of candelas in the areal image of the source by measuring the luminance it produces when illuminating a perfectly diffusing surface placed at a given distance from it. One candela placed 1 m in front of a perfectly reflecting and diffusing surface gives it a luminance of 0.1 mL. When a small source placed at a distance x meters in front of a perfectly diffusing surface of reflectance r produces a luminance of \bar{B} mL, the source has intensity $10 \bar{B}x^2/r$ cd, and by equation (1b) such a source image in the plane of the pupil would give a retinal illuminance of $10 \bar{B}x^2/r$ lm/sterad or $10^7 \bar{B}x^2/r$ trolands.

The procedure for measuring the retinal illuminance in a Maxwellian view system where all the radiation passing through the instrument enters the pupil is:

1. Place a perfectly diffusing surface of reflectivity r at a distance x m beyond the source image that would ordinarily be in the plane of eye's pupil.
2. Measure the luminance of the surface with one of the standard instruments. Let it be \bar{B} mL.
3. Calculate the retinal illuminance using the formula

$$E \text{ (trolands)} = 10^7 \bar{B}x^2/r \quad (3)$$

Focus

Our discussion has so far concerned itself only with the retinal illuminance achievable with a Maxwellian view arrangement. Unless a uniform field of view is desired, a target is placed in the beam. In the typical experiment it is intended to stimulate the subject's retina in a certain way, for example, by illuminating certain areas and not others. The almost universally adopted way of creating a desired pattern of light stimulation on the retina, is to place in the object plane, conjugate to the retina, a target which is then imaged on the retina by the optics of the eye. (Rare exceptions to this approach may occur, for example, when setting up interference fringes on the retina, Westheimer, 1960.) In the arrangement of Fig. 1(a), a reasonable place for the target is just in front of or behind the lens L . Unfortunately, the resulting retinal image will only be clear when the eye is accommodated on L . This restriction is no longer necessary, however, if we re-arrange the optical system by splitting the lens L into two components, L_1 and L_2 . In Fig. 1(a) it is seen that the function of L is to image the source in the pupil. Suppose this is done in two stages, first by using a lens, L_1 , to produce a parallel beam of the light from the source, and then by having another lens, L_2 , image this parallel beam in the pupil (Fig. 2). Because the beam in its passage between L_1 and L_2 is parallel, these two lenses may be separated by any desired distance without altering their effect of imaging the source in the pupil. But the target position can now be chosen so that the subject's accommodative state may be anything within a wide limit. For example, if the subject is emmetropic and is not to accommodate, the target would be placed at the first focal point of lens L_2 , which then images it at infinity, and this is its situation relative to the subject's eye. If, on the other hand, the subject is 1 diopter myopic (or is emmetropic and is to accommodate 1 diopter) the target's place with respect to lens L_2 must be such that this lens images it at a distance of 1 m in front of the eye. Suppose the lens has focal length f_2 . Let x be the distance of the object from the first focal point of lens L_2 , and x' be the distance of the image from the second focal point of lens L_2 (which in our present arrangement is the eye's entrance pupil); then by Newton's lens formula we have the relationship

$$xx' = -f_2^2 \quad \text{or} \quad 1/x' = -x/f_2^2$$

*The author was introduced to this technique by W. A. H. Rushton, who had originally described in the literature (Rushton, 1956a) a related approach.

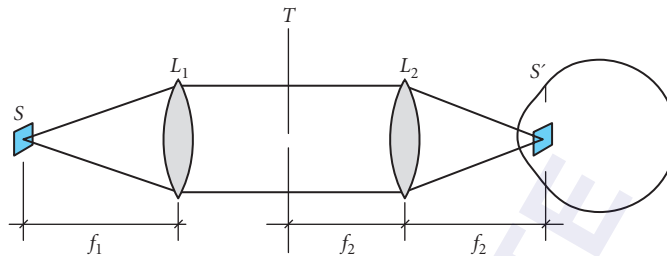


FIGURE 2 More versatile arrangement for Maxwellian view. The source S is placed at the first principal focal plane of lens L_1 and the target T is now trans-illuminated by a parallel beam. A laser may be substituted for S and L_1 . For the image of the source to fall in the plane of the eye's entrance pupil, the latter should be in the second principal focal plane of lens L_2 . If the observer is emmetropic, the target T should be placed at the first principal focal plane of lens L_2 .

Now $1/x'$ is the reciprocal of the image distance from the eye, i.e., it is the desired accommodative state of an emmetropic eye (or the far point of an unaccommodated ametropic eye) in diopters. It follows that a linear relationship exists between x the displacement of the target from the first primary focal point of lens L_2 , and the target focus with respect to the eye in diopters. Each diopter of target focus corresponds to a distance $x = f_2^2 (x, f_2 \text{ in meters})$. For example, if the lens L_2 has a focal length of 0.20 m, each diopter of target focus as far as the eye is concerned is represented by 0.04 m displacement of the target from the first focal point of the lens L_2 . The displacement is towards L_2 if the subject is accommodating, or is myopic, and is away from L_2 if he is hyperopic. The device clearly makes a good optometer, and was described in an arrangement for this purpose by Schmidt-Rimpler in 1877 (Emsley, 1939).

Field of View

Another important variable in setting up a vision experiment, in addition to the level of retinal illuminance and the state of focus, which were discussed above, is the area of retina illuminated. We shall follow the approach adopted throughout this paper of specifying retinal distances in angular measure in the eye's object space. The angle should be measured at the center of the entrance pupil, which is about 3 mm behind the corneal vertex; the specification of retinal distance is then more invariant with such factors as accommodation and out-of-focusness than if another reference point, e.g. nodal point, is used. For a typical human eye, 1 min of arc in the object space corresponds to a retinal distance of just under 5 μ , so that there arc about 3.4° in each millimeter of retinal distance.

If the target is placed in the plane of lens L_2 in Fig. 2, and if its diameter is c , its angular subtense is approximately c/f_2 rad. Here, c and f_2 should be measured in the same units. To convert to minutes of arc, it should be remembered that there are 3.44×10^3 min of arc in each radian. In a hypothetical situation where $f_2 = 20$ cm, and it is desired to present the eye with a stimulus 60 min of arc in diameter, its linear dimension must then be $c = 20 \times 60/3.44 \times 10^3$ cm = 0.3 cm. Clearly, if small targets have to be presented to the eye, the focal length of the lens L_2 should be long so that the linear size of the targets does not become too small. For example, if a lens of focal length 34.4 cm is chosen, each minute of arc of target size is 0.1 mm.

The above arguments hold identically if the target is not in the plane of L_2 but is placed anywhere in the parallel beam between lenses L_1 and L_2 . The important property of the principle of placing the

second focal point of the lens L_2 in the plane of eye's entrance pupil, the so-called telecentric principle, is now apparent: for purposes of changing focus (Section "Focus") the target may be moved in the space between the two lenses L_1 and L_2 , but such a move does not change either the retinal illuminance or the image size.

The maximum image size that can be achieved in the arrangement of Fig. 2 is one in which the apertures of lenses L_1 and L_2 remain unobscured. If the aperture of the smaller of the two lenses is c , the maximum field size is c/f_2 rad. If, for example, lens L_2 is filled and has an aperture of 50 mm and focal length 20 cm, the largest field is $5/20 = 1/4$ rad. $\approx 14^\circ$ in diameter.

Pupil Size

The effect of pupil size on the quantity of light entering the eye has already been discussed. The Stiles-Crawford effect has an influence, of course; whether it needs to be taken into consideration is a decision that depends on the particular application since, for example, the *S-C* effect does not show up in rod vision.

The pupil aperture also affects the quality of the retinal image. While in an important way everything to be said here will be modified by the optical aberrations of the particular eye to which it is applied, there still remains a substantial body of generally valid facts based on the diffraction theory of perfect optical instruments. Only under exceptional circumstances can an eye be expected to be as good as this, but our considerations allow us to set a bound to the best performance of any eye. In the present study no attempt will be made to go beyond this approach.

Incoherent Target* It used to be taught that the resolution limit of an eye with an entrance pupil diameter of a mm for light of wave length λ is $1.22 \lambda/a$ rad of visual angle in object space. This is at best an incomplete description of the situation, based on the fact that the diffraction-limited point-spread function, the Fraunhofer diffraction pattern, then has a radius (to its first zero) of $1.22 \lambda/a$ rad of visual angle. A fuller analysis of the performance limitation of the eye is afforded by Fourier theory. The fundamental description of the eye's imagery, which the above statement expresses in the space domain (a point object is imaged in the Fraunhofer diffraction spread pattern), may also be phrased in terms of the spatial frequency domain. A grating with sinusoidal intensity pattern with a period α min of arc visual angle may also be said to have repetition rate or spatial frequency of $1/\alpha$ c/min of arc. In the spatial frequency domain, the way of describing the performance limitation of an optical system is to state the demodulation experienced by a sinusoidal grating target in the process of being imaged. When this quantity is plotted against spatial frequency we have what is now called the modulation transfer function. Such curves for a hypothetical aberration-free eye with various entrance pupils are shown in Fig. 3. Their main feature is a progressive reduction in the modulation of the image (when the object grating retains constant contrast) as the spatial frequency is increased up to the cut-off spatial frequency, beyond which there is no indication in the image that the object contains any features of those repetition frequencies.

The above statement is based on the acceptable premise that optical imagery obeys the principles of linearity. When the target is self-luminous, as it is assumed to be the case with ordinary viewing, this is indeed the case. The following is a summary in mathematical language of the basic concepts of image formation for an incoherent target.

Let the two dimensions in the eye's object space be α, β ; distances involving α refer to the horizontal direction, and β the vertical direction. Let distances in the two corresponding directions of the eye's entrance pupil be y_α, y_β . Let spatial frequencies in these two directions be denoted by $\omega_\alpha, \omega_\beta$. For example a grating made up of vertical bars with period 10 min of arc has spatial frequency $\omega_\alpha = 0.1$ c/min of arc.

* An excellent introduction to the following discussions can be found in O'Neill (1963) or in Hopkins (1953).

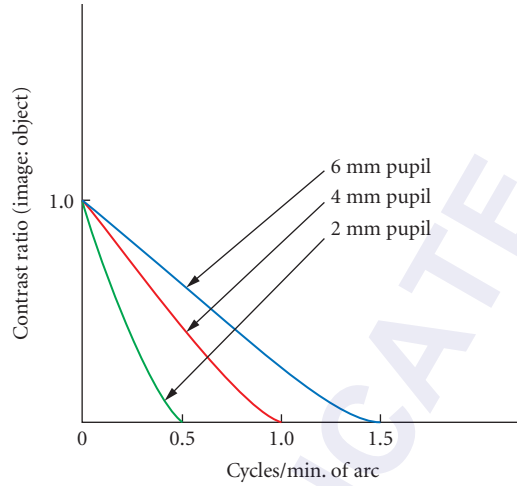


FIGURE 3 Contrast transfer functions (ratios of contrast in images and corresponding objects for sinusoidal light distributions) of optical system like that of the human eye but free from aberrations with round pupil of various diameters in monochromatic light of wavelength $\lambda = 555 \text{ m}\mu$.

Let the pupil aperture be described by the function $g(y_\alpha, y_\beta)$. For example, if the pupil is rectangular with a horizontal diameter of a , and vertical diameter b , and is quite unobstructed,

$$g(y_\alpha, y_\beta) = \begin{cases} 1 & \frac{-a}{2} \leq y_\alpha \leq \frac{+a}{2} \\ & \frac{-b}{2} \leq y_\beta \leq \frac{+b}{2} \\ 0 & \text{elsewhere.} \end{cases}$$

If it is round,

$$g(y_\alpha, y_\beta) = \begin{cases} 1 & y_\alpha^2 + y_\beta^2 \leq \left(\frac{a}{2}\right)^2 \\ 0 & \text{elsewhere.} \end{cases}$$

The following theorems from diffraction theory are applicable:

1. Let the pupil function be expressed in terms of the variables $\omega_\alpha = y_\alpha/\lambda$ and $\omega_\beta = y_\beta/\lambda$ where λ is the wavelength of light. In this way we define the function $\hat{S}(\omega_\alpha, \omega_\beta) = g(y_\alpha, y_\beta)$. It is the Fourier transform of $\hat{s}(\alpha, \beta)$, the amplitude distribution of the disturbance in the image of a point object, that is to say

$$\hat{s}(\alpha, \beta) = \int \int_{-\infty}^{+\infty} g(y_\alpha, y_\beta) e^{-2\pi i(y_\alpha \alpha + y_\beta \beta)} dy_\alpha dy_\beta \quad (4)$$

2. The intensity distribution in the image of a point object, $s(\alpha, \beta)$ is the product of the amplitude distribution $\hat{s}(\alpha, \beta)$ and its complex conjugate $\hat{s}^*(\alpha, \beta)$. Thus

$$s(\alpha, \beta) = \hat{s}(\alpha, \beta) \cdot \hat{s}^*(\alpha, \beta) \quad (5)$$

3. The Fourier transform of the intensity distribution in the image of a point source $s(\alpha, \beta)$, is the modulation transfer function, $\tau(\omega_\alpha, \omega_\beta)$. Thus

$$\tau(\omega_\alpha, \omega_\beta) = \int_{-\infty}^{+\infty} s(\alpha, \beta) e^{-2\pi i(\alpha\omega_\alpha + \beta\omega_\beta)} d\alpha d\beta \quad (6)$$

4. The modulation transfer function is the self-convolution of the pupil function

$$\tau(\omega_\alpha, \omega_\beta) = g(y_\alpha, y_\beta) \otimes g(y_\alpha, y_\beta) \quad (7)$$

In practice, the following steps can lead from the knowledge of the aperture function of an optical system to the detailed knowledge of the image intensity distribution:

1. The modulation transfer function is given by the self-convolution (auto-correlation) of the pupil aperture function, equation (7).
2. The product of the Fourier transform of the object intensity distribution with the modulation transfer function gives the Fourier transform of the image intensity distribution. In the space domain, the equivalent statement is that the image intensity distribution is the convolution of the object intensity distribution, with the intensity distribution in the image of a point source. These relationships, which are derived and discussed in considerable detail in the relevant texts (O'Neil, Born and Wolf) are summarized in Table 1.

The fundamental imaging property of an instrument like the eye is, therefore, contained in the modulation transfer function $\tau(\omega_\alpha, \omega_\beta)$ in the Fourier domain, and in the image intensity distribution of a point object $s(\alpha, \beta)$ in the space domain. This holds not only for purely diffraction-limited instruments but also for those afflicted with aberrations. In such cases the modulation transfer function will be lower (less transmission of spatial frequencies) and the point-spread function larger.

The application of this theory to imaging in the eye of incoherently illuminated objects has been discussed in a previous study (Westheimer, 1964).

Coherent Illumination The whole situation is radically altered if the target is transilluminated by light which has its origin in a single source, so that the radiation emitted from different portions of the target is capable of mutually interfering. This is realized in Maxwellian viewing when the source S

TABLE 1 Incoherent Imagery

	Object Space α, β Visual Angle	Pupil Plane y_α, y_β Distance	Fourier Domain $\omega_\alpha, \omega_\beta$ (c/min of arc)	Image Space α, β (min of arc)
Object	$o(\alpha, \beta)$		$O(\omega_\alpha, \omega_\beta)$	
Pupil aperture function		$g(y_\alpha, y_\beta)$		
Amplitude function in image of point object			$\hat{S}(\omega_\alpha, \omega_\beta)$	$\hat{s}(\alpha, \beta)$
Transfer function			$\tau(\omega_\alpha, \omega_\beta) = g(y_\alpha, y_\beta) \otimes g(y_\alpha, y_\beta)$	
Intensity function in image of point object				$s(\alpha, \beta) = \hat{s}(\alpha, \beta) \cdot \hat{s}^*(\alpha, \beta)$
Image intensity function			$I(\omega_\alpha, \omega_\beta) = \tau(\omega_\alpha, \omega_\beta) \cdot O(\omega_\alpha, \omega_\beta)$	$i(\alpha, \beta) = o(\alpha, \beta) \hat{s}^*(\alpha, \beta)$

See footnote to Table 2.

in Fig. 1(a) is very small, and monochromatic. Then the light after being collimated by lens L_1 may be regarded as perfectly coherent, and the following theory holds rigidly. (It also holds for targets illuminated by lasers which emit perfectly coherent light.†)

The following is a convenient way of visualizing the situation. Since the wavefront impinging on the target is plane, there will be formed in the plane of the eye's pupil the Fraunhofer diffraction pattern produced at the passage of the light through the target. The amplitude distribution in this diffraction image is the Fourier transform of the target amplitude function.

If $\delta(\alpha, \beta)$ is the distribution of amplitude in the object, the distribution of amplitude in the disturbance in the plane of the pupil is given by

$$O(\omega_a, \omega_\beta) = \iint_{-\infty}^{+\infty} \delta(\omega_a, \omega_\beta) e^{-2\pi i(\omega_a \alpha + \omega_\beta \beta)} d\alpha d\beta \quad (8)$$

Naturally, only that portion of $O(\omega_a, \omega_\beta)$ that is transmitted by the pupillary aperture can be effective in the formation of the image amplitude distribution $i(\alpha, \beta)$ which is again the Fourier transform of the pupil amplitude function. Thus

$$i(\alpha, \beta) = \iint_{\text{pupil}} O(\omega_a, \omega_\beta) e^{-2\pi i(\omega_a \alpha + \omega_\beta \beta)} d\omega_a d\omega_\beta \quad (9)$$

Finally, the image intensity distribution $i(\alpha, \beta)$ is the product of the image amplitude distribution with its complex conjugate:

$$i(\alpha, \beta) = i(\alpha, \beta) \cdot i^*(\alpha, \beta)$$

It is seen that the limitation to the fidelity of reproduction of the original target on the retina (aside from any aberrations that may be present) is the fact that the pupil aperture is necessarily finite, while for any finite target the diffraction image, or Fourier transform, that is formed in the pupillary plane will be infinite in extent. All portions of the diffraction image of the target that fall outside the pupil will be cut off. The diffraction limitation during Maxwellian viewing lies therefore in equation (9) which states that the range of integration in the Fourier re-synthesis of the image extends only over the pupil instead of going to infinity. Now there is a very simple relationship between the target's amplitude spatial frequencies (ω_a, ω_β) , and the corresponding positions in the diffraction image in the pupil. Each spatial frequency ω is situated at a distance $y = \lambda\omega$ from the geometrical image of the source in the pupil. The pupil aperture function $g(y_\alpha, y_\beta)$, written in terms of the variables $\omega_a = y_\alpha/\lambda$, $\omega_\beta = y_\beta/\lambda$, thus, gives us the amplitude transmission function $\tau(\omega_a, \omega_\beta)$. It is this function which, when multiplied with $O(\omega_a, \omega_\beta)$, the Fourier transform of the object amplitude function, gives $I(\omega_a, \omega_\beta)$ the Fourier transform of the image amplitude function. All this is also summarized in Table 2.

Two examples follow. Suppose that the disturbance issuing from the target plane has an amplitude distribution that is sinusoidal with spatial frequency ω , i.e.

$$O(\alpha, \beta) = 1 + \cos(2\pi\alpha\omega)$$

The associated amplitude distribution in the plane of the pupil is three delta functions, i.e. there is a concentration of disturbance in 3 spots in the pupillary plane and nowhere else. The three are:

$$(1) \quad \omega_a = 0, \text{ i.e. } y_a = \omega_a \lambda = 0, \text{ i.e. at the center of the pupil and}$$

$$(2 \text{ and } 3) \quad \omega_a = \pm\omega, \text{ i.e. } y_a = \pm\omega\lambda$$

† A laser may be substituted for the source and lens L_1 . It is to be noted that all the light from the laser is then focused in the eye's pupil. Obviously, precautions are necessary here, just as much as in using a laser and letting the eye's optical system focus its beam on the retina.

TABLE 2 Incoherent Imagery

	Object Space α, β Visual Angle	Pupil Plane y_α, y_β Distance	Spatial Frequency Domain ω (c/min of arc)	Image Space
Object amplitude function	$o(\alpha, \beta)$		$O(\omega_\alpha, \omega_\beta)$	
Aperture function		$g(y_\alpha, y_\beta)$		
Complex amplitude transmission function			$\tau(\omega_\alpha, \omega_\beta)$	
Image amplitude function			$I(\omega_\alpha, \omega_\beta) = \tau(\omega_\alpha, \omega_\beta) \cdot O(\omega_\alpha, \omega_\beta)$	$\hat{i}(\alpha, \beta)$
Image intensity function				$i(\alpha, \beta) = i(\alpha, \beta) i^*(\alpha, \beta)$

The relationship between variables is as follows:

- Where quantities in 1st and 3rd columns are entered between a pair of horizontal lines, they are Fourier transforms:

$$o(\omega_\alpha, \omega_\beta) = \iint_{-\infty}^{+\infty} O(\alpha, \beta) e^{-2\pi i(\alpha\omega_\alpha + \beta\omega_\beta)} d\alpha d\beta$$

where α , and β are visual angles in minutes of arc, $\omega_\alpha, \omega_\beta$ are spatial frequency in c/min of arc.

- Distances y_α, y_β in plane of entrance pupil measured from its center are related to the spatial frequencies in the same orientation by the equations

$$y_\alpha = \omega_\alpha \lambda \quad y_\beta = \omega_\beta \lambda$$

λ being the wavelength of light used. Functions in the second column (pupil space) are transferred directly to the third column, spatial frequency, by this change in variables.

The above example shows that the representation of a sinusoidal grating in the frequency domain consists of a d.c. term and a pair of a.c. terms which are physically placed in the pupillary plane in positions according to their spatial frequency. Frequencies higher than the one present at the edge of the pupil can clearly not be available to make up the image. If the pupil aperture has diameter a , this limit of sinusoidal components of amplitude is given by

$$\omega_\alpha = a/2\lambda$$

Another example: suppose the target is a trans-illuminated circular disk of diameter α_0 rad, centered on the optical system. This case has been dealt with in detail in a previous study (Westheimer, 1959). The amplitude distribution in the diffraction image in the pupil is here the familiar Airy disk, with the first zero at a distance from the center given by

$$y = 1.22 \lambda / \alpha_0$$

y and λ being in the same units. For a disk 2.0 min of arc in diameter (0.00058 rad) and light of a representative wavelength, $y = 1.2$ mm. Thus, if the pupil diameter is smaller than 2.4 mm (or if an artificial pupil of smaller diameter than 2.4 mm is placed in front of the eye) only the center portion of the diffraction image of the disk is admitted into the eye. This constitutes a severe restriction of the “bandwidth” of spatial frequencies making up the original disk. The apparent anomaly of the situation is that most of the amplitude of the disturbance (>90 per cent) from the point source passing through the target disk, is contained in the central part of the diffraction image. Yet much of the “information” about the target—whether or not it is a sharply bounded disk and the position and exact nature of its edge—is contained in the higher spatial frequencies. They are imaged in the periphery of the eye’s entrance pupil and could thus be screened out needlessly by an experimenter who, unaware of their information content, might stop them from entering the eye by placing too small an artificial pupil in front of an eye seeing the target in Maxwellian view.

It should be emphasized that the process of coherent imaging is linear in *amplitude* of the electromagnetic radiation, not in *intensity* as it is with incoherent targets. When a target is lit by coherent light, it is helpful to consider the amplitude of the disturbance emerging from the target plane as being focused in the pupil plane. The amplitude of the disturbance there, which is the Fourier

transform of that in the target plane, represents the spatial frequency content of the latter displayed at distances from the geometrical image of the source proportional to the spatial frequencies. Since the pupil necessarily cuts out all portions of the diffraction pattern further out than its radius, the disturbance entering the eye will now make up an image amplitude distribution on the retina that is restricted in bandwidth of its spatial frequencies. In addition, intraocular effects and aberrations may also superimpose inhomogeneities and phase shifts. Ultimately the amplitude distribution of the disturbance in the plane of the retina is subjected to a non-linear transformation, viz. squaring, and the resulting intensity distribution constitutes the retinal image. Once its origin is properly understood, the Stiles-Crawford effect can be allowed for in this series of transformations. The place of the non-linearity in coherent imagery—it occurs as the last step in the transformation—prevents the formulation of any general statements about the retinal image. All particular cases of target configurations, pupil openings (not to mention aberrations and defocusing) will give rise to their own kinds of imaging characteristics. For incoherent target illumination, the non-linearity occurs at a more convenient juncture in the sequence, for once the point-spread function, or in the Fourier domain, the modulation transfer function, has been fixed by the particulars of the pupil size, aberrations and focus errors, the image for any target can be readily determined.

When it is coherent, the amplitude of the electromagnetic disturbance emerging from the target plane is being spatially Fourier analyzed (displayed as a diffraction image in the plane of the pupil), filtered, and then resynthesized as a distribution of amplitude of electromagnetic disturbance in the retinal plane. Ultimately the amplitude is squared at each point in the image and this constitutes the image intensity pattern. It is seen that filtering the spatial frequency in the domain of the amplitude of the disturbance will lead to results which, when looked at in the realm of the intensity pattern, may be non-linear. As an example, consider the case of an amplitude distribution that is sinusoidal. When squared this becomes an intensity distribution that is also sinusoidal but has twice the frequency. Suppose that now a constant (d.c.) level is added in the realm of amplitude. The laws of linearity are observed, so that the new amplitude distribution is the sum of the a.c. and d.c. terms. But when this is squared to yield the intensity distribution, a result is obtained which is quite different from that obtained when the amplitude distribution due to the a.c. term alone is squared.

Partial Coherence

The theory of coherent imagery (Beran and Parrant, 1964; Born and Wolf, 1959) developed above is applicable when the source S is essentially a monochromatic point source. (A laser in place of source S and lens L_1 realizes this situation in practice.) It is then possible for light coming from a point on the target to interfere with all the light coming from all the other parts of the target. When the source is larger than that, a point on the target is transilluminated by light from a range of points on the source, and the capacity to subsequently interfere is limited. In the limit, when the source is quite extensive, the target may be regarded as incoherently illuminated. These effects are described by the theory of partial coherence, the chief theorem of which, the van Cittert-Zernike theorem, permits us to define the target region beyond which no significant degree of coherence exists.

Applied to the situation illustrated in Fig. 2, the van Cittert-Zernike theorem states that the degree of coherence in the target plane is distributed in the manner of the Fraunhofer pattern that would be produced in the hypothetical case due to diffraction at the aperture of the source S . Let the source be circular with linear diameter y_0 , let the ratio of the focal lengths of the lens L_2 to that of lens L_1 (the magnification of the source in the plane of the pupil) be M . Then the separation of two points in the source for the first reduction of their mutual coherence to zero (i.e. the radius of the Fraunhofer diffraction pattern) expressed in terms of α , the angular distance at the eye, is

$$\alpha = 1.22\lambda/My_0$$

It is seen that for α to be, say, 1 min of arc, My_0 should be about 2.2 mm, i.e. the geometrical image of the source in the plane of the pupil should be 2.2 mm. This is an outer limit, since two points on the source separated by, say, half the radius of the Airy disk due to the source, would have a mutual coherence that is already small.

It is thus possible to lay down the following guidelines on whether a particular experimental situation calls for the application of the theory of coherent or incoherent imagery.

When the source has an image in the plane of the pupil larger than, say, 2.5 mm in diameter, the target may be regarded as incoherently illuminated. An artificial pupil can with advantage be employed and the approach to the light distribution in the retinal image would be by way of the modulation transfer function of the eye.

On the other extreme, one may wish to utilize the characteristics of coherent imagery for a particular purpose, such as spatial filtering. How small does the geometrical image of the source have to be for a satisfactory degree of coherence to be achieved? We have mentioned that the laser will do this, but all the essential ingredients are also present in a tungsten filament or mercury or xenon arc, provided that the angular subtense of the source at lens L_1 is small enough and the light has been restricted to a narrow enough wavelength band. If one wishes to ensure that the light coming from a 10 min of arc region of the visual field is at least, say, 90 per cent coherent, the geometrical image of the source in the pupil should be no larger than

$$My_0 = 0.32\lambda/2.9 \times 10^{-3} \quad \text{i.e.} \quad 0.06 \text{ mm}$$

When an interference filter is used to make the beam monochromatic, the resultant retinal illuminance cannot be very high. Nevertheless, a tungsten filament with a luminance of 10^3 cd/cm² can still give a few hundred trolands of quasi-monochromatic retinal illuminance with good coherence properties. With the use of a high-pressure mercury arc this can be increased by at least two orders of magnitude.

References

- Aguilar, M. and Stiles, W. S. (1954). Saturation of the rod mechanism of the retina at high levels of stimulation. *Optica Acta* 1, 59–65.
- Beran, M. J. and Parrent, G. B. (1964). *Theory of Partial Coherence*. Prentice Hall, Englewood Cliffs, N.J.
- Born M. and Wolf, E. (1959). *Principles of Optics*. Pergamon Press, Oxford.
- Emsley, H. H. (1939). *Visual Optics*. 2nd Edition, Hatten Press, London.
- Hopkins, H. H. (1953). On the diffraction theory of optical images. *Proc. R. Soc. A* 217, 408–432.
- Le Grand, Y. (1948). *Optique Physiologique*. Tome 2ieme, Ed. Rev, d'Optique.
- Maxwell, J. C. (1860)—On the theory of compound colours, and the relations of the colours of the spectrum. *Phil. Trans. R. Soc.* 150, 57.
- O'Neill, E. L. (1963). *Introduction to Statistical Optics*. Addison-Wesley, Reading, Mass.
- Rushton, W. A. H. (1956a). The difference spectrum and the photo-sensitivity of the rhodopsin in the living eye. *J. Physiol.* 134, 11–29.
- Rushton, W. A. H. (1956b). The rhodopsin density in the human rods. *J. Physiol.* 134, 30–46.
- Walsh, J. W. T. (1958). *Photometry*. 3rd Edition, Constable, London.
- Westheimer, G. (1959). Retinal light distributions for circular apertures in Maxwellian View, *J. Opt. Soc. Am.* 49, 41.
- Westheimer, G. (1960). Modulation thresholds for sinusoidal light distribution on the retina. *J. Physiol.* 152, 67–74.
- Westheimer, G. (1964). Pupil size and visual resolution. *Vision Res.* 4, 39–45.

6.3 POSTSCRIPT (2008)

These supplementary remarks are informed by developments in the two intersecting topics which this article, written 45 years ago, covers: physical optics and dioptrics of the eye.

The Maxwellian view is a nontraditional way of generating the retinal image. In ordinary viewing, the energy over the whole bundle of rays—or the wavefront—reaching the entrance pupil is uniform.

The recent great advances in adaptive optics, that is, the modifying of the wavefront for purposes of neutralizing deviations caused by the eye's dioptric defects, deal only with the phase, not amplitude (Chap. 15 in this volume, and Chap. 5, "Adaptive Optics," by Robert Q. Fugate in Vol. V).

But it is precisely the amplitude of the electromagnetic disturbance across the entering wavefront that is in play in Maxwellian viewing. A simple way of formulating the eye's imaging is the statement that Fourier transformation relates the complex amplitude distributions of the electromagnetic disturbance in the planes of the exit pupil and the retinal image. The image intensity distribution is then the product of the amplitude with its complex conjugate—in practice squaring, once phase effects have been factored in. Hence in principle one can configure a desired retinal image by manipulating the incoming wavefront as it reaches the eye. In the most elementary case of Maxwellian viewing, the disturbance in the pupil plane is merely the diffraction pattern generated when a coherent beam passes through a transparent target comes to a focus there. But spatial amplitude and phase distributions can be created by other means. Limiting it to two points by a double prisms, for example, is a time-honored manner of producing interference fringes on the retina. If, a technology were available that allows both amplitude and phase modulation of a uniform plane wavefront incident on the eye—instead of only the phase as currently implemented in "adaptive optics"—an unprecedented control of the retinal image would be achieved in the absence of any physical objects. A start of this has already been made in attempts to measure the performance of the retinal mosaic by not just interference fringes but also in two-point separation acuity.¹

The original paper discussed the always difficult problem of partial coherence. This is only appropriate, if for no other reason than that resolution measurements as a function of the degree of coherence had been performed in the first instance by visual observation². In the early days of this research, generating a highly coherent light signal was quite difficult. It involved strictly monochromatic light, usually produced by mercury arcs limited to a single spectral line by narrow-band interference filters, focused on a pin-hole aperture only a fraction of a millimeter in diameter. The possibility of coherence being only partial needed to be kept in mind, particularly since its mathematical handling involves delicate problems.

All this changed radically with the advent of the laser. This advance in technology took place in the interval between two papers devoted to measuring retinal resolution bypassing the eye's optics and illustrates how it allowed researchers to concentrate on the biology rather than the physics of the process.^{3,4} As the result of the easy access to fully coherent beams, partial coherence is now of more theoretical than practical interest.

Utilizing lasers and sources other than incandescent ones has opened up a question that at one time occupied center stage in vision research: photon statistics at the absolute visual threshold. Hecht, Shlaer and Pirenne's epochal work,⁵ done while the quantum theory of light was still novel, pointed to the physical basis of threshold fluctuation through the contention that a Poisson distribution was the governing rule when the number of absorbed photons was only a handful. Since then the subject of photon statistics has been opened more widely with the utilization of photon emission both from lasers^{6,7} and from oscilloscope screens,⁸ each of which are used in vision research.

In the process of estimating the retinal image by way of Fourier transformation of the complex amplitude distribution of the electromagnetic disturbance in the plane of the entrance pupil, the question cannot be skirted of how to factor in the Stiles-Crawford effect⁹ (see Chaps. 8 and 9). The original measurement involved the reduction in luminous efficiency for a narrow bundle as it was moved from the center to the edge of the pupil. The most trivial explanation of absorption in the eye media was quickly eliminated by the fact that the phenomenon differs between rod and cone vision.

The problem arises when computing the light distribution in the image formed by the eye's optics as it is in turn is funneled into the receptors to form the immediate stimulus for the visual process leading to photoisomerizations of photopigment molecules in rods and cones. Most calculations make no particular assumptions about the uptake of radiation in retinal receptors, yet their properties in channeling and accepting radiation cannot be ignored, especially if that involves conditioning the impinging wavefront. There are two possible interpretations of the role of the Stiles-Crawford effect at the level just prior to the phototransduction:

- a. The effect is an apodization phenomenon, where the entering wavefront near the edge of the pupil somehow is diminished in amplitude, perhaps as a result of waveguide properties of receptors (Chap. 9) or acceptance-lobe limitations at the molecular absorption levels. In such a case, the

calculation from pupil aperture to image distributions would have to fold in amplitude variations across the wavefront; in other words, the amplitude of the wavefront is in effect diminished toward the pupil periphery as the square root of the measured Stiles-Crawford effect, the latter measurements having been obtained in the intensity, not the amplitude domain of the electromagnetic disturbance.

- b. The effect is a diminution of the Poynting vector of the electromagnetic disturbance as a function of its direction with respect to the orientation of the receptor cells. This would be implemented by physical light screening at the retina, for example, inter-receptor absorbing pigments.

The case for (a) has been made well enough by arguments in favor of retinal receptors acting as waveguides (Chap. 8) and resultant improvement in image quality through reducing the weight of the wavefront in the pupil periphery more prone to be afflicted by ocular aberrations. It is strengthened by the fact that the Stiles-Crawford effect is much more prominent in cone vision where light level is not an issue but image quality is, than in rod vision.¹⁰

But the alternative (b) cannot yet be dismissed. There is evidence that the photopigment pools reached by light entering at opposite edges of the pupil are not identical.¹¹ The situation has been laid out in greater detail elsewhere.^{12,13}

An effective retinal image, that is, a pattern of intra-receptor light absorption, based on proposition (a) will differ from one based on (b) and would lead to a rod-cone difference in spatial visual stimulus patterns.

References

1. J. Liang and G. Westheimer, "Method for Measuring Visual Resolution at the Retinal Level," *J. Opt. Soc. Am. A* **10**(8):1691–1696 (1993).
2. A. Arnulf, O. Dupuy, and F. Flamant, "Étude expérimentale de la variation de la limite de résolution en fonction de la cohérence," *Revue d'Optique* **32**:529–552 (1953).
3. G. Westheimer, "Modulation Thresholds for Sinusoidal Light Distributions on the Retina," *J. Physiol. (London)* **152**(1):67–74 (1960).
4. F. W. Campbell and D. G. Green, "Optical and Retinal Factors Affecting Visual Resolution," *J. Physiol.* **181**:576–593 (1965).
5. S. Hecht, S. Shlaer, and M. H. Pirenne, "Energy, Quanta and Vision," *J. Physiol.* **25**:819–840 (1942).
6. E. Wolf and C. L. Mehta, "Determination of the Statistical Properties of Light from Photoelectric Measurements," *Physical Rev. Letters* **13**:1–3 (1964).
7. R. J. Glauber, "Optical Coherence and Photon Statistics," *Quantum Optics and Electronics—Les Houches Lectures 1964*, C. DeWitt, (ed.), Gordon & Breach, New York, 1965.
8. M. C. Teich, P. R. Prucnal, G. Vannucci, M. E. Breton, and W. J. McGill, "Multiplication Noise in the Human Visual System at Threshold. 3. The role of non-Poisson quantum fluctuations," *Biol. Cybernetics* **44**:157–165 (1982).
9. W. S. Stiles and B. H. Crawford, "The Luminous Efficiency of Rays Entering the Eye Pupil at Different Points," *Proc. Roy. Soc. (London) B* **112**:428–450 (1933).
10. W. S. Stiles, "The Directional Sensitivity of the Retina and the Spectral Sensitivities of the Rods and Cones," *Proc. Roy. Soc. (London) B* **127**:64–105 (1939).
11. W. Makous, "A Transient Stiles-Crawford Effect," *Vis. Res.* **8**:1271–1284 (1968).
12. G. Westheimer, "Specifying and Controlling the Optical Image on the Retina," *Progress in Retinal and Eye Research* **25**:19–42 (2006).
13. G. Westheimer, "Directional Sensitivity of the Retina: Seventy-Five Years of Stiles-Crawford Effect," *Proc. Roy. Soc. (London) B* **275**:2777–2789 (2008).

This page intentionally left blank.

DO NOT DUPLICATE

OCULAR RADIATION HAZARDS*

David H. Sliney

Consulting Medical Physicist

Fallston, Maryland

Retired, U.S. Army Center for Health Promotion and Preventive Medicine

Laser/Optical Radiation Program

Aberdeen Proving Ground, Maryland

7.1 GLOSSARY

Bunsen-Roscoe law of photochemistry. The reciprocal relation between irradiance (dose-rate) and exposure duration to produce a constant radiant exposure (dose) to produce an effect.

Coroneo effect. See *Ophthalmoheliosis*.

Erythema. Skin reddening (acute damage) such as that caused by UV radiation; sunburn.

Macula lutea. Yellowed-pigment central region (macular area) of the retina of the human eye.

Maxwellian view. A method for directing a convergent beam into the eye's pupil to illuminate a large retinal area.

Ophthalmoheliosis. The concentration of UV radiation upon the nasal region of the cornea (the limbus) and nasal equatorial region of the lens, thus contributing to the development of pterygium on the cornea and cataract in the lens beginning in these regions.

Photokeratitis. A painful inflammation of the cornea of the eye.

Photoretinitis. Retinal injury due to viewing the sun or other high-intensity light source arising from photochemical damage due to short, visible wavelengths (also termed *photic maculopathy*).

Pterygium. A fleshy growth that invades the cornea (the clear front “window” of the eye). It is an abnormal process in which the conjunctiva (a membrane that covers the white of the eye) grows into the cornea. A pterygium may be small or grow large enough to interfere with vision and commonly occurs on the inner (nasal) corner of the eye.

Solar retinopathy. Retinal injury caused from staring at the sun—photochemical damage produced by short, visible (blue) light (also referred to as *solar maculopathy*, since it normally occurs in the macula lutea).

*Note: The opinions or assertions herein are those of the author and should not be construed as official policies of the U.S. Department of the Army or Department of Defense.

7.2 INTRODUCTION

Optical radiation hazards to the eye vary greatly with wavelength and also depend upon the ocular exposure duration. Depending upon the spectral region the cornea, conjunctiva, lens, and/or retina may be at risk when exposed to intense light sources such as lasers and arc sources and even the sun. Photochemical damage mechanisms dominate in the ultraviolet (UV) end of the spectrum, and thermal damage mechanisms dominate at longer wavelengths in the visible and infrared (IR) spectral regions. Natural aversion responses to very bright light sources limit ocular exposure, as when one glances at the sun, but artificial sources may be brighter than the solar disk (e.g., lasers) or have different spectral or geometrical characteristics that overcome natural avoidance mechanisms, and the eye can suffer injury. Guidelines for the safe exposure of the eye from conventional light sources as well as lasers have evolved in recent decades. Because of special concerns about laser hazards, safety standards and regulations for lasers exist worldwide. However, the guidelines for exposure to conventional light sources are more complex because of the likelihood of encountering more than one type of ocular hazard from the broader emission spectrum.

The adverse effects associated with viewing lasers or other bright light sources such as the sun, arc lamps, and welding arcs have been studied for decades.¹⁻⁴ During the last three decades, guidelines for limiting exposure to protect the eye have evolved.⁴⁻¹² The guidelines were fostered to a large extent by the growing use of lasers and the quickly recognized hazard posed by viewing laser sources. Injury thresholds for acute injury in experimental animals for corneal, lenticular, and retinal effects have been corroborated for the human eye from accident data. Exposure limits are based upon this knowledge.³ These safe exposure criteria were based not only upon studies aimed at determining thresholds in different ocular structures, but also upon studies of injury mechanisms and how the different injury mechanisms scaled with wavelength, exposure duration, and the area of irradiated tissue. The exposure guidelines also had to deal with competing damage to different ocular structures such as the cornea, lens, and retina, and how these varied with wavelength.³⁻¹¹ Although laser guidelines could be simplified to consider only the type of damage dominating for a single wavelength and exposure duration, the guidelines for incoherent, broadband sources were necessarily more complex.¹¹ To understand the guidelines for safe exposure, it is first necessary to consider the principal mechanisms of injury and which ocular structures may be injured under different exposure conditions.

7.3 INJURY MECHANISMS

For most exposure durations ($t > 1$ s), optical radiation injury to ocular structures is dominated by either photochemically or thermally initiated events taking place during or immediately following the absorption of radiant energy. At shorter durations, nonlinear optical interaction mechanisms may play a role.¹² Following the initial insult; biological repair responses may play a significant role in determining the final consequences of the event. While inflammatory, repair responses are intended to reduce the sequelae; in some instances, this response could result in events such as scarring, which could have an adverse impact upon biological function.³

Action Spectra

Photochemical and thermal effects scale differently with wavelength, exposure duration, and irradiated spot size. Indeed, experimental biological studies in humans and animals make use of these different scaling relationships to distinguish which mechanisms are playing a role in observed injury. Photochemical injury is highly wavelength dependent; thermal injury depends upon exposure duration.

Exposure Duration and Reciprocity

The Bunsen-Roscoe law of photochemistry describes the *reciprocity* of exposure rate and duration of exposure, which applies to any photochemical event. The product of the dose-rate (in watts per square centimeter) and the exposure duration (in seconds) is the exposure dose (in joules

per square centimeter at the site of absorption) that may be used to express an injury threshold. Radiometrically, the irradiance E in watts per square centimeter, multiplied by the exposure duration t , is the radiant exposure H in joules per square centimeter, as shown in Eq. (1).

$$H = E \cdot t \quad (1)$$

This reciprocity helps to distinguish photochemical injury mechanisms from thermal injury (burns) where heat conduction requires a very intense exposure within seconds to cause photocoagulation; otherwise, surrounding tissue conducts the heat away from the absorption site (e.g., from a retinal image). Biological repair mechanisms and biochemical changes over long periods and photon saturation for extremely short periods will lead to reciprocity failure.

Thermal injury is a *rate process* that is dependent upon the time-temperature history resulting from the volumic absorption of energy across the spectrum. The thermochemical reactions that produce coagulation of proteins and cell death require critical temperatures for detectable biological injury. The critical temperature for an injury gradually decreases with the lengthening of exposure duration. This approximate decrease in temperature varies over many orders of magnitude in time and decreases approximately as a function of the exposure duration t raised to the -0.25 power [i.e., $f(t^{-0.25})$].

As with any photochemical reaction, the *action spectrum* should be known.¹³ The action spectrum describes the relative effectiveness of different wavelengths in causing a photobiological effect. For most photobiological effects (whether beneficial or adverse), the full width at half-maximum is less than 100 nm and a long-wavelength cutoff exists where photon energy is insufficient to produce the effect. This is not at all characteristic of thermal effects; in which the effect occurs over a wide range of wavelengths where optical penetration and tissue absorption occur. For example, significant radiant energy can penetrate the ocular media and be absorbed in the retina in the spectral region between 400 and nearly 1400 nm; the absorbed energy can produce retinal thermal injury.

Although the action spectra for acute photochemical effects upon the cornea,^{3,14-16} upon the lens^{3,16} and upon the retina^{3,17} have been published, the action spectra of some other effects appear to be quite imprecise or even unknown.¹⁴ The action spectrum for neuroendocrine effects mediated by the eye is still only approximately known. This points to the need to specify the spectrum of the light source of interest as well as the irradiance levels if one is to compare experimental results from different experimental studies.

Although both E and H may be defined over the entire optical spectrum, it is necessary to employ an action spectrum for photochemical effects. The International Commission on Illumination (CIE) photopic $V(\lambda)$, UV hazard $S(\lambda)$, and blue-light hazard $B(\lambda)$ curves of Fig. 1 are examples of action spectra that may be used to spectrally weight the incident light. With modern computer spreadsheet programs, one can readily spectrally weight a lamp's spectrum by a large variety of photochemical action spectra. These computations all take the following form:

$$E_{\text{eff}} = \Sigma E_{\lambda} \cdot A(\lambda) \cdot \Delta(\lambda) \quad (2)$$

where $A(\lambda)$ may be any action spectrum (unitless) of interest. One can then compare different sources to determine the relative effectiveness of the same irradiance from several lamps for a given action spectrum.

7.4 TYPES OF INJURY

The human eye is actually quite well-adapted to protect itself against the potential hazards from optical radiation (UV, visible, and IR radiant energy) from most environmental exposures encountered from sunlight. However, if ground reflections are unusually high, as when snow is on the ground, reflected UV radiation may produce "snow blindness," that is, UV photokeratitis. Another example occurs during a solar eclipse, if one stares at the solar disc for more than 1 to 2 min without eye protection, the result may be an "eclipse burn" of the retina.^{3,17} Lasers may injure the eye, and under unusual situations other artificial light sources, such as mercury-quartz-discharge lamps or arc lamps, may pose a potential ocular hazard. This is particularly true when the normal defense

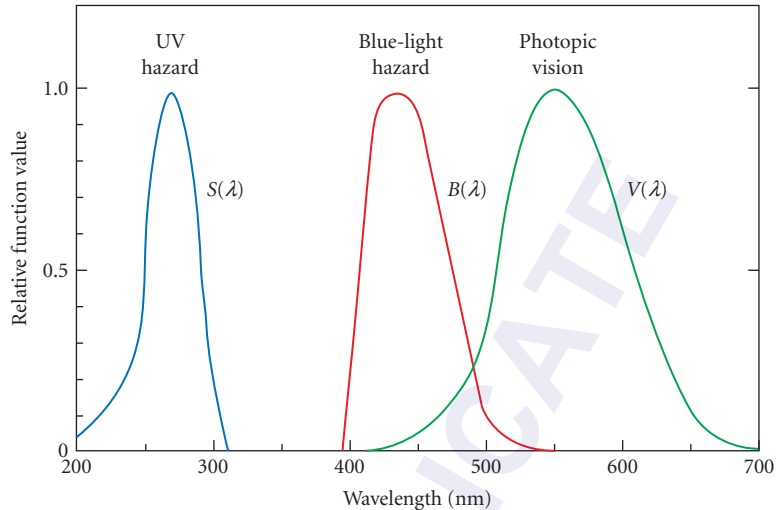


FIGURE 1 Action spectra. The ACGIH UV hazard function $S(\lambda)$ describes approximately the relative spectral risk for photokeratitis, and the blue-light hazard function $B(\lambda)$ describes the spectral risk for photoreinitis. For comparison, the CIE spectral sensitivity (standard observer) function $V(\lambda)$ for the human eye is also provided.

mechanisms of squinting, blinking, and aversion to bright light are overcome. Ophthalmic, clinical exposures may increase this risk if the pupil is dilated or the head is stabilized.

There are at least six separate types of hazards to the eye from lasers and other intense optical sources, and these are shown in Fig. 2 (along with the spectral range of responsible radiation):³

1. Ultraviolet photochemical injury to the cornea (photokeratitis); also known as “welder’s flash” or “snow blindness” (180 to 400 nm)^{3,15,16}
2. Ultraviolet photochemical injury to the crystalline lens (cataract) of the eye (295 to 325 nm, and perhaps to 400 nm)¹⁶
3. Blue-light photochemical injury to the retina of the eye (principally, 400 to 550 nm; unless aphakic, 310–550 nm)^{3,4,17}
4. Thermal injury (photocoagulation) to the retina of the eye (400 to nearly 1400 nm) and thermoacoustic injury at very short laser exposure durations^{3,4,12,18}
5. Near-IR thermal hazards to the lens (~800–3000 nm)^{3,19}
6. Thermal injury (burns) of the cornea of the eye (~1300 nm to 1 mm)³

The potential hazards to the eye and skin are illustrated in Fig. 2, which shows how effects are generally dominant in specific CIE photobiological spectral bands, UV-A, -B, and -C and IR-A, -B, and -C.²⁰ Although these photobiological bands are useful shorthand notations, they do not define fine lines between no effect and an effect in accordance with changing wavelengths.

Photokeratitis

This painful but transient (1 to 2 d) photochemical effect normally occurs only over snow or open arc. The relative position of the light source and the degree of lid closure can greatly affect the proper calculation of this UV exposure dose. For assessing risk of photochemical injury, the spectral distribution of the light source is of great importance.

Nonionizing radiation band	UV-C	UV-B	UV-A	Visible	IR-A	IR-B	IR-C	
Wavelength (nm)	100	280	315	400	760	1400	3000	10 ⁶
Adverse effects	Photokeratitis			Retinal burns		Corneal burns		
		Cataract →				Cataracts		
		Erythema		Color vision Night vision Degradation				
				Thermal skin burns				
Skin penetration of radiation (depth)								

FIGURE 2 The separate types of hazards to the eye from lasers and other intense optical sources, along with the spectral range of responsible radiation.³

Cataract

Ultraviolet cataract can be produced from exposure to UV radiation (UVR) of wavelengths between 295 and 325 nm (and even to 400 nm). Action spectra can only be obtained from animal studies, which have shown that anterior, cortical, and posterior subcapsular cataract can be produced by intense exposure delivered over a period of days. Human cortical cataract has been linked to chronic, lifelong UV-B radiation exposure. Although the animal studies and some epidemiological studies suggest that it is primarily UV-B radiation in sunlight, and not UV-A, that is most injurious to the lens, biochemical studies suggest that UV-A radiation may also contribute to accelerated aging of the lens. It should be noted that the wavelengths between 295 and 325 nm are also in the wavelength region where solar UVR increases significantly with ozone depletion. Because there is an earlier age of onset of cataract in equatorial zones, UVR exposure has frequently been one of the most appealing of a number of theories to explain this latitudinal dependence. The UVR guideline, established by the American Conference of Governmental Industrial Hygienists (ACGIH) and the International Commission on Non-Ionizing Radiation Protection (ICNIRP) is also intended to protect against cataract.

Despite the collection of animal and epidemiological evidence that exposure of the human eye to UVR plays an etiological role in the development of cataract, this role in cataractogenesis continues to be questioned by others. Indeed, more recent studies would support a stronger role for ambient temperature as the primary environmental etiological factor in one type of cataract—nuclear cataract.¹⁹

Because the only direct pathway of UVR to the inferior germinative area of the lens is from the extreme temporal direction, it has been speculated that side exposure is particularly hazardous. For any greatly delayed health effect, such as cataract or retinal degeneration, it is critical to determine the actual dose distribution at critical tissue locations. A factor of great practical importance is the actual UVR that reaches the germinative layers of any tissue structure. In the case of the lens, the germinative layer where lens fiber cell nuclei are located is of great importance. The DNA in these cells is normally well-shielded by the parasol effect of the irises. However, Coroneo²¹ has suggested that the focusing of very peripheral rays by the temporal edge of the cornea, those which do not even reach the retina, can enter the pupil and reach the equatorial region as shown in Fig. 3. He terms this effect, which can also produce a concentration of UVR at the nasal side of the limbus and lens, “ophthalmoheliosis.” He also noted the more frequent onset of cataract in the nasal quadrant of the lens and the formation of pterygium in the nasal region of the cornea. Figure 4 shows the percentage of cortical cataract that

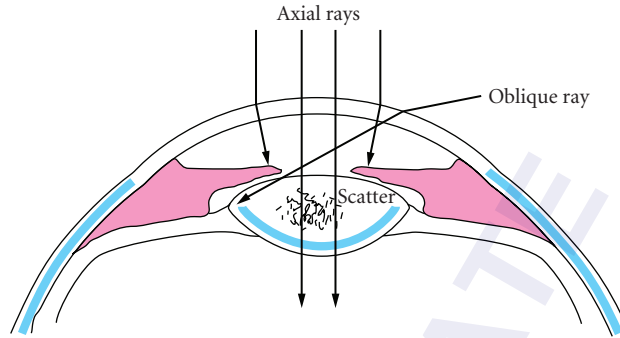


FIGURE 3 The Coroneo effect. Very oblique, temporal rays can be focused near the equator of the lens in the inferior nasal quadrant.

actually first appears in each quadrant of the lens from the data of Barbara Klein and her colleagues in their study of a population in the U.S. Midwest, in Beaver Dam, Wisconsin.²² This relationship is highly consistent with the Coroneo hypothesis.

Pterygium and Droplet Keratopathies

The possible role of UVR in the etiology of a number of age-related ocular diseases has been the subject of many medical and scientific papers. However, there is still a debate as to the validity of these arguments. Although photokeratitis is unquestionably caused by UVR reflection from the snow,^{3,23,24} pterygium and droplet keratopathies are less clearly related to UVR exposure.²⁴ Pterygium, a fatty growth over the conjunctiva that may extend over the cornea, is most common in ocean island residents (where both UVR and wind exposure is prevalent). Ultraviolet radiation

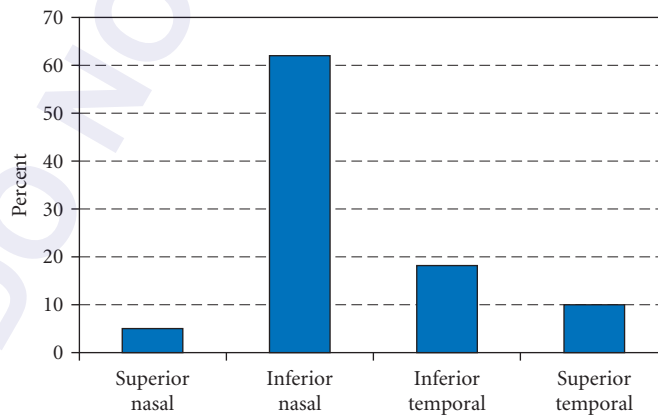


FIGURE 4 Distribution of cortical cataract by segment. The percentage of cortical cataract that actually first appears in each quadrant of the lens.²²

is a likely etiological factor,^{24,25} and the Coroneo effect may also play a role.^{21,25} To better answer these epidemiological questions, far better ocular dosimetry is required. Epidemiological studies can arrive at erroneous conclusions if assignments of exposure are seriously in error, and assumptions regarding relative exposures have been argued to be incorrect.²⁵ Before one can improve on current epidemiological studies of cataract (or determine the most effective UVR protective measures), it is necessary to characterize the actual solar UVR exposure to the eye.

Photoretinitis

The principal retinal hazard resulting from viewing bright continuous-wave (CW) light sources is photoretinitis [e.g., *solar retinitis* with an accompanying scotoma (“blind spot”), which results from staring at the sun]. Solar retinitis was once referred to as “eclipse blindness” and associated “retinal bum.” At one time, solar retinitis was thought to be a thermal injury, but it has been since shown conclusively (1976) to result from a photochemical injury related to exposure of the retina to shorter wavelengths in the visible spectrum (i.e., violet and blue light).^{3,17} For this reason, it has frequently been referred to as the “blue-light” hazard. The action spectrum for photoretinitis peaks at about 445 nm in the normal phakic eye; however, if the crystalline lens has been removed (as in cataract surgery), the action spectrum continues to increase for shorter wavelengths into the UV spectrum, until the cornea blocks shorter wavelengths, and the new peak shifts down to nearly 305 nm.

As a consequence of the Bunsen-Roscoe law of photochemistry, blue-light retinal injury (photoretinitis) can result from viewing either an extremely bright light for a short time or a less bright light for longer periods. The approximate retinal threshold for a nearly monochromatic source at 440 nm is approximately 22 J/cm², hence a retinal irradiance of 2.2 W/cm² delivered in 10 s, or 0.022 W/cm² delivered in 1000 s will result in the same threshold retinal lesion. Eye movements will reduce the hazard, and this is particularly important for sources subtending an angle of less than 11 milliradians (mrad), because saccadic motion even during fixation is of the order of 11 mrad.²⁷

Infrared Cataract

Infrared cataract in industrial workers appears only after lifetime exposures of the order of 80 to 150 mW/cm². Although thermal cataracts were observed in glassblowers, steelworkers, and others in metal industries at the turn of the century, they are rarely seen today. These aforementioned irradiances are almost never exceeded in modern industry, where workers will be limited in exposure duration to brief periods above 10 mW/cm², or they will be (should be) wearing eye protectors. Good compliance in wearing eye protection occurs at higher irradiances if only because the worker desires comfort of the face.³

7.5 RETINAL IRRADIANCE CALCULATIONS

Retinal irradiance (exposure rate) is directly related to source radiance (brightness). It is not readily related to corneal irradiance.³ Equation (3) gives the general relation, where E_r is the retinal irradiance in watts per square centimeter, L_s is the source radiance in watts per square centimeter per steradian (sr), f is the effective focal length of the eye in centimeters, d_e is the pupil diameter in centimeters, and t is the transmittance of the ocular media:

$$E_r = \frac{(\pi \cdot L_s \cdot \tau \cdot d_e^2)}{(4f^2)} \quad (3)$$

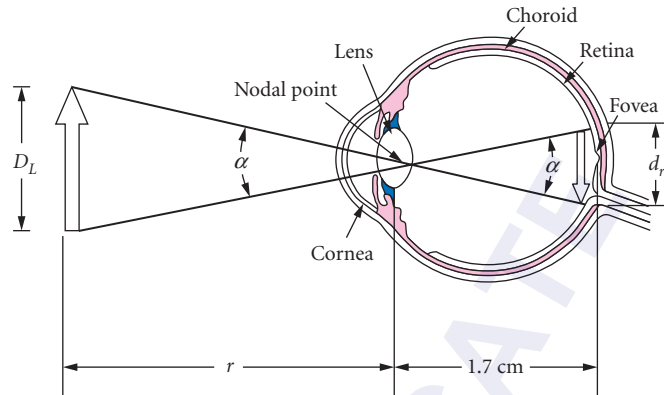


FIGURE 5 Retinal image size related to source size. The retinal image size can be calculated based upon the equal angular subtense of the source and the retinal image at the eye's nodal point (17 mm in front of the retina).

Equation (3) is derived by considering the equal angular subtense of the source and the retinal image at the eye's nodal point (Fig. 5). The detailed derivation of this equation is given elsewhere³ and in Chap. 37 "Radiometry and Photometry for Vision Optics," by Yoshi Ohno in Vol. II. The transmittance t of the ocular media in the visible spectrum for younger humans (and most animals) is as high as 0.9 (i.e., 90 percent).³ If one uses the effective focal length f of the adult human eye (Gullstrand eye), where $f = 1.7$ cm, one has

$$E_r = 0.27 \cdot L_s \cdot \tau \cdot d_e^2 \quad (4)$$

All of the preceding equations assume that the iris is pigmented, and the pupil acts as a true aperture. In albino individuals, the iris is not very effective, and some scattered light reaches the retina. Nevertheless, imaging of a light source still occurs, and Eq. (4) is still valid if the contribution of scattered light (which falls over the entire retina) is added.

7.6 EXAMPLES

As an example, a typical cool-white fluorescent lamp has an illumination of 200 lx, and the total irradiance is 0.7 mW/cm². By spectral weighting, the effective blue-light irradiance is found to be 0.15 mW/cm². From measurement of the solid angle subtended by the source at the measurement distance,^{3,11} the blue-light radiance is 0.6 mW/(cm²·sr) and the radiance is 2.5 mW/(cm² · sr). When fluorescent lamps are viewed through a plastic diffuser, the luminance and blue-light radiance are reduced.

As another example, a 1000-watt tungsten halogen bulb has a greater radiance, but the percentage of blue light is far less. A typical blue-light radiance is 0.95 W/(cm² · sr) compared with a total radiance of 58 W/(cm² · sr). The luminance is 2600 candelas (cd)/cm²—a factor of 3000 times brighter than the cool-white fluorescent. However, because the retinal image is small, eye movements spread the exposure over a much larger area.

Once the source radiance, L or L_B , is known, the retinal irradiance E_r is calculated by Eq. (4). The preceding examples illustrate the importance of considering the size of a light source and the impact of eye movements in any calculation of retinal exposure dose. If one were exposed to a focal beam

of light (e.g., from a laser or LED) that was brought to focus in the anterior chamber (the aqueous humor), the pupil plane, or the lens, the light beam would diverge past this focal point and could be incident upon the retina as a relatively large image. This type of retinal illumination is frequently referred to as *Maxwellian view* and does not occur in nature. The retinal irradiance calculation in this case would be determined by the depth of the focal spot in the eye; the closer to the retina, the smaller the retinal image and the greater the irradiance.

Because the iris alters its diameter (and pupil center) dynamically, one must be alert to vignetting. The pupil aperture also diminishes with age. Near-sighted, or myopic, individuals tend to have larger pupils, and far-sighted, or hyperopic, individuals tend toward smaller pupils. Accommodation also results in a decrease in pupil diameter.

7.7 EXPOSURE LIMITS

A number of national and international groups have recommended occupational or public exposure limits (ELs) for optical radiation (i.e., UV, light, and IR radiant energy). Although most such groups have recommended ELs for UV and laser radiation, only one group has recommended for some time ELs for visible radiation (i.e., light). This one group is well-known in the field of occupational health—the American Conference of Governmental Industrial Hygienists (ACGIH).^{5,6} The ACGIH refers to its ELs as *threshold-limit values* (TLVs); these are issued yearly, so there is an opportunity for a yearly revision. The current ACGIH TLVs for light (400 to 760 nm) have been largely unchanged for the last decade—aside from an increase in the retinal thermal limits in 2008. In 1997, with some revisions, the International Commission on Non-Ionizing Radiation Protection (ICNIRP) recommended these as international guidelines.¹¹ The ICNIRP guidelines are developed through collaboration with the World Health Organization (WHO) by jointly publishing criteria documents that provide the scientific database for the exposure limits.⁴ The ACGIH TLVs and ICNIRP ELs are generally applied in product safety standards of the International Commission on Illumination (CIE), the International Electrotechnical Commission (IEC), the International Standardization Organization (ISO), and consensus standards from the American National Standards Institute and other groups. They are based in large part on ocular injury data from animal studies and data from human retinal injuries resulting from viewing the sun and welding arcs. All of the guidelines have an underlying assumption that outdoor environmental exposures to visible radiant energy is normally not hazardous to the eye except in very unusual environments such as snow fields, deserts, or out on the open water.

Applying the UV Limits

To apply the UV guideline, one must obtain the average spectrally weighted UV irradiance E_{eff} at the location of exposure. The spectrally weighted irradiance is

$$E_{\text{eff}} = \sum E_{\lambda} \cdot S(\lambda) \cdot \Delta(\lambda) \quad (5)$$

where the summation covers the full spectral range of $S(\lambda)$, which is the normalized “Ultraviolet Hazard Function” (a normalized action spectrum for photokeratitis). The maximum duration of exposure to stay within the limit t_{max} is determined by dividing the daily EL of $3 \text{ mJ} \cdot \text{cm}^{-2}$ by the measured effective irradiance to obtain the duration in seconds, as noted in Eq. (6):

$$t_{\text{max}} = \frac{3 \times 10^{-3} \text{ J} \cdot \text{cm}^{-2}}{E_{\text{eff}}} \quad (6)$$

In addition to the $S(\lambda)$ envelope action spectrum-based EL, there has always been one additional criterion to protect the lens and to limit the dose rate to both lens and the skin from very high irradiances.

Initially, this was based only upon a consideration to conservatively protect against thermal effects. This was later thought essential not only to protect against thermal damage, but to also hedge against a possible unknown photochemical damage in the UV-A to the lens.

Sharply defined photobiological action spectra apply to the ultraviolet hazard function $[S(\lambda)]$, the blue-light hazard function $[B(\lambda)]$, and the retinal thermal hazard function $[R(\lambda)]$. These are shown in Fig. 1.

Guidelines for the Visible

The ocular exposure limits for intense visible and IR radiation exposure of the eye (from incoherent radiation) are several, because they protect against either photochemical or thermal effects to the lens or retina.

The two primary hazards that must be assessed in evaluating an intense visible light source are (1) the photoretinitis (blue-light) hazard and (2) the retinal thermal hazard. Additionally, lenticular exposure in the near IR may be of concern. It is almost always true that light sources with a luminance less than 1 cd/cm² (10⁴ cd/m²) will not exceed the limits, and this is generally a maximal luminance for comfortable viewing. Although this luminance value is not considered a safety limit, and may not be sufficiently conservative for a violet LED, it is frequently provided as a quick check to determine the need for further hazard assessment.^{3,5}

The retinal thermal criteria, based upon the action spectrum $R(\lambda)$, applies to pulsed light sources and to intense sources. The longest viewing duration of potential concern is 10 s, because pupillary constriction and eye movements limit the added risk from greater exposure durations. The retinal thermal hazard EL is therefore not specified for longer durations.¹² Indeed, in 2008 ACGIH limited the changing value to 0.25 s. These retinal thermal limits are

$$\sum L_{\lambda} \cdot R(\lambda) \cdot \Delta\lambda \leq 5/\alpha \cdot t^{0.25} \text{ W} \cdot \text{cm}^{-2} \cdot \text{sr}^{-1}, \quad 1 \mu\text{s} < t < 10 \text{ s} \quad (\text{ICNIRP, 1997}) \quad (7)$$

$$\sum L_{\lambda} \cdot R(\lambda) \cdot \Delta\lambda \leq 640/t^{0.25} \text{ W} \cdot \text{cm}^{-2} \cdot \text{sr}^{-1}, \quad 1 \mu\text{s} < t < 0.625 \text{ ms} \quad (\text{ACGIH, 2008}) \quad (8a)$$

$$\sum L_{\lambda} \cdot R(\lambda) \cdot \Delta\lambda \leq 16/t^{0.75} \text{ W} \cdot \text{cm}^{-2} \cdot \text{sr}^{-1}, \quad 0.625 \text{ ms} < t < 0.25 \text{ s} \quad (\text{ACGIH, 2008}) \quad (8b)$$

$$\sum L_{\lambda} \cdot R(\lambda) \cdot \Delta\lambda \leq 45 \text{ W} \cdot \text{cm}^{-2} \cdot \text{sr}^{-1}, \quad t > 0.25 \text{ s} \quad (\text{ACGIH, 2008}) \quad (8c)$$

The blue-light photoretinitis hazard criteria were based upon the work of Ham et al.^{6,11,17} The limit for time t is expressed as a $B(\lambda)$ spectrally weighted radiance:

$$\sum L_{\lambda} \cdot B(\lambda) \cdot t \cdot \Delta\lambda \leq 100 \text{ J} \cdot \text{cm}^{-2} \cdot \text{sr}^{-1} \quad (8)$$

Applying IR Limits

There are two criteria in the near-IR region. To protect the lens against IR cataract, the EL that is applicable to IR-A and IR-B radiant energy (i.e., 780 to 3000 nm) specifies a maximal irradiance for continued exposure of 10 mW/cm² (average) over any 1000-s period, but not to exceed an irradiance of $1.8t^{-0.75}$ W/cm² (for times of exposure less than 1000 s). This is based upon the fact that IR cataract in workers appears only after lifetime exposure of 80 to 150 mW/cm.^{2,3,6}

The second IR EL is to protect against retinal thermal injury from low-luminance IR illumination sources. This EL is for very special applications, where near-IR illuminators are used for night surveillance applications or IR LEDs are used for illumination or signaling. These illuminators have a very low visual stimulus and therefore would permit lengthy ocular exposure with dilated pupils. Although the retinal thermal limit [based upon the $R(\lambda)$ function] for intense, visible, broadband sources is not provided for times greater than 0.25 to 10 s because of pupillary constriction and the like, the

retinal thermal hazard—for other than momentary viewing—will only realistically occur when the source can be comfortably viewed, and this is the intended application of this special-purpose EL. The IR illuminators are used for area illumination for nighttime security where it is desirable to limit light trespass to adjacent housing. If directly viewed, the typical illuminator source may be totally invisible, or it may appear as a deep cherry red source that can be comfortably viewed. The EL is proportional to $1/\alpha$ and is simply limited to:

$$L_{\text{NIR}} = \sum L_{\lambda} \cdot R(\lambda) \cdot \Delta\lambda \leq 3.2/\alpha \cdot t^{0.25} \text{ W} \cdot \text{cm}^{-2} \cdot \text{sr}^{-1} \quad t < 810 \text{ s} \quad (\text{ACGIH}, 2008) \quad (10a)$$

$$L_{\text{NIR}} = \sum L_{\lambda} \cdot R(\lambda) \cdot \Delta\lambda \leq 0.6/\alpha \text{ W} \cdot \text{cm}^{-2} \cdot \text{sr}^{-1} \quad t < 810 \text{ s} \quad (\text{ACGIH}, 2008) \quad (10b)$$

It must be emphasized that this criterion is not applied to white-light sources, because bright light produces an aversion response and Eq. (8) or (9) would apply instead.

7.8 DISCUSSION

Exceeding the Exposure Limits

When deriving the ELs for minimal-image-size visible and infrared laser radiations, the ICNIRP generally includes a factor of 10 to reduce the 50-percent probability of retinal injury by a factor of 10. This is not a true “safety factor,” because there is a statistical distribution of damage, and this factor was based upon several considerations. These included the difficulties in performing accurate measurements of source radiance or corneal irradiance, the measurement of the source angular subtense, as well as histological studies showing retinal changes occurring at the microscopic level at levels of approximately 2 below the ED-50 value.³ In actual practice, this means that an exposure at two to three times the EL would not be expected to actually cause a physical retinal injury. At five times the EL, one would expect to find some injuries in a population of exposed subjects. The ELs are guidelines for controlling human exposure and should not be considered as fine lines between safe and hazardous exposure. By employing benefit-versus-risk considerations, it would be appropriate to have some relaxed guidelines; however, to date, no standards group has seen the need to do this.

Laser Hazards

The very high radiance (brightness) of a laser (MW and $\text{TW cm}^{-2} \text{sr}^{-1}$) is responsible for the laser’s great value in material processing and laser surgery, but it also accounts for its significant hazard to the eye (Fig. 6). When compared with a xenon arc or the sun, even a small He–Ne alignment laser is typically 10 times brighter (Fig. 7). A collimated beam entering the relaxed human eye will experience an increased irradiance of about 10^5 (i.e., $1 \text{ W} \cdot \text{cm}^{-2}$ at the cornea becomes 100 kW/cm^2 at the retina). Of course, the retinal image size is only about 10 to 20 μm , considerably smaller than the diameter of a human hair. So you may wonder: “So what if I have such a small lesion in my retina? I have millions of cone cells in my retina.” The retinal injury is always larger because of heat flow and acoustic transients, and even a small disturbance of the retina can be significant. This is particularly important in the region of central vision, referred to by eye specialists as the *macula lutea* (yellow spot), or simply the *macula*. The central region of the macula, the fovea centralis, is responsible for your detailed 20/20 vision. Damage to this extremely small (about 150- μm diameter) central region can result in severe vision loss even though 98 percent of the retina is unscathed. The surrounding retina is useful for movement detection and other tasks but possesses limited visual acuity (after all, this is why your eye moves across a line of print, because your retinal area responsible for detailed vision has a very small angular subtense). Outside the retinal hazard region (400 to 1400 nm), the cornea—and even the lens—can be damaged by laser beam exposure.

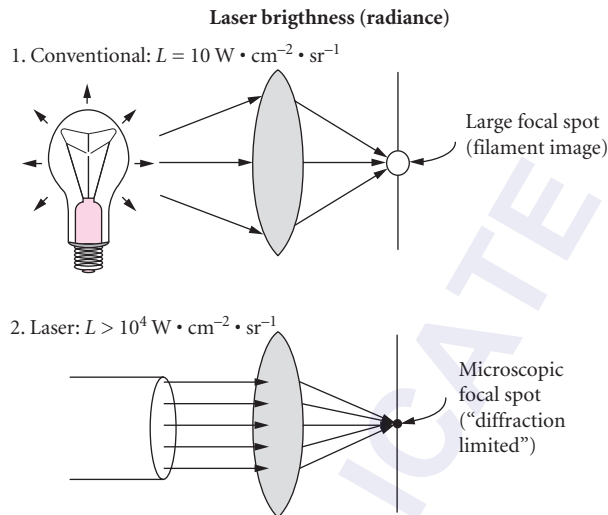


FIGURE 6 The radiance of a light source determines the irradiance in the focal spot. Hence, the very high radiance of a laser permits one to focus laser radiation to a very small image of very high irradiance.

Laser Safety Standards In the United States, the American National Standard, ANSI Z136.1-2000, *The Safe Use of Lasers*, is the national consensus standard for laser safety in the user environment. It evolved through several editions since 1973. Maximum permissible exposure (MPE) limits are provided as sliding scales, with wavelength and duration for all wavelengths from 180 nm to 1 mm and for exposure durations of 100 fs to 30 ks (8-h workday). Health and safety specialists want the simplest expression of the limits, but some more mathematically inclined scientists and engineers on standards committees argue for sophisticated formulas to express the limits (which belie the real level of biological uncertainty). These exposure limits (which are identical to those of ACGIH) formed the basis for the U.S. Federal Product Performance Standard (21 CFR 1040).³ The latter standard regulates only laser manufacturers. On the international scene, the International Electrotechnical Commission Standard IEC 60825-1 Ed. 2 (2007) originally grew out of an amalgam of the ANSI standard for user control measures and exposure limits and the U.S. Federal product classification regulation. The ICNIRP, which has a special relationship with the WHO in developing criteria documents on laser radiation, now recommends exposure limits for laser radiation.¹⁰ All of the aforementioned standards are basically in agreement.

Laser safety standards existing worldwide group all laser products into four general hazard classes and provide safe measures for each hazard class (classes 1 to 4).^{4,5} The U.S. Federal Product Performance Standard (21 CFR 1040) requires all commercial laser products to have a label indicating the hazard class. Once one understands and recognizes the associated ocular hazards, the safety measures recommended in these standards are quite obvious (e.g., beam blocks, shields, baffles, eye protectors). The ocular hazards are generally of primary concern. Many laser products sold in the USA and worldwide, will also be certified to meet the corresponding standard of IEC 60825-1:2007.

Laser Accidents A graduate student in a physical chemistry laboratory is aligning a Nd:YAG-pumped optical parametric oscillator (OPO) laser beam to direct it into a gas cell to study photodissociation parameters for a particular molecule. Leaning over a beam director, he glances down over an upward, secondary beam and approximately 80 μJ enters his left eye. The impact produces a microscopic hole in his retina, a small hemorrhage is produced over his central vision, and he sees only red in his left eye. Within an hour, he is rushed to an eye clinic where an ophthalmologist tells him

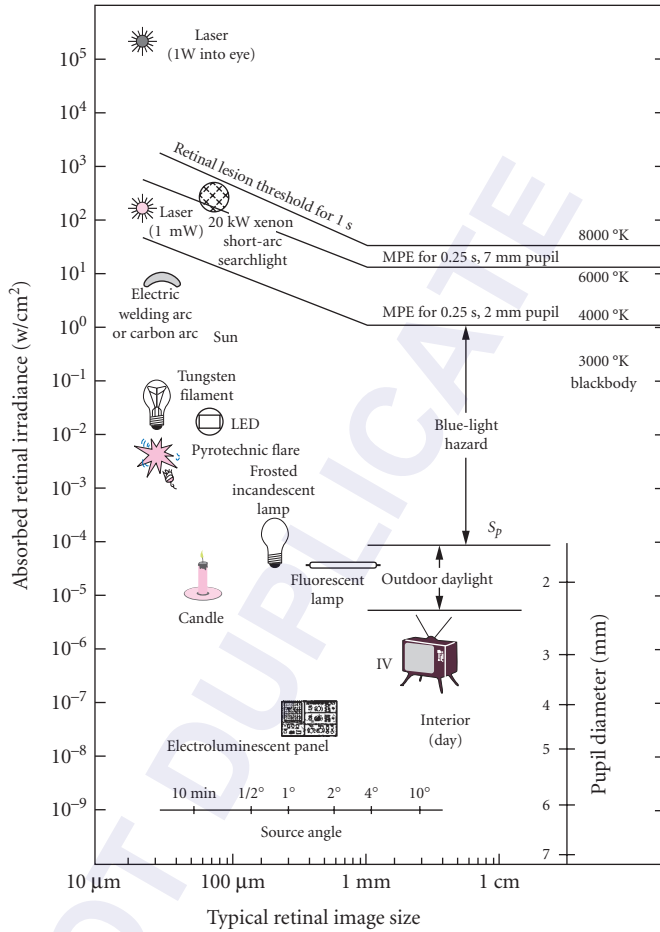


FIGURE 7 Relative retinal irradiances for staring directly at various light sources. Two retinal exposure risks are depicted: (1) retinal thermal injury, which is image-size dependent and (2) photochemical injury, which depends on the degree of blue light in the source’s spectrum. The horizontal scale indicates typical image sizes. Most intense sources are so small that eye movements and heat flow will spread the incident energy over a larger retinal area. The range of photoretinitis (the “blue-light hazard”) is shown to be extending from the normal outdoor light levels and above. The xenon arc lamp is clearly the most dangerous of nonlaser hazards.

he has only 20/400 vision. In another university, a physics graduate student attempts to realign the internal optics in a Q-switched Nd:YAG laser system—a procedure normally performed by a service representative that the student had witnessed several times before. A weak secondary beam reflected upward from a Brewster window enters the man’s eye, producing a similar hemorrhagic retinal lesion with a severe loss of vision. Similar accidents occur each year and frequently do not receive publicity because of litigation or for administrative reasons.^{3,26} Scientists and engineers who work with open-beam lasers really need to realize that almost all such lasers pose a very severe hazard to the eye if eye protection is not worn or if other safety measures are not observed!²⁻⁴

A common element in most laser laboratory accidents is an attitude that “I know where the laser beams are; I do not place my eye near to a beam; safety goggles are uncomfortable; therefore, I do not need to wear the goggles.” In virtually all accidents, eye protectors were available, but not worn. The probability that a small beam will intersect a 3- to 5-mm pupil of a person’s eye is small to begin with, so injuries do not always happen when eye protectors are not worn. However, it is worthwhile to consider the following analogy. If these individuals were given an air rifle with 100 BBs to fire, were placed in a cubical room that measures 4 m on each side and is constructed of stainless-steel walls, and were told to fire all of the BBs in any directions they wished, how many would be willing to do this without heavy clothing and eye protectors?!! Yet, the probability of an eye injury is similar. In all fairness, there are laser goggles that are comfortable to wear. The common complaints that one cannot see the beam to align it are mere excuses that are readily solved with some ingenuity once you accept the hazard. For example, image converters and various fluorescent cards (such as are used to align the Nd: YAG 1064-nm beam) can be used for visible lasers as well.

Laser Eye Protectors In the ANSI Z136.1 standard, there is broad general guidance for the user to consider factors such as comfort and fit, filter damage threshold, and periodic inspection, as well as the critical specification of wavelength and optical density (OD). However, there has never been any detailed U.S. specifications for standard marking and laboratory proofing (testing) of filters. By contrast, the approach in Germany (with DIN standards) and more recently in Europe (with CEN standards) has been to minimize the decision making by the user and place heavy responsibility upon the eyewear manufacturer to design, test, and follow standardized marking codes to label the eye protector. Indeed, the manufacturer had been required to use third-party test houses at some considerable expense to have type eyewear tested. Each approach has merit, and a new standard in the United States is now being developed for testing for standardized marking that is understandable to the wearer. The most important test of a protective filter is to measure the OD under CW, Q-switched, and mode-locked pulse irradiation conditions to detect saturable absorption-reversible bleaching.^{28,29}

The marking of laser eye protection in an intelligible fashion to ensure that the user will not misunderstand and select the wrong goggle has been a serious issue, and there appears to be no ideal solution. Several eye injuries appear to have been caused by a person choosing the wrong protector. This is particularly likely in an environment where multiple and different laser wavelengths are in use, as in a research laboratory or a dermatological laser setting. A marking of OD may be clearly intelligible to optical physicists, but this is seldom understandable to physicians and industrial workers. Probably the best assurance against misuse of eyewear has been the application of customized labeling by the user. For example, “Use only with the model 12A,” or “Use only for alignment of mode-locked YAG,” or “Use only for port-wine stain laser.” Such cautions supplement the more technical marking. The terms ruby, neodymium, carbon dioxide, and similar labels in addition to the wavelength can reduce potential confusion with a multi-wavelength laser. Marking of broadband eye protection, such as welding goggles, does not pose a problem to the wearer. The welding shade number (e.g., WG-8), even if misunderstood by the welder, poses little risk. The simple suggestion to select the filter based upon visual comfort prevents the welder from being injured; wearing too light a shade such that the BL for blue light is exceeded is virtually impossible because of the severe disability glare that would result. The same is not true for laser eye protectors: One could wear a goggle that protects against a visible wavelength but that doesn’t protect against the invisible 1064-nm wavelength. Improved, standardized marking is certainly needed! Recently, a new standard in the ANSI Z136 series of standards was published: ANSI Z136.7:2008 (American National Standard for Testing and Labeling of Laser Protective Equipment), which provides more reasonable guidance.²⁸

Lamp Safety Standards

Lamp safety is dealt with in ANSI/IESNA RP27.3-2005 [*Recommended Practice for Photobiological Safety for Lamps—Risk Group Classification & Labeling*, and in CIE Standard S-009:2006(IEC62471:2006) *Photobiological Safety for Lamps and Lamp Systems*. In addition, LEDs are also included in one laser safety standard: *American National Standard for the Safe Use of Lasers in Optical Fiber Communications Systems (OFCS)*], ANSI Z136.2 (1997), since LEDs used in this special application may be viewed by

optical aids such as a microscope or eye loupe, and such viewing actually happens in realistic splicing conditions and the instrumentation to distinguish the narrow bandwidth of the laser from the wider, nominal 50-nm bandwidth of the LED is a needless burden. All current LEDs have a maximum radiance of about $12 \text{ W} \cdot \text{cm}^{-2} \cdot \text{sr}^{-1}$ and can in no way be considered hazardous, but for a period an IEC laser standard 60825-1 included LEDs (an inclusion dropped in 2007). Some observers of the IEC group accused the group of having an ulterior motive of providing test houses and safety delegates with greater business!

7.9 REFERENCES

1. L. R. Solon, R. Aronson, and G. Gould, "Physiological Implications of Laser Beams," *Science* **134**:1506–1508 (1961).
2. D. H. Sliney and B. C. Freasier, "The Evaluation of Optical Radiation Hazards," *Applied Optics* **12**(1):1–24 (1973).
3. D. H. Sliney and M. L. Wolbarsht, *Safety with Lasers and Other Optical Sources*. Plenum Publishing, New York, 1980.
4. World Health Organization (WHO), *Environmental Health Criteria No. 23, Lasers and Optical Radiation*, joint publication of the United Nations Environmental Program, the International Radiation Protection Association, and the World Health Organization, Geneva, 1982.
5. American Conference of Governmental Industrial Hygienists (ACGIH), *TLVs and BEIs Based on the Documentation of the Threshold Limit Values for Chemical Substances and Physical Agents and Biological Exposure Indices*, American Conference of Governmental Industrial Hygienists, Cincinnati, OH, 2008.
6. ACGIH, *Documentation for the Threshold Limit Values*, American Conference of Governmental Industrial Hygienists, Cincinnati, OH, 2007.
7. A. S. Duchene, J. R. A. Lakey, and Michael H. Repacholi (eds.), *IRPA Guidelines on Protection Against Non-Ionizing Radiation*, MacMillan, New York, 1991.
8. D. H. Sliney, J. Mellerio, V. P. Gabel, and K. Schulmeister, "What is the Meaning of Threshold in Laser Injury Experiments? Implications for Human Exposure Limits," *Health Phys.*, **82**(3):335–347, 2002.
9. American National Standards Institute (ANSI), *Safe Use of Lasers*, ANSI Z136.1-2007, Laser Institute of America, Orlando, FL, 2007.
10. International Commission on Non-Ionizing Radiation Protection (ICNIRP), "Guidelines on Limits for Laser Radiation of Wavelengths between 180 nm and 1,000 μm ," *Health Phys.* **71**(5):804–819 (1996); update (in press).
11. ICNIRP, "Guidelines on limits of Exposure to Broad-Band Incoherent Optical Radiation (0.38 to 3 μm)," *Health Phys.* **73**(3):539–554 (1997).
12. W. P. Roach et al., "Proposed Maximum Permissible Exposure Limits for Ultrashort Laser Pulses," *Health Phys.* **77**:61–68 (1999).
13. T. P. Coohill, "Photobiological Action Spectra—What Do They Mean?" *Measurements of Optical Radiation Hazards*, CIE/ICNIRP, Munich, pp. 27–39 (1998).
14. D. H. Sliney, "Radiometric Quantities and Units Used in Photobiology and Photochemistry: Recommendations of the Commission Internationale de l'Éclairage (International Commission on Illumination)," *Photochem. Photobiol.* **83**:425–432 (2007).
15. J. A. Zuclich, "Ultraviolet-Induced Photochemical Damage in Ocular Tissues," *Health Phys.* **56**(5):671–682 (1989).
16. D. G. Pitts et al., "Ocular Effects of Ultraviolet Radiation from 295 to 365 nm," *Invest. Ophthalmol. Vis. Sci.* **16**(10):932–939 (1977).
17. W. T. Ham, Jr., "The Photopathology and Nature of the Blue-Light and Near-UV Retinal Lesions Produced by Lasers and Other Optical Sources," *Laser Applications in Medicine and Biology*, M. L. Wolbarsht (ed.), Plenum Publishing, New York, 1989.
18. W. T. Ham et al., "Evaluation of Retinal Exposures from Repetitively Pulsed and Scanning Lasers," *Health Phys.* **54**(3):337–344 (1988).
19. D. H. Sliney, "Physical Factors in Cataractogenesis-Ambient Ultraviolet Radiation and Temperature," *Invest. Ophthalmol. Vis. Sci.* **27**(5):781–789 (1986).

20. Commission Internationale de l'éclairage (International Commission on Illumination), *International Lighting Vocabulary*, CIE Publication No. 17.4, CIE, Geneva, Switzerland (1987).
21. M. T. Coroneo et al., "Peripheral Light Focussing by the Anterior Eye and the Ophthalmohelioses," *Ophthalmic Surg.* **22**:705–711 (1991).
22. B. E. K. Klein, R. Klein, and K. L. Linton, "Prevalence of Age-Related Lens Opacities in a Population, the Beaver Dam Eye Study," *Am. J. Pub. Hlth.* **82**(12):1658–1662 (1992).
23. D. H. Sliney, "Eye Protective Techniques for Bright Light," *Ophthalmology* **90**(8):937–944 (1983).
24. P. J. Dolin, "Assessment of the Epidemiological Evidence that Exposure to Solar Ultraviolet Radiation Causes Cataract," *Doc. Ophthalmol.* **88**:327–337 (1995).
25. D. H. Sliney, "Geometrical Gradients in the Distribution of Temperature and Absorbed Ultraviolet Radiation in Ocular Tissues," *Dev. Ophthalmol.* **35**:40–59 (2002).
26. D. H. Sliney, "Ocular Injuries from Laser Accidents," *SPIE Proceedings of Laser-Inflicted Eye Injuries: Epidemiology, Prevention, and Treatment*, San Jose, CA, 1996, pp. 25–33.
27. J. W. Ness et al., "Retinal Image Motion during Deliberate Fixation: Implications to Laser Safety for Long Duration Viewing," *Health Phys.* **72**(2):131–142 (2000).
28. ANSI, "American National Standard for Testing and Labeling of Laser Protective Equipment," ANSI Z136.7-2008 (2008).
29. D. H. Sliney and W. J. Marshall (eds.), "LIA Guide for the Selection of Laser Eye Protection," *Laser Institute of America*, 2008.

BIOLOGICAL WAVEGUIDES

Vasudevan Lakshminarayanan

*School of Optometry and Departments of Physics and Electrical Engineering
University of Waterloo
Waterloo, Ontario, Canada*

Jay M. Enoch

*School of Optometry
University of California at Berkeley
Berkeley, California*

8.1 GLOSSARY

Apodization. Modification of the pupil function. In the present context, the Stiles-Crawford directional effect can be considered as an apodization, that is, variable transmittance across the pupil.

Bessel functions. These special functions, first defined by the mathematician Daniel Bernoulli and generalized by Friedrich Bessel, are canonical solutions of a particular differential equation called the Bessel differential equation. Bessel's equation arises when finding separable solutions to Laplace's equation and the Helmholtz equation in cylindrical or spherical coordinates systems. Bessel functions are therefore especially important for many problems of wave propagation, static potentials, such as propagation of electromagnetic waves in cylindrical waveguides, heat conduction in cylindrical objects, and the like. A specific form of the Bessel functions are called Hankel functions. These functions are very useful to describe waveguiding.

Born approximation. If an electromagnetic wave is expressed as the sum of an incident wave and the diffracted secondary wave, the scattering of the secondary wave is neglected. This neglect represents what is known as Born's first-order approximation.

Bruch's membrane. Inner layer of the choroid.

Cilia. Hair.

Cochlea. The spiral half of the labyrinth of the inner ear.

Coleoptile. The first true leaf of a monocotyledon.

Copepod. A subclass of crustacean arthropods.

Cotyledons. The seed leaf—the leaf or leaves that first appear when the seed germinates.

Ellipsoid. The outer portion of the inner segment of a photoreceptor (cone or rod). It is located between the myoid and the outer segment. It often contains oriented mitochondria for metabolic activity.

Fluence. A measure of the time-integrated energy flux usually given in units of Joules/cm². In biology, it is used as a measure of light exposure.

Helmholtz's reciprocity theorem. Also known as the reversion theorem. It is a basic statement of the reversibility of optical path.

Henle's fibers. Slender, relatively uniform fibers surrounded by more lucent Müller cell cytoplasm in the retina. They arise from cone photoreceptor bodies and expand into cone pedicles.

Hypocotyl. The part of the axis of the plant embryo that lies below the cotyledons.

Interstitial matrix. Space separating photoreceptors. It is of a lower refractive index than the photoreceptors and, hence, can be considered as "cladding."

Mesocotyl. The node between the sheath and the cotyledons of seedling grasses.

Mode. The mode of a dielectric waveguide is an electromagnetic field that propagates along the waveguide axis with a well-defined phase velocity and that has the same shape in any arbitrary plane transverse along that axis. Modes are obtained by solving the source-free Maxwell's equations with suitable boundary conditions.

Myoid. It is the inner portion of the inner segment of a photoreceptor (cone or rod), located between the ellipsoid and the external limiting membrane of the retina.

Organ of Corti. The structure in the cochlea of the mammals which is the principal receiver of sound; contains hair cells.

Pedicle. Foot of the cone photoreceptor; it is attached to a narrow support structure.

Phototropism. Movement stimulated by light.

Spherule. End bulb of a rod photoreceptor.

V-parameter. Waveguide parameter [Eq. (3)]. This parameter completely describes the optical properties of a waveguide and depends upon the diameter of the guide, the wavelength of light used, and the indices of refraction of the inside (core) and outside (cladding) of the guide.

Equation (1). This is a statement of Snell's law.

n_1 refractive index of medium 1; light is incident from this medium

n_2 refractive index of medium 2; light is refracted into this medium

θ_1 angle of incidence of light measured with respect to the surface normal in medium 1

θ_2 angle of refraction of light in medium 2 measured with respect to the surface normal

Equation (2). This is an expression for the numerical aperture.

n_1 refractive index of the inside of the dielectric cylinder (the "core")

n_2 refractive index of the surround (the "cladding")

n_3 refractive index of medium light is incident from before entering the cylinder

θ_L limiting angle of incidence

Equation (3). Expression for the waveguide V-parameter.

d diameter of the waveguide cylinder

λ wavelength of incident light

n_1 refractive index of core

n_2 refractive index of cladding

8.2 INTRODUCTION

In this chapter, we will be studying two different types of biological waveguide model systems in order to relate the models' structures to the waveguide properties initiated within the light-guide construct. In these models, we will be discussing how these biological waveguides follow the Stiles-Crawford effect of the first kind, wavelength sensitivity of transmission, and Helmholtz's reciprocity theorem of optics. Furthermore, we will be investigating the different sources of biological waveguides such as vertebrate photoreceptors, cochlear hair cells (similar to cilia), and fiber-optic plant tissues, as well as considering their applications and comparing their light-guiding effects to the theoretical models. The emphasis will be on waveguiding in vertebrate retinal photoreceptors, though similar considerations apply to other systems.

8.3 WAVEGUIDING IN RETINAL PHOTORECEPTORS AND THE STILES-CRAWFORD EFFECT

Photoreceptor optics is defined as the science that investigates the effects of the optical properties of the retinal photoreceptors—namely, their size, shape, refractive index, orientation, and arrangement—on the absorption of light by the photopigment, as well as the techniques and instrumentation necessary to carry out such studies. It also explores the physiological consequences of the propagation of light within photoreceptor cells. The *Stiles-Crawford effect of the first kind*¹ (SCE I), discovered in 1933, represents a major breakthrough in our understanding of retinal physiology and the modern origin of the science of photoreceptor optics. The SCE I refers to the fact that visual sensitivity in normal eyes is greatest for light entering near the center of the eye pupil, and the response falls off roughly symmetrically from this peak. The SCE I underscores the fact that the retina is a complex optical processing system whose properties play a fundamental role in visual processing. The individual photoreceptors (rods and cones) behave as light collectors which capture the incident light and channel the electromagnetic energy to sites of visual absorption, the photolabile pigments, where the transduction of light energy to physicochemical energy takes place. The photoreceptors act as classic fiber-optic elements and the retina can be thought of as an enormous fiber bundle (a typical human retina has about 130×10^6 receptors). Some aspects of the behavior of this fiber optic bundle can be studied using the SCE I function, a psychophysical measurement of the directional sensitivity of the retina. The SCE I is an important index reflecting waveguide properties of the retina and is a psychophysical measurement of the directional sensitivity of the retina. The SCE I can be used to elucidate various properties, and the photodynamic nature, of the photoreceptors. The SCE I can also be used to indicate the stage and degree of various retinal abnormalities. The reader is referred to the literature for detailed discussion of various aspects of photoreceptor optics.²⁻⁵ A collection of classic reprints on the Stiles-Crawford effect and photoreceptor optics has been published by the Optical Society of America.⁶ The books by Dowling⁷ and Rodieck,⁸ for example, give good basic information on the physiology of the retina as well as discussions of the physical and chemical events involved in the early stages of visual processing. Additionally, there is a color effect, known as the Stiles-Crawford effect of the second kind (SCE II). Here, if the locus of entry of a beam in the entrance pupil of the eye is changed, it altered the perceived hue and saturation of the retinal image. The SCE II will not be dealt with in this chapter and the reader is referred to the literature.^{9,10} Added features of the SCE are addressed in the following chapter.

8.4 WAVEGUIDES AND PHOTORECEPTORS

It appears as though there are two major evolutionary lines associated with the development of visual receptors: the ciliary type (vertebrates) and the rhabdomeric type (invertebrates). Besides being different in structure, the primary excitation processes also differ between these two types of detectors. Ciliary type (modified hair) detectors are found in all vertebrate eyes. They have certain common features with other sensory systems, for example, the hair cells of the auditory and vestibular systems. It is possible to say that all vertebrate visual detectors have waveguides associated with the receptor element. The incident light is literally guided or channeled by the receptor into the outer segment where the photosensitive pigments are located. In other words, the fiber transmits radiant energy from one point (the entrance pupil or effective aperture of the receptor) to other points (the photolabile pigment transduction sites). Given the myriad characteristics found in photoreceptors in both vertebrate and invertebrate eyes (including numerous superposition and apposition eyes), what is common in all species is the presence of the fiber-optic element. It appears as though fiber-optic properties evolved twice in the vertebrate and invertebrate receptors. This fact emphasizes the importance of these fiber-optic properties.

Like any optical device, the photoreceptor as a waveguide can accept and contain light incident within a solid angle about its axis. Therefore, all waveguides are directional in their acceptance of light and selection of only a certain part of the incident electromagnetic wavefront for transmittance to the transducing pigment.

The detailed analysis of optics and image formation in various forms of vertebrate and invertebrate eyes is beyond the scope of this chapter (see Refs. 11 and 12 as well as Chap. 1 in this volume for excellent reviews). However, all such eyes must share a common characteristic. That is, each detecting element must effectively collect light falling within its bound and each must “view” a slightly different aspect of the external world, the source of visual signal. Because of limitations of ocular imagery, there must be a degree of overlap between excitation falling on neighboring receptors; there is also the possibility of optical cross talk. It has been postulated that there is an isomorphic correspondence between points on the retina and a direction in (or points in) real (the outside world) space. A number of primary and secondary sources of light stimulus for the organism are located in the outside world. The organism has to locate those stimuli falling within its field of view and to relate those sensed objects to its own egocentric (localization) frame of reference. Let us assume that there is a fixed relationship between sense of direction and retinal locus. The organism not only needs to differentiate between different directions in an orderly manner, but it must also localize objects in space. This implies that a pertinent incident visual signal is identifiable and used for later visual processing. Further, directionality-sensitive rods and cones function most efficiently only if they are axially aligned with the aperture of the eye, that is, the source of the visual signal. Thus, it becomes necessary to relate the aperture of the receptor as a fiber-optic element relative to the iris aperture. The retinal receptor-waveguide-pupillary aperture system should be thought of as an integrated unit that is designed for optimization of light capture of quanta from the visual signal in external space and for rejection of stray light noise in the eye/image. A most interesting lens-aperture-photoreceptor system in the invertebrate world is that of a copepod, *Copilia*. The *Copilia* has a pair of image-forming eyes containing two lenses, a corneal lens and another lens cylinder (with an associated photoreceptor), separated by a relatively enormous distance. Otherwise, the photoreceptor system is similar to the insect eye. *Copilia* scans the environment by moving the lens cylinder and the attached photoreceptor. The rate of scan varies from about five scans per second to one scan per two seconds. This is like a mechanical television camera. There are other examples of scanning eyes in nature.¹²

The retina in the vertebrate eye is contained within the white, translucent, diffusing, integrating-sphere-like eye. It is important to realize that the retinal receptor “field of view” is not at all limited by the pupillary aperture. The nonimaged light may approach the receptor from virtually any direction. In Fig. 1 the larger of the two converging cones of electromagnetic energy is meant to portray (in a

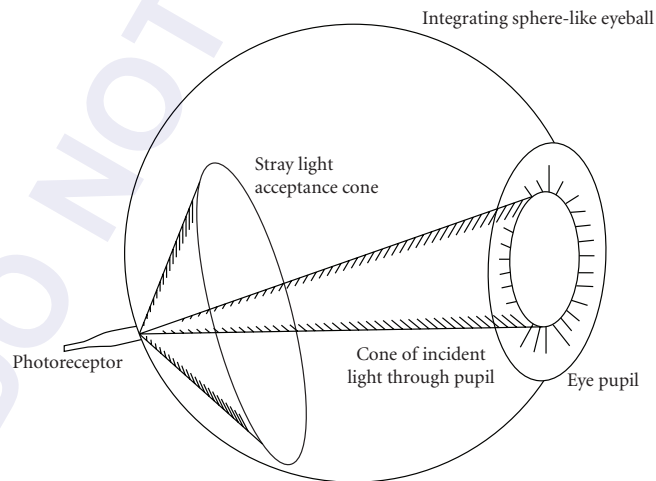


FIGURE 1 Schematic drawing showing light incidence at a single photoreceptor. Light incident on the retina has a maximum angle of incidence of about $\pm 10^\circ$ for a receptor at the posterior pole of the eye. Stray light may enter the receptor over a much larger acceptance angle (shown by the larger cone).

limited sense) the extent of the solid angle over which energy may impinge on the receptor. While stray light from any small part of this vast solid angle may be small, the total, when integrated over the sphere surrounding the receptor, may be quite a significant sum. To prevent such noise from destroying image quality, the organism has evolved specific mechanisms, including screening dark pigments in the choroid and pigment epithelium, morphological designs, photomechanical changes, aligned photolabile pigment and fiber-optic properties, and in a number of species, a tapetum.

Since the pupillary aperture is the source of the pertinent visual signal, the primary purpose of the retinal fiber-optic element, the receptor, is to be a directionally selective mechanism, which accepts incident light from only within a rather narrow solid angle. Given that the receptor has a rather narrow, limiting aperture, exhibits directionality, and favors absorption of light passing along its axis, this optical system will be most effective if it is aligned with the source of the signal, the pupillary aperture. This argument applies equally for cones and rods.

In order for a material to serve as a fiber-optic element, the material must be reasonably transparent and have a refractive index that is higher than the immediately surrounding medium. To reduce cross talk in an array of such fiber-optic elements, they must be separated by a distance approximately equal to or greater than a wavelength of light by a medium having a lower index of refraction than the fiber. The inner and outer segments of retinal photoreceptors and the separating interstitial matrix satisfy such condition. The index of refraction of biological materials is generally dependent on the concentration of large molecules of proteins, lipids, and lipoproteins. The receptor has a high concentration of such solids. These materials are in low concentrations in the surrounding interstitial matrix space. As examples, Sidman¹³ has reported refractive index values of 1.4 for rod outer segments, 1.419 for foveal cone outer segments, and a value between 1.334 and 1.347 for the interstitial matrix in the primate retina. Enoch and Tobey¹⁴ found a refractive index difference of 0.06 between the rod outer segment and the interstitial matrix in frogs. Photoreceptors are generally separated by about 0.51 μm . A minimum of about 0.5 μm (1 wavelength in the green part of the spectrum) is needed. Separation may be disturbed in the presence of pathological processes. This can alter waveguiding effects, transmission as a function of wavelength, directionality, cross talk, stray light, resolution capability, and the like.

8.5 PHOTORECEPTOR ORIENTATION AND ALIGNMENT

As noted above, Stiles and Crawford¹ found that radiant energy entering the periphery of the dilated pupil was a less effective stimulus than a physically equal beam entering the pupil center. Under photopic foveal viewing conditions, it was found that a beam entering the periphery of a pupil had to have a radiance of from 5 to 10 times that of a beam entering the pupil center to have the same subjective brightness. Stiles and Crawford plotted a parameter η , defined as the ratio of the luminance of the standard beam (at pupil center) to that of the displaced beam at photometric match point as a function of the entry portion of the displaced beam. The resultant curve has essentially a symmetric falloff in sensitivity (Fig. 2). This cannot be explained as being due to preretinal factors and this directional sensitivity of the retina is known as the Stiles-Crawford Effect of the first kind (SCE I). Later a chromatic effect (alteration of perceived hue and saturation of displaced beam) was also discovered that is called the *Stiles-Crawford Effect of the second kind*.^{9,10} SCE I is unaffected by polarized light. Normative values of the parameters of this function describing the effect (as well as various mathematical descriptions) can be found in Applegate and Lakshminarayanan.¹⁵

The directionality of photoreceptors depends upon the acceptance angles of individual photoreceptors and on the variability of their orientation within the tested photoreceptor population. MacLeod¹⁶ developed a selective adaptation technique to study the variability of photoreceptor orientations and concluded that foveal cones were aligned with great precision (similar results were shown, for example, Burns et al.¹⁷ using reflectometry methods). Roorda and Williams¹⁸ measured the orientations of individual cones and found the average disarray to be 0.17 mm in the pupil plane and point out that this disarray accounts for less than 1 percent of the breadth of the overall tuning function.

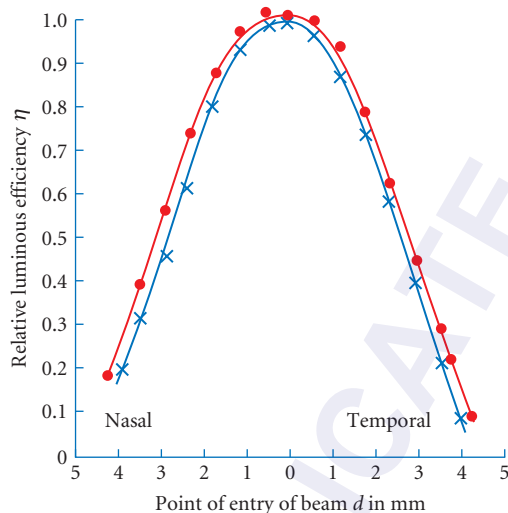


FIGURE 2 Psychophysically measured Stiles-Crawford function for left and right eyes of an observer. Nasal and temporal refer to direction with respect to pupil center. The observer's task is to match the brightness of a beam whose point of entry is varied in the pupil to a fixed beam coming in through the center of the pupil.

The directionality of photoreceptors is due to their structure that makes them act as optical fibers. Therefore, measurement of photoreceptor directionality is an important tool for testing the physical properties of photoreceptors *in vivo*. Much valuable information has been obtained from studies of patients with pathologies or conditions that affect the photoreceptor (and adjoining structures) layers. However clinical studies are limited by the fact that psychophysical methods are time consuming and require excellent cooperation from the subject. Therefore, different reflectometric methods have been developed.

Krauskopf¹⁹ was the first to show a direct correspondence between the psychophysical SCE and changes in reflectivity when the entry and exit pupils are moved. There are three major methods to measure the directional properties of photoreceptors by reflectometry. They are: (a) Moving entry and exit pupils method—here the angles of incidence and reflection at the retina are changed in tandem. This method has been applied to design the photoreceptor alignment reflectometer²⁰ and a custom built Scanning Laser Ophthalmoscope.²¹ (b) Moving exit pupil method—in this method the entrance pupil remains fixed and this design was applied to the reflectometers of van Blokland^{22,23} and Burns et al.²⁴ (c) Moving entrance pupil method²⁵—Roorda and Williams¹⁸ applied this method to sample fields as small as a single cone by means of the Rochester adaptive optics ophthalmoscope. (See Chap. 15 in this volume for a discussion of adaptive optics in vision science).

Histological and x-ray diffraction studies show that in many species, including human, photoreceptors tend to be oriented not toward the center of the eye, but toward an anterior point. This pointing tendency is present at birth. The directional characteristics of an individual receptor are difficult to assess, and cannot, as yet, be directly related to psychophysical SCE I data. If one assumes that there is a retinal mechanism controlling alignment and that the peak of the SCE I represents the central alignment tendencies of the receptors contained within the sampling area under consideration, it is possible to investigate where the peak of the receptor distribution is pointed at different loci tested across the retina and to draw conclusions relative to the properties of an alignment system or mechanism. Results of the SCE I functions studied over a wide area of the retina show that the overall pattern of receptor orientation is toward the approximate center of the exit pupil of the eye.

This alignment is maintained even up to 35° in the periphery. There is a growing evidence that receptor orientation is extremely stable over time in normal observers.²⁶

Photoreceptor orientation disturbance and realignment of receptors following retinal disorders have also been extensively studied. A key result of these studies is that with remission of the pathological process, the system is often capable of recovering orientation throughout life. Another crucial conclusion from these studies is that alignment is *locally controlled* in the retina (e.g., Ref. 27). The stability of the system, as well as corrections in receptor orientation after pathology, implies the presence of at least one mechanism (possibly phototropic) for alignment. Even though the center-of-the-exit-pupil-of-the-eye receptor alignment characteristic is predominant, rare exceptions to this have been found. These individuals exhibit an approximate center-of-the-retinal-sphere-pointing receptor alignment tendency (Fig. 3).²⁸ Though certain working hypotheses have been made,²⁹ what maintains the receptor is an open unanswered question and is the subject of current research that is beyond the scope of this chapter. Eckmiller,³⁰ for example, has proposed that the alignment of photoreceptors is accomplished by a feedback controlled bending of the cell at the myoid. Some unspecified local changes are thought to occur within the myoid that activate molecular motors and induce local movements by cytoskeletal elements, and bending can be accomplished by a differential change in the myoid at different positions around the cell perimeter. To summarize, the resultant directionality must reflect the properties of the media, the waveguide properties of receptors, the alignment properties of photolabile pigments, and so forth. Mechanical tractional effects also affect receptor alignment (e.g., Refs. 31 and 32). Other factors could include stress due to microscars, traction or shear forces due to accommodation, inertial forces during saccades, local growth for example, in high myopia, and the like. The SCE I function obtained experimentally depends upon all these factors as well as the sum of photoreceptor alignment properties sampled in the retinal test area. The response is dominated by the receptor units most capable of responding in the test area and is affected by the level of light adaptation.

What are the functional visual consequences of SCE I? Recent studies have dealt with this question (e.g., Ref. 33). SCE is important in photometry and retinal image quality (see Chap. 9). Calculations for change in effective retinal illuminance were first given by Martin³⁴ and expanded on by Atchison et al.³⁵ Baron and Enoch,³⁶ used a half-sensitivity, half-width measurement of retinal directional sensitivity as the basis for integrating a parabolic approximation of the SCE over the pupillary area. It is possible to define an “effective” troland which takes into account the SCE compensated entrance pupil area. These studies show that increasing the pupil size from 2- to 8-mm diameter provides an increase in effective retinal illuminance of about 9 times rather than the 16 times increase in pupil area.

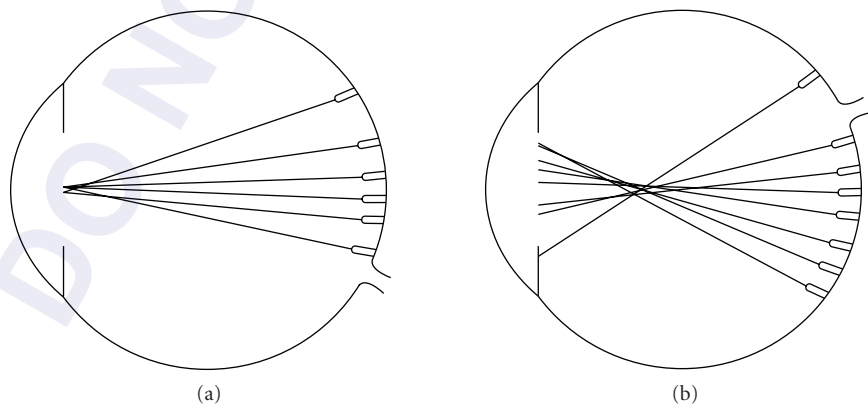


FIGURE 3 Inferred photoreceptor alignment from SCE measurements. (a) The more usual center-of-the-exit-pupil pointing tendency. (b) Center-of-the-retinal-sphere pointing tendency (rare).

The influence of the SCE on visual functions can be investigated by using filters based on the apodization model of the SCE.^{37,38} Rynders et al.³⁹ developed SCE neutralizing filters, but provided few details. Following up, Scott et al.⁴⁰ have described construction of practical filters to achieve neutralization of the SCE. When light from the edge of a large pupil plays a significant role in degrading retinal image quality, attenuating this light by the effective SCE I apodization will, of course, improve image quality. In general, the magnitude of this improvement will depend on the magnitude of the SCE and the level of blur caused by marginal rays. The presence of aberrations influences the impact of the SCE apodization on defocused image quality. Since defocus of the opposite sign to that of the aberration can lead to well-focused marginal rays, attenuating these well-focused rays with apodization will not improve image quality.⁴¹ In addition to generally attenuating the impact of defocus on image contrast, this apodization will effectively remove some of the phase reversals (due to zero crossings of the modulation transfer function) created by positive defocus. The displacement of phase reversals to higher spatial frequencies will have a direct impact on defocused visual resolution. Spatial tasks in which veridical phase perception is critical will be significantly more resistant to positive defocus because of the SCE. However, Atchison et al.^{41,42} suggest that SCE has little effect on image quality when dealing with well centered pupils. Zhang et al.³³ have used a wave optics model to examine the effects of SCE and aberrations and have compared the results with psychophysical measurements of the effect of defocus on contrast sensitivity and perceived phase reversals. They found that SCE apodization had the biggest effect on defocused image quality only when defocus and spherical aberration have the same sign and conclude that SCE can significantly improve defocused image quality and defocused vision particularly for tasks that require veridical phase perception. It has been argued that measurements of the *transverse chromatic aberration (TCA)* can be affected due to SCE. Various studies show that the effect of the SCE on the amount of TCA varies strongly across individuals, and between eyes in the same individual. In conclusion, it is thought that aberrations and SCE modify the amount and direction of TCA, and SCE does not necessarily reduce the impact of TCA.⁴³ Longitudinal chromatic aberration, on the other hand, is only slightly affected by SCE I. It has also been shown that decentering the SCE produces an appreciable shift in subjective TCA for large pupil sizes.⁴²

The wavelength dependence of the SCE has been explained by using a model of fundus reflectance developed by van de Kraats et al.⁴⁴ Berendschot et al.⁴⁵ showed that there is a good fit between the model (based on geometrical optics) and experimental data, if self-screening and backscattered choroidal light are included in the model.

Is the retina an optimal instrument? To answer this question, Marcos and Burns⁴⁶ studied both wavefront aberration and cone directionality. They concluded that cone directionality apodization does not always occur at the optically best pupillary region and that in general ocular optics and cone alignment do not develop toward an optimal optical design.

8.6 INTRODUCTION TO THE MODELS AND THEORETICAL IMPLICATIONS

First, we need to explore the theoretical constructs or models of the biological waveguide systems based on the photoreceptors to get a full appreciation of how these models transmit light through their structure. There are two models of biological waveguide systems that will be discussed here. The first model, which is the retinal layer of rods and cones, is defined as beginning approximately at the *external limiting membrane (ELM)* (Fig. 4), and it is from that boundary that the photoreceptors are believed to take their orientation relative to the pupillary aperture. In any local area within that layer, the photoreceptors are roughly parallel cylindrical structures of varying degrees of taper, which have diameters of the order of magnitude of the wavelength of light, and which are separated laterally by an organized lower refractive index substance, called the *interphotoreceptor matrix (IPM)*. The photoreceptors thus form a highly organized array of optical waveguides whose longitudinal axes are, beginning at the ELM, aligned with the central direction of the pupillary illumination. The IPM serves as a cladding for those waveguides. We will assume that the ELM approximates the inner bound of the receptor fiber bundle for light that is incident upon the retina in the physiological

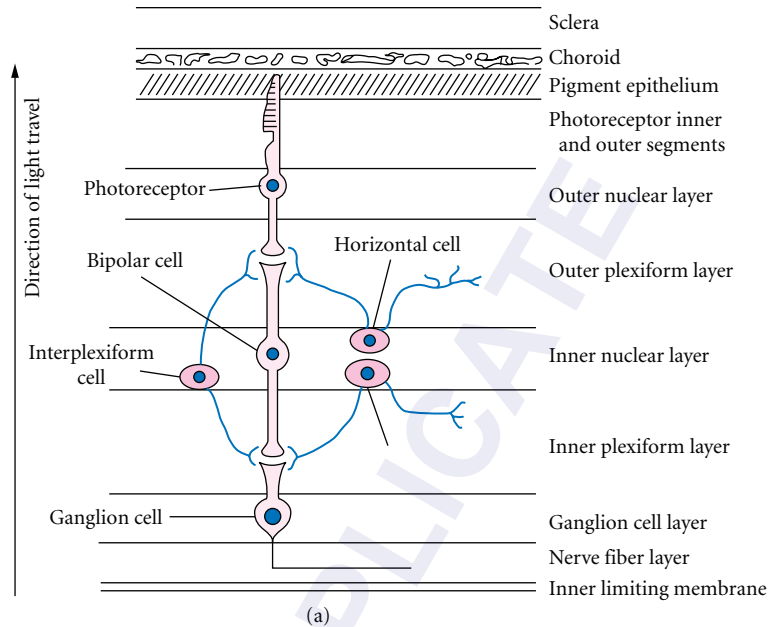


FIGURE 4 Schematic diagram of the retina and the photoreceptor. (a) Retinal layers. Note: the direction of light travel from the pupil is shown by the arrow.

direction.⁴⁷ Portions of the rod and cone cells that lie anterior to the ELM are not necessarily aligned with the pupillary aperture, and might be struck on the side, rather than on the axis, by the incident illumination. Some waveguiding is observed for a short distance into this region. However, these portions of the rod and cone cells anterior to the ELM would only be weakly excited by the incoming light.

Figure 5 shows the three-segment model that is generally used, denoting the idealized photoreceptor. All three segments are assumed to have a circular cross section, with uniform myoid and outer segment and a smoothly tapered ellipsoid. For receptors with equal myoid and outer segment radii, the ellipsoid is untapered. In this model, the sections are taken to be homogeneous, isotropic, and of higher refractive index than the surrounding near homogeneous and isotropic medium, which is the IPM. Both the myoid and the ellipsoid are assumed to be lossless, or nonabsorbing, since the cytochrome pigments in the mitochondria of the ellipsoids are not considered in this model. Absorption in the outer segments by the photolabile pigment molecules aligned within the closely packed disk membranes is described by the *Naperian absorption coefficient* (α). The aperture imposed at the ELM facilitates calculation of the equivalent illumination incident upon the receptor.

A number of assumptions and approximations have been made for these two models described above:

1. The receptor cross section is approximately circular and its axis is a straight line. This appears to be a reasonable assumption for the freshly excised tissue that is free of ocular diseases or pathology.
2. Ellipsoid taper is assumed to be smooth and gradual. Tapers that deviate from these conditions may introduce strong mode coupling and radiation loss.
3. Individual segments are assumed to be homogeneous. This is a first-order approximation, as each of the photoreceptor sections has different inclusions, which produce local inhomogeneities.

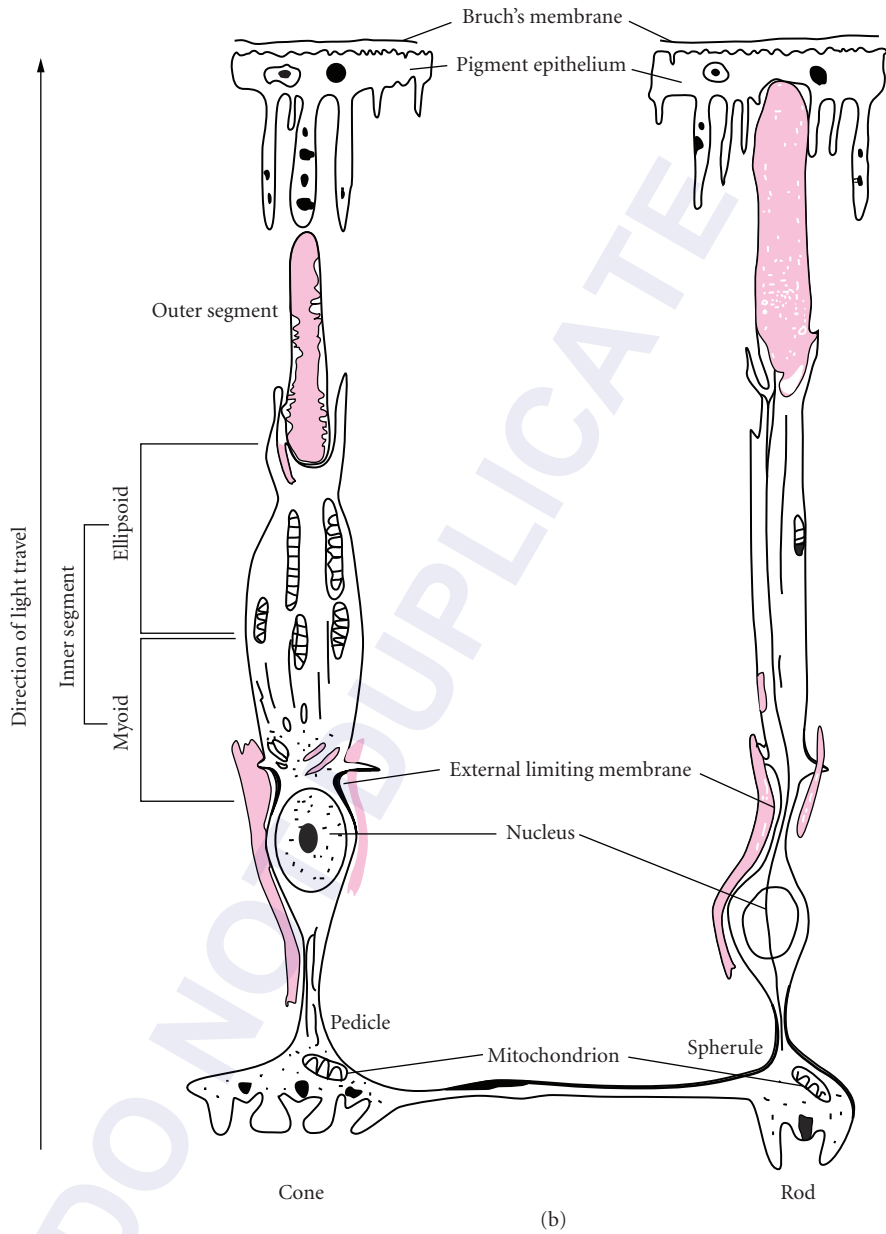


FIGURE 4 (Continued) (b) Structure of the cone and rod photoreceptor cells.

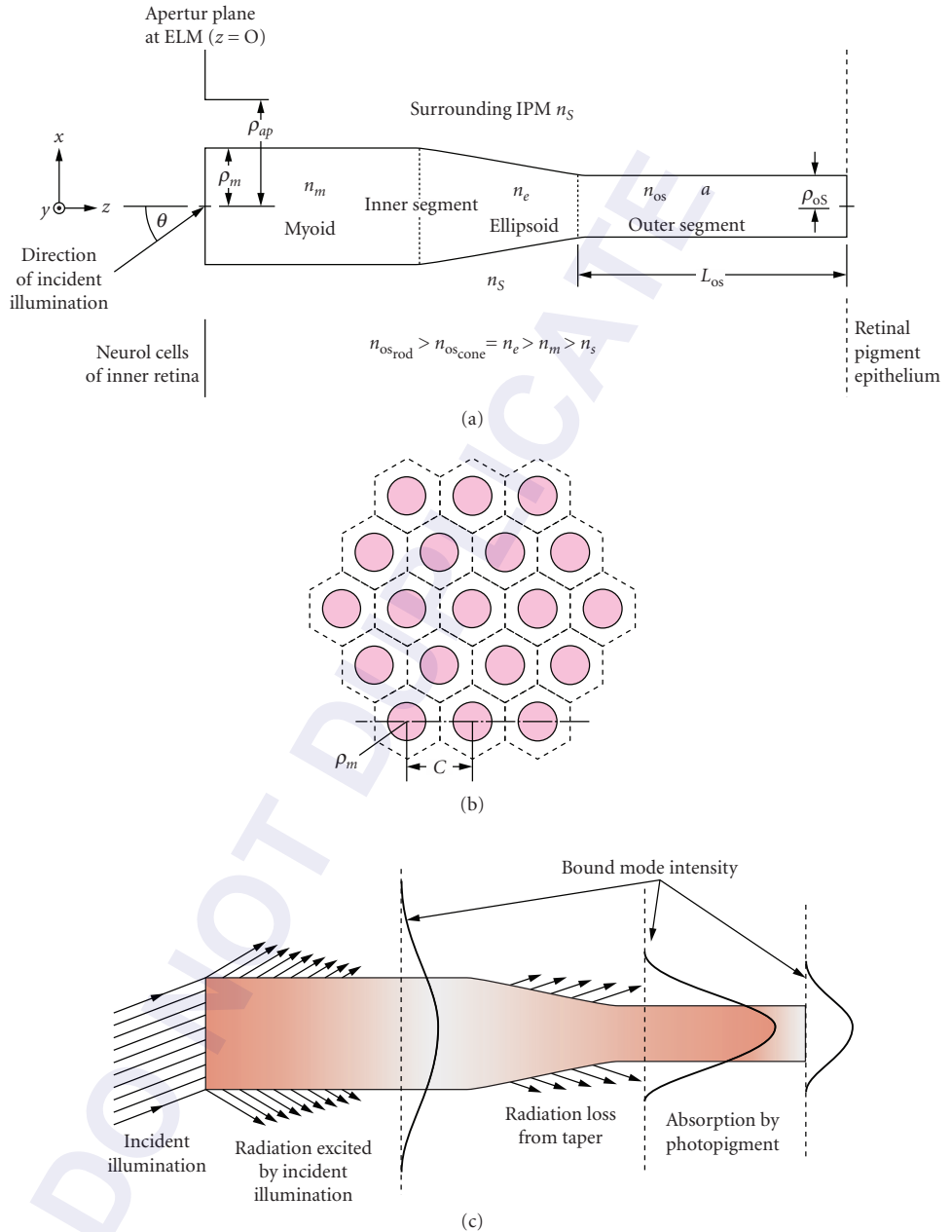


FIGURE 5 (a) Generalized three-segment model of an idealized photoreceptor. (b) View of retinal mosaic from the ELM (external limiting membrane), which is the aperture plane of the biological waveguide model for photoreceptors. The solid circles represent myoid cross sections. (c) Illustration of optical excitation of retinal photoreceptor showing radiation loss, funneling by ellipsoid taper, and absorption in the outer segment. Regions of darker shading represent greater light intensity.

4. Individual segments are assumed to be isotropic as a first-order approximation. However, in reality, the ellipsoid is packed with highly membranous mitochondria whose dimensions approach the wavelength of visible light, and which are shaped differently in different species. The mitochondrion-rich ellipsoid introduces light scattering and local inhomogeneity. The outer segment is composed of transversely stacked disks, having a repetition period an order of magnitude smaller than the wavelength of light; in addition, these outer segments contain photolabile pigment.
5. Linear media are assumed. Nonlinearities may arise from a number of sources. These could include the following:
 - Absorption of light is accompanied by bleaching of the photopigments and the production of absorbing photoproducts that have different absorption spectra from those of the pigments. Bleaching is energy and wavelength dependent and the bleached pigment is removed from the pool of available pigment, thus affecting the absorption spectrum. Under physiological conditions in the intact organism where pigment is constantly renewed, the relatively low illumination levels that result in a small percent bleach can be expected to be consistent with the assumed linearity. However, with very high illumination levels where renewal is virtually nonexistent, researchers must utilize special procedures to avoid obtaining distorted spectra (e.g., albino retina).
 - Nonlinearity may arise in the ellipsoid owing to illumination-dependent activity of the mitochondria, since it has been shown that mitochondria of ellipsoid exhibit metabolically linked mechanical effects. If such effects occurred in the retina, then illumination-driven activity could alter inner segment scattering.
 - Some species exhibit illumination-dependent photomechanical responses that result in gross movement of retinal photoreceptors. That movement is accompanied by an elongation or contraction of the myoid in the receptor inner segment and has been analyzed in terms of its action as an optical switch.⁴⁸ In affected species, rod myoids elongate and cone myoids shorten in the light, while rod myoids contract and the cone myoids elongate in the dark. In the light, then, the cones lie closer to the outer limiting membrane and away from the shielding effect of the pigment epithelial cell processes, and contain absorbed pigment, while the pigment granules migrate within the pigment epithelial processes.⁴⁹
6. The refractive indices of the segment are assumed to be independent of wavelength, which completely contradicts the definition of the *refractive index*, which is $n = cv$ where v (velocity) = f (frequency) $\cdot \lambda$ (wavelength), and v is proportional to wavelength because f is constant for all media. Therefore, a change in the velocity dependent on wavelength also changes the refractive index in real life.
7. The medium to the left of the external limiting membrane is assumed to be homogeneous. However, in reality, although the ocular media and inner (neural) retina are largely transparent, scattering does occur within the cornea, aqueous, lens, vitreous, and inner retina. Preparations of excised retina will be devoid of the ocular media, but the inner retina will be present. In turn, the preparations would account for light scattering upon light incidence.
8. It is assumed that the medium surrounding the photoreceptors is homogeneous and nonabsorbing. Again, this is an approximation. Microvilli originating in the Müller cells extend into the space between the inner segments, and microfibrils of the *retina pigment epithelium* (RPE) surround portions of the outer segments. Also, the assumption of homogeneity of the surround neglects the reflection and backscatter produced by the RPE, choroid, and sclera. Note that preparations of excised retinæ that are used for observing waveguiding are usually (largely) devoid of interdigitating RPE.
9. In some species the outer segment is tapered. Gradual taper can be accommodated in the model. Sharp tapers, however, must be dealt with separately.
10. The exact location of the effective input plane to the retinal fiber optics bundle is not known; however, in early studies it was assumed that the ELM was the effective inner bound of the retinal fiber bundle. To date, though, universal agreement does not exist regarding the location of the effective input plane; some investigators assume it to be located at the outer segment level.

11. In species with double cones (e.g., goldfish), the two components may often not be of equal length or contain identical photolabile pigments.

Next, we will evaluate the electromagnetic validity of the preceding models based on the assumptions discussed above. An additional assumption is that meaningful results relating to the *in situ* properties of a retinal photoreceptor can be obtained by means of calculations performed upon a single photoreceptor model. Involved here are considerations regarding optical excitation, optical coupling between neighboring receptors, and the role of scattered light within the retina. These assumptions may be summarized as follows:

1. The illumination incident upon the single photoreceptor through the aperture in Fig. 5 is taken to be representative of that illumination available to an individual receptor located in the retinal mosaic. Figure 5*b* illustrates the cross section of a uniform receptor array as seen from the ELM. For a uniformly illuminated, infinitely extended array, the total illumination available to each receptor would be that falling on a single hexagonal area.
2. The open, or unbounded, transverse geometry of the dielectric waveguide supports two classes of modes, referred to as bound, guided, or trapped modes, and unbound, unguided, or radiation modes. Different behavior is exhibited by the two mode species. As their names suggest, the bound modes carry power along the waveguide, whereas the unbound or radiation modes carry power away from the guide into the radiation field (Fig. 5*c*). From a ray viewpoint, the bound modes are associated with rays that undergo total internal reflection and are trapped within the cylinder; unbound modes are associated with rays that refract or tunnel out of the cylinder. These modes are obtained as source-free solutions of Maxwell's equations and are represented in terms of Bessel and Hankel functions (for a cylindrical geometry). These mode shapes show a high degree of symmetry. A full discussion of waveguide theory can be found in Snyder and Love.⁵⁰
3. Coherent illumination of the photoreceptors is assumed. Here, the equivalent aperture of the photoreceptor is defined to be the circle whose area is equal to that of the hexagon seen in Fig. 5*b* specified by the appropriate intercellular spacing and inner segment diameter.

The earliest analyses of light propagation within the retinal photoreceptors employed geometric optics. As early as 1843, Brucke discussed the photoreceptor's trapping of light via total internal reflection (as quoted by von Helmholtz),⁵¹ but stopped just short of attributing a directionality to vision. In fact, Hannover had observed waveguide modes in photoreceptors,⁵² but interpreted them as being cellular fine structure (Fig. 6). Interest resumed after the discovery of the Stiles-Crawford effect of the

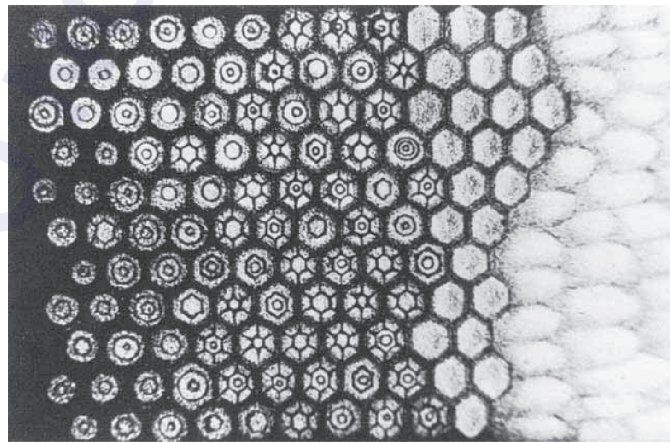


FIGURE 6 Drawing of frog outer segment viewed end-on. (From Ref. 21.)

first kind (SCE I). O'Brien⁵³ suggested a light-funneling effect produced by the tapered ellipsoids, and Winston and Enoch⁵⁴ employed geometric optics to investigate the light-collecting properties of the retinal photoreceptor.

Recall the behavior of rays incident upon the interface separating two dielectric media (see Fig. 7a). Rays incident from medium 1 at an angle θ_1 to the interface normal, will be refracted at an angle θ_2 in medium 2; the angles are related by Snell's law where:

$$n_1 \sin \theta_1 = n_2 \sin \theta_2 \quad (1)$$

where n_1 and n_2 are the refractive indices of media 1 and 2, respectively. As shown in Fig. 6a, the refracted rays are bent toward the interface, since the incidence is from a medium of greater refractive index toward a medium of lesser refractive index. In this situation, the limiting angle of incidence angle θ_L equals the critical angle of the total internal reflection, for which the refracted ray travels along the interface (angle $\theta_2 = 90^\circ$); from Snell's law, the critical angle = $\arcsin(n_2/n_1)$. Rays incident at angles smaller than critical angle will undergo partial reflection and refraction; rays incident at angles equal to or larger than critical angle will undergo total internal reflection.

Let us now consider the pertinent geometry involved with a dielectric cylinder. A circular dielectric cylinder of refractive index n_1 is embedded in a lower-refractive-index surround n_2 . Light is incident upon the fiber end face at an angle θ to the z or cylinder axis from a third medium of refractive index $n_3 < n_1$ (Fig. 7b). We consider the meridional rays, that is, rays which intersect the fiber axis. There are

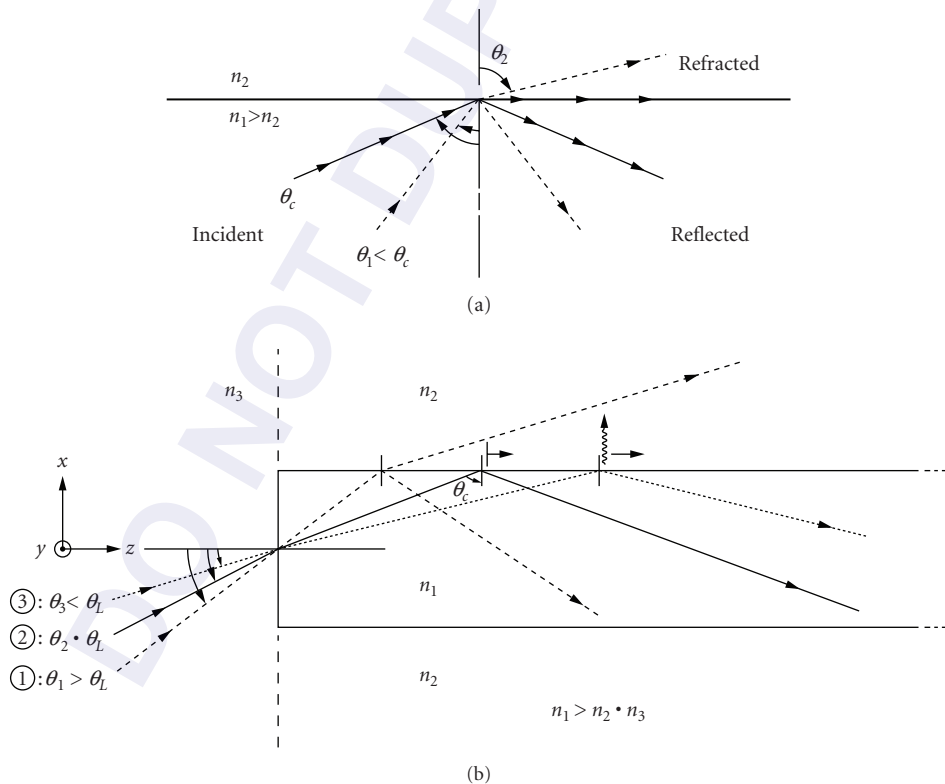


FIGURE 7 (a) Reflection and refraction of light rays incident upon an interface separating two different media. (b) Ray optics analysis of light incident on an idealized cylindrical photoreceptor (see text for details).

three possibilities for meridional rays, shown by rays 1 through 3. Once inside the cylinder, ray 1 is incident upon the side wall at an angle smaller than the critical angle, is only partially reflected, and produces, at each reflection, a refracted ray that carries power away from the fiber. Ray 2 is critically incident and is thus totally reflected; this ray forms the boundary separating meridional rays that are refracted out of the fiber from meridional rays that are trapped within the fiber. Ray 3 is incident upon the side wall at an angle greater than the critical angle, and is trapped within the fiber by total internal reflection.

Ray optics thus predicts a limiting incidence angle θ_L , related to the critical angle via Snell's law, given by

$$n_3 \sin \theta_L = [(n_1)^2 - (n_2)^2]^{1/2} \quad (2)$$

for which incident meridional rays may be trapped within the fiber. Only those rays incident upon the fiber end face at angles $\theta < \text{limiting angle } \theta_L$ are accepted; those incident at greater angles are quickly attenuated owing to refraction. The numerical aperture (NA) of the fiber is defined to be $n_3 \sin \theta_L$.











Application of meridional ray optics thus yields an acceptance property wherein only rays incident at $\theta < \theta_L$ are trapped within the fiber. One thing to note is that the greater the angle of incidence is, the faster the attenuation of light will be through the fiber. This result is consistent with results obtained with electromagnetic waveguide modal analysis, which show that the farther away from a cutoff a mode is, the more tightly bound to the cylinder it is, and the closer its central ray is to the cylinder axis.

Geometric optics predicts that an optical fiber can trap and guide light; however, it does not provide an accurate quantitative description of phenomena on fibers having diameters of the order of the wavelength of visible light, where diffraction effects are important. Also, taking a simplified electromagnetic analysis of the photoreceptor waveguide model, illumination is assumed to be monochromatic, and the first Born approximation⁵⁵—which ignores the contribution of waves that are scattered by the apertures—is used to match the incident illumination in the equivalent truncating aperture to the modes of the lossless guide that represents the myoid. Modes excited on the myoid are then guided to the tapered ellipsoid. Some of the initially excited radiation field will have leaked away before the ellipsoid is reached. Each segment is dealt with separately; thus, the ellipsoid will guide or funnel the power received from the myoid to the outer segment. Since each segment is dealt with separately, and because the modes of those segments are identical to those of an infinitely extended cylinder of the same radius and refractive index, the dielectric structure is considered to have an infinite uniform dielectric cylinder of radius p and refractive index n_1 embedded in a lower-refractive-index surround n_2 .

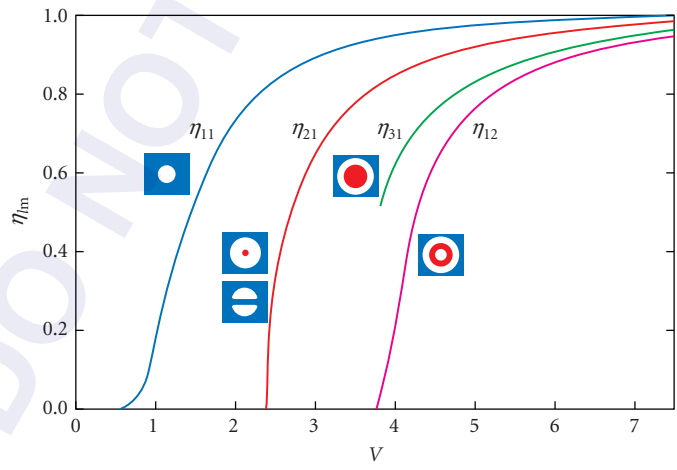
8.7 QUANTITATIVE OBSERVATIONS OF SINGLE RECEPTORS

The first observations of waveguide modal patterns in photoreceptors were made by Enoch.^{47,56–58} Given that vertebrate photoreceptors act as waveguides, in order to treat the optics of a single receptor theoretically, it is modeled as a cylindrical waveguide of circular symmetry, with two fixed indices of refraction inside and outside the guide. Note that light is propagated along the guide in one or more of a set of modes, each associated with a characteristic three-dimensional distribution of energy density in the guide. In the ideal guide their distributions may be very nonuniform but will show a high degree of symmetry. Modes are grouped into orders, related to the number of distinct modal surfaces on which the electromagnetic fields go to zero. Thus, in Fig. 8, the lowest (first) order mode (designated HE_{11}) goes to zero only at an infinite distance from the guide axis. Note that for all modes the fields extend to infinity in principle. At the terminus of the guide the fiber effectively becomes an antenna, radiating light energy into space. If one observes the end of an optical waveguide with a microscope, one sees a more or less complex pattern of light which roughly approximates a cross section through the modes being transmitted.

Modal patterns

Designations	Negative form	Cutoff
HE ₁₁		0.0
TE ₀₁ , TM ₀₁ HE ₂₁		2.405 2.4+
(TE ₀₁ or TM ₀₁) + HE ₂₁		2.4+
HE ₁₂		3.812
EH ₁₁ HE ₃₁		3.812 3.8+
HE ₁₂ + (HE ₃₁ or EH ₁₁)		3.8+
HE ₁₃ + EH ₁₁		3.8+
EH ₂₁ or HE ₄₁		5.2
TE ₀₂ , TM ₀₂ , HE ₂₂		5.52
(TE ₀₂ or TM ₀₂) + HE ₂₂		5.52

(a)



(b)

FIGURE 8 (a) Some commonly observed modal patterns found in human photoreceptors; their designations and V values. (b) The fraction of transmitted energy as a function of the V parameter for some commonly observed modal patterns.

The optical properties of the ideal guide are completely determined by a parameter called the *waveguide parameter* V , defined by

$$V = (\pi d/\lambda) [(n_1)^2 - (n_0)^2]^{1/2} \quad (3)$$

where d is the diameter of the guide, λ the wavelength in vacuum, and n_1 and n_0 are the indices of refraction of the inside (core) and outside (cladding) of the fiber, respectively. For a given guide, V is readily changed by changing wavelength as shown above. With one exception, each of the various modes is associated with a cutoff value, V_c , below which light cannot propagate indefinitely along the fiber in that mode. Again, the lowest order, HE_{11} , mode does not have a cutoff value. V may be regarded as setting the radial scale of the energy distribution; for $V \gg V_c$, the energy is concentrated near the axis of the guide with relatively little transmitted outside the core. As V approaches V_c , the energy spreads out with proportionally more being transmitted outside the guide wall. A quantity η (not to be confused with the relative luminous efficiency defined for SCE I) has been defined that represents the fraction of transmitted energy which propagates inside the guide wall. For a given mode or combination of modes and given V , η can be evaluated from waveguide theory. Figure 8a shows schematic representations of the mode patterns commonly observed in vertebrate receptors, cochlear hair cells, cilia (various), and the like, together with V_c for the corresponding modes. Figure 8b shows values of η as a function of V for some of the lower order modes.⁵⁹

It is evident that an understanding of receptor waveguide behavior requires a determination of V for a range of typical receptors. An obvious straightforward approach is suggested, where if one could determine V at any one wavelength, it could be found for other wavelengths by simple calculation. How is this done? If one observes a flat retinal preparation illuminated by a monochromator so that the mode patterns are sharp, then by rapidly changing the wavelength throughout the visible spectrum, some of the receptors will show an abrupt change of the pattern from one mode to another. This behavior suggests that for these receptors V is moving through the value of V_c for the more complex mode resulting in a change in the dominant mode being excited. To be valid, this approach requires that the fraction of energy within the guide (η) should fall abruptly to zero at cutoff. In experimental studies, it has been shown that η does go to zero at cutoff for the second-order modes characterized by bilobed or thick annular patterns as shown in Fig. 8. These modes include HE_{12} , TE_{01} , TM_{01} , $\text{TM}_{01}(\text{TE}_{01} + \text{HE}_{21})$, and $(\text{TM}_{01} + \text{HE}_{21})$. Figure 9 displays some waveguide modes obtained in humans and monkeys, as well as variation with wavelength. In turn, it has been discovered that these second-order modal patterns are by far the most frequently observed in small-diameter mammalian receptors. Also, a detailed understanding of receptor-waveguide behavior requires accurate values for guide diameters. Another fact to note concerning parameter V is that the indices of refraction of the entire dielectric cylinder are fairly constant with no pronounced change in waveguide behavior as determined from past experimental studies. For reinforcement, considerable power is carried outside the fiber for low (above cutoff) values of V . In the limit as $V \rightarrow$ infinity, the mode's total power is carried within the fiber.

The most complete study of the electrodynamic properties of visible-light interaction with the outer segment of the vertebrate rod based on detailed, first-principles computational electromagnetics modeling using a direct time integration of Maxwell's equation in a two-dimensional space grid for both transverse-magnetic and transverse-electric vector field modes was presented by Picket-May, Taflou, and Troy.⁶⁰ Detailed maps of the standing wave within the rod were generated (Fig. 10). The standing-wave data were Fourier analyzed to obtain spatial frequency spectra. Except for isolated peaks, the spatial frequency spectra were found to be essentially independent of the illumination wavelength. This finding seems to support the hypothesis that the electrodynamic properties of the rod contribute little if at all to the wavelength specificity of optical absorption.⁶¹ Frequency-independent structures have found major applications in broad band transmission and reception of radio-frequency and microwave signals. As Picket-May et al.⁶⁰ point out, it is conceivable that some engineering application of frequency-independent, retinal-rod-like structures may eventually result for optical signal processing.

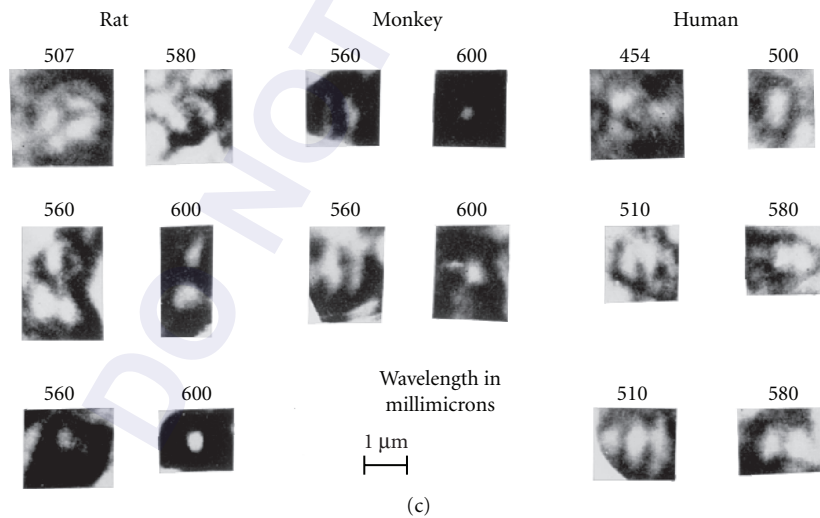
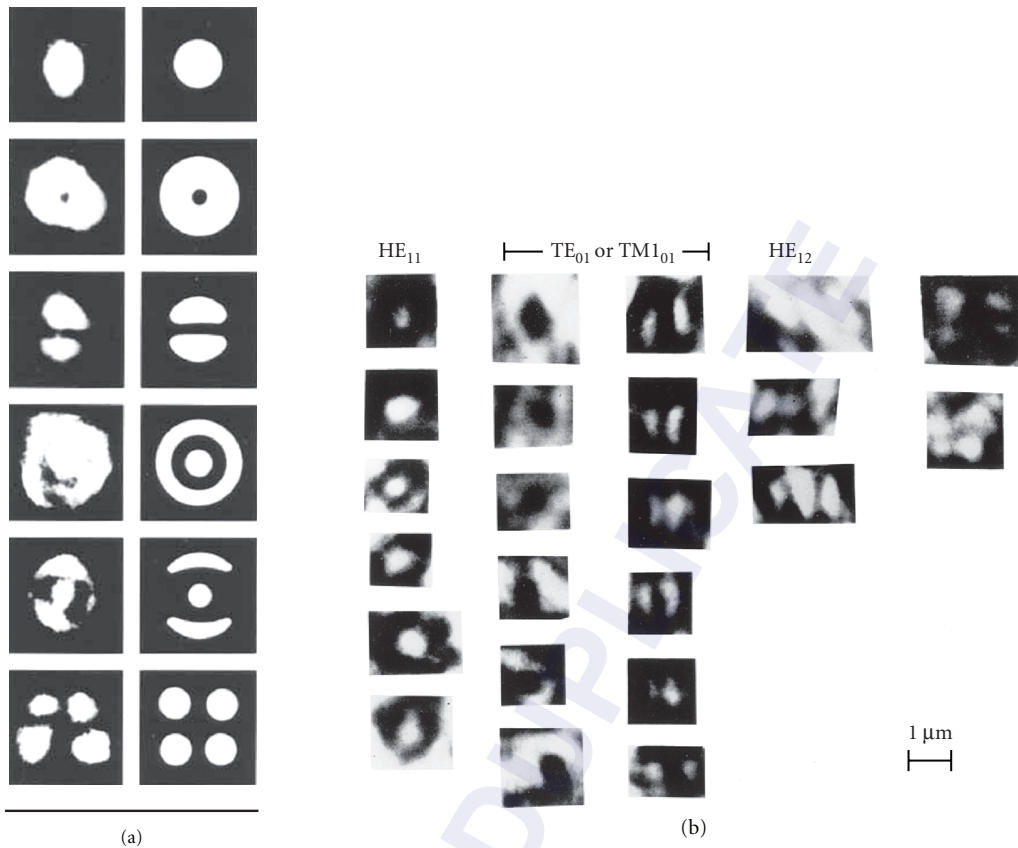


FIGURE 9 (a) Human waveguide modal patterns and their theoretical, idealized equivalents. The bar represents 10 μm . (b) Monkey waveguide modal patterns. The bar represents 1 μm . (c) Changes in modal patterns induced by wavelength in three species.



FIGURE 10 Computed figure showing the magnitude of the optical standing wave within a retinal rod at three different wavelengths (475, 505, and 714 nm). White areas are standing wave peaks, dark areas standing wave nulls. (From Ref. 26.)

8.8 WAVEGUIDE MODAL PATTERNS FOUND IN MONKEY/HUMAN RETINAL RECEPTORS

When a given modal pattern is excited, the wavelength composition of the energy associated with the propagation of that pattern in the retinal receptor (acting as a waveguide) may differ from the total distribution of energy at the point of incidence; for it is important to realize that the energy which is not transmitted is not necessarily absorbed by the receptor. It may be reflected and refracted out of the receptor. A distinct wavelength separation mechanism is seen to exist in those retinal

receptors observed. In some instances, one type of modal pattern was seen having a dominant hue. Where multiple colors were seen in a given receptor unit, those colors in some instances were due to changes in modal pattern with wavelength. The appearance of multiple colors, observed when viewing the outer segments, occurred in some receptors where mode coupling or interactive effects were present. There are differences in the wavelength distribution of the energy transmitted to the terminal end of the receptor. In addition, there are differences in the spatial or regional distribution of that energy within the outer segments as a function of wavelength due to the presence of different modal patterns and interactions governed by parameter V .⁶²

If the same well-oriented central foveal receptors of humans are illuminated at three angles of obliquity of incident light, other changes are observed. It is observed subjectively that increased obliquity results in more red being transmitted, or less yellow and green being seen. Why does this occur? This phenomenon is dependent on the angle of incidence of light upon the receptor and the spectral sensitivity of the outer segment. First, the angle of incidence of light upon the receptor determines the amount of light transmitted through the receptor, as shown. Therefore, different values of V result in different amounts of light being collected for a fixed angle of incidence, as well as in different shapes, hence, in different widths, of the respective selectivity curves. Also, as the angle of incidence approaches the limiting angle, there is a decrease in the amount of bound-mode power of light transmitted through the receptor. In turn, absorption is greater for the larger-diameter receptors which, having higher V values, are more efficient at collecting light. So for higher amounts of light transmitted through the receptor with less contrast of different hues, the receptor has a high V value (large myoid diameter) and the light incident upon it is closer to normal to the myoid's surface versus its limiting angle. Another fact to note is that the absorption peak occurs at decreasing wavelengths for decreasing receptor radius, and the absorption curves broaden for decreasing radius of outer segment.

Also, since wavelength transmissivity, as well as total transmissivity, changes with obliquity of incidence, it may be further assumed that receptors that are not oriented properly in the retina will not respond to stimuli in the same manner as receptors that are oriented properly. In other words, disturbance in the orientation of receptors should result in some degree of anomalous color vision as well as visual acuity.^{55,63–65}

It has been suggested that Henle fibers found in the retina may act as optical fibers directing light toward the center of the fovea.⁶⁶ Henle fibers arise from cone cell bodies, and turn horizontally to run radially outward from the center of the fovea, parallel to the retinal surface. At the foveal periphery, the Henle fibers turn again, perpendicular to the retinal surface, and expand into cone pedicles. Pedicles of foveal cones form an annulus around the fovea. Cone pedicles are located between 100 to 1000 μm from the center of the fovea and the Henle fibers connecting them to cone photoreceptor nuclear regions can be up to 300 μm long. If this hypothesis is correct, then Henle fiber-optic transmission could increase foveal irradiance, channeling short wavelength light from perifoveal pedicles to the central fovea and increasing the fovea's risk of photochemical damage from intense light exposure. This property has also been suggested as a reason for foveomacular retinitis burn without direct solar observation and could also be considered as a previously unsuspected risk in perifoveal laser photocoagulation therapy.

More recently, Franze et al.⁶⁷ investigated intact retinal tissue and individual Muller cells, which are radial glial cells spanning the entire thickness of the retina. Transmission and reflection confocal microscopy of retinal tissue *in vivo* and *in vitro* showed that these cells provide a low scattering passage for light from the retinal surface to the photoreceptors. In addition, using a modified dual beam laser trap, these researchers were able to demonstrate that these individual cells act as optical fibers. Their parallel array in the retina is analogous to fiberoptic plates used for low-distortion image transfer. It should be noted that these cells have an extended funnel shape, a higher refractive index than their surrounding tissue and are oriented along the direction of light propagation. Calculation of the V parameter showed that it varied from 2.6 to 2.9 (at 700 nm) for the different parts along the Muller cell, which is sufficiently high for low loss propagation of a few modes in the structure at this long wavelength. In the more intermediate wavelengths, the V varies from 3.6 to 4.0. Even though the refractive index and diameter of the cells change along their length, the V parameter and thus, their light guiding capability is very nearly constant. Unlike photoreceptors, muller cells do not have

the smooth cylindrical shape and have complex side branching processes. Their inclusion, using an “effective index” actually increases the V parameter. Therefore, even though Muller cells have a complex morphology, they can function as optical waveguides for visible light.

Modern theories of the Stiles-Crawford Effect I (SCE I) have used waveguide analysis to describe the effect. The most noteworthy is that of Snyder and Pask.⁶⁸ Snyder and Pask hypothesized an ideal “average” foveal cone with uniform dimensions and refractive indices. The cone is thought to be made up of two cylinders representing inner and outer segments. The response of a single cone is assumed to be a monotonically increasing function of absorbed light. The psychophysical response is assumed to be a linear combination of single-cone responses. Values of various parameters in the associated equations were adjusted to fit the wavelength variation of the SCE I reported by Stiles.⁶⁹ The model predicts the magnitude and trends of the experimental data (Fig. 11). The predictions, however, are highly dependent on the choice of refractive indices. Even small variations ($\sim 0.5\%$) produce enormous changes in the results. The model is based not upon the absorption of light, but upon transmitted modal power. Other waveguide models include those of, for example, Alpern,⁷⁰ Starr,⁷¹ and Wijngaard et al.⁷² These papers deal with SCE II, the color effect. Additionally, Goyal et al.⁷³ have proposed an inhomogeneous waveguide model for cones. Vohnsen and his colleagues^{74,75} have modified the existing models by making an important assumption, namely that instead of assuming that the light that propagates as rays in the eye to/from photoreceptors are propagated as waves. This implies that diffraction of both the incoming beam and the light emanated by the photoreceptors has to be included in the analysis. The mode field is approximated as a Gaussian. This is an important modification, since diffraction effects in relation to waveguide coupling have not been considered in previous models. Using this model it is possible to derive analytical expressions for the directionality parameters. The model also approximates the overall wavelength variation of directionality.

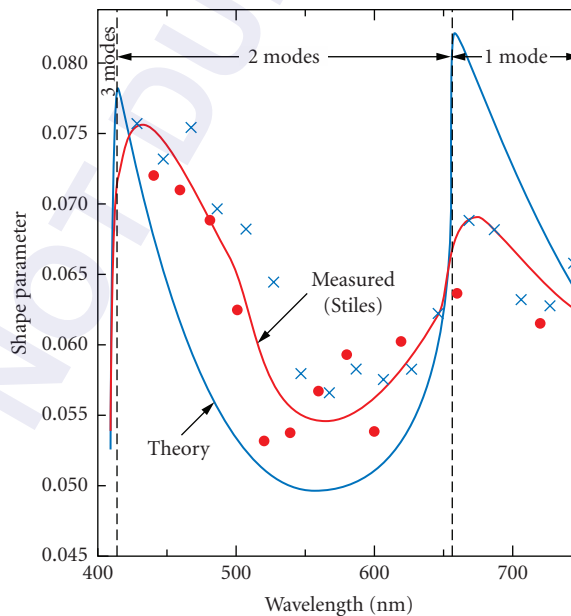


FIGURE 11 The idealized cone model of Snyder and Pask. (From Ref. 31) Here the magnitude of the shape parameter (the curvature of the SCE function) is plotted as a function of the wavelength. The blue curve is the theoretical prediction and the red curve is plotted through Stiles’s experimental points. Note the qualitative agreement between theory and experiment.

When phase variations (say, due to aberrations) of the field incident at the photoreceptor aperture are considered, the model predicts that the finite width of photoreceptors not only leads to a slight broadening of the point spread function, but also reduces the impact of aberrations on the visual sensation that is produced. This may play a role in accommodation, since a defocused image couples less light.

A caveat when considering waveguide models is that attention must be given to the physical assumptions made. These, along with uncertainties in the values of various physical parameters (e.g., refractive index), can seriously affect the results and hence the conclusions drawn from these models.

A formalism for modal analysis in absorbing photoreceptor waveguides has been presented by Lakshminarayanan and Calvo.⁷⁶⁻⁷⁸ Here, an exact expression for the fraction of total energy confined in a waveguide was derived and applied to the case of a cylindrically symmetric waveguide supporting both the zeroth- and first-order sets of excited modes. The fraction of energy flow confined within the cylinder can be decomposed into two factors, one containing a Z -dependence (Z is direction of cylinder axis) as a damping term and the second as a Z -dependent oscillatory term. It should be noted that because of non-negligible absorption, Bessel and Hankel functions with complex arguments have to be dealt with. Using this formalism, if it is found that, using physiological data, the attenuation of the fraction of confined power increased with distance, the energy is completely absorbed for a penetration distance $Z > 60 \mu\text{m}$.

More interesting results are obtained when two neighboring waveguides separated by a distance of the order of the wavelength (a situation found in the foveal photoreceptor mosaic) are considered. An incoherent superposition of the intensity functions in the two waveguides is assumed. The results are shown in Fig. 12. The radial distribution of transmitted energy for two values of core radius ($0.45 \mu\text{m}$

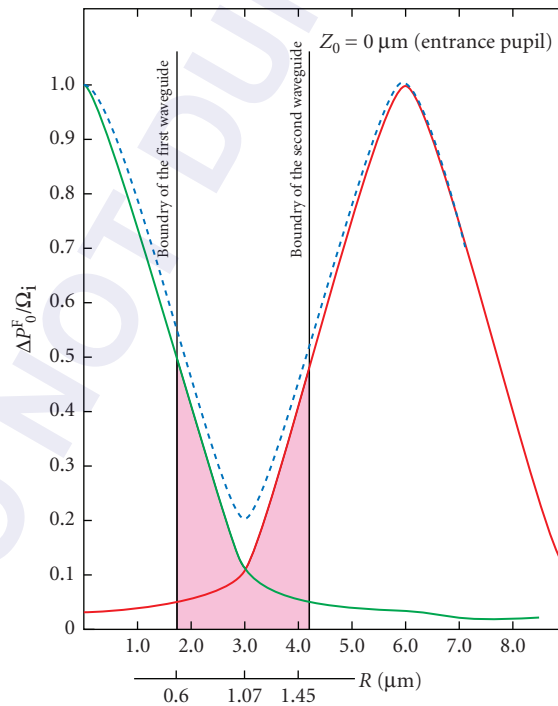


FIGURE 12 Effect of two neighboring absorbing waveguide photoreceptors on the fraction of confined power per unit area. The pink area denotes the region in which the energy of both waveguides is “cooperative.”

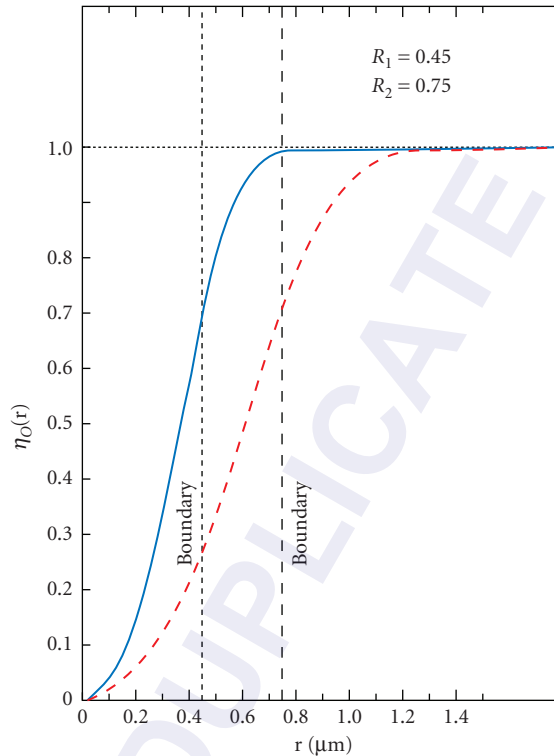


FIGURE 13 Flow of transmitted energy in a monomode absorbing waveguide for two value of the core radius (blue lines— $0.45 \mu\text{m}$; red line— $0.75 \mu\text{m}$) (see text for details).

and $0.75 \mu\text{m}$) is shown in Fig. 13. It is seen that 70 percent of the energy is still confined within the waveguide boundary, with some percentage flowing into the cladding, producing an “efficient” waveguide with radius of the order of $0.75 \mu\text{m}$ and $0.95 \mu\text{m}$ in the two cases. This effect together with the incoherent superposition implies a mechanism wherein as light propagates along the waveguide, there is a two-way transverse distribution of energy from both sides of the waveguide boundary, providing a continuous energy flow while avoiding complete attenuation of the transmitted signal. Analysis also showed that for specific values of the modal parameters, some energy can escape in a direction transverse to the waveguide axis at specific points of the boundary. This mechanism, a result of the Z -dependent oscillatory term, also helps to avoid complete attenuation of the transmitted signal.

The spatial impulse response of a single-variable cross-section photoreceptor has been analyzed by Lakshminarayanan and Calvo.⁷⁹ The modal field propagating under strict confinement conditions can be shown to be written in terms of a superposition integral which describes the behavior of a spatially invariant system. Using standard techniques of linear systems analysis, it is possible to obtain the modulation transfer functions of the inner and outer segments. This analysis neglects scattering and reflection at the inner-outer segment interface. The transfer function depends on both the modal parameter defining the waveguide regime and spatial frequency (Fig. 14). It is seen that both inner and outer segments behave as low-pass filters, although for the outer segment a wide frequency range is obtained, especially for the longer wavelengths. Using a different analysis, Stacey and Pask^{80,81} have shown similar results and conclude that photoreceptors themselves contribute to visual acuity

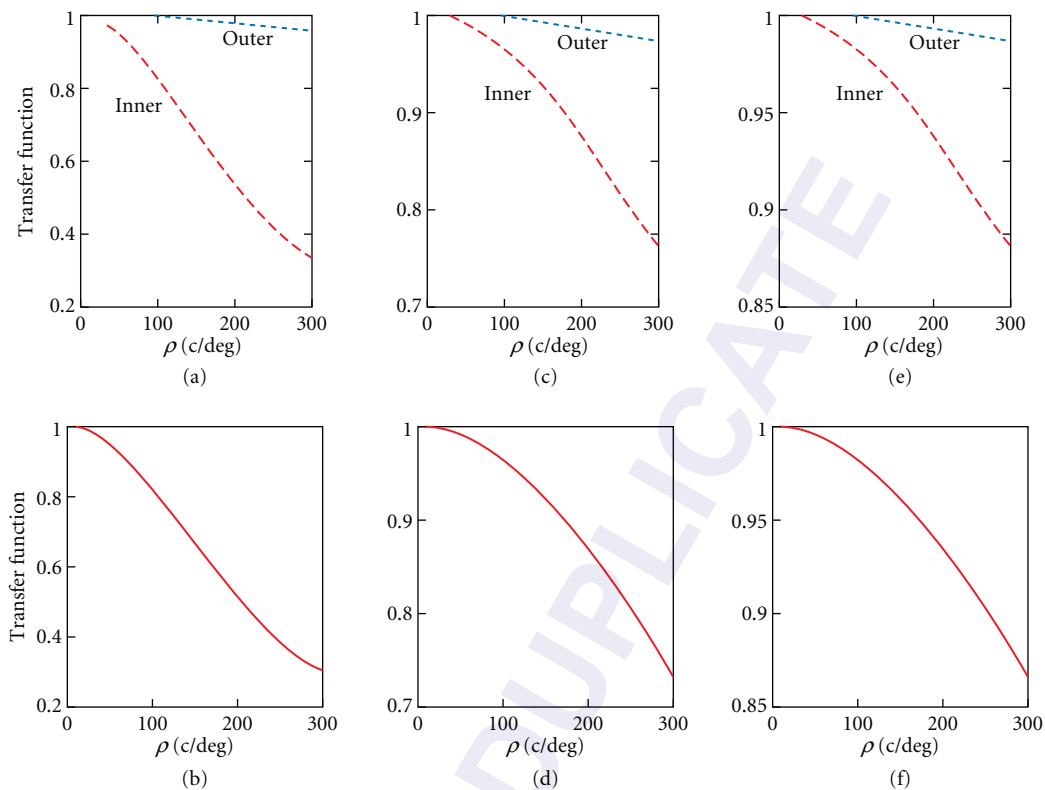


FIGURE 14 Transfer function of the inner and outer segments (top) and the total transfer function of the photoreceptor for three wavelengths: 680, 550, and 460 nm. The values of the modal parameters are fixed. Horizontal axis is in units of cycles/degree.

through a wavelength-dependent response at each spatial frequency and that the response is dependent on the coherence of the source. They have used this model to study the Campbell effect. (When measuring human visual acuity with an incoherent source, Campbell,⁸² noticed that the response decreased markedly when an artificial pupil placed in front of the eye was displaced perpendicular to the test fringes. There was no observed reduction when the pupil was displaced parallel to the fringes.) They conclude⁸³ that the waveguiding properties of photoreceptors make them sensitive to obliquely incident exciting waves and this provides some support for the hypothesis that both the SCE and the Campbell effect are manifestations of the same underlying waveguide nature of the photoreceptors.^{83,84}

8.9 LIGHT GUIDE EFFECT IN COCHLEAR HAIR CELLS AND HUMAN HAIR

The vibration of the organ of Corti in response to sound is becoming evident at the cellular level using the optical waveguide property of cochlear hair cells to study the motion of single hair cells. It has been reported that conspicuous, well-defined bright spots can be seen on the transilluminated organ of Corti. The spots are arranged in the mosaic fashion of hair cells known from surface views of the organ of Corti.⁸⁵ Thus, the hair cells appear to act as light guides. Also, a definite property of

biological waveguide requires the extracellular matrix or fluid to be lower in refractive index as compared to the fiber (in this case, cochlear hair cell), and this property has been elucidated by phase-contrast micrographs of cochlear hair cells.⁸⁶

Cochlear hair cells are transparent, nearly cylindrical bodies with a refractive index higher than the surrounding medium. Inner hair cells are optically denser than the surrounding supporting cells, and outer hair cells are optically denser than the surrounding cochlear fluid. Therefore, in optical studies, hair cells may be considered as optical fibers and the organ of Corti as a fiber-optics array. Even though the significance of the function of the organ of Corti in audition is understood, the role of waveguiding in hair cells is not clearly established.

To check for light-guide action, a broad light is usually used so that many hair cells show up simultaneously as light guides. For broad beams the conditions of light entrance into individual cells are not specified. Light can enter hair cells along the cylindrical wall. Part of such light is scattered at optical discontinuities inside the cells. To a large extent scattered light is trapped in cells by total internal reflection. Thus, to maximize contrast for transillumination with broad beams, no direct light is allowed to enter the microscope. Areas with hair cells appear dark. However, the heads of hair cells appear brightly illuminated owing to the light-collecting effect of hair cells on scattered light. Also, the diameter of cochlear hair cells is from 5 to 9 μm , that is, 10 to 20 times the wavelength of light. In fiber optics, fibers with diameters larger than 10 to 20 times the wavelength are considered multimodal because they support a large number of characteristic modes of light-energy propagation. Fibers with diameters smaller than 10 to 20 times the wavelength support fewer modes, and discrete modes of light-energy propagation can be launched more readily. Because of this division, cochlear hair cells fall in between typically single-mode and multimode optical fibers.⁸⁷ Also, an important parameter for the evaluation of fibers with small diameters is the parameter V . Recall that the smaller parameter V , the smaller the number of discrete modes a given fiber is capable of supporting. In addition, parameter V is greater for cochlear hair cells than for retinal receptor cells mainly because they have larger diameters than retinal rods and cones (see above for parameter V).

Gross features of surface preparations of the organ of Corti can be recognized with the binocular dissecting microscope. When viewed, the rows of hair cells can be seen as strings of bright pearls. Bright spots in the region of hair cells are due to light transmission through hair cells. Bright spots disappear when hair cells are removed.

Helmholtz's reciprocity theorem of optics⁵⁵ (reversibility of light path) also applies to hair cells where exit radiation pattern is equivalent to entrance radiation pattern. There is no basic difference in light-guide action when the directions of illumination and observation are reversed compared to the normal directions. Obviously, Helmholtz's theorem applies to visual photoreceptors too, since they are just modified cilia.

It is also seen that the detection of light guidance is shown to strongly correlate with the viewing angle, where incident light hits the numerical aperture of the hair cell, which is 30 to 40°. When the incident light is normal to the numerical aperture of the hair cell, higher transmittance of light and contrast is observed; however, at more oblique angles to the numerical aperture of the hair cell, less contrast is observed because the angle of incident light is approaching the limiting angle, which is a limit to total internal reflection and more refraction or scattering of rays occurs through the fiber.

A useful tool for studying the quality and preservation of preparations of organ of Corti utilizes the concept that hair cells show little change in appearance under the light microscope for up to 20 to 40 min after the death of an animal. Consequently, poor fixation and preservation of the organ of Corti can be recognized immediately at low levels of magnification.

To study the modal light intensity distributions on single hair cells, researchers use a monocular microscope. Hair cells exercise relatively little mode selection, which is to be expected because of their relatively large diameters (large parameter V). As the monocular microscope is focused up and down individual hair cells, it is seen that the intensity distributions in different cross sections of hair cells may differ. Focusing above the exit of hair cells shows that the near-field radiation pattern may differ from the intensity distribution on the cell exit. Furthermore, other studies have shown that the modal distribution of light intensity can be changed drastically for individual cells by changing the angle of light incidence relative to the axis of hair cells. Also, several low-order patterns on the exit end of single hair cells are easily identifiable and are designated according to the customary notation.

As a note, since cochlear hair cells show relatively little mode selection, it is not always possible to launch a desired mode pattern on a given cell with ordinary light. In order to restrict the spectrum of possible modes, it is advantageous to work with monochromatic, coherent light, that is, laser light. However, the identification of patterns relative to the cells was found to be difficult. This is due to the fact that laser stray light is highly structured (speckled), which makes it difficult to focus the microscope and to recognize cell boundaries. Light guidance in cochlear hair cells may be considered to be a general research tool for cochlear investigations in situations where the organ of Corti can be transilluminated and where relatively low levels of magnification are desirable.

Furthermore, biological waveguide effects have also been shown to occur in human hair,⁸⁸ which are similar to cochlear hair cell waveguide effects. There is a reduction in light intensity with increased hair length, but by far the most important factor was hair color. Little or no light was transmitted by brown hair, while gray hair acts as a natural fiber optic which can transmit light to its matrix, the follicular epithelium, and to the dermis. Interestingly, whether light transmission down hairs affects skin and hair needs investigation. As pointed out by Enoch and Lakshminarayanan,⁸⁹ many major sensory systems incorporate cilia in their structure. In the eye, the invagination of the neural groove of the surface ectoderm of the embryo leads to the development and elaboration of the nervous system. The infolding and incorporation of surface tissue within the embryo accounts for the presence of ciliated cells within nervous tissue neurons, which after maturation form photoreceptors containing cilia, lining the outer wall of the collapsed brain ventricle that will form the retina (see Oyster⁹⁰ for a description of the embryology and growth of the retina). A study of the evolution of the relationships between cilia and specific sensory systems would be of great interest.

8.10 FIBER-OPTIC PLANT TISSUES

Dark-grown (etiolated) plant tissues are analogous to multiple fiber-optic bundles capable of coherent transfer of light over at least 20 mm; hence, they can transmit a simple pattern faithfully along their length. Each tissue examined can accept incident light at or above certain angles relative to normal as is expected of any optical waveguide. The peak of the curve that describes this angular dependence, the acceptance angle, and the overall shape of the curve appear to be characteristic of a given plant species and not of its individual organs.⁹¹ Knowledge of this optical phenomenon has permitted localization of the site of photoreception for certain photomorphogenic responses in etiolated oats. Before the discussion of the fiber-optic properties of etiolated plant tissues, it is important to have a general idea of what plant structures are presently being discussed to get a full understanding of the structure to function relationship (Fig. 15). To define, the *epicotyl* is a region of shoot above the point of attachment of the cotyledons, often bearing the first foliage leaves. The *hypocotyl* is the part of the embryo proper below the point of attachment of the cotyledons; it will form the first part of the stem of the young plant. The *coleoptile* is a cylindrical sheath that encloses the first leaves of seedlings of grasses and their relatives.

Total internal reflection of light should be independent of fluence rate since this phenomenon is simply a special case of light scattering. The amount of light emerging from the tissue segment was measured first at the highest fluence rate available and then again when the incident light beam had been attenuated with one or more calibrated neutral density filters over a range of 4.5 absorbance units. The amount of light axially transmitted at a given fluence rate (expressed as a percentage of the amount transmitted at the highest fluence rate for individual tissue segments) decreased log-linearly as filters of increasing absorbance were used for oat, mung bean, and corn tissues. Therefore, light guiding in etiolated tissues is clearly fluence-rate-independent over at least 4.5 orders of magnitude.⁹²

Light guiding through segments of etiolated root, stem, and coleoptiles, plus primary leaves from several species, all display the spectral dependence expected of any light scattering agent, with low transmission in the blue relative to that in the far-red regions of the spectrum. As tissue length increases, the difference between transmissions in the blue and far-red regions of the spectrum increases.

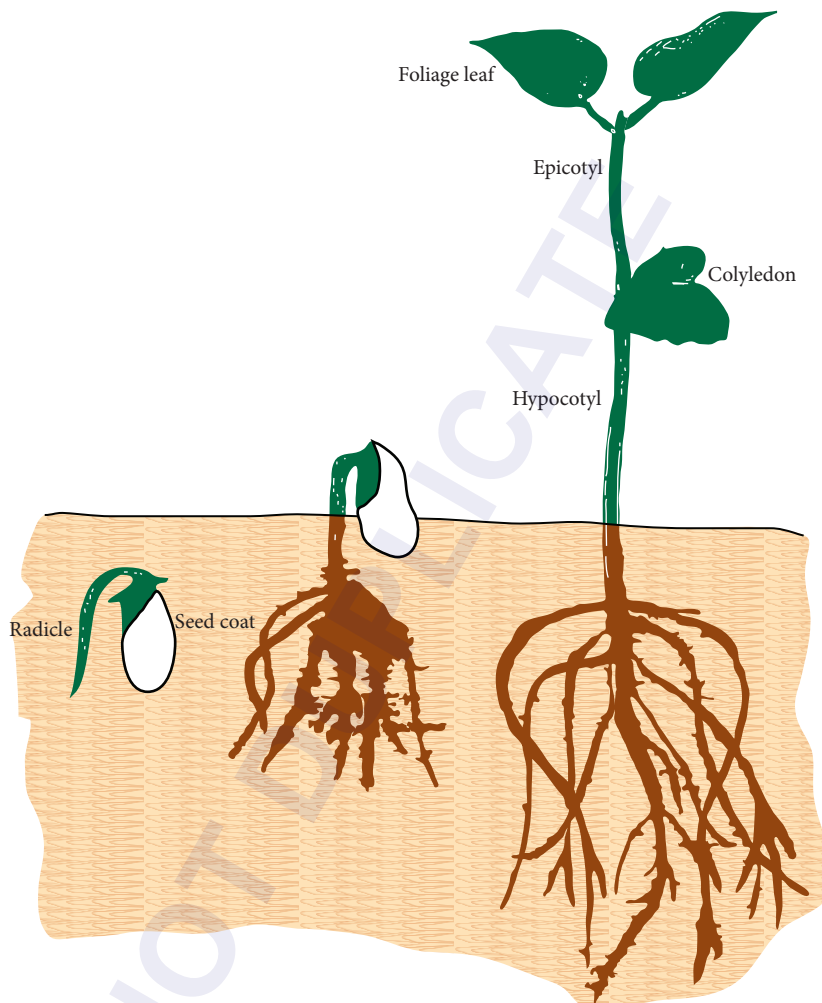


FIGURE 15 General structure of plant tissues.

It has been found that the spectra obtained when the incident beam of white light is applied to the side of the intact plant are indistinguishable from those obtained when the cut end of a tissue is irradiated. Spectra taken on tissue segments excised from different regions of the plant may show traces of certain pigments present only in those regions. Whereas there are no differences between spectra of the apical and basal portions of oat mesocotyls, the hook regions of mung bean hypocotyls contain small amounts of carotenoids, which is indicated as a small dip in the transmittance curve that is absent in spectra from the lower hypocotyl or hypocotyl-root transition zone. Oat coleoptilar nodes, which join the mesocotyl to the coleoptile, decrease light transmission from one organ to another when light is passed from the coleoptile to the mesocotyl or from the mesocotyl to the coleoptile, respectively, but show a spectral dependence which is independent of the direction of light transmission based on Helmholtz's reciprocity theorem of light. Older etiolated leaves do not transmit red light axially relative to the cylindrical coleoptile which sheaths them. Consequently, the presence or absence of the primary leaves does not alter the spectrum of the etiolated coleoptile.

Plant tissues mainly consist of vertical columns of cells. Light internally reflected through them passes predominantly along these columns of cells rather than across columns when the incident beam is at or above the tissue acceptance angle.⁹¹ As optical fibers, the tissues are about 10 percent as effective as glass rods of comparable diameter, and about 1 percent as effective as commercial fiber optics.⁹² Within a given cell, light could be traveling predominantly in the cell wall, the cytoplasm, the vacuole, or any combination of these three compartments. Observations of the cut ends of tissue segments that are light-guiding indicate that the cell wall does not transmit light axially but that the cell interior does. Also, plant cells contain pigments that are naturally restricted to the vacuole or to the cytoplasmic organelles. Mung bean hypocotyls contain carotenoids (absorption maxima near 450 nm), which are located in the etioplasts, organelles found only in the cytoplasm, whereas etiolated beet and rhubarb stems are composed of cells with clear cytoplasm and vacuoles containing betacyanin and/or anthocyanin (absorption maxima near 525 to 530 nm). Spectra of the light guided through these tissues show that light guiding occurs both through the cytoplasm and the vacuole. The relative size of each compartment, the pigment concentration, and the refractive indices involved will largely determine the amount of light that a given compartment will transmit. As the seedlings are exposed to light and become green, troughs due to absorption by chlorophylls (near 670 nm) appear. Completely greened tissue has very little light-guiding capacity except in the far-red region of the spectrum.

Also, the red to far-red (R/FR) ratio determines the fraction of phytochrome in the active state in the plant cell.⁹³ This ratio plays an important role in the regulation of growth. Similarly, the blue to red (B/R) ratio may be important for responses mediated by phytochrome alone when transmission is low in the red relative to the blue, or for those mediated by both phytochrome and a pigment absorbing blue light. The R/FR ratio of the light emitted from several curved, etiolated plant tissues decreases as the tissue length increases. However, the patterns of change in the B/R ratio differ with each tissue examined. All B/R ratios presumably decrease along the first 20 mm of tissues but in mung hypocotyl (stem) this ratio then increases, whereas in oat mesocotyl it increases and then decreases further. Changes in the R/FR and B/R ratios also occur in a given length of tissue as it becomes green. The R/FR ratio does not change in tissues that remain white (i.e., oat mesocotyl) but decreases dramatically in those that synthesize pigments (i.e., oat coleoptile or mung bean, beet, and rhubarb stems). The B/R ratios, unlike the R/FR ratios, tend to increase as the tissues become green, reflecting a smaller decrease in transmission at 450 nm than 660 nm as pigments accumulate.

All of these results demonstrate certain properties of light-guiding effects in plants: the fluence-rate independence of light guiding in plants; an increase in tissue length showing a more pronounced difference between transmission in the blue and far-red regions of the spectra; and the effect of the greening status on the pertinent R/FR and B/R ratios for several species of plants grown under a variety of conditions. Certainly, both the B/R and the R/FR ratios change dramatically in light that is internally reflected through plants as etiolated tissue length is increased and as greening occurs. Finally, the sensitivity of etiolated plant tissues to small fluences of light and the inherent capacity of plants to detect and respond to small variations in the R/FR ratio indicate that light guiding occurring in etiolated plants beneath the soil and the concomitant spectral changes that occur along the length of the plant may be of importance in determining and/or fine-tuning the plant's photomorphogenic response.^{94,95}

8.11 SPONGES

Recently, researchers at Bell Laboratories studied the fiber optical features of a deep sea organism, namely, the glass sponge (*Euplectella*, "venus flower basket").⁹⁶ The spicules of these sponges have some remarkable fiber optic properties. The skeleton of the hexactinellid class of sponges is constructed from amorphous hydrated silica. The sponge has a lattice of fused spicules that provide extended structural support. The spicules are 5 to 15 cm long and 40 to 70 μm in diameter. The spicules have a characteristic layered morphology and a cross-sectional variation in composition: a pure silica core of about 2 μm in diameter that enclosed an organic filament, a central cylinder with maximum organic material and a striated shell that has a gradually decreasing organic content.

The researchers conducted interferometric index profiling and found that corresponding to the three regions, there was a core with a high refractive index that is comparable to or higher than that of vitreous silica, a cylinder of lower refractive index that surrounds the core and an oscillating pattern with progressively increasing refractive index at the outer part of the spicule. These embedded spicules act as single or a few-mode waveguides. When light was coupled to free standing spicules, they functioned as multimode fibers.

These fibers are similar to commercial telecommunication fibers in that they are made of the same material and have comparable dimensions. They also function as efficient as single-mode, few-mode, or multi-mode, fibers depending upon optical launch conditions.

In addition, it is possible that these spicules in addition to providing structural anchorage support, could also act as a fiber optic network for distributing light in the deep sea environment. The fascinating aspect of these fibers is that they are made under ambient conditions, as opposed to commercial manufactured fibers which require very high temperatures (about 1700°C). Organic ligands at the exterior of the fiber seem to protect it and provide an effective crack-arresting mechanism and the fibers seem to be doped with specialized impurities that improve the refractive index profile and hence the waveguiding properties.

8.12 SUMMARY

In short, we have pointed out some similar relationships between the theoretical waveguide models and the actual biological waveguides such as vertebrate photoreceptors, cochlear hair cells of the inner ear, human hair (cilia), and fiber-optic plant tissues. The similar relationships between the theoretical model of waveguides and actual biological waveguides were demonstrated as (a) almost entirely following the Stiles-Crawford effect of the first kind (SCE-I) for vertebrate photoreceptors, (b) exhibiting wavelength sensitivity of transmission, and (c) conforming to Helmholtz's reciprocity theorem of optics. Of special interest is the fact that visual photoreceptors exhibit self-orientation. Similar phototropic self-orientation is exhibited by growing and fully-grown plants (e.g., sunflowers). For theoretical or possible applications, the light-guiding properties of biological waveguides can be used not only to determine the structure of such guides, but also to provide the basis for determining the functioning of other structures related to these guides. Lastly, these studies which give us greater insight into low temperature, biologically inspired processes could possibly result in better fiber optic materials and networks for commercial applications such as telecommunications.

8.13 REFERENCES

1. W. S. Stiles and B. H. Crawford, "The Luminous Efficiency of Rays Entering the Eye Pupil at Different Points," *Proc. R. Soc. Lond. B.* **112**:428–450 (1933).
2. A. W. Snyder and R. Menzel, *Photoreceptor Optics*, Springer-Verlag, New York, 1975.
3. J. M. Enoch and F. L. Tobey Jr., *Vertebrate Photoreceptor Optics*, vol. 23, *Springer Series in Optical Science*. Springer-Verlag, Berlin, 1981.
4. J. M. Enoch and V. Lakshminarayanan, "Retinal Fiber Optics," in W. N. Charman (ed.), *Vision and Visual Dysfunction: Visual Optics and Instrumentation*, vol. 1, MacMillan Press, London, 1991, pp. 280–309.
5. V. Lakshminarayanan, "Waveguiding in Retinal Photoreceptors: An Overview," *Proc. SPIE* **3211**:182–192 (1998).
6. V. Lakshminarayanan, J. M. Enoch, and A. Raghuram, *The Stiles-Crawford Effects*, Classic Reprints on CD-ROM, vol. 4, Optical Society of America, Washington, D.C. (2003).
7. J. Dowling, *The Retina: An Approachable Part of the Brain*, Harvard University Press, Cambridge, MA, 1987.
8. R. W. Rodieck, *The First Steps in Seeing*, Sinauer, Sunderland, MA, 1998.

9. W. S. Stiles, "The Luminous Efficiency of Monochromatic Rays Entering the Eye Pupil at Different Points and a New Color Effect," *Proc. R. Soc. Lond. B* **123**:90–118 (1937).
10. J. M. Enoch and W. S. Stiles, "The Color Change of Monochromatic Light with Retinal Angle of Incidence," *Optica Acta* **8**:329–358 (1961).
11. M. F. Land, "The Optics of Animal Eyes," *Contemp. Physiol.* **29**:435–455 (1988).
12. M. F. Land, "Optics and Vision in Invertebrates," in H. Autrum (ed.), *Handbook of Sensory Physiology*, vol. VII/6B, Springer-Verlag, Berlin, 1981, pp. 471–592.
13. R. Sidman, "The Structure and Concentration of Solids in Photoreceptor Cells Studied by Refractometry and Interference Microscopy," *J. Biophys. Biochem. Cytol.* **3**:15–30 (1957).
14. J. M. Enoch and F. L. Tobey Jr., "Use of the Waveguide Parameter V to Determine the Differences in the Index of Refraction between the Rat Rod Outer Segment and Interstitial Matrix," *J. Opt. Soc. Am.* **68**(8):1130–1134 (1978).
15. R. A. Applegate and V. Lakshminarayanan, "Parametric Representation of Stiles-Crawford Function: Normal Variation of Peak Location and Directionality," *J. Opt. Soc. Am. A* **10**(7):1611–1623 (1993).
16. D. I. A. MacLeod, "Directionally Selective Light Adaptation: A Visual Consequence of Receptor Disarray?" *Vis. Res.* **14**:369–378 (1974).
17. S. A. Burns, S. Wu, F. C. DeLori, and A. E. Elsner, "Variations in Photoreceptor Directionality across the Central Retina," *J. Opt. Soc. Am. A* **14**:2033–2040 (1997).
18. A. Roorda and D. Williams, "Optical Fiber Properties of Human Cones," *J. Vision* **2**:404–412 (2002).
19. J. Krauskopf, "Some Experiments with a Photoelectric Ophthalmoscope," *Excerpta Medica International Congress Series* **125**:171–181 (1965).
20. J. M. Gorrand and F. C. DeLori, "A Reflectometric Technique for Assessing Photoreceptor Alignment," *Vis. Res.* **35**:999–1010 (1995).
21. P. J. deLint, T. T. J. M. Berendschott, and D. van Norren, "Local Photoreceptor Alignment Measured with a Scanning Laser Ophthalmoscope," *Vis. Res.* **37**:243–248 (1997).
22. G. J. van Blokland and D. van Norren, "Intensity and Polarization of Light Scattered at Small Angles from the Human Fovea," *Vis. Res.* **26**:485–494 (1986).
23. G. J. van Blokland, "Directionality and Alignment of the Foveal Receptors Assessed with Light Scattering from the Fundus in Vivo," *Vis. Res.* **26**:495–500 (1986).
24. S. A. Burns, S. Wu, J. C. He, and A. Elsner, "Variations in Photoreceptor Directionality Across the Central Retina," *J. Opt. Soc. Am. A* **14**:2033–2040 (1997).
25. S. Marcos and S. A. Burns, "Cone Spacing and Waveguide Properties from Cone Directionality Measurements," *J. Opt. Soc. Am. A* **16**:2437–2447 (1999).
26. J. M. Enoch, J. S. Werner, G. Haegerstrom-Portnoy, V. Lakshminarayanan, and M. Rynders, "Forever Young: Visual Functions Not Affected or Minimally Affected by Aging: A Review," *J. Gerontol. Biol. Sci.* **54**(8):B336–351 (1999).
27. E. Campos, J. M. Enoch, and C. R. Fitzgerald, "Retinal Receptive Field Like Properties and the Stiles-Crawford Effect in a Patient with a Traumatic Choroidal Rupture," *Doc. Ophthalmol.* **45**:381–395 (1978).
28. V. Lakshminarayanan, J. M. Enoch, and S. Yamade, "Human Photoreceptor Orientation: Normals and Exceptions," in *Advances in Diagnostic Visual Optics*, A. Fiorentini, D. L. Guyton, and I. M. Siegel, (eds.), Springer Verlag, Heidelberg, Germany, 1987, pp. 28–32.
29. V. Lakshminarayanan, *The Stiles Crawford effect in Aniridia*, Ph.D. Dissertation, University of California, Berkeley (1985).
30. M. S. Eckmiller, "Defective Cone Photoreceptor Cytoskeleton, Alignment, Feedback and Energetics can Lead to Energy Depletion in Macular Degeneration," *Prog. Retin Eye Res.* **23**:495–522 (2004).
31. V. Lakshminarayanan, J. E. Bailey, and J. M. Enoch, "Photoreceptor Orientation and Alignment in Nasal Fundus Ectasia," *Optom. Vis. Sci.* **74**(12):1011–1018 (1997).
32. S. S. Choi, J. M. Enoch, and M. Kono, "Evidence for Transient Forces/Strains at the Optic Nervehead in Myopia: Repeated Measurements of the Stiles Crawford Effect of the First Kind (SCE-1) over Time," *Ophthalm. Physiol. Optics* **24**:194–206 (2004).
33. X. Zhang, Ye. Ming, A. Bradley, and L. Thibos, "Apodization of the Stiles Crawford Effect Moderates the Visual Impact of Retinal Image Defocus," *J. Opt. Soc. Am. A* **16**(4):812–820 (1999).

34. L. Martin, *Technical Optics*, vol. 2, Pitman, London, UK, 1953.
35. D. Atchison, D. Scott, and G. Smith, "Pupil Photometric Efficiency and Effective Center," *Ophthalm. Physiol. Optics* **20**:501–503 (2001).
36. W. S. Baron and J. M. Enoch, "Calculating Photopic Illuminance," *Amer. J. Optometry Physiol. Optics* **59**(4): 338–341 (1982).
37. H. Metcalf, "Stiles–Crawford apodization," *J. Opt. Soc. Am.* **55**:72–74 (1965).
38. J. P. Carrol, "Apodization Model of the Stiles–Crawford Effect," *J. Opt. Soc. Am.* **70**:1155–1156 (1980).
39. M. Rynders, L. Thibos, A. Bradley, and N. Lopez-Gil, "Apodization Neutralization: A New Technique for Investigating the Impact of the Stiles–Crawford Effect on Visual Function," in: *Basic and Clinical Applications of Vision Science*, V. Lakshminarayanan, (ed.), Kluwer, Dordrecht, The Netherlands, 1997, pp. 57–61.
40. D. H. Scott, D. A. Atchison, and P. A. Pejeski, "Description of a Method for Neutralizing the Stiles–Crawford Effect," *Ophthalm. Physiol. Optics* **21**(2):161–172 (2001).
41. D. A. Atchison, A. Joblin, and G. Smith, "Influence of Stiles–Crawford Apodization on Spatial Visual Performance," *J. Opt. Soc. Am. A* **15**(9):2545–2551 (1998).
42. D. A. Atchison, D. H. Scott, A. Joblin, and F. Smith, "Influence of Stiles–Crawford Apodization on Spatial Visual Performance with Decentered Pupils," *J. Opt. Soc. Am. A* **18**(6):1201–1211 (2001).
43. S. Marcos, S. Burns, E. Moreno-Barriuso, and R. Navarro, "A New Approach to the Study of Ocular Chromatic Aberrations," *Vis. Res.* **39**:4309–4323 (2000).
44. J. van de Kraats, T.J.M. Berendschot, and D. van Norren, "The Pathways of Light Measured in Fundus Reflectometry," *Vis. Res.* **36**:2229–2247 (1996).
45. T. T. J. Berenschot, J. van de Kraats, and D. van Norren, "Wavelength Dependence of the Stiles–Crawford Effect Explained by Perception of Backscattered Light from the Choroid," *J. Opt. Soc. Am. A* **18**(7):1445–1451 (2001).
46. S. Marcos and S. A. Burns, "On the Symmetry between Eyes of Wavefront Aberration and Cone Directionality," *Vis. Res.* **40**:2437–4580 (2000).
47. J. M. Enoch, "Optical Properties of Retinal Receptors," *J. Opt. Soc. Am.* **53**:71–85 (1963).
48. W. H. Miller and A. W. Snyder, "Optical Function of Myoids," *Vis. Res.* **12**(11):1841–1848 (1972).
49. B. Burnside and C. King-Smith, "Retinomotor Movements," in L. Squire, *New Encyclopedia of Neuroscience*, Elsevier, in press 2008.
50. A. W. Snyder and J. D. Love, *Optical Waveguide Theory*, Chapman and Hall, London, 1983.
51. H. von Helmholtz, *Treatise in Physiological Optics*, Dover, New York, 1962, p. 229.
52. A. Hannover, *Vid. Sel. Naturv. Og Math. Sk.* **X** (1843).
53. B. O'Brien, "A Theory of the Stiles Crawford Effect," *J. Opt. Soc. Am.* **36**(9):506–509 (1946).
54. R. Winston and J. M. Enoch, "Retinal Cone Receptors as an Ideal Light Collector," *J. Opt. Soc. Am.* **61**(8):1120–1122 (1971).
55. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon Press, Oxford, 1984.
56. J. M. Enoch, "Waveguide Modes in Retinal Receptors," *Science* **133**(3461):1353–1354 (1961).
57. J. M. Enoch, "Nature of the Transmission of Energy in the Retinal Photoreceptor," *J. Opt. Soc. Am.* **51**:1122–1126 (1961).
58. J. M. Enoch, "Optical Properties of the Retinal Receptor," *J. Opt. Soc. Am.* **53**:71–85 (1963).
59. K. Kirschfeld and A. W. Snyder, "Measurement of Photoreceptors Characteristic Waveguide Parameter," *Vis. Res.* **16**(7):775–778 (1976).
60. M. J. May Picket, A. Taflove, and J. B. Troy, "Electrodynamics of Visible Light Interaction with the Vertebrate Retinal Rod," *Opt. Letters* **18**:568–570 (1993).
61. A. Fein and E. Z. Szuts, *Photoreceptors: Their Role in Vision*, Cambridge University Press, Cambridge, UK, 1982.
62. J. M. Enoch, "Receptor Amblyopia," *Am. J. Ophthalmol.* **48**:262–273 (1959).
63. V. C. Smith, J. Pokorny, and K. R. Diddie, "Color Matching and the Stiles–Crawford Effect in Observers with Early Age-Related Macular Changes," *J. Opt. Soc. Am. A.* **5**:2113–2121 (1988).
64. V. C. Smith, J. Pokorny, J. T. Ernest, and S. J. Starr, "Visual Function in Acute Posterior Multifocal Placoid Pigment Epitheliopathy," *Am. J. Ophthalmol.* **85**:192–199 (1978).

65. F. W. Campbell and A. H. Gregory, "The Spatial Resolving Power of the Human Retina with Oblique Incidence," *J. Opt. Soc. Am.* **50**:831 (1960).
66. M. A. Mainster, "Henle Fibers May Direct Light toward the Center of the Fovea," *Lasers and Light in Ophthalmol.* **2**:79–86 (1988).
67. K. Franze, J. Grosche, S. N. Skatchkov, S. Schinkinger, C. Foja, D. Schiold, O. Uckermann, K. Travis, A. Reichenbach, and J. Guck, "Muller Cells are Living Optical Fibers in the Vertebrate Retina," *Proc. Natl. Acad. Sciences (USA)* **104**(20):8287–8292 (2007).
68. A. W. Snyder and C. Pask, "The Stiles Crawford Effect—Explanation and Consequences," *Vis. Res.* **13**(6):1115–1137 (1973).
69. W. S. Stiles, "The Directional Sensitivity of the Retina and the Spectral Sensitivities of the Rods and Cones," *Proc. R. Soc. Lond. B* **127**:64–105 (1939).
70. M. Alpern, "A Note on Theory of the Stiles Crawford Effects," in J. D. Mollon and L. T. Sharpe (eds.), *Color Vision: Physiology and Psychophysics*, Academic Press, London, 1983, pp. 117–129.
71. S. J. Starr, *Effects of Luminance and Wavelength on the Stiles Crawford Effect in Dichromats*, Ph.D. dissertation, Univ. of Chicago, Chicago, 1977.
72. W. Wijngaard, M. A. Bouman, and F. Budding, "The Stiles Crawford Color Change," *Vis. Res.* **14**(10):951–957 (1974).
73. I. C. Goyal, A. Kumar, and A. K. Ghatak, "Stiles Crawford Effect: An Inhomogeneous Model for Human Cone Receptor," *Optik* **49**:39–49 (1977).
74. B. Vohnsen, I. Iglesias, and P. Artal, "Guided Light and Diffraction Model of Human-Eye Photoreceptors," *J. Opt. Soc. Am. A* **22**(11):2318–2328 (2005).
75. B. Vohnsen, "Photoreceptor Waveguides and Effective Retinal Image Quality," *J. Opt. Soc. Am. A* **24**:597–607 (2007).
76. V. Lakshminarayanan and M. L. Calvo, "Initial Field and Energy Flux in Absorbing Optical Waveguides II. Implications," *J. Opt. Soc. Am.* **A4**(11):2133–2140 (1987).
77. M. L. Calvo and V. Lakshminarayanan, "Initial Field and Energy Flux in Absorbing Optical Waveguides I. Theoretical formalism," *J. Opt. Soc. Am.* **A4**(6):1037–1042 (1987).
78. M. L. Calvo and V. Lakshminarayanan, "An Analysis of the Modal Field in Absorbing Optical Waveguides and Some Useful Approximations," *J. Phys. D: Appl. Phys.* **22**:603–610 (1989).
79. V. Lakshminarayanan and M. L. Calvo, "Incoherent Spatial Impulse Response in Variable-Cross Section Photoreceptors and Frequency Domain Analysis," *J. Opt. Soc. Am.* **A12**(10):2339–2347 (1995).
80. A. Stacey and C. Pask, "Spatial Frequency Response of a Photoreceptor and Its Wavelength Dependence I. Coherent Sources," *J. Opt. Soc. Am.* **A11**(4):1193–1198 (1994).
81. A. Stacey and C. Pask, "Spatial Frequency Response of a Photoreceptor and Its Wavelength Dependence II. Partial Coherent Sources," *J. Opt. Soc. Am.* **A14**:2893–2900 (1997).
82. F. W. Campbell, "A Retinal Acuity Directional Effect," *J. Physiol. (London)* **144**:25P–26P (1958).
83. C. Pask and A. Stacey, "Optical Properties of Retinal Photoreceptors and the Campbell Effect," *Vis. Res.* **38**(7):953–961 (1998).
84. J. M. Enoch, "Retinal Directional Resolution" in J. Pierce and J. Levine (eds.), *Visual Science*, Indiana University Press, Bloomington, IN, 1971, pp. 40–57.
85. H. Engstrom, H. W. Ades, and A. Anderson, *Structural Pattern of the Organ of Corti*, Almquist and Wiskell, Stockholm, 1966.
86. R. Thalmann, L. Thalmann, and T. H. Comegys, "Quantitative Cytochemistry of the Organ of Corti. Dissection, Weight Determination, and Analysis of Single Outer Hair Cells," *Laryngoscope* **82**(11):2059–2078 (1972).
87. N. S. Kapany and J. J. Burke, *Optical Waveguides*, New York, Academic Press, 1972.
88. J. Wells, "Hair Light Guide" (letter), *Nature* **338**(6210):23 (1989).
89. J. M. Enoch and V. Lakshminarayanan, "Biological Light Guides" (letter), *Nature* **340**(6230):194 (1989).
90. C. Oyster, *The Human Eye: Structure and Function*, Sinauer, Sunderland, MA, 1999.
91. D. F. Mandoli and W. R. Briggs, "Optical Properties of Etiolated Plant Tissues," *Proc. Natl. Acad. Sci. USA* **79**:2902–2906 (1982).

92. D. F. Mandoli and W. R. Briggs, "The Photoreceptive Sites and the Function of Tissue Light-Piping in Photomorphogenesis of Etiolated Oat Seedlings," *Plant Cell Environ.* **5**(2):137–145 (1982).
93. W. L. Butler, H. C. Lane, and H. W. Siegelman, "Nonphotochemical Transformations of Phytochrome, In Vivo," *Plant Physiol.* **38**(5):514–519 (1963).
94. D. C. Morgan, T. O'Brien, and H. Smith, "Rapid Photomodulation of Stem Extension in Light Grown *Sinapis Alba*: Studies on Kinetics, Site of Perception and Photoreceptor," *Planta* **150**(2):95–101 (1980).
95. W. S. Hillman, "Phytochrome Conversion by Brief Illumination and the Subsequent Elongation of Etiolated *Pisum* Stem Segments," *Physiol. Plant* **18**(2):346–358 (1965).
96. V. C. Sundar, A. D. Yablon, J. L. Grazul, M. Ilan, and J. Aizenberg, "Fiber Optical Features of a Glass Sponge," *Nature*, **424**:899–900 (2003).

This page intentionally left blank.

DO NOT DUPLICATE

THE PROBLEM OF CORRECTION FOR THE STILES-CRAWFORD EFFECT OF THE FIRST KIND IN RADIOMETRY AND PHOTOMETRY, A SOLUTION

Jay M. Enoch

*School of Optometry
University of California at Berkeley
Berkeley, California*

Vasudevan Lakshminarayanan

*School of Optometry and Departments of Physics and Electrical Engineering
University of Waterloo
Waterloo, Ontario, Canada*

9.1 GLOSSARY*

Aniridia. (Greek “an” = without; “iridia” = iris.) One or both eyes of a patient who is either born without an iris (usually this is a bilateral anomaly), or, for whatever reason the individual has no iris (e.g., as a result of disease, or trauma, or surgery, etc.). Note: In some cases, on examination, a very tiny remnant or very small iris nubbin may be present in congenital cases.

Aniseikonia. (Greek “anis” = unequal; “eikon or ikon” = image.) The result of aniseikonia, or the presence of unequal image sizes in the two eyes, is the formation of characteristic spatial distortions in binocular vision. If this anomaly is large enough, there is a loss of fusion of the two images.

Aphakia. (Greek “a” = not or without; “phakia” = lens.) The state of an eye after surgical removal of an eye lens. In some cases, the eye lens may become displaced rather than removed (this used to be seen commonly with an older form of cataract surgery), or for example, as a complication of certain diseases such as Marfan’s syndrome. Today, the excised eye lens is commonly replaced with a plastic “intraocular lens” (see Chap. 21).

*Added useful Glossaries: please also refer to glossaries of Chap. 8, “Biological Waveguides” by Vasudevan Lakshminarayanan and Jay M. Enoch and Chap. 6, “The Maxwellian View with an Addendum on Apodization” by Gerald Westheimer in this volume and to Chap. 37, “Radiometry and Photometry for Vision Optics” by Yoshi Ohno in Vol. II.

Keratoconus. This is an anomaly of the cornea which results in a cone-like forward-projection of the central or near-central cornea (usually slightly decentered). It is also often accompanied by progressive thinning of the area within and surrounding the cone.

Stiles-Crawford Effect of the first kind (SCE-1). The directional sensitivity of the retina. SCE-1 is the result largely of photoreceptor waveguide properties (see Chap. 8). We test vision for beams of light entering the eye at different points in the entrance pupil of the eye, and then refracted to the retina by the optical system of the eye. *The SCE-1 affects perceived brightness of luminous objects.*

Stiles-Crawford Effect of the second kind (SCE-2). In addition to perceived brightness effects, there occur alterations in perceived hue and saturation of objects viewed. This is associated with the angle of incidence of the viewed beam at the retinal receptor (e.g., see Refs. 45 and 46).

Optical Stiles-Crawford effect. Distribution of reverse path (reflected) luminous energy distribution assessed in the entrance pupil of the eye by reverse path. It is radiant energy which had been transmitted to and through the retina and has been reflected back through the ocular media after having been reradiated after reflection by the photoreceptor waveguides. It is assumed that these retinal receptors are preferentially aligned with the exit pupil of the eye. Also present are back-reflected scattered and nonguided light, for example, see Refs. 25–27.

9.2 INTRODUCTION

The Troland

In the vision science literature, the stimulus to vision is defined as the product of (1) the luminance of the object of regard assessed in the plane of the entrance pupil of the eye, expressed in candelas/m², times (2) the area of the entrance pupil of the eye expressed in mm². These values are expressed in either photopic or scotopic “trolands.” The unit is named after the late Dr. Leonard Troland.^{1–3}

It Is Time for a Reasoned Reassessment of “the Troland,” the Unit of Retinal Illuminance

Leonard Troland, 1889 to 1932, died in a tragic accident as a relatively young man. He left behind a respectable body of research which led to the modern approach to assessment of retinal illuminance, the accepted unit used to specify the stimulus to vision.^{1–3} Please also refer to Dr. Yoshi Ohno’s discussion in Chap. 37 in the section on “Radiometry and Photometry for Vision Optics,” in Vol. II. He discusses current practices. We also call attention to the recent paper by Cheung et al. which addresses the standard units of optical radiation.⁴

There remains, of course, the larger issue of considering whether the currently employed means of estimating/assessing the stimulus to vision (often expressed today as “the photopic or scotopic troland”) might be improved and/or modified in order to bring this form of analysis, and measurement units employed, more effectively into parallel with other modern measures of radiometry and photometry.

Please note the spectral response of the eye, when exposed to brighter/more intense stimuli exciting the cone-dominant-*photopic*-visual-response-system (V_λ), differs from the rod-receptor-dominant-*scotopic*-visual-response-system (V'_λ). The latter determination(s) pertain to less intense visual stimuli. Readers are referred to the book by Wyszecki and Stiles for further discussion.⁵

An “equivalent troland” also has been defined. This unit makes use of specification of the SCE-1 in the entrance pupil where the stimulus has 50 percent effectiveness relative to the level of response determined at the SCE-1 peak location.^{6–9} The more complete SCE-1 function, etc., can be inferred using data obtained with this approach.

The Photometric Efficiency Factor

For completeness, we call attention to another approach to integration of and utilization of SCE-1 data for photometric and other purposes. This is based on the photometric efficiency (PE) factor.^{10,11} The PE is defined by the following ratio:

$$\frac{\text{Effective light collected by pupil}}{\text{Actual light collected through the pupil}} \quad (1)$$

This is an interesting concept, but it is not easy to compute.¹⁰ What is sought here is determination of “an equivalent entrance pupil size” where all points in the entrance pupil of the eye would theoretically contribute equally to the visual stimulus. Particularly for photopic vision, this equivalent pupil, would be smaller than the measured entrance pupil of the eye.^{10,11} By using this approach, an integrated “equivalent entrance pupil” of given diameter and area may be determined. In addition, an adjusted SCE-1 centrum^{10,11} is provided by developing appropriate weighting factors for the asymmetries encountered. Please remember, for different size pupils, the center of the entrance pupil may change by as much as 0.5 mm with dilation of the pupil.¹² For the same dilation of the eye pupil, usually this centrum is quite stable.

Current Practice

The troland,¹⁻³ as it is ordinarily computed (also see Chap. 37, “Radiometry and Photometry for Vision Optics,” by Yoshi Ohno in Vol. II), does not take into consideration the Stiles-Crawford effects, nor factors associated with existing blur of the retinal image, nor uncorrected blur encountered when large entrance pupil areas are included in the experimental design (using a large natural, or a dilated pupil of the eye, or a decentered eye pupil). It is in our collective interest to enhance the quality of determinations, to improve the precision of measurements, to reconsider units used, and to correct these deficiencies.

The Argument Presented in This Chapter

Here we address a number of aspects of these issues, and offer an approach to the solution of resultant photometric problems. We suggest incorporation of the integrated Stiles-Crawford effect of the first kind (SCE-1)¹³⁻¹⁶ in photometric analyses *when degradation of images caused by peripheral portions of the optical elements of the eye (the cornea and the eye lens) are/have been effectively neutralized*. Even though this is still not a perfect solution to the specification of the visual stimulus, it does represent an advance relative to techniques employed today. This will allow us to move effectively toward a finer estimate of the magnitude of the visual stimulus seen either by the experimental subject or patient tested.

9.3 THE PROBLEM AND AN APPROACH TO ITS SOLUTION

Statement of the Problem

When one seeks to determine effective “visual stimulation induced by radiant energy” within the visible spectrum (i.e., *light*, the stimulus to vision),⁹ and integrates all energy entering within the entrance pupil of the eye,¹³ for example, one needs to consider the Stiles-Crawford effect of the first kind (SCE-1).¹³⁻¹⁶ SCE-1 is also known as the directional sensitivity of the retina (see Chap. 8). The issues considered here have rarely been treated adequately in experimental procedures. *If one integrates data derived from point-by-point SCE-1 determinations (obtained by measuring responses obtained when using tiny apertures imaged at a number of different locations in the entrance pupil of the eye)*

the resultant integrated function may not predict properly the visual result obtained when different diameter areas in the entrance pupil of the eye are sampled (please see further arguments presented in this chapter).

Most modern assessments of the integrated SCE-1 are based on point-by-point measurements of SCE-1 determined experimentally at a number of discrete points across the entrance pupil of the eye in a number of meridians (commonly tested in two meridians, i.e., horizontal and vertical). These are then fitted by a parabolic function, proposed in 1937, by *both* Stiles (using a log base 10 relationship)^{14,15} and Crawford (employing a \log_e or \ln base e equation),¹⁶ and that function is then integrated around the entrance pupil. Other functions have been proposed, for example, a Gaussian function by Safir and Hyams,¹⁷ Burns et al.¹⁸ as well as the relationship used in figures in this chapter by Enoch.^{19,20} A review of various analytical expressions used to fit SCE-1 data is discussed in Applegate and Lakshminarayanan.²¹

In his dissertation, Lakshminarayanan recorded SCE-1 functions in a unique patient. This was a young lady born aniridic (without an iris), and (importantly) she had developed only marginal nystagmus (i.e., approximately $\frac{1}{4}$ mm amplitude nystagmus—in others, it is almost always of greater amplitude). She had quite normal visual acuity. In this eye, not only was the SCE-1 recorded, but, in her case, it was possible to assess better the associated side lobes of SCE-1.^{22–24} Assuming (1) her data correspond to characteristic population (central lobe) SCE-1 patterns in normal observers, (2) then none of the above cited equations adequately conforms to the SCE-1 coupled with *the side lobes recorded*.

Notes: (1) In most experiments measuring the SCE-1 performed in recent years, very small-size projections of the aperture stop of the test apparatus have been imaged in the plane of the observer's entrance pupil (i.e., they are usually less than 1 mm in diameter). (2) This discussion is limited to monocular testing. (3) Measured SCE-1 functions obtained vary with wavelength (see Fig. 8). (4) We do not consider here recent discussions of reflected, re-emitted, and projected light gathered from the retinal photoreceptor waveguides which can be assessed by reverse path irradiation/illumination. This relationship is commonly termed “the optical SCE-1” (see Glossary and Refs. 25–27).

Confounds

Fundamentally, when an observer views a display or a scene in an eye with a relatively large natural eye pupil, or with a dilated iris resulting in an enlarged pupil (for purposes considered here, let us say this occurs when the entrance pupil *diameter* is greater than 3 mm), one must consider *both the SCE-1 and additional blur effects* present in the peripheral portions of the cornea and the eye lens. These two factors affect perceived brightness of a stimulus!

For simplicity, assume we are testing at photopic levels, the refraction of the individual has been appropriately corrected for spherical and astigmatic errors (these are first- and second-order aberrations), and the observer is properly accommodated (or corrected visually) for the object viewed. When point-to-point assessments for SCE-1 in the entrance pupil of the eye are properly determined and analyzed, *the resultant uncertainty occurring is largely assignable to the remaining/residual blur of the retinal image* (see data and figures). And as will be seen, this uncertainty is associated with blur induced by peripheral corneal and eye lens aberrations. In discussions of related topics, the senior author has referred to this blur factor as *a dirty variable*, that is, this is an uncontrolled or poorly controlled factor both in vision research and in clinical practice. To further simplify this discussion, we do not address chromatic aberrations (these factors can be addressed/considered separately), nor the use of binocular stimuli. Similarly, here, we do not consider a variety of anomalous optical conditions encountered in clinical settings (e.g., keratoconus, cataracts, or aphakia, etc.), or decentered and/or misshapen irises/pupils, or intraocular lenses centered or decentered, refractive surgery and its complications, movements/decentration(s) of contact lenses, binocular anomalies, etc. (see Chaps. 13, 14, 16, 20, and 21).

Both the magnitude of blur encountered and the nature of a blurred image vary with accommodation of the eye, and, of course, with mismatches between the retinal plane and the plane of focus of the retinal image formed within the eye. We restate the point that the center of the entrance pupil of the eye varies naturally with pupillary dilation, and is virtually always decentered to some extent relative to the optical axis of the eye. Further, the optical axis of the cornea as best determined is not the

same as that of the eye lens (e.g., Refs. 28–30). Simply, the optical system of the eye is not a perfectly centric system. We also know that dilation does not affect the SCE-1 function, *per se* (e.g., Ref. 31). For radiometric and/or photometric purposes, we seek the very best result possible depending on the level of accuracy required for the task.

A Bit of History

When Walter “Stanley” Stiles and Brian “Hewson” Crawford first reported the SCE effect 75 years ago in 1933,¹³ their results were based on data obtained from their own eyes. And it is apparent that, by chance, their eyes manifested little blur of the retinal image, and hence, the complications one often encounters during integration of SCE-1 in the entrance pupil of the eye were largely not encountered in their data.¹³ In such relatively rare individuals, the point-by-point estimates obtained through SCE-1 studies predict well the results of integrating such data across the entrance pupil of the eye.^{19,20} (For example, see the example given under “Sample Data from Enoch” in Sec. 9.4 for experimental subject B.W.)^{19,20}

On the other hand, if there is meaningful degradation of the retinal image caused by peripheral corneal and lenticular aberrations of the eye resulting in image blur, this factor, on its own, can alter the effective integrated perceived luminance of the visual stimulus (e.g., Refs. 18, 19, 32–37). The results found need not be the same between eyes or individuals, nor as predictable as SCE-1.

Collectively, these points have been demonstrated in number of studies, particularly in those of Enoch and coworkers^{18,19,32,33} and Drum.^{34,35}

Today, A Solution of the Problem Is Possible by Using (in Addition) Adaptive Optics (AO) Techniques

What is needed is a joint test of the Stiles-Crawford function using a suitable SCE-1 apparatus (e.g., Ref. 38) *and* utilization of one of a number of modern adaptive optics (AO) devices (see Chap. 15) capable of determining and correcting the incident wavefront in order to deliver a corrected or a nearly corrected image at the retinal test locus sampled. By such an approach, we have the ability to determine more accurately the magnitude of the visual stimulus for the designated entrance pupil size. This is achieved by correcting blur (and higher-order aberrations), and adjusting or correcting for the SCE-1, and properly integrating the resultant data. Earlier, except for use of small entrance pupil diameters or imaged apertures in experiments (see Chap. 6), such corrections were not always possible. Note, devices utilizing adaptive optics are moving toward incorporation of means of correction for chromatic aberrations in their designs or overall experimental techniques.

If a Nonmonochromatic Stimulus to Vision Is Employed

For a more complete treatment of stimuli which are not monochromatic, one might use an achromatic approach/technique proposed by Powell³⁹ (see Chap. 15), or other special lens(es) such as designed and used by the late Gunter Wyszecki to correct the chromatic aberrations of the human eye (e.g., Ref. 5). Using the Wyszecki lens, JME made such adjustments in some of his earlier research, and found this approach quite useful. Please note, the Wyszecki lens must be aligned carefully with the optic axis.

A Reasonable Estimate Is Possible without Full Application of Adaptive Optics

By correcting refraction, *and* spherical aberration in annular zones within the eye pupil, and by using monochromatic stimuli, in 1956, Enoch was able to approximate a more complete correction of blur in the eyes of his experimental subjects.^{19,20}

A rather similar research project to that of Enoch was conducted by Drum just a few years later.^{34,35}

These added issues relate to image formation in the retinal plane, and the assumption is made (often implied, but rarely stated) that energy is transferred from a single plane into the rod and cone photoreceptor waveguides. Please realize that the living cell is a highly complex entity (e.g., Ref. 43). Cell morphology and dimensions (particularly cones) often alter quite a bit across small/modest distances on the retina. Clearly, diffraction and interference effects, as well as rapid changes occurring in the spatial retinal image itself, are occurring across small distances about the entrance to the photoreceptor waveguides. For example, it is useful to consider image alterations occurring over small distances about a plane of focus (both longitudinally and laterally) of the retinal image. As but one example, consider the contents of a paper by Bachinski and Bekkefi,⁴⁰ and data recorded by Enoch and Fry⁴¹ in greatly enlarged and simplified rod and cone receptor models assessed in the microwave spectrum. Without addressing such issues further, here we suggest that a somewhat better estimate of the retinal stimulus to vision *per se* is achievable by incorporating the integrated SCE-1 function in a blur-corrected or blur-minimized retinal image. We argue that this provides a better assessment of the visual stimulus than that provided in what has become now a “classical” approach defined by the late Leonard Troland many years ago, and which is still used for assessing retinal illuminance.^{2,3} See Chap. 37, “Radiometry and Photometry for Vision Optics,” by Yoshi Ohno in Vol. II. So saying, there is more that needs consideration relative to this topic in future.

9.4 SAMPLE POINT-BY-POINT ESTIMATES OF SCE-1 AND INTEGRATED SCE-1 DATA

Sample Data from Stiles and Crawford, 1933

In this section, data are presented which demonstrate reasonably the magnitude and nature of the errors induced in calibration and measurement without correction of the defined errors encountered. While these are not extreme errors, *per se*, these are real effects, they are readily correctable, and should be considered in studies of vision.

Data here are taken from the dissertation of Enoch^{19,20} and the function utilized for integrations of “relative directional sensitivity” here indicated as “ η ” (plotted as a \log_{10} scale) are defined as

$$\eta = 0.25(1 + \cos 9.5 \theta)^2 \quad (2)$$

η equals 1.0 at the peak of the SCE-1 curve; θ is the angle of incidence of the ray of light at the retinal plane. A 1-mm beam displacement at the entrance pupil of the eye alters the angle of incidence of the incident ray of light by 2.5° of oblique incidence at the retinal plane (based on computations using the constants of the Gullstrand schematic eye).²³

This particular equation provides a bit better fit of SCE-1 data^{13–16} than the Stiles’ 1937 (log base 10) relationship¹⁴ or Crawford’s 1937 (ln base e) equation¹⁶ for a parabola. It is used here to simplify transfer of the figures employed below. Because the improvements made when fitting SCE-1 data using this relationship were modest, the authors of this chapter have not used this format for some years, and, rather, have used the more broadly employed Stiles 1937 formulation¹⁴ in its stead

$$\log_{10} \eta = -\rho r^2 \quad (3)$$

In this equation “ η ” is the measured relative visual sensitivity of the stimulus to vision (in trolands, \log_{10} scale); “ ρ ” is a constant defining the curvature of this parabolic function; and “ r ,” in mm, is the distance of the point being sampled from the peak of the measured SCE-1 function in the plane of the entrance pupil of the eye.

SCE-1 data shown in Fig. 1 were taken from Stiles and Crawford¹³ and were fitted by Eq. (2).

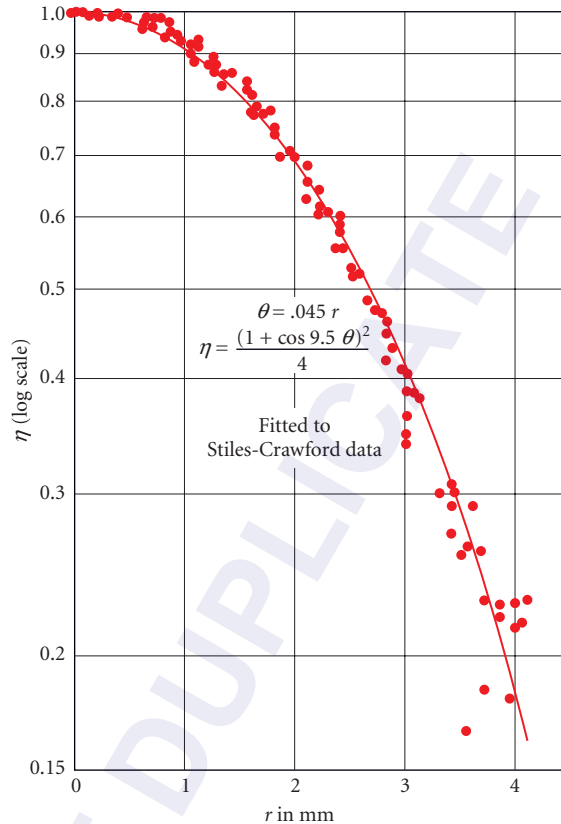


FIGURE 1 The 1933 SCE-1 data of W. S. Stiles and B. H. Crawford were adjusted such that the maxima of measured relative visual sensitivity functions were located at $r = 0.0$ mm in their entrance pupils, and were assigned a relative sensitivity value of 1.0 (at/as the peak of these data functions). Radial distance settings, “ r ,” of the test beam in the entrance pupils of their eyes from the peaks of SCE-1 are displayed on the abscissa. The ordinate displays η (\log_{10} scale). Data used in this illustration originated in the paper of W. S. Stiles and B. H. Crawford, 1933.¹³ (This illustration was copied from Fig. 1, Enoch, 1956, Dissertation; and Fig. 1, Enoch, 1958, J.O.S.A.^{19,20} These figures in this paper are reproduced with permission of the author and the publisher of J.O.S.A.)

In their initial study, Stiles and Crawford, 1933,¹³ were seeking to measure precisely the area of the eye pupil (i.e., their instrument had been designed as a straightforward pupillometer). They assumed, at the outset, that energy passing through each part of the eye pupil contributed equally and proportionately to the retinal image and visual sensitivity. They discovered that the device gave results which were not consistent with their *à priori* assumption, and they properly interpreted their results to indicate that there was evidence for the presence of directional sensitivity of the retina, which became known as the Stiles-Crawford effect of the first kind (SCE-1).

Their instrument sought to assess the integrated visual response resulting from irradiating the entire entrance pupil, and for different size entrance pupils of the eye. *Since they had little image blur in their own eyes*, the integrated result gave little evidence that when the full pupil was measured, the resultant additivity (of the contributions of different parts of the pupil) might also be different in eyes where there were meaningful peripheral corneal and eye lens aberrations resulting in additional degradation (blur) of the retinal image. That is, they did not encounter other than near perfect additivity from integrating point-by-point SCE-1 determinations and comparing them to full eye-pupil assessments. Thus, we need to differentiate between (1) simple addition of the contributions of the entire eye pupil, versus (2) that sum adjusted for SCE-1, and (3) *the result obtained in eyes where there are meaningful aberrations in the periphery of the optical elements of the eye.*

Sample Data from Enoch

Data were obtained from three well-trained graduate student subjects.^{19,20} Subject B.W. had an eye with fine optical properties. Both subjects R. V. and A. M. had greater peripheral ocular aberrations in their measured eyes than B. W.; with subject A. M. exhibiting somewhat greater image degradation than observer R. V. Note these are all photopic determinations. By avoiding scotopic measurements here, we simplify the argument.⁴² Except for overall refractive correction, the single aberration measured independently as part of this dissertation (1955 to 1956) was spherical aberration. The latter was achieved by using a technique similar to one employed by Ivanoff.⁴³ That is, the pupil was divided into discrete nonoverlapping annular zones, and the required refractive correction was altered as needed within each annular zone. Table 1 provides values useful if one is performing the described integrations.

The SCE-1 was measured separately for both (1a) white light (obtained by using a ribbon filament lamp run at 16 amperes, with the beam passing through a heat filter), and paired neutral density filters, or (1b) by appropriately adding an interference filter having a near monochromatic green light band of light (peak $\lambda = 552 \text{ m}\mu$) into the same optical system.

A second variable was introduced in this experiment. There was the traverse of the test beam across the entrance pupil of the eye employing (2a) an instrument design where the ribbon-lamp-filament was imaged in the plane of the entrance pupil of the eye (a classic Maxwellian view instrument design) (see Chap. 6), and, alternatively (2b) the light source was imaged directly on the retina of the observer having passed through the same number of surfaces and by using the same dimension beam image in the entrance pupil of the eye, and the same viewed stimulus target. This was regarded as a non-Maxwellian-view illumination system. Each of these cases utilized a projected aperture which subtended slightly less than 1-mm diameter image in the entrance pupil of the eye. Differences in measured SCE-1 results encountered for these different test conditions were very modest for all three subjects. Direct comparisons were made using both (a) bipartite fields, and, separately, (b) flicker photometry.

TABLE 1 Area in the Entrance Pupil of the Human Eye Based on Gullstrand Schematic Eye

Ent. Pupil Area (Abscissa)*	Solve for Radius and Diameter of Entrance Pupil	
	Radius	Diameter
10 mm ²	1.78 mm	3.57 mm
20	2.52	5.05
30	3.09	6.18
40	3.57	7.14
50	3.99	7.98
60	4.37	8.74
70	4.72	9.44

* The Abscissas in Figs. 3 to 6 are plotted as \log_{10} entrance pupil area.

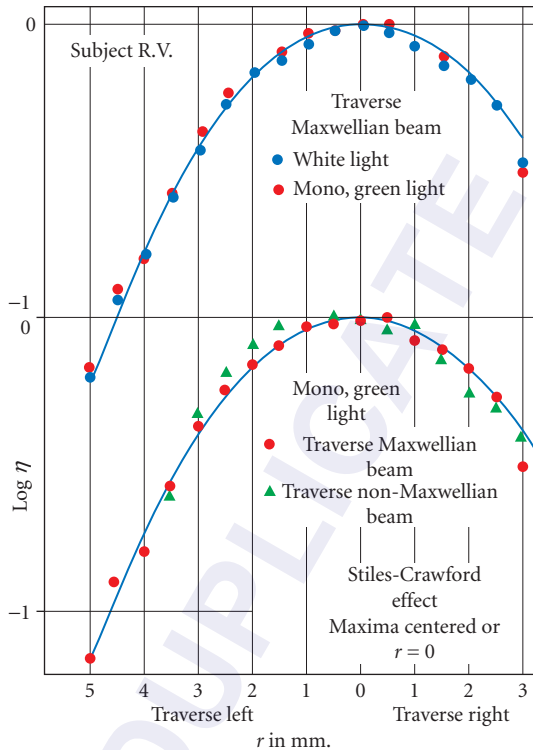
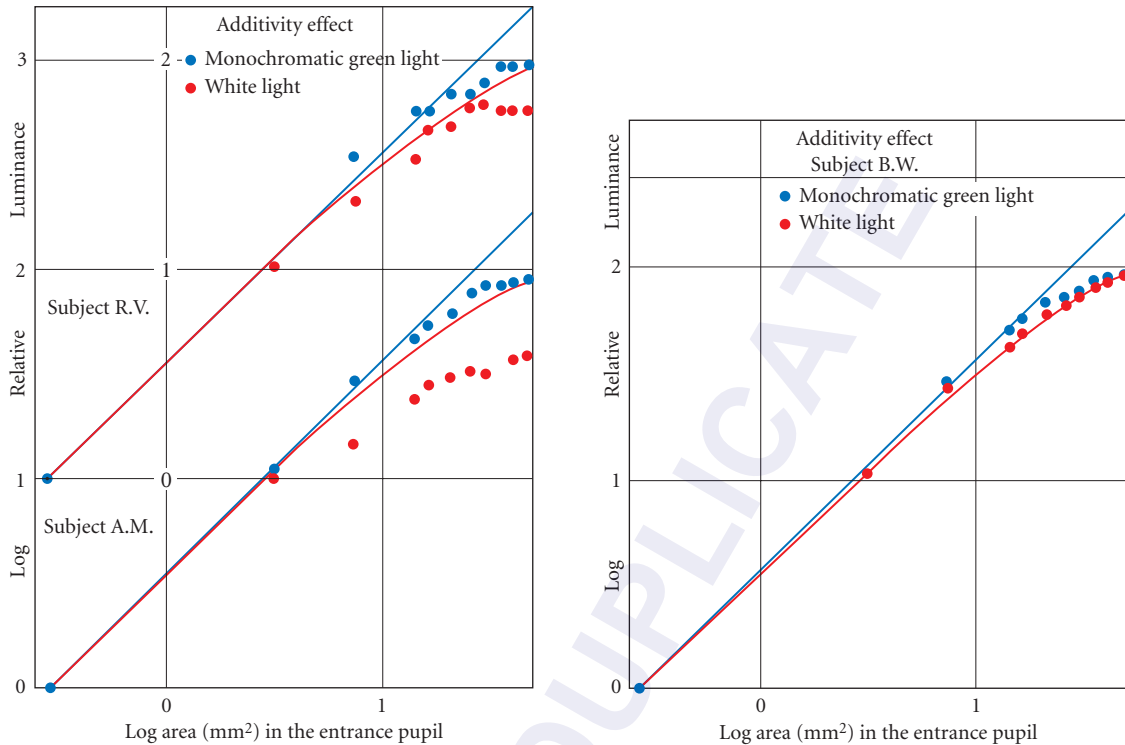


FIGURE 2 SCE-1 data for subject R. V. are presented as an example of the many data sets acquired. His SCE-1 maxima were set at $r = 0.0$ mm on the abscissa. *Top*: SCE-1 data obtained for white light and monochromatic green light (552 m μ max.) are shown. *Bottom*: Both data sets were obtained using monochromatic light. Two different optical conditions were employed. That is, in one experiment, a Maxwellian beam traversed in the eye pupil, and in a second set of sessions, a non-Maxwellian beam was used to traverse the entrance pupil of the eye.^{19,20} Also please refer to Chap. 6. These data are presented with Eq. (2) fitted to these data. (This illustration was taken from Fig. 14, Enoch, 1956, *Dissertation*;¹⁹ and Fig. 9, Enoch, 1958, *J.O.S.A.*^{19,20} These figures are reproduced with permission of the author and the publisher of *J.O.S.A.*)

Figure 2 shows sample SCE-1 test data for subject R. V.^{19,20} As is obvious, there were only small differences between measured SCE-1 functions for the stated conditions. These data are presented as an example of the many data sets acquired. The peaks of his SCE-1 maxima were set at $r = 0.0$ mm on the abscissa.

In Figs. 3 and 4, two different predictive additivity (integrated) plots are presented. The straight line on these double logarithmic plots assumes each point within the entrance pupil of the observer's eye contributes equally to the perceived visual stimulus. The modestly curved line represents the integrated SCE-1 function [here based on Eq. (2) above]. In Fig. 3, data for subjects A. M. and R. V. are presented, and in Fig. 4, data for subject B. W. appears. In Fig. 3 data for one observer were displaced by 1 log unit for clarity. For monochromatic green light [wavelength 522 m μ , (blue circles)] data obtained from all subjects show quite good agreement with the integrated SCE-1 function plot in



FIGURES 3 AND 4 These figures demonstrate the equivalent of integrated or “additivity” diagrams for excitation entering defined entrance pupil areas. The straight line would predict perfect additivity of luminous stimuli entering the entrance pupil of the eye tested of the designated subject. The curved line incorporates a correction for the integration of SCE-1 data about the entrance pupil of the eye. In Fig. 3 data for Subjects R. V. and A. M. are presented; in Fig. 4 data are presented for Subject B. W. (This illustration was taken from Figs. 18 and 19, Enoch, 1956, Dissertation 19; and Figs. 12,13, Enoch, 1958, J.O.S.A.20 These figures are reproduced with permission of the author and the publisher of J.O.S.A.)

both Figs. 3 and 4 (blue circles). The performance of subjects R. V. and A. M. clearly was less good for white light (red circles). On the other hand, data obtained from the eye of B. W. (Fig. 4) shows remarkably good “additivity” for both monochromatic and white light stimuli. We can infer that ocular aberrations in B. W.’s eye exhibited lower amplitudes than those encountered in the eyes of A. M. and R. V. That is, the optical quality of the periphery of B. W.’s cornea and eye lens was better than that found in the peripheral portions of the ocular components of the other two observers. This difference is the result of factors contributing to both monochromatic and chromatic aberrations.

In each case, near monochromatic green stimuli more closely matched the integrated SCE-1 curve than did white light, and in the case of B. W. both data sets were in close approximation to the predicted result.^{19,20}

The added chromatic aberrations apparently made little difference to B. W. (or he retained his focus about mid-spectrum?), while the chromatic aberrations noticeably affected performance of the two subjects, R. V. and A. M.

In Fig. 5, the additivity effect obtained when placing an added -1.00 diopter spherical (D.S.) lens (before the test eye of each subject) is shown for a limited number of different test aperture sizes projected into the entrance pupil of the eye, and this result is contrasted with the addition of a $+2.00$ D.S. lens placed in front of the same eye (with the -1.00 D.S. lens removed). The former lens (the -1.00 D.S.) makes that eye more hyperopic “far-sighted,” while the $+2.00$ D.S. makes the same eye more myopic “near-sighted.” The test eye could at least partially self-correct for the -1.00 D.S.

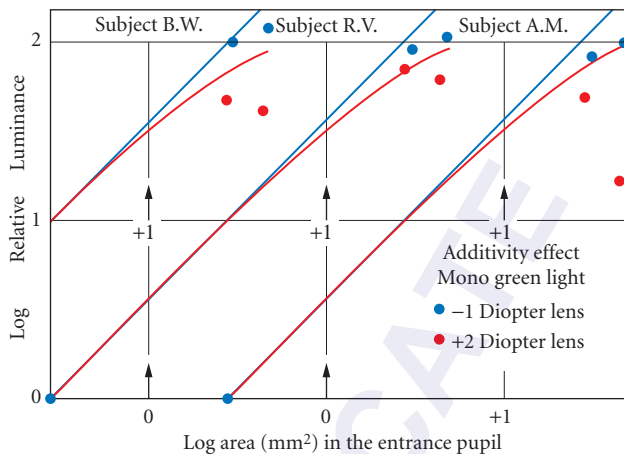


FIGURE 5 A comparison was made between the SCE additivity effect obtained when a -1.00 D.S. lens was placed in the apparatus before the eye of each observer's test eye, and (separately) for addition of a $+2.00$ D.S. lens. Two differing relatively large entrance pupil areas were selected for this test. These lenses slightly altered the area of the imaged aperture stop in the observer's entrance pupil. Hence, the $+2$ D.S. lens is shown as subtending a slightly smaller area in the entrance pupil (on the abscissa) than the -1.00 D.S. lens. (This illustration was taken from Fig. 22, Enoch, 1956, *Dissertation*;¹⁹ and Fig. 16, Enoch, 1958, *J.O.S.A.*²⁰ These figures are reproduced with permission of the author and the publisher of *J.O.S.A.*)

lens by accommodating on the test target. That is, these were nonpresbyopic test subjects. That is, the pupil-dilating drops employed had only slight/modest effect on accommodation in these observers. All tests were monocular.

When a flicker photometric test method was employed,^{19,20} *blur of the retinal image proved to have much less of an influence upon brightness matches made by the subjects* (note, observers were instructed to compare the apparent brightness of the central parts of the flickering fields). When this same technique was combined with use of monochromatic $\lambda = 552$ m μ stimuli, the matches of all three subjects proved to be similar and closely matched the integrated SCE-1 estimated curve. That is, all three subjects exhibited near perfect additivity approximately matching the SCE-1 predicted values (Fig. 6).^{19,20}

Thus, by eliminating and/or minimizing blur effects on perceived brightness due to aberrations (i.e., both monochromatic and chromatic aberrations!), *the integrated SCE function* becomes a good estimate of visual performance for any given pupil aperture utilized. Use of such methods, now (today) are made more readily achievable by using modern adaptive optics techniques for nonchromatic image formation, and, thus, for visual stimulus control. Taking the defined steps allows a superior estimate to be made of perceived brightness, and for the specification of the visual stimulus.

Finally, please note, the factor rho (or ρ) is not a constant across the retina (i.e., it varies with eccentricity from the point of fixation (Fig. 7)).^{12,14,15,44,45} It also varies within the visual spectrum with test wavelength utilized (as implied in Fig. 8).¹² Here we do not consider the Stiles-Crawford effect of the second kind (SCE-2), that is, the alteration of perceived hue and saturation associated with variation of the angle of incidence at the retina.^{14-16, 44-46} Separately, the SCE-1 is also affected to some degree by the distribution of and density of yellow (blue-absorbing) pigment in the eye lens.⁴⁷ The absorbance of this pigment also varies with wavelength. This factor is addressed in Fig. 8. See Chap. 3 in Ref. 12.

As noted above, Drum addressed some similar issues from a somewhat different point of view.^{34,35} He also clearly demonstrated additivity of the SCE-1 to exist for all assessed tests and conditions. The simple fact is that SCE-1 and associated integrated results are rather robust under quite a variety of test conditions.

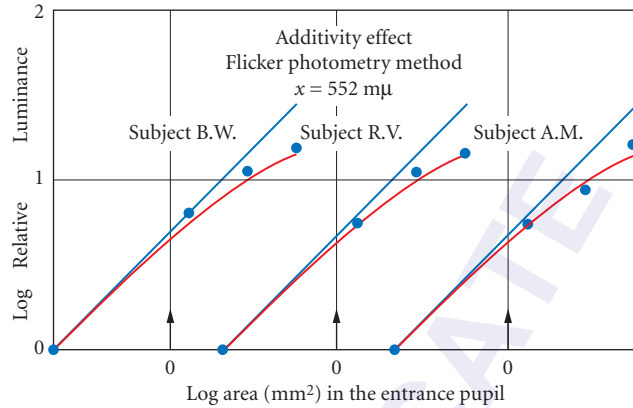


FIGURE 6 These are integrated additivity data for each of the three trained subjects. The flicker photometry method was used. Effects of peripheral corneal and eye lens blur of the observers were minimized by testing with (1) monochromatic light, and by employing (2) a photometric matching technique less influenced by blurred imagery. (This illustration was taken from Fig. 25, Enoch, 1956, Dissertation;¹⁹ and Fig. 19, Enoch, 1958, J.O.S.A.²⁰ These figures are reproduced with permission of the author and the publisher of J.O.S.A.)

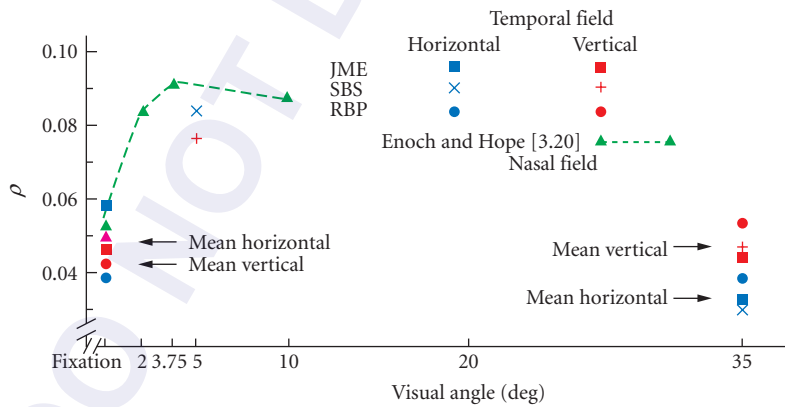


FIGURE 7 This figure addresses the variation of rho, ρ , with eccentricity from the point of fixation (assume this locus corresponds to the center of the fovea). This seemingly complex figure combines data from more than one paper on Stiles' ρ factor [Eq. (3)]. Tests were performed at fixation, in the para-foveal area, and separately at a test locus located 35° from the point of fixation. Horizontal and vertical in this figure refer to SCE-1 tests conducted in these two meridians within the entrance pupils of human eyes. (This figure is reproduced from Fig. 3.8 in Enoch and Tobey, 1980, (see Ref. 12, p. 99), and Enoch and Hope, 1973.^{12,45} It is printed with the permission of the authors and publishers.)

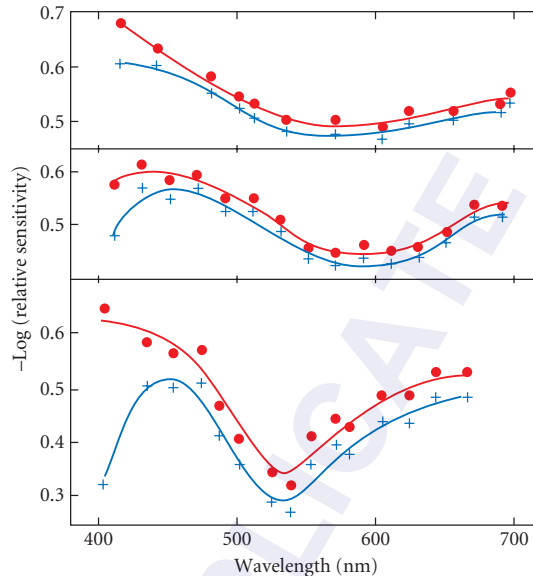


FIGURE 8 These data relate to values obtained from measured human foveal SCE-1 functions (see Ref. 12, Fig. 3.13 located on p. 109).^{15,47} The upper curves present values of (-) log relative sensitivity for different subjects measured at a number of different wavelengths. The upper data curves presented were affected by the presence of yellow eye lens pigment(s), and the lower curves have been corrected for these pigment effects. Please be aware that the density of the yellow lens pigments increases with age. (*This illustration has origin in data of Stiles, 1939,¹⁵ with additions obtained from the paper by Vos and Van Os, 1975.⁴⁷ These figures are reproduced with permission of the authors and the publishers.*)

9.5 DISCUSSION

When Do Matters Considered Here Warrant Consideration?

Like in all psychophysical studies, the demand for precision is greater the more sensitive the determination, and the more complex the problem encountered in the observer or patient. We think we would all agree that the approach argued here is indicated whenever large or dilated pupils are used, particularly for tests of photopic vision! If a pupil is not well centered, or even reasonably asymmetric, or relatively nonresponsive, or otherwise not normal, one needs to rule out any unnecessary variable. Similarly, if the SCE-1 is abnormal, or displaced from its natural centrum, then it is clearly advisable to pursue the matter further. Note, given retinal photoreceptor packing properties, that it is almost unheard of for rods to be normally aligned if cones are found to be disturbed in their orientations in a given retinal area, and vice versa.

Zwack and his coworkers (e.g., Ref. 48) have studied the effects of laser exposures on retinal survival, recovery, and/or lack thereof mainly in snakes. In their experiments on animals, Zwack et al. used (1) subthreshold laser exposures, (2) threshold burns, (3) above threshold burns, and (4) more severe exposures. These studies were backed up by assessments of individual humans who have been exposed to comparable laser burns. Such laser burns result in immediate and substantial disturbances in photoreceptor alignments (often extending for some distance about the burn site), resultant early and late sensitivity losses. Snake eyes have large photoreceptors visible directly through the animal's natural pupil *in vivo*, and are available for study. Snakes are relatively nonreactive when exposed

to high intensity laser light when sedated, and show relatively rapid recovery *if* there is a potential for recovery. A meaningful area of the retina is affected. These authors⁴⁸ have carefully documented recovery and failure to recover in their studies.

One must remember that alignment of photoreceptors with the center of the exit pupil of the eye is maintained normally throughout life, except in cases where there has been some retinal disorder or disease.^{49,50} In a number of cases, the retina can recover if the inducing anomaly is corrected or self corrects (remits).^{51,52} So saying, residual scar tissue, local adhesions, or other impediments to recovery can alter the result. If ever doubts exist, suitable tests are indicated.

As has been inferred above, and can be seen in Fig. 3 to 6, in the normal observer the smooth, integrated SCE-1 function curve differs only modestly from the linear additivity curve until there is about a 4- to 5-mm-diameter entrance pupil.

Once again, looking at Figs. 3 and 4, uncertainty enters with peripheral blur of the optical components of the eye. The effects are greater for white light than monochromatic light. Subject A. M. apparently manifests less aberrations in his peripheral eye lens and cornea than observer V. R., and subject B. W. exhibits virtually no effect due to aberrations. At the outset, it was pointed out that image blur or degradation enters as a “dirty variable.” Without correction of these degrading peripheral lens blur factors, one introduces uncertainty. If one adds refractive error to the mix, the situation is made worse (see Fig. 5). If refractive and peripheral blur are corrected or at least meaningfully minimized, then a SCE correction applied to such data predict quite well the visual stimulus presented to the observer (see Fig. 6)!

Among other very interesting results, Makous and his coworkers,^{53–58} and Applegate and Lakshminarayanan,²¹ correctly point out, as can be inferred from comments made above, for small pupil apertures, certainly for 3-mm-diameter entrance pupils (1.5-mm radial distances) and perhaps a bit more, the two functions differ only slightly if refraction is well corrected, and the system is centered on the subject’s entrance pupil. And Makous et al. have also addressed and thoughtfully considered issues regarding effects of coherence and noncoherence of imagery, loci of adaptation, quantum activation rates in photoreceptors, effects of cell obliquity on excitation processes, etc.

9.6 TELEOLOGICAL AND DEVELOPMENTAL FACTORS

In a sense, it is also useful to consider these visual functions from a teleological point of view. It is apparent that the ultimate purpose(s) of the SCE-1 and associated waveguide properties of retinal photoreceptors are to enhance the detection of critical visual signal, both photopic and scotopic, and to help suppress intraocular stray light “noise” present in the “integrating-sphere-like” eye (e.g., Ref. 12). The feedback mechanisms controlling receptor alignment (under normal circumstances) and their directional properties *together* serve to favor the detection of directed visual signal content passing through the center of the exit pupil in eyes of vertebrates. In many invertebrate species the eyes serve similarly, for example, octopus, etc.; or in a different (but effectively comparable) fashion in the vast number of other invertebrate species. These optical and anatomical features serve to enhance greatly visual processes, and, as such, play critical roles in vision and survival of species. This argument emphasizes the great importance exerted by such mechanisms in evolutionary processes.

Related to such matters, the attention of the reader is called to a *remarkable recent paper* by Detlev Arendt et al. in *Science* in 2004.⁵⁹ These authors located a primitive and ancient form of invertebrate aquatic worm which had the usual paired invertebrate eye structures *as well as vertebrate-type, cylindrically shaped photoreceptors* containing cilia located in its brain in that area which controlled circadian rhythms. And the latter cells were also shown to contain a cone-type of opsin!⁵⁹

9.7 CONCLUSIONS

There is a need to optimally correct and control the quality of the images formed in the eye, and/or the eye plus associated optical apparatuses. This is needed in order to define better the observer’s stimulus to vision, and to understand, in a superior way, the observer’s visual responses. This will include (1) a satisfactory correction of refraction, that is, lower-order aberrations, (2) correction

of higher-order monochromatic aberrations, and also (3) correction of chromatic aberrations. The totality can be aided by utilization of modern adaptive optics (AO) techniques. In radiometric and photometric studies, there is also need to include a factor which corrects for the Stiles-Crawford effect of the first kind (SCE-1), and that need increases with pupil diameter, particularly for photopic vision.

Note, because we do not wholly understand the nature of the stimulus to accommodation, we must be careful not to be too overly aggressive in seeking to eliminate *all* image blur.

And, as pointed out in the introductory remarks, it is time to reconsider the definition and units used to describe/define the troland.

Here, we only have considered monocular visual corrections. That is, in this discussion, we have not considered issues associated with maintenance of sustained and comfortable binocular vision (see Chap. 13) both in instrument design as well as in assessment and corrections of vision. For effective binocular results, it is clear that issues associated with maintaining an observer's comfort while performing extended binocular visual tasks need to be carefully addressed in both design and assessment of vision roles/function, including careful attention being paid to equating the sizes of the two retinal images (i.e., countering the effects of aniseikonia), and fusion of those images. That is, in design, we always need to address those additional factors affecting ocular motility and binocular fusion (Chap. 13).

The list of references appended to this chapter includes a number of citations which address additional aspects of issues considered here.⁶⁰⁻⁶⁸

9.8 REFERENCES

1. J. Howard, "The History of OSA; Profile in Optics, Leonard Thompson Troland, Optics and Photonic News (O.P.N.)," *Optical Society of America* **19**(6):20-21 (2008).
2. L. T. Troland, "Theory and Practice of the Artificial Pupil," *Psych. Review* **22**(3):167-176 (1915).
3. L. T. Troland, "On the Measurement of Visual Sensation Intensities," *J. Exp. Psych.* **2**(1):1-34 (1917).
4. J. Y. Cheung, C. J. Chunnillall, E. R. Woolliams, N. P. Fox, J. R. Mountford, J. Wang, and P. J. Thomas, "The Quantum Candela: A Redefinition of the Standard Units for Optical Radiation," *J. Mod. Opt.* **54**(2-3):373-396 (2007).
5. G. Wyszecki and W. S. Stiles, *Color Science: Concepts, and Methods, Quantitative Data and Formulae*, 2nd (ed.), Wiley, New York, 1982.
6. W. S. Baron and J. M. Enoch, "Calculating Photopic Illuminance," *Am. J. Optom. Physiol. Optics* **59**(4):338-341 (1982).
7. J. M. Enoch, "Vision; Physiology," *Modern Ophthalmology*, vol. 1, Arnold Sorsby, (ed.), 1st (ed.), Butterworths, Washington DC, 1963, sec. 1, Chap. 3, pp. 202-289.
8. J. M. Enoch and Harold E. Bedell, "Specification of the Directionality of the Stiles-Crawford Function," *Am. J. Optom. Physiol. Optics* **56**:341-344 (1979).
9. *Handbook of the Illumination Engineering Society* (any issue), Appendix, Conversion Factors, p. A-1.
10. D. Atchison, D. H. Scott, and G. Smith, "Pupil Photometric Efficiency and Effective Centre," *Ophthalmol. Physiol. Optics* **20**(6):501-503 (2000).
11. L. C. Martin, *Technical Optics*, 1st (ed.), vol. 2, Pitman, London, 1954.
12. J. M. Enoch, F. L. Tobey, Jr., (eds.), *Vertebrate Photoreceptor Optics. Springer Series in Optical Sciences*, vol. 23, Springer-Verlag, Berlin, Heidelberg, New York, 1981, ISBN 3-540-10515-8 Berlin, etc., and ISBN 0-387-10515-8 New York, etc.
13. W. S. Stiles and B. H. Crawford, "The Luminous Efficiency of Light Entering the Eye Pupil at Different Points," *Proc. Roy. London, Ser. B* **112**:428-450 (1933).
14. W. S. Stiles, "The Luminous Efficiency of Monochromatic Rays Entering the Eye Pupil at Different Points and a New Color Effect," *Proc. Roy. Soc. London, Ser. B* **123**:(830) 90-118 (1937).
15. W. S. Stiles, "The Directional Sensitivity of the Retina and the Spectral Sensitivities of the Rods and Cones," *Proc. Roy. Soc. London Ser. B* **127**:64-105 (1939).

16. B. H. Crawford, "The Luminous Efficiency of Light Entering the Eye Pupil at Different Points and Its Relation to Brightness Threshold Measurements," *Proc. Roy. Soc. London, Ser. B* **124**:(834) 81–96 (1937).
17. A. Safir, L. Hyams, and J. Philpott, "The Retinal Directional Effect: A Model Based on the Gaussian Distribution of Cone Orientations," *Vis. Res.* **11**:819–831 (1971).
18. S. Marcos and S. Burns, "Cone Spacing and Waveguide Properties from Cone Directionality Measurements," *J. Opt. Soc. Am. A* **16**:995–1004 (1999).
19. J. M. Enoch, *Summated Response of the Retina to Light Entering Different Parts of the Pupil*, Dissertation, Professor Glenn A. Fry, Advisor, Ohio State University, 1956.
20. J. M. Enoch, "Summated Response of the Retina to Light Entering Different Parts of the Pupil," *J. Opt. Soc. Am.* **48**:392–405 (1958).
21. R. A. Applegate and V. Lakshminarayanan, "Parametric Representation of the Stiles-Crawford Functions: Normal Variation of Peak Location and Directionality," *J. Opt. Soc. Am. A* **10**:1611–1623 (1993).
22. Vasudevan Lakshminarayanan, *The Stiles-Crawford Effect in Anirida*, Ph.D. Dissertation, Jay M. Enoch, Advisor, University of California, Berkeley, 1985.
23. J. M. Enoch, V. Lakshminarayanan, and S. Yamade, "The Stiles-Crawford Effect (SCE) of the First Kind: Studies of the SCE in an Aniridic Observer," *Perception* **15**:777–784 (1986) (the W.S. Stiles Memorial issue).
24. V. Lakshminarayanan, J. M. Enoch, and S. Yamade, "Human Photoreceptor Orientation: Normals and Exceptions," *Advances in Diagnostic Visual Optics*, A. Fiorentini, D. L. Guyton, and I. M. Siegel. (eds.), Heidelberg, Springer-Verlag, 1987, pp. 28–32.
25. J. C. He, S. Marcos, and S. A. Burns, "Comparison of Cone Directionality Determined by Psychophysical and Reflectometric Techniques," *J. Opt. Soc. Am. A* **16**:2363–2369 (1999).
26. M. J. Kanis, *Foveal Reflection Analysis in a Clinical Setting*, Dissertation, Utrecht University, Faculty of Medicine, the Netherlands. p. 155, 2008. ISBN: 978-90-39348536, Dr. Dirk van Nooren, Advisor.
27. W. Gao, B. Cense, Y. Zhang, R. S. Jonnal, and D. T. Miller, "Measuring Retinal Contributions to the Optical Stiles-Crawford Effect with Optical Coherence Tomography," *Opt. Express* **16**:6486–6501 (2008).
28. H. von Helmholtz, *Helmholtz's Treatise on Physiological Optics*, 3d German (ed.), vol. 1, English Translation by J.P.C. Southall, Rochester, NY, Optical Society of America, 1924.
29. C. Cui and V. Lakshminarayanan, "The Choice of Reference Axis in Ocular Wavefront Aberration Measurement," *J. Opt. Soc. Am. A* **15**:2488–2496 (1998).
30. C. Cui and V. Lakshminarayanan, "The Reference Axis in Corneal Refractive Surgeries—Visual Axis or the Line of Sight?" *J. Mod. Opt.* **50**:1743–1749 (2003).
31. L. Ronchi, "Influence d'un mydriatique sur l'effet Stiles-Crawford," *Optica Acta* **2**(1):47–49 (1955).
32. A. M. Laties and J. M. Enoch, "An Analysis of Retinal Receptor Orientation: I. Angular Relationship of Neighboring Photoreceptors," *Invest. Ophthalmol.* **10**(1):69–77 (1971).
33. J. M. Enoch and A. M. Laties, "An Analysis of Retinal Receptor Orientation: II. Predictions of Psychophysical Tests," *Invest. Ophthalmol.* **10**(12):959–970 (1971).
34. B. Drum, *Additivity of the Stiles-Crawford Effect for a Fraunhofer Image*, Dissertation, Ohio State University, 1973, Carl R. Ingling, Advisor.
35. B. Drum, "Additivity of the Stiles-Crawford Effect for a Fraunhofer Image," *Vis. Res.* **15**:291–298 (1975).
36. G. Bocchino, "Studio della variazione della luminosità di un cannocchiale a variare della pupilla d'uscita," *Ottica I*:136–142 (1936).
37. G. T. di Francia and W. Sbrolli, "Sulla legge integrale dell'effetto Stiles-Crawford," *Atti della Fond. G. Ronchi* **2**:100–104 (1947).
38. J. M. Enoch and G. M. Hope, "An Analysis of Retinal Receptor Orientation: III Results of Initial Psychophysical Tests," *Invest. Ophthalmol.* **11**(9):765–782 (1972).
39. I. Powell, "Lenses for Correcting Chromatic Aberration of the Eye," *Appl. Opt.* **20**:4152–4155 (1981).
40. M. P. Bachynski and G. Bekefi, "Study of Optical Diffraction Images at Microwave Frequencies," *J. Opt. Soc. Am.* **47**:428–438 (1957).
41. J. M. Enoch and G. A. Fry, "Characteristics of a Model Retinal Receptor Studied at Microwave Frequencies," *J. Opt. Soc. Am.* **48**(12):899–911 (1958).
42. J. M. Enoch, H. E. Bedell, and E. C. Campos, "Local Variations in Rod Receptor Orientation," *Vis. Res.* **18**(1): 123–124 (1978).

43. A. Ivanoff, *Les Aberrations de l'Oeil, Leur Role dans l'Accommodation*, (The relation between pupil efficiencies for small and extended pupils of entry.) Editions de la Revue d'Optique, Paris, 1953.
44. J. M. Enoch and G. M. Hope, "Directional Sensitivity of the Foveal and Parafoveal Retina," *Invest. Ophthalmol.* **12**:497–503 (1973).
45. J. M. Enoch and W. S. Stiles, "The Colour Change of Monochromatic Light with Retinal Angle of Incidence," *Optica Acta* **8**:329–358 (1961).
46. P. L. Walraven and M. A. Bouman, "Relation Between Directional Sensitivity and Spectral Response Curves in Human Cone Vision," *J. Opt. Soc. Am* **60**:780–784 (1960).
47. J. J. Vos and F. L. Van Os, "The Effect of Lens Density on the Stiles-Crawford Effect," *Vis. Res.* **15**:749–751 (1975).
48. H. Zwick, P. Edsall, B. E. Stuck, E. Wood, R. Elliott, R. Cheramie, and H. Hacker, "Laser Induced Photoreceptor Damage and Recovery in the High Numerical Aperture Eye of the Garter Snake," *Vis. Res.* **48**:486–493 (2008).
49. M. Rynders, T. Grosvenor, and J. M. Enoch, "Stability of the Stiles-Crawford Function in a Unilateral Amblyopic Subject over a 38 year Period: A Case Study," *Optom. Vis. Sci.* **72**(3):177–185 (1995).
50. J. M. Enoch, J. S. Werner, G. Haegerstrom-Portnoy, V. Lakshminarayanan, and M. Rynders, "Forever Young: Visual Functions Not Affected or Minimally Affected By Aging," *J. Gerontology: Biological Sciences* **55A**(8): B336–B351 August (1999).
51. E. C. Campos, H. E. Bedell, J. M. Enoch, and C. R. Fitzgerald, "Retinal Receptive Field-like Properties and Stiles-Crawford Effect in a Patient with a Traumatic Choroidal Rupture," *Doc. Ophthalmol.* **45**:381–395 (1978).
52. J. M. Enoch, C. R. Fitzgerald, and E. C. Campos, *Quantitative Layer-by-Layer Perimetry: An Extended Analysis*, Grune and Stratton, New York, 1981.
53. W. Makous and J. Schnapf, "Two Components of the Stiles-Crawford Effect: Cone Aperture and Disarray," (Abstract) *Program A.R.V.O.* 1973, p. 88.
54. M. J. McMahan and D. I. A. MacLeod, "Retinal Contrast Losses and Visual Resolution with Obliquely Incident Light," *J. Opt. Soc. Am. A* **18**(11):2692–2703 (2001).
55. J. Schnapf and W. Makous, "Individually Adaptable Optical Channels in Human Retina," (Abstract) *Program A.R.V.O.* 1974, p. 26.
56. B. Chen and W. Makous, "Light Capture in Human Cones," *J. Physiol. (London.)* **414**:89–108 (1989).
57. W. Makous, "Fourier Models and the Loci of Adaptation," *J. Opt. Soc. Am. A* **14**:2323–2345 (see p. 2332 pertinent to this discussion), (1997).
58. W. Makous, "Scotopic vision," in John Werner and L. M. Chapula (eds.), *The Visual Neurosciences*, Boston, MIT Press, 2004, (this paper corrects an erroneous table in the prior reference) pp. 838–850.
59. D. Arendt, K. Tessmar-Raible, H. Snyman, A. Dorresteijn, and J. Wittbrodt, "Ciliary Photoreceptors with a Vertebrate-type Opsin in an Invertebrate-brain," *Science* **306**(29 October):869–871 (2004). See the interesting discussion of this paper by Elizabeth Pennisi, pp. 796–797.
60. L. Lundström and P. Unsbo, "Transformation of Zernike Coefficients: Scaled, Translated, and Rotated Wavefronts with Circular and Elliptical Pupils," *J. Opt. Soc. Am. A* **24**(3):569–577 (2007).
61. R. A. Applegate, W. J. Donnelly III, J. D. Marsack, D. E. Koenig, and K. Pesudovs, "Three-Dimensional Relationship between High-order Root-mean-square Wavefront Error, Pupil Diameter, and Aging," *J. Opt. Soc. Am.* **23**(3):578–587 (2007).
62. X. Zhang, M. Ye, A. Bradley, and L. Thibos, "Apodization by the Stiles-Crawford Effect Moderates the Visual Impact of Retinal Image Defocus," *J. Opt. Soc. Am. A* **16**:812–820 (1999).
63. J. M. Enoch, "Retinal Directional Resolution," International Conference on Visual Science, Bloomington, Indiana (April 1968), in *Visual Science* by J. Pierce and J. Levene, (eds.), Indiana University Press, Bloomington, Indiana, pp. 40–57 (1971).
64. H. Metcalf, "Stiles-Crawford Apodization," *J. Opt. Soc. Am.* **55**:72–74 (1965).
65. J. P. Carroll, "Apodization of the Stiles-Crawford Effect," *J. Opt. Soc. Am.* **70**:1155–1156 (1980).
66. D. A. Palmer, "Stiles-Crawford Apodization and the Stiles-Crawford Effect," *J. Opt. Soc. Am. A* **2**:1371–1374 (1985).

67. L. L. Sloan, "Size of Pupil as a Variable Factor in Measurements of the Threshold: An Experimental Study of the Stiles-Crawford Phenomenon," (Abstract) *J. Opt. Soc. Am.* **30**:271, (June, 1940), Paper: *Arch. Ophthalmol.* **24** (New Series, N.S.) July-December, pp. 258–275 (1940).
68. E. J. Fernández, A. Unterhuber, B. Považay, B. Hermann, P. Artal, and W. Drexler, "Chromatic Aberration Correction of the Human Eye for Retinal Imaging in the Near Infrared," *Optics Express* **14**(13):6213–6225, (June 26, 2006).

DO NOT DUPLICATE

COLORIMETRY

David H. Brainard

*Department of Psychology
University of Pennsylvania
Philadelphia, Pennsylvania*

Andrew Stockman

*Department of Visual Neuroscience
UCL Institute of Ophthalmology
London, United Kingdom*

10.1 GLOSSARY

Chromaticity coordinates. Tristimulus values normalized to sum to unity.

CIE. Commission Internationale de l'Éclairage or International Commission on Illumination. Organization that develops standards for color and lighting.

Color-matching functions (CMFs). Tristimulus values of the equal-energy spectrum locus.

Color space transformation matrix. Multiply a vector of tristimulus values for one color space by such a matrix to obtain tristimulus values in another color space.

Cone coordinates. Tristimulus values of a light with respect to the cone fundamentals.

Cone fundamentals. Estimates of the cone spectral sensitivities at the cornea. Equivalently, the CMFs that would result if primaries that uniquely stimulated the three cones could be and were used.

Linear model. Set of spectral functions that may be scaled and added to approximate other spectral functions. For example, the spectral power distributions of three monitor primaries are a linear model for the set of lights that can be emitted by the monitor.

Metamers. Two physically different lights that match in appearance to an observer.

Photopic luminosity function. Measure of luminous efficiency as a function of wavelength under photopic (i.e., rod-free) conditions.

Primary lights. Three independent lights (real or imaginary) to whose scaled mixture a test light is matched (actually or hypothetically). They must be independent in the sense that no combination of any two can match the third.

Standard observer. The standard observer is the hypothetical individual whose color-matching behavior is represented by a particular set of CMFs.

Tristimulus values. The tristimulus values of a light are the intensities of the three primary lights required to match it.

Visual angle. The angle subtended by an object in the external field at the effective optical center of the eye. Colorimetric data are typically specified for centrally fixated 2° or 10° fields of view.

10.2 INTRODUCTION

Scope

The goal of colorimetry is to incorporate properties of the human color vision system into the measurement and numerical specification of visible light. Thanks in part to the inherent simplicity of the initial stages of visual coding, this branch of color science has been quite successful. We now have effective quantitative representations that predict when two lights will appear identical to a human observer and a good understanding of how these matches are related to the spectral sensitivities of the underlying cone photoreceptors. Although colorimetric representations do not directly predict color sensation,¹⁻³ they do provide the foundation for the scientific study of color appearance. Moreover, colorimetry can be applied successfully in practical applications. Foremost among these is perhaps color reproduction.⁴⁻⁶

As an illustrative example, Fig. 1 shows an image processing chain. Light from an illuminant reflects from a collection of surfaces. This light is recorded by a color camera and stored in digital form. The digital image is processed by a computer and rendered on a color monitor. The reproduced image is viewed by a human observer. The goal of the image processing is to render an image with the same color appearance at each image location as the original. Although exact reproduction is not always possible with this type of system, the concepts and formulas of colorimetry do provide a reasonable solution.^{4,7} To develop this solution, we will need to consider how to represent the spectral properties of light, the relation between these properties and color camera responses, the representation of the restricted set of lights that may be produced with a color monitor, and the way in which the human visual system encodes the spectral properties of light. We will treat each of these topics in this chapter, with particular emphasis on the role played by the human visual system.

Reference Sources

A number of excellent references are available that provide detailed treatments of colorimetry and its applications. Wyszecki and Stiles' comprehensive book⁸ is an authoritative reference and

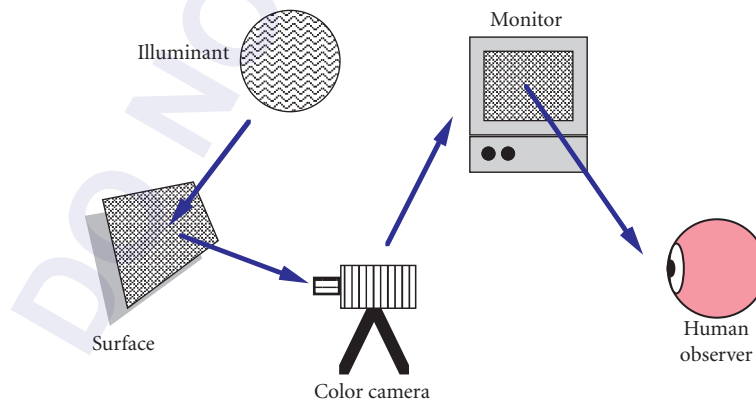


FIGURE 1 A typical image processing chain. Light reflects from a surface or collection of surfaces. This light is recorded by a color camera and stored in digital form. The digital image is processed by a computer and rendered on a color monitor. The reproduced image is viewed by a human observer.

provides numerous tables of standard colorimetric data. Smith and Pokorny⁹ provide a treatment complementary to the one developed here. Several publications of the Commission Internationale de l'Éclairage (International Commission on Illumination, commonly referred to as the CIE) describe current international technical standards for colorimetric measurements and calculations.¹⁰ The most recent CIE proposal is for a set of physiologically relevant color-matching functions or cone fundamentals based mainly on the results of human psychophysical measurements.¹¹ Other sources cover colorimetry's mathematical foundations,^{12,13} its history,^{14–16} its applications,^{2,5,17,18} and its relation to neural mechanisms.^{19–21} Chapters 3, 5, 11, and 22 in this volume, and Chap. 37, “Radiometry and Photometry for Vision Optics,” by Yoshi Ohno in Vol. II of this *Handbook* are also relevant.

Chapter Overview

The rest of this chapter is organized into three main sections. Section 10.3, “Fundamentals of Colorimetry,” reviews the empirical foundation of colorimetry and introduces basic colorimetric methods. In this section, we adhere to notation and development that is now fairly standard in the field.

Section 10.4, “Color Coordinate Systems,” discusses practicalities of using basic colorimetric ideas and reviews standard coordinate systems for representing color data.

Desktop computers can easily handle all standard colorimetric calculations. In Sec. 10.5 we introduce vector and matrix representations of colorimetric data and formulas. This development enables direct translation between colorimetric concepts and computer calculations. Matrix algebra is now being used increasingly in the colorimetric literature.^{4,22,23}

Section 10.6 uses the vector and matrix formulation developed in Sec. 10.5 to treat some advanced topics.

The appendix (Sec. 10.7) reviews the elementary facts of matrix algebra required for this chapter. Numerous texts treat the subject in detail.^{24–27} Many software packages (e.g., MATLAB, S-Plus, R) provide extensive support for numerical matrix algebra.

10.3 FUNDAMENTALS OF COLORIMETRY

Introduction

We describe the light reaching the eye from an image location by its spectral power distribution. The spectral power distribution generally specifies the radiant power density at each wavelength in the visible spectrum. For human vision, the visible spectrum extends roughly between 400 and 700 nm (but see subsection “Sampling the Visible Spectrum” in Sec. 10.5). Depending on the viewing geometry, measures of radiation transfer other than radiant power may be used. These measures include radiance, irradiance, exitance, and intensity. The distinctions between these measures and their associated units as well as equivalent photometric measures are treated in Chaps. 34, 36, and 37 of Vol. II on this *Handbook* and are not considered here.

Color and color perception are limited at the first stage of vision by the spectral properties of the layer of light-sensitive *photoreceptors* that cover the rear surface of the eye (upon which an inverted image of the world is projected by the eye's optics). These photoreceptors transduce arriving photons to produce the patterns of electrical signals that eventually lead to perception. Daytime (photopic) color vision depends mainly upon the three classes of cone photoreceptor, each with different spectral sensitivity. These are referred to as long-, middle-, and short-wavelength-sensitive cones (L, M, and S cones), according to the part of the visible spectrum to which they are most sensitive (see Fig. 6). Night-time (scotopic) vision, by contrast, depends on a single class of photoreceptor, the rod.

TABLE 1 Glossary of Conventional Colorimetric Terms and Notation

Chromaticity coordinates	x, y , or in terms of the tristimulus values $X/(X+Y+Z)$ and $Y/(X+Y+Z)$, respectively (or r, g for <i>RGB</i> space, or l, m for <i>LMS</i> space).
Color-matching functions or CMFs	$\bar{x}(\lambda)$, $\bar{y}(\lambda)$, and $\bar{z}(\lambda)$. Tristimulus values of the equal-energy spectrum locus.
Cone fundamentals	$\bar{l}(\lambda)$, $\bar{m}(\lambda)$, and $\bar{s}(\lambda)$ in CMF notation, or often $L(\lambda)$, $M(\lambda)$, and $S(\lambda)$. These are the CMFs that would result if primaries that uniquely stimulated the three cones could be used.
Photopic luminosity function	Photometric measure of luminous efficiency as a function of wavelength under photopic (<i>i.e.</i> , rod-free) conditions: $V(\lambda)$ or $\bar{y}(\lambda)$.
Primary lights	R, G, B , the three independent primaries (real or imaginary) to which the test light is matched (actually or hypothetically). They must be independent in the sense that no combination of two can match the third.
Standard observer	The standard observer is the hypothetical individual whose color-matching behavior is represented by a particular set of mean CMFs.
Tristimulus values	R, G, B , the amounts of the three primaries required to match a given stimulus.
Visual angle	The angle subtended by an object in the external field of view at the effective optical center of the eye. Colorimetric data are typically for centrally fixated 2 or 10° fields of view.

Conventional Colorimetric Terms and Notation

Table 1 provides a glossary of conventional colorimetric terms and notation. We adhere to these conventions in our initial development, in this section and in Sec.10.4. See also Table 3.1 of Ref. 28, and compare with the matrix algebra glossary in Table 2.

Trichromacy and Univariance

Normal human vision is *trichromatic*. With some important provisos (see subsection “Conditions for Trichromatic Color Matching” in Sec. 10.3), observers can match a test light of any spectral composition to an appropriately adjusted mixture of just three other lights. Consequently, colors can be defined by three variables: the intensities of the three primary lights with which they match. These are called *tristimulus values*.

The range of colors that can be produced by the additive combination of three lights is simulated in Fig. 2. Overlapping red, green, and blue lights produce regions that appear cyan, purple, yellow, and white. Other, intermediate, colors can be produced by varying the relative intensities of the three lights.

Human vision is *trichromatic* because there are only three classes of cone photoreceptor in the eye, each of which responds univariantly to the rate of photon absorption.^{29,30} *Univariance* refers to the fact that the effect of a photon, once absorbed, is independent of wavelength. What varies with wavelength is the probability that a photon is in fact absorbed, and this variation is described by the photoreceptor’s spectral sensitivity. Photoreceptors are, in effect, sophisticated photon counters the outputs of which vary according to the rate of absorbed photons. Changes in the absorption rate can result from a change in photon wavelength or from a change in the number of incident photons. This confound means that individual photoreceptors are effectively color blind. Normal observers are able to see color by comparing the outputs of the three, individually color-blind, cone types.

TABLE 2 Glossary of Notation Used in Matrix Algebra Development

		Link to Conventional Notation
λ	Wavelength	
N_λ	Number of wavelength samples	
\mathbf{b}	Spectral power distribution; basis vector	
\mathbf{B}	Linear model basis vectors	
\mathbf{a}	Linear model weights	
N_b	Linear model dimension	
\mathbf{p}	Primary spectral power distribution	
\mathbf{P}	Linear model for primaries	
\mathbf{t}	Tristimulus coordinates	$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix}$
\mathbf{T}	Color-matching functions (\bar{x} , \bar{y} , and \bar{z} are rows of \mathbf{T})	$\begin{bmatrix} \bullet & \bullet & \bullet & \bar{x} & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bar{y} & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bar{z} & \bullet & \bullet & \bullet \end{bmatrix}$
\mathbf{r}	Cone (or sensor) coordinates	$\begin{bmatrix} L \\ M \\ S \end{bmatrix}$
\mathbf{R}	Cone (or sensor) sensitivities (\bar{l} , \bar{m} , and \bar{s} are rows of \mathbf{R})	$\begin{bmatrix} \bullet & \bullet & \bullet & \bar{l} & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bar{m} & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bar{s} & \bullet & \bullet & \bullet \end{bmatrix}$
ν	Luminance	$[Y]$
\mathbf{V}	Luminous efficiency function (V_λ is the single row of vector \mathbf{v})	$[\bullet \bullet \bullet V_\lambda \bullet \bullet \bullet]$
\mathbf{M}	Color space transformation matrix	



FIGURE 2 Additive color mixing. Simulated overlap of projected red, green, and blue lights. The additive combination of red and green is seen as yellow, red and blue as purple, green and blue as cyan, and red, green, and blue as white.

Color Matching

Trichromacy, together with other critical properties of color matching described in subsection “Critical Properties of Color Matching” in Sec. 10.3 mean that the color-matching behavior of an individual can be characterized as the intensities of three independent *primary lights* that are required to match a series of monochromatic spectral lights spanning the visible spectrum. Two experimental methods have been used to measure color matches: the maximum saturation method and Maxwell’s method. Most standard color-matching functions have been obtained using the maximum saturation method, though it is arguably inferior.

Maximum Saturation Method The maximum saturation method was used by Wright³¹ and Guild³² to obtain the matches that form the basis of the CIE 1931 color-matching functions (see subsection “CIE 1931 2° Color Matching Functions” in Sec. 10.4). In this method, the observer is presented with a half field illuminated by a monochromatic test light of variable wavelength λ as illustrated in Fig. 3a and an abutting half field illuminated by a mixture of red (R), green (G), and blue (B) primary lights.

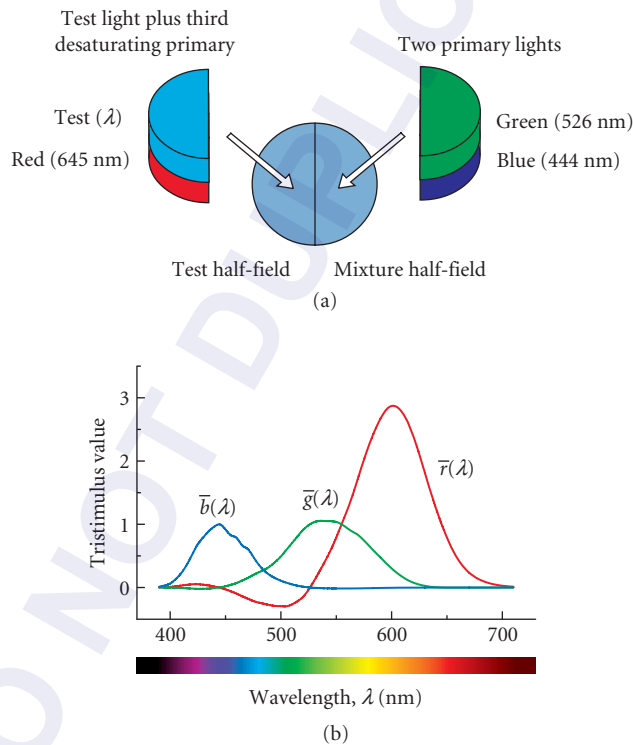


FIGURE 3 (a) Maximum saturation method of color matching. A monochromatic test field of wavelength λ can be matched using a mixture of red (645 nm), green (526 nm), and blue (444 nm) primary lights, one of which must usually be added to the test field to complete the match. (b) Color-matching functions. The amounts of each of the three primaries required to match equal energy monochromatic lights spanning the visible spectrum are known as the red $\bar{r}(\lambda)$, green $\bar{g}(\lambda)$, and blue $\bar{b}(\lambda)$, CMFs. These are shown as the red, green, and blue lines respectively. A negative sign means that primary must be added to the target to complete the match. (Based on Fig. 2.6 of Stockman and Sharpe.²¹ The data are from Stiles and Burch.³³)

(Note that in this section of the chapter, bold uppercase symbols denote primary lights, not matrices.) Often the primary lights are chosen to be monochromatic, although this is not necessary. For each test wavelength λ , the observer adjusts the intensities and arrangement of the three primary lights to make a match between the half field containing the test light and the adjacent half field. Generally, one of the primary lights is admixed with the test, while the other two are mixed together in the adjacent half field. Figure 3*b* shows the mean $\bar{r}(\lambda)$, $\bar{g}(\lambda)$, and $\bar{b}(\lambda)$ color-matching functions (hereafter abbreviated as CMFs) obtained by Stiles and Burch³³ for primary lights of 645, 526, and 444 nm. Notice that one of the CMFs is usually negative. There is no “negative light.” Negative values mean that the primary in question has been added to the test light in order to make a match. Matches using real primaries result in negative values because the primaries do not uniquely stimulate single cone photoreceptors, the spectral sensitivities of which overlap throughout the visible spectrum (see Fig. 6). Although color-matching functions are generally plotted as functions of wavelength, it is helpful to keep in mind that they represent matches, not light spectral power distributions.

The maximum saturation match between \mathbf{E}_λ , a monochromatic constituent of the equal unit energy stimulus of wavelength λ , and the three primary lights (\mathbf{R} , \mathbf{G} , and \mathbf{B}) is denoted by

$$\mathbf{E}_\lambda \sim \bar{r}(\lambda)\mathbf{R} + \bar{g}(\lambda)\mathbf{G} + \bar{b}(\lambda)\mathbf{B} \quad (1)$$

where $\bar{r}(\lambda)$, $\bar{g}(\lambda)$, and $\bar{b}(\lambda)$ are the three CMFs, and where negative CMF values indicate that the corresponding primary was mixed with the test to make the perceptual match. CMFs are usually defined for a stimulus, \mathbf{E} , which has equal unit energy throughout the spectrum. However, in practice the spectral power of the test light used in most matching experiments is varied with wavelength. In particular, longer-wavelength test lights are typically chosen to be intense enough to saturate the rods so that rods do not participate in the matches (see, e.g., Ref. 34). CMFs and the spectral power distributions of lights are always measured and tabulated as discrete functions of wavelength, typically defined in steps of 1, 5, or 10 nm.

We use the symbol \sim in Eq. (1) to indicate that two lights are a perceptual match. Perceptual matches are to be carefully distinguished from physical matches, which are denoted by the $=$ symbol. Of course, when two lights are a physical match, they must also be a perceptual match. Two lights that are a perceptual match but not a physical match are referred to as metameric color stimuli or metamers. The term metamerism is often used to refer to the fact that two physically different lights can appear identical.

The color-matching functions are defined for equal energy monochromatic test lights. More generally any test light, whether monochromatic or not, may be matched in the color-matching experiment. As noted above, we refer to the primary weights R , G , and B required to match any light as its tristimulus values. As with CMFs, tristimulus values may be negative, indicating that the corresponding primary is mixed with the test to make the match. Once the matching primaries are specified, the tristimulus values of a light provide a complete description of its effect on the human cone-mediated visual system, subject to the caveats discussed below. In addition, knowledge of the color-matching functions is sufficient to compute the tristimulus values of any light (see subsection “Tristimulus Values for Arbitrary Lights” in Sec. 10.3).

Conditions for Trichromatic Color Matching There are a number of qualifications to the empirical generalization that it is possible for observers to match any test light by adjusting the intensities of just three primaries. Some of these qualifications have to do with ancillary restrictions on the experimental conditions (e.g., the size of the bipartite field and the overall intensity of the test and matching lights). The other qualifications have to do with the choice of primaries and certain conventions about the matching procedure. First the primaries must be chosen so that it is not possible to match any one of them with a weighted superposition of the other two. Second, the observer sometimes wishes to increase the intensity of one or more of the primaries above its maximum value. In this case, we must allow him to scale the intensity of the test light down. We follow the convention of saying that the match was possible and scale up the reported primary weights by the same factor. Third, as discussed in more detail above, the observer sometimes wishes to decrease the intensity of one or more of the primaries below zero. This is always the case when the test light is a spectral light unless its wavelength is equivalent to one of the primaries. In this case, we must allow the observer

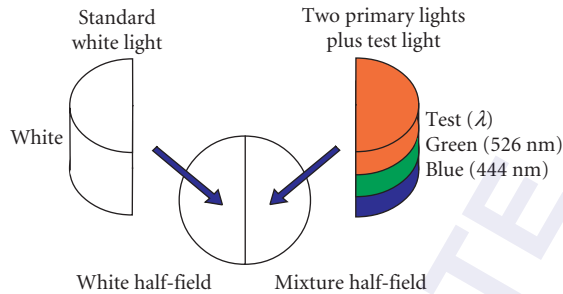


FIGURE 4 Maxwell's method of color matching. A monochromatic test field of wavelength λ replaces the primary light to which it is most similar, and a match is made to the white standard by adjusting the intensities of the two remaining primaries and the test field. (Based on Fig. 3 of Stockman.²⁰⁶)

to superimpose each such primary on the test light rather than on the other primaries. We follow the convention of saying that the match was possible but report with negative sign the intensity of each transposed primary.

With these qualifications, matching with three primaries is always possible for small fields. For larger fields, spatial inhomogeneities may make it impossible to produce a match simultaneously across the entire field (see subsections "Specificity of CMFs" and "Tristimulus Values for Arbitrary Lights" in Sec. 10.3).

Maxwell's Matching Method It is of methodological interest to note that the maximum saturation method is not the only way to instrument the color-matching experiment. Indeed the first careful quantitative measurements of color matching and trichromacy were made by Maxwell.³⁵ In Maxwell's method, which is illustrated in Fig. 4, the matched fields always appear white, so that at the match point the eye is always in the same state of adaptation whatever the test wavelength (in contrast to the maximum saturation method in which the chromaticity of the match varies with wavelength). In the experiment, the subject is first presented with a white standard half-field, and is asked to match it with the three primary lights. The test light then replaces the primary light to which it is most similar and the match is repeated. Grassmann's laws are invoked to convert the two empirical matches to the form of Eq. (1).

Critical Properties of Color Matching Color-matching data are usually obtained for monochromatic test lights. Such data are useful in general only if they can be used to predict matches for other lights with arbitrary spectral power distributions, and by extension the matches that would be made for other sets of primary lights. For this to be possible, the color-matching experiment must exhibit a number of critical properties. We review these properties briefly below. Given that they hold, it is possible to show that tristimulus values provide a complete representation for the spectral properties of light as these affect human vision. Krantz provides a detailed formal treatment.¹²

Grassmann's laws Grassmann's laws describe several of the key properties of color matching. They are:^{8,12}

1. **Symmetry:** If light **X** matches light **Y**, then **Y** matches **X**.
2. **Transitivity:** If light **X** matches light **Y** and **Y** matches light **Z**, then **X** matches **Z**.
3. **Proportionality:** If light **X** matches light **Y**, then $n\mathbf{X}$ matches $n\mathbf{Y}$ (where n is a constant of proportionality).
4. **Additivity:** If **W** matches **X** and **Y** matches **Z**, then the combination of **W** and **Y** matches the combination of **X** and **Z** (and similarly the combination of **X** and **Y** matches **W** and **Z**).

These laws have been tested extensively and hold well.^{8,19} To a first approximation, color matching can be considered to be linear and additive.^{12,36}

Uniqueness of color matches The tristimulus values of a light should be unique. This is equivalent to the requirement that only one weighted combination of the apparatus primaries produces a match to any given test light. The uniqueness of color matches ensures that tristimulus values are well-defined. In conjunction with transitivity, uniqueness also guarantees that two lights that match each other will have identical tristimulus values. It is generally accepted that, apart from variability, trichromatic color matches are unique for color normal observers.

Persistence of color matches The above properties concern color matching under a single set of viewing conditions. By viewing conditions, we refer to the properties of the image surrounding the bipartite field and the sequence of images viewed by the observer before he made the match. An important property of color matching is that lights that match under one set of viewing conditions continue to match when the viewing conditions are changed. This property is referred to as the persistence or stability of color matches.^{8,19} It holds to good approximation (but see subsection “Limits of Color Matching Data” in Sec. 10.4). The importance of the persistence law is that it allows a single set of tristimulus values to be used across viewing conditions.

Consistency across observers Finally, for the use of tristimulus values to have general validity, it is important that there should be agreement about matches across observers. For the majority of the population, there is good agreement about which lights match. We discuss individual differences in color matching in section “Limits of Color-Matching Data.”

Specificity of CMFs Color-matching data are specific to the conditions under which they were measured, and strictly to the individual observers in whom they were measured. By applying the data to other conditions and using them to predict other observer’s matches, some errors will inevitably be introduced.

An important consideration is the area of the retina within which the color matches were made. Standard color matching data (see section “Color-Matching Functions” in Sec. 10.4) have been obtained for centrally viewed fields with diameters of either 2° or 10° of visual angle. The visual angle refers to the angle subtended by an object in the external field at the effective optical center of the eye. The size of a circular matching field used in colorimetry is defined as the angular difference subtended at the eye between two diametrically opposite points on the circumference of the field. Thus, matches are defined according to the retinal size of the matching field *not* by its physical size. A 2° diameter field is known as a *small field*, whereas a 10° one as a *large field*. (One degree of visual angle is roughly equivalent to the width of the fingernail of the index finger held at arm’s length.) Color matches vary with retinal size and position because of changes in macular pigment density and photopigment optical density with visual angle (see section “Limits of Color-Matching Data”).

Standardized CMFs are mean data that are also known as *standard observer* data, in the sense that they are assumed to represent the color-matching behavior of a hypothetical typical human observer. The color matches of individual observers, however, can vary substantially from the mean matches represented by standard observer CMFs. Individual differences in lens pigment density, macular pigment density, photopigment optical density, and in the photopigments themselves can all influence color matches (see section “Limits of Color-Matching Data”).

Tristimulus Values for Arbitrary Lights Given that additivity holds for color matches, the tristimulus values, R , G , and B for an arbitrarily complex spectral radiant power distribution $P(\lambda)$ can be obtained from the $\bar{r}(\lambda)$, $\bar{g}(\lambda)$, and $\bar{b}(\lambda)$ CMFs by:

$$R = \int P(\lambda)\bar{r}(\lambda)d\lambda, \quad G = \int P(\lambda)\bar{g}(\lambda)d\lambda, \quad \text{and} \quad B = \int P(\lambda)\bar{b}(\lambda)d\lambda \quad (2)$$

Since spectral power distributions and CMFs are usually discrete functions, the integration in Eq. (2) is usually replaced by a sum.

Transformability of CMFs The $\bar{r}(\lambda)$, $\bar{g}(\lambda)$, and $\bar{b}(\lambda)$ CMFs shown in Fig. 3 are for monochromatic RGB (red-green-blue) primaries of 645, 526, and 444 nm. These CMFs can be transformed to other sets of real primary lights, and to CMFs for imaginary primary lights, such as the CIE **X**, **Y**, and **Z** primaries, or to CMFs representing the *LMS* cone spectral sensitivities (*cone fundamentals*). These transformations are illustrated in Fig. 5.

Each transformation of CMFs is accomplished by multiplying the CMFs, viewed as a column vector at each wavelength, by a 3×3 matrix. For now we simply assert this result, as our key point here is to note that such transformation is possible to enable a discussion of commonly used tristimulus representations. See Sec. 10.5 or Sec. 3.2.5 of Ref. 8 for more details about transformations between primaries.

The primaries selected by the CIE produced $\bar{x}(\lambda)$, $\bar{y}(\lambda)$, and $\bar{z}(\lambda)$ CMFs that are always positive. The $\bar{y}(\lambda)$ CMF is also the luminosity function (see section “Brightness Matching and Photometry” and also Chap. 11) thus incorporating luminosity information into the CMFs, and linking colorimetry and photometry.

The primaries that yield the cone fundamentals $\bar{l}(\lambda)$, $\bar{m}(\lambda)$, and $\bar{s}(\lambda)$ as CMFs are three imaginary primary lights that would uniquely stimulate each of the three classes of cones. Although $\bar{l}(\lambda)$, $\bar{m}(\lambda)$, and $\bar{s}(\lambda)$ cannot be obtained directly from color matches, they are strongly constrained by color-matching data since they should be a linear transformation of any other set of CMFs. Derivation of cone fundamentals is discussed in the section “Cone Fundamentals.”

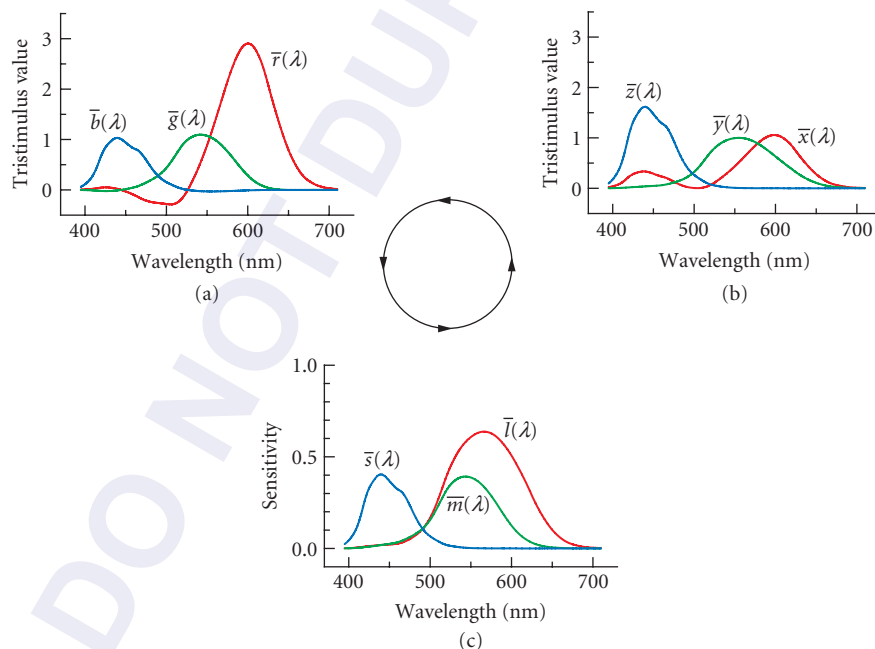


FIGURE 5 CMFs can be linearly transformed from one set of primaries to another. Illustrated here are CMFs for **R**, **G**, and **B** primaries (a), for the imaginary **X**, **Y**, and **Z** primaries (b), and the cone fundamental **L**, **M**, and **S** primaries (c). The CMFs shown in (a) and (b) are Judd-Vos modified CIE 1931 *RGB* and *XYZ* functions, respectively (see subsection “Judd-Vos Modified 2° Color-Matching Functions” in Sec. 10.4) and those shown in (c) are the Smith-Pokorny cone fundamentals (see section “Cone Fundamentals”). (Based on Fig. 4 of Stockman.²⁰⁶)

10.4 COLOR COORDINATE SYSTEMS

Overview

For the range of conditions where the color-matching experiment obeys the properties described in the previous sections, tristimulus values (or cone coordinates) provide a complete and efficient representation of human color vision. When two lights have identical tristimulus values, they are indistinguishable to the visual system and may be substituted for one another. When two lights have tristimulus values that differ substantially, they can be distinguished by an observer with normal color vision.

The relation between spectral power distributions and tristimulus values depends on the choice of primaries used in the color-matching experiment. In this sense, the choice of primaries in colorimetry is analogous to the choice of unit (e.g., foot versus meter) in the measurement of length. We use the terms *color coordinate system* and *color space* to refer to a representation derived with respect to a particular choice of primaries. We will also use the term *color coordinates* as synonym for tristimulus values.

Although the choice of primaries determines a color space, specifying primaries alone is not sufficient to compute tristimulus values. Rather, it is the color-matching functions that characterize the properties of the human observer with respect to a particular set of primaries. As noted in section “Fundamentals of Colorimetry” above and developed in detail in Sec. 10.5, “Matrix Representations and Calculations,” knowledge of the color-matching functions allows us to compute tristimulus values for arbitrary lights, as well as to derive color-matching functions with respect to other sets of primaries. Thus in practice we can specify a color space either by its primaries or by its color-matching functions.

A large number of different color spaces are in common use. The choice of which color space to use in a given application is governed by a number of considerations. If all that is of interest is to use a three-dimensional representation that accurately predicts the results of the color-matching experiment, the choice revolves around the question of finding a set of color-matching functions that accurately capture color-matching performance for the set of observers and viewing conditions under consideration. From this point of view, color spaces that differ only by an invertible linear transformation are equivalent. But there are other possible uses for color representation. For example, one might wish to choose a space that makes explicit the responses of the physiological mechanisms that mediate color vision. We discuss a number of commonly used color spaces based on CMFs, cone fundamentals, and transformations of the cone fundamentals guided by assumptions about color vision after the photoreceptors.

Many of the CMFs and cone fundamentals are available online in tabulated form at URL <http://www.cvrl.org/>.

Stimulus Spaces

A stimulus space is the color space determined by the primaries of a particular apparatus. For example, stimuli are often specified in terms of the excitation of three monitor phosphors. Stimulus color spaces have the advantage that they provide a direct description of the physical stimulus. On the other hand, they are nonstandard and their use hampers comparison of data collected in different laboratories. A useful compromise is to transform the data to a standard color space, but to provide enough side information to allow exact reconstruction of the stimulus. Often this side information can be specification of the apparatus primaries.

Color-Matching Functions

Several sets of standard CMFs are available for the central 2° or the central 10° of vision. For the central 2° (the small-field matching conditions), they are the CIE 1931 CMFs,³⁷ the Judd-Vos modified 1931 CMFs,^{38,39} and the Stiles and Burch CMFs.³³ For the central 10° (the large-field matching conditions), they are the 10° CMFs of Stiles and Burch,³⁴ and the related 10° CIE 1964 CMFs. CIE functions are available as $\bar{r}(\lambda)$, $\bar{g}(\lambda)$, and $\bar{b}(\lambda)$ for the real primaries **R**, **G**, and **B**, or as $\bar{x}(\lambda)$, $\bar{y}(\lambda)$, and $\bar{z}(\lambda)$ for the imaginary primaries **X**, **Y**, and **Z**. The latter are more commonly used in applied colorimetry.

CIE 1931 2° Color-Matching Functions In 1931, the CIE integrated a body of empirical data to determine a standard set of CMFs.^{37,40} The notion was that the CIE 1931 color-matching functions would characterize the results of a color-matching experiment performed on an “average” or “standard” color-normal human observer known as the CIE 1931 standard observer. They are available in both $\bar{r}(\lambda)$, $\bar{g}(\lambda)$, and $\bar{b}(\lambda)$ and $\bar{x}(\lambda)$, $\bar{y}(\lambda)$, and $\bar{z}(\lambda)$ form.

The empirical color-matching data used to construct the 1931 standard observer were those of Wright⁴¹ and Guild,³² which provided only the ratios of the three primaries required to match spectral test lights. Knowledge of the absolute radiances of the matching primaries is required to generate CMFs, but this was unavailable. The CIE reconstructed this information by assuming that a linear combination of the three unknown CMFs was equal to the 1924 CIE $V(\lambda)$ function.^{37,42} In addition to uncertainties about the validity of this assumption,⁴³ the $V(\lambda)$ curve that was used as the standard is now known not to provide an accurate description of typical human performance; it is far too insensitive at short wavelengths (see Fig. 2.13 of Ref. 44).

More generally, there is now considerable evidence that the color-matching functions standardized by the CIE in 1931 differ from those of the average human observer^{21,33,34,38,39} and the CIE has recently recommended¹¹ a new set of color-matching functions based on estimates of the cone photoreceptor spectral sensitivities and the Stiles and Burch 10° CMFs.³⁴ A large body of extant data is available only in terms of the CIE 1931 system, however, and many colorimetric instruments are designed around it. Therefore it seems likely that the CIE 1931 system will continue to be of practical importance for some time. Its inadequacy at short-wavelengths is well-known, and is often taken into account in colorimetric and photometric applications.

Judd-Vos Modified 2° Color-Matching Functions In 1951, Judd reconsidered the 1931 CMFs and came to the conclusion that they could be improved.³⁸ He increased the sensitivity of $V(\lambda)$ used to reconstruct the CIE CMFs below 460 nm, and derived a new set of CMFs [see Table 1 (5.5.2) of Ref. 8, which were later slightly modified by Vos,³⁹ see his Table 1].

The modifications to the $V(\lambda)$ function introduced by Judd had the unwanted effect of producing CMFs that are relatively insensitive near 460 nm (where they were unchanged). Although this insensitivity can be roughly characterized as being consistent with a high macular pigment density,^{33,45,46} the CMFs are somewhat artificial and thus removed from real color matches. Nevertheless, in practice the Judd-Vos modifications lead to a set of CMFs that are probably more typical of the average human observer than the original CIE 1931 color-matching functions. These functions were never officially standardized. However, they are widely used in practice, especially in vision science, because they are the basis of a number of estimates of the human cone spectral sensitivities, including the recent versions of the Smith-Pokorny cone fundamentals.⁴⁷

Stiles and Burch (1955) 2° CMFs The assumption used to construct the CIE 1931 standard observer, namely that $V(\lambda)$ is a linear combination of the CMFs is now unnecessary, since current instrumentation allows CMFs to be measured in conjunction with absolute radiometry. The Stiles and Burch 2° CMFs³³ are an example of directly measured functions. Though referred to by Stiles as “pilot” data, these CMFs are the most extensive set of directly measured color-matching data for 2° vision available, being averaged from matches made by 10 observers. Even compared in relative terms, there are real differences between the CIE 1931 and the Stiles and Burch³³ 2° color-matching data in the range between 430 and 490 nm. These CMFs are seldom used.

Stiles and Burch (1959) 10° CMFs The most comprehensive set of color-matching data are the large-field, centrally viewed 10° CMFs of Stiles and Burch.³⁴ Measured in 49 subjects from approximately 390 to 730 nm (and in nine subjects from 730 to 830 nm), these data are probably the most secure set of existing CMFs. Like the Stiles and Burch 2° functions,³³ the 10° functions represent directly measured CMFs, and so do not depend on measures of $V(\lambda)$. These CMFs are the basis of the Stockman and Sharpe⁴⁶ cone fundamentals (see section “Cone Fundamentals”) and thus the recent CIE proposal for a set of physiologically relevant CMFs.¹¹

1964 10° Color-Matching Functions In 1964, the CIE standardized a second set of CMFs appropriate for larger field sizes. These CMFs take into account the fact that human color matches depend on

the size of the matching fields. The CIE 1964 10° color-matching functions are an attempt to provide a standard observer for these larger fields. The use of 10° color-matching functions is recommended by the CIE when the sizes of the regions under consideration are larger than 4°.¹⁰ The large field CIE 1964 CMFs are based mainly on the 10° CMFs of Stiles and Burch³⁴ and to a lesser extent on the arguably inferior and possibly rod-contaminated 10° CMFs of Speranskaya.⁴⁸ These functions are available as $\bar{r}(\lambda)$, $\bar{g}(\lambda)$, and $\bar{b}(\lambda)$ and $\bar{x}(\lambda)$, $\bar{y}(\lambda)$, and $\bar{z}(\lambda)$.

While the CIE 1964 CMFs are similar to the 10° CMFs of Stiles and Burch functions, they differ in several ways that compromise their use as the basis for cone fundamentals.⁴⁶ The CIE¹¹ has now recommended a new set of 10° color-matching functions that are more tightly coupled to estimates of the cone spectral sensitivities and are based on the original Stiles and Burch 10° data.

Cone Fundamentals

An important goal in color science since the establishment of trichromatic color theory,^{49–52} has been the determination of the linear transformation between $\bar{r}(\lambda)$, $\bar{g}(\lambda)$, and $\bar{b}(\lambda)$ and the three cone spectral sensitivities, $\bar{l}(\lambda)$, $\bar{m}(\lambda)$, and $\bar{s}(\lambda)$.

A match between the test and mixture fields in a color-matching experiment is a match at the level of the cone photoreceptors. The response of each cone class to the mixture of primaries equals the response of that cone class to the test light. Put more formally, the following equations must hold for each unit energy test light:

$$\begin{aligned}\bar{l}_R \bar{r}(\lambda) + \bar{l}_G \bar{g}(\lambda) + \bar{l}_B \bar{b}(\lambda) &= \bar{l}(\lambda) \\ \bar{m}_R \bar{r}(\lambda) + \bar{m}_G \bar{g}(\lambda) + \bar{m}_B \bar{b}(\lambda) &= \bar{m}(\lambda) \\ \bar{s}_R \bar{r}(\lambda) + \bar{s}_G \bar{g}(\lambda) + \bar{s}_B \bar{b}(\lambda) &= \bar{s}(\lambda)\end{aligned}\quad (3)$$

where \bar{l}_R , \bar{l}_G , and \bar{l}_B are, respectively, the L-cone sensitivities to the **R**, **G**, and **B** primary lights, \bar{m}_R , \bar{m}_G , and \bar{m}_B are the M-cone sensitivities to the primary lights, and \bar{s}_R , \bar{s}_G , and \bar{s}_B are the S-cone sensitivities.

Since the S cones are now known to be insensitive to long wavelengths, it can be assumed that \bar{s}_R is effectively zero for a long-wavelength **R** primary. There are therefore eight unknowns required, and we can rewrite Eq. (3) as a linear transformation:

$$\begin{pmatrix} \bar{l}_R & \bar{l}_G & \bar{l}_B \\ \bar{m}_R & \bar{m}_G & \bar{m}_B \\ 0 & \bar{s}_G & \bar{s}_B \end{pmatrix} \begin{pmatrix} \bar{r}(\lambda) \\ \bar{g}(\lambda) \\ \bar{b}(\lambda) \end{pmatrix} = \begin{pmatrix} \bar{l}(\lambda) \\ \bar{m}(\lambda) \\ \bar{s}(\lambda) \end{pmatrix}\quad (4)$$

Moreover, since we are often more concerned about the relative $\bar{l}(\lambda)$, $\bar{m}(\lambda)$, and $\bar{s}(\lambda)$ cone spectral sensitivities, rather than their absolute values, the eight unknowns become five:

$$\begin{pmatrix} \bar{l}_R/\bar{l}_B & \bar{l}_G/\bar{l}_B & 1 \\ \bar{m}_R/\bar{m}_B & \bar{m}_G/\bar{m}_B & 1 \\ 0 & \bar{s}_G/\bar{s}_B & 1 \end{pmatrix} \begin{pmatrix} \bar{r}(\lambda) \\ \bar{g}(\lambda) \\ \bar{b}(\lambda) \end{pmatrix} = \begin{pmatrix} k_l \bar{l}(\lambda) \\ k_m \bar{m}(\lambda) \\ k_s \bar{s}(\lambda) \end{pmatrix}\quad (5)$$

Note that the constants k_l , k_m , and k_s remain unknown. Their values are typically chosen to scale the three cone fundamentals to meet some side criterion: for example, so that $k_l \bar{l}(\lambda)$, $k_m \bar{m}(\lambda)$, and $k_s \bar{s}(\lambda)$ peak at unity. Smith and Pokorny⁵³ assume that $k_l \bar{l}(\lambda) + k_m \bar{m}(\lambda)$ sum to $V(\lambda)$, the luminous efficiency function. Care should be taken when drawing conclusions that depend on the scaling chosen.

The five unknowns in the left of Eq. (5) can be estimated by fitting linear combinations of CMFs to cone spectral sensitivity measurements made in dichromatic observers and in normal observers under special conditions that isolate the responses of single cone types. They can also be estimated by comparing

color matches made by normal and dichromatic observers. Estimates from dichromats depend on the “loss,” “reduction,” or “König” assumption that dichromatic observers lack one of the three cone types, but retain two that are identical in spectral sensitivity to the normal counterparts.^{35,54} The identity of the two remaining cone types means that dichromats accept all color matches set by normal trichromats. The loss hypothesis now has a firm empirical foundation, because it has become possible to sequence and identify the photopigment opsin genes of normal, dichromatic and monochromatic observers.^{55,56} As a result, individuals who conform to the loss assumption can be selected by genetic analysis. Thanks to the longer wavelength part of the visible spectrum being effectively dichromatic, because of the insensitivity of the S cones to longer wavelength lights, the unknown value, \bar{s}_G / \bar{s}_B , can also be derived directly from normal color-matching data (see Refs. 57 and 58 for details).

Several authors have estimated LMS cone spectral sensitivities using the loss hypothesis.^{8,53,59–66} Figure 6 shows estimates by Smith and Pokorny⁵³ and Stockman and Sharpe.⁴⁶ The Smith-Pokorny estimates are a transformation of the Judd-Vos corrected CIE 1931 functions (see earlier). The Stockman-Sharpe estimates are a transformation of the Stiles and Burch 10° (see earlier) adjusted to 2° (see Ref. 21 for further information).

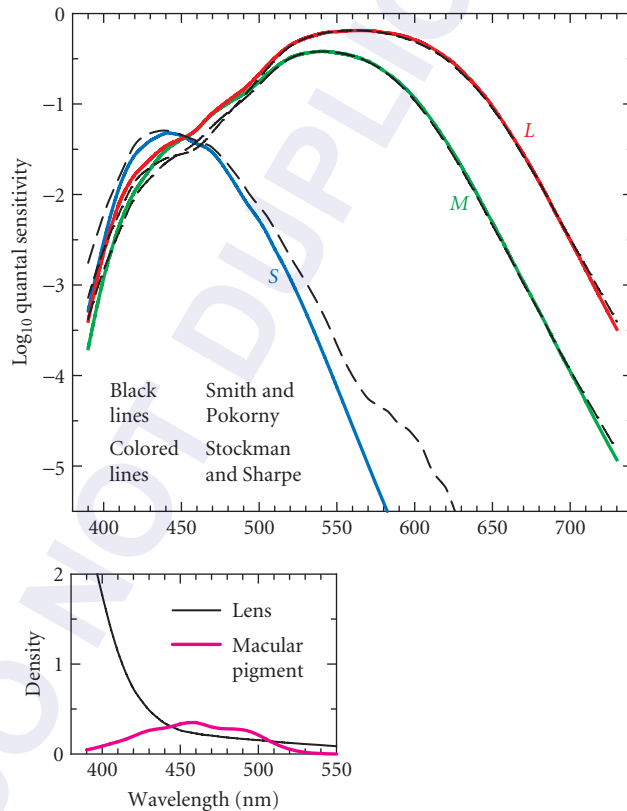


FIGURE 6 S-, M-, and L-cone spectral sensitivity estimates of Stockman and Sharpe⁴⁶ (colored lines) compared with the estimates of Smith and Pokorny⁵³ (dashed black lines). The lower inset shows the lens pigment optical density spectrum (black line) and the macular pigment optical density spectrum (magenta line) from Stockman and Sharpe.⁴⁶ Note the logarithmic vertical scale—commonly used in such plots to emphasize small sensitivities. (Based on Fig. 5 of Stockman.²⁰⁶)

Limits of Color-Matching Data

Specifying a stimulus using tristimulus values depends on having an accurate set of color-matching functions. The CMFs and cone fundamentals discussed in preceding sections are designed to be representative of a standard observer under typical viewing conditions. A number of factors limit the precision to which a standard color space can predict the individual color matches. We describe some of these factors below. Wyszecki and Stiles⁸ provide a more detailed treatment.

For most applications, standard calculations are sufficiently precise. However, when high precision is required, it is necessary to tailor a set of color-matching functions to the individual and observing conditions of interest. Once such a set of color-matching functions or cone fundamentals is available, the techniques described in other sections may be used to compute corresponding color coordinates.

Standard sets of color-matching functions are summaries or means of color-matching results for a number of color-normal observers. There is small but systematic variability between the matches set by individual observers, and this variability limits the precision to which standard color-matching functions may be taken as representative of any given color-normal observer. A number of factors underlie the variability in color matching. Stiles and Burch carefully measured color-matching functions for 49 observers using 10° fields.^{33,34} Webster and MacLeod analyzed individual variation in these color-matching functions.⁶⁷ They identified five primary factors that drive the variation in individual color matches. These are macular pigment density, lens pigment density, photopigment optical density, amount of rod intrusion into the matches, and variability in the absorption spectra of the L, M, and S cone photopigments.

Macular Pigment Density Light must pass through the ocular media before reaching the photoreceptors. At the fovea this includes the macula lutea, which contains macular pigment. This pigment absorbs lights of shorter wavelengths covering a broad spectral region centered on 460 nm (see inset of Fig. 6). There are large individual differences in macular pigment density, with peak densities at 460 nm ranging from 0.0 to about 1.2.^{68–70}

Lens Pigment Density Light is focused on the retina by the cornea and the yellow pigmented crystalline lens. The lens pigment absorbs light mainly of short wavelengths (see inset of Fig. 6). Individual differences in lens pigment density range by as much as ± 25 percent of the mean density in young observers (<30 years old).⁷¹ Lens pigment also increases with age,^{72,73} resulting in systematic differences in color-matching functions between populations of different ages.⁷⁴

Photopigment Optical Density The axial optical density of the photopigment in the photoreceptor outer segment depends on several factors, including the underlying photopigment extinction or absorbance spectra, outer segment length, and the photopigment concentration within the outer segment. All these factors can vary between individuals,^{75–82} and within individuals. Photoreceptor outer segment length, and thus axial photopigment optical density, decreases with retinal eccentricity.^{83,84} Although changes in photopigment optical density are typically neglected, they can become important under circumstances where very intense adapting fields (which dilute the photopigment by bleaching) are employed or where fixation is eccentric. See section “Adjusting Cone Spectral Sensitivities for Individual Differences” for corrections that account for changes in photopigment optical density.

Variability in Photopigment λ_{\max} Genetic and behavioral evidence shows that there are multiple versions of the human L- and M-cone photopigments.^{56,85–89} This multiplicity is known as cone polymorphism. The most common genetic polymorphism is the substitution of alanine for serine at position 180 of the L-cone photopigment gene. This substitution produces a shift in the L-cone photopigment spectral sensitivity of several nanometers, with the A180 variant shifted toward shorter wavelengths relative to the S180 variant (see Ref. 89 for a review of shift estimates). In applications where precise knowledge of an individual’s cone spectral sensitivities is important, genotyping can now help provide key information.^{46,90} Some individuals possess more than one variant of the L- or M- cone photopigment gene.^{55, 91–93}

Color-Deficient Observers A class of color-deficient individuals, known as anomalous red-green trichromats, are trichromatic but set color matches substantially different from color-normal observers. Anomalous red-green trichromacy is caused by the spectral sensitivity of either the L- or the M-cone photopigment being shifted from its normal location to an anomalous position that lies closer to the location of the spectral sensitivity function of the remaining normal M- or L-cone photopigment (for a review, see Ref. 89). These shifts result from the inheritance of hybrid LM- or ML-cone photopigment opsin genes, which are fusion genes produced by intragenic crossing over, containing the coding sequences of both L- and M-cone pigment genes. Measurements of the absorbance spectrum peaks of the hybrid pigments made *in vitro*^{87,94} and *in vivo*^{95,96} reveal a wide range of possible anomalous spectra that lie between the normal L- and M-cone spectra. The peak absorbances of the LM hybrid pigments cluster within about 8 nm of the peak absorbance of the normal M-cone pigment, while those of the ML hybrid pigments cluster within about 12 nm of the peak absorbance of the normal L-cone pigment (see Table 1 of Ref. 97). In protanomalous trichromats, one of the two polymorphic variants of the normal L-cone pigment has been replaced with a hybrid LM pigment, whereas in deuteranomalous trichromats one of the two polymorphic variants of the normal M-cone pigment has been replaced with a hybrid ML pigment.

Our development of colorimetric calculations in Sec. 10.5 can be used to tailor color specification in a particular application for color anomalous individuals, if their color-matching functions are known. Estimates of the cone sensitivities of color anomalous observers are available.^{47,98} Estimates of the A180 and S180 variants of the Stockman and Sharpe 2° functions are tabulated in Table A of Ref. 99. Details of how to adjust cone fundamentals for different λ_{\max} values are discussed in section “Adjusting Cone Spectral Sensitivities for Individual Differences” (see also section “Photopigment Optical Density Spectra” of Ref. 21).

Some individuals require only two primaries in the color-matching experiment (i.e., they are dichromats) or in rare cases only one primary (i.e., they are monochromats). Dichromats, like anomalous trichromats, are referred to as color deficient. Monochromats are, however, truly color blind (except for rod-cone interactions at mesopic levels in single cone monochromats¹⁰⁰). Most forms of monochromacy and dichromacy can be understood by assuming that the individual lacks one or more of the normal three types of cone photopigment.^{35,101} Individuals who lack the L-, M-, or S-cone photopigments are known, respectively, as protanopes, deuteranopes, or tritanopes. Protanopes and deuteranopes are much more common than tritanopes.⁸⁹ Some protanopes and deuteranopes have only one of the two normal longer wavelength cone photopigments, and so are true loss dichromats. Some, however, have a single hybrid ML- or LM-cone photopigment, which is intermediate in spectral position between M and L, while others have two cone photopigments with identical or nearly identical spectral sensitivities. For dichromatic and monochromatic individuals with normal cone photopigments (i.e., those without hybrid photopigments), the use of standard color coordinates will produce acceptable results, since a match for all three cone types will also be a match for any subset of these types. In very rare cases, an individual has no cones at all and his vision is mediated entirely by rods. His visual matches can be predicted by the CIE scotopic luminosity function [see Table I (4.3.2) of Ref. 8].

Simple standard tests exist for identifying color-deficient and color-anomalous individuals. These include the Ishihara pseudoisochromatic plates,¹⁰² the Farnsworth 100 hue test,¹⁰³ and the Rayleigh match.¹⁰⁴ For coverage of the available clinical tests see Ref. 105. Genetic analysis may also be used to identify the variants of cone pigments likely to be expressed by a given individual.⁴⁶

Retinal Inhomogeneity Most standard colorimetric systems are based on color-matching experiments where the bipartite field was either 2° or 10° in diameter and viewed foveally. The distribution of photoreceptors is not homogeneous across the retina, however, and both macular pigment and photopigment optical density decline with eccentricity. Thus, CMFs that are accurate for the fovea do not necessarily describe color matching in the extra fovea. The CIE 1964 10° XYZ color-matching functions are designed for situations where the colors being judged subtend a large visual angle. Stockman and Sharpe⁴⁶ provide both 2° and 10° cone fundamentals.

Another consideration is that the absence of S cones in approximately the central 25-min diameter of vision makes color matches confined to that small region tritanopic.^{106–109}

Rod Intrusion Both outside the fovea and at low light levels, rods can play a role in color-matching. Under conditions where rods play a role, there is a shift in the color-matching functions due to the contribution of rod signals. Wyszecki and Stiles⁸ discuss approximate methods for correcting standard sets of color-matching functions when rods intrude into color vision.

Chromatic Aberrations By some standards, even the small (roughly 2°) fields used as the basis of most color coordinate systems are rather coarse. The optics of the eye contain chromatic aberrations which cause different wavelengths of light to be focused with different accuracy. These aberrations can cause a shift in the color-matching functions if the stimuli being matched have fine spatial structure. Two stimuli which are metameric at low spatial frequencies may no longer be so at high spatial frequencies. Such effects can be quite large.¹¹⁰⁻¹¹² It is possible to correct color coordinates for chromatic aberration if enough side information is available. Such correction is rare in practice but can be important for stimuli with fine spatial structure. Some guidance is available from the literature.^{111,113} Another strategy available in the laboratory is to correct the stimulus for the chromatic aberration of the eye.¹¹⁴

Adjusting Cone Spectral Sensitivities for Individual Differences

Adjustments from Corneal to Photoreceptor Sensitivities Cone spectral sensitivities and CMFs are measured with respect to light entering the observers' cornea. However, between the cornea and photoreceptor, the light passes through the pigmented crystalline lens, and in the fovea through the macula lutea. Both of these filters markedly reduce the observers' sensitivity to short-wavelength lights (see Fig. 6).

In the first part of this section, we describe how to adjust the cone spectral sensitivities back to their values at the photoreceptor. A related adjustment is to correct cone spectral sensitivities and CMFs for individual differences in lens and macular pigment densities.

The calculation of photoreceptor sensitivities is straightforward given the lens [$d_{\text{lens}}(\lambda)$] and macular [$d_{\text{mac}}(\lambda)$] density spectra, as well as the respective scaling constants, k_{lens} and k_{mac} , by which each should be multiplied. Beginning with the *quantal* corneal spectral sensitivity of, for example, the L cones [$\bar{l}(\lambda)$], the filtering by the lens pigment [$k_{\text{lens}} d_{\text{lens}}(\lambda)$] and the macular pigment [$k_{\text{mac}} d_{\text{mac}}(\lambda)$] is removed:

$$\log_{10}[\bar{l}_r(\lambda)] = \log_{10}[\bar{l}(\lambda)] + k_{\text{lens}} d_{\text{lens}}(\lambda) + k_{\text{mac}} d_{\text{mac}}(\lambda) \quad (6)$$

to give $\bar{l}_r(\lambda)$, the spectral sensitivity of the cones at the photoreceptor. The mean or standard $d_{\text{lens}}(\lambda)$ and $d_{\text{mac}}(\lambda)$ spectra that are assumed appropriate for the Stockman and Sharpe 2° cone fundamentals are tabulated in Table 2 of their paper.⁴⁶ These densities correspond to a macular density of 0.35 at 460 nm, and a lens density of 1.765 at 400 nm. For the standard 2° observer, the values of k_{lens} and k_{mac} are set to 1, but should be adjusted appropriately for individual observers or groups of observers with different lens and macular densities. For the mean 10° observer of Stockman and Sharpe, the values of k_{lens} and k_{mac} are assumed to be 1 and 0.27, respectively.

To calculate back from photoreceptor to corneal sensitivities, the filtering by the lens and macular pigments is added back:

$$\log_{10}[\bar{l}(\lambda)] = \log_{10}[\bar{l}_r(\lambda)] - k_{\text{lens}} d_{\text{lens}}(\lambda) - k_{\text{mac}} d_{\text{mac}}(\lambda) \quad (7)$$

Again, k_{lens} and k_{mac} should be adjusted as appropriate.

Macular pigment density can be estimated psychophysically from the differences between spectral sensitivities measured centrally and peripherally (in the macular-free area). Note, however, that such estimates can be affected by other changes between the two locations, such as photopigment optical density (see Fig. 2.5 of Ref. 21). Relative estimates of lens density can be obtained psychophysically by measuring spectral sensitivities in a macular-free area of the retina, and then comparing them with

mean spectral sensitivity data. Typically, rod spectral sensitivities are measured¹¹⁵ and then compared with mean rod data, such as the data for 50 observers measured by Crawford⁷² to obtain the mean standard rod spectral sensitivity function, $V'(\lambda)$. Absolute lens density estimates can be obtained by comparing spectral sensitivities with photopigment spectra. See Ref. 21 for discussion.

Adjustments for Photopigment Optical Density As noted above, decreases and increases in photopigment optical density result in a narrowing or broadening, respectively, of the cone spectral sensitivity curves. Corrections for these changes are most easily applied to the cone fundamentals.

The photopigment optical density or absorbance spectra $[\bar{l}_{OD}(\lambda)]$ can be calculated from photoreceptor spectral sensitivity $[\bar{l}_r(\lambda)]$ given the value of D_{peak} , the peak optical density of the photopigment, thus:

$$\bar{l}_{OD}(\lambda) = \frac{-\log_{10}[1 - \bar{l}_r(\lambda)]}{D_{\text{peak}}} \quad (8)$$

Note that $\bar{l}_r(\lambda)$ should be scaled before applying Eq. (8) for $\bar{l}_{OD}(\lambda)$ to peak at 1. Stockman and Sharpe⁴⁶ assume L-, M-, and S-cone D_{peak} of 0.5, 0.5, and 0.4, respectively, for their mean 2° observer, and values of 0.38, 0.38, and 0.3 for their mean 10° observer.

The spectral sensitivity at the photoreceptor, $\bar{l}_r(\lambda)$, can be calculated from the normalized photopigment optical density spectrum, $\bar{l}_{OD}(\lambda)$, by the inversion of Eq. (8) (see Ref. 116):

$$\bar{l}_r(\lambda) = 1 - 10^{-D_{\text{peak}} \bar{l}_{OD}(\lambda)} \quad (9)$$

Calculations from corneal spectral sensitivities to retinal photopigment optical densities ignore changes in spectral sensitivity that may result from the structure of the photoreceptor or other ocular structures and pigments (unless they are incorporated in estimates of the lens or macular pigment density spectra).

Photopigment optical density can be estimated from the differences between spectral sensitivities or color matches obtained when the concentration of the photopigment is dilute and those obtained when it is in its normal concentration. This can be achieved psychophysically by comparing data obtained under bleached versus unbleached conditions or for obliquely versus axially presented lights. Bleaching measurements yield mean peak optical density values in the range 0.3 to 0.6, those that depend on oblique presentation in the range 0.7 to 1.0 and objective measures in the range 0.35 to 0.57. See Ref. 21 for discussion.

Adjustments for Changes in Photopigment λ_{max} Adjustments in the spectral position of the photopigment spectra can be affected by shifting them along an appropriate spectral scale before applying Eq. (7) to restore the prereceptor filtering, the appropriate scale being one that preserves the shapes of photopigment spectra, in general, as λ_{max} changes (i.e., their shape should be invariant). An early proposal was by Dartnall¹¹⁷ who proposed a “nomogram” or template shape for photopigment spectra that was invariant when shifted along frequency or wavenumber ($1/\lambda$, in units of cm^{-1}) scale. Shape invariance, however, is better preserved when spectra are plotted as a function of log frequency or log wavenumber $[\log(1/\lambda)]$,^{118–120} which is equivalent to log wavelength $[\log(\lambda)]$ or normalized frequency $(\lambda_{\text{max}}/\lambda)$. Barlow¹²¹ has also proposed an abscissa of the fourth root of wavelength $(\sqrt[4]{\lambda})$. A recent nomogram proposed by Govardovskii et al.¹²² is seeing considerable use. See also Eq. (8) of Ref. 46 for human photopigment nomograms. Linear wavelength scales (λ) should not be used to shift pigment templates unless the spectral shift is quite small.

Opponent and Contrast Spaces

Cone coordinates are useful because they make explicit the responses of the initial physiological mechanisms thought to mediate color vision. A number of investigators have begun to use representations that attempt to represent the responses of subsequent postreceptor mechanisms. Two basic

ideas underlie these representations. The first is the general opponent processing model described in companion chapter (Chap. 11) in this volume. We call representations based on this idea opponent color spaces. The second idea is that stimulus contrast is more relevant than stimulus magnitude.¹²³ We call spaces that are based on this second idea modulation or contrast color spaces. Some color spaces are both opponent and contrast color spaces.

Cone contrast space To derive coordinates in the cone contrast color space, the stimulus is first expressed in terms of its cone coordinates. The cone coordinates of a white point are then chosen. Usually these are the cone coordinates of a uniform adapting field or the spatio-temporal average of the cone coordinates of the entire image sequence. The cone coordinates of the white point are subtracted from the cone coordinates of the stimulus and the resulting differences are normalized by the corresponding cone coordinates of the white point.

The DKL color space Derrington, Krauskopf, and Lennie¹²⁴ introduced an opponent modulation space that is now widely used. This space is closely related to the chromaticity diagram suggested by MacLeod and Boynton¹²⁵ (see also Ref. 126). To derive coordinates in the DKL color space, the stimulus is first expressed in cone coordinates. As with cone contrast space, the cone coordinates of a white point are then subtracted from the cone coordinates of the stimulus of interest. The next step is to reexpress the resulting difference as tristimulus values with respect to a new choice of primaries that are thought to isolate the responses of post-receptoral mechanisms.^{127,128} The three primaries are chosen so that modulating two of them does not change the response of the photopic luminance mechanism (see section “Brightness Matching and Photometry”). The color coordinates corresponding to these two primaries are often called the constant B and constant R and G coordinates. Modulating the constant R and G coordinates of a stimulus modulates only the S cones. Modulating the constant B coordinate modulates both the L and M cones but keeps the S-cone response constant. Because the constant R and G coordinates are not allowed to change the response of the photopic luminance mechanism, the DKL color space is well-defined only if the S cones do not contribute to luminance. The third primary of the space is chosen so that it has the same relative cone coordinates as the white point. The coordinate corresponding to this third primary is called the luminance coordinate. Flitcroft¹¹⁰ and Brainard¹²⁹ provide detailed treatments of the DKL color space.

Caveats The basic ideas underlying the use of opponent and modulation/contrast color spaces seem to be valid. On the other hand, there is not a general agreement about how signals from cones are combined into opponent channels, about how this combination depends on adaptation, or about how adaptation affects signals originating in the cones. Since a specific model of these processes is implicit in any opponent or modulation/contrast color space, coordinates in these spaces must be treated carefully. This is particularly true of contrast spaces, where the relation between the physical stimulus and coordinates in the space depends on the choice of white point. As a consequence, radically different stimuli can have identical coordinates in a contrast space. For example, 100 percent contrast monochromatic intensity gratings are all represented by the same coordinates in contrast color spaces, independent of their wavelength. Nonetheless, such stimuli appear very different to human observers. Identity of coordinates in a contrast color space does not imply identity of appearance across different choices of white points. See Ref. 129 for more extended discussion.

Visualizing Color Data

A challenge facing today’s color scientist is to produce and interpret graphical representations of color data. Because color coordinates are three-dimensional, it is difficult to plot them on a two-dimensional page. Even more difficult is to represent a dependent measure of visual performance as a function of color coordinates. We discuss several approaches.

Three-Dimensional Approaches One strategy is to plot the three-dimensional data in perspective. In many cases the projection viewpoint may be chosen to provide a clear view of the regularities of

interest in the data. In Fig. 7a the spectrum locus is shown in the LMS tristimulus space. The three-dimensional structure of the data may be emphasized by the addition of various monocular depth cues to such figures, such as shading or drop lines. A number of computer graphics packages now provide facilities to aid in the preparation of three-dimensional perspective plots. Often these programs allow variation of the viewpoint and automatic inclusion of monocular depth cues.

Computer display technology also provides promise for improved methods of viewing three-dimensional data. For example, it is now possible to produce computer animations that show plots that vary over time. Such plots have the potential for representing multidimensional data in a manner that is more comprehensible to a human viewer than a static plot. Other interesting possibilities include the use of stereo depth cues and color displays. Online publication is making the use of such technologies more widely available for archival purposes.

Another approach to showing the three-dimensional structure of color data is to provide multiple two-dimensional views, as in a draftsman's sketch. This is illustrated in Fig. 7.

Chromaticity Diagrams A second strategy for plotting color data is to reduce the dimensionality of the data representation. One common approach is through the use of chromaticity coordinates. Chromaticity coordinates are defined so that any two lights with the same relative color coordinates have identical chromaticity coordinates. That is, the chromaticity coordinates of a light are invariant with respect to intensity scaling. Because chromaticity coordinates have one fewer degree of freedom than color coordinates, they can be described by just two numbers and plotted in a plane. We call a plot of chromaticity coordinates a chromaticity diagram. A chromaticity diagram eliminates all information about the intensity of a stimulus.

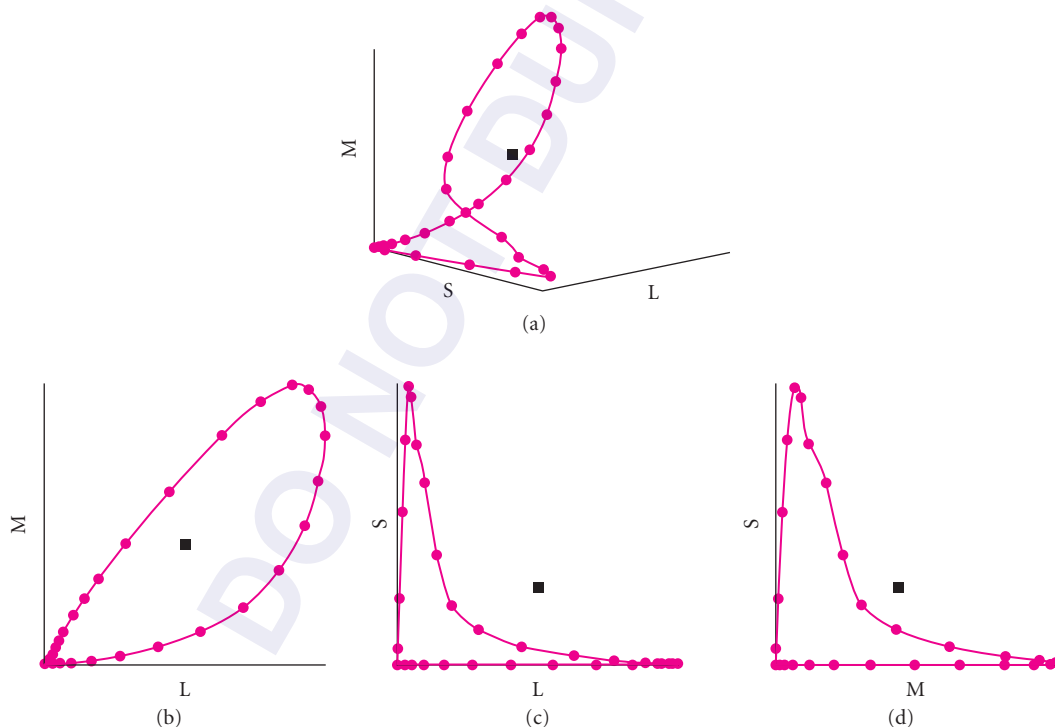


FIGURE 7 Three-dimensional views of color data. The figure shows the color coordinates of an equal energy spectrum in color space defined by the human cone sensitivities (connected closed circles) and the color coordinates of CIE daylight D65 (closed squares). (a) The data in perspective. (b, c, and d) Three two-dimensional views of the same data.

There are many ways to normalize color coordinates to produce a set of chromaticity coordinates. In general, the chromaticity coordinates $[r(\lambda), g(\lambda), \text{and } b(\lambda)]$ of the spectrum locus are related to the CMFs $[\bar{r}(\lambda), \bar{g}(\lambda), \text{and } \bar{b}(\lambda)]$ as follows:

$$\begin{aligned} r(\lambda) &= \frac{\bar{r}(\lambda)}{\bar{r}(\lambda) + \bar{g}(\lambda) + \bar{b}(\lambda)} \\ g(\lambda) &= \frac{\bar{g}(\lambda)}{\bar{r}(\lambda) + \bar{g}(\lambda) + \bar{b}(\lambda)} \quad \text{and} \\ b(\lambda) &= \frac{\bar{b}(\lambda)}{\bar{r}(\lambda) + \bar{g}(\lambda) + \bar{b}(\lambda)} \end{aligned} \quad (10)$$

Given $r(\lambda) + g(\lambda) + b(\lambda) = 1$, only $r(\lambda)$ and $g(\lambda)$ are typically plotted, since $b(\lambda)$ is $1 - [r(\lambda) + g(\lambda)]$. For the special case of the 1931 CMFs, we have:

$$\begin{aligned} x(\lambda) &= \frac{\bar{x}(\lambda)}{\bar{x}(\lambda) + \bar{y}(\lambda) + \bar{z}(\lambda)} \quad \text{and} \\ y(\lambda) &= \frac{\bar{y}(\lambda)}{\bar{x}(\lambda) + \bar{y}(\lambda) + \bar{z}(\lambda)} \end{aligned} \quad (11)$$

Figure 8 shows the spectrum locus in the 1931 CIE x, y chromaticity space with an approximate representation of the colors associated with each coordinate.

Neither r, g nor x, y chromaticity diagrams provide a strong visual connection between the data representation and the underlying cone mechanisms. For this reason, there is increasing use of chromaticity diagrams defined by the cone fundamentals. Figure 9 shows the spectrum locus plotted in l, m chromaticity coordinates.

A useful property of most chromaticity diagrams is that the chromaticity coordinates of the mixture of two lights is always a weighted combination of chromaticity coordinates of the individual lights. This is easily verified for the CIE 1931 chromaticity diagram by algebraic manipulation. Thus the chromaticity of a mixture of lights will plot somewhere on the chord connecting the chromaticities of the individual lights. Wyszecki and Stiles⁸ review a number of standard chromaticity diagrams not discussed here.

Implicit in the use of chromaticity coordinates is the assumption that scalar multiplication of the stimuli does not affect the visual performance being plotted. If the overall intensity of the stimuli matter, then the use of chromaticity coordinates can obscure important regularities. For example, the shape of color discrimination contours (see “Color Discrimination” in Sec. 10.6 and the Chap. 11) depends on how the overall intensity of the stimuli covaries with their chromaticities. Yet these contours are often plotted on a chromaticity diagram. This practice can lead to misinterpretation of the discrimination data. We recommend that plots of chromaticity coordinates be treated with some caution.

Functions of Wavelength Color data are often represented as functions of wavelength. The wavelength spectrum parameterizes a particular path through the three-dimensional color space. The exact path depends on how overall intensity covaries with wavelength. For an equal energy spectrum, the path is illustrated by Fig. 7.

Wavelength representations are particularly useful in situations where knowing the value of a function for the set of monochromatic stimuli provides a complete characterization of performance. Color-matching functions, for example, are usefully plotted as functions of wavelength because these functions may be used to predict the tristimulus values of any light. Plots of detection threshold

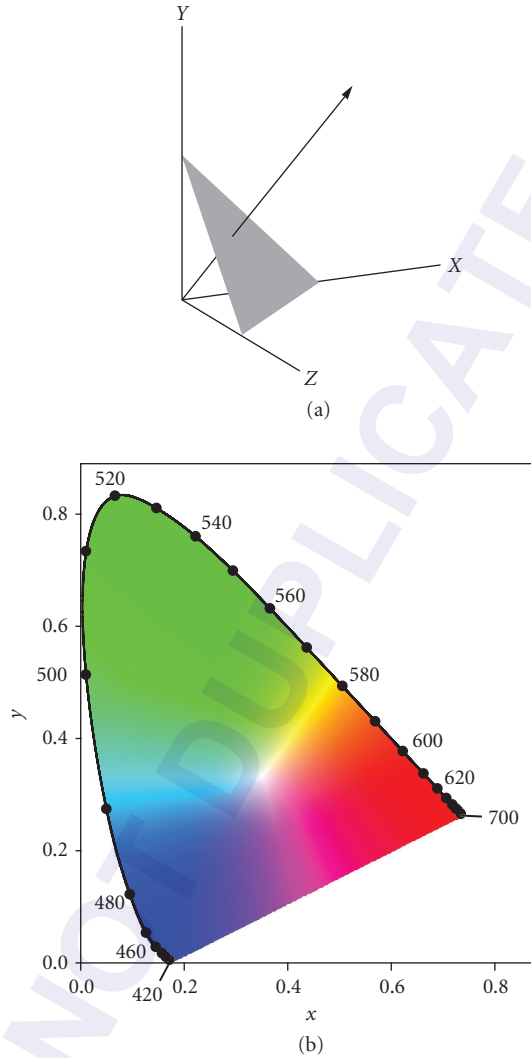


FIGURE 8 CIE 1931 xy chromaticity diagram. (a) A perspective view of the CIE 1931 XYZ tristimulus space. The ray shows a locus of points with constant chromaticity coordinates. The actual chromaticity coordinates for each ray are determined by where the ray intersects the plane described by the equation $X + Y + Z = 1$. This plane is indicated. The X and Y tristimulus values at the point of intersection are the x and y chromaticity coordinates for the ray. (b) The chromaticity coordinates of an equal energy spectrum with the interior colored to provide a rough indication of the color appearance of a stimulus of each chromaticity when viewed in a neutral context.

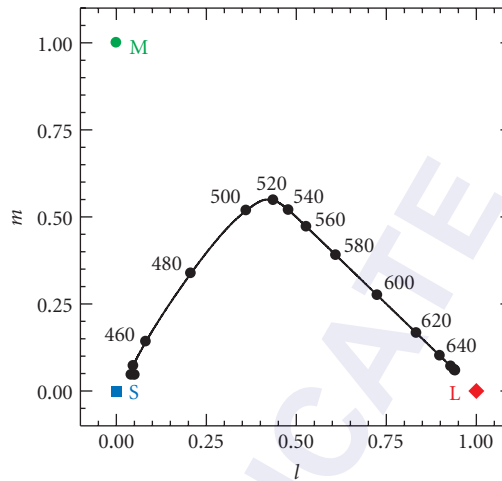


FIGURE 9 Spectrum locus (continuous line) and selected wavelengths (filled circles) plotted in the Stockman and Sharpe⁴⁶ $2^\circ l, m$ cone chromaticity space. The L- (red diamond), M- (green circle), and S- (blue square) cone fundamentals plot at (1,0), (0,1), and (0,0), respectively.

versus wavelength, on the other hand, cannot be used to predict the detection threshold for arbitrary lights.¹³⁰ Just as the chromaticity diagram tends to obscure the potential importance of manipulating the overall intensity of light, wavelength representations tend to obscure the potential importance of considering mixtures of monochromatic lights.

Colorimetric Measurements

To apply the formulas described in this chapter, it is often necessary to measure the colorimetric properties of stimuli. The most general approach is to measure the full spectral power distribution of the stimuli. Often, however, it is not necessary to know the full spectral power distribution; knowledge of the tristimulus values (in some standard color space) is sufficient. For example, the color space transformations summarized in Table 3 depend on the full spectral power distributions of the primaries only through their tristimulus values.

Specialized instruments, called colorimeters, can measure tristimulus values directly. These instruments typically operate using the same principles as photometers with the exception that they have three calibrated filters rather than just one. Each filter mimics the spectral shape of one of the color-matching functions. Wyszecki and Stiles discuss the basics of colorimeter design.⁸ Colorimeters are generally less expensive than radiometers and are thus an attractive option when full spectral data are not required.

Two caveats are worth noting. First, it is technically difficult to design filters that exactly match a desired set of color-matching functions. Generally, commercial colorimeters are calibrated so that they give accurate readings for stimuli with broadband spectral power distributions. For narrow band stimuli (e.g., the light emitted by the red phosphor of many color monitors) the reported readings may be quite inaccurate. Second, most colorimeters are designed according to the CIE 1931 standard. This may not be an optimal choice for the purpose of predicting the matches of an average human observer.

TABLE 3 Color Space Transformations

Spectral Functions Known			
Source	Destination	M	Notes
Primaries \mathbf{P}_1	CMFs \mathbf{T}_2	$\mathbf{M} = \mathbf{T}_2 \mathbf{P}_1$	T is any set of CMFs. Use regression to find M.
CMFs \mathbf{T}_1	Primaries \mathbf{P}_2	$\mathbf{M} = (\mathbf{T}_1 \mathbf{P}_2)^{-1}$	
Primaries \mathbf{P}_1	Primaries \mathbf{P}_2	$\mathbf{M} = (\mathbf{TP}_2)^{-1}(\mathbf{TP}_1)$	
CMFs \mathbf{T}_1	CMFs \mathbf{T}_2	$\mathbf{T}_2 = \mathbf{MT}_1$	
One Space Specified in Terms of Other			
Known Tristimulus Coordinates		How to Construct M	
Source primaries known in destination space.		Put them in columns of M.	
Source CMFs known in destination space.		Put them in rows of \mathbf{M}^{-1} .	
Destination primaries known in source space.		Put them in columns of \mathbf{M}^{-1} .	
Destination CMFs known in source space.		Put them in rows of M.	

CMFs stands for color matching functions.
The table summarizes how to form the matrix M that transforms color coordinates between two spaces.

10.5 MATRIX REPRESENTATIONS AND CALCULATIONS

Introduction

In the remainder of the chapter we move away from the conventional representation of colorimetric data and formulae as continuous functions of wavelength to their representation as vectors and matrices. Matrix algebra greatly simplifies the implementation of colorimetry on digital computers. Although a discrete representation provides only samples of the underlying function of wavelength, the information loss caused by this sampling can be made arbitrarily small by sampling at smaller intervals.

Notation for Matrix Calculations The conventional notation used in colorimetry does not lend itself easily to matrix and vector representations, and at the risk of some confusion between the notation used in Secs. 10.3 “Fundamentals of Colorimetry” and 10.4 “Color Coordinate Systems” and that used here and in Sec. 10.6 “Topics,” we now switch notational conventions. Table 2 provides a glossary of the major symbol usage for the matrix formulation. The following notational conventions are used: (a) scalars are denoted with italic symbols, (b) vectors are denoted with lowercase bold symbols, and (c) matrices are denoted with uppercase bold symbols. Symbols used in the appendix are generic.

Stimulus Representation

Vector Representation of Spectral Functions Suppose that spectral power density has been measured at N_λ discrete sample wavelengths $\lambda_1 \dots \lambda_{N_\lambda}$, each separated by an equal wavelength step $\Delta\lambda$. As shown in Fig. 10, we can represent the measured spectral power distribution using an N_λ dimensional column vector \mathbf{b} . The n th entry of \mathbf{b} is simply the measured power density at the n th sample wavelength multiplied by $\Delta\lambda$. Note that the values of the sample wavelengths $\lambda_1 \dots \lambda_{N_\lambda}$ and wavelength step $\Delta\lambda$ are not explicit in the vector representation. These values must be provided as side information when they are required for a particular calculation. In colorimetric applications, sample wavelengths are typically spaced evenly throughout the visible spectrum at steps of 1, 5, or 10 nm. We follow the convention that the entries of \mathbf{b} incorporate $\Delta\lambda$, however, so that we need not represent $\Delta\lambda$ explicitly when we approximate integrals over wavelength.

Manipulation of Light Intensity scaling is an operation that changes the overall power of a light at each wavelength without altering the relative power between any pair of wavelengths.

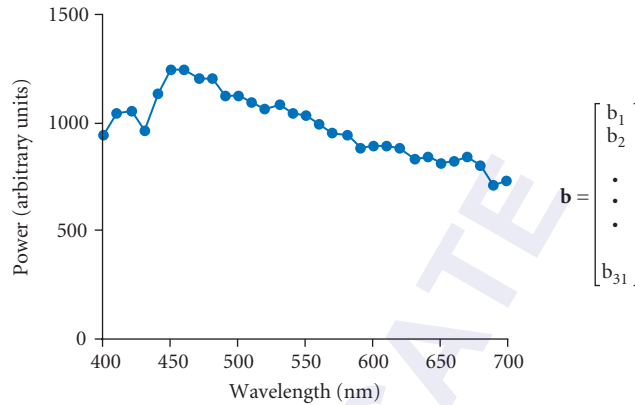


FIGURE 10 The vector representation of functions of wavelength. The plot shows a spectral power distribution measured at 10-nm intervals between 400 and 700 nm. Each point on the plot represents the power at a single sample wavelength. The vector \mathbf{b} on the right depicts the vector representation of the same spectral power distribution. The n th entry of \mathbf{b} is simply the measured power density at the n th sample wavelength times $\Delta\lambda$. Thus b_1 is derived from the power density at 400 nm, b_2 is derived from the power density at 410 nm, and b_{31} is derived from the power density at 700 nm.

The superposition of two lights is an operation that produces a new light whose power at each wavelength is the sum of the power in the original lights at the corresponding wavelength. The effects of both manipulations may be expressed using matrix algebra.

We use scalar multiplication to represent intensity scaling. If a light \mathbf{b}_1 is scaled by a factor a , then the result \mathbf{b} is given by the equation $\mathbf{b} = \mathbf{b}_1 a$. The expression $\mathbf{b}_1 a$ represents a vector whose entries are obtained by multiplying the entries of the vector \mathbf{b}_1 by the scalar a . Similarly, we use vector addition to represent superposition. If we superimpose two lights \mathbf{b}_1 and \mathbf{b}_2 , then the result \mathbf{b} is given by the equation $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$. The expression $\mathbf{b}_1 + \mathbf{b}_2$ represents a vector whose entries are obtained by adding the entries of the vectors \mathbf{b}_1 and \mathbf{b}_2 component by component. Figures 11 and 12 depict both of these operations.

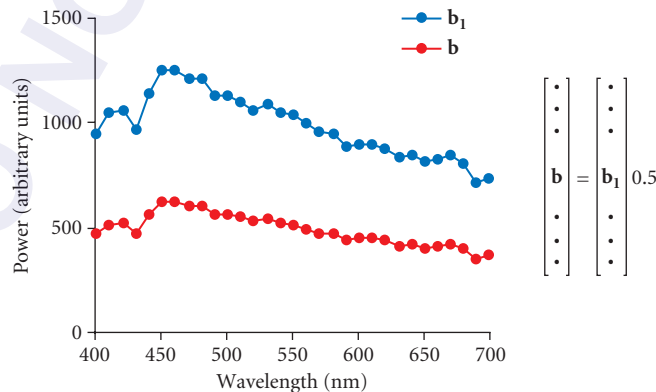


FIGURE 11 Representation of intensity scaling. Suppose that light \mathbf{b} is created by reducing the power in light \mathbf{b}_1 by a factor of 0.5 at each wavelength. The result is shown graphically in the plot. The vector representation of the same relation is given by the equation $\mathbf{b} = \mathbf{b}_1 0.5$.

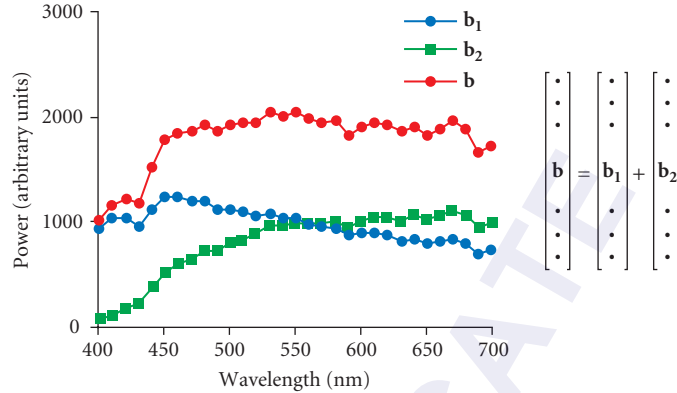


FIGURE 12 Representation of superposition. Suppose that light \mathbf{b} is created by superimposing two lights \mathbf{b}_1 and \mathbf{b}_2 . The result is shown graphically in the plot. The vector representation of the same relation is given by the equation $\mathbf{b} = \mathbf{b}_1 + \mathbf{b}_2$.

Linear Models for Spectral Functions Intensity scaling and superposition may be used in combination to produce a wide range of spectral functions. Suppose that we have a set of N_b lights that we can individually scale and superimpose. Let the vectors $\mathbf{b}_1 \dots \mathbf{b}_{N_b}$ represent the spectral power distributions of these lights. In this case, we can produce any spectral power distribution \mathbf{b} that has the form

$$\mathbf{b} = \mathbf{b}_1 a_1 + \dots + \mathbf{b}_{N_b} a_{N_b} \quad (12)$$

Now suppose we know that a spectral function \mathbf{b} is constrained to have the form of Eq. (12) where the vectors $\mathbf{b}_1 \dots \mathbf{b}_{N_b}$ are known. Then we can specify \mathbf{b} completely by providing the values of the scalars $a_1 \dots a_{N_b}$. If the number of primaries N_b is less than the number of sample wavelengths N_λ , then this specification is more efficient (i.e., requires fewer numbers) than specifying the entries of \mathbf{b} directly. We say that the spectral functions that satisfy Eq. (12) are described by (or lie within) a linear model. We call N_b the dimension of the linear model. We call the vectors $\mathbf{b}_1 \dots \mathbf{b}_{N_b}$ the basis vectors for the model. We call the scalars $a_1 \dots a_{N_b}$ required to construct any particular spectral function the model weights for that function.

Matrix Representation of Linear Models Equation (12) can be written using vector and matrix notation. Let \mathbf{B} be an N_λ by N_b dimensional matrix whose columns are the basis vectors $\mathbf{b}_1 \dots \mathbf{b}_{N_b}$. We call \mathbf{B} the basis matrix for the linear model. The composition of the basis matrix is shown pictorially on the left of Fig. 13. Let \mathbf{a} be an N_b dimensional vector whose entries are the weights $a_1 \dots a_{N_b}$. Figure 13 also depicts the vector \mathbf{a} . Using \mathbf{B} and \mathbf{a} we can re-express Eq. (12) as the matrix multiplication

$$\mathbf{b} = \mathbf{B}\mathbf{a} \quad (13)$$

The equivalence of Eqs. (12) and (13) may be established by direct expansion of the definition of matrix multiplication (see App. A). A useful working intuition for matrix multiplication is that the effect of multiplying a matrix times a vector (e.g., $\mathbf{B}\mathbf{a}$) is to produce a new vector (e.g., \mathbf{b}) that is a weighted superposition of the columns of the matrix (e.g., \mathbf{B}), where the weight for column \mathbf{b}_i is given by the i th weight, a_i , in \mathbf{a} .

Use of Linear Models When we know that a spectral function is described by a linear model, we can specify it by using the weight vector \mathbf{a} . The matrix \mathbf{B} , which is determined by the basis vectors, specifies the side information necessary to convert the vector \mathbf{a} back to the discrete wavelength

$$\mathbf{B} = \begin{bmatrix} \mathbf{b}_1 & \mathbf{b}_2 & \dots & \mathbf{b}_{N_b} \end{bmatrix} \quad \mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{N_b} \end{bmatrix} \quad \mathbf{b} = \mathbf{B} \mathbf{a}$$

FIGURE 13 Vector representation of linear models. The matrix \mathbf{B} represents the basis vectors of the linear model. The vector \mathbf{a} represents the model weights required to form a particular spectral power distribution \mathbf{b} . The relation between \mathbf{b} , \mathbf{a} , and \mathbf{B} is given by Eq. (13) and is depicted on the right of the figure.

representation \mathbf{b} . When we represent spectral functions in this way, we say that we are representing the functions within the specified linear model.

Representing spectral functions within a small-dimensional linear model places strong constraints on the form of functions. As the dimension of the model grows, linear models can represent progressively wider classes of spectral power distributions. In many cases of interest, there is prior information that allows us to assume that spectra are indeed described by a linear model. A common example of this situation is the light emitted from a computer controlled color monitor. Such a monitor produces different spectral power distributions by scaling the intensity of the light emitted by three different types of phosphor (see Chap. 22). Thus the emitted light lies within a three-dimensional linear model whose basis vectors are given by the emission spectra of the monitor's phosphors. Linear model constraints also turn out to be useful for describing naturally occurring surface and illuminant spectra.

Note that representing spectral functions within linear models is a generalization of, rather than an alternative to, the more traditional wavelength representation. To understand this, we need to only note that we can choose the basis vectors of the linear model to be discrete delta functions centered at each sample wavelength. We refer to this special choice of basis vectors as the identity basis or wavelength basis. We refer to the corresponding linear model as the identity model. For the identity model, the basis matrix \mathbf{B} is the (square) N_λ by N_λ identity matrix, where N_λ is the number of sample wavelengths. The identity matrix contains ones along its main diagonal and zeros elsewhere. Multiplying the identity matrix times any vector simply results in the same vector. From Eq. (13), we can see that when \mathbf{B} is the identity matrix, the representation of any light \mathbf{b} within the linear model is simply $\mathbf{b} = \mathbf{a}$.

Sampling the Visible Spectrum To use a discrete representation for functions of wavelength, it is necessary to choose a sampling range and sampling increment. Standard practice varies considerably. The Commission Internationale de l'Éclairage (International Commission on Illumination, commonly referred to as the CIE) provides recommendations on the sampling of the visible spectrum.¹⁰ For many applications, using 5-nm increments between 380 and 780 nm is sufficient, and coarser sampling at 10 nm between 400 and 700 nm is not uncommon. In cases where a subset of the spectral data required for a calculation is not available, interpolation or extrapolation may be used to estimate the missing values.

Vector Representation of Colorimetric Data

The Basic Color-Matching Experiment In the maximum saturation experiment color matching (see Fig. 3), observer's task is to adjust the spectral power distribution \mathbf{b}_m of the matching light on one side of the bipartite field to match the appearance of the spectral power distribution \mathbf{b}_i of the

test light on the other side of the field. As described above, the matching light's spectral power distribution is described completely by a three-dimensional linear model whose basis vectors are the primary lights' spectral power distributions. The tristimulus values of a test light are precisely the linear model weights required to form the matching light. We denote the primary spectral power distributions by the vectors $\mathbf{p}_1 \dots \mathbf{p}_3$. The associated linear model matrix \mathbf{P} contains these vectors in its three columns. We denote the tristimulus values of a light using the three-dimensional vector \mathbf{t} . Thus we can use the tristimulus values of any test light \mathbf{b} to reconstruct a matching light $\mathbf{P}\mathbf{t}$ such that:

$$\mathbf{b} \sim \mathbf{P}\mathbf{t} \quad (14)$$

We emphasize that in general that $\mathbf{P}\mathbf{t}$ will not be equal to \mathbf{b} .

Grassmann's Laws Revisited In vector notation, the proportionality law states

$$\begin{aligned} \text{if} & \quad \mathbf{b}_1 \sim \mathbf{b}_2 \\ \text{then} & \quad \mathbf{b}_1 a \sim \mathbf{b}_2 a \end{aligned} \quad (15)$$

where a is a scalar that represents any intensity scaling. The additivity law states

$$\begin{aligned} \text{if} & \quad \mathbf{b}_1 \sim \mathbf{b}_2 \quad \text{and} \quad \mathbf{b}_3 \sim \mathbf{b}_4 \\ \text{then} & \quad \mathbf{b}_1 + \mathbf{b}_3 \sim \mathbf{b}_2 + \mathbf{b}_4 \end{aligned} \quad (16)$$

The proportionality law allows us to determine the relation between the tristimulus values of a light and the tristimulus values of a scaled version of that light. Suppose that $\mathbf{b} \sim \mathbf{P}\mathbf{t}$. Applying the proportionality law, we conclude that for any scalar a , we have $\mathbf{b}a \sim (\mathbf{P}\mathbf{t})a$. Because matrix multiplication is associative, we can conclude that

$$\begin{aligned} \text{if} & \quad \mathbf{b} \sim \mathbf{P}\mathbf{t} \\ \text{then} & \quad \mathbf{b}a \sim \mathbf{P}(\mathbf{t}a) \end{aligned} \quad (17)$$

This means that the tristimulus values of a light $\mathbf{b}a$ may be obtained by scaling the tristimulus values \mathbf{t} of the light \mathbf{b} . A similar argument shows that the additivity law determines the relation between the tristimulus values of two lights and the tristimulus values of their superposition

$$\begin{aligned} \text{if} & \quad \mathbf{b}_1 \sim \mathbf{P}\mathbf{t}_1 \quad \text{and} \quad \mathbf{b}_2 \sim \mathbf{P}\mathbf{t}_2 \\ \text{then} & \quad \mathbf{b}_1 + \mathbf{b}_2 \sim \mathbf{P}(\mathbf{t}_1 + \mathbf{t}_2) \end{aligned} \quad (18)$$

Implication of Grassmann's laws If the tristimulus values of the basis vectors for a linear model are known, then Grassmann's laws allow us to determine the tristimulus values of any light within the linear model. Let $\mathbf{t}_1 \dots \mathbf{t}_{N_b}$ be the tristimulus values corresponding to the model basis vectors and let \mathbf{T}_B be the 3 by N_b matrix whose columns are $\mathbf{t}_1 \dots \mathbf{t}_{N_b}$. For any light \mathbf{b} within the linear model, we can write that $\mathbf{b} = \mathbf{B}\mathbf{a}$ where \mathbf{a} now represents the vector of weights to be applied to each basis vector to produce \mathbf{b} . By expanding this matrix product and applying Eqs. (17) and (18), it is possible to show that the tristimulus values of \mathbf{b} are given by the matrix product

$$\mathbf{t} = \mathbf{T}_B \mathbf{a} \quad (19)$$

Equation (19) is very important. It tells how to compute the tristimulus values for any light within a linear model from the tristimulus values for each of the basis vectors. Thus a small number of color matches (one for each of the basis vectors) allow us to predict color matches for a large number of lights, that is, any light within the linear model.

Color-Matching Functions and Cone Fundamentals Let \mathbf{T} be the matrix of tristimulus values for the basis vectors of the identity model (i.e., equal energy monochromatic lights). In this case \mathbf{T} has

dimensions 3 by N_λ , where N_λ is the number of sample wavelengths. Each column of \mathbf{T} is the tristimulus values for a monochromatic light. Within the identity model, the representation of any light \mathbf{b} is simply \mathbf{b} itself. From Eq. (19) we conclude directly that the tristimulus values for any light are given by

$$\mathbf{t} = \mathbf{Tb} \quad (20)$$

Once we know the tristimulus values for a set of monochromatic lights centered at each of the sample wavelengths, we can use Eq. (20) to compute the tristimulus values of any light. Equation (20) is the matrix algebra version of Eq. (2).

We can regard each of the rows of \mathbf{T} as a function of wavelength, and in doing so we can identify these as the standard color-matching functions obtained with respect to the primaries used in the matching experiment.

We can similarly represent the spectral sensitivity functions of the three classes of cones by the rows of a 3 by N_λ matrix \mathbf{R} . Let \mathbf{r} be a three-dimensional vector whose entries are the cone quantal absorption rates of an arbitrary light represented by \mathbf{b} . We can compute the absorption rates through the matrix equation

$$\mathbf{r} = \mathbf{Rb} \quad (21)$$

This computation accomplishes the wavelength-by-wavelength multiplication and summation for each cone class.

We use the term cone coordinates to refer to the vector \mathbf{r} . We can relate cone coordinates to tristimulus values in a straightforward manner. Suppose that in a color-matching experiment performed with primaries \mathbf{P} we find that a light \mathbf{b} has tristimulus values \mathbf{t} . From our mechanistic explanation of color matching in terms of the cones, we have that

$$\mathbf{r} = \mathbf{Rb} = \mathbf{RPt} \quad (22)$$

Recall that the matrix \mathbf{P} holds the power spectrum of the three primaries used in a matching experiment. If we define the matrix $\mathbf{M}_{T,R} = (\mathbf{RP})$, we see that the tristimulus values of a light are related to its cone coordinates by the linear transformation

$$\mathbf{r} = \mathbf{M}_{T,R}\mathbf{t} \quad (23)$$

By comparing Eq. (20) with Eq. (22) and noting that these equations hold for any light \mathbf{b} , we derive

$$\mathbf{R} = \mathbf{M}_{T,R}\mathbf{T} \quad (24)$$

Equation (24) has the key implication, which we took advantage of in Sec. 10.4 “Color Coordinate Systems” [Eqs. (3) and (4)], that the color-matching experiment determines the cone sensitivities up to a free linear transformation, the matrix $\mathbf{M}_{T,R}$ of Eq. (24).

Transformations between Color Spaces

Because of the large number of color spaces currently in use, the ability to transform data between various color spaces is of considerable practical importance. The derivation of such transformations depends on what is known about the source and destination color spaces. Below we discuss cases where both the source and destination color space are derived from the same underlying observer (i.e., when the source and destination color spaces both predict identical color matches). Table 3 summarizes these transformations. When the source and destination color spaces are characterized by a different underlying observer (e.g., if they are based on different CMFs) the transformation is more difficult and often cannot be done exactly. We discuss possible approaches in section “Color Cameras and Other Visual Systems.”

Source Primaries and Destination Color-Matching Functions Known Let \mathbf{P}_1 be the matrix representing the spectral power distributions of a set of known primaries, with one primary in each column. Let \mathbf{T}_2 be the matrix representing a known set of color-matching functions (e.g., the CIE 1931 XYZ color-matching functions), with one function in each of its three rows. We would like to determine a transformation between the color coordinate system specified by \mathbf{P}_1 and that specified by \mathbf{T}_2 . For example, linearized frame buffer values input to a computer-controlled color monitor may be thought of as tristimulus values in a color space defined by the monitor's phosphor emission spectra. The transformation we seek allows computation of the CIE 1931 tristimulus values from the linearized frame buffer values.

We start by using Eq. (20) to compute the tristimulus values, with respect to \mathbf{T}_2 , for the three primary lights specified by \mathbf{P}_1 . Each of these primaries is contained in a column of \mathbf{P}_1 , so that we may perform this calculation directly through the matrix multiplication

$$\mathbf{M}_{BT} = \mathbf{T}_2 \mathbf{P}_1 \quad (25)$$

Let the matrix \mathbf{P}_2 represent the destination primaries. We do not need to know these explicitly, only that they exist. The meaning of Eq. (25) is that

$$\mathbf{P}_1 \sim \mathbf{P}_2 \mathbf{M}_{BT} \quad (26)$$

where we have generalized the symbol “ \sim ” to denote a column-by-column visual match for the matrices on both sides of the relation. This relation follows because the columns of \mathbf{M}_{BT} specify how the destination primaries should be mixed to match the source primaries. Equation (26) tells us that we can substitute the three lights represented by the columns of $\mathbf{P}_2 \mathbf{M}_{BT}$ for the three lights represented by the columns of \mathbf{P}_1 in any color matching experiment. In particular, we may make this substitution for any light \mathbf{b} with tristimulus values \mathbf{t}_1 in the source color coordinate system. We have

$$\mathbf{b} \sim \mathbf{P}_1 \mathbf{t}_1 \sim \mathbf{P}_2 \mathbf{M}_{BT} \mathbf{t}_1 \quad (27)$$

By inspection, this tells us that the three-dimensional vector

$$\mathbf{t}_2 = \mathbf{M}_{BT} \mathbf{t}_1 \quad (28)$$

is the tristimulus values of \mathbf{b} in the destination color coordinate system.

Equation (28) provides us with the means to transform tristimulus values from a coordinate system where the primaries are known to one where the color-matching functions are known. The transformation matrix \mathbf{M}_{BT} required to perform the transformation depends only on the known primaries \mathbf{P}_1 and the known color-matching functions \mathbf{T}_2 . Given these, \mathbf{M}_{BT} may be computed directly from Eq. (25).

Source Color-Matching Functions and Destination Primaries Known A second transformation applies when the color-matching functions in the source color space and the primaries in the destination color space are known. This will be the case, for example, when we wish to render a stimulus specified in terms of CIE 1931 tristimulus values on a calibrated color monitor.

Let \mathbf{T}_1 represent the known color-matching functions and \mathbf{P}_2 represent the known primaries. By applying Eq. (28) we have that the relation between source tristimulus values and the destination tristimulus values is given by $\mathbf{t}_1 = \mathbf{M}_{BT} \mathbf{t}_2$. This is a system of linear equations that we may solve to find an expression for \mathbf{t}_2 in terms of \mathbf{t}_1 . In particular, as long as the matrix \mathbf{M}_{BT} is nonsingular, we can convert tristimulus values using the relation

$$\mathbf{t}_2 = \mathbf{M}_{T,P} \mathbf{t}_1 \quad (29)$$

where we define

$$\mathbf{M}_{T,P} = (\mathbf{M}_{BT})^{-1} = (\mathbf{T}_1 \mathbf{P}_2)^{-1} \quad (30)$$

Source and Destination Primaries Known A third transformation applies when the primaries of both the source and destination color spaces are known. One application of this transformation is to generate matching stimuli on two different calibrated monitors.

Let \mathbf{P}_1 and \mathbf{P}_2 represent the two sets of primaries. Let \mathbf{T} represent a set of color-matching functions for any human color coordinate system. (There is no requirement that the color-matching functions be related to either the source or the destination primaries. For example, the CIE 1931 XYZ color-matching functions might be used.) To do the conversion, we simply use Eq. (28) to transform from the color coordinate system described by \mathbf{P}_1 to the coordinate system described by \mathbf{T} . Then we use Eq. (29) to transform from the coordinates system described by \mathbf{T} to the coordinate system described by \mathbf{P}_2 . The overall transformation is given by

$$\mathbf{t}_2 = \mathbf{M}_{P_2} \mathbf{t}_1 = (\mathbf{M}_{T,P_2})(\mathbf{M}_{P_1,T}) \mathbf{t}_1 = (\mathbf{TP}_2)^{-1} (\mathbf{TP}_1) \mathbf{t}_1 \quad (31)$$

It should not be surprising that this transformation requires the specification of a set of color-matching functions. These color-matching functions are the only source of information about the human observer in the transformation equation.

Source and Destination Color-Matching Functions Known Finally, it is sometimes of interest to transform between two color spaces that are specified in terms of their color-matching functions. An example is transforming between the space defined by the Stiles and Burch 10° color-matching functions³⁴ and the space defined by the Stockman and Sharpe 10° cone fundamentals.⁴⁶

Let \mathbf{T}_1 and \mathbf{T}_2 represent the source and destination color-matching functions. Our development above assures us that there is some three-by-three transformation matrix, call it $\mathbf{M}_{T,T}$, that transforms color coordinates between the two systems. Recall that the columns of \mathbf{T}_1 and \mathbf{T}_2 are themselves tristimulus values for corresponding monochromatic lights. Thus $\mathbf{M}_{T,T}$ must satisfy

$$\mathbf{T}_2 = \mathbf{M}_{T,T} \mathbf{T}_1 \quad (32)$$

This is a system of linear equations where the entries of $\mathbf{M}_{T,T}$ are the unknown variables. This system may be solved using standard regression methods. Once we have solved for $\mathbf{M}_{T,T}$, we can transform tristimulus values using the equation

$$\mathbf{t}_2 = \mathbf{M}_{T,T} \mathbf{t}_1 \quad (33)$$

The transformation specified by Eq. (33) will be exact as long as the two sets of color-matching functions \mathbf{T}_1 and \mathbf{T}_2 characterize the performance of the same observer. One sometimes wishes, however, to transform between two color spaces that are defined with respect to different observers. For example, one might want to convert CIE 1931 XYZ tristimulus values to Judd-Vos modified tristimulus values. Although the regression procedure described here will still produce a transformation matrix in this case, the result of the transformation is not guaranteed to be correct.⁴ We will return to this topic in section “Color Cameras and Other Visual Systems.”

Interpreting the Transformation Matrix It is useful to interpret the rows and columns of the matrices derived in the preceding sections. Let \mathbf{M} be a matrix that maps the color coordinates from a source color space to a destination color space. Both source and destination color spaces are associated with a set of primaries and a set of color-matching functions. From our derivations above, we can conclude that the columns of \mathbf{M} are the coordinates of the source primaries in the destination color space [see Eq. (25)] and the rows of \mathbf{M} provide the destination color-matching functions with respect to the linear model whose basis functions are the primaries of source color space (see subsection “Use of Linear Modes” in Sec. 10.5). Similarly, the columns of \mathbf{M}^{-1} are the coordinates of the destination primaries in the source color-matching space and the rows of \mathbf{M}^{-1} are the source color-matching functions with respect to the linear model whose basis functions are the primaries of the destination color space. Thus in many cases it is possible to construct the matrix \mathbf{M} without full knowledge of the spectral functions. This can be of practical importance. For example, monitor manufacturers often specify the CIE 1931 tristimulus values of their monitors’ phosphors. In addition, colorimeters that measure tristimulus values directly are often more readily available than spectral radiometers.

Transforming Primaries and Color-Matching Functions We have shown that color coordinates in any two color spaces may be related by applying a linear transformation \mathbf{M} . The converse is also true. If we pick any nonsingular linear transformation \mathbf{M} and apply it to a set of color coordinates we have defined a new color space that will successfully predict color matches. The color-matching functions for this new space will be given by $\mathbf{T}_2 = \mathbf{M}\mathbf{T}_1$. A set of primaries for the new space will be given by $\mathbf{P}_2 = \mathbf{P}_1\mathbf{M}^{-1}$. These derived primaries are not unique. Any set of primaries that match the constructed primaries will also work.

The fact that new color spaces can be constructed by applying linear transformations has an important implication for the study of color. If we restrict attention to what we may conclude from the color-matching experiment, we can only determine the psychological representation of color up to a free linear transformation. There are two attitudes one can take toward this fact. The conservative attitude is to refrain from making any statements about the nature of color vision that depend on a particular choice of color space. The other is to appeal to experiments other than the color-matching experiment to choose a privileged representation.

10.6 TOPICS

Surfaces and Illuminants

As shown in Fig. 1, the light reaching the eye is often formed when light from an illuminant reflects from a surface. Illuminants and surfaces are of interest in color reproduction applications involving inks, paints, and dyes and in lighting design applications.

Reflection Model Illuminants are specified by their spectral power distributions. We will use the vector \mathbf{e} to represent illuminant spectral power distributions. In general, the interaction of light with matter is quite complex. For many applications, however, a rather simple model is acceptable. Using this model, we describe a surface by its surface reflectance function. The surface reflectance function specifies, for each sample wavelength, the fraction of illuminant power that is reflected to the observer. We will use the vector \mathbf{s} to represent surface reflectance spectra. Each entry of \mathbf{s} gives the reflectance measured at a single sample wavelength. Thus the spectral power distribution \mathbf{b} of the reflected light is given by the wavelength-by-wavelength product of the illuminant spectral power distribution and the surface reflectance function.

The most important consideration neglected in this formulation is viewing geometry. The relation between the radiant power emitted by a source of illumination, the material properties of a surface, and the radiant power reaching an observer can depend strongly on the viewing geometry. In our formulation, these geometrical factors must be incorporated implicitly into the specification of the illuminant and surface properties, so that any actual calculation is specific to a particular viewing geometry. Moreover, the surface reflectance must be understood as being associated with a particular image location, rather than with a particular object. A topic of current research in computer graphics is to find accurate and efficient ways to specify illuminants and surfaces for rendering,^{131,132} and parallel work in human vision seeks to understand how the reflectance of spatially complex objects is perceived.^{133–136} A second complexity that we neglect is fluorescence.

Computing the Reflected Light The relation between the surface reflectance function and the reflected light spectral power distribution is linear if the illuminant spectral power distribution is held fixed. We form the N_λ by N_λ diagonal illuminant matrix \mathbf{E} whose diagonal entries are the entries of \mathbf{e} . This leads to the relation $\mathbf{b} = \mathbf{E}\mathbf{s}$. By substituting into Eq. (20), we arrive at an expression for the tristimulus values of the light reflected from a surface:

$$\mathbf{t} = (\mathbf{TE})\mathbf{s} \quad (34)$$

where \mathbf{T} holds the colour matching functions in its three rows. The matrix (\mathbf{TE}) in this equation plays exactly the same role as the color-matching functions do in Eq. (20). Any result that holds for spectral

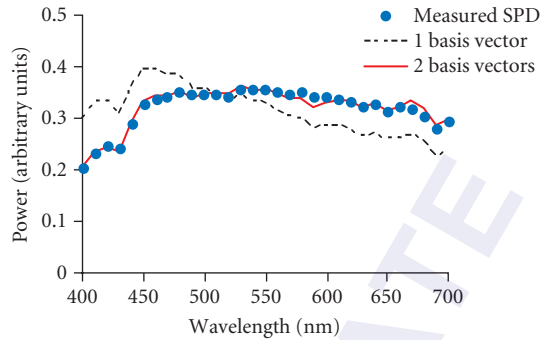


FIGURE 14 The figure shows a daylight spectral power distribution and its approximation using the CIE linear model for daylight. For this particular illuminant, only two basis functions were required to provide a very good fit.

power distributions may thus be directly extended to a result for surface reflectance functions when the illuminant is known and held fixed (see subsection “Color Coordinates of Surfaces” in Sec. 10.6).

Linear Model Representations for Surfaces and Illuminants Judd, MacAdam, and Wyszecki¹³⁷ measured the spectral power distributions of a large number of naturally occurring daylights. They determined that a four-dimensional linear model provided a good description of their spectral measurements. Consideration of their results and other daylight measurements led the CIE to standardize a three-dimensional linear model for natural daylights.¹⁰ Figure 14 depicts a daylight spectral power distribution measured by the first author and its approximation using the first two basis vectors of the CIE linear model for daylight.

Cohen¹³⁸ analyzed the reflectance spectra of a large set of Munsell papers^{139,140} and concluded that a four-dimensional linear model provided a good approximation to the entire data set. Maloney¹⁴¹ reanalyzed these data, plus a set of natural spectra measured by Krinov¹⁴² and confirmed Cohen’s conclusion. More recently, reflectance measurements of the spectra of additional Munsell papers and natural objects^{143,144} have been described by small-dimensional linear models. Figure 15 shows a measured surface reflectance spectrum (red cloth, measured by the first author) and its approximation using Cohen’s four-dimensional linear model.

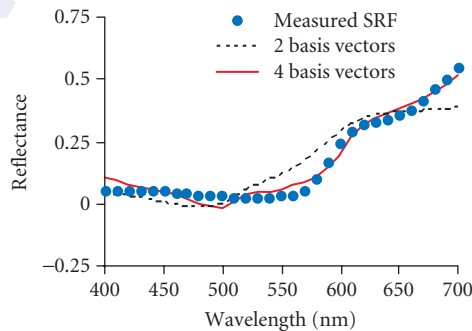


FIGURE 15 The figure shows a measured surface reflectance function and a fit to it using Cohen’s four-dimensional linear model.¹³⁸

It is not yet clear why natural illuminant and surface spectra are well-approximated by small-dimensional linear models nor how general this conclusion is. Maloney¹⁴¹ provides some speculations. None the less, the assumption that natural spectra do lie within small-dimensional linear models seems reasonable in light of the currently available evidence. This assumption makes possible a number of interesting practical calculations, as we illustrate in some of the following sections.

Determining a Linear Model from Raw Spectral Data Given a set of spectral measurements, it is possible, for any integer N_b , to find the N_b dimensional linear model that best approximates the spectral data set (in a least squares sense). Suppose that the data set consists of N_{meas} spectra, each of which is represented at N_λ sample wavelengths. Let \mathbf{X} be an N_λ by N_{meas} data matrix whose columns represent the individual spectral measurements. The goal of the calculation is to determine the N_λ by N_b matrix \mathbf{B} and an N_b by N_{meas} matrix of coefficients \mathbf{A} such that the linear model approximation $\tilde{\mathbf{X}} = \mathbf{BA}$ is the best least squares approximation to the data matrix \mathbf{X} over all possible choices of \mathbf{B} and \mathbf{A} .

The process of finding the matrix \mathbf{B} is called one mode components analysis.^{145,146} It is very closely related to the principle components analysis technique discussed in most multivariate statistics texts.¹⁴⁷ One mode 1 components analysis may be accomplished numerically through the use of the singular value decomposition.^{27,148} We define the singular value decomposition in Sec. 10.7 “Appendix—Matrix Algebra.” To see how the singular value decomposition is used to determine an N_b dimensional linear model for \mathbf{X} , consider Fig. 16. Figure 16a depicts the singular value decomposition of an N_λ by N_{meas} matrix \mathbf{X} for the case $N_{\text{meas}} > N_\lambda$, where the two matrices \mathbf{D} and \mathbf{V}^T have been collapsed. This form makes it clear that each column of \mathbf{X} is given by a linear combination of the columns of \mathbf{U} . Furthermore, for each column of \mathbf{X} , the weights needed to combine the columns of \mathbf{U} are given by the corresponding column of the matrix \mathbf{DV}^T . Suppose we choose an N_b dimensional linear model \mathbf{B} for the data in \mathbf{X} by extracting the first N_b columns of \mathbf{U} . In this case, it should be clear that we can form an approximation $\tilde{\mathbf{X}}$ to the data \mathbf{X} as shown in Fig. 16b. Because the columns of \mathbf{U} are orthogonal, the matrix \mathbf{A} consists of the first N_b rows of \mathbf{DV}^T . The accuracy of the approximation depends on how important the columns of \mathbf{U} excluded from \mathbf{B} were to the original expression for \mathbf{X} . Under certain assumptions, it can be shown that choosing \mathbf{B} as above produces a linear model that minimizes the squared error of the approximation, for any choice of N_b .¹⁴⁵ Thus computing the singular value decomposition of \mathbf{X} allows us to find a good linear model of any desired dimension for $N_b < N_\lambda$. Computing linear models from data is quite feasible on modern desktop computers.

Although the above procedure produces the linear model that provides the best least squares fit to a data set, there are a number of additional considerations that should go into choosing a linear

$$\begin{aligned} \left[\begin{array}{c} \mathbf{X} \end{array} \right] &= \left[\begin{array}{c} \mathbf{U} \end{array} \right] \left[\begin{array}{c} \mathbf{DV}^T \end{array} \right] \\ &\text{(a)} \\ \left[\begin{array}{c} \tilde{\mathbf{X}} \end{array} \right] &= \left[\begin{array}{c} \mathbf{B} \end{array} \right] \left[\begin{array}{c} \mathbf{A} \end{array} \right] \\ &\text{(b)} \end{aligned}$$

FIGURE 16 (a) The figure depicts the singular value decomposition (SVD) of an N_λ by N_{meas} matrix \mathbf{X} for the case $N_{\text{meas}} > N_\lambda$. In this view we have collapsed the two matrices \mathbf{D} and \mathbf{V}^T . To determine an N_b dimensional linear model \mathbf{B} for the data in \mathbf{X} we let \mathbf{B} consist of the first N_b columns of \mathbf{U} . (b) The linear model approximation of the data is given by $\tilde{\mathbf{X}} = \mathbf{BA}$, where \mathbf{A} consists of the first N_b rows of \mathbf{DV}^T .

model. First, we note that the choice of linear model is not unique. Any nonsingular linear combination of the columns of \mathbf{B} will produce a linear model that provides an equally good account of the data. Second, the least squares error measure gives more weight to spectra with large amplitudes. In the case of surface spectra, this means that the more reflective surfaces will tend to drive the choice of basis vectors. In the case of illuminants, the more intense illuminants will tend to drive the choice. To avoid this weighting, the measured spectra are sometimes normalized to unit length before performing the singular value decomposition. The normalization equalizes the effect of the relative shape of each spectrum in the data set.¹⁴¹ Third, it is sometimes desired to find a linear model that best describes the variation of a data set around its mean. To do this, the mean of the data set should be subtracted before performing the singular value decomposition. When the mean of the data is subtracted, one mode components analysis is identical to principle components analysis. Finally, there are circumstances where the linear model will be used not to approximate spectra but rather to approximate some other quantity (e.g., color coordinates) that depend on the spectra. In this case, more general techniques, closely related to those discussed here, may be used.¹⁴⁹

Approximating a Spectrum with Respect to a Linear Model Given an N_b dimensional linear model \mathbf{B} , it is straightforward to find the representation of any spectrum with respect to the linear model. Let \mathbf{X} be a matrix representing the spectra of functions to be approximated. These spectra do not need to be members of the data set that was used to determine the linear model. To find the matrix of coefficients \mathbf{A} such that $\tilde{\mathbf{X}} = \mathbf{BA}$ best approximates \mathbf{X} we use simple linear regression. Regression routines to solve this problem are provided as part of any standard matrix algebra software package.

Digital Image Representations If in a given application illuminants and surfaces may be represented with respect to small-dimensional linear models, then it becomes feasible to use point-by-point representations of these quantities in digital image processing. In typical color image processing, the image data at each point are represented by three numbers at each location. These numbers are generally tristimulus values in some color space. In calibrated systems, side information about the color-matching functions or primary spectral power distributions that define the color space is available to interpret the tristimulus values. It is straightforward to generalize this notion of color images by allowing the images to contain N_b numbers at each point and allowing these numbers to represent quantities other than tristimulus values.⁷ For example, in representing the image produced by a printer, it might be advantageous to represent the surface reflectance at each location.¹⁵⁰ If the gamut of printed reflectances can be represented within a small-dimensional linear model, then representing the surface reflectance functions with respect to this model would not require much more storage than a traditional color image.⁷ The basis functions for the linear model only need be represented once, not at each location. But by representing reflectances rather than tristimulus values, it becomes possible to compute what the tristimulus values reflected from the printed image would be under any illumination. We illustrate the calculation in the next section. Because of the problem of metamerism (see subsection “Computing the Reflected Light”), this calculation is not possible if only the tristimulus values are represented in the digital image. To avoid this limitation, hyperspectral images record full spectra at each image location.^{151–153}

Simulation of Illuminated Surfaces Consider the problem of producing a signal on a monitor that has the same tristimulus values as a surface under a variety of different illuminants. The solution to this problem is straightforward and is useful in a number of applications. These include rendering digitally archived paintings,^{154,155} generating stimuli for use in psychophysics,¹⁵⁶ and producing photorealistic computer-generated imagery.¹⁷ We show the calculation for the data at a single image location. Let \mathbf{a} be a representation of the surface reflectance with respect to an N_b dimensional linear model \mathbf{B} . Let \mathbf{E} represent the illuminant spectral power distribution in diagonal matrix form. Let \mathbf{T} represent the color-matching functions for a human observer, and \mathbf{P} represent the primary phosphor spectral power distributions for the monitor on which the surface will be rendered. We wish to determine tristimulus values \mathbf{t} with respect to the monitor primaries so that the light emitted from the monitor will appear identical to the light reflected from the simulated surface under the

simulated illuminant. From Eqs. (13) (cast as $\mathbf{s} = \mathbf{Ba}$), (34), (29), and (30) we can write directly the desired rendering equation

$$\mathbf{t} = [(\mathbf{TP})^{-1}(\mathbf{TE})\mathbf{B}]\mathbf{a} \quad (35)$$

The rendering matrix $[(\mathbf{TP})^{-1}(\mathbf{TE})\mathbf{B}]$ has dimensions 3 by N_b and maps the surface weights directly to monitor tristimulus values. It is quite general, in that we may use it for any calibrated monitor and any choice of linear models. It does not depend on the particular surface being rendered and may be computed once for an entire image. Because the rendering matrix is of small dimension, rendering of this sort is feasible, even for very large images. As discussed in subsection “Transformations between Color Spaces” in Sec. 10.5, it may be possible to determine the matrix $\mathbf{M}_{TP} = (\mathbf{TP})^{-1}$ directly. A similar shortcut is possible for the matrix $(\mathbf{TE})\mathbf{B}$. Each column of this matrix is the tristimulus values of one linear model basis vector under the illuminant specified by the matrix \mathbf{E} .

Color Coordinates of Surfaces Our discussion thus far has emphasized describing the color coordinates of lights. In many applications of colorimetry, it is desirable to describe the color properties of reflective objects. One efficient way to do this, as described above, is to use linear models to describe the full surface reflectance functions. Another possibility is to specify the color coordinates of the light reflected from the surface under standard illumination. This method allows the assignment of tristimulus values to surfaces in an orderly fashion.

The CIE has standardized several illuminant spectral power distributions that may be used for this purpose (see the next section). Using the procedures defined above, one can begin with the spectral power distribution of the illuminant and the surface reflectance function and from there calculate the desired color-matching coordinates.

The relative size of the tristimulus values assigned to a surface depend on its spectral reflectance function and on the illuminant chosen for specification. To factor the intensity of the illuminant out of the surface representation, the CIE specified a normalization of the color coordinates for use with 1931 XYZ tristimulus values. This normalization consists of multiplying the computed tristimulus values by the quantity $100/Y_0$, where Y_0 is the Y tristimulus value for the illuminant.

The tristimulus values of a surface provide enough information to match the surface when it is viewed under the illuminant used to compute those coordinates. It is important to bear in mind that two surfaces that have the same tristimulus values under one illuminant do not necessarily share the same tristimulus values under another illuminant. A more complete description can be generated using the linear model approach described above.

Standard Sources of Illumination The CIE has standardized a number of illuminant spectral power distributions.¹⁵⁷ These were designed to be typical of various common viewing conditions and are useful as specific choices of illumination when the illuminant cannot be measured directly. CIE Illuminant A is designed to be representative of tungsten-filament illumination. CIE Illuminant D65 is designed to be representative of average daylight. Other CIE standard daylight illuminants may be computed using the CIE principle components of daylight as basis vectors and the formulas specified by the CIE.¹⁰ Spectra representative of fluorescent lamps and other artificial sources are also available.^{8,10}

Metamerism

Recovering Spectral Power Distributions from Tristimulus Values It is not possible in general to recover a spectral power distribution from its tristimulus values. If some prior information about the spectral power distribution of the color signal is available, however, then recovery may be possible. Such recovery is of most interest in applications where direct spectral measurements are not possible and where knowing the full spectrum is important. For example, the effect of lens chromatic aberrations on cone quantal absorption rates depends on the full spectral power distribution.¹¹⁰

Suppose the spectral power distribution of interest is known to lie within a three-dimensional linear model. We may write $\mathbf{b} = \mathbf{Ba}$, where the basis matrix \mathbf{B} has dimensions N_λ by 3. Let \mathbf{t} be the

tristimulus values of the light with respect to a set of color-matching functions \mathbf{T} . We can conclude that $\mathbf{a} = (\mathbf{TB})^{-1}\mathbf{t}$, which implies

$$\mathbf{b} = \mathbf{B}(\mathbf{TB})^{-1}\mathbf{t} \quad (36)$$

When we do not have a prior constraint that the signal belongs to a three-dimensional linear model, we may still be able to place some linear model constraint, of dimension higher than three, on the spectral power distribution. For example, when we know that the signal was produced by the reflection of daylight from a natural object, it is reasonable to assume that the color signal lies within a linear model of dimension that may be as low as nine.¹⁵⁸ In this case, we can still write $\mathbf{b} = \mathbf{Ba}$, but we cannot apply Eq. (36) directly because the matrix (\mathbf{TB}) will be singular. To deal with this problem, we can choose a reduced linear model with only three dimensions. We then proceed as outlined above, but substitute the reduced model for the true model. This will lead to an estimate for the actual spectral power distribution \mathbf{b} . If the reduced linear model provides a reasonable approximation to \mathbf{b} , the estimation error may be quite small. The estimate will have the property that it is a metamer of \mathbf{b} . The techniques described above for finding linear model approximations may be used to choose an appropriate reduced model.

Finding Metamers of a Light It is often of interest to find metamers of a light. We discuss two approaches here. Wyszecki and Stiles⁸ treat the problem in considerable detail.

Using a linear model If we choose any three-dimensional linear model we can combine Eq. (36) with the fact that $\mathbf{t} = \mathbf{Tb}$ [Eq. (20)] to compute a pair of metameric spectral power distributions \mathbf{b} and $\hat{\mathbf{b}}$

$$\hat{\mathbf{b}} = \hat{\mathbf{B}}(\hat{\mathbf{T}}\hat{\mathbf{B}})^{-1}\hat{\mathbf{T}}\mathbf{b} \quad (37)$$

Each choice of $\hat{\mathbf{B}}$ will lead to a different metamer. Figure 17 shows a number of metameric spectral power distributions generated in this fashion.

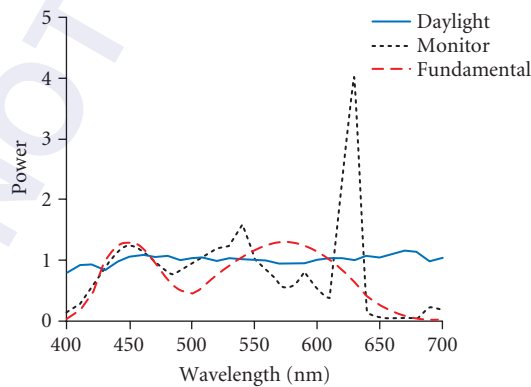


FIGURE 17 The figure shows three metameric color signals with respect to the CIE 1931 standard observer. The three metamers were computed using Eq. (37). The initial spectral power distribution \mathbf{b} (not shown) was an equal energy spectrum. Three separate linear models were used: one that describes natural daylights, one typical of monitor phosphor spectral power distributions, and one that provides Cohen's "fundamental metamer."²⁰⁷

Metameric blacks Another approach to generating metamers is to note that there will be some spectral power distributions \mathbf{b}_0 that have the property $\mathbf{T}\mathbf{b}_0 = \mathbf{0}$. Wyszecki referred to such distributions as metameric blacks, since they have the same tristimulus values as no light at all.¹⁵⁹ Grassmann's laws imply that adding a metameric black \mathbf{b}_0 to any light \mathbf{b} yields a metamer of \mathbf{b} . Given a linear model \mathbf{B} with dimension greater than three it is possible to find a second linear model \mathbf{B}_0 such that (a) all lights that lie in \mathbf{B}_0 also lie in \mathbf{B} and (b) all lights in \mathbf{B}_0 are metameric blacks. We determine \mathbf{B}_0 by finding a linear model for the null space of the matrix \mathbf{TB} . The null space of a matrix consists of all vectors that are mapped to $\mathbf{0}$ by the matrix. Finding a basis for the null space of a matrix is a standard operation in numerical matrix algebra. If we have a set of basis vectors \mathbf{N}_0 for the null space of \mathbf{TB} , we can form $\mathbf{B}_0 = \mathbf{B}\mathbf{N}_0$. This technique provides a way to generate a large list of metamers for any given light \mathbf{b} . We choose a set of weights \mathbf{a} at random and construct $\mathbf{b}_0 = \mathbf{B}_0\mathbf{a}$. We then add \mathbf{b}_0 to \mathbf{b} to form a metamer. To generate more metamers, we simply repeat with new choices of weight vector \mathbf{a} .

Surface and Illuminant Metamerism The formal similarity between Eq. (20) (which gives the relation between spectral power distributions and tristimulus values) and Eq. (34) (which gives the relation between surface reflectance functions and tristimulus values when the illuminant is known) makes it clear that our discussion of metamerism can be applied to surface reflectance spectra. Two physically different surfaces will appear identical if the tristimulus values of the light reflected from them is identical. This fact can be used to good purpose in some color reproduction applications. Suppose that we have a sample surface or textile whose color we wish to reproduce. It may be that we are not able to reproduce the sample's surface reflectance function exactly because of various limitations in the available color reproduction technology. If we know the illuminant under which the reproduction will be viewed, we may be able to determine a reproducible reflectance function that is metameric to that of the desired sample. This will give us a sample whose color appearance is as desired. Applications of this sort make heavy use of the methods described earlier to determine metamers.

But what if the illuminant is not known or if it is known to vary? In this case there is an additional richness to the topic of determining metamers. We can pose the problem of finding surface reflectance functions that will be metameric to a desired reflectance under multiple specified illuminants or under all of the illuminants within some linear model. The general methods developed here have been extended to analyze this case.^{158,160} Similar issues arise in lighting design, where we desire to produce an artificial light whose color-rendering properties match those of a specified light (such as natural daylight). When wavelength by wavelength matching of the spectra is not feasible, it may still be possible to find a spectrum so that the light reflected from surfaces within a linear model is identical for the two light sources. Because of the symmetric role of illuminants and surfaces in reflection, this problem may be treated by the same methods as used for surface reproduction.

Color Cameras and Other Visual Systems

We have treated colorimetry from the point of view of specifying the spectral information available to a human observer. We have developed our treatment, however, in such a way that it may be applied to handle other visual systems. Suppose that we wish to define color coordinates with respect to some arbitrary visual system with N_{device} photosensors. This visual system might be an artificial system based on a color camera or scanner, a nonhuman biological visual system, or the visual system of a color anomalous human observer. We assume that the sensitivities of the visual system's photosensors are known up to a linear transformation. Let $\mathbf{T}_{\text{device}}$ be an N_{device} by N_{λ} matrix whose entries are the sensitivities of each of the device's sensors at each sample wavelength. We can compute the responses of these sensors to any light \mathbf{b} . Let $\mathbf{t}_{\text{device}}$ be a vector containing the responses of each sensor type to the light. Then we have $\mathbf{t}_{\text{device}} = \mathbf{T}_{\text{device}}\mathbf{b}$. We may use $\mathbf{t}_{\text{device}}$ as the device color coordinates of \mathbf{b} .

Transformation between Color Coordinates of Different Visual Systems Suppose that we have two different visual systems and we wish to transform between the color coordinates of each.

A typical example might be trying to compute the CIE 1931 XYZ tristimulus values of a light from the responses of a color camera. Let N_s be the number of source sensors, with sensitivities specified by \mathbf{T}_s . Similarly, let N_d be the number of destination sensors with sensitivities specified by \mathbf{T}_d . For any light \mathbf{b} we know that the source device color coordinates are given by $\mathbf{t}_s = \mathbf{T}_s \mathbf{b}$ and the destination device color coordinates $\mathbf{t}_d = \mathbf{T}_d \mathbf{b}$. We would like to transform between \mathbf{t}_s and \mathbf{t}_d without direct knowledge of \mathbf{b} .

If we can find an N_d by N_s matrix \mathbf{M} such that $\mathbf{T}_d = \mathbf{M} \mathbf{T}_s$ then it is easy to show that the matrix \mathbf{M} may be used to compute the destination device color coordinates from the source device color coordinates through $\mathbf{t}_d = \mathbf{M} \mathbf{t}_s$. We have already considered this case (in a less general form) in subsection “Transformations between Color Spaces.” The extension here is that we allow the possibility that the dimensions of the two color coordinate systems differ. When a linear transformation between \mathbf{T}_s and \mathbf{T}_d exists, it can be found by standard regression methods.

Horn demonstrated that when no exact linear transformation between \mathbf{T}_s and \mathbf{T}_d exists, it is not in general possible to transform between the two sets of color coordinates.⁴ The reason for this is that there will always exist a pair of lights that have the same color coordinates for the source device but different color coordinates for the destination device. The transformation will therefore be incorrect for at least one member of this pair. When no exact linear transformation exists, it is still possible to make an approximate transformation. One approach is to use linear regression to find the best linear transformation \mathbf{M} between the two sets of color-matching functions in a least squares sense. This transformation is then applied to the source color coordinates as if it were exact.⁴ Although this is an approximation, in many cases the results will be acceptable. In the absence of prior information about the spectral power distribution of the original light \mathbf{b} , it is a sensible approach.

A second possibility is to use prior constraints on the spectral power distribution of the light to guide the transformation.²² Suppose that we know that the light is constrained to lie within an N_b dimensional linear model \mathbf{B} . Then we can find the best linear transformation \mathbf{M} between the two matrices $\mathbf{T}_s \mathbf{B}$ and $\mathbf{T}_d \mathbf{B}$. This transformation may then be used to transform the source color coordinates to the destination color coordinates. It is easy to show that the transformation will be exact if $\mathbf{T}_d \mathbf{B} = \mathbf{M} \mathbf{T}_s \mathbf{B}$. Otherwise it is a reasonable approximation that takes the linear model constraint into account.

A number of recent papers present more elaborated methods for color correction.^{161–163}

Computational Color Constancy An interesting application is the problem of estimating surface reflectance functions from color coordinates. This problem is of interest for two reasons. First, it appears that human color vision makes some attempt to perform this estimation, so that our percept of color is more closely associated with object surface properties than with the proximal properties of the light reaching the eye. Second, an artificial system that could estimate surface properties would have an important cue to aid object recognition. In the case where the illuminant is known, the problem of estimating surface reflectance properties is the same as the problem of estimating the color signal, because the illuminant spectral power distribution can simply be incorporated into the sensor sensitivities. In this case the methods outlined in the preceding section for estimating color signal spectral properties can be used.

The more interesting case is where both the illuminant and the surface reflectance are unknown. In this case, the problem is more difficult. Considerable insight has been gained by applying linear model constraints to both the surface and illuminant spectral power distributions. A number of approaches have been developed for recovering surface reflectance functions or otherwise achieving color constancy.^{158,164–172} Each approach differs (1) in the additional assumptions that are made about the properties of the image and (2) in the sophistication of the model of illuminant surface interaction and scene geometry used. A thorough review of all of these methods is beyond the scope of this chapter. It is instructive, however, to review one of the simpler methods, that of Buchsbaum.¹⁶⁵ See Ebner¹⁷³ for a discussion of many current algorithms.

Buchsbaum assumed that in any given scene, the average reflectance function of the surfaces in the scene is known. This is commonly called the “gray world” assumption. He also assumed that the illuminant was diffuse and constant across the scene and that the illuminants and surfaces in the scene are described by linear models with the same dimensionality as the number of sensors.

Let S_{avg} be the spectral power distribution of the known average surface, represented in diagonal matrix form. Then it is possible to write the relation between the space average of the sensor responses and the illuminant as

$$\mathbf{t}_{\text{avg}} = \mathbf{T} \mathbf{S}_{\text{avg}} \mathbf{B}_e \mathbf{a}_e \quad (38)$$

where \mathbf{a}_e is a vector containing the weights of the illuminant within the linear model representation \mathbf{B}_e . Because we assume that the dimension $N_e = N_p$, the matrix $\mathbf{T} \mathbf{S}_{\text{avg}} \mathbf{B}_e$ will be square and typically may be inverted. From this we recover the illuminant as $\mathbf{e} = \mathbf{B}_e (\mathbf{T} \mathbf{S}_{\text{avg}} \mathbf{B}_e)^{-1} \mathbf{t}_{\text{avg}}$. If we let \mathbf{E} represent the recovered illuminant in matrix form, then at each image location we can write

$$\mathbf{t} = \mathbf{T} \mathbf{E} \mathbf{B}_s \mathbf{a}_s \quad (39)$$

where \mathbf{a}_s is a vector containing the weights of the surface within the linear model representation \mathbf{B}_s . Proceeding exactly as we did for the illuminant, we may recover the surface reflectance from this equation.

Although Buchsbaum's method depends on rather strong assumptions about the nature of the scene, subsequent algorithms have shown that these assumptions can be relaxed.^{22,166,172} Several authors treat the relation between computational color constancy and the study of human vision.^{174–178}

Color Discrimination

Measurement of Small Color Differences Our treatment so far has not included any discussion of the precision to which observers can judge identity of color appearance. To specify tolerances for color reproduction, it would be helpful to know how different the color coordinates of two lights must be for an observer to reliably distinguish between them. A number of techniques are available for measuring human ability to discriminate between colored lights. We review these briefly here as an introduction to the topic of uniform color spaces. A more extended discussion of color discrimination and its relation to postreceptoral mechanisms is presented in the companion chapter (Chap. 11).

One experimental method, employed in seminal work by MacAdam,^{179,180} is to examine the variability in individual color matches. That is, if we have observers set matches to the same test stimulus, we will discover that they do not always set exactly the same values. Rather, there will be some trial-to-trial variability in the settings. MacAdam and others^{181,182} used the sample covariance of the individual match tristimulus values as a measure of observers' color discrimination.

A second approach is to use more direct psychophysical methods (see Chap. 3) to measure observers' thresholds for discriminating between pairs of colored lights. Examples include increment threshold measurements for monochromatic lights¹⁸³ and thresholds measured systematically in a three-dimensional color space.^{184,185}

Measurements of small color differences are often summarized with isodiscrimination contours. An isodiscrimination contour specifies the color coordinates of lights that are equally discriminable from a common standard light. Figure 18 shows an illustrative isodiscrimination contour. Isodiscrimination contours are often modeled as ellipsoids^{184,185} and the figure is drawn to the typical ellipsoidal shape. The well-known MacAdam ellipses¹⁷⁹ are an example of representing discrimination data using the chromaticity coordinates of a cross-section of a full three-dimensional isodiscrimination contour (see the legend of Fig. 18). Chapter 11 provides a more extensive discussion of possible models of discrimination contours. Under some experimental conditions, the measured contour may not be ellipsoidal.

CIE Uniform Color Spaces Figure 19 shows chromaticity plots of theoretical isodiscrimination contours. A striking feature of the plots is that the size and shape of the contours depends on the standard stimulus. For this reason, it is not possible to predict whether two lights will be discriminable solely on the basis of the Euclidean distance between their color coordinates. The heterogeneity of the isodiscrimination contours must also be taken into account.

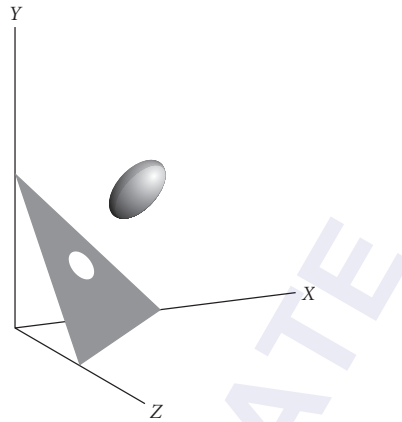


FIGURE 18 Isodiscrimination contour. The plotted ellipsoid shows a hypothetical isodiscrimination contour in the CIE XYZ color space. This contour represents color discrimination performance for the standard light whose color coordinates are located at the ellipsoid's center. Isodiscrimination contours such as the one shown are often summarized by a two-dimensional contour plotted on a chromaticity diagram (see Fig. 19). The two-dimensional contour is obtained from a cross section of the full contour, and its shape can depend on which cross section is used. This information is not available directly from the two-dimensional plot. A common criterion for choice of cross section is isoluminance. The ellipsoid shown in the figure is schematic and does not represent actual human performance.

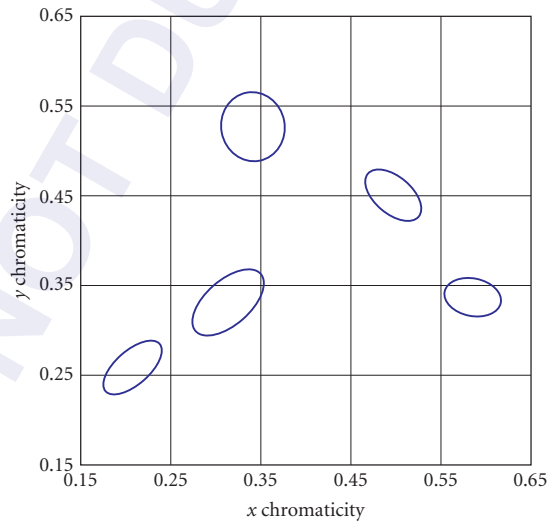


FIGURE 19 Isodiscrimination contours plotted in the chromaticity diagram. These were computed using the CIE $L^*a^*b^*$ uniform color space and ΔE_{ab}^* difference metric. They provide an approximate representation of human performance. For each standard stimulus, the plotted contour represents the color coordinates of lights that differ from the standard by $15 \Delta E_{ab}^*$ units but that have the same luminance as the standard. The choice of $15 \Delta E_{ab}^*$ units magnifies the contours compared to those that would be obtained in a threshold experiment. The contours shown are a projection of isodiscrimination contours computed for isoluminant color differences. The luminance of the white point used in the CIELAB computations was set at 1000 cd/m^2 , while the discriminations were around stimuli with a luminance of 500 cd/m^2 .

The CIE¹⁰ provides formulas that may be used to predict the discriminability of colored lights. The most recent recommendations are based on the CIE 1976 $L^*a^*b^*$ (CIELAB) color coordinates. These are obtained by a nonlinear transformation from CIE 1931 XYZ color coordinates. The transformation stretches the XYZ color space so that the resulting Euclidean distance between color coordinates provides an approximation to the how well lights may be discriminated. The $L^*a^*b^*$ system is referred to as a uniform color space. There is also a CIE 1976 $L^*u^*v^*$ (CIELUV) system, but this is now less widely used than the $L^*a^*b^*$ system and its derivatives.

Transformation to CIELAB coordinates The CIE 1976 $L^*a^*b^*$ color coordinates of a light may be obtained from its CIE XYZ coordinates according to the equations

$$L^* = \begin{cases} 116 \left(\frac{Y}{Y_n} \right)^{1/3} - 16 & \frac{Y}{Y_n} > 0.008856 \\ 903.3 \left(\frac{Y}{Y_n} \right) & \frac{Y}{Y_n} \leq 0.008856 \end{cases} \quad (40)$$

$$a^* = 500 \left[f \left(\frac{X}{X_n} \right) - f \left(\frac{Y}{Y_n} \right) \right]$$

$$b^* = 200 \left[f \left(\frac{Y}{Y_n} \right) - f \left(\frac{Z}{Z_n} \right) \right]$$

where the function $f(s)$ is defined as

$$f(s) = \begin{cases} (s)^{1/3} & s > 0.008856 \\ 7.787(s) + \frac{16}{116} & s \leq 0.008856 \end{cases} \quad (41)$$

The quantities X_n , Y_n , and Z_n in the equations above are the tristimulus values of a white point. Little guidance is available as to how to choose an appropriate white point. In the case where the lights being judged are formed when an illuminant reflects from surfaces, the tristimulus values of the illuminant may be used. In the case where the lights being judged on a computer-controlled color monitor, the sum of the tristimulus values of the three monitor phosphors stimulated at their maximum intensity may be used.

Distance in CIELAB space The Euclidean distance between the $L^*a^*b^*$ coordinates of two lights provides a rough guide to their discriminability. The symbol ΔE_{ab}^* is used to denote distance in the uniform color space and is defined as

$$\Delta E_{ab}^* = \sqrt{(\Delta L^*)^2 + (\Delta a^*)^2 + (\Delta b^*)^2} \quad (42)$$

where the various quantities on the right represent the differences between the corresponding coordinates of the two lights. Roughly speaking, a ΔE_{ab}^* value of 1 corresponds to a color difference that can just be reliably discerned by a human observer under optimal viewing conditions. A ΔE_{ab}^* value of 3 is sometimes used as an acceptable tolerance in industrial color reproduction applications.

The CIE color difference measure ΔE_{ab}^* provides only an approximate guide to the discriminability between two lights. There are a number of reasons why this is so. The first is that the relatively simple nonlinear transformation between CIE XYZ and CIE $L^*a^*b^*$ coordinates does not completely capture the empirical data on color discrimination between two samples. In part this is because the formulae were designed to predict not only discrimination data but also certain suprathreshold judgments of color appearance.¹⁸⁶ Second, color discrimination thresholds depend heavily on factors other than the tristimulus values. These factors include the adapted state of the observer,¹⁸³ the spatial and temporal structure of the stimulus,^{187–189} and the task demands placed on the observer.^{190–193} Therefore, the complete specification of a uniform color space must incorporate these factors. The CIE has now recommended a more involved method of computing small color differences from the CIE $L^*a^*b^*$ coordinates that attempts to provide better prediction of small color differences.^{10,194} The resultant computed difference is referred to as ΔE_{00}^* . The details of the computation of ΔE_{00}^* are provided and discussed in a CIE technical report.¹⁹⁴ The reader considering using ΔE_{00}^* is encouraged to study Wyszecki and Stiles's^{8(pp. 584–586)} insightful discussion of color vision models.

Effect of Errors in Color-Matching Functions

Given that there is some variation between different standard estimates of color-matching functions, between the color-matching functions of different individuals, and between the color-matching functions that mediate performance for different viewing conditions, it is of interest to determine whether the magnitude of this variation is of practical importance. There is probably no general method for making this determination, but here we outline one approach.

Consider the case of rendering a set of illuminated surfaces on a color monitor. If we know the spectral power distribution of the monitor's phosphors, it is possible to compute the appropriate weights on the monitor phosphors to produce a light metameric to each illuminated surface. The computed weights will depend on the choice of color-matching functions. Once we know the weights, however, we can find the $L^*a^*b^*$ coordinates of the emitted light. This suggests the following method to estimate the effect of differences in color-matching functions. First we compute the $L^*a^*b^*$ coordinates of surfaces rendered using the first set of color-matching functions. Then we compute the corresponding coordinates when the surfaces are rendered using the second set of color-matching functions. Finally, we compute the ΔE_{ab}^* difference between corresponding sets of coordinates. If the ΔE_{ab}^* are large, then the differences between the color-matching functions are important for the rendering application.

We have performed this calculation for a set of 462 measured surfaces^{139,140} rendered under CIE Illuminant D_{65} . The two sets of color-matching functions used were the 1931 CIE XYZ color-matching functions and the Judd-Vos modified XYZ color-matching functions. The monitor phosphor spectral power distributions were measured by the first author. The results are shown in Fig. 20. The plot shows a histogram of the differences. The median difference is 1.2 units. This difference is quite close to discrimination threshold and for many applications, the differences between the two sets of color-matching functions will probably not be of great consequence.

Brightness Matching and Photometry

The foundation of colorimetry is the human observer's ability to judge identity of color appearance. It is sometimes of interest to compare certain perceptual attributes of lights that do not, as a whole, appear identical. In particular, there has been a great deal of interest in developing formulas that predict when two lights with different relative spectral power distributions will appear equally bright. Colorimetry provides a partial answer to this question, since two lights that match in appearance must appear equally bright. Intuitively, however, it seems that it should be possible to set the relative intensities of any two lights so that they match in brightness.

In a heterochromatic brightness matching experiment, observers are asked to scale the intensity of a matching light until its brightness matches that of an experimentally controlled test light. Although

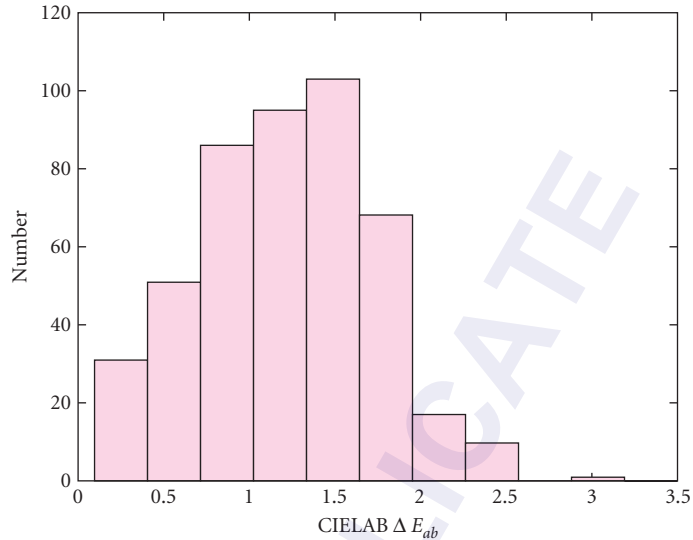


FIGURE 20 Effect of changes in color-matching functions. The plot shows a histogram of the ΔE_{ab}^* differences between two sets of lights, each of which is a monitor rendering of the same set of illuminated surfaces. The two renderings were computed using different sets of color-matching functions. The white point used in the CIELAB transformations was the XYZ coordinates of the illuminant used to compute the renderings, CIE D65.

observers can perform the heterochromatic brightness matching task, they often report that it is difficult and their matches tend to be highly variable.⁸ Moreover, the results of brightness-matching experiments are not additive.^{43,195} For photometry to be as practicable as radiometry, the measured luminous efficiency of any mixture of lights must equal the sum of the luminous efficiencies of the component lights. Such additivity is known as obedience to Abney's law.^{196,197} For this reason, more indirect methods for equating the overall effectiveness of lights at stimulating the visual system have been developed.^{8,195,198–203} The most commonly used method is that of flicker photometry. In a flicker photometric experiment, two lights of different spectral power distributions are presented alternately at the same location. At moderate flicker rates (about 20 Hz), subjects are able to adjust the overall intensity of one of the lights to minimize the apparent flicker. The intensity setting that minimizes apparent flicker is taken to indicate that the two lights match in their effectiveness as visual stimuli. Two lights equated in this way are said to be equiluminant or to have equal luminance.

Because experiments for determining when lights have the same luminance obey linearity properties similar to Grassmann's laws, it is possible to determine a luminous efficiency function that allows the assignment of a luminance value to any light. A luminous efficiency function specifies, for each sample wavelength, the relative contribution of that wavelength to the overall luminance. We can represent a luminous efficiency function as an N_λ dimensional row vector \mathbf{v} . Each entry of the matrix specifies the relative luminance of light at the corresponding sample wavelength. The luminance ν of an arbitrary spectral power distribution \mathbf{b} may be computed by the equation

$$\nu = \mathbf{vb} \quad (43)$$

The CIE has standardized four luminous efficiency functions by definition. The most commonly used of these is the standard photopic luminous efficiency function $V(\lambda)$. This is identical by definition to the 1931 XYZ color-matching function $\bar{y}(\lambda)$. For lights that subtend more than 4° of visual angle, a luminous efficiency function $V_{10}(\lambda)$ given by the 1964 10° XYZ color-matching functions is

preferred. More recently, the Judd-Vos modified color matching function has been made a supplemental standard.²⁰⁴ A final standard luminous efficiency function is available for use at low light levels when the rods are the primary functioning photoreceptors. A new luminous efficiency function will be incorporated into the new CIE proposal for a set of physiologically relevant color-matching functions. The notation V_λ or $V(\lambda)$ is often used in the literature to denote luminous efficiency functions. Note that Eq. (43) allows the computation of luminance in arbitrary units. Ref. 8 discusses standard measurement units for luminance.

It is important to note that luminance is a construct derived from flicker photometric and related experiments. As such, it does not directly predict when two lights will be judged to have the same brightness. The relation between luminance and brightness is quite complicated.^{8,205}

It is also worth noting that there is considerable individual variation in flicker photometric judgments, even among color-normal observers. For this reason, it is a common practice in psychophysical experiments to use flicker photometry to determine isoluminant stimuli for individual subjects with the stimuli of interest.

10.7 APPENDIX—MATRIX ALGEBRA

This appendix provides a brief introduction to matrix algebra. The development emphasizes the aspects of matrix algebra that are used in this chapter and is somewhat idiosyncratic. In addition, we do not prove any of the results we state. Rather, our intention is to provide the reader unfamiliar with matrix algebra with enough information to make this chapter self-contained.

Basic Notions

Vectors and Matrices A vector is a list of numbers. We use lowercase bold letters to represent vectors. We use single subscripts to identify the individual entries of a vector. The entry a_i refers to the i th number in \mathbf{a} . We call the total number of entries in a vector its dimension.

A matrix is an array of numbers. We use uppercase bold letters to represent matrices. We use dual subscripts to identify the individual entries of a matrix. The entry a_{ij} refers to the number in the i th row and j th column of \mathbf{A} . We sometimes refer to this as the ij th entry of \mathbf{A} . We call the number of rows in a matrix its row dimension. We call the number of columns in a matrix its column dimension. We generally use the symbol N to denote dimensions.

Vectors are a special case of matrices where either the row or the column dimension is 1. A matrix with a single column is often called a column vector. A matrix with a single row is often called a row vector. By convention, all vectors used in this chapter should be understood to be column vectors unless explicitly noted otherwise.

It is often convenient to think of a matrix as being composed of vectors. For example, if a matrix has dimensions N_r by N_c , then we may think of the matrix as consisting of N_c column vectors, each of which has dimension N_r .

Addition and Multiplication A vector may be multiplied by a number. We call this scalar multiplication. Scalar multiplication is accomplished by multiplying each entry of the vector by the number. If \mathbf{a} is a vector and b is a number, then $\mathbf{b} = \mathbf{ab} = \mathbf{ba}$ is a vector whose entries are given by $c_j = ba_j$.

Two vectors may be added together if they have the same dimension. We call this vector addition. Vector addition is accomplished through entry-by-entry addition. If \mathbf{a} and \mathbf{b} are vectors with the same dimension, the entries of $\mathbf{c} = \mathbf{a} + \mathbf{b}$ are given by $c_j = a_j + b_j$.

Two matrices may be added if they have the same row and column dimensions. We call this matrix addition. Matrix addition is also defined as entry by entry addition. Thus if \mathbf{A} and \mathbf{B} are matrices with the same dimension, the entries of $\mathbf{C} = \mathbf{A} + \mathbf{B}$ are given by $c_{ij} = a_{ij} + b_{ij}$. Vector addition is a special case of matrix addition.

A column vector may be multiplied by a matrix if the column dimension of the matrix matches the dimension of the vector. If \mathbf{A} has dimensions N_r by N_c and \mathbf{b} has dimension N_c , then $\mathbf{c} = \mathbf{A}\mathbf{b}$ is an N_r dimensional vector. The i th entry of \mathbf{c} is related to the entries of \mathbf{A} and \mathbf{b} by the equation:

$$c_i = \sum_{j=1}^{N_c} a_{ij} b_j \quad (44)$$

It is also possible to multiply a matrix \mathbf{B} , by another matrix, \mathbf{A} on the left, if the column dimension of \mathbf{A} matches the row dimension of \mathbf{B} . If \mathbf{A} has dimensions N_r by N and \mathbf{B} has dimensions N by N_c , then $\mathbf{C} = \mathbf{A}\mathbf{B}$ is an N_r by N_c dimensional matrix. The ik th entry of \mathbf{C} is related to the entries of \mathbf{A} and \mathbf{B} by the equation:

$$c_{ik} = \sum_{j=1}^N a_{ij} b_{jk} \quad (45)$$

By comparing Eqs. (44) and (45) we see that multiplying a matrix by a matrix is a shorthand for multiplying several vectors by the same matrix. Denote the N_c columns of \mathbf{B} by $\mathbf{b}_1, \dots, \mathbf{b}_{N_c}$ and the N_c columns of \mathbf{C} by $\mathbf{c}_1, \dots, \mathbf{c}_{N_c}$. If $\mathbf{C} = \mathbf{A}\mathbf{B}$, then $\mathbf{c}_k = \mathbf{A}\mathbf{b}_k$ for $k = 1, \dots, N_c$.

It is possible to show that matrix multiplication is associative. Suppose we have three matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} whose dimensions are such that the matrix products $(\mathbf{A}\mathbf{B})$ and $(\mathbf{B}\mathbf{C})$ are both well-defined. Then $(\mathbf{A}\mathbf{B})\mathbf{C} = \mathbf{A}(\mathbf{B}\mathbf{C})$. We often write \mathbf{ABC} to denote either product. Matrix multiplication is not commutative. Even when both products are well-defined, it is not in general true that \mathbf{BA} is equal to \mathbf{AB} .

Matrix Transposition The transpose of an N_r by N_c matrix \mathbf{A} is an N_c by N_r matrix \mathbf{B} whose ij th entry is given by $b_{ij} = a_{ji}$. We use the superscript “ T ” to denote matrix transposition: $\mathbf{B} = \mathbf{A}^T$. The identity $(\mathbf{A}\mathbf{B})^T = \mathbf{B}^T\mathbf{A}^T$ always holds.

Special Matrices and Vectors A diagonal matrix \mathbf{D} is an N_r by N_c matrix whose entries d_{ij} are zero if $i \neq j$. That is, the only nonzero entries of a diagonal matrix lie along its main diagonal. We refer to the nonzero entries of a diagonal matrix as its diagonal entries.

A square matrix is a matrix whose row and column dimensions are equal. We refer to the row and column dimensions of a square matrix as its dimension.

An identity matrix is a square diagonal matrix whose diagonal entries are all one. We use the symbol \mathbf{I}_N to denote the N by N identity matrix. Using Eq. (45) it is possible to show that for any N_r by N_c matrix \mathbf{A} , $\mathbf{A}\mathbf{I}_{N_c} = \mathbf{I}_{N_r}\mathbf{A} = \mathbf{A}$.

An orthogonal matrix \mathbf{U} is a square matrix that has the property that $\mathbf{U}^T\mathbf{U} = \mathbf{U}\mathbf{U}^T = \mathbf{I}_N$, where N is the dimension of \mathbf{U} .

A zero vector is a vector whose entries are all zero. We use the symbol $\mathbf{0}_N$ to denote the N dimensional zero vector.

Linear Models

Linear Combinations of Vectors Equation (44) is not particularly intuitive. A useful way to think about the effect of multiplying a vector \mathbf{b} by matrix \mathbf{A} is as follows. Consider the matrix \mathbf{A} to consist of N_c column vectors $\mathbf{a}_1, \dots, \mathbf{a}_{N_c}$. Then from Eq. (44) we have that the vector $\mathbf{c} = \mathbf{A}\mathbf{b}$ may be obtained by the operations of vector addition and scalar multiplication by

$$\mathbf{c} = \mathbf{a}_1 b_1 + \dots + \mathbf{a}_{N_c} b_{N_c} \quad (46)$$

where the numbers b_1, \dots, b_{N_c} are the entries of \mathbf{b} . Thus the effect of multiplying a vector by a matrix is to form a weighted sum of the columns of the matrix. The weights that go into forming the sum are the entries of the vector. We call any expression that has the form of the right hand side of Eq. (46) a linear combination of the vectors $\mathbf{a}_1, \dots, \mathbf{a}_{N_c}$.

Independence and Rank Consider a collection of vectors $\mathbf{a}_1, \dots, \mathbf{a}_{N_c}$. If no one of these vectors may be expressed as a linear combination of the others, then we say that the collection is independent. We define the rank of a collection of vectors $\mathbf{a}_1, \dots, \mathbf{a}_{N_c}$ as the largest number of linearly independent vectors that may be chosen from that collection. We define the rank of a matrix \mathbf{A} to be the rank of the vectors $\mathbf{a}_1, \dots, \mathbf{a}_{N_c}$ that make up its columns. It may be proved that the rank of a matrix is always less than or equal to the minimum of its row and column dimensions. We say that a matrix has full rank if its rank is exactly equal to the minimum of its row and column dimensions.

Linear Models When a vector \mathbf{c} has the form given in Eq. (46), we say that \mathbf{c} lies within a linear model. We call N_c the dimension of the linear model. We call the vectors $\mathbf{a}_1, \dots, \mathbf{a}_{N_c}$ the basis vectors for the model. Thus an N_c dimensional linear model with basis vectors $\mathbf{a}_1, \dots, \mathbf{a}_{N_c}$ contains all vectors \mathbf{c} that can be expressed exactly using Eq. (46) for some choice of numbers b_1, \dots, b_{N_c} . Equivalently, the linear model contains all vectors \mathbf{c} that may be expressed as $\mathbf{c} = \mathbf{A}\mathbf{b}$ where the columns of the matrix \mathbf{A} are the vectors $\mathbf{a}_1, \dots, \mathbf{a}_{N_c}$ and \mathbf{b} is an arbitrary vector. We refer to all vectors within a linear model as the subspace defined by that model.

Simultaneous Linear Equations

Matrix and Vector Form Matrix multiplication may be used to express a system of simultaneous linear equations. Suppose we have a set of N_r simultaneous linear equations in N_c unknowns. Call the unknowns b_1, \dots, b_{N_c} . Conventionally we would write the equations in the form

$$\begin{aligned} a_{11}b_1 + \dots + a_{1N_c}b_{N_c} &= c_1 \\ a_{21}b_1 + \dots + a_{2N_c}b_{N_c} &= c_2 \\ &\dots \\ a_{N_r1}b_1 + \dots + a_{N_rN_c}b_{N_c} &= c_{N_r} \end{aligned} \quad (47)$$

where the a_{ij} and c_i represent known numbers. From Eq. (44) it is easy to see that we may rewrite Eq. (47) as a matrix multiplication

$$\mathbf{A}\mathbf{b} = \mathbf{c} \quad (48)$$

In this form, the entries of the vector \mathbf{b} represent the unknowns. Solving Eq. (48) for \mathbf{b} is equivalent to solving the system of simultaneous linear equations of Eq. (47).

Solving Simultaneous Linear Equations A fundamental topic in linear algebra is finding solutions for systems of simultaneous linear equations. We will rely on several basic results in this chapter, which we state here.

When the matrix \mathbf{A} is square and has full rank, it is always possible to find a unique matrix \mathbf{A}^{-1} such that $\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}_{N_r}$. We call the matrix \mathbf{A}^{-1} the inverse of the matrix \mathbf{A} . The matrix \mathbf{A}^{-1} is also square and has full rank. Algorithms exist for determining the inverse of a matrix and are provided by software packages that support matrix algebra.

When the matrix \mathbf{A} is square and has full rank, a unique solution \mathbf{b} to Eq. (48) exists. This solution is given simply by the expression $\mathbf{b} = \mathbf{A}^{-1}\mathbf{c}$. When the rank of \mathbf{A} is less than its row dimension, then there will not in general be an exact solution to Eq. (48). There will, however, be a unique vector \mathbf{b} that is the best solution in a least squares sense. We call this the least squares solution to Eq. (48). Finding the least squares solution to Eq. (48) is often referred to as linear regression. Algorithms exist for performing linear regression and are provided by software packages that support matrix algebra.

A generalization of Eq. (48) is the matrix equation

$$\mathbf{A}\mathbf{B} = \mathbf{C} \quad (49)$$

where the entries of the matrix \mathbf{B} are the unknowns. From our interpretation of matrix multiplication as a shorthand for multiple multiplications of a vector by a matrix, we can see immediately that this type of equation may be solved by applying the above analysis in a columnwise fashion. If \mathbf{A} is square and has full rank, then we may determine \mathbf{B} uniquely as $\mathbf{A}^{-1}\mathbf{C}$. When the rank of \mathbf{A} is less than its row dimension, we may use regression determine a matrix \mathbf{B} that satisfies Eq. (49) in a least squares sense.

It is also possible to solve matrix equations of the form $\mathbf{BA} = \mathbf{C}$ where the entries of \mathbf{B} are again the unknowns. An equation of this form may be converted to the form of Eq. (49) by applying the transpose identity (see subsection “Matrix Transposition”). That is, we may find \mathbf{B} by solving the equation $\mathbf{A}^T\mathbf{B}^T = \mathbf{C}^T$ if \mathbf{A}^T meets the appropriate conditions.

Null Space When the rank of a matrix \mathbf{A} is greater than its row dimension N_r , it is possible to find nontrivial solutions to the equation

$$\mathbf{A}\mathbf{b} = \mathbf{0}_{N_r} \tag{50}$$

Indeed, it is possible to determine a linear model such that all vectors contained in the model satisfy Eq. (50). This linear model will have dimension equal to the difference between N_r and the rank of the matrix \mathbf{A} . The subspace defined by this linear model is called the null space of the matrix \mathbf{A} . Algorithms to find the basis vectors of a matrix’s null space exist and are provided by software packages that support matrix algebra.

Singular Value Decomposition

The singular value decomposition allows us to write any N_r by N_c matrix \mathbf{X} in the form

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T \tag{51}$$

where \mathbf{U} is an N_r by N_r orthogonal matrix, \mathbf{D} is an N_r by N_c “diagonal” matrix, and \mathbf{V} is an N_c by N_c orthogonal matrix.¹⁴⁸ The diagonal entries of \mathbf{D} are guaranteed to be nonnegative. Some of the diagonal entries may be zero. By convention, the entries along this diagonal are arranged in decreasing order. We illustrate the singular value decomposition in Fig. 21. The singular value decomposition has a large number of uses in numerical matrix algebra. Routines to compute it are generally provided as part of software packages that support matrix algebra.

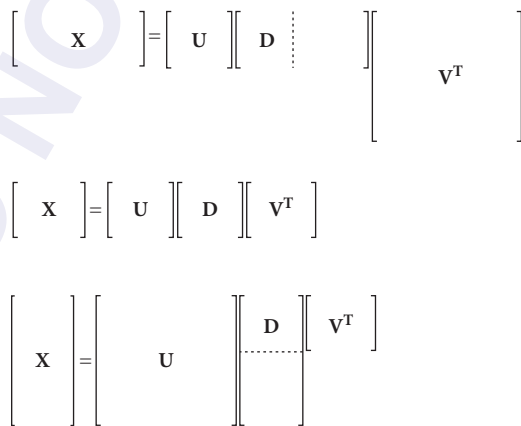


FIGURE 21 The figure depicts the singular value decomposition (SVD) of an N by M matrix \mathbf{X} for three cases: $N_c > N_r$, $N_c = N_r$, and $N_c < N_r$.

10.8 REFERENCES

1. J. von Kries, "Chromatic Adaptation," 1902, reprinted in *Sources of Color Vision*, D. L. MacAdam, (ed.), MIT Press, 1970, Cambridge, MA, pp. 109–119.
2. R. M. Evans, *An Introduction to Color*, Wiley, New York, 1948.
3. G. Wyszecki, "Color Appearance," in *Handbook of Perception and Human Performance: Sensory Processes and Perception*, K. R. Boff, L. Kaufman, and J. P. Thomas, (eds.), John Wiley & Sons, New York, 1986, pp. 9.1–9.56.
4. B. K. P. Horn, "Exact Reproduction of Colored Images," *Computer Vision, Graphics and Image Processing* **26**:135–167 (1984).
5. R. W. G. Hunt, *The Reproduction of Colour*, 4th ed., Fountain Press, Tolworth, England, 1987.
6. M. D. Fairchild, *Color Appearance Models*, Addison-Wesley, Reading, MA, 1998.
7. D. H. Brainard and B. A. Wandell, "Calibrated Processing of Image Color," *Color Research and Application* **15**:266–271 (1990).
8. G. Wyszecki and W. S. Stiles, *Color Science—Concepts and Methods, Quantitative Data and Formulae*, 2nd ed., John Wiley & Sons, New York, 1982.
9. V. C. Smith and J. Pokorny, "Color Matching and Color Discrimination," in *The Science of Color*, 2nd ed., S. K. Shevell, (ed.), Optical Society of America; Elsevier Ltd, Oxford, 2003, pp. 103–148.
10. CIE, *Colorimetry*, 3rd edition, Publication 15.2004, Bureau Central de la CIE, Vienna, 2004.
11. CIE, *Fundamental Chromaticity Diagram with Physiological Axes—Part 1, Publication 170–1*, Bureau Central de la CIE, Vienna, 2006.
12. D. H. Krantz, "Color Measurement and Color Theory: I. Representation Theorem for Grassmann Structures," *Journal of Mathematical Psychology* **12**:283–303 (1975).
13. P. Suppes, D. H. Krantz, R. D. Luce, and A. Tversky, *Foundations of Measurement*, Academic Press, San Diego, 1989, Vol. II.
14. D. L. MacAdam, *Sources of Color Science*, MIT Press, Cambridge, MA, 1970.
15. W. D. Wright, "The Origins of the 1931 CIE System," in *Human Color Vision*, 2nd ed., P. K. Kaiser and R. M. Boynton, (eds.), Optical Society of America, Washington, D.C., 1996, pp. 534–543.
16. J. D. Mollon, "The Origins of Modern Color Science," in *The Science of Color*, 2nd ed., S. K. Shevell, (ed.), Optical Society of America; Elsevier Ltd, Oxford, 2003, pp. 1–39.
17. R. Hall, *Illumination and Color in Computer Generated Imagery*, Springer-Verlag, New York, 1989.
18. D. B. Judd and G. Wyszecki, *Color in Business, Science, and Industry*, John-Wiley and Sons, New York, 1975.
19. G. S. Brindley, *Physiology of the Retina and the Visual Pathway*, 2nd ed., Williams and Wilkins, Baltimore, 1970.
20. R. M. Boynton, "History and Current Status of a Physiologically Based System of Photometry and Colorimetry," *Journal of the Optical Society of America A* **65**:1609–1621 (1996).
21. A. Stockman and L. T. Sharpe, "Cone Spectral Sensitivities and Color Matching," in *Color Vision: From Genes To Perception*, K. R. Gegenfurtner and L. T. Sharpe, (eds.), Cambridge University Press, Cambridge, MA, 1999, pp. 53–87.
22. B. A. Wandell, "The Synthesis and Analysis of Color Images," *IEEE Transactions on Pattern Analysis and Machine Intelligence* **9**:2–13 (1987).
23. S. Westland and C. Ripamonti, *Computational Colour Science using MATLAB*, John Wiley & Sons, 2004.
24. B. Noble and J. W. Daniel, *Applied Linear Algebra*, 2nd ed., Prentice-Hall, Inc., Englewood Cliffs, NJ, 1977.
25. G. H. Golub and C. F. van Loan, *Matrix Computations*, Johns Hopkins University Press, Baltimore, 1983.
26. W. K. Pratt, *Digital Image Processing*, John Wiley & Sons, New York, 1978.
27. W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling, *Numerical Recipes in C: The Art of Scientific Computing*, Cambridge University Press, Cambridge, MA, 1988.
28. G. Wyszecki and W. Stiles, *Color Science—Concepts and Methods, Quantitative Data and Formulas*, John Wiley & Sons, New York, 1967.
29. W. S. Stiles, "The Physical Interpretation of the Spectral Sensitivity Curve of the Eye," in *Transactions of the Optical Convention of the Worshipful Company of Spectacle Makers*, Spectacle Maker's Company, London, 1948, pp. 97–107.

30. D. E. Mitchell and W. A. H. Rushton, "Visual Pigments in Dichromats," *Vision Research* **11**:1033–1043 (1971).
31. W. D. Wright, "A Re-determination of the Trichromatic Coefficients of the Spectral Colours," *Transactions of the Optical Society* **30**:141–164 (1928–1929).
32. J. Guild, "The Colorimetric Properties of the Spectrum," *Philosophical Transactions of the Royal Society of London A* **230**:149–187 (1931).
33. W. S. Stiles, "Interim Report to the Commission Internationale de l'Éclairage Zurich, 1955, on the National Physical Laboratory's Investigation of Colour-matching" (with an Appendix by W. S. Stiles & J. M. Burch) *Optica Acta* **2**:168–181 (1955).
34. W. S. Stiles and J. M. Burch, "NPL Colour-matching Investigation: Final Report (1958)," *Optica Acta* **6**:1–26 (1959).
35. J. C. Maxwell, "On the Theory of Compound Colours and the Relations of the Colours of the Spectrum," *Philosophical Transactions of the Royal Society of London* **150**:57–84 (1860).
36. J. G. Grassmann, "Theory of Compound Colors," in MacAdam, D.L., (ed.), *Sources of Color Vision*, MIT Press: Cambridge, MA, 1970. (Originally published in *Annalen der Physik und Chemie*, **89**:69–84 1853).
37. CIE, *Commission Internationale de l'Éclairage Proceedings, 1931*, Cambridge University Press, Cambridge, MA, 1932.
38. D. B. Judd, "Report of U.S. Secretariat Committee on Colorimetry and Artificial Daylight," *Proceedings of the Twelfth Session of the CIE* **1**:11 (1951).
39. J. J. Vos, "Colorimetric and Photometric Properties of a 2° Fundamental Observer," *Color Research and Application* **3**:125–128 (1978).
40. ISO/CIE, *CIE Standard Colorimetric Observers*, Reference Number 10527, International Organization for Standardization, Geneva, 1991.
41. I. E. Abdou and W. K. Pratt, "Quantitative Design and Evaluation of Enhancement/Thresholding Edge Detectors," *Proceedings of the IEEE* **67**:753–763 (1979).
42. CIE, *Commission Internationale de l'Éclairage Proceedings, 1924*, Cambridge University Press, Cambridge, MA, 1926.
43. H. G. Sperling, "An Experimental Investigation of the Relationship Between Colour Mixture and Luminous Efficiency," in *Visual Problems of Colour*, vol. 1, Her Majesty's Stationery Office, London, 1958, pp. 249–277.
44. A. Stockman, L. T. Sharpe, and C. C. Fach, "The Spectral Sensitivity of the Human Short-wavelength Cones," *Vision Research* **39**:2901–2927 (1999).
45. V. C. Smith, J. Pokorny, and Q. Zaidi, "How do Sets of Color-matching Functions Differ?," in *Colour Vision*, J. D. Mollon and L. T. Sharpe, (eds.), Academic Press, London, 1983, pp. 93–105.
46. A. Stockman and L. T. Sharpe, "Spectral Sensitivities of the Middle- and Long-Wavelength Sensitive Cones Derived from Measurements in Observers of Known Genotype," *Vision Research* **40**:1711–1737 (2000).
47. P. DeMarco, J. Pokorny, and V. C. Smith, "Full-spectrum Cone Sensitivity Functions for X-chromosome-linked Anomalous Trichromats," *Journal of the Optical Society of America* **9**:1465–1476 (1992).
48. N. I. Speranskaya, "Determination of Spectrum Color Coordinates for Twenty-seven Normal Observers," *Optics and Spectroscopy* **7**:424–428 (1959).
49. O. N. Rood, *Modern Chromatics, with Applications to Art and Industry*, D. Appleton & Co., New York, 1879.
50. T. Young, "On the Theory of Light and Colours," *Philosophical Transactions of the Royal Society of London* **92**:12–48 (1802).
51. H. L. F. von Helmholtz, "On the Theory of Compound Colours," *Philosophical Magazine Series, 4* **4**:519–534 (1852).
52. J. C. Maxwell, "Experiments on Colours, as Perceived by the Eye, with Remarks on Colour-blindness," *Transactions of the Royal Society of Edinburgh* **21**:275–298 (1855).
53. V. Smith and J. Pokorny, "Spectral Sensitivity of the Foveal Cone Photopigments between 400 and 500 nm," *Vision Research* **15**:161–171 (1975).
54. J. C. Maxwell, "On the Theory of Colours in Relation to Colour-blindness. A letter to Dr. G. Wilson," *Transactions of the Royal Scottish Society of Arts* **4**:394–400 (1856).
55. J. Nathans, T. P. Piantanida, R. L. Eddy, T. B. Shows, and D. S. Hogness, "Molecular Genetics of Inherited Variation in Human Color Vision," *Science* **232**:203–210 (1986).
56. J. Nathans, D. Thomas, and D. S. Hogness, "Molecular Genetics of Human Color Vision: the Genes Encoding Blue, Green and Red Pigments," *Science* **232**:193–202 (1986).

57. M. M. Bongard and M. S. Smirnov, "Determination of the Eye Spectral Sensitivity Curves from Spectral Mixture Curves," *Doklady Akademiia nauk S.S.S.R.* **102**:1111–1114 (1954).
58. A. Stockman, D. I. A. MacLeod, and N. E. Johnson, "Spectral Sensitivities of the Human Cones," *Journal of the Optical Society of America A* **10**:2491–2521 (1993).
59. A. König and C. Dieterici, "Die Grundempfindungen in Normalen und anomalen Farben Systemen und ihre Intensitäts-Verthielung im Spectrum," *Z Psychol Physiol Sinnesorg* **4**:241–347 (1893).
60. P. J. Bouma, "Mathematical Relationship between the Colour Vision System of Trichromats and Dichromats," *Physica* **9**:773–784 (1942).
61. D. B. Judd, "Standard Response Functions for Protanopic and Deuteranopic Vision," *Journal of the Optical Society of America* **35**:199–221 (1945).
62. D. B. Judd, "Standard Response Functions for Protanopic and Deuteranopic Vision," *Journal of the Optical Society of America* **39**:505 (1949).
63. J. J. Vos and P. L. Walraven, "On the Derivation of the Foveal Receptor Primaries," *Vision Research* **11**:799–818 (1971).
64. O. Estévez, "On the Fundamental Database of Normal and Dichromatic Color Vision," Ph.D., Amsterdam University, 1979.
65. J. J. Vos, O. Estevez, and P. L. Walraven, "Improved Color Fundamentals Offer a New View on Photometric Additivity," *Vision Research* **30**:937–943 (1990).
66. A. Stockman, D. I. A. MacLeod, and J. A. Vivien, "Isolation of the Middle- and Long-wavelength Sensitive Cones in Normal Trichromats," *Journal of the Optical Society of America A* **10**:2471–2490 (1993).
67. M. A. Webster and D. I. A. MacLeod, "Factors Underlying Individual Differences in the Color Matches of Normal Observers," *Journal of the Optical Society of America A* **5**:1722–1735 (1988).
68. G. Wald, "Human Vision and the Spectrum," *Science* **101**:653–658 (1945).
69. R. A. Bone and J. M. B. Sparrock, "Comparison of Macular Pigment Densities in the Human Eye," *Vision Research* **11**:1057–1064 (1971).
70. P. L. Pease, A. J. Adams, and E. Nuccio, "Optical Density of Human Macular Pigment," *Vision Research* **27**:705–710 (1987).
71. D. van Norren and J. J. Vos, "Spectral Transmission of the Human Ocular Media," *Vision Research* **14**:1237–1244 (1974).
72. B. H. Crawford, "The Scotopic Visibility Function," *Proceedings of the Physical Society of London B* **62**:321–334 (1949).
73. F. S. Said and R. A. Weale, "The Variation with Age of the Spectral Transmissivity of the Living Human Crystalline Lens," *Gerontologia* **3**:213–231 (1959).
74. J. Pokorny, V. C. Smith, and M. Lutze, "Aging of the Human Lens," *Applied Optics* **26**:1437–1440 (1987).
75. M. Alpern, "Lack of Uniformity in Colour Matching," *Journal of Physiology* **288**:85–105 (1979).
76. H. Terstiege, "Untersuchungen zum Persistenz- und Koeffizientensatz," *Die Farbe* **16**:1–120 (1967).
77. S. S. Miller, "Psychophysical Estimates of Visual Pigment Densities in Red-green Dichromats," *Journal of Physiology* **223**:89–107 (1972).
78. P. E. King-Smith, "The Optical Density of Erythrolabe Determined by a New Method," *Journal of Physiology* **230**:551–560 (1973).
79. P. E. King-Smith, "The Optical Density of Erythrolabe Determined by Retinal Densitometry Using the Self-screening Method," *Journal of Physiology* **230**:535–549 (1973).
80. V. Smith and J. Pokorny, "Psychophysical Estimates of Optical Density in Human Cones," *Vision Research* **13**:1099–1202 (1973).
81. S. A. Burns and A. E. Elsner, "Color Matching at High Luminances: Photopigment Optical Density and Pupil Entry," *Journal of the Optical Society of America A* **10**:221–230 (1993).
82. T. T. J. M. Berendschot, J. van der Kraats, and D. van Norren, "Foveal Cone Mosaic and Visual Pigment Density in Dichromats," *Journal of Physiology* **492**:307–314 (1996).
83. C. A. Curcio, K. R. Sloan, R. E. Kalina, and A. E. Hendrickson, "Human Photoreceptor Topography," *Journal of Comparative Neurology* **292**:497–523 (1990).

84. A. E. Elsner, S. A. Burns, and R. H. Webb, "Mapping Cone Photopigment Optical Density," *Journal of the Optical Society of America A* **10**:52–58 (1993).
85. M. Neitz, J. Neitz, and G. H. Jacobs, "Spectral Tuning of Pigments Underlying Red-green Color Vision," *Science* **252**:971–973 (1991).
86. S. C. Merbs and J. Nathans, "Absorption Spectra of Human Cone Photopigments," *Nature* **356**:433–435 (1992).
87. S. L. Merbs and J. Nathans, "Absorption Spectra of the Hybrid Pigments Responsible for Anomalous Color Vision," *Science* **258**:464–466 (1992).
88. J. Winderickx, D. T. Lindsey, E. Sanocki, D. Y. Teller, A. G. Motulsky, and S. S. Deeb, "Polymorphism in Red Photopigment Underlies Variation in Colour Matching," *Nature* **356**:431–433 (1992).
89. L. T. Sharpe, A. Stockman, H. Jägle, and J. Nathans, "Opsin Genes, Cone Photopigments, Color Vision, and Color Blindness," in *Color Vision: From Genes To Perception*, K. R. Gegenfurtner and L. T. Sharpe, (eds.), Cambridge University Press, Cambridge, MA, 1999, pp. 3–51.
90. D. H. Brainard, A. Roorda, Y. Yamauchi, J. B. Calderone, A. Metha, M. Neitz, J. Neitz, D. R. Williams, and G. H. Jacobs, "Functional Consequences of the Relative Numbers of L and M Cones," *Journal of the Optical Society of America A* **17**:607–614 (2000).
91. M. Drummond-Borg, S. S. Deeb, and A. G. Motulsky, "Molecular Patterns of X Chromosome-linked Color Vision Genes among 134 Men of European Ancestry," *Proceedings of the National Academy of Sciences* **86**:983–987 (1989).
92. J. Neitz, M. Neitz, and G. H. Jacobs, "More than 3 Different Cone Pigments among People with Normal Color Vision," *Vision Research* **33**:117–122 (1993).
93. J. P. Macke and J. Nathans, "Individual Variation in the Size of the Human Red and Green Pigment Gene Array," *Investigative Ophthalmology and Visual Science* **38**:1040–1043 (1997).
94. A. B. Asenjo, J. Rim, and D. D. Oprian, "Molecular Determinants of Human Red/Green Color Discrimination," *Neuron* **12**:1131–1138 (1994).
95. L. T. Sharpe, A. Stockman, H. Jägle, H. Knau, and J. Nathans, "L, M, and L-M Hybrid Cone Photopigments in Man: Deriving λ_{\max} from Flicker Photometric Spectral Sensitivities," *Vision Research* **39**:3513–3525 (1999).
96. M. Neitz, J. Neitz, and G. H. Jacobs, "Genetic Basis of Photopigment Variations in Human Dichromats," *Vision Research* **35**:2095–2103 (1995).
97. A. Stockman, L. T. Sharpe, S. Merbs, and J. Nathans, "Spectral Sensitivities of Human Cone Visual Pigments Determined *in vivo* and *in vitro*," in *Vertebrate Phototransduction and the Visual Cycle, Part B. Methods in Enzymology, Vol. 316*, K. Palczewski, (ed.), Academic Press, New York, 2000, pp. 626–650.
98. J. J. Kremers, T. Usui, H. P. Scholl, and L. T. Sharpe, "Cone Signal Contributions to Electrograms in Dichromats and Trichromats," *Investigative Ophthalmology and Visual Science* **40**:920–930 (1999).
99. A. Stockman, H. Jägle, M. Pirzer, and L. T. Sharpe, "The Dependence of Luminous Efficiency on Chromatic Adaptation," *Journal of Vision* **8**:1–26 (2008).
100. A. Reitner, L. T. Sharpe, and E. Zrenner, "Is Colour Vision Possible with Only Rods and Blue-sensitive Cones?," *Nature* **352**:798–800 (1991).
101. A. König, and C. Dieterici, "Die Grundempfindungen und ihre Intensitäts-Vertheilung im Spectrum," *Sitzungsberichte Akademie der Wissenschaften*, Berlin **1886**:805–829 (1886).
102. S. Ishihara, *Tests for Colour-Blindness*, Kanehara Shuppen Company, Ltd., Tokyo, 1977.
103. D. Farnsworth, "The Farnsworth-Munsell 100 Hue and Dichotomous Tests for Color Vision," *Journal of the Optical Society of America* **33**:568–578 (1943).
104. L. Rayleigh, "Experiments on Colour," *Nature* **25**:64–66 (1881).
105. J. B. Birch, *Diagnosis of Defective Colour Vision*, Oxford University Press, Oxford, 1993.
106. A. König, "Über den Menschlichen Sehpurpur und seine Bedeutung für das Sehen," *Academie der Wissenschaften Sitzungsberichte* **30**:577–598 (1894).
107. E. N. Willmer, "Further Observations on the Properties of the Central Fovea in Colour-blind and Normal Subjects," *Journal of Physiology* **110**:422–446 (1950).
108. L. C. Thomson and W. D. Wright, "The Colour Sensitivity of the Retina within the Central Fovea of Man," *Journal of Physiology* **105**:316–331 (1947).

109. D. Williams, D. I. A. MacLeod, and M. Hayhoe, "Foveal Tritanopia," *Vision Research* **21**:1341–1356 (1981).
110. D. I. Flitcroft, "The Interactions between Chromatic Aberration, Defocus and Stimulus Chromaticity: Implications for Visual Physiology and Colorimetry," *Vision Research* **29**:349–360 (1989).
111. D. R. Williams, N. Sekiguchi, W. Haake, D. H. Brainard, and O. Packer, "The Cost of Trichromacy for Spatial Vision," in *From Pigments to Perception*, B. B. Lee and A. Valberg, (eds.), Plenum Press, New York, 1991, pp. 11–22.
112. D. H. Marimont and B. A. Wandell, "Matching Color Images—the Effects of Axial Chromatic Aberration," *Journal of the Optical Society of America A* **11**:3113–3122 (1994).
113. L. N. Thibos, M. Ye, X. X. Zhang, and A. Bradley, "The Chromatic Eye—a New Reduced-eye Model of Ocular Chromatic Aberration in Humans," *Applied Optics* **31**:3594–3667 (1992).
114. I. Powell, "Lenses for Correcting Chromatic Aberration of the Eye," *Applied Optics* **20**:4152–4155 (1981).
115. K. H. Ruddock, "The Effect of Age Upon Colour Vision. II. Changes with Age in Light Transmission of the Ocular Media," *Vision Research* **5**:47–58 (1965).
116. A. Knowles and H. J. A. Dartnall, "The Photobiology of Vision," in *The Eye*, H. Davson, (ed.), Academic Press, New York, 1977, pp. 321–533.
117. H. J. A. Dartnall, "The Interpretation of Spectral Sensitivity Curves," *British Medical Bulletin* **9**:24–30 (1953).
118. R. J. W. Mansfield, "Primate Photopigments and Cone Mechanisms," in *The Visual System*, A. Fein and J. S. Levine, (eds.), Alan R. Liss, New York, 1985.
119. E. F. MacNichol, Jr., "A Unifying Presentation of Photopigment Spectra," *Vision Research* **26**:1543–1556 (1986).
120. T. D. Lamb, "Photoreceptor Spectral Sensitivities: Common Shape in the Long-Wavelength Spectral Region," *Vision Research* **35**:3083–3091 (1995).
121. H. B. Barlow, "What Causes Trichromacy? A Theoretical Analysis using Comb-filtered Spectra," *Vision Research* **22**:635–643 (1982).
122. V. I. Govardovskii, N. Fyhrquist, T. Reuter, D. G. Kuzmin, and K. Donner, "In Search of the Visual Pigment Template," *Visual Neuroscience* **17**:509–528 (2000).
123. J. Walraven, C. Enroth-Cugell, D. C. Hood, D. I. A. MacLeod, and J. L. Schnapf, "The Control of Visual Sensitivity. Receptor and Postreceptor Processes," in *Visual Perception: The Neurophysiological Foundations*, L. Spillman and J. S. Werner, (eds.), Academic Press, San Diego, 1990, pp. 53–101.
124. A. M. Derrington, J. Krauskopf, and P. Lennie, "Colour-opponent Cells in the Dorsal Lateral Geniculate Nucleus of the Macaque," *Journal of Physiology* **329**:22–23 (1982).
125. D. I. A. MacLeod and R. M. Boynton, "Chromaticity Diagram Showing Cone Excitation by Stimuli of Equal Luminance," *Journal of the Optical Society of America* **69**:1183–1186 (1979).
126. R. Luther, "Aus dem Gebiet der Farbreizmetrik," *Zeitschrift für Technische Physik* **8**:540–558 (1927).
127. J. Krauskopf, D. R. Williams, and D. W. Heeley, "Cardinal Directions of Color Space," *Vision Research* **22**:1123–1131 (1982).
128. J. Krauskopf, D. R. Williams, M. B. Mandler, and A. M. Brown, "Higher Order Color Mechanisms," *Vision Research* **26**:23–32 (1986).
129. D. H. Brainard, "Cone Contrast and Opponent Modulation Color Spaces," in *Human Color Vision*, 2nd ed., P. K. Kaiser and R. M. Boynton, Optical Society of America, Washington, D.C., 1996, pp. 563–579.
130. A. B. Poirson and B. A. Wandell, "The Ellipsoidal Representation of Spectral Sensitivity," *Vision Research* **30**:647–652 (1990).
131. R. Ramamoorthi and P. Hanrahan, "A Signal-processing Framework for Reflection," *ACM Transactions on Graphics* **23**:1004–1042 (2004).
132. R. O. Dror, A. S. Willsky, and E. H. Adelson, "Statistical Characterization of Real-world Illumination," *Journal of Vision* **4**:821–837 (2004).
133. R. W. Fleming, R. O. Dror, and E. H. Adelson, "Real-world Illumination and the Perception of Surface Reflectance Properties," *Journal of Vision* **3**:347–368 (2003).
134. R. W. Fleming, A. Torralba, and E. H. Adelson, "Specular Reflections and the Perception of Shape," *Journal of Vision* **4**:798–820 (2004).

135. B. Xiao and D. H. Brainard, "Surface Gloss and Color Perception of 3D Objects," *Visual Neuroscience* **25**: 371–385 (2008).
136. I. Motoyoshi, S. Nishida, L. Sharan, and E. H. Adelson, "Image Statistics and the Perception of Surface Qualities," *Nature* **447**:206–209 (2007).
137. D. B. Judd, D. L. MacAdam, and G. W. Wyszecki, "Spectral Distribution of Typical Daylight as a Function of Correlated Color Temperature," *Journal of the Optical Society of America* **54**:1031–1040 (1964).
138. J. Cohen, "Dependency of the Spectral Reflectance Curves of the Munsell Color Chips," *Psychonomic Science* **1**:369–370 (1964).
139. K. L. Kelly, K. S. Gibson, and D. Nickerson, "Tristimulus Specification of the Munsell Book of Color from Spectrophotometric Measurements," *Journal of the Optical Society of America* **33**:355–376 (1943).
140. D. Nickerson, "Spectrophotometric Data for a Collection of Munsell Samples," U.S. Department of Agriculture, 1957.
141. L. T. Maloney, "Evaluation of Linear Models of Surface Spectral Reflectance with Small Numbers of Parameters," *Journal of the Optical Society of America A* **3**:1673–1683 (1986).
142. E. L. Krinov, "Surface Reflectance Properties of Natural Formations," National Research Council of Canada: Technical Translation **TT-439** (1947).
143. J. P. S. Parkkinen, J. Hallikainen, and T. Jaaskelainen, "Characteristic Spectra of Munsell Colors," *Journal of the Optical Society of America* **6**:318–322 (1989).
144. T. Jaaskelainen, J. Parkkinen, and S. Toyooka, "A Vector-subspace Model for Color Representation," *Journal of the Optical Society of America A* **7**:725–730 (1990).
145. J. R. Magnus and H. Neudecker, *Matrix Differential Calculus with Applications in Statistics and Econometrics*, Wiley, Chichester, 1988.
146. T. W. Anderson, *An Introduction to Multivariate Statistical Analysis*, 2nd ed., John Wiley and Sons, New York, 1971.
147. R. A. Johnson and D. W. Wichern, *Applied Multivariate Statistical Analysis*, Prentice-Hall, Englewood Cliffs, NJ, 1988.
148. J. M. Chambers, *Computational Methods for Data Analysis*, John Wiley and Sons, New York, 1977.
149. D. H. Marimont and B. A. Wandell, "Linear Models of Surface and Illuminant Spectra," *Journal of the Optical Society of America A* **9**:1905–1913 (1992).
150. B. A. Wandell and D. H. Brainard, "Towards Cross-media Color Reproduction," Presented at the OSA Applied Vision Topical Meeting, San Francisco, CA, 1989.
151. C. A. Parraga, G. Brelstaff, T. Troscianko, and I. R. Moorehead, "Color and Luminance Information in Natural Scenes," *Journal of the Optical Society of America A* **15**:563–569 (1998).
152. P. Longère and D. H. Brainard, "Simulation of Digital Camera Images from Hyperspectral Input," in *Vision Models and Applications to Image and Video Processing*, C. J. van den Branden Lambrecht (ed.), Kluwer Academic, Boston, 2001, pp. 123–150.
153. D. H. Foster, S. M. C. Nascimento, and K. Amano, "Information Limits on Neural Identification of Colored Surfaces in Natural Scenes," *Visual Neuroscience* **21**:1–6 (2004).
154. D. Saunders and A. Hamber, "From Pigments to Pixels: Measurement and Display of the Colour Gamut of Paintings," Proceedings of the SPIE: *Perceiving, Measuring, and Using Color* **1250**:90–102 (1990).
155. R. Berns, "Rejuvenating Seurat's Palette Using Color and Imaging Science: a Simulation," in *Seurat and the Making of La Grande Jatte*, R. L. Herbert, (ed.), University of California Press, 2004, pp. 214–227.
156. D. H. Brainard, D. G. Pelli, and T. Robson, "Display Characterization," in *Encyclopedia of Imaging Science and Technology*, J. Hornak, (ed.) Wiley, 2002, pp. 72–188.
157. ISO/CIE, *CIE Standard Colorimetric Illuminants*, Reference Number 10526, International Organization for Standardization, Geneva, 1991.
158. D. H. Brainard, B. A. Wandell, and W. B. Cowan, "Black Light: How Sensors Filter Spectral Variation of the Illuminant," *IEEE Transactions on Biomedical Engineering* **36**:140–149 (1989).
159. G. Wyszecki, "Evaluation of Metameric Colors," *Journal of the Optical Society of America* **48**:451–454 (1958).
160. S. A. Burns, J. B. Cohen, and E. N. Kuznetsov, "Multiple Metatmers: Preserving Color Matches Under Diverse Illuminants," *Color Research and Application* **14**:16–22 (1989).

161. G. D. Finlayson and M. S. Drew, "The Maximum Ignorance Assumption with Positivity," Presented at the 4th IS&T/SID Color Imaging Conference, Scottsdale, AZ, 1996, pp. 202–204.
162. J. A. S. Viggiano, "Minimal-knowledge Assumptions in Digital Still Camerage Characterization. I. Uniform Distribution, Toeplitz Correlation," Presented at the 9th IS&T/SID Color Imaging Conference, Scottsdale, AZ, 2001, pp. 332–336.
163. X. Zhang and D. H. Brainard, "Bayesian Color-correction Method for Non-colorimetric Digital Image Sensors," Presented at the 12th IS&T/SID Color Imaging Conference, Scottsdale, AZ, 2004, pp. 308–314.
164. M. H. Brill, "A Device Performing Illuminant-invariant Assessment of Chromatic Relations," *Journal of Theoretical Biology* **71**:473–478 (1978).
165. G. Buchsbaum, "A Spatial Processor Model for Object Colour Perception," *Journal of the Franklin Institute* **310**:1–26 (1980).
166. L. T. Maloney and B. A. Wandell, "Color Constancy: A Method for Recovering Surface Spectral Reflectances," *Journal of the Optical Society of America A* **3**:29–33 (1986).
167. H. C. Lee, "Method for Computing the Scene-illuminant Chromaticity from Specular Highlights," *Journal of the Optical Society of America A* **3**:1694–1699 (1986).
168. B. Funt and J. Ho, "Color from Black and White," Presented at the International Conference on Computer Vision, Tampa, FL, 1988, pp. 2–8.
169. B. V. Funt and M. S. Drew, "Color Constancy Computation in Near-Mondrian Scenes Using a Finite Dimensional Linear Model," Presented at the IEEE Conference on Computer Vision and Pattern Recognition, Ann Arbor, MI, 1988, pp. 544–549.
170. D. A. Forsyth, "A Novel Approach to Colour Constancy," Presented at the International Conference on Computer Vision, Tampa, FL, 1988, pp. 9–18.
171. G. D. Finlayson, P. H. Hubel, and S. Hordley, "Color by correlation," Presented at the IS&T/SID Fifth Color Imaging Conference, Scottsdale, AZ, 1997, pp. 6–11.
172. D. H. Brainard and W. T. Freeman, "Bayesian Color Constancy," *Journal of the Optical Society of America A* **14**:1393–1411 (1997).
173. M. Ebner, *Color Constancy*, Wiley, Chichester, UK, 2007.
174. M. D'Zmura and P. Lennie, "Mechanisms of Color Constancy," *Journal of the Optical Society of America A* **3**:1662–1672 (1986).
175. A. C. Hurlbert, "Computational Models of Color Constancy," in *Perceptual Constancy: Why Things Look As They Do*, V. Walsh and J. Kulikowski, (eds.), Cambridge University Press, Cambridge, MA, 1998, pp. 283–322.
176. L. T. Maloney, "Physics-Based Approaches to Modeling Surface Color Perception," in *Color Vision: From Genes to Perception*, K. T. Gegenfurtner and L. T. Sharpe, (eds.), Cambridge University Press, Cambridge, MA, 1999, pp. 387–416.
177. D. H. Brainard, J. M. Kraft, and P. Longère, "Color Constancy: Developing Empirical Tests of Computational Models," in *Colour Perception: Mind and the Physical World*, R. Mausfeld and D. Heyer, (eds.), Oxford University Press, Oxford, 2003, pp. 307–334.
178. D. H. Brainard, P. Longere, P. B. Delahunt, W. T. Freeman, J. M. Kraft, and B. Xiao, "Bayesian Model of Human Color Constancy," *Journal of Vision* **6**:1267–1281 (2006).
179. D. L. MacAdam, "Visual Sensitivities to Color Differences in Daylight," *Journal of the Optical Society of America* **32**:247–274 (1942).
180. D. L. MacAdam, "Colour Discrimination and the Influence of Colour Contrast on Acuity" *Documenta Ophthalmologica* **3**:214–233 (1949).
181. G. Wyszecki, "Matching Color Differences," *Journal of the Optical Society of America* **55**:1319–1324 (1965).
182. G. Wyszecki and G. H. Fielder, "New Color-matching Ellipses," *Journal of the Optical Society of America* **61**:1135–1152 (1971).
183. W. S. Stiles, "Color vision: The Approach Through Increment Threshold Sensitivity," *Proceedings National Academy of Sciences (USA)* **45**:100–114 (1959).
184. C. Noorlander and J. J. Koenderink, "Spatial and Temporal Discrimination Ellipsoids in Color Space," *Journal of the Optical Society of America* **73**:1533–1543 (1983).
185. A. B. Poirson, B. A. Wandell, D. Varner, and D. H. Brainard, "Surface Characterizations of Color Thresholds," *Journal of the Optical Society of America A* **7**:783–789 (1990).

186. A. R. Robertson, "The CIE 1976 Color-difference Formula," *Color Research and Application* **2**:7–11 (1977).
187. H. de Lange, "Research into the Dynamic Nature of the Human Fovea-cortex Systems with Intermittent and Modulated Light. I. Attenuation Characteristics with White and Coloured Light," *Journal of the Optical Society of America* **48**:777–784 (1958).
188. H. de Lange, "Research into the Dynamic Nature of the Human Fovea-cortex Systems with Intermittent and Modulated Light. II. Phase Shift in Brightness and Delay in Color Perception," *Journal of the Optical Society of America* **48**:784–789 (1958).
189. N. Sekiguchi, D. R. Williams, and D. H. Brainard, "Efficiency for Detecting Isoluminant and Isochromatic Interference Fringes," *Journal of the Optical Society of America A* **10**:2118–2133 (1993).
190. E. C. Carter and R. C. Carter, "Color and Conspicuousness," *Journal of the Optical Society of America* **71**:723–729 (1981).
191. A. B. Poirson and B. A. Wandell, "Task-Dependent Color Discrimination," *Journal of the Optical Society of America A* **7**:776–782 (1990).
192. A. L. Nagy and R. R. Sanchez, "Critical Color Differences Determined with a Visual Search Task," *Journal of the Optical Society of America A* **7**:1209–1217 (1990).
193. H. E. Smithson, S. S. Khan, L. T. Sharpe, and A. Sotckman, "Transitions between Color Categories Mapped with a Reverse Stroop Task," *Visual Neuroscience* **23**:453–460 (2006).
194. CIE, *Improvement to Industrial Colour-Difference Evaluation*, Publication 142, Bureau Central de la CIE, Vienna, 2001.
195. G. Wagner and R. M. Boynton, "Comparison of Four Methods of Heterochromatic Photometry," *Journal of the Optical Society of America* **62**:1508–1515 (1972).
196. W. Abney and E. R. Festing, "Colour Photometry," *Philosophical Transactions of the Royal Society of London* **177**:423–456 (1886).
197. W. Abney, *Researches in Colour Vision*, Longmans, Green, London, 1913, p. 418.
198. A. Dresler, "The Non-additivity of Heterochromatic Brightness," *Transactions of the Illuminating Engineering Society* (London) **18**:141–165 (1953).
199. R. M. Boynton and P. Kaiser, "Vision: The Additivity Law Made to Work for Heterochromatic Photometry with Bipartite Fields," *Science* **161**:366–368 (1968).
200. S. L. Guth, N. J. Donley, and R. T. Marrocco, "On Luminance Additivity and Related Topics," *Vision Research* **9**:537–575 (1969).
201. Y. Le Grand, "Spectral Luminosity," in *Visual Psychophysics, Handbook of Sensory Physiology*, D. Jameson and L. H. Hurvich, (eds.), Springer-Verlag, Berlin, 1972, pp. 413–433.
202. J. Pokorny, V. C. Smith, and M. Lutze, "Heterochromatic Modulation Photometry," *Journal of the Optical Society of America A* **6**:1618–1623 (1989).
203. P. Lennie, J. Pokorny, and V. C. Smith, "Luminance," *Journal of the Optical Society of America A* **10**:1283–1293 (1993).
204. CIE, *CIE 1988 2° Spectral Luminous Efficiency Function for Photopic Vision*, Publication 86, Bureau Central de la CIE, Vienna, 1990.
205. J. Pokorny and V. C. Smith, "Colorimetry and Color Discrimination," in *Handbook of Perception and Human Performance*, K. R. Boff, L. Kaufman, and J. P. Thomas, (eds.), John Wiley & Sons, 1986.
206. A. Stockman, "Colorimetry," in *The Optics Encyclopedia: Basic Foundations and Practical Applications*, T. G. Brown, K. Creath, H. Kogelnik, M. A. Kriss, J. Schmit, and M. J. Weber, (eds.), Wiley-VCH, Berlin, 2003, pp. 207–226.
207. J. B. Cohen and W. E. Kappauf, "Metameric Color Stimuli, Fundamental Metamers, and Wyszecki's Metameric Blacks: Theory, Algebra, Geometry, Application," *American Journal of Psychology* **95**:537–564 (1982).

COLOR VISION MECHANISMS

Andrew Stockman

*Department of Visual Neuroscience
UCL Institute of Ophthalmology
London, United Kingdom*

David H. Brainard

*Department of Psychology
University of Pennsylvania
Philadelphia, Pennsylvania*

11.1 GLOSSARY

Achromatic mechanism. Hypothetical psychophysical mechanisms, sometimes equated with the luminance mechanism, which respond primarily to changes in intensity. Note that achromatic mechanisms may have spectrally opponent inputs, in addition to their primary nonopponent inputs.

Bezold-Brücke hue shift. The shift in the hue of a stimulus toward either the yellow or blue invariant hues with increasing intensity.

Bipolar mechanism. A mechanism, the response of which has two mutually exclusive types of output that depend on the balance between its two opposing inputs. Its response is nulled when its two inputs are balanced.

Brightness. A perceptual measure of the apparent intensity of lights. Distinct from luminance in the sense that lights that appear equally bright are not necessarily of equal luminance.

Cardinal directions. Stimulus directions in a three-dimensional color space that isolate two of the three “cardinal mechanisms.” These are the isolating directions for the L+M, L–M, and S–(L+M) mechanisms. Note that the isolating directions do not necessarily correspond to mechanism directions.

Cardinal mechanisms. The second-site bipolar L–M and S–(L+M) chromatic mechanisms and the L+M luminance mechanism.

Chromatic discrimination. Discrimination of a chromatic target from another target or background, typically measured at equiluminance.

Chromatic mechanism. Hypothetical psychophysical mechanisms that respond to chromatic stimuli, that is, to stimuli modulated at equiluminance.

Color appearance. Subjective appearance of the hue, brightness, and saturation of objects or lights.

Color-appearance mechanisms. Hypothetical psychophysical mechanisms that mediate color appearance, especially as determined in hue scaling or color valence experiments.

Color assimilation. The phenomenon in which the hue of an area is perceived to be closer to that of the surround than to its hue when viewed in isolation. Also known as the von Bezold spreading effect.

Color constancy. The tendency of objects to retain their color appearance despite changes in the spectral characteristics of the illuminant, or, more generally, despite changes in viewing context.

Color contrast. The change in the color appearance of an area caused by the presence of a colored surround. The color change, unlike assimilation, is usually complementary to the surround color.

Color-discrimination mechanisms. Hypothetical psychophysical mechanisms that determine performance in chromatic detection or discrimination tasks. Assumed in some models to correspond to cone-opponent mechanisms.

Color spaces. Representations of lights either in terms of the responses of some known or hypothetical mechanisms thought to underlie the perception of color (such as cone or postreceptoral mechanisms), or in terms of the projection of the lights onto stimulus-based vectors (such as monochromatic primaries or mechanism-isolating vectors).

Color valence. A measure of the color of a light in terms of the amount of a cancelling light required to null one of the hue sensations produced by that light. Thus, if a light appears red it is cancelled by light that appears green, and the amount of this green light is its red valence. In opponent-colors theory, color appearance depends on the relative red-green and blue-yellow valences.

Cone contrast. The contrast (or relative change in quantal or energy catch) presented to each cone photoreceptor: $\Delta L/L$, $\Delta M/M$, and $\Delta S/S$.

Cone contrast space. A color space where the position along each axis represents the contrast of one cone class.

Cone mechanisms. Hypothetical psychophysical mechanisms, the performances of which are limited at the cone photoreceptors.

Cone-opponent mechanism. Hypothetical psychophysical mechanisms with opposed cone inputs.

Derrington Krauskopf Lennie (DKL) space. Color space, the axes of which are the stimulus strengths in each of the three cardinal mechanism directions. Closely related to the spaces proposed by Lüther⁸⁵ and MacLeod and Boynton.⁸⁶ In some accounts of this space the axes are defined in a different way, in terms of the three vectors that isolate each of the three cardinal mechanisms.

Detection surface or contour. Detection thresholds measured in many directions in color space form a detection surface. Confined to a plane, they form a contour. The terms threshold surface and threshold contour are synonymous with detection surface and detection contour, respectively.

Field method. A method in which the observer's sensitivity for detecting or discriminating a target is measured as a function of some change in context or in the adapted state of the mechanism of interest.

First-site adaptation. Adaptation, usually assumed to be cone-class specific, occurring at or related to the photoreceptor level.

Habituation. Loss of sensitivity caused by prolonged adaptation to chromatic and/or achromatic stimulus modulations, also known as contrast adaptation.

Incremental cone-excitation space. A color space in which the axes represent the deviations of each of the three classes of cones from a background. Deviations can be negative (decrements) as well as increments.

Intensity. Generic term to denote variation in stimulus or modulation strength when chromatic properties are held constant. In the particular context of modulations around a background, the vector length of a modulation may be used as a measure of intensity.

Invariant hue. A stimulus produces an invariant hue if that hue is independent of changes to stimulus intensity. Generally studied in the context of monochromatic stimuli.

Isolating direction. Direction in a color space that isolates the response of a single mechanism.

Linear visual mechanisms. Hypothetical mechanisms that behave linearly, usually with respect to the cone isomerization rates, but in some models with respect to the cone outputs after von Kries adaptation or contrast coding.

Luminance. A measure of the efficiency (or effectiveness) of lights often linked to the assumed output of the achromatic mechanism.

Mechanism direction. Stimulus color direction along which a specified mechanism is most sensitive. Note that the mechanism direction is not, in general, the same as the isolating direction for the same mechanism.

Noise masking. Threshold elevations caused by superimposing targets in noise.

Nonlinear visual mechanisms. Hypothetical mechanisms that behave nonlinearly either with respect to the cone inputs or with respect to their own (assumed) inputs.

Opponent-colors theory. A color theory that accounts for color appearance in the terms of the perceptual opposition of red and green (R/G), blue and yellow (B/Y), and dark and light (W/B).

Pedestal effects. Changes in sensitivity that occur when a target is superimposed on another stimulus, called the pedestal, which may have either identical or different spatio-chromatic-temporal characteristics to the target.

Second-site desensitization. Adaptation or sensitivity losses that act on the outputs of second-site cone-opponent and achromatic mechanisms, and thus on the combined cone signals processed by each mechanism.

Test method. A method in which the sensitivity for detecting or discriminating a target is measured as a function of some target parameter, such as wavelength, size, or temporal frequency.

Threshold surface or contour. Synonyms for detection surface or contour.

Unique hues. Hues that appear perceptually unmixed, such as unique blue and unique yellow (which appear neither red nor green).

Unipolar mechanism. A mechanism that responds to only one pole of bipolar cone-opponent excursions, thought to be produced by half-wave rectification of bipolar signals.

Univariant mechanism. A mechanism, in which the output varies unidimensionally, irrespective of the characteristics of its inputs.

von Bezold spreading. See Color assimilation.

von Kries adaptation. Reciprocal sensitivity adjustment in response to changing light levels assumed to occur independently within each of the three cone mechanisms.

Weber's law. $\Delta I/I = \text{constant}$. The sensitivity to increments (ΔI) is inversely proportional to the adaptation level (I).

11.2 INTRODUCTION

The first stage of color vision is now well understood (see Chap. 10). When presented in the same context under photopic conditions, pairs of lights that produce the same excitations in the long-, middle-, and short-wavelength-sensitive (L-, M-, and S-) cones match each other exactly in appearance. Moreover, this match survives changes in context and changes in adaptation, provided that the changes are applied equally to both lights. Crucially, however, while the match survives such manipulations, the shared appearance of the lights does not. Substantial shifts in color appearance can be caused both by changes in context and by changes in chromatic adaptation. The identity of lights matched in this way reflects univariance at the cone photoreceptor level, whereas their changed appearance reflects the complex activity of postreceptoral mechanisms acting on the outputs of the cone photoreceptors. Figure 1 shows examples of how color contrast and color assimilation can affect the color appearance of pairs of lights that are physically identical.

In addition to affecting color appearance, postreceptoral mechanisms play a major role in determining the discriminability of color stimuli. Indeed, measurements of color thresholds (detection and discrimination) are critical in guiding models of postreceptoral mechanisms. Models of color discrimination are also important in industrial applications, for instance, in the specification tolerances for color reproduction (see subsection “CIE Uniform Color Spaces” in Sec. 10.6, Chap. 10).

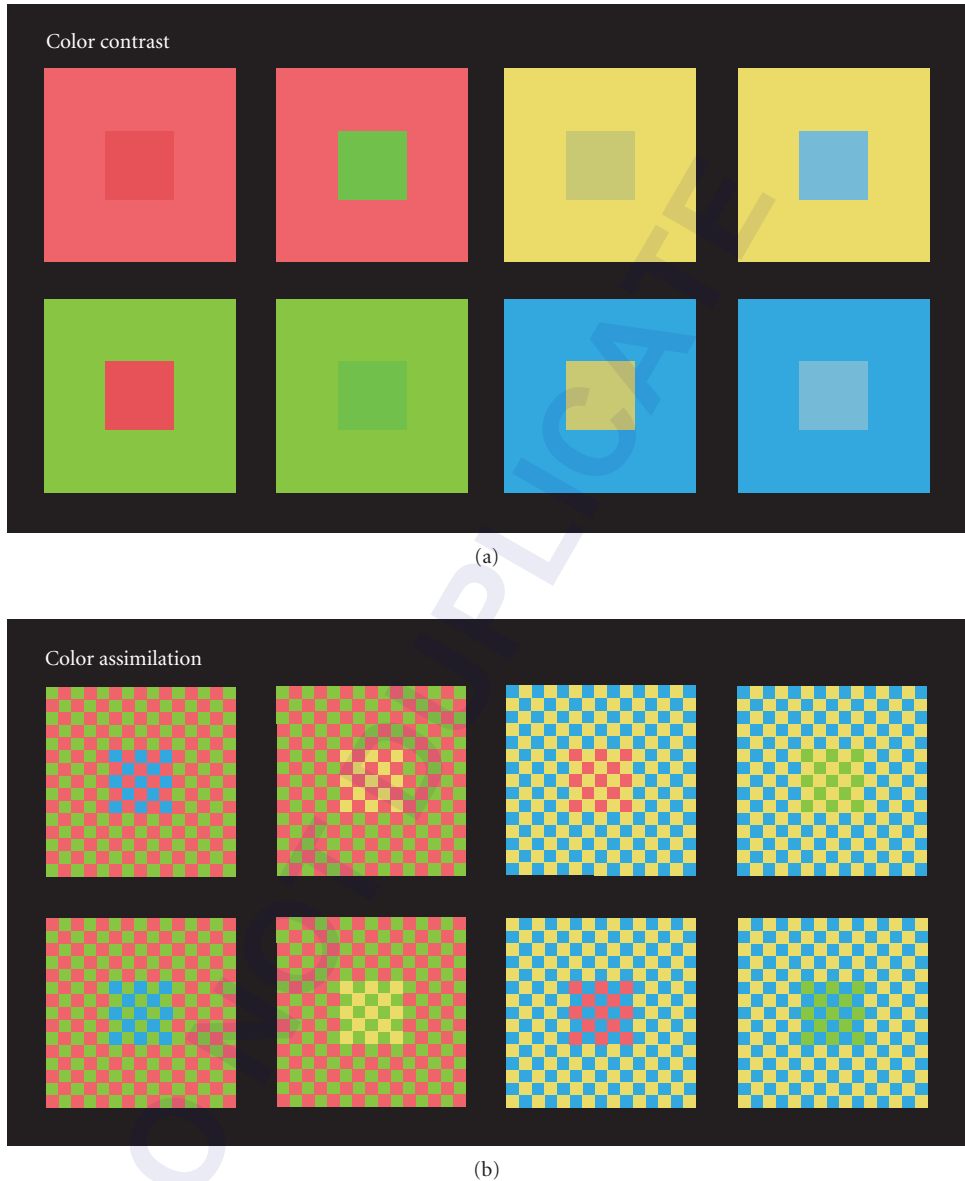


FIGURE 1 (a) Color contrast: The pairs of smaller squares in each of the four vertical columns are physically the same, but their color appearances are very different. The differences arise because of the surrounding areas, which induce complementary color changes in the appearance of the central squares.⁴⁷³ Comparable examples of color contrast have been produced by Akiyoshi Kitaoka,⁴⁷⁴ which he attributed to Kasumi Sakai.⁴⁷⁵ (b) Color assimilation or the von Bezold spreading effect:⁴⁷⁶ The tiny squares that make up the checkerboard patterns in each of the four columns are identical, except in the square central areas. In those central areas, one of the checkerboard colors has been replaced by a third color. The replacement color is the same in the upper and lower patterns, but the colors of the checkers that it replaces are different. The result is that the replacement color is surrounded by a different color in the upper and lower patterns. Although the replacement color is physically the same in each column, it appears different because of the color of the immediately surrounding squares. Unlike color contrast, the apparent color change is toward that of the surrounding squares.

The Mechanistic Approach

This chapter is about color vision after the photoreceptors. In the development, we adopt a mechanistic approach. The idea is to model color vision as a series of stages that act on the responses of the cones. Within the mechanistic approach, the central questions are: how many stages are needed, what are the properties of the mechanisms at each stage, and how are the mechanisms' outputs linked to measured performance? We focus on psychophysical (perceptual) data. Nonetheless, we are guided in many instances by physiological and anatomical considerations. For reviews of color physiology and anatomy, see, for example, Gegenfurtner and Kiper,¹ Lennie and Movshon², and Solomon and Lennie.³ A useful online resource is Webvision at <http://webvision.med.utah.edu/>.

The distinction between color encoded at the photoreceptors and color encoded by postreceptoral mechanisms was anticipated by two theories that have dominated color vision research since the late nineteenth century. First, in the Young-Helmholtz trichromatic theory,^{4,5} color vision is assumed to depend on the univariant responses of the three fundamental color mechanisms (see Chap. 10). Color vision is therefore trichromatic. Trichromacy allows us to predict which mixtures of lights *match*, but it does not address how those matches *appear*, nor the discriminability or similarity of stimuli that do not match.

Second, in Hering's^{6,7} opponent colors theory, an early attempt was made to explain some of the phenomenological aspects of color appearance, and, in particular, the observation that under normal viewing conditions some combinations of colors, such as reddish-blue, reddish-yellow, and greenish-yellow, are perceived together, but others, such as reddish-green or yellowish-blue, are not. This idea is illustrated in Fig. 2. Hering proposed that color appearance arises from the action of three signed mechanisms that represent opposing sensations of red versus green, blue versus yellow, and light versus dark.^{6,7} A consequence of this idea is that opposing or opponent pairs of sensations are exclusive, since they cannot both be simultaneously encoded. In this chapter, we will use the term "color-appearance mechanisms" to refer to model-constructs designed to account for the appearance of stimuli, and in particular the opponent nature of color appearance.

Early attempts to reconcile trichromacy with the opponent phenomenology of color appearance suggested that the color-appearance mechanisms reflect a postreceptoral stage (or "zone") of color processing that acts upon the outputs of the three Young-Helmholtz cone mechanisms. Modern versions of the two-stage theory explicitly incorporate the cone's characteristics as a first stage as well as

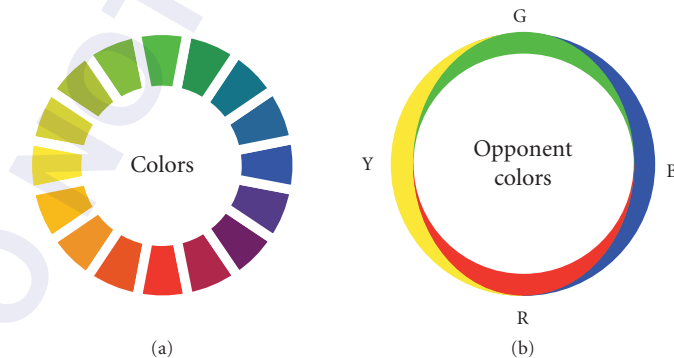


FIGURE 2 Hering's opponent-colors diagram. A diagrammatic representation of opponent-colors theory. The ring on the left (a) shows a range of colors changing in small steps from green at the top clockwise to blue, red, yellow, and back to green. The ring on the right (b) shows the hypothetical contributions of each of the color-opponent pairs [red (R) vs. green (G), and blue (B) vs. yellow (Y)] to the appearance of the corresponding colors in (a). In accordance with opponent-colors theory, the opposed pairs of colors are mutually exclusive. (Redrawn from Plate 1 of Ref. 7).

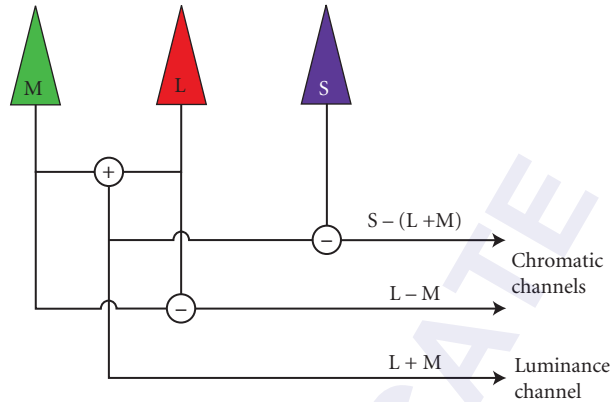


FIGURE 3 Model of the early postreceptoral stages of the visual system. The signals from the three cone types, S, M, and L, are combined to produce an achromatic or luminance channel, $L+M$, and two cone-opponent channels, $L-M$ and $S-(L+M)$. Note that there is assumed to be no S-cone input to the luminance channel. (Based on Fig. 7.3 of Ref. 15).

a second stage at which signals from the separate cone classes interact (e.g., Refs. 8–14). A familiar version of the two-zone model from Boynton¹⁵ with chromatic, $L-M$ and $S-(L+M)$, and achromatic, $L+M$, postreceptoral mechanisms is shown in Fig. 3.

Interestingly, the particulars of many modern two-stage models were formulated not to account for color appearance, but rather to explain threshold measurements of the detection and discrimination of visual stimuli. As Fig. 3 indicates, the opponent mechanisms in these models take on a simple form, represented as elementary combinations of the outputs of the three cone classes. We refer to opponent mechanisms that are postulated to explain threshold data as “color-discrimination mechanisms,” to distinguish them from color-appearance mechanisms postulated to explain appearance phenomena. Here we use the omnibus term color-discrimination to refer both to detection (where a stimulus is discriminated from a uniform background) and discrimination (where two stimuli, each different from a background, are discriminated from each other.)

The distinction between color-appearance and color-discrimination mechanisms is important, both conceptually and in practice. It is important conceptually because there is no a priori reason why data from the two types of experiments (appearance and threshold) need be mediated by the same stages of visual processing. Indeed, as we will see below, the theory that links measured performance to mechanism properties is quite different in the two cases. The distinction is important in practice because the mechanism properties derived from appearance and discrimination data are not currently well reconciled.

The discrepancy between color-discrimination and color-appearance mechanisms has been commented on by recent authors,^{11,14,16–22} but the discrepancy is implicit in early versions of the three-stage Müller zone theories,^{23,24} Judd’s version of which²⁴ was discussed again some years later in a more modern context (see Fig. 6. of Ref. 25). It is remarkable that models with separate opponent stages for the two types of data were proposed well before the first physiological observation of cone opponency in fish²⁶ and primate.²⁷

Figure 4 illustrates a modern version of Judd’s three-stage Müller zone theory, which is described in more detail in the subsection “Three Stage Zone Models” in Sec. 11.6. The figure shows the spectral sensitivities of each of the three stages. The spectral sensitivities of Stage 1 correspond to the cone spectral sensitivities of Stockman and Sharpe,²⁸ those of Stage 2 to the spectral sensitivities of

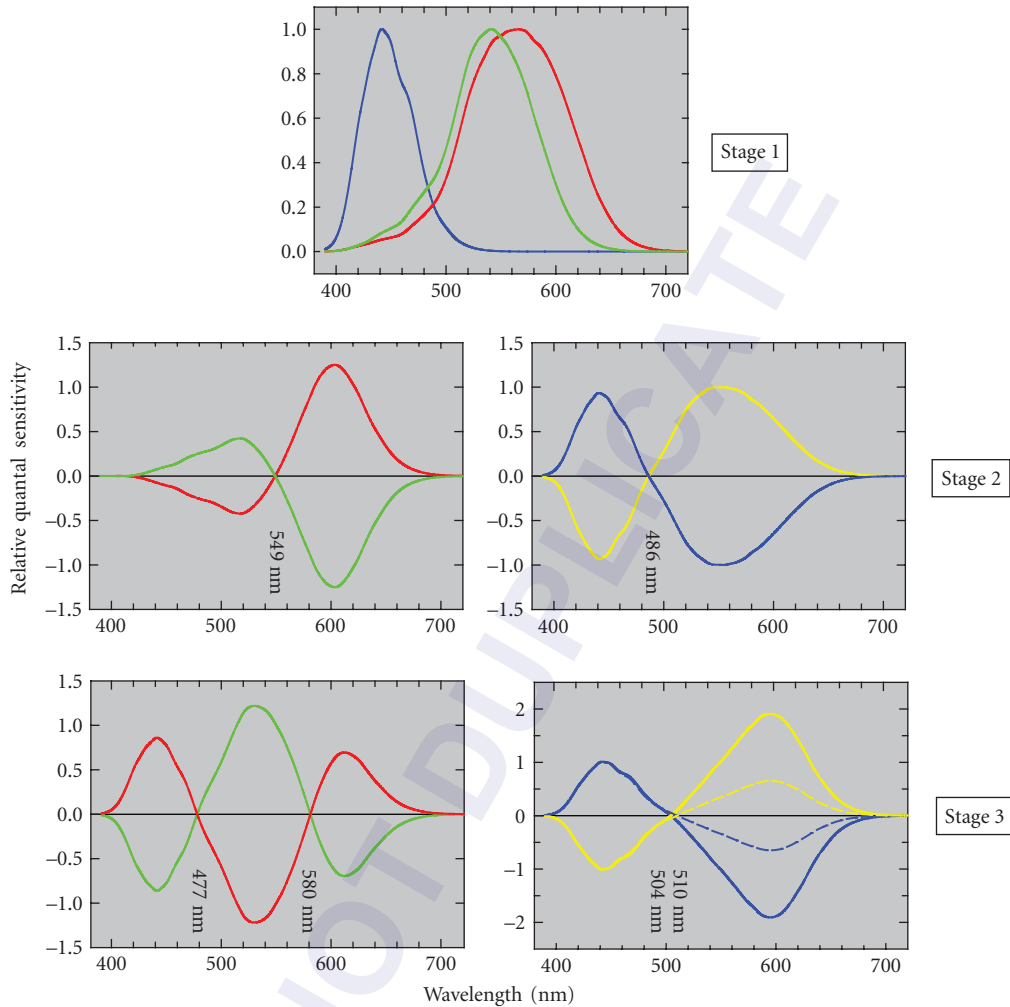


FIGURE 4 Version of the three-stage Müller zone model with updated spectral sensitivities. The panels shows the assumed spectral sensitivities of the color mechanisms at Stages 1 (upper panel), 2 (middle panels), and 3 (lower panels). *Stage 1*: L- (red line), M- (green line), and S- (blue line) cone fundamental spectral sensitivities.²⁸ *Stage 2*: L-M (red line), M-L (green line), S-(L+M) (blue line), and (L+M)-S (yellow line) cone-opponent mechanism spectral sensitivities. *Stage 3*: R/G (red line), G/R (green line), B/Y (blue line), Y/B (yellow line) color-opponent spectral sensitivities. Our derivation of the cone-opponent and color-opponent spectral sensitivities is described in the subsection “Three-Stage Zone Models” in Sec. 11.6. The dashed lines in the lower right panel are versions of the B/Y and Y/B color-opponent spectral sensitivities adjusted so that the Y and B spectral sensitivity poles are equal in area. The wavelengths of the zero crossings of the Stage 2 and Stage 3 mechanisms are given in the figure. The spectral sensitivities of the achromatic mechanisms have been omitted.

color-discrimination mechanisms as suggested by threshold data, and those of Stage 3 to the spectral sensitivities of color-appearance mechanisms as suggested by appearance data.

Figure 4 sets the scene for this chapter, in which we will review the theory and data that allow derivation of the properties of color-discrimination and color-appearance mechanisms, and discuss the relation between the two. According to some commentators, one of the unsolved mysteries of color vision is how best to understand the relation between mechanisms referred to as Stages 2 and 3 of Fig. 4.

Nomenclature

One unnecessary complication in the literature is that discrimination and appearance mechanisms are frequently described using the same names. Thus, the terms red-green (R/G), blue-yellow (B/Y), and luminance are often used to describe both types of mechanisms. We will attempt in this chapter to maintain a distinct nomenclature for distinct mechanisms.

It is now accepted that cones should be referred to as long-, middle-, and short-wavelength-sensitive (L-, M-, and S-), rather than red, green, and blue, because the color descriptions correspond neither to the wavelengths of peak cone sensitivity nor to the color sensations elicited by the excitation of single cones.³⁰ However, it is equally misleading to use color names to refer to color-discrimination mechanisms. Stimulation of just one or other side of such a mechanism does not necessarily give rise to a simple color sensation. Indeed, current models of opponent color-discrimination mechanisms have the property that modulating each in isolation around an achromatic background produces in one case a red/magenta to cyan color variation and in the other a purple to yellow/green variation.^{31,32} Consequently, the perception of blue, green, and yellow, and to a lesser extent red, requires the modulation of *both* cone-opponent discrimination mechanisms (see subsection “Color Appearance and Color Opponency” in Sec. 11.5). We therefore refer to chromatic color-discrimination mechanisms according to their predominant cone inputs: L–M and S–(L+M). Although this approach has the unfortunate consequence that it neglects to indicate smaller inputs, usually from the S-cones (see subsection “Sensitivity to Different Directions of Color Space” in Sec. 11.5.), it has the advantage of simplicity and matches standard usage in much of the literature. We refer to the nonchromatic color discrimination mechanism as L+M. Note also that this nomenclature is intended to convey the identity and sign of the predominant cone inputs to each mechanism, but not the relative weights of these inputs.

In contrast, the perception of pure or “unique” red, green, yellow, and blue is, by construction of the theory, assumed to result from the responses of a single opponent color-appearance mechanism, the response of the other mechanism or mechanisms being nulled or in equilibrium (see subsection “Opponent-Colors Theory” in Sec. 11.5). We refer to opponent color-appearance mechanisms as R/G and B/Y, according to the color percepts they are assumed to generate. We refer to the nonopponent appearance mechanism as brightness.

Guiding Principles

Behavioral measurements of color vision reflect the activity of an inherently complex neural system with multiple sites of processing that operate both in series and in parallel. Moreover, these sites are essentially nonlinear. The promise of the mechanistic approach lies in two main areas. First, in terms of developing an overall characterization of postreceptoral color vision, the hope is that it will be possible to identify broad regularities in the behavior of the system that can be understood in terms of models that postulate a small number of relatively simple mechanism constructs. Second, in terms of using psychophysics to characterize the behavior of particular neural sites, and to link behavior to physiology, the hope is that specific stimulus conditions can be identified for which the properties of the site of interest dominate the measured performance. In the conceptual limit of complexity, where a mechanistic model explicitly describes the action of every neuron in a given visual pathway, these models can, in principle, predict performance. But as a practical matter, it

remains unclear the degree to which parsimonious mechanistic models derived from behavioral measurements will succeed. Given the complexity of the underlying neural system, it is apparent that the mechanistic approach is ambitious.

In this context, several points are worth bearing in mind. First, the concept of a psychophysical mechanism will have the most utility when it can be shown to have an existence that extends beyond the particular stimulus and task conditions from which it was derived. Thus, we regard a mechanism as a theoretical construct whose usefulness depends on the range of data it can explain parsimoniously. This is a broader consideration than those sometimes used to determine the value of a mechanism.^{33,34} In Secs. 11.3 and 11.4, we review two broad mechanism concepts that satisfy this criterion: opponency and adaptation.

Second, while some psychophysical techniques may emphasize the contribution of particular stages of the system, it is, with the exception of color matching, a simplification to ignore the contributions of earlier and/or later stages. For example, an often made but usually implicit assumption is that the cone quantal absorption rates are transmitted directly to postreceptoral mechanisms, as if the photoreceptors had no role other than to pass on a linear copy of their inputs. An interesting future direction for mechanistic models is to account for interactions between different stages of processing more fully.

Finally, in spite of our concerns, we have written this chapter primarily from the point of view of psychophysics. We readily acknowledge that over the past 50 years, results from molecular genetics, anatomy, and physiology have helped propel many psychophysical theories from intelligent speculation to received wisdom, particularly in the cases of cone spectral sensitivities and early retinal visual processing. Physiology and anatomy have also provided important evidence about color coding at different levels of the early visual system (e.g., Refs. 35–37), crucial information that is hard to obtain psychophysically. Obversely, as the increasing complexity of the visual system at higher levels of processing begins to limit the utility of a mechanistic psychophysical approach, it also constrains how much can be understood from a knowledge of the properties of the component neurons in the processing chain, the majority of which have yet to be characterized.

Chapter Organization

The rest of this chapter is organized as follows: In the next two sections, we describe how mechanism concepts allow us to understand important features of both color discrimination and color appearance. Our treatment here is illustrative. We review a small selection of experimental data, and emphasize the logic that leads from the data to constraints on the corresponding models. In Sec. 11.3 we discuss discrimination data; in Sec. 11.4 we turn to appearance data. Two broad mechanistic concepts emerge from this initial review: those of color opponency and of adaptation. The initial review provides us with a basic model for postreceptoral color vision. Although this model is clearly oversimplified, it serves as the point of departure for much current work, and understanding it is crucial for making sense of the literature. The remainder of the chapter (Secs. 11.5 and 11.6) is devoted to a discussion of advanced issues that push the boundaries of the basic model.

11.3 BASICS OF COLOR-DISCRIMINATION MECHANISMS

What Is a Mechanism?

Given the central role of the mechanism concept in models of color vision, one might expect that this concept is clearly and precisely defined, with a good consensus about its meaning. Our experience, however, is that although there are precise definitions of what constitutes a mechanism,^{33,34} these are generally specific to a particular model and that the term is often used fairly loosely. In keeping with this tradition, we will proceed to discuss mechanisms through examples without attempting a rigorous definition.

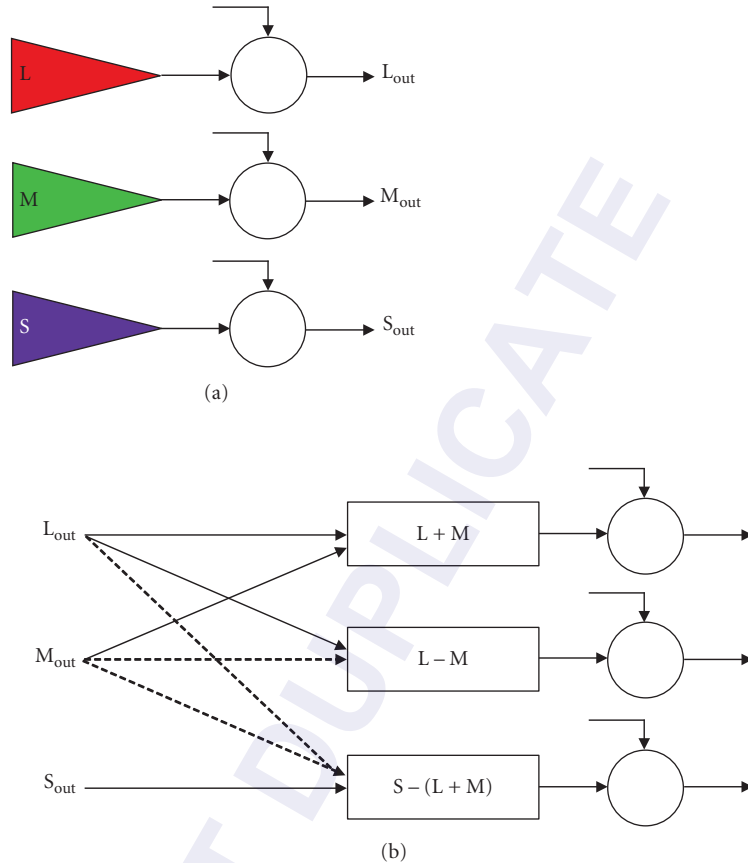


FIGURE 5 Basic mechanisms. (a) First-stage cone mechanisms: L, M, and S. The cone outputs are subject to some form of gain control (open circles), which can, in principle, be modified by signals from the same or from different cone mechanisms. (b) Second-stage color-discrimination mechanisms: $L + M$, $L - M$, and $S - (L + M)$. The inputs to these mechanisms are the adapted cone outputs from the cone mechanisms. As with the cone mechanisms, the outputs of the second-stage mechanism are subject to some gain control, which can be modified by signals from the same or from different second-stage mechanisms. Dashed arrows indicate inhibitory outputs.

Figure 5 illustrates the broad mechanism concept. Figure 5a shows the L, M, and S cones. Each of these cones is a mechanism and satisfies several key properties. First, we conceive of each cone mechanism as providing a spatiotemporal representation of the retinal image, so that the mechanism output may be regarded as a spatial array that changes over time. Second, at a single location and time, the output of each cone mechanism is univariate and conveys only a single scalar quantity. This means that the output of a single cone mechanism confounds changes in the relative spectrum of the light input with the overall intensity of the spectral distribution. Third, the relation between mechanism input and output is subject to adaptation. This is indicated in the figure by the open circle along the output pathway for each cone mechanism. Often the adaptation is characterized as a gain control, or a gain control coupled with a term that subtracts steady-state input (e.g., Refs. 38 and 39).

Key to the mechanism approach is that the behavior of the mechanism can be probed experimentally with adaptation held approximately fixed. This may be accomplished, for example, by presenting brief, relatively weak, flashes on spatially uniform backgrounds. In this case, the state of adaptation is taken to be determined by the background alone; that is, it is assumed that the weak flashes do not measurably perturb the state of adaptation and thus that the state of adaptation is independent of the flash.

The arrows pointing into each circle in Fig. 5 indicate that the signals that control the state of adaptation can be quite general, and need to be specified as part of the model. In the classical (von Kries) mechanism concept, however, these signals are assumed to arise entirely within the array of signals from the same cone mechanism.³³ For example, the gain (or adaptation) applied to the L-cone signals at a location would be taken to be determined entirely by the array of L-cone responses, and be independent of the responses of M and S cones. Adaptation that acts on the output of the cones is referred to as “first-site” adaptation.

Figure 5*b* shows a simple model of three second-stage color-discrimination mechanisms with the same basic wiring diagram as Fig. 3. These mechanisms have a similar structure to the cone mechanisms, in the sense that they are assumed to provide a spatial array of univariate responses. In addition, as with the cone mechanisms, the second-stage mechanisms can adapt. An important difference between these second-stage mechanisms and the cone mechanisms is the nature of their input. Where the cone mechanisms transduce light into a neural signal, the second-stage mechanisms take the adapted cone responses as their input. As shown in the Fig. 5, each of the postreceptoral mechanisms can be thought of as computing a weighted combination of cone signals. The three mechanisms are labeled by the manner in which they combine cone inputs. The first mechanism, (L+M), which is often referred to as the luminance mechanism, adds inputs from the L and M cones. The second mechanism, (L–M), takes a difference between L and M cone signals. And the third mechanism, S–(L+M), takes a difference between S-cone signals and summed L- and M-cone signals.

Although not necessary as part of the mechanism definition, making the input to the mechanism be a linear function of the output of the preceding stage simplifies the model, and enables it to make stronger predictions; whether the cone combinations are actually linear is an empirical question. Adaptation that acts on the outputs of the second-stage postreceptoral mechanisms is usually referred to as “second-site” adaptation or desensitization.

The cone mechanisms used in color vision models represent the action of a well-defined neural processing stage (the cone photoreceptors). The connection between postreceptoral mechanisms and neurons is not as tight. Although the cone inputs to different classes of retinal ganglion cells are similar to those used in models, these models often do not incorporate important features of real ganglion cells, such as their spatial receptive field structure (but see subsection “Multiplexing Chromatic and Achromatic Signals” in Sec. 11.5), cell-to-cell variability in cone inputs, and the fact that ganglion cells come in distinct ON and OFF varieties that each partially rectify their output (but see “Unipolar vs. Bipolar Chromatic Mechanisms” in Sec. 11.6).

Psychophysical Test and Field Methods

Broadly speaking, two principal psychophysical techniques have been used to investigate the properties of color-discrimination mechanisms.^{40–42} In the “test” (or “target”) sensitivity method, the observer’s sensitivity for detecting or discriminating a target is measured as a function of some target parameter, such as its wavelength, size, or temporal frequency. As noted above, implicit in the use of this method is the assumption that the presentation of the target does not substantially alter the properties or sensitivity of the detection mechanism, so that targets are usually presented against a background and kept near visual threshold.

In the “field” sensitivity method, the observer’s sensitivity for detecting or discriminating a target is measured as a function of some change in the adaptive state of the mechanism. Field methods complement test methods by explicitly probing how contextual signals control adaptation. On the assumption that the control of adaptation occurs solely through signals within the mechanism mediating detection, the spectral properties of that mechanism can be investigated by, for example,

superimposing the target on a steady adapting field and changing the adapting field chromaticity and/or radiance, by habituating the observer to backgrounds temporally modulated in chromaticity and/or luminance just prior to the target presentation, or by superimposing chromatic and/or luminance noise on the target.

Both test and field methods have obvious limitations. In the test method, it is difficult to ensure that a target is detected by a single mechanism when the target parameters are varied, with the result that multiple mechanisms typically mediate most sets of test sensitivity measurements. As we discuss below, assigning different portions of chromatic detection contours to different mechanisms is problematic (see subsection “Sensitivity to Different Directions of Color Space” in Sec. 11.5). In the field method, it is often easier to ensure that a target is detected by a single mechanism. However, the assumption that adaptation is controlled entirely by signals from within the mechanism mediating detection is a strong one, and interpretation of results becomes much more complicated under conditions where this assumption is not secure. More generally, without an explicit and accurate model about the properties of the mechanism, both test and field sensitivity data may be uninterpretable.

How Test Measurements Imply Opponency

In the next three sections, we introduce color-discrimination mechanisms. These are characterized by opponent recombination of cone signals at a second site (as shown in Fig. 5*b*), by adaptation in cone specific pathways (first-site adaptation, Fig. 5*a*), and by adaptation after the opponent recombination (second-site adaptation, Fig. 5*b*).

We begin with evidence from test sensitivity measurements that some sort of opponent recombination occurs. In the canonical test sensitivity experiment, a small test stimulus is presented against a uniform background. We denote the L-, M-, and S-cone excitations of the background by L_b , M_b , and S_b . If we fix the background, the test stimulus may be characterized by how much the cone excitations it produces deviate from the background. Denote these deviations by ΔL , ΔM , and ΔS , so that the overall cone excitations of the test are given by $L = L_b + \Delta L$, $M = M_b + \Delta M$, and $S = S_b + \Delta S$. Note that the deviations ΔL , ΔM , and ΔS may be positive or negative.

For any background, we can consider a parametric family of test stimuli whose L-, M-, and S-cone deviations are in the same proportion. That is, we can define a test *color direction* by a triplet of deviations ΔL_d , ΔM_d , and ΔS_d , normalized so that $\sqrt{\Delta L_d^2 + \Delta M_d^2 + \Delta S_d^2} = 1$. Test stimuli that share the same color direction have the form $\Delta L = c\Delta L_d$, $\Delta M = c\Delta M_d$, and $\Delta S = c\Delta S_d$. We refer to the constant c as the *intensity* of the test stimulus along the given color direction. Figure 6*a* illustrates the test color direction and intensity concept for two vectors in the ΔL , ΔM plane.

If the background and color direction of a test stimulus are held fixed, an experimenter can vary the test intensity and determine the psychophysical threshold for detection. This is the lowest intensity at which the observer can just see the test, and its value can be found experimentally using a variety of procedures.⁴³ The experiment can then be repeated for different choices of color direction, and the threshold intensity can be determined for each one.

Figure 6*b* plots one way in which the data from a threshold experiment might come out. For simplicity, we assume that ΔS_d was set to zero, and show data only in the ΔL , ΔM plane. Each point in the plot represents the results of a threshold measurement for one color direction. The set of threshold points together trace out a “detection contour” (or “threshold contour”), since each point on the contour leads to an equal level of detection performance (the “threshold”). The two points that lie directly on the ΔL axis represent threshold for increments and decrements that drive only the L cones and leave the M cones silent; that is, these targets produce zero change in the M cones. Similarly two points that lie directly on the ΔM axis represent increments and decrements that produce only changes in the M cones. The other points on the contour show thresholds obtained when L- and M-cone signals are covaried in various ratios.

Understanding how threshold contours inform us about color-discrimination mechanisms is a key idea in the theory we present in this chapter. It is useful to begin by asking how the contour would come out if the visual system had a single color-discrimination mechanism consisting of, say,

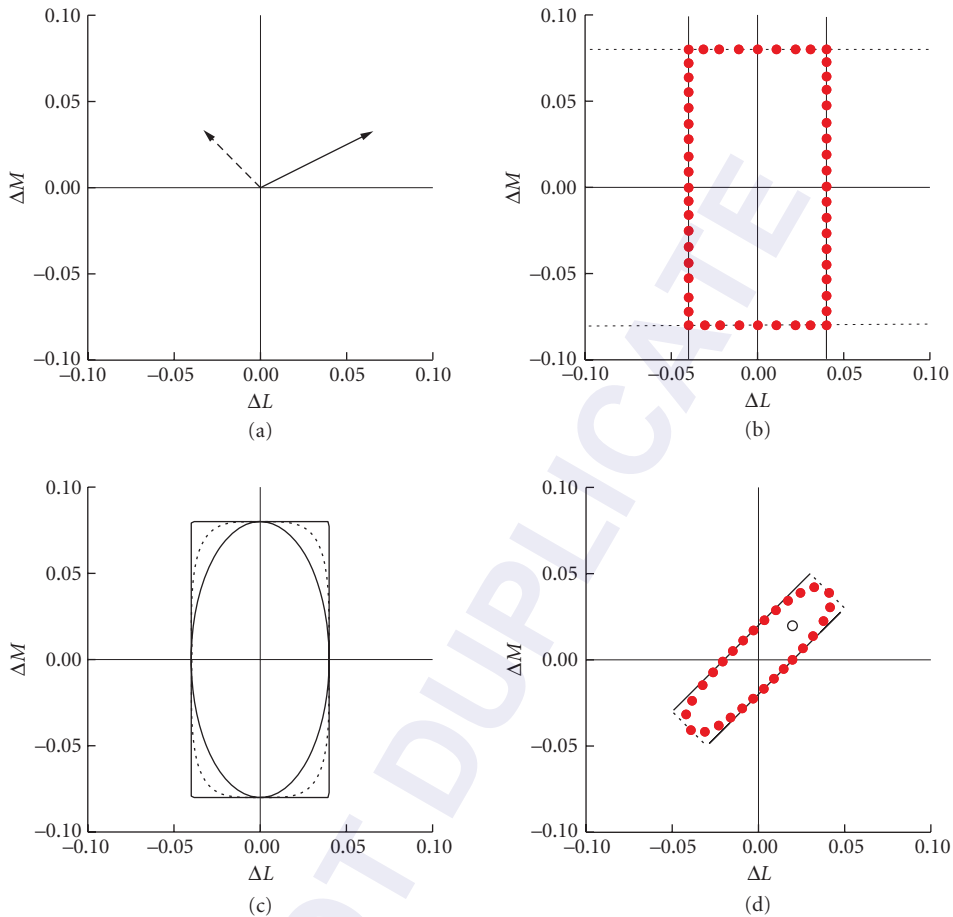


FIGURE 6 Basic thresholds. (a) Two vectors plotted in the ΔL , ΔM plane that represent two lights of different color directions and different intensities. (b) Incremental and decremental thresholds determined by an L-cone mechanism (vertical black lines) or by an M-cone mechanism (black, horizontal dotted lines). The points define the contour expected if threshold is determined independently and with no interactions between L- and M-cone detection mechanisms. (c) The joint detection contours when the L- and M-cone detection mechanisms show different degrees of summation. The inner, middle, and outer contours are for summation exponents of $k = 2, 4$, and 1000 , respectively. (d) Idealized detection contours (red solid points) for thresholds determined by chromatic L–M (solid lines at 45°) and achromatic L+M (dotted lines at -45°) mechanisms. The single open circle shows the threshold expected in the $\Delta L_d = \Delta M_d$ direction if the threshold was determined by the cone mechanisms (i.e., by the thresholds along the axes $\Delta L_d = 0$ and $\Delta M_d = 0$).

only the L cones. In this case, only the ΔL component of the test modulation would affect performance, and the threshold contour would consist of two lines, both parallel to the ΔM axis. One line would represent threshold for incremental tests with a positive ΔL value, while the other would represent threshold for decremental tests with a negative ΔL value. This hypothetical detection contour is shown as two solid vertical lines in Fig. 6b. If there were no other mechanisms, the threshold contour would simply continue along the extensions of the solid vertical lines. The contour would not be closed, because modulations along the M-cone isolating direction produce ΔL values of zero and thus are invisible to the L-cone mechanism.

We can also consider a visual system with only an M-cone discrimination mechanism. By reasoning analogous to that used for the L-cone case above, the threshold contour for this system would consist of the two horizontal dotted lines also shown in the figure.

Finally, we can consider a system with independent L- and M-cone mechanisms and a threshold determined when either its ΔL or ΔM component reached the threshold for the corresponding mechanism. This would lead to a closed rectangular detection contour formed by the intersecting segments of the solid vertical and dotted horizontal lines shown in Fig. 6*b*. The threshold data plotted as solid red circles in the panel correspond to this simple model.

When signals from both L- and M-cone mechanisms mediate detection, the measured detection contour would be expected not to reach the corners of the rectangular contour shown in Fig. 6*b*. This is because even if the L- and M-cone mechanisms are completely independent, detection is necessarily probabilistic, and in probabilistic terms the likelihood that L, M, or both L and M together will signal the test is greater than the likelihood of either cone mechanism doing so alone even when the mechanisms are independent, especially when L and M are both near threshold. This “probability summation” will reduce thresholds in the corners where L and M signals are of similar detectability, rounding them. A number of simple models of how mechanisms outputs might combine in the presence of neural noise predict such rounding, with the exact shape of the predicted contour varying across models.

A convenient parametric form for a broad class of summation models is that test threshold is reached when the quantity $\Delta = \sqrt[k]{\sum_{i=1}^N |\Delta_i|^k}$ reaches some criterion value. In this expression, N is the number of mechanism whose outputs are being combined, Δ_i is the output of the i th mechanism, and the exponent k determines the form of summation. When $k = 2$, the quantity Δ represents the Euclidean vector length of the mechanism outputs. When $k \rightarrow \infty$, Δ represents the output of the mechanism whose output is greatest.

Figure 6*c* shows expected threshold contours for the output of the L and M mechanisms for three values of k . The outer rectangular contour (solid lines) shows the contour for $k = 1000$. Here there is effectively no summation, and the contour has the rectangular shape shown in Fig. 6*b*. The inner contour (solid line) shows the case for $k = 2$. Here the contour is an ellipse. The middle contour (dotted line) shows the case for $k = 4$. Determining the value of k that best fits experimental data is a topic of current interest, as the answer turns out to have important implications for how to reason from sensitivity contours to mechanism properties (see subsection “Sensitivity to Different Directions of Color Space” in Sec. 11.5). Here, however, the key point is that we expect actual experimental data to show more rounded threshold contours than the one depicted in Fig. 6*b*. Under conditions where sensitivity is determined by the output of the L- and M-cone mechanisms, the expected shape of the threshold contour is a closed form whose major axes are aligned with the ΔL and ΔM axes.

Figure 6*d* shows an idealized representation of how the data actually come out when thresholds are measured against a neutral background (see also Fig. 20). Actual data deviate clearly from the predictions shown in Fig. 6*c*, which were based on the assumption that the color-discrimination mechanisms are the responses of the L and M cones. Instead, the data lie near an elongated ellipsoid whose axes are rotated intermediate to the ΔL and ΔM axes. The deviations from the predictions of the Fig. 6*c* are large and robust. In particular, note the location of the open circle in the figure. This circle shows an upper bound on threshold in the $\Delta L_d = \Delta M_d$ color direction; if the color-discrimination mechanisms were the L and M cones and they operated without summation, then threshold would be given by the open circle. If there was also summation between L- and M-cone mechanisms, then threshold would lie between the open circle and the origin, with the exact location depending on the degree of summation. Thresholds in this direction vastly exceed the bound shown by the open circle. This observation demonstrates unequivocally that the outputs of different cone mechanisms must be postreceptorally recombined.

A natural interpretation of the threshold contour shown in Fig. 6*d* is indicated by the rotated rectangle shown on the plot. The parallel solid lines represent the threshold contour that would be observed if detection were mediated by a single mechanism that computed as its output the difference between the L- and M-cone excitations to the test, and if threshold depended on $|\Delta L - \Delta M|$ reaching some criterion threshold level. Similarly, the dotted lines represent the threshold contours of a mechanism that sums the L- and M-cone excitations of the test. The idealized threshold data

shown is thus consistent with the first two color-discrimination mechanisms illustrated in Fig. 5b, with some amount of summation between mechanism outputs accounting for the quasi-ellipsoidal shape of the overall contour.

Logic similar to that shown, applied to data where ΔS is varied, leads us to postulate a third color-discrimination mechanism whose output is given by a weighted opposition of ΔS on the one hand and ΔL and ΔM on the other. Figure 21, later, shows detection contours in the equiluminant plane partially determined by the S-(L+M) mechanism.

Although it is straightforward to predict detection contours if the spectral properties of the underlying detection mechanisms and the interactions between them are known, it is harder, as we shall see below, to infer unequivocally the exact mechanism properties (e.g., the exact weights on signals from each cone class) from the contours.

First-Site Adaptation

Weber's law and contrast coding The measurements described in the previous section assess sensitivity under conditions where the state of adaptation was held fixed. We now turn to how input signals from the cones depend on the conditioning or adapting background. The key idea here is that the cone excitations are converted to a contrast representation. Here we use L cones as an example and take the L-cone contrast of the test, C_L , to be given by $C_L = \Delta L/L_b$, where ΔL is (as above) the difference between the L-cone excitations produced by the test and background, and L_b is the L-cone excitation produced by the background. Similar expressions apply for the M and S cones.

The conversion of raw cone excitations (photoisomerizations) to a contrast code implies an important adaptive function. The range of cone excitation rates encountered in the environment can be greater than 10^6 . This range greatly exceeds the limitations imposed by the individual neurons in the visual pathway, many of which have dynamic ranges of no more than about 10^2 from the level of noise—their spontaneous firing rate in the absence of light stimulation—to their response ceiling.^{44,45} If there were no adaptation, the cone visual response would often convey no useful information, either because it was too small to rise above the noise at the low end of the stimulus range or because it was saturated at the maximum response level and thus unable to signal changes. For adaptation to protect the cone-mediated visual system from overload as the light level increases, the primary mechanisms of sensitivity regulation are likely to be early in the visual pathways, most likely in the cone photoreceptors themselves.^{46,48} Indeed, the molecular mechanisms of adaptation acting within the photoreceptor are now fairly well understood.^{49–52}

The idea that signals leaving the cones are effectively converted to a contrast code makes a specific prediction about how thresholds should depend on the background. Suppose we consider test stimuli that only stimulate the L cones, and that we manipulate the L-cone component of the background, L_b . Because of the cone-specific adaptation, increasing L_b will produce a proportional decrease in the contrast signal fed to the postreceptoral detection mechanisms. On the assumption that the noise limiting discrimination behavior remains constant across the change of background (see “Sites of Limiting Noise” in Sec. 11.3), this decrease means that the differential signal, ΔL , required to bring a test stimulus to threshold will increase in proportion to L_b . This sort of behavior is often described as obedience to Weber's law; predictions of Weber's law models are shown by the blue line in Fig. 7a. Figure 7a shows a log-log plot of increment threshold ΔL against background component L_b . The predicted thresholds increase with a slope of unity on the log-log plot.

A feature of contrast-coding predictions that is clearly not realistic is that the threshold ΔL should decrease toward zero in complete darkness. A generalized and more realistic form of contrast coding is given by $C_L = \Delta L/(L_b + L_o)$, where L_o is a constant. The dashed line in Fig. 7a shows a prediction of this form, which has the feature that thresholds approach a constant nonzero value as L_b decreases below L_o . The value of L_o may be thought of as the hypothetical excitation level—sometimes called a “dark light”—produced internally within the visual system in the absence of light that limits performance at low background levels.

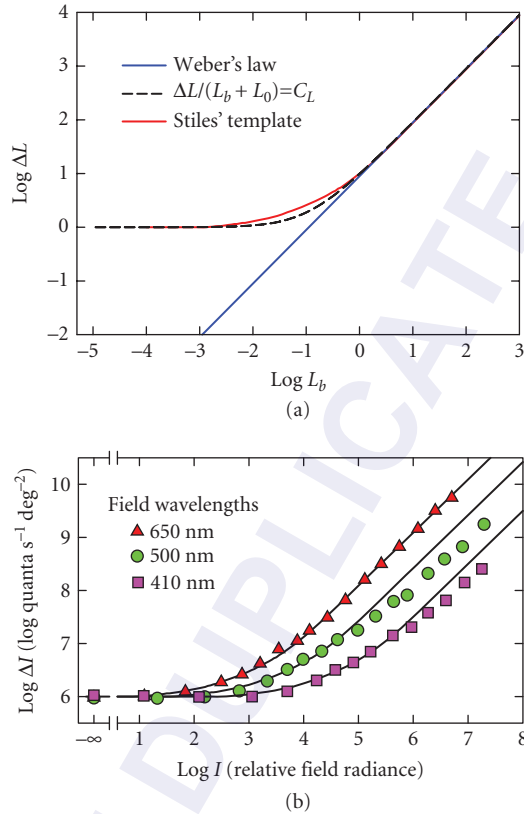


FIGURE 7 Weber's law and increment thresholds. (a) Predictions of three first-site adaptation models for test lights detected by an L-cone mechanism as a function of background intensity. (b) Increment threshold data attributed to the L-cone mechanism (π_5 , see subsection "Stiles' π Mechanisms" in Sec. 11.5) replotted from Fig. 2 (Observer: EP) of Sigel and Pugh⁵⁴ obtained using a 200-ms duration, 667-nm target presented on background fields of 650 nm (red triangles), 500 nm (green circles), and 410 nm (purple squares). The solid lines aligned with each data set at low background radiances are Stiles' standard template [Table 1(7.4.3) of Ref. 53]. The horizontal position of the data is arbitrary.

Stiles' template shape [see Table 1(7.4.3) of Ref. 53], shown as the red line in Fig. 7a, has a form similar to that predicted by generalized contrast coding (dashed line). This template was derived to account for increment threshold data obtained under a variety of conditions (see subsections "Sensitivity to Spectral Lights" and "Stiles' π Mechanisms" in Sec. 11.5). Figure 7b shows increment threshold data measured by Sigel and Pugh,⁵⁴ for three different background wavelengths under conditions where the L-cone mechanism is thought to dominate detection. These data fall approximately along Stiles' template, although deviations are clear at high background radiances for two of the three spectral backgrounds shown. The deviations are consistent with contributions of postreceptoral mechanisms to detection performance.^{54,55} Data like these speak to the difficulties of distinguishing the contributions of different mechanisms from even the simplest data. Stiles, for example, accounted for the same sort of deviations by proposing an additional cone mechanism.³³

The signals reaching the second site Although increment threshold data suggest a roughly Weber-type gain control, they do not by themselves determine whether the signal transmitted to the second site is of the contrast form $\Delta L/(L_b + L_o)$, as we have written above, or of the simpler form $L/(L_b + L_o)$, where to recap $L = \Delta L + L_b$. We cannot, in other words, determine from detection thresholds alone whether the adapted signal generated by the background is partially or wholly subtracted from the signal transmitted to the postreceptoral mechanisms.

To understand the ambiguity, define the gain g as $g = 1/(L_b + L_o)$. In a gain control only model (with no background subtraction), the signal transmitted from the L cones in response to the background plus test would be $L' = gL$, where L represents the combined cone excitations to the background plus test. Similarly, under this model, the gain adjusted response to the background alone would be $L'_b = gL_b$. If we assume that threshold requires the difference in response to background plus test on the one hand and background alone on the other to reach some criterion level, then the threshold will depend on $(L' - L'_b) = g(L - L_b) = (L - L_b)/(L_b + L_o) = \Delta L/(L_b + L_o) = C$, which is the contrast form. Thus, predictions about detection threshold are independent of whether a constant is subtracted from the signals leaving the cones.

To deduce the need for a subtractive term (i.e., to differentiate contrast-like coding from gain control alone) from threshold data requires discrimination experiments, in which threshold is measured not just for detecting a test against the background, but also for detecting the test (often referred to as a probe in this context) presented against an incremental flash, both of which are presented against a steady background. In these experiments, the threshold intensity for the probe depends on the flash intensity, and changes in this dependence with the background may be used to infer the nature of the adaptive processes. We will not review these experiments or their analysis here; other papers provide detailed descriptions (e.g., Refs. 56–59). The conclusion drawn by these authors is that to account for steady-state adaptation a subtractive term is required in addition to multiplicative gain control.

Second-Site Adaptation

The final piece of the basic model of color-discrimination mechanisms is a second stage of adaptation that modifies the output of the second-stage mechanisms (Fig. 5*b*). First-site adaptation is postulated to remove the effects of steady, uniform backgrounds by converting signals to a contrast representation; second-site adaptation is postulated to be driven by contrast signals in the background (see subsection “Field Sensitivities” in Sec. 11.5). Before proceeding, however, we should add two caveats. First, even with uniform, steady fields, first-site adaptation is usually incomplete. Were it complete, detection thresholds expressed as contrasts would be independent of background chromaticity, which they clearly are not (see Fig. 8 and subsection “Field Sensitivities” in Sec. 11.5). Second, first-site adaptation is not instantaneous. If it were, then in the extreme case of complete first-site adaptation being local to the photoreceptor as well as being instantaneous, neither local nor global changes in contrast would be transmitted to the second-site. In the less extreme case of complete, instantaneous adaptation being more spatially extended, then global changes in contrast would not be transmitted—with the result that large uniform field or Ganzfeld flicker should be invisible, which is not the case.⁶⁰ To understand the basic logic that connects threshold experiments to second-site contrast adaptation or desensitization, we first consider the effects of steady backgrounds and then the effects of habituation.

Second-site desensitization by steady fields Because first-site adaptation is incomplete, steady chromatic backgrounds can desensitize second-site mechanisms. The effect of this type of desensitization on detection contours is illustrated in Fig. 8. The colored circles in Fig. 8*a* show segments of detection contours in the $\Delta L/L, \Delta M/M$ plane of cone contrast space for two background chromaticities. The contours have a slope of 1, which is consistent with detection by a chromatic L–M mechanism with equal and opposite L- and M-cone contrast weights. If the background chromaticity is changed and adaptation at the first site follows Weber’s law, then the contours plotted in cone contrast units should not change. What happens in practice, when for example, a field is changed

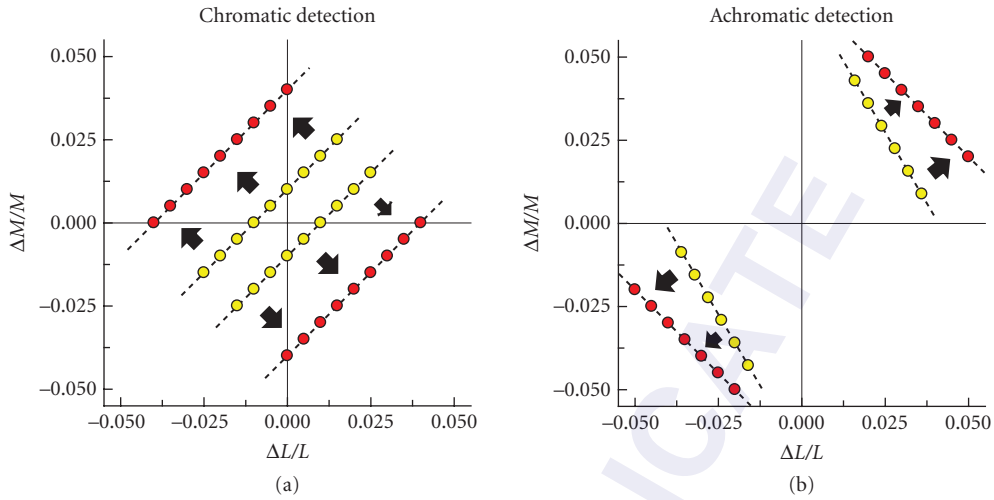


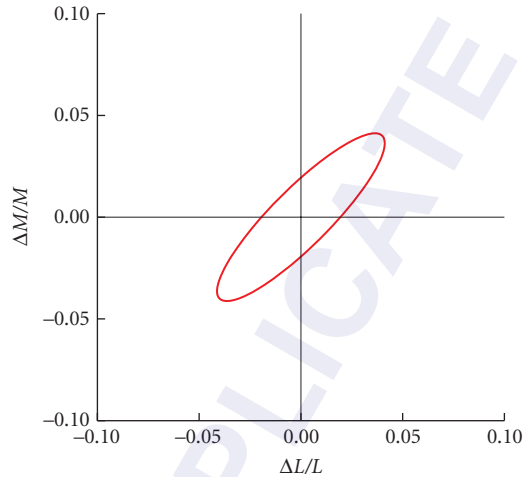
FIGURE 8 Changes in contrast threshold caused by second-site adaptation to steady fields. (a) Hypothetical contours for detection mediated by the L–M mechanism plotted in cone contrast space before (yellow circles) and after (red circles) second-site desensitization (caused by, e.g., a change in background chromaticity from yellow to red). The L–M mechanism is assumed to have equal and opposite L-cone contrast weights at its input, and this equality is assumed to be unaffected by the desensitization. (b) Hypothetical detection contours for detection by the L+M mechanism before (yellow circles) and after (red circles) a change in background chromaticity from yellow to red. Initially, the L+M mechanism is assumed to have an L-cone contrast weight 1.7 times greater than the M-cone contrast weight. Changing the field to red suppresses the L-cone contrast weight, causing a rotation of the contour, as indicated by the arrows and the red circles.⁶³

in chromaticity from spectral yellow to red, is that the contours move outward,⁶¹ as indicated by the arrows and the red circles in Fig. 8a. (Actual data of this type are shown in Fig. 28.)

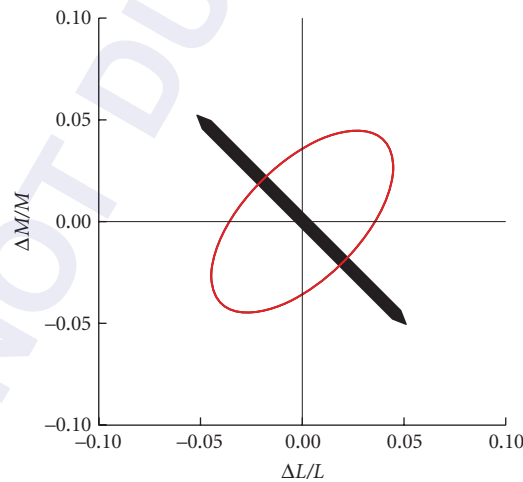
The constant slope of the L–M detection contours with changing background chromaticity is consistent with Weber’s law operating at the first site. The loss of sensitivity in excess of Weber’s law is most naturally interpreted as a second-site desensitization that depends on background chromaticity and acts on the joint L–M signals. Experimental evidence also supports the idea that first-site adaptation is in the Weber regime under conditions where supra-Weber desensitization is observed.⁶²

The effect of second-site chromatic adaptation on the L+M luminance mechanism is different. The yellow circles show partial detection contours with a slope of -1.7 , which is consistent with detection by an achromatic L+M mechanism with an L-cone contrast weight 1.7 times greater than the M-cone contrast weight (this choice is arbitrary, but, in general, the L-cone weight is found to be greater than the M-cone weight, see subsection “Luminance” in Sec. 11.5). If the background chromaticity is again changed from spectral yellow to red, the detection contours rotate, as indicated by the arrows and the yellow circles.⁶³ In this case, the rotation is due to a suppression by the long-wavelength field of the L-cone contribution to luminance detection relative to the M-cone contribution.^{63,64} Given first-site adaptation that follows Weber’s law independently in the L- and the M-cones, first-site adaptation should affect L- and M-cone inputs equally in all postreceptoral mechanisms. The fact that the L–M measurements show an equal effect on the L- and M-cone contrast inputs, whereas the L+M measurements show a differential effect, is most easily explained by positing adaptation effects at a second-site. The idea that the cone inputs to luminance and chromatic mechanisms undergo adaptation that is specific to each type of mechanism was framed by Ahn and MacLeod.⁶⁵ The same idea is also explicit in the earlier work of Stromeyer, Cole, and Kronauer.^{61,63} The effects of steady fields on the L+M mechanism are discussed further in subsection “Achromatic Direction and Chromatic Adaptation” in Sec. 11.5.

Second-site habituation The effects of steady backgrounds at the second site are comparatively small because the background signals are attenuated by adaptation at the first site. Second-site effects can be enhanced by temporally modulating the chromaticity or luminance of the background as in an habituation experiment. Figure 9a shows an idealized threshold contour for detection thresholds



(a)



(b)

FIGURE 9 Habituation predictions. (a) Hypothetical threshold contour under conditions where detection is mediated by second-site mechanisms. The contour is plotted in a cone contrast representation, and shows the contour in the $\Delta L/L$, $\Delta M/M$ contrast plane. The S-cone contrast is taken to be zero, so that thresholds are mediated by the L+M and L-M mechanisms. The contour was computed on the assumption that the gain of the L-M mechanism was four times that of the L+M mechanism, using a summation model with an exponent of 2. (b) Corresponding contour after habituation has reduced the gain of the L-M mechanism by a factor of 2 and left the L+M mechanism gain unchanged. The habituating direction is indicated by the black arrow. The result is the threshold contour becomes more elongated along the negative diagonal in the cone contrast plot.

in the $\Delta L/L$, $\Delta M/M$ plane obtained against a uniform background. This is the same type of data shown in Fig. 6*d*, but plotted on cone contrast axes. Detection is assumed to be mediated by L–M and L+M mechanisms with L–M being four times more sensitive to cone contrast than L+M. The elliptical contour reflects a summation exponent of 2, which we assume here for didactic purposes. If the same experiment is repeated after the subject has been habituated to contrast modulated in the L–M direction (as indicated by the black arrow Fig. 9*b*), the basic second-site model predicts that thresholds will be elevated for test stimuli that are modulated in that direction. On the other hand, for test stimuli modulated in the L+M direction, thresholds should be unaffected. Finally, for stimuli in intermediate directions, thresholds will be elevated to the extent that the L–M opponent mechanism contributes to detection. Figure 9*b* shows the predicted threshold contour, computed on the assumption that the L–M habituation cuts the gain of the L–M mechanism in half while leaving that of the L+M mechanism unaffected.

The data shown in Fig. 9 may be replotted by taking the difference of the post- and prehabituation threshold stimuli, for each color direction. Figure 10*a* shows an habituation effect plot of this sort. In this plot, the distance from the origin to the contour in each direction shows the threshold increase produced by the L–M habituation for test stimuli modulated in that color direction. The radius is large in the L–M test direction, and drops to zero for the L+M test direction.

Figure 10*b* and *c* show idealized habituation effect plots for habituation in the L+M direction and to an intermediate color direction, as shown by the black arrows. What the figure makes clear is that the effect of habituation on detection threshold should have a characteristic signature for each habituation direction.

Figure 11 replots the parts of Fig. 10 on axes that represent not cone contrast but rather stimulus directions that isolate the L+M and L–M discrimination mechanisms. These axes have been normalized so that a unit step along each axis direction corresponds to a threshold stimulus (without any habituation) for the corresponding mechanism. We discuss in more detail in subsection “Color Data Representations” in Sec. 11.5 the advantages and disadvantages of using an opponent representation of this type, but introduce it here to allow direct comparison of the idealized predictions of the basic model to the data reported by Krauskopf, Williams, and Heeley.¹⁴ These authors made measurements of how habituation in different color directions affected threshold contours. Their data for observer DRW, as shown in Fig. 29 in subsection “Habituation or Contrast Adaptation Experiments” in Sec. 11.5, agree qualitatively with the predictions shown in Fig. 11, and provide fairly direct evidence to support habituation at a second site. As with other aspects of the basic model, we return below to discuss more subtle features of the data that deviate from its predictions.

Sites of Limiting Noise

Before proceeding to color-appearance mechanisms, we pause to discuss a fundamental theoretical issue. This is the question of how the properties of a mechanism at an early stage of the visual processing chain could directly mediate observed detection thresholds, given that the signals from such a site must pass through many additional stages of processing before a perceptual decision is made. This question parallels one often asked in the context of colorimetry, namely, how it is that the very first stage of light absorption can mediate the overall behavior observed in color matching (see Chap. 10). Both questions have the same underlying answer, namely, that information lost at an early stage of visual processing cannot be restored at subsequent stages. The difference is the nature of the information loss. In the case of color matching, the information loss occurs because of the univariant nature of phototransduction. In the case of color thresholds, the information loss occurs because of noise in the responses of visual mechanisms. These sources of noise, relative to the signal strength, set the limits on visual detection and discrimination.

The link between thresholds and noise is well developed under the rubric of the theory of signal detection.^{43,66–68} Here we introduce the basic ideas. Consider a test stimulus of intensity c in the $\Delta L_d = 1$, $\Delta M_d = 0$ color direction. Only the L-cone mechanism responds to this stimulus, so for this initial example we can restrict attention to the L-cone response. If we present this same stimulus over many trials, the response of the L-cone mechanism will vary from trial-to-trial. This is indicated

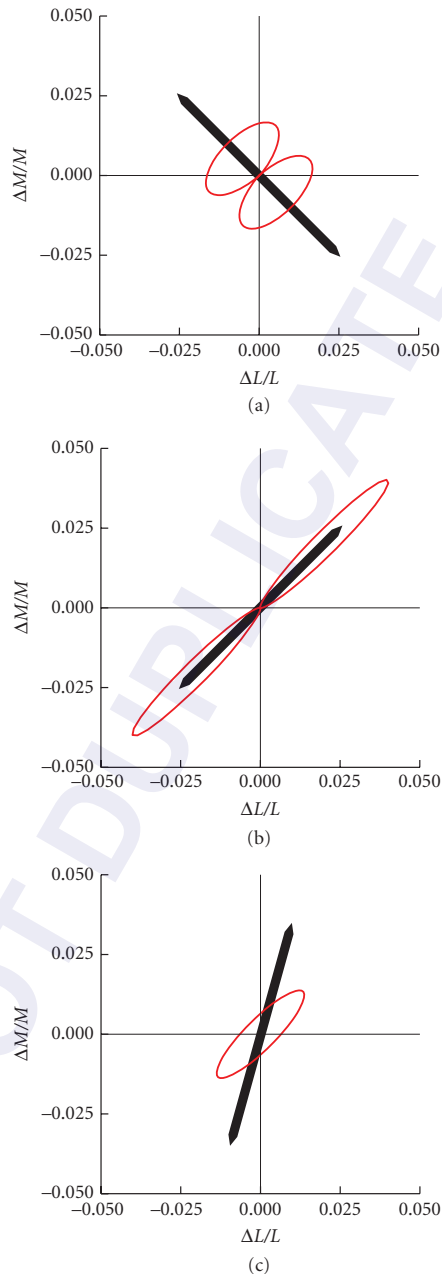


FIGURE 10 Changes in contrast threshold caused by habituation. It is conventional to report the results of a habituation experiment using a plot that shows the change in threshold caused by habituation. (a) Depiction of the hypothetical data shown in Fig. 9 in this fashion. For each direction in the cone contrast plot, the distance between the origin and the contour shows the increase in threshold caused by habituation for tests in that color direction. For the example shown, there is no change in threshold in the L+M contrast direction, and the contour returns to the origin for this direction. (b) Changes in threshold that would be produced by a habituating stimulus that decreases the gain of the L+M mechanism by a factor of 2 and leaves the gain of the L–M mechanism unchanged. (c) Changes in threshold that would be produced by a habituating stimulus that decreases the gain of both mechanisms equally by a factor of 1.33. The habituating directions are indicated by the black arrows.

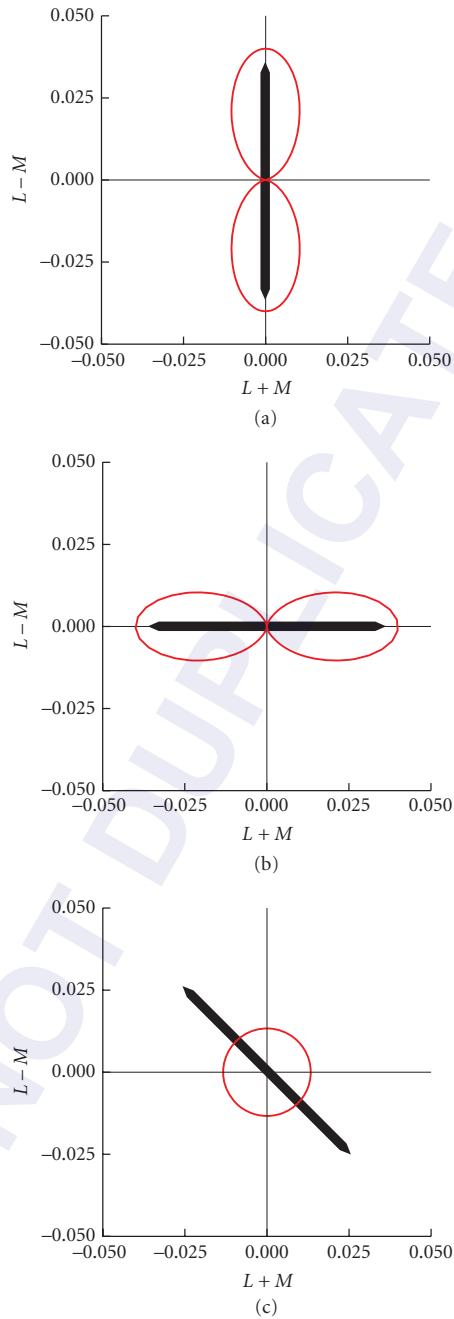


FIGURE 11 Changes in L+M and L-M sensitivity caused by habituation. Panels (a) to (c) replot the threshold changes shown in the corresponding panels of Fig. 10. Rather than showing the effects in cone-contrast space, here the results are shown in a space where the axes correspond to the L+M and L-M color directions. In addition, the units of each of these rotated axes have been chosen so that detection threshold for tests along the axis is 1. This representation highlights the canonical pattern of results expected when habituation along each axis reduces sensitivity for only the corresponding mechanism, and where habituation along the 45° direction reduces sensitivity equally for the two mechanisms. The habituating directions are indicated by the black arrows.

by the probability distribution (solid red line) in Fig. 12a. Although the mean response here is 0.06, sometimes the response is larger and sometimes smaller. Similarly, the mechanism response to the background alone will also fluctuate (probability distribution shown as dotted blue line).

It can be shown^{43,69} that optimal performance in a detection experiment occurs if the observer sets an appropriate criterion level and reports “background alone” if the response falls below this criterion

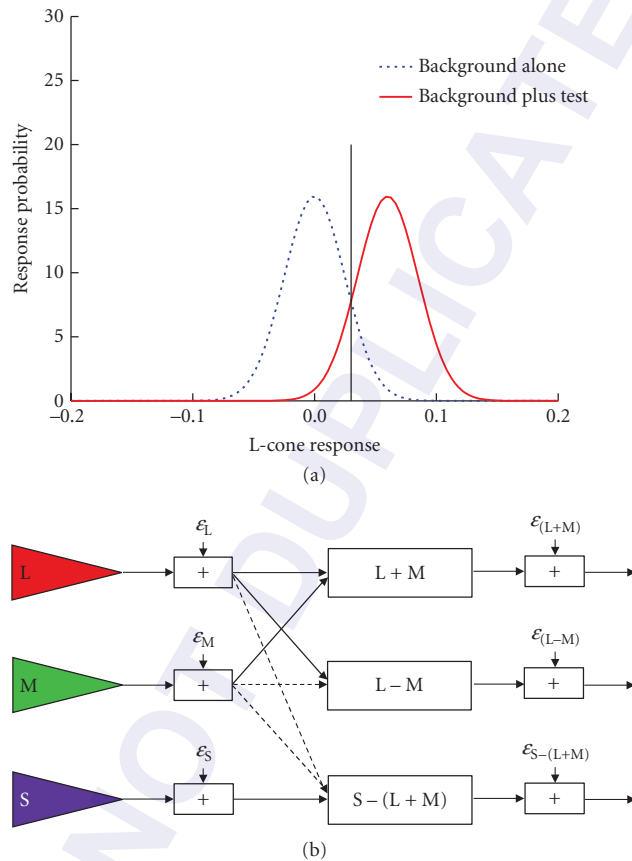


FIGURE 12 Thresholds and noise. (a) The theory of signal detection’s account of how noise determines thresholds, for a single mechanism. When the background alone is presented, the mechanism’s responses vary somewhat from trial-to-trial because of additive response noise. The response variation is illustrated by the probability distribution centered on zero and shown by the dotted blue line. When the background and test are presented together, the average response increases but there is still trial-to-trial variability. This is illustrated by the probability distribution shown by the solid red line. When the two types of trial are presented equally often and the magnitudes of the costs associated with correct and incorrect responses are equal, it can be shown that the observer maximizes percent correct by reporting that the test was present whenever the mechanism response exceeds the criterion value shown in the figure by the vertical line. Different costs and prior probabilities simply move the location of this criterion. Even when this optimal criterion is used, there will still be some incorrect trials, those on which the response to background alone exceeds the criterion and those on which the response to background and test falls below the criterion. Threshold is reached when the test is intense enough that the observer is correct on a sufficient percentage of trials. How much the average response to background plus test must exceed that to background alone is determined by the magnitude of the noise. (b) The two-stage model of color-discrimination mechanisms, drawn in a manner that emphasizes that noise (ϵ) may be added both to the output of the first-site mechanisms and to the output of the second-site mechanisms. A full signal detection theoretic model of detection takes into account the response gain at each stage (which determines the separation of the average response to background and to background plus test) and the magnitude of noise at each stage. It also replaces the simple criterion shown in (a) with an optimal multivariate classifier stage.⁷⁰

and “background plus test” otherwise. The plot shows the location of the optimal criterion (vertical line) for the special case when background-plus-test and background-alone events are presented equally often and when the benefits associated with both sorts of correct response (hits and correct rejections) have the same magnitude as the costs associated with both sorts of incorrect response (false alarms and misses). Note that even the optimal strategy will not lead to perfect performance unless the intensity of the test is large enough that the two distributions become completely separated. In this sense, the magnitude of response noise (which determines the widths of the distributions), relative to the effect of the test on the mean mechanism response, is what limits performance. More generally, the theory of signal detection allows computation of the probability of the observer performing correctly in the detection experiment as a function of the noise magnitude and test intensity.

With this idea in mind, we can turn to understanding how threshold measurements can reveal properties of early mechanisms, despite later processing by subsequent mechanisms. Figure 12*b* shows the basic model of color-discrimination mechanisms from Fig. 5, but without the adaptation stages shown explicitly. What is shown here, however, is that noise (ϵ) is added to the response of each mechanism. At the first stage, the cone responses are noisy. The noisy cone responses are then transformed at the second-site, and more noise is added. Thus the responses of the second-site mechanisms are subject both to noise that propagates from the first site, in addition to noise injected directly at the second site. For simplicity, we consider the noise added to each mechanism at each site to be independent; note that even with this assumption the noise at the opponent site will be correlated across the three color-discrimination mechanisms as a result of the way the first-site noise is transformed.

This model allows us to compute the expected shape of the threshold contours for a visual system that makes optimal use of the output of the second-site mechanisms. The computation is more involved than the simple unidimensional signal detection example illustrated in Fig. 12*a*, because optimal performance makes use of the output of all three second-site mechanisms and must take into account the correlated nature of the noise at this stage. Nonetheless, the ideas underlying the computation are the same as in the simple example, and the theory that allows the computation is well worked out.⁷⁰ We can thus ask how the expected discrimination contours, computed on the basis of the output at the second site, vary with the properties of the noise added at each stage.

Figure 13 shows threshold contours computed for three choices of noise. Figure 13*a* shows a case where the standard deviation of the noise injected at the first site is 5 times larger than that at the second site. The derived contour has the properties expected of detection by cone mechanisms, with the major axes of the contour aligned with the L- and M-cone contrast axes (see Fig. 6*b* and *c*). Figure 13*c* on the other hand, shows the opposite situation where the noise injected at the second site is 5 times larger than that at the first site. Simply changing the noise magnitude has major effects on the observed contour; the shape of the contour shown in Fig. 13*c* is the expected result for detection mediated by second-site mechanisms (see Fig. 6*d*). Finally, Fig. 13*b* shows an intermediate case, with equal noise injected at the two sites. Here the result is intermediate between the two other cases.

This example lets us draw a number of important conclusions. First, the result in Fig. 13*a* shows explicitly how, in a two-stage visual system, threshold contours can reveal properties of the first-stage mechanism. The necessary condition is that the noise added before the transformation to the second stage be large enough to dominate the overall performance. When a situation like this holds, we say that the first stage represents the *site of limiting noise*. As long as subsequent processing does not add significant additional noise, it is the site of limiting noise whose properties will be reflected by threshold measurements. Although Fig. 13*a* shows an example where the cones were the site of limiting noise, the fact that detection contours often have the shape shown in Fig. 13*c* suggests that often the second site limits performance, and that subsequent processing at sites past the second site do not add significant additional noise. As an aside, we note that much of threshold psychophysics is concerned with arranging auxiliary stimulus manipulations that manipulate mechanism properties so as to move the site of limiting noise from one stage of processing to another, or to shift it across mechanisms within a stage (e.g., to change the relative contribution of L+M and L–M mechanisms to performance), with the goal of enabling psychophysical probing of individual mechanisms at different stages of processing. Generally, the stimulus manipulations are thought to

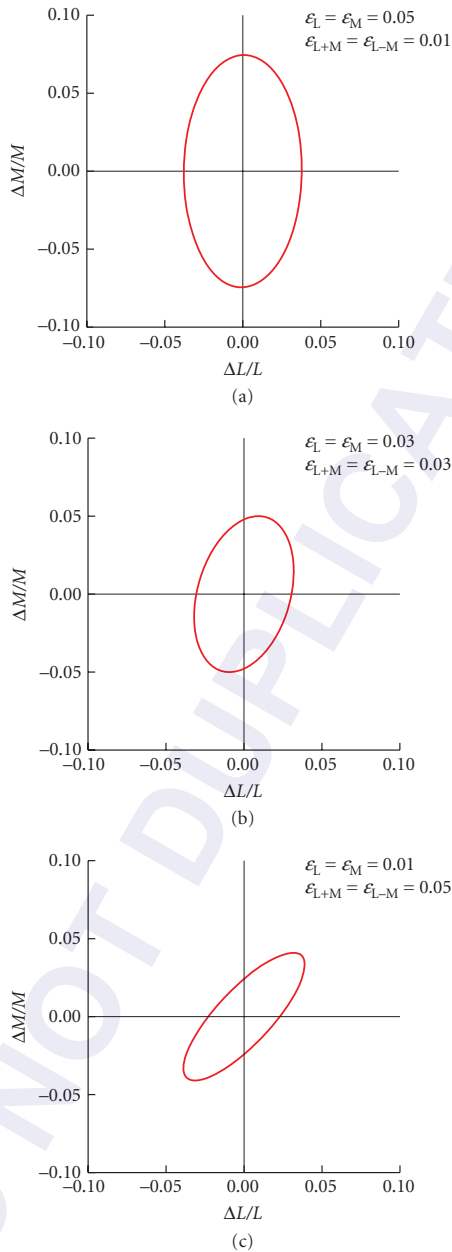


FIGURE 13 Threshold predictions from signal detection theoretical model. The figure shows threshold contours predicted using the model shown in Fig. 12b, for various choices of noise amplitude. In all cases, only the L+M and L-M mechanisms were considered. The gain on L-cone contrast was set at 2, while the gain on M-cone contrast was set at 1. The gain of the L+M mechanism was set at 1, while that of the L-M mechanism was set at 4. The noise added at each site was assumed to be Gaussian, and for each test color direction and choice of noise magnitude, the test intensity that led to optimal classification performance of 75 percent was determined. (a) Resulting contour when the noise at the first site is five times that of the noise at the second site. In this case, the detection contour has the shape expected for the first-site mechanisms, because little additional information is lost at the second-site mechanisms. (c) Contour when the magnitude of the noise at the second site is five times that of the first. In this case, the properties of the second-site mechanisms dominate measured performance. (b) Contour when the noise at the two sites is of comparable magnitude, and here the result is itself intermediate between that expected from the properties of first- and second-site mechanisms.

act by changing the gains applied to mechanism outputs. This changes the effect of the stimulus relative to the noise, rather than the noise magnitude itself.

A second important point can be drawn from Fig. 13*b*. It is not unreasonable to think that noise at different stages of visual processing will be of a comparable magnitude, and Fig. 13*b* shows that in such cases observed performance will represent a mixture of effects that arise from mechanisms properties at different stages. This is a more formal development of a point we stressed in the introduction, namely, that in spite of our attempt to attribute performance under any given experimental conditions primarily to the action of a small number of mechanisms, this is likely to be a highly simplified account. Eventually, it seems likely that models that explicitly account for the interactions between mechanisms at different stages will be needed. The theory underlying the production of Fig. 13, which explicitly accounts for the role of noise at each stage, would be one interesting approach toward such models.

Finally, note that the development of Fig. 13 is in the tradition of ideal-observer models, where performance is derived from an explicit model of a set of stages of visual processing, followed by the assumption that further processing is optimal.^{71–74} Although the assumption of optimal processing is unlikely to be exactly correct, ideal-observer models are useful because they can provide theoretical clarity and because cases where the data deviate from the predictions of such a model highlight phenomena that require further exploration and understanding. For example, the rounded shape of threshold contours, which we initially described using a descriptive vector length account, emerges naturally from the ideal-observer analysis. Indeed, for the ideal observer the predicted shape of the threshold contours is essentially ellipsoidal (summation exponent of 2). This prediction is another reason why careful assessment of the observed shape of threshold contours is of theoretical interest.

In subsection “Noise-Masking Experiments” in Sec. 11.5, we consider the effects of adding external noise to the system.

11.4 BASICS OF COLOR-APPEARANCE MECHANISMS

We now turn to consider how a different class of observations can inform us about mechanisms. These are measurements that assess not the observer’s ability to detect or discriminate stimuli, but rather measurements of how stimuli look to observers. More specifically, we will consider measurements of color appearance.

As with threshold measurements, color-appearance measurements may be divided into two classes, those that assess the appearance of test stimuli as some function of those stimuli (test measurements) and those that assess how the color appearance of a test stimulus depends on the context in which it is viewed (field measurements). We review how both sorts of measurements can inform a model of color-appearance mechanisms. This treatment parallels our development for color-discrimination mechanisms in order to highlight similarities and differences between models that account for the two types of data.

Appearance Test Measurements and Opponency

As described in the introduction, one of the earliest suggestions that the output of the cone mechanisms were recombined at a second site came from Hering’s observation that certain pairs of color sensations (i.e., red and green, blue and yellow) are mutually exclusive. Hurvich and Jameson elaborated this observation into a psychophysical procedure known as hue cancellation (see subsection “Spectral Properties of Color-Opponent Mechanisms” in Sec. 11.5). Hue cancellation can be performed separately for the red/green opponent pair and the blue/yellow opponent pair.

In the red/green hue-cancellation experiment, the observer is presented with a test stimulus, which is usually monochromatic, and asked to judge whether its appearance is reddish or greenish. The results may be used to derive the properties of an opponent mechanism, under the “linking

hypothesis” that when the mechanism’s response is of one sign the stimulus appears reddish, while when it is of the opposite sign the stimulus appears greenish. That is, color-appearance measurements are used to derive mechanism properties on the assumption that the mechanism output provides an explicit representation of the sensation experienced by the observer.

To understand the hue-cancellation procedure in more detail, imagine that the experimenter picks a fixed monochromatic reference light that in isolation appears greenish to the observer. The experimenter then chooses a series of monochromatic test lights of unit intensity, each of which appears reddish in isolation. The observer judges mixtures of each test light and the reference light. If the mixture appears reddish, the intensity of the greenish reference light is increased. If the mixture appears greenish, the intensity of the greenish reference light is decreased. Repeating this procedure allows the experimental determination of a balance or equilibrium point for the red-green mechanism, where the result of the mixture appears neither reddish nor greenish. (Depending on the test stimulus, the balanced mixture will appear either as yellow or blue or achromatic.) Colored stimuli that appear neither reddish nor greenish are referred to as “unique yellow” or “unique blue.” The intensity of the green reference light in the balanced mixture may be taken as a measure of the amount of redness in the test stimulus, in the same way that the amount of weight added to one side of a balance scale indexes the weight on the other side when the balance is even (see Krantz’s “force table” analogy.⁷⁵)

A similar procedure may be used to determine the amount of greenness in test stimuli that have a greenish component, by choosing an adjustable reference stimulus that appears reddish. Moreover, the units of redness and greenness may be equated by balancing the two reference stimuli against each other. Figure 14*a* and *b* show prototypical results of such measurements for spectral test lights, as well as results for an analogous procedure used to determine the amount of blueness and yellowness in the same set of spectral test lights. Such data are referred to as “chromatic valence” data.

Krantz^{75,76} showed that if the amount of redness, greenness, blueness, and yellowness in any test light is represented by the output of two opponent color-appearance mechanisms, and if these mechanisms combine cone signals linearly, then the spectral hue-valence functions measured in the cancellation experiment must be a linear transformation of the cone spectral sensitivities. Although the linearity assumption is at best an approximation (see subsection “Linearity of Color-Opponent Mechanisms” in Sec. 11.5), Krantz’s theorem allows derivation of the cone inputs from the data shown in Fig. 14. Such fits are shown in Fig. 14*a* and *b*, and lead to the basic wiring diagram shown in Fig. 14*c*. This shows two opponent color-appearance mechanisms, labeled R/G and B/Y. The cone inputs to these mechanisms are similar to those for the opponent color-discrimination mechanisms, with one striking difference: the S cones make a strong R contribution to the R/G appearance mechanism. By contrast, the S-cone input to the L–M discrimination mechanism is small at most. Another difference, which we address later, is that the M-cone contribution to B/Y in some models is of the same sign as the S-cone contribution (see “Three-Stage Zone Models” in Sec. 11.6).

Although the hue-cancellation procedure may be used to derive properties of opponent-appearance mechanisms, it is less suited to deriving the properties of a third mechanism that signals the brightness of colored stimuli. This aspect of color appearance may be assessed using heterochromatic brightness matching, where subjects are asked to equate the brightness of two stimuli, or scaling methods, where observers are asked to directly rate how bright test stimuli appear (see subsection “Luminance and Brightness” in Sec. 11.5).

Appearance Field Measurements and First-Site Adaptation

It is also possible to make field measurements using appearance judgments. A common method is known as asymmetric matching, which is a simple extension of the basic color-matching experiment (see Chap. 10). In the basic color-matching experiment, the observer adjusts a match stimulus to have the same appearance as a test stimulus when both are seen side-by-side in the same context. The asymmetric matching experiment adds one additional variable, namely that the test and match are each seen in separate contexts. As illustrated in Fig. 1, when an identical test patch is presented

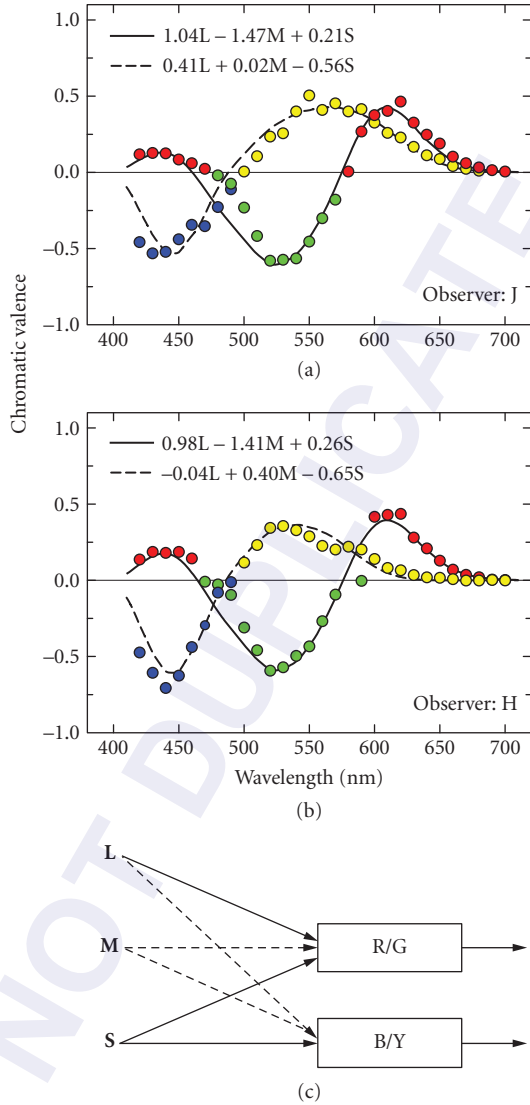


FIGURE 14 Chromatic valence data and wiring. (a) and (b) Valence data (colored symbols) replotted from Figs. 4 and 5 of Jameson and Hurvich²⁷⁹ for observers “J” (a) and “H” (b) fitted with linear combinations of the Stockman and Sharpe²⁸ cone fundamentals (solid and dashed curves). The best-fitting cone weights are noted in the key. (c) Wiring suggested by this (a) and subsequent color valence data. In the diagram, the sign of the M-cone contribution to B/Y is shown as negative (see subsection “Spectral Properties of Color-Opponent Mechanisms” in Sec. 11.5).

against two different surrounds, it can appear quite different. In the asymmetric matching experiment, the observer compares test patches seen in two different contexts (e.g., against different adapting backgrounds) and adjusts one of them until it matches the other in appearance. The cone coordinates of the match test patches typically differ, with the shift being a measure of the effect of changing from one context to the other.

Figure 15a plots asymmetric matching data from an experiment of this sort reported by MacAdam.⁷⁷ The filled circles show the L- and S-cone coordinates of a set of test stimuli, seen in a region of the retina adapted to a background with the chromaticity of a typical daylight. The open circles show a subset of the corresponding asymmetric matches seen in a region of the retina adapted to a typical tungsten illuminant. The change in adapting background produces a change in test appearance, which observers have compensated for in their adjustments.

Asymmetric matching data can be used to test models of adaptation. For example, one can ask whether simple gain control together with subtractive adaptation applied to the cone signals can account for the matches. Let $(L_1, M_1, S_1), (L_2, M_2, S_2), \dots, (L_N, M_N, S_N)$ represent the L-, M-, and S-cone coordinates of N test stimuli seen in one context. Similarly, let $(L'_1, M'_1, S'_1), (L'_2, M'_2, S'_2), \dots, (L'_N, M'_N, S'_N)$ represent the cone coordinates of the matches set in the other context. If the data can be accounted for by cone-specific gain control together with subtractive adaptation, then we should be able to find gains g_L, g_M, g_S and subtractive terms L_0, M_0, S_0 such that $L'_i \approx g_L L_i - L_0$, $M'_i \approx g_M M_i - M_0$, and $S'_i \approx g_S S_i - S_0$ for all of the N matches. The asterisks plotted in Fig. 15a show predictions of this gain control model in the L-S cone plane. Figure 15b, c, and d show predicted shifts against the measured shifts for all three cone classes. The model captures the broad trends in the data, although it misses some of the detail. The first-site adaptation captured by this model is often taken to be the same adaptation revealed by sensitivity experiments.⁷⁸

Appearance Field Measurements and Second-Site Adaptation

Asymmetric matching may also be used to investigate the effect of contrast adaptation on color appearance. In a series of experiments, Webster and Mollon^{32,79} investigated the effect of contrast adaptation (or habituation) on the color appearance of suprathreshold lights (see also Refs. 80 and 81). Observers habituated to a background field that was sinusoidally modulated at 1 Hz along various color directions. The habituating field was interleaved with a stimulus presented at the same location. The observer's task was to adjust the chromaticity and luminance of a match stimulus so that its appearance matched that of the test field. Webster and Mollon found that habituation affected the appearance of the test stimulus.

The motivation behind these experiments was to test whether or not the changes in color appearance caused by contrast adaptation were consistent with the type of second-site adaptation that characterizes sensitivity experiments, namely independent adaptation of three cardinal mechanisms with cone inputs: L-M, L+M, and S-(L+M). The prediction for this case is that the appearance effects should always be greatest along the color axis corresponding to the mechanism most adapted by the habituating stimulus, and least along the axis corresponding to the mechanism least adapted by the stimulus—and this pattern should be found whatever the chromatic axis of the habituating stimulus. Thus the *relative* changes in chromaticity and/or luminance caused by habituation should be roughly elliptical with axes aligned with the mechanism axes.

The results showed that the sensitivity changes did not consistently align with the axes of the cardinal mechanisms. Instead, the largest selective sensitivity losses in the equiluminant plane, for two out of three subjects, aligned with the adapting axis, while the smallest losses aligned with an axis approximately 90° away. Comparable results were found when the habituating and test stimuli varied in color and luminance in either the plane including the L-M and L+M+S axes or the plane including the S and L+M+S axes. Thus, contrast adaptation produces changes in color appearance that are, in general, selective for the habituating axis rather than selective for the purported underlying mechanism axes. Nevertheless, greater selectivity was found for habituation along the three cardinal axes, which suggests their importance. See Ref. 32 for further details.

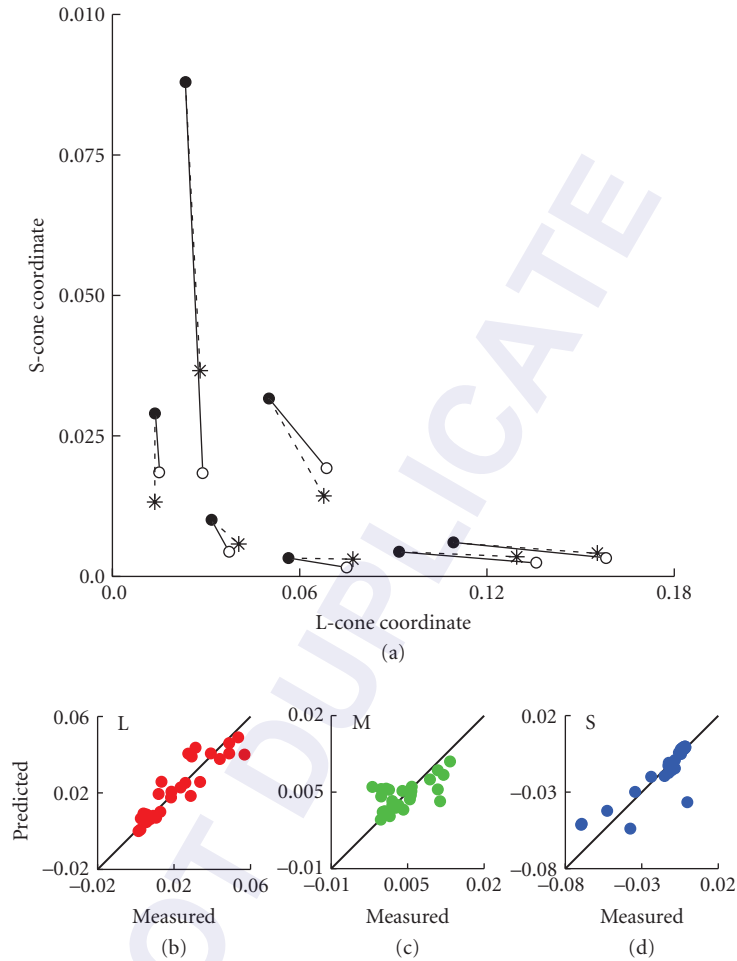


FIGURE 15 Analysis of asymmetric matching data. (a) Subset of the data from Observer DLM reported by MacAdam.⁷⁷ The closed circles show L- and S-cone coordinates of tests viewed on part of the retina adapted to a background with the chromaticity of daylight. The connected open circles show the asymmetric matches to these tests, when the matching stimulus was viewed on part of the retina adapted to the chromaticity of a tungsten illuminant. The star symbols show the predictions of a model that supposes first-site gain control and subtractive adaptation, operating separately within each cone mechanism. (b), (c), and (d) Model predictions for the full set of 29 matches set by Observer DLM. Each panel shows data for one cone class. The x axis shows the measured difference between the test and match, and represents the size of the asymmetric matching effect. The y axis shows the corresponding predicted difference. If the model were perfect, all of the data would lie along the positive diagonals. The model captures the broad trend of the data, but misses in detail.

Webster and Mollon³² conclude that the changes in color appearance following habituation are inconsistent with models of color vision that assume adaptation in just three independent postreceptoral channels. Given, however, that the properties of the color-appearance mechanisms as revealed by test methods differ from those of color-discrimination mechanisms, this result is hardly surprising. Also note that comparisons between sensitivity measurements made for test stimuli near detection threshold and appearance measurements made for suprathreshold tests are complicated by the presence of response nonlinearities at various stages of processing. See Refs. 78 and 82 for discussion.

11.5 DETAILS AND LIMITS OF THE BASIC MODEL

Sections 11.3 and 11.4 introduce what we will refer to as a “basic model” of color-discrimination and color-appearance mechanisms. There is general agreement that the components incorporated in this model are correct in broad outline. These are (1) first-site, cone specific, adaptation roughly in accord with Weber’s law for signals of low temporal frequency, (2) a recombination of cone signals that accounts for opponent and nonopponent postreceptoral second-site signals, (3) second-site desensitization and habituation, and (4) the necessity for distinct mechanistic accounts of discrimination and appearance.

In the earlier sections, we highlighted example data that strongly suggest each feature of the basic model. Not surprisingly, those data represent a small subset of the available experimental evidence brought to bear on the nature of postreceptoral color vision. In this section we provide a more detailed review, with the goal of highlighting both data accounted for by the basic model, as well as data that expose the model’s limitations. Current research in the mechanistic tradition is aimed at understanding how we should elaborate the basic model to account parsimoniously for a wider range of phenomena, and at the end of this chapter we outline approaches being taken toward this end.

Sections “Test Sensitivities” through “Color Appearance and Color Opponency” follow the same basic organization as the introductory sections. We consider discrimination and appearance data separately, and distinguish between test and field methods. We begin, however, with a discussion of color data representations.

Color Data Representations

Over the past 50 years, color vision research has been driven in part by the available technology. Conventional optical systems with spectral lights and mechanical shutters largely restricted the volume of color space that could be easily investigated to the 1/8 volume containing positive, incremental cone modulations. When such data are obtained as a function of the wavelength of monochromatic stimuli, a natural stimulus representation is to plot the results as a function of wavelength (e.g., Fig. 14*a* and *b*). This representation, however, does not take advantage of our understanding of the first stage of color processing, transduction of light by the cone photoreceptors. Moreover, with the availability of color monitors and other three-primary devices, stimuli are now routinely generated as a combination of primaries, which makes increments, decrements, and mixtures readily available for experimentation. Spectral plots are not appropriate for such stimuli, which instead are represented within some tristimulus space (see Chap. 10).

Our estimates of the human cone spectral sensitivities have become increasingly secure over the past 30 years (see Chap. 10). This allows stimulus representations that explicitly represent the stimulus in terms of the L-, M-, and S-cone excitations and which therefore connect data quite directly to the first-stage color mechanisms. Such representations are called “cone-excitation spaces,” and an example of this sort of representation is provided in Fig. 7 of Chap. 10.

An extension of cone-excitation space is one in which the cone excitations are converted to increments and decrements of each cone type relative to a background. Thus, the three axes are the changes in cone excitations: $\pm\Delta L$, $\pm\Delta M$, and $\pm\Delta S$. We introduced such spaces in Fig. 6; they are useful when

the stimulus is most naturally conceived as a modulation relative to a background, and where the properties of the background per se are of less interest. We will use the shorthand “incremental cone spaces” to refer to this type of representation, although clearly they allow expression of both increments and decrements.

Closely related to incremental cone spaces are “cone contrast spaces,” in which the incremental/decremental changes in cone excitations produced by the target are divided by the cone excitations produced by the background.^{61,83} Here, the three axes are dimensionless cone contrasts: $\Delta L/L_b$, $\Delta M/M_b$, and $\Delta S/S_b$. If Weber’s law holds independently for each cone type then $\Delta L/L_b$, $\Delta M/M_b$, and $\Delta S/S_b$ all remain constant. Thus a particular advantage of this space is that it factors out the gross effects of first-site cone-specific adaptation, which tends to follow Weber’s law at higher intensities (see subsection “First-Site Adaptation” in Sec. 11.3). Plots in cone contrast space help to emphasize desensitization that occurs after the receptors, and, to the extent that Weber’s law holds, provide an explicit representation of the inputs to postreceptoral mechanisms. (We introduced a cone contrast representation in Fig. 8 exactly for this reason.)

At lower adaptation levels, when adaptation falls short of Weber’s law (see Fig. 7) or under conditions where Weber’s law does not hold (such as at high spatial and temporal frequencies; see, e.g., Ref. 84), cone contrast space is less useful. However, for targets of low temporal and spatial frequencies, Weber’s law has been found to hold for chromatic detection down to low photopic levels.⁶²

A final widely used color space is known as the “Derrington-Krauskopf-Lennie” (DKL) space^{36,85,86} in which the coordinates represent the purported responses of the three second-site color-discrimination mechanism, L+M, L–M, and S–(L+M). In this context, modulation directions that change the response of one of these mechanisms while leaving the response of the other two fixed are referred to as “cardinal directions.”¹⁴ The DKL representation was introduced in Fig. 11 and is further illustrated in Fig. 16. See Ref. 87 for further discussion of color spaces and considerations of when each is most appropriately used.

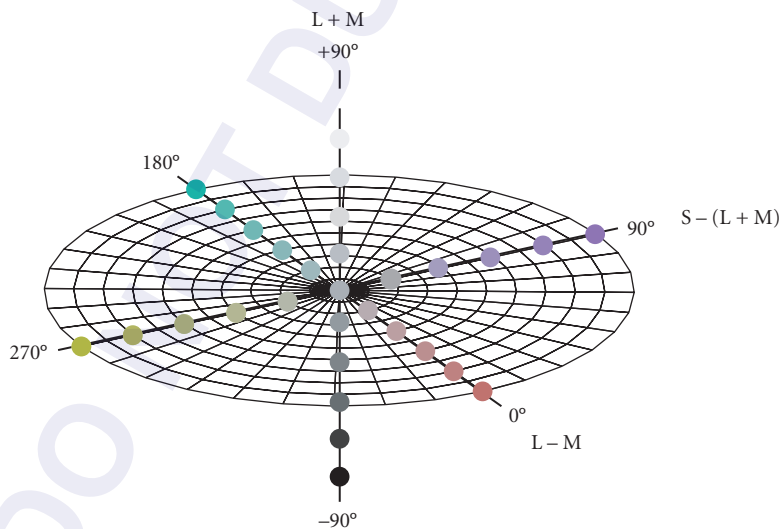


FIGURE 16 Derrington-Krauskopf-Lennie (DKL) color space. The grid corresponds to the equiluminant plane, which includes the L–M (0° – 180°) and S–(L+M) (90° – 270°) cardinal mechanisms axes. The vertical axis is the achromatic L+M axis (-90° – $+90^\circ$). The colors along each axis are approximate representations of the appearances of lights modulated along each cardinal axis from an achromatic gray at the center. Notice that the unique hues do not appear. The axes are labeled according to the mechanisms that are assumed to be uniquely excited by modulations along each cardinal axis. (Figure provided by Caterina Ripamonti.)

Although the conceptual ideas underlying the various color spaces discussed here are clear enough, confusion often arises when one tries to use one color space to represent constructs from another. The confusion arises because there are two distinct ways to interpret the axes of a three-dimensional color space. The first is in terms of the responses of three specified mechanisms, while the second is in terms of how the stimulus is decomposed in terms of the modulation directions that isolate the three mechanisms. To understand why these are different, note that to compute the response of a particular mechanism, one only needs to know its properties (e.g., that it is an L–M mechanism). But, in contrast, the stimulus direction that isolates that same mechanism is the one that silences the responses of the other two mechanisms. The isolating stimulus direction depends not on the properties of the mechanism being isolated but on those of the other two. For this reason, it is important to distinguish between the *mechanism direction*, which is the stimulus direction that elicits the maximum positive response of a given mechanism per unit intensity, and the *isolating direction* for the same mechanism, which produces no response in the other two mechanisms. For the triplet of mechanisms L+M, L–M, and S–(L+M), the mechanism and isolating directions align (by definition) in DKL space, but this property does not hold when these directions are plotted in the antecedent incremental cone space or cone contrast space. Note that in some descriptions of DKL space, the axes are defined in terms of the cone modulations that isolate mechanisms rather than in terms of the mechanism responses. These conflicting definitions of DKL space have led to a good deal of confusion. See Chap. 10 and Refs. 87–89 for further discussion of transformations between color spaces and the relation between mechanism sensitivities, mechanism directions, and isolating directions.

Caveats Crucial for the use of any color space based on the cone mechanisms either for generating visual stimuli or for interpreting the resulting data is that the underlying cone spectral sensitivities are correct. Our estimates of these are now secure enough for many applications, but there are uncertainties that may matter for some applications, particularly at short wavelengths. These are compounded by individual differences (see Chap. 10), as well as issues with some of the sets of color-matching functions used to derive particular estimates of cone spectral sensitivities (again see Chap. 10, and also Ref. 90 for further discussion).

Such concerns are magnified in the use of postreceptoral spaces like the DKL space, because as well as the cone spectral sensitivities needing to be correct, the rules by which they are combined to produce the spectral sensitivities of the postreceptoral mechanisms must also be correct. For example, to silence the L+M luminance mechanism it is necessary to specify exactly the relative weights of the M and L cones to this mechanism. However, these weights are uncertain. As discussed in subsection “Luminance” in Sec. 11.5, the M- and L-cone weights show large individual differences. Furthermore, the assumption typically used in practice, when constructing the DKL space, is that luminance is a weighted sum of M and L cones and that the S cones therefore do not contribute to luminance. This may not always be the case.^{91–94} Luminous efficiency also varies with chromatic adaptation, the spatial and temporal properties of the stimulus, and the experimental task used to define luminous efficiency. Lastly, the standard 1924 CIE $V(\lambda)$ luminous efficiency function, which determines candelas/m², and has been used by several groups to define the spectral sensitivity of the luminance mechanism, seriously underestimates luminous efficiency at shorter wavelengths (see Ref. 90).

One final note of caution is that it is useful to keep in mind that independently representing and manipulating the responses of mechanisms at a particular stage of visual processing only makes sense if there are three or fewer mechanisms. If there are more than three photopic mechanisms at the second stage of color processing (see also subsection “Low-Level and Higher-Order Color-Discrimination Mechanisms” in Sec. 11.6), independently manipulating all of them is not possible because of the inherent three-dimensionality of the first stage.

These caveats are not meant to discourage entirely the use of physiologically motivated color spaces for data representation; on the contrary we believe such representations offer important advantages if used judiciously. It is crucial, however, to bear in mind that such spaces usually assume particular theories of processing, and that interpreting data represented within them should be done in the context of a clear understanding of the limitations of the corresponding theories.

Test Sensitivities

The most direct method of studying the chromatic properties of color mechanisms is the test method. Variations of this method have been tailored to investigate the properties of different visual mechanisms, including cone, chromatic, and achromatic (or luminance) mechanisms.

Sensitivity to Spectral Lights Perhaps the most influential chromatic detection data have been the two-color threshold data of Stiles,^{41,42,95–98} so called because the threshold for detecting a target or test field of one wavelength is measured on a larger adapting or background field usually of a second wavelength (or mixture of wavelengths). These data were collected within the context of an explicit model of color mechanisms, referred to as “ π -mechanisms,” and how they adapt. Enoch⁹⁹ provides a concise review of Stiles’ work.

An important limitation of Stiles’ model is that it did not allow for opponency: the π -mechanisms had all-positive spectral sensitivities and interactions between them were required to be summative. In some instances, the deviations caused by opponency were accommodated in Stiles’ model by the postulation of new π -mechanisms, such as π_5' (see Refs. 54 and 55). As data in the two-color tradition accumulated, this limitation led to the demise of the π -mechanism framework and the development of the basic model we introduced in Sec. 11.3. Nonetheless, any correct model must be consistent with the experimental data provided by the two-color threshold technique, and we review some key aspects in the context of current thinking.

Two-color threshold test spectral sensitivity measurements obtained on steady backgrounds as a function of target wavelength have characteristic peaks and troughs that depend on the background wavelength (see Fig. 1 of Ref. 98). These undulations were interpreted as the envelope of the spectral sensitivities of the underlying cone mechanisms (or “ π -mechanisms,” see Fig. 23), but crucially they are now known also to reflect opponent interactions between mechanisms.¹⁰⁰

The shape of test spectral sensitivity functions depend not only on background wavelength but also on the type of target used in the detection task. Detection by second-site opponent color-discrimination mechanisms is generally favored by targets comprised mainly of low temporal and spatial frequency components, whereas detection by achromatic mechanisms, which sum cone signals, is favored by targets with higher temporal and spatial frequency components.^{83,101,102} Stiles typically used a flashed target of 1° in visual diameter and 200 ms in duration, which favors detection by chromatic mechanisms. When Stiles’ two-color threshold experiments are carried out using brief, 10-ms duration targets, the results are more consistent with detection being mediated by achromatic mechanisms.¹⁰³

Detection by chromatically opponent mechanisms is most evident in test spectral sensitivity functions measured using targets of low temporal and/or spatial frequency. Measured on neutral fields, chromatic detection is characterized by peaks in the spectral sensitivity curves at approximately 440, 530, and 610 nm that are broader than the underlying cone spectral sensitivity functions and separated by pronounced notches.^{102,104–108} The so-called “Sloan notch”¹⁰⁹ corresponds to the loss of sensitivity when the target wavelength is such that the target produces no chromatic signal (e.g., when the L- and M-cone inputs to the L–M chromatic mechanism are equal), so that detection is mediated instead by the less-sensitive achromatic mechanism. By contrast, the broad peaks correspond to target wavelengths at which the target produces a large chromatic signal.^{101,102,105–108,110,111} Examples of detection spectral sensitivities with strong contributions from opponent mechanisms are shown in Fig. 17a and Fig. 18a. Those shown as circles in Fig. 17 and Fig. 18 were measured on a white background. The other functions in Fig. 17 were measured on backgrounds of 560 nm (green squares), 580 nm (yellow triangles), and 600 nm (orange triangles).

The nature of the mechanisms responsible for the spectral sensitivity data in Fig. 17a can be seen clearly by replotting the data as cone contrasts, as illustrated in Fig. 17b for the three chromatic backgrounds (from Ref. 89). As indicated by the straight lines fitted to each set of data, the detection contours have slopes close to one in cone contrast space, which is consistent with detection by L–M chromatic mechanisms with equal cone contrast weights (see subsection “Sensitivity to Different Directions of Color Space” in Sec. 11.5). Note also that the contours move outward with increasing field wavelength. This is consistent with second-site desensitization (see Fig. 8). The dependence

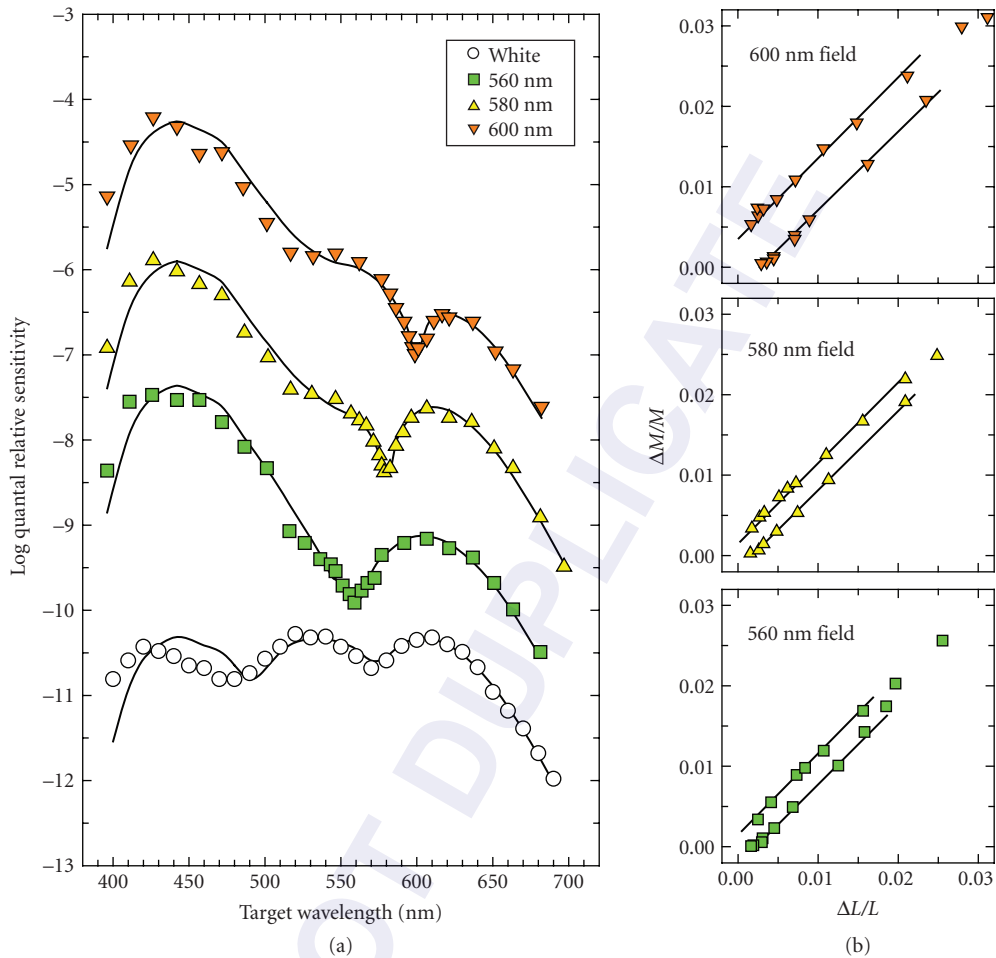


FIGURE 17 Sloan notch and chromatic adaptation. The three data sets in (a) are spectral sensitivity data replotted from Fig. 1 of Thornton and Pugh¹⁰⁷ measured on 10^{10} quanta sec^{-1} deg^{-2} fields of 560 nm (green squares), 580 nm (yellow triangles), and 600 nm (orange inverted triangles) using a long duration (one period of a 2-Hz cosine), large (3° diameter, Gaussian-windowed) target. The lowest data set in (a) are data replotted from Fig. 4 (Observer JS) of Sperling and Harwerth¹⁰⁵ measured on 10,000 td, 5500 K, white field (open circles) using a 50-ms duration, 45-min diameter target. The data have been modeled (solid lines) by assuming that the spectral sensitivities can be described by $b|aL - M| + c|M - aL| + d|S - e(L + 0.5M)|$, where L, M, and S are the quantal cone fundamentals²⁸ normalized to unity peak, and $a-e$ are the best-fitting scaling factors. These fits are comparable to the ones originally carried out by Thornton and Pugh¹⁰⁷ using different cone fundamentals. Notice that the Sloan notch (on colored backgrounds) coincides with the field wavelength. The detection contours in (b) have been replotted from Fig. 18.4B of Eskew, McLellan, and Giulianini.⁸⁹ The data are the spectral sensitivities measured on 560 nm (green squares), 580 nm (yellow triangles), and 600 nm (orange inverted triangles) fields already shown in (a), but transformed and plotted in cone-contrast space. Only data for target wavelengths ≥ 520 nm are shown. As indicated by the straight lines fitted to each set of data, the detection contours have slopes close to one in cone contrast space, which is consistent with detection by L–M chromatic mechanisms with equal L- and M-cone contrast weights.

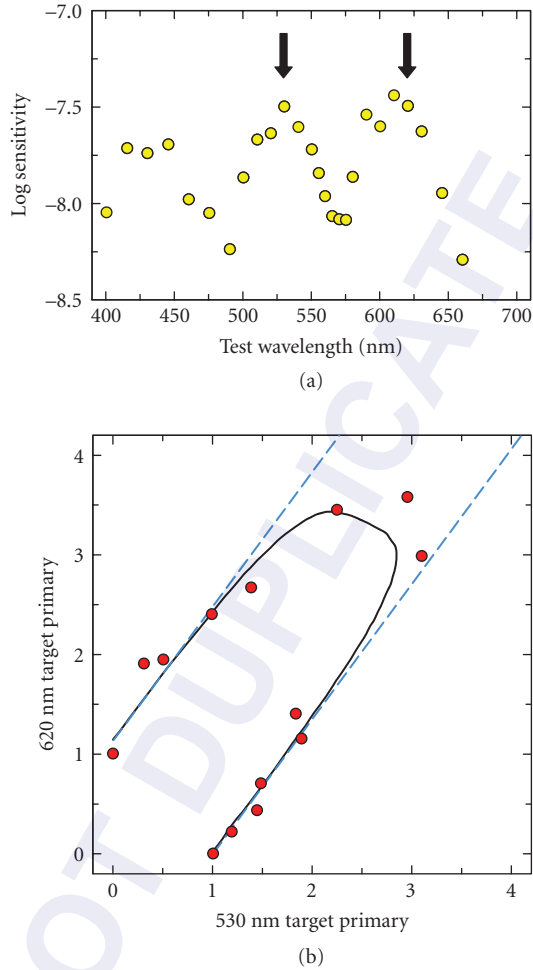


FIGURE 18 Test additivity. (a) Spectral sensitivity for the detection of a low frequency target measured on a 6000 td xenon (white) background (yellow circles). The shape is characteristic of detection by chromatic channels. The arrows indicate the wavelengths of the two primaries used in the test additivity experiment shown in (b). (b) Detection thresholds measured on the same field for different ratios of the two primaries. Data and model fits replotted from Fig. 2 of Thornton and Pugh.¹⁶⁴ The solid curve is the fit of a multi-mechanism model comprising an achromatic (L+M) and chromatic mechanisms (L-M and M-L with the same but opposite weights). The sensitivities of the chromatic mechanisms are shown by the dashed blue line. (See Ref. 164 for details.)

of the spectral position of the Sloan notch on background wavelength is discussed in subsection “Chromatic Adaptation and the Sloan Notch” in Sec. 11.5.

Luminance

Luminous efficiency The target parameters (and the experimental task) can also be chosen to favor detection by achromatic mechanisms. For example, changing the target parameters from a 200-ms duration, 1° diameter flash to a 10 ms, 0.05° diameter flash presented on a white background changes the spectral sensitivity from a chromatic one to an achromatic one, with chromatic peaks at approximately 530- and 610-nm merging to form a broad peak near 555 nm.¹⁰² Measurements that favor detection by achromatic mechanisms have been carried out frequently in the applied field of photometry as a means of estimating scotopic (rod), mesopic (rod-cone), or photopic (cone) “luminous efficiency” functions.

Luminous efficiency was introduced by the CIE (Commission Internationale de l’Éclairage) to provide a perceptual analog of radiance that could be used to estimate the visual effectiveness of lights. The measurement of luminous efficiency using tasks that favor achromatic detection is a practical solution to the requirement that luminous efficiency should be additive. Additivity, also known as obedience to Abney’s law,^{112,113} is necessary in order for the luminous efficiency of a light of arbitrary spectral complexity to be predictable from the luminous efficiency function for spectral lights, $V(\lambda)$.

Some of the earlier luminous efficiency measurements^{114–117} incorporated into the original 1924 CIE photopic luminous efficiency function, $V(\lambda)$, used techniques that are now known to fail the additivity requirement. Techniques that satisfy the additivity requirement, and depend mainly on achromatic or luminance mechanisms, include heterochromatic flicker photometry (HFP), heterochromatic modulation photometry (HMP), minimally distinct border (MDB), and minimum motion (MM). Details of these techniques can be found in several papers.^{28,53,90,118–124}

In visual science, $V(\lambda)$ or its variants, has often been assumed to correspond to the spectral sensitivity of the human postreceptoral achromatic or “luminance” mechanism, which is assumed to add positively weighted inputs from the L and M cones.¹²⁵ As a result, estimates of luminous efficiency have taken on an important theoretical role in studies of the opponent mechanisms, because they are used to produce stimulus modulations that produce no change in the purported luminance mechanism. In this context, the $V(\lambda)$ standard is often overinterpreted, since it obscures a number of factors that affect actual measurements of luminous efficiency. These include (1) the strong dependence of luminous efficiency on the mean state of chromatic adaptation (see subsection “Achromatic Detection and Chromatic Adaptation” in Sec. 11.5), (2) the sizeable individual differences in luminous efficiency that can occur between observers (see below), and (3) the fact that luminous efficiency is affected both by the spatial properties of the target, as well by where on the retina it is measured.

A more fundamental problem is that the CIE photopic 1924 $V(\lambda)$ function seriously underestimates luminous efficiency at short wavelengths, because of errors made in its derivation. Consequently, $V(\lambda)$ is seldom used in visual science, except for the derivation of troland values (with the result that luminance is often underestimated at short wavelengths). Attempts to improve $V(\lambda)$ ^{126,127} have been less than satisfactory,^{28,90} but are continuing in the form of a new estimate referred to as $V^*(\lambda)$.^{128,129} The $V^*(\lambda)$ (green line) and the CIE 1924 $V(\lambda)$ (green circles) luminous efficiency functions for centrally viewed fields of 2° in visual diameter, and the $V_{10}^*(\lambda)$ (orange line) and the CIE 1964 $V_{10}(\lambda)$ (orange squares) luminous efficiency functions for centrally viewed fields of 10° in visual diameter can be compared in Fig. 19. The functions differ mainly at short wavelengths. The $V^*(\lambda)$ and $V_{10}^*(\lambda)$ functions have been recommended by the CIE for use in “physiologically relevant” colorimetric and photometric systems.¹³⁰

In general, luminous efficiency functions can be approximated by a linear combination of the L-cone [$\bar{l}(\lambda)$] and M-cone [$\bar{m}(\lambda)$] spectral sensitivities, thus: $V(\lambda) = a\bar{l}(\lambda) + \bar{m}(\lambda)$, where a is the L-cone weight relative to the M-cone weight. Luminous efficiency shows sizeable individual differences even after individual differences in macular and lens pigmentation have been taken into account. For example, under neutral adaptation, such as daylight D65 adaptation, the L-cone weight in a group of 40 subjects, in whom 25-Hz HFP spectral sensitivity data were measured on a 3 log td white D65 background, varied from 0.56 to 17.75, while the mean HFP data were consistent with an L-cone weight of 1.89.^{128,129}

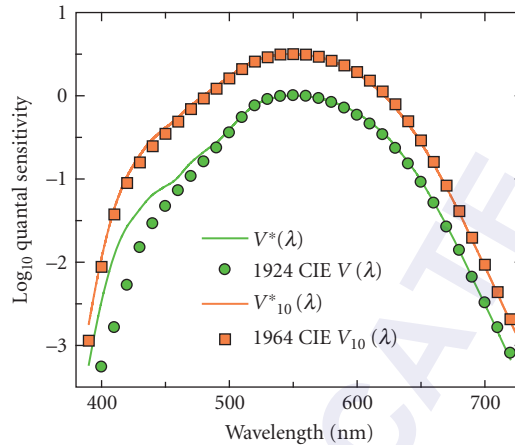


FIGURE 19 Luminous efficiency functions for 2° and 10° central viewing conditions. The lower pair of curves compare the CIE 1924 2° $V(\lambda)$ function (green circles) with the 2° $V^*(\lambda)$ function^{128,129} (green line). The upper pair of curves compare the CIE 1964 10° $V_{10}(\lambda)$ function (orange squares) with the 10° $V^*(\lambda)$ function (orange line). The main discrepancies are found at short wavelengths.

Cone numerosity On the assumption that the variation in L:M cone contribution to luminous efficiency is directly related to the variation in the ratio of L:M cone numbers in the retina, several investigators have used luminous efficiency measurements to estimate relative L:M cone numerosity.^{125,131–141} It is important to note, however, that luminous efficiency depends on factors other than relative L:M cone number, and is affected, in particular, by both first-site and second-site adaptation.^{63,64,142} Indeed, the L:M contribution to luminous efficiency could in principle have little or nothing to do with the relative numbers of L and M cones, but instead reflect the relative L- and M-cone contrast gains set in a manner independent of cone number. Such an extreme view seems unlikely to hold, however. Population estimates of relative L:M cone numerosity obtained using luminous efficiency, correlate reasonably well with estimates obtained using direct imaging,³⁰ flicker ERG,^{143,144} and other methods.^{134,137–139,145–148} More precise delineation of the relation between luminous efficiency and relative L:M cone numerosity awaits further investigation. Note that an important piece of this enterprise is to incorporate measurements of individual variation in L- and M-cone spectral sensitivity.^{128,144,146}

Multiple luminance signals As developed later in subsection “Multiplexing Chromatic and Achromatic Signals” in Sec. 11.5, cone-opponent, center-surround mechanisms that are also spatially opponent, such as parvocellular or P-cells, encode not only chromatic information, but also luminance information (see Fig. 38). The fact that P-cells transmit a luminance signal suggests that there may be two distinct “luminance” signals at the physiological level: one generated by P-cells and another by magnocellular or M-cells, a possibility which has been discussed in several papers by Ingling and his coworkers.^{149–152} The need for a P-cell-based luminance system is also suggested by the fact that the retinal distribution of M-cells is too coarse to support the observed psychophysical spatial acuity for luminance patterns¹⁵³ (but see Ref. 154). Others have also suggested that multiple mechanisms might underlie luminance spectral sensitivity.¹²⁴

Luminance mechanisms based on M-cell and P-cells should have different spatial and temporal properties. Specifically, the mechanism more dependent on M-cell activity, which we refer to as $(L+M)_M$, should be more sensitive to high temporal frequencies, whereas the mechanism more dependent on

P-cell activity, which we refer to as $(L+M)_p$, should be more sensitive to high spatial frequencies. If, as has been suggested,^{149–152} $(L+M)_p$ is multiplexed with L–M, we might also expect the interactions between $(L+M)_p$ and L–M to be different from those between $(L+M)_M$ and L–M. Similarly, chromatic noise may have a greater effect on $(L+M)_p$ than on $(L+M)_M$. Note, however, that if the stimulus temporal frequency is high enough for the spatially opponent surround to become synergistic with the center (i.e., when the surround delay is half of the flicker period), the P-cell will, in principle, become sensitive to any low spatial frequency luminance components of the stimulus.^{151,155}

Psychophysical evidence for multiple luminance mechanisms is relatively sparse. Given the likely difference in the M- and L-cone weights into $(L+M)_M$, the conventional luminance mechanism, and into L–M, the chromatic mechanism [from which $(L+M)_p$ is assumed to derive], the spectral sensitivities for spatial acuity tasks, which favor $(L+M)_p$, and for flicker tasks, which favor $(L+M)_M$, should be different. However, little difference in spectral sensitivity is actually found.^{122,156–158} Ingling and Tsou¹⁵² speculate that this similarity is found despite there being different luminance mechanisms, because the spatial acuity task may depend on P-cell centers rather than on centers and surrounds. Averaged over many cells, the P-cell *center* spectral sensitivity may be similar to the M-cell sensitivity, particularly if the latter depends on relative cone numerosity (see previous subsection).

Webster and Mollon¹⁷ looked at the effect of correlated chromatic and luminance contrast adaptation on four different measures of “luminous efficiency.” They found that lightness settings and minimum-motion settings for 1-Hz counter-phase gratings were strongly biased by contrast adaptation, but that flicker settings and minimum-motion settings for 15-Hz counter-phase gratings were unaffected. On the face of it, these results are consistent with the idea that different luminance mechanisms mediate low- and high-temporal frequency tasks.

We speculate that some of the discrepancies found between the detection contours measured in the L,M plane might be related to detection being mediated by different luminance mechanisms under different conditions. Conditions chosen to favor the conventional luminance channel (e.g., the presence of chromatic noise and/or the use of targets of high temporal frequency) are likely to favor detection by $(L+M)_M$. Perhaps, under these conditions the $(L+M)_M$ luminance contours appear as distinct segments, because there is relatively little threshold summation with L–M (e.g., $k = 4$). Examples of such segments are shown by the dashed lines in Fig. 20c. Conditions *not* specially chosen to favor the conventional luminance channel (e.g., the use of targets of high spatial and/or low temporal frequencies) may favor instead detection by $(L+M)_p$. Perhaps, under these conditions the luminance $(L+M)_p$ contours form elliptical contours with L–M ($k = 2$). Examples of such contours are shown in Fig. 20d. Noorlander, Heuts, and Koenderink,⁸³ however, reported that elliptical contours were found across all spatial and temporal frequencies.

Sensitivity to Different Directions of Color Space Before considering results obtained using trichromatic test mixtures, we first consider earlier results obtained using bichromatic mixtures.

Sensitivity to bichromatic test mixtures The use of spectral lights restricts the range of measurements to the spectrum locus in color space, thus ignoring the much larger and arguably environmentally more relevant volume of space that includes desaturated colors, achromatic lights, and extra-spectral colors (such as purples). Measuring detection sensitivity for mixtures of spectral targets probes the inner volume of color space. Studying these “bichromatic mixtures,” moreover, can provide a direct test of the additivity of a mechanism. For example, under conditions that isolate a univariant linear mechanism, pairs of lights should combine additively, so that the observer should be more sensitive to the combination than to either light alone. For a cone-opponent mechanism, on the other hand, pairs of lights that oppositely polarize the mechanism should combine subadditively, so that the observer is less sensitive to the mixture than to either component (as illustrated in Fig. 6d). If different mechanisms detect the two targets, then other potential interactions include probability summation or gating inhibition (winner takes all). Boynton, Ikeda, and Stiles¹⁰⁰ found a complex pattern of interactions between various test fields that included most of the possible types of interactions.

Several workers have demonstrated subadditivity for mixtures of spectral lights that is consistent with either L–M or S–(L+M) cone opponency.^{100,106,120,159–165} Clear examples of subadditivity were

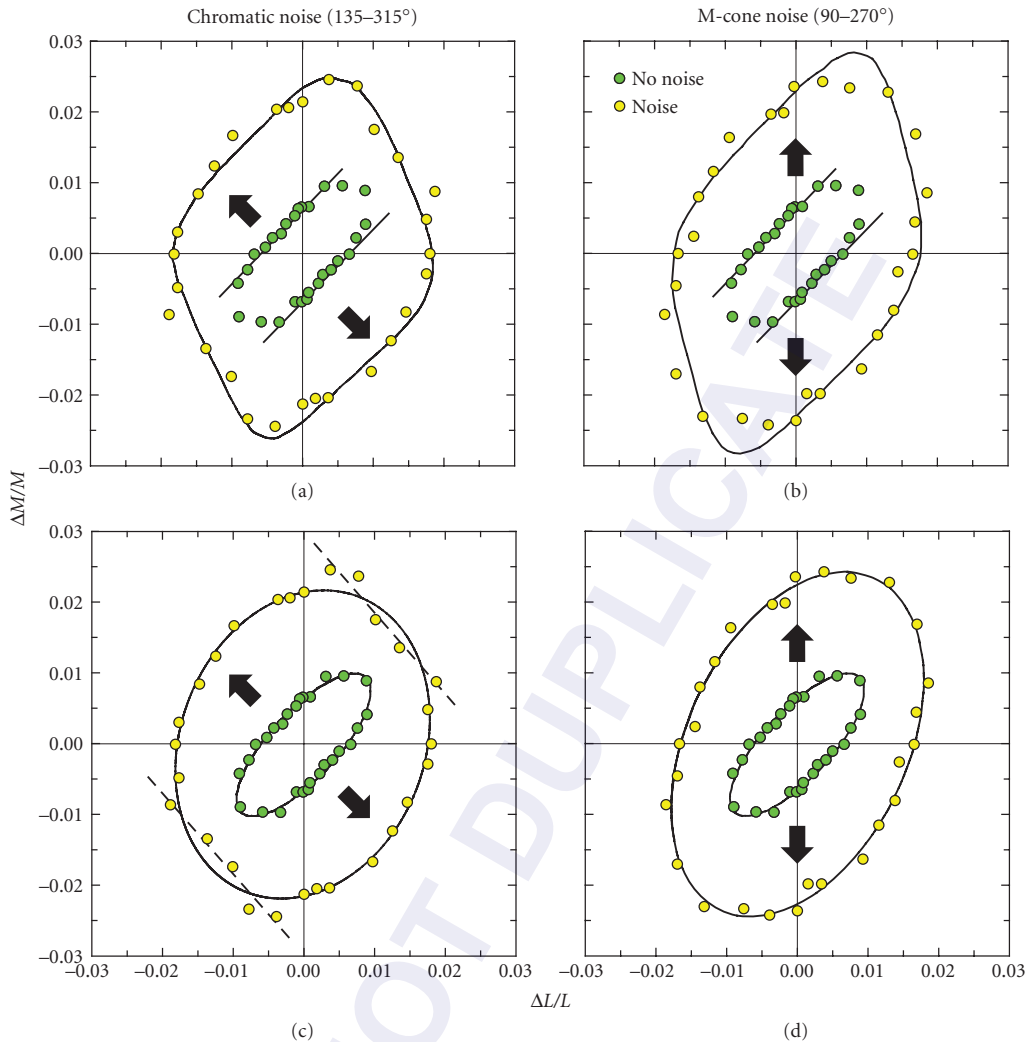


FIGURE 20 Detection contours in the $\Delta L/L$, $\Delta M/M$ plane. Detection thresholds for a 1-cpd Gabor replotted from Fig. 5 of Giulianini and Eskew¹⁷⁵ obtained with (yellow circles) or without (green circles) noise made up of superposed flickering rings modulated along the $\Delta L/L = -\Delta M/M$ (chromatic) axis (a and c) or the $\Delta M/M$ (M-cone) axis (b and d). The arrows indicate the noise directions. (a) and (b). Model fits (solid lines) by Giulianini and Eskew¹⁷⁵ in which L–M chromatic detection is assumed to be mediated by a mechanism with relative cone contrast weights of 0.70 $\Delta L/L$ and 0.72 $\Delta M/M$ of opposite sign, while L+M achromatic detection is mediated by a mechanism with relative weights of 0.90 $\Delta L/L$ and 0.43 $\Delta M/M$ of the same sign. Without chromatic noise, the detection data are consistent with detection solely by the chromatic mechanism (as indicated by the parallel lines) except along the luminance axis. With added chromatic noise, the sensitivity of the chromatic mechanism is reduced, revealing the achromatic mechanism in the first and third quadrants. In the presence of noise, the data are consistent with detection by both chromatic and achromatic mechanisms (as indicated by the rounded parallelogram). The combination exponent for the two mechanisms, k , was assumed to be 4 (see subsection “How Test Measurements Imply Opponency” in Sec. 11.3). (c) and (d) Alternative fits to the data using elliptical contours, which provide a surprisingly good fit to the detection data, except in the achromatic direction with chromatic noise (c). The ellipse fitted to the yellow circles in (c) has the formula $3332x^2 + 2021y^2 + 873xy = 1$ (major axis 73.16°), while that fitted to the green circles has the formula $27,191x^2 + 24,842y^2 + 43,005xy = 1$ (major axis 46.56°), where $x = \Delta L/L$ and $y = \Delta M/M$. Those data with noise that are poorly accounted for by the ellipse have been fitted with a straight line of best-fitting slope -1.14 (dashed line). The ellipse fitted to the yellow circles in (d) has the formula $3648x^2 + 1959y^2 + 1992xy = 1$ (major axis 65.15°), while that fitted to the green circles has the formula $26,521x^2 + 23,587y^2 + 38,131xy = 1$ (major axis 47.20°).

obtained by Thornton and Pugh,¹⁶⁴ who presented targets chosen to favor chromatic detection on a white (xenon-arc) background. Figure 18a shows the test spectral sensitivity obtained on the white field characteristic of chromatic detection. The spectral positions of the test mixture primaries used in the additivity experiment are indicated by the arrows. Figure 18b shows the results of the text mixture experiment. Detection contours aligned with the parallel diagonal contours (dashed blue lines) is consistent with chromatic detection by an L–M mechanism, whereas detection along the 45° vector (parallel to the contours) is consistent with achromatic detection by an L+M mechanism. Thornton and Pugh¹⁶⁴ also used test primaries of 430 and 570 nm, which were chosen to favor detection by the S cones and by luminance (L+M), respectively. Although less clear than the L–M case shown here, they found evidence for inhibition between the 430- and 570-nm target consistent with detection by an S–(L+M) cone-opponent mechanism.

Guth and Lodge¹¹ and Ingling and Tsou¹² both used models that incorporated opponency to analyze bichromatic and spectral thresholds. Both found that the threshold data could be accounted for by two opponent-color mechanisms [which can be approximated to the basic cone-opponent discrimination mechanisms, L+M and S–(L+M)], and a single nonopponent mechanism (L+M). Both also developed models that could better account for “suprathreshold” data, such as color valence data, by modifying their basic models. These modifications took the form of an increase in the contribution of the S–(L+M) mechanism and either inhibition between the two cone-opponent mechanisms¹¹ or addition of an S-cone contribution to the L-cone side of the L–M mechanism.¹² These suprathreshold models are related to the color-appearance models discussed in subsection “Color Appearance and Color Opponency” in Sec. 11.5. It is worth noting that Guth and Lodge acknowledge that modulations of their cone-opponent discrimination mechanisms do not produce strictly red–green or yellow–blue appearance changes, but rather yellowish–red versus blue–green and greenish–yellow versus violet. That is, their threshold discrimination mechanisms do not account for unique hues (see subsection “Spectral Properties of Color-Opponent Mechanisms” in Sec. 11.5).

Kranda and King-Smith¹⁰⁶ also modeled bichromatic and spectral threshold data, but they required four mechanisms: an achromatic mechanism (L+M), two L–M cone-opponent mechanisms of opposite polarities (L–M and M–L), and an S-cone mechanism (S). The main discrepancy from the previous models is that they found no clear evidence for an opponent inhibition of S by L+M. Two opposite polarity, L–M mechanisms were required because Kranda and King-Smith assumed that only the positive lobe contributed to detection, whereas previous workers had assumed that both negative and positive lobes could contribute. This presaged the unipolar versus bipolar mechanism debate (see “Unipolar versus Bipolar Chromatic Mechanisms” in Sec. 11.6).

With the exception of Boynton, Ikeda, and Stiles,¹⁰⁰ most of the preceding studies used incremental spectral test lights, which confined the measurements to the eighth volume of color space in which all three cone signals increase. To explore other directions in color space, both incremental and decremental stimuli must be used.

Detection contours in the L,M plane The first explicit use of incremental and decremental stimuli in cone contrast space was by Noorlander, Heuts, and Koenderink,⁸³ who made measurements in the L,M plane. They fitted their detection contours with ellipses, and found that the alignment of the minor and major axes depended on temporal and spatial frequency. At lower frequencies, the minor and major axes aligned with the chromatic ($\Delta M/M = -\Delta L/L$) and achromatic ($\Delta M/M = \Delta L/L$) directions, respectively. By contrast, at higher frequencies, the reverse was the case. Such a change is consistent with relative insensitivity of chromatic mechanisms to higher temporal and spatial frequencies (see subsection “Spatial and Temporal Contrast Sensitivity Functions” in Sec. 11.5).

If the individual detection contours are truly elliptical in two-dimensional color space or ellipsoidal in three, then the vectors of the underlying linear detection mechanisms cannot be unambiguously defined. This ambiguity arises because ellipses or ellipsoids can be linearly transformed into circles or spheres, which do not have privileged directions.¹⁶⁵ Consequently, whether or not detection contours are elliptical or ellipsoidal has become an important experimental question. Knoblauch and Maloney¹⁶⁶ specifically addressed this question by testing the shapes of detection contours measured in two color planes. One plane was defined by modulations of the red and green phosphors of their display (referred to by them as the $\Delta R, \Delta G$ plane), and the other by the modulations of the red and green phosphors together and the blue phosphor of their display (referred to as the $\Delta Y, \Delta B$ plane). The

detection contours they obtained did not differ significantly in shape from an elliptical contour.¹⁶⁶ Note, however, that a set of multiple threshold contours, each ellipsoidal but measured under related conditions, can be used to identify mechanisms uniquely, given a model about how the change in conditions affects the contours (see, e.g., Ref. 167; discussed in subsection “Spatial and Temporal Contrast Sensitivity Functions” in Sec. 11.5).

Deviations from an ellipse, if they do occur, are most likely to be found under experimental conditions that maximize the probability that different portions of the contour depend on detection by different mechanisms. Such a differentiation is unlikely to be found under the conditions used by Knoblauch and Maloney.¹⁶⁶ At lower temporal frequencies (1.5-Hz sine wave or a Gaussian-windowed, $\sigma = 335$ ms, 1.5-Hz cosine wave), the detection contour in the ΔR , ΔG phosphor space is likely to be dominated by the much more sensitive L–M mechanisms, whereas the contour in ΔY , ΔB space is likely to depend upon detection by all three mechanisms [L–M, L+M, and S–(L+M)]. A comparable conclusion about the shape of threshold contours was reached by Poirson, Wandell, Varner, and Brainard,¹⁶⁸ who found that color thresholds could be described well by ellipses, squares, or parallelograms. Again, however, they used long duration (a Gaussian envelope lasting 1 s, corresponding to +5 to –5 standard deviations) and large (2° diameter) target that would favor detection by L–M.

The relatively high sensitivity of the L–M mechanism to foveally viewed, long-duration flashed targets^{169–171} means that the L–M mechanism dominates the detection contours unless steps are taken to expose the L+M luminance mechanisms, such as by the use of flickering targets,^{63,172} very small targets,¹⁷³ or chromatic noise (see subsection “Noise-Masking Experiments” in Sec. 11.5). Detection threshold contours measured in the L,M plane that expose both L–M and L+M contours are arguably better described by a parallelogram with rounded corners rather than an ellipse, with the parallel sides of positive slope corresponding to detection by the L–M chromatic mechanism and those of negative slope corresponding to detection by the L+M luminance mechanism.^{61,63,173,174} Examples of detection contours that appear nonelliptical include parts of Figs. 3 and 4 of Ref. 63; Fig. 3, 5, and 7 of Ref. 175; and Fig. 2 of Ref. 176.

Figure 20 shows the detection data replotted from Fig. 5 of Ref. 175. The data are shown twice: once in the upper panels (a, b) and again in the lower panels (c, d). The data shown as green circles were measured without noise. The data shown as yellow circles in Fig. 20a and c were measured in the presence of chromatic noise, while those in Fig. 20b and d were measured in the presence of M-cone noise. The noise directions are indicated by the black arrows.

If, for the sake of argument, we accept that the mechanism vectors of the detection mechanisms can be unambiguously defined from detection contours (as illustrated in Fig. 20a and b), it is possible to estimate their cone weights and therefore their spectral sensitivities. The chromatic L–M detection contours in Fig. 20a and b obtained without noise (green circles) have been fitted with lines of parallel slopes of approximately 1.0, which indicates that the L–M chromatic mechanism responds to the linear difference of the L- and M-cone *contrast* signals; $|c\Delta L/L - d\Delta M/M| = \text{constant}$, where $c = d$.^{61,63,173,174} The equality of the L- and M-cone weights, unlike the underlying relative L- and M-cone numerosities (e.g., Ref. 30), shows a remarkable lack of individual variability across observers (see Ref. 89). Similarly, color appearance settings (see subsection “Color Appearance and Color Opponency” in Sec. 11.5) in the form of unique yellow settings show much less variability than would be expected if unique yellow depended on the relative number of L and M cones in the retina.¹⁴⁶

As suggested by the yellow circles in Fig. 20a and c, the use of a high frequency flickering target in chromatic noise exposes more of the L+M contour in the first and third quadrants. Unlike the L–M contour, the L+M contour has a negative slope, which indicates that the L+M achromatic chromatic mechanism responds to the linear sum of the L- and M-cone *contrast* signals; $|a\Delta L/L + b\Delta M/M| = \text{constant}$, where a is typically greater than b .^{61,63,89,173,174}

Mechanism interactions The shapes of the detection contours when two or more postreceptoral mechanisms mediate threshold depends upon the way in which those mechanisms interact (see Fig. 6c). In principle, such interactions can range from complete summation, through probability summation to exclusivity (winner takes all) and inhibition (see Ref. 100). The usual assumption is that L–M, S–(L+M), and L+M are stochastically independent mechanisms that combine by probability summation. That is, if two or three mechanisms respond to the target, then the threshold will be lower than if only one responds to it. How much the threshold is lowered depends upon the

underlying frequency-of-seeing curve for each mechanism. Steep frequency-of-seeing curves result in a relatively small drop in threshold, whereas shallow curves result in a larger drop. Several groups have mimicked probability summation by using the summation rule introduced in “How Test Measurements Imply Opponency” in Sec. 11.3. The summation exponent can be varied to mimic the effects of probability summation for psychometric functions with different slopes (see Refs. 174 and 89). A closely related geometrical description of the threshold contours was used by Sankeralli and Mullen.¹⁷² For the fits in Fig. 20*a* and *b*, the summation exponent was assumed to be 4, which produces a contour shaped like a rounded parallelogram.^{89,175}

The superiority of a rounded parallelogram over an ellipse (summation exponent of 2) is not always clear. For example, Fig. 20*c* and *d* show the same detection data as the *a* and *b*. In *c* and *d*, the data have been fitted with ellipses. As can be seen, each set of data is well fitted by an ellipse, except perhaps for a small part of the contour measured in the presence of chromatic noise shown in Fig. 20*c*. In this region, we have fitted a straight line plotted as the dashed line. Changing the direction of the noise from L–M chromatic to M cone between Fig. 20*c* and *d* rotates the ellipse slightly anticlockwise. The L–M data along the opponent diagonals clearly follow an elliptical rather than a straight contour in both cases. As previously noted, when an elliptical representation underlies the data, it does not uniquely determine the inputs to the underlying mechanisms.

Support for the summation exponent of 4 used to fit the data in Fig. 20*a* and *b* is also lacking from estimates of the slopes of the underlying psychometric functions, which are consistent with much smaller exponents.^{171,177,178}

Detection in the other planes An investigation of mechanisms other than L–M and L+M (as well as information about whether or not there are S-cone inputs into those mechanisms) requires the measurement of responses outside the L,M plane of color space. Two systematic investigations have been carried out that include S-cone modulations. Cole, Hine, and McIlhagga¹⁷⁴ made measurements on a roughly 1000 log td white background using a Gaussian-blurred 2°, 200-ms spot that favors chromatic detection. Their data were consistent with three independent mechanisms: an L–M mechanism with equal and opposite L- and M-cone inputs but no S-cone input; an L+M mechanisms with unequal inputs and a small positive S-cone input, and an S–(L+M) mechanism. Sankeralli and Mullen¹⁷² used three types of gratings with spatiotemporal properties chosen to favor either the chromatic L–M (1 cpd, 0 Hz), chromatic S–(L+M) (0.125 cpd, 0 Hz), or the achromatic L+M mechanism (1 cpd, 24 Hz). They confirmed that L–M had equal and opposite L- and M-cone inputs, but found a small S-cone input of about 2 percent added to either the L or the M cones for different observers. In addition, they found a small 5 percent S-cone input into the luminance mechanism that opposed the L+M input (see also Ref. 92). The S–(L+M) mechanism was assumed to have balanced opposed inputs.¹⁷² Eskew, McLellan, and Giulianini⁸⁹ provide a separate analysis of both these sets of data (see Table 18.1 of Ref. 89).

Figure 21 shows chromatic detection data measured in the L–M,S plane by Eskew, Newton, and Giulianini.¹⁷⁹ The central panel shows detection thresholds with (filled circles) and without (open squares) L–M chromatic masking noise. Detection by the relatively insensitive S-cone chromatic mechanisms becomes apparent when detection by L–M chromatic mechanism is prevented by chromatic noise. The nearly horizontal portions of the black contour reflect detection by S–(L+M), while the vertical portions reflect detection by L–M (the model fit is actually of the unipolar mechanisms $|L-M|$, $|M-L|$, $|S-(L+M)|$, and $|(L+M)-S|$, see Ref. 34 for details; also “Unipolar versus Bipolar Chromatic Mechanisms” in Sec. 11.6).

The question of whether or not the S cones make a small contribution to the L–M detection mechanisms remains controversial. Boynton, Nagy, and Olsen¹⁸⁰ found that S cones show summation with L–M (+S with L–M, and –S with M–L) in chromatic difference judgments. Stromeyer et al.¹⁸¹ showed that S-cone flicker facilitated detection and discrimination of L–M flicker, with +S facilitating L–M and –S facilitating M–L. Yet, the estimated S-cone contrast weight was only approximately 1/60th of the L- and M-cone contrasts weight into L–M (see p. 820 of Ref. 181). In general, in detection measurement the S-cone signal adds to redness (L–M) rather than greenness (M–L)—as expected from color-opponent theory (see Sec. 11.4), but the size of the contribution is much less than expected. Eskew and Kortick,²⁰ however, made both hue equilibria and detection measurements and estimated that the S-cone contrast weight was about 3 percent that of the L- and

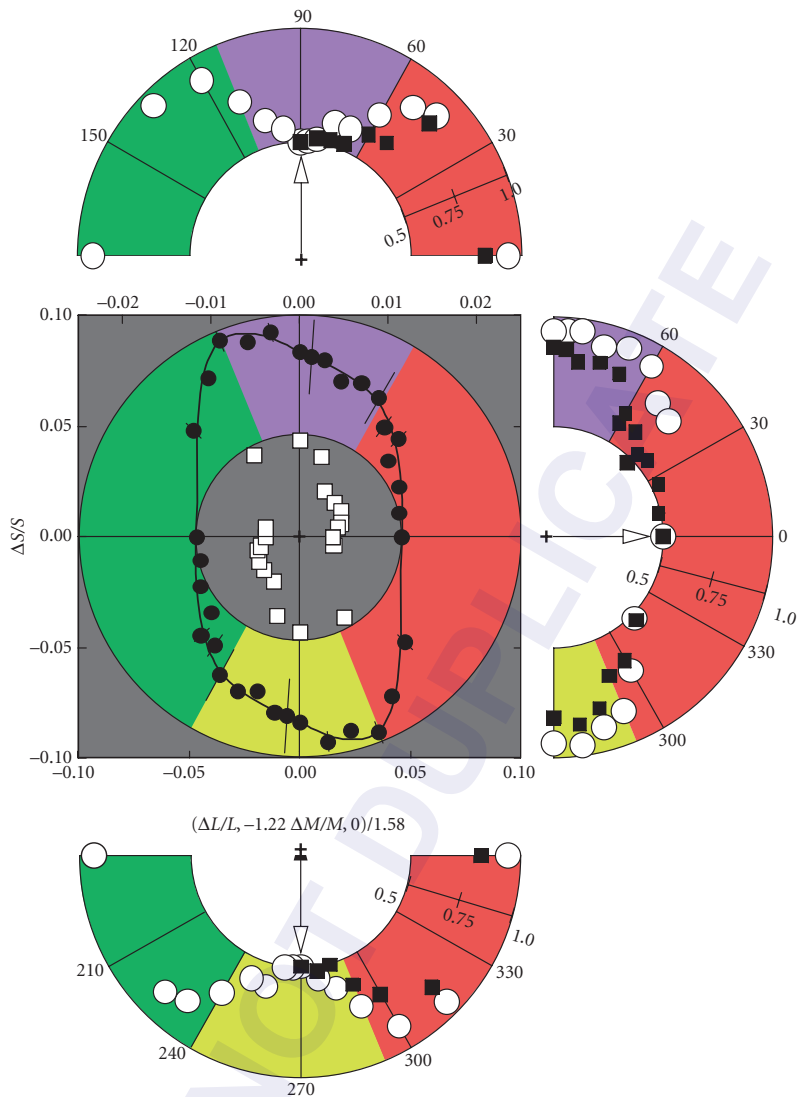


FIGURE 21 Detection and discrimination contours in the equiluminant plane. This figure is a reproduction of Fig. 5 from Eskew.³⁴ The central panel shows thresholds for the detection of Gabor patches with (filled circles) and without (open squares) L–M chromatic masking noise (0°/180° color direction). For clarity, the horizontal scale for the no-noise thresholds (open squares) has been expanded as indicated by the upper axis. Data are symmetric about the origin due to the symmetric stimuli. Detection by the relatively insensitive S-cone chromatic mechanisms becomes apparent when detection by the L–M chromatic mechanism(s) is suppressed by L–M chromatic noise. The black contour fitted to the filled circles is the prediction of a model in which detection is mediated by probability summation between four symmetric unipolar linear mechanisms $|L-M|$, $|M-L|$, $|S-(L+M)|$, and $|(L+M)-S|$, (see subsection “Unipolar versus Bipolar Chromatic Mechanisms” Sec. 11.6). The three outlying semicircular polar plots are discriminability data (filled squares), in which the angular coordinate corresponds to the same stimuli, at the same angles, as in the detection data. The radial coordinate gives the discriminability between 0.5 (chance) and 1.0 (perfect discriminability) of the given test color angle relative to the standard color angle indicated by the arrow. The open circles are predictions of a Bayesian classifier model that takes the outputs of the four mechanisms fitted to the detection data as its inputs. Colored regions indicate bands of poorly discriminated stimuli, and are redrawn on the detection plot for comparison. These results suggest that just four univariant mechanisms are responsible for detection and discrimination under these conditions. See Ref. 34 for further details. [Figure adapted by Eskew³⁴ from Fig. 4 (Observer JRN) of Ref. 179 Reprinted from *The Senses: A Comprehensive Reference, Volume 2: Vision II*, Fig. 5 from R. T. Eskew, Jr., “Chromatic Detection and Discrimination,” pp. 109–117. Copyright (2008), with permission from Elsevier.]

M-cone contrast weights for both tasks. Yet, although test methods support an S-cone contribution to L–M, field methods do not (see Ref. 34). Changes in the S-cone adaptation level have little effect on L–M detection¹⁸² or on L–M mediated wavelength discrimination.¹⁸³

Given the uncertainties about cone spectral sensitivities and luminous efficiency at short wavelengths, as well as the substantial individual differences in prereceptor filtering between observers (see, e.g., Ref. 90), evidence for small S-cone inputs should be treated with caution.

Spatial and Temporal Contrast Sensitivity Functions The test method can also be used to determine the temporal and spatial properties of visual mechanisms. These properties are traditionally characterized by temporal and spatial “contrast sensitivity functions” (CSFs). These are sometimes incorrectly called “modulation transfer functions” (MTFs)—incorrectly because MTFs require phase as well as amplitude specifications.

In a temporal CSF determination, the observer’s sensitivity for detecting sinusoidal flicker is measured as a function of temporal frequency. Generally, the chromatic mechanisms have a lowpass temporal frequency response with typically greater temporal integration and a poorer response to higher temporal frequencies than the more band-pass achromatic mechanism.^{94,102,184–190} A lowpass CSF means that the temporal response is approximately independent of frequency at low frequencies but falls off at higher frequencies, whereas a band-pass CSF means that the response peaks at an intermediate frequency and falls off at both lower and higher frequencies.

On the grounds that different chromatic flicker frequencies can be discriminated near threshold, Metha and Mullen¹⁹⁰ suggest that at least two chromatic mechanisms with different temporal CSFs must underlie the measured chromatic CSE, one with a band-pass CSF and the other with a lowpass CSF. They argue that two mechanisms must operate, because a single *univariant* flicker mechanism will confound changes in contrast and flicker frequency, making frequency identification impossible. However, the assumption that flicker is encoded univariantly may be flawed. At low frequencies, flicker may be encoded as a moment-by-moment variation that matches the flicker frequency.¹⁹¹

In a spatial CSF determination, the observer’s sensitivity for detecting sinusoidal gratings is measured as a function of their spatial frequency. Generally, the chromatic mechanisms have a lowpass spatial frequency response with greater spatial integration and a poorer response to higher spatial frequencies than the usually more band-pass achromatic mechanisms.^{102,192–197}

Spatial CSFs are degraded by sensitivity losses introduced by the optics of the eye. In one study, laser interference fringes were used to project gratings directly on the retina, thereby effectively bypassing the MTF of the eye’s optics.^{198,199} However, because the red and green equiluminant interference fringes could not be reliably aligned to produce a steady artifact-free chromatic grating, they were drifted in opposite directions at 0.25 Hz. Thus, the gratings continuously changed from a red-green chromatic (equiluminant) grating to a yellow-black luminance (equichromatic) grating. The subject’s task was to set the contrast threshold separately for the different spatial phases, which was apparently possible.¹⁹⁹ One concern about this method, however, is that the close juxtaposition of chromatic and luminance gratings in time and space might produce facilitation or interference (see subsection “Pedestal Experiments” in Sec. 11.5). For S-cone detection, a single grating was used.

The results for equiluminant, chromatic (red circles) and equichromatic, achromatic (blue circles) fringes are shown for three observers in the upper panels of Fig. 22. As expected, there is a steeper fall-off in high spatial-frequency sensitivity for chromatic gratings than for luminance gratings. By estimating the sensitivity losses caused by known factors such as photon noise, ocular transmission, cone-aperture size, and cone sampling, Sekiguchi, Williams, and Brainard¹⁹⁸ were able to estimate an ideal observer’s sensitivity at the retina. Any additional losses found in the real data are then assumed to be due to neural factors alone. Estimates of the losses due to neural factors are shown in the lower panels of Fig. 22. Based on these estimates, the neural mechanism that detects luminance gratings has about 1.8 times the bandwidth of the mechanisms that detect chromatic gratings. Interestingly, the neural mechanisms responsible for the detection of equiluminant L–M gratings and S-cone gratings have strikingly similar (bandpass) contrast sensitivities, despite the very different spatial properties of their submosaics.²⁰⁰ This commonality may have important implications for color processing (see “Three-Stage Zone Models” in Sec. 11.6).

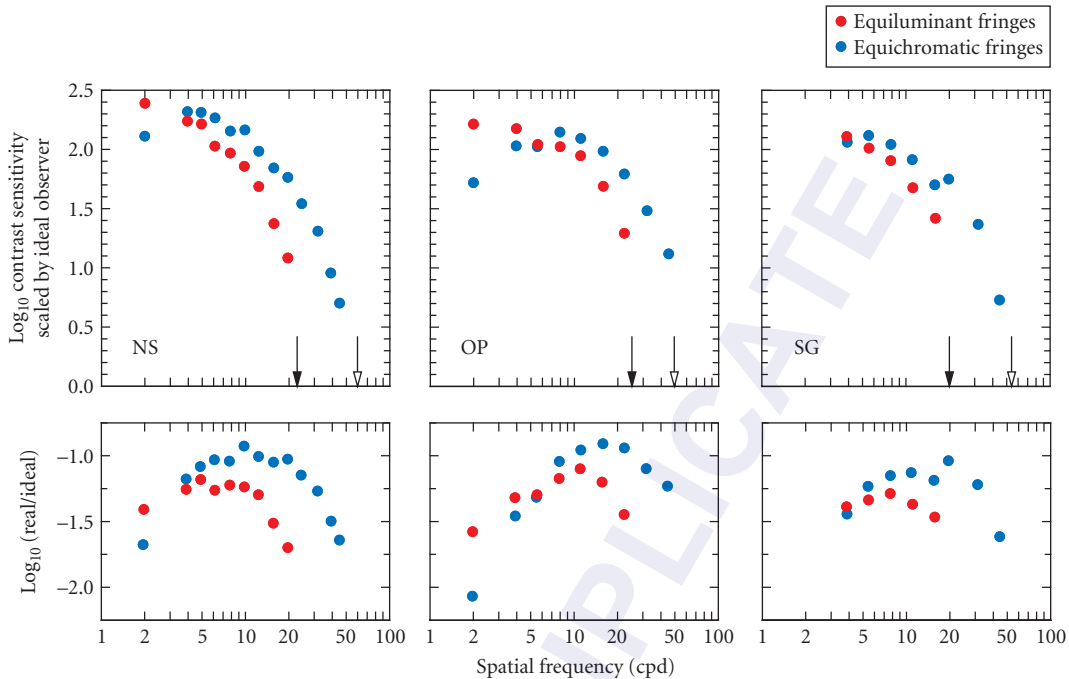


FIGURE 22 Spatial equiluminant and equichromatic contrast sensitivity functions. Top panels: Foveal CSFs for equiluminant (red circles) and equichromatic (blue circles) interference fringes for three observers. Filled and open arrows represent the foveal resolution limit for the equiluminant and the equichromatic stimuli, respectively. Bottom panels: The ratio of the ideal to the real observer's contrast sensitivity for detecting equiluminant (red circles) and equichromatic (blue circles) stimuli. (Based on Figure 9 of Ref. 199.)

Poirson and Wandell¹⁶⁷ measured spatial CSFs along the two cardinal directions and along intermediate directions in color space; that is, they effectively measured threshold contours for a series of grating spatial frequencies. They fitted the contours with an ellipsoidal model, but added a constraint of pattern-color separability; that is, they assumed that the effect of pattern operated independently in each of three second-site mechanisms. This constraint resolves the mechanism ambiguity inherent in fitting ellipsoidal contours to the individual data, and allowed Poirson and Wandell to estimate the sensitivities of three underlying mechanisms by jointly analyzing the contours across spatial frequency. They found two opponent and one nonopponent mechanisms, with sensitivities consistent with the cardinal mechanisms, L+M, L-M, and S-(L+M) (see Fig. 6 of Ref. 167).

We next consider results obtained using the field method.

Field Sensitivities

Stiles' π -Mechanisms Stiles used the field sensitivity method to define the spectral sensitivities of the π -mechanisms, which are illustrated in Fig. 23. The tabulated field sensitivities can be found in Table 2 (7.4.3) of Ref. 53 and Table B of Ref. 33. The field sensitivities are the reciprocals of the field radiances (in log quanta $s^{-1} \text{ deg}^{-2}$) required to raise the threshold of each isolated π -mechanism by one log unit above its absolute threshold. The tabulated data are average data for four subjects: three females aged 20 to 30 and one male aged 51.

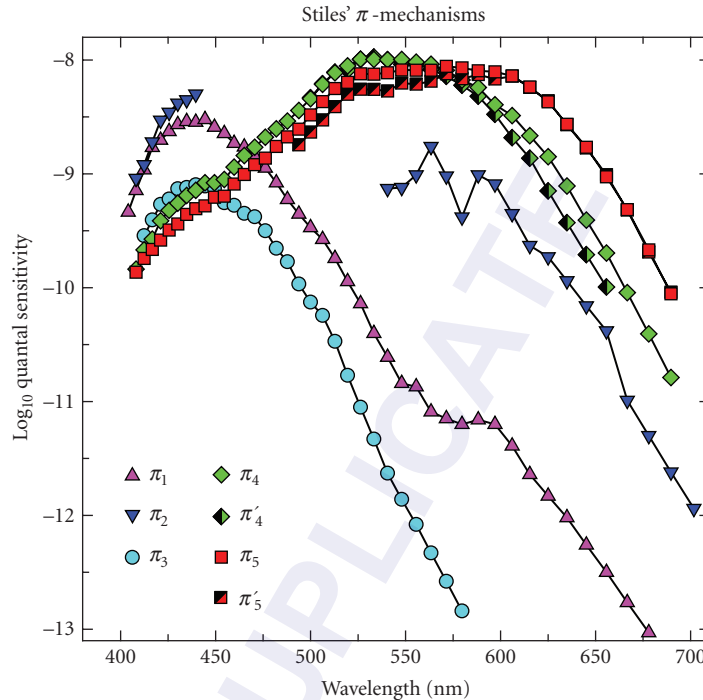


FIGURE 23 Stiles' π -mechanisms. Field spectral sensitivities of the seven photopic π -mechanisms: π_1 (purple triangles), π_2 (dark blue inverted triangles), π_3 (light blue circles), π_4 (green diamonds), π_4' (green half-filled diamonds), π_5 (red squares), π_5' (red half-filled squares). The field sensitivities are the reciprocals of the field radiances (in log quanta $\text{s}^{-1} \text{deg}^{-2}$) required to raise the threshold of each "isolated" π -mechanism by one log unit above its absolute threshold. The results were averaged across four subjects: three females aged 20 to 30 and one male aged 51. Details of the experimental conditions can be found on p. 12 of Stiles' book.³³ (Data from Table 2(7.4.3) of Wyszecki and Stiles.⁵³)

Through field sensitivity measurements, and measurements of threshold versus background radiance for many combinations of target and background wavelength, Stiles identified seven photopic π -mechanisms: three predominately S-cone mechanisms (π_1 , π_2 , and π_3), two M-cone (π_4 and π_4'), and two L-cone (π_5 and π_5'). Although it has been variously suggested that some of the π -mechanisms might correspond to cone mechanisms,^{97,201–203} it is now clear that all reflect some form of postreceptoral cone interaction. The surfeit of π -mechanisms and their lack of obvious correspondence to cone mechanisms or simple postreceptoral mechanisms led to the demise of Stiles' model. Its usefulness, however, was prolonged in the work of several investigators, who used failures of Stiles' model as a way of investigating and understanding the properties of chromatic mechanisms (see "Field Additivity" in this section). But, as noted in the parallel discussion in the context of test sensitivity, the potential ongoing value of the π -mechanism field sensitivities is that they provide an empirical database that may be used to constrain current models.

Achromatic Detection and Chromatic Adaptation Achromatic spectral sensitivity (or luminous efficiency) is strongly dependent on chromatic adaptation.^{63,64,142,204–210} Figure 24a shows the changes in luminous efficiency caused by changing background field wavelength from 430 to 670 nm and Fig. 24b shows the changes caused by mixing 478- and 577-nm adapting fields in different luminance

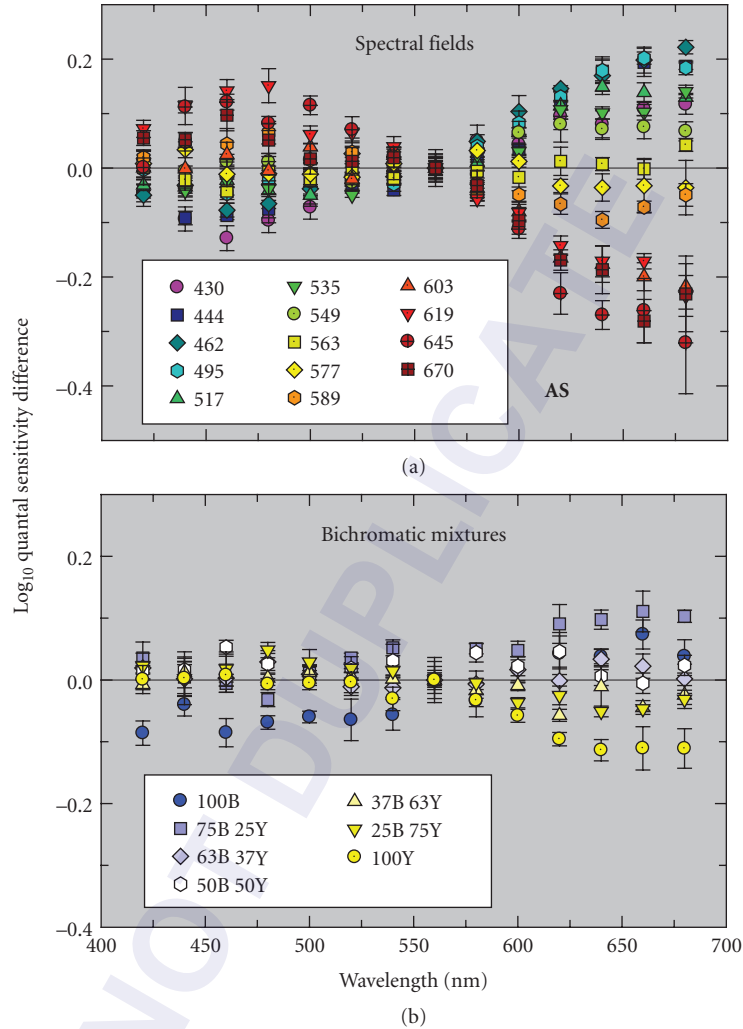


FIGURE 24. Dependence of luminous efficiency on chromatic adaptation. Changes in quantal luminous efficiency caused by changes in adapting field chromaticity plotted relative to the mean luminous efficiency. Data from Observer S1 of Stockman, Jägle, Pirzer and Sharpe.¹²⁹ Measurements were made using 25-Hz heterochromatic flicker photometry (HFP) on 1000-td backgrounds. (a) HFP changes on fourteen 1000-td spectral adapting fields from 430 to 670 nm and (b) changes on seven 1000 td bichromatic, 478(B) + 577(Y)-nm adapting field mixtures that varied from 100 percent B to 100 percent Y.

ratios. These changes are consistent with selective adaptation reducing L-cone sensitivity relative to M on longer wavelength fields, and reducing M-cone sensitivity relative to L on shorter wavelength fields. Between about 535 and 603 nm the selective changes are roughly consistent with Weber's law ($\Delta I/I = \text{constant}$). Consequently, plotted in cone contrast space, the sensitivity contour for the achromatic mechanism should remain of constant slope between field wavelengths of 535 and 603 nm, *unlike* the example shown in Fig. 8b. For background wavelengths shorter than 535 nm,

on the other hand, the changes in spectral sensitivity are consistent with a relative loss of M-cone sensitivity in excess of that implied by Weber's law (for details, see Ref. 129). This is consistent with other evidence that also shows that the luminance contribution of the L- or M-cone type more sensitive to a given chromatic field can be suppressed by more than Weber's law implies.^{63,64,142} Such supra-Weber suppression, illustrated in Fig. 8*b* for the L-cone case, causes the luminance sensitivity contours in cone contrast space to rotate away from the axis of the more suppressed cone.

Multiple cone inputs There is now good psychophysical evidence that the standard model of the luminance channel as an additive channel with just fast L- and M-cone inputs (referred to here as $+fL+fM$) is incomplete. According to this model, pairs of sinusoidally alternating lights that are "luminance-equated" and detected solely by the luminance channel should appear perfectly steady or nulled whatever their chromaticities (equating alternating lights to produce a null is known as heterochromatic flicker photometry or HFP). And, indeed, such nulls are generally possible, but sometimes only if moderate to large phase delays are introduced between the pairs of flickering lights.^{184,211–215} Moreover, these large phase delays are often accompanied by substantial frequency-dependent changes in the spectral sensitivity of flicker detection and modulation sensitivity.^{216–218} From the phase delays and sensitivity changes, it is possible to infer that, in addition to the faster $+fM+fL$ signals, sluggish $+sM-sL$, $+sL-sM$, and $-sS$ signals also contribute to the luminance-signal nulls in a way that depends on adapting chromaticity and luminance.^{94,215–223}

Figure 25 shows examples of large phase delays that are found for M- and L-cone-detected flickering lights on four increasingly intense 658 nm fields. In these plots, 0° means that the two lights cancelled when they were physically in opposite phase (i.e., when they were alternated), while $\pm 180^\circ$ means that they cancelled when they were actually in the same phase. The differences between the plotted M- and L-cone phase delays therefore show the delays between the M- and L-cone signals introduced *within* the visual system. As can be seen, some of the phase delays are substantial even at moderately high temporal frequencies, particularly for the M-cone signals. Such delays are inconsistent with the standard model of luminance, which, except for phase differences caused by the selective adaptation of the L cones by the long-wavelength field, predicts that no phase adjustments should be required. Notice also the abrupt changes in phase *between* intensity levels 3 and 4.

The continuous lines in Fig. 25 are fits of a vector model, in which it was assumed that sluggish and fast signals both contribute to the resultant M- and L-cone signals, but in different ratios depending on the cone type and intensity level. At levels 1 to 3, the dominant slow and fast signals are $+sL-sM$ and $+fL+fM$, respectively, as illustrated in the lower neural circuit diagram of Fig. 25, whereas at level 4 they are $-sL+sM$ and $+fL+fM$, as illustrated in the upper neural circuit diagram. According to the model, the polarities of the both the slow M- and the slow L-cone signals reverse between levels 3 and 4.

Note that although the sluggish signals $+sM-sL$ and $+sL-sM$ are spectrally opponent, they produce an achromatic percept that can be flicker photometrically cancelled with luminance flicker, rather than producing an R/G chromatic percept.

Chromatic Adaptation and the Sloan Notch Thornton and Pugh¹⁰⁷ measured test spectral sensitivity functions on fields of 560, 580, and 600 nm using a target that strongly favors chromatic detection. As illustrated in Fig. 17*a*, they found that the local minimum in spectral sensitivity, or Sloan notch, coincides with the field wavelength (i.e., it occurs when the target and background wavelengths are the same; i.e., homochromatic). As discussed in subsection "Sensitivity to Spectral Lights" in Sec. 11.5, the notch is thought to correspond to the target wavelength that produces a minimum or null in the L–M chromatic channel, so that the less-sensitive, achromatic channel takes over detection. For the notch to occur at the homochromatic target wavelength means that adaptation to the field has shifted the zero crossing or null in the cone-opponent mechanism to the field wavelength. Such a shift is consistent with reciprocal (von Kries) adaptation occurring independently in both the M and the L cones. This type of adaptation is plausible given that Weber's law behavior for detection will have been reached on the fairly intense fields used by Thornton and Pugh.³³ If first-site adaptation had fallen short of the proportionality implied by Weber's Law, then the Sloan notch would not have shifted as far as the homochromatic target wavelength.

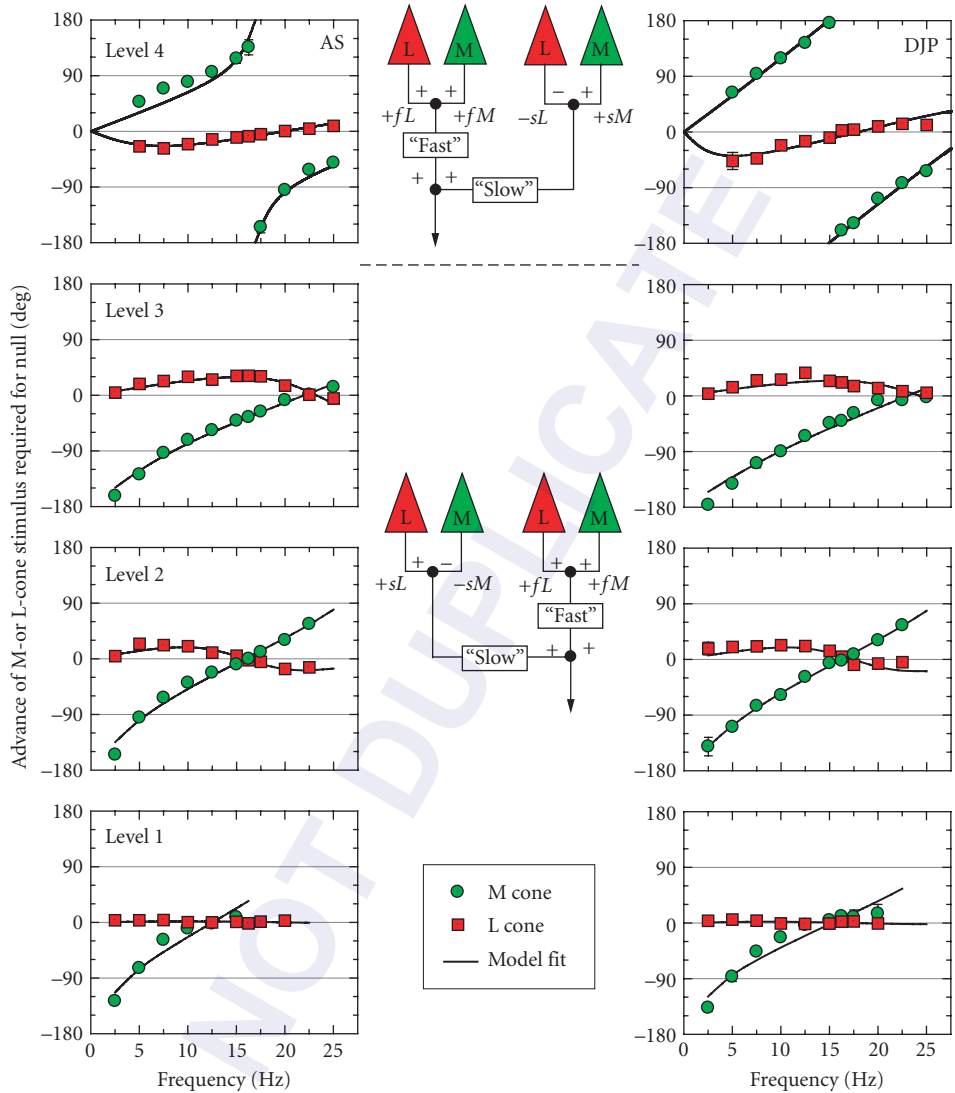


FIGURE 25 Large phase delays in the “luminance” channel. Phase advances of M-cone (green circles) or L-cone (red squares) stimuli required to flicker-photometrically null a 656-nm target for observers AS (left panels) and DJP (right panels) measured on 658-nm backgrounds of 8.93 (*Level 1*), 10.16 (*Level 2*), 11.18 (*Level 3*), or 12.50 (*Level 4*) \log_{10} quanta $s^{-1} \text{ deg}^{-2}$. The M-cone stimuli were alternating pairs of L-cone-equated 540 and 650 nm targets; and the L-cone stimuli were pairs of M-cone-equated 650- and 550-nm targets. The continuous lines are fits of a model in which the L- and M-cone signals are assumed to be the resultant of a fast signal (f) and a delayed slow signal (s) of the same or opposite sign. The predominant signals at each level are indicated by the wiring diagrams in the central insets (see Ref. 223).

As noted in Sec. 11.3, the idea that the zero or null point of an opponent second site always corresponds to the background field wavelength is due to von Kries adaptation is clearly an oversimplification. Were that the case, chromatic detection contours measured in cone contrast units would be independent of field wavelength, which they clearly are not (see subsection “Detection Contours and Field Adaptation” in Sec. 11.5). Similarly, chromatic sensitivity would be independent of field chromaticity at intensities at which Weber’s law had been reached, which is also not the case in field additivity experiments (see next). In general, if first-site adaptation was precisely inversely proportional to the amount of background excitation, second-site opponent desensitization would play little or no role in adaptation to steady fields. It would, of course, still play an important role in contrast adaptation; that is, in adaptation to excursions around the mean adapting level.

Field Additivity By using mainly spectral backgrounds (but sometimes with fixed-wavelength auxiliary backgrounds), Stiles restricted his measurements to the spectrum locus of color space. Some of the most revealing failures of his model, however, occur in the interior volume of color space when two spectral backgrounds are added together in different intensity ratios. Given the expectation in Stiles’ model that π -mechanisms should behave univariantly, bichromatic field mixtures should raise the threshold of a π -mechanism simply in accordance with the total quantum catch generated by those fields in that mechanism. For a linear mechanism, the effects of mixed fields, in other words, should be additive whenever only a single π -mechanism is involved. However, failures of field additivity can be subadditive, the mixed fields have less than the expected effect, or superadditive, when they have more than the expected effect. Subadditivity is an indication of opponency, and cases of subadditivity are thus another way that field experiments provide evidence of opponency.

Arguably the most interesting failures of field additivity occur under conditions that isolate Stiles’ S-cone mechanisms, π_1 and π_3 , as illustrated in Fig. 26. On a blue, 476-nm background, the threshold for a 425 nm target rises faster than Weber’s law predicts (leftmost data set, dashed vs continuous lines), this is superadditivity and is consistent with S-cone saturation.^{224,225} However, when a yellow, 590-nm background is added to the blue background, additivity typically fails. For lower 476-nm radiances (open circles, dashed vs continuous lines), superadditivity is found again,²²⁶ whereas for higher 476-nm radiances (open and filled squares, open triangles) subadditivity is found.^{227,228} Comparable failures of field additivity are found for the L-cone mechanism or π_3 ,^{55,188} but less clearly for the M-cone mechanisms or π_4 .^{188,196,229,230}

First- and Second-Site Adaptation Pugh and Mollon²³¹ developed a model to account for the failures of field additivity and other “anomalies” in the S-cone pathway, which has been influential. They assumed that the S-cone signal can be attenuated by gain controls at two sites. Gain at the first site is cone-specific, being controlled in the case of π_1 and π_3 solely by excitation in the S cones. By contrast, gain at the cone-opponent second-site is controlled by the net difference between signals generated by the S cones and those generated by the two other cone types. Excursions at the second site either in the S-cone direction or in the L+M-cone direction reduce the gain of that site, and thus attenuate the transmitted S-cone signal. Sensitivity at the second site is assumed to be greatest when the S-cone and the L+M-cone signals are balanced. Superadditivity occurs when attenuation at the first and second sites work together, for example on a blue spectral background. Subadditivity occurs when the addition of a background restores the balance at the second site, for example when a yellow background is added to a blue one. The original version of the Pugh and Mollon²³¹ model is formalized in their Eqs. (1) to (3). Another version of the model in which the second-site “gain” is implemented as a bipolar static nonlinearity (see Fig. 6 of Ref. 227) is illustrated in Fig. 27. The nonlinearity is a sigmoidal function that compresses the response and thus reduces sensitivity if the input is polarized either in the +S or $-(L+M)$ direction or in the L+M or $-S$ direction. The presence of a compressive nonlinearity in the output of mechanisms is now generally accepted as an important extension of the basic model. Such nonlinearities form a critical piece of models that aim to account for discrimination as well as detection data (see subsection “Pedestal Experiments” in Sec. 11.5).

There is an inconsistency in the Pugh and Mollon scheme for steady-state adaptation. In order for first-site adaptation to be consistent with Weber’s law ($\Delta I/I = \text{constant}$), sensitivity adjustment is assumed to be reciprocal—any increase in background radiance is compensated for by a reciprocal

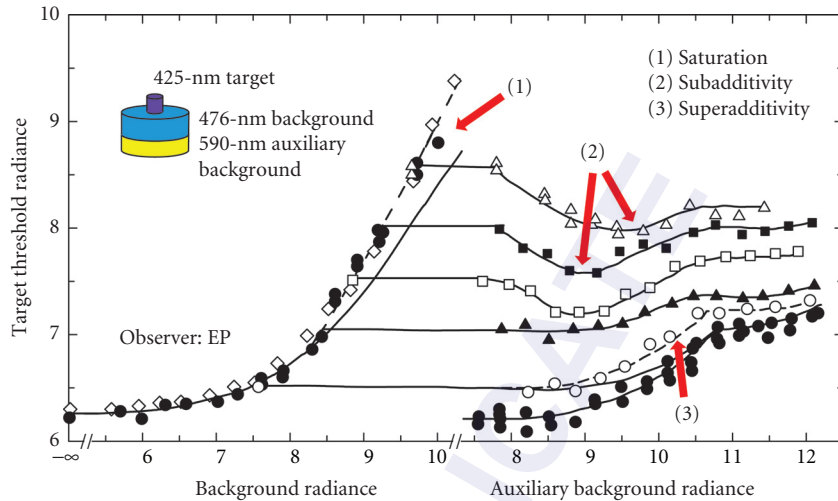


FIGURE 26 Field additivity failures. S-cone threshold data replotted from Fig. 1A of Pugh and Larimer.²²⁸ The data are increment thresholds for a 425 nm, 200-ms duration, 1° diameter foveal target presented in the centre of 476-nm fields (filled circles and open diamonds, left curves), 590-nm fields (filled circles, right curves), and various bichromatic mixtures of the two. (A 590-nm field of at least $8.0 \log_{10}$ quanta $s^{-1} \text{ deg}^{-2}$ was always added to the 476-nm field.) In the bichromatic mixture experiments a series of increasingly intense 590-nm fields, the intensities of which are indicated along the abscissa on the right, were added to a fixed intensity 476-nm field, the intensity of which corresponds to the abscissa associated with the dashed curve at target threshold radiance values corresponding to the data points at the lowest auxiliary background radiances. The shape of solid curves associated with the filled circles and open diamonds on the left and with the filled circles on the right is the standard Stiles threshold-versus-increment (tvi) template shape. The data provide examples of saturation (1), when the threshold elevation grows faster than Weber's law,^{224,225} subadditivity (2), when the threshold elevation of the bichromatic mixture falls short of the additivity prediction,²²⁷ and superadditivity (3), when the threshold elevation exceeds the additivity prediction.²²⁶ Figure design based on Figs. 1A and 2 of Pugh and Larimer.²²⁸ Units of radiance are \log quanta $\text{sec}^{-1} \text{ deg}^{-2}$.

decrease in gain. Consequently, the signal that leaves the first site (and potentially reaches the second site) should be independent of background radiance (see also “First-Site Adaptation” in Sec. 11.3). However, for the Pugh and Mollon model (and comparable models of first- and second-site adaptation) to work requires that the cone signals reaching the second site continue to grow with background radiance. To achieve this in their original model, they decrease the S-cone gain at the second site in proportion to the background radiance raised to the power n (where n is 0.71, 0.75, or 0.82 depending on the subject) rather than in proportion to the background radiance as at the first site ($n = 1$). Thus, the S-cone signal at the input to the second site continues to grow with background radiance, even though the signal at the output of the first site is constant. This scheme requires that the cone signal at second site somehow bypasses first-site adaptation, perhaps depending on signals at the border of the field⁸⁹ or on higher temporal frequency components, which are not subject to Weber's law.⁵²

A now mainly historical issue is how much adaptation occurs at the cone photoreceptors. Although some evidence has suggested that relatively little adaptation occurs within photoreceptors until close to bleaching levels,^{232,233} other compelling evidence suggests that significant adaptation occurs at much lower levels.^{46–48} Measures of adaptation in monkey horizontal cells show that light adaptation is well advanced at or before the first synapse in the visual pathway, and begins at levels as low as 15 td.^{234,235} Moreover, psychophysically, local adaptation has been demonstrated to occur with the resolution of single cones.^{236–238} The available evidence supports a receptor site of cone adaptation.

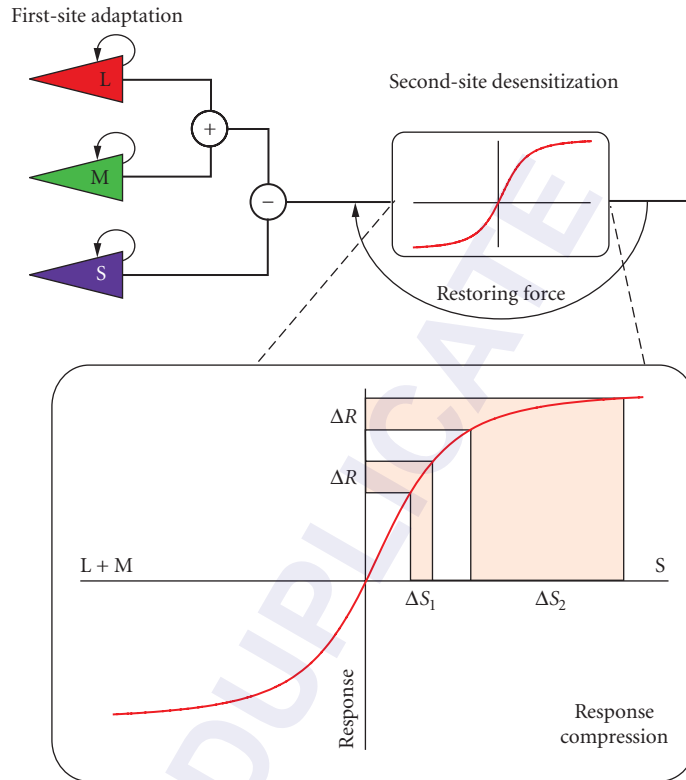


FIGURE 27 Pugh and Mollon cone-opponent model. Version of the Pugh and Mollon model²³¹ in which the cone-opponent, second-site desensitization is implemented at a bipolar static nonlinearity.²²⁷ The nonlinearity is a sigmoidal function that compresses the response as the input is polarized either in the +S [or $-(L+M)$] direction or in the L+M or ($-S$) direction. The effect of the nonlinearity can be understood by assuming that the system needs a criterion change in response (ΔR) in order for a change in the input to be detected. At smaller excursions along the +S direction, a smaller change in the S-cone signal, ΔS_1 , is required to produce the criterion response than the change in signal, ΔS_2 , required at greater excursions. Sensitivity should be greatest when the polarization is balanced, and the system is in the middle of its response range.

Detection Contours and Field Adaptation Stromeyer, Cole, and Kronauer⁶¹ measured cone contrasts in the L,M-cone plane (the plane defined by the $\Delta L/L$ and $\Delta M/M$ cone contrast axes, which includes chromatic L–M and achromatic L+M excursions) on spectral green, yellow, and red adapting backgrounds. Changing background wavelength moved the parallel chromatic contours for the detection of 200-ms flashes either outward (reducing sensitivity) or inward (increasing sensitivity). In general, background field wavelength strongly affects chromatic detection with sensitivity being maximal on yellow fields, declining slightly on green fields, and declining strongly on red fields. Since the slopes of the contours in contrast space are unaffected by these changes, the results again show that chromatic adaptation does not alter the relative weights of the M- and L-cone contrast inputs to the chromatic mechanism, but does change the sensitivity of the chromatically opponent second site. These data should be compared with the predictions shown in Fig. 8a. S-cone stimulation did not affect the L–M chromatic detection contours.⁶¹

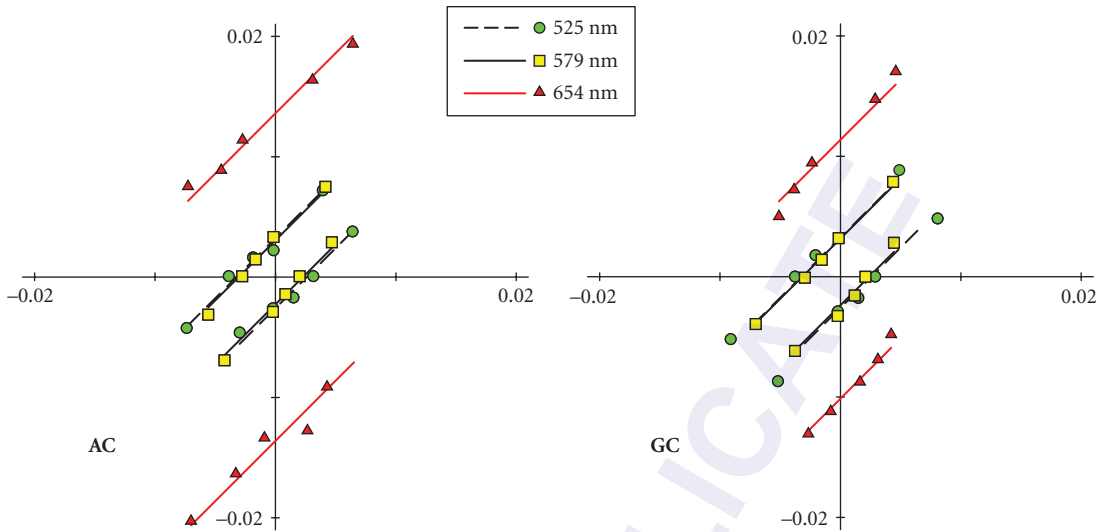


FIGURE 28 Second-site desensitization of L–M by steady fields. L–M detection contours in cone contrast space for two subjects measured by Chaparro, Stromeyer, Chen, and Kronauer⁶² on 3000-td fields of 525 nm (green circles), 579 nm (yellow squares), and 654 nm (red triangles) replotted from their Fig. 4. $\Delta M/M$ is plotted along the ordinate and $\Delta L/L$ along the abscissa. The straight contours fitted to thresholds have slopes of 1, consistent with detection by an L–M mechanism with equal and opposite L- and M-cone contrast weights. Constant relative L- and M-cone weights are a characteristic of von Kries first-site adaptation. The displacement of the contour from the origin on the 654 nm field is characteristic of second-site adaptation.

Figure 28 shows a clear example of the effect of changing the background wavelength of 3000 td fields from 525 to 579 to 654 nm from Fig. 4 of Ref. 62. Consistent with Stromeyer, Cole, and Kronauer,⁶¹ the L–M chromatic detection contours have positive unity slopes, and move outward from the origin as the wavelength lengthens to 654 nm.

The relation between these detection contours and spectral sensitivity data is nicely illustrated by replotting the spectral sensitivity data as cone contrasts. Figure 17*b*, shows the spectral sensitivity data of Thornton and Pugh¹⁰⁷ replotted in cone contrast space by Eskew, McLellan, and Giulianini.⁸⁹ The replotted spectral sensitivity data have the same characteristics as the detection contours shown in Fig. 28.

Habituation or Contrast Adaptation Experiments The extent of second-site desensitization produced by steady backgrounds is fundamentally limited by first-site adaptation. A natural extension of the field sensitivity method is to temporally modulate the chromaticity or luminance of the background around a mean level, so partially “bypassing” all but the most rapid effects of first-site adaptation and enhancing the effects of second-site desensitization. Moreover, if the spectral properties of the postreceptoral detection mechanisms are known, one mechanism can be selectively desensitized by adapting (or habituating) to modulations along a cardinal axis, which are invisible or “silent” to the other mechanisms (see Ref. 239 for a review of silent substitution methods and their origins).

In an influential paper, Krauskopf, Williams, and Heeley¹⁴ used habituation (the loss of sensitivity following prolonged exposure, in this case, to a temporally varying adapting stimulus) to investigate the properties of what they referred to as the “cardinal mechanisms” (see also Ref. 240). They used a 50 td, 2° diameter field made up of 632.5-, 514.5-, and 441.6-nm laser primaries. Following habituation to 1-Hz sinusoidal modulations (a 30-s initial habituation, then 5-s top-ups between test presentations) that were either chromatic (in the L–M, or S directions at equiluminance), or achromatic (L+M+S), or intermediate between these three directions, the detection sensitivity for Gaussian target pulses ($\sigma = 250$ ms) was measured in various directions of color space.

Figure 29 shows the results for observer DRW. The red circles show the losses of sensitivity and the black circles the increase in sensitivity following habituation to stimuli modulated in the directions of the orange arrows shown in all panels. These results should be compared with the predictions shown in Fig. 11. In their original paper, Krauskopf, Williams, and Heeley¹⁴ concluded that their results showed that habituation occurred independently along each of the three cardinal mechanism axes: L-M, S, and L+M. Thus, the sensitivity losses to habituation along one axis were mainly confined to targets modulated along the same axis, whereas the losses to habituation along axes intermediate to the cardinal ones involved losses in both mechanisms.

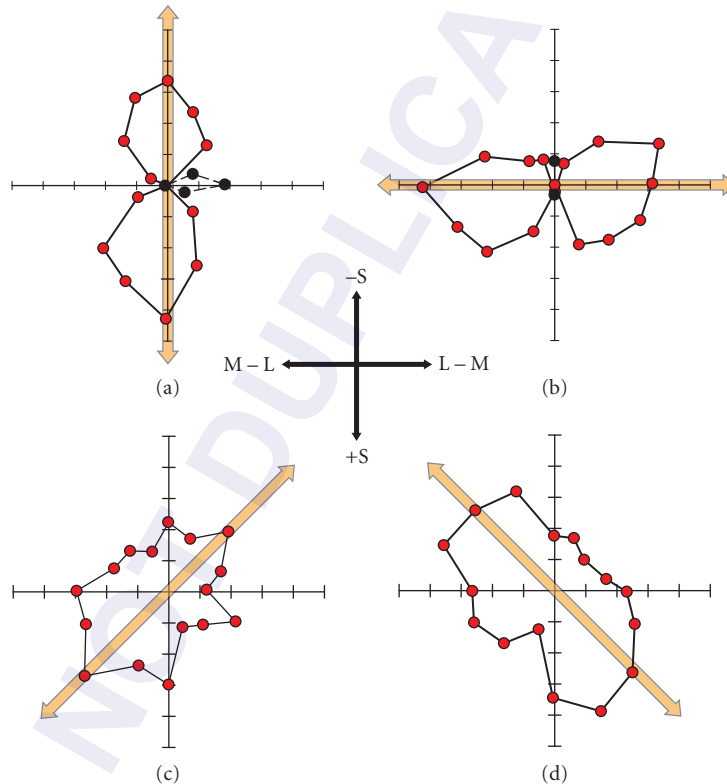


FIGURE 29 Habituation. Data from Fig. 6 of Krauskopf, Williams, and Heeley¹⁴ plotted in the canonical format shown in Fig. 11. These data are for stimuli that silence the L+M mechanism and reveal the properties of the L-M and S-(L+M) mechanisms plotted in a color space in which the horizontal axis corresponds to stimuli that only stimulate the L-M mechanism and the vertical axis corresponds to stimuli that only stimulate the S-(L+M) mechanism (by modulating S). (a) and (b) Results for habituation stimuli oriented along the horizontal and vertical axes, respectively. (c) and (d) Threshold changes for habituation stimuli along intermediate directions. The pattern of results is in qualitative agreement with the predictions developed in Figs. 9 to 11, and provides evidence for second-site adaptation in two opponent-discrimination mechanisms. The red circles plot increases in threshold following habituation, whereas the black circles plot cases where threshold was slightly decreased, rather than increased, by habituation.

Three years later, after a more rigorous analysis using Fourier methods, Krauskopf, Williams, Mandler, and Brown²⁴¹ concluded instead that the same data showed that habituation could desensitize multiple “higher-order” mechanisms tuned to intermediate directions in the equiluminant plane between the L–M and S cardinal directions. Yet, in terms of the magnitude of the selective losses, those along the intermediate axes are very much second-order effects compared with those found along the cardinal L–M or S axes. These second-order effects are not particularly compelling evidence for the existence of higher-order mechanisms, because they might also be due to more mundane causes, such as the visual distortion of the habituating stimuli or low-level mechanism interactions. The strong conclusion that can be drawn from the habituating experiments is that selective losses are largest when the habituating stimuli are modulated along the S, L–M, or L+M+S cardinal axes. The evidence for higher-order mechanisms along intermediate axes is secondary. Other workers have suggested that instead of there being multiple mechanisms, the cardinal mechanisms adapt to decorrelate the input signals so that their mechanism vectors rotate.^{242–244}

Noise-Masking Experiments Another field method used to investigate visual mechanisms is the introduction of masking noise to raise detection threshold. The spectral properties of the underlying mechanisms can then be investigated by varying either the chromaticity and/or luminance of the target (which, if the noise is held constant, is strictly speaking a test method) or by varying the chromaticity and/or luminance of the noise. As with habituating stimuli, the noise stimuli can be fashioned to excite one or other receptor or postreceptor mechanism. However, whereas habituation is assumed to have little effect at the first site, noise masking can potentially raise thresholds at both the first and the second sites. Although it is possible to selectively excite different first-stage cone mechanisms by modulating the noise along cone-isolating vectors or to selectively excite different second-stage color-discrimination mechanisms by modulating the noise along cardinal vectors, each noise stimulus inevitably excites both first- and second-stage mechanisms. As described in section “Sites of Limiting Noise” in Sec. 11.3 with respect to internal noise (see Figs. 12 and 13), changing the balance of noise at the first and second sites has the potential of rotating the detection contours in color space. Consequently, the discovery of detection contours for some noise vectors that do not align with the assumed low-level mechanism axes does not necessarily imply the existence of high-order mechanisms.

Gegenfurtner and Kiper²⁴⁵ measured thresholds for a 1.2-cpd Gabor patch ($\sigma = 0.8^\circ$ in space and 170 ms in time) in the presence of noise modulated in different directions around white in the L,M plane of DKL space. The target was either modulated in the L–M or L+M+S directions or in an intermediate direction along which the target color appeared to change between bright red and dark green. For all three directions masking was maximal when the noise was in the target direction, and minimal when it was in the orthogonal direction, even if the target direction did not lie along a cardinal axis. These findings suggested to the authors that there must be multiple mechanisms, some tuned to axes other than the cardinal ones.

The mechanisms derived by Gegenfurtner and Kiper, however, combined cone inputs nonlinearly. In a linear mechanism, the effectiveness of the noise should depend on the cosine between the cardinal direction and the noise direction, but Gegenfurtner and Kiper²⁴⁵ found mechanisms that were more narrowly tuned. Curiously, when square rather than Gabor patches were used, evidence for multiple mechanisms was not found.

The central findings of Gegenfurtner and Kiper have been contradicted in other papers, some of which used comparable techniques. Sankeralli and Mullen²⁴⁶ found evidence for only three mechanisms in $\Delta L/L$, $\Delta M/M$, $\Delta S/S$ cone contrast space using a 1-cpd semi-Gabor patch ($\sigma = 1.4^\circ$ vertically, and 4° horizontally with sharp edges in space, and $\sigma = 170$ ms in time). Mechanism tuning was found to be linear in cone contrast units with a cosine relationship between the noise effectiveness and the angle between the mechanism and noise vectors. Giulianini and Eskew¹⁷⁵ working in the $\Delta L/L$, $\Delta M/M$ cone contrast plane using 200-ms circular Gaussian blobs ($\sigma = 1^\circ$) or horizontal Gabor patches (1 cpd, $\sigma = 1^\circ$) similarly found evidence for only two mechanisms in that plane. Their results were consistent with a linear L–M mechanism with equal cone contrast weights, and a linear L+M mechanism with greater L-than M-cone weights.

Some of Giulianini and Eskew's¹⁷⁵ results from their Fig. 5 are shown in Fig. 20, for noise masks in the chromatic direction (Fig. 20a and c) and in the M-cone direction (Fig. 20b and d). Their interpretation

of the data is shown by the solid contours in Fig. 20*a* and *b*. Their model is consistent with detection by two mechanisms with roughly constant L–M contours across the two noise directions, but L+M contours that rotate slightly away from the M-cone axis in the presence of M-cone noise. However, as our elliptical fits in Fig. 20*c* and *d* show, other interpretations of their data are also plausible; most of the data in the presence of chromatic and M-cone noise can be described by ellipses, but the major axes of the ellipses do not align with the noise directions.

Subsequently, Eskew, Newton, and Giulianini¹⁷⁹ extended these measurements to include S-cone modulations in the equiluminant plane, and again found no evidence for higher-order mechanisms in detection or near-threshold discrimination tasks. Some of their results are shown in Fig. 21. These data seem more consistent with a rounded parallelogram than an ellipse. All three studies just described, in contradiction to Gegenfurtner and Kiper,²⁴⁵ found that the directions at which noise masking was maximal more closely aligned with the expected mechanism directions than with the noise directions. However, see also Hansen and Gegenfurtner.²⁴⁷

In a related study, D’Zmura and Knoblauch²⁴⁸ did find evidence for multiple mechanisms within the equiluminant plane of DKL space. They used “sectored” noise made up of noise centered on the signal direction combined with noise of various strengths in the orthogonal direction. Contrast thresholds for “yellow,” “orange,” “red,” and “violet” signals ($\sigma = 2.4^\circ$ in space and 160 ms in time) were independent of the strength of the orthogonal noise, which suggests the existence of linear broad-band mechanisms tuned to each of the signal directions (for which the orthogonal direction is a null direction). Like Gegenfurtner and Kiper²⁴⁵ multiple mechanisms were found, but they were broadly rather than narrowly tuned. D’Zmura and Knoblauch²⁴⁸ argued that narrow tuning can result from off-axis looking in which subjects switch among multiple linear broad-band detection channels to reduce the noise and maximize the detectability of the signal.

An important caveat about D’Zmura and Knoblauch’s experiments is that they used the 1924 $V(\lambda)$ function to equate luminance (see p. 3118 of Ref. 248). Because $V(\lambda)$ substantially underestimates luminous efficiency at short wavelengths, modulations of their blue phosphor will produce a luminance signal in their nominally equiluminant plane, which could complicate the interpretation of their results. For example, nominally S-cone modulations might have superimposed and cancelling L+M modulations, and some signal directions might be luminance-detected (and therefore artifactually independent of chromatic noise in the orthogonal direction). Nonetheless, Monaci et al.²⁴⁹ using a related technique, but a better definition of luminous efficiency, came to the same conclusion.

Most recently, Giulianini and Eskew²⁵⁰ used noise masking to investigate the properties of the S–(L+M) mechanism, and, in particular, whether or not the mechanism responds linearly to noises in different directions of cone contrast space. Though the effects of noise on the L–M mechanism are linear, for both polarities of the S–(L+M) mechanism (i.e., for +S and –S detection) its effects are nonlinear.²⁵⁰

Chromatic Discrimination

Basic Experiments Chromatic discrimination is the ability to distinguish the difference between two test lights that differ in chromaticity. The prototypical set of chromatic discrimination data are the discrimination ellipses of MacAdam²⁵¹ originally plotted as CIE x, y chromaticity coordinates, and thus in a color space that is not especially conducive to understanding the underlying mechanisms. These ellipses were determined from the variability of repeated color matches made at 25 different reference chromaticities in the CIE color space. Le Grand²⁵² reanalyzed MacAdam’s data and found that much of the variability between the ellipses could be accounted for by considering variability along the tritan (S) dimension and along the cone-opponent (L–M) dimension. Interpretation of MacAdam’s data is complicated by the fact that each ellipse was measured for a different state of adaptation, namely that determined by the mean match whose variability was being assessed.

Boynton and Kambe²⁵³ made a series of discrimination measurements in the equiluminant plane along the important S and L–M color directions. They found that discrimination along the S axis depended on only S-cone excitation level; specifically that $\Delta S/S + S_0 = \text{constant}$, where S is the background excitation and S_0 is an “eigengrau” or “darklight” term. At high S-cone excitation levels, as S_0 becomes less significant, Weber’s law holds (i.e., $\Delta S/S$ is approximately constant).^{254–256} A similar

conclusion was reached by Rodieck²⁵⁷ in his reanalysis of MacAdam's²⁵¹ data (see Figs. XXIII-10 and XXIII-11 of Ref. 257). When plotted in a normalized cone excitation diagram, chromaticity ellipses for individual subjects have similar shapes regardless of the reference chromaticity.²⁵⁸

Chromatic Discrimination Near Detection Threshold Under conditions that isolate the chromatic mechanisms, near-threshold lights can be indistinguishable over fairly extended spectral or chromaticity ranges. On a white background using 0.75° diameter, 500-ms duration test discs, Mullen and Kulikowski¹⁹⁷ found that there were four indistinguishable wavelength ranges that corresponded to the threshold lights appearing orange (above 585 nm), pale yellow (566–585 nm), green (c. 490–566 nm), and blue (c. 456–490 nm), and that there was weaker evidence for a fifth range that appeared violet (< c. 456 nm). Importantly, too, these boundaries were independent of the wavelength *differences* between pairs of lights, which suggests that each mechanism is univariant. These results suggest the existence of four or five distinct unipolar mechanisms at threshold.¹⁹⁷ Unlike the bipolar mechanisms (L–M) and S–(L+M), a unipolar mechanism responds with only positive sign. The bipolar (L–M) mechanism, for example, may be thought of as consisting of two unipolar mechanisms |L–M| and |M–L|, each of which is half-wave rectified (see “Unipolar versus Bipolar Chromatic Mechanisms” in Sec. 11.6).

Eskew, Newton, and Giulianini¹⁷⁹ looked at discrimination between near-threshold Gabor patches modulated in various directions in the equiluminant plane. Noise masking was used to desensitize the more sensitive L–M mechanism and thus reveal more of the S-cone mechanism, and perhaps higher-order mechanisms. Despite the presence of noise they found evidence for only four color regions within which lights were indiscriminable. These regions corresponded to detection by the four poles of the classical cone-opponent mechanisms |L–M|, |M–L|, |S–(L+M)|, and |(L+M)–S|. Some of their results are shown in Fig. 21. The four colors shown in the three outer arcs of the figure correspond to the four regions within which lights were indiscriminable. The data in each semicircular polar plot (filled squares) show the discriminability of a given chromatic vector from the standard vector shown by the arrow. The open circles are model predictions (see legend). There is no evidence here for higher-order color mechanisms. In a later abstract, Eskew, Wang, and Richters²⁵⁹ reported evidence also from threshold level discrimination measurements for a fifth unipolar mechanism, which from hue-scaling measurements they identified as generating a “purple” hue percept (with the other four mechanisms generating blue, green, yellow, and red percepts).

Thus, the work of both Mullen and Kulikowski¹⁹⁷ and Eskew et al.^{179,259} may be consistent with five threshold-level color mechanisms. By contrast, Krauskopf et al.²⁴¹ found evidence for multiple mechanisms. They measured the detection and discrimination of lights that were 90° apart in the equiluminant DKL plane and modulated either along cardinal or intermediate axes. No difference was found between data from the cardinal and noncardinal axes, which should not be the case if only cardinal mechanisms mediate discrimination. The authors argued that stimuli 90° apart modulated along the two cardinal axes should always be discriminable at threshold, because they are always detected by different cardinal mechanisms. In contrast, stimuli 90° apart modulated along intermediate axes will only be discriminable if they are detected by different cardinal mechanisms, which will not always be the case. Care is required with arguments of this sort, because angles between stimuli are not preserved across color space transformations (i.e., they depend on the choice of axes for the space and also on the scaling of the axes relative to each). More secure are conclusions based on an explicit model of the detection and discrimination process. Note that this caveat applies not just here but also to many studies that rely on intuitions about the angular spacing of stimuli relative to each other.

Clearly, colors that are substantially suprathreshold and stimulate more than one color mechanism should be distinguishable. However, the range of suprathreshold univariance can be surprisingly large. Calkins, Thornton, and Pugh¹¹¹ made chromatic discrimination and detection measurements in the red-green range from 530 to 670 nm both at threshold and at levels up to 8 times threshold. Measurements were made with single trough-to-trough periods of a 2-Hz raised-cosine presented on a yellow field of 6700 td and 578 nm. The observers' responses were functionally monochromatic either for wavelengths below the Sloan notch (between 530 and 560 nm) or for wavelengths above the Sloan notch (between 600 and 670 nm); that is, any two wavelengths within each range were indistinguishable at some relative intensity. The action spectra for indiscriminability above threshold

and for detection at threshold were consistent with an L–M opponent mechanism. Discrimination for pairs of wavelengths either side of the Sloan notch was very good, and between stimuli near the notch and stimuli further away from the notch was also good. The good discriminability near the notch is consistent with a diminution of $|L-M|$ as the signals become balanced and with the intrusion of another mechanism, which between 575 and 610 nm has a spectral sensitivity comparable with the photopic luminosity function [although whether or not this corresponds to luminance or the L+M lobe of S–(L+M) is unclear]. For full details, see Ref. 111.

The results reviewed in this section are especially important because they suggest that each color mechanism is not only univariant, but behaves like a “labeled line.” As Calkins et al. note on their p. 2365 “Phenomenologically speaking, the discrimination judgments [in the monochromatic regions] are based upon variation in perceptual brightness.”

Pedestal Experiments Discrimination experiments, sometimes called “pedestal experiments” can be used to investigate the properties of mechanisms over a range of contrasts both above and below the normal detection threshold. In a typical pedestal experiment, an observer is presented with two identical but brief “pedestal” stimuli, separated either in space or in time, upon one of which a test stimulus is superimposed. The observer’s task is to discriminate in which spatial or temporal interval the test stimulus occurred. By presenting the test and pedestal stimuli to the same mechanism (sometimes called the “uncrossed” condition), it is possible to investigate the properties of chromatic and luminance mechanisms separately. By presenting the test and pedestal stimuli to different mechanisms (sometimes called the “crossed” condition), it is possible to investigate interactions between them. The presentation of the pedestals causes a rapid shift away from the mean adapting state. The transient signals produced by these shifts will be affected by instantaneous or rapid adaptation mechanisms, such as changes in the response that depend on static transducer functions, but not by more sluggish adaptation mechanisms.

When pedestal and test stimuli are both uncrossed and of the *same* sign, the discrimination of the test stimulus typically follows a “dipper” function as the pedestal contrast is increased. As the pedestal contrast is raised from below its contrast threshold, the threshold for detecting the test stimulus falls by a factor of about two or three, reaching a minimum just above the pedestal threshold, after which further increases in pedestal contrast hurts the detection of the test.^{260–265} Facilitation has been explained by supposing there is a transducer (input-output) function that accelerates at lower contrasts, thus causing the difference in output generated by the pedestal and test-plus-pedestal to increase with pedestal contrast (e.g., Refs. 263, 266). Alternatively, it has been explained by supposing that the pedestal, because it is a copy of the test, reduces detection uncertainty.²⁶⁷ The dipper function is found both for uncrossed luminance and uncrossed chromatic stimuli, and has similar characteristics for the two mechanisms.^{169,268,269} When pedestals and test stimuli are uncrossed but of the *opposite* sign, increases in pedestal contrast first raise threshold—consistent with subthreshold cancellation between the pedestal and test. As the pedestal approaches its own threshold, however, discrimination threshold abruptly decreases, after which further increases in pedestal contrast hurt detection of the test for both uncrossed luminance and uncrossed chromatic conditions.¹⁶⁹ The increase followed by decrease in threshold has been referred to as a “bumper” function.²⁷⁰ The dipper and bumper functions found under uncrossed conditions show clear evidence of subthreshold summation or cancellation. Figure 30a shows examples of the dipper (upper right and lower left quadrants) and bumper functions (upper left and lower right quadrants) for uncrossed L–M chromatic pedestals and targets from Cole, Stromeyer, and Kronauer.¹⁶⁹

Presenting the pedestals and test stimuli to different mechanisms in the crossed conditions is a potentially powerful way of investigating how the luminance and chromatic mechanisms interact. For example, if the two mechanisms interact before the hypothesized accelerating transducer, then the crossed and uncrossed dipper and bumper functions should be comparable and should also show evidence for subthreshold interactions. If the chromatic and luminance mechanisms are largely independent until later stages of visual processing, then subthreshold effects should be absent, but given the uncertainty model suprathreshold facilitation might still be expected.

The results obtained with crossed test and pedestal stimuli, however, remain somewhat equivocal. Luminance square-wave gratings or spots have been shown to facilitate chromatic discrimination.^{194,271}

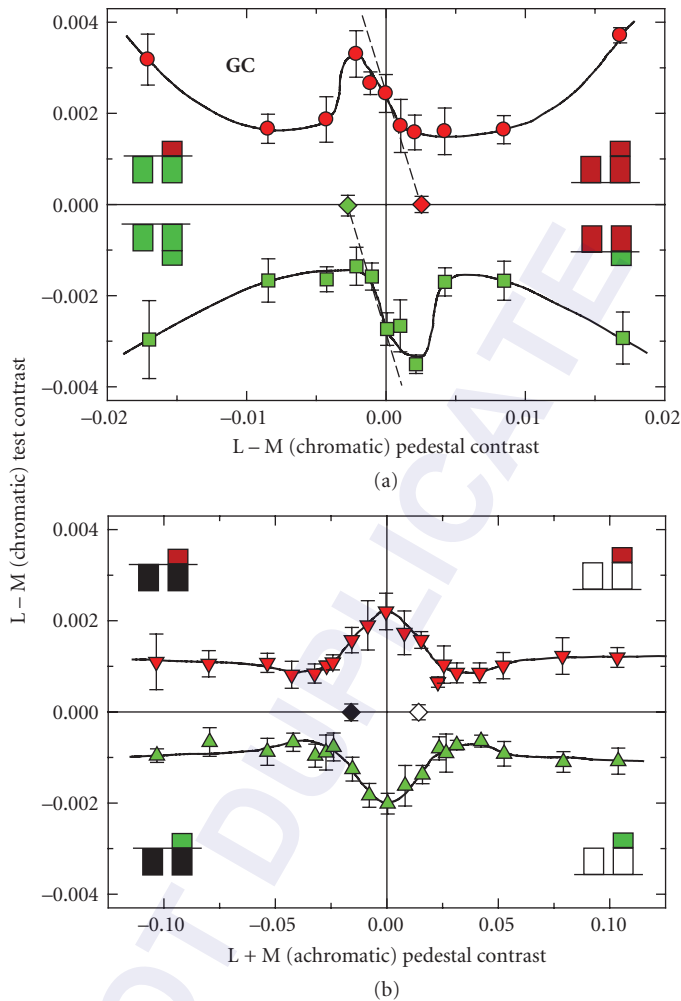


FIGURE 30 Pedestal effects. (a) Detection of chromatic L-M tests on chromatic L-M pedestals. As indicated by the icons in each quadrant, the data points are the thresholds for detecting red chromatic targets (red circles) on either red (upper right quadrant) or green (upper left quadrant) chromatic pedestals, or for detecting green chromatic targets (green squares) on either red (lower right quadrant) or green (lower left quadrant) pedestals. The diamonds show the thresholds for detecting the red (red diamond) or green (green diamond) chromatic pedestals alone. The dashed lines show predictions that the discrimination of the presence of the test depends only on the intensity of the test plus the pedestal (which implies that pedestal alone is ineffectual). (b) Detection of chromatic L-M tests on achromatic L+M pedestals. As indicated by the icons, the data are the thresholds for detecting red chromatic targets (red inverted triangles) on incremental (upper right quadrant) or decremental (upper left quadrant) luminance pedestals, or for detecting green chromatic targets (green triangles) on incremental (lower right quadrant) or decremental (lower left quadrant) luminance pedestals. The diamonds show the thresholds for detecting the incremental (open diamond) or decremental (filled diamond) luminance pedestals alone. (Data replotted from Figs. 4 and 5 of Cole, Stromeyer, and Kronauer.¹⁶⁹ Design of the figure based on Fig. 18.10 of Eskew, McLellan, and Giulianini.⁸⁹)

Using sine-wave gratings, De Valois and Switkes²⁶⁵ and Switkes, Bradley, and De Valois²⁶⁸ found that crossed facilitation and masking was asymmetric: for luminance pedestals and chromatic tests they found facilitation and no masking, but for chromatic pedestals and luminance tests they found no facilitation and strong masking. This curious asymmetry has not always been replicated in subsequent measurements. Using spots, Cole, Stromeyer, and Kronauer¹⁶⁹ obtained equivalent results for crossed luminance pedestals and chromatic tests and for crossed chromatic pedestals and luminance tests. They found in both cases that increasing the contrast of the crossed pedestal had little effect on test threshold until the pedestal approached or exceeded its own threshold by a factor of two or more, and that further increases in pedestal contrast had little or no masking effect.¹⁶⁹ Figure 30*b* shows examples of their crossed luminance and chromatic results. Similarly, Mullen, and Losada²⁶⁹ found facilitation for both types of crossed pedestals and tests when the pedestal was suprathreshold, but unlike Cole, Stromeyer, and Kronauer¹⁶⁹ they found crossed masking for higher pedestal contrasts. Chen, Foley, and Brainard,²⁷² however, obtained results that were more consistent with De Valois and Switkes²⁶⁵ and Switkes, Bradley, and De Valois.²⁶⁸ They used spatial Gabors (1 cpd, 1° standard deviation) with a Gaussian temporal presentation (40-ms standard deviation, 160-ms total duration), and investigated the effects of targets and pedestals presented in various directions of color space. They found masking for all target-pedestal pairs and facilitation for most target-pedestal pairs. The exceptions were that (1) equiluminant pedestals did not facilitate luminance targets, and (2) tritan (blue/yellow) pedestals facilitated only tritan targets.

So, what does all this mean? The fact that the crossed facilitation, when it is found, occurs mostly with suprathreshold pedestals^{169,268,269} suggests that the interaction occurs after the stages that limit the thresholds in the chromatic and luminance mechanisms. This is surprising, given that center-surround L–M neurons respond to both chromatic and luminance gratings (see subsection “Multiplexing Chromatic and Achromatic Signals” in Sec. 11.5). This early independence is supported by findings with gratings that the crossed facilitation is independent of the spatial phase of the gratings.^{268,269} However, the facilitation is not found with crossed dichoptic pedestals and tests. The finding that crossed facilitation survives if the luminance pedestal is replaced with a thin ring¹⁶⁹ suggests that it might be due to uncertainty reduction. However, yes-no psychometric functions and receiver operating characteristics show clearly that the facilitation of chromatic detection by luminance contours is inconsistent with the uncertainty model.¹⁷⁷ To explain their results, Chen, Foley, and Brainard²⁷³ presented a model made up of three postreceptoral channels with separate excitatory inputs (producing facilitation), but each with divisive inhibition aggregated from all channels (producing masking).

Krauskopf and Gegenfurtner²⁵⁵ performed a series of experiments in which chromatic discrimination was measured after an abrupt change in chromaticity away from an equal-energy-white adapting light (see also Refs. 274, 275). The change in chromaticity was effected by simultaneously presenting four 36-min diameter discs for 1 s. Three of the discs were of the same chromaticity, but the fourth was slightly offset from the others. The observer’s task was to identify which disc was different. The task can be thought of as the presentation of three stimuli defined by a “pedestal” color vector from the white point with the fourth stimulus being defined by the vector sum of the pedestal vector plus a “test” vector.⁸⁹ The experiments were carried out in the DKL equiluminant plane at a luminance of 37 cd/m². The discrimination of a test vector along one of the cardinal axes was found to deteriorate if the pedestal vector was along the same axis, but was relatively unaffected if the pedestal vector was along the other axis. Color discrimination, therefore, depended primarily on the L–M and S–(L+M) mechanisms nominally isolated by modulations along the cardinal axes.²⁵⁵ A more detailed determination of the shapes of chromatic discrimination ellipses for pedestal and test vectors along 16 color directions was also consistent with discrimination being determined by the cardinal mechanisms, but there was also evidence for another mechanism orientated along a blue-green to orange axis, from 135 to 315° in the DKL space.²⁵⁵ This might reflect the existence of a higher-level mechanism, but it might also be consistent with other explanations. Maintaining the phosphors at 37 cd/m² [a photometric measure defined by the discredited CIE $V(\lambda)$ function that underestimates luminous efficiency in the blue and violet] results in modulation of the blue phosphor at nominal “equiluminance” producing a luminance (L+M) axis close to the blue-green axis. The effect of such a signal is hard to predict, but could affect the sensitivity of the S–(L+M) mechanism by

cancelling S, or it could be detected by L+M. Another possibility is that the evidence for an additional mechanism might instead reflect a small S-cone contribution to L–M, which would rotate the L–M mechanism axis.⁸⁹

Color Appearance and Color Opponency

When viewed in isolation (i.e., in the “aperture” mode), monochromatic lights have characteristic color appearances. Across the visible spectrum hues vary from violet through blue, blue-green or cyan, green, yellow-green, yellow, orange, and red.²⁷⁶ Yet, however compelling these colors might seem, they are private perceptions. Consequently, to investigate them, we must move away from the safer and more objective psychophysical realms of color matching, color detection, and color discrimination, and rely instead on observers’ introspections about their perceptual experiences (see Ref. 277 for a related discussion of “Class A” and “Class B” observations). Experiments that depend upon introspection are an anathema to some psychophysicists, yet they represent one of the few methods of gaining insights into how color is processed beyond the early stages of the visual system. As long as the limitations of the techniques are recognized, such experiments can be revealing. In addition, it is of interest to understand color appearance per se, and to do so requires the use of subjective methods in which observers report their experience.

Physiological measurements of the spectral properties of neurons in the retina and LGN are generally inconsistent with the spectral properties of color-appearance mechanisms, which suggests that color appearance is likely to be determined by cortical processes (see, e.g., Ref. 18). Relative to the input at the photoreceptors, the relevant neural signals that eventually determine color appearance will have undergone several transformations, likely to be both linear and nonlinear, before they reach the cortex. Despite this many models of color appearance are based on the assumption that appearance can be accounted for by mechanisms that linearly combine cone inputs (see Sec. 11.4). As we will see later, this is at best an approximation, at least for some aspects of color appearance. One reason for the failures of linearity may be that the suprathreshold test stimuli used to probe appearance are often adapting stimuli that intrude significantly on the measurements. That is, unlike in detection experiments where near-threshold tests can reasonably be assumed not to perturb the adapted state of the visual system very much, nonlinear adaptive processes cannot be ignored even for nominally steady-state color-appearance measurements, except perhaps under some contrived experimental conditions (see subsection “Linearity of Color-Opponent Mechanisms” in Sec. 11.5).

Opponent-Colors Theory Color appearance is often conceptualized in terms of Hering’s opponent-color theory; that is, in terms of the perceptual opposition of red and green (R/G), blue and yellow (B/Y), and dark and light. In many versions of the model, the opponent mechanisms are assumed, either explicitly or implicitly, to combine their inputs from cones in a linear fashion.^{11,12,18,19,29,76,278–282} Stage 3 of Fig. 4 represents a version of the linear model derived by the simple combinations of cone-opponent mechanisms (see subsection “Three-Stage Zone Models” in Sec. 11.6 for details).

The assumption of linearity constrains and simplifies the predicted properties of the color-opponent mechanisms. For example, under this assumption each opponent mechanism is completely characterized by its spectral response. As a function of wavelength, the R/G opponent mechanism responds R at both short and long wavelengths, corresponding to the red constituent in the appearance of both short- and long-wavelength lights, and responds G at middle wavelengths. The B/Y mechanism responds B at shorter wavelengths and Y at longer wavelengths. Each response (R, G, Y, or B) is univariant, and the responses of opposed poles (R vs. G, and B vs. Y) are mutually exclusive. The colors of lights detected solely by the same pole of a given mechanism cannot be distinguished at threshold, but lights detected by different poles can. In the model, color appearance depends on the relative outputs of the color-opponent appearance mechanisms. The individual mechanisms signal one opponent color or the other (e.g., redness or greenness) in different strengths, or a null.

The wavelength at which the opponent response of a mechanism changes polarity is a zero crossing for that mechanism, at which its response falls to zero. These zero crossings correspond to the unique hues. Unique blue and yellow are the two zero crossings for the R/G mechanism, and unique green and red are the two zero crossings for the Y/B mechanism. Unique red does not appear in plots

of valence against wavelength because it is extra-spectral (i.e., it has no monochromatic metamer and can only be produced by a mixture of spectral lights). Another extra-spectral color is white, which corresponds to an “equilibrium color” for both the R/G and B/Y mechanisms. A property of opponent-color mechanisms is that both the R/G and the Y/B mechanisms do not respond to a suitably chosen broadband white stimulus (such as an equal energy white). Thus, the sum of the outputs of the positive and negative lobes of each mechanism in response to such stimuli should be zero.

Tests of the assumptions of linearity of opponent-colors mechanisms are discussed below (see subsection “Linearity of Color-Opponent Mechanisms” in Sec. 11.5). See also Hurvich²⁸³ for more details of opponent colors theory.

Spectral Properties of Color-Opponent Mechanisms Several techniques have been used to investigate the spectral properties of the color-opponent mechanisms. Most rely on the assumption that opponent-color mechanisms can be nulled or silenced by lights or combinations of lights that balance opposing sensations. The null may be achieved experimentally by equating two lights so that their combination appears to be in equilibrium (see section “Appearance Test Measurements and Opponency” in Sec. 11.4). A closely related alternative is to find single spectral lights that appear to be in equilibrium for a particular mechanism (e.g., that appear neither red nor green, or neither blue nor yellow). Another related technique is hue scaling, in which lights are scaled according to how blue, green, yellow, or red they appear. Thus, lights that are rated 100 percent blue, green, yellow, or red are equilibrium or unique hues.

Hue scaling Despite the variety of color terms that could be used to describe the colors of lights, observers require only four—red, yellow, green, and blue—to describe most aspects of color appearance.^{6,7} Other colors such as orange can be described as reddish-yellow, cyan as bluish-green, purple as reddish-blue, and so on. The need for just four color terms is consistent with opponent-colors theory. By asking observers to scale how blue, green, yellow, and blue spectral lights appear, hue scaling can be used to estimate the spectral response curves of the opponent mechanisms.^{21,284–286} A cyan, for example, might be described as 50-percent green, 50-percent blue, and an orange as 60-percent red and 40-percent yellow. Figure 31 shows hue scaling data obtained by De Valois et al.,²¹ who instead of using spectral lights used equiluminant modulations in DKL space, the implications of which are discussed below. Hue-scaling results have been reported to be consistent with hue-cancellation valence functions,²⁸⁶ but inconsistencies have also been reported.²⁸⁷

Color-opponent response or valence functions As introduced in section “Appearance Test Measurements and Opponency,” the phenomenological color-opponency theory of Hering⁶ was given a more quantitative basis by the hue-cancellation technique.^{279,280,288} Jameson and Hurvich (see Fig. 14 for some of their data) determined the amount of particular standard lights that had to be added to a spectral test light in order to cancel the perception of either redness, greenness, blueness, or yellowness that was apparent in the test light, and so produced “valence” functions. Figure 32 shows additional R/G valence data from Ingling and Tsou¹² and B/Y valence data from Werner and Wooten.²⁸⁶ Individual lobes of the curves show the amounts of red, green, blue, or yellow light required to cancel the perception of its opposing color in the test light. The zero crossings are the “unique” colors of blue, green, and yellow.

These types of judgments are, of course, highly subjective because the observer must abstract a particular color quality from color sensations that vary in more than one perceptual dimension. In valence settings, the color of the lights at equilibrium depends on the wavelength of the spectral test light. Lights that are in red/green equilibrium vary along a blue-achromatic-yellow color dimension, while those in yellow/blue equilibrium vary along a red-achromatic-green color dimension. To reduce the need for the observer to decide when a light is close enough to the neutral point to accept the adjustment, some authors have used techniques in which the observer judges simply whether a presented light appears, for example, reddish or greenish, to establish the neutral points of each mechanism.²⁸⁹

Unique hues and equilibrium colors The equilibrium points of the R/G color opponent mechanism correspond to the unique blue and unique yellow hues that appear neither red nor green,

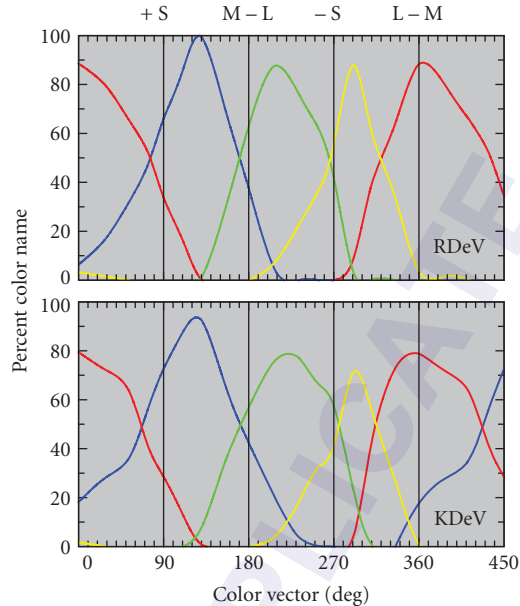


FIGURE 31 Hue scaling. Data replotted from Fig. 3 (Observers: RDeV and KDeV) of De Valois, De Valois, Switkes, and Mahon.²¹ The data show the percentage of times that a stimulus with a given color vector in the equiluminant plane of DKL space was called red (red lines), yellow (yellow lines), green (green lines), or blue (blue lines). The vertical lines indicate the cardinal axes at $90/450^\circ$ (+S), 180° (M-L), 270° (-S), and $0/360^\circ$ (L-M). With the exception perhaps of the red function for RDeV, the other hue naming functions do not align with the cardinal axes.

whereas those of the B/Y mechanism correspond to the unique green and unique red hues that appear neither blue nor yellow. These hues are “unique” in the sense that they appear phenomenologically unmixed.²⁹⁰

Measurements of the spectral positions of the unique hues have been made in several studies.^{281,290–296} Kuehni²⁹⁷ provides a useful summary of the unique hue settings from 10 relatively recent studies. For unique yellow, the mean and standard deviation after weighting by the number of observers in each study are 577.8 and 2.9 nm, respectively; for unique green, 527.2 and 14.9 nm; and for unique blue, 476.8 and 5.2 nm. Studies that generate colors either on monitors^{22,289,298} or in print,^{299,300} and therefore use desaturated, nonspectral colors, are not necessarily comparable with the results obtained using spectral lights, since perceived hue depends upon saturation.³⁰¹ The mean unique yellow, green, and blue wavelengths taken from Table of Kuehni²⁹⁷ *excluding* those obtained by extrapolation from desaturated, nonspectral equilibrium lights are 576.3, 514.8, and 478.0 nm, respectively. Although unique hues show sizable individual differences, the settings within observers are generally quite precise.^{302,303} Dimmick and Hubbard²⁹⁰ review historical estimates.

Dimmick and Hubbard²⁹⁰ reported, nearly 70 years ago, that unique blue and yellow and unique red and green are not complementaries (see also Ref. 304). By extending unique hue settings to include nonspectral colors it is possible to determine the equilibrium vectors of the opponent mechanisms in two dimensions or the equilibrium planes in three dimensions. Valberg³⁰⁵ determined unique blues, greens, yellows, and reds as a function of saturation. Although not commented on in the original paper, two important features of his results were that (1) the unique red and green vectors plotted

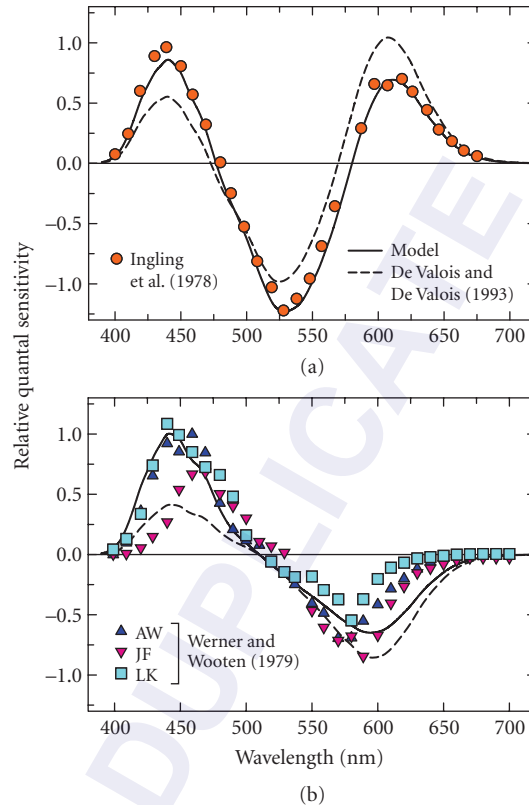


FIGURE 32 Color valence data. (a) R/G valence data (red circles) replotted from Fig. 2A of Ingling, Russell, Rea, and Tsou.³¹⁴ (b) B/Y valence data for observers AW (dark blue triangles), JF (purple inverted triangles) and LK (light blue squares) replotted from Fig. 9 of Werner and Wooten.²⁸⁶ The continuous lines are the spectral sensitivity predictions of Stage 3 of the Müller-zone model outlined in the text (see section “Three-Stage Zone Models” in Sec. 11.6). Each set of valence data has been scaled to best fit the spectral sensitivity predictions. For comparison, the spectral sensitivity predictions of the model proposed by De Valois and De Valois¹⁸ are also shown as dashed lines, in each case scaled to best fit the Müller-zone model predictions.

in a CIE x,y chromaticity diagram were not colinear, and (2) unique yellow and blue vectors, though continuous through the white point, were slightly curved. Burns et al.¹⁶ in a comparable experiment confirmed that the unique green and red vectors were not colinear, and that the unique blue vector was curved, but found that the unique yellow and green vectors were roughly straight.

Chichilnisky and Wandell²⁸⁹ located equilibrium boundary planes by presenting stimuli varying in both intensity and chromaticity on various backgrounds, and then asking subjects to classify them into red-green, blue-yellow, and white-black opponent-color categories. They found that the opponent-color classification boundaries were not coplanar. Wuerger, Atkinson, and Cropper²² determined the null planes for the four unique hues using a hue-selection task in which subjects selected the equilibrium hue from an array of colors. They concluded that unique green and unique red planes in a space defined by the color monitor phosphors were not coplanar, consistent with the previous findings, but

that unique blue and unique yellow could form a single plane. However, they did not try to fit their data with curved surfaces, as might be suggested by previous work.

In summary, unique red and green are not colinear in two-dimensional color space or coplanar in three, and the unique blue “vector” or “plane” is curved. As we discuss below, failures of colinearity or coplanarity imply that either a bipolar color-opponent mechanism is not a single mechanism or it is a single but nonlinear mechanism, while curved vectors and planes imply failures of additivity within single mechanisms.

Linearity of Color-Opponent Mechanisms

Tests of linearity For the spectral sensitivities of color-opponent mechanisms to be generalizable implicitly or explicitly, the color-opponent mechanisms must be linear and additive. This assumption is made in many accounts of zero crossings and/or in descriptions of valence functions as linear combinations of the cone fundamentals.^{11,18,19,24,29,76,278,279,306–308}

Larimer, Krantz, and Cicerone^{307,309} specifically tested whether the equilibrium colors behaved additively under dark adapted conditions. They tested two predictions of additivity: (1) that equilibrium colors should be intensity-invariant (the so-called “scalar invariance law”), and (2) that mixtures of equilibrium hues should also be in equilibrium (“the additivity law”). Another consequence of these predictions is that the chromatic response or valence function should be a linear combination of the cone fundamentals.⁷⁵ For R/G opponency, Larimer, Krantz, and Cicerone³⁰⁷ found that the spectral position of the blue and yellow equilibrium hues were intensity-invariant over a range of 1 to 2 log units, and that mixtures of red/green equilibrium hues remained in red/green equilibrium. These results suggest that R/G color opponency is additive. Additivity for the R/G mechanism has also been reported in other studies^{310,311} and is consistent with the R/G valence functions being linear combinations of the cone spectral sensitivities.^{12,286,312}

Several studies, however, suggest that linearity fails for the R/G mechanism. The curvilinear unique blue vector^{16,305} implies that mixtures of lights that individually appear unique blue and therefore in R/G equilibrium will not be in equilibrium when they are mixed—as is required for additivity. Similarly, the nonplanar red-green color categories²⁸⁹ also imply failures of additivity. Ayama, Kaiser, and Nakatsue³¹³ carried out an additivity test for red valence and found that while additivity was found for some pairs of wavelength within separate R valence lobes (400 paired with 440 nm and 610 paired with 680 nm) failures were found between lobes (400 paired with 680 nm) for 3 out of 4 observers. Ingling et al.³¹⁴ found that the short-wavelength shape of the R valence lobe was dependent on the technique used to measure it. If the redness of the violet light was assessed by matching rather than by cancellation, the estimate of redness was as much as 30 times less. Ingling, Barley, and Ghani²⁸⁷ analyzed previous R/G hue cancellation and hue-scaling data and found them to be inconsistent with the linear model.

By contrast, most evidence suggests that the B/Y color-opponent mechanism is nonlinear. Larimer, Krantz, and Cicerone³⁰⁹ found that equilibrium green was approximately intensity-invariant, but that equilibrium red was not, becoming more bluish-red as the target intensity was increased. Moreover, although mixtures of B/Y equilibrium hues remained in equilibrium as their overall intensity changed, the equilibria showed an intensity-dependence. Other results also suggest that B/Y color opponency is nonlinear. The unique green and red vectors are not colinear in two dimensions of color space^{16,305} and the unique green and unique red planes are not coplanar.^{22,289} These failures of colinearity and coplanarity suggest that the B and the Y valence function have different spectral sensitivities, but the same neutral or balance point. The failure of additivity is consistent with the fact that the B/Y valence functions cannot be described by a linear combination of the cone fundamentals; instead some form of nonlinearity must be introduced.^{286,309} Elzinga and de Weert³¹⁵ found failures of intensity-invariance for equilibrium mixtures and attributed the failures to a nonlinear (power) transform of the S-cone input. Ayama and Ikeda³¹⁶ found additivity failures for combinations of several wavelength pairs. Ingling, Barley, and Ghani²⁸⁷ found that previous B/Y hue cancellation and hue-scaling data were inconsistent with the linear model. Knoblauch and Shevell³¹⁷ found that the B/Y nonlinearity might be related to the signs of the cone signals changing with luminance.

Interestingly, the B/Y mechanism behaves linearly in protanopes and deutanopes, which suggests that the nonlinearity might depend in some way on L–M.^{318,319}

Why linearity? This question of whether color-opponent mechanisms are linear or nonlinear was originally formulated when it was still considered plausible that cone signals somehow fed directly into color-opponent channels without modification. Given our current knowledge of receptor adaptation and postreceptor circuitry, that hope now seems somewhat optimistic. Receptor adaptation is an essentially nonlinear process, so that unless the adaptive states of all three cone types are held approximately constant (or they somehow change together as perhaps in the case of invariant hues), linearity will fail.

Tests that have found linearity have used elaborate experimental strategies to avoid the effects of adaptation. Larimer, Krantz, and Cicerone,^{307,309} for example, used dark-adapted conditions and presented the stimuli for only 1 second once every 21 seconds. This protracted procedure succeeded in achieving additivity in the case of the red/green equilibrium experiments but not in the case of the blue/yellow ones.

It would be a mistake to conclude from results like these that color opponency and color appearance are, in general, additive with respect to the cone inputs. Experiments like those of Larimer, Krantz, and Cicerone^{307,309} demonstrate additivity, but only under very specific conditions. While the results are important, because they provide information about the properties of the “isolated” postreceptor mechanisms, they are not necessarily relevant to natural viewing conditions under which the mechanisms may well behave nonlinearly with respect to the cone inputs.

Bezold-Brücke Effect and Invariant Hues Hue depends upon light intensity. Although a linear explanation of such hue changes has been suggested (see, e.g., Fig. 4 on p.74 of Hurvich²⁸³), it now seems clear that these changes reflect underlying nonlinearities in the visual response. These nonlinearities arise in part because the cone outputs are nonlinear functions of intensity, but they may also reflect inherent nonlinearities within the color-appearance mechanisms themselves and the way in which they combine cone signals. The dependence of hue on intensity is revealed clearly by the Bezold-Brücke effect,^{320,321} which is illustrated in Fig. 33 using prototypical data from Purdy.³²² As the intensity of a spectral light is increased, those with wavelengths shorter than 500-nm shift in appearance toward that of an invariant blue hue of approximately 474 nm, while those with wavelengths

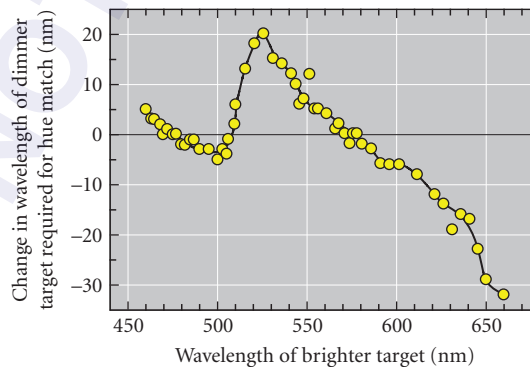


FIGURE 33 Bezold-Brücke effect. Data replotted from Fig. 1 of Purdy³²² illustrating the changes in apparent hue as the light level increases, a phenomenon known as the Bezold-Brücke hue shift. The graph shows the change in the wavelength of a 100-td target from that of the 1000-td target required to match the hue of the 1000-td target plotted as a function of the wavelength of the more intense target.

longer than about 520 nm shift in appearance toward that of the invariant yellow of approximately 571 nm. Intermediate spectral colors shift to an invariant green of approximately 506 nm.

The invariant hues are usually assumed to coincide with the unique hues of color-opponent theory,²⁸³ yet when the coincidence is tested, small discrepancies are found in some^{284,322,323} but not all^{307,309} studies (see also Ref. 296). Although invariant hues are often interpreted as being the zero crossings of cone-opponent mechanisms,¹³⁶ they are also consistent with models that incorporate just receptor adaptation.^{323,324} As Vos³²³ pointed out, if hue invariance and unique hues reflect different underlying processes, their approximate agreement may serve an important purpose.

Color Appearance and Chromatic Adaptation In the Bezold-Brücke effect, the level of adaptation is varied by increasing or decreasing the intensity of spectral lights of fixed wavelength. Under such conditions, as noted in the previous section, the appearances of some wavelengths are invariant with intensity. If, however, the state of chromatic adaptation is varied by superimposing the spectral lights on a chromatic background, their color appearance will change. In general, the appearances of targets superimposed on an adapting background move away from the appearance of the adapting wavelength. Thus, a target that appears yellow when viewed alone will appear reddish if superimposed on a middle-wavelength adapting background or greenish if superimposed on a long-wavelength one.^{325–327} Similarly, the appearance of a superimposed target of the *same* wavelength as the adapting background will move toward the appearance of the equilibrium yellow hue. Indeed, on intense adapting fields of 560, 580, and 600 nm, Thornton and Pugh¹⁰⁷ showed that the spectral position of equilibrium yellow for a superimposed target coincided with the adapting field wavelength. However, this agreement breaks down on longer wavelength fields. Long-wavelength flashes presented on long-wavelength fields appear red, not equilibrium yellow.^{65,328,329}

Historically, changes in color appearance with adaptation were thought to be consistent with von Kries adaptation,³³⁰ in which the gain of each photoreceptor is attenuated in proportion to the background adaptation. However, von Kries adaptation cannot completely account for the changes in color appearance^{77,331–333} (but see Ref. 334). In particular, asymmetric color matches, in which targets presented on different adapting backgrounds are adjusted to match in hue, do not survive proportional changes in overall intensity—as they should if the von Kries coefficient law holds.³³⁵

Jameson and Hurvich³³⁶ proposed that the adapting field alters color appearance in two ways. First, by changing the cone sensitivities in accordance with von Kries adaptation, and second by an additive contribution to the appearance of the incremental mixture. For example, a red background tends to make an incremental target appear greener because it selectively adapts the L- cones, but it also tends to make it appear redder because it adds to the target. Color-appearance models of this sort remain in vogue (see Ref. 337). Figure 15 shows a case where a model with both gain control at the cones and an additive term accounts for much but not all of the variance in an asymmetric matching data set.

The details of this two-stage theory have been contentious, despite the fact that all authors are in essential agreement that much of the effect of the background is removed from the admixed background and test. Walraven,^{325,338} using 1.5° diameter targets superimposed on 7° diameter backgrounds, determined the equilibrium mixtures of incremental 540 and 660 nm targets on various 660 nm backgrounds. He found that the equilibria were consistent with von Kries adaptation produced by the background, but that the color appearance of the background was completely discounted, so that it made no additive contribution to the appearance of the incremental target. By contrast, Shevell, using thin annuli or transiently presented 150-ms targets that were contiguous with the background, found that the background is not discounted completely and instead makes a small additive contribution to the target.^{326,327,339} Walraven³³⁸ suggested that Shevell's results reflected the fact that his conditions hindered the easy separation of the incremental target from the background. Yet, later workers came to similar conclusions as Shevell.^{340–342} However, although the background is not entirely discounted, its additive contribution to color appearance is much less than would be expected from its physical admixture.^{327,343} Thus, a discounting mechanism must be playing some role. Indeed, the failure of complete discounting is small compared to the magnitude of the additive term, so that the failure is very much a second-order effect.

Color Appearance and Chromatic Detection and Discrimination Several attempts have been made to link explicitly the cone-opponent effects evident in field sensitivity and field adaptation detection experiments with changes in color appearance measured under the same conditions.

On a white, xenon field, Thornton and Pugh¹⁶⁴ showed that the L–M cone-opponent detection implied by the Sloan notch is consistent with the suprathreshold color opponency implied by red–green hue equilibria. Thus, the wavelength of the notch, which is assumed to occur when the L–M cone-opponent signals are equal, coincided with an R/G equilibrium hue. A comparable convergence was found for 430- and 570-nm mixtures and the suprathreshold yellow–blue equilibrium locus.

S-cone thresholds on bichromatic mixtures of spectral blue and yellow fields are subadditive (see Fig. 26). According to the Pugh and Mollon²³¹ cone-opponent model, the greatest sensitization should occur when the opposing signals at the second site are in balance; that is, when the S- and L+M-cone signals are equal and opposite. In an attempt to link these cone-opponent effects to the color opponency of Hering, two groups have suggested that the background mixture that yields the lowest threshold, which is therefore the balance point of the cone-opponent at the second site, should also be the mixture that is in blue–yellow equilibrium (i.e., the mixture that appears neither yellow nor blue).

Pugh and Larimer²²⁸ investigated the detection sensitivity of π_1/π_2 on field mixtures that appeared to be in yellow–blue equilibrium. They reasoned that if such mixtures are also null stimuli for the cone-opponent S–(L+M) second site, then detection sensitivity should depend only on first-site adaptation. Their results were consistent with their hypothesis, since field mixtures in yellow–blue equilibrium never produced superadditivity, a characteristic assumed to be due to second-site adaptation (see “Field Additivity” in Sec. 11.5). However, Polden and Mollon²²⁷ looked specifically at the relationship between the maximum sensitization and the contribution of the longer wavelength component of the field mixture and found that it approximately followed a $V(\lambda)$ spectral sensitivity. By contrast, the comparable equilibrium hue settings for long-wavelength fields fell much more steeply with wavelength for fields longer than 590 nm, implying that different underlying spectral sensitivities operate in the two cases.

Rinner and Gegenfurtner³⁴⁴ looked at the time course of adaptation for color appearance and discrimination and identified three processes with half lives of less than 10 ms, 40 to 70 ms, and 20 s. The slow and intermediate adaptation processes were common to color appearance and discrimination, but the fast process affected only color appearance. However, in an extensive study, Hillis and Brainard⁷⁸ compared the effects of chromatic adaptation on color discrimination and asymmetric color matching. Pedestal discrimination measurements made as a function of pedestal intensity on five different chromatic backgrounds were used to estimate the response-intensity curves on each background. These response-intensity curves were then used to predict pairs of light that should match on different pairs of backgrounds. The agreement they found between the predicted asymmetric matches and measured ones suggests that color appearance and discriminability on different uniform backgrounds are controlled by the same underlying mechanism.⁷⁸ A follow-up study reached the same conclusion about the effect of habituation on “unstructured” spatiotemporal contrast,⁸² although this study did not probe performance with test stimuli likely to tap purported higher-order chromatic mechanisms.

Overall, the results of these experiments are somewhat equivocal, but the majority of experiments find some correspondence between the effects of adaptation on color detection and discrimination on the one hand, and its effects on color appearance on the other. Given, however, that the discrimination and appearance processes have very different spectral sensitivities, this correspondence must break down under some conditions of chromatic adaptation, at least if the signals that control the adaptation of mechanisms depend strongly on the output of those same mechanisms. Of note here is that Hillis and Brainard did find clear dissociations of the effect of adaptation on discrimination and on appearance when the stimuli contained sufficient spatial structure to be perceived as illuminated surfaces.³⁴⁵ Elucidation of the nature of this dissociation awaits further investigation.

Color Appearance and Habituation The results of Webster and Mollon^{32,79} were used as an example in subsection “Appearance Field Measurements and Second-Site Adaptation” in Sec. 11.4. Briefly, they

found that contrast adaptation produced changes in appearance that were selective for the habituating axis whether that axis was in the equiluminant plane or not. This finding is a potentially important difference between the effects of habituation on color appearance and its effects on detection, since only the appearance data provides clear evidence for higher-order mechanisms that are sensitive to both chromatic and luminance modulations. The detection data measured after habituation suggest that the chromatic and luminance mechanisms behave independently.²⁴¹ Determining conclusively whether this comparison across studies reflects a dissociation between habituation effects on thresholds and appearance, or instead represents a difference between the processing of near-threshold and suprathreshold test stimuli, could be resolved by applying the logic developed by Hillis and Brainard^{78,82} to data collected for the full set of stimulus conditions studied by Webster and Mollon.^{32,79}

Webster and Mollon^{32,79} suggested two models that might account for their data. In the first, they replaced the three color mechanisms tuned to the three cardinal directions of color space with many mechanisms with tuning that varied according to a Gaussian distribution around each of the three cardinal directions. They found that they could account for the individual data by varying the standard deviations of the Gaussian distributions. In the second, they assumed that there were just three mechanisms, but that their tuning could be modified by adaptation. With adaptation, inhibition was assumed to build up between the mechanisms to decorrelate the input signals and reduce redundancy (see also Refs. 242–244). These models can also be applied to chromatic detection data.

Luminance and Brightness Brightness matching functions are broader than luminous efficiency functions measured using techniques such as HFP or MDB that produce additive results (see subsection “Sensitivity to Spectral Lights” in Sec. 11.5) and are thought to measure the sensitivity of the luminance channel.^{115,118,122,346}

Figure 34 shows examples of a brightness matching function measured by direct brightness matching (red line) and luminous efficiency measured by flicker photometry (black line) replotted from Wagner and Boynton.¹²² The brightness matching function is relatively more sensitive in the blue and orange spectral regions and slightly less sensitive in the yellow. These differences are usually attributed to a chromatic contribution to brightness but not to luminance.^{11,122,347}

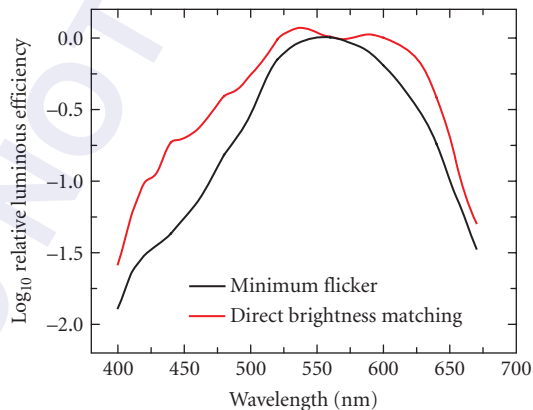


FIGURE 34 Luminous efficiency and brightness. Mean matching data for three subjects replotted from Figs. 6 and 8 of Wagner and Boynton¹²² obtained either using flicker photometry (black line) or using direct brightness matching (red line). The discrepancies are consistent with a chromatic contribution to brightness matching but not to flicker photometric matching.^{120,122}

Mechanisms of Color Constancy

The literature on mechanisms of color appearance is closely tied to a phenomenon known as “color constancy.” The need for color constancy stems from two observations. First, our normal use of color appearance is to associate colors with objects. Second, the stimulus reflected to the eye from a fixed object depends not only on the object’s intrinsic reflectance properties but also on the spectrum of the illumination. When the illuminant changes, so does the spectrum of the reflected light. If the visual system did not compensate for this change to stabilize the appearance of objects, it would not be possible to use color appearance as a reliable guide to object properties.

Empirically, color constancy is studied in much the same way as color appearance, using techniques of asymmetric matching, hue cancellation, or hue scaling. Papers published under the rubric of color constancy, however, have tended to employ stimulus conditions designed to model scenes consisting of illuminated objects and to emphasize context changes induced by manipulating the illuminant, rather than the simpler test-stimulus-on-background configurations favored for working out mechanistic accounts of color appearance.^{348–353} The evidence supports the empirical generalization that when the illuminant is varied in a manner typical of the variations that occur in the natural environment, the visual system adapts to compensate for the physical change in reflected light, and thus to stabilize color appearance. That is, constancy in the face of naturally occurring illumination changes is good.^{352,353} At the same time, it is important to note that constancy is not perfect: objects do not appear exactly the same across illumination changes. Moreover, other scene manipulations, such as those that change the reflectance of objects near the one of interest, also affect color appearance, even in rich scenes.³⁵⁴ Figure 1 shows examples of this latter type of effect for simple stimulus configurations; these may be regarded as failures of constancy. Several recent reviews treat empirical studies of color constancy in more detail.^{355–357}

Considerable theoretical attention has been devoted to modeling human color constancy. Much of this work is in the computational tradition, in which the theorist starts not with experimental data but rather by asking how constancy could in principle be achieved (or approximated), given the information about object reflectance properties that is available in the retinal image.^{358–363}

With a few exceptions, work in the computational tradition does not connect explicitly with the mechanistic account of processing that we have developed in this chapter. Rather the work seeks to elucidate the sources of information in an image that could contribute to constancy. This is then used both to guide experiments³⁵⁴ as well as in models whose purpose is to explicitly predict color appearance.^{364–367} A particular promise of this approach is that it generalizes naturally as scene complexity increases, allowing (for example) predictions of the effects of manipulations in complex three-dimensional scenes using the principles developed for simpler scene configurations.^{368–370}

A limitation of computational models is that they do not, in and of themselves, provide much insight about how neural mechanisms might act to provide constancy. Not surprisingly, then, there is also an important line of theoretical work that attempts to understand constancy in terms of the action of the first- and second-site mechanisms of adaptation that we have reviewed in this chapter.

Land’s retinex theory^{371–374} may be understood as a computational model that incorporates first-site cone-specific gain control to stabilize the representation of object color. Several authors have shown that for natural illuminant changes, such first-site gain control provides excellent, although not perfect, compensation for the effect of illumination changes on the light reflected from natural objects as long as the gains are set correctly for the scene illumination.^{375–378} Moreover, asymmetric matching data from experiments conducted using rich, naturalistic stimuli are well fit by postulating first-site gain control.³⁵² Webster and Mollon³⁷⁹ (see also Ref. 380) extended this general line of thinking by showing that following first-site cone-specific adaptation with a second-site process that is sensitive to image contrast can improve the degree of compensation for the physical effect of illuminant changes. As we have reviewed in subsections “Appearance Field Measurements and Second Site Adaptation” in Sec. 11.4 and “Color Appearance and Habituation” in Sec. 11.5, it has also been established that contrast adaptation affects color appearance,^{32,79,379,381} which provides a connection between human performance and the observation that contrast adaptation can help achieve color constancy.

A number of authors^{382–384} have emphasized that it is clarifying, in the context of understanding constancy and more generally adaptation, to distinguish two questions. First, what parameters of visual processing (e.g., first-site gains, subtractive terms, second-site gains) are affected by contextual factors. That is, what can adapt? Second, what are the particular contextual features that act to set each of the adaptable parameters? Our understanding of these aspects is less developed. For example, it is not obvious how theories based on data obtained with spatially uniform adapting fields should be generalized to predict the adapted state for image contexts that have rich spatial and temporal structure. Simple hypotheses for generalization include the idea that a weighted average of the cone coordinates of the image might play the role of the uniform field, as might the most luminous image region. Ideas of this sort are implicit, for example, in Land's retinex theory. Explicit tests of hypotheses of this nature for rich scenes, however, have not revealed simple scene statistics sufficient to predict the visual system's state of adaptation.³⁵⁴ Indeed, understanding what image statistics are used by the visual system to achieve constancy, and connecting the extraction of these statistics explicitly to mechanistic sites of adaptation is a primary focus of current research on constancy. In this regard, the computational work may provide critical guidance, as at the heart of each computational theory is an analysis of where in the retinal image the most valuable information about the illumination may be found. Zaidi³⁸⁵ and Golz and MacLeod,³⁸⁶ for example, translate analyses of image statistics that carry information about the illuminant into mechanistic terms.

Color and Contours

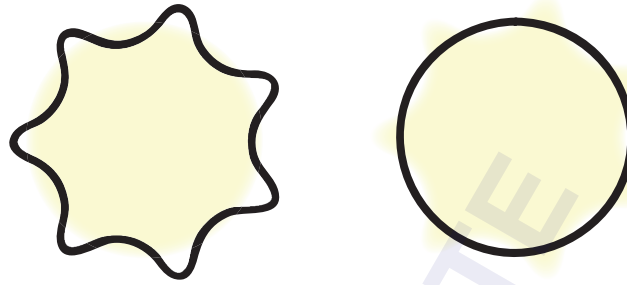
In the natural world, differently colored regions are usually delimited by contours across which there is also a luminance difference. If the luminance difference is removed, the contour becomes indistinct, and it becomes more difficult to discriminate the chromatic differences between the two regions. For example, Boynton, Hayhoe, and MacLeod³⁸⁷ showed that color discrimination is severely impaired if the difference between two regions is dependent on only the S cones, and Eskew and Boynton³⁸⁸ showed that the contour remains indistinct even for chromaticity differences up to twice threshold. For two regions separated by a border defined by S cones (i.e., if the regions are tritan pairs), "melting" of the border³⁸⁹ and "chromatic diffusion" between the regions³⁸⁸ have been reported.

Figure 35a shows two examples of chromatic regions filling in an area bordered by a luminance contour. These are versions of the "Boynton illusion" (see p. 287 of Ref. 15 and <http://www.yorku.ca/eye/boynton.htm>). The yellow areas, which are discriminated from the white background mainly by S cones, appear to fill-in the black borders as you move further away from the figure. The filling-in illustrates the tendency of luminance contours to constrain the spatial extent of signals mediated by the chromatic pathways.

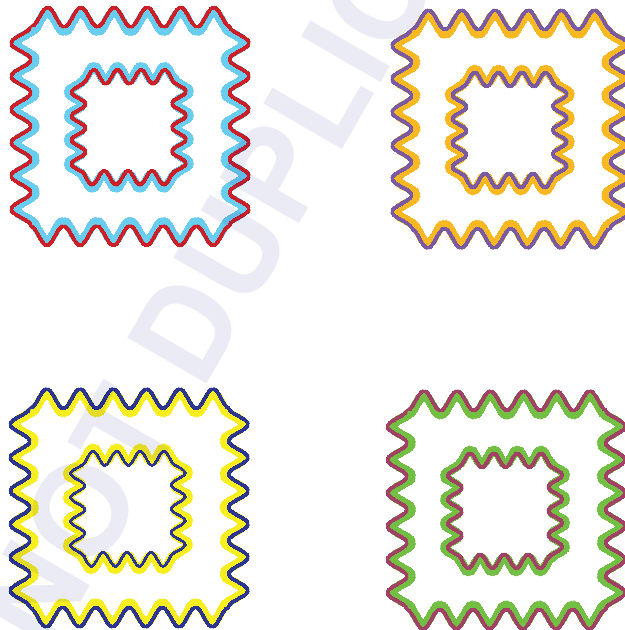
Figure 35b shows four examples of a related illusion, which is known as the watercolor illusion or effect.^{390–392} If an area is delineated by a darker outer chromatic contour flanked by a brighter inner chromatic contour, then the brighter color spreads faintly into the inner area.³⁹³ This faint coloration resembles a "watercolor" wash.

The "Gap" Effect and Luminance Pedestals The gap effect is a well-known phenomenon described by Boynton, Hayhoe, and MacLeod.³⁸⁷ A small gap introduced between two juxtaposed test fields improves their discriminability if the fields differ only in S-cone excitation, impairs it if they differ only in luminance, and can result in a small improvement if they differ only in L- and M-cone excitation at constant luminance (i.e., in L–M excitation).^{387,394} The small gap that separates the tritan pair and improves discriminability can be produced either by a luminance difference or by an L–M chromatic difference.³⁹⁵

The improvement in chromatic discriminability caused by gaps can be related to the improvement in chromatic sensitivity when a luminance pedestal or contour is coincident with the chromatic stimulus^{169,177,194,269,271,272}—the so-called "crossed" pedestal facilitation described above (see subsection "Pedestal Experiments" in Sec. 11.5). Moreover, as also noted above, the crossed facilitation survives if the luminance pedestal is replaced by a thin ring¹⁶⁹ (see also Ref. 177). Montag³⁹⁴ compared the gap



(a)



(b)

FIGURE 35 Boynton and watercolor illusions. (a) Two examples of the Boynton illusion. In the left picture, a star-shaped black outline is superimposed on a circular yellow area, while in the right picture a circular black outline is superimposed on a star-shaped yellow area. As you move away from the picture, the yellow areas appear to fill the black outlines. These are variations of an illusion attributed to Robert M. Boynton by Peter Kaiser, see: <http://www.yorku.ca/eye/boynton.htm>. (b) Four examples of the watercolor illusion.^{390–392} Each square is outlined by a pair of juxtaposed sinuous contours. The similar colors of the inner contour of each larger square and the outer contour of each smaller square fill in the intervening space with a faint color (which looks like a watercolor wash). Inspired by a version of the watercolor illusion designed by Akiyoshi Kitaoka, see: <http://www.psy.ritsumei.ac.jp/~akitaoka/watercolorillusionsamples.jpg>.

effect with facilitation produced by adding thin lines of the same orientation to gratings at half-cycle intervals. He found that the detection of S-cone gratings was facilitated by placing dark lines at the midpoints between the peaks and troughs of the gratings (i.e., where the luminance of the grating is equal to the mean luminance of the stimulus—zero crossings of the spatially varying factor in the stimulus), but that the facilitation declined to zero as the lines were moved toward the peaks and troughs. A comparable but smaller effect was found for equiluminant L- and M-cone modulated gratings, while the detection of luminance gratings was always impaired. Montag suggested that the gap effect and pedestal enhancement may be complementary effects.³⁹⁴ Indeed, Gowdy, Stromeyer, and Kronauer³⁹⁶ reported that the facilitation of L–M detection by luminance pedestals is markedly enhanced if the luminance grating is square-wave rather than sinusoidal. They argued that the square-wave produces an abrupt border across which the L–M mechanism is somehow able to compare chromatic differences.

Color Appearance and Stabilized Borders The importance of contours on color appearance is made clear in experiments in which the retinal image is partially or wholly stabilized. When the retinal image is wholly stabilized, the color and brightness of the stimulus fades until it is no longer seen.³⁹⁷ The color appearance of partially stabilized images in which stabilized and unstabilized borders are combined can be entirely independent of the local quantal absorptions. Krauskopf³⁹⁸ found that when the border between a disc and a surrounding annulus is stabilized, and the outer border of the annulus unstabilized, the disc “disappears” and its color is filled in by the color of the annulus. When, for example, the border between a green, 529-nm annulus, and an orange, 640-nm disc, is stabilized, the disc appears green. These changes are not restricted to the stabilized disc. Unstabilized targets presented on stabilized discs also change their color appearance as the color of the disc fills in to the color of the annulus.³⁹⁹ For example, if the annulus is yellow and the stabilized disc is either red or green, both discs fill in to take on the yellow color of the annulus. Unstabilized yellow targets presented in the center of the stabilized discs, however, take on the color appearance complementary to actual color of the disc, so that the yellow target on a green disc appears red, and the yellow target on a red disc appears green. Such changes are consistent with the local contrast signals at the target edges determining the filled-in appearance of the yellow target (see Experiment 4 of Ref. 400).

It has also been reported that stabilized boundaries can produce colors that appear reddish-green and yellowish-blue, and so violate the predictions of the opponent-colors theory that opposed colors (red and green, or yellow and blue) cannot be perceived simultaneously. These “forbidden colors” were produced by presenting two vertical stripes side-by-side, with their common border stabilized, but their outer borders unstabilized. When the juxtaposed stripes were red and green, most observers reported that the merged stripes appeared reddish-green or greenish-red, and when the stripes were blue and yellow most observers reported that they appeared bluish-yellow⁴⁰¹ (see Experiment 3 of Ref. 400). Forbidden colors were also reported by some subjects in another study, but only when the stripes were equiluminant.⁴⁰²

The results obtained using stabilized and unstabilized borders add to the view that color appearance need not be a local retinal phenomenon. A crucial question, then, is whether the changes in color appearance caused by image stabilization also influence phenomena that are ostensibly more low-level, and thus more likely to be local, such as detection. Two experiments have addressed this question. Nerger, Piantanida, and Larimer⁴⁰³ found that when a red disk was surrounded by a yellow annulus, stabilizing the edge between the two fields on the retina caused the yellow to fill in, making the disk appear yellow too. This filling-in affected the color appearance of small tests added to the disk, but it did not affect their increment threshold. The results for S-cone flicker detection, somewhat surprisingly, suggest that filling-in can change S-cone flicker sensitivity. Previously, Wisowaty and Boynton,⁴⁰⁴ using flickering tritan pairs to isolate the S-cone response, had found that yellow background adaptation reduces the S-cone modulation sensitivity compared to no background adaptation. Piantanida⁴⁰⁵ showed that the reduction in S-cone flicker sensitivity caused by a yellow field could also be caused by a dark field that only appeared yellow because its outer border with a yellow annulus was stabilized. In a related experiment, Piantanida and Larimer,⁴⁰⁰ compared S-cone modulation sensitivities on yellow and green fields, and found that the S-cone modulation sensitivity was

lower on the yellow field. Again, this sensitivity difference depended upon the field appearance rather than upon its spectral content. Thus, the reduced sensitivity was found not only on the yellow field but also on a green field that appeared yellow because its border with a surrounding yellow annulus was stabilized. Similarly, the increased sensitivity was found not only on the green field but also on a yellow field that appeared green because its border with a surrounding green annulus was stabilized. The differences between studies may indicate that different sites of limiting noise act to limit detection in the different tasks.

Contours and Aftereffects Daw⁴⁰⁶ showed that the saliency of colored after-images depended upon the presence of an aligned monochrome image. Thus, adaptation to a colored image produces a clear afterimage if the after-image is aligned with a monochrome image with the same contours, but produces a poor after-image if the two are misaligned. Figure 36 demonstrates this effect.

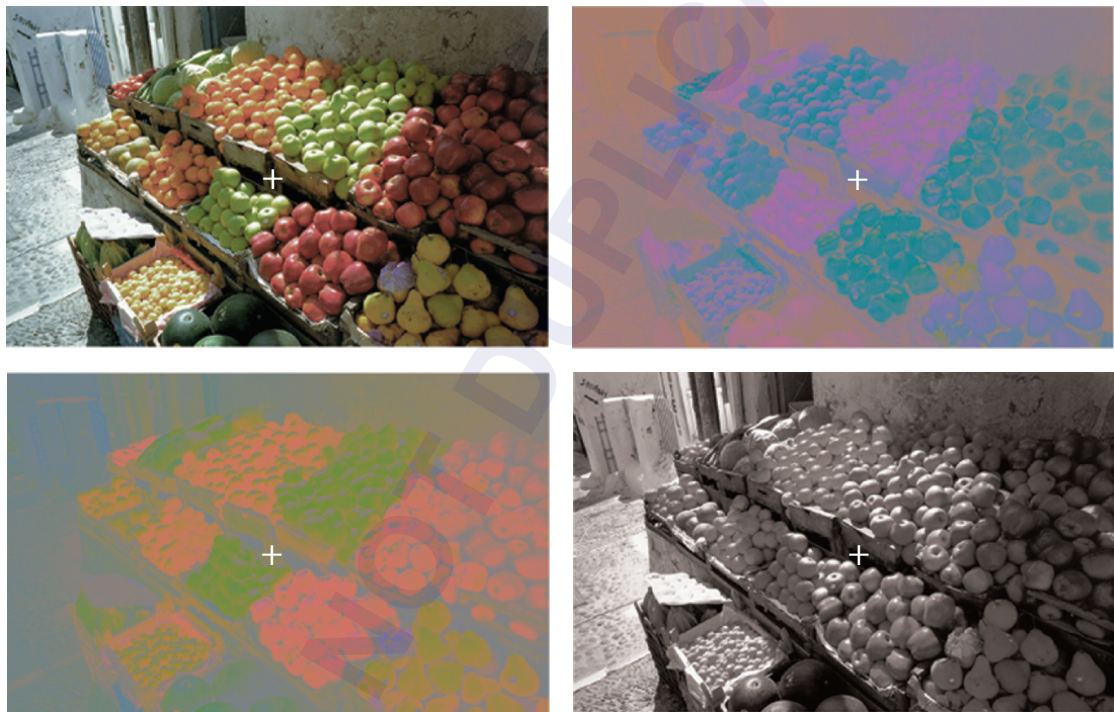


FIGURE 36 Color and contour. Four different images of a fruit stall are shown. The full color image (top left) has been decomposed into its chromatic (bottom left) and luminance components (bottom right). The image on the top right is the complementary or inverse of the chromatic image. Notice that the details of the image are much better preserved in the luminance image than in the chromatic image. The chromatic information, although important for object discrimination and identification, is secondary to the perception of form. Color fills-in the picture delineated by the luminance information. This can be demonstrated by fixating the cross in the center of the complementary chromatic image (top right) for several seconds, and then diverting your gaze quickly to the cross in the center of the luminance image (bottom right). You should see a correctly colored version of the picture. Notice that if you shift your gaze slightly, so that the after-image no longer aligns precisely with the luminance image, the color disappears. You can also try adapting to the chromatic image (bottom left) and repeating the experiment. The effects are stronger with projected CRT versions of the images and successive presentation, which can be seen at <http://www.cvrl.org>. Color scene after Fig. 1.14A of Sharpe et al.⁴⁷⁷ from an original by Minolta Corp. A version of this aftereffect demonstration was first described in 1962 by Nigel Daw.⁴⁰⁶ Other examples of this illusion along with instructions about how to produce comparable images can be found at: http://www.johnsadowski.com/color_illusion_tutorial.html.

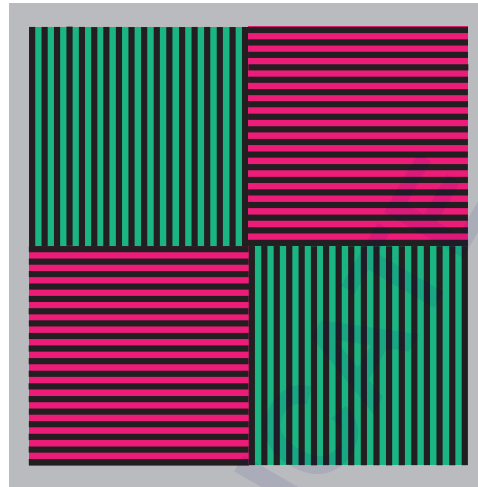
Impressive demonstrations of the influence of contours on after-images have recently been produced by van Lier and Vergeer⁴⁰⁷ and van Lier, Vergeer, and Anstis.⁴⁰⁸ Their demonstrations show that colored stimuli can produce different after-image colors at the same retinal location, depending on the positioning of contours presented with the after-image. In one version, observers adapt to a plaid pattern made up of blue, green, orange, and purple squares, and then view, together with the after-image, either horizontal or vertical contours aligned with the borders of the squares. The color of the after-image changes with the orientation of the contours, because the “mean” after-image along the horizontal rows and along the vertical columns is different, thanks to the arrangement of the adapting squares in the plaid. The demonstration can be seen at: <http://www-psy.ucsd.edu/~sanstis/SAai.html>.

McCollough Effect The McCollough effect is a well-known orientation-contingent color after-effect.⁴⁰⁹ After prolonged viewing of colored “adapting” gratings of, say, vertical or horizontal orientation, neutral, black-and-white test gratings of the same orientation take on the complementary hue of the adapting grating. For example, following adaptation to red horizontal and green vertical gratings, horizontal and vertical black-and-white gratings appear, respectively, slightly greenish and reddish (see Fig. 37).

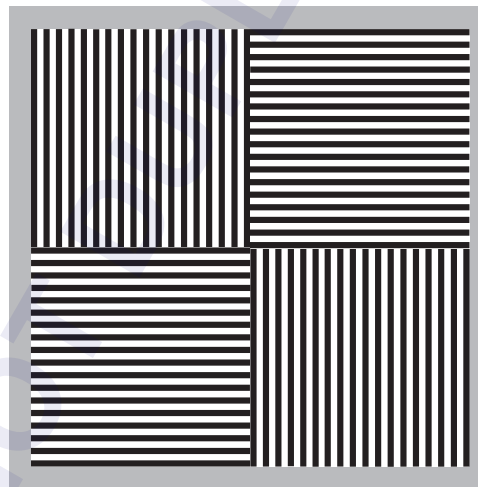
The McCollough effect can persist undiminished for many days provided that test gratings are not viewed.⁴¹⁰ This prolonged effect is distinct from colored afterimages, such as the effect demonstrated in the previous subsection, which decline relatively quickly. Vul, Krizay, and MacLeod⁴¹¹ have recently identified two processes in the generation of the McCollough effect. The first process behaves like a leaky integrator, and produces an effect that declines exponentially with a time constant of roughly 30 s. The second process behaves like a perfect integrator, and produces an effect that shows no decay.⁴¹¹ Several lines of evidence suggest that the McCollough effect is generated relatively early in the visual pathway. First, the McCollough effect exhibits little interocular transfer.^{409,412,413} Second, it can be produced by alternating red-stripped and green-stripped adapting patterns at rates as high as 50 Hz, frequencies at which the colored patterns cannot be consciously perceived.⁴¹⁴ Third, it is dependent on the spectral content of the inducing stimuli, not their color appearance.⁴¹⁵ These findings suggest that the “adaptation” has an effect that occurs at a relatively early visual site that precedes binocularity, is able to respond high-frequency flicker, and retains some veridical information about the spectral content of lights. However, the contingency of the McCollough effect on grating orientation suggests that the site cannot be earlier than primary visual cortex, V1, where orientation selectivity is first clearly established.^{416,417} There is an extensive literature on the McCollough effect (e.g., Refs. 418–423).

Multiplexing Chromatic and Achromatic Signals Information about color and contour may be transmitted—at least in part—by the same postreceptoral pathways. The ways in which these mixed signals are demultiplexed may give rise to some of the color and border phenomena described in the previous subsections. The evidence for the color-luminance multiplexing comes mainly from knowledge of the spatial properties of neuronal receptive fields in the early visual pathway. Cone-opponent, center-surround mechanisms that are also spatially opponent encode not only chromatic information, which is dependent upon the difference between the spectral sensitivities of their center and surround, but also “achromatic” information, which is dependent on the sum of their spectral sensitivities. For color, the center and surround behave synergistically producing a low-pass response to spatial variations in chromaticity, but for luminance they behave antagonistically producing a more band-pass response to spatial variations in luminance.^{149–151} This type of multiplexing can, in principle, apply to any mechanism that is both chromatically and spatially opponent, but it is most often considered in the context of P-cells or midget ganglion cells, which make up as much as 80 percent of the primate ganglion cell population.⁴²⁴ These cells are chromatically opponent with opposed L- and M-cone inputs,³⁶ but also respond to spatial differences in luminance.^{425,426} In the fovea, they may be chromatically opponent simply by virtue of having single L- or M-cone centers^{427,428} and mixed surrounds.^{429–431} How segregated the L- and M-cone inputs are to P-cell surrounds and to P-cell centers in the periphery remains controversial.^{432–440}

The multiplexing of color and luminance signals may just be an unwanted and unused artifact of having both spatial and chromatic opponency in P-cells. Indeed, Rodieck⁴⁴¹ has argued that the L–M



(a)



(b)

FIGURE 37 McCollough effect. View the upper colored image (a) for several minutes letting your gaze fall on different colored areas for several seconds at a time. Look next at the lower monochrome image (b).

opponent system is mediated instead by a population of so-called Type II cells with coincident centers and surrounds, which are chromatically but not spatially opponent. For the multiplexing of color and luminance signals in cells with concentric center-surrounds to be unambiguously useful to later visual stages, the signals must be decoded. Several decoding strategies have been suggested, most of which depend on the chromatic signal being spatially low-pass filtered by the center-surround interaction, and the luminance signal being spatially band-pass filtered. Because of this filtering, signals between adjacent mechanisms will, in general, change slowly if the signals are chromatic and rapidly if they are achromatic. Thus, the chromatic and luminance signals can be decoded by spatially

low-pass or band-pass filtering, respectively, across adjacent mechanisms.^{149,150} In such a scheme, low spatial-frequency luminance information and high spatial-frequency chromatic (edge) information is lost. The high-frequency luminance or edge information can be used both in “filling-in” the low spatial frequency luminance information, and to define missing chromatic edges. Several well-known visual illusions, such as the Craik-O’Brien-Cornsweet illusion (see Ref. 442), are consistent with filling-in, while others, such as the Boynton illusion (see Fig. 35a), are consistent with luminance edges defining chromatic borders.

Simple mechanisms for decoding the luminance and chromatic signals that difference or sum center-surround chromatically opponent neurons have been proposed.^{153,443–445} Such mechanisms are, in fact, mathematically related to decoding using matched spatial filters.⁴⁴⁵ In one version, spatially superimposed opponent cells with different cone inputs are either summed or differenced.¹⁵³ The band-pass luminance signal is decoded by summing $+Lc-Ms$ and $+Mc-Ls$ to give $+(L+M)c-(L+M)s$ and by summing $-Mc+Ls$ and $-Lc+Ms$ to give $-(L+M)c+(L+M)s$, both of which are also bandpass and, in principle, achromatic (where c = center, s = surround). The low-pass chromatic signal is decoded by differencing $+Lc-Ms$ and $-Mc+Ls$ to give $(+L-M)c,s$ and by differencing $-Lc+Ms$ and $+Mc-Ls$ to give $(+M-L)c,s$, both of which have spatially coincident centers and surrounds. These combination schemes are illustrated in Fig. 38. As Kingdom and Mullen⁴⁴⁶ point out, superimposed centers with

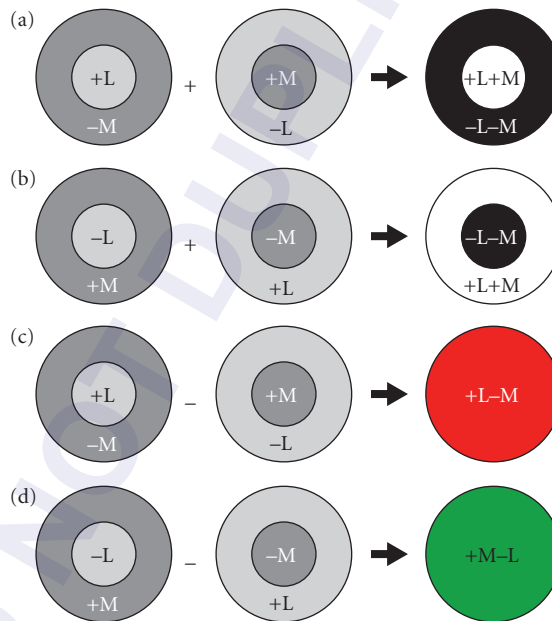


FIGURE 38 Double-duty L–M center-surround opponency. Achromatic and chromatic information can be demultiplexed from LM-cone-opponent center-surround mechanisms by summing or differencing different types of mechanisms. Spatially opponent On and Off center achromatic mechanisms can be produced by summing, respectively, either (a) L and M On-center mechanisms, or (b) L and M Off-center mechanisms. Spatially nonopponent L–M and M–L chromatic mechanisms can be produced by differencing, respectively, either (c) L and M On-center mechanisms, or (d) L and M Off-center mechanisms. See Fig. 24 of Lennie and D’Zmura⁴⁷⁸ and Fig. 2 of Kingdom and Mullen.⁴⁴⁶

different cone inputs must be spatially offset, simply because two cones cannot occupy the same space. They found that such offsets produce significant crosstalk between the decoded luminance and chromatic outputs. However, this problem is mitigated somewhat by the blurring of the image by the eye's optics, and, in the case of chromatic detection, by neural blurring (see Fig. 7 of Ref. 200).

In another version of the decoding mechanism, De Valois and De Valois¹⁸ added S-cone inputs to generate spectral sensitivities more consistent with opponent-colors theory (see subsection "Spectral Properties of Color-Opponent Mechanisms" in Sec. 11.5). Billock,⁴⁴⁵ by contrast, opposed P-cells of the same sign in both the center and surround, thus directly generating double-opponent cells. This can be achieved by differencing (e.g., $+L-M$ in the center differenced from $+L-M$ in the surround) or by summing (e.g., $+L-M$ in the center summed with $-L+M$ in the surround).

The importance of multiplexing in spatially and chromatically opponent mechanisms, and, in particular, in P-cells, remains controversial. Clearly, the P-cells, which make up 80 percent of the ganglion cells and have small receptive fields, must be important for spatial, achromatic vision, but their importance for color vision, while likely, is not yet firmly established. Speculative arguments have been based on the phylogenetically recent duplication of the M/L-cone photopigment gene, which is thought to have occurred about 30 to 40 million years ago after the divergence of Old and New World primates.⁴⁴⁷ The typical argument is that the color opponency in P-cells is parasitic on an "ancient" system that before the duplication detected only luminance contrast.⁴⁴⁸ However, this raises the possibility that the "decoding" of color information might be done simply to remove it and so improve image fidelity rather than to provide a color signal.

Many psychophysical experiments have attempted to elucidate the properties of the P-cells by employing equiluminant stimuli.⁴⁴⁹ The idea is that these stimuli silence the M-cells and thus reveal the properties of the P-cells. If we accept this logic, there is still, of course, an important limitation to these experiments: they do not, by their design, provide any information about how the P-cells respond to luminance modulation. Thus, these experiments can provide only a partial characterization of the P-cells' response.

11.6 CONCLUSIONS

In Sections 11.3 and 11.4, we introduced basic mechanistic models for color discrimination and color appearance. Section 11.5 reviewed evidence that both supports the basic model and indicates areas where it is deficient. Here we summarize what we see as the current state of mechanistic models of color vision.

Low-Level and Higher-Order Color-Discrimination Mechanisms

Test methods provide clear evidence for the existence of three low-level "cardinal" postreceptoral mechanisms: $L-M$, $S-(L+M)$, and $L+M$ in detection and discrimination experiments (see subsection "Test Sensitivities" in Sec. 11.5). This simple picture is complicated by evidence that $L-M$ and $L+M$ may have small S-cone inputs (see subsection "Sensitivity to Different Directions of Color Space" in Sec. 11.5), but nevertheless the weight of the evidence suggests that just three bipolar cardinal mechanisms (or six unipolar mechanisms—see next section "Unipolar versus Bipolar Chromatic Mechanisms") are sufficient to understand the preponderance of the data. Consistent with the existence of cardinal mechanisms, near-threshold color-discrimination experiments, with one exception,²⁴¹ require only four or five unipolar detection mechanisms.^{179,197,260} Field methods (see subsection "Field Sensitivities" in Sec. 11.5) also confirm the importance of the three cardinal mechanisms. This agreement across methods provides strong psychophysical evidence for the independent existence of these mechanisms.

Test methods provide little or no evidence for the existence of "higher-order" (or noncardinal) mechanisms. Evidence for higher-order mechanisms comes almost entirely from field methods carried

out mainly in the DKL equiluminance plane using detection and habituation,²⁴¹ detection and noise masking,^{245,248,450} discrimination and habituation,^{241,255} texture segmentation and noise,⁴⁵¹ and image segmentation and noise.⁴⁵² In contrast, three experiments using detection and noise carried out with stimuli specified in cone contrast space found little or no evidence for higher-order mechanisms.^{175,179,246}

In principal, the choice of space should not influence the results. However, the choice of space used to represent the data does tend to affect the ensemble of stimuli studied, in the sense that stimuli that uniformly sample color directions when they are represented in one space do not necessarily do so when they are represented in another. Another concern is that experiments carried out at nominal equiluminance are liable to luminance artifacts, since the elimination of luminance signals is, in practice, difficult to achieve. Aside from limits on calibration precision or inaccuracies introduced by the assumption that the 1924 CIE $V(\lambda)$ function (i.e., candelas/m²) defines equiluminance, a recurring difficulty in producing truly equiluminant stimuli is that the luminance spectral sensitivity varies not only with chromatic adaptation and experimental task, but also shows large individual differences (see subsection “Luminance” in Sec. 11.5). Moreover, the luminance channel almost certainly has cone-opponent inputs (see Fig. 25). Finally, even when the luminance mechanism is silenced, the low-level L–M mechanism may have a small S-cone input, which will rotate the mechanism axis away from the cardinal axis.

These concerns notwithstanding, the results of field methods also suffer from an inherent ambiguity. Just because noise masking or habituation alters the detection spectral sensitivity does not necessarily reveal transitions between multiple higher-order mechanisms. Such changes can also result if the simple model used to characterize mechanisms fails to capture, for example, low-level interactions and nonlinearities. Indeed, Zaidi and Shapiro²⁴³ have suggested that such “failures” might actually be an adaptive orthogonalization among cardinal mechanisms that reduces sensitivity to the adapting stimulus and improves sensitivity to the orthogonal direction. Moreover, noise masking can produce rotated detection contours by changing the balance of noise at the first and second sites (see subsection “Sites of Limiting Noise” in Sec. 11.3) without the need to invoke high-order mechanisms.

Of course, it would be absurd to suppose that higher-order color mechanisms do not exist in some form. They are clearly found in more “cognitive” experiments such as Stroop or reverse-Stroop experiments, which show clearly the effects of color categorization (e.g., Ref. 453). The central question with respect to color discrimination, however, is how such mechanisms can be revealed and investigated psychophysically using conventional threshold techniques. For example, if the site of limiting noise precedes higher-order mechanisms for most stimulus conditions, this may not be possible. In that case, the study of higher-order mechanisms will require the use of appearance or other experimental tasks in addition to knowledge of the input characteristics.

On the other hand, some discrimination tasks do provide clear evidence for higher-order mechanisms. Zaidi and Halevy⁴⁵⁴ showed that the discrimination thresholds for an excursion in color space away from a background, the chromaticity of which was modulated at 0.46 Hz along a circle in the DKL equiluminant-color plane, depended on the direction of the background color change. Discrimination was consistently worse if the excursion was in the same direction as the background color change, than if it was in the opposite direction. Since this effect was independent of the actual background chromaticity, the effect is inconsistent with there being a limited number of cardinal mechanisms. More generally, Flanagan, Cavanagh, and Favreau³¹ used the tilt after-effect to investigate the spectral properties of the orientation-selective mechanisms thought to underlie the after-effect. They found that the orientation-selectivity occurs independently in each of the three cardinal mechanisms, but they also found evidence for secondary orientation-selective mechanisms at other directions in color space. In addition, they found evidence for another orientation-selective mechanism that was nonselective for color directions in any color direction. Other experiments that suggest high-order mechanisms include visual search,⁴⁵⁵ color appearance,^{32,79,456} and the motion coherence of plaid patterns.⁴⁵⁷ A unified understanding of when and why higher-order mechanisms are revealed experimentally remains an important question for the field.

Unipolar versus Bipolar Chromatic Mechanisms

The strong implication that cone-opponent and color-opponent mechanisms may be unipolar rather than bipolar came first from physiology. The inhibitory range of LGN opponent cells

(i.e., the extent of decreases in firing rate from the resting or spontaneous level) is seldom more than 10 to 20 impulses per second, compared with an excitatory range of several hundred (i.e., the extent of increases in firing rate from the resting level).⁴⁵⁸ Consequently, cells of both polarities (i.e., L–M and M–L) are required in a push-pull relationship to encode fully both poles of a cone-opponent response. By the cortex, the low spontaneous firing rates of cortical neurons makes cells effectively half-wave rectifiers, because they cannot encode inhibitory responses.^{459,460} Thus, the bipolar L–M, M–L, S–(L+M), and (L+M)–S cone-opponent mechanisms must become unipolar |L–M|, |M–L|, |S–(L+M)|, and |(L+M)–S| cone-opponent mechanisms, while the bipolar R/G and B/Y mechanism become unipolar R, G, B, and Y mechanisms.

Behavioral evidence for unipolar mechanisms at the cone-opponent level is sparse in the case of the L–M mechanism, because the opposite polarity L–M and M–L detection mechanisms have similar cone weights (see subsection “Detection contours in the L, M Plane” in Sec. 11.5), so in many ways they behave like a unitary bipolar mechanism.³⁴ Nevertheless, there is some evidence that the |L–M| and |M–L| poles behave differently. In particular, habituations to sawtooth chromaticity modulations can selectively elevate thresholds for |L–M| or |M–L| targets.^{14,461} Sankeralli and Mullen⁴⁶² raised detection threshold using unipolar masking noise and found that masking was most evident when the test and noise had the same polarity (thus, |L–M| noise selectively masks an |L–M| test, and |M–L| noise selectively masks an |M–L| test). Other evidence suggests that |L–M| and |M–L| are separate mechanisms at the fovea.^{396,463} Evidence for separate mechanisms is also apparent in the peripheral retina, where the |L–M| mechanism predominates.^{176,178,464}

There is better evidence for the existence of unipolar S-cone mechanisms. Sawtooth chromaticity modulations selectively elevate thresholds for |S–(L+M)| or |(L+M)–S| targets,^{14,461} as do unipolar |S–(L+M)| or |(L+M)–S| masks.⁴⁶² Other evidence shows that (1) the spatial summation for S-cone incremental and decremental flashes differs,⁴⁶⁵ (2) the longer wavelength field sensitivity for transient tritanopia is different for incremental and decremental flashes, which indicates different L- and M-cone inputs for the two S cone polarities,⁴⁶⁶ and (3) habituation or transient adapting flashes have differential effects on S increment and S decrement thresholds.^{467,468} There is also ample physiological and anatomical evidence that S-cone ON and OFF pathways are distinct (e.g., Ref. 2). For a recent review, see Ref. 34.

Eskew³⁴ takes the more philosophical stance that a psychophysical mechanism should be defined as a univariant mechanism, the output of which is a “labeled line.”^{469,470} This means that a bipolar cone-opponent mechanism cannot by definition be a unitary mechanism, since it can exist in two discriminable states (positive and negative).

Discrepancies between Color-Discrimination and Color-Appearance Mechanisms

Much evidence indicates that color-discrimination and color-appearance mechanisms are distinct, at least given the linking hypotheses currently used to connect mechanism properties to performance.^{14,16–22} As with color-discrimination mechanisms, much color appearance data may be accounted for by postulating three postreceptoral mechanisms. But the detailed properties of these mechanisms differ from those of the color-discrimination mechanisms.

In general, excitation of single cone-opponent mechanisms along cardinal directions does not give rise to the perceptions expected of isolated color-opponent mechanisms. Thus, the modulation of L–M around white produces a red/magenta to cyan color variation, whereas modulation of S around white produces a purple to yellow/green modulation.^{31,32}

In specific tests of the relationship between color-discrimination and color-appearance mechanisms, unique hues, which correspond to the null vector of color-appearance mechanisms (see subsection “Spectral Properties of Color-Opponent Mechanisms” in Sec. 11.5), have been compared with the null vectors of color-discrimination mechanisms. Webster et al.²⁹⁸ plotted pairs of unique colors of different saturation in DKL space. They found that only the unique red pair aligned with one of the cone-opponent axes, the L–M axis. The blue and yellow pair aligned roughly along an intermediate axis, but the green pair was not colinear with the red pair. Thus, unique blue and yellow, in particular,

must reflect the joint activity of L–M and S–(L+M). Chichilnisky and Wandell²⁸⁹ and Wuerger, Atkinson, and Cropper²² found that their four unique hue planes did not correspond to the null planes of color-discrimination mechanisms. Using a hue-naming technique, De Valois et al.²¹ found that hue names are also not easily related to modulations of color-discrimination mechanisms (see Fig. 31).

The inconsistencies between color-discrimination and color-appearance mechanisms strongly suggest that the two are fundamentally different. But can they be considered to be mechanisms at different levels of a common serial processing stream or are they even more distinct? And if they do act serially, how might they be related? The simplified three-stage zone model described in the next section suggests that they could be parts of a common serial process and, at least for R/G, very simply related.

Three-Stage Zone Models

As noted in the Introduction (see section “The Mechanistic Approach” in Sec. 11.2), linear three-stage Müller zone models,²³ in which the second stage is roughly cone-opponent and the third stage color-opponent, have been proposed by Judd²⁴ and more recently by Guth²⁹ and De Valois, and De Valois.¹⁸

The De Valois and De Valois¹⁸ three-stage zone model is an interesting example based on physiological and anatomical assumptions about the relative numerosity of cone inputs to center-surround-opponent neurons. In their indiscriminate-surround model, neurons are assumed to have single cone inputs to their centers and mixed cone inputs to their surrounds (in the ratio of 10L:5M:S). Consequently, both the L–M and M–L cone-opponent stages have –S inputs, as a result of which the third color-opponent stage with an S input to R/G is arguably not needed (see Guth⁴⁷¹ and De Valois and De Valois⁴⁷²). [The dangers of assuming that cone weights can be simply related to relative cone numerosity were discussed before in the context of luminous efficiency (see subsection “Luminance” in Sec. 11.5).] The spectral sensitivities of the De Valois and De Valois¹⁸ third-stage color-opponent mechanisms are shown in Fig. 32, as dashed black lines. Despite their physiologically based approach, the De Valois and De Valois¹⁸ zone model inevitably derives from earlier psychophysical models, partly because there are only a limited number of ways in which the cone fundamentals can be linearly combined to produce plausible cone-opponent and color-opponent spectral sensitivities. As well as three-stage models, there are also many examples of two-stage models, in which the second stage is designed to account either for color-discrimination or for color-appearance data in isolation.^{8–13,15}

As an exercise, we next provide an illustrative example of a linear three-stage zone model based on a few very simple assumptions about the signal transformations at each stage. Our goal was to see if we could derive plausible estimates of the zero crossings of opponent-colors theory without resorting to speculative physiological models or psychophysical data fitting.

First zone At the first stage, we assume the Stockman and Sharpe²⁸ cone fundamentals: $\bar{l}(\lambda)$, $\bar{m}(\lambda)$, and $\bar{s}(\lambda)$, the spectral sensitivities of which are labelled Stage 1 in Fig. 4. These functions are normalized to unity peak.

Second zone At the second stage, we assume classical cone opponency. In the case of L–M (and M–L), we assign equal cone weights, thus yielding $\bar{l}(\lambda) - \bar{m}(\lambda)$, and its opposite-signed pair $\bar{m}(\lambda) - \bar{l}(\lambda)$. This is consistent with the evidence from psychophysical detection experiments for equal L- and M-cone contrast weights into L–M (see subsection “Sensitivity to Different Directions of Color Space” in Sec. 11.5). Note that the zero crossing of this cone-opponent L–M mechanism is 549 nm, which is far from unique yellow. For the zero crossing to be near the 580 nm—the unique yellow assumed at the next stage—the relative M:L cone weight would have to be increased from 1 to 1.55 (see, e.g., Fig. 7.4 of Ref. 15).

In the case of S–(L+M) and (L+M)–S, we assign half the weight to M as to L (in accordance with many other models) and then scale both by 0.69 to give equal weights (in terms of peak spectral sensitivity) to S and L+0.5M, thus yielding $\bar{s}(\lambda) - 0.69[\bar{l}(\lambda) + 0.5\bar{m}(\lambda)]$ and its opposite signed pair $0.69[\bar{l}(\lambda) + 0.5\bar{m}(\lambda)] - \bar{s}(\lambda)$. The zero crossing of this mechanism is 486 nm, which is closer to unique blue (a zero crossing of R/G) than to unique green (a zero crossing of Y/B). For a linear Y/B to have a zero crossing near unique green, the B pole must have a contribution from M or L. The cone weights

into the cone-opponent mechanisms assumed at the second level are consistent with psychophysical measurements of color detection and discrimination estimates (e.g., Table 18.1 of Ref. 89).

The spectral sensitivities of the cone-opponent pairs are labelled Stage 2 in Fig. 4. (The L–M sensitivities have also been scaled by 2.55 in accordance with the proposals for the next stage.)

Third zone At the third stage, we sum the outputs of the second-stage mechanisms (or equivalently oppose ones of different signs), to give four color-opponent mechanisms (see also Refs. 29 and 18). The L–M cone-opponent input to the third stage is weighted by 2.55, so that R/G is zero for an equal-quantum white light, thus:

Red [L–M summed with S–(L+M)]

$$2.55[\bar{l}(\lambda) - \bar{m}(\lambda)] + (\bar{s}(\lambda) - 0.69[\bar{l}(\lambda) + 0.5\bar{m}(\lambda)]) = 1.86\bar{l}(\lambda) - 2.90\bar{m}(\lambda) + \bar{s}(\lambda)$$

(bipolar) or $|1.86\bar{l}(\lambda) - 2.90\bar{m}(\lambda) + \bar{s}(\lambda)|$ (unipolar)

Green [M–L and (L+M)–S summed]

$$2.55[\bar{m}(\lambda) - \bar{l}(\lambda)] + (0.69[\bar{l}(\lambda) + 0.5\bar{m}(\lambda)] - \bar{s}(\lambda)) = -1.86\bar{l}(\lambda) + 2.90\bar{m}(\lambda) - \bar{s}(\lambda)$$

(bipolar) or $|-1.86\bar{l}(\lambda) + 2.90\bar{m}(\lambda) - \bar{s}(\lambda)|$ (unipolar)

R/G is required to be zero for an equal-quantal white, because otherwise, according to opponent-colors theory, such whites would appear colored (see Ref. 283). Initially, we used the same weights for the cone-opponent inputs to B/Y and Y/B, thus:

Blue [M–L summed with S–(L+M)]

$$2.55[\bar{m}(\lambda) - \bar{l}(\lambda)] + (\bar{s}(\lambda) - 0.69[\bar{l}(\lambda) + 0.5\bar{m}(\lambda)]) = -3.24\bar{l}(\lambda) + 2.21\bar{m}(\lambda) + \bar{s}(\lambda)$$

(bipolar) or $|-3.24\bar{l}(\lambda) + 2.21\bar{m}(\lambda) + \bar{s}(\lambda)|$ (unipolar)

Yellow [L–M summed with (L+M)–S]

$$2.55[\bar{l}(\lambda) - \bar{m}(\lambda)] + (0.69[\bar{l}(\lambda) + 0.5\bar{m}(\lambda)] - \bar{s}(\lambda)) = 3.24\bar{l}(\lambda) - 2.21\bar{m}(\lambda) - \bar{s}(\lambda)$$

(bipolar) or $|3.24\bar{l}(\lambda) - 2.21\bar{m}(\lambda) - \bar{s}(\lambda)|$ (unipolar)

The spectral sensitivities of the color-opponent mechanisms are labelled Stage 3 in Fig. 4. The R/G mechanism yields reasonable estimates of unique blue (477 nm) and unique yellow (580 nm), and the B/Y mechanism yields a reasonable estimate of unique green (504 nm).

One potential problem, however, is that the opposing poles of B/Y and Y/B are unbalanced (as they also are in the De Valois and De Valois model), so that B/Y and Y/B will produce a nonzero response to an equal-quantum white. Given that $Y > B$, the white field would be expected to appear yellowish. This imbalance can be corrected by decreasing the relative contributions of Y perhaps after the half-wave rectification of B/Y into unipolar B and Y mechanisms. Alternatively, it can be corrected by decreasing the weights of the L- and M-cone inputs into B/Y and into Y/B. For this illustrative example, we choose the latter correction and accordingly scale the weights of the L- and M-inputs into Y/B by 0.34 to give a zero response to white, so that:

Blue [M–L summed with S–(L+M)] becomes

$$-1.10\bar{l}(\lambda) + 0.75\bar{m}(\lambda) + \bar{s}(\lambda) \text{ (bipolar) or } |-1.10\bar{l}(\lambda) + 0.75\bar{m}(\lambda) + \bar{s}(\lambda)| \text{ (unipolar)}$$

Yellow [L–M summed with (L+M)–S] becomes

$$1.10\bar{l}(\lambda) - 0.75\bar{m}(\lambda) - \bar{s}(\lambda) \text{ (bipolar) or } |1.10\bar{l}(\lambda) - 0.75\bar{m}(\lambda) - \bar{s}(\lambda)| \text{ (unipolar)}$$

The spectral sensitivities of the Y/B and B/Y color-opponent pair are shown in the lower right panel of Fig. 4 as dashed lines. The correction shifts unique green to 510 nm. Such post hoc adjustments are inevitably somewhat speculative, however. Moreover, given that B/Y is clearly nonlinear (see subsection “Linearity of Color-Opponent Mechanisms” in Sec. 11.5), any mechanism that adjusts the balance of B/Y is also likely to be nonlinear.

Nonetheless, this example demonstrates that plausible estimates of the zero crossings of third-stage color opponency can be generated by making a few very simple assumptions about the transformations at the cone-opponent and color-opponent stages. But, how close are the R/G and B/Y spectral sensitivities to color-appearance data at wavelengths removed from the zero crossings? Figure 32 compares the R/G and B/Y spectral sensitivities with color valence data. Figure 32*a* shows as red circles R/G valence data from Ingling et al.³¹⁴ scaled to best (least-squares) fit the R/G spectral sensitivity. As can be seen, the agreement between R/G and the valence data is remarkably good. Figure 32*b*, shows B/Y valence data for three subjects (AW, JF, and LK, dark blue triangles, purple inverted triangles, and light blue squares, respectively) from Werner and Wooten²⁸⁶ each scaled to best (least-squares) fit B/Y. Consistent with their own failed attempts to find linear combinations of cone fundamentals to describe their B/Y valence data,²⁸⁶ our B/Y spectral sensitivity agrees only approximately with the data for AW and LK, and poorly with the data for JF.

Also shown in Fig. 32 are the spectral sensitivities of R/G and B/Y color-opponent mechanisms proposed by De Valois and De Valois¹⁸ scaled to best fit our functions. Their R/G and B/Y functions agree poorly with our functions and with the valence data. In contrast, our R/G spectral sensitivity agrees well with the proposal by Ingling et al.,³¹⁴ which was based on the valence data shown in Fig. 32*a* (see also Ref. 12).

In summary, simple assumptions about signal transformations at the second and third stages yield a reasonable estimate of the spectral sensitivity of the R/G color-opponent mechanism. Perhaps, then, the R/G mechanism does reflect a cortical summing of the chromatic responses of double-duty L–M center-surround neurons and S–(L+M) neurons as has been suggested before.^{18,29} In contrast, the B/Y mechanism cannot be accounted for so simply, and may reflect much more complex nonlinear transformations (see subsection “Linearity of Color-Opponent Mechanisms” in Sec. 11.5). These differences between R/G and B/Y could be linked to the idea that the R/G mechanism represents the opportunistic use of an existing or perhaps slightly modified neural process following the relatively recent emergence of separate L- and M-cone photopigment genes,⁴⁴⁷ whereas the B/Y mechanism is some more ancient neural process.⁴⁴⁸

Figure 39 shows a speculative “wiring” diagram of the three-stage Müller zone model.

Final Remarks

Although the three-stage zone model provides a way to integrate the separate two-stage models required for discrimination and for appearance, it is important to note several limitations.

First, as formulated, it does not account for the nonlinearities in unique hue and hue-scaling judgments described in “Linearity of Color-Opponent Mechanisms” in Sec. 11.5. It is possible that explicitly incorporating nonlinearities in the contrast-response functions of each stage could remedy this, but this remains an open question.

Second, the model as outlined just above is for a single state of adaptation. We know that adaptive processes act at both the first and second stages and, with the three-stage models, the effects of these would be expected to propagate to the third stage and thus affect both discrimination and appearance in a common fashion. Although some data, reviewed in “Color Appearance and Chromatic Detection and Discrimination” in Sec. 11.5, suggest that this is the case, further investigation on this point is required. In addition, test on pedestal data for crossed conditions suggest that somewhere in the processing chain there may be additional cross mechanism adaptive effects. Such effects have also been suggested to play a role in appearance data.³⁸¹

Third, the model as formulated does not explicitly account for effects of the temporal and spatial structure of the stimuli. These components could certainly be added, but whether this can be done in a parsimonious fashion that does justice to the empirical phenomena is not clear. If our mechanistic understanding of color vision is to be applied to natural viewing, then its extension to handle the

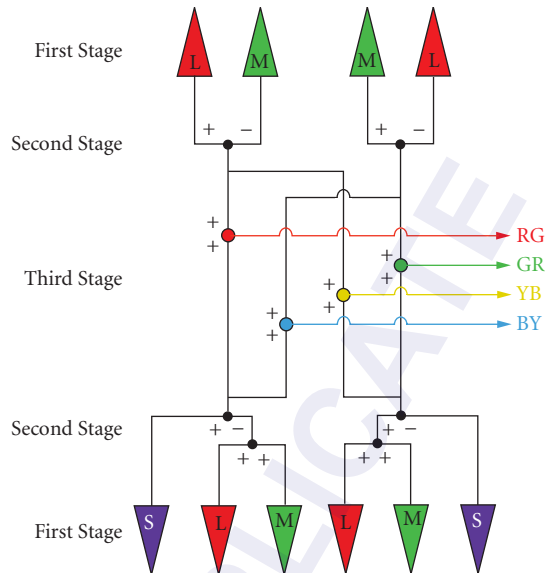


FIGURE 39 Three-stage Müller zone model. *First stage:* L-, M-, and S-cone photoreceptors (top and bottom). *Second stage:* L–M and M–L cone opponency (top) and S–(L+M) and (L+M)–S cone opponency (bottom). *Third stage:* Color opponency (center) is achieved by summing the various cone-opponent second-stage outputs.

complexity of natural retinal stimulation must be a high priority. At the same time, this is a daunting problem because of the explosion in stimulus parameters and the difficulties in controlling them adequately that occur when arbitrary spatiotemporal patterns are considered for both test and field.

Finally, we do not know whether models of this sort can provide a unified account of data across a wider range of tasks than simple threshold and appearance judgments.

Despite these unknowns and limitations, the type of three-stage model described here provides a framework for moving forward. It remains to be seen whether a model of this type will eventually provide a unified account of a wide range data, or whether what will be required, as was the case with Stiles' π -mechanism model which preceded it, is a reconceptualization of the nature of the psychophysical mechanisms and/or the linking hypotheses that connect them to behavioral data.

11.7 ACKNOWLEDGMENTS

The first author acknowledges a significant intellectual debt to Rhea Eskew not only for his writings on this subject, but also for the many discussions over many years about color detection and discrimination. He also acknowledges the support of Sabine Apitz, his wife, without whom this chapter could not have been written; and the financial support of the Wellcome Trust, the BBSRC, and Fight for Sight. The second author acknowledges the support of NIH R01 EY10016. The authors are grateful to Bruce Henning for helpful discussions, proofreading and help with Fig. 20, and Caterina Ripamonti for help with Fig. 16. They also thank Rhea Eskew, Bruce Henning, Ken Knoblauch, and Caterina Ripamonti for helpful comments on the manuscript.

11.8 REFERENCES

1. K. R. Gegenfurtner and D. C. Kiper, "Color Vision," *Annual Review of Neuroscience* **26**:181–206(2003).
2. P. Lennie and J. A. Movshon, "Coding of Color and Form in the Geniculostriate Visual Pathway," *Journal of the Optical Society of America A* **10**:2013–2033 (2005).
3. S. G. Solomon and P. Lennie, "The Machinery of Colour Vision," *Nature Reviews Neuroscience* **8**:276–286 (2007).
4. T. Young, "On the Theory of Light and Colours," *Philosophical Transactions of the Royal Society of London* **92**:20–71 (1802).
5. H. von Helmholtz, *Handbuch der Physiologischen Optik*, Hamburg and Leipzig, Voss, 1867.
6. E. Hering, *Zur Lehre vom Lichtsinne. Sechs Mittheilungen an die Kaiserliche Akademie der Wissenschaften in Wien*, Carl Gerold's Sohn, Wien, 1878.
7. E. Hering, *Grundzüge der Lehre vom Lichtsinn*, Springer, Berlin, 1920.
8. L. M. Hurvich and D. Jameson, "An Opponent-Process Theory of Color Vision," *Psychological Review* **64**:384–404 (1957).
9. R. M. Boynton, "Theory of Color Vision," *Journal of the Optical Society of America* **50**:929–944 (1960).
10. P. L. Walraven, "On the Mechanisms of Colour Vision," Institute for Perception RVO-TNO, The Netherlands, 1962.
11. S. L. Guth and H. R. Lodge, "Heterochromatic Additivity, Foveal Spectral Sensitivity, and a New Color Model," *Journal of the Optical Society of America* **63**:450–462 (1973).
12. C. R. Ingling, Jr. and H. B.-P. Tsou, "Orthogonal Combination of the Three Visual Visual Channels," *Vision Research* **17**:1075–1082 (1977).
13. S. L. Guth, R. W. Massof, and T. Benzschawel, "Vector Model for Normal and Dichromatic color vision," *Journal of the Optical Society of America* **70**:197–212 (1980).
14. J. Krauskopf, D. R. Williams, and D. W. Heeley, "Cardinal Directions of Color Space," *Vision Research* **22**:1123–1131 (1982).
15. R. M. Boynton, *Human Color Vision*, Holt, Rinehart and Winston, New York, 1979.
16. S. A. Burns, A. E. Elsner, J. Pokorny, and V. C. Smith, "The Abney Effect: Chromaticity Coordinates of Unique and Other Constant Hues," *Vision Research* **24**:479–489(1984).
17. M. A. Webster and J. D. Mollon, "Contrast Adaptation Dissociates Different Measures of Luminous Efficiency," *Journal of the Optical Society of America A* **10**:1332–1340 (1993).
18. R. L. De Valois and K. K. De Valois, "A Multi-Stage Color Model," *Vision Research* **33**:1053–1065 (1993).
19. I. Abramov and J. Gordon, "Color Appearance: On Seeing Red-or Yellow, or Green, or Blue," *Annual Review of Psychology* **45**:451–485 (1994).
20. R. T. Eskew, Jr. and P. M. Kortick, "Hue Equilibria Compared with Chromatic Detection in 3D Cone Contrast Space," *Investigative Ophthalmology and Visual Science (supplement)* **35**:1555 (1994).
21. R. L. De Valois, K. K. De Valois, E. Switkes, and L. Mahon, "Hue Scaling of Isoluminant and Cone-Specific Lights," *Vision Research* **37**:885–897 (1997).
22. S. M. Wuerger, P. Atkinson, and S. Cropper, "The Cone Inputs to the Unique-Hue Mechanisms," *Vision Research* **45**:25–26 (2005).
23. G. E. Müller, "Über die Farbenempfindungen," *Zeitschrift für Psychologie und Physiologie der Sinnesorgane, Ergänzungsband* **17**:1–430 (1930).
24. D. B. Judd, "Response Functions for Types of Vision According to the Müller Theory, Research Paper RP1946," *Journal of Research of the National Bureau of Standards* **42**:356–371 (1949).
25. D. B. Judd, "Fundamental Studies of Color Vision from 1860 to 1960," *Proceedings of the National Academy of Science of the United States of America* **55**:1313–1330 (1966).
26. G. Svaetichin and E. F. MacNichol, Jr., "Retinal Mechanisms for Chromatic and Achromatic Vision," *Annals of the New York Academy of Sciences* **74**:385–404 (1959).
27. R. L. De Valois, I. Abramov, and G. H. Jacobs, "Analysis of Response Patterns of LGN Cells," *Journal of the Optical Society of America* **56**:966–977 (1966).

28. A. Stockman and L. T. Sharpe, "Spectral Sensitivities of the Middle- and Long-Wavelength Sensitive Cones Derived from Measurements in Observers of Known Genotype," *Vision Research* **40**:1711–1737 (2000).
29. S. L. Guth, "A Model for Color and Light Adaptation," *Journal of the Optical Society of America A* **8**:976–993 (1991).
30. H. Hofer, J. Carroll, J. Neitz, M. Neitz, and D. R. Williams, "Organization of the Human Trichromatic Cone Mosaic," *Journal of Neuroscience* **25**:9669–9679 (2005).
31. P. Flanagan, P. Cavanagh, and O. E. Favreau, "Independent Orientation-Selective Mechanisms for the Cardinal Directions of Colour Space," *Vision Research* **30**:769–778 (1990).
32. M. A. Webster and J. D. Mollon, "The Influence of Contrast Adaptation on Color Appearance," *Vision Research* **34**:1993–2020 (1994).
33. W. S. Stiles, *Mechanisms of Colour Vision*, Academic Press, London, 1978.
34. R. T. Eskew, Jr, "Chromatic Detection and Discrimination," in *The Senses: A Comprehensive Reference, Volume 2: Vision II*, T. D. Albright and R. H. Masland, eds., Academic Press Inc., San Diego, 2008, pp. 101–117.
35. D. A. Baylor, B. J. Nunn, and J. L. Schnapf, "The Photocurrent, Noise and Spectral Sensitivity of Rods of the Monkey Macaca Fascicularis," *Journal of Physiology* **357**:575–607 (1984).
36. A. M. Derrington, J. Krauskopf, and P. Lennie, "Chromatic Mechanisms in Lateral Geniculate Nucleus of Macaque," *Journal of Physiology* **357**:241–265 (1984).
37. A. M. Derrington and P. Lennie, "Spatial and Temporal Contrast Sensitivities of Neurones in the Lateral Geniculate Nucleus of Macaque," *Journal of Physiology* **357**:219–240 (1984).
38. N. Graham and D. C. Hood, "Modeling the Dynamics of Adaptation: the Merging of Two Traditions," *Vision Research* **32**:1373–1393 (1992).
39. D. C. Hood, "Lower-Level Visual Processing and Models of Light Adaptation," *Annual Review of Psychology* **49**:503–535 (1998).
40. W. S. Stiles, "The Directional Sensitivity of the Retina and the Spectral Sensitivity of the Rods and Cones," *Proceedings of the Royal Society of London. Series B: Biological Sciences* **B127**:64–105 (1939).
41. W. S. Stiles, "Incremental Thresholds and the Mechanisms of Colour Vision," *Documenta Ophthalmologica* **3**:138–163 (1949).
42. W. S. Stiles, "Further Studies of Visual Mechanisms by the Two-Colour Threshold Technique," *Coloquio sobre problemas opticos de la vision* **1**:65–103 (1953).
43. D. M. Green and J. A. Swets, *Signal Detection Theory and Psychophysics*, John Wiley & Sons, New York, 1966.
44. H. B. Barlow and W. R. Levick, "Threshold Setting by the Surround of Cat Retinal Ganglion Cells," *Journal of Physiology* **259**:737–757 (1976).
45. R. Shapley and C. Enroth-Cugell, "Visual Adaptation and Retinal Gain Controls," in *Progress in Retinal Research*, N. Osborne and G. Chader, eds., Pergamon Press, New York, 1984, pp. 263–346.
46. R. M. Boynton and D. N. Whitten, "Visual Adaptation in Monkey Cones: Recordings of Late Receptor Potentials," *Science* **170**:1423–1426 (1970).
47. J. M. Valetton and D. van Norren, "Light Adaptation of Primate Cones: An Analysis Based on Extracellular Data," *Vision Research* **23**:1539–1547 (1983).
48. D. A. Burkhardt, "Light Adaptation and Photopigment Bleaching in Cone Photoreceptors *in situ* in the Retina of the Turtle," *Journal of Neuroscience* **14**:1091–1105 (1994).
49. V. Y. Arshavsky, T. D. Lamb, and E. N. Pugh, Jr, "G Proteins and Phototransduction," *Annual Review of Physiology* **64**:153–187 (2002).
50. R. D. Hamer, S. C. Nicholas, D. Tranchina, T. D. Lamb, and J. L. P. Jarvinen, "Toward a Unified Model of Vertebrate Rod Phototransduction," *Visual Neuroscience* **22**:417–436 (2005).
51. E. N. Pugh, Jr, S. Nikonov, and T. D. Lamb, "Molecular Mechanisms of Vertebrate Photoreceptor Light Adaptation," *Current Opinion in Neurobiology* **9**:410–418 (1999).
52. A. Stockman, M. Langendörfer, H. E. Smithson, and L. T. Sharpe, "Human Cone Light Adaptation: From Behavioral Measurements to Molecular Mechanisms," *Journal of Vision* **6**:1194–1213 (2006).
53. G. Wyszecki and W. S. Stiles, *Color Science: Concepts and Methods, Quantitative Data and Formulae*, 2nd ed., Wiley, New York, 1982.

54. C. Sigel and E. N. Pugh, Jr., "Stiles's π_5 Color Mechanism: Tests of Field Displacements and Field Additivity Properties," *Journal of the Optical Society of America* **70**:71–81 (1980).
55. B. A. Wandell and E. N. Pugh, Jr., "Detection of Long-Duration Incremental Flashes by a Chromatically Coded Pathway," *Vision Research* **20**:625–635 (1980).
56. E. H. Adelson, "Saturation and Adaptation in the Rod System," *Vision Research* **22**:1299–1312(1982).
57. W. S. Geisler, "Effects of Bleaching and Backgrounds on the Flash Response of the Cone System," *Journal of Physiology* **312**:413–434 (1981).
58. D. C. Hood and M. A. Finkelstein, "Sensitivity to Light," in *Handbook of Perception and Human Performance*, K. Boff, L. Kaufman and J. Thomas, eds., Wiley, New York, 1986, pp. 5-1–5-66.
59. M. M. Hayhoe, N. I. Benimof, and D. C. Hood, "The Time Course of Multiplicative and Subtractive Adaptation Processes," *Vision Research* **27**:1981–1996 (1987).
60. D. H. Kelly, "Effects of Sharp Edges in a Flickering Field," *Journal of the Optical Society of America* **49**:730–732 (1959).
61. C. F. Stromeyer, III, G. R. Cole, and R. E. Kronauer, "Second-Site Adaptation in the Red-Green Chromatic Pathways," *Vision Research* **25**:219–237 (1985).
62. A. Chaparro, C. F. Stromeyer, III, G. Chen, and R. E. Kronauer, "Human Cones Appear to Adapt at Low Light Levels: Measurements on the Red-Green Detection Mechanism," *Vision Research* **35**:3103–3118 (1995).
63. C. F. Stromeyer, III, G. R. Cole, and R. E. Kronauer, "Chromatic Suppression of Cone Inputs to the Luminance Flicker Mechanisms," *Vision Research* **27**:1113–1137 (1987).
64. A. Eisner and D. I. A. MacLeod, "Flicker Photometric Study of Chromatic Adaptation: Selective Suppression of Cone Inputs by Colored Backgrounds," *Journal of the Optical Society of America* **71**:705–718 (1981).
65. S. J. Ahn and D. I. A. MacLeod, "Link-Specific Adaptation in the Luminance and Chromatic Channels," *Vision Research* **33**:2271–2286 (1991).
66. Y. W. Lee, *Statistical Theory of Communication*, John Wiley & Sons, Inc., New York, 1968.
67. H. L. van Trees, *Detection, Estimation, and Modulation Theory*, Part I, Wiley, New York, 1968.
68. J. Nachmias, "Signal Detection Theory and its Applications to Problems in Vision," in *Visual Psychophysics, Handbook of Sensory Physiology*, Vol. VII/4, D. Jameson and L. H. Hurvich, eds., Springer-Verlag, Berlin, 1972, pp. 56–77.
69. C. H. Coombs, R. M. Dawes, and A. Tversky, *Mathematical Psychology*, Prentice-Hall, Englewood Cliff, New Jersey, 1970.
70. R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification and Scene Analysis*, 2nd ed., John Wiley & Sons, New York, 2001.
71. H. B. Barlow, "Retinal and Central Factors in Human Vision Limited by Noise," in *Vertebrate Photoreception*, B. B. P. Fatt, ed., Academic Press, London, 1977, pp. 337–358.
72. A. B. Watson, H. B. Barlow, and J. G. Robson, "What Does the Eye See Best?," *Nature* **31**:419–422 (1983).
73. W. S. Geisler, "Physical Limits of Acuity and Hyperacuity," *Journal of the Optical Society of America A* **1**:775–782 (1984).
74. W. S. Geisler, "Sequential Ideal-Observer Analysis of Visual Discriminations," *Psychological Review* **96**:267–314 (1989).
75. D. H. Krantz, "Color Measurement and Color Theory: I. Representation Theorem for Grassmann Structures," *Journal of Mathematical Psychology* **12**:283–303 (1975).
76. D. H. Krantz, "Color Measurement and Color Theory: II Opponent-Colors Theory," *Journal of Mathematical Psychology* **12**:304–327 (1975).
77. D. L. MacAdam, "Chromatic Adaptation," *Journal of the Optical Society* **46**:500–513 (1956).
78. J. M. Hillis and D. H. Brainard, "Do Common Mechanisms of Adaptation Mediate Color Discrimination and Appearance? Uniform Backgrounds," *Journal of the Optical Society America A* **22**:2090–2106 (2005).
79. M. A. Webster and J. D. Mollon, "Changes in Colour Appearance Following Post-Receptor Adaptation," *Nature* **349**:235–238 (1991).
80. J. P. Moxley and S. L. Guth, "Hue Shifts Caused by Post-Receptor Adaptation," *Investigative Ophthalmology and Visual Science (supplement)* **20**:206 (1981).

81. S. L. Guth and J. P. Moxley, "Hue Shifts Following Differential Postreceptor Achromatic Adaptation," *Journal of the Optical Society of America* **72**:301–303 (1982).
82. J. M. Hillis and D. H. Brainard, "Do Common Mechanisms of Adaptation Mediate Color Discrimination and Appearance? Contrast Adaptation," *Journal of the Optical Society of America A* **24**:2122–2133 (2007).
83. C. Noorlander, M. J. G. Heuts, and J. J. Koenderink, "Sensitivity to Spatiotemporal Combined Luminance and Chromaticity Contrast," *Journal of the Optical Society of America* **71**:453–459 (1981).
84. D. H. Kelly, "Flicker," in *Visual Psychophysics, Handbook of Sensory Physiology*, Vol. VII/4, D. Jameson and L. H. Hurvich, eds., Springer-Verlag, Berlin, 1972, pp. 273–302.
85. R. Luther, "Aus dem Gebiet der Farbreizmetrik," *Zeitschrift für technische Physik* **8**:540–558 (1927).
86. D. I. A. MacLeod and R. M. Boynton, "Chromaticity Diagram Showing Cone Excitation by Stimuli of Equal Luminance," *Journal of the Optical Society of America* **69**:1183–1186 (1979).
87. D. H. Brainard, "Cone Contrast and Opponent Modulation Color Spaces," in *Human Color Vision*, P. K. Kaiser, and R. M. Boynton, eds., Optical Society of America, Washington, D.C., 1996, pp. 563–579.
88. K. Knoblauch, "Dual Bases in Dichromatic Color Space," in *Colour Vision Deficiencies XII*, B. Drum, ed., Kluwer Academic Publishers, Dordrecht, 1995, pp. 165–176.
89. R. T. Eskew, Jr., J. S. McLellan, and F. Giulianini, "Chromatic Detection and Discrimination," in *Color Vision: From Genes to Perception*, K. Gegenfurtner and L. T. Sharpe, eds., Cambridge University Press, Cambridge, 1999, pp. 345–368.
90. A. Stockman and L. T. Sharpe, "Cone Spectral Sensitivities and Color Matching," in *Color Vision: From Genes to Perception*, K. Gegenfurtner and L. T. Sharpe, eds., Cambridge University Press, Cambridge, 1999, pp. 53–87.
91. B. Drum, "Short-Wavelength Cones Contribute to Achromatic Sensitivity," *Vision Research* **23**:1433–1439 (1983).
92. A. Stockman, D. I. A. MacLeod, and D. D. DePriest, "An Inverted S-cone Input to the Luminance Channel: Evidence for Two Processes in S-cone Flicker Detection," *Investigative Ophthalmology and Visual Science (supplement)* **28**:92 (1987).
93. J. Lee and C. F. Stromeyer, III, "Contribution of Human Short-Wave Cones to Luminance and Motion Detection," *Journal of Physiology* **413**:563–593 (1989).
94. A. Stockman, D. I. A. MacLeod, and D. D. DePriest, "The Temporal Properties of the Human Short-Wave Photoreceptors and Their Associated Pathways," *Vision Research* **31**:189–208 (1991).
95. W. S. Stiles, "Separation of the 'Blue' and 'Green' Mechanisms of Foveal Vision by Measurements of Increment Thresholds," *Proceedings of the Royal Society of London. Series B: Biological Sciences* **B 133**:418–434 (1946).
96. W. S. Stiles, "The Physical Interpretation of the Spectral Sensitivity Curve of the Eye," in *Transactions of the Optical Convention of the Worshipful Company of Spectacle Makers*, Spectacle Maker's Company, London, 1948, pp. 97–107.
97. W. S. Stiles, "Color Vision: the Approach Through Increment Threshold Sensitivity," *Proceedings of the National Academy of Science of the United States of America* **45**:100–114 (1959).
98. W. S. Stiles, "Foveal Threshold Sensitivity on Fields of Different Colors," *Science* **145**:1016–1018 (1964).
99. J. M. Enoch, "The Two-Color Threshold Technique of Stiles and Derived Component Color Mechanisms," in *Handbook of Sensory Physiology*, Vol. VII/4, D. Jameson and L. H. Hurvich, eds., Springer-Verlag, Berlin, 1972 pp. 537–567.
100. R. M. Boynton, M. Ikeda, and W. S. Stiles, "Interactions Among Chromatic Mechanisms as Inferred from Positive and Negative Increment Thresholds," *Vision Research* **4**:87–117 (1964).
101. P. E. King-Smith, "Visual Detection Analysed in Terms of Luminance and Chromatic Signals," *Nature* **255**:69–70 (1975).
102. P. E. King-Smith and D. Carden, "Luminance and Opponent-Color Contributions to Visual Detection and Adaptation and to Temporal and Spatial Integration," *Journal of the Optical Society of America* **66**:709–717 (1976).
103. B. A. Wandell and E. N. Pugh, Jr., "A Field Additive Pathway Detects Brief-Duration, Long-Wavelength Incremental Flashes," *Vision Research* **20**:613–624 (1980).

104. W. S. Stiles and B. H. Crawford, "The Liminal Brightness Increment as a Function of Wavelength for Different Conditions of the Foveal and Parafoveal Retina," *Proceedings of the Royal Society of London. Series B: Biological Sciences* **B113**:496–530 (1933).
105. H. G. Sperling and R. S. Harwerth, "Red-Green Cone Interactions in Increment-Threshold Spectral Sensitivity of Primates," *Science* **172**:180–184 (1971).
106. K. Kranda and P. E. King-Smith, "Detection of Colored Stimuli by Independent Linear Systems," *Vision Research* **19**:733–745 (1979).
107. J. E. Thornton and E. N. Pugh, Jr., "Red/Green Color Opponency at Detection Threshold," *Science* **219**:191–193 (1983).
108. M. Kalloniatis and H. G. Sperling, "The Spectral Sensitivity and Adaptation Characteristics of Cone Mechanisms Under White Light Adaptation," *Journal of the Optical Society of America A* **7**:1912–1928 (1990).
109. L. L. Sloan, "The Effect of Intensity of Light, State of Adaptation of the Eye, and Size of Photometric Field on the Visibility Curve," *Psychological Monographs* **38**:1–87 (1928).
110. C. R. Ingling, Jr., "A Tetrachromatic Hypothesis for Human Color Vision," *Vision Research* **9**:1131–1148 (1969).
111. D. J. Calkins, J. E. Thornton, and E. N. Pugh, Jr., "Monochromatism Determined at a Long-Wavelength/Middle-Wavelength Cone-Antagonistic Locus," *Vision Research* **13**:2349–2367 (1992).
112. W. d. W. Abney and E. R. Festing, "Colour Photometry," *Philosophical Transactions of the Royal Society of London* **177**:423–456 (1886).
113. W. d. W. Abney, *Researches in Colour Vision*, Longmans, Green, London, 1913.
114. H. E. Ives, "Studies in the Photometry of Lights of Different Colours. I. Spectral Luminosity Curves Obtained by the Equality of Brightness Photometer and Flicker Photometer under Similar Conditions," *Philosophical Magazine Series 6* **24**:149–188 (1912).
115. W. W. Coblenz and W. B. Emerson, "Relative Sensibility of the Average Eye to Light of Different Color and Some Practical Applications," *Bulletin of the Bureau of Standards* **14**:167–236 (1918).
116. E. P. Hyde, W. E. Forsythe, and F. E. Cady, "The Visibility of Radiation," *Astrophysical Journal*, **48**:65–83 (1918).
117. K. S. Gibson and E. P. T. Tyndall, "Visibility of Radiant Energy," *Scientific Papers of the Bureau of Standards* **19**:131–191 (1923).
118. A. Dresler, "The Non-Additivity of Heterochromatic Brightness," *Transactions of the Illuminating Engineering Society* **18**:141–165 (1953).
119. R. M. Boynton and P. Kaiser, "Vision: The Additivity Law Made to Work for Heterochromatic Photometry with Bipartite Fields," *Science* **161**:366–368 (1968).
120. S. L. Guth, N. V. Donley, and R. T. Marrocco, "On Luminance Additivity and Related Topics," *Vision Research* **9**:537–575 (1969).
121. Y. Le Grand, "Spectral Luminosity," in *Visual Psychophysics, Handbook of Sensory Physiology*, Vol. VII/4, D. Jameson and L. H. Hurvich, eds., Springer-Verlag, Berlin, 1972, pp. 413–433.
122. G. Wagner and R. M. Boynton, "Comparison of Four Methods of Heterochromatic Photometry," *Journal of the Optical Society of America* **62**:1508–1515 (1972).
123. J. Pokorny, V. C. Smith, and M. Lutze, "Heterochromatic Modulation Photometry," *Journal of the Optical Society of America A* **6**:1618–1623 (1989).
124. P. Lennie, J. Pokorny, and V. C. Smith, "Luminance," *Journal of the Optical Society of America A* **10**:1283–1293 (1993).
125. V. C. Smith, and J. Pokorny, "Spectral Sensitivity of the Foveal Cone Photopigments between 400 and 500 nm," *Vision Research* **15**:161–171 (1975).
126. D. B. Judd, "Report of U.S. Secretariat Committee on Colorimetry and Artificial Daylight," in *Proceedings of the Twelfth Session of the CIE, Stockholm*, Bureau Central de la CIE, Paris, 1951, pp. 1–60.
127. J. J. Vos, "Colorimetric and Photometric Properties of a 2-deg Fundamental Observer," *Color Research and Application* **3**:125–128 (1978).
128. L. T. Sharpe, A. Stockman, W. Jagla, and H. Jägle, "A Luminous Efficiency Function, $V'(\lambda)$, for Daylight Adaptation," *Journal of Vision* **5**:948–968 (2005).

129. A. Stockman, H. Jägle, M. Pirzer, and L. T. Sharpe, "The Dependence of Luminous Efficiency on Chromatic Adaptation," *Journal of Vision* **8**(16):1, 1–26 (2008).
130. CIE, *Fundamental Chromaticity Diagram with Physiological Axes—Part 1. Technical Report 170-1*, Central Bureau of the Commission Internationale de l'Éclairage, Vienna, 2007.
131. H. De Vries, "Luminosity Curves of Trichromats," *Nature* **157**:736–737 (1946).
132. H. De Vries, "The Heredity of the Relative Numbers of Red and Green Receptors in the Human Eye," *Genetica* **24**:199–212 (1948).
133. R. A. Crone, "Spectral Sensitivity in Color-Defective Subjects and Heterozygous Carriers," *American Journal of Ophthalmology* **48**:231–238 (1959).
134. W. A. H. Rushton and H. D. Baker, "Red/Green Sensitivity in Normal Vision," *Vision Research* **4**:75–85 (1964).
135. A. Adam, "Foveal Red-Green Ratios of Normals, Colorblinds and Heterozygotes," *Proceedings Tel-Hashomer Hospital: Tel-Aviv* **8**:2–6 (1969).
136. J. J. Vos and P. L. Walraven, "On the Derivation of the Foveal Receptor Primaries," *Vision Research* **11**:799–818 (1971).
137. M. Lutze, N. J. Cox, V. C. Smith, and J. Pokorny, "Genetic Studies of Variation in Rayleigh and Photometric Matches in Normal Trichromats," *Vision Research* **30**:149–162(1990).
138. R. L. P. Vimal, V. C. Smith, J. Pokorny, and S. K. Shevell, "Foveal Cone Thresholds," *Vision Research* **29**:61–78 (1989).
139. J. Kremers, H. P. N. Scholl, H. Knau, T. T. J. M. Berendschot, and L. T. Sharpe, "L/M-Cone Ratios in Human Trichromats Assessed by Psychophysics, Electroretinography and Retinal Densitometry," *Journal of the Optical Society of America A* **17**:517–526 (2000).
140. K. R. Dobkins, A. Thiele, and T. D. Albright, "Comparisons of Red-Green Equiluminance Points in Humans and Macaques: Evidence for Different L:M Cone Ratios between Species," *Journal of the Optical Society of America A* **17**:545–556 (2000).
141. K. L. Gunther, and K. R. Dobkins, "Individual Differences in Chromatic (red/green) Contrast Sensitivity are Constrained by the Relative Numbers of L- versus M-cones in the Eye," *Vision Research* **42**:1367–1378 (2002).
142. A. Stockman, D. I. A. MacLeod, and J. A. Vivien, "Isolation of the Middle- and Long-Wavelength Sensitive Cones in Normal Trichromats," *Journal of the Optical Society of America A* **10**:2471–2490 (1993).
143. J. Carroll, C. McMahon, M. Neitz, and J. Neitz, "Flicker-Photometric Electroretinogram Estimates of L:M Cone Photoreceptor Ratio in Men with Photopigment Spectra Derived from Genetics," *Journal of the Optical Society of America A* **17**:499–509 (2000).
144. J. Carroll, J. Neitz, and M. Neitz, "Estimates of L:M Cone Ratio from ERG Flicker Photometry and Genetics," *Journal of Vision* **2**:531–542 (2002).
145. M. F. Wesner, J. Pokorny, S. K. Shevell, and V. C. Smith, "Foveal Cone Detection Statistics in Color-Normals and Dichromats," *Vision Research* **31**:1021–1037 (1991).
146. D. H. Brainard, A. Roorda, Y. Yamauchi, J. B. Calderone, A. Metha, M. Neitz, J. Neitz, D. R. Williams, and G. H. Jacobs, "Functional Consequences of the Relative Numbers of L and M Cones," *Journal of the Optical Society of America A* **17**:607–614 (2000).
147. J. Albrecht, H. Jägle, D. C. Hood, and L. T. Sharpe, "The Multifocal Electroretinogram (mfERG) and Cone Isolating Stimuli: Variation in L- and M-cone driven signals across the Retina," *Journal of Vision* **2**:543–558 (2002).
148. L. T. Sharpe, E. de Luca, T. Hansen, H. Jägle, and K. Gegenfurtner, "Advantages and Disadvantages of Human Dichromacy," *Journal of Vision* **6**:213–223 (2006).
149. C. R. Ingling, Jr., and E. Martinez-Uriegas, "The Relationship between Spectral Sensitivity and Spatial Sensitivity for the Primate r-g X-Channel," *Vision Research* **23**:1495–1500(1983).
150. C. R. Ingling, Jr. and E. Martinez, "The Spatio-Chromatic Signal of the r-g Channels," in *Colour Vision: Physiology and Psychophysics*, J. D. Mollon and L. T. Sharpe, eds., Academic Press, London, 1983, pp. 433–444.
151. C. R. Ingling, Jr. and E. Martinez-Uriegas, "The Spatiotemporal Properties of the R-G X-cell Channel," *Vision Research* **25**:33–38 (1985).
152. C. R. Ingling, Jr. and H. B.-P. Tsou, "Spectral Sensitivity for Flicker and Acuity Criteria," *Journal of the Optical Society of America A* **5**:1374–1378 (1988).

153. P. Lennie and M. D’Zmura, “Mechanisms of Color Vision,” *CRC Critical Reviews in Neurobiology* **3**:333–400 (1988).
154. P. K. Kaiser, B. B. Lee, P. R. Martin, and A. Valberg, “The Physiological Basis of the Minimally Distinct Border Demonstrated in the Ganglion Cells of the Macaque Retina,” *Journal of Physiology* **422**:153–183 (1990).
155. P. Gouras and E. Zrenner, “Enhancement of Luminance Flicker by Color-Opponent Mechanisms,” *Science* **205**:587–589 (1979).
156. J. L. Brown, L. Phares, and D. E. Fletcher, “Spectral Energy Thresholds for the Resolution of Acuity Targets,” *Journal of the Optical Society of America* **50**:950–960 (1960).
157. J. Pokorny, C. H. Graham, and R. N. Lanson, “Effect of Wavelength on Foveal Grating Acuity,” *Journal of the Optical Society of America* **58**:1410–1414 (1968).
158. C. R. Ingling, S. S. Grigsby, and R. C. Long, “Comparison of Spectral Sensitivity Using Heterochromatic Flicker Photometry and an Acuity Criterion,” *Color Research & Application* **17**:187–196 (1992).
159. M. Ikeda, “Study of Interrelations between Mechanisms at Threshold,” *Journal of the Optical Society of America A* **53**:1305–1313 (1963).
160. S. L. Guth, “Luminance Addition: General Considerations and Some Results at Foveal Threshold,” *Journal of the Optical Society of America* **55**:718–722 (1965).
161. S. L. Guth, “Nonadditivity and Inhibition Among Chromatic Luminances at Threshold,” *Vision Research* **7**:319–328 (1967).
162. S. L. Guth, J. V. Alexander, J. I. Chumbly, C. B. Gillman, and M. M. Patterson, “Factors Affecting Luminance Additivity at Threshold Among Normal and Color-Blind Subjects and Elaborations of a Trichromatic-Opponent Color Theory,” *Vision Research* **8**:913–928 (1968).
163. M. Ikeda, T. Uetsuki, and W. S. Stiles, “Interrelations among Stiles π Mechanisms,” *Journal of the Optical Society of America* **60**:406–415 (1970).
164. J. E. Thornton and E. N. Pugh, Jr., “Relationship of Opponent-Colors Cancellation Measures to Cone Antagonistic Signals Deduced from Increment Threshold Data,” in *Colour Vision: Physiology and Psychophysics*, J. D. Mollon and L. T. Sharpe, eds., Academic Press, London, 1983, pp. 361–373.
165. A. B. Poirson and B. A. Wandell, “The Ellipsoidal Representation of Spectral Sensitivity,” *Vision Research* **30**:647–652 (1990).
166. K. Knoblauch and L. T. Maloney, “Testing the Indeterminacy of Linear Color Mechanisms from Color Discrimination Data,” *Vision Research* **36**:295–306 (1996).
167. A. B. Poirson and B. A. Wandell, “Pattern-Color Separable Pathways Predict Sensitivity to Simple Colored Patterns” *Vision Research* **36**:515–526 (1996).
168. A. B. Poirson B. A. Wandell, D. C. Varner, and D. H. Brainard, “Surface Characterizations of Color Thresholds,” *Journal of the Optical Society of America A* **7**:783–789 (1990).
169. G. R. Cole, C. F. Stromeyer, III, and R. E. Kronauer, “Visual Interactions with Luminance and Chromatic Stimuli,” *Journal of the Optical Society of America A* **7**:128–140 (1990).
170. A. Chaparro, C. F. Stromeyer, III, E. P. Huang, R. E. Kronauer, and R. T. Eskew, Jr, “Colour is What the Eye Sees Best,” *Nature* **361**:348–350 (1993).
171. R. T. Eskew, Jr., C. F. Stromeyer, III, and R. E. Kronauer, “Temporal Properties of the Red-Green Chromatic Mechanism,” *Vision Research* **34**:3127–3137 (1994).
172. M. J. Sankeralli and K. T. Mullen, “Estimation of the L-, M-, and S-cone Weights of the Postreceptoral Detection Mechanisms,” *Journal of the Optical Society of America A* **13**:906–915 (1996).
173. A. Chaparro, C. F. Stromeyer, III, R. E. Kronauer, and R. T. Eskew, Jr., “Separable Red-Green and Luminance Detectors for Small Flashes,” *Vision Research* **34**:751–762 (1994).
174. G. R. Cole, T. Hine, and W. McIlhagga, “Detection Mechanisms in L-, M-, and S-cone Contrast Space,” *Journal of the Optical Society of America A* **10**:38–51 (1993).
175. F. Giulianini and R. T. Eskew, Jr., “Chromatic Masking in the ($\Delta L/L$, $\Delta M/M$) Plane of Cone-Contrast Space Reveals only Two Detection Mechanisms,” *Vision Research* **38**:3913–3926(1998).
176. J. R. Newton and R. T. Eskew, Jr., “Chromatic Detection and Discrimination in the Periphery: A Postreceptoral Loss of Color Sensitivity,” *Visual Neuroscience* **20**:511–521 (2003).
177. R. T. Eskew, Jr., C. F. Stromeyer, III, C. J. Picotte, and R. E. Kronauer, “Detection Uncertainty and the Facilitation of Chromatic Detection by Luminance Contours,” *Journal of the Optical Society of America A* **8**: 394–403 (1991).

178. C. F. Stromeyer, III, J. Lee, and R. T. Eskew, Jr., "Peripheral Chromatic Sensitivity for Flashes: a Post-Receptorial Red-Green Asymmetry," *Vision Research* **32**:1865–1873 (1992).
179. R. T. Eskew, Jr., J. R. Newton, and F. Giulianini, "Chromatic Detection and Discrimination Analyzed by a Bayesian Classifier," *Vision Research* **41**:893–909 (2001).
180. R. M. Boynton, A. L. Nagy, and C. X. Olson, "A Flaw in Equations for Predicting Chromatic Differences," *Color Research and Application* **8**:69–74 (1983).
181. C. F. Stromeyer, III, A. Chaparro, C. Rodriguez, D. Chen, E. Hu, and R. E. Kronauer, "Short-Wave Cone Signal in the Red-Green Detection Mechanism," *Vision Research* **38**:813–826 (1998).
182. C. F. Stromeyer, III and J. Lee, "Adaptational Effects of Short Wave Cone Signals on Red-Green Chromatic Detection," *Vision Research* **28**:931–940 (1988).
183. J. D. Mollon and C. R. Cavonius, "The Chromatic Antagonisms of Opponent-Process Theory are not the Same as Those Revealed in Studies of Detection and Discrimination," in *Colour Deficiencies VIII, Documenta Ophthalmologica Proceedings Series*, 46, G. Verriest, ed., Nijhoff-Junk, Dordrecht, 1987, pp. 473–483.
184. H. De Lange, "Research into the Dynamic Nature of the Human Fovea-Cortex Systems with Intermittent and Modulated Light. II. Phase Shift in Brightness and Delay in Color Perception," *Journal of the Optical Society of America* **48**:784–789 (1958).
185. D. Regan and C. W. Tyler, "Some Dynamic Features of Colour Vision," *Vision Research* **11**:1307–1324 (1971).
186. D. H. Kelly and D. van Norren, "Two-Band Model of Heterochromatic Flicker," *Journal of the Optical Society of America* **67**:1081–1091 (1977).
187. D. J. Tolhurst, "Colour-Coding Properties of Sustained and Transient Channels in Human Vision," *Nature* **266**:266–268 (1977).
188. C. E. Sternheim, C. F. Stromeyer, III, and M. C. K. Khoo, "Visibility of Chromatic Flicker upon Spectrally Mixed Adapting Fields," *Vision Research* **19**:175–183 (1979).
189. V. C. Smith, R. W. Bowen, and J. Pokorny, "Threshold Temporal Integration of Chromatic Stimuli," *Vision Research* **24**:653–660 (1984).
190. A. B. Metha and K. T. Mullen, "Temporal Mechanisms Underlying Flicker Detection and Identification for Red-Green and Achromatic Stimuli," *Journal of the Optical Society of America A* **13**:1969–1980 (1996).
191. A. Stockman, M. R. Williams, and H. E. Smithson, "Flicker-Clicker: Cross Modality Matching Experiments," *Journal of Vision* **4**:86a (2004).
192. G. J. C. van der Horst, C. M. M. de Weert, and M. A. Bouman, "Transfer of Spatial Chromaticity-Contrast at Threshold in the Human Eye," *Journal of the Optical Society of America* **57**:1260–1266 (1967).
193. G. J. C. van der Horst and M. A. Bouman, "Spatio-Temporal Chromaticity Discrimination," *Journal of the Optical Society of America* **59**:1482–1488 (1969).
194. R. Hilz and C. R. Cavonius, "Wavelength Discrimination Measured with Square-Wave Gratings," *Journal of the Optical Society of America* **60**:273–277 (1970).
195. E. M. Granger and J. C. Heurtley, "Visual Chromaticity-Modulation Transfer Function," *Journal of the Optical Society of America* **63**:1173–1174 (1973).
196. C. F. Stromeyer, III and C. E. Sternheim, "Visibility of Red and Green Spatial Patterns upon Spectrally Mixed Adapting Fields," *Vision Research* **21**:397–407 (1981).
197. K. T. Mullen and J. J. Kulikowski, "Wavelength Discrimination at Detection Threshold," *Journal of the Optical Society of America A* **7**:733–742 (1990).
198. N. Sekiguchi, D. R. Williams, and D. H. Brainard, "Efficiency in Detection of Isoluminant and Isochromatic Interference Fringes," *Journal of the Optical Society of America A* **10**:2118–2133 (1993).
199. N. Sekiguchi, D. R. Williams, and D. H. Brainard, "Aberration-Free Measurements of the Visibility of Isoluminant Gratings," *Journal of the Optical Society of America A* **10**:2105–2117 (1993).
200. D. R. Williams, N. Sekiguchi, and D. H. Brainard, "Color, Contrast Sensitivity, and the Cone Mosaic," *Proceedings of the National Academy of Science of the United States of America* **90**:9770–9777 (1993).
201. E. N. Pugh, Jr. and C. Sigel, "Evaluation of the Candidacy of the π -Mechanisms of Stiles for Color-Matching Fundamentals," *Vision Research* **18**:317–330 (1978).
202. O. Estévez, "On the Fundamental Database of Normal and Dichromatic Color Vision," PhD thesis, Amsterdam University, Amsterdam 1979.

203. H. J. A. Dartnall, J. K. Bowmaker, and J. D. Mollon, "Human Visual Pigments: Microspectrophotometric Results from the Eyes of Seven Persons," *Proceedings of the Royal Society of London. Series B: Biological Sciences* **B 220**:115–130 (1983).
204. H. De Vries, "The Luminosity Curve of the Eye as Determined by Measurements with the Flicker Photometer," *Physica* **14**:319–348 (1948).
205. M. Ikeda and M. Urakubo, "Flicker HRTF as Test of Color Vision," *Journal of the Optical Society of America* **58**:27–31 (1968).
206. L. E. Marks and M. H. Bornstein, "Spectral Sensitivity by Constant CFF: Effect of Chromatic Adaptation," *Journal of the Optical Society of America* **63**:220–226 (1973).
207. P. E. King-Smith and J. R. Webb, "The Use of Photopic Saturation in Determining the Fundamental Spectral Sensitivity Curves," *Vision Research* **14**:421–429 (1974).
208. A. Eisner, "Comparison of Flicker-Photometric and Flicker-Threshold Spectral Sensitivities While the Eye is Adapted to Colored Backgrounds," *Journal of the Optical Society of America* **72**:517–518 (1982).
209. W. H. Swanson, "Chromatic Adaptation Alters Spectral Sensitivity at High Temporal Frequencies," *Journal of the Optical Society of America A* **10**:1294–1303 (1993).
210. C. F. Stromeyer, III, A. Chaparro, A. S. Tolia, and R. E. Kronauer, "Colour Adaptation Modifies the Long-Wave Versus Middle-Wave Cone Weights and Temporal Phases in Human Luminance (but not red-green) Mechanism," *Journal of Physiology* **499**:227–254 (1997).
211. W. B. Cushman and J. Z. Levinson, "Phase Shift in Red and Green Counter-Phase Flicker at High Frequencies," *Journal of the Optical Society of America* **73**:1557–1561 (1983).
212. P. L. Walraven and H. J. Leebeek, "Phase Shift of Sinusoidally Alternating Colored Stimuli," *Journal of the Optical Society of America* **54**:78–82 (1964).
213. D. T. Lindsey, J. Pokorny, and V. C. Smith, "Phase-Dependent Sensitivity to Heterochromatic Flicker," *Journal of the Optical Society of America A* **3**:921–927 (1986).
214. W. H. Swanson, J. Pokorny, and V. C. Smith, "Effects of Temporal Frequency on Phase-Dependent Sensitivity to Heterochromatic Flicker," *Journal of the Optical Society of America A* **4**:2266–2273 (1987).
215. V. C. Smith, B. B. Lee, J. Pokorny, P. R. Martin, and A. Valberg, "Responses of Macaque Ganglion Cells to the Relative Phase of Heterochromatically Modulated Lights," *Journal of Physiology* **458**:191–221 (1992).
216. A. Stockman, D. J. Plummer, and E. D. Montag, "Spectrally-Opponent Inputs to the Human Luminance Pathway: Slow +M and -L Cone Inputs Revealed by Intense Long-Wavelength Adaptation," *Journal of Physiology* **566**:61–76 (2005).
217. A. Stockman and D. J. Plummer, "Spectrally-Opponent Inputs to the Human Luminance Pathway: Slow +L and -M Cone Inputs Revealed by Low to Moderate Long-Wavelength Adaptation," *Journal of Physiology* **566**:77–91 (2005).
218. A. Stockman, E. D. Montag, and D. J. Plummer, "Paradoxical Shifts in Human Colour Sensitivity Caused by Constructive and Destructive Interference between Signals from the Same Cone Class," *Visual Neuroscience* **23**:471–478 (2006).
219. A. Stockman, E. D. Montag, and D. I. A. MacLeod, "Large Changes in Phase Delay on Intense Bleaching Backgrounds," *Investigative Ophthalmology and Visual Science (supplement)* **32**:841 (1991).
220. A. Stockman and D. J. Plummer, "The Luminance Channel Can be Opponent?," *Investigative Ophthalmology and Visual Science (supplement)* **35**:1572 (1994).
221. C. F. Stromeyer, III, P. D. Gowdy, A. Chaparro, S. Kladakis, J. D. Willen, and R. E. Kronauer, "Colour Adaptation Modifies the Temporal Properties of the Long- and Middle-Wave Cone Signals in the Human Luminance Mechanism," *Journal of Physiology* **526**:177–194 (2000).
222. A. Stockman, "Multiple Cone Inputs to Luminance," *Investigative Ophthalmology and Visual Science (supplement)* **42**:S320 (2001).
223. A. Stockman and D. J. Plummer, "Long-Wavelength Adaptation Reveals Slow, Spectrally Opponent Inputs to the Human Luminance Pathway," *Journal of Vision* **5**:702–716 (2005).
224. J. D. Mollon and P. G. Polden, "Saturation of a Retinal Cone Mechanism," *Nature* **259**:243–246 (1977).
225. C. F. Stromeyer, III, R. E. Kronauer, and J. C. Madsen, "Response Saturation of Short-Wavelength Cone Pathways Controlled by Color-Opponent Mechanisms," *Vision Research* **19**:1025–1040 (1979).
226. E. N. Pugh, Jr., "The Nature of the π_1 Mechanism of W. S. Stiles," *Journal of Physiology* **257**:713–747 (1976).

227. P. G. Polden and J. D. Mollon, "Reversed Effect of Adapting Stimuli on Visual Sensitivity," *Proceedings of the Royal Society of London. Series B: Biological Sciences* **B 210**:235–272 (1980).
228. E. N. Pugh, Jr. and J. Larimer, "Test of the Identity of the Site of Blue/Yellow Hue Cancellation and the Site of Chromatic Antagonism in the π_1 Pathway," *Vision Research* **19**:779–788 (1980).
229. D. B. Kirk, "The Putative π_4 Mechanism: Failure of Shape Invariance and Field Additivity," *Investigative Ophthalmology and Visual Science (supplement)* **26**:184 (1985).
230. A. Reeves, "Field Additivity of Stiles's Pi-4 Color Mechanism," *Journal of the Optical Society of America A* **4**:525–529 (1987).
231. E. N. Pugh, Jr. and J. D. Mollon, "A Theory of the π_1 and π_3 Color Mechanisms of Stiles," *Vision Research* **20**:293–312 (1979).
232. J. L. Schnapf, B. J. Nunn, M. Meister, and D. A. Baylor, "Visual Transduction in Cones of the Monkey *Macaca Fascicularis*," *Journal of Physiology* **427**:681–713 (1990).
233. D. C. Hood and D. G. Birch, "Human Cone Receptor Activity: The Leading Edge of the A-Wave and Models of Receptor Activity," *Visual Neuroscience* **10**:857–871 (1993).
234. B. B. Lee, D. M. Dacey, V. C. Smith, and J. Pokorny, "Horizontal Cells Reveal Cone Type-Specific Adaptation in Primate Retina," *Proceedings of the National Academy of Sciences of the United States of America* **96**:14611–14616 (1999).
235. B. B. Lee, D. M. Dacey, V. C. Smith, and J. Pokorny, "Dynamics of Sensitivity Regulation in Primate Outer Retina: The Horizontal Cell Network," *Journal of Vision* **3**:513–526 (2003).
236. G. J. Burton, "Evidence for Nonlinear Response Processes in the Human Visual System from Measurements on the Thresholds of Spatial Beat Frequencies," *Vision Research* **13**:1211–1225 (1973).
237. D. I. A. MacLeod, and S. He "Visible Flicker from Invisible Patterns," *Nature* **361**:256–258 (1993).
238. D. I. A. MacLeod, D. R. Williams, and W. Makous, "A Visual Nonlinearity Fed by Single Cones," *Vision Research* **32**:347–363 (1992).
239. O. Estévez and H. Spekreijse, "The 'Silent Substitution' Method in Visual Research," *Vision Research* **22**:681–691 (1982).
240. T. Benzschawel and S. L. Guth, "Post-Receptor Chromatic Mechanisms Revealed by Flickering vs Fused Adaptation," *Vision Research* **22**:69–75 (1982).
241. J. Krauskopf, D. R. Williams, M. B. Mandler, and A. M. Brown, "Higher Order Color Mechanisms," *Vision Research* **26**:23–32 (1986).
242. H. B. Barlow and P. Foldiak, "Adaptation and Decorrelation in the Cortex," in *The Computing Neuron*, R. Durbin, C. Miall, and G. J. Mitchison, eds., Addison-Wesley, Wokingham, 1989, pp. 54–72.
243. Q. Zaidi and A. G. Shapiro, "Adaptive Orthogonalization of Opponent-Color Signals," *Biological Cybernetics* **69**:415–428 (1993).
244. J. J. Atick, Z. Li, and A. N. Redlich, "What Does Post-Adaptation Color Appearance Reveal about Cortical Color Representation?," *Vision Research* **33**:123–129 (1993).
245. K. Gegenfurtner and D. C. Kiper, "Contrast Detection in Luminance and Chromatic Noise," *Journal of the Optical Society of America A* **9**:1880–1888 (1992).
246. M. J. Sankeralli and K. T. Mullen, "Postreceptoral Chromatic Detection Mechanisms Revealed by Noise Masking in Three-Dimensional Cone Contrast Space," *Journal of the Optical Society of America A* **14**:2633–2646 (1997).
247. T. Hansen and K. Gegenfurtner, "Higher Level Chromatic Mechanisms for Image Segmentation," *Journal of Vision* **6**:239–259 (2006).
248. M. D'Zmura and K. Knoblauch, "Spectral Bandwidths for the Detection of Colour," *Vision Research* **38**:3117–3128 (1998).
249. G. Monaci, G. Menegaz, S. Süsstrunk, and K. Knoblauch, "Chromatic Contrast Detection in Spatial Chromatic Noise," *Visual Neuroscience* **21**:291–294 (2005).
250. F. Giulianini and R. T. Eskew, Jr., "Theory of Chromatic Noise Masking Applied to Testing Linearity of S-cone Detection Mechanisms," *Journal of the Optical Society of America A* **24**:2604–2021 (2007).
251. D. L. MacAdam, "Visual Sensitivities to Color Differences in Daylight," *Journal of the Optical Society of America* **32**:247–274 (1942).
252. Y. Le Grand, "Les Seuils Différentiels de Couleurs dans la Théorie de Young," *Revue d'Optique* **28**:261–278 (1949).

253. R. M. Boynton and N. Kambe, "Chromatic Difference Steps of Moderate Size Measured Along Theoretically Critical Axes," *Color Research and Application* **5**:13–23 (1980).
254. A. L. Nagy, R. T. Eskew, Jr., and R. M. Boynton, "Analysis of Color-Matching Ellipses in a Cone-Excitation Space," *Journal of the Optical Society of America A* **4**:756–768 (1987).
255. J. Krauskopf and K. Gegenfurtner, "Color Discrimination and Adaptation," *Vision Research* **11**:2165–2175 (1992).
256. J. Romero, J. A. Garcia, L. Jiménez del Barco, and E. Hita, "Evaluation of Color-Discrimination Ellipsoids in Two-color Spaces," *Journal of the Optical Society of America A* **10**:827–837 (1993).
257. R. W. Rodieck, *The Vertebrate Retina*, Freeman, San Francisco, 1973.
258. R. M. Boynton, R. T. Eskew, Jr., and A. L. Nagy, "Similarity of Normalized Discrimination Ellipses in the Constant-Luminance Chromaticity Plane," *Perception* **15**:755–763 (1986).
259. R. T. Eskew, Jr., Q. Wang, and D. P. Richters, "A Five-Mechanism Model of Hue Sensations," *Journal of Vision* **4**:315 (2004).
260. T. N. Cornsweet, and H. M. Pinsker, "Luminance Discrimination of Brief Flashes Under Various Conditions of Adaptation," *Journal of Physiology* **176**:294–310 (1965).
261. F. W. Campbell and J. J. Kulikowski, "Orientation Selectivity of the Human Visual System," *Journal of Physiology* **187**:437–445 (1966).
262. J. Nachmias and E. C. Kocher, "Discrimination of Luminance Increments," *Journal of the Optical Society of America* **60**:382–389 (1970).
263. J. Nachmias and R. V. Sansbury, "Grating Contrast: Discrimination May Be Better Than Detection," *Vision Research* **14**:1039–1042 (1974).
264. J. M. Foley and G. Legge, "Contrast Detection and Near-Threshold Discrimination in Human Vision," *Vision Research* **21**:1041–1053 (1981).
265. K. K. De Valois and E. Switkes, "Simultaneous Masking Interactions between Chromatic and Luminance Gratings," *Journal of the Optical Society of America* **73**:11–18 (1983).
266. C. F. Stromeyer, III and S. Klein, "Spatial Frequency Channels in Human Vision as Asymmetric (edge) Mechanisms," *Vision Research* **14**:1409–1420 (1974).
267. D. G. Pelli, "Uncertainty Explains Many Aspects of Visual Contrast Detection and Discrimination," *Journal of the Optical Society of America A* **2**:1508–1532 (1985).
268. E. Switkes, A. Bradley, and K. K. De Valois, "Contrast Dependence and Mechanisms of Masking Interactions among Chromatic and Luminance Gratings," *Journal of the Optical Society of America A* **5**:1149–1162 (1988).
269. K. T. Mullen and M. A. Losada, "Evidence for Separate Pathways for Color and Luminance Detection Mechanisms," *Journal of the Optical Society of America A* **11**:3136–3151 (1994).
270. R. W. Bowen and J. K. Cotten, "The Dipper and Bumper: Pattern Polarity Effects in Contrast Discrimination," *Investigative Ophthalmology and Visual Science (supplement)* **34**:708 (1993).
271. R. Hiltz, G. Huppmann, and C. R. Cavonius, "Influence of Luminance Contrast on Hue Discrimination," *Journal of the Optical Society of America* **64**:763–766 (1974).
272. C.-C. Chen, J. M. Foley, and D. H. Brainard, "Detection of Chromoluminance Patterns on Chromoluminance Pedestals I: Threshold Measurements," *Vision Research* **40**:773–788 (2000).
273. C.-C. Chen, J. M. Foley, and D. H. Brainard, "Detection of Chromoluminance Patterns on Chromoluminance Pedestals II: Model," *Vision Research* **40**:789–803 (2000).
274. B. A. Wandell, "Measurement of Small Color Differences," *Psychological Review* **89**:281–302 (1982).
275. B. A. Wandell, "Color Measurement and Discrimination," *Journal of the Optical Society of America A* **2**:62–71 (1985).
276. A. C. Beare, "Color-Name as a Function of Wavelength," *The American Journal of Psychology* **76**:248–256 (1963).
277. G. S. Brindley, *Physiology of the Retina and the Visual Pathway*, Williams and Wilkins, Baltimore, 1970.
278. E. Schrödinger, "Über das Verhältnis der Vierfarben zur Dreifarben-theorie," *Sitzungsberichte. Abt. 2a, Mathematik, Astronomie, Physik, Meteorologie und Mechanik. Akademie der Wissenschaften in Wien, Mathematisch-Naturwissenschaftliche Klasse* **134**:471 (1925).
279. D. Jameson and L. M. Hurvich, "Some Quantitative Aspects of an Opponent-Colors Theory. I. Chromatic Responses and Spectral Saturation," *Journal of the Optical Society of America* **45**:546–552 (1955).

280. L. M. Hurvich and D. Jameson, "Some Quantitative Aspects of an Opponent-Colors Theory. II. Brightness, Saturation, and Hue in Normal and Dichromatic Vision," *Journal of the Optical Society of America* **45**:602–616 (1955).
281. D. Jameson and L. M. Hurvich, "Some Quantitative Aspects of an Opponent-Colors Theory. III. Changes in Brightness, Saturation, and Hue with Chromatic Adaptation," *Journal of the Optical Society of America* **46**:405–415 (1956).
282. L. M. Hurvich and D. Jameson, "Some Quantitative Aspects of an Opponent-Colors Theory. IV. A Psychological Color Specification System," *Journal of the Optical Society of America* **46**:1075–1089 (1956).
283. L. H. Hurvich, *Color Vision*, Sinauer, Sunderland, Massachusetts, 1981.
284. R. M. Boynton and J. Gordon, "Bezold-Brücke Hue Shift Measured by a Color-Naming Technique," *Journal of the Optical Society of America* **55**:78–86 (1965).
285. J. Gordon and I. Abramov, "Color Vision in the Peripheral Retina. II. Hue and Saturation," *Journal of the Optical Society of America* **67**:202–207 (1977).
286. J. S. Werner and B. R. Wooten, "Opponent Chromatic Mechanisms: Relation to Photopigments and Hue Naming," *Journal of the Optical Society of America* **69**:422–434 (1979).
287. C. R. Ingling, Jr., J. P. Barley, and N. Ghani, "Chromatic Content of Spectral Lights," *Vision Research* **36**:2537–2551 (1996).
288. A. Brückner, "Zur Frage der Eichung von Farbensystemen," *Zeitschrift für Sinnesphysiologie* **58**:322–362 (1927).
289. E. J. Chichilnisky and B. A. Wandell, "Trichromatic Opponent Color Classification," *Vision Research* **39**:3444–3458 (1999).
290. F. L. Dimmick and M. R. Hubbard, "The Spectral Location of Psychologically Unique Yellow, Green, and Blue," *American Journal of Psychology* **52**:242–254 (1939).
291. M. Ayama, T. Nakatsue, and P. E. Kaiser, "Constant Hue Loci of Unique and Binary Balanced Hues at 10, 100, and 1000 Td," *Journal of the Optical Society of America A* **4**:1136–1144 (1987).
292. B. E. Shefrin and J. S. Werner, "Loci of Spectral Unique Hues Throughout the Life-Span," *Journal of the Optical Society of America A* **7**:305–311 (1990).
293. G. Jordan and J. D. Mollon, "Rayleigh Matches and Unique Green," *Vision Research* **35**:613–620 (1995).
294. J. L. Nerger, V. J. Volbrecht, and C. J. Ayde, "Unique Hue Judgments as a Function of Test Size in the Fovea and at 20-deg Temporal Eccentricity," *Journal of the Optical Society of America A* **12**:1225–1232 (1995).
295. V. J. Volbrecht, J. L. Nerger, and C. E. Harlow, "The Bimodality of Unique Green Revisited," *Vision Research* **37**:404–416 (1997).
296. R. W. Pridmore, "Unique and Binary Hues as Functions of Luminance and Illuminant Color Temperature, and Relations with Invariant Hues," *Vision Research* **39**:3892–3908 (1999).
297. R. G. Kuehni, "Variability in Unique Hue Selection: A Surprising Phenomenon," *Color Research and Application* **29**:158–162 (2004).
298. M. A. Webster, E. Miyahara, G. Malkoc, and V. E. Raker, "Variations in Normal Color Vision. II. Unique Hues," *Journal of the Optical Society of America A* **17**:1545–1555 (2000).
299. R. G. Kuehni, "Determination of Unique Hues Using Munsell Color Chips," *Color Research and Application* **26**:61–66 (2001).
300. M. A. Webster, S. M. Webster, S. Bharadwaj, R. Verma, J. Jaikumar, G. Madan, and E. Vaithilingham, "Variations in Normal Color Vision. III. Unique Hues in Indian and United States Observers," *Journal of the Optical Society of America A* **19**:1951–1962 (2002).
301. W. d. W. Abney, "On the Change of Hue of Spectrum Colors by Dilution with White Light," *Proceedings of the Royal Society of London* **A83**:120–127 (1910).
302. K. Richter, "Antagonistische Signale beim Farbensehen und ihr Zusammenhang mit der empfindungsgemässen Farbordnung," PhD thesis, University of Basel, 1969.
303. J. D. Mollon and G. Jordan, "On the Nature of Unique Hues," in *John Dalton's Colour Vision Legacy*, I. M. D. C. C. Dickinson, ed., Taylor and Francis, London, 1997, pp. 381–392.
304. J. H. Parsons, *An Introduction to Colour Vision*, 2nd ed., Cambridge University Press, Cambridge, 1924.
305. A. Valberg, "A Method for the Precise Determination of Achromatic Colours Including White," *Vision Research* **11**:157–160 (1971).

306. D. Jameson and L. M. Hurvich, "Opponent-Response Functions Related to Measured Cone Photopigments," *Vision Research* **58**:429–430 (1968).
307. J. Larimer, D. H. Krantz, and C. M. Cicerone, "Opponent-Process Additivity—I: Red/Green Equilibria," *Vision Research* **14**:1127–1140 (1974).
308. C. R. Ingling, Jr., "The Spectral Sensitivity of the Opponent-Color Channels," *Vision Research* **17**:1083–1089 (1977).
309. J. Larimer, D. H. Krantz, and C. M. Cicerone, "Opponent-Process Additivity—II: Yellow/Blue Equilibria and Nonlinear Models," *Vision Research* **15**:723–731 (1975).
310. J. G. W. Raaijmakers and C. M. M. de Weert, "Linear and Nonlinear Opponent Color Coding," *Perception and Psychophysics* **18**:474–480 (1975).
311. Y. Ejima and Y. Takahashi, "Bezold-Brücke Hue Shift and Nonlinearity in Opponent-Color Process" *Vision Research* **24**:1897–1904 (1984).
312. S. Takahashi and Y. Ejima, "Spatial Properties of Red-Green and Yellow-Blue Perceptual Opponent-Color Response," *Vision Research* **24**:987–994 (1984).
313. M. Ayama, P. K. Kaiser, and T. Nakatsue, "Additivity of Red Chromatic Valence," *Vision Research* **25**:1885–1891 (1985).
314. C. R. Ingling, Jr., P. W. Russel, M. S. Rea, and B. H.-P. Tsou, "Red-Green Opponent Spectral Sensitivity: Disparity between Cancellation and Direct Matching Methods," *Science* **201**:1221–1223 (1978).
315. C. H. Elzinga and C. M. M. de Weert, "Nonlinear Codes for the Yellow/Blue Mechanism," *Vision Research* **24**:911–922 (1984).
316. M. Ayama and M. Ikeda, "Additivity of Yellow Chromatic Valence," *Vision Research* **26**:763–769 (1985).
317. K. Knoblauch and S. K. Shevell, "Relating Cone Signals to Color Appearance: Failure of Monotonicity in Yellow/Blue," *Visual Neuroscience* **18**:901–906 (2001).
318. M. Ikeda and M. Ayama, "Non-linear Nature of the Yellow Chromatic Valence," in *Colour Vision: Physiology and Psychophysics*, J. D. Mollon and L. T. Sharpe, eds., Academic Press, London, 1983, pp. 345–352.
319. K. Knoblauch, L. Sirovich, and B. R. Wooten, "Linearity of Hue Cancellation in Sex-Linked Dichromacy," *Journal of the Optical Society America A* **2**:136–146 (1985).
320. W. von Bezold, "Über das Gesetz der Farbenmischung und die physiologischen Grundfarben," *Annalen der Physiologie und Chemie* **150**:221–247 (1873).
321. E. W. Brücke, "Über einige Empfindungen im Gebiete der Sehnerven," *Sitzungsberichte der Akademie der Wissenschaften in Wien, Mathematisch-Naturwissenschaftliche Klasse, Abteilung 3* **77**:39–71 (1878).
322. D. M. Purdy, "Spectral Hue as a Function of Intensity," *American Journal of Psychology* **63**:541–559 (1931).
323. J. J. Vos, "Are Unique and Invariant Hues Coupled?," *Vision Research* **26**:337–342 (1986).
324. P. L. Walraven, "On the Bezold-Brücke Phenomenon," *Journal of the Optical Society of America* **51**:1113–1116 (1961).
325. J. Walraven, "Discounting the Background—The Missing Link in the Explanation of Chromatic Induction," *Vision Research* **16**:289–295 (1976).
326. S. K. Shevell, "The Dual Role of Chromatic Backgrounds in Color Perception," *Vision Research* **18**:1649–1661 (1978).
327. S. K. Shevell, "Color Perception under Chromatic Adaptation: Equilibrium Yellow and Long-Wavelength Adaptation," *Vision Research* **22**:279–292 (1982).
328. P. Whittle, "The Brightness of Coloured Flashes on Backgrounds of Various Colours and Luminances," *Vision Research* **13**:621–638 (1973).
329. C. M. Cicerone, D. H. Krantz, and J. Larimer, "Opponent-Process Additivity—III: Effect of Moderate Chromatic Adaptation," *Vision Research* **15**:1125–1135 (1975).
330. J. von Kries, "Influence of Adaptation on the Effects Produced by Luminous Stimuli," in *Sources of Color Science (1970)*, D. L. MacAdam, ed., MIT Press, Cambridge, MA, 1905, pp. 120–1126.
331. R. W. Burnham, R. W. Evans, and S. M. Newhall, "Influence on Color Perception of Adaptation to Illumination," *Journal of the Optical Society of America* **42**:597–605 (1952).
332. H. V. Walters, "Some Experiments on the Trichromatic Theory of Vision," *Proceedings of the Royal Society of London* **131**:27–50 (1942).

333. D. L. MacAdam, "Influence of Chromatic Adaptation on Color Discrimination and Color Perception," *Die Farbe* **4**:133–143 (1955).
334. E. J. Chichilnisky and B. A. Wandell, "Photoreceptor Sensitivity Changes Explain Color Appearance Shifts Induced by Large Uniform Backgrounds in Dichoptic Matching," *Vision Research* **35**:239–254 (1995).
335. L. H. Hurvich and D. Jameson, "Further Developments of a Quantified Opponent Colors Theory," in *Visual Problems of Colour*, Volume 2, Her Majesty's Stationery Office, London, 1958, pp. 691–723.
336. D. Jameson and L. M. Hurvich, "Sensitivity, Contrast, and Afterimages," in *Visual Psychophysics*, Vol. VII/4, *Handbook of Sensory Physiology*, D. Jameson and L. H. Hurvich, eds., Springer-Verlag, Berlin, 1972, pp. 568–581.
337. S. K. Shevell, "Color Appearance," in *The Science of Color*, (2nd ed.), S. K. Shevell, ed., Elsevier, Oxford, 2003, pp. 149–190.
338. J. Walraven, "No Additive Effect of Backgrounds in Chromatic Induction," *Vision Research* **19**:1061–1063 (1979).
339. S. K. Shevell, "Unambiguous Evidence for the Additive Effect in Chromatic Adaptation," *Vision Research* **20**:637–639 (1980).
340. B. Drum, "Additive Effect of Backgrounds in Chromatic Induction," *Vision Research* **21**:959–961 (1981).
341. E. H. Adelson, "Looking at the World through a Rose-Colored Ganzfeld," *Vision Research* **21**:749–750 (1981).
342. J. Larimer, "Red/Green Opponent Colors Equilibria Measured on Chromatic Adapting Fields: Evidence for Gain Changes and Restoring Forces," *Vision Research* **21**:501–512 (1981).
343. J. Wei and S. K. Shevell, "Color Appearance under Chromatic Adaptation Varied Along Theoretically Significant Axes in Color Space," *Journal of the Optical Society of America A* **12**:36–46 (1995).
344. O. Rinner and K. Gegenfurtner, "Time Course of Chromatic Adaptation for Color Appearance and Discrimination," *Vision Research* **40**:1813–1826 (2000).
345. J. M. Hillis and D. H. Brainard, "Distinct Mechanisms Mediate Visual Detection and Identification," *Current Biology* **17**:1714–1719 (2007).
346. P. K. Kaiser, "Minimally Distinct Border as a Preferred Psychophysical Criterion in Heterochromatic Photometry," *Journal of the Optical Society of America* **61**:966–971 (1971).
347. A. Kohlrausch, "Theoretisches und Praktisches zur heterochromen Photometrie," *Pflügers Archiv für die gesamte Physiologie des Menschen und der Tiere* **200**:216–220 (1923).
348. H. Helson and V. B. Jeffers, "Fundamental Problems in Color Vision. II. Hue, Lightness, and Saturation of Selective Samples in Chromatic Illumination," *Journal of Experimental Psychology* **26**:1–27 (1940).
349. R. W. Burnham, R. M. Evans, and S. M. Newhall, "Prediction of Color Appearance with Different Adaptation Illuminations," *Journal of the Optical Society of America* **47**:35–42 (1957).
350. J. J. McCann, S. P. McKee, and T. H. Taylor, "Quantitative Studies in Retinex Theory: A Comparison between Theoretical Predictions and Observer Responses to the 'Color Mondrian' Experiments," *Vision Research* **16**:445–458 (1976).
351. L. E. Arend and A. Reeves, "Simultaneous Color Constancy," *Journal of the Optical Society of America A* **3**:1743–1751 (1986).
352. D. H. Brainard, W. A. Brunt, and J. M. Speigle, "Color Constancy in the Nearly Natural Image. 1. Asymmetric Matches," *Journal of the Optical Society of America A* **14**:2091–2110 (1997).
353. D. H. Brainard, "Color Constancy in the Nearly Natural Image. 2. Achromatic Loci," *Journal of the Optical Society of America A* **15**:307–325 (1998).
354. J. M. Kraft and D. H. Brainard, "Mechanisms of Color Constancy under Nearly Natural Viewing," *Proceedings of the National Academy of Sciences of the United States of America*, **96**:307–312 (1999).
355. D. H. Brainard, "Color Constancy," in *The Visual Neurosciences*, L. Chalupa and J. Werner, eds., MIT Press, Cambridge, MA, 2004, pp. 948–961.
356. H. E. Smithson, "Sensory, Computational, and Cognitive Components of Human Color Constancy," *Philosophical Transactions of the Royal Society of London B* **360**:1329–1346 (2005).
357. S. K. Shevell and F. A. A. Kingdom, "Color in Complex Scenes," *Annual Review of Psychology* **59**:143–166 (2008).

358. G. Buchsbaum, "A Spatial Processor Model for Object Colour Perception," *Journal of the Franklin Institute* **310**:1–26 (1980).
359. L. T. Maloney and B. A. Wandell, "Color Constancy: A Method for Recovering Surface Spectral Reflectances," *Journal of the Optical Society of America A* **3**:29–33 (1986).
360. D. H. Brainard and W. T. Freeman, "Bayesian Color Constancy," *Journal of the Optical Society of America A* **14**:1393–1411 (1997).
361. B. V. Funt, M. S. Drew, and J. Ho, "Color Constancy from Mutual Reflection," *International Journal of Computer Vision* **6**:5–24 (1991).
362. M. D'Zmura and G. Iverson, "Color Constancy. III. General Linear Recovery of Spectral Descriptions for Lights and Surfaces," *Journal of the Optical Society of America A* **11**:2389–2400 (1994).
363. G. D. Finlayson, P. H. Hubel, and S. Hordley, "Color by Correlation," in *Proceedings of the IS&T/SID Fifth Color Imaging Conference*, Scottsdale, AZ, 1997, pp. 6–11.
364. A. C. Hurlbert, "Computational Models of Color Constancy," in *Perceptual Constancy: Why Things Look As They Do*, V. Walsh and J. Kulikowski, eds., Cambridge University Press, Cambridge, 1998, pp. 283–322.
365. L. T. Maloney and J. N. Yang, "The Illuminant Estimation Hypothesis and Surface Color Perception," in *Colour Perception: From Light to Object*, R. Mausfeld and D. Heyer, eds., Oxford University Press, Oxford, 2001, pp. 335–358.
366. D. H. Brainard, J. M. Kraft, and P. Longère, "Color Constancy: Developing Empirical Tests of Computational Models," in *Colour Perception: Mind and the Physical World*, R. Mausfeld and D. Heyer, eds., Oxford University Press, Oxford, 2003, pp. 307–334.
367. D. H. Brainard, P. Longere, P. B. Delahunt, W. T. Freeman, J. M. Kraft, and B. Xiao, "Bayesian Model of Human Color Constancy," *Journal of Vision* **6**:1267–1281 (2006).
368. H. Boyaci, L. T. Maloney, and S. Hersh, "The Effect of Perceived Surface Orientation on Perceived Surface Albedo in Binocularly Viewed Scenes," *Journal of Vision* **3**:541–553 (2003).
369. H. Boyaci, K. Doerschner, and L. T. Maloney, "Perceived Surface Color in Binocularly Viewed Scenes with Two Light Sources Differing in Chromaticity," *Journal of Vision* **4**:664–679 (2004).
370. M. Bloj, C. Ripamonti, K. Mitha, S. Greenwald, R. Hauck, and D. H. Brainard, "An Equivalent Illuminant Model for the Effect of Surface Slant on Perceived Lightness," *Journal of Vision* **4**:735–746 (2004).
371. E. H. Land and J. J. McCann, "Lightness and Retinex Theory," *Journal of the Optical Society of America* **61**:1–11 (1971).
372. E. H. Land, "The Retinex Theory of Color Vision," *Scientific American* **237**:108–128 (1977).
373. E. H. Land, "Recent Advances in Retinex Theory," *Vision Research* **26**:7–21 (1986).
374. A. Hurlbert, "Formal Connections between Lightness Algorithms," *Journal of the Optical Society of America A* **3**:1684–1694 (1986).
375. G. West and M. H. Brill, "Necessary and Sufficient Conditions for von Kries Chromatic Adaptation to Give Color Constancy," *Journal of Mathematical Biology* **15**:249–250 (1982).
376. J. A. Worthey and M. H. Brill, "Heuristic Analysis of von Kries Color Constancy," *Journal of the Optical Society of America A* **3**:1708–1712 (1986).
377. D. H. Brainard and B. A. Wandell, "Analysis of the Retinex Theory of Color Vision," *Journal of the Optical Society of America A* **3**:1651–1661 (1986).
378. D. H. Foster and S. M. C. Nascimento, "Relational Colour Constancy from Invariant Cone-Excitation Ratios," *Proceedings of the Royal Society of London. Series B: Biological Sciences* **257**:115–121 (1994).
379. M. A. Webster and J. D. Mollon, "Colour Constancy Influenced by Contrast Adaptation," *Nature* **373**:694–698 (1995).
380. Q. Zaidi, B. Spehar, and J. DeBonet, "Color Constancy in Variegated Scenes: Role of Low-level Mechanisms in Discounting Illumination Changes," *Journal of the Optical Society of America A* **14**:2608–2621 (1997).
381. M. D'Zmura and B. Singer, "Contrast Gain Control," in *Color Vision: From Molecular Genetics to Perception*, K. Gegenfurtner and L. T. Sharpe, eds., Cambridge University Press, Cambridge, 1999, pp. 369–385.
382. W. S. Stiles, "Mechanism Concepts in Colour Theory," *Journal of the Colour Group* **11**:106–123 (1967).
383. D. Krantz, "A Theory of Context Effects Based on Cross-Context Matching," *Journal of Mathematical Psychology* **5**:1–48 (1968).

384. D. H. Brainard and B. A. Wandell, "Asymmetric Color-matching: How Color Appearance Depends on the Illuminant," *Journal of the Optical Society of America A* **9**:1433–1448(1992).
385. Q. Zaidi, "Identification of Illuminant and Object Colors: Heuristic-Based Algorithms," *Journal of the Optical Society of America A* **15**:1767–1776 (1998).
386. J. Golz and D. I. A. MacLeod, "Influence of Scene Statistics on Colour Constancy," *Nature* **415**:637–640 (2002).
387. R. M. Boynton, M. M. Hayhoe, and D. I. A. MacLeod, "The Gap Effect: Chromatic and Achromatic Visual Discrimination as Affected by Field Separation," *Optica Acta* **24**:159–177(1977).
388. R. T. Eskew, Jr. and R. M. Boynton, "Effects of Field Area and Configuration on Chromatic and Border Discriminations," *Vision Research* **27**:1835–1844 (1987).
389. B. W. Tansley and R. M. Boynton, "A Line, Not a Space, Represents Visual Distinctness of Borders Formed by Different Colors" *Science* **191**:954–957 (1976).
390. B. Pinna, "Un Effetto di Colorazione," in *Il laboratorio e la città. XXI Congresso degli Psicologi Italiani*, V. Majer, M. Maeran, and M. Santinello, eds., Società Italiana di Psicologia, Milano, 1987, p. 158.
391. B. Pinna, G. Brelstaff, and L. Spillmann, "Surface Color from Boundaries: A New 'Watercolor' Illusion," *Vision Research* **41**:2669–2676 (2001).
392. B. Pinna, J. S. Werner, and L. Spillmann, "The Watercolor Effect: A New Principle of Grouping and Figure-ground Organization," *Vision Research* **43**:43–52 (2003).
393. F. Devinck, P. B. Delahunt, J. L. Hardy, L. Spillmann, and J. S. Werner, "The Watercolor Effect: Quantitative Evidence for Luminance-Dependent Mechanisms of Long-Range Color Assimilation," *Vision Research* **45**:1413–1424 (2005).
394. E. D. Montag, "Influence of Boundary Information on the Perception of Color," *Journal of the Optical Society of America A* **14**:997–1006 (1997).
395. R. T. Eskew, Jr. "The Gap Effect Revisited: Slow Changes in Chromatic Sensitivity as Affected by Luminance and Chromatic Borders," *Vision Research* **29**:717–729 (1989).
396. P. D. Gowdy, C. F. Stromeyer, III, and R. E. Kronauer, "Facilitation between the Luminance and Red-Green Detection Mechanisms: Enhancing Contrast Differences Across Edges," *Vision Research* **39**:4098–4112 (1999).
397. L. A. Riggs, F. Ratliff, J. C. Cornsweet, and T. C. Cornsweet, "The Disappearance of Steadily Fixated Visual Test Objects," *Journal of the Optical Society of America* **43**:495–501 (1953).
398. J. Krauskopf, "Effect of Retinal Image Stabilization on the Appearance of Heterochromatic Targets," *Journal of the Optical Society of America* **53**:741–744 (1963).
399. A. L. Yarbus, *Eye Movements and Vision*, Plenum Press, New York, 1967.
400. T. P. Piantanida and J. Larimer, "The Impact of Boundaries on Color: Stabilized Image Studies," *Journal of Imaging Technology* **15**:58–63 (1989).
401. H. D. Crane and T. Piantanida, "On Seeing Reddish Green and Yellowish Blue," *Science* **221**:1078–1080 (1983).
402. V. A. Billock, G. A. Gleason, and B. H. Tsou, "Perception of Forbidden Colors in Retinally Stabilized Equiluminance Images: An Indication of Softwired Cortical Color Opponency," *Journal of the Optical Society of America A* **10**:2398–2403 (2001).
403. J. L. Nerger, T. P. Piantanida, and J. Larimer, "Color Appearance of Filled-in Backgrounds Affects Hue Cancellation, but not Detection Thresholds," *Vision Research* **33**:165–172(1993).
404. J. J. Wisowaty and R. M. Boynton, "Temporal Modulation Sensitivity of the Blue Mechanism: Measurements Made Without Chromatic Adaptation," *Vision Research* **20**:895–909 (1980).
405. T. P. Piantanida, "Temporal Modulation Sensitivity of the Blue Mechanism: Measurements Made with Extraretinal Chromatic Adaptation," *Vision Research* **25**:1439–1444(1985).
406. N. W. Daw, "Why After-Images are Not Seen in Normal Circumstances," *Nature* **196**:1143–1145(1962).
407. R. van Lier and M. Vergeer, "Filling in the Afterimage after the Image," *Perception (ECVP Abstract Supplement)* **36**:200–201 (2007).
408. R. van Lier, M. Vergeer, and S. Anstis, "'Mixing-In' Afterimage Colors," *Perception (ECVP Abstract Supplement)* **37**:84 (2008).
409. C. McCollough, "Color Adaptation of Edge Detectors in the Human Visual System," *Science* **149**:1115–1116 (1965).

410. P. D. Jones and D. H. Holding, "Extremely Long-Term Persistence of the McCollough Effect," *Journal of Experimental Psychology: Human Perception & Performance* **1**:323–327 (1975).
411. E. Vul, E. Krizay, and D. I. A. MacLeod, "The McCollough Effect Reflects Permanent and Transient Adaptation in Early Visual Cortex," *Journal of Vision* **8**:4, 1–12 (2008).
412. G. M. Murch, "Binocular Relationships in a Size and Color Orientation Specific Aftereffect," *Journal of Experimental Psychology* **93**:30–34 (1972).
413. R. L. Savoy, " 'Extinction' of the McCollough Effect does not Transfer Interocularly," *Perception & Psychophysics* **36**:571–576 (1984).
414. E. Vul and D. I. A. MacLeod, "Contingent Aftereffects Distinguish Conscious and Preconscious Color Processing," *Nature Neuroscience* **9**:873–874 (2006).
415. P. Thompson and G. Latchford, "Colour-Contingent After-Effects Are Really Wavelength-Contingent," *Nature* **320**:525–526 (1986).
416. D. H. Hubel and T. N. Wiesel, "Receptive Fields, Binocular Interaction and Functional Architecture in the Cat's Visual Cortex," *Journal of Physiology* **160**:106–154 (1962).
417. D. H. Hubel and T. N. Wiesel, "Receptive Fields and Functional Architecture of Monkey Striate Cortex," *Journal of Physiology* **195**:215–243 (1968).
418. C. F. Stromeyer, III, "Form-Color Aftereffects in Human Vision," in *Perception, Handbook of Sensory Physiology*, Vol. VIII, R. Held, H. W. Leibowitz, and H. L. Teuber, eds., Springer-Verlag, Berlin, 1978, pp. 97–142.
419. D. Skowbo, B. N. Timney, T. A. Gentry, and R. B. Morant, "McCollough Effects: Experimental Findings and Theoretical Accounts," *Psychological Bulletin* **82**:497–510 (1975).
420. H. B. Barlow, "A Theory about the Functional Role and Synaptic Mechanism of Visual After-effects," in *Visual Coding and Efficiency*, C. B. Blakemore, ed., Cambridge University Press, Cambridge, 1990, pp. 363–375.
421. P. Dodwell and G. K. Humphrey, "A Functional Theory of the McCollough Effect," *Psychological Review* **97**:78–89 (1990).
422. G. K. Humphrey and M. A. Goodale, "Probing Unconscious Visual Processing with the McCollough Effect," *Consciousness and Cognition* **7**:494–519 (1998).
423. C. McCollough, "Do McCollough Effects Provide Evidence for Global Pattern Processing?," *Perception & Psychophysics* **62**:350–362 (2000).
424. V. H. Perry, R. Oehler, and A. Cowey, "Retinal Ganglion Cells that Project to the Dorsal Lateral Geniculate Nucleus in the Macaque Monkey," *Neuroscience* **12**:1101–1123 (1984).
425. T. N. Wiesel and D. Hubel, "Spatial and Chromatic Interactions in the Lateral Geniculate Body of the Rhesus Monkey," *Journal of Neurophysiology* **29**:1115–1156 (1966).
426. R. L. De Valois and P. L. Pease, "Contours and Contrast: Responses of Monkey Lateral Geniculate Nucleus Cells to Luminance and Color Figures," *Science* **171**:694–696 (1971).
427. H. Kolb and L. Dekorver, "Midget Ganglion Cells of the Parafovea of the Human Retina: A Study by Electron Microscopy and Serial Reconstructions," *Journal of Comparative Neurology* **303**:617–636 (1991).
428. D. J. Calkins, S. J. Schein, Y. Tsukamoto, and P. Sterling, "M and L Cones in Macaque Fovea Connect to Midget Ganglion Cells by Different Numbers of Excitatory Synapses," *Nature* **371**:70–72 (1994).
429. B. B. Boycott, J. M. Hopkins, and H. G. Sperling, "Cone Connections of the Horizontal Cells of the Rhesus Monkey's retina," *Proceedings of the Royal Society of London. Series B: Biological Sciences* **B229**:345–379 (1987).
430. W. Paulus and A. Kröger-Paulus, "A New Concept of Retinal Colour Coding," *Vision Research* **23**:529–540 (1983).
431. P. Lennie, P. W. Haake, and D. R. Williams, "The Design of Chromatically Opponent Receptive Fields," in *Computational Models of Visual Processing*, M. S. Landy and J. A. Movshon, eds., MIT Press, Cambridge, MA, 1991, pp. 71–82.
432. R. C. Reid and R. M. Shapley, "Spatial Structure of Cone Inputs to the Receptive Fields in Primate Lateral Geniculate Nucleus," *Nature* **356**:716–718 (1992).
433. B. B. Lee, J. Kremers and T. Yeh, "Receptive Fields of Primate Retinal Ganglion Cells Studied with a Novel Technique," *Visual Neuroscience* **15**:161–175 (1998).
434. D. M. Dacey, "Parallel Pathways for Spectral Coding in Primate Retina," *Annual Review of Neuroscience* **23**:743–775 (2000).

435. K. T. Mullen and F. A. A. Kingdom, "Differential Distributions of Red-Green and Blue-Yellow Cone Opponency across the Visual Field," *Visual Neuroscience* **19**:109–118 (2002).
436. L. Diller, O. S. Packer, J. Verweij, M. J. McMahon, D. R. Williams, and D. M. Dacey, "L and M Cone Contributions to the Midget and Parasol Ganglion Cell Receptive Fields of Macaque Monkey Retina," *Journal of Neuroscience* **24**:1079–1088 (2004).
437. S. G. Solomon, B. B. Lee, A. J. White, L. Rüttiger, and P. R. Martin, "Chromatic Organization of Ganglion Cell Receptive Fields in the Peripheral Retina," *Journal of Neuroscience* **25**:4527–4539 (2005).
438. C. Vakrou, D. Whitaker, P. V. McGraw, and D. McKeefry, "Functional Evidence for Cone-Specific Connectivity in the Human Retina," *Journal of Physiology* **566**:93–102 (2005).
439. P. Buzas, E. M. Blessing, B. A. Szmajda, and P. R. Martin, "Specificity of M and L Cone Inputs to Receptive Fields in the Parvocellular Pathway: Random Wiring with Functional Bias," *Journal of Neuroscience* **26**:11148–11161 (2006).
440. P. R. Jusuf, P. R. Martin, and U. Grünert, "Random Wiring in the Midget Pathway of Primate Retina," *Journal of Neuroscience* **26**:3908–3917 (2006).
441. R. W. Rodieck, "What Cells Code for Color?," in *From Pigments to Perception. Advances in Understanding Visual Processes*, A. Valberg, and B. B. Lee, eds., Plenum, New York, 1991.
442. T. Cornsweet, *Visual Perception*, Academic Press, New York, 1970.
443. P. Lennie, "Recent Developments in the Physiology of Color Vision," *Trends in Neurosciences* **7**:243–248 (1984).
444. E. Martinez-Uriegas, "A Solution to the Color-Luminance Ambiguity in the Spatiotemporal Signal of Primate X Cells," *Investigative Ophthalmology and Visual Science (supplement)* **26**:183 (1985).
445. V. A. Billock, "The Relationship between Simple and Double Opponent Cells," *Vision Research* **31**:33–42 (1991).
446. F. A. A. Kingdom and K. T. Mullen, "Separating Colour and Luminance Information in the Visual System," *Spatial Vision* **9**:191–219 (1995).
447. J. Nathans, D. Thomas, and S. G. Hogness, "Molecular Genetics of Human Color Vision: The Genes Encoding Blue, Green and Red Pigments," *Science* **232**:193–202 (1986).
448. J. D. Mollon, "'Tho' She Kneel'd in That Place Where They Grew...': The Uses and Origins of Primate Colour Vision," *Journal of Experimental Biology* **146**:21–38 (1989).
449. M. S. Livingstone and D. H. Hubel, "Segregation of Form, Color, Movement, and Depth: Anatomy, Physiology, and Perception," *Science* **240**:740–749 (1988).
450. D. T. Lindsey and A. M. Brown, "Masking of Grating Detection in the Isoluminant Plane of DKL Color Space," *Visual Neuroscience* **21**:269–273 (2004).
451. A. Li and P. Lennie, "Mechanisms Underlying Segmentation of Colored Textures," *Vision Research* **37**:83–97 (1997).
452. T. Hansen and K. Gegenfurtner, "Classification Images for Chromatic Signal Detection," *Journal of the Optical Society of America A* **22**:2081–2089 (2005).
453. H. E. Smithson, S. Khan, L. T. Sharpe, and A. Stockman, "Transitions between Colour Categories Mapped with Reverse Stroop Interference and Facilitation," *Visual Neuroscience* **23**:453–460 (2006).
454. Q. Zaidi and D. Halevy, "Visual Mechanisms That Signal the Direction of Color Changes," *Vision Research* **33**:1037–1051 (1986).
455. M. D'Zmura, "Color in Visual Search," *Vision Research* **31**:951–966 (1991).
456. J. Krauskopf, Q. Zaidi, and M. B. Mandler, "Mechanisms of Simultaneous Color Induction," *Journal of the Optical Society of America A* **3**:1752–1757 (1986).
457. J. Krauskopf, H. J. Wu, and B. Farell, "Coherence, Cardinal Directions and Higher-order Mechanisms," *Vision Research* **36**:1235–1245 (1996).
458. A. Valberg, "Unique Hues: An Old Problem for a New Generation," *Vision Research* **41**:1645–1657 (2001).
459. J. A. Movshon, I. D. Thompson, and D. J. Tolhurst, "Spatial Summation in the Receptive Fields of Simple Cells in the Cat's Striate Cortex," *Journal of Physiology* **283**:53–77 (1978).
460. H. Spitzer and S. Hochstein, "Complex-Cell Receptive Field Models," *Progress in Neurobiology* **31**:285–309 (1988).

461. J. Krauskopf and Q. Zaidi, "Induced Desensitization," *Vision Research* **26**:759–762 (1986).
462. M. J. Sankeralli and K. T. Mullen, "Bipolar or Rectified Chromatic Detection Mechanisms?," *Visual Neuroscience* **18**:127–135 (2001).
463. P. D. Gowdy, C. F. Stromeyer, III, and R. E. Kronauer, "Detection of Flickering Edges: Absence of a Red-Green Edge Detector," *Vision Research* **39**:4186–4191 (1999).
464. M. Sakurai and K. T. Mullen, "Cone Weights for the Two Cone-Opponent Systems in Peripheral Vision and Asymmetries of Cone Contrast Sensitivity," *Vision Research* **46**:4346–4354 (2006).
465. A. Vassilev, M. S. Mihaylova, K. Racheva, M. Zlatkova, and R. S. Anderson, "Spatial Summation of S-cone ON and OFF Signals: Effects of Retinal Eccentricity," *Vision Research* **43**:2875–2884 (2003).
466. J. S. McClellan and R. T. Eskew, Jr., "ON and OFF S-cone Pathways Have Different Long-Wave Cone Inputs" *Vision Research* **40**:2449–2465 (2000).
467. A. G. Shapiro and Q. Zaidi, "The Effects of Prolonged Temporal Modulation on the Differential Response of Color Mechanisms," *Vision Research* **32**:2065–2075 (1992).
468. Q. Zaidi, A. G. Shapiro, and D. C. Hood, "The Effect of Adaptation on the Differential Sensitivity of the S-cone Color System," *Vision Research* **32**:1297–1318 (1992).
469. N. V. S. Graham, *Visual Pattern Analyzers*, Oxford University Press, New York, 1989.
470. A. B. Watson and J. G. Robson, "Discrimination at Threshold: Labelled Detectors in Human Vision," *Vision Research* **21**:1115–1122 (1981).
471. S. L. Guth, "Comments on 'A Multi-Stage Color Model'" *Vision Research* **36**:831–833 (1996).
472. R. L. De Valois and K. K. De Valois, "On 'A Three-Stage Color Model,'" *Vision Research* **36**:833–836 (1996).
473. M. E. Chevreul, *De la loi du Contraste Simultané des Couleurs*, Pitois-Levreault, Paris, 1839.
474. A. Kitaoka, "Illusion and Color Perception," *Journal of the Color Science Association of Japan* **29**:150–151 (2005).
475. K. Sakai, "Color Representation by Land's Retinex Theory and Belsey's Hypothesis," Gradual thesis, *Department of Psychology*, Ritsumeikan University, Japan, 2003.
476. W. von Bezold, *Die Farbenlehre in Hinblick auf Kunst und Kunstgewerbe*, Westermann, Braunschweig, 1874.
477. L. T. Sharpe, A. Stockman, H. Jägle, and J. Nathans, "Opsin Genes, Cone Photopigments, Color Vision and Colorblindness," in *Color Vision: From Genes to Perception*, K. Gegenfurtner and L. T. Sharpe, eds., Cambridge University Press, Cambridge, 1999, pp. 3–51.
478. P. Lennie, "Parallel Visual Pathways: A Review," *Vision Research* **20**:561–594 (1980).

ASSESSMENT OF REFRACTION AND REFRACTIVE ERRORS AND THEIR INFLUENCE ON OPTICAL DESIGN

B. Ralph Chou

*School of Optometry
University of Waterloo
Waterloo, Ontario, Canada*

12.1 GLOSSARY

Definitions

Accommodation. The increase of refractive power of the eye by changing the shape of the crystalline lens that enables focusing of the eye on a near object.

Amblyopia. Reduced visual acuity that is not improved with corrective lenses and occurs in the absence of anatomical or pathological anomalies in the eye.

Ametropia. A refractive condition of the unaccommodated eye in which light from optical infinity does not focus on the retina.

Aniseikonia. A relative difference in the perceived size and/or shape of the images in the two eyes.

Anisometropia. Unequal refractive state in the two eyes, which may result in aniseikonia.

Aphakia. Absence of the crystalline lens from the eye.

Astigmatism. Refractive state of the eye in which rays from a point object form two line images at different distances from the retina.

Back vertex power. Reciprocal of the distance in meters between the pole of the back surface of a lens and the axial position of the image formed of an infinitely distant object. The unit of BVP is the **diopter** (or reciprocal meter).

Base curve of a contact lens. The radius of curvature of the ocular surface of a contact lens.

Base curve of a spectacle lens. The power of the flattest meridian on the front surface of a spectacle lens. Alternatively, for a spectacle lens design, the reference surface power for a series of lenses of different back vertex powers.

Emmetropia. Refractive state of an unaccommodated eye in which light from an infinitely distant object is focused on the retina.

Far point. A point in the object space of an ametropic eye which is conjugate to the retina.

Presbyopia. Reduced ability to accommodate, usually due to age-related changes in the crystalline lens.

Prismatic effect. The deviation of a ray of light as it passes through an optical system, considered as if due to a prism of known deviating power that replaces the optical system. The unit of prismatic effect is the **prism diopter**, which is 1 cm of displacement per meter of travel of the ray along the optical axis.

Rimless mounting. A system of mounting spectacle lenses in which there is no frame; all parts for resting the spectacles on the nose and ears are attached directly to the lenses.

Visual acuity. Clinical measure of the minimum angle of resolution of the eye, with or without corrective lenses.

Working distance. Distance of the object of regard in front of the eye, usually considered 40 cm for reading and other near tasks.

Symbols

f	spectacle lens focal length
F	spectacle lens power
F_v	back vertex power of a spectacle lens
F_x	effective back vertex power of a spectacle lens displaced a distance x
K	keratometry reading of corneal curvature
L	axial length of the eye
P	power of intraocular implant lens
x	displacement distance of a spectacle lens
P_H	prismatic effect in a lens of power F at a point x cm from the optical center

Equations

Equation (1) is the lens effectivity equation which is used to adjust the power of a spectacle lens when it is moved a distance x meter from its original position.

F_x	effective back vertex power of a spectacle lens displaced a distance x
F_v	original back vertex power of the spectacle lens
x	displacement distance of the spectacle lens

Equation (2) estimates the spectacle correction needed after cataract extraction.

F	back vertex power of the postoperative spectacle lens
F_{old}	back vertex power of the preoperative spectacle lens

Equation (3) is the SRKII formula for power of an intraocular lens needed to produce emmetropia in an eye after cataract removal.

P	IOL power
$A1$	a constant
K	keratometer reading of the central corneal curvature
L	axial length of the eye in mm

Equation (4) is Prentice's Rule.

P_H	prismatic effect
x	distance in centimeters of the point on the lens through which the line of sight passes from the optical center of the lens
F	the back vertex power of the lens

Equation (5) is the exact expression for spectacle magnification of a spectacle lens.

M	magnification
d_v	distance between the back vertex of the lens and the cornea

- F back vertex power of the spectacle lens
- t axial thickness of the spectacle lens in meters
- n index of refraction of the lens
- F_1 front surface power of the lens

Equation (6) is the approximate formula for spectacle magnification.

12.2 INTRODUCTION

At the beginning of the twenty-first century, much of the information upon which we rely comes to us through vision. Camera view finders, computer displays, liquid crystal displays at vehicle controls, and optical instruments are only a few examples of devices in which refractive error and its correction can have important implications for both their users and optical designers. The options for correction of refractive error are many and varied; each has its advantages and disadvantages when the use of optical devices is considered. Most observers are binocular, and whether one or both eyes are fully or optimally corrected may affect one's ability to use a given optical device or instrument. In addition, the optical and physiological changes associated with aging of the eye may also have a profound effect on visual performance.

12.3 REFRACTIVE ERRORS

As described by Charman¹ in Chap. 1, the eye consists of two refractive elements, the cornea and the crystalline lens, which are separated by the watery aqueous humor, and the gel-like vitreous humor, which fills the rest of the eyeball between the crystalline lens and the retina (Fig. 1). The cornea provides approximately one-third of the optical power of the eye and the crystalline lens the remainder. Light entering the cornea from optical infinity ideally is focused through the dioptics of the eye onto the retina, where it stimulates the photoreceptors and triggers the cascade of neurophysiological processes leading to the visual percept. An eye in which the retinal image is in sharp focus is described as emmetropic. The shape of the crystalline lens can be changed through the action of intraocular muscles in the process called *accommodation*. Accommodation increases the power of the eye so that light from a “near” object at a finite distance in front of the eye is sharply focused on the *retina*. The neural mechanisms governing accommodation affect the neural control of convergent and divergent eye movement. This ensures that as the eyes fixate on a near object, both eyes form images that overlap, giving rise to a single binocular percept. The crystalline lens continues to grow throughout life and physical changes within it gradually cause a loss of flexibility which reduces the amplitude of accommodation. The clinical onset of *presbyopia* (old eye) is marked by the inability to maintain focus on a target viewed at distance of 40 cm in front of the eye. This normally occurs at the age of 40 to 50 years.

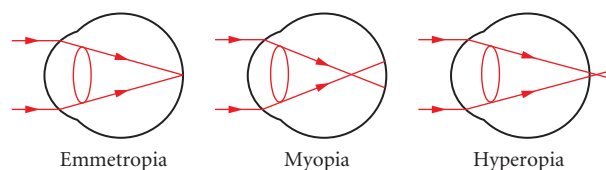


FIGURE 1 Focusing collimated light in emmetropic, myopic, and hyperopic eyes.

The failure to accurately focus light from a remote object onto the retina is referred to as *ametropia*.¹ Ametropia or refractive error is present in many eyes. Uncorrected ametropia has been identified as one of the leading treatable causes of blindness around the world.² The degree of ametropia is quantified by the optical power F of the spectacle lens that “corrects” the focus of the eye by bringing light from the distant object to focus at the retina. The dioptric power of the lens is given by the reciprocal of its focal length f : $F = 1/f$. The unit of focal power is the *diopter*, abbreviated as D, which has dimensions of reciprocal meters. Most spectacle lenses can be regarded as thin lenses,³ and the power of superimposed thin lenses can be found as the algebraic sum of the powers of the individual lenses. This is convenient in ophthalmic applications such as measurement of refractive error and the use of “trial” corrections made with combinations of loose or trial case lenses.

Types of Refractive Error

An unaccommodated eye which focuses light from a distant object onto the retina is *emmetropic* and requires no corrective lenses (Fig. 1).

A “near sighted” or *myopic* eye focuses light in front of the retina (Fig. 1). The myope sees closer objects in focus, but distant objects are blurred. Myopia is corrected by lenses of negative power that diverge light.

A “far sighted” or *hyperopic* eye focuses light behind the retina (Fig. 1). The hyperope can exercise accommodation to increase the power of the eye and focus the image of a distant object clearly on the retina. Nearer objects can also be seen clearly by exercising more accommodation, but the increased effort to maintain accommodation and a single binocular percept may result in symptoms of “eye-strain” and blurred near point vision. If accommodative fatigue occurs, the observer’s ability to maintain a clearly focused image of the distant object is lost, and blur at both distance and near results. Measurement of the refractive error of a hyperope often results in a finding of apparent emmetropia when accommodation is being exercised (latent hyperopia) and varying degrees of manifest hyperopia when control of accommodation is lost. Hyperopia is corrected by lenses of positive power that converge light.

The foregoing discussion assumes that the dioptics of the eye bring collimated light to a point focus. However, most eyes are *astigmatic*, that is, the optical components form two line foci that are usually perpendicular to one another. *Astigmatism* results in blur for all object distances, and the greater the astigmatism, the greater the blur. If the orientation of the line foci is in a direction away from the horizontal and vertical, the effects of ocular astigmatism on vision may be very serious, hindering development of the visual system.

Astigmatism is corrected with a *spherocylindrical lens* (Fig. 2), which usually has one spherical and one toric surface. Optically, the cylinder component brings the two line foci of the eye together,

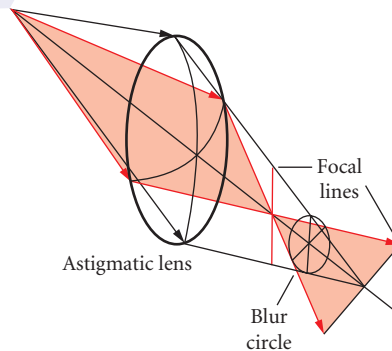


FIGURE 2 The astigmatic pencil formed by a spherocylindrical lens.

and the sphere component places the superimposed line foci on the retina. In clinical practice, the lens prescription is written in terms of *sphere*, *cylinder*, and *axis*. The position of the cylinder axis is specified using the trigonometric coordinate system centered on the eye as seen by the examiner, with 0° to the examiner's right or the patient's left, in what is referred to as the TABO or Standard Notation.³ For example, the prescription $+1.00 - 2.50 \times 030$ corresponds to a +1.00 D spherical lens (sometimes designated +1.00 DS) superimposed on a -2.50 D cylindrical lens with its axis at 30° (sometimes designated -2.00 DC \times 030). Occasionally, for clinical reasons, astigmatism may be corrected by the *spherical equivalent lens*, a spherical lens that causes the line foci to straddle the retina. The power of the spherical equivalent is calculated as the sum of the sphere plus one-half of the cylinder in the spectacle prescription. The spherical equivalent of the above prescription would be -0.25 D.

Spectacle lenses are usually manufactured in quarter diopter steps. By ophthalmic convention, spectacle prescriptions are written with a plus or minus sign and to two decimal places. When the power is less than ± 1.00 D, a 0 is written before the decimal point, for example, +0.25 D.

Refractive errors are usually one of these simple types because the curvatures of the various refractive surfaces of the eye vary slowly over the area of the entrance pupil.^{4,5} In some instances, however, *irregular astigmatism* may occur and no lens of any type can produce sharp vision. These are usually cases of eye injury or disease, such as distortion of the cornea in advanced keratoconus, corneal scarring after severe trauma, corneal dystrophy, disruption of the optics of the crystalline lens due to cataract, and traumatic damage to the lens. Tilting or displacement of the crystalline lens and distortion of the cornea following rigid contact lens wear may also lead to irregular refraction.⁶ Williams⁷ has described a case of irregular refraction due to tilting of the retina.

12.4 ASSESSMENT OF REFRACTIVE ERROR

In ophthalmic practice, the measurement of refractive error or ocular refraction is an iterative process in which a series of tests is employed to refine the results until a satisfactory end point is reached. The usual starting point is an objective test where no patient interaction is required, followed by subjective tests that require patient responses. Clinical tests are generally accurate to within ± 0.25 D, which is comparable to the depth of field of the human eye.⁸

Objective Tests

Objective tests require only that the patients fix their gaze upon a designated target, normally located at a distance of 6 m from the eye so that with accurate fixation there is a negligible amount of accommodation being exercised. Three common tests are keratometry, direct ophthalmoscopy, and retinoscopy.

Keratometry measures the curvature of the front surface of the cornea. A bright object, the *keratometer mire*, is reflected by the cornea (1st Purkinje image) and its magnification determined from the lateral displacement of a doubled image of the mire.⁹ The curvature of the corneal surface (K reading) in its principal meridians is read from a calibrated dial along with the axis of any corneal astigmatism. The total power of the cornea may be estimated from the K reading, and the corneal astigmatism may be used to predict the total amount and orientation of the astigmatism of the whole eye. Various rules of thumb have been developed (e.g., Javal's rule) to estimate the total astigmatism from the keratometric readings of normal eyes. This is due to the observation that large degrees of ocular astigmatism are almost always due to the cornea. In an *aphakic* eye, from which the crystalline lens has been removed, the corneal astigmatism is the ocular astigmatism.

Although primarily designed to view the retina of the living eye, the *direct ophthalmoscope* can also be used to estimate refractive error. Light is directed from a source in the ophthalmoscope through the pupil to illuminate the patient's retina. The examiner views the resulting image through lenses of varying power until a clear view of the retina is obtained. The algebraic sum of the lens power

and refractive error of the examiner is the patient's refractive error. When astigmatism is taken into account, this result is a crude estimate of the spherical equivalent of the patient's ocular refraction and is not sufficiently accurate to prescribe a corrective lens power. It may be helpful in assessing the refractive error of incommunicative patients when other methods cannot be employed.

The most common objective measurement of ocular refraction is retinoscopy. The retinoscope is designed to view the red reflex in the pupil as light is reflected from the retina. This is the same red reflex that often is seen in flash photographs of faces. The retinoscope is held at a known *working distance* from the eye, and the examiner observes the red reflex through the peephole in the retinoscope mirror as the retinoscope beam is moved across the patient's eye. Vignetting of the beam by the pupils of the patient and the examiner results in a movement of the reflex across the pupil. The direction and speed of the reflex motion depends on the patient's refractive error.¹⁰

Most refractionists use either a slightly divergent or collimated retinoscope beam at a working distance of 50 or 66 cm from the eye. If the red reflex moves in the same direction as the beam, this is a "with" movement and indicates that the retina is conjugate to a point closer than the working distance. As lenses with increasing plus power are interposed in 0.25 D increments, the with movement becomes faster until the end point is reached where a lens neutralizes the refractive error and the reflex moves infinitely fast or instantly fills the pupil when the beam touches the edge of the pupil. The next lens in the sequence should cause the reflex motion to reverse direction. The patient's refractive error is the value of the end point lens minus the working distance correction (+2.00 D for 50 cm working distance and +1.50 D for 66 cm). When the red reflex shows an "against" movement (it moves in the opposite direction to the retinoscope beam), the retina is conjugate to a point farther than the working distance.

Neutralization is achieved by interposing lenses of increasing minus power until the reflex shows a with movement. The end point lens is one incremental step back in power; with this lens, the reflex instantly fills the pupil when the beam touches the edge of the pupil. The patient's refractive error is calculated by correcting for the examiner's working distance.

The shape and motion of the reflex can be used to determine the axis of astigmatism. Once the two principal meridians of the eye are identified, the refractive error is measured in each separately. The sphere, cylinder, and axis components of the ocular refraction can then be calculated.

Retinoscopy provides an estimate of the refractive error very quickly, but requires some skill. A streak-shaped beam can facilitate the process in many patients, but retinoscopic findings may still vary from the true ocular refraction in some patients.

Automated equipment has been developed to determine the refractive error and corneal curvature without depending on the examiner's skill or judgment. Autokeratometers and autorefractors are frequently used to provide a starting point for the subjective refraction. More recently, aberrometers, instruments that measure ocular aberrations as well as the paraxial ocular refraction, have entered the ophthalmic market.

Subjective Techniques

The final refinement of the refractive findings in an oculo-visual examination is usually obtained through subjective refraction. Although many techniques have been developed over the years, in practice only two or three are used. These methods rely on patient interaction to determine the sphere, cylinder, and axis components of the refractive error.

Subjective methods rely on the use of a visual acuity chart and trial lenses that can be placed in a trial frame worn by the patient, or a phoropter, a mechanical device that contains a range of powers of spherical and cylindrical lenses, prisms and accessory filters and lenses. The visual acuity chart may be either printed or projected and contains letters, numbers, or pictographs arranged in a sequence of sizes from large at the top to small at the bottom. The target size is rated according to the distance at which the finest detail within it would subtend 1 min of arc. Both recognition and resolution of detail within the target are required to read the smallest characters on the chart. The *visual acuity* is a representation of the minimum angle of resolution of the eye viewing the chart with or without corrective lenses. For example, an eye with a 1.5 min minimum angle of resolution should just read

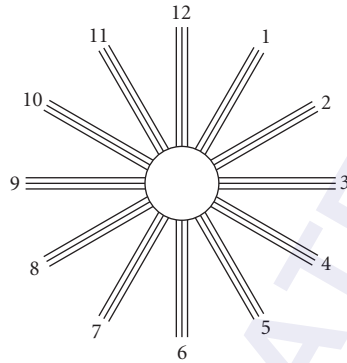


FIGURE 3 The astigmatic dial.

a target rated for 9 m or 30 ft when viewed at a distance of 6 m or 20 ft. The visual acuity would be recorded as 6/9 in metric, or 20/30 in imperial units. The refractive error can be estimated as the quotient of the denominator in imperial units divided by 100—in this example the refractive error would be approximately ± 0.30 D spherical equivalent. We usually expect ametropic patients with healthy eyes to have corrected visual acuity of 6/4.5 (20/15) or 6/6 (20/20). Rabbetts¹¹ has an excellent review of the design and use of visual acuity charts.

Subjective refraction usually begins with the patient's retinoscopic finding or previous spectacle correction in the phoropter or trial frame. Sometimes it may be preferable or necessary to add plus lenses in 0.25 D steps before the unaided eye to blur or fog the patient, then reduce lens power or add minus power incrementally until the best possible visual acuity is achieved with a spherical lens power (unfogging).

The astigmatism is measured next. If the patient has an equivalent sphere lens in front of the eye, vision is fogged while viewing the astigmatic dial, which is a chart with radial lines arranged in 30° intervals (Fig. 3). The patient identifies the line that appears clearest or darkest; the line corresponds to the line focus nearest to the retina and its orientation allows for calculation of the axis of the astigmatism. Since the chart looks like an analog clock dial, the hour value times 30 gives the approximate axis value, where 12 o'clock has the value 0 for the TABO axis notation. Cylinder power is added with this axis position until all lines appear equally dark or clear.

Refinement of the astigmatic correction is achieved with the *Jackson cross-cylinder* check test. This uses a spherocylindrical lens with equal but opposite powers in its principal meridians; normally a ± 0.25 D or ± 0.50 D lens is used. The Jackson cross-cylinder lens is designed to rotate around one of its principal meridians to check cylinder power and around an axis 45° from that meridian to check cylinder axis orientation. The patient looks through the spherocylinder correction and Jackson cross-cylinder lens and notes changes in clarity of small letters on the acuity chart as the cross cylinder is flipped and adjustments are made first to the axis position, then the cylinder power until both cross-cylinder orientations give the same degree of clarity. The sphere is then adjusted for best corrected visual acuity. Further adjustments may be made to equalize the visual acuity between the two eyes before the final corrective power is confirmed with trial lenses and frame.

Presbyopia

Presbyopia or “old eye” is the refractive condition in which the ability to accommodate or focus the eyes on a near target is reduced. This is an age-related phenomenon due to the gradually increasing rigidity of the crystalline lens. It generally becomes clinically significant in the mid to late 40s; absolute presbyopia occurs by the late 50s when there is no accommodation. The onset of presbyopia is

earlier among hyperopes whose first symptom is the inability to focus on distant objects that were previously seen clearly without corrective lenses. Myopes can compensate for early presbyopia by removing their corrective lenses to view near objects. Presbyopia is corrected by adding positive power (the reading addition or add) to the distance correction until the patient is able to achieve the desired near visual acuity for reading. Single vision reading glasses or bifocal lenses may be prescribed for near work.

The added power may be determined either objectively using techniques such as dynamic retinoscopy or subjectively using a variety of techniques.¹²

12.5 CORRECTION OF REFRACTIVE ERROR

The Art of Prescribing

Ametropia is usually corrected by placing a lens in front of the eye that will focus light from a distant object onto the axial position that is conjugate to the retina of the ametropic eye. This position is the *far point* or *punctum remotum* of the eye. Near objects can then be seen clearly through the lens by exercising accommodation.

The fact that humans usually have two functioning eyes and the frequent occurrence of astigmatism in ametropia require eye care practitioners to consider a given patient's binocular status and visual needs in addition to the refractive errors of the eyes. Prescribing corrective lenses is therefore as much an art as a science.

The spatial distortion inherent in astigmatic corrections can cause adaptation problems, whether or not the patient is well adapted to spectacle correction. As a result, most practitioners will try to minimize the amount of cylindrical correction that is needed for clear comfortable vision. In some cases, particularly in first corrections, the spherical equivalent may be preferred over even a partial correction of the astigmatism to facilitate adaptation to the lenses.

Younger ametropes may find vision slightly more crisp when they exercise a small degree of accommodation when looking at a distant target. Beginning refractionists often err in "over minus-ing" patients who accommodate during the refraction. The methods of subjective refraction are designed to minimize this tendency. While absolute presbyopes may be expected to require a full distance correction, and one expects myopic individuals will normally require full correction of their ametropia, younger hyperopes may often show a manifest error that is considerably lower than what can be found in cycloplegic refraction. They may require a partial correction until they learn to relax their accommodation when looking at distant objects.

Anatomical and neuromuscular anomalies of extraocular muscles that control eye movement may result in ocular deviations, both latent and manifest, that affect the ability of the visual system to provide a clear binocular visual percept for all target distances in all directions of gaze through corrective lenses. Differential *prismatic effects* may lead to double vision when the patient looks in certain directions or at certain distances in front of the eyes. The practitioner may need to modify the lens prescription in order to provide comfortable vision.

As presbyopia develops, a reading addition is required to maintain single clear binocular vision during near visual tasks like reading. The power of the add usually is chosen so that the patient uses about one-half of the available accommodation when reading. Addition power is determined using a standardized testing distance of 40 cm, then adjusted using trial lenses and the patient's preferred working distance for near vision tasks. The design of the lenses (bifocal, multifocal, invisible bifocal, single vision readers) to be prescribed largely depends on the wearer's visual tasks, intended use of the lenses, and need for eye protection.

It should be noted that presbyopes are not the only ones who can benefit from using bifocals and reading glasses. Since these corrections can reduce the demand for accommodation, they may be prescribed for younger patients who have anomalous accommodation and complain of eyestrain, visual discomfort, and even double vision. Finally, some individuals with low degrees of myopia often find that they are most comfortable when reading without any spectacle correction at all.

Spectacle Correction

Optical Considerations Spectacles have been in use for at least 700 years, although it was only after the late nineteenth century that they were generally used to correct distance vision. Prior to that time, the lenses were self-selected mostly to aid in reading, and rarely for distance correction. The now familiar *bent* shape of modern spectacle lenses was not generally adopted until the early 1900s.

Spectacle lenses are often treated as thin lenses and it is often assumed that their exact position in front of the eyes is not critical. These assumptions apply to the majority of patients who have low-to-moderate refractive errors. It is only when we consider lenses for correction of high refractive error that thick lens optics must be used.

Spectacle lenses are specified by their *back vertex power*, which is the reciprocal of the distance from the pole of the ocular or back surface of the lens to its second focal point. The powers of phoropter and trial case lenses are also given as back vertex powers. The *focimeter* or *lensometer* is an instrument used to measure the back vertex power of spectacle lenses.¹³

The back vertex power of a corrective lens depends on the *vertex distance*, which is the distance between the back surface of the lens and the cornea. The vertex distance is typically between 12 and 15 mm. When the vertex distance of the phoropter or trial frame used in determining a patient's refractive error differs from the actual vertex distance of the corrective lenses placed in the spectacle frame, the back vertex power of the lenses can be modified using the formula

$$F_x = \frac{F_v}{1 - xF_v} \quad (1)$$

where x is the distance in meters between the vertex distance of the refraction and the vertex distance of the spectacles ($x > 0$ when the spectacles are moved closer to the eye). It can be shown that changes in vertex distance of 1 or 2 mm are insignificant unless the lens power is greater than 8.00 D.¹⁴

Of all the parameters that are considered in designing spectacle lenses, the most significant is the *base curve*, which defines the lens form. The base curve is adjusted so that the lens forms images of distant objects near the far point of the wearer's eye in all directions of gaze. Ideally the image surface coincides with the far point sphere of the eye, so that the same spectacle prescription is required in all fields of gaze.¹⁵ In practice, this means that the lens design will minimize unwanted oblique astigmatism and control curvature of field. Many corrected curve or best form lens designs have been developed, each with its own assumptions about vertex distance, field size, thickness, material, and relative importance of the Seidel aberrations, using spherical curves in the principal meridians of the lens. High plus power lenses must be made with aspheric curves to achieve adequate performance. Astigmatic lenses are made in minus cylinder form, with a spherical front surface and toric back.

The introduction of computer-aided design and computer numerically controlled (CNC) surfacing machines has led to the free form surfacing technology of today's ophthalmic lens industry. By specifying the deviation of both lens surfaces from reference spherical surfaces at many points, a diamond cutter can be used to form any desired surface shape. Double aspheric lenses which may provide optimal visual performance for wearers of any spectacle correction can be made readily with this technology.

Spectacle Lens Materials Spectacle lenses of all types are available in many glass and plastic materials. The final choice of material for a given pair of lenses depends on the diameter and thickness of the lenses (which determine their weight), the need for impact and scratch resistance, and the sensitivity of the wearer to color fringes caused by transverse chromatic aberration. The glass lenses of the twentieth century have been largely replaced by plastics, the most common being CR39 with a refractive index of 1.498. Polycarbonate (index 1.586) is frequently used because of its high impact resistance for protective lenses for occupational and sports use and for patients with high refractive errors. This material requires scratch-resistant coatings because of its soft surface and it suffers from chromatic aberration. Trivex (index 1.53) is a recently introduced plastic material that is more impact resistant than CR39 and can be used in *rimless mountings* and high power prescriptions where polycarbonate shows poorer performance. High index urethane lenses have largely replaced the high index glass lenses of the late 1900s with refractive indices of 1.6 to 1.74 available. Many

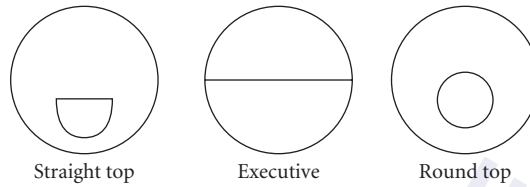


FIGURE 4 Common bifocal styles.

lenses are supplied with antireflection coatings to enhance visual performance and cosmetic appearance; the coating can significantly reduce impact resistance.^{16–18} Tints can be applied for cosmetic, protective, or vision enhancement purposes.¹⁹

Presbyopic Corrections Most presbyopic patients prefer a bifocal correction in which the top part of the spectacle lens provides the distance correction and an area in the lower part, the *segment*, contains the correction for near vision. The difference between the segment and distance power is the add of the bifocal. Adds typically range from +1.00 D to +2.50 D, with the power increasing with age. A small minority of patients prefer separate pairs of glasses for distance use and reading.

The most common bifocal segments are the *straight top*, *Executive*, and *round top* designs shown in Fig. 4. Although many segment designs were developed in the early 1900s, only a few remain commercially available. All plastic bifocals are one piece, that is, the distance and near powers are achieved by changing the surface curvature on one side of the lens. Glass Executive bifocals are also one piece, but the straight top and round top styles are made by fusing a high index glass button into a countersink in the front surface of a lower index glass lens. They are referred to as fused bifocals. Bifocal segments are normally placed so that the top of the segment (the segment line) is at the edge of the lower eyelid. This allows comfortable near vision with minimal downgaze and keeps the segment as inconspicuous as possible. The style and width of the bifocal segment used is mainly determined by the wearer's visual needs at near, particularly with regard to width of the near visual field, as well as cosmetic appearance.

Patients with adds greater than +1.75 D often find that objects at a distance of 50 to 60 cm appear blurred through both the distance and segment portions of their bifocals. A trifocal lens that contains an intermediate segment between the distance and near parts of the lens allows for clear vision at arm's-length distance. Trifocals are positioned with the top line at or just below the lower edge of the pupil.

Special occupational multifocal lenses can be prescribed for presbyopes who need to see intermediate or near objects overhead. These lenses contain an intermediate or near add in the upper part of the lens in addition to the usual bifocal or trifocal segment. They are often referred to as occupational trifocal or quadrifocal lenses, respectively. Such lenses are recommended for presbyopic mechanics, electricians, plumbers, carpenters, painters, and librarians, among others. The *progressive addition lens* or *PAL* is an increasingly popular alternative to bifocal and trifocal lenses. The PAL design (Fig. 5) features a continuously varying power in the lens from the distance visual point, the pupil position when looking at distance, to the near visual point that is used to read. Originally a single aspheric progressive surface on the front of the lens was used, but more recent designs that take advantage of freeform surfacing technology have incorporated the power progression on either or both surfaces. Clear vision is achieved not only at the distance and near visual points, but also through a narrow optical corridor or umbilicus that joins them.

Outside of the distance and near vision zones and the corridor, the changes in surface shape result in distorted blurred images. The myriad of PAL designs in the ophthalmic lens market represents the many solutions to designing the progressive surface(s) to optimize single clear binocular vision through as much of the lens as possible. Some manufacturers claim to include preference for head or eye movements to view objects to the side, differences in prismatic effect between the eyes, and other factors in their designs.²⁰ Although marketed as lineless bifocals, PALs have proven popular because they hide the fact the wearer is old enough to need a reading prescription and provide clear vision at

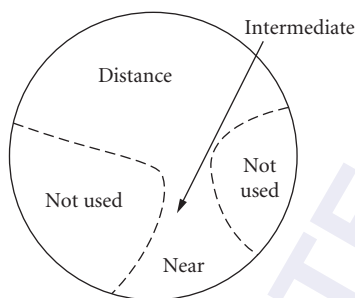


FIGURE 5 A progressive addition lens (PAL) has a large area of distance prescription power and an area of addition power connected by a narrow corridor. Areas to either side of the corridor and reading area are less usable due to blur induced by the progressive power surface.

all distances once the patient has learned to use them. The lens must be properly positioned in front of the eyes and the patient instructed on how to use it. Most patients adapt quickly to their PALs and prefer them to conventional bifocals despite the extra cost.

Recently PAL technology has been adapted for both patients approaching presbyopia as well as advanced presbyopes, who work with computer displays. These office lenses focus on arm's length and near working distances to provide a wider intermediate field of view than can be achieved in the corridor of a PAL. Prepresbyopes with complaints of eyestrain and presbyopes suffering "bifocal neck" while working at a computer display may find that these lens designs relieve their symptoms. Since the lenses have the intermediate portion in front of the pupil in the straight ahead gaze position, the head need not be tilted backward to enable the wearer to see the screen. Musicians, hairdressers, and other workers who require wide clear intermediate fields may also find these lenses useful.

Contact Lenses

A contact lens can be considered as a spectacle lens with a zero vertex distance. The spectacle lens power must be adjusted for the change in vertex distance from the spectacle plane; in practice lens powers under 5 D are not adjusted for vertex distance since the change in power is less than 0.25 D. The contact lens power will be more positive or less negative than the spectacle lens power. Since the contact lens is supported by the cornea, the fit of the contact lens must be assessed to ensure that the back surface of the contact lens is supported without significantly affecting the shape of the cornea. The tear layer between the cornea and the contact lens may affect the fit as well as the power of the contact lens.

Most contact lenses are prescribed for cosmetic reasons and many patients ask for them in sports activity. Patients with very high refractive errors often find that contact lenses provide better quality of vision, particularly with regard to magnification effects and extent of the visual field, than spectacles. They are also much more comfortable than when wearing thick, heavy spectacle lenses.

The main contraindications for wearing contact lenses are dry eye and allergies.

Rigid Contact Lenses The first contact lenses were made from polymethylmethacrylate (PMMA). This plastic has very good optical and physical properties, but is impermeable to oxygen. Contact lenses had to be designed so that the overall diameter and the curvature of the back surface (the base curve of the contact lens) allowed a lens to move over the eye with the blink and permit tear fluid to be exchanged under the lens. This allowed enough oxygen to reach the cornea to maintain its physiology. Poorly fitting PMMA lenses could cause many problems related to disrupted corneal physiology due to hypoxia, including hypoesthesia, edema, surface irregularities, and overwear syndrome.

Deformation of the cornea could also occur, especially in corneal astigmatism, if the base curve did not provide a good mechanical fit to the toric corneal surface.

Modern rigid lenses made of gas permeable materials²¹ have greatly reduced many of the problems caused by PMMA. These biocompatible materials transmit oxygen and can be made larger and thinner than PMMA lenses without adversely affecting corneal physiology. With less movement, the lenses are more comfortable, although initially there is a period of adaptation. Adaptation usually occurs over a few weeks as the patient starts wearing the lenses for a few hours each day, then gradually increases the wearing time. Eventually most patients wear their contact lenses between 10 and 14 hours a day; some patients may attempt extended wear, when the lenses are worn continuously for several days. Once adapted, the patient must continue to wear the lenses on a schedule that maintains the eyes' adaptation.

Rigid lenses do not conform to the shape of the underlying cornea. Thus, the tears filling the space between the ocular surface of the contact lens and the cornea form a *tear lens*. The sum of the power of the tear lens plus the power of the contact lens itself must equal the vertex adjusted power of the spectacle correction. Very thin gas permeable lenses tend to flex and change shape on the eye because of lid pressure and the tear lens. Most contact lens fitters will use trial lenses to estimate how the lens behaves on the eye and thus choose a base curve and lens power that will provide both a good physical fit and optimal vision correction.

Because the tears and cornea have almost the same index of refraction, the optical power of the cornea is replaced by that of the front surface of the tear layer beneath the lens. A spherical base curve therefore effectively masks the corneal astigmatism that often contributes most of the refractive astigmatism in ametropia. As a result, many patients with moderate astigmatism can be managed satisfactorily with spherical contact lenses, provided that an acceptable fit of the base curve to the cornea can be achieved. If the astigmatism is too great, or if a suitable physical fit cannot be achieved, rigid contact lenses with toric front and/or back surfaces can be used. In these cases, the lenses must be oriented using prism ballast in which prism power is ground into the lens so that the thicker base edge is at the bottom of the lens—its orientation is maintained by gravity.

Rigid lenses, particularly those made of PMMA, will change the shape of the cornea.²² The curvature change may vary from day to day, making it difficult to determine a spectacle prescription for use when the patient is not wearing the contact lenses. This is described as spectacle blur. In extreme cases, the warped cornea may have a permanent irregular astigmatism. *Orthokeratology* is a clinical approach to reducing or eliminating refractive error by corneal molding with rigid contact lenses.²³ The procedure involves a series of contact lenses with successively flatter base curves to reduce corneal curvature. After several months, the patient experiences a mild reduction of 1 or 2 D of myopia. This is a temporary (i.e., several hours) effect. Retainer lenses must be worn for several hours each day to prevent the cornea from rebounding to its original shape.

Most rigid lenses are made with a handling tint, often a very light blue, brown, or green, to facilitate locating a lens dropped on the floor or countertop. The tint should not affect the wearer's colour perception. The right lens of a pair is normally identified with a dot engraved on the front surface. Lenses must be cleaned and disinfected with special care solutions after each wearing. The fit of the lens should be checked every year. The expected service life of a rigid gas permeable contact lens is 2 to 3 years.

Hydrogel Lenses Soft or hydrogel contact lenses were introduced in the 1970s and today they dominate the contact lens market. A hydrogel is a polymer that is able to absorb a significant amount of water.²¹ Although brittle and fragile when dry, a hydrogel becomes soft and flexible when hydrated. The first commercially successful lens, the Bausch & Lomb Soflens®, was made of hydroxyethylmethacrylate (HEMA), which contains 38.6 percent water by weight when fully hydrated. Manufacturers have since developed HEMA derivatives with water content up to 77 percent.²⁴ A hydrogel lens that is saturated with ophthalmic saline solution buffered to the pH of the tears will smoothly drape over the cornea if an appropriate base curve has been selected. The base curve of the lens is determined from the value of K , the corneal curvature measured by a keratometer. A trial lens of known power with this base curve is placed on the eye and an over refraction is performed to determine the final contact lens power. A properly sized hydrogel lens will cover the entire cornea

and overlap the sclera by about 1 mm; the peripheral flattening of the cornea is accommodated by flattening of the edge zone of the ocular surface of the contact lens. The optic zone of the lens, where the base curve is found, covers most of the cornea.

Most soft contact lenses are spherical with small amounts of astigmatism being corrected with the equivalent sphere of the prescription. Because a hydrogel lens conforms more to the shape of the cornea than a rigid lens, higher amounts of astigmatism must be corrected with a toric lens. Maintaining the orientation of a toric lens on the eye is a challenge, since interaction of the lens with the underlying cornea, the tear layer, and the action of the lids will influence how it sits on the eye. Various methods have been used by manufacturers to stabilize lens position.

Soft lenses can be worn on a daily, disposable/frequent replacement basis. In daily wear, the lenses are worn up to 12 to 14 hours a day and removed each night. The lenses must be cleaned and disinfected before they are worn again. Cleaning with digital rubbing of the lens surfaces removes mucous deposits while disinfection eliminates infectious agents. These lenses can be worn for about 1 year. Disposable or frequent replacement lenses are worn between 1 day and several weeks on a daily basis, then discarded. They are prescribed for patients who experience rapid build-up of protein deposits or who have limited access to cleaning and sterilizing solutions. Extended wear lenses are worn up to a week at a time before being removed for cleaning and disinfection. Their higher water content and reduced thickness increase oxygen permeability but make them more susceptible to damage when handled. Extended wear lenses are usually replaced monthly.

One important advantage of hydrogel over rigid lens materials is that the lens is comfortable to wear with little or no adaptation time. Tear exchange is not as important since hydrogel transmits oxygen to the cornea. Thin lenses with high water content transmit significantly more oxygen. The lenses do not dislodge easily, making them ideal for sports activity; however, it has been reported that hydrogel lenses worn in chlorinated pool water exhibit significantly more microbial colonization than lenses that were never worn in the pool. This may increase the risk of bacterial keratitis if the lenses are worn in water sports.²⁵ Lenses worn in water sports should be discarded as soon as possible after the activity is finished.²⁶

In the last few years, silicone hydrogel lenses have entered the market. Silicone hydrogel has very high oxygen transmissivity, but is stiffer than HEMA-based hydrogel. Initially there is more awareness of the lens in situ; however, the oxygen permeability of these lenses has largely eliminated many of the complications of earlier hydrogel lenses.

Most soft contact lenses are either clear or have a slight handling tint. Cosmetic lenses with an overall tint can be used to make light irides appear in different colors. Lenses with arrays of tinted dots can be used to change the color appearance of dark irides with varying degrees of success. Novelty tints can be used to change the appearance of the eye (e.g., slit pupils, “unnatural” pupil color); however, there is a controversy over whether such lenses should only be supplied through eye care practitioners because of serious risks to eye health.²⁷ Similar types of lens can be used to disguise scarred or disfigured eyes due to damage or disease.²⁸

Contact Lenses for Presbyopia Presbyopia has become an important clinical problem as patients who began wearing their contact lenses in the 1970s and 1980s enter their 40s or 50s. The quality of vision and convenience of a largely spectacle-free life are compromised. Reading glasses worn over the contact lenses is a simple approach, but defeats the purpose of having contacts in the first place.

Monovision is the practice of fitting one eye for distance correction and the other for near. This approach reduces the need for spectacles, but relies on the patient’s ability to suppress the vision of one or the other eye when looking at a given target. The resulting reduction of binocular vision may create more serious problems involving depth and space perception.

Bifocal contact lenses have met with varying degrees of success. Alternating vision rigid contact lenses have a segment ground into the lens. When the patient looks at a near object the lower eyelid moves the lens over the cornea to bring the segment in front of the pupil. The theory of this design is seldom found in practice.

The simultaneous vision lenses are more successful. Both the distance portion and the segment are within the pupil so that superimposed clear and blurred images of distant and near objects are seen. There is a resultant loss of clarity and contrast. A diffractive contact lens design uses concentric

rings cut into the base curve of the lens to provide a near correction while the front surface curvature provides the distance correction simultaneously.²⁹

For a more detailed discussion of contact lens optics and technology, see Chap. 20.

Refractive Surgery

The desire of many ametropic patients to be free of both spectacles and contact lenses has driven the development of refractive surgery to reduce or neutralize ametropia. Earlier procedures that used incisions in the cornea to alter its shape (radial keratotomy) have since been largely abandoned due to complications including corneal perforation³⁰ and the possibility of rupture of the compromised cornea.³¹ More recently, the excimer laser has been used in *photorefractive keratectomy (PRK)* and *laser in-situ keratomileusis (LASIK)* to reshape the cornea. These procedures are not without risks of compromised vision and complications include stromal haze, regression of refractive effect, infection, and optical and/or mechanical instability of the corneal flap in LASIK.³⁰ Almost all patients achieve visual acuity of at least 20/40 (6/12) and most achieve postoperative visual acuity of at least 20/20 (6/6). LASIK results in increased higher-order ocular aberrations; a new wavefront-guided LASIK procedure improves on this result.³² A more detailed description of PRK and LASIK can be found in Chap. 16.

Extremely high myopia (over 12.00 D) may be treated by clear lens exchange in which the crystalline lens is extracted and an intraocular lens of suitable power is implanted. Phakic lens implants have also been used in either the anterior chamber or posterior chamber.

Low myopes can achieve visual acuity of better than 20/40 (6/12) when treated with intrastromal corneal rings (ICR). Circular arcs are implanted concentrically to the optical axis of the cornea in the stroma of the peripheral cornea to flatten the corneal curvature mechanically. The procedure is at least partly reversed by removal of the rings.

None of these procedures is guaranteed to result in perfect vision. All have some deleterious effect on quality of vision (glare, loss of contrast, increased higher-order aberrations³⁰) and none eliminates the need for glasses for certain visually demanding tasks. As these patients become presbyopic, they will still require some form of reading prescription.

Aphakia and Pseudophakia

Cataract is the general term for a crystalline lens that loses transparency because of aging or injury. Cortical cataract is seen in the outer layers of the crystalline lens, the cortex, and may comprise general haze, punctuate or wedge shaped opacities, and bubbles. In earlier stages, cortical cataract does not greatly affect visual acuity that is measured with high contrast targets, but there may be a significant reduction of contrast sensitivity. Patients may complain of glare in bright sunlight and difficulty reading. Nuclear cataract often appears as a gradual yellowing of the core of the crystalline lens. This is thought to arise from photochemical changes in the lens crystalline proteins triggered by long-term chronic exposure to ultraviolet radiation.³³ Visual consequences include reduced color discrimination and visual acuity.

Treatment of cataract is by surgical removal of the crystalline lens. *Intracapsular cataract extraction (ICCE)* is the removal of the lens and its surrounding lens capsule. ICCE has been largely replaced by *extracapsular cataract extraction (ECCE)* in which the lens is removed, but the capsule remains in the eye. Phacoemulsification is a procedure for breaking up the lens with an ultrasonic probe to facilitate ECCE. A very small incision near the edge of the cornea is required for ECCE, whereas a much larger circumferential incision is required for ICCE. Consequently, the potential for optical and physical postoperative complications in ECCE is much smaller than for ICCE. Approximately one in three patients who undergo ECCE will experience postoperative opacification of the lens capsule. This is remedied by an in-office YAG laser capsulotomy to open a hole in the posterior capsule to restore a clear optical path to the retina.³⁴ An eye which has had its lens removed is described as *aphakic*.

Postsurgical Correction Since the crystalline lens contributes about one-third of the optical power of the eye (see Chap. 1), the aphakic eye requires a high plus optical correction whose approximate power is

$$F = +11.25 + 0.62 F_{\text{old}} \quad (2)$$

where F_{old} is the preoperative equivalent spectacle correction.³⁵ This can be provided in the form of spectacles, contact lenses, or intraocular lens implants.

Aphakic spectacle corrections, sometimes referred to as cataract lenses, are heavy, thick, and cosmetically unattractive. They also have significant optical problems including distortion, field curvature, and oblique astigmatism that cannot be minimized using best-form principles, high magnification, and ring scotoma arising from prismatic effect in the periphery.³⁶ Aspheric lenses and lenticular designs can minimize some of these effects, but aphakic spectacle corrections remain problematic.

Contact lenses provide much better visual performance for aphakic patients. Cosmetic appearance is improved and most of the problems due to the aberrations and prismatic effects of spectacle lenses are eliminated. The magnification is closer to that of the phakic eye, making spatial adaptation and hand-eye coordination easier.³⁷ The optical advantages are somewhat offset by the difficulties of fitting these lenses due to the thick optical center and the challenge of providing sufficient oxygen to maintain corneal physiology.

Except in rare cases, most patients with clinically significant cataracts now have their natural crystalline lenses replaced with *intraocular lenses (IOLs)*. They are described as being *pseudophakic*. In general, posterior chamber IOLs that are implanted in the lens capsule after ECCE, provide better optical performance than anterior chamber lenses that are fixed in the anterior chamber angle or clipped to the iris following ICCE. The power of the IOL is determined from the keratometric measurement of corneal curvature and the axial length of the eye, as determined by ultrasound A-scan. There are many formulas for estimating the power of the IOL. One of the most widely used, the SRK II formula is

$$P = A1 - 0.9 K - 2.5 L \quad (3)$$

where P is the IOL power needed for emmetropia, $A1$ is a constant that varies according to the axial length of the eye L in millimeters, and K is the keratometer reading of central corneal curvature.³⁸ IOLs for astigmatism and bifocal IOLs have also been produced and recently designs incorporating wavefront correction have been introduced. IOL formulas have also been developed for use with patients who have undergone LASIK.³⁹ A more detailed discussion of the optics and clinical considerations of IOLs is found in Chap. 21.

12.6 BINOCULAR FACTORS

Convergence and Accommodation

In most individuals the two eyes have similar refractive errors. When looking at a distant object, the amount of defocus is therefore similar in both eyes and accommodation to view a near object results in retinal images that appear similar in size and clarity. (It is assumed that both eyes will accommodate by the same amount.) The visual system has an easier task to fuse the two retinal images into a single binocular percept.

Ideally, corresponding points in the two retinal images would be fused into a single perceived image. However, binocular fusion does not require exactly matching images. When a single point in the retina of one eye is stimulated, an area of retina in the other eye surrounding the corresponding point will give rise to a single fused percept when any point in this area is stimulated. This is called *Panum's area*. The size of Panum's area varies with location in the retina. It is very small in the macula and increases toward the periphery of the retina. Disparities between the two points are used by the visual system to determine stereoscopic depth.

Control of the extraocular muscles is neurologically linked to the process of accommodation. As accommodation increases, the eyes converge so that near objects can be seen as a single binocular percept. When the patient accommodates while wearing corrective lenses, the lines of sight pass through off-axis points of the lens and prismatic effect will alter the amount of convergence needed to maintain fusion. Myopes experience less convergence demand while hyperopes experience more; the amount of this difference from the emmetropic convergence demand can be estimated for each eye using Prentice's rule:

$$P_H = xF \quad (4)$$

where P_H is the prismatic effect in *prism diopters* along the horizontal direction, x is the distance in centimeters along the horizontal meridian of the line of sight from the optical center of the lens, and F is the power of the corrective lens in the horizontal meridian.

The foregoing discussion assumes that the ametropia is spherical. If the ametropia is astigmatic, it is open to question whether accommodation is driven to focus on one or the other principal meridian or on the circle of least confusion (spherical equivalent). Most writers assume that focus is to the spherical equivalent in the interval of Sturm.

A more extensive discussion of binocular factors is found in Chap. 13.

Anisometropia and Aniseikonia

Anisometropia is considered clinically significant when the spherical equivalent refractive error of the two eyes differs by 1.00 D or more.⁴⁰ Whether corrected or uncorrected, anisometropia can lead to significant problems with binocular vision for the patient.

Uncorrected Ametropia In uncorrected ametropia, depending on the refractive errors of the two eyes and the distance of the object of regard, one retinal image or the other may be in focus, or neither, since both eyes accommodate by the same amount. If the patient is a young child, the eye that more frequently has a clear retinal image will develop normally; however, the eye that tends to have a blurred image more of the time is likely to develop *amblyopia*. This is often the case where one eye is hyperopic and the other emmetropic, or if both eyes are hyperopic but to different degrees. If one eye is myopic and the other emmetropic, the child may learn to use the myopic eye for near objects and the emmetropic one for distant objects, in which case there is no amblyopia, but binocular fusion and stereopsis may not develop normally.

Corrected Ametropia When anisometropia is corrected, accommodation brings both retinal images of a near object into focus; however, differences in prismatic effect in off-axis gaze and in retinal image size may lead to problems with binocular vision, particularly fusion. This is particularly a problem when the patient looks down to read. Although the visual system has the capacity to fuse images with substantial horizontal disparities, it is extremely limited when disparities in the vertical meridian are encountered. In downward gaze the lines of sight pass through a spectacle lens approximately 1.0 cm below the optical center. Many patients experience visual discomfort when the vertical prismatic imbalance between the eyes is about 1 prism diopter (this unit is often written as Δ) and vertical diplopia (double vision) when it is over 2Δ . While a younger patient can avoid this problem by tipping the chin down so that the lines of sight pass through the paraxial zone of the lens in the vertical meridian, presbyopic anisometropes would need either bifocals or PALs with slab-off prism in the near portion of the lens with more minus or less plus distance power or a separate pair of single vision reading glasses.⁴⁰ Bifocal contact lenses, monovision, or reading glasses with contact lenses may avoid the problem altogether, but may cause other complications with the patient's vision.

Aniseikonia Anisometropia corrected with spectacles almost invariably results in *aniseikonia*, a condition in which binocular fusion is difficult or impossible because of the disparity in size of the retinal or cortical images from the two eyes.

The spectacle magnification M of a spectacle lens is given by the formula

$$M = \frac{1}{1 - d_v F} \frac{1}{1 - \frac{t}{n} F_1} \quad (5)$$

where d_v is the distance between the back vertex of the spectacle lens and the cornea, F is the back vertex power of the lens, t is the axial thickness of the lens, n its refractive index, and F_1 its front surface power. The first factor is the power factor, the magnification due to a thin lens of power F at distance d_v in front of the cornea and the second is the shape factor, the magnification due to an afocal lens with front surface power F_1 , index n , and thickness t . For thin lenses the equation can be written as

$$M = 1 + d_v F + \frac{t}{n} F_1 \quad (6)$$

using a thin lens approximation.⁴¹

Aniseikonia can be measured with a space eikonometer or estimated from the refractive error. The percentage difference in image size is between 1 and 1.5 percent per diopter of anisometropia and can be attributed to differences in refractive power of the eyes rather than axial length. When relative spectacle magnification is calculated on the basis of axial versus refractive ametropia, aniseikonia must arise from differences between the eyes in refractive ametropia. Size lenses, afocal lenses with magnification, can be used to determine the amount of change in spectacle magnification of the corrective lenses needed to neutralize the aniseikonia. The spectacle magnification formula can then be used to calculate the new lens parameters needed to modify the shape factor of the lens.⁴¹

Aniseikonia is relatively rare; however, cataract surgery and corneal refractive surgery often induces it particularly when only one eye is operated. Clinicians require more understanding of this clinical phenomenon than is assumed in their training. The reader is referred to Remole and Robertson's book on this subject⁴² as well as the Web site www.opticaldiagnostics.com/info/aniseikonia.html⁴³ for further information.

12.7 CONSEQUENCES FOR OPTICAL DESIGN

Designers of optical equipment must take into account that many potential users of their devices will have ametropia, either corrected or uncorrected, and perhaps presbyopia as well. While each refractive condition has its challenges, the combination can make for a great deal of discomfort and inconvenience when using optical equipment. For example, digital point-and-shoot cameras have LCD panels to display images as well as control menus and exposure information. The panel is also used to aim the camera. Most users will hold the device approximately 40 cm from the eyes—viewing the display is uncomfortable and difficult, particularly if the user is wearing bifocal or PAL spectacles. It is also important to determine whether the device is to be used monocularly or binocularly. In the latter case, the optical design must allow for convergence of the eyes and the possible change in convergence demand if the user has to accommodate. A further complication is the size of the exit pupil of the device: as patients age, their pupil size decreases and results in vignetting of the exit beam. Human visual performance factors must be considered carefully in the design of future optical devices.

12.8 REFERENCES

1. W. N. Charman, "Optics of the Eye," Chap. 1, in V. Lakshminarayanan and Enoch J. M. (eds.), *OSA Handbook of Optics*, Vol. 3, McGraw-Hill, New York, 2009.
2. World Health Organization. *Elimination of Avoidable Visual Disability due to Refractive Errors*, WHO/PBL/00.79, WHO, Geneva, 2001.

3. M. Jalie, *Ophthalmic Lenses and Dispensing*, 3rd ed., Butterworth Heinemann, Boston, 2008.
4. W. Long, "Why Is Ocular Astigmatism Regular?" *Am. J. Optom. Physiol. Opt.* **59**(6):20–522 (1982).
5. W. Charman and G. Walsh, "Variations in the Local Refractive Correction of the Eye Across its Entrance Pupil," *Opt. Vis. Sci.* **66**(1):34–40 (1989).
6. C. Sheard, "A Case of Refraction Demanding Cylinders Crossed at Oblique Axes, Together with the Theory and Practice Involved," *Ophth. Rec.* **25**:558–567 (1916).
7. T. D. Williams, "Malformation of the Optic Nerve Head," *Am. J. Optom. Physiol. Opt.* **55**(10):706–718 (1978).
8. R. B. Rabbetts, *Clinical Visual Optics*, 3rd ed., Butterworth Heinemann, Boston, 1998, pp. 288–289.
9. R. B. Rabbetts, *Clinical Visual Optics*, 3rd ed., Butterworth Heinemann, Boston, 1998, pp. 380–389.
10. R. B. Rabbetts, *Clinical Visual Optics*, 3rd ed., Butterworth Heinemann, Boston, 1998, pp. 330–350.
11. R. B. Rabbetts, *Clinical Visual Optics*, 3rd ed., Butterworth Heinemann, Boston, 1998, pp. 19–61.
12. J. M. Newman, "Analysis, Interpretation and Prescription for the Ametropias and Heterophorias," Chap. 22, in W. J. Benjamin (ed.), *Borish's Clinical Refraction*, 2nd ed., Butterworth Heinemann, Boston, 2006, pp. 1002–1009.
13. M. Jalie, *Ophthalmic Lenses and Dispensing*, 3rd ed., Butterworth Heinemann, 2008, pp. 24–27.
14. W. F. Long, "Paraxial Optics of Vision Correction," Chap. 4, in W. N. Charman (ed.), *Vision and Visual Dysfunction*, Vol. 1, *Visual Optics and Instrumentation*, MacMillan, London, 1991, p. 55.
15. M. Jalie, *The Principles of Ophthalmic Lenses*, 3rd ed., Association of Dispensing Opticians, London, 1984, pp. 413–468.
16. B. R. Chou and J. K. Hovis, "Durability of Coated CR-39 Industrial Lenses," *Optom. Vis. Sci.* **80**(10):703–707 (2003).
17. B. R. Chou, A. Gupta, and J. K. Hovis, "The Effect of Multiple Antireflective Coatings and Centre Thickness on Resistance of Polycarbonate Spectacle Lenses to Penetration by Pointed Missiles," *Optom. Vis. Sci.* **82**(11):964–969 (2005).
18. B. R. Chou and J. K. Hovis, "Effect of Multiple Antireflection Coatings on Impact Resistance of Hoya Phoenix Spectacle Lenses," *Clin. Exp. Optom.* **89**(2):86–89 (2006).
19. D. G. Pitts and B. R. Chou, "Prescription of Absorptive Lenses," Chap. 25, in W. J. Benjamin and I. M. Borish (eds.), *Borish's Clinical Refraction*, 2nd ed., W.B. Saunders Company: Cambridge, MA, 2006, pp. 1153–1187.
20. M. Jalie, *Ophthalmic Lenses and Dispensing*, 3rd ed., Butterworth Heinemann, Boston, 2008, pp. 169–195.
21. M. F. Refojo, "Chemical Composition and Properties," Chap. 2, in M. Guillon and M. Ruben (eds.), *Contact Lens Practice*, Chapman & Hall, London, 1994, pp. 29–36.
22. G. E. Lowther, "Induced Refractive Changes," Chap. 44, in M. Guillon and M. Ruben (eds.), *Contact Lens Practice*, Chapman & Hall, London, 1994.
23. L. G. Carney, "Orthokeratology," Chap 37, in M. Guillon and M. Ruben (eds.), *Contact Lens Practice*, Chapman & Hall, London, 1994.
24. M. Callender, B. R. Chou, and B. E. Robinson (eds.), *Contact Lenses and Solutions Available in Canada Univ. Waterloo Contact Lens J.* **34**(1):5–70 (2007).
25. J. Choo, K. Vuu, P. Bergenske, K. Burnham, J. Smythe, and P. Caroline, "Bacterial Populations on Silicone Hydrogel and Hydrogel Contact Lenses after Swimming in a Chlorinated Pool," *Optom. Vis. Sci.* **82**(2):134–137 (2005).
26. D. Lam and T. B. Edrington, "Contact Lenses and Aquatics," *Contact Lens Spectrum*, **21**(5):2–32 (2006).
27. Food and Drug Administration. FDA Reminds Consumers of Serious Risks of Using Decorative Contact Lenses without Consulting Eye Care Professionals. *FDA News* (Online) October 27, 2006. Accessed June 15, 2008, www.fda.gov/bbs/topics/NEWS/2006/NEW01499.html.
28. M. J. A. Port, "Cosmetic and Prosthetic Contact Lenses," Chap 20.6, in A. J. Phillips and L. Speedwell (eds.), *Contact Lenses*, 4th ed., Butterworth Heinemann, Oxford, 1997, pp. 752–754.
29. A. L. Cohen, "Diffractive Bifocal Lens Designs," *Optom. Vis. Sci.*, **70**(6):461–468 (1993).
30. J. P. G. Bergmanson and E. J. Farmer, "A Return to Primitive Practice? Radial Keratotomy Revisited," *Contact Lens Anterior Eye* **22**(1):2–10(1999).
31. P. Vinger, W. Mieler, J. Oestreicher, and M. Easterbrook, "Ruptured Globe Following Radial and Hexagonal Keratotomy," *Arch. Ophthalmol.* **114**:129–134 (1996).

32. M. J. Lee, S. M. Lee, H. J. Lee, W. R. Wee, J. H. Lee, and M. K. Kim, "The Changes of Posterior Corneal Surface and High-order Aberrations after Refractive Surgery in Moderate Myopia," *Korean J. Ophthalmol*, **21**(3): 131–136 (2007).
33. D. G. Pitts, L. L. Cameron, J. G. Jose, S. Lerman, E. Moss, S. D. Varma, S. Zigler, S. Zigman, and J. Zulich, "Optical Radiation and Cataracts," Chap. 2, in M. Waxler and V. M. Hitchins (eds.), *Optical Radiation and Visual Health*. CRC Press, Boca Raton, FL, 1986, pp. 23–25.
34. B. Noble and I. Simmons, *Complications of Cataract Surgery: A Manual*, Butterworth Heinemann, Boston, 2001, pp. 83–85.
35. H. H. Emsley, "Optics of Vision," Vol. 1, *Visual Optics*, 5th ed. Butterworths, London, 1952, p. 111.
36. T. E. Fannin and T. P. Grosvenor, *Clinical Optics*, 2nd ed, Butterworth Heinemann, Boston, 1996, pp. 328–345.
37. T. E. Fannin and T. P. Grosvenor, *Clinical Optics*, 2nd ed., Butterworth Heinemann, Boston, 1996, p. 396.
38. D. R. Sanders, J. Retzlaff, and M. C. Kraff, "Comparison of the SRK II Formula and other Second Generation Formulas," *J. Cataract Refract Surg*. **14**:136–141 (1988).
39. L. Wang, M. A. Booth, and D. D. Koch, "Comparison of Intraocular Lens Power Calculation Methods in Eyes that Have Undergone Laser-assisted in-situ Keratomileusis," *Trans. Am. Ophthalmol Soc.*, **102**:189–197 (2004).
40. T. E. Fannin and T. P. Grosvenor, *Clinical Optics*, 2nd ed., Butterworth Heinemann, Boston, 1996, pp. 294–299.
41. T. E. Fannin and T. P. Grosvenor, *Clinical Optics*, 2nd ed., Butterworth Heinemann, Boston, 1996, pp. 300–323.
42. A. Remole and K. M. Robertson, *Aniseikonia and Anisophoria*, Runestone Publishing, Waterloo, Ontario, Canada, 1996.
43. About Aniseikonia. *Optical Diagnostics*. (Online) Accessed June 25, 2008, www.opticaldiagnostics.com/info/aniseikonia.html.

This page intentionally left blank.

DO NOT DUPLICATE

BINOCULAR VISION FACTORS THAT INFLUENCE OPTICAL DESIGN

Clifton Schor

*School of Optometry
University of California
Berkeley, California*

13.1 GLOSSARY

Accommodation. Change in focal length or optical power of the eye produced by change in power of the crystalline lens as a result of contraction of the ciliary muscle. This capability decreases with age.

Ametrope. An eye with a refractive error.

Aniseikonia. Unequal perceived image sizes from the two eyes.

Baseline. Line intersecting the entrance pupils of the two eyes.

Binocular disparity. Differences in the perspective views of the two eyes.

Binocular fusion. Act or process of integrating percepts of retinal images formed in the two eyes into a single combined percept.

Binocular parallax. Angle subtended by an object point at the nodal points of the two eyes.

Binocular rivalry. Temporal alternation of perception of portions of each eye's visual field when the eyes are stimulated simultaneously with targets composed of dissimilar colors or different contour orientations.

Center of rotation. A pivot point within the eye about which the eye rotates to change direction of gaze.

Concomitant. Equal amplitude synchronous motion or rotation of the two eyes in the same direction.

Conjugate. Simultaneous motion or rotation of the two eyes in the same direction.

Convergence. Inward rotation of the two eyes.

Corresponding retinal points. Regions of the two eyes, which when stimulated result in identical perceived visual directions.

Cyclopean eye. A descriptive term used to symbolize the combined view points of the two eyes into a single location midway between them.

Cyclovergence. Unequal torsion of the two eyes.

Disconjugate. Simultaneous motion or rotation of the two eyes in opposite directions.

Divergence. Outward rotation of the two eyes.

Egocenter. Directional reference point for judging direction relative to the head from a point midway between the two eyes.

Emmetrope. An eye with no refractive error.

Entrance pupil. The image of the aperture stop formed by the portion of an optical system on the object side of the stop. The eye pupil is the aperture stop of the eye.

Extraretinal cues. Nonvisual information used in space perception.

Eye movement. A rotation of the eye about its center of rotation.

Fixation. The alignment of the fovea with an object of interest.

Focus of Expansion. The origin of velocity vectors in the optic flow field.

Frontoparallel plane. Plane that is parallel to the face and orthogonal to the primary position of gaze.

Haploopia. Perception of a single target by the two eyes.

Heterophoria. Synonymous with phoria.

Horopter. Locus of points in space whose images are formed on corresponding points of the two retinas.

Hyperopia. An optical error of the eye in which objects at infinity are imaged behind the retinal surface while the accommodative response is zero.

Midsagittal plane. Plane that is perpendicular to and bisects the baseline. The plane vertically bisecting the midline of the body.

Motion parallax. Apparent relative displacement or motion of one object or texture with respect to another that is usually produced by two successive views by a moving observer of stationary objects at different distances.

Myopia. An optical error of the eye in which objects at infinity are imaged in front of the retinal surface while the accommodative response is zero.

Nonconcomitant. Unequal amplitudes of synchronous motion or rotation of the two eyes in the same direction.

Optic flow. Pattern of retinal image movement.

Percept. That which is perceived.

Perception. The act of awareness through the senses, such as vision.

Perspective. Variations of perceived size, separation, and orientation of objects in 3-D space from a particular viewing distance and vantage point.

Phoria. An error of binocular alignment revealed when one eye is occluded.

Primary position of gaze. Position of the eye when the visual axis is directed straight ahead and is perpendicular to the frontoparallel plane.

Retinal disparity. Difference in the angles formed by two targets with the entrance pupils of the eyes.

Shear. Differential displacement during motion parallax of texture elements along an axis perpendicular to meridian of motion.

Skew movement. A vertical vergence or vertical movement of the two eyes in opposite directions.

Stereopsis. The perception of depth stimulated by binocular disparity.

Strabismus. An eye turn or misalignment of the two eyes during attempted binocular fixation.

Tilt. Amount of angular rotation in depth about an axis in the frontoparallel plane—synonymous with orientation.

Torsion. Rotation of the eye around the visual axis.

Vantage point. Position in space of the entrance pupil through which 3-D space is transformed by an optical system with a 2-D image or view plane.

Visual field. The angular region of space or field of view limited by the entrance pupil of the eye, the zone of functional retina, and occlusion structures such as the nose and orbit of the eye.

Visual plane. Any plane containing the fixation point and entrance pupils of the two eyes.

Yoked movements. Simultaneous movement or rotation of the two eyes in the same direction.

13.2 COMBINING THE IMAGES IN THE TWO EYES INTO ONE PERCEPTION OF THE VISUAL FIELD

The Visual Field

Optical designs that are intended to enhance vision usually need to consider limitations imposed by having two eyes and the visual functions that are made possible by binocular vision. When developing an optical design it is important to tailor it to specific binocular functions that you wish to enhance or at least not limit binocular vision. Binocular vision enhances visual perception in several ways. First and foremost, it expands the visual field.¹ The horizontal extent of the visual field depends largely on the placement and orientation of eyes in the head. Animals with laterally placed eyes have panoramic vision that gives them a full 360° of viewing angle. The forward placement of our eyes reduces the visual field to 190°. Eye movements let us expand our field of view. The forward placement of the eyes adds a large region of binocular overlap so that we can achieve another binocular function, stereoscopic depth perception. The region of binocular overlap is 114° and the remaining monocular portion is 37° for each eye. Each eye sees slightly more of the temporal than nasal-visual field and they also see more of the ipsilateral than the contralateral side of a binocularly viewed object.

Visual perception is not uniform throughout this region. The fovea or central region of the retina is specialized to allow us to resolve fine detail in the central 5° of the visual field. Eye movements expand the regions of available space. They allow expansion of the zone of high resolution, and accommodation expands the range of distances within which we can have high acuity and clear vision from viewing distances ranging from as near as 5 cm in the very young to optical infinity.² The peripheral retina is specialized to aid us in locomotor tasks such as walking so that we can navigate safely without colliding with obstacles. It is important that optical designs allow users to retain the necessary visual to perform tasks aided by the optical device.

Perceived Space

Space perception includes our perception of direction, distance, orientation, and shape of objects; object trajectories with respect to our location; sensing body orientation, location, and motion in space; and heading. These percepts can be derived from a 2-D image projection on the two retinas of a 3-D object space. Reconstruction of a 3-D percept from a 2-D image requires the use of information within the retinal image that is geometrically constrained by the 3-D nature of space. Three primary sources of visual or vision-related information are used for this purpose and they include monocular information, binocular information, and extraretinal or motor information. Monocular visual or retinal cues include familiarity with size and shape of objects, linear perspective and shape distortions, texture density, shading, partial image occlusion or overlap, size expansion and optic flow patterns, and motion parallax.³ Binocular information is also available, including stereoscopic depth sensed from combinations of horizontal and vertical disparity, and motion in depth.⁴ Extraretinal cues include accommodation, convergence, and gaze direction of the eyes. Many of these cues are redundant and provide ways to check their consistency, and to sense their errors that can be corrected by adaptably reweighting their emphasis or contribution to the final 3-D percept.

Monocular Cues Static monocular cues include familiarity with size and shape of objects, linear perspective and shape distortions, texture density, partial image occlusion or overlap.³ Linear perspective refers to the distortions of the retinal image that result from the projection or imaging of a 3-D object located at a finite viewing distance onto a 2-D surface such as the retina. Texture refers to a repetitive pattern of uniform size and shape such as a gravel bed. Familiarity with the size and shape of a target allows us to sense its distance and orientation. As targets approach us and change their orientation, even though the retinal image expands and changes shape, the target appears to maintain a constant perceived size and rigid shape. Retinal image shape changes are interpreted as changes in object orientation produced by rotations about three axes. Texture density can indicate which axes the object is rotated about. Rotation about the vertical axis produces tilt and causes a compression of texture along

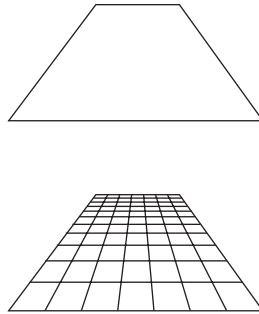


FIGURE 1 When the edges of an object are parallel, they meet at a point in the retinal image plane called the vanishing point.

the horizontal axis. Rotation about the horizontal axis produces slant and it causes a compression of texture along the vertical axis. Rotation about the z axis produces roll or torsion and causes a constant change in orientation of all texture elements. Combinations of these rotations cause distortions or foreshortening of images; however, if the distortions are decomposed into three rotations, the orientation of rigid objects can be computed. The gradient of compressed texture varies inversely with target distance. If viewing distance is known, the texture density gradient is a strong cue to the amount of slant about a given axis.

Perspective cues arising from combinations of observer's view point, object distance, and orientation can contribute to perception of distance and orientation. Perspective cues utilize edge information that is extrapolated until it intersects another extrapolated edge of the same target (Fig. 1).⁵ When the edges of an object are parallel, they meet at a point in the retinal image plane called the vanishing point. Normally we do not perceive the vanishing point directly, and it needs to be derived by extrapolation. It is assumed that the line of sight directed at the vanishing point is parallel to the edges of the surface such as a roadway that extrapolates to the same vanishing point. As with texture density gradients, perspective distortion of the retinal image increases with object proximity. If the viewing distance is known, the amount of perspective distortion is an indicator of the amount of slant or tilt about a given axis.

Depth ordering can be computed from image overlap. This is a very powerful cue and can override any of the other depth cues, including binocular disparity. A binocular version of form overlap is called DaVinci stereopsis where the overlap is compared between the two eyes' views.⁶ Because of their lateral separation, each eye sees a slightly different view of 3-D objects. Each eye perceives more of the ipsilateral side of the object than the contralateral side. Thus for a given eye, less of the background is occluded by a near object on the ipsilateral than the contralateral side. This binocular discrepancy is sufficient to provide a strong sense of depth in the absence of binocular parallax.

Kinetic Cues Kinetic monocular motion cues include size expansion, optic flow patterns, and motion parallax. Distance of targets that we approach becomes apparent from the increased velocity of radial flow of the retina (loom).⁷ Motion of texture elements across the whole image contributes to our perception of shape and location of objects in space. The relative motion of texture elements in the visual field is referred to as optic flow and this flow field can be used to segregate the image into objects at different depth planes as well as to perceive the shape or form of an object. Motion parallax is the shear of optic flow in opposite directions resulting from the translation of our view point.⁸ Movies that pan the horizon yield a strong sense of 3-D depth from the motion parallax they produce between near and far parts of 3-D objects. Two sequential views resulting from a lateral shift in the viewpoint (motion parallax) are analogous to two separate and simultaneous views from two eyes with separate viewpoints (binocular disparity). Unlike horizontal disparity, the shear

between sequential views by itself is ambiguous, but once we know which way our head is translating, we can correctly identify the direction or sign of depth.⁹ The advantage of motion parallax over binocular disparity is that we can translate in any meridian, horizontal or vertical, and get depth information, while stereo only works for horizontal disparities.⁸

Motion perception helps us to determine where we are in space over time. Probably the most obvious application is to navigation and heading judgments (Fig. 2). Where are we headed, and can

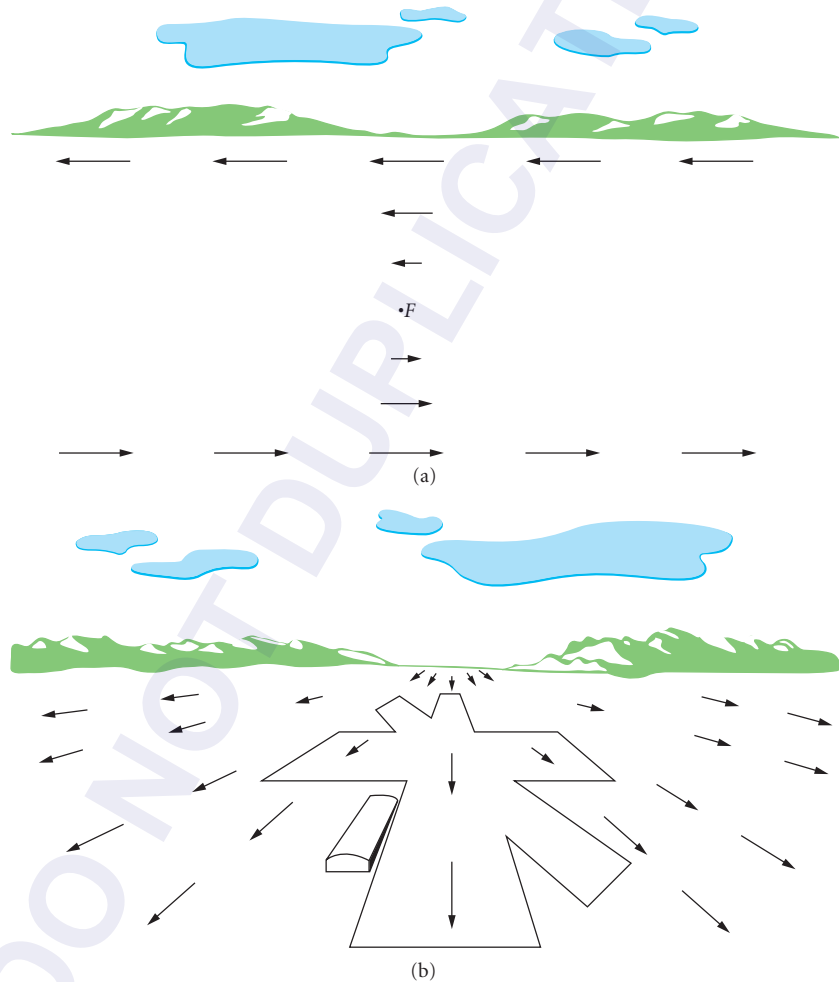


FIGURE 2 (a) Motion Parallax. Assume that an observer moving toward the left fixates a point at F . Objects nearer than F will appear to move in a direction opposite to that of the movement of the observer; objects farther away than F will appear to move in the same direction as the observer. The length of the arrows signifies that the apparent velocity of the optic flow is directly related to the distance of objects from the fixation point. (b) Motion perspective. The optical flow in the visual field as an observer moves forward. The view is that as seen from an airplane in level flight. The direction of apparent movement in the terrain below is signified by the direction of the motion vectors (arrows); speed of apparent motion is indicated by the length of the motion vectors. The expansion point in the distance from which motion vectors originate is the heading direction.

we navigate to a certain point in space? Generally we perceive a focus of expansion in the direction we are heading as long as our eyes and head remain fixed.^{7,10–13} These types of judgments occur whenever we walk, ride a bicycle, or operate a car. Navigation also involves tasks to intercept or avoid moving objects. For example, you might want to steer your car around another car or avoid a pedestrian walking across the crosswalk. Or you may want to intercept something such as catching or striking a ball. Another navigation task is anticipation of time of arrival or contact. This activity is essential when pulling up to a stop sign or traffic light. If you stop too late you enter the intersection and risk a collision. These judgments are also made while walking down or running up a flight of stairs. Precise knowledge of time to contact is essential to keep from stumbling. Time to contact can be estimated from the distance to contact divided by the speed of approach. It can also be estimated from the angular subtense of an object divided by its rate of expansion.¹⁴ The task is optimal when looking at a surface that is in the frontoparallel plane, and generally we are very good at it. It looks like we use the rate of expansion to predict time of contact since we make unreliable judgments of our own velocity.¹⁵

Another class of motion perception is self- or ego-motion.¹⁶ It is related to the vestibular sense of motion and balance and it tells us when we are moving. Self-motion responds to vestibular stimulation as well as visual motion produced by translational and rotational optic flow of large fields. Finally, we can use binocular motion information to judge the trajectory of a moving object and determine if it will strike us in the head or not. This motion in depth results from a comparison of the motion in one eye to the other. The visual system computes the direction of motion from the amplitude ratio and relative direction of horizontal retinal image motion (Fig. 3).¹⁷ If horizontal motion in the two

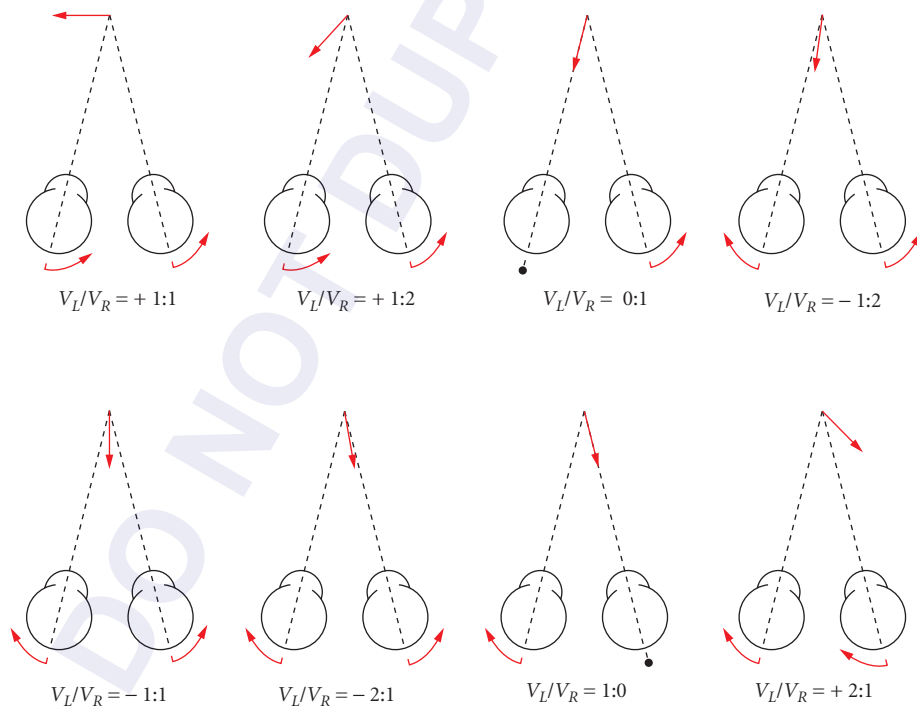


FIGURE 3 Motion in depth. Relative velocities of left and right retinal images for different target trajectories. When the target moves along a line passing between the eyes, its retinal images move in opposite directions in the two eyes; when the target moves along a line passing wide of the head, the retinal images move in the same direction, but with different velocities. The ratio (V_L/V_R) of left- and right-eye image velocities provides an unequivocal indication of the direction of motion in depth relative to the head.

eyes is equal in amplitude and direction, it corresponds to an object moving in the frontoparallel plane. If the horizontal retinal image motion is equal and opposite, the object is moving toward us in the midsagittal plane and will strike us between the eyes. If the motion has any disconjugacy at all, it will strike us. A miss is indicated by yoked, albeit noncomitant, motion.

Binocular Cues and Extraretinal Information from Eye Movements Under binocular viewing conditions, we perceive a single view of the world as though seen by a single cyclopean eye, even though the two eyes receive slightly different retinal images. The binocular percept is the average of monocularly sensed shapes and directions (allelotopia). The benefit of this binocular combination is to allow us to sense objects as single with small amounts of binocular disparity so that we can interpret depth from the stereoscopic sense. Interpretation of the effects of prism and magnification upon perceived direction through an instrument needs to consider how the visual directions of the two eyes are combined. If we only had one eye, direction could be judged from the nodal point of the eye, a site where viewing angle in space equals visual angle in the eye, assuming the nodal point is close to the radial center of the retina. However, two eyes present a problem for a system that operates as though it has only a single cyclopean eye. The two eyes have view points separated by approximately 6.5 cm. When the two eyes converge accurately on a near target placed along the midsagittal plane, the target appears straight ahead of the nose, even when one eye is occluded. In order for perceived egocentric direction to be the same when either eye views the near target monocularly, there needs to be a common reference point for judging direction. This reference point is called the cyclopean locus or egocenter, and is located midway on the interocular axis. The egocenter is the percept of a reference point for judging visual direction with either eye alone or under binocular viewing conditions.

Perceived Direction

Direction and distance can be described in polar coordinates as the angle and magnitude of a vector originating at the egocenter. For targets imaged on corresponding retinal points, this vector is determined by the location of the retinal image and by the direction of gaze that is determined by the average position of the two eyes (conjugate eye position). The angle the two retinal images form with the visual axes is added to the conjugate rotational vector component of binocular eye position (the average of right- and left-eye position). This combination yields the perceived egocentric direction. Convergence of the eyes, which results from disconjugate eye movements, has no influence on perceived egocentric direction. Thus, when the two eyes fixate near objects to the left or right of the midline in asymmetric convergence, only the conjugate component of the two eyes' positions contributes to perceived direction. These facets of egocentric direction were summarized by Hering¹⁸ as five laws of visual direction, and they have been restated by Howard.¹⁹ The laws are mainly concerned with targets imaged on corresponding retinal regions (i.e., targets on the horopter).

We perceive space with two eyes as though they were merged into a single cyclopean eye. This merger is made possible by a sensory linkage between the two eyes that is facilitated by the anatomical superposition or combination of homologous regions of the two retinas in the visual cortex. The binocular overlap of the visual fields is very specific. Unique pairs of retinal regions in the two eyes (corresponding points) must receive images of the same object so that these objects can be perceived as single and at the same depth as the point of fixation. This requires that retinal images must be precisely aligned by the oculomotor system with corresponding retinal regions in the two eyes. As described in Sec. 13.8, binocular alignment is achieved by yoked movements of the eyes in the same direction (version) and by movements of the eyes in opposite direction (vergence). Slight misalignment of similar images from corresponding points is interpreted as depth. Three-dimensional space can be derived geometrically by comparing the small differences between the two retinal images that result from the slightly different vantage points of the two eyes caused by their 6.5-cm separation. These disparities are described as horizontal, vertical, and torsional as well as distortion or shear differences between the two images. The disparities result from surface shape,

depth, and orientation with respect to the observer as well as direction and orientation (torsion) of the observer's eyes.²⁰ These disparities are used to judge the layout of 3-D space and to sense the solidness or curvature of surfaces. Disparities are also used to break through camouflage in images such as seen in tree foliage.

Binocular Visual Direction—Corresponding Retinal Points and the Horopter Hering¹⁸ defined binocular correspondence by retinal locations in the two eyes, which when stimulated, resulted in a percept in identical visual directions. For a fixed angle of convergence, projections of corresponding points along the equator and midline of the eye converge upon real points in space. In other cases, such as oblique eccentric locations of the retina, corresponding points have visual directions that do not intersect in real space. The horopter is the locus in space of real objects or points whose images can be formed on corresponding retinal points. It serves as a reference throughout the visual field for the same depth or disparity as at the fixation point. To appreciate the shape of the horopter, consider a theoretical case in which corresponding points are defined as homologous locations on the two retinas. Thus corresponding points are equidistant from their respective foveas. Consider binocular matches between the horizontal meridians or equators of the two retinas. Under this circumstance, the visual directions of corresponding points intersect in space at real points that define the longitudinal horopter. This theoretical horopter is a circle whose points will be imaged at equal eccentricities from the two foveas on corresponding points except for the small arc of the circle that lies between the two eyes.²² While the theoretical horopter is always a circle, its radius of curvature increases with viewing distance. This means that its curvature decreases as viewing distance increases. In the limit, the theoretical horopter is a straight line at infinity that is parallel to the interocular axis. Thus a surface representing zero disparity has many different shapes depending on the viewing distance. The consequence of this spatial variation in horopter curvature is that the spatial pattern of horizontal disparities is insufficient information to specify depth magnitude or even depth ordering or surface shape. It is essential to know viewing distance to interpret surface shape and orientation from depth related retinal image disparity. Viewing distance could be obtained from extraretinal information such as the convergence state of the eyes, or from retinal information in the form of vertical disparity as described further.

The empirical or measured horopter differs from the theoretical horopter in two ways.²² It can be tilted about a vertical axis and its curvature can be flatter or steeper than the Vieth-Muller circle. This is a circle that passes through the fixation point in the midsagittal plane and the entrance pupils of the two eyes. Points along this circle are imaged at equal retinal eccentricities from the foveas of the two eyes. The tilt of the empirical horopter can be the result of a horizontal magnification of the retinal image in one eye. Image points of a frontoparallel surface subtend larger angles in the magnified images. A similar effect occurs when no magnifier is worn and the fixation plane is tilted toward one eye. The tilt causes a larger image to be formed in one eye than the other. Thus magnification of a horizontal row of vertical rods in the frontoparallel plane causes the plane of the rods to appear tilted to face the eye with the magnifier.

How does the visual system distinguish between a surface fixated in eccentric gaze that is in the frontoparallel plane (parallel to the face), but tilted with respect to the Vieth-Muller circle, and a plane fixated in forward gaze, that is tilted toward one eye? Both of these planes project identical patterns of horizontal retinal image disparity, but they have very different physical tilts in space. In order to distinguish between them, the visual system needs to know the horizontal gaze eccentricity. One way is to register the extraretinal eye position signal, and the other is to compute gaze eccentricity from the vertical disparity gradient associated with eccentrically located targets. It appears that the latter method is possible, since tilt of a frontoparallel plane can be induced by a vertical magnifier before one eye (the induced effect). The effect is opposite to the tilt produced by placing a horizontal magnifier before the same eye (the geometric effect). If both a horizontal and vertical magnifier are placed before one eye, in the form of an overall magnifier, no tilt of the plane occurs. Interestingly, in eccentric gaze, the proximity to one eye causes an overall magnification of the image in the nearer eye, and this could be sufficient to allow observers to make accurate frontoparallel settings in eccentric gaze (Fig. 4). Comparison of vertical and horizontal magnification seems to be sufficient to disambiguate tilt from eccentric viewing.

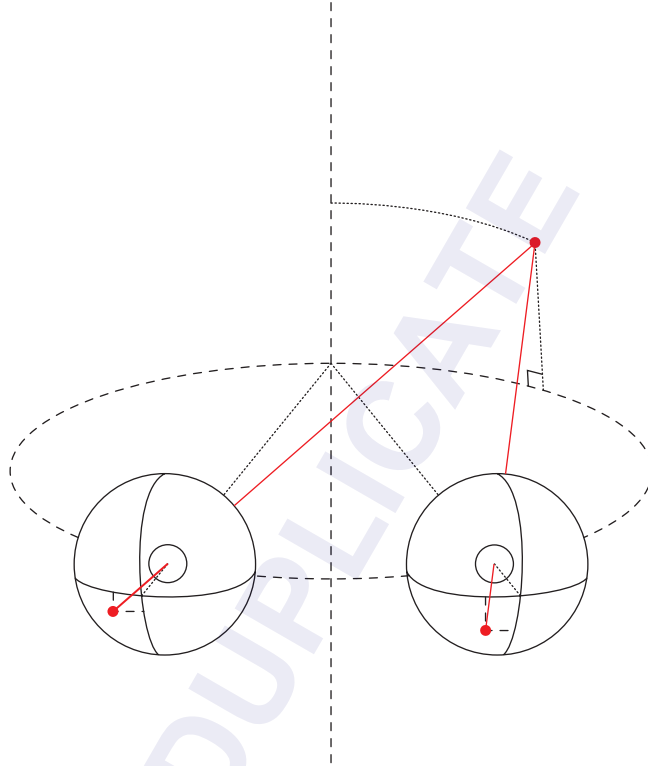


FIGURE 4 The horopter. The theoretical horopter or Vieth-Muller circle passes through the fixation point and two entrance pupils of the eyes. The theoretical vertical horopter is a vertical line passing through the fixation point in the midsagittal plane. Fixation on a tertiary point in space produces vertical disparities because of greater proximity to the ipsilateral eye. (Reprinted by permission from Tyler and Scott.²¹)

Flattening of the horopter from a circle to an ellipse results from nonuniform magnification of one retinal image. If the empirical horopter is flatter than the theoretical horopter, corresponding retinal points are more distant from the fovea on the nasal than temporal hemiretina. Curvature changes in the horopter can be produced with nonuniform magnifiers such as prisms (Fig. 5). A prism magnifies more at its apex than its base. If base-out prism is placed before the two eyes, the right half of the image is magnified more in the left eye than in the right eye and the left half of the image is magnified more in the right eye than in the left eye. This causes flat surfaces to appear more concave and the horopter to be less curved or more convex.

The Vertical Horopter The theoretical vertical point horopter for a finite viewing distance is limited by the locus of points in space where visual directions from corresponding points will intersect real objects.²¹ These points are described by a vertical line in the midsagittal plane that passes through the Vieth-Muller circle. Eccentric object points in tertiary gaze (points with both azimuth and elevation) lie closer to one eye than the other eye. Because they are imaged at different vertical eccentricities from the two foveas, tertiary object points cannot be imaged on theoretically corresponding retinal points. However, all object points at an infinite viewing distance can be imaged on corresponding retinal regions and at this infinite viewing distance, the vertical horopter becomes a plane.

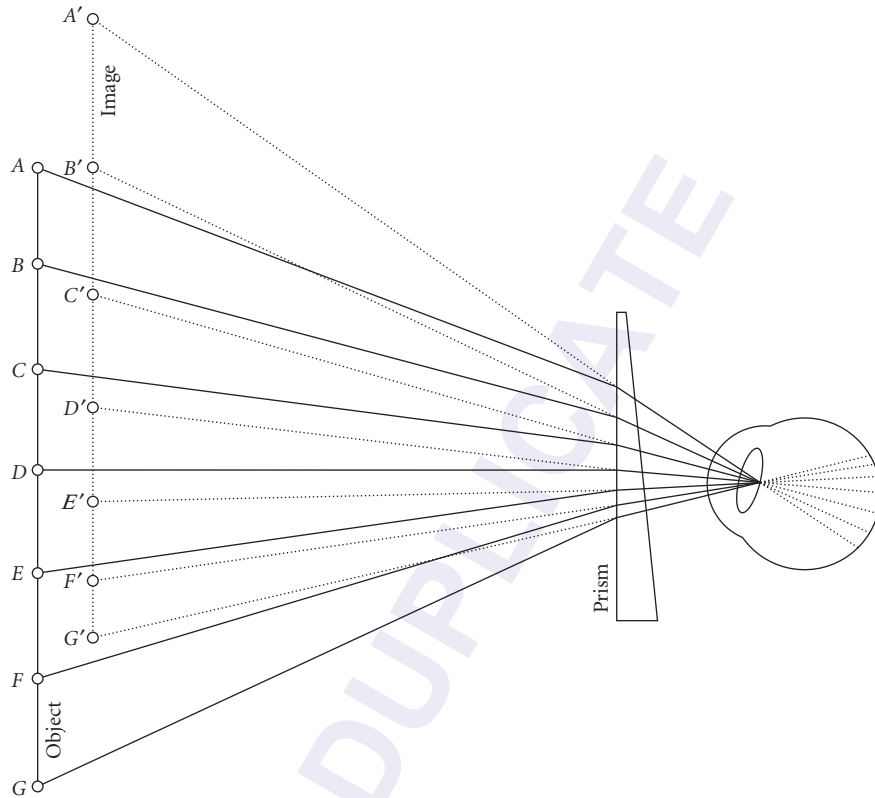


FIGURE 5 Figure showing nonuniform magnification of a prism. Disparateness of retinal images, producing stereopsis.

The empirical vertical horopter is declinated (top slanted away from the observer) in comparison to the theoretical horopter. Helmholtz²³ reasoned that this was because of a horizontal shear of the two retinal images which causes a real vertical plane to appear inclined toward the observer. Optical infinity is the only viewing distance that the vertical horopter becomes a plane. It is always a vertical line at finite viewing distances. Targets that lie away from the midsagittal plane always subtend a vertical disparity due to their unequal proximity and retinal image magnification in the two eyes. The pattern of vertical disparity varies systematically with viewing distance and eccentricity from the midsagittal plane. For a given target height, vertical disparity increases with horizontal eccentricity from the midsagittal plane and decreases with viewing distance. Thus a horizontally extended target of constant height will produce a vertical disparity gradient that increases with eccentricity. The gradient will be greater at near than far viewing distances. It is possible to estimate viewing distance from the vertical disparity gradient or from vertical disparity at a given point if target eccentricity is known. This could provide a retinal source of information about viewing distance to allow the visual system to scale disparity to depth and to compute depth ordering. Several investigators have shown that modifying vertical disparity can influence the magnitude of depth (scaling) and surface slant, such that vertical disparity is a useful source of information about viewing distance that can be used to scale disparity and determine depth ordering.

Stereopsis

Three independent variables involved in the calculation of stereodepth are retinal image disparity, viewing distance, and the separation in space of the two viewpoints (i.e., the baseline or interpupillary distance) (Fig. 6). In stereopsis, the relationship between the linear depth interval between two objects and the retinal image disparity that they subtend is approximated by the following expression

$$\Delta d = \frac{\eta \times d^2}{2a}$$

where η is retinal image disparity in radians, d is viewing distance, $2a$ is the interpupillary distance and Δd is the linear depth interval. $2a$, d , and Δd are all expressed in the same units (e.g., meters). The formula implies that in order to perceive depth in units of absolute distance (e.g., meters), the visual system utilizes information about the interpupillary distance and the viewing distance. Viewing distance could be sensed from the angle of convergence²⁴ or from other retinal cues such as oblique or vertical disparities. These disparities occur naturally with targets in tertiary directions from the point of fixation.²⁵⁻³⁰

The equation illustrates that for a fixed retinal image disparity, the corresponding linear depth interval increases with the square of viewing distance and that viewing distance is used to scale the

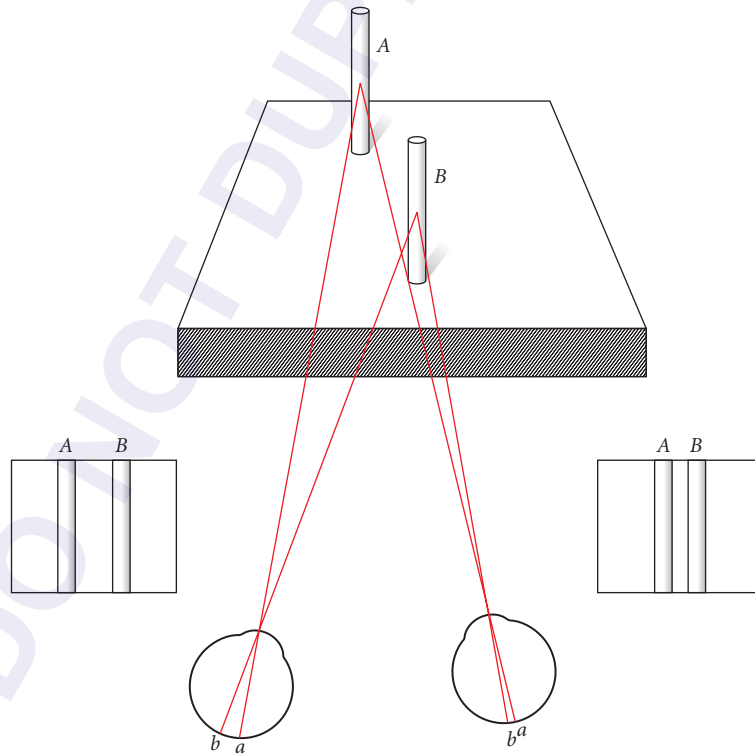


FIGURE 6 The differences in the perspective views of the two eyes produce binocular that can be used to perceive depth.

horizontal disparity into a linear depth interval. When objects are viewed through base-out prisms that stimulate additional convergence, perceived depth should be reduced by underestimates of viewing distance. Furthermore, the pattern of zero retinal image disparities described by the curvature of the longitudinal horopter varies with viewing distance. It can be concave at near distances and convex at far distances in the same observer.²² Thus, without distance information, the pattern of retinal image disparities across the visual field is insufficient to sense either depth ordering (surface curvature) or depth magnitude.²⁵ Similarly, the same pattern of horizontal disparity can correspond to different slants about a vertical axis presented at various horizontal gaze eccentricities.²² Convergence distance and direction of gaze are important sources of information used to interpret slant from disparity fields associated with slanting surfaces.³¹ Clearly, stereodepth perception is much more than a disparity map of the visual field.

Binocular Fusion and Suppression

Even though there is a fairly precise point-to-point correspondence between the two eyes for determining depth in the fixation plane, images in the two eyes can be combined into a single percept when they are anywhere within a small range or retinal area around corresponding retinal points. Thus a point in one eye can be combined perceptually with a point imaged within a small area around its corresponding retinal location in the other eye. These ranges are referred to as Panum's fusional area (PFA) and they serve as a buffer zone to eliminate diplopia for small disparities near the horopter. PFA allows for the persistence of single binocular vision in the presence of constant changes in retinal image disparity caused by various oculomotor disturbances. For example, considerable errors of binocular alignment (>15 arc min) may occur during eye tracking of dynamic depth produced either by object motion or by head and body movements.³² Stereopsis could exist without singleness but the double images near the fixation plane would be a distraction. The depth of focus of the human eye serves a similar function. Objects that are nearly conjugate to the retina appear as clear as objects focused precisely on the retina. The buffer for the optics of the eye is much larger than the buffer for binocular fusion. The depth of focus of the eye is approximately $0.75 D$. Panum's area, expressed in equivalent units, is only 0.08 meter angles or approximately one tenth the magnitude of the depth of focus. Thus we are more tolerant of focus errors than we are of convergence errors.

Fusion is only possible when images have similar size, shape, and contrast polarity. This similarity ensures that we only combine images that belong to the same object in space. Natural scenes contain many objects that are at a wide range of distances from the plane of fixation. Some information is coherent, such as images formed within Panum's fusional areas. Some information is fragmented, such as partially occluded regions of space resulting in visibility to only one eye. Finally, some information is uncorrelated because it is either ambiguous or in conflict with other information, such as the superposition of separate diplopic images arising from objects seen by both eyes behind or in front of the plane of fixation. One objective of the visual system is to preserve as much information from all three sources as possible to make inferences about objective space without introducing ambiguity or confusion of space perception. In some circumstances conflicts between the two eyes are so great that conflicting percepts are seen alternately every 4 s (binocular rivalry suppression) or in some cases one image is permanently suppressed such as when you look through a telescope with one eye while the other remains open. The view through the telescope tends to dominate over the view of the background seen by the other eye. For example, the two ocular images may have unequal clarity or blur such as in asymmetric convergence, or large unfusible disparities originating from targets behind or in front of the fixation plane may appear overlapped with other large diplopic images. Fortunately, when dissimilar images are formed within a fusible range of disparity that perception of the conflicting images is suppressed.

Four classes of stimuli evoke what appear to be different mechanisms of interocular suppression. The first is unequal contrast or blur of the two retinal images which causes interocular blur suppression. The second is physiologically diplopic images of targets in front or behind the singleness horopter which result in suspension of one of the redundant images.³³ The third is targets of different shape

presented in identical visual directions. Different size and shape result an alternating appearance of the two images referred to as either binocular retinal rivalry or percept rivalry suppression. The fourth is partial occlusion that obstructs the view of one eye such that the portions of the background are only seen by the unoccluded eye and the overlapping region of the occluder that is imaged in the other eye is permanently suppressed.

13.3 DISTORTION OF SPACE BY MONOCULAR MAGNIFICATION

The magnification and translation of images caused by optical aids will distort certain cues used for space perception while leaving other cues unaffected. This will result in cue misrepresentation of space as well as cue conflicts that will produce errors of several features of space perception. These include percepts of direction, distance, trajectory or path of external moving targets, heading and self-motion estimates, as well as shape and orientation (curvature, slant, and tilt) of external objects. Fortunately, if the optical distortions are long standing, the visual system can adapt and make corrections to their contributions to space perception.

Magnification and Perspective Distortion

Most optical systems magnify or minify images, and this can produce errors in perceived direction. When objects are viewed through a point other than the optical center of a lens, they appear displaced from their true direction. This is referred to as a prismatic displacement. Magnification will influence most of the monocular cues mentioned above except overlap. Magnification will produce errors in perceived distance and produce conflicts with cues of texture density and linear perspective. For example, the retinal image of a circle lying on a horizontal ground plane has an elliptical shape when viewed at a remote distance, and a more circular shape when viewed from a shorter distance. When a remote circle is viewed through a telescope, the uniformly magnified image appears to be closer and squashed in the vertical direction to form an ellipse (shape distortion). Uniform magnification also distorts texture-spacing gradients such that the ground plane containing the texture is perceived as inclined or tilted upward toward vertical. Normally, the change in texture gradient is greatest for near objects (Fig. 7). When a distant object is magnified, it appears nearer but its low texture density gradient is consistent with a plane inclined toward the vertical. Magnification also affects slant derived from perspective cues in the same way. For this reason, a tiled rooftop appears to have a steeper pitch when viewed through binoculars or a telescope, the pitcher-to-batter distance appears reduced in telephoto view from the outfield, and the spectators may appear larger than the players in front of them in telephoto lens shots.

Magnification and Stereopsis

Magnification increases both retinal image disparity and image size. The increased disparity stimulates greater stereo depth while increased size stimulates reduced perceived viewing distance.³⁴ Greater disparity will increase perceived depth while reduced perceived viewing distance will scale depth to a smaller value. Because stereoscopic depth from disparity varies with the square of viewing distance, the perceptual reduction of viewing distance will dominate the influence of magnification on stereoscopic depth. Thus, perceived depth intervals sensed stereoscopically from retinal image disparity appears smaller with magnified images because depth derived from binocular retinal image disparity is scaled by the square of perceived viewing distance. Objects appear as flat surfaces with relative depth but without thickness. This perceptual distortion is referred to as cardboarding. Most binocular optical systems compensate for this reduction in stereopsis by increasing the baseline or separation of the two ocular objectives. For example, the prism design of binoculars folds the optical path into

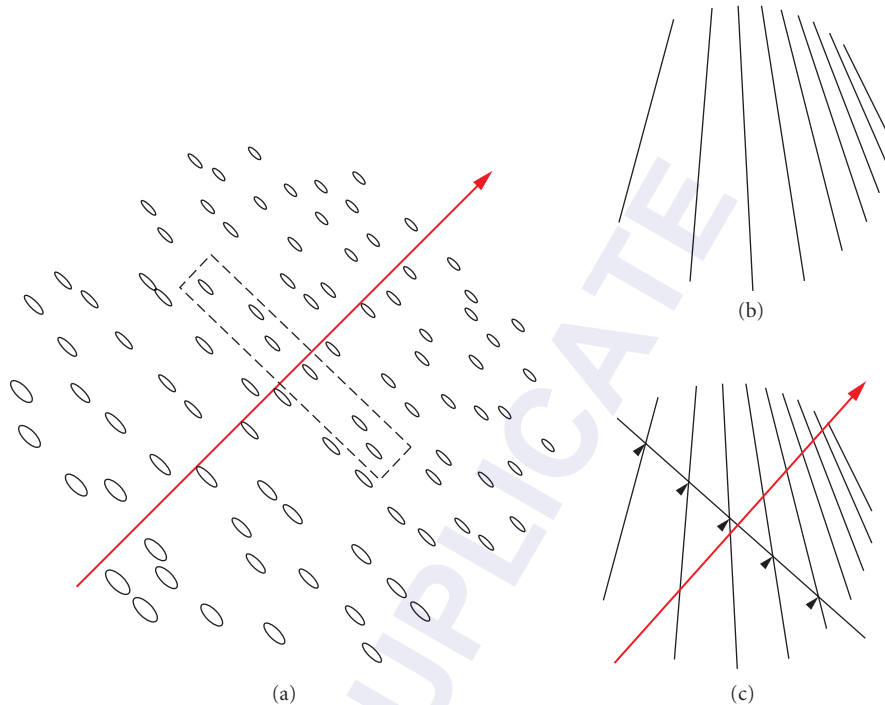


FIGURE 7 Texture gradients. The tilt of a surface is the direction in which it is slanted away from the gaze normal of the observer. (a) If the surface bears a uniform texture, the projection of the axis of tilt in the image indicates the direction in which the local density of the textures varies most, or, equivalently, it is perpendicular to the direction that the texture elements are most uniformly distributed (dotted rectangle). Either technique can be used to recover the tilt axis, as illustrated. Interestingly, however, the tilt axis in situations like (b) can probably be recovered most accurately using the second method, i.e., searching for the line that is intersected by the perspective lines at equal intervals. This method is illustrated in (c). (Reprinted by permission from Stevens.³⁵)

a compact package and also expands the distance between the objective lenses creating an expanded interpupillary distance (telestereoscope). This increase of effective separation between the two eyes' views exaggerates disparity and objects can appear in hyperstereoscopic depth if it exceeds the depth reduction caused by perceived reductions in viewing distance. The depth can be predicted from the formula in the preceding section that computes disparity from viewing distance, interpupillary distance, and physical linear depth interval.

Magnification and Motion Parallax—Heading

Magnification will produce errors in percepts derived from motion cues. Motion will be exaggerated by magnification and perceived distance will be shortened. As with binocular disparity, depth associated with motion parallax will be underestimated by the influence of magnification on perceived distance. Heading judgments in which gaze is directed away from the heading direction will be in error by the amount the angle of eccentric gaze is magnified. Thus if your car is headed north and your gaze is shifted 5° to the left, heading errors will be the percentage the 5° gaze shift is magnified by the optical system and you will sense a path that is northeast. Heading judgments that are constantly

changing, such as along a curved path, might be affected in a similar way, if the gaze lags behind the changing path of the vehicle. In this case, changes of perceived heading along the curved path would be exaggerated by the magnifier. Unequal magnification of the two retinal images in anisometropia will cause errors in sensing the direction of motion in depth by changing the amplitude ratio of the two moving images.

Bifocal Jump

Most bifocal corrections have separate optical centers for the distance lens and the near addition. One consequence is that objects appear to change their direction or jump as the visual axis crosses the top line of the bifocal segment (Fig. 8). This bifocal jump occurs because prism is introduced as the eye rotates downward behind the lens away from the center of the distance correction and enters the bifocal segment. When the visual axis enters the optical zone of the bifocal, a new prismatic displacement is produced by the displacement of the visual axis from the optical center of the bifocal addition. All bifocal additions are positive so that the prismatic effect is always the same; objects in the upper part of the bifocal appear displaced upward compared to views just above the bifocal. This results in bifocal jump. A consequence of bifocal jump is that a region of space seen near the upper edge of the bifocal is not visible. It has been displaced out of the field of view by the vertical prismatic effect of the bifocal. The invisible or missing part of the visual field equals the prismatic jump which can be calculated from Prentice's rule [distance (cm) between the upper edge and center of the bifocal times the power of the add]. Thus for a 2.5-D bifocal add with a top to center distance of 0.8 cm, approximately 1° of the visual field is obscured by bifocal jump. This loss may seem minor; however, it is very conspicuous when paying attention to the ground plane while walking. There is a clear horizontal line of discontinuity that is exaggerated by the unequal speed of optic flow above and below the bifocal boundary. Optic flow is exaggerated by the magnification of the positive bifocal

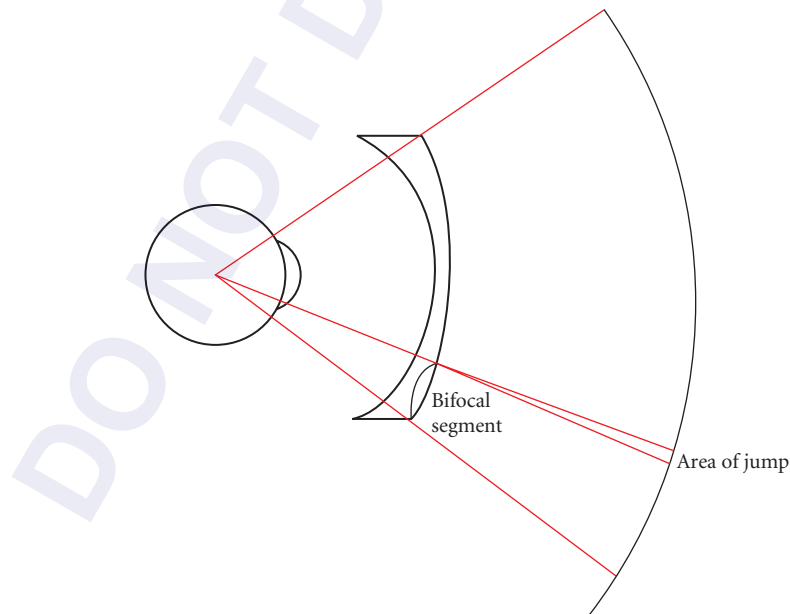


FIGURE 8 Jump figure. The missing region of the visual field occurs at the top edge of the bifocal segment as a result of prism jump.

segment. It can produce errors of estimated walking speed and some confusion if you attend to the ground plane in the lower field while navigating toward objects in the upper field. This can be a difficult problem when performing such tasks as walking down steps. It is best dealt with by maintaining the gaze well above the bifocal segment while walking, which means that the user is not looking at the ground plane in the immediate vicinity of the feet. Some bifocal designs eliminate jump by making the pole of the bifocal and distance portion of the lens concentric. The bifocal line starts at the optical centers of both lenses. This correction always produces some vertical displacement of the visual field since the eyes rarely look through the optical center of the lens; however, there is no jump.

Discrepant Views of Objects and Their Images

Translation errors in viewing a scene can lead to errors in space perception. A common problem related to this topic is the interpretation of distance and orientation of objects seen on video displays or in photographs. The problem arises from discrepancies between the camera's vantage point with respect to a 3-D scene, and the station point of the observer who is viewing the 2-D projected image after it has been transformed by the optics of the camera lens. Errors of observer distance and tilt of the picture plane contribute to perceptual errors of perspective distortion, magnification, and texture density gradients expected from the station point. The station point represents the position of the observer with respect to the view plane or projected images. To see the picture in the correct perspective, the station point should be the same as the location of the camera lens with respect to the film plane. Two violations of these requirements are errors in viewing distance and translation errors (i.e., horizontal or vertical displacements). Tilt of the view screen can be decomposed into errors of distance and translation errors. In the case of the photograph, the correct distance for the station point is the camera lens focal length multiplied by the enlargement scaling factor. Thus, when a 35-mm negative taken through a 55-mm lens is enlarged 7 times to a 10-in print, it should be viewed at a distance of 38.5 cm or 15 in (7×5.5 cm).

13.4 DISTORTION OF SPACE PERCEPTION FROM INTEROCULAR ANISO-MAGNIFICATION (UNEQUAL BINOCULAR MAGNIFICATION)

Optical devices that magnify images unequally or translate images unequally in the two eyes produce errors of stereoscopic depth and perceived direction as well as eye alignment errors that the oculomotor system is not accustomed to correcting. If these sensory and motor errors are longstanding, the visual system adapts perceptual and motor responses to restore normal visual performance. Patients who receive optical corrections that have unequal binocular magnification will adapt space percepts in 10 days to 2 weeks.

Lenses and Prisms Magnification produced by lenses is uniform compared to the nonuniform magnification produced by prisms, where the magnification increases toward the apex of the prism. Uniform magnification differences between the horizontal meridians of the two eyes produce errors in perceived surface tilt about a vertical axis (geometric effect). A frontoparallel surface appears to face the eye with the more magnified image. This effect is offset by vertical magnification of one ocular image that causes a frontoparallel surface to appear to face the eye with the less magnified image. Thus perceptual errors of surface tilt occur mainly when magnification is greater in the horizontal or vertical meridian. These meridional magnifiers are produced by cylindrical corrections for astigmatism. Nonuniform magnification caused by horizontal prism such as base-out or base-in (Fig. 5), cause surfaces to appear more concave or convex, respectively. Finally, cylindrical corrections for astigmatism cause scissor or rotational distortion of line orientations (Fig. 9). When the axes of astigmatism are not parallel in the two eyes, cyclodisparities are introduced. These disparities

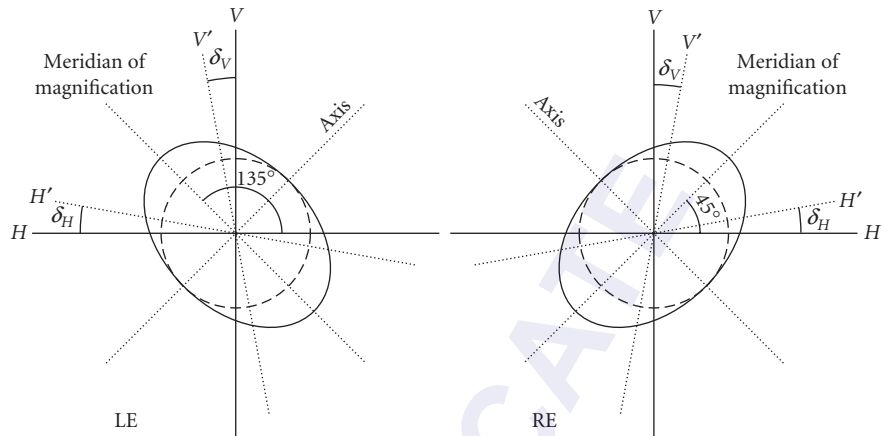


FIGURE 9 Magnification ellipse. Meridional magnification of images in the two eyes in oblique meridians cause “scissors” or rotary deviations of images of vertical (and horizontal) lines and affect stereoscopic spatial localization in a characteristic manner. (Reprinted from Ogle and Boder.³⁶)

produce slant errors and cause surfaces to appear inclined or declinated. These spatial depth distortions are adapted to readily so that space appears veridical; however, some people are unable to cope with the size differences of the two eyes and suffer from an anomaly referred to as aniseikonia. While aniseikonic errors are tolerated by many individuals, they are subject of spatial distortions which, in given situations, may cause meaningful problems.

Aniseikonia

Unequal size of the two retinal images in anisometropia can precipitate a pathological mismatch in perceived size of the two ocular images (aniseikonia). Aniseikonia is defined as a relative difference in perceived size and/or shape of the two ocular images. Most difficulty with aniseikonia occurs when anisometropia is corrected optically with spectacles. Patients experience a broad range of symptoms including spatial disorientation, depth distortion, diplopia, vertigo, and asthenopia. These symptoms result from several disturbances including spatial distortion and interference with binocular sensory and motor fusion.

Two optical distortions produced by spectacle corrections of anisometropia are called axial and lateral aniseikonia. Axial aniseikonia describes size distortions produced by differences in the magnification at the optical centers of ophthalmic spectacle corrections for anisometropia. The magnification difference produced by the spectacle correction depends in part on the origin of the refractive error. Retinal images in the uncorrected axial and refractive ametropias are not the same. In axial myopia, the uncorrected retinal image of the elongated eye is larger than that of an emmetropic eye of similar refractive power. An uncorrected refractive myope has a blurred image of similar size to an emmetropic eye with the same axial length. Knapp’s law states that image size in an axial ametropic eye can be changed to equal the image size in an emmetropic eye having the same refractive power by placing the ophthalmic correction at the anterior focal plane of the eye. The anterior focal plane of the eye is usually considered to correspond approximately to the spectacle plane. However, this is only a very general approximation. The anterior focal point will vary with power of the eye and in anisometropia, it will be unequal for the two eyes. Assuming the anterior focal point is the spectacle plane, then the abnormally large image of an uncorrected axial-myopic eye can be reduced to near normal with a spectacle correction, whereas a contact lens correction would leave the retinal image of the myopic eye enlarged. If the ametropia is refractive in nature (i.e., not axial), then the uncorrected

retinal images are equal in size and a contact lens correction will produce less change in image size than a spectacle correction. If the refractive myope is corrected with spectacles, the retinal image becomes smaller than that in an emmetropic eye of similar refractive power. This change in image size of refractive anisometropes corrected with spectacle lenses is one of several factors leading to ocular discomfort in aniseikonia.

Paradoxically, spectacle correction of axial myopes often produces too drastic a reduction in ocular image size. The ocular image of the axially myopic eye, corrected with spectacle lenses, appears smaller than that of the less ametropic eye in anisometropia.^{37–42} As mentioned earlier, ocular image size depends not only upon retinal image size but also upon the neural factor.²² For example, in cases of axial myopia, the elongation of the eye produces stretching of the retina and pigment epithelium. Accordingly, retinal images of normal size occupy a smaller proportion of the stretched retina than they would in a retina with normal retinal element density. Indeed, anisometropes with axial myopia report minification of the ocular image in their more myopic eye after the retinal images have been matched in size.^{41,43,44} Thus, in the uncorrected state, myopic anisometropic individuals have optical magnification and neural minification. As a result, a spectacle correction would overcompensate for differential size of ocular images.

Surprisingly, not all patients with a significant amount of aniseikonia (4 to 6%) complain of symptoms. This could be due in part to suppression of one ocular image. However, normally we are able to fuse large size differences up to 15 percent of extended targets.⁴⁵ Sinusoidal gratings in the range from 1 to 5 cpd can be fused with grating differing in spatial frequency by 20 to 40 percent^{46–48} and narrow band pass filtered bars can be fused with 100- to 400-percent difference in interocular size.^{49,50} These enormous size differences demonstrate the remarkable ability of the visual system to fuse large size differences in axial aniseikonia. However, in practice, it is difficult to sustain fusion and stereopsis with greater than 8- to 10-percent binocular image size difference.

With extended targets, magnification produces an additional lateral displacement due to cumulative effects such that nonaxial images are enlarged and imaged on noncorresponding points. Given our tolerances for axial disparity gradients in size, it appears that most of the difficulty encountered in aniseikonia is lateral. Von Rohr⁵¹ and Erggelet⁵² considered this prismatic effect as the most important cause of symptomatic discomfort in optically induced aniseikonia.

Differential magnification of the two eyes images can be described as a single constant; however, the magnitude of binocular-position disparity increases proportionally with the eccentricity that an object is viewed from the center of a spectacle lens. Prentice's rule approximates this prismatic displacement from the product of distance from the optical center of the lens in centimeters and the dioptric power of the lens. In cases of anisometropia corrected with spectacle lenses, this prismatic effect produces binocular disparity along the meridian of eccentric gaze. When these disparities are horizontal they result in distortions of stereoscopic depth localization, and when they are vertical, they can produce diplopia and eye strain. Fortunately, there is an increasing tolerance for disparity in the periphery provided by an increase in the binocular sensory fusion range with retinal eccentricity.²² In cases of high anisometropia, such as encountered in monocular aphakia, spectacle corrections can cause diplopia near the central optical zone of the correction lens. In bilateral aphakia, wearing glasses or contact lenses does not interfere with the normal motor or sensory fusion range to horizontal disparity; however, stereothresholds are elevated.⁵³ In unilateral aphakia, contact lenses and intraocular lens implants support a normal motor fusion range; however, stereopsis is far superior with the interocular lens implants.⁵³ Tolerance of size difference errors varies greatly with individuals and their work tasks. Some may tolerate 6 to 8 percent errors while others may complain with only 0.25- to 0.5-percent errors.

Interocular Blur Suppression with Anisometropia

There is a wide variety of natural conditions that present the eyes with unequal image contrast. These conditions include naturally occurring anisometropia, unequal amplitudes of accommodation, and asymmetric convergence on targets that are closer to one eye than the other. This blur can be eliminated in part by a limited degree of differential accommodation of the two eyes⁵⁴ and by

interocular suppression of the blur. The latter mechanism is particularly helpful for a type of contact lens patient who can no longer accommodate (presbyopes) and prefer to wear a near contact lens correction over one eye and a far correction over the other (monovision) rather than wearing bifocal spectacles. For most people (66%), all of these conditions result in clear, nonblurred, binocular percepts with a retention of stereopsis⁵⁵ albeit with the stereothreshold elevated by approximately a factor of two. Interocular blur suppression is reduced for high contrast targets composed of high spatial frequencies.⁵⁵ There is an interaction between interocular blur suppression and binocular rivalry suppression. Measures of binocular rivalry reveal a form of eye dominance defined as the eye that is suppressed least when viewing dichoptic forms of different shape. When the dominant eye for rivalry and aiming or sighting is the same, interocular suppression is more effective than when dominance for sighting and rivalry are crossed (i.e., in different eyes).⁵⁶

When the contrast of the two eyes' stimuli is reduced unequally, stereoacuity is reduced more than if the contrast of both targets is reduced equally. Lowering the contrast in one eye reduces stereo acuity twice as much as when contrast is lowered in both eyes. This contrast paradox only occurs for stereopsis and is not observed for other disparity-based phenomenon including binocular sensory fusion. The contrast paradox occurs mainly with low spatial frequency stimuli (<2.0 cycles/degree). Contrast reduction caused by lens defocus is greater at high spatial frequencies, and perhaps this is why stereopsis is able to persist in the presence of interocular blur caused by anisometropia.

13.5 DISTORTIONS OF SPACE FROM CONVERGENCE RESPONSES TO PRISM

Horizontal vergence responses to prism influence estimates of viewing distance.^{24,57-59} The resulting errors in perceived distance will influence depth percepts that are scaled by perceived viewing distance such as stereopsis and motion parallax. Time of arrival or contact with external objects can be estimated from the distance to contact divided by the speed of approach. If distance is underestimated while wearing convergent or base-out (temporal) prism, then time to contact, computed from velocity and distance, will also be underestimated. In addition, prism stimuli to convergence will produce conflicts in distance estimates from vergence and retinal image size. When vergence responds to prism, retinal image size remains constant instead of increasing as it would in response to real changes in viewing distance. The conflict is resolved by a perceived decrease in perceived size. This illusion is referred to as convergence micropsia. An average magnitude of shrinkage is 25 percent for an accommodation and convergence increase equivalent to 5 D or 5 meter angles.⁶⁰ Very large prism stimuli to vergence present major conflicts between the fixed stimulus to accommodation and the new vergence stimulus. The oculomotor system compromises by partially responding to the vergence stimulus. If physical size of a target in a virtual display has been adjusted to correspond to an expected full vergence response, (increase retinal image size with convergence stimuli), but a smaller vergence response actually occurs, the perceptual compensation for size is based on the smaller vergence response. As a result, a paradoxical increase of perceived size accompanies convergence (convergence macropsia).⁶¹⁻⁶³ The incompletely fused images are perceived at the distance of convergence alignment of the eyes rather than at the intended distance specified by physical size and disparity. Thus, aniseikonic and anisophoric differences in the two eyes, which are made worse by anisometropia, can cause difficulty when performing binocular tasks associated with the use of optical instruments.

13.6 EYE MOVEMENTS

There are two general categories of eye movements. One class stabilizes the motion of the retinal image generated by our own body and head movements and the other class stabilizes the retinal image of an object moving in space. In the former case, the whole retinal image is moving in unison

whereas in the latter case only the retinal image of the moving object moves relative to the stationary background and when that image is tracked it remains stabilized on the fovea during target inspection. In all, there are five types of eye movements (VOR, OKN, saccades, pursuits, and vergence) and they are classified as either stabilizing or tracking movements.

Two classes of reflex eye movements perform the image stabilization task. The first class of stabilizing eye movement compensates for brief head and body rotation and is called the vestibulo-ocular reflex (VOR). During head movements in any direction, the semicircular canals of the vestibular labyrinth signal how fast the head is rotating and the oculomotor system responds reflexively to this signal by rotating the eyes in an equal and opposite rotation. During head rotation, this reflex keeps the line of sight fixed in space and visual images stationary on the retina. Body rotation in darkness causes a repetitious sequence of eye movements which move the line of sight slowly in the direction opposite to head rotation (slow phase) followed by a rapid (saccadic) movement in the direction the head moves (fast phase). This reflex is almost always active and without it we would be unable to see much of anything due to constant smear of the retinal image. The vestibulo-ocular reflex is extremely accurate for active head motion because it has an adaptive plasticity that makes it easily calibrated in response to altered visual feedback.^{64,65} The VOR may be inhibited or suppressed by the visual fixation mechanism⁶⁶ as a result of the system's attempt to recalibrate to the head-stationary image.

Optokinetic reflex is an ocular following response that stabilizes whole field motion when we move about in a visual scene. Unlike the vestibulo-ocular reflex, this optokinetic reflex requires a visible retinal image whereas the vestibulo-ocular reflex operates in total darkness. Rotation or translation of the visual field causes the eyes to execute a slow following and fast reset phase sequence of eye movements that have the same pattern as the VOR. The optokinetic reflex supplements the vestibulo-ocular reflex in several ways. The vestibulo-ocular reflex responds to acceleration and deceleration of the head but not to constant velocity. In contrast, optokinetic reflex responds to constant retinal image velocity caused by constant body rotation or translation. The vestibulo-ocular reflex controls initial image stabilization and optokinetic reflex maintains the stabilization.

Tracking eye movements place the retinal image of the object of regard on the fovea and keep it there, even if the object moves in space. The class of eye movements that places the image on the fovea is called a saccade. This is a very fast eye movement that shifts the image in a step-like motion. Saccades can be made voluntarily, in response to visual, tactile, and auditory stimuli, and even in darkness to willed directions of gaze. Their chief characteristic is they are fast, reaching velocities of nearly 1000 deg/s. Once the target is on the fovea, slow following eye movements called pursuits maintain the image of the target on the fovea. They can keep images stabilized that are formed of objects moving as fast as 30 deg/s in humans. Generally, pursuits cannot be made voluntarily in the absence of a visual moving stimulus. However, they can respond to a moving sound or tactile sensation in darkness.

13.7 COORDINATION AND ALIGNMENT OF THE TWO EYES

If you observe the motion of the two eyes, you will notice that they either move in the same or opposite direction. When they move in the same direction they are said to be yoked, or conjugate, like a team of horses. These are also called version movements. At other times the eyes move in opposite directions, especially when we converge from far to near. These are called disjunctive or vergence movements. Both conjugate and disjunctive movements can be either fast or slow depending if the eyes are shifting fixation or maintaining foveation (stable fixation) of a target. A century ago, Ewald Hering described the movements of the two eyes as equal and symmetrical (Hering's law). Our eyes are yoked when we are born but the yoking can be modified in response to anisometric spectacles that produce unequal image sizes and require us to make unequal eye movements to see bifoveally. This image inequality leads to the binocular anomalies of aniseikonia (unequal perceived image size) and anisophoria (noncomitant phoria or a phoria that varies with direction of gaze).

Targets are brought into alignment with corresponding points by movements of the eyes in opposite directions. These vergence movements have 3° of freedom (horizontal, vertical, and torsional). Horizontal vergence can rotate the eyes temporalward (divergence) or nasalward (convergence). Since near targets only require convergence of the eyes, there is a limited need for divergence other than to overcome postural errors of the eyes. Consequently, the range of divergence is less than half the range of convergence. The eyes are able to diverge approximately 4° and converge approximately 8 to 12° in response to disparity while focusing accurately on a target at a fixed viewing distance. Vergence movements serve the same function as conjugate movements of the eyes which is to preserve the alignment of the two retinal images on corresponding retinal points. Vergence is required when changing viewing distance, whereas versions accomplish the task when changing direction of gaze.

The stimulus to vergence depends upon both the viewing distance and the distance between the two eyes (interpupillary distance). It is equal to the angle subtended by a target in space at the entrance pupils of the two eyes (binocular parallax). As will be described in Sec. 13.9, "Prism-Induced Errors of Eye Alignment," the stimulus to vergence can be changed by stimulating accommodation with lenses, or changing disparity (binocular parallax) with prisms. The exact stimulus value of lenses and prisms to vergence depends upon the placement of the prism relative to the eye and this relationship will be described in Sec. 13.8, "Effects of Lenses and Prism on Vergence and Phoria." Vergence is measured in degrees of ocular rotation and prism diopters. The meter angle is also used in order to get a comparable metric for accommodation and convergence. Meter angles equal the reciprocal of the viewing distance in meters, and the value is approximately equal to the accommodative stimulus in diopters.

$$MA = \frac{1}{\text{target distance in meters}}$$

Meter angles (MA) can be converted into prism diopters simply by multiplying MA by the interpupillary distance (IPD) in cm. Prism diopters equal the 100 tan vergence angle in degrees. Vergence eye movements respond to four classes of stimuli described nearly a century ago by Maddox.

Intrinsic Stimuli to Vergence

There are four classes of stimuli to align the eyes with horizontal vergence at various viewing distances. The first is an intrinsic innervation that develops postnatally. During the first 6 weeks of life, our eyes diverge during sleep. This divergence posture is known as the anatomical position of rest and it represents the vergence angle in the absence of any innervation. As we develop during the first 1.5 months of life, the resting vergence angle is reduced to nearly zero by the development of tonic vergence which reduces the resting position from the anatomical position of rest to the physiological position of rest. The anatomical position of rest refers to eye alignment in the absence of all innervation, such as in deep anesthesia. The physiological position of rest refers to the angle that vergence assumes when we are awake and alert, in the absence of any stimulus to accommodation or binocular fusion or perceived distance. Usually, this physiological position of rest differs from the vergence demand, (i.e., the amount of vergence needed to fuse a distance target). The discrepancy between the vergence demand at a far viewing distance and the physiological position of rest is called the phoria. The phoria is an error of binocular alignment, but the error is only manifest during monocular viewing or any condition that disrupts binocular vision. If binocular vision is allowed (not disrupted) other components of vergence reduce the error nearly to zero. Adaptation of the phoria continually calibrates eye alignment and keeps the phoria constant throughout life.

There are two components of phoria adaptation that compensate for eye misalignment. One is a conjugate process that is called prism adaptation and it refers to our ability to adjust our phoria in response to a uniform prism. This prism is constant throughout the visual field and it produces a conjugate change in phoria in all directions of gaze. The conjugate adaptation is very rapid and can be completed in several minutes in a normal patient. The eyes are able to adapt horizontal vergence several degrees over a 1 h time span and larger amounts for unlimited adaptation periods. The slow-anomalous fusional movements reported by Bagolini⁶⁷ for strabismus patients are similar in time course to adaptive responses to prism and they indicate that the potential range of adaptation must exceed 10° .

The second component of phoria adaptation is a nonconjugate process that adjusts the phoria separately for each direction and distance of gaze. For example, anisometric spectacle corrections disrupt the normal alignment pattern needed for fusion. A person wearing a plus (converging) lens over the right eye needs to over converge in left gaze and over diverge in right gaze. Our oculomotor system is able to make these noncomitant adjustments such that the phoria varies with direction and distance of gaze consistent to compensate for the problems encountered with anisometric spectacle corrections. This disconjugate adaptation occurs in a matter of hours and it persists when the spectacle lenses are removed. It provides a means to align the eyes without retinal cues of binocular disparity. The adapted variations of phoria are guided by a synergistic link with the eye positions in which they were learned. For example when an anisometric person first wears a plus lens spectacle over the right eye, the right eye is hypophoric in upgaze and hyperphoric in downgaze. After wearing this spectacle correction for several hours the vergence errors are adapted to so that there is no longer a phoria or vergence error while wearing the spectacles. However, when the spectacles are removed, the right eye is hyperphoric in upgaze and hypophoric in downgaze. This demonstrates a synergistic link between vertical eye position and vertical phoria where eye position guides the vertical vergence alignment of the two eyes, even when only one eye is open and the other is occluded.

Binocular Parallax

Horizontal vergence also responds to misalignment of similar images from corresponding retinal points (retinal image disparity). It reduces the disparity to a value within Panum's fusional area and causes images to appear fused or single. This fusional or disparity vergence is analogous to optical reflex accommodation. It serves as a rapid fine adjustment mechanism that holds or maintains fixation of a given target and corrects for minor changes of disparity caused by small movements of the head or the target being inspected. It has a range of approximately 3° of convergence and divergence. Rapid adaptations of the phoria account for increased ranges of fusion in response to disparity.

Vertical and cyclovergence eye movements also respond to binocular disparity. Vertical vergence can rotate one eye up or down relative to the other eye (skew movements) over a range of 2° . These movements respond to vertical disparity that is always associated with near targets in tertiary directions of gaze (targets at oblique eccentricities). Vertical disparity is a consequence of spatial geometry. Tertiary targets are closer to one eye than the other so that they subtend unequal image sizes in the two eyes. Unequal vertical eccentricities resulting from unequal image magnification requires a vertical vergence to align tertiary targets onto corresponding retinal regions. At near distances greater than 10 in, vertical disparities are very small and only require a limited range of movement.

Torsional or cyclovergence can twist one eye about its visual axis relative to the orientation of the other eye. Normally, torsional vergence keeps the horizontal meridians of the retinas aligned when the eyes are converged on a near target in the upper or lower visual field. Torsional vergence stimulated by cyclodisparities is limited to a range of approximately 10° . Like vertical vergence, torsional vergence responds to unequal orientations or torsional disparities that result from the perspective distortions of near targets viewed by the two eyes. Ocular torsion is mainly stimulated by lateral head roll and amounts to approximately 10 percent of the amount of head roll, up to a maximum of 6° (60° head roll). The head roll can produce vertical disparities of the near fixation targets if the eyes are not accurately converged. With the head tilted toward the right shoulder, a divergence fixation error can cause the right eye to be depressed from the target and a convergence fixation error will cause the right eye to be elevated with respect to the fixation target. However, no vertical fixation errors result from head roll if the eyes are accurately converged on the fixation target.

Defocus and Efforts to Clear Vision

Stimulation of either convergence or accommodation causes synergistic changes in convergence, accommodation, and pupil size. This synergy is called the near triad. During monocular occlusion, when the eyes focus from optical infinity to 1 m (1 D) they converge approximately 2 to 3° without

disparity cues. This accommodative/vergence response brings the eyes into approximate binocular alignment, and any residual alignment errors produce binocular disparity that is corrected by fusional vergence. The activation of convergence when accommodation is stimulated while one eye is covered is called accommodative vergence. The gain or strength of this interaction is called the accommodative convergence/accommodation (AC/A) ratio which describes the amount of convergence activated per diopter of accommodation. Similarly, if convergence is stimulated while wearing pinhole pupils, accommodation also responds without blur cues. This reflex is called convergence accommodation and the gain or strength of this interaction is the convergence-accommodation/convergence (CA/C) ratio. Normally these ratios are tuned appropriately for the geometry of our eye separation so that efforts of either accommodation or convergence bring the eyes into focus and alignment at nearly the same viewing distance. Occasionally, because of uncorrected refractive errors, or large abnormal phorias, or because of optical distortions caused by instruments, there is an imbalance of accommodation and vergence causing them to respond to different distances. The mismatch may result in fatigue and associated symptoms if the visual task is extended. The resulting mismatch needs to be corrected by fusional vergence and optical reflex accommodation. If the mismatch is beyond the range of the compensation, then vision remains blurred or double.

Variations of Cross-Coupling with Direction and Distance of Gaze

In straight-ahead gaze (zero azimuth), targets in the midsagittal plane produce equal magnitudes of disparity and defocus when they are expressed in meter angles [$MA = 1/\text{viewing distance (m)}$] and diopters [$D = 1/\text{viewing distance (m)}$] respectively. However, when targets are viewed in asymmetric vergence [i.e., with some combination of horizontal version (azimuth) and vergence], conflicts arise between stimuli for accommodation and convergence as a consequence of 3-D spatial geometry. Figure 10 compares a plan view of the isoaccommodation circle to the isovergence circle. The isoaccommodation circle describes the locus of points that subtend a constant average accommodation

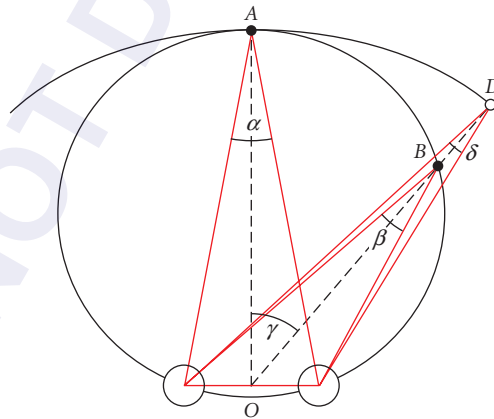


FIGURE 10 Diagram showing isovergence circle (large complete circle), isoaccommodation circle (shown as an arc) and various azimuth angles. The two small circles represent the right and the left eye. The point A in the figure represents the spatial location that corresponds to matched stimuli for accommodation and convergence. α = vergence angle at point A (angle made by the intersection of two lines of sight), β = vergence angle at point B, and $\alpha = \beta$ along the isovergence circle. γ = lateral gaze angle. Note that at point D on the isoaccommodation circle, the convergence and accommodation demand is lower than at point B.

stimulus to the two eyes with increasing azimuth. This circle has a radius equal to the viewing distance from the cyclopean eye to the object of regard in any horizontal direction of gaze lying in a common visual plane. The isovergence circle passes through the fixation point and two centers of eye rotation.⁶⁸ Its distance to the cyclopean eye varies with both viewing distance and direction of gaze (azimuth). The circle has a radius equal to the interpupillary distance (PD) divided by twice the sine of the angle of convergence (theta).

$$\text{Isovergence radius} = \frac{\text{PD}}{2 * \sin \theta}$$

The only spatial location that corresponds to matched stimuli for accommodation and convergence is straight ahead where these two circles intersect (Fig. 10). For points lying to the left or right, the stimulus to accommodation is always greater than the stimulus to convergence. This conflict is resolved by having the normal cross-link interaction between accommodative-convergence and accommodation (AC/A ratio) set for an ideal value for 50° azimuth.⁶⁹ This low ratio results in a small divergence error of vergence at lesser azimuths that is easily overcome with fusional convergence. The normal convergence accommodation/convergence interaction (CA/C ratio) is ideal for straight ahead. This produces a lag of accommodation in small lag of accommodation that increases with gaze eccentricity.

Perceived Distance

The eyes align and focus in response to perceived distance. The term proximal refers to perceived distance. Perceived distance is classified as a spatiotopic cue because it results in a conscious percept of distance and direction. Our perception of space is made up of an infinite number of different targets in different directions and distances. Our eyes can only select and fuse one of these target configurations at a time. One of the main tasks of the oculomotor system is to select one fixation target from many alternatives. This selection process usually results in a shift of binocular fixation from one distance and direction to another. Proximal vergence serves as a course adjustment mechanism and accounts for nearly 100 percent of the eye movement when gaze shifts are too large for us to sense them from blur and disparity. We perform these large gaze shifts without needing to see the new intended target clearly or even with both eyes.⁷⁰ We only need to have an estimate of target location using monocular cues for distance and direction. The majority of binocular shifts in fixation are stimulated by proximal vergence. Once the shift has been made another component of vergence (fusional vergence) takes over to maintain the response at the new target location. Errors in perceived distance can result in errors of eye alignment and focus. Mislocalization of large print or symbols in a head mounted display, or of the distance of a microscope slide can cause the eyes to accommodate excessively and produce a blurred image of a target that is imaged at optical infinity.

The Zone of Clear and Single Binocular Vision

These components of vergence add linearly to form a range of accommodative and disparity stimuli that the observer is able to respond to with single and clear binocular vision (Fig. 11).² The range is determined by the convergence and divergence fusional vergence range from the resting phoria position, plus the linkage between accommodation and vergence described in the section on “Defocus and Efforts to Clear Vision” that increases with accommodative effort until the amplitude of accommodation is reached. Geometrically matched stimuli for convergence and accommodation that correspond to common viewing distances normally lie within this range of clear single binocular vision unless there is an alignment error such as an eye turn or strabismus. Mismatches between stimuli for accommodation and convergence produced by lenses or prisms or optical instruments can be compensated for as long as the combined stimuli for accommodation and convergence lie within this stimulus operating zone for binocular vision. The ability to respond to unmatched stimuli for convergence and accommodation does not result from decoupled independent responses of the

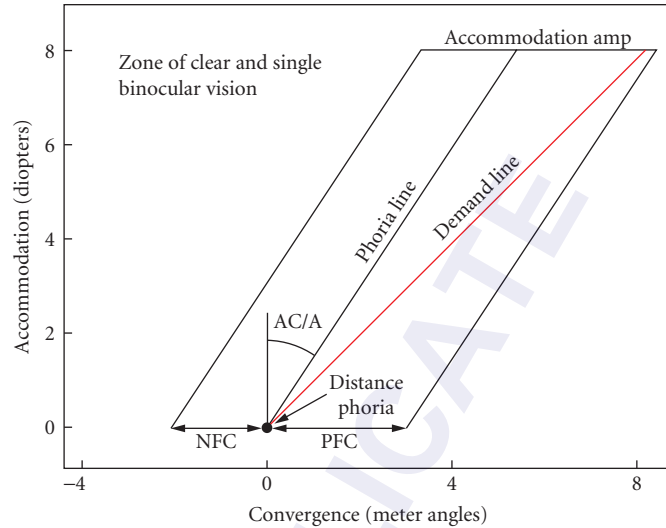


FIGURE 11 Figure of the ZCSBV. Schematic representation of the zone of clear single binocular vision in terms of five fundamental variables (amplitude of accommodation, distance phoria, AC/A ratio positive fusional convergence, negative fusional convergence).

two systems. Accommodation and convergence always remain cross-coupled, but fusional vergence and optical reflex accommodation are able to make responses to overcome differences in phorias produced by accommodative vergence and prisms.

13.8 EFFECTS OF LENSES AND PRISM ON VERGENCE AND PHORIA

Errors of horizontal and vertical vergence can result from differences in the two ocular images caused by differential magnification and translation. Optical power differences will cause magnification differences. Translation of the image is caused by prism or by decentration of the line of sight from the optical center of the lens system.

Magnification-Induced Errors of Eye Alignment

Magnification of the ocular image is influenced by several factors including index of refraction, lens thickness, front surface power, overall power of the lens, and the ocular distance from the eyepiece. For a spectacle lens the magnification is described by the product of a shape factor and a power factor. Magnification described by the shape factor (M_s) equals

$$M_s = \frac{1}{1 - F_1 t/n}$$

and by the power factor (M_p) equals

$$M_p = \frac{1}{1 - P_v h}$$

where F_1 is the front surface power, t is the lens thickness in centimeters, n is the index of refraction, h is the distance of the spectacle plane from the cornea in centimeters, and P_v is the thin lens power. Percent angular magnification caused by the power factor is approximately 1.5 percent per diopter power. The percent magnification ($M\%$) can be approximated by the following formula:²²

$$M\% = \frac{F_1 t}{n + P_v h}$$

Differences between the lenses before the two eyes will produce image size differences, depending on the source of refractive error. Two basic sources are an abnormal axial length (axial ametropia) or abnormal corneal and lenticular power (refractive ametropia). Magnification differences between the two eyes caused by interocular differences in axial length can be eliminated largely by correcting axial errors with spectacle lenses located in the primary focal plane (the spectacle plane) and placing the lenses not too near to the eyes. Magnification differences between the two eyes caused by unequal powers of the two corneas can be prevented by correcting refractive errors with contact lenses instead of spectacle lenses (Knapp's law). If Knapp's law is violated by correcting a refractive myope with spectacle lenses placed in the primary focal plane of the eye (16.67 mm), the myopic eye will have a smaller image than an emmetropic eye with a normal corneal power. Similarly, if an axial myope is corrected with a contact lens, the myopic eye will have a larger image than an emmetropic eye with a normal axial length. However, the perceived image size in axial refractive errors also depends on whether the retina is stretched in proportion to the axial length. If it is, then the contact lens correction for an axial ametropia would produce a perceived image size equal to that of an emmetropic eye with a normal axial length. This analysis only applies to perceived image size and it is an attempt to correct perceived image size differences between the two eyes (aniseikonia). However, it does not address the unequal ocular movements stimulated by anisometric spectacle corrections (anisophoria) resulting, as described in Sec. 13.10 on gaze, from the separation of the entrance pupil and center of rotation.

Anisophoria is defined by Friedenwald⁷¹ as a heterophoria which varies in magnitude with direction of gaze. It is produced by spectacle corrections for anisometropia independent of whether refractive errors are axial or refractive in nature and only concerns the refractive power of the spectacle lens correction. The lateral prismatic effect occurs even when size distortions of an axial refractive error have been properly corrected with a spectacle lens. A prismatic deviation will always exist whenever the visual axis of the eye and optical axis of a lens do not coincide. As a result, the initial movement of an ametropic eye that is corrected with spectacle lenses will be too small in hypermetropia and too large in myopia. These optical factors produce differential prism in anisometropia corrected with spectacles. Initially, anisometropes corrected with spectacles will make vergence errors during versional movements to eccentric targets because the eyes follow Hering's law of equal innervation to the extraocular muscles. Interocular differences in power of spectacle lenses can stimulate vergence eye movements in all meridians of anisometropia (vertical, horizontal, and oblique). This optically induced anisophoria can be compensated for easily in the horizontal meridian by disparity vergence eye movements. The vertical errors will be the most troublesome because of the limited range and slow speed of vertical vergence.^{72,73} Difficulty in compensating for optically induced anisophoria is believed to contribute to the symptoms of aniseikonia corrected with spectacle lenses.^{71,74} Anisophoria may explain why both refractive and axial anisometropes prefer contact lens corrections that do not produce anisophoria. Contact lenses move with the eye, and the line of sight always passes through the optical center such that fixation errors of targets at infinity always equal the correct motor error to guide the oculomotor system during gaze shifts. Optically induced anisophoria can be very large. For example, Ogle²² computes a 4.2 prism diopter phoria produced by a +4-D anisometropia spectacle correction during a 20° eye movement. This constitutes a 12-percent error. Normally the oculomotor system is capable of recalibration of Hering's law to compensate for optically induced anisophoria. For example, anisometric patients wearing bifocal spectacle corrections do not exhibit the predicted change in vertical phoria while viewing in downgaze.^{75,76} Adjustments of heterophoria in all directions of gaze occur in approximately 2.5 h for differential afocal magnification of up to 10 percent.^{77,78} Similar oculomotor adjustments occur normally in asymmetric convergence whereby the abducted eye makes larger saccades than the adducted eye in order to compensate for geometrically induced aniseikonia.⁷⁹ Unequal prism before the two eyes produced by a 10-percent magnification difference,

such as observed in monocular aphakia corrected with spectacles or spectacle and contact lens combinations,⁸⁰ is extremely difficult to adapt to. Presumably a 10-percent magnification difference sets an upper limit to the degree of anisometropia that will not disrupt binocular versional eye movements and binocular alignment. It is possible that individuals who have symptoms of aniseikonia when corrected with spectacle lenses, are simply unable to make the necessary adjustments in conjugacy to overcome their optically induced anisophoria. Indeed Charnwood⁸¹ observed that aniseikonia patients preferred prism incorporated into their spectacle corrections to reduce anisophoria and Field⁸² reported that nearly all of 25 patients with clinically significant aniseikonia had oculomotor imbalance. Clearly anisophoria is an important by-product of anisometropia corrected with spectacles.

Torsional vergence errors can also result from optical corrections that have cylindrical corrections for astigmatism. The power of a cylindrical lens varies with the square of the cosine of the angular separation from the major power meridian. This power variation produces a meridian-dependent magnification of the ocular image described as a magnification ellipse. The ellipse produces rotational deviations of lines at different orientations toward the meridian of maximum magnification. The rotary deviation between a line and its image meridian (δ) depends on the angle between the orientation of the line (γ) and the major axis of the magnification ellipse (θ). The rotary deviation is approximated by the following equation:

$$100 \tan \delta = 1/2 f \sin 2(\gamma\theta)$$

where f equals the percent of magnification produced by the cylinder.²² This equation demonstrates that most of the shear or scissor distortion of lines will occur in meridians between the major and minor axes of the magnification ellipse. These rotary deviations produce lines imaged on noncorresponding retinal meridians or torsional disparities because the axis of astigmatism is normally not parallel in the two eyes. Rather, astigmatism tends to be mirror symmetric about oblique meridians such that an axis of 45° in one eye is often paired with an axis of 135° in the other eye. As a result, the shear or rotary deviations of vertical meridians tend to be in opposite directions. The resulting cyclodisparities produce perceptual declination errors of stereo surface slant (tilt of a surface about a vertical or horizontal axis). Cyclodisparities about horizontal axes stimulate cyclovergence movements of the two eyes (torsion of the two eyes in opposite directions). If the disparities are long lasting, the cyclovergence response adapts and persists even in the absence of binocular stimulation.⁸³

13.9 PRISM-INDUCED ERRORS OF EYE ALIGNMENT

Errors of binocular eye alignment also result from translation of one ocular image with respect to the other. The translation can be caused by prism or by displacement of the optical center of the lens from the line of sight. When an object is viewed through a prism, it appears in a new direction that is deviated toward the prism apex and the eye must move in that direction to obtain foveal fixation. The amount the eye must rotate to compensate for the prismatic displacement of targets at finite viewing distances decreases as the distance between the eye and prism increases. This is because the angle that the displaced image subtends at the center of eye rotation decreases with eye-prism separation. This reduction is referred to as prism effectivity.⁸⁴ It does not occur for targets at infinity since the target distance is effectively the same to both the prism and center of rotation. The effective prism magnitude at the eye (Z_e) is described by the following equation:

$$Z_e = \frac{Z}{1 - C_{\text{rot}} \times V}$$

where Z is prism value, C is center of rotation to prism separation in meters, and V is the reciprocal of the object to lens separation in meters which corresponds to the dioptric vergence leaving the prism.⁸⁴

Binocular disparities induced by a prism or lens determines the direction of the phoria or vergence error. For example, base-out prism produces a crossed disparity and divergence error (exophoria) which is corrected under binocular conditions by fusional convergence. Minus lenses stimulate accommodation and accommodative convergence, which under binocular conditions produce a convergence error (esophoria) that is corrected by fusional divergence. Similarly, base-in prism produces an uncrossed disparity and convergence error (esophoria) which is corrected under binocular conditions by fusional divergence, and plus lenses stimulate a relaxation of both accommodation and accommodative convergence, which produces a divergence error (exophoria) that is corrected by fusional convergence.

The prismatic translation caused by optical displacement is described by Prentice's rule, where the prism displacement (Z) equals the power of the lens in diopters (P) times the displacement of the lens center from the line of sight (d) in centimeters ($Z = dP$). The angular deviation of any paraxial ray that passes through a lens at a distance from the optical center is described by this rule. Prentice's rule can also be applied to cylindrical lenses where the magnification occurs in a power meridian. The axis of a cylindrical lens is perpendicular to its power meridian. The prism induced can be calculated from the vector component of the displacement that is perpendicular to the axis of the cylinder. This distance (Hp) equals the sine of the angle (θ) between the displacement direction and the axis of the cylinder times the magnitude of the displacement (h)

$$Hp = h \sin \theta$$

Thus, according to Prentice's rule, prism induced by displacing a cylindrical lens equals

$$Z = Ph \sin \theta$$

The magnitude of prism induced by lens displacement also depends on the separation between the lens and the center of rotation of the eye. Effective prism displacement (Z_e) decreases as the distance of the lens from the eye increases according to the following equation:⁸⁴

$$Z_e = \frac{(Zc + Z)}{1 - C_{\text{rot}} \times V}$$

where Zc is $100 \tan c$, c is the angle formed between the object point and the displaced point from the optical center (major reference point) and the straight ahead line, V is the reciprocal of the image distance to the lens in meters which corresponds to the dioptric vergence leaving the prism, Z is the prism calculated from Prentice's rule, C_{rot} is separation between the lens and center of rotation of the eye. The reduction of prismatic displacement with eye to lens separation is referred to as prism effectivity with lenses.

When the two eyes view targets through lenses that are displaced from the optical axes, a disparity stimulus is introduced. Most typically, when the optical centers of a binocular optical system are separated by an amount greater or less than the interpupillary distance, disparate horizontal prism is introduced that requires the eyes to make a horizontal vergence adjustment to achieve binocular alignment. For example, when looking through a pair of plus lenses that are farther apart than the interpupillary distance, each eye is viewing through a point that is nasalward from the optical axis. This introduces base-out prism that stimulates convergence of the two eyes. This prism vergence stimulus may produce conflicts with the accommodative stimulus that is not affected by lens displacement so that the eyes need to focus for one distance and converge for another. Accommodation and convergence are linked synergistically so that their responses are closely matched to similar target distances. There is a limited degree of disparity vergence compensation available to overcome any discrepancy between the two systems, and these efforts can result in headaches and other ocular discomfort.

Anisometropic ophthalmic lens prescriptions with bifocal corrections in the lower segment of the lens can also introduce vertical prism while viewing near objects through the reading addition. The amount of vertical prism is predicted by Prentice's rule applied to the amount of anisometropia. Thus, for a 2-D anisometropic individual, a patient looking downward 1 cm through their bifocal has a vertical phoria of two prism diopters. The vertical phoria can be corrected with a slab-off

prism (a base-down prism that is cemented onto the bifocal over the lens with the most plus power). However, this procedure is rarely necessary, since the patient usually is able to adapt the vertical phoria to vary nonconjugately with vertical eye position as described above.

13.10 HEAD AND EYE RESPONSES TO DIRECTION (GAZE CONTROL)

When objects attract our attention in the peripheral visual field, we are able to fixate or shift our central gaze toward them using head and eye movements. The eyes move comfortably over a range of 15° starting from primary position. Larger gaze shifts require combinations of head and eye movements. Normally, the eyes begin the large movements, then the head moves toward the target while the eye movements are reduced by the VOR, and when the head stops the eyes refine their alignment. The reduction of eye motion during head rotation results from a gaze stabilization reflex that allows us to keep our eyes fixed on a stationary object as we move our body or head. This reflex operates with head motion sensors called the vestibular organs.

The manner in which we shift our gaze influences optical designs. Optical systems usually are mounted on the head. They need to provide a large enough viewing angle to accommodate the normal range of eye movements, and to remain stable on the head as it moves. They also need to provide peak resolution in all directions that the oculomotor system is able to shift gaze while the head remains stationary. Gaze shifts can be inaccurate because of interactions between spectacle refractive corrections and movements of the eyes.

The accuracy of eye movements is affected by stationary optical systems such as spectacles that are mounted in front of the eyes as opposed to optical systems such as contact lenses that move with the eye. The problem occurs when viewing distant targets because the entrance pupil of the eye and its center of rotation are at different distances from the spectacle plane. The entrance pupil is at a distance of about 15 mm from the spectacle lens and the ocular center of rotation lies about 27 mm behind the spectacle plane. Thus the center of rotation and entrance pupil are separated by about 12 mm and a target in space subtends different angles at these two points. The angle at the entrance pupil determines the visual error and the angle at the center of rotation determines the motor error. Whether a spectacle lens is worn or not, this difference alters retinal image feedback during eye rotations for any target imaged at a finite viewing distance from the eye. Images at infinity don't have this problem because they subtend equal angles at both the entrance pupil and center of rotation. Thus eye movements cause motion parallax between targets viewed at near and far viewing distances and some animals like the chameleon use this parallax to judge distance.⁸⁵

Problems occur when the spectacle lens forms an image of an object at infinity at some finite distance from the eye. This image of the distant target subtends different angles at the entrance pupil and center of rotation because of their 12-mm separation. Thus, perceived visual angle and motor error are unequal. Depending on whether the image is formed in front (myopic error) or behind (hyperopic error) the eye, this discrepancy requires that a fixational eye movement must be smaller or larger, respectively, than the perceived eccentricity or retinal angular error of the target sensed prior to the eye movement. The amount of eye rotation can be computed from the equation describing prism effectivity with lenses (Z_e):⁸⁴

$$Z_e = \frac{Zc}{1 - C_{\text{rot}} \times V}$$

where Zc equals $100 \tan c$ (c equals angular eccentricity of the target measured at the lens plane), and V equals the dioptric vergence of the image distance from the lens plane, and C_{rot} equals the lens to center of rotation separation in meters. When different power spectacle lenses are worn before the two eyes to correct unequal refractive errors (anisometropia), a target that would normally require equal movements of the two eyes now requires unequal binocular movements. The oculomotor system learns to make these unequal movements and it retains this ability even when one eye is occluded after wearing anisometric spectacle corrections for only 1 hour.

During head movements, magnification of the retinal image by spectacles increases retinal image motion during head rotation and causes the stationary world to appear to move opposite to the direction the head rotates. This can produce vertigo as well as stimulate a change in the vestibular stabilization reflex. When a magnifier is worn for only 1 h, the reflex can be exaggerated to compensate for the excessive retinal image motion. Following this adaptation, when the optical aid is removed, the exaggerated VOR causes the world to appear to move in the same direction that the head rotates. If you wear glasses, you can demonstrate this adaptation to yourself by removing your glasses and shaking your head laterally. If you wear negative lenses to correct for myopia, the minification they produce has caused a reduction in image stabilization and without your glasses, the world will appear to move opposite to the direction of head shake.

13.11 FOCUS AND RESPONSES TO DISTANCE

Our visual sense of space requires that we are able to accurately perceive objects at different distances. Single and clear views of targets at a finite viewing distance require that the visual system align the similar images of the target on corresponding retinal regions and that the optics of the eye adjust the focus to accommodate the near target (accommodation). The range of accommodation varies with age, starting at birth at about 18.5 D and decreasing 1 D every 3 years until approximately age 55 when it is no longer available (absolute presbyopia). The amplitude of accommodation available to an individual of a given age is remarkably stable across the population if extremes of refractive are not considered.

Spectacle corrections and other optical systems can influence both aspects of the near response. Accommodation declines with age until eventually, the near point is beyond the near working distance (functional presbyopia). In this situation, near targets can be imaged farther away from the eye with plus lenses so that they lie within the available range of remaining accommodation. The accommodative system is only able to maintain 2/3 of its full amplitude for long periods of time and when this range is exceeded by the optical demands of working distance, bifocals are typically prescribed. The magnitude of the bifocal correction depends upon the accommodative amplitude and the near working distance. For a typical 16-in or 40-cm near working distance that requires 2.5 D of accommodation, a bifocal correction is needed when the amplitude of accommodation is less than 3.75 D. Typically this amplitude is reached at about 45 years of age.

If spectacles are worn, the amplitude of accommodation is influenced by the sign (plus or minus) of the spectacle correction. Accommodation demand equals the difference in the dioptric value of the far point and the near point of the eye. These dioptric values can be described relative to the front surface of the eye or to the spectacle plane. The reciprocal of a specific viewing distance to the spectacle plane describes the spectacle accommodative demand, and this value is independent of the lens correction. However, the ocular accommodative demand (amount the eye needs to change power) required to focus targets at this distance does depend on the power and sign of the spectacle correction. Ocular accommodation is computed by the difference between the dioptric vergence subtended by the near point target at the cornea, after it is imaged by the spectacle lens, and the vergence subtended at the cornea by the far point of the eye. For myopes corrected with minus spectacle lenses, ocular accommodation through spectacles is greater than the spectacle accommodative demand and the reverse is true for hyperopes. For example, the ocular accommodative response to a target placed 20 cm in front of the spectacle plane that is 14 mm from the eye is 4.67 D from an emmetrope, 4 D for a myope wearing a -6-D corrective spectacle lens, and 5.56 D for a hyperope wearing a +6 D corrective spectacle lens. Thus, near targets present unequal accommodative stimuli to the two eyes of anisometropes corrected with spectacles. Although the accommodative response is consensual (equal in the two eyes) the difference in the two stimuli will go unnoticed as long as it lies within the depth of focus of the eye (approximately 0.25 D). However, in the example given above, a 6-D anisometrope would have about a 1-D difference in the stimulus to accommodation presented by a target at a 20-cm viewing distance. Finally, this lens effectivity causes spectacle lens wearers who are hyperopic to become presbyopic before myopes.

13.12 VIDEO HEAD SETS, HEAD'S UP DISPLAYS AND VIRTUAL REALITY: IMPACT ON BINOCULAR VISION

A variety of head mounted visual display systems (HMD) have been developed to assist performance of a variety of tasks. Telepresence is an environment in which images from remote cameras are viewed on video screens to perform tasks at distal sites. In other applications, information in the form of symbology is added to a natural image to augment reality, such as in a heads up aviation application, and in still another application, virtual environments are simulated and displayed for entertainment or instructional value. The head mounted video systems usually display information in one of three configurations. These are monocular or monoscopic that present information to one eye, biocular where both eyes view an identical image, and binocular where slight differences in the two eyes' images stimulate stereoscopic depth perception. The biocular view avoids image conflicts between the two eyes that arise in the monoscopic systems, and stereopsis in the binocular systems helps to sort out or distinguish elevation and distance information in the perspective scene. Stereopsis also provides the user with fine depth discrimination and shape sensitivity for objects located within an arm's length. Each of these systems attempts to present coherent cues to space perception, however, due to physical limitations this is not always possible and cue conflicts result.

Distance Conflicts

The current designs of HMD devices introduce three stimulus conflicts. They present conflicting information about target distance to accommodation and convergence, conflicting visual-vestibular information, and conflicts concerning the spatial location of information to the two eyes. For example, the symbology presented to aviators in heads-up displays is set for infinity but large symbols appear to be near and this can initiate an inappropriate accommodative response. However, given sufficient time, subjects will accommodate correctly.⁸⁶ Currently all HMD units present a fixed optical viewing distance that is not always matched to the convergence stimulus in biocular units or the many states of convergence stimulated in binocular units. This situation is exacerbated by the presence of uncorrected refractive errors that can cause images to be blurred in the case of myopia or to force extra accommodation in the case of hyperopia, or to present unequal blur and accommodative stimuli in the case of anisometropia. Normally, observers with uncorrected refractive errors are able to adjust target distance to clear images but the HMD has a fixed target distance that is not adjustable. Units that provide self-adjusting focus run the risk of users introducing larger errors or mismatching the clarity of the two ocular images.

As described above, accommodation and convergence are tightly coupled to respond to approximately the same viewing distance. There is some flexibility to compensate for small differences between stimuli to accommodation and convergence but long periods of usage, particularly those requiring fusional divergence, to overcome an esophoria caused by excessive accommodation can produce eye strain. Tolerance for accommodative errors can be increased by extending the depth of focus with small exit pupils or with screens that have reduced resolution. It is also helpful to set the viewing distance at an intermediate value of approximately 2 m⁸⁷ to allow uncorrected mild myopes to see clearly while allowing a 3° range of stereo depth disparities. In biocular units where disparity is not varied, the optical distance can be set closer to 1 m to further compensate for uncorrected myopia and to place images near arm's length. Prolonged viewing at close distances can cause instrument myopia, a form of excessive accommodation, and postimmersion short-term aftereffects such as accommodative spasm.⁸⁸ Accommodative fatigue also can result after 30 min of stereo display use.⁸⁹ If refractive corrections are worn, the optimal image distance will depend on the tasks performed. Simulated computer screens should be at arm's length whereas views for driving simulators should be remote to provide more realism in the display.

Errors of perceived distance can occur when vergence distance is not matched to image size or perceptual task. These errors can result in perceptual image expansion (macropsia) or constriction (micropsia), particularly in binocular displays that stimulate a range of convergence values. These effects can be minimized by limiting the range of binocular disparities to a few degrees for near

objects and to use a biocular display when presenting views of distant objects simultaneously with foreground detail. Perspective cues along with motion parallax, overlap or superposition, and texture density gradients provide powerful depth cues for remote or complex scenes containing a wide range of depth information and binocular disparity adds little to these displays. Stereopsis is most useful for facilitating manual dexterity tasks with fine depth discrimination between objects and their 3-D shape located within an arm's length of the user. Static displays that don't utilize motion parallax will benefit more from the addition of stereoscopic cues.

Some HMDs and other virtual environments combine real and virtual information with a see-through visor or with remote rear projection systems that surround the viewer panoramically as in the Computer-Automatic Virtual Environment.⁹⁰ These mixed environments can introduce multiple and sometimes conflicting accommodative stimuli. When real variable focus and virtual fixed focus images are at different image distances, both cannot be cleared at the same time, though they may appear to be in similar locations in space. Conflicting near virtual targets can interfere with recognition of real distant targets.⁹¹ Several optical solutions have been suggested to minimize these focus conflicts.⁵⁴ The depth of focus can be extended with pinhole pupils. A pinhole exit pupil projected to the eye pupil in a Maxwellian view system will preserve field of view. The virtual image is viewed through the Maxwellian pupil imaged via a beam splitter while the real environment is viewed through the natural pupil. Eye tracking might be necessary in Maxwellian systems to shift the exit pupil to follow eye movements.⁹² Monovision is another possible solution in which a near lens is worn over one eye and a far lens over the other. A target that is in focus for one eye will be blurred for the other. The visual system is able to suppress the blurred image and retain the focused images of both eyes and also retain a large functional range of stereopsis.⁹³ Monovision has been used successfully in presbyopes with contact lens corrections, and it has been shown to work on prepresbyopes in most individuals.⁵⁵ The limitation of monovision is that it provides clear vision only for the two focal depths, plus and minus the depth of focus. The range of the two clear vision distance might be extended by aniso-accommodation; the ability of the two eyes to focus independently of each other while remaining binocular.⁵⁴ Another solution is a chromatic bifocal that takes advantage of the longitudinal chromatic aberration of the eye. Wavelengths of light at opposite ends of the visual spectrum have a 2-D difference in dioptric power for a given object distance. Given that nearer objects tend to be in the lower visual field, different monochromatic wavelengths combined with field segregation could be used in the virtual image control panel. Users viewing near objects could view long-wavelength images in the lower half of the virtual control panel. Thus, both real and red virtual images would stimulate an increase in accommodation. For distant real objects, the user could attend to the upper short wavelength images in the virtual display panel and real and blue virtual targets would stimulate a relaxation of accommodation. A simpler solution is to place a bifocal lens segment in the lower portion of the see-through window, bringing near real images to the same focal distance as the images seen on the virtual display.

Visual-Vestibular Conflicts

Virtual systems that are worn on the head or remain fixed in space both present conflicting visual-vestibular information. Normally, when the head rotates or translates, the eyes rotate in the opposite direction to keep the images of objects in space stabilized on the retina. The relationship between retinal image motion and head and eye motion is disrupted by the virtual display. In the case of an external virtual projection system like CAVE (Computer-Augmented Virtual Environment),⁹⁰ simulated movements of the visual scene that correspond to body movements are not accompanied by the appropriate vestibular signals. For example, a driving simulator displays changes in the road view while steering the virtual car. When making turns, the visual field appears to rotate but there is no corresponding vestibular signal. As a result, drivers of virtual systems tend to oversteer the turns because they underestimate the curved path of the automobile. Perceptual disturbances during similar phenomenon in flight simulators are referred to as simulator sickness. In contrast head mounted video displays can present a fixed visual stimulus while the head rotates such that vestibular signals are not accompanied by the appropriate visual motion stimuli. These conflicts can also

result in symptoms of nausea, vertigo, and general malaise even when very small head movements occur.⁹⁴⁻⁹⁶ Visual-vestibular conflicts also occur with telepresence systems that present images that are either magnified or minified as a result of differences between the camera distance to objects and the viewer distance to the viewscreen. As a result, head movements that cause camera movements will stimulate excessive retinal image motion with magnification or too little retinal image motion with image minification. The mismatch will cause adaptive adjustments of the VOR and possible simulator sickness.

The oculomotor system responds to this visual-vestibular conflict during steady fixation by suppressing the VOR and adaptively reducing its gain. Thus, a user viewing a 2-h motion picture on an HMD while making a normal range of head movements during conversation with other persons will adapt their VOR to be less responsive to head motion. This adaptation can produce postimmersion symptoms where the adapted VOR remains suppressed such that head movements without the HMD will not stimulate a sufficient vestibular ocular movement to stabilize the retinal image. The user will perceive the world to move with his head and this can cause nausea, vertigo, and malaise. Attempts to move images with the head using a head tracking system that updates image position during head rotation have been unsuccessful⁹⁷ because of small delays in image processing that desynchronize expected image motion and head movements. It is possible that an image could be blanked during large head movements and stabilized with a head tracker during small head movements. The blanking would be similar to the normal suppression of vision during saccadic eye movements. Other solutions are to use see-through and see-around virtual displays so that real objects in the periphery are always visible during normal head movements so that some visual and vestibular stimuli remain coherent.

Spatial Location Conflicts

No existing HMD has the wide field of view of the human visual system in a natural environment. The horizontal field of view in currently available HMDs ranges from 22.5° to 155°.⁹⁸ The larger fields of view in binocular systems are expanded by partially overlapping the two monocular images. Note that the entire monocular and binocular portions of the field of view lie within the normal 114 to 120° binocular overlap of the normal visual system. Partial overlap in HMDs can be produced by displacing the right field to the right and left field to the left. Each eye has a monocular flanking component in its temporal visual field. It is also possible to reverse the displacement so that the monocular components are in the nasal-visual field of each eye (Fig. 12). These two organizations are consistent with two different occlusion aperture configurations. The former has the same organization as a binocular field seen on a near screen in front of a background. Thus the total image is seen as a discontinuous surface in depth with the central binocular portion at a proximal distance compared to the monocular portion that is seen farther away.⁶ The latter case has the same organization as a single continuous surface viewed through an aperture (porthole effect). The former system is more likely to cause binocular rivalry in the monocular regions of the field than is the latter approach. This suppression can cause a moonlike crescent shape (luning) in the flanking monocular region that corresponds to the border of the image seen by the contralateral eye and this blocks visibility of the monocular field crescent seen by the ipsilateral eye. The luning can be reduced by blurring the edges of each eye's image borders or by using the crossed view configuration. Interestingly, the latter crossed-view approach is preferable in head mounted displays.⁹⁹

Visual suppression and binocular rivalry suppression are also stimulated in monocular image displays. Rivalry occurs between the image presented to one eye and the natural scene presented to the other. Even if one eye is occluded, the dark image of the patched eye will rival with the visible image of the open eye. This rivalry can be minimized by presenting the patched eye with a uniform field with the same mean luminance and color as the open eye image. This can be accomplished easily by placing a diffusing filter over the nonseeing eye. In cases where the two eyes view independent images, it might be possible to control which eye is suppressed by moving the image in the seeing eye and keeping it stationary in the nonseeing eye. This technique is based on the observation that motion is a dominant feature for binocular rivalry.

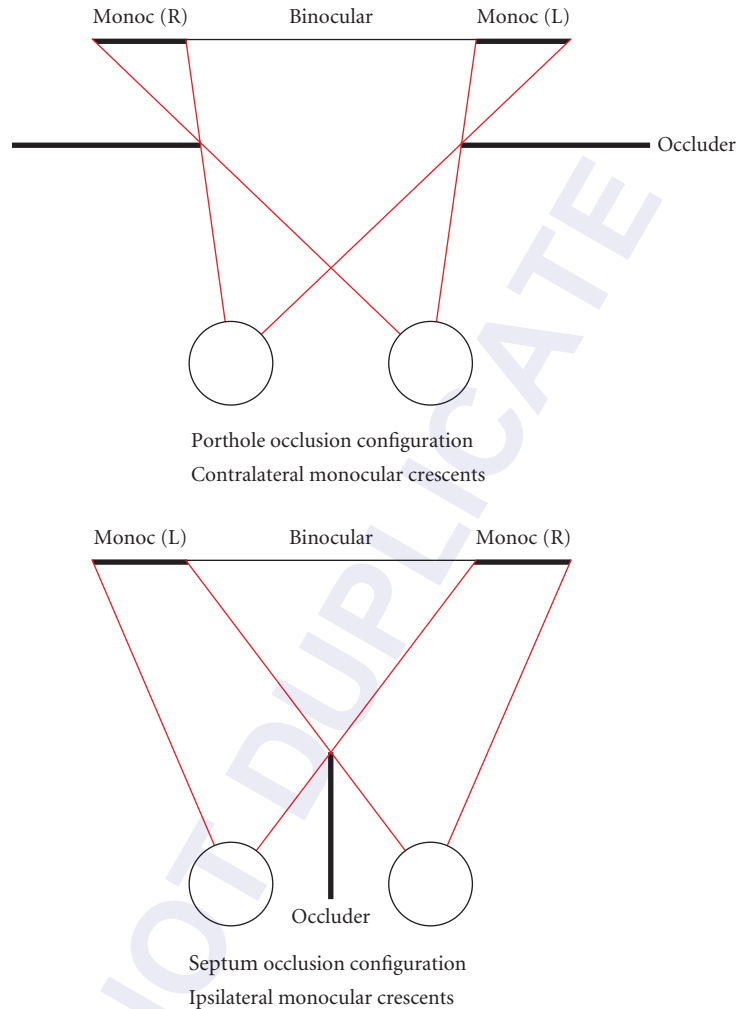


FIGURE 12 Occlusion figure with port hole and septum. Partial overlap of the visual field simulates obstructed views by a port-hole aperture or septum in the mid-sagittal plane.

Optical Errors

The success of any binocular viewing instrument depends in large part on the precision of alignment and quality of distortion free optics. Alignment errors can produce vertical, horizontal, and torsional disparities. Of these, the vertical errors are the most troublesome. There is limited range for vertical vergence and slow response times.⁷⁷ The oculomotor system can adapt its vertical alignment posture (vertical phoria) if the vertical disparity is constantly present; however, if the HMD is worn intermittently then no adaptation will occur. Horizontal errors, especially in the convergence direction, can be overcome easily by vergence eye movements. Torsional errors cause overall torsional disparities which can cause perception of slant about the horizontal axis if the virtual scene is superimposed

by see-through optics on the natural visual environment. However, if an opaque visor is used so that only the virtual display is visible, slant distortion will not be perceived with cyclotorsion disparity.¹⁰⁰ This is because the visual system only interprets slant when there is horizontal shear disparity between vertical lines but no vertical shear disparity between horizontal lines. When both forms of shear are present, this is interpreted as resulting from eye rotation rather than real object slant. Cyclodisparity in large field displays will stimulate cyclovergence responses that can have postimmersion aftereffects in the form of cyclophoria.⁸³ Errors in the baseline separation of the entrance pupils of the optical system can produce divergent disparities if it is set wider than the interpupillary distance. The oculomotor system has limited divergence ability and this stimulus can cause headaches and visual discomfort. Finally, low quality optics that have distortions such as barrel, pincushion, and keystone will introduce horizontal, vertical, and torsional disparities in both biocular and binocular displays if the distortions are not identical for the two eyes or if the visual fields are displaced slightly in the two eyes. Displacements could result from alignment errors or from partial field overlap used to expand the field of view. The disparities produced by these distortions can cause idiosyncratic warpage of the virtual image and vertical disparities that stress the oculomotor system.

13.13 REFERENCES

1. G. L. Walls, *The Vertebrate Eye and Its Adaptive Radiation*, Hafner, New York, 1967.
2. H. W. Hofstetter, "Graphical Analysis," Chapter 13 in *Vergence Eye Movements: Clinical and Basic Aspects*, C. M. Schor and K. Ciuffreda (eds.), Butterworth, Boston, 1983.
3. W. Dember and J. Warm, *Psychology of Perception*, Holt, Rinehart, and Winston, New York, 1979.
4. B. Cumming, "Motion—In Depth," Chapter 12 in *Visual Detection of Motion*, A. T. Smith and R. J. (eds.), Snowden Academic Press, San Diego, CA, 1994.
5. W. Epstein, (ed.), *Stability and Constancy in Visual Perception: Mechanisms and Processes*, John Wiley and Sons, New York, 1977.
6. K. Nakayama and S. Shimojo, "DaVinci Stereopsis: Depth and Subjective Occluding Contours from Unpaired Image Points," *Vis. Res.* **30**:1811–1825 (1990).
7. M. G. Harris, "Optic Retinal Flow," Chapter 2 in *Visual Detection of Motion*, A. T. Smith and R. J. Snowden (eds.), Academic Press, San Diego, CA, 1994.
8. B. Rogers and B. J. Graham, "Similarities between Motion Parallax and Stereopsis in Human Depth Perception," *Vis. Res.* **22**:261–270 (1982).
9. M. E. Ono, J. Rivest, and H. Ono, "Depth Perception as a Function of Motion Parallax and Absolute-Distance Information," *J. Exp. Psychol: Human Perception and Performance.* **12**:331–337 (1986).
10. L. R. Harris, "Visual Motion Caused by Movements of the Eye, Head and Body," Chapter 14 in *Visual Detection of Motion*, A. T. Smith and R. J. Snowden (eds.), Academic Press, San Diego, CA, 1994.
11. D. Regan and K. I. Beverley, "Visually Guided Locomotion: Psychophysical Evidence for a Neural Mechanism Sensitive to Flow Patterns," *Science* **205**:311–313 (1979).
12. W. H. Warren and D. J. Hannon, "Direction of Self-Motion is Perceived from Optical Flow," *Nature* **336**:162–163 (1988).
13. J. A. Crowell and M. S. Banks, "Ideal Observer for Heading Judgments," *Vis. Res.* **36**:471–490 (1996).
14. D. N. Lee, "A Theory of Visual Control of Braking Based on Information about Time to Collision," *Perception* **5**:437–459 (1976).
15. S. P. McKee and S. Watamaniuk, "The Psychophysics of Motion Perception," Chapter 4 in *Visual Detection of Motion*, A. T. Smith and R. J. Snowden (eds.), Academic Press, San Diego, CA, 1994.
16. T. Brandt, J. Dichgans, and W. Koenig, "Differential Central and Peripheral Contributions to Self-Motion Perception," *Exp. Brain Res.* **16**:476–91 (1973).
17. K. I. Beverley and D. Regan, "Evidence for the Existence of Neural Mechanisms Selectively Sensitive to the Direction of Movement in Space," *J. Physiol. Lond.* **235**:17–29 (1973).
18. E. Hering, *Spatial Sense and Movements of the Eye* (in German), C. A. Radde (trans.), American Academy of Optometry, Baltimore, 1879.

19. I. P. Howard, *Human Visual Orientation*, John Wiley and Sons, New York, 1982.
20. R. van Ee and C. J. Erkelens, "Stability of Binocular Depth Perception with Moving Head and Eyes," *Vis. Res.* **36**:3827–3842 (1996).
21. C. W. Tyler and A. B. Scott, "Binocular Vision," in *Physiology of the Human Eye and Visual System*, R. Records (ed.), Harper and Row, Hagerstown, 1979.
22. K. N. Ogle, *Researches in Binocular Vision*, Hafner, New York, 1964.
23. H. V. Helmholtz, *Handbuch der Physiologischen Optik*, 3d German ed. (1962) English translation by J. P. C. Southall, (trans.), 1909.
24. J. Foley, "Binocular Distance Perception," *Psychol. Rev.* **87**:411–434 (1980).
25. J. Garding, J. Porrill, J. E. W. Mayhew, and J. P. Frisby, "Stereopsis, Vertical Disparity and Relief Transformations," *Vis. Res.* **35**(5):703–722 (1995).
26. J. E. W. Mayhew and H. C. Longuet-Higgins, "A Computational Model of Binocular Depth Perception," *Nature* **297**:376–378 (1982).
27. B. Gillam and B. Lawergren, "The Induced Effect, Vertical Disparity, and Stereoscopic Theory," *Perception and Psychophysics* **34**:121–130 (1983).
28. L. Liu, S. B. Stevenson, and C. M. Schor, "A Polar Coordinate System for Describing Binocular Disparity," *Vis. Res.* **34**(9):1205–1222 (1994).
29. B. J. Rogers and M. F. Bradshaw, "Vertical Disparities, Differential Perspective and Binocular Stereopsis," *Nature* **361**:253–255 (1993).
30. G. Westheimer and M. W. Pettet, "Detection and Processing of Vertical Disparity by the Human Observer," *Proc. R. Soc. Lond. B Biol. Sci.* **250**(1329):243–247 (1992).
31. B. T. Backus, M. S. Banks, R. van Ee, and J. A. Crowell, "Horizontal and Vertical Disparity, Eye Position, and Stereoscopic Slant Perception," *Vis. Res.* **39**:1143–70 (1999).
32. R. M. Steinman and H. Collewijn "Binocular Retinal Image Motion during Active Head Rotation," *Vis. Res.* **20**:415–429 (1980).
33. D. Cline, H. W. Hofstetter, and J. R. Griffin, *Dictionary of Visual Science*, 4th ed. Chilton, Radnor, PA, 1989.
34. G. Westheimer, "Effect of Binocular Magnification Devices on Stereoscopic Depth Resolution," *J. Opt. Soc. Am.* **46**:278–280 (1956).
35. K. Stevens, "Surface Perception from Local Analysis of Texture and Contour," Ph.D. thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, 1979.
36. K. N. Ogle and P. Boder, "Distortion of Stereoscopic Spatial Localization," *J. Opt. Soc. Am.* **38**:723–733 (1948).
37. C. Bourdy, "Aniseikonia and Dioptric Elements of the Eye," *J. Am. Opt. Assn.* **39**:1085–1093 (1968).
38. R. S. Arner, "Eikonometer Measurements in Anisometropes with Spectacles and Contact Lenses," *J. Am. Optom. Assn.* **40**:712–715 (1969).
39. L. Rose and A. Levenson, "Anisometropia and Aniseikonia," *Am. J. Optom.* **49**:480–484 (1972).
40. S. Awaya and G. K. von Noorden, "Aniseikonia Measurement by Phase Difference Haploscope in Myopia, Anisometropia, and Unilateral Aphakia (with Special Reference to Knapp's Law and Comparison Between Correction with Spectacle Lenses and Contact lenses)," *J. Jpn. Contact Lens Soc.* **13**:131 (1971).
41. A. Bradley, J. Rabin, and R. D. Freeman, "Nonoptical Determinants of Aniseikonia," *Invest. Ophthalmol. Vis. Sci.* **24**:507–512 (1983).
42. B. Winn, C. A. Ackerley, F. K. Brown, J. Murray, J. Prars, and M. F. St. John, "Reduced Aniseikonia in Axial Anisometropia with Contact Lens Correction," *Ophthalmic Physiol. Opt.* **8**:341–344 (1987).
43. R. M. Burnside and C. Langley, "Anisometropia and the Fundus Camera," *Am. J. Ophthalmol.* **58**:588–594 (1964).
44. E. Engle, "Incidence and Magnitude of Structurally Imposed Retinal Image Size Differences," *Percept. Motor Skills* **16**:377–384 (1963).
45. B. Julesz, *Foundations of Cyclopean Perception*, Chicago University of Chicago Press, Chicago, 1971.
46. C. B. Blakemore, "A New Kind of Stereoscopic Vision," *Vis. Res.* **10**:1181–1199 (1970).
47. A. Fiorentini and L. Maffei, "Binocular Depth Perception without Geometrical Cues," *Vis. Res.* **11**:1299–1305 (1971).

48. C. W. Tyler and E. Sutter, "Depth from Spatial Frequency Difference: an Old Kind of Stereopsis?" *Vis. Res.* **19**:858–865 (1979).
49. C. M. Schor, I. C. Wood, and J. Ogawa, "Spatial Tuning of Static and Dynamic Local Stereopsis," *Vis. Res.* **24**:573–578 (1984).
50. C. M. Schor and T. Heckmen, "Interocular Differences in Contrast and Spatial Frequency: Effects on Stereopsis and Fusion," *Vis. Res.* **29**:837–847 (1989).
51. M. von Rohr, *Die Brille als Optisches Instrument*, Wilhelm Englemann, Leipzig, pp. 172, 1911.
52. H. Ergleget, "Über Brillenwirkungen," *Zschr. F. Ophth. Optik* **3**:170–183 (1916).
53. J. Nolan, A. Hawkswell and S. Becket, "Fusion and Stereopsis in Aphakia," in *Orthoptics: Past Present and Future*, S. Moore, J. Mein, and L. Stockbridge (eds.), Shatten Int. Cont. Med. Book Co., New York, pp. 523–529 1975.
54. L. Marran, and C. M. Schor, "Lens Induced Aniso-Accommodation," *Vis. Res.* **38**(22):3601–3619 (1998).
55. C. M. Schor, L. Landsman, and P. Erickson, "Ocular Dominance and the Interocular Suppression of Blur in Monovision," *Am. J. Optom. and Physiol. Optics* **64**:723–736 (1987).
56. M. J. Collins and A. Goode, "Interocular Blur Suppression and Monovision," *Acta Ophthalmologica* **72**: 376–380 (1994).
57. H. Wallach and K. J. Frey "Adaptation in Distance Perception Based on Oculomotor Cues," *Perception and Psychophysics* **11**:77–83 (1972).
58. D. A. Owens and H. W. Leibowitz, "Accommodation, Convergence and Distance Perception in Low Illumination," *Am. J. Optom. Physiol. Opt.* **57**:540–50 (1980).
59. S. M. Ebenholtz, "Hysteresis Effects in the Vergence Control System: Perceptual Implications," in D. F. Fisher, R. A. Monty, and J. W. Senders (eds.), *Eye Movements: Visual Perception and Cognition*, Erlbaum, Hillsdale, NJ, 1981.
60. D. W. McCreedy, "Size, Distance Perception and Accommodation Convergence Micropsia. A Critique," *Vis. Res.* **5**:189–206 (1965).
61. T. Morita and N. Hiruma, "Effect of Vergence Eye Movements for Size Perception of Binocular Stereoscopic Images," in *Proc. Third International Display Workshops*, Kobe, Japan, SID, 1996.
62. S. R. Ellis, U. J. Bucher, and B. M. Menges, "The Relationship of Binocular Convergence and Errors in Judged Distance to virtual Objects," in *Proc. International Federation of Automatic Control*, Boston, MA (1995).
63. Y. Y. Yeh, "Visual and Perceptual Issues in Stereoscopic Color Displays," Chapter 4, in *Stereo Computer Graphics and Other True 3D Technologies*, D. F. McAllister (ed.), Princeton University Press, Princeton, NJ, 1993.
64. A. Gonshor and G. Melvill-Jones, "Extreme Vestibulo-Ocular Adaptation Induced by Prolonged Optical Reversal of Vision," *J. Physiol. (Lond.)* **256**:381–414 (1976).
65. G. M. Gauthier and D. A. Robinson, "Adaptation of the Human Vestibular Ocular Reflex to Magnifying Lenses," *Brain Res.* **92**:331–335 (1975).
66. R. M. Burde, "The Extraocular Muscles," in *Adler's Physiology of the Eye: Clinical Applications*, R. A. Moses (ed.), St. Louis, CV Mosby, pp 84–183 (1981).
67. B. Bagolini, "Part II Sensorial–Motor Anomalies in Strabismus (Abnormal Movements)," *Doc. Ophthalmol.* **41**:23–41 (1976).
68. R. K Lunenburg, *Mathematical Analysis of Binocular Vision*, Princeton University Press, Princeton, NJ, (1948).
69. D. Nguyen, I. Vedamurthy, and C. M. Schor, "Cross-Coupling Between Accommodation and Convergence Is Optimized for a Broad Range of Directions and Distances of Gaze," *Vis. Res.* **48**(7):893–903 (2008).
70. C. M. Schor, J. Alexander, L. Cormack, and S. Stevenson, "A Negative Feedback Control Model of Proximal Convergence and Accommodation," *Ophth. and Physiol. Optics.* **12**:307–318 (1992).
71. J. S. Friedenwald, "Diagnosis and Treatment of Anisophoria," *Arch. Ophth.* **15**:283–304 (1936).
72. A. L. Perlmutter and A. Kertesz, "Measurement of Human Vertical Fusional Response," *Vis. Res.* **18**:219–223 (1978).
73. W. A. Houtman, T. H. Roze, and W. Scheper, "Vertical Motor Fusion," *Doc. Ophthalmol.* **44**:179–185 (1977).
74. A. Romole "Dynamic versus Static Aniseikonia," *Aust. J. Optom.* **67**:108–113 (1984).
75. V. Ellerbrock and G. Fry, "Effects Induced by Anisometropic Corrections," *Am. J. Optom.* **19**:444–459 (1942).

76. P. L. Cusick and H. W. Hawn, "Prism Compensation in Cases of Anisometropia," *Arch Ophthalmol.* **25**:651–654 (1941).
77. D. B. Henson and B. Dharamski, "Oculomotor Adaptation to Induced Heterophoria and Anisometropia," *Invest. Ophthalmol. Vis. Sci.* **22**:234–240 (1982).
78. C. M. Schor, J. Gleason, and D. Horner, "Selective Nonconjugate Binocular Adaptation of Vertical Saccades and Pursuits," *Vis. Res.* **30**(11):1827–1845 (1990).
79. L. C. Morrison, "Stereoscopic Localization with the Eyes Asymmetrically Converged," *Am. J. Optom.* **54**:556–566 (1977).
80. J. M. Enoch, "A Spectacle-Contact Lens Combination used as a Reverse Galilean Telescope in Unilateral Aphakia," *Am. J. Optom.* **45**:231–240 (1968).
81. J. R. B. Charnwood, *An Essay on Binocular Vision*, London, Hatton, pp. 73–88, 1950.
82. H. B. Field "A Comparison of Ocular Imagery," *Arch. Ophthalmol.* **29**:981–988 (1943).
83. J. S. Maxwell and C. M. Schor, "Adaptation of Ocular Torsion in Relation to Head Position," *Vis. Res.* **36**(8):1195–1205 (1999).
84. M. P. Keating, *Geometric, Physical and Visual Optics*, Butterworths, Boston, MA, 1988.
85. M. F. Land, "Fast-Focus Telephoto Eye," *Nature* **373**:658 (1995).
86. C. M. Schor and L. Task, "Effects of Overlay Symbolology in Night Vision Goggles on Accommodation and Attention Shift Reaction Time," *Aviation, Space and Environmental Medicine* **67**:1039–1047 (1996).
87. D. A. Southard, "Transformations for Stereoscopic Visual Stimuli," *Computers and Graphics* **16**:401–410 (1992).
88. M. Rosenfield and K. J. Ciuffreda, "Cognitive Demand and Transient Nearwork-Induced Myopia," *Optom. Vis. Sci.* **71**:381–385 (1994).
89. T. Inoue and H. Ohzu "Accommodation and Convergence When Looking at Binocular 3D Images," in *Human Factors in Organizational Design and Management III*, K. Noro and O. Brown, Jr. (eds), Elsevier Science Publishers, Amsterdam, pp. 249–252, 1990.
90. R. V. Kenyon, T. A. DeFanti, and J. T. Sandin, "Visual Requirements for Virtual-Environment Generation," *Society for Information Display Digest* **3**:357–360 (1995).
91. J. Norman and S. Ehrlich, "Visual Accommodation and Virtual Image Displays: Target Detection and Recognition," *Hum. Factors* **28**:135–151 (1986).
92. G. C. de Wit and R. A. E. W. Beek, "Effects of a Small Exit Pupil in a Virtual Reality Display System," *Opt. Eng.* **36**:2158–2162 (1997).
93. P. Erickson and C. M. Schor, "Visual Functions with Presbyopic Contact Lens Corrections," *Optometry and Visual Science* **67**:22–28 (1990).
94. R. L. Hughes, L. R. Chason, and J. C. H. Schwank, *Psychological Considerations in the Design of Helmet-Mounted Displays and Sights: Overview and Annotated Bibliography*, U.S. Government Printing Office, AMRL-TR-73-16, National Technical Information Service and U.S. Dept. of Commerce, Washington D.C., 1973.
95. P. A. Howarth, "Empirical Studies of Accommodation, Convergence and HMD Use," in *Proc. Hoso-Bunka Foundation Symposium: The Human Factors in 3D Imaging*, Hoso-Bunka Foundation, Tokyo, Japan.
96. P. A. Howarth and P. J. Costello "The Occurrence of Virtual Simulation Sickness Symptoms When an HMD was Used as a Personal Viewing System," *Displays* **18**:107–116 (1997).
97. T. Piantanida, D. K. Bowman, J. Larimer, J. Gille, and C. Reed, "Studies of the Field of View/Resolution Trade-Off in Virtual-Reality Systems," *Hum. Vis. Visual Proc. Digital Display III, Proc. SPIE* **1666**:448–456 (1992).
98. E. T. Davis, "Visual Requirements in HMDs: What Can We See and What Do We Need to See," Chapter 8 in *Head Mounted Displays*, J. E. Melzer and K. Moffett, (eds.) Optical and Electro-Optical Engineering Series, McGraw-Hill, New York, 1997.
99. J. E. Melzer and K. W. Moffitt, "Ecological Approach to Partial Binocular Overlap," *Large Screen Projection, Avionic and Helmet-Mounted Displays, Proc. SPIE* **1456**:124 (1991).
100. I. P. Howard and H. Kaneko, "Relative Shear Disparities and the Perception of Surface Inclination," *Vis. Res.* **34**:2505–2517 (1994).

OPTICS AND VISION OF THE AGING EYE

John S. Werner

*Department of Ophthalmology & Vision Science
University of California, Davis
Sacramento, California*

Brooke E. Scheffrin

*Department of Psychology
University of Colorado
Boulder, Colorado*

Arthur Bradley

*School of Optometry
Indiana University
Bloomington, Indiana*

14.1 GLOSSARY

Aphakia. An eye lacking a lens, usually as a result of cataract surgery.

Age-related macular degeneration (AMD). A disease that results in loss of vision in the central retina or macula due to damage to the photoreceptors and retinal pigment epithelium of the retina, often secondary to changes in the choroid (blood vessel layer) and Bruch's membrane. There are two major forms, called "wet," or exudative, and "dry," or atrophic. The dry form is more common, but the wet form is responsible for the majority of cases involving severe vision loss.

Cataract. An opacity or clouding of the lens resulting in a reduction in the amount of light transmitted to the retinal image and an increase in the light scatter. This is the leading cause of blindness worldwide, but can be treated successfully by removal of the lens.

Diabetic retinopathy. A disease of the eye caused by diabetes that leads to a proliferation of blood vessels. These blood vessels may swell and leak. When the resultant scar tissue shrinks the surrounding vitreous humor, it may cause the retina to detach, leading to loss of vision in the affected region of retina.

Glaucoma. A group of diseases, often but not always, associated with high intraocular pressure, resulting in death of ganglion cells.

Intraocular lens (IOL). A prosthetic lens usually implanted in the eye to replace the natural lens following its removal. IOLs may be placed in front of the pupillary plane (anterior chamber lenses) or

behind the pupil (posterior chamber lenses). They are manufactured from a variety of materials, the most common of which are polymethylmethacrylate, silicone, and hydrogels.

Phakia. An eye with a natural lens, from the Greek *phakos* or lens.

Presbyopia. Greek for “old sight,” is a condition in which an individual has lost the ability to accommodate sufficiently for near work. This results in the need for “reading glasses” around the age of 40 or 50 years.

Pseudophakia. An eye with an artificial or intraocular lens implanted to replace the natural lens following its removal.

14.2 INTRODUCTION

Senescence occurs in the optics of all human eyes and the resulting changes in the spatial and spectral distribution of light comprising the retinal image can have a profound impact on vision. These optical changes have enormous economic consequences for health care and society, and although some effective treatments exist, optical quality in the senescent eye remains inferior to that of younger eyes. Before describing the optical and visual changes that characterize normal aging, it is useful to consider the scope of the problem in light of demographic changes in the world’s aging population.

14.3 THE GRAYING OF THE PLANET

The proportion of the world’s population that is elderly is higher now than at any other point in recorded history. In the last two decades, the number of individuals above age 65 years has increased at more than twice the rate as that of individuals below 65 years.¹ At present, every month, another 1 million of the world’s population enters its sixth decade of life. As a result, the number of individuals age 60 and above is likely to change from one in ten today to one in five by 2050, and one in three by the year 2150. As illustrated by Fig. 1, the global growth rate for individuals over age 80 is occurring

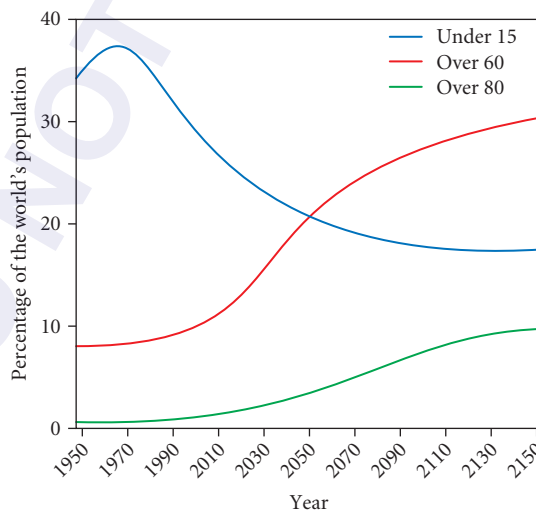


FIGURE 1 Percentage of the global population projected from 1950 to 2150 for three age groups based on statistics provided by the United Nations.²

at an even faster rate. Indeed, the number of centenarians in the United States has doubled every 10 years since 1960. At the same time, there is a corresponding decline in the proportion of the population below 15 years of age.

Many factors have contributed to this remarkable demographic shift in the age distribution, including a change in the fertility rate, improved hygiene, nutrition, disease control, and education. These trends are generally more apparent in less developed countries, and the demographic shift in the elderly population is now more rapid in these parts of the world. For example, in 1995 there was a global increase of 12 million persons at or above age 60 and nearly 80 percent of it was due to a change in less developed countries. The countries with the fewest resources at present will thus have less time to cope with the social and economic consequences of an elderly population. It is not clear how long these present demographic shifts will continue into the future and when they will reach equilibrium.

Some Implications of Living Longer

Some theoretical calculations place the maximum human life span at about 120 years, but there are reasons to think that this may not apply to the eye. The demographic shifts in longevity of the eye are apparently not keeping pace with that of the rest of the body. Research to Prevent Blindness estimates that the number of blind individuals will double between 1999 and 2030. If true, the social, health, and economic consequences will be profound. One eminent authority put it this way: “Mathematically speaking, it is permissible to say that those who prolong life promote blindness” (Ref. 3, p. 306).

Health Implications and the Problem of Defining Normal Aging Current projections for the future are that as the baby boomers (the cohort of 76 million Americans born between 1946 and 1964) start to turn age 65 in the year 2011, they will be better educated, have better health, and more financial resources than did their parents when they retired. All indications are that retiring baby boomers will read more and want to spend more time in outdoor leisure activities than their parents did at comparable ages. For the age group over 65 years of age, there are also some indications that the rates of disability and disease are decreasing compared to previous generations, but this is not so for vision. Older adults experience a higher rate of visual disability and ocular disease than younger individuals, and the rate accelerates after about 40 years of age. For example, glaucoma is about 8 times more common for individuals above age 65 than in the general population. In the Framingham study,⁴ 50 percent of participants age 75 and over had at least one of four visual “diseases:” age-related macular degeneration, cataract, diabetic retinopathy, or glaucoma.

Many studies of “normal” aging eyes fail to include significant numbers of eyes older than 80 years (a notable exception being Haegerstrom-Portnoy⁵). Thus, most of the aging effects described in this chapter fail to describe very old (80 to 100 years) eyes. In addition, most studies of aging attempt to characterize the visual system with many factors controlled, such as by testing with optimal refraction rather than habitual refraction. These controls are essential for separating the factors underlying normal aging, but they necessarily overestimate the quality of vision experienced outside the laboratory.

There is little agreement about where to draw the line between senescent degradation in vision and some forms of ocular pathology. For many researchers, normal aging is defined by exclusion—the changes in vision that occur when overt ocular disorder and disease are not detected. An alternative distinction between aging and disease is based on the amount of tissue changed or when visual loss exceeds a criterion. Both of these definitions are used to define cataract, a loss of transparency of the lens resulting from cumulative changes in products of oxidation that can dramatically elevate intra-ocular light scatter. There are several different types of cataract that can degrade the optical quality of images formed on the retina and consequently undermine the quality of the “neural” image transmitted by the retina. The biochemical events leading up to these lenticular changes occur continuously across the life span so the presence of this disorder is largely a matter of where the criterion is set on the continuum of normal aging.

Most individuals will still be able to approach their golden years without frank visual pathology. However, no individual should expect to have the quality of vision experienced in youth. Indeed, perhaps 20 years prior to retirement, most individuals already will have experienced the need for

“reading glasses” to correct near vision due to the senescent reduction in or loss of accommodation. Some problems of optical origin can be corrected or compensated by optical means, but additional losses in vision may be expected due to neural changes for which there are no known remedies. Some current research has focused on factors that might contribute to normal aging of the eye that could be controlled by life style choices such as diet⁶ or exposure to ultraviolet⁷ and short-wave visible⁸ radiation. Whatever the cause, any impairment of our dominant sensory system will inevitably have behavioral consequences. In some cases, this may be an inconvenience such as slower reading or the need for better lighting, while in other cases it may lead to more serious problems such as accidents, falls, and greater dependence on others for daily activities. For many, senescent changes in vision will become the dominant factor limiting their quality of life.⁹

Economic and Social Implications of Aging Eyes There is a gradual reduction in accommodative amplitude throughout adulthood resulting in a complete loss of accommodation during the sixth decade of life. Combining this information with the fact that each day 10,000 people in the United States reach the age of 50, there is little wonder that ophthalmic corrections for reading and near work will account for approximately half the revenues in the multibillion dollar ophthalmic lens industry.

Visual impairment secondary to diseases of the eye also creates a sizeable economic impact on our society. For example, visual impairment and blindness for individuals above age 40 have an estimated cost of \$51 billion per annum in the United States alone.¹⁰ In addition, approximately 2.8 million cataract surgeries are performed each year in the United States at an annualized cost of about \$25 billion. The costs of presbyopic spectacle correction may be even higher. These numbers are almost certain to increase due to the increasing proportion of the population that is elderly.

Additional economic consequences of our aging population will be created because the elderly generally require more special services, housing needs, and medical care, often as a result of visual impairment or blindness. For example, vision impairment increases the risks of falls by sixfold.¹¹ Senescent changes in vision are also likely to impose new burdens on public transportation systems in view of the fact that about half the world’s elderly population today live in cities, and that proportion is expected to increase to about three-fourths, at least in more developed countries. New demands will be made for housing these elderly individuals in order to promote independent living.¹² Being able to care for oneself is not only important for self-esteem, but it also has enormous economic impact. Architects and lighting designers will have to take into account age-related changes to the visual system in planning for new environments inhabited by the elderly. Current standards for lighting are almost entirely based on college-age students, although there have been new calls for explicit attention to the visual needs of aged individuals.¹³ The goal of such standards would be to find conditions that can compensate for sensory losses, to be described later, while not introducing conditions that promote “disability glare.” Such lighting considerations might promote greater safety in the home to facilitate independent mobility and support communication through computers and reading and other visual tasks involving hobbies and work.

14.4 SENESCENCE OF THE EYE’S OPTICS

Senescent changes influencing image formation are prevalent in all optical components of the eye. The following sections describe some of these changes in detail.

Senescent Changes in the Optical Components of the Eye

Tears The anterior surface of the eye is bathed by tears that produce an optically smooth air-tear interface mirroring the global shape of the cornea. The refractive index of the tears approaches that of the underlying corneal epithelium, therefore, the optical impact of the rough corneal epithelium surface is essentially eliminated by index matching. Because the tear film is so thin ($3\ \mu\text{m}$)¹⁴ it is vulnerable to evaporation effects, and it is inherently unstable. Stability and uniformity are enhanced by a

surface monolayer of lipids and blinking. Insufficient blinking and/or inadequate tear volume or lipid coating can cause the uniform tears to “break-up” into a patchwork of tears leaving regions of dehydrated mucins and the optically rough epithelium exposed, resulting in increased optical scatter and aberrations¹⁵ and reduced visual performance.¹⁶ The presence of dry eye is an indicator of inadequate tears, and the prevalence of dry eye syndrome increases with age.¹⁷ With increasing age, tear production by the accessory and main lacrimal glands decreases,³ the lipid layer gets thinner,¹⁸ corneal innervation and sensitivity decline,¹⁹ and in combination, these factors lead to reduced tear production^{20,21} and blink rates. All of these factors may contribute to reduced optical quality of tears in older eyes, which may be most pronounced during activities such as reading which lowers blink rates.

Cornea The human cornea consists of three discrete layers separated by two basement membranes. The outermost epithelial layer, composed of basal, wing, and surface cells, acts as a barrier to the external environment and is subject to a rapid rate of renewal, with the result that its cells are replaced every 5 to 7 days.²² The middle layer, or stroma, is a matrix of collagen fibers as well as fibroblast-like cells called keratocytes and accounts for approximately 90 percent of the corneal thickness. The stroma is sandwiched between two membranes. Anterior and posterior to the stroma are Bowman’s and Decemet’s membrane, respectively. The transparency of the stroma is due to the uniform density and subwavelength spacing of collagen fibers, and any disruption of this order caused by such factors as increased hydration will cause light scatter and loss of transparency. Because of the tendency of the stroma to imbibe water there must be a way to regulate corneal hydration. The extent of corneal hydration is largely controlled by energy-dependent-ion coupled fluid transport mechanisms located at the innermost layer of the cornea, the corneal endothelium. The endothelium is made up of a single layer of cells. At birth, the corneal endothelium has about 4500 cells/mm², but the number is reduced to about 2000 to 2500 cells/mm² by age 70 or 80. Although the surviving cells become larger to maintain approximately constant total tissue volume, a smaller total number of fluid transport mechanisms remain in the corneal endothelium. As a result, proper corneal hydration is sometimes a problem in the elderly. The rate of decline in endothelial cell density is accelerated by cataract surgeries,²³ but for the typical eye sufficient endothelial cells remain throughout life.

The cornea is the most powerful refractive element in the human eye, accounting for about 40 to 45 D of optical power. Although it has a lower refractive index than the lens, the difference of refractive indices, which is directly proportional to dioptric power, is highest at the air-tear/cornea interface. The cornea is not spherical; its reduced curvature in the periphery helps reduce the spherical aberration of the cornea. The refractive index of the cornea²⁴ and curvature of its anterior surface²⁵ appear to undergo little change with age, but with aging the axis of corneal astigmatism loses the “with the rule” bias seen in younger eyes.^{26,27}

As can be seen from Fig. 2, the cornea absorbs most of the ultraviolet (UV) radiation below about 320 nm. Absorption of this radiation in sufficient intensity can damage the cornea (e.g., actinic keratitis), but long-term effects are rare because damage is typically confined to the corneal epithelium which is capable of rapid renewal. Stem cells in the corneal margins continue to generate epithelial cells, which then migrate across the cornea to generate the basal layer of the epithelium. This may be why the absorption spectrum of the cornea undergoes only small changes with advancing age.²⁸ Cumulative absorption of UV light may, however, lead to corneal abnormalities such as pterygium, a condition that involves the advancing growth of abnormal tissue from the border between the cornea and sclera, the limbus, onto the clear cornea. This increases in prevalence with increasing age.²⁹

Anterior and Posterior Chambers The cornea and lens are separated by a fluid-filled space. This fluid (the aqueous humor) maintains the metabolism of the cornea and lens, as the transparency requirements of these structures necessitate that they have no blood supply. The much larger volume between the cornea and iris is the “anterior” chamber, and the smaller volume between the iris and the lens is the “posterior” chamber. The aqueous humor is produced by the ciliary body around the margins of the posterior chamber, flows to the anterior chamber, and is drained through the trabecular meshwork and into Schlemm’s canal around the radial margins of the anterior chamber. The optical role of the aqueous, therefore, is twofold. It enables the avascular lens and corneal endothelium to function (and thus retain lens and corneal transparency) while maintaining its own

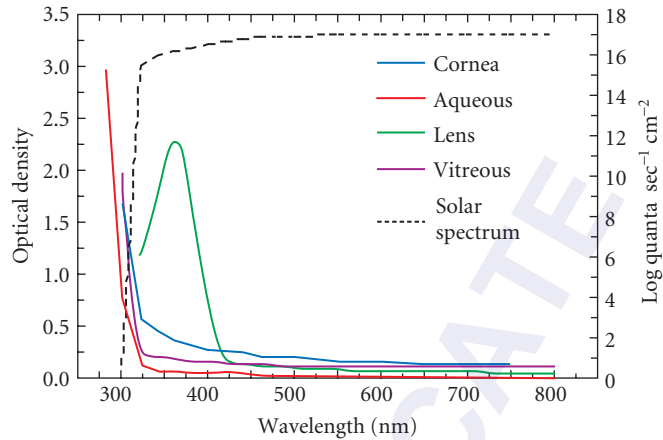


FIGURE 2 Optical density (decadic logarithm of the reciprocal of transmittance) of human ocular media is plotted as function of wavelength (left ordinates; from Boettner and Wolter³⁰). Log quanta in sunlight reaching the surface of the earth from an angle of 30° from the horizon is plotted as a function of wavelength (right ordinates; from Brandhorst et al.³¹). (After Werner.³²)

transparency. The optical density of the aqueous changes negligibly with age³⁰ and any damage sustained by the aqueous resulting from absorption of light is likely to be unimportant for senescence because it is rapidly replaced.²² Although there are pathological conditions that directly affect aqueous transparency (e.g., uveitis in response to infection), they are primarily transient and not related directly to age. There is a senescent reduction in both the production and drainage of aqueous,³³ and an increased likelihood of failure of the homeostatic production and drainage system. As drainage is unable to keep up with production, intraocular pressure rises, and glaucoma can ensue (see section “Glaucomas”). One optical side effect of glaucoma can occur when fluid in the corneal epithelial layer builds up between the basal epithelial cells,³⁴ causing colored diffraction halos similar to those generated by nocturnal eye closure, contact lens wear, or prolonged exposure to fresh water during swimming. A healthy epithelium maintains its high transparency by tight junctions between adjacent cells; however, as fluid builds up between the regular array of cells, classic wavelength-dependent diffraction occurs which is visible as colored halos. Glaucoma is generally associated with corneal epithelial edema only when pressures become very high, which is more typical of acute angle closure glaucoma. Interestingly, a much rarer condition in older eyes, ocular hypotony, can generate loss of corneal transparency because of stromal edema resulting from IOP being too low.³⁵ Both cases emphasize the critical role played by intraocular pressure in maintaining corneal hydration and thus corneal transparency.

Pupil The dilator and sphincter muscles of the iris control the size of the pupillary aperture that, for young adults, varies in diameter with the level of the prevailing illumination from approximately 2 to 8 mm. As in any optical system, the size of the aperture has important consequences for image formation. The smallest pupillary diameters can decrease optical quality through diffraction effects. The optimal pupil diameter is about 2.4 mm from the point of view of classic image quality.³⁶ Traditionally, larger pupils at lower light levels are considered to represent a compromise between improved sensitivity (more light entering the eye) and diminished image quality (more aberrations in peripheral optics). However, this view ignores the very significant role of photon noise on image quality.³⁷ Including photon noise in image calculations, Laughlin³⁸ argued that pupil dilation as light levels drop ensures optimal image quality. Employing a similar hypothesis (pupil size is adjusted to optimize retinal image quality), pupil miosis that accompanies senescence may also be considered as an adaptation to poorer optical quality in older eyes.

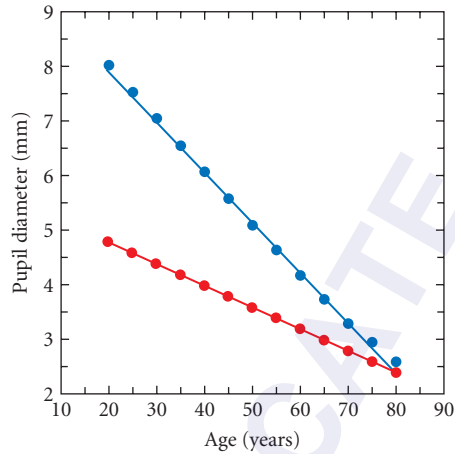


FIGURE 3 Mean diameter of the human eye pupil is plotted as a function of age. Blue and red symbols show dark- and light-adapted values, respectively. (Data from Kornzweig.³⁹)

Maximal pupil diameter occurs in adolescence and then progressively decreases with increasing age. As shown in Fig. 3, age-related change in diameter differs for dark-adapted and light-adapted conditions. Changes in pupil diameter with light level provide a mechanism of adaptation through adjustment of visual sensitivity with illumination. However, the change from maximum to minimum size influences retinal illuminance by only about a factor of 10, so the pupil plays only a small role in adaptation mechanisms that allow us to function over a range of ambient illumination covering about 12 log units.

Age-related changes in pupil size are noticeable in the elderly but actually begin in the early adult years. For example, when fully dark adapted, there is essentially a monotonic decrease in pupil diameter beyond 10 years of age.⁴⁰ Spatial vision depends on light level, and although a smaller pupillary diameter allows less incident light to fall onto the retina, it is unlikely that a smaller pupil accounts for significant age-related changes in spatial vision. For example, it has been demonstrated that age-related differences in contrast sensitivity under photopic and mesopic ambient light levels cannot be explained solely by differences in pupillary diameter.⁴¹ In fact, there actually are some benefits associated with an age-related reduction in pupil diameter. The smaller pupil of older persons diminishes age-related changes in higher-order aberrations⁴² and the optical transfer function of the eye at low luminance levels.⁴³ In addition, the smaller pupil in the elderly reduces the diameter of the retinal blur circle, thereby increasing the depth of focus by about 0.5 D between the ages of 30 and 70 years.⁴⁴ More significantly, since blur circle size (radians) = defocus (Diopters) * pupil diameter (m),⁴⁵ as pupil size is halved, blur circle size is also halved. For example, if defocus was 4 D, then the impact of halving the pupil would be the same as halving the defocus, or a 2-D change in image quality. Thus halving of the photopic pupil and the more than factor of 3 reduction in mesopic pupil diameters in older eyes will greatly reduce the level of blur. Senescent pupil miosis may be nature's attempt at minimizing the impact of presbyopia.

The Lens The lens consists of only one cell type, but its thickness increases throughout the entire life span due to the accumulation of cells to form a series of concentric layers.⁴⁶ This results in a continuous axial thickening and reduction in the depth of the anterior chamber.⁴⁷ The oldest cells are thus found in the center or nucleus (indeed these are the oldest cells in the body in that they are born early in gestation and are never replaced), while the youngest cells comprise the outer or cortical layers. As the cells of the lens age, they lose their organelles and become essentially metabolically inert.

Molecules that absorb UV and short-wave visible light are present in the newborn lens and their numbers increase continuously over the life span.⁴⁸ Thus, the curve for the lens shown in Fig. 2 can be, to a first approximation, scaled multiplicatively with age.⁴⁹ A number of studies have measured the optical density of the human lens as a function of age. Some have described these changes in density by a linear function of age^{49,50} while other evidence demonstrates a bilinear function due to an acceleration in the aging function after about 60 years.⁵¹ The precise form of the aging function is difficult to determine because of substantial individual differences at every age (on the order of 1.0 log unit at 400 nm). However, it should be noted that the average infant lens transmits about 25 times more light at 400 nm to the retina than the average 70-year-old eye!

Transparency/cataract Cataract refers to an opacity of the lens that decreases vision by a somewhat arbitrary criterion amount. An acceleration in the population mean rate of change in lens optical density beyond approximately 60 years of age could be due to the inclusion of subjects having cataract.⁵² However, it is not clear where normal aging ends and cataract begins because this increase in optical density continues throughout life. In his comprehensive monograph on this subject, Young⁵³ summarized a substantial literature which lead him to conclude that: “No discontinuity between senescent and cataractous changes can be detected at the molecular level in the human lens....The most sophisticated techniques of modern biophysical and biochemical analysis have so far failed to uncover any feature of the cataractous lens that suggests cataract is anything more than an advanced stage of normal aging” (Ref. 53, p. 56).

It should be added that there are several forms of cataract depending on the part of the lens affected (nuclear, cortical, posterior subcapsular) or etiology (congenital, age-related, chemically induced, trauma induced, disease induced). Losses of lenticular transparency can sometimes occur through occupational hazards associated with infrared light exposure such as in “glassblowers cataract” due to the conduction of radiant heat that is absorbed by the cornea and, to a lesser extent, the iris.⁵⁴ Most often, changes in the transparency of the lens are associated with cumulative exposure to ultraviolet radiation. Weale⁵⁵ has proposed that the specific form of cataract depends on genetic factors while the common etiology is most often cumulative exposure to solar radiation (see section “Life-Span Environmental Radiation Damage”).

Currently, the only effective treatment for cataract requires surgical removal of the lens and in more than 90 percent of such surgeries performed in the United States an intraocular lens (IOL) is implanted in its place. A properly selected IOL obviates the need to wear thick spectacle corrections or contact lenses for the hyperopia that would otherwise ensue. Most of these lenses now contain UV-absorbing chromophores, but this is a relatively recent design improvement.

Accommodation and presbyopia Accommodation refers to the ability of the eye to adjust its power to focus the images of objects located at varying distances in the environment. Contraction of the ciliary muscle located within the ciliary body of the eye allows the anterior, and to a lesser degree the posterior, surface of the lens to change shape thereby altering the dioptric power of the lens (see recent review by Charman⁵⁶). The amplitude of accommodation refers to the *dioptric* difference between the farthest and nearest distances in which an object remains in focus. Age-related loss in accommodative amplitude (Fig. 4) has been attributed to various causes including reduced ciliary body function, changes in the geometric relationships of the lens and the surrounding zonules and ciliary body based on progressive circumferential enlargement of the crystalline lens, but primarily it is due to changes in the elastic properties of the lens and to a lesser extent the capsule.^{56,57} The red symbols in Fig. 4 are from the Basel longitudinal study of individuals followed for 20 years, with different subjects starting the study at different ages.⁵⁸ These data agree well with the classic data of Donders⁵⁹ shown by blue symbols.

As the amplitude of accommodation decreases, a near correction for tasks such as reading may be required. This typically occurs within the fifth decade of life and is known as presbyopia (meaning “old sight”). Because the accommodative amplitude changes continuously, the definition of presbyopia and its onset are somewhat arbitrary. The two sets of data in Fig. 4, although separated in time by about a century, are from similar geographic locations. Age of onset of presbyopia depends on many factors including geographical latitude and average ambient temperature.³

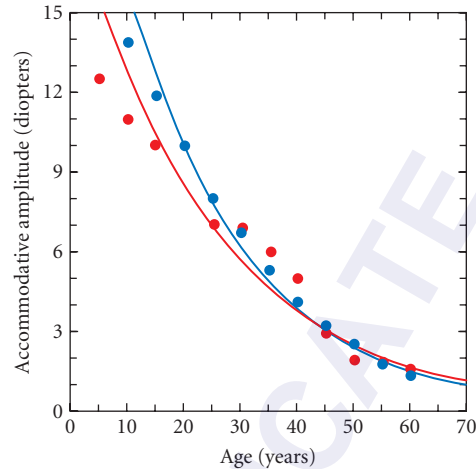


FIGURE 4 Amplitude of accommodation (diopters) is plotted as a function of age. Note that the tasks used for measurement involved both accommodation and depth of focus (which the clinical literature often calls pseudoaccommodation). This depends on both optical and neural factors. Red symbols (curve) from Brückner et al.⁵⁸ and blue symbols (curve) from Donders.⁵⁹ Curves are based on best-fitting exponential functions.

Vitreous Separating the lens from the retina is the vitreous chamber. This space, which occupies about 80 percent of the eye's volume, is filled with a colorless gelatinous substance called the vitreous humor. It is made up of about 99-percent water mixed into a gel-like mesh of collagen fibers and hyaluronic acid. Its role is primarily structural and not metabolic. In order to function it must retain its volume and remain transparent. As the vitreous ages, the gel is gradually replaced by unbound water. This liquefaction of the vitreous is responsible for the increased movement of the vitreous during and following eye movements in the aging eye.⁶¹ Local shrinkage of the vitreous gel is responsible for the increased prevalence of posterior vitreous detachment which can be present in up to 65 percent of eyes older than 65 years.⁶⁰ The liquefaction and shrinkage can lead to increased numbers and movement of floaters (visible shadows and diffraction phenomena associated with clusters of fibers and other optical inhomogeneities in the vitreous) consistent with increased movement of the vitreous in old eyes.⁶¹

Retina The optical role of the retina is complex. Its primary optical role is to absorb photons in the photopigment. However, the secondary optical role reflects needs for transparency, opacity, and waveguiding. Because the retina is "inverted," light must pass through the neural retina prior to reaching the photopigment. Thus, the neural retina must remain highly transparent. For example, it is for this reason that retinal ganglion cell axons are generally unmyelinated until they have exited the eye. Even though most of the retinal tissue is transparent, the blood supply to the postreceptoral neurons remains opaque. The lateral displacement of the foveal postreceptoral neurons removes the need for any retinal blood supply in the central fovea. Thus, most eyes have a foveal avascular zone with an accompanying increase in transparency and presumably image quality in the central fovea.

External to the layer of photopigment (which resides in the outer segments of the photoreceptors) lies a second layer of pigment, melanin, in the pigmented epithelium cells (melanin and hemoglobin in the choroid are also effective pigments). It has been suggested that melanin in the RPE has a primary antioxidant role protecting the retina from age-related oxidative damage.⁶² Optically, this pigment plays the same role as the black paint used to coat the inside of a camera; it reduces the

reflection of photons that are not absorbed by the photopigment. For example, at the peak sensitivity of the photopigment, about 2/3 of the photons entering the outer segment are absorbed by the photopigment (0.5 optical density), and the outer segments cover about 70 percent of the retinal area, and thus, slightly less than 50 percent of photons arriving at the retina will be absorbed. This percentage is lower in the peripheral retina where the outer segments are shorter, and will generally be lower due to photopigment bleaching (which never exceeds 50% under natural conditions), and of course will be lower at wavelengths other than the peak. The effectiveness of the photopigment and the pigment epithelium combine to improve image quality of the eye but make it difficult for clinicians to obtain high intensity images of the fundus as little light is reflected back out of the eye.

A third major pigment in the retina, most dense around the fovea, is a yellow pigment called the macular pigment (MP). It is located in the receptor fiber and inner plexiform layers⁶³ where it selectively absorbs short-wavelength light en route to the photoreceptors. Figure 5 shows the spectral absorbance of the foveal macular pigment for an average observer. The peak density is at 460 nm, near the peak sensitivity of the short-wave-sensitive cones (440 nm).

Several different methods (fundus photography, entoptic visualization, and psychophysics) indicate that the peak density of the MP is at the fovea and follows an exponential decay to negligible levels at approximately 5 to 10° retinal eccentricity.⁶⁵ There are large individual differences in the absorption of the macular pigment (from about 0 to 1.0 log unit at 460 nm), but there is little age dependency after early childhood.^{66,67}

While the function of the MP is still uncertain, it has been suggested that the presence of this pigment improves the optical image quality at the retina and may aid in maintaining the health of the macula. Because MP selectively absorbs short-wave light, it has been proposed that it helps to minimize the effects of chromatic aberration⁶⁸ thereby improving polychromatic image quality in the central retina. With regard to the health of the eye, MP is made up of carotenoids that tend to neutralize photosensitized molecules⁶⁹ and inhibit free radical reactions that are toxic to retinal cells. Thus, it has been suggested that the MP may protect the retina from actinic damage that may lead to age-related changes in visual sensitivity and to the development of age-related macular degeneration.⁸

A fourth, rhodopsin-based, visual pigment has recently been observed in some retinal ganglion cells. It has peak sensitivity around 460 nm, and it has been suggested that it may be responsible for light-induced melatonin suppression and thus maintenance of the circadian rhythm. In a recent

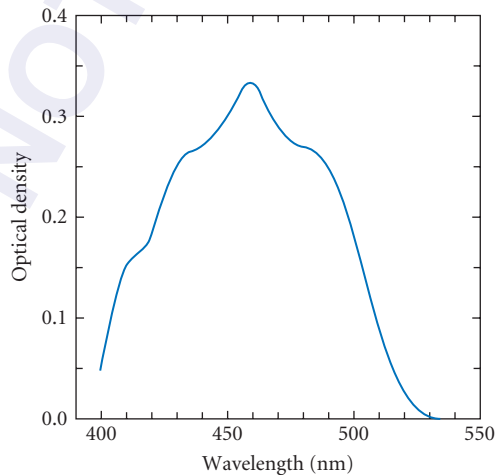


FIGURE 5 Optical density of the human macular pigment is plotted as a function of wavelength. (Data from Vos.⁶⁴)

analysis, Charman⁵⁶ showed that the combined effects of senescent miosis and lens yellowing would reduce the photon catch of this pigment in old eyes to only 10 percent of that of younger eyes and thus may contribute to age-related changes in circadian rhythms.

In addition to the transparency and pigmentation properties described above, the retinal photoreceptors have important optical waveguide properties. Photoreceptors must contain large numbers of photopigment molecules because large numbers of photons need to be captured, and photopigment regeneration is slow (half-life of 2 and 7 min for cones and rods, respectively). Also, in order to adequately sample the retinal image spatially, photoreceptors must be small. Typical photoreceptors contain long tubes filled with pigment with their axis aligned approximately perpendicular to the image plane (e.g., outer segments of foveal cones can be 60 μm long and only 2 μm in diameter). Such geometry creates fiber-optic or waveguide properties in these cells. These waveguide properties play an important role in reducing the impact of highly aberrated and scattered light so that they maximize the probability of photon catch for light entering in and around the pupil center [the Stiles-Crawford effect (SCE) see Chap. 8], and in so doing, enhance the optical quality of the retinal image.⁷⁰ A recent report⁷¹ has demonstrated that retinal Müller cells may also act as optical fibers funneling light through the retina to the photoreceptors to avoid scattering by the neural and vascular cells of the retina. No significant senescent changes occur in the SCE implying stability in photoreceptor orientation.^{72,73}

There are several senescent changes in the optics of the retina. The depth of the foveal pit decreases with increasing age,⁷⁴ while the pigmented epithelial cells accumulate lipofuscin throughout life as a by-product of phagocytosis of photoreceptor outer segments. Over age 70, lipofuscin and melanolipofuscin granules may occupy up to 20 percent to 33 percent of the free cytoplasmic space of the cell.⁷⁵ Also, between the ages of 10 and 90 years, there is a 2.5-fold decrease in melanin concentration within the retinal pigmented epithelial cells.⁷⁶ All of these changes can alter the final destination and distribution of photons within the retina. In addition, these cellular alterations are manifest in an age-related increase in lipofuscin autofluorescence,⁷⁷ which may serve as a possible marker for cellular aging, oxidative damage, and dysfunction associated with retinal disease.⁷⁸ The combined impact of cellular senescence on the retinal optics is unknown, but it likely contributes to the increased scatter observed in older eyes.⁷⁹

Eye Size The globe of the eye grows rapidly in the first year of life and then more slowly until about 5 years of age when it reaches an asymptotic sagittal length of about 23.5 mm. While the globe of the eye changes relatively little over the adult years, the lens continues to grow. These and other structural changes alter the refractive state of the eye throughout the life span. The coordinated growth of eye size and optical power in many eyes enables a focused retinal image without need for accommodation or refractive treatments. Such eyes are said to be emmetropic, but emmetropia is not the norm during all stages of development and senescence. Infant eyes are almost always hyperopic but during the first 2 years of life they shift toward emmetropia. For many older children and adults, a further slow shift toward myopia relegates many to an adult life of near sightedness (e.g., recent estimates indicate that myopia rates can be as high as 71 and 84 percent in Hong Kong⁸⁰ and Taiwan,⁸¹ respectively), whereas in later life there is an overall tendency for the refractive error of the eye to shift in a hyperopic direction. Although the prevalence of emmetropia can be high in modern populations, both eyes tend to maintain highly correlated growth resulting in similar refractive states in the two eyes. Recent reports, however, indicate that this correlated growth pattern begins to fail as the eyes age,⁸² and can result in significant interocular differences in refractive error (anisometropia) among the population older than 70 years,⁸³ compared to children and younger adults.

Senescent Changes in Retinal Image Quality

Many factors act to alter the retinal image throughout the life span. In this section, their combined effects are described along with one aspect of the retinal image (chromatic aberration) that does not seem to be age dependent.

Intraocular Scatter The spatial distribution of light on the retina from a point source is described by the point spread function (PSF). The fraction of light transmitted from a point source in the forward direction is the total transmittance, while the fraction transmitted within a cone of a particular angle a , is the direct transmittance. The value of a is arbitrary, but typically is about 1° . The difference between the total and direct transmittances produces a “tail” in the PSF and is called stray light. However, because the photoreceptors are directionally selective (see Chap. 9), the optical PSF is not necessarily equivalent to the functional PSF. Westheimer and Liang⁸⁴ have used psychophysical methods to measure the PSF of younger and older observers for a broadband source. The total scattered light from $7'$ to 90° is about 5 times larger in a 69-year-old observer compared to a young adult. Put another way, 90 percent of the light from a point source falls within a region of $7'$ for a young subject, but only 50 percent of the light falls within $7'$ for the older subject. (This analysis of the PSF does not take into account differences in the density of the ocular media which would further increase the difference between young and old.)

A number of studies have shown that large particle scatter by the human ocular media increases across the life span (e.g., Wolf⁸⁵). The relation between light scatter and age has been described as nonlinear due to an acceleration in the function in middle age.⁸⁶ A more recent study⁸⁷ measured the integral of the scatter function of the eye from 2.2 to 90° . The results show little or no increase in scatter from 16 to 40 years of age but a rapid increase afterward (i.e., a factor of 5 increase at 65 years compared to young adults).

Stray light in the human eye is not strongly related to wavelength,^{88,89} implying that it is not due to scatter by small particles (smaller than the wavelength of light). A small wavelength dependence, greater at long than middle wavelengths (opposite to Rayleigh scatter, where scatter $= k\lambda^{-4}$), has been attributed to fundal reflections, especially in lightly pigmented eyes.⁹⁰ A recent study⁹¹ has proposed that intraocular scatter is a combination of small particle Rayleigh scatter and a red dominated scatter produced by light entering through the sclera. This latter source must pass through the choroidal blood and thus is red biased. The combination of these effects in lightly pigmented eyes produces scatter levels that vary little with wavelength. This can be important in evaluating stray light in the elderly because pigmentation levels in the RPE decrease with age. The lens seems to be the main source of age-related change in intraocular stray light.⁸⁵ The cornea produces light scatter as well, but the amount does not change significantly with age.³⁰

Scattered light in the elderly may also be due to the age-related accumulation of two fluorophores in the lens.⁹² These are molecules that emit light at one wavelength (e.g., 520 nm) when irradiated by another (e.g., 410 nm). As these emissions are scattered randomly, they would be expected to degrade visual resolution. Elliott et al.⁹³ demonstrated a small effect of fluorescence on resolution of low contrast targets.

Intraocular stray light reduces the contrast of the retinal image. Normal levels of scattered light do not affect color detection thresholds significantly, but high levels can affect chromatic saturation and estimates of color constancy.⁹⁴ The primary effect of intraocular stray light is to impair visual performance when lighting varies dramatically within a visual scene, such as when driving at night, because even small amounts of scatter from a high intensity source can dramatically reduce contrast of the retinal images of nearby low intensity targets. The effects of light scatter are sometimes referred to as disability glare. Discomfort glare, on the other hand, describes a subject's visual discomfort when high intensity glare sources are present. The level of discomfort glare shows large intersubject variability, but it is generally correlated with the amount of disability glare. The older eye is not necessarily more sensitive to glare, but it experiences more of it due to greater scatter.^{79,95} Low luminance and low contrast visual acuity exhibit marked deficits with aging, which is further exacerbated by glare (veiling auto headlights, and so on). With the exception of cataract surgery, there are no treatments for increased ocular scatter, and therefore older people suffer the visual consequences of scatter or modify their behavior to avoid its most serious consequences, for example, changing night driving behavior.

Monochromatic Aberrations Almost 95 percent of the wavefront error measured monochromatically is caused by the lower-order aberrations of defocus and astigmatism,⁹⁶ and even in well-refracted eyes these lower-order aberrations continue to dominate.⁹⁷ However, there has been considerable interest generated recently in the higher-order aberrations of the eye.

Although the typical older eye exhibits appreciable levels of higher-order aberration, this results less from a specific ocular surface than from a decreased ability of lenticular aberrations to compensate for corneal aberrations—in contrast to younger eyes. With aging there is a change in meridian of the corneal astigmatism, the cornea develops negative spherical aberration, and the corneal coefficients for coma change as well. The changes are not accompanied by compensatory lenticular changes, resulting in greater whole-eye aberrations.^{27,98} Thus, senescent changes in sign and not the magnitude of both the cornea and lens aberrations are the primary contributor to increased ocular aberrations seen in older eyes. The age dependence of this trend, however, is not completely clear for two reasons. First, population data show large variability between eyes, and most studies clearly show that the distribution of aberrations in older eyes is mostly overlapping with the distribution of younger eyes.^{26,99} Thus, many older eyes have young eye levels of higher-order aberrations. Second, although most studies generally indicate best optical quality in young adult eyes (20 to 30 age group), one shows that optical quality is maximum at around age 40.⁹⁹ Finally, there is a shortage of data from eyes above 65 years.

A quantitative description of the quality of an optical system is provided by its modulation transfer function (MTF), the relation between contrast attenuation and spatial frequency. The eye's MTF can be derived from *in vivo* measurements of the optical point spread function on the retina. Guirao et al.⁴³ derived the MTFs from carefully collected point spread functions of 60 observers ranging in age from 20 to 70 years while controlling for such important factors as pupil size and the refractive state of the eye. In Fig. 6, the blue and red curves represent MTFs for observers with mean ages of 24 and 63 years, respectively. Modulation is scaled logarithmically to highlight differences at middle and high spatial frequencies that are obscured on a linear plot. Optical performance decreases with increasing spatial frequency for both age groups, and there is an age-related loss in performance at all spatial frequencies. Results similar to these have also been reported by Artal et al.¹⁰⁰ The results in Fig. 6 may be attributed to optical aberrations but not ocular scatter because it is factored out in the measurement procedure. Inclusion of intraocular scatter would tend to produce a more or less uniform “halo” in the retinal image that further reduces the contrast of stimuli in the elderly eye.

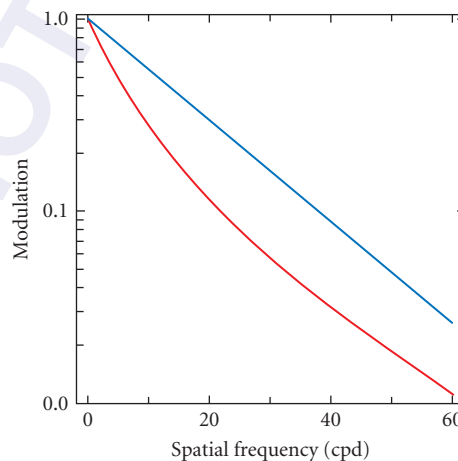


FIGURE 6 Modulation transfer functions for average individuals at 23 (blue curve) and 63 (red curve) years of age for a 3-mm pupil, 543-nm stimulus, calculated from equation and parameters from Guirao et al.⁴³

Perhaps the most important finding is that a smaller pupil can compensate for the age-related increase in monochromatic aberrations. Indeed senescent miosis may provide a complete compensation for the failed correction of corneal aberrations by lenticular aberrations. That is, one compensation mechanism evolves as the other disappears, and this does not include additional compensation mechanisms at a neural level.¹⁰¹ The bottom line is that the major loss in *optimal* retinal image quality with age will not be due to increased monochromatic or chromatic aberrations.

Chromatic Aberration The refractive index of the ocular media is wavelength dependent so that the spectral components of a broadband stimulus will not be imaged in the same plane within the eye. This is an example of chromatic aberration. Chromatic aberration can be appreciated from the colored fringes that can be observed in many telescopes and microscopes. Refractive indices of the optical elements of the eye decline as wavelength increases moving the long wavelength images farther from the cornea. The degree of chromatic aberration can be quantified in terms of the lens power required to image a particular wavelength on the retina when the eye is focused for a reference wavelength (change in refractive error as a function of wavelength). With a 578-nm reference, the aberration varies from -1.70 D at 405 nm to 0.40 D at 691 nm.¹⁰² Howarth et al.¹⁰³ (1988) measured the longitudinal chromatic aberration for four younger (27 to 33 years) and six older (48 to 72 years) participants using three different methods. All methods failed to show a change in axial chromatic aberration with age.

Interestingly, one suggestion for treating presbyopia employs a hyperchromatic lens design which can double the effective chromatic aberration of the eye. Experimental studies have shown that such a lens can increase the depth of field of the presbyopic eye.¹⁰⁴ Conversely, diffractive IOL designs have been proposed that can cancel the refractive chromatic aberration of the human pseudophakic eye.^{105,106}

Summary of Optical Senescence

Senescent changes in the eye's optics are widespread and varied, but the overall effects are dominated by three important changes. First, the reflexive autofocus feature of the eye (accommodation), which allows objects at any distance to be well focused on the retina, slowly declines with age and disappears altogether during the sixth decade of life. From the fifth decade onward, the inability to focus objects in our three-dimensional world is a major source of retinal image degradation (due to spherical defocus). Second, due to senescent changes in lens and other ocular media transparency,¹⁰⁷ and pupil miosis, retinal image illuminance is chronically lower in older eyes, and this is particularly true for short wavelength light sources. Interestingly, the widely available low transparency yellow filters used in many sunglass prescriptions that mimic the older eye transmission characteristics indicate that this aging effect will have little impact in high light environments. The third significant and widespread optical change with age is the increased optical scatter (forward scatter by the eye's optics, and increased scatter by the retina and through the sclera due to reduced pigmentation). Scatter will always have the greatest impact at night when environmental light levels can vary by huge amounts from one location to another. In combination, the reduced transparency and the increased scatter will compromise vision at night. Also, the loss of accommodation and increased scatter will compromise retinal image quality when reading at near, particularly with high gloss materials that reflect light specularly, and thus can create local high intensity regions on the page.

14.5 SENESCENT CHANGES IN VISION

Age-related changes in visual performance are due to both optical and neural factors. Changes in performance that cannot be explained solely by optical factors are assumed to be due to senescent changes in neural pathways. An understanding of neuro-visual senescence thus depends on an understanding of the optics of the aging eye.

Senescent Changes in Sensitivity under Scotopic and Photopic Conditions

Visible radiation in natural viewing conditions is normally limited to the band between 400 and 700 nm. As can be inferred from Fig. 2, the limit at short wavelengths is due to the absorption characteristics of the ocular media; the limit at long wavelengths is due to the absorption spectrum of the photopigments contained in the rod and cone receptors.

Various studies using psychophysical procedures have demonstrated age-related losses in the sensitivity of scotopic and photopic mechanisms. Dark adaptation functions (absolute threshold versus time in the dark) have two scalloped-shaped branches, the first branch being mediated by cone photoreceptors and the second mediated by rods. A number of early studies have shown that the asymptotes of each of these branches become elevated as a function of age. In addition, it has been noted that the rate of photopic and scotopic dark adaptation decreases with increasing age.¹⁰⁸ This latter result has been confirmed by more recent studies.^{109,110} It was not clear from some of the early studies whether these effects were entirely due to age-related changes in retinal illuminance caused by, among other factors, smaller pupils (see subsection “Pupils” in Sec. 14.4). Most, but not all,¹¹¹ of the studies that have taken age-related changes in pupil size and ocular media density into account have demonstrated age-related losses in scotopic sensitivity.^{112,113} Based on those studies that have shown losses in sensitivity that cannot be explained solely by preretinal changes, it can be inferred that losses in sensitivity reflect neural changes that take place at or beyond the level of the photoreceptors.

Several studies indicate that changes at the receptor level may partially account for age-related changes in light and dark adaptation and for overall sensitivity under scotopic and photopic conditions. In vivo measurements of the photopigment kinetics indicate that a change in the rate of dark adaptation is due, at least in part, to a decrease in the rate of photopigment regeneration.^{109,114} Other changes at the receptor level involve the relative ability of “aging” rods to efficiently capture light quanta. Although there are substantial losses in numbers of rod photoreceptors with age,¹¹⁵ the amount of rhodopsin photopigment in the retina decreases only slightly, if at all, with increasing age.¹¹⁶ This implies that individual surviving rods must contain more photopigment in the older retina. However, some outer segments of these aged photoreceptors undergo morphological changes such that they appear to be more convoluted in shape.¹¹⁷ These relatively deformed outer segments may contribute to reductions in scotopic sensitivity because rhodopsin molecules now possess less than optimal orientations for absorbing light.

Photopic vision is mediated by three different classes of cone photoreceptor, with each class of receptor having peak sensitivity at either short- (S), middle- (M), or long- (L) wavelengths. Several studies have shown that the sensitivity of the S cones decreases with age.^{118–120} In addition, the results of Werner and Steele (Fig. 7) indicate that M and L cones also sustain significant sensitivity losses with age, and that the rate of loss is similar for all three cone types.

One complication in specifying the rates of age-related sensitivity loss for individual cone types is that the magnitudes of these losses depend on the level of light adaptation. For example, it has been shown for an S-cone mechanism that age-related differences are greatest near the absolute threshold and gradually diminish at higher light levels.¹²¹ For this reason, Werner et al.¹²² have measured the absolute threshold of cone mechanisms at the fovea, and at locations 4 and 8° eccentric to the fovea, along the horizontal meridian in the temporal retina. Their study revealed linear declines in sensitivity with age for all three cone mechanisms that could not be accounted for by age-related changes in preretinal structures of the eye. Thus, losses in sensitivity that occur during our life span reflect, in part, neural changes taking place in both scotopic and photopic mechanisms.

Color Vision

The ability to discriminate hue declines with age beginning in early adulthood, as assessed by color matching^{123,124} and arrangement tasks.^{125–127} These losses in discrimination may be due to age-related changes in retinal illuminance caused by smaller pupil size and increased short-wave absorption by the lens. Results from the Farnsworth-Munsell 100-hue test indicate that losses in discrimination

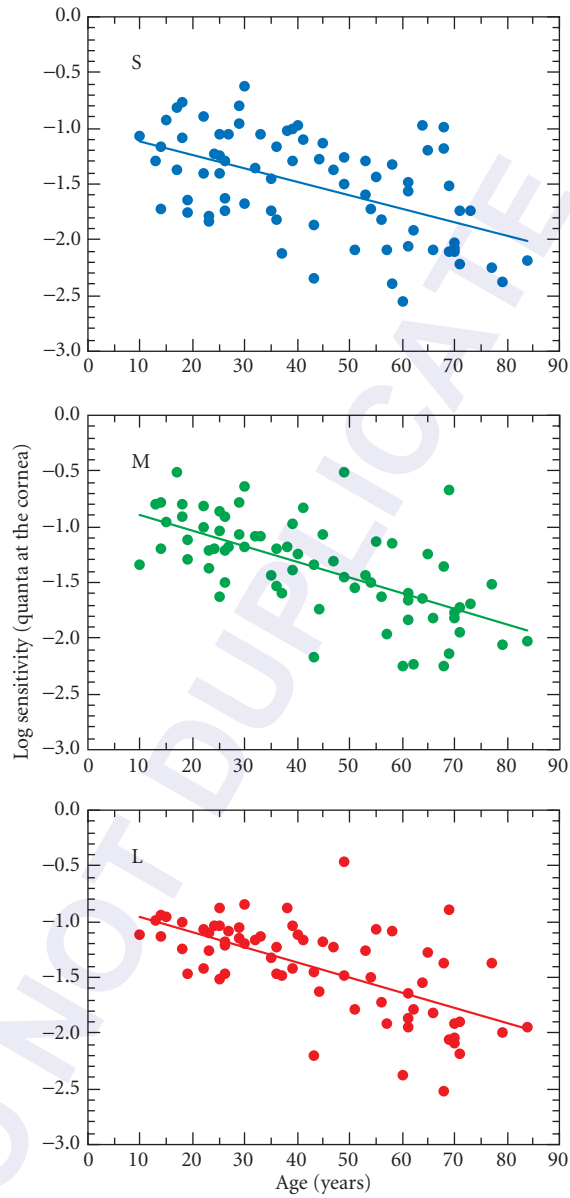


FIGURE 7 Log relative sensitivity (on a quantal basis) of S-, M-, and L-cone mechanisms is plotted as a function of age. (Data from Werner and Steele.¹²⁰)

mediated by S cones occur earlier in life than do discriminations mediated by M and L cones. However, this effect is due, at least in part, to the construction of the FM-100 test.¹²⁸ When chromatic discriminations are measured under conditions that equate the stimuli individually for retinal illuminance, losses occur not only for S-cone mediated discriminations,¹²⁹ but throughout color space for both spectral¹³⁰ and nonspectral stimuli.¹³¹ These latter results indicate that age-related losses in chromatic discrimination are due, in part, to neural processes.

Senescence of the lens results in a reduction in retinal illuminance and a change in the spectral distribution of stimuli arriving at the retina. Concurrently, there is a reduction in sensitivity of the photoreceptors and numerous morphological changes in the visual pathways. Nevertheless, there is a remarkable degree of stability in color perception across the life span. Color naming of broadband reflective samples, the wavelengths of unique blue and yellow, saturation, and the achromatic locus are all relatively invariant with normal aging.¹³² The proportion of different hue responses (e.g., blue or green) to a large set of simulated Munsell samples is independent of individual ocular media density.¹³³ Stability of color perception despite many changes in the retinal stimulus implies that the visual system continuously renormalizes itself to maintain constancy of perception. Thus, an elderly person may call the same stimulus “white” as he or she did 70 years ago, even though it must be based upon a markedly different retinal stimulus and strength of retinal signals.¹³⁴ To characterize the renormalization implicit in these findings, the chromaticity of the achromatic point was measured before and after cataract surgery.¹³⁵ There was a shift following cataract surgery (removal of a brunescient lens) that was initially toward yellow in color space, but over the course of months, it drifted back in the direction of the achromatic point before surgery. This long-term renormalization is as it should be; otherwise, the white of the young would be the yellow of the old. Such adjustments may occur by calibrating the average responses in color mechanisms according to the average color in scenes.

Spatial Vision

Visual acuity, the minimum angle of resolution, is the most time-honored measure of human visual performance, and in many clinical settings it is the only measured visual ability. It can be measured quickly and reliably with a variety of test patterns such as Snellen letters or Landolt C's. Traditionally, visual acuity is measured with letters having high contrast and moderate luminances. On average, such data suggest some stability in acuity until about age 50 and then a linear decrease with increasing age.¹³⁶ The results have been quite variable between studies, and the failures to screen subjects for disease and to optimize their refractions for the test distance make it difficult to know the true course of age-related change in visual acuity. Figure 8 shows results from several recent studies plotted in terms of log minimum angle of resolution (left axis) or equivalent Snellen and metric acuities (right axis). Data from Owsley et al.¹³⁷ resemble earlier studies, although participants were screened for visual disease and refracted for the test distance. However, with more stringent screening and optimized refraction, the change in visual acuity follows a somewhat different course according to more recent studies.^{138,139} As shown in Fig. 8, these latter studies show better overall visual acuity with a progressive decline after about age 25 years. Conversely, when subjects are not screened for visual disease and are tested with their habitual refractive correction, their acuity is considerably worse than shown in Fig. 8.⁵ Which data sets represent “normal” aging is unclear, but the better acuity values indicate that individuals can often reliably resolve finer details, or higher spatial frequencies, than suggested by earlier research, for example, 80 year olds can have 20/20 acuity.

A more general characterization of spatial vision is provided by the contrast sensitivity function (CSF), defined by the reciprocal of the minimum contrast required to detect sinusoidal gratings that vary in spatial frequency (the number of cycles per degree of visual angle, cpd). Figure 9 shows CSFs for nonflickering stimuli interpolated from a large data set by linear regression for hypothetical 20 (blue symbols) and 75-year-old (red symbols) observers. These functions could be computed in this way because the age-related changes in contrast sensitivity appear to be linear with age. A comparison of the photopic CSFs (circles) show little change in sensitivity with age at low spatial frequencies, but notable declines in sensitivity are found at middle and high spatial frequencies. These results are consistent with other studies that measured contrast sensitivity in subjects refracted for the test

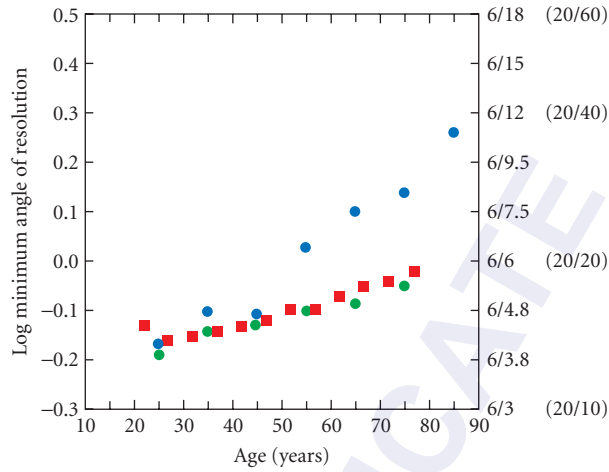


FIGURE 8 Log minimum angle of resolution (left) and Snellen equivalent acuity (right) plotted as a function of age. Blue circles are from Owsley et al.,¹³⁷ red squares and green circles are from Frisén & Frisén¹³⁸ and Elliott et al.,¹³⁹ respectively.

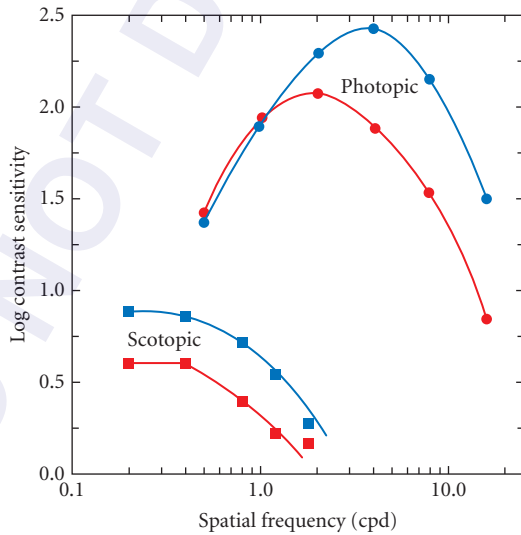


FIGURE 9 Log contrast sensitivity is plotted as a function of spatial frequency for static sinusoidal gratings. Blue and red circles represent average observers age 20 and 75 years, respectively, determined from regression equations for stimuli centered on the fovea. (From Owsley et al.¹³⁷) Blue and red squares represent average observers age 20 and 75 years, respectively, determined from regression equations for stimuli centered at 8° nasal. (From Scheffrin et al.¹⁴⁴)

distance and free from ocular disease.^{140–143} Currently it is unclear whether age-related differences in contrast sensitivity, especially at high spatial frequencies, become greater if the temporal modulation of luminance varying sinusoidal gratings is increased.^{137,142,143}

Evidence from a number of studies suggests that preneural and neural factors are required to explain the decline in contrast sensitivity with age. Preretinal factors such as smaller pupils, increases in ocular media density and intraocular scatter can partially account for these results. However, losses in contrast sensitivity are still demonstrated by studies that have controlled for pupil size¹⁴¹ and have attempted to mimic age-related reductions in retinal illuminance due to smaller pupils and increases in lens density.¹⁴³ In addition, age-related losses of contrast sensitivity at high spatial frequencies are found when measured with laser interferometry,¹⁴⁵ or when ocular aberrations, including dynamic higher-order aberrations, are corrected with adaptive optics.¹⁴⁶ It has been suggested that age-related changes in visual pathways account for less than half of the total decline in photopic contrast sensitivity at medium and high spatial frequencies with the remainder of the loss in sensitivity due to preretinal factors.¹⁴⁷ This point needs to be more thoroughly investigated.

Measurement of the CSF under scotopic conditions shows a pattern of loss quite different from photopic conditions, as illustrated by the squares in Fig. 9. These data were obtained following 30 min dark adaptation, using an artificial pupil and stimulus wavelength to minimize age-related variation in retinal illuminance.¹⁴⁴ The largest sensitivity loss to horizontally oriented sinusoidal gratings is at low spatial frequencies, a result that implies age-related change in neural mechanisms.

In Fig. 9, extrapolation of the CSFs to zero at high spatial frequencies provides a measure of visual resolution. At both photopic and scotopic light levels, age-related changes in the high frequency cutoff are noted even though participants were included only if they had retinæ without pathology when examined by ophthalmoscopy.

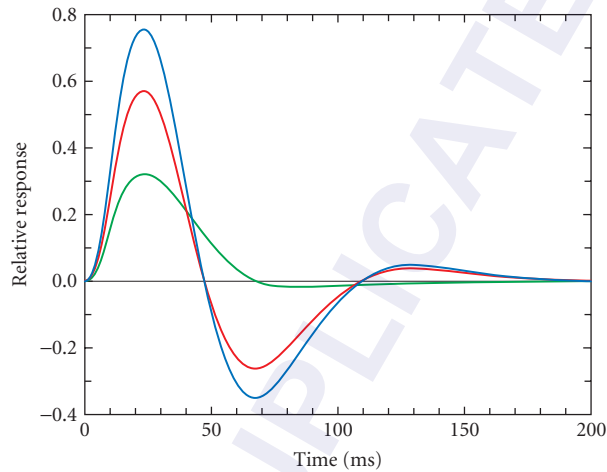
Finally, spatial visual mechanisms have been characterized in terms of their chromatic sensitivity.¹⁴⁸ To be certain that the spatial variations are strictly chromatic, it is necessary to control for chromatic aberration and individual differences in luminosity functions. With these controls in place, Hardy et al.¹⁴⁹ used stimuli modulated along separate S-cone and L-M-cone axes. They found a significant difference in chromatic contrast sensitivity between younger and older observers at all spatial frequencies and for both chromatic axes. The difference in sensitivity is greater for S-axis stimuli than for L and M varying stimuli, prior to adjustments based upon ocular media density. Contrast sensitivities measured with controls for variation in ocular media density show that differences between young and old are reduced, at least for S-cone modulation. Nevertheless, when stimuli are equated at the retina, significant differences between younger and older observers remain. The difference is similar across spatial frequency and for the two chromatic axes. From these results, it may be concluded that both optical and neural factors contribute to age-related losses in chromatic contrast sensitivity.

Temporal Vision

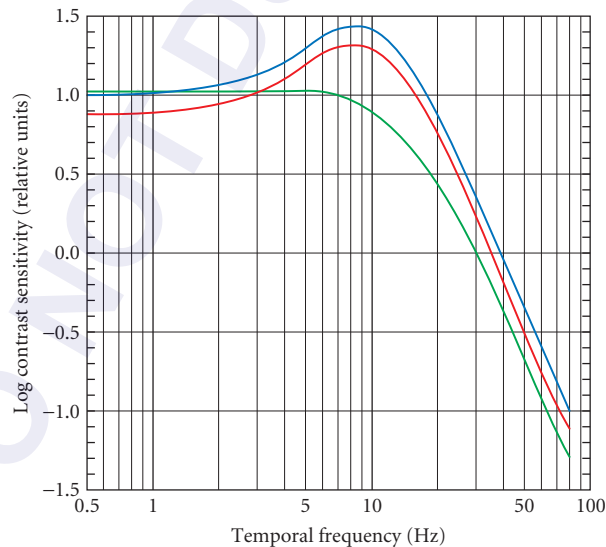
The critical fusion frequency (CFF), the lowest frequency of flicker that is indiscriminable from a steady light, decreases with age.¹⁵⁰ Part of this loss in temporal resolution is due to retinal illuminance differences between young and old. Lower luminance levels are associated with a loss in temporal resolution. While the CFF measures the high temporal resolution limit, a more complete description of temporal processing is provided by the temporal contrast sensitivity function (tCSF), the amount of contrast needed for detection as a function of temporal frequency. Tyler¹⁵¹ reported small age-related declines in temporal sensitivity for long-wavelength sinusoidal modulation, with the magnitude of the loss increasing with temporal frequency. The decrease in sensitivity began around age 20 years. Similar results were reported by Kim and Mayer,¹⁵² but sensitivity losses were not observed until about the midforties, after which there was a loss in temporal sensitivity of about 0.78 decilog per decade. It was argued in both studies that these losses in sensitivity are greater than those expected from age-related changes in optical factors alone. Based on simultaneous recordings of pattern electroretinograms (PERG) and visually evoked potentials (VEPs) in young and elderly observers, Porciatti et al.¹⁵³ have suggested that at least some of the loss in temporal sensitivity with age is due to changes at postretinal sites.

Another approach to understanding age-related changes in temporal sensitivity is to measure the impulse response function (IRF), the theoretical response to a pulse of short duration.¹⁵⁴ This approach

is important theoretically because the IRF allows predictions of how the visual system should respond to any time-varying stimulus; it is related to the tCSF by its Fourier transform. The IRF was measured for 62 observers ranging in age from 16 to 85 years using a double-pulse method.¹⁵⁵ The IRFs with luminance modulation follow an exponentially damped sine function, with a fast excitatory response followed by an inhibitory phase.¹⁵⁶ Average IRFs and tCSFs are presented in Fig. 10 by red and blue curves. Note that while the amplitude of the IRF is reduced for older subjects, the timing is similar to



(a)



(b)

FIGURE 10 (a) Calculated impulse response functions for theoretical 20-year-old (blue curve) and 80-year-old observer with normal (red curve) or reduced second-phase (inhibitory) amplitude (green curve). (b) Corresponding temporal contrast sensitivity functions calculated from the impulse response functions. (After Shinomori and Werner.¹⁵⁵)

that for younger subjects. For 9 of 25 observers over 55 years of age, however, IRFs are relatively weak and show that the second (inhibitory) phase is reduced. Consequently, the IRF for these observers is quite slow and long, as shown by the IRF plotted as a green curve. The loss of inhibitory phase changes the tCSF from band-pass to low-pass as shown by the green tCSF in Fig. 10. It is not clear why older observers fall into two groups, as all observers met stringent inclusion criteria. Control conditions demonstrated that age-related changes in the IRF under these conditions cannot be ascribed to optical factors. Thus, in most cases, the human visual system maintains a stable speed of response to a flash until at least about 80 years of age, even while the response amplitude decreases with age.

The Visual Field

The extent of the visual field narrows over adulthood, by several degrees each decade until about age 50 to 60 years and then somewhat more rapidly with increasing age.¹⁵⁷ Iwase et al.¹⁵⁸ suggested that the “volume” of the visual field decreases starting at about 30 to 40 years of age and continues for the remainder of the life span.

Senescent constriction of and reduced light sensitivity within the visual field are due in part to age-related reductions in pupil diameter and increases in light absorbance by the ocular media. Johnson et al.¹⁵⁹ performed static perimetry using experimental conditions that notably lessen the effects of lenticular senescence and reduced pupil size. They demonstrated age-related losses in sensitivity that varied from approximately 0.5 to 1 dB per decade within the central 30° visual field. Their results suggest that at least part of the visual field loss must be ascribed to neural factors.

Another approach to perimetry is to define the functional or “useful field of view,” or the area of the visual field over which information can be processed within a single fixation. The size of the useful field of view decreases when observers are engaged in additional tasks.¹⁶⁰ Ball et al.¹⁶¹ compared sensitivity of groups of younger and older observers within the central 30°. Stimuli were presented in the presence and absence of distracting stimuli to measure the useful field of view. The older group was more impaired by the distracting stimuli, particularly at greater retinal eccentricities. This restriction of the functional visual field is apparently not due to optical or retinal factors. A more recent study of individuals ages 58 to 102 years⁵ shows a dramatic impact of focal attention tasks on visual field changes with age, as shown in Fig. 11. For a typical light detection task, there is only a modest change even up to 90 years, but for a detection task with fixation attention demand, the visual field drops dramatically with age.

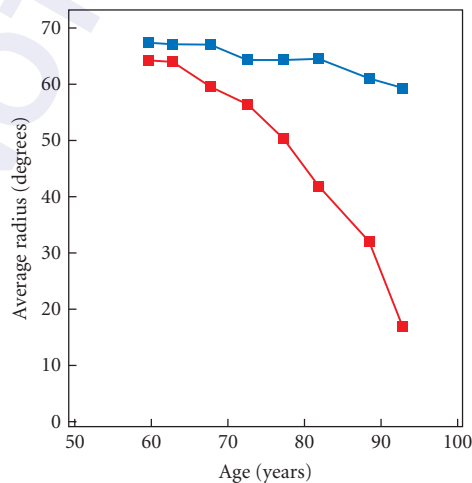


FIGURE 11 The median average radius of the visual field is plotted as a function of age for a standard light detection task (blue symbols) and with concurrent foveal attentional demand (red symbols). (After Haegerstrom-Portnoy et al.⁵)

Depth and Stereovision

Our ability to perceive the distances of objects depends upon several different kinds of information. One of the most important depth cues is stereopsis. Stereoscopic depth information is provided by retinal disparity that occurs when a single object is imaged at different retinal locations between the two eyes. There is a paucity of investigations that have examined age-related changes in stereovision. This is unfortunate because the use of retinal disparity to extract three-dimensional depth could be used to probe cortical changes that may occur with aging.^{162,163} In addition, based on the increased number of falls and other accidents that occur among the elderly, measures of stereoscopic depth perception across adulthood might be valuable.

Jani¹⁶⁴ used a standardized test of stereovision (flashlight diastereo test) to determine the percentage of 1207 individuals who failed the test at various ages. This percentage decreased from his 0 to 9 year group and remained stable from 10 to 19 through about 40 to 49 years. With further increases in age, there was a monotonic increase in the percentage of individuals failing the test. A more sensitive indicator of stereovision requires measurement of thresholds. Bell et al.¹⁶⁵ using a Verhoeff stereopter (based on the furthest distance at which the relative distance between three rods can be detected), showed that stereo resolution decreases slowly from 30 to 40 years of age and then more rapidly after about age 40. It is not clear whether these changes in stereovision are secondary to changes in the spatial CSF. This was controlled by using random-dot stereograms which varied in their binocular correlation.¹⁶⁶ Older observers required significantly higher binocular correlation than younger observers to perceive stereodepth.

Some Visual Functions That Change Minimally with Age

Senescent changes in visual function are widespread as might be expected given the optical and early neural changes that take place. It is perhaps surprising, therefore, that a number of visual functions change relatively little with age. The stability of color perception across the life span has already been described, but a number of other visual functions also appear resistant to aging.¹⁶⁷ For example, vernier acuity tasks have been used with subjects between 10 and 94 years of age. They were asked to align two or three lines or several points. This can be done with great precision (hence the name hyperacuity), and it does not change with age. The directional sensitivity of the retina (measured by the Stiles-Crawford effect of the first kind) also does not appear much affected by senescence. This task depends on the integrity of both physical and physiological mechanisms.

These and perhaps other (undiscovered) constancies in visual function across adulthood may provide insights into the types of neural changes that accompany age. In addition, such results might also serve as useful indices to distinguish between normal aging and changes in the same functions due to disease.

14.6 AGE-RELATED OCULAR DISEASES AFFECTING VISUAL FUNCTION

There are several ocular and retinal disorders that are so strongly correlated with increasing age that “age” has become a part of their name, for example, “age-related cataracts,” “age-related macular degeneration,” etc. Age, of course, represents only the passage of time, not a process, but the name does accurately imply that these problems are more likely to occur in the elderly.

Life-Span Environmental Radiation Damage

Unlike most of the human body, the eye is unique in that it contains highly transparent tissue that allows a large proportion of incident radiation to penetrate about 2.5 cm into the body. When considering radiation damage, transparency has the advantage that the energy passing through transparent

material is not absorbed and thus can do no damage. However, all of the quanta passing through the transparent portions of the eye will ultimately be absorbed deep within the eye, and thus, unlike the rest of the body in which only the surface is vulnerable to environmental radiation hazards, tissue deep within the eye can be damaged by radiation. Therefore, environmental radiation has the potential to contribute to the aging process of ocular tissue deep within the eye. The retina is additionally vulnerable to radiation hazards because flux densities at the retina can exceed those at the cornea due to the focusing characteristics of the eye's optics. The eye's susceptibility to the accumulated effect of chronic radiation damage is further exacerbated by the need for transparency across the visible spectrum and the fact that critical cells within the eye (e.g., lens fiber cells and photoreceptors) cannot regenerate. Not surprisingly, the primary ocular targets of life-span environmental radiation damage are the lens and the retina.

With the exception of those employed in professions with high levels of UV radiation (created by some welding techniques, polymer curing, electronic circuitry imaging, fluorescent testing, etc.), the primary source of environmental UV radiation is the sun. The spectral distribution of sunlight reaching different parts of the earth's surface will differ somewhat from that shown in Fig. 2 due to such factors as latitude, altitude, cloud cover, and aerosol content as well as time of day. It is important to note that just as Rayleigh scatter makes the sky blue, the atmosphere scatters an even larger proportion of UV light, and thus the sky itself is a significant source of UV radiation.

While it is straightforward to document surface UV light, it is very difficult to calculate ocular exposure because the eye is typically shaded by the brow and upper lid and thus is rarely directly irradiated by sunlight. Instead, environmental exposure to solar radiation will come primarily from reflections off surfaces or from scatter by the atmosphere. For example, while the sunlight stays constant when moving from a snowless to a snow-covered field, ocular UV irradiance increases substantially because of the high reflectivity of the snow in the UV range.¹⁶⁸ The combined effects of high altitude (little atmosphere to absorb UV) and snow (fresh snow reflects more than 90 percent of the incident UV) can rapidly create solar keratitis in any skier who fails to wear UV-filtering eye protection.

The impact of solar radiation on the eye depends on the proportion of this radiation that is absorbed. At very short wavelengths (UVB), it is the cornea (and primarily the surface epithelium) that absorbs most of the UV, and thus is at risk for UV light damage. However, due to the rapid regeneration of corneal epithelial cells (the entire corneal epithelium is regrown every 5 to 7 days), chronic UV damage is not easily accumulated. At longer UV wavelengths, corneal transparency increases to about 50 percent, and almost all of the remaining 50 percent that passes through the cornea is then absorbed by the lens. Thus, the lens is also vulnerable to UV-damage, and particularly to chronic UV damage because cells in the lens are not replaced or recycled.

Epidemiological evidence provides support for the view that UV radiation is associated with cataract formation. For example, there is a clear correlation between incidence of cataract and other diseases (e.g., pterygium and skin cancer) and exposure to sunlight. In addition, the age of cataract decreases and its incidence increases from the poles to the equator. Young's⁵³ summary of this literature indicates prevalence rates of about 14 percent for countries above 36°, 22 percent for mid-latitude (28 to 36°) countries and 28 percent for low-latitude (10 to 26°) countries. While epidemiological studies are complex due to the many factors that vary among groups, UV radiation appears to be the only factor that changes with latitude *and* has been shown to play a significant role in aging of the lens. A similar latitude dependency appears in the age of onset of presbyopia³ indicating that both of the primary age-dependent changes in lens function (presbyopia and cataract) may be regulated by levels of UV exposure.

Because of the opacity of the cornea and lens to UV radiation, relatively little UV reaches the adult retina, however, the retina is vulnerable to photochemical damage caused by UV light as well as short-wave-visible light (the "blue light hazard"). A study comprised of pseudophakic patients who had undergone bilateral cataract extraction and implantation of two IOLs differing in the amount of UV absorption suggested that increased UV exposure in the eye with the UV transparent IOL reduced visual sensitivity of short-wave-sensitive cones. It was estimated that 5 years of exposure to light through an IOL without UV absorbing chromophores produced the equivalent of more than 30 years of normal aging.⁷ Thus, the lens pigmentation responsible for UV absorption may be an adaptation to protect the retina from photochemical effects that are damaging.

Cataract

Cataracts are opacities of the lens that interfere with vision by reducing the amount of light transmitted while at the same time increasing its scatter within the eye. There are a number of different forms of cataract, and each form exhibits its own characteristic effects on vision and has different stages of development. For example, cataracts can form within the lens nucleus and cortex, and they may originate in either the axial or equatorial regions of the lens. Cataracts, depending on their type, can cause a host of visual problems including loss of visual acuity, changes in refractive error, anisikonia, perceived spatial distortions, disability glare, and the formation of multiple images. Often times more than one form is present in a given eye, and the lens of each eye may develop cataract at different rates.

In the United States, cataract is found in about one in seven individuals above the age of 40, and the proportion increases with age. The rate is much higher and typically occurs at an earlier age in individuals with diabetes. It is known that UV and infrared light can cause cataracts. Most individuals exhibit changes in their crystalline lenses with advancing age and senescent cataract may be an extension of this process.⁵³ In developing countries, cataracts occur much earlier in life (e.g., 14 to 18 years earlier in India than the United States), perhaps due to nutritional factors or greater ultraviolet exposure (as many developing countries have lower latitudes).

Cataract is the leading worldwide cause of blindness. Fortunately, it is treated with a high success rate by surgically removing the cataractous lens. In most cases in developed countries, an intraocular lens (IOL) is implanted. Cataract surgery in the west is one of the most successful therapies known, with visual improvement in nearly all cases. Nearly three million cataract surgeries were performed in the United States in 2007 and over 14 million worldwide.

Cataracts generate both forward and backward scatter of light, and quantification of both can explain about 50 percent of the variability in the visual acuity of older Americans.¹⁶⁹ It is the forward scatter, and not the backscatter seen by the clinician, that interferes with vision. Wide angle forward scatter would be expected to generate a uniform drop in contrast sensitivity, but in patients with cataract Hess and Woo¹⁷⁰ observed reductions in contrast sensitivity that were larger at high and middle spatial frequencies (not unlike those seen with retinal image defocus). Their results indicate that forward scatter produced by cataracts is not wide angle, but, as predicted by the larger than wavelength size of the particles formed in a cataract,¹⁷¹ much of the forward scattered light is distributed over a small angle (Mie as opposed to Rayleigh scattering).

Age-Related Macular Degeneration

Age-related macular degeneration (AMD) is generally a progressive, bilateral disease that compromises the integrity and functioning of the central retina. More than 13 million persons over the age of 50 in the United States have signs of macular degeneration. An individual in the early stages of AMD may experience declines in visual acuity, losses in sensitivity for central photopic¹⁷² and scotopic mechanisms,^{173,174} an inability to make fine distinctions between colors,¹⁷⁵ and metamorphopsia, a distortion in the appearance of objects in the visual scene. As the disease progresses, individuals afflicted with this condition undergo further declines in visual performance including the onset of central scotomas, defined areas of absolute or relative vision loss within the visual field, making it difficult to recognize readily faces of individuals encountered. However, because the disease does not spread past the macula or central retina, individuals can still navigate about in the environment and even “regain” some function through training to use unaffected adjacent portions of the macular region. About 1.2 million persons in the United States have reached sight-threatening late stages of AMD. This latter group is composed primarily of Caucasians aged 60 years and above, as the disease affects other ethnic groups with lower frequency. The numbers of affected individuals are increasing rapidly with the aging of the population, affecting more than 30 percent of individuals over the age of 75.

AMD is divided clinically into the exudative (wet) and nonexudative (dry) forms of the disease. In the more devastating wet form of this disease, choroidal neovascular membranes grow through

cracks in Bruch's membrane, the layer separating the choroid from the retinal pigment epithelium (RPE). Amongst other things, the RPE is responsible for regulating and providing nourishment to and removing waste products from the retina. Death of retinal cells occurs when the retina is separated for lengthy time periods from the underlying RPE as can happen when the leaky blood vessels comprising the neovascular membranes begin to exude their contents into the potential space between the RPE and photoreceptor layer of the retina. In addition, the presence of these subretinal neovascular membranes lead to scar formation beneath the macula, and this process of scarring is the primary cause of vision loss in the wet form of AMD. The wet form only comprises 10 to 15% of all cases of AMD but accounts for 90 percent of legal blindness from AMD.

In contrast to the wet form of AMD that may produce a sudden loss in vision, dry, or atrophic, AMD results in a slow progressive vision loss that usually takes place over many years. The dry form of AMD is characterized by the formation of drusen, deposits of accumulated waste products from the retina and atrophy of the RPE cells. Presently, it is unclear what causes retinal atrophy in the dry form of AMD, but it is thought to reflect ischemia from the choriocapillaris and/or a metabolic disturbance of the receptor/RPE complex leading to apoptosis.

Currently, clinical advances have been made to retard the progress of the wet form of AMD through medications and laser therapy. Fluorescein and indocyanine green angiography have been effective tools for both diagnosing wet AMD in its early stages of the disease process and defining the subretinal sites of neovascular membranes and exudative leakage from these blood vessels. By taking advantage of this valuable knowledge, laser photocoagulation is now used to destroy these vascular membranes thereby retarding the progress of the disease. For subfoveal lesions, intravitreal injection of anti-VEGF drugs can stabilize or improve vision in 95 percent of cases. Beyond these forms of treatment, visual rehabilitation has concentrated on the use of low vision aides and training AMD patients to use parafoveal and peripheral portions of their retina that are less affected by the disease.

Diabetic Retinopathy

Diabetic retinopathy is also a common disorder associated with aging. Juvenile-onset (Type I) diabetes is generally diagnosed during the first two decades of life, whereas adult-onset (Type II) diabetes is diagnosed later in life. Type II diabetes is associated with increasing age, obesity, and a family history of the disease; Mexican Americans and African Americans are affected relatively more often than Caucasian Americans. A patient with diabetes (either Type I or Type II) develops diabetic retinal changes over time; at least 40 percent of diabetic patients will eventually develop detectable retinopathy. The incidence of diabetic retinopathy is related to both blood sugar control and the duration of the disease, with most estimates suggesting clinical evidence of retinopathy occurs about 10 years after the onset of hyperglycemia. Approximately 8000 people are blinded each year by this disease.¹⁷⁶

There are several stages to diabetic retinopathy. Early in the disease there is patchy loss of the endothelial cells lining the retinal capillaries which leads to increased permeability of these blood vessels and the formation of microaneurysms. The appearances of small punctate hemorrhages, microaneurysms, and proteinaceous deposits within the retina are hallmarks of the early or "background" stage of diabetic retinopathy. As the disease progresses, blood vessels may become occluded resulting in retinal ischemia. In response to disrupted blood flow, the ischemic retina releases growth factors which induce the proliferation of new blood vessels or increased permeability of the damaged retinal vessels. These new blood vessels formed in the "proliferative" stage of diabetic retinopathy are fragile and can spontaneously bleed, leading to loss of vision from vitreous hemorrhage. The permeable damaged vessels can cause exudation and swelling of the macula leading to vision loss. In the latter stages of the disease there is death of ganglion cells from intraretinal hypoxia and contraction of scar tissue leading to retinal traction and subsequent retinal detachment in areas adjacent to relatively high concentrations of neovascularization. Tractional retinal detachment may cause visual distortion and blindness, and requires retinal surgery to minimize vision loss.

Initial treatment of diabetic retinopathy involves controlling blood sugar levels via dietary changes, oral agents, and/or insulin to facilitate the uptake of glucose into body tissue. Poorly controlled blood

sugar levels can lead to fluctuations in vision due to changes in the refractive power of the crystalline lens and the early formation of cataracts. The Diabetes Control and Complications Trial¹⁷⁷ demonstrated that tight control of blood glucose reduced the incidence of diabetic retinopathy. Additionally, fluctuations in blood glucose levels have been shown to influence normal functioning of S-cone pathways¹⁷⁸ perhaps accounting for decreased chromatic discrimination along a tritan axis demonstrated by diabetic individuals.¹⁷⁹

In the latter stages of diabetic retinopathy, laser photocoagulation is used in an attempt to maintain the integrity of the central retina. Laser photocoagulation is utilized to control retinal neovascularization associated with proliferative disease and reduced vision secondary to macular edema. These treatments have been shown to decrease severe vision loss by approximately 50 percent. Because laser photocoagulation to control neovascularization is a destructive form of treatment, newer strategies are being investigated that target specific growth factors causing neovascularization.

Glaucomas

The glaucomas are a group of chronic diseases that result in the progressive degeneration of retinal ganglion cells whose axons comprise the optic nerve. In most but not all cases of glaucoma, elevated intraocular pressure (IOP) is believed to be one of the primary risk factors for ganglion cell death. Ganglion cell death occurs through the process of apoptosis, brought on, it is believed, by the loss of trophic factors normally transported to the ganglion cell body by retrograde axoplasmic flow. In glaucoma, cell death is thought to occur following a disruption in axoplasmic flow of the retinal ganglion cell axons comprising the optic nerve; the site of this disruption is believed to be at the optic nerve head, where the insult may be mechanical¹⁸⁰ due to elevated IOP and/or a disruption in the vascular supply to the nerve head.¹⁸¹

Aqueous humor serves as the optically clear “nutritional supply” to the metabolically active but avascular structures of the lens and cornea. IOP is determined by a balance between production of aqueous humor (by the ciliary body, a secretory tissue that lies directly behind the iris) and drainage through the trabecular meshwork, a circumferential drainage system located just anterior to the iris at an angle formed by the cornea and iris root. Abnormalities in the trabecular meshwork are thought to underlie the elevated IOP in many forms of “open-angle” glaucoma.

Of the approximately 67 million persons affected by glaucoma worldwide, about equal numbers suffer from open-angle and angle-closure types, with the former being more prevalent in the United States and Europe and the latter being more prevalent in Asia and the Pacific Rim. Glaucoma accounts for about 10 percent of blindness in the United States. The frequencies of glaucomatous disorders increase with increasing age, particularly after the age of 40. By age 65 it affects about 10 percent of African Americans and about 3 percent of Caucasian Americans.

Angle-closure glaucoma typically occurs when a forward protruding iris covers the trabecular meshwork thus impeding the flow of aqueous humor from the eye. In this form of glaucoma, there can be sudden dramatic increases in IOP resulting in an abnormally hydrated cornea that leads to diminished visual acuity and accompanying ocular pain. In contrast, open-angle glaucoma is a slow progressive condition in which the IOP is too high for a particular eye resulting in ganglion cell death. Visual loss due to glaucoma occurs generally in the peripheral visual field and affects central vision only late in the disease process. Because the loss is gradual and usually occurs in areas overlapping “normal” vision in the fellow eye, patients are generally unaware of their visual loss until a great deal of damage has occurred. In addition to IOP, other risk factors associated with glaucoma include age, ethnicity, family history, corneal thickness, and myopia. Because increased IOP is not invariably associated with glaucoma, a complete eye examination is needed to make an accurate diagnosis. Often times these examinations will reveal pathognomonic changes such as an abnormal appearance of the optic nerve head, and on visual field testing, specific patterns of visual loss.

Because damage to optic nerve fibers is irreversible and the progress of glaucoma can, in most instances, be successfully controlled, early detection is of great importance. Decline in visual performance at earlier stages in the disease process can be detected through several forms of visual field testing targeting vulnerable ganglion cells.^{182,183}

Current treatment regimens for glaucoma involve the use of pharmacological agents and surgical procedures aimed at controlling IOP. Drug therapy includes the use of topical and systemic agents used singly or in combination to inhibit the formation of aqueous or promote increased drainage of aqueous from the eye. Procedures involving lasers or incisional surgery typically aid in aqueous outflow and are used when drug therapy, by itself, fails to adequately retard the progress of the disease. In addition, ongoing research is currently investigating possible therapies involving neuroprotective agents designed to reduce ganglion cell death after mechanical or chemical damage and to promote new growth.¹⁸⁴

14.7 THE AGING WORLD FROM THE OPTICAL POINT OF VIEW: PRESBYOPIC CORRECTIONS

Presbyopia refers to a loss of accommodation resulting in the near point, the closest distance on which the eye can focus, receding from the eye. This is a continuous process (see Fig. 4) and generally requires some optical aid once accommodative amplitude declines to less than 3 D. Every human suffers from this senescent loss of accommodation, and thus, optical substitutes for accommodation are widely available. Because there is no cure for this condition, and no treatment is perfect, many different treatment options have been developed. Historically, spectacle lenses have been the only option for treatment of presbyopia, but recent attempts have been made to provide presbyopes with contact lenses, IOLs, and other refractive surgery corrections to expand the depth of field of the aging eye. Most of these strategies make no claim to reinstate true accommodation, and thus, can be considered examples of “pseudoaccommodation,” in that the eye’s depth of field increases without any dynamic change in optical power. Reduced pupil size, increased aberrations, bifocal optics, and residual myopia can all lead to improved near vision without any accommodation.

During the last couple of years, several IOLs and a unique surgical procedure have all made claims that they actually restore accommodation; however, almost all of the evidence reported in support of these claims (e.g., improved near vision) could be generated by pseudoaccommodation. Thus, at this time, there is no effective treatment that can restore accommodation in presbyopic humans, but considerable research effort continues with this goal in mind.

Spectacles

To obtain a well-focused retinal image when viewing a near object, presbyopes that are emmetropic or corrected ametropes require a “near” optical correction. Since the invention of the bifocal lens by Benjamin Franklin, several ingenious approaches to this problem have been developed. With traditional bifocal lenses, the individual looks through different regions of a lens, each region having a different optical power, in order to see clearly objects at near or far. Extensions of this concept include trifocal lenses having three powers confined to three distinct zones, and aspheric lens designs in which the optical power continually changes in a prescribed manner from the top to the inferior portion of the lens. Peripheral zones of the latter designs often exhibit large amounts of astigmatism but are considered cosmetically superior. In each of these designs, the patient is able to select the preferred optical power by adjusting head and eye movements such that different regions of the lens are used to create the retinal image.

Contact Lenses

The principles of bifocal and multifocal corrections have been incorporated into contact lenses and intraocular lenses (see next section). Like their spectacle counterparts, bifocal and aspheric contact lenses have been designed to provide the presbyope with an increased depth of field. However, unlike the case of spectacle lenses, the patient cannot select a different region of the lens with which

to view different target distances. As the eye and head move, so does the contact lens. Thus, multifocal contact and intraocular lenses include two or more optical powers within the same region of the lens, and the retinal image is generated by two or more powers simultaneously. The obvious consequence is that, in a simple bifocal lens, if 50 percent of the light is well focused on the retina, then 50 percent will be defocused. The challenge for the patient is one of seeing the focused portion of the image with the simultaneously present defocused image. Such lenses are often called “simultaneous vision” lenses.

Some designs of bifocal contact lenses consist of concentric annuli that alternate in power between near and far optical corrections while other bifocal designs include only two regions, one with a distance correction and one with a near correction. A third type of simultaneous vision lens attempts to provide increased depth of field by employing a highly aberrated “multifocal” design. A fourth type of bifocal lens employs a diffractive optical element (DOE) and is often referred to as a “diffractive” lens.¹⁸⁵ Although most designs provide improved near vision compared to a standard monofocal lens,^{186,187} this improvement is generally accompanied by a reduction in best-corrected distance vision.

Instead of including two different corrections in the same eye, an alternate approach called “monovision” consists of different monofocal optical corrections in each eye, one for near and the other for far. Thus, unlike simultaneous bifocal vision which adds the focused and defocused images linearly (light adds) at the retina, monovision patients employ nonlinear binocular neural summation¹⁸⁹ in the cortex which biases vision to the higher contrast and better focused image. Thus, the defocused image present in one eye when viewing a distant or near target generally goes unnoticed. This failure to observe the defocused image has led to the notion that the blurred image is somehow suppressed, by a process called “blur suppression.”¹⁹⁰

It should be noted that the retinal blur produced by monovision and simultaneous bifocal contact lens corrections leads to losses in contrast sensitivity, especially at high spatial frequencies. The magnitude of vision loss can be quite small at distance and near but large at intermediate distances where neither image is well focused.^{187,191} Stereopsis can also be affected adversely by the unilateral blur and aniseikonia that accompany monovision corrections. In addition, patients with simultaneous vision corrections that may provide reasonable distance and near vision may experience “ghosting” caused by the defocused portion of the light.¹⁹²

Intraocular Lenses

Following extraction of a cataractous lens, a refractive error of about 17 D remains to be corrected. A person without their natural lens is called an aphake. In most cases, an intraocular lens (IOL) is inserted, and the wearer is said to be pseudophakic. IOLs have traditionally been constructed from polymethylmethacrylate, but a variety of other materials have been developed (e.g., silicone, hydrogels) so that the lens can be folded and inserted through smaller incisions than are required to insert a hard plastic lens. Many IOLs in the past 15 years have incorporated UV-absorbing chromophores, but older IOLs were virtually transparent to UV radiation from 320 to 400 nm.

Bifocal and multifocal corrections have been incorporated into IOLs. While the optical principles are mostly the same as presbyopic contact lens corrections, some different problems may emerge because of the placement of the optic and maintenance of its position. For example, IOL tilt is associated with astigmatism. However, decentering a bifocal IOL will generate a lateral shift of the defocused images. These off-set defocused images are described as ghosts which are also observed by presbyopes with bifocal contact lens corrections.

A simple segmented bifocal design might have a distance correction in the center and a near correction in the lens surround. With the correct pupil size, light might be equally distributed between the distance and near optics, but as pupil size changes, this balance will be lost. For example, if the pupil is reduced in size such that it matched the size of the central distance portion, the bifocal lens becomes a monofocal correction. To avoid this pupil size dependence, both contact lenses and IOLs have been designed with multiple zones. One type of multizone lens employs standard refractive optics but has alternating near and distance rings across the lens ensuring that both near and distant

optics will contribute to the retinal image irrespective of the pupil size, and thus, the increased depth of field will be present at all times. A second type of multizone bifocal lens, a diffractive optical element, necessarily has multiple rings as part of its design,^{193,194} and thus is inherently resistant to pupil size changes. These diffractive lenses control their refractive power by the size of the rings and control the balance between near and distance by the size of the step between successive zones. In addition to enabling bifocal optics, these diffractive lenses have the added potential advantage of correcting some of the eye's chromatic aberration because, unlike refractive lenses (and human eyes) which will always have more power at short wavelengths, diffractive lenses have more power at longer wavelengths.

As with bifocal contact lenses, reduced image contrast in any eye using a bifocal or multifocal IOL can be expected due to the superposition of this second out-of-focus image.¹⁹⁵ The consequence of these problems may be exacerbated in older patients due to their lower contrast sensitivity. Problems with glare and halos around light sources are greater in patients with bilateral multifocal compared to bilateral monofocal IOLs, but the former patients sometimes achieve less dependence on additional spectacle corrections.¹⁹⁶

Noncataract Related Refractive Surgeries

There are three issues of ocular senescence related to refractive surgeries (see Chap. 16). First, as postrefractive emmetropes, many young ex-myopes will enjoy a spectacle or contact lens free life. However, as presbyopia progresses, these corrected myopes will now need reading aids for near work. Of course, if left uncorrected, many of these myopes would avoid the need for reading glasses or bifocals by simply removing their distance corrections. Second, some presbyopic patients opt for monovision refractive surgeries, essentially exchanging bilateral myopia for anisometropia. Third, some young adults opting for clear lens extraction or phakic IOLs to treat their high myopia will all experience a form of early and total presbyopia, and thus they will need a near reading aid. There is some concern that these early IOL surgeries will generate accelerated loss of endothelial cells which may lead to corneal failure in the distant future.

Accommodation Restoration

All of the standard treatments for presbyopia attempt to increase the depth of field without any attempt to restore the lens' ability to adjust its power. That is, pseudoaccommodation is used to replace true accommodation. In the last few years, however, attempts have been made to restore accommodation in presbyopic eyes and to enable true accommodation in pseudophakic eyes (see Ref. 197 for recent comprehensive review). These approaches are based upon the assumption that some useful ciliary muscle strength persists into senescence. There have been two general approaches tried so far. In one, the sclerotic senescent lens is replaced with a lens made from a pliable elastic material (e.g., silicone) that will change shape to achieve a change in power. However, physiological incompatibilities leading to capsule opacification, and problems in the balance between power and change in power, remain.¹⁹⁸ The second approach employs fixed power optical elements that generate a change in whole eye power by movement. Interestingly, although natural accommodation generates increased positive power by a combination of changing lens shape and position, the position changes required by this type of accommodating IOL are much larger than those exhibited by the natural lens. One other technique involves surgical implantation of scleral "expansion rings" to expand the equatorial diameter of the presbyopic eye with the goal of regaining accommodation. The theory of accommodation and loss of accommodation with age behind this surgical procedure is not generally accepted to be correct,^{197,199} and careful study of one satisfied patient found no evidence of accommodative amplitude higher than those of untreated age-matched controls.²⁰⁰

IOL designs that employ structural changes in the eye created by contraction of the ciliary muscle to change the position or shape of the lens are called "accommodating" IOLs. Designs that employ lens movement to change power have already made it to the market place, whereas designs that employ pliable polymers continue to be tested in the lab.

Two general strategies have been employed to generate accommodating IOLs that work via lens movements. One type uses a single plus lens supported by a hinged haptic, and when the equatorial diameter of the eye decreases as the ciliary muscle contracts, lateral forces on the outer edge of the haptic, combined with potential forward movement of the vitreous, will push the lens forward. When the positive powered IOL and positive powered cornea are closer together, the net result is more plus power. A second design strategy employs a dual optic lens containing a high powered plus lens and a compensatory negative lens. Again, as the ciliary muscle contracts, the forces produced are intended to increase the separation of the two optical components of the dual optic lens, and in so doing, increase its power.^{197,201} Optical theory predicts that the latter design has more potential for changing power than that produced by movement of the simple lens.²⁰² Experimental and clinical studies struggle to show evidence of accommodation.^{197,203} Some evidence for lens movement has been demonstrated in response to “pharmacologically induced accommodation” in which lens position or eye refraction is compared in eyes in which ciliary muscle contraction is either blocked (e.g., with Cyclopentolate) or activated (e.g., with pilocarpine).^{197,204} Disappointingly, stimulus driven accommodation restoration has yet to be demonstrated,^{202,205} and future studies need to ensure that they carefully discriminate between genuine accommodation and pseudoaccommodation.

This is a rapidly evolving field in which new designs are expected to emerge in the next few years. Lens designs that employ elastic IOLs that can change shape in response to ciliary muscle contraction offer the greatest potential,²⁰⁶ but these face very significant challenges. The translating designs that have made it to market seem to generate little if any true accommodation.²⁰³ It is important to note that most people in their 40s with perhaps only 2 or 2.5 D of accommodation can function without a reading add, therefore accommodating IOLs only need to generate about this level of accommodation for most people to function.

14.8 CONCLUSIONS

The proportion of the population above age 65 years is increasing, and there is a parallel rise in the proportion of the population that is blind or in need of refractive corrections. Senescent deterioration in the eye’s optics can lead to a wide range of changes in visual function due to reductions in retinal illuminance and degradations in image quality. Age-related changes in vision also result from neural losses. A few visual functions are spared, but substantial changes in the spatial, spectral, and temporal analysis of the retinal image have been well documented. The demographic shift in the proportion of the population that is elderly has not yet reached an equilibrium, so the social and economic consequences can be expected to be even more profound in the future. At this point in history, in spite of technical, optical, and surgical advances, the optical and visual quality of all eyes will deteriorate with advancing age, and for some, age-related pathology will precipitate serious disability due to dramatic loss of visual capability.

14.9 ACKNOWLEDGMENTS

This chapter was supported by NIA grant AG04058.

14.10 REFERENCES

1. U.S. Senate Committee on Aging, “Aging America,” (U.S. Department of Health and Human Services, Washington, D.C., 1985–86).
2. United Nations Population Division, “World Population Projections to 2150,” (New York, NY, 1998).
3. R. A. Weale, *A Biography of the Eye* (Lewis, London, 1982).
4. H. M. Leibowitz, D. E. Krueger, L. R. Maunder, R. C. Milton, M. M. Kini, H. A. Kahn, R. J. Nickerson, et al., “The Framingham Eye Study Monograph,” *Surv. Ophthalmol. (Supplement)* **24**:335–610 (1980).

5. G. Haegerstrom-Portnoy, M. E. Schneck, and J. A. Brabyn, "Seeing into Old Age: Vision Function beyond Acuity," *Optom. Vis. Sci.* **76**:141–158 (1999).
6. B. R. Hammond, E. J. Johnson, R. M. Russell, N. I. Krinsey, K. J. Yeum, R. B. Edwards, and D. M. Shodderly, "Dietary Modification of Human Macular Pigment Density," *Invest. Ophthalmol. Visual Sci.* **38**:1795–1801 (1997).
7. J. S. Werner, V. G. Steele, and D. S. Pfoff, "Loss of Human Photoreceptor Sensitivity Associated with Chronic Exposure to Ultraviolet Radiation," *Ophthalmology* **96**:1552–1558 (1989).
8. G. Haegerstrom-Portnoy, "Short-Wavelength-Sensitive-Cone Sensitivity Loss with Aging: A Protective Role for Macular Pigment?," *J. Opt. Soc. Am.* **A5**:2140–2144 (1988).
9. L. M. Verbrugge and D. L. Patrick, "Seven Chronic Conditions: Their Impact on U.S. Adults' Activity Levels and Use of Medical Services," *Am. J. Public Health* **85**:173–182 (1995).
10. S. Book, "Vision Loss," (Alliance for Aging Research, 2007).
11. R. S. Ramrattan, R. C. W. Wolfs, S. Panda-Jonas, J. B. Jonas, D. Bakker, H. A. Pols, A. Hofman, and P. T. V. M. de Jong, "Prevalence and Causes of Visual Field Loss in the Elderly and Associations with Impairment in Daily Functioning," *Arch. Ophthalmol.* **119**(12):1788–1794 (2001).
12. S. Brink, "Greying of Our Communities Worldwide," *Ageing Internat.* **23**:13–31 (1997).
13. H. Bouma and K. Sagawa, "Ageing Effects on Vision," (Commission Internationale de L'Eclairage, Publication No. 133, Warsaw, Poland, 1999), pp. 368–370.
14. P. E. King-Smith, B. A. Fink, R. M. Hill, K. W. Koelling, and J. M. Tiffany, "The Thickness of the Tear Film," *Curr. Eye Res.* **29**(4–5):357–368 (2004).
15. N. L. Himebaugh, A. R. Wright, A. Bradley, C. G. Begley, and L. N. Thibos, "Use of Retroillumination to Visualize Optical Aberrations Caused by Tear Film Break-Up," *Optom. Vis. Sci.* **80**(1):69–78 (2003).
16. R. Tutt, A. Bradley, C. Begley, and L. N. Thibos, "Optical and Visual Impact of Tear Break-Up in Human Eyes," *Invest. Ophthalmol. Visual Sci.* **41**(13):4117–4123 (2000).
17. D. A. Schaumberg, D. A. Sullivan, J. E. Buring, and M. R. Dana, "Prevalence of Dry Eye Syndrome among U.S. women," *Am. J. Ophthalmol.* **136**(2):318–326 (2003).
18. S. Patel and I. Wallace, "Tear Meniscus Height, Lower Punctum Lacrimale, and the Tear Lipid Layer in Normal Aging," *Optom. Vis. Sci.* **83**(10):731–739 (2006).
19. M. C. Acosta, M. L. Alfaro, F. Borrás, C. Belmonte, and J. Gallar, "Influence of Age, Gender and Iris Color on Mechanical and Chemical Sensitivity of the Cornea and Conjunctiva," *Exp. Eye Res.* **83**(4):932–938 (2006).
20. W. D. Mathers, J. A. Lane, and M. B. Zimmerman, "Tear Film Changes Associated with Normal Aging," *Cornea* **15**(3):229–234 (1996).
21. D. A. Dartt, "Dysfunctional Neural Regulation of Lacrimal Gland Secretion and its Role in the Pathogenesis of Dry Eye Syndromes," *Ocul. Surf.* **2**:76–91 (2004).
22. R. W. Young, "The Bowman Lecture, 1982. Biological Renewal. Applications to the eye," *Trans. Ophthalmol. Soc. U.K.* **102** (Pt 1):42–75 (1982).
23. J. Richard, L. Hoffart, F. Chavane, B. Ridings, and J. Conrath, "Corneal Endothelial Cell Loss after Cataract Extraction by Using Ultrasound Phacoemulsification versus a Fluid-Based System," *Cornea* **27**(1):17–21 (2008).
24. M. Millodot and I. A. Newton, "A Possible Change of Refractive Index with Age and Its Relevance to Chromatic Aberration," *Graefes Arch. Clin. Exp. Ophthalmol.* **201**:159–167 (1976).
25. M. W. Morgan, "Changes in Visual Function in the Aging Eye," in *Vision and Aging: General and Clinical Perspectives*, J. A. A. Rosenbloom, and M. W. Morgan (eds.) (Fairchild Publications, New York, 1986), pp. 121–134.
26. A. Guirao, M. Redondo, and P. Artal, "Optical Aberrations of the Human Cornea and Function of Age," *J. Opt. Soc. Am.* **A17**(10):1697–1702 (2000).
27. P. Artal, E. Berrio, A. Guirao, and P. Piers, "Contribution of the Cornea and Internal Surfaces to the Change of Ocular Aberrations with Age," *J. Opt. Soc. Am.* **A19**(1):137–143 (2002).
28. S. Lerman, "Biophysical Aspects of Corneal and Lenticular Transparency," *Curr. Eye Res.* **3**:3–14 (1984).
29. J. Lu, Z. Wang, P. Lu, X. Chen, W. Zhang, K. Shi, Y. Kang, L. Ke, and R. Chen, "Pterygium in an Aged Mongolian Population: A Population-Based Study in China," *Eye* (2007).
30. E. A. Boettner and J. R. Wolter, "Transmission of the Ocular Media," *Invest. Ophthalmol.* **1**:776–783 (1962).

31. H. Brandhorst, J. Hickey, H. Curtis, and E. Ralph, "Interim Solar Cell Testing Procedures for Terrestrial Applications," (Lewis Research Center, Cleveland, OH, 1975).
32. J. S. Werner, "The Damaging Effects of Light on the Eye and Implications for Understanding Changes in Vision Across the Life Span," in *The Changing Visual System: Maturation and Aging in the Central Nervous System*, P. Bagnoli and W. Hodós (eds.) (Plenum Press, New York, 1991), pp. 295–309.
33. B. T. Gabelt and P. L. Kaufman, "Changes in Aqueous Humor Dynamics with Age and Glaucoma," *Prog. Retin. Eye Res.* **24**:612–637 (2005).
34. A. Caldicott and W. N. Charman, "Diffraction Haloes Resulting from Corneal Oedema and Epithelial Cell Size," *Ophthalm. Physiol. Opt.* **22**(3):209–213 (2002).
35. M. P. Hatton, V. L. Perez, and C. H. Dohlman, "Corneal Oedema in Ocular Hypotony," *Exp. Eye Res.* **78**(3):549–552 (2004).
36. F. W. Campbell and R. W. Gubisch, "Optical Quality of the Human Eye," *J. Physiol. (London)* **186**:558–578 (1966).
37. M. S. Banks, W. S. Geisler, and P. J. Bennett, "The Physical Limits of Grating Visibility," *Vis. Res.* **27**(11):1915–1924 (1987).
38. S. B. Laughlin, "Retinal Information Capacity and the Function of the Pupil," *Ophthalm. Physiol. Opt.* **12**(2):161–164 (1992).
39. A. L. Kornzweig, "Physiological Effects of Age on the Visual Process," *Sight-Saving Rev.* **24**:130–138 (1954).
40. I. E. Loewenfeld, "Pupillary Changes Related to Age," in *Topics in Neuro-ophthalmology* H. S. Thompson, R. Daroff, L. Frisén, J. S. Glaser, and M. D. Sanders (eds.) (Williams and Wilkins, Baltimore, MD, 1979), pp. 124–150.
41. M. E. Sloane, C. Owsley, and S. L. Alvarez, "Aging, Senile Miosis and Spatial Contrast Sensitivity at Low Luminance," *Vis. Res.* **28**:1235–1246 (1988).
42. R. A. Applegate, I. Donnelly, W. J., J. D. Marsack, D. E. Koenig, and K. Pesudovs, "Three-Dimensional Relationship between High-Order Root-Mean-Square Wavefront Error, Pupil Diameter, and Aging," *J. Opt. Soc. Am. A* **24**:578–587 (2007).
43. A. Guirao, C. González, M. Redondo, E. Geraghty, S. Norrby and P. Artal, "Average Optical Performance of the Human Eye as a Function of Age in a Normal Population," *Invest. Ophthalmol. Visual Sci.* **40**:203–213 (1999).
44. R. A. Weale, *The Senescence of Human Vision* (Oxford University Press, Oxford, U.K., 1992).
45. G. Smith, "Angular Diameter of Defocus Blur Discs," *Am. J. Optom. Physiol. Opt.* **59**(11):885–889 (1982).
46. R. C. Augusteyn, "Growth of the Lens: In Vitro Observations," *Clin. Exp. Optom.* **91**(3):226–239 (2008).
47. D. A. Atchison, E. L. Markwell, S. Kasthurirangan, J. M. Pope, G. Smith, and P. G. Swann, "Age-Related Changes in Optical and Biometric Characteristics of Emmetropic Eyes," *J. Vis.* **8**:1–20 (2008).
48. L. M. Bova, M. H. Sweeney, J. F. Jamie, and R. J. Truscott, "Major Changes in Human Ocular UV Protection with Age," *Invest. Ophthalmol. Visual Sci.* **42**(1):200–205 (2001).
49. J. S. Werner, "Development of Scotopic Sensitivity and the Absorption Spectrum of the Human Ocular Media," *J. Opt. Soc. Am.* **72**:247–258 (1982).
50. R. A. Weale, "Age and the Transmittance of the Human Crystalline Lens," *J. Physiol. (London)* **395**:577–587 (1988).
51. J. Pokorny, V. C. Smith, and M. Lutze, "Aging of the Human Lens," *Appl. Opt.* **26**:1437–1440 (1987).
52. P. A. Sample, F. D. Esterson, R. N. Weinreb, and R. M. Boynton, "The Aging Lens: In Vivo Assessment of Light Absorption in 84 Eyes," *Invest. Ophthalmol. Visual Sci.* **29**:1306–1311 (1988).
53. R. Young, *Age-Related Cataract* (Oxford University Press, New York, 1991).
54. J. J. Vos and D. van Norren, "Thermal Cataract, from Furnaces to Lasers," *Clin. Exp. Optom.* **87**(6):372–376 (2004).
55. R. A. Weale, "Senile Cataract. The Case Against Light," *Ophthalmology* **90**:420–423 (1983).
56. W. N. Charman, "The Eye in Focus: Accommodation and Presbyopia," *Clin. Exp. Optom.* **91**(3):207–225 (2008).
57. B. K. Pierscionek and R. A. Weale, "Presbyopia—A Maverick of Human Aging," *Arch. Gerontol. Geriatr.* **20**(3):229–240 (1995).
58. R. Brückner, E. Batschelet, and F. Hugenschmidt, "The Basel Longitudinal Study on Aging (1955–1978)," *Doc. Ophthalmol.* **64**:235–310 (1987).

59. F. C. Donders, *On the Anomalies of Accommodation and Refraction of the Eye* (New Sydenham, London, 1864).
60. C. Oyster, *The Human Eye, Structure and Function* (Sinauer Associates Inc., Sunderland, 1999).
61. K. A. Walton, C. H. Meyer, C. J. Harkrider, T. A. Cox, and C. A. Toth, "Age-Related Changes in Vitreous Mobility as Measured by Video B Scan Ultrasound," *Exp. Eye Res.* **74**(2):173–180 (2002).
62. Z. Wang, J. Dillon, and E. R. Gaillard, "Antioxidant Properties of Melanin in Retinal Pigment Epithelial Cells," *Photochem. Photobiol.* **82**(2):474–479 (2006).
63. D. M. Snodderly, J. D. Auran, and F. C. Delori, "The Macular Pigment. II. Spatial Distribution in Primate Retinas," *Invest. Ophthalmol. Visual Sci.* **25**:674–685 (1984).
64. J. J. Vos, "Tabulated Characteristics of a Proposed 2° Fundamental Observer," (Institute for Perception—TNO, Soesterberg, The Netherlands, 1978).
65. D. Moreland and P. Bhatt, "Retinal Distribution of Macular Pigment," in *Colour Vision Deficiencies VII*, G. Verriest (ed.) (Dr. W. Junk, The Hague, 1984), pp. 127–132.
66. J. S. Werner, S. K. Donnelly, and R. Kliegl, "Aging and Human Macular Pigment Density; Appended with Translations from the Work of Max Sultze and Ewald Hering," *Vis. Res.* **27**:257–268 (1987).
67. R. A. Bone, J. T. Landrum, L. Fernandez, and S. L. Tarsis, "Analysis of the Macular Pigment by HPLC: Retinal Distribution and Age Study," *Invest. Ophthalmol. Visual Sci.* **29**:843–849 (1988).
68. V. M. Reading and R. A. Weale, "Macular Pigment and Chromatic Aberration," *J. Opt. Soc. Am.* **A64**:231–234 (1974).
69. K. Kirschfield, "Carotenoid Pigments: Their Possible Role in Protecting Against Photooxidation in Eyes and Photoreceptor Cells," *Proc. R. Soc. London* **B216**:71–85 (1982).
70. X. Zhang, M. Ye, A. Bradley, and L. Thibos, "Apodization by the Stiles-Crawford Effect Moderates the Visual Impact of Retinal Image Defocus," *J. Opt. Soc. Am.* **A16**(4):812–820 (1999).
71. K. Franze, J. Grosche, S. N. Skatchkov, S. Schinkinger, C. Foja, D. Schild, O. Uckermann, K. Travis, A. Reichenbach, and J. Guck, "Muller Cells are Living Optical Fibers in the Vertebrate Retina," *Proc. Natl. Acad. Sci. USA* **104**(20):8287–8292 (2007).
72. M. C. Rynders, T. Grosvenor, and J. M. Enoch, "Stability of the Stiles-Crawford Function in a Unilateral Amblyopic Subject over a 38-Year Period: A Case Study," *Optom. Vis. Sci.* **72**(3):177–185 (1995).
73. P. J. DeLint, J. J. Vos, T. T. Berendschot, and D. van Norren, "On the Stiles-Crawford Effect with Age," *Invest. Ophthalmol. Visual Sci.* **38**(6):1271–1274 (1997).
74. J. M. Gorrand and F. C. Delori, "Reflectance and Curvature of the Inner Limiting Membrane at the Foveola," *J. Opt. Soc. Am.* **A16**(6):1229–1237 (1999).
75. L. Feeney-Burns, E. S. Hilderbrand, and S. Eldridge, "Aging Human RPE: Morphometric Analysis of Macular, Equatorial, and Peripheral Cells," *Invest. Ophthalmol. Visual Sci.* **25**(2):195–200 (1984).
76. T. Sarna, J. M. Burke, W. Korytowski, M. Rozanowska, C. M. Skumatz, A. Zareba, and M. Zareba, "Loss of Melanin from Human RPE with Aging: Possible Role of Melanin Photooxidation," *Exp. Eye Res.* **76**(1):89–98 (2003).
77. F. C. Delori, C. K. Dorey, G. Staurenghi, O. Arend, D. G. Goger, and J. J. Weiter, "In Vivo Fluorescence of the Ocular Fundus Exhibits Retinal Pigment Epithelium Lipofuscin Characteristics," *Invest. Ophthalmol. Visual Sci.* **36**:718–729 (1995).
78. J. R. Sparrow and M. Boulton, "RPE Lipofuscin and its Role in Retinal Pathobiology," *Exp. Eye Res.* **80**:595–606 (2005).
79. J. J. Vos, "On the Cause of Disability Glare and Its Dependence on Glare Angle, Age and Ocular Pigmentation," *Clin. Exp. Optom.* **86**(6):363–370 (2003).
80. W. S. Goh and C. S. Lam, "Changes in Refractive Trends and Optical Components of Hong Kong Chinese Aged 19–39 Years," *Ophthalm. Physiol. Opt.* **14**(4):378–382 (1994).
81. L. L. Lin, Y. F. Shih, C. K. Hsiao, and C. J. Chen, "Prevalence of Myopia in Taiwanese Schoolchildren: 1983 to 2000," *Ann. Acad. Med. Singapore* **33**(1):27–33 (2004).
82. R. A. Weale, "On the Age-Related Prevalence of Anisometropia," *Ophthalm. Res.* **34**(6):389–392 (2002).
83. J. R. Lavery, J. M. Gibson, D. E. Shaw, and A. R. Rosenthal, "Refraction and Refractive Errors in an Elderly Population," *Ophthalm. Physiol. Opt.* **8**(4):394–396 (1988).

84. G. Westheimer and J. Liang, "Influence of Ocular Light Scatter on the Eye's Optical Performance," *J. Opt. Soc. Am.* **A12**:1417–1424 (1995).
85. E. Wolf, "Glare and Age," *Arch. Ophthalmol.* **64**:502–514 (1960).
86. J. k. IJspeert, P. W. T. de Waard, T. J. T. P. van den Berg, and P. T. V. M. de Jong, "The Intraocular Straylight Function in 129 Healthy Volunteers; Dependence on Angle, Age and Pigmentation," *Vis. Res.* **36**:699–707 (1990).
87. M. L. Hennelly, J. L. Barbur, D. F. Edgar, and E. G. Woodward, "The Effect of Age on the Light Scattering Characteristics of the Eye," *Ophthalm. Physiol. Opt.* **18**:197–203 (1998).
88. B. R. Wooten and G. A. Geri, "Psychophysical Determination of Intraocular Light Scatter as a Function of Wavelength," *Vis. Res.* **27**:1291–1298 (1987).
89. D. Whitaker, R. Steen, and D. B. Elliott, "Light Scatter in the Normal Young, Elderly, and Cataractous Eye Demonstrates Little Wavelength Dependency," *Optom. Vis. Sci.* **70**(11):963–968 (1993).
90. T. J. T. P. van den Berg, J. K. IJspeert, and P. W. T. de Waard, "Dependence of Intraocular Straylight on Pigmentation and Light Transmission Through the Ocular Wall," *Vis. Res.* **31**:1361–1367 (1991).
91. J. E. Coppens, L. Franssen, and T. J. van den Berg, "Wavelength Dependence of Intraocular Straylight," *Exp. Eye Res.* **82**(4):688–692 (2006).
92. S. Lerman and R. F. Borkman, "Spectroscopic Evaluation and Classification of the Normal, Aging and Cataractous Lens," *Ophthalm. Res.* **8**:335–353 (1976).
93. D. B. Elliott, K. C. H. Yang, K. Dumbleton, and A. P. Cullen, "Ultraviolet-Induced Lenticular Fluorescence: Intraocular Straylight Affecting Visual Function," *Vis. Res.* **33**:1827–1833 (1993).
94. J. L. Barbur, A. J. Harlow, and C. Williams, "Light Scattered in the Eye and Its Effect on the Measurement of the Colour Constancy Index," in *Colour Vision Deficiencies XIII*, C. R. Cavonius, (ed.), (Kluwer Academic Publishers, Dordrecht, 1997), pp. 439–438.
95. J. J. Vos, "Age Dependence of Glare Effects and Their Significance in Terms of Visual Ergonomics," (Lighting Research Institute, New York, 1995).
96. J. Porter, A. Guirao, I. G. Cox, and D. R. Williams, "Monochromatic Aberrations of the Human Eye in a Large Population," *J. Opt. Soc. Am.* **A18**(8):1793–1803 (2001).
97. L. N. Thibos, X. Hong, A. Bradley, and X. Cheng, "Statistical Variation of Aberration Structure and Image Quality in a Normal Population of Healthy Eyes," *J. Opt. Soc. Am.* **A19**(12):2329–2348 (2002).
98. J. L. Alio, P. Schimchak, H. P. Negri, and R. Montes-Mico, "Crystalline Lens Optical Dysfunction Through Aging," *Ophthalmology* **112**(11):2022–2029 (2005).
99. I. Brunette, J. M. Bueno, M. Parent, H. Hamam, and P. Simonet, "Monochromatic Aberrations as a Function of Age, from Childhood to Advanced Age," *Invest. Ophthalmol. Visual Sci.* **44**(12):5438–5446(2003).
100. P. Artal, M. Ferro, I. Miranda, and R. Navarro, "Effects of Aging in Retinal Image Quality," *J. Opt. Soc. Am.* **A10**:1656–1662 (1993).
101. R. I. Calver, M. J. Cox, and D. B. Elliott, "Effect of Aging on the Monochromatic Aberrations of the Human Eye," *J. Opt. Soc. Am.* **A16**:2069–2078 (1999).
102. L. Thibos, M. Ye, X. Zhang, and A. Bradley, "Chromatic Eye: A New Reduced-Eye Model of Ocular Chromatic Aberration in Humans," *Appl. Opt.* **31**:3594–3600 (1992).
103. P. P. Howarth, X. X. Zhang, A. Bradley, D. L. Still, and L. N. Thibos, "Does the Chromatic Aberration of the Eye Vary with Age?," *J. Opt. Soc. Am.* **A5**:2087–2092 (1988).
104. H. D. Whitefoot and W. N. Charman, "Hyperchromatic Lenses as Potential Aids for the Presbyope," *Ophthalm. Physiol. Opt.* **15**(1):13–22 (1995).
105. D. A. Atchison, M. Ye, A. Bradley, M. J. Collins, X. Zhang, H. A. Rahman, and L. N. Thibos, "Chromatic Aberration and Optical Power of a Diffractive Bifocal Contact Lens," *Optom. Vis. Sci.* **69**(10):797–804 (1992).
106. N. Lopez-Gil and R. Montes-Mico, "New Intraocular Lens for Achromatizing the Human Eye," *J. Cataract Refract. Surg.* **33**(7):1296–1302 (2007).
107. J. van de Kraats and D. van Norren, "Optical Density of the Aging Human Ocular Media in the Visible and the UV," *J. Opt. Soc. Am.* **A24**(7):1842–1857 (2007).
108. R. G. Domey, R. A. McFarland, and E. Chadwick, "Dark Adaptation as a Function of Age and Time: II. A Derivation," *J. Gerontol.* **15**:267–279 (1960).

109. C. D. Coile and H. D. Baker, "Foveal Dark Adaptation, Photopigment Regeneration, and Aging," *Vis. Neurosci.* **8**:27–29 (1992).
110. G. R. Jackson, C. Owsley, and G. McGwin, "Aging and Dark Adaptation," *Vis. Res.* **39**:3975–3982 (1999).
111. E. Pulos, "Changes in Rod Sensitivity through Adulthood," *Invest. Ophthalmol. Visual Sci.* **30**:1738–1742 (1989).
112. J. F. Sturr, L. Zhang, H. A. Taub, D. J. Hannon, and M. M. Jackowski, "Psychophysical Evidence for Losses in Rod Sensitivity in the Aging Visual System," *Vis. Res.* **37**:475–481 (1997).
113. B. E. Scheffrin, M. L. Bieber, R. McLean, and J. S. Werner, "The Area of Complete Scotopic Spatial Summation Enlarges with Age," *J. Opt. Soc. Am.* **A15**:340–348 (1998).
114. J. E. E. Keunen, D. V. Norren, and G. J. V. Meel, "Density of Foveal Cone Pigments at Older Age," *Invest. Ophthalmol. Visual Sci.* **28**:985–991 (1987).
115. C. A. Curcio, C. L. Millican, K. A. Allen, and R. E. Kalina, "Aging of the Human Photoreceptor Mosaic: Evidence for Selective Vulnerability of Rods in Central Retina," *Invest. Ophthalmol. Visual Sci.* **34**(12):3278–3296(1993).
116. J. J. Plantner, H. L. Barbour, and E. L. Kean, "The Rhodopsin Content of the Human Eye," *Curr. Eye Res.* **7**:1125–1129 (1988).
117. J. Marshall, J. Grindle, P. L. Ansell, and B. Borwein, "Convolution in Human Rods: An Ageing Process," *Br. J. Ophthalmol.* **63**:181–187 (1979).
118. A. Eisner, S. A. Fleming, M. L. Klein, and W. M. Mauldin, "Sensitivities in Older Eyes with Good Acuity: Cross-Sectional Norms," *Invest. Ophthalmol. Visual Sci.* **28**:1824–1831 (1987).
119. C. A. Johnson, A. J. Adams, J. D. Twelker, and J. M. Quigg, "Age-Related Changes of the Central Visual Field for Short-Wavelength-Sensitive Pathways," *J. Opt. Soc. Am.* **A5**:2131–2139(1988).
120. J. S. Werner and V. G. Steele, "Sensitivity of Human Foveal Color Mechanisms Throughout the Life Span," *J. Opt. Soc. Am.* **A5**:2122–2130 (1988).
121. B. E. Scheffrin, J. S. Werner, M. Plach, N. Utlaut, and E. Switkes, "Sites of Age-Related Sensitivity Loss in a Short-Wave Cone Pathway," *J. Opt. Soc. Am.* **A9**:355–363 (1992).
122. J. S. Werner, B. E. Scheffrin, and M. L. Bieber, "Senescence of Foveal and Parafoveal Cone Sensitivities and Their Relations to Macular Pigment Density," *J. Opt. Soc. Am.* **A17**:1918–1932 (2000).
123. R. Lakowski, "Is the Deterioration of Colour Discrimination with Age due to Lens or Retinal Changes?" *Die Farbe* **11**:69–86 (1962).
124. Y. Ohta and H. Kato, "Colour Perception Changes with Age," *Mod. Probl. Ophthalmol.* **17**:345–352 (1976).
125. G. Verriest, "Further Studies on Acquired Deficiency of Color Discrimination," *J. Opt. Soc. Am.* **53**:185–195 (1963).
126. V. C. Smith, J. Pokorny, and A. S. Pass, "Color-Axis Determination on the Farnsworth-Munsell 100-Hue Test," *Am. J. Ophthalmol.* **100**:176–182(1985).
127. K. Knoblauch, F. Saunders, M. Kusuda, R. Hynes, M. Podor, K. E. Higgins, and F. M. de Monasterio, "Age and Illuminance Effects in the Farnsworth-Munsell 100-Hue Test," *Appl. Opt.* **26**:1441–1448 (1987).
128. J. Birch, *Diagnosis of Defective Colour Vision* (Oxford University Press, New York, 1993).
129. B. E. Scheffrin, K. Shinomori, and J. S. Werner, "Contributions of Neural Pathways to Age-Related Losses in Chromatic Discrimination," *J. Opt. Soc. Am.* **A12**:1233–1241 (1995).
130. K. Shinomori, B. E. Scheffrin, and J. S. Werner, "Age-Related Changes in Wavelength Discrimination," *J. Opt. Soc. Am.* **A18**:310–318 (2001).
131. J. M. Kraft and J. S. Werner, "Aging and the Saturation of Colors: 1. Colorimetric Purity Discrimination," *J. Opt. Soc. Am.* **A16**:223–230 (1999).
132. J. S. Werner, "Visual Problems of the Retina during Ageing: Compensation Mechanisms and Colour Constancy across the Life Span," *Prog. Retin. Eye Res.* **15**:621–645 (1996).
133. J. L. Hardy, C. M. Frederick, P. Kay, and J. S. Werner, "Color Naming, Lens Aging, and Grue: What the Optics of the Aging Eye Can Teach Us about Color Language," *Psychol. Sci.* **44**:321–327 (2005).
134. J. S. Werner and B. E. Scheffrin, "Loci of Achromatic Points Throughout the Life Span," *J. Opt. Soc. Am.* **A10**:1509–1516 (1993).
135. P. B. Delahunt, M. A. Webster, L. Ma, and J. S. Werner, "Long-Term Renormalization of Chromatic Mechanisms Following Cataract Surgery," *Vis. Neurosci.* **21**:301–307 (2004).

136. D. G. Pitts, "The Effects of Aging on Selected Visual Functions: Dark Adaptation, Visual Acuity, Stereopsis, and Brightness Contrast," in *Aging and Human Visual Functions*, R. Sekuler, D. Kline, and K. Dismukes, (eds.), (A.R. Liss, New York, 1982), pp. 131–159.
137. C. Owsley, R. Sekuler, and D. Siemsen, "Contrast Sensitivity Throughout Adulthood," *Vision Res.* **23**:689–699 (1983).
138. L. Frisén and M. Frisén, "How Good is Normal Acuity? A Study of Letter Acuity Thresholds as a Function of Age," *Graefes Arch. Clin. Exp. Ophthalmol.* **215**:149–157(1981).
139. D. B. Elliott, K. C. H. Yang, and D. Whitaker, "Visual Acuity Changes Throughout Adulthood in Normal, Healthy Eyes: Seeing beyond 6/6," *Optom. Vis. Sci.* **72**:186–191 (1995).
140. G. Derefeldt, G. Lennerstrand, and B. Lundh, "Age Variations in Normal Human Contrast Sensitivity," *Acta Ophthalmol.* **57**:679–690 (1979).
141. K. E. Higgins, M. J. Jafe, R. C. Caruso, and F. M. deMonasterio, "Spatial Contrast Sensitivity: Effects of Age, Test Retest, and Psychophysical Method," *J. Opt. Soc. Am.* **A5**:2173–2180(1988).
142. U. Tulunay-Keesey, J. N. VerHoeve, and C. Terkla-McGrane, "Threshold and Suprathreshold Spatiotemporal Response Throughout Adulthood," *J. Opt. Soc. Am.* **A5**:2191–2200(1988).
143. D. Elliott, D. Whitaker, and D. MacVeigh, "Neural Contribution to Spatiotemporal Contrast Sensitivity Decline in Healthy Ageing Eyes," *Vis. Res.* **30**:541–547 (1990).
144. B. E. Scheffrin, S. J. Tregear, L. O. Harvey Jr., and J. S. Werner, "Senescent Changes in Scotopic Contrast Sensitivity," *Vis. Res.* **39**:3728–3736 (1999).
145. J. D. Morrison and C. McGrath, "Assessment of the Optical Contributions to the Age-Related Deterioration in Vision," *Q. J. Exp. Psychol.* **70**:249–269 (1985).
146. S. L. Elliott, S. S. Choi, N. Doble, J. L. Hardy, J. W. Evans, and J. S. Werner, "Role of High-Order Aberrations in Senescent Changes in Spatial Vision," *J. Vis.* **9**:1–16 (2009).
147. K. B. Burton, C. Owsley, and M. E. Sloane, "Aging and Neural Spatial Contrast Sensitivity: Photopic Vision," *Vis. Res.* **33**:939–946 (1993).
148. A. Fiorentini, V. Porciatti, M. C. Morrone, and D. C. Burr, "Visual Ageing: Unspecific Decline of the Responses to Luminance and Colour," *Vis. Res.* **36**(21):3557–3566 (1996).
149. J. L. Hardy, P. B. Delahunt, K. Okajima, and J. S. Werner, "Senescence of Spatial Chromatic Contrast Sensitivity. I. Detection under Conditions Controlling for Optical Factors," *J. Opt. Soc. Am.* **A22**:49–59 (2005).
150. N. W. Coppinger, "The Relationship between Critical Flicker Frequency and Chronological Age for Varying Levels of Stimulus Brightness," *J. Gerontol.* **10**:48–52 (1955).
151. C. W. Tyler, "Two Processes Control Variations in Flicker Sensitivity over the Life Span," *J. Opt. Soc. Am.* **A6**:481–490 (1989).
152. C. B. Y. Kim and M. J. Mayer, "Flicker Sensitivity in Healthy Aging Eyes. II. Cross-Sectional Aging Trends from 18 through 77 Years of Age," *J. Opt. Soc. Am.* **A11**:1958–1969 (1994).
153. V. Porciatti, D. C. Burr, C. Morrone, and A. Fiorentini, "The Effects of Ageing on the Pattern Electroretinogram and Visual Evoked Potential in Humans," *Vis. Res.* **32**:1199–1209(1992).
154. M. Ikeda, "Temporal Summation of Positive and Negative Flashes in the Visual System," *J. Opt. Soc. Am.* **A55**:1527–1534 (1965).
155. K. Shinomori and J. S. Werner, "Senescence of the Temporal Impulse Response to a Luminous Pulse," *Vis. Res.* **43**:617–627 (2003).
156. D. C. Burr and M. C. Morrone, "Impulse-Response Functions for Chromatic and Achromatic Stimuli," *J. Opt. Soc. Am.* **A10**:1706–1713 (1993).
157. E. Wolf, "Studies on the Shrinkage of the Visual Field with Age," *Highway Res. Rec.* **167**:1–7 (1967).
158. A. Iwase, Y. Kitazawa, and Y. Ohno, "On Age-Related Norms of the Visual Field," *Jpn. J. Ophthalmol.* **32**:429–437 (1988).
159. C. A. Johnson, A. J. Adams, and R. A. Lewis, "Evidence for a Neural Basis of Age-Related Visual Field Loss in Normal Observers," *Invest. Ophthalmol. Visual Sci.* **30**:2056–2064 (1989).
160. M. Ikeda and T. Takeuchi, "Influence of Foveal Load on the Functional Visual Field," *Percept. Psychophys.* **18**:255–260 (1965).
161. K. K. Ball, B. L. Beard, D. L. Roenker, R. L. Miller, and D. S. Griggs, "Age and Visual Search: Expanding the Useful Field of View," *J. Opt. Soc. Am.* **A5**:2210–2221 (1988).

162. K. O. Devaney and H. A. Johnson, "Neuron Loss in the Aging Visual Cortex of Man," *J. Gerontol.* **35**:836–841 (1980).
163. P. D. Spear, "Neural Bases of Visual Deficits during Aging," *Vis. Res.* **33**:2589–2609 (1993).
164. S. N. Jani, "The Age Factor in Stereopsis Screening," *Am. J. Optom.* **43**:653–657 (1966).
165. M. D. Bell, E. Wolf, and C. D. Bernholz, "Depth Perception as a Function of Age," *Aging Human Develop.* **3**:77–81 (1972).
166. S. Laframboise, D. DeGuise, and J. Faubert, "Effect of Aging on Stereoscopic Interocular Correlation," *Optom. Vis. Sci.* **83**:589–593 (2006).
167. J. M. Enoch, J. S. Werner, G. Haegerstrom-Portnoy, V. Lakshminarayanan, and M. Rynders, "Forever Young: Visual Functions not Affected or Minimally Affected by Aging," *J. Gerontol.: Bio. Sci.* **A54**:B336–B351 (1999).
168. D. Sliney and M. Wolbarsht, *Safety with Lasers and Other Optical Sources* (Plenum Press, New York, 1980).
169. W. J. Donnelly, 3rd, K. Pesudovs, J. D. Marsack, E. J. Sarver, and R. A. Applegate, "Quantifying Scatter in Shack-Hartmann Images to Evaluate Nuclear Cataract," *J. Refract. Surg.* **20**(5):S515–S522 (2004).
170. R. Hess and G. Woo, "Vision Through Cataracts," *Invest. Ophthalmol. Visual Sci.* **17**(5):428–435 (1978).
171. M. J. Costello, S. Johnsen, K. O. Gilliland, C. D. Freel, and W. C. Fowler, "Predicted Light Scattering from Particles Observed in Human Age-Related Nuclear Cataracts Using Mie Scattering Theory," *Invest. Ophthalmol. Visual Sci.* **48**(1):303–312 (2007).
172. A. E. Elsner, S. A. Burns, and J. J. Weiter, "Cone Photopigment in Older Subjects: Decreased Optical Density in Early Age-Related Macular Degeneration," *J. Opt. Soc. Am.* **A19**(1):215–222 (2002).
173. G. Haegerstrom-Portnoy and B. Brown, "Two-Color Increment Thresholds in Early Age Related Maculopathy," *Clin. Vis. Sci.* **4**:165–172 (1989).
174. J. S. Sunness, R. W. Massof, M. A. Johnson, N. M. Bressler, S. B. Bressler, and S. L. Fine, "Diminished Foveal Sensitivity may Predict the Development of Advanced Age-Related Macular Degeneration," *Ophthalmology* **96**:375–381 (1989).
175. R. A. Applegate, A. J. Adams, J. C. Cavender, and F. Zisman, "Early Color Vision Changes in Age-Related Maculopathy," *Appl. Opt.* **26**:1458–1462 (1987).
176. National Society to Prevent Blindness, "Prevent Blindness America" (Schaumburg, IL, 1994).
177. Diabetes Control and Complications Trial Research Group, "The Effect of Intensive Treatment of Diabetes on the Development and Progression of Long-Term Complications in Insulin-Dependent Diabetes Mellitus," *New Eng. J. Med.* **329**:977–986 (1993).
178. V. J. Volbrecht, M. E. Schneck, A. J. Adams, J. A. Linfoot, and E. Ai, "Diabetic Short-Wavelength Sensitivity: Variations with Induced Changes in Blood Glucose Level," *Invest. Ophthalmol. Visual Sci.* **35**:1243–1246 (1994).
179. R. Lakowski, P. A. Aspinall, and P. R. Kinnear, "Association between Colour Vision Losses and Diabetes Mellitus," *Ophthalm. Physiol. Opt.* **4**:145–159 (1972/73).
180. Y. Yang, J. C. Downs, C. Girken, L. Sakata, A. Bellezza, H. Thompson, and C. F. Burgoyne, "3-D Histomorphometry of the Normal and Early Glaucomatous Monkey Optic Nerve Head: Lamina Cribosa and Peripapillary Scleral Position and Thickness," *Invest. Ophthalmol. Visual Sci.* **48**:4597–4607 (2007).
181. J. Flammer and S. Orgül, "Optic Nerve Blood-Flow Abnormalities in Glaucoma," *Prog. Retin. Eye Res.* **17**:267–289 (1998).
182. C. A. Johnson and S. J. Samuels, "Screening for Glaucomatous Visual Field Loss with Frequency-Doubling Perimetry," *Invest. Ophthalmol. Visual Sci.* **38**:413–425 (1997).
183. H. Sun, M. W. Dul, and W. H. Swanson, "Linearity Can Account for the Similarity among Conventional, Frequency-Doubling, and Gabor-Based Perimetric Tests in the Glaucomatous Macula," *Optom. Vis. Sci.* **83**(7):455–465 (2006).
184. S. J. Chew and R. Ritch, "Neuroprotection: The Next Breakthrough in Glaucoma? Proceedings of the 3rd Annual Optic Nerve Rescue and Restoration Think Tank," *J. Glaucoma* **6**:263–266 (1997).
185. A. L. Cohen, "Diffractive Bifocal Lens Designs," *Optom. Vis. Sci.* **70**(6):461–468 (1993).
186. M. A. Bullimore and R. J. Jacobs, "Subjective and Objective Assessment of Soft Bifocal Contact Lens Performance," *Optom. Vis. Sci.* **70**(6):469–475 (1993).
187. A. Bradley, H. Abdul Rahman, P. S. Soni, and X. Zhang, "Effects of Target Distance and Pupil Size on Letter Contrast Sensitivity with Simultaneous Vision Bifocal Contact Lenses," *Optom. Vis. Sci.* **70**(6):476–481 (1993).

188. S. Sanislo, D. Wicker, and D. G. Green, "Contrast Sensitivity Measurements with the Echelon Diffractive Bifocal Contact Lens as Compared to Bifocal Spectacles," *Clao J.* **18**(3):161–164 (1992).
189. G. E. Legge and G. S. Rubin, "Binocular Interactions in Suprathreshold Contrast Perception," *Percept. Psychophys.* **30**(1):49–61 (1981).
190. C. Schor, L. Landsman, and P. Erickson, "Ocular Dominance and the Interocular Suppression of Blur in Monovision," *Am. J. Optom. Physiol. Opt.* **64**(10):723–730 (1987).
191. R. Legras, V. Hornain, A. Monot, and N. Chateau, "Effect of Induced Anisometropia on Binocular Through-focus Contrast Sensitivity," *Optom. Vis. Sci.* **78**(7):503–509 (2001).
192. A. Back, T. Grant, and N. Hine, "Comparative Visual Performance of Three Presbyopic Contact Lens Corrections," *Optom. Vis. Sci.* **69**(6):474–480 (1992).
193. J. T. Holladay, H. Van Dijk, A. Lang, V. Portney, T. R. Willis, R. Sun, and H. C. Oksman, "Optical Performance of Multifocal Intraocular Lenses," *J. Cataract Refract. Surg.* **16**(4):413–422 (1990).
194. M. Simpson, "Diffractive Multifocal Intraocular Lens Image Quality," *Appl. Opt.* **31**(19):3621–3626 (1992).
195. E. Peli and A. Lang, "Appearance of Images Through a Multifocal Intraocular Lens," *J. Opt. Soc. Am.* **A18**(2):302–309 (2001).
196. J. C. Javitt, F. Wang, D. J. Trentacost, M. Rowe, and N. Tarantino, "Outcomes of Cataract Extraction with Multifocal Intraocular Lens Implantation," *Ophthalmology* **104**:589–599 (1997).
197. A. Glasser, "Restoration of Accommodation: Surgical Options for Correction of Presbyopia," *Clin. Exp. Optom.* **91**(3):279–295 (2008).
198. O. Nishi and K. Nishi, "Accommodation Amplitude after Lens Refilling with Injectable Silicone by Sealing the Capsule with a Plug in Primates," *Arch. Ophthalmol.* **116**(10):1358–1361 (1998).
199. P. L. Kaufman, "Scleral Expansion Surgery for Presbyopia," *Ophthalmology* **108**(12):2161–2162 (2001).
200. L. A. Ostrin, S. Kasthurirangan, and A. Glasser, "Evaluation of a Satisfied Bilateral Scleral Expansion Band Patient," *J. Cataract Refract. Surg.* **30**(7):1445–1453 (2004).
201. S. D. McLeod, V. Portney, and A. Ting, "A Dual Optic Accommodating Foldable Intraocular Lens," *Br. J. Ophthalmol.* **87**(9):1083–1085 (2003).
202. S. D. McLeod, "Optical Principles, Biomechanics, and Initial Clinical Performance of a Dual-Optic Accommodating Intraocular Lens (an American Ophthalmological Society thesis)," *Trans. Am. Ophthalmol. Soc.* **104**:437–452 (2006).
203. O. Findl and C. Leydolt, "Meta-Analysis of Accommodating Intraocular Lenses," *J. Cataract Refract. Surg.* **33**(3):522–527 (2007).
204. A. Langenbucher, B. Seitz, S. Huber, N. X. Nguyen, and M. Kuchle, "Theoretical and Measured Pseudophakic Accommodation after Implantation of a New Accommodative Posterior Chamber Intraocular Lens," *Arch. Ophthalmol.* **121**(12):1722–1727 (2003).
205. J. S. Cumming, S. G. Slade, and A. Chayet, "Clinical Evaluation of the Model AT-45 Silicone Accommodating Intraocular Lens: Results of Feasibility and the Initial Phase of a Food and Drug Administration Clinical Trial," *Ophthalmology* **108**(11):2005–2009; discussion 2010 (2001).
206. G. Marchini, P. Mora, E. Pedrotti, F. Manzotti, R. Aldigeri, and S. A. Gandolfi, "Functional Assessment of two Different Accommodative Intraocular Lenses Compared with a Monofocal Intraocular Lens," *Ophthalmology* **114**(11):2038–2043 (2007).

ADAPTIVE OPTICS IN RETINAL MICROSCOPY AND VISION

Donald T. Miller

*School of Optometry
Indiana University
Bloomington, Indiana*

Austin Roorda

*School of Optometry
University of California
Berkeley, California*

15.1 GLOSSARY

Adaptive optics. Opto-electro-mechanical method for dynamic closed-loop compensation of wavefront aberrations that reside between the object and image.

Bimorph mirror. Class of wavefront corrector distinguished by its piezoelectric material that is sandwiched between a continuous top electrode (reflective) and a bottom, patterned electrode array.

Contrast sensitivity. Reciprocal of the minimum perceptible contrast.

Deformable mirror. Mirror whose surface profile is controlled by applying voltages to a set of adjacent electrodes or actuators.

Discrete-actuator deformable mirror. Class of wavefront corrector distinguished by its continuous reflective surface and array of underlying actuators that alter the surface shape.

Higher-order aberrations. Aberrations higher than second order.

Influence function. Deflection of the corrector surface when a unit voltage is applied to a single actuator.

Low-order aberrations. Aberrations of first and second order.

Membrane mirror. Class of wavefront corrector distinguished by its edge-clamped, flexible, reflective membrane that is sandwiched between a transparent top electrode and an underlying array of electrodes.

Ophthalmoscope. Instrument for viewing and recording images of the back of the eye.

Optical coherence tomography. Interferometric-based imaging modality that optically sections tissue with micron-scale axial resolution. Axial resolution is inversely related to the spectrum of the imaging light source.

Photoreceptor. Cell that initiates the visual process in the retina by transducing the captured photons into an encoded neural signal.

- Point spread function.** Response (image) of an optical system to a point source of light (object).
- Psychophysics.** A branch of science that deals with the relationship between physical stimuli and sensory response.
- Retinal microscopy.** Observation of retinal features at the cellular level in the living eye.
- Scanning laser ophthalmoscope.** Imaging modality that raster scans a point source across the object and filters the reflected light with a confocal pinhole for enhanced contrast and axial sectioning.
- Segmented corrector.** Class of wavefront corrector distinguished by its array of small closely spaced reflective or transmissive elements that induce local piston and tilt on the incident wavefront.
- Shack-Hartmann wavefront sensor.** An objective wavefront sensor that measures the wave aberrations of the eye with principle components being a light source beacon, lenslet array, and areal CCD.
- Spectacles.** Ophthalmic appliance that corrects prism, defocus (sphere), and astigmatism (cylinder).
- Strehl ratio.** Ratio of the maximum light intensity in the aberrated point spread function to that in the aberration-free point spread function with the same pupil diameter. Values range from 0 to 1.
- Visual acuity.** Capacity for seeing distinctly the details of an object. In units of the Snellen fraction, average clinical visual acuity varies between 20/15 and 20/20.
- Wave aberration.** Optical deviations of a wavefront that degrade image quality and that can vary temporally and spatially across the pupil of the optical system.
- Zernike polynomials.** Specific set of polynomials frequently used to mathematically represent the wave aberration of the eye.

15.2 INTRODUCTION

Vision is a most acute human sense. Yet one need only to examine the first step in the visual process—the formation of an image on the retina—to realize that it is often far from perfect. Image quality at the retina is limited fundamentally by the wave aberrations intrinsic to the cornea and crystalline lens and diffraction due to the finite size of the eye's pupil. Together, aberrations and diffraction limit not only what the eye sees looking out, but also determine the smallest internal structures that can be observed when looking into the eye with a microscope. Conventional corrective methods such as spectacles, contact lenses, and refractive surgery provide a static amelioration of low-order sphere (defocus) and cylinder (astigmatism). However, ocular image quality can be significantly improved by dilating the pupil to minimize diffraction and correcting the aberrations across the larger pupil. The presence of aberrations in the human eye¹ has been recognized for some time. Indeed, Hermann von Helmholtz, commenting on the eye 150 years ago, put it as follows: “Now, it is not too much to say that if an optician wanted to sell me an instrument which had all these defects, I should think myself quite justified in blaming his carelessness in the strongest terms, and giving him back his instrument.”²

Spectacles have been used to correct defocus since at least as early as the thirteenth century. They began to be used to correct astigmatism shortly after 1801, when Thomas Young discovered that condition in the human eye. However, it has only been in the last decade and a half—two centuries later—that additional, higher-order aberrations have been corrected. There are three reasons for this long period of seemingly little progress.

First, the most significant and troublesome aberrations in the eye are defocus and astigmatism. Correcting these two usually improves vision to an acceptable level. Second, for most of this period defocus and astigmatism were the only aberrations of the eye that could be easily measured. Third, even when higher-order aberrations could be measured by the cumbersome techniques available at that time, there was no simple or economical way to correct them. In 1961, for instance, M. S. Smirnov³ measured many of the higher-order aberrations in the human eye for the first time. But his psychophysical method required 1 to 2 h of measurements per eye, followed by an additional 10 to 12 h of calculations. On the strength of his results, Smirnov surmised that it would be possible to manufacture custom lenses to compensate for the higher-order aberrations of individual eyes.

Recent technological advances, notably in adaptive optics (AO),⁴ have provided solutions to the problems Smirnov encountered. The concept of AO was proposed in 1953 by astronomer Horace Babcock⁵ as a means of compensating for the wavefront distortion induced by atmospheric turbulence.⁶ Transforming AO into a noninvasive vision tool fundamentally consists of measuring the aberrations of the eye and then compensating for them. But the technique must be able to deal with the problem of variation of ocular and retinal tissue for a given person, as well as from person to person, and it must also take into account human safety.

Prior to the development of AO for the eye, several techniques showed promise for enhancing vision and improving the resolution of retinal imaging. These included use of interference fringes to probe the visual system without optical blurring;⁷ a postdetection imaging technique coined speckle interferometry to recover high spatial frequencies of the cone mosaic;⁸ and precise correction of defocus and astigmatism in conjunction with quick flashes of quasi-monochromatic light to image individual cone photoreceptors.⁹ Success with the latter was confined to eyes of relatively low higher-order aberrations. In 1989, the first wavefront corrector was applied to the eye, a 13-actuator membrane mirror that corrected the astigmatism in one subject's eye using his conventional prescription for spectacles.¹⁰ Correction of additional aberrations was not possible as there was no effective means at that time for measuring the aberrations.

These early techniques shared the common problem that the aberrations beyond defocus and astigmatism were largely unknown. To address this limitation, Liang et al.¹¹ developed a wavefront sensor based on the principles of Shack-Hartmann wavefront sensing (SHWS) and further refined it for measuring aberrations up through the first 10 radial orders (corresponding to the first 65 Zernike modes).¹² The SHWS was a watershed. It provided—for the first time—a rapid, automated, noninvasive, and objective measure of the wave aberrations of the eye. The method required no feedback from the patient; its measurements could be collected in a fraction of a second; and its near-infrared illumination was less intrusive to the subject than earlier techniques had been. That same year, the first AO retina camera was developed based on a SHWS and a 37-actuator deformable mirror.¹³ The system included an additional channel for collecting images of the living retina and measuring the visual performance of subjects. With a pupil 6 mm in diameter, the system substantially corrected aberrations up through the fourth-Zernike order, where most of the significant aberrations in the eye reside. Thus, Liang et al. demonstrated for the first time that AO allows coma, spherical aberration, and other irregular aberrations to be corrected in the eye. The closed-loop speed of their system was on the order of tens of seconds and was incapable of tracking temporal fluctuations in the eye's aberrations. Their system may be more accurately characterized as an active rather than adaptive optics system.

In little more than a decade since AO for the eye has experienced exponential growth.¹⁴ Essentially all facets of AO have undergone substantive development for use with the eye, including its wavefront sensor, wavefront corrector, and control algorithm. In addition, AO has been integrated into a variety of retina camera architectures to increase their resolution and sensitivity including the three principle types: flood-illuminated ophthalmoscope,^{13,15–18} scanning laser ophthalmoscope (SLO),^{19–24} and optical coherence tomography (OCT).^{25–35} AO for the eye has also entered the commercial arena. The SHWS for measuring ocular aberrations has been commercially available for several years; commercial AO systems for vision testing for the last couple of years; and commercial AO retina cameras now entering the market.

While early AO retina cameras focused on imaging cone photoreceptors owing to their high reflectance and contrast, many other cellular structures have since been imaged in the living human retina including nerve fiber bundles, individual blood cells flowing through the smallest retinal capillaries, fibers of Henle, retinal pigment epithelium cells, and astrocytes. Ganglion cells and their axons have also been visualized with AO in primates when used in conjunction with a contrast-enhancing fluorescent dye. There is also an ever expanding list of fundamental and clinical research applications that are underway and which have been made possible by the increased resolution and sensitivity afforded by AO.

Several reviews of AO for vision science already exist,^{36–40} including a textbook devoted to the topic.⁴¹ Nevertheless, the field continues to grow rapidly and much has changed since the first *OSA Handbook* chapter on this topic was written in 2000.⁴² Our aim here is to summarize the current understanding of the ocular aberration properties (pertinent for designing AO systems for the eye), detail the status of AO technologies for the eye, and review some of the latest scientific and clinical uses of this powerful technology.

15.3 PROPERTIES OF OCULAR ABERRATIONS

Aberrations of the eye vary spatially across the eye's pupil, vary over time, and vary with field location at the retina (isoplanatism). These fundamental properties (spatial distribution, temporal distribution, and field dependence) dictate the performance requirements of AO for effective correction of ocular aberrations. Here we summarize the most pertinent aspects of each for AO.

The spatial distribution of ocular aberrations varies considerably among subjects, and often does not correspond to the classical aberrations of manufactured optical systems (e.g., Seidel aberrations). Instead, Zernike polynomials are routinely used for representing ocular aberrations. The first-order Zernike polynomials represent tilt, the second order correspond to defocus and astigmatism, the third order are coma and trefoil, the fourth order are spherical aberration, secondary astigmatism, and quadrafoil, and so on. Third order and higher are collectively termed higher-order aberrations or irregular aberrations. Zernikes are not the most efficient polynomial representation of the eye's aberrations, but are only slightly less compact than the principal modes that derive from a principal components analysis.^{43,44} Zernikes are also mathematically simpler and have well-established properties. Certainly the long history of Zernike polynomials in astronomy for representing atmospheric turbulence⁴⁵ influenced their early use by the vision community. The now ubiquitous use of Zernike polynomials is perhaps best reflected by the community's establishment and universal acceptance of a single naming convention for the Zernike coefficients and polynomials as described by the OSA/VSIA Standards Taskforce.⁴⁶

Today detailed measurements of the spatial distribution of wave aberrations have been made in large populations of normal, healthy adult eyes with large pupils,^{43,44,47} including the pooling of data collected at 10 laboratories (2560 eyes).⁴⁸ More recent studies have expanded this to include the effect of age, accommodation, refractive state, refractive surgery, and disease.^{49–57}

The spatial characteristics of the ocular aberrations most important for AO are spatial fidelity (frequency composition of the aberrations) and magnitude [peak-to-valley (PV) wavefront error]. Both are captured by the Zernike representation—spatial fidelity by the Zernike order and magnitude by the Zernike coefficient, and both are strongly influenced by pupil size. As an example, Fig. 1*a* shows the wavefront variance decomposed by Zernike order for a population of about 100 normal subjects, most in their early twenties.⁴⁴ Data points are shown for three pupil sizes: 4.5, 6.0, and 7.5 mm, which span the range over which AO is typically applied to the human eye. For comparison, the maximum physiological pupil size for young subjects is nominally 8 mm. In Fig. 1*a*, the power decreases monotonically (approximately linear on a logarithmic scale) with second-order aberrations dominating the total wavefront error. Power also decreases monotonically with decreasing pupil size. The 7.5-mm pupil represents the most demanding condition for AO. For this pupil size, correction of Zernike polynomials up through at least 10th order is necessary to reach diffraction-limited imaging [$\lambda/14$ root-mean-square (RMS) error] in approximately 95 percent of the population, that is, two times the standard deviation of the \log_{10} (wavefront variance).

Figure 1*b* shows the corresponding PV wavefront error that encompasses 95 percent of the population as a function of pupil size (4.5 to 7.5 mm). Three curves are shown that correspond to three different second-order states. As expected, the PV error increases monotonically with pupil size and depends strongly on the second-order state. The PV error for the 7.5-mm pupil ranges from 7 to 11 μm depending on the second-order state. The largest error of 11 μm represents the most demanding condition for AO correction. Note that for this study, each subject was meticulously refracted with trial lenses. Typically, AO systems operate under less ideal conditions in which the second-order aberrations are only coarsely corrected or not at all. Under these more realistic conditions, second-order aberrations can readily surpass the PV error shown in the figure by several times. As such, careful design of an AO system must take into account the anticipated refractive state of the population and the expected manner in which trial lenses or translating lenses will be applied.

The population measurements in Fig. 1 represent essentially static wave aberrations and therefore do not capture the temporal behavior of the ocular media. Two early studies investigated the temporal dynamics with attention to their impact on AO performance.^{58,59} Both studies employed a SHWS to track the dynamics. Hofer et al. measured aberration changes up to approximately 5 to 6 Hz, and

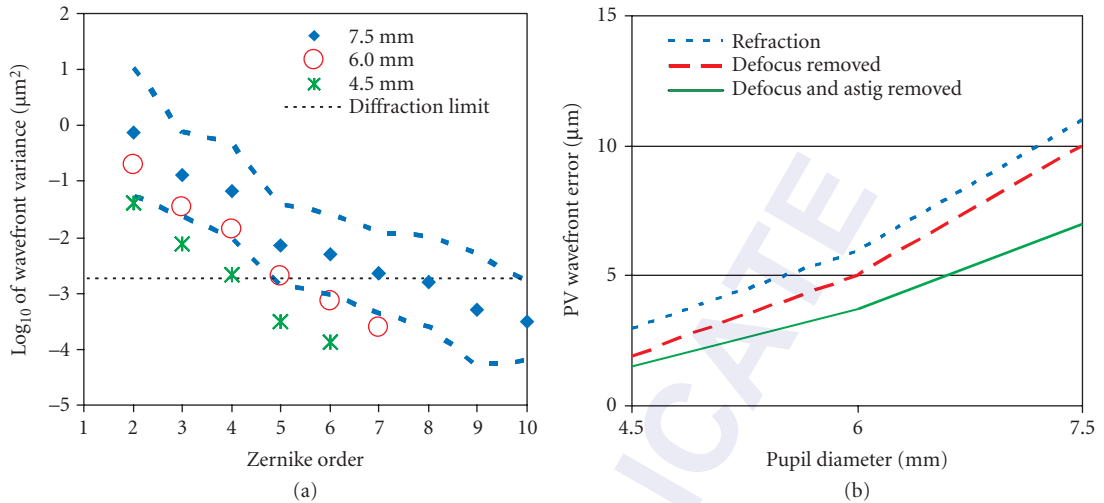


FIGURE 1 Spatial properties of ocular aberrations in a large population of 100 normal eyes.⁴⁴ (a) \log_{10} of the wavefront variance after a conventional refraction using trial lenses is plotted as a function of Zernike order and pupil size (4.5, 6.0, and 7.5 mm). Diamond and corresponding dashed curves represent the mean and mean \pm two times the standard deviation of the \log_{10} (wavefront variance), respectively, for a 7.5-mm pupil. Star and open circle correspond to 4.5- and 6.0-mm pupils. Thin, horizontal dashed line corresponds to $\lambda/14$ RMS for $\lambda = 0.6 \mu\text{m}$. (b) PV wavefront error that encompasses 95 percent of the population is plotted as a function of pupil diameter. Three second-order states are shown: (i) residual aberrations after a conventional refraction using trial lenses (short blue dashed line) (ii) all aberrations present with zeroed defocus (long red dashed line), and (iii) all aberrations present with zeroed defocus and astigmatism (solid green line).

Lois-Diaz et al. up to approximately 30 Hz. Figure 2a shows representative temporal traces of the total RMS wavefront error and several Zernike terms for one subject and a model eye over a 5 s interval. Figure 2b shows corresponding average temporal power spectra which reveal the frequency content of the traces. As depicted in the figure, temporal fluctuations are found in all of the eye's aberrations, not just defocus, even when the eye's accommodation is paralyzed. In both studies, the reported temporal fluctuations decreased at $f^{-4/3}$ or equivalently four dB/octave, where f is the temporal frequency. This decrease is evident in the temporal power spectrum shown in Fig. 2b. The vast majority of the aberration power lies below 1 to 2 Hz, suggesting effective AO correction needs only a temporal bandwidth of a couple of hertz. For comparison, the 1 to 2 Hz is approximately two orders smaller than that of atmospheric turbulence for ground-based telescopes.

Finally, aberrations of the eye also vary with field location at the retina. That is, image quality (or more precisely the point spread function or optical transfer function) at one point on the retina is not the same as at another point when the two points are separated sufficiently. This difference in image quality stems from the fact that the ocular aberrations originate at the cornea and crystalline lens rather than at the pupil of the eye. Rays originating from different field locations on the retina will therefore take slightly different paths through the ocular media, accumulating different phase delays or equivalently different aberrations. In general there is a local region about any point on the retina in which the path differences are sufficiently small that the ocular aberrations are effectively constant. This region is termed the isoplanatic patch.⁶⁰ The diameter of the patch depends on the properties of the eye, but also on the definition of image quality, as for example, use of the stringent Maréchal criterion, $\lambda/14$ - μm RMS (Strehl ratio = 0.8) yields a narrower isoplanatic diameter than the more relaxed criterion of 1 rad² (Strehl ratio = 0.37). Which criterion to choose depends on the application, though it remains to be seen what criterion is optimal for common AO retinal imaging uses as well as AO psychophysical uses.

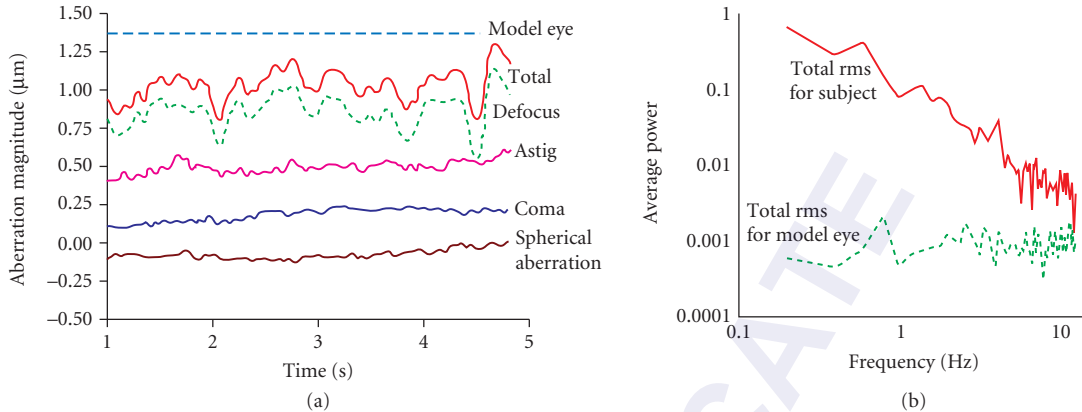


FIGURE 2 Temporal properties of ocular aberrations.⁵⁸ (a) Temporal traces of the total RMS wavefront error and Zernike terms: defocus, astigmatism, coma, and spherical aberration for one subject. A trace of the total RMS wavefront error for an artificial eye is also shown and reflects the sensitivity of the instrument. (b) Average temporal power spectra are shown for the fluctuations in the total RMS wavefront error for a model eye and one subject whose accommodation was paralyzed. Aberrations were computed for a 4.7-mm pupil size. (Reproduced with permission from Ref. 58.)

Anecdotal evidence in the AO literature indicates the isoplanatic patch diameter for the human eye is a couple of degrees. Experiments to directly measure the patch size have only recently been conducted. A detailed survey of efforts in this area coupled with a more thorough analysis were recently published by Bedggood et al.⁶¹ Figure 3 shows a representative result from this study, plotting isoplanatic patch diameter as a function of the residual wavefront RMS at the patch edge. As expected, the patch size increases as the criterion relaxes, that is, larger residual RMS. As an example, AO systems

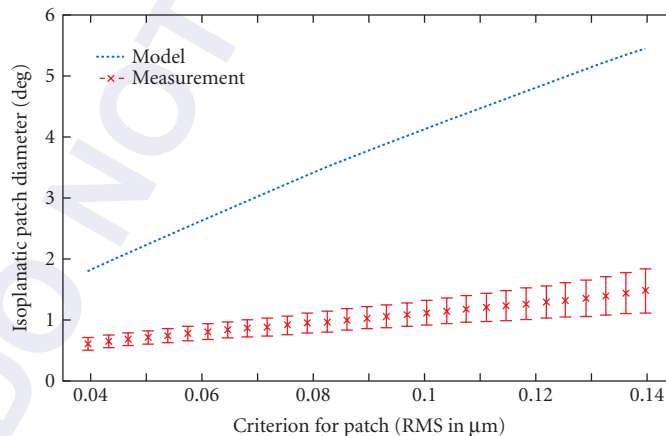


FIGURE 3 Isoplanatic properties of ocular aberrations. Isoplanatic patch diameter is plotted versus residual wavefront RMS at the patch edge after perfect correction of ocular aberrations at the patch center. Data points are average measurements across seven healthy subjects, aged 22 to 39, with error bars at ± 1 standard deviation. Predicted patch diameter is also shown using the Liou Brennan schematic eye.⁶² Pupil size for both is 6 mm. (Reproduced with permission from Ref. 61.)

typically correct the eye (with a large pupil) to within a wavefront RMS error of approximately $0.1 \mu\text{m}$. Using this value as a realistic criterion for the residual RMS at the patch edge in Fig. 3, the corresponding isoplanatic diameter is approximately 1.1° , consistent with anecdotal evidence and other measurements reported in the literature. For comparison, this patch diameter is three orders of magnitude larger than that encountered with ground-based telescopes. Interestingly, optical models of the eye (e.g., the Liou Brennan schematic eye in Fig. 3) are not sufficiently developed yet to predict experimental findings, both in terms of the patch size and criterion dependence. As evident in Fig. 3, the Liou Brennan model overestimates the patch diameter by several times.

15.4 IMPLEMENTATION OF AO

The principle components of an AO system are the wavefront sensor, wavefront corrector, and control system. Figure 4 shows a simplified schematic of AO applied to the eye and highlights the principle components. The effectiveness of AO depends fundamentally on its ability to measure, correct,

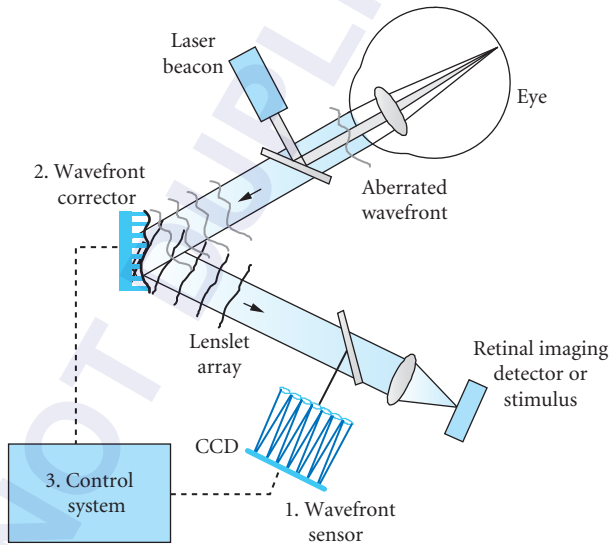


FIGURE 4 Concept schematic of AO applied to the eye. The three principle components of AO are (1) the wavefront sensor, (2) the wavefront corrector, and (3) the control system. The wavefront sensor (shown here as a Shack-Hartmann implementation) works by focusing a laser beam onto the subject's retina. This is analogous to the artificial-guide-star approach now being used to measure the wave aberrations of atmospheric turbulence.⁶⁵ The reflection from the retinal spot is distorted as it passes back through the refracting media of the eye. A two-dimensional lenslet array, placed conjugate with the eye's pupil, samples the exiting wavefront forming an array of images of the retinal spot. A CCD sensor records the displacement of the spots, from which first local wavefront slopes and then global wavefront shape are determined. Wavefront compensation is realized with a deformable mirror. The mirror lies in a plane conjugate with the subject's pupil and the lenslet array of the SHWS. For retinal imaging, a light source (not shown) illuminates the retina, some of which is reflected out of the eye, reflects from the wavefront corrector, and forms an aerial image at the retinal imaging detector. For visual testing, the retinal imaging detector is replaced with a visual stimulus that is viewed by the subject through the AO system.

and track the ocular aberrations, the extent to which depends heavily on the spatial and temporal properties of the aberrations as described in the previous section. Measurement is accomplished by the wavefront sensor, correction by the wavefront corrector, and tracking by the control system in conjunction with the sensor and corrector. The status of each of these components and their requirements for effective use with the eye are described in order below.

Wavefront sensor

There is an abundance of optical sensors for wavefront measuring with many representing highly mature devices. These include the Shack-Hartman sensor, shearing interferometer, pyramid sensor, curvature sensor, phase diversity, laser ray tracing, phase-shifting interferometer, and common path interferometer. While several of these have been applied to the eye and others continue to attract interest, to date all operational AO systems for the eye have been constructed around the SHWS. As such, we confine our attention to this sensor type, though much of which will be applicable to the other sensors. The discussion that follows will reference Fig. 5, which shows an enlarged schematic of the SHWS (lenslet and CCD array) along with typical values of several key parameters.

A first-order design of the SHWS entails evaluation of several key properties: (1) spatial sampling, (2) dynamic range, (3) accuracy, and (4) signal-to-noise. These are primarily set by the characteristics of the lenslet and CCD arrays of the sensor, and therefore selection of these two components is critical. Figure 5*b* shows typical values in use. Spatial sampling of the sensor is specified by the number of lenslets that fills the pupil of the eye. The necessary number of lenslets for reliable representation of the ocular aberrations depends on the composition of the ocular aberrations, pupil size of the eye, and to a lesser extent the SH wavelength. It should be noted that while we assume here the lenslets to be uniformly distributed (as is the case of SH sensors in current AO systems), several nonuniform distributions have also been evaluated.⁶⁴

Theoretical modeling in conjunction with measured aberrations of the eye indicates a one-to-one relationship between the number of Zernike modes and number of lenslets necessary to reliably reconstruct these modes.⁶⁵ Therefore, the number of required lenslets that fill the pupil must be at

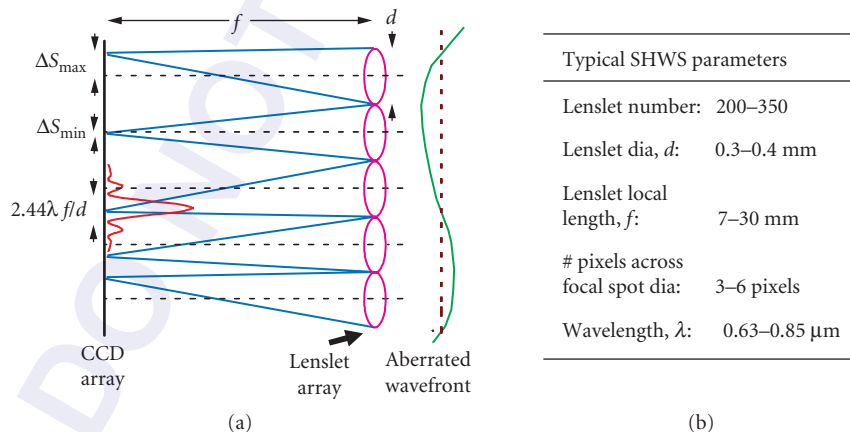


FIGURE 5 (a) A two-dimensional lenslet array, placed conjugate with the eye's pupil, samples the exiting wavefront forming an array of images of the retinal spot. A CCD sensor records the displacement of the spots. Maximum and minimum detectable displacements of the spots are indicated by ΔS_{\max} and ΔS_{\min} . (b) Key SHWS parameters and typical values for use with the eye. Lenslet number refers to the number that fills the pupil of the eye. Lenslet diameter is referenced to the eye pupil. Pixels refer to the detector unit of the CCD array. λ is the wavelength of the SHWS laser beacon.

least equal to the number of Zernike modes, which is related to the maximum Zernike order by $(N+1)(N+2)/2 - 3$. N is the maximum Zernike order and piston, tip, and tilt are not included. Based on this, reliable representation up through 10th Zernike order requires at least 63 lenslets, that is, equals $(10+1)(10+2)/2 - 3$, that uniformly cover the pupil.

Figure 1a provides additional information—at least for normal healthy eyes—as to the number of Zernike orders necessary for reliably representing ocular aberrations. As an example, for the stringent condition of diffraction-limited imaging through a large 7.5-mm pupil in 95 percent of the population, Fig. 1a indicates Zernike representation up through at least 10th order is necessary. For this scenario, at least 63 lenslets are therefore required to reliably sample the aberrations. Fewer lenslets are necessary for a smaller pupil or a less stringent image quality metric. In practice, SHWSs routinely sample the pupil much higher. There is little cost to do so, since there is an over abundance of light for the sensor. Oversampling makes the measurements more robust (to pupil edge effects, eye motion, drying of the precorneal tear film, and system noise) and applicable to eyes with a wider range of aberrations.

Dynamic range refers to the maximum wavefront slope $\Delta\theta_{\max}$ that can be measured by the SHWS (depicted in Fig. 5a). $\Delta\theta_{\max}$ is typically expressed as $\Delta S_{\max}/f$, where ΔS_{\max} is the maximum displacement of the lenslet spot (= lenslet radius) and f the lenslet focal length. Since spherical and cylindrical refractive errors typically consume much of the dynamic range of the sensor, it is often more useful to convert ΔS_{\max} to a maximum measurable defocus (diopters), expressed by $D = \Delta\theta_{\max} / (\text{pupil radius})$. For example, a lenslet diameter of 0.4 mm and a lenslet focal length of 24 mm, gives a $\Delta\theta_{\max}$ and D of 8.3 mrad and 2.45 diopters, respectively, for a 6.8-mm eye pupil. Note that D represents the dioptric range at the wavefront sensor. The corresponding range at the eye is often different depending on the eye-to-sensor magnification. In general, the dynamic range can be increased by increasing the lenslet diameter, decreasing the lenslet focal length, or decreasing the pupil diameter. The pupil diameter, however, is usually not adjustable, being defined early in the design process of the retina camera.

Accuracy refers to the minimum wavefront slope, $\Delta\theta_{\min}$, that can be measured by the SHWS. $\Delta\theta_{\min}$ can be expressed as $\Delta S_{\min}/f$, where ΔS_{\min} is the minimum detectable displacement of the lenslet spot (see Fig. 5a). ΔS_{\min} depends on the pixel size of the SHWS detector (typically a CCD array), diameter of lenslet spot at the detector (which is a function of the lenslet diameter and focal length), accuracy of the centroiding and thresholding algorithms, and signal-to-noise of the captured SHWS spot image. In general, subpixel accuracy is routinely achieved.⁶⁶ The above example (lenslet diameter = 0.4 mm and focal length = 24 mm) with a CCD pixel size of 26 μm yields approximately 4.4 pixels across the focal spot and $\Delta\theta_{\min} \ll 1.1$ mrad. While the SHWS parameters can be manipulated to increase accuracy, dynamic range of the sensor depends inversely on the lenslet diameter and focal length, thereby creating a clear trade-off between dynamic range and accuracy. Additional discussion on this topic can be found in Ref. 66.

Signal-to-noise of the captured SHWS spot images is limited by the amount of light that can be safely directed into the eye, number of pixels per lenslet focal spot, quantum efficiency of the SHWS CCD, and throughput efficiency the SHWS. Unlike astronomical AO systems, however, that operate under light-starved conditions and covet each photon to achieve the approximately 100 photons/lenslet required to close the AO loop, vision AO systems operate under comparatively relaxed light requirements. A typical vision AO system has a minimum of approximately 500,000 detected photons/lenslet, this with hundreds of lenslets that sample the pupil. Such high light levels permit use of relatively low quantum efficient CCDs and tens of pixels that sample the core of the lenslet focal spot.

Wavefront Corrector

The effectiveness of AO to correct ocular aberrations is typically limited by the performance of the wavefront corrector, which is also by far the most expensive component. This device dynamically imparts an ideally conjugate aberration profile onto the passing wavefront, thus canceling the original aberrations. Numerous types of wavefront correctors have been employed in AO systems for the eye. Historically, correctors have been selected somewhat arbitrarily for the eye, that is, driven

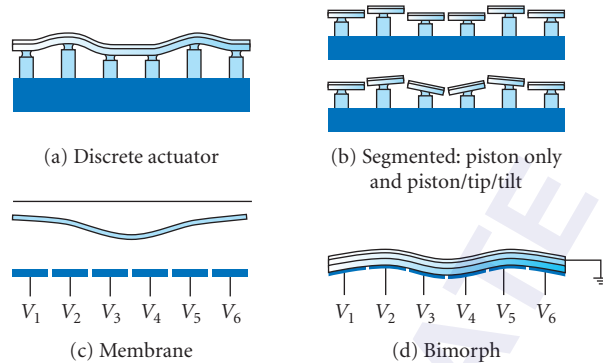


FIGURE 6 Four main classes of wavefront correctors. (a) Discrete actuator deformable mirrors consist of a continuous, reflective surface, and an array of actuators, each capable of producing a local deformation in the surface. (b) Piston-only segmented correctors consist of an array of small planar mirrors whose axial motion (piston) is independently controlled. LC-SLMs modulate the wavefront in a similar fashion, but rely on changes in refractive index rather than the physical displacement of a mirror surface. Piston/tip/tilt segmented correctors add independent tip and tilt motion to the piston-only correctors. (c) Membrane mirrors consist of a grounded, flexible, reflective membrane sandwiched between a transparent top electrode, and an underlying array of patterned electrodes, each of which is capable of producing a global deformation in the surface. (d) Bimorph mirrors consist of a layer of piezoelectric material sandwiched between a continuous top electrode and a bottom, patterned electrode array. A top mirrored layer is added to the top continuous electrode. An applied voltage causes a deformation of the top mirrored surface. (Reproduced with permission from Ref. 41, chap. 4, Fig. 4.2.)

primarily by availability and cost rather than performance, with the expectation that image quality will improve, the extent of which was empirically determined. More recently, several extensive theoretical studies have evaluated the performance of general wavefront corrector classes^{67,68} as well as specific commercial devices.^{69,70} The four main classes of wavefront correctors are depicted in Fig. 6. These include discrete actuator, membrane, bimorph, and segmented piston and piston/tip/tilt wavefront correctors. Regardless of class, performance is governed primarily by the corrector's actuator number, stroke, and influence function. Stroke refers to the dynamic range of the corrector and limits the largest wavefront error that can be corrected. Larger stroke provides better correction of large-magnitude aberrations. Actuator number across the eye's pupil and actuator influence function (localized deflection of the mirror surface that results when a single actuator is pushed or pulled) determine the fidelity of the correction. More actuators and a more localized influence provide better correction of high-spatial-frequency aberrations (i.e., aberration, that are of higher order). Other important corrector parameters include temporal response, surface reflectivity, diameter of mirror surface, and cost.

Wavefront correctors first applied to the eye were developed primarily for compensation of atmospheric turbulence. A common example being macroscopic discrete actuator deformable mirrors (DMs), such as those manufactured by Xinetics, Inc.⁷¹ Specifically, their actuator number, stroke, influence function, and speed were tailored to the spatial and temporal properties of the atmosphere⁴ rather than that of the eye.^{43,44,58,59} The dynamic range of such devices often limited the compensation of ocular aberrations, especially second-order contributions, and their kilohertz response was overkill for the 1 to 2 Hz fluctuations of the ocular aberrations. While wavefront correctors represent a small fraction of the total cost of ground-based telescopes in which they are employed, they represent a significant fraction of the total cost of most retina cameras. Atmospheric wavefront correctors are also generally bulky, with large mirror surfaces (~several centimeters or more) that require long focal length relay optics to magnify the pupil of the eye. A smaller corrector comparable to the dilated pupil of the eye (4 to 8 mm) can substantially reduce the instrument size, which is attractive for research retina cameras and mandatory for commercial ones.

Alternative wavefront corrector technologies, which are more cost effective and smaller, have become commercially available and span the four corrector classes depicted in Fig. 6. Large stroke (up to $16\ \mu\text{m}$) bimorph mirrors having 13 to 35 actuators have been investigated by several groups.^{15–17,28,30,72} An even larger stroke ($50\ \mu\text{m}$) corrector, a magnetic membrane mirror with 52 actuators was recently evaluated by Fernandez et al.⁷³ Microelectromechanical systems (MEMS) promises batch fabrication of low cost, compact wavefront correctors. Bulk micromachined membrane MEMS mirrors, employing 37 electrodes and $3.5\text{-}\mu\text{m}$ stroke have been successfully applied to the eye.^{74,75} Although bimorph and membrane mirrors have a large dynamic range for correction of low-order aberrations, the effective range drops rapidly with increasing order of the aberration mode.

Surface micromachined devices are another class of MEMS mirror whose mode of operation is comparable to discrete actuator DMs. Early performance results with a surface micromachined MEMS DM was reported by Doble et al.⁷⁶ Today, this mirror type is employed in numerous AO systems for the eye.^{20,21,24,28,30,32} Liquid spatial light modulators (LC-SLMs) are an alternative wavefront corrector technology. Transmissive, pixelated designs with 69 and 127 pixels were examined by Thibos et al.⁷⁷ and Vargas et al.,⁷⁸ respectively. Prieto et al.⁷⁹ and Fernandez et al.²⁹ investigated an optically addressed LC-SLM. This reflective device has high spatial resolution (480×480 piston-only pixels) and operates with low control voltages ($\sim 5\ \text{V}$). LC-SLMs in general, however, are limited to 2π phase modulation of polarized light. Extended ranges can be achieved with phase wrapping.

Theoretical studies to predict the performance of wavefront correctors and guide their use with the eye have been undertaken with representative results from two such studies given in Fig. 7. Figure 7a shows the predicted corrected Strehl for three general types of wavefront correctors (discrete actuator, piston/tip/tilt, and piston) in conjunction with the measured ocular aberrations of Fig. 1. In this study, the effect of actuator stroke was separated from that of the actuator number and influence function by assuming sufficient stroke in the Fig. 7a plot. For the discrete actuator corrector in the figure, at least a 13×13 array of actuators is required to achieve diffraction-limited imaging (Strehl = 0.8) in 95 percent of the population, that is, ± 2 standard deviations. Piston/tip/tilt correctors require somewhat less and piston correctors require substantively more. As shown in Fig. 7a, 100 to 150 piston

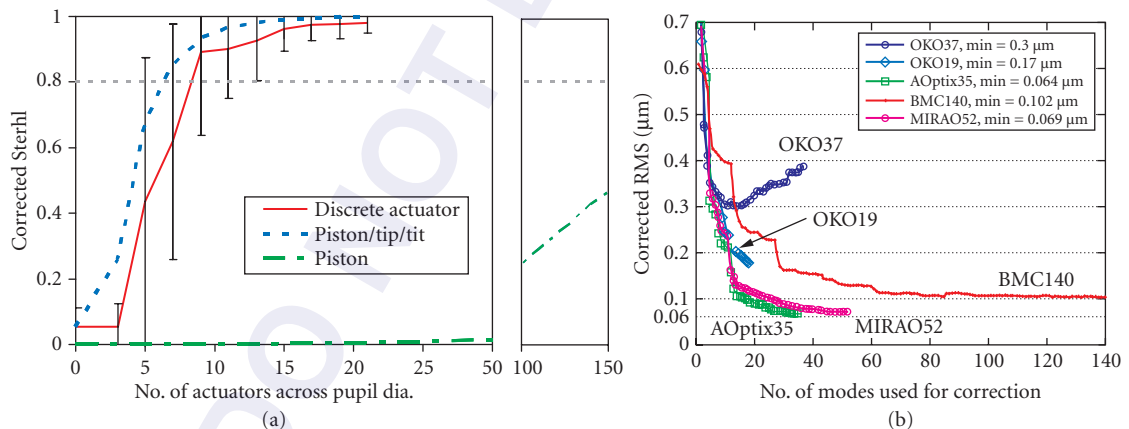


FIGURE 7 (a) Predicted corrected Strehl ratio for three general types of wavefront correctors (discrete actuator, piston/tip/tilt, and piston) as a function of actuator number across the pupil diameter.⁶⁸ Pupil diameter at the eye is 6 mm, wavelength is $0.6\ \mu\text{m}$, and the ocular aberrations are given in Fig. 1. Sufficient corrector stroke is assumed. The error bars for the single representative curve correspond to ± 2 standard deviations. (b) Predicted corrected RMS for several commercial correctors is shown as a function of number of corrector modes used for correction.⁷⁰ Model includes thresholding to account for the finite mirror stroke. Devices are from OKO Technologies (OKO19 and OKO37), Boston Micromachines Corp. (BMC140), AOptix Technologies (AOptix35), and Imagine Eyes (MIRA052). Pupil diameter at the eye is 6 mm and the randomly generated ocular aberrations are based on the population statistics of Fig. 1.

segments (maximum sampling density of the study) across the pupil diameter are required to reach a corrected Strehl of only 0.4. As for the required stroke of the corrector, it should be at least one-half the PV error of the aberrations to be corrected, for example, the PV error in Fig. 1*b*. For example, in Fig. 1, the PV error for the 7.5 mm pupil ranges from 7 to 11 μm depending on the second-order state. Thus the required stroke must be at least one half of this (3.5 to 5.5 μm). Note that the required stroke can increase significantly if a larger portion of the second-order aberrations is corrected with AO rather than for example trial lenses. As an example, a 1-D spherical error across a 7.5-mm pupil increases the required stroke by an additional 3.5 μm based on the equation: stroke = $1/4$ (refractive error in diopters) (pupil radius in millimeters)².

Figure 7*b* shows results from another study in which performance was predicted for several commercially available correctors. The modeling took into account the actual influence functions, finite stroke of each corrector (thresholding), and number of corrector modes included in the correction. The commercial devices evaluated included three membrane mirrors (OKO19, OKO37, and MIRA052), one bimorph mirror (AOptix35), and a surface micromachined MEMS deformable mirror (BMC140). Of these only two (AOptix35 and MIRA052) approach diffract limited performance.

In practice, AO retina cameras with a single wavefront corrector have not achieved sufficient correction to yield diffraction-limited imaging across large pupils (≥ 6 mm), more so in the presence of typical levels of second-order aberrations. Current correctors—while having improved noticeably in recent years—struggle to provide the necessary high spatial fidelity and high dynamic range to fully compensate both low- and higher-order aberrations of the eye. To overcome this barrier, other approaches have been explored, including a multipass configuration that traverses the same corrector more than once⁸⁰ and a cascade of two wavefront correctors, one of high spatial fidelity and the other of large stroke.^{34,81,82}

Control System

The control system performs the critical step of rapidly and repeatedly converting the raw output of the wavefront sensor into voltage commands that shape the corrector surface, that is, the deformation created by the sum response of the corrector actuators. The control loop takes into account the number, arrangement, and interaction of lenslets and actuators; the influence function and stroke of actuators; temporal response of the AO components; and a desired performance criteria, for example, minimize the RMS wavefront error. The control loop can be characterized by distinct spatial and temporal properties, which are each discussed in order further.

Spatial control starts with the layout of lenslets and actuators. For AO astronomical telescopes, the lenslet-actuator arrangement typically follows a Southwell or Fried configuration in which each lenslet is centered on one actuator or centered between four adjacent actuators. Both geometries yield one lenslet per actuator. This ratio maximizes the number of photons per lenslet, critical for the photon starved conditions under which astronomical AO systems operate, while preserving sensitivity to actuator motion. The relative coarse sampling of the lenslets, however, makes the AO sensitive to misalignment of the lenslet and actuator arrays, and the wavefront sensor cannot detect certain corrector modes, as for example waffle mode. As discussed in section “Wavefront Sensor,” the eye provides substantially more light than astronomical applications, allowing many more lenslets per actuator. Typical AO systems for the eye sample each actuator with two to ten lenslets. This oversampling greatly increases the tolerance to alignment errors and permits detection of all corrector modes.

Spatial (static) control of the wavefront corrector is normally realized through a multiplication of two matrices, one representing the wavefront sensor output (slope measurements) and the other a reconstructor matrix representing the interaction of each actuator with each lenslet. At its most basic level, the wavefront sensor output is an areal image of focus spots arranged in a quasi-regular pattern and resting on a noise floor, one spot formed behind each lenslet. The critical information in the image is the displacement of the focus spots relative to predetermined reference positions. Determining these displacements is nontrivial owing to noise in the image, and the finite width

and often irregular shape of the focus spots. Collectively, these reduce the accuracy to locate the spot centers. Centroiding based on center of mass is a common solution to this problem and is effective as it takes into account the full energy distribution of the spot. Following the mathematical notation by L. Chen,⁸³ the horizontal and vertical centroid positions for the k th lenslet in the array and a unit voltage applied to the m th actuator in the corrector can be expressed as

$$\Delta x_{k,m} = \frac{\sum_{i,j \in k} x_{i,j} I_{i,j,m}}{\sum_{i,j \in k} I_{i,j,m}} - x_{k,ref} \quad \text{and} \quad \Delta y_{k,m} = \frac{\sum_{i,j \in k} y_{i,j} I_{i,j,m}}{\sum_{i,j \in k} I_{i,j,m}} - y_{k,ref} \quad (1)$$

i and j represent coordinates for a pixelated detector, such as a CCD array. $I_{i,j,m}$ is the intensity distribution of the focus spots; $x_{k,ref}$ and $y_{k,ref}$ are reference coordinates, for example, those that minimize the wavefront error. The corresponding horizontal and vertical wavefront slopes are determined from the centroids by

$$s_{km_x} = \frac{\Delta x_{k,m}}{f} \quad \text{and} \quad s_{km_y} = \frac{\Delta y_{k,m}}{f} \quad (2)$$

where f is the focal length of the lenslets (see Fig. 5). In matrix notation, the wavefront slope output of the sensor can be expressed as

$$S_m = (s_{1m_x}, s_{1m_y}, s_{2m_x}, s_{2m_y}, \dots, s_{km_x}, s_{km_y}, \dots, s_{Km_x}, s_{Km_y})^T \quad (3)$$

where K is the total number of lenslets. Systematically applying a unit voltage to each of the M actuators in the corrector and recording the resulting displacement of the focus spots S_m , a $2K \times M$ matrix is constructed called the influence function matrix:

$$A = (S_1, \dots, S_m, \dots, S_M) \quad (4)$$

Because A is generated by systematically applying unit voltages to the actuators, A is mathematically related to the applied voltages, $V = (v_1, v_2, \dots, v_M)^T$, and the resulting wavefront slopes, $S = (s_{1_x}, s_{1_y}, \dots, s_{k_x}, s_{k_y}, \dots, s_{K_x}, s_{K_y})^T$, by

$$S = A V \quad (5)$$

Equation (5) can be solved for V by first finding an inverse for A . Since A is typically singular (e.g., a nonsquare matrix), no inverse exists and mathematical techniques such as singular value decomposition (SVD) are commonly used to find an appropriate pseudoinverse, denoted A^+ :

$$V = A^+ S \quad (6)$$

A^+ is often referred to as a reconstructor matrix as a single matrix multiplication with S reconstructs the actuator voltages necessary to correct the measured aberrations. The single matrix multiplication is fast and efficient, and for SVD represents a least squares solution.

The control model described above assumes the AO operates in a linear fashion. While this is overly simplistic for actual AO systems, the AO control loop is much more robust than suggested, converging even when nonlinearities are present owing to its closed-loop operation. Equation (6) refers to what is called the direct slope reconstruction method. Another approach that is sometimes employed is the modal reconstruction method in which wavefront slope measurements are first converted to Zernike aberration modes that in turn are converted to applied voltages. While this approach is slower and less efficient, it has the advantage of allowing explicit control of individual modes, including which ones are sent to the corrector.

Equation (6) does not capture the closed-loop performance of AO as it neglects the temporal dynamics of the ocular aberrations and the time delays associated with measuring the wavefront slopes, and computing and applying the control voltages. The resulting error can be mathematically

expressed as $\phi_{\text{resid}}(x, t) = \phi_{\text{aberr}}(x, t) - \phi_{\text{correct}}(x, t)$, where $\phi_{\text{aberr}}(x, t)$ is the ocular aberrations at time t , $\phi_{\text{correct}}(x, t)$ the applied correction based on slope measurements acquired at an earlier time $t - \tau$, and $\phi_{\text{resid}}(x, t)$ the uncorrected portion at time t . τ is the closed-loop period of the AO. Note that $\phi_{\text{correct}}(x, t)$ is directly related to the voltages, V , that are applied iteratively to the corrector and expressed as $V(t) = V(t - \tau) - G V_{\text{estimate}}(t)$. $V(t - \tau)$ are voltages already on the corrector as determined from the previous loop iteration, $V_{\text{estimate}}(t)$ are adjustments to the voltages based on Eq. (6), and $V(t)$ represent the updated voltages. G is the variable loop gain (adjusted between 0 and 1) that controls the extent to which $V_{\text{estimate}}(t - \tau)$ is applied and represents the trade-off between stability and sensitivity of the AO system, for example, higher gain increases sensitivity to temporal fluctuations in the aberrations, but at the expense of decreased system stability. Typical values of G for AO vision science instruments are between 0.1 to 0.4.

Predicting the temporal performance of an AO system is critical for optimizing system design for the eye and is accomplished by modeling the system as a cascade of transfer functions, each representing an independent component of the system. Figure 8 shows a block diagram representation of an AO control loop with input/output parameters $M(s)$, $X(s)$, and $R(s)$ that depict the time-varying aberrations of the eye, residual aberrations after correction, and wavefront corrector voltages. s is a complex frequency and defined as $i2\pi f$, where f is the temporal frequency and $i^2 = -1$. Noise in the system is assumed negligible. For analysis, the AO system can be decomposed as a linear cascade of independent transfer functions, in this case the four AO stages that most often limit temporal performance in AO vision science systems: SHWS exposure, sensor readout and computation of V , integral compensator, and finally holding mirror position for one AO loop period. Common expressions for the four transfer functions are given by

$$H_{\text{exp}}(s) = \frac{1 - e^{-sT_1}}{sT_1} \quad (7)$$

$$H_{\text{delay}}(s) = e^{-st_1} \quad (8)$$

$$H_{\text{comp}}(s) = \frac{G}{1 - e^{-sT_2}} \quad (9)$$

$$H_{\text{zoh}}(s) = \frac{1 - e^{-sT_2}}{sT_2} \quad (10)$$

where T_1 is the exposure duration of the SHWS detector, T_2 is the sampling period of the AO loop, and t_1 is the total delay to readout the sensor and compute V .

Using the notation in Fig. 8, the open-loop transfer function (single iteration), $H_{\text{OL}}(s)$, is equal to $M(s)/R(s)$ and corresponds to the multiplication of the four transfer functions, that is, $H_{\text{exp}}(s) H_{\text{delay}}(s) H_{\text{comp}}(s) H_{\text{zoh}}(s)$. While this quantifies the extent to which fluctuations in the aberrations are

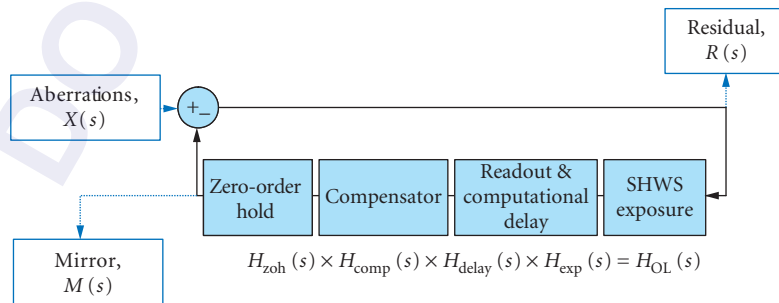


FIGURE 8 Block diagram of the AO system for modeling temporal performance.

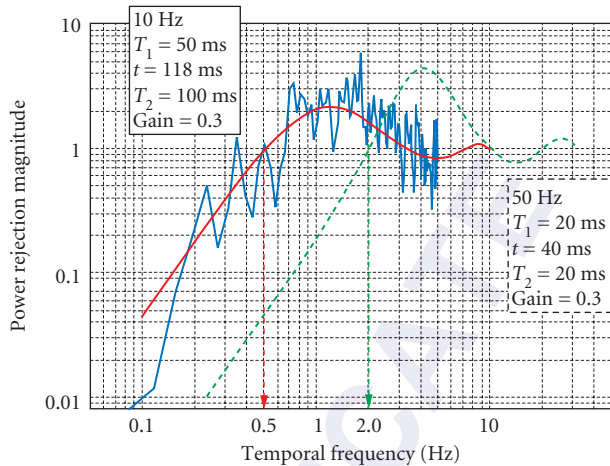


FIGURE 9 Experimental and theoretical curves for the power rejection magnitude of a representative AO system for vision science. Experimental curve (jagged blue line) was obtained on one eye using a gain of 0.3 for a 6.8-mm pupil. The corresponding theoretical curve (solid, red) was based on the actual system parameters as highlighted in the leftmost box. The second theoretical curve (dashed, green) predicts the performance of an AO system that is 5 times faster. System parameters are highlighted in the rightmost box. Cutoff frequencies for the two configurations are 0.5 and 2.0 Hz.

reliably transferred to the corrector, a more useful transfer function is the error transfer function, $H_{\text{error}}(s)$, defined as $R(s)/X(s)$. $H_{\text{error}}(s)$ quantifies the extent to which fluctuations in the aberrations are reduced by the AO. $H_{\text{error}}(s)$ is related to $H_{\text{OL}}(s)$ by the expression $1/[1 + H_{\text{OL}}(s)]$. As an example, Fig. 9 shows the power rejection magnitude (defined as $|H_{\text{error}}(s)|^2$) for two representative AO configurations, one being 5 times faster than the other. Both configurations employ double buffering, a common technique realized with a frame-transfer CCD camera, in which one SHWS frame is processed while the next frame is being acquired. The power rejection magnitude permits determining the cutoff frequency of the AO system, being defined as the frequency at which the power rejection magnitude first reaches 1. For the two configurations in Fig. 9, cutoff frequencies occur at 0.5 and 2.0 Hz, with the former consistent with the experimental measurements shown. Note that a cutoff of 0.5 Hz indicates that temporal frequencies below 0.5 Hz are (at least partially) corrected by the system, while those above 0.5 Hz are (at least partially) amplified. Because the vast majority of the aberration power in the eye lies below 1 to 2 Hz (see Fig. 2 and related discussion), the 2.0-Hz bandwidth AO configuration should be adequately fast and should provide somewhat improved performance over that of the 0.5-Hz system.

15.5 APPLICATION OF AO TO THE EYE

The sole function of AO in ophthalmic applications is to compensate or control the aberrations in the optics that lie between the eye's fundus and the rest of the world. As such, AO can be implemented into any vision or imaging application that employs these optics. This section will be split into applications where the goal is to get a better view of the retina and applications where the goal is to reduce or manipulate the blur of light delivered to the retina.

Applications Involving Imaging of the Retina

AO's role in a retina camera is to achieve diffraction-limited imaging through the highest possible numerical aperture in the eye. Most imaging applications that use AO correct over pupil sizes of 5 mm or greater, thereby taking advantage of the reduced diffraction. For example, a corrected pupil of 6 mm has a Rayleigh resolution limit about $2\ \mu\text{m}$, which is on the scale of the smallest cells in the retina (rods and foveal cones, for example).

Flood-Illuminated AO Ophthalmoscope The first ophthalmic imaging modality that demonstrated a high-order aberration correction with AO was a flood-illuminated ophthalmoscope.¹³ In this system (Fig. 10), the camera looks at the retina through an optical system where the deformable mirror is placed in the path in a plane that is conjugate to the pupil of the eye. A low coherent collimated light source (SLD, superluminescent diode) is generally used to project a beacon onto the retina. The scattered light passes through the system, off the DM and to the wavefront sensor (LA & CCD). The AO system operates in a closed loop to correct the eye's aberration. At the same time, a separate light source is used to illuminate a patch of retina limited by the field stop (FS) at the location of the beacon. The scattered light emerges, its wavefront is compensated by the DM, and the light is directed to the science camera after reflection from a dichroic beamsplitter (DBS2) to form a sharp retinal image. The camera can be moved to focus on different layers of the retina. Keeping track of the retinal and pupil conjugate positions is critical and they are labeled throughout the schematic.

The flood-illuminated AO ophthalmoscope has been used effectively to measure basic optical properties of photoreceptors. A combination of AO imaging with retinal densitometry was used to provide the first maps of the arrangement of the three cone classes in the human retina (see Fig. 11).^{84,85} It was

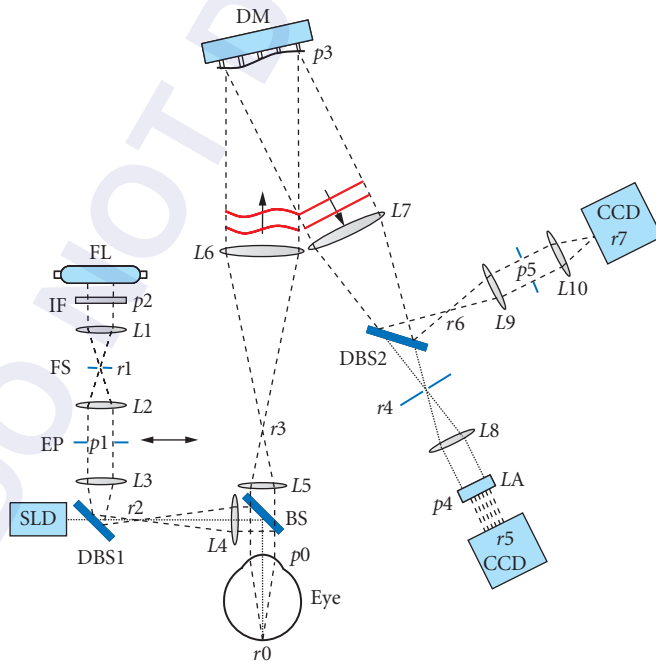


FIGURE 10 Schematic of a flood-illuminated AO retina camera.

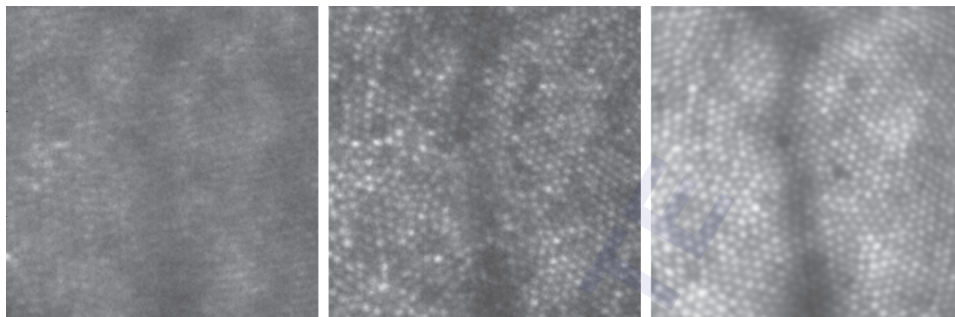


FIGURE 11 The three images are of the same location in a human retina about 1° ($300\ \mu\text{m}$) from the foveal center. The left image is taken after careful correction of defocus and astigmatism. The center image is a single frame taken with higher-order aberrations correction. The right image is a registered sum of about 20 frames. The small spots are single cone photoreceptors and in the registered frame virtually all cone photoreceptors are resolved.

also used to measure the directional properties of individual cones,⁸⁶ measure the locus of fixation relative to the maximum density of foveal cones,⁸⁷ measure changes in reflectance of cones over a diurnal period,⁸⁸ reveal new information about photoreceptor packing and structure in eyes with retinal degenerations,^{89–91} and detect fast scattering changes of cones in response to visual stimulation.^{92,93}

Adaptive Optics Scanning Laser Ophthalmoscope The scanning laser ophthalmoscope (SLO)⁹⁴ is a special application of the scanning laser microscope, invented in 1955 by Marvin Minsky.⁹⁵ The main difference with an SLO is that the eye is used as the objective lens and the retina is always the sample being imaged. As such, the lateral and axial resolution of the SLO imaging is limited by the quality of the objective lens (cornea and lens), which is why AO can provide a significant benefit.

The main motivation to employ AO in an SLO is because of its confocal property. Unlike flood-illuminated cameras, confocal systems block scattered light from features that are outside the focal plane, thereby improving the contrast of features in the focal plane of interest. Moreover, by moving the focal plane, one can optically “section” the retinal tissue (see Fig. 12 for example). These features offer significant benefits for ophthalmoscopy because the retina is a thick, multilayered structure, which scatters throughout its depth.

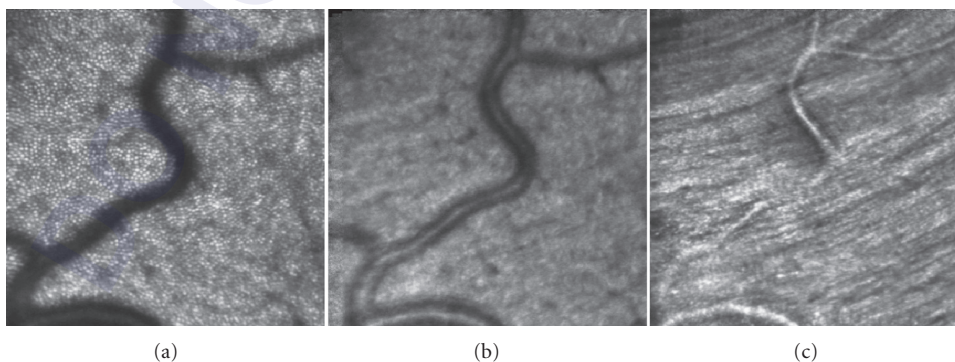


FIGURE 12 Confocal optical sectioning with an AO SLO. The sequence of AO-SLO images shows a variation in retinal structure as the focal plane of the instrument is systematically shifted from the (a) photoreceptor layer, (b) the blood vessels, and (c) the nerve fiber layer.

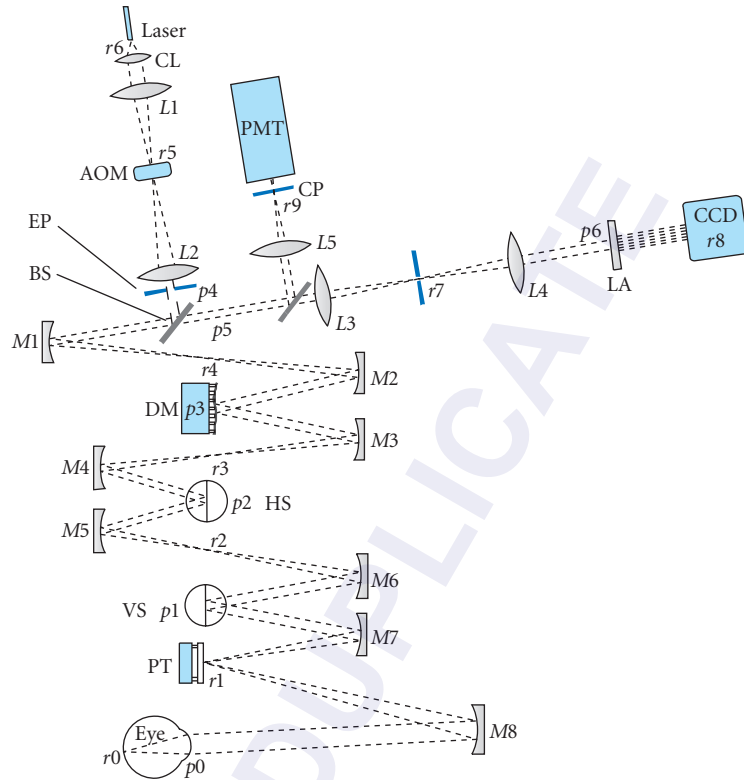


FIGURE 13 Schematic of an adaptive optics scanning laser ophthalmoscope. Key: CL—collimating lens; AOM—acousto-optic modulator; EP—entrance pupil; BS1—beamsplitter 1; DM—deformable mirror; HS—horizontal scanning mirror; VS—vertical scanning mirror; PT—pupil tracking mirror; LA—lenslet array; CP—confocal pinhole; PMT—photomultiplier tube. Pupil and retinal conjugates are labeled p and r , respectively. Mirrors and lenses are labeled $M\#$ and $L\#$ along the optical path. Telescope lens/mirror-pairs for relaying the pupil through the path are $L1-L2$, $L3-L4$, $M1-M2$, $M3-M4$, $M5-M6$, and $M7-M8$.

An SLO image is acquired over time as the scattered light is recorded from a focused spot on the retina as it scans in a raster pattern across the retina. The image is formed digitally in a frame grabber, which combines horizontal and vertical position information from the scanning mirrors with the digitized values of the analog intensity stream from the detector, to form an extended field image.

The concept of applying AO to an SLO was first proposed by Dreher et al.¹⁰ but was not fully implemented with a wavefront sensor and a mirror capable of correcting high-order aberrations until Roorda et al.¹⁹ The layout of a representative AO-SLO instrument is shown in Fig. 13.

In a confocal scanning laser imaging system, the effective PSF is computed as

$$\text{PSF}_{\text{SLO}} = \text{PSF}_{\text{in}} \times (\text{PSF}_{\text{out}} \otimes D) \quad (11)$$

where D represents the confocal aperture. When the confocal aperture approaches a size equal to the radius of the airy disk of the collection path of the system, the effective PSF is simply the product of the ingoing and outgoing PSFs. Under these optimal conditions, the lateral resolution exceeds conventional imaging by about 40 percent, as assessed by the full width at half maximum (FWHM)

of the PSF. If the confocal aperture is large, then image quality is governed only by PSF_{in} , effective optical sectioning disappears, and lateral resolution is the same as a flood-illuminated imaging system. Axial resolution in confocal SLOs can be defined and computed in several ways. The standard way to determine axial resolution is to measure the detected intensity from a flat, diffusely scattering surface as it moves relative to the focal plane of the SLO.¹⁰ The FWHM of the resulting intensity distribution is a measure of the axial resolution of the SLO.⁹⁶

The AO SLO employs a series of telescopes to relay the light to the elements that act on the beam. For example, wavefront sensing ($p6$ in Fig. 13), wavefront correcting ($p3$), and tip and tilt adjustments ($p1$ and $p2$) for raster scanning the beam all need to be done in pupil-conjugate planes. At present, each of these actions is done with a separate component and so a telescope is required to relay conjugate images of the pupil to the various components. Maintaining conjugacy of the pupil planes is critical in an AO SLO. If, for example, the eye's pupil was not conjugate to the scanning mirror, the beam would scan across the pupil, rather than pivoting about the pupil, and neither the wavefront sensor nor deformable mirror would see a stable aberration pattern. To avoid back reflections, the telescopes employ focusing mirrors rather than lenses. A consequence of using mirrors in the optical path is that they have to be used off-axis, which generates unwanted aberrations in the optical path. However, with the help of optical design software, each mirror in the system can be adjusted to compensate for the coma generated by the previous reflection, and the remaining astigmatism can be corrected by a cylindrical lens placed somewhere in the path.⁹⁷

Wavefront sensing and compensation is done in a unique way in an AO SLO. First, AO is effective in both directions; to focus the light to a sharp point on the retina and to take the scattered light from the eye and focus it to a sharp point at the confocal pinhole. In the AO SLO, the wave aberration is measured with the same light that is used to form the image. This is possible because, although the light is being scanned in a raster pattern on the retina, the light is descanned on the return path, rendering it stationary again. Thus, the wavefront sensor sees the light from the retina as though it were coming from a single spot, which makes an aberration measurement possible. This method has several advantages. First, it automatically implements the method of Hofer et al. in which scanning on the retina is employed to remove speckle from the short-exposure retinal images.⁵⁸ A second advantage is that the average aberration is measured over the entire field of view of the system, thereby ensuring more uniformity in the correction over the field that is being imaged. The final advantage is that the wavefront sensing and the imaging portions of the system use the same light path and light source, which reduces noncommon path errors and eliminates noncommon aberrations between the wavefront sensor and imaging camera due to chromatic effects.

Adaptive Optics Optical Coherence Tomography OCT is a noninvasive, interferometric imaging modality that provides significantly higher-axial resolution and sensitivity than the flood-illuminated ophthalmoscope and SLO modalities discussed previously. The term OCT was coined in 1991 in reference to the first optical B-scan (xz plane) images collected of the in vitro human retina using interferometry.⁹⁸ Note that this was preceded by substantial work in the 1980s on a technique called low-coherence interferometry that provided depth-resolved A-scan (z) images of the retina using essentially the same interferometric principles.^{99,100} Since the early 1990s, OCT technology and knowledge have grown rapidly. This has led to an increasingly large and diverse array of OCT designs that fall into two broad categories: time domain and spectral (or Fourier) domain. These domains refer to the temporal and spectral detection of the OCT signal, respectively. Details of time domain and spectral domain OCT, associated underlying theory, and other applications can be found in Chap. 18.

In an AO-OCT combination, OCT provides the high-axial resolution, and AO the complementary high-transverse resolution. Together, the AO-OCT combination provides a powerful imaging tool whose 3D resolution and sensitivity in the eye substantially surpass that of any current retinal imaging modality. Current state-of-the-art AO-OCT instruments have an isotropic 3D resolution approaching $3 \times 3 \times 3 \mu\text{m}^3$ in retinal tissue. Most OCT designs have already been combined with AO and have demonstrated both an increase in transverse resolution and sensitivity beyond that realized with OCT alone. AO-OCT combinations include time-domain en face (xy) flood-illumination OCT using an areal CCD,²⁵ time-domain tomographic scanning (xz) ultrahigh-resolution OCT,²⁶ time-domain en face scanning OCT,^{31,33} high-resolution spectral-domain OCT,^{27,28,30–32} ultrahigh-resolution spectral-domain OCT,^{29,34,35,82} and swept

source OCT.¹⁰¹ AO-OCT instruments have been evaluated with a variety of wavefront correctors that cover the four types depicted in Fig. 6: discrete actuator, membrane, bimorph, and LC-SLM wavefront correctors, all in combination with a SHWS.

As a representative example, Fig. 14 shows a schematic of the Indiana AO-OCT instrument based on ultrahigh-resolution spectral-domain OCT. Use of a single mode fiber beam splitter provides true confocality [see Eq. (11)]. The sample arm is characteristic of an AO-SLO (see Fig. 13 and corresponding text), raster scanning a focused spot across the retina. The slower acquisition of OCT, however, precludes use of a resonant scanner, which is used with the SLO. For the AO OCT shown in the figure, improved aberration correction is realized by cascading two wavefront correctors, a bimorph that corrects low-order, large-magnitude ocular aberrations and a discrete actuator deformable mirror for higher-order, small-magnitude aberrations. As an added benefit, the large stroke of the bimorph reduces the need for trial lens correction of individual refractive error, which is important for OCT because trial lenses produce unequal dispersion in the sample and reference channels. A third advantage of two wavefront correctors is that the focus of the AO correction can be varied within the retina for optimal lateral resolution at a particular retinal layer of interest.

Figure 15 illustrates the optical benefit of AO OCT compared to the widely used OCT clinical instrument, the Stratus OCT3. The inset shows essentially the same cross-sectional patch of retina imaged with both instruments. While the cellular structures are clearly too small for the clinical instrument to resolve, the same structures (in particular the lateral and axial extent of individual cone photoreceptor outer segments) are observed with AO OCT with focus at the photoreceptor layer. An additional benefit of AO—readily apparent in the inset—is the much smaller speckle size due to the larger pupil through which the image is captured, in this case more than four times larger in diameter. The larger pupil also increases substantially the light collection capacity of the instrument.

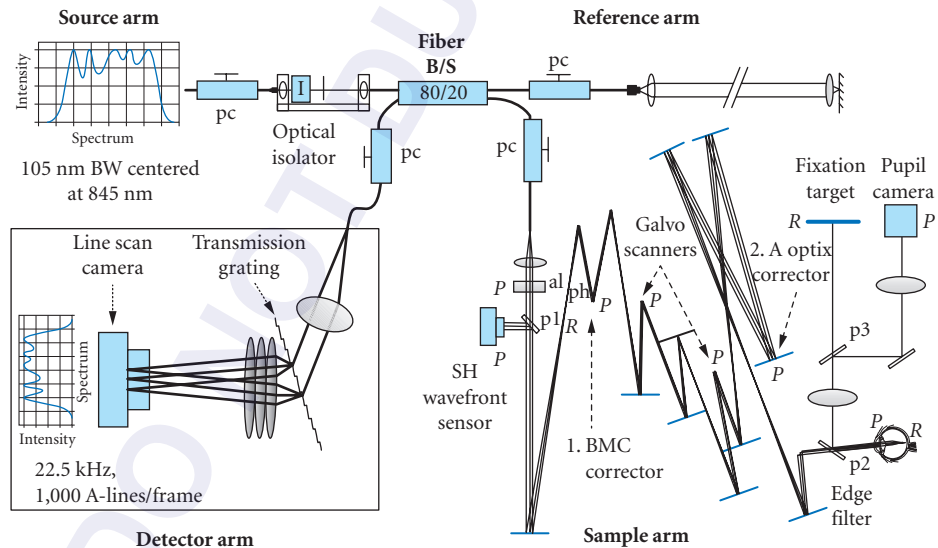


FIGURE 14 Schematic of an AO-OCT instrument based on ultrahigh-resolution spectral-domain OCT.⁸² The AO system is located in the sample arm and consists of a Shack-Hartmann wavefront sensor, a 140-element Boston Micromachines MEMS deformable mirror, and a 37-element AOptix bimorph mirror. The MEMS corrects high-order, small-magnitude ocular aberrations and complements the low-order, large-magnitude correction of the AOptix DM. The AOptix DM is positioned close to the eye to prevent refractive errors of the eye from distorting the circular beam at the SHWS.³⁰ Key: al—customized achromatizing lens; B/S—beam splitter; p1–3—pellicles; pc—polarization controller; ph—pinhole; P/R—pupil and retina conjugate planes. Footprint of the instrument is approximately 2.5×4 ft.

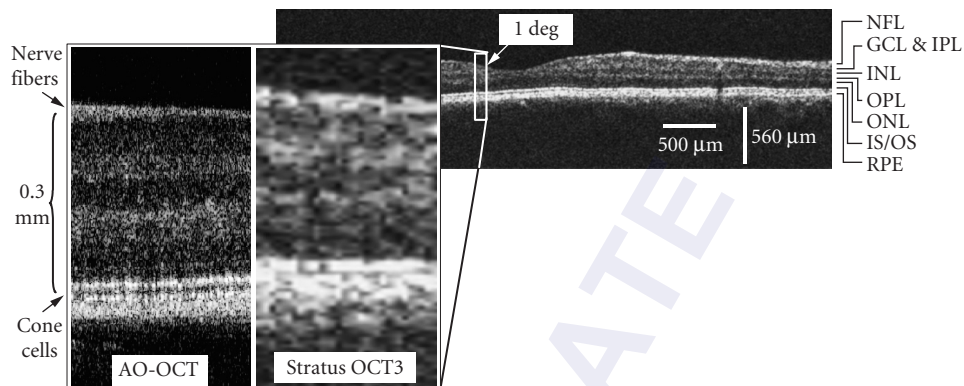


FIGURE 15 Cross-sectional images of the retina in the same subject obtained with (a) AO high-resolution spectral-domain OCT and (b) clinical OCT. Images were acquired at 1° retinal eccentricity. Key: NFL—nerve fiber layer; GCL—IPL—ganglion cell layer; IPL and OPL—inner and outer plexiform layers; INL and ONL—inner and outer nuclear layers; IS/OS—inner segment/outer segment junction; RPE—retinal pigment epithelium.

For example, a 4 times increase in pupil diameter (e.g., 6 mm compared to 1.5 mm) leads to a nominal 12-dB upper limit increase in sensitivity. The actual increase is lower depending on the scattering properties of the retina and the throughput efficiency of the sample channel, which contains many more elements than conventional OCT.

Current AO-OCT systems have been used successfully to create 3D retinal images of structures previously only visible with histology, including the foveal microvasculature, bundles within the retinal nerve fiber layer, the fibers of Henle, the 3D photoreceptor mosaic, drusen in age-related macular degeneration, and the tiny pores of the lamina cribrosa of the optical nerve in normal and glaucoma patients. AO OCT is beginning to be employed for studying a number of other clinical conditions, including age-related macular degeneration, hereditary retinal dystrophies, and optic neuropathies.

Contrast and Resolution Image contrast is a key requirement for visualizing features in retinal images. There are three elements to achieving an image with sufficient contrast. The first is to transfer as much of the available contrast from the object to the image. This is achieved by making the modulation transfer function (MTF) as high as possible and where the role of AO is crucial. Additional increases in the MTF are obtained by judicious selection of the imaging modality. AO SLO, for example, uses confocal imaging to increase the lateral and axial MTF.^{96,102} OCT uses coherence-gated detection to increase the MTF in the axial direction. The second key is to make the most of the signal that reaches the image plane. Proper selection of detectors and detection modalities can be used to achieve high signal-to-noise ratios with low levels of light. Flood-illuminated cameras use cooled scientific grade CCD cameras with low readout noise, AO SLOs use sensitive photomultiplier tubes and avalanche photodiodes with carefully designed amplification electronics,¹⁰³ and OCT uses heterodyning techniques to amplify weak signals. The third element is to have contrast in the object itself. Not only do retinal features need to scatter light, but they need to scatter light in such a way as to reveal their intrinsic structure. Cone photoreceptors, for example, have very high intrinsic contrast and lend themselves very well to imaging. The fiber optic property of cones⁸⁶ makes them appear as an array of regularly spaced point sources of nearly 100 percent contrast. But most other features in the retina scatter very little light and have less intrinsic contrast. The index of refraction of a ganglion cell body, for example, is well matched to its surrounding media and therefore scatters very little light. Under these situations, manipulation of the imaging wavelength and/or polarization might be used to improve contrast. Blood vessels, for example, are best viewed with shorter wavelengths¹⁰⁴ and the contrast of the photoreceptors and other retinal features can be enhanced by careful control of the incident and detected polarization.¹⁰⁵

A common method to achieve contrast in many imaging systems is to use fluorescence, both intrinsic and extrinsic. Recently, an AO SLO at the University of Rochester has been equipped with a fluorescence channel where they imaged the autofluorescent lipofuscin to reveal the complete and contiguous array of retinal pigment epithelium cells in human and monkeys. By using extrinsic contrast agents in a monkey, they could image retinal capillaries (fluorescein angiography), and ganglion cell bodies and their dendrites.²¹

Applications Involving AO-controlled Light Delivery to the Retina

The pioneers of AO for the human eye appreciated that the benefits of AO work in both directions.¹³ Human vision has always been limited by diffraction for small pupils and aberrations for large pupils. Consequently, the optimal pupil size for sharpest vision in a typical human eye often lies somewhere between 2 and 4 mm,^{12,106–110} depending on image quality criteria and the specific individual. Even for the optimal pupil size, the image quality remains very poor by comparison with its diffraction limited condition and is far from what could be achieved if the aberrations of the largest anatomical pupil size were corrected.

Before AO, the only method to generate high frequency, high-contrast images on the retina was using interferometry. Two mutually coherent point sources, laterally displaced in the pupil would produce sinusoidal interference fringes of 100 percent contrast at spatial frequencies governed by the separation of the two points.^{111–113} These techniques have been used effectively to understand some aspects of the limits of human vision and can produce higher-contrast images than any AO system. But, aside from a few modifications,¹¹⁴ the stimuli were limited largely to sinusoidal gratings. The contrast of images projected through the AO system is still limited by the MTF of large pupils, but can generate complex images with more contrast and resolution than ever before. Moreover, an AO system need not be used only to correct the aberration, but can be used to manipulate them as well.

Challenges for AO Systems That Project Light into the Eye All AO systems face the challenge of how to ensure that the AO is working properly. In astronomy, the performance of the AO system is gauged by the quality of the images of distant stars, which are effective point sources. In the retinal image, where no point sources exist, two metrics are typically used to gauge how well an AO system is working, the residual wavefront error, and the relative improvements of the retinal image. Neither metric provides absolute, objective measures of system performance. In postprocessing, multiframe blind deconvolution can produce reasonable estimates of the residual point spread function,¹¹⁵ but these cannot be used as feedback in closed-loop operation. When projecting light into the eye, it is even more difficult because feedback on image quality is not available in most cases, and we can only rely on residual wavefront error and the subject's perception or visual performance.

The challenge is made worse because the stimulus is viewed along a different optical path, it extends over a finite field angle, and often employs wavelengths or a span of wavelength that are different than the wavefront sensing beacon. The main challenges that one faces in developing an AO system for vision science are:

Longitudinal chromatic aberration (LCA) The eye has a lot of chromatic aberration. Refractive power from 400 to 900 nm will change by about 2.6 D^{116–118} and the foci span an axial focus range in the eye approaching 1 mm. By comparison, the thickness of the neural retina is about 250 μm . This means that if one measures and corrects aberrations with IR light and tests vision with visible light, then they must adjust for the focus of the visual target. Fortunately, the LCA of the eye is fairly consistent between individuals so one can estimate where the correct focus must be.^{116,119} But there is no element of the AO closed loop that will be able to provide feedback on the fine tuning of the visible wavelength. This forces the experimenter to use the subject's perceptual responses as feedback for alignment of the system.

Transverse chromatic aberration (TCA) Images of different wavelengths projected onto the retina will be displaced laterally because of TCA. Consider, for example, a collimated white light

beam entering the eye. Relative image shifts between both ends of the white light spectrum can be significant, depending on where the beam enters the eye. TCA arises because the chief ray of the beam does not correspond with the achromatic axis of the eye. Image displacements of ± 3 arcmin can occur with pupil shifts of ± 1 mm.^{120–121} So, whenever a broadband visual stimulus is used off the achromatic axis of the eye, the spread can significantly degrade image quality and defeat the benefits of AO.

Alignment A final challenge in achieving the best AO performance for stimuli delivered to the retina is the alignment of the stimulus itself. The stimulus channel may be subject to noncommon path errors (discussed earlier) and therefore will present new aberrations that the AO system cannot detect or correct. Since there is no feedback from this optical path, the experimenter must be obsessive about its optical design and alignment. Perceived sharpness of the image is not sufficient to judge image quality as the human visual system is limited by its retina to detection of spatial frequencies at 60 cycles/degree at best, whereas the optics can transmit frequencies over 200 cycles/degree.¹¹³ Moreover, perceived sharpness depends on context and on whether you are preadapted to a blurred or sharp stimulus.¹²²

Conventional AO Vision Systems The most basic implementation of an AO system for vision testing is one in which a subject is given a view of a target through a deformable mirror that is conjugate to the optics of the eye. The optics are no different than for a AO flood-illuminated ophthalmoscope (Fig. 10), except that there is no illumination source and the retina camera is replaced by the stimulus. The first experiments were designed simply to test the extent to which the correction of aberration improved vision.^{13,123} After correction of aberration and removal of chromatic aberration (by using monochromatic light), improvements of 20 times were reported for contrast sensitivity of a 32 cycles/degree gratings through a 6-mm pupil.¹²³ Of course, the largest improvements were realized for the largest pupils, for which diffraction is low and aberrations are high. As for visual acuity, improvements were significant, but absolute performance was limited because of the low luminances available in the display. The latter result begs the question about the practical visual benefit of correcting all aberrations. Large pupils, for which the visual benefits are greatest, are only reached under dim light conditions—the same conditions for which the retina starts to limit vision.^{124–126} In a practical task (detection of the contrast threshold for a 20/60 letter) the visual benefit for all luminances after taking typical pupil sizes into account was improved with AO, but the improvements were relatively modest.¹²⁷

AO-Controlled Stimulus Delivery in an AO SLO An alternate way to test vision with AO-corrected stimuli is to project the image directly onto the retina. By modulating the laser power during the raster scan, complex AO-controlled stimuli can be projected directly on the retina and be seen by the subject. Moreover, the same modulation patterns will appear in the retinal image.¹²⁸ This method for delivering stimuli to the retina offers some important advantages. First, the scanning beam and the stimulus can employ the same light source and are therefore not affected by TCA or LCA. Second, the recorded image can be used as a metric to determine the quality of the projected image as well as where the stimulus was presented, both axially and laterally. In some cases, two lasers might be used, one primarily for imaging (e.g., an infrared channel which is used to guide the location of the stimulus) and a second for delivery of the stimulus (e.g., small spot stimulus of a visible wavelength). In this situation, LCA and TCA do need to be taken into account so the operator knows exactly where the visible stimulus is landing.²³ Applications that employ stimulus delivery in an SLO have been used on numerous occasions,^{94,129,130} but the application of the technique with AO SLO is relatively new and has been implemented on only a few systems.

What Will Vision AO Systems Be Used for?

The role of aberrations and visual performance How much can vision improve? Experiments with AO SLO have reported visual acuities better than 20/8 (letter size for 80% correct judgments of the orientation of a tumbling E).^{128,131} Ultimately, acuity is limited by other factors in the visual system,

such as cone photoreceptor spacing, eye movements, and cortical processing limits. By bypassing the optics and simultaneously providing unprecedented views of retina structure as well as eye movements, AO SLO and other modalities provide the best tools to help us understand these limits. AO SLO studies of the role of eye movements on visual acuity are currently underway.

Measuring activity of individual cones The eye was never designed to allow stimulation of single cones. Nevertheless, recording sensations from single cones can help reveal how cones are connected to neurons further downstream in the human visual system. There are several applications where stimulation on the scale of single cones is useful. Hofer et al. for example, used an AO system to deliver brief flashes of light to individual cones, and recorded subject's perceived color sensation on each trial.¹³² A simple hypothesis, termed the "elemental sensation hypothesis" was that color percepts would fall into three categories, red, green, and blue, corresponding to the three cone types in the eye. Instead, viewers required up to eight different color names, requiring a rethinking of how color sensations are developed.¹³³ From a clinical perspective, recording the sensitivity of individual cones can provide important relationships between retinal structure and function. In an elegant set of experiments, Makous et al. recorded light sensations in brief flashes to confirm that individuals with a rare genetic mutation leading to apparent loss in a subset of photoreceptors had small, cone-sized gaps across their visual field.¹³⁴ The most recent application involves the direct and targeting stimulation of individual cones in a monkey while simultaneously recording electrical activity in downstream neurons.¹³⁵ Although neuroscientists have been recording activity of single cells in the visual system for decades, AO allows, for the first time, stimulation of single cones as the input and promises to be a useful tool for mapping connections between the retina and brain.

Using AO to generate aberrations There are many cases where the presence of aberrations may be beneficial to human vision. Considering that most people are not really handicapped by their ocular aberrations, it is sensible to consider how one could work with tolerable levels of aberration to increase visual performance in some other way. The most promising application for controlling aberrations is for the relief of presbyopia, which refers to the age-related loss in ability to deform ones natural lens to focus on near objects. Although aberrations cannot restore accommodation, they can increase the eyes depth of focus and relieve these effects. AO vision systems are currently being used to test the benefits of specific aberration profiles, which can later be implemented into contact lenses or intraocular lenses.^{136–138}

15.6 ACKNOWLEDGMENTS

The authors thank Nathan Doble of IrisAO, Inc. for assistance with the Fig. 1 data. Financial support was provided by the National Eye Institute grants 1R01 EY018339 and 5R01 EY014743 to Donald T. Miller and EY014375 to Austin Roorda. The authors are also supported in part by the National Science Foundation Science and Technology Center for Adaptive Optics, managed by the University of California at Santa Cruz under cooperative agreement No. AST-9876783.

15.7 REFERENCES

1. For reviews, see W. N. Charman, "Optics of the Eye," *Handbook of Optics*, (eds.), McGraw-Hill, New York, 2008. See also P. Artal, J. M. Bueno, A. Guirao, and P. M. Prieto, "Aberration Structure of the Eye," *Adaptive Optics for Vision Science: Principles, Practices, Design and Applications*, J. Porter, H. Queener, J. Lin, K. Thorn, and A. Awwal (eds.), Wiley, New York, 2006. See also R. R. Krueger, R. A. Applegate, and S.M. MacRae, *Wavefront Customized Visual Correction: The Quest for Super Vision II*, Slack, Inc., Thorofare, New Jersey, 2004.
2. H. von Helmholtz, *Popular Scientific Lectures*, M. Kline (ed.), Dover Publications, New York, 1962.
3. M. S. Smirnov, "Measurement of the Wave Aberration of the Human Eye," *Biophysics* 7:766–795 (1962).
4. R. K. Tyson, *Principles of Adaptive Optics*, 2nd ed., Academic, 1998.

5. H. W. Babcock, "The Possibility of Compensating Astronomical Seeing," *Pub. of the Astronomical Soc. of the Pacific* **65**:229–236 (1953).
6. L. A. Thompson, "Adaptive Optics in Astronomy," *Phys. Today* 24–31 (December, 1994).
7. F. W. Campbell and D. G. Green, "Optical and Retinal Factors Affecting Visual Resolution," *J. Physiol. (London)* **181**:576–593 (1965).
8. P. Artal and R. Navarro, "High-Resolution Imaging of the Living Human Fovea: Measurement of the Intercenter Cone Distance by Speckle Interferometry," *Opt. Lett.* **14**:1098–1100 (1989).
9. D. T. Miller, D. R. Williams, G. M. Morris, and J. Liang, "Images of the Cone Mosaic in the Living Human Eye," *Vis. Res.* **36**:1067–1079 (1996).
10. A. W. Dreher, J. F. Bille, and R. N. Weinreb, "Active Optical Depth Resolution Improvement of the Laser Tomographic Scanner," *Appl. Opt.* **24**:804–808 (1989).
11. J. Liang, B. Grimm, S. Goelz, and J. F. Bille, "Objective Measurement of the Wave Aberrations of the Human Eye Using a Shack-Hartmann Wavefront Sensor," *J. Opt. Soc. Am. A* **11**:1949–1957 (1994).
12. J. Liang and D. R. Williams, "Aberrations and Retinal Image Quality of the Normal Human Eye," *J. Opt. Soc. Am. A* **14**, 2873–2883 (1997).
13. J. Liang, D. R. Williams, and D. T. Miller, "Supernormal Vision and High Resolution Retinal Imaging through Adaptive Optics," *J. Opt. Soc. Am. A* **14**:2884–2892 (1997).
14. D. R. Williams and C. Max, Figure F.1 in "Foreword," *Adaptive Optics for Vision Science: Principles, Practices, Design and Applications*, J. Porter, H. Queener, J. Lin, K. Thorn, and A. Awwal (eds.), Wiley, New York, xviii (2006).
15. V. Larichev, P. V. Ivanov, N. G. Iroshnikov, V. I. Shmalhauzen, and L. J. Otten, "Adaptive System for Eye-Fundus Imaging," *Quantum Electron.* **32**:902–908 (2002).
16. N. Ling, Y. Zhang, X. Rao, X. Li, C. Wang, Y. Hu, and W. Jiang, "Small Table-Top Adaptive Optical Systems for Human Retinal Imaging," in *High-Resolution Wavefront Control: Methods, Devices, and Applications IV*, J. D. Gonglewski, M. A. Vorontsov, M. T. Gruneisen, S. R. Restaino, and R. K. Tyson (eds.), *Proc. SPIE* **4825**: 99–108 (2002).
17. M. Glanc, E. Gendron, F. Lacombe, D. Lafaille, J. F. Le Gargasson, and P. Lena, "Towards Wide-Field Imaging with Adaptive Optics," *Opt. Commun.* **230**:225–238 (2004).
18. J. Rha, R. S. Jonnal, K. E. Thorn, J. Qu, Y. Zhang, and D. T. Miller, "Adaptive Optics Flood-Illumination Camera for High Speed Retinal Imaging," *Opt. Express* **14**:4552–4569 (2006).
19. A. Roorda, F. Romero-Borja, W. J. Donnelly, H. Queener, T. J. Hebert, and M. C. W. Campbell, "Adaptive Optics Scanning Laser Ophthalmoscopy," *Opt. Express* **10**:405–412 (2002).
20. Y. Zhang, S. Poonja, and A. Roorda, "MEMS-Based Adaptive Optics Scanning Laser Ophthalmoscopy," *Opt. Lett.* **31**:1268–1270 (2006).
21. D. C. Gray, W. Merigan, J. I. Wolfing, B. P. Gee, J. Porter, A. Dubra, T. H. Twietmeyer et al., "In Vivo Fluorescence Imaging of Primate Retinal Ganglion Cells and Retinal Pigment Epithelial Cells," *Opt. Express* **14**:7144–7158 (2006).
22. D. X. Hammer, R. D. Ferguson, C. E. Bigelow, N. V. Iftimia, T. E. Ustun, and S. A. Burns, "Adaptive Optics Scanning Laser Ophthalmoscope for Stabilized Retinal Imaging," *Opt. Express* **14**:3354–3367 (2006).
23. K. Grieve, P. Tiruveedhula, Y. Zhang, and A. Roorda, "Multi-Wavelength Imaging with the Adaptive Optics Scanning Laser Ophthalmoscope," *Opt. Express* **14**:12230–12242 (2006).
24. S. A. Burns, R. Tumar, A. E. Elsner, R. D. Ferguson, and D. X. Hammer, "Large-Field-of-View, Modular, Stabilized, Adaptive-Optics-Based Scanning Laser Ophthalmoscope," *J. Opt. Soc. Am. A* **24**:1313–1326 (2007).
25. D. T. Miller, J. Qu, R. S. Jonnal, and K. Thorn, "Coherence Gating and Adaptive Optics in the Eye," in *Coherence Domain Optical Methods and Optical Coherence Tomography in Biomedicine VII*, V. Valery, V. Tuchin, J. A. Izatt, J. G. Fujimoto (eds.), *Proc. SPIE* **4956**:65–72 (2003).
26. B. Hermann, E. J. Fernández, A. Unterhuber, H. Sattmann, A. F. Fercher, W. Drexler, P. M. Prieto, and P. Artal, "Adaptive-Optics Ultrahigh-Resolution Optical Coherence Tomography," *Opt. Lett.* **29**:2142–2144 (2004).
27. Y. Zhang, J. Rha, R. S. Jonnal, and D. T. Miller, "Adaptive Optics Parallel Spectral Domain Optical Coherence Tomography for Imaging the Living Retina," *Opt. Express* **13**:4792–4811 (2005).
28. R. J. Zawadzki, S. Jones, S. S. Olivier, M. Zhao, B. A. Bower, J. A. Izatt, S. S. Choi, S. Laut, and J. S. Werner, "Adaptive-Optics Optical Coherence Tomography for High-Resolution and High-Speed 3D Retinal In Vivo Imaging," *Opt. Express* **13**:8532–8546 (2005).

29. E. J. Fernández, B. Považay, B. Hermann, A. Unterhuber, H. Sattmann, P. M. Prieto, R. Leitgeb, P. Ahnelt, P. Artal, and W. Drexler, "Three-Dimensional Adaptive Optics Ultrahigh-Resolution Optical Coherence Tomography Using a Liquid Crystal Spatial Light Modulator," *Vis. Res.* **45**:3432–3444 (2005).
30. Y. Zhang, B. Cense, J. Rha, R. S. Jonnal, W. Gao, R. J. Zawadzki, J. S. Werner, S. Jones, S. Olivier, and D. T. Miller, "High-Speed Volumetric Imaging of Cone Photoreceptors with Adaptive Optics Spectral Domain Optical Coherence Tomography," *Opt. Exp.* **14**:4380–4394 (2006).
31. D. Merino, C. Dainty, A. Bradu, and A. G. Podoleanu, "Adaptive Optics Enhanced Simultaneous En-Face Optical Coherence Tomography and Scanning Laser Ophthalmoscopy," *Opt. Exp.* **14**:3345–3353 (2006).
32. C. E. Bigelow, N. V. Iftimia, R. D. Ferguson, T. E. Ustun, B. Bloom, and D. X. Hammer, "Compact Multimodal Adaptive-Optics Spectral-Domain Optical Coherence Tomography Instrument for Retinal Imaging," *J. Opt. Soc. Am. A* **24**:1327–1336 (2007).
33. M. Pircher, R. J. Zawadzki, J. W. Evans, J. S. Werner, and C. K. Hitzenberger, "Simultaneous Imaging of Human Cone Mosaic with Adaptive Optics Enhanced Scanning Laser Ophthalmoscopy and High-Speed Transversal Scanning Optical Coherence Tomography," *Opt. Lett.* **33**:22–24 (2008).
34. R. J. Zawadzki, B. Cense, Y. Zhang, S. S. Choi, D. T. Miller, and J. S. Werner, "Ultrahigh-Resolution Optical Coherence Tomography with Monochromatic and Chromatic Aberration Correction," *Opt. Express* **16**:8126–8143 (2008).
35. E. J. Fernández, B. Hermann, B. Považay, A. Unterhuber, H. Sattmann, B. Hofer, P. K. Ahnelt, and W. Drexler, "Ultrahigh Resolution Optical Coherence Tomography and Pancorrection for Cellular Imaging of the Living Human Retina," *Opt. Express* **16**:11083–11094 (2008).
36. D. R. Williams, J. Liang, D. T. Miller, and A. Roorda, "Wavefront Sensing and Compensation for the Human Eye," *Adaptive Optics Engineering Handbook*, R. K. Tyson (ed.), Marcel Dekker, New York (2000).
37. A. Roorda and D. R. Williams, "Retinal Imaging Using Adaptive Optics," in *Wavefront Customized Visual Correction: The Quest for Supervision II*, R. R. Krueger, R. A. Applegate, and S. M. MacRae (eds.), Slack, Incorporated, Thorofare, 2004, pp. 43–51.
38. J. Carroll, D. C. Gray, A. Roorda, and D. R. Williams, "Recent Advances in Retinal Imaging with Adaptive Optics," *Optics & Photonics News* 36–42 (January 2005).
39. N. Doble, "High-Resolution, In Vivo Retinal Imaging Using Adaptive Optics and its Future Role in Ophthalmology," *Expert Rev. Med. Devices*, **2**:205–216 (2005).
40. A. Roorda, K. Venkateswaran, F. Romero-Borja, D. R. Williams, J. Carroll, and H. Hofer, "Adaptive Optics Ophthalmoscopy," in *Atlas of Posterior Segment Imaging*, D. Huang, P. K. Kaiser, C. Y. Lowder, and E. Traboulsi (eds.), Elsevier Science, 2005.
41. J. Porter, H. Queener, J. Lin, K. Thorn, and A. Awwal (eds.), *Adaptive Optics for Vision Science: Principles, Practices, Design and Applications*, Wiley, New York (2006).
42. D. T. Miller, "Adaptive Optics in Retinal Microscopy and Vision," *OSA Handbook of Optics Vol. III.*, M. Bass, J. M. Enoch, E. W. Van Stryland, and W. L. Wolfe (eds.), McGraw-Hill, New York, (2000).
43. J. Porter, A. Guirao, I. G. Cox, and D. R. Williams, "Monochromatic Aberrations of the Human Eye in a Large Population," *J. Opt. Soc. Am. A* **18**:1793–1803 (2001).
44. L. N. Thibos, X. Hong, A. Bradley, and X. Cheng, "Statistical Variation of Aberration Structure and Image Quality in a Normal Population of Healthy Eyes," *J. Opt. Soc. Am. A* **19**:2329–2348 (2002).
45. M. C. Roggemann and B. Welch, *Imaging through Turbulence*, M. J. Weber (ed.), CRC Press, Boca Raton, 1996.
46. L. N. Thibos, R. A. Applegate, J. T. Schwiegerling, R. Webb, and VISA Standards Taskforce Members, "Standards for Reporting the Optical Aberrations of Eyes," *J. Refract. Surg.* **18**:S652–S660 (2002).
47. L. Wang and D. D. Koch, "Ocular Higher-Order Aberrations in Individuals Screened for Refractive Surgery," *J. Cataract Refract. Surg.* **29**:1896–1903 (2003).
48. T. O. Salmon and C. van de Pol, "Normal-Eye Zernike Coefficients and Root-Mean-Square Wavefront Errors," *J. Cataract Refract. Surg.* **32**:2064–2074 (2006).
49. M. J. Collins, C. F. Wildsoet, and D. A. Atchison, "Monochromatic Aberrations and Myopia," *Vis. Res.* **35**:1157–1163 (1995).
50. X. Hong and L. Thibos, "Longitudinal Evaluation of Optical Aberrations Following Laser in Situ Keratomileusis Surgery," *J. Refract. Surg.* **16**:S647–S650 (2001).

51. X. Cheng, A. Bradley, X. Hong, and L. Thibos, "Relationship between Refractive Error and Monochromatic Aberrations of the Eye," *Optom. Vis. Sci.* **80**:43–49 (2003).
52. S. Amano, Y. Amano, S. Yamagami, T. Miyai, K. Miyata, T. Samejima, and T. Oshika, "Age-Related Changes in Corneal and Ocular Higher-Order Wavefront Aberrations," *Am. J. Ophthalmol.* **137**:988–992 (2004).
53. J. McLellan, S. Marcos, and S. A. Burns, "Age-Related Changes in Monochromatic Wave Aberrations of the Human Eye," *IOVS* **42**:1390–1395 (2001).
54. H. Cheng, J. K. Barnett, A. S. Vilupuru, J. D. Marsack, S. Kasthurirangan, R. A. Applegate, and A. Roorda, "A Population Study on Changes in Wave Aberrations with Accommodation," *J. Vis.* **4**:272–280 (2004).
55. J. Porter, G. Yoon, S. MacRae, I. Cox, and D.R. Williams, "Aberrations Induced in Customized Laser Refractive Surgery due to Shifts between Natural and Dilated Pupil Center Locations," *J. Cataract Refract. Surg.* **32**:21–32 (2006).
56. H. Radhakrishnan and W. N. Charman, "Age-Related Changes in Ocular Aberrations with Accommodation," *J. Vis.* **7**:11–21 (2007).
57. S. Pantanelli, S. MacRae, T. M. Jeong, and G. Yoon, "Characterizing the Wave Aberration in Eyes with Keratoconus or Penetrating Keratoplasty Using a High Dynamic Range Wavefront Sensor," *Ophthalmology* **114**:2013–2021 (2007).
58. H. Hofer, P. Artal, B. Singer, J. L. Aragón, and D. R. Williams, "Dynamics of the Eye's Wave Aberration," *J. Opt. Soc. Am. A* **18**:497–506 (2001).
59. L. Diaz-Santana, C. Torti, I. Munro, P. Gasson, and C. Dainty, "Benefit of Higher Closed-Loop Bandwidths in Ocular Adaptive Optics," *Opt. Express* **11**:2597–2605 (2003).
60. D. L. Fried, "Anisoplanatism in Adaptive Optics," *J. Opt. Soc. Am.* **72**:52–61 (1982).
61. P. Bedggood, M. Daaboul, R. Ashman, G. Smith, and A. Metha, "Characteristics of the Human Isoplanatic Patch and Implications for Adaptive Optics Retinal Imaging," *J. Biomed. Opt.* **13**:024008 (2008).
62. H.-L. Liou and N. A. Brennan, "Anatomically Accurate, Finite Model Eye for Optical Modeling," *J. Opt. Soc. Am. A* **14**:1684–1694 (1997).
63. R. Q. Fugate, D. L. Fried, G. A. Ameer, B. R. Boeke, S. L. Browne, P. H. Roberts, R. E. Ruane, G. A. Tyler, and L. M. Wopat, "Measurement of Atmospheric Wavefront Distortion Using Scattered Light from a Laser Guide-Star," *Nature* **353**:144–146 (1991).
64. L. Llorente, S. Marcos, C. Dorronsoro, and S. A. Burns, "Effect of Sampling on Real Ocular Aberration Measurements," *J. Opt. Soc. Am. A* **24**:2783–2796 (2007).
65. G. Yoon, "Wavefront Sensing and Diagnostic Uses," in *Adaptive Optics for Vision Science: Principles, Practices, Design and Applications*, J. Porter, H. Queener, J. Lin, K. Thorn, and A. Awwal (eds.), Wiley, New York, 2006.
66. D. R. Neal, J. Copland, and D. Neal, "Shack-Hartmann Wavefront Sensor Precision and Accuracy," in *Advanced Characterization Techniques for Optical, Semiconductor, and Data Storage Components*, A. Duparré and B. Singh (eds.), Proc. SPIE **4779**:148–160 (2002).
67. D. T. Miller, L. N. Thibos, and X. Hong, "Requirements for Segmented Correctors for Diffraction-Limited Performance in the Human Eye," *Opt. Express* **13**:275–289 (2005).
68. N. Doble, D. T. Miller, G. Yoon, and D. R. Williams, "Requirements for Discrete Actuator and Segmented Wavefront Correctors for Aberration Compensation in Two Large Populations of Human Eyes," *Appl. Opt.* **46**:4501–4514 (2007).
69. G. T. Kennedy and C. Paterson, "Correcting the Ocular Aberrations of a Healthy Adult Population Using Microelectromechanical (MEMS) Deformable Mirrors," *Opt. Commun.* **271**:278–284 (2007).
70. T. Farrell, E. Daly, E. Dalimier, and C. Dainty, "Task-Based Assessment of Deformable Mirrors," *Proc. SPIE* **6467**:64670F (2007).
71. B. R. Oppenheimer, D. L. Palmer, R. G. Dekany, A. Sivaramakrishnan, M. A. Ealey, and T. R. Price, "Investigating a Xinetics Inc. Deformable Mirror," *Proc. SPIE* **3126**:569–579 (1997).
72. D. A. Horsley, H. Park, S. P. Laut, J. S. Werner, "Characterization of a Bimorph Deformable Mirror Using Stroboscopic Phase-Shifting Interferometry," *Sens. Actuators A* **134**:221–230 (2006).
73. E. J. Fernández, L. Vabre, B. Hermann, A. Unterhuber, B. Povazay, and W. Drexler, "Adaptive Optics with a Magnetic Deformable Mirror: Applications in the Human Eye," *Opt. Express* **14**:8900–8917 (2006).
74. E. J. Fernández and P. Artal, "Membrane Deformable Mirror for Adaptive Optics: Performance Limits in Visual Optics," *Opt. Exp.* **11**:1056–1069 (2003).

75. E. J. Fernández, I. Iglesias, and P. Artal, "Closed-Loop Adaptive Optics in the Human Eye," *Opt. Lett.* **26**:746–748 (2001).
76. N. Doble, G. Yoon, L. Chen, P. Bierden, B. Singer, S. Olivier, and D. R. Williams, "The Use of a Microelectromechanical Mirror for Adaptive Optics in the Human Eye," *Opt. Lett.* **27**:1579–1581 (2002).
77. L. N. Thibos and A. Bradley, "Use of Liquid-Crystal Adaptive Optics to Alter the Refractive State of the Eye," *Optom. Vis. Sci.* **74**:581–587 (1997).
78. F. Vargas-Martín, P. M. Prieto, and P. Artal, "Correction of the Aberrations in the Human Eye with a Liquid-Crystal Spatial Light Modulator: Limits to Performance," *J. Opt. Soc. Am. A* **15**:2552–2562 (1998).
79. P. M. Prieto, E. J. Fernández, S. Manzanera, and P. Artal, "Adaptive Optics with a Programmable Phase Modulator: Applications in the Human Eye," *Opt. Exp.* **12**:4059–4071 (2004).
80. R. H. Webb, M. J. Albanese, Y. Zhou, T. Bifano, and S. A. Burns, "Stroke Amplifier for Deformable Mirrors," *Appl. Opt.* **43**:5330–5333 (2004).
81. D. C. Chen, S. M. Jones, D. A. Silva, and S. S. Olivier, "High-Resolution Adaptive Optics Scanning Laser Ophthalmoscope with Dual Deformable Mirrors," *J. Opt. Soc. Am. A* **24**:1305–1312 (2007).
82. B. Cense, E. Koperda, J. M. Brown, O. P. Kocaoglu, W. Gao, R. S. Jonnal, and D. T. Miller, "Volumetric Retinal Imaging with Ultrahigh-Resolution Spectral-Domain Optical Coherence Tomography and Adaptive Optics: Comparison of Two Broadband Light Sources" *Opt. Exp.* **17**:4095–4111 (2009).
83. L. Chen, "Control Algorithms," in *Adaptive Optics for Vision Science: Principles, Practices, Design and Applications*, J. Porter, H. Queener, J. Lin, K. Thorn, and A. Awwal (eds.), Wiley, New York, 2006.
84. A. Roorda and D. R. Williams, "The Arrangement of the Three Cone Classes in the Living Human Eye," *Nature* **397**:520–522 (1999).
85. H. Hofer, J. Carroll, J. Neitz, M. Neitz, and D. R. Williams, "Organization of the Human Trichromatic Cone Mosaic," *J. Neurosci.* **25**:9669–9679 (2005).
86. A. Roorda and D. R. Williams, "Optical Fiber Properties of Individual Human Cones," *J. Vis.* **2**:404–412 (2002).
87. N. M. Putnam, H. Hofer, N. Doble, L. Chen, J. Carroll, and D. R. Williams, "The Locus of Fixation and the Foveal Cone Mosaic," *J. Vis.* **5**:632–639 (2005).
88. A. Pallikaris, D. R. Williams, and H. Hofer, "The Reflectance of Single Cones in the Living Human Eye," *Invest Ophthalmol. Vis. Sci.* **44**:4580–4592 (2003).
89. J. I. Wolfing, M. Chung, J. Carroll, A. Roorda, and D. R. Williams, "High-Resolution Retinal Imaging of Cone-Rod Dystrophy," *Ophthalmology* **113**:1014–1019 (2006).
90. J. Carroll, M. Neitz, H. Hofer, J. Neitz, and D. R. Williams, "Functional Photoreceptor Loss Revealed with Adaptive Optics: An Alternate Cause of Color Blindness," *Proc. Natl. Acad. Sci. U. S. A.* **101**, 8461–8466 (2004).
91. S. S. Choi, N. Doble, J. L. Hardy, S. M. Jones, J. L. Keltner, S. S. Olivier, and J. S. Werner, "In Vivo Imaging of the Photoreceptor Mosaic in Retinal Dystrophies and Correlations with Visual Function," *Invest Ophthalmol. Vis. Sci.* **47**:2080–2092 (2006).
92. K. Grieve and A. Roorda, "Intrinsic Signals from Human Cone Photoreceptors," *Invest Ophthalmol. Vis. Sci.* **49**:713–719 (2008).
93. R. S. Jonnal, J. Rha, Y. Zhang, B. Cense, H. Gao, and D. T. Miller, "In Vivo Functional Imaging of Human Cone Photoreceptors," *Opt. Exp.* **15**:16141–16160 (2007).
94. R. H. Webb, G. W. Hughes, and O. Pomerantzeff, "Flying Spot TV Ophthalmoscope," *Appl. Opt.* **19**:2991–2997 (1980).
95. M. Minsky, "Memoir on Inventing the Confocal Scanning Laser Microscope," *Scanning* **10**:128–138 (1988).
96. F. Romero-Borja, K. Venkateswaran, A. Roorda, and T. J. Hebert, "Optical Slicing of Human Retinal Tissue *In Vivo* with the Adaptive Optics Scanning Laser Ophthalmoscope," *Appl. Opt.* **44**:4032–4040 (2005).
97. W. J. Donnelly, "Improving Imaging in the Confocal Scanning Laser Ophthalmoscope," M.S. dissertation, (University of Houston, Houston, TX, 2001).
98. D. Huang, E. A. Swanson, C. P. Lin, J. S. Schuman, W. G. Stinson, W. Chang, M. R. Hee, et al., "Optical Coherence Tomography," *Science* **254**:1178–1181 (1991).
99. A. F. Fercher and E. Roth, "Ophthalmic Laser Interferometry," in *Optical Instrumentation for Biomedical Laser Applications*, G. J. Mueller (ed.), *Proc. SPIE* **658**:48–51 (1986).
100. A. F. Fercher, K. Mengedoh, and W. Werner, "Eye Length Measurement by Interferometry with Partially Coherent Light," *Opt. Lett.* **13**:186–188 (1988).

101. D. T. Miller, B. Cense, Y. Zhang, W. Gao, J. Jiang, and A. Cable, "Retinal Imaging at 850 nm with Swept Source Optical Coherence Tomography and Adaptive Optics," *Invest. Ophthalmol. Vis. Sci.* **48**:E-Abstract 2769 (2007).
102. C. J. R. Sheppard and M. Gu, "The Significance of 3-D Transfer Functions in Confocal Scanning Microscopy," *J. Microsc.* **165**:377–390 (1991).
103. Y. Zhang and A. Roorda, "Photon Signal Detection and Evaluation in the Adaptive Optics Scanning Laser Ophthalmoscope," *J. Opt. Soc. Am. A Opt. Image Sci. Vis.* **24**:1276–1283 (2007).
104. J. A. Martin and A. Roorda, "Pulsatility of Parafoveal Capillary Leukocytes," *Exp. Eye Res.* **8**(3):356–360 (2008).
105. H. Song, Y. Zhao, X. Qi, Y. T. Chui, and S. A. Burns, "Stokes Vector Analysis of Adaptive Optics Images of the Retina," *Opt. Lett.* **33**:137–139 (2008).
106. F. W. Campbell and D. G. Green, "Optical and Retinal Factors Affecting Visual Resolution," *J. Physiol.* **181**:576–593 (1965).
107. F. W. Campbell and R. W. Gubisch, "Optical Quality of the Human Eye," *J. Physiol.* **186**:558–578 (1966).
108. H. C. Howland and B. Howland, "A Subjective Method for the Measurement of Monochromatic Aberrations of the Eye," *J. Opt. Soc. Am.* **67**:1508–1518 (1977).
109. G. Walsh, W. N. Charman, and H. C. Howland, "Objective Technique for the Determination of Monochromatic Aberrations of the Human Eye," *J. Opt. Soc. Am. A* **1**:987–992 (1984).
110. W. J. Donnelly III and A. Roorda, "Optimal Pupil Size in the Human Eye for Axial Resolution," *J. Opt. Soc. Am. A* **20**:2010–2015 (2003).
111. M. A. Arnulf and M. O. Dupuy, "La Transmission des Contrastes par le Système Optique de L'oeil et les Seuils de Contrastes Retinines," *C. R. Acad. Sci. (Paris)* **250**:2757–2759 (1960).
112. G. Westheimer, "Modulation Thresholds for Sinusoidal Light Distributions on the Retina," *J. Physiol.* **152**:67–74 (1960).
113. D. R. Williams, "Aliasing in Human Foveal Vision," *Vis. Res.* **25**:195–205 (1985).
114. J. Liang and G. Westheimer, "Method for Measuring Visual Resolution at the Retinal Level," *J. Opt. Soc. Am. A Opt. Image Sci. Vis.* **10**:1691–1696 (1993).
115. J. C. Christou, A. Roorda, and D. R. Williams, "Deconvolution of Adaptive Optics Retinal Images," *J. Opt. Soc. Am. A* **21**:1393–1401 (2003).
116. A. G. Bennett and R. B. Rabbetts, *Clinical Visual Optics*, 2nd ed. (Butterworths, London, 1989).
117. L. N. Thibos, M. Ye, X. Zhang, and A. Bradley, "The Chromatic Eye: A New Reduced-Eye Model of Ocular Chromatic Aberration in Humans," *Appl. Opt.* **31**:3594–3600 (1992).
118. E. J. Fernández, A. Unterhuber, B. Povazay, B. Hermann, P. Artal, and W. Drexler, "Chromatic Aberration Correction of the Human Eye for Retinal Imaging in the Near Infrared," *Opt. Exp.* **14**:6213–6225 (2006).
119. L. N. Thibos, A. Bradley, D. L. Still, X. X. Zhang, and P. A. Howarth, "Theory and Measurement of Ocular Chromatic Aberration," *Vis. Res.* **30**:33–49 (1990).
120. P. Simonet and M. C. W. Campbell, "The Optical Transverse Chromatic Aberration on the Fovea of the Human Eye," *Vis. Res.* **30**(2):187–206 (1990).
121. L. N. Thibos, "Calculation of the Influence of Lateral Chromatic Aberration on Image Quality Across the Visual Field," *J. Opt. Soc. Am. A* **4**(8):1673–1680 (1987).
122. M. A. Webster, J. S. Werner, and D. Field, "Adaptation and the Phenomenology of Perception," in *Fitting the Mind to the World: Adaptation and Aftereffects in High Level Vision*, C. Clifford and G. Rhodes (eds.), Oxford University Press, 2005, pp. 241–277.
123. G. Y. Yoon and D. R. Williams, "Visual Performance after Correcting the Monochromatic and Chromatic Aberrations of the Eye," *J. Opt. Soc. Am. A* **19**:266–275 (2002).
124. D. H. Kelly, "Spatial Frequency Selectivity in the Retina," *Vis. Res.* **15**:665–672 (1975).
125. A. M. Rohaly and G. Buchsbaum, "Global Spatiochromatic Mechanism Accounting for Luminance Variations in Contrast Sensitivity Functions," *J. Opt. Soc. Am. A* **6**:312–317 (1989).
126. A. Stockman and L. T. Sharpe, "Into the Twilight Zone: The Complexities of Mesopic Vision and Luminous Efficiency," *Ophthalmic Physiol. Opt.* **26**:225–239 (2006).
127. E. Dalimier, C. Dainty, and J. L. Barbur, "Effects of Higher-Order Aberrations on Contrast Acuity as a Function of Light Level," *J. Mod. Opt.* **55**:791–803 (2008).

128. S. Poonja, S. Patel, L. Henry, and A. Roorda, "Dynamic Visual Stimulus Presentation in an Adaptive Optics Scanning Laser Ophthalmoscope," *J. Refract. Surg.* **21**:S575–S580 (2005).
129. G. T. Timberlake, M. A. Mainster, R. H. Webb, G. W. Hughes, and C. L. Trempe, "Retinal Localization of Scotomata by Scanning Laser Ophthalmoscopy," *Invest. Ophthalmol. Vis. Sci.* **22**:91–97 (1982).
130. M. A. Mainster, G. T. Timberlake, R. H. Webb, and G. W. Hughes, "Scanning Laser Ophthalmoscopy. Clinical Applications," *Ophthalmology* **89**:852–857 (1982).
131. E. A. Rossi, P. Weiser, J. Tarrant, and A. Roorda, "Does the Correction of Higher-Order Aberrations Improve Visual Performance in Myopia?," *J. Vis.* **7**(8):14; 1–14 (2007).
132. H. Hofer, B. Singer, and D. R. Williams, "Different Sensations from Cones with the Same Pigment," *J. Vision* **5**:444–454 (2005).
133. D. H. Brainard, D. R. Williams, and H. Hofer, "Trichromatic Reconstruction from the Interleaved Cone Mosaic: Bayesian Model and the Color Appearance of Small Spots," *J. Vis.* **8**(5):15; 1–23 (2008).
134. W. Makous, J. Carroll, J. I. Wolfing, J. Lin, N. Christie, and D. R. Williams, "Retinal Microscotomas Revealed with Adaptive-Optics Microflashes," *Invest. Ophthalmol. Vis. Sci.* **47**:4160–4167 (2006).
135. L. C. Sincich, Y. Zhang, P. Tiruveedhula, J. C. Horton, and A. Roorda, "Resolving Single Cone Inputs to Visual Receptive Fields," *Nature Neurosci.* in press (2009).
136. P. A. Piers, E. J. Fernández, S. Manzanera, S. Norrby, and P. Artal, "Adaptive Optics Simulation of Intraocular Lenses with Modified Spherical Aberration," *Invest. Ophthalmol. Vis. Sci.* **45**:4601–4610 (2004).
137. P. A. Piers, S. Manzanera, P. M. Prieto, N. Gorceix, and P. Artal, "Use of Adaptive Optics to Determine the Optimal Ocular Spherical Aberration," *J. Cataract Refract. Surg.* **33**:1721–1726 (2007).
138. H. Guo, D. A. Atchison, and B. J. Birt, "Changes in Through-Focus Spatial Visual Performance with Adaptive Optics Correction of Monochromatic Aberrations," *Vis. Res.* **48**:1804–1811 (2008).

REFRACTIVE SURGERY, CORRECTION OF VISION, PRK AND LASIK

L. Diaz-Santana

*Department of Optometry and Visual Science
City University
London, United Kingdom*

Harilaos Ginis

*Institute of Vision and Optics
University of Crete, Greece*

16.1 GLOSSARY

Definitions

Acellular. Containing no cells.

Accommodation. The ability of the crystalline lens to change its optical power allowing the eye to bring into focus the images of objects at different distances.

Collagen. The main structural protein found in animal connective tissues, yielding gelatin when boiled.

Endothelium. The tissue that forms a single layer of cells lining various organs and cavities of the body, especially the blood vessels, heart, and lymphatic vessels. It is formed from the embryonic mesoderm.

Epithelium. The thin tissue forming the outer layer of a body's surface.

Fibroblast. Type of cell that synthesizes and maintains the extracellular matrix of many animal tissues. Fibroblasts provide a structural framework (stroma) for many tissues, and play a critical role in wound healing. They are the most common cells of connective tissue in animals.

Hydrophilic. Having a tendency to mix with, dissolve in, or be wetted by water.

Keratocyte. A flattened connective tissue cell lying between the fibrous tissue lamellae which constitute the cornea. The keratocyte has branching processes that communicate with those of other keratocytes.

Lamella. A thin layer, membrane, scale, or platelike tissue.

Nomogram. A graphical calculating device, a two-dimensional diagram, designed to allow the approximate graphical computation of a function.

Refraction. Measurement of the focusing characteristics of an eye or eyes.

Equations

Equation (1). This is the field of a monochromatic scalar propagating wave as

$U(r)$ modulus of the wave

r position

$\phi(r)$ wavefront

k wave number

π 3.14159

ν optical frequency

t time

Equation (2). This is the definition of the wave number.

π 3.14159

ν optical frequency

c speed of light

λ optical wavelength

Equation (3). This is the wavefront aberration as a polynomial expansion.

ω wavefront aberration

x horizontal coordinate

y vertical coordinate

a_n n th weighting value or coefficient

P_n n th polynomial from an orthogonal set

Equation (4). This relates the back focal length of a correcting lens with the vertex distance and the eye's far point.

k distance from the eye to its far point

f'_c back focal length of the correcting lens

d correcting lens' vertex distance

Equation (5). This is the power of a correcting lens placed at a finite distance from the eye.

F'_c optical power of the correcting lens

K $1/k$

k distance from the eye to its far point

d correcting lens' vertex distance

Equation (6). This is an empirical equation to describe the dependence of the ablation rate on the fluence.

A ablation rate

F fluence

F_{thr} ablation threshold

m $0.3 \mu\text{m}/\text{pulse}$

16.2 INTRODUCTION

Optics of the Human Eye

The human eye can be thought of as a photographic camera; an optical element(s) forms a real image of an object onto a photosensitive surface, the retina. The eye, like any other optical system, is formed from several different materials (tissues), each with a different refractive index and transmittance. The total optical power of the eye is about +60 diopters (D) of which +48 D are provided by its frontal part, the cornea. This tissue has refractive index of $n = 1.376$. The cornea is in contact with the air, giving it strong optical power.¹⁻³

The eye is self-contained in an almost spherical “bag” of opaque white tissue, except on the cornea, called the sclera. The cornea is approximately oval in shape with average dimensions of 12.6-mm horizontal diameter and 11.7-mm vertical diameter. It absorbs most of the UV radiation with a peak at 270 nm, but transmits radiation with wavelengths from 310 to 2500 nm approximately.¹⁻³

Immediately behind the cornea is a chamber known as the anterior chamber filled with a watery fluid: the aqueous humour. This fluid replenishes nutrients and takes away metabolic wastes from the posterior cornea, crystalline lens, and anterior vitreous. It also helps to keep intraocular pressure stable. Its refractive index is $n = 1.336$.

The iris marks the end of the anterior chamber. The volume behind it is known as the posterior chamber. The iris is the aperture stop of the eye. In a normal young eye the iris can expand or contract the pupil diameter from 2 or 3 mm to 7 or 8 mm. This reduces or increases the amount of light passing through the system and increases or decreases the depth of focus of the eye. The iris can be easily identified as it gives the eye its characteristic colour: blue, green, brown, and so on. In later life the ability of the iris to expand is severely reduced.²⁻³

Next to the iris is the crystalline lens of the eye. It has a complex layered structure and is remarkably transparent. It completely lacks a blood supply and innervations. It continues to grow in size and mass throughout life. It has a gradient refractive index which varies from about 1.406 in the inner core to approximately 1.386 at the outer core. It is a flexible structure and can change shape in order to provide a fine focusing mechanism.² By combining the cornea and the lens into a single optical system it is possible to treat the system as a thin lens in a medium. The nodal point of such a system lies at 17 mm from the retina and its posterior principal point lies on the cornea, 5.6 mm from the nodal point. A schematic eye representing this system is shown in Fig. 1.

The posterior chamber is filled with the vitreous. With a volume of almost 4 mL, it represents about the 80 percent of the total volume of the eye. It is a transparent collagen-based substance with a refractive index of 1.336. The vitreous plays an important role in several metabolic eye processes like oxygenation, depository for metabolic wastes, such as lactic acid, and as a medium for active transport of different substances throughout the eye. Immediately behind the vitreous lies the retina, the surface where the image is formed. The retina is a very complex multilayered structure, that not only captures the light, but some of the initial visual processes occur there. The central part of the retina, the fovea, contains the highest density of photoreceptors and is the region where the image formed by the optics of the eye must be of highest quality.^{2,3}

The optics of the eye are covered in greater detail in Chap. 1.

The Cornea

The first optical surface of the eye is the cornea, and hence the interface between air and tissue produces the highest refractive power, and makes the cornea the most important optical surface in the eye. The most common refractive surgery procedures achieve a change in the overall optical power of the eye by changing the shape of the cornea. For this reason we will give more attention to its cellular structure.

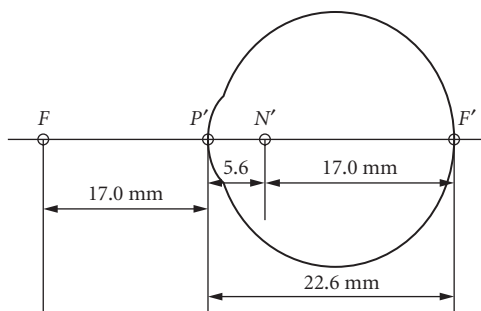


FIGURE 1 Schematic eye treating the whole optical system as a single thin lens. (Adapted from Ref. 1.)

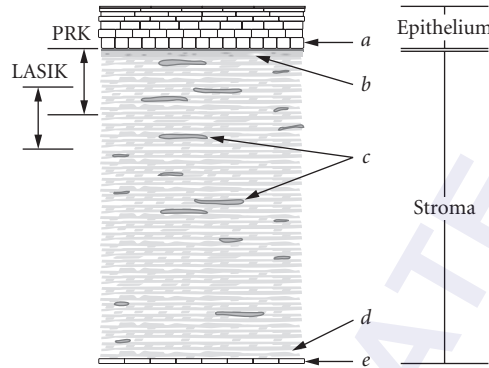


FIGURE 2 Schematic representation of the corneal layers.

The cornea is a remarkably transparent and mechanically stable tissue. The cornea features five layers with distinct roles as seen in Fig. 2. The anterior surface of the cornea is covered by the corneal epithelium (*a* in Fig. 2). The corneal epithelium is the outermost surface of the cornea and it is the surface that is wetted by the tear film.

The corneal epithelium (stratified squamous epithelium), has a thickness of approximately $50\ \mu\text{m}$. It is a multicellular layer, comprising fast regenerating cells that provide a barrier to water, ions, and bacteria. The epithelium consists of progressively flatter cells, so that its outer surface is capable of smoothing (to some extent) minor underlying irregularities or roughness. The epithelium is not a permanent tissue, in the sense that it is continuously regenerated. Consequently, in order to produce permanent changes to the corneal geometry, any modifications must be applied to the underlying layers.

Underneath the epithelium there is a layer of dense collagen (*b* in Fig. 2) having a thickness of about $12\ \mu\text{m}$, called Bowman's layer. This layer has no cells and it has been hypothesized that this structure provides mechanical strength to the cornea. The main part of the cornea representing approximately 90 percent of its thickness is the corneal stroma. The corneal stroma consists mainly of elongated bands of Type I and V collagen arranged in a lamellar configuration. Consecutive lamellae are organized in random orientations perpendicular to the optical axis to avoid anisotropy in the transverse orientations. The space around the collagen fibers is filled with a hydrophilic glycosaminoglycan and water. In the corneal stroma there are relatively sparse keratocytes (*c* in Fig. 2) that are the fibroblastic cells of the cornea. Keratocyte activity sustains proper collagen spacing and collagen fibril diameter in the cornea. In the case of injury, keratocytes produce collagen at accelerated rates in an attempt to restore stromal integrity. Often this new collagen is substantially less organized than normal collagen and is characterized by decreased transparency.

Descemet's membrane (also known as the posterior limiting lamina) (*d* in Fig. 2) of the cornea is an acellular membrane made mostly of collagen. The corneal endothelium is a monolayer of specialized cells that lines the posterior surface of the cornea. The corneal endothelium acts like a fluid pump across the posterior surface of the cornea that actively maintains the cornea in the relatively dehydrated state that maintains dense collagen organization and therefore transparency of the cornea.

Refractive Errors

The function of the optics of the eye is to form an image of the outside world onto the retina. In an *emmetropic* eye and in the absence of accommodation, the image of an object at infinity is sharply focused onto the retina. When the unaccommodated eye fails to produce a well-focused image of a distant object on the retina, it is said to be *ametropic* or to have a refractive error. The degradation in

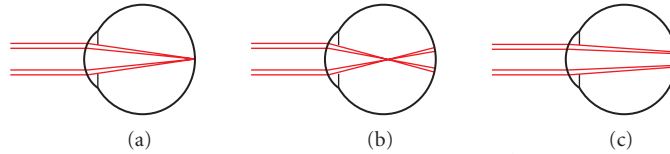


FIGURE 3 Schematic representation of refractive errors. (a) Normal emmetropic eye, (b) myopic eye, and (c) hyperopic eye.

the image seen by an ametropic eye is produced purely by the optics of the eye, and hence, in principle, it is possible to correct it by optical means alone.

The definition of refractive errors or ametropias characterize the excess or lack of optical power in the eye. In these definitions it is understood that the crystalline lens is meant to be relaxed, or *unaccommodated*. The case of an eye free of refractive errors or emmetropic is shown in Fig. 3a; the image of an object at infinity is formed on the retina. In ametropia this is no longer the case, an object at infinity appears blurred on the retina. Ametropias are divided in two main types: spherical ametropia and astigmatism.

The reader should be aware that *presbyopia* is not considered to be an ametropia. Presbyopia is the decrease in our ability to accommodate as we grow old. This process starts from youth and continues till the age of about 60 when our ability to accommodate disappears completely. In presbyopia, the crystalline lens does not change shape any longer and we are left with a single focal plane. In addition to presbyopia, an eye can have any of the refractive errors or ametropias that will be defined in the following sections.

Spherical Ametropias In spherical ametropia, the optical power of the eye is symmetric around its optical axis, but there is a mismatch between the optical power of the eye and its axial length. It is divided in two cases: *myopia* and *hypermetropia*. Myopia occurs when the optical power of the eye is too large for its length and the image of an object at infinity is formed in front of the retina. This is shown in Fig. 3b. Hypermetropia, sometimes referred to as hyperopia, occurs when the optical power of the eye is too small for its length so that the image plane lies behind the retina. This situation is shown in Fig. 3c.

The point conjugated with the retina in the absence of accommodation is called the *far point*. In an unaccommodated emmetropic eye the far point lies at infinity. Hence, accommodation is only required when seeing objects closer than infinity. In myopia, the far point is placed at a finite distance from the eye. The larger the myopia the closer the far point will be. Accommodation can only bring into focus objects closer to the eye than the far point, as accommodation can only increase the optical power of the eye. In a myopic eye any object beyond the far point will always be out of focus.

In contrast, the hypermetropic eye must always accommodate. The far point lies behind the retina, and hence power must always be added in order to bring an image into focus. The young hypermetrope can accommodate, and this is usually done unconsciously. Hence hypermetropia can pass unnoticed for many years.

Astigmatism Myopia and hyperopia affect the optical power of the eye equally in all meridians. In *astigmatism* this is not the case. An astigmatic optical surface has two perpendicular principal meridians; the curvature varying from a maximum in one meridian to a minimum in the other. The astigmatism of such a surface is defined as the difference between the optical powers at each one of these meridians. Hence astigmatism as a refractive error has the same meaning as that of on-axis astigmatism normally found in geometrical optics. Oblique astigmatism, arises when the angle of incidence of a pencil of rays is not parallel to its optical axis. This can also arise in the eye, but it is not a refractive error and hence it will not be dealt in this chapter.

The simplest possible case of an astigmatic system is that of a cylindrical lens in which one meridian has zero curvature. A lens with a toroidal surface is another example, and is commonly used to make spectacle lenses. It is generated by the revolution of an arch AB around an axis CD not passing through the arch's center of curvature C_c .

In general astigmatism can be regarded as a combination of spherical and cylindrical elements. Optically this is equivalent to an sphero-cylindrical lens with cylindrical power equal to the difference

between the two principal powers. Additionally, in the case of toroidal surfaces the spherical power is that of the weakest principal meridian.

The angle of astigmatism is measured on a plane perpendicular to the eye's optical axis. It is defined as the angle at which the meridian with the maximum curvature lies. In optometric practice this angle is measured counterclockwise starting from the horizontal, when looking into the eye. The same system is used in both eyes. For horizontal astigmatism the angle used is 180, not 0. The ° symbol is not used to avoid confusion with the number zero, so that, for example, 13° is not mistaken with 130.

Astigmatism is present in practically all human eyes. There are several reasons for this. Firstly, the cornea is seldom perfectly spherical. It usually presents certain degree of astigmatism. Likewise, the crystalline lens may present astigmatic surfaces. Moreover, the mutual alignment between crystalline lens and cornea is not perfect. There is always certain amount of tilt or decentration between these elements.

Ocular Wavefronts The eye is a biological system, and as such it is not perfectly symmetric. The refractive errors discussed so far are only a first approximation. The eye presents other more irregular errors in its refractive power. Older texts refer to these errors as *irregular astigmatism*. The more modern literature and current research papers refer to these errors as *higher-order aberrations*. Strictly speaking refractive error is a measure of the error in optical power in the eye and has units of diopters. Instead, higher-order aberrations are a measure of phase error and can have units of microns or radians. The term higher-order aberrations originates from the description of an optical wavefront.

The field of a monochromatic scalar propagating wave may be written explicitly as

$$\mathbf{u}(r, t) = \text{Re}[U(r)\exp(-ik\phi(r))\exp(-i2\pi\nu t)] \quad (1)$$

where $U(r)$ and $k\phi(r)$ are, respectively, the modulus and phase of the wave at position r , k is the wave number

$$k = 2\pi \frac{\nu}{c} = \frac{2\pi}{\lambda} \quad (2)$$

and $\text{Re}[\cdot]$ is short notation meaning the "real part of." In the definition of the wave number ν is the optical frequency, λ is the optical wavelength, and c is the speed of light in free space.

The time dependence of the optical wave of fixed wave number in Eq. (1) is known a priori. Consequently the complex function $U(r) = U(r)\exp[-ik\phi(r)]$ is enough to describe the optical disturbance. The function $\phi(r)$ is typically referred to as the wavefront of the optical perturbation and the optical perturbation $U(r)$ is referred to as the optical field, the field, or the propagating wave.

A perfect optical system, in the presence of a point source produces at the exit pupil plane a perfectly spherical wavefront converging to an ideal focusing point. Any departures from this perfectly spherical shape are known as *aberrations* and the function describing them is usually represented by $\omega(x, y)$. If the reference spherical wavefront is known, then the wavefront aberration $\omega(x, y)$ not only contains all the information about the aberrations in the system, but it can be used to describe the phase of the propagating optical wave, keeping in mind that the reference sphere has been subtracted from this function.

The quantity $\omega(x, y)$ is a real function. In Cartesian coordinates, it can be expressed as a polynomial expansion of the form

$$\omega(x, y) = \sum_{n=0}^{\infty} a_n P_n(x, y) \quad (3)$$

where P_n are the expansion polynomials and a_n are the weighting values for each polynomial. In principle $\{P_n\}$ can be any polynomial base, however, in vision sciences the Zernike polynomials have gained popularity and are widely used. The American National Standards Institute (ANSI) has produced a standard for reporting optical aberrations of eyes which uses Zernike polynomials.⁴

When using Zernike polynomials, aberrations are split into higher- and lower-order aberrations. Table 1 shows the first 14 Zernike polynomials according to the ANSI standard (having ignoring the piston or constant term, as it is not possible to measure it in the eye). The first five terms are the lower-order

TABLE 1 The First 14 Zernike Polynomials According to the ANSI Standard;⁴ ρ and θ Are Radial Coordinates on the Pupil Plane

No.	Polynomial	Description
1	$2\rho\sin\theta$	x tilt or horizontal prism
2	$2\rho\cos\theta$	y tilt or vertical prism
3	$\sqrt{6}\rho^2\sin 2\theta$	45° astigmatism
4	$\sqrt{3}(2\rho^2-1)$	Focus or sphere
5	$\sqrt{6}\rho^2\cos 2\theta$	0° astigmatism
6	$\sqrt{8}\rho^3\sin 3\theta$	Trefoil
7	$\sqrt{8}(3\rho^3-2\rho)\sin\theta$	x Coma
8	$\sqrt{8}(3\rho^3-2\rho)\cos\theta$	y Coma
9	$\sqrt{8}\rho^3\cos 3\theta$	Trefoil
10	$\sqrt{10}\rho^4\sin 4\theta$	-
11	$\sqrt{10}(4\rho^4-3\rho^2)\sin(2\theta)$	-
12	$\sqrt{5}(6\rho^4-6\rho^2+1)$	Spherical aberration
13	$\sqrt{10}(4\rho^4-3\rho^2)\cos(2\theta)$	-
14	$\sqrt{10}\rho^4\cos 4\theta$	-

aberrations and correspond to prism, sphere or defocus, and astigmatism. Aberrations described by any polynomials of order higher than two are referred to as higher-order aberrations.

Study of the wavefront error and higher-order aberrations on the eye has been possible thanks to the development of objective techniques to measure the ocular wavefront aberrations. The most popular of these techniques is the use of a Shack-Hartmann sensor, based the work by Liang et al.⁵ Another popular technique is *laser ray tracing* which was developed simultaneously by Moreno-Barriuso et al.⁶ and by Molebny.⁷ Several other techniques exist; see for instance Refs. 8 to 10.

Currently, there is no consensus about how to estimate refractive error from wavefront data. A number of different metrics have been proposed.¹¹ Also a methodology to represent wavefront maps as refractive power maps has recently been proposed by Iskander et al.¹² The details of these works fall outside the scope of this chapter and will not be addressed here.

Correcting Refractive Errors

Refractive correction must be done for the unaccommodated eye and turn the compound system formed by the eye and its correction into an emmetropic “eye”. That is, a well-corrected eye must be able to form a sharp image of an object at infinity on the retina in the absence of accommodation. We have already defined the far point of the eye as the optical conjugate of the retina in the absence of accommodation. Hence, refractive correction must shift that point to infinity. Alternatively, applying the reversibility of the optical path, a correcting lens must create an image of an object at infinity on the far point of the uncorrected eye.

Figure 4a shows the optical principle of the correction of a myopic eye. The dashed line shows the optical rays converging to the far point (F'_p) in the absence of correction. The red lines show the

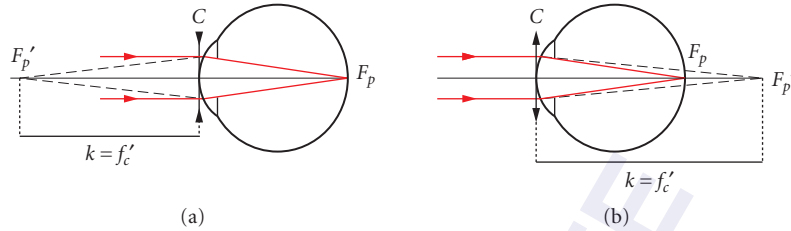


FIGURE 4 Optical principle for the correction of ametropias. (a) Myopia and (b) hypermetropia. (Adapted from Ref. 13)

light path after refractive correction. A lens of negative optical power is placed in front of the eye. Its back focal length (f'_c) is equal to the distance between the eye and its far point (k). Under these circumstances a pencil of rays coming from an object at infinity will diverge after passing through the lens, and a virtual image will form on the eye's far point F'_p . The optical system of the eye will now be able to bring this image into focus on the retina F_p .

The principle of correction for a hypermetropic eye is exactly the same as above, but in this case a correcting lens with positive power is required. This situation is shown Fig. 4b. The correcting positive lens creates the image of an object at infinity behind the eye, on its far point. As in the previous case, the unaccommodated eye can now form a sharp image of the object on the retina.

The inverse of the distance to the far point $1/k \equiv K$ is then equal to the optical power F'_p of the lens required to correct the refractive error and we have

$$k = f'_c \Rightarrow K = F'_c$$

The diagrams in Fig. 4 are both for the case when the correction lens is in contact with the eye, like the case of a contact lens. When correction is achieved with spectacle correction one must take into account the distance between the lens and the eye. This distance is always measured from the back vertex of the lens (the one closest to the eye), and is called the vertex distance and denoted by d . Diagrams illustrating this situation are shown in Fig. 5. It is clear that in this case

$$k = f'_c - d \quad (4)$$

and the power of the correcting lens must be adjusted as follows:

$$F'_c = \frac{K}{1 + dK} \quad (5)$$

It should be clear from Fig. 5 that if a refractive correction is measured with a lens at a distance d , changing this distance to d' will have an effect on the power of the correcting lens. The larger the power of

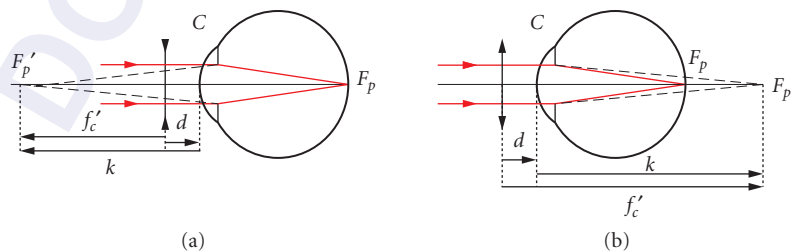


FIGURE 5 Optical principle for the correction of ametropias with a spectacle lens. (a) Myopia and (b) hypermetropia. (Adapted from Ref. 13.)

the lens the larger this impact. However, since the far point of the eye must coincide with the focal point of the correcting lens, it is straightforward to adjust the prescription. If the lens is to be moved by $d - d'$ mm closer to the eye's far point, then the focal length f_p' of the correcting lens must be reduced by $d - d'$ mm.

This is particularly important for refractive surgery, because refractive errors are typically prescribed using trial lenses placed 10 to 14 mm away from the cornea. The most common type of refractive surgery is aimed at changing the optical power of the cornea itself, so this distance must be taken into consideration.

An astigmatic eye can be treated as if it had two separate far points, each one corresponding to one of the principal meridians. The correction of astigmatism is achieved treating each principal meridian independently using the same principles described above. Namely, using an astigmatic lens, the focal point of each meridian of the lens must coincide with the far point of the corresponding meridian in the eye.

Correction of higher-order aberrations has been proposed by several authors.^{14–18} Partial correction of aberrations has also been proposed as a way to improve visual performance.¹⁹ Although these studies have shown an improvement in contrast sensitivity after correction of higher-order aberrations, this has only been possible in laboratory conditions. Moreover, Dalimier et al. found that aberration correction is most effective when the target was viewed in bright light. For lower light levels, the improvement in optical quality was masked out by a decrease in neural sensitivity. That is, for lower light levels, resolution is limited by neural mechanisms and not by optical quality.

These results mean that correction of higher-order aberrations, as an everyday solution, may not be practical for the majority of people. During the day, the pupil constricts reducing the benefit that correction could bring. At night the pupil dilates, but neural sensitivity is reduced. Higher-order aberration correction may only benefit patients with pathologically abnormal corneas where they play a dominant role in image degradation. Nevertheless, understanding of important visual mechanisms like the triggering of accommodation may largely benefit from carefully designed studies where aberration-free images are projected on the retina.

Refractive corrections are covered in greater detail in Chap. 12.

16.3 REFRACTIVE SURGERY MODALITIES

As mentioned earlier, refractive surgical procedures are attempting to alter the optical power of the eye with permanent means in order to compensate for the eye's refractive error. There are three major categories of refractive procedures: *Intraocular lens (IOL) implantation*, *corneal incisions/implants*, and *laser corneal ablation*.

Intraocular Lenses, Corneal Incisions, and Implants

Intraocular Lenses Intraocular lenses (IOLs) were first developed for the substitution of the crystalline lens after cataract extraction. The first lenses were made of appropriate glass and were substantially heavy.^{20,21} Modern lenses are made of appropriate biocompatible polymers (such as silicone or acrylic). Some types are designed to be implanted in the eye while the natural lens is present and therefore to compensate for possible refractive errors. These IOLs (phakic IOLs^{22–24}) are either implanted in the posterior chamber (between the iris and the crystalline lens) or in the anterior chamber supported by the iris. The latter is the most common type currently used. The main advantage of IOL implantation is that the correction is reversible in the sense that the implant can be removed. Generally high corrections can be achieved.

Corneal Incisions/Implants A different approach for correcting ametropias and, historically, one of the first to be widely applied corneal incisions procedures. Corneal incisions are performed (typically using a diamond knife) either in a radial or in a circumferential manner. Depth of the incisions is typically adjusted to be 90 percent of the local thickness of the cornea. The rationale of corneal incisions is that local relaxation of mechanical stress results to a change in corneal geometry. Radial keratotomy is one of the most common incisional techniques (Fig. 6).

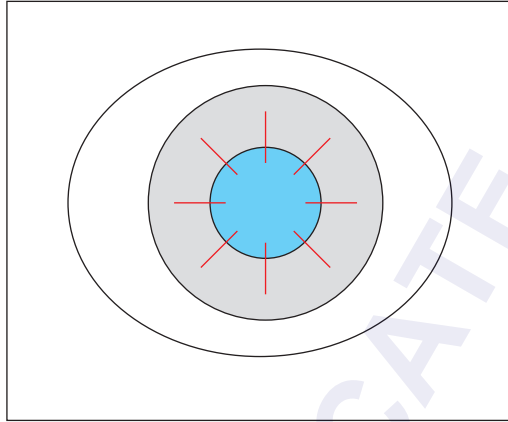


FIGURE 6 Schematic representation of radial incisions for the correction of myopia (radial keratotomy—RK).

The main drawbacks of radial keratotomy are the low predictability,²⁵ the large amount of induced aberrations,²⁶ and the compromised mechanical strength of the eye globe.²⁷ The low predictability of incisional surgery is associated with the fact that the long-term redistribution of stress in the corneal surface is not stable over time but is varying during the healing process. The final outcome must be examined several months after the incisions when healing has reached a stable condition.

The variability of the healing process may be attributed to individual factors such as the particular anatomy of the given cornea, age, and other biochemical and biological factors. In order to reduce this variability the desired refractive change is achieved by using empirical nomograms with different number, length, and positioning of the incisions that additionally take into consideration the factors mentioned above. Circumferential (arcuate) incisions have been employed for the correction of astigmatism usually in pairs across the steep meridian of the corneal astigmatism.

Another modality that has been employed to correct low refractive errors is the implantation of intracorneal ring segments.²⁸ These semicircular arcs are implanted in circumferential tunnels (created either mechanically or by means of femtosecond laser photodisruption) in the corneal stroma (Fig. 7).

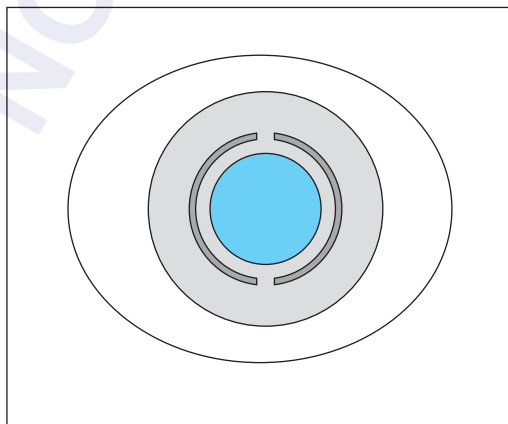


FIGURE 7 Schematic representation of intracorneal ring segments (ICRS) implanted around the pupil.

These segments are made of stiff material like polymethyl methacrylate (PMMA) and their rigid shape provides a new set of boundary conditions for corneal geometry. Their refractive effect is achieved by elevating the circumferential portion of the cornea over their implantation tunnel. The fact that they are rigid has been utilized in order to provide mechanical support to corneas with compromised mechanical stability such as in keratoconus²⁹ or keratectasia.³⁰

Laser Corneal Refractive Procedures

The majority of the refractive procedures performed today are based on controlled ablation of corneal tissue in order to modify its central curvature. Among all the different approaches, the most critical differentiation is the depth at which the ablation is performed. In surface ablations (PRK, Epi-LASIK, LASEK) tissue removal involves the Bowman's layer and superficial stroma after appropriate removal of the corneal epithelium. In LASIK, tissue is removed from deeper stroma after creation of a hinged flap. Surface ablations are considered less invasive than LASIK in the sense that for any given correction a thinner part of the stroma is affected.

Photorefractive Keratectomy (PRK) PRK was the first laser procedure to be widely available.^{31,32} It involves removal of corneal epithelium with a blunt instrument such as a spatula or a brush. After the Bowman's layer is exposed, laser ablation is performed using an appropriate superposition of a number of laser pulses. A schematic representation of this procedure is shown in Fig. 8. After

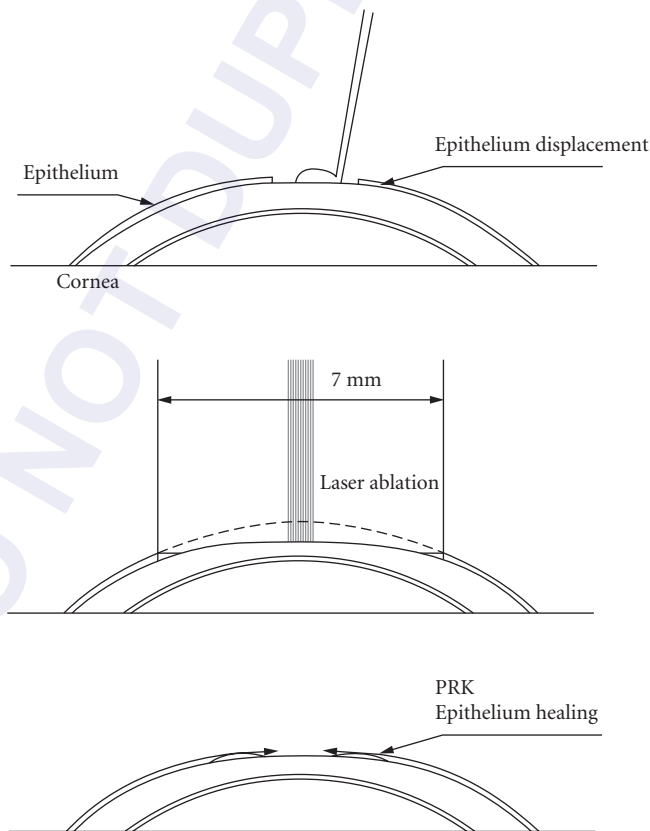


FIGURE 8 Schematic representation of photorefractive keratectomy.

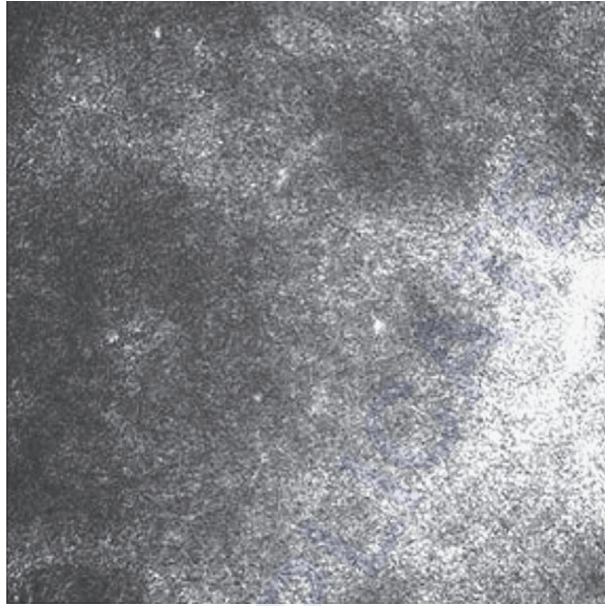


FIGURE 9 Confocal microscope image of the subepithelial scar tissue after PRK. Side is 400 μm , depth is at the epithelial—stromal interface (50 μm).

ablation, the corneal epithelium regenerates slowly and within 3 to 5 days completely covers the reshaped corneal surface. During this period of epithelial healing (that is characterized by pain or discomfort to the patient) the cornea is covered by a contact lens to reduce discomfort. Antibiotics are also administered as drops during this period of time to prevent infection of the exposed cornea. Additionally to antibiotics, anti-inflammatory treatment often in the form of steroids is prescribed. Time duration of the steroid treatment is not standard. Several months may be necessary depending on the physician's preferred scheme.

The technique is surgically simple and straightforward and is considered safe and effective for the treatment of low to moderate refractive errors.³³ Since treatment of the cornea with the excimer laser results to trauma and a transient inflammation, there is a healing response that not only modulates the refractive outcome, but additionally influences the transparency of the cornea due to the low organization of newly formed collagen. Figure 9 shows a confocal microscope image of corneal scar. This postoperative deterioration of corneal transparency (known as corneal haze) is not expressed with significant severity in the majority of the patients.³⁴

In order to avoid scar tissue formation several different approaches have been developed, collectively called advanced surface ablations. These techniques are variants of the PRK where either special pharmaceutical treatment is applied in order to control the production of new collagen or the epithelial layer is preserved in order to cover the treated area immediately after laser ablation. The most common pharmaceutical is mitomycin C.³⁵ Mitomycin C is an anticancer drug known also to effectively inhibit keratocyte proliferation and therefore haze formation. However, there is skepticism about its use by many practitioners because its long-term effects and safety have not yet been documented. In order to reduce postoperative exposure to inflammation factors the LASEK procedure was developed.

LASEK/Epi-LASIK In these procedures, the epithelium is preserved as a whole sheet and is repositioned in order to cover the treated portion of the cornea. Presumably, this would reduce

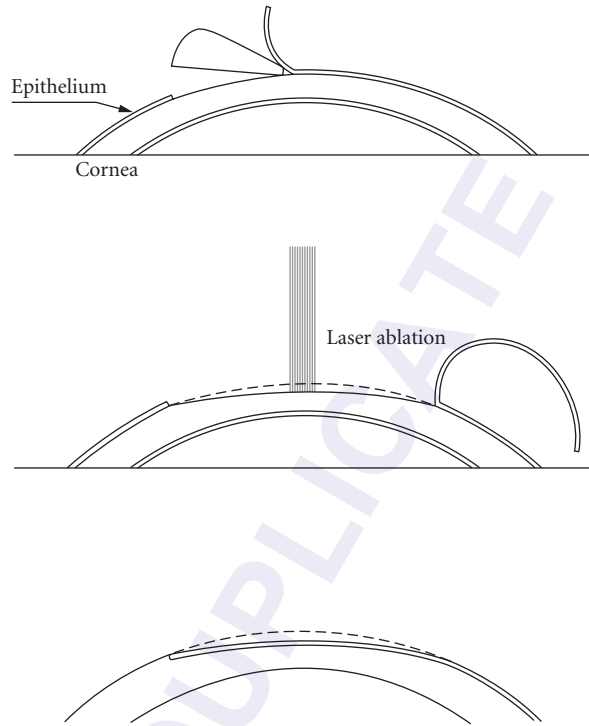


FIGURE 10 Schematic representation of the Epi-LASIK technique.

exposure of keratocytes to biochemical factors in the tears that trigger and enhance the healing response. In LASEK (laser subepithelial keratomileusis),³⁶ the epithelium is loosened by the application of an alcohol solution for a short-time interval (typically 20% alcohol in water for 30 s). After the treatment with alcohol, the epithelium can be delaminated using a blunt spatula in a manner that preserves its continuity. After laser ablation, the epithelium is repositioned onto the cornea. This epithelium is replaced in 3 to 5 days by new epithelial cells migrating from the periphery.

Epi-LASIK (epithelial laser in situ keratomileusis), is a similar technique (Fig. 10) where the epithelium is mechanically separated from the underlying Bowman's layer by means of a special device called epikeratome.^{37,38}

Although promising, these techniques have not been proved to be more effective than PRK in terms of postoperative risk for haze development. Both LASEK and Epi-LASIK being essentially surface ablations are mostly preferred for low to moderate refractive errors.

Laser In Situ Keratomileusis (LASIK) The most common laser technique for the correction of moderate and high refractive errors is laser in situ keratomileusis (LASIK). LASIK was invented in the late 1980s by Pallikaris et al.^{39,40} and has gained popularity due to its faster rehabilitation and less postoperative pain. In LASIK (Fig. 11) a thin (typically 120 μm) flap is created exposing deeper stromal layers while preserving epithelium and the Bowman's layer.

The flap is created using a special instrument called "microkeratome." The microkeratomes feature an oscillating blade and an applanator that flattens a portion of the cornea where the cut is performed. By this instrument, a flap of constant thickness can be produced where this thickness is controlled by the gap between the applanator and the knife's edge.

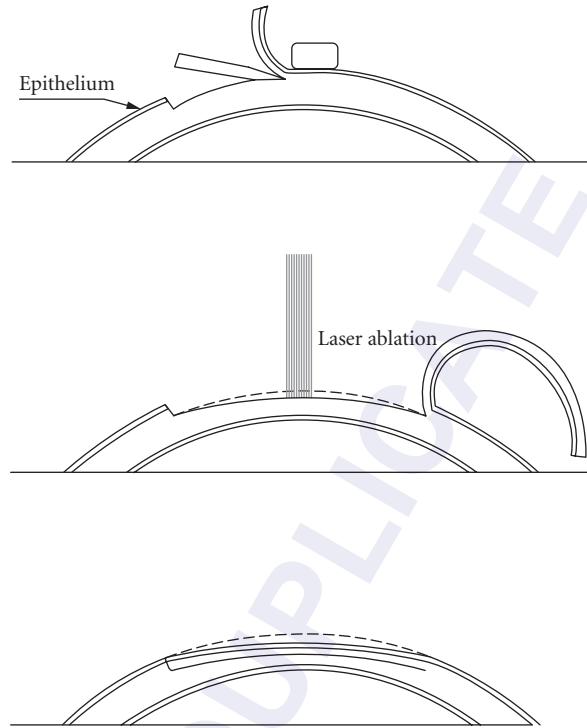


FIGURE 11 Schematic representation of laser in situ keratomileusis (LASIK).

After laser ablation, the flap is repositioned on the reshaped corneal stromal bed. This, effectively isolates the treated area from inflammation factors of the environment and therefore corneal transparency is less affected by the healing response. The creation of the flap along with the ablation may influence the mechanical stability of the cornea. There are different safety criteria related to how much ablation can be performed where the most common is that the remaining corneal stromal thickness under the ablation should be no thinner than $300\ \mu\text{m}$.

Ablation Profiles After the corneal stroma is exposed with either of the above-mentioned techniques, a sequence of overlapping laser pulses is delivered in order to excise corneal tissue in a controlled manner and modify the curvature of the corneal surface. The desired ablation profile is achieved by the appropriate superposition of—typically—a few tens of thousands of pulses, depending on the diameter of the treated zone, the desired refractive correction, and the type of ametropia.

Typically, only the portion of the cornea in front of the pupil is treated in order to minimize the tissue that needs to be excised. For this purpose, the diameter of the pupil in scotopic conditions is measured and the ablation profile is calculated for a zone preferably slightly larger than the scotopic pupil diameter.

For the correction of myopia, the refractive power of the cornea has to be reduced and therefore the ablation profile is designed to excise more tissue centrally than peripherally and therefore flatten the central cornea. Similarly, the correction of hyperopia requires steepening of the central cornea by the removal of peripheral tissue. The removed tissue has the shape of a positive or negative meniscus

lens respectively. In the case of astigmatism, the ablation profile has a toric shape introducing a relative change in dioptric power across two perpendicular meridians.

The ablation profiles typically attempt the modification of corneal curvature in a circular zone centred on the line-of-sight that has diameter of the order of 6 to 8 mm. This zone is known as the optical zone of the profile. The edges of the optical zone are treated with additional pulses in order to smooth out the border between corrected and noncorrected portion of the cornea creating a transition zone. The width of the transition zone is of the order of 1 mm and depends on the magnitude of the discontinuity that needs to be smoothed. It is clear that the wider the optical zone the more the tissue needs to be excised for any given dioptric correction. The exact shape that needs to be removed in any case is calculated by geometrical considerations known as the Munnerlyn's equations.^{41,42}

Laser systems usually incorporate further empirical adjustments to compensate for biological factors such as the biomechanical response of the cornea and/or the effect of healing (especially at the border of the optical zone). Ablation profiles can be either calculated directly from the spherocylindrical refractive error or can be generated based on measurements of corneal topography or wavefront aberration of the eye. The advent of topography and wavefront-guided ablation profiles generated an interest in creating supernormal (aberration-free) optical systems for the human eye.⁴³ However, it is now understood that the outcome variability that originates from individual factors such as the healing response does not allow the precise sculpturing of corneal surface in order to eliminate high-order aberrations in a predictable manner.³⁵

In the next section we turn our attention to the characteristics of the lasers used to perform the refractive surgery procedures and how they interact with the living tissue.

16.4 LASER ABLATION

Excimer lasers are pulsed ultraviolet (UV) lasers where the amplification medium is a biatomic system existing only at an electronically excited state.⁴⁴ Excited dimers (e.g., ArF) can be formed during electric discharges through gas mixtures of the corresponding elements and they do not exist in their ground electronic state (Fig. 12).^{44,45} Excited dimers have a typical lifetime of a few nanoseconds. Their spontaneous dissociation is accompanied by the emission of a UV photon. In excimer lasers, typically excitation is provided by a transverse electric discharge (perpendicular to the optical axis) through a mixture of gases of pressure similar to that of atmospheric pressure. The dissociation of

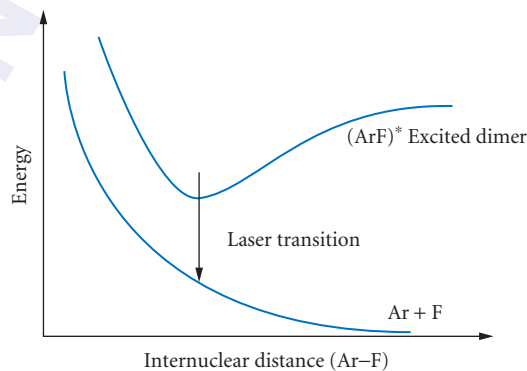


FIGURE 12 Potential energy for ArF in electronically excited and ground state.

the excited dimers renders the medium transparent to its emission wavelength. Typically, excimer lasers are characterized by high gain and low feedback therefore the beam has low coherence and generally has considerable inherent divergence.

The most common excimer lasers and their respective emission wavelengths are XeF (351 nm), XeCl (308 nm), KrF (248 nm), and ArF (193 nm). Typically, for surgical systems the pulse duration is of the order of 20-ns FWHM and the total energy per pulse is of the order of tens or hundreds of millijoule. The electric efficiency for a typical ArF laser system is of the order of 2 percent.⁴⁵ Applications of excimer lasers include material's processing, micromachining, and the excision of tissue.

Corneal Photoablation

The first experiments for the evaluation of excimer lasers as means to excise corneal tissue, were performed in the late 1980s.^{46–48} These studies concluded that the most suitable wavelength was 193 nm (ArF excimer laser) due to its relatively limited thermal damage to the underlying tissue, the accuracy of the excision, and the smoothness of the resulting surface.⁴⁸ The interaction of high-power pulsed UV radiation with corneal tissue results from the process of ablative photodecomposition (or simply photoablation) described just below.

The UV pulse is absorbed by a thin surface layer. Precise measurement of the absorption coefficient in real conditions has not been accomplished; however, the penetration depth is estimated to be a few microns.⁴⁹ As energetic photons (6.4 eV) are absorbed, the molecules that form the extracellular matrix are decomposed (Fig. 13a). This process is primarily photochemical.⁴⁸

According to a model proposed by Dougherty et al. to describe the role of corneal hydration in the photoablation process,⁵⁰ radiation with a wavelength of 193 nm is primarily absorbed by collagen, while water plays a more passive role absorbing a small fraction of the radiation and—mostly—affecting the thermal and mechanical properties of the tissue. As radiation is absorbed and molecular fragments are formed, more photons are absorbed by the fragments and this results to a portion of tissue from the surface to be effectively vaporized. For typical conditions (pulse energy/duration) the pressure near the surface may reach levels of the order of 100 bar.⁵¹ The ablation products expand rapidly and are ejected perpendicularly to the surface with supersonic velocities.⁵²

Based on the description above, it is reasonable to expect that the deeper a collagen fiber the larger the fragments and the less the kinetic energy of each fragment. Therefore, there is a critical depth (depending on the pulse energy density) for which fragments are energetic enough to be expelled from the material. Under this depth, the collagen has undergone direct photochemical damage but not to the extent required for ablation. For energy densities below a certain threshold photoablation does not occur. This threshold has been experimentally determined to be approximately 40 mJ/cm² for normal cornea and typical pulse durations.⁵³

Ablation Rate

At any given conditions (pulse duration, pulse repetition rate) the ablation rate is defined as the depth of tissue excised by each pulse and is generally expressed in $\mu\text{m}/\text{pulse}$. As mentioned above, ablation rate depends on the incident energy density. Given that for a particular laser system the pulse duration is constant, the quantity that is empirically used to quantify the incident laser energy density is the surface density (fluence) typically expressed in mJ/cm^2 .

The empirical equation to describe the dependence of the ablation rate on fluence is logarithmic⁵³

$$A = m \cdot \ln \frac{F}{F_{\text{thr}}} \quad (6)$$

where A is the ablation rate in $\mu\text{m}/\text{pulse}$, F is the fluence, F_{thr} the ablation threshold, (in this equation it is equal to 50 mJ/cm^2) and m is an empirical constant equal to 0.3 $\mu\text{m}/\text{pulse}$

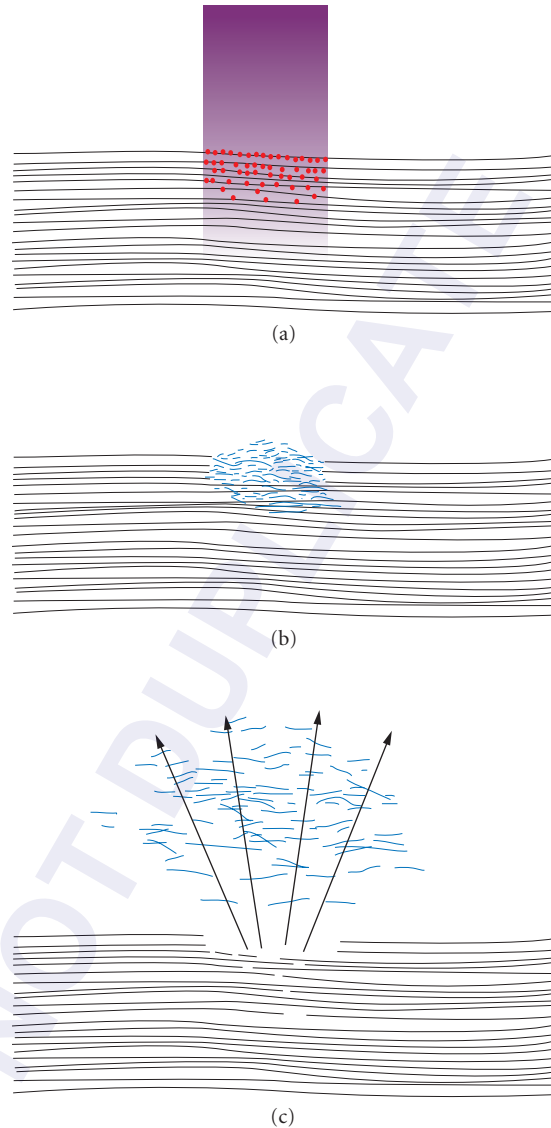


FIGURE 13 Schematic representation of the corneal photoablation process.

(Fig. 14). Typically, fluences employed in refractive surgery applications range⁵⁴ from 120 to 225 mJ/cm².

In practice, parameters such as pulse duration, beam profile, the presence of a gas stream for ablation debris aspiration, and the repetition rate may influence the ablation rate. These parameters as well as the temperature and humidity in the operating room are factors that are standardized in refractive surgery in order to minimize the variability of the delivered ablation. As mentioned earlier, the appropriate ablation profiles (that are designed to modify corneal geometry in a predetermined manner) are created with an appropriate superposition of pulses. Typical repetition rates are of the

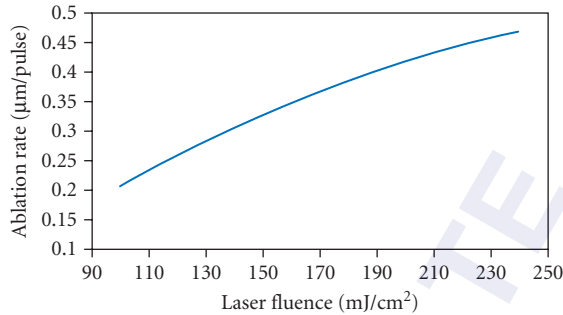


FIGURE 14 Corneal ablation rate as a function of fluence.

order of 400 Hz. For the removal of tissue necessary for the correction of moderate myopia, the total number of pulses depend on all of the parameters mentioned above and are of the order of a few tens of thousands of pulses. Besides standardization of laser and environmental parameters during ablation, the accuracy of tissue excision is enhanced by the utilization of empirical nomograms⁵⁵ where other parameters are taken into account such as attempted refractive correction, diameter of the treatment zone, patient age, corneal radius of curvature, corneal thickness, and others.

The main variable that may influence the effective ablation rate is corneal hydration.⁵⁰ Standardization of the procedures minimizes the variability associated to corneal hydration changes during surgery. It has been suggested that hydration could be measured spectroscopically and ablation parameters could be adjusted dynamically to compensate for possible changes of hydration.⁵⁶

Thermal, Photochemical, and Photoacoustic Effects

As mentioned above, photoablation can be primarily seen as a fast photochemical decomposition of the structural elements of the cornea (the extracellular matrix) during which high-pressure gradients are developed removing fragments from the surface of the irradiated tissue.^{46–48} In parallel to photoablation there are phenomena that their potential implications in the acute postoperative inflammation and long-term healing response of the cornea is not fully understood.

A significant part of the laser energy is used for the photochemical decomposition of the extracellular matrix. However, there is a part of the laser energy directly dissipated as heat (e.g., for depths greater than the ablation rate where energy density is subthreshold). Moreover, the highly energetic fragments deposit part of their energy to the surface as they are expanding. As a result of these two mechanisms each laser pulse results to a certain amount of heat dissipated on the free surface of the ablated cornea.

The temperature rise associated with this heat dissipation is highly variable and depends on the particular conditions (repetition rate, fluence, beam diameter pulse superposition pattern, etc.). In practical conditions the temperature rise has been reported to range from 0 to 8°C.^{57,58} One of the parameters taken into account in pulse superposition pattern design is the effective reduction of repetition rate in the sense that irradiation of any point with consecutive pulses is avoided. In this manner, minimal changes in surface temperature of the cornea can be achieved.⁵⁹

In principle, ArF excimer laser radiation is not capable of producing damage to the DNA of keratocytes as the penetration depth is smaller than the dimensions of one cell.⁶⁰ In that sense if the nucleus of a cell is exposed to this UV radiation a significant portion of the cell is already ablated and therefore the cell is not viable. It is the secondary radiation (fluorescence) that bears a potential risk for mutagenesis.^{60,61} Fluorescence is primarily emitted at wavelengths in the region between 260 and 320 nm.⁶¹ Photoreactivation experiments with yeast cell cultures showed a significant amount of DNA repair activity after 193-nm excimer laser irradiation. Although the role of

secondary radiation in the healing response is not yet fully understood⁶² there seems to be minimal concern among the physicians as clinical experience of more than 20 years has not shown that the potential risk is significant.

Additionally, it has been hypothesised⁵¹ that the mechanical stress of the cornea due to the pressure gradients associated to photoablation may result to structural changes in collagen and to the formation of scar tissue and loss of corneal transparency.

16.5 ACKNOWLEDGMENTS

The authors would like to thank Loukia Leonidou for assistance in preparing the figures in this chapter.

16.6 REFERENCES

1. D. W. Meltzer, ed., *Optics, Refraction and Contact Lenses*, Basic and Clinical Science Course, Sec. 3, American Academy of Ophthalmology, 1994.
2. W. M. Hart, ed., *Adler's Physiology of the Eye, Clinical Application*, 9th ed., Mosby Year Book, St. Louis, 1993.
3. R. M. Boynton, *Human Color Vision*, Optical Society of America, 1992.
4. *American National Standard for Ophthalmics, ANSI-Z80.28-2004, Methods for Reporting Optical Aberrations of Eyes*, American National Standards Institute, Inc., 2004.
5. J. Liang, B. Grimm, S. Goelz, and J. F. Bille, "Objective Measurement of Wave Aberration of the Human Eye with the Use of a Hartmann-Shack Wavefront Sensor," *J. Opt. Soc. Am. A*. **11**:1949–1957 (1994).
6. E. Moreno-Barriuso, S. Marcos, R. Navarro, and S. A. Burns, "Comparing Laser Ray Tracing, Spatially Resolved Refracometer and Hartmann–Shack Sensor to Measure the Ocular Wave Aberration," *Optom. Vis. Sci.* **78**:152–156 (2001).
7. V. Molebny, "Principles of Ray-tracing Aberrometry," in *1st International Congress of Wavefront Sensing and Aberration-Free Refractive Correction*, vol. 16, pp. 572–575 (2000).
8. F. Diaz-Doutón, J. Pujol, M. Arjona, and S. O. Luque, "Curvature Sensor for Ocular Wavefront Measurement," *Opt. Lett.* **31**:2245–2247 (2006).
9. C. Torti, S. Gruppeta, and L. Diaz-Santan., "Wavefront Curvature Sensing for the Human Eye," *J. Mod. Opt.* **55**:691–702 (2008).
10. I. Iglesias, R. Ragazzoni, Y. Julien, and P. Artal, "Extended Source Pyramid Wavefront Sensor for the Human Eye," *Opt. Express* **10**:419–428 (2002).
11. Jason D. Marsack, Larry N. Thibos, and Raymond A. Applegate, "Metrics of Optical Quality Derived from Wave Aberrations Predict Visual Performance," *J. Vis.* **4**(4):322–328 (2004).
12. D. R. Iskander, B. A. Davis, M. J. Collins, and R. Franklin, "Objective Refraction from Monochromatic Wavefront Aberrations via Zernike Power Polynomials," *Ophthalmic. Physiol. Opt.* **27**(3):245–255 (2007).
13. A. G. Bennett and R. B. Rabbetts, *Clinical Visual Optics*, Butterworths, London, 1984.
14. J. Liang and D. R. Williams, "Supernormal Vision and High Resolution Retinal Imaging through Adaptive Optics," *J. Opt. Soc. Am. A*. **14**:2884–2892 (1997).
15. P. Artal, F. Vargas, and I. Iglesias, "Detection of the Wavefront Aberration in the Human Eye and Its Partial Correction with a Liquid Crystal Light Modulator," in *International Workshop: Adaptive Optics for Industry and Medicine*, Shatura, Rusia, 1997.
16. M. B. McDonald, "Summit—Autonomus CustomCornea Laser in situ Keratomileusis Outcomes," in *1st International Congress of Wavefront Sensing and Aberration-Free Refractive Correction*, vol. 16, pp. 617–618, 2000.
17. H. Hofer, L. Chen, G. Y. Yoon, B. Singer, Y. Yamauchi, and D. R. Williams, "Improvement in Retinal Image Quality with Dynamic Correction of the Eye's Aberrations," *Opt. Exp.* **8**:631–643 (2001).
18. E. Dalimier, C. Dainty, and J. L. Barbur, "Effects of Higher-Order Aberrations on Contrast Acuity as a Function of Light Level," *J. Mod. Opt.* **55**(4):791–803 (2008).

19. Patricia A. Piers, S. Manzanera, P. M. Prieto, N. Gorceix, and P. Artal, "Use of Adaptive Optics to Determine the Optimal Ocular Spherical Aberration," *J. Cataract Refract. Surg.* **33**:1721–1726 (2007).
20. H. Ridley, "Intraocular Acrylic Lenses after Cataract Extraction," *Lancet* **1**(6699):118–121 (1952).
21. H. Ridley, "The Treatment of Cataract," *Practitioner* **178**(1067):525–533 (1957).
22. C. F. Lovisolo and D. Z. Reinstein, "Phakic Intraocular Lenses," *Surv. Ophthalmol.* **50**(6):549–587 (2005).
23. L. Espandar, J. J. Meyer, and M. Moshirfar, "Phakic Intraocular Lenses," *Curr. Opin. Ophthalmol.* **19**(4):349–356 (2008).
24. D. H. Chang and E. A. Davis, "Phakic Intraocular Lenses," *Curr. Opin. Ophthalmol.* **17**(1):99–104 (2006).
25. D. R. Sanders, M. R. Deitz, and D. Gallagher, "Factors Affecting Predictability of Radial Keratotomy," *Ophthalmology* **92**(9):1237–1243 (1985).
26. R. A. Applegate, H. C. Howland, R. P. Sharp, A. J. Cottingham, and R. W. Yee, "Corneal Aberrations and Visual Performance after Radial Keratotomy," *J. Refract. Surg.* **14**(4):397–407 (1998).
27. P. F. Vinger, W. F. Mieler, J. H. Oestreicher, and M. Easterbrook, "Ruptured Globes Following Radial and Hexagonal Keratotomy Surgery," *Arch. Ophthalmol.* **114**(2):129–134 (1996).
28. T. E. Burris, "Intrastromal Corneal Ring Technology: Results and Indications," *Curr. Opin. Ophthalmol.* **9**(4):9–14 (1998).
29. G. D. Kymionis, C. S. Siganos, N. S. Tsiklis, A. Anastasakis, S. H. Yoo, A. I. Pallikaris, N. Astyrakakis, and I. G. Pallikaris, "Long-term Follow-up of Intacs in Keratoconus," *Am. J. Ophthalmol.* **143**(2):236–244 (2007).
30. G. D. Kymionis, N. S. Tsiklis, A. I. Pallikaris, G. Kounis, V. F. Diakonis, N. Astyrakakis, and C. S. Siganos, "Long-term Follow-up of Intacs for Post-Lasik Corneal Ectasia," *Ophthalmology* **113**(11):1909–1917 (2006).
31. C. R. Munneryn, S. J. Koons, and J. Marshall, "Photorefractive Keratectomy: A Technique for Laser Refractive Surgery," *J. Cataract Refract. Surg.* **14**(1):46–52 (1988).
32. M. B. McDonald, J. M. Frantz, S. D. Klyce, R. W. Beuerman, R. Varnell, C. R. Munneryn, T. N. Clapham, B. Salmeron, and H. E. Kaufman, "Central Photorefractive Keratectomy for Myopia. The Blind Eye Study," *Arch. Ophthalmol.* **108**(6):799–808 (1990).
33. C. W. Jr. Flowers, P. J. McDonnell, and S. D. McLeod, "Excimer Laser Photorefractive Keratectomy," *Ophthalmol. Clin. North. Am.* **14**(2):274–283 (2001).
34. P. Fagerholm, "Phototherapeutic Keratectomy: 12 Years of Experience," *Acta Ophthalmol Scand.* **81**(1):19–32 (2003).
35. W. B. Trattler and S. D. Barnes, "Current Trends in Advanced Surface Ablation," *Curr. Opin. Ophthalmol.* **19**(4):330–334 (2008).
36. M. Camellin, "Laser Epithelial Keratomileusis for Myopia," *J. Refract. Surg.* **19**(6):666–670 (2003).
37. I. G. Pallikaris, V. J. Katsanevaki, M. I. Kalyvianaki, and I. I. Naoumidi, "Advances in Subepithelial Excimer Refractive Surgery Techniques: Epi-Lasik," *Curr. Opin. Ophthalmol.* **14**(4):207–212 (2003).
38. I. G. Pallikaris, M. I. Kalyvianaki, V. J. Katsanevaki, and H. S. Ginis, "Epi-Lasik: Preliminary Clinical Results of an Alternative Surface Ablation Procedure," *J. Cataract Refract. Surg.* **31**(5):879–885 (2005).
39. I. G. Pallikaris, M. E. Papatzanaki, E. Z. Stathi, O. Frenschock, and A. Georgiadis, "Laser in situ Keratomileusis," *Lasers Surg. Med.* **10**(5):463–468 (1990).
40. I. G. Pallikaris, M. E. Papatzanaki, D. S. Siganos, and M. K. Tsilimbaris, "A Corneal Flap Technique for Laser in situ Keratomileusis. Human Studies," *Arch. Ophthalmol.* **109**(12):1699–1702 (1991).
41. C. R. Munneryn, S. J. Koons, and J. Marshall, "Photorefractive Keratectomy: A Technique for Laser Refractive Surgery," *J. Cataract Refract Surg.* **Jan**; **14**(1):46–52 (1988).
42. P. J. McDonnell, H. Moreira, J. Garbus, T. N. Clapham, J. D'Arcy, and C. R. Munneryn, "Photorefractive Keratectomy to Create Toric Ablations for Correction of Astigmatism," *Arch. Ophthalmol.* **May**; **109**(5):710–3 (1991).
43. T. Seiler, M. H. Dastjerdi, "Customized Corneal Ablation," *Curr. Opin. Ophthalmol.* **Aug**; **13**(4):256–260 (2002).
44. D. Basting, ed., *Excimer Laser Technology: Laser Sources, Optics, Systems and Applications*, Lambda Physik, Gottingen, 2001.
45. V. M. Borisov, I. E. Bragin, A. Yu. Vinokhodov, and V. A. Vodchits, "Pumping Rate of Electric-Discharge Excimer Lasers," *Quantum Electronics* **25**(6):507–510 (1995).

46. R. Srinivasan and E. Sutcliffe, "Dynamics of the Ultraviolet Laser Ablation of Corneal Tissue," *Am. J. Ophthalmol.* **103**(3 pt 2):470–471 (1987).
47. R. Srinivasan, "Ablation of Polymers and Biological Tissue by Ultraviolet Lasers," *Science* **234**(4776):559–65 (1986).
48. S. L. Trokel, R. Srinivasan, and B. Braren, "Excimer Laser Surgery of the Cornea," *Am. J. Ophthalmol.* **96**(6):710–715 (1983).
49. G. B. Petit and M. N. Ediger, "Corneal-tissue Absorption Coefficients for 193- and 213-nm Ultraviolet Radiation," *App. Opt.* **35**:3386–3391 (1995).
50. P. J. Dougherty, K. L. Wellish, and R. K. Maloney, "Excimer Laser Ablation Rate and Corneal Hydration," *Am. J. Ophthalmol.* **118**:169–176 (1994).
51. O. Kermani and H. Lubatschowski, "Structure and Dynamics of Photo-Acoustic Shock-waves in 193 nm Excimer Laser Photo-ablation of the Cornea," *Fortschr Ophthalmol.* **88**(6):748–753 (1991).
52. D. W. Hahn, M. N. Ediger, and G. H. Pettit, "Dynamics of Ablation Plume Particles Generated during Excimer Laser Corneal Ablation," *Lasers Surg. Med.* **16**(4):384–389 (1995).
53. T. Seiler and P. J. McDonnell, "Excimer Laser Photorefractive Keratectomy," *Surv. Ophthalmol.* **40**:89–118 (1995).
54. Surgical Lasers Operation Manuals and Data sheets. (Various.)
55. M. H. Feltham and R. J. Wolfe, "Some Variables to Consider to Avoid the Need for Lasik Surgical Enhancements," *Clin. Exp. Optom.* **83**(2):76–81 (2000).
56. I. G. Pallikaris, H. S. Giniis, G. A. Gounis, D. Anglos, T. G. Papazoglou, and L. P. Naoumidis, "Corneal Hydration Monitored by Laser-Induced Breakdown Spectroscopy," *J. Refract. Surg.* **14**(6):655–660 (1998).
57. C. Maldonado-Codina, P. B. Morgan, and N. Efron, "Thermal Consequences of Photorefractive Keratectomy," *Cornea* **20**(5):509–515 (2001).
58. T. Bende, T. Seiler, and J. Wollensak, "Side Effects in Excimer Corneal Surgery: Corneal Thermal Gradients," *Graefes Arch. Clin. Exp. Ophthalmol.* **226**:277–280 (1988).
59. M. Vetrugno, A. Maino, E. Valenzano, and E. Cardia, "Corneal Temperature Changes during Photorefractive Keratectomy using the Laserscan 2000 Flying Spot Laser," *J. Refract. Surg.* **17**(4):454–459 (2001).
60. T. Seiler, T. Bende, K. Winckler, and J. Wollensak, "Side Effects in Excimer Corneal Surgery. DNA Damage as a result of 193 nm Excimer Laser Radiation," *Graefes Arch. Clin. Exp. Ophthalmol.* **226**(3):273–276 (1988).
61. N. Muller-Stolzenburg, S. Schrunder, J. Helfmann, H. J. Buchwald, and G. I. Muller, "Fluorescence Behavior of the Cornea with 193 nm Excimer Laser Irradiation," *Fortschr Ophthalmol.* **87**(6):653–658 (1990).
62. Z. Z. Nagy, J. Toth, A. Nagymihaly, and I. Suveges, "The Role of Ultraviolet-b in Corneal Healing Following Excimer Laser in situ Keratomileusis," *Pathol. Oncol. Res.* **8**(1):41–46 (2002).

This page intentionally left blank.

DO NOT DUPLICATE

THREE-DIMENSIONAL CONFOCAL MICROSCOPY OF THE LIVING HUMAN CORNEA

Barry R. Masters

*Department of Biological Engineering
Massachusetts Institute of Technology
Cambridge, Massachusetts*

17.1 GLOSSARY

Adaptive optics. An optical system that consists of a wavefront sensor that can measure the wavefront of the light, that is, a Shack-Hartmann sensor, a deformable mirror that can alter the wavefront of the light, and a closed loop-feedback control system to minimize the wavefront error. An example is the use of adaptive optics in a scanning laser ophthalmoscope, in which the wave front distortions due to the cornea and the ocular lens are corrected by the adaptive optics.

Confocal microscope. A microscope that is based on spatial filtering with two pinholes or slit apertures located in conjugate planes. The word confocal means that both the point source of light and the point detector are cofocused on the same focal volume within the specimen. A point source of light (source aperture) illuminates a diffraction limited focal volume within the specimen. The scattered or emitted light from this focal volume is focused on a detection pinhole located in a plane that is conjugate to the source aperture. A large area detector is placed behind the photodetector pinhole.

Diffraction limited imaging. The theoretical resolution of an optical microscope is limited by the diffraction of light. Abbe defined this diffraction limit for microscopic resolution. The diffraction limited resolution is proportional to numerical aperture of the microscope objective, that is, its ability to collect light, and inversely proportional to the wavelength of the light from the specimen. A microscope that can image at the theoretical limit is a diffraction limited imaging system. Today, there are optical microscopes that exceed the Abbe diffraction limit.

Laser scanning confocal microscope. A confocal microscope in which the laser beam is scanned over the back focal plane of the microscope objective with the simultaneous scanning of the diffraction limited beam over the specimen. The scattered or emitted light from the focal volume is collected by an epi-illumination system, descanned, and passes through an aperture placed in front of a large area photodetector.

Linear optics. The response of the illuminated material to the electromagnetic field is linear with the amplitude of the field. Examples are one-photon absorption, scattering, and fluorescence.

Microlens Nipkow disk confocal microscope. This is a major improvement over the Nipkow disk tandem scanning confocal microscope that has an illumination efficiency of only 1 to 2 percent. In this new device a second Nipkow disk which contains about 20,000 microlenses is added to the

Nipkow disk microscope. The lower Nipkow disk contains another 20,000 pinholes that are arranged in the same pattern as the microlenses on the upper disk. The lower pinhole disk is located in the focal plane of the upper microlens disk. This arrangement increases the illumination efficiency to 50 percent and results in a high sensitivity even for weak reflecting specimens. The same pinholes are used for illumination (laser excitation light) and spatial filtering of the fluorescence light, and a dichroic mirror placed between the two disks separates the excitation light from the fluorescence that is focused on the detector. Both disks rotate together on the same axis to form a real-time scanned image of the specimen.

Nipkow disk tandem-scanning confocal microscope. A real-time, direct-view, tandem-scanning microscope based on a rotating Nipkow disk, contains sets of conjugate pinholes arranged in a set of spirals about the axis of rotation. On a given diameter of the disk there is one pinhole for the illumination and one pinhole for detection of the scattered or emitted light from the specimen. As the disk rotates, many conjugate sets of illumination and detection pinholes (arranged in conjugate planes on opposite sides of the Nipkow disk) operate in parallel to simultaneously scan the specimen with a diffraction limited illumination spot, and to detect the descanned scattered or emitted light.

Nonlinear optics. The response of the illuminated material to the electromagnetic field is nonlinear with the amplitude of the field. At high light intensities, the incident field modifies the optical properties of the material such that the light waves can interact and exchange energy and momentum. Examples of nonlinear optics are two-photon excitation fluorescence, second-harmonic generation, and third-harmonic generation.

Numerical aperture. The numerical aperture (NA) of a microscope objective is a measure of the range of angles that the lens can accept or emit light. The NA of a microscope objective is defined as: $NA = n \sin \theta$, where n is the index of refraction of the medium between the objective and the specimen, and θ is the half-angle of the maximum cone of light that can enter the objective.

Optical aberrations. They represent the failure of an optical system to produce a perfect image. Modern microscope objectives are made to minimize five categories of optical aberrations: spherical aberrations (results in a lack of sharp focus), coma (results when light is focused at points off the optical axis), astigmatism, field curvature, and distortion. In addition to these monochromatic Seidel aberrations, there are also axial chromatic aberrations and lateral chromatic aberrations in which different colors of light focus at different positions.

Scanning slit confocal microscope. A real-time, confocal microscope in which a slit of incident light is scanned over the back focal plane of the microscope objective. The reflected and the emitted light from the focal volume within the specimen is collected by the microscope objective, descanned by the scanning system, and detected by a wide-area detector.

Second-harmonic generation (SHG). A nonlinear process in which an incident wave of frequency ω interacts with a material that does not contain an inversion symmetry, and generates a new signal at the frequency 2ω .

Slit-lamp microscope. A microscope that is based on the principle of focal illumination. A slit of light is used to illuminate a volume within the eye. A biomicroscope collects the scattered and reflected light from a volume that intersects the illumination volume. The intersection of the illumination volume and the collection volume is in focus.

Specular microscope. A modification of the slit lamp in which one half of the microscope objective is used for illumination of the specimen and one half of the microscope objective is used to collect the scattered and the reflected light. For a plane reflecting surface when the angle between the normal to the surface and the illumination beam of light is equal to the angle between the normal to the surface and the reflected beam of light, the condition of specular reflection exists. The microscope that images this specular reflection is called a specular microscope. It is used to image the endothelial cells due to the strong specular reflection at the endothelial cell and the aqueous humor interface.

Third-harmonic generation (THG). A nonlinear process in which an incident wave of frequency ω interacts with any material, and generates a new signal at the frequency 3ω .

17.2 INTRODUCTION

This chapter presents the optical principles of the confocal microscope and its applications to the design of clinical confocal microscopes for imaging the living human cornea. A series of technical advances resulted in new optical designs from the early slit lamp with its focal plane illumination to the modern scanning slit and laser scanning clinical confocal microscopes. A carefully edited selection of papers, books, and commercial Web sites is presented as a further resource of clinical images, instrument designs, and clinical protocols. The emphasis of this chapter is to critically examine the limitations of the various designs of clinical instruments, and to present the author's perspective on the state of the art instruments and the proposed future instrument development that will incorporate the suggested improvements as well as those based on the emerging field of nonlinear optical microscopy,^{1,2} specifically second- and third-generation harmonic microscopy. It is further suggested that future clinical biomicroscopes will incorporate correlative microscopy, that is, both linear reflected confocal microscopy together with nonlinear higher-harmonic microscopy.

17.3 THEORY OF CONFOCAL MICROSCOPY

The principles of the confocal microscope in both the transmission mode and in the epi-illumination mode were clearly stated and illustrated in the 1961 patent granted to Marvin Minsky.³ To facilitate easy access to the key papers and the important patents in the field of confocal microscopy and its many biological applications SPIE Press published an edited collection of these documents.³ Further theoretical and practical details on instrument design, theories of transverse, and axial optical resolution are contained in a modern book.⁴

A widefield optical microscope has no capability for axial optical sectioning; therefore, when the specimen is a thick, highly scattering tissue the images at the focal plane are degraded by out-of-focus contributions to the image. The confocal microscope solves this problem with the method of spatial filtering.³ A point source of light from an incoherent or a coherent source illuminates a point in the specimen via the microscope objective. The back-scattered or the fluorescence light from that point is collected by the microscope objective and passes an aperture that is in a conjugate plane to the point source of illumination. Only light from the focal plane of the illumination spot passes the conjugate aperture and is detected by the detector located directly behind the detection aperture. Light from the out-of-focal plane cannot pass the detection aperture and therefore it is not detected. This confocal principle is shown in Fig. 1. The use of spatial filtering provides an enhanced axial resolution (z axis of the microscope) and is the source of the axial optical sectioning in the confocal microscope.

In order to acquire two-dimensional images of the focal plane it is necessary to scan the light beam over the back focal plane of the microscope objective and to descann the scattered or emitted light that is collected from the microscope objective. Typically two sets of oscillating galvanometer driven mirrors, one for the horizontal scan and one for the vertical scan, are used for laser scanning and descanning.

17.4 THE DEVELOPMENT OF CONFOCAL INSTRUMENTS

The late David Maurice is honored by the vision research community for his outstanding contributions to the development of optical ophthalmic instrumentation.⁵ The development of the scanning slit confocal microscope is based on a series of linked technical advances that include parallel invention as well as unique creative invention, and the story of the development of live cell imaging in ophthalmology is described in a recent book.⁶

Two groups of confocal microscopes were developed for ophthalmology: the first was a laser scanning ophthalmoscope for imaging the retina was developed by Webb and his collaborators⁷ and the

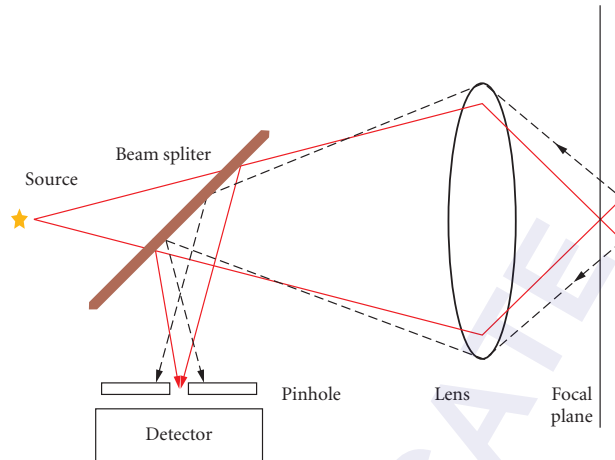


FIGURE 1 The principle of the confocal microscope. This schematic diagram shows the depth discrimination capability of a confocal microscope and how it discriminates against reflected light from out of the focal plane. The dashed vertical line to the right of the focal plane represents an out-of-focus plane. Only the light rays (solid lines) from the focal plane pass the pinhole and are detected. (Reproduced with permission from Ref. 6.)

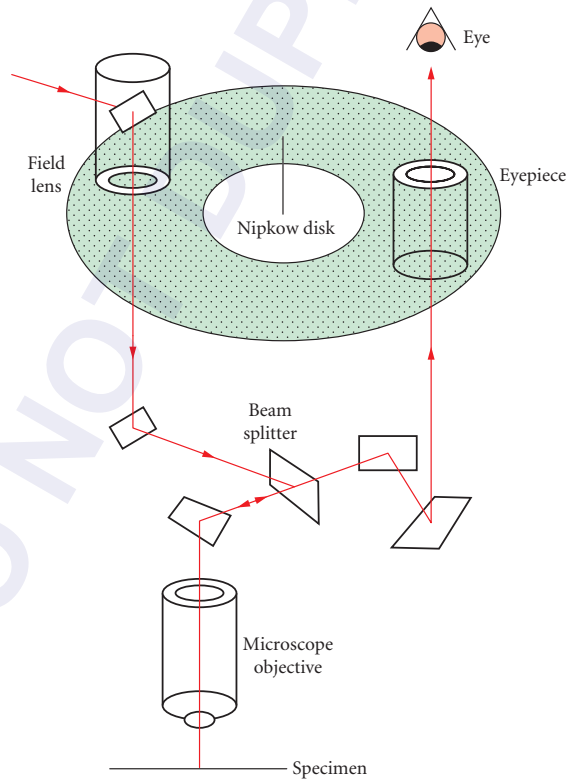


FIGURE 2 The principle of Nipkow disk confocal microscope. This schematic shows the real-time, direct-view, tandem-scanning Nipkow disk confocal microscope. The light source from the upper left is a mercury arc lamp or a tungsten filament lamp. (Reproduced with permission from Ref. 6.)

second group of confocal microscopes that typically used an incoherent light source was developed for imaging the cornea.⁸

Several early types of *in vivo* confocal microscopes were based on the spinning Nipkow disk that contained spiral arrays of multiple sets of conjugate pinholes, one pinhole was used for illumination, and a conjugate pinhole was used for spatial filtering in front of the detector.^{3,6} The design of the Nipkow disk based confocal microscope is shown in Fig. 2.

The Nipkow disk-based confocal microscope has a very poor illumination efficiency (1 to 2 percent) and therefore weak reflecting corneal cell layers such as the basal epithelial cell layer cannot be imaged. The solution to this problem was made by Akira Ichihara and his coworkers at the Yokogawa Institute Corporation in Japan with their invention of a microlens Nipkow disk confocal microscope.⁶ The addition of an upper disk that contained 20,000 microlenses which focused the light on the 20,000 pinholes on the lower disk. The key point is that the lower pinhole disk is aligned and located in the focal plane of the upper disk. The resulting illumination efficiency is about 50 percent. Figure 3 shows the principle of the microlens Nipkow disk confocal microscope.

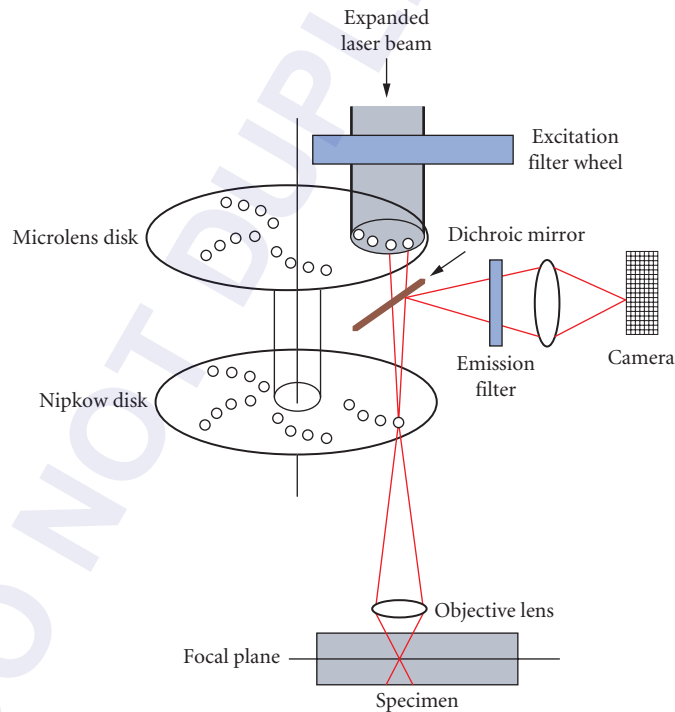


FIGURE 3 The principle of microlens Nipkow disk confocal microscope. Schematic drawing of the Yokogawa microlens Nipkow disk confocal microscope. The two disks rotate on a common axis. The upper microlens disk contains 20,000 microlenses arranged in a series of spirals. The lower Nipkow disk contains 20,000 pinholes arranged in a series of spirals. The lower disk containing the pinholes is in the focal plane of the microlenses in the upper disk. A dichroic mirror placed between the two disks passes the incident laser beam, but reflects the beam of emitted light to the emission filter and then it is focused on the camera. (Reproduced with permission from Ref. 6.)

17.5 THE SCANNING SLIT AND LASER SCANNING CLINICAL CONFOCAL MICROSCOPES

Based on the Svishchev confocal microscope design of a two-sided oscillating mirror for scanning a slit of light across the back focal plane of a high NA, water immersion microscope objective, a clinical confocal microscope with an incoherent halogen lamp light source was designed and constructed.⁹ Figure 4 shows the principle of the two-sided, oscillating mirror in the scanning slit confocal microscope.

The advantage of this scanning slit confocal microscope is that the slit widths of the two conjugate slits, one for the illumination side and one for the detection side, are completely adjustable. They could be opened to increase the reflected signal from weak scattering corneas, or could be reduced to increase the resolution of the microscope.¹⁰ Later commercial designs used fixed slits.

This scanning slit confocal microscope was placed in the clinic and was used to study a variety of normal and clinical conditions of the in vivo human cornea. Figure 5 shows the principle of the clinical scanning slit confocal microscope based on the Svishchev oscillating two-sided mirror design.

For example, this scanning slit confocal microscope was successfully used to study the basal epithelial cell layer, and its adjacent wing cell layer. In the clinic the confocal microscope was used to relocate the same cells within the cornea over different periods of time, this provided a method to track cell migration over time. Since the microscope objective contacted the anterior surface of the cornea there is the need for both an anesthetic and the application of a drop of an index matching gel to the tip of the microscope objective.

In recent years the scanning slit confocal microscope has been made into a commercial product that is now in the fourth generation.¹¹ One useful feature of the commercial version, the Confoscan 4, is the 20x microscope objective that works under noncontact conditions; therefore, there is no need

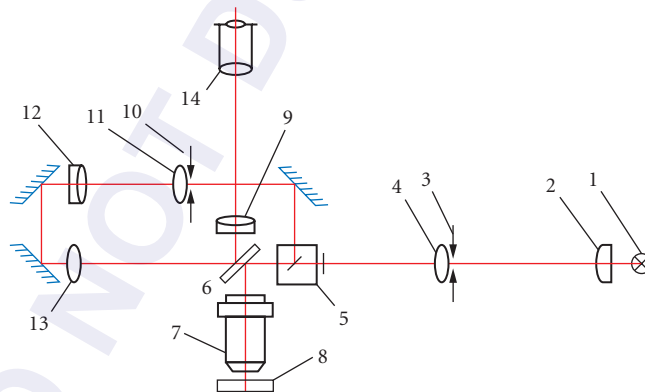


FIGURE 4 The principle of the Svishchev two-sided, oscillating mirror, scanning slit confocal microscope. The light source, 1, is projected by a condenser lens, 2, onto the first slit, 3; the light passes through a prism cube, 5; an image of the first slit is scanned over the back focal plane of the microscope objective, 7, by the two-sided oscillating mirror, 6, which descans the reflected light from the focal plane of the specimen. The second slit, 10, is conjugate with the first slit (confocal) and excludes the light that is not in the specimen's focal plane, 8. The two-sided oscillating mirror, 6, performs three functions: it scans the image of the slit 3 over the back focal plane of the microscope objective, descans the beam from the object, and rescans the beam for observation in the ocular, 14. The parts labeled 2, 4, 9, 11, 12, and 13, are lenses. (Reproduced with permission from Ref. 6.)

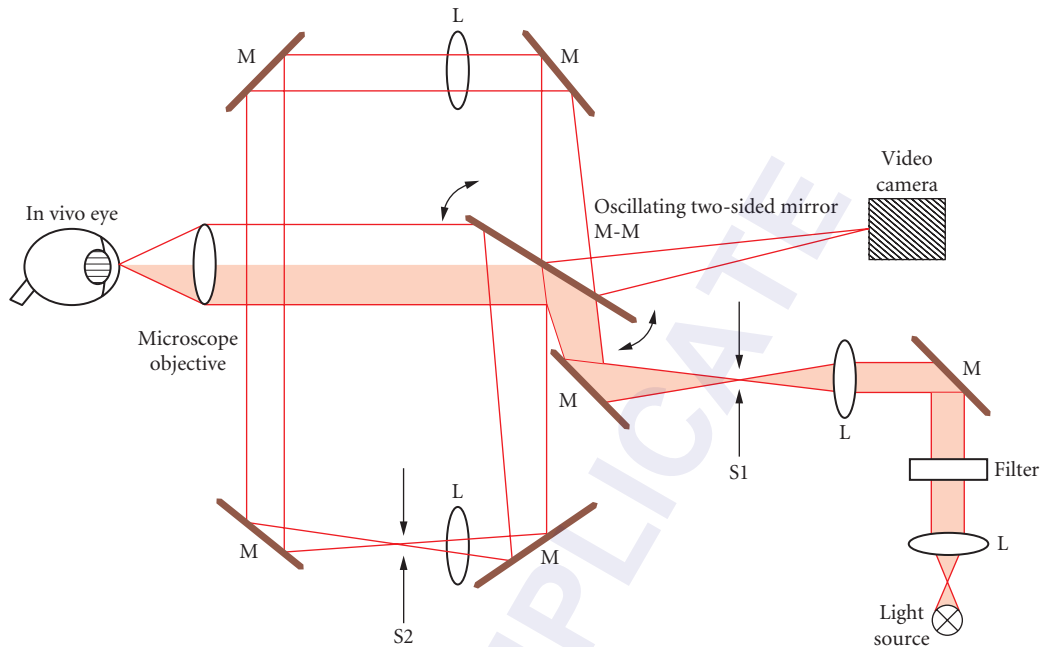


FIGURE 5 The optical principles of the clinical, video-rate, scanning-slit, in vivo, confocal microscope with two adjustable slits: S1 for illumination and S2 for imaging. An oscillating two-sided mirror M-M scans S1 over the back focal plane of the objective, descans the collected light and directs it to S2, then rescans the beam, which is imaged on the photocathode of a video camera. The dark ray path indicated the illumination light rays. The light ray path indicates the reflected light that is collected by the microscope objective. One half of the NA of the microscope objective is used for illumination, and one half of the NA is used for collection of the reflected light form the specimen. (Reproduced with permission from Ref. 6.)

for anesthetic nor for an index matching gel on the tip of the objective. Numerous publications based on the slit scanning confocal microscope are replete with confocal images of normal corneas and the diseased and altered corneas.¹²⁻¹⁶ In addition, the Nidek Web site features many examples of confocal microscopic images of the cornea.¹¹

Modern clinical confocal microscopes can acquire stacks of confocal images through the full thickness of the human cornea and produce three-dimensional visualizations of the a volume from the anterior surface to the posterior endothelial cell layer. This approach of three-dimensional visualization of the living cornea was first demonstrated on rabbit corneas by Masters and coworkers.¹⁷⁻¹⁹

An alternative to the commercial scanning slit confocal microscope is a laser scanning confocal microscope designed for retinal imaging by Heidelberg Engineering that can be modified with a Rostock Cornea Module (RCM) microscope that permits confocal microscopy of the cornea.²⁰ The Heidelberg Retinal Tomograph (HRT) is a confocal laser scanning microscope with a diode laser light source at 670 nm. The optical principle of this confocal microscope is shown in Fig. 1.

The RCM unit has the ability to change the focal plane within the cornea. A single-use planar cap is fitted over the tip of the module and it is in contact with the anterior surface of the cornea. The plastic cap stabilizes the distance between the module and the cornea. An index matching gel is placed between the instrument tip and cornea. Between the movable plastic cap is an index matching gel that serves to increase the numerical aperture of the objective which results in increased resolution and contrast. When it is desirable to use the confocal microscope in the noncontact mode, the water immersion microscope objective is removed together with its contact cap and is replaced with a long

focal length dry objective. Typically a noncontact Nikon 50X, NA 0.45 microscope objective is used to image the cornea in the noncontact mode. The specifications from the manufacturer include a focus range of 4 mm, the lateral image size of 380 by 380 μm , the lateral optical resolution of about 1 μm , and the depth of focus of about 1 to 2 μm .

17.6 CLINICAL APPLICATIONS OF CONFOCAL MICROSCOPY

The clinical utility of a clinical confocal microscope for the examination of the *in vivo* human cornea depends on the optical and mechanical design parameters of the confocal microscope, the optical condition of the subject's cornea, and finally, on the skill and the experience of the clinician who is making the examination.

A general introduction to the principles and the practice of clinical confocal microscopy is a rich resource of clinical protocols. A review on the confocal microscopy of the cornea by Böhnke and Masters provides the accumulated clinical experience of many years of patient examination.²¹ This practical discussion of clinical techniques includes the following topics: technical setup of the confocal microscope, the technique of biomicroscopy with the confocal microscope, the suggested clinical examination procedure, patient data archiving and retrieval, and the dissemination, documentation, and publication of microscopic findings. This paper also contains an atlas of confocal images from the normal cornea, the effects of aging, the effects of contact lens wear, the cornea with known pathologies, and the postsurgical cornea.

Creative scientists, engineers, and clinicians have worked together in the last decades to develop new clinical confocal microscopes that have successfully been incorporated into clinical practice. Optics and optical engineering are in the forefront of the design and development of medical devices for both treatment and diagnostics providing rich career opportunities for highly trained people.

Two books are available with many clinical examples of the application of clinical confocal microscopy to the cornea.^{22,23} These books are replete with confocal images from a wide variety of normal and pathological corneas and cover the effects of photorefractive surgery on the structure of the cornea. Another source of clinical corneal images is the World Wide Web. I include two commercial sites of companies that manufacture clinical confocal instruments.^{11,20} Their Web sites include many examples of confocal images of the cornea in a variety of normal and pathological conditions.

As an exemplar of the use of the clinical confocal microscopy to investigate new clinical findings the following study is presented.²⁴ This paper illustrates the confounding effect of acquiring clinical confocal images of the human cornea with inappropriate optical resolution. The authors posed the following question: what are the long-term effects of contact lens wear on the human cornea? They answered this question by first acquiring confocal images through the full thickness of the *in vivo* human cornea, and then by a frame by frame analysis to determine morphological changes within the cornea that are correlated with long-term contact lens wear and do not occur in subjects that do not wear contact lenses.

The scanning slit confocal microscope with a 50X/1.0 NA water immersion objective was used to investigate the corneal morphology in long-term contact lens wearers. The authors investigated 13 patients with a history of up to 26 years of soft contact lens wear, 11 patients with a history of up to 25 years of rigid gas permeable contact lens wear, and a control group of 29 normal subjects without a history of contact lens wear. For contact lens wearers epithelial microcystic changes and alterations of endothelial cell morphology were found as described previously. The significant new finding was there were highly reflective panstromal microdot (submicron) deposits in the entire thickness of the stroma for the contact lens wearers. It was concluded that this newly observed stromal microdot degeneration scales with the years of contact lens wear and may be the early stage of a significant corneal disease. Further correlative microscopic studies involved electron microscopy of *ex vivo* human corneas, and spectroscopic studies of the microdots. Initially, other groups were not able to observe the stromal microdot deposits because they used a Nipkow disk confocal microscope with a low power and a low NA microscope objective. The low power, low NA microscope objective did not

provide the necessary resolution to image the stromal microdot deposits. This study illustrates the requirement of the appropriate optical resolution in studies with the clinical confocal microscope, as well as correlative electron microscopy to validate the clinical conclusions.

17.7 PERSPECTIVES

In the linear domain confocal microscopes can be developed that use different colored light-emitting diodes as the light source. The degree of penetration of the incident light is a function of the wavelength and the amount of absorbing and scattering material in a particular cornea. The selection of various colors of illumination light would provide the clinician with improved image contrast since the illumination could be tuned to optimize either the penetration depth or the contrast.

The use of noncontact confocal microscopes such as the Heidelberg Retina Tomograph II with the Rostock Corneal Module is a major improvement over the previous designs with an applanating microscope objective that flattens the corneal surface. The noncontact confocal microscope is less likely to transmit microbes from patient to patient and causes less trauma to the cornea. The Nidek ConfoScan4 clinical confocal microscope also operates in a noncontact mode.

Another problem that still requires further development to reach a solution is a method to determine the position of the focal plane within the cornea. The current use of movable lenses together with accurate position sensors do not provide the true position of the focal plane within the cornea with respect to the corneal surface.

Typically, the technical specifications of commercial clinical confocal microscopes do not include information on the point-spread function, or other metrics of the transverse and the optical resolution of the microscope with a stated frequency of illumination, and a stated numerical aperture of the microscope objective. Measurements of the image, a series of optical sections, through a subdiffraction of gold particle or a 0.1- μm fluorescent bead are readily made and would provide useful information.

Microscope images obtained through thick, scattering, and absorbing tissues such as a cornea containing scar tissues result in further degradation of the image quality. Even when a gel is used to approximately match the refractive index of the cornea and the microscope objective there are induced optical aberrations that degrade the image quality. For example, the refractive index varies within the cornea, both in depth and in the imaging plane of the microscope. One option to mitigate this problem is to incorporate adaptive optics (a closed loop-feedback system that corrects the wavefront distortions due to the cornea). Adaptive optics have successfully been incorporated into scanning-laser ophthalmoscopes for retinal imaging.²⁵

All of the clinical confocal microscopes in use today are based on linear optics and generate contrast from the absorption and scattering of the incident light. There are many published studies of nonlinear higher-harmonic generation microscopy that were performed on *ex vivo* corneas, but there is a paucity of *in vivo* animal studies and no *in vivo* human studies have been reported.²⁶ Second- and third-harmonic generation can be explained by the theory of the nonlinear susceptibility in which for the time domain the polarization is expanded in a power series of the sum of products of the linear susceptibility and the electric field, the second-order susceptibility and the square of the electric field, and the third-order susceptibility and the cube of the electric field and higher-order terms. Second-harmonic generation (SHG) is described by the second-order susceptibility and third-harmonic generation (THG) is described by the third-order susceptibility. THG has a broad potential for tissue imaging since it occurs in all materials, including dielectric materials with inversion symmetry. In general THG is a weak process; however, it is dipole allowed. THG microscopy has an emerging potential for cell and tissue imaging. As new techniques are being developed for the surface-enhanced THG at interfaces, the technique may become useful for imaging the cornea. Another emerging development is the use of enhanced THG microscopy with nanogold particles to utilize the surface-plasmon-resonance effect.

Both SHG and THG imaging do not exhibit the saturation effects or the photobleaching effects that are associated with multiphoton excitation fluorescence microscopy. These microscopic techniques can

penetrate into millimeters of tissue and provide submicron three-dimensional optical sectioning. Collagen has a noncentrosymmetric structure and can convert the incident ultrashort laser pulse to its second harmonic, thus providing a noninvasive imaging modality for the cornea. The backward-scattered SHG signal may provide a sensitive clinical method to study corneal haze following photorefractive surgery.

Structures within the cornea, for example, collagen fibrils that are a few nanometers in diameter, can only be imaged with an electron microscope. When imaged with a confocal microscope or with a SHG or THG microscope the collagen fibrils, and all subdiffraction limit structures appear at the diffraction limit of the microscope. This optical phenomena is rarely explained in many papers on nonlinear microscopy of the cornea.

Finally, the image quality of a clinical corneal microscope is dependent on the microscope objective. Olympus Corporation has developed a new 25X, NA 1.05 water immersion microscope objective that is optimized for nonlinear microscopy based on femtosecond pulses of laser light. This lens that produced tightly focused light in the focal volume, a high collection efficiency due to the high NA, and a long free-working distance may enhance the performance of corneal microscope that operates in the nonlinear optical domain.

17.8 SUMMARY

The goals of biomicroscopy of the living human cornea are to improve our understanding of the cellular structure and function in the normal cornea and in diseased, aging, and surgically altered corneas. Biomicroscopy of the anterior segment of the eye has provided the clinician with a series of improved instruments from the early slit lamps to a variety of confocal microscopes; with each new design there were improvements in both the image resolution and the contrast. Modern clinical confocal microscopes provide transverse resolution and axial resolution of several microns. Both the scanning slit confocal microscope and the laser scanning confocal microscope provide high-resolution images with high contrast and are ideal to investigate nerve regeneration in the postsurgical cornea, the formation of scar tissue, and other morphological anomalies due to disease, contact lens wear, overnight lid closure, or surgical procedures such as photorefractive surgery. Routine confocal microscopy of the cornea reveals the cellular and subcellular structure of all the cells in the full thickness of the cornea. Nonlinear microscopy, especially higher-harmonic generation microscopy has the potential for molecular imaging with the benefit that it is based on a scattering process and, therefore, there is no net energy deposition within the ocular tissue. Future clinical biomicroscopes can incorporate correlative microscopy, that is, both linear reflected confocal microscopy and nonlinear higher-harmonic microscopy.

17.9 ACKNOWLEDGMENTS

The author and Professor M. Böhnke shared the 1999 Alfred Vogt-Prize for Ophthalmology (the highest award in Switzerland for scientific research in ophthalmology) from the Alfred Vogt-Stiftung zur Förderung der Augenheilkunde Zürich, for their work: "Confocal Microscopy of the Cornea." The author thanks Professor M. Böhnke, formerly of the Department of Ophthalmology, University of Bern, for many years of collaboration, and Dr. Andreas A. Thae for his collaboration in the development of the clinical confocal microscope.

17.10 REFERENCES

1. B. R. Masters and P. T. C. So, *Handbook of Biomedical Nonlinear Optical Microscopy*, Oxford University Press, New York, 2008.

2. R. W. Boyd, *Nonlinear Optics*, 3rd ed., Academic Press, San Diego, CA, 2008.
3. B. R. Masters, *Selected Papers on Confocal Microscopy*, Milestone Series MS 131, SPIE Optical Engineering Press, Bellingham, WA, 1996.
4. T. R. Corle and G. S. Kino, *Confocal Scanning Optical Microscope and Related Imaging Systems*, Academic Press, San Diego, CA, 1996.
5. B. R. Masters, "David Maurice's Contributions to Optical Ophthalmic Instrumentation: Roots of the Scanning Slit Confocal Microscope," *Experimental Eye Research* **78**:315–326 (2004).
6. B. R. Masters, *Confocal Microscopy and Multiphoton Excitation Microscopy: The Genesis of Live Cell Imaging*, SPIE Optical Engineering Press, Bellingham, WA, 2006.
7. R. H. Webb, "Scanning Laser Ophthalmoscope," In B. R. Masters, (ed.), *Noninvasive Diagnostic Techniques in Ophthalmology*, Springer Verlag, New York, pp. 438–540, 1990.
8. B. R. Masters, *Noninvasive Diagnostic Techniques in Ophthalmology*, Springer Verlag, New York, 1990.
9. B. R. Masters and A. A. Thaer, "Real-Time Scanning Slit Confocal Microscopy of the In Vivo Human Cornea," *Applied Optics*, **33**(4):695–701 (1994).
10. B. R. Masters, "Scanning Slit Confocal Microscopy of the In Vivo Cornea," *Optical Engineering* **34**(3):684–692 (1995).
11. <http://usa.nidek.com>; accessed May 26, 2009.
12. B. R. Masters and A. A. Thaer, "In Vivo Human Corneal Confocal Microscopy of Identical Fields of Subepithelial Nerve Plexus, Basal Epithelial and Wing Cells at Different Times," *Microscopy Research and Technique* **29**:350–356 (1994).
13. B. R. Masters and A. A. Thaer, "In Vivo, Real-Time Confocal Microscopy of the Continuous Wing Cell Layer Adjacent to the Basal Epithelium in the Human Cornea: a New Benchmark for In Vivo Corneal Microscopy," *Bioimages* **3**(1):7–11 (1995).
14. B. R. Masters and A. A. Thaer, "Real-Time Confocal Microscopy of In Vivo Human Corneal Nerves," *Bioimages*, **4**(3):129–134 (1997).
15. B. R. Masters and M. Böhnke, "Three-Dimensional Confocal Microscopy of the Human Cornea in Vivo," *Ophthalmic Research* **33**(3):125–135 (2001).
16. B. R. Masters and M. Böhnke, "Three-Dimensional Confocal Microscopy of the Living Eye," *Annual Review of Biomedical Engineering*, Annual Reviews, Palo Alto, CA, **4**:69–91 (2002).
17. B. R. Masters and S. W. Paddock, "Three-Dimensional Reconstruction of the Rabbit Cornea by Confocal Scanning Optical Microscopy and Volume Rendering," *Applied Optics* **29**:3816–3822 (1990).
18. B. R. Masters and M. A. Farmer, "Three-Dimensional Confocal Microscopy and Visualization of the In Situ Cornea," *Computerized Medical Imaging and Graphics*, **17**(3):211–219 (1993).
19. B. R. Masters, "Three-Dimensional Confocal Microscopy of the Living In Situ Rabbit Cornea," *Optics Express* **3**(9):351–355 (1998).
20. <http://www.heidelbergengineering.com>; accessed May 26, 2009.
21. M. Böhnke and B. R. Masters, "Confocal Microscopy of the Cornea," *Progress in Retina & Eye Research* **18**(5): 553–628 (1999).
22. L. Mastropasqua and M. Nubile, *Confocal Microscopy of the Cornea*, SLACK Incorporated, Thorofare, NJ, 2002.
23. R. F. Guthoff, C. Baudouin, and J. Stave, *Atlas of Confocal Laser Scanning In-Vivo Microscopy in Ophthalmology*, Springer Verlag, Berlin, 2006.
24. M. Böhnke and B. R. Masters, "Long Term Contact Lens Wear Induces a Corneal Degeneration with Micro-dot Deposits in the Corneal Stroma," *Ophthalmology* **104**:1887–1896 (1997).
25. J. Porter, H. Queener, J. Lin, K. Thorn, and A. Awwal, (eds.), *Adaptive Optics for Vision Science*, Wiley-Interscience, Hoboken, NJ, 2006.
26. B. R. Masters, "Correlation of Histology and Linear and Nonlinear Microscopy of the Living Human Cornea." *Journal of Biophotonics* **2**:127–139 (2009).

This page intentionally left blank.

DO NOT DUPLICATE

DIAGNOSTIC USE OF OPTICAL COHERENCE TOMOGRAPHY IN THE EYE

Johannes F. de Boer

*Department of Physics
VU University, Amsterdam
Rotterdam Ophthalmic Institute
Rotterdam, The Netherlands*

18.1 GLOSSARY

A-line. Depth scan along the beam direction in an OCT measurement.

Axial resolution. Resolution along the beam direction in an OCT measurement.

Birefringence. Optical property of a material where the refractive index of the material depends on the polarization state of light.

Bruchs membrane. The innermost layer of the choroid, connected to the retinal pigmented epithelium.

Capillary. The smallest blood vessels, approximately 5 to 10 μm in diameter.

Coherence length. The length along the beam path over which light is correlated.

Diattenuation or dichroism. Optical property of a material where the light absorption of the material depends on the polarization state of light.

Fourier transform. Operation that transforms one function of a variable into another, where the new function describes which frequencies are present in the original function.

Fovea. The area in the back of the inside of the eye located in the center of the macula, responsible for sharp central vision.

Choroid. The vascular layer of the eye lying between the retina and the sclera.

Jones matrix formalism. In the Jones matrix formalism the electromagnetic wave before and after interaction is described by Jones vectors with two complex-valued entries, and their transformation is described by the complex-valued 2×2 Jones matrix.

Michelson interferometer. Most common configuration for optical interferometry where the source light is split by a beam splitter into a sample and a reference arm. Reflected light is combined in the detection arm.

Mueller matrix formalism. In the Mueller matrix formalism the polarization properties of light are described by Stokes vectors with four real-valued entries, and their transformation is described by the real-valued 4×4 Mueller matrix.

Optical circulator. An optical circulator is a three-port device that allows light to travel in only one direction—from port 1 to port 2, then from port 2 to port 3.

Optical coherence tomography (OCT). Optical interferometric technique to create cross-sectional images of scattering media.

Optical frequency domain imaging (OFDI) or swept source OCT (SS-OCT). OCT configuration in which the reference arm length is stationary and the interference in the detection arm is spectrally resolved by scanning a narrow band source over a wavelength range.

Optic nerve head. The circular area in the back of the inside of the eye where the optic nerve connects to the retina.

Phase retardation. The retardation in phase that originates from the difference in refractive index for polarized light in a birefringent material.

Poincaré sphere. Graphical representation of the Stokes vector on a sphere.

Relative intensity noise. Noise associated with the instability in the power level of a laser or light source.

Retinal pigmented epithelium. The pigmented cell layer between the choroid and the photoreceptors.

Semiconductor optical amplifier. Semiconductor device that amplifies light by several orders of magnitude.

Shot noise. Fundamental noise generated by statistical fluctuations in a stream of particles (photons, electrons).

Signal to noise ratio. Ratio of a signal power to the noise power corrupting the signal.

Spectral or Fourier domain OCT. OCT configuration in which the reference arm length is stationary and the interference in the detection arm is spectrally resolved.

Stokes vector. Four element vector that describes the polarization properties of light.

Temporal coherence. The feature that light at different time points exhibit correlation.

Time domain OCT. OCT configuration in which the reference arm length is mechanically scanned.

Unitary matrix. A unitary matrix is an n by n complex matrix U satisfying the condition $UU^* = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$, where $\begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}$ is the identity matrix and U^* is the conjugate transpose (also called the Hermitian adjoint) of U .

Vitreous humor. The clear gel that fills the space between the lens and the retina.

White noise. Noise that has a flat frequency distribution (all frequencies are equally present in the noise).

18.2 INTRODUCTION

Optical coherence tomography (OCT) is a low coherence interferometric method for imaging of biological tissue.^{1,2} OCT creates cross-sectional images of tissue by measuring the location (lateral and depth) where light is backscattered. OCT has become a very powerful imaging tool in ophthalmology, since it provides high resolution cross-sectional images of the retina and the cornea that approach the resolution of histology. The development of objective criteria to characterize ophthalmic diseases based on OCT images is an active area of research.

OCT is analogous to ultrasound imaging, where a sound pulse is sent into tissue and the time delay is measured of reflected sound waves. The time delay is converted to a distance using the propagation speed of sound in tissue. The speed of light (3×10^8 m/s) is 5 orders of magnitude higher

than the speed of sound (1.5×10^3 m/s), and therefore it is not possible to measure directly the small time delay of a reflected optical pulse. This time delay can be measured however by interfering the light reflected from the tissue with a reference in a Michelson type interferometer. This principle will be explained in more detail further.

For more than a decade after its inception, the dominant implementation of OCT has been time domain OCT (TD-OCT), in which the length of a reference arm in a Michelson type interferometer is rapidly scanned. TD-OCT has been used extensively in ophthalmology, where cross-sectional OCT images of the retina have provided useful information regarding the presence or progression of specific ocular diseases.³ The acquisition rate of clinical and preclinical TD-OCT systems is limited by sensitivity and the maximum permissible incident power on the eye⁴ to about 400 depth profiles/s, preventing comprehensive screening of large retinal areas. Over the past 5 years, the dominant implementation of OCT in ophthalmology has become spectral or Fourier domain OCT (SD/FD-OCT). As will be discussed, SD/FD-OCT has a significant sensitivity advantage over TD-OCT. The sensitivity improvement of SD-OCT allows for dramatically increased acquisition speeds by approximately a factor of 100 without compromising image quality. Three-dimensional data sets can therefore be rapidly acquired, opening the possibility of comprehensive screening. The first SD/FD-OCT implementation was reported in 1995.⁵ In SD-OCT the reference arm is kept stationary, and the depth information is obtained by a Fourier transform of the spectrally resolved interference fringes measured with a spectrometer in the detection arm of a Michelson interferometer.

An alternative method to SD/FD-OCT that provides the same sensitivity advantage is optical frequency domain imaging (OFDI) or swept source OCT.⁶⁻⁹ In OFDI, the spectrally resolved interference fringes are not measured by a spectrometer in the detection arm of a Michelson interferometer, but by rapidly tuning a narrowband source over a wavelength range. A detector in the reference arm measures the spectrally resolved interference as a function of time.

OCT images provide structural information based on the scattering properties of the tissue under investigation. OCT can be enhanced by functional extensions such as Doppler sensitivity to detect flow, and polarization sensitivity to detect birefringent structures (structures for which the refractive index depends on the polarization state of light). In this chapter, the aforementioned techniques will be discussed and explained in more detail.

18.3 PRINCIPLE OF OCT: TIME DOMAIN OCT

Optical coherence tomography is a very sensitive interferometric method to measure the location of structures that reflect light with an axial resolution of 1 to 10 μm . OCT exploits the low temporal coherence properties of a broadband light source. For a low temporal coherence light source, electromagnetic waves will not interfere when the time delay between two beams is larger than the coherence time. The product of the coherence time and the speed of light gives the coherence length. This property can be used to measure the path length difference between two arms of a Michelson interferometer. When the path length difference is larger than the coherence length, light will not interfere. The path length difference for which the two arms match in length within the coherence length of the source produces an interference pattern. Figure 1 shows the principle of an OCT system with a scanning mirror in the reference arm.

In the sample arm is a structure with three reflecting layers. As the reference arm mirror is translated, the detector in the detection arm measures an interference pattern each time the reference arm path length matches the path length to one of the reflecting layers in the sample arm. In OCT the envelope of the interference pattern is detected and a single scan of the reference arm mirror provides a single axial (along the beam direction) reflectivity profile of the sample. A single depth profile is also called an A-line. Repeating many depth scans (A-lines) at lateral locations creates a cross-sectional image of, for example, the retina, as shown in Fig. 2.

In a more mathematical description the intensity at the detector is given by

$$\langle I(\Delta z) \rangle = \langle I_r \rangle + \langle I_s \rangle + \langle E_r^* E_s(\Delta z) \rangle + \langle E_r E_s^*(\Delta z) \rangle \quad (1)$$

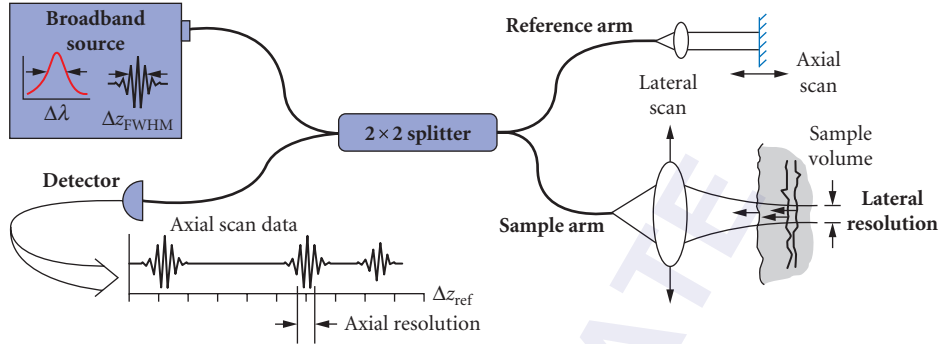


FIGURE 1 Basic OCT configuration. Light from a broadband source is split into a sample and a reference arm. The reference arm mirror scans in depth, creating interference fringes when the reference arm path length matches the sample arm path length to one of the three reflecting layers in the sample.

where $E_{r,s}$, $E_{r,s}^*$ are the electric field component and its complex conjugate, Δz is the path length difference between the two arms of the interferometer, and subscripts r and s denote fields reflected from the reference and sample arm, respectively. The angular brackets denote time averaging. The first two terms on the right-hand side of Eq. (1) give the intensity reflected from the reference and sample arm, respectively.

The last two terms of Eq. (1) correspond to the interference between reference and sample arm light. In our analysis, the electric field amplitude is represented by a complex analytic function,¹⁰ $E(z)$, with

$$E(z) = \int \tilde{e}(k) \exp(-ikz) dk \quad (2)$$

where $\tilde{e}(k)$ is the field amplitude as a function of free space wave number $k = 2\pi/\lambda$, with

$$\tilde{e}(k) = 0 \quad \text{if} \quad k < 0 \quad (3)$$

From the Wiener-Khinchine theorem, it follows,

$$\langle \tilde{e}(k) \tilde{e}(k') \rangle = S(k) \delta(k - k') \quad (4)$$

which defines $\tilde{e}(k)$ in terms of the source power spectral density $S(k)$. Using the Wiener-Khinchine theorem [Eq. (4)], the interference term in Eq. (1) is given by

$$I_{r,s}(z, \Delta z) = E_r^* E_s + E_r E_s^* \propto \sqrt{R(z)} \int \cos(2k\Delta z) S(k) dk \quad (5)$$



FIGURE 2 OCT image of a human retina. Arrow indicates the position of a blood vessel. (Reproduced from Ref. 11 with permission from the Optical Society of America.)

with z the depth in the tissue and Δz the optical path length difference between sample and reference arms, $\Delta z = z_r - z_s$. We will assume a Gaussian power spectral density for the source

$$S(k) \propto \exp\left[-\left(\frac{k-k_0}{\kappa}\right)^2\right] \quad (6)$$

with the FWHM spectral bandwidth of the source given by $\kappa 2\sqrt{\ln 2}$. The integration over k in Eq. (5) can now easily be performed

$$I_{r,s}(z, \Delta z) \propto \sqrt{R(z)} \cos(2k_0 \Delta z) \exp[-(\Delta z / \Delta l)^2] \quad (7)$$

with the FWHM of the interference fringes envelope given $\Delta l 2\sqrt{\ln 2}$

$$\Delta l = \frac{1}{\kappa} = \frac{\lambda_0^2 \sqrt{\ln 2}}{\pi \Delta \lambda} \quad (8)$$

and $\Delta \lambda$ the spectral FWHM of the source in wavelength. Equation (7) is an important result; it describes the interference fringe intensity that is observed for a reflective structure in the sample arm when the reference arm is scanned. The right-hand side of Eq. (7) consists of the reflected amplitude at depth z ($\sqrt{R(z)}$), the Doppler shift or carrier frequency generated by the variation of the optical path length difference between sample and reference arm ($\cos(2k_0 \Delta z)$), and the interference fringes envelope ($\exp(-(\Delta z / \Delta l)^2)$), which describes the width of the interference pattern, respectively. The depth resolution is determined by the width Δl of the interference pattern, and Eq. (8) shows the relation between axial or depth resolution and the spectral width of the source in OCT. The broader the source (larger $\Delta \lambda$), the smaller the coherence length Δl , that is, the better the axial resolution. The availability of sources with very broad spectral widths over the past decade has increased the axial (depth) resolution of OCT from $15 \mu\text{m}$ to approximately 2 to $3 \mu\text{m}$.¹² The lateral resolution in OCT is determined by the lateral size of the focused spot, in analogy to confocal microscopy.

18.4 PRINCIPLE OF OCT: SPECTRAL DOMAIN OCT

Spectral domain OCT is an alternative implementation of the time domain OCT technology. The first SD/FD-OCT implementation was reported in 1995.⁵ Figure 3 shows the basic configuration of an SD/FD-OCT system.

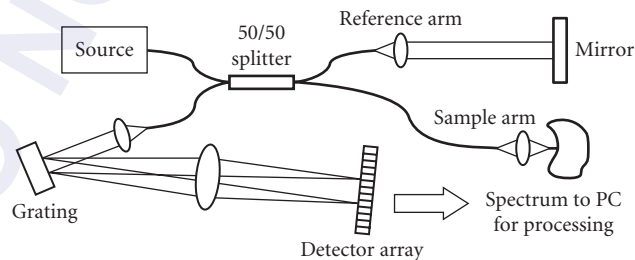


FIGURE 3 Basic configuration of an SD/FD-OCT system. Light from a broadband light source is split into a sample and a reference arm. The mirror in the reference arm is stationary. After reflection, the sample and reference arm light is recombined and interferes. In the detection arm, the light is dispersed by a grating and each wavelength is focused on an element in a detector array. (Reproduced from Ref. 13 with permission from the Biophysical Society.)

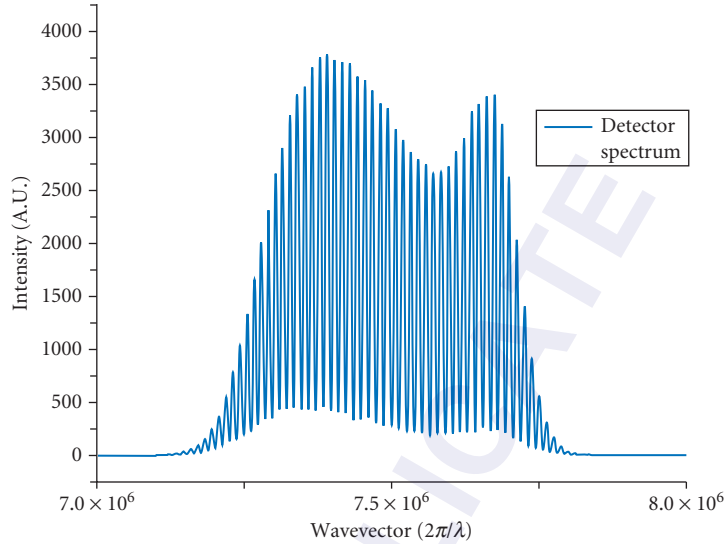


FIGURE 4 Spectrally resolved interference pattern as detected by the spectrometer in the detection arm of an SD-OCT system with a single reflector in the sample arm. (Reproduced from Ref. 13 with permission from the Biophysical Society.)

The source light is split into a reference and sample arm. At recombination in the detection arm, the light interferes. In the detection arm, the light is spectrally dispersed by a grating, and each wavelength is focused on an element in a linear detection array by a lens. The path length difference between the sample and reference arm is determined by the spectrally resolved interference detected by the linear detection array. In Fig. 4, an example of such a spectrally resolved interference pattern is shown for a single reflector in the sample arm. The power reflected from the sample and reference arm was approximately equal, giving nearly full constructive and destructive interference fringes.

The interference pattern shown in Fig. 4 can be understood as follows: for wavelengths for which the path length difference between sample and reference arm is exactly a multiple of the wavelength we observe constructive interference, and for wavelengths for which the path length difference between sample and reference arm is $(N + 0.5)\lambda$ we observe destructive interference. Mathematically, the interference spectrum is described by

$$I(k) = I_r(k) + I_s(k) + 2\sqrt{I_s(k)I_r(k)} \sum_n \alpha_n \cos(kz_n) \quad (9)$$

where $I_r(k)$ and $I_s(k)$ are the wavelength-dependent intensities reflected from reference and sample arms, respectively, and k is the wave number. The last term on the right-hand side of Eq. (9) represents the interference between light returning from the reference and sample arms. α_n is the square root of the sample reflectivity at depth z_n , and z_n is the path length difference between sample and reference arm. The summation sums over all reflective structures at location z_n in the sample arm. Equation (9) shows that the path length difference is encoded on the spectrum by the modulation term $\cos(kz_n)$. The modulation frequency as a function of wave vector is determined by the path length difference z_n . The larger the path length difference, the higher the modulation frequency. This also illustrates that the maximum path length difference (depth range) that can be measured is determined by the ability to measure high modulation frequencies as a function of wave vector. The spectral resolution of the spectrometer plays a crucial role in measuring these high modulation frequencies. In SD-OCT, the depth range z is inversely proportional to the spectral resolution $\delta\lambda$ of the spectrometer and given by¹⁴

$$z = \lambda_0^2 / 4n\delta\lambda \quad (10)$$

A depth profile or depth scan is obtained from the spectrum through a Fourier transformation [see Eq. (18) later in this chapter].

SD-OCT provides a significant advantage in sensitivity or equivalently signal to noise ratio (SNR), that despite reports as early as 1997^{14,15} has taken about half a decade to be recognized fully by the OCT community.^{7,16,17} The first demonstration of SD-OCT for in vivo retinal imaging in 2002¹⁸ was followed by a full realization of the sensitivity advantage by video rate in vivo retinal imaging,¹⁰ including high speed three-dimensional volumetric imaging,¹⁹ ultra high resolution video rate imaging,^{20,21} and Doppler blood flow determination in the human retina.^{22,23} The superior sensitivity of SD-OCT, combined with the lack of need for a fast mechanical scanning mechanism, has opened up the possibility of much faster scanning without loss of image quality and provided a paradigm shift from point sampling to volumetric mapping of biological tissue in vivo. The technology has been particularly promising in ophthalmology.^{24,25} Later in this chapter a detailed analysis of the sensitivity advantage of SD-OCT is presented.

18.5 PRINCIPLE OF OCT: OPTICAL FREQUENCY DOMAIN IMAGING

In SD-OCT the spectrally resolved interference fringes are detected by a spectrometer in the detection arm. The same information can be obtained by scanning a narrowband source over a large wavelength range, and detecting the light with a single detector. In this case, the wavelength is encoded as a function of time. This technique, optical frequency domain imaging (OFDI) or swept source OCT (SS-OCT), gives the same sensitivity advantage as SD-OCT.^{7,8} Figure 5 gives a basic

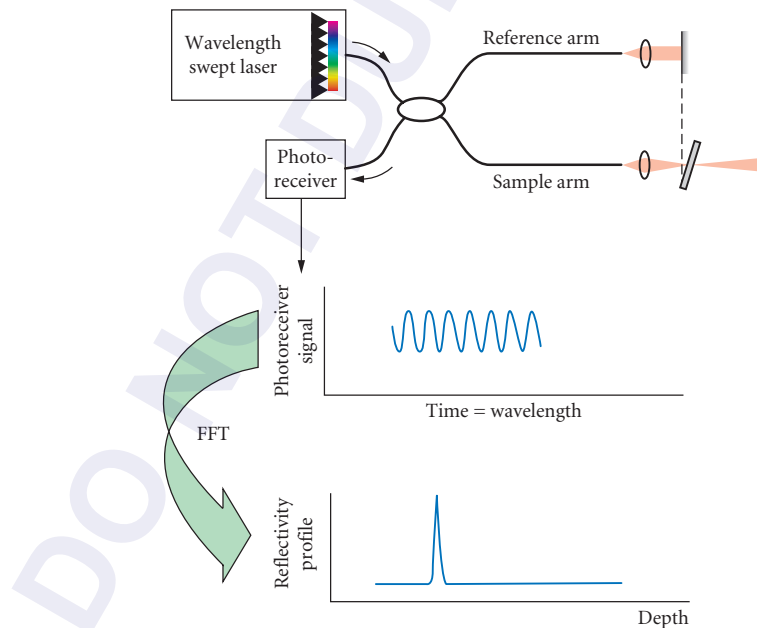


FIGURE 5 Basic OFDI configuration. A narrowband light source rapidly tunes over a wavelength range. Sample and reference arm light interferes in the detection arm. The light is detected by the photoreceiver as a function of time, where time is equivalent to wavelength or wave vector, giving the spectrally resolved interference fringes. A Fourier transform gives a reflectivity profile of the sample.

configuration of an OFDI or SS-OCT system. A narrowband laser source rapidly scans over a large wavelength range. The photo receiver detects the interference between sample and reference arm light as a function of time, where time is equivalent to wavelength or wave vector. A single wavelength scan gives the wavelength-resolved interference, which after a Fourier transform provides a single reflectivity profile of the sample.

OFDI has been demonstrated at 1300 nm in a variety of tissues. In ophthalmology, OFDI has been demonstrated at a wavelength of 1050 nm²⁶ and 840 nm.²⁷ Particularly 1050 nm is an interesting wavelength for ophthalmology. It is well known that longer wavelengths penetrate deeper into biological tissue. This suggests that 1300 nm would provide deeper penetration in retinal tissue; however, the water absorption of 1300 nm is too large to image the retina through the vitreous humor. At 1050 nm there is a window where the water absorption is lower than at 1300 nm. The better penetration of 1050 nm over 800 nm into the retinal layers and the choroid is demonstrated in Fig. 6 and Ref. 28 in a comparison between an OFDI system at 1050 nm and an SD-OCT system at 840 nm. The better visualization below of the retinal pigmented epithelium layer and into the choroid could be particularly useful for the diagnosis of age-related macular degeneration (AMD).²⁹

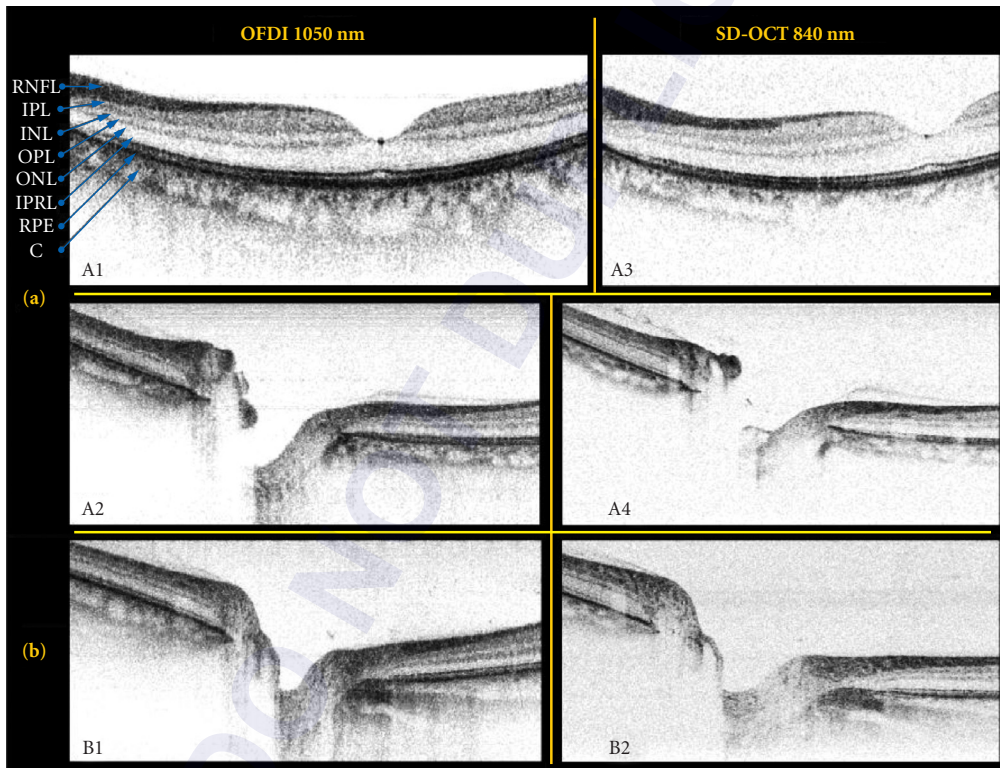


FIGURE 6 Comparison of two imaging systems (OFDI at 1050 nm and SD-OCT at 840 nm). **A1** and **A2**: OFDI images at fovea and optic nerve head, respectively, from volunteer A, 36-year-old Asian male. **A3** and **A4**: SD-OCT images from the same volunteer at similar tissue locations. **B1** and **B2**: OFDI and SD-OCT images, respectively, obtained from volunteer B, 41-year-old Caucasian male. OFDI images exhibit considerably deeper penetration in tissue than SD-OCT images in all the data sets. The OFDI image (**A1**) shows the anatomical layered structure: RNFL, retinal nerve fiber layer; IPL, inner plexiform layer; INL, inner nuclear layer; OPL, outer plexiform layer; ONL, outer nuclear layer; IPRL, interface between the inner and outer segments of the photoreceptor layer; RPE, retinal pigmented epithelium; and C, choriocapillaris and choroid. (Reproduced from Ref. 26 with permission from the Optical Society of America.)

18.6 SD-OCT VERSUS OFDI

The best choice for a particular OCT technology in ophthalmology depends on a number of factors. Currently, SD-OCT is the preferred implementation at wavelengths around 800 nm, due to the availability of high performance silicon array detectors and the difficulty to create robust rapid tunable sources in this wavelength range. At wavelengths of 1050 and 1300 nm, OFDI is the preferred implementation due to the availability of a key component for a rapid tunable source at these wavelengths, a high performance semiconductor optical amplifier (SOA), and the relatively poor performance of InGaAs detector arrays. Other considerations that play a role are the relative ease with which to construct a broadband light source compared to a broadband rapid tunable source, favoring SD-OCT. On the other hand, SD-OCT suffers from sensitivity decay as a function of depth, limiting the effective depth range to about 1 mm, that will be detailed further below. Experimental OFDI systems have in general demonstrated an effective depth range that is 2 to 4 times better than SD-OCT. In cases where a large effective depth range is required, for example, imaging the cornea, or imaging patients with deep optic nerve head cupping, OFDI has an advantage.

18.7 SENSITIVITY ADVANTAGE OF SD-OCT OVER TD-OCT

In the standard time-domain (TD) implementation of OCT, the position of the reference mirror in the interferometer is rapidly scanned by mechanical means in order to obtain a depth profile (A-line) within a sample. In SD/FD-OCT, no mechanical scanning of the reference arm is required. Instead, the cross-spectral density at the detection arm of the interferometer is measured by means of a spectrometer.^{5,14,30} Although this method has long been proposed and demonstrated, only recently have there been efforts to explicitly show that SD-OCT can produce a better detection sensitivity than the time domain method.^{7,16,17} SD-OCT does not require scanning of the reference arm length and therefore has the potential for faster image acquisition rates. High speed line scan cameras in the spectrometer with an integration time as small as 34 and even 16 μ s permit the acquisition of A-lines at a rate of 30 to 60 kHz. Recent work has experimentally demonstrated a 148-fold (21.7 dB) sensitivity improvement of SD-OCT¹⁰ and a near shot noise limited performance. Shot noise is a fundamental noise source that cannot be improved upon.

In essence the SNR advantage of SD-OCT over TD-OCT is based on the significant reduction of noise obtained by replacing the single-element detector with a multielement array detector. In a TD-OCT system, each wavelength is uniquely encoded as a frequency, and shot noise has a white noise characteristic. In a single detector TD-OCT system, the shot noise generated by the power density at one particular wavelength is present at all frequencies, and therefore adversely affects the SNR at all other wavelengths. By spectrally dispersing each wavelength to a separate detector, the cross shot noise term is eliminated in both hybrid and fully parallel SD-OCT systems.¹⁷

18.8 NOISE ANALYSIS OF SD-OCT USING CHARGE COUPLED DEVICES (CCDs)

The core of an SD-OCT system is the spectrometer in the detection arm. In general, light detection in the spectrometer is achieved by a charge coupled device (CCD) like a line array. Alternative implementations have been suggested that require individual processing of pixel charges.³¹ Here we will treat the noise analysis based on a simple integrating CCD. To facilitate the noise analysis in the

case of detection by CCDs, the signal and noise terms will be expressed in charge squared (e^2). In SD-OCT the signal S_{SD} given by,^{16,17}

$$S_{SD} = \frac{\eta^2 e^2 P_{ref} P_{sample} \tau_i^2}{E_v^2} (e^2) \quad (11)$$

with η the quantum efficiency of the detector, e the electron charge, P_{ref} and P_{sample} , respectively, the reference arm and sample arm power at the detection arm fiber tip, τ_i the integration time, and E_v the photon energy.

In the noise analysis, it is generally assumed that the reference arm power is much larger than the sample arm power, and the shot noise and relative intensity noise (RIN) contributions of the sample arm power to the noise are neglected.³² The read-out and dark noise, shot noise, and RIN contributions to the overall noise in electrons squared per read-out cycle are then given by, respectively,^{16,17}

$$\sigma_{noise}^2 = \sigma_{r+d}^2 + \frac{\eta e^2 P_{ref} \tau_i}{E_v} + \left(\frac{\eta e P_{ref}}{E_v} \right)^2 \tau_i \tau_{coh} (e^2) \quad (12)$$

with σ_{r+d}^2 the sum of read-out noise and dark noise, and $\tau_{coh} = \sqrt{2 \ln 2 / \pi} \lambda_0^2 / c \delta \lambda$ the coherence time, with c the speed of light.³³ A close look at Eq. (12) shows that the noise consists of a constant term (σ_{r+d}^2), a term that linearly increases with the reference arm power (shot noise) and a term that quadratically increases with the reference arm power (RIN). The signal [Eq. (11)] increases linearly with the reference arm power. The signal to noise ratio (SNR) determines the performance of the system. Since the shot noise is the limiting noise term that cannot be improved upon, the optimal signal to noise performance is achieved when shot noise dominates both read-out noise and RIN.³² When shot noise dominates the noise, both the signal [Eq. (11)] and the noise [Eq. (12)] increase linearly with the reference arm power and the SNR is constant. Shot noise dominates read-out noise and dark noise when $\eta e^2 P_{ref} \tau_i / \sigma_{r+d}^2 E_v > 1$, and shot noise dominates RIN when their ratio is larger than 1, that is, $E_v / \eta P_{ref} \tau_{coh} > 1$. The optimal reference arm power is found when read-out noise and dark noise are equal to the RIN.³⁴

$$\sigma_{r+d}^2 = \left(\frac{\eta e P_{ref}}{E_v} \right)^2 \tau_i \tau_{coh} \quad (13)$$

Thus for a system to operate close to shot-noise-limited performance, shot noise should dominate thermal and RIN noise at the optimal reference arm power

$$P_{ref} = \frac{\sigma_{r+d} E_v}{\eta e \sqrt{\tau_i \tau_{coh}}} \quad (14)$$

At this optimal reference arm power, the inequalities describing shot noise dominance over read-out noise and RIN, respectively, reduce to the same equation

$$\frac{e \sqrt{\tau_i}}{\sigma_{rd} \sqrt{\tau_{coh}}} > 1 \quad (15)$$

In general, one would like to choose the integration time τ_i as short as possible, and the coherence time τ_{coh} as long as possible. The coherence time is inversely related to the spectral resolution of the spectrometer which in turn relates linearly to the maximum depth range of the system. In conclusion, the parameter that most determines the system performance is the read-out and dark noise of the detector σ_{rd} .

18.9 SIGNAL TO NOISE RATIO AND AUTOCORRELATION NOISE

In the shot noise limit, the following expressions are found for the SNR in the time domain³² and the spectral domain,¹⁷ where the expression for the spectral domain can also be derived by the ratio of Eqs. (11) and (12)

$$\text{SNR}_{\text{TD}} = \frac{\eta P_{\text{sample}}}{E_v \text{BW}} \quad \text{SNR}_{\text{SD}} = \frac{\eta P_{\text{sample}} \tau_i}{E_v} \quad (16)$$

where η is the spectrometer efficiency, P_{sample} is the sample arm power returning to the detection arm, BW is the electronic detection bandwidth in a time domain system, τ_i is the detector integration time, and E_v is the photon energy. The electronic signal bandwidth BW is centered at the carrier frequency given by the Doppler shift [see Eq. (7)] and the bandwidth is proportional to the spectral width of the source and the velocity of the reference arm mirror. A detailed numerical comparison between the SNR in the time and spectral domain showed a more than 2 orders better SNR in the spectral domain.¹⁰ Unlike the SNR in the time domain, Eq. (16) demonstrates that in the shot noise limit, SNR_{SD} is independent of the spectral width of the source. This implies that the axial resolution can be increased at no penalty to the SNR, provided that the full spectral width of the source can be imaged onto an array detector. However, this result should be interpreted with some care. The sample arm power returning to the detection arm is assumed to come from a single reflecting surface. In tissue, however, the reflected power comes from multiple structures along a depth profile. The SNR for a particular position along the depth profile is given on average by the total power reflected by all structures within the coherence length of the source. As the resolution increases (the coherence length decreases), the total reflected power within the coherence length decreases. As a consequence, the SNR at a particular position along the depth profile will reduce as the resolution increases by increasing the source optical bandwidth.

In SD-OCT the depth information is obtained by a Fourier transform of the spectrally resolved interference fringes. The detected interference signal at the spectrometer may be expressed as¹⁴

$$I(k) = I_r(k) + I_s(k) + 2\sqrt{I_s(k)I_r(k)} \sum_n \alpha_n \cos(kz_n) \quad (17)$$

where $I_r(k)$ and $I_s(k)$ are the wavelength-dependent intensities reflected from reference and sample arms, respectively, and k is the wave number. The last term on the right-hand side of Eq. (17) represents the interference between light returning from reference and sample arms. α_n is the square root of the sample reflectivity at depth z_n . Depth information is retrieved by performing an inverse Fourier transform of Eq. (16), yielding the following convolution¹⁴

$$\left| \text{FT}^{-1}[I(k)] \right|^2 = \Gamma^2(z) \otimes \left\{ \delta(0) + \sum_n \alpha_n^2 \delta(z - z_n) + \sum_n \alpha_n^2 \delta(z + z_n) + O\left[\frac{I_s^2}{I_r^2}\right] \right\} \quad (18)$$

with $\Gamma(z)$ representing the envelope of the coherence function. The first term in the braces on the right-hand side describes the autocorrelation signal from the reference arm and has magnitude unity. The second and third terms are due to interference between light returning from reference and sample arms and form two images, where each has magnitude on the order of I_s/I_r . These two terms provide mirror images, where one is retained. The final term, with magnitude on the order of I_s^2/I_r^2 , describes autocorrelation noise due to interference within the sample arm.^{14,18} I_s and I_r represent the total intensity reflected from sample and reference arm, respectively. Sample autocorrelation noise is generated by the interference of light reflected at different depth locations in the sample. Equation (18) indicates that the relative contribution of sample autocorrelation noise can be reduced by increasing the reference arm power with respect to the signal. Decreasing the detector

integration time permits an increase in the reference arm power without saturating the detector, decreasing the ratio I_s^2/I_r^2 and consequently reducing the contribution of autocorrelation noise in ultrahigh-speed SD-OCT.

18.10 SHOT-NOISE-LIMITED DETECTION

In this section we take a closer look at the different noise components described in Eq. (12) for an actual system. In an SD-OCT system the sample arm was blocked, and only the reference arm light was detected by the spectrometer. One thousand spectra were recorded at an acquisition rate of 29.3 kHz. The read-out and shot noise are shown in Fig. 7.

The noise was determined by calculating the variance at each camera pixel over 1000 consecutive spectra. Dark noise measurements were taken with the source light off. Only light returning from the reference arm was used to measure the shot noise and RIN in the system. The shot noise and RIN are given by the second and third term on the right-hand side of Eq. (12), expressed in number of electrons squared. Taking into account only the shot and dark + read-out noise, the measured variance is proportional to

$$\sigma^2(\lambda) \sim \frac{P_{\text{ref}}(\lambda)}{E_v} + \sigma_{r+d}^2 \quad (19)$$

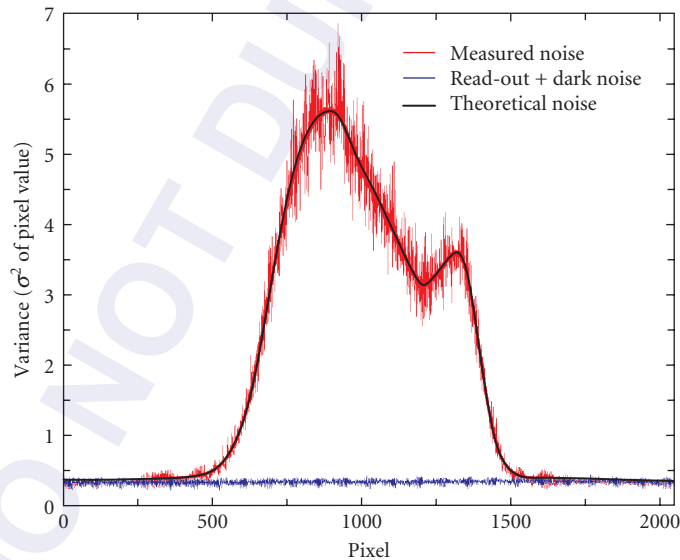


FIGURE 7 Noise analysis of an SD-OCT system. The variance (vertical axis) was calculated as a function of wavelength (pixel number, horizontal axis) for 1000 consecutive spectra. Both the read-out + dark noise (blue curve, no illumination) and the shot noise (red curve, illumination by reference arm only) were determined. The theoretical shot noise curve was fit using Eq. (19) to the measured noise. The excellent fit of the theoretical shot noise curve to the measured noise demonstrates that RIN was significantly smaller than shot noise and did not contribute. (Reproduced from Ref. 19 with permission from the Optical Society of America.)

The first term on the right-hand side of Eq. (19) is the shot noise contribution which is linearly proportional to the reference arm power, and the second term is the dark and read-out contribution to the noise. Equation (19) was fit to the measurements, limiting the fit to the central 700 pixels. Relative intensity noise (RIN) was not dominant in this setup, as experimentally demonstrated by the excellent fit of only the shot noise and dark and read-out noise to the measured noise in Fig. 7, and theoretically since the maximum power per pixel (4.6 nW) at a 34.1 μs integration time does not meet the criteria for RIN dominated noise.¹⁰ This demonstrates shot-noise-limited performance of an SD-OCT system.

18.11 DEPTH DEPENDENT SENSITIVITY

In SD-OCT, signal sensitivity is strongly dependent on depth within an image. To characterize system sensitivity as a function of ranging depth, 1000 A-lines were acquired at an acquisition speed of 34.1 $\mu\text{s}/\text{A-line}$ for 9 different positions of a weak reflector in the sample arm. The reflected sample arm power was 1.18 nW for all reflector positions. The noise floor decayed by 5 dB between a depth of 500 μm and 2 mm, and the peak signal dropped by 21.7 dB over the first 2 mm. Due to fixed-pattern noise, the true noise floor could not be determined between 0 and 500 μm . After zero-padding before the FFT to correct wavelength-mapping errors from pixel to wave vector, a 16.7 dB loss in peak signal was noted across the first 2 mm, whereas the noise level dropped by only 0.4 dB between 500 μm (35.1 dB) and 2 mm (34.7 dB) (Fig. 8). The zero-padding method produced a nearly constant noise level and improved the signal by more than 5 dB at the greatest depths in the scan. Although zero-padding did not change the local SNR, this method eliminated the shoulders that are present at larger scan depths.⁷ The decay in both the signal and the noise level across the entire scan

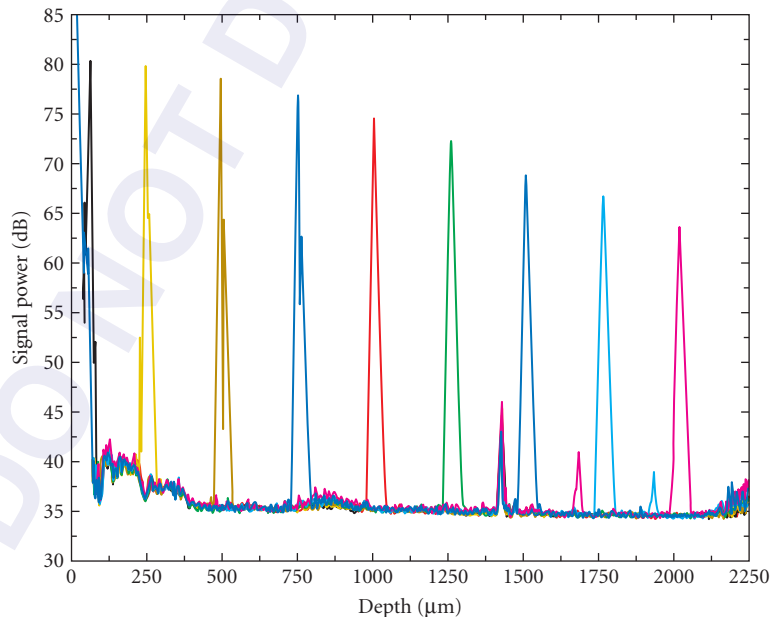


FIGURE 8 The depth dependent loss in signal sensitivity from a weak reflector. The signal decayed 16.7 dB between 0 and 2 mm. The peaks at 1.4, 1.6, and 1.85 mm are fixed-pattern noise. (Reproduced from Ref. 19 with permission from the Optical Society of America.)

length of 2.4 mm have been theorized to amount to 4 dB as a result of the finite pixel width.¹⁶ As demonstrated by the experimental data, the noise level decayed by less than 4 dB over the entire scan length, which we attribute to the statistical independence of the shot noise between neighboring pixels of the array. Thus, the finite pixel width does not introduce a decay of the noise level.

The finite spectrometer resolution introduces a sensitivity decay³⁴ similar to that introduced by the finite pixel size.¹⁶ Convolution of the finite pixel size with the Gaussian spectral resolution yields the following expression for the sensitivity reduction, R , as a function of imaging depth, z ³⁴

$$R(z) = \frac{\sin^2(\pi z/2d)}{(\pi z/2d)^2} \exp\left[-\frac{\pi^2 \omega^2}{8 \ln 2} \left(\frac{z}{d}\right)^2\right] \quad (20)$$

where d is the maximum scan depth and ω is the ratio of the spectral resolution to the sampling interval. Equation (20) was fit to the signal decay data presented in Fig. 8 with ω as a free parameter, and the result is shown in Fig. 9. Due to its proximity to the autocorrelation peak, the first data point was not included in the fit. The value for ω obtained from the fit was 1.85, demonstrating that the working spectral resolution was 0.139 nm.

The SNR was determined by the ratio of the peak at 250 μm (79.8 dB) and the noise level. Due to the fixed-pattern noise at 250 μm , the noise level was determined to be 35.2 dB by extrapolation of the linear region between 0.5 and 2 mm. The resulting SNR of 44.6 dB for 1.18 nW returning to the detection arm was 2.2 dB below the theoretical value given by Eq. (16) of 46.8 dB, for an integration time of 34.1 μs , a central wavelength of 840 nm, and a spectrometer efficiency of 28 percent. With 600 μW of power incident on an ideal reflector in the sample arm, the measured power returning to the detection arm was 284 μW . The sum of the SNR at 1.18 nW (44.6 dB) and the 10 log ratio of maximum (284 μW) over measured (1.18 nW) power (53.8 dB) gives a sensitivity of 98.4 dB.

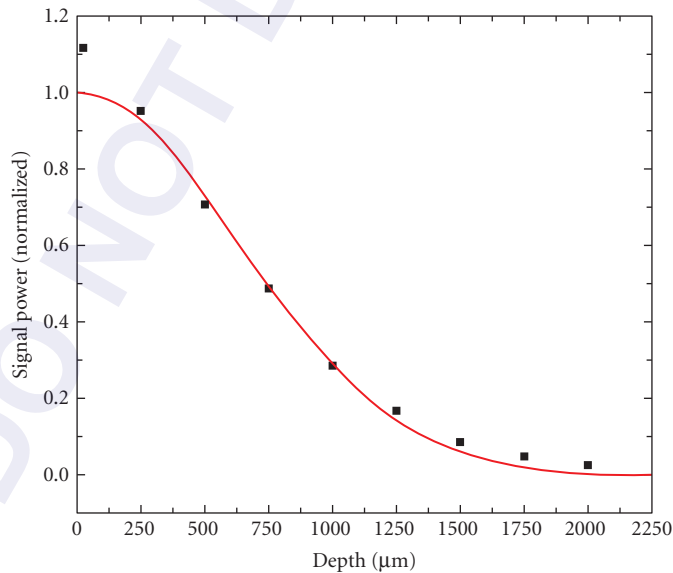


FIGURE 9 Decay of sensitivity across the measurement range. Symbols: Peak intensities of data presented in Fig. 8. Solid line: Fit of Eq. (20) to the data points. (Reproduced from Ref. 19 with permission from the Optical Society of America.)

18.12 MOTION ARTIFACTS AND FRINGE WASHOUT

As OCT utilizes lateral point-scanning, motion of the sample or scanning beam during the measurement causes SNR reduction and image degradation in SD-OCT and OFDI.³⁵ Yun et al. theoretically investigated axial and lateral motion artifacts in continuous wave (CW) SD-OCT and swept-source OFDI, and experimentally demonstrated reduced axial and lateral motion artifacts using a pulsed source and a swept source in endoscopic imaging of biological tissue.^{35,36} Stroboscopic illumination in full field OCT was demonstrated, resulting in reduced motion artifacts for in vivo measurement.³⁷ In ophthalmic applications of SD-OCT, SNR reduction caused by high speed lateral scanning of the beam over the retina may be dominant over axial patient motion. Using pulsed illumination reduces lateral motion artifacts and provides a better SNR for in vivo high-speed human retinal imaging.³⁸

18.13 OFDI AT 1050 NM

An alternative technique to SD-OCT is optical frequency domain imaging (OFDI).⁸ In OFDI, a rapidly tuned laser source is used and the spectrally resolved interference fringes are recorded as a function of time in the detection arm of the interferometer. Published results in healthy volunteers have shown that OFDI has better immunity to sensitivity degradation due to lateral and axial eye motion, and has an effective ranging depth that is 2 to 2.5 times better than SD-OCT (depth dependent sensitivity decay of 6 dB over 2 to 2.5 mm).^{8,26,27}

More importantly, recent research at 1050 nm spectral range has demonstrated a better retinal penetration depth,^{26,28,39} particularly important for detecting retinal abnormalities at or below the retinal pigment epithelium (RPE). A wavelength of 1050 nm has less attenuation from scattering in opaque media, commonly seen in cataract patients.⁴⁰ Although the water absorption at 1050 nm is higher than in the 850 nm region, this is partially compensated by the approximately 3 times higher maximum permissible exposure according to the ANSI standards (1.9 mW at 1050 nm).⁴

Figure 10a depicts a schematic of a laser source for OFDI in a linear cavity configuration.^{8,26} The gain medium was a commercially available, bi-directional semiconductor optical amplifier

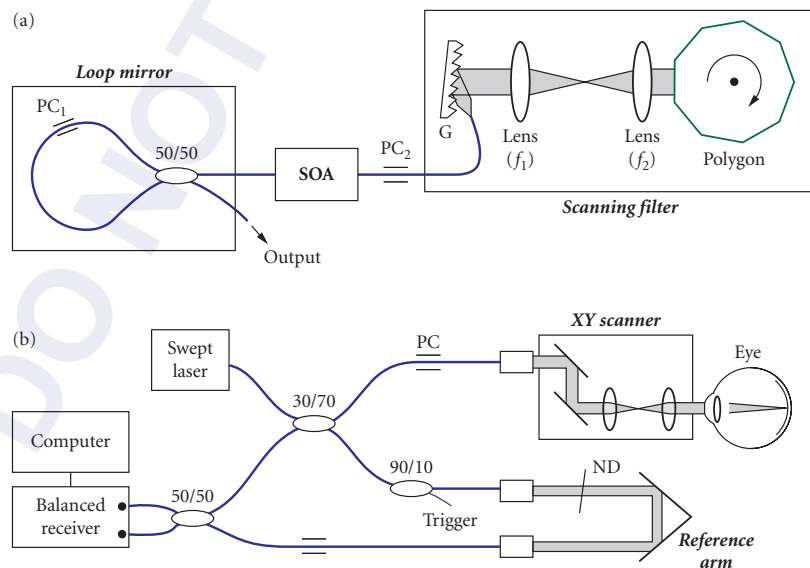


FIGURE 10 Experimental setup: (a) wavelength-swept laser and (b) OFDI system. (Reproduced from Ref. 26 with permission of the Optical Society of America.)

(QPhotonics, Inc., QSOA-1050) driven at an injection current level of 400 mA. One port of the amplifier was coupled to a wavelength-scanning filter⁴¹ that comprises a diffraction grating (1200 lines/mm), a telescope ($f_1 = 100$ mm, $f_2 = 50$ mm), and a polygon mirror scanner (Lincoln Lasers, Inc., 40 facets). The design bandwidth and free spectral range of the filter were approximately 0.1 and 61 nm, respectively. The amplifier's other port was spliced to a loop mirror made of a 50/50 coupler. Sweep repetition rates of up to 36 kHz were possible with 100 percent duty cycle.

Figure 10*b* depicts the complete ophthalmic OFDI system.²⁶ The effective ranging depth was 2.4 mm (depth dependent sensitivity decay of 6 dB over 2.4 mm) due to the finite coherence length of the laser output.

The OFDI system acquired data as the focused sample beam was scanned over an area of 6 (horizontal) by 5.2 mm (vertical) across the macular region in the retina. Each image frame in the three-dimensional volume was constructed from a thousand A-line scans. Given three-dimensional tomographic data of the eye's posterior segment, integrating the pixel values along the entire depth axis readily produces a two-dimensional fundus-type reflectivity image.^{42,43} Figure 11*a* depicts an integrated reflectivity

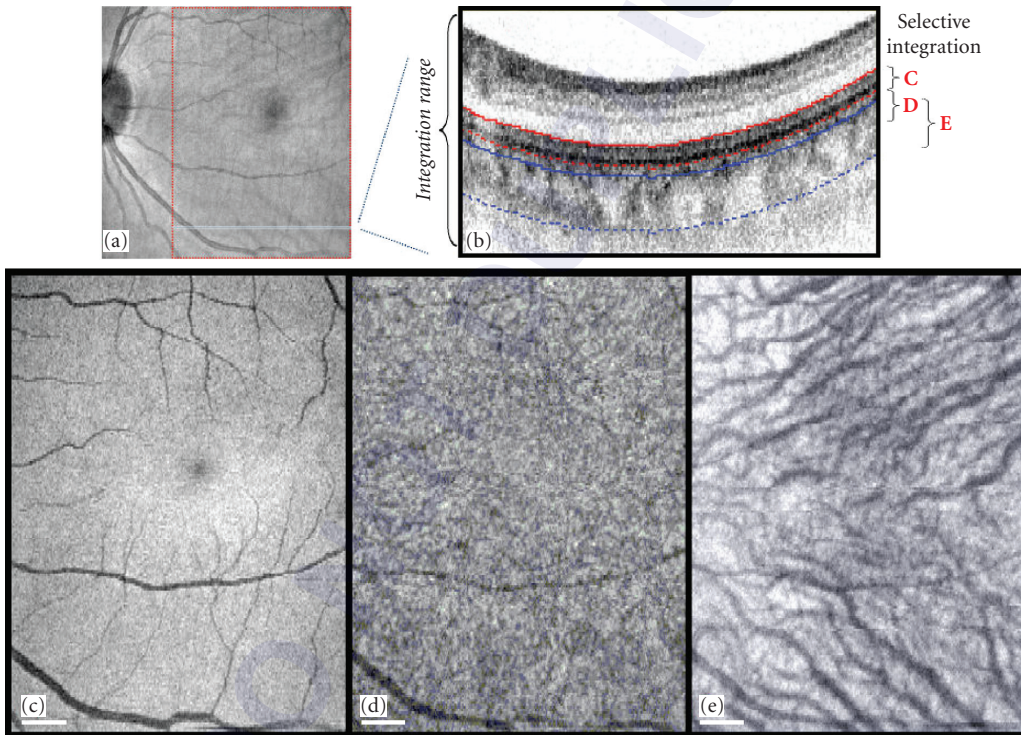


FIGURE 11 The retinal and choroidal vasculature extracted from a three-dimensional OFDI data set. (a) Two-dimensional reflectance image (5.3×5.2 mm²) obtained with the conventional full-range integration method. Higher (lower) reflectivity is represented by white (black) in the grayscale. (b) Illustration of the axial-sectioning integration method, with the different integration regions labeled C,D,E corresponding to the following fundus-type reflectivity images, respectively: (c) retinal reflectivity image showing the shadow of retinal vasculature (3.8×5.2 mm²), (d) reflectivity image obtained from the upper part of the choroid, and (e) reflectivity image from the center of the choroid revealing the choroidal vasculature. Shadows of retinal vasculature are also visible in (d) and (e). Scale bars: 0.5 mm. (Figure reproduced from Ref. 26 with permission from the Optical Society of America.)

image generated from the entire OFDI image sequence. The image visualizes the optic nerve head, fovea, retinal vessels, and the faint outline of the deep choroidal vasculature; however, the depth information is completely lost. To overcome this limitation of the conventional method, we integrated only selective regions based on anatomical structures. For example, to visualize the retinal vasculature with maximum contrast, we used automatic image segmentation techniques⁴³ and integrated the reflectivity in the range between IPRL and RPE (marked by red lines and labeled C in Fig. 11*b*), where the shadow or loss of signal created by the retinal vessels above appeared most distinctly.⁴² Integrating over the entire retina, including the vessel, often results in a lower contrast in the vasculature because retinal blood vessels produce large signals by strong scattering. Figure 11*c* depicts the fundus-type reflectivity image (shadow) of the retinal vessels produced with this depth-sectioning method. The choriocapillary layer contains abundant small blood vessels and pigment cells. Using a thin integration region in the upper part of choroid (labeled D in Fig. 11*b*), we also obtained an image of the choriocapillary layer (Fig. 11*d*). To obtain an image of the complete choroid region, we used the bottom integration region (marked by blue lines and labeled E in Fig. 11*b*). The choroidal vasculature is clearly visualized in the resulting reflectivity image (Fig. 11*e*).

Figure 12 demonstrates OFDI at 1050 nm in an AMD patient. The left and right panel show images at the same location pre- and posttreatment with anti-VEGF therapy. The advantage of 1050 nm is a better penetration below the retinal pigmented epithelium (RPE), providing detailed information of what is believed to be type I vascularization between the RPE and Bruchs membrane (presumed sub-RPE CNV).

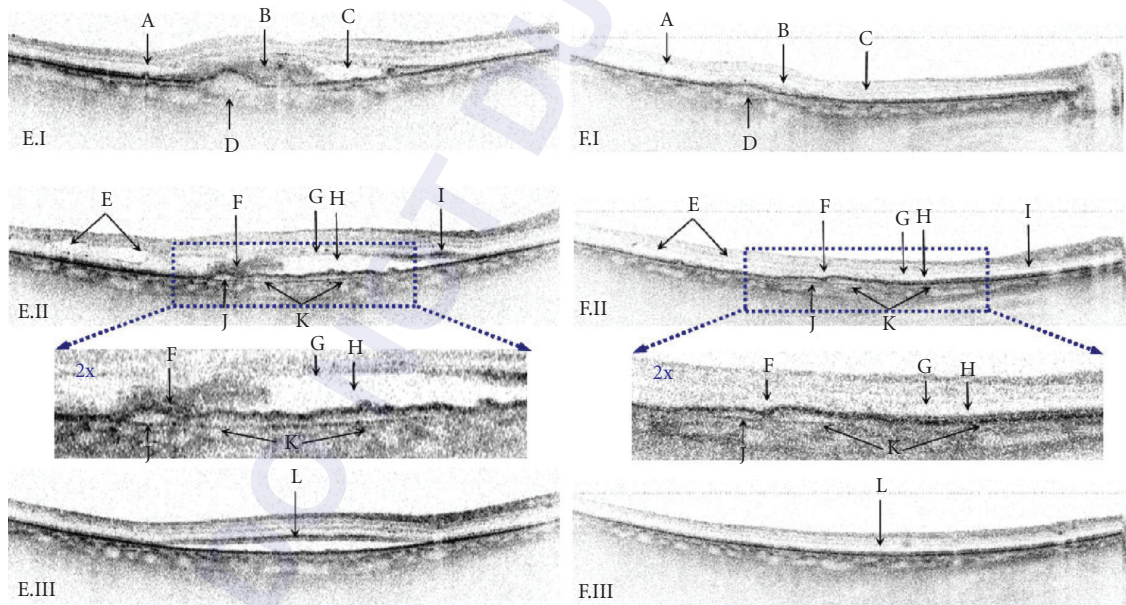


FIGURE 12 OFDI images at 1050 nm of an AMD patient pretreatment: E.I-III, and posttreatment: F.I-III. Features—A: drusen, B: blood clot, C: subretinal fluid, D: RPE detachment, E: cystic changes, F: blood clot, G: weak scattering in photoreceptors, H: subretinal fluid, I: strong scattering in photoreceptors, J: RPE detachment, K: presumed sub-RPE CNV, and L: strong scattering from photoreceptors in the periphery of the subretinal fluid. (Reproduced from Ref. 29 with permission from the Association for Research in Vision and Ophthalmology.)

18.14 FUNCTIONAL EXTENSIONS: DOPPLER OCT AND POLARIZATION SENSITIVE OCT

Optical coherence tomography is an interferometric technique capable of noninvasive high-resolution cross-sectional imaging by measuring the intensity of light reflected from within tissue.² This results in a noncontact imaging modality that provides images similar in scale and geometry to histology. Just as different stains can be used to enhance the contrast in histology, various extensions of OCT allow for visualization of features not readily apparent in traditional OCT. For example, optical Doppler tomography⁴⁴ can enable depth-resolved imaging of flow by observing differences in phase between successive depth scans.^{45–47} Polarization-sensitive OCT (PS-OCT) utilizes depth-dependent changes in the polarization state of detected light to determine the light polarization changing properties of a sample.^{48–53}

18.15 DOPPLER OCT AND PHASE STABILITY

In the past, phase-resolved optical Doppler tomography (ODT) based on time domain OCT (TD-OCT) has proven able to make high-resolution, high-velocity-sensitivity cross-sectional images of in vivo blood flow.^{45–47,54–59} ODT measurements of blood flow in the human retina have been demonstrated,^{60,61} yet the accuracy and sensitivity was compromised by a slow A-line rate and patient motion artifacts, which can introduce phase inaccuracy and obscure true retinal topography. Originally, Doppler shifts were detected by the shift in the carrier frequency in TD-OCT systems, where a trade-off between Doppler sensitivity and spatial resolution had to be made.^{54–56,62} The detection of Doppler shifts improved significantly with the method pioneered by Zhao et al.^{45,46} In this method two sequential A-lines are acquired at the same location. Phase-resolved detection of the interference fringes permits the determination of small phase shifts between the interferograms. The phase difference $\Delta\phi$ divided by the time laps ΔT between the sequential A-lines gives a Doppler shift $\Delta\omega = \Delta\phi / \Delta T$, associated with the motion of the scattering particle. Figure 13 gives a graphical example of this method.

Combining optical Doppler tomography with the superior sensitivity and speed of SD-OCT has allowed a significant improvement in detecting Doppler signals in vivo. In the first combination of these technologies, velocity of a moving mirror and capillary tube flow was demonstrated,⁶³ followed by in vivo demonstration of retinal blood flow.^{22,23}

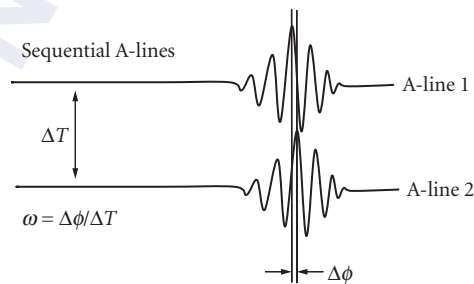


FIGURE 13 Principle of fast Doppler OCT. OCT depth profiles are acquired sequentially at the same location. Small motion of the scattering object results in a phase shift of the interference fringes. The phase difference divided by the time lapse between the sequential A-lines gives a Doppler shift due to the motion of the scattering object.

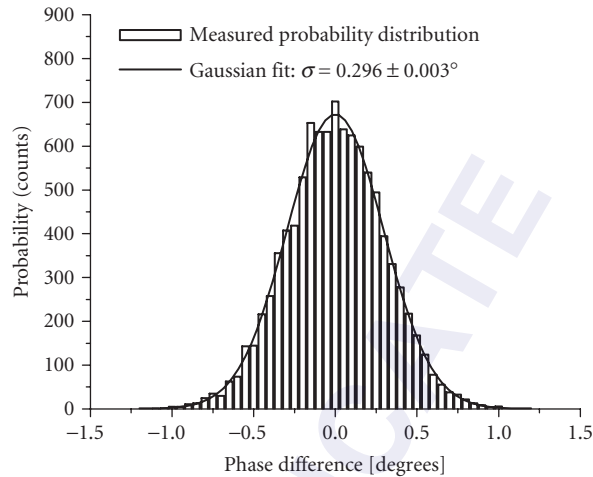


FIGURE 14 Probability distribution of the measured phase difference between adjacent A-lines, with a stationary reflector in the sample arm. Bars: Counted phase difference for 9990 A-lines. Bin size = 0.05° . Solid line: Gaussian fit to the distribution, with a measured standard deviation of $0.296 \pm 0.003^\circ$. (Reproduced from Ref. 23 with permission from the Optical Society of America.)

In SD-OCT, a phase-sensitive image is generated by simply determining the phase difference between points at the same depth in adjacent A-lines. The superior phase stability of SD-OCT, due to the absence of moving parts, is demonstrated in Fig. 14. The data was acquired with a stationary mirror in the sample arm, without scanning the incident beam. Ideally, interference between sample and reference arm light should have identical phase at the mirror position for all A-lines. This condition underlies the assumption that any phase difference between adjacent A-lines is solely due to motion within the sample. The actual phase varies in a Gaussian manner about this ideal, as demonstrated in Fig. 14, where we present the measured probability distribution of phase differences with a standard deviation of $0.296 \pm 0.003^\circ$. This value is over 25 times lower than previously quantified figures for time domain optical Doppler tomography systems,^{59,64} and at an acquisition speed of 29 kHz corresponds to a minimum detectable Doppler shift of ± 25 Hz. With a time difference of $34.1 \mu\text{s}$ between acquired A-lines, phase wrapping occurs at Doppler shifts greater than 15 kHz. Thus, the system dynamic range described by the ratio of maximum to minimum detectable Doppler shifts before phase wrapping occurs is a factor of 600.

In vivo images of structure and Doppler flow were acquired at 29 frames per second (1000 A-lines per frame) and subsequently processed. The images presented in Fig. 15 is 1.6 mm wide and has been cropped in depth to $580 \mu\text{m}$, from their original size of 1.7 mm. The layers of the retina visible in the intensity image have been identified and described previously,³ with the thick, uppermost layer being the nerve fiber layer, and the thinner, strongly scattering deep layer being the retinal pigmented epithelium. One can see the pulsatility of blood flow in the artery (a), while the flow in the vein (v) is less variable (see Ref. 23). At the lower left-center of the image, it is possible to distinguish blood flow deep within the retina (d). With reference to the intensity image, one can see that this blood flow is being detected *below* the retinal pigmented epithelium, and we believe this is the first time that optical Doppler tomography imaging techniques have been able to observe and localize blood flow within the choroid. To the left of the large vessel on the right-hand side of the image, note the appearance of a very small vessel (c). The diameter of this vessel is slightly under $10 \mu\text{m}$.

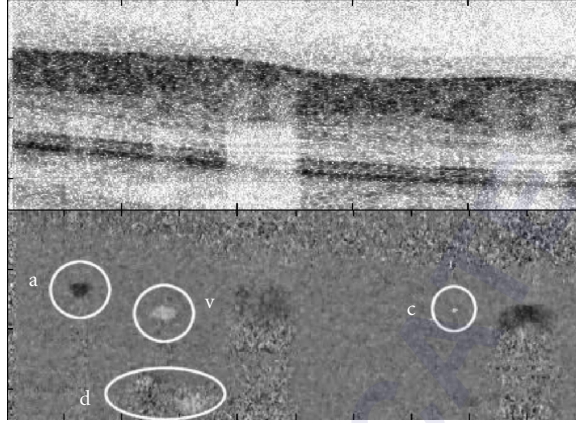


FIGURE 15 Movie of structure (top panel) and bi-directional flow (bottom panel) acquired in vivo in the human eye at a rate of 29 frames per second. The sequence contained 95 frames, totaling 3.28 s (see Ref. 23). Image size is 1.6 mm wide by 580 μm deep. a: artery; v: vein; c: capillary; and d: choroidal vessel. (*Reproduced from Ref. 23 with permission from the Optical Society of America.*)

18.16 POLARIZATION SENSITIVE OCT

Polarization-sensitive OCT (PS-OCT) utilizes depth-dependent changes in the polarization state of detected light to determine the light polarization changing properties of a sample.^{48–53} These material properties, including birefringence, dichroism, and optic axis orientation, can be determined by studying the depth evolution of the Stokes parameters,^{49–52,65–69} or by using the changing reflected polarization states to first determine Jones or Mueller matrices.^{53,70–74} PS-OCT provides additional contrast to identify tissue structures. Nearly all linear tissue structures, like collagen, nerve fibers, and muscle fiber exhibit birefringence. PS-OCT has been used in a wide variety of applications. In dermatology the reduced birefringence in tumor tissue was observed due to the destruction of the extracellular collagen matrix by the tumor.⁷⁵ The reduced birefringence due to thermally denatured collagen was observed in burns.⁶⁷ In ophthalmology, the birefringence of the retinal nerve fiber was measured and found to depend on the location around the optic nerve head.^{76,77} The onset and progression of caries lesions was associated with changes in birefringence of dental tissue.⁷⁸

We will first treat the Jones formalism to describe polarization properties and use the Jones formalism to determine the polarization properties of tissue. The Jones formalism provides a convenient mathematical description of polarized light and polarization effects.⁷⁹ The complex electric field vector \mathbf{E} can be decomposed into a pair of orthonormal basis vectors to yield

$$\begin{aligned}\mathbf{E} &= E_p \hat{\mathbf{e}}_p + E_\perp \hat{\mathbf{e}}_\perp \\ E_p &= a_p e^{-i\delta_p} \quad E_\perp = a_\perp e^{-i\delta_\perp}\end{aligned}\tag{21}$$

where $\hat{\mathbf{e}}_p$ and $\hat{\mathbf{e}}_\perp$ are unit vectors along the horizontal and vertical, a_p and a_\perp are the amplitude along the horizontal and vertical, and δ_p and δ_\perp are the phase along the horizontal and vertical, respectively. In this case, the vibrational ellipse of the electric field can be reformulated as

$$\mathbf{E}_{\text{vib}}(t) = a_p \cos(\omega t + \delta_p) \hat{\mathbf{e}}_p + a_\perp \cos(\omega t + \delta_\perp) \hat{\mathbf{e}}_\perp\tag{22}$$

with ω the angular frequency and t time. The overall irradiance, or intensity, of the beam of light then can be expressed as the scalar quantity

$$I = a_p^2 + a_\perp^2 \quad (23)$$

It is worth noting that while the overall irradiance of a beam is not dependent on its polarization state, it is possible to measure irradiance along a particular orientation (e.g., the intensity of a beam in the horizontally polarized direction).

Linear polarization states occur for phase differences $\Delta\delta = \delta_p - \delta_\perp = m\pi$, where $m \in \mathbb{Z}$, as the vibrational ellipse collapses to a line described by

$$E_{\text{vib}}(t) = (a_p \hat{e}_p + (-1)^m a_\perp \hat{e}_\perp) \cos(\omega t + \delta_p) \quad (24)$$

The orientation of the linear polarization state depends on the ratio of amplitudes a_p and a_\perp . The polarization state of light is horizontal or vertical when $a_\perp = 0$ or $a_p = 0$, respectively, and oriented at $\pm 45^\circ$ if $|a_p| = |a_\perp|$. An orientation angle θ can then be defined according to the relations

$$\begin{aligned} a_p &= a \cos\theta & a_\perp &= a \sin\theta \\ a &= \sqrt{a_p^2 + a_\perp^2} & \theta &= \tan^{-1} \frac{a_p}{a_\perp} \end{aligned} \quad (25)$$

where a can be thought of as an overall amplitude of the electric field and $\Delta\delta = 0$.

Circular polarization states are obtained when $|a_p| = |a_\perp|$ and $\Delta\delta = \frac{\pi}{2} + n\pi$, which is evident through the form of the resultant vibrational ellipse

$$E_{\text{vib}}(t) = a_p (\cos(\omega t + \delta_p) \hat{e}_p \pm \sin(\omega t + \delta_p) \hat{e}_\perp) \quad (26)$$

This describes a circle, the handedness (left- or right-circular) of which is determined by the sign between the orthogonal components. Circular polarization states differ only in their phase difference compared to linearly polarized light at 45° , and so the phase difference $\Delta\delta$ between orthogonal electric field components reflects the ratio between circular and linear components of the polarization state. Figure 16 gives a graphical representation of the different vibrational ellipses

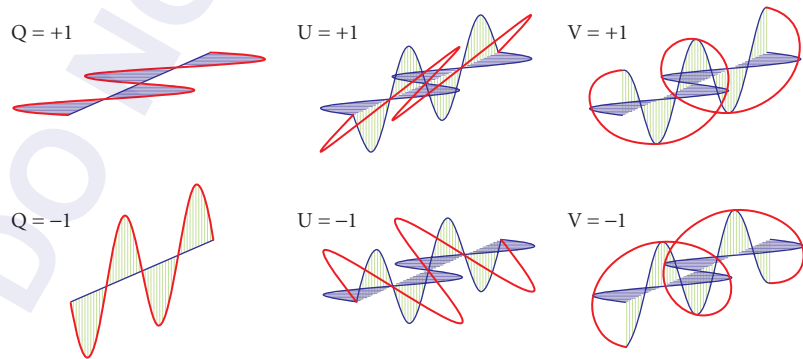


FIGURE 16 Vibrational ellipses for various polarization states $Q = 1$ and $Q = -1$ correspond to horizontal and vertical linear polarized light. $U = 1$ and $U = -1$ correspond to linear polarized light at $+45^\circ$ and -45° , and $V = 1$ and $V = -1$ correspond to circular polarized light, respectively.

The electric field decomposition in Eq. (21) can be rewritten as a complex 2-vector such that

$$\mathbf{E} = \begin{bmatrix} E_p \\ E_{\perp} \end{bmatrix} = \begin{bmatrix} a_p e^{-i\delta_p} \\ a_{\perp} e^{-i\delta_{\perp}} \end{bmatrix} = a e^{-i\delta_p} \begin{bmatrix} \cos\theta \\ e^{i\Delta\delta} \sin\theta \end{bmatrix} \quad (27)$$

while the time-invariant electric field vector \mathbf{E} , also known as a Jones vector, does depend on the amplitude and exact phase of the electric field components, it should be noted that the polarization state itself is completely determined by the orientation angle θ and the phase difference $\Delta\delta$.

Just as two vectors of length n can be related using a matrix of dimension $n \times n$, two polarization states can be related using a complex 2×2 matrix known as a Jones matrix. The polarization properties of any nondepolarizing optical system can be described using a Jones matrix. The transmitted polarization state \mathbf{E}' as a result of an optical system represented by a Jones matrix \mathbf{J} acting on an incident polarization state \mathbf{E} can be determined by

$$\mathbf{E}' = \begin{bmatrix} E'_p \\ E'_{\perp} \end{bmatrix} = \begin{bmatrix} J_{11} & J_{12} \\ J_{21} & J_{22} \end{bmatrix} \begin{bmatrix} E_p \\ E_{\perp} \end{bmatrix} = \mathbf{J}\mathbf{E} \quad (28)$$

Subsequent transmission of \mathbf{E}' through an optical system \mathbf{J}' results in a polarization state $\mathbf{E}'' = \mathbf{J}'\mathbf{E}' = \mathbf{J}'(\mathbf{J}\mathbf{E}) = \mathbf{J}'\mathbf{J}\mathbf{E}$. As a result, the combined polarization effect of a cascade of optical elements, $\mathbf{J}_1, \mathbf{J}_2, \dots, \mathbf{J}_n$, can be described by the product $\mathbf{J} = \mathbf{J}_n \cdots \mathbf{J}_2 \mathbf{J}_1$.

The Jones matrix for a birefringent material that induces a phase retardation η between electric field components parallel and orthogonal to a polarization state characterized by an orientation angle θ and a circularity related to ϕ is given by⁸⁰

$$\mathbf{J}_b = \begin{bmatrix} e^{i\eta/2} C_{\theta}^2 + e^{-i\eta/2} S_{\theta}^2 & (e^{i\eta/2} - e^{-i\eta/2}) C_{\theta} S_{\theta} e^{-i\phi} \\ (e^{i\eta/2} - e^{-i\eta/2}) C_{\theta} S_{\theta} e^{i\phi} & e^{i\eta/2} S_{\theta}^2 + e^{-i\eta/2} C_{\theta}^2 \end{bmatrix} \quad (29)$$

where $C_{\theta} = \cos\theta$ and $S_{\theta} = \sin\theta$. The Jones matrix of a dichroic material with attenuation ratios of P_1 and P_2 for electric field components parallel and orthogonal, respectively, to a polarization state given by an orientation angle Θ and a circularity Φ has the form⁸⁰

$$\mathbf{J}_d = \begin{bmatrix} P_1 C_{\Theta}^2 + P_2 S_{\Theta}^2 & (P_1 - P_2) C_{\Theta} S_{\Theta} e^{-i\Phi} \\ (P_1 - P_2) C_{\Theta} S_{\Theta} e^{i\Phi} & P_1 S_{\Theta}^2 + P_2 C_{\Theta}^2 \end{bmatrix} \quad (30)$$

As birefringence seems to be the primary polarization property exhibited by biological tissue, most PS-OCT analysis methods concentrate on determination of phase retardation. Schoenberger et al.⁸¹ analyzed system errors introduced by the extinction ratio of polarizing optics and chromatic dependence of wave retarders, and errors due to dichroism. System errors can be kept small by careful design of the system with achromatic elements, but can never be completely eliminated. In principle, dichroism is a more serious problem when interpreting results as solely due to birefringence. However, Mueller matrix ellipsometry measurements have shown that the error due to dichroism in the eye is relatively small,^{82,83} and earlier PS-OCT work shows that dichroism is of minor importance in rodent muscle.⁵² Despite this, a method for simultaneous determination of sample birefringence and dichroism is desirable, especially one that can be applied to systems with the unrestricted use of optical fiber and fiber components.

The nondepolarizing polarization properties of an optical system can be completely described by its complex Jones matrix, \mathbf{J} , which transforms an incident polarization state, described by a complex electric field vector, $\mathbf{E} = [H, V]^T$, to a transmitted state, $\mathbf{E}' = [H', V']^T$. A Jones matrix can be decomposed in the form $\mathbf{J} = \mathbf{J}_R \mathbf{J}_P = \mathbf{J}_P \mathbf{J}_R$.⁸⁰ Birefringence, described by \mathbf{J}_R , can be parameterized by three variables: a degree of phase retardation η about an axis defined by two angles, γ and δ .

Diattenuation, described by J_p , is defined as $d = (P_1^2 - P_2^2) / (P_1^2 + P_2^2)$ and can be parameterized by four variables, where P_1 and P_2 are the attenuation coefficients parallel and orthogonal, respectively, to an axis defined by angles Γ and Δ . These seven independent parameters, along with an overall common phase $e^{i\psi}$, account for all four complex elements of a general Jones matrix \mathbf{J} . Assuming that birefringence and diattenuation arise from the same fibrous structures in biological tissue and thus share a common axis $\delta = \Delta$ and $\gamma = \Gamma$,⁷² the number of independent parameters is reduced by two to five. In order to determine these five parameters, the sample needs to be probed by two unique polarization states. One incident and reflected polarization state yield three relations involving the two orthogonal amplitudes and the relative phase between them.⁵² Therefore, it is possible to use the six relationships defined by two unique pairs of incident and reflected polarization states to exactly solve for the Jones matrix of a sample.

In general terms, a PS-OCT system sends polarized light from a broadband source into the sample and reference arms of an interferometer, and reflected light from both arms is recombined and detected. Figure 17 shows an example of a fiber-based PS-OCT system. The source light is sent to a fiber-based polarization controller and a polarizer. The polarization of the incident state on the sample is modulate by a polarization modulator that allows switching between two polarization states. Light passes a circulator, which sends the light to a fiber-based beam splitter that splits the light into a sample and a reference arm. Upon reflection, the light is recombined in the beam splitter and the circulator sends the light to the detection arm where the light is split into two orthogonal polarization states before detection by two separate detectors. The optical path from source to detector can be described by three Jones matrices. Define J_{in} as the Jones matrix representing the optical path from the polarized light source (the polarization modulator) to the sample surface, J_{out} as that going from the sample surface to the detectors, and J_s as the round-trip Jones matrix for light propagation through a sample.⁷⁴ This nomenclature can be applied to all PS-OCT systems, ranging from bulk-optic systems^{48,49,51–53,65} to those with fibers placed such that they are traversed in a round-trip manner,⁷³ to time-domain^{66, 68} and spectral-domain⁶⁹ PS-OCT systems with the unrestricted use of optical fiber and nondiattenuating fiber components, and even for retinal systems,⁷⁶ where the polarization effects of the cornea can be included in J_{in} and J_{out} . The electric field of light reflected from the

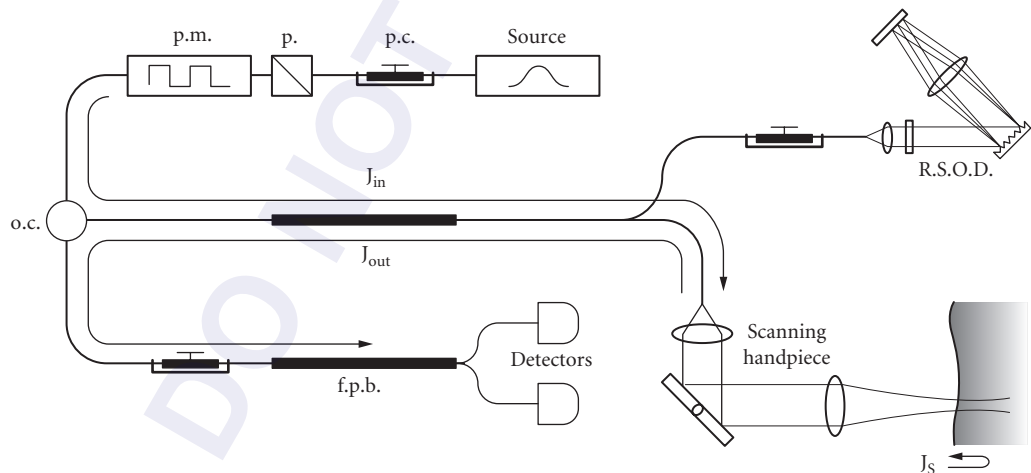


FIGURE 17 Schematic of the fiber-based PS-OCT system (p.c., polarization controller; p, polarizer; pm, polarization modulator; oc, optical circulator; RSOD, rapid scanning optical delay; and fpb, fiber polarizing beamsplitter). J_{in} , J_{out} , and J_s are the Jones matrix representations for the one-way optical path from the polarization modulator to the scanning handpiece, the one-way optical path back from the scanning handpiece to the detectors, and the round-trip path through some depth in the sample, respectively. (Reprinted from Ref. 74 with permission from the Optical Society of America.)

sample surface, \mathbf{E} , can be expressed as $\mathbf{E} = e^{i\psi} \mathbf{J}_{\text{out}} \mathbf{J}_{\text{in}} \mathbf{E}_{\text{source}}$, where ψ represents a common phase and $\mathbf{E}_{\text{source}}$ represents the electric field of light coming from the polarized source. Likewise, the electric field of light reflected from some depth within the tissue may be described by $\mathbf{E}' = e^{i\psi'} \mathbf{J}_{\text{out}} \mathbf{J}_{\text{in}} \mathbf{E}_{\text{source}}$. These two measurable polarization states can be related to each other such that $\mathbf{E}' = e^{i\Delta\psi} \mathbf{J}_T \mathbf{E}$, where $\mathbf{J}_T = \mathbf{J}_{\text{out}} \mathbf{J}_S \mathbf{J}_{\text{out}}^{-1}$ and $\Delta\psi = \psi' - \psi$. Thus, \mathbf{J}_T can be determined by comparing the polarization state reflected from the surface and from a certain depth in the sample. Now we need to determine how to extract \mathbf{J}_S from \mathbf{J}_T .

If the optical system representing \mathbf{J}_{out} is nondiattenuating, \mathbf{J}_{out} can be treated as a unitary matrix with unit determinant after separating out a common attenuation factor. \mathbf{J}_{out} is assumed to be nondiattenuating since optical fibers are virtually lossless. \mathbf{J}_S can be decomposed into a diagonal matrix $\mathbf{J}_C = [P_1 e^{i\eta/2}, 0; 0, P_2 e^{-i\eta/2}]$, containing complete information regarding the amount of sample diattenuation and phase retardation, surrounded by unitary matrices \mathbf{J}_A with unit determinant that define the sample optic axis, $\mathbf{J}_S = \mathbf{J}_A \mathbf{J}_C \mathbf{J}_A^{-1}$. \mathbf{J}_T can be rewritten such that $\mathbf{J}_T = \mathbf{J}_{\text{out}} \mathbf{J}_S \mathbf{J}_{\text{out}}^{-1} = \mathbf{J}_{\text{out}} (\mathbf{J}_A \mathbf{J}_C \mathbf{J}_A^{-1}) \mathbf{J}_{\text{out}}^{-1} = \mathbf{J}_U \mathbf{J}_C \mathbf{J}_U^{-1}$, where $\mathbf{J}_U = \mathbf{J}_{\text{out}} \mathbf{J}_A$. Since unitary matrices with unit determinant form the special unitary group $\text{SU}(2)$,⁸⁴ \mathbf{J}_U must also be a unitary matrix with unit determinant by closure and can be expressed in the form

$$\mathbf{J}_U = e^{i\beta} \begin{bmatrix} C_\theta e^{i(\phi-\theta)} & -S_\theta e^{i(\phi+\theta)} \\ S_\theta e^{-i(\phi+\theta)} & C_\theta e^{-i(\phi-\theta)} \end{bmatrix} \quad (31)$$

\mathbf{J}_T can be obtained from the measurements by combining information from two unique incident polarization states, $[H'_1, H'_2; V'_1, V'_2] = e^{i\Delta\psi_1} \mathbf{J}_T [H_1, e^{i\alpha} H_2; V_1, e^{i\alpha} V_2]$, where $\alpha = \Delta\psi_2 - \Delta\psi_1$. The polarization properties of interest can be obtained by equating the two expressions for \mathbf{J}_T to yield

$$e^{i\Delta\psi_1} \begin{bmatrix} P_1 e^{i\eta/2} & 0 \\ 0 & P_2 e^{-i\eta/2} \end{bmatrix} = \begin{bmatrix} C_\theta & S_\theta \\ -S_\theta & C_\theta \end{bmatrix} \begin{bmatrix} e^{-i\phi} & 0 \\ 0 & e^{i\phi} \end{bmatrix} \begin{bmatrix} H'_1 & H'_2 \\ V'_1 & V'_2 \end{bmatrix} \quad (32)$$

$$\begin{bmatrix} H_1 & e^{i\alpha} H_2 \\ V_1 & e^{i\alpha} V_2 \end{bmatrix}^{-1} \begin{bmatrix} e^{i\phi} & 0 \\ 0 & e^{-i\phi} \end{bmatrix} \begin{bmatrix} C_\theta & -S_\theta \\ S_\theta & C_\theta \end{bmatrix}$$

In principle, parameters θ , ϕ , and α can be solved for with the condition that the off-diagonal elements of the matrix product on the right-hand side of Eq. (32) are equal to zero. In practice, real solution cannot always be found, as measurement noise can induce nonphysical transformations between incident and transmitted polarization states. To account for this, Eq. (32) can be solved by optimizing parameters θ , ϕ , and α to minimize the sum of the magnitudes of the off-diagonal elements. In principle, this can be achieved using two unique incident polarization states to probe the same volume of a sample. However, when two orthogonal incident polarization states are used,⁷³ birefringence cannot be retrieved under all circumstances.⁸⁵ A better choice is to use two incident polarization states perpendicular in a Poincaré sphere representation to guarantee that polarization information can always be extracted.^{66–69,76,86} The degree of phase retardation can easily be extracted through the phase difference of the resulting diagonal elements, and the diattenuation by their magnitudes. It should be noted that these phase retardation values range from $-\pi$ to π , and can therefore be unwrapped to yield overall phase retardations in excess of 2π .

18.17 PS-OCT IN OPHTHALMOLOGY

Ophthalmological application of OCT has arguably driven a great deal of its development, and probably represents the most researched clinical application of the technology to date. PS-OCT in particular has been used to measure the birefringence of the human retinal nerve fiber layer in vivo^{48,76,77,87–89} for potential early detection of glaucoma, the world's second leading cause of blindness.

Glaucoma causes damage to the retinal ganglion cells, resulting in a thinning of the retinal nerve fiber layer (RNFL). In addition, nerve fiber layer tissue loss may be preceded by changes in birefringence as ganglion cells become necrotic and axons in the RNFL are replaced by a less organized and amorphous tissue composed of glial cells. When glaucoma is detected at an early stage, further loss of vision can be prevented by treatment. The visual field test is the current standard method of detecting loss of peripheral vision in glaucoma. However, measurements show that up to 40 percent of nerves are irreversibly damaged before loss of peripheral vision can be clinically detected. PS-OCT has the potential to detect changes to the RNFL at an earlier time point through changes in its birefringence and thickness.

Ophthalmic studies can be performed using systems similar to that used by Cense et al.,⁷⁶ in which a slit lamp has been adapted for use with PS-OCT. Figure 18 is a typical example of a structural-intensity time-domain OCT image of the retina in the left eye of a healthy volunteer obtained with a circular scan with a radius of 2.1 mm around the optic nerve head (ONH). The image measures 13.3 mm wide and 0.9 mm deep and is shown at an expanded aspect ratio in depth for clarity. Structural layers such as the RNFL, the interface between the inner and outer segments of the photoreceptors, and the retinal pigmented epithelium can be seen.

The addition of polarization sensitivity allows for localized quantitative assessment of the thickness and birefringence of the RNFL. Figure 19 shows two examples of combined thickness and birefringence measurements, one of a region temporal to the ONH, the other of a region superior to the ONH. The depth of the RNFL can be determined by a decrease in backscattered intensity from the RNFL to the inner plexiform layer. The birefringence of the RNFL can then be estimated from a linear least-square fit of the measured double-pass phase retardation through the determined depth. Two main observations can be drawn from such graphs: the retinal layers directly below the RNFL are minimally birefringent, and that the thickness and birefringence of the RNFL are not constant. These observations can also be seen in Fig. 20, which overlays the thickness and birefringence determined

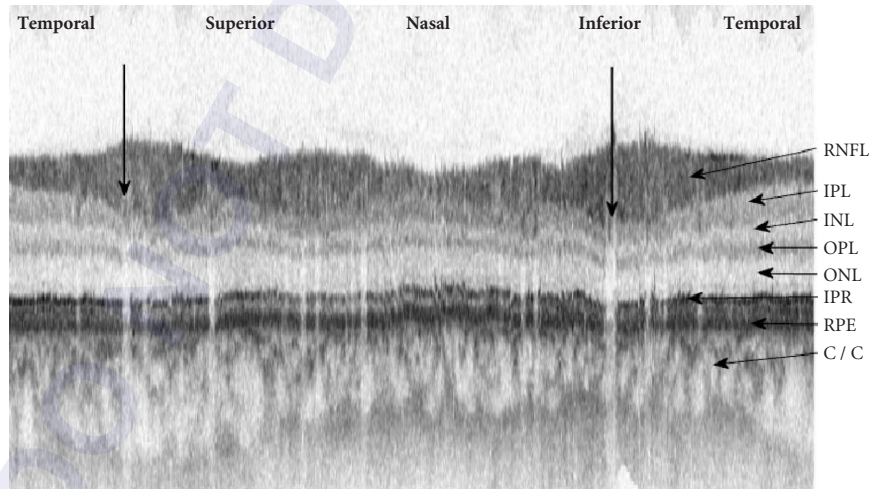


FIGURE 18 A realigned OCT intensity image created with a 2.1-mm radius circular scan around the ONH. The dynamic range of the image is -36 dB. Black pixels represent strong reflections. The image measures 13.3 mm wide and 0.9 mm deep. Visible structures: retinal nerve fiber layer (RNFL); inner plexiform layer (IPL); inner nuclear layer (INL); outer plexiform layer (OPL); outer nuclear layer (ONL); interface between the inner and outer segments of the photoreceptor layer (IPR); retinal pigmented epithelium (RPE); and choriocapillaris and choroid (C/C). Vertical arrows: locations of the two largest blood vessels. Other smaller blood vessels appear as vertical white areas in the image. (Reprinted from Ref. 77 with permission from the Association for Research in Vision and Ophthalmology.)

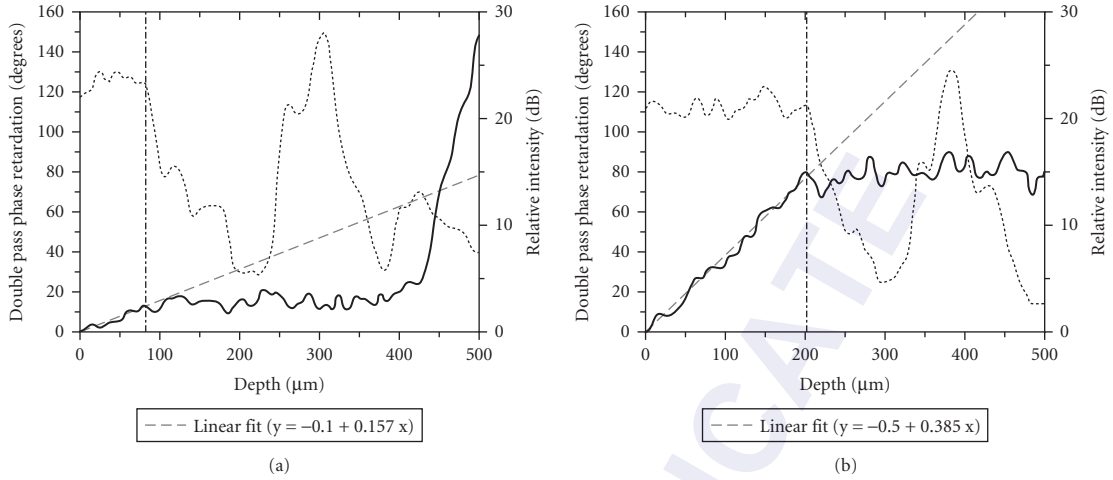


FIGURE 19 Thickness (dotted line) and birefringence (solid line) plots of an area temporal (a) and superior (b) to the ONH. DPPR data belonging to the RNFL is fit with a least-square linear fit. The slope in the equation represents the DPPR/UD or birefringence. The vertical line indicates the boundary of the RNFL, as determined from the intensity and DPPR data. (a) The increase in DPPR at a depth beyond 450 μm is caused by either a relatively low signal-to-noise ratio, or by the presence of a highly birefringent material—for instance, collagen in the sclera. (Reprinted from Ref. 77 with permission from the Association for Research in Vision and Ophthalmology.)

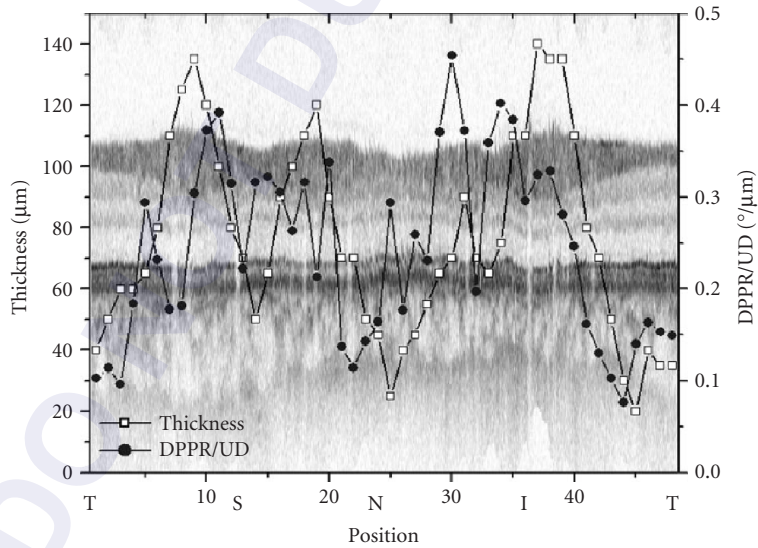


FIGURE 20 A typical example of combined RNFL thickness and birefringence measurements along a circular scan around the ONH. The intensity image is plotted in the background. The RNFL is relatively thicker superiorly (S) and inferiorly (I). A similar development can be seen in the birefringence plot. The birefringence is relatively higher in the thicker areas, whereas it is lower in the thinner temporal (T) and nasal (N) areas. (Reprinted from Ref. 77 with permission from the Association for Research in Vision and Ophthalmology.)

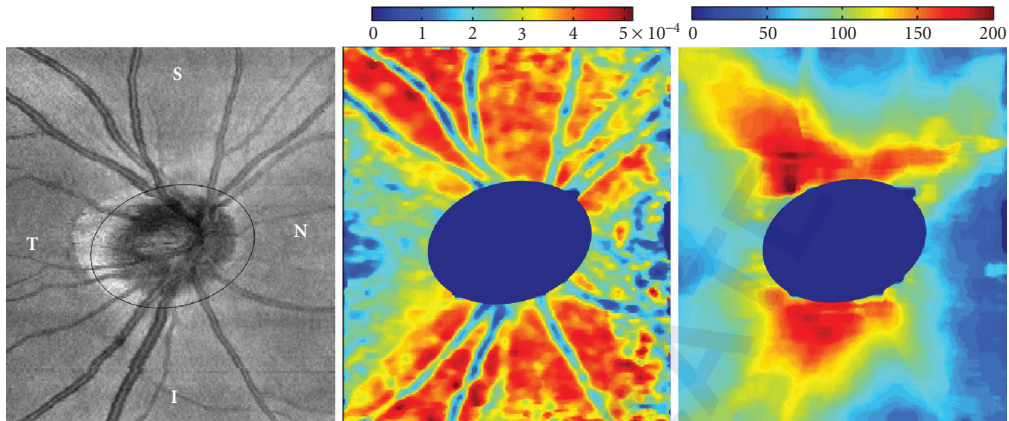


FIGURE 21 OCT scan ($4.24 \times 5.29 \text{ mm}^2$) of the retina of a normal volunteer, centered on the ONH. (a) Integrated reflectance map showing a normal temporal crescent (white area temporal to the ONH); (b) birefringence map; and (c) RNFL thickness map (color bar scaled in microns). The circle on the left indicates the excluded area in the birefringence and thickness maps as corresponding to ONH. (S = superior, N = nasal, I = inferior, T = temporal). (Reprinted from Ref. 90 with permission from the International Society for Optical Engineering.)

as in Fig. 19 on a circular scan across the ONH. The plots indicate that the RNFL is thickest and most birefringent superiorly and inferiorly to the ONH.

En face maps of RNFL thickness and birefringence can be generated from data obtained with recently developed spectral-domain ophthalmic PS-OCT systems.⁹⁰ A three-dimensional volume ($4.24 \times 5.29 \times 1.57 \text{ mm}^3$) of the retina of a normal volunteer (right eye) was scanned at a rate of 29 fps with 1000 A-lines/frame, and contains 190 frames (B-scans) acquired in 6.5 s. The integrated reflectance, birefringence, and retinal nerve fiber layer (RNFL) thickness maps are shown in Fig. 21, confirming previous findings that the RNFL birefringence is not uniform across the retina. Superior, nasal, inferior, and temporal areas of the retina around the ONH are indicated by the letters S, N, I, and T. The integrated reflectance map, obtained by simply integrating the logarithmic depth profiles, illustrates the blood vessel structure around the ONH. The RNFL thickness map is scaled in microns (color bar on the top of the image) indicating a RNFL thickness of up to 200 μm . The central dark-blue area corresponds to the position of the ONH that was excluded from both the thickness and the birefringence maps. A typical bow-tie pattern can be seen for the distribution of the RNFL thickness around the ONH, showing a thicker RNFL superior and inferior to the ONH. The birefringence map illustrates a variation of the birefringence values between 0 and 5.16×10^{-4} , and it clearly demonstrates that the RNFL birefringence is not uniform across the retina; it is smaller nasal and temporal and larger superior and inferior to the ONH.

Given that measurements of the thickness and birefringence of the RNFL can be acquired with the speed and accuracy demonstrated, further research into changes in these parameters with glaucoma can be performed. Experiments—for instance, a longitudinal study with PS-OCT on patients at high risk for development of glaucoma—will either confirm or reject the hypothesis. In addition, PS-OCT can enhance the specificity in determining RNFL thickness in structural OCT images by using changes in tissue birefringence to determine the border between the RNFL and the ganglion cell layer.

18.18 RETINAL IMAGING WITH SD-OCT

Many examples of retinal imaging with SD-OCT are available in the literature. Commercial systems are being introduced into the market at this point. Below, two typical examples of high



FIGURE 22 High resolution SD-OCT image of a human retina in vivo, centered on the optic nerve head. The image is magnified in depth by a factor of 2. Image size: 4.662×1.541 mm.

quality SD-OCT images are presented, acquired with a system providing a depth resolution of $3 \mu\text{m}$ in tissue. Both figures consist of approximately 1000 A-lines or depth profiles. Figure 22 shows a cross-section centered on the optic nerve head. Figure 23 shows an image centered on the fovea and the corresponding en face image (Fig. 24) generated from the three-dimensional data set.

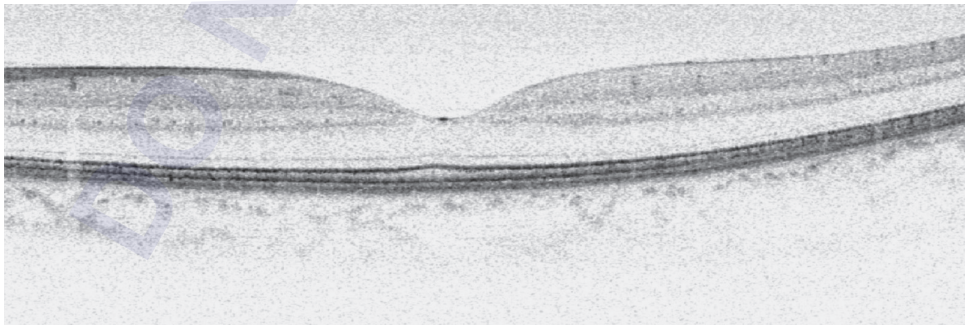


FIGURE 23 High resolution SD-OCT image of a human retina in vivo, centered on the fovea. The image is magnified in depth by a factor of 2. Image size: 4.973×0.837 mm.

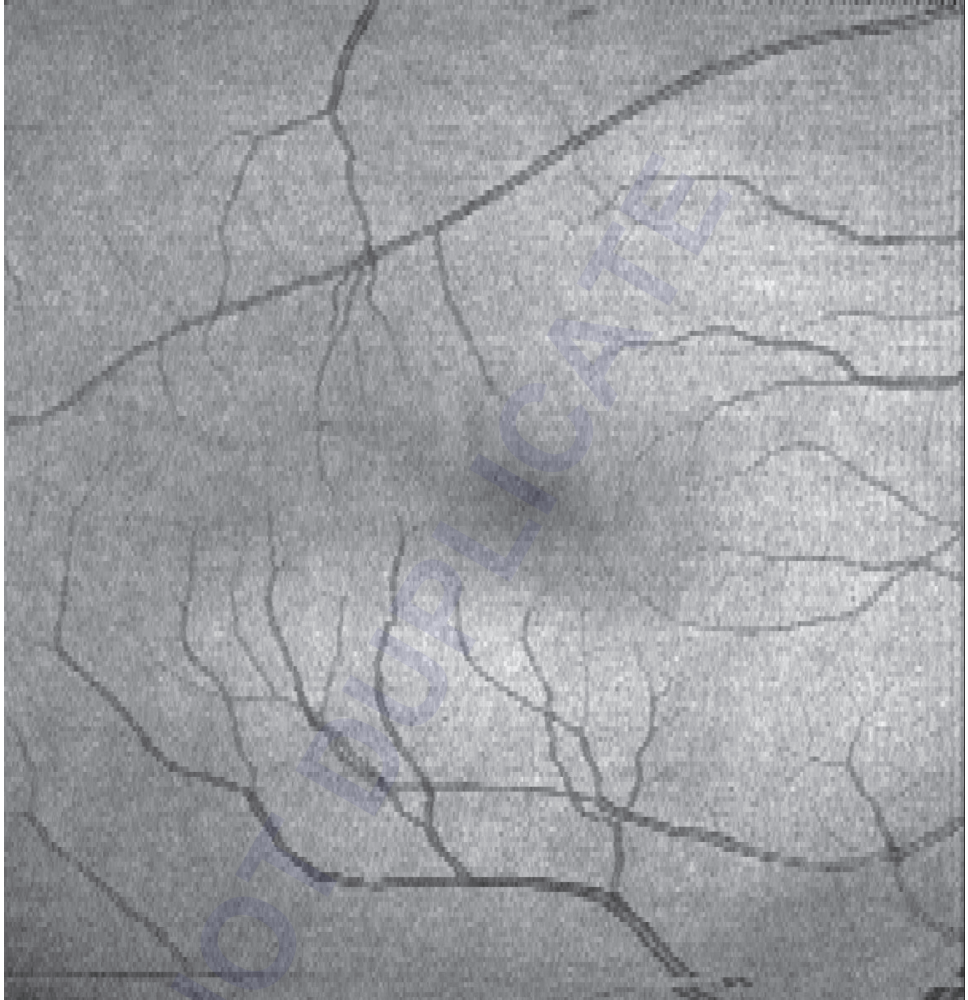


FIGURE 24 En face reconstruction of the fovea region of the retina from a three-dimensional volumetric SD-OCT data set of 200 images. Image size: 4.97 mm \times 5.18 mm.

18.19 CONCLUSION

Spectral domain or frequency domain OCT (SD/FD-OCT) has become the preferred method for retinal imaging owing to its high imaging speed,^{10,19} enhanced signal to noise ratio (SNR),^{7,15–17} and the availability of broadband sources permitting ultrahigh resolution retinal imaging.^{20,21} However, the state-of-the-art spectrometers are hindering further improvements (1) with limited detection efficiency (~25 percent)¹⁹ and (2) the obtainable spectral resolution causes approximately a 6-dB sensitivity drop over a 1-mm depth range.¹⁹ Furthermore, rapid scanning of the probe beam in SD-OCT has the adverse effect of fringe washout, which causes SNR to decrease.³⁶ Fringe washout can be addressed by pulsed illumination.³⁸ A competing technique such as optical frequency domain imaging

(OFDI), the dominant implementation of Fourier domain OCT technologies at 1.3 μm ,⁹ has the advantage of larger depth range and better immunity to motion artifacts. OFDI has recently been demonstrated in the 800 and 1050 nm range,^{26,27,91} but has not reached the superior resolution of SD-OCT.^{20,21}

18.20 ACKNOWLEDGMENTS

This research was supported in part by research grants from the National Institutes of Health (1R24 EY12877, R01 EY014975, and RR19768), Department of Defense (F4 9620-01-1-0014), CIMIT, and a gift from Dr. and Mrs. J. S. Chen to the optical diagnostics program of the Wellman Center of Photomedicine. The author would like to thank a number of graduate students and post doctoral fellows who have contributed to the results presented in this chapter, Barry Cense, Nader Nassif, Brian White, Hyle Park, Jang Woo You, Mircea Mujat, Hyungsik Lim, Martijn de Bruin, Daina Burnes, and Yueli Chen. Special thanks to Teresa Chen, MD, my invaluable collaborator at the Massachusetts Eye and Ear Infirmary, without whom all this work would not have been possible.

18.21 REFERENCES

1. A. F. Fercher, K. Mengedoht, and W. Werner, "Eye-Length Measurement by Interferometry with Partially Coherent-Light," *Optics Letters*, 1988, **13**(3):186–188.
2. D. Huang, E. A. Swanson, C. P. Lin, J. S. Schuman, W. G. Stinson, W. Chang, M. R. Hee, et al., "Optical Coherence Tomography," *Science*, 1991, **254**(5035):1178–1181.
3. W. Drexler, H. Sattmann, B. Hermann, T. H. Ko, M. Stur, A. Unterhuber, C. Scholda, et al., "Enhanced Visualization of Macular Pathology with the Use of Ultrahigh-Resolution Optical Coherence Tomography," *Archives of Ophthalmology*, 2003, **121**(5):695–706.
4. American National Standards Institute, *American National Standard for Safe Use of Lasers Z136.1*, 2000, Orlando.
5. A. F. Fercher, C. K. Hitzenberger, G. Kamp, and S. Y. Elzaiat, "Measurement of Intraocular Distances by Backscattering Spectral Interferometry," *Optics Communications*, 1995, **117**(1–2):43–48.
6. B. Golubovic, B. E. Bouma, G. J. Tearney, and J. G. Fujimoto, "Optical Frequency-Domain Reflectometry Using Rapid Wavelength Tuning of a Cr4+:Forsterite Laser," *Optics Letters*, 1997, **22**(22):1704–1706.
7. M. A. Choma, M. V. Sarunic, C. H. Yang, and J. A. Izatt, "Sensitivity Advantage of Swept Source and Fourier Domain Optical Coherence Tomography," *Optics Express*, 2003, **11**(18):2183–2189.
8. S. H. Yun, G. J. Tearney, J. F. de Boer, N. Iftimia, and B. E. Bouma, "High-Speed Optical Frequency-Domain Imaging," *Optics Express*, 2003, **11**(22):2953–2963.
9. S. H. Yun, G. J. Tearney, B. J. Vakoc, M. Shishkov, W. Y. Oh, A. E. Desjardins, M. J. Suter, et al., "Comprehensive Volumetric Optical Microscopy In Vivo," *Nature Medicine*, 2006, **12**(12):1429–1433.
10. L. Mandel and E. Wolf, *Optical Coherence and Quantum Optics*, 1995, Cambridge, England: Cambridge University Press.
11. N. Nassif, B. Cense, B. H. Park, S. H. Yun, T. C. Chen, B. E. Bouma, G. J. Tearney, and J. F. de Boer, "In Vivo Human Retinal Imaging by Ultrahigh-Speed Spectral Domain Optical Coherence Tomography," *Optics Letters*, 2004, **29**(5):480–482.
12. W. Drexler, U. Morgner, R. K. Ghanta, F. X. Kartner, J. S. Schuman, and J. G. Fujimoto, "Ultrahigh-Resolution Ophthalmic Optical Coherence Tomography," *Nature Medicine*, 2001, **7**(4):502–507.
13. T. Akkin, C. Joo, and J. F. de Boer, "Depth-Resolved Measurement of Transient Structural Changes during Action Potential Propagation," *Biophysical Journal*, 2007, **93**(4):1347–1353.
14. G. Hausler and M. W. Lindner, "Coherence Radar and Spectral Radar—New Tools for Dermatological Diagnosis," *Journal of Biomedical Optics*, 1998, **3**(1):21–31.
15. T. Mitsui, "Dynamic Range of Optical Reflectometry with Spectral Interferometry," *Japanese Journal of Applied Physics Part 1-Regular Papers Short Notes & Review Papers*, 1999, **38**(10):6133–6137.

16. R. Leitgeb, C. K. Hitzenberger, and A. F. Fercher, "Performance of Fourier Domain vs. Time Domain Optical Coherence Tomography," *Optics Express*, 2003, **11**(8):889–894.
17. J. F. de Boer, B. Cense, B. H. Park, M. C. Pierce, G. J. Tearney, and B. E. Bouma, "Improved Signal-to-Noise Ratio in Spectral-Domain Compared with Time-Domain Optical Coherence Tomography," *Optics Letters*, 2003, **28**(21):2067–2069.
18. M. Wojtkowski, R. Leitgeb, A. Kowalczyk, T. Bajraszewski, and A. F. Fercher, "In Vivo Human Retinal Imaging by Fourier Domain Optical Coherence Tomography," *Journal of Biomedical Optics*, 2002, **7**(3):457–463.
19. N. A. Nassif, B. Cense, B. H. Park, M. C. Pierce, S. H. Yun, B. E. Bouma, G. J. Tearney, T. C. Chen, and J. F. de Boer, "In Vivo High-Resolution Video-Rate Spectral-Domain Optical Coherence Tomography of the Human Retina and Optic Nerve," *Optics Express*, 2004, **12**(3):367–376.
20. B. Cense, N. Nassif, T. C. Chen, M. C. Pierce, S. H. Yun, B. H. Park, B. E. Bouma, G. J. Tearney, and J. F. de Boer, "Ultrasound-High-Speed Retinal Imaging Using Spectral-Domain Optical Coherence Tomography," *Optics Express*, 2004, **12**(11):2435–2447.
21. M. Wojtkowski, V. J. Srinivasan, T. H. Ko, J. G. Fujimoto, A. Kowalczyk, and J. S. Duker, "Ultrasound-High-Speed, Fourier Domain Optical Coherence Tomography and Methods for Dispersion Compensation," *Optics Express*, 2004, **12**(11):2404–2422.
22. R. A. Leitgeb, L. Schmetterer, W. Drexler, A. F. Fercher, R. J. Zawadzki, and T. Bajraszewski, "Real-time Assessment of Retinal Blood Flow with Ultrafast Acquisition by Color Doppler Fourier Domain Optical Coherence Tomography," *Optics Express*, 2003, **11**(23):3116–3121.
23. B. R. White, M. C. Pierce, N. Nassif, B. Cense, B. H. Park, G. J. Tearney, B. E. Bouma, T. C. Chen, and J. F. de Boer, "In Vivo Dynamic Human Retinal Blood Flow Imaging Using Ultra-High-Speed Spectral Domain Optical Doppler Tomography," *Optics Express*, 2003, **11**(25):3490–3497.
24. M. Wojtkowski, T. Bajraszewski, I. Gorczynska, P. Targowski, A. Kowalczyk, W. Wasilewski, and C. Radzewicz, "Ophthalmic Imaging by Spectral Optical Coherence Tomography," *American Journal of Ophthalmology*, 2004, **138**(3):412–419.
25. T. C. Chen, B. Cense, M. C. Pierce, N. Nassif, B. H. Park, S. H. Yun, B. R. White, B. E. Bouma, G. J. Tearney, and J. F. de Boer, "Spectral Domain Optical Coherence Tomography—Ultrasound-High-Speed, Ultra-High Resolution Ophthalmic Imaging," *Archives of Ophthalmology*, 2005, **123**(12):1715–1720.
26. E. C. W. Lee, J. F. de Boer, M. Mujat, H. Lim, and S. H. Yun, "In Vivo Optical Frequency Domain Imaging of Human Retina and Choroid," *Optics Express*, 2006, **14**(10):4403–4411.
27. H. Lim, M. Mujat, C. Kerbage, E. C. W. Lee, Y. Chen, T. C. Chen, and J. F. de Boer, "High-Speed Imaging of Human Retina In Vivo with Swept-Source Optical Coherence Tomography," *Optics Express*, 2006, **14**(26):12902–12908.
28. A. Unterhuber, B. Povazay, B. Hermann, H. Sattmann, A. Chavez-Pirson, and W. Drexler, "In Vivo Retinal Optical Coherence Tomography at 1040 nm-Enhanced Penetration into the Choroid," *Optics Express*, 2005, **13**(9):3252–3258.
29. D. M. de Bruin, D. L. Burnes, J. Loewenstein, Y. Chen, S. Chang, T. C. Chen, D. D. Esmaili, and J. F. de Boer, "In Vivo Three-Dimensional Imaging of Neovascular Age-Related Macular Degeneration Using Optical Frequency Domain Imaging at 1050 nm," *Investigative Ophthalmology and Visual Science*, 2008, **49**:4545–4552.
30. A. F. Fercher, W. Drexler, C. K. Hitzenberger, and T. Lasser, "Optical Coherence Tomography—Principles and Applications," *Reports on Progress in Physics*, 2003, **66**(2):239–303.
31. A. B. Vakhtin, K. A. Peterson, W. R. Wood, and D. J. Kane, "Differential Spectral Interferometry: An Imaging Technique for Biomedical Applications," *Optics Letters*, 2003, **28**(15):1332–1334.
32. W. V. Sorin and D. M. Baney, "A Simple Intensity Noise-Reduction Technique for Optical Low-Coherence Reflectometry," *IEEE Photonics Technology Letters*, 1992, **4**(12):1404–1406.
33. L. Mandel and E. Wolf, "Measures of Bandwidth and Coherence Time in Optics," *Proceedings of the Physical Society of London*, 1962, **80**(516):894–897.
34. S. H. Yun, G. J. Tearney, B. E. Bouma, B. H. Park, and J. F. de Boer, "High-Speed Spectral-Domain Optical Coherence Tomography at 1.3 μm Wavelength," *Optics Express*, 2003, **11**(26):3598–3604.
35. S. H. Yun, G. J. Tearney, J. F. de Boer, and B. E. Bouma, "Motion Artifacts in Optical Coherence Tomography with Frequency-Domain Ranging," *Optics Express*, 2004, **12**(13):2977–2998.
36. S. H. Yun, G. J. Tearney, J. F. de Boer, and B. E. Bouma, "Pulsed-Source and Swept-Source Spectral-Domain Optical Coherence Tomography with Reduced Motion Artifacts," *Optics Express*, 2004, **12**(23):5614–5624.

37. G. Moneron, A. C. Boccara, and A. Dubois, "Stroboscopic Ultrahigh-Resolution Full-Field Optical Coherence Tomography," *Optics Letters*, 2005, **30**(11):1351–1353.
38. J. W. You, T. C. Chen, M. Mujat, B. H. Park, and J. F. de Boer, "Pulsed Illumination Spectral-Domain Optical Coherence Tomography for Human Retinal Imaging," *Optics Express*, 2006, **14**(15):6739–6748.
39. S. Bourquin, A. Aguirre, I. Hartl, P. Hsiung, T. Ko, J. Fujimoto, T. Birks, W. Wadsworth, U. Bting, and D. Kopf, "Ultrahigh Resolution Real Time OCT Imaging Using a Compact Femtosecond Nd:Glass Laser and Nonlinear Fiber," *Optics Express*, 2003, **11**(24):3290–3297.
40. M. E. J. van Velthoven, M. H. van der Linden, M. D. de Smet, D. J. Faber, and F. D. Verbraak, "Influence of Cataract on Optical Coherence Tomography Image Quality and Retinal Thickness," *British Journal of Ophthalmology*, 2006, **90**(10):1259–1262.
41. S. H. Yun, C. Boudoux, G. J. Tearney, and B. E. Bouma, "Highspeed Wavelength-Swept Semiconductor Laser with a Polygonscanner-Based Wavelength Filter," *Optics Letters*, 2003, **28**:1981–1983.
42. S. L. Jiao, R. Knighton, X. R. Huang, G. Gregori, and C. A. Puliafito, "Simultaneous Acquisition of Sectional and Fundus Ophthalmic Images with Spectral-Domain Optical Coherence Tomography," *Optics Express*, 2005, **13**(2):444–452.
43. M. Mujat, R. C. Chan, B. Cense, B. H. Park, C. Joo, T. Akkin, T. C. Chen, and J. F. de Boer, "Retinal Nerve Fiber Layer Thickness Map Determined from Optical Coherence Tomography Images," *Optics Express*, 2005, **13**(23):9480–9491.
44. X. J. Wang, T. E. Milner, and J. S. Nelson, "Characterization of Fluid-Flow Velocity by Optical Doppler Tomography," *Optics Letters*, 1995, **20**(11):1337–1339.
45. Y. H. Zhao, Z. P. Chen, C. Saxer, S. H. Xiang, J. F. de Boer, and J. S. Nelson, "Phase-Resolved Optical Coherence Tomography and Optical Doppler Tomography for Imaging Blood Flow in Human Skin with Fast Scanning Speed and High Velocity Sensitivity," *Optics Letters*, 2000, **25**(2):114–116.
46. Y. H. Zhao, Z. P. Chen, C. Saxer, Q. M. Shen, S. H. Xiang, J. F. de Boer, and J. S. Nelson, "Doppler Standard Deviation Imaging for Clinical Monitoring of In Vivo Human Skin Blood Flow," *Optics Letters*, 2000, **25**(18):1358–1360.
47. V. Westphal, S. Yazdanfar, A. M. Rollins, and J. A. Izatt, "Realtime, High Velocity-Resolution Color Doppler Optical Coherence Tomography," *Optics Letters*, 2002, **27**(1):34–36.
48. M. R. Hee, D. Huang, E. A. Swanson, and J. G. Fujimoto, "Polarization-Sensitive Low-Coherence Reflectometer for Birefringence Characterization and Ranging," *Journal of the Optical Society of America B-Optical Physics*, 1992, **9**(6):903–908.
49. J. F. de Boer, T. E. Milner, M. J. C. van Gemert, and J. S. Nelson, "Two-Dimensional Birefringence Imaging in Biological Tissue by Polarization-Sensitive Optical Coherence Tomography," *Optics Letters*, 1997, **22**(12):934–936.
50. J. F. de Boer, S. M. Srinivas, A. Malekafzali, Z. P. Chen, and J. S. Nelson, "Imaging Thermally Damaged Tissue by Polarization Sensitive Optical Coherence Tomography," *Optics Express*, 1998, **3**(6):212–218.
51. M. J. Everett, K. Schoenenberger, B. W. Colston, and L. B. Da Silva, "Birefringence Characterization of Biological Tissue by Use of Optical Coherence Tomography," *Optics Letters*, 1998, **23**(3):228–230.
52. M. G. Ducros, J. F. de Boer, H. E. Huang, L. C. Chao, Z. P. Chen, J. S. Nelson, T. E. Milner, and H. G. Rylander, "Polarization Sensitive Optical Coherence Tomography of the Rabbit Eye," *IEEE Journal of Selected Topics in Quantum Electronics*, 1999, **5**(4):1159–1167.
53. G. Yao and L. V. Wang, "Two-Dimensional Depth-Resolved Mueller Matrix Characterization of Biological Tissue by Optical Coherence Tomography," *Optics Letters*, 1999, **24**(8):537–539.
54. X. J. Wang, T. E. Milner, Z. P. Chen, and J. S. Nelson, "Measurement of Fluid-Flow-Velocity Profile in Turbid Media by the Use of Optical Doppler Tomography," *Applied Optics*, 1997, **36**(1):144–149.
55. Z. P. Chen, T. E. Milner, D. Dave, and J. S. Nelson, "Optical Doppler Tomographic Imaging of Fluid Flow Velocity in Highly Scattering Media," *Optics Letters*, 1997, **22**(1):64–66.
56. J. A. Izatt, M. D. Kulkarni, S. Yazdanfar, J. K. Barton, and A. J. Welch, "In Vivo Bidirectional Color Doppler Flow Imaging of Picoliter Blood Volumes Using Optical Coherence Tomography," *Optics Letters*, 1997, **22**(18):1439–1441.
57. A. M. Rollins, S. Yazdanfar, J. K. Barton, and J. A. Izatt, "Realtime In Vivo Color Doppler Optical Coherence Tomography," *Journal of Biomedical Optics*, 2002, **7**(1):123–129.
58. Z. H. Ding, Y. H. Zhao, H. W. Ren, J. S. Nelson, and Z. P. Chen, "Real-Time Phase-Resolved Optical Coherence Tomography and Optical Doppler Tomography," *Optics Express*, 2002, **10**(5):236–245.

59. V. X. D. Yang, M. L. Gordon, B. Qi, J. Pekar, S. Lo, E. Seng-Yue, A. Mok, B. C. Wilson, and I. A. Vitkin, "High Speed, Wide Velocity Dynamic Range Doppler Optical Coherence Tomography (Part I): System Design, Signal Processing, and Performance," *Optics Express*, 2003, **11**(7):794–809.
60. S. Yazdanfar, A. M. Rollins, and J. A. Izatt, "Imaging and Velocimetry of the Human Retinal Circulation with Color Doppler Optical Coherence Tomography," *Optics Letters*, 2000, **25**(19):1448–1450.
61. S. Yazdanfar, A. M. Rollins, and J. A. Izatt, "In Vivo Imaging of Human Retinal Flow Dynamics by Color Doppler Optical Coherence Tomography," *Archives of Ophthalmology*, 2003, **121**(2):235–239.
62. M. D. Kulkarni, T. G. van Leeuwen, S. Yazdanfar, and J. A. Izatt, "Velocity-Estimation Accuracy and Frame-Rate Limitations in Color Doppler Optical Coherence Tomography," *Optics Letters*, 1998, **23**(13):1057–1059.
63. R. Leitgeb, L. F. Schmetterer, M. Wojtkowski, C. K. Hitzenberger, M. Sticker, and A. F. Fercher, "Flow Velocity Measurements by Frequency Domain Short Coherence Interferometry," *Proceedings of SPIE*, 2002, 4619.
64. J. F. de Boer, C. E. Saxer, and J. S. Nelson, "Stable Carrier Generation and Phase-Resolved Digital Data Processing in Optical Coherence Tomography," *Applied Optics*, 2001, **40**(31):5787–5790.
65. C. K. Hitzenberger, E. Gotzinger, M. Sticker, M. Pircher, and A. F. Fercher, "Measurement and Imaging of Birefringence and Optic Axis Orientation by Phase Resolved Polarization Sensitive Optical Coherence Tomography," *Optics Express*, 2001, **9**(13):780–790.
66. C. E. Saxer, J. F. de Boer, B. H. Park, Y. H. Zhao, Z. P. Chen, and J. S. Nelson, "High-Speed Fiber-Based Polarization-Sensitive Optical Coherence Tomography of In Vivo Human Skin," *Optics Letters*, 2000, **25**(18):1355–1357.
67. B. H. Park, C. Saxer, S. M. Srinivas, J. S. Nelson, and J. F. de Boer, "In Vivo Burn Depth Determination by High-Speed Fiberbased Polarization Sensitive Optical Coherence Tomography," *Journal of Biomedical Optics*, 2001, **6**(4):474–479.
68. M. C. Pierce, B. H. Park, B. Cense, and J. F. de Boer, "Simultaneous Intensity, Birefringence, and Flow Measurements with High-Speed Fiber-Based Optical Coherence Tomography," *Optics Letters*, 2002, **27**(17):1534–1536.
69. B. H. Park, M. C. Pierce, B. Cense, S. H. Yun, M. Mujat, G. J. Tearney, B. E. Bouma, and J. F. de Boer, "Real-time Fiberbased Multifunctional Spectral-Domain Optical Coherence Tomography at 1.3 μm ," *Optics Express*, 2005, **13**(11):3931–3944.
70. S. L. Jiao, G. Yao, and L. H. V. Wang, "Depth-Resolved two-dimensional Stokes Vectors of Backscattered Light and Mueller Matrices of Biological Tissue Measured with Optical Coherence Tomography," *Applied Optics*, 2000, **39**(34):6318–6324.
71. S. L. Jiao and L. H. V. Wang, "Jones-matrix Imaging of Biological Tissues with Quadruple-Channel Optical Coherence Tomography," *Journal of Biomedical Optics*, 2002, **7**(3):350–358.
72. S. L. Jiao and L. H. V. Wang, "Two-dimensional Depth-Resolved Mueller Matrix of Biological Tissue Measured with Double-Beam Polarization-Sensitive Optical Coherence Tomography," *Optics Letters*, 2002, **27**(2):101–103.
73. S. L. Jiao, W. R. Yu, G. Stoica, and L. H. V. Wang, "Optical-fiber-Based Mueller Optical Coherence Tomography," *Optics Letters*, 2003, **28**(14):1206–1208.
74. B. H. Park, M. C. Pierce, B. Cense, and J. F. de Boer, "Jones Matrix Analysis for a Polarization-Sensitive Optical Coherence Tomography System Using Fiber-Optic Components," *Optics Letters*, 2004, **29**(21):2512–2514.
75. J. Strasswimmer, M. C. Pierce, B. H. Park, V. Neel, and J. F. de Boer, "Polarization-Sensitive Optical Coherence Tomography of Invasive Basal Cell Carcinoma," *Journal of Biomedical Optics*, 2004, **9**(2):292–298.
76. B. Cense, T. C. Chen, B. H. Park, M. C. Pierce, and J. F. de Boer, "In Vivo Depth-Resolved Birefringence Measurements of the Human Retinal Nerve Fiber Layer by Polarization-Sensitive Optical Coherence Tomography," *Optics Letters*, 2002, **27**(18):1610–1612.
77. B. Cense, T. C. Chen, B. H. Park, M. C. Pierce, and J. F. de Boer, "Thickness and Birefringence of Healthy Retinal Nerve Fiber Layer Tissue Measured with Polarization-Sensitive Optical Coherence Tomography," *Investigative Ophthalmology and Visual Science*, 2004, **45**(8):2606–2612.
78. D. Fried, J. Xie, S. Shafi, J. D. B. Featherstone, T. M. Breunig, and C. Le, "Imaging Caries Lesions and Lesion Progression with Polarization Sensitive Optical Coherence Tomography," *Journal of Biomedical Optics*, 2002, **7**(4):618–627.
79. R. C. Jones, "A New Calculus for the Treatment of Optical Systems I. Description and Discussion of the Calculus," *Journal of the Optical Society of America A*, 1941, **31**(7):488–493.

80. J. J. Gil, and E. Bernabeu, "Obtainment of the Polarizing and Retardation Parameters of a Nondepolarizing Optical System from the Polar Decomposition of its Mueller Matrix," *Optik*, 1987, **76**(2):67–71.
81. K. Schoenenberger, B. W. Colston, D. J. Maitland, L. B. Da Silva, and M. J. Everett, "Mapping of Birefringence and Thermal Damage in Tissue by Use of Polarization-Sensitive Optical Coherence Tomography," *Applied Optics*, 1998, **37**(25):6026–6036.
82. G. J. van Blokland, "Ellipsometry of the Human Retina In Vivo: Preservation of Polarization," *Journal of the Optical Society of America A*, 1985, **2**:72–75.
83. H. B. K. Brink and G. J. van Blokland, "Birefringence of the Human Foveal Area Assessed In Vivo with Mueller-Matrix Ellipsometry," *Journal of the Optical Society of America A*, 1988, **5**(1):49–57.
84. W. K. Tung, *Group Theory in Physics*, 1985, Singapore World Scientific.
85. B. H. Park, M. C. Pierce, and J. F. de Boer, "Comment on Optical-Fiber-Based Mueller Optical Coherence Tomography," *Optics Letters*, 2004, **29**(24):2873–2874.
86. B. H. Park, M. C. Pierce, B. Cense, and J. F. de Boer, "Realtime Multifunctional Optical Coherence Tomography," *Optics Express*, 2003, **11**(7):782–793.
87. B. Cense, H. C. Chen, B. H. Park, M. C. Pierce, and J. F. de Boer, "In Vivo Birefringence and Thickness Measurements of the Human Retinal Nerve Fiber Layer Using Polarization-Sensitive Optical Coherence Tomography," *Journal of Biomedical Optics*, 2004, **9**(1):121–125.
88. M. Pircher, E. Gotzinger, R. Leitgeb, H. Sattmann, O. Findl, and C. K. Hitzenberger, "Imaging of Polarization Properties of Human Retina In Vivo with Phase Resolved Transversal PS-OCT," *Optics Express*, 2004, **12**(24):5940–5951.
89. E. Gotzinger, M. Pircher, and C. K. Hitzenberger, "High Speed Spectral Domain Polarization Sensitive Optical Coherence Tomography of the Human Retina," *Optics Express*, 2005, **13**(25):10217–10229.
90. M. Mujat, B. H. Park, B. Cense, T. C. Chen, and J. F. de Boer, "Autocalibration of Spectral-Domain Optical Coherence Tomography Spectrometers for In Vivo Quantitative Retinal Nerve Fiber Layer Birefringence Determination," *Journal of Biomedical Optics*, 2007, **12**(4).
91. H. Lim, J. F. de Boer, B. H. Park, E. C. W. Lee, R. Yelin, and S. H. Yun, "Optical Frequency Domain Imaging with a Rapidly Swept Laser in the 815-870 nm Range," *Optics Express*, 2006, **14**(13):5937–5944.

GRADIENT INDEX OPTICS IN THE EYE

Barbara K. Pierscionek

*Department of Biomedical Sciences
University of Ulster
Coleraine, United Kingdom*

19.1 GLOSSARY

Cataract. Opacification in the eye lens caused by changes in protein shape, structure, or interaction that causes light to be scattered or absorbed.

Crystallin proteins. The major structural proteins in the eye lens.

Equatorial plane. The plane of the eye lens that contains the lens equator.

Gradient index or GRIN. The property of a lens or medium that has a gradually varying refractive index.

Homogenous index. The property of a lens or medium that has a constant refractive index.

Isoindicial contours. Contours of constant refractive index.

Magnetic resonance imaging (MRI). Noninvasive method used for imaging body organs that utilizes the interaction of a static magnetic field, radio frequency pulses, and magnetic gradients to form images.

Myopia. Shortsightedness.

Optic fiber. A fiber made of glass or plastic that guides light along its length by total internal reflection.

Optical fiber perform. The precursor to an optical fiber: a rod with the refractive index distribution of the desired optical fiber that is stretched to form the fiber.

Raman microspectroscopy. A technique used to detect constituents of a system from the spectrum of molecular vibrations created by specific interactions of monochromatic light with matter. These interactions cause scattering of photons at lower frequencies than that of the incident light.

Ray tracing. Monitoring the path of rays as they traverse a lens or optical system. This is often used to measure a particular property of the lens or system.

Reflectometric sensor. A device that measures optical or material properties of a medium from the proportion of light that the medium reflects.

Refractive index. The property of a medium that contributes to its capacity to bend light and that is related to the density of the material.

Refractive power. The capacity of a lens or optical system to bend and focus light.

Sagittal plane. The vertical plane of the eye lens that contains the optic axis.

19.2 INTRODUCTION

This chapter considers the common forms of gradient index in optical media. The lens of the eye has a gradient of index and this is examined in detail; the form of the index profile as well as the magnitude of refractive index are described. The refractive index distribution has been measured or derived from a number of animal eyes as well as from the human eye. The different types of gradients are presented and compared.

19.3 THE NATURE OF AN INDEX GRADIENT

The refractive power of a lens is determined by its shape, the refractive index of its medium, and the refractive index of the medium that surrounds the lens. In terms of shape, the more curved the lens the greater is the refraction of light. The relationship between refractive index and refractive power is not as direct: light is not necessarily refracted to a greater degree in a higher-index medium. A ray of light will travel in a straight line whether it travels in air (refractive index = 1), in water (refractive index = 1.33), or in glass (refractive index around 1.5 to 1.6). The effect of refractive index on refraction occurs when there is a change in refractive index at the interface between two media or a gradual variation in index within a medium. The greater the difference between refractive indices at an interface or the steeper the gradient of refractive index in a medium, the greater the degree of refraction.

Gradient index (GRIN) deals with media which have a varying refractive index. The effect of varying or gradient index occurs in nature; changes in temperature or pressure can induce fluctuations or gradations in the refractive index of air. One of the most commonly cited examples of this is seen in the shiny reflective marks on a road that can be seen on a hot day. The heat of the road warms the air closest to it and gradients of temperature and of refractive index are formed. The higher the temperature, the less dense the air and the lower the refractive index. The index gradient has the effect of deflecting light that hits the surface of the road at certain angles so that it appears to have been reflected from the surface of the road.

The steepness of the index gradient controls the degree of bending of light in the GRIN medium and this is used in the creation of GRIN lenses and devices. A refractive contribution from the medium is particularly effective when the surface of the given structure is flat and cannot contribute to the refractive power.

19.4 SPHERICAL GRADIENTS

In a lens with a spherical refractive index distribution, the index varies with distance from the center of the lens. The most famous early GRIN lens was Maxwell's fish-eye lens¹ that has a refractive index distribution with spherical symmetry about a point and can be described as having the form

$$n(r) = \frac{n_0}{1 + (r/a)^2} \quad (1)$$

where r is the distance from the center, $n(r)$ is the refractive index at r , n_0 is the refractive index at the center of the lens, and a is a constant.

The special property of a fish-eye lens is that all rays emanating from a given object point pass through the same image point.² This produces sharp imagery but only for those points on the surface and within the lens. There is no evidence that fish possess such lenses. The name possibly derives from the challenge posed by the Irish Academy of Sciences: to produce a medium with a refractive index distribution that was capable of forming images with the least depth possible. As the eye of the fish is

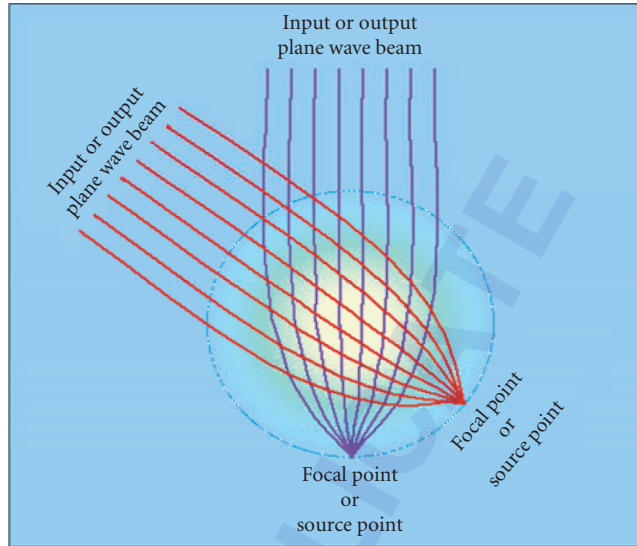


FIGURE 1 Light paths through a Luneberg lens in surround media index matched to the lens surface.

known to be relatively flat, Maxwell's solution to the problem was called the fish-eye.³ The limitation of this lens is that sharp imagery of object points only occurs for those points that lie within or on the surface of the lens.^{2,4}

This restriction does not apply to the Luneberg lens,⁵ which, like the Maxwell fish-eye lens, has a refractive index distribution that is spherically symmetrical. This lens is a sphere with a refractive index distribution described by

$$n(r) = (2 - r^2)^{1/2} \quad (2)$$

where r is the distance from the center of the lens and $n(r)$ is the refractive index at r .

It has no intrinsic optical axis but focuses parallel rays, incident on any part of its surface, to a point on the opposite surface of the lens (Fig. 1).

The restriction of the Luneberg lens is that with this lens sharp images are only created for objects at infinity. In order to ensure that rays enter in a parallel bundle, the index of the surrounding media must match that of the lens surface. In air this would require an outer surface refractive index of 1, which is difficult to achieve for visible light but can be utilized for frequencies in the microwave range.² It has wide ranging applications in microwave antennas and as a radar device. A very useful modern application is in the design of broadband airborne antennas. These allow passengers on an aircraft to access the Internet, email, satellite television, and video.

19.5 RADIAL GRADIENTS

When the refractive index distribution is cylindrically symmetrical, that is, the refractive index varies with distance from a fixed line, it is described as a radial index gradient. A simple type of radial index lens is the Wood lens, named after its designer.⁶ This lens has a radial refractive index gradient and flat surfaces.² In such a lens, parallel rays from objects at infinity are refracted purely by the

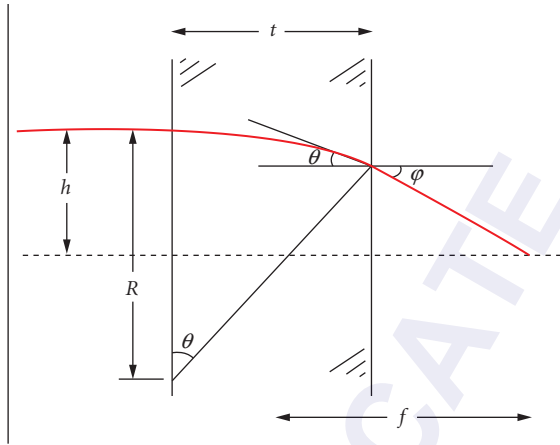


FIGURE 2 Refraction of a ray through a Wood lens. $R, f \gg h, t$. [From N. Morton, "Gradient Refractive Index Lenses," *Phys. Educ.* **19**:86–90 (1984).]

index distribution within the lens (Fig. 2). When the refractive index decreases with distance from the axis, the lens acts as a converging lens; when the refractive index increases with distance from the axis, the lens diverges light. With a Wood lens, the object and image planes are external to the lens.²

Optical fiber cores with radial index gradients are, cross-sectionally, akin to a Wood lens. However, because optical fibers are elongated in the axial direction and have very small diameters, there is an additional periodicity to the ray path. In a graded index optical fiber core, the index distribution follows a general form⁷

$$n(r) = n_0[1 - 2(r/a)^\alpha \Delta]^{1/2} \tag{3}$$

where r is the radial distance from the center, $n(r)$ is the refractive index at r , n_0 is the refractive index along the axis of the fiber, a is the fiber radius, and Δ is the fractional refractive index difference between the center and edge of the fiber core.

For $\alpha = 2$ and $\Delta \ll 1$, the profile is parabolic and approximated as

$$n(r) = n_0[1 - (r/a)^2 \Delta] \tag{4}$$

In an optical fiber, the core is surrounded by a single index cladding to minimize light losses along the length of the fiber and to protect the quality of the fiber core. Within the core, the rays follow a sinusoidal path along the axis of the fiber (as seen in Fig. 3).

Optical fibers have many uses particularly when the light needs to be sent some distance from the source or an image needs to be collected from a remote or hard to access location. Optical fibers have replaced wires in communication systems. They are now used routinely in medical instruments and applications, such as endoscopes and sensors.

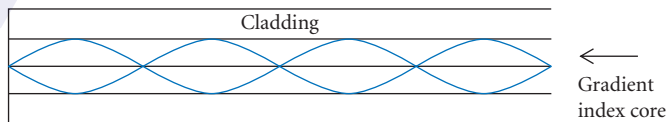


FIGURE 3 Light rays follow a sinusoidal path when traveling in the gradient index core of a GRIN fiber.

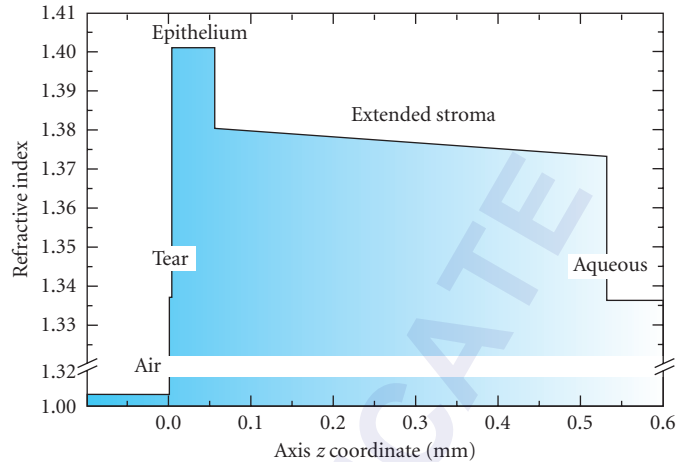


FIGURE 4 Variations in corneal refractive index along the optical axis. (From Ref. 10.)

19.6 AXIAL GRADIENTS

Refractive index can also vary as a function of distance from a reference plane. This is most often the plane perpendicular to the optic axis and in such case the refractive index varies along the optic axis. Such a gradient can be formed by thin films and has important use in the design of multilayer antireflection coatings.^{8,9} Southwell⁹ has shown that shape of the index gradient along the axis is a very important determinant that can be used to minimize adverse reflections. Measurement of axial gradients has application in the study of biological tissues, many of which are layered. The lamellar stroma, which constitutes the bulk of the cornea of the eye and has been considered to have a homogenous refractive index, has recently been modeled as an axial index structure¹⁰. The axial index model has been found to produce more spherical aberration, for larger pupil sizes, than homogenous index models. Any refractive index differences along the axis of the cornea are likely to be very small and the model introduced very slight variations (Fig. 4).

19.7 THE EYE LENS

The refractive index of the eye lens has a radial as well as an axial gradient. This is because of the three-dimensional structure as well as the nature of the growth mode of the lens. The lens grows by tissue accretion in layers on the lens surface, with no concomitant loss of existing cells, resulting in a layered structure. The refractive index varies across the layers resulting in a radial gradient, in the equatorial plane and an axial gradient along the optic axis. In spherical lenses (seen in some fish), the index gradient varies along a radial distance from the center in all directions.

The shapes and sizes of eye lenses, like any other part of the eye, vary widely across species. This variation bears no relationship to the size of the animal nor to the size of the eye. For example, the lens of a rabbit is larger than that of a human, while the rat's lens is highly curved and fills the eye almost entirely. With such a variety of lenticular shapes and sizes, it is not surprising that the refractive index distribution also varies between species and is not related to the size of the animal or to its eye but to its visual needs.

The refractive index gradients in human and several types of animal lenses have been measured or derived using various techniques. These have included ray tracing, optic fiber sensing, or calculation from protein and water concentration measurements and MRI scanning. Earlier studies on dispersion have provided refractive index measurements on inner and outer sections of lenses from a number of vertebrates.^{11,12}

19.8 FISH

The fish lens contributes more to the refractive power of the eye than lenses from land living animals because the corneal power in the fish eye is significantly reduced by the very small difference in refractive index between the cornea and water. Fish lenses are spherical or nearly spherical with steep refractive index gradients and high maximal index values to provide the necessary refractive power to the eye whilst minimizing spherical aberration.^{13–20} The refractive index gradients in fish lenses have been fitted with elliptical,¹⁴ parabolic,^{16–19} and sixth-order polynomial functions (Table 1).²⁰

In spite of the differences in the size of fish and their lenses, the maxima of refractive index in the lens centers are very close in value and, with the exception of the Blue Eye Trevally lens, similar values are also found at the lens edge (Table 1). The edge index value in the trevally lens is substantially higher than that found in the other fish lenses. It is also notable that the trevally has an enormous eye relative to its size and its lens is the largest of any lens studied: a diameter of 26.6 mm and a weight of over 9.5 g.¹⁷ It is a deep-sea fish, unlike the other fish listed in Table 1. From a structural aspect, almost a third of the trevally lens protein comprises γ -crystallin.¹⁷ This crystallin protein contributes more to the refractive index of the lens than the other crystallins²¹ and is capable of closer protein packing with less surrounding water.²² Although the proportion of γ -crystallin is comparatively high, it is not high enough to alone account for the refractive index magnitude. This may also be attributed to the very high proportion of insoluble protein that was found on extraction.¹⁷ This does not necessarily indicate that the proteins are in an insoluble state in the intact trevally lens but may suggest a propensity for proteins in this lens to aggregate with no loss of transparency. This will depend on protein types and relative distributions as well as protein-protein and protein-water interactions. It is not known whether the high proportions of γ -crystallin and insoluble protein found in the trevally lens apply to the lenses of other fish.

Ageing changes in the fish lens have not been studied but the effect on the index gradient with size and therefore growth has been investigated. Kroger et al.¹⁹ reported that the refractive index in the center of the African Cichlid fish lens increases with lens size. The resultant steepening of the index gradient with growth causes a decrease in focal length. There is also a concomitant decrease in spherical aberration that is attributed to a smoothing of the index gradient.¹⁹ It has been suggested that depressions in the index profile may have a discrete function in producing multiple focal points.¹⁹

Jagger¹⁵ compared models of fish lenses and found that a higher-order polynomial function to describe the index gradient provided results that were most akin to measurements for fish lenses in terms of image quality, spherical, and chromatic aberration. It has also been pointed out that there may be internal asymmetries in the index distribution in fish lenses that would render erroneous measurements based on ray tracing assumptions.²³

TABLE 1 Details of Refractive Index Profiles in a Variety of Fish Lenses

Fish	Refractive Index (max)	Refractive Index (min)	Wavelength (nm)	Gradient Shape
Goldfish ¹⁴	1.55–7	1.35–8	633	Elliptical
African Cichlid ¹⁶	1.540	1.361	633	Parabolic
Blue Eye Trevally ¹⁷	1.546	1.408	633	Parabolic
Trout ¹⁸	1.538	1.3716	550	Parabolic
Black Dory ²⁰	1.543	1.368	590	Sixth-order polynomial

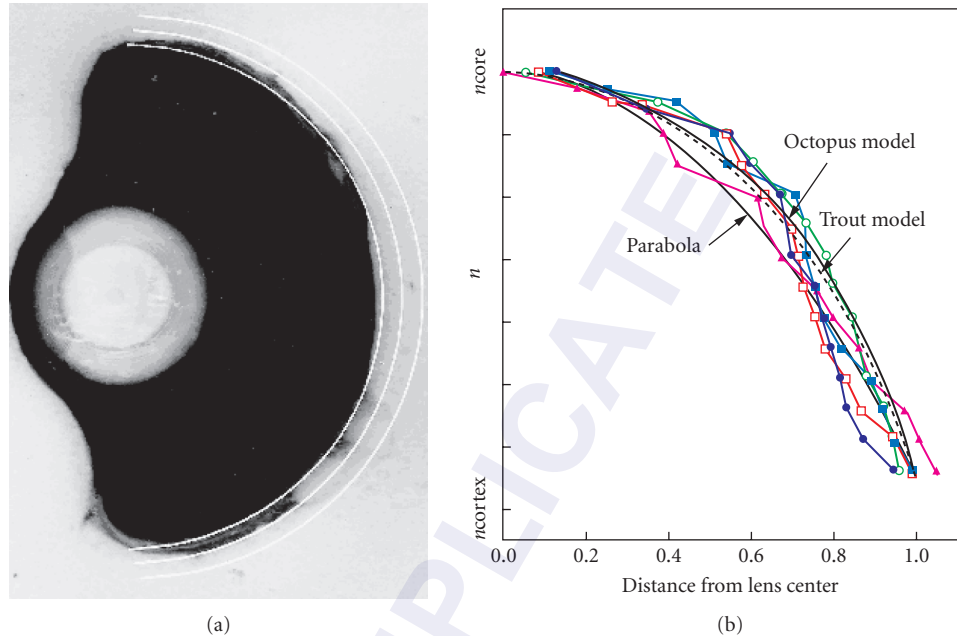


FIGURE 5 (a) Octopus eye showing lens and (b) refractive index gradient for five octopus lenses as measured in the equatorial plane. Model and experimental index profiles are shown together with the index profile for the trout lens. (From Ref. 24.)

19.9 OCTOPUS

The eye of the octopus has no cornea; the refractive power of the eye is provided by the lens (Fig. 5a).²⁴ It was found to have concentric isoindicial contours and an index gradient that could be fitted to a 10th order polynomial (Fig. 5b). The maximum refractive index was 1.509 decreasing to 1.357 for 550 nm.

19.10 RAT

The refractive index distribution in the rat lens was measured using ray tracing^{25,26} following the method for determining the refractive index in optical fiber performs devised by Chu.²⁷ This method enabled the calculation of the gradient refractive index profile using, as parameters, the distance of the parallel incident rays from the optic axis and the corresponding angles formed by the exiting rays and the optic axis.

The studies of Campbell and Hughes²⁵ and Campbell²⁶ showed that the rat has the steepest refractive index gradient of all species yet measured, rising from around 1.38 at the edge of the lens to a maximum of around 1.5 at the lens center in a profile that can be fitted to a parabolic form. This rise occurs over a very short distance of 2 mm.²⁶ A refractive index value of 1.50 is approaching the refractive index range for glass and, in terms of protein concentration, is close to the refractive index of pure protein.²⁸ With such a high protein concentration, the rat lens cannot be deformed, and hence the rat eye has no accommodative capacity. The high magnitude of refractive index in the rat lens is

partially caused by the relatively high proportion of the protein γ -crystallin. In terms of shape, the rat lens is not quite spherical but tends to have a longer equatorial diameter than sagittal radius with little variation between the anterior and posterior curvatures.²⁶ The radial and axial gradient profiles will not be the same, the latter slightly less steep than the former. The curved nature of the rat lens, its high magnitude of refractive index, and the steep index gradient gives the lens a very high refractive power rendering the rat eye highly myopic. This accords with the animal's visual needs which are for very close distances and do not require a range of focus.

19.11 GUINEA PIG

The refractive index distribution of the guinea pig lens, measured with an optical fiber reflectometric sensor, has been found to be parabolic.²⁹ This was for a single lens, and the age of the animal was not known. Hence magnitudes of refractive index: a maximum value of 1.429 at the lens center and a minimum value of 1.389 at the edge can only be taken as applicable for that particular lens. Calculated protein content for the guinea pig lens was 48 percent²⁹ and the lens shape is almost spherical with a sagittal width of around 5 mm and an equatorial radius of 5.4 mm. The refractive index distribution will therefore be similar in the axial and radial directions. Because the guinea pig lens is relatively small, its refractive index gradient is quite steep and this adds to its refractive power. The steep distribution of refractive index may be linked to the structural proteins and the uncommon crystallin protein (ζ -crystallin) that comprises about 10 percent of the total soluble proteins.³⁰ This protein differs from the other major crystallins in structure^{30,31} and in its response to the effects of heat and the protective action of α -crystallin.³² ζ -crystallin requires a tenfold greater amount (w/w) of α -crystallin than the other crystallin proteins to prevent it forming aggregates and disturbing the smoothness of the refractive index gradient.³² The greater level of protection required by ζ -crystallin could indicate a propensity to exist in regions of relatively higher protein concentration. The relatively steep refractive index gradient and the high curvature of the guinea pig lens are needed to meet the refractive requirements of the guinea pig: clear vision at close distances.

19.12 RABBIT

The refractive index gradient in the rabbit lens was measured in the equatorial and horizontal planes on thin slices.³³ Isoindicial curves were found as concentric circles in the equatorial plane and as eccentric ellipses in the horizontal plane, the eccentricity arising because of differences in curvature between the anterior and posterior lens surfaces. The three-dimensional refractive index distribution was represented by an elliptical paraboloid

$$n = n_0 - [Z^2/A^2 + (X^2 + Y^2)/B^2] \quad (5)$$

where n_0 is the refractive index at the lens center, X is the lateral axis (through the middle of the horizontal plane), Y is the vertical axis (through the middle of the equatorial plane), Z is the optic axis, A and B are constants (A differs in the anterior and posterior sections of the lens).

The refractive index for 589 nm was reported to vary from 1.464 to 1.431 at the lens center and 1.393 to 1.384 at the surface, which is not very different from the magnitudes of refractive index in the guinea pig lens. However, the rabbit lens is larger than that of the guinea pig (around 10 mm thickness in the horizontal and sagittal planes) and so the index gradient in the rabbit lens would not be as steep providing a lower contribution to refractive power than does the index gradient in the guinea pig lens.

19.13 CAT

The lens of the cat is shaped as a biconvex ellipsoid³⁴ and its refractive index follows a parabolic shape with a maximum of 1.4813 at the center and minimum of 1.3877 in the periphery (values averaged from five cat lenses) for a wavelength 550 nm. The parabolic shape was found in both the sagittal (axial) and equatorial planes and the index variations follow isoindicial contours that retain the geometric shape of the lens.³⁴

Jagger³⁴ cites Matthiessen (1880),³⁵ who showed the index in a homogenous lens would give equivalent power to one with a variation of index (based on a parabolic gradient). This rule states that the difference between the total index and the core (maximum) index in a grin lens is equal to the difference between the core index and cortical (edge) index

$$n(\text{total}) - n(\text{core}) = n(\text{core}) - n(\text{cortex})$$

so

$$n(\text{total}) = 2 \times n(\text{core}) - n(\text{cortex})$$

According to the rule of Matthiessen,³⁵ the equivalent to a lens with a parabolic gradient index, with a maximum value of 1.48 and an edge value of 1.38, would require a homogenous index lens to have a refractive index of 1.58. This is higher than the refractive index of pure protein.²⁸

19.14 BOVINE

The optical properties of the bovine lens have been investigated in a number of studies and the refractive index has been measured using ray tracing^{36,37} and fiber optic reflectometric sensing methods.³⁸ In addition, the refractive index in the bovine lens has been measured in both equatorial and sagittal planes and in lenses from a wide age range from foetal³⁹ to approximately 15 years of age.³⁶

The refractive index, measured in the equatorial plane using a ray tracing technique for a wavelength of 633 nm^{36,37} was found to have some variations depending on the age of the lens. In all lenses the gradient followed the same general shape: that of a second-order polynomial, but the index magnitude increased across the gradient with age. The size and therefore the radius of the lens in the equatorial plane had to be incorporated into the equation to account for size and age differences. The general form of the refractive index gradient covering the wide range of lenses investigated is

$$n(r) = a_r - N_r/\rho^2 + C \quad (6)$$

where r is the radial distance from the center of the lens, $n(r)$ is the refractive index at r (mm), a_r is the refractive index at the center of the oldest lens investigated (approximately 15 years of age),³⁶ N_r is a non-normalized gradient function, ρ is the equatorial radius of the lens (mm), and C is a correction factor for age-related differences in the refractive index at the lens center.

With

$$N_r = 0.702 - 0.471r + 0.186r^2 \quad (7)$$

and

$$C = -0.0025 - 0.00583r + 0.00131r^2 \quad (8)$$

the final form of the equation is

$$n(r) = 1.4575 - 0.702/\rho^2 + 0.471r/\rho^2 - 0.186r^2/\rho^2 - 0.00583r + 0.00131r^2 \quad (9)$$

In a subsequent study the refractive index along the optic axis of the bovine lens was measured using a fiber optic sensor and a wavelength of 670 nm.³⁸ The lens age range was not as great as in the ray tracing study of the equatorial plane index: weights varied from 1.676 to 2.646 g corresponding to ages of around 1 to 8 years, respectively.⁴⁰ A second-order polynomial was again found to fit the profiles but here slight differences were found in the anterior and posterior sections of the lens. The anterior axial profile was found to fit the following equation

$$n(x) = 1.3790 + 0.018484x - 0.0014925x^2 \quad (10)$$

and the posterior axial profile could be described by

$$n(x) = 1.4102 + 0.009647x - 0.0010741x^2 \quad (11)$$

where x is the distance from the anterior pole (mm) and $n(x)$ is the refractive index at r .

It is not possible to definitively state whether the bovine lens does have contours of refractive index that follow the external shape of the lens but the forms of the gradients in the sagittal and equatorial planes are similar and using dimensional data from the lens it is possible to do approximate calculations to see if the magnitudes of refractive index following a given contour are close. If so, the assumption of concentric contours of refractive index may apply to the bovine lens. From the dimensions of bovine lenses⁴¹ a lens of about 1.5 g weight has an equatorial radius of 7.5 mm, an anterior sagittal thickness of 5 mm and a posterior sagittal thickness of 6 mm.

Taking a contour at 2 mm from the anterior pole, this will be at a radial distance of 4.5 along the equatorial plane (2 mm is 40 percent of the anterior sagittal thickness and 4.5 mm is 40 percent of the equatorial radius) and 2.4 mm from the posterior pole (8.6 mm from the anterior pole). Putting these values into Eqs. (9), (10), and (11) gives refractive index values of 1.4099 (~1.41) along the anterior sagittal thickness, 1.416 along the equatorial plane, and 1.414 along the posterior sagittal thickness, respectively. Taking into account differences in measurement technique, the fact that measurements of radial and axial gradients were conducted on different lenses, and the recognition that in biological tissues there will always be individual variations in lens shape and size, the above calculations can only provide an estimate. This notwithstanding, the values are close and suggest that in the bovine lens isoindicial contours of refractive index may follow the surface shape.

The bovine lens is one of few in which age-related changes to the GRIN optics have been investigated. The changes with age to the index distribution should take into account the fact that the lens continues to grow, and are more correctly described as changes with time and not exclusively occurring as a result of senescence. The lens grows by successive tissue layering and, within each layer, the cells stretch from anterior to posterior pole. As there is an asymmetry between the anterior and posterior surface curvatures, the cells must stretch more and more asymmetrically with growth, and a greater component of the cellular length is found in the posterior part of the lens. How the cellular contents, mainly the proteins, are distributed along the cell length is not known. Nevertheless, as the anterior and posterior refractive index gradients meet in the equatorial plane, the refractive index at that point must be the same. Similarly, refractive indices at the anterior and posterior poles should be close in value, as these are from the same cell layers, and large difference in refractive index between the poles would not meet the refractive needs of the eye. Fiber optic sensing has shown that there is, indeed, very little difference in the refractive index at the anterior and posterior poles of lenses that have a relatively lower weight and were therefore likely to be younger.³⁸ However, whilst the refractive index does not appear to alter with age at the anterior pole, at the posterior pole a statistically significant correlation with age was found.³⁸ The consequence of this is that in very old bovine lenses, some differences in polar refractive index may occur. This is most likely to be a consequence of the growth mode and increasing asymmetry of shape that may affect and result in greater differences in anterior and posterior refractive index gradients in advanced age.

As both axial and radial gradients were measured with similar wavelength sources (633 nm for the radial gradient and 670 nm for the axial) wavelength differences would have a negligible impact on the results. Comparison of ray tracing and fiber optic sensing have shown consistent findings.^{36,38}

A single study on development of the foetal bovine lens³⁹ has shown that the refractive index gradient starts to take on its parabolic form around the middle of the gestational period of 9.3 months. Prior to this, the refractive index distribution has no defined shape (Fig. 6).

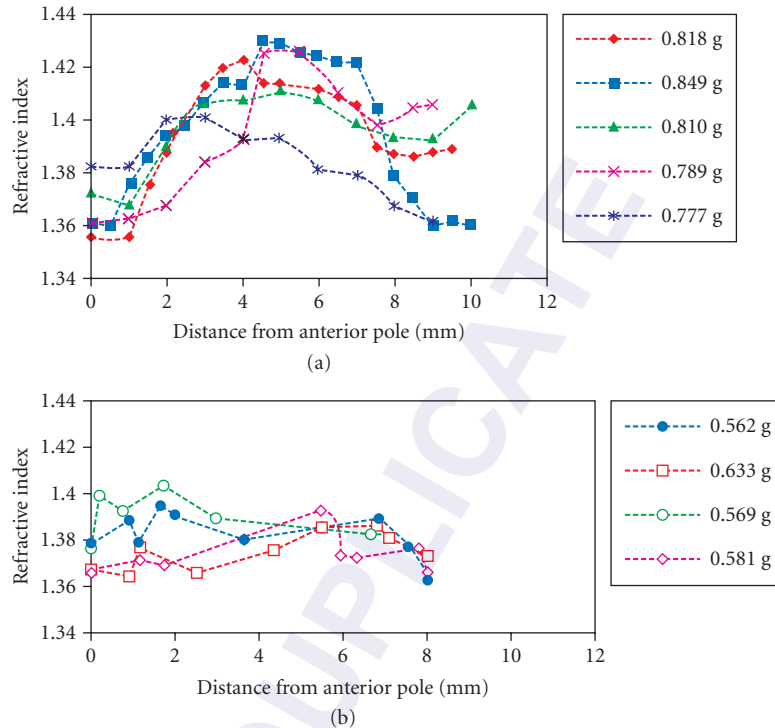


FIGURE 6 Refractive index profiles along optic axes of foetal bovine lenses. (a) Greater than 0.7 g and (b) lesser than 0.7 g wet weight. (From Ref. 39.)

19.15 PIG

The refractive index in the porcine lens is similar in shape to that found in the bovine lens.^{42,43} It has been fitted with a second-order polynomial, a shape that has been found from ray tracing measurements of refractive index⁴² (as shown as Fig. 7) and from protein concentration profiles determined using Raman microspectroscopy.⁴⁴

Jones and Pope⁴³ using magnetic resonance imaging (MRI), found that parabolic, power function, and higher-order polynomial fits could be used to provide a very good description of the index profile. They reported that the power function provided the best fit but this was only by a very slight margin. Comparing the results of Pierscionek et al.⁴² who measured the refractive index profile in the equatorial plane with those of Jones and Pope,⁴³ whose measurements were in the sagittal plane, there is very little difference in profile shapes. This suggests that there may be isoindicial contours of refractive index in the porcine lens that follow the surface shape. The porcine lens index gradient has been measured for wavelengths 633 and 532⁴² and 589 nm (converted to 543.5 nm).⁴³ The index magnitude was reported to vary from an average of 1.354 at the edge of the lens up to 1.396 at the center for 633 nm⁴² from 1.365 to 1.368 at the edge up to an average of 1.4374 at the center for 543.5 nm,⁴³ and from 1.358 at the edge to 1.404 at the center for 532 nm. The central index value found by Jones and Pope⁴³ concurs with measurements made by Sivak and Mandelman.¹²

The higher values of refractive index found by Jones and Pope⁴³ compared to Pierscionek et al.⁴² could have reflected differences in ages of animals; it is not possible to determine whether the two studies^{42,43} had used lenses of similar ages. Age-related variations in refractive index have not been studied in the porcine lens.

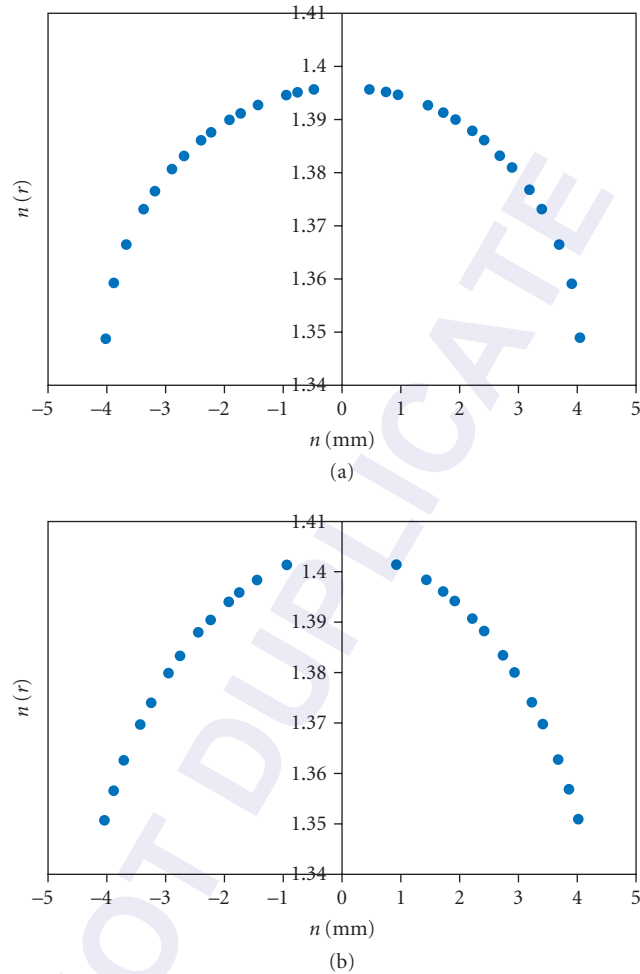


FIGURE 7 Refractive index profiles measured in the equatorial plane of porcine lenses shown against radial distance from the center of the lens for (a) 633 nm and (b) 532 nm light. (From Ref. 42.)

19.16 HUMAN/PRIMATE

The refractive index gradient of the human lens differs from index gradients in the lenses of most other species. The gradient of refractive index is only found in the cortex of the lens, which is generally taken to represent the outer one-third of the lens, approximately, in terms of radial distance. In the inner, nuclear region, the refractive index is relatively constant around a value of 1.39 to 1.41^{12,45,46} (depending on wavelength) and remains so with age as long as there are no opacifications or cataractous changes.⁴⁷ This supports findings on protein⁴⁸ and water gradients.^{49,50} Changes with age occur in the steepness of the cortical index gradient and have been proposed as an answer to why the human lens, in spite of continued growth and increased curvature^{51,52} does not become myopic.^{29,53} A recent study using MRI has suggested that there is a decrease, with age, in the nuclear

refractive index⁵⁴ and that this may explain the fact that eye does not become myopic. However, the results reported are not consistent with what is known about the protein/total water proportions in the human lens and their constancy with age.²⁹ It is most likely that Moffat et al.⁵⁴ were measuring the proportions of free water, which increases in the human lens with age.⁵⁵

An overall decrease in the nuclear refractive index, would be physiologically unsound, as it would require an imbibing of water or a leakage of protein from the center of the lens. Given that the nuclear refractive index in the human lens is not a part of the gradient index but is constant, an overall decrease in magnitude would have a negligible effect on refractive power of the eye.

The constancy of the nuclear refractive index and the index gradient in the cortex require a higher-order polynomial to fit the index distribution in the human lens. A fourth-order polynomial was found to give a highly significant fit to the refractive index gradient in the equatorial plane of the human lens (Fig. 8).²⁹

This form for the radial refractive index gradient (i.e., in the equatorial plane) has been found in studies that used ray tracing (for 633 nm)⁴⁵ and for measurements taken with a fiber optic sensor (for 670 nm).⁴⁶ The two techniques have also shown results with similar magnitudes for the nuclear refractive index.^{45,46}

The refractive index along the optic axis similarly has a gradient in the cortical region with a constant index nucleus.^{46,47} However, the assumptions of concentric, isoindicial contours that follow the surface shape of the lens, may not apply to younger lenses (from the third decade of life).⁴⁶ As with the bovine lens, and consistent with the asymmetry of shape in the sagittal plane, there is a slight difference in the refractive index gradients in the anterior and posterior sections of the lens and this is also affected by lens age.²⁹

The difficulty in ascribing a single formulaic description to the refractive index gradient in the human lens is not only complicated by changes to the cortical gradient with age, and by differences in the axial and radial gradients. There is also the confounding factor of accommodation. The human lens, in individuals under the age of around 50, is capable of altering its shape to meet the focusing requirements of the eye. This capacity is greater in younger eyes and it is not certain, what accommodative state a young lens is in when it is removed from the eye. It may not be unreasonable to expect that when the lens is released from the tension applied by the ciliary muscle, in situ, it takes on its most accommodated, that is, most curved shape. This, however, assumes that the lens is an elastic body, with the fully accommodated state: the one to which it returns when applied stretching forces mediated by the ciliary muscle are removed. This assumption has and continues to court controversy and so it remains uncertain what accommodative state the lens is in when removed from the eye. Moreover, given that the capacity for accommodation diminishes with age, the accommodative state of the in vitro lens will vary depending on the age of the lens. How the refractive index gradient alters with accommodation is not known, although some studies have shown that the accommodative process is caused by a change

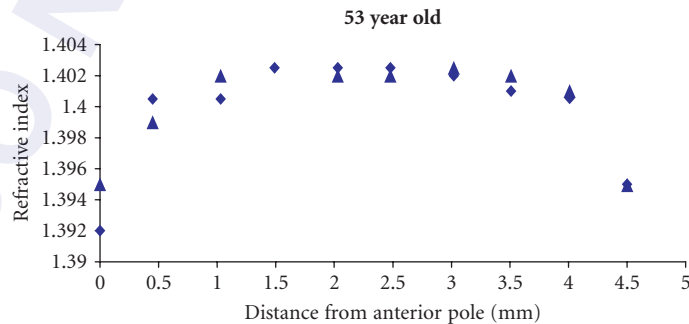


FIGURE 8 Variations in refractive index (for 670 nm) along the optic axis of a 53-year-old human lens plotted against the normalized distance from the anterior pole. (From Ref. 29.)

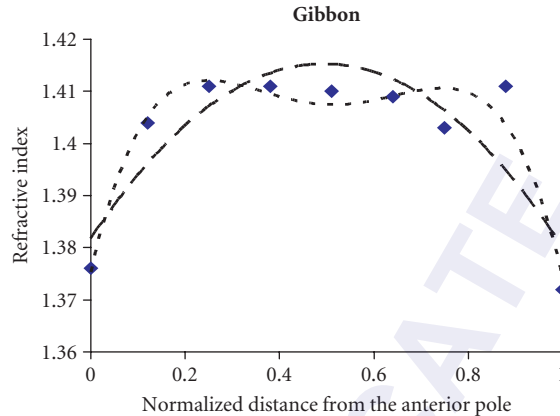


FIGURE 9 Variations in refractive index (for 670 nm) along the optic axis of the gibbon lens plotted against the normalized distance from the anterior pole. (From Ref. 29.)

in the length of the nucleus^{51,56–58} with no change in the thickness of the cortex. This would suggest that the cortical refractive index gradient remains unchanged, the only alteration to the profile being an elongation to the constant index nuclear section.

There is very little information about the refractive index gradients in the lenses of lower-order primates. Results from a single gibbon lens show the same type of profile as for the human lens with a similar index magnitude (Fig. 9).²⁹

It is not certain whether the same confounding factor of accommodation applies to the optics and refractive index gradient in the gibbon lens. However, studies on the rhesus monkey have shown that the species does accommodate.⁵⁹ As the gibbon belongs to a primate family that is more closely related to humans and its visual needs, from close range feeding to swinging from trees suggests that this species is likely to have an accommodating lens.

From the refractive index gradients, the primate lenses are the only ones that do not have a parabolic gradient but the refractive index gradient flattens to a region of constant refractive index in the nucleus. This constant, maximum value is lower than for other species thus far studied rendering these lenses softer and more pliable. This is consistent with the need for accommodation but also means that the refractive index gradient in a primate lens contributes less to the refractive power than the parabolic gradients, with higher-index magnitudes, found in other species.

19.17 FUNCTIONAL CONSIDERATIONS

Gradient index media are found in nature and the concept is exploited in the manufacture of lenses and optical devices for a wide range of applications. The fundamental difference between the biological and the manufactured lenses is that the former is not static in its optical qualities. Within any given species, the refractive index gradient will change with age because the lens continues to grow throughout life and it may change its shape in those species that possess accommodating lenses. Moreover, when considering the eye lens and its optical properties, it should not be considered in isolation. The eye lens functions as part of the optical system of the eye. The rays incident on its anterior surface will have been refracted by the cornea and traversed the anterior chamber. The shapes of the cornea and anterior chamber depths vary from species to species and individual variations within a species can be found. The humours that surround the lens in the eye will also impact on the contribution to refraction by the refractive index gradient. The more

closely the edge of the gradient is matched to the surrounding media, the more contribution is made by the refractive index gradient. No biological lens has yet been shown to have an edge index value matched to the refractive index of the aqueous and vitreous (around 1.336) and, with the techniques of ray tracing and optic fiber sensing, this is the most difficult part of the refractive index gradient to measure.

19.18 SUMMARY

The major types of index gradient found in lenses, optical devices, and created by nature have been described. These can be broadly categorized into gradients that vary parallel to the path of light, perpendicular to the path of light, or with radial distance from the center of the structure (found in spherical lenses). The eye lens is an excellent example of a GRIN structure with an index that varies both along and perpendicular to the optic axis. The refractive index distribution of the eye lens is related to its biological organization and to the properties of the major structural entities, the crystallin proteins. The concentration of proteins affects the magnitude of refractive index and the distributions of the different classes influence the shape of the index gradient. In spite of the complex structure of proteins and the changing interactions between protein subclasses and between proteins and water, it has been possible to measure the GRIN structure of the eye lens and in most species examined this appears to have a similar form. Transposing this form to an intraocular implant that cannot only mimic the GRIN structure but also provide some accommodative capacity is the subject of further research.

19.19 REFERENCES

1. J. C. Maxwell, *Cambridge Dublin Math Journal* **8**:188 (1854).
2. E. W. Marchand, *Gradient Index Optics*, Academic Press, New York, 1978.
3. A. Frank, F. Leyvraz, and K. B. Wolf, "Hidden Symmetry and Potential Group on the Maxwell Fish-eye," *J. Math. Phys.* **31**:2757–2768 (1990).
4. M. Born and E. Wolf, *Principles of Optics*, Pergamon, New York, 1959.
5. R. K. Luneberg, *Mathematical Theory of Optics*, Brown University, Providence, Rhode Island, 1944, pp. 189–213.
6. R. W. Wood, *Physical Optics*, Macmillan, New York, 1905, p. 72.
7. J. C. Palais, *Optic Fiber Communications*, Prentice Hall, Englewood Cliffs, NJ, 1988, p. 95.
8. P. H. Berning, "Use of Equivalent Films in the Design of Infrared Multilayer Antireflection Coatings," *J. Opt. Soc. Am.* **52**:431–436 (1962).
9. W. H. Southwell, "Gradient-index Antireflection Coatings," *Opt. Lett.* **8**:584 (1983).
10. S. Barbero, "Refractive Power of a Multilayer Rotationally Symmetric Model of the Human Cornea and Tear Film," *J. Opt. Soc. Am. A* **23**:1578–1585 (2006).
11. D. A. Palmer and J. Sivak, "Crystalline Lens Dispersion," *J. Opt. Soc. Am.* **71**:780–782 (1981).
12. J. G. Sivak and T. Mandelman, "Chromatic Dispersion of the Ocular Media," *Vis. Res.* **22**:997–1003 (1982).
13. R. D. Fernald and S. E. Wright, "Maintenance of Optical Quality during Crystalline Lens Growth," *Nature* **301**:618–620 (1983).
14. D. Axelrod, D. Lerner, and P. J. Sands, "Refractive Index within the Lens of a Goldfish Eye Determined from the Paths of Thin Laser Beams," *Vis. Res.* **28**:57–65 (1988).
15. W. S. Jagger, "The Optics of the Spherical Fish Lens," *Vis. Res.* **32**:1271–1284 (1992).
16. R. H. H. Kroger, M. C. W. Campbell, R. Munger, and R. D. Fernald, "Refractive Index Distribution and Spherical Aberration in the Crystalline Lens of the African Cichlid Fish *Haplochromis Burtoni*," *Vis. Res.* **34**:1815–1822 (1994).
17. B. K. Pierscionek and R. C. Augusteyn, "Refractive Index and Protein Distributions in the Blue Eye Trevally Lens," *J. Am. Optom. Assoc.* **66**:739–743 (1995).

18. W. S. Jagger and P. J. Sands, "A Wide-angle Gradient Index Optical Model of the Crystalline Lens and Eye of the Rainbow Trout," *Vis. Res.* **36**:2623–2639 (1996).
19. R. H. H. Kroger, M. C. W. Campbell, and R. D. Fernald, "The Development of the Crystalline Lens is Sensitive to Visual Input in the African Cichlid Fish, *Haplochromis Burtoni*," *Vis. Res.* **41**:549–559 (2001).
20. L. F. Garner, G. Smith, S. Yao, and R. C. Augusteyn, "Gradient Refractive Index of the Crystalline Lens of the Black Oreo Dory (*Alloctytus Niger*): Comparison of Magnetic Resonance Imaging (MRI) and Laser Ray-trace Methods," *Vis. Res.* **41**:973–979 (2001).
21. B. K. Pierscionek, G. Smith, and R. C. Augusteyn, "The Refractive Increments of Bovine α , β , and γ -Crystallins," *Vis. Res.* **27**:1539–1541 (1987).
22. C. Slingsby, B. Norledge, A. Simpson, O. A. Bateman, G. Wright, H. P. C. Driessen, P. F. Lindley, D. S. Moss, and B. Bax, "X-ray Diffraction and Structure of Crystallins," *Prog. Retin Eye Res.* **16**:3–29 (1996).
23. P. E. Malkki and R. H. H. Kroger, "Visualization of Chromatic Correction of Fish Lenses by Multiple Focal Lengths," *J. Opt. A; Pure Appl. Opt.* **7**:691–700 (2005).
24. W. S. Jagger and P. J. Sands, "A Wide-angle Gradient Index Optical Model of the Crystalline Lens and Eye of the Octopus," *Vis. Res.* **39**:2841–2852 (1999).
25. M. C. W. Campbell and A. Hughes, "An Analytic Gradient Index Schematic Lens and Eye for the Rat that Predicts Aberrations for Finite Pupil," *Vis. Res.* **21**:1129–1148 (1981).
26. M. C. W. Campbell, "Measurement of Refractive Index in an Intact Crystalline Lens," *Vis. Res.* **24**:409–415 (1984).
27. P. L. Chu, "Non Destructive Measurements of Index Profile of an Optical-fiber Perform," *Electronic Letts.* **13**:736–738 (1977).
28. R. Barer and S. Joseph, "Refractometry of Living Cells. Part 1. Basic Principles," *Quart. J. Microscop. Sci.* **95**:399–423 (1954).
29. B. K. Pierscionek, "Species Variations in the Refractive Index of the Eye Lens and Patterns of Change with Ageing," in O. R. Ioseliani, (ed.), *Focus on Eye Research*, Nova Science Publishers Inc., New York, 2005, pp. 91–116.
30. Q. L. Huang, P. Russell, S. H. Stone, and J. S. Zigler Jr., "Zeta-crystallin a Novel Lens Protein from the Guinea Pig," *Curr. Eye Res.* **6**:725–732 (1987).
31. T. Borras and A. Rodokanaki, "Primary Structure of Zeta-crystallin Protein from Guinea Pig. Its Similarity to the Enzyme Alcohol Dehydrogenase," *Lens Eye Toxic. Res.* **6**:795–805 (1989).
32. P. V. Rao, J. Horwitz, and J. S. Zigler Jr., "Chaperone-like Activity of α -Crystallin. The Effect of NADPH on its Interaction with Zeta-crystallin," *J. Biol. Chemistry* **269**:13266–13272 (1994).
33. S. Nakao, S. Fujimoto, R. Nagata, and K. Iwata, "Model of Refractive-Index Distribution in the Rabbit Crystalline Lens," *Opt. Soc. Am.* **58**:1125–1130 (1968).
34. W. S. Jagger, "The Refractive Structure and Optical Properties of the Isolated Crystalline Lens of the Cat," *Vis. Res.* **30**:723–738 (1990).
35. L. Mattheissen, "Untersuchungen uber den Aplanatismus und die Periscopie der Kristalllinse in den Augen der Fische," *Pfluegers Archiv.* **21**:287–307 (1880).
36. B. K. Pierscionek, "Growth and Ageing Effects on the Refractive Index in the Equatorial Plane of the Bovine Lens," *Vis. Res.* **29**:1759–1766 (1989).
37. B. K. Pierscionek, D. Y. C. Chan, J. P. Ennis, G. Smith, and R. C. Augusteyn, "A Non-destructive Method of Constructing Three-dimensional Gradient Index Models for Crystalline Lenses: I. Theory and Experiment," *Am. J. Optom. Physiol. Opt.* **65**:481–491 (1988).
38. B. K. Pierscionek, "The Refractive Index along the Optic Axis of the Bovine Lens," *Eye* **9**:776–782 (1995).
39. B. K. Pierscionek, A. Belaidi, and H. H. Bruun, "Optical Development in the Foetal Bovine Lens," *Exp. Eye Res.* **77**:639–641 (2003).
40. O. Hockwin, "Age Changes of Lens Metabolism," *Altern Entwickl* **1**:95–129 (1971).
41. B. K. Pierscionek and R. C. Augusteyn, "Growth Related Changes to Functional Parameters in the Bovine Lens," *Biochim. Biophys. Acta* **1116**:283–290 (1992).
42. B. K. Pierscionek, A. Belaidi, and H. H. Bruun, "Refractive Index Gradient in the Porcine Lens for 532 and 633 nm Light," *Eye* **19**:375–381 (2005).
43. C. E. Jones and J. M. Pope, "Measuring Optical Properties of an Eye Lens Using Magnetic Resonance Imaging," *Magn. Reson. Imaging* **22**:211–220 (2004).

44. C. L. de Korte, A. F. van der Steen, J. M. Thijssen, J. J. Duindam, C. Otto, and G. J. Puppels, "Relation Between Local Acoustic Parameters and Protein Distribution in Human and Porcine Eye Lenses," *Exp. Eye Res.* **59**:617–627 (1994).
45. B. K. Pierscionek and D. Y. C. Chan, "The Refractive Index Gradient of the Human Lens," *Optom. Vis. Sci.* **66**:822–829 (1989).
46. B. K. Pierscionek, "Refractive Index Contours in the Human Lens," *Exp. Eye Res.* **64**:887–893 (1997).
47. B. K. Pierscionek, "Variations in Refractive Index and Absorbance of 670 nm Light with Age and Cataract Formation in Human Lenses," *Exp. Eye Res.* **60**:407–414 (1995).
48. P. P. Fagerholm, B. T. Philipson, and B. Lindstrom, "Normal Human Lens, the Distribution of Protein," *Exp. Eye Res.* **33**:615–620 (1981).
49. A. Huizinga, A. C. C. Bot, F. F. M. de Mul, G. F. J. M. Vrensen, and J. Greve, "Local Variation in Absolute Water Content of Human and Rabbit Eye Lenses Measured by Raman Microspectroscopy," *Exp. Eye Res.* **48**:487–496 (1989).
50. I. Siebinga, G. F. J. M. Vrensen, F. F. M. de Mul, and J. Greve, "Age-related Changes in Local Water and Protein Content of Human Eye Lenses Measured by Raman Microspectroscopy," *Exp. Eye Res.* **53**:233–239 (1991).
51. N. Brown, "The Change in Shape and Internal Form of the Lens of the Eye on Accommodation," *Exp. Eye Res.* **15**:441–459 (1973).
52. N. Brown, "The Change in Lens Curvature with Age," *Exp. Eye Res.* **19**:175–183 (1974).
53. B. K. Pierscionek, "Presbyopia and the Effect of Refractive Index," *Clin. Exp. Optom.* **73**:26–36 (1990).
54. B. A. Moffat, D. A. Atchison, and J. M. Pope, "Age-related Changes in Refractive Index Distribution and Power of the Human Lens as Measured by Magnetic Resonance Micro-imaging in vitro," *Vis. Res.* **42**:1683–1693 (2002).
55. D. Lahm, L. K. Lee, and F. A. Bettelheim, "Age Dependence of Freezable and Nonfreezable Water Content of Normal Human Lenses," *Invest. Ophthalmol. Vis. Sci.* **26**:1162–1165 (1985).
56. B. Patnaik, "A Photographic Study of Accommodative Mechanisms: Changes in Lens Nucleus during Accommodation," *Invest. Ophthalmol.* **6**:601–611 (1963).
57. J. F. Koretz, C. A. Cook, and P. L. Kaufman, "Accommodation and Presbyopia in the Human Eye. Changes in the Anterior Segment and Crystalline Lens with Focus," *Invest. Ophthalmol. Vis. Sci.* **38**:569–578 (1997).
58. M. Dubbelman, G. L. van der Heijde, H. A. Weeber, and G. F. Vrensen, "Changes in the Internal Structure of the Human Crystalline Lens with Age and Accommodation," *Vis. Res.* **43**:2363–2375 (2003).
59. J. F. Koretz, A. M. Bertasso, M. W. Neider, B. True-Gabelt, and P. L. Kaufman, "Slit-lamp Studies of the Rhesus Monkey Eye. II. Changes in Crystalline Lens Shape, Thickness and Position during Accommodation and Aging," *Exp. Eye Res.* **45**:317–326 (1987).

This page intentionally left blank.

DO NOT DUPLICATE

Edward S. Bennett

*College of Optometry
University of Missouri
St. Louis, Missouri*

20.1 GLOSSARY

Accommodative demand. This represents the difference between distance and near refractive corrections.

Aspheric lens design. An aspherical surface, typically elliptical, often used on the back surface to align with an aspheric anterior cornea. In some presbyopic designs, either a back surface with a hyperbolic or high rate of flattening is used to provide plus power or an aspheric front surface is used for the same purpose.

Axial ametropia. Ametropia resulting from the axial elongation of the globe.

Axial edge lift. Pertains to the distance between the extension of the tangent of the base curve radius and the absolute edge; it is measured parallel to the lens optical axis or the vertical distance from the lens edge to an extension of the base curve radius.

Base curve radius (BCR). The radius of curvature of the central posterior optical section in millimeters.

Back vertex power (BVP). Commonly used to express the power of a contact lens, the power is determined from a fixed position with the concave surface of the lens against the lensometer lens stop.

Bitoric design. Gas permeable (GP) lens designs incorporating two base curve radii and multiple powers for correction of high corneal astigmatism, typically 2.50 D or greater.

Contact lens power. The dioptric power of the contact lens, usually expressed as back vertex power.

Edge clearance. The actual distance from the lens edge to the cornea. This value is less than the calculated edge lift due to the asphericity of the anterior cornea.

Effective power. The power of the contact lens at the corneal plane.

Fluoro-silicone/acrylate (F-S/A). The current generation of rigid gas-permeable contact lenses consisting primarily of fluorine, silicone, and methyl methacrylate.

Front vertex power (FVP). The power is determined from a fixed position with the convex surface of the lens against the lensometer lens stop.

High-order aberrations. Aberrations induced by the optics of the eye, notably spherical aberration and coma.

“K.” This refers to the flatter keratometry reading. GP lenses are generally selected with a base curve radius slightly steeper or flatter than “K.”

Lacrimal lens. The tear lens between a GP contact lens and the anterior corneal surface. It is a plus power if the lens exhibits apical clearance; it is minus if the lens exhibits apical bearing.

Multifocal contact lenses. Contact lenses that provide more than one correcting power within the lens; these are often used for the correction of presbyopia.

Optical zone diameter (OZD). This represents the linear equivalent of the base curve radius and provides the visual optics for both GP and soft lenses.

Overall diameter. This represents the linear, chord, edge-to-edge measurement of the lens in millimeters.

Oxygen permeability (Dk). The current term to describe the potential of a contact lens material to transmit oxygen which considers both, solubility and diffusion.

Prism (in contact lenses). Prism is produced by varying the thickness from the superior to the inferior region of a contact lens while maintaining the same front and back surface curvatures.

Radial edge lift. Pertains to the distance between the extension of the tangent of the base curve radius and the absolute edge; it is measured from the lens edge perpendicular to an extension of the base curve radius.

Refractive ametropia. Ametropia resulting from abnormal refractive component(s) of the eye.

Relative spectacle magnification (RSM). This compares the corrected ametropic retinal image to that of a standard emmetropic schematic eye. Ametropia can be purely axial, resulting from the axial elongation of the globe, refractive, resulting from abnormal refractive component(s) of the eye or a combination of both factors.

Residual astigmatism. With spherical GP lenses this pertains to the total ocular astigmatism (defined as the refractive cylinder minus the corneal (often keratometric) cylinder. With spherical soft lenses, this pertains to the refractive cylinder.

Rigid GP contact lenses. Rigid lens materials, smaller in diameter with higher optical quality than soft lenses, which typically incorporate methyl methacrylate, silicone, and fluorine within the polymer matrix.

Silicone-hydrogel. The current generation of soft lens materials incorporating silicone within the matrix for higher oxygen permeability.

Soft toric lenses. Soft lens designs that correct for refractive astigmatism.

Spectacle magnification. This pertains to the ratio of retinal image size of the corrected ametropic eye to the retinal image size of the same eye uncorrected.

20.2 INTRODUCTION

The optics pertaining to contact lenses have several similarities to those of spectacle lenses. However, the primary differences pertain to such factors as the lesser thickness of contact lenses, the fact they are in direct contact with the tear film—not 11 to 15 mm away from the cornea—the differences present between soft and rigid gas permeable (GP) lens materials, as well as differences in accommodation and convergence effects on the eye. It is evident that, as a result of their optical properties, contact lenses are unique among the forms of correction available to wearers as well as, in many cases, a modality that exceeds other visual correction options in both the quality of vision and the visual freedom present to the wearer. The goal of this chapter is to emphasize the basic contact lens optics commonly used in eye care practices.

20.3 CONTACT LENS MATERIAL, COMPOSITION, AND DESIGN PARAMETERS

Contact Lens Materials

Contact lens materials can be considered either rigid or soft. The first modern rigid lens was the corneal polymethylmethacrylate (PMMA) lens which was popular from the late 1940s to the early 1980s. This material was essentially limited to the PMMA plastic and was, therefore, not oxygen permeable. GP lenses incorporating silicone (for oxygen permeability) as well as stability agents and wetting agents (the latter to help offset the hydrophobic nature of silicone) was introduced in 1979. However, the introduction of fluoro-silicone/acrylate (F-S/A) lenses in the mid 1980s was a major breakthrough and almost all GP lenses in use today are F-S/A materials. The addition of fluorine both assisted in surface wettability via encouraging the tear film mucin to form an even layer on the lens surface but also increased oxygen permeability—through increased oxygen solubility in fluorine—while allowing the silicone content to be reduced. GP materials have many benefits including high quality of vision resulting from the excellent optical quality of these materials and their ability to mold the front surface of the cornea and automatically correct for astigmatism.

Soft (or hydrogel) contact lenses were introduced in 1971. Soft lenses differ from GP lenses as a result of their ability to bind substantial amounts of water.¹ A loose or less tight cross-linking structure allows these materials to range in water content from 25 to greater than 70 percent when equilibrated in normal (0.9 percent) saline.² The water content can be expressed by the following equation:³

$$\text{Water content} = \frac{\text{Weight of lens water}}{\text{Total weight of hydrated lens}} \times 100$$

Soft lens materials are quite porous relative to GP lenses resulting in a reduction in optical quality which can impact vision. This pore size has been estimated at about 8 Å for low water content materials and 20 to 30 Å for high water content lenses.⁴ However, as a result of their larger overall diameter resulting in less movement with the blink, soft lenses have the advantage of being more initially comfortable and easier to adapt to for lens wearers. As soft lenses allow very little oxygen via tear exchange with the blink and the oxygen transmission is achieved primarily through diffusion through the lens, they have not represented a highly successful long-term option for extended wear. However, recently silicone hydrogel lenses have been introduced with oxygen permeability values as much as 5 to 10 times greater than hydrogel lenses and have increased from 0 percent of new fits in the United States in 2002 to over 35 percent in 2006.⁵ Silicone hydrogel lens materials combine the high oxygen permeability of silicone with other material characteristics of conventional hydrogel materials that facilitate fluid transport and lens movement.⁶ The surfaces of these lenses are treated to make them more hydrophilic and, therefore, more biocompatible. These materials are FDA approved for up to 30 days of continuous lens wear.

Curvature, Diameter, and Thickness

To understand the optics associated with contact lenses, it is important to have a good working knowledge of contact lens design. A brief description of each parameter is discussed below and summarized in Table 1.

Base Curve Radius (BCR) The base curve radius of a contact lens is the radius of curvature of the central posterior optical section in millimeters⁷ (Fig. 1). This has also been termed the back central optical radius (BCOR). For rigid GP lens materials the base curve radius is typically selected to be approximately the same curvature as the anterior cornea such that the lens will “align” with the anterior corneal surface. With soft or hydrogel lenses, the base curve radius is typically about 1.0 mm flatter than the anterior corneal curvature as they are much larger in diameter and simply drape over

TABLE 1 Contact Lens Design Parameters

Base Curve Radius	The radius of curvature of the central posterior optical section in millimeters.
Secondary Curve Radius/Width	In a tricurve or tetracurve GP lens design, the secondary curve radius is the curve adjacent to the central base curve radius and is typically flatter than the base curve radius. The width of this curve (often 0.2 to 0.5 mm) represents the linear size of this curve.
Intermediate Curve Radius/Width	In a tetracurve GP lens design, the intermediate curve radius is the curve adjacent to the secondary curve radius and is flatter than the secondary curve radius. The width of this curve (often 0.2 to 0.5 mm) represents the linear size of this curve.
Peripheral Curve Radius/Width	In a tricurve or tetracurve GP lens design, the peripheral curve radius is the outermost curve radius. It is the flattest curve and its parameter is most responsible for the resultant edge clearance (distance from edge to cornea). The width of this curve (often 0.2 to 0.5 mm) represents the linear size of this curve.
Transition Zone Blend	As the junction between two peripheral curves can be rather sharp, laboratories typically blend or smooth these junctions.
Overall Diameter	This represents the linear, chord, edge-to-edge measurement of the lens in millimeters.
Optical Zone Diameter	This represents the linear equivalent of the base curve radius and provides the visual optics for both GP and soft lenses.
Center Thickness	This is the thickness as measured at the exact center of the lens.
Edge Thickness	This is the thickness as measured at the edge of the lens. This value should be similar to the center thickness. As high minus lenses have thick edges, a plus lenticular is often manufactured to thin the anterior edge; likewise for low minus and plus power lenses, a minus lenticular is generated to increase the edge thickness.
Power	Lens power is the dioptric power of the contact lens, usually expressed as back vertex power.

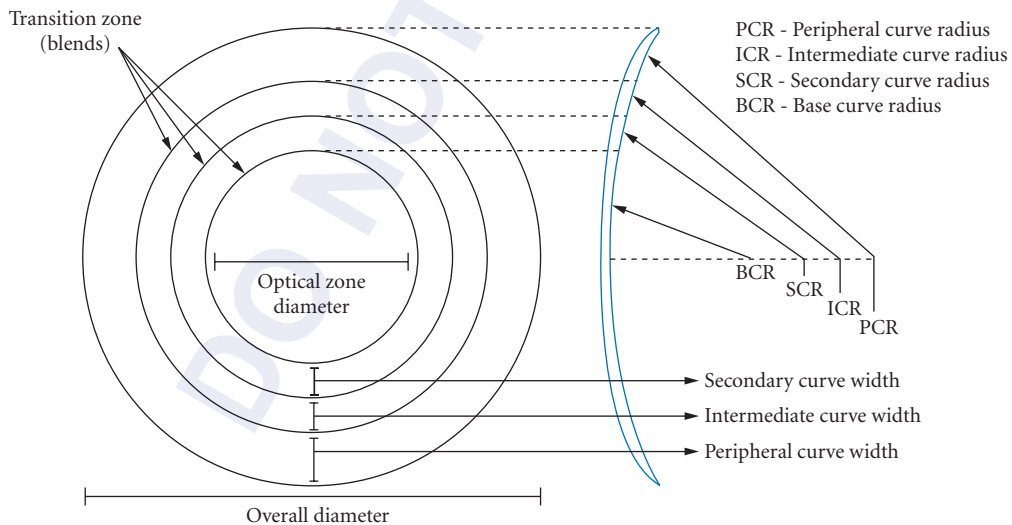


FIGURE 1 Cross-section of a tetracurve rigid GP contact lens showing overall diameter, optical zone diameter, secondary curve, intermediate curve, and peripheral curve radii and widths and transition zone blends. (From Ref. 30.)

the cornea and align with the adjacent sclera.⁸ For GP lenses base curve radius is often specified in millimeters; however, it has also been denoted by its dioptric equivalent. The conversion formula is

$$\text{BCR (in diopters)} = \frac{n' - n}{r}$$

where n' = index surrounding the lens (although the tear layer is 1.336, the keratometer is calibrated for 1.3375 and this value is used)
 n = refractive index of air or 1.0
 r = base curve radius in meters

For example, if the base curve radius was verified to be 7.50 mm with the radiuscope, the dioptric equivalent would be $1.3375 - 1.0/0.0075$ or 45 D. BCR values in diopters and their converted values in millimeters are provided in Table 2.

Overall Diameter (OAD)/Optical Zone Diameter (OZD) The linear, chord, edge-to-edge measurement of the lens in millimeters is the overall diameter. GP lenses typically have OAD values between 9.0 to 10.0 mm. Soft or hydrogel lenses typically have OAD values of approximately 14.0 mm. The optical zone diameter is the linear equivalent of the base curve radius and is the central section of the lens which provides usable optics for vision; the lens periphery has a primary function of optimizing tear exchange and the fitting relationship. By convention, the lens OZD for GP lenses is the overall diameter minus the posterior peripheral curve widths (often the OAD—approximately 1.0 to 1.5 mm). The OZD for soft lenses is the overall lens diameter minus twice the anterior peripheral curve width in millimeters.

Posterior Peripheral Curve Systems GP lenses have a standard curvature (i.e., the base curve radius) with typically two-to-three peripheral curve radii. These peripheral curve radii are progressively

TABLE 2 Keratometer Diopter Conversion to Millimeters

BCR (D)	BCR (mm)	BCR (D)	BCR (mm)
40.00	8.44	46.25	7.30
40.25	8.39	46.50	7.26
40.50	8.33	46.75	7.22
40.75	8.28	47.00	7.18
41.00	8.23	47.25	7.14
41.25	8.18	47.50	7.11
41.50	8.13	47.75	7.07
41.75	8.08	48.00	7.03
42.00	8.04	48.25	7.00
42.25	7.99	48.50	6.96
42.50	7.94	48.75	6.92
42.75	7.89	49.00	6.89
43.00	7.85	49.25	6.85
43.25	7.80	49.50	6.82
43.50	7.76	49.75	6.78
43.75	7.72	50.00	6.75
44.00	7.67	50.25	6.72
44.25	7.63	50.50	6.68
44.50	7.58	50.75	6.65
44.75	7.54	51.00	6.62
45.00	7.50	51.25	6.59
45.25	7.46	51.50	6.55
45.50	7.42	51.75	6.52
45.75	7.38	52.00	6.49
46.00	7.34		

flatter from the base curve to the outermost curve. As the cornea is aspheric and flattens from center to periphery, a moncurve GP lens would likely be too steep in the periphery and not allow tear exchange—important for oxygen exchange and debris removal—to occur. Therefore, if the base curve radius was 7.80 mm, one common philosophy recommends a secondary curve radius (SCR) equal to 1 mm flatter than the base curve radius, and a peripheral curve radius 2 mm flatter than the SCR.⁹ Both of these curves would be 0.3 mm wide. This is termed a tricurve design. A tetracurve design would include a secondary curve, an intermediate curve, and a peripheral curve as diagrammed in Fig. 1. Some designs (to be discussed) are manufactured such that they have a gradual increase in flattening in the periphery, termed a peripheral aspheric design. Soft lens posterior curvature is either moncurve, bicurve (typically one peripheral curve on average about 0.5 mm wide and as flat as 12 mm in curvature), or aspheric. The nonmonocurve designs are intended to align the lens periphery with the peripheral limbal or scleral topography.⁷

Center Thickness Contact lenses are relatively thin lens designs (relative to spectacles) with lenses as thin as 0.035 mm for a high minus power ultrathin soft lens to as high as 0.5 to 0.6 mm for high plus power soft and GP lenses, especially if prism ballast is incorporated into the lens to stabilize it on the eye. Minus power GP lenses tend to vary from 0.09 mm (ultrathin) to 0.19 mm for low minus standard thickness designs. Low plus power lenses typically have center thicknesses in the 0.20 to 0.25 mm range and high plus power lenses have center thicknesses in the 0.26 to 0.60 mm range. Center thickness impacts both lens mass and oxygen transmission. Plus power lenses have greater mass and an anterior center of gravity making centration on the eye more challenging. In addition, oxygen transmission (i.e., Dk/t where “Dk” is the oxygen permeability of the material and “t” is the thickness—either t_{avg} for average thickness or t_c for center thickness) reduces as the center thickness increases. However, a minus power GP lens that is too thin will bend or flex on the eye and will negate one of the benefits of this modality, the ability to mold the anterior corneal surface into a sphere and correct astigmatism.

Edge Thickness Edge thickness varies for GP and soft lenses but values of 0.06 to 0.12 mm are not uncommon. With most GP designs, the apex is tapered and rounded toward the posterior surface. As high minus power (typically ≥ 5 D) lenses tend to have thick edges and both low minus power (typically ≤ 1.50 D) and all plus power lenses have thin edges, special manufacturing procedures are utilized to maintain a uniform edge thickness. With GP lenses this is typically a plus lenticular (a manufacturing process that thins the edge) for high minus power lenses and a minus lenticular (a manufacturing process that results in a thicker edge) for low minus and plus power lenses. This uniform edge thickness is important for allowing the upper lid to lift up the lens with the blink while not being so thick that the lid would push the lens inferiorly.

20.4 CONTACT LENS POWER

Contact lens power typically differs from spectacle lens power. For example, a spectacle lens power of $-2.00 -1.00 \times 180$ would indicate that the horizontal meridian of the spectacles has a power of -2.00 D and the vertical meridian has a power of -3.00 D. If this patient has no residual astigmatism (i.e., the anterior corneal astigmatism is also equal to -1.00×180), the back surface of a spherical power GP lens will mold the front surface of the cornea into a sphere; therefore, the astigmatism is corrected. A spherical soft lens, however, simply drapes over the eye and does not correct for refractive astigmatism and, in theory, the -1.00×180 correction in the spectacles would not be corrected by such a soft lens and a special design, a soft toric lens (to be discussed) would be indicated.

The Contact Lens as a Thick Lens

Contact lenses are manufactured at a much lesser center thickness (typically 0.08 to 0.5 mm) than spectacle lenses. Nevertheless, they still need to be considered as thick lenses when determining the overall power. Therefore, instead of simply adding the refractive power of the front surface to the

back surface power (i.e., $F_T = F_1 + F_2$), the curvatures of both surfaces, center thickness and index of refraction of the lens material need to be considered. The thick lens formula used for spectacle lenses could be considered.¹⁰ It is as follows:

$$F_T = F_1 + F_2 - (t/n)F_1F_2$$

where F_T = equivalent or true refractive power in diopters

$F_1 = (n' - n)/r_1$ = refractive power of the anterior lens surface

$F_2 = (n - n')/r_2$ = refractive power of the posterior lens surface

t = center thickness of the lens in meters

n' = refractive index of the lens material

n = refractive index of the medium surrounding the lens

r_1 = radius of curvature of posterior surface, in meters

r_2 = radius of curvature of anterior surface, in meters

However, as the positions of the principal planes vary with the design of the lens and its two surface powers, it is impractical to use the equivalent, or true, power of a contact lens as the absolute definition of refractive power. It is, instead, important to use a fixed position to measure refractive power.

Back Vertex Power and Front Vertex Power

A fixed position from which refractive power can be determined is via the use of either back vertex power (BVP) or front vertex power (FVP). The equations of these powers are as follows:

$$\text{BVP} = F_1 + \left[\frac{F_2}{1 - (t/n')F_2} \right]$$

$$\text{FVP} = \frac{F_1}{1 - (t/n')F_1} + F_2$$

The measurement of BVP and FVP is very straightforward on a lensometer. BVP is obtained with the back surface of the lens (concave surface) against the lens stop (see Fig. 2). Front surface



FIGURE 2 Determination of back surface power via the use of a lensometer.

TABLE 3 Differences between Back Vertex Power (BVP) and Front Vertex Power (FVP) in Air for GP Lenses ($N = 1.46$)

BVP (D)	FVP (D)	Center Thickness (mm)	Power Difference (D)
-15.00	-14.86	0.12	-0.14
-12.00	-11.88	0.12	-0.12
-9.00	-8.91	0.12	-0.09
-6.00	-5.94	0.12	-0.06
-3.00	-2.96	0.14	-0.04
Plano	Plano	0.18	0
+3.00	+2.95	0.21	+0.05
+6.00	+5.87	0.24	+0.13
+9.00	+8.75	0.30	+0.25
+12.00	+11.59	0.36	+0.41
+15.00	+14.41	0.41	+0.59
+20.00	+19.01	0.50	+0.99

power is obtained with the front surface of the lens (i.e., convex surface) against the lens stop. The focal lengths (in meters) are also easy to determine as they are simply the reciprocals of the vertex powers in diopters. For example, a BVP of -10 D would have a focal length of $1/10$ or 0.1 m (or 10 cm).

The difference between BVP and FVP is often negligible, especially in all minus and low plus power lenses. This is the result of the smaller center thicknesses that are present in these lens powers. However, with moderate-to-high plus power lenses—notably aphakic powers—this difference becomes significant (see Table 3).

As it can be shown, BVP is always greater than FVP. This has been a source of confusion in the past when it was not uncommon to specify (and/or receive from the fabricating laboratory) aphakic lenses via FVP. However, if the practitioner ordered the lenses based on BVP, the amount of resulting error could be considerable.¹¹ Fortunately, the contact lens industry has—with rare exception—utilized back surface power for all contact lenses in recent years.

Effective Power

A primary difference between spectacle power and contact lens power—notably in high ametropia—pertains to the effective power at the corneal plane. A spectacle lens is often positioned 10 to 15 mm in front of the cornea. This distance from the spectacle lens to the corneal apex is termed “vertex distance.” A contact lens, however, by virtue of being positioned adjacent to the cornea has a vertex distance of zero.

As a rule, as a plus power lens is brought closer to the cornea, its effective power decreases, such that its refractive power must be increased in order to maintain a constant amount of power relative to the eye. Likewise, when a minus power lens is brought closer to the eye, its effective power increases, such that its refractive power must be decreased in order to maintain a constant amount of power relative to the eye.¹⁰ When correcting an individual with a myopic refractive error, a contact lens will require a longer focal length, therefore resulting in less minus power than a spectacle lens. Likewise, when correcting an individual with a hyperopic refractive error, a contact lens will require a smaller focal length, therefore resulting in a higher plus refractive power than with a spectacle lens. Effective power differences between spectacle and contact lens corrections become clinically significant at about ± 4 D (see Table 4).

TABLE 4 Effective Spectacle Lens Power at the Corneal Plane (12 mm Vertex Distance)

Minus Lenses (in Diopters)			
Spectacle Lens Power	Effective Lens Power	Spectacle Lens Power	Effective Lens Power
-4.00	-3.75	+4.00	+4.25
-4.50	-4.25	+4.50	+4.75
-5.00	-4.75	+5.00	+5.25
-5.50	-5.25	+5.50	+6.00
-6.00	-5.50	+6.00	+6.50
-6.50	-6.00	+6.50	+7.00
-7.00	-6.50	+7.00	+7.75
-7.50	-7.00	+7.50	+8.25
-8.00	-7.25	+8.00	+8.75
-8.50	-7.75	+8.50	+9.50
-9.00	-8.25	+9.00	+10.00
-9.50	-8.50	+9.50	+10.75
-10.00	-9.00	+10.00	+11.25
-10.50	-9.50	+10.50	+12.00
-11.00	-9.75	+11.00	+12.75
-11.50	-10.00	+11.50	+13.25
-12.00	-10.50	+12.00	+14.00
-12.50	-10.75	+12.50	+14.75
-13.00	-11.25	+13.00	+15.50
-13.50	-11.50	+13.50	+16.00
-14.00	-12.00	+14.00	+16.75
-14.50	-12.50	+14.50	+17.50
-15.00	-12.75	+15.00	+18.25
-15.50	-13.00	+15.50	+19.00
-16.00	-13.50	+16.00	+19.75
-16.50	-13.75	+16.50	+20.50
-17.00	-14.00	+17.00	+21.25
-17.50	-14.50	+17.50	+22.25
-18.00	-14.75	+18.00	+23.00

There are two methods of determining the power at the corneal plane. The first is simply to use the following equation

$$F(c) = \frac{F(sp)}{1 - dF(sp)}$$

where $F(c)$ = power at the corneal plane
 $F(sp)$ = power at the spectacle plane
 d = vertex distance (in m)

Whereas, the vertex distance should be carefully measured in high ametropia patients, an average value of 12 mm is commonly used and will be used in the calculations presented in this chapter. If a patient has a refractive power of $-7.00 - 1.50 \times 180$, and the power in each meridian (i.e., -7.00 at 180 and -8.50 at 090) is "vertexed" to the corneal plane via this formula, the resulting powers would be -6.46 D and -7.71 D or $-6.46 - 1.25 \times 180$, which can be rounded off to $-6.50 - 1.25 \times 180$. If the individual was hyperopic and had a refractive error of $+7.00 - 1.50 \times 180$, and the power in each meridian (i.e., $+7.00$ at 180 and $+5.50$ D at 090) was vertexed to the corneal plane, the resulting powers would be $+7.64$ D and $+5.89$ D or $+7.64 - 1.75 \times 180$ or rounded off to $+7.75 - 1.75 \times 180$.

It can be observed that in compound myopic astigmatism, the spherical and cylinder power components decrease when referred to the corneal plane. In compound hyperopic astigmatism, the spherical and cylindrical power components increase or become greater when referred to the corneal plane.

The second method of effective power calculation would refer back to the diagrams in Figs. 3 and 4 via the use of focal lengths. If the individual's refractive error is $-7.00 - 150 \times 180$, the focal lengths of the two primary meridians are as follows:

$$\text{Horizontal meridian: } f_{180} = \frac{1}{F_{180}} = \frac{1}{7 \text{ D}} = 0.1429 \text{ m or } 142.9 \text{ mm (Fig. 3a)}$$

$$\text{Vertical meridian: } f_{090} = \frac{1}{F_{090}} = \frac{1}{8.50 \text{ D}} = 0.1176 \text{ m or } 117.6 \text{ mm (Fig. 3b)}$$

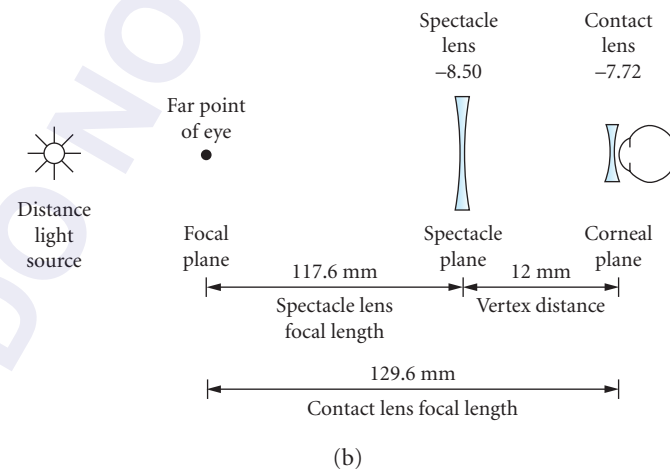
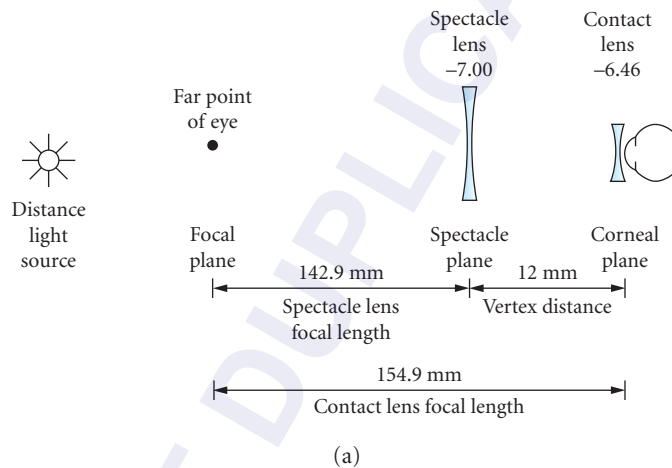


FIGURE 3 Determination of the effective power for the myopic patient in both the (a) horizontal meridian and (b) vertical meridian.

The focal lengths of the correction at the corneal plane will be 12 mm more; therefore, the final power per meridian is

$$\text{Horizontal meridian: } F_{180} = \frac{1}{0.1549 \text{ m}} = -6.46 \text{ D}$$

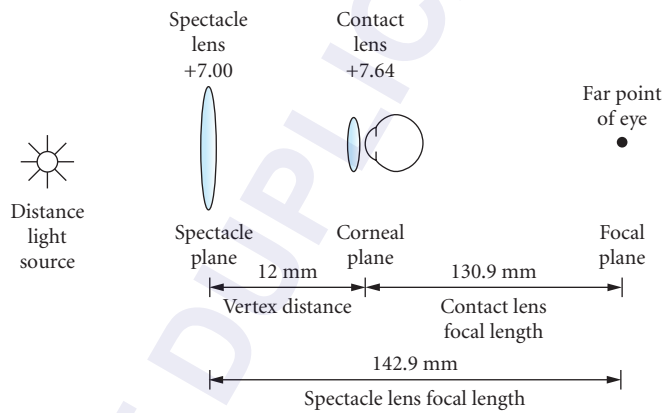
$$\text{Vertical meridian: } F_{090} = \frac{1}{0.1296 \text{ m}} = -7.72 \text{ D}$$

The corneal plane refraction is, therefore, $-6.46 - 1.26 \times 180$ or rounded off to $-6.50 - 1.25 \times 180$.

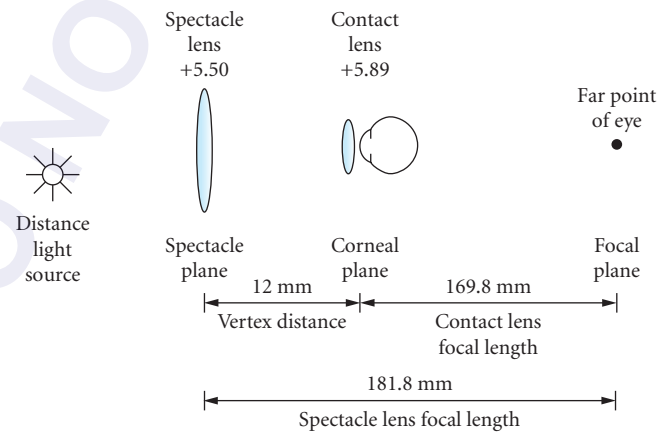
If the individual's refractive error is $+7.00 - 1.50 \times 180$, the focal lengths of the two primary meridians are as follows:

$$\text{Horizontal meridian: } f_{180} = \frac{1}{F_{180}} = \frac{1}{7 \text{ D}} = 0.1429 \text{ m or } 142.9 \text{ mm (Fig. 4a)}$$

$$\text{Vertical meridian: } f_{090} = \frac{1}{F_{090}} = \frac{1}{5.5 \text{ D}} = 0.1818 \text{ m or } 181.8 \text{ mm (Fig. 4b)}$$



(a)



(b)

FIGURE 4 Determination of the effective power for the hyperopic patient in both the (a) horizontal meridian and (b) vertical meridian.

The focal lengths of the correction at the corneal plane will be 12 mm less; therefore, the final power per meridian is

$$\text{Horizontal meridian: } F_{180} = \frac{1}{0.1309 \text{ m}} = +7.64 \text{ D}$$

$$\text{Vertical meridian: } F_{90} = \frac{1}{0.1698 \text{ m}} = +5.89 \text{ D}$$

The corneal plane refraction is, therefore, $+7.64 - 1.75 \times 180$ or rounded off to $+7.75 - 1.75 \times 180$.

Lacrimal Lens Considerations

Whereas effective power impacts both soft and GP lenses with patients who have a spherical refractive value equal to 4 D or greater, as soft lenses merely drape over the cornea, no compensation for the lacrimal or tear lens between the cornea and the contact lens is necessary. With GP lenses, however, a lacrimal or tear lens between cornea and contact lens often exists and needs to be factored into the final contact lens power determination. The baseline values to be considered when determining GP power pertain to the flatter keratometry reading and the spherical refractive value (at the corneal plane). The base curve radius selected for most GP wearers is very similar in curvature to the flatter keratometry reading; therefore, very little—if any—lacrimal lens is present. If fitted in alignment with the flatter keratometry reading, this is typically referred to as an “On *K*” lens-to-cornea fitting relationship. If the lens is fitted more curved or steeper in base curve radius than the baseline corneal curvature value it is termed “steeper than *K*.” If the lens is fitted less curved or flatter than the baseline corneal curvature value, it is termed “flatter than *K*.” The ideal GP lens-to-cornea fitting relationship would consist of an alignment or even tear film between lens and cornea as easily observed when fluorescein dye is applied to the tear film. Such a pattern is shown in Fig. 5. A base curve radius selected excessively steeper than “*K*” resulting in direct contact of the lens against the central cornea is shown in Fig. 6. This apical bearing relationship can ultimately result in central corneal distortion as well as excessive lens awareness due to the large amount of edge clearance present. A base curve radius selected excessively steeper than “*K*” can result in very little tear exchange peripherally, possibly resulting in edge sealoff and corneal edema. (See Fig. 7.)

In theory, as a result of the fact that the cornea is aspheric and gradually flattens in curvature from center to periphery, a slightly flatter than “*K*” base curve radius is often recommended. Of course the rate of flattening or corneal eccentricity can vary considerably and this will impact the fitting relationship and the base curve radius that will ultimately result in success. It is also important to remember

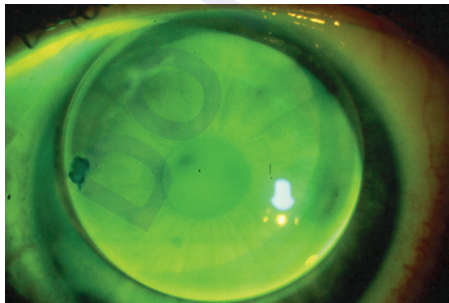


FIGURE 5 An alignment GP lens-to-cornea fitting relationship as demonstrated with fluorescein application.

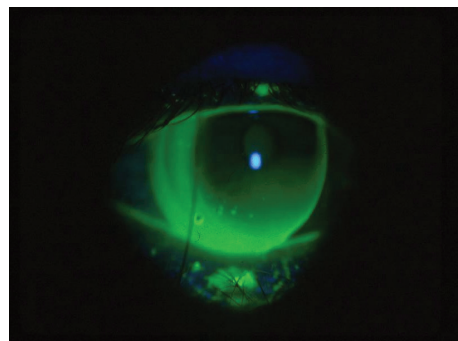


FIGURE 6 An apical bearing or flat GP lens-to-cornea fitting relationship as demonstrated with fluorescein application.

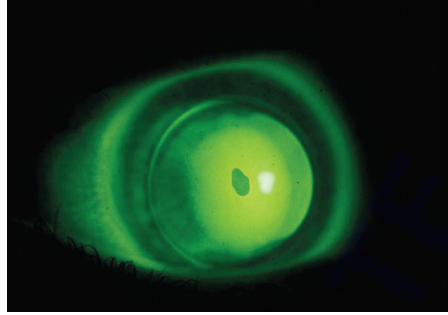


FIGURE 7 An apical clearance or steep GP lens-to-cornea fitting relationship as demonstrated with fluorescein application.

that a lens that is fitted flatter than “K” will induce a minus power lacrimal lens; for example, if a lens was fitted 0.50 D flatter than “K,” a -0.50 D lacrimal lens would be present and $+0.50$ D would need to be incorporated into the contact lens power to compensate for the lacrimal lens. Likewise, a lens that is fitted steeper than “K” will induce a plus power lacrimal lens; for example, if a lens was fitted 0.50 D steeper than “K,” a $+0.50$ D lacrimal lens would be present and -0.50 D would need to be incorporated into the contact lens power to compensate for the lacrimal lens.

The actual power can be derived in several ways. For example, the GP lenses were to be ordered empirically (i.e., from calculations only, not via the application of a diagnostic lens), the following formula can be used for predicting the power to be ordered

$$F_{cl} = F_{cp} - LLP$$

where F_{cl} = power of the contact lens
 F_{cp} = power of the refraction at the corneal plane
 LLP = lacrimal lens power

Using this formula, if the patient had the following refraction and keratometry values

Refraction: $-3.00 - 0.75 \times 180$

Keratometry: 43.50 at 180; 44.25 at 090

and a base curve radius was selected 0.75 D flatter than “K” or 42.75 (7.89 mm), the following GP contact lens power would be predicted:

$$F_{cl} = -3.00 - (-)0.75 = -2.25 \text{ D}$$

If the base curve radius was selected to be 0.50 D steeper than “K,” the predicted contact lens power would be

$$F_{cl} = -3.00 - 0.50 = -3.50$$

This example is shown in Fig. 8. An easy acronym to remember when factoring the tear lens into the contact lens power is “SAM-FAP” (i.e., “steep-add-minus” and “flat-add-plus”).

The most accurate method of determining contact lens power is via the use of a diagnostic lens. The formula for deriving the power when a diagnostic lens is placed on the eye is very simple

$$F_{cl} = F_{dl} + OR$$

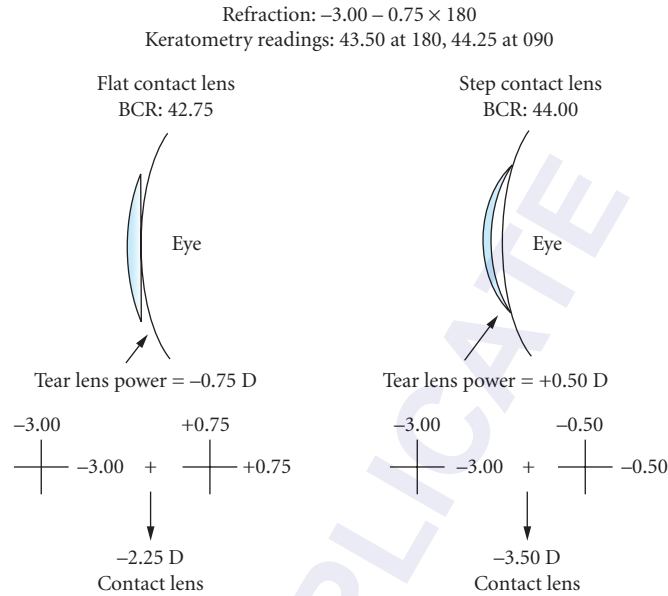


FIGURE 8 The “SAM (steep-add-minus)-FAP (flat-add-plus)” acronym as it applies to power determination based on the tear lens power.

where F_{cl} = final contact lens power
 F_{dl} = power of the diagnostic lens
 OR = spherical overrefraction

For example, if a patient has the following spectacle refraction and keratometry values

Refraction: $-6.00 - 0.50 \times 180$

Keratometry: 43.50 at 180; 44.00 at 090

If a diagnostic GP lens with a power of -3.00 D and a base curve radius equal to 43.25 (7.80 mm) was selected, the predicted contact lens power would equal

$$\begin{aligned}
 F_{cl} &= F_{cp} - LLP \\
 &= -5.50 \text{ (refraction at corneal plane)} - (-)0.25 \\
 &= -5.25 \text{ D}
 \end{aligned}$$

Once the contact lens has been evaluated and an overrefraction performed, it can be determined if the actual contact lens power to be ordered is equal to the predicted. If the spherical overrefraction was -2.25 D, this would be the case as

$$\begin{aligned}
 F_{cl} &= F_{dl} + \text{OR} \quad \text{or} \\
 &= -3.00 + (-)2.25 \\
 &= -5.25 \text{ D}
 \end{aligned}$$

As this example clearly demonstrates, the two most important factors when determining the power of a GP lens are effective power (if the ametropia is ≥ 4 D) and lacrimal lens power. If the contact lens power, as determined by overrefraction, deviates from the predicted power, the former is typically more accurate as inaccuracies in refraction and keratometry can impact the predicted value. The final lens order for this patient could be as follows:

Power: -5.25 D
 OAD/OZD: 9.40/8.00 mm
 BCR: 7.80 mm
 SCR/W: 8.80/0.3 mm
 PCR/W: 10.80/0.3 mm
 CT: 0.12 mm
 Material: Boston ES (fluoro-silicone/acrylate)

It is also important to understand how the lens design impacts the sagittal depth relationship between contact lens and cornea. As the optical zone diameter represents the linear equivalent to the base curve radius, if an alignment lens-to-cornea fitting relationship exists but the patient's pupil diameter is quite large, necessitating a larger optical zone diameter to optimize vision, the new lens should exhibit an apical clearance or steep fitting relationship due to the increase in sagittal depth. Essentially the steepest section of the contact lens, the base curve radius, has become larger. Therefore, the base curve radius should be designed to be flatter to compensate for this change in optical zone diameter; typically a 0.25 D change for every 0.5 mm change in OZD should maintain an alignment lens-to-cornea fitting relationship.

Residual Astigmatism

A factor that has to be considered when determining whether a patient should be fit into GP lenses is their residual astigmatism. This is defined as¹²

Total ocular astigmatism – astigmatism as measured by keratometry

In the previous examples, the refractive and anterior corneal astigmatism (as measured with a keratometer) were the same; therefore, no residual astigmatism would be predicted. If, however, a patient had the following refractive information:

Refraction: $-2.00 - 1.50 \times 180$
 Keratometry: 42.50 at 180; 43.00 at 090

the predicted residual astigmatism would equal the refractive (i.e., total ocular) astigmatism minus the keratometric astigmatism or

$$-1.50 \times 180 - (-) 0.50 \times 180 = -1.00 \times 180$$

As a GP lens only corrects anterior corneal astigmatism, this residual astigmatism would not be corrected and could result in symptoms of blurred vision and eyestrain. Typically, when there is -0.75 D or greater of residual astigmatism, a spherical GP lens is not recommended. Fortunately, this is quite uncommon.

Soft Lens Power

For spherical soft lenses, the primary factor to consider would pertain to effective power considerations when the ametropia equals or exceeds 4 D. As the back surface of a soft lens essentially conforms to—or drapes over—the anterior corneal surface, the power of the lacrimal lens is approximately

plano or afocal.¹³ This conformation of lens to the topography of the cornea has been termed “flexure.” Changes in base curve radius, diameter, or any other parameters do not appear to induce any lacrimal lens power. If the lens dehydrates during wear, the refractive index of the material will increase which will increase the magnitude of both plus and minus corrections, and the base curve radius will steepen and the lens will increase in minus power.¹⁴ In addition, it has been reported that as the cornea steepens, the amount of back vertex power is predicted to increase for both plus and minus lenses.¹⁰ Total flexure or conformation to the cornea may not be realistic with soft lenses, notably for plus power lenses.¹⁵

Another factor that impacts soft lens power as well as other design parameters pertains to the amount of expansion that occurs when the lens is hydrated (i.e., from the dry to wet state). According to Douthwaite¹³ a -3.25 D 59 percent water content soft lens with a 13 mm overall diameter, 8.60 mm base curve radius, and center thickness of 0.085 mm will have—prior to hydration—a -5.82 D power, 6.10 mm base curve radius, 9.22 mm overall diameter, and 0.060 mm center thickness. Therefore, as a result of all of these factors (dehydration, incomplete flexure, and the parameter changes when hydrated), it cannot be assumed that the labeled lens power of a soft contact lens will represent the true power of the lens on the eye.

Toric Lens Powers

Soft Toric Lenses As spherical soft lenses do not correct for the refractive (ocular) astigmatism, lenses incorporating this correction on the front surface are termed “soft toric lenses.” As it has been estimated that about 45 percent of interested contact lens patients have refractive astigmatism greater than or equal to 0.75 D¹⁶ if the criterion for significant astigmatism is a minimum of 0.75 D, then close to half of the spectacle-wearing population are potential candidates for soft toric lenses.¹⁷

The astigmatism is typically corrected onto the vertical meridian of the front surface and the lenses use several types of stabilizing techniques that will minimize any rotation of the lens on the eye. Traditionally these designs were prism ballasted inferiorly to weight the lens. Several newer techniques have resulted in the necessary thickness differential between the central, superior, and inferior sections of the lens without prism use.

The cylinder axis of the soft toric lens should align with the axis of the refractive astigmatism of the eye. For example, if a patient has a refractive error of $-2.00 - 0.75 \times 180$, this should be the powers of the soft toric lens fit to this patient. The cylinder axis is typically identified by a laser mark; often marks are also present 20° to the left and to the right of the axis. If the axis shifts or rotates with the blink, a new lens should be attempted which compensates for this rotation. The acronym LARS (left-add; right subtract) is valuable in these cases (Fig. 9). If the laser mark shifts 15° to the observer’s left, a new lens with an axis at 15° should be attempted; likewise if the lens shifts 15° to the observer’s right, a lens with an axis of 165° should be attempted.

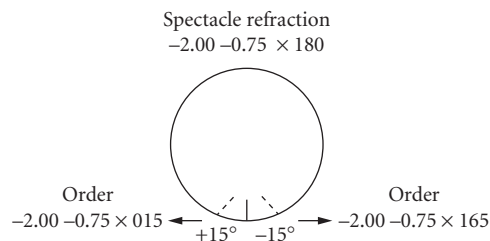


FIGURE 9 The “LARS” (i.e., left-add; right-subtract) concept as it pertains to determining axis for a soft toric lens.

If lens rotation is minimal, particularly with high spherical refractive errors (i.e., ≥ 5 D) and low astigmatic corrections (i.e., ≤ 1.25 D), little-to-no impact on vision should result. On average, relatively stable soft toric lenses usually average 0 to 10° of nasal rotation.¹⁸ However, when the rotation is significant, a “crossed cylinder” problem is presented. When the axis of cylinder of a correcting contact lens does not align with the best refractive correction, an overrefraction will result in cylinder at an axis oblique to that of the axis of best correction.¹⁰ Residual cylinder will be present in the overrefraction and this value will increase as axis rotation increases and with higher refractive cylinder corrections. The extent of the shift in position of the resultant plus cylinder axis for two cylinders of equal power but opposite sign has been calculated by Pascal.¹⁹ The formula $(90 + a)/2$ is used, where “ a ” is the angular discrepancy between the two combined cylinders. Therefore, the resultant plus cylinder axis is 45° from the midpoint between the axes of the two combined cylinders. The axis of the resultant plus cylinder of the crossed cylinder combination will appear on the side of the axis of the original plus cylinder opposite to that of the original minus cylinder. If the axis of ocular astigmatism is, for example, 90° and a minus cylinder of equal power is placed at axis 70° , or 20° away, the resultant plus cylinder will appear 55° [i.e., $(90 + 20)/2$] from 070 on the opposite side, or “x125.” The resultant plus cylinder will require a minus cylinder correction.

GP Bitoric Lens Powers Spherical GP wearing patients who exhibit a high amount of corneal astigmatism (typically ≥ 2.50 D) will not show the alignment fluorescein pattern demonstrated in Fig. 5 but will, instead, show a “dumbbell-shaped” pattern as shown in Fig. 10. In a high with-the-rule astigmatic patient, excessive edge clearance will exist in the steeper vertical meridian and bearing along the flatter horizontal meridian. As this fitting relationship often results in lens decentration on the eye, more edge awareness and corneal desiccation (dryness-induced epithelial cell loss), a lens that has a toric back surface with curvatures approximating the corneal curvatures in the two principle meridians is indicated. This design—termed bitoric lens design—essentially consists of two curvatures and two corresponding powers. A common method of selecting base curve radii is the Mandell-Moore bitoric lens guide philosophy.²⁰ This is provided in Table 5. For example, if the keratometry readings were

42.00 at 180 and 45.00 at 090

based on the Mandell-Moore base curve radius selection guide for three diopters of corneal cylinder, the flatter base curve radius would equal 0.25 D flatter than flat “ K ” or 41.75 D and the steeper base curve radius would equal 0.75 D flatter than steep “ K ” or 44.25 D. The base curve radii are selected to be flatter than their respective corneal curvatures due to the flattening of the cornea from center to periphery. The greater amount of flattening in the steeper meridian creates a slight amount of

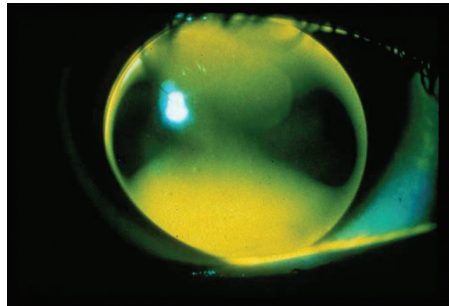


FIGURE 10 An astigmatic or “dumbbell-shaped” fluorescein pattern of a spherical GP lens on a highly with-the-rule astigmatic patient.

TABLE 5 Mandell-Moore Base Curve Selection

Corneal Cylinder	Flat Meridian	Steep Meridian
2.0 D	“On K”	0.50 D Flatter than “K”
2.5 D	0.25 D Flatter than “K”	0.50 D Flatter than “K”
3.0 D	0.25 D Flatter than “K”	0.75 D Flatter than “K”
3.5 D	0.25 D Flatter than “K”	0.75 D Flatter than “K”
4.0 D	0.25 D Flatter than “K”	0.75 D Flatter than “K”
5.0 D	0.25 D Flatter than “K”	0.75 D Flatter than “K”

toricity which has been found to enhance tear exchange. Determining the lens powers requires two tear lens calculations (vertical and horizontal meridians), not one as with conventional GP lenses. For example, if the patient’s refraction in the above case was $Pl - 3.00 \times 180$, the powers to be ordered for each meridian would consist of the following

$$F_{cph} - LLP \quad \text{and} \quad F_{cpv} - LLP$$

where F_{cph} = refractive power at the corneal plane in horizontal meridian and F_{cpv} = refractive power at the corneal plane in vertical meridian

$$F_{cph} - LLP = Pl - (-0.25) = +0.25 \text{ D}$$

$$F_{cpv} - LLP = -3.00 - (-0.75) = -2.25 \text{ D}$$

Therefore, the final powers of this bitoric lens would be $+0.25/-2.25$. This can also be diagrammed via the use of optical crosses showing the refractive power per meridian in combination with the correction for the lacrimal lens power (i.e., using “SAM-FAP”). This is shown in Fig. 11.

Another method of determining the base curve radii and powers without the use of optical crosses is to use the Mandell-Moore guide. An example is provided in Fig. 12.

An alternative method, which appears to be equally popular and effective, is to fit the patient with bitoric diagnostic lenses.²¹ In this case, a spherical overrefraction is typically all that is required and this value is added to the diagnostic lens power in each meridian. If, for example, in the above case, a bitoric diagnostic lens is selected that has base curve radii of 41.50/44.50 and powers of $Pl/-3.00 \text{ D}$, and the overrefraction is $+0.50 \text{ DS}$, the final powers would be

$$\text{Plano} + (+) 0.50 = +0.50 \text{ D}$$

$$-3.00 + (+) 0.50 = -2.50 \text{ D}$$

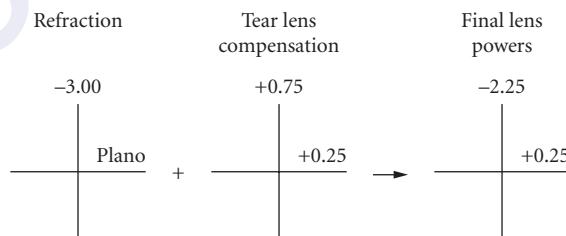


FIGURE 11 The use of optical crosses to show tear lens compensation in determining the final powers of a bitoric lens.

Mandell-Moore Bitoric Lens Guide - Per Eye

1. Keratometry

2. Spectacle Rx (Minus Cyl Form)

	Flattest K	Sphere Power	Steepest K	Sph+Cyl Power
3. Enter K	<input type="text" value="42.00"/>		<input type="text" value="45.00"/>	
4. Enter Spectacle Power		<input type="text" value="plano"/>		<input type="text" value="-3.00"/>
5. Vertex Adjust Line 4		<input type="text" value="plano"/>		<input type="text" value="-3.00"/>
6. Insert Fit Factor	<input type="text" value="(-) 0.25"/>	<input type="text" value="(+) 0.25"/>	<input type="text" value="(-) 0.75"/>	<input type="text" value="(+) 0.75"/>
Add Lines	3&6	5&6	3&6	5&6
7. Final CL Rx	<input type="text" value="41.75"/>	<input type="text" value="+0.25"/>	<input type="text" value="44.25"/>	<input type="text" value="-2.25"/>
	Base Curve	Power	Base Curve	Power

Bitoric Lens Fit Factor

Corneal Cyl	Fit Flat Meridian	Fit Steep Meridian
2.0 Diopters	On K (0 D)	0.50D Flatter
2.5 Diopters	0.25D Flatter	0.50D Flatter
3.0 Diopters	0.25D Flatter	0.75D Flatter
3.5 Diopters	0.25D Flatter	0.75D Flatter
4.0 Diopters	0.25D Flatter	0.75D Flatter
5.0 Diopters	0.25D Flatter	0.75D Flatter

FIGURE 12 The Mandell-Moore bitoric fitting guide.

In this case there is no residual astigmatism. In these examples the bitoric lens has a spherical power effect in that it can rotate on the eye with no impact on vision as the corresponding vergence will be present to correct for the change in corneal curvature. (See Fig. 13.)

If the patient has significant residual astigmatism, a sphero-cylindrical overrefraction will be necessary and the overrefraction in a meridian will be added to the corresponding power in the bitoric diagnostic lens. For example, if a patient had the following refractive information

Refraction: +1.00 - 4.00 × 180

Keratometry: 42.00 at 180; 45.00 at 090

Diagnostic bitoric lens: 41.50/44.50 Pl/-3.00D

Overrefraction: +1.50 - 1.00 × 180

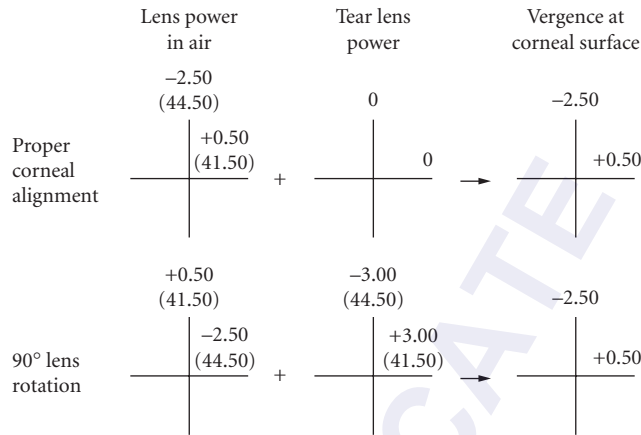


FIGURE 13 An example of how a bitoric lens can have spherical power effect optics. If the lens rotates 90°, the compensating vergences will be present.

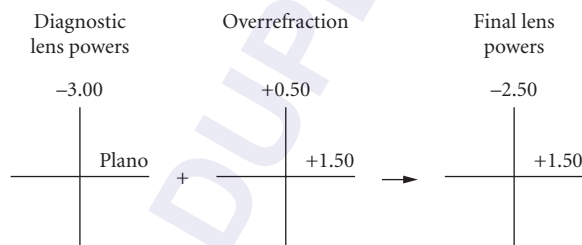


FIGURE 14 Determining the final powers of a bitoric lens in which a sphero-cylindrical overrefraction simply pertains to adding the overrefraction in a specific meridian to the power of the diagnostic lens in that meridian.

The final powers are determined by adding the overrefraction in the 180 meridian (+1.50) to the diagnostic lens power (Plano) and then adding the overrefraction in the 090 meridian [$+1.50 + (-)1.00 = +0.50$] to the diagnostic lens power in that meridian (-3.00 D). This is shown in Fig. 14.

20.5 OTHER DESIGN CONSIDERATIONS

Aspheric Lens Designs

An aspherical surface has been defined as a surface which is not spherical.²² In contact lens design terminology, it was originally indicated by Feinbloom²³ that “(a) an ellipsoid represents a better approximation of the form of the surface of the cornea of the human eye, and (b) contact lenses with inner elliptical surfaces represent a marked improvement in the comfort and wearing time by the patient.” This elliptical corneal shape consists of a central spherical zone which progressively flattens as the limbus is approached. The progressive flattening is representative of ellipses—and conic sections—in general.²⁴

The shapes of conic sections are described mathematically by their eccentricity (e). As it pertains to conic sections, eccentricity can be defined as deviating from a circular path. A spherical central

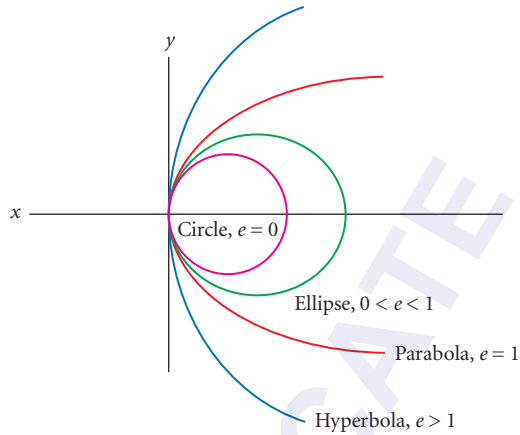


FIGURE 15 Different conic sections which have the same prolate apical radius but differ in eccentricity or e -value.

area of a given radius of curvature is combined in a model system with varying rates of paracentral flattening in the same scale as that of the cornea that it represents. The “ e ” increases, as does the rate of paracentral flattening. This can be demonstrated graphically by plotting the conic sections on the same set of axes with the vertex of each figure having the same radius of curvature at the origin of the coordinates (see Fig. 15). A circle has an “ e ” value of 0 with its radius of curvature the same for all points on the curve and the meridians all equal in length. An ellipse has an “ e ” value between 0 and 1; a parabola has an “ e ” value equal to 1; a hyperbola has an “ e ” value greater than 1. As will be discussed, most of the aspheric lens designs in common use mimic the average corneal eccentricity which is a prolate ellipse with an “ e ” value of approximately 0.45.^{10,24} A prolate ellipse is produced by elongating the horizontal meridian of a circle (Fig. 16). If “ a ” is the major meridian of the ellipse and “ b ” the minor meridian, the eccentricity is calculated by the following equation:

$$e = \frac{\sqrt{1 - b^2}}{a}$$

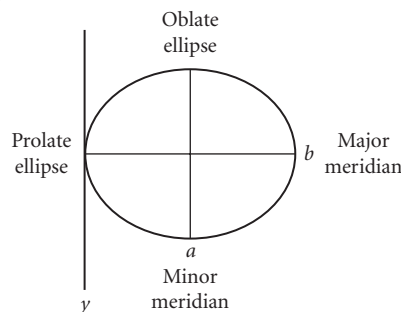


FIGURE 16 The major meridian of a typical ellipse is the “ a ” axis; rotation around this axis will produce a prolate ellipsoid, characterized by increasing radius of curvature outward from the origin. An oblate ellipsoid is produced by rotation around the minor axis “ b .” This is characterized by a decrease in radius of curvature outward from the origin.

Several GP aspheric lens designs are in common use. Some are manufactured with a spherical base curve radius and an aspheric periphery. In many of these designs, the aspheric section is tangential to the base curve radius, therefore creating an aspheric continuum or continuous curve. Several aspheric designs have a totally aspheric posterior surface; however, the posterior optical zone and the periphery are two different curves. One such design that has enjoyed some success in the United States is the Envision lens (Bausch & Lomb). Originally developed by Hanita Contact Lenses in Israel, this design has an elliptical posterior optical zone of 0.4 and a hyperbolic periphery that is tangential to the posterior optical zone.

Aspheric lens designs have become increasingly popular in the correction of presbyopia. Both center-near and center-distance power soft lenses are currently used which involve continuous change in power from the lens axis to the peripheral section of the lens, therefore creating a multifocal effect. These designs have the benefit of progressive power change resulting in distance, intermediate, and near image correction. This power change is typically on the front surface of the lens and these designs have been termed “simultaneous vision” lenses as the different powers are in front of the pupil at the same time; therefore, some compromise in vision may exist with these designs at any desired viewing distance. This type of presbyopic correction has been found to be preferred by 76 percent of individuals who wore both this correction and monovision (i.e., one eye optimally corrected for distance vision; the other eye optimally corrected for near vision).²⁵ GP aspheric multifocal designs have enjoyed greater success as a result of the rigidity, optical quality, and ability to provide high add powers. Although some designs are front surface aspherics with an increase in near power from center to periphery, most designs in common use today are center-distance back surface designs although this varies from region to region worldwide. These designs have much greater eccentricity than single vision designs, ranging from high “e” value elliptical to hyperbolic back surface designs.^{26,27} This type of design has been found to result in similar visual performance to progressive addition spectacle lenses while performing significantly better than both monovision and soft lens multifocals.^{28,29} Although the rapid increase in plus power required to provide optimal near vision for patients with high add requirements can compromise distance vision, recently introduced GP multifocal designs have addressed this issue by incorporating some of the add power in a paracentral annular region on the front surface of the lens (see Fig. 17). These designs also rely on some upward shifting of the lens—or translation—when the patient views inferiorly to take advantage of the paracentral/peripheral near power.

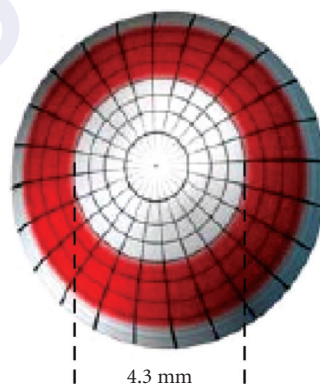


FIGURE 17 The Essentials CSA lens design with the near add power on the front surface shown in red.

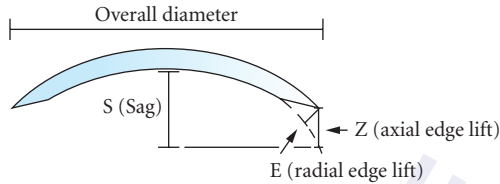


FIGURE 18 Axial versus radial edge lift.

Posterior Peripheral Curve Design Systems and Calculations

As previously discussed, there are several terms used to describe the peripheral curve radii and width, including secondary curve, intermediate curve, peripheral curve, and aspheric peripheries. Edge lift or edge clearance are terms used to indicate the distance from the GP lens edge to the cornea. Edge lift pertains to the distance between the extension of the tangent of the base curve radius and the absolute edge after the addition of peripheral curves. Radial edge lift (REL) is measured normal to the base curve radius extension (i.e., the distance from the lens edge perpendicular to an extension of the base curve radius.⁹ (See Fig. 18.) Axial edge lift is measured parallel to the lens optical axis or the vertical distance from the lens edge to an extension of the base curve radius. Axial edge lift is often used and values of 0.10 to 0.15 mm are typically recommended.³⁰ In a tricurve design, for example, the peripheral curve often contributes approximately two-thirds of the overall axial edge lift. AEL values can be determined via the tables presented in Musset and Stone.³¹

Table 6 provides some representative AEL values for different designs. A more accurate estimate of the distance between the lens edge and the cornea is the actual edge clearance. As this does not pertain to an extension of the base curve radius, this value would be less than the axial edge lift. An ideal value for axial edge clearance has been found to be 0.08 mm.³² Although this has traditionally been a difficult parameter to measure, Douthwaite³³ provides a very good overview of his recent efforts in calculating axial edge clearance and his text should be referred to for anyone desiring more information on this topic.

TABLE 6 Recommended Secondary and Peripheral Curve Radii

Axial Edge Lift = 0.10 mm

Secondary Curve Width = 0.3 mm

Peripheral Curve Width = 0.3 mm

Overall Diameter = 9.0; optical zone diameter = 7.8 mm

Base Curve Radius (BCR): varies from 7.5 to 8.3 mm

If secondary curve radius (SCR)/width contributes 0.04 and the peripheral curve radius (PCR)/width contributes 0.06 to the overall axial edge lift, the following values would be calculated:

BCR (mm)	SCR (mm)	PCR (mm)
7.50	9.00	10.20
7.70	9.30	10.70
7.90	9.60	11.20
8.10	10.00	11.70
8.30	10.30	12.20

Source: From Ref. 32.

For practical purposes, if the peripheral curve is designed or modified to be flatter or wider, the resulting edge clearance is greater. Flattening the base curve radius will also “lift” the periphery away from the cornea and increase edge clearance. Reducing the optical zone diameter while keeping the overall diameter constant will increase edge clearance as the steepest part of the lens (i.e., the base curve radius) is reduced and the overall peripheral curve width is being increased. Conversely, edge clearance can be reduced by steepening the peripheral curve radius, reducing the peripheral curve width, steepening the base curve radius, or increasing the optical zone diameter.

Aberrations and Contact Lenses

As is well established, there are several aberrations induced by the optics of the eye which can influence the quality of vision. Lower-order aberrations include tilt or prism, astigmatism, and defocus. The most common high-order aberrations found in human eyes are spherical aberration and coma. Whereas in spectacle lens designs, the primary problems pertain to minimizing the effects of oblique astigmatism and curvature of field to optimize quality of vision in peripheral gaze, this is not a significant problem with contact lenses as the lens moves with the eye in all areas of gaze. However, spherical aberration and coma can be problematic with contact lens wear, notably when the lens decenters.²² The latter can be especially important. If a lens is decentered 0.5 mm, an amount of coma approximately equivalent to the amount of spherical aberration is induced. If this decentration is 1 mm, a magnitude of coma equal to twice the level of spherical aberration is induced.³⁴ A well-centered lens is especially important as recent research has found that coma is, in fact, more visually compromising than spherical aberration.³⁵

The introduction of sensitive, sophisticated aberrometers has been beneficial in evaluating the relationship of contact lenses and aberrations. Several studies have evaluated the optical quality of eyes wearing different types of contact lenses. Certain types of soft lenses (i.e., manufactured via cast-molding or spin-casting) induced more high-order aberrations, such as coma and spherical aberration as measured via aberrometry.³⁶ In comparing both soft and GP lenses it was found that, whereas both soft and GP lenses induce more aberrations for the eyes that have low wavefront aberrations, soft lens wear tends to induce more higher-order aberrations and GP lens wear tends to reduce higher-order aberrations.^{37,38} CRT lenses used to dramatically reshape the cornea to reduce myopia when worn at night have been found to increase higher-order aberrations, especially spherical aberration.^{39,40}

To correct for wavefront aberrations it is important to slow the propagation of the wavefront in areas of the pupil where it is advanced and speed up the propagation of the wavefront in areas where it's retarded.⁴¹ This can be accomplished in a contact lens with localized changes in lens thickness. Essentially the lens must be thinner in areas where the eye it is correcting shows the greatest delays in wavefront and thickest where the eye's wavefront demonstrates its greatest advancement. This is shown in Fig. 19 in which the introduction of an aberration-correcting contact lens in both myopic (Fig. 19a) and comatic (Fig. 19b) wavefronts is demonstrated.⁴¹

Several contact lens designs, especially soft lenses, have been recently introduced with some aberration-correction ability. Keratoconic eyes reportedly have higher-order aberration levels as much as 5.5 times higher than normal eyes;⁴² however, custom soft lenses have been found to provide a threefold reduction in aberration on keratoconic eyes.^{43,44} A few designs have attempted to introduce a constant level of spherical aberration at all lens powers that is equal in magnitude but opposite in sign to the population mean for human eyes.⁴² Some designs achieve this over part of the dioptric power range (i.e., PureVision from Bausch & Lomb); the Biomedics Premier lens from CooperVision achieved this goal over the entire dioptric range.^{42,45} The Biomedics XC lens from CooperVision has been introduced with the strategy of having zero spherical aberration at all lens powers and, therefore, will neither add to or correct the eye's own spherical aberration. These lenses have the advantages of not introducing large amounts of spherical aberration often provided by standard spherical soft lenses and as these lenses have no spherical aberration, they will not introduce coma if they decenter.

It's evident that—although documented clinical success has not been definitely established with these designs—the future looks promising for aberration-controlling lenses. Likewise, manufacturers will introduce designs that either center well or decenter by a fixed amount. The introduction of

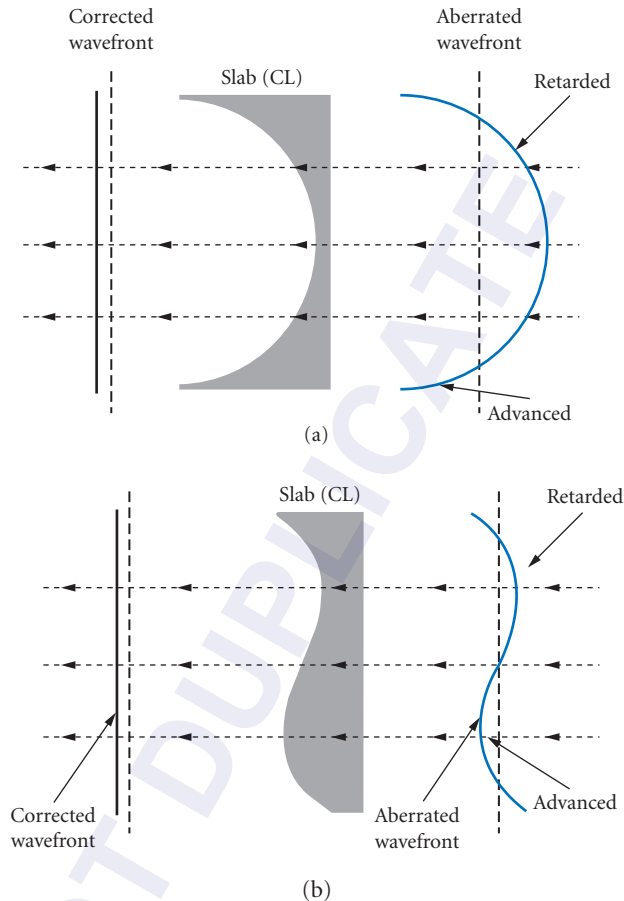


FIGURE 19 Schematic showing the correction of a (a) myopic and (b) comatic wavefront by the introduction of an aberration-correcting contact lens. (From Ref. 41.)

contact lens-only aberrometers, such as the ClearWave (AMO Wavefront Sciences), allows for a greater ability to monitor the aberration characteristics of these designs and ultimately introduce designs that are more effective in correcting for the aberrations of the eye.

20.6 CONVERGENCE AND ACCOMMODATION EFFECTS

Convergence

Spectacles which are well centered for distance vision but which are also used for all directions of gaze can induce a prismatic effect when the gaze is not straightahead. When viewing at near, a spectacle-wearing hyperopic patient will experience a base-out effect and a spectacle wearing myopic patient will experience a base-in effect. The fact that the optical center of a contact lens remains (or ideally should remain) centered in front of the eye in different directions of gaze results in several benefits

for the contact lens wearer via reducing or eliminating the following prismatic effects common to spectacle wear including:^{10,46}

1. Decreased near convergence demand for bilateral hyperopic patients and increased near convergence demand for bilateral hyperopic patients.
2. Base right prism and base left prism for bilateral myopic patients, and base left prism and base right prism for bilateral hyperopic patients, in right and left gaze, respectively.
3. Vergence demand alterations in anisometropia or antimetropia, required for right and left gaze.
4. Vertical prismatic effects in up and down gaze, and imbalances in down gaze resulting from anisometropia or antimetropia.

There are a few potential problems, however, when changing a patient from spectacles to contact lenses. If the patient requires lateral prism correction to maintain alignment and fusional ability, this is not possible in contact lenses. For the young child this is an important consideration because of the cosmesis and freedom of vision provided by contact lenses. However, only after successful binocular vision training and/or surgical intervention should contact lenses be prescribed. Likewise, patients requiring correction for a vertical deviation will (almost always) not achieve that correction in a contact lens. It is possible to prescribe base down prism in a contact lens; however, this is unlikely to solve the problem unless prescribed binocularly as asthenopic complaints are likely due to the prism differential between the two eyes.

Finally, for a given distance, the contact lens-wearing hyperopic patient exhibits less convergence and the myope exerts more convergence than with spectacles. The binocular myopic patient loses the “base-in” effect received from spectacles and, therefore, if exophoria is present and the patient has an abnormally long or remote near point of convergence, eyestrain and even diplopia may result. This same problem may result for the esophoric hyperopic patient who loses the “base-out” effect provided by spectacles and—if borderline fusional divergence ability was initially present—the decreased convergence demand at near induced by contact lens wear may result in compromised binocularity. Fortunately, as will be discussed, accommodative vergence tends to compensate for the differences in vergence demands between the two modes of correction.

Accommodative Demands

The review of effective power considerations and, specifically, vertex distance, becomes especially appropriate as it pertains to accommodative demand. As will be discussed, the vergence entering the eye differs with light coming from a distance object versus a near object. Accommodative demand represents the difference between distance and near refractive corrections.¹⁰

Accommodative demand can best be described via reviewing the differences in demand between a hyperopic and myopic patient, both viewing at a 40 cm distance. Figure 20 shows the vergence of light for a +5 D spectacle-wearing patient, using a vertex distance of 15 mm. A vergence of 0 would arrive at the spectacle lens from a distant object. Via use of the effective power equation, the power at the corneal plane would equal

$$\begin{aligned}
 F(c) &= \frac{F(sp)}{1 - dF(sp)} \\
 &= + \frac{5.00}{1 - 0.015(+5.00)} \\
 &= +5.41 \text{ D}
 \end{aligned}$$

This value can also be obtained by simply subtracting the difference in focal lengths. The spectacle lens has a focal length of $1/+5.00$ or 0.2 m (or 200 mm). At the corneal plane relative to the far point of the eye, the focal length is reduced by 15 mm or equals 185 mm (0.185 m). $1/0.185 \text{ m} = +5.41 \text{ D}$

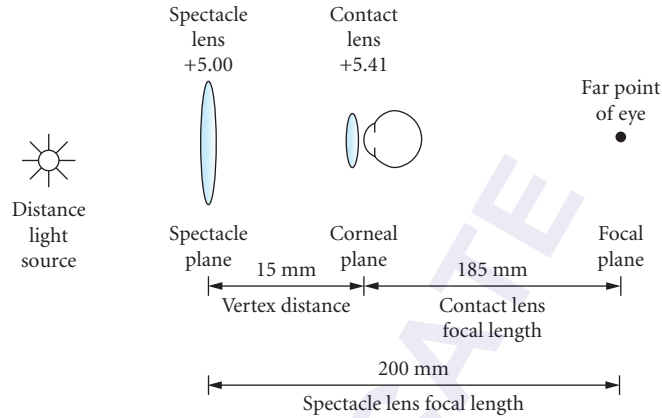


FIGURE 20 The power and corresponding focal lengths of a +5.00 D spectacle lens at the corneal plane when viewing a distant object.

which would equal the vergence of light at the corneal plane. For the emmetrope, it can be assumed that the vergence of light at the cornea is zero.

When viewing at a distance of 40 cm, a vergence of -2.50 D (i.e., $1/-0.40$) would enter the spectacle lens (Fig. 21). A vergence of $+2.50$ D (i.e., $-2.50 + (+)5.00$) would exit the spectacle lens. To determine the power at the corneal plane, the effective power equation can be used.

$$\begin{aligned}
 F(c) &= \frac{F(sp)}{1 - dF(sp)} \\
 &= + \frac{2.50}{1 - 0.015(+2.50)} \\
 &= +2.60 \text{ D}
 \end{aligned}$$

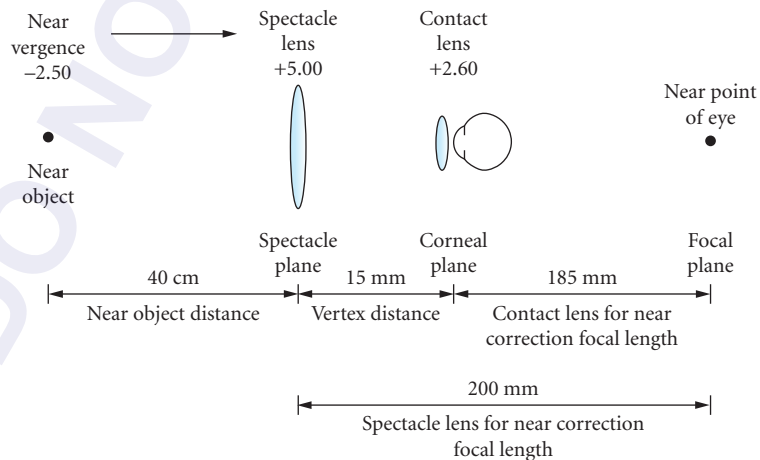


FIGURE 21 The power and corresponding focal lengths of a +5.00 D spectacle lens at the corneal plane when viewing a near object.

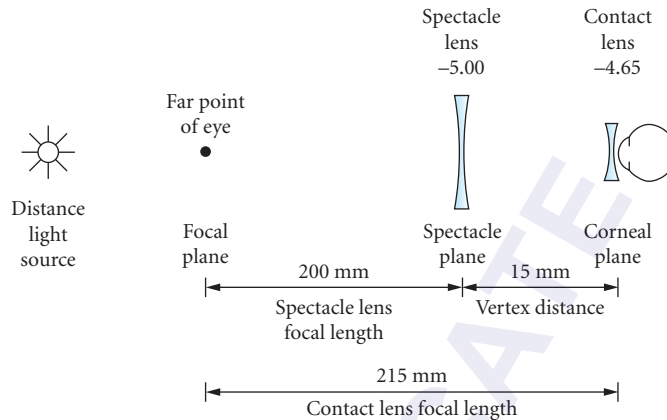


FIGURE 22 The power and corresponding focal lengths of a -5.00 D spectacle lens at the corneal plane when viewing a distant object.

Likewise, this value can also be determined by subtracting the difference in focal lengths. $1/2.50$ or 0.4 m (400 mm) provides the focal length of the power exiting the spectacle lens and, considering the 15 mm reduction at the corneal plane, the power at the corneal plane would be $1/0.385$ or $+2.60$ D. Therefore, the corneal accommodative demand for the 5 D hyperope would equal the distance accommodative demand minus the near accommodative demand or $+5.41 - (+)2.60 = +2.81$ D. This can be compared to the emmetrope who would have a corneal accommodative demand of $1/0.415$ (i.e., 40 cm plus the 15 mm vertex distance) or 2.41 D. Therefore, the hyperopic patient requires more accommodation than the emmetropic patients when viewing a near object with spectacles.

Figure 22 shows the vergence of light for a -5 D spectacle-wearing patient, using a vertex distance of 15 mm. A vergence of 0 would arrive at the spectacle lens from a distant object. Via use of the effective power equation, the power at the corneal plane would equal

$$\begin{aligned} F(c) &= \frac{F(\text{sp})}{1 - dF(\text{sp})} \\ &= \frac{-5.00}{1 - 0.015(-5.00)} \\ &= -4.65 \text{ D} \end{aligned}$$

This value can also be obtained by simply subtracting the difference in focal lengths. The spectacle lens has a focal length of $1/-5.00$ or -0.2 m (or -200 mm). At the corneal plane relative to the far point of the eye, the focal length is increased by 15 mm or equals -215 mm (-0.215 m). $1/-0.215$ m = -4.65 D which would equal the vergence of light at the corneal plane.

When viewing at a distance of 40 cm, a vergence of -2.50 D would enter the spectacle lens (Fig. 23). A vergence of -7.50 D (i.e., $-2.50 + (-)5.00$) would exit the spectacle lens. To determine the power at the corneal plane, the effective power equation can be used

$$\begin{aligned} F(c) &= \frac{F(\text{sp})}{1 - dF(\text{sp})} \\ &= \frac{-7.50}{1 - 0.015(-7.50)} \\ &= -6.74 \text{ D} \end{aligned}$$

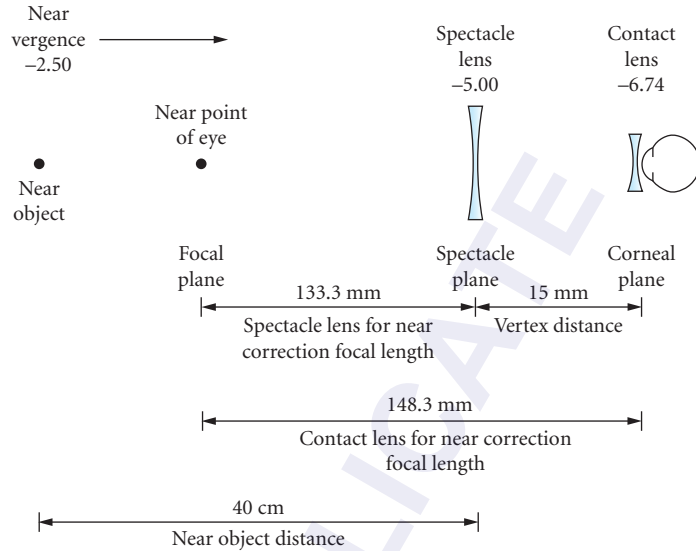


FIGURE 23 The power and corresponding focal lengths of a -5.00 D spectacle lens at the corneal plane when viewing a near object.

Likewise, this value can also be determined by subtracting the difference in focal lengths. $1/(-7.50$ or -0.1333 m (-133.3 mm) provides the focal length of the power exiting the spectacle lens and, considering the increase of 15 mm at the corneal plane, the power at the corneal plane would be $1/(-0.1483$ or -6.74 D. Therefore, the corneal accommodative demand for the 5 D myope would equal the distance accommodative demand minus the near accommodative demand or $-4.60 - (-)6.74 = +2.14$ D. This can be compared to the emmetrope who would have a corneal accommodative demand of 2.41 D; therefore, the myopic patient requires less accommodation than the emmetropic patient when viewing a near object with spectacles. The differences in corneal accommodative demand for myopic and hyperopic patients has been well documented.^{46–49} Table 7 shows the difference in corneal accommodative demand at the corneal plane as it pertains to myopic and hyperopic spectacle lens powers.¹⁰

As contact lenses are positioned at the corneal plane, they induce accommodative demands equivalent to that of emmetropia. Therefore, the emerging presbyopic hyperope who changes from spectacles to contact lenses may find that the near symptoms are reduced. However, the emerging

TABLE 7 Accommodative Demands at the Corneal Plane with Spectacle Lens Wear

Difference in Corneal Plane Accommodative Demand Compared to Emmetropia	Back Vertex Power of Hyperopic Spectacle Lens (D)	Back Vertex Power of Myopic Spectacle Lens (D)
+/- 0.25	+3.25	-3.87
+/- 0.50	+6.00	-8.37
+/- 0.75	+8.62	-13.75
+/- 1.00	+10.87	-20.87

*Assume a near distance of 40 cm and a vertex distance of 15 m.
Source: From Ref. 10.

presbyopic highly myopic patient should be advised that a separate near correction (or a contact lens multifocal) may be necessary to satisfy all vision demands.

Overall, it can be concluded that hyperopic-correcting spectacle lenses will require an increase in convergence over that required when uncorrected whereas myopic-correcting lenses reduce the convergence required.⁵⁰ However, as indicated in this section, myopic patients accommodate more and hyperopic patients accommodate less when wearing a contact lens correction versus spectacles. Therefore, the overall accommodation-convergence ratio is only minimally impacted.

20.7 PRISMATIC EFFECTS

Prism-Ballasted Contact Lenses

Prism within a contact lens is produced by varying the thickness from the superior to the inferior regions while maintaining the same front and back surface curvatures. There are several types of contact lenses for which prism base down is incorporated within the design. The most popular design, including such prism, pertains to segmented, translating bifocal GP lenses. These lenses include a distance zone in the upper half of the lens, a near zone in the bottom half of the lens and (with some designs) an intermediate power zone between the near and distance segments. Typically, 1.5 to 3.0Δ base down is incorporated within the lens to stabilize it such that the thicker inferior edge will interact with the inferior lid margin on inferior gaze and push the lens up (termed “translation”) such that the wearer will be viewing through the inferior near power zone when reading. Prism has also been used for stabilizing soft toric lenses as well as front toric GP lenses (i.e., used to correct residual astigmatism), although front toric GPs are little used due to the application of soft toric lenses in most of these cases and soft toric lenses tend to use other—reduced mass—options for stabilization.

The amount of prism within a contact lens can be computed using the following formula:¹⁰

$$P = \frac{100(n-1)(BT-AT)}{BAL}$$

where P = prismatic power in prism diopters (Δ)
 n = refractive index of prismatic lens
 BT = base thickness of prismatic component of lens
 AT = apex thickness of prismatic component of lens
 BAL = length of base-apex line or the diameter of the contact lens long the base-apex line of prism

For example, if the diameter of the lens was 9.0 mm, the thickness at the apex was 0.11 mm and the thickness at the base was 0.40 mm, and the refractive index of the material was 1.46, the amount of prism within this lens would be

$$\begin{aligned} P &= \frac{100(n-1)(BT-AT)}{BAL} \\ &= \frac{100(1.45-1)(0.40-0.11)}{9.0} \\ &= 1.48 \text{ or approximately } 1.50\Delta \text{ base down} \end{aligned}$$

Refractive power of a prismatic lens does vary along the base-apex line via the aforementioned formula for deriving back vertex power. This is the result of the fact that lens thickness is increased toward the base with an absence of change in surface curvature. Via Prentice’s rule (to be discussed), back surface power becomes less minus/more plus as thickness increases toward the base of the prism.

Unintentional (Induced) Prism

Prentice's rule can also be used for determining the prism induced by lens decentration on the eye. This is due to the fact that a contact lens is considered to be separated from the tear lens by a thin layer of air.⁵¹ Prentice's rule is

$$P = Fd$$

where P = prism induced in prism diopters
 F = dioptric power of lens
 d = decentration in centimeters

For example, if a +5.00 D lens decenters 2 mm, the resulting induced prism would be $+6.00 \text{ D} \times 0.2 = 1.2\Delta$ base down. If the other lens is well centered, this may be sufficient vertical imbalance to cause symptoms of eyestrain. Fortunately, soft lenses typically do not decenter on the eye due to their large size and the fact they drape over the cornea. GP lenses do move with the blink but typically, if decentration is present, it is essentially the same for both eyes. In cases of anisometropia, however, any decentration could result in asthenopic complaints. For example, if the patient was wearing a +2 D lens in the right eye and a +7 D lens in the left eye and both lenses decentered 2 mm inferiorly, the amount of induced prism would be OD: $+2.00 \times 0.2 = 0.4\Delta$; OS: $+7.00 \times 0.2 = 1.4\Delta$. The resulting difference would be 1.0Δ . Typically patients will not be symptomatic with less than 1Δ ; however, in this case, asthenopic complaints are possible.

20.8 MAGNIFICATION

It has been well established that when a spectacle lens is moved toward the eye, the retinal image decreases in size for a plus lens and increases in size for a minus lens. Therefore, with contact lenses both of these effects would tend to be advantageous. The hyperopic patient with a high prescription—and most certainly the aphakic patient—are typically very satisfied to find that objects have returned to nearly a normal size. Myopic patients are pleased to find that everything looks larger with contact lenses. Several methods have been used to assess magnification effects, including magnification of correction and relative spectacle magnification.

Magnification of Correction (Spectacle Magnification)

Spectacle magnification (SM) is the ratio of retinal image size of the corrected ametropic eye to the retinal image size of the same eye uncorrected.¹⁰ This has been used as an index of how corrective lenses alter retinal image size as compared by the patient before and after corrective lenses are placed on the eye. Spectacle magnification formulas include a power factor and a shape factor^{10,47,52–54}

$$\text{SM} = \underbrace{\frac{1}{1} - h(\text{BVP})}_{\text{Power factor}} \times \underbrace{\frac{1}{1} - (tn')F_1}_{\text{Shape factor}}$$

where SM = spectacle magnification or magnification of correction
 BVP = back vertex power of correcting lens (D)
 h = stop distance, from plane of correcting lens to ocular entrance pupil in meters
 (i.e., vertex distance +3 mm)
 t = center thickness of correcting lens (m)
 n' = refractive index of correcting lens
 F_1 = front surface power of correcting lens (D)

With both contact lenses and spectacle lenses the shape factor is nearly always greater than 1.0, representing magnification, due to their respective convex anterior surfaces (i.e., F_1 is a positive number). However, for contact lenses, the lacrimal lens needs to be considered. Therefore, for a contact lens, the following shape factor should be used:

$$\text{Shape factor} = \frac{1}{1} - (t/n')F_1 \times \frac{1}{1} - (t_L/n_L)F_L$$

where t = center thickness of correcting lens (m)
 n' = refractive index of correcting lens
 F_1 = front surface power of correcting lens (D)
 t_L = center thickness of lacrimal lens in meters
 n_L = refractive index of lacrimal lens (1.3375)
 F_L = front surface power of lacrimal lens in keratometric diopters

The power factor for contact lenses is essentially the same as that for spectacle lenses with the exception that, as the vertex distance is 0, the stop distance would be only equal to 3 mm.

In high myopia, a significant minification ($SM < 1.0$) of the retinal image occurs with spectacle correction; this minification is greatly reduced with contact lens wear. Therefore, the highly myopic patient changing from spectacles to contact lenses may comment that their vision is clearer with contact lenses.

The change in spectacle magnification when changing from spectacle to contact lens correction can be compared with the following formula:

$$\frac{\text{Contact lens power factor}}{\text{Spectacle lens power factor}} = 1 - h(\text{BVP})$$

where h = stop distance, from plane of correcting lens to ocular entrance pupil in meters
 (i.e., vertex distance +3 mm) and
 BVP = back vertex distance of spectacle lens (in D)

For the highly myopic (i.e., -10 D) patient, approximately 15 percent minification will be present with spectacles and only about 5 percent with contact lenses. For the aphakic patient, the magnification will be 25 to 30 percent with spectacles which would be reduced to 5 to 8 percent with contact lenses. Nevertheless, if the patient is a unilateral aphake wearing one contact lens, this image difference may still result in binocular vision problems. Therefore, intraocular lenses, with a stop distance of 0, is optimum in maintaining similar image sizes between the two eyes.

Relative Spectacle Magnification

Retinal image size can also be estimated via comparing the corrected ametropic retinal image to that of a standard emmetropic schematic eye. Relative spectacle magnification (RSM) is the ratio derived for spectacle lenses.^{10,47,52-54} Ametropia can be purely axial, resulting from the axial elongation of the globe, refractive, resulting from abnormal refractive component(s) of the eye or a combination of both factors.

Relative spectacle magnification for axial ametropia can be derived from the following equation

$$\text{RSM} = \frac{1}{1 + g(\text{BVP})}$$

where RSM = relative spectacle magnification for axial ametropia
 g = distance in meters from anterior focal point of eye to correcting lens; $g = 0$ if lens is 15.7 mm in front of the eye
 BVP = back vertex power of refractive correction (D)

with axial ametropia, g is equal to 0 or close to it for most spectacle lenses; therefore, the relative spectacle magnification is very close to 0. (See Fig. 24.)⁵¹ In theory, anisometropia of axial origin is best corrected with spectacles and, in fact, anisokonia problems may be encountered with contact lens wear. However, it has been suggested that in axial myopic eyes the retina will stretch to compensate

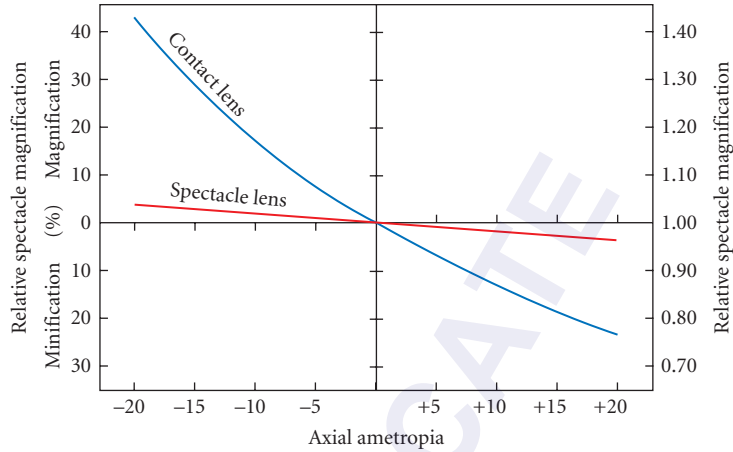


FIGURE 24 Relative spectacle magnification for axial ametropia when corrected by spectacles and contact lenses. (From Ref. 51.)

resulting in greater spacing between retinal receptors and receptive fields that effectively reduce the neurological image.⁵¹

Relative spectacle magnification for refractive ametropia can be derived from the following equation:

$$\text{RSM} = \frac{1}{1 - d(\text{BVP})}$$

where RSM = relative spectacle magnification for axial ametropia

d = stop distance in meters from correcting lens to entrance pupil (i.e., vertex distance +3 mm)

BVP = back vertex power of refractive correction (D)

This equation is the same for the shape factor for spectacle magnification. It can be observed in Fig. 25 that the RSM for contact lens wear is close to 1.00 but, conversely, it is quite minified for highly myopic lenses and very magnified for high plus lenses.⁵⁵

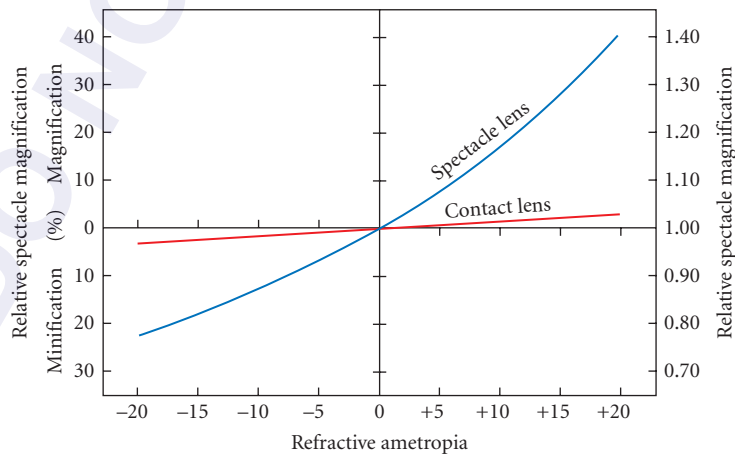


FIGURE 25 Relative spectacle magnification for refractive ametropia when corrected by spectacles and contact lenses. (From Ref. 51.)

20.9 SUMMARY

This chapter emphasizes the many optical benefits of contact lenses versus spectacles. In particular, the elimination of undesirable aberrations and induced prism when viewing away from the optical center and a lesser magnified (or minified) image size, it is important to emphasize that contact lenses will result in an increase in accommodative demand for the emerging presbyopic myope when changing from spectacles. In addition, children with binocular vision problems who benefit from the lateral prism in spectacles will not benefit from contact lens wear and may, in fact, be very poor candidates. This chapter also emphasizes contact lens optics that are in common use today. Determining the accommodative demand, induced prism and magnification are important in some contact lens wearers. Determination of contact lens power is important in all contact lens wearers.

20.10 ACKNOWLEDGMENTS

The author would like to acknowledge the contributions of Teresa Mathew and Jason Bechtoldt.

20.11 REFERENCES

1. O. Wichterle and D. Lim, "Hydrophilic Gels for Biological Use," *Nature (London)* **185**:117–118 (1960).
2. M. F. Refojo, "The Chemistry of Soft Hydrogel Lens Materials," in M. Ruben, ed., *Soft Contact Lenses*, New York, John Wiley & Sons, 19–39 (1978).
3. M. Refojo, "The Relationship of Linear Expansion to Hydration of Hydrogel Contact Lenses," *Cont. Lens Intraoc. Lens Med. J.* **1**:153–162 (1976).
4. I. Fatt, "Water Flow and Pore Diameter in Extended Wear Gel Materials," *Am. J. Optom. Physiol. Opt.* **55**:294–301 (1978).
5. P. B. Morgan, C. A. Woods, D. Jones, et al., "International Contact Lens Prescribing in 2006," *Contact Lens Spectrum* **22**(1):34–38 (2007).
6. D. F. Sweeney, N. A. Carnt, R. Du Toit, et al., "Silicone Hydrogel Lenses for Continuous Wear," in E. S. Bennett and B. A. Weissman, eds., *Clinical Contact Lens Practice*, Philadelphia, Lippincott Williams & Wilkins, 693–717 (2005).
7. A. Cannella and J. A. Bonafini, "Polymer Chemistry," in E. S. Bennett and B. A. Weissman, eds., *Clinical Contact Lens Practice*, Philadelphia, Lippincott Williams & Wilkins, 233–242 (2005).
8. B. A. Weissman and K. M. Gardner, "Power and Radius Changes Induced in Soft Contact Lens Systems by Flexure," *Am. J. Optom. Physiol. Opt.* **61**:239 (1984).
9. E. S. Bennett, "Silicone/Acrylate Lens Design," *Int. Contact Lens. Clin.* **11**:547 (1984).
10. W. J. Benjamin, "Optical Phenomena of Contact Lenses," in E. S. Bennett and B. A. Weissman, eds., *Clinical Contact Lens Practice*, Philadelphia, Lippincott Williams & Wilkins, 111–163 (2005).
11. R. B. Mandell, "Optics," in R. B. Mandell, ed., *Contact Lens Practice* (4th ed.), Illinois, Springfield, Charles. C. Thomas, 954–980 (1988).
12. M. W. Ford and J. Stone, "Optics and Contact Lens Design," in A. J. Phillips and L. Speedwell, eds., *Contact Lenses* (5th ed.), London, Elsevier, 129–158 (2007).
13. W. A. Douthwaite, "The Contact Lens," in W. A. Douthwaite, ed., *Contact Lens Optics and Lens Design* (3d ed.), London, Elsevier, 27–55 (2006).
14. M. W. Ford, "Computation of the Back Vertex Powers of Hydrophilic Lenses," *Paper read at the interdisciplinary Conference on Contact Lenses, Department of Ophthalmic Optics and Visual Science, City University, London* (1976).
15. B. A. Weissman, "Loss of Power with Flexure of Hydrogel Plus Lenses," *Am. J. Optom. Physiol. Opt.* **63**:166–169 (1986).
16. B. A. Holden, "The Principles and Practice of Correcting Astigmatism with Soft Contact Lenses," *Aust. J. Optom.* **58**:279–299 (1975).

17. E. S. Bennett, P. Blaze, and M. R. Remba, "Correction of Astigmatism," in E. S. Bennett and V. A. Henry, eds., *Clinical Manual of Contact Lenses* (2d ed.), Philadelphia, Lippincott Williams & Wilkins, 351–409 (2000).
18. W. J. Benjamin, "Practical Optics of Contact Lens Prescription, in E. S. Bennett and B. A. Weissman, eds., *Clinical Contact Lens Practice* (2d ed.), Philadelphia, Lippincott Williams & Wilkins, 165–195 (2005).
19. J. L. Pascal, "Cross Cylinder Tests—Meridional Balance Technique," *Opt. J. Rev. Optom.* **87**:31–35 (1950).
20. R. B. Mandell and C.F. Moore, "A Bitoric Guide That is Really Simple," *Contact Lens Spectrum* **3**:83–85 (1988).
21. E. S. Bennett, "Astigmatic Correction," in E. S. Bennett and M. M. Hom, eds., *Manual of Gas Permeable Contact Lenses* (2d ed.), Elsevier, St. Louis, 286–323 (2004).
22. W. A. Douthwaite, "Aspherical Surface," in W. A. Douthwaite, ed., *Contact Lens Optics and Lens Design* (3d ed.), Elsevier, London, 91–125 (2006).
23. W. Feinbloom, "Corneal Contact Lens having Inner Ellipsoidal Surface. US Patent No. 3,227,507, January 4, 1966.
24. G. L. Feldman and E. S. Bennett, "Aspheric Lens Designs," in E. S. Bennett and B. A. Weissman, eds., *Clinical Contact Lens Practice*, Philadelphia, JB Lippincott, 16-1–16-10 (1990).
25. K. Richdale, G. L. Mitchell, and K. Zadnik, "Comparison of Multifocal and Monovision Soft Contact Lens Corrections in Patients with Low-Astigmatic Presbyopia," *Optom. Vis. Sci.* **83**(5):266–273 (2006).
26. T. B. Edrington and J. T. Barr, "Aspheric What?," *Contact Lens Spectrum* **17**(5):50 (2003).
27. E. S. Bennett, "Contact Lens Correction of Presbyopia," *Clin. Exp. Optom.* **91**(3):265–278 (2008).
28. A. S. Rajagopalan, E. S. Bennett, and V. Lakshminarayanan, "Visual Performance of Subjects Wearing Presbyopic Contact Lenses," *Optom. Vis. Sci.* **83**:611–615 (2006).
29. A. S. Rajagopalan, E. S. Bennett, and V. Lakshminarayanan, "Contrast Sensitivity with Presbyopic Contact Lenses," *J. Modern Optics* **54**:7–9:1325–1332 (2007).
30. E. S. Bennett, "Lens Design, Fitting, and Troubleshooting," in E. S. Bennett and R. M. Grohe, eds., *Rigid Gas-Permeable Contact Lenses*, New York, Professional Press, 189–224 (1986).
31. A. Musset and J. Stone, "Contact Lens Design Tables," Butterworths, London, 79–108 (1981).
32. M. Townsley, "New Knowledge of the Corneal Contour," *Contacto* **14**(3):38–43 (1970).
33. W. A. Douthwaite, "Contact Lens Design," in W. A. Douthwaite, ed., *Contact Lens Optics and Lens Design* (3d ed.), Elsevier, London, 165–201 (2006).
34. P. S. Kollbaum and A. Bradley, "Correcting Aberrations with Contact Lenses," *Contact Lens Spectrum* **22**(11):24–32 (2007).
35. A. Bradley, Personal communication, May, 2009.
36. H. Jiang, D. Wang, L. Yang, P. Xie, and J. C. He, "A Comparison of Wavefront Aberrations in Eyes Wearing Different Types of Soft Contact Lenses," *Optom. Vis. Sci.* **83**(10):769–774 (2006).
37. F. Lu, X. Mao, J. Qu, D. Xu, and J. C. He, "Monochromatic Wavefront Aberrations in the Human Eye with Contact Lenses," *Optom. Vis. Sci.* **80**(2):135–141 (2003).
38. X. Hong, N. Himebaugh, and L. N. Thibos, "On-Eye Evaluation of Optical Performance of Rigid and Soft Contact Lenses," *Optom. Vis. Sci.* **78**(12):872–880 (2001).
39. C. E. Joslin, S. M. Wu, T. T. McMahon, and M. Shahidi, "Higher-Order Wavefront Aberrations in Corneal Refractive Therapy," *Optom. Vis. Sci.* **80**(12):805–811 (2003).
40. D. A. Berntsen, J. T. Barr, and G. L. Mitchell, "The Effect of Overnight Contact Lens Corneal Reshaping on Higher-Order Aberrations and Best-Corrected Visual Acuity," *Optom. Vis. Sci.* **82**(6):490–497 (2005).
41. P. S. Kollbaum and A. Bradley, "Correcting Aberrations with Contact Lenses: Part 2," *Contact Lens Spectrum* **22**(12):31–34 (2007).
42. S. Pantanelli and S. MacCrae, "Characterizing the Wave Aberration in Eyes with Keratoconus or Penetrating Keratoplasty Using a High Dynamic Range Wavefront Sensor," *Ophthalmology* **114**(11):2013–2021 (2007).
43. R. Sabesan, M. Jeong, L. Carvalho, I. G. Cox, D. R. Williams, and G. Yoon, "Vision Improvement by Correcting Higher-Order Aberrations with Customized Soft Contact Lenses in Keratoconic Eyes," *Opt. Lett.* **32**(8):1000–1002 (2007).
44. J. D. Marsack, K. E. Parker, K. Pesudovs, W. J. Donnelly, and R. A. Applegate, "Uncorrected Wavefront Error and Visual Performance During RGP Wear in Keratoconus," *Optom. Vis. Sci.* **84**(6):463–470 (2007).

45. P. Kollbaum and A. Bradley, "Aspheric Contact Lenses: Fact and Fiction," *Contact Lens Spectrum* **20**(3):34–38 (2005).
46. M. Alpern, "Accommodation and Convergence with Contact Lenses," *Am. J. Optom.* **26**:379–387 (1949).
47. G. Westheimer, "The Visual World of the New Contact Lens Wearer," *J. Am. Optom. Assoc.* **34**:135–138 (1962).
48. J. Neumueller, "The Effect of the Ametropic Distance Correction Upon the Accommodation and Distance Correction," *Am. J. Optom.* **15**:120–128 (1938).
49. J. S. Hermann and R. Johnson, "The Accommodation Requirement in Myopia," *Arch. Ophthalmol.* **76**:47–51 (1966).
50. W. A. Douthwaite, "Basic Visual Optics," in W. A. Douthwaite, ed., *Contact Lens Optics and Lens Design* (3d ed.), Elsevier, London, 1–26 (2006).
51. R. B. Mandell, "Contact Lens Optics," in R. B. Mandell, ed., *Contact Lens Practice* (4th ed.), Illinois, Springfield, Charles C. Thomas, 954–979 (1988).
52. A. G. Bennett, "Optics of Contact Lenses," (4th ed.), Association of Dispensing Opticians, London, (1966).
53. J. F. Neumiller, "The Optics of Contact Lenses," *Am. J. Optom. Arch. Am. Acad. Optom.* **45**:786–796(1968).
54. S. Duke-Elder and D. Abrams, "Optics," in S. Duke-Elder and D. Abrams, eds., Section I, *Ophthalmic Optics and Refraction*, volume V of Duke-Elder S, ed., *System of Ophthalmology*, C.V. Mosby, Missouri, St. Louis, 25–204 (1970).
55. R. B. Mandell, "Corneal Topography," in R. B. Mandell, ed., *Contact Lens Practice* (4th ed.), Illinois, Springfield, Charles C. Thomas, 107–135 (1988).

INTRAOCULAR LENSES

Jim Schwiegerling

*Department of Ophthalmology
University of Arizona
Tucson, Arizona*

21.1 GLOSSARY

Accommodating intraocular lens. An artificial lens implanted in the eye that changes power and/or position in response to contraction of the ciliary muscle.

Accommodation. The ability to change the power of the eye and focus on objects at different distances.

Aphakia. A condition of the eye without a crystalline lens or intraocular lens.

Aspheric intraocular lens. An artificial implanted lens that has at least one aspheric surface. Typically used to compensate spherical aberration of the cornea.

Capsulorhexis. Removal of the lens capsule. Usually during cataract surgery, the anterior portion of the capsule is removed to allow access to the crystalline lens within the capsule.

Cataract. An opacification that occurs within the crystalline lens that reduces the quality of vision.

Chromophore. A molecule doped into the lens material that absorbs specific wavelength bands such as ultraviolet or blue light.

Dysphotopsia. Stray light effects encountered following implantation of intraocular lenses. Positive dysphotopsia causes streaks or glints of light seen in the peripheral vision. Negative dysphotopsia causes dark bands or shadows to appear in peripheral vision.

Haptic. A structure that aids in alignment and support of intraocular lenses within the eye. Typically two or three arms emanate from the edge of the lens. Plate haptics are rectangular flanges that protrude from the sides of the lens.

Intraocular lens. An artificial lens that is implanted into the eye to modify the eye's optical power.

Lens capsule. An elastic bag that encapsulates the crystalline lens. Most modern cataract procedures leave the capsule in place and only remove the crystalline lens within.

Limbus. The boundary between the cornea and the sclera, the white of the eye.

Multifocal intraocular lens. An artificial lens that incorporates two or more powers into the same lens. The purpose of these lenses is to allow objects at different distances to be in focus simultaneously.

Phacoemulsification. A technique for removing a cataractous crystalline where an ultrasonic probe is used to shatter the lens and then the lens bits are suctioned out of the eye.

Phakic lens. An artificial lens implanted into the eye while leaving the crystalline lens intact. Typically this procedure is used to alleviate refractive error and is not a treatment for cataracts.

Posterior capsule opacification. A potential side effect of the intraocular lens implant. The lens capsule can adhere to the implant and an opacification develops on the posterior side of the intraocular lens. Laser pulses are typically used to open a hole in the opacified capsule to restore vision.

Presbyopia. The state of the eye where accommodation has been completely lost due to stiffening and enlargement of the crystalline lens.

Pseudophakia. The state of having an artificial lens implanted within the eye.

Sclera. The white of the eye.

21.2 INTRODUCTION

The eye forms the optical system of the human visual system. A variety of references provide a good general introduction to the essential components of the eye and their function.^{1,2} Here, a brief overview of the main optical elements and their mechanism of action will be provided. The eye consists of two separated lenses that ideally form an image on the retina, the array of photosensitive cells lining the back surface of the eyeball. The eye's first lens is the cornea, which is the clear membrane on the external portion of the eye. The cornea is a meniscus lens and has a power of about 43 diopters (D). The iris resides several millimeters behind the cornea. The iris is heavily pigmented to block light transmission and serves as the aperture stop of the eye. The diameter of the opening in the iris varies with light level, object proximity, and age. The crystalline lens is the second optical element of the eye. It lies immediately behind the iris and provides the focusing mechanism of the visual system. Together, the cornea and iris form images on the external environment onto the retina. At the retina, the light is converted to a neural signal and transmitted to the brain where the signals are interpreted into our perceived image of the surrounding scene.

The crystalline lens has several notable features that provide for a flexible optical system capable of focusing to a wide range of distances. First, the crystalline lens has a gradient index structure with its refractive index varying in both the axial and the radial directions. The maximum index of refractive occurs in the center or nucleus of the lens. The index then gradually decreases toward the lens periphery as well as the front and back surfaces. In addition, the crystalline lens is also flexible, allowing the lens to change shape and consequently power in response to muscular forces. The crystalline lens resides in an elastic bag called the capsule. In the young lens, this capsule causes a contraction of the periphery of the lens and an increase in the central thickness of the lens. As a result of this effect, the radii of curvature of the front and back surfaces of the lens decrease, leading to an increase in the optical power of the lens. In its maximally flexible state, the capsule can form the crystalline lens such that it gains a power up to a maximum 34 D.

The lens and capsule are suspended behind the iris by a series of ligaments called the zonules. One end of a zonule attaches to the perimeter of the lens/capsule and the other end attaches to the ciliary muscle. The ciliary muscle is a ring-shaped muscle that resides behind the iris. In its relaxed state, the ciliary muscle dilates, increasing tension on the zonules and exerting an outward radial force on the perimeter of the lens and capsule. This force causes the crystalline lens shape to flatten. The center thickness of the lens reduces and the radii of curvature of its front and back surfaces increase. The net effect of the lens flattening is a reduction in its power to about 20 D. Constriction of the ciliary muscle causes a reduction in the zonule tension and the elastic capsule causes the lens to thicken and increase its power. Accommodation is the ability to change the power of the crystalline lens through constriction or relaxation of the ciliary muscle. Accommodation allows the eye to bring objects at a variety of distances into focus on the retina. If the eye is focused at infinity when the ciliary muscle is in its relaxed state, then the fully accommodated crystalline lens allows objects at 70 mm from the eye to be brought into proper focus.

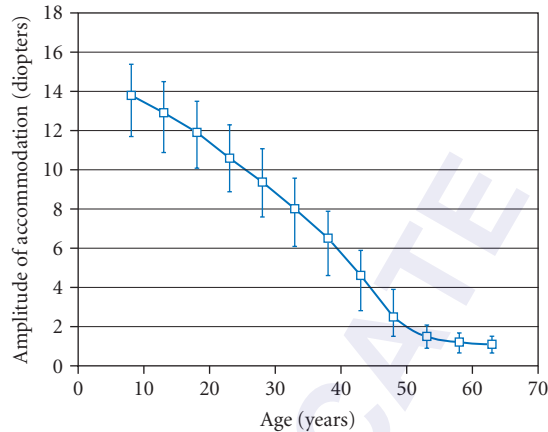


FIGURE 1 Loss of the amplitude of accommodation with age.

The eye is meant to provide a lifetime of vision. However, as with all parts of the body, changes to the eye's performance occur with aging. One effect is the continual decrease in the amplitude of accommodation that occurs almost immediately after birth. Figure 1 shows the accommodative amplitude as a function of age.³ The preceding description of lens function described a lens that could change its power from 20 to 34 D. With age, this range gradually reduces. The changes are usually not noticeable until the fifth decade of life. At this point, the range of accommodation is usually 20 D in the relaxed state to 23 D in the maximally accommodated state. Under these conditions, if the eye is focused at infinity for a relaxed ciliary muscle, then with maximum exertion, objects at 33 cm from the eye can be brought into proper focus. This distance is a typical distance reading material is held from the eye. Objects closer than 33 cm cannot be brought into focus and ocular fatigue can rapidly set in from continual full exertion of the ciliary muscle. Moving into the sixth decade of life, the range of accommodation continues to reduce, leading to the inability to focus on near objects. This lack of focusing ability is called presbyopia. Reading glasses or bifocal spectacles are often used to aid the presbyope in seeing near objects. These lenses provide additional power to the eye that can no longer be provided by the crystalline lens. Presbyopia is caused by the continual growth of the crystalline lens.⁴ The lens consists of a series of concentric layers and new layers are added throughout life. This growth has two effects. First, the lens size and thickness increases with age, leading to a lens that is stiffer and more difficult to deform. Second, the diameter of the lens increases, leading to a reduced tension on the zonules. The combination of these two effects results in the reduction in accommodative range with age and eventually presbyopia.

A second aging effect that occurs in the eye is senile cataract. Cataract is the gradual opacification of the crystalline lens that affects visual function.⁵ An increase in lens scatter initially causes a yellowing of the lens. As these scatter centers further develop large opacities occur within or at the surface of the lens. The progression of cataract initially leads to a reduction in visual acuity and losses in contrast sensitivity due to glare and halos. The endpoint of cataract is blindness, with the opacifications completely blocking light from reaching the retina. It is estimated that 17 million people worldwide are blind due to cataract.⁶ However, the distribution of these blind individuals is skewed toward less developed countries. The prevalence of blindness due to cataracts is 0.014 percent in developed countries and 2 percent in undeveloped countries. The most common cause of lens opacification is chronic long-term exposure to the UV-B radiation in sunlight. Other less common incidences include traumatic, pharmacological, and congenital cataracts. Risk factors such as smoking and poor diet can lead to earlier emergence of cataract. However, since exposure to sunlight occurs nearly daily, there is high likelihood of developing cataracts within a lifetime. Due to this prevalence, treatments for cataracts have long been sought after.

TABLE 1 Prevalence of Cataract in the Population of U.S. Residents

Age	Prevalence (%)
43–54	1.6
55–64	7.2
65–74	19.6
75–85	43.1

Currently, there exist no means for preventing cataract formation. Reducing risk factors such as UV-B exposure through limiting sun exposure and using sunglasses, quitting smoking, and healthy diet can delay the onset of cataracts, but still the opacifications are likely to emerge. Table 1 shows the prevalence of cataract in the United States.⁷

Treatment of cataracts is the main means for restoring lost vision caused by the opacifications. Variations in access to these treatments is a reason for the skew in blindness caused by cataracts in developed and undeveloped countries. In developed countries access to cataract treatment is high and usually provided when the cataractous lens has caused relatively mild loss in visual performance. In fact, cataract surgery is the most widely performed surgery in the United States.⁶ In undeveloped countries, access may be nonexistent and severely limited leading to a high rate of blindness due to complete opacification of the lens.

21.3 CATARACT SURGERY

In restoring vision loss from cataracts, the opacified lens is removed to allow light to enter the eye again. This concept has been known in various evolutionary forms for thousands of years.⁸ However, since the crystalline lens forms about one-third of the overall power of the eye, removing the lens leaves the patient severely farsighted. This condition is known as aphakia. Aphakia requires patients to wear high-power spectacle lenses to replace the roughly 20 D of power lost from the crystalline lens. These high-powered lenses are problematic for a variety of reasons. These effects include a 20 to 35 percent magnification of the image and peripheral distortion, since spectacles with spherical surface cannot be corrected for oblique astigmatism at such high positive powers and a ring scotoma due to a gap in the field of view between light entering just inside and just outside the edge of the spectacle lens.^{9,10} The advent of intraocular lenses (IOLs) in 1949 served to “complete” cataract surgery.^{11,12} IOLs are artificial lenses that are implanted into the eye following the removal of a cataractous lens and are designed to provide the needed optical power to allow normal image formation to occur within the eye. Placement of the IOL within the eye allows aphakic spectacles and their problems can be avoided.

Modern cataract treatments include extracapsular cataract extraction (ECCE) and more recently a variation of this technique known as phacoemulsification. ECCE was first performed over 250 years ago, and has been widely used to remove the opacified lens.¹³ In this procedure, a large incision is made in the limbus, which is the boundary of the cornea. A capsulorhexis, which is the removal of the anterior portion of the capsule, is performed and then the crystalline lens is removed, leaving the remaining portion of the capsule intact within the eye. The corneal incision site is then sutured to seal the eye. Phacoemulsification was developed in 1967 by Charles Kelman.¹⁴ In this technique, a small incision is made at the limbus. Tools inserted through this small opening are used to make a capsulorhexis. With the crystalline lens now exposed, a hollow ultrasonic probe is slipped through the opening. The vibrations of the probe tip cause the crystalline lens to shatter and suction is used to aspirate the lens fragments through the hollow center of the probe. The entire lens is removed in this fashion leaving an empty capsule. The incision in phacoemulsification is sufficiently small so that no sutures are needed to seal the wound. Phacoemulsification is similar to modern surgical techniques of orthoscopic and laposcopic surgery. Incisions are kept small to promote rapid healing and limit the risk of infection.

21.4 INTRAOCULAR LENS DESIGN

Sir Harold Ridley, a British surgeon specializing in ophthalmology, was dissatisfied with the existing form of cataract surgery since it left patients aphakic and dependent on poorly performing spectacle lenses for visual function.¹² Ridley gained a key insight into “fixing” cataract surgery while treating Royal Air Force pilots from World War II. The canopies of British fighter planes were constructed with a newly developed polymer, polymethyl methacrylate (PMMA), also known as Perspex, plexiglass, or acrylic. In some instances, when the gunfire shattered the canopies, shards of the PMMA would become embedded into the cornea of the pilot. Ridley, in treating these aviators, noted that the material caused little or no inflammatory response of the surrounding ocular tissue. The typical response of the body to a foreign body is to encapsulate and reject the material. However, the PMMA did not elicit this response and sparked the concept of forming a lens from a “biocompatible” material that could be implanted following extraction of the cataractous lens. In 1949, Ridley implanted the first IOL into a 45-year-old woman. Following surgery, the patient was 14 D nearsighted (as opposed to 20 D farsighted in aphakic patients). While the power of the implanted IOL was clearly wrong, Ridley successfully demonstrated that an artificial lens could be implanted into the eye that evokes a marked power change. The error in the power of the IOL stemmed from modeling the implant after the natural crystalline lens. The approximate radii of curvature of the crystalline lens were used for the implant, but the higher-refractive index of the PMMA was not taken into account. Ridley’s second surgery resulted in a similar degree of nearsightedness, but by the third surgery, the IOL power had been suitably refined, leaving the patient mildly nearsighted.

Ridley’s IOL invention is revolutionary, as it restored normal vision following cataract surgery. However, the medical establishment did not immediately accept the value of this invention. Instead, many leading figures in ophthalmology at the time thought the surgery to be reckless and unwarranted. Ridley met with disdain, criticism, and the threat of malpractice. The result of this animosity was a long delay in the evolution of the implant technology as well as a limited use of the lenses by the surgical community. The shift to acceptance of the technology did not begin to occur until the 1970s. In 1981, the Food and Drug Administration finally approved IOLs for implantation following cataract surgery in the United States. Today, phacoemulsification followed by IOL implantation is the most widely performed and successful surgery in the United States. Virtually no aphakic subjects exist anymore and the term “cataract surgery” now implicitly implies implantation of the IOL.

Intraocular Lens Power

As Ridley’s first two surgeries illustrate, calculation of the IOL power is important to minimize the error in refractive power of the eye following surgery. Ridley ventured into unknown territory in designing the initial IOLs. Accurate determination of the IOL power remains the key component to a successful surgery. Cataract surgeons typically target slight nearsightedness following implantation of the IOL. Conventional IOLs leave the recipient presbyopic because the lens power and position are fixed within the eye. Unlike the crystalline lens, the conventional IOL does not change shape or position and consequently accommodation is lost. Extensive research into accommodating IOLs is currently being pursued and a description of emerging technologies is given further in the chapter. With the fixed-focus IOL, surgeons need to determine the proper power prior to surgery, which can be challenging for several reasons. First, the eye must be tested with the existing crystalline lens in place. The power of the crystalline lens must be deduced from the overall power of the eye and the shape of the cornea. The crystalline lens power must be subtracted from the overall ocular power since the crystalline lens is removed by the surgery. Next, the power of the IOL must be estimated based on its anticipated position within the eye following implantation. Finally, the power of the IOL must also be adjusted for any inherent refractive error in the eye. In other words, the IOL power and the cornea must work in conjunction to give a sharp image on the retina. To further complicate this calculation, the patient is receiving the implant because of cataract, so subjective measurement of refractive error may be confounded by the opacities in the lens.

The IOL power calculation formulas fall into two categories. Some formulas are based on multiple regression analysis of the variables associated with the eye and implant, while other formulas are based on a theoretical prediction of the IOL power. The SRK formula was originally the most widely used regression-type formula to predict IOL power. The SRK formula is given by^{15,16}

$$\phi_{\text{IOL}} = A - 0.9 K - 2.5 L \quad (1)$$

where ϕ_{IOL} is the power of the IOL in diopters, A is the A-constant of the lens, K is the corneal keratometry, and L is the axial length of the eye. The A-constant is a value provided by the IOL manufacturer that is most closely associated with the position of the IOL within the eye. Modern IOLs are about 1 mm thick and about one-fourth the thickness of the natural crystalline lens. Consequently, there is a variation in the positioning of the IOL behind the iris that is dependent on the shape of the implant and surgical technique. The A-constant is empirically derived to account for these factors. Surgeons often customize the A-constants of lenses routinely implanted based on their specific outcomes.

Keratometry is a measure of corneal power. Keratometry is calculated measuring the radius of curvature R_a in mm of the anterior corneal surface. A power is then calculated as

$$K = 1000 \frac{(n_k - 1)}{R_a} \quad (2)$$

where n_k is the keratometric index of refraction and K is in units of diopters. The keratometric index of refraction is an effective corneal index that attempts to incorporate the power associated with the posterior surface of cornea into the keratometry measurement. Values of 1.3315 to 1.3375 are routinely used for n_k in clinical devices used for measuring keratometry. Measurement of the posterior corneal radius of curvature has been difficult until recently, so historically keratometry measurements have been made based on an estimate of the contribution of this surface to the overall power of the cornea. Note that n_k is lower than the true refractive index of the cornea, 1.376, to account for the negative power induced by the posterior cornea. As an example, the radii of curvature of the anterior and posterior surfaces of the cornea are approximately 7.8 and 6.5 mm, respectively.² The refractive index of the cornea is 1.376 and of the aqueous humor on the back side of the cornea is 1.337. A typical thickness of the cornea is 0.55 mm. With these assumptions, the keratometry of the cornea, as given by Eq. (2), is

$$K = 1000 \frac{(1.3315 - 1)}{7.8} = 42.50 \text{ D} \quad (3)$$

The true power of the cornea is given by

$$\phi_{\text{cornea}} = 1000 \left[\frac{(1.376 - 1)}{7.8} + \frac{(1.337 - 1.376)}{6.5} - \left(\frac{0.55}{1.376} \right) \left(\frac{(1.376 - 1)}{7.8} \right) \left(\frac{(1.337 - 1.376)}{6.5} \right) \right] = 42.32 \text{ D} \quad (4)$$

Clearly, there are some small differences in the predicted keratometric and the corneal power. This discrepancy is one potential source for error in the calculation of the implant power. As can be seen from Eq. (1), a 0.9 D error in keratometry leads to a 1 D error in the calculated IOL power.

Traditionally, the axial length, L , of the eye is measured with A-scan ultrasonography. An ultrasonic probe, typically with a frequency of 10 MHz, is placed lightly in contact with the cornea to avoid deforming the surface. Echoes from the crystalline lens surfaces and the retina are then measured. The time of flight of these echoes is directly related to the speed of sound in the various materials of the eye. Table 2 summarizes typical values of these speeds used by commercial A-scan units. The longitudinal resolution of a 10-MHz ultrasound is 200 μm .^{17,18} Consequently, Eq. (1) suggests this level of accuracy would lead to a limit of the accuracy of the IOL power of ± 0.25 D.

TABLE 2 Speed of Sound in Various Ocular and IOL Materials

Material	Speed of Sound
Cornea	1,640 m/s
Aqueous	1,532 m/s
Lens—Normal	1,640 m/s
Lens—Cataractous	1,629 m/s
Vitreous	1,532 m/s
PMMA IOL	2,760 m/s
Silicone IOL	1,000 m/s

Partial coherence interferometry has recently become a popular technique for measuring axial length.¹⁹ In this technique, a Michelson interferometer is used to measure distances within the eye. A short-coherence-length narrow-band infrared source is shone into the eye. The return beam is interfered with light from the test arm. Since the beam has low-coherence, fringes only appear when the path lengths in the eye and the test arm are nearly identical. By changing the length of the reference arm and observing when fringes appear, the length of the eye can be accurately measured. The axial resolution of this technique is on the order of 12 μm , leading to a resolution of the IOL power of 0.015 D.

The first regression-type formulas are accurate for eyes that fall within “normal” parameters for axial length and corneal power. However, as the more extreme sizes and shapes are approached, the accuracy of formulas degrades markedly. The regression-type formulas were further refined to better handle these extreme cases.²⁰ While regression analysis is still used for these formulas, a piecewise linear approximation is used to better approximate the full range of variation of human eyes. The SRK II is an example of an evolved formula. The IOL power is defined as

$$\phi_{\text{IOL}} = A_1 - 0.9 K - 2.5 L \quad (5)$$

where

$$\begin{aligned} A_1 &= A + 3 & \text{for } L < 20.0 \text{ mm} \\ A_1 &= A + 2 & \text{for } 20.0 \leq L < 21.0 \\ A_1 &= A + 1 & \text{for } 21.0 \leq L < 22.0 \\ A_1 &= A & \text{for } 22.0 \leq L < 24.5 \\ A_1 &= A - 0.5 & \text{for } L \geq 24.5 \end{aligned}$$

Figure 2 compares the SRK and SRK II formulas.

An alternative to regression-based formulas for IOL power prediction are theoretical formulas.²¹ Theoretical formulas use geometrical optics to make predictions about the implant power, given knowledge of the axial length, corneal power, and position of the IOL following implantation. Figure 3 shows a simplified model of the eye. The cornea is modeled as a single surface of power K . In the aphakic eye shown in Fig. 3a, the cornea forms an image of a distant object at a distance $1000 n_{\text{aq}}/K$ behind the corneal vertex for K in diopters. When the IOL is implanted into the eye, it sits at a location called the effective lens position (ELP). The ELP is deeper into the eye than the iris plane for posterior chamber IOLs. Figure 3b shows the IOL implanted into the aphakic of Fig. 3a. The image formed by the cornea becomes a virtual object for the IOL located at a distance

$$S = 1000 n_{\text{aq}} / K - \text{ELP} \quad (6)$$

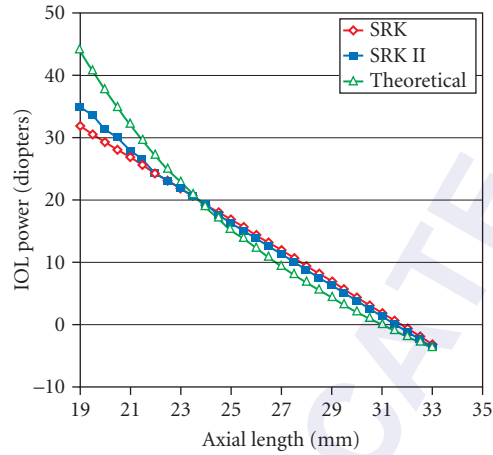


FIGURE 2 A comparison of the predicted IOL power for the SRK and SRK II formulas assuming $A = 118$ and $K = 43$ D. Also shown is a theoretical prediction of the IOL power assuming $n_{\text{aq}} = 1.337$, $K = 43$ D, and $\text{ELP} = 5$ mm.

from the plane of the IOL. The IOL, in turn, needs to image this virtual object onto the retina. If the distance from the IOL to the retina is given by S' , then

$$\frac{n_{\text{aq}}}{S'} - \frac{n_{\text{aq}}}{S} = \phi_{\text{IOL}} = n_{\text{aq}} \left[\frac{1}{L - \text{ELP}} - \frac{1}{1000n_{\text{aq}}/K - \text{ELP}} \right] \quad (7)$$

which can be rewritten as

$$\phi_{\text{IOL}} = \frac{n_{\text{aq}}}{K} \left[\frac{1000n_{\text{aq}} - KL}{(L - \text{ELP})(1000n_{\text{aq}}/K - \text{ELP})} \right] \quad (8)$$

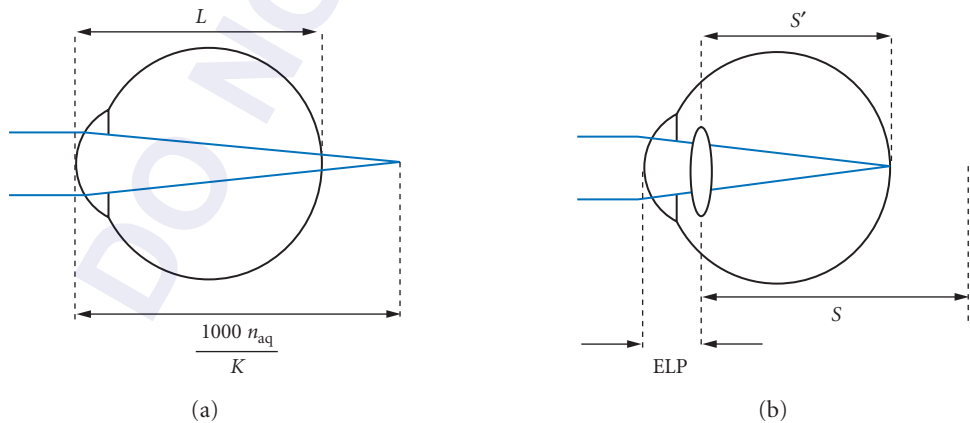


FIGURE 3 (a) In the aphakic eye, the image is formed at a distance $1000n_{\text{aq}}/K$, where n_{aq} is the refractive index of the aqueous and K is the power of the cornea. (b) When the IOL is inserted at the effective lens position (ELP), the aphakic image is reformed onto the retina.

Equation (8) is the theoretical IOL power based on thin lens analysis of the pseudophakic eye. Figure 2 also shows the predicted IOL power based on Eq. (8).

A variety of theoretical models exist.^{22–25} They are all based on the preceding Eq. (8). Subtle differences exist, however, in the manner each of the variables is calculated. For example, the axial length provided by ultrasonography is slightly short because the echo occurs at the inner limiting membrane of the retina. Light absorption occurs at the photoreceptor layer that is approximately 20 μm deeper into the eye. The ELP is also affected by the anterior chamber depth, which is the distance from the cornea vertex to the iris plane. Eyes with flatter corneas or smaller diameter corneas can have shallower anterior chamber depths, causing a reduction in the ELP. Holladay has worked extensively to determine the location within the eye of available IOLs and to relate the A-constant to the ELP as well as what Holladay calls the surgeon factor (SF).^{26–31} The SF is the distance from iris plane to the front principal plane of the IOL. The theoretical models in most cases provide excellent prediction of the IOL power. The typical error in power is less than 1 D in over 80 percent of patients.²⁰ Surgeons tend to err on the side of nearsightedness to allow the patient to have some plane in their field of vision conjugate to the retina. A farsighted patient would only have a virtual object in focus on their retina and consequently they tend to be less happy since additional refractive correction is required for all object distances. The current generation of IOL power calculations, however, has difficulty in predicting accurate IOL powers in patients who have undergone corneal refractive surgeries such as radial keratotomy (RK) and laser in situ keratomileusis (LASIK). These surgical procedures modify the shape of the cornea to correct for refractive error. However, this shape change can have marked effect on the measurement of the keratometry. The relationship between the shape of the front and back surfaces of the cornea and in the case of LASIK, the corneal thickness has been artificially modified by the refractive procedure. As a result, the keratometric index of refraction is no longer valid and gross errors in IOL prediction can occur as a result. A variety of techniques have been suggested for handling this subset of cataract patients.^{32–42} Only a limited number of patients currently fall into this category. However, as the refractive surgery generation ages, an increasing percentage of cataract patients will require specialized techniques for predicting IOL power and much research is still needed in this area to achieve satisfactory results.

Intraocular Lens Aberrations

Since IOLs are singlets, only limited degrees of freedom exist to perform aberration correction. The human cornea typically has positive spherical aberration.⁴³ The natural crystalline lenses typically compensates for the cornea with inherent negative spherical aberration leaving the whole eye with low levels of positive spherical aberration.⁴⁴ The crystalline lens creates negative spherical aberration through its gradient index structure and aspheric anterior and posterior surfaces. Traditionally, IOLs have been made with spherical surfaces due to the ease of manufacturability. Smith and Lu analyzed the aberrations induced by IOLs using the thin lens Seidel aberration formulas described by Welford.^{45,46} Spherical aberration is given by

$$\frac{y^4 \phi_{\text{IOL}}^3}{4n_{\text{aq}}^2} \left[\frac{n_{\text{aq}}^2 (2n_{\text{aq}} + n_{\text{IOL}})}{n_{\text{IOL}} (n_{\text{aq}} - n_{\text{IOL}})^2} X^2 + \frac{4n_{\text{aq}} (n_{\text{aq}} + n_{\text{IOL}})}{n_{\text{IOL}} (n_{\text{aq}} - n_{\text{IOL}})} XY + \frac{2n_{\text{aq}} + 3n_{\text{IOL}}}{n_{\text{IOL}}} Y^2 + \frac{n_{\text{IOL}}^2}{(n_{\text{aq}} - n_{\text{IOL}})^2} \right] \quad (9)$$

where y is the height of the marginal ray at the IOL, ϕ_{IOL} is the power of the implant, $n_{\text{aq}} = 1.337$ is the refractive index of the aqueous, and n_{IOL} is the refractive index of the IOL. Typical values for n_{IOL} range from 1.41 to 1.55 depending on the material. The shape factor X is defined as

$$X = \frac{C + C'}{C - C'} \quad (10)$$

where C is the curvature of the anterior surface of the IOL and C' is the curvature of the posterior surface. The conjugate factor Y is defined as

$$Y = \frac{S + S'}{S - S'} \quad (11)$$

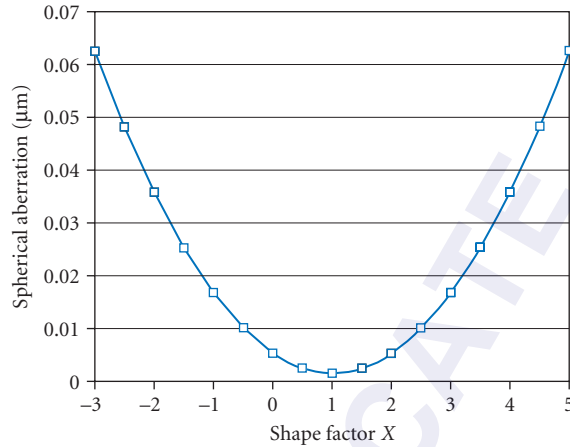


FIGURE 4 The spherical aberration as a function of shape factor for an IOL.

where S is the distance from the IOL to the virtual object formed by the cornea and S' is the distance from the IOL to the retina. In the eye, a 43-diopter cornea forms an image of a distant object about 31 mm behind the corneal vertex. The IOL typically resides in a plane 5 mm into the eye. Under these conditions, $S = 31 - 5 = 26$ mm. If the eye has an axial length of $L = 24$ mm, then $S' = 19$ mm. In this example, the conjugate factor $Y = 6.43$. Figure 4 shows the spherical aberration as a function of lens shape factor for a 20 D-PMMA lens with the preceding eye parameters and $\gamma = 2$ mm. Two key features of this plot are apparent. First, the values of spherical aberration are strictly positive. This is true for any positive power lens with spherical surfaces. Since the corneal spherical aberration is typically positive as well, the spherical-surfaced IOL cannot compensate for this aberration, as the natural crystalline lens tends to do. Instead, the total ocular spherical aberration of the pseudophakic eye can only be minimized through the appropriate choice of IOL shape factor X . In the example in Fig. 4, a shape factor $X = +1.0$ minimizes the spherical aberration. This shape corresponds to a plano-convex lens with the flat side toward the retina. The shape factor $X = +1$ is ideal for a perfectly aligned cornea and IOL. However, natural variations in the tilt and decentration of the IOL and the orientation of the visual axis within the eye make this perfect alignment difficult. Typical values for the tilt and decentration of the IOL are less than 2.6° and 0.4 mm, respectively.⁴⁷ Atchison has shown that a shape factor of $X = +0.5$, double-convex lens with a less curved radius toward the retina, is less sensitive to these errors and modern spherical IOLs tend to target this shape factor.^{48,49}

Aspheric Lenses

Recently, IOLs have shifted toward incorporating an aspheric surface for aberration control.⁵⁰ Whereas, conventional spherical-surfaced IOLs seek to minimize the aberrations induced by the implant, aspheric IOLs can be used to compensate spherical aberration induced by the cornea. The corneal spherical aberration stems mainly from its aspheric shape. The cornea over the entrance pupil of the eye is well represented by a conic section of radius R and conic constant Q . The sag of the cornea z is then described as

$$z = \frac{1}{Q+1} \left[R - \sqrt{R^2 - (Q+1)r^2} \right] \quad (12)$$

TABLE 3 Radius R and Conic Constant Q of the Human Cornea

References	R (mm)	Q	Technique
Keily ⁵¹	7.72	-0.26	Photokeratoscopy
Guillon ⁵²	7.85	-0.15	Photokeratoscopy
Dubbelman ⁵³	7.87	-0.18	Scheimpflug
Dubbelman ⁵⁴	7.87 (male) 7.72 (female)	-0.13	Scheimpflug

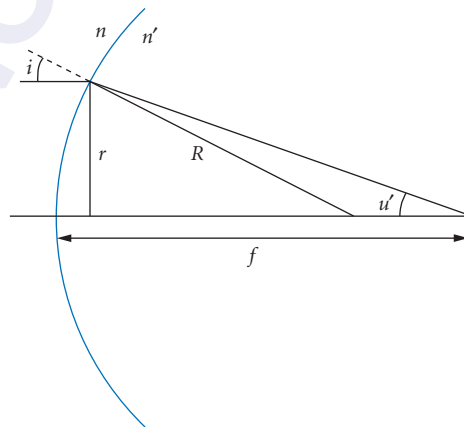
where r is the radial distance from the optical axis. Table 3 summarizes the values of R and Q of the anterior cornea found in several studies of corneal topography.⁵¹⁻⁵⁴ There is also a tendency for the cornea to become more spherical with age^{54,55} leading to an increase in corneal spherical aberration.⁵⁶

The spherical aberration of a surface can be found from paraxial raytracing values. The spherical aberration wavefront coefficient W_{040} of a spherical surface is given by⁴⁶

$$W_{040} = -\frac{1}{8}r(ni)^2 \left(\frac{u'}{n'} - \frac{u}{n} \right) \quad (13)$$

where n is the refractive index, i is the paraxial angle of incidence of the marginal ray on the surface, and u the angle the marginal ray makes with respect to the optical axis. The primed terms denote variables following the surface, whereas the unprimed terms denote variable prior to the surface. Figure 5 shows the layout for calculating the spherical aberration of the anterior cornea for a distant object. In this case, $u = 0$, $n = 1.0$, u' is the marginal ray angle following the surface, $n' = 1.376$ is the refractive index of the cornea, and i is the paraxial angle of incidence of the marginal ray. From the figure it is evident that

$$i = \frac{r}{R} \quad (14)$$


FIGURE 5 Variables in the calculation of the spherical aberration of the anterior cornea.

within the paraxial small angle assumption. Furthermore, the focal length f of the surface is given by

$$f = \frac{n'}{n'-1}R \quad (15)$$

which leads to a marginal ray angle of

$$u' = -\frac{(n'-1)r}{n'R} \quad (16)$$

Hopkins⁵⁷ showed that a correction factor can be added to the spherical surface aberration term to account for the asphericity of a conic surface. This correction factor is given by

$$\Delta W_{040} = \frac{Q(n'-1)r^4}{8R^3} \quad (17)$$

Combining Eq. (13) to Eq. (17), the total spherical aberration of the anterior cornea is given by

$$W_{040} + \Delta W_{040} = \frac{(n'-1)(1+n'^2Q)r^4}{8n'^2R^3} \quad (18)$$

If the preceding technique is repeated for the posterior corneal surface, then its spherical aberration typically is 2 orders of magnitude smaller than the spherical aberration of the anterior surface. Consequently, the contribution of the posterior cornea to the spherical aberration of the aphakic eye can be ignored. While Eq. (18) represents the wavefront coefficient of the corneal spherical aberration, the literature has somewhat adopted a different format. Typically, corneal spherical aberration has been given in terms of the Zernike polynomial expansion coefficient c_{40} for a normalization radius $r = 3$ mm.⁵⁸ In this case, the corneal spherical aberration is given by

$$c_{40} = W_{040}(r=3) + \Delta W_{040}(r=3) = \frac{1}{6\sqrt{5}} \frac{81(n'-1)(1+n'^2Q)}{8n'^2R^3} = \frac{0.15}{R^3} + \frac{0.284Q}{R^3} \quad (19)$$

The value of c_{40} is linear with respect to the corneal conic constant, as shown in Fig. 6. Table 4 summarizes the values of c_{40} from various studies in the literature.⁵⁹⁻⁶² Since the aberrations of the various elements of the eye add, the implanted aspheric IOL should compensate for the corneal spherical

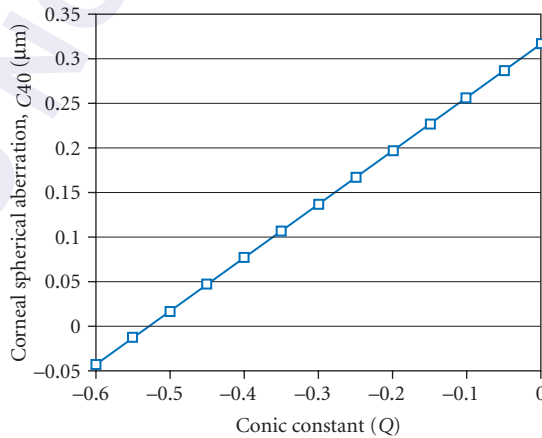


FIGURE 6 Corneal spherical aberration as a function of conic constant Q for a corneal radius $R = 7.8$ mm.

TABLE 4 Corneal Spherical Aberration c_{40}

References	c_{40}
Holladay ⁵⁰	$0.270 \pm 0.200 \mu\text{m}$
Wang ⁵⁹	$0.280 \pm 0.086 \mu\text{m}$
Belluci ⁶⁰	$0.276 \pm 0.036 \mu\text{m}$
Guirao ⁶¹	$0.320 \pm 0.120 \mu\text{m}$

aberration given by Eq. (19). Examples of aspheric IOLs include the Tecnis Z9000 (Advanced Medical Optics, Irvine, CA) designed to correct $0.27 \mu\text{m}$ of corneal spherical aberration, the AcrySof SN60WF (Alcon Laboratories, Fort Worth, TX) designed to correct for $0.20 \mu\text{m}$ of corneal spherical aberration, and the SofPort Advanced Optics IOL (Bausch & Lomb, Rochester, NY) designed such that the implant has zero spherical aberration.

Toric Lenses

Astigmatism is another important consideration following cataract surgery. IOLs have been traditionally rotationally symmetric, again due to manufacturing ease. However, astigmatism may be present in the aphakic eye. This astigmatism is due to a toric shape in the cornea and, since the optical system is not rotationally symmetric, the astigmatism appears along the visual axis. Recently, toric IOLs have become available, but their use has been limited.⁶³ It is more common to reduce corneal astigmatism at the time of surgery with procedures known as corneal relaxing incisions or limbal relaxing incisions. In these procedures, arcuate incisions are made in the peripheral cornea or limbus, respectively. These incisions have the effect of making the cornea more rotationally symmetric after the incisions heal. Consequently, corneal astigmatism is reduced and a conventional rotationally symmetric IOL can be implanted. Corneal or limbal relaxing incisions are typically performed at the time of cataract surgery. An alternative to relaxing incisions is to perform a refractive surgery procedure such as LASIK to reduce corneal astigmatism. This procedure uses an excimer laser to ablate corneal tissue, leaving the postoperative cornea more rotationally symmetric. Toric IOLs are available and similar power calculations must be taken into account.⁶⁴⁻⁶⁷ FDA-approved toric IOLs are available in the United States from Staar Surgical (Monrovia, CA) and Alcon Laboratories (Fort Worth, TX).

Testing IOLs

Clinical testing of implanted IOLs provides insight into the function of the lenses, but can also be difficult due to confounding factors such as cornea clarity, retina function, variations in surgical positioning of the implant, and neural processing. Furthermore, large numbers of subjects and long follow-up times make these types of studies expensive and slow. Several methods, including the Ronchi test and speckle interferometry, have been used to test isolated IOLs.^{68,69} Alternatively, model eyes can be used to test the optical performance of IOLs in a configuration similar to their implanted state.^{62,70-78} Model eyes allow for rapid evaluation of the lenses that avoids the confounding factors described above. In testing the optical performance of IOLs with a model eye, the geometry of how the lens is used within the eye is taken into account. Aberrations, which affect the IOL optical performance, are dependent on the vergence of light entering the implant. In general, the model eye consists of a lens that represents an artificial cornea and a wet cell, which is a planar-faced vial containing saline into which the IOL is mounted. The wet cell is placed behind the artificial cornea such that the corneal lens modifies incident plane waves and creates a converging beam onto the IOL. The vergence striking the IOL is meant to be representative of the vergence seen by an implanted IOL. An artificial pupil can also be introduced to simulate performance for different pupil sizes. The performance of the lens can then be evaluated at the image plane of the eye model. Several different types of artificial corneas have been used in eye models. The international standard ISO11979-2 recommends using a model cornea that is virtually free of aberrations in conjunction with the light

source used.⁷⁹ The standard provides an example model cornea made from a Melles-Griot LAO 034 lens. This lens is a commercially available cemented achromat and consequently is well corrected for both chromatic and spherical aberration. With this model cornea, the spherical aberration of the whole model eye is created solely from the IOL. This model eye specification predated the advent of aspheric IOLs. This type of model eye is therefore only suitable for earlier spherical surface IOLs to evaluate how well these implants minimize their inherent spherical aberration. However, this type of model cornea is not ideal for modern aspheric IOL designs. These newer designs are made to work in conjunction with the cornea and compensate for the corneal spherical aberration. Model corneas with representative levels of corneal spherical and chromatic dispersion have been proposed and used to test these advanced IOL designs.⁶²

Measurement of the optical transfer function (OTF) or its modulus the modulation transfer function (MTF) are routine tests for measuring the optical quality of IOLs.^{70–78} The MTF of an optical system describes the amount of contrast that is passed through the system. If a high contrast sinusoidal target is imaged by an optical system, then the contrast of the resultant image is reduced. In general, the contrast tends to decrease more severely with higher spatial frequency (i.e., finer spacing between the bars of the sinusoidal target). The MTF of a model eye is conveniently calculated by measuring the line spread function (LSF) of the model eye. The LSF, as its name implies, is simply the image of a narrow slit formed by the model eye. The image $i(x, y)$ of the vertically oriented slit on the retina is given by

$$i(x, y) = \text{rect}\left(\frac{x}{d}\right) * \text{PSF}(x, y) \quad (20)$$

where d is the width of the unaberrated slit on the retina and $\text{PSF}(x, y)$ is the point spread function of the eye. The Fourier transform of Eq. (12) gives

$$\mathfrak{F}\{i(x, y)\} = d \text{sinc}(d\xi) \delta(\eta) \text{OTF}(\xi, \eta) \Rightarrow \text{OTF}(\xi, 0) = \frac{\mathfrak{F}\{i(x, y)\}}{d \text{sinc}(d\xi)} \quad (21)$$

where $\mathfrak{F}\{ \}$ denotes the Fourier transform and $\text{sinc}(\xi) = \sin(\pi\xi)/(\pi\xi)$. Equation (21) says that the optical transfer function (OTF) is the ratio of the image spectrum and a sinc function. Note that the denominator of Eq. (13) goes to 0 for $\xi = 1/d$. Under this condition, the OTF approaches infinity. Consequently, the size of d must be made sufficiently small to ensure that the spatial frequencies ξ fall within a desirable range.

Multifocal Lenses

A variety of strategies have been employed to help alleviate presbyopia resulting from cataract surgery. The ideal situation would be to implant an IOL that had the capability to either change its power as the young crystalline does or change its position in response to contraction of the ciliary muscle. The latter case provides accommodation by changing the overall power of the eye through a shift in the separation between the cornea and IOL. These strategies and similar variations are being aggressively pursued as “accommodating” IOLs discussed in more detail in the following section. Clinical results demonstrating the ability of these IOLs to restore some degree of accommodation have been somewhat mixed, but improvement is likely as the accommodative mechanism becomes better understood and the lens designs incorporate these actions.

Multifocal IOLs represent a successful bridging technology between conventional IOLs and a true accommodating IOL. Multifocal optics have two or more distinct powers within their aperture. These lenses take advantage of simultaneous vision. In other words, both in-focus and out-of-focus images are simultaneously presented to the retina. The role of the brain then is to filter out the blurred component and interpret the sharp component providing suitable vision for two distinct distances. For example, suppose the required monofocal (single power) IOL power to correct distance vision in a patient undergoing refractive surgery is 20 D. If a multifocal lens containing dual powers of 20 and 24 D is implanted instead, the patient will have simultaneous vision following surgery. (Note that a 4 D add power in the

IOL plane is approximately equivalent to a 3 D add in the spectacle plane.) If the patient views a distant object, the distance portion of the IOL will form a sharp image on the retina, while the near portion of the lens will create a blurred image of the distant scene on the retina. Similarly, if the patient now views the text of a book, the near portion of the IOL will form a sharp image of the text on the retina, while the distance portion of the IOL will now create a blurred image of the page. In both cases, if the deleterious effects of the blurred portion of the retinal image are sufficiently low, then successful interpretation of the in-focus portion can be made. As a result of simultaneous vision, some contrast is lost in the in-focus image. Multifocal optics therefore represent a trade-off in visual performance. Monofocal lenses provide sharp, high contrast images for distance and horribly blurred low-contrast images for near objects. Multifocal lenses provide reasonably high contrast images for both distance and near vision, but there is some loss in contrast compared to monofocal distance vision. When designing multifocal lenses, control of the out-of-focus portion of the retinal image and understanding the conditions under which the lenses are used are important to optimizing the design. Proper understanding of these issues can provide high-quality multifocal images and lead to happy, spectacle-free patients.

Multifocal IOLs require two or more powers to be simultaneously present in the lens.⁸⁰ Zonal refractive lenses introduce multiple powers into the lens by having distinct regions or zones that refract light differently. The width or area of these zones is large compared to the wavelength of light. As a result of these large zones, the bending of the light rays is determined strictly by refraction. The local surface curvature of the lens determines the local lens power, and different powers can be incorporated simply by having regions with different curvatures. A simple example of a zonal refractive lens would be to bore a 2 mm hole in the center of a 20 D lens and then fill the region with a 2 mm region cookie cut from a 24 D lens. These resulting lens would act as a 20 D lens with a 4 D add in the center, providing multifocal optics. One clear problem arises in this design. The junction between the two lens regions will be abrupt and discontinuous leading the stray light effects. Zonal refractive lenses blend the transition between regions. This blend can be rapid, or the transition can be made slowly to introduce regions of intermediate power. Clearly, this concept can be extended to more than two discrete regions, so that multiple annular regions of alternating power can be created. Furthermore, the regions do not even need to be concentric circles. Instead, tear shaped, wedged, and even swirled regions have been demonstrated. The ReZoom IOL and its predecessor the Array IOL (Advanced Medical Optics, Irvine, CA) are examples of zonal refractive lenses. Figure 7 illustrates the power profile of such a lens. As can be expected with simultaneous vision, there exists a degradation in the optical quality of the retinal image with zonal refractive lenses. For example, starburst and halo patterns can result from these types of lenses.^{81,82} However, these downsides provide roughly two-thirds of patients with near visual acuity of 20/40 or better, and the differences between zonal refractive and conventional monofocal lenses on complex tasks such as driving is small.^{82,83} These lenses provide improved near vision at the expense of the quality of overall vision.

The second optical phenomenon that is exploited to create multifocal optics is diffraction.^{84–86} Multifocal optics with diffractive structures are often misunderstood because they move away from the geometrical picture of light rays bending at the surface of the lens. Instead, these lenses take

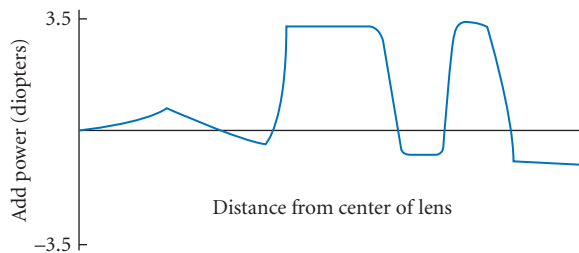


FIGURE 7 Power profile of the Array lens. The lens is formed from multiple concentric regions that oscillate between the base lens power and the higher add power.

advantage of the wave nature of light. In diffractive IOLs, concentric annular zones are created on the face of the lens. The j th zone occurs at a radius

$$r_j = \sqrt{2j\lambda_o F} \quad (22)$$

where λ_o is the design wavelength and F is the focal length of the add power. At the junction of each zone, an abrupt step appears. Both the height of the step and the dimensions of the zones control the degree of multifocality of the lens. The dimensions of the zones are related to the desired add power and in general, the spacing between zones gets progressively smaller from the lens center to its edge. The height of the step at the boundary of each zone determines how much light is put into the add portion. In typical multifocal diffractive lenses, this step height is chosen so that the peaks of one diffractive zone line up with the troughs of the next larger diffractive zone immediately following the lens. As these waves propagate to the retina, the waves from the various diffractive zones mix and there are two distinct regions of constructive interference that correspond to the two main foci of the multifocal lens. The optical phase profile $\phi(r)$ of a diffractive lens is given by⁸⁷

$$\phi(r) = 2\pi\alpha \left(j - \frac{r^2}{2\lambda_o F} \right) \quad r_j \leq r < r_{j+1} \quad (23)$$

where α represents a fraction of the 2π phase delay. This phase pattern is superimposed onto one of the curved surfaces of the IOL. Note that if the change of variable $\rho = r^2$ is used, then the phase profile becomes periodic in ρ . Figure 8 shows the phase profile in both r - and ρ -space for a design wavelength of 555 nm, $\alpha = 0.5$, and $F = 250$ mm.

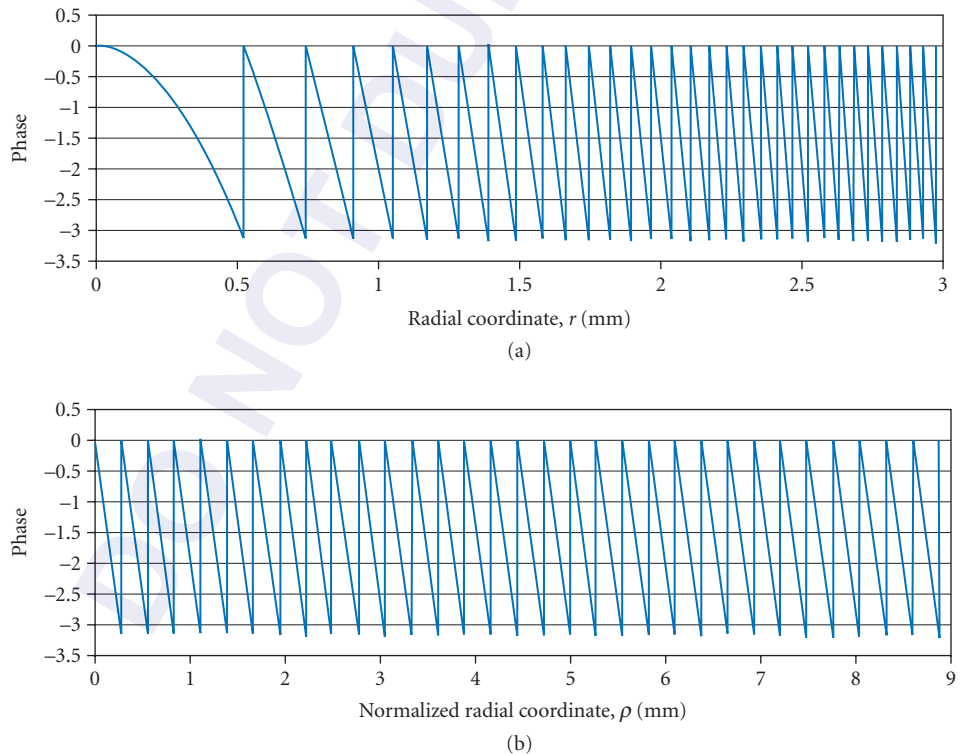


FIGURE 8 The phase profile of a diffractive lens as a function of r (a) and ρ (b). Note that (b) is periodic.

The periodic phase profile can be represented as a Fourier series such that

$$\exp[i\phi(\rho)] = \sum_m c_m \exp\left[-im2\pi\left(\frac{\rho}{2\lambda_o F}\right)\right] \quad (24)$$

where the coefficients c_m are given by

$$c_m = \frac{1}{2\lambda_o F} \int_0^{2\lambda_o F} \exp\left[-\frac{i\pi\alpha\rho}{\lambda_o F}\right] \exp\left[\frac{im2\pi\rho}{2\lambda_o F}\right] d\rho \quad (25)$$

Carrying out the integration in Eq. (25) gives

$$\exp[i\phi(r)] = \sum_m \exp[i\pi(m-\alpha)] \exp\left[-\frac{i\pi r^2}{\lambda_o(F/m)}\right] \text{sinc}[m-\alpha] \quad (26)$$

Note that each term in the series in Eq. (26) acts like a lens of focal length F/m . The sinc function dictates how much energy is distributed into each foci. The diffraction efficiency η_m is defined as

$$\eta_m = \text{sinc}^2[m-\alpha] \quad (27)$$

The diffraction efficiency describes the percent of going into each diffraction order. In the case where $m = 0$, the focal length of this order becomes infinite (i.e., zero power). The underlying refractive carrier of the IOL provides all of the power of the implant. Consequently, the refractive power of the IOL is chosen to correct the eye for distance vision. For $m = 1$, energy is distributed into the +1 diffraction order. The power of the IOL for this diffraction order is the underlying refractive power of the IOL plus the power $+1/F$ provided by the diffractive structure. The amount of energy going into each of these diffractive orders is η_0 and η_1 , respectively. Similarly, energy is sent to higher-diffractive orders that have efficiency η_m and add a power m/F to the underlying refractive power of the IOL. Table 5 summarizes the add power and the diffraction efficiency for a case where $\alpha = 0.5$. Clearly, most of the energy goes into the 0 and +1 diffraction orders.

The Tecnis ZM900 (Advanced Medical Optics, Irvine, CA) is an example of a diffractive lens with $\alpha = 0.5$ and $F = 250$ mm. The Acri.LISA IOL (Carl Zeiss Meditec, Oberkochen, Germany) is another example of a diffractive IOL. This IOL is currently only available outside the United States. This lens distributes about twice the energy into the zero-order as it does into the +1 order. This effect biases distance vision and can be achieved by making $\alpha = 0.414$. The ReSTOR IOL (Alcon Laboratories, Fort Worth, TX) is another variation of a diffractive lens.⁸⁸ These lenses have a diffractive structure over the central 3 mm of the lens and are purely refractive in their periphery. The refractive portion of the lens provides distance vision, while the diffractive portion of the lens provides both near and distance vision. The step heights between annular zones of the lens gradually decrease toward the edge of the diffractive zone. This decrease is called apodization. Traditionally in optics, apodization has referred to a variable transmission of the pupil of an optical system. However, apodization in this regard describes the gradual change in diffraction efficiency that occurs with these types of lenses. The net result of these types of lenses is that the energy sent to the distance and near foci is nearly equal for

TABLE 5 Add Power and Diffraction Efficiency for $\alpha = 0.5$ and $F = 250$ mm

Diffraction Order, m	Diffraction Efficiency, η_m	Add Power
-1	4.5%	-4 D
0	40.5%	0 D
+1	40.5%	4 D
+2	4.5%	8 D
+3	1.6%	12 D

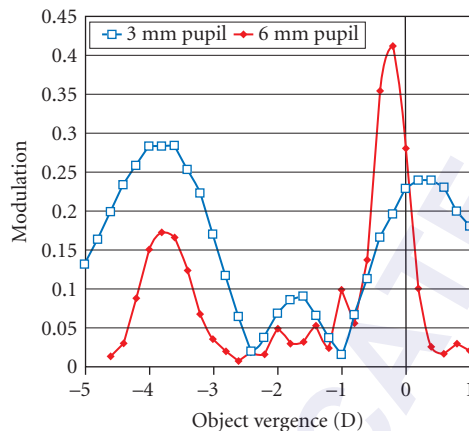


FIGURE 9 MTF measurements of the ReSTOR apodized diffractive IOL as a function of object vergence: 3- and 6-mm pupil.

small pupils, but is systematically shifted to a distance bias as the pupil size expands. Figure 9 shows the through-focus MTF of the ReSTOR lens for a 3- and 6-mm pupil. The horizontal axis of these plots is in units of vergence or the reciprocal of the object distance in units of inverse meters. As with zonal refractive lenses, there is some loss in visual performance with diffractive lenses.^{89–91} This loss is offset by roughly 78 percent of patients achieving near visual acuity of 20/40 or better.⁹¹ Biasing distance vision through modification of the step height or apodization tends to improve performance over conventional equal-split diffractive lenses for large pupil diameters.^{92,93} Stray light effects, flares, and haloes can be seen with diffractive lenses as well. Apodization of the IOL tends to markedly dampen these stray light effects for large pupils.⁹³ Finally, diffractive lenses tend to provide superior visual performance when compared to zonal refractive lenses.^{94,95}

Accommodating Lenses

Accommodating IOLs are the natural progression of implant technology. As shown in Fig. 1, accommodation is lost by about the age of 50. Implantation of IOLs may restore vision in cataract patients, but presbyopia remains. Multifocal IOLs have been developed to address the need of presbyopic pseudophakic patients, but there is always a trade-off with these types of lenses. Simultaneous vision by its nature reduces contrast and image quality, and designing multifocal IOLs is always a balance between providing functional near vision and minimizing artifacts attributed to multifocality. Accommodating lenses avoid the issues of simultaneous vision and provide high quality imaging at a variety of object distances. These types of lenses would function in a manner similar to the natural crystalline lens. Several accommodating IOLs are currently available; however, their performance has been marginal at best.⁹⁶ Advances in accommodating lens technology are rapidly occurring and this technology will likely replace multifocal IOLs once a suitable solution has been found. Furthermore, cataract surgery may no longer be the requirement for IOL implantation. Once solutions have been developed that offer reasonable levels of accommodation (likely allowing focusing from distance to 50 cm), a procedure known as refractive lens exchange (RLE) or clear lens extraction (CLE) is likely to rise in popularity.⁹⁷ In these procedures the noncataractous crystalline lens is removed and replaced with an IOL. These techniques are one option for patients with severe myopia, but surgeons are hesitant to perform these procedures on healthy, albeit presbyopic lenses with low degrees of refractive error. However, there was much resistance to implanting IOLs in the first place and once RLE or CLE demonstrates safe results that restore accommodation, the hesitancy will likely wane.

Two mechanisms are available for creating an accommodating IOL. The trigger for the lens action is the constriction of the ciliary muscle. The first mechanism is an axial shift of the position of the lens. Equation (8) describes the IOL power given the corneal curvature K , axial length L , and the effective lens position, ELP. Differentiating Eq. (8) with respect to the ELP gives

$$\frac{d\phi_{\text{IOL}}}{d(\text{ELP})} = \frac{(2K \cdot \text{ELP} - KL - 1000n_{\text{aq}})(KL - 1000n_{\text{aq}})n_{\text{aq}}}{(\text{ELP} - L)^2(K \cdot \text{ELP} - 1000n_{\text{aq}})^2} \quad (28)$$

Assuming a corneal power $K = 43$ D, and axial length $L = 24$ mm, $n_{\text{aq}} = 1.337$, $\text{ELP} = 5$ mm in the unaccommodated state and that the change in ELP is small compared to the total ELP, then the change in IOL power $\Delta\phi_{\text{IOL}}$ is approximately

$$\Delta\phi_{\text{IOL}} \approx 1.735 \Delta\text{ELP} \quad (29)$$

Equation (29) shows that under these conditions, the 1 mm of movement of the IOL toward the cornea gives about 1.75 D of accommodation. Examples of IOLs that use axial movement are the FDA-approved Crystalens (Bausch & Lomb, Rochester, NY) and 1CU (HumanOptics AG, Erlangen, Germany) available in Europe. Both lenses have plate-type haptics or flat flanges emanating from the side of the lens. These flanges are hinged such that the lens vaults toward the cornea when the ciliary muscle constricts. Still there is much improvement needed in these technologies. The Synchrony accommodating IOL (Visiogen, Irvine, CA) is a variation on the axial movement concept that uses two lenses that move relative to one another instead of a single lens for the IOL.

The second mechanisms for an accommodating IOL would be a lens that changes its power in response to the constriction of the ciliary muscle. Examples of these types of lenses include the Fluid Vision IOL (Powervision, Belmont, CA) which pumps fluid from a peripheral reservoir into a lens with a deformable membrane surface to change its power, and the AkkoLens (AkkoLens International, Delft, The Netherlands) that uses two cubic phase plates that translate laterally to achieve a power change.⁹⁸ These accommodating IOLs represent emerging technologies but clinical demonstration of their capabilities is still needed.

Phakic Lenses

Phakic IOLs, as their name implies, are lenses that are implanted while the crystalline lens remains clear and intact. These lenses are used for treating nearsightedness instead of cataracts. Since the crystalline lens remains in place, phakic IOLs need to be implanted in the remaining space. There are two suitable locations for this implantation. The first is the anterior chamber, which is the space between the posterior cornea and the iris. Two techniques have been used to support anterior chamber IOLs. The first technique is to wedge the haptics into the angle where the peripheral cornea meets the iris. The second technique is to use specially designed haptics that clip to the peripheral iris to maintain the lens position. The main disadvantage of anterior chamber phakic IOLs is the risk to the corneal endothelium. The corneal endothelium is a thin layer of cells that coats the posterior surface of the cornea. These cells regulate the nutrition and health of the cornea and cannot be replaced if damaged. Mounting an IOL in the anterior chamber risks abrasions of these endothelial cells. Phakic IOLs can also be supported by the sulcus, which means the phakic IOL is placed on the area directly behind the iris, but in front of the crystalline lens. This second type of phakic IOL positioning requires a vaulting of the implant away from the surface of the crystalline lens since contact between the artificial and natural lenses are likely to lead to cataract formation. Examples of phakic IOLs include the Verisyse (Advanced Medical Optics, Santa Ana, CA), which is an iris-supported lens, and the Visian ICL (Staar Surgical, Monrovia, CA), which is a sulcus-supported lens. Excellent results have been achieved with phakic IOLs in extreme nearsightedness.⁹⁹ Phakic IOLs for the treatment of farsightedness and astigmatism as well as multifocal lenses are likely to be available in the United States in the near future.

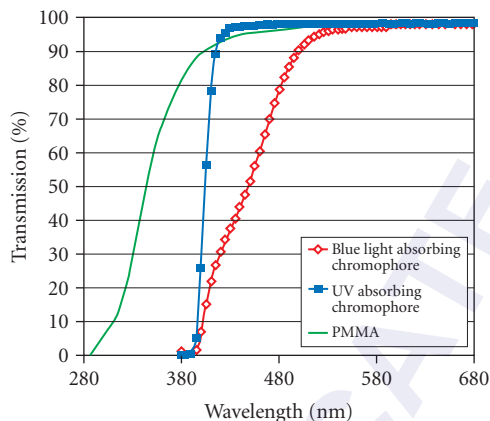


FIGURE 10 The transmission of a PMMA, and blue-light filtering, and UV-absorbing IOLs for ultra-violet and visible wavelengths.

Chromophores

The crystalline lens is the main absorber of wavelengths below 400 nm in the eye.^{2,100} Removal of the crystalline lens causes a dramatic increase in the level of ultraviolet (UV) radiation that reaches the retina. These highly energetic wavelengths can cause phototoxicity, leading to damage the retinal neural structure and the underlying retinal pigment epithelium. Aphakic individuals are particularly susceptible to this type of damage since they no longer have the appropriate UV filtering. IOLs, in general, do little to mitigate the UV exposure. PMMA has significant transmission for wavelengths between 300 and 400 nm. Mainster recognized the risks to retinal health posed by transparent IOLs^{101,102} and most IOLs incorporated a UV-absorbing chromophore by 1986.¹⁰³ These UV-absorbing chromophores typically have a cutoff wavelength of 400 nm, leaving high transmission of visible wavelengths, but negligible transmission of UV wavelengths. Blue-green wavelengths can also cause retinal phototoxicity, although much higher dosages are required for damage when compared to UV wavelengths. Several companies have introduced blue-light absorbing IOLs (AF-1 UY, Hoya, Tokyo, Japan and the Acrysof Natural, Alcon laboratories, Fort Worth, TX). These IOLs seek to have a transmission similar to the middle-aged human crystalline lens. Consequently, they partially absorb a portion of the spectrum between 400 and 500 nm. Figure 10 shows the transmission of a 3-mm slab of PMMA immersed in saline, along with the transmissions of IOLs containing UV-absorbing and blue-light absorbing IOLs.

Blue-light absorbing IOLs are somewhat controversial. Critics have suggested that scotopic and color vision are degraded by the lenses.^{104,105} However, these claims have not been supported by theoretical analyses and clinical measures of in-eye lens performance.^{106–113} Supporters of blue-light filtering lenses advocate their use to promote retinal health. Some laboratory studies have demonstrated protection to cell cultures afforded by the blue-light chromophores, but long-term clinical data remains unavailable.¹¹⁴ Blue light has been shown to mediate the circadian rhythm¹¹⁵ and an intermediate violet-absorbing chromophore that balances retinal protection and the circadian clock has been suggested.^{116,117} Again, long-term validation of the benefits of such a lens are not available.

21.5 INTRAOCULAR LENS SIDE EFFECTS

A variety of optical deficits can arise following the implantation of intraocular lenses (IOLs). While artificial lenses replace cataractous lenses that have degraded markedly in optical quality, these IOLs can also introduce some visual side effects. These side effects include glare, halos, streaks, starbursts,

shadows, and haze. While the vast majority of pseudophakic patients are free from these effects under most circumstances and tolerate the problem in the few situations when they do arise, a small fraction of patients suffer from problems that are comparable or worse than their preoperative state. Understanding the cause of these optical phenomena allows for IOL manufacturers to improve their designs and minimize potential problems in future recipients.

Dysphotopsia

Dysphotopsia is the introduction of unwanted patterns onto the retina. These unwanted patterns are superimposed over the true retinal image and can cause degradation in visual performance. In pseudophakic dysphotopsia, the design and material of the artificial lens is typically responsible for redirecting the unwanted light to the retina. Dysphotopsia comes in two forms: positive and negative. Positive dysphotopsia is the introduction of bright artifacts onto the retina. These artifacts include arcs, streaks, rings, and halos and may only be present under certain lighting conditions or for certain locations of glare sources in the peripheral field. Negative dysphotopsia, conversely, is the blockage of light from reaching certain portions of the retina. Shadows and dark spots are perceived usually in the temporal field and again this phenomenon is affected by lighting conditions and position. Both positive and negative dysphotopsias affect visual performance because the unwanted images obscure the relevant retinal image formed directly by the IOL. Consequently, understanding their cause and eliminating these undesirable effects will lead to improve performance and satisfaction of patients.

Posterior Capsule Opacification

Another side effect of IOL implantation is posterior capsule opacification (PCO). Implantation of IOLs into the capsular bag provides a stable platform for the lens. However, in about 25 percent of patients, PCO can result. In PCO, the capsular bag gradually opacifies, and postcataract patients perceive many of the same symptoms seen prior to surgery, namely a gradual loss in acuity, glare, and blurred vision. YAG posterior capsulotomy is typically performed to open a hole in the capsule allowing light to pass through again. The large incidence of PCO causes dissatisfaction in patients and adds expense due to the need for additional equipment, procedures, and time to rectify the situation. Consequently, a reduction in PCO would be beneficial and strong efforts have been made to modify IOL designs and materials to drastically reduce the incidence of PCO. Acrylic materials and square edge designs have both appear to reduce the incidence of PCO.^{118,119} However, dysphotopsia increased with introduction of these changes. To gain the benefit of reduced PCO, alterations were needed to acrylic lens edge designs.

Holladay et al.¹²⁰ compared the edge glare caused by sharp and rounded edge designs using a nonsequential raytracing technique. Nonsequential raytracing is routinely used for analyzing stray light and illumination effects. Typically, hundreds of thousands of rays are launched from a light source and allowed to refract, reflect, and scatter from various elements. Following the raytracing, the concentration of these rays in a given plane can be evaluated to determine the illumination pattern of the stray light. Holladay et al. found that both square and round edge IOLs produce stray light, but that only the square-edge design concentrated the light into a well-formed arc on the retina. Round-edge designs tended to disperse the stray light over a much larger portion of the retina, suggesting its visual consequences fall below a perceptible threshold. Clinical reports of dysphotopsia support that square-edge designs with a smooth finish are a likely culprit for many of the positive dysphotopsia problems.^{121,122} Meacock et al.¹²³ performed a prospective study of 60 patients split between acrylic lenses with textured and nontextured edges. By prospectively analyzing the two groups, determination about the advantages of textured edges is assessed. One month post-operatively 67 percent of the nontextured IOL patients and 13 percent of textured-IOL patients had glare symptoms. The textured edges provided a statistically significant reduction in glare symptoms. Franchini et al.^{124,125} used nonsequential raytracing methods to compare different types of edge design on positive dysphotopsia. They found similar result to Holladay et al.¹²⁰ in that the untextured square-edge design produced a

ring pattern on the retina, while a round edge distributed the light over a larger portion of the retina. This group also analyzed an “OptiEdge” design where the edge is beveled in an attempt to minimize stray light effects. In this case, they found an arc pattern is formed on the retina, but this arc in general has intensity far below the smooth edge pattern. Furthermore, Franchini et al. modeled a frosted square-edge design and found that the circular pattern retina is reduced, but the light is distributed over the retina and could possibly reduce contrast in the perceived image. The results of these efforts demonstrate that tailoring the shape and texture of the edge of the IOL can lead to an implant that reduces both PCO and dysphotopsia.

21.6 SUMMARY

The addition of IOL implantation following cataract removal revolutionized the procedure and enabled the restoration of high quality vision. While the technique was slow to be accepted by the medical establishment, the benefits of IOLs won over even the harshest critics to become the standard of care. IOLs are entering a postmodern era. Conventional IOLs are reliably fabricated and advances in materials and design have led to small incision surgeries that minimize risks to the patient and allowing outpatient procedures. Power calculations and the devices used to measure ocular dimensions have continued to improve, narrowing the residual error following implantation. Manufacturers are now turning to more subtle improvements to the technology. Aspheric surfaces are being incorporated to minimize ocular spherical aberration. Multifocal IOLs are providing simultaneous distance and near vision, enabling a spectacle-free lifestyle following surgery. Accommodating IOLs are emerging to fully restore accommodation. Furthermore, manufacturers are expanding the role of IOLs into the treatment of refractive error with phakic IOLs. Finally, chromophores are incorporated into IOLs to promote the health of the retina for the remainder of the patient’s life. IOLs have been extraordinarily successful and a true benefit to the dim world of cataracts.

21.7 REFERENCES

1. W. J. Smith, *Modern Optical Engineering*, 4th ed., McGraw-Hill, New York, 2000.
2. D. A. Atchison and G. Smith, *Optics of the Human Eye*, Butterworth-Heinemann, Oxford, 2000.
3. A. Duane, “Normal Values of the Accommodation at all Ages,” *Trans. Sec. Ophthalmol. AMA* **53**:383–391 (1912).
4. B. K. Pierscionek, “Reviewing the Optics of the Lens and Presbyopia,” *Proc. SPIE* **3579**:34–39 (1998).
5. National Eye Institute, “Cataract—What You Should Know,” *NIH Publication* 03-201 (2003).
6. J. C. Javitt, F. Wang, and S. K. West, “Blindness due to Cataract: Epidemiology and Prevention,” *Annu. Rev. Public Health* **17**:159–177 (1996).
7. B. E. K. Klein, R. Klein, and K. L. P. Linton. “Prevalence of Age-Related Lens Opacities in a Population,” *Ophthalmology* **99**:546–552 (2003).
8. K. B. Kansupada and J. W. Sassani, “Sushruta: The Father of Indian Surgery and Ophthalmology,” *Doc. Ophthalmology* **93**:159–167 (1997).
9. M. Jalie, *The Principles of Ophthalmic Lenses*, 4th ed., Association of Dispensing Opticians, London, 1984.
10. R. J. Schechter, “Optics of Intraocular Lenses,” in *Duane’s Clinical Ophthalmology*, eds., W. Tasman and E. A. Jaeger, Lippincott, Philadelphia, 2004.
11. D. J. Apple and J. Sims, “Harold Ridley and the Invention of the Intraocular Lens,” *Surv. Ophthalmol.* **40**:279–292 (1996).
12. D. J. Apple, *Sir Harold Ridley and His Fight for Sight*, Slack, New Jersey, 2006.
13. I. Obuchowska and Z. Mariak, “Jacques Daviel—The Inventor of the Extracapsular Cataract Extraction Surgery,” *Klin. Oczna.* **107**:567–571 (2005).

14. C. D. Kelman, "Phacoemulsification and Aspiration. A New Technique of Cataract Removal. A Preliminary Report," *Am. J. Ophthalmol.* **64**:23–35 (1967).
15. J. Retzlaff, "A New Intraocular Lens Calculation Formula," *Am. Intra-Ocular Implant Soc. J.* **6**:148–152 (1980).
16. D. R. Sanders and M. C. Kraff, "Improvement of Intraocular Lens Power Calculation Using Empirical Data," *Am. Intra-Ocular Implant Soc. J.* **6**:268–270 (1980).
17. R. D. Binkhorst, "The Accuracy of Ultrasonic Measurement of Axial Length of the Eye," *Ophthalmic. Surg.* **13**:363–365 (1981).
18. T. Olsen, "The Accuracy of Ultrasonic Determination of Axial Length in Pseudophakic Eyes," *Acta Ophthalmol.* **67**:141–144 (1989).
19. W. Drexler, O. Findl, R. Menapace, G. Rainer, C. Vass, C. K. Hitzenberger, and A. F. Fercher, "Partial Coherence Interferometry: A Novel Approach to Biometry in Cataract Surgery," *Am. J. Ophthalmol.* **126**: 524–534 (1998).
20. D. R. Sanders, J. R. Retzlaff, and M. C. Kraff, "Comparison of the SRK II Formula and Other Second Generation Formulas," *J. Cataract Refract. Surg.* **14**:136–141 (1988).
21. T. Olson, "Calculation of Intraocular Lens Power: A Review," *Acta Ophthalmol. Scand.* **85**:472–485 (2007).
22. O. Pomerantzef, M. M. Pankratov, and G. Wang, "Calculation of an IOL from the Wide-Angle Optical Model of the Eye," *Am. Intra-Ocular Implant Soc. J.* **11**:37–43 (1985).
23. J. T. Holladay, T. C. Prager, T. Y. Chandler, K. H. Musgrove, J. W. Lewis, and R. S. Ruiz, "A Three-Part System for Refining Intraocular Lens Power Calculations," *J. Cataract Refract. Surg.* **14**:17–23 (1988).
24. J. A. Retzlaff, D. R. Sanders, and M. C. Kraff, "Development of the SRK/T Intraocular Lens Implant Power Calculation Formula," *J. Cataract Refract. Surg.* **16**:333–340 (1990).
25. K. J. Hoffer, "The Hoffer Q Formula: A Comparison of Theoretic and Regression Formulas," *J. Cataract Refract. Surg.* **19**:700–712 (1993).
26. J. T. Holladay and K. J. Maverick, "Relationship of the Actual Thick Intraocular Lens Optic to the Thin Lens Equivalent," *Am. J. Ophthalmol.* **126**:339–347 (1998).
27. J. T. Holladay, "International Intraocular Lens and Implant Registry," *J. Cataract Refract. Surg.* **26**:118–134 (2000).
28. J. T. Holladay, "International Intraocular Lens and Implant Registry," *J. Cataract Refract. Surg.* **27**:143–164 (2001).
29. J. T. Holladay, "International Intraocular Lens and Implant Registry," *J. Cataract Refract. Surg.* **28**:152–174 (2002).
30. J. T. Holladay, "International Intraocular Lens and Implant Registry," *J. Cataract Refract. Surg.* **29**:176–197 (2003).
31. J. T. Holladay, "International Intraocular Lens and Implant Registry," *J. Cataract Refract. Surg.* **30**:207–229 (2004).
32. L. C. Celikkol, G. Pavlopoulos, B. Weinstein, G. Celikkol, and S. T. Feldman, "Calculation of Intraocular Lens Power after Radial Keratotomy with Computerized Videokeratography," *Am. J. Ophthalmol.* **120**:739–749 (1995).
33. H. V. Gimbel and R. Sun, "Accuracy and Predictability of Intraocular Lens Power Calculation after Laser in situ Keratomileusis," *J. Cataract Refract. Surg.* **27**:571–576 (2001).
34. J. H. Kim, D. H. Lee, and C. K. Joo, "Measuring Corneal Power for Intraocular Lens Power Calculation after Refractive Surgery," *J. Cataract Refract. Surg.* **28**:1932–1938 (2002).
35. N. Rosa, L. Capasso, and A. Romano, "A New Method of Calculating Intraocular Lens Power after Photorefractive Keratectomy," *J. Refract. Surg.* **18**:720–724 (2002).
36. A. A. Stakheev, "Intraocular Lens Calculation for Cataract after Previous Radial Keratotomy," *Ophthal. Physiol. Opt.* **22**:289–295 (2002).
37. C. Argento, M. J. Cosentino, and D. Badoza, "Intraocular Lens Power Calculation after Refractive Surgery," *J. Refract. Surg.* **29**:1346–1351 (2003).
38. L. Chen, M. H. Mannis, J. J. Salz, F. J. Garcia-Ferrer, and J. Ge, "Analysis of Intraocular Lens Power Calculation in Pose-Radial Keratotomy Eyes," *J. Cataract Refract. Surg.* **29**:65–70 (2003).
39. E. Jarade and K. F. Tabbara, "New Formula for Calculating Intraocular Lens Power after Laser in situ Keratomileusis," *J. Cataract Refract. Surg.* **30**:1711–1715 (2004).

40. G. Ferrara, G. Cennamo, G. Marotta, and E. Loffredo, "New Formula to Calculate Corneal Power after Refractive Surgery," *J. Refract. Surg.* **20**:465–471 (2004).
41. G. Savini, P. Barboni, and M. Zanini, "Intraocular Lens Power Calculation after Myopic Refractive Surgery," *Ophthalmology* **113**:1271–1282 (2006).
42. S. T. Awwad, S. Dwarakanathan, W. Bowman, D. Cavanagh, S. M. Verity, V. V. Mootha, and J. P. McCulley, "Intraocular Lens Power Calculation after Radial Keratotomy: Estimating the Refractive Corneal Power," *J. Refract. Surg.* **33**:1045–1050 (2007).
43. L. Wang, E. Dai, D. D. Koch, and A. Nathoo, "Optical Aberrations of the Human Anterior Cornea," *J. Cataract Refract. Surg.* **29**:1514–1521 (2003).
44. J. Porter, A. Guirao, I. G. Cox, and D. R. Williams, "Monochromatic Aberrations of the Human Eye in a Large Population," *J. Opt. Soc. Am. A* **18**:1793–1803 (2001).
45. G. Smith and C. -W. Lu, "The Spherical Aberration of Intraocular Lenses," *Ophthal. Physiol. Opt.* **8**:287–294 (1988).
46. W. T. Welford, *Aberrations of Optical Systems*, Adam Hilger, Bristol, 1986.
47. A. Castro, P. Rosales, and S. Marcos, "Tilt and Decentration of Intraocular Lenses In Vivo from Purkinje and Scheimpflug Imaging," *J. Cataract Refract. Surg.* **33**:418–429 (2007).
48. D. Atchison, "Optical Design of Intraocular Lenses. III. On-Axis Performance in the Presence of Lens Displacement," *Optom. Vis. Sci.* **66**:671–681 (1989).
49. D. Atchison, "Refractive Errors Induced by Displacement of Intraocular Lenses Within the Pseudophakic Eye," *Optom. Vis. Sci.* **66**:146–152 (1989).
50. J. T. Holladay, P. A. Piers, G. Koranyi, M. van der Mooren, and S. Norrby, "A New Intraocular Lens Design to Reduce Spherical Aberration of Pseudophakic Eyes," *J. Refract. Surg.* **18**:683–692 (2002).
51. P. M. Kiely, G. Smith, and L. G. Carney, "The Mean Shape of the Human Cornea," *Optica Acta* **29**:1027–1040 (1982).
52. M. Guillon, D. P. M. Lyndon, and C. Wilson, "Corneal Topography: A Clinical Model," *Ophthal. Physiol. Opt.* **6**:47–56 (1986).
53. M. Dubbelman, H. A. Weeber, and R. G. L. van der Heijde, H. J. Völker-Dieben, "Radius and Asphericity of the Posterior Corneal Surface Determined by Corrected Scheimpflug Photography," *Acta Ophthalmol. Scand.* **80**:379–383 (2002).
54. M. Dubbelman, V. A. Sicam, and R. G. L. van der Heijde, "The Shape of the Anterior and Posterior Surface of the Aging Human Cornea," *Vis. Res.* **46**:993–1001 (2006).
55. A. Guirao, M. Redondo, and P. Artal, "Optical Aberrations of the Human Cornea as a Function of Age," *J. Opt. Soc. Am. A* **17**:1697–1702 (2000).
56. T. Oshika, S. D. Klyce, R. A. Applegate, H. Howland, and M. Danasoury, "Comparison of Corneal Wavefront Aberrations after Photorefractive Keratectomy and Laser In Situ Keratomileusis," *Am. J. Ophthalmol.* **127**:1–7 (1999).
57. H. H. Hopkins, *Wave Theory of Aberrations*, Clarendon, Oxford, 1950.
58. J. Schwiegerling, J. G. Greivenkamp, and J. M. Miller, "Representation of Videokeratoscopic Height Data with Zernike Polynomials," *J. Opt. Soc. Am. A* **12**:2105–2113 (1995).
59. L. Wang, E. Dai, D. D. Koch, and A. Nathoo, "Optical Aberrations of the Human Anterior Cornea," *J. Cataract Refract. Surg.* **29**:1514–1521 (2003).
60. R. Bellucci, S. Morselli, and P. Piers, "Comparison of Wavefront Aberrations and Optical Quality of Eyes Implanted with Five Different Intraocular Lenses," *J. Refract. Surg.* **20**:297–306 (2004).
61. A. Guirao, J. Tejedor, and P. Artal, "Corneal Aberrations before and after Small-Incision Cataract Surgery," *Invest. Ophthalmol. Vis. Sci.* **45**:4312–4319 (2004).
62. S. Norrby, P. Piers, C. Campbell, and M. van der Mooren, "Model Eyes for the Evaluation of Intraocular Lenses," *Appl. Opt.* **46**:6595–6605 (2007).
63. D. F. Chang, "When Do I Use LRIs versus Toric IOLs," in *Mastering Refractive IOLs: The Art and Science*, ed., D. F. Chang, Slack, New Jersey, 2008.
64. A. Langenbucher and B. Seitz, "Computerized Calculation Scheme for Bitoric Eikonic Intraocular Lenses," *Ophthal. Physiol. Opt.* **23**:213–220 (2003).

65. A. Langenbucher, S. Reese, T. Sauer, and B. Seitz, "Matrix-Based Calculation Scheme for Toric Intraocular Lenses," *Ophthalm. Physiol. Opt.* **24**:511–519 (2004).
66. A. Langenbucher and B. Seitz, "Computerized Calculation Scheme for Toric Intraocular Lenses," *Acta Ophthalmol. Scand.* **82**:270–276 (2004).
67. A. Langenbucher, N. Szentmary, and B. Seitz, "Calculating the Power of Toric Phakic Intraocular Lenses," *Ophthalmic. Physiol. Opt.* **27**:373–380 (2007).
68. L. Carretero, R. Fuentes, and A. Fimia, "Measurement of Spherical Aberration of Intraocular Lenses with the Ronchi Test," *Optom. Vis. Sci.* **69**:190–192 (1992).
69. G. Bos, J. M. Vanzo, P. Maufof, and J. L. Gutzwiller, "A New Interferometric to Assess IOL Characteristics," *Proc. SPIE* **2127**:56–61 (1994).
70. R. E. Fischer and K. C. Liu, "Advanced Techniques for Optical Performance Characterization of Intraocular Lenses," *Proc. SPIE* **2127**:14–25 (1994).
71. R. Sun and V. Portney, "Multifocal Ophthalmic Lens Testing," *Proc. SPIE* **2127**:82–87 (1994).
72. B. G. Broome, "Basic Considerations for IOL MTF Testing," *Proc. SPIE* **2127**:2–15 (1994).
73. J. S. Chou, L. W. Blake, J. M. Fridge, and D. A. Fridge, "MTF Measurement System that Simulates IOL Performances in the Human Eye," *Proc. SPIE* **2393**:271–279 (1995).
74. E. Keren and A. L. Rotlex, "Measurement of Power, Quality, and MTF of Intraocular and Soft Contact Lenses in Wet Cells," *Proc. SPIE* **2673**:262–273 (1996).
75. D. Tognetto, G. Sanguinetti, P. Sirotti, P. Cecchini, L. Marcucci, E. Ballone, and G. Ravalico, "Analysis of the Optical Quality of Intraocular Lenses," *Invest. Ophthalmol. Vis. Sci.* **45**:2686–2690 (2004).
76. P. A. Piers, N. E. S. Norrby, and U. Mester, "Eye Models for the Prediction of Contrast Vision in Patients with New Intraocular Lens Designs," *Opt. Lett.* **29**:733–735 (2004).
77. R. Rawer, W. Stork, C. W. Spraul, and C. Lingenfelder, "Imaging Quality of Intraocular Lenses," *J. Cataract Refract. Surg.* **31**:1618–1631 (2005).
78. P. G. Gobbi, F. Fasce, S. Bozza, and R. Brancato, "Optomechanical Eye Model with Imaging Capabilities for Objective Evaluation of Intraocular Lenses," *J. Cataract Refract. Surg.* **32**:643–651 (2006).
79. International Standard 11979-2, "Ophthalmic Implants—Intraocular Lenses. Part 2: Optical Properties and Test Methods," ISO, Geneva, 1999.
80. T. Avitabile and F. Marano, "Multifocal Intraocular Lenses," *Curr. Opin. Ophthalmol.* **12**:12–16 (2001).
81. J. D. Hunkeler, T. M. Coffman, J. Paugh, A. Lang, P. Smith, and N. Tarantino, "Characterization of Visual Phenomena with the Array Multifocal Intraocular Lens," *J. Cataract Refract. Surg.* **28**:1195–1204 (2002).
82. H. N. Sen, A. U. Sarikkola, R. J. Uusitalo and L. Llaatikainen, "Quality of Vision after AMO Array Multifocal Intraocular Lens Implantation," *J. Cataract Refract. Surg.* **31**:2483–2493 (2004).
83. K. A. Featherstone, J. R. Bloomfield, A. J. Lang, M. J. Miller-Meeks, G. Woodworth, R. F. Steinert, "Driving Simulation Study: Bilateral Array Multifocal versus Bilateral AMO Monofocal Intraocular Lenses," *J. Cataract Refract. Surg.* **25**:1254–1262 (1999).
84. M. Larsson, C. Beckman, A. Nyström, S. Hård and J. Sjöstrand, "Optical Properties of Diffractive, Bifocal Intraocular lenses," *Proc. SPIE* **1529**:63–70 (1991).
85. A. Issacson, "Global Status of Diffraction Optics as the Basis for an Intraocular Lens," *Proc. SPIE* **1529**:71–79 (1991).
86. M. J. Simpson, "Diffractive Multifocal Intraocular Lens Image Quality," *Appl. Opt.* **31**:3621–3626 (1992).
87. D. Faklis and G. M. Morris, "Spectral Properties of Multiorder Diffractive Lenses," *Appl. Opt.* **34**:2462–2468 (1995).
88. J. A. Davison and M. J. Simpson, "History and Development of the Apodized Diffractive Intraocular Lens," *J. Cataract Refract. Surg.* **32**:849–858 (2006).
89. H. V. Gimbel, D. R. Sanders, and M. G. Raanan, "Visual and Refractive Results of Multifocal Intraocular Lenses," *Ophthalmology* **98**:881–887 (1991).
90. C. T. Post, "Comparison of Depths of Focus and Low-Contrast Acuties for Monofocal versus Multifocal Intraocular Lens Patients at 1 Year," *Ophthalmology* **99**:1658–1663 (1992).
91. R. L. Lindstrom, "Food and Drug Administration Study Update. One Year Results from 671 Patients with the 3M Multifocal Intraocular Lens," *Ophthalmology* **100**:91–97 (1993).

92. G. Schmidinger, C. Simander, I. Dejaco-Ruhswurm, C. Skorpik, and S. Pieh, "Contrast Sensitivity Function in Eyes with Diffractive Bifocal Intraocular Lenses," *J. Cataract Refract. Surg.* **31**:2076–2083 (2005).
93. J. Choi and J. Schwiegerling, "Optical Performance Measurement and Night Driving Simulation of ReSTOR, ReZoom, and Tecnis Multifocal Intraocular Lenses in a Model Eye," *J. Refract. Surg.* **24**:218–222 (2008).
94. W. W. Hutz, B. Eckhardt, B. Rohrig, and R. Grolmus, "Reading Ability with 3 Multifocal Intraocular Lens Models," *J. Cataract Refract. Surg.* **32**:2015–2021 (2006).
95. U. Mester, W. Hunold, T. Wesendahl, and H. Kaymak, "Functional Outcomes after Implantation of Tecnis ZM900 and Array SA40 Multifocal Intraocular Lenses," *J. Cataract Refract. Surg.* **33**:1033–1040 (2007).
96. O. Findl and C. Leydolt, "Meta-Analysis of Accommodating Intraocular Lenses," *J. Cataract Refract. Surg.* **33**:522–527 (2007).
97. M. Packer, I. H. Fine, and R. S. Hoffman, "The Crystalline Lens as a Target for Refractive Surgery," in *Refractive Lens Surgery*, eds., I. H. Fine, M. Packer, and R. S. Hoffman, Springer-Verlag, Berlin, 2005.
98. A. N. Simonov and G. Vdovin, "Cubic Optical Elements for an Accommodative Intraocular Lens," *Opt. Exp.* **14**:7757–7775 (2006).
99. I. Brunette, J. M. Bueno, M. Harissi-Dagher, M. Parent, M. Podtetenov, and H. Hamam, "Optical Quality of the Eye with the Artisan Phakic Lens for the Correction of High Myopia," *Optom. Vis. Sci.* **80**:167–174 (2003).
100. E. A. Boettner and J. R. Wolter, "Transmission of the Ocular Media," *Invest. Ophthalmol.* **1**:776–783 (1962).
101. M. A. Mainster, "Spectral Transmittance of Intraocular Lenses and Retinal Damage from Intense Light Sources," *Am. J. Ophthalmol.* **85**:167–170 (1978).
102. M. A. Mainster, "Solar Retinitis, Photoc Maculopathy and the Pseudophakic Eye," *J. Am. Intraocul. Implant Soc.* **4**:84–86 (1978).
103. M. A. Mainster, "The Spectra, Classification, and Rationale of Ultraviolet-Protective Intraocular Lenses," *Am. J. Ophthalmol.* **102**:727–732 (1986).
104. M. A. Mainster and J. R. Sparrow, "How Much Blue Light Should an IOL Transmit," *Br J Ophthalmol.* **87**:1532–1529 (2003).
105. M. A. Mainster, "Blue-Blocking Intraocular Lenses and Pseudophakic Scotopic Sensitivity," *J. Cataract Refract. Surg.* **32**:1403–1406 (2006).
106. K. Niwa, Y. Yoshino, F. Okuyama, and T. Tokoro, "Effects of Tinted Intraocular Lens on Contrast Sensitivity," *Ophthalm. Physiol. Opt.* **16**:297–302 (1996).
107. S. M. Raj, A. R. Vasavada, and M. A. Nanavaty, "AcrySof Natural SN60AT versus AcrySof SA60AT Intraocular Lens in Patients with Color Vision Defects," *J. Cataract Refract. Surg.* **31**:2324–2328 (2005).
108. A. Rodriguez-Galietero, R. Montes-Mico, G. Munoz, and C. Albarran-Diego, "Blue-Light Filtering Intraocular Lens in Patients with Diabetes: Contrast Sensitivity and Chromatic Discrimination," *J. Cataract Refract. Surg.* **31**:2088–2092 (2005).
109. A. R. Galietero, R. M. Mico, G. Munoz, and C. A. Diego, "Comparison of Contrast Sensitivity and Color Discrimination after Clear and Yellow Intraocular Lens Implantation," *J. Cataract Refract. Surg.* **31**:1736–1740 (2005).
110. J. Schwiegerling, "Blue-Light-Absorbing Lenses and Their Effect on Scotopic Vision," *J. Cataract Refract. Surg.* **32**:141–144 (2006).
111. R. J. Cionni and J. H. Tsai, "Color Perception with AcrySof Natural and AcrySof Single-Piece Intraocular Lenses under Photopic and Mesopic Conditions," *J. Cataract Refract. Surg.* **23**:236–242 (2006).
112. N. Kara-Junior, J. L. Jardim, E. O. Leme, M. Dall'Col, and R. S. Junior, "Effect of the AcrySof Natural Intraocular Lens on Blue-Yellow Perimetry," *J. Cataract Refract. Surg.* **32**:1328–1330 (2006).
113. V. C. Greenstein, P. Chiosi, P. Baker, W. Seiple, K. Holopigian, R. E. Braunstein, and J. R. Sparrow, "Scotopic Sensitivity and Color Vision with a Blue-Light-Absorbing Intraocular Lens," *J. Cataract Refract. Surg.* **33**:667–672 (2007).
114. J. R. Sparrow, A. S. Miller, and J. Zhou, "Blue Light-Absorbing Intraocular Lens and Retinal Pigment Epithelium Protection In Vitro," *J. Cataract Refract. Surg.* **30**:873–878 (2004).
115. G. C. Brainard, J. P. Hanifin, J. M. Greeson, B. Byrne, G. Glickman, E. Gerner, M. D. Rollag, "Action Spectrum for Melatonin Regulation in Humans: Evidence for a Novel Circadian Photoreceptor," *J. Neurosci.* **21**:6405–6412 (2001).

116. M. A. Mainster, "Violet and Blue Light Blocking Intraocular Lenses: Photoprotection versus Photoreception," *Br. J. Ophthalmol.* **90**:784–792 (2006).
117. J. Kraats and D. Norren, "Sharp Cutoff Filters in Intraocular Lenses Optimize the Balance between Light Reception and Light Protection," *J. Cataract Refract. Surg.* **33**:879–887 (2007).
118. E. J. Hollick, D. J. Spalton, P. G. Ursell, M. V. Pande, S. A. Barman, J. F. Boyce, and K. Tilling, "The Effect of Polymethylmethacrylate, Silicone, and Polyacrylic Intraocular Lenses on Posterior Capsular Opacification 3 Years after Cataract Surgery," *Ophthalmology* **106**:49–54 (1999).
119. Q. Peng, N. Visessook, D. J. Apple, S. K. Pandey, L. Werner, M. Escobar-Gomez, R. Schoderbek, K. D. Solomon, and A. Guindi, "Surgical Prevention of Posterior Capsule Opacification. Part 3: Intraocular Lens Optic Barrier Effect as a Second Line of Defense," *J. Cataract Refract. Surg.* **26**:198–213 (2000).
120. J. T. Holladay, A. Lang, and V. Portney, "Analysis of Edge Glare Phenomena in Intraocular Lens Edge Designs," *J. Cataract Refract. Surg.* **25**:748–752 (1999).
121. S. Masket, "Truncated Edge Design, Dysphotopsia, and Inhibition of Posterior Capsule Opacification," *J. Cataract Refract. Surg.* **26**:145–147 (2000).
122. J. A. Davison, "Positive and Negative Dysphotopsia in Patients with Acrylic Intraocular Lenses," *J. Cataract Refract. Surg.* **26**:1346–1355 (2000).
123. W. R. Meacock, D. J. Spalton, and S. Khan, "The Effect of Texturing the Intraocular Lens Edge on Postoperative Glare Symptoms: A Randomized, Prospective, Double-Masked Study," *Arch. Ophthalmol.* **120**:1294–1298 (2002).
124. A. Franchini, B. Z. Gallarati, and E. Vaccari, "Computerized Analysis of the Effects of Intraocular Lens Edge Design on the Quality of Vision in Pseudophakic Patients," *J. Cataract Refract. Surg.* **29**:342–347 (2003).
125. A. Franchini, B. Z. Gallarati, and E. Vaccari, "Analysis of Stray-Light Effects Related to Intraocular Lens Edge Design," *J. Cataract Refract. Surg.* **30**:1531–1536 (2004).

This page intentionally left blank.

DO NOT DUPLICATE

DISPLAYS FOR VISION RESEARCH

William Cowan

Department of Computer Science
University of Waterloo
Waterloo, Ontario, Canada

22.1 GLOSSARY

$D(\lambda)$	spectral reflectance of the diffuser
e_p	phosphor efficiency
I_B	beam current of a CRT
M_{ij}	coefficient ij in the regression matrix linking ΔX_i to ΔR_i
m_j	vector indicating the particular sample measured
N_p	number of monochrome pixels in a color pixel
n_j	vector indicating the particular sample measured
$R(\lambda)$	spectral reflectance of the faceplate of an LCD
R_j	input coordinates of the CRT: usually R_1 is the voltage sent to the red gun; R_2 is the voltage sent to the green gun; and R_3 is the voltage sent to the blue gun
V	voltage input to a CRT
V_A	acceleration voltage of a CRT
V_0	maximum voltage input to a CRT
v_a	voltage applied to gun a of a CRT; $a = R, G, B$
v_B	maximum scanning velocity for the electron beam
v_h	horizontal velocity of the beam
v_v	vertical velocity of the beam
X_{ai}	tristimulus values of light emitted from monochrome pixel of color a
X_i	tristimulus values of light emitted from a color pixel
X_i	tristimulus value i , that is, $X = X_1$; $Y = X_2$; and $Z = X_3$
X_{0i}	tristimulus value i of a reference color
x_p	horizontal interpixel spacing, $x_p = v_h \tau_p$
y_p	vertical interpixel (interline) spacing

γ	exponent in power law expressions of gamma correction
δ	interpolation parameter for one-dimensional characterizations
ΔR_i	change in input coordinate i
ΔX_i	change in tristimulus value i
ϵ	interpolation parameter for inverting one-dimensional characterizations
V_{\max}	maximum frequency of CRT input amplifiers
μ	index of color measurements taken along a curve in color space
Φ	power of the emitted light
$\Phi_{a\lambda}^{(AMB)}(V_a)$	spectral power distribution of light emitted from monochrome pixel of color a as a result of ambient light falling on an LCD
$\Phi_{a\lambda}^{(BL)}$	spectral power distribution of light emitted from monochrome pixel of color a as a result of the backlight
$\Phi_{a\lambda}(V_a)$	spectral power distribution of light emitted by monochrome pixel of color a , depending on the voltage with which the pixel is driven
$\Phi_{a\lambda}^{(R)}$	spectral power distribution of light reflected from the faceplate of an LCD
$\Phi(x, y)$	power of the emitted light as a function of screen position
Φ_λ	spectral power distribution of light emitted by all monochrome pixels in a color pixel
Φ_λ	spectral power of the emitted light
$\Phi_\lambda^{(AMB)}$	light output caused by ambient light falling on the faceplate of an LCD
$\Phi_\lambda^{(BL)}$	light output caused by the backlight of an LCD
Φ_0	maximum light power output from a CRT
$\Phi_{0\lambda}^{(AMB)}$	ambient light falling on the faceplate of an LCD
$\Phi_{0\lambda}^{(BL)}$	light emitted by the backlight-diffuser element of an LCD
$\Phi_\lambda(v_R, v_G, v_B)$	spectral power distribution of the light emitted when voltages v_R, v_G, v_B are applied to its inputs
$\tau_a(\lambda)$	spectral transmittance of the filter on monochrome pixel of color a
τ_d	phosphor decay time
τ_f	time spent scanning a complete frame
τ_l	time spent scanning a complete line, including horizontal flyback
τ_p	time the beam spends crossing a single pixel
$\tau(V)$	voltage-dependent transmittance of the light modulating element

22.2 INTRODUCTION

Complex images that are colorimetrically calibrated are needed for a variety of applications, from color prepress to psychophysical experimentation. Unfortunately, such images are extremely difficult to produce, especially using traditional image production technologies such as photography or printing. In the last decade technical advances in computer graphics have made digital imaging the dominant technology for all such applications, with the color television monitor the output device of choice. This chapter describes the operational characteristics and colorimetric calibration techniques for color television monitors, with its emphasis on methods that are likely to be useful for psychophysical experimentation. A short final section describes a newer display device, the color liquid crystal display, which is likely to become as important in the decade to come as the color monitor is today.

22.3 OPERATIONAL CHARACTERISTICS OF COLOR MONITORS

The color television monitor is currently the most common display device for digital imaging applications, especially when temporally varying images are required. Its advantages include good temporal stability, a large color gamut, well-defined standards, and inexpensive manufacture. A wide variety of different CRT types is now available, but all types are derived from a common technological base. This section describes the characteristics shared by most color monitors, including design, controls, and operation. Of course, a chapter such as this can only skim the surface of such a highly evolved technology. There is a vast literature available for those who wish a greater depth of technical detail. Fink et al.¹ provides a great deal of it, with a useful set of references for those who wish to dig even deeper. Colorimetry is discussed in Chap. 10 in this volume.

Output devices that go under a variety of names, color television receivers, color video monitors, color computer displays, and so on, all use the same display component. This chapter describes only the display component. Technically, it is known as a color cathode ray tube (CRT), the term that is used exclusively in the remainder of the chapter.

In this chapter several operational characteristics are illustrated by measurements of CRT output. These measurements, which are intended only to be illustrative, were performed using a Tektronix SR690, a now obsolete CRT produced for the broadcast monitor market. While the measurements shown are typical of CRTs I have measured, they are intended to be neither predictors of CRT performance nor ideals to be pursued. In fact, color CRTs vary widely from model to model, and any CRT that is to be used in an application where colorimetry is critical should be carefully characterized before use.

Color CRT Design and Operation

The color CRT was derived from the monochrome CRT, and shares many features of its design. Thus, this section begins by describing the construction and operation of the monochrome CRT. New features that were incorporated to provide color are then discussed.

Monochrome CRTs A schematic diagram of a monochrome CRT is shown in Fig. 1. The envelope is a sealed glass tube from which all the air has been evacuated. Electrons are emitted from the cathode, which is heated red hot. The flux of electrons in the beam, the beam current I_b , is determined by a control grid. A variety of magnetic and/or electrostatic electrodes then focus, accelerate, and deflect the beam. The beam strikes a layer of phosphor on the inside of the CRT faceplate, depositing

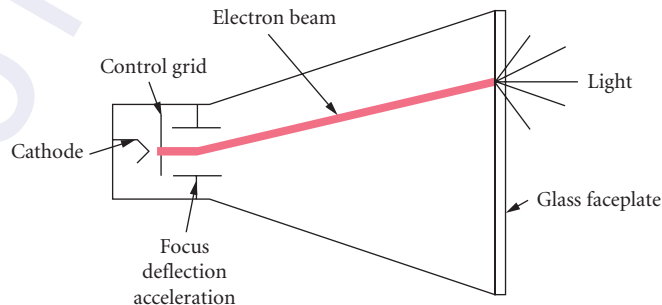


FIGURE 1 Schematic diagram of a monochrome CRT, showing the path of the electron beam and the location of the phosphor on the inside of the faceplate.

power $I_B V_A$, where V_A is the acceleration voltage. Some of this power is converted to light by the phosphor. The power Φ in the emitted light is given by

$$\Phi = \int \Phi_\lambda \cdot \frac{hc}{\lambda} d\lambda$$

Φ_λ , the spectral power distribution of the emitted light is determined, up to a single multiplicative factor, by the chemical composition of the phosphor. The multiplicative factor is usually taken to be a linear function of the beam power. Thus, the efficiency of the phosphor e_p given by

$$e_p = \Phi / I_B V_A$$

is independent of the beam current. Power not emitted as light becomes heat, with two consequences.

1. If the beam remains in the same spot on the screen long enough, the phosphor coating will heat up and boil the phosphor from the back of the faceplate leaving a hole in any image displayed on the screen.
2. Anything near the phosphor, such as the shadowmask in a color CRT, heats up.

A stationary spot on the screen produces an area of light. The intensity of the light is generally taken to have a gaussian spatial profile. That is, if the beam is centered at (x_0, y_0) , the spatial profile of the light emitted is given by

$$\Phi(x, y) \propto \exp \left[-\frac{1}{2\sigma_x^2}(x - x_0)^2 - \frac{1}{2\sigma_y^2}(y - y_0)^2 \right]$$

The dependence of this spatial profile on the beam current is the subject of active development in the CRT industry. In most applications the beam is scanned around the screen, making a pattern of illuminated areas. The brightness of a given illuminated area depends on the beam current when the electron beam irradiates the area, with the beam current determined by the voltage applied to the control grid.

Shadowmask Color CRTs The basic technology of the monochrome CRT is extended to produce color. The standard method for producing color is to take several, usually three, monochrome images that differ in color and mix them additively to form a gamut of colors. (Even very unconventional technologies, such as the Tektronix bichromatic liquid crystal shutter technology, produce color by additive mixture.) This section describes the basic principles of shadowmask CRT technology, which is the dominant technology in color video.

Geometry of the shadowmask CRT A color CRT must produce three images to give a full range of color. To do so usually requires a tube with three guns. The three electron beams are scanned in exactly the way that monochrome beams are, but arrive at the screen traveling in slightly different directions. Near the faceplate is a screen called the shadowmask. It is made of metal, with a regular pattern of holes. Electrons in the beams either hit the screen to be conducted away or pass ballistically through the holes. Because the three beams are traveling in different directions, they diverge once they have passed the hole, striking the back of the faceplate in different places. The phosphor on the back of the faceplate is not uniform, but is distributed in discrete areas that radiate red, green, or blue light. The geometry is arranged so that electrons from the red gun hit red-emitting phosphor, electrons from the green gun hit green-emitting phosphor, and electrons from the blue gun hit blue-emitting phosphor. This geometry is illustrated in Fig. 2. Several different geometries of shadowmask tube exist:

1. *Delta guns.* The three guns are arranged at the corners of an equilateral triangle, irradiating phosphor dots arranged in a triad.

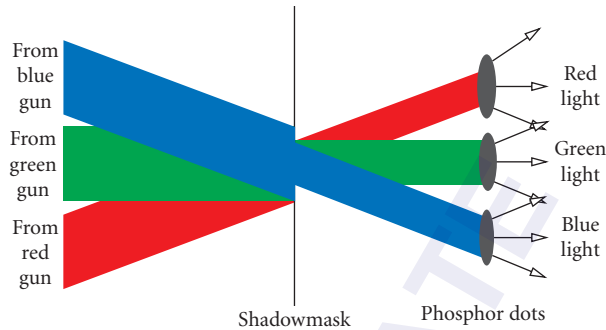


FIGURE 2 Electron beam/shadowmask/phosphor geometry in a shadowmask CRT.

2. *In-line guns.* The three guns are arranged in a line, irradiating phosphor dots side by side. The lines of phosphor dots are offset from line to line, so that the dot pattern is identical to that of the delta gun.
3. *Trinitron.* This is an in-line gun configuration, but the phosphor is arranged in vertical stripes. The holes in the shadowmask are rectangular, oriented vertically.

Other types of technology are also possible, though none is in widespread use. Beam index tubes, for example, dispense with the shadowmask, switching among the red, green, and blue signals as the electron beam falls on the different phosphor dots.

The most fundamental colorimetric property determining the colors produced by a shadowmask CRT is the light emitted by the phosphors. Only those colors that are the additive mixture of the phosphor colors (CRT primaries) can be produced. Because the color and efficiency of the phosphors are important determinants of perceived picture quality in broadcast applications phosphor chemistry undergoes continuous improvement, and varies from CRT to CRT. The emission spectra of the phosphors of a “typical” shadowmask CRT are shown in Fig. 3. The chromaticities of these phosphors are:

	x	y	z
Red phosphor	0.652	0.335	0.013
Green phosphor	0.298	0.604	0.098
Blue phosphor	0.149	0.064	0.787

Common problems in shadowmask CRTs The shadowmask is a weak point in color CRT design, so that color CRTs suffer from a variety of problems that potentially degrade their performance. The three most common problems are doming, blooming, and shadowmask magnetization. Doming occurs when the shadowmask heats because of the large energy density in the electrons it stops. It then heats up and expands, often not uniformly. The result is a distortion of its shape called “doming.” This geometrical distortion means that the registration of holes and dots is disturbed, and what should be uniform colors become nonuniform. Trinitron tubes have a tensioning wire to reduce this problem. It is often visible as a horizontal hairline running across the whole width of the tube about a third of the way from the top or bottom of the tube, depending on which side of the tube is installed up.

Blooming occurs when too much energy is deposited in the electron beam. Then electrons arrive at the screen in the right position but their entire kinetic energy is not immediately absorbed by the phosphor. They then can move laterally and deposit energy on nearby dots which can be both the wrong color and outside the intended boundary of the bright area. Colors become desaturated, and edges of areas become blurred.

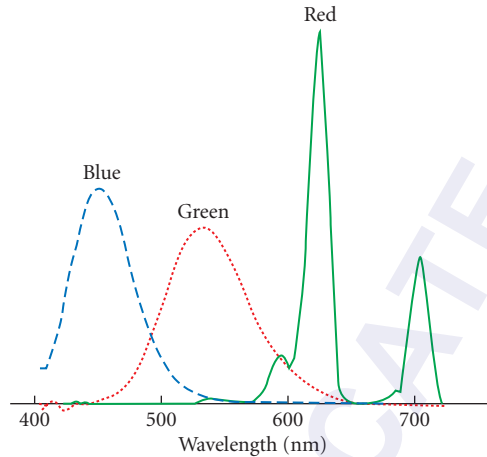


FIGURE 3 Spectral power distributions of the light output by a typical set of phosphors. The long wavelength peak of the red phosphor is sometimes missed in spectroradiometric measurements.

Magnetic fields build up in the shadowmask when the electromagnetic forces produced by electrons, assisted by the heat build-up, create magnetized areas in the shadowmask. These areas are nonuniform and usually produce large regions of nonuniformity in color or brightness on the screen. CRTs usually have automatic degaussing at power-up to remove these magnetic fields. If this is insufficient to remove the field buildup, inexpensive degaussing tools are available.

CRT Electronics and Controls The CRT receives voltage signals at its inputs, one for monochrome operation, three for color operation. The input signals must be suitably amplified to control the beam current using the grid electrode. The amplification process is described first, followed by the controls that are usually available to control the amplification.

Amplification Two properties of the amplification are important to the quality of the displayed image: bandwidth and gamma correction.

The bandwidth of the amplifiers determines the maximum frequency at which the beam current can be modulated. In most applications the beam is moved about the screen to produce the image, so that the maximum frequency translates into a maximum spatial frequency signal. Suppose, unrealistically, that the amplifiers have a sharp cutoff frequency v_{\max} and that the beam is scanned at velocity v_B . Then the maximum spatial frequency of which the CRT is capable is v_{\max}/v_B . In some types of CRT, particularly those designed for vector use, settling time is more important than bandwidth. It can be similarly related to beam velocity to determine the sharpness of spatial transitions in the image.

For colorimetric purposes, the relationship between the voltage applied at the CRT input and the light emitted from the screen is very important. It is determined by the amplification characteristics of the input amplifiers and of the CRT itself, and creating a good relationship is part of the art of CRT design. When the relationship must be known for colorimetric purposes, it is usually determined by measurement and handled tabularly. However, for explanatory purposes it is often written in the form

$$\Phi = \Phi_0 (V/V_0)^\gamma$$

where V is the voltage input to the CRT, normalized to its maximum value V_0 . The exponent, which is conventionally written as γ gives this amplification process its name, gamma correction.

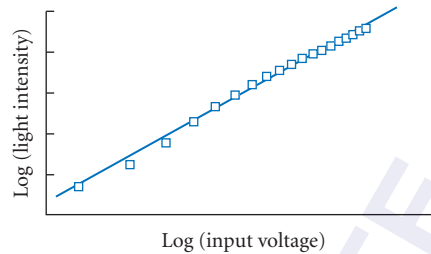


FIGURE 4 Graph of \ln (light intensity) against \ln (input voltage), sometimes used to determine the gamma correction exponent. Note the systematic deviation from a straight line.

Figure 4 shows some measured values, and a log-linear regression line drawn through them. Note the following features of this data.

1. The line is close to linear.
2. There are regions of input voltage near the top and bottom where significant deviations from linearity occur.
3. The total dynamic range (the ratio between the highest and lowest outputs) is roughly 100, a typical value. This quantity depends on the setting of brightness (black level) and contrast used during the measurement.

There will be additional discussion of gamma correction in the section on monitor setup.

Controls that affect amplification in monochrome CRTs Monochrome monitors have several controls that adjust several properties of the electron beam. Some controls are external to the monitor, some are internal, usually in the form of small potentiometers mounted on the circuit boards. The particular configuration depends on the monitor. In terms of decreasing likelihood of being external they are: brightness (black level), contrast, focus, underscan/overscan, pedestal, gain, horizontal and vertical size. A brief description of the action of each control follows. It is important, however, to remember that while names of controls tend to be constant from CRT to CRT, the action performed by each control often varies.

Brightness (black level) This control adjusts the background level of light on the monitor screen. It is designed for viewing conditions where ambient light is reflected from the monitor faceplate. It usually also varies the gamma exponent to a higher value when the background level (black level) is increased. A typical variation of the light intensity/input voltage relationship when brightness is varied is shown in Fig. 5.

Contrast This control varies the ratio between the intensity of the lightest possible value and the darkest possible value. High contrast is usually regarded as a desirable attribute of a displayed image, and the contrast control is usually used to produce the highest contrast that is consistent with a sharp image. A typical variation of the light intensity/input voltage relationship when this control is varied is shown in Fig. 6.

Focus This control varies the size of the electron beam. A more tightly focused electron beam produces sharper edges, but the beam can be focused too sharply, so that flat fields show artifactual spatial structure associated with beam motion. Focus is usually set with the beam size just large enough that no intensity minimum is visible between the raster lines on a uniform field.

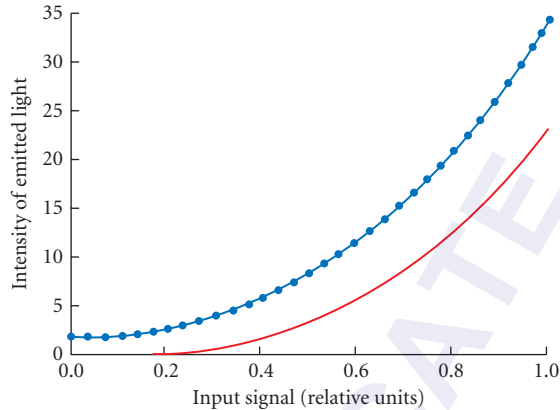


FIGURE 5 Variation of the light intensity/input voltage relationship when the brightness (black level) control is varied. The lower curve shows the relationship with brightness set near its minimum; the upper one with brightness set somewhat higher.

Pedestal and gain These controls, which are almost always internal, are similar to brightness and contrast, but are more directly connected to the actual amplifiers. Pedestal varies the level of light output when the input voltage is zero. Gain varies the rate at which the light output from the screen increases as the input voltage increases.

Controls specific to color CRTs Color monitors have a standard set of controls similar to those of monochrome monitors. Some of these, like brightness and contrast, have a single control applied simultaneously to each of the color components. Others, like gain and pedestal, have three controls, one for each of the color channels. There are several aspects of color that need to be controlled, however, and they are discussed in the paragraphs that follow.

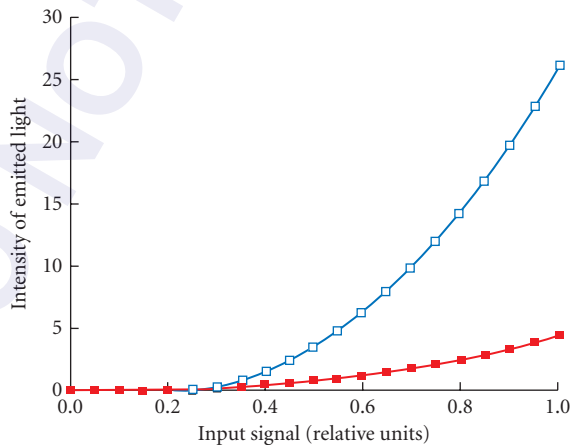


FIGURE 6 Variation of the light intensity/input voltage relationship when the contrast is varied. The lower curve shows the relationship with contrast set near its minimum; the upper one with contrast near its maximum.

Purity Purity is an effect associated with beam/shadowmask geometry. It describes the possibility that the electron beams can cause fluorescence in inappropriate phosphors. There is no standard set of controls for adjusting purity. Generally, there are magnets in the yoke area whose position can be adjusted to control purity, but this control is very difficult to perform. There will be no need to alter them under normal conditions. Purity is most influenced by stray magnetic fields, and can often be improved by moving the CRT.

White balance It is important for most monitor applications that when the red, green, and blue guns are turned on equally, calibration white (usually either D_{6500} or D_{9200}) appears on the screen. This should be true at all intensities. Thus, controls that alter the voltage input/light output relationship for each channel should be available. At a minimum, there will be the pedestal and gain for each channel.

Degauss Above we mentioned magnetic fields that build up in the shadowmask. There is generally a set of wires that run around the edge of the faceplate. At power-up a degaussing signal is sent through the wires and the ensuing magnetic field degausses the shadowmask. The degauss control can produce the degauss signal at any time.

CRT Operation The CRT forms an image on its screen by scanning the electron beam from place to place, modulating the beam current to change the brightness from one part of the image to another. A variety of scan patterns are possible, divided into two categories. Scan patterns that are determined by the context of the image (in which, for example, the beam moves following lines in the image) are used in vector displays, which were once very common but now less so. Scan patterns that cover the screen in a regular pattern independent of image content are used in raster displays: the scan pattern is called the raster. A variety of different rasters are possible; the one in most common use is a set of horizontal lines, drawn from top to bottom of the screen.

Raster generation Almost all raster CRTs respond to a standard input signal that creates a raster consisting of a set of horizontal lines. This section describes the path of the electron beam as it traverses the CRT screen; a later section discusses the signal configurations that provide the synchronization necessary to drive it.

Frames and fields The image on the screen is scanned out as a set of lines. Each line is scanned from left to right, as seen from a position facing the screen. Between the end of one line and the beginning of the next the beam returns very quickly to the left side of the screen. This is known as horizontal retrace or flyback. The successive lines are scanned from top to bottom of the screen. One field consists of a scan from top to bottom of the screen. Between the end of one field and the beginning of the next the beam returns very quickly to the top of the screen. This is known as vertical retrace or flyback.

One frame consists of a scan of all the lines in an image. In the simpler type of display a field is identical to a frame, and all the lines of the image are scanned out, from top to bottom of the display. The scan pattern is shown in Fig. 7. It is called noninterlaced.

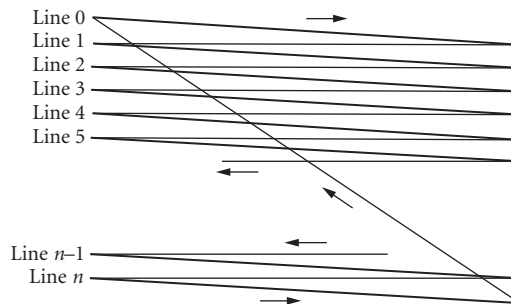


FIGURE 7 Scan pattern for a noninterlaced raster.

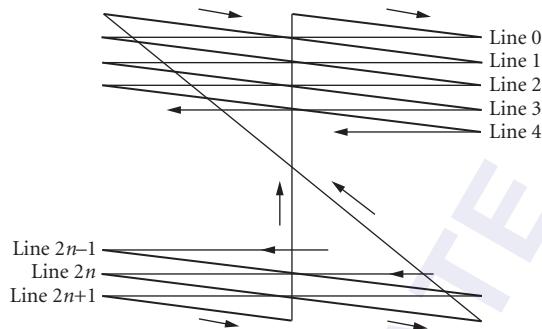


FIGURE 8 Scan pattern for an interlaced raster, two fields per frame. Line 3 eventually scans to line $2n - 1$ by way of the odd lines; line 4 scans to line $2n$.

A more complicated type of raster requires more than one field for each frame. The most usual case has two fields per frame, an even field and an odd field. During the even field, the even-numbered lines of the image are scanned out with spaces between them. This is followed by vertical retrace. Then, during the odd field the odd-numbered lines of the image are scanned out into the spaces left during the even-field scan. A second vertical retrace completes the frame. This scan pattern, known as interlaced, is shown in Fig. 8. The purpose of interlace is to decrease the visible flicker within small regions of the screen. It works well for this purpose, provided there are no high-contrast horizontal edges in the display. If they are present they appear to oscillate up and down at the frame rate, which is usually 30 Hz. This artifact can be very visible and objectionable, particularly to peripheral vision. Interlace having more than two fields per frame is possible, but uncommon.

In viewing Figs. 7 and 8 note that the vertical scan is produced by scanning the beam down continuously. Thus the visible lines are actually sloped down from left to right while the retrace, which is much faster than the horizontal scan, is virtually unsloped. The method used to put the odd field of the interlaced scan between the lines of the even field is to scan half a line at the end of the even field followed by half a line at the beginning of the odd field. Thus an interlaced raster has an odd number of lines.

Relationship of the raster to the CRT input signal The CRT receives a serial input signal containing a voltage that controls the intensity for each location in the image. This signal must be synchronized with the raster in order to make sure that each pixel is displayed in the right location. This section describes the relationship between the raster and the CRT input.

Horizontal scanning The input signal for one line is shown in Fig. 9. It consists of data interspersed with blank periods, called the horizontal blanking intervals. During the horizontal blanking interval the beam is stopped at the end of the line, scanned quickly back for the beginning of the next line, then accelerated before data for the next line begins. In the data portion, 0.0 V indicates black and 1.0 V indicates white. The signal shown would produce a line that is dim on the left, where the line starts, and bright on the right where it ends.

A second signal, the horizontal synchronization signal, has a negative-going pulse once per line, positioned during the horizontal blanking interval. This pulse is the signal for the CRT to begin the horizontal retrace.

When the synchronization signal is combined with the data signal the third waveform in Fig. 9 is produced. The intervals in the horizontal blanking interval before and after the synchronization are known as the front and back porch. They are commonly used to set the voltage level corresponding to the black level, a process known as black clamping. Black clamping reduces the effect of low-frequency noise on the color of the image, but requires good synchronization between the precise timing of the

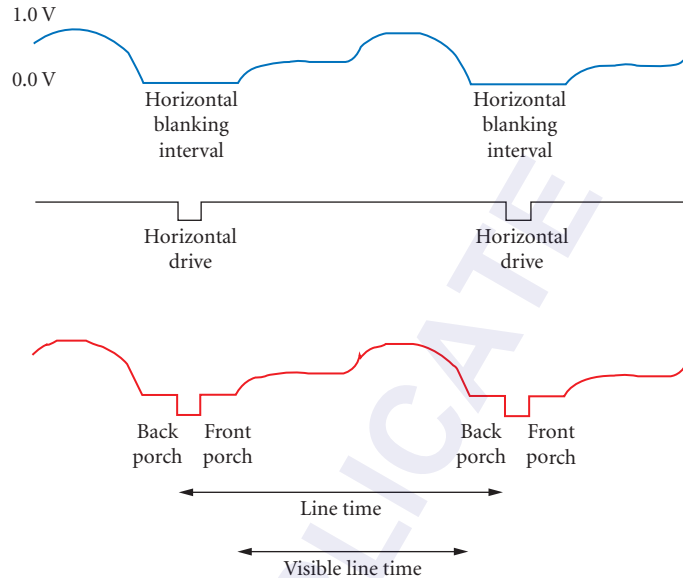


FIGURE 9 A schematic input signal that would generate a single line of raster, including the end of the preceding line and the beginning of the succeeding line. The top trace shows the signal that produces the image, including the blank between lines; the middle trace shows the synchronization signal with the horizontal drive pulse located during the blank; the bottom trace shows the synchronization signal and the picture signal in a single waveform.

signal and the raster production of the CRT. The ability of the monitor to hold synchronization with particular timing in this signal is an important parameter of the monitor electronics.

Vertical scanning Figure 10 shows the input signal for one complete field. There is data for each line, separated by the horizontal blanking period, which is too short to be visible in the figure. Separating the data portion of each field is the vertical blanking interval, during which the beam is scanned back to the top of the screen.

The synchronization signal consists of a vertical drive pulse signaling the time at which the beam should be scanned to the top of the screen. This pulse is long compared to the horizontal drive pulses so that it can easily be separated from them when the synchronization signals are combined into composite synch. In composite synch positive-going pulses in positions conjugate to the positions of the horizontal drive pulses are added during the vertical drive signal. These pulses are designed to keep the phase of the horizontal oscillator from wandering during vertical drive. When they are interpreted incorrectly, as was the case in some monitors early in the history of digital electronics, the result is a small shear in the first few lines at the top of the screen as the horizontal oscillator gets back into phase.

The composite synch signal can be added to the data, as shown in the fourth illustration in Fig. 10. Most real systems have this type of synchronization, which was designed for an era when signals were to be sent long distances by broadcast or wire. Today, in computer graphics applications, we often find that electronic circuitry in the source carefully folds the two synch signals into composite synch and the composite synch signal into the data signal. This signal is carried along a short piece of wire to the receiver, where electronic circuitry strips the signals apart. Bad synchronization is often caused by the malfunction of this superfluous circuitry, but it is unlikely that this situation will change in the immediate future.

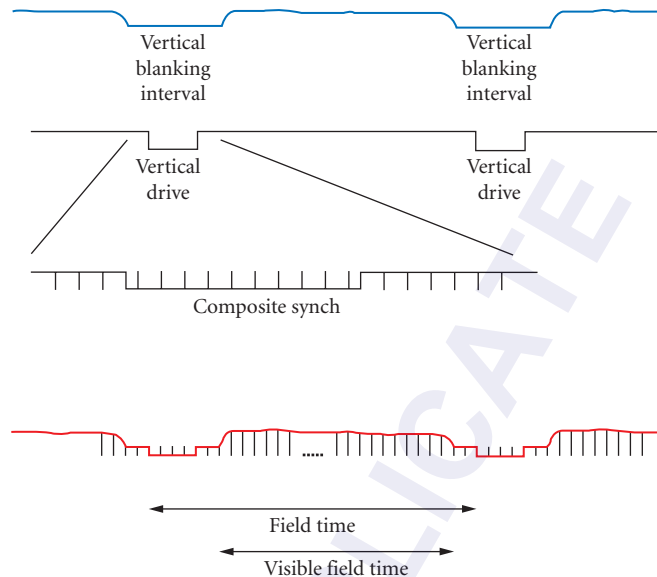


FIGURE 10 A schematic input signal that would generate a single field of a raster, including the end of the preceding field and the beginning of the succeeding field. The top trace shows the picture signal, with the vertical blank shown and the horizontal blank omitted. The second trace shows the vertical synchronization signal with the vertical drive pulse. The third trace shows the horizontal synchronization signal added to the vertical signal to give composite synch. The bottom trace shows composite synch added to the picture signal.

At the input to a color CRT, three or four signals like the ones described above are provided. The four signals are three input signals containing pixel intensities and blanking intervals, plus a fourth signal carrying the composite synchronization signal. To get three signals the synchronization signal is combined with one or more of the pixel signals. When only one pixel signal has synchronization information it is almost always the green one.

Controls that affect the raster Several controls affect the placement and size of the raster on the CRT. They are described in this section.

Horizontal/vertical size and position These controls provide continuous variation in the horizontal and vertical sizes of the raster, and in the position on the CRT faceplate where the origin of the raster is located.

Underscan/overscan Most CRTs provide two standard sizes of raster, as shown in Fig. 11. This control toggles between them. In the underscan position the image is smaller than the cabinet port, so that the whole image is visible, surrounded by a black border. In the overscan position the image is slightly larger than the cabinet port, so that no edges of the image are visible. There is a standard amount of overscan, which is used in home television receivers.

Convergence A shadowmask color CRT has, in effect, three rasters, one for each of the primary colors. It is essential that these rasters be positioned and sized exactly the same in all parts of the image; otherwise, spurious colored fringes appear at the edges in the image. To achieve this it is essential

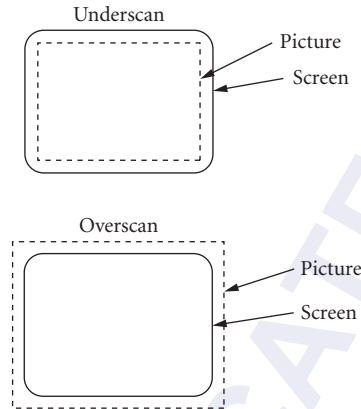


FIGURE 11 Schematic illustration of overscan and underscan. The rectangle with rounded corners is the cabinet mask that defines the viewing area; the rectangle with square corners is the displayed image.

that all three electron beams should be at the same place at the same time. For example, if the green gun is slightly to the left of the red and blue guns, a white line of a black background will have a green fringe on its left side and a magenta fringe on its right side. Good convergence is relatively easy to obtain with in-line gun configurations, so the usual practice is to adjust convergence at the factory using ring-shaped magnets placed over the yoke of the CRT, then to glue them permanently in place with epoxy cement. Satisfactory readjustment is very difficult to achieve.

Delta gun configurations have, by reputation, less stable convergence. Consequently, there are often controls available to the user for adjustment of convergence. The Tektronix 690SR was an extreme case, with a pullout drawer containing 52 potentiometers for controlling convergence on different areas of the screen. However, these controls, as might be expected, are not independent, so that “fine tuning” the convergence is very difficult, even when the controls are conveniently located.

Magnetic fields are generally the main culprit when convergence is objectionable, since the paths of electrons are curved in magnetic fields. Small fields from power transformers or from other electronic equipment or even the earth’s magnetic field can be the problem. Thus, before considering redoing the convergence on a new monitor, try rotating it and/or moving other equipment around in the room. Another useful trick: some low-cost CRTs have poorly positioned power transformers, and convergence can be improved by moving the power transformer far outside the CRT cabinet.

Operational Characteristics of Color CRTs

Two factors influence the light output that is produced in response to input signals. One is conventional: the input signal must be sequenced and timed precisely to match what is expected by the CRT. This problem is handled by input signal standards of long standing in the television industry. Synchronization standards are described first, followed by colorimetric standards. The second is technical: the construction of CRTs produces certain variations, even when the input signal is constant. The section ends with a discussion of this variation, specifying the temporal and spatial characteristics of the light output from a CRT and describing reasonable expectations for stability and uniformity in CRT-produced images.

Timing and Synchronization Standards Figures 9 and 10 show how the information needed to specify a CRT image are arranged in the input signal. For practical use they need to be specified numerically, thereby allowing CRT manufacturers to ensure that their products will respond appropriately to

the input signals they will encounter in service. Two standards have been created by the Electronic Industries Association (EIA) for this purpose. They are actually standards for television studios, prescribing how signals should be distributed in closed-circuit applications, and specifically between television cameras and studio monitors. In practice they are more widely used. Among other applications, they specify the interface between the tuner and the CRT in television receivers, and the output of digital graphics systems allowing studio monitors to be used as display devices.

The two standards are RS-170,² which was designed for lower bandwidth applications, and RS-343,³ which was designed for higher bandwidth applications. Each gives minimum and maximum timings for each part of the input signal in terms of several parameters that are allowed to vary from application to application. These parameters are then part of CRT input specification, allowing all timing characteristics to be deduced. The most important parameter is the line rate, the number of lines displayed per second. Most CRTs are intolerant of large variations in this parameter. Another important parameter is the field rate, the number of fields displayed per second. Older CRTs were quite intolerant of variations in this parameter, but most modern CRTs can handle signals with a very wide variety of field rates.

RS-170 and RS-343 are monochrome standards. When used for color they are tripled, with the input signal assumed to be in three parallel monochrome signals. As mentioned above, the synchronization signal is either placed on a fourth, or incorporated into a single color signal, usually the green one. However, although this practice is almost universal there is no official color standard for the RGB input to a monitor. This lack can have unfortunate consequences. For example, the NTSC color signal naturally decodes into three signals with peak-to-peak voltages of 0.7 V. Thus RGB monitors were built with inputs expecting this range. RS-170 and RS-343, on the other hand, specify a peak-to-peak voltage of 1.0 V. Early digital graphics systems were built to provide exact RS-170 and RS-343 output. These systems, naturally, overdrove standard monitors badly.

Colorimetric Standards In broadcast applications the image transmitter should be able to specify the precise color that will be displayed on the receiver's CRT. This requirement can be supplied by a colorimetric standard. The NTSC color standard was agreed upon for use in the North American broadcast television industry. It is a complete color standard, specifying phosphor chromaticities, color representation on the carrier signal, signal bandwidth, gamma correction, color balance, and so on. Thus, if the NTSC standard were followed in both transmitter and receiver, home television would provide calibrated colors. It is not followed exactly, however, since both television manufacturers and broadcasters have discovered that there are color distortions that viewers prefer to colorimetrically precise color. Furthermore, it is not useful for high-quality imagery since the low bandwidth it allocates to the chromatic channels produces edges that are incompatible with good image quality.

Spatial and Temporal Characteristics of Emitted Light The light emitted from a CRT is not precisely uniform, but suffers from small- and large-scale variations. The small-scale variations arise because the image is actually created by a spot that moves over the entire display surface in a time that is intended to be short compared to temporal integration times in the human visual system. In fact, a short enough glimpse of the CRT screen reveals only a single point of light; a longer one reveals a line as the point moves during the viewing time. These patterns are designed to be relatively invisible under normal viewing conditions, but often need to be considered when CRTs are used for vision experimentation, or when radiometric measurements are made. They are controlled largely by the input signal, so they are relatively constant from CRT to CRT.

The small-scale variations are described in this section; the large-scale variations, which occur as spatial nonuniformity and temporal instability in the emitted light, are discussed in the following section.

Spatial characteristics The electron beam scans the screen horizontally, making an image that consists of a set of horizontal lines. During the scan its intensity is modulated, so that the line varies in brightness—vertical edges, for example, being created by coordinated modulation in a series of horizontal lines. The sharpness of such a line depends on two factors: the ability of the video amplifiers to produce an abrupt change in intensity and the size of the spot of light on the screen, which is essentially the same as the cross section of the electron beam. Video amplifiers vary greatly in bandwidth

from one model of CRT to another, and adjusting their performance is beyond the reach of virtually all CRT users. Unfortunately, interactions between adjacent horizontal pixels are not even linear,⁴⁻⁶ and the measurements needed to determine the nonlinearity of a given CRT are extremely demanding. Thus, compensating for amplifier bandwidth or even measuring its effects is beyond the scope of this chapter.

By contrast, the sharpness of a horizontal line is independent of the video amplifiers, but depends only on the spot size and its spatial profile, which is well modeled as a two-dimensional gaussian. The width of the gaussian depends on the focus electrodes, and is user-adjustable on most CRTs. The shrinking raster technique for CRT setup (discussed in subsection “CRT Setup for Image Display” discussed later in this chapter) determines a spot size that has a specific relationship to the interline spacing. Assuming this setup, it is possible to make reasonable assumptions about the contrast of images on a CRT,⁷ and these assumptions can be extended to the small-scale spatial structure of arbitrary images. Note that many CRTs used primarily as display terminals are overfocused compared to the shrinking raster criterion because such overfocusing allows the display of smaller text. Overfocusing is easily detected as visible raster lines, usually in the form of closely spaced dark horizontal lines when a uniform field is displayed.

Temporal characteristics Because of the scan pattern the light emitted from any portion of the screen has a complicated temporal dependence. The next few paragraphs describe several levels of this dependence, assuming, for simplicity, that the scan is noninterlaced. Similar results for an interlaced scan are easy to derive. The relevant variables are

τ_d	phosphor decay time
τ_p	time the beam spends crossing a single pixel
v_h	horizontal velocity of the beam
x_p	horizontal interpixel spacing, $x_p = v_h \tau_p$
τ_l	time spent scanning a complete line, including horizontal flyback
v_v	vertical velocity of the beam
y_p	vertical interpixel (interline) spacing, $y_p = v_v \tau_l$
τ_f	time spent scanning a complete frame (identical to the time spent scanning a complete field)

These temporal factors usually change the colorimetry of the CRT. When one or more is changed, usually because the video source has been reconfigured, the color output for a given input to the CRT usually changes. Thus, a recalibration should be done after any such change. And it is probably practical to recommend that the system be used with the new video timing parameters for a while before the recalibration is done, since tuning or touch-ups will require further recalibration. The intensity of the light emitted from a vanishingly small area of the CRT screen is a function that is zero until the beam traverses the point t_s , then decays exponentially afterward,

$$\Phi(t) = \sum \Phi_0 \theta(t - t_s) \exp(-(t - t_s) / \tau_d)$$

τ_d ranges between $10^{-3} \tau_l$ and τ_f , and usually varies from one phosphor to another. Most often the green phosphor has the largest τ_d , the blue one the smallest. Broadcast monitors tend to have small τ_d 's since they are designed to display moving imagery; data display CRTs tend to have large τ_d 's since they are intended to display static imagery. Occasionally a CRT primary (most often red) is a mixture of two phosphors with different decay times. In such cases, the chromaticity of the light emitted by the primary changes over time, though the change is not usually visible.

If a second pixel is n_h pixels to the right of a given pixel and n_v lines below it, the light emitted by the second pixel lags behind the light emitted by the first pixel by

$$n_h \tau_p + n_v \tau_l \approx d_h / v_h + d_v / v_v$$

where d_h and d_v are the horizontal and vertical screen distances.

Commonly, a detector measures a finite area of the CRT screen. Then intensity of the detected light is a closely spaced group of exponentials followed by a long gap, then another closely spaced group, and so on. Each peak in the closely spaced group is spaced about τ_p from the previous one, occurring when each line of the raster crosses the detection area. Each group is spaced about τ_f from the previous one, occurring each time a new field repaints the detected target. The time constant of this composite signal is τ_d , the decay time of the phosphor for each individual peak in the signal.

Stability and Uniformity of CRT Output The small-scale structure of the emitted light, discussed above, determines the specific form of visual stimuli produced by a CRT. The large-scale structure determines the scope of applicability of measurements made on a particular area of a CRT at a particular time.

Temporal stability How constant is the light emitted by a CRT that receives the same input signal? Figures 12 and 13 show the results of colorimetric measurements taken over a 24-hour period. In each figure the calibration bars shows a 2 percent variation compared to the average value. The variation decreases considerably in the latter part of the graph. It is the period from 5 p.m. until 9 the next morning, showing that variations in the building power are the largest source of temporal instability. Furthermore, this variation affected all guns similarly, as is shown by the smaller variation in Fig. 13, which shows the chromaticity coordinates, thereby factoring out overall changes in intensity.

Measurements like this one are important for establishing the precision at which calibration is sensible: it isn't worthwhile to calibrate to a greater precision than the variation in the colorimetry over the period between calibrations. While this CRT has excellent temporal stability, significantly better than most other light sources, the same is not true for all CRTs. Some vary in output by as much as 20 to 30 percent over times as short as a few minutes.

The measurements shown in Figs. 12 and 13 give an idea of the variation of light output over periods of hours or days, and illustrate the precision it is possible to expect from a recently calibrated CRT. CRTs also vary on a time scale of years, but the effects are not documented. Anecdotal reports suggest the following. First, electronic components change with age, so properties that depend on the

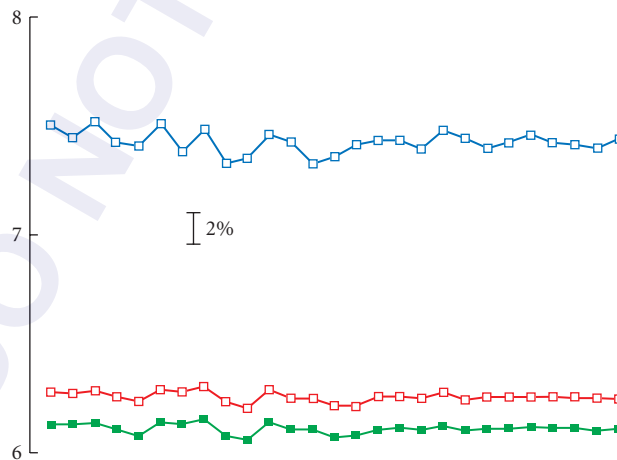


FIGURE 12 Variation of light output from a color CRT over 24 hours of continuous operation. This graph shows the three tristimulus values when a neutral color is displayed. The latter part of the graph is at night when almost all other equipment in the building is turned off.

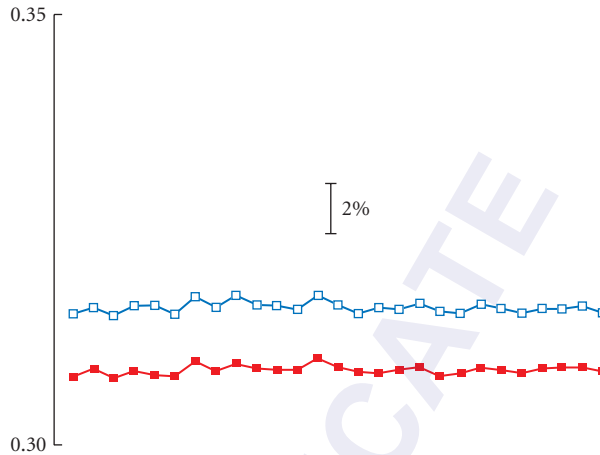


FIGURE 13 Variation of light output from a color CRT over 24 hours of continuous operation. This graph shows the chromaticity coordinates corresponding to Fig. 12. They show less variation than the tristimulus values which covary considerably.

CRT electronics, such as video amplifier gain and bandwidth, change with age, almost always for the worse. Second, chemical properties do not change with age, so properties such as the spectral power of light emitted by a specific phosphor do not change. One anecdotal report⁸ describes no variation in chromaticity coordinates of phosphor emission spectra over several years of CRT operation. It is possible, however, that phosphor concentrations diminish as tubes age, probably because the phosphor slowly evaporates. Such an effect would reduce the intensity of light emitted from the phosphor without changing its chromaticity. The magnitude of this effect is controversial.

Spatial uniformity The light emitted by a specific input signal varies a surprising amount from one area of the screen to another. Figure 14 shows the variation of luminance at constant input voltage as

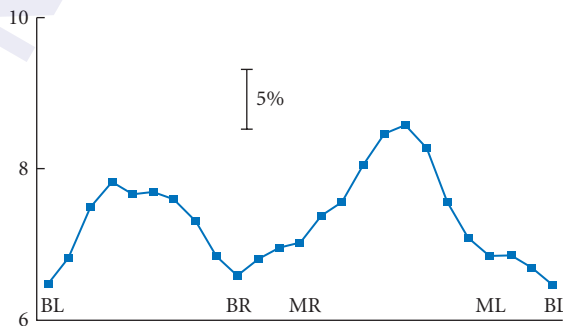


FIGURE 14 Variation of light output when different parts of a CRT screen are measured from a fixed point. Horizontal lines mark variations of about 5 percent.

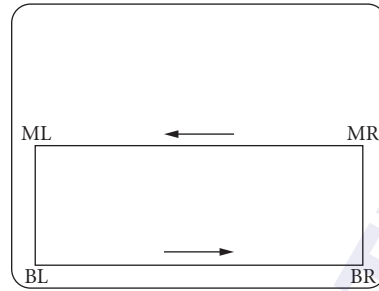


FIGURE 15 The measurement path used for the measurements shown in Fig. 14.

we measure different areas of the screen from a fixed measurement point. (Figure 15 shows the location on the screen of the measurement path.) Note that the light intensity decreases as the measurement point moves away from the center either horizontally or vertically, and is lowest in the corners. Two effects work together to create this effect. As the beam scans away from the center of the tube, it meets the shadowmask at more and more oblique angles, making the holes effectively smaller. Because of the curvature of the tube and the finite distance of the observer, the edges and corners are viewed at angles off the normal to the tube. Light is emitted in a non-Lambertian distribution, preferring directions closer to the normal to the tube face. The effects in Fig. 14 occur in all CRTs. How large they are, however, depends strongly on the type and setup of the monitor. Correcting for this nonuniformity is usually impractical. Doing so requires very extensive measurement.⁹

Closer examination of measured light shows, however, that many experimental situations are not hampered by this nonuniformity. Usually the chromaticity coordinates are very close to being constant, even though the luminance varies greatly. General intuition about color, as well as some recent experiments,¹⁰ shows that humans are quite insensitive to smooth luminance gradients, even when they are as large as 20 percent. This fact, combined with commonsense layout of experimental displays (making them symmetrical with respect to the center of the screen, for example) overcomes spatial nonuniformity without extensive measurement. Important to mention in this respect is the difficulty of creating good stereoscopic viewing conditions on a single monitor. First, there is only a single area of the screen that is ideal for the position of the center of an image. Second, unless one image is horizontally inverted, two images from opposite sides of the screen combined into a single image present drastically different luminance gradients to the two eyes.

Setup and Viewing Environments for Color CRTs

Many adjustments of the CRT electronics change its performance substantially. Some, such as the purity and convergence, have specific “correct” settings and are not designed to be user-adjustable. Most CRTs with in-line guns, for example, have the ring magnets that adjust the convergence glued into place at the factory. Other controls, such as contrast and brightness, are user-adjustable and should be changed when viewing conditions change if an optimal image is to be produced. These controls are provided in the expectation that the CRT will be displaying an image that is broadcast to a large number of CRTs that are viewed in very different visual environments. There is, therefore, a correct way to adjust these controls; the next few paragraphs describe the basic adjustments to be done. The procedure is based on a technical manual produced by the Canadian Broadcasting Corporation.¹¹ This manual also provides recommended viewing conditions for critical assessment of displayed images. The procedures and viewing conditions are expected to form the basis of an SMPTE standard for CRT image display.

When CRTs are used for visual experimentation, of course, they are often displaying specifically controlled images in unusual viewing conditions, often total darkness. For such applications the

adjustment procedure described below is unlikely to be interesting, and extreme values of the controls are likely to be desired. For example, an experiment conducted in total darkness is likely to need the black level (brightness) set so that there is no background light emitted from the screen. Or an experiment to measure thresholds is likely to benefit from whatever value of the contrast control minimizes the gain of the CRT at the intensity levels where the measurement is performed. (The objective is to minimize intensity quantization of the stimulus.) In fact, modern CRTs with computer-controllable contrast and black level might be used with different control settings for different trials of a single experiment.

CRT Setup for Image Display When a CRT is set up for image display, four adjustments are usually made: focus, brightness, contrast, and color balance and tracking. These adjustments form a rough order, with changes to one often requiring changes to succeeding ones. The following procedures are simplified from Benedikt.¹¹

Focus Focus should be adjusted by the shrinking raster method. A uniform gray field is displayed on the CRT, and the focus adjustment used to shrink the beam size until raster lines are clearly visible. The beam size is then increased until the raster lines just barely disappear. The smallest beam size for which no horizontal raster lines are visible indicates the correct setting for the focus control.

Brightness The brightness or black-level control is set so that zero input to the CRT produces a visual impression of black. This setting *must* be performed in lighting conditions that are exactly those in which the images will be viewed, and with the observer at exactly the viewing distance at which the images will be viewed. With no signal input to the CRT, and the image set to underscan if possible, the brightness control is increased until the image area is noticeably gray. It is then reduced to the highest setting at which the image area looks black.

Contrast An input that has all three guns fully on is used for this adjustment. With such a color displayed, the contrast control is adjusted until the luminance of the screen is the maximum luminance desired. This setting should be performed using either a luminance meter or a color comparator, a CRT adjustment device that allows an observer to view the CRT as half of a bipartite field, the other half containing a reference white at the correct luminance. In any case, it is essential that the CRT should not bloom at the contrast setting in use. Blooming occurs when a too-intense electron beam spreads after it has passed through the shadowmask, stimulating more than one phosphor. It reduces the purity of the image, with the visual consequence that bright saturated colors are washed out. Narrow lines of red, green, and blue at full intensity can be used to check for blooming.

Color balance and tracking A white with red, green, and blue inputs equal at maximum intensity and a gray with red, green, and blue inputs equal at half intensity are used to set the color balance. In addition, colorimetric capability is needed. Visual colorimetry using a luminous white reference is usual in the broadcast industry, but this measurement can also be made instrumentally. For visual colorimetry the white reference is usually supplied by a color comparator. Adjust the red, green, and blue gain controls until the displayed white matches the reference in chromaticity. Then adjust the red, green, and blue screen controls until the displayed gray matches the reference in chromaticity. It may now be necessary to readjust the brightness and contrast. Do so, then repeat the color balance and tracking adjustment until all the three adjustments are simultaneously satisfactory.

Viewing Environments In applications where image quality is critical, it is necessary to control the viewing environment very closely. The following viewing conditions are typical of those used in the broadcast television industry.

1. The luminance of reference white is about 70 cd/m².
2. The observer views the screen from a direction normal to the screen, with the screen-observer distance between 4 and 6 times the screen height.

3. The CRT should be surrounded by a neutral matte area at least 8 times the screen area. The surround should have the chromaticity of the CRT reference white, and a luminance of about 10 cd/m².
4. A narrow matte-black mask should frame the CRT image.
5. All room lighting should be as close as possible to the chromaticity of the CRT reference white.

22.4 COLORIMETRIC CALIBRATION OF VIDEO MONITORS

When users of color CRTs talk precisely about colorimetric calibration, a careful distinction between calibration and characterization is usually made. Measuring the input voltage/output color relationship well enough that it is possible to predict the output color from the input voltage, or to discover the input voltage needed to produce a given output color is a CRT characterization. The function mapping voltage to color characterizes the colorimetric performance of the CRT. Adjusting the CRT so that its characterization function matches the characterization function of a standard CRT is a CRT calibration. Talking loosely, however, calibration usually covers both characterization and calibration, and most CRT users are actually more interested in characterizations than in calibrations. Thus, this chapter describes several methods for creating colorimetric characterizations of CRTs. Some such characterization is an essential part of any calibration, but omits the detailed internal adjustment of a CRT, which requires significant electronic expertise if it is to be done safely.

There is no single characterization method that suits all needs. Thus, this chapter provides a variety of methods, each of which does some jobs well and others badly. They can be divided into three basic types.

1. *Exhaustive characterization methods (ECM)* Useful when good precision is needed over the complete monitor gamut. The same or similar algorithms can be used for the characterization of other output devices, such as printers. These methods tend to be both computationally and radiometrically expensive.
2. *Local characterization methods (LCM)* Useful when a precise characterization is needed for only a small portion of the monitor's output range.
3. *Model-dependent characterization methods (MDCM)* Useful when a characterization of moderate precision is needed for the complete monitor gamut. These methods tend to be specific to a given monitor and useful for only a small set of monitor setup methods, but they are computationally and radiometrically inexpensive. In addition they can be done so that the perceptual effects of mischaracterizations remain small, even when the colorimetric effects are large.

Here is a small set of criteria for deciding which type of method is best for a specific application:

1. High precision (better than 1 to 3 percent is impractical by any method discussed in this chapter)—ECM or LCM
2. Complete gamut characterization needed—ECM or MDCM
3. Minimal memory available for characterization tables—LCM or MDCM
4. No or slow floating point available—ECM or LCM
5. Fast inverse needed (we call the transformation from *RGB* to *XYZ* the characterization, the transformation from *XYZ* to *RGB* the inverse)—ECM, LCM, or MDCM
6. Forgiving behavior with out-of-gamut colors—LCM or MDCM
7. Photometry available but not radiometry—MDCM
8. Change of small number of parameters when monitor or monitor setup changes—LCM or MDCM

This list is not comprehensive; each type of characterization is actually a family of methods and there may be an effective method of a particular type even when general considerations seem to rule out that type. Undoubtedly more methods will be discovered as time goes on, leading to a broadening of these rules, and perhaps even a useful splitting of the characterization types into subtypes.

Important considerations when deciding on the characterization method to be used are measurement precision and monitor stability. There is no point in using a characterization method more precise than measurement precision, or more precise than the stability of the CRT over the range of conditions in which the characterization is to be used.

The characterization methods described below make use of a variety of measurements, but measurement methodology is not discussed. For reference, the methods used are (1) spectroradiometry, measurement of spectral power distributions; (2) colorimetry, measurement of tristimulus values; and (3) photometry, measurement of luminance.

The phrase “monitor coordinates” is used throughout this section. It denotes a set of controllable CRT inputs to which the experimenter has immediate access. For example, in applications where the CRT is computer controlled, it indicates the *RGB* values in the color lookup tables of the frame buffer that supplies the video signal to the CRT. The characterization techniques discussed below are independent of the monitor coordinates used. The single exception is the set of model-dependent characterization methods.

Exhaustive Characterization Methods

These are the most costly characterization methods in terms of time and computer storage. They are also the most precise and the most general. Their utility depends only on the stability of the monitor. (If the monitor doesn't have enough stability for exhaustive methods to be usable it is unlikely to have enough stability to use any other method.)

General Description The idea behind exhaustive characterization is to measure all the colors a monitor can produce and store them in a large table. When a color of given monitor coordinates (*RGB*) is displayed, its tristimulus values are determined by looking them up in the table. When a color of given tristimulus coordinates is desired, the table is searched and that set of monitor coordinates closest to the desired value is chosen for display. (Dithering among nearby colors is also possible if more precision is desired.) This method is also useful for any kind of device, not just for monitors. Thus, software designed for monitors can be reused when, for example, printers must be characterized. The obvious drawback to this method is the number of measurements that must be made. For example, with 24-bit color (8 bits per gun), over 16 million colors can be produced. If each measurement takes 1 s, the measurement process consumes over 4500 h, or almost 200 days, measuring around the clock. The solution to this problem is to sample the set of realizable colors, measure the samples, then interpolate between them. Thus, what we call exhaustive methods are only relatively exhaustive. Handling the practical problems is discussed in the following sections.

Sampling Densities and Interpolation Algorithms How many measurements should be made, and which colors should be selected for measurement? The answer depends on the nature of the sampling algorithm to be used, and is an open research question for general sampling. Thus, while it is possible to discuss the issues, practical decisions depend on the experience of the user.

Sampling procedures Most sampling algorithms sample the monitor coordinates linearly in the monitor coordinates. For example, if $512 = 2^9$ samples are to be taken for a color CRT, $8 = 2^3$ values on each of the red, green, and blue guns would be used, linearly interpolating the range of output voltages. The 512 samples would be the cartesian product of each set of gun values. Thus, if there were 256 possible values for each gun, running from 0 to 255, the 8 values chosen would be 0, 36, 72, 109, 145, 182, 218, and 255. If, as is often the case, saturation exists at the low and/or high ends of the input range, the full range is not used, and these values are scaled to the usable range. Behind this choice is the idea that a well-designed output device should have device coordinates that are close

to perceptually uniform. If so, linear sampling in device coordinates roughly approximates even sampling in a perceptually uniform space. The density of the sampling is related to the interpolation algorithm, which must approximate the exact value to within the desired precision of the characterization. This necessarily entails a trade-off between measurement time and online computational complexity when the characterization is being used: the more dense the sampling, the simpler the interpolation algorithm, and vice versa.

Interpolation Most often, linear interpolation is used. In this case, since we are interpolating in a cubic volume the interpolation is trilinear. Thus, given R_j (RGB), what are the tristimulus values?

1. We assume the existence of a table $X_i(R_j(n_j))$, consisting of the three tristimulus values X_i measured when the three guns are driven by the three voltages $R_j(n_j)$. Each three-vector n_j labels a different sample.
2. Find the sample m_j that has $R_j(m_j)$ less than and closest to the values to be estimated R_j .
3. Take the eight samples (m_0, m_1, m_2) , $(m_0 + 1, m_1, m_2)$, $(m_0, m_1 + 1, m_2)$, \dots , and $(m_0 + 1, m_1 + 1, m_2 + 1)$ as the vertices of a polyhedron.
4. Interpolate in R_0 on the four sides running from (m_0, m_1, m_2) to $(m_0 + 1, m_1, m_2)$, from $(m_0, m_1 + 1, m_2)$ to $(m_0 + 1, m_1 + 1, m_2)$, from $(m_0, m_1, m_2 + 1)$ to $(m_0 + 1, m_1, m_2 + 1)$, and from $(m_0, m_1 + 1, m_2 + 1)$ to $(m_0 + 1, m_1 + 1, m_2 + 1)$. The interpolation algorithm for $X_i(*, l_1, l_2)$ is given by

$$X_i(*, l_1, l_2) = \frac{R_0 - R_0(m_0)}{R_0(m_0 + 1) - R_0(m_0)} X_i(R_0(m_0 + 1), R_1(l_1), R_2(l_2)) \\ + \frac{R_0(m_0 + 1) - R_0}{R_0(m_0 + 1) - R_0(m_0)} X_i(R_0(m_0), R_1(l_1), R_2(l_2))$$

where l_j is either m_j or $m_j + 1$.

5. Treat the four values as the corners of a polygon. Interpolate in R_1 along the two sides running from (m_1, m_2) to $(m_1 + 1, m_2)$ and from $(m_1, m_2 + 1)$ to $(m_1 + 1, m_2 + 1)$. The interpolation algorithm for $X_i(*, *, l_2)$ is given by

$$X_i(*, *, l_2) = \frac{R_1 - R_1(m_1)}{R_1(m_1 + 1) - R_1(m_1)} X_i(*, R_1(m_1 + 1), R_2(l_2)) \\ + \frac{R_1(m_1 + 1) - R_1}{R_1(m_1 + 1) - R_1(m_1)} X_i(*, R_1(m_1), R_2(l_2))$$

6. Treat these two values as the endpoints of a line segment, and interpolate in R_2 . The final value is given by

$$X_i = \frac{R_2 - R_2(m_2)}{R_2(m_2 + 1) - R_2(m_2)} X_i(*, *, R_2(m_2 + 1)) \\ + \frac{R_2(m_2 + 1) - R_2}{R_2(m_2 + 1) - R_2(m_2)} X_i(*, *, R_2(m_2))$$

The above equations implement trilinear interpolation within a (possibly distorted) cube. It has recently been observed that tetrahedral interpolation has some desirable properties that are lacking in trilinear interpolation. It is accomplished by subdividing the cube into five or six tetrahedra, the corners of which coincide with the corners of the cube. Barycentric coordinates are then used to determine the tetrahedron in which the displayed color lies, and to interpolate within that

tetrahedron. More complex interpolation methods are also possible for nonuniform measurement sampling.¹²

Such interpolation methods implement the characterization directly in terms of the measured values. To test the adequacy of the interpolation it is important to choose regions of color space where curvature of the input/output response is expected to be high and to test the interpolation against exact measurements in that region of color space. If the interpolated values do not match the measured values there are two possible solutions: increase the measurement sampling density or improve the interpolation function. The first is expensive in terms of measurement time and computer memory; the second in terms of online calculation time. (There is a middle way, where a complicated algorithm is used to interpolate among coarsely sampled measurements for the purpose of generating a table large enough that linear interpolation can be done online. This has been used for hardcopy devices but not for monitors.) If the interpolation function is to be generalized it is possible to use higher-order interpolating functions, like splines, or to linearize the measurement space by transforming it before interpolating. For example, it is possible to build a characterization function that provides the logarithms of tristimulus values in terms of the logarithms of input coordinates. Any reasonable powerful generalization is bound to be computationally expensive. Much remains to be done to improve the sampling and interpolation process.

Inverses Calculating inverses for multidimensional tabular functions is not straightforward. Speed of computation can be optimized by creating an inverse table, giving values of R_j for regularly sampled values of X_j . The easiest way to construct this table is not by measurement, but by calculation (off-line) from the characterization table. Once the table is available, linear interpolation can be done online to provide voltages corresponding to given tristimulus values. Two possible methods exist for calculating inverses from tables.

1. Newton's method of finding zeros using derivatives to calculate an error function is possible, but can have bad numerical instabilities on tabular data that contains measurement error.
2. Examine the table that drives the forward transform to find the cell that contains the desired tristimulus values. Then subdivide the cell, using trilinear interpolation to determine new values for the corners. Determine which subcell contains the desired tristimulus values, and continue subdividing until the solution is sufficiently precise. This method is robust, but too expensive computationally to use online.

A promising alternative, not yet in wide use, takes advantage of barycentric coordinates in a tetrahedral lattice to find the appropriate cell quickly.

Nonlinear interpolation schemes can be used in the forward transformation when producing tables for linear inverse mappings. Presumably tables for nonlinear interpolation schemes in the inverse mapping can also be produced.

Three issues about inverses stand out.

1. Nonlinear schemes have not been sufficiently explored for their potential to be clear.
2. The straightforward application of well-understood computational methodology can provide substantial improvements over present methods.
3. Whatever inversion method is employed, out-of-gamut colors remain a serious unsolved problem.

Out-of-Gamut Colors What should be done when a program is asked to display a set of tristimulus values that are outside the color gamut of the monitor on which the color is to be displayed? The answer, of course, depends on the nature of the application. If the application demands drastic action when a request for an out-of-gamut color occurs, the solution is easy. The interpolation algorithm returns an illegal value of R_j . Then the application can display an error color in that area, or exit with an error. If, instead, the application must display a reasonable color within the monitor gamut, there is no natural solution. Solutions like projecting onto the surface of the monitor gamut have been used with success for some applications. Unfortunately, however, they are frequently computationally expensive and visually unsatisfactory.

Local Characterization Methods

In many applications only a small region of the color gamut of the monitor must be characterized, but that region must be characterized very precisely. A typical example is a threshold experiment in which there are a small number of reference stimuli. Only colors close to the reference stimuli need be characterized precisely since they are the only colors likely to arise in the experiment. Methods that are specialized for the characterization of small parts of the color gamut, such as the threshold stimulus set, are the subject of this section. They are the methods most often appropriate for vision experimentation.

General Description Local characterization methods try to take advantage of simplifications that arise because the set of colors to be characterized is small. For example, a small set of colors all very close to a reference color can be characterized by a linear approximation to the global characterization function. Such simplifications offer many desirable properties such as high precision with minimal measurement (individual colors), linear characterizations and inverses (local regions of color), and simple inverses over extended color ranges (one-dimensional color spaces). To realize this precision, colorimetric measurements of the required precision must be available and easily usable. Detailed descriptions of three limited-gamut characterization schemes follow. Others are certainly possible, and may be worked out by analogy.

Individual colors It is often the case that a small number of distinct colors is needed. Under such circumstances the best method available is to perform a colorimetric measurement on each color. This is easy if the colors are arbitrary but need to be known precisely. Then they can be chosen by their *RGB* coordinates, with colorimetric measurements used to establish their tristimulus values. It is less easy if a small number of colors of specified tristimulus values are needed. The best way to deal with the latter problem is to use online colorimetric measurement, adjusting *RGB* values until the required tristimulus values are produced. This should be repeated several times throughout the time during which the characterization is needed to establish the stability of the monitor while the characterization is in use.

This type of characterization leads most directly into the most difficult question that CRT characterization must face: how should a characterization handle spatial and angular variations of light emitted by the CRT? Suppose, for example, an experiment uses only two color stimuli, a red one and a green one. Suppose further that the red one always appears at one point on the screen, while the green one always appears at a different point. Clearly, a characterization method that performs a colorimetric measurement on each stimulus yields maximum precision for minimum effort. Now, should the green color be measured in its particular location and the red one be measured in its location, or should both be measured in a common location, probably the screen center? And how should possible variations in the color of a stimulus when the other is turned off and on be handled? There is no universal best way of resolving such questions. Each must be resolved in a way that suits the particular display application. Note that this problem does not arise only for individual color characterizations. Indeed, it arises with any characterization whatsoever. The fact that every other aspect of individual color characterizations is so simple makes it particularly noticeable in this case.

Local regions of color Somewhat more complicated than individual color characterizations are ones in which it is necessary to characterize a small part of the color gamut surrounding a particular color. This might arise, for example, in a matching experiment, where all matches are close to a fixed reference stimulus. In such cases a small region of the color gamut must be characterized very precisely, with an emphasis on exact presentation of differences of color between the reference color and other colors in the region. The procedure is to make a precise colorimetric measurement of the reference color. Call the result \mathbf{X}_0 with components X_{0p} , which is the result of measuring a color having input voltages R_{0p} . Then, for small changes in input voltage ΔR_p , measure the corresponding changes in tristimulus values ΔX_p . When the results are plotted they show a region in which the changes in tristimulus value are linearly related to changes in input coordinate with nonlinear effects growing in importance near the edge of the region. The size of the nonlinear effects, when compared

to the required prevision of the characterization, determines the size of the region that can be characterized linearly. Within this region the tristimulus values are given by

$$X_i = X_{0i} + \sum_{j=1}^3 M_{ij} \Delta R_j$$

The matrix M_{ij} is determined by a multilinear regression of the data, with each component of the tristimulus values regressed simultaneously against all three input coordinates. Nonlinear effects and interaction terms are not included in the regression. The matrix entry M_{ij} , which arises as a regression coefficient, is the change of X_i for a unit change in R_j . Thus, for example, M_{23} describes how much the Y tristimulus value varies with changes in the B gun of the monitor. This type of characterization is probably the most common in color research. It is most important to remember to determine the limits of such a linear characterization, or at least to determine that the limits are outside the region in which the characterization is to be used.

One-dimensional color spaces Often it is necessary to characterize a one-dimensional set of colors, that is, a line, not necessarily straight, in color space. Here is a method for doing so easily. Make a set of colorimetric measurements spaced more or less evenly along the curve. Use the variable μ for the measurement number, numbering from one end of the line to the other. Plot the measured tristimulus values and the input voltages as functions of μ . The measurement points should be dense enough that intermediate values can be approximated by linear interpolation. Now any set of RGB values must correspond to a value of μ , not necessarily integral. How is this value determined?

1. Each measured μ has a corresponding set of RGB values.
2. Find two consecutive sets such that the desired RGB values are intermediate between the RGB values of the sets. Call the μ value of the sets μ_0 and $\mu_0 + 1$. Thus, mathematically,

$$R_i(\mu_0) \leq R_i \leq R_i(\mu_0 + 1) \quad \forall i$$

3. Now calculate “how far” between μ_0 and $\mu_0 + 1$ the desired RGB lies. The value for gun j is δ_j where

$$\delta_j = \frac{R_j - R_j(\mu_0)}{R_j(\mu_0 + 1) - R_j(\mu_0)}$$

4. If the three δ values are close together this method works. Otherwise, the line of colors must be measured more densely.
5. Use $\delta = (\delta_1 + \delta_2 + \delta_3)/3$ as the distance between μ_0 and the desired RGB value. This distance can then be used to interpolate in the tristimulus values.
6. The interpolated result is the characterized result.

$$X_i = X_i(\mu_0) + \delta(X_i(\mu_0 + 1) - X_i(\mu_0))$$

This method requires no special precautions if the line is relatively straight in RGB space. It is important to use the interpolation to calculate the tristimulus values of some measured colors to establish the precision of the measurement, or to decide if denser sampling is necessary.

Inverses The major reason for using this family of characterization methods is the computational simplicity of the characterizations and their inverses. Neither extensive memory nor expensive computation is needed to realize them.

Individual colors Since there is a small number of discrete colors, they can be arranged in a table. To get the characterization, you look up RGB in the table and read off the tristimulus values. To get the inverse, look up the tristimulus values in the output side of the table and the corresponding input gives RGB for the desired output. If the tristimulus values to be inverted are not in the table, then the color is not in the characterized set and it cannot be realized.

Local regions of color To derive the inverse for local regions of color, write the characterization equation as

$$X_i = X_{0i} = \Delta X_i = \sum_{j=1}^3 M_{ij} \Delta R_j$$

Then, inverting the matrix M_{ij} to get M_{ji}^{-1} ,

$$\Delta R_j = \sum_{i=1}^3 M_{ji}^{-1} \Delta X_i$$

which can be written explicitly as a solution for the RGB values needed to generate a color of given tristimulus values. That is,

$$R_j = R_{0j} + \sum_{i=1}^3 M_{ji}^{-1} (X_i - X_{0i})$$

It is important, after doing the above calculation, to check that the *RGB* values so determined lie within the region for which the characterization offers satisfactory precision. If not, however, unlike other methods, the inverse so determined is usually a reasonable approximation to the desired color.

Instead of inverting the matrix off-line, the three linear equations in R_j can be solved online.

One-dimensional color spaces For one-dimensional color sets the interpolation equations provide the inverse provided that *RGB* values and tristimulus values are reversed. To avoid confusion ϵ is used in place of δ .

1. To calculate ϵ find a consecutive pair of tristimulus values that contain the desired values. Call the first of the pair μ_0 .
2. Then calculate ϵ_i using

$$\epsilon_i = \frac{X_i - X_i(\mu_0)}{X_i(\mu_0 + 1) - X_i(\mu_0)}$$

If the three ϵ_i values are not close together, then the tristimulus values designate a color that is not in the set, or the wrong interval has been found in a line of colors that must be very convoluted, and the set must be measured more densely.

3. Use $\epsilon = (\epsilon_1 + \epsilon_2 + \epsilon_3)/3$ as the distance between μ_0 and the desired tristimulus values.
4. Then the *RGB* values corresponding to the desired tristimulus values are

$$R_j = R_j(\mu_0) + \epsilon(R_j(\mu_0 + 1) - R_j(\mu_0))$$

The results of this calculation should be checked to ensure that the *RGB* values do indeed lie in the color set. This is the second check that is performed, the first being a check for a small spread in the three ϵ_i values. The checks are needed if the set is convoluted, and are always passed if the set is simple.

One way of checking whether the measurement density is adequate is to cycle an *RGB* value through the characterization and its inverse. If it fails to meet the required precision, then the measurement density is probably insufficient.

Out-of-Gamut Colors Local gamut characterization methods are attractive because of their easy-to-use inverses. Their treatment of out-of-gamut colors is also appealing.

Individual colors Colors not in the set of individual colors do not have inverses, and are not found in a search of tristimulus values. Because this characterization method ignores colors near to or between the measured colors, no method is possible for representing colors that are not in the

measured set. However, if the search is done with finite precision, colors close to a measured color are aliased onto it in the search. The precision of the search determines the region of colors that is aliased onto any measured color.

Local regions of color When a color is outside the local region to which the characterization applies, this method returns the input coordinates that would apply to a linear characterization. The values are close to those produced by an exact characterization, but outside the bounds of the specified precision. The error increases as the color gets farther away from the characterized region, but slowly in most cases. Thus, this method clearly indicates when a color is outside the characterized region and fails gracefully outside the region, giving values that are close, if not exact.

One-dimensional color spaces When a set of *RGB* values or tristimulus values lies outside the one-dimensional space, the fact is indicated during the calculation. The calculation can be carried through anyway; the result is a color in the space that is close to the given color. Exactly how close, and how the characterization defines the “closest” color in the one-dimensional space, depends on the details of how the measurement samples are selected, the curvature of the color set, and the appropriate experimental definition of “close.” If this method comes to have wider utility, investigation of sampling methods might be warranted to find the “best” sampling algorithms. In the meantime, it is possible to say that for reasonably straight sets colors close to the set are mapped onto colors in the set that are colorimetrically near the original color.

Model-Dependent Characterization Methods

The two types of characterization methods described above are independent of any particular monitor properties. In fact, they can be used for any color-generating device. The third type of characterization makes use of a parametrized model of CRT color production. A characterization is created by measuring the parameters that are appropriate to the CRT being characterized. The next few sections describe the most common model for color CRTs, which is tried and true. Some CRTs may require it to be modified, but any modifications are likely to be small. The emphasis in this description is the set of assumptions on which the model is based, since violations of the assumptions require modifications to the model.

General Description The standard model of a color CRT has the following parts.

1. Any displayed color is the additive mixture of three component colors. The component colors are generally taken to be the light emitted by a single phosphor.
2. The spectral power distribution of the light in each component is determined by a single input signal, *R*, *G*, or *B*, and is independent of the other two signals.
3. The relative spectral power distribution of the light in each component is constant. Hence, the chromaticity of the component color is constant.
4. The intensity of the light in each component is a power function of the appropriate input voltage.

Taken together, these parts form a mathematical model of CRT colorimetry.

The standard model has been described many times. For a concise presentation see Cowan,¹³ for a historical one see Tannenbaum.¹⁴

Gun independence The light emitted when the input coordinates are the three voltages v_R , v_G , and v_B (generically called v_a) is $\Phi_\lambda(v_R, v_G, v_B)$. It turns out to be convenient to form slightly different components than the standard model does, following not the physical description of how a monitor operates, but the logical description of what is done to realize a color on the monitor. The first step when generating a color is to create the *RGB* input from the separate *R*, *G*, and *B* inputs. Imagine turning on each input by itself, and assume that the color when all guns are turned on together is the

additive mixture of the colors produced when the guns are turned on individually. This assumption is called “gun independence.” In terms of tristimulus values it implies the condition

$$X_i = X_{Ri} + X_{Gi} + X_{Bi}$$

(Usually CRTs realize gun independence by exciting different phosphors independently,

$$\Phi_\lambda(\mathbf{v}_R, \mathbf{v}_G, \mathbf{v}_B) = \Phi_\lambda(\mathbf{v}_R, 0, 0) + \Phi_\lambda(0, \mathbf{v}_G, 0) + \Phi_\lambda(0, 0, \mathbf{v}_B)$$

This condition is, in fact, stronger than the gun-independence assumption, and only gun independence is needed for characterization.)

Gun independence was described by Cowan and Rowell,¹⁵ along with a “shotgun” method for testing it. The tristimulus values of many colors were measured, then predictions based on the assumption of gun independence were tested. Specifically, gun independence implies consistency relationships that must hold within the set of measurements. Cowan and Rowell showed that a particular CRT had a certain level of gun independence, but didn’t attempt to make measurements of many CRTs. Clearly, it is worth making measurements on a variety of monitors to determine how widely gun independence occurs, and how large violations of it are when they do occur. For CRTs without gun independence there are two ways to take corrective action:

1. Use a characterization method that is not model-dependent.
2. Modify the monitor to create gun independence. For most CRTs there are good, inexpensive methods for improving gun independence. They will be more widely available in the future as characterization problems caused by failures of gun independence become more widely known.

Extreme settings of the monitor controls can cause violations of gun independence in otherwise sound monitors. The worst culprit is usually turning the brightness and/or contrast up so high that blooming appears at high input levels.

Phosphor constancy When gun independence holds, it is necessary to characterize only the colors that arise when a single gun is turned on, since the rest can be derived from them. Usually, it is assumed that the colors produced when a single gun is turned on have constant chromaticity. Thus, for example, the tristimulus values of colors produced by turning on the red gun alone take the form

$$X_{Ri}(\mathbf{v}_R) = E_R(\mathbf{v}_R) \cdot x_{Ri}$$

where x_{Ri} the chromaticity coordinates of the emitted light, $x_{R'}$, $y_{R'}$, and $z_{R'}$ are independent of the voltage input. Thus the tristimulus values of the color produced depend on input voltage only through the intensity $E_R(\mathbf{v}_R)$.

The engineering that usually lies behind phosphor constancy is the following. A CRT should be designed so that the beam current in a given electron gun is independent of the beam current in any other gun. The deflection and shadowmask geometries should be designed and adjusted so that the electrons from any gun fall only on a single phosphor type. The physical properties of the phosphor should guarantee that the phosphor emits light of which the chromaticity is independent of the intensity. Meeting these conditions ensures phosphor constancy. It is possible, but unusual, to have phosphor constancy under other conditions, if it happens that effects cancel out appropriately.

The measurement of phosphor constancy is described by Cowan and Rowell.¹⁵ Here is a short description of how to make the required measurements for one of the guns.

1. Turn the gun on at a variety of input voltages, and make a colorimetric measurement at each input voltage.
2. Calculate chromaticity coordinates for each measurement, which are independent of input voltage if phosphor constancy is present.

3. Constancy of chromaticity should also obtain in the presence of varying input to the other guns. To consider this case, turn the two other guns on, leaving the gun to be measured off.
4. Measure the tristimulus values for a baseline value.
5. Then turn the gun to be measured on to a variety of input voltages, making a colorimetric measurement at each.
6. Subtract the baseline tristimulus values from each measurement, then calculate the chromaticity coordinates. These should be constant, independent of the input to the measured gun, and independent of the input to the other two guns as well.

One of the most common monitor characterization problems is a measured lack of phosphor constancy. This is a serious problem, since the functions derived from the CRT model cannot then be easily inverted. Sometimes the measured lack of phosphor constancy is real, in which case there is no solution but to use a different type of characterization procedure. More often, the measurement is caused by poor CRT setup or viewing conditions. This occurs as follows, first for viewing conditions. Ambient light reflected from the screen of the monitor is added to the emitted light. This light is independent of input voltage, and adds equally to all colors. Imagine a phosphor of constant chromaticity (x, y) , at a variety of input voltages. The tristimulus values of light emitted from the screen are

$$X_i = X_{0i} + E(\mathbf{v}) \cdot x_i$$

where X_{0i} are the tristimulus values of the ambient light reflected from the screen. Note that subtracting the tristimulus values measured when the input voltage is zero from each of the other measurements gives phosphor constancy. (Failure to do the subtraction gives phosphor chromaticities that tend toward white as the input voltage decreases.) This practice should be followed when ambient light is present. Psychophysical experiments usually go out of their way to exclude ambient light, so the above considerations are usually not a problem. But under more normal viewing conditions ambient light is always present.

When the black level/brightness control is set above its minimum there is a background level of light emitted from the screen even when the input voltage is zero. This light should be constant. Thus, it can be handled in exactly the same way as ambient light, and the two can be lumped together in a single normalizing measurement.

Now let's put together this gun independence and phosphor constancy. Gun independence means that the tristimulus values $X_i(\mathbf{v}_R, \mathbf{v}_G, \mathbf{v}_B)$ can be treated as the sum of three tristimulus values that depend on only a single gun:

$$X_i(\mathbf{v}_R, \mathbf{v}_G, \mathbf{v}_B) = \sum_{a=R,G,B} X_{ai}(\mathbf{v}_a)$$

Phosphor constancy means that the tristimulus values from a single gun, once the background has been subtracted away, have constant chromaticity:

$$X_i(\mathbf{v}_R, \mathbf{v}_G, \mathbf{v}_B) = (X_0, Y_0, Z_0) + \sum_{a=R,G,B} E_a(\mathbf{v}_a) \cdot x_{ai}$$

To complete the characterization a model is needed for the phosphor output intensity $E_a(\mathbf{v}_a)$.

Phosphor output models Several different methods exist for modeling the phosphor output. The most general, and probably the best, is to build a table for each gun. To do so requires only photometric measurement, and the responsivity of the photometer need not be known. (This property of the model is very convenient. Colorimetric measurements are fairly easy to perform at light levels in the top 70 percent of the light levels produced by a CRT, but hard to perform for dark colors. In this model the top of the output range can be used when doing the colorimetric measurements required to obtain phosphor chromaticities. Then, photometric measurements can be performed easily using photodiodes to calibrate the output response at the low end of the CRT output.) Choose the largest

input voltage to be used, and call it $v_{a\max}$. Measure its light output with the photometer. Then measure the light output from a range of smaller values v_a . The ratio of the light at lower input to the light at maximum input is the relative excitation $e_a(v_a)$. Store the values in a table. When excitations for intermediate voltage values are needed, find them by linear interpolation.

$$e_a(v_a) = \mu e_a(v_{a1}) + (1 - \mu)e_a(v_{a2})$$

where the interpolation parameter μ is

$$\mu = \frac{v_a - v_{a2}}{v_{a1} - v_{a2}}$$

The memory requirements of the table can be reduced by using more complicated functions, but at the price of complexity of computation in the model.

It is also possible to approximate the measurements using a parameterized function. The most commonly used form is the gamma correction equation

$$e_a(v_a) = e_{a\max} (v_a / v_{a\max})^{\gamma_a}$$

with two parameters γ_a and $e_{a\max}$. Here, $v_{a\max}$ is the maximum input voltage, $e_{a\max}$ is the maximum relative excitation, usually taken to be 1.0, and γ_a is called the gamma correction exponent. It is determined empirically by regressing the logarithm of the measured excitation against the logarithm of the input voltage.

Normalization The excitation that enters the characterization equation $E_a(v_a)$ is the product of the voltage-dependent relative excitation $e_a(v_a)$ and a normalization coefficient N_a :

$$E_a(v_a) = N_a \cdot e_a(v_a)$$

The normalization coefficients complete the characterization model.

Summary equation The characterization is summarized by the single equation

$$X_i = X_{0i} + \sum_{a=R,G,B} N_a \cdot e_a(v_a) \cdot x_{ai}$$

which provides the tristimulus values for any set of *RGB* input coordinates.

Conditions for Use The assumptions discussed above are the conditions under which the model developed here can be used successfully. The assumptions should be checked carefully before this characterization method is used. The departure of the CRT from the ideal CRT, defined by perfect adherence to these conditions, is a measure of the imprecision in the characterization procedure. Note that when the CRT fails to meet the conditions, which occurs most frequently near the edges of its gamut, the erroneous colors are usually plausible. Thus, the model can be in error by a surprisingly large amount and still produce satisfactory characterizations in noncritical applications.

Partial models It is fairly clear that there are partial models meeting some of the conditions, and that they provide useful approaches to characterization. For example, suppose that the phosphor chromaticities vary with input voltage. They can be stored in a table indexed by input voltage (exact values are produced by interpolation) and used in a variant of the characterization equation:

$$X_i = X_{0i} + \sum_{a=R,G,B} N_a \cdot e_a(v_a) \cdot x_{ai}(v_a)$$

Such a generalization works well for the transformation from input voltages to tristimulus values, but the inverse transformation is virtually unusable.

Measurement of Parameters The characterization equation requires the measurement of a set of parameters that varies from CRT to CRT. Thus, they must be measured for the CRT to be characterized. Different types of measurements—spectral, colorimetric, and photometric—are needed, as described below.

Ambient light and black level X_{0i} A colorimetric measurement of the screen with all input voltages set to zero produces this value. It should be measured under the exact conditions in which the characterization will be used.

Phosphor chromaticities x_{ai} A colorimetric measurement of the screen with a single gun turned on produces these values. Careful subtraction of the ambient light and black level is important. Measuring at a variety of input voltages produces:

1. A measure of phosphor constancy
2. The range of input voltages over which phosphor constancy holds well enough for use of the characterization model
3. A value for the phosphor chromaticities, produced by averaging over the chromaticities in the characterization range

Gamma correction functions Measurement of the relationship between the input voltages and the excitation functions requires photometric measurement at a large variety of input voltages. If interpolation into a table is the method of choice, the sampling of input voltages must be dense enough that the interpolation method chosen, usually linear, is close enough to the exact values to yield the desired accuracy. If a functional relationship, like the power function, is chosen the validity of the function chosen must be inspected very carefully. Whatever method is chosen it is essential that the ambient light/black level be subtracted from the measurements. If it is not, it is counted three times when the contributions from the three guns are added together.

Normalization coefficients The normalization coefficients are determined through the characterization equation

$$X_i = \sum_{a=R,G,B} N_a \cdot e_a(v_a) \cdot x_{ai}$$

They may be determined in several ways, falling into two categories: measurement methods and comparison methods. A typical measurement method assumes a single color of known tristimulus values. This knowledge can be produced by a colorimetric measurement or by varying *RGB* to match a known sample such as the reference field of a color comparator. The tristimulus values are substituted into the equation above, giving three equations for the three normalization coefficients. Solving the equations gives the normalization coefficients. Note that the relative excitations and the phosphor chromaticities are both unitless. Thus, the normalization coefficients carry whatever units are used to measure the tristimulus values. The units enter naturally when the linear equations are solved.

The method for solving these linear equations utilizes the inverse transform. A matrix \mathbf{M} , the entries of which are the phosphor chromaticities

$$\mathbf{M}_{ai} = x_{ai}$$

is defined, and its inverse, \mathbf{M}_{ia}^{-1} determined. Multiplying it into the characterization equation gives

$$\sum_{i=1}^3 \mathbf{M}_{ia}^{-1} X_i = N_a \cdot e_a(v_a)$$

Then N_a is calculated from

$$N_a = e_a(v_a) \cdot \sum_{i=1}^3 \mathbf{M}_{ia}^{-1} X_i$$

Note that $e_a(v_a)$ is known since the input voltage that matches the standard, or of the color that has been measured, is known.

The second method is used when it is possible to match a known standard. For example, it might be possible to find voltages that make colors produced by individual guns equiluminous. Call the voltages v_{YR} , v_{YG} , and v_{YB} , and substitute them into the $Y = X_2$ characterization equation

$$N_R e_R(v_{YR}) x_{R2} = N_G e_G(v_{YG}) x_{G2} = N_B e_B(v_{YB}) x_{B2}$$

This equation is easily solved for the ratios N_G/N_R and N_B/N_R , leaving a single constant, an overall normalizing coefficient, undetermined. This constant, which carries the units in which the tristimulus values are measured, can be determined only by an absolute measurement (or, equivalently, a color match to an absolute standard). However, because vision is insensitive to variations in overall intensity, the precise value of the overall normalizing coefficient is unimportant in many applications.

Inverse Transformations One of the most attractive features of model-dependent characterizations is the simplicity of their inverses. It is simple to invert the characterization equation

$$X_i = \sum_{a=R,G,B} N_a \cdot e_a(v_a) \cdot x_{ai}$$

which determines the tristimulus values of a color in terms of the input voltages that cause it to be displayed, into the inverse equation, which specifies input voltages that create a color of given tristimulus values. Define the matrix of phosphor chromaticities

$$M_{ai} = x_{ai}$$

Then determine its inverse M_{ia}^{-1} by conventional means. It is used to solve the characterization function

$$\sum_{i=1}^3 M_{ia}^{-1} X_i = N_a \cdot e_a(v_a)$$

Then, after the normalization coefficients are divided out

$$e_a(v_a) = \frac{1}{N_a} \sum_{i=1}^3 M_{ia}^{-1} X_i$$

It is then necessary to find the input voltages that give the appropriate excitations e_a , which requires inversion of the excitation function. Thus,

$$v_a = e_a^{-1} \left(\frac{1}{N_a} \sum_{i=1}^3 M_{ia}^{-1} X_i \right)$$

If the relative excitation function is specified in a closed form, the inverse can be calculated using only elementary algebra. Otherwise, it is necessary to develop a tabular inverse, which is most easily done by one-dimensional interpolation.

Violations of phosphor constancy When phosphor constancy does not hold, the inverse cannot be calculated using the methods given above. Numerical inversion techniques, which are beyond the scope of this section, can be used, but it is probably a better idea to use a different calibration method.

Out-of-Gamut Colors Most inversion methods fail internally if they are asked to produce the input coordinates needed to display an out-of-gamut color. For example, tabular inverses fail to find an interval or cell in which to interpolate. Because the functions in model-dependent characterizations are well defined beyond the gamut, these methods generally do not fail, but calculate values that

lie outside the domain of input coordinates. Thus, input coordinate values calculated from inverses of model-dependent characterizations should be checked, because writing out-of-range values to color lookup tables produces unpredictable results.

Absolute Characterization versus Characterization for Interaction

The characterization methods discussed above are designed to provide absolute characterizations, which are needed for color coordinates defined in terms of an external standard. Probably just as often characterization is needed to support specific interaction techniques. Suppose, for example, an experiment calls for the observer to have three CIELAB controls, one changing L^* , one changing a^* , and one changing b^* . To make this possible a characterization of the CRT in terms of CIE tristimulus values must be done to provide a basis for the transformation from tristimulus values to $L^*a^*b^*$. This type of characterization situation has several characteristics that are quite common:

1. At any given time there is a known color on the screen and the problem is to calculate a color slightly different from it. This situation arises because the control is sampled frequently enough that large color changes are produced only cumulatively over many color updates.
2. The new color coordinates must be calculated quickly and frequently.
3. Errors in the calculated color increments are errors in quantities that are already small. Thus, they can be considerably larger (in percentage terms) than would be tolerable for an absolute characterization.
4. The inverse transformation, tristimulus values to RGB is usually needed.

This combination of requirements is most easily met when the characterization has a linear inverse, and as little calculation as possible in the main loop. Exhaustive characterization methods have the necessary ingredients as they stand, provided an inverse table has been built, as do local methods. But nonlinear transformations, like the one from $L^*a^*b^*$ to tristimulus values violate the conditions imposed above, and it is necessary to linearize them. They can then be combined with whatever characterization is in use to linearize the entire transformation from input coordinates to color specification.

Linearization Consider a transformation f that takes a vector \mathbf{x} into a vector \mathbf{y} . An example is the transformation between $L^*a^*b^*$ and XYZ . Suppose two vectors that are known to transform into each other are known: \mathbf{x}_0 ($L_0^*a_0^*b_0^*$) corresponds to the set \mathbf{y}_0 ($X_0Y_0Z_0$). Now, if one vector undergoes a small change, how does the other change? If the transformation is suitably smooth (and the functions of color science are adequately smooth), small changes of one variable can be expressed in terms of small changes of the other by simple matrix multiplication

$$\Delta y_i = \sum_j M_{ij} \Delta x_j$$

The entries in the matrix \mathbf{M} are the partial derivatives of f with respect to \mathbf{x} , evaluated at \mathbf{x}_0 . That is,

$$M_{ij} = \left. \frac{\partial f_i}{\partial x_j} \right|_{\mathbf{x}=\mathbf{x}_0}$$

Even though the computation of the derivatives may be expensive, they need to be calculated only infrequently, and the calculation can be scheduled when the system has idle resources available.

Practical Comments

CRT characterization is essentially a practical matter. It is almost always done because a superordinate goal requires it. Thus, it is appropriate to end this chapter with several practical remarks that make characterization easier and less costly.

Do as little characterization as possible! Characterization is usually done within a fixed budget of time and equipment. Identifying the least characterization that provides all the information needed for the application allows the extra resources to be used to improve precision and to check measurements.

Characterize frequently! CRTs age and drift; colleagues turn knobs unaware that they change a critical characterization. Discovering that a CRT does not match its characterization function invalidates all data collected or experiments run since the characterization was last checked. There are only two defenses. First, measure the CRT as often as is feasible. Second, devise visual checks that can be run frequently and that will catch most characterization errors. It is possible to use luminance-matching techniques like minimally distinct border or minimum motion to check characterizations in a few seconds.

Understand the principles underlying the characterization methods and alter them to fit the specific stimuli needed for the application. In particular, CRTs have many ways of being unstable. The best defense is to characterize the CRT with stimuli that are as close as possible to those that are used in the application: same size, same position on the screen, measurement apparatus at the observer's eye point, same background color and intensity, and so on. Such adherence is the best defense against inadvertently finding a new type of CRT instability that invalidates months of work.

22.5 AN INTRODUCTION TO LIQUID CRYSTAL DISPLAYS

The purpose of this section is to provide an introduction to the operational principles that provide color output from liquid crystal displays (LCDs). No actual colorimetric measurements are shown, because LCDs are evolving too fast for them to be of any lasting interest. Instead, schematic illustrations are given, showing features that the author considers likely to be present in future displays.

At the time of writing, colorimetric quality of LCDs is uneven. Nonetheless, since LCDs are semiconductor components and since they are now being produced in large quantities for low-cost computers, it is reasonable to expect a rapid decrease in cost combined with an increase in quality over the next decade. As the succeeding sections should make clear, they have a variety of interesting properties which will often make them preferred to CRTs as visual stimulators. Only by actively following their development will it be possible to use them at the earliest opportunity.

Two properties make liquid crystal devices particularly interesting to the visual scientist. First, they are usable as self-luminous, transmissive, or reflective media. They combine this flexibility with the range of temporal and spatial control possible on a CRT. This property is likely to make them particularly valuable for studies of color appearance and color constancy. Second, they have the ability to produce temporal and spatial edges that are much sharper than those producible on a CRT. This property has been called "acutance" in the display literature.¹⁶ High acutance is potentially a serious problem for image display, since the gaussian blur that occurs at edges in CRT images helps to hide artifacts of the CRT raster. But high acutance also presents an opportunity since many small-scale irregularities that cannot be controlled on the CRT can be controlled on the LCD.

Operational Principles of Monochrome LCDs

This section describes the operation of monochrome LCDs. The components of color LCDs are so similar to those of monochrome ones that they are most easily explained based on an understanding of monochrome LCDs.

Overview The important components of a monochrome LCD are shown schematically in Fig. 16. In the dark, the main source of light is the backlight. Its luminous flux passes through a diffuser. The uniformity of the backlight/diffuser pair determines the spatial uniformity of the display surface. Light from the diffuser then passes through a light-modulating element that consists of two crossed

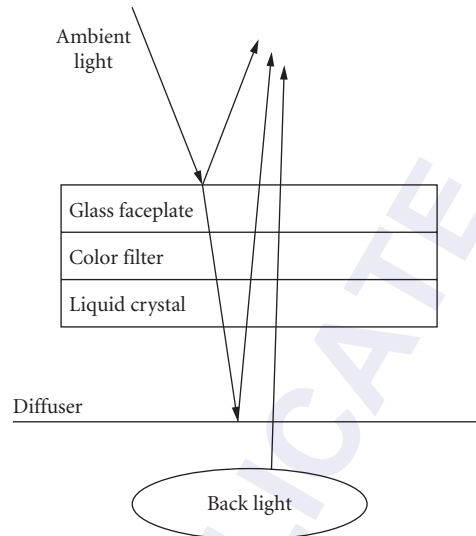


FIGURE 16 Illustration of the typical components of a single LCD pixel, seen side on. The arrows show the three light paths that contribute to the light emitted from the pixel.

polarizers and an intervening space, usually a few microns thick, filled with the liquid crystal. The liquid crystal rotates the polarization of light passing through it by an amount that depends on the local electric field. Thus, the electric field can be varied to vary the amount of light passing through the light-modulating element. The electric field is produced by a capacitor that covers the area of the pixel. The different varieties of LCD differ mainly in the electronics used to drive the capacitor. At present, the most common type is bilevel, with the capacitor controlled by a transistor that is either on or off, thus defining two light levels, white and black. More interesting are multilevel devices, in which an analog electric field allows more or less continuous variation of light output. Currently available analog devices are capable of 16 levels of gray scale with a pixel that is about 150 microns square.

In bilevel devices the array of pixels is very similar to the layout of memory cells in a random access semiconductor memory. Thus, the display itself can be considered a write-only random-access frame buffer. Most LCDs do not, at present, make use of the random-access capability, but accept input via the RS-170 analog signal standard used for CRTs. The multilevel variant is similar, except that an analog value is stored at each location instead of 0 or 1.

An interesting feature of this design is the ability of ambient light to enter the display. It then passes through the light-modulating element and is scattered with random phase by the diffuser. From there it passes back through the light-modulating element to the viewer. Thus, under high-ambient conditions there are two contributions to light emitted from the display: one from the backlight that is linear with respect to the transmittance of the light-modulating element and one from ambient light that is quadratic with respect to the transmittance of the light-modulating element. More specifically, the contribution from the backlight, $\Phi_{\lambda}^{(BL)}$ is given by

$$\Phi_{\lambda}^{(BL)} = \tau(V)\Phi_{0\lambda}^{(BL)}$$

where $\Phi_{0\lambda}^{(BL)}$ is the light emitted from the display when the light-modulating element is in its maximally transmissive state and $\tau(V)$ is the voltage-dependent transmittance of the light-modulating element. Note that the light-modulating process is assumed to be spectrally neutral. (In fact, chemists

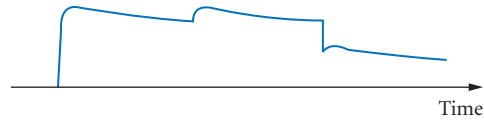


FIGURE 17 Illustration of the time dependence of the light output from an LCD pixel that is turned on for two refreshes, then turned to a lower level for the third refresh.

attempt to make the liquid crystal as spectrally neutral as possible.) The contribution from ambient light $\Phi_{\lambda}^{(AMB)}$ is

$$\Phi_{\lambda}^{(AMB)} = \tau^2(V)(1 - R(\lambda))D(\lambda)\Phi_{0\lambda}^{(AMB)}$$

where $\Phi_{0\lambda}^{(AMB)}$ is the ambient light incident on the display, $R(\lambda)$ is the reflectance of the faceplate of the display, and $D(\lambda)$ is the proportion of light reflected back by the diffuser. The second contribution is particularly interesting because the display is essentially characterized by a variable reflectance $\tau^2(V)(1 - R(\lambda))D(\lambda)$. Thus it is possible to create reflective images that are computer-controlled pixel by pixel. Note the importance of the front surface reflection in this pixel model. Light that is not transmitted at the front surface produces glare; light that is transmitted creates the image. Consequently, it is doubly important that the faceplate be treated to minimize reflection.

Pixel Structure The most interesting characteristic of the LCD is the shape and layout of its pixels, which are unlike the CRT pixel both spatially and temporally.

The time dependence of the light output from an LCD pixel is shown schematically in Fig. 17. Each refresh of the pixel produces a sharp rise followed by a slow fall. If the light output decreases from one frame to the next, the decrease occurs abruptly at the beginning of the pixel, without the exponential tail characteristic of CRTs with long persistence phosphors. The turnoff is limited only by the resistance of the circuit that removes charge from the capacitor. The slow fall is produced by loss of charge from the capacitor: the circuit is essentially a sample and hold and the flatness of the pixel profile is determined by the quality of the hold.

The spatial structure of the LCD display surface is shown schematically in Fig. 18. The pixels are nonoverlapping, with dark lines between them because an interpixel mask is used to hide regions

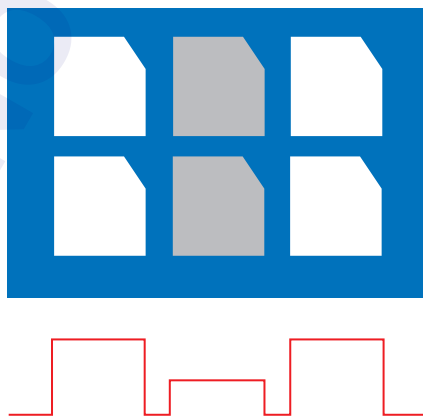


FIGURE 18 Illustration of the layout and shape of pixels on an LCD. The lower panel shows the spatial dependence of emitted light, on a cross section through the center of a row of pixels. The two middle pixels are turned on to a lower level than the outside ones.

where the electric fields of neighboring capacitors overlap. The light emitted is uniform over the pixel area, so that the intensity profile is uniform. Usually there is a small amount taken off the corner of the pixel to allow space for the electronic circuits that turn the pixel off and on. The pitch of the display, the distance from one pixel to the center of the next pixel, ranges from 100 to 500 μm .

Problems Current high resolution, analog, color LCDs suffer from a variety of problems, which may be expected to diminish rapidly in importance as manufacturing technology improves.

Thickness variation Uniformity of output depends critically on precise control of the thickness of the liquid crystal from pixel to pixel.

Heat Current displays are very temperature-dependent, changing their performance as they warm up and if anything warm, like a human hand, is held in contact with them.

Gray scale The present mass-produced LCDs have poor gray-scale capability: each pixel can be only off or on. The manufacturing processes needed to produce gray scale, which produce about 30 levels of gray scale on prototype displays, need considerable development.

Electronics A typical display depends on thousands of address wires, connected to the display at the pitch of the pixels. These wires are delicate and subject to failure, which produces a row or column of inactive pixels, called a line-out.

The Color LCD

The color LCD is a straightforward generalization of the monochrome LCD. Colored filters are placed in the light path to color the light emitted from each pixel. Three differently colored filters, and sometimes four, are used to create several interpenetrating colored images that combine additively to form a full-color image. The repetition pattern of the different images is regular, and it is usually possible to group a set of differently colored pixels into a single-colored pixel with three or four primaries, as shown in Fig. 19. (Terminology is important but nonstandard. In this chapter a single-colored pixel is called a monochrome pixel; a set of differently colored monochrome pixels that combines additively to produce a continuum of color at a given point is called a color pixel.)

Geometry of the Color Pixel There are a variety of different geometrical arrangements of monochrome pixels within a color pixel. Several are illustrated in Fig. 19. It is not currently known which

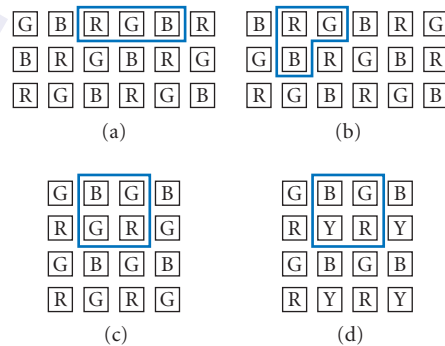


FIGURE 19 Schematic diagram of a variety of different monochrome pixel geometries: (a) and (b) show two triad geometries; (c) and (d) show two quad geometries. The color pixel is outlined in gray.

arrangement is best for which type of visual information. When the viewer is far enough back from the display the arrangement does not matter, only the relative number of each primary color. One strategy that is being tried is to add a fourth color, white or yellow, in a quad arrangement in order to increase brightness. Such a display offers interesting possibilities for experimentation with mesopic stimuli and for diagnosis and investigation of color-vision defects.

It is obvious that every geometry has directions in which the primaries form unwanted color patterns. Such patterns do not arise in CRT imagery because there is a random correlation between the dot triads that produce color and the pixel locations. Whether there is a regular geometry, as yet untried, that removes the patterns, whether software techniques similar to antialiasing can eliminate their visual effects, or whether a randomized primary arrangement is needed to remove the patterns is a problem as yet unresolved. (It should also be remembered that some stimuli may best be created by graphical techniques that take advantage of the existence of the patterns.) Whatever the solution, however, it must address the interrelationship between spatial and color aspects of the stimulus in the visual system of the perceiver.

Colorimetry of the Color Pixel The dominant factor in the colorimetry of LCDs is the interaction between the spectral power distribution of the backlight and the transmittances of the filters. Typically backlight and filter design deals with a trade-off between brightness and color gamut. The higher the excitation purity of the filter the more saturated the primary but the less bright the display. An additional factor considered in this trade-off is heating of the liquid crystal, which is greater when the filters have high excitation purity. Choice of a backlight that has as much as possible of its spectral power in wavelengths passed by the filters will improve performance.

The colorimetric properties of an LCD can be summed up in a few equations. The first is that the color of a pixel is the additive mixture of the colors of its monochrome pixel components. The sum of the spectral powers is

$$\Phi_{\lambda} = \sum_{a=1}^{N_p} \Phi_{a\lambda}(V_a)$$

where $\Phi_{a\lambda}$ is the spectral power emitted by monochrome pixel a , which depends on the voltage V_a applied to it. Similarly, the tristimulus values are the sum of the tristimulus values of the monochrome pixels

$$X_i = \sum_{a=1}^{N_p} X_{ai}$$

where X_{ai} is the set of tristimulus values for primary a . Note that “gun independence” is assumed.

The spectral power contributed by monochrome pixel a is the sum of three components: one from the backlight, $\Phi_{a\lambda}^{(BL)}(V_a)$, one from ambient light reflected from the front surface of the LCD, $\Phi(R)_{a\lambda}$, and one from ambient light reemitted from the display, $\Phi_{a\lambda}^{(AMB)}(V_a)$.

$$\Phi_{a\lambda}(V_a) = \Phi_{a\lambda}^{(BL)}(V_a) + \Phi_{a\lambda}^{(R)} + \Phi_{a\lambda}^{(AMB)}(V_a)$$

The contribution from the backlight is

$$\Phi_{a\lambda}^{(BL)}(V_a) = \tau_a(\lambda)\tau(V_a)\Phi_{0\lambda}^{(BL)}$$

where $\Phi_{0\lambda}^{(BL)}$ is the spectral power distribution of the backlight, $\tau(V_a)$ is the voltage-dependent transmittance of the liquid crystal/polarizer sandwich, and $\tau_a(\lambda)$ is the transmittance of the color filter. The function $\tau(V)$ is independent of the primary because the display uses the same light-modulating element for each primary. The function $\tau_a(\lambda)$ is independent of the applied voltage since the light-modulating element is spectrally neutral. In the absence of ambient light, this term is the only contribution to the light emitted by the LCD. The colorimetric properties are then identical to those of a CRT with primaries of the same chromaticity, which is the chromaticity of the product $\tau_a(\lambda)\Phi_{0\lambda}^{(BL)}$, though the equivalent of the gamma function $\tau(V_a)$ is certain to have a different functional form.

The contribution from light reflected from the front surface is

$$\Phi_{a\lambda}^{(R)} = R(\lambda)\Phi_{0\lambda}^{(AMB)}$$

where $R(\lambda)$ is the reflectance of the front surface and $\Phi_{0\lambda}^{(AMB)}$ is the ambient light incident on the surface. Usually, the reflectance is spectrally neutral (independent of wavelength) and the reflected light has the same color as the incident ambient light.

The contribution from light reemitted from the display is

$$\Phi_{a\lambda}^{(R)} = (1 - R(\lambda))D(\lambda)\tau_a^2(\lambda)\tau^2(V_a)\Phi_{0\lambda}^{(AMB)}$$

where $D(\lambda)$ is the reflectance of the diffuser. This contribution can be modeled as a light of the same chromaticity as the spectral power distribution $(1 - R(\lambda))D(\lambda)\tau_a^2(\lambda)\Phi_{0\lambda}^{(AMB)}$, with its intensity modulated by the function $\tau^2(V_a)$. Note that this chromaticity is not in general the same as the chromaticity produced by light from the backlight. (It will usually be more saturated than the backlight component since it passes through the color filter twice.) Note also that the voltage dependence of the intensity of this light differs from that of the backlight component. Thus, the chromaticity of the primary changes as the voltage changes when ambient light is present, and the effect cannot be subtracted off, as can the ambient component.

Controls and Input Standards LCDs will have standard sets of controls and input signals only when the technology is much more mature than it is at present. Currently, the most common input seems to be analog video, as used for CRTs. While this allows a one-for-one substitution of an LCD for a CRT, it seems quite inappropriate for computer applications. Specifically, information is extracted from the frame buffer, a random-access digital device, laboriously sequenced into the serial analog video signal, the reextracted for presentation for the random-access LCD. Thus, it seems likely that digital random-access input standards are likely to supersede video signals for LCD input, and that the display will be more tightly coupled to the image storage system than is common with CRTs.

Temporal and Spatial Variations in Output

LCD fabrication technology is changing too quickly for quantitative limits to spatial and temporal variability to be predictable. Nonetheless, it seems likely that certain qualitative properties are inherent in all LCD designs. They are discussed in general terms in the next few sections.

Short-Time Temporal Variation The light output over a short time (about 100 ms) is shown schematically in Fig. 17. The details shown in that figure, the rise and decay times and the turnoff time, are bound to change as fabrication techniques improve. Currently, the turnoff time is not very good for most LCDs, and images leave shadows that decay visibly for several seconds after they are removed. That defect is not inherent in LCD electronics and it should be expected to disappear as circuitry improves. The other possible defect is the ripple in the light output, which is nonetheless much smaller than the ripple of CRT light output. The size of the ripple is determined by the quality of the sample-and-hold circuit that maintains charge on the capacitor, and is likely to decrease as LCDs improve.

The interaction between ripple and turnoff time is quite different for an LCD than for a CRT. To decrease the ripple on a CRT it is necessary to increase the decay time of the phosphors; degradation of the turnoff time cannot be avoided. For an LCD, on the other hand, the turnoff is active, and independent of the ripple. Thus, it is possible to improve the ripple without degrading the turnoff. As a result, future LCDs are likely to be very useful for generating stimuli that require precise temporal control.

Long-Time Temporal Variation Even at the current state of development, the long-time stability of LCDs is quite good, with reports showing variations of about 2 percent on a time scale of several hours. This performance is comparable to good-quality CRTs and incandescent light sources.

The main contributor to instability is heat. LCDs are only stable once they are warmed up, and even small changes in cooling configuration can cause appreciable changes in light output. Thus, good stability of light output requires well-controlled temperature and a long warm-up time.

Small-Scale Spatial Variation Small-scale variation is determined by the spatial structure of the pixel, which is illustrated in Fig. 18. The important feature is the sharp edge of the pixel: unlike CRT pixels, there is no blending at the edges of adjacent pixels. This feature makes the creation of some stimulus characteristics easy—sharp vertical and horizontal edges, for example—and makes the creation of other stimulus parameters very difficult—rounded corners, for example. Many graphical techniques and image-quality criteria have evolved to take advantage of display characteristics that are peculiar to the CRT; there is likely to be a substantial investment in devising techniques that are well suited to the very different pixel profile of the LCD.

An interesting future possibility arises because of the similarity of LCD manufacturing technology to that of random-access memories. Consequently, it is likely that future LCDs will be limited in resolution not by pixel size but by input bandwidth. If so, it would be sensible to have logical pixels within the control system that consist of many physical pixels on the display, with enough processing power on the display itself to translate commands referring to logical pixels into drive signals for physical pixels. In fact, a physical pixel could even belong to more than one logical pixel. If such a development occurs, considerable control over the pixel profile will be possible, which may greatly extend the range of spatial variation that images can possess.

Large-Scale Spatial Variation Because LCDs have no large-scale structural features, like the beam deflection in the CRT, only manufacturing tolerances should induce spatial variation. The main source of manufacturing variability at the time of writing is the physical size of the pixels and the electronic components—capacitors, especially—that control them. Ideally, such variations would be random, so that spatial variation would add only gaussian noise to the displayed image.

Much more serious is angular variation of emitted light. The light-emitting element uses effects that have strong directional dependence, and it is sometimes even possible to find angles at which an LCD reverses contrast compared to perpendicular viewing. Although this effect is inherent in LC technology, it is reasonable to hope that display improvements will reduce the angular variability below its currently unacceptable level.

22.6 ACKNOWLEDGMENTS

The author would particularly like to express his gratitude to the late Gunter Wyszecki who provided impetus for his interest in this subject, and who provided intellectual and logistical support during the time in which the author carried out most of the work on which this chapter is based. A large number of other people have provided useful insights and technical assistance over the years, including, but not restricted to: John Beatty, Ian Bell, David Brainard, Pierre Jolicoeur, Nelson Rowell, Marueen Stone, and Brian Wandell. The author would also like to thank the editor, David Williams, for patience above and beyond the call of duty during the very difficult time when I was putting this material into the form in which it appears here.

22.7 REFERENCES

1. D. G. Fink, K. B. Benson, C. W. Rhodes, M. O. Felix, L. H. Hoke Jr, and G. M. Stamp, "Television and Facsimile Systems," *Electronic Engineers Handbook*, D. G. Fink and D. Christiansen (eds.), 3d ed., McGraw-Hill, New York, 1989, pp. 20-1–20-127.
2. Electronic Industries Association, *Electrical Performance Standards—Monochrome Television Studio Facilities*, RS-170, Washington, D.C., 1957.

3. Electronic Industries Association, *Electrical Performance Standards—Monochrome Television Studio Facilities*, RS-343, Washington, D.C., 1969.
4. N. P. Lyons and J. E. Farrell, "Linear Systems Analysis of CRT Displays," *Society for Information Display Annual Symposium: Digest of Technical Papers* **20**:220–223 (1989).
5. J. B. Mulligan and L. S. Stone, "Half-toning Methods for the Generation of Motion Stimuli," *Journal of the Optical Society of America* **A6**:1217–1227 (1989).
6. A. C. Naiman and W. Makous, "Spatial Non-Linearities of Grayscale CRT Pixels," *Proceedings of SPIE: Human Vision, Visual Processing, and Digital Display III* **1666**:41–56 (1992).
7. B. MacIntyre and W. B. Cowan, "A Practical Approach to Calculating Luminance Contrast on a CRT," *ACM Transaction on Graphics* **11**:336–347 (1992).
8. P. M. Tannenbaum, 1986 (personal communication).
9. D. H. Brainard, "Calibration of a Computer-Controlled Monitor," *Color Research and Application* **14**:23–34 (1989).
10. W. T. Wallace and G. R. Lockhead, "Brightness of Luminance Distributions with Gradient Changes," *Vision Research* **27**:1589–1602 (1987).
11. F. Benedikt, "A Tutorial on Color Monitor Alignment and Television Viewing Conditions," Development Report 6272, Canadian Broadcasting Corporation, Montreal, 1986.
12. I. E. Bell and W. Cowan, "Characterizing Printer Gamuts Using Tetrahedral Interpolation," *Color Imaging Conference: Transforms and Transportability of Color*, Scottsdale, Arizona, 1993.
13. W. B. Cowan, "An Inexpensive Scheme for Calibration of a Color Monitor in Terms of CIE Standard Coordinates," *Computer Graphics* **17**(3):315–321 (1983).
14. P. M. Tannenbaum, "The Colorimetry of Color Displays: 1950 to the Present," *Color Research and Application* **11**:S27–S28 (1986).
15. W. B. Cowan and N. L. Rowell, "On the Gun Independence and Phosphor Constancy of Color Video Monitors," *Color Research and Application* **11**:S34–S38 (1986).
16. L. M. Biberman, "Image Quality," *Perception of Displayed Information*, L. M. Biberman (ed.), Plenum, New York, 1973.

This page intentionally left blank.

DO NOT DUPLICATE

VISION PROBLEMS AT COMPUTERS

Jeffrey Anshel

*Corporate Vision Consulting
Encinitas, California*

James E. Sheedy

*College of Optometry
Pacific University
Forest Grove, Oregon*

23.1 GLOSSARY

Accommodation. In regard to the visual system, accommodation is the focusing ability of the eye.

Acuity. A measure of the ability of the eye to resolve fine detail, specifically to distinguish that two points separated in space are distinctly separate.

Afterimage. An optical illusion that refers to an image continuing to appear in one's vision after the exposure to the original image has ceased.

Anisometropia. A visual condition in which there is a significant refractive difference between the two eyes.

Astigmatism. A visual condition in which the light entering the eye is distorted such that it does not focus at one single point in space.

Binocularity. The use of two eyes at the same time, where the usable visual areas of each eye overlap to produce a three-dimensional perception.

Brightness. The subjective attribute of light to which humans assign a label between very dim and very bright (brilliant). Brightness is perceived, not measured. Brightness is what is perceived when lumens fall on the rods and cones of the eye's retina. The sensitivity of the eye decreases as the magnitude of the light increases, and the rods and cones are sensitive to the luminous energy per unit of time (power) impinging on them.

Cataracts. A loss of clarity of the crystalline lens within the eye which causes partial or total blindness.

Cathode ray tube (CRT). A glass tube that forms part of most video display terminals. The tube generates a stream of electrons that strike the phosphor coated display screen and cause light to be emitted. The light forms characters on the screen.

Color convergence. Alignment of the three electron beams in the CRT that generate the three primary screen colors—red, green, and blue—used to form images on screen. In a misconverged image, edges will have color fringes (e.g., a white area might have a blue fringe on one side).

Color temperature. A way of measuring color accuracy. Adjusting a monitor's color-temperature control, for example, may change a bluish white to a whiter white.

Convergence. That visual function of realigning the eyes to attend an object closer than optical infinity. The visual axes of the eyes continually point closer to each other as the object of viewing gets closer to the viewer.

Contrast. The difference in color and light between parts of an image.

Diplopia (double vision). That visual condition where the person experiences two distinct images while looking at one object. This results from the breakdown of the coordination skills of the person.

Disability glare. A type of glare that causes objects to appear to have lower contrast and is usually caused by scattering of light within the media in the eye.

Discomfort glare. Glare that produces ocular discomfort including eye fatigue, eyestrain, and irritation.

Dot matrix. A pattern of dots that forms characters (text) or constructs a display image (graphics) on the VDT screen.

Dot pitch. The distance between two phosphor dots of the same color on the screen.

Electromagnetic radiation. A form of energy resulting from electric and magnetic effects which travels as invisible waves.

Ergonomics. The study of the relationship between humans and their work. The goal of ergonomics is to increase worker's comfort, productivity, and safety.

Eyesight. The process of receiving light rays into the eyes and focusing them onto the retina for interpretation.

Eyestrain (asthenopia). Descriptive terms for symptoms of visual discomfort. Symptoms include burning, itching, tiredness, aching, watering, blurring, etc.

Farsightedness (hyperopia). A visual condition where objects at a distance are more easily focused as opposed to objects up close.

Font. A complete set of characters including typeface, style, and size used for screen or printer displays.

Focal length. The distance from the eye to the viewed object needed to obtain clear focus.

Foot-candle. The amount of illumination inside the surface of an imaginary 1 ft radius sphere would be receiving if there were a uniform point source of 1 cd in the exact center of the sphere. One foot-candle \approx 10.764 lux.

Glare. The loss in visual performance or visibility, or the annoyance of discomfort, produced by a luminance in the visual field greater than the illuminance to which the eyes are adapted.

Hertz (Hz). Cycles per seconds. Used to express the refresh rate of video displays.

Illuminance. The luminous flux incident on a surface per unit area. The unit is the lux or lumen per square meter. The foot-candle (fc) or lumen per square foot is also used. An illuminance photometer measures the luminous flux per unit area at the surface being illuminated without regard to the direction from which the light approaches the sensor.

Interlaced. An interlaced monitor scans the odd lines of an image first followed by the even lines. This scanning method does not successfully eliminate flicker on computer screens.

Lag. In optometric terms, the measured difference between the viewed object and the actual focusing distance.

LASIK. A surgical procedure that alters the curvature of the cornea (front surface of the eye) which reduces the amount of nearsightedness or astigmatism.

LCD (Liquid crystal display). A display technology that relies on polarizing filters and liquid-crystal cells rather than phosphors illuminated by electron beams to produce an on-screen image. To control the intensity of the red, green, and blue dots that comprise pixels, an LCD's control circuitry applies varying charges to the liquid-crystal cells through which polarized light passes on its way to the screen.

Light. The radiant energy that is capable of exciting the retina and producing a visual sensation. The visible wavelengths of the electromagnetic spectrum extend from about 380 to 770 nm. The unit of light energy is the lumen.

Luminous flux. The visible power or light energy per unit of time. It is measured in lumens. Since light is visible energy, the lumen refers only to visible power.

Luminous intensity. The luminous flux per solid angle emitted or reflected from a point. The unit of measure is the lumen per steradian, or candela (cd). (The steradian is the unit of measurement of a solid angle.)

Luminance. The luminous intensity per unit area projected in a given direction. The unit is the candela per square meter, which is still sometimes called a nit. The foot-lambert (fL) is also in common use. Luminance is the measurable quantity which most closely corresponds to brightness.

Lux (see foot-candle). A unit of illuminance and luminous emittance. It is used in photometry as a measure of the intensity of light.

MHz (MegaHertz). A measurement of frequency in millions of cycles per second.

Myopia (nearsightedness). The ability to see objects clearly only at a close distance.

Musculoskeletal. Relating to the muscles and skeleton of the human body.

Nearpoint. The nearest point of viewing, usually within arms length.

Noninterlaced. A noninterlaced monitor scans the lines of an image sequentially, from top to bottom. This method provides less visible flicker than interlaced scanning.

Ocular motility. Relating to the movement abilities of the eyes.

Parabolic louver. A type of light fixture that is designed to direct light in a limited and narrowed direction.

Perception. The understanding of sensory input (vision, hearing, touch, etc.).

Pixel. The smallest element of a display screen that can be independently assigned color and intensity.

Phosphor. A substance that emits light when stimulated by electrons.

Photophobia. A visual condition of excessive sensitivity to light.

Presbyopia. A reduction in the ability to focus on near objects caused by the decreased flexibility in the lens, usually noticed around the age of 40 years old or later.

Polarity. The arrangement of the light and dark images on the screen. Normal polarity has light characters against a dark background; reverse polarity has dark characters against a light background.

Refractive. Having to do with the bending of light rays, usually in producing a sharp optical image.

Refresh rate. The number of times per second that the screen phosphors must be painted to maintain proper character display.

Resolution. The number of pixels, horizontally and vertically, that make up a screen image. The higher the resolution, the more detailed the image.

Resting point of accommodation (RPA). The point in space where the eyes naturally focus when at rest.

Specular reflection. The perfect, mirror-like reflection of light from a surface, in which light from a single incoming direction is reflected into a single outgoing direction.

Suppression. The “turning off” of the image of one eye by the brain, most often to avoid double vision or reduce excess stress.

SVGA (Super video graphics array). A video adapter capable of higher resolution pixels and/or colors) than the $320 \times 200 \times 256$ and $640 \times 480 \times 16$ which IBM's VGA adapter is capable of producing. SVGA enables video adapters to support resolutions of 1024 by 768 pixels and higher with up to 16.7 million simultaneous colors (known as true color).

VDT (Video display terminal). An electronic device consisting of a monitor unit (e.g., cathode ray tube) with which to view input into a computer.

Vision. A learned awareness and perception of visual experiences (combined with any or all other senses) that results in mental or physical action. Not simply eyesight.

Visual stress. The inability of a person to visually process light information in a comfortable, efficient manner.

Vision therapy. Treatment (by behavioral optometrists) used to develop and enhance visual abilities.

23.2 INTRODUCTION

Vision and eye problems are the most commonly reported symptoms among workers at computer monitors. The percentage of computer workers who experience visually related symptoms is generally reported to be from 70 up to almost 90 percent.¹ The current concern for, and popularization of these problems is likely related to the rapid introduction of this new technology into the workplace, the job restructuring associated with it, and a greater societal concern about discomfort associated with work.

The problems are largely symptoms of discomfort. There is little evidence of permanent physiological change or damage associated with extended work at a computer monitor. The vision and eye problems caused by working at a computer monitor have been collectively named “Computer Vision Syndrome (CVS)” by the American Optometric Association. The most common symptoms of CVS obtained from a survey of optometrists in order of frequency are: eyestrain, headaches, blurred vision, dry or irritated eyes, neck and/or backaches, photophobia, double vision, and colored afterimages.

A syndrome is defined as a “group of signs and symptoms that indicate a disease or disease process.”² While CVS is not technically a “disease,” it can also be properly characterized as a syndrome due to the fact that it is “a complex of symptoms indicating the existence of undesirable condition or quality.” Thus, CVS is not a single disease entity, as the definition seems to require. However, it is a group of signs and symptoms related to a specific occupational usage of the eyes.

The causes for these symptoms are a combination of individual vision problems, poor office ergonomics, or work habits. Many individuals have marginal vision disorders that do not cause symptoms of less demanding visual tasks. On the other hand, there are numerous aspects of the computer monitor and the computer work environment that make it a more demanding visual task than others—therefore, more individuals are put beyond their threshold for experiencing symptoms.

23.3 WORK ENVIRONMENT

Lighting

Improper lighting is likely the largest environmental factor contributing to visual discomfort. The room lighting is a particular problem for computer workers because the horizontal gaze angle of computer work exposes the eyes to numerous glare sources such as overhead lights and windows. Other deskwork is performed with a downward gaze angle and the glare sources are not in the field of view.

Glare is the loss in visual performance or visibility, or the annoyance of discomfort, produced by a luminance in the visual field greater than the illuminance to which the eyes are adapted. There are generally four types of glare: distracting glare, discomfort glare, disabling glare, and blinding glare.

Distracting glare results from light being reflected from the surface of an optical medium and is usually below 3000 lumens. This form of glare most often results in an annoyance to the viewer and leads to eye fatigue. Discomfort glare ranges from 3000 to 10,000 lumens and produces ocular discomfort, including eye fatigue, eyestrain, and irritation. Disabling glare, also known as veiling glare, results when light reaches 10,000 lumens and can actually block visual tasks. This type of glare causes objects to appear to have lower contrast and is caused by scattering of light within the media in the

TABLE 1 Display and Glare Source Luminances

Visual Object	Luminance (cd/m ²)
Dark background display	20–25
Light background display	80–120
Reference material with 750 lumens/m ²	200
Reference material with auxiliary light	400
Blue sky (window)	2,500
Concrete in sun (window)	6,000–1,2000
Fluorescent light (poor design)	1,000–5,000
Auxiliary lamp (direct)	1,500–10,000

eye. Blinding glare results from incident light reflecting from smooth shiny surfaces such as water and snow. It can block vision to the extent that the wearer becomes visually compromised and recovery time is needed to be fully comfortable once again.

The threshold luminance ratios and locations of visual stimuli that cause glare discomfort have been determined,³ but the physiological basis for glare discomfort is not known. Since large luminance disparities in the field of view can cause glare discomfort, it is best to have a visual environment in which luminances are relatively equal. Primarily because of glare discomfort, the Illumination Engineering Society of North America (IESNA) established maximum luminance ratios that should not be exceeded.⁴ The luminance ratio should not exceed 1:3 or 3:1 between the task and visual surrounding within 25°, nor should the ratio exceed 1:10 or 10:1 between the task and more remote visual surroundings. Many objects in the field of view can cause luminance ratios in excess of those recommended by IESNA. Table 1 shows common luminance levels of objects in an office environment. Relative to the display luminance, several objects can greatly exceed the IESNA recommended ratios. A major advantage of using light background displays compared to dark background displays is that it enables better conformity with office luminance levels.

Good lighting design can reduce discomfort glare. Light leaving the fixture can be directed so that it goes straight down and not into the eyes of the room occupants. This is most commonly accomplished with parabolic louvers in the fixture. A better solution is indirect lighting in which the light is bounced off the ceiling—resulting in a large low luminance source of light for the room. Traditional office lighting recommendations have been suggested at about 100 fc (1000 lux) but computerized offices most often require less light, in the range of 50 fc (500 lux) to better balance with the light emanating from the computer display. Proper treatment of light from windows is also important. Shades or blinds should be employed to give flexibility to control outdoor light.

Screen Reflections

In most cases, the most bothersome source of cathode ray tube (CRT) screen reflections is from the phosphor on the inner surface of the glass. (“Monitor Characteristics” are discussed later.) The phosphor is the material that emits light when struck by the electron beam—it also passively reflects room light. The reflections coming from the phosphor are diffuse. Most computer monitors have a frosted glass surface; therefore, the light reflected from the glass is primarily diffuse with a specular component.

Since most of the screen reflections are diffuse, a portion of any light impinging on the screen, regardless of the direction from which it comes, is reflected into the eyes of the user and causes the screen to appear brighter. This means that black is no longer black, but a shade of gray. The diffuse reflections from the phosphor and from the glass reduce the contrast of the text presented on the screen. Instead of viewing black characters on a white background, gray characters are presented on a white background. Calculations⁵ show that these reflections significantly reduce contrast—from 0.96 to 0.53 [contrast = $(L_t - L_b)/(L_t + L_b)$, where L_t and L_b are luminances of the task and background] under common condition. Decreased contrast increases demand upon the visual system.

A common method of treating reflections is with an antireflection filter placed on the computer monitor. The primary purpose of the filter is to make the blacks blacker, thereby increasing contrast. The luminance of light that is emitted by the screen (the desired image) is decreased by the transmittance factor of the glass filter, whereas light that is reflected from the computer screen (undesired light) is decreased by the square of the filter transmittance, since it must pass through it twice. For a typical filter of 30 percent transmittance, the luminances of black and white are reduced respectively to 9 percent and 30 percent of their values without the filter. This results in an increase in the contrast.

Some antireflection filters have a circular polarized element in them. This feature results in circular polarized lights being transmitted to the computer screen. The resulting reflection on the screen changes the rotation of the polarization of the light so that it is blocked from coming back out through the filter. Since only specular reflected light maintains its polarized properties after reflection, this polarizing feature provides added benefit for only the specular reflections. If significant specular reflections are present, it is beneficial to obtain an antireflection filter with circular polarization.

The influx of liquid crystal displays (LCDs) into the workplace has addressed some of the issues regarding reflections by using a matte surface to diffuse the reflections. However, the matte surface may be disappearing in the future; the trend in high performance LCDs is back to a glossy surface, reintroducing the problem of reflections on the surface of the display. Whichever surface, matte or glossy, users can still experience problems with glare and contrast reduction or image washout. A 2003 survey⁶ of LCD users supports this, with 39 percent reported having a glare problem on the LCD and 85 percent of LCD users reacted favorably when an antireflection filter was used with their LCD, stating that they were bothered by glare on their display and preferred working with a glare reduction filter on the display.

Independent, scientific testing of a midtransmission level antireflection computer filter to the same international standard that computer monitors must comply with, ISO 9241-7, has shown that this filter can actually improve a monitor's performance against this standard for reflection reduction and contrast improvement. The significance of this testing is that the quality of the antireflection coatings and the level of absorption technology are important considerations. There are many products on the market today that claim to be antireflection computer filters. These filters offer very low quality, if any, antireflection performance and little to no absorption technology. It is because of these lower performance products that ergonomic specialists, when considering options for reducing reflections and glare on an electronic display, often ridicule antireflection computer filters.

It is important to discuss the counterargument to antireflection computer filters—that monitors today don't need them because they already have antireflection treatments. While this is the case for many computer displays, the amount of glare and reflection reduction can be misleading. A few key points to keep in mind regarding computer displays and antireflection treatments:

- Some simply change the state of the reflection from specular to diffuse through silica coatings or etching the surface of the display (a matte finish).
- Some use a spin coating process that may only reduce first surface reflections down about 1 to 2 percent. Good quality antireflection computer filters reduce first surface reflections to less than 1 percent.
- Flat screens are better, reduce off-axis angle reflections, but still do little for normal incidence reflections.
- Few have absorptive coatings to improve contrast.
- LCDs have many advantages over the older CRT technology, matte surfaces to reduce reflections, but they still have issues with glare and contrast reduction.

Moreover, the trend of going to glossy displays means increased reflections.

Monitor Characteristics

One obvious method of improving the visual environment for computer workers is to improve the legibility of the display at which they are working. The visual image on computer monitors is compromised in many ways compared to most paper tasks.

A major aspect of computer displays is that the images are composed of pixels—the “picture element” that creates the “dots” of an image. A greater pixel density will be able to display greater detail. Most older computer screens had pixel densities in the range of from 80 to 110 dots/in (dpi: 31 to 43 dots/cm). This pixel density can be compared to impact printers (60 to 100 dpi), laser printers (300 dpi), and advanced laser printers (400 to 1200 dpi). Computer monitors today have dpi specifications that are comparable to laser printers—greater pixel density would result in visually discernible differences, as for printers. Besides pixel density there are many other factors that affect display quality, such as pixel definition (measured by area under the modulation transfer function), font type, font size, letter spacing, line spacing, stroke width, contrast, color, gray scale, and refresh rate.⁷ Several studies have shown that increased screen resolution characteristics increase reading speed and/or decrease visual symptoms.⁵ Although technological advances have resulted in improved screen resolutions, there is still substantial room for improvement.

To adequately cover all aspects of computer monitor characteristics, it is necessary to discuss both the CRT and the LCD technologies used to generate these images.

The CRT is a light-emitting device. It was first demonstrated over 100 years ago and it is a fairly simple device to describe, but one requiring great precision to manufacture. It consists of a glass bottle under a high vacuum with a layer of phosphorescent material at one end and an electron gun at the other. The electron gun creates a stream of electrons that are accelerated toward the phosphor by a very high voltage. When the electrons strike the phosphor, it glows at the point of impact. A coil of wire called the yoke is wrapped around the neck of the CRT. The yoke actually consists of a horizontal coil and a vertical coil, each of which generates a magnetic field when a voltage is applied. These fields deflect the beam of electrons both horizontally and vertically, thereby illuminating the entire screen rather than a pinpoint. To create an image, the electron beam is modulated (turned on and off) by the signal from the video card as it sweeps across the phosphor from left to right. When it reaches the right side of the tube, the beam is turned off or “blanked” and moved down one line and back to the left side. This process occurs repeatedly until the beam reaches the bottom right-hand side and is blanked and moved to the top left.

The LCD or “flat panel” as it is sometimes called, is a totally different way of generating a computer image. Rather than being light emissive as the CRT is, the LCD is light transmissive. Each dot on the screen acts like an electrically controlled shutter that either blocks or passes a very high intensity backlight (or backlights) that is always on at full brightness. Whereas the dots comprising the image of a CRT-based display vary in size depending on the displayed resolution (if you display VGA resolution, the entire screen is lit with $640 \times 480 = 307,200$ dots while if you display UXGA, the screen is filled by $1600 \times 1200 = 1,920,000$ dots), an LCD panel has a fixed number of pixels. Each pixel is made up of three subpixels, a red, a green, and a blue. Since the backlight is white, the colors are achieved by means of tiny red, green, and blue filters.

Although the properties of a liquid crystal material were discovered in the late 1800s, it was not until 1968 before the first LCD display was demonstrated. Simply, light from the backlight is polarized, passed through the liquid crystal material (which twists the light 90°) and passed through a second polarizer that is aligned with the first one. This results in a structure that blocks the backlight. The liquid crystal material has a unique property, however, that causes its molecules to line up with an applied voltage. The 90° twist goes to 0 and the polarizer sandwich becomes transparent to the backlight. By applying more or less voltage, the light can be varied from 0 (black) to 100 percent (white). Modern displays can make the transition from full off to full on in 256 steps, which is 8 bits ($2^8 = 256$). Each pixel can display 256 shades each of red, green, and blue for a total palette of $256 \times 256 \times 256 = 16.7$ million colors.

It is interesting that many monitor manufacturers will use that 16.7 million-color number very prominently in their advertising. However, in this chapter we are considering the eyes and visual system. Just how many colors does the human visual system perceive? Well, as near as we can tell (and there are many estimates here), the human visual system perceives between 7 and 8 million colors (in the most generous estimates). So, here we have a display that produces 16.7 million colors, but we can only see about half of those! For this reason, high performance monochrome monitors are often measured in “just noticeable differences,” or JNDs. However, color combinations can be a visual issue. It is important to maintain high contrast (color contrast, as well) between the letters and the background of the display, regardless of which colors are used.

So, which is the better technology? This can be very subjective, but the LCD exhibits no linearity or pincushion (straight lines looking curved) distortion and no flicker, all of which can contribute to computer vision syndrome. Furthermore, the LCD is always perfectly focused as long as it is set to display at its native resolution (the fixed number of pixels that may be addressed by the video card regardless of the resolution being displayed). When displaying a nonnative resolution, the LCD invokes a digital processor called a scaler to add enough pixels to fill the screen without distorting the image. This is a very complex procedure and the final result is a slightly softened image that many people find tiring. A good rule of thumb is to run an LCD monitor at its native resolution unless there is a compelling reason not to.

The electron beam in a CRT display is constantly scanning or drawing the image over and over. When it hits a specific spot on the screen, light is emitted, but as soon as the beam moves on, the spot starts to fade. It is critical that the beam scans all the points on the screen and returns to the first point in a short enough time that the eye does not perceive the dimming. The number of times an image is “painted” on the CRT is called the refresh rate or vertical sync frequency and it must be set to a value ≥ 75 Hz to avoid flicker. Some people are hypersensitive to flicker and need an even higher refresh rate. The LCD display does not flicker by design.

Workstation Arrangement

It is commonly stated that: “the eyes lead the body.” Since working at the computer is a visually intensive task, our body will do what is necessary to get the eyes in the most comfortable position—often at the expense of good posture and causing musculoskeletal ailments such as sore neck and back.

The most common distance at which people view printed material is 40 cm from the eyes. This is the distance at which eye doctors routinely perform near visual testing and for which most bifocal or multifocal glasses are designed. Most commonly, the computer screen is farther from the eyes—from 50 to 70 cm. This distance is largely dictated by other workstation factors, such as desk space and having room for a keyboard. The letter sizes on the screen are commensurately larger compared to printed materials to enable the longer viewing distance.

Because we are now using a vertical, rather than horizontal work surface, the height of the computer screen is a very important aspect of the workstation arrangement. A person can adapt to the vertical location of their task by changing their gaze angle (the elevation of the eyes in the orbit) and/or by changing the extension/flexion of the neck. Typically, users will alter their head position rather than eye position to adjust to a different viewing height. This can cause awkward posture and result in neck and/or backache. The eyes work best with a depression of from 10 to 20°;⁵ therefore, this is the preferred viewing angle. As a practical guide, it is often recommended that the computer user set their monitor just below their straight-ahead gaze, with the top of the display tilted back about 10°.

Whatever is viewed most often during daily work should be placed straight in front of the worker when seated at the desk. This applies to the computer and/or the reference documents—whichever is viewed most frequently. Although this seems self-evident, many computer workers situate their work so that they are constantly looking off to one side. If using hard copy on a regular basis, the material should be placed so that it is close in proximity to the monitor. Whether it is below the monitor (an “in-line” setup) or adjacent to it, the closer proximity will allow easier eye movements between the screen and the hard copy. In addition, the hard copy should reside on the same side as the handedness of the user: right side for a right-handed person. This will allow the person to either write or manipulate the hard copy without stretching into an awkward posture.

Work Habits

Since the computer has essentially taken over as the most essential piece of office equipment, office workers are spending an excessive number of hours viewing their display screens. To simplify a method to break up their viewing habits, one author has devised a method called “the 3 B’s”:

blink, breathe, and break. The rationale for blinking will be discussed in the next section on dry eyes. Breathing is an important aspect of computer use simply because we are sedentary for extended hours and our physical activity is reduced. This reduces blood flow and creates a condition for shallow breathing, which in turn tires the computer user making them lethargic and less productive.

Regarding breaks, this author recommends the “20/20/20” rule: Every 20 min, take just 20 s, and look 20 ft away. This short, more frequent break will allow the accommodative and convergence systems to relax and presumably regain some level of performance.

23.4 VISION AND EYE CONDITIONS

Many individuals have marginal vision disorders that do not cause symptoms on less demanding visual work, but will cause symptoms when the individual performs a demanding visual task. Given individual variation in visual systems and work environments, individual assessment is required to solve any given person’s symptoms. Following are the major categories of eye and vision disorders that cause symptoms among computer users.

Dry Eyes

Computer users commonly experience symptoms related to dry eyes. These symptoms include irritated eyes, dry eyes, excessive tearing, burning eyes, itching eyes, and red eyes. Contact lens wearers also often experience problems with their contact lenses while working at a computer display—related to dry eyes. In response to dry eye and ocular irritation, reflex tearing sometimes occurs and floods the eyes with tears. This is the same reflex that causes tearing when we cry or are exposed to foreign particles in our eye.

Computer workers are at greater risk for experiencing dry eye because the blink rate is significantly decreased and, because of the higher gaze angle, the eyes are wide open with a large exposed ocular surface. Patel et al.⁸ measured blink rate by direct observation on a group of 16 subjects. The mean blink rate during conversation was 18.4 blinks/min and during computer use it was 3.6 blinks/min—more than a fivefold decrease. Tsubota and Nakamori⁹ measured blink rates on 104 office workers. The mean blink rates were 22 blinks/min under relaxed conditions, 10 blinks/min while reading a book on the table, and 7 blinks/min while viewing text on a computer screen. Although both book reading and computer work result in significantly decreased blink rates, a difference between them is that computer work usually requires a higher gaze angle, resulting in an increased rate of tear evaporation. Tsubota and Nakamori⁹ measured a mean exposed ocular surface of 2.2 cm² while subjects were relaxed and 1.2 cm² while working at a computer. Since the primary route of tear elimination is through evaporation, and the amount of evaporation is a roughly linear function of ocular aperture area, the higher gaze angle when viewing a computer display results in faster tear loss. Even though reading a book and display work both reduce the blink rate, the computer worker is more at risk because of the higher gaze angle.

Other factors may also contribute to dry eye. For example, the office air environment is often low in humidity and can contain contaminants. Working under an air vent also increases evaporation. The higher gaze angle results in a greater percentage of blinks that are incomplete, resulting in poor tear film and higher evaporation rates.

Dry eyes are common in the population, reported as being 15 percent of a clinical population. For example, postmenopausal women, arthritis sufferers, post-LASIK patients, those taking some systemic medications, and contact lens wearers are more prone to dry eyes. All of these conditions should be explored as causative for dry eye symptoms. For many people with marginal dry eye problems, work at a computer will cause their dry eye problem to become clinically significant.

Resolving dry eye symptoms may take many approaches, depending on the source of the problem. The longstanding, traditional approach of simply adding an “artificial tear” to the eyes is losing favor because it has shown to be ineffective. Newer drops are available that create a better quality tear film.

In addition, a prescription drop that incorporates a mild antibiotic to reduce inflammation has shown to be effective in some cases. There is also a swell of attention to oral supplementation to stimulate better tear formation. Regardless of which approach the eye care practitioner takes, resolving dry eyes for the computer user should be a priority in the treatment of computer vision issues.

Refractive Error

Many patients simply need an accurate correction of refractive error (myopia, hyperopia, or astigmatism). The blur created by the refractive error makes it difficult to easily acquire visual information at the computer display. This reduces work efficiency and induces fatigue. A 2004 study by Daum¹⁰ indicates that even a slightly inaccurate vision prescription at the computer can have a significant negative impact on a worker's productivity. Results indicate that uncorrected vision, even if no symptoms appear, can affect employee productivity and accuracy. In addition, a vision miscorrection of as little as 0.50 D can cause a decrease in productivity of 9 percent and accuracy by 38 percent.

Patients with hyperopia must exert added accommodative effort to see clearly at near distances. Therefore, many hyperopic computer patients require refractive correction that they may not require for a less visually demanding job. Patients with 2.00 to 3.50 D of myopia who habitually read without their glasses often have visual and musculoskeletal difficulties at the computer because, to see clearly without their glasses, they must work too closely to the computer screen. These patients often require a partial correction of their myopia for proper function at the computer display.

Some computer-using patients develop a late-onset myopia from 0.25 to 1.00 D. The preponderance of evidence supports the contention that near workplaces some people are greater risk for the development of myopia. However, there is no evidence that work at a computer monitor causes myopia more than other forms of extended near-point work. The fact is, however, that our society is dictating that we do more and more of our work on the computer, which translates into long hours viewing the display screen. Take for example some school districts that require all of their students to perform all of their work on computers. The consequences of this type of visual stress may take years to realize.

Accommodation

Accommodation is the mechanism by which the eye changes its focus to look at near objects. Accommodation is accomplished by constriction of the ciliary muscle that surrounds the crystalline lens within the eye. The amplitude of accommodation (dioptric change in power) decreases with age. Presbyopia (see later) is the age-related condition that results when the amplitude of accommodation is no longer adequate to meet near visual needs. (See Chaps. 12 and 14.)

Many workers have reduced amplitude of accommodation for their age, or accommodative *infacility* (inability to change the level of accommodation quickly and accurately). These conditions result in a blur at near working distances and/or discomfort. Extended near work also commonly results in *accommodative hysteresis*—that is, the accommodative mechanism becomes “locked” into the near focus and it takes time (a few minutes to hours) to fully relax for distance focus. This is effectively a transient myopia that persists after extended near work. Some theories on myopia development look at this transient myopia as a step in the progression to developed myopia. Accommodative disorders are diagnosed in approximately one-third of the prepresbyopic segment of a clinical population. Accommodative disorders in the younger or prepresbyopic patient are usually treated with near prescription spectacles that enable the worker to relax accommodation. Vision training can also work for some patients.

Most work with display screens are performed with a higher gaze than other near-point work and accommodative amplitude is shown to be reduced with elevation of the eye. Relative to a 40° downward viewing angle, elevations of 20° downward, straight ahead, and 20° upward result in average accommodative decreases of 1.22 D, or 2.05 D, and 2.00 D, respectively, in a group of 80 prepresbyopes

with an average age of 26.4 years.¹¹ The higher gaze angles at most computer workstations result in a viewing condition for which the amplitude of accommodation is reduced—thus placing a greater strain on accommodation than near tasks performed at lower gaze angles.

Binocular Vision

Approximately 95 percent of people keep both eyes aligned on the object of regard. Those individuals who habitually do not keep their eyes aligned have *strabismus*. Even though most people can keep their eyes aligned when viewing an object, many individuals have difficulty maintaining this ocular alignment and experience symptoms such as fatigue, headaches, blur, double vision, and general ocular discomfort (see Chap. 13).

Binocular fusion is the sensory process by which the images from each eye are combined to form a single percept. When the sensory feedback loop is opened (e.g., by blocking an eye), the eyes assume their *position of rest* with respect to one another. If the position of rest is outward or diverged, the patient has *exophoria*. If it is inward or converged, the condition is *esophoria*. Clinically, the phoria is measured by occluding one of the eyes and measuring the eye alignment while occluded.

If the patient has a phoric deviation, as is the case for most people, then a constant neuromuscular effort is required to keep the eyes aligned. Whether a person experiences symptoms depends on the amount of the misalignment, the ability of the individual to overcome that misalignment, and the task demands. The symptoms associated with excessive phoria deviations can be eyestrain, double vision, headaches, eye irritation, and general fatigue.

Eye alignment at near viewing distances is more complex than at far viewing distances because of the ocular convergence required to view near objects and because of the interaction between the ocular convergence and accommodative mechanisms. Treatment of these conditions can include refractive or prismatic correction in spectacles, or vision training.

Anisometropia

One of the possible causes of binocular vision dysfunction is a difference in the correction for each eye. If there is a significant refractive difference between the two eyes, the condition of *anisometropia* exists. While this can happen developmentally and genetically, it can also be induced with surgical intervention of cataracts. If one eye has a cataract procedure performed and surgery on the other eye is not needed, then there may be a difference in the refractive errors, creating an anisometropic condition.

This condition is more significant if the refractive error of each eye is corrected with spectacles. Wearing a correction for two variable prescriptions can create an image-size differential between the eyes. One solution is to fit the patient with contact lenses. Wearing the prescription directly on the eye will reduce the magnification differential and allow for equal-sized images.

Presbyopia

Presbyopia is the condition in which the normal age-related loss of accommodation results in an inability to comfortably maintain focus on near objects. This usually begins at about the age of 40. The usual treatment for presbyopia is to prescribe reading glasses or multifocal lenses that have a distance vision corrective power (if needed) in the top of the lens and the near vision corrective power in the bottom of the lens. The most common lens designs for correcting presbyopia are bifocals and progressive addition lenses (PALs). As usually designed and fitted, these lenses work well for the most common everyday visual tasks and provide clear vision at 30 cm with a downward gaze angles of about 25°. The computer screen is typically farther away (from 50 to 70 cm) and higher (from 10 to 20° of ocular depression). A presbyope who tries to wear their usual multifocal correction at the computer will either not see the screen clearly or will need to assume an awkward

posture, resulting in neck and back strain. This is because the zone of intermediate vision is most often too narrow for full-time computer use. Many, if not most, presbyopic computer workers require a separate pair of spectacles for their computer work. Several newer lenses are designed specifically for people with occupations that require intermediate viewing distances. (See Chaps. 12 and 14.)

23.5 REFERENCES

1. National Institute for Occupational Safety and Health, "Potential Health Hazards of Video Display Terminals," *DHHS (NIOSH) Publication No. 81-129*, National Institute of Occupational Safety and Health, Cincinnati, 1981.
2. *The American Heritage Dictionary*, 4th ed., Houghton-Mifflin, Boston, MA, 2000.
3. S. K. Guth, "Prentice Memorial Lecture: The Science of Seeing—A Search for Criteria," *American Journal of Optometry and Physiological Optics* **58**:870–885 (1981).
4. "VDT Lighting. IES Recommended Practice for Lighting Offices Containing Computer Visual Display Terminals," Illumination Engineering Society of North America, New York, 1990.
5. J. E. Sheedy, "Vision at Computer Displays," *Vision Analysis*, Walnut Creek, CA, 1995.
6. T. Allan, "Glare Screen Home Usage Test Report," Decision Analysis, Inc. Study 2003–0559 (2005).
7. ANSI/HFS 100, *American National Standard for Human Factors Engineering of Visual Display Terminal Workstations*, Human Factors Society, Santa Monica, CA, 1999.
8. S. Patel, R. Henderson, L. Bradley, B. Galloway, and L. Hunter, "Effect of Visual Display Unit Use on Blink Rate and Tear Stability," *Optometry and Vision Science* **68**(11):888–892 (1991).
9. K. Tsubota and K. Nakmori, "Dry Eyes and Video Display Terminals. Letter to Editor," *New England Journal of Medicine* **328**:524 (1993).
10. Daum KM, Clore KA, Simms SS, Vesely JW, Dwilczek DD, Spittle BM, and Good GW. "Productivity Associated with Visual Status of Computer Users," *Optometry* **75**:33–47 (2004).
11. P. H. Ripple, "Accommodative Amplitude and Direction of Gaze," *American Journal of Ophthalmology* **35**:1630–1634 (1952).

HUMAN VISION AND ELECTRONIC IMAGING

Bernice E. Rogowitz

*IBM T. J. Watson Research Center
Hawthorne, New York*

Thrasyvoulos N. Pappas

*Department of Electrical and Computer Engineering
Northwestern University
Evanston, Illinois*

Jan P. Allebach

*Electronic Imaging Systems Laboratory
School of Electrical and Computer Engineering
Purdue University
West Lafayette, Indiana*

24.1 INTRODUCTION

The field of *electronic imaging* has made incredible strides over the past decade as increased computational speed, bandwidth, and storage capacity have made it possible to perform image computations at interactive speeds. This means larger images, with more spatial, temporal, and chromatic resolution, can be captured, compressed, transmitted, stored, rendered, printed, and displayed. It also means that workstations and PCs can accommodate more complex image and data formats, more complex operations for analyzing and visualizing information, more advanced interfaces, and richer image environments, such as virtual reality. This, in turn, means that image technology can now be practically used in an expanding world of applications, including video, home photography, internet catalogues, digital libraries, art, and scientific data analysis. These advances in technology have been greatly influenced by research in human perception and cognition, and in turn, have stimulated new research into the vision, perception, and cognition of the human observer. Some important topics include spatial, temporal, and color vision, attentive and preattentive vision, pattern recognition, visual organization, object perception, language, and memory.

The study of the interaction between human vision and electronic imaging is one of the key growth areas in *imaging science*. Its scope ranges from printing and display technologies to image processing algorithms for image rendering and compression, to applications involving interpretation, analysis, visualization, search, design, and aesthetics. Different electronic imaging applications call on different human capabilities. At the bottom of the visual food chain are the visual phenomena mediated by the threshold sensitivity of low-level spatial, temporal, and color mechanisms. At the next level are perceptual effects, such as color constancy, suprathreshold pattern, and texture analysis.

Moving up the food chain, we find cognitive effects, including memory, semantic categorization, and visual representation, and moving to the next level, we encounter aesthetic and emotional aspects of visual processing.

Overview of This Chapter

In this chapter, we review several key areas in the two-way interaction between human vision and technology. We show how technology advances in electronic imaging are increasingly driven by methods, models, and applications of vision science. We also show how advances in vision science are increasingly driven by the rapid advances in technologies designed for human interaction. An influential force in the development of this new field has been the Society for Imaging Science and Technology (IS&T)/Society of Photographic and Instrumentation Engineers (SPIE) Conference on Human Vision and Electronic Imaging, which had its origins in 1988 and since 1989 has grown annually as a forum for multidisciplinary research in this area.^{1–14} This chapter has been strongly influenced by the body of research that has been presented at these conferences, and reflects their unique perspectives.

A decade ago, the field of *human vision and electronic imaging* focused on the threshold sensitivity of the human visual system (HVS) as it relates to display technology and still-image compression. Today, the scope of human vision and electronic imaging is much larger and keeps expanding with the electronic imaging field. We have organized this chapter to reflect this expanding scope. We begin with early vision approaches in Sec. 24.2, showing how the explicit consideration of human spatial, temporal, and color sensitivity has affected the development of algorithms for compression and rendering, as well as the development of image quality metrics. These approaches have been most successful in algorithmically evaluating the degree to which artifacts introduced by compression or rendering processes will be detected. Section 24.3 considers how image features are detected, processed, and perceived by the human visual system. This work has been influential in the design of novel gaze-dependent compression schemes, new visualization and user interface designs, and perceptually based algorithms for image retrieval systems. Section 24.4 moves higher up the food chain to consider emotional and aesthetic evaluations. This work has influenced the development of virtual environments, high-definition TV, and tools for artistic appreciation and analysis. Section 24.5 concludes the chapter paper, and suggestions for further reading are given in Sec. 24.6.

24.2 EARLY VISION APPROACHES: THE PERCEPTION OF IMAGING ARTIFACTS

There are many approaches to characterizing image artifacts. Some approaches are based purely on physical measurement. These include measuring key image or system parameters to ensure that they fall within established tolerances, or comparing an original with a rendered, displayed, or compressed version, on a pixel-by-pixel basis, using a metric such as mean square error. These approaches have the advantage of being objective and relatively easy to implement; on the other hand, it is difficult to generalize from these data. Another approach is to use human observers to judge perceived quality, either by employing a panel of trained experts or by running psychological experiments to measure the perception of image characteristics. These approaches have the advantage of considering the human observer explicitly; on the other hand, they are costly to run and the results may not generalize. A third approach is to develop metrics, based on experiments measuring human visual characteristics, that can stand in for the human observer as a means of estimating human judgments. These perceptual models are based on experiments involving the detection, recognition, and identification of carefully controlled experimental stimuli. Since these experiments are designed to reveal the behavior of fundamental mechanisms of human vision, their results are more likely to generalize. This section reviews several models and metrics based on early human vision which have been developed for evaluating the perception of image artifacts.

Early vision models are concerned with the processes mediating the threshold detection of spatial, temporal, and color stimuli. In the early days of television design, Schade¹⁵ and his colleagues introduced the notion of characterizing human threshold sensitivity in terms of the response to spatial-frequency patterns, thereby beginning a long tradition of work to model human spatial and temporal sensitivity using the techniques of linear systems analysis. The simplest model of human vision is a threshold contrast sensitivity function (CSF). A curve representing our sensitivity to spatial modulations in luminance contrast is called a spatial CSF. A curve representing our sensitivity to temporal modulations in luminance is a temporal CSF. These band-pass curves vary depending on many other parameters, including, for example, the luminance level, the size of the field, temporal modulation, and color. In early applications of visual properties to electronic imaging technologies, these simple threshold shapes were used. As the field has progressed, the operational model for early vision has become more sophisticated, incorporating interactions between spatial, temporal, and color vision, and making more sophisticated assumptions about the underlying processes. These include, for example, the representation of the visual system as a set of band-pass spatial-frequency filters,^{16,17} the introduction of near-threshold contrast masking effects, and nonlinear processing.

One key area where these models of early vision have been applied is in the evaluation of image quality. This includes the psychophysical evaluation of image quality, perceptual metrics of image distortion, perceptual effects of spatial, temporal, and chromatic sampling, and the experimental comparison of compression, sampling, and halftoning algorithms. This work has progressed hand-in-hand with the development of vision-based algorithms for still image and video compression, image enhancement, restoration and reconstruction, image halftoning and rendering, and image and video quantization and display.

Image Quality and Compression

The basic premise of the work in perceptual image quality is that electronic imaging processes, such as compression and halftoning, introduce distortions. The more visible these distortions, the greater the impairment in image quality. The human vision model is used to evaluate the degree to which these impairments will be detected. Traditional metrics of image compression do not incorporate any models of human vision and are based on the mean squared error (i.e., the average squared difference between the original and compressed images). Furthermore, they typically fail to include a calibration step, or a model of the display device, thereby providing an inadequate model of the information presented to the eye.

The first perceptually based image quality metrics used the spatial contrast sensitivity function as a model of the human visual system.^{18,19} The development of perceptual models based on multiple spatial-frequency channels greatly improved the objective evaluation of image quality. In these models, the visual system is treated as a set of spatial-frequency-tuned channels, or as Gabor filters with limited spatial extent distributed over the visual scene. The envelope of the responses of these channels is the contrast sensitivity function. These multiple channel models provide a more physiologically representative model of the visual system, and more easily model interactions observed in the detection of spatially varying stimuli. One important interaction is contrast masking, where the degree to which a target signal is detected depends on the spatial-frequency composition of the masking signal.

In 1992, Daly introduced the *visual differences predictor*,^{6,20} a multiple-channel model for image quality that models the degree to which artifacts in the image will be detected, and thus will impair perceived image quality. This is a spatial-frequency model which incorporates spatial contrast masking and light adaptation. At about the same time, Lubin²¹ proposed a similar metric that also accounts for sensitivity variations due to spatial frequency and masking. It also accounts for fixation depth and image eccentricity in the observer's visual field. The output of such metrics is either a map of detection probabilities or a point-by-point measure of the distance between the original and degraded image normalized by the HVS sensitivity to error at each spatial frequency and location. These detection probabilities or distances can be combined into a single number that represents the overall picture quality. While both models were developed for the evaluation of displays and high quality imaging systems, they have been adapted for a wide variety of applications.

Perceptually based image compression techniques were developed in parallel with perceptual models for image quality. This is not surprising since quality metrics and compression algorithms are closely related. The image quality metric is trying to characterize the human response to an image; an image compression algorithm is either trying to minimize some distortion metric for a given bit rate, or trying to minimize the bit rate for a given distortion. In both cases, a perceptually based distortion metric can be used.

In 1989, three important papers introduced the notion of “perceptually lossless” compression. (See Refs. 22–24.) In this view, the criterion of importance in image compression was the degree to which an image could be compressed without the user’s perceiving a difference between the compressed and original images. Safranek and Johnston²⁵ presented the perceptual subband image coder (PIC) which incorporated a perceptual model and achieved perceptually lossless compression at lower rates than state-of-the-art perceptually lossy schemes. The Safranek-Johnston coder used an empirically derived perceptual masking model that was obtained for a given CRT display and viewing conditions. As with the quality metrics we discussed earlier, the model determines the HVS sensitivity to errors at each spatial frequency and location, which is called the *just noticeable distortion level* (JND). In subsequent years, perceptual models were used to improve the results of traditional approaches to image and video compression, such as those that are based on the discrete cosine transform (DCT). (See Refs. 26 to 28.)

Perceptually based image quality metrics have also been extended to video. In 1996, Van den Branden Lambrecht and Verscheure²⁹ described a video quality metric that incorporates spatio-temporal contrast sensitivities as well as luminance and contrast masking adjustments. In 1998, Watson extended his DCT-based still image metric, proposing a video quality metric based on the DCT.^{30,31} Since all the current video coding standards are based on the DCT, this metric is useful for optimizing and evaluating these coding schemes without significant additional computational overhead.

In recent years, significant energy has been devoted to comparing these methods, models, and results, and to fine-tuning the perceptual models. The Video Quality Experts Group (VQEG), for example, has conducted a cross-laboratory evaluation of perceptually based video compression metrics. They found that different perceptual metrics performed better on different MPEG-compressed image sequences, but that no one model provided a clear advantage, including the nonperceptual measure, peak signal-to-noise ratio (Ref. 32). We believe that the advantages of the perceptual metrics will become apparent when this work is extended to explicitly compare fundamentally different compression schemes over different rates, image content, and channel distortions. An important approach for improving low-level image quality metrics is to provide a common set of psychophysical data to model. Modelfest (Ref. 33) is a collaborative modeling effort where researchers have volunteered to collect detection threshold data on a wide range of visual stimuli, under carefully controlled conditions, in order to provide a basis for comparing the predictions of early vision models. This effort should lead to a converged model for early vision that can be used to develop image quality metrics and perceptually based compression and rendering schemes. (Comprehensive reviews of the use of perceptual criteria for the evaluation of image quality can be found in Refs. 34 and 35.)

Image Rendering, Halftoning, and Other Applications

Another area where simple models of low-level vision have proven to be quite successful is image halftoning and rendering. All halftoning and rendering algorithms make implicit use of the properties of the HVS; they would not work if it were not for the high spatial-frequency cutoff of the human CSF. In the early 1980s, Allebach introduced the idea of halftoning based on explicit human visual models and models of the display device. However, the development of halftoning techniques that rely on such models came much later. In 1991, Sullivan et al.³⁶ used a CSF model to design halftoning patterns of minimum visibility and Pappas and Neuhoff³⁷ incorporated printer models in error diffusion. In 1992, methods that use explicit visual models to minimize perceptual error in image halftoning were independently proposed by Analoui and Allebach,³⁸ Mulligan and Ahumada,³⁹ and Pappas and Neuhoff.⁴⁰ Allebach et al. extended this model-based approach to color in 1993⁴¹ and to video in 1994.⁴² Combining the spatial and the chromatic properties of early vision, Mulligan⁴³ took advantage of the low spatial-frequency sensitivity of chromatic channels to hide high spatial-frequency halftoning artifacts.

Low-level perceptual models and techniques have also been applied to the problem of target detection in medical images (e.g., Ref. 44) and to the problem of embedding digital watermarks in electronic images (Refs. 45 and 46).

Early Color Vision and Its Applications

The trichromacy of the human visual system has been studied for over a hundred years, but we are just in the infancy in applying knowledge of early color vision to electronic imaging systems. From an engineering perspective, the CIE 1931 model, which represents human color-matching as a linear combination of three color filters, has been the most influential model of human color vision. This work allows the determination of whether two patches will have matching colors, when viewed under a constant illuminant. The CIE 1931 color space, however, is not perceptually uniform. That is, equal differences in chromaticity do not correspond to equal perceived differences. To remedy this, various transformations of this space have been introduced. Although CIE $L^*a^*b^*$ and CIE $L^*u^*v^*$ are commonly used as perceptually uniform spaces, they still depart significantly from this objective. By adding a nonlinear gain control at the receptor level, Guth's ATD model^{47,48} provided a more perceptually uniform space and could model a wide range of perceptual phenomena. Adding a spatial component to a simplified version of Guth's model, Granger⁴⁹ demonstrated how this approach could be used to significantly increase the perceived similarity of original and displayed images. The observation that every pixel in an image provides input both to mechanisms sensitive to color variations and to mechanisms sensitive to spatial variations has provided other important contributions. For example, Uriegas^{50,51} used the multiplexing of color and spatial information by cortical neurons to develop a novel color image compression scheme.

Understanding how to represent color in electronic imaging systems is a very complicated problem, since different devices (e.g., printers, CRT or TFT/LCD displays, film) have different mechanisms for generating color, and produce different ranges, or gamuts, of colors. Several approaches have been explored for producing colors on one device that have the same appearance on another. Engineering research in this field goes under the name of *device independent color*, and has recently been energized by the need for accurate color rendering over the Internet. (See Ref. 52.)

Limitations of Early Vision Models for Electronic Imaging

The goal of the early vision models is to describe the phenomena of visual perception in terms of simple mechanisms operating at threshold. Pursuing this Ockham's razor approach has motivated the introduction of masking models and nonlinear summation models, and has allowed us to extend these simple models to describe the detection and recognition of higher-level patterns, such as textures and multiple-sinewave plaids. These models, however, eventually run out of steam in their ability to account for the perception of higher-level shapes and patterns. For example, the perception of a simple dot cannot be successfully modeled in terms of the response of a bank of linear spatial-frequency filters.

In image quality, early vision approaches consider an image as a collection of picture elements. They measure and model the perceived fidelity of an image based on the degree to which the pixels, or some transformation of the pixels, have changed in their spatial, temporal, and color properties. But what if two images have identical perceived image fidelity, but one has a lot of blur and little blockiness while the other has little blur and a lot of blockiness? Are they equivalent? Also, since it is not always possible to conceal these artifacts, it is important to understand how their effects combine and interact in order to minimize their objectionable and annoying effects. This has led to the development of new suprathreshold scaling methods for evaluating the perceived quality of images with multiple suprathreshold artifacts (e.g., Refs. 53 and 54), and has led to the development of new dynamic techniques to study how the annoyance of an artifact depends on its temporal position in a video sequence (Ref. 55).

In color vision, a key problem for electronic imaging is that simply capturing and reproducing the physical color of individual color pixels and patches in an image is not sufficient to describe

the perceived colors in the image. Hunt,⁵⁶ for example, has shown that surrounding colors affect color appearance. Furthermore, as McCann⁵⁷ has demonstrated using spatially complex stimuli, the perception of an image depends, not simply on the local luminance and color values, and not simply on nearby luminance and color values, but on a more global consideration of the image and its geometry.

24.3 HIGHER-LEVEL APPROACHES: THE ANALYSIS OF IMAGE FEATURES

Higher-level approaches in vision are dedicated to understanding the mechanisms for perceiving more complex features such as dots, plaids, textures, and faces. One approach is to build up from early spatial-frequency models by positing nonlinear mechanisms that are sensitive to two-dimensional spatial variations, such as t-junctions (e.g., Ref. 58). Interesting new research by Webster and his colleagues⁵⁹ identifies higher-level perceptual mechanisms tuned to more complex visual relationships. Using the same types of visual adaptation techniques commonly used in early vision experiments to identify spatial, temporal, or color channel properties, they have demonstrated adaptation to complex visual stimuli, including, for example, the adaptation to complex facial distortions and complex color distributions.

Another interesting approach to understanding the fundamental building blocks of human vision is to consider the physical environment in which it has evolved, with the idea in mind that the statistics of the natural world must somehow have guided, or put constraints on, its development. This concept, originally introduced by Field,^{60,61} has led to extensive measurements of the spatial and temporal characteristics of the world's scenes, and to an exploration of the statistics of the world's illuminants and reflective surfaces. One of the major findings in this field is that the spatial amplitude spectra of natural scenes falls off as $f^{-1.1}$. More recently, Field et al.⁶² have related this function to the sparseness of spatial-frequency mechanisms in human vision. This low-level description of the global frequency content in natural scenes has been tied to higher-level perception by Rogowitz and Voss.⁶³ They showed that when the slope of the fall-off in the amplitude spectrum is in a certain range (corresponding to a fractal dimension of from 1.2 to 1.4), observers see nameable shapes in images like clouds. This approach has also been applied to understanding how the statistics of illuminants and surfaces in the world constrains the spectral sensitivities of visual mechanisms (Ref. 64), thereby influencing the mechanisms of color constancy (see also Refs. 65 and 66).

The attention literature provides another path for studying image features by asking which features in an image or scene naturally attract attention and structure perception. This approach is epitomized by Triesman's work on "preattentive" vision⁶⁷ and by Julesz' exploration of "textons."^{68,69} Their work has had a very large impact in the field of electronic imaging by focusing attention on the immediacy with which certain features in the world, or in an image, are processed. Their work both identifies image characteristics that attract visual attention and guide visual search, and also provides a paradigm for studying the visual salience of image features. An important debate in this area is whether features are defined bottom up—that is, generated by successive organizations of low-level elements such as edges, as suggested by Marr⁷⁰—or whether features are perceived immediately, driven by top-down processes, as argued by Stark.⁷¹

This section explores the application of feature perception and extraction in a wide range of electronic imaging applications, and examines how the demands of these new electronic tasks motivate research in perception. One important emerging area is the incorporation of perceptual attention or "importance" in image compression and coding. The idea here is that if we could identify those aspects of an image that attracted attention, we could encode this image using higher resolution for the areas of importance and save bandwidth in the remaining areas. Another important area is image analysis, where knowing which features are perceptually salient could be used to index, manipulate, analyze, compare, or describe an image. An important emerging application in this area is digital libraries, where the goal is to find objects in a database that have certain features, or that are similar to a target object. Another relevant application area is visualization, where the goal is to develop methods for visually representing structures and features in the data.

Attention and Region of Interest

A new idea in electronic imaging is to analyze images to identify their “regions of interest,” features that draw our attention and have a special saliency for interpretation. If we could algorithmically decompose an image into its salient features, we could develop compression schemes that would, for example, compress more heavily those regions that did not include perceptual features, and devote extra bandwidth to regions of interest. Several attempts have been made to algorithmically identify regions of interest. For example, Leray et al.⁷² developed a neural network model that incorporates both the response of low-level visual mechanisms and an attentional mechanism that differentially encodes areas of interest in the visual image, and used this to develop a compression scheme.

Stelmach and Tam⁷³ measured the eye movements of people examining video sequences and concluded that there wasn’t enough consistency across users to motivate developing a compression scheme for TV broadcast based on user eye movements. Geissler and Perry,⁷⁴ however, demonstrated that by guiding a user’s eye movements to a succession of target locations, and only rendering the image at high resolution at these “foveated” regions, the entire image appeared to have full resolution. If the system knew in advance where the eye movements would be directed, it could adjust the resolution at those locations, on the fly, thus saving bandwidth. Finding those perceptually relevant features, however, is a very difficult problem. Yarbus⁷⁵ showed that there is not a single stereotypical way that an image is viewed; the way the user’s saccades are placed on the picture depends on the task. For example, if the goal is to identify how many people there are in the picture, the set of eye movements differs from those where the goal is to examine what the people are wearing.

A most important voice in this discussion is that of Stark.^{76,77} Noton and Stark’s classic papers^{71,78} set the modern stage for measuring the eye movement paths, or “scanpaths,” of human vision. They used this methodology to explore the role of particular visual stimuli in driving attention (bottom up) versus the role of higher-level hypothesis-testing and scene checking goals (top down) in driving the pattern of visual activity. Stark and his colleagues have contributed both to our basic understanding of the processes that drive these scanpaths, and to integrating this knowledge into the development of better electronic imaging systems. Recent evidence for top-down processing in eye movements has been obtained using an eye tracking device that can be worn while an observer moves about in the world (Ref. 79). As observers perform a variety of everyday tasks, they perform what the authors call “planful” eye movements, eye movements that occur in the middle of a task, in anticipation of the subject’s interaction with an object in an upcoming task.

Image Features, Similarity, and Digital Libraries Applications

In digital libraries applications, the goal is to organize and retrieve information from a database of images, videos, graphical objects, music, and sounds. In these applications, the better these objects are indexed and organized, the more easily they can be searched. A key problem, therefore, is to identify features and attributes of importance, and to provide methods for indexing and searching for objects based on these features. Understanding which features are meaningful, or what makes objects similar, however, is a difficult problem. Methods from signal processing and computer vision can be used to algorithmically segment images and detect features. Since these databases are being designed for humans to search and navigate, it is important to understand which features are salient and meaningful to human observers, how they are extracted from complex objects, and how object features are used to judge image and object similarity. Therefore, these algorithms often incorporate heuristics gleaned from perceptual experiments regarding human color, shape, and pattern perception. For example, the work by Petkovic and his colleagues⁸⁰ includes operators that extract information about color, shape, texture, and composition. More recently, knowledge about human perception and cognition has been incorporated more explicitly. For example, Frese et al.⁸¹ developed criteria for image similarity based on a multiscale model of the human visual system, and used a psychophysical model to weight the parameters of their model. Another method for explicitly incorporating human perception has been to study how humans explicitly judge image similarity and use

this knowledge to build better search algorithms. Rogowitz et al.,⁸² for example, asked observers to judge the similarity of a large collection of images and used a multidimensional scaling technique to identify the dimensions along which natural objects were organized perceptually. Mojsilović et al.⁸³ asked observers to judge the similarity of textured designs, then built a texture retrieval system based on the perceptual results.

In this expanding area, research opportunities include the analysis of human feature perception, the operationalization of these behaviors into algorithms, and the incorporation of these algorithms into digital library systems. This includes the development of perceptual criteria for image retrieval, the creation of perceptual image similarity metrics, methods for specifying perceptual metadata for characterizing these objects, and perceptual cues for navigating through large multimedia databases. An important emerging opportunity lies in extending these methods to accommodate more complex multidimensional objects, such as three-dimensional objects and auditory patterns.

Visualization

One consequence of the expanding computer age is the creation of terabytes and terabytes of data, and an interest in taking advantage of these data for scientific, medical, and business purposes. These can be data from satellite sensors, medical diagnostic sensors, business applications, simulations, or experiments, and these data come in many different varieties and forms. The goal of visualization is to create visual representations of these data that make it easier for people to see patterns and relationships, identify trends, develop hypotheses, and gain understanding. To do so, the data are mapped onto visual (and sometimes auditory) dimensions, producing maps, bar charts, sonograms, statistical plots, etc. A key perceptual issue in this area is how to map data onto visual dimensions in a way that preserves the structure in the data without creating visual artifacts. Another key perceptual issue is how to map the data onto visual dimensions in a way that takes advantage of the natural feature extraction and pattern-identification capabilities of the human visual system. For example, how can color and texture be used to draw attention to a particular range of data values, or a departure from a model's predictions? (See Ref. 84 for an introduction to these ideas.)

One important theme in this research is the return to Gestalt principles of organization for inspiration about the analysis of complex visual information. The Gestalt psychologists identified principles by which objects in the visual world are organized perceptually. For example, objects near each other (proximity) or similar to each other (similarity) appear to belong together. If these fundamental principles could be operationalized, or if systems could be built that take advantage of these basic rules of organization, it could be of great practical importance. Some recent papers in this area include Kubovy's⁸⁵ theoretical and experimental studies, experiments by Hon et al.⁸⁶ on the interpolation and segmentation of sampled contours, and an interactive data visualization system based on Gestalt principles of perceptual organization (Ref. 87).

User Interface Design

As more and more information becomes available, and computer workstations and web browsers support more interactivity and more color, we move into a new era in interface design. With so many choices available, we now have the luxury to ask how best to use color, geometry, and texture to represent information. A key vision paper in this area has been Boynton's paper⁸⁸ on "the eleven colors which are almost never confused." Instead of measuring color discriminability, Boynton asked people to sort colors, to get at a higher-level color representation scheme, and found that the responses clustered in 11 categories, each organized by a common color name. That is, although in side-by-side comparison, people can discriminate millions of colors, when the task is to categorize colors, the number of perceptually distinguishable color categories is very small. In a related experiment, Derefeldt and Swartling⁸⁹ tried to create as many distinct colors as possible for a digital map application; they were only able to identify 30. More recently, Yendrikovski⁹⁰ extracted thousands of images from the web and used a *K*-means clustering algorithm to group the colors of the pixels. He also found a small

number of distinct color categories. This suggests that the colors of the natural world group into a small number of categories that correspond to a small set of nameable colors. This result is important for many applications that involve selecting colors to represent semantic entities, such as using colors to represent different types of tumors in a medical visualization. It is also important for image compression or digital libraries, suggesting that the colors of an image can be encoded using a minimal number of color categories. Several groups have begun developing user-interface design tools that incorporate perceptually based intelligence, guiding the designer in choices of colors, fonts, and layouts based on knowledge about color discriminability, color deficiency, color communication, and legibility. Some examples include Hedin and Derefeldt⁹¹ and Rogowitz et al.⁹²

24.4 VERY HIGH-LEVEL APPROACHES: THE REPRESENTATION OF AESTHETIC AND EMOTIONAL CHARACTERISTICS

The higher-level applications discussed above, such as visualization and digital libraries, push the envelope in electronic imaging technology. They drive developers to address the content of the data, not just the pixel representation. As the technology improves and becomes more pervasive, new application areas are drawn into the web, including applications that use visual media to convey artistic and emotional meaning. For these applications, it is not just the faithfulness of the representation, or its perceptual features, that are of concern, but also its naturalness, colorfulness, composition, and appeal. This may seem an enormous leap, until we realize how natural it is to think of the artistic and emotional aspects of photographic imaging. Interest in using electronic media to convey artistic and emotional goals may just reflect the maturation of the electronic imaging field.

Image Quality

Although the word “quality” itself connotes a certain subjectivity, the major focus of the research in image quality, discussed so far, has been aimed at developing objective, algorithmic methods for characterizing the judgments of human observers. The research in Japan on high-definition TV has played a major role in shaping a more subjective approach to image quality. This approach is based on the observation that the larger format TV, because it stimulates a greater area of the visual field, produces images which appear not only more saturated, but also more lively, realistic, vivid, natural, and compelling (Ref. 93). In related work on the image quality of color images, DeRidder et al.⁹⁴ found that, in some cases, observers preferred somewhat more saturated images, even when they clearly perceived them to be somewhat less natural. The observers’ responses were based on aesthetic characteristics, not judgments of image fidelity.

Virtual Reality and Presence

Virtual reality is a set of technologies whose goal is to immerse the user in a visual representation. Typically, the user will wear a head-tracking device to match the scene to the user’s viewpoint and stereo glasses whose two interlaced views provide a stereoscopic view. A tracked hand device may also be included so that the user can interact with this virtual environment. Although these systems do not provide adequate visual resolution or adequate temporal resolution to produce a realistic visual experience, they do give the sense of “presence,” a sense of being immersed in the environment. This new environment has motivated considerable perceptual research in a diverse set of areas, including the evaluation of depth cues, the visual contributions to the sense of presence, the role of auditory cues in the creation of a realistic virtual environment, and cues for navigation. This new environment also provides an experimental apparatus for studying the effects of large-field visual patterns and the interaction of vision and audition in navigation.

Color, Art, and Emotion

Since the beginning of experimental vision science, there have been active interactions between artists and psychologists. This is not surprising since both groups are interested in the visual impression produced by physical media. As artists use new media, and experiment with new visual forms, perceptual psychologists seek to understand how visual effects are created. For example, Chevreul made fundamental contributions to color contrast by studying why certain color dyes in artistic textiles didn't appear to have the right colors. The Impressionists and the Cubists read the color vision literature of their day, and in some cases, the artist and vision scientist were the same person. For example, Seurat studied the perception of color and created color halftone paintings. Artists have continued to inspire vision scientists to this day. For example, Papathomas,⁹⁵ struck by the impressive visual depth illusions in Hughes's wall sculptures, used this as an opportunity to deepen our knowledge of how the visual system constructs the sensation of depth. Koenderink⁹⁶ used three-dimensional sculptures to study how the human observer constructs an impression of three-dimensional shape from the set of two-dimensional views, where there is an infinite number of valid reconstructions. Tyler⁹⁷ brought psychophysical testing to the issue of portrait composition, and discovered that portrait painters systematically position one of the subject's eyes directly along the central median of the painting.

Another approach in this area is to understand the artistic process in order to emulate it algorithmically. The implementation, thus, is both a test of the model and a method for generating objects of an artistic process. For example, Burton⁹⁸ explored how young children express themselves visually and kinesthetically in drawings, and Brand⁹⁹ implemented his model of shape perception by teaching a robot to sculpt. Some of these goals are esoteric, but practical applications can be found. For example, Dalton¹⁰⁰ broadened the dialogue on digital libraries by exploring how algorithms could search for images that are the same, except for their artistic representation.

Also, as more and more artists, fabric designers, and graphic artists exploit the advances of electronic imaging systems, the more important it becomes to give them control over the aesthetic and emotional dimensions of electronic images. The MIT Media Lab has been a major player in exploring the visual, artistic, and emotional parameters of visual stimuli. In 1993, Feldman and Bender¹⁰¹ conducted an experiment where users judged the affective aspects (e.g., "energy," "expressiveness") of color pairs, demonstrating that even such seemingly vague and subjective dimensions could be studied using psychophysical techniques. They found that color energy depended on the color distance between pairs in a calibrated Munsell color space, where complementary hues, or strongly different luminances or saturations, produced greater emotional energy. These authors and their colleagues have gone on to develop emotion-based color palettes and software tools for artistic design.

24.5 CONCLUSIONS

The field of human vision and electronic imaging has developed over the last decade, evolving with technology. In the 1980s, the bandwidth available to workstations allowed only the display of green text, and in those days people interested in perceptual issues in electronic imaging studied display flicker, spatial sampling, and gray-scale/resolution trade-offs. As the technology has evolved, almost all desktop displays now provide color, and the bandwidth allows for high-resolution, 24-bit images. These displays, and the faster systems that drive them, allow the development of desktop image applications, desktop digital libraries applications, and desktop virtual reality applications. Application developers interested in art, visualization, data mining, and image analysis can now build interactive image solutions, integrating images and data from multiple sources worldwide. This rapid increase in technology enables new types of imaging solutions, solutions that allow the user to explore, manipulate, and be immersed in images; and these in turn pose new questions to researchers interested in perceptually based imaging systems.

Some of these new questions include, for example: How do we measure the image quality of virtual reality environments? How do we measure presence, and what are the factors (including interactions of multiple media) that contribute to it? How do we model the interactive visualization of

multivariate data? How do we provide environments that allow people to search through a sea of images? How do we create tools that artists and designers can use? How do we use image technology to perform remote surgery? How do we build user interfaces that sense our mood? With each advance in technology, the systems we build grow increasingly richer, and the questions we ask about human observers become more complex. In particular, they increasingly require a deeper understanding of higher levels of human perception and cognition. Low-level models of retinal and striate cortex function are the foundation, but upon this base we need to consider other dimensions of the human experience: color perception, perceptual organization, language, memory, problem solving, aesthetic appreciation, and emotional response.

As electronic imaging becomes ever more prevalent in the home and the workplace, the two-way interaction between human vision and electronic imaging technology will continue to grow in scope and importance. Research in human vision and its interactions with the other senses will continue to shape the solutions that are developed to meet the requirements of the imaging industry, which in turn will continue to motivate our understanding of the complex process of vision that is the window to the world around us.

24.6 ADDITIONAL INFORMATION ON HUMAN VISION AND ELECTRONIC IMAGING

A number of archival journals and conferences have provided a venue for dissemination and discussion of research results in the field of human vision or in the field of electronic imaging. For human vision, these include the annual meeting of the Optical Society of America, the annual meeting of the Association for Research in Vision and Ophthalmology, and the European Conference on Visual Perception (ECPV). For electronic imaging, these include the IEEE Signal Processing Society's International Conference on Image Processing, the annual meeting of the Society for Information Display (SID), and several conferences sponsored or cosponsored by the Society for Imaging Science and Technology (IS&T), especially the International Conference on Digital Printing Technologies (NIP), the Color Imaging Conference (CIC), co-sponsored with SID, the Conference on Picture and Image Processing (PICS), and several conferences that are part of the Symposium on Electronic Imaging cosponsored by the Society of Photographic and Instrumentation Engineers (SPIE) and IS&T.

The major forum for human vision *and* electronic imaging is the Conference on Human Vision and Electronic Imaging cosponsored by the International Imaging Society (IS&T) and the Society for Photographic and Imaging Engineers (SPIE).

There are a number of excellent books on various aspects of human perception and electronic imaging, including *Human Color Vision* by Kaiser and Boynton,¹⁰² *Foundations of Vision* by Wandell,¹⁰³ *Visual Perception* by Cornsweet,¹⁰⁴ *Spatial Vision* by De Valois and De Valois,¹⁰⁵ *The Senses*, edited by Barlow and Molon,¹⁰⁶ *The Artful Eye*, edited by Gregory, Harris, Rose, and Heard,¹⁰⁷ *How the Mind Works* by Pinker,¹⁰⁸ *Information Visualization: Optimizing Design For Human Perception*, edited by Colin Ware,¹⁰⁹ *Digital Image and Human Vision*, edited by Watson,¹¹⁰ *Computational Models of Visual Processing*, edited by Landy and Movshon,¹¹¹ and *Early Vision and Beyond*, edited by Papathomas et al.¹¹²

Finally, we cite a number of review articles and book chapters.^{34,35,113}

24.7 REFERENCES

1. G. W. Hughes, P. E. Mantey, and B. E. Rogowitz (eds.), *Image Processing, Analysis, Measurement, and Quality*, Proc. SPIE, vol. 901, Los Angeles, California, Jan. 13–15, 1988.
2. B. E. Rogowitz (ed.), *Human Vision, Visual Processing, and Digital Display*, Proc. SPIE, vol. 1077, Los Angeles, California, Jan. 18–20, 1989.
3. B. E. Rogowitz and J. P. Allebach (eds.), *Human Vision and Electronic Imaging: Models, Methods, and Applications*, Proc. SPIE, vol. 1249, Santa Clara, California, Feb. 12–14, 1990.

4. M. H. Brill (ed.), *Perceiving, Measuring, and Using Color*, *Proc. SPIE*, vol. 1250, Santa Clara, California, Feb. 15–16, 1990.
5. B. E. Rogowitz, M. H. Brill, and J. P. Allebach (eds.), *Human Vision, Visual Processing, and Digital Display II*, *Proc. SPIE*, vol. 1453, San Jose, California, Feb. 27–Mar. 1, 1991.
6. B. E. Rogowitz (ed.), *Human Vision, Visual Processing, and Digital Display III*, *Proc. SPIE*, vol. 1666, San Jose, California, Feb. 10–13, 1992.
7. J. P. Allebach and B. E. Rogowitz (eds.), *Human Vision, Visual Processing, and Digital Display IV*, *Proc. SPIE*, vol. 1913, San Jose, California, Feb. 1–4, 1993.
8. B. E. Rogowitz and J. P. Allebach (eds.), *Human Vision, Visual Processing, and Digital Display V*, *Proc. SPIE*, vol. 2179, San Jose, California, Feb. 8–10, 1994.
9. B. E. Rogowitz and I. P. Allebach (eds.), *Human Vision, Visual Processing, and Digital Display VI*, *Proc. SPIE*, vol. 2411, San Jose, California, Feb. 6–8, 1995.
10. B. E. Rogowitz and J. P. Allebach (eds.), *Human Vision and Electronic Imaging*, *Proc. SPIE*, vol. 2657, San Jose, California, Jan. 29–Feb. 1, 1996.
11. B. E. Rogowitz and T. N. Pappas (eds.), *Human Vision and Electronic Imaging II*, *Proc. SPIE*, vol. 3016, San Jose, California, Feb. 10–13, 1997.
12. B. E. Rogowitz and T. N. Pappas (eds.), *Human Vision and Electronic Imaging III*, *Proc. SPIE*, vol. 3299, San Jose, California, Jan. 26–29, 1998.
13. B. E. Rogowitz and T. N. Pappas (eds.), *Human Vision and Electronic Imaging TV*, *Proc. SPIE*, vol. 3644, San Jose, California, Jan. 25–28, 1999.
14. B. E. Rogowitz and T. N. Pappas (eds.), *Human Vision and Electronic Imaging V*, *Proc. SPIE*, vol. 3959, San Jose, California, Jan. 24–27, 2000.
15. O. H. Schade, “Optical and Photoelectric Analog of the Eye,” *Journal of the Optical Society of America* **46**:721–739 (1956).
16. F. W. Campbell and J. G. Robson, “Application of Fourier Analysis to the Visibility of Gratings,” *Journal of Physiology* **197**:551–566 (1968).
17. N. Graham and J. Nachmias, “Detection of Grating Patterns Containing Two Spatial Frequencies: A Comparison of Single-Channel and Multiple-Channels Models,” *Vision Research* **11**:252–259 (1971).
18. H. Snyder, “Image Quality and Observer Performance in Perception of Displayed Information,” in *Perception of Displayed Information*, L. M. Bieberman (ed.), Plenum Press, New York, 1973.
19. P. G. J. Barten, “The SQRI Method: A New Method for the Evaluation of Visible Resolution on a Display,” *Proceedings of the Society for Information Display*, **28**:253–262 (1987).
20. S. Daly, “The Visible Differences Predictor: An Algorithm for the Assessment of Image Fidelity,” in A. B. Watson (ed.), *Digital Images and Human Vision*, MIT Press, Cambridge, MA, pp. 179–206, 1993.
21. J. Lubin, “The Use of Psychophysical Data and Models in the Analysis of Display System Performance,” in A. B. Watson (ed.), *Digital Images and Human Vision*, MIT Press, Cambridge, MA, pp. 163–178, 1993.
22. V. Ramamoorthy and N. S. Jayant, “On Transparent Quality Image Coding Using Visual Models,” *Proc. SPIE*, **1077**:146–154, 1989.
23. S. Daly, “The Visible Difference Predictor: An Algorithm for the Assessment of Image Fidelity,” *Proc. SPIE*, **1077**:209–216, 1989.
24. A. B. Watson, “Receptive Fields and Visual Representations,” *Proc. SPIE*, **1077**:190–197, 1989.
25. R. J. Safranek and J. D. Johnston, “A Perceptually Tuned Sub-Band Image Coder with Image Dependent Quantization and Post-Quantization Data Compression,” in *Proc. ICASSP-89*, vol. 3, Glasgow, Scotland, pp. 1945–1948, May 1989.
26. H. A. Peterson, A. J. Ahumada, and A. B. Watson, “Improved Detection Model for DCT Coefficient Quantization,” *Proc. SPIE*, **1913**:191–201, 1993.
27. A. B. Watson, “DCT Quantization Matrices Visually Optimized for Individual Images,” *Proc. SPIE*, **1913**:202–216, 1993.
28. D. A. Silverstein and S. A. Klein, “DCT Image Fidelity Metric and its Application to a Text Based Scheme for Image Display,” *Proc. SPIE*, **1913**:229–239, 1993.

29. C. J. Van den Branden Lambrecht and O. Verscheure, "Perceptual Quality Measure Using a Spatio-Temporal Model of the Human Visual System," in V. Bhaskaran, F. Sijstermans, and S. Panchanathan (eds.), *Digital Video Compression: Algorithms and Technologies, Proc. SPIE*, vol. 2668, San Jose, California, pp. 450–461, Jan./Feb. 1996.
30. A. B. Watson, "Toward a Perceptual Visual Quality Metric," *Proc. SPIE*, **3299**:139–147, 1998.
31. A. B. Watson, Q. J. Hu, J. F. McGowan, and J. B. Mulligan, "Design and Performance of a Digital Video Quality Metric," *Proc. SPIE*, **3644**:168–174, 1999.
32. P. J. Corriveau, A. A. Webster, A. M. Rohaly, and J. M. Liebert, "Video Quality Experts Group: The Quest for Valid and Objective Methods," *Proc. SPIE*, **3959**:129–139, 2000.
33. T. Carney, C. W. Tyler, A. B. Watson, W. Makous, B. Beutter, C. Chen, A. M. Norcia, and S. A. Klein, "Modelfest: Year One Results and Plans for Future Years," *Proc. SPIE*, **3959**:140–151, 2000.
34. M. P. Eckert and A. P. Bradley, "Perceptual Quality Metrics Applied to Still Image Compression," *Signal Processing* **70**:177–200 (1998).
35. T. N. Pappas and R. J. Safranek, "Perceptual Criteria for Image Quality Evaluation," in *Handbook of Image and Video Processing*, A. C. Bovik (ed.), Academic Press, New York, pp. 669–684, 2000.
36. J. Sullivan, L. Ray, and R. Miller, "Design of Minimum Visual Modulation Halftone Patterns," *IEEE Transactions on Systems, Man, and Cybernetics* **21**:33–38 (Jan./Feb. 1991).
37. T. N. Pappas and D. L. Neuhoff, "Printer Models and Error Diffusion," *IEEE Trans. Image Process* **4**:66–79 (1995).
38. M. Analoui and J. P. Allebach, "Model Based Halftoning Using Direct Binary Search," *Proc. SPIE*, **1666**:96–108, 1992.
39. J. B. Mulligan and A. J. Ahumada, "Principled Halftoning Based on Human Vision Models," *Proc. SPIE*, **1666**:109–121, 1992.
40. T. N. Pappas and D. L. Neuhoff, "Least Squares Model Based Halftoning," *Proc. SPIE*, **1666**:165–176, 1992.
41. T. J. Flohr, B. W. Kolpatzik, R. Balasubramanian, D. A. Carrara, C. A. Bouman, and J. P. Allebach, "Model Based Color Image Quantization," *Proc. SPIE*, **1913**:270–281, 1993.
42. C. B. Atkins, T. J. Flohr, D. P. Hilgenberg, C. A. Bouman, and J. P. Allebach, "Model Based Color Image Sequence Quantization," *Proc. SPIE*, **2179**:318–326, 1994.
43. J. B. Mulligan, "Digital Halftoning Methods for Selectively Partitioning Error into Achromatic and Chromatic Channels," *Proc. SPIE*, **1249**:261–270, 1990.
44. M. P. Eckstein, A. J. Ahumada, and A. B. Watson, "Image Discrimination Models Predict Signal Detection in Natural Medical Image Backgrounds," *Proc. SPIE*, **1316**:58–69, 1997.
45. I. J. Cox and M. L. Miller, "Review of Watermarking and the Importance of Perceptual Modeling," *Proc. SPIE*, **1316**:92–99, 1997.
46. C. I. Podilchuk and W. Zeng, "Digital Image Watermarking Using Visual Models," *Proc. SPIE*, **1316**:100–111, 1997.
47. S. L. Guth, "Unified Model for Human Color Perception and Visual Adaptation," *Proc. SPIE*, **1077**:370–390, 1989.
48. S. L. Guth, "Unified Model for Human Color Perception and Visual Adaptation II," *Proc. SPIE*, **1913**:440–448, 1993.
49. E. M. Granger, "Uniform Color Space as a Function of Spatial Frequency," *Proc. SPIE*, **1913**:449–461, 1993.
50. E. M. Uriegas, J. D. Peters, and H. D. Crane, "Comparison of Digital Color Images Based the Model of Spatiochromatic Multiplexing of Human Vision," *Proc. SPIE*, **2179**:400–406, 1994.
51. E. M. Uriegas, J. D. Peters, and H. D. Crane, "Spatiotemporal Multiplexing: A Color Image Representation for Digital Processing and Compression," *Proc. SPIE*, **2657**:412–420, 1996.
52. J. Gille, J. Luszcz, and J. O. Larimer, "Error Diffusion Using the Web-Safe Colors: How Good Is it Across Platforms?" *Proc. SPIE*, **3299**:368–375, 1998.
53. H. de Ridder and G. M. Majoor, "Numerical Category Scaling: An Efficient Method for Assessing Digital Image Coding Impairments," *Proc. SPIE*, **1249**:65–77, 1990.
54. J. A. Roufs and M. C. Boschman, "Methods for Evaluating the Perceptual Quality of VDUs," *Proc. SPIE*, **1249**:2–11, 1990.
55. D. E. Pearson, "Viewer Response to Time Varying Video Quality," *Proc. SPIE*, **3299**:2–15, 1998.

56. R. W. G. Hunt, *The Reproduction of Colour in Photography, Printing and Television*. Fountain Press, England, 1987.
57. J. J. McCann, "Color Imaging System and Color Theory: Past, Present and Future," *Proc. SPIE*, **3299**:38–46, 1998.
58. C. Zetzsche and E. Barth, "Image Surface Predicates and the Neural Encoding of Two Dimensional Signal Variations," *Proc. SPIE*, **1249**:160–177, 1990.
59. M. A. Webster and O. H. MacLin, "Visual Adaptation and the Perception of Distortions in Natural Images," *Proc. SPIE*, **3299**:264–273, 1998.
60. D. J. Field, "What the Statistics of Natural Images Tell us about Visual Coding," *Proc. SPIE*, **1077**:269–276, 1989.
61. D. J. Field, "Relations Between the Statistics of Natural Images and the Response Properties of Cortical Cells," *Journal of the Optical Society of America A* **4**:2379–2394 (1987).
62. D. J. Field, B. A. Olshausen, and N. Brady, "Wavelets, Blur and the Sources of Variability in the Amplitude Spectra of Natural Scenes," *Proc. SPIE*, **2657**:108–119, 1996.
63. B. E. Rogowitz and R. Voss, "Shape Perception and Low Dimension Fractal Boundary Contours," *Proc. SPIE*, **1249**:387–394, 1990.
64. L. T. Maloney, "Photoreceptor Spectral Sensitivities and Color Correction," *Proc. SPIE*, **1250**:103–110, 1990.
65. G. D. Finlayson, "Color Constancy and a Changing Illumination," *Proc. SPIE*, **2179**:352–363, 1994.
66. D. H. Brainard and W. T. Freeman, "Bayesian Method for Recovering Surface and Illuminant Properties from Photosensor Responses," *Proc. SPIE*, **2179**:364–376, 1994.
67. A. Treisman and G. Gelade, "A Feature Integration Theory of Attention," *Cognitive Psychology* **12**:97–136 (1980).
68. B. Julesz, "AI and Early Vision II," *Proc. SPIE*, **1077**:246–268, 1989.
69. B. Julesz, "Textons, the Elements of Texture Perception and Their Interactions," *Nature* **290**:91–97 (1981).
70. D. Marr, *Vision*. W. H. Freeman Company, San Francisco, CA, 1982.
71. D. Noton and L. Stark, "Eye Movements and Visual Perception," *Scientific American* **224**(6):34–43 (1971).
72. P. Leray, F. Guyot, P. Marchal, and Y. Burnod, "CUBICORT: Simulation of the Visual Cortical System for Three Dimensional Image Analysis, Synthesis, and Hypercompression for Digital TV, HDTV and Multimedia," *Proc. SPIE*, **2179**:247–258, 1994.
73. L. B. Stelmach and W. J. Tam, "Processing Image Sequences Based on Eye Movements," *Proc. SPIE*, **2179**:90–98, 1994.
74. W. Geisler and J. S. Perry, "Real Time Foveated Multiresolution System for Low Bandwidth Video Communication," *Proc. SPIE*, **3299**:294–305, 1997.
75. A. Yarbus, *Eye Movements and Vision*. Plenum Press, New York, 1967.
76. A. M. Liu, G. K. Tharp, and L. Stark, "Depth Cue Interaction in Telepresence and Simulated Telemanipulation," *Proc. SPIE*, **1666**:541–547, 1992.
77. L. Stark, H. Yang, and M. Azzariti, "Symbolic Binding," *Proc. SPIE*, **3959**:254–267, 2000.
78. D. Noton and L. Stark, "Scanpaths in Saccadic Eye Movements While Viewing and Recognizing Patterns," *Vision Research* **11**:929–942 (1971).
79. J. B. Pelz, R. L. Canosa, D. Kucharczyk, J. S. Babcock, A. Silver, and D. Konno, "Portable Eye Tracking: A Study of Natural Eye Movements," *Proc. SPIE*, **3959**:566–582, 2000.
80. M. Flickner, H. Sawhney, W. Niblack, J. Ashley, Q. Huang, B. Dom, M. Gorkani, J. Hafner, D. Lee, D. Petkovic, D. Steele, and P. Yanke, "Query by Image and Video Content; the QBIC System," *Computer* **28**:23–32 (Sept. 1995).
81. T. Frese, C. A. Bouman, and J. P. Allebach, "Methodology for Designing Image Similarity Metrics Based on Human Visual System Models," *Proc. SPIE*, **3016**:472–483, 1997.
82. B. E. Rogowitz, T. frees, J. R. Smith, C. A. Bouman, and E. B. Kalin, "Perceptual Image Similarity Experiments," *Proc. SPIE*, **3299**:576–590, 1998.
83. A. Mojsilovic, J. Kovacevic, J. Hu, R. J. Safranek, and S. K. Ganapathy, "Retrieval of Color Patterns Based on Perceptual Dimensions of Texture and Human Similarity Rules," *Proc. SPIE*, **3644**:441–452, 2000.
84. B. E. Rogowitz and L. Treinish, "Data Visualization: The End of the Rainbow," *IEEE Spectrum*, 52–59 (Dec. 1998).

85. M. Kubovy, "Gestalt Laws of Grouping Revisited and Quantified," *Proc. SPIE*, **3016**:402–408, 1997.
86. A. K. Hon, L. T. Maloney, and M. S. Landy, "Influence Function for Visual Interpolation," *Proc. SPIE*, **3016**:409–419, 1997.
87. B. E. Rogowitz, D. A. Rabenhorst, J. A. Garth, and E. B. Kalin, "Visual Cues for Data Mining," *Proc. SPIE*, **2657**:275–300, 1996.
88. R. M. Boynton, "Eleven Colors that are Almost Never Confused," *Proc. SPIE*, **1077**:322–331, 1989.
89. G. A. M. Derefeldt and T. Swartling, "How to Identify upto 30 Colors without Training: Color Concept Retrieval by Free Color Naming," *Proc. SPIE*, **2179**:418–428, 1994.
90. S. N. Yendrikhoviski, "Computing Color Categories," *Proc. SPIE*, **3959**:356–364, 2000.
91. C. E. Hedin and G. A. M. Derefeldt, "Palette: A Color Selection Aid for VDU Displays," *Proc. SPIE*, **1250**:165–176, 1990.
92. B. E. Rogowitz, D. A. Rabenhorst, J. A. Garth, and E. B. Kalin, "Visual Cues for Data Mining," *Proc. SPIE*, **2657**:275–300, 1996.
93. H. Kusaka, "Apparent Depth and Size of Stereoscopically Viewed Images," *Proc. SPIE*, **1666**:476–482, 1992.
94. H. de Ridder, F. J. J. Blommaert, and E. A. Fedorovskaya, "Naturalness and Image Quality: Chroma and Hue Variation in Color Images of Natural Images," *Proc. SPIE*, **2411**:51–61, 1995.
95. T. V. Pappathomas, "See How They Turn: False Depth and Motion in Hughes's Reverspectives," *Proc. SPIE*, **3959**:506–517, 2000.
96. J. J. Koenderink, A. J. van Doorn, A. M. L. Kappers, and J. T. Todd, "Directing the Mental Eye in Pictorial Perception," *Proc. SPIE*, **3959**:2–13, 2000.
97. C. W. Tyler, "Eye Placement Principles in Portraits and Figure Studies over the Past Two Millennia," *Proc. SPIE*, **3299**:431–438, 1998.
98. E. Burton, "Seeing and Scribbling: A Computer Representation of the Relationship Between Perception and Action in Young Children's Drawings," *Proc. SPIE*, **3016**:314–323, 1998.
99. M. Brand, "Computer Vision for a Robot Sculptor," *Proc. SPIE*, **3016**:508–516, 1997.
100. J. C. Dalton, "Image Similarity Modes and the Perception of Artistic Representations of Natural Images," *Proc. SPIE*, **3016**:517–525, 1997.
101. U. Feldman, N. Jacobson, and W. R. Bender, "Quantifying the Experience of Color," *Proc. SPIE*, **1913**:537–547, 1993.
102. P. K. Kaiser and R. M. Boynton, *Human Color Vision*. Optical Society of America, Washington, DC.
103. B. A. Wandell, *Foundations of Vision*. Sinauer, Sunderland, MA, 1995.
104. T. N. Cornsweet, *Visual Perception*. Academic Press, New York, 1970.
105. R. L. D. Valois and K. K. D. Valois, *Spatial Vision*. Oxford University Press, New York, 1990.
106. H. B. Barlow and J. D. Molon (eds.), *The Senses*. Cambridge University Press, Cambridge, 1982.
107. R. L. Gregory, J. Harris, D. Rose, and P. Heard (eds.), *The Artful Eye*. Oxford University Press, Oxford, 1995.
108. S. Pinker, *How the Mind Works*. Norton, New York, 1999.
109. C. Ware (ed.), *Information Visualization: Optimizing Design For Human Perception*. Morgan Kaufmann Publishers, San Francisco, CA, 1999.
110. A. B. Watson (ed.), *Digital Image and Human Vision*. MIT Press, Cambridge, MA, 1993.
111. M. S. Landy and J. A. Movshon (eds.), *Computational Models of Visual Processing*. MIT Press, Cambridge, MA, 1991.
112. T. V. Pappathomas, C. Chubb, A. Gorea, and E. Kowler (eds.), *Early Vision and Beyond*. MIT Press, Cambridge, MA, 1995.
113. B. E. Rogowitz, "The Human Visual System: A Guide for the Display Technologist," in *Proceedings of the SID*, vol. 24/3, pp. 235–252, 1983.

This page intentionally left blank.

DO NOT DUPLICATE

VISUAL FACTORS ASSOCIATED WITH HEAD-MOUNTED DISPLAYS

Brian H. Tsou

*Air Force Research Laboratory
Wright Patterson AFB, Ohio*

Martin Shenker

*Martin Shenker Optical Design, Inc.
White Plains, New York*

25.1 GLOSSARY

Aniseikonia. Unequal right and left retinal image sizes. (Also see Chaps. 12 and 13.)

C, F lines. Hydrogen lines at the wavelengths of 656 and 486 nm, respectively.

D (diopter). A unit of lens power; the reciprocal of focal distance in m. (Also see Chap. 12.)

Prism diopter. A unit of angle; 100 times the tangent of the angle. (Also see Chaps. 12 and 13.)

25.2 INTRODUCTION

Some virtual reality applications require a head-mounted display (HMD). Proper implementation of a binocular and lightweight HMD is quite a challenge; many aspects of human-machine interaction must be considered when designing such complex visual interfaces. We learned that it is essential to couple the knowledge of visual psychophysics with sound human/optical engineering techniques when making HMDs. This chapter highlights some visual considerations necessary for the successful integration of a wide field-of-view HMD.

25.3 COMMON DESIGN CONSIDERATIONS AMONG ALL HMDS

A wide field-of-view HMD coupled with a gimballed image-intensifying (see Chap. 31 of Vol. II for details) or infrared (see Chap. 33 of Vol. II for details) sensor enables helicopter pilots to fly reconnaissance, attack, or search-and-rescue missions regardless of weather conditions. Before any HMD

can be built to specifications derived from mission requirements, certain necessary display performance factors should be considered. These factors include functionality, comfort, usability, and safety. A partial list of subfactors that affect vision follows:

Functionality	Comfort	Usability	Safety
Field of view	Eye line-of-sight/focus	Eye motion box	Center-of-mass
Image quality	Alignment	Range of adjustment	Weight

Safety

An analytic approach to considering these factors has been proposed.¹ In general, factors can be weighted according to specific needs. At the same time, safety—especially in the context of a military operation—cannot be compromised.² The preliminary weight and center-of-mass safety requirements for helicopters³ and jets⁴ have been published.

It is extremely difficult to achieve the center-of-mass requirement. Hanging eye lenses in front of one’s face necessitates tipping the center-of-mass forward. A recent survey disclosed that pilots routinely resort to using counterweights to balance the center-of-mass of popular night-vision goggles⁵ in order to obtain stability and comfort, even at the expense of carrying more weight.

Light yet sturdy optics and mechanical materials have already made the HMD lighter than was previously possible. In the future, miniature flat-panel displays (see Chap. 22 in this volume, Chaps. 17 and 19 in Vol. II, and Chap. 8 in Vol. V for details) and microelectromechanical system (MEMS) technology⁶ might further lessen the total head-borne weight. Previous experience⁷ shows that both the field-of-view and exit pupil have a significant impact on HMD weight. Although we can calculate¹ the eye motion box [see Eq. (1)] that defines the display exit pupil requirement, we still do not have a quantitative field-of-view model for predicting its optimum value.

Usability

In order to avoid any vignetting, the display exit pupil size should match the eye motion box. The eye motion box is the sum of lateral translation of the eye and the eye pupil projection for viewing the full field of view (FOV). We assume the radius of eye rotation is 10 mm.

Eye motion box in mm

$$= 2 \times 10 \times \sin(\text{FOV}/2) + \text{eye pupil diameter} \times \cos(\text{FOV}/2) + \text{helmet slippage} \tag{1}$$

The range of pertinent optomechanical HMD adjustments for helicopter aviators⁸ is listed as minimum-maximum in millimeters in the following table:

	Interpupillary Distance	Eye to Top of Head	Eye to Back of Head
Female	53–69	107.7–141.7	152.5–190.0
Male	56–75	115.1–148.2	164.4–195.4

Comfort

Levy⁹ shows that, in total darkness, the eye line-of-sight drifts, on average, about 5° down from the horizon. Figure 1 shows the geometric relationship of the visual axis to the horizon. Levy argues that the physiological position of rest of the eye orbit is downward pointing. Menozzi et al.¹⁰ confirm that

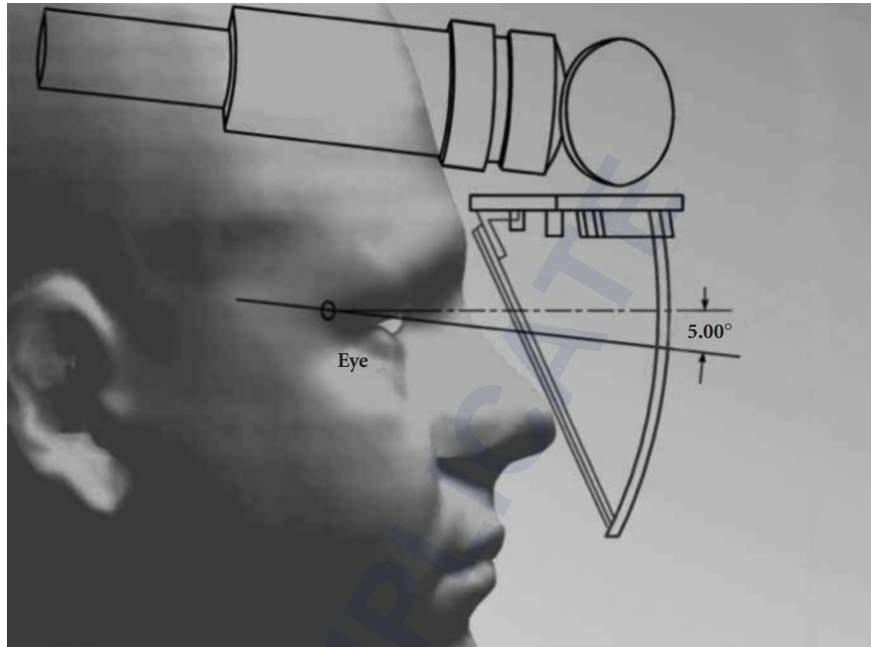


FIGURE 1 Geometric layout of the HMD showing the visual axis being 5° down from the horizon, as discussed in the text. Also shown are the schematic depictions of miniature CRT, folding mirror, beam-splitter, and spherical combiner. Its prescription is given in Table 1.

maximum comfort is in fact achieved when line of sight is pointed downward. The viewing comfort was determined by the perceived muscle exertion during gaze in a given direction for a given time. The median line of sight for the most comfortable gaze at 1 m is around 10° downward.

Our laboratory has also examined how normal ocular-motor behavior might improve the comfort of using any HMD. In a short-term wear study,¹¹ Gleason found that a range of eyepiece lens power of AN/AVS-6 Aviator Night Vision Imaging Systems (ANVIS) produced comparable visual acuity independent of luminance and contrast. The single best-overall eyepiece lens power produced visual acuity equal to, or better than, that of subject-adjusted eyepiece lens power, producing visual acuity within 2 percent of optimal. Infinity-focused eyepieces made visual acuity worse, reducing it by 10 percent. In his long-term wear (4 hours) study, -1.5 diopter (D) eyepiece lens power caused half of the subjects ($n = 12$) to complain of blurred or uncomfortable vision. These studies indicate that those users, who are optically corrected to a “most plus-best binocular visual acuity” endpoint (see Sec. 25.6, “Appendix”), achieve satisfactory comfort and near optimal binocular visual acuity for extended ANVIS viewing when an eyepiece lens power of approximately -0.75 D is added to the clinical refraction. This result, consistent with an earlier report by Home and Poole,¹² may extend to other binocular visual displays.

Alignment tolerances¹³ are the maximum permissible angular deviation between the optical axes of a binocular device that displays separate targets to the two eyes (see Chap. 13, “Binocular Vision Factors That Influence Optical Design,” in this volume for general issues). The ocular targets may be either identical or dichoptic (for stereoscopic presentation) but will be treated differently depending on whether there is a requirement for superimposition of the targets onto a direct view of the environment. If the direct view is either blocked out or too dark to see (e.g., flying at night), then it is classified as a nontransparent or closed HMD. A closed HMD resembles immersive virtual reality simulation where the real world is not needed.

Closed HMD Normal subjects have a substantial tolerance for angular deviation between the images in the two eyes for the three degrees of freedom of eye rotation: horizontal, vertical, and cyclorotational. These limits are especially of interest to optometrists and ophthalmologists. The optic axes of spectacle lenses must be aligned with the visual axes of the eyes of the wearer. Displacement of a spectacle lens from a centered position on the visual axis of the eye introduces a prismatic deviation that is entirely analogous to the misalignment of binocular devices. The American National Standards Institute (ANSI) has published ANSI 280.1–1987¹⁴ and permits a maximum deviation of $\frac{2}{3}$ prism diopter (23 arc minutes) horizontally and $\frac{1}{3}$ prism diopter (11.5 arc minutes) vertically. This standard can be adopted for closed binocular HMD. These tolerance values are summarized as follows:

Recommended Tolerances for Closed HMD	
	Tolerances (arc minutes)
Horizontal	±23
Vertical	±11.5
Cyclorotational	±12

Regarding the cyclorotational tolerances, both ophthalmic lenses and binocular displays can produce rotational deviation of the images about the line of fixation. ANSI spectacle standards do not exist for this alignment axis. Earlier research (dealing with distortion of stereoscopic spatial localization resulting from meridional aniseikonia at oblique axes) shows that there is considerable tolerance for cyclotorsional rotations of the images.¹⁵ However, the effects of cyclorotational misalignment are complicated by the fact that these rotations may result in either an oculomotor (cyclofusional) or sensory (stereoscopic) response. If the eyes do not make complete compensatory cyclorotations, a declination error (rotational disparity) will exist. Declination errors result in an apparent inclination (rotation about a horizontal axis) of the display; that is, the display will appear with an inappropriate pitch orientation.

Sensitivity to declination (stereoscopic response to rotational disparity) is quite acute. Ogle¹⁵ reports normal threshold values of ± 6 arc minutes, which corresponds to an object inclination of 5° at 3 m distance. If this threshold value were adopted for the cyclotorsion tolerance measurement, it would be overly conservative, as it would deny any reflex cyclofusional capacity to compensate for misalignment errors. A 97 percent threshold¹⁵ (± 12 arc minutes) is suggested for adoption as the cyclofusional tolerance.

Transparent HMD Transparent HMD optically superimposes a second image upon the directly viewed image using a beam combiner (see Sec. 25.6, “Appendix”). The difference in alignment angles between the direct and superimposed images equals the binocular disparity between target pairs.¹⁶ Horizontal binocular disparity is the stimulus to binocular stereoscopic vision. Consequently, sufficient lateral misalignment will lead to the perception of a separation in depth between the direct and superimposed images.

The exact nature of the sensory experience arising from lateral misalignment depends on the magnitude of the relative disparity. Although the threshold disparity for stereopsis is generally placed at about 10 arc seconds for vertical rods,¹⁷ it can approach 2 arc seconds under optimal conditions using greater than 25 arc-minutes-long vertical line targets.¹⁸ The normative value of threshold stereopsis in the standardized Howard-Dolman testing apparatus is 16 arc seconds. This defines 100 percent stereoacuity for clinical-legal requirements.

This extremely small stereothreshold (~ 2 arc seconds for optimal conditions) establishes an unpractically small tolerance for the design of displays. However, this precision is only required if the relative depths of the direct and superimposed fields need to correspond and the direct and superimposed fields are dissimilar (nonfusible). An example would be the need to superimpose a fiduciary marker over a particular location in the direct field, which is important to military operations. If this is not a requirement, the binocular disparity between the direct and superimposed fields merely needs to not exceed the range for pyknostereopsis in order to permit sensory averaging of the fields. When two

views are overlaid with small disparity separations (<20 arc minutes), they are depth averaged to form a single depth plane (called pyknostereopsis).^{19,20} This sets an alignment tolerance of ± 20 arc minutes (i.e., a maximum binocular disparity of 20 arc minutes). This limit is also within the disparity range required for maximum efficacy of oculomotor alignment of the eyes.²¹ Because of the poor sensitivity of the human visual system to absolute disparity,²² this amount of oculomotor misalignment will not influence the perception of absolute distances of objects in the display.

Vertical disparities do not result in local stereopsis, but they can have an effect on the overall stereoscopic orientation of objects in the binocular field. This response to vertical disparity is known as the *induced effect*. The sensitivity to vertical disparity for extended targets is similar to that of the horizontal disparity.²³ The other factor to consider when establishing vertical tolerances is the dimension of Panum's fusional area. Images falling within the Panum's area will fuse even if these images do not match exactly on the corresponding points between the two eyes. The vertical dimension of Panum's area can be as much as 2.5 times smaller than the horizontal dimension.²⁴

The recommended cyclorotational tolerance is based on Ogle's¹⁵ normal threshold value of ± 6 arc minutes, as discussed in the previous section. Because reflex cyclofusional would be restricted in transparent HMD by the direct field, no relaxation for compensatory eye movements has been included in this tolerance recommendation.

In formulating the final recommendations, the range of pyknostereopsis was reduced by half to introduce a margin for individual differences. Summary recommendations for tolerance in transparent HMD are provided as follows:

Recommended Tolerances for Transparent HMD

	Tolerance (arc minutes)
Horizontal	± 10
Vertical	± 4
Cyclorotational	± 6

Functionality

A larger field of view equates with larger optics and, hence, a heavier piece of equipment. As early as 1962, optical engineers partially overlapped the individual telescopes to achieve a wider field of view without greatly increasing weight.²⁵ The impact of partial overlap on image quality, however, was not fully addressed until recently.^{26,27} These psychophysical studies on binocular summation, rivalry, and stereopsis support a 40° binocular overlap in the middle of the field of view as both necessary and sufficient, with a divergent configuration of the two optical axes preferred over a convergent configuration (please note the visual axes are still parallel or convergent). Quantitative flight data (eye gaze) and pilots' subjective assessment are consistent with laboratory predictions.²⁸

Still, how large should the total field of view be? Everyday experience tells us that it should be as large as our visual field. On the other hand, most studies consistently show there is no significant human performance gain as the field of view extends beyond the 40° to 60° range.²⁹ Restricting head movements or manipulating field-of-view aspect ratio³⁰ does not seem to affect the performance except for the smaller field of view (from 15° to 30°). These studies suggest that (1) we are basically foveated beings,³¹⁻³³ and (2) unless compensated by large target size, our peripheral sensitivity is very poor.³⁴ Nevertheless, quantitative models for the useful field-of-view^{35,36} requirement are urgently needed so that HMD field of view is optimized to suit each application and scenario.

Image quality is perhaps the most important parameter for the acceptance of HMD since it is the aspect most noticeable by the users. Various image quality metrics have been devised (see Chap. 1 of this volume for reviews). The fundamental limitation on image quality in any display is either the loss of resolution, or lack thereof. Historically, the poor resolution of the miniature image source has hampered HMD progress; soon, the expected higher pixel density should significantly improve perceived image quality. For example, with high-definition TV's approximate 2000 pixels, a 50° field-of-view

HMD will have 20 cy/deg resolution or a respectable 20/30. The loss of resolution in visual optics is primarily caused by the defocus and lateral chromatic smearing of an image point.³⁷ When the eye is centered in the exit pupil of a display, defocus is due to residual field curvature and astigmatism. As the eye is decentered (or off aperture) in the exit pupil, residual amounts of uncorrected spherical aberration and coma in the display cause additional astigmatism. Astigmatism is defined as the variation in focus as a function of orientation in the exit pupil. Recently, Mouroulis and Zhang³⁸ have shown that MTFa (the area under the modulation transfer function over a specified frequency range of from 5 to 24 cy/deg) gives good correlation with subjective sharpness judgment for all target orientations. There is another convenient approximation to use for image quality. Figure 2 shows plots of the variation in modulation transfer function (MTF) (see Chap. 1 of Vol. II for reviews) for a 3-mm aberration-free pupil out to a frequency of 1 arc minute per line pair (or 20/10) as its image is defocused from the best focus position. A 3-mm pupil has been chosen because the eye is generally considered to be diffraction-limited in image quality near this value; we conclude that for “each tenth of diopter defocus 1 arc minute of resolution is lost” (see Fig. 2). We arrived at the same conclusion using Greivendamp’s schematic eye.³⁹ Greivendamp et al. showed their model correlates quite well with Snellen visual acuity. It should be noted that, for small defocus errors (<0.5 D), a real eye can see slightly better than our calculation predicts.⁴⁰ This may be attributable to the residual spherical aberration in the real eye. While this spherical aberration lowers the MTF for the exact in-focus image, it increases the depth-of-focus over which the maximum resolution is visible. Regardless, our rule of thumb is quite useful for the purpose of design evaluation.

Similarly, we approximate the residual chromatic smear as the sum of the lateral color plus the off-aperture manifestation of the residual longitudinal color by using a pair of wedges with zero deviation at a middle wavelength but a specific chromatic spread between the C and F spectral lines. The loss of MTF for a 3-mm, otherwise aberration-free exit pupil is shown in Fig. 3. It can be seen that these defects have a very significant effect on the resolution capabilities at both the high and low spatial frequencies. Therefore, a rule of thumb of “ $\frac{2}{3}$ arc minute of resolution is lost for each arc minute of chromatic spread between C and F lines” is applicable.

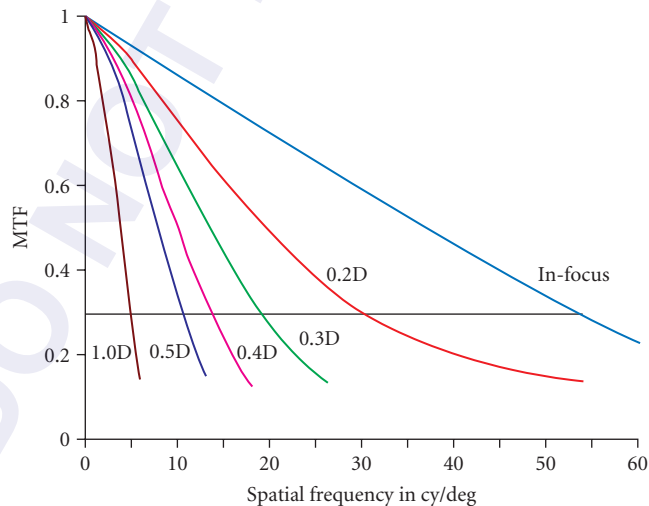


FIGURE 2 Loss of MTF in terms of resolution with defocus. The diopter values are the amount of defocus. At 30 percent modulation, the resolutions are about 1.1 arc minutes per line pair for in-focus, 2 for 0.2 D defocus, 3 for 0.3 D, 4.5 for 0.4 D, 6 for 0.5 D, and 12 for 1.0 D, respectively.

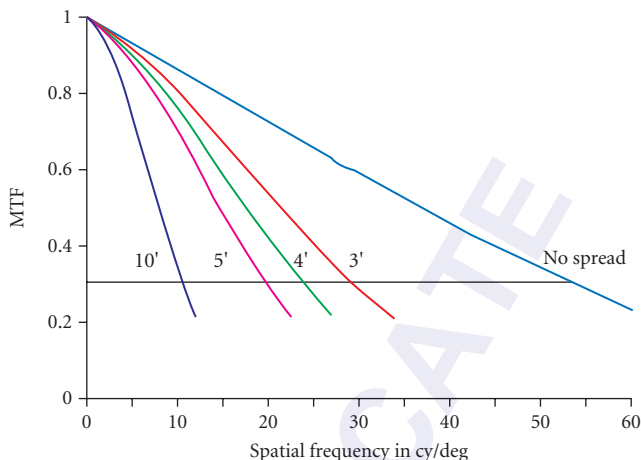


FIGURE 3 Loss of MTF with chromatic smear. The arc minute values represent smear between C and F lines. At 30 percent modulation, the resolutions are about 1.1 arc minutes per line pair for no chromatic spread, 2 for 3 arc minutes spread, 2.5 for 4', 3.1 for 5', and 6 for 10', respectively.

25.4 CHARACTERIZING HMD

Traditionally, we design and evaluate the image at the focal plane of photographic lenses with a plane, aberration-free wavefront, entering the entrance pupil of the objective. This is likely the arrangement in which the photographic lens will be utilized. The design of visual optics is usually performed in a similar manner with a perfect wavefront at the exit pupil of the display (which is defined as the entrance pupil of the eye in the design process). However, the use and the real performance of such a system is best defined, from a visual point of view, by ray tracing from a perfect object and evaluating the nature of the light exiting the system. There are at least two characteristics of this light that are of interest to the display designer.

The first of these is the direction of a ray that is defined as intersecting the center of the eye pupil that is viewing a particular point on the object. This datum is used for evaluation of the system magnification (i.e., mapping and distortion). If the system is supposed to be displaying an image projected to infinity (collimated display), then the apparent direction of a point in the field of view should not change as the eyes move relative to the display. Any deviation of this ray direction across the field of view will bring about an apparent change of the display distortion as well. This variation of the display distortion is often referred to as “dynamic distortion” or “swimming.” Next, if all the rays over the small (3 to 4 mm) entrance pupil of the eye are in focus, a relatively well-corrected display will demonstrate essentially diffraction-limited performance. However, if the eye cannot accommodate for the defocus, then this can cause a significant loss of modulation in the image.

Thus, our approach to evaluating the display image quality is to assess the instantaneous image quality available to the observer as he/she looks at different field angles from a particular point within the eye motion box. We define a small eye pupil at various positions in the larger exit pupil of the display; then we compute the apparent sagittal and tangential foci for various points in the field of view, as seen from this particular eye point. Both the exit pupils and the field of view are defined in terms of a clock face, as shown in Fig. 4. The displacement of the small eye pupil is in the 12:00 direction

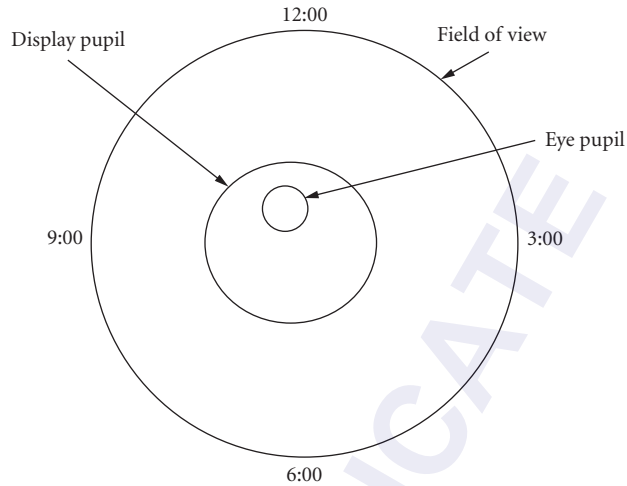


FIGURE 4 Spatial and directional relationships between display exit pupil and eye pupil with respect to field of view.

within the larger exit pupil of the display. This radial direction is the same as the direction of the off-axis point in the field of view. We make the following further definitions:

6:00 is the direction in the field-of-view opposite to the direction in which the eye pupil is displaced or off-aperture within the large exit pupil of the display.

3:00 and 9:00 are directions in the field perpendicular to the direction in which the eye pupil is off-aperture in the large exit pupil of the display. These values are identical.

We simply trace a matrix of object points on the object and calculate the tangential and sagittal foci in diopters for each of these points. We are thus assuming that the aberration variation across the 3- to 4-mm eye pupil is small and that the only residual defect is defocus. The three possibilities are as follows:

If the tangential and sagittal dioptric values are equal and the defocus is such that the observer can accommodate for this apparent distance, then the resolution potential of the display at this point in the field reverts to the diffraction limited potential of the aperture of his/her eye.

If the tangential and sagittal values are unequal but he/she can accommodate to either of the foci, then the image quality loss is defined by a focus halfway between the two foci.

If the tangential and/or sagittal values are such that the observer cannot accommodate for the defocus, then the loss of image quality is defined by this total defocus value.

In this manner, a particular optical design can be evaluated for image quality. An indication that the design is not sufficiently corrected is when focus conditions at various points in the field of view change rapidly for small displacements in the exit pupil. Our method for evaluating the potential display image quality is completely independent of their scale. It is valid for the 1.5-m focal length tilted mirrors used in the collimated displays for flight simulators, the 150- to 200-mm focal length refractive collimators used in heads-up displays, the 50- to 75-mm focal length eyepieces used in binocular magnifiers, and the 20-mm focal length erecting eyepiece systems used in HMD. However, for the HMD, the loss of image quality is likely to be the most severe. We will illustrate this method with a 50° field-of-view full-overlap binocular HMD, using only four elements to minimize weight. To conserve expense, all elements are simple spheres. The display has an overall focal length of 24.6 mm and is

TABLE 1 Prescription for a 50° Field-of-View Display Optics

Surface	Radius (mm)	Thickness (mm)	Glass
Object	∞	0.519	Air
1	∞	2.54	BK7
2	20.559	5.715	Air
3	∞	8.89	SK16
4	-26.75	33.02	Air
5	119.907	2.54	SF58
6	30.384	11.43	LAF2
7	-48.552	7.62	Air
8	105.5	5.08	SK16
9	-105.5	95.325	Air
10	∞	-38.1	Mirror
11	98.298	76.2	Mirror
Stop	∞	0	
Image	∞		

dimensioned to provide a 19-mm exit pupil with a horizontal field of view of 50° provided by a 23-mm image source format. Our electronic driver bandwidth can only adequately display 1000 video lines, resulting in a resolution of slightly better than 20/60 (10 cy/deg). This resolution is deemed minimal in order to support flying nap-of-the-earth (terrain following) or contour flight.⁴¹ The prescription is given in Table 1.

Table 2 lists the optical defocus values exhibited by this display example. A negative value indicates that the light exiting the display is convergent and corresponds to images focused beyond infinity. Objective measurement of the display system confirms the defocus table. The display focus is then set to 0.87 D (see subsection entitled “Comfort”). Table 3 lists the visual defocus values according to the

TABLE 2 Optical Defocus versus Pupil Displacements

Field Direction	Field Angle (deg)	Pupil Displaced							
		0 mm		2 mm		4 mm		6 mm	
		12-6	3-9	12-6	3-9	12-6	3-9	12-6	3-9
12:00	25	-0.4*	+0.2†	-0.5	+0.1	-1.0	-0.1	-1.7	-0.5
12:00	20	-0.6	0	-0.6	-0.1	-1.2	-0.3	-2.0	-0.6
12:00	15	-0.4	-0.1	-0.6	-0.1	-1.2	-0.3	-2.0	-0.7
12:00	10	-0.2	0	-0.4	-0.1	-1.1	-0.4	-2.0	-0.7
12:00	5	-0.1	0	-0.3	-0.1	-1.1	-0.4	-2.1	-0.8
	0	0	0	-0.3	-0.1	-1.3	-0.4	-2.3	-0.9
6:00	5	-0.1	0	-0.5	-0.1	-1.5	-0.5	-2.6	-1.0
6:00	10	-0.2	0	-0.7	-0.2	-1.8	-0.6	-2.9	-1.1
6:00	15	-0.4	-0.1	-0.9	-0.2	-2.0	-0.6	-3.1	-1.1
6:00	20	-0.6	0	-1.1	-0.2	-2.1	-0.5	-2.9	-0.9
6:00	25	-0.4	+0.2	-0.7	+0.1	-1.5	-0.2	-1.7	-0.5
3:00	25	+0.2	-0.4	-0.2	-0.5	-1.0	-0.7	-1.8	-1.1
3:00	20	0	-0.6	-0.3	-0.7	-1.2	-1.0	-2.1	-1.4
3:00	15	-0.1	-0.4	-0.4	-0.5	-1.3	-0.9	-2.2	-1.3
3:00	10	0	-0.2	-0.4	-0.3	-1.3	-0.7	-2.3	-1.1
3:00	5	0	-0.1	-0.3	-0.2	-1.3	-0.5	-2.3	-1.0
	0	0	0	-0.3	-0.1	-1.3	-0.4	-2.3	-0.9

*Negative diopters indicate convergent rays or farther than optical infinity.

†Positive diopters indicate divergent rays or closer than optical infinity.

TABLE 3 Visual Defocus versus Pupil Displacements

Field Direction	Field Angle (deg)	Pupil Displaced							
		0 mm		2 mm		4 mm		6 mm	
		12–6	3–9	12–6	3–9	12–6	3–9	12–6	3–9
12:00	25	-0.3*	+0.3 [†]	-0.3	+0.3	-0.9	0	-1.3	0
12:00	20	-0.3	+0.3	-0.3	+0.3	-0.2	0	-1.3	0
12:00	15	-0.3	+0.3	-0.2	+0.2	-0.2	0	-1.3	0
12:00	10	-0.1	+0.1	-0.2	+0.2	-0.1	0	-1.3	0
12:00	5	0	0	-0.1	+0.1	-0.1	0	-1.3	0
	0	0	0	-0.1	+0.1	-0.3	0	-2.3	-0.9
6:00	5	0	0	-0.2	+0.2	-1.0	0	-2.6	-1.0
6:00	10	-0.1	+0.1	-0.3	+0.3	-1.2	0	-2.9	-1.1
6:00	15	-0.3	+0.3	-0.7	0	-1.5	0	-3.1	-1.1
6:00	20	-0.3	+0.3	-0.9	0	-1.5	0	-2.9	-0.9
6:00	25	-0.3	+0.3	-0.4	+0.4	-1.3	0	-1.2	0
3:00	25	+0.3	-0.3	+0.2	-0.2	-0.2	0	-1.8	-1.1
3:00	20	+0.3	-0.3	+0.2	-0.2	-1.2	-1.0	-2.1	-1.4
3:00	15	+0.2	-0.2	+0.1	-0.1	-0.4	0	-2.2	-1.3
3:00	10	+0.1	-0.1	0	0	-0.6	0	-2.3	-1.1
3:00	5	0	0	-0.1	+0.1	-0.8	0	-2.3	-1.0
	0	0	0	-0.1	+0.1	-0.8	0	-2.3	-0.9

*Negative diopters indicate convergent rays or farther than where eye accommodates.

[†]Positive diopters indicate divergent rays or closer than where eye accommodates.

eye's ability to accommodate. Most observers cannot relax their accommodation beyond infinity (or > -0.87 D in our example) and thus suffer a loss in resolution corresponding to the full amount of defocus. The subjective evaluations agree with the predictions. Namely, while performance is satisfactory for the center of the field when the eye is positioned in the center of the display exit pupil, the performance falls off radically as the eye is decentered in the pupil. In addition, there is a significant falloff of performance across the field of view. We suspect the image quality would not be acceptable if the display resolution were much higher than the 20/60 used.⁴²

25.5 SUMMARY

We have described various visual issues dealing with a 50° binocular HMD breadboard. Both optical axes point downward by 5° and converge and intersect at a point about 1 m in front. Preliminary evaluation of the breadboard in the laboratory has not revealed any problems with respect to usability and comfort. However, the optics are undercorrected for both spherical aberration and astigmatism as shown by the defocus calculations. The astigmatism causes the variation of focus across the field of view when viewed from the center of the eye motion box. The spherical aberration causes the variation of focus when the eye pupil is decentered in the eye motion box. We found that using the defocus to predict display resolution is both valuable and intuitive in assessing HMD image quality during design and initial evaluation.

25.6 APPENDIX

There has been a concern⁴³ with respect to the transparent HMD and heads-up display that virtual images appear closer than the directly viewed image, inducing an undesired accommodation of the viewer. After investigating the issue, we did not find any support for misaccommodation. First, let us define the following terms:

Refractive error is defined by the eye's far point, which is the farthest distance to which the eye can focus. Eye focus is the dioptric distance at which the eye is focused, making objects at that distance optically conjugate with the eye's retina. Therefore, eye focus equals the dioptric sum of refractive error plus accommodation.

Accommodation is the eye's ability to focus closer than the far point by increasing its optical power through the action of the ciliary muscles. The hyperope's far point is beyond optical infinity, and therefore negative in value. Hyperopes must accommodate to focus to optical infinity. Myopia has a positive refractive error—that is, the eye is optically too strong to focus on distance objects. Do not confuse this with the clinical convention of associating myopia with the negative lens power of the correcting spectacles.

How is refractive error⁴⁴ determined? Optometers and clinical refraction optometers (e.g., infrared,^{45–48} laser,^{49,50} and stigmatoscopy⁵¹) measure eye focus, not accommodation. Optometers cannot distinguish between nonaccommodating 1-D myopes, emmetropes accommodating 1 D, or 1-D hyperopes accommodating 2 D since all are focused 1 m away. Optometers approximate refractive error (i.e., far point) only when accommodation is minimized under cycloplegia. Cycloplegia is the pharmacologically induced paralysis of accommodation.

Thus, clinical refraction does not determine refractive error in the strict sense. Clinical refraction determines the most-plus lens power resulting in best binocular visual acuity based on a subjective blur criterion (see Chap. 12, "Assessment of Refraction and Refractive Errors and Their Influence on Optical Design" and Chap. 13, "Binocular Vision Factors That Influence Optical Design," in this volume for general issues). Therefore, its theoretical endpoint places the distal end of the eye's depth-of-field (see Chap. 1, "Optics of the Eye," in this volume for reviews) at the eye chart, and does not place eye focus at optical infinity. Since eye charts are usually located about 6 m (+0.17 D) away and the eye's depth-of-field is +0.43 D for detecting blur of small sharp-edged high-contrast targets,⁵² the clinical emmetrope is about +0.6 D myopic. This concept is supported by Morgan's⁵³ report that subjective tests tend to overcorrect by +0.75 D, and Bannon et al.'s⁵¹ report that apparent emmetropes focused from +0.5 to +0.75 D closer than a fixated distant target. The presumed benefit of a myopic endpoint is that the eye's entire depth of field is maximized in real space when distant objects are viewed.

There are two potential causes for misaccommodation with transparent HMD: near-awareness of the virtual image and accommodative trap by the transparent combiner. While proximal accommodation would be the result of the former, the latter may be related to the Mandelbaum's phenomenon. However, Jones⁵⁴ found that unless the exit pupil is abnormally constrained (<3 mm), proximal accommodation does not affect the accommodative accuracy for targets in virtual image displays. Therefore, those deleterious proximal effects in virtual image displays may be avoided by utilization of an exit pupil in binocular HMD sufficiently large to encompass a natural pupil (e.g., display exit pupil matches eye motion box).

Contrary to accepted opinion, Gleason⁵⁵ found that Mandelbaum's phenomenon is not caused by accommodative capture. Almost fifty years ago, Mandelbaum⁵⁶ reported that distant scenery appears blurred when viewed through a window screen. The blur subsided with side-to-side head motion and returned after a delay when head motion ceased. These collective observations are called Mandelbaum's phenomenon. Mandelbaum hypothesized that the eye reflexively focused on the intervening screen instead of the desired distant scenery. Indeed, the assumed disruption of voluntary accommodation by intervening objects is now called the Mandelbaum effect.⁵⁷ Gleason investigated Mandelbaum's phenomenon using objective and continuous measurements of accommodation. In each trial, subjects looked directly at the distant E and were instructed to keep the E as clear as possible. After 7 s, a mesh screen was placed in the line of sight at the distance of either causing the greatest Mandelbaum phenomenon or at 40 cm in front. The 40-cm distance was near the subjects' dark focus. At 15 s, subjects were instructed to focus on the screen. At 23 s, subjects were instructed to refocus on the distant E. After the trial, the experimenter shook the screen side to side to verify that Mandelbaum's phenomenon was seen. Subjects were then questioned about what they saw during the trial. After each maximum Mandelbaum's phenomenon trial, subjects reported the distant letter blurred within 1 s of screen interjection; the letter remained blurred until the experimenter shook the screen at trial's end. These observations indicate that Mandelbaum's

phenomenon occurred during the trial. Mandelbaum's phenomenon was weak or nonexistent in 40-cm trials.

A narrow range of screen distances (0.7 to 0.9 D) produced maximum Mandelbaum's phenomenon. These results nearly match Mandelbaum's report of from 0.6 to 1.1 D. The screen distance of maximum Mandelbaum's phenomenon was at 0.8 D. At this distance, the screen strands subtended 0.8 arc minute and the nominal screen gap width was 5.5 arc minutes. Dark focus averaged 1.7 D, matching Leibowitz and Owens^{58,59} reports of 1.7 and 1.5 D for 124 and 220 college-age observers, respectively.

Gleason found that eye focus was essentially equal for all viewing tasks requiring fixation of the distant letter. If Mandelbaum's phenomenon were caused by the screen's capture of accommodation, then accommodation would shift eye focus toward the screen. But the eye focus is essentially unchanged before and after the introduction of the intervening screen, indicating that the screen did not influence accommodation. Eye focus shifted to the screen only after the instruction to focus on the screen. Moreover, eye focus returned to its initial level after the instruction to refocus on the distant letter. Clearly, the intervening screen did not disrupt volitional accommodation, either. Even more interesting is the fact that accommodative accuracy was not diminished by the perceptual blur as eye focus shifted from the screen back to the distant letter. In fact, it turns out that clinical emmetropia promotes Mandelbaum's phenomenon. Since clinical emmetropes are actually myopic (+0.6 D nearsighted), when they look at the mesh screen (placed at +0.8 D), the distant scenery (0.0 D) is defocused by 0.6 D, but the screen is only defocused by 0.2 D. Thus, the more blurred distant scenery may be suppressed for the benefit of the screen. This imbalance of retinal blur increases with uncorrected clinical myopia and, consequently, less time is needed to initiate Mandelbaum's phenomenon. On the other hand, an uncorrected hyperope can actually focus to the distant scenery, making the nearby screen blurrier. Now, the screen is suppressed instead of the distant scenery, precluding Mandelbaum's phenomenon. In summary, all five subjects in Gleason's study saw Mandelbaum's phenomenon without a concomitant shift in accommodation, which eliminates the misaccommodation concern over transparent HMD.

25.7 ACKNOWLEDGMENTS

The author (BT) thanks Tom Metzler, Mike Richey, Trang Bui, Luc Biberman, and Walter Hollis for the opportunities to be involved with Army HMD programs. We thank Jay Enoch, Neil Charman, and Vince Billock for their critical reviews of the manuscript. We also thank Sarah Tsou for her editorial support. Finally, we acknowledge the significant contributions by Jerry Gleason toward the advancements of HMD designs.

25.8 REFERENCES

1. B. H. Tsou, "System Design Considerations for a Visually Coupled System," in S. R. Robinson (ed.), *The Infrared and Electro-Optics System Handbook: Vol. 8. Emerging Systems and Technologies*. SPIE Optical Engineering Press, Bellingham, WA, 1993, pp. 515–536.
2. D. A. Fulghum, "Navy Orders Contractors to Stop Work on Ejectable Night Vision Helmets," *Aviation Week and Space Technology* **133**(23):67–68 (Dec. 3, 1990).
3. B. J. McEntire and D. F. Shanahan, "Mass Requirements for Helicopter Aircrew Helmets," *Advisory Group for Aerospace Research and Development (AGARD) Conference Proceedings* **597**:20.1–20.6 (Nov. 7–9, 1996).
4. C. E. Perry and J. R. Buhrman, "HMD Head and Neck Biomechanics," in J. E. Melzer and K. Moffitt (eds.), *Head Mounted Displays: Designing for the User*, McGraw-Hill, New York, 1997, pp. 147–172.
5. B. McLean, S. Shannon, J. McEntire, and S. Armstrong, "Counterweights Used with ANVIS," Report USAARL-96-30, US Army Aeromedical Research Laboratory, Fort Rucker, AL, 1990.
6. K. Y. Lau, "Microscanner Raster-Scanning Display: A Spyglass for the Future," *Optics and Photonics News* **10**(5):47–50, 84 (May 1999).

7. M. Shenker, "Optical Design Criteria for Binocular Helmet-Mounted Displays," *SPIE* **778**:70–78 (1987).
8. S. M. Donelson and C. C. Gordon, "1988 Anthropometric Survey of US Army Personnel: Pilot Summary Statistics," Report Natick-TR-91-040, US Army Natick Research, Development and Engineering Center, Natick, MA, 1991.
9. J. Levy, "Physiological Position of Rest and Phoria," *Am. J. Ophthalmol.* **68**(4):706–713 (1968).
10. M. Menozzi, A. Buol, H. Krueger, and Ch. Mieke, "Direction of Gaze and Comfort: Discovering the Relation for the Ergonomic Optimization of Visual Tasks," *Ophthalmic Physiol. Opt.* **14**:393–399 (1994).
11. G. A. Gleason and J. T. Riegler, "The Effects of Eyepiece Focus on Visual Acuity Through ANVIS Night Vision Goggles During Short- and Long-Term Wear," unpublished report performed for A. F. Workunit 71842604, 1997. (Part of the report was first presented at 1996 Human Factors Society Annual Meeting.)
12. R. Home and J. Poole, "Measurement of the Preferred Binocular Dioptric Settings at a High and Low Light Level," *Opt. Acta* **24**(1):97–98 (1977).
13. R. Jones, "Binocular Fusion in Helmet Mounted Displays," unpublished report performed for A. F. Workunit 71842601, 1996. (Part of the report was first presented at 1992 Society for Information Display Annual Meeting.)
14. American National Standards Institute, 1430 Broadway, New York, NY 10018.
15. K. N. Ogle, *Researches in Binocular Vision*, Hafner Publishing Co., New York, 1964, pp. 283–290.
16. G. A. Fry, "Measurement of the Threshold of Stereopsis," *Optometric Weekly* **33**:1029–1033 (1942).
17. G. Sigmar, "Observations on Vernier Stereo Acuity with Special Consideration to Their Relationship," *Acta Ophthalmol.* **48**:979–998 (1970).
18. E. E. Anderson and F. Weymouth, "Visual Perception and the Retinal Mosaic," *Am. J. Physiol.* **64**:561–591 (1923).
19. R. A. Schumer, "Mechanisms in Human Stereopsis," Ph.D. dissertation, Stanford University, Palo Alto, CA, 1979.
20. C. W. Tyler, "Sensory Processing of Binocular Disparity," in Schor and Ciuffreda (eds.), *Vergence Eye Movements: Basic and Clinical Aspects*, Butterworths, Boston, 1983, pp. 199–295.
21. A. M. Norcia, E. E. Sutter, and C. W. Tyler, "Electrophysiological Evidence for the Existence of Coarse and Fine Disparity Mechanisms in Humans," *Vis. Res.* **25**:1603 (1985).
22. C. J. Erkelens and H. Collewijn, "Motion Perception During Dichoptic Viewing of Moving Random Dot Stereograms," *Vis. Res.* **25**:345–353 (1985).
23. K. N. Ogle, *Researches in Binocular Vision*, Hafner Publishing Co., New York, 1964, pp. 173–199.
24. C. M. Schor and C. W. Tyler, "Spatial Temporal Properties of Panum's Fusional Area," *Vis. Res.* **21**:683–692 (1981).
25. M. Shenker, J. LaRussa, P. Yoder, and W. Scidmore, "Overlapping-Monoculars—An Ultrawide-Field Viewing System," *Appl. Opt.* **1**:399–402 (1962).
26. S. S. Grigsby and B. H. Tsou, "Grating and Flicker Sensitivity in the Near and Far Periphery: Naso-Temporal Asymmetries and Binocular Summation," *Vis. Res.* **34**:2841–2848 (1994).
27. S. S. Grigsby and B. H. Tsou, "Visual Processing and Partial-Overlap Head-Mounted Displays," *J. Soc. for Information Display* **2/2**:69–73 (1994).
28. T. H. Bui, R. H. Vollmerhausen, and B. H. Tsou, "Overlap Binocular Field-of-View Flight Experiment," *1994 Society for Information Display International Symposium Digest of Technical Papers* **25**:306–308 (1994).
29. M. J. Wells, M. Venturino, and R. K. Osgood, "The Effect of Field-of-View Size on Performance at a Simple Simulated Air-to-Air Mission," *SPIE* **1116**:126–137 (1989).
30. B. H. Tsou and B. M. Rogers-Adams, "The Effect of Aspect Ratio on Helmet-Mounted Display Field-of-View," *Report of the 30th Meeting of Air Standardization Coordination Committee (ASCC) Working Party 61: Aerospace Medical and Life Support Systems Symposium: Aeromedical Aspects of Vision* **4**:136–146. Defense and Civil Institute of Environmental Medicine, Toronto, Canada, 1990.
31. A. T. Bayhill and D. Adler, "Most Naturally Occurring Human Saccades Have Magnitudes of 15 Degrees or Less," *Invest. Ophthalmol.* **14**:468–469 (1975).
32. G. H. Robinson, "Dynamics of the Eye and Head During Movement Between Displays: A Qualitative and Quantitative Guide for Designers," *Human Factors* **21**:343–352 (1979).
33. J. Wolfe, P. O'Neil, and S. C. Bennett, "Why Are There Eccentricity Effects in Visual Search? Visual and Attentional Hypotheses," *Percept. Psychophys.* **60**:140–156 (1998).

34. S. P. McKee and K. Nakayama, "The Detection of Motion in the Peripheral Visual Field," *Vis. Res.* **24**(1):25–32 (1984).
35. K. K. Ball, B. L. Beard, D. L. Roenker, R. L. Miller, and D. S. Griggs, "Age and Visual Search: Expanding the Useful Field of View," *J. Opt. Soc. Am.* **A5**:2210–2219 (1988).
36. P. Havig and B. H. Tsou, "Is There a Useful Field of Visual Search?" [Abstract]. *Optical Society of America Annual Meeting*, Baltimore, MD, Oct. 4–9, 1998.
37. L. N. Thibos and A. Bradley, "Modeling the Refractive and Neuro-Sensor Systems of the Eye," in P. Mouroulis (ed.), *Visual Instrumentation Handbook*, McGraw-Hill, New York, 1999, pp. 4.19–4.20.
38. P. Mouroulis and H. Zhang, "Visual Instrument Image Quality Metrics and the Effects of Coma and Astigmatism," *J. Opt. Soc. Am.* **A9**(1):34–12 (1992).
39. J. E. Greivenkamp, J. Schwiegerling, J. M. Miller, and M. D. Mellinger, "Visual Acuity Modeling Using Optical Raytracing of Schematic Eyes," *Am. J. Ophthalmol.* **120**(2):227–240 (1995).
40. W. N. Charman and J. A. M. Jennings, "The Optical Quality of the Monochromatic Retinal Image As a Function of Focus," *Br. J. Physiol. Opt.* **31**:119–134 (1976).
41. D. Greene, "Night Vision Pilotage System Field-of-View (FOV)/Resolution Tradeoff Study Flight Experiment," Report NV-1-26, Center for Night Vision and Electro-Optics, Ft. Belvoir, VA, 1988.
42. G. E. Legge, K. T. Mullen, G. C. Woo, and F. W. Campbell, "Tolerance to Visual Defocus," *J. Opt. Soc. Am.* **A4**(5):851–863 (1987).
43. J. H. Iavecchia, H. P. Iavecchia, and S. N. Roscoe, "Eye Accommodation to Head-Up Virtual Images," *Human Factors*, **130**(6):687–702 (1988).
44. G. A. Gleason and J. T. Riegler, "Do Clinical Refraction Result in Best Visual Acuity?" unpublished report performed for A. F. Workunit 71842604, 1997. (Part of the report was first presented at 1996 Association for Research in Vision and Ophthalmology Annual Meeting.)
45. T. N. Cornsweet and H. D. Crane, "Servo-Controlled Infrared Optometer," *J. Opt. Soc. Am.* **60**:548–554 (1970).
46. H. D. Crane and C. M. Steele, "Accurate Three Dimensional Eyetracker," *Appl. Opt.* **17**:691–705 (1978).
47. F. W. Campbell and J. G. Robson, "High-Speed Infra-Red Optometer," *J. Opt. Soc. Am.* **49**:268–272 (1959).
48. N. A. McBrien and M. Millodot, "Clinical Evaluation of the Canon Autorefractor R-1," *Am. J. Optom. Physiol. Opt.* **62**:786–792 (1985).
49. R. T. Hennessy and H. W. Leibowitz, "Laser Optometer Incorporating the Badal Principle," *Behav. Res. Meth. and Instrum.* **4**(5):237–239 (1972).
50. A. Morrell and W. N. Charman, "A Bichromatic Laser Optometer," *Am. J. Optom. Physiol. Opt.* **64**:790–795 (1987).
51. R. E. Bannon, F. H. Cooley, H. M. Fisher, and R. T. Textor, "The Stigmatoscopy Method of Determining the Binocular Refractive Status," *Am. J. Optom. Arch. Am. Acad. Optom.* **27**:371 (1950). (Cited in I. M. Borish, *Clinical Refraction*, 3d ed., Professional Press, Chicago, 1975, pp. 783–784.)
52. F. W. Campbell, "The Depth of Field of the Human Eye," *Opt. Acta* **4**:157–164 (1957).
53. M. W. Morgan, "The Meaning of Clinical Tests in the Light of Laboratory Measurements," *The Ind. Opt.* **2**:6 (1947). (Cited in I. M. Borish, *Clinical Refraction*, 3d ed., Professional Press, Chicago, 1975, p. 356.)
54. R. Jones, "Proximal Accommodation in Virtual Displays," *Society for Information Display International Symposium Digest of Technical Papers* **24**:195–198 (1998).
55. G. A. Gleason and R. V. Kenyon, "Mandelbaum's Phenomenon Is Not Caused by Accommodative Capture," unpublished report performed for A. F. Workunit 71842604, 1998. (Part of the report was first presented at 1997 Association for Research in Vision and Ophthalmology Annual Meeting.)
56. J. Mandelbaum, "An Accommodative Phenomenon," *AMA Arch. Ophthalmol.* **63**:923–926 (1960).
57. D. A. Owens, "The Mandelbaum Effect: Evidence for an Accommodative Bias Toward Intermediate Viewing Distance," *J. Opt. Soc. Am.* **65**:646–652 (1975).
58. H. W. Leibowitz and D. A. Owens, "Anomalous Myopias and the Dark Focus of Accommodation," *Science* **189**:646–648 (1975).
59. H. W. Leibowitz and D. A. Owens, "New Evidence for the Intermediate Position of Relaxed Accommodation," *Doc. Ophthalmol.* **46**:133–147 (1978).

INDEX

Index note: The *f* after a page number refers to a figure, the *n* to a note, and the *t* to a table.

- Aberrations (in general):
and contact lenses, **20.24–20.25, 20.25f**
and forward scattering of light, **1.21**
in intraocular lenses, **21.9–21.10, 21.10f**
and Stiles-Crawford effect, **8.8**
- Aberrations (human eye), **1.3, 16.6**
absence of, **1.12–1.14, 1.13f, 1.14f**
and AO-controlled light delivery
to generate aberrations, **15.24**
longitudinal chromatic aberration, **15.22**
transverse chromatic aberration,
15.22–15.23
- astigmatism, **15.2**
- chromatic, **1.19–1.20, 1.19f**
and accommodation, **1.34**
age-related, **14.14**
correcting color coordinates for, **10.17**
correction of, **1.25–1.26, 1.26f**
longitudinal, **15.22**
and macular pigment, **1.9**
transverse, **15.22–15.23**
with visual instruments, **1.28**
- control of, **1.3**
correction of, **1.25–1.26, 1.26f, 10.17**
- defocus, **15.2**
and depth of focus, **1.28**
higher-order, **16.6, 16.7**
and idiosyncratic peculiarities, **1.6**
- image quality
for aberration-free eye, **1.12–1.14**
calculating, **1.21–1.22**
- monochromatic, **1.4, 1.14–1.19**
age-related, **14.12–14.14, 14.13f**
correction of, **1.25, 1.26, 1.26f**
off-axis, **1.18–1.19, 1.18f**
on the visual axis, **1.15–1.18**
and observed optical performance, **1.23–1.26**
properties of, **15.4–15.7, 15.5f, 15.6f**
- Aberrations (human eye) (*Cont.*):
and pupil diameter, **1.8**
and refractive surgery, **1.15**
- Aberrometers, **1.15, 1.23, 12.6**
- Ablation rate (refractive correction),
16.16–16.18, 16.18f
- Abney's law, obedience to, **10.44, 11.37**
- Absorption, lenticular, **1.9**
- Accommodating intraocular lenses,
14.29–14.30, 21.1, 21.14, 21.18–21.19
- Accommodation, **1.3, 1.29–1.36, 12.3, 21.2**
accuracy of, **1.32–1.34, 1.33f, 1.34f**
in aging eyes, **1.35–1.36, 1.35f, 14.4, 21.3**
application to instrumentation, **1.34–1.35**
and change in lens, **1.5**
with computer work, **23.10–23.11**
defined, **12.1, 13.1, 16.1, 21.1, 23.1**
dynamics of, **1.31–1.32, 1.31f**
and extraocular muscle movement, **12.16**
fluctuations in, **1.32**
with head-mounted displays, **13.31, 25.10–25.12**
and ocular aberration, **1.17–1.18**
and presbyopia, **14.8, 14.9f, 14.29–14.30**
and refractive errors, **12.15–12.16**
and spherical ametropias, **16.5**
stability of, **1.32**
vergence input, **1.29, 1.30, 1.34**
- Accommodation demand, **13.30**
- Accommodative demand:
with contact lenses, **20.26–20.30,**
20.27f–20.29f, 20.29t
defined, **20.1**
- Accommodative miosis, **1.30**
- Acellular, **16.1**
- Achromatic detection:
and chromatic adaptation, **11.47–11.49,**
11.48f
measurements favoring, **11.37**

- Achromatic mechanism (color vision),
11.1, 11.37
- Achromatic signals, multiplexing of chromatic signals and (color vision), 11.76–11.79, 11.78f
- Action spectra (ocular radiation), 7.2, 7.3, 7.4f
- Adaptation:
in aging eyes, 14.15
in contrast detection, 2.25–2.27
multiplicative and subtractive, 2.26
phoria, 13.21–13.22
in wearing contact lenses, 12.12
- Adaptive optics (AO) (in retinal microscopy and vision), 1.25, 15.1–15.24
- AO-controlled light delivery to the retina, 15.22–15.24
alignment, 15.23
in an AO SLO, 15.23
conventional AO vision systems, 15.23
to generate aberrations, 15.24
longitudinal chromatic aberration, 15.22
measuring activity of individual cones, 15.24
transverse chromatic aberration, 15.22–15.23
uses of, 15.23–15.24
- for correcting monochromatic aberrations, 1.25
and correction for SCE-1, 9.5
defined, 15.1, 17.1
history of, 15.2–15.3
imaging of the retina, 15.16–15.22
contrast and resolution, 15.21–15.22
flood-illuminated AO ophthalmoscope, 15.16–15.17, 15.16f, 15.17f
optical coherence tomography, 15.19–15.21, 15.20f, 15.21f
scanning laser ophthalmoscope, 15.17–15.19, 15.18f, 15.19f, 17.9
implementation of, 15.7–15.15, 15.7f
control system, 15.12–15.15
wavefront corrector, 15.9–15.12, 15.10f, 15.11f
wavefront sensor, 15.8–15.9, 15.8f
properties of ocular aberrations, 15.4–15.7, 15.5f, 15.6f
wavefront corrector, 15.9–15.12, 15.10f, 15.11f
wavefront sensor, 15.8–15.9, 15.8f
- Adaptive optics retina cameras, 15.3, 15.12
(*See also* Ophthalmoscopes)
- Additivity:
in color matching, 10.8
of color opponency and color appearance, 11.67
with equilibrium colors, 11.66
field, 11.51, 11.52f
obedience to Abney's law, 10.44, 11.37
Stiles-Crawford effect, 9.10f–9.11f
- Adjustment tasks:
judgment tasks vs., 3.2
psychophysical measurement of, 3.4–3.5
magnitude production, 3.5
matching, 3.4–3.5
nulling, 3.4
threshold, 3.4, 3.5f
- Aesthetics and emotional characteristics (electronic imaging), 24.9–24.10
color, art, and emotion, 24.10
image quality, 24.9
virtual reality and presence, 24.9
- After-effect (color vision):
and McCollough effect, 11.76
and orientation of contours, 11.75–11.76, 11.75f
and orientation-selectivity, 11.80
- Afterimage, 23.1
- Against-the-rule astigmatism, 1.6
- Age-related changes in vision, 1.5, 1.7, 1.35–1.36, 14.1–14.30
in accommodation, 1.35–1.36, 1.35f, 14.4, 21.3
in bovine lenses, 19.10
changes in the eye causing, 21.3–21.4
in color vision, 14.15, 14.17
and demographic changes in world, 14.2–14.4, 14.2f
in depth and stereovision, 14.22
economic/social implications of, 14.4
and fluorescence, 1.21
and health implications of living longer, 14.3–14.4
index of diffusion for, 1.23
and level of aberration, 1.18
minimal, 14.22
ocular diseases, 14.22–14.27
age-related macular degeneration, 14.24–14.25
cataract, 14.24
diabetic retinopathy, 14.25–14.26
glaucomas, 14.26–14.27
life-span environmental radiation damage, 14.22–14.23

- Age-related changes in vision (*Cont.*):
- in optics, 14.4–14.14
 - accommodation and presbyopia, 14.8, 14.9*f*
 - anterior and posterior chambers, 14.5, 14.6
 - cornea, 14.5, 14.6*f*
 - eye size, 14.11
 - lens, 14.7–14.8
 - pupil, 14.6–14.7, 14.7*f*
 - retina, 14.9–14.11
 - tears, 14.4–14.5
 - transparency/cataract, 14.8
 - presbyopic corrections for, 14.27–14.30
 - accommodation restoration, 14.29–14.30
 - contact lenses, 14.27–14.28
 - intraocular lenses, 14.28–14.29
 - noncataract related refractive surgeries, 14.29
 - spectacles, 14.27
 - in pupil diameter, 1.8, 1.18
 - in retinal image quality, 14.11–14.14
 - chromatic aberration, 14.14
 - intraocular scatter, 14.12
 - monochromatic aberrations, 14.12–14.14, 14.13*f*
 - in retinal reflectance, 1.11
 - and RMS wavefront error, 1.15–1.17, 1.16*t*
 - in scattering, 1.20
 - in sensitivity, 14.15, 14.16*f*
 - in spatial vision, 14.17–14.19, 14.18*f*
 - in temporal vision, 14.19–14.21, 14.20*f*
 - in transmittance, 1.9
 - in visual acuity, 4.13
 - in visual field, 14.21, 14.21*f*
 - Age-related macular degeneration (AMD), 14.24–14.25
 - defined, 14.1
 - OFDI at 1050 nm in, 18.17, 18.17*f*
 - OFDI in diagnosis of, 18.8
 - Airy diffraction pattern, 1.12
 - Airy's disk, 4.1, 4.4*f*, 4.5
 - Alignment:
 - of eye (*see* Eye alignment)
 - of photoreceptors, 8.5–8.8, 8.6*f*, 8.7*f*
 - A-line, 18.1
 - Allelotropia, 13.7
 - Amacrine cells, 2.10, 2.11
 - Amblyopia, 2.34, 2.35, 12.1, 12.16
 - American Conference of Governmental Industrial Hygienists (ACGIH), 7.9
 - American National Standard Institute (ANSI), 7.9, 7.11
 - Ametropias, 1.6–1.7, 12.4, 16.4–16.5, 16.5*f*
 - astigmatism, 16.5–16.6
 - correcting, 12.16, 16.8*f*
 - defined, 12.1, 13.1
 - spherical, 16.5
 - uncorrected, 12.16
 - Amplification, in color CRTs, 22.6–22.9, 22.7*f*
 - Angle α , 1.20
 - Angle of incidence, obliquity of, 8.20
 - Angular resolution, in diffraction-limited eye, 1.12
 - Aniridia, 9.1
 - Aniseikonia, 1.41–1.42
 - binocular factors in, 12.16–12.17
 - defined, 9.1, 12.1, 13.1, 25.1
 - distortion from interocular anisomagnification, 13.17–13.18
 - Anisometropia, 1.7, 1.41, 1.42
 - binocular factors in, 12.16–12.17
 - defined, 12.1, 23.1
 - interocular blur suppression with, 13.18–13.19
 - as problem with computer work, 23.11
 - Anisophoria, optically induced, 13.26
 - Anterior chamber, 14.5, 14.6, 16.3
 - Aphakia, 12.14, 14.28
 - correction of, 12.14–12.15
 - defined, 9.1, 12.1, 14.1, 21.1
 - Apodization, 6.14, 21.17, 21.18
 - defined, 8.1
 - Stiles-Crawford effect, 8.8
 - Aqueous humor, 1.3*f*, 14.5, 14.6, 14.26, 16.3
 - Aspheric intraocular lenses, 21.1, 21.10–21.12, 21.11*f*, 21.12*t*
 - Aspheric lens design, 20.1, 20.20–20.23, 20.21*f*–20.23*f*
 - Asthenopia (eyestrain), 23.2
 - Astigmatic dial, 12.7*f*
 - Astigmatism, 1.6, 1.7, 1.7*f*, 1.18*f*, 12.4, 1.6.5–16.6
 - and binocular instrumentation, 1.41–1.42
 - and contact lenses, 20.15

- Astigmatism (*Cont.*):
 correction of, 16.9
 cylindrical correction, 13.16, 13.17
 with hydrogel contact lenses, 12.13
 with spectacle lenses, 12.9, 15.2
 with spherical contact lenses, 12.12
 defined, 12.1, 23.1
 determining axis of, 12.6
 following cataract surgery, 21.13
 irregular, 16.6
 measuring, 12.7
 off-axis, 1.18, 1.18*f*
- Asymmetric color matching, 11.27, 11.29, 11.30*f*
 and changes in overall intensity, 11.68
 and chromatic adaptation, 11.69
 and color constancy, 11.71
- Atrophic (dry) age-related macular degeneration, 14.1
- Attention, in human vision, 24.6, 24.7
- Autocorrelation noise, 18.11–18.12
- Autokeratometers, 12.6
- Autorefractors, 12.6
- Axial ametropia, 20.1
- Axial aniseikonia, 13.17, 13.18
- Axial edge lift, 20.1
- Axial gradients, 19.5
- Axial resolution, 18.1
- Back central optical radius (BCOR), 20.3
- Back vertex power (BVP), 12.1, 20.1
 contact lenses, 20.7–20.8, 20.8*t*
 spectacle lenses, 12.9
- Base curve:
 of contact lenses, 12.1, 12.12
 of spectacle lenses, 12.1, 12.9
- Base curve radius (BCR):
 for contact lenses, 20.3, 20.4*f*, 20.5
 defined, 20.1
- Baseline, 13.1
- Beam separators, 5.10
- Beamsplitters, 5.9–5.10, 5.9*f*
- Beer-Lambert law, 2.7
- Bessel functions, 8.1, 8.13, 8.22
- Best-correction, 1.6
- Bezold-Brücke hue shift, 11.1, 11.67–11.68, 11.67*f*
- Bichromatic test mixtures, sensitivity to, 11.39, 11.40
- Bifocal jump, with monocular magnification, 13.15–13.16, 13.15*f*
- Bifocal lenses, 12.10*f*, 14.27
 contact lenses, 12.13, 14.28
 intraocular lenses, 14.28
 for presbyopia, 12.8, 12.10
 vertical prism introduced by, 13.28–13.29
- Bimorph mirror, 15.1, 15.10, 15.10*f*, 15.11
- Binocular cues, 13.3, 13.7
- Binocular disparities, 2.40, 13.1, 13.4–13.5
 induced by prisms or lenses, 13.28
 and perceived direction, 13.7–13.8
- Binocular field, 1.3
 horizontal angular extent of, 1.38*f*
 stereopsis in, 1.38–1.42
- Binocular fusion, 13.1, 13.12–13.13
- Binocular instrumentation:
 and chromostereopsis, 1.20
 differential focusing for, 1.7
 tolerances in, 1.41–1.42
- Binocular parallax, 13.1, 13.22
- Binocular rivalry, 13.1, 13.19
- Binocular rivalry suppression, 13.12, 13.33
- Binocular stereoscopic discrimination, 2.40–2.41, 2.41*f*
- Binocular vision factors, 13.1–13.35
 in computer work, 23.11
 coordination of eyes, 13.20–13.25
 binocular parallax, 13.22
 cross-coupling and direction/distance of gaze, 13.23–13.24, 13.23*f*
 defocus and efforts to clear vision, 13.22–13.23
 intrinsic stimuli to vergence, 13.21–13.22
 perceived distance, 13.24
 zone of clear and single binocular vision, 13.24–13.25, 13.25*f*
 distortion by monocular magnification, 13.13–13.16
 bifocal jump, 13.15–13.16, 13.15*f*
 from convergence responses to prism, 13.19
 discrepant views of objects/images, 13.16
 motion parallax, 13.14–13.15
 perspective distortion, 13.13, 13.14*f*
 stereopsis, 13.13, 13.14
 distortion from interocular anisomagnification, 13.16–13.19
 aniseikonia, 13.17–13.18
 interocular blur suppression with anisometropia, 13.18–13.19
 lenses and prisms, 13.16–13.17, 13.17*f*

- Binocular vision factors (*Cont.*):
- eye alignment, 13.20–13.25
 - magnification induced errors of, 13.25–13.27
 - prism induced errors of, 13.27–13.29
 - eye movements, 13.19–13.20
 - focus and responses to distance, 13.30
 - fusion and suppression, 13.12–13.13
 - gaze control, 13.29–13.30
 - and head mounted visual display systems, 13.31–13.35
 - distance conflicts, 13.31–13.32
 - optical errors, 13.34–13.35
 - spatial location conflicts, 13.33, 13.34f
 - visual-vestibular conflicts, 13.32–13.33
 - lens effects on vergence and phoria, 13.25–13.27
 - perceived direction, 13.7–13.10
 - corresponding retinal points, 13.8
 - horopter, 13.8–13.9, 13.9f, 13.10f
 - vertical horopter, 13.9
 - perceived space, 13.3–13.7
 - binocular cues, 13.7
 - extraretinal information for eye movements, 13.7
 - kinetic cues, 13.4–13.7, 13.5f, 13.6f
 - monocular cues, 13.3–13.7
 - prisms
 - distortion from interocular anisomagnification, 13.16–13.17
 - effects on vergence and phoria, 13.25–13.27
 - errors of alignment with, 13.25–13.29
 - in refractive errors, 12.15–12.17
 - aniseikonia, 12.16–12.17
 - anisometropia, 12.16–12.17
 - convergence and accommodation, 12.15–12.16
 - stereopsis, 13.11–13.12, 13.11f
 - visual field, 13.3
- Binocularly, 23.1
- Biological waveguides, 8.1–8.29
- in cochlear hair cells and human hair, 8.24–8.26
 - in fiber-optic plant tissues, 8.26–8.28, 8.27f
 - models of, 8.8–8.15, 8.13f, 8.14f
 - physical assumptions in, 8.22
 - and propagation of light, 8.21, 8.22
 - retinal layer of rods and cones, 8.8–8.9, 8.9f, 8.10f, 8.12–8.15
 - three-segment model, 8.9, 8.11–8.15
- Biological waveguides (*Cont.*):
- and photoreceptors, 8.3–8.5
 - modal patterns in monkey/human receptors, 8.19–8.24, 8.22f
 - orientation and alignment, 8.5–8.8, 8.6f, 8.7f
 - photoreceptor optics, 8.3
 - quantitative observations of single receptors, 8.15, 8.16f, 8.17, 8.18f, 8.19f
 - in sponges, 8.28–8.29
 - and Stiles-Crawford effect of the first kind, 8.3
- Bipolar cells (retina), 2.9, 2.10
- Bipolar mechanisms (color vision), 11.1, 11.80–11.81
- Birefringence, 1.10, 18.1, 18.20, 18.22, 18.25, 18.27
- Bitoric lenses, 20.1, 20.17–20.20, 20.18f–20.20f
- Blindness:
- from cataract, 14.24, 21.3
 - from diabetic retinopathy, 14.25
 - economic impact of, 14.4
 - from glaucoma, 14.26
 - and increasing life span, 14.3
 - from uncorrected ametropia, 12.4
- Bloch's law, 2.28
- Blue to red (B/R) ratio, 8.28
- Blue-light photochemical injury, 7.4, 7.10
- Blur, 1.28
 - in correction for SCE-1, 9.4–9.6, 9.14
 - with monovision, 14.28
 - suppression of, 13.18–13.19, 14.28
 - (*See also* Astigmatism; Defocus)
- Born approximation, 8.1, 8.15
- Bovine lenses, 19.8–19.11, 19.11f
- Bowman's layer, 16.4
- Bowman's membrane (cornea), 14.5
- Boynton illusion, 11.72, 11.73f, 11.78
- Brightness:
- in color CRTs, 22.7, 22.19
 - defined, 11.1, 11.8, 23.1
 - and luminous efficiency, 11.70, 11.70f
 - in monochrome CRTs, 22.8f
 - in visual acuity tests, 4.8
- Brightness matching, 10.43–10.45
- Bruch's membrane, 8.1, 14.25, 18.1
- Bumper function, 11.59, 11.60f
- Bunsen-Roscoe law of photochemistry, 7.1–7.3, 7.7

- C line, 25.1
- Calibration:
- of color CRTs, 22.20–22.21 (*See also* Characterization, of CRTs)
 - errors in, and SCE-1 correction, 9.6–9.13
 - in Maxwellian viewing systems, 5.17–5.18, 5.17*f*
- Cameras:
- adaptive optics retina cameras, 15.3, 15.12 (*See also* Ophthalmoscopes)
 - color, 10.38–10.40
- Capillary, 18.1
- Capsulorhexis, 21.1
- Cardinal directions (color vision), 11.1, 11.56
- Cardinal mechanisms (color vision), 11.1, 11.54, 11.56
- evidence for, 11.79
 - and sensitivity losses, 11.29
- Cat lenses, 19.9
- Cataract, 1.21, 14.8
- in aging eyes, 14.24, 21.3–21.4
 - defined, 14.1, 19.1, 21.1, 23.1
 - infrared, 7.7, 7.10
 - postsurgical correction of, 12.14
 - from radiation, 7.4–7.6, 7.6*f*, 14.23
 - treatment of, 12.14
- Cataract lenses, 12.15
- Cataract surgery, 14.4, 14.24
- intraocular lenses in, 21.4, 21.5 (*See also* Intraocular lenses)
 - renormalization following, 14.17
- Cathode ray tubes (CRTs), 23.1
- color (*see* Color cathode ray tubes)
 - and Computer Vision Syndrome, 23.7, 23.8
 - monochrome, 22.3–22.4, 22.3*f*
 - controls for, 22.7, 22.8*f*
 - design and operation of, 22.3–22.4, 22.3*f*
 - standards for, 22.14
 - in optical systems, 5.16
 - screen reflections with, 23.5
- Cauchy's equation, 1.19, 1.20
- Cavanagh, P., 11.80
- Center of rotation (eye), 1.42, 13.1
- Center thickness, for contact lenses, 20.6
- Central visual processing, 2.12–2.14, 2.13*f*
- Characterization:
- of CRTs, 22.20–22.34
 - absolute vs. for interaction, 22.33
 - calibration vs., 22.20
 - choice of method for, 22.20–22.21
- Characterization, of CRTs (*Cont.*):
- exhaustive, 22.21–22.23
 - local, 22.24–22.27
 - model-dependent, 22.27–22.33
 - of head-mounted displays, 25.7–25.10, 25.8*f*, 25.9*t*–25.10*t*
- Chief ray (retinal imaging), 4.3
- Choroid, 18.1
- Chromatic aberration, 1.19–1.20, 1.19*f*
- and accommodation, 1.34
 - age-related, 14.14
 - correcting color coordinates for, 10.17
 - correction of, 1.25–1.26, 1.26*f*
 - longitudinal, 15.22
 - and macular pigment, 1.9
 - transverse, 15.22–15.23
 - with visual instruments, 1.28
- Chromatic adaptation, 11.35*f*
- achromatic detection and, 11.47–11.49, 11.48*f*
 - color appearance and, 11.68
 - and luminous efficiency, 11.37
 - second-site, 11.17–11.18
 - and the Sloan notch, 11.49, 11.51
- Chromatic contrast detection, 2.29–2.31, 2.30*f*
- Chromatic CSFs, 2.29–2.31
- Chromatic detection:
- and color appearance, 11.69
 - on neutral fields, 11.34
- Chromatic discrimination, 11.57–11.62
- and chromatic adaptation, 11.69
 - and color appearance, 11.69
 - defined, 11.1
 - and gap effect, 11.72, 11.74
 - near detection threshold, 11.58–11.59
 - pedestal experiments, 11.59, 11.60*f*, 11.61–11.62
- Chromatic mechanisms, 11.1, 11.80–11.81
- Chromatic signals, multiplexing of achromatic signals and, 11.76–11.79, 11.78*f*
- Chromatic spatial CSFs, 2.29
- Chromatic temporal CSFs, 2.31
- Chromatic valence data, 11.27, 11.28*f*
- Chromaticity coordinates, 10.1, 10.4*t*, 10.20–10.21, 10.22*f*
- Chromaticity diagrams, 10.20–10.21, 10.22*f*
- Chromophores, 21.1, 21.20, 21.20*f*
- Chromostereopsis, 1.20
- CIE 1931 2° color-matching functions, 10.6, 10.12

- CIE 1964 10° color-matching functions, 10.12–10.13, 10.16
- CIE luminous efficiency functions, 10.44–10.45, 11.37, 11.38*f*
- CIELAB, 10.42–10.43
- Cilia, 8.1
- Ciliary body, 1.3*f*
- Ciliary muscle, 21.2
- Ciliary photoreceptors, 8.3
- Ciliary ring, 1.30–1.31
- Circular polarization (OCT), 18.21
- Clear lens extraction (CLE), 21.18
- Closed HMDs, 25.3, 25.4
- Cochlea, 8.1, 8.26
- Cochlear hair cells, light guide effect in, 8.24–8.26
- Coherence length, 18.1
- Coherent coupling, 1.27–1.28
- Coherent illumination, pupil size and, 6.9–6.12, 6.11*t*
- Coleoptile, 8.1, 8.26
- Collagen, 16.1, 16.4
- Color:
- in electronic imaging, 24.5–24.6, 24.8–24.9
 - and visual acuity, 4.10
- Color appearance:
- and chromatic adaptation, 11.68
 - and chromatic detection and discrimination, 11.69
 - and color constancy, 11.71
 - and color opponency, 11.62–11.66
 - color-opponent response or valence functions, 11.63, 11.65*f*
 - hue scaling, 11.63, 11.64*f*
 - opponent-colors theory, 11.62–11.63
 - spectral properties of color-opponent mechanisms, 11.63
 - unique hues and equilibrium colors, 11.63–11.66
 - defined, 11.2
 - and habituation, 11.69–11.70
 - phenomenological aspects of, 11.5
 - and postreceptor mechanisms, 11.3, 11.4*f*, 11.5
 - and stabilized borders, 11.74–11.75
 - trichromacy *vs.*, 11.5
- Color assimilation, 11.2, 11.4*f*
- Color balance and tracking, in color CRTs, 22.19
- Color cameras, 10.38–10.40
- Color cathode ray tubes (color CRTs), 22.1–22.34
- colorimetric calibration/characterization, 22.20–22.34
- absolute *vs.* characterization for interaction, 22.33
 - choice of method for, 22.20–22.21
 - exhaustive methods, 22.21–22.23
 - local methods, 22.24–22.27
 - model-dependent methods, 22.27–22.33
- design and operation of, 22.3–22.13
- electronics and controls, 22.6–22.13, 22.7*f*–22.13*f*
 - and monochrome CRTs, 22.3–22.4, 22.3*f*
 - shadowmask color CRTs, 22.4–22.6, 22.5*f*, 22.6*f*
- operational characteristics of, 22.13–22.18
- colorimetric standards, 22.14
 - spatial characteristics of emitted light, 22.14–22.15
 - spatial uniformity, 22.17–22.18, 22.17*f*, 22.18*f*
 - stability of output, 22.16–22.17, 22.16*f*, 22.17*f*
 - temporal characteristics of emitted light, 22.14–22.16
 - timing and synchronization standards, 22.13–22.14
- setup for image display, 22.18–22.19
- brightness, 22.19
 - color balance and tracking, 22.19
 - contrast, 22.19
 - focus, 22.19
 - viewing environments, 22.19–22.20
- Color constancy, 10.39–10.40, 11.2
- Color constancy mechanisms, 11.71–11.72
- Color contrast, 11.2, 11.4*f*
- Color convergence, 23.1
- Color coordinate systems, 10.11–10.24
- adjusting cone spectral sensitivities, 10.17–10.18
 - colorimetric measurements, 10.23–10.24
 - color-matching functions, 10.11–10.13
 - cone fundamentals, 10.13–10.14
 - limits of color-matching data, 10.15–10.17
 - opponent and contrast spaces, 10.18–10.19
 - stimulus spaces, 10.11
 - visualizing color data, 10.19–10.23, 10.19*f*, 10.22*f*, 10.23*f*

- Color coordinates:
 - in color-deficient observers, 10.16
 - correcting for chromatic aberrations, 10.17
 - of different visual systems, 10.38–10.39
 - of surfaces, 10.36
 - transformation to CIELAB, 10.42–10.43
- Color data representations, 11.31–11.33, 11.32*f*
- Color direction, 11.12
- Color discrimination, 10.40–10.43, 10.41*f*, 14.15, 14.17
- Color LCDs, 22.37–22.40
 - colorimetry of color pixels, 22.38–22.39
 - controls and input standards, 22.39
 - geometry of color pixels, 22.37–22.38, 22.37*f*
 - spatial variations in output, 22.40
 - temporal variations in output, 22.39–22.40
- Color matching, 10.6–10.10
 - in color-deficient observers, 10.16
 - consistency across observers, 10.9
 - critical properties of, 10.8–10.10
 - errors in, 10.43, 10.43*f*, 10.44
 - Grassmann's laws, 10.8–10.9
 - maximum saturation method, 10.6–10.7, 10.6*f*
 - Maxwell's method, 10.8, 10.8*f*
 - persistence of, 10.9
 - trichromatic, 10.7–10.8
 - tristimulus values for arbitrary lights, 10.9
 - uniqueness in, 10.9
- Color opponency:
 - and color appearance, 11.62–11.66
 - and “forbidden” colors, 11.74
 - implied by test measurements, 11.12–11.15, 11.13*f*, 11.26–11.27, 11.28*f*
 - and multiplexing of color and luminance signals, 11.76–11.79
 - third-level, zero crossings of, 11.83, 11.84
- Color space transformation matrix, 10.1, 10.31
- Color spaces, 10.11, 10.18–10.19, 11.31–11.33
- contrast, 10.19
- defined, 11.2
- different directions of, 11.39, 11.41–11.46, 11.44*f*, 11.46*f*
 - detection contours in L, M plane, 11.40*f*, 11.41–11.42
 - detection in planes other than L,M, 11.43–11.45, 11.44*f*
 - mechanism interactions, 11.42–11.43
 - spatial and temporal CSFs, 11.45–11.46, 11.46*f*
- Color spaces (*Cont.*):
 - DKL, 10.19, 11.32, 11.32*f*, 11.33
 - opponent, 10.19
 - specifying, 10.11
 - stimulus, 10.11
 - transformations between, 10.24*f*, 10.29–10.32
 - uniform, 10.40, 10.42
- Color temperature, 23.2
- Color valence, 11.2, 11.63, 11.65*f*, 11.66
- Color vision, age-related changes in, 14.15, 14.17
- Color vision mechanisms, 11.1–11.85
 - basic model details and limits, 11.31
 - chromatic discrimination, 11.57–11.62
 - near detection threshold, 11.58–11.59
 - pedestal experiments, 11.59, 11.60*f*, 11.61–11.62
 - color and contours, 11.72–11.79, 11.73*f*–11.75*f*
 - color appearance and stabilized borders, 11.74–11.75
 - contours and after-effects, 11.75–11.76, 11.75*f*
 - gap effect and luminance pedestals, 11.72, 11.74
 - McCollough effect, 11.76, 11.77*f*
 - multiplexing chromatic and achromatic signals, 11.76–11.79, 11.78*f*
- color appearance and color opponency, 11.62–11.66
 - color-opponent response (valence functions), 11.63, 11.65*f*
 - hue scaling, 11.63, 11.64*f*
 - opponent-colors theory, 11.62–11.63
 - spectral properties of color-opponent mechanisms, 11.63
 - unique hues and equilibrium colors, 11.63–11.66
- color constancy, 11.71–11.72
- color data representations, 11.31–11.33, 11.32*f*
- color-appearance mechanisms, 11.26–11.31
 - color-discrimination mechanisms vs., 11.5–11.8, 11.6*f*, 11.7*f*, 11.81–11.82
 - field measurements and first-site adaptation, 11.27, 11.29, 11.30*f*
 - field measurements and second-site adaptation, 11.29, 11.31
 - test measurements and opponency, 11.26–11.27, 11.28*f*

- Color vision mechanisms (*Cont.*):
- color-discrimination mechanisms, 11.9–11.26
 - color-appearance mechanisms vs., 11.5–11.8, 11.6f, 11.7f, 11.81–11.82
 - field method, 11.11–11.12
 - first-site adaptation, 11.15–11.17, 11.16f
 - opponency implied by test measurements, 11.12–11.15, 11.13f
 - psychophysical test method, 11.9, 11.11, 11.12
 - second-site adaptation, 11.17–11.22, 11.18f, 11.19f
 - sites of limiting noise, 11.20, 11.23–11.26, 11.23f, 11.25f
 - field sensitivities, 11.46–11.57
 - achromatic detection and chromatic adaptation, 11.47–11.49, 11.48f
 - chromatic adaptation and the Sloan notch, 11.49, 11.51
 - detection contours and field adaptation, 11.53–11.54, 11.54f
 - field additivity, 11.51, 11.52f
 - first- and second-site adaptation, 11.51, 11.52, 11.53f
 - habituation or contrast adaptation experiments, 11.54–11.56, 11.55f
 - multiple cone inputs, 11.49, 11.50f
 - noise-masking experiments, 11.56–11.57
 - Stiles' π -mechanisms, 11.46, 11.47, 11.47f
 - guiding principles of, 11.8–11.9
 - linearity of color-opponent mechanisms, 11.66–11.70
 - Bezold-Brücke effect and invariant hues, 11.67–11.68, 11.67f
 - color appearance and chromatic adaptation, 11.68
 - color appearance and chromatic detection/discrimination, 11.69
 - color appearance and habituation, 11.69–11.70
 - luminance and brightness, 11.70, 11.70f
 - tests of linearity, 11.66–11.67
 - low-level and higher-order mechanisms, 11.79–11.80
 - and mechanism concept, 11.9–11.11, 11.10f
 - nomenclature for, 11.8
 - test sensitivities, 11.34–11.46
- Color vision mechanisms, test sensitivities (*Cont.*):
- to different directions of color space, 11.39, 11.41–11.43, 11.44f, 11.45–11.46, 11.46f
 - luminance, 11.37–11.39, 11.38f, 11.40f
 - to spectral lights, 11.34, 11.35f, 11.36f, 11.37
 - three-stage zone models, 11.82–11.85, 11.85f
 - unipolar vs. bipolar chromatic mechanisms, 11.80–11.81
- Color-appearance mechanisms, 11.26–11.31
- color-discrimination mechanisms vs., 11.5–11.8, 11.6f, 11.7f, 11.81–11.82
 - defined, 11.2
 - field measurements
 - and first-site adaptation, 11.27, 11.29, 11.30f
 - and second-site adaptation, 11.29, 11.31
 - test measurements and opponency, 11.26–11.27, 11.28f
- Color-discrimination mechanisms, 11.9–11.26, 11.10f
- color-appearance mechanisms vs., 11.5–11.8, 11.6f, 11.7f, 11.81–11.82
 - defined, 11.2
 - field method, 11.11–11.12
 - first-site adaptation, 11.15–11.17, 11.16f
 - low-level and higher-order, 11.79–11.80
 - opponency implied by test measurements, 11.12–11.15, 11.13f
 - psychophysical test method, 11.9, 11.11, 11.12
 - second-site adaptation, 11.17–11.22, 11.18f, 11.19f
 - sites of limiting noise, 11.20, 11.23–11.26, 11.23f, 11.25f
- Colorimeters, 10.23
- Colorimetry, 10.1–10.45
- brightness matching and photometry, 10.43–10.45
 - color cameras, 10.38–10.40
 - color coordinate systems, 10.11–10.24
 - adjusting cone spectral sensitivities, 10.17–10.18
 - colorimetric measurements, 10.23–10.24
 - color-matching functions, 10.11–10.13
 - cone fundamentals, 10.13–10.14
 - limits of color-matching data, 10.15–10.17
 - opponent and contrast spaces, 10.18–10.19
 - stimulus spaces, 10.11
 - visualizing color data, 10.19–10.23, 10.19f, 10.22f, 10.23f

- Colorimetry (*Cont.*):
- color discrimination, 10.40–10.43, 10.41f
 - color matching, 10.6–10.10
 - conventional terms/notation, 10.4, 10.4t, 10.5
 - errors in color-matching functions, 10.43, 10.43f, 10.44
 - image processing chain, 10.2f
 - matrix representations/calculations, 10.24–10.32
 - stimulus representation, 10.24–10.27, 10.25f, 10.26f
 - transformations between color spaces, 10.24t, 10.29–10.32
 - vector representation of data, 10.25f, 10.27–10.29, 10.27f
 - metamerism, 10.36–10.38, 10.37f
 - scope of, 10.3
 - standards for color CRTs, 22.14
 - surfaces and illuminants, 10.32–10.36, 10.33f, 10.34f
 - trichromacy, 10.4–10.6
 - univariance, 10.4
 - visual systems, 10.38–10.40
- Color-matching functions (CMFs), 10.11–10.13
- defined, 10.1, 10.4t
 - limits of data, 10.15–10.17
 - and luminosity function, 10.10
 - and maximum saturation method, 10.7
 - online tabulation of, 10.11
 - specificity of, 10.9
 - standards for, 10.12–10.13
 - tailored to individuals, 10.15
 - transformation of, 10.10, 10.10f
- Color-opponent mechanisms:
- linearity of, 11.66–11.70
 - Bezold-Brücke effect and invariant hues, 11.67–11.68, 11.67f
 - color appearance and chromatic adaptation, 11.68
 - color appearance and chromatic detection/discrimination, 11.69
 - color appearance and habituation, 11.69–11.70
 - luminance and brightness, 11.70, 11.70f
 - tests of linearity, 11.66–11.67
 - and opponent-colors theory, 11.3
 - spectral properties of, 11.63
- Commission Internationale de l'Éclairage (CIE), 7.9, 10.1, 10.3, 10.27
- Complex cells (cortical neurons), 2.14
- Computer Vision Syndrome (CVS), 23.1–23.12
- disorders and eye conditions, 23.9–23.12
 - accommodation, 23.10–23.11
 - anisometropia, 23.11
 - binocular vision, 23.11
 - dry eyes, 23.9
 - presbyopia, 23.11–23.12
 - refractive error, 23.10
 - and work environment, 23.4–23.9
 - lighting, 23.4–23.5, 23.5t
 - monitor characteristics, 23.6–23.8
 - screen reflections, 23.5–23.6
 - work habits, 23.8–23.9
 - workstation arrangement, 23.8
- Computer-Automatic Virtual Environment (CAVE), 13.32
- Concomitant eye motion, 13.1
- Cone contrast spaces, 11.2, 11.32, 11.41–11.42
- Cone contrasts, 11.2, 11.32
- Cone coordinates, 10.1, 10.11 (*See also* Tristimulus values)
- Cone fundamentals, 10.13–10.14
- adjusting cone spectral sensitivities, 10.17–10.18
 - defined, 10.1, 10.4t
 - online tabulation of, 10.11
 - primaries yielding, 10.10
 - Smith-Pokorny, 10.12
 - Stockman and Sharpe, 10.12, 10.16
 - tailored to individuals, 10.15
- Cone magno pathway, 2.9, 2.10f
- Cone mechanisms, 11.10f, 11.11
- defined, 11.2
 - and stimulus direction, 11.33
 - and test measurements, 11.12–11.14
- Cone parvo pathway, 2.9, 2.10f
- Cone pathways, 2.9, 2.10, 2.10f
- Cone polymorphism, 10.15
- Cone-excitation spaces, 11.31–11.32
- Cone-opponent mechanisms, 11.8
- and bichromatic test mixtures, 11.39
 - defined, 11.2
 - sensitization in, 11.69
 - unipolar vs. bipolar, 11.80–11.81

- Cones (cone photoreceptors):
 adaptation at, 11.52
 and age-related photopic vision changes, 14.15
 alignment of, 8.4
 densities and distributions of, 2.7
 in dichromatic observers, 10.14
 directional sensitivity of, 8.5
 ideal “average,” 8.21, 8.21*f*
 inner segments, 2.6*f*
 light-collection area of, 2.8
 linear density of, 2.6*f*
 and maximum saturation color matching, 10.7
 measuring activity of, 15.24
 optical waveguide properties of, 14.11
 outer segments of, 2.6*f*
 pedicles of, 8.20
 photocurrent responses of, 2.8*f*
 in retinal layer of rods and cones model of biological waveguides, 8.8–8.9, 8.9*f*, 8.10*f*, 8.12–8.15
 S-cone flicker sensitivity, 11.74–11.75
 spatial distribution of, 2.6
 spectral sensitivities of, 10.3, 10.13, 11.8
 adjusting for individual differences, 10.17–10.18
 estimates of, 10.14, 10.14*f*
 loss hypothesis for, 10.14
 in normal and dichromatic observers, 10.13
 time constant of photopigment regeneration, 2.7
 types and functions of, 2.4, 2.6
- Confocal microscopes, 17.1–17.10
 clinical, 17.3, 17.6–17.9
 clinical applications using, 17.8–17.9
 defined, 17.1
 development of, 17.3, 17.5
 laser scanning, 17.3, 17.7–17.8
 Nipkow disk, 17.4*f*, 17.5, 17.5*f*
 scanning slit, 17.3, 17.6–17.9, 17.6*f*, 17.7*f*
 spatial filtering with, 17.3
 Svishchev, 17.6, 17.6*f*
 theory of confocal microscopy, 17.3, 17.4*f*
- ConfoScan 4 microscope, 17.6, 17.7, 17.9
- Conjugate (version) eye movements/position, 1.42, 1.43, 13.1, 13.7, 13.20
- Contact lenses, 20.1–20.34
 accommodation with, 20.26–20.30, 20.27*f*–20.29*f*, 20.29*t*
 anisophoria and, 13.26
 for aphakic patients, 12.15
 base curve of, 12.1
 convergence with, 20.25–20.26
 correction with, 12.11–12.14
 hydrogel lenses, 12.12–12.13
 for presbyopia, 12.13–12.14, 14.27–14.28
 rigid lenses, 12.11–12.12
 design considerations, 20.20–20.25
 aberrations, 20.24–20.25, 20.25*f*
 aspheric lenses, 20.20–20.23, 20.21*f*–20.23*f*
 posterior peripheral curve, 20.23–20.24, 20.24*t*
 design parameters, 20.2–20.6, 20.4*t*–20.5*t*
 base curve radius (BCR), 20.3, 20.4*f*, 20.5
 center thickness, 20.6
 edge thickness, 20.6
 optical zone diameter (OZD), 20.5
 overall diameter (OAD), 20.5
 posterior peripheral curve systems, 20.5, 20.6
 magnification, 20.31–20.33
 relative spectacle, 20.32–20.33, 20.33*f*
 spectacle, 20.31–20.32
 materials for, 20.3
 power of, 20.6–20.20
 back and front vertex power, 20.7–20.8, 20.8*t*
 and contact lens as thick lens, 20.6–20.7
 effective power, 20.8–20.12, 20.9*t*, 20.10*f*, 20.11*f*
 lacrima lens consideration, 20.12–20.15, 20.12*f*–20.14*f*
 residual astigmatism, 20.15
 of soft lenses, 20.15–20.16
 of toric lenses, 20.16–20.20, 20.16*f*–20.20*f*, 20.18*t*
 prismatic effects with, 20.30–20.31
 prism-ballasted contacted lenses, 20.30
 unintentional (induced) prism, 20.31
 (See also *specific types of lenses*)
- Contours (color vision), 11.72–11.79, 11.73*f*–11.75*f*
 and after-effects, 11.75–11.76, 11.75*f*
 color appearance and stabilized borders, 11.74–11.75

- Contours (color vision) (*Cont.*):
 detection surface/contour, 11.12
 defined, 11.2
 and directions of color spaces, 11.40*f*
 in equiluminant plane, 11.44*f*
 and field adaptation, 11.53–11.54, 11.54*f*
 in the L,M plane, 11.40*f*, 11.41–11.42
 discrimination contours, 11.44*f*
 gap effect and luminance pedestals, 11.72, 11.74
 McCollough effect, 11.76, 11.77*f*
 multiplexing chromatic and achromatic signals, 11.76–11.79, 11.78*f*
 threshold surface/contour, 11.12–11.15, 11.13*f*
 defined, 11.3
 and loss of information, 11.20
 and noise, 11.20, 11.23–11.26, 11.23*f*, 11.25*f*
 and second-site adaptation to steady fields, 11.18*f*
- Contrast:
 in color CRTs, 22.19
 defined, 23.2
 in monochrome CRTs, 22.7, 22.8*f*
 in retinal imaging, 15.21–15.22
 in vision experiments, 3.4
 and visual acuity, 4.12, 4.12*f*
- Contrast coding:
 first-site adaptation, 11.15–11.16
 Weber's law and, 11.15–11.16, 11.16*f*
- Contrast (modulation) color spaces, 10.19
- Contrast constancy, 2.33
- Contrast detection, 2.19–2.31
 adaptation and inhibition, 2.25–2.27, 2.25*f*
 chromatic, 2.29–2.31, 2.30*f*
 eye movements, 2.21
 as function of spatial frequency, 2.20*f*
 optical transfer function, 2.21–2.22, 2.22*f*
 optical/retinal inhomogeneity, 2.24, 2.25*f*
 receptors, 2.22–2.23
 spatial channels, 2.23–2.24
 temporal, 2.27–2.29, 2.27*f*, 2.29*f*
- Contrast discrimination, 2.31–2.32, 2.31*f*
- Contrast discrimination functions, 2.31–2.32, 2.31*f*
- Contrast estimation, 2.33
- Contrast masking, 2.32–2.33, 2.32*f*
- Contrast masking functions, 2.32, 2.32*f*
- Contrast sensitivity, 15.1
- Contrast sensitivity functions (CSFs), 2.20–2.23, 2.20*f*, 2.22*f*, 24.3
 and adaptation, 2.25–2.27, 2.25*f*
 age-related changes in, 14.17, 14.19
 chromatic, 2.29–2.31
 in color vision, 11.45–11.46, 11.46*f*
 at different retinal eccentricities, 2.24, 2.25*f*
 of ideal observer, 2.24
 for motion detection, 2.36, 2.37
 spatial, 11.45–11.46, 11.46*f*
 spatio-temporal, 2.28–2.29
 for stereopsis, 2.40, 2.41
 temporal, 2.27–2.29, 2.29*f*, 11.45
- Contrast-transfer function, resolving capacity of the eye and, 4.5
- Convergence, 13.21
 in color CRTs, 22.12, 22.13
 with contact lenses, 20.25–20.26
 defined, 13.1, 23.2
 and egocentric direction, 13.7
 and refractive errors, 12.15–12.16
- Convergence accommodation/convergence interaction (CA/C ratio), 13.24
- Convergence micropsia, 13.19
- Convergence responses to prism, 13.19
- Copepod, 8.1
- Copilia*, 8.4
- Cornea, 1.3*f*, 1.4–1.6, 16.3–16.4, 21.2
 absorption of ultraviolet light at, 1.9
 aging-related changes in, 14.5, 14.6*f*
 asphericity in, 1.5*f*
 and cataract (*see* Cataract)
 confocal microscopy of (*see* Confocal microscopes)
 as entrance pupil, 1.8
 injury to, 7.4
 keratometry of, 21.6
 laser ablation, 16.11–16.19
 ablation profiles, 16.14–16.15
 ablation rate, 16.16–16.18, 16.18*f*
 corneal photoablation, 16.16, 16.17*f*
 Epi-LASIK, 16.12, 16.13, 16.13*f*
 LASEK, 16.12, 16.13
 LASIK, 16.13–16.14, 16.14*f*
 photorefractive keratectomy (PRK), 16.11–16.12, 16.11*f*, 16.12*f*
 thermal, photochemical, and photoacoustic effects, 16.18–16.19
 optical axis of, 9.5
 optical power of, 16.2

- Cornea (*Cont.*):
 and refraction in the eye, 12.3
 refractive index of, 14.5
 refractive surgery modalities
 corneal incisions/implants, 16.9–16.11
 laser corneal procedures, 16.11–16.15
 and rigid contact lenses, 12.12
 spatially modulated excimer laser ablation of,
 1.25
 spectral sensitivities at, 10.18–10.19
 stray light from, 1.20
 Corneal endothelium, 16.4
 Corneal epithelium, 16.4
 Corneal hydration, 16.18
 Corneal incisions/implants, 16.9–16.11, 16.10f
 Corneal photoablation, 16.16–16.19, 16.17f
 ablation rate, 16.16–16.18, 16.18f
 thermal, photochemical, and photoacoustic
 effects, 16.18–16.19
 Corneal relaxing incisions, 21.13
 Corneal stroma, 16.4
 Coroneo effect, 7.5, 7.6f, 7.7
 Corresponding retinal points, 13.1, 13.8
 Cortical cataract, 12.14, 14.8
 Cortical neurons, 2.14
 Cotyledons, 8.1, 8.26, 8.27f
 Craik-O'Brien-Cornsweet illusion, 11.78
 Critical fusion frequency (CFF), 14.19
 Cross-coupling, direction/distance of gaze and,
 13.23–13.24, 13.23f
 Crystallin proteins, 19.1
 Crystalline lens, 1.3f, 16.3, 21.2
 and accommodation, 1.30–1.31
 aging-related changes in, 14.7–14.8, 14.17
 cataract of, 12.14, 21.3 (*See also* Cataract;
 Intraocular lenses)
 and cone spectral sensitivities, 10.17
 contours of refractive index in, 1.5f
 distribution of refractive index in, 1.5–1.6,
 1.5f
 and fluctuations in accommodation, 1.32
 fluorescence in, 1.21
 gradient index structure of, 1.18, 21.2
 growth of, 1.4
 and refraction in the eye, 12.3
 refractive index gradient of, 1.3, 19.12–19.13,
 19.13f, 19.14f
 scattered light from, 1.20
 UV absorption in, 1.9
 Cyclopean eye, 13.1
 Cyclopean locus, 13.7
 Cyclophoria, 13.35
 Cycloplegia, 1.25, 25.11
 Cyclovergence, 13.1, 13.22, 13.27
 Cytochrome oxidase “blobs,” 2.14
 Da Vinci stereopsis, 13.4
 Dark focus, 1.33, 1.34, 1.34f
 De Valois De Valois zone model, 11.82
 Decision tasks (in experiments), 3.2
 Defocus:
 and efforts to clear vision, 13.22–13.23
 with head-mounted displays, 25.9–25.10,
 25.9t–25.10t
 modulation transfer with, 1.29
 with monovision, 14.28
 off-axis, 1.18, 1.19
 Rayleigh units of, 1.29
 spectacles for, 15.2
 and visual acuity, 4.9, 4.9f, 4.10f
 with visual instruments, 1.28
 (*See also* Blur)
 Deformable mirrors, 15.1, 15.10, 15.10f
 Degauss, in color CRTs, 22.9
 Delta gun CRTs, 22.4, 22.13
 Depth dependent sensitivity (OCT),
 18.13–18.14, 18.13f, 18.14f
 Depth ordering, 13.4, 13.5, 13.10–13.12
 Depth perception, 2.39–2.40
 age-related changes in, 14.22
 stereodepth, 13.11–13.12
 Depth range:
 with OFDI, 18.16
 in SD-OCT, 18.6–18.7
 Depth-of-focus (DOF):
 in human eye, 1.28–1.29
 and pupil diameter, 1.8, 1.29, 1.30f
 Derrington Krauskopf Lennie (DKL) space,
 11.32, 11.32f, 11.33
 defined, 11.2
 equiluminant plane
 hue scaling, 11.63
 multiple mechanisms in, 11.57, 11.58
 Descemet's membrane, 14.5, 16.4
 Detection surface or contour (color vision), 11.12
 defined, 11.2
 and directions of color spaces, 11.40f
 in equiluminant plane, 11.44f
 and field adaptation, 11.53–11.54, 11.54f
 in the L,M plane, 11.40f, 11.41–11.42

- Detection tasks, 2.15, 3.2
- Detection threshold:
 - color vision
 - chromatic discrimination near, 11.58–11.59
 - two-stage model of, 11.23^f
 - and discrimination thresholds, 4.6–4.7
- Detectors, light, 5.21, 5.22^t
- Deuteranopes, 10.16
- DeVries-Rose law, 2.26
- Diabetic retinopathy, 14.1, 14.25–14.26
- Diaphany, 1.20
- Diattenuation, 18.1
- Dichroism, 18.1
- Dichromatic vision, 10.13, 10.14, 10.16
- Diffraction:
 - in multifocal optics, 21.15–21.18
 - and retinal image quality, 1.12–1.14, 1.21
- Diffraction limit, 4.1
- Diffraction limited imaging, 17.1
- Diffraction theory, 6.1
- Diffraction contact lenses, 14.28
- Diffusers, 5.10–5.11, 5.10^f
- Digital imaging:
 - color representations in, 10.35
 - displays for vision research
 - color cathode ray tubes (color CRTs), 22.1–22.34
 - liquid crystal displays (LCDs), 22.34–22.40
- Digital Libraries, 24.7–24.8, 24.10
- Diopter (D), 12.1, 12.4, 13.21, 25.1
- Dioptric errors of focus, 1.13, 1.14, 1.28
- Diplopia (double vision), 23.2
- Dipper function, 11.59, 11.60^f
- Direct ophthalmoscope, 12.5–12.6
- Direction, perception of, 13.7–13.10
 - corresponding retinal points, 13.8
 - horopter, 13.8–13.9, 13.9^f, 13.10^f
 - vertical horopter, 13.9
- Direction discrimination, 2.37
- Direction selectivity, 2.14ⁿ
- Disability glare, 14.4, 14.12, 23.2
- Discomfort glare, 14.12, 23.2
- Disconjugate eye movements, 13.1
- Discrete actuator deformable mirrors, 15.1, 15.10, 15.10^f
- Discrimination contours (color vision), 11.44^f
- Discrimination experiments (*see* Pedestal experiments)
- Discrimination tasks, 2.15, 3.2
- Discrimination thresholds, 4.6–4.7
- Diseases of the eye:
 - age-related, 14.3
 - age-related macular degeneration, 14.1, 14.24–14.25
 - cataract, 14.24
 - diabetic retinopathy, 14.25–14.26
 - glaucomas, 14.26–14.27
 - life-span environmental radiation damage, 14.22–14.23
 - visual impairment secondary to, 14.4
 - (*See also specific diseases*)
- Disjunctive eye movements (*see* Vergence eye movements)
- Displays for vision research:
 - color cathode ray tubes (color CRTs), 22.1–22.34
 - colorimetric calibration of, 22.20–22.34
 - design and operation of, 22.3–22.13
 - operational characteristics of, 22.13–22.18
 - setup for image display, 22.18–22.19
 - viewing environments, 22.19–22.20
 - liquid crystal displays (LCDs), 22.34–22.40
 - color LCDs, 22.37–22.40
 - monochrome, operational principles of, 22.34–22.37
- Disruptive movements (*see* Vergence eye movements)
- Distal stimuli (human vision), 4.2
- Distance, perception of, 13.24
- Distance conflicts, in head mounted display systems, 13.31–13.32
- Distortion, with head-mounted displays, 25.7
- Divergence, 13.21
 - defined, 13.1
 - in first 6 weeks of life, 13.21
- DKL color space, 10.19
- Doppler OCT, 18.18–18.19, 18.18^f–18.20^f
- Dot matrix, 23.2
- Dot pitch, 23.2
- Double-pass methods (retinal image quality), 1.22–1.23
- Drift, 1.44
- Droplet keratopathies, 7.6, 7.7
- Dry age-related macular degeneration, 14.1
- Dry eyes:
 - with increasing age, 14.5
 - as problem with computer work, 23.9

- Dysphotopsia:
 defined, 21.1
 with intraocular lenses, 21.21
- Eccentricity, 2.3, 2.7, 2.24, 2.25*f*
- Eclipse burn of the retina, 7.3
- Edge clearance, 20.1
- Edge image, in diffraction-limited eye, 1.13
- Edge thickness, for contact lenses, 20.6
- Effective power:
 defined, 20.1
 for spectacle and contact lenses, 20.8–20.12,
 20.9*t*, 20.10*f*, 20.11*f*
- “Effective” troland, 8.7
- Ego-center, 13.2, 13.7
- Ego-motion, 13.8, 13.9
- Electromagnetic radiation, 23.2
- Electronic imaging:
 human vision and (*see* Human vision and
 electronic imaging)
 in ophthalmoscopic methods, 1.23
- Ellipse blindness, 7.7
- Ellipsoid, 8.1
- Emmetropia, 1.6, 1.32, 12.4, 16.4, 16.5*f*
 age-related, 14.11
 defined, 6.1, 12.1, 13.2
 and focus of collimated light, 12.3*f*
 and Maxwellian viewing, 5.4–5.5
- Emmetropization, 1.7
- Emotional characteristics in electronic
 imaging (*see* Aesthetic and emotional
 characteristics)
- Empirical horopter, 13.8
- Empty field myopia, 1.33
- Endothelium:
 corneal, 16.4
 defined, 16.1
- Entrance pupil:
 and correction for SCE-1, 9.4
 defined, 13.2
 in Maxwellian viewing, 5.7
 in reflectometry, 8.6–8.7
 and retinal illuminance, 1.12
 and retinal irradiance, 2.4
 and stimulus specification,
 4.2–4.3, 4.3*f*
- Epicotyl, 8.26, 8.27*f*
- Epi-LASIK (epithelial laser in situ
 keratomileusis), 16.12, 16.13, 16.13*f*
- Epithelial layer (cornea), 14.5
- Epithelium:
 corneal, 16.4
 defined, 16.1
- Equatorial plane, 19.1
- Equilibrium colors:
 additivity with, 11.66
 and unique hues, 11.63–11.66
 white as, 11.63
- Equilibrium level (accommodation),
 1.33
- Equiluminance, difficulty in producing,
 11.80
- Equiluminant plane:
 detection and discrimination contours in,
 11.44*f*
 of DKL space
 hue scaling, 11.57
 multiple mechanisms in, 11.57, 11.58
- Equivalent reflectance, retinal, 1.11
- “Equivalent” troland, 9.2
- Equivalent veiling luminance, 1.20
- Equivalent-sphere correction, 1.6
- Ergonomics, 23.2
- Erythema, 7.1
- Estimation tasks, 2.15–2.16
- Etiolated plant tissues, 8.26–8.28
- Excimer lasers, 16.15–16.16
- Exhaustive characterization methods (ECM)
 (color CRTs), 22.21–22.23
 interpolation, 22.22–22.23
 inverses, 22.23
 out-of-gamut colors, 22.23
 sampling, 22.21–22.22
- Exit pupil, 2.3
 measurement of power at, 5.17
 in reflectometry, 8.6–8.7, 8.7*f*
 and retinal illuminance, 1.12
 and Stiles-Crawford effects, 1.11
- Experimental conditions, 3.2
- Exposure limits (ELs) (radiation), 7.9–7.11
 exceeding, 7.11
 for infrared, 7.10–7.11
 for laser light, 7.12
 for ultraviolet, 7.9–7.10
 for visible light, 7.10
- External limiting membrane (ELM), 8.8, 8.9,
 8.10*f*, 8.12
- Extracapsular cataract extraction (ECCE),
 12.14, 21.4
- Extraocular muscles, 1.42, 12.8, 12.16

- Extraretinal cues:
 defined, 13.2
 for eye movements, 13.7
 in space perception, 13.3
- Exudative (wet) age-related macular degeneration, 14.1, 14.24–14.35
- Eye alignment:
 binocular parallax, 13.22
 cross-coupling and direction/distance of gaze, 13.23–13.24, 13.23f
 defocus and efforts to clear vision, 13.22–13.23
 intrinsic stimuli to vergence, 13.21–13.22
 magnification induced errors of, 13.25–13.27
 in Maxwellian viewing, 5.8
 perceived distance, 13.24
 prism induced errors of, 13.27–13.29
 zone of clear and single binocular vision, 13.24–13.25, 13.25f
- Eye movements, 1.42–1.45, 13.19–13.20
 analysis of, 24.7
 characteristics of, 1.43–1.44
 in contrast detection, 2.21
 coordination of
 binocular parallax, 13.22
 cross-coupling and direction/distance of gaze, 13.23–13.24, 13.23f
 defocus and efforts to clear vision, 13.22–13.23
 intrinsic stimuli to vergence, 13.21–13.22
 perceived distance, 13.24
 zone of clear and single binocular vision, 13.24–13.25, 13.25f
- defined, 13.2
 extraretinal information for, 13.7
 gaze control, 13.29–13.30
 and optical flow, 2.39
 stability of fixation, 1.44, 1.45f
 types of, 13.19–13.20
- Eye protectors, for laser hazards, 7.14
- Eye relief, 1.7
- Eyes:
 bovine lenses, 19.8–19.11, 19.11f
 cat lenses, 19.9
 fish lenses, 19.6, 19.6f
 gibbon lenses, 19.14
 guinea pig lenses, 19.8
 of invertebrate aquatic worms, 9.14
 model, 21.13
 octopus lenses, 19.7, 19.7f
- Eyes (*Cont.*):
 pig lenses, 19.11, 19.12f
 porcine lenses, 19.11, 19.12f
 primate eye lens, 19.13f, 19.14, 19.14f
 rabbit lenses, 19.8
 rat lenses, 19.7–19.8
 (*See also* Human eye)
- Eyesight, defined, 23.2
- Eyestrain (asthenopia), 23.2
- F line, 25.1
- Far point, 12.1, 12.8, 16.5
- Farsightedness, 23.2 (*See also* Hyperopia)
- Fiber optic bundles, human receptors as, 8.3–8.5
- Fiber optic features:
 of plant tissues, 8.26–8.28, 8.27f
 of sponges, 8.28–8.29
- Fibroblast, 16.1
- Field additivity (color vision), 11.51, 11.52f
- Field measurements (color vision):
 and first-site adaptation, 11.27, 11.29, 11.30f
 and second-site adaptation, 11.29, 11.31
- Field method (color vision), 11.11–11.12
 defined, 11.2
 evidence for higher-order mechanisms from, 11.79–11.80
- Field of view, 8.4
 for head-mounted displays, 25.2, 25.5, 25.7–25.8, 25.9f
 in Maxwellian viewing, 6.6–6.7, 6.6f
- Field quality, in optical systems, 5.11
- Field sensitivities (color vision), 11.46–11.57
 achromatic detection and chromatic adaptation, 11.47–11.49, 11.48f
 chromatic adaptation and the Sloan notch, 11.49, 11.51
 detection contours and field adaptation, 11.53–11.54, 11.54f
 field additivity, 11.51, 11.52f
 first- and second-site adaptation, 11.51, 11.52, 11.53f
 habituation or contrast adaptation experiments, 11.54–11.56, 11.55f
 multiple cone inputs, 11.49, 11.50f
 noise-masking experiments, 11.56–11.57
 Stiles' π -mechanisms, 11.46, 11.47, 11.47f
 "Field" sensitivity method (color vision), 11.11–11.12
- Fields, in color CRTs, 22.9–22.10, 22.9f, 22.10f

- Filters, intensity, 5.14, 5.14*t*
- First-site adaptation (color vision), 11.11, 11.15–11.17, 11.16*f*
- defined, 11.2
 - field measurements and, 11.27, 11.29, 11.30*f*
 - and field sensitivities, 11.51, 11.52, 11.53*f*
 - incompleteness of, 11.17
 - signals reaching second site, 11.17
 - Weber's law and contrast coding, 11.15–11.16, 11.16*f*
- Fish lenses, 19.6, 19.6*f*
- Fish-eye lenses, 19.2–19.3
- Fixation, 1.42
- defined, 13.2
 - stability of, 1.44, 1.45*f*
- Flexure, 20.16
- Flicker photometrics, 10.34
- Flood-illuminated AO ophthalmoscope, 15.3, 15.16–15.17, 15.16*f*, 15.17*f*
- Fluence, 8.26
- corneal photoablation, 16.16, 16.17, 16.19*f*
 - defined, 8.1
- Fluorescence:
- lenticular, 1.21
 - stray light as result of, 1.21
- Fluoro-silicone/acrylate (F-S/A), 20.1
- Focal length, 23.2
- Focimeter, 12.9
- Focus:
- in color CRTs, 22.19
 - in Maxwellian viewing, 5.4–5.5, 5.5*f*, 6.5–6.6
 - in monochrome CRTs, 22.7
 - and responses to distance, 13.30
- Focus of expansion, 13.2
- Font, 23.2
- Foot-candle, 23.2
- “Forbidden” colors, 11.74
- Fourier approach to optics, 4.8
- Fourier domain OCT (*see* Spectral domain OCT)
- Fourier Theory of Optics:
- defined, 6.1
 - in Maxwellian viewing, 6.10–6.12
- Fourier transform, 2.23, 18.1
- Fovea, 1.3, 1.3*f*, 1.9, 1.15, 2.5*f*, 14.9
- and cone spectral sensitivities, 10.17
 - defined, 4.1, 18.1
 - disparities in, 2.40
 - function of, 13.3
 - receptor size in, 2.23
- Foveal avascular zone, 14.9
- Foveal pit, 14.11
- Foves, 2.24
- Frames, in color CRTs, 22.9–22.10, 22.9*f*, 22.10*f*
- Fraunhofer diffraction pattern, 6.10, 6.12
- Free (newtonian) viewing, 5.2–5.4, 5.3*f*
- limitations of, 5.4
 - retinal illuminance, 5.2–5.3, 5.3*f*
 - the troland, 5.3–5.4
- Fringe washout (OCT), 18.15
- Front vertex power (FVP):
- contact lenses, 20.7–20.8, 20.8*t*
 - defined, 20.1
- Frontoparallel plane, 13.2
- Gain control, in CRTs, 22.8
- Ganglion cells, 2.10–2.11
- contrast sensitivity functions of, 2.12*f*
 - density of, 2.11
 - and LGN, 2.12
 - linear density of, 2.6*f*
 - midget, 2.10–2.11, 2.10*n*
 - off-center, 2.10
 - on-center, 2.10
 - parasol, 2.10*n*
 - P-cells, 11.76, 11.79
 - physiology of, 2.11
 - transfer functions of, 2.11–2.12
- Ganzfeld, 1.12
- Gap effect (color vision), 11.72, 11.74
- Gas permeable (GP) contact lenses, 12.12, 20.3
- aberrations and, 20.24
 - aspheric, 20.22
 - base curve radius for, 20.3, 20.5
 - bitoric, power of, 20.17–20.20
 - center thickness of, 20.6
 - edge thickness of, 20.6
 - lacrimar or tear lens with, 20.12–20.15
 - OAD/OZD of, 20.5
 - posterior peripheral curve systems of, 20.5, 20.6
 - and residual astigmatism, 20.15
- Gaussian noise, 2.17, 2.18
- Gaze control, 13.29–13.30
- Gaze eccentricity, 13.8
- Gaze holding, 1.42
- Gaze shifting, 1.42
- Generalized pupil function, 2.3
- Geometric optics, 8.15
- Gibbon lenses, 19.14, 19.14*f*

- Glare:
 with computer work, 23.4–23.5, 23.5*t*
 defined, 23.2
 in human eye, 1.20
- Glass lenses, for spectacles, 12.9
- Glaucoma, 14.3, 18.25–18.27
 in aging eyes, 14.26–14.27
 cause of, 14.6
 defined, 14.1
- Gradient index (GRIN), 19.1, 19.2
- Gradient index optics, 19.1–19.15
 axial gradients, 19.5
 bovine lenses, 19.8–19.11, 19.11*f*
 cat lenses, 19.9
 eye lens, 19.5–19.6
 fish lenses, 19.6, 19.6*f*
 functional considerations, 19.14–19.15
 guinea pig lenses, 19.8
 human/primate lenses, 19.12–19.14, 19.13*f*,
 19.14*f*
 nature of index gradient, 19.21–19.15
 octopus lenses, 19.7, 19.7*f*
 pig lenses, 19.11, 19.12*f*
 rabbit lenses, 19.8
 radial gradients, 19.3–19.5, 19.4*f*, 19.5*f*
 rat lenses, 19.7–19.8
 spherical gradients, 19.2–19.3
- Grassmann's laws, 10.8–10.9, 10.28
- Grating acuity tasks, 2.34, 2.34*f*, 2.35*f*
- “Gray world” assumption, 10.39
- Guinea pig lenses, 19.8
- Habituation (color vision), 11.54–11.56,
 11.55*f*
 and color appearance, 11.29, 11.69–11.70
 defined, 11.2
 second-site, 11.19–11.20, 11.19*f*, 11.21*f*,
 11.22*f*
- Half-bleaching constant, 2.7
- Halftoning, low-level vision models in, 24.4
- Halos, 1.21
 in aging eyes, 14.13
 with glaucoma, 14.6
- Hankel functions, 8.1
 in modal analysis, 8.22
 for waveguides, 8.13
- Haploopia, 13.2
- Haptic, 21.1
- Heading judgments, magnification and,
 13.14–13.15
- Head-mounted displays (HMDs), 25.1–25.12
 accommodation with, 25.10–25.12
 binocular vision factors in design of,
 13.31–13.35
 distance conflicts, 13.31–13.32
 optical errors, 13.34–13.35
 spatial location conflicts, 13.33, 13.34*f*
 visual-vestibular conflicts, 13.32–13.33
 characterizing, 25.7–25.10, 25.8*f*,
 25.9*t*–25.10*t*
 common design considerations for, 25.2–25.7
 comfort, 25.2–25.5, 25.3*f*
 functionality, 25.5–25.6, 25.6*f*, 25.7*f*
 safety, 25.2
 usability, 25.2
- Heat sinks, 5.15
- Heidelberg Retinal Tomograph (HRT),
 17.7, 17.9
- Helmholtz's reciprocity theorem, 8.1, 8.25, 8.27
- Henle fibers, 8.2, 8.20
- Hering's law, 13.20, 13.26
- Hertz (HZ), 23.2
- Heterochromatic flicker photometry (HFP),
 11.37, 11.49
- Heterochromatic modulation photometry
 (HMP), 11.37
- Heterophoria, 13.2, 13.26
- Higher-order aberrations, 16.6, 16.7
 correction of, 16.9
 defined, 15.1, 20.2
- Higher-order mechanisms (color vision), 11.80
- Homogenous index, 19.1
- Horizontal cells (retina), 2.9
- Horizontal scanning, in color CRTs, 22.10–22.11,
 22.11*f*
- Horizontal/vertical size and position, in color
 CRTs, 22.12
- Horoapter, 2.40, 13.2, 13.8–13.9, 13.9*f*, 13.10*f*
- Hue cancellation, 11.26–11.27
- Hue scaling, 11.63, 11.64*f*
- Human eye, 1.1–1.45, 1.3*f*
 aberrations in, 1.3 (*See also* Aberrations
 [human eye])
 accommodation response, 1.29–1.36
 accuracy of, 1.32–1.34
 age-dependent changes in, 1.35–1.36
 application to instrumentation, 1.34–1.35
 dynamics of, 1.31–1.32
 stability of, 1.32
 vergence input, 1.34

- Human eye (*Cont.*):
- age-related changes in, 1.5, 1.7, 14.4–14.14
(*See also* Adaptive optics)
 - accommodation, 1.35–1.36, 1.35f
 - accommodation and presbyopia, 14.8, 14.9f
 - anterior and posterior chambers, 14.5, 14.6
 - cornea, 14.5, 14.6f
 - eye size, 14.11
 - fluorescence, 1.21
 - index of diffusion for, 1.23
 - lens, 14.7–14.8
 - and level of aberration, 1.18
 - pupil, 14.6–14.7, 14.7f
 - pupil diameter, 1.8, 1.18
 - retina, 14.9–14.11
 - retinal reflectance, 1.11
 - RMS wavefront error, 1.15–1.17, 1.16f
 - scattering, 1.20
 - tears, 14.4–14.5
 - transmittance, 1.9
 - transparency/cataract, 14.8
- ametropia, 1.6–1.7
- cornea, 16.3, 16.4, 16.4f
- damage from light exposure, 5.18–5.19, 5.18f
- depth-of-focus, 1.8, 1.28–1.29, 1.30f
- models of, 1.36–1.38
- paraxial, 1.36–1.37, 1.37f
 - wide-angle, 1.38
- monochromatic ocular aberrations, 1.4, 1.14–1.19
- off-axis, 1.18–1.19, 1.18f
 - on the visual axis, 1.15–1.18
- movements of, 1.42–1.45
- characteristics of, 1.43–1.44
 - stability of fixation, 1.44, 1.45f
- ocular parameters, 1.4–1.6, 1.4f, 1.5f
- ocular radiometry, 1.11–1.12
- optical components of, 2.2
- optics of, 16.2–16.3, 16.3f
- pupil diameter, 1.8–1.9
- refractive elements in, 12.3
- retinal illuminance, 1.11–1.12
- retinal image quality, 1.12–1.28
- in aberration-free eye, 1.12–1.14
 - calculation from aberration data, 1.21–1.22
 - chromatic aberration, 1.19–1.20, 1.19f
 - comparison between methods, 1.23
- Human eye, retinal image quality (*Cont.*):
- effects of aberration correction, 1.25–1.26
 - intraocular scattered light, 1.20–1.21
 - lenticular fluorescence, 1.21
 - monochromatic ocular aberrations, 1.14–1.19
 - observed optical performance, 1.23–1.25
 - ophthalmoscopic (double-pass) methods, 1.22–1.23
 - psychophysical comparison method, 1.22 and pupil diameter, 1.8
 - on the visual axis, 1.21–1.28
 - with visual instruments, 1.27–1.28
- retinal reflectance, 1.11
- size of, 14.11
- stereopsis, 1.38–1.42
- aniseikonia, 1.41–1.42
 - stereoscopic and related instruments, 1.41
 - tolerances in binocular instrumentation, 1.41–1.42
- Stiles-Crawford effect, 1.10–1.11
- structure of, 21.2–21.3
- transmittance, 1.9–1.11
- visual vs. optical axis orientation in, 1.3
(*See also specific topics*)
- Human hair, light guide effect in, 8.24–8.26
- Human vision and electronic imaging, 24.1–24.11
- analysis of image features, 24.6–24.9
 - attention and region of interest, 24.7
 - Digital Libraries applications, 24.7–24.8
 - user interface design, 24.8–24.9
 - visualization, 24.8
- information sources for, 24.11
- perception of imaging artifacts, 24.2–24.6
- early color vision, 24.5
 - embedding digital watermarks, 24.5
 - image quality and compression, 24.3–24.4
 - limitations of early vision models, 24.5–24.6
 - rendering and halftoning, 24.4
 - target detection in medical images, 24.5
- representation of aesthetic and emotional characteristics, 24.9–24.10
- color, art, and emotion, 24.10
 - image quality, 24.9
 - virtual reality and presence, 24.9
- Hydrogel contact lenses (*see* Soft contact lenses)
- Hydrophilic, defined, 16.1

- Hyperacuity, 2.34, 2.35, 4.13–4.16, 4.14f
 across ages, 14.22
 defined, 4.1
 and superresolution, 4.15
- Hypermetropia (hyperopia), 1.6, 16.5, 16.5f, 16.8
- Hyperopia, 1.18, 12.4
 defined, 13.2
 and focus of collimated light, 12.3f
- Hypocotyl, 8.2, 8.26, 8.27f
- Ideal observer, in judgment experiments, 3.6
- Ideal-observer models (color vision), 11.26
- Ideal-observer theory, 2.16–2.19, 2.19f, 2.24
- Identification tasks, 2.15, 2.16
- Illuminance:
 defined, 23.2
 measurement of, 5.17, 5.18
 retinal, 1.11–1.12
 SI units of, 5.3
- Illuminants:
 metamerism for, 10.38
 reflected from surfaces, 10.32–10.36
- Image artifacts:
 and perceived fidelity, 24.5
 perception of (*see* Perception, of imaging artifacts)
- Image compression, early techniques for, 24.4
- Image features in vision, 24.6–24.9
 attention and region of interest, 24.7
 Digital Libraries applications, 24.7–24.8
 user interface design, 24.8–24.9
 visualization, 24.8
- Image formation, 2.2–2.4
- Image overlap, 13.4
- Image quality:
 for head-mounted displays, 25.2, 25.5–25.7, 25.7f
 in human eye (*see* Retinal image quality)
 subjective approach to, 24.9
- Image sampling, by photoreceptors, 2.4–2.9
- Imaging science, 24.1
- Impulse response function (IRF), 14.19–14.21, 14.20f
- Incoherent light sources, 6.1
- Incoherent targets, pupil size and, 6.7–6.9, 6.8f, 6.9t
- Increment contrast, 2.31, 2.32
- Increment thresholds (color vision), 11.16f
- Incremental cone-excitation space, 11.2, 11.31, 11.32
- Infants:
 corneal endothelium in, 14.5
 lenses in, 14.8
 size of eye globe in, 14.11
- Influence function, 15.1
- Information theory, visual resolution and, 4.15–4.16
- Information-processing model, 2.15–2.16, 2.15f
- Infrared (IR):
 damage from, 7.2
 exposure limits for, 7.10–7.11
- Infrared cataract, 7.7, 7.10
- Inhibition, 2.25–2.27, 2.25f
- Injuries, radiation, 7.3–7.7, 7.5f
 cataract, 7.5–7.6, 7.6f
 droplet keratopathies, 7.6, 7.7
 infrared cataract, 7.7
 mechanisms of, 7.2–7.3
 action spectra, 7.2, 7.4f
 exposure duration and reciprocity, 7.2–7.3
 photokeratitis, 7.4
 photoretinitis, 7.7
 pterygium, 7.6, 7.7
- In-line gun CRTs, 22.5
- Input standards, for color LCDs, 22.39
- Instrument myopia, 1.34–1.35
- Integrating spheres, 5.10, 5.11
- Intensity (color vision):
 defined, 11.2
 dependence of hue on, 11.67–11.68, 11.67f
 of test stimulus, 11.12
- Intensity of light, controlling, 5.13–5.15, 5.14t
- Intensity scaling (light), 10.24, 10.25, 10.25f
- Interferometers, 5.24
- Interlaced monitors, 23.2
- International Commission on Illumination (*see* Commission Internationale de l'Éclairage)
- International Commission on Non-Ionizing Radiation Protection (ICNIRP), 7.9, 7.12
- International Electrotechnical Commission (IEC), 7.9, 7.12
- International Standardization Organization (ISO), 7.9
- Interocular aniso-magnification, distortion from, 13.16–13.19
 aniseikonia, 13.17–13.18
 interocular blur suppression with anisometropia, 13.18–13.19
 lenses and prisms, 13.16–13.17, 13.17f
- Interphotoreceptor matrix (IPM), 8.8, 8.9

- Interpupillary distance (IPD), 1.39–1.41, 1.40*f*
- Interstitial matrix, 8.2
- Intracapsular cataract extraction (ICCE), 12.14
- Intracorneal ring segments, 16.10–16.11, 16.10*f*
- Intraocular lenses (IOLs), 12.15, 16.9, 21.1–21.22
- accommodating, 14.29–14.30, 21.18–21.19
 - and aging of the eye, 21.3–21.4
 - aspheric, 21.10–21.12, 21.11*f*, 21.12*t*
 - for cataract correction, 14.8
 - and cataract surgery, 21.4
 - defined, 14.1–14.2, 21.1
 - design of, 21.5–21.20
 - accommodating lenses, 21.18–21.19
 - aspheric lenses, 21.10–21.12, 21.11*f*, 21.12*t*
 - chromophores, 21.20, 21.20*f*
 - IOL aberrations, 21.9–21.10, 21.10*f*
 - IOL power, 21.5–21.9, 21.7*t*, 21.8*f*
 - multifocal lenses, 21.14–21.18, 21.15*f*, 21.16*f*, 21.16*t*, 21.18*f*
 - phakic lenses, 21.19
 - testing IOLs, 21.13–21.14
 - toric lenses, 21.13, 21.12*f*
 - effective lens position (ELP), 21.7–21.9
 - multifocal, 21.14–21.18, 21.15*f*, 21.16*f*, 21.16*t*, 21.18*f*
 - phakic, 21.19
 - for presbyopic correction, 14.28–14.29
 - side effects of, 21.20–21.22
 - dysphotopsia, 21.21
 - posterior capsule opacification, 21.21–21.22
 - and structure of the eye, 21.2–21.3
 - toric, 21.13, 21.12*f*
- Intraocular pressure (IOP), 14.2, 14.26, 14.27, 16.3
- Intraocular scatter, age-related, 14.12
- Invariant hues:
 - and Bezold-Brücke effect, 11.67–11.68, 11.67*f*
 - defined, 11.2
- Invertebrate aquatic worm, eye structures in, 9.14
- Iris, 1.3*f*, 16.3, 21.2
- Irregular astigmatism, 16.6
- Isoaccommodation circle, 13.23, 13.23*f*
- Isochromatic CFSs, 2.30, 2.30*f*
- Isodiscrimination contours, 10.40
- Isoidential contours, 19.1
- Isolating direction (color vision), 11.2, 11.33
- Isoluminant CFSs, 2.30, 2.30*f*
- Isoluminant patterns, 2.30
- Isoplanatic patch, 15.5–15.7, 15.6*f*
- Isovergence circle, 13.23, 13.23*f*
- Jackson cross-cylinder check test, 12.7
- Javal's rule, 12.5
- Jones matrix formalism, 18.1, 18.20, 18.22–18.24
- Judd's three-stage Müller zone theory, 11.6, 11.7*f*, 11.8
- Judd-Vos modified 2° color-matching functions, 10.12, 10.45
- Judgment tasks:
 - adjustment tasks *vs.*, 3.2
 - psychophysical measurement of, 3.6–3.8
 - ideal observer, 3.6
 - rating scale, 3.6–3.8, 3.7*f*
 - response time, 3.8
 - two-alternative forced choice (2afc), 3.8
 - yes-no, 3.6
- Just noticeable distortion level (JND), 24.4
- “K,” 20.2
- Keratoconus, 9.2, 20.24
- Keratocytes, 16.1, 16.4
- Keratometer mire, 12.5
- Keratometry, 12.5, 21.6
- Kinetic cues, for perceived space, 13.4–13.7, 13.5*f*, 13.6*f*
- Knapp's law, 13.17, 13.26
- Kohler illumination, 5.6
- La Jolla colorimeters, 5.10
- Lacrimal lens:
 - and contact lens power, 20.12–20.15, 20.12*f*–20.14*f*
 - defined, 20.2
- Lag, 23.2
- Lamellas, 16.1, 16.4
- Lamp safety standards, 7.14–7.15
- Landolt C test, 4.8
- Landolt ring target, 2.34
- Large field (color matching):
 - defined, 10.9
 - standards for, 10.11, 10.12–10.13
- LASEK (laser subepithelial keratomileusis), 16.12, 16.13

- Laser ablation, **16.11–16.19**
 ablation profiles, **16.14–16.15**
 ablation rate, **16.16–16.18, 16.18f**
 corneal photoablation, **16.16, 16.17f**
 Epi-LASIK, **16.12, 16.13, 16.13f**
 LASEK, **16.12, 16.13**
 LASIK, **16.13–16.14, 16.14f**
 photorefractive keratectomy (PRK),
16.11–16.12, 16.11f, 16.12f
 thermal, photochemical, and photoacoustic
 effects, **16.18–16.19**
- Laser photocoagulation, **14.26**
- Laser ray tracing, **16.7**
- Laser scanning confocal microscopes, **17.1, 17.3, 17.7–17.8**
- Laser scanning ophthalmoscopes, **17.3**
- Lasers:
 eye injury from, **7.3, 7.4**
 hazards related to, **7.11–7.14, 7.12f, 7.13f**
 accidents, **7.12–7.14**
 and eye protectors, **7.14**
 safety standards, **7.12**
 and speckle fields, **5.21**
 types of, **5.20t**
- LASIK (laser-assisted keratomileusis) surgery,
1.6, 12.14, 16.13–16.14, 16.14f, 23.2
- Lateral aniseikonia, **13.17**
- Lateral chromatic aberration (*See* Transverse
 chromatic aberration)
- Lateral geniculate nucleus (LGN), **2.12–2.14**
 anatomy of, **2.12–2.13, 2.13f**
 physiology of, **2.13–2.14**
- Le Grand eye model, **2.3**
- LEDs, as light source, **5.15**
- Lens capsule, **21.1, 21.2, 21.21–21.22**
- Lens pigment density:
 and color matches, **10.9**
 and variations in color matching, **10.15**
- Lenses:
 bovine, **19.8–19.11, 19.11f**
 cat, **19.9**
 computer-assisted surfacing technologies
 for, **12.9**
 for correction of refractive error, **12.4**
 distortion from interocular aniso-
 magnification, **13.16–13.17**
 effects on vergence and phoria, **13.25–13.27**
 fish, **19.6, 19.6f**
 fish-eye, **19.2–19.3**
 gibbon, **19.14, 19.14f**
- Lenses (*Cont.*):
 guinea pig, **19.8**
 of human eye, **1.3f, 19.5–19.6** (*See also*
 Cornea; Crystalline lens; Gradient index
 optics)
 lacrimal, **20.2, 20.12–20.15, 20.12f–20.14f**
 for Maxwellian viewing systems, **5.4**
 octopus, **19.7, 19.7f**
 in optical systems, **5.11**
 porcine, **19.11, 19.12f**
 primate, **19.13f, 19.14, 19.14f**
 rabbit, **19.8**
 rat, **19.7–19.8**
 spherocylindrical, **12.4, 12.5**
 stimulus value of, **13.21**
 for vision correction (*see* Contact lenses;
 Intraocular lenses; Spectacles)
 Wyszceki, **9.5**
- Lensometer, **12.9**
- Lenticular absorption, **1.9**
- Lenticular fluorescence, **1.21**
- Leray, P., **24.7**
- Letter acuity task, **2.34**
- Life-span environmental radiation damage,
14.22–14.23
- Light, defined, **4.1, 23.3**
- Light output:
 in color CRTs
 spatial characteristics of, **22.14–22.15**
 spatial uniformity of, **22.17–22.18, 22.17f, 22.18f**
 stability of, **22.16–22.17, 22.16f, 22.17f**
 in color LCDs
 spatial variations in, **22.40**
 temporal variations in, **22.39–22.40**
- Light sources:
 coherent, **5.19, 5.21**
 control of (*see* Optical generation of visual
 stimulus)
 in Maxwellian viewing, **6.2**
 in vision laboratories, **5.19, 5.19t–5.20t**
- Lighting:
 for aging individuals, **14.4**
 for computer work, **23.4–23.5, 23.5t**
- Limbal relaxing incisions, **21.13**
- Limbus, **21.1**
- Limiting noise, site of, **11.24**
- Linear colorimetry models, **10.26–10.27, 10.27f**
- Linear models (matrix algebra), **10.1, 10.46–10.47**

- Linear optics, 17.1
- Linear perspective, 13.3
- Linear polarization (OCT), 18.21
- Linear visual mechanisms, 11.3
- Linearity of color-opponent mechanisms, 11.66–11.70
 - Bezold-Brücke effect and invariant hues, 11.67–11.68, 11.67*f*
 - color appearance
 - and chromatic adaptation, 11.68
 - and chromatic detection and discrimination, 11.69
 - and habituation, 11.69–11.70
 - luminance and brightness, 11.70, 11.70*f*
 - tests of linearity, 11.66–11.67
- Linearization, in optical systems, 5.15
- Line-spread function (LSF), 1.21
 - in diffraction-limited eye, 1.13
 - in ophthalmoscopic methods, 1.23
- Liou Brennan schematic eye, 15.6*f*, 15.7
- Liquid crystal displays (LCDs), 22.34–22.40
 - color, 22.37–22.40
 - colorimetry of color pixels, 22.38–22.39
 - controls and input standards, 22.39
 - geometry of color pixels, 22.37–22.38, 22.37*f*
 - spatial variations in output, 22.40
 - temporal variations in output, 22.39–22.40
 - and Computer Vision Syndrome, 23.7–23.8
 - defined, 23.2
 - monochrome, operational principles of, 22.34–22.37, 22.35*f*, 22.36*f*
 - in optical systems, 5.16
 - reflections from, 23.6
- Liquid spatial light modulators (LC-SLMs), 15.11
- Local characterization methods (LCM) (color CRTs), 22.24–22.27
 - individual colors, 22.24–22.26
 - inverses, 22.25
 - local regions of color, 22.24–22.26
 - one-dimensional color spaces, 22.25, 22.26
 - out-of-gamut colors, 22.25
- Localization, 4.14, 4.15, 8.4
- Longitudinal (axial) chromatic aberration (LCA), 1.19, 1.20, 1.28, 8.8, 15.22
- Longitudinal horopter, 13.8
- Long-range motion discrimination mechanism, 2.38
- Loss hypothesis, 10.14
- Low-level mechanisms (color vision), 11.80
- Low-order aberrations, 15.1
- Luminance, 11.37–11.39
 - cone numerosity, 11.38
 - defined, 2.29*n*, 11.3, 23.3
 - luminous efficiency, 11.37
 - luminous efficiency functions, 11.37, 11.38*f*
 - multiple cone inputs, 11.49, 11.50*f*
 - multiple luminance signals, 11.38–11.39, 11.40*f*, 11.70, 11.70*f*
 - and visual acuity, 4.11, 4.11*f*, 4.12
- Luminance mechanism, 11.11
- Luminance pedestals, 11.72, 11.74
- Luminance ratios, for computer work, 23.5
- Luminosity function, 10.10, 10.16
- Luminous efficiency, 11.37
 - and brightness, 11.70, 11.70*f*
 - and chromatic adaptation, 11.48*f*
 - variations in, 11.33
- Luminous efficiency functions, 10.13, 10.44–10.45, 11.33, 11.37, 11.38*f*
- Luminous flux, 23.3
- Luminous intensity:
 - defined, 23.3
 - in stimulus specification, 4.3–4.4
- Luneberg lens, 19.3
- Lux, 23.3
- MacAdam ellipses, 10.40
- Mach-Zender interferometers, 5.24
- Macula lutea (yellow spot), 7.1, 7.11
- Macular pigment (MP), 1.9, 1.11, 14.9–14.10
- Macular pigment density:
 - adjustments for individual differences, 10.17
 - and color matches, 10.9
 - and variations in color matching, 10.15
- Magnetic resonance imaging (MRI), 19.1
- Magnification, with contact lenses, 20.31–20.33
 - relative spectacle, 20.32–20.33, 20.33*f*
 - spectacle, 20.31–20.32
- Magnification and translation of images:
 - distortion by monocular magnification, 13.13–13.16
 - bifocal jump, 13.15–13.16, 13.15*f*
 - from convergence responses to prism, 13.19
 - discrepant views of objects/images, 13.16
 - motion parallax, 13.14–13.15
 - perspective distortion, 13.13, 13.14*f*
 - stereopsis, 13.13, 13.14

- Magnification and translation of images (*Cont.*):
 distortion from interocular aniso-
 magnification, 13.16–13.19
 aniseikonia, 13.17–13.18
 interocular blur suppression with
 anisometropia, 13.18–13.19
 lenses and prisms, 13.16–13.17, 13.17*f*
 errors of eye alignment induced by,
 13.25–13.27
- Magnitude (optical aberrations), 15.4
- Magnitude estimation, in psychophysical
 measurement, 3.8
- Magnitude production, in adjustment
 experiments, 3.5
- Magnocellular laminae, 2.12, 2.13, 2.13*f*
- Mandelbaum's phenomenon, 25.11–25.12
- Mandell-Moore bitoric lens guide, 20.17–20.19,
 20.18*t*, 20.19*f*
- Maréchal criterion, 1.15
- Masks (in experiments), 3.2
- Matching:
 in adjustment experiments, 3.4, 3.5
 color (*see* Color matching)
- Matrix algebra, 10.45–10.48
 addition and multiplication, 10.45–10.46
 colorimetry, 10.24–10.32
 stimulus representation, 10.24–10.27,
 10.25*f*, 10.26*f*
 transformations between color spaces,
 10.29–10.32
 vector representation of data, 10.25*f*,
 10.27–10.29, 10.27*f*
 glossary of notation in, 10.5*t*
 linear models, 10.46–10.47
 matrix transposition, 10.46
 simultaneous linear equations,
 10.47–10.48
 singular value decomposition, 10.48
 special matrices and vectors, 10.46
 vectors and matrices, 10.45
- Maximum saturation method (color
 matching), 10.6–10.7, 10.6*f*
- Maxwellian view (viewing), 5.4–5.8, 5.5*f*,
 6.1–6.14
 advantages of, 5.7
 control of focus, 5.4–5.5, 5.5*f*
 defined, 7.1
 field of view, 6.6–6.7, 6.6*f*
 focus, 6.5–6.6
 interferometers, 5.24
- Maxwellian view (viewing) (*Cont.*):
 partial coherence, 6.12–6.14
 positioning of pupil in, 5.9
 pupil size, 6.7–6.12
 coherent illumination, 6.9–6.12, 6.11*t*
 incoherent target, 6.7–6.9, 6.8*f*, 6.9*t*
 retinal conjugate plane, 5.5*f*
 retinal illuminance, 5.6–5.7, 6.3–6.5
 size, 5.5–5.6
 spatial frequency content of stimuli, 5.8
 two-channel, 5.21, 5.23–5.24, 5.23*f*
- Maxwell's method, for color matching,
 10.8, 10.8*f*
- McCollough effect, 11.76, 11.77*f*
- Mean-sphere correction, 1.6
- Measured horopter, 13.8
- Mechanism direction (color vision), 11.3, 11.33
- Mechanist approach (color vision) (*see* Color
 vision mechanisms)
- Melanin, 14.9–14.10
- Membrane mirrors (AO), 15.1, 15.10–15.12,
 15.10*f*
- Mesocotyl, 8.2
- Metamerism, 10.7, 10.36–10.38, 10.37*f*
- Metamers, 10.1, 10.7
- Method of constant stimuli, 3.9
- MHz (MegaHertz), 23.3
- Michaelis-Menton function, 2.14
- Michelson contrast, in vision experiments, 3.4
- Michelson formula, 4.8
- Michelson interferometer, 18.1
- Microelectromechanical systems (MEMS),
 15.11
- Microkeratomes, 16.13, 16.14
- Microlens Nipkow disk confocal microscope,
 17.1–17.2
- Micromachined membrane MEMS mirrors, 15.11
- Microsaccades, 1.44
- Midget ganglion cells, 2.10–2.11, 2.10*n*
- Midsagittal plane, 13.2
- Miniature eye movements, 1.44
- Minimally distinct border (MDB), 11.37
- Minimum motion (MM), 11.37
- Mode:
 defined, 8.2
 and energy density distribution, 8.15
 in hair cells, 8.25–8.26
 in human photoreceptors, 8.16*f*, 8.18*f*
 in monkey/human retinal receptors,
 8.19–8.24, 8.22*f*

- Model eyes, 21.13
- Model-dependent characterization methods (MDCM) (color CRTs), 22.27–22.33
 conditions for use, 22.30
 gun independence, 22.27–22.28
 inverse transformations, 22.32
 measurement of parameters, 22.31
 normalization coefficients, 22.31–22.32
 out-of-gamut colors, 22.32–22.33
 partial models, 22.30
 phosphor constancy, 22.28–22.29, 22.32
 phosphor output models, 22.29–22.30
- Modulation (contrast) color spaces, 10.19
- Modulation thresholds, 1.22
- Modulation transfer:
 with defocus, 1.29
 TCA effect on, 1.20
- Modulation transfer function (MTF):
 aberration-derived, 1.22
 in aging eyes, 14.13
 in diffraction-limited eye, 1.13, 1.13f
 double-pass, 1.23
 and observed optical performance, 1.24–1.25
 and off-axis image quality, 1.26–1.27
 in ophthalmoscopic methods, 1.23
 and optical quality of IOLs, 21.14
 and scattered light, 1.21
 of visual instruments, 1.28
 in young adult eyes, 1.24f
- Modulators, intensity, 5.14
- Monitors:
 and Computer Vision Syndrome, 23.6–23.8
 interlaced, 23.2
 noninterlaced, 23.3
 in optical systems, 5.16
- Monkey photoreceptors, 8.19–8.24, 8.22f
- Monochromatic ocular aberrations, 1.4, 1.14–1.19
 age-related, 14.12–14.14, 14.13f
 correction of, 1.25, 1.26, 1.26f
 off-axis, 1.18–1.19, 1.18f
 on the visual axis, 1.15–1.18
- Monochromatic vision, 10.16
- Monochrome CRTs, 22.3–22.4, 22.3f
 controls for, 22.7, 22.8f
 design and operation of, 22.3–22.4, 22.3f
 standards for, 22.14
- Monochrome LCDs, operational principles of, 22.34–22.37
- Monocular cues, for perceived space, 13.3–13.7
- Monocular field:
 horizontal angular extent of, 1.38f
 and stereopsis, 1.38–1.42
- Monocular magnification, distortion by, 13.13–13.16
 bifocal jump, 13.15–13.16, 13.15f
 from convergence responses to prism, 13.19
 discrepant views of objects/images, 13.16
 motion parallax, 13.14–13.15
 perspective distortion, 13.13, 13.14f
 stereopsis, 13.13, 13.14
- Mono vision:
 with contact lenses, 14.28
 defined, 12.13
 with head mounted visual displays, 13.32
- Motion artifacts (OCT), 18.15
- Motion detection/discrimination, 2.36–2.40
 optic flow fields, 2.39f
 thresholds, 2.37f
- Motion parallax, 13.4, 13.5
 defined, 13.2
 with monocular magnification, 13.14–13.15
- Motion perception, 13.6–13.7
- Moving entry and exit pupils method (photo-receptor directionality), 8.6
- Mueller matrix formalism, 18.1, 18.20, 18.22
- Müller cells, 8.20–8.21, 14.11
- Müller zone model, 11.65f, 11.84, 11.85f
- Müller zone theories, 11.6, 11.7f
- Multifocal lenses:
 contact lenses, 14.28, 20.2
 intraocular lenses, 21.1, 21.14–21.18, 21.15f, 21.16f, 21.16t, 21.18f
- Multiplicative adaptation, 2.26
- Multizone intraocular lenses, 14.28–14.29
- Munnerlyn's equations, 16.15
- Musculoskeletal, 23.3
- Myoid, 8.2
- Myopia (nearsightedness), 1.6, 1.18, 12.4, 16.5
 correction of, 13.17–13.18, 16.7–16.8
 defined, 13.2, 19.1, 23.3
 and early presbyopia, 12.8
 empty field, 1.33
 and focus of collimated light, 12.3f
 instrument, 1.34–1.35
 night (twilight), 1.34
 refractive surgery for, 12.14
- Naperian absorption coefficient, 8.9
- Near triad, 1.30

- Nearpoint, **23.3**
- Nearsightedness (*see* Myopia)
- Neuronal receptive fields, resolving capacity of the eye and, **4.6**
- Newtonian viewing (*see* Optical generation of visual stimulus, free (newtonian) viewing)
- Night (twilight) myopia, **1.34**
- Nipkow disk tandem-scanning confocal microscope, **17.2, 17.4f, 17.5, 17.5f**
- Noise (color vision), threshold contours for, **11.20, 11.23–11.26, 11.23f, 11.25f**
- Noise masking (color vision):
 defined, **11.3**
 experiments in, **11.56–11.57**
- Nomogram, **16.1**
- Nonconcomitant movement of eyes, **13.2**
- Nonexudative (dry) age-related macular degeneration, **14.1, 14.24–14.35**
- Noninterlaced monitors, **23.3**
- Nonlinear optics, **17.2, 17.3**
- Nonlinear visual mechanisms, **11.3**
- Nuclear cataract, **7.5, 12.14, 14.8**
- Nulling, in adjustment experiments, **3.4**
- Numerical aperture, **17.2**
- Nyquist limit, **2.6, 2.7, 2.11, 2.22–2.23**
- Obedience to Abney's law, **10.44, 11.37** (*See also* Additivity)
- Objective amplitude of accommodation, **1.32**
- Objective tasks, **2.15n**
- Oblique effect (pattern discrimination), **2.35**
- Octopus lenses, **19.7, 19.7f**
- Ocular motility, **23.3**
- Ocular parameters, **1.4–1.6**
- Ocular radiation hazards, **7.1–7.15**
 examples of, **7.8–7.9**
 exposure limits, **7.9–7.11**
 exceeding, **7.11**
 guidelines for visible light, **7.10**
 IR, **7.10–7.11**
 UV, **7.9–7.10**
 injury mechanisms, **7.2–7.3**
 action spectra, **7.2, 7.4f**
 exposure duration and reciprocity, **7.2–7.3**
 lamp safety standards, **7.14–7.15**
 laser hazards, **7.11–7.14, 7.12f, 7.13f**
 accidents, **7.12–7.14**
 and eye protectors, **7.14**
 safety standards, **7.12**
 retinal irradiance calculations, **7.7–7.8, 7.8f**
- Ocular radiation hazards (*Cont.*):
 types of injury, **7.3–7.7, 7.5f**
 cataract, **7.5–7.6, 7.6f**
 droplet keratopathies, **7.6, 7.7**
 infrared cataract, **7.7**
 photokeratitis, **7.4**
 photoretinitis, **7.7**
 pterygium, **7.6, 7.7**
- Ocular radiometry, **1.11–1.12**
- Ocular wavefronts, **16.6–16.7, 16.7t**
- Off-axis image quality, **1.18–1.19, 1.18f, 1.26–1.27, 1.27f**
- Olympus water immersion microscope, **17.10**
- One mode components analysis, **10.34**
- Ophthalmoheliosis, **7.1, 7.5**
- Ophthalmoscopes:
 defined, **15.1**
 flood-illuminated AO, **15.3, 15.16–15.17**
 scanning laser, **15.2, 15.3, 15.17–15.19, 15.18f, 15.19f**
- Ophthalmoscopic (double-pass) methods (retinal image quality), **1.22–1.23**
- Opponent color spaces, **10.18–10.19**
- Opponent-colors theory, **11.3, 11.5, 11.6f, 11.62–11.63**
- Optic disc, **1.3f**
- Optic fiber, defined, **19.1**
- Optic flow, **13.4**
 and bifocal jump, **13.15, 13.16**
 defined, **13.2**
- Optic flow fields, **2.38–2.39, 2.39f**
- Optic nerve head, **18.2**
- Optical aberrations:
 categories of, **17.2**
 in human eyes (*see* Aberrations [human eye])
- Optical axis, **1.3f**
- Optical circulator, **18.2**
- Optical coherence tomography (OCT), **15.3, 15.19, 18.1–18.30**
 at 1050 nm, **18.15–18.17, 18.15f–18.17f**
 autocorrelation noise, **18.11–18.12**
 combined with adaptive optics, **15.19–15.21, 15.20f, 15.21f**
 defined, **15.2, 18.2**
 depth dependent sensitivity, **18.13–18.14, 18.13f, 18.14f**
 Doppler OCT, **18.18–18.19, 18.18f–18.20f**
 fringe washout, **18.15**
 motion artifacts, **18.15**

- Optical coherence tomography (OCT) (*Cont.*):
 optical frequency domain imaging,
 18.7–18.9, 18.7f, 18.8f
 polarization sensitive OCT, 18.18, 18.20–18.27,
 18.21f, 18.23f, 18.25f–18.27f
 shot-noise-limited detection, 18.12–18.13
 signal to noise ratio, 18.11
 spectral domain OCT, 18.5–18.7, 18.5f,
 18.6f, 18.9
 noise analysis of, 18.9–18.10, 18.12f
 retinal imaging with, 18.27–18.29, 18.28f,
 18.29f
 sensitivity advantage of, 18.9
 Stratus OCT 3, 15.20–15.21, 15.21f
 time domain OCT, 18.3–18.5, 18.4f
- Optical design:
 binocular vision factors in, 13.1–13.35
 and refractive errors, 12.1–12.17
 assessment of, 12.5–12.8
 binocular factors, 12.15–12.17
 correction of, 12.8–12.15
 types of, 12.4–12.5
- Optical errors, in head mounted display
 systems, 13.34–13.35
- Optical fiber perform, 19.1
- Optical fibers, 5.10, 5.11, 19.4
- Optical frequency domain imaging (OFDI),
 18.3, 18.7–18.9, 18.7f, 18.8f
 at 1050 nm, 18.15–18.17
 defined, 18.2
 SD-OCT vs., 18.9
- Optical generation of visual stimulus,
 5.1–5.25
 building an optical system, 5.8–5.18
 alternating of source and retinal planes,
 5.8–5.9
 calibration, 5.17–5.19, 5.17f
 combining lights, 5.9–5.11, 5.9f, 5.10f
 controlling intensity, 5.13–5.15
 controlling wavelength, 5.11–5.13, 5.12t
 field quality, 5.11
 generating complex patterns, 5.16, 5.16f
 lenses, 5.11
 turning field on/off, 5.13, 5.13t
 coherent radiation, 5.19, 5.21
 detectors, 5.21, 5.22t
 free (newtonian) viewing, 5.2–5.4, 5.3f
 limitations of, 5.4
 retinal illuminance, 5.2–5.3, 5.3f
 the troland, 5.3–5.4
- Optical generation of visual stimulus (*Cont.*):
 light exposure and optical safety, 5.18–5.19,
 5.18f
 light sources, 5.19, 5.19t–5.20t
 Maxwellian viewing, 5.4–5.8, 5.5f
 advantages of, 5.7
 control of focus, 5.4–5.5, 5.5f
 interferometers, 5.24
 positioning of pupil in, 5.9
 retinal conjugate plane, 5.5f
 retinal illuminance, 5.6–5.7
 size, 5.5–5.6
 spatial frequency content of stimuli, 5.8
 two-channel, 5.21, 5.23–5.24, 5.23f
 size of stimulus, 5.2
- Optical power (of cornea), 16.2
- Optical safety, 5.18–5.19, 5.18f
- Optical Stiles-Crawford effect, 9.2
- Optical transfer function (OTF), 2.21–2.22,
 2.22f
 defined, 4.1
 for human eye, 1.21–1.22
 in ophthalmoscopic methods, 1.23
 and optical quality of IOLs, 21.14
 with visual instruments, 1.27
- Optical zone diameter (OZD):
 for contact lenses, 20.5
 defined, 20.2
- Opticaldiagnostics.com, 12.17
- Optics:
 of aging eyes, 14.4–14.14
 accommodation and presbyopia,
 14.8, 14.9f
 anterior and posterior chambers, 14.5, 14.6
 cornea, 14.5, 14.6f
 eye size, 14.11
 lens, 14.7–14.8
 pupil, 14.6–14.7, 14.7f
 retina, 14.9–14.11
 tears, 14.4–14.5
 transparency/cataract, 14.8
 of photoreceptors, 8.3, 8.10f
 (*See also* Adaptive optics)
- Optokinetic nystagmus, 1.44
- Optokinetic reflex, 13.20
- Organ of Corti, 8.2, 8.24–8.26
- Orthokeratology, 12.12
- Overall diameter (OAD):
 for contact lenses, 20.5
 defined, 20.2

- Overscan, in color CRTs, 22.12
 Oxygen permeability (Dk), 20.2
- Panum's fusional area (PFA), 12.15, 13.12, 25.5
 Parabolic louver, 23.3
 Parasol ganglion cells, 2.10*n*
 Paraxial models (of human eye), 1.36–1.37, 1.37*f*
 Partial coherence, in Maxwellian viewing, 6.12–6.14
 Partial coherence interferometry, 21.7
 Parvocellular laminae, 2.12, 2.13, 2.13*f*
 Pattern discrimination, 2.35–2.36
 Pedestal contrast, 2.31
 Pedestal control, in CRTs, 22.8
 Pedestal effects (color vision), 11.60*f*
 and chromatic discrimination, 11.69
 in crossed conditions, 11.59
 defined, 11.3
 and gap effect, 11.72, 11.74
 in uncrossed conditions, 11.59, 11.61
 Pedestal experiments, 3.2, 11.59–11.62, 11.60*f*
 Pedicles (cones), 8.2, 8.20
 Percept rivalry suppression, 13.14
 Perception:
 of color, 14.17
 defined, 13.2, 23.3
 of direction, 13.7–13.10
 corresponding retinal points, 13.8
 horopter, 13.8–13.9, 13.9*f*, 13.10*f*
 vertical horopter, 13.9
 of distance, 13.24
 Gestalt principles in, 24.8
 of imaging artifacts in vision, 24.2–24.6
 early color vision, 24.5
 embedding digital watermarks, 24.5
 image quality and compression, 24.3–24.4
 limitations of early vision models, 24.5–24.6
 rendering and halftoning, 24.4
 target detection in medical images, 24.5
 measuring (*see* Psychophysical measurement)
 of size, 13.19
 of space, 13.3–13.7
 binocular cues, 13.7
 distortion of, 13.16–13.19
 extraretinal information for eye movements, 13.7
 of space (*Cont.*):
 kinetic cues, 13.4–13.7, 13.5*f*, 13.6*f*
 monocular cues, 13.3–13.7
 and visual effects created by artists, 24.10
 Percepts:
 3-D, 13.3
 defined, 13.2
 Perceptual Subband Image Coder (PIC), 24.4
 Perceptually based image compression, 24.4
 Performance, measuring (*see* Psychophysical measurement)
- Perimetry, 14.21
 Peripheral field, 1.3
 and retinal illuminance, 1.12
 TCA in, 1.20
 Peripheral retina, 13.3, 14.10
 Perspective:
 defined, 13.2
 distorted, with monocular magnification, 13.13, 13.14*f*
 Phacoemulsification, 12.14, 21.2, 21.4, 21.5
 Phakia, 14.2
 Phakic lenses:
 defined, 21.2
 intraocular, 16.9, 21.19
 Phase discrimination, 2.36
 Phase retardation, 18.2
 Phase stability, Doppler OCT and, 18.18–18.19, 18.18*f*–18.20*f*
 Phase transfer function (PTF):
 in diffraction-limited eye, 1.13
 and off-axis image quality, 1.26, 1.27
 in ophthalmoscopic methods, 1.23
 Phoria, 13.21–13.22
 defined, 13.2
 effect of lenses and prisms on, 13.25–13.27
 Phoropter, 12.6
 Phosphor, 23.3
 Photocoagulation, 7.4
 Photographic recording, in ophthalmoscopic methods, 1.23
 Photokeratitis, 1.9, 7.4*f*
 defined, 7.1
 from radiation, 7.4
 UV, 7.3
 Photometric efficiency (PE) factor, 9.3
 Photometry, 9.1–9.15, 10.10, 10.43–10.45, 11.37 (*See also* Chapter 37 in Volume II)
 Photon noise, 2.4
 Photon-flux irradiance (retina), 2.4, 2.7

- Photophobia, 23.3
- Photopic luminosity function, 10.1, 10.4*t*
- Photopic luminous efficiency function, 10.44
- Photopic retinal illuminance, 2.4
- Photopic troland, 9.2
- Photopic vision:
 - accommodation response, 1.32–1.34, 1.33*f*, 1.34*f*
 - age-related changes in, 14.15, 14.16*f*
 - color, 10.3
 - Stiles-Crawford effect, 1.10, 1.10*f*
- Photopigment opsin genes, 10.14
 - in color-deficient observers, 10.16
 - and variations in color matching, 10.15
- Photopigment optical density:
 - and adjustment of cone spectral sensitivities, 10.17
 - adjustments for individual differences, 10.18
 - and color matches, 10.9
- Photopigments:
 - absorption spectra of, 2.7
 - adjustments for individual differences, 10.18
 - and color matches, 10.9
 - in color-deficient observers, 10.16
 - time constant of regeneration, 2.7
 - variability in, 10.15
- Photoreceptors, 10.3, 10.4
 - biological waveguide models of, 8.8–8.9, 8.9*f*, 8.10*f*, 8.12–8.15
 - and color appearance, 11.62
 - and contrast detection, 2.22–2.23
 - defined, 15.2
 - directional sensitivity of, 8.5
 - dynamic range of, 2.9
 - image sampling by, 2.4–2.9
 - inhomogeneous, 10.16
 - length of outer segments, 2.7–2.8
 - modal patterns in, 8.16*f*, 8.18*f*, 8.19–8.24, 8.22*f*
 - for neighboring waveguides, 8.22
 - radial distribution of transmitted energy, 8.22–8.23, 8.23*f*
 - Snyder and Pask cone model, 8.21*f*
 - transfer function, 8.24*f*
 - optical waveguide properties of, 14.11
 - optics of, 8.3, 8.10*f*
 - orientation and alignment of, 8.5–8.8, 8.6*f*, 8.7*f*
 - quantitative observations of, 8.15, 8.16*f*, 8.17, 8.18*f*, 8.19*f*
- Photoreceptors (*Cont.*):
 - and resolving capacity of the eye, 4.5, 4.6*f*
 - schematic diagram of, 8.9*f*
 - temporal response properties of, 2.8–2.9
 - types and functions of, 2.4, 2.6
 - waveguiding in, 8.3–8.5
 - (*See also* Cones; Rods)
- Photorefractive keratectomy (PRK), 12.14, 16.11–16.12, 16.11*f*, 16.12*f*
- Photoretinitis, 7.1, 7.7
- Phototropism, 8.2
- Physiological optics, 2.2
- Pig lenses, 19.11, 19.12*f*
- Pigments:
 - macular, 1.9, 1.11, 14.9–14.10
 - of the retina, 14.9–14.11
- Pixel, 23.3
- Plant tissues, fiber-optic, 8.26–8.28
- Plastic lenses, for spectacles, 12.9
- Poincaré sphere, 18.2
- Point-spread function (PSF), in human eye, 1.21
 - defined, 4.1, 15.2
 - in diffraction-limited eye, 1.12–1.14
 - direct measurements of, 2.3–2.4
 - and image formation, 2.3
 - and resolving capacity of the eye, 4.4, 4.5
- Polarity, 23.3
- Polarization, of light entering eye, 1.10
- Polarization sensitive OCT (PS-OCT), 18.18, 18.20–18.27, 18.21*f*, 18.23*f*, 18.25*f*–18.27*f*
- Polycarbonate lenses, for spectacles, 12.9
- Polymethylmethacrylate (PMMA) contact lenses, 12.11–12.12, 20.3
- Porcine lenses, 19.11, 19.12*f*
- Posterior capsule opacification (PCO):
 - defined, 21.2
 - with intraocular lenses, 21.21–21.22
- Posterior chamber, 14.5, 14.6, 16.3
- Posterior limiting lamina (*see* Descemet's membrane)
- Posterior peripheral curve (contact lenses), 20.5, 20.6, 20.23–20.24, 20.24*t*
- Posterior subcapsular cataract, 14.8
- Prentice's rule, 12.16, 13.15, 13.18, 13.28, 20.31
- Presbyopia (old eye/old sight), 1.7, 12.3, 16.5
 - accommodation restoration for, 14.29–14.30
 - as age-related, 14.8, 14.9*f*
 - assessment of, 12.7–12.8
 - clinical onset of, 12.3, 12.7–12.8

- Presbyopia (old eye/old sight) (*Cont.*):
 correction of, 14.27–14.30
 accommodation restoration, 14.29–14.30
 bifocals for, 12.8
 contact lenses for, 12.13–12.14,
 14.27–14.28
 hyperchromatic lens design for, 14.14
 intraocular lenses for, 14.28–14.29
 noncataract related refractive surgeries for,
 14.29
 spectacles for, 12.10–12.11, 14.27
 defined, 12.1, 14.2, 21.2, 23.3
 as problem with computer work, 23.11–23.12
 and UV radiation exposure, 14.23
- Primary lights (primaries), 10.4
 additive mixing of, 10.4, 10.5f
 CMFs for, 10.10, 10.11
 in color matching, 10.6, 10.7
 defined, 10.1, 10.4t
 for DKL color space, 10.19
 imaginary, 10.10
 in Maxwell's matching method, 10.8, 10.8f
 perceptual vs. physical matches of, 10.7
 for stimulus spaces, 10.11
 in trichromatic color matching, 10.7, 10.8
 and vector representations, 10.27–10.29
- Primary position, 1.42, 13.2
- Primary visual cortex, 2.12, 2.14
- Primate eye lens, 19.13f, 19.14, 19.14f
- Prism(s):
 in contact lenses, 20.2
 distortion from interocular aniso-magnifica-
 tion, 13.16–13.17, 13.17f
 effects on vergence and phoria, 13.25–13.27
 errors of alignment with, 13.27–13.29
 nonuniform magnification of, 13.10f
 stimulus value of, 13.21
- Prism diopter, 12.2, 25.1
- Prismatic effects, 12.8
 and bifocal jump, 13.15
 with contact lenses, 20.30–20.31
 prism-ballasted contacted lenses, 20.30
 unintentional (induced) prism, 20.31
 with corrected ametropia, 12.16
 defined, 12.2
- Prism-ballasted contacted lenses, 20.30
- Progressive addition lenses (PALs), 12.10–
 12.11, 12.11f
- Proportionality (color matching), 10.8
- Protanopes, 10.16
- Proximal (psychic) convergence, 1.35
- Proximal stimuli (human vision), 4.2
- Pseudoaccommodation, 14.27, 14.29, 14.30
- Pseudophakia, 12.15, 14.28
 correction of, 12.14–12.15
 defined, 14.2, 21.2
- Psychic convergence, 1.35
- Psychophysical comparison method (retinal
 image quality), 1.22
- Psychophysical measurement, 3.1–3.10
 of adjustment tasks, 3.4–3.5
 magnitude production, 3.5
 matching, 3.4–3.5
 nulling, 3.4
 threshold, 3.4, 3.5f
 definitions related to, 3.2–3.3
 of judgment tasks, 3.6–3.8
 ideal observer, 3.6
 rating scale, 3.6–3.8, 3.7f
 response time for, 3.8
 two-alternative forced choice (2afc), 3.8
 yes-no, 3.6
 magnitude estimation in, 3.8
 professional tips for, 3.10
 stimulus sequencing in, 3.9
 method of constant stimuli, 3.9
 sequential estimation methods, 3.9
 of visual acuity, 4.6–4.7
 visual stimuli, 3.3–3.4
- Psychophysical test method (color vision), 11.9,
 11.11, 11.12
- Psychophysics, 4.1, 15.2
- Pterygium, 7.1, 7.6, 7.7
- Punctum remotum*, 12.8
- Pupil:
 center of, 1.20
 defined, 6.1
 diameter of (*see* Pupil diameter/size)
 interpupillary distance, 1.39–1.41, 1.40f
 in Maxwellian viewing, 5.9
 in ophthalmoscopic methods, 1.23
 and retinal illuminance, 1.12
 true area of, 1.10
 and visual acuity, 4.9
 and visual instruments, 1.27–1.28
 (*See also* Entrance pupil; Exit pupil)
- Pupil conjugate plane:
 in Maxwellian viewing, 5.8–5.9
 rotating a mirror in, 5.16, 5.16f
 shutters at, 5.13, 5.13t

- Pupil diameter/size, 1.8–1.9, 1.8*f*
 age-related changes in, 1.18, 14.6–14.7, 14.7*f*, 14.14
 change in, 16.3
 and correction for refractive error, 1.23, 1.24
 and correction for SCE-1, 9.4, 9.13
 and depth-of-focus, 1.29, 1.30*f*
 in Maxwellian viewing, 6.3–6.5, 6.7–6.12
 coherent illumination, 6.9–6.12, 6.11*t*
 incoherent target, 6.7–6.9, 6.8*f*, 6.9*t*
 optimal, 15.22
 and RMS wavefront error, 1.15–1.16
 and Stiles-Crawford effect, 1.10
 variation in parameters as function of, 1.13
- Pupil function, OTF and, 1.21
- Purity, in color CRTs, 22.9
- Pursuit movements, 1.43
- QUEST, 3.9
- Rabbit lenses, 19.8
- Radial edge lift, 20.2
- Radial gradients, 19.3–19.5, 19.4*f*, 19.5*f*
- Radial keratotomy, 12.14, 16.9–16.10, 16.10*f*
- Radiation damage, life-span, 14.22–14.23
- Radiation hazards (*see* Ocular radiation hazards)
- Radiometry (*See also* Chapter 37 in Volume II):
 for color matching, 10.12
 correction for SCE-1 in, 9.1–9.15
 ocular, 1.11–1.12
- Raman microspectroscopy, 19.1
- Random-dot stereograms, 2.40
- Raster CRTs, 22.9–22.13, 22.13*f*
- Raster scanning, to despeckle light sources, 5.21
- Rat lenses, 19.7–19.8
- Rating scale judgments, 3.6–3.8, 3.7*f*
- Ray optics, 8.15
- Ray tracing, 19.1
- Rayleigh units, 1.13, 1.29
- Receiver operating characteristic (ROC) curves, 3.7–3.8, 3.7*f*
- Reciprocity theorem of optics (Helmholtz), 8.25, 8.27
- Red to far-red (R/FR) ratio, 8.28
- Reduced eye model, 1.36, 1.37, 2.3, 2.4
- Reflectance, retinal, 1.11
- Reflection:
 from computer screens, 23.5–23.6
 in fiber-optic plant tissues, 8.26, 8.28
 upon interfaces separating different media, 8.14–8.15, 8.14*f*
- Reflection model, 10.32
- Reflective apertures, 5.10
- Reflectometric sensors, 19.1
- Refraction:
 in human eye, 12.3, 16.1
 and retinal image quality, 1.21
 subjective, 12.6–12.7
 upon interfaces separating different media, 8.14–8.15, 8.14*f*
- Refractive, defined, 23.3
- Refractive ametropia, 20.2
- Refractive error correction, 12.8–12.15, 16.7–16.19
 ametropias, 16.8*f*
 aphakia and pseudophakia, 12.14–12.15
 with contact lenses, 12.11–12.14
 hydrogel, 12.12–12.13
 for presbyopia, 12.13–12.14
 rigid, 12.11–12.12
 with laser ablation, 16.11–16.19
 ablation rate, 16.16–16.18, 16.18*f*
 corneal photoablation, 16.16, 16.17*f*
 refractive surgery modalities, 16.11–16.15, 16.11*f*–16.14*f*
 thermal, photochemical, and photoacoustic effects, 16.18–16.19
- prescribing, 12.8
 with refractive surgery, 12.14, 16.9–16.15
 corneal incisions/implants, 16.9–16.11, 16.10*f*
 intraocular lenses, 16.9
 laser ablation modalities, 16.11*f*–16.14*f*
 modalities for, 16.11–16.15
 with spectacles, 12.9–12.11, 12.10*f*, 12.11*f*
- Refractive errors, 12.1–12.17, 12.3*f*, 12.4*f*, 16.4–16.19, 16.5*f*
 assessment of, 12.5–12.8
 objective tests, 12.5–12.6
 presbyopia, 12.7–12.8
 subjective techniques, 12.6–12.7
- astigmatism, 16.5–16.6
- binocular factors, 12.15–12.17
 aniseikonia, 12.16–12.17
 anisometropia, 12.16–12.17
 convergence and accommodation, 12.15–12.16

- Refractive errors (*Cont.*):
 consequences for optical design, 12.18
 correction of (*see* Refractive error correction)
 ocular wavefronts, 16.6–16.7, 16.7*t*
 as problem with computer work, 23.10
 spherical ametropias, 16.5
 types of, 12.4–12.5
- Refractive index, 8.12
 of cornea, 14.5
 defined, 19.1
 in lens of human eye, 1.5–1.6, 1.5*f*
 (*See also* Gradient index optics)
- Refractive lens exchange (RLE), 21.18
- Refractive power, defined, 19.1
- Refractive surgery, 1.15, 1.25, 12.14, 16.9–16.15
 change in cornea with, 16.3
 corneal incisions/implants, 16.9–16.11, 16.10*f*
 intraocular lenses, 16.9
 laser corneal procedures, 16.11–16.15
 ablation profiles, 16.14–16.15
 Epi-LASIK, 16.12, 16.13, 16.13*f*
 LASEK, 16.12, 16.13
 LASIK, 16.13–16.14, 16.14*f*
 photorefractive keratectomy (PRK), 16.11–16.12, 16.11*f*, 16.12*f*
 for presbyopic correction, 14.29
- Refresh rate, 23.3
- Region of interest, in human vision, 24.7
- Regular astigmatism, 1.6
- Relative directional sensitivity, 9.6
- Relative intensity noise (RIN):
 defined, 18.2
 with SD-OCT, 18.10
- Relative spectacle magnification (RSM), 20.2, 20.32–20.33
- Rendering, low-level vision models in, 24.4
- Residual astigmatism:
 and contact lens power, 20.15
 defined, 20.2
- Resolution (visual acuity):
 defined, 23.3
 and information theory, 4.15–4.16
 optics of resolving capacity, 4.4–4.5, 4.4*f*
 contrast-transfer function, 4.5
 point-spread function, 4.4, 4.5
 spatial-frequency coordinates, 4.5
 in retinal imaging, 15.21
 and superresolution, 4.15
 thresholds of, 4.6–4.7, 4.15
- Response time, as measure of performance, 3.8
- Response to direction (gaze control), 13.29–13.30
- Resting point of accommodation (RPA), 23.3
- Resting positions (eyes), 13.21
- Resting state (accommodation), 1.33
- Retina, 1.3*f*, 21.2
 aging-related changes in, 14.9–14.11
 anatomy of, 2.9–2.11
 AO-controlled light delivery to, 15.22–15.24
 alignment, 15.23
 in an AO SLO, 15.23
 conventional AO vision systems, 15.23
 to generate aberrations, 15.24
 longitudinal chromatic aberration, 15.22
 measuring activity of individual cones, 15.24
 transverse chromatic aberration, 15.22–15.23
 uses of, 15.23–15.24
 control of alignment in, 8.7
 direction-corresponding points, 13.8
 fovea, 2.24
 injury to, 7.4
 neural pathways in, 2.9, 2.10*f*
 nonfoveal areas of, 2.24
 OCT image of, 18.4*f*
 optic flow fields, 2.38–2.39, 2.39*f*
 physiology of, 2.11–2.12
 pigments of, 14.9–14.11
 and refraction in the eye, 12.3
 and resolving capacity of the eye, 4.5–4.6, 4.6*f*
 schematic diagram of, 8.9*f*, 8.10*f*
 UV light damage to, 14.23
- Retina cameras, AO, 15.3, 15.12 (*See also* Ophthalmoscopes)
- Retina pigment epithelium (RPE), 8.12
- Retinal burn, 7.7
- Retinal conjugate plane, in Maxwellian viewing, 5.5*f*, 5.8–5.9
- Retinal disorders, photoreceptor orientation/
 realignment after, 8.7
- Retinal disparity, 13.2
- Retinal eccentricity, visual acuity and, 4.10–4.11, 4.11*f*
- Retinal illuminance, 1.11–1.12
 in free (newtonian) viewing, 5.2–5.4, 5.3*f*
 maximum permissible, 5.18*f*
 in Maxwellian viewing, 5.6–5.7, 6.3–6.5
 in normal viewing, 6.3
 and pupil diameter, 1.8, 1.9

- Retinal image:
- AO ophthalmic applications, **15.16–15.22**
 - contrast and resolution, **15.21–15.22**
 - flood-illuminated AO ophthalmoscope, **15.16–15.17, 15.16f, 15.17f**
 - optical coherence tomography, **15.19–15.21, 15.20f, 15.21f**
 - scanning laser ophthalmoscope, **15.17–15.19, 15.18f, 15.19f**
 - light spread in, **4.4–4.5**
 - relating actual object and, **4.2**
- Retinal image disparity, **13.22**
- Retinal image quality, **1.12–1.28**
- in aberration-free eye, **1.12–1.14, 1.13f, 1.14f**
 - aging-related changes in, **14.11–14.14**
 - chromatic aberration, **14.14**
 - intraocular scatter, **14.12**
 - monochromatic aberrations, **14.12–14.14, 14.13f**
 - chromatic aberration, **1.19–1.20**
 - computing, **2.3**
 - intraocular scattered light, **1.20–1.21**
 - lenticular fluorescence, **1.21**
 - monochromatic ocular aberrations, **1.14–1.19**
 - off-axis, **1.18–1.19, 1.18f, 1.26–1.27, 1.27f**
 - on the visual axis, **1.15–1.18**
 - in peripheral field, **1.3**
 - and pupil diameter, **1.8**
 - variation with field location, **15.5**
 - on the visual axis, **1.21–1.28**
 - calculation from aberration data, **1.21–1.22**
 - comparison between methods, **1.23**
 - effects of aberration correction, **1.25–1.26**
 - observed optical performance, **1.23–1.25**
 - ophthalmoscopic (double-pass) methods, **1.22–1.23**
 - psychophysical comparison method, **1.22**
- Retinal imaging, with spectral domain OCT, **18.27–18.29, 18.28f, 18.29f**
- Retinal irradiance, **2.4, 7.7–7.8, 7.8f**
- Retinal layer of rods and cones model of
- biological waveguides, **8.8–8.9, 8.9f, 8.10f, 8.12–8.15**
 - assumptions and approximations for, **8.9, 8.12–8.13**
 - electromagnetic validity of, **8.13**
- Retinal microscopy:
- adaptive optics in, **15.1–15.24**
 - AO-controlled light delivery to the retina, **15.22–15.24**
 - control system, **15.12–15.15**
 - history of, **15.2–15.3**
 - imaging of the retina, **15.16–15.22**
 - implementation of, **15.7–15.15**
 - in ophthalmic applications, **15.15–15.24**
 - properties of ocular aberrations, **15.4–15.7, 15.5f, 15.6f**
 - wavefront corrector, **15.9–15.12, 15.10f, 15.11f**
 - wavefront sensor, **15.8–15.9, 15.8f**
 - defined, **15.2**
- Retinal nerve fiber layer (RNFL), **18.25–18.27**
- Retinal neurons, **2.5f**
- ganglion cells, **2.10–2.11**
 - information transmission by, **2.11**
- Retinal pigmented epithelium (RPE), **14.25, 18.2**
- Retinal processing, **2.9–2.12**
- Retinal reflectance, **1.11**
- layers of occurrence for, **1.23**
 - and ophthalmoscopic methods, **1.23**
 - and stray light, **1.20**
- Retinex theory (Land), **11.71, 11.72**
- Retinopathy:
- diabetic, **14.1, 14.25–14.26**
 - solar, **7.1, 7.3**
- Retinoscopy, **12.5–12.6**
- Rhabdomeric photoreceptors, **8.3**
- Rigid contact lenses:
- correction of refractive errors with, **12.11–12.12**
 - early types of, **20.2, 20.3**
 - (*See also* Gas permeable contact lenses)
- Rimless mounting, **12.2, 12.9**
- RMS (root mean square) contrast, **3.4**
- Rod amacrine cells, **2.5f**
- Rod bipolar cells, **2.5f**
- Rod pathway, **2.9, 2.10, 2.10f**
- Rods, **2.5f, 10.3**
- and age-related scotopic vision changes, **14.15**
 - alignment of, **8.4**
 - amacrine cells, **2.10**
 - and color matching, **10.17**
 - in color-deficient observers, **10.16**
 - diameter of, **2.6f**

- Rods (*Cont.*):
 directional sensitivity of, 8.5
 function of, 2.4
 linear density of, 2.6f
 and maximum saturation color matching, 10.7
 optical standing waves in, 8.17, 8.19f
 optical waveguide properties of, 14.11
 photocurrent responses of, 2.8f
 in retinal layer of rods and cones model of biological waveguides, 8.8–8.9, 8.9f, 8.10f, 8.12–8.15
 spatial distribution of, 2.6
 spectral sensitivities of, 10.18
 time constant of photopigment regeneration, 2.7
- Root-mean-square (RMS) wavefront error, 1.15–1.17, 1.16f
- Rostock Cornea Module (RCM) microscope, 17.7, 17.9
- Runs (trial sequences), 3.9
- Saccades, 1.42, 1.43, 1.43f, 13.20
- Saccadic suppression, 1.43
- Safety, of head-mounted displays, 25.2
- Safety standards, for laser hazards, 7.12
- Sagittal plane, 19.1
- Scanning laser ophthalmoscope (SLO), 15.3, 15.17, 17.9
 applying adaptive optics to, 15.17–15.19, 15.18f, 15.19f, 15.23
 AO-controlled stimulus delivery, 15.17–15.19
 visual acuity with, 15.23–15.24
 defined, 15.2
- Scanning slit confocal microscopes, 17.2, 17.3, 17.6–17.9, 17.6f, 17.7f
- Scattered light:
 with cataracts, 14.24
 intraocular, 1.20–1.21
 and OTF calculation, 1.21
 and psychophysical comparison method, 1.22
- Schematic eye model, 1.36, 1.37
- Schlemm's canal, 14.5
- Sclera, 16.3, 21.2
- Scotoma, 7.7
- Scotopic troland, 9.2
- Scotopic vision, 10.3
 age-related changes in, 14.15, 14.16f
 in color-deficient observers, 10.16
 Stiles-Crawford effect, 1.10, 1.10f
- Screen reflections, with computer work, 23.5–23.6
- Second-harmonic generation (SHG), 17.2, 17.9–17.10
- Second-site adaptation (color vision), 11.17–11.22, 11.18f, 11.19f
 desensitization by steady fields, 11.17, 11.18, 11.18f
 field measurements and, 11.29, 11.31
 and field sensitivities, 11.51, 11.52, 11.53f
 habituation, 11.19–11.20, 11.19f, 11.21f, 11.22f
- Second-site desensitization (color vision):
 defined, 11.3
 by steady fields, 11.17, 11.18, 11.18f, 11.54f
- Segmented corrector, 15.2
- Segmented piston and piston/tip/tilt wavefront correctors, 15.10, 15.10f
- Seidel aberrations, 12.9
- Self-motion, 13.8, 13.9
- Semiconductor optical amplifiers, 18.2
- Senescent changes in vision (*see* Age-related changes in vision)
- Sequential stimulus sequencing, 3.9
- Shack-Hartmann wavefront sensing (SHWS), 15.3, 15.8–15.9, 15.8f
- Shack-Hartmann wavefront sensors, 15.2, 16.7
- Shadowmask color CRTs, 22.4–22.6
 common problems in, 22.5, 22.6
 geometry of, 22.4–22.5, 22.5f, 22.6f
- Shear, 13.2
- Short-range motion discrimination mechanism, 2.38
- Shot noise, 18.9
 defined, 18.2
 with SD-OCT, 18.10
- Shot-noise-limited detection, 18.12–18.13
- Signal detection theory, 2.17, 11.20, 11.23f, 11.24
- Signal to noise ratio (SNR), 18.11
 defined, 18.2
 for in-vivo high speed human retinal imaging, 18.15
 and preferred frequency, 2.24
 with SD-OCT, 18.7, 18.10

- Silicone hydrogel contact lenses, **20.2, 20.3**
- Simple cells (cortical neurons), **2.14**
- Simplified schematic eye model, **1.36, 1.37**
- Simultaneous vision contact lenses,
12.13–12.14, 14.28, 20.22
- Site of limiting noise, **11.24**
- Size, perceived, **13.19**
- Skew movement, **13.2**
- Slit-lamp microscopes, **17.2**
- Slit-scanning arrangements, in
ophthalmoscopic methods, **1.23**
- Sloan notch, **11.34, 11.35f, 11.37**
and chromatic adaptation, **11.49, 11.51**
and chromatic discrimination, **11.57, 11.58**
- Small field (color matching):
defined, **10.9**
standards for, **10.11, 10.12**
- Snakes, effect of laser exposures in, **9.13**
- Snellen letters, **4.1, 4.8**
- Snell's law, biological waveguides and,
8.14, 8.15
- Snow blindness, **1.9, 7.3, 7.4**
- Soft (hydrogel) contact lenses:
aberrations and, **20.24**
base curve radius for, **20.3, 20.5**
center thickness of, **20.6**
correction of refractive errors with,
12.12–12.13
edge thickness of, **20.6**
OAD/OZD of, **20.5**
posterior peripheral curve systems of, **20.6**
power of, **20.15–20.16**
toric lenses
defined, **20.2**
power of, **20.16–20.17**
- Solar keratitis, **14.23**
- Solar retinitis, **7.7**
- Solar retinopathy, **7.1, 7.3**
- Space, perception of, **13.3–13.7**
binocular cues, **13.7**
distortion of, **13.16–13.19**
extraretinal information for eye movements,
13.7
kinetic cues, **13.4–13.7, 13.5f, 13.6f**
monocular cues, **13.3–13.7**
- Spatial channels, in contrast detection,
2.23–2.24
- Spatial equichromatic CSFs, **11.45–11.46, 11.46f**
- Spatial equiluminant CSFs, **11.45–11.46, 11.46f**
- Spatial fidelity (optical aberrations), **15.4**
- Spatial filtering:
with confocal microscopes, **17.3**
to despeckle light sources, **5.21**
- Spatial frequency, **4.1**
- Spatial frequency channels (in vision), **3.2**
- Spatial impulse response, of photoreceptors, **8.23**
- Spatial location conflicts, in head mounted
display systems, **13.33, 13.34f**
- Spatial sinusoids, **2.20**
- Spatial vision, age-related changes in, **14.7, 14.17–14.19, 14.18f**
- Spatial-frequency coordinates, resolving
capacity of the eye and, **4.5**
- Spatially-modulated excimer laser ablation of
cornea, **1.25**
- Spatio-temporal CSFs, **2.29**
- Speckle fields, **5.19, 5.21**
- Spectacle blur, **12.12**
- Spectacle lenses:
anisophoria and, **13.26**
base curve of, **12.1, 12.9**
effective power of, **20.8, 20.9t**
and gaze control, **13.29–13.30**
materials for, **12.9–12.10**
optical power of, **12.4**
power of, **20.6**
- Spectacle magnification (SM), **20.2, 20.31–20.32**
- Spectacles, **15.2**
correction by (*See also* Spectacle lenses)
for anisometropia, **12.16, 13.17–13.18**
aphakia, **12.15**
for presbyopia, **14.27**
for refractive errors, **12.9–12.11, 12.10f, 12.11f**
defined, **15.2**
and response to distance, **13.30**
- Spectral color sensitivities:
test method for, **11.34, 11.35f, 11.36f, 11.37**
in three-stage Müller zone model, **11.7f**
two-color measurements, **11.34**
- Spectral (Fourier) domain OCT, **18.5–18.7, 18.5f, 18.6f, 18.9**
defined, **18.2**
noise analysis of, **18.9–18.10, 18.12f**
OFDI vs., **18.9**
retinal imaging with, **18.27–18.29, 18.28f, 18.29f**
sensitivity advantage of, **18.9**
- Spectral power density, vector representation of,
10.24, 10.25f

- Spectral power distribution, **10.3**
 in colorimetric measurements, **10.23**
 and tristimulus values, **10.11, 10.36–10.37**
- Specular microscopes, **17.2**
- Specular reflection, **23.3**
- Speed discrimination, **2.37**
- Spherical ametropias, **1.6–1.7, 16.5**
 hypermetropia, **16.5**
 myopia, **16.5**
- Spherical gradients, **19.2–19.3**
- Spherocylindrical lenses, **12.4, 12.5**
- Spherule, **8.2**
- Sponges, **8.28–8.29**
- Square root law, **2.26**
- Stability of fixation, **1.44, 1.45f**
- Standard Notation (lenses), **12.5**
- Standard observer:
 CIE 1931, **10.12**
 CIE 1964, **10.13**
 defined, **10.1, 10.4t**
 limitations of, **10.15**
 standardized CMFs as, **10.9**
- “Star” patterns (human vision), **1.21**
- Steradian, **6.3**
- Stereoacuity, **2.41**
- Stereopsis, **1.38–1.42, 2.40**
 acuity, **1.39–1.40**
 age-related changes in, **14.22**
 aniseikonia, **1.41–1.42**
 Da Vinci, **13.4**
 defined, **13.2**
 factors affecting, **2.40, 2.41**
 with head mounted visual displays, **13.31, 13.32**
 with monocular magnification, **13.13, 13.14**
 with monovision, **14.28**
 in perception of depth, **13.11–13.12, 13.11f**
 stereoscopic and related instruments, **1.41**
 tolerances in binocular instrumentation, **1.41–1.42**
- Stereoscopic instruments, **1.41**
- Stiles and Burch (1955) 2° color-matching functions, **10.12**
- Stiles and Burch (1959) 10° color-matching functions, **10.12, 10.13**
- Stiles’ π -mechanisms, **11.16, 11.16f, 11.46, 11.47, 11.47f**
 and contrast coding, **11.16, 11.16f**
 and field additivity, **11.51**
 limitation of, **11.34**
- Stiles-Crawford effect, **1.29, 2.8**
 of the first kind, **1.21, 8.5** (*See also* Stiles-Crawford effect of the first kind)
 correction for, **9.1–9.15**
 just prior to phototransduction, **6.14**
 and Maxwellian viewing, **6.2**
 optical, **9.2**
 and pupil size, **6.7**
 of the second kind, **8.3, 8.5, 9.2**
 and human eye, **1.11**
 and waveguide models, **8.21**
- Stiles-Crawford Effect of the first kind (SCE I/SCE-1), **1.21, 8.5**
 and biological waveguides, **8.3, 8.6f**
 correction for, **9.1–9.15**
 adaptive optics techniques for, **9.5**
 application of approach, **9.13–9.14**
 confounds in, **9.4–9.5**
 history of, **9.5**
 with nonmonochromatic stimulus to vision, **9.5**
- Photometric efficiency factor, **9.3**
- sample point-by-point SCE-1 estimates, **9.6–9.13, 9.7f, 9.9f–9.13f**
- teleological and developmental factors in, **9.14**
 and trolands, **9.2, 9.3**
 without adaptive optics, **9.5–9.6**
 defined, **9.2**
 and human eye, **1.10, 1.10f, 1.11**
 and photoreceptor directionality, **8.6–8.8**
 and waveguide analysis, **8.21**
- Stimulus color spaces (colorimetry), **10.11**
- Stimulus sequencing, in psychophysical measurement, **3.9**
- Stimulus specification (visual acuity), **4.2–4.4, 4.3f**
- Stokes parameters, **18.20**
- Stokes vector, **18.2**
- Strabismus, **13.2**
- Strehl ratio, **1.23–1.24, 1.24f, 1.28, 15.2**
- Stroma (cornea), **14.5**
- Subadditivity, **11.51, 11.52f**
- Subjective amplitude of accommodation, **1.32**
- Subjective refraction, **12.6–12.7**
- Subjective tasks, **2.15n**
- Subtractive adaptation, **2.26**
- Superadditivity, **11.51, 11.52f**
- Superposition (colorimetry), **10.25, 10.26f**
- Superresolution (visual acuity), **4.1, 4.15**

- Suppression (binocular vision), **13.12–13.13**
of blur, **13.18–13.19**
defined, **23.3**
with head mounted visual displays, **13.33**
- Suprathreshold models (color vision), **11.41**
- Surface micromachined MEMS devices, **15.11**
- Surfaces:
illuminants reflected from, **10.32–10.36**
metamerism for, **10.38**
- Surround, visual acuity and, **4.12, 4.13f**
- SVGA (super video graphics array), **23.3**
- Svishchev confocal microscope, **17.6, 17.6f**
- Swept Source OCT (SS-OCT), **18.2** (*See also*
Optical frequency domain imaging)
- Symmetry (color matching), **10.8**
- Synchronization standards, for color CRTs,
22.13–22.14
- TABO, **12.5, 12.7**
- Taflove, A., **8.17**
- “Target” sensitivity method (color vision)
(*see* Test sensitivity method (color vision))
- Tear lens, **12.12**
- Tears, aging-related changes in, **14.4–14.5**
- Telecentric principle, **6.6–6.7**
- Temporal coherence, **18.2**
- Temporal contrast detection, **2.27–2.29, 2.27f, 2.29f**
- Temporal despeckling, **5.21**
- Temporal vision, age-related changes in,
14.19–14.21, 14.20f
- Test (“target”) sensitivity method (color
vision), **11.11, 11.12, 11.34–11.46**
defined, **11.3**
and different directions of color space, **11.39, 11.41–11.43, 11.44f, 11.45–11.46, 11.46f**
luminance, **11.37–11.39, 11.38f, 11.40f**
and spectral lights, **11.34, 11.35f, 11.36f, 11.37**
- Tests of visual acuity, **4.7–4.8, 4.7f**
for infants, **4.8**
traditional visual acuity chart, **4.6**
- Theoretical horopter, **13.8, 13.9f**
- Third-harmonic generation (THG), **17.2, 17.9–17.10**
- Three-dimensional color data, **10.19–10.20, 10.20f**
- Three-segment model of biological waveguides,
8.9, 8.11–8.15, 8.11f
assumptions and approximations for, **8.9, 8.12–8.13**
electromagnetic validity of, **8.13**
- Three-stage zone models (color vision), **11.6, 11.82–11.85, 11.85f**
De Valois and De Valois model, **11.82**
Müller model, **11.6, 11.7f, 11.65f, 11.84, 11.85f**
- Threshold experiments, **11.17**
- Threshold surface or contour (color vision),
11.12–11.15, 11.13f
defined, **11.3**
and loss of information, **11.20**
and noise, **11.20, 11.23–11.26, 11.23f, 11.25f**
and second-site adaptation to steady fields,
11.18f
- Threshold-limit values (TLVs), **7.9**
- Thresholds, **3.3f**
in adjustment experiments, **3.4**
defined, **3, 3.2**
detection (color vision)
chromatic discrimination near,
11.58–11.59
two-stage model of, **11.23f**
discrimination, **4.6–4.7**
estimation of, **3.9**
increment (color vision), **11.16f**
modulation, **1.22**
for motion detection/discrimination, **2.37f**
in psychophysics, **3.2–3.3**
suprathreshold models (color vision), **11.41**
of visual resolution, **4.6–4.7, 4.15**
- Tilt, **13.2**
- Time, visual acuity and, **4.12, 4.13f**
- Time domain OCT, **18.2–18.5, 18.4f**
- Timing standards, for color CRTs, **22.13–22.14**
- Tonic level of accommodation, **1.33–1.34**
- Toric lenses:
for astigmatism, **12.13**
bitoric, **20.17–20.20, 20.18f–20.20f, 20.18t**
intraocular, **21.13, 21.12f**
power of, **20.16–20.20, 20.16f, 20.17f**
prism with, **20.30**
- Torsion, **13.2, 13.8, 13.22**
- Torsional vergence, **13.22, 13.27**
- Transfer function (human vision), **2.3**
defined, **6.1**
of photoreceptors, **8.23–8.24, 8.24f**
temporal and spatial components of, **2.11, 2.12f**
- Transitivity (color matching), **10.8**
- Transmittance, in the human eye, **1.9–1.11, 1.9f**

- Transparency, aging-related changes in, **14.8**
(*See also* Cataract)
- Transparent HMDs, **25.4–25.5**
- Transverse (lateral) chromatic aberration (TCA), **1.19, 1.20, 15.22–15.23**
measurement of, **8.8**
with visual instruments, **1.28**
- Tremor, **1.44**
- Trials, experimental, **3.2**
- Trichromacy (trichromatic vision), **10.4–10.6**
color appearance *vs.*, **11.5**
in color matching, **10.7–10.9**
in color-deficient observers, **10.16**
- Trifocal lenses, **12.10, 14.27**
- Trinitron CRTs, **22.5**
- Tristimulus values (cone coordinates), **10.4**
for arbitrary lights, **10.9**
in colorimetric measurements, **10.23**
and consistency across observers, **10.9**
defined, **10.1, 10.4t**
in maximum stimulus method, **10.7**
negative, **10.7**
and spectral power distribution, **10.11**
spectral power distributions from, **10.36–10.37**
tailored to individuals, **10.15**
uniqueness of, **10.9**
in vector representations, **10.27–10.29**
- Tritanopes, **10.16**
- Trivex lenses, **12.9**
- The troland, **5.3–5.4, 6.3**
defined, **2.7n**
“effective,” **8.7**
“equivalent,” **9.2**
limitations of, **9.3**
as unit of retinal illuminance, **9.2**
- Twilight myopia, **1.34**
- Two-alternative forced choice (2afc), **3.8**
- Two-channel Maxwellian viewing, **5.21, 5.23–5.24, 5.23f**
- Two-color threshold experiments, **11.34, 11.37**
- Ultraviolet (UV) radiation, **14.23**
cataract from, **7.5–7.6**
damage from, **7.2**
exposure limits for, **7.9–7.10**
pterygium and droplet keratopathies from, **7.6, 7.7**
- Unaccommodated eye, **16.5**
- Underscan, in color CRTs, **22.12, 22.13f**
- Uniform color spaces, **10.40, 10.42**
- Unintentional (induced) prism, with contact lenses, **20.31**
- Unipolar chromatic mechanisms, **11.80–11.81**
- Unipolar mechanism (color vision), **11.3**
- Unique hues, **11.27**
in color discrimination *vs.* color appearance tests, **11.81–11.82**
defined, **11.3**
and equilibrium colors, **11.63–11.66**
zero crossings, **11.62–11.63**
- Unitary matrix, **18.2**
- Univariance, **10.4**
- Univariant mechanism (color vision):
defined, **11.3**
flicker, **11.45**
suprathreshold, **11.58**
- User interfaces, human image features and, **24.8–24.9**
- Van Cittert-Zernike theorem, **6.12**
- Vanishing point, **13.4, 13.4f**
- Vantage point, **13.2**
- Vergence (disjunctive) eye movements, **1.42, 1.44, 13.20, 13.21**
- Vergence system, **1.29, 1.30, 1.32, 1.34, 1.43–1.44**
binocular parallax, **13.22**
defocus and effort to clear vision, **13.22–13.23**
effect of lenses and prisms on, **13.25–13.27**
with head mounted visual displays, **13.31–13.32**
intrinsic stimuli to, **13.21–13.22**
and variations of cross-coupling, **13.23–13.24**
zone of clear and single binocular vision, **13.24–13.25, 13.25f**
- Vernier acuity, **4.1**
- Vernier acuity task, **2.34, 2.35f**
- Version eye movements (*see* Conjugate eye movements/position)
- Vertex distance, **12.9, 16.8**
- Vertical horopter, **13.9**
- Vertical scanning, in color CRTs, **22.11–22.12, 22.12f**
- Vestibulo-ocular reflex (VOR), **13.20**
- Vestibulo-ocular responses, **1.42**
- Video display terminal (VDT), **23.3**

- Video head sets (*see* Head-mounted displays)
- Video monitors, in optical systems, 5.16
- Video Quality Experts Group (VQEG), 24.4
- Vieth-Müller circle, 13.8, 13.9, 13.9f
- Vieth-Müller horopter, 2.40
- Viewing environments:
 for color CRTs, 22.19–22.20
 for computer work, 23.4–23.9
 lighting, 23.4–23.5, 23.5t
 monitor characteristics, 23.6–23.8
 screen reflections, 23.5–23.6
 work habits, 23.8–23.9
 workstation arrangement, 23.8
- Virtual reality:
 head-mounted displays for, 25.1–25.12
 characterizing, 25.7–25.10, 25.8f, 25.9t–25.10t
 design considerations, 25.2–25.7
 research in development of, 24.9
- Visible light, hazards from, 7.10
- Vision, defined, 23.4
- Vision experiments/research:
 displays for, 22.1–22.40
 color cathode ray tubes (color CRTs), 22.1–22.34
 liquid crystal displays (LCDs), 22.34–22.40
 using Maxwellian view in (*see* Maxwellian view)
 (*See also* Psychophysical measurement)
- Vision laboratories, control of light sources in (*see* Optical generation of visual stimulus)
- Vision therapy, 23.4
- Visual acuity, 2.33–2.35, 2.34f, 2.35f, 4.1–4.13, 12.6, 12.7
 age-related changes in, 14.17, 14.19
 and color, 4.10
 and contrast, 4.12, 4.12f
 defined, 4.1, 12.2, 15.2, 23.1
 and defocus, 4.9, 4.9f, 4.10f
 factors affecting, 4.9–4.13
 hyperacuity, 4.13–4.16, 4.14f
 defined, 4.1
 and superresolution, 4.15
 and luminance, 4.11, 4.11f, 4.12
 practice effect on, 4.12
 and pupil, 4.9
 resolution
 and information, 4.15–4.16
 and superresolution, 4.15
- Visual acuity (*Cont.*):
 resolving capacity of the eye, 4.4–4.5, 4.4f
 contrast-transfer function, 4.5
 point-spread function, 4.4, 4.5
 spatial-frequency coordinates, 4.5
 and retinal eccentricity, 4.10–4.11, 4.11f
 retinal limitations, 4.5–4.6, 4.6f
 and stage of development/aging, 4.13
 stimulus specification, 4.2–4.4, 4.3f
 and surround, 4.12, 4.13f
 tests of visual acuity, 4.7–4.8, 4.7f
 and time, 4.12, 4.13f
 visual resolution threshold determination, 4.6–4.7
- Visual angles:
 angular extent, 5.2
 defined, 2.3, 4.1, 10.1, 10.4t
 degrees of, 10.9
 and image formation, 2.3
 steradian of, 6.3
- Visual axis, 2.3
- Visual differences predictor, 24.3
- Visual displays, measuring performance of (*see* Psychophysical measurement)
- Visual fields:
 age-related changes in, 14.21, 14.21f
 binocular overlap of, 13.7
 defined, 13.2
 perception of, 13.3
- Visual instruments:
 and accommodation response in eye, 1.34–1.35
 and chromostereopsis, 1.20
 exit pupil of, 1.27–1.28
 retinal image quality with, 1.27–1.28
 stereoscopic, 1.41
- Visual performance, 2.1–2.41
 aging-related changes in, 14.14–14.22
 color vision, 14.15, 14.17
 depth and stereovision, 14.22
 minimal, 14.22
 sensitivity under scotopic and photopic conditions, 14.15, 14.16f
 spatial vision, 14.17–14.19, 14.18f
 temporal vision, 14.19–14.21, 14.20f
 visual field, 14.21, 14.21f
 binocular stereoscopic discrimination, 2.40–2.41, 2.41f

- Visual performance (*Cont.*):
- central visual processing, 2.12–2.14
 - anatomy of LGN, 2.12–2.13, 2.13f
 - physiology of LGN, 2.13–2.14
 - contrast detection, 2.19–2.31
 - adaptation and inhibition, 2.25–2.27
 - chromatic, 2.29–2.31, 2.30f
 - eye movements, 2.21
 - as function of spatial frequency, 2.20f
 - optical transfer function, 2.21–2.22, 2.22f
 - optical/retinal inhomogeneity, 2.24, 2.25f
 - receptors, 2.22–2.23
 - spatial channels, 2.23–2.24
 - temporal, 2.27–2.29, 2.27f, 2.29f
 - contrast discrimination, 2.31–2.32, 2.31f
 - contrast estimation, 2.33
 - contrast masking, 2.32–2.33, 2.32f
 - ideal-observer theory, 2.16–2.19, 2.19f
 - image formation, 2.2–2.4
 - image sampling by photoreceptors, 2.4–2.9
 - information-processing model for, 2.15–2.16, 2.15f
 - motion detection/discrimination, 2.36–2.40
 - optic flow fields, 2.39f
 - thresholds, 2.37f
 - pattern discrimination, 2.35–2.36
 - retinal processing, 2.9–2.12
 - anatomy of retina, 2.9–2.11
 - physiology of retina, 2.11–2.12
 - visual acuity, 2.33–2.35, 2.34f, 2.35f
- Visual plane, 13.2
- Visual resolution threshold, determination of, 4.6–4.7
- Visual stimuli:
- in decision tasks, 3.2
 - optical generation of (*see* Optical generation of visual stimulus)
 - in psychophysical experiments, 3.3–3.4
 - spatial frequency content of, 5.8
- Visual stress, 23.4
- Visual systems (colorimetry), 10.38–10.40
- Visual tasks, 2.15
- Visualization, in electronic imaging, 24.8
- Visual-vestibular conflicts, in head mounted display systems, 13.32–13.33
- Vitreous humor/gel, 1.3f, 14.9, 16.3, 18.2
- Von Bezold spreading, 11.3, 11.4f
- Von Kries adaptation, 11.3, 11.49, 11.51, 11.68
- Watercolor illusion/effect, 11.72, 11.73f
- Wave aberration, 15.2
- Wavefront aberration, 1.15–1.19, 1.16f, 1.16t, 1.17f
 - and defocus of visual instruments, 1.28
 - and final retinal image quality, 1.21–1.22
- Wavefront correctors (AO systems for the eye), 15.3, 15.7f, 15.9–15.12, 15.10f, 15.11f
- Wavefront sensor (AO systems for the eye), 15.3, 15.7f, 15.8–15.9, 15.8f
- Waveguide parameter (*V*-parameter), 8.2, 8.17, 8.20–8.21
- Waveguides, biological (*see* Biological waveguides)
- Wavelength:
- color data as functions of, 10.21, 10.23, 10.23f
 - in optical systems, 5.11–5.13, 5.12t
- Weber contrast, in vision experiments, 3.4
- Weber's law, 11.16f
 - and cone contrast spaces, 11.32
 - and contrast coding, 11.15–11.16, 11.16f
 - defined, 11.3
 - first-site adaptation, 11.15, 11.16f
 - second-site adaptation, 11.18
- Webvision, 11.5
- Welder's flash, 1.9, 7.4
- Wet age-related macular degeneration, 14.1, 14.24
- White balance, in color CRTs, 22.9
- White noise, 18.2
- Wide-angle models (of human eye), 1.38
- With-the-rule astigmatism, 1.6
- Wood lens, 19.3–19.4, 19.4f
- Working distance (retinoscopy), 12.2, 12.6
- Workstation arrangement, for computer work, 23.8
- World Health Organization (WHO), 7.9
- Wyszecki lens, 9.5
- Yes-no judgments, 3.6
- Yoked eye movements, 13.2, 13.20
- Young-Helmholtz trichromatic theory, 11.5
- Young's fringes, 1.22

- Z-dependent oscillary term, 8.22, 8.23
- Zernike coefficient, for RMS wavefront error,
1.15–1.19, 1.16*t*, 1.17*f*
- Zernike polynomials:
defined, 15.2
for representing ocular aberrations, 15.4,
15.5*f*, 16.6–16.7, 16.7*t*
- ZEST, 3.9
- Zonal refractive lenses (IOLs), 21.15
- Zone models (color vision), 11.5
three-stage, 11.6, 11.7*f*, 11.8, 11.82–11.85,
11.85*f*
two-stage, 11.5, 11.6*f*
- Zonules, 1.3*f*, 21.2

DO NOT DUPLICATE

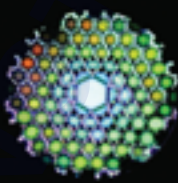
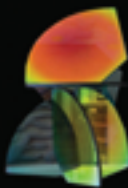
Third Edition

Sponsored by the Optical Society of America

HANDBOOK OF OPTICS

Volume IV

*Optical Properties of Materials, Nonlinear Optics,
Quantum Optics*



Editor-in-Chief:
Michael Bass

Associate Editors:
Casimer M. DeCusatis
Jay M. Enoch
Vasudevan Lakshminarayanan
Guifang Li
Carolyn MacDonald
Virendra N. Mahajan
Eric Van Stryland

OSA[®]

HANDBOOK OF OPTICS

DO NOT DUPLICATE

ABOUT THE EDITORS

Editor-in-Chief: Dr. Michael Bass is professor emeritus at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Associate Editors:

Dr. Casimer M. DeCusatis is a distinguished engineer and technical executive with IBM Corporation, Poughkeepsie, New York.

Dr. Jay M. Enoch is dean emeritus and professor at the School of Optometry at the University of California, Berkeley.

Dr. Vasudevan Lakshminarayanan is professor of Optometry, Physics, and Electrical Engineering at the University of Waterloo, Ontario, Canada.

Dr. Guifang Li is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Dr. Carolyn MacDonald is professor and chair of physics at the University at Albany, SUNY, and the director of the Albany Center for X-Ray Optics.

Dr. Virendra N. Mahajan is a distinguished scientist at The Aerospace Corporation.

Dr. Eric Van Stryland is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

HANDBOOK OF OPTICS

Volume IV
Optical Properties of Materials,
Nonlinear Optics, Quantum Optics

THIRD EDITION

Sponsored by the
OPTICAL SOCIETY OF AMERICA

Michael Bass Editor-in-Chief
*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

Guifang Li Associate Editor
*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

Eric Van Stryland Associate Editor
*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*



New York Chicago San Francisco Lisbon London Madrid
Mexico City Milan New Delhi San Juan Seoul
Singapore Sydney Toronto

This page intentionally left blank.

DO NOT DUPLICATE

COVER ILLUSTRATIONS

Left: Photograph of a femtosecond optical parametric oscillator pumped in the blue by the second harmonic of a Ti:sapphire laser and operating in the orange. The oscillator can deliver femtosecond pulses across the entire visible range from the blue-green to yellow-red by simple rotation of the nonlinear crystal. (*Courtesy of Radiantis, S. L., Barcelona, Spain.*) See Chapter 19.

Middle: Photograph of a thin-film-based sculpture showing the beautiful colors of thin films seen in reflection and transmission. The variety of properties one may achieve with optical thin films is demonstrated in this photo by the different colors of reflected and transmitted light seen as a result of different coating design and angle of incidence. See Chapter 7.

Right: This is an optical micrograph of the end face of a hollow core photonic crystal fiber with super continuum white light launched at the far end. It shows the separation of colors according to the lifetimes of Mie resonances in the hollow channels. This illustrates nonlinear optical phenomena as discussed in several chapters of this volume, but also uses fibers as discussed in Chapter 11 of Vol. V.

This page intentionally left blank.

DO NOT DUPLICATE

CONTENTS

Contributors	xiii
Brief Contents of All Volumes	xv
Editors' Preface	xxi
Preface to Volume IV	xxiii
Glossary and Fundamental Constants	xxv

Part 1. Properties

Chapter 1. Optical Properties of Water	<i>Curtis D. Mobley</i>	1.3
<hr/>		
1.1	Introduction / 1.3	
1.2	Terminology, Notation, and Definitions / 1.3	
1.3	Radiometric Quantities Useful in Hydrologic Optics / 1.4	
1.4	Inherent Optical Properties / 1.9	
1.5	Apparent Optical Properties / 1.12	
1.6	The Optically Significant Constituents of Natural Waters / 1.13	
1.7	Particle Size Distributions / 1.15	
1.8	Electromagnetic Properties of Water / 1.16	
1.9	Index of Refraction / 1.18	
1.10	Measurement of Absorption / 1.20	
1.11	Absorption by Pure Sea Water / 1.21	
1.12	Absorption by Dissolved Organic Matter / 1.22	
1.13	Absorption by Phytoplankton / 1.23	
1.14	Absorption by Organic Detritus / 1.25	
1.15	Bio-Optical Models for Absorption / 1.27	
1.16	Measurement of Scattering / 1.29	
1.17	Scattering by Pure Water and by Pure Sea Water / 1.30	
1.18	Scattering by Particles / 1.30	
1.19	Wavelength Dependence of Scattering: Bio-Optical Models / 1.35	
1.20	Beam Attenuation / 1.40	
1.21	Diffuse Attenuation and Jerlov Water Types / 1.42	
1.22	Irradiance Reflectance and Remote Sensing / 1.46	
1.23	Inelastic Scattering and Polarization / 1.47	
1.24	Acknowledgments / 1.50	
1.25	References / 1.50	
Chapter 2. Properties of Crystals and Glasses	<i>William J. Tropf, Michael E. Thomas, and Eric W. Rogala</i>	2.1
<hr/>		
2.1	Glossary / 2.1	
2.2	Introduction / 2.3	
2.3	Optical Materials / 2.4	
2.4	Properties of Materials / 2.5	
2.5	Properties Tables / 2.36	
2.6	References / 2.77	
Chapter 3. Polymeric Optics	<i>John D. Lytle</i>	3.1
<hr/>		
3.1	Glossary / 3.1	
3.2	Introduction / 3.1	

- 3.3 Forms / 3.2
- 3.4 Physical Properties / 3.2
- 3.5 Optical Properties / 3.5
- 3.6 Optical Design / 3.7
- 3.7 Processing / 3.11
- 3.8 Coatings / 3.17
- 3.9 References / 3.18

Chapter 4. Properties of Metals *Roger A. Paquin* 4.1

- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.2
- 4.3 Summary Data / 4.11
- 4.4 References / 4.70

Chapter 5. Optical Properties of Semiconductors
*David G. Seiler, Stefan Zollner, Alain C. Diebold,
and Paul M. Amirtharaj* 5.1

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.3
- 5.3 Optical Properties / 5.8
- 5.4 Measurement Techniques / 5.56
- 5.5 Acknowledgments / 5.83
- 5.6 Summary and Conclusions / 5.83
- 5.7 References / 5.91

Chapter 6. Characterization and Use of Black Surfaces for Optical Systems *Stephen M. Pompea and Robert P. Breault* 6.1

- 6.1 Introduction / 6.1
- 6.2 Selection Process for Black Baffle Surfaces in Optical Systems / 6.10
- 6.3 The Creation of Black Surfaces for Specific Applications / 6.13
- 6.4 Environmental Degradation of Black Surfaces / 6.16
- 6.5 Optical Characterization of Black Surfaces / 6.18
- 6.6 Surfaces for Ultraviolet and Far-Infrared Applications / 6.21
- 6.7 Survey of Surfaces with Optical Data / 6.34
- 6.8 Paints / 6.35
- 6.9 Conclusions / 6.59
- 6.10 Acknowledgments / 6.59
- 6.11 References / 6.60
- 6.12 Further Readings / 6.67

Chapter 7. Optical Properties of Films and Coatings
Jerzy A. Dobrowolski 7.1

- 7.1 Introduction / 7.1
- 7.2 Theory and Design of Optical Thin-Film Coatings / 7.5
- 7.3 Thin-Film Manufacturing Considerations / 7.10
- 7.4 Measurements on Optical Coatings / 7.12
- 7.5 Antireflection Coatings / 7.15
- 7.6 Two-Material Periodic Multilayers Theory / 7.32
- 7.7 Multilayer Reflectors—Experimental Results / 7.39
- 7.8 Cutoff, Heat-Control, and Solar-Cell Cover Filters / 7.53
- 7.9 Beam Splitters and Neutral Filters / 7.61
- 7.10 Interference Polarizers and Polarizing Beam Splitters / 7.69

- 7.11 Bandpass Filters / 7.73
- 7.12 High Performance Optical Multilayer Coatings / 7.96
- 7.13 Multilayers for Two or Three Spectral Regions / 7.98
- 7.14 Phase Coatings / 7.101
- 7.15 Interference Filters with Low Reflection / 7.104
- 7.16 Reflection Filters and Coatings / 7.106
- 7.17 Special Purpose Coatings / 7.113
- 7.18 References / 7.114

Chapter 8. Fundamental Optical Properties of Solids *Alan Miller* 8.1

- 8.1 Glossary / 8.1
- 8.2 Introduction / 8.3
- 8.3 Propagation of Light in Solids / 8.4
- 8.4 Dispersion Relations / 8.14
- 8.5 Lattice Interactions / 8.16
- 8.6 Free Electron Properties / 8.21
- 8.7 Band Structures and Interband Transitions / 8.24
- 8.8 References / 8.32

Chapter 9. Photonic Bandgap Materials *Pierre R. Villeneuve* 9.1

- 9.1 Glossary / 9.1
- 9.2 Introduction / 9.2
- 9.3 Maxwell's Equations / 9.2
- 9.4 Three-Dimensional Photonic Crystals / 9.4
- 9.5 Microcavities in Three-Dimensional Photonic Crystals / 9.6
- 9.6 Microcavities in Photonic Crystals with Two-Dimensional Periodicity / 9.8
- 9.7 Waveguides / 9.12
- 9.8 Conclusion / 9.17
- 9.9 References / 9.18

Part 2. Nonlinear Optics

Chapter 10. Nonlinear Optics *Chung L. Tang* 10.3

- 10.1 Glossary / 10.3
- 10.2 Introduction / 10.4
- 10.3 Basic Concepts / 10.5
- 10.4 Material Considerations / 10.19
- 10.5 Appendix / 10.21
- 10.6 References / 10.23

Chapter 11. Coherent Optical Transients *Paul R. Berman and Duncan G. Steel* 11.1

- 11.1 Glossary / 11.1
- 11.2 Introduction / 11.2
- 11.3 Optical Bloch Equations / 11.3
- 11.4 Maxwell-Bloch Equations / 11.6
- 11.5 Free Polarization Decay / 11.7
- 11.6 Photon Echo / 11.11
- 11.7 Stimulated Photon Echo / 11.15
- 11.8 Phase Conjugate Geometry and Optical Ramsey Fringes / 11.19
- 11.9 Two-Photon Transitions and Atom Interferometry / 11.22
- 11.10 Chirped Pulse Excitation / 11.25
- 11.11 Experimental Considerations / 11.26
- 11.12 Conclusion / 11.28
- 11.13 References / 11.28

Chapter 12. Photorefractive Materials and Devices <i>Mark Cronin-Golomb and Marvin Klein</i>	12.1
<hr/>	
12.1 Introduction / 12.1	
12.2 Materials / 12.10	
12.3 Devices / 12.28	
12.4 References / 12.38	
12.5 Further Reading / 12.45	
Chapter 13. Optical Limiting <i>David J. Hagan</i>	13.1
<hr/>	
13.1 Introduction / 13.1	
13.2 Basic Principles of Passive Optical Limiting / 13.4	
13.3 Examples of Passive Optical Limiting in Specific Materials / 13.9	
13.4 References / 13.13	
Chapter 14. Electromagnetically Induced Transparency <i>Jonathan P. Marangos and Thomas Halfmann</i>	14.1
<hr/>	
14.1 Glossary / 14.1	
14.2 Introduction / 14.2	
14.3 Coherence in Two- and Three-Level Atomic Systems / 14.4	
14.4 The Basic Physical Concept of Electromagnetically Induced Transparency / 14.5	
14.5 Manipulation of Optical Properties by Electromagnetically Induced Transparency / 14.10	
14.6 Electromagnetically Induced Transparency, Driven by Pulsed Lasers / 14.15	
14.7 Steady State Electromagnetically Induced Transparency, Driven by CW Lasers / 14.16	
14.8 Gain without Inversion and Lasing without Inversion / 14.18	
14.9 Manipulation of the Index of Refraction in Dressed Atoms / 14.19	
14.10 Pulse Propagation Effects / 14.20	
14.11 Ultraslow Light Pulses / 14.22	
14.12 Nonlinear Optical Frequency Conversion / 14.24	
14.13 Nonlinear Optics at Maximal Atomic Coherence / 14.28	
14.14 Nonlinear Optics at the Few Photon Level / 14.32	
14.15 Electromagnetically Induced Transparency in Solids / 14.33	
14.16 Conclusion / 14.36	
14.17 Further Reading / 14.36	
14.18 References / 14.37	
Chapter 15. Stimulated Raman and Brillouin Scattering <i>John Reintjes and Mark Bashkansky</i>	15.1
<hr/>	
15.1 Introduction / 15.1	
15.2 Raman Scattering / 15.1	
15.3 Stimulated Brillouin Scattering / 15.43	
15.4 References / 15.54	
15.5 Additional References / 15.60	
Chapter 16. Third-Order Optical Nonlinearities <i>Mansoor Sheik-Bahae and Michael P. Hasselbeck</i>	16.1
<hr/>	
16.1 Introduction / 16.1	
16.2 Quantum Mechanical Picture / 16.4	
16.3 Nonlinear Absorption and Nonlinear Refraction / 16.7	
16.4 Kramers-Kronig Dispersion Relations / 16.9	
16.5 Optical Kerr Effect / 16.11	
16.6 Third-Harmonic Generation / 16.14	
16.7 Stimulated Scattering / 16.14	
16.8 Two-Photon Absorption / 16.19	
16.9 Effective Third-Order Nonlinearities; Cascaded $\chi^{(1)}\chi^{(2)}$ Processes / 16.20	
16.10 Effective Third-Order Nonlinearities; Cascaded $\chi^{(2)}\chi^{(2)}$ Processes / 16.22	

- 16.11 Propagation Effects / 16.24
 16.12 Common Experimental Techniques and Applications / 16.26
 16.13 References / 16.31

Chapter 17. Continuous-Wave Optical Parametric Oscillators 17.1
Majid Ebrahim-Zadeh

- 17.1 Introduction / 17.1
 17.2 Continuous-Wave Optical Parametric Oscillators / 17.2
 17.3 Applications / 17.21
 17.4 Summary / 17.29
 17.5 References / 17.31

Chapter 18. Nonlinear Optical Processes for Ultrashort Pulse Generation 18.1
Uwe Siegner and Ursula Keller

- 18.1 Glossary / 18.1
 18.2 Abbreviations / 18.3
 18.3 Introduction / 18.3
 18.4 Saturable Absorbers: Macroscopic Description / 18.5
 18.5 Kerr Effect / 18.11
 18.6 Semiconductor Ultrafast Nonlinearities: Microscopic Processes / 18.15
 18.7 References / 18.23

Chapter 19. Laser-Induced Damage to Optical Materials 19.1
Marion J. Soileau

- 19.1 Introduction / 19.1
 19.2 Practical Estimates / 19.2
 19.3 Surface Damage / 19.2
 19.4 Package-Induced Damage / 19.4
 19.5 Nonlinear Optical Effects / 19.5
 19.6 Avoidance of Damage / 19.5
 19.7 Fundamental Mechanisms / 19.6
 19.8 Progress in Measurements of Critical NLO Parameters / 19.9
 19.9 References / 19.11

Part 3. Quantum and Molecular Optics

Chapter 20. Laser Cooling and Trapping of Atoms 20.3
Harold J. Metcalf and Peter van der Straten

- 20.1 Introduction / 20.3
 20.2 General Properties Concerning Laser Cooling / 20.4
 20.3 Theoretical Description / 20.6
 20.4 Slowing Atomic Beams / 20.11
 20.5 Optical Molasses / 20.13
 20.6 Cooling Below the Doppler Limit / 20.17
 20.7 Trapping of Neutral Atoms / 20.21
 20.8 Applications / 20.26
 20.9 References / 20.39

Chapter 21. Strong Field Physics 21.1
Todd Ditmire

- 21.1 Glossary / 21.1
 21.2 Introduction and History / 21.2
 21.3 Laser Technology Used in Strong Field Physics / 21.4
 21.4 Strong Field Interactions with Single Electrons / 21.5
 21.5 Strong Field Interactions with Atoms / 21.10

- 21.6 Strong Field Interactions with Molecules / 21.22
- 21.7 Strong Field Nonlinear Optics in Gases / 21.27
- 21.8 Strong Field Interactions with Clusters / 21.31
- 21.9 Strong Field Physics in Underdense Plasmas / 21.36
- 21.10 Strong Field Physics at Surfaces of Overdense Plasmas / 21.46
- 21.11 Applications of Strong Field Interactions with Plasmas / 21.52
- 21.12 References / 21.55

Chapter 22. Slow Light Propagation in Atomic and Photonic Media *Jacob B. Khurgin* 22.1

- 22.1 Glossary / 22.1
- 22.2 Introduction / 22.2
- 22.3 Atomic Resonance / 22.2
- 22.4 Bandwidth Limitations in Atomic Schemes / 22.9
- 22.5 Photonic Resonance / 22.9
- 22.6 Slow Light in Optical Fibers / 22.13
- 22.7 Conclusion / 22.15
- 22.8 References / 22.16

Chapter 23. Quantum Entanglement in Optical Interferometry
*Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver,
and Jonathan P. Dowling* 23.1

- 23.1 Introduction / 23.1
- 23.2 Shot-Noise Limit / 23.4
- 23.3 Heisenberg Limit / 23.6
- 23.4 “Digital” Approaches / 23.7
- 23.5 N00n State / 23.9
- 23.6 Quantum Imaging / 23.13
- 23.7 Toward Quantum Remote Sensing / 23.14
- 23.8 References / 23.15

Index I.1

CONTRIBUTORS

- Paul M. Amirtharaj** *Sensors and Electron Devices Directorate, U.S. Army Research Laboratory, Adelphi, Maryland* (CHAP. 5)
- Mark Bashkansky** *Optical Sciences Division, Naval Research Laboratory, Washington, D.C.* (CHAP. 15)
- Paul R. Berman** *Physics Department, University of Michigan, Ann Arbor, Michigan* (CHAP. 11)
- Robert P. Breault** *Breault Research Organization, Inc., Tucson, Arizona* (CHAP. 6)
- Mark Cronin-Golomb** *Department of Biomedical Engineering, Tufts University, Medford, Massachusetts* (CHAP. 12)
- Alain C. Diebold** *College of Nanoscale Science and Engineering, University at Albany, Albany, New York* (CHAP. 5)
- Todd Ditmire** *Texas Center for High Intensity Laser Science, Department of Physics, The University of Texas at Austin, Austin, Texas* (CHAP. 21)
- Jerzy A. Dobrowolski** *Institute for Microstructural Sciences, National Research Council of Canada, Ottawa, Ontario, Canada* (CHAP. 7)
- Jonathan P. Dowling** *Hearne Institute for Theoretical Physics, Department of Physics and Astronomy, Louisiana State University, Baton Rouge, Louisiana* (CHAP. 23)
- Majid Ebrahim-Zadeh** *ICFO—Institut de Ciències Fotoniques, Mediterranean Technology Park, Barcelona, Spain, and Institutio Catalana de Recerca i Estudis Avancats (ICREA), Passeig Lluís Companys, Barcelona, Spain* (CHAP. 17)
- David J. Hagan** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 13)
- Thomas Halfmann** *Institute of Applied Physics, Technical University of Darmstadt, Darmstadt, Germany* (CHAP. 14)
- Michael P. Hasselbeck** *Department of Physics and Astronomy, University of New Mexico, Albuquerque, New Mexico* (CHAP. 16)
- Sean D. Huver** *Hearne Institute for Theoretical Physics, Department of Physics and Astronomy, Louisiana State University, Baton Rouge, Louisiana* (CHAP. 23)
- Ursula Keller** *Institute of Quantum Electronics, Physics Department, Swiss Federal Institute of Technology (ETH), Zurich, Switzerland* (CHAP. 18)
- Jacob B. Khurgin** *Department of Electrical and Computer Engineering, Johns Hopkins University, Baltimore, Maryland* (CHAP. 22)
- Marvin Klein** *Intelligent Optical Systems, Inc., Torrance, California* (CHAP. 12)
- Hwang Lee** *Hearne Institute for Theoretical Physics, Department of Physics and Astronomy, Louisiana State University, Baton Rouge, Louisiana* (CHAP. 23)
- John D. Lytle** *Advanced Optical Concepts, Santa Cruz, California* (CHAP. 3)
- Jonathan P. Marangos** *Quantum Optics and Laser Science Group, Blackett Laboratory, Imperial College, London, United Kingdom* (CHAP. 14)
- Harold J. Metcalf** *Department of Physics, State University of New York, Stony Brook, New York* (CHAP. 20)
- Alan Miller** *Scottish Universities Physics Alliance, School of Engineering and Physical Sciences, Heriot-Watt University, Edinburgh, Scotland* (CHAP. 8)
- Curtis D. Mobley** *Applied Electromagnetics and Optics Laboratory, SRI International, Menlo Park, California* (CHAP. 1)
- Roger A. Paquin** *Advanced Materials Consultant, Tucson, Arizona, and Optical Sciences Center, University of Arizona, Tucson* (CHAP. 4)
- Stephen M. Pompea** *National Optical Astronomy Observatory, Tucson, Arizona* (CHAP. 6)
- John Reintjes** *Optical Sciences Division, Naval Research Laboratory, Washington, D.C.* (CHAP. 15)

- Eric W. Rogala** *Raytheon Missile Systems, Tucson, Arizona* (CHAP. 2)
- David G. Seiler** *Semiconductor Electronics Division, National Institute of Standards and Technology, Gaithersburg, Maryland* (CHAP. 5)
- Mansoor Sheik-Bahae** *Department of Physics and Astronomy, University of New Mexico, Albuquerque, New Mexico* (CHAP. 16)
- Uwe Siegner** *Institute of Quantum Electronics, Physics Department, Swiss Federal Institute of Technology (ETH), Zurich, Switzerland* (CHAP. 18)
- Marion J. Soileau** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 19)
- Duncan G. Steel** *Physics Department, University of Michigan, Ann Arbor, Michigan* (CHAP. 11)
- Chung L. Tang** *School of Electrical and Computer Engineering, Cornell University, Ithaca, New York* (CHAP. 10)
- Michael E. Thomas** *Applied Physics Laboratory, Johns Hopkins University, Laurel, Maryland* (CHAP. 2)
- William J. Tropf** *Applied Physics Laboratory, Johns Hopkins University, Laurel, Maryland* (CHAP. 2)
- Peter van der Straten** *Debye Institute, Department of Atomic and Interface Physics, Utrecht University, Utrecht, The Netherlands* (CHAP. 20)
- Pierre R. Villeneuve** *Department of Physics, Massachusetts Institute of Technology, Cambridge, Massachusetts* (CHAP. 9)
- Christoph F. Wildfeuer** *Hearne Institute for Theoretical Physics, Department of Physics and Astronomy, Louisiana State University, Baton Rouge, Louisiana* (CHAP. 23)
- Stefan Zollner** *Freescale Semiconductor, Inc., Hopewell Junction, New York* (CHAP. 5)

BRIEF CONTENTS OF ALL VOLUMES

VOLUME I. GEOMETRICAL AND PHYSICAL OPTICS, POLARIZED LIGHT, COMPONENTS AND INSTRUMENTS

PART 1. GEOMETRICAL OPTICS

Chapter 1. General Principles of Geometrical Optics *Douglas S. Goodman*

PART 2. PHYSICAL OPTICS

Chapter 2. Interference *John E. Greivenkamp*

Chapter 3. Diffraction *Arvind S. Marathay and John F. McCalmont*

Chapter 4. Transfer Function Techniques *Glenn D. Boreman*

Chapter 5. Coherence Theory *William H. Carter*

Chapter 6. Coherence Theory: Tools and Applications *Gisele Bennett, William T. Rhodes, and J. Christopher James*

Chapter 7. Scattering by Particles *Craig F. Bohren*

Chapter 8. Surface Scattering *Eugene L. Church and Peter Z. Takacs*

Chapter 9. Volume Scattering in Random Media *Aristide Dogariu and Jeremy Ellis*

Chapter 10. Optical Spectroscopy and Spectroscopic Lineshapes *Brian Henderson*

Chapter 11. Analog Optical Signal and Image Processing *Joseph W. Goodman*

PART 3. POLARIZED LIGHT

Chapter 12. Polarization *Jean M. Bennett*

Chapter 13. Polarizers *Jean M. Bennett*

Chapter 14. Mueller Matrices *Russell A. Chipman*

Chapter 15. Polarimetry *Russell A. Chipman*

Chapter 16. Ellipsometry *Rasheed M. A. Azzam*

PART 4. COMPONENTS

Chapter 17. Lenses *R. Barry Johnson*

Chapter 18. Afocal Systems *William B. Wetherell*

Chapter 19. Nondispersive Prisms *William L. Wolfe*

Chapter 20. Dispersive Prisms and Gratings *George J. Zissis*

Chapter 21. Integrated Optics *Thomas L. Koch, Frederick J. Leonberger, and Paul G. Suchoski*

Chapter 22. Miniature and Micro-Optics *Tom D. Milster and Tomasz S. Tkaczyk*

Chapter 23. Binary Optics *Michael W. Farn and Wilfrid B. Veldkamp*

Chapter 24. Gradient Index Optics *Duncan T. Moore*

PART 5. INSTRUMENTS

Chapter 25. Cameras *Norman Goldberg*

Chapter 26. Solid-State Cameras *Gerald C. Holst*

Chapter 27. Camera Lenses *Ellis Betensky, Melvin H. Kreitzer, and Jacob Moskovich*

Chapter 28. Microscopes *Rudolf Oldenbourg and Michael Shribak*

Chapter 29. Reflective and Catadioptric Objectives *Lloyd Jones*

- Chapter 30. Scanners *Leo Beiser and R. Barry Johnson*
- Chapter 31. Optical Spectrometers *Brian Henderson*
- Chapter 32. Interferometers *Parameswaran Hariharan*
- Chapter 33. Holography and Holographic Instruments *Lloyd Huff*
- Chapter 34. Xerographic Systems *Howard Stark*
- Chapter 35. Principles of Optical Disk Data Storage *Masud Mansuripur*

VOLUME II. DESIGN, FABRICATION, AND TESTING; SOURCES AND DETECTORS; RADIOMETRY AND PHOTOMETRY

PART 1. DESIGN

- Chapter 1. Techniques of First-Order Layout *Warren J. Smith*
- Chapter 2. Aberration Curves in Lens Design *Donald C. O'Shea and Michael E. Harrigan*
- Chapter 3. Optical Design Software *Douglas C. Sinclair*
- Chapter 4. Optical Specifications *Robert R. Shannon*
- Chapter 5. Tolerancing Techniques *Robert R. Shannon*
- Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.*
- Chapter 7. Control of Stray Light *Robert P. Breault*
- Chapter 8. Thermal Compensation Techniques *Philip J. Rogers and Michael Roberts*

PART 2. FABRICATION

- Chapter 9. Optical Fabrication *Michael P. Mandina*
- Chapter 10. Fabrication of Optics by Diamond Turning *Richard L. Rhorer and Chris J. Evans*

PART 3. TESTING

- Chapter 11. Orthonormal Polynomials in Wavefront Analysis *Virendra N. Mahajan*
- Chapter 12. Optical Metrology *Zacarias Malacara and Daniel Malacara-Hernández*
- Chapter 13. Optical Testing *Daniel Malacara-Hernández*
- Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant*

PART 4. SOURCES

- Chapter 15. Artificial Sources *Anthony LaRocca*
- Chapter 16. Lasers *William T. Silfvast*
- Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman*
- Chapter 18. High-Brightness Visible LEDs *Winston V. Schoenfeld*
- Chapter 19. Semiconductor Lasers *Pamela L. Derry, Luis Figueroa, and Chi-shain Hong*
- Chapter 20. Ultrashort Optical Sources and Applications *Jean-Claude Diels and Ladan Arissian*
- Chapter 21. Attosecond Optics *Zenghu Chang*
- Chapter 22. Laser Stabilization *John L. Hall, Matthew S. Taubman, and Jun Ye*
- Chapter 23. Quantum Theory of the Laser *János A. Bergou, Berthold-Georg Englert, Melvin Lax, Marian O. Scully, Herbert Walther, and M. Suhail Zubairy*

PART 5. DETECTORS

- Chapter 24. Photodetectors *Paul R. Norton*
- Chapter 25. Photodetection *Abhay M. Joshi and Gregory H. Olsen*
- Chapter 26. High-Speed Photodetectors *John E. Bowers and Yih G. Wey*
- Chapter 27. Signal Detection and Analysis *John R. Willison*
- Chapter 28. Thermal Detectors *William L. Wolfe and Paul W. Kruse*

PART 6. IMAGING DETECTORS

- Chapter 29. Photographic Films *Joseph H. Altman*
- Chapter 30. Photographic Materials *John D. Baloga*

- Chapter 31. Image Tube Intensified Electronic Imaging *C. Bruce Johnson and Larry D. Owen*
 Chapter 32. Visible Array Detectors *Timothy J. Tredwell*
 Chapter 33. Infrared Detector Arrays *Lester J. Kozlowski and Walter F. Kosonocky*

PART 7. RADIOMETRY AND PHOTOMETRY

- Chapter 34. Radiometry and Photometry *Edward F. Zalewski*
 Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection *James M. Palmer*
 Chapter 36. Radiometry and Photometry: Units and Conversions *James M. Palmer*
 Chapter 37. Radiometry and Photometry for Vision Optics *Yoshi Ohno*
 Chapter 38. Spectroradiometry *Carolyn J. Sher DeCusatis*
 Chapter 39. Nonimaging Optics: Concentration and Illumination *William Cassarly*
 Chapter 40. Lighting and Applications *Anurag Gupta and R. John Koshel*

VOLUME III. VISION AND VISION OPTICS

- Chapter 1. Optics of the Eye *Neil Charman*
 Chapter 2. Visual Performance *Wilson S. Geisler and Martin S. Banks*
 Chapter 3. Psychophysical Methods *Denis G. Pelli and Bart Farell*
 Chapter 4. Visual Acuity and Hyperacuity *Gerald Westheimer*
 Chapter 5. Optical Generation of the Visual Stimulus *Stephen A. Burns and Robert H. Webb*
 Chapter 6. The Maxwellian View with an Addendum on Apodization *Gerald Westheimer*
 Chapter 7. Ocular Radiation Hazards *David H. Sliney*
 Chapter 8. Biological Waveguides *Vasudevan Lakshminarayanan and Jay M. Enoch*
 Chapter 9. The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution *Jay M. Enoch and Vasudevan Lakshminarayanan*
 Chapter 10. Colorimetry *David H. Brainard and Andrew Stockman*
 Chapter 11. Color Vision Mechanisms *Andrew Stockman and David H. Brainard*
 Chapter 12. Assessment of Refraction and Refractive Errors and Their Influence on Optical Design *B. Ralph Chou*
 Chapter 13. Binocular Vision Factors That Influence Optical Design *Clifton Schor*
 Chapter 14. Optics and Vision of the Aging Eye *John S. Werner, Brooke E. Scheffrin, and Arthur Bradley*
 Chapter 15. Adaptive Optics in Retinal Microscopy and Vision *Donald T. Miller and Austin Roorda*
 Chapter 16. Refractive Surgery, Correction of Vision, PRK, and LASIK *L. Diaz-Santana and Harilaos Ginis*
 Chapter 17. Three-Dimensional Confocal Microscopy of the Living Human Cornea *Barry R. Masters*
 Chapter 18. Diagnostic Use of Optical Coherence Tomography in the Eye *Johannes F. de Boer*
 Chapter 19. Gradient Index Optics in the Eye *Barbara K. Pierscionek*
 Chapter 20. Optics of Contact Lenses *Edward S. Bennett*
 Chapter 21. Intraocular Lenses *Jim Schwiegerling*
 Chapter 22. Displays for Vision Research *William Cowan*
 Chapter 23. Vision Problems at Computers *Jeffrey Anshel and James E. Sheedy*
 Chapter 24. Human Vision and Electronic Imaging *Bernice E. Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allebach*
 Chapter 25. Visual Factors Associated with Head-Mounted Displays *Brian H. Tsou and Martin Shenker*

VOLUME IV. OPTICAL PROPERTIES OF MATERIALS, NONLINEAR OPTICS, QUANTUM OPTICS

PART 1. PROPERTIES

- Chapter 1. Optical Properties of Water *Curtis D. Mobley*
 Chapter 2. Properties of Crystals and Glasses *William J. Tropf, Michael E. Thomas, and Eric W. Rogala*
 Chapter 3. Polymeric Optics *John D. Lytle*
 Chapter 4. Properties of Metals *Roger A. Paquin*

- Chapter 5. Optical Properties of Semiconductors *David G. Seiler, Stefan Zollner, Alain C. Diebold, and Paul M. Amirtharaj*
- Chapter 6. Characterization and Use of Black Surfaces for Optical Systems *Stephen M. Pompea and Robert P. Breault*
- Chapter 7. Optical Properties of Films and Coatings *Jerzy A. Dobrowolski*
- Chapter 8. Fundamental Optical Properties of Solids *Alan Miller*
- Chapter 9. Photonic Bandgap Materials *Pierre R. Villeneuve*

PART 2. NONLINEAR OPTICS

- Chapter 10. Nonlinear Optics *Chung L. Tang*
- Chapter 11. Coherent Optical Transients *Paul R. Berman and Duncan G. Steel*
- Chapter 12. Photorefractive Materials and Devices *Mark Cronin-Golomb and Marvin Klein*
- Chapter 13. Optical Limiting *David J. Hagan*
- Chapter 14. Electromagnetically Induced Transparency *Jonathan P. Marangos and Thomas Halfmann*
- Chapter 15. Stimulated Raman and Brillouin Scattering *John Reintjes and Mark Bashkansky*
- Chapter 16. Third-Order Optical Nonlinearities *Mansoor Sheik-Bahae and Michael P. Hasselbeck*
- Chapter 17. Continuous-Wave Optical Parametric Oscillators *Majid Ebrahim-Zadeh*
- Chapter 18. Nonlinear Optical Processes for Ultrashort Pulse Generation *Uwe Siegner and Ursula Keller*
- Chapter 19. Laser-Induced Damage to Optical Materials *Marion J. Soileau*

PART 3. QUANTUM AND MOLECULAR OPTICS

- Chapter 20. Laser Cooling and Trapping of Atoms *Harold J. Metcalf and Peter van der Straten*
- Chapter 21. Strong Field Physics *Todd Ditmire*
- Chapter 22. Slow Light Propagation in Atomic and Photonic Media *Jacob B. Khurgin*
- Chapter 23. Quantum Entanglement in Optical Interferometry *Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver, and Jonathan P. Dowling*

VOLUME V. ATMOSPHERIC OPTICS, MODULATORS, FIBER OPTICS, X-RAY AND NEUTRON OPTICS

PART 1. MEASUREMENTS

- Chapter 1. Scatterometers *John C. Stover*
- Chapter 2. Spectroscopic Measurements *Brian Henderson*

PART 2. ATMOSPHERIC OPTICS

- Chapter 3. Atmospheric Optics *Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman*
- Chapter 4. Imaging through Atmospheric Turbulence *Virendra N. Mahajan and Guang-ming Dai*
- Chapter 5. Adaptive Optics *Robert Q. Fugate*

PART 3. MODULATORS

- Chapter 6. Acousto-Optic Devices *I-Cheng Chang*
- Chapter 7. Electro-Optic Modulators *Georgianne M. Purvinis and Theresa A. Maldonado*
- Chapter 8. Liquid Crystals *Sebastian Gauza and Shin-Tson Wu*

PART 4. FIBER OPTICS

- Chapter 9. Optical Fiber Communication Technology and System Overview *Ira Jacobs*
- Chapter 10. Nonlinear Effects in Optical Fibers *John A. Buck*
- Chapter 11. Photonic Crystal Fibers *Philip St. J. Russell and G. J. Pearce*
- Chapter 12. Infrared Fibers *James A. Harrington*
- Chapter 13. Sources, Modulators, and Detectors for Fiber Optic Communication Systems *Elsa Garmire*
- Chapter 14. Optical Fiber Amplifiers *John A. Buck*
- Chapter 15. Fiber Optic Communication Links (Telecom, Datacom, and Analog) *Casimer DeCusatis and Guifang Li*

- Chapter 16. Fiber-Based Couplers *Daniel Nolan*
 Chapter 17. Fiber Bragg Gratings *Kenneth O. Hill*
 Chapter 18. Micro-Optics-Based Components for Networking *Joseph C. Palais*
 Chapter 19. Semiconductor Optical Amplifiers *Jay M. Wiesenfeld and Leo H. Spiekman*
 Chapter 20. Optical Time-Division Multiplexed Communication Networks *Peter J. Delfyett*
 Chapter 21. WDM Fiber-Optic Communication Networks *Alan E. Willner, Changyuan Yu, Zhongqi Pan, and Yong Xie*
 Chapter 22. Solitons in Optical Fiber Communication Systems *Pavel V. Mamyshev*
 Chapter 23. Fiber-Optic Communication Standards *Casimer DeCusatis*
 Chapter 24. Optical Fiber Sensors *Richard O. Claus, Ignacio Matias, and Francisco Arregui*
 Chapter 25. High-Power Fiber Lasers and Amplifiers *Timothy S. McComb, Martin C. Richardson, and Michael Bass*

PART 5. X-RAY AND NEUTRON OPTICS

Subpart 5.1. Introduction and Applications

- Chapter 26. An Introduction to X-Ray and Neutron Optics *Carolyn MacDonald*
 Chapter 27. Coherent X-Ray Optics and Microscopy *Qun Shen*
 Chapter 28. Requirements for X-Ray Diffraction *Scott T. Misture*
 Chapter 29. Requirements for X-Ray Fluorescence *George J. Havrilla*
 Chapter 30. Requirements for X-Ray Spectroscopy *Dirk Lützenkirchen-Hecht and Ronald Frahm*
 Chapter 31. Requirements for Medical Imaging and X-Ray Inspection *Douglas Pfeiffer*
 Chapter 32. Requirements for Nuclear Medicine *Lars R. Furenlid*
 Chapter 33. Requirements for X-Ray Astronomy *Scott O. Rohrbach*
 Chapter 34. Extreme Ultraviolet Lithography *Franco Cerrina and Fan Jiang*
 Chapter 35. Ray Tracing of X-Ray Optical Systems *Franco Cerrina and Manuel Sanchez del Rio*
 Chapter 36. X-Ray Properties of Materials *Eric M. Gullikson*

Subpart 5.2. Refractive and Interference Optics

- Chapter 37. Refractive X-Ray Lenses *Bruno Lengeler and Christian G. Schroer*
 Chapter 38. Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region *Malcolm R. Howells*
 Chapter 39. Crystal Monochromators and Bent Crystals *Peter Siddons*
 Chapter 40. Zone Plates *Alan Michette*
 Chapter 41. Multilayers *Eberhard Spiller*
 Chapter 42. Nanofocusing of Hard X-Rays with Multilayer Laue Lenses *Albert T. Macrander, Hanfei Yan, Hyon Chol Kang, Jörg Maser, Chian Liu, Ray Conley, and G. Brian Stephenson*
 Chapter 43. Polarizing Crystal Optics *Qun Shen*

Subpart 5.3. Reflective Optics

- Chapter 44. Reflective Optics *James Harvey*
 Chapter 45. Aberrations for Grazing Incidence Optics *Timo T. Saha*
 Chapter 46. X-Ray Mirror Metrology *Peter Z. Takacs*
 Chapter 47. Astronomical X-Ray Optics *Marshall K. Joy and Brian D. Ramsey*
 Chapter 48. Multifoil X-Ray Optics *Ladislav Pina*
 Chapter 49. Pore Optics *Marco Beijersbergen*
 Chapter 50. Adaptive X-Ray Optics *Ali Khounsary*
 Chapter 51. The Schwarzschild Objective *Franco Cerrina*
 Chapter 52. Single Capillaries *Donald H. Bilderback and Sterling W. Cornaby*
 Chapter 53. Polycapillary X-Ray Optics *Carolyn MacDonald and Walter Gibson*

Subpart 5.4. X-Ray Sources

- Chapter 54. X-Ray Tube Sources *Susanne M. Lee and Carolyn MacDonald*
 Chapter 55. Synchrotron Sources *Steven L. Hulbert and Gwyn P. Williams*
 Chapter 56. Laser Generated Plasmas *Alan Michette*

- Chapter 57. Pinch Plasma Sources *Victor Kantsyrev*
Chapter 58. X-Ray Lasers *Greg Tallents*
Chapter 59. Inverse Compton X-Ray Sources *Frank Carroll*

Subpart 5.5. X-Ray Detectors

- Chapter 60. Introduction to X-Ray Detectors *Walter Gibson and Peter Siddons*
Chapter 61. Advances in Imaging Detectors *Aaron Couture*
Chapter 62. X-Ray Spectral Detection and Imaging *Eric Lifshin*

Subpart 5.6. Neutron Optics and Applications

- Chapter 63. Neutron Optics *David Mildner*
Chapter 64. Grazing-Incidence Neutron Optics *Mikhail Gubarev and Brian Ramsey*

DO NOT DUPLICATE

EDITORS' PREFACE

The third edition of the *Handbook of Optics* is designed to pull together the dramatic developments in both the basic and applied aspects of the field while retaining the archival, reference book value of a handbook. This means that it is much more extensive than either the first edition, published in 1978, or the second edition, with Volumes I and II appearing in 1995 and Volumes III and IV in 2001. To cover the greatly expanded field of optics, the *Handbook* now appears in five volumes. Over 100 authors or author teams have contributed to this work.

Volume I is devoted to the fundamentals, components, and instruments that make optics possible. Volume II contains chapters on design, fabrication, testing, sources of light, detection, and a new section devoted to radiometry and photometry. Volume III concerns vision optics only and is printed entirely in color. In Volume IV there are chapters on the optical properties of materials, nonlinear, quantum and molecular optics. Volume V has extensive sections on fiber optics and x ray and neutron optics, along with shorter sections on measurements, modulators, and atmospheric optical properties and turbulence. Several pages of color inserts are provided where appropriate to aid the reader. A purchaser of the print version of any volume of the *Handbook* will be able to download a digital version containing all of the material in that volume in PDF format to one computer (see download instructions on bound-in card). The combined index for all five volumes can be downloaded from www.HandbookofOpticsOnline.com.

It is possible by careful selection of what and how to present that the third edition of the *Handbook* could serve as a text for a comprehensive course in optics. In addition, students who take such a course would have the *Handbook* as a career-long reference.

Topics were selected by the editors so that the *Handbook* could be a desktop (bookshelf) general reference for the parts of optics that had matured enough to warrant archival presentation. New chapters were included on topics that had reached this stage since the second edition, and existing chapters from the second edition were updated where necessary to provide this compendium. In selecting subjects to include, we also had to select which subjects to leave out. The criteria we applied were: (1) was it a specific application of optics rather than a core science or technology and (2) was it a subject in which the role of optics was peripheral to the central issue addressed. Thus, such topics as medical optics, laser surgery, and laser materials processing were not included. While applications of optics are mentioned in the chapters there is no space in the *Handbook* to include separate chapters devoted to all of the myriad uses of optics in today's world. If we had, the third edition would be much longer than it is and much of it would soon be outdated. We designed the third edition of the *Handbook of Optics* so that it concentrates on the principles of optics that make applications possible.

Authors were asked to try to achieve the dual purpose of preparing a chapter that was a worthwhile reference for someone working in the field and that could be used as a starting point to become acquainted with that aspect of optics. They did that and we thank them for the outstanding results seen throughout the *Handbook*. We also thank Mr. Taisuke Soda of McGraw-Hill for his help in putting this complex project together and Mr. Alan Tourtlotte and Ms. Susannah Lehman of the Optical Society of America for logistical help that made this effort possible.

We dedicate the third edition of the *Handbook of Optics* to all of the OSA volunteers who, since OSA's founding in 1916, give their time and energy to promoting the generation, application, archiving, and worldwide dissemination of knowledge in optics and photonics.

Michael Bass, Editor-in-Chief

Associate Editors:

Casimer M. DeCusatis

Jay M. Enoch

Vasudevan Lakshminarayanan

Guifang Li

Carolyn MacDonald

Virendra N. Mahajan

Eric Van Stryland

This page intentionally left blank.

DO NOT DUPLICATE

PREFACE TO VOLUME IV

Volume IV is a compendium of articles on properties (Chapters 1 to 9), nonlinear optics (Chapters 10 to 19), and quantum and molecular optics (Chapters 20 to 23). As with the rest of the *Handbook*, articles were chosen for their archival nature. Clearly, optical properties of materials fit into the archival category well. This volume devotes a large number of pages to explain and describe the optical properties of water, crystals and glasses, metals, semiconductors, solids in general, thin films and coatings including optical blacks, and photonic bandgap materials. These articles have been updated to include new materials and understanding developed since the previous edition including, among other things, advances in thin-film materials. Nonlinear optics is a mature field, but with many relatively new applications, much of them driven by advances in optical materials. Areas covered here are frequency conversion via second-order nonlinearities including optical parametric oscillators, third-order nonlinearities of two-photon absorption and nonlinear refraction, as well as stimulated Raman and Brillouin scattering, photorefractive materials and devices, coherent optical transients, electromagnetically induced transparency, optical limiting, and laser-induced damage. Nonlinear optical processes for ultrashort pulses is included here and has been a major part of the revolution in sources for obtaining laser pulses now down to attoseconds; however, other chapters on these ultrashort pulses are included in Volume II. Clearly, advances in fiber optic telecommunications have been greatly impacted by nonlinear optics, thus much work in this field is included in the fiber optics chapters in Volume V. The new chapter on laser-induced damage is a much needed addition to the *Handbook* covering a problem from the earliest days of the laser. Chapters on quantum optics in general cover some more modern aspects of optics that have become archival: laser cooling and trapping, where multiple Nobel prizes have recently been awarded; high-field physics that result from the availability of the extreme irradiance produced by lasers; slow light, topics related to being able to slow and even stop light propagation in materials; and correlated states or quantum entanglement, the unusual behavior of quantum systems where optics has played a pivotal role in its understanding as well as some interesting applications in secure communication/cryptography. The chapter on the quantum theory of lasers is, however, included in Volume II. We thank all of the many authors who gave their input to this volume of the *Handbook of Optics*.

Guifang Li and Eric Van Stryland
Associate Editors

This page intentionally left blank.

DO NOT DUPLICATE

GLOSSARY AND FUNDAMENTAL CONSTANTS

Introduction

This glossary of the terms used in the *Handbook* represents to a large extent the language of optics. The symbols are representations of numbers, variables, and concepts. Although the basic list was compiled by the author of this section, all the editors have contributed and agreed to this set of symbols and definitions. Every attempt has been made to use the same symbols for the same concepts throughout the entire *Handbook*, although there are exceptions. Some symbols seem to be used for many concepts. The symbol α is a prime example, as it is used for absorptivity, absorption coefficient, coefficient of linear thermal expansion, and more. Although we have tried to limit this kind of redundancy, we have also bowed deeply to custom.

Units

The abbreviations for the most common units are given first. They are consistent with most of the established lists of symbols, such as given by the International Standards Organization ISO¹ and the International Union of Pure and Applied Physics, IUPAP.²

Prefixes

Similarly, a list of the numerical prefixes¹ that are most frequently used is given, along with both the common names (where they exist) and the multiples of ten that they represent.

Fundamental Constants

The values of the fundamental constants³ are listed following the sections on SI units.

Symbols

The most commonly used symbols are then given. Most chapters of the *Handbook* also have a glossary of the terms and symbols specific to them for the convenience of the reader. In the following list, the symbol is given, its meaning is next, and the most customary unit of measure for the quantity is presented in brackets. A bracket with a dash in it indicates that the quantity is unitless. Note that there is a difference between units and dimensions. An angle has units of degrees or radians and a solid angle square degrees or steradians, but both are pure ratios and are dimensionless. The unit symbols as recommended in the SI system are used, but decimal multiples of some of the dimensions are sometimes given. The symbols chosen, with some cited exceptions, are also those of the first two references.

RATIONALE FOR SOME DISPUTED SYMBOLS

The choice of symbols is a personal decision, but commonality improves communication. This section explains why the editors have chosen the preferred symbols for the *Handbook*. We hope that this will encourage more agreement.

Fundamental Constants

It is encouraging that there is almost universal agreement for the symbols for the fundamental constants. We have taken one small exception by adding a subscript B to the k for Boltzmann's constant.

Mathematics

We have chosen i as the imaginary almost arbitrarily. IUPAP lists both i and j , while ISO does not report on these.

Spectral Variables

These include expressions for the wavelength λ , frequency ν , wave number σ , ω for circular or radian frequency, k for circular or radian wave number and dimensionless frequency x . Although some use f for frequency, it can be easily confused with electronic or spatial frequency. Some use $\tilde{\nu}$ for wave number, but, because of typography problems and agreement with ISO and IUPAP, we have chosen σ ; it should not be confused with the Stefan-Boltzmann constant. For spatial frequencies we have chosen ξ and η , although f_x and f_y are sometimes used. ISO and IUPAP do not report on these.

Radiometry

Radiometric terms are contentious. The most recent set of recommendations by ISO and IUPAP are L for radiance [$\text{Wcm}^{-2}\text{sr}^{-1}$], M for radiant emittance or exitance [Wcm^{-2}], E for irradiance or incidence [Wcm^{-2}], and I for intensity [Wsr^{-2}]. The previous terms, W , H , N , and J , respectively, are still in many texts, notably Smith⁴ and Lloyd⁵ but we have used the revised set, although there are still shortcomings. We have tried to deal with the vexatious term *intensity* by using *specific intensity* when the units are $\text{Wcm}^{-2}\text{sr}^{-1}$, *field intensity* when they are Wcm^{-2} , and *radiometric intensity* when they are Wsr^{-1} .

There are two sets of terms for these radiometric quantities, which arise in part from the terms for different types of reflection, transmission, absorption, and emission. It has been proposed that the *ion* ending indicate a process, that the *ance* ending indicate a value associated with a particular sample, and that the *ivity* ending indicate a generic value for a "pure" substance. Then one also has reflectance, transmittance, absorptance, and emittance as well as reflectivity, transmissivity, absorptivity, and emissivity. There are now two different uses of the word emissivity. Thus the words *exitance*, *incidence*, and *sterance* were coined to be used in place of emittance, irradiance, and radiance. It is interesting that ISO uses radiance, exitance, and irradiance whereas IUPAP uses radiance exitance [*sic*], and irradiance. We have chosen to use them both, i.e., emittance, irradiance, and radiance will be followed in square brackets by exitance, incidence, and sterance (or vice versa). Individual authors will use the different endings for transmission, reflection, absorption, and emission as they see fit.

We are still troubled by the use of the symbol E for irradiance, as it is so close in meaning to electric field, but we have maintained that accepted use. The spectral concentrations of these quantities, indicated by a wavelength, wave number, or frequency subscript (e.g., L_λ) represent partial differentiations; a subscript q represents a photon quantity; and a subscript v indicates a quantity normalized to the response of the eye. Thereby, L_v is luminance, E_v illuminance, and M_v and I_v luminous emittance and luminous intensity. The symbols we have chosen are consistent with ISO and IUPAP.

The refractive index may be considered a radiometric quantity. It is generally complex and is indicated by $\tilde{n} = n - ik$. The real part is the relative refractive index and k is the extinction coefficient. These are consistent with ISO and IUPAP, but they do not address the complex index or extinction coefficient.

Optical Design

For the most part ISO and IUPAP do not address the symbols that are important in this area.

There were at least 20 different ways to indicate focal ratio; we have chosen FN as symmetrical with NA; we chose f and efl to indicate the effective focal length. Object and image distance, although given many different symbols, were finally called s_o and s_i since s is an almost universal symbol for distance. Field angles are θ and ϕ ; angles that measure the slope of a ray to the optical axis are u ; u can also be $\sin u$. Wave aberrations are indicated by W_{ijk} , while third-order ray aberrations are indicated by σ_i and more mnemonic symbols.

Electromagnetic Fields

There is no argument about \mathbf{E} and \mathbf{H} for the electric and magnetic field strengths, Q for quantity of charge, ρ for volume charge density, σ for surface charge density, etc. There is no guidance from Refs. 1 and 2 on polarization indication. We chose \perp and \parallel rather than p and s , partly because s is sometimes also used to indicate scattered light.

There are several sets of symbols used for reflection transmission, and (sometimes) absorption, each with good logic. The versions of these quantities dealing with field amplitudes are usually specified with lower case symbols: r , t , and a . The versions dealing with power are alternately given by the uppercase symbols or the corresponding Greek symbols: R and T versus ρ and τ . We have chosen to use the Greek, mainly because these quantities are also closely associated with Kirchhoff's law that is usually stated symbolically as $\alpha = \epsilon$. The law of conservation of energy for light on a surface is also usually written as $\alpha + \rho + \tau = 1$.

Base SI Quantities

length	m	meter
time	s	second
mass	kg	kilogram
electric current	A	ampere
temperature	K	kelvin
amount of substance	mol	mole
luminous intensity	cd	candela

Derived SI Quantities

energy	J	joule
electric charge	C	coulomb
electric potential	V	volt
electric capacitance	F	farad
electric resistance	Ω	ohm
electric conductance	S	siemens
magnetic flux	Wb	weber
inductance	H	henry
pressure	Pa	pascal
magnetic flux density	T	tesla
frequency	Hz	hertz
power	W	watt
force	N	newton
angle	rad	radian
angle	sr	steradian

Prefixes

<i>Symbol</i>	<i>Name</i>	<i>Common name</i>	<i>Exponent of ten</i>
F	exa		18
P	peta		15
T	tera	trillion	12
G	giga	billion	9
M	mega	million	6
k	kilo	thousand	3
h	hecto	hundred	2
da	deca	ten	1
d	deci	tenth	-1
c	centi	hundredth	-2
m	milli	thousandth	-3
μ	micro	millionth	-6
n	nano	billionth	-9
p	pico	trillionth	-12
f	femto		-15
a	atto		-18

Constants

c	speed of light vacuo [299792458 ms ⁻¹]
c_1	first radiation constant = $2\pi^2 h = 3.7417749 \times 10^{-16}$ [Wm ²]
c_2	second radiation constant = $hc/k = 0.014838769$ [mK]
e	elementary charge [$1.60217733 \times 10^{-19}$ C]
g_n	free fall constant [9.80665 ms ⁻²]
h	Planck's constant [$6.6260755 \times 10^{-34}$ Ws]
k_B	Boltzmann constant [1.380658×10^{-23} JK ⁻¹]
m_e	mass of the electron [$9.1093897 \times 10^{-31}$ kg]
N_A	Avogadro constant [6.0221367×10^{23} mol ⁻¹]
R_∞	Rydberg constant [10973731.534 m ⁻¹]
ϵ_0	vacuum permittivity [$\mu_0^{-1}c^{-2}$]
σ	Stefan-Boltzmann constant [5.67051×10^{-8} Wm ⁻¹ K ⁻⁴]
μ_0	vacuum permeability [$4\pi \times 10^{-7}$ NA ⁻²]
μ_B	Bohr magneton [$9.2740154 \times 10^{-24}$ JT ⁻¹]

General

B	magnetic induction [Wbm ⁻² , kgs ⁻¹ C ⁻¹]
C	capacitance [f, C ² s ² m ⁻² kg ⁻¹]
C	curvature [m ⁻¹]
c	speed of light in vacuo [ms ⁻¹]
c_1	first radiation constant [Wm ²]
c_2	second radiation constant [mK]
D	electric displacement [Cm ⁻²]
E	incidence [irradiance] [Wm ⁻²]
e	electronic charge [coulomb]
E_v	illuminance [lux, lmm ⁻²]
E	electrical field strength [Vm ⁻¹]
E	transition energy [J]
E_g	band-gap energy [eV]
f^g	focal length [m]
f_f^g	Fermi occupation function, conduction band
f_v^g	Fermi occupation function, valence band

FN	focal ratio (f /number) [—]
g	gain per unit length [m^{-1}]
g_{th}	gain threshold per unit length [m^{-1}]
H	magnetic field strength [Am^{-1} , $\text{Cs}^{-1} \text{m}^{-1}$]
h	height [m]
I	irradiance (see also E) [Wm^{-2}]
I	radiant intensity [Wsr^{-1}]
I	nuclear spin quantum number [—]
I	current [A]
i	$\sqrt{-1}$
$\text{Im}()$	imaginary part of
J	current density [Am^{-2}]
j	total angular momentum [$\text{kg m}^2 \text{s}^{-1}$]
$J_1()$	Bessel function of the first kind [—]
k	radian wave number $=2\pi/\lambda$ [rad cm^{-1}]
k	wave vector [rad cm^{-1}]
k	extinction coefficient [—]
L	sterance [radiance] [$\text{Wm}^{-2} \text{sr}^{-1}$]
L_v	luminance [cdm^{-2}]
L	inductance [h, $\text{m}^2 \text{kg C}^2$]
L	laser cavity length
L, M, N	direction cosines [—]
M	angular magnification [—]
M	radiant exitance [radiant emittance] [Wm^{-2}]
m	linear magnification [—]
m	effective mass [kg]
MTF	modulation transfer function [—]
N	photon flux [s^{-1}]
N	carrier (number) density [m^{-3}]
n	real part of the relative refractive index [—]
\tilde{n}	complex index of refraction [—]
NA	numerical aperture [—]
OPD	optical path difference [m]
P	macroscopic polarization [C m^{-2}]
$\text{Re}()$	real part of [—]
R	resistance [Ω]
r	position vector [m]
S	Seebeck coefficient [VK^{-1}]
s	spin quantum number [—]
s	path length [m]
S_o	object distance [m]
S_i	image distance [m]
T	temperature [K, C]
t	time [s]
t	thickness [m]
u	slope of ray with the optical axis [rad]
V	Abbe reciprocal dispersion [—]
V	voltage [V , $\text{m}^2 \text{kg s}^{-2} \text{C}^{-1}$]
x, y, z	rectangular coordinates [m]
Z	atomic number [—]

Greek Symbols

α	absorption coefficient [cm^{-1}]
α	(power) absorptance (absorptivity)

ϵ	dielectric coefficient (constant) [—]
ϵ	emittance (emissivity) [—]
ϵ	eccentricity [—]
ϵ_1	Re (ϵ)
ϵ_2	Im (ϵ)
τ	(power) transmittance (transmissivity) [—]
ν	radiation frequency [Hz]
ω	circular frequency = $2\pi\nu$ [rads ⁻¹]
ω	plasma frequency [Hz]
λ	wavelength [μm , nm]
σ	wave number = $1/\lambda$ [cm ⁻¹]
σ	Stefan Boltzmann constant [Wm ⁻² K ⁻¹]
ρ	reflectance (reflectivity) [—]
θ, ϕ	angular coordinates [rad, °]
ξ, η	rectangular spatial frequencies [m ⁻¹ , r ⁻¹]
ϕ	phase [rad, °]
ϕ	lens power [m ⁻²]
Φ	flux [W]
χ	electric susceptibility tensor [—]
Ω	solid angle [sr]

Other

\Re	responsivity
$\exp(x)$	e^x
$\log_a(x)$	log to the base a of x
$\ln(x)$	natural log of x
$\log(x)$	standard log of x : $\log_{10}(x)$
Σ	summation
Π	product
Δ	finite difference
δx	variation in x
dx	total differential
∂x	partial derivative of x
$\delta(x)$	Dirac delta function of x
δ_{ij}	Kronecker delta

REFERENCES

1. Anonymous, *ISO Standards Handbook 2: Units of Measurement*, 2nd ed., International Organization for Standardization, 1982.
2. Anonymous, *Symbols, Units and Nomenclature in Physics*, Document U.I.P. 20, International Union of Pure and Applied Physics, 1978.
3. E. Cohen and B. Taylor, "The Fundamental Physical Constants," *Physics Today*, 9 August 1990.
4. W. J. Smith, *Modern Optical Engineering*, 2nd ed., McGraw-Hill, 1990.
5. J. M. Lloyd, *Thermal Imaging Systems*, Plenum Press, 1972.

William L. Wolfe
 College of Optical Sciences
 University of Arizona
 Tucson, Arizona

PART

1

PROPERTIES

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

OPTICAL PROPERTIES OF WATER

Curtis D. Mobley

*Applied Electromagnetics and Optics Laboratory
SRI International
Menlo Park, California*

1.1 INTRODUCTION

This article discusses the optical properties of three substances: pure water, pure sea water, and natural water. Pure water (i.e., water molecules only) without any dissolved substances, ions, bubbles, or other impurities, is exceptionally difficult to produce in the laboratory. For this and other reasons, definitive direct measurements of its optical properties at visible wavelengths have not yet been made. Pure sea water—pure water plus various dissolved salts—has optical properties close to those of pure water. Neither pure water nor pure sea water ever occur in nature. Natural waters, both fresh and saline, are a witch's brew of dissolved and particulate matter. These solutes and particulates are both optically significant and highly variable in kind and concentration. Consequently, the optical properties of natural waters show large temporal and spatial variations and seldom resemble those of pure water.

The great variability of the optical properties of natural water is the bane of those who desire precise and easily tabulated data. However, it is the connections between the optical properties and the biological, chemical, and geological constituents of natural water and the physical environment that define the critical role of optics in aquatic research. For just as optics utilizes results from the biological, chemical, geological, and physical subdisciplines of limnology and oceanography, so do those subdisciplines incorporate optics. This synergism is seen in such areas as bio-optical oceanography, marine photochemistry, mixed-layer dynamics, laser bathymetry, and remote sensing of biological productivity, sediment load, or pollutants.

1.2 TERMINOLOGY, NOTATION, AND DEFINITIONS

Hydrologic optics is the quantitative study of the interactions of radiant energy with the earth's oceans, estuaries, lakes, rivers, and other water bodies. Most past and current research within hydrologic optics has been within the subfield of oceanic optics, in particular the optics of deep ocean waters, as opposed to coastal or estuarine areas. This emphasis is reflected in our uneven understanding of the optical properties of various water types.

Although the optical properties of different water bodies can vary greatly, there is an overall similarity that is quite distinct from, say, the optical properties of the atmosphere. Therefore, hydrologic and atmospheric optics have developed considerably different theoretical formulations, experimental methodologies, and instrumentation as suited to each field's specific scientific issues. Chapter 3, "Atmospheric Optics," by Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman in Vol. V discusses atmospheric optics. The text by Mobley¹ gives a comprehensive treatment of hydrologic optics.

Radiative transfer theory is the framework that connects the optical properties of water with the ambient light field. A rigorous mathematical formulation of radiative transfer theory as applicable to hydrologic optics has been developed by Preisendorfer² and others. Preisendorfer found it convenient to divide the optical properties of water into two classes: inherent and apparent. *Inherent optical properties* (IOPs) are those properties that depend only upon the medium and therefore are independent of the ambient light field within the medium. The two fundamental IOPs are the absorption coefficient and the volume scattering function (VSF). Other IOPs include the attenuation coefficient and the single-scattering albedo. *Apparent optical properties* (AOPs) are those properties that depend both on the medium (the IOPs) and on the geometric (directional) structure of the ambient light field and that display enough regular features and stability to be useful descriptors of the water body. Commonly used AOPs are the irradiance reflectance, the average cosines, and the various attenuation functions (*K* functions). (All of these quantities are defined below.) The radiative transfer equation provides the connection between the IOPs and the AOPs. The physical environment of a water body—waves on its surface, the character of its bottom, the incident radiance from the sky—enters the theory via the boundary conditions necessary for solution of the radiative transfer equation.

The IOPs are easily defined but they can be exceptionally difficult to measure, especially in situ. The AOPs are generally much easier to measure, but they are difficult to interpret because of the confounding environmental effects. (A change in the sea surface wave state or in the sun's position changes the radiance distribution, and hence the AOPs, even though the IOPs are unchanged.)

Hydrologic optics employs standard radiometric concepts and terminology, although the notation adopted by the International Association for Physical Sciences of the Ocean (IAPSO³) differs somewhat from that used in other fields. Table 1 summarizes the terms, units, and symbols for those quantities that have proven most useful in hydrologic optics. These quantities are defined and discussed in Secs. 1.3 to 1.5. Figure 1 summarizes the relationships among the various inherent and apparent optical properties. In this figure, note the central unifying role of radiative transfer theory. Note also that the spectral absorption coefficient and the spectral volume scattering function are the fundamental inherent optical properties in the sense that all inherent optical properties are derivable from those two. Likewise, spectral radiance is the parent of all radiometric quantities and apparent optical properties. The source term *S* in the radiative transfer equation accounts both for true internal sources such as bioluminescence and for radiance appearing at the wavelength of interest owing to inelastic scattering from other wavelengths.

Most radiative transfer theory assumes the radiant energy to be monochromatic. In this case the associated optical properties and radiometric quantities are termed *spectral* and carry a wavelength (λ) argument or subscript [e.g., the spectral absorption coefficient $a(\lambda)$ or a_λ or the spectral downward irradiance $E_d(\lambda)$]. Spectral radiometric quantities have the SI unit nm^{-1} added to the units shown in Table 1 [e.g., $E_d(\lambda)$ has units $\text{W m}^{-2} \text{nm}^{-1}$]. Many radiometric on the other hand, respond to a fairly wide bandwidth, which complicates the comparison of data and theory.

1.3 RADIOMETRIC QUANTITIES USEFUL IN HYDROLOGIC OPTICS

Consider an amount ΔQ of radiant energy incident in a time interval Δt centered on time t , onto a surface of area ΔA located at (x, y, z) . The energy arrives through a set of directions contained in a solid angle $\Delta\Omega$ about the direction (θ, ϕ) normal to the area ΔA and is produced by photons in

TABLE 1 Terms, Units, and Symbols for Quantities Commonly Used in Hydrologic Optics

Quantity	SI Units	IAPSO Recommended Symbol*	Historic Symbol† (if different)
Fundamental quantities			
Most of the fundamental quantities are not defined by IAPSO, in which case common usage is given.			
geometric depth below water surface	m	z	
polar angle of photon travel	radian or degree	θ	
wavelength of light (in vacuo)	nm	λ	
cosine of polar angle	dimensionless	$\mu \equiv \cos \theta$	
optical depth below water surface	dimensionless	τ	
azimuthal angle of photon travel	radian or degree	ϕ	
scattering angle	radian or degree	ψ, γ or Θ	θ
solid angle	sr	Ω or ω	Ω
Radiometric quantities			
The quantities as shown represent broadband measurements. For narrowband (monochromatic) measurements add the adjective "spectral" to the term, add nm^{-1} to the units, and add a wavelength index λ to the symbol [e.g., spectral radiance, L_λ or $L(\lambda)$] with units $\text{W m}^{-2} \text{sr}^{-1} \text{nm}^{-1}$. PAR is always broadband.			
(plane) irradiance	W m^{-2}	E	H
downward irradiance	W m^{-2}	E_d	$H(-)$
upward irradiance	W m^{-2}	E_v	$H(+)$
net (vertical) irradiance	W m^{-2}	\bar{E}	\bar{H}
scalar irradiance	W m^{-2}	E_0	h
downward scalar irradiance	W m^{-2}	E_{0d}	$h(-)$
upward scalar irradiance	W m^{-2}	E_{0u}	$h(+)$
radiant intensity	W sr^{-1}	I	J
radiance	$\text{W m}^{-2} \text{sr}^{-1}$	L	N
radiant excitation	W m^{-2}	M	W
photosynthetically available radiation	$\text{photons s}^{-1} \text{m}^{-2}$	PAR or E_{PAR}	
quantity of radiant energy	J	Q	U
radiant power	W	Φ	P
Inherent optical properties			
absorptance	dimensionless	A	
absorption coefficient	m^{-1}	a	
scatterance	dimensionless	B	
scattering coefficient	m^{-1}	b	s
backward scattering coefficient	m^{-1}	b_b	b
forward scattering coefficient	m^{-1}	b_f	f
attenuance	dimensionless	C	
attenuation coefficient	m^{-1}	c	α
(real) index of refraction	dimensionless	n	
transmittance	dimensionless	T	
volume scattering function	$\text{m}^{-1} \text{sr}^{-1}$	β	σ
scattering phase function	sr^{-1}	$\tilde{\beta}$	P
single-scattering albedo	dimensionless	ω_0 or $\tilde{\omega}$	
Apparent optical properties			
(vertical) attenuation coefficients			
of downward irradiance $E_d(z)$	m^{-1}	K_d	$K(-)$
of total scalar irradiance $E_0(z)$	m^{-1}	K_0	k

(Continued)

TABLE 1 Terms, Units, and Symbols for Quantities Commonly Used in Hydrologic Optics (*Continued*)

Quantity	SI Units	IAPSO Recommended Symbol ^a	Historic Symbol [†] (if different)
Apparent optical properties			
of downward scalar irradiance $E_{0d}(z)$	m^{-1}	K_{0d}	$k(-)$
of upward scalar irradiance $E_{0u}(z)$	m^{-1}	K_{0u}	$k(+)$
of PAR	m^{-1}	K_{PAR}	
of upward irradiance $E_u(z)$	m^{-1}	K_u	$K(+)$
of radiance $L(z, \theta, \phi)$	m^{-1}	$K(\theta, \phi)$	
irradiance reflectance (ratio)	dimensionless	R	$R(-)$
average cosine of light field	dimensionless	$\bar{\mu}$	
of downwelling light	dimensionless	$\bar{\mu}_d$	$D(-)=1/\bar{\mu}_d$
of upwelling light	dimensionless	$\bar{\mu}_u$	$D(+)=1/\bar{\mu}_u$
distribution function	dimensionless		$D=1/\bar{\mu}$

^aReferences 1 and 3.[†]Reference 2.

a wavelength interval $\Delta\lambda$ centered on wavelength λ . Then an *operational* definition of the *spectral radiance* is

$$L(x, y, z, t, \theta, \phi, \lambda) \equiv \frac{\Delta Q}{\Delta t \Delta A \Delta \Omega \Delta \lambda} \quad \text{Js}^{-1} \text{m}^{-2} \text{sr}^{-1} \text{nm}^{-1}$$

In practice, one takes Δt , ΔA , $\Delta \Omega$, and $\Delta \lambda$ small enough to get a useful resolution of radiance over the four parameter domains but not so small as to encounter diffraction effects or fluctuations from photon shot noise at very low light levels. Typical values are $\Delta t \sim 10^{-3}$ to 10^3 s (depending on whether or not one wishes to average out sea surface wave effects), $\Delta A \sim 10^{-3} \text{m}^2$, $\Delta \Omega \sim 10^{-2}$ sr, and $\Delta \lambda \sim 10$ nm. In the conceptual limit of infinitesimal parameter intervals, the spectral radiance is defined as

$$L(x, y, z, t, \theta, \phi, \lambda) \equiv \frac{\partial^4 Q}{\partial t \partial A \partial \Omega \partial \lambda} \quad \text{Js}^{-1} \text{m}^{-2} \text{sr}^{-1} \text{nm}^{-1}$$

Spectral radiance is the fundamental radiometric quantity of interest in hydrologic optics: it specifies the positional (x, y, z), temporal (t), directional (θ, ϕ) and spectral (λ) structure of the light field. For typical oceanic environments, horizontal variations (on a scale of tens to thousands of meters) of inherent and apparent optical properties are much less than variations with depth, and it is usually assumed that these properties vary only with depth z . Moreover, since the time scales for changes in IOPs or in the environment (seconds to seasons) are much greater than the time required for the radiance field to reach steady state (microseconds) after a change in IOPs or boundary conditions, time-independent radiative transfer theory is adequate for most hydrologic optics studies. The spectral radiance therefore usually is written $L(z, \theta, \phi, \lambda)$. The exceptions are applications such as time-of-flight lidar.

There are few conventions on the choice of coordinate systems. Oceanographers usually measure the depth z positive downward from $z = 0$ at the mean water surface. In radiative transfer theory it is convenient to let (θ, ϕ) denote the direction of photon travel (especially when doing Monte Carlo simulations). When displaying data it is convenient to let (θ, ϕ) represent the direction in which the instrument was pointed (the viewing direction) in order to detect photons traveling in

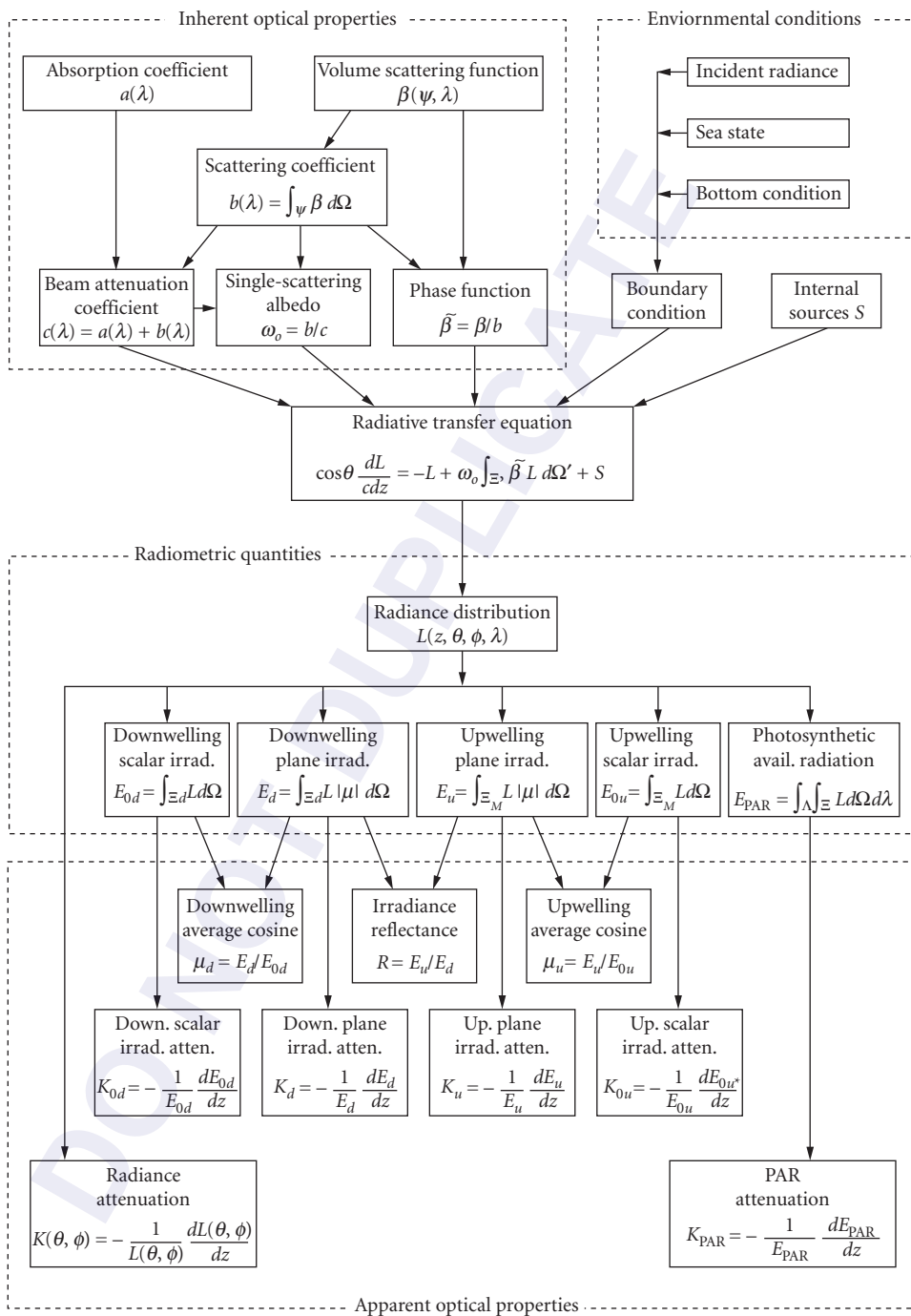


FIGURE 1 Relationships between the various quantities commonly used in hydrologic optics.

the opposite direction. Some authors measure the polar angle θ from the zenith (upward) direction, even when z is taken as positive downward; others measure θ from the $+z$ axis (nadir, or downward, direction). In the following discussion Ξ_u denotes the hemisphere of *upward* directions (i.e., the set of directions (θ, ϕ) such that $0 \leq \theta \leq \pi/2$ and $0 \leq \phi \leq 2\pi$ if θ is measured from the zenith direction) and Ξ_d denotes the hemisphere of *downward* directions. The element of solid angle is $d\Omega = \sin\theta d\theta d\phi$ with units of steradian. The solid angle measure of the set of directions Ξ_u or Ξ_d is $\Omega(\Xi_u) = \Omega(\Xi_d) = 2\pi$ sr.

Although the spectral radiance completely specifies the light field, it is seldom measured both because of instrumental difficulties and because such complete information often is not needed for specific applications. The most commonly measured radiometric quantities are various *irradiances*.

Consider a light detector constructed so as to be equally sensitive to photons of a given wavelength λ traveling in any direction (θ, ϕ) within a *hemisphere* of directions.⁴ If the detector is located at depth z and is oriented facing *upward*, so as to collect photons traveling *downward*, then the detector output is a measure of the *spectral downward scalar irradiance* at depth z , $E_{0d}(z, \lambda)$. Such an instrument is summing radiance over all the directions (elements of solid angle) in the downward hemisphere Ξ_d . Thus $E_{0d}(z, \lambda)$ is related to $L(z, \theta, \phi, \lambda)$ by

$$E_{0d}(z, \lambda) = \int_{\Xi_d} L(z, \theta, \phi, \lambda) d\Omega \quad \text{W m}^{-2} \text{ nm}^{-1}$$

The symbolic integral over Ξ_d can be evaluated as a double integral over θ and ϕ after a specific coordinate system is chosen.

If the same instrument is oriented facing *downward*, so as to detect photons traveling *upward*, then the quantity measured is the *spectral upward scalar irradiance* $E_{0u}(z, \lambda)$:

$$E_{0u}(z, \lambda) = \int_{\Xi_u} L(z, \theta, \phi, \lambda) d\Omega \quad \text{W m}^{-2} \text{ nm}^{-1}$$

The *spectral scalar irradiance* $E_0(z, \lambda)$ is just the sum of the downward and upward components:

$$E_0(z, \lambda) \equiv E_{0d}(z, \lambda) + E_{0u}(z, \lambda) = \int_{\Xi} L(z, \theta, \phi, \lambda) d\Omega \quad \text{W m}^{-2} \text{ nm}^{-1}$$

Here $\Xi = \Xi_d \cup \Xi_u$ is the set of all directions; $\Omega(\Xi) = 4\pi$ sr. $E_0(z, \lambda)$ is useful² because it is proportional to the spectral radiant energy density ($\text{J m}^{-3} \text{ nm}^{-1}$) at depth z .

Now consider a detector designed⁴ so that its sensitivity is proportional to $|\cos\theta|$, where θ is the angle between the photon direction and the normal to the surface of the detector. This is the ideal response of a “flat plate” collector of area ΔA , which when viewed at an angle θ to its normal appears to have an area of $\Delta A |\cos\theta|$. If such a detector is located at depth z and is oriented facing *upward*, so as to detect photons traveling *downward*, then its output is proportional to the *spectral downward plane irradiance* $E_d(z, \lambda)$ (usually called *spectral downwelling irradiance*). This instrument is summing the downwelling radiance weighted by the cosine of the photon direction:

$$E_d(z, \lambda) = \int_{\Xi_d} L(z, \theta, \phi, \lambda) |\cos\theta| d\Omega \quad \text{W m}^{-2} \text{ nm}^{-1}$$

Turning this instrument upside down gives the *spectral upward plane irradiance* (*spectral upwelling irradiance*) $E_u(z, \lambda)$:

$$E_u(z, \lambda) = \int_{\Xi_u} L(z, \theta, \phi, \lambda) |\cos\theta| d\Omega \quad \text{W m}^{-2} \text{ nm}^{-1}$$

E_d and E_u are useful because they give the energy flux (power per unit area) across the horizontal surface at depth z owing to downwelling and upwelling photons, respectively.

The *spectral net irradiance* at depth z , $\bar{E}(z, \lambda)$ is the difference in the downwelling and upwelling plane irradiances:

$$\bar{E}(z, \lambda) = E_d(z, \lambda) - E_u(z, \lambda)$$

Photosynthesis is a quantum phenomenon (i.e., it is the *number* of available photons rather than the *amount* of radiant energy that is relevant to the chemical transformations). This is because if a photon of, say, $\lambda = 350$ nm, is absorbed by chlorophyll it induces the same chemical change as does a photon of $\lambda = 700$ nm, even though the 350-nm photon has twice the energy of the 700-nm photon. Only a part of the photon energy goes into photosynthesis; the excess is converted to heat or reradiated. Moreover, chlorophyll is equally able to absorb and utilize a photon regardless of the photon's direction of travel. Therefore, in studies of phytoplankton biology the relevant measure of the light field is the *photosynthetically available radiation*, PAR or E_{PAR} , defined by

$$\text{PAR}(z) \equiv \int_{350 \text{ nm}}^{700 \text{ nm}} \frac{\lambda}{hc} E_0(z, \lambda) d\lambda \quad \text{photons s}^{-1} \text{ m}^{-2}$$

where $h = 6.6255 \times 10^{-34}$ J s is Planck's constant and $c = 3.0 \times 10^{17}$ nm s⁻¹ is the speed of light. The factor λ/hc converts the energy units of E_0 (watts) to quantum units (photons per second). Bio-optical literature often states PAR values in units of mol photons s⁻¹ m⁻² or einst s⁻¹ m⁻². Morel and Smith⁵ found that over a wide variety of water types from very clear to turbid, with corresponding variations in the spectral nature of the irradiance, the conversion factor for energy to quanta varied by only ± 10 percent about the value 2.5×10^{18} photons s⁻¹ W⁻¹ (4.2 $\mu\text{einst s}^{-1} \text{ W}^{-1}$).

For practical reasons related to instrument design, PAR is sometimes estimated using the spectral downwelling plane irradiance and the visible wavelengths only:

$$\text{PAR}(z) \approx \int_{400 \text{ nm}}^{700 \text{ nm}} \frac{\lambda}{hc} E_d(z, \lambda) d\lambda \quad \text{photons s}^{-1} \text{ m}^{-2}$$

However, it is now recognized^{6,7} that the use of E_d rather than E_0 can lead to errors of 20 to 100 percent in computations of PAR. Omission of the 350–400-nm band is less troublesome since those wavelengths are rapidly absorbed near the water surface, except in very clear waters.

1.4 INHERENT OPTICAL PROPERTIES

Consider a small volume ΔV of water of thickness Δr as seen by a narrow collimated beam of monochromatic light of spectral radiant power $\Phi_i(\lambda)$ W nm⁻¹ as schematically illustrated in Fig. 2. Some part $\Phi_a(\lambda)$ of the incident power $\Phi_i(\lambda)$ is absorbed within the volume of water. Some part $\Phi_s(\psi, \lambda)$ is scattered out of the beam at an angle ψ , and the remaining power $\Phi_t(\lambda)$ is transmitted through the volume with no change in direction. Let $\Phi_s(\lambda)$ be the total power that is scattered into all directions. Furthermore, assume that no inelastic scattering occurs (i.e., assume that no photons undergo a change in wavelength during the scattering process). Then by conservation of energy,

$$\Phi_i(\lambda) = \Phi_a(\lambda) + \Phi_s(\lambda) + \Phi_t(\lambda)$$

The *spectral absorptance* $A(\lambda)$ is the fraction of incident power that is absorbed within the volume:

$$A(\lambda) \equiv \frac{\Phi_a(\lambda)}{\Phi_i(\lambda)}$$

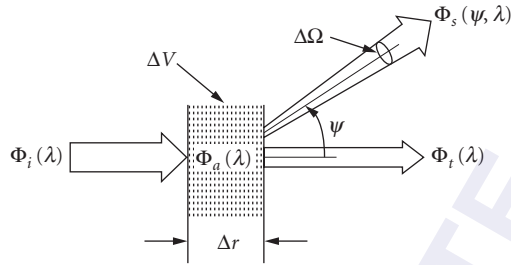


FIGURE 2 Geometry used to define inherent optical properties.

Likewise the *spectral scatterance* $B(\lambda)$ is the fractional part of the incident power that is scattered out of the beam,

$$B(\lambda) \equiv \frac{\Phi_s(\lambda)}{\Phi_i(\lambda)}$$

and the *spectral transmittance* $T(\lambda)$ is

$$T(\lambda) \equiv \frac{\Phi_t(\lambda)}{\Phi_i(\lambda)}$$

Clearly, $A(\lambda) + B(\lambda) + T(\lambda) = 1$. A quantity easily confused with the absorptance $A(\lambda)$ is the absorbance $D(\lambda)$ (also called *optical density*) defined as⁸

$$D(\lambda) \equiv \log_{10} \frac{\Phi_i(\lambda)}{\Phi_s(\lambda) + \Phi_t(\lambda)} = -\log_{10}[1 - A(\lambda)]$$

$D(\lambda)$ is the quantity actually measured in a spectrophotometer.

The inherent optical properties usually employed in hydrologic optics are the spectral absorption and scattering coefficients which are, respectively, the spectral absorptance and scatterance *per unit distance* in the medium. In the geometry of Fig. 2, the *spectral absorption coefficient* $a(\lambda)$ is defined as

$$a(\lambda) \equiv \lim_{\Delta r \rightarrow 0} \frac{A(\lambda)}{\Delta r} \quad \text{m}^{-1}$$

and the *spectral scattering coefficient* $b(\lambda)$ is

$$b(\lambda) \equiv \lim_{\Delta r \rightarrow 0} \frac{B(\lambda)}{\Delta r} \quad \text{m}^{-1}$$

The *spectral beam attenuation coefficient* $c(\lambda)$ is defined as

$$c(\lambda) \equiv a(\lambda) + b(\lambda)$$

Hydrologic optics uses the term attenuation rather than extinction.

Now take into account the angular distribution of the scattered power, with $B(\psi, \lambda)$ being the fraction of incident power scattered out of the beam through an angle ψ into a solid angle

$\Delta\Omega$ centered on ψ as shown in Fig. 2. Then the angular scatterance per unit distance and unit solid angle, $\beta(\psi, \lambda)$, is

$$\beta(\psi, \lambda) \equiv \lim_{\Delta r \rightarrow 0} \lim_{\Delta\Omega \rightarrow 0} \frac{B(\psi, \lambda)}{\Delta r \Delta\Omega} = \lim_{\Delta r \rightarrow 0} \lim_{\Delta\Omega \rightarrow 0} \frac{\Phi_s(\psi, \lambda)}{\Phi_i(\lambda) \Delta r \Delta\Omega} \quad \text{m}^{-1} \text{sr}^{-1}$$

The spectral power scattered into the given solid angle $\Delta\Omega$ is just the spectral radiant intensity scattered into direction ψ times the solid angle: $\Phi_s(\psi, \lambda) = I_s(\psi, \lambda) \Delta\Omega$. Moreover, if the incident power $\Phi_i(\lambda)$ falls on an area ΔA , then the corresponding incident irradiance $E_i(\lambda) = \Phi_i(\lambda) / \Delta A$. Noting that $\Delta V = \Delta r \Delta A$ is the volume of water that is illuminated by the incident beam gives

$$\beta(\psi, \lambda) = \lim_{\Delta V \rightarrow 0} \frac{I_s(\psi, \lambda)}{E_i(\lambda) \Delta V}$$

This form of $\beta(\psi, \lambda)$ suggests the name *spectral volume scattering function* and the physical interpretation of scattered intensity per unit incident irradiance per unit volume of water; $\beta(\psi, \lambda)$ also can be interpreted as the differential scattering cross section per unit volume. Integrating $\beta(\psi, \lambda)$ over all directions (solid angles) gives the total scattered power per unit incident irradiance and unit volume of water or, in other words, the spectral scattering coefficient:

$$b(\lambda) = \int_{\Xi} \beta(\psi, \lambda) d\Omega = 2\pi \int_0^\pi \beta(\psi, \lambda) \sin \psi \, d\psi$$

The last equation follows because scattering in natural waters is azimuthally symmetric about the incident direction (for unpolarized sources and for randomly oriented scatterers). This integration is often divided into forward scattering, $0 \leq \psi \leq \pi/2$, and backward scattering, $\pi/2 \leq \psi \leq \pi$, parts. The corresponding spectral forward and backward scattering coefficients are, respectively,

$$b_f(\lambda) \equiv 2\pi \int_0^{\pi/2} \beta(\psi, \lambda) \sin \psi \, d\psi$$

$$b_b(\lambda) \equiv 2\pi \int_{\pi/2}^\pi \beta(\psi, \lambda) \sin \psi \, d\psi$$

The preceding discussion has assumed that no inelastic (transpectral) scattering processes are present. However, transpectral scattering does occur in natural waters attributable to fluorescence by dissolved matter or chlorophyll and to Raman or Brillouin scattering by the water molecules themselves (see Sec. 1.23). Power lost from wavelength λ by scattering into wavelength $\lambda' \neq \lambda$ appears in the above formalism as an increase in the spectral absorption.⁹ In this case, $a(\lambda)$ accounts for "true" absorption (e.g., conversion of radiant energy into heat) as well as for the loss of power at wavelength λ by inelastic scattering to another wavelength. The gain in power at λ' appears as a source term in the radiative transfer formalism.

Two more inherent optical properties are commonly used in hydrologic optics. The *spectral single-scattering albedo* $\omega_0(\lambda)$ is

$$\omega_0(\lambda) \equiv \frac{b(\lambda)}{c(\lambda)}$$

The single-scattering albedo is the probability that a photon will be scattered (rather than absorbed) in any given interaction; hence, $\omega_0(\lambda)$ is also known as the *probability of photon survival*. The *spectral volume scattering phase function*, $\tilde{\beta}(\psi, \lambda)$ is defined by

$$\tilde{\beta}(\psi, \lambda) \equiv \frac{\beta(\psi, \lambda)}{b(\lambda)} \quad \text{sr}^{-1}$$

Writing the volume scattering function $\beta(\psi, \lambda)$ as the product of the scattering coefficient $b(\lambda)$ and the phase function $\beta(\psi, \lambda)$ partitions $\beta(\psi, \lambda)$ into a factor giving the *strength* of the scattering, $b(\lambda)$ with units of m^{-1} , and a factor giving the *angular distribution* of the scattered photons, $\beta(\psi, \lambda)$ with units of sr^{-1} .

1.5 APPARENT OPTICAL PROPERTIES

The quantity

$$\bar{\mu}_d(z, \lambda) \equiv \frac{\int_{\Xi_d} L(z, \theta, \phi, \lambda) |\cos \theta| d\Omega}{\int_{\Xi_u} L(z, \theta, \phi, \lambda) d\Omega} \equiv \frac{E_d(z, \lambda)}{E_{0d}(z, \lambda)}$$

is called the *spectral downwelling average cosine*. The definition shows that $\bar{\mu}_d(z, \lambda)$ is the average value of the cosine of the polar angle of all the photons contributing to the downwelling radiance at the given depth and wavelength. The *spectral upwelling average cosine* is defined analogously:

$$\bar{\mu}_u(z, \lambda) \equiv \frac{E_u(z, \lambda)}{E_{0u}(z, \lambda)}$$

The average cosines are useful one-parameter measures of the directional structures of the downwelling and upwelling light fields. For example, if the downwelling light field (radiance distribution) is colimated in direction (θ_0, ϕ_0) so $L(\theta, \phi) = L_0 \delta(\theta - \theta_0) \delta(\phi - \phi_0)$, where δ is the Dirac δ function, then $\bar{\mu}_d = |\cos \theta_0|$. If the downwelling radiance is completely diffuse (isotropic), $L(\theta, \phi) = L_0$ and $\bar{\mu}_d = \frac{1}{2}$. Typical values of the average cosines for waters illuminated by the sun and sky are $\bar{\mu}_d \approx \frac{3}{4}$ and $\bar{\mu}_u \approx \frac{3}{8}$. Older literature generally refers to *distribution functions*, D_d and D_u , rather than to average cosines. The distribution functions are just reciprocals of the average cosines:

$$D_d(z, \lambda) = \frac{1}{\bar{\mu}_d(z, \lambda)} \quad \text{and} \quad D_u(z, \lambda) = \frac{1}{\bar{\mu}_u(z, \lambda)}$$

The *spectral irradiance reflectance* (or *irradiance ratio*) $R(z, \lambda)$ is the ratio of spectral upwelling to downwelling plane irradiances:

$$R(z, \lambda) \equiv \frac{E_u(z, \lambda)}{E_d(z, \lambda)}$$

$R(z, \lambda)$ just beneath the sea surface is of great importance in remote sensing (see Sec. 1.22).

Under typical oceanic conditions for which the incident lighting is provided by the sun and sky, the various radiances and irradiances all decrease approximately exponentially with depth, at least when far enough below the surface (and far enough above the bottom, in shallow water) to be free of boundary effects. For example, it is convenient to write the depth dependence of $E_d(z, \lambda)$ as

$$E_d(z, \lambda) \equiv E_d(0, \lambda) \exp\left[-\int_0^z K_d(z', \lambda) dz'\right] \equiv E_d(0, \lambda) \exp[-\bar{K}_d(\lambda)z]$$

where $K_d(z, \lambda)$ is the *spectral diffuse attenuation coefficient* for spectral downwelling plane irradiance and $\bar{K}_d(\lambda)$ is the average value of $K_d(z, \lambda)$ over the depth interval 0 to z . Solving for $K_d(z, \lambda)$ gives

$$K_d(z, \lambda) = -\frac{d \ln E_d(z, \lambda)}{dz} = -\frac{1}{E_d(z, \lambda)} \frac{dE_d(z, \lambda)}{dz} \quad \text{m}^{-1}$$

The distinction between *beam* and *diffuse* attenuation coefficients is important. The beam attenuation coefficient $c(\lambda)$ is defined in terms of the radiant power lost from a single, narrow, collimated beam of photons. The downwelling diffuse attenuation coefficient $K_d(z, \lambda)$ is defined in terms of the decrease with depth of the ambient downwelling irradiance $E_d(z, \lambda)$, which comprises photons heading in all downward directions (a diffuse or uncollimated light field). $K_d(z, \lambda)$ clearly depends on the directional structure of the ambient light field and so is classified as an apparent optical property. Other diffuse attenuation coefficients, e.g., K_u , K_{0d} , K_{0u} , K_{PAR} , and $K(\theta, \phi)$ are defined in an analogous manner, using the corresponding radiometric quantities.

In *homogeneous* waters, these “ K functions” depend only weakly on depth and therefore can serve as convenient, if imperfect, descriptors of the water body. Smith and Baker¹⁰ have pointed out other reasons why K functions are useful:

1. The K 's are defined as ratios and therefore do not require absolute radiometric measurements.
2. The K 's are strongly correlated with chlorophyll concentration (i.e., they provide a connection between biology and optics).
3. About 90 percent of the diffusely reflected light from a water body comes from a layer of water of depth $1/K_d(0, \lambda)$ (i.e., K_d has implications for remote sensing).
4. Radiative transfer theory provides several useful relations between the K 's and other quantities of interest, such as absorption and beam attenuation coefficients, the irradiance reflectance, and the average cosines.
5. Instruments are available for routine measurement of the K 's.

It must be remembered, however, that in spite of their utility K functions are apparent optical properties—a change in the environment (e.g., solar angle or sea state) changes their value, sometimes by a negligible amount but sometimes greatly. However, numerical simulations by Gordon¹¹ show how with a few additional but easily made measurements measured values of $K_d(z, \lambda)$ and $\bar{K}_d(\lambda)$ can be “normalized” to remove the effects of solar angle and sea state. The normalized K_d and \bar{K}_d are equal to the values that would be obtained if the sun were at the zenith and the sea surface were calm. If this normalization is performed, the resulting $K_d(z, \lambda)$ and $\bar{K}_d(\lambda)$ can be regarded as *inherent* optical properties for all practical purposes. It is strongly recommended that Gordon's procedure be routinely followed by experimentalists.

1.6 THE OPTICALLY SIGNIFICANT CONSTITUENTS OF NATURAL WATERS

Dissolved Substances

Pure sea water consists of pure water plus various dissolved salts, which average about 35 parts per thousand (%) by weight. These salts increase scattering above that of pure water by 30 percent (see Table 10 in Sec. 1.17). It is not well established what, if any, effect these salts have on absorption, but it is likely that they increase absorption somewhat at ultraviolet wavelengths.

Both fresh and saline waters contain varying concentrations of dissolved organic compounds. These compounds are produced during the decay of plant matter and consist mostly of various humic and fulvic acids.⁸ These compounds are generally brown in color and in sufficient concentrations can color the water yellowish brown. For this reason the compounds are generically referred to as *yellow matter*, *Gelbstoff*, or *gilvin*. Yellow matter absorbs very little in the red, but absorption increases rapidly with decreasing wavelength. Since the main source of yellow matter is decayed terrestrial vegetation, concentrations are generally greatest in lakes, rivers, and coastal waters influenced by river runoff. In such waters yellow matter can be the dominant absorber at the blue end of the spectrum. In mid-ocean waters absorption by yellow matter is usually small compared to absorption by other constituents, but some yellow matter is likely to be present as the result of decaying phytoplankton, especially at the end of a bloom.

Particulate Matter

Particulate matter in the oceans has two distinct origins: biological and physical. The organic particles of optical importance are created as bacteria, phytoplankton, and zooplankton grow and reproduce by photosynthesis or by eating their neighbors. Particles of a given size are destroyed by breaking apart after death, by flocculation into larger aggregate particles, or by settling out of the water column. Inorganic particles are created primarily by weathering of terrestrial rocks and soils. These particles can enter the water as wind-blown dust settles on the sea surface, as rivers carry eroded soil to the sea, or as currents resuspend bottom sediments. Inorganic particles are removed from the water by settling, aggregating, or dissolving. This particulate matter usually is the major determiner of both the absorption and scattering properties of natural waters and is responsible for most of the temporal and spatial variability in these optical properties.

Organic Particles These occur in many forms.

Viruses Natural marine waters contain virus particles¹² in concentrations of 10^{12} to 10^{15} particles m^{-3} . These particles are generally much smaller (2–200 nm) than the wavelength of visible light, and it is not known what, if any, direct effect viruses have on the optical properties of sea water.

Colloids Nonliving colloidal particles in the size range 0.4–1.0 μm are found¹³ in typical number concentrations of $10^{13} m^{-3}$ and colloids of size $\leq 0.1 \mu m$ are found¹⁴ in abundances of $10^{15} m^{-3}$. Some of the absorption traditionally attributed to dissolved matter may be due to colloids, some of which strongly resemble fulvic acids in electron micrographs.¹⁴

Bacteria Living bacteria in the size range 0.2–1.0 μm occur in typical number concentrations of 10^{11} – $10^{13} m^{-3}$. It only recently has been recognized^{15–17} that bacteria can be significant scatterers and absorbers of light, especially at blue wavelengths and in clean oceanic waters where the larger phytoplankton are relatively scarce.

Phytoplankton These ubiquitous microscopic plants occur with incredible diversity of species, size, shape, and concentration. They range in cell size from less than 1 μm to more than 200 μm , and some species form even larger chains of individual cells. It has long been recognized that phytoplankton are the particles primarily responsible for determining the optical properties of most oceanic waters. Their chlorophyll and related pigments strongly absorb light in the blue and red and thus when concentrations are high determine the spectral absorption of sea water. These particles are generally much larger than the wavelength of visible light and are efficient scatterers, especially via diffraction, thus influencing the scattering properties of sea water.

Organic detritus Nonliving organic particles of various sizes are produced, for example, when phytoplankton die and their cells break apart. They may also be formed when zooplankton graze on phytoplankton and leave behind cell fragments and fecal pellets. Even if these detrital particles contain pigments at the time of their production, they can be rapidly photo-oxidized and lose the characteristic absorption spectrum of living phytoplankton, leaving significant absorption only at blue wavelengths.

Large particles Particles larger than 100 μm include zooplankton (living animals with sizes from tens of micrometers to two centimeters) and fragile amorphous aggregates¹⁸ of smaller particles (“marine snow,” with sizes from 0.5 mm to tens of centimeters). Such particles occur in highly variable numbers from almost none to thousands per cubic meter. Even at relatively large concentrations these large particles tend to be missed by optical instruments that randomly sample only a few cubic centimeters of water or that mechanically break apart the aggregates. However, these large particles can be efficient diffuse scatterers of light and therefore may significantly affect the optical properties (especially backscatter) of large volumes of water, e.g., as seen by remote sensing instruments. Although such optical effects are recognized, they have not been quantified.

Inorganic Particles These generally consist of finely ground quartz sand, clay minerals, or metal oxides in the size range from much less than 1 μm to several tens of micrometers. Insufficient attention has been paid to the optical effects of such particles in sea water, although it is recognized that inorganic particles are sometimes optically more important than organic particles. Such situations can occur both in turbid coastal waters carrying a heavy sediment load and in very clear oceanic waters which are receiving wind-blown dust.¹⁹

At certain stages of its life, the phytoplankton coccolithophore species *Emiliania huxleyi* is a most remarkable source of crystalline particles. During blooms *E. huxleyi* produces and sheds enormous numbers of small (2–4 μm) calcite plates; concentrations of 3×10^{11} plates m^{-3} have been observed.²⁰ Although they have a negligible effect on light absorption, these calcite plates are extremely efficient light scatterers: irradiance reflectances of $R = 0.39$ have been observed²⁰ at blue wavelengths during blooms (compared with $R = 0.02$ to 0.05 in the blue for typical ocean waters, discussed in Sec. 1.22). Such coccolithophore blooms give the ocean a milky white or turquoise appearance.

1.7 PARTICLE SIZE DISTRIBUTIONS

In spite of the diverse mechanisms for particle production and removal, observation shows that a single family of particle size distributions often suffices to describe oceanic particulate matter in the optically important size range from 0.1 to 100 μm . Let $N(x)$ be the number of particles per unit volume with size greater than x in a sample of particles; x usually represents equivalent spherical diameter computed from particle volume, but also can represent particle volume or surface area. The Junge (also called *hyperbolic*) cumulative size distribution²¹ is then

$$N(x) = k \left(\frac{x}{x_0} \right)^{-m}$$

where k sets the scale, x_0 is a reference size, and $-m$ is the slope of the distribution when $\log N$ is plotted versus $\log x$; k , x_0 , and m are positive constants.

Oceanic particle size distributions usually have m values between 2 and 5, with $m = 3$ to 4 being typical; such spectra can be seen in McCave,²² Fig. 7. It often occurs that oceanic particle size spectra are best described by a segmented distribution in which a smaller value of m is used for x less than a certain value and a larger value of m is used for x greater than that value. Such segmented spectra can be seen in Bader,²¹ and in McCave,²² Fig. 8.

The quantity most relevant to optics, e.g., in Mie scattering computations for polydisperse systems, is not the *cumulative* size distribution $N(x)$, but rather the *number* size distribution $n(x)$. The number distribution is defined such that $n(x) dx$ is the number of particles in the size interval from x to $x + dx$. The number distribution is related to the cumulative distribution by $n(x) = |dN(x)/dx|$, so that for the Junge distribution

$$n(x) = kmx_0^{-m} x^{-m-1} \equiv Kx^{-s}$$

where $K \equiv kmx_0^{-m}$ and $s \equiv m + 1$; s is commonly referred to as the slope of the distribution. Figure 3 shows the number distribution of biological particles typical of open ocean waters; note that a value of $s = 4$ gives a reasonable fit to the plotted points.

It should be noted, however, that the Junge distribution sometimes fails to represent oceanic conditions. For example, during the growth phase of a phytoplankton bloom the rapid increase in population of a particular species may give abnormally large numbers of particles in a particular size range. Such bloom conditions therefore give a “bump” in $n(x)$ that is not well modeled by the simple Junge distribution. Moreover, Lambert et al.²³ found that a log-normal distribution sometimes better described the distributions of inorganic particles found in water samples taken from near the bottom at deep ocean locations. These particles were principally aluminosilicates in the 0.2- to 10.0- μm size range but included quartz grains, metal oxides, and phytoplankton skeletal parts such as coccolithophore plates. Based on the sampling location it was assumed that the inorganic particles were resuspended sediments. Lambert et al. found that the size distributions of the individual particle types (e.g., aluminosilicates or metal oxides) obeyed log-normal distributions which “flattened out” below 1 μm . For particles larger than $\sim 1 \mu\text{m}$, log-normal and Junge distributions gave nearly

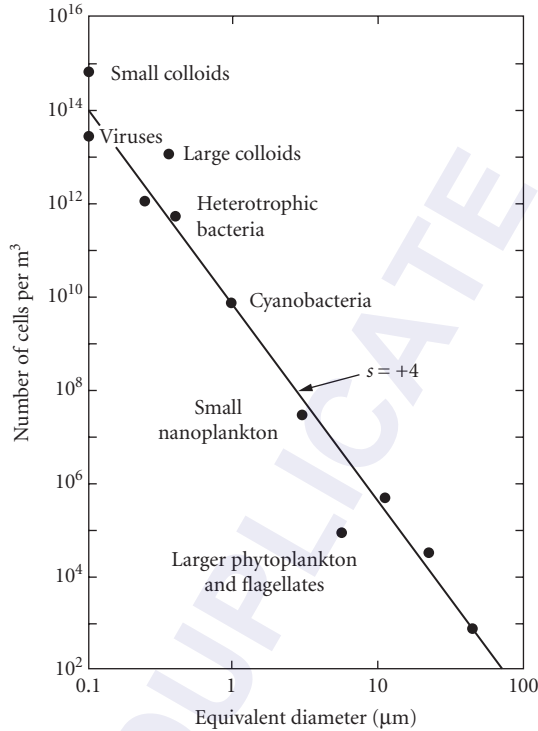


FIGURE 3 Number size distribution typical of biological particles in the open ocean. (Based on Stramski and Kiefer,¹⁷ with permission.)

equivalent descriptions of the data. Biological particles were not as well described by the log-normal distribution, especially for sizes greater than 5 μm.

1.8 ELECTROMAGNETIC PROPERTIES OF WATER

In studies of electromagnetic wave propagation at the level of Maxwell's equations it is convenient to specify the bulk electromagnetic properties of the medium via the electrical permittivity ϵ , the magnetic permeability μ , and the electrical conductivity σ . Since water displays no significant magnetic properties, the permeability can be taken equal to the free-space (in vacuo) value at all frequencies: $\mu = \mu_0 = 4\pi \times 10^{-7} \text{ N A}^{-2}$. Both ϵ and σ depend on the frequency ω of the propagating electromagnetic wave as well as on the water temperature, pressure, and salinity. Low-frequency ($\omega \rightarrow 0$) values for the permittivity are of order $\epsilon \approx 80\epsilon_0$, where $\epsilon_0 = 8.85 \times 10^{-12} \text{ A}^2 \text{ s}^2 \text{ N}^{-1} \text{ m}^{-2}$ is the free-space value. This value decreases to $\epsilon \approx 1.8\epsilon_0$ at optical frequencies. Extensive tabulations of ϵ/ϵ_0 as a function of temperature and pressure are given for pure water in Archer and Wang.²⁴ The low-frequency conductivity ranges from $\sigma \approx 4 \times 10^{-6} \text{ siemen m}^{-1}$ for pure water to $\sigma \approx 4.4 \text{ siemen m}^{-1}$ for sea water.

The effects of ϵ , μ , and σ on electromagnetic wave propagation are compactly summarized in terms of the complex index of refraction, $m = n - ik$, where n is the real index of refraction,

k is the dimensionless electrodynamic absorption coefficient, and $i = \sqrt{-1}$; n and k are collectively called the *optical constants* of water (a time dependence convention of $\exp(+i\omega t)$ is used in deriving wave equations from Maxwell's equations). The explicit dependence of m on ϵ , μ , and σ is given by²⁵

$$\begin{aligned} m^2 &= \mu\epsilon c^2 - i \frac{\mu\sigma c^2}{\omega} \\ &= (n - ik)^2 = n^2 - k^2 - i2nk \end{aligned}$$

where $c = (\epsilon_0\mu_0)^{-1/2}$ is the speed of light in vacuo. These equations can be used to relate n and k to the bulk electromagnetic properties. The optical constants are convenient because they are directly related to the scattering and absorbing properties of water. The real index of refraction $n(\lambda)$ governs scattering both at interfaces (via the laws of reflection and refraction) and within the medium (via thermal or other fluctuations of $n(\lambda)$ at molecular and larger scales). The spectral absorption coefficient $a(\lambda)$ is related to $k(\lambda)$ by²⁵

$$a(\lambda) = \frac{4\pi k(\lambda)}{\lambda}$$

Here λ refers to the in vacuo wavelength of light corresponding to a given frequency ω of electromagnetic wave.

Figure 4 shows the wavelength dependence of the optical constants n and k for pure water. The extraordinary feature seen in this figure is the narrow “window” in $k(\lambda)$, where $k(\lambda)$ decreases by over nine orders of magnitude between the near ultraviolet and the visible and then quickly rises again in the near infrared. This behavior in $k(\lambda)$ gives a corresponding window in the spectral absorption coefficient $a(\lambda)$ as seen in Table 2. Because of the opaqueness of water outside the near-UV to near-IR wavelengths, hydrologic optics is concerned only with this small part of the electromagnetic spectrum. These wavelengths overlap nicely with the wavelengths of the sun's maximum energy output and with a corresponding window in atmospheric absorption, much to the benefit of life on earth.

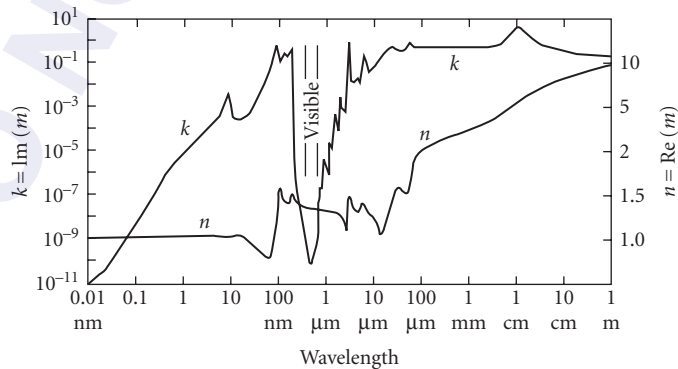


FIGURE 4 The optical constants of pure water. The left axis gives $k = \text{Im}(m)$ and the right axis gives $n = \text{Re}(m)$ where m is the complex index of refraction. (Redrawn from Zoloratev and Demin,²⁶ with permission.)

TABLE 2 Absorption Coefficient a of Pure Water As a Function of Wavelength λ ^a

λ	a (m ⁻¹)	λ	a (m ⁻¹)
0.01 nm	1.3×10^1	700 nm	0.650
0.1	6.5×10^2	800	2.07
1	9.4×10^4	900 nm	7.0
10	3.5×10^6	1 μ m	3.3×10^1
100	5.0×10^7	10	7.0×10^4
200	3.07	100 μ m	6.5×10^4
300	0.141	0.001 m	1.3×10^4
400	0.0171	0.01	3.6×10^3
500	0.0257	0.1	5.0×10^1
600 nm	0.244	1 m	2.5

^aData for $200 \text{ nm} \leq \lambda \leq 800 \text{ nm}$ taken from Table 6. Data for other wavelengths computed from Fig. 4.

1.9 INDEX OF REFRACTION

Seawater

Austin and Halikas²⁷ exhaustively reviewed the literature on measurements of the real index of refraction of sea water. Their report contains extensive tables and interpolation algorithms for the index of refraction (relative to air), $n(\lambda, S, T, p)$, as a function of wavelength ($\lambda = 400$ to 700 nm), salinity ($S = 0$ to 43‰), temperature ($T = 0$ to 30°C), and pressure ($p = 10^5$ to 10^8 Pa , or 1 to 1080 atm). Figure 5 illustrates the general dependence of n on these four parameters: n decreases with increasing wavelength or temperature and increases with increasing salinity or pressure. Table 3 gives the values of n for the extreme values of each parameter. The extreme values of n , 1.329128 and 1.366885, show that n varies by less than 3 percent over the entire parameter range relevant to hydrologic optics. Table 4 gives selected values of $n(\lambda, T)$ for fresh water ($S = 0$) and for typical sea water ($S = 35\text{‰}$) at atmospheric pressure ($p = 10^5 \text{ Pa}$). The values in Table 4 can be multiplied by 1.000293 (the index of refraction of dry air at STP and $\lambda = 538 \text{ nm}$) if values relative to vacuum are

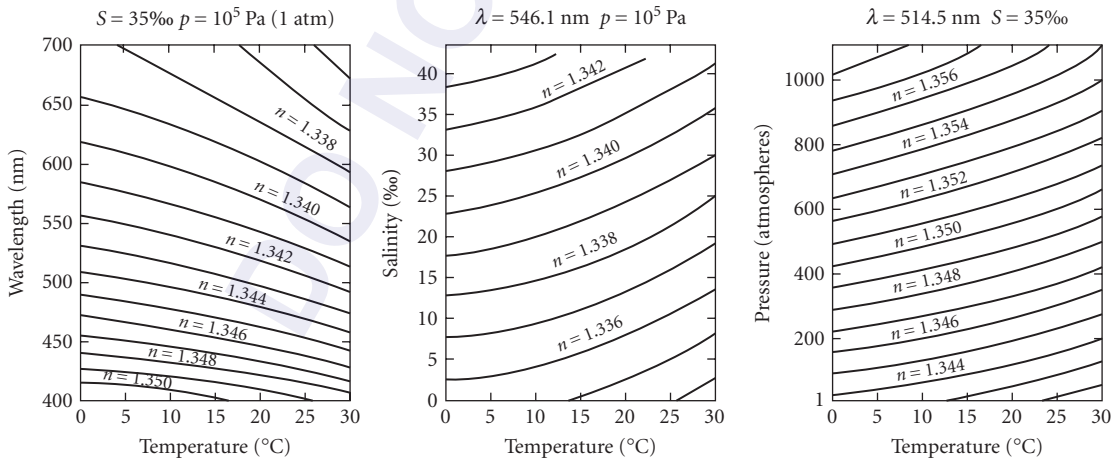


FIGURE 5 Real index of refraction of water for selected values of pressure, temperature, and salinity. (Adapted from Austin and Halikas.²⁷)

TABLE 3 Index of Refraction of Water n for the Extreme Values of Pressure p , Temperature T , Salinity S , and Wavelength λ Encountered in Hydrologic Optics*

p (Pa)	T (°C)	S (‰)	λ (nm)	n
1.01×10^5	0	0	400	1.344186
1.01	0	0	700	1.331084
1.01	0	35	400	1.351415
1.01	0	35	700	1.337906
1.01	30	0	400	1.342081
1.01	30	0	700	1.329128
1.01	30	35	400	1.348752
1.01	30	35	700	1.335316
1.08×10^8	0	0	400	1.360076
1.08	0	0	700	1.346604
1.08	0	35	400	1.366885
1.08	0	35	700	1.352956
1.08	30	0	400	1.356281
1.08	30	0	700	1.342958
1.08	30	35	400	1.362842
1.08	30	35	700	1.348986

*Reproduced from Austin and Halikas.²⁷**TABLE 4** Index of Refraction of Fresh Water and of Sea Water at Atmospheric Pressure for Selected Temperatures and Wavelengths*

Temp (°C)	Fresh water ($S = 0$) wavelength (nm)							
	400	420	440	460	480	500	520	540
0	1.34419	1.34243	1.34092	1.33960	1.33844	1.33741	1.33649	1.33567
10	1.34390	1.34215	1.34064	1.33933	1.33817	1.33714	1.33623	1.33541
20	1.34317	1.34142	1.33992	1.33860	1.33745	1.33643	1.33551	1.33469
30	1.34208	1.34034	1.33884	1.33753	1.33638	1.33537	1.33445	1.33363
Temp (°C)	Wavelength (nm)							
	560	580	600	620	640	660	680	700
0	1.33494	1.33424	1.33362	1.33305	1.33251	1.33200	1.33153	1.33108
10	1.33466	1.33399	1.33336	1.33279	1.33225	1.33174	1.33127	1.33084
20	1.33397	1.33328	1.33267	1.33210	1.33156	1.33105	1.33059	1.33016
30	1.33292	1.33223	1.33162	1.33106	1.33052	1.33001	1.32955	1.32913
Temp (°C)	Sea water ($S = 35\text{‰}$) wavelength (nm)							
	400	420	440	460	480	500	520	540
0	1.35141	1.34961	1.34804	1.34667	1.34548	1.34442	1.34347	1.34263
10	1.35084	1.34903	1.34747	1.34612	1.34492	1.34385	1.34291	1.34207
20	1.34994	1.34814	1.34657	1.34519	1.34401	1.34295	1.34200	1.34116
30	1.34875	1.34694	1.34539	1.34404	1.34284	1.34179	1.34085	1.34000
Temp (°C)	Wavelength (nm)							
	560	580	600	620	640	660	680	700
0	1.34186	1.34115	1.34050	1.33992	1.33937	1.33885	1.33836	1.33791
10	1.34129	1.34061	1.33997	1.33938	1.33882	1.33830	1.33782	1.33738
20	1.34039	1.33969	1.33904	1.33845	1.33791	1.33739	1.33690	1.33644
30	1.33925	1.33855	1.33790	1.33731	1.33676	1.33624	1.33576	1.33532

*Data extracted from Austin and Halikas.

TABLE 5 Index of Refraction Relative to Water, n , of Inorganic Particles Found in Sea Water

Substance	n
Quartz	1.16
Kaolinite	1.17
Montmorillonite	1.14
Hydrated mica	1.19
Calcite	1.11/1.24

desired. Millard and Seaver²⁸ have developed a 27-term formula that gives the index of refraction to part-per-million accuracy over most of the oceanographic parameter range.

Particles

Suspended particulate matter in sea water often has a bimodal index of refraction distribution. Living phytoplankton typically have “low” indices of refraction in the range 1.01 to 1.09 relative to the index of refraction of seawater. Detritus and inorganic particles generally have “high” indices in the range of 1.15 to 1.20 relative to seawater.²⁹ Typical values are 1.05 for phytoplankton and 1.16 for inorganic particles.

Table 5 gives the relative index of refraction of terrigenous minerals commonly found in river runoff and wind-blown dust. Only recently has it become possible to measure the refractive indices of individual phytoplankton cells.³⁰ Consequently, little is yet known about the dependence of refractive index on phytoplankton species, or on the physiological state of the plankton within a given species, although it appears that the dependence can be significant.³¹

1.10 MEASUREMENT OF ABSORPTION

Determination of the spectral absorption coefficient $a(\lambda)$ for natural waters is a difficult task for several reasons. First, water absorbs only weakly at near-UV to blue wavelengths so that very sensitive instruments are required. More importantly, scattering is never negligible so that careful consideration must be made of the possible aliasing of the absorption measurements by scattering effects. In pure water at wavelengths of $\lambda = 370$ to 450 nm, molecular scattering provides 20 to 25 percent (Table 10) of the total beam attenuation, $c(\lambda) = a(\lambda) + b(\lambda)$. Scattering effects can dominate absorption at all visible wavelengths in waters with high particulate loads. Additional complications arise in determining the absorption of pure water because of the difficulty of preparing uncontaminated samples.

Many techniques have been employed in attempts to determine the spectral absorption coefficient for pure water, $a_w(\lambda)$; these are reviewed in Smith and Baker.³² The most commonly employed technique for routine determination of $a(\lambda)$ for oceanic waters consists of filtering a sample of water to retain the particulate matter on a filter pad. The spectral absorption of the particulate matter, $a_p(\lambda)$, is then determined in a spectrophotometer. The absorption of pure water, $a_w(\lambda)$, must be added to $a_p(\lambda)$ to obtain the total absorption of the oceanic water sample. Even though this technique for determining absorption has been in use for many years, the methodology is still evolving^{33–35} because of the many types of errors inherent in the $a_p(\lambda)$ measurements (e.g., inability of filters to retain all particulates, scattering effects within the sample cell, absorption by dissolved matter retained on the filter pad, and decomposition of pigments during the filtration process).

Moreover, this methodology for determining total absorption assumes that absorption by dissolved organic matter (yellow substances) is negligible, which is not always the case. If the absorption by yellow matter, $a_y(\lambda)$, is desired, then the absorption of the *filtrate* is measured, and $a_y(\lambda)$ is taken to be $a_{\text{filtrate}}(\lambda) - a_w(\lambda)$. Several novel instruments under development³⁶⁻³⁸ show promise for circumventing the problems inherent in the filter-pad technique as well as for making in situ measurements of total absorption which at present is difficult.³⁹

1.11 ABSORPTION BY PURE SEA WATER

Table 2 showed the absorption for pure water over the wavelength range from 0.01 nm (x-rays) to 1 m (radio waves). As is seen in the table, only the near-UV to near-IR wavelengths are of interest in hydrologic optics. Smith and Baker³² made a careful but indirect determination of the *upper bound* of the spectral absorption coefficient of pure sea water, $a_w(\lambda)$, in the wavelength range of oceanographic interest, $200 \text{ nm} \leq \lambda \leq 800 \text{ nm}$. Their work assumed that for the clearest natural waters (1) absorption by salt or other dissolved substances was negligible, (2) the only scattering was by water molecules and salt ions, and (3) there was no inelastic scattering (i.e., no fluorescence or Raman scattering). With these assumptions the inequality (derived from radiative transfer theory)

$$a_w(\lambda) \leq K_d(\lambda) - \frac{1}{2}b_m^{\text{sw}}(\lambda)$$

holds. Here $b_m^{\text{sw}}(\lambda)$ is the spectral scattering coefficient for pure sea water; $b_m^{\text{sw}}(\lambda)$ was taken as known (Table 10). Smith and Baker then used measured values of the diffuse attenuation function $K_d(\lambda)$ from very clear waters (e.g., Crater Lake, Oregon, U.S.A., and the Sargasso Sea) to estimate $a_w(\lambda)$. Table 6 gives their self-consistent values of $a_w(\lambda)$, $K_d(\lambda)$, $b_m^{\text{sw}}(\lambda)$.

The Smith and Baker absorption values are widely used. However, it must be remembered that the values of $a_w(\lambda)$ in Table 6 are upper bounds; the true absorption of pure water is likely to be somewhat lower, at least at violet and blue wavelengths.⁴⁰ Smith and Baker pointed out that there are uncertainties because K_d , an apparent optical property, is influenced by environmental conditions. They also commented that at wavelengths below 300 nm, their values are “merely an educated guess.” They estimated the accuracy of $a_w(\lambda)$ to be within +25 and -5 percent between 300 and 480 nm and +10 to -15 percent between 480 and 800 nm. Numerical simulations by Gordon¹¹ indicate that a more restrictive inequality,

$$a_w(\lambda) \leq \frac{K_d(\lambda)}{D_0(\lambda)} - 0.62b_m^{\text{sw}}(\lambda)$$

could be used. Here $D_0(\lambda)$ is a measurable distribution function [$D_0(\lambda) > 1$] that corrects for the effects of sun angle and sea state on $K_d(\lambda)$ (discussed earlier). Use of the Gordon inequality could reduce the Smith and Baker absorption values by up to 20 percent at blue wavelengths. And finally, the Smith and Baker measurements were not made in optically pure water but rather in the “clearest natural waters.” Even these waters contain a small amount of dissolved and particulate matter which will contribute something to both absorption and scattering.

There is evidence⁴¹ that absorption is weakly dependent on temperature, at least in the red and near infrared ($\partial a/\partial T \sim 0.0015 \text{ m}^{-1} \text{ }^\circ\text{C}^{-1}$ at $\lambda = 600 \text{ nm}$ and $\partial a/\partial T \sim 0.01 \text{ m}^{-1} \text{ }^\circ\text{C}^{-1}$ at $\lambda = 750 \text{ nm}$) and perhaps also slightly dependent on salinity; these matters are under investigation.

TABLE 6 Spectral Absorption Coefficient of Pure Sea Water, a_w , As Determined by Smith and Baker (Values of the molecular scattering coefficient of pure sea water, b_m^{sw} , and of the diffuse attenuation coefficient K_d used in their computation of a_w are also shown.)*

λ (nm)	a_w (m^{-1})	b_m^{sw} (m^{-1})	K_d (m^{-1})	λ (nm)	a_w (m^{-1})	b_m^{sw} (m^{-1})	K_d (m^{-1})
200	3.07	0.151	3.14	500	0.0257	0.0029	0.0271
210	1.99	0.119	2.05	510	0.0357	0.0026	0.0370
220	1.31	0.0995	1.36	520	0.0477	0.0024	0.0489
230	0.927	0.0820	0.968	530	0.0507	0.0022	0.0519
240	0.720	0.0685	0.754	540	0.0558	0.0021	0.0568
250	0.559	0.0575	0.588	550	0.0638	0.0019	0.0648
260	0.457	0.0485	0.481	560	0.0708	0.0018	0.0717
270	0.373	0.0415	0.394	570	0.0799	0.0017	0.0807
280	0.288	0.0353	0.306	580	0.108	0.0016	0.109
290	0.215	0.0305	0.230	590	0.157	0.0015	0.158
300	0.141	0.0262	0.154	600	0.244	0.0014	0.245
310	0.105	0.0229	0.116	610	0.289	0.0013	0.290
320	0.0844	0.0200	0.0944	620	0.309	0.0012	0.310
330	0.0678	0.0175	0.0765	630	0.319	0.0011	0.320
340	0.0561	0.0153	0.0637	640	0.329	0.0010	0.330
350	0.0463	0.0134	0.0530	650	0.349	0.0010	0.350
360	0.0379	0.0120	0.0439	660	0.400	0.0008	0.400
370	0.0300	0.0106	0.0353	670	0.430	0.0008	0.430
380	0.0220	0.0094	0.0267	680	0.450	0.0007	0.450
390	0.0191	0.0084	0.0233	690	0.500	0.0007	0.500
400	0.0171	0.0076	0.0209	700	0.650	0.0007	0.650
410	0.0162	0.0068	0.0196	710	0.839	0.0007	0.834
420	0.0153	0.0061	0.0184	720	1.169	0.0006	1.170
430	0.0144	0.0055	0.0172	730	1.799	0.0006	1.800
440	0.0145	0.0049	0.0170	740	2.38	0.0006	2.380
450	0.0145	0.0045	0.0168	750	2.47	0.0005	2.47
460	0.0156	0.0041	0.0176	760	2.55	0.0005	2.55
470	0.0156	0.0037	0.0175	770	2.51	0.0005	2.51
480	0.0176	0.0034	0.0194	780	2.36	0.0004	2.36
490	0.0196	0.0031	0.0212	790	2.16	0.0004	2.16
				800	2.07	0.0004	2.07

*Reproduced from Smith and Baker,³² with permission.

1.12 ABSORPTION BY DISSOLVED ORGANIC MATTER

Absorption by yellow matter is reasonably well described by the model⁴²

$$a_y(\lambda) = a_y(\lambda_0) \exp[-0.014(\lambda - \lambda_0)]$$

over the range $350 \text{ nm} \leq \lambda \leq 700 \text{ nm}$. Here λ_0 is a reference wavelength usually chosen to be $\lambda_0 = 440 \text{ nm}$ and $a_y(\lambda_0)$ is the absorption due to yellow matter at the reference wavelength. The value of $a_y(\lambda)$ of course depends on the concentration of yellow matter in the water. The exponential decay constant depends on the relative proportion of specific types of yellow matter; other studies have found exponents of -0.014 to -0.019 (Roesler et al.,⁴³ table 1). Both total concentration and proportions are highly variable. Table 7 gives measured values of $a_y(440)$ for selected waters. Because of the variability in yellow matter concentrations, the values found in Table 7 have little general

TABLE 7 Measured Absorption Coefficient at $\lambda = 440$ nm Due to Yellow Matter, $a_y(440 \text{ nm})$, for Selected Waters*

Water Body	$a_y(440 \text{ nm})$ (m^{-1})
Oceanic waters	
Sargasso Sea	≈ 0
off Bermuda	0.01
Gulf of Guinea	0.024–0.113
oligotrophic Indian Ocean	0.02
mesotrophic Indian Ocean	0.03
eutrophic Indian Ocean	0.09
Coastal and estuarine waters	
North Sea	0.07
Baltic Sea	0.24
Rhone River mouth, France	0.086–0.572
Clyde River estuary, Australia	0.64
Lakes and rivers	
Crystal Lake, Wisconsin, U.S.A.	0.16
Lake George, Australia	0.69–3.04
Lake George, Uganda	3.7
Carrao River, Venezuela	12.44
Lough Napeast, Ireland	19.1

*Condensed from Kirk,⁸ with permission.

validity even for the particular water bodies sampled, but they do serve to show representative values and the range of influence of yellow matter in determining the total absorption. Although the above model allows the determination of spectral absorption by yellow matter if the absorption is known at one wavelength, no model yet exists that allows for the direct determination of $a_y(\lambda)$ from given concentrations of yellow matter constituents.

1.13 ABSORPTION BY PHYTOPLANKTON

Phytoplankton cells are strong absorbers of visible light and therefore play a major role in determining the absorption properties of natural waters. Absorption by phytoplankton occurs in various photosynthetic pigments of which the chlorophylls are best known to nonspecialists. Absorption by chlorophyll itself is characterized by strong absorption bands in the blue and in the red (peaking at $\lambda \approx 430$ and 665 nm, respectively, for chlorophyll *a*), with very little absorption in the green. Chlorophyll occurs in all plants, and its concentration in milligrams of chlorophyll per cubic meter of water is commonly used as the relevant optical measure of phytoplankton abundance. (The term “chlorophyll concentration” usually refers to the sum of chlorophyll *a*, the main pigment in phytoplankton cells, and the related pigment pheophytin *a*.) Chlorophyll concentrations for various waters range from 0.01 mg m^{-3} in the clearest open ocean waters to 10 mg m^{-3} in productive coastal upwelling regions to 100 mg m^{-3} in eutrophic estuaries or lakes. The globally averaged, near-surface, open-ocean value is in the neighborhood of 0.5 mg m^{-3} .

The absorbing pigments are not evenly distributed within phytoplankton cells but are localized into small “packages” (chloroplasts) which are distributed nonrandomly throughout the cell. This localized distribution of pigments means⁸ that the spectral absorption by a phytoplankton cell or by a collection of cells in water is “flatter” (has less-pronounced peaks and reduced overall absorption) than if the pigments were uniformly distributed throughout the cell or throughout the water. This so-called “pigment packaging effect” is a major source of both inter- and intraspecies variability in spectral absorption by phytoplankton. This is because the details of the pigment

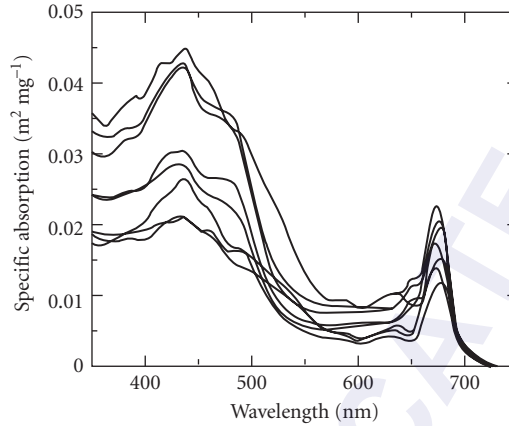


FIGURE 6 Chlorophyll-specific spectral absorption coefficients for eight species of phytoplankton. (Redrawn from Sathyendranath et al.,⁴⁴ with permission.)

packaging within cells depend not only on species but also on a cell's size and physiological state (which in turn depends on environmental factors such as ambient lighting and nutrient availability). Another source of variability in addition to chlorophyll *a* concentration and packaging is changes in pigment composition (the relative proportions of accessory pigments, namely, chlorophylls *b* and *c*, pheopigments, biliproteins, and carotenoids) since each pigment displays a characteristic absorption curve.

A qualitative feel for the nature of phytoplankton absorption can be obtained from Fig. 6 which is based on absorption measurements from eight different single-species laboratory phytoplankton cultures.⁴⁴ Measured spectral absorption coefficients for the eight cultures, $a_i(\lambda)$, $i = 1$ to 8, were first reduced by subtracting $a_i(737)$ to remove the effects of absorption by detritus and cell constituents other than pigments: the assumption is that pigments do not absorb at $\lambda = 737$ nm and that the residual absorption is wavelength independent (which is a crude approximation). The resulting curves were then normalized by the chlorophyll concentrations of the respective cultures to generate the *chlorophyll-specific spectral absorption* curves for phytoplankton, $a_i^*(\lambda)$.

$$a_i^*(\lambda) = \frac{a_i(\lambda) - a_i(737)}{C_i} \quad \frac{\text{m}^{-1}}{\text{mg m}^{-3}} = \text{m}^2 \text{mg}^{-1}$$

which are plotted in Fig. 6.

Several general features of phytoplankton absorption are seen in Fig. 6:

1. There are distinct absorption peaks at $\lambda \approx 440$ and 675 nm.
2. The blue peak is one to three times as high as the red one (for a given species) due to the contribution of accessory pigments to absorption in the blue.
3. There is relatively little absorption between 550 and 650 nm, with the absorption minimum near 600 nm being 10 to 30 percent of the value at 440 nm.

Similar analysis by Morel⁴⁵ yielded the average specific absorption curve shown in Fig. 7. Morel's curve is an average of spectra from 14 cultured phytoplankton species. The Morel curve is qualitatively the same as the curves of Fig. 6 and is as good a candidate as any for being called a "typical" phytoplankton specific absorption curve. The $a^*(\lambda)$ values of Fig. 7 are tabulated in Table 8 for reference.

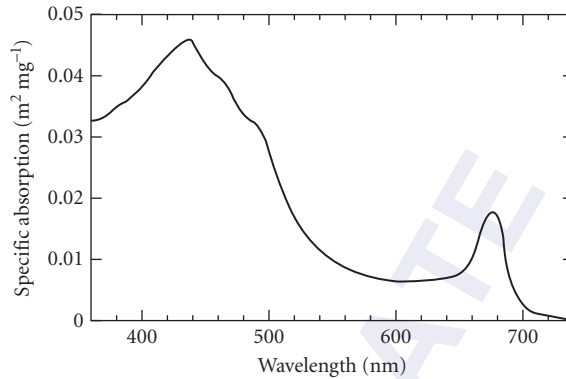


FIGURE 7 Average chlorophyll-specific spectral absorption coefficient for 14 species of phytoplankton. (Redrawn from Morel,⁴⁵ with permission.)

1.14 ABSORPTION BY ORGANIC DETRITUS

Only recently has it become possible to determine the relative contributions of living phytoplankton and nonliving detritus to the total absorption by particulates. Iturriaga and Siegel⁴⁶ used micro-spectrophotometric techniques capable of measuring the spectral absorption of individual particles as small as 3 μm diameter to examine the absorption properties of particulates from Sargasso sea waters. Roesler et al.⁴³ employed a standard filter-pad technique with measurements made before

TABLE 8 Average Chlorophyll-Specific Spectral Absorption Coefficient a^* for 14 Species of Phytoplankton As Plotted in Fig. 7 (The standard deviation is $\sim 30\%$ of the mean except in the vicinity of 400 nm, where it is $\sim 50\%$.)

λ (nm)	a^* ($\text{m}^2 \text{mg}^{-1}$)	λ (nm)	a^* ($\text{m}^2 \text{mg}^{-1}$)	λ (nm)	a^* ($\text{m}^2 \text{mg}^{-1}$)
400	0.0394	500	0.0274	600	0.0053
405	0.0395	505	0.0246	605	0.0053
410	0.0403	510	0.0216	610	0.0054
415	0.0417	515	0.0190	615	0.0057
420	0.0429	520	0.0168	620	0.0059
425	0.0439	525	0.0151	625	0.0061
430	0.0448	530	0.0137	630	0.0063
435	0.0452	535	0.0125	635	0.0064
440	0.0448	540	0.0115	640	0.0064
445	0.0436	545	0.0106	645	0.0066
450	0.0419	550	0.0098	650	0.0071
455	0.0405	555	0.0090	655	0.0084
460	0.0392	560	0.0084	660	0.0106
465	0.0379	565	0.0078	665	0.0136
470	0.0363	570	0.0073	670	0.0161
475	0.0347	575	0.0068	675	0.0170
480	0.0333	580	0.0064	680	0.0154
485	0.0322	585	0.0061	685	0.0118
490	0.0312	590	0.0058	690	0.0077
495	0.0297	595	0.0055	695	0.0046
				700	0.0027

*Data courtesy of A. Morel.⁴⁵

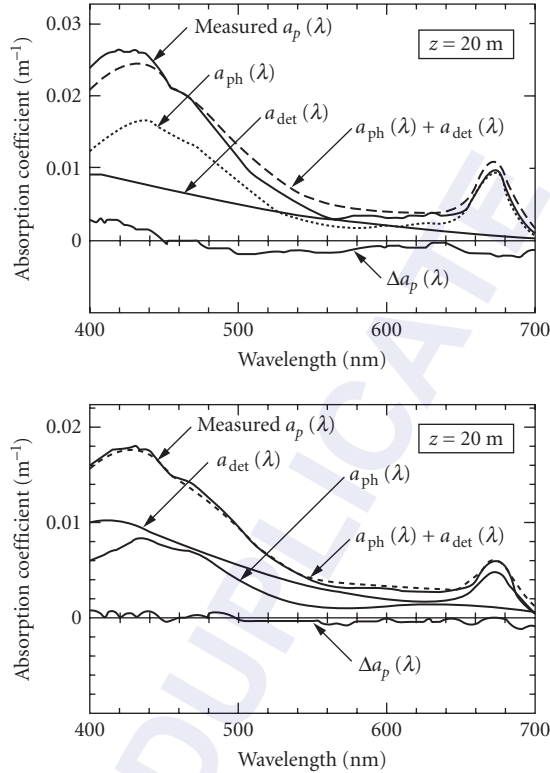


FIGURE 8 Examples of the relative contributions of absorption by phytoplankton $a_{\text{ph}}(\lambda)$, and by organic detritus $a_{\text{det}}(\lambda)$, to the total particulate absorption $a_p(\lambda)$, from Sargasso Sea waters. (Redrawn from Iturriaga and Siegel,⁴⁶ with permission.)

and after pigments were chemically extracted to distinguish between absorption by pigmented and nonpigmented particles from fjord waters in the San Juan Islands, Washington, U.S.A. Each of these dissimilar techniques applied to particles from greatly different waters found very similar absorption spectra for nonpigmented organic particles derived from phytoplankton.

Figure 8 shows the microspectrophotometrically determined contributions of absorption by phytoplankton, $a_{\text{ph}}(\lambda)$, and of absorption by detritus, $a_{\text{det}}(\lambda)$, to the independently measured (by the filter-pad technique) total particulate absorption, $a_p(\lambda)$, for two depths at the same Atlantic location. The small residual, $\Delta a_p(\lambda) = a_p(\lambda) - a_{\text{ph}}(\lambda) - a_{\text{det}}(\lambda)$ shown in the figure is attributed either to errors in the determination of the phytoplankton and detrital parts (particles smaller than $\sim 3 \mu\text{m}$ were not analyzed) or to contamination by dissolved organic matter of the filter-pad measurements of total particulate absorption. Note that at the shallow depth the phytoplankton are relatively more important at blue wavelengths whereas the detritus is slightly more important at the deeper depth. There is no generality in this result (other locations showed the reverse)—it merely illustrates the variability possible in water samples taken only 60 vertical meters apart.

The important feature to note in Fig. 8 is the general shape of the spectral absorption curve for detritus. Roesler et al. found essentially identical curves in their determination of $a_{\text{det}}(\lambda)$. These curves are reminiscent of the absorption curves for yellow matter and, indeed, Roesler et al. found that the model

$$a_{\text{det}}(\lambda) = a_{\text{det}}(400) \exp[-0.011(\lambda - 400)]$$

provides a satisfactory fit to detrital absorption curves. Other studies have found coefficients of -0.006 to -0.014 (Roesler et al.,⁴³ table 1) instead of -0.011 .

1.15 BIO-OPTICAL MODELS FOR ABSORPTION

Depending on the concentrations of dissolved substances, phytoplankton, and detritus, the total spectral absorption coefficient of a given water sample can range from almost identical to that of pure water to one which shows orders-of-magnitude greater absorption than pure water, especially at blue wavelengths. Figure 9 shows some $a(\lambda)$ profiles from various natural waters. Figure 9a shows absorption profiles measured in phytoplankton-dominated waters where chlorophyll concentrations ranged from $C = 0.2$ to 18.4 mg m^{-3} . In essence, the absorption is high in the blue because of absorption by phytoplankton pigments and high in the red because of absorption by the water. Figure 9b shows the absorption at three locations where $C \approx 2 \text{ mg m}^{-3}$ but where the scattering coefficient b varied from 1.55 to 3.6 m^{-1} indicating that nonpigmented particles were playing an important role in determining the shape of $a(\lambda)$. Figure 9c shows curves from waters rich in yellow matter. One of the goals of bio-optics is to develop predictive models for absorption curves such as those seen in Fig. 9.

Case 1 waters are waters in which the concentration of phytoplankton is high compared to non-biogenic particles.⁴⁷ Absorption by chlorophyll and related pigments therefore plays a major role in determining the total absorption coefficient in such waters, although detritus and dissolved organic matter derived from the phytoplankton also contribute to absorption in case 1 waters. Case 1 water can range from very clear (oligotrophic) water to very turbid (eutrophic) water, depending on the phytoplankton concentration. *Case 2 waters* are “everything else,” namely, waters where inorganic particles or dissolved organic matter from land drainage dominate so that absorption by pigments is relatively less important in determining the total absorption. (The cases 1 and 2 classifications must not be confused with the Jerlov water types 1 and 2, discussed later.) Roughly 98 percent of the world’s open ocean and coastal waters fall into the case 1 category, and therefore almost all bio-optical research has been directed toward these phytoplankton-dominated waters. However, near-shore and estuarine case 2 waters are disproportionately important to human interests such as recreation, fisheries, and military operations.

Prieur and Sathyendranath⁴⁸ developed a pioneering bio-optical model for the spectral absorption coefficient of case 1 waters. Their model was statistically derived from 90 sets of spectral absorption data taken in various case 1 waters and included absorption by phytoplankton pigments,

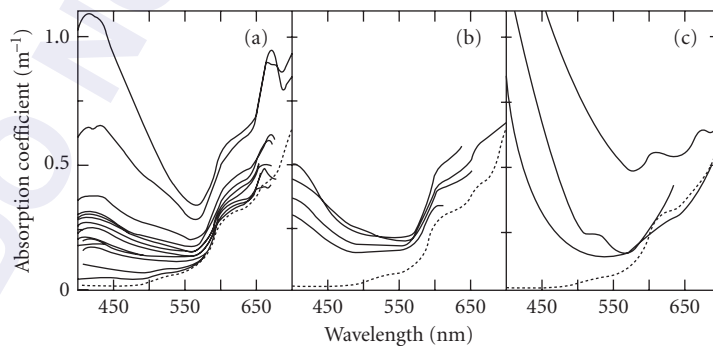


FIGURE 9 Examples of spectral absorption coefficients $a(\lambda)$ for various waters. Panel (a) shows $a(\lambda)$ for waters dominated by phytoplankton, panel (b) is for waters with a high concentration of nonpigmented particles, and panel (c) is for waters rich in yellow matter. (Based on Prieur and Sathyendranath,⁴⁸ with permission.)

by nonpigmented organic particles derived from deceased phytoplankton, and by yellow matter derived from decayed phytoplankton. The contribution of phytoplankton to the total absorption was parametrized in terms of the chlorophyll concentration C (i.e., chlorophyll a plus pheophytin a). The contributions of nonpigmented particles and of yellow matter were parametrized in terms of both the chlorophyll concentration and the total scattering coefficient at $\lambda = 550$ nm, $b(550)$. The essence of the Prieur-Sathyendranath model is contained in a more recent and simpler variant given by Morel:⁶

$$a(\lambda) = [a_w(\lambda) + 0.06a_c^{*'}(\lambda)C^{0.65}][1 + 0.2\exp(-0.014(\lambda - 440))] \quad (1)$$

Here $a_w(\lambda)$ is the absorption coefficient of pure water and $a_c^{*'}(\lambda)$ is a nondimensional, statistically derived chlorophyll-specific absorption coefficient; $a_w(\lambda)$ and $a_c^{*'}(\lambda)$ values are given in Table 9 [these $a_w(\lambda)$ values are slightly different than those of Table 6]. When C is expressed in mg m^{-3} and λ is in nm, the resulting $a(\lambda)$ is in m^{-1} .

Another simple bio-optical model for absorption has been developed independently by Kopelevich.⁴⁹ It has the form⁵⁰

$$a(\lambda) = a_w(\lambda) + C[a_c^0(\lambda) + 0.1\exp[-0.015(\lambda - 400)]]$$

where $a_c^0(\lambda)$ is the chlorophyll-specific absorption coefficient for phytoplankton ($\text{m}^2 \text{mg}^{-1}$), and $a_w(\lambda)$ and C are defined as for Eq. (1). The Kopelevich model as presently used⁴⁹ takes $a_w(\lambda)$ from Smith and Baker³² (Table 6) and takes $a_c^0(\lambda)$ from Yentsch.⁵¹

Although these and similar bio-optical models for absorption are frequently used, caution is advised in their application. Both models assume that the absorption by yellow matter covaries with that due to phytoplankton; i.e., each implies that a fixed percentage of the total absorption at a given wavelength always comes from yellow matter. The general validity of this assumption is doubtful even for open ocean waters: Bricaud et al.⁴² show data (Fig. 5) for which $a(375)$, used as an index for yellow matter concentration, is uncorrelated with chlorophyll concentration even in oceanic regions uninfluenced by freshwater runoff. Gordon⁵² has developed a model that avoids assuming any relation between yellow matter and phytoplankton. However, his model becomes singular as $C \rightarrow 0.01 \text{ mg m}^{-3}$ and cannot be expected to work well for C much less than 0.1 mg m^{-3} . The Kopelevich model has the chlorophyll contribution proportional to C , whereas the Morel model has $C^{0.65}$. The exponent of 0.65 is probably closer to reality, since it reflects a change in the relative contributions to absorption by phytoplankton and by detritus as the chlorophyll concentration changes (absorption by detritus is relatively more important at low chlorophyll concentrations⁵²). Moreover, the chlorophyll-specific absorption curve of Yentsch⁵¹ used in the Kopelevich

TABLE 9 Absorption by Pure Sea Water, a_w , and the Nondimensional Chlorophyll-Specific Absorption Coefficient $a_c^{*'}$ Used in the Prieur-Sathyendranath-Morel Model for the Spectral Absorption Coefficient $a(\lambda)$ *

λ (nm)	a_w (m^{-1})	$a_c^{*'}$	λ (nm)	a_w (m^{-1})	$a_c^{*'}$	λ (nm)	a_w (m^{-1})	$a_c^{*'}$
400	0.018	0.687	500	0.026	0.668	600	0.245	0.236
410	0.017	0.828	510	0.036	0.618	610	0.290	0.252
420	0.016	0.913	520	0.048	0.528	620	0.310	0.276
430	0.015	0.973	530	0.051	0.474	630	0.320	0.317
440	0.015	1.000	540	0.056	0.416	640	0.330	0.334
450	0.015	0.944	550	0.064	0.357	650	0.350	0.356
460	0.016	0.917	560	0.071	0.294	660	0.410	0.441
470	0.016	0.870	570	0.080	0.276	670	0.430	0.595
480	0.018	0.798	580	0.108	0.291	680	0.450	0.502
490	0.020	0.750	590	0.157	0.282	690	0.500	0.329
						700	0.650	0.215

*Condensed with permission from Prieur and Sathyendranath,⁴⁸ who give values every 5 nm.

model is based on laboratory cultures of phytoplankton, whereas the later work by Prieur and Sathyendranath used in situ observations to derive the a_c^* (λ) values of Table 9—an additional point in favor of Eq. (1). Either of these bio-optical models is useful but clearly imperfect. They may (or may not) give correct *average* values, but they give no information about the *variability* of $a(\lambda)$. It can be anticipated that the simple models now available will be replaced, perhaps by models designed for specific regions and seasons, as better understanding of the variability inherent in spectral absorption is achieved.

1.16 MEASUREMENT OF SCATTERING

Scattering in natural waters is caused both by small scale ($\ll \lambda$) density fluctuations attributable to random molecular motions and by the ubiquitous large ($> \lambda$) organic and inorganic particles. Scattering by water molecules (and salt ions, in seawater) determines the minimum values for the scattering properties. However, as is the case for absorption, the scattering properties of natural waters are greatly modified by the particulate matter that is always present.

Scattering measurements are even more difficult than absorption measurements. The conceptual design of an instrument for measuring the volume scattering function $\beta(\psi, \lambda)$ is no more complicated than Fig. 2 and the defining equation $\beta(\psi, \lambda) = I_s(\psi, \lambda)/[E_i(\lambda) \Delta V]$: a collimated beam of known irradiance E_i illuminates a given volume of water ΔV and the scattered intensity I_s is measured as a function of scattering angle and wavelength. However, the engineering of instruments capable of the in situ determination of $\beta(\psi, \lambda)$ is quite difficult. The magnitude of the scattered intensity typically increases by five or six orders of magnitude in going from $\psi = 90^\circ$ to $\psi = 0.1^\circ$ for a given natural water sample, and scattering at a given angle ψ can vary by two orders of magnitude among water samples. The required dynamic range of an instrument is therefore great. Corrections must be made for absorption within the sample volume and also along the incident and scattered beam paths for in situ instruments. The rapid change in $\beta(\psi, \lambda)$ at small scattering angles requires very precise alignment of the optical elements, but rolling ships seem designed to knock things out of alignment. Because of these design difficulties only a few one-of-a-kind instruments have been built for in situ measurement of the volume scattering function, and measurements of $\beta(\psi, \lambda)$ are not routinely made. Petzold⁵³ gives the details of two such instruments, one for small scattering angles ($\psi = 0.085, 0.17, \text{ and } 0.34^\circ$) and one for larger angles ($10^\circ \leq \psi \leq 170^\circ$); these are the instruments used to obtain the data presented in Sec. 1.18. Other instruments are referenced in Kirk⁸ and in Jerlov.²⁹

Commercial instruments are available for laboratory measurement of $\beta(\psi, \lambda)$ at fixed scattering angles (e.g., ψ every 5° from $\sim 20^\circ$ to $\sim 160^\circ$). These instruments are subject to their own problems, such as degradation of samples between the times of collection and measurement. Moreover, measurements of $\beta(\psi, \lambda)$ over a limited range of ψ are not sufficient for determination of $b(\lambda)$ by integration. In practice, the scattering coefficient $b(\lambda)$ is usually determined by the conservation of energy relation $b(\lambda) = c(\lambda) - a(\lambda)$ after measurements of beam attenuation and absorption have been made.

Both in situ and laboratory instruments sample ($\sim \text{cm}^3$) volumes of water and therefore may fail to detect the presence of optically significant large aggregates (marine snow) if such particles are too few in number to be reliably captured in the sample volume. However, such particles can affect the scattering properties of large volumes of water (e.g., as seen in remote sensing or underwater visibility studies).

Measurements at near forward ($\psi < 1^\circ$) and near backward ($\psi > 179^\circ$) angles are exceptionally difficult to make, yet the behavior of $\beta(\psi, \lambda)$ at these extreme angles is of considerable interest. Accurate determination of β at small angles is crucial to the determination of b by integration since typically one-half of all scattering takes place at angles of less than a few degrees. Scattering at small angles is important in underwater imaging and it is of theoretical interest for its connections to scattering theory, particle optical properties, and particle size distributions. The behavior of β very near $\psi = 180^\circ$ is important in laser remote-sensing applications.

Spinrad et al.⁵⁴ and Padmabandu and Fry⁵⁵ have reported measurements at very small angles on suspensions of polystyrene spheres but no such measurements have been published for natural water samples. The Padmabandu and Fry technique is notable in that it allows the measurement of β at $\psi = 0^\circ$ exactly by use of the coupling of two coherent beams in a photorefractive crystal to measure the phase shift that corresponds to 0° scattering. Measurement of $\beta(0, \lambda)$ is of theoretical interest because of its relation to attenuation via the optical theorem. Enhanced backscatter has been reported⁵⁶ in suspensions of latex spheres; a factor-of-two increase in scattered intensity between $\psi = 179.5$ and 180.0° is typical. Whether or not such backscattering enhancement ever occurs in natural waters is a subject of heated debate.

1.17 SCATTERING BY PURE WATER AND BY PURE SEA WATER

Morel⁵⁷ has reviewed in detail the theory and observations pertaining to scattering by pure water and by pure sea water. In pure water random molecular motions give rise to rapid fluctuations in the number of molecules in a given volume ΔV , where ΔV is small compared to the wavelength of light but large compared to atomic scales (so that the liquid within the volume is adequately described by statistical thermodynamics). The Einstein-Smoluchowski theory of scattering relates these fluctuations in molecular number density to associated fluctuations in the index of refraction, which give rise to scattering. In sea water the basic theory is the same but random fluctuations in the concentrations of the various ions (Cl^- , Na^+ , etc.) give somewhat greater index of refraction fluctuations, and hence greater scattering. The net result of these considerations is that the volume scattering function for either pure water or for pure sea water has the form

$$\beta_w(\psi, \lambda) = \beta_w(90^\circ, \lambda_0) \left(\frac{\lambda_0}{\lambda} \right)^{4.32} (1 + 0.835 \cos^2 \psi) \quad \text{m}^{-1} \text{sr}^{-1} \quad (2)$$

which is reminiscent of the form for Rayleigh scattering. The wavelength dependence of $\lambda^{-4.32}$ rather than λ^{-4} (for Rayleigh scattering) results from the wavelength dependence of the index of refraction. The 0.835 factor is attributable to the anisotropy of the water molecules. The corresponding phase function is

$$\tilde{\beta}_w(\psi) = 0.06225(1 + 0.835 \cos^2 \psi) \quad \text{sr}^{-1}$$

and the total scattering coefficient $b_w(\lambda)$ is given by

$$b_w(\lambda) = 16.06 \left(\frac{\lambda_0}{\lambda} \right)^{4.32} \beta_w(90^\circ, \lambda_0) \quad \text{m}^{-1} \quad (3)$$

Table 10 gives values of $\beta_w(90^\circ, \lambda)$ and $b_w(\lambda)$ for selected wavelengths for both pure water and pure sea water ($S = 35$ to 39‰). Note that the pure sea water values are about 30 percent greater than the pure water values at all wavelengths. Table 11 shows the dependence of $b_w(546)$ on pressure, temperature, and salinity. Note that molecular scattering decreases as decreasing temperature or increasing pressure reduce the smallscale fluctuations.

1.18 SCATTERING BY PARTICLES

Heroic efforts are required to obtain water of sufficient purity that a Rayleigh-like volume scattering function is observed. As soon as there is a slight amount of particulate matter in the water—always the case for even the clearest natural water—the volume scattering function becomes highly peaked in the forward direction, and the scattering coefficient increases by at least a factor of 10.

TABLE 10 The Volume Scattering Function at $\psi = 90^\circ$, $\beta_w(90^\circ, \lambda)$, and the Scattering Coefficient $b_w(\lambda)$ for Pure Water and for Pure Sea Water ($S = 35\text{--}39\%$)*

λ (nm)	Pure water		Pure sea water	
	$\beta_w(90^\circ)$ ($10^{-4} \text{ m}^{-1} \text{ sr}^{-1}$)	b_w^\dagger (10^{-4} m^{-1})	$\beta(90^\circ)$ ($10^{-4} \text{ m}^{-1} \text{ sr}^{-1}$)	b_w^\dagger (10^{-4} m^{-1})
350	6.47	103.5	8.41	134.5
375	4.80	76.8	6.24	99.8
400	3.63	58.1	4.72	75.5
425	2.80	44.7	3.63	58.1
450	2.18	34.9	2.84	45.4
475	1.73	27.6	2.25	35.9
500	1.38	22.2	1.80	28.8
525	1.12	17.9	1.46	23.3
550	0.93	14.9	1.21	19.3
575	0.78	12.5	1.01	16.2
600	0.68	10.9	0.88	14.1

*Reproduced from Morel,⁵⁷ with permission.†Computed from $b(\lambda) = 16.0\beta(90^\circ, \lambda)$.

Even for the most numerous oceanic particles (e.g., colloids at a concentration of 10^{15} m^{-3}) the average distance between particles is greater than ten wavelengths of visible light. For the optically most significant phytoplankton the average separation is thousands of wavelengths. Moreover, these particles usually are randomly distributed and oriented. Ocean water therefore can be treated as a very dilute suspension of random scatterers and consequently the intensity of light scattered by an ensemble of particles is given by the sum of the intensities due to the individual particles. Coherent scattering effects are negligible except perhaps at extremely small scattering angles.⁵⁸ An overview of scattering by particles is given in Chap. 7 "Scattering by Particles" by Craig F. Bohren in Vol. I.

The contribution of the particulate matter to the total volume scattering function $\beta(\psi, \lambda)$ is obtained from

$$\beta_p(\psi, \lambda) \equiv \beta(\psi, \lambda) - \beta_w(\psi, \lambda)$$

Here the subscript p refers to particles, and w refers to pure water (if β is measured in fresh water) or pure sea water (for oceanic measurements). Figure 10 shows several particle volume scattering functions determined from in situ measurements of $\beta(\psi, \lambda)$ in a variety of waters ranging from very clear to very turbid. The particles cause at least a four-order-of-magnitude increase in scattering between $\psi \approx 90^\circ$ and $\psi \approx 1^\circ$. The contribution of molecular scattering to the total is therefore completely negligible except at backscattered directions ($\psi \geq 90^\circ$) in the clearest natural waters. The top curve in Fig. 10 is shown for small scattering angles in Fig. 11. The scattering function shows

TABLE 11 Computed Scattering Coefficient b of Pure Water ($S = 0$) and of Pure Sea Water ($S = 35\%$) at $\lambda = 546 \text{ nm}$ as a Function of Temperature T and Pressure p (Numbers in the body of the table have units of m^{-1} .*)

T ($^\circ\text{C}$)	$p = 10^5 \text{ Pa}$ (1 atm)		$p = 10^7 \text{ Pa}$ (100 atm)		$p = 10^8 \text{ Pa}$ (1000 atm)	
	$S = 0$	$S = 35\%$	$S = 0$	$S = 35\%$	$S = 0$	$S = 35\%$
0	0.00145	0.00195	0.00140	0.00192	0.00110	0.00167
10	0.00148	0.00203	0.00143	0.00200	0.00119	0.00176
20	0.00149	0.00207	0.00147	0.00204	0.00125	0.00183
40	0.00150	0.00213	0.00149	0.00212	0.00136	0.00197

*Data extracted from the more extensive table of Shifrin,⁵⁸ with permission.

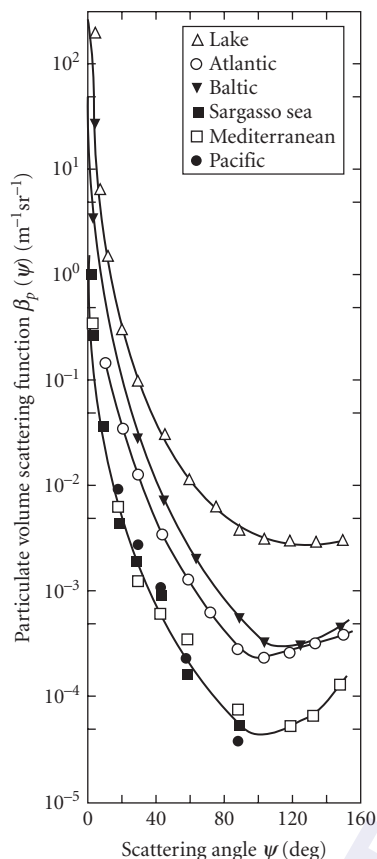


FIGURE 10 Particulate volume scattering functions $\beta_p(\psi, \lambda)$ determined from *in situ* measurements in various waters; wavelengths vary. (Redrawn from Kullenberg,⁵⁹ with permission.)

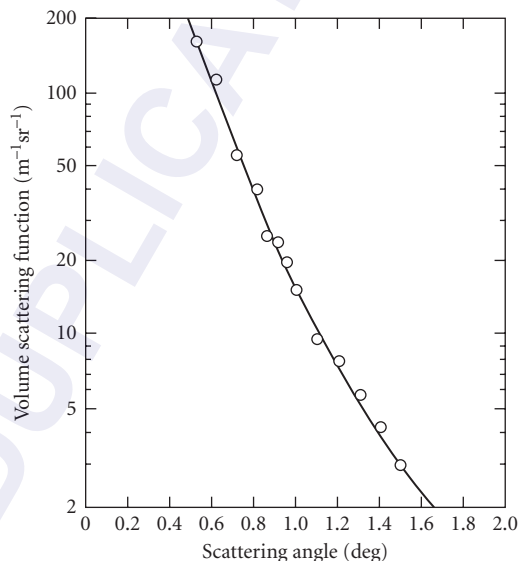


FIGURE 11 Detail of the forward scattering values of the “lake” volume scattering function seen in the top curve of Fig. 10. (Redrawn from Preisendorfer.²)

no indication of “flattening out” even at angles as small as 0.5° . Note that the scattering function increases by a factor of 100 over only a one-degree range of scattering angle.

Highly peaked forward scattering such as that seen in Figs. 10 and 11 is characteristic of diffraction-dominated scattering in a polydisperse system. Scattering by refraction and reflection from particle surfaces becomes important at large scattering angles ($\psi \geq 15^\circ$). Mie scattering calculations are well able to reproduce observed volume scattering functions given the proper optical properties and size distributions. Early efforts along these lines are seen in Kullenberg⁵⁹ and in Brown and Gordon.⁶⁰ Brown and Gordon were unable to reproduce observed backscattering values using measured particle size distributions. However, their instruments were unable to detect submicrometer particles. They found that the Mie theory properly predicted backscattering if they assumed the presence of numerous, submicrometer, low-index-of-refraction particles. It is reasonable to speculate that bacteria and the recently discovered colloidal particles are the particles whose existence was inferred by Brown and Gordon. Recent Mie scattering calculations⁶¹ have used three-layered spheres to model the structure of phytoplankton (cell wall, chloroplasts, and cytoplasm core) and have used polydisperse mixtures of both organic and inorganic particles.

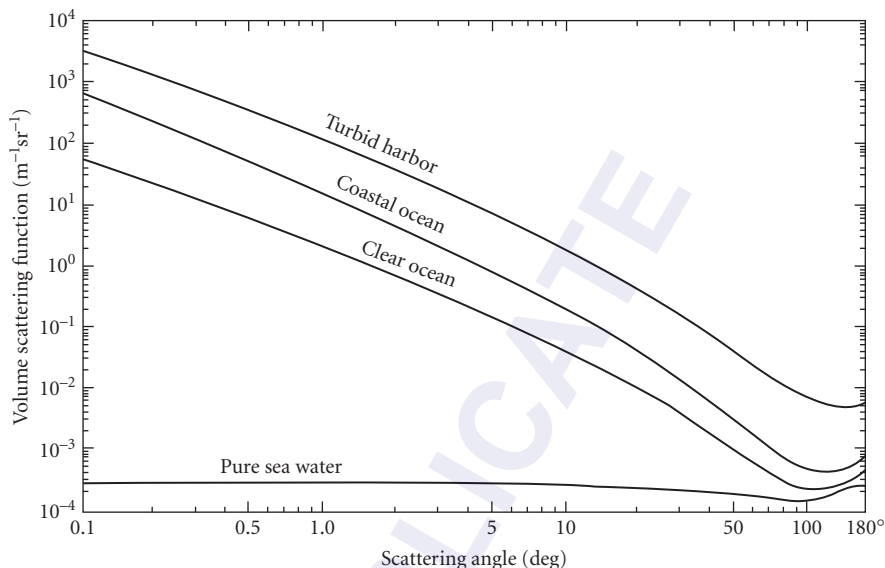


FIGURE 12 Measured volume scattering functions from three different natural waters and the computed volume scattering function for pure sea water, all at $\lambda = 514$ nm. (Redrawn from Petzold.⁵³)

The most carefully made and widely cited scattering measurements are found in Petzold.⁵³ Figure 12 shows three of his $\beta(\psi, \lambda)$ curves displayed on a log-log plot to emphasize the forward scattering angles. The instruments had a spectral response centered at $\lambda = 514$ nm with an FWHM of 75 nm. The top curve was obtained in the very turbid water of San Diego Harbor, California; the center curve comes from near-shore coastal water in San Pedro Channel, California; and the bottom curve is from very clear water in the Tongue of the Ocean, Bahama Islands. The striking feature of these volume scattering functions (and those of Fig. 10) from very different waters is the similarity of their shapes. Although the scattering coefficients b of the curves in Fig. 12 vary by a factor of 50 (Table 13), the uniform shapes suggest that it is reasonable to define a “typical” particle phase $\tilde{\beta}_p(\psi, \lambda)$. This has been done⁶² with three sets of Petzold’s data from waters with a high particulate load (one set being the top curve of Fig. 12), as follows: (1) subtract $\beta_w(\psi, \lambda)$ from each curve to get three particle volume scattering functions $\beta_p^i(\psi, \lambda)$, $i = 1, 2, 3$; (2) compute three particle phase functions via $\tilde{\beta}_p^i(\psi, \lambda) = \beta_p^i(\psi, \lambda) / b^i(\lambda)$; (3) average the three particle phase functions to define the typical particle phase function $\tilde{\beta}_p(\psi, \lambda)$. Table 12 displays the three Petzold volume scattering functions plotted in Fig. 12, the volume scattering function for pure sea water, and the average particle phase function computed as just described. This particle phase function satisfies the normalization $2\pi \int_0^\pi \tilde{\beta}_p(\psi, \lambda) \sin \psi d\psi = 1$ if a behavior of $\tilde{\beta}_p \sim \psi^{-m}$ is assumed for $\psi < 0.1^\circ$ (m is a positive constant between zero and two, determined from $\tilde{\beta}_p$ at the smallest measured angles), and a trapezoidal rule integration is used for $\psi \geq 0.1^\circ$, with linear interpolation used between the tabulated values. This average particle phase function is adequate for many radiative transfer calculations. However, the user must remember that significant deviations from the average can be expected in nature (e.g., in waters with abnormally high numbers of either large or small particles), although the details of such deviations have not been quantified.

Table 13 compares several inherent optical properties for pure sea water and for the three Petzold water samples of Fig. 12 and Table 12. These data show how greatly different even clear ocean water is from pure sea water. Note that natural water ranges from absorption-dominated ($\omega_0 = 0.247$) to scattering-dominated ($\omega_0 = 0.833$) at $\lambda = 514$ nm. The ratio of backscattering to total scattering is typically a few percent in natural water. However, there is no clear relation between b_b/b and the water type, at least for the Petzold data of Table 13. This lack of an obvious relation is likely the result of

TABLE 12 Volume Scattering Functions $\beta(\psi, \lambda)$ for Three Oceanic Waters and for Pure Sea Water and a Typical Particle Phase $\beta_p(\psi, \lambda)$, All at $\lambda = 514$ nm

Scattering Angle (deg)	Volume scattering functions ($\text{m}^{-1} \text{sr}^{-1}$)				
	Clear Ocean*	Coastal Ocean*	Turbid Harbor*	Pure Sea Water†	Particle Phase Function‡ (sr^{-1})
0.100	5.318×10^1	6.533×10^2	3.262×10^3	2.936×10^{-4}	1.767×10^3
0.126	4.042	4.577	2.397	2.936	1.296
0.158	3.073	3.206	1.757	2.936	9.502×10^2
0.200	2.374	2.252	1.275	2.936	6.991
0.251	1.814	1.579	9.260×10^2	2.936	5.140
0.316	1.360	1.104	6.764	2.936	3.764
0.398	9.954×10	7.731×10^1	5.027	2.936	2.763
0.501	7.179	5.371	3.705	2.936	2.012
0.631	5.110	3.675	2.676	2.936	1.444
0.794	3.591	2.481	1.897	2.936	1.022
1.000	2.498	1.662	1.329	2.936	7.161×10^1
1.259	1.719	1.106	9.191×10^1	2.935	4.958
1.585	1.171	7.306×10^0	6.280	2.935	3.395
1.995	7.758×10^{-1}	4.751	4.171	2.934	2.281
2.512	5.087	3.067	2.737	2.933	1.516
3.162	3.340	1.977	1.793	2.932	1.002
3.981	2.196	1.273	1.172	2.930	6.580×10^0
5.012	1.446	8.183×10^{-1}	7.655×10^0	2.926	4.295
6.310	9.522×10^{-2}	5.285	5.039	2.920	2.807
7.943	6.282	3.402	3.302	2.911	1.819
10.0	4.162	2.155	2.111	2.896	1.153
15.0	2.038	9.283×10^{-2}	9.041×10^{-1}	2.847	4.893×10^{-1}
20.0	1.099	4.427	4.452	2.780	2.444
25.0	6.166×10^{-3}	2.390	2.734	2.697	1.472
30.0	3.888	1.445	1.613	2.602	8.609×10^{-2}
35.0	2.680	9.063×10^{-3}	1.109	2.497	5.931
40.0	1.899	6.014	7.913×10^{-2}	2.384	4.210
45.0	1.372	4.144	5.858	2.268	3.067
50.0	1.020	2.993	4.388	2.152	2.275
55.0	7.683×10^{-4}	2.253	3.288	2.040	1.699
60.0	6.028	1.737	2.548	1.934	1.313
65.0	4.883	1.369	2.041	1.839	1.046
70.0	4.069	1.094	1.655	1.756	8.488×10^{-3}
75.0	3.457	8.782×10^{-4}	1.345	1.690	6.976
80.0	3.019	7.238	1.124	1.640	5.842
85.0	2.681	6.036	9.637×10^{-3}	1.610	4.953
90.0	2.459	5.241	8.411	1.600	4.292
95.0	2.315	4.703	7.396	1.610	3.782
100.0	2.239	4.363	6.694	1.640	3.404
105.0	2.225	4.189	6.220	1.690	3.116
110.0	2.239	4.073	5.891	1.756	2.912
115.0	2.265	3.994	5.729	1.839	2.797
120.0	2.339	3.972	5.549	1.934	2.686
125.0	2.505	3.984	5.343	2.040	2.571
130.0	2.629	4.071	5.154	2.152	2.476
135.0	2.662	4.219	4.967	2.268	2.377
140.0	2.749	4.458	4.822	2.384	2.329
145.0	2.896	4.775	4.635	2.497	2.313

(Continued)

TABLE 12 Volume Scattering Functions $\beta(\psi, \lambda)$ for Three Oceanic Waters and for Pure Sea Water and a Typical Particle Phase $\tilde{\beta}_p(\psi, \lambda)$, All at $\lambda = 514$ nm (*Continued*)

Scattering Angle (deg)	Volume scattering functions ($\text{m}^{-1} \text{sr}^{-1}$)				Particle Phase Function [‡] (sr^{-1})
	Clear Ocean [*]	Coastal Ocean [*]	Turbid Harbor [*]	Pure Sea Water [†]	
150.0	3.088	5.232	4.634	2.602	2.365
155.0	3.304	5.824	4.900	2.697	2.506
160.0	3.627	6.665	5.142	2.780	2.662
165.0	4.073	7.823	5.359	2.847	2.835
170.0	4.671	9.393	5.550	2.896	3.031
175.0	4.845	9.847	5.618	2.926	3.092
180.0	5.109	1.030×10^{-3}	5.686	2.936	3.154

^{*}Data reproduced from Petzold.⁵³

[†]Computed from Eq. (2) and Table 10.

[‡]Courtesy of H. R. Gordon; see also Ref. 62.

differing particle types in the three waters. Since refraction and reflection are important processes at large scattering angles, the particle indices of refraction are important in determining b_p . Total scattering is dominated by diffraction and so particle composition has little effect on b values. The last row of Table 13 gives the angle ψ such that one-half of the total scattering occurs at angles between 0 and ψ . This angle is rarely greater than 10° in natural waters.

1.19 WAVELENGTH DEPENDENCE OF SCATTERING: BIO-OPTICAL MODELS

The strong $\lambda^{-4.32}$ wavelength dependence of pure-water scattering is not seen in natural waters. This is because scattering is dominated by diffraction from polydisperse particles that are usually much larger than the wavelength of visible light. Although diffraction depends on the particle size-to-wavelength ratio, the presence of particles of many sizes diminishes the wavelength effects that are seen in diffraction by a single particle. Moreover, diffraction does not depend on particle composition. However, some wavelength dependence is to be expected, especially at backward scattering angles where refraction, and hence particle composition, is important. Molecular scattering also contributes something to the total scattering and can even dominate the particle contribution at backscatter angles in clear water.⁶³

Morel⁶⁴ presents several useful observations on the wavelength dependence of scattering. Figure 13 shows two sets of volume scattering functions, one from the very clear waters of the Tyrrhenian Sea and one from the turbid English Channel. Each set displays $\beta(\psi, \lambda)/\beta(90^\circ, \lambda)$ for $\lambda = 366, 436$, and

TABLE 13 Selected Inherent Optical Properties for the Waters Presented in Fig. 12 and in Table 12 (All values are for $\lambda = 514$ nm except as noted.)

Water	a (m^{-1})	b (m^{-1})	c (m^{-1})	ω_0	b_0/b	ψ for $\frac{1}{2}b$ (deg)
Pure sea water	0.0405 [*]	0.0025 [†]	0.043	0.058	0.500	90.00
Clear ocean	0.114 [‡]	0.037	0.151 [§]	0.247	0.044	6.25
Coastal ocean	0.179 [‡]	0.219	0.398 [§]	0.551	0.013	2.53
Turbid harbor	0.366 [‡]	1.824	2.190 [§]	0.833	0.020	4.68

^{*}Value obtained by interpolation in Table 6.

[†]Value obtained by interpolation in Table 10.

[‡]Estimated by Petzold⁵³ from $c(530 \text{ nm}) - b(514 \text{ nm})$.

[§]Measured by Petzold at $\lambda = 530$ nm.

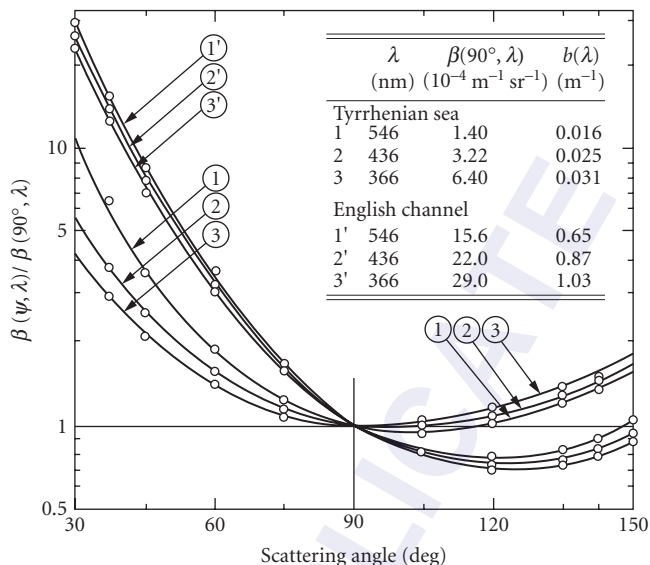


FIGURE 13 Wavelength dependence of total volume scattering functions measured in very clear (Tyrrhenian Sea) and in turbid (English Channel) waters. (Redrawn from Morel.⁶⁴)

546 nm. The clear water shows a definite dependence of the shape of $\beta(\psi, \lambda)$ on λ whereas the particle-rich turbid water shows much less wavelength dependence. In each case the volume scattering function of shortest wavelength is most nearly symmetric about $\psi = 90^\circ$, presumably because symmetric molecular scattering is contributing relatively more to the total scattering at short wavelengths.

Figure 14 shows a systematic wavelength dependence of particle volume scattering functions. Figure 14a shows average values of $\beta_p(\psi, 366 \text{ nm})/\beta_p(\psi, 546 \text{ nm})$ for N samples as labeled in

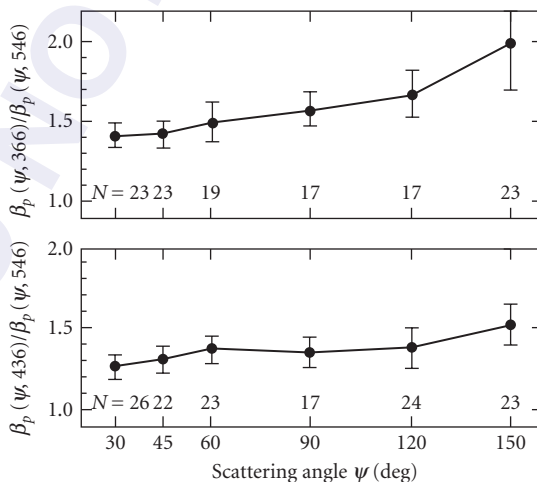


FIGURE 14 Wavelength dependence of particulate volume scattering functions. N is the number of samples. (Redrawn from Morel.⁶⁴)

TABLE 14 Exponents n Required to Fit the Data of Fig. 14 Assuming That $\beta_p(\psi, \lambda) = \beta_p(\psi, 546)(546/\lambda)^n$

Wavelength λ (nm)	Scattering angle ψ		
	30°	90°	150°
366	0.84	1.13	1.73
436	0.99	1.33	1.89

the figure. The vertical bars are one standard deviation of the observations. Figure 14b shows the ratio for $\lambda = 436$ to 546 nm. These ratios clearly depend both on wavelength and scattering angle. Assuming that $\beta_p(\psi, \lambda)$ has a wavelength dependence of

$$\beta_p(\psi, \lambda) = \beta_p(\psi, 546) \left(\frac{546 \text{ nm}}{\lambda} \right)^n$$

the data of Fig. 14 imply values for n as seen in Table 14. As anticipated, the wavelength dependence is strongest for backscatter ($\psi = 150^\circ$) and weakest for forward scatter ($\psi = 30^\circ$).

Kopelevich^{49,65} has statistically derived a two-parameter model for spectral volume scattering functions (VSFs). This model separates the contributions by “small” and “large” particles to the particulate scattering. Small particles are taken to be mineral particles less than 1 μm in size and having an index of refraction of $n = 1.15$; large particles are biologic particles larger than 1 μm in size and having an index refraction of $n = 1.03$. The model is defined by

$$\beta(\psi, \lambda) = \beta_w(\psi, \lambda) + v_s \beta_s^*(\psi) \left(\frac{550 \text{ nm}}{\lambda} \right)^{1.7} + v_\ell \beta_\ell^*(\psi) \left(\frac{550 \text{ nm}}{\lambda} \right)^{0.3} \quad (4)$$

with the following definitions:

- $\beta_w(\psi, \lambda)$ the VSF of pure sea water, given by Eq. (2) with $\lambda_0 = 550$ nm and an exponent of 4.30
- v_s the volume concentration of small particles, with units of cm^3 of particles per m^3 of water, i.e., parts per million (ppm)
- v_ℓ the analogous volume concentration of large particles
- $\beta_s^*(\psi)$ the small-particle VSF per unit volume concentration of small particles, with units of $\text{m}^{-1} \text{sr}^{-1} \text{ppm}^{-1}$
- $\beta_\ell^*(\psi)$ the analogous large-particle concentration-specific VSF

The concentration-specific VSFs for small and large particles are given in Table 15. Equation (4) can be evaluated as if the two parameters v_s and v_ℓ are known; the ranges of values for oceanic waters are $0.01 \leq v_s \leq 0.20$ ppm $0.01 \leq v_\ell \leq 0.40$ ppm. However, these two parameters are themselves parametrized in terms of the total volume scattering function measured at $\lambda = 550$ nm for $\psi = 1$ and 45° :

$$\begin{aligned} v_s &= -1.45 \times 10^{-4} \beta(1^\circ, 550 \text{ nm}) + 10.2 \beta(45^\circ, 550 \text{ nm}) - 0.002 \\ v_\ell &= 2.2 \times 10^{-2} \beta(1^\circ, 550 \text{ nm}) - 1.2 \beta(45^\circ, 550 \text{ nm}) \end{aligned} \quad (5)$$

Thus $\beta(\psi, \lambda)$ can also be determined from two measurements of the total VSF.

The mathematical form of the Kopelevich model reveals its underlying physics. Large particles give diffractive scattering at very small angles; thus $\beta_\ell^*(\psi)$ is highly peaked for small ψ and the wavelength dependence of the large particle term is weak ($\lambda^{-0.3}$). Small particles contribute more to

TABLE 15 The Concentration-Specific Volume Scattering Functions for Small (β_s^*) and Large (β_ℓ^*) Particles As a Function of the Scattering Angle ψ for Use in the Kopelevich Model for Spectral Volume Scattering Functions, Eq. (4)*

ψ (deg)	$\beta_s^* \left(\frac{\text{m}^{-1} \text{sr}^{-1}}{\text{ppm}} \right)$	$\beta_\ell^* \left(\frac{\text{m}^{-1} \text{sr}^{-1}}{\text{ppm}} \right)$	ψ (deg)	$\beta_s^* \left(\frac{\text{m}^{-1} \text{sr}^{-1}}{\text{ppm}} \right)$	$\beta_\ell^* \left(\frac{\text{m}^{-1} \text{sr}^{-1}}{\text{ppm}} \right)$
0	5.3	140	45	9.8×10^{-2}	6.2×10^{-4}
0.5	5.3	98	60	4.1	3.8
1	5.2	46	75	2.0	2.0
1.5	5.2	26	90	1.2	6.3×10^{-5}
2	5.1	15	105	8.6×10^{-3}	4.4
4	4.6	3.6	120	7.4	2.9
6	3.9	1.1	135	7.4	2.0
10	2.5	0.20	150	7.5	2.0
15	1.3	5.0×10^{-2}	180	8.1	7.0
30	0.29	2.8×10^{-3}	$b^* =$	$1.34 \text{ m}^{-1}/\text{ppm}$	$0.312 \text{ m}^{-1}/\text{ppm}$

*Reproduced from Kopelevich.⁴⁹

scattering at large angles and thus have a more symmetric VSF and a stronger wavelength dependence ($\lambda^{-1.7}$). This model gives a reasonably good description of VSFs observed in a variety of waters (Shifrin,⁵⁸ fig. 5.20).

Several simple models are available for the scattering coefficient $b(\lambda)$. A commonly employed bio-optical model for $b(\lambda)$ is that of Gordon and Morel⁶⁶ (see also Ref. 6):

$$b(\lambda) = b_w(\lambda) + \left(\frac{550 \text{ nm}}{\lambda} \right) 0.30C^{0.62} \quad \text{m}^{-1} \quad (6)$$

Here $b_w(\lambda)$ is given by Eq. (3) and Table 10. λ is in nm and C is the chlorophyll concentration in mg m^{-3} . A related bio-optical model for the backscatter coefficient $b_b(\lambda)$ is found in Morel.⁴⁵

$$b_b(\lambda) = \frac{1}{2} b_w(\lambda) + \left[0.002 + 0.02 \left(\frac{1}{2} - \frac{1}{4} \log C \right) \left(\frac{500 \text{ nm}}{\lambda} \right) \right] 0.30C^{0.62}$$

The $\left(\frac{1}{2} - \frac{1}{4} \log C \right)$ factor gives $b_b(\lambda)$ a λ^{-1} wavelength dependence in very clear ($C = 0.01 \text{ mg m}^{-3}$) water and no wavelength dependence in very turbid ($C = 100 \text{ mg m}^{-3}$) water. These empirically derived models are intended for use only in case 1 waters.

A feeling for the accuracy of the $b(\lambda)$ model of Eq. (6) can be obtained from Fig. 15, which plots measured $b(550 \text{ nm})$ values versus chlorophyll concentration C in both case 1 and case 2 waters. Note that even when the model is applied to the case 1 waters from which it was derived, the predicted $b(550 \text{ nm})$ value easily can be wrong by a factor of 2. If the model is misapplied to case 2 waters, the error can be an order of magnitude. Note that for a given C value, $b(550 \text{ nm})$ is higher in case 2 waters than in case 1 waters, presumably because of the presence of additional particles that do not contain chlorophyll.

Integration over ψ of the Kopelevich $b(\psi, \lambda)$ model of Eq. (4) yields another model for $b(\lambda)$:

$$b(\lambda) = 0.0017 \left(\frac{500 \text{ nm}}{\lambda} \right)^{4.3} + 1.34v_s \left(\frac{550 \text{ nm}}{\lambda} \right)^{1.7} + 0.312v_\ell \left(\frac{550 \text{ nm}}{\lambda} \right)^{0.3} \quad \text{m}^{-1}$$

where v_s and v_ℓ are given by Eq. (5). Kopelevich claims that the accuracy of this model is ~ 30 percent.

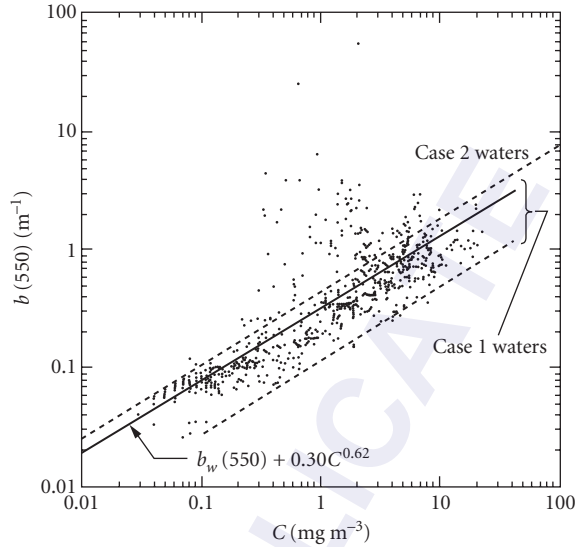


FIGURE 15 Measured scattering coefficients at $\lambda = 550$ nm, $b(550)$, as a function of chlorophyll concentration C . Case 1 waters lie between the dashed lines. Case 2 waters lie above the upper dashed line, which is defined by $b(550) = 0.45C^{0.62}$. The solid line is the model of Eq. (6). (Redrawn from Gordon and Morel,⁶⁶ with permission.)

A bio-optical model related to the Kopelevich model is found in Haltrin and Kattawar⁵⁰ (their notation):

$$b(\lambda) = b_w(\lambda) + b_{ps}^0(\lambda)P_s + b_{pl}^0(\lambda)P_l$$

Here $b_w(\lambda)$ is given by

$$b_w(\lambda) = 5.826 \times 10^{-3} \left(\frac{400}{\lambda} \right)^{4.322}$$

which is essentially the same as Eq. (3) and the data in Table 10. The terms $b_{ps}^0(\lambda)$ and $b_{pl}^0(\lambda)$ are the specific scattering coefficients for small and large particles, respectively, and are given by

$$b_{ps}^0(\lambda) = 1.1513 \left(\frac{400}{\lambda} \right)^{1.7} \quad \text{m}^2 \text{ g}^{-1}$$

$$b_{pl}^0(\lambda) = 0.3411 \left(\frac{400}{\lambda} \right)^{0.3} \quad \text{m}^2 \text{ g}^{-1}$$

P_s and P_l are the concentrations in g m^{-3} of small and large particles, respectively. These quantities are parametrized in terms of the chlorophyll concentration C , as shown in Table 16. This work also presents a model for backscattering:

$$b_b(\lambda) = \frac{1}{2} b_w(\lambda) + B_s b_{ps}^0(\lambda)P_s + B_l b_{pl}^0(\lambda)P_l$$

TABLE 16 Parameterization of Small (P_s) and Large (P_ℓ) Particle Concentrations in Terms of the Chlorophyll Concentration C for Use in the Kopelevich-Haltrin-Kattawar Models for $b(\lambda)$ and $b_b(\lambda)$ *

C (mg m^{-3})	P_s (g m^{-3})	P_ℓ (g m^{-3})
0.00	0.000	0.000
0.03	0.001	0.035
0.05	0.002	0.051
0.12	0.004	0.098
0.30	0.009	0.194
0.60	0.016	0.325
1.00	0.024	0.476
3.00	0.062	1.078

*Reproduced from Haltrin and Kattawar,⁵⁰ with permission

Here $B_s = 0.039$ is the backscattering probability for small particles and $B_\ell = 0.00064$ is the backscattering probability for large particles.

The bio-optical models for scattering just discussed are useful but very approximate.

The reason for the frequent large discrepancies between model predictions and measured reality likely lies in the fact that scattering depends not just on particle concentration (as parameterized in terms of chlorophyll concentration), but also on the particle index of refraction and on the details of the particle size distribution which are not well parameterized in terms of the chlorophyll concentration alone. Whether or not the Kopelevich model or its derivative Haltrin-Kattawar form which partition the scattering into large and small particle components is in some sense better than the Gordon-Morel model is not known at present. Another consequence of the complexity of scattering is seen in the next section.

1.20 BEAM ATTENUATION

The spectral beam attenuation coefficient $c(\lambda)$ is just the sum of the spectral absorption and scattering coefficients: $c(\lambda) = a(\lambda) + b(\lambda)$. Since both $a(\lambda)$ and $b(\lambda)$ are highly variable functions of the nature and concentration of the constituents of natural waters so is $c(\lambda)$. Beam attenuation near $\lambda = 660$ nm is the only inherent optical property of water that is easily, accurately, and routinely measured. This wavelength is used both for engineering reasons (the availability of a stable LED light source) and because absorption by yellow matter is negligible in the red. Thus the quantity

$$c_p(660 \text{ nm}) \equiv c(660 \text{ nm}) - a_w(660 \text{ nm}) - b_w(660 \text{ nm}) \equiv c(660 \text{ nm}) - c_w(660 \text{ nm})$$

is determined by the nature of the suspended particulate matter. The particulate beam attenuation $c_p(660 \text{ nm})$ is highly correlated with total particle volume concentration (usually expressed in parts per million), but it is much less well correlated with chlorophyll concentration.⁶⁷ The particulate beam attenuation can be used to estimate the total particulate load (often expressed as g m^{-3}).⁶⁸ However, the dependence of the particulate beam attenuation on particle properties is not simple. Spinrad⁶⁹ used Mie theory to calculate the dependence of the volume-specific particulate beam attenuation (particulate beam attenuation coefficient c_p in m^{-1} per unit suspended particulate volume in parts per million) on the relative refractive index and on the slope s of an assumed Junge size distribution for particles in the size range from 1 to 80 μm ; the result is shown in Fig. 16. Although the details of the figure are sensitive to the choice of upper and lower size limits in the Mie calculations, the qualitative behavior of the curves is generally valid and supports the statements made in the closing paragraph of Sec. 1.19.

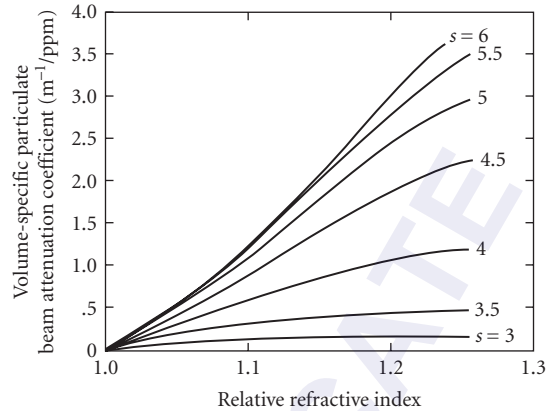


FIGURE 16 Computed relationship between volume-specific particulate beam attenuation coefficient, relative refractive index, and slope s of a Junge number size distribution. (Reproduced from Spinrad,⁶⁹ with permission.)

Because of the complicated dependence of scattering and hence of beam attenuation on particle properties, the construction of bio-optical models for $c(\lambda)$ is not easy. The reason is that chlorophyll concentration alone is not sufficient to parametrize scattering.⁷⁰ Figure 17 illustrates this insufficiency. The figure plots vertical profiles of $c(665 \text{ nm})$, water density (proportional to the oceanographic variable σ_t), and chlorophyll concentration (proportional to fluorescence by chlorophyll and related pigments). Note that the maximum in beam attenuation at 46 m depth coincides with the interface (pycnocline) between less dense water above and more dense water below. Peaks in beam attenuation are commonly observed at density interfaces because particle concentrations are often greatest there. The maximum in chlorophyll concentration occurs at a depth of 87 m. The chlorophyll concentration depends not just on the number or volume of chlorophyll-bearing particles but also on their photoadaptive state, which depends on nutrient availability and ambient lighting. Thus chlorophyll concentration cannot be expected to correlate well with total scattering or with particulate beam attenuation $c_p(\lambda)$.

Voss⁷¹ has developed an empirical model for $c(\lambda)$ given a measurement of c at $\lambda = 490 \text{ nm}$:

$$c(\lambda) = c_w(\lambda) + [c(490 \text{ nm}) - c_w(490 \text{ nm})][1.563 - 1.149 \times 10^{-3} \lambda]$$

where λ is in nm and c is in m^{-1} . The attenuation coefficient for pure sea water, $c_w = a_w + b_w$, is given by the Smith-Baker data of Table 6. This model was statistically derived from data of global extent. Testing of the model with independent data usually gave errors of less than 5 percent, although occasional errors of ~ 20 percent were found.

Voss also determined a least-squares fit of $c(490 \text{ nm})$ to the chlorophyll concentration. The result

$$c(490 \text{ nm}) = 0.39C^{0.57} \quad (7)$$

is similar in form to the chlorophyll dependence of the $a(\lambda)$ and $b(\lambda)$ models seen in Eqs. (1) and (6), respectively. Figure 18 shows the spread of the data points used to determine Eq. (7). Note that for a given value of C there is an order-of-magnitude spread in values of $c(490 \text{ nm})$. The user of Eq. (7) or of the models for $b(\lambda)$ must always keep in mind that large deviations from the predicted values will be found in natural waters.

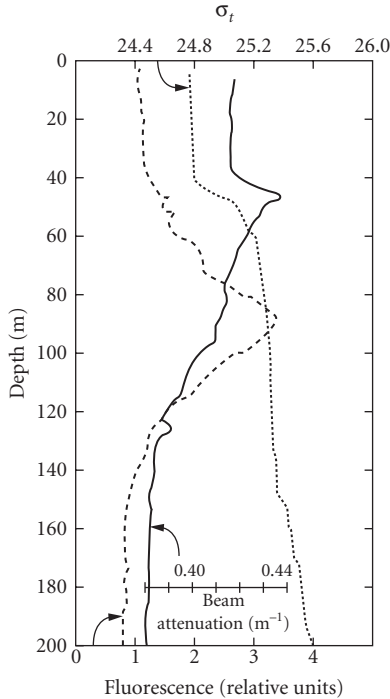


FIGURE 17 Example from the Pacific Ocean water of the depth dependence of beam attenuation (solid line), water density (σ_t , dashed line), and chlorophyll concentration (fluorescence, dotted line). (Reproduced from Kitchen and Zaneveld,⁷⁰ with permission.)

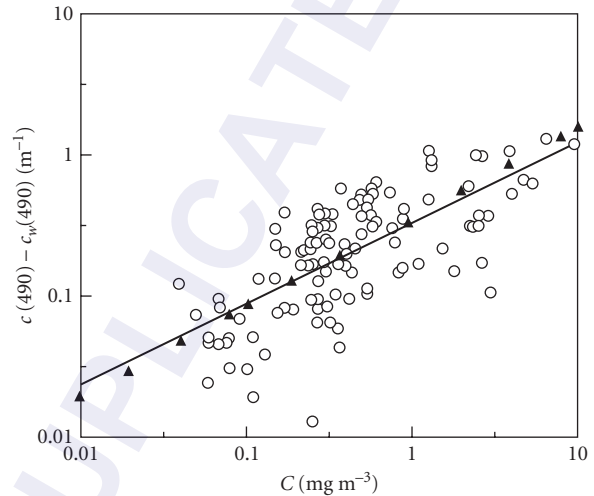


FIGURE 18 Particulate beam attenuation at 490 nm (open circles) as a function of chlorophyll concentration C as used to determine Eq. (7) which is given by the solid line. Solid triangles give values as predicted by the sum of Eqs. (1) and (6). (Redrawn from Voss,⁷¹ with permission.)

1.21 DIFFUSE ATTENUATION AND JERLOV WATER TYPES

As seen in Fig. 1 and in Table 1 there is a so-called diffuse attenuation coefficient for any radiometric variable. The most commonly used diffuse attenuation coefficients are those for downwelling plane irradiance, $K_d(z, \lambda)$, and for PAR, $K_{PAR}(z)$. Although the various diffuse attenuation coefficients are conceptually distinct, in practice they are often numerically similar and they all asymptotically approach a common value at great depths in homogeneous water.² The monograph by Tyler and Smith⁷² gives tabulations and plots of $E_d(z, \lambda)$, $E_u(z, \lambda)$ and the associated $K_d(z, \lambda)$, $K_u(z, \lambda)$, and $R(z, \lambda)$ measured in a variety of waters.

Observation shows that $K_d(z, \lambda)$ varies systematically with wavelength over a wide range of waters from very clear to very turbid. Moreover, $K_d(z, \lambda)$ is often rather insensitive to environmental effects⁷³ except for extreme conditions⁷⁴ (such as the sun within 10° of the horizon) and in most cases correction can be made¹¹ for the environmental effects that are present in K_d . K_d therefore is regarded as a quasi-inherent optical property whose variability is governed primarily by changes in the inherent optical properties of the water body and not by changes in the external environment.

Jerlov²⁹ exploited this benign behavior of K_d to develop a frequently used classification scheme for oceanic waters based on the spectral shape of K_d . The *Jerlov water types* are in essence a classification based on water clarity as quantified by $K_d(z_s, \lambda)$, where z_s is a depth just below the sea surface.

This classification scheme can be contrasted with the case 1 and case 2 classification described earlier, which is based on the nature of the suspended matter within the water. The Jerlov water types are numbered I, IA, IB, II, and III for open ocean waters, and 1 through 9 for coastal waters. Type I is the clearest and type III is the most turbid open ocean water. Likewise, for coastal waters type 1 is clearest and type 9 is most turbid. The Jerlov types I to III generally correspond to case 1 water since phytoplankton predominate in the open ocean. Types 1 to 9 correspond to case 2 waters where yellow matter and terrigenous particulates dominate the optical properties. A rough correspondence between chlorophyll concentration and Jerlov oceanic water type is given by⁴⁵

$$C: \quad 0-0.01 \sim 0.05 \sim 0.1 \sim 0.5 \sim 1.5-2.0 \text{ mg m}^{-3}$$

$$\text{water type:} \quad \text{I} \quad \text{IA} \quad \text{IB} \quad \text{II} \quad \text{III}$$

Austin and Petzold⁷⁵ reevaluated the Jerlov classification using an expanded database and slightly revised the $K_d(\lambda)$ values used by Jerlov in his original definition of the water types. Table 17 gives the revised values for $K_d(\lambda)$ for the water types commonly encountered in oceanography. These values are recommended over those found in Jerlov.²⁹ Figure 19 shows the percent transmittance of $E_d(\lambda)$ per meter of water for selected Jerlov water types. Note how the wavelength of maximum transmittance shifts from blue in the clearest open ocean water (type I) to green (types III and 1) to yellow in the most turbid, yellow-matter-rich coastal water (type 9).

Austin and Petzold also presented a simple model that allows the determination of $K_d(\lambda)$ at all wavelengths from a value of K_d measured at any single wavelength. This model is defined by

$$K_d(\lambda) = \frac{M(\lambda)}{M(\lambda_0)} [K_d(\lambda_0) - K_{\text{dw}}(\lambda_0)] + K_{\text{dw}}(\lambda)$$

Here λ_0 is the wavelength at which K_d is measured and K_{dw} refers to values for pure sea water. $K_{\text{dw}}(\lambda)$ and the statistically derived coefficients $M(\lambda)$ are given in Table 18. (These K_{dw} values differ slightly from those seen in Table 6.) This model is valid in waters where $K_d(490) \leq 0.16 \text{ m}^{-1}$ which corresponds to a chlorophyll concentration of $C \leq 3 \text{ mg m}^{-3}$.

TABLE 17 Downwelling Irradiance Diffuse Attenuation Coefficients $K_d(\lambda)$ Used to Define the Jerlov Water Types As Determined by Austin and Petzold* (All quantities in the body of the table have units of m^{-1} .)

λ (nm)	Jerlov water type					
	I	IA	IB	II	III	1
350	0.0510	0.0632	0.0782	0.1325	0.2335	0.3345
375	0.0302	0.0412	0.0546	0.1031	0.1935	0.2839
400	0.0217	0.0316	0.0438	0.0878	0.1697	0.2516
425	0.0185	0.0280	0.0395	0.0814	0.1594	0.2374
450	0.0176	0.0257	0.0355	0.0714	0.1381	0.2048
475	0.0184	0.0250	0.0330	0.0620	0.1160	0.1700
500	0.0280	0.0332	0.0396	0.0627	0.1056	0.1486
525	0.0504	0.0545	0.0596	0.0779	0.1120	0.1461
550	0.0640	0.0674	0.0715	0.0863	0.1139	0.1415
575	0.0931	0.0960	0.0995	0.1122	0.1359	0.1596
600	0.2408	0.2437	0.2471	0.2595	0.2826	0.3057
625	0.3174	0.3206	0.3245	0.3389	0.3655	0.3922
650	0.3559	0.3601	0.3652	0.3837	0.4181	0.4525
675	0.4372	0.4410	0.4457	0.4626	0.4942	0.5257
700	0.6513	0.6530	0.6550	0.6623	0.6760	0.6896

*Reproduced from Austin and Petzold,⁷⁵ with permission.

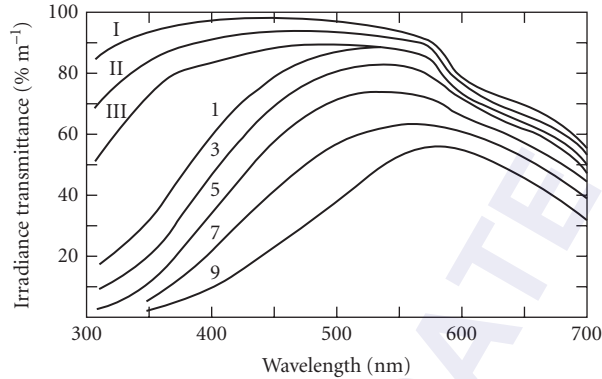


FIGURE 19 Percentage transmittance per meter of water of downwelling irradiance E_d as a function of wavelength for selected Jerlov water types. (Reproduced from Jerlov,²⁹ with permission.)

Unlike the beam attenuation coefficient $c(\lambda)$, the diffuse attenuation $K_d(z, \lambda)$ is highly correlated with chlorophyll concentration. The reason is seen in the approximate formula¹¹

$$K_d(\lambda) \approx \frac{a(\lambda) + b_b(\lambda)}{\cos \theta_{sw}}$$

where θ_{sw} is the solar angle measured within the water. Since $a(\lambda) \gg b_b(\lambda)$ for most waters, $K_d(\lambda)$ is largely determined by the absorption properties of the water, which are fairly well parametrized by the chlorophyll concentration. Beam attenuation on the other hand is proportional to the total scattering which is not well parametrized by chlorophyll concentration. Observations show⁷⁶ that the beam attenuation at 660 nm is not in general correlated with diffuse attenuation.

A bio-optical model for $K_d(\lambda)$ is given by Morel:⁴⁵

$$K_d(\lambda) = K_{dw}(\lambda) + \chi(\lambda)C^{e(\lambda)}$$

Here $K_{dw}(\lambda)$ is the diffuse attenuation for pure sea water, and $\chi(\lambda)$ and $e(\lambda)$ are statistically derived functions that convert the chlorophyll concentration C in mg m^{-3} into K_d values in m^{-1} . Table 19

TABLE 18 Values of the Coefficient $M(\lambda)$ and of the Downwelling Diffuse Attenuation Coefficient for Pure Sea Water, $K_{dw}(\lambda)$, for Use in the Austin and Petzold Model for $K_d(\lambda)$ *

λ (nm)	M (m^{-1})	K_{dw} (m^{-1})	λ (nm)	M (m^{-1})	K_{dw} (m^{-1})	λ (nm)	M (m^{-1})	K_{dw} (m^{-1})
350	2.1442	0.0510	470	1.1982	0.0179	590	0.4840	0.1578
360	2.0504	0.0405	480	1.0955	0.0193	600	0.4903	0.2409
370	1.9610	0.0331	490	1.0000	0.0224	610	0.5090	0.2892
380	1.8772	0.0278	500	0.9118	0.0280	620	0.5380	0.3124
390	1.8009	0.0242	510	0.8310	0.0369	630	0.6231	0.3296
400	1.7383	0.0217	520	0.7578	0.0498	640	0.7001	0.3290
410	1.7591	0.0200	530	0.6924	0.0526	540	0.7300	0.3559
420	1.6974	0.0189	540	0.6350	0.0577	660	0.7301	0.4105
430	1.6108	0.0182	550	0.5860	0.0640	670	0.7008	0.4278
440	1.5169	0.0178	560	0.5457	0.0723	680	0.6245	0.4521
450	1.4158	0.0176	570	0.5146	0.0842	690	0.4901	0.5116
460	1.3077	0.0176	580	0.4935	1.1065	700	0.2891	0.6514

*Condensed with permission from Austin and Petzold,⁷⁵ who give values every 5 nm.

TABLE 19 Values of the Coefficients $\chi(\lambda)$ and $e(\lambda)$ and of the Downwelling Diffuse Attenuation Coefficient for Pure Sea Water, $K_{dw}(\lambda)$, for Use in the Morel Model for $K_d(\lambda)$ *

λ (nm)	$\chi(\lambda)$	$e(\lambda)$	$K_{dw}(\lambda)$ (m^{-1})	λ (nm)	$\chi(\lambda)$	$e(\lambda)$	$K_{dw}(\lambda)$ (m^{-1})
400	0.1100	0.668	0.0209	550	0.0410	0.650	0.0640
410	0.1125	0.680	0.0196	560	0.0390	0.640	0.0717
420	0.1126	0.693	0.0183	570	0.0360	0.623	0.0807
430	0.1078	0.707	0.0171	580	0.0330	0.610	0.1070
440	0.1041	0.707	0.0168	590	0.0325	0.618	0.1570
450	0.0971	0.701	0.0168	600	0.0340	0.626	0.2530
460	0.0896	0.700	0.0173	610	0.0360	0.634	0.2960
470	0.0823	0.703	0.0175	620	0.0385	0.642	0.3100
480	0.0746	0.703	0.0194	630	0.0420	0.653	0.3200
490	0.0690	0.702	0.0217	640	0.0440	0.663	0.3300
500	0.0636	0.700	0.0271	650	0.0450	0.672	0.3500
510	0.0578	0.690	0.0384	660	0.0475	0.682	0.4050
520	0.0498	0.680	0.0490	670	0.0515	0.695	0.4300
530	0.0467	0.670	0.0518	680	0.0505	0.693	0.4500
540	0.0440	0.660	0.0568	690	0.0390	0.640	0.5000
				700	0.0300	0.600	0.6500

*Condensed with permission from Morel,⁴⁵ who gives values every 5 nm.

gives the K_{dw} , χ , and e values used in the Morel model. This model is applicable to case 1 waters with $C \leq 30 \text{ mg m}^{-3}$, although the χ and e values are somewhat uncertain for $\lambda > 650 \text{ nm}$ because of sparse data available for their determination. Some feeling for the accuracy of the Morel $K_d(\lambda)$ model can be obtained from Fig. 20 which shows predicted (the line) and observed $K_d(450)$ values as a function of C . Errors can be as large as a factor of 2 in case 1 waters (dots) and can be much larger if the model is misapplied to case 2 waters (open circles). The Morel model allows the determination of $K_d(\lambda)$ if C is measured; the Austin and Petzold model determines $K_d(\lambda)$ from a measurement at one wavelength.

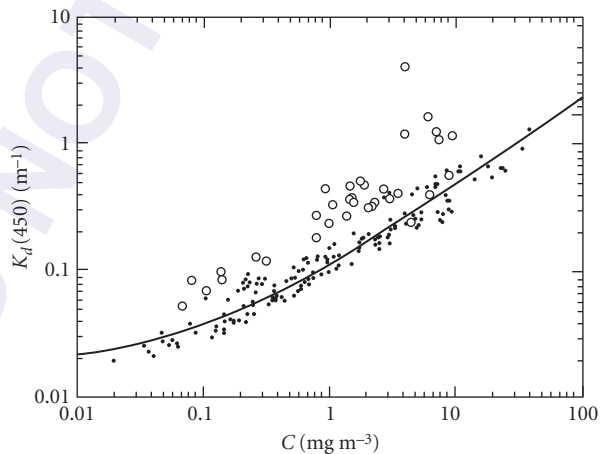


FIGURE 20 Measured K_d values at 450 nm as a function of chlorophyll concentration C . Dots are measurements from case 1 waters; open circles are from case 2 waters. The solid line gives $K_d(450)$ as predicted by the Morel bio-optical model. (Redrawn from Morel,⁴⁵ with permission.)

TABLE 20 Approximate Depth of the Euphotic Zone, z_{eu} , in Homogeneous Case 1 Water As a Function of Chlorophyll Concentration C .*

C (mg m ⁻³)	z_{eu} (m)	C (mg m ⁻³)	z_{eu} (m)
0.0	183	1	39
0.01	153	2	29
0.03	129	3	24
0.05	115	5	19
0.1	95	10	14
0.2	75	20	10
0.3	64	30	8
0.5	52		

*Data extracted from Morel,⁴⁵ with permission.

Morel⁴⁵ also presents a very simple bio-optical model for $\bar{K}_{\text{PAR}}(0, z_{\text{eu}})$ the value of $K_{\text{PAR}}(z)$ averaged over the euphotic zone $0 \leq z \leq z_{\text{eu}}$:

$$\bar{K}_{\text{PAR}}(0, z_{\text{eu}}) = 0.121C^{0.428}$$

where C is the mean chlorophyll concentration in the euphotic zone in mg m⁻³ and \bar{K}_{PAR} is in m⁻¹. The euphotic zone is the region where there is sufficient light for photosynthesis to take place; it extends roughly to the depth where $E_{\text{PAR}}(z)$ is 1 percent of its surface value [i.e., $E_{\text{PAR}}(z_{\text{eu}}) = 0.01 E_{\text{PAR}}(0)$]. Table 20 gives z_{eu} as a function of C as determined by the Morel model.

1.22 IRRADIANCE REFLECTANCE AND REMOTE SENSING

The spectral irradiance reflectant $R(\lambda) \equiv E_u(\lambda)/E_d(\lambda)$ is an important apparent optical property. Measurements of $R(z, \lambda)$ within the water have been used⁷⁷ to estimate water quality parameters such as the chlorophyll concentration, the particle backscattering coefficient, and the absorption coefficient of yellow matter. More importantly, $R(\lambda)$ just below the water surface can be related to the radiance leaving the water;⁷⁸ this radiance is available for detection by aircraft- or satellite-borne instruments. Understanding the dependence of $R(\lambda)$ upon the constituents of natural waters is therefore one of the central problems in remote sensing of water bodies.

Figure 21 illustrates the variability of $R(\lambda)$ in natural waters. Figure 21a shows $R(\lambda)$ in percent for various case 1 waters. For low-chlorophyll concentrations $R(\lambda)$ is highest at blue wavelengths, hence the blue color of clean ocean water. As the chlorophyll concentration increases, the maximum in $R(\lambda)$ shifts to green wavelengths. The enhanced reflectance near $\lambda = 685$ nm is due to chlorophyll fluorescence. Also note the exceptionally high values measured²⁰ within a coccolithophore bloom; $R(\lambda)$ is high there because of the strong scattering by the numerous calcite particles (see Sec. 1.6). Figure 21b shows $R(\lambda)$ from waters dominated by suspended sediments (i.e., by nonpigmented particles). For high-sediment concentrations $R(\lambda)$ is nearly flat from blue to yellow wavelengths, and therefore the water appears brown. Figure 21c is from waters with high concentrations of yellow substances; the peak in $R(\lambda)$ lies in the yellow. With good reason the term “ocean color” is often used as a synonym for $R(\lambda)$.

One of the main goals of oceanic remote sensing is the determination of chlorophyll concentrations in near-surface waters because of the fundamental role played by phytoplankton in the global ecosystem. Gordon et al.⁷⁸ define the *normalized water-leaving radiance* $[L_w(\lambda)]_N$ as the radiance that would leave the sea surface if the sun were at the zenith and the atmosphere were absent; this

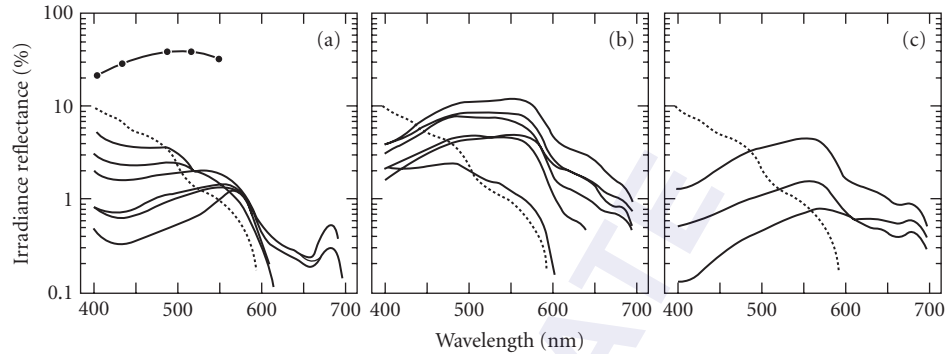


FIGURE 21 Measured spectral irradiance reflectances $R(\lambda)$ from various waters. Panel (a) is from case waters with different quantities of phytoplankton; the dotted line is $R(\lambda)$ for pure sea water. The heavy dots give values measured within a coccolithophore bloom.²⁰ Panel (b) is from case 2 waters dominated by suspended sediments and panel (c) is from case 2 waters dominated by yellow matter. (Redrawn from Sathyendranath and Morel,⁷⁹ with permission.)

quantity is fundamental to remote sensing. They then show that $[L_w]_N$ is directly proportional to R and that R is proportional to $b_b/(a + b_b)$ (i.e., to b_b/K_d). Although a or K_d are reasonably well modeled in terms of chlorophyll concentration C in case 1 waters, b_b is not well described in terms of C . Thus poor agreement is to be expected between observed values of $[L_w(\lambda)]_N$ and values predicted by a model parametrized in terms of C . This is indeed the case as is seen in Fig. 22a and b which shows observed and predicted $[L_w(443 \text{ nm})]_N$ and $[L_w(550 \text{ nm})]_N$ values as a function of chlorophyll concentration. Based on these figures there seems to be little hope of being able to reliably retrieve C from a remotely sensed $[L_w(\lambda)]_N$ value. However, in spite of the noise seen in Fig. 22a and b, the ratios of normalized water-leaving radiances for different wavelengths can be remarkably well-behaved functions of C . Figure 22c shows predicted (the line) and observed (dots) values of $[L_w(443 \text{ nm})]_N/[L_w(550 \text{ nm})]_N$; the agreement between prediction and observation is now rather good. Thus, measurement of $[L_w(\lambda)]_N$ at two (carefully chosen) wavelengths along with application of a bio-optical model for their ratio can yield a useful estimate of chlorophyll concentration. Such models are the basis of much remote sensing.

1.23 INELASTIC SCATTERING AND POLARIZATION

Although the basic physics of inelastic scattering and polarization is well understood, only recently has it become computationally practicable to incorporate these effects into predictive numerical models of underwater radiance distributions. For this reason as well as because of the difficulty of making needed measurements, quantitative knowledge about the significance of inelastic scattering and polarization in the underwater environment is incomplete.

Inelastic scattering processes are often negligible in comparison to sunlight or artificial lights as sources of underwater light of a given wavelength. However, in certain circumstances transpectral scatter is the dominant source of underwater light at some wavelengths. Figure 23 illustrates just such a circumstance.

Figure 23 shows a measured depth profile from the Sargasso Sea of irradiance reflectance at $\lambda = 589 \text{ nm}$ (yellow-orange light); note that $R(589 \text{ nm})$ increases with depth. Because of the fairly high absorption of water at this wavelength [$a_w(589 \text{ nm}) = 0.152 \text{ m}^{-1}$] most of the yellow-orange component of the incident solar radiation is absorbed near the surface and monochromatic radiative transfer theory shows that the reflectance should approach a value of $R(589 \text{ nm}) \approx 0.04$ for the

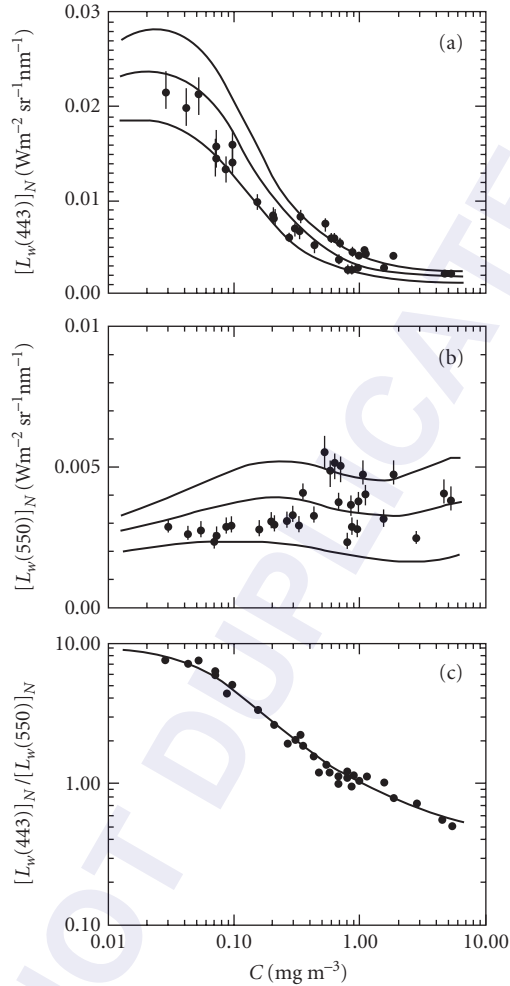


FIGURE 22 In panels (a) and (b) the solid lines are values of the normalized water-leaving radiances $[L_w(\lambda)]_N$ at $\lambda = 443$ and 550 nm, respectively, as predicted by various models that relate $[L_w(\lambda)]_N$ to the chlorophyll concentration C . The dots are measured values of $[L_w(\lambda)]_N$. Panel (c) shows the predicted (line) and observed ratio of the $[L_w(\lambda)]_N$ values of panels (a) and (b). (Redrawn from Gordon et al.,⁷⁸ with permission.)

water body of Fig. 22. Calculations by Marshall and Smith⁸⁰ explain the paradox. Light of blue-green wavelengths ($\lambda \sim 500$ nm) can penetrate to great depth in the clear Sargasso Sea water [$a_w(500 \text{ nm}) \approx 0.026 \text{ m}^{-1}$]. Some of this light is then Raman scattered from blue-green to yellow-orange wavelengths providing a source of yellow-orange light at depth. Moreover, since the phase function for Raman scattering is symmetric in the forward and backward hemispheres Raman scattered photons are equally likely to be heading upward [and thus contribute to $E_u(589 \text{ nm})$] or downward [and thus contribute to $E_d(589 \text{ nm})$] even though most of the blue-green light at depth is heading downward [e.g., $E_d(500 \text{ nm}) \gg E_u(500 \text{ nm})$].

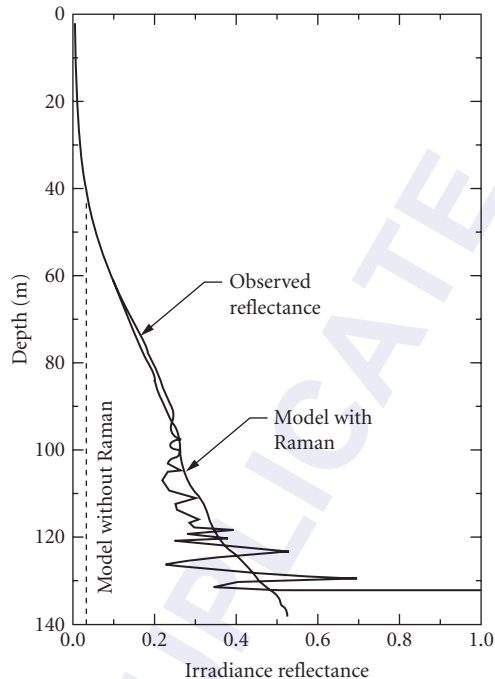


FIGURE 23 Observed irradiance reflectance R at 589 nm (light line) and values predicted by a model including Raman scattering (heavy line) and omitting Raman scattering (dashed line). (Redrawn from Marshall and Smith,⁸⁰ with permission.)

Thus as depth increases and Raman scattering becomes increasingly important relative to transmitted sunlight as a source of ambient yellow-orange light, $E_u(589 \text{ nm})$ and $E_d(589 \text{ nm})$ become more nearly equal and the irradiance reflectance $R(589 \text{ nm})$ increases. Such an increase is not seen at blue-green wavelengths since E_d transmitted from the surface remains much greater than E_u at great depths. Since Raman scattering is by the water molecules themselves this process is present (and indeed relatively more important) even in the clearest waters. Another oceanographic effect of Raman scattering occurs in the filling of Fraunhofer lines in the solar spectrum as seen underwater; this matter is just now coming under detailed investigation.⁸¹

Fluorescence by chlorophyll or other substances can be significant if the fluorescing material is present in sufficient quantity. Chlorophyll fluoresces strongly near $\lambda = 685 \text{ nm}$; this source of red light is responsible⁸² for the enhanced reflectance near 685 nm noted in Fig. 21a. The spectral signature of fluorescence is a useful tool for analyzing many of the constituents of natural waters.⁸³

Relatively little attention has been paid to the state of polarization of underwater light fields.⁸⁴ Some use of polarized light has been made in enhancing underwater visibility⁸⁵ and it is well established that many oceanic organisms sense polarized light when navigating.⁸⁶ Voss and Fry⁸⁷ measured the Mueller matrix for ocean water and Quinby-Hunt et al.⁸⁸ have studied the propensity of certain phytoplankton to induce circular polarization in unpolarized or linearly polarized light. Kattawar and Adams⁸⁹ have shown that errors of up to 15 percent can occur in calculations of underwater radiance if scalar (unpolarized) radiative transfer theory is used instead of vector (polarized) theory.

1.24 ACKNOWLEDGMENTS

This paper was written while the author held a National Research Council Resident Research Associateship (Senior Level) at the Jet Propulsion Laboratory, California Institute of Technology. This associateship was supported by the Ocean Biochemistry Program at NASA Headquarters. Final proofing was supported by SRI. Karen Baker, Annick Bricaud, Howard Gordon, Richard Honey, Rodolpho Iturriaga, George Kattawar, Scott Pegau, Mary Jane Perry, Collin Roesler, Shubha Sathyendranath, Richard Spinrad, and Kenneth Voss all made helpful comments on a draft of the paper; their efforts are greatly appreciated.

1.25 REFERENCES

1. C. D. Mobley, *Light and Water: Radiative Transfer in Natural Waters*, Academic Press, San Diego, 1994, 592 pp.
2. R. W. Preisendorfer, *Hydrologic Optics*, 6 volumes, U.S. Dept. of Commerce, NOAA, Pacific Marine Environmental Lab., Seattle, 1976, 1757 pp. Available from National Technical Information Service, 5285 Port Royal Road, Springfield, Virginia 22161.
3. A. Morel and R. C. Smith, "Terminology and Units in Optical Oceanography," *Marine Geodesy* 5(4):335 (1982).
4. N. Højerslev, "A Spectral Light Absorption Meter for Measurements in the Sea," *Limnol. Oceanogr.* 20(6):1024 (1975).
5. A. Morel and R. C. Smith, "Relation between Total Quanta and Total Energy for Aquatic Photosynthesis," *Limnol. Oceanogr.* 19(4):591 (1974).
6. A. Morel, "Light and Marine Photosynthesis: A Spectral Model with Geochemical and Climatological Implications," *Prog. Oceanogr.* 26:263 (1991).
7. J. T. O. Kirk, "The Upwelling Light Stream in Natural Waters," *Limnol. Oceanogr.* 34(8):1410 (1989).
8. J. T. O. Kirk, *Light and Photosynthesis in Aquatic Ecosystems*, Cambridge Univ. Press, New York, 1983, 410 pp.
9. R. W. Preisendorfer and C. D. Mobley, "Theory of Fluorescent Irradiance Fields in Natural Waters," *J. Geophys. Res.* 93(D9):10831 (1988).
10. R. C. Smith and K. Baker, "The Bio-Optical State of Ocean Waters and Remote Sensing," *Limnol. Oceanogr.* 23(2):247 (1978).
11. H. Gordon, "Can the Lambert-Beer Law Be Applied to the Diffuse Attenuation Coefficient of Ocean Water?," *Limnol. Oceanogr.* 34(8):1389 (1989).
12. C. A. Suttle, A. M. Chan, and M. T. Cottrell, "Infection of Phytoplankton by Viruses and Reduction of Primary Productivity," *Nature* 347:467 (1990).
13. I. Koike, S. Hara, K. Terauchi, and K. Kogure, "Role of Submicrometer Particles in the Ocean," *Nature* 345:242 (1990).
14. M. L. Wells and E. D. Goldberg, "Occurrence of Small Colloids in Sea Water," *Nature* 353:342 (1991).
15. R. W. Spinrad, H. Glover, B. B. Ward, L. A. Codispoti, and G. Kullenberg, "Suspended Particle and Bacterial Maxima in Peruvian Coastal Water during a Cold Water Anomaly," *Deep-Sea Res.* 36(5):715 (1989).
16. A. Morel and Y.-H. Ahn, "Optical Efficiency Factors of Free-Living Marine Bacteria: Influence of Bacterioplankton upon the Optical Properties and Particulate Organic Carbon in Oceanic Waters," *J. Marine Res.* 48:145 (1990).
17. D. Stramski and D. A. Kiefer, "Light Scattering by Microorganisms in the Open Ocean," *Prog. Oceanogr.* 28:343 (1991).
18. A. Alldredge and M. W. Silver, "Characteristics, Dynamics and Significance of Marine Snow," *Prog. Oceanogr.* 20:41 (1988).
19. K. L. Carder, R. G. Steward, P. R. Betzer, D. L. Johnson, and J. M. Prospero, "Dynamics and Composition of Particles from an Aeolian Input Event to the Sargasso Sea," *J. Geophys. Res.* 91(D1):1055 (1986).
20. W. M. Balch, P. M. Holligan, S. G. Ackleson, and K. J. Voss, "Biological and Optical Properties of Mesoscale Coccolithophore Blooms in the Gulf of Maine," *Limnol. Oceanogr.* 36(4):629 (1991).

21. H. Bader, "The Hyperbolic Distribution of Particle Sizes," *J. Geophys. Res.* **75**(15):2822 (1970).
22. I. N. McCave, "Particulate Size Spectra, Behavior, and Origin of Nepheloid Layers over the Nova Scotian Continental Rise," *J. Geophys. Res.* **88**(C12):7647 (1983).
23. C. E. Lambert, C. Jehanno, N. Silverberg, J. C. Brun-Cottan, and R. Chesselet, "Log-Normal Distributions of Suspended Particles in the Open Ocean," *J. Marine Res.* **39**(1):77 (1981).
24. D. G. Archer and P. Wang, "The Dielectric Constant of Water and Debye-Hückel Limiting Law Slopes," *J. Phys. Chem. Ref. Data* **19**:371 (1990).
25. M. Kerker, *The Scattering of Light and Other Electromagnetic Radiation*, Academic Press, New York, 1969, 666 pp.
26. V. M. Zoloratev and A. V. Demin, "Optical Constants of Water over a Broad Range of Wavelengths, 0.1 Å-1 m," *Opt. Spectrosc. (U.S.S.R.)* **43**(2):157 (Aug. 1977).
27. R. W. Austin and G. Halikas, "The Index of Refraction of Seawater," SIO ref. no. 76-1, Scripps Inst. Oceanogr., San Diego, 1976, 121 pp.
28. R. C. Millard and G. Seaver, "An Index of Refraction Algorithm over Temperature, Pressure, Salinity, Density, and Wavelength," *Deep-Sea Res.* **37**(12):1909 (1990).
29. N. G. Jerlov, *Marine Optics*, Elsevier, Amsterdam, 1976, 231 pp.
30. R. W. Spinrad and J. F. Brown, "Relative Real Refractive Index of Marine Micro-Organisms: A Technique for Flow Cytometric Measurement," *Appl. Optics* **25**(2):1930 (1986).
31. S. G. Ackleson, R. W. Spinrad, C. M. Yentsch, J. Brown, and W. Korjef-Bellows, "Phytoplankton Optical Properties: Flow Cytometric Examinations of Dilution-Induced Effects," *Appl. Optics* **27**(7): 1262 (1988).
32. R. C. Smith and K. S. Baker, "Optical Properties of the Clearest Natural Waters," *Appl. Optics* **20**(2):177 (1981).
33. T. T. Bannister, "Estimation of Absorption Coefficients of Scattering Suspensions Using Opal Glass," *Limnol. Oceanogr.* **33**(4, part 1):607 (1988).
34. B. G. Mitchell, "Algorithms for Determining the Absorption Coefficient of Aquatic Particles Using the Quantitative Filter Technique (QFT)," *Ocean Optics X*, R. W. Spinrad (ed.), *Proc. SPIE* **1302**:137 (1990).
35. D. Stramski, "Artifacts in Measuring Absorption Spectra of Phytoplankton Collected on a Filter," *Limnol. Oceanogr.* **35**(8):1804 (1990).
36. J. R. V. Zaneveld, R. Bartz, and J. C. Kitchen, "Reflective-Tube Absorption Meter," *Ocean Optics X*, R. W. Spinrad (ed.), *Proc. SPIE* **1302**:124 (1990).
37. E. S. Fry, G. W. Kattawar, and R. M. Pope, "Integrating Cavity Absorption Meter," *Appl. Optics* **31**(12):2025 (1992).
38. W. Doss and W. Wells, "Radiometer for Light in the Sea," *Ocean Optics X*, R. W. Spinrad (ed.), *Proc. SPIE* **1302**:363 (1990).
39. K. J. Voss, "Use of the Radiance Distribution to Measure the Optical Absorption Coefficient in the Ocean," *Limnol. Oceanogr.* **34**(8):1614 (1989).
40. F. Sogandares, Z.-F. Qi, and E. S. Fry, "Spectral Absorption of Water," presentation at the Optical Society of America Annual Meeting, San Jose, Calif., 1991.
41. W. S. Pegau and J. R. V. Zaneveld, "Temperature Dependent Absorption of Water in the Red and Near Infrared Portions of the Spectrum," *Limnol. Oceanogr.* **38**(1):188 (1993).
42. A. Bricaud, A. Morel, and L. Prieur, "Absorption by Dissolved Organic Matter of the Sea (Yellow Substance) in the UV and Visible Domains," *Limnol. Oceanogr.* **26**(1):43 (1981).
43. C. S. Roesler, M. J. Perry, and K. L. Carder, "Modeling in situ Phytoplankton Absorption from Total Absorption Spectra in Productive Inland Marine Waters," *Limnol. Oceanogr.* **34**(8):1510 (1989).
44. S. Sathyendranath, L. Lazzara, and L. Prieur, "Variations in the Spectral Values of Specific Absorption of Phytoplankton," *Limnol. Oceanogr.* **32**(2):403 (1987).
45. A. Morel, "Optical Modeling of the Upper Ocean in Relation to Its Biogenous Matter Content (Case 1 Waters)," *J. Geophys. Res.* **93**(C9):10749 (1988).
46. R. Iturriaga and D. Siegel, "Microphotometric Characterization of Phytoplankton and Detrital Absorption Properties in the Sargasso Sea," *Limnol. Oceanogr.* **34**(8):1706 (1989).
47. A. Morel and L. Prieur, "Analysis of Variations in Ocean Color," *Limnol. Oceanogr.* **22**(4):709 (1977).
48. L. Prieur and S. Sathyendranath, "An Optical Classification of Coastal and Oceanic Waters Based on the Specific Spectral Absorption Curves of Phytoplankton Pigments, Dissolved Organic Matter, and Other Particulate Materials," *Limnol. Oceanogr.* **26**(4):671 (1981).

49. O. V. Kopelevich, "Small-Parameter Model of Optical Properties of Sea Water," *Ocean Optics*, vol 1, *Physical Ocean Optics*, A. S. Monin (ed.), Nauka Pub., Moscow, 1983, chap. 8 (in Russian).
50. V. I. Haltrin and G. Kattawar, "Light Fields with Raman Scattering and Fluorescence in Sea Water," Tech. Rept., Dept. of Physics, Texas A&M Univ., College Station, 1991, 74 pp.
51. C. S. Yentsch, "The Influence of Phytoplankton Pigments on the Color of Sea Water," *Deep-Sea Res.* 7:1 (1960).
52. H. Gordon, "Diffuse Reflectance of the Ocean: Influence of Nonuniform Pigment Profile," *Appl. Optics* 31(12):2116 (1992).
53. T. J. Petzold, "Volume Scattering Functions for Selected Ocean Waters," SIO Ref. 72-78, Scripps Inst. Oceanogr., La Jolla, 1972 (79 pp). Condensed in *Light in the Sea*, J. E. Tyler (ed.), Dowden, Hutchinson & Ross, Stroudsburg, 1977, chap. 12, pp. 150-174.
54. R. W. Spinrad, J. R. V. Zaneveld, and H. Pak, "Volume Scattering Function of Suspended Particulate Matter at Near-Forward Angles: A Comparison of Experimental and Theoretical Values," *Appl. Optics* 17(7):1125 (1978).
55. G. G. Padmabandu and E. S. Fry, "Measurement of Light Scattering at 0° by Small Particle Suspensions," *Ocean Optics X*, R. W. Spinrad (ed.), *Proc. SPIE* 1302:191 (1990).
56. Y. Kuga and A. Ishimaru, "Backscattering Enhancement by Randomly Distributed Very Large Particles," *Appl. Optics* 28(11):2165 (1989).
57. A. Morel, "Optical Properties of Pure Water and Pure Sea Water," *Optical Aspects of Oceanography*, N. G. Jerlov and E. S. Nielsen (eds.), Academic Press, New York, 1974, chap. 1, pp 1-24.
58. K. S. Shifrin, *Physical Optics of Ocean Water*, AIP Translation Series, Amer. Inst. Physics, New York, 1988, 285 pp.
59. G. Kullenberg, "Observed and Computed Scattering Functions," *Optical Aspects of Oceanography*, N. G. Jerlov and E. S. Nielsen (eds.), Academic Press, New York, 1974, chap. 2, pp 25-49.
60. O. B. Brown and H. R. Gordon, "Size-Refractive Index Distribution of Clear Coastal Water Particulates from Scattering," *Appl. Optics* 13:2874 (1974).
61. J. C. Kitchen and J. R. V. Zaneveld, "A Three-Layer Sphere, Mie-Scattering Model of Oceanic Phytoplankton Populations," presented at Amer. Geophys. Union/Amer. Soc. Limnol. Oceanogr. Annual Meeting, New Orleans, 1990.
62. C. D. Mobley, B. Gentili, H. R. Gordon, Z. Jin, G. W. Kattawar, A. Morel, P. Reinersman, K. Stamnes, and R. H. Starn, "Comparison of Numerical Models for Computing Underwater Light Fields," *Appl. Optics* 32(36):7484 (1993).
63. A. Morel and B. Gentili, "Diffuse Reflectance of Ocean Waters: Its Dependence on Sun Angle As Influenced by the Molecular Scattering Contribution," *Appl. Optics* 30(30):4427 (1991).
64. A. Morel, "Diffusion de la lumière par les eaux de mer. Résultats expérimentaux et approach théorique," NATO AGARD lecture series no. 61, *Optics of the Sea*, chap. 3.1, pp. 1-76 (1973), G. Halikas (trans.), Scripps Inst. Oceanogr., La Jolla, 1975, 161 pp.
65. O. V. Kopelevich and E. M. Mezhericher, "Calculation of Spectral Characteristics of Light Scattering by Sea Water," *Izvestiya, Atmos. Oceanic Phys.* 19(2):144 (1983).
66. H. R. Gordon and A. Morel, "Remote Assessment of Ocean Color for Interpretation of Satellite Visible Imagery, A Review," *Lecture Notes on Coastal and Estuarine Studies*, vol. 4, Springer-Verlag, New York, 1983, 114 pp.
67. J. C. Kitchen, J. R. V. Zaneveld, and H. Pak, "Effect of Particle Size Distribution and Chlorophyll Content on Beam Attenuation Spectra," *Appl. Optics* 21(21):3913 (1982).
68. J. K. Bishop, "The Correction and Suspended Particulate Matter Calibration of Sea Tech Transmissometer Data," *Deep-Sea Res.* 33:121 (1986).
69. R. W. Spinrad, "A Calibration Diagram of Specific Beam Attenuation," *J. Geophys. Res.* 91(C6):7761 (1986).
70. J. C. Kitchen and J. R. V. Zaneveld, "On the Noncorrelation of the Vertical Structure of Light Scattering and Chlorophyll *a* in Case 1 Water," *J. Geophys. Res.* 95(C11):20237 (1990).
71. K. J. Voss, "A Spectral Model of the Beam Attenuation Coefficient in the Ocean and Coastal Areas," *Limnol. Oceanogr.* 37(3):501 (1992).

72. J. E. Tyler and R. C. Smith, *Measurements of Spectral Irradiance Underwater*, Gordon and Breach, New York, 1970, 103 pp.
73. K. S. Baker and R. C. Smith, "Quasi-Inherent Characteristics of the Diffuse Attenuation Coefficient for Irradiance," *Ocean Optics VI*, S. Q. Duntley (ed.), *Proc. SPIE* **208**:60 (1979).
74. C. D. Mobley, "A Numerical Model for the Computation of Radiance Distributions in Natural Waters with Wind-Blown Surfaces," *Limnol. Oceanogr.* **34**(8):1473 (1989).
75. R. W. Austin and T. J. Petzold, "Spectral Dependence of the Diffuse Attenuation Coefficient of Light in Ocean Water," *Opt. Eng.* **25**(3):471 (1986).
76. D. A. Siegel and T. D. Dickey, "Observations of the Vertical Structure of the Diffuse Attenuation Coefficient Spectrum," *Deep-Sea Res.* **34**(4):547 (1987).
77. S. Sugihara and M. Kishino, "An Algorithm for Estimating the Water Quality Parameters from Irradiance Just below the Sea Surface," *J. Geophys. Res.* **93**(D9):10857 (1988).
78. H. R. Gordon, O. B. Brown, R. E. Evans, J. W. Brown, R. C. Smith, K. S. Baker, and D. C. Clark, "A Semianalytic Model of Ocean Color," *J. Geophys. Res.* **93**(D9):10909 (1988).
79. S. Sathyendranath and A. Morel, "Light Emerging from the Sea—Interpretation and Uses in Remote Sensing," *Remote Sensing Applications in Marine Science and Technology*, A. P. Cracknell (ed.), D. Reidel, Dordrecht, 1983, chap. 16, pp. 323–357.
80. B. R. Marshall and R. C. Smith, "Raman Scattering and In-Water Ocean Optical Properties," *Appl. Optics* **29**:71 (1990).
81. G. W. Kattawar and X. Xu, "Filling-in of Fraunhofer Lines in the Ocean by Raman Scattering," *Appl. Optics* **31**(30):6491 (1992).
82. H. Gordon, "Diffuse Reflectance of the Ocean: The Theory of Its Augmentation by Chlorophyll *a* Fluorescence at 685 nm," *Appl. Optics* **18**:1161 (1979).
83. J. J. Cullen, C. M. Yentsch, T. L. Cucci, and H. L. MacIntyre, "Autofluorescence and Other Optical Properties As Tools in Biological Oceanography," *Ocean Optics IX*, M. A. Blizard (ed.), *Proc. SPIE* **925**:149 (1988).
84. A. Ivanoff, "Polarization Measurements in the Sea," *Optical Aspects of Oceanography*, N. G. Jerlov and E. S. Nielsen (eds.), Academic Press, New York, 1974, chap. 8, pp. 151–175.
85. G. D. Gilbert and J. C. Pernicka, "Improvement of Underwater Visibility by Reduction of Backscatter with a Circular Polarization Technique," *SPIE Underwater Photo-Optics Seminar Proc*, Santa Barbara, Oct. 1966.
86. T. H. Waterman, "Polarization of Marine Light Fields and Animal Orientation," *Ocean Optics X*, M. A. Blizard (ed.), *Proc. SPIE* **925**:431 (1988).
87. K. J. Voss and E. S. Fry, "Measurement of the Mueller Matrix for Ocean Water," *Appl. Optics* **23**:4427 (1984).
88. M. S. Quinby-Hunt, A. J. Hunt, K. Lofftus, and D. Shapiro, "Polarized Light Studies of Marine *Chlorella*," *Limnol. Oceanogr.* **34**(8):1589 (1989).
89. G. W. Kattawar and C. N. Adams, "Stokes Vector Calculations of the Submarine Light Field in an Atmosphere-Ocean with Scattering According to a Rayleigh Phase Matrix: Effect of Interface Refractive Index on Radiance and Polarization," *Limnol. Oceanogr.* **34**(8):1453 (1989).

This page intentionally left blank.

DO NOT DUPLICATE

PROPERTIES OF CRYSTALS AND GLASSES

William J. Tropf and Michael E. Thomas

*Applied Physics Laboratory
Johns Hopkins University
Laurel, Maryland*

Eric W. Rogala

*Raytheon Missile Systems
Tucson, Arizona*

2.1 GLOSSARY

A, B, C, D, E, G	dispersion equation constants
a, b	partial dispersion equation constants
a, b, c	crystal axes
B	inverse dielectric constant ($=1/\epsilon = 1/n^2$)
B	bulk modulus
C	heat capacity
c	speed of light
c	elastic stiffness tensor
D	electric displacement
D	dispersion
d	piezoelectric coefficient
$d_{ij}^{(2)}$	nonlinear optical coefficient
E	Young's modulus
E	energy
E	electric field
e	strain
G	shear modulus
G	thermal optical constant
g	degeneracy
Hi	Hilbert transform
h	heat flow
k	Index of absorption
k_B	Boltzmann constant
ℓ	phonon mean free path
L	length
MW	molecular weight
m	Integer, mass
$N()$	occupation number density
n	refractive index

\tilde{n}	complex refractive index = $n + ik$
\mathbf{P}	electric polarization
P	relative air pressure = air pressure/one atmosphere (dimensionless)
P	pyroelectric constant
$P_{x,y}$	relative partial dispersion
\mathbf{p}	elasto-optic tensor
\mathbf{q}	piezo-optic tensor
r	electro-optic coefficient
r	amplitude reflection coefficient
r_{ij}	electro-optic coefficient
$S(\)$	oscillator strength
\mathbf{s}	elastic compliance tensor
T	temperature
t	amplitude transmission coefficient
U	enthalpy
u	atomic mass unit
V	Volume
v	velocity of sound
x	displacement
x	variable of integration
Z	formulas per unit cell
α	linear thermal expansion coefficient
α	intensity (power) absorptance
α_m	macroscopic polarizability
α, β, γ	crystal angles
β	power absorption coefficient
$\gamma(\)$	line width
γ	Gruneisen parameter
ϵ	relative dielectric constant, permittivity
ϵ	emittance
θ_D	Debye temperature
κ	thermal conductivity
λ	wavelength
$\Lambda(\)$	complex function
μ	permeability
ν	wave number ($\omega/2\pi c$)
ν_d	Abbe number (constringence)
ρ	density
ρ	intensity reflectance
σ	stress
τ	intensity (power) transmittance
χ	susceptibility
$\chi^{(2)}$	second-order susceptibility
Ω	solid angle
ω	frequency

Subscripts

ABS	absorptance
abs	absolute
bb	blackbody
C	656.3 nm
d	587.6 nm
EXT	extinctance

F	486.1 nm
i	integers
P	constant pressure
r	relative
SCA	scatterance
V	constant volume
0	vacuum or constant terms (or $T = 0$)
1, 2, 3	principal axes

2.2 INTRODUCTION

Nearly every nonmetallic crystalline and glassy material has potential use in transparent optics. If a nonmetal is sufficiently dense and homogeneous, it will have good optical properties. Generally, a combination of desirable optical properties, good thermal and mechanical properties, and cost and ease of manufacture dictate the number of readily available materials. In practice, glasses dominate the available optical materials for several important reasons. Glasses are easily made of inexpensive materials, and glass manufacturing technology is mature and well-established. The resultant glass products can have very high optical quality and meet most design requirements.

Common glasses, however, are composed of low-atomic-weight oxides and therefore will not transmit beyond about 2.5 μm or below 0.3 μm . Some crystalline materials transmit at wavelengths longer (e.g., heavy-metal halides and chalcogenides) or shorter (e.g., fluorides) than common, oxide-based glasses. Crystalline materials are also used for situations that require the material to have very low scatter, high thermal conductivity, or exceptionally high hardness and strength, especially at high temperature. Other applications of crystalline optical materials make use of their directional properties, particularly those of noncubic (uni- or biaxial) crystals. Phase matching (e.g., in wave mixing) and polarization (e.g., in wave plates) are example applications. For these reasons, crystalline solids are used for a wide variety of specialized applications.

Polycrystalline materials form an intermediate class. Typically engineered from powders, polycrystalline materials approach the properties of crystals with reduced cost, particularly when made into complex shapes. Most polycrystalline materials are made with large grains ($\gg \lambda$) to reduce scatter, and the properties of crystalline materials in this chapter are obtained from samples with large grains or single crystals. Recent work has produced materials with sub-micrometer grains ($\ll \lambda$) for applications such as laser gain media and highly durable windows.

This chapter gives the physical, mechanical, thermal, and optical properties of selected crystalline, polycrystalline, and glassy materials. Crystals are chosen based on availability of property data and usefulness of the material. Unfortunately, for many materials, property data are imprecise, incomplete, or not wholly applicable to optical-quality material. Glasses are more accurately and uniformly characterized, but their optical property data are usually limited to wavelengths below 1.06 μm . Owing to the preponderance of glasses, only a small, representative fraction of available glasses are included below. SI-derived units, as commonly applied in material characterization, are used.

Property data are accompanied with brief explanations and useful functional relationships. We have updated the previous version of this chapter¹ by extracting property data from standard compilations²⁻¹² as well as recent literature. Unfortunately, property data are often sparse. For example, refractive index data may be available for only a portion of the transparent region or the temperature dependence of the index may not be known. Strength of many materials is poorly characterized. Thermal conductivity is frequently unavailable and other thermal properties are usually sketchy.

For brevity, we have reduced the number of materials in this chapter. Many seldom-used or no-longer-produced materials have been eliminated. New materials, principally used for laser gain media or wave mixing were added. Glasses have been updated to their new arsenic- and lead-free analogs.

2.3 OPTICAL MATERIALS

Crystalline and amorphous (including glass) materials are differentiated by their structural (crystallographic) order. The distinguishing structural characteristic of amorphous substances is the absence of long-range order; the distinguishing characteristic of crystals is the presence of long-range order. This order, in the form of a periodic structure, can cause directional-dependent (anisotropic) properties that require a more complex description than needed for isotropic, amorphous materials. The periodic structure of crystals are used to classify them into six crystal systems,* and further arrange them into 14 (Bravais) space lattices, 32 point groups, and 230 space groups based on the characteristic symmetries found in a crystal.

Glass is by far the most widely used optical material, accounting for more than 90 percent of all optical elements manufactured. Traditionally, glass has been the material of choice for optical systems designers, owing to its high transmittance in the visible-wavelength region, high degree of homogeneity, ease of molding, shaping, and machining, relatively low cost, and the wide variety of index and dispersion characteristics available.

Under proper conditions, glass can be formed from many different inorganic mixtures. Primary glass-forming compounds include oxides, halides, and chalcogenides with the most common mixtures being the oxides of silicon, boron, and phosphorous used for glasses transmitting in the visible spectrum. By varying the chemical composition of glasses (glasses are not fixed stoichiometrically), the properties of the glass can be modified. For optical applications, glass compositions are altered to vary the refractive index, dispersion, and thermo-optic coefficient, allowing glasses to be used in combination to make optics with excellent performance over wide spectral and temperature ranges.

Historically, emphasis was placed on optimizing optics for the visible spectrum. The approach was to produce a wide variety of glass and selecting the best combination of glasses to optimize a design. Early glass technologists found that adding BaO offered a high-refractive-index glass with lower than normal dispersion, B_2O_3 offered low index and very low dispersion, and by replacing oxides with fluorides, glasses could be obtained with very low index and very low dispersion. Later, others developed very high index glasses with relatively low dispersions by introducing rare-earth elements, especially lanthanum, to glass compositions. Other compounds are added to silica-based glass mixtures to help with chemical stabilization, typically the alkaline earth oxides and in particular Al_2O_3 to improve the resistance to attack by water.

Advances in detector and source technology pushed the design spectrum into both the UV and IR. To extend the transmission range of glasses into the ultraviolet, a number of fluoride and fluorophosphate glasses were developed. Nonoxide glasses are used for infrared applications requiring transmission beyond the transmission limit of typical optical glasses (cf., 2.4 to 2.7 μm for an absorption coefficient of $<1 \text{ cm}^{-1}$). These materials include chalcogenides such as As_2S_3 glass and heavy-metal fluorides such as ZrF_4 -based glasses.

With today's computing power, optical designs are becoming far more elegant and simple, and require fewer optical glass types. Advances in polishing and use of aspheres have introduced a new design variable, reducing the dependence on a continuous range material properties. An increased awareness of the environment has led to the elimination of many hazardous materials from optical glasses. These factors lead to a major reduction of the glass library available from manufacturers.

Crystalline materials include naturally occurring minerals and manufactured crystals. Both single crystals and polycrystalline forms are available for many materials. Polycrystalline optical materials are typically composed of many small (cf., 50 to 200 μm) individual crystals ("grains") with random orientations and *grain boundaries* between them. These grain boundaries are a form of material defect arising from the lattice mismatch between individual grains. Polycrystalline materials are made by diverse means such as pressing powders (usually with heat applied), sintering, or chemical vapor deposition. Single crystals are typically grown from dissolved or molten material using a variety of techniques. Usually, polycrystalline materials offer greater hardness and strength at the expense of increased scatter and reduced thermal conductivity.

*Cubic (or isometric), hexagonal (including rhombohedral), tetragonal, orthorhombic, monoclinic, and triclinic are the crystallographic systems.

Uniformity of the refractive index throughout an optical element is a prime consideration in selecting materials for high-performance lenses, elements for coherent optics, laser harmonic generation, and acousto-optical devices. In general, highly pure, single crystals achieve the best uniformity, followed by glasses (especially those specially manufactured for homogeneity), and lastly polycrystalline materials. Similarly, high-quality single crystals have very low scatter, typically one-tenth that of the best glasses.

Applications requiring optical elements with direction-dependent properties, such as polarizers and phase-matching materials for harmonic generation, frequently use single crystals.

Purity of starting materials is a prime factor in determining the quality of the final product. High material quality and uniformity of processing is required to avoid impurity absorption, index non-uniformity, voids, cracks, and bubbles, and excess scatter. Practical manufacturing techniques limit the size of optics produced from a given material (glasses are typically limited by the moduli, i.e., deformation caused by the weight of the piece). Some manufacturing methods, such as hot pressing, also produce significantly lower quality material when size becomes large (especially when the thinnest dimension is significantly increased). Cost of finished optical elements is a function of size, raw material cost, and the difficulty of machining, polishing, and coating the material. Any one of these factors can dominate cost.

Information on general optical glass properties and classification can be found on select manufacturer's websites¹³⁻¹⁵ and in the literature.^{16,17} General information on the manufacturing methods for glasses and crystalline materials is available in several sources.^{10,12,18,19} Information on cutting and polishing of optical elements can be found in the literature.^{3,20,21} Further information on applications of optical material are found in this section of this *Handbook* and a variety of sources.^{2,18,22,23}

2.4 PROPERTIES OF MATERIALS

Symmetry Properties

The description of the properties of solids depends on structural symmetry. The structural symmetry of crystalline materials dictates the number of appropriate directional-dependent terms that describe a property.²⁴ *Neumann's principle* states that physical properties of a crystal must possess at least the symmetry of the point group of the crystal. Amorphous (glassy) materials, having no long-range symmetry, are generally considered isotropic,* and require the least number of property terms. Macroscopic material properties are best described in tensor notation; the rank of the property tensors can range from zero (scalar property) up to large rank. Table 1 summarizes common material properties and the rank of the tensor that describes them.

The rank of the property tensor and the symmetry of the material (as determined by the point group for crystals) determine the number of terms needed to describe a property. Table 2 summarizes both the number of terms needed to describe a property and the number of unique terms as a function of property tensor rank and point group. Another important term in the definition of tensor characteristics is the *principal values* of a second-rank tensor property. The principal values of a property are those values referenced to (measured along) the crystal axes as defined in Table 2. For example, a second-rank tensor of a triclinic crystal has nine nonzero coefficients, and because of symmetry, six independent coefficients. These six coefficients can be separated into (1) three principal values of the quantity (e.g., thermal expansion coefficient, dielectric constant, refractive index, and stress), and (2) the three angles describing the orientation of the crystal axes (α , β , γ).

Properties of materials depend on several fundamental constants that are listed in Table 3.²⁵ These constants are defined as follows.

*Glass properties are dependent on cooling rate. Nonuniform cooling will result in density variation and residual stress which cause anisotropy in all properties. Such nonisotropic behavior is different from that of crystals in that it arises from thermal gradients rather than periodic structure order. Hence, the nature of anisotropy in glasses varies from sample to sample.

TABLE 1 Tensor Characteristics and Definitions of Properties

Tensor rank	Property	Symbol	Units	Relationship
0	Enthalpy (energy)	U	J/mole	—
	Temperature	T	K	—
	Heat capacity	C	J/(mole · K)	$C = \partial U / \partial T$
1	Displacement	x	m	—
	Heat flow	h	W/m ²	—
	Electric field	\mathbf{E}	V/m	—
	Electric polarization	\mathbf{P}	C/m ²	$\mathbf{P} = \epsilon_0 \chi \mathbf{E}$
	Electric displacement	\mathbf{D}	C/m ²	$\mathbf{D} = \epsilon_0 \mathbf{E} + \mathbf{P}$
	Pyroelectric constant	p	C/(m ² · K)	$\Delta \mathbf{P} = p \Delta T$
	2	Stress	σ	Pa
Strain		e	—	—
Thermal expansion		α	K ⁻¹	$e = \alpha T$
Thermal conductivity		k	W/(m · K)	$h = -k (\partial T / \partial x)$
Dielectric constant (relative permittivity)		ϵ	—	$\mathbf{D} = \epsilon_0 \epsilon_r \mathbf{E}$
Inverse dielectric tensor		\mathbf{B}	—	$\mathbf{B} = \epsilon_r^{-1}$
Susceptibility		χ	—	$\epsilon_r = \chi + 1$
3	Piezoelectric coefficient (modulus)	d	m/V \equiv C/N	$\mathbf{P} = d \cdot \sigma$
	(converse piezoelectric effect)	d	m/V \equiv C/N	$e = d \cdot \mathbf{E}$
	Electro-optic coefficient (linear)	r	m/V	$\Delta \mathbf{B} = r \cdot \mathbf{E}$
4	Second-order susceptibility	$\chi^{(2)}$	m/V	$\mathbf{P} = \epsilon_0 \chi^{(2)} \mathbf{E}_1 \mathbf{E}_2$
	Elastic stiffness	\mathbf{c}	Pa	$\sigma = \mathbf{c} \cdot \mathbf{e}; \mathbf{c} = 1/\mathbf{s}$
	Elastic compliance	\mathbf{s}	Pa ⁻¹	$e = \mathbf{s} \cdot \sigma; \mathbf{s} = 1/\mathbf{c}$
	Elasto-optic tensor	\mathbf{p}	—	$\Delta \mathbf{B} = \mathbf{p} \cdot \mathbf{e}; \mathbf{p} = \mathbf{q} \cdot \mathbf{c}$
	Piezo-optic tensor	\mathbf{q}	Pa ⁻¹	$\Delta \mathbf{B} = \mathbf{q} \cdot \sigma; \mathbf{q} = \mathbf{p} \cdot \mathbf{s}$

Optical Properties: Introduction

Refractive Index Important optical properties, definitions, formulas, and basic concepts are derived from a classical description of propagation based on the macroscopic Maxwell's equations.²² The standard wave equation for the electric field \mathbf{E} is obtained from Faraday, Gauss, and Ampere laws in Maxwell's form:*

$$\nabla^2 \mathbf{E} = (-i\omega\mu\sigma - \omega^2\mu\epsilon)\mathbf{E} \quad (1)$$

where σ , ϵ , and μ are the frequency-dependent conductivity, permittivity, and permeability, respectively. These quantities are scalars in an isotropic medium. To simplify notation, a generalized permittivity is sometimes defined as:

$$\epsilon_c(\omega) = \epsilon_r(\omega) \left[1 + i \frac{\sigma(\omega)}{\omega\epsilon(\omega)} \right] \quad (2)$$

where $\epsilon_c(\omega)$ is a generalized relative permittivity of dielectric constant (i.e., with ϵ_0 removed from $\epsilon = \epsilon_0 \epsilon_c$) that includes contributions from free charges [via the conductivity $\sigma(\omega)$] and bound charges [via the relative permittivity, $\epsilon_r(\omega)$]. Assuming a nonmagnetic material ($\mu_r = 1$), and using the preceding generalized dielectric constant, the plane-wave solution to the wave equation is

$$\mathbf{E}(z, \omega) = \mathbf{E}(0) \exp \left[i(\omega z / c) \sqrt{\epsilon_c(\omega)} \right] \quad (3)$$

*The definition of the dielectric constant and refractive index in this section is based on a harmonic field of the form $\exp(-i\omega t)$. Other definitions lead to different sign conventions (e.g., $\vec{n} = n - ik$) and care must be taken to ensure consistency.

TABLE 2 Crystal Classes and Symmetries

Crystal System	Crystal Axes	Space lattice		Point group		Space Group Nos.	Tensor coefficients*			
		Types	Symmetry	Schönflies	Internat ¹		Rank 1	Rank 2	Rank 3	Rank 4
Triclinic (—)	$a \neq b \neq c$	P	$\bar{1}$	C_1	$\bar{1}$	1	3 (3)	9 (6)	18 (18)	36 (21)
	$\alpha \neq \beta \neq \gamma$			C_1	$\bar{1}$	2	0	9 (6)	0	36 (21)
Monoclinic (2-fold axis)	$a \neq b \neq c$	P, I (or C)	2/m	C_2	2	3-5	1 (1)	5 (4)	8 (8)	20 (13)
	$\alpha = \beta = 90^\circ$			C_s	m	6-9	2 (2)	5 (4)	10 (10)	20 (13)
	$\gamma = 90^\circ$			C_{2h}	2/m	10-15	0	5 (4)	0	20 (13)
Orthorhombic (3 \perp 2-fold axes)	$a \neq b \neq c$	P, I, C, F	mmm	D_2	222	16-24	0	3 (3)	3 (3)	12 (9)
	$\alpha = \beta = \gamma = 90^\circ$			C_{2v}	2mm	25-46	1 (1)	3 (3)	5 (5)	12 (9)
				D_{2h}	mmm	47-74	0	3 (3)	0	12 (9)
				C_4	4	75-80	1 (1)	3 (2)	7 (4)	16 (7)
				C_4	$\frac{4}{2}$	81-82	0	3 (2)	7 (4)	16 (7)
				C_{2h}	4/m	83-88	0	3 (2)	0	16 (7)
				D_4	422	89-98	0	3 (2)	2 (1)	12 (6)
				C_{4v}	4mm	99-110	1 (1)	3 (2)	5 (3)	12 (6)
				D_{2d}	$\frac{4}{2}$	111-122	0	3 (2)	3 (2)	12 (6)
				D_{4h}	4/mmm	123-142	0	3 (2)	0	12 (6)
Hexagonal (3-fold axis)	$a = b \neq c$	P, R	$\bar{3}$	C_3	$\frac{3}{2}$	143-146	1 (1)	3 (2)	13 (6)	24 (7)
	$\alpha = \beta = 90^\circ$			C_{3i}	$\frac{3}{2}$	147-148	0	3 (2)	0	24 (7)
	$\gamma = 120^\circ$			D_3	32	149-155	0	3 (2)	5 (2)	18 (6)
				D_{3d}	$\frac{3}{2}$	156-161	1 (1)	3 (2)	8 (4)	18 (6)
				D_{3h}	$\frac{3}{2}$	162-167	0	3 (2)	0	18 (6)
				C_6	$\frac{6}{2}$	168-173	1 (1)	3 (2)	7 (4)	12 (5)
				C_{3h}	6	174	0	3 (2)	6 (2)	12 (5)
				C_{6h}	6/m	175-176	0	3 (2)	0	12 (5)
				D_6	622	177-182	0	3 (2)	2 (1)	12 (5)
				C_{6v}	6mm	183-186	1 (1)	3 (2)	5 (3)	12 (5)
Cubic (isometric) (Four 3-fold axes)	$a = b = c$	P, I, F	m $\bar{3}$	D_{3h}	62M	187-190	0	3 (2)	3 (1)	12 (5)
	$\alpha = \beta = \gamma = 90^\circ$			D_{6h}	6/mmm	191-194	0	3 (2)	0	12 (5)
				T	23	195-199	0	3 (1)	3 (1)	12 (3)
				T_h	m $\bar{3}$	200-206	0	3 (1)	0	12 (3)
				O	432	207-224	0	3 (1)	0	12 (3)
				T_d	43m	215-220	0	3 (1)	3 (1)	12 (3)
Isotropic	Amorphous	—	$\infty\infty m$	—	—	221-230	0	3 (1)	0	12 (3)
							0	3 (1)	0	12 (2)

*Values are the number of nonzero coefficients in (equilibrium) property tensors and the values in parentheses are the numbers of independent coefficients in the tensor. Note that the elasto-optic and piezo-optic tensors have lower symmetry than the rank-4 tensors defined in this table and therefore have more independent coefficients than shown. Second-, third-, and fourth-rank tensors are given in the usual reduced index format (see text).

TABLE 3 Fundamental Physical Constants (2006 CODATA Values²⁵)

Constant	Symbol	Value	Unit
Atomic mass unit (amu)	u	$1.660\,538\,782(83) \times 10^{-27}$	kg
Avogadro constant	N_A	$6.022\,141\,79(30) \times 10^{23}$	mole ⁻¹
Boltzmann constant	k_B	$1.380\,6504(24) \times 10^{-23}$	J/K
Elementary charge	e	$1.602\,176\,487(40) \times 10^{-19}$	C
Permeability of vacuum	μ_0	$4\pi \cdot 10^{-7} = 12.566\,370\,614 \times 10^{-7}$	N/A ² or H/m
Permittivity of vacuum	ϵ_0	$8.854\,187\,187 \times 10^{-12}$	F/m
Planck constant	h	$6.626\,068\,96(33) \times 10^{-34}$	J · s
Speed of light	c	299 792 458	m/s

In optics, it is frequently convenient to define a *complex refractive index*, \bar{n} , the square root of the complex dielectric constant (henceforth, the symbol ϵ will be used for the relative complex dielectric constant):

$$\begin{aligned}\bar{n} &= n + ik = \sqrt{\epsilon} = \sqrt{\epsilon' + i\epsilon''} \\ \epsilon' &= n^2 - k^2; & \epsilon'' &= 2nk \\ n^2 &= \frac{1}{2} \left[\epsilon' + \sqrt{\epsilon'^2 + \epsilon''^2} \right] \\ k^2 &= \frac{1}{2} \left[-\epsilon' + \sqrt{\epsilon'^2 + \epsilon''^2} \right] = \left(\frac{\epsilon''}{2n} \right)^2\end{aligned}\quad (4)$$

where n is the (real) *index of refraction* and k is the *index of absorption* (or imaginary part of the complex refractive index). (The index of absorption is also called the absorption constant, index of extinction, or some other combination of these terms.) Using this definition of the complex index of refraction and the solution of the wave equation [Eq. (3)], the optical power density (proportional to $1/2|\mathbf{E}|^2$ from Poynting's vector) is

$$\text{Power density} = \frac{1}{2} n \sqrt{\mu_0/\epsilon_0} |\mathbf{E}(z)|^2 = \frac{1}{2} n \sqrt{\mu_0/\epsilon_0} |\mathbf{E}(0)|^2 \exp(-2\omega k z/c) \quad (5)$$

where the exponential function represents the attenuation of the wave. The meanings of n and k are clear: n contributes to phase effects (time delay or variable velocity) and k contributes to attenuation by absorption. In practice, attenuation is conveniently described by a power absorption coefficient, β_{ABS} , which describes the internal transmittance over a distance z , i.e.,

$$\tau = \frac{|\mathbf{E}(z)|^2}{|\mathbf{E}(z=0)|^2} = e^{-2\omega k z/c} = e^{-\beta_{\text{ABS}} z} \quad (6a)$$

and β_{ABS} (with units of reciprocal length, usually cm⁻¹) is

$$\beta_{\text{ABS}} = 2\omega k/c = 4\pi\nu k \quad (6b)$$

where ν is the wave number in reciprocal length, usually cm⁻¹.

Kramers–Krönig and Sum Rule Relationships The principal of causality—that a material cannot respond until acted upon—when applied to optics, produces important symmetry properties and relationships that are very useful in modeling and analyzing optical properties. As a consequence of these symmetry properties, the real and imaginary parts of the dielectric constant (and of the complex index of refraction) are Hilbert transforms of each other. The Hilbert transform, $\text{Hi}[\Lambda(\omega)]$, of the complex function $\Lambda(\omega)$ is defined as (the symbol P denotes the principal value of the integral)

$$\text{Hi}[\Lambda(\omega)] = \frac{i}{\pi} P \int_{-\infty}^{\infty} \frac{\Lambda(\omega')}{\omega - \omega'} d\omega' \quad (7)$$

and the relationships between the components of the dielectric constant or index of refraction are

$$\begin{aligned}\epsilon' - 1 &= \text{Hi}[\epsilon''] \\ \epsilon'' &= \text{Hi}^{-1}[\epsilon' - 1] \\ n - 1 &= \text{Hi}[k] \\ k &= \text{Hi}^{-1}[n - 1]\end{aligned}\quad (8)$$

These are the *Kramers–Krönig relationships* (abbreviated KK). Usually the Hilbert transforms for the refractive index are written in single-sided form:

$$n(\omega) - 1 = \frac{2}{\pi} P \int_0^{\infty} \frac{\omega' k(\omega')}{\omega'^2 - \omega^2} d\omega' \quad (9a)$$

with the inverse transform given by

$$k(\omega) = \frac{2\omega}{\pi} P \int_0^{\infty} \frac{n(\omega') - 1}{\omega^2 - \omega'^2} d\omega' \quad (9b)$$

These are fundamental relationships of any causal system.

A number of useful integral relationships, or *sum rules* result from Fourier transforms and the Kramers–Krönig relationship.²⁶ For example, the real part of the refractive index satisfies

$$\int_0^{\infty} [n(\omega) - 1] d\omega = 0 \quad (10a)$$

and

$$\begin{aligned}n(\omega=0) - 1 &= \frac{2}{\pi} \int_0^{\infty} \frac{k(\omega')}{\omega'} d\omega' \\ &= \frac{c}{\pi} \int_0^{\infty} \frac{\beta_{\text{ABS}}(\omega')}{\omega'^2} d\omega'\end{aligned}\quad (10b)$$

The dielectric constant and the refractive index also have the following symmetry properties:

$$\begin{aligned}\epsilon(\omega) &= \epsilon^*(-\omega) \\ \bar{n}(\omega) &= \bar{n}^*(-\omega)\end{aligned}\quad (10c)$$

Practical models of the dielectric constant or refractive index must satisfy these fundamental symmetry properties and integral relationships.

Optical Properties: Origin and Models

Intrinsic optical properties of a material are determined by three basic physical processes: electronic transitions, lattice vibrations, and free-carrier effects.^{5-7,27,28} However, the dominant physical process depends on the material and spectral region of interest. All materials have contributions to the complex index of refraction from electronic transitions. Insulators and semiconductors also require the characterization of the lattice vibrations (or phonons) to fully understand the optical properties. Transparency of semiconductors, particularly those with small bandgaps, is additionally influenced by free-carrier effects. The strength of the free-carrier influence on transmission and absorption depends on the free-carrier concentration; thus free-carrier effects dominate the optical properties of metals in the visible and infrared.

In the range of transparency of a bulk material, more subtle effects such as multiphonon processes (see later discussion; also refer to Chap. 16, “Third-Order Nonlinearities,” by Mansoor Sheik-Bahae

and Michael P. Hasselbeck), impurity and defect absorption, and scattering become the important loss mechanisms. Intrinsic atomic (Rayleigh) scattering is a very weak effect, but is important in long-path optical fibers and ultraviolet materials (refer to Chap. 9, “Volume Scattering in Random Media,” by Aristide Dogariu and Jeremy Ellis in Vol. I). Extrinsic scattering, caused by density (local composition) variations, defects, or grains in polycrystalline solids, is typically much larger than intrinsic scattering in the visible and infrared spectral regions. Impurity and defect (electronic or vibrational) absorption features can be of great concern depending on the spectral region, incident radiation intensity, or material temperature required by the application.

Figure 1 illustrates the frequency dependence of n and k for an insulating polar crystal.²⁹ The value of $n(\omega, T)$ is essentially the sum of the contributions of all electronic and lattice vibration resonances, and is dominated by those with fundamental oscillation frequencies above ω . Figure 1a indicates regions of validity for the popular Sellmeier model (see discussion under “Electronic Transitions”). The frequency dependence of the imaginary part of the index of refraction $k(\omega, T)$ requires consideration of

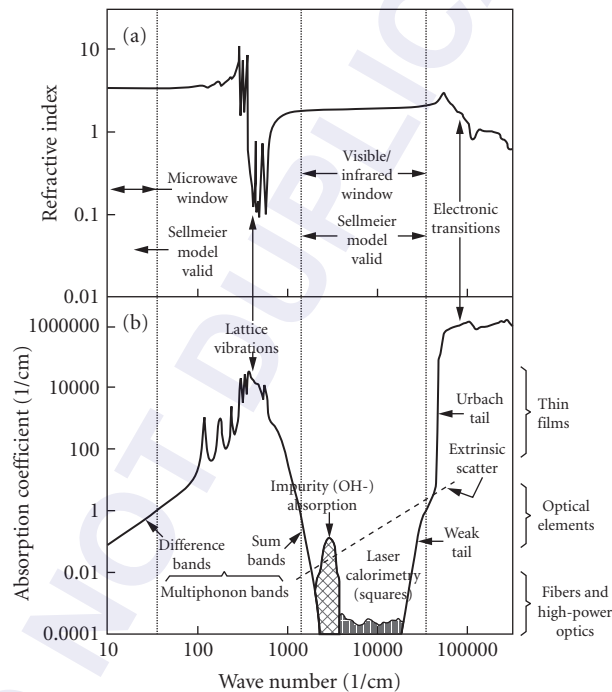


FIGURE 1 The wave number (frequency) dependence of the complex refractive index of Yttria²⁹ (a) shows the real part of the refractive index, $n(\omega)$. The real part is high at low frequency and monotonically increases, becomes oscillatory in the lattice vibration (phonon) absorption bands, increases monotonically (normal dispersion) in the optical transparent region, and again becomes oscillatory in the electronic absorption region; (b) shows the imaginary part of the refractive index, in terms of the absorption coefficient, $\beta(\nu) = 4\pi\nu k(\nu)$. The absorption coefficient is small in the transparent regions and very high in the electronic and vibrational (phonon) absorption bands. The optical transparent region is bounded by the “Urbach tail” absorption at high frequency and by multiphonon absorption at low frequency (wave number). In between, loss is primarily due to impurities and scatter. (Reprinted by permission of McGraw-Hill, Inc.)

not only the dominant physical processes but also higher-order processes, impurities, and defects as illustrated in Fig. 1*b*. The spectral regions of the fundamental resonances are opaque. The infrared edge of transparency is controlled by multiphonon transitions. Transparent regions for insulators are divided in two regions: microwave and visible/infrared.

Lattice Vibrations Atomic motion, or lattice vibrations, accounts for many material properties, including heat capacity, thermal conductivity, elastic constants, and optical and dielectric properties. Lattice vibrations are quantized; the quantum of lattice vibration is called a *phonon*. In crystals, the number of lattice vibrations is equal to three times the number of atoms in the primitive unit cell (see further discussion); three of these are acoustic vibrations (translational modes in the form of sound waves), the remainder are optical vibrations (or modes or phonons) for materials with more than one type atom per unit cell. For most practical temperatures, only the acoustic phonons are thermally excited because optical phonons are typically of much higher frequency, hence acoustic modes play a dominant role in thermal and elastic properties. There are three types of optical modes: infrared-active, Raman-active, and optically inactive. Infrared-active modes, typically occurring in the spectral region from 100 to 1000 cm^{-1} , are those that (elastically) absorb light (photon converted to phonon) through an interaction between the electric field and the light and the dipole moment of the crystal. Raman modes* (caused by phonons that modulate the polarizability of the crystal to induce a dipole moment) weakly absorb light through an inelastic mechanism (photon converted to phonon and scattered photon) and are best observed with intense (e.g., laser) light. Optically inactive modes have no permanent or induced dipole moment and therefore do not interact with light. Optically active modes can be simultaneously infrared- and Raman-active and are experimentally observed by either infrared or Raman spectroscopy, as well as by x-ray or neutron scattering.

Crystal symmetries reduce the number of unique lattice vibrations (i.e., introduce vibrational degeneracies). Group theory analysis determines the number of optical modes of each type for an ideal material. Defects and impurities will increase the number of observed infrared-active and Raman-active modes in a real material. As structural disorder increases (nonstoichiometry, defects, variable composition), the optical modes broaden and additional modes appear. Optical modes in noncrystalline (amorphous) materials such as glasses are very broad compared to those of crystals.

Lattice vibration contributions to the static dielectric constant, $\epsilon(0)$, are determined from the longitudinal-mode (LO) and transverse-mode (TO) frequencies of the optical modes using the Lyddane-Sachs-Teller (LST) relationship³⁰ as extended by Cochran and Cowley³¹ for materials with multiple optical modes:

$$\epsilon(0) = \epsilon(\infty) \prod_j \frac{\omega_j^2(\text{LO})}{\omega_j^2(\text{TO})} \quad (11)$$

where ω is frequency (also expressible in wave numbers) and $\epsilon(\infty)$ is the high-frequency (electronic transition) contribution to the dielectric constant (not the dielectric constant at infinite frequency). This relationship holds individually for each principal axis. The index j denotes infrared-active lattice vibrations with minimum value usually found from group theory (discussed later). This LST relationship has been extended to include the frequency dependence of the real dielectric constant (in transparent, i.e., nonabsorbing, spectral regions):

$$\epsilon(\omega) = \epsilon(\infty) \prod_j \frac{\omega_j^2(\text{LO}) - \omega^2}{\omega_j^2(\text{TO}) - \omega^2} \quad (12a)$$

A modified form of this fundamental equation, including absorption, is used by Gervais and Piriou³² and others to model the complete dielectric constant in the infrared:

$$\epsilon(\omega) = \epsilon(\infty) \prod_j \frac{\omega_j^2(\text{LO}) - \omega^2 - i\omega\gamma_j(\text{LO})}{\omega_j^2(\text{TO}) - \omega^2 - i\omega\gamma_j(\text{TO})} \quad (12b)$$

*Brillouin scattering is a term applied to inelastic scattering of photons by acoustic phonons.

where γ is the line width of the longitudinal and transverse modes as denoted by the symbol in parentheses. This form of the dielectric constant is known as the *semiquantum* four-parameter model.

Frequently, the infrared dielectric constant is modeled in a three-parameter classical oscillator form (or Maxwell–Helmholtz–Drude³³ dispersion formula), namely

$$\varepsilon(\omega) = \varepsilon(\infty) + \sum_j^{\text{jmax}} \frac{S_j \omega_j^2(\text{TO})}{\omega_j^2(\text{TO}) - \omega^2 - i\omega\gamma_j} \quad (13)$$

where $S_j (= \Delta\varepsilon_j)$ is the oscillator strength and γ_j is the full width of the j th mode. This model assumes no coupling between modes and provides a good representation of the dielectric constant, especially if the modes are weak [small separation between $\omega(\text{TO})$ and $\omega(\text{LO})$] and uncoupled. The static dielectric constant $\varepsilon(0)$ is merely the sum of the high-frequency dielectric constant $\varepsilon(\infty)$ and strengths S_j of the individual, IR-active modes. This formulation has been widely used to model infrared dispersion. This model can also be used to represent the high-frequency dielectric constant using additional modes [i.e., $\varepsilon(\infty)$ is replaced by 1, the dielectric constant of free space, plus the contribution of electronic modes; see later discussion in this chapter]. Strengths and line widths for the classical dispersion model can be derived from the values of the four-parameter model.³⁴

Both the classical and four-parameter dispersion models satisfy the Kramers–Krönig relationship [Eqs. (7) and (8)] and are therefore physically realizable. The frequencies $\omega(\text{TO})$ and $\omega(\text{LO})$ arise from the interaction of light with the material, correspond to solutions of the Maxwell wave equation, $\mathbf{V} \cdot \mathbf{D} = 0$ (no free charges), and are obtained from measurements. A transverse field corresponds to an electromagnetic wave with \mathbf{E} perpendicular to the wave vector. At higher frequencies [$\omega > \omega(\text{TO})$], the external electric field counters the internal polarization field of the material until the real part of the dielectric constant is zero, hence $\mathbf{D} = 0$ at $\omega(\text{LO})$. The longitudinal frequency $\omega(\text{LO})$ is always greater than the transverse frequency $\omega(\text{TO})$. The separation of $\omega(\text{LO})$ and $\omega(\text{TO})$ is a measure of the strength (S_j or $\Delta\varepsilon_j$) of the optical mode. Raman modes are, by their nature, very weak and therefore $\omega(\text{LO}) \approx \omega(\text{TO})$, hence Raman modes do not contribute appreciably to dielectric properties.

In Eq. (12a), the $\omega(\text{TO})$ frequencies correspond to the poles and the $\omega(\text{LO})$ frequencies correspond to the zeros of the dielectric constant. The real part of the dielectric constant continuously rises with frequency [except for a discontinuity at $\omega(\text{TO})$] and the dielectric constant is real and negative [i.e., highly absorbing, Eq. (3)] between the transverse and longitudinal frequencies. When damping is added to the dielectric constant model to represent the response of real materials [e.g., Eq. (12b) or (13)], the transverse and longitudinal frequencies become the maxima and minima of the dielectric constant. With damping, a negative dielectric constant is not a necessary condition for absorption, and the material will also absorb outside the region bounded by $\omega(\text{TO})$ and $\omega(\text{LO})$. Furthermore, damping allows the real part of the dielectric constant (and also the refractive index) to decrease with frequency near $\omega(\text{TO})$. This dispersive condition is called *anomalous dispersion* and can only occur in absorptive regions.

As an example of lattice vibrations, consider the simple case of crystalline sodium chloride (NaCl) which has four molecules (eight atoms) per unit cell. Since the sodium chloride structure is face-centered cubic, a primitive cell has one molecule or two atoms. The number of unique vibrations is further reduced by symmetry: for sodium chloride, the lattice vibrations consist of one (triply degenerate) acoustic mode and one (triply degenerate) optical mode. Many metals have one atom per primitive unit cell and therefore have no optical modes. Group IV cubic materials (diamond, silicon, germanium) have two atoms per unit cell with one (triply degenerate) optical mode that is Raman-active only. Therefore, to first order, these nonpolar materials are transparent from their bandgap to very low frequencies. In practice, multiphonon vibrations, defects, impurities, and free carriers introduce significant absorption.

Figure 2 shows a typical crystal infrared spectrum and the corresponding classical oscillator³⁵ and four parameter model³⁶ fits to the data. Parameters for both the extended LST [Eq. (12)] and classical oscillator [Eq. (13)] vibration models are given at the end of this chapter for many materials.

Electronic Transitions Electronic transitions in a solid begin at the material's bandgap. This point generally marks the end (upper frequency) of a material's useful transparency. Above the bandgap, the

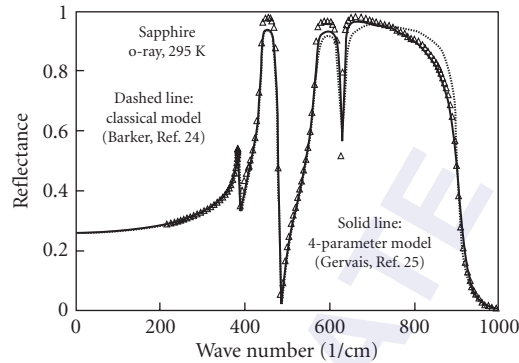


FIGURE 2 The infrared spectrum of sapphire (Al_2O_3) showing the fit to data of the three-parameter classic oscillator model of Barker (dashed line³⁵) and the four-parameter model of Gervais (solid line³⁶). The four-parameter model better fits experimental data (triangles) in the 650 to 900- cm^{-1} region.

material is highly reflective. The large number of possible electronic transitions produces broad-featured spectra. However, electronic structure is fundamental to understanding the nature of the bonds forming the solid and thus many of the material properties. In three dimensions, the band structure becomes more complicated, because it varies with direction just as lattice vibrations do. This point is illustrated in Fig. 3 for the case of the electronic k -space diagram for sapphire.³⁷ Also included in this figure is the corresponding electronic density of states that determines the strength of the absorption.

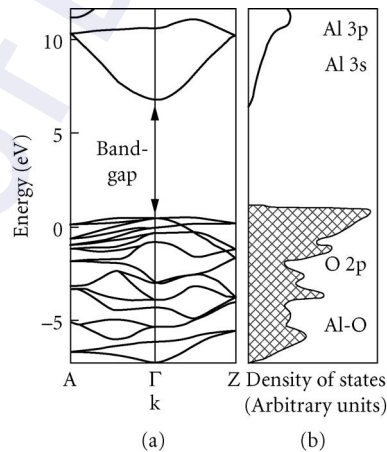


FIGURE 3 Electronic energy band diagram for sapphire at room temperature³⁷ (a) shows the complex k -space energy levels of the electrons of sapphire. The arrow denotes the direct bandgap transition and (b) shows the density of electronic states as a function of energy level. The many direct and indirect electronic transitions give rise to a broad electronic absorption with few features. (Reprinted by permission of the American Ceramic Society.)

In polar insulators, no intraband electronic transitions are allowed, and the lowest frequency electronic absorption is often caused by creation of an *exciton*, a bound electron-hole pair. The photon energy required to create this bound pair is lower than the bandgap energy. Excitons have many properties similar to that of a hydrogen atom. The absorption spectrum of an exciton is similar to that of hydrogen and occurs near the bandgap of the host material. The bond length between the electron-hole pair, hence the energy required to create the exciton, depends on the host medium. Long bond lengths are found in semiconductors (low-energy exciton) and short bond lengths are found in insulating materials.

Other lower-frequency transitions are caused by interband transitions between the anion valence band and the cation conduction band. For sapphire (see Fig. 3), the lowest energy transitions are from the upper valence band of oxygen ($2p^6$) to the conduction band of aluminum ($3s + 3p$). There are typically many of these transitions which appear as a strong, broad absorption feature. Higher energy absorption is caused by surface and bulk plasmons (quanta of collective electronic waves), and still higher energy absorption is attributable to promotion of inner electrons to the conduction band and ultimately liberation of electrons from the material (photoemission).

Classical electronic polarization theory produces a model of the real part of the dielectric constant similar to the model used for the real dielectric constant of lattice vibrations [Eq. (13)]. General properties of the real dielectric constant (and real refractive index) can be deduced from this model. Bound electrons oscillate at a frequency proportional to the square root of the binding energy divided by the electronic mass. Oscillator strength is proportional to the inverse of binding energy. This means that insulators with light atoms and strong bonding have large bandgaps, hence good UV transmission (cf., LiF). Furthermore, high bonding energy (hence high-energy bandgap) means low refractive index (e.g., fluorides).

When both electronic and lattice vibrational contributions to the dielectric constant are modeled as oscillations, the dielectric constant in the transparent region between electronic and vibrational absorption is (mostly) real and the real part takes the form

$$\varepsilon'(\omega) - 1 = n^2(\omega) - 1 = \sum_j \frac{S_j \omega_j^2}{\omega_j^2 - \omega^2} \quad (14)$$

which is the widely used *Sellmeier dispersion formula*.³⁸ The sum includes both electronic (UV) and vibrational (IR or *polar* or *ionic*) contributions. Most other dispersion formulas (such as the Schott glass power series) are recast, series expanded, or simplified forms of the Sellmeier model. The refractive index of most materials with good homogeneity can be modeled to a few parts in 10^5 over their entire transparent region with a Sellmeier fit of a few terms. The frequency (ω_j) term(s) in a Sellmeier fit are not necessarily TO modes, but are correlated to the strong TOs nearest the transparent region, with adjustments to the constants made to account for weaker modes, multiphonon effects, and impurities. The Sellmeier model works well because it is (1) based on a reasonable physical model and (2) adjusts constants to match data. The relationship between the Sellmeier equation and other dispersion formulas is discussed later.

The *Urbach tail* model is successfully applied to model the frequency and temperature dependence of the ultraviolet absorption edge in a number of materials, particularly those with a direct bandgap, over several orders of magnitude in absorption. Urbach³⁹ observed an exponential absorption edge in silver halide materials (which have an indirect bandgap). Further development added temperature dependence in the form

$$\beta_{\text{ABS}}(E, T) = \beta_0 \exp \left[-\sigma(T) \frac{(E_0 - E)}{k_B T} \right] \quad (15a)$$

where

$$\sigma(T) = \sigma_0(T) \frac{2k_B T}{E_p} \tanh \frac{E_p}{2k_B T} \quad (15b)$$

where E_0 is the bandgap energy (typically the energy of an exciton) at $T = 0$ K, T temperature in kelvins, and E_p a characteristic phonon energy. The interpretation of the Urbach tail is a broadening of the electronic bandgap by phonon interactions, and several detailed theories have been proposed.^{28,40}

Below the Urbach tail, absorption continues to decrease exponentially, albeit much slower than predicted by the Urbach formula. This region of slowly decreasing absorption is sometimes called the “weak tail,” and has been observed in both semiconductor⁴¹ and crystalline materials.^{42,43} Typically, the weak tail begins at the point when the absorption coefficient falls to 0.1 cm^{-1} .

Unfortunately, the optical properties of the electronic transitions, which drive optoelectronic device performance, can not be modeled in a straightforward manner like the vibrational transition models because most electronic transitions are coupled together to form broad homogeneous absorption bands. Thus, the shape of the density-of-states function that determines the spectral response is seldom in the simple functional form of a classical oscillator model.

Adachi⁴⁴ has developed models that successfully represent the complex permittivity of electronic transitions as a function of both frequency and temperature by considering the effect of various electronic interband transitions from valence to conduction band. These transitions peak at various critical point locations (e.g., Γ and L) and include energy level shifts caused by spin-orbit coupling. Various density-of-state models are used, including the classical oscillator to represent several closely spaced transitions.

The Adachi model accounts for first-order electronic transitions, and represents the complex refractive index over a wide frequency and temperature range. Model predictions agree well with experimental complex index of refraction data in transparent spectral regions ending at the bandgap. Model parameters are available for CdTe, GaAs, GaP, GaSb, InAs, InP, InSb, ZnTe, AlGaAs, InGaAsP, and others.⁴⁴⁻⁵¹

The electronic transitions have a very different functional form than the vibrational modes, yet a Sellmeier model, an approximate form to the classical oscillator model, is commonly used to represent the real part of the relative permittivity in transparent spectral regions ending at the bandgap. However, a polynomial expansion of the electronic susceptibility functions shows similar structure to an expansion of the Sellmeier equation. It is expected that when high precision is needed (e.g., $<10^{-5}$) in the refractive index, the more physically correct (Adachi) model should be used. This is certainly the case for absorption near the bandgap.

Free Carriers Free carriers, such as electrons in metals, or electrons and holes in semiconductors, also affect the optical properties of materials. For insulators or wide-bandgap semiconductors (i.e., bandgap greater than 0.5 eV) with a low number of free carriers at room temperature (low conductivity), the effect of free carriers on optical absorption is small [see Eq. (2)]. For nonmetals, the free-carrier concentration grows with temperature so that even an “insulator” has measurable conductivity (and free-carrier absorption) at very high temperature. Commonly used optical materials such as silicon and germanium have a significant increase in free-carrier absorption at moderately high temperature as illustrated in Fig. 4.⁵²

Free-carrier effects can be modeled as an additional contribution to the dielectric constant model. For example, the classical model [Eq. (13)] takes the form

$$\epsilon(\omega) = \epsilon(\infty) + \sum_j^{\max} \frac{S_j \omega_j^2(\text{TO})}{\omega_j^2(\text{TO}) - \omega^2 - i\omega\gamma_j} - \frac{\omega_p^2}{\omega^2 + i\omega\gamma_c} \quad (16)$$

where ω_p is the plasma frequency, proportional to the square root of the free-carrier density, and γ_c is the damping frequency (i.e., determines the effective width of the free-carrier influence). Such a model is well known to accurately predict the far-infrared ($\geq 10 \mu\text{m}$) refractive index of metals⁵³ and also has been used to model the free-carrier contribution to the optical properties of semiconductors.

Multiphonon Absorption and Refraction Absorption at the infrared edge of insulators is principally caused by anharmonic terms in the lattice potential leading to higher harmonics of the lattice resonances. This phenomenon is called multiphonon absorption because the frequencies are harmonics

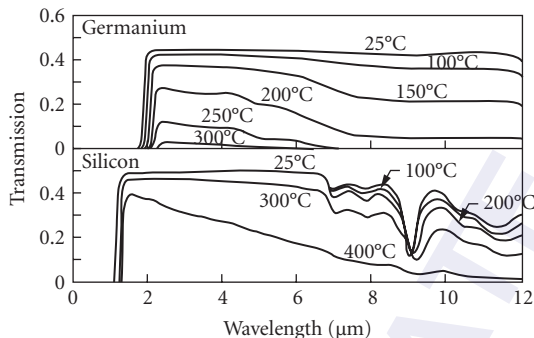


FIGURE 4 Decreased transmission of germanium and silicon with temperature is attributable to an increase in free-carrier concentration resulting in increased absorption.⁵² The absorption is greater at longer wave lengths; see Eq. (16). (Reprinted by permission of the Optical Society of America.)

of the characteristic lattice phonons (vibrations). For absorption in the infrared, each successively higher multiple of the fundamental frequency is weaker (and broader) leading to decreasing absorption beyond the highest fundamental absorption frequency (maximum transverse optical frequency). At about three times $\omega(\text{TO})$ the absorption coefficient becomes small and a material with thickness of 1 to 10 μm is reasonably transparent. The infrared absorption coefficient of materials (especially highly ionic insulators) can be characterized by an exponential absorption coefficient,⁵⁴ β_{ABS} of the form

$$\beta_{\text{ABS}} = \beta_0 \exp\left(-\gamma \frac{\omega}{\omega_0}\right) \quad (17a)$$

where β_0 is a constant (dimensions same as the absorption coefficient, typically cm^{-1}), γ is a dimensionless constant (typically found to be near 4), ω_0 is frequency or wave number of the maximum transverse optical frequency (units are cm^{-1} for wave numbers; values are given in the property data tables), and ω is the frequency or wave number of interest. This formula works reasonably well for ionic materials at room temperature for the range of absorption coefficients from 0.001 to 10 cm^{-1} .

In the classical (continuum) limit, the temperature dependence of multiphonon absorption has a T^{-n} dependence, where n is the order of the multiphonon process,⁵⁵ that is $n \approx \omega/\omega_0$. At low temperature, there is no temperature dependence since only transitions from the ground state occur. Once the temperature is sufficiently high (e.g., approaching the Debye temperature), transitions that originate from excited states become important and the classical temperature dependence is observed. Bendow et al.⁵⁶ has developed a simple model of the temperature dependence of multiphonon absorption based on a Bose–Einstein distribution of states:

$$\beta_{\text{ABS}}(\omega, T) = \beta_0 \frac{[N(\omega_0, T) + 1]^{\omega/\omega_0}}{N(\omega, T) + 1} \exp(-A\omega/\omega_0) \quad (17b)$$

where $N(\omega, T)$ is the phonon occupation density from Bose–Einstein statistics.

Thomas et al.⁵⁷ have successfully developed a semiempirical, quantum mechanical model of (sum band) multiphonon absorption based on the Morse interatomic potential and a gaussian function for the phonon density of states. Use of the Morse potential leads to an exact solution to the Schrödinger equation and includes anharmonic effects to all orders. The model contains parameters derived from room-temperature measurements of absorption, is computationally efficient, and has been applied to many ionic substances. Figure 5 shows a typical result compared to experimental data. Model parameters are available for oxides,⁵⁸ diamond,⁵⁸ glasses,⁵⁹ and semiconductors.⁶⁰

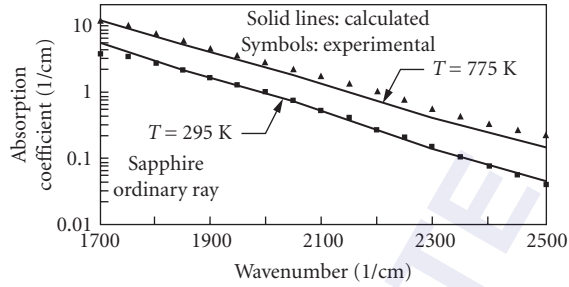


FIGURE 5 Temperature-dependent change of absorption in insulators is principally confined to the absorption edges, especially the infrared multiphonon absorption edge. This figure shows measured and predicted absorption coefficients at the infrared edge of transparency for the ordinary ray of crystalline sapphire. Increasing temperature activates higher multiphonon processes, resulting in a rapid increase in absorption. The multiphonon model of Thomas et al.⁵⁷⁻⁵⁹ accurately predicts the frequency- and temperature-dependence of infrared absorption in highly ionic materials such as oxides and halides.

Multiphonon absorption modeling also contributes to the real refractive index. Although multiphonon contributions to the real index are small compared to one-phonon contributions, they are important for two cases in the infrared: (1) when the refractive index must be known beyond two decimal places, or (2) at high temperature. In the first case, multiphonon contributions to the index are significant over a large spectral region. In the second case, the contribution of multiphonon modes to the real refractive index grows rapidly at high temperature because of the T^{n-1} dependence of the n th mode strength.

Absorption in the Transparent Region In the transparent region, away from the electronic and vibrational resonances, absorption is governed by impurities and defects. The level of absorption is highly dependent on the purity of the starting materials, conditions of manufacture, and subsequent machining and polishing. For example, OH impurities are common in oxides,* occurring at frequencies below the fundamental (nonbonded) OH vibration at 3735 cm^{-1} . OH can be removed by appropriate heat treatment.

Low-level absorption coefficient measurements are typically made by laser calorimetry or photoacoustic techniques. Data are available for a number of materials in the visible,^{42,61} at $1.3\text{ }\mu\text{m}$,⁶² and at 2.7 and $3.8\text{ }\mu\text{m}$.⁶³

Optical Properties: Applications

Dielectric Tensor and Optical Indicatrix Many important materials are nonisotropic (e.g., the crystals Al_2O_3 , SiO_2 , and MgF_2), and their optical properties are described by tensor relationships (see earlier section, “Symmetry Properties”). The dielectric constant, ϵ , a second-rank tensor, relates the electric field \mathbf{E} to the electric displacement \mathbf{D} :

$$\begin{pmatrix} \mathbf{D}_x \\ \mathbf{D}_y \\ \mathbf{D}_z \end{pmatrix} = \epsilon_0 \cdot \begin{pmatrix} \epsilon_{xx} & \epsilon_{xy} & \epsilon_{xz} \\ \epsilon_{yx} & \epsilon_{yy} & \epsilon_{yz} \\ \epsilon_{zx} & \epsilon_{zy} & \epsilon_{zz} \end{pmatrix} \cdot \begin{pmatrix} \mathbf{E}_x \\ \mathbf{E}_y \\ \mathbf{E}_z \end{pmatrix} \quad (18a)$$

*The OH vibrational impurity absorption in oxides is known for Al_2O_3 , ALON, MgAl_2O_4 , MgO, SiO_2 , Y_2O_3 , and Yb_2O_3 .

From the symmetry of properties, this is a symmetric tensor with $\epsilon_{ab} = \epsilon_{ba}$. Usually, the dielectric constant components are given as principal values, that is, those values along the unit cell of the appropriate crystal class. In this case, the principal dielectric constants are

$$\begin{aligned}\epsilon_{x'x'} &\equiv \epsilon_1 \\ \epsilon_{y'y'} &\equiv \epsilon_2 \\ \epsilon_{z'z'} &\equiv \epsilon_3\end{aligned}\tag{18b}$$

where the primes on the subscripts denote principal values (i.e., along the crystallographic axes, possible in a nonorthogonal coordinate system) and the subscripts 1, 2, and 3 denote reduced notation for these values (see "Elastic Properties"). The relationship between dielectric constant and refractive index, Eq. (4), means there are similarly three principal values for the (complex) refractive index. Also, the components of the dielectric tensor are individually related to the corresponding components of the refractive index. (Subscripts a , b , and c or x , y , and z as well as others may be used for the principal values of the dielectric constant or refractive index.)

Three important cases arise:

1. Isotropic and cubic materials have only one dielectric constant, ϵ , (hence one refractive \bar{n}). Therefore $\epsilon_1 = \epsilon_2 = \epsilon_3 = \epsilon \bar{n}^2$.
2. Hexagonal (including trigonal) and tetragonal crystals have two principal dielectric constants, ϵ_1 and ϵ_3 (hence two refractive indices, \bar{n} and \bar{n}_3). Therefore $\epsilon_1 = \epsilon_2 \neq \epsilon_3$. Such materials are called *uniaxial*; the unique crystallographic axis is the c axis, which is also called the *optical axis*. One method of denoting the two unique principal axes is to state the orientation of the electric field relative to the optical axis. The dielectric constant for the $\mathbf{E} \perp c$ situation is ϵ_1 or ϵ_{\perp} . This circumstance is also called the *ordinary ray*, and the corresponding symbol for the refractive index is $\bar{n}_1, \bar{n}_{\perp}$ (for ordinary ray), and \bar{n}_o or ω (primarily in the older literature). The dielectric constant for $\mathbf{E} \parallel c$ situation is ϵ_3 or ϵ_{\parallel} . This condition is called the *extraordinary ray*, and the corresponding symbol for the refractive index is $\bar{n}_3, \bar{n}_e, \bar{n}_e$ (for extraordinary ray), and \bar{n}_e or ϵ (again, primarily in the older literature). Crystals are called *positive uniaxial* when $n_e - n_o > 0$, and *negative uniaxial* otherwise. Since the dispersions of the ordinary and extraordinary wave are different, a crystal can be positive uniaxial in one wavelength region and negative uniaxial in another (AgGaS₂ is an example).
3. Orthorhombic, monoclinic, and triclinic crystals have three principal dielectric constants, ϵ_1 , ϵ_2 , and ϵ_3 (hence, three refractive indices, \bar{n}_1, \bar{n}_2 , and \bar{n}_3). Therefore $\epsilon_1 \neq \epsilon_2 \neq \epsilon_3 \neq \epsilon_1$. These crystals are called *biaxial*. Confusion sometimes arises from the correlation of the principal dielectric constants with the crystallographic orientation owing to several conventions in selecting the crystal axes. [The optical indicatrix (see following discussion) of a biaxial material has two circular sections that define optical axes. The orientation of these axes are then used to assign a positive- or negative-biaxial designation.]

The existence of more than one dielectric constant or refractive index means that, for radiation with arbitrary orientation with respect to the crystal axes, two plane-polarized waves, of different speed, propagate in the crystal. Hence, for light propagating at a random orientation to the principal axes, a uniaxial or biaxial crystal exhibits two effective refractive indices different from the individual principal values. The refractive index of the two waves is determined from the *optical indicatrix* or *index ellipsoid*, a triaxial ellipsoidal surface defined by

$$\frac{x_1^2}{n_1^2} + \frac{x_2^2}{n_2^2} + \frac{x_3^2}{n_3^2} = 1\tag{19a}$$

where the x_1 , x_2 , and x_3 are the principal axes of the dielectric constant. The indicatrix is illustrated in Fig. 6: for a wave normal in an arbitrary direction (OP), the two waves have refraction indices equal to the axes of the ellipse perpendicular to the wave normal (OA and OB). The directions

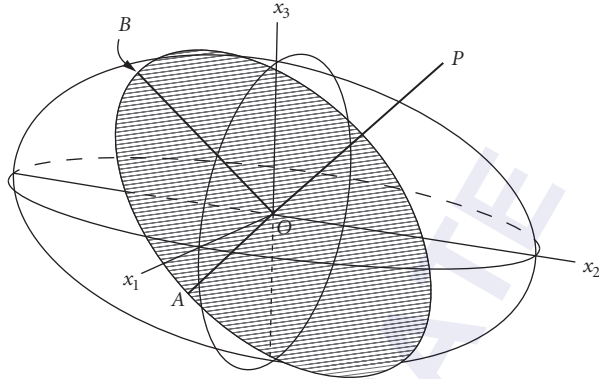


FIGURE 6 The optical indicatrix or index ellipsoid used to determine the effective refractive index for an arbitrary wave normal in a crystal. The axes of the ellipsoid correspond to the principal axes of the crystal, and the radii of the ellipsoid along the axes are the principal values of the refractive indices. For propagation along an arbitrary wave normal (OP), the effective refractive indices are the axes of the ellipse whose normal is parallel to the wave normal. In the illustrated case, the directions OA and OB define the effective refractive indices.²⁴ (Reprinted by permission of Oxford University Press.)

represented by OA and OB are the vibrational planes of the electric displacement vector \mathbf{D} of the two waves. When the wave normal is parallel to an optic axis, the two waves propagate with principal refractive indices. For uniaxial materials and the wave normal parallel to the x_3 (or z or c crystallographic or optical) axis, the vibrational ellipsoid is circular and the two waves have the same refractive index (n_o), and there is no double refraction. Equation (19a) can also be written

$$B_1 x_1^2 + B_2 x_2^2 + B_3 x_3^2 = 1 \quad (19b)$$

where $B_i = 1/\bar{n}_i^2 = 1/\epsilon_i$ is called the *inverse dielectric tensor*. The inverse dielectric tensor is used in defining electro-optic, piezo-optic, and elasto-optic relationships (see Table 1).

Fresnel formulas, found elsewhere in this *Handbook*, use the complex refractive index to calculate reflection and transmission at a material boundary as a function of incident angle, polarization, and (crystal) orientation (for uni- and biaxial crystals).

Total Power Law Incident light on a material is reflected, transmitted, or absorbed. Scattering is a term used to describe diffuse reflectance (surface scatter) and diffuse transmittance (bulk scatter). Conservation of energy dictates that the fractional amount reflected ρ , absorbed α_{ABS} , transmitted τ , and scattered α_{SCA} total to unity, hence

$$1 = \rho(\Omega_i, \omega) + \alpha_{\text{SCA}}(\Omega_i, \omega) + \alpha_{\text{ABS}}(\Omega_i, \omega) + \tau(\Omega_i, \omega) \quad (20)$$

where these time-averaged quantities, illustrated in Fig. 7, are

$$\rho(\Omega_i, \omega) = \frac{\Phi_r(\Omega_i, \omega)}{\Phi_i(\omega)} = \text{total integrated reflectance} \quad (21a)$$

$$\alpha_{\text{SCA}}(\Omega_i, \omega) = \frac{\Phi_s(\Omega_i, \omega)}{\Phi_i(\omega)} = \text{total integrated scatterance} \quad (21b)$$

$$\alpha_{\text{ABS}}(\Omega_i, \omega) = \frac{\Phi_a(\Omega_i, \omega)}{\Phi_i(\omega)} = \text{total integrated absorptance} \quad (21c)$$

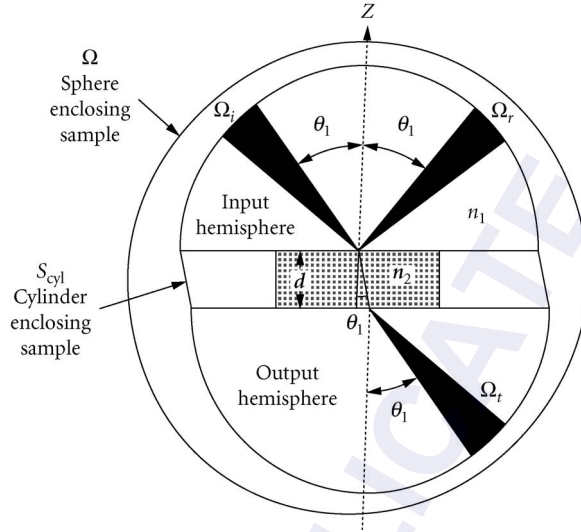


FIGURE 7 Geometry of incident, transmitted, and reflected beams for a plane transparent slab of thickness d . The power equals the reflected, refracted, and absorbed components, assuming no scatter.

and

$$\tau(\Omega_i, \omega) = \frac{\Phi_t(\Omega_i, \omega)}{\Phi_i(\omega)} = \text{total integrated transmittance} \quad (21d)$$

Notice that these quantities are functions of the angle of incidence and frequency only.

The sum of total integrated scatterance and total integrated absorbance can be defined as the total integrated extintance α_{EXT}

$$\alpha_{\text{EXT}}(\Omega_i, \omega) = \alpha_{\text{ABS}}(\Omega_i, \omega) + \alpha_{\text{SCA}}(\Omega_i, \omega) \quad (22)$$

and the total power law becomes

$$1 = \rho(\Omega_i, \omega) + \alpha_{\text{EXT}}(\Omega_i, \omega) + \tau(\Omega_i, \omega) \quad (23)$$

Another useful quantity is emittance, which is defined as

$$\varepsilon(\Omega_i, \omega) = \frac{\Phi_e(\Omega_i, \omega)}{\Phi_{\text{bb}}(\omega)} = \text{total integrated emittance} \quad (24a)$$

where Φ_{bb} is the blackbody function representing the spectral emission of a medium which totally absorbs all light at all frequencies. When $\Phi_i(\omega) = \Phi_{\text{bb}}(\omega)$, then the total integrated emittance equals the total integrated absorbance:

$$\varepsilon(\Omega_i, \omega) = \alpha_{\text{ABS}}(\Omega_i, \omega) \quad (24b)$$

Dispersion Formulas for Refractive Index The dielectric constant and refractive index are functions of frequency, hence wavelength. The frequency or wavelength variation of refractive index is called *dispersion*. Dispersion is an important property for optical design (i.e., correction of chromatic aberration) and in the transmission of information (i.e., pulse spreading). Other optical properties are derived from the change in refractive index with other properties such as temperature (*thermo-optic* coefficient), stress or strain (*piezo-optic* or *elasto-optic* coefficients), or applied field (*electro-optic* or *piezo-electric* coefficients). Since the dielectric constant is a second-order tensor with three principal values, the coefficients defined here are also tensor properties (see Table 1).

Precise refractive index measurements give values as functions of wavelength. Frequently, it is desirable to have a functional form for the dispersion of the refractive index (i.e., for calculations and value interpolation). There are many formulas used for representing the refractive index. One of the most widely used is the Sellmeier (or Drude or Maxwell-Helmholtz-Drude) dispersion model [Eq. (14)], which arises from treating the absorption like simple mechanical or electrical resonances. Sellmeier³⁸ proposed the following dispersion formula in 1871–1872 (although Maxwell had also considered the same derivation in 1869). The usual form of this equation for optical applications gives refractive index as a function of wavelength rather than wave number, frequency, or energy. In this form, the Sellmeier equation is

$$n^2(\lambda) - 1 = \sum_{i=1} \frac{A_i \cdot \lambda^2}{\lambda^2 - \lambda_i^2} \quad (25a)$$

The Sellmeier formula has the appropriate physical basis to accurately represent the refractive index throughout the transparent region in the simplest manner. The Sellmeier constants have physical meaning, particularly for simple substances. Most other dispersion formulas are closely related to (or are a disguised form of) the Sellmeier equation. Many of these other dispersion formulas are unable to cover as wide a spectral region, and unlike the Sellmeier form, do not lend themselves to extrapolation outside the interval of available measurements. For these reasons, we strongly urge that the Sellmeier model be universally used as the standard representation of the refractive index.

Modifications of the Sellmeier terms that include composition variation^{64,65} and temperature dependence⁶⁶ have been applied to successfully model refractive index. The variation of the Sellmeier A_i and λ_i constants is usually modeled as linearly dependent on the mole fraction of the components and temperature.

An often-used, slight modification of this formula puts the wavelength of the shortest wavelength resonance at zero ($\lambda_i = 0$), that is, the first term is a constant. This constant term represents contributions to refractive index from electronic transitions at energies far above the bandgap. Sellmeier terms with small λ_i (representing electronic transitions) can be expanded as a power series,

$$\frac{A_i \cdot \lambda^2}{\lambda^2 - \lambda_i^2} = A_i \cdot \sum_{i=0}^{\infty} (\lambda_i^2 / \lambda^2)^j = A_i + \frac{A_i \cdot \lambda_i^2}{\lambda^2} + \frac{A_i \cdot \lambda_i^4}{\lambda^4} + \dots \quad (25b)$$

and the terms with large λ_i (representing vibrational transitions) are expanded as:

$$\frac{A_i \cdot \lambda^2}{\lambda^2 - \lambda_i^2} = -A_i \cdot \sum_{j=1}^{\infty} (\lambda^2 / \lambda_i^2)^j = -A_i \cdot \frac{\lambda^2}{\lambda_i^2} - A_i \cdot \frac{\lambda^4}{\lambda_i^4} - \dots \quad (25c)$$

The first term of this expansion is occasionally used to represent the long-wavelength (infrared) contribution to the index of refraction (instead of the Sellmeier term).

A generalized form of the short-wavelength approximation to the Sellmeier equation is the Cauchy formula, developed in 1836. This was the first successful attempt to represent dispersion by an equation:

$$n = A_0 + \sum_{i=1} \frac{A_i}{\lambda^{2i}} \quad \text{or} \quad n^2 = A'_0 + \sum_{i=1} \frac{A'_i}{\lambda^{2i}} \quad (26)$$

Power series approximations to the Sellmeier equation are expressed in many forms. One common form is the Schott glass formula used for glasses:

$$n^2 = A_0 + A_1\lambda^2 + A_2\lambda^{-2} + A_3\lambda^{-4} + A_4\lambda^{-6} + A_5\lambda^{-8} + \dots \quad (27)$$

For typical high-quality glasses, this equation is accurate to $\pm 3 \cdot 10^{-6}$ in the visible (400–765 nm) and within $\pm 5 \cdot 10^{-6}$ from 365 to 1014 nm. A comparison of the Schott power series formula with a three-term Sellmeier formula showed equivalent accuracy of the range of the Schott fit, but that the Sellmeier model was accurate over a much wider wavelength range.⁶⁷ A number of other power series dispersion formulas (e.g., Ketteler–Neumann⁶⁸) are occasionally used.

Frequently, Sellmeier terms are written in altered fashion such as this form used by Li⁶⁹ and others:

$$\frac{A_i}{(\lambda^2 - \lambda_i^2)} = \frac{A_i(\lambda^2 / \lambda_i^2)}{(\lambda^2 - \lambda_i^2)} - \frac{A_i}{\lambda_i^2} \quad (28a)$$

which is the equivalent of the combination of two Sellmeier terms, one located at zero wavelength and the other at λ_i . The Zernike formula⁷⁰ also uses a term in this form. Another way to modify Sellmeier terms is to convert the wavelength of the resonances to wave number or energy [see Eq. (14)].

Another common formula for the index of refraction is the Hartmann⁴³ or Cornu equation:

$$n = A + \frac{B}{\lambda - \lambda_0} \quad (28b)$$

This equation is more distantly related to the Sellmeier formulation. Note that a two-term Sellmeier formula (with $\lambda_1 = 0$) can be written as

$$(n^2 - n_0^2) \cdot (\lambda^2 - \lambda_0^2) = (n - n_0) \cdot (\lambda - \lambda_0) \cdot (n + n_0) \cdot (\lambda + \lambda_0) = \text{constant} \quad (28c)$$

and the Hartmann formula can be written as a hyperbola

$$(n - n_0) \cdot (\lambda - \lambda_0) = \text{constant} \quad (28d)$$

Note that in a limited spectral region, the difference terms of the Sellmeier formula of Eq. (28c) vary much more rapidly than do the sum terms, hence the Hartmann and Sellmeier forms will have nearly the same shape in a limited spectral range.

Other equations that combine Sellmeier and power series terms (cf., Wemple formula) are often used. One such formulation is the Herzberger equation, first developed for glasses⁷¹ and later applied to infrared crystalline materials:⁷²

$$n = A + \frac{B}{(\lambda^2 - 0.028)} + \frac{C}{(\lambda^2 - 0.028)^2} + D\lambda^2 + E\lambda^4 \quad (28e)$$

where the choice of the constant $\lambda_0^2 = 0.028$ is arbitrary in that it is applied to all materials.

The Pikhtin–Yas'kov formula⁷³ is nearly the same as the Sellmeier form with the addition of another term representing a broadband electronic contribution (which is equal to the real part of one of the terms in Adachi's model⁴⁴):

$$n^2 - 1 = \frac{A}{\pi} \ln \frac{E_1^2 - (\hbar\omega)^2}{E_0^2 - (\hbar\omega)^2} + \sum_i \frac{G_i}{E_i^2 - (\hbar\omega)^2} \quad (28f)$$

This formulation has been applied to some semiconductor materials. The unique term arises from assuming that the imaginary part of the dielectric constant is a constant between energies E_0 and E_1 , and that infinitely narrow resonances occur at E_i . The formula is then derived by applying the Kramers–Krönig relationship to this model.

Dispersion Formulas Used for Glasses Dispersion in glasses (and most other materials) is accurately represented by a Sellmeier model [Eq. (25a)] with three terms, two representing electronic (short wavelength) absorption and one representing lattice vibration (long wavelength) absorption contribution to the refractive index. Typically for glasses, the *Abbe number*, or constringence v_d is also given. The Abbe number is a measure of dispersion in the visible and is defined as $v_d = (n_d - 1)/(n_F - n_C)$, where n_d , n_F , and n_C are refractive indices at 587.6, 486.1, and 656.3 nm, respectively. The quantity $(n_F - n_C)$ is known as the *principle dispersion*. A relative partial dispersion $P_{x,y}$ can be calculated at any wavelengths x and y from

$$P_{x,y} = \frac{(n_x - n_y)}{(n_F - n_C)} \quad (29a)$$

For so-called “normal” glasses, the partial dispersions obey a linear relationship, namely

$$P_{x,y} = a_{x,y} + b_{x,y} v_d \quad (29b)$$

where $a_{x,y}$ and $b_{x,y}$ are empirical constant characteristics of normal glasses. However, for correction of secondary spectrum in an optical system (i.e., achromatization for more than two wavelengths), it is necessary to employ a glass that does not follow the glass line. Glass manufacturers usually list $\Delta P_{x,y}$ for a number of wavelength pairs as defined by:

$$P_{x,y} = a_{x,y} + b_{x,y} v_d + \Delta P_{x,y} \quad (29c)$$

The deviation term $\Delta P_{x,y}$ is a measure of the dispersion characteristics differing from the normal glasses. Schott glasses F2 and K7 define the normal glass line.

In the transparent region, refractive index decreases with increasing wavelength, and the magnitude of $dn/d\lambda$ is a minimum between the electronic and vibrational absorptions. The wavelength (λ_0) of minimum $dn/d\lambda$, called the “zero”-dispersion point, is given by

$$\left. \frac{d^2 n(\lambda)}{d\lambda^2} \right|_{\lambda=\lambda_0} = 0 \quad (30)$$

which is the desired operating point for high-bandwidth, information-carrying optical fibers, as well as the optimum wavelength for single-element refractive optical systems. For glassy silica fibers, the zero dispersion point is 1.272 μm .

The equation for the dispersion in a standard single-mode fiber is

$$D_{\text{chromatic}} = \frac{S_0}{4} \left(\lambda - \frac{\lambda_0^4}{\lambda^3} \right) \quad (31)$$

where the operating wavelength is λ , the zero-dispersion wavelength is λ_0 , and S_0 is the slope of the dispersion curve. The dispersion units are picoseconds per wavelength bandwidth (nm) per unit length of fiber (km).

One approach to reducing both dispersion and loss is to use a material with a wide transparent region, that is, widely separated electronic and vibrational absorptions, hence the interest in materials such as heavy-metal fluoride glasses for fiber applications. Recent work in photonic bandgap (PBG) fibers and micro-structured optical fibers demonstrates that light can be guided in hollow cores. Index-guided fibers can operate below 800 nm but the effective mode area is small. PBG fibers have been theorized to work below 800 nm with an order of magnitude larger effective mode area.⁷⁴ Hence, dispersionless transmission of short pulses in the near visible range is achieved.^{75–77}

Thermo-Optic and Photoelastic Coefficients Temperature is one of the main factors influencing the refractive index of solids. The thermo-optical coefficients $\partial n/\partial T$ (or $\partial \epsilon/\partial T$) can be estimated from a derivation of the Clausius–Mossotti relationship:⁷⁸

$$\frac{1}{(\epsilon-1)(\epsilon+2)}\left(\frac{\partial \epsilon}{\partial T}\right) = -\alpha \left[1 - \frac{V}{\alpha_m} \left(\frac{\partial \alpha_m}{\partial V} \right)_T \right] + \frac{1}{3\alpha_m} \left(\frac{\partial \alpha_m}{\partial T} \right)_V \quad (32)$$

where α_m is the macroscopic polarizability. The first two terms are the principal contributors in ionic materials: a positive thermal expansion coefficient results in a negative thermo-optic coefficient and a positive change in polarizability with volume results in a positive thermo-optic coefficient. In ionic materials with a low melting point, thermal expansion is high and the thermo-optic coefficient is negative (typical of alkali halides); when thermal expansion is small (indicated by high melting point, hardness, and high elastic moduli), the thermo-optic coefficient is positive, dominated by the volume change in polarizability (typical of the high-temperature oxides).

Thermal expansion has no frequency dependence but polarizability does. At frequencies (wavelengths) near the edge of transparency, the polarizability (and $\partial \alpha_m/\partial V$) rises, and $\partial n/\partial T$ becomes more positive (or less negative). Figure 8 shows the variation of refractive index for sodium chloride as a function of frequency and temperature.

An approach similar to that for the thermo-optic coefficient, above, can be used to estimate the photoelastic constants of a material. The simplest photoelastic constant is that produced by uniform pressure, that is, dn/dP . More complex photoelastic constants are tensors whose components define the effect of individual strain (elasto-optic coefficients) or stress (piezo-optic coefficients) tensor terms. Bendow et al.⁷⁹ calculate dn/dP and the elasto-optic coefficients for a number of cubic crystals and compare the results to experiment.

Thermo-Optic Coefficients for Glasses The effect of temperature on the refractive index is dependent upon several factors, including the material and wavelength. The refractive index also changes as a function of air pressure. A distinction is therefore made between absolute refractive index, which is measured in vacuum, and refractive index relative to normal air, which is defined as the refractive index at normal air pressure (1 atmosphere = 101.325 kPa). The standard temperature

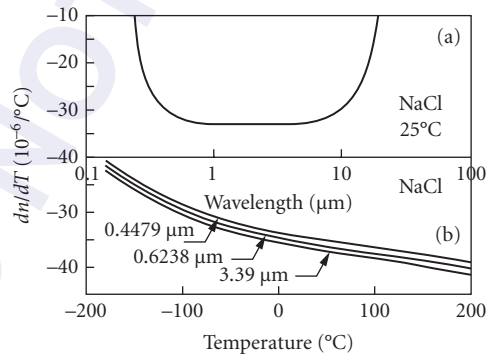


FIGURE 8 The thermo-optic coefficient (dn/dT) of sodium chloride (NaCl): (a) shows the wavelength dependence of room-temperature thermo-optic coefficient. The thermo-optic coefficient is nearly constant in the transparent region, but increases significantly at the edges of the transparent region and (b) illustrates the temperature dependence of the thermo-optic coefficient. The thermo-optic coefficient decreases (becomes more negative) with increasing temperature, primarily as the density decreases.

for refractive index measurements is 22°C, and manufacturer's quoted index of refraction values are always relative to normal air pressure.

The change in the absolute refractive index as a function of temperature is modeled as a first-order or linear change with temperature,

$$n_{\text{abs}}(\lambda, T) = n_{\text{abs}}(\lambda, T_0) + \frac{dn_{\text{abs}}(\lambda, T)}{dT} \Delta T \quad (33)$$

In the visible, the absolute refractive index as a function of wavelength can be well-approximated by a single-term Sellmeier model, Eq. (25a), and the temperature derivative of refractive index is then

$$\frac{dn_{\text{abs}}}{dT} = \frac{1}{2n_{\text{abs}}} \frac{d(n_{\text{abs}}^2 - 1)}{dT} = \frac{n_{\text{abs}}^2 - 1}{2n_{\text{abs}}} \left[\frac{1}{A} \frac{dA}{dT} + \frac{2\lambda_k}{\lambda^2 - \lambda_k^2} \frac{d\lambda_k}{dT} \right] \quad (34)$$

a function of the temperature dependence of the mode strength (A) and its location (λ_k). The mode strength is the product of the strength (polarizability) of the individual absorption oscillators and the density of those oscillators (i.e., is a function of the volume expansion).

The change in n with respect to T is then modeled by using constants D_0 , D_1 , and D_2 to represent the temperature dependence of mode strength, and constants E_0 and E_1 to represent the temperature dependence of mode location, namely,

$$\frac{dn_{\text{abs}}(\lambda, T)}{dT} = \frac{n^2(\lambda, T_0) - 1}{2n(\lambda, T_0)} \left[D_0 + 2D_1 \Delta T + 3D_2 \Delta T^2 + \frac{E_0 + 2E_1 \Delta T}{\lambda^2 - \lambda_k^2} \right] \quad (35)$$

T_0 is the reference temperature, generally 20°C. T is the temperature of interest, and ΔT is the difference, $T - T_0$. The wavelength of interest is λ . The terms D_0 , D_1 , D_2 , E_0 , E_1 , and λ_k are all constants dependent upon the glass type and are found by fitting to measurements. The relative index values found in catalogs can be used in Eq. (35) without loss of accuracy. This is why the abs subscript was dropped from the n terms on the right-hand side of the equation. Using Eq. (35) in Eq. (33) allows for the calculation of the change in the refractive index, as well as the new refractive index as a function of temperature.

The refractive index relative to air as a function of temperature, and the change itself as a function of temperature can be computed from the relationships¹³

$$n_{\text{rel}}(\lambda, T) = \frac{n_{\text{abs}}(\lambda, T)}{n_{\text{air}}(\lambda, T, P)} \quad (36a)$$

$$\frac{dn_{\text{rel}}(\lambda, T)}{dT} = \frac{\frac{dn_{\text{abs}}(\lambda, T)}{dT} - n_{\text{rel}}(\lambda, T) \frac{dn_{\text{air}}(\lambda, T, P)}{dT}}{n_{\text{air}}(\lambda, T, P)} \quad (36b)$$

The approximation that $n_{\text{rel}}(\lambda, T_0)$ can be substituted for $n_{\text{rel}}(\lambda, T)$ maintains sufficient accuracy, and simplifies the calculation.

The widely accepted equations for n_{air} and dn_{air}/dT of dry air are

$$n_{\text{air}}(\lambda, T, P) = 1 + \frac{n_{\text{air}}(\lambda, T_{15^\circ\text{C}}, P_0) - 1}{1 + 3.4785 \cdot 10^{-3}(T - 15^\circ\text{C})} P \quad (37a)$$

$$\frac{dn_{\text{air}}(\lambda, T, P)}{dT} = \frac{-0.00367 n_{\text{air}}(\lambda, T, P) - 1}{1 + 0.00367 T} \quad (37b)$$

where P is the atmospheric pressure divided by the reference pressure ($P_0 = 1$ atmosphere).

The index of refraction in air at 15°C and 1 atmosphere of pressure is given by

$$n_{\text{air}}(\lambda, T_{15^\circ\text{C}}, P_0) = 1 + \left[6432.8 + \frac{2949810\lambda^2}{146\lambda^2 - 1} + \frac{25540\lambda^2}{41\lambda^2 - 1} \right] \times 10^{-8} \quad (37c)$$

The wavelength λ is in units of micrometers, and temperature is in units of degrees Celsius. P is the relative air pressure and is dimensionless. This completes the equations necessary to compute the absolute and relative index of refraction at any temperature and pressure, provided the proportionality constants in Eq. (33) and the parameters of Eq. (35) are available for the glass of interest. These equations are typically valid for temperature ranges of -100°C to 140°C , and a wavelength range of 0.3650 to 1.014 μm .

Athermal Glasses The environmental temperature affects glass in terms of thermal expansion coefficient [α , see Eq. (46a)] and the variation of the index of refraction, dn_{rel}/dT . Both effects deteriorate the wavefront due to the changed optical path and the change in the index of refraction. To mitigate the effect of temperature variations, the glass should ideally compensate for any thermal expansion effects by the refractive index change. Since all glasses exhibit a positive thermal expansion coefficient (increase in size with increased temperature), the ideal glass would have a negative temperature coefficient. This can be expressed in terms of the thermo-optical constant, G :

$$G = \alpha(n_{\text{rel}}(\lambda, T) - 1) + \frac{dn_{\text{rel}}(\lambda, T)}{dT} \quad (38)$$

Any glass with G near zero (e.g., $|G| < 1 \cdot 10^{-6} \text{ K}^{-1}$) is referred to as athermal. Not all glasses with negative relative index temperature coefficients are athermal. Schott glasses N-PK51, N-FK51A, and N-PK52A show athermal behavior.

Nonlinear Optical Coefficients One of the most important higher-order optical coefficients is the nonlinear (or second-order) susceptibility. With the high electric fields generated by lasers, the nonlinear susceptibility gives rise to important processes such as second-harmonic generation, optical rectification, parametric mixing, and the linear electro-optic (Pockels) effect. The second-order susceptibility, $\chi^{(2)}$, is related to the polarization vector \mathbf{P} by

$$\mathbf{P}_i(\omega) = \epsilon_0 [\chi_{ij} \mathbf{E}_j + g\chi_{ijk}^{(2)} \mathbf{E}_j(\omega_1) \mathbf{E}_k(\omega_2)] \quad (39)$$

where g is a degeneracy factor arising from the nature of the electric fields applied. If the two frequencies are equal, the condition of optical rectification and second-harmonic generation (SHG) arises, and $g = 1/2$. When the frequencies of \mathbf{E}_j and \mathbf{E}_k are different, parametric mixing occurs and $g = 1$. If \mathbf{E}_k is a dc field, the situation is the same as the linear electro-optic (or Pockels) effect, and $g = 2$. The value of nonlinear susceptibility is a function of the frequencies of both the input fields and the output polarization ($\omega = \omega_1 \pm \omega_2$).

The nonlinear susceptibility is a third-order (3 by 3 by 3) tensor. A nonlinear optical coefficient $d^{(2)}$, frequently used to describe these nonlinear properties, is equal to one-half of the second-order nonlinear susceptibility, that is, $d^{(2)} = 1/2\chi^{(2)}$. Nonlinear optical coefficients are universally written in reduced (matrix) notation, $d_{ij}^{(2)}$, where the index $i = 1, 2, \text{ or } 3$ and the index j runs from 1 to 6.⁸⁰ (Both the piezo-electric coefficient and the nonlinear optical coefficient are given the symbol d , and the resulting confusion is enhanced because both coefficients have the same units.) The relationship between the electro-optic coefficient r and the nonlinear optical coefficient $d^{(2)}$ is

$$r_{ij} = \frac{2gd_{ij}^{(2)}}{\epsilon^2} \quad (40)$$

Units of the second-order nonlinear optical coefficient are m/V (or pm/V, where pm = 10^{-12}) in mks units.

TABLE 4 Typical Nonlinear Optical Coefficients

Crystal	Nonlinear Optical Coefficient (pm/V)
β -BaB ₂ O ₄	$d_{11} = 1.60$
KH ₂ PO ₄	$d_{36} = 0.39$
LiB ₃ O ₅	$d_{32} = 1.21$
LiNbO ₃	$d_{31} = 5.07$
LiIO ₃	$d_{31} = 3.90$
KTiOPO ₄	$d_{31} = 5.85$
Urea	$d_{14} = 1.17$

Typical values of the nonlinear optical coefficients are listed in Table 4.⁸¹ Additional nonlinear optical coefficients are given in reviews (also refer to Chap. 10, “Nonlinear Optics” by Chung L. Tang).⁸⁰⁻⁸³

Scatter Scatter is both an intrinsic and extrinsic property. Rayleigh, Brillouin, Raman, and stoichiometric (index variation) contributions to scatter have been derived in simple form and used to estimate scatter loss in several fiber-optic materials.⁸⁴ Rayleigh scattering refers to elastic scatter from features small compared to the wavelength of light. In highly pure and defect-free optical crystals, Rayleigh scatter is caused by atomic scale inhomogeneities (much smaller than the wavelength) in analogy to Rayleigh scatter from molecules in the atmosphere. In most materials, including glasses, Rayleigh scatter is augmented by extrinsic contributions arising from localized density variations (which also limit the uniformity of the refractive index). Attenuation in high-quality optical materials is frequently limited by Rayleigh scatter rather than absorption.

Brillouin and Raman scatter are forms of inelastic scattering from acoustic and optical phonons (vibrations). The frequency of the scattered light is shifted by the phonon frequency. Creation of phonons results in longer-wavelength (low-frequency) scattered light (Stokes case) and annihilation of phonons results in higher-frequency scattered light (anti-Stokes case). Rayleigh, Brillouin, and Raman scatter all have λ^{-4} wavelength dependence. Polycrystalline and translucent materials have features such as grain boundaries and voids whose size is larger than the wavelength of light. This type of scatter is often called Mie scatter because the scattering features are larger than the wavelength of the light. Mie scatter typically has a measured λ^{-m} dependence where the parameter m typically lies between 1 and 2.^{85,86} Rayleigh and Mie scatter may arise from either surface roughness or bulk nonuniformities.

Other Properties of Materials

Naming of Crystals and Glasses All materials are characterized by name(s) for identification, a chemical formula (crystalline materials) or approximate composition (glasses, amorphous substances), and a density (ρ , in kg/m³). Crystalline materials are further identified by crystal class, space group, unit-cell lattice parameters, molecular weight (of a formula unit in atomic mass units, amu), and number of formula units per unit cell (Z). (See standard compilations of crystallographic data.^{87,88})

Material Designation and Composition Crystals are completely identified by both the chemical formulation and the space group. Chemical formulation alone is insufficient for identification because many substances have several structures (called *polymorphs*) with different properties. Properties in the data tables pertain only to the specific structure listed. Materials in the data tables having several stable polymorphs at room temperature include SiO₂ (eight polymorphs), C (diamond, graphite, and amorphous forms), and SiC and ZnS (both have cubic and hexagonal forms).

The space group also identifies the appropriate number of independent terms (see Tables 1 and 2) that describe a physical quantity. Noncubic crystals require two or more values to fully describe thermal expansion, thermal conductivity, refractive index, and other properties. Often, scalar quantities are

given in the literature when a tensor characterization is needed. Such a characterization may be adequate for polycrystalline materials, but is unsatisfactory for single crystals that require knowledge of directional properties.

Historically, optical glasses have been identified by traditional names derived from their composition and their dispersion relative to their index of refraction. This dates back to the 1880s when low dispersion crown glasses (typically with Abbe number $v_d > 50$) and higher dispersion flint glasses (typically, $v_d < 50$) dominated the glass map.⁸⁹ It took Ernst Abbe, Carl Zeiss, and the work of Otto Schott to expand the glass by developing two new glass formers, fluorine and boron.^{90,91} In addition, Schott's work with BaO network modifiers pioneered new variations in index of refraction within the same groups. This gave rise to classifications of "heavy" and "light" (index of refraction designation), in addition to crown and flint (Abbe number classification).

The next major advancement in glass technology came in the 1930s when rare earth elements and network-former phosphorous became readily available.^{92,93} These new glass types were named by adding the compounds chemical symbol to the traditional nomenclature of crown or flint, heavy or light. This led to various groups evolving on the Abbe diagram, groups whose position gave strong indication of the chemical composition of the glass. This distinction and borders in the Abbe diagram remains till the present day. However the traditional nomenclature is rapidly ending as manufacturers, primarily in Japan, are adopting new naming conventions for their glasses.

The primary reason for this nomenclature change is advances in manufacturing techniques, raw material supply and purity, and the evolution of new materials have all led to new compositional families of glasses that cross traditional borders. So a given glass family name such as PSK may have glasses without any phosphorous. Furthermore, different manufacturers producing a glass with identical index of refraction and Abbe numbers may in fact have entirely different physical and chemical properties (e.g., Schott N-LaSF31A, Hikari's E-LaSF08, and Hoya's TaFD30.) A single name would be disingenuous.

The main classical groups are listed in Table 8, and a representative glass from each is shown. Figure 9 shows the classic Abbe diagram with traditional nomenclature and the newer style showing glass systems based more on constituent materials. A thorough discussion of the various glass systems is found in the literature.¹⁷

Regardless of the categorizing and naming of glasses and glass systems or groups, a specific glass identifier scheme, as defined in military standard MIL-G-174, remains in use today. This convention uses a six-digit number representing the first three digits of $(n_d - 1)$ and the first three digits of v_d . Each manufacturer also has its own designator, usually based on traditional or new names, that uniquely identifies each glass. For example, the glass with code 517624 has the following manufacturer's designations:

Manufacturer	Designation for Glass 517642
Schott	N-BK7
Corning	B-16-64
Pilkington	BSC-517642B
Hoya	BSC-7
Ohara	S-BSL 7 (closest value 516641)

Properties of glass are primarily determined from the compositions, but also depend on the manufacturing process, specifically the thermal history. In fact, refractive index specifications in a glass catalog should be interpreted as those obtained with a particular annealing schedule. Annealing removes stress (and minimizes stress-induced birefringence) and minimizes the effect of thermal history, producing high refractive-index uniformity. Special (*precision*) annealing designed to maximize refractive index homogeneity may, however, increase refractive index slightly above a nominal (catalog) value. A complete discussion of the annealing process can be found in literature.^{16,94}

In general, both glass composition and thermal processing are proprietary (so the compositions given are only illustrative). However, manufacturers' data sheets on individual glasses can provide detailed and specific information on optical and mechanical properties. Also, the data sheets supply useful details on index homogeneity, climate resistance, stain resistance, and chemical (acid and

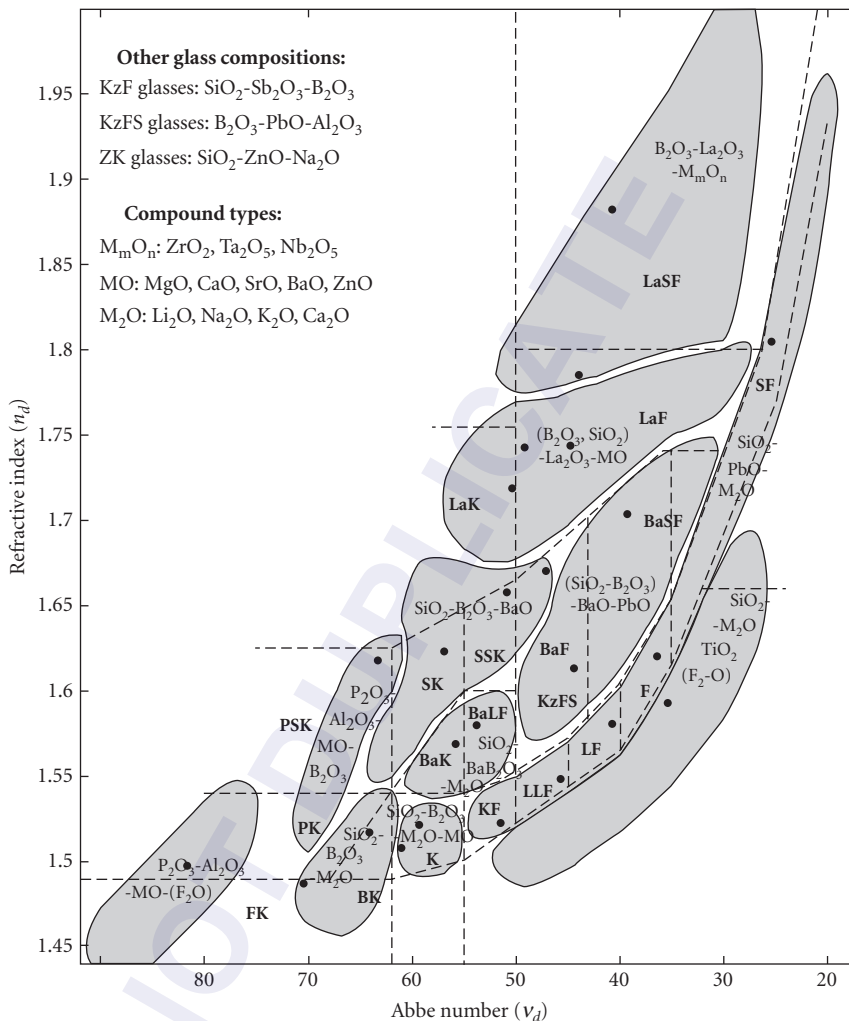


FIGURE 9 The Abbe diagram (refractive index, n_d versus Abbe number, v_d) showing the composition of common glass systems, the Schott glass classifications, and the sample glasses (dots) included in this chapter. (Reprinted by permission of Springer-Verlag.)

alkaline) resistance for a particular glass type. For very detailed work or demanding applications, a glass manufacturer can supply a *melt data sheet* providing accurate optical properties for a specific glass lot.

A driving force in the recent breaking of traditional glass manufacture practices has been a growing movement for environmentally responsible glass. The primary elements of contention were arsenic and lead, and other toxic elements. For example, thorium oxide use has been discontinued, and cadmium oxide is restricted to colored glass manufacture. The principal concerns lie not in the end user, despite the negative connotation lead and lead-based products have incurred the past few decades. In fact, lead-containing glass is still available and labeled as such. The concern rests in the manufacture and disposal of hazardous waste. No guarantee could be made that the workplace would be safe, and the waste disposed of properly.

The replacements for PbO are typically TiO₂, Nb₂O₅, ZrO₂, and WO₃. The replacement for As₂O₃ as a refining agent has been Sb₂O₃. For approximately 50 percent of the glasses containing As₂O₃ and no PbO, the use of Sb₂O₃ has not changed the properties of the glass. Using substitute materials for Glasses formerly containing PbO and As₂O₃ has modified the properties of the glass, sometimes significantly. One by-product of lead removal has been lighter weight. Manufacturers have begun to denote the reformulated glasses by specific designations. Schott uses “N-” to prefix the new glass materials, Hikari uses “E-”, and Ohara uses “S-”. References to environmental concerns are included in the reference section.^{10,95}

Unit Cell Parameters, Molecular Weight, and Density The structure and composition of crystals can be used to calculate density. This calculated (theoretical or x-ray) density should closely match that of optical-quality materials. Density, ρ , is mass divided by volume:

$$\rho = \frac{Z \cdot (\text{MW}) \cdot u}{a \cdot b \cdot c \sqrt{\sin^2 \alpha + \sin^2 \beta + \sin^2 \gamma - 2(1 - \cos \alpha \cdot \cos \beta \cdot \cos \gamma)}} \quad (41)$$

where Z is the number of formula units in a crystal unit cell, MW is the molecular weight of a formula unit in amu, u is weight of an amu (Table 3), a , b , and c are unit cell axes lengths, and α , β , and γ are unit cell axes angles.

Typically, pure amorphous materials have lower density than the corresponding crystalline materials. Density of glasses and other amorphous materials is derived from measurements.

Elastic Properties Elastic properties of materials can be described with a hierarchy of terms. On the atomic scale, interatomic force constants or potential energies can be used to predict the vibrational modes, thermal expansion, and elastic properties of a material. On the macroscopic scale, elastic properties are described using elastic moduli (or constants) related to the directional properties of a material. The tensor relationships between stress (σ , a second-order tensor) and strain (e , a second-order tensor) are

$$\begin{aligned} \sigma_{ij} &= c_{ijkl} \cdot e_{kl} \\ e_{ij} &= s_{ijkl} \cdot \sigma_{kl} \end{aligned} \quad (42a)$$

where the fourth-rank tensors c_{ijkl} and s_{ijkl} are named elastic stiffness c and elastic compliance s , respectively. This is the tensor form of Hooke's law. Each index (i , j , k , and l) has three values (i.e., x , y , and z), hence the c and s tensors have 81 terms.

The stiffness and compliance tensors are usually written in a matrix notation made possible by the symmetry relationship of the stress and strain tensors. Symmetry reduces the number of independent terms in the stiffness and compliance tensors from 81 to 36. The usual notation for the reduced (matrix) notation form of the stiffness and compliance tensors is

$$\begin{aligned} \sigma_i &= c_{ij} \cdot e_j \\ e_i &= s_{ij} \cdot \sigma_j \end{aligned} \quad (42b)$$

where the indices, an abbreviation of the ij or kl components, run from 1 to 6. Table 5 shows the conversion from tensor to matrix notation. Thus, the stiffness and compliance tensors are written as 6 by 6 matrices which can again be shown to be symmetric, given 21 independent terms. Virtually all data will be found in matrix notation. These tensors (matrices) that relate stress and strain are sometimes called *second-order* stiffness and compliance. Higher-order tensors are used to describe nonlinear elastic behavior (i.e., third-order stiffness determines the stress tensor from the square of the strain tensor).

Stiffness and compliance tensors are needed to completely describe the linear elastic properties of a crystal. Even a completely amorphous material has two independent constants that describe the relationship between stress and strain. Usually, the elastic properties of materials are expressed in terms of engineering (or technical) moduli: Young's modulus (E), shear modulus (or modulus

TABLE 5 Matrix Notation for Stress, Strian, Stiffness, and Compliance Tensors

Tensor-to-matrix index conversion		Tensor-to-matrix element conversion	
Tensor Indices ij or kl	Matrix Indices m or n	Notation	Condition
11	1	$\sigma_m = \sigma_{ij}$	$m = 1, 2, 3$
22	2	$e_m = \sigma_{ij}$ $\sigma_m = \sigma_{ij}$	$m = 4, 5, 6$
33	3	$e_m = 2e_{ij}$	all m, n
23 or 32	4	$c_{mn} = c_{ijkl}$	$m, n = 1, 2, 3$
13 or 31	5	$s_{mn} = s_{ijkl}$	$m = 1, 2, 3$ and $n = 4, 5, 6$
12 or 21	6	$s_{mn} = 4s_{ijkl}$	$m = 4, 5, 6$ and $n = 1, 2, 3$

of rigidity G), bulk modulus (B , compressibility⁻¹), and Poisson's ratio (ν). For example, Young's modulus is defined as the ratio of the longitudinal tension to the longitudinal strain for tension, a quantity which is anisotropic (i.e., directionally dependent) for all crystal classes (but is isotropic for amorphous materials). Therefore the engineering moduli only accurately describe the elastic behavior of isotropic materials. The engineering moduli also approximately describe the elastic behavior of polycrystalline materials (assuming small, randomly distributed grains). Various methods are available to estimate the engineering moduli of crystals.

Values of the engineering moduli for crystalline materials given in the data tables are estimated from elastic moduli using the Voigt and Reuss methods (noncubic materials) or the Haskin and Shtrikman method⁹⁶ (cubic materials) to give shear (G) and bulk (B) moduli. Young's modulus (E) and Poisson's ratio (ν) are then calculated assuming isotropy using the following relationships:

$$E = \frac{9 \cdot G \cdot B}{G + 3 \cdot B} \quad \nu = \frac{3 \cdot B - 2 \cdot G}{6 \cdot B + 2 \cdot G} \quad (43)$$

Hardness and Strength Hardness is an empirical and relative measure of a material's resistance to wear (mechanical abrasion). Despite the qualitative nature of the result, hardness testing is quantitative, repeatable, and easy to measure. The first measure of hardness was the Mohs scale which compares the hardness of materials to one of 10 minerals. Usually, the Knoop indent test is used to measure hardness of optical materials. The test determines the resistance of a surface to penetration by a diamond indenter with a fixed load (usually 50 to 200 g). The Knoop hardness number (usually measured in $\text{kg}/\text{mm}^2 = 9.8 \text{ MPa}$) is the indenter mass (proportional to load) divided by the area of the indent. Figure 10 compares the Mohs and Knoop scales.⁹

Materials with Knoop values less than $100 \text{ kg}/\text{mm}^2$ are very soft, difficult to polish, and susceptible to handling damage. Knoop hardness values greater than 750 are quite hard. Typical glasses have hardness values of 350 to $600 \text{ kg}/\text{mm}^2$. Hardness qualitatively correlates to Young's modulus and to strength. Hardness of crystals is dependent on the orientation of the crystal axes with respect to the tested surface. Coatings can significantly alter hardness.

Strength is a measure of a material's resistance to fracture (or onset of plastic deformation). Strength is highly dependent on material flaws and therefore on the method of manufacture, as well as the method of measurement. For optical materials, strength is most conveniently measured in flexure; tensile strengths are typically 50 to 90 percent of those measured in flexure. Because of high variability in strength values, quoted strength values should only be used as a guide for comparison of materials. Strength of crystals also is dependent on the orientation of the crystal axes with respect to the applied stress. Applications requiring high strength to avoid failure should use large safety margins (typically a factor of four) over average strength whenever possible.

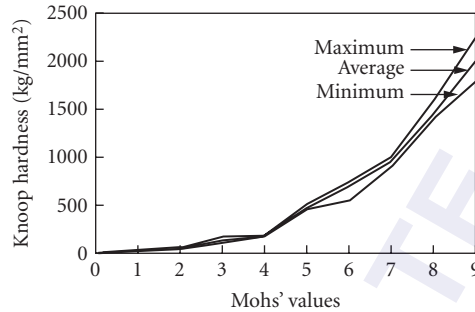


FIGURE 10 A comparison of Mohs and Knoop hardness scales.¹² The Mohs scale is qualitative, comparing the hardness of a material to one of 10 minerals: talc (Moh \equiv 1), gypsum (\equiv 2), calcite (\equiv 3), fluorite (\equiv 4), apatite (\equiv 5), orthoclase (\equiv 6), quartz (\equiv 7), topaz (\equiv 8), sapphire (\equiv 9), and diamond (Moh \equiv 10). The Knoop scale is determined by area of a mark caused by an indenter; the Knoop value is the indenter mass divided by the indented area. The mass of the indenter (load) is usually specified; 200 or 500 g are typical. (Reprinted by permission of Ashlee Publishing Company.)

Fracture toughness is another measure of strength, specifically, a material's ability to resist crack propagation. Fracture toughness measures the applied stress required to enlarge a flaw (crack) of given size and has units of $\text{MPa} \cdot \text{m}^{1/2}$. Values for representative materials are given in Table 6.

Characteristic Temperatures Characteristic temperatures of crystalline materials are those of melting (or vaporization or decomposition) and phase transitions. Of particular importance are the phase-transition temperatures. These temperatures mark the boundaries of a particular structure. A phase transition can mean a marked change in properties. One important phase-transition point is the Curie temperature of ferroelectric materials. Below the Curie temperature, the material is ferroelectric; above this temperature it is paraelectric. The Curie temperature phase transition is particularly significant because the change in structure is accompanied by drastic changes in some properties such as the static dielectric constant which approaches infinity as temperature nears the Curie temperature. This transition is associated with the lowest transverse optical frequency (the soft mode) approaching zero [hence the static dielectric constant approaches infinity from the Lyddane-Sachs-Teller relationship, Eq. (11)].

TABLE 6 Fracture Toughness of Some Materials

Fracture toughness		Fracture toughness	
Material	$\text{MPa} \cdot \text{m}^{1/2}$	Material	$\text{MPa} \cdot \text{m}^{1/2}$
Al_2O_3	3	MgAl_2O_4	1.5
ALON	1.4	MgF_2	1.0
AlN	3	PbTiO_3	1.1
N-BK7 glass	1.1	SF58 glass	0.38
C, diamond	2.0	Si	0.95
CaF_2	0.5	fused SiO_2	0.8
CaLa_2S_4	0.68	Y_2O_3	0.7
F2 glass	0.55	ZnS	0.5 (crystal) 0.8 (CVD)
GaP	0.9	ZnSe	0.33
Ge	0.66	$\text{ZrO}_2\text{:Y}_2\text{O}_3$	2.0
LaK10 glass	0.95	Zerodur	0.9

The term *glass* applies to a material that retains an amorphous state upon solidification. More accurately, glass is an undercooled, inorganic liquid with a very high viscosity at room temperature and is characterized by a gradual softening with temperature and a hysteresis between glass and crystalline properties. The gradual change in viscosity with temperature is characterized by several temperatures, especially the glass transition temperature and the softening-point temperature. The glass transition temperature defines a second-order phase transition analogous to melting. At this temperature, the temperature dependence of various properties changes (in particular, the linear thermal expansion coefficient) as the material transitions from a liquid to glassy state. Glasses can crystallize if held above the transition temperature for sufficient time. The annealing point is defined as the temperature resulting in a glass viscosity of 10^{13} poise and at which typical glasses can be annealed within an hour or so. In many glasses, the annealing point and the glass transition temperature are close. In an optical system, glass elements need to be kept 150 to 200°C below the glass transition temperature to avoid significant surface distortion. At the *softening temperature*, viscosity is $10^{7.6}$ poise and glass will rapidly deform under its own weight; glasses are typically molded at this temperature. Glasses do not have a true melting point; they become progressively softer (more viscous) with increased temperature. Other amorphous materials may not have a well-defined glass transition; instead they may have a conventional melting point. Glasses that crystallize at elevated temperature also have a well-defined melting point.

Glass-ceramics have been developed which are materials with both glasslike and crystalline phases. In particular, low-thermal-expansion ceramics comprise a crystalline phase with a negative thermal expansion and a vitreous phase with a positive thermal expansion. Combined, the two phases result in very high dimensional stability. Typically, the ceramics are made like other glasses, but after stresses are removed from a blank, a special heat-treatment step forms the nuclei for the growth of the crystalline component of the ceramic. Although not strictly ceramics, similar attributes can be found in some two-phase glasses.

Heat Capacity and Debye Temperature Heat capacity, or specific heat, a scalar quantity, is the change in thermal energy with a change in temperature. Units are typically J/(gm · K). Debye developed a theory of heat capacity assuming that the energy was stored in acoustical phonons. This theory, which assumes a particular density of states, results in a Debye molar heat capacity [units J/(mole · K)] of the form

$$C_V(T) = 9mN_A k_B \left(\frac{T}{\theta_D} \right)^3 \int_0^{\theta_D/T} \frac{x^4 e^x}{(e^x - 1)^2} dx \quad (44a)$$

where C_V is the molar heat capacity in J/(mole · K) per unit volume, θ_D is the Debye temperature, m the number of atoms per formula unit, N_A is Avogadro's number, and k_B is Boltzmann's constant. At low temperatures ($T \rightarrow 0$ K), heat capacity closely follows the T^3 law of Debye theory

$$C_V(T) = \frac{12\pi^4}{5} mN_A k_B \left(\frac{T}{\theta_D} \right)^3 = 1943.76 \cdot m \left(\frac{T}{\theta_D} \right)^3 \text{ J/(mole} \cdot \text{K)} \quad (T \ll \theta_D) \quad (44b)$$

and the high-temperature (classical) limit is

$$C_V(T) = 3mN_A k_B = 24.943 \cdot m \text{ J/(mole} \cdot \text{K)} \quad (T > \theta_D) \quad (44c)$$

If heat capacity data are fit piecewise to the Debye equation, a temperature-dependent θ_D can be found. Frequently, a Debye temperature is determined from room-temperature elastic constants, and is therefore different from the low-temperature value. The Debye temperature given in the tables is, when possible, derived from low-temperature heat capacity data.

The Debye equations can be used to estimate heat capacity C_V over the entire temperature range, typically to within 5 percent of the true value using a single Debye temperature value. Usually, however, C_p , the constant pressure heat capacity, rather than C_V , is desired. At low temperatures, thermal

expansion is small, and $C_p \approx C_v$. At elevated temperature, the relationship between C_v and C_p is given by the thermodynamic relationship

$$C_p(T) = C_v(T) + 9\alpha^2 TVB \quad (45)$$

where T is temperature (K), V is the molar volume (m^3/mole), α is the thermal expansion coefficient, and B is the bulk modulus ($\text{Pa} = \text{Nt}/\text{m}^2$).

Molar heat capacity can be converted to usual units by dividing by the molecular weight (see the physical property tables for crystals) in g/mole . Since molar heat capacity approaches the value of $24.943 \text{ J}/(\text{mole} \cdot \text{K})$ the heat capacity per unit weight is inversely proportional to molecular weight at high temperature (i.e., above the Debye temperature).

Thermal Expansion The linear thermal expansion coefficient α is the fractional change in length with a change in temperature as defined by

$$\alpha(T) = \frac{1}{L} \frac{dL}{dT} \quad (46a)$$

and units are $1/\text{K}$. The units of length are arbitrary. Thermal expansion is a second-rank tensor; nonisometric crystals have a different thermal expansion coefficient for each principal direction. At low temperature, thermal expansion is low, and the coefficient of thermal expansion approaches zero as $T \rightarrow 0 \text{ K}$. The expansion coefficient generally rises with increasing temperature; Fig. 11 shows temperature dependence of the expansion coefficient for several materials. Several compilations of data exist.^{97,98}

The volume expansion coefficient α_v is the fractional change in volume with an increase in temperature. For a cubic or isotropic material with a single linear thermal expansion coefficient

$$\alpha_v(T) = \frac{1}{V} \frac{dV}{dT} = 3\alpha \quad (46b)$$

can be used to estimate the temperature change of density.

The Grüneisen relationship relates the thermal expansion coefficient to molar heat capacity

$$\alpha(T) = \frac{\gamma C_v(T)}{3B_0 V_0} \quad (46c)$$

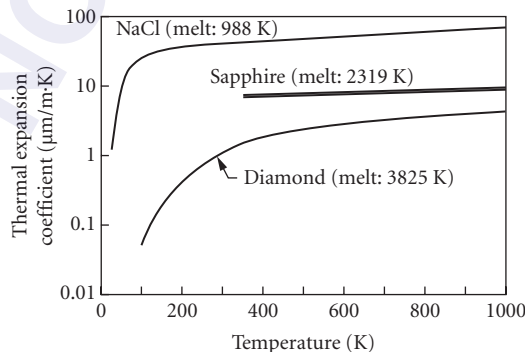


FIGURE 11 Thermal expansion of several materials. Expansion arises from the anharmonicity of the interatomic potential. At low temperature, expansion is very low and the expansion coefficient is low. As temperature increases, the expansion coefficient rises, first quickly, then less rapidly.

where B_0 is the bulk modulus at $T = 0$ K, V_0 is the volume at $T = 0$ K, and γ is the Grüneisen parameter. This relationship shows that thermal expansion has the same temperature dependence as the heat capacity. Typical values of γ lie between 1 and 2.

Thermal Conductivity Thermal conductivity κ determines the rate of heat flow through a material with a given thermal gradient. Conductivity is a second-rank tensor with up to three principal values. This property is especially important in relieving thermal stress and optical distortions caused by rapid heating or cooling. Units are $W/(m \cdot K)$

Kinetic theory gives the following expression for thermal conductivity, κ

$$\kappa = \frac{1}{3} C_V v \ell \quad (47)$$

where v is the phonon (sound) velocity and ℓ is the phonon mean free path. At very low temperature ($T < \theta_D/20$), the temperature dependence of thermal conductivity is governed by C_V , which rises as T^3 [see Eq. (36b)]. At high temperature ($T > \theta_D/10$), the phonon mean free path is limited by several mechanisms. In crystals, scattering by other phonons usually ℓ . In the high-temperature limit, the phonon density rises proportional to T and thermal conductivity is inversely proportional to T . Figure 12 illustrates the temperature dependence of thermal conductivity of several crystalline materials.

Thermal conductivity in amorphous substances is quite different compared to crystals. The phonon mean free path in glasses is significantly less than in crystals, limited by structural disorder. The mean free path of amorphous materials is typically the size of the fundamental structural units (e.g., 10 Å) and has little temperature dependence; hence the temperature dependence of thermal conductivity is primarily governed by the temperature dependence of heat capacity. At room temperature, the thermal conductivity of oxide glasses is a factor of 10 below typical oxide crystals. Figure 12 compares the thermal conductivity of fused silica to crystalline silica (quartz).

Thermal conductivity of crystals is highly dependent on purity and order. Mixed crystals, second-phase inclusion, nonstoichiometry, voids, and defects can all lower the thermal conductivity

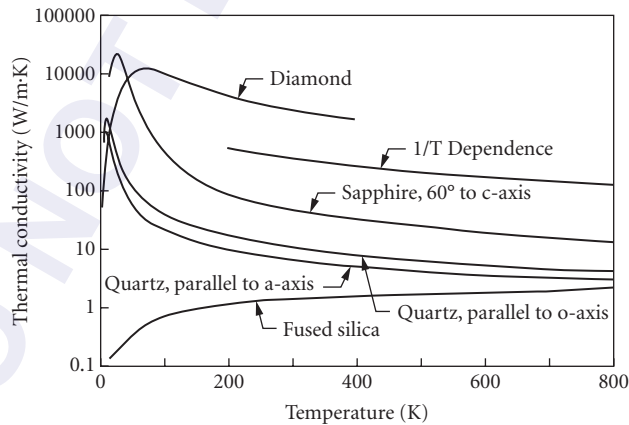


FIGURE 12 Thermal conductivity of several materials. Conductivity initially rises rapidly as the heat capacity increases. Peak thermal conductivity of crystals is high due to the long phonon mean free path of the periodic structure which falls with increasing temperature as the phonon free path length decreases ($\approx 1/T$). The phonon mean free path of amorphous materials (e.g., fused silica) is small and nearly independent of temperature, hence thermal conductivity rises monotonically and approaches the crystalline value at high temperature.

of a material. Values given in the data table are for the highest-quality material. Thermal conductivity data are found in compilations⁹⁸ and reviews.⁹⁹

Correlation of Properties All material properties are correlated to a relatively few factors, for example, constituent atoms, the bonding between the atoms, and the structural symmetry. The binding forces, or chemical bonds, play a major role in properties. Tightly bonded materials have high moduli, high hardness and strength, and high Debye temperature (hence high room-temperature thermal conductivity). Strong bonds also mean lower thermal expansion, lower refractive index, higher-frequency optical vibrational modes (hence less infrared transparency), and higher energy bandgaps (hence more ultraviolet transparency). Increased mass of the constituent atoms lowers the frequency of both electronic and ionic resonances. Similarly, high structural symmetry can increase hardness and eliminate (to first order) ionic vibrations (cf., diamond, silicon, and germanium).

Combinations of Properties A given material property is influenced by many factors. For example, the length of a specimen is affected by stress (producing strain), by electric fields (piezo-electric effect), and by temperature (thermal expansion). The total strain is then a combination of these linear effects and can be written as

$$\Delta e_{ij} = \left(\frac{\partial e_{ij}}{\partial \sigma_{kl}} \right)_{E,T} \Delta \sigma_{kl} + \left(\frac{\partial e_{ij}}{\partial E_k} \right)_{\sigma,T} \Delta E_k + \left(\frac{\partial e_{ij}}{\partial T} \right)_{\sigma,E} \Delta T \quad (48)$$

$$\Delta e_{ij} = s_{ijkl} \Delta \sigma_{kl} + d_{ijk} \Delta E_k + \alpha_{ij} \Delta T$$

Often, these effects are interrelated, and frequently dependent on measurement conditions. Some properties of some materials are very sensitive to measurement conditions (the subscripts in the preceding equation denote the variable held constant for each term). For example, if the measurement conditions for the elastic contribution were adiabatic, stress will cause temperature to fall, which, in turn, decreases strain (assuming positive thermal expansion coefficient). Thus the elastic contribution to strain is measured under isothermal (and constant \mathbf{E} field) conditions so as not to include the temperature effects already included in the thermal expansion term.

The conditions of measurement are given a variety of names that may cause confusion. For example, the mechanical state of “clamped,” constant volume, and constant strain all refer to the same measurement condition which is paired with the corresponding condition of “free,” “unclamped,” constant pressure, or constant stress. In many cases, the condition of measurement is not reported and probably unimportant (i.e., different conditions give essentially the same result). Another common measurement condition is constant \mathbf{E} field (“electrically free”) or constant \mathbf{D} field (“electrically clamped”).

Some materials, particularly ferroelectrics, have large property variation with temperature and pressure, hence measurement conditions may greatly alter the data. The piezoelectric effect contributes significantly to the clamped dielectric constant of ferroelectrics. The difference between the isothermal clamped and free dielectric constants is

$$\epsilon_0 \cdot (\epsilon_i^c - \epsilon_i^\sigma) = -d_{ij} d_{ik} c_{jk}^E \quad (49)$$

where d_{ij} is the piezo-electric coefficient and c_{jk}^E is the (electrically free) elastic stiffness. If the material structure is centrosymmetric, all components of d_{ij} vanish, and the two dielectric constants are the same.

2.5 PROPERTIES TABLES

The following tables summarize the basic properties for representative crystals and glasses. In general, the presented materials are (1) of general interest, (2) well-characterized (within the limitations imposed by general paucity of data and conflicting property values), and (3) represent a wide

range of representative types and properties. Few materials can be regarded as well-characterized. Crystalline materials are represented by alkali halides, oxides, chalcogenides, and a variety of crystals with nonlinear optical, ferroelectric, piezoelectric, and photorefractive properties.

A number of materials in the previous edition have been removed for the sake of brevity. Properties of the following crystals and glasses are found in this chapter:¹

β -AgI (iodyrite)	InAs	PbS
AlAs	InP	PbSe
AlN	KCl	Se
BN	KF	β -SiC
BP	LaF ₃	α -SiC
CaLa ₂ S ₄	NaBr	Te
CsBr	NaF (valliaumite)	Tl ₃ AsSe ₃ (TAS)
CsCl	NaI	TlCl
CuCl (nantokite)	PbF ₂	ZnO (zincite)
TiK1	TiF1	PK2
BaLK1	KzF6	ZK1
PSK3	BaK1	LgSK2
SK4	SSK4	BaSF10
TaC2	TaD5	ULTRAN 30 (548743)
CORTRAN 9753	ZBL	ZBLA
ZBT	HBL	HBLA
HBT		

The physical property tables define the composition, density, and structure (of crystalline materials). Table 7 gives data for 71 crystalline materials. Table 8 compares glasses from different manufacturers and selects 24 representative “optical” glasses intended for visible and near-infrared use (typically to 2.5 μm). Table 9 gives physical property data for 17 specialty glasses and substrate materials. The specialty glasses include fused silica and germania, calcium aluminate, fluoride, germanium-, and chalcogenide-based glasses, many of which are intended for use at longer wavelengths. The three substrate materials, Pyrex, Zerodur, and ULE, are included because of their widespread use for mirror blanks.

Mechanical properties for crystals are given in two forms, room-temperature elastic constants (or moduli) for crystals (Tables 10 through 16), and engineering moduli, flexure strength, and hardness for both crystals (Table 17) and glasses (Table 18). Engineering moduli for crystalline materials should only be applied to the polycrystalline forms of these materials. Accurate representation of the elastic properties of single crystals requires the use of elastic constants in tensor form. Strength is highly dependent on manufacture method and many have significant sample-to-sample variability. These characteristics account for the lack of strength data. For these reasons, the provided strength data is intended only as a guide. Glasses and glass-ceramics flexure strengths typically range from 30 to 200 MPa, although glass fibers with strength exceeding 1000 MPa have been reported.

Thermal properties are given in Tables 19 and 20 for crystals and glasses, respectively. Characteristic temperatures (Debye, phase change, and melt for crystals; glass transitions, soften, and melt temperatures for glasses), heat capacity, thermal expansion, and thermal conductivity data are included. Directional thermal properties of crystals are given when available. Only room-temperature properties are reported except for thermal conductivity of crystals, which is also given for temperatures above and below ambient, if available.

Optical properties are summarized in Tables 21 and 22 for crystals and glasses, respectively. These tables give the wavelength boundaries of the optical transparent region (based on a 1 cm^{-1} absorption coefficient), characteristic refractive index (\bar{n}_{∞} , the asymptote of the electronic contribution to the refractive index for crystals), or n_d and v_d for glasses), and values of dn/dT at various wavelengths. Tables 23 and 24 give dispersion formulas for crystals and glasses, respectively.

TABLE 7 Composition, Structure, and Density of Crystals

Material	Crystal System and Space Group	Unit Cell Dimension (Å)	Molecular Weight (amu)	Formulae/ Unit Cell	Density (g/cm ³)
Ag ₃ AsS ₃ (proustite)	Hexagonal R3c (C _{3v} ⁶) #161	a = 10.756 c = 8.652	494.72	6	5.686
AgBr (bromyrite)	Cubic Fm3m (O _h ⁵) #225	5.7745	187.77	4	6.477
AgCl (cerargyrite)	Cubic Fm3m (O _h ⁵) #225	5.547	143.32	4	5.578
AgGaS ₂ (AGS)	Tetragonal 142d (D _{2d} ¹²) #122	a = 5.757 c = 10.304	241.72	4	4.701
AgGaSe ₂	Tetragonal 142d (D _{2d} ¹²) #122	a = 5.992 c = 10.886	335.51	4	5.702
Al ₂ O ₃ (sapphire, alumina)	Hexagonal R3c (D _{3d} ⁶) #167	a = 4.759 c = 12.989	101.96	6	3.987
Al ₂₃ O ₂₇ N ₅ , ALON	Cubic Fd3m (O _h ⁷) #227	7.948	1122.59	1	3.713
Ba ₃ [B ₃ O ₆] ₂ , BBO	Hexagonal R3 (C ₃ ⁴) #146	a = 12.532 c = 12.726	668.84	6	3.850
BaF ₂	Cubic Fm3m (O _h ⁵) #225	6.2001	175.32	4	4.886
BaTiO ₃	Tetragonal P4 ₂ /mmn (D _{4h} ¹⁴) #136	a = 3.9920 c = 4.0361	233.19	1	6.020
BeO (bromellite)	Hexagonal P6 ₃ mc (C _{6v} ⁴) #186	a = 2.693 c = 4.395	25.012	2	3.009
Bi ₁₂ GeO ₂₀ , BGO	Cubic I23 (T ³) # 197	10.143	2900.39	2	9.231
Bi ₁₂ SiO ₂₀ , BSO (selenite)	Cubic I23 (T ³) # 197	10.1043	2855.84	2	9.194
BiB ₃ O ₆ , BIBO	Monoclinic C2 (C ₂ ³) #5	a = 7.1203 b = 4.9948 c = 6.5077 γ = 105.59°	337.41	2	5.027
C (diamond)	Cubic Fd3m (O _h ⁷) #227	3.56696	12.011	8	3.516
CaCO ₃ (calcite)	Hexagonal R3c (D _{3d} ⁶) #167	a = 4.9898 c = 17.060	100.09	6	2.711
CaF ₂ (fluorite)	Cubic Fm3m (O _h ⁵) #225	5.46295	78.07	4	3.181
CaMoO ₄ (powellite)	Tetragonal I4 ₁ /a (C _{4h} ⁶) #88	a = 5.23 c = 11.44	200.04	4	4.246
CaWO ₄ (scheelite)	Tetragonal I4 ₁ /a (C _{4h} ⁶) #88	a = 5.243 c = 11.376	287.92	4	6.115
CdGeAs ₂	Tetragonal 142d (D _{2d} ¹²) #122	a = 5.9432 c = 11.216	334.89	4	5.615
CdS (greenockite)	Hexagonal P6 ₃ mc (C _{6v} ⁴) #186	a = 4.1367 c = 6.7161	144.48	2	4.821
CdSe	Hexagonal P6 ₃ mc (C _{6v} ⁴) #186	a = 4.2972 c = 7.0064	191.37	2	5.672

TABLE 7 Composition, Structure, and Density of Crystals (*Continued*)

Material	Crystal System and Space Group	Unit Cell Dimension (Å)	Molecular Weight (amu)	Formulae/ Unit Cell	Density (g/cm ³)
CdTe	Cubic F $\bar{4}3m$ (T $_d^2$) #216	6.4830	240.01	4	5.851
CsLiB ₆ O ₁₀ , CLBO	Tetragonal I $\bar{4}2d$ (D $_{2d}^{12}$) #122	a = 10.494 c = 8.939	492.70	4	2.461
CsI	Cubic Pm $\bar{3}m$ (O $_h^5$) #221	4.566	259.81	1	4.532
CuGaS ₂	Tetragonal I $\bar{4}2d$ (D $_{2d}^{12}$) #122	a = 5.351 c = 10.480	197.40	4	4.369
GaAs	Cubic F $\bar{4}3m$ (T $_d^2$) #216	5.65325	144.64	4	5.317
GaN	Hexagonal P6 $_3$ mc (C $_{6v}^4$) #186	a = 3.186 c = 5.178	83.73	2	6.109
GaP	Cubic F $\bar{4}3m$ (T $_d^2$) #216	5.4495	100.70	4	4.133
Ge	Cubic Fd $\bar{3}m$ (O $_h^7$) #227	5.65741	72.64	8	5.329
KBr	Cubic Fm $\bar{3}m$ (O $_h^5$) #225	6.600	119.00	4	2.749
KH ₂ PO ₄ , KDP	Tetragonal I $\bar{4}2d$ (D $_{2d}^{12}$) #122	a = 7.4529 c = 6.9751	136.09	4	2.338
KI	Cubic Fm $\bar{3}m$ (O $_h^5$) #225	7.065	166.00	4	3.127
KNbO ₃	Orthorhombic Bmm2 (C $_{2v}^{14}$) #38	a = 5.6896 b = 3.9692 c = 5.7256	180.00	2	4.623
KTaO ₃	Cubic Pm $\bar{3}m$ (O $_h^5$) #221	3.9885	268.04	1	7.015
KTiOPO ₄ , KTP	Orthorhombic Pna2 ₁ (C $_{2v}^9$) #33	a = 12.8164 b = 6.4033 c = 10.5897	197.94	8	3.026
LiB ₃ O ₅ , LBO	Orthorhombic Pna2 ₁ (C $_{2v}^9$) #33	a = 8.4473 b = 7.3788 c = 5.1395	119.37	4	2.475
LiCaAlF ₆ , LiCAF	Hexagonal P $\bar{3}1c$ (D $_{3d}^2$) #163	a = 5.008 c = 9.643	187.99	2	2.981
LiF	Cubic Fm $\bar{3}m$ (O $_h^5$) #225	4.0173	25.939	4	2.657
α -LiIO ₃	Hexagonal P6 $_3$ (C $_6^6$) #173	a = 5.4815 c = 5.1709	181.84	2	4.488
LiNbO ₃	Hexagonal R3c (C $_{3v}^6$) #161	a = 5.1483 c = 13.8631	147.85	6	4.629
LiYF ₄ , YLF	Tetragonal I4 $_1/a$ (C $_{4h}^6$) #88	a = 5.175 c = 10.74	171.84	4	3.968
MgAl ₂ O ₄ (spinel)	Cubic Fd $\bar{3}m$ (O $_h^7$) #227	8.084	142.27	8	3.577

(Continued)

TABLE 7 Composition, Structure, and Density of Crystals (*Continued*)

Material	Crystal System and Space Group	Unit Cell Dimension (\AA)	Molecular Weight (amu)	Formulae/ Unit Cell	Density (g/cm^3)
MgF ₂ (sellaite)	Tetragonal P4 ₂ /mnm (D _{4h} ¹⁴) #136	a = 4.623 c = 3.053	62.302	2	3.171
MgO (periclase)	Cubic Fm3m (O _h ⁵) #225	4.2117	40.304	4	3.583
NaCl (halite, rock salt)	Cubic Fm3m (O _h ⁵) #225	5.63978	58.44	4	2.164
[NH ₄] ₂ CO (urea, carbamide)	Tetragonal I4 ₂ m (D _{2d} ³) #113	a = 5.661 c = 4.712	60.055	2	1.321
NH ₄ H ₂ PO ₄ , ADP	Tetragonal I4 ₂ d (D _{2d} ¹²) #122	a = 7.4991 c = 7.5493	115.03	4	1.800
PbMoO ₄ (wulfenite)	Tetragonal I4 ₁ /a (C _{4h} ⁶) #88	a = 5.4312 c = 12.1065	367.16	4	6.829
PbTe (altaite)	Cubic Fm3m (O _h ⁵) #225	6.443	334.80	4	8.314
PbTiO ₃	Tetragonal P4 ₂ /mnm (D _{4h} ¹⁴) #136	a = 3.8966 c = 4.1440	303.07	1	7.998
RbTiOPO ₄ , RTP	Orthorhombic Pna2 ₁ (C _{2v} ⁹) #33	a = 12.948 b = 6.494 c = 10.551	244.31	8	3.654
Si	Cubic Fd3m (O _h ⁷) #227	5.43085	28.0855	8	2.329
SiO ₂ (α -quartz)	Hexagonal P3 ₂ 21 (D ₃ ⁶) #154	a = 4.9136 c = 5.4051	60.084	3	2.648
SrF ₂	Cubic Fm3m (O _h ⁵) #225	5.7996	125.62	4	4.277
SrMoO ₄	Tetragonal I4 ₁ /a (C _{4h} ⁶) #88	a = 5.380 c = 11.97	247.58	4	4.746
SrTiO ₃	Cubic Pm3m (O _h ⁵) #221	3.9049	183.49	1	5.117
TeO ₂ (paratellurite)	Tetragonal P4 ₁ 2 ₁ 2 (D ₄ ⁴) #92	a = 4.810 c = 7.613	159.60	4	6.019
TiO ₂ (rutile)	Tetragonal P4 ₂ /mnm (D _{4h} ¹⁴) #136	a = 4.5937 c = 2.9618	79.866	2	4.244
TlBr	Cubic Pm3m (O _h ⁵) #221	3.9846	284.29	1	7.462
Tl[0.46Br, 0.54I], KRS-5	Cubic Pm3m (O _h ⁵) #221	4.108	307.79	1	7.372
Tl[0.7Cl, 0.3Br], KRS-6	Cubic Pm3m (O _h ⁵) #221	3.907	253.17	1	7.049
Y ₃ Al ₅ O ₁₂ , YAG	Cubic Ia3d (O _h ¹⁰) #230	12.008	593.62	8	4.554
Y ₂ O ₃ (yttria)	Cubic Ia3 (T _h ⁷) #206	10.603	225.81	16	5.033
YVO ₄	Tetragonal I4 ₁ /amd (D _{4h} ¹⁹) #141	a = 7.1192 c = 6.2898	203.84	2	4.247

TABLE 7 Composition, Structure, and Density of Crystals (*Continued*)

Material	Crystal System and Space Group	Unit Cell Dimension (Å)	Molecular Weight (amu)	Formulae/ Unit Cell	Density (g/cm ³)
ZnGeP ₂	Tetragonal I4 ₂ d (D _{2d} ¹²) #122	a = 5.466 c = 10.722	199.97	4	4.146
β-ZnS (zincblende)	Cubic F4 ₃ m (T _d ²) #216	5.4094	97.445	4	4.089
α-ZnS (wurtzite)	Hexagonal P6 ₃ mc (C _{6v} ⁴) #186	a = 3.8218 c = 6.2587	97.445	2	4.088
ZnSe	Cubic F4 ₃ m (T _d ²) #216	5.6685	144.34	4	5.264
ZnTe	Cubic F4 ₃ m (T _d ²) #216	6.1034	192.98	4	5.638
ZrO ₂ :0.12Y ₂ O ₃ (cubic zirconia)	Cubic Fm3m (O _h ⁵) #225	5.148	121.98	4	5.939

TABLE 8 Optical Glass Reference Table [Cross-referenced glass examples, including 6-digit glass code and density (ρ, g/cm³)]

Common Name	Manufacturer			
	Schott	Ohara	Hoya	Pilkington
Fluor crown	FK	FSL	FC, FCD ¹	
	N-FK5 (487704) ρ = 2.45	S-FSL5 (487702) ρ = 2.46	FC5 (487704) ρ = 2.45	—
Phosphate crown	PK	FPL	PC, PCS ²	
	N-PK52A (497816) ρ = 3.70	S-FPL51 (497816) ρ = 3.62	FCD1 (497816) ρ = 3.70	—
Zinc crown	ZK		ZnC	ZC
	N-ZK7 (508612) ρ = 2.49	—	—	ZC508612A ρ = 2.50
Borosilicate crown	BK	BSL	BSC	BSC
	N-BK7 (517642) ρ = 2.51	S-BSL7(516614) ρ = 2.52	BSC7 (517642) ρ = 2.52	BSC517642B ρ = 2.50
Crown	K	NSL	C	HC ³
	N-K5 (522595) ρ = 2.59	S-NSL5 (522598) ρ = 2.49	—	HC522595A ρ = 2.54
Crown flint	KF	NSL	CF	
	N-KF9 (523515) ρ = 2.50	—	—	—
Very light flint	LLF	TIL	FEL ⁴	ELF ⁴
	LLF1 (548457) ρ = 2.94	S-TIL1 (548458) ρ = 2.49	E-FEL1 (548458) ρ = 2.54	—
Dense phosphate crown	PSK	BAL, PHM	PCD	
	N-PSK53A (618634) ρ = 3.57	S-PHM52 (618634) ρ = 3.67	PCD4 (618634) ρ = 3.52	—
Barium crown	BaK	BAL	BaC	MBC ⁵
	N-BaK4 (569560) ρ = 3.05	S-BAL14 (569563) ρ = 2.89	BAC4 (569560) ρ = 2.85	MBC569561A ρ = 3.06

(Continued)

TABLE 8 Optical Glass Reference Table [Cross-referenced glass examples, including 6-digit glass code and density (ρ , g/cm³)] (Continued)

Common Name	Manufacturer			
	Schott	Ohara	Hoya	Pilkington
Barium light flint	BaLF	BAL	BaFL ⁶	
	N-BaLF4 (580540)	S-BAL3 (571530)	BaFL3 (571530)	—
Light flint	$\rho = 3.11$	$\rho = 2.98$	$\rho = 2.98$	
	LF	TIL	FL	LF
	LF5 (581409)	S-TIL25 (581407)	E-FL5 (581409)	LF581409A
Fluoro flint	$\rho = 3.22$	$\rho = 2.59$	$\rho = 2.59$	$\rho = 3.23$
	TiF	TIM	FF	—
	TiF1 (511510)	S-FTM16 (593353)	*FF5 (593354)	—
Special short flint	$\rho = 2.479$ (obsolete)	$\rho = 2.64$	$\rho = 2.64$	
	KzFS	BAH, BAM	ADF ⁷	
	N-KzFS4 (613445)	S-NBM51 (613443)	E-ADF10 (613443)	—
Dense crown	$\rho = 3.00$	$\rho = 2.93$	$\rho = 3.04$	
	SK	BAL, BSM	BaCD ⁸	DBC ⁸
	N-SK10 (623570)	S-BSM10 (623570)	E-BACD10 (623569)	DBC623569A
Flint	$\rho = 3.67$ (obsolete)	$\rho = 3.60$	$\rho = 3.66$	$\rho = 3.66$
	F	TIM	F	DF ⁹
	N-F2 (620364)	S-TIM1 (626357)	E-F2 (620363)	DF620364A
Very dense crown	$\rho = 2.65$	$\rho = 2.71$	$\rho = 2.67$	$\rho = 3.61$
	SSK	BSM	BaCED ¹⁰	DBC
	N-SSK5 (658509)	S-BSM25 (658509)	BaCED5 (658509)	DBC658509A
Barium flint	$\rho = 3.71$	$\rho = 3.50$	$\rho = 3.64$	$\rho = 3.56$
	BaF	BAH, BAF	BaF	
	N-BaF10 (670471)	S-BAH10 (670473)	BAF10 (670472)	—
Barium dense flint	$\rho = 3.75$	$\rho = 3.48$	$\rho = 3.61$	
	BaSF	BAH	BaFD	
	N-BaSF64 (704394)	—	BAFD15 (702402)	—
Lanthanum crown	$\rho = 3.20$		$\rho = 2.99$	
	LaK	LAL, YGH	LaC, LaCL ¹¹	LAC
	N-LaK10 (720506)	S-LAL10 (720502)	LAC10 (720503)	LAC720504B
Tantalum crown	$\rho = 3.69$	$\rho = 3.86$	$\rho = 3.87$	$\rho = 3.86$
	LaK	LAL, YGH	TaC	
	N-LaK33A (754523)	S-YGH51 (755523)	TAC6 (755523)	—
Niobium flint	$\rho = 4.22$	$\rho = 4.40$	$\rho = 4.27$	
	LaF	LAM	NbF, NbFD	
	N-LaF35 (743494)	S-LAM60 (743493)	NbF1 (743492)	—
Light lanthanum flint	$\rho = 4.12$	$\rho = 4.06$	$\rho = 4.17$	
	N-LaF33 (786441)	S-LAH51 (786442)	NBFD11 (786439)	—
	$\rho = 4.36$	$\rho = 4.40$	$\rho = 4.43$	
Lanthanum flint	LaF	LAM	LaF	LaF
	N-LaF2 (744449)	S-LAM2 (744448)	LaF2 (744449)	LAF744447B
	$\rho = 4.30$	$\rho = 4.32$	$\rho = 4.39$	$\rho = 4.32$
Lanthanum flint	LaF	LAM	LaFL	LaF
	LaF2 (744447)	S-LAM59 (697485)	LAFL2 (697485)	—
	$\rho = 4.34$ (obsolete)	$\rho = 3.77$	$\rho = 4.05$	

TABLE 8 Optical Glass Reference Table [Cross-referenced glass examples, including 6-digit glass code and density (ρ , g/cm³)] (Continued)

Common Name	Manufacturer			
	Schott	Ohara	Hoya	Pilkington
Dense flint	SF	TIM, TIH	FD, FDS ¹²	EDF, DEDF, LDF
	SF2 (648338)	S-TIM22 (648338)	E-FD2 (648338)	EDF648339A
	$\rho = 3.86$	$\rho = 2.79$	$\rho = 2.77$	$\rho = 3.86$
	N-SF10 (728285)	S-TIH10 (728285)	E-FD10 (728283)	DEDF728284A
	$\rho = 3.05$	$\rho = 3.06$	$\rho = 3.07$	$\rho = 4.27$
	N-SF6 (805254)	S-TIH6 (805254)	FD60 (805255)	LDF805254A
Lanthanum dense flint	$\rho = 3.37$	$\rho = 3.37$	$\rho = 3.36$	$\rho = 3.37$
	LaSF	LAH	TaF ¹³ , TaFD ¹⁴	—
	N-LaSF44 (804465)	S-LAH65 (804466)	TaF3 (804465)	—
	$\rho = 4.44$	$\rho = 4.76$	$\rho = 4.65$	—
Antimony flint	N-LaSF31A (883408)	S-LAH58 (883408)	TaFD30 (883408)	—
	$\rho = 5.51$	$\rho = 5.52$	$\rho = 5.51$	—
	KzF	—	SbF	—
Fused silica (UV)	KzFN1 (551496)	—	SbF1 (551495)	—
	$\rho = 2.71$ (obsolete)	—	$\rho = 2.72$	—
	Lithosil-Q	—	—	—
	Lithosil-Q (458678)	—	—	—
	$\rho = 2.20$	—	—	—

¹FCD = dense fluor crown; ²PCS = special phosphate crown; ³HC = hard crown; ⁴FEL = ELF = extra light flint; ⁵MBC = medium barium crown; ⁶BaFL = light barium flint; ⁷ADF = abnormal dispersion flint; ⁸BACD = DBC = dense barium crown; ⁹DF = dense flint; ¹⁰BaCED = extra dense barium crown; ¹¹LaCL = light lanthanum crown; ¹²FDS = special dense flint; ¹³TaF = tantalum flint; ¹⁴TaFD = tantalum dense flint.

TABLE 9 Physical Properties of Specialty Glasses and Substrate Materials

Glass type	Density (g/cm ³)	Typical composition
Fused silica (SiO ₂) (e.g., Corning 7940 or Schott Lithosil-Q)	2.202	100%SiO ₂
Fused germania (GeO ₂)	3.604	100%GeO ₂
BS-39B (Barr and Stroud)	3.1	50%CaO, 34%Al ₂ O ₃ , 9%MgO
CORTRAN 9754 (Corning)	3.581	33%GeO ₂ , 20%CaO, 37%Al ₂ O ₃ , 5%BaO, 5%ZnO
IRG 2 (Schott)	5.00	Germanium glass
IRG 9 (Schott)	3.63	Fluorophosphate glass
IRG 11 (Schott)	3.12	Calcium aluminate glass
IRG 100 (Schott)	4.67	Chalcogenide glass
HTF-1 (Ohara) [443930]	3.94	Heavy metal fluoride glass
ZBLAN	4.52	56%ZrF ₄ , 14%BaF ₂ , 6%LaF ₃ , 4%AlF ₃ , 20%NaF
Arsenic trisulfide (As ₂ S ₃)	3.198	100%As ₂ S ₃
Arsenic triselenide (As ₂ Se ₃)	4.69	100%As ₂ Se ₃
AMTRI-1/TI-20	4.41	55%Se, 33%Ge, 12%As
AMTIR-3/TI-1173	4.70	60%Se, 28%Ge, 12%Sb
Pyrex (e.g., Corning 7740)	2.23	81%SiO ₂ , 13%B ₂ O ₃ , 4%Na ₂ O, 2%Al ₂ O ₃ [two-phase glass]
Zerodur (Schott)	2.53	56%SiO ₂ , 25%Al ₂ O ₃ , 8%P ₂ O ₅ , 4%Li ₂ O, 2%TiO ₂ 2%ZrO ₂ , ZnO/MgO/Na ₂ O/As ₂ O ₃ [glass ceramic]
ULE (Corning 7971)	2.205	92.5%SiO ₂ , 7.5%TiO ₂ [glass ceramic]

TABLE 10 Room-Temperature Elastic Constants of Cubic Crystals

Material	Stiffness (GPa)			Compliance (TPa ⁻¹)			Ref.
	c_{11}	c_{12}	c_{44}	s_{11}	s_{12}	s_{44}	
AgBr	56.3	32.8	7.25	31.1	-11.5	138	83, 100
AgCl	59.6	36.1	6.22	31.1	-11.7	161	83, 101
ALON	393	108	119	2.89	-0.62	8.40	102
BaF ₂	91.1	41.2	25.3	15.3	-4.8	39.5	83
Bi ₁₂ GeO ₂₀ (BGO)	125.0	32.4	24.9	8.96	-1.84	40.4	103
Bi ₁₂ SiO ₂₀ (BSO)	129.8	29.7	24.7	8.42	-1.57	40.2	104
C (diamond)	1077	124.7	557	0.95	-0.099	1.73	83
CaF ₂	165	46	33.9	6.94	-1.53	29.5	83
CdTe	53.5	36.9	20.2	42.6	-17.4	49.4	83, 105
CsI	24.5	6.6	6.31	46.1	-9.7	158	83
GaAs	118	53.5	59.4	11.75	-3.66	16.8	83
GaP	141	62.4	71.2	9.70	-2.97	14.0	83
Ge	129	48	67.1	9.73	-2.64	14.9	83
KBr	34.5	5.5	5.10	30.3	-4.2	196	83
KI	27.4	4.3	3.70	38.2	-5.2	270	83
KTaO ₃	431	103	109	2.7	-0.63	9.2	83
LiF	112	46	63.5	11.6	-3.35	15.8	83
MgAl ₂ O ₄	282.9	155.4	154.8	5.79	-2.05	6.46	106
MgO	297.8	95.1	155.8	3.97	-0.96	6.42	106
NaCl	49.1	12.8	12.8	22.9	-4.8	78.3	83
PbTe	107	8	13.2	9.46	-0.64	75.8	83, 107
Si	165	63	79.1	7.68	-2.12	12.6	83
SrF ₂	124	44	31.8	9.86	-2.57	31.5	83
SrTiO ₃	315.6	102.7	121.5	3.77	-0.93	8.23	108
TlBr	37.6	14.8	7.54	34.2	-9.6	133	83
Tl[Br:I], (KRS-5)	34.1	13.6	5.79	38.0	-10.8	173	83
Tl[Cl:Br], (KRS-6)	39.7	14.9	7.23	31.9	-8.8	139	83
Y ₃ Al ₅ O ₁₂ (YAG)	328.1	106.4	113.7	3.62	-0.89	8.80	109
Y ₂ O ₃	223.7	112.4	74.6	6.73	-2.25	13.4	110
ZnS	102	64.6	44.6	19.5	-7.6	22.5	83
ZnSe	86.4	51.5	40.2	21.0	-7.9	24.9	83
ZnTe	71.5	40.8	31.1	23.9	-8.5	32.5	83
ZrO ₂ :Y ₂ O ₃	405.1	105.3	61.8	2.77	-0.57	16.18	111

TABLE 11 Room-Temperature Elastic Constants of Tetragonal Crystals (Point groups 4mm, $\bar{4}2m$, 422, and 4/mmm)

Material	c or s	Subscript of stiffness (GPa) or compliance (TPa ⁻¹)						Ref.
		11	12	13	33	44	66	
AgGaS ₂	c	87.9	58.4	59.2	75.8	24.1	30.8	83, 112
	s	26.2	-7.7	-14.5	35.9	41.5	32.5	
AgGaSe ₂	c	80.1	51.6	52.6	70.7	21.2	24.7	113
	s	26.9	-8.2	-13.9	34.9	47.2	40.5	
BaTiO ₃	c ^E	275	179	152	165	54.4	113	83
	s ^E	8.05	-2.35	-5.24	15.7	18.4	8.84	
	c ^E	211	107	114	160	56.2	127	114
	s ^E	8.01	-1.57	4.60	12.8	17.8	7.91	

TABLE 11 Room-Temperature Elastic Constants of Tetragonal Crystals (Point groups $4mm$, $\bar{4}2m$, 422 , and $4/mmm$) (Continued)

Material	c or s	Subscript of stiffness (GPa) or compliance (TPa ⁻¹)						Refs.
		11	12	13	33	44	66	
CdGeAs ₂	c	94.5	59.6	59.7	83.4	42.1	40.8	83, 115
	s	21.6	-7.04	-10.4	26.9	23.8	24.5	
CsLiB ₆ O ₁₀	c	62.79	-16.74	23.86	37.4	33.31	26.7	116
	s	38.33	25.75	-40.88	78.90	44.82	37.45	
CuGaS ₂	c							
	s							
KH ₂ PO ₄ (KDP)	c	71.2	-5.0	14.1	56.8	12.6	6.22	83
	s	14.9	1.9	-4.2	19.7	78.4	161	
MgF ₂	c	138	88	62	201	56.5	95.6	83
	s	12.6	-7.2	-1.65	6.01	17.7	10.5	
[NH ₄] ₂ CO (Urea)	c	21.7	8.9	24	53.2	6.26	0.45	117
	s	95	16	-50	64	160	2220	
NH ₄ H ₂ PO ₄ (ADP)	c	67.3	5.0	19.8	33.7	8.57	6.02	83
	s	18.3	2.2	-12.0	43.7	117	166	
PbTiO ₃	c ^E	235	101	98.8	105	65.1	104	118
	s ^E	7.06	-0.40	-6.27	21.3	15.5	9.62	
TeO ₂	c	56.12	51.55	23.03	105.71	26.68	66.14	119, 120
	s	114.5	-104.3	-2.3	10.5	37.5	15.1	
TiO ₂	c	269	177	146	480	124	192	83
	s	6.80	-4.01	-0.85	2.60	8.06	5.21	
YVO ₄	c	244.51	48.93	81.09	313.7	24.18	16.18	121
	s	4.54	-0.57	-1.03	3.72	20.76	61.6	
ZnGeP ₂	c	(87)	(66)	(64)	(81)	(20)	(23)	122
	s	(33.5)	(-14.2)	(-15.2)	(36.4)	(50)	(43.5)	

TABLE 12 Room-Temperature Elastic Constants of Tetragonal Crystals (Point groups 4 , $\bar{4}$, and $4/m$)

Material	c or s	Subscript of stiffness (GPa) or compliance (TPa ⁻¹)							Ref.
		11	12	13	16	33	44	66	
CaMoO ₄	c	144	65	47	-13.5	127	36.8	45.8	83, 123
	s	9.90	-4.2	2.1	4.4	9.48	27.1	24.4	
CaWO ₄	c	141	61	41	-17	125	33.7	40.7	83, 124
	s	10.5	-4.7	-1.9	6.4	9.3	29.7	30.2	
PbMoO ₄	c	107.2	61.9	52.0	-15.8	93.2	26.4	34.8	83, 125
	s	20.8	-11.8	-5.0	14.9	16.3	37.9	42.3	
SrMoO ₄	c	117.3	58.7	46.8	-10.8	103.8	34.9	46.6	83, 126
	s	13.2	-5.7	-3.3	4.5	12.6	28.7	23.7	
YLiF ₄	c	121	60.9	52.6	-7.7	156	40.9	17.7	83, 127
	s	12.8	-6.0	-2.3	8.16	7.96	24.4	63.6	

TABLE 13 Room-Temperature Elastic Constants of Hexagonal Crystals (Point groups 6, $\bar{6}$, 6/m, 622, 6mm, $\bar{6}2m$ and 6/mmm)

Material	c or s	Subscript of stiffness (GPa) or compliance (TPa ⁻¹)					Ref.
		11	12	13	33	44	
BeO	c	470	168	119	494	153	83
	s	2.52	-0.80	-0.41	2.22	6.53	
CdS	c	88.4	55.4	48.0	95.2	15.0	83
	s	20.5	-9.9	-5.3	15.9	66.7	
CdSe	c	74.1	45.2	38.9	84.3	13.4	83
	s	23.2	-11.2	-5.5	16.9	74.7	
GaN	c	296	130	158	267	241	128
	s	5.10	-0.92	-2.48	6.68	4.15	
LiIO ₃	c ^E	81.24	31.84	9.25	52.9	17.83	129
	s ^E	14.7	-5.6	-1.6	19.5	56.1	
ZnS	c	122	58	43	138	28.7	83
	s	11.0	-4.5	-2.1	8.6	34.8	

TABLE 14 Room-Temperature Elastic Constants of Hexagonal (Trigonal) Crystals (Point groups 32, 3m, $\bar{3}m$)

Material	c or s	Subscript of stiffness (GPa) or compliance (TPa ⁻¹)						Ref.
		11	12	13	14	33	44	
Ag ₃ AsS ₃	c	59.5	31.7	29.6	0.18	39.8	9.97	83, 130
	s	28.6	-7.3	-15.9	-0.6	48.8	100	
Al ₂ O ₃	c	496	159	114	-23	499	146	83
	s	2.35	-0.69	-0.38	0.47	2.18	7.0	
β -Ba ₃ B ₆ O ₁₂ (BBO)	c	123.8	60.3	49.4	12.3	53.3	7.8	131
	s	25.63	-14.85	-9.97	-63.97	37.21	331.3	
CaCO ₃ (Calcite)	c	144	54.2	51.2	-20.5	84.3	33.5	83
	s	11.4	-4.0	-4.5	9.5	17.3	41.4	
LiCaAlF ₆ (LiCAF)	c	118	42	54	± 19	107	50	132
	s	12.9	-3.4	-4.8	± 6.2	-14.2	24.7	
LiNbO ₃	c	202	55	72	8.5	244	60.2	83
	s	5.81	-1.12	-1.38	-0.98	4.93	16.9	
α -SiO ₂	c	86.6	6.74	12.4	-17.8	106.4	58.0	83
	s	12.8	-1.75	-1.30	4.47	9.73	20.0	

TABLE 15 Room-Temperature Elastic Constants of Orthorhombic Crystals

Material	c or s	Subscript of stiffness (GPa) or compliance (TPa ⁻¹)									Ref.
		11	12	13	22	23	33	44	55	66	
KNbO ₃	c ^E	224	102	182	273	130	245	75	28.5	95	133
	s ^E	11.3	-0.3	-8.2	4.9	-2.4	11.5	13.3	35.1	10.5	
KTiOPO ₄ (KTP)	c	166	37	54	164	51	181	56	54	45	134
	s	6.8	-1	-1.8	6.8	-1.6	6.5	17.9	18.5	22.2	
LiB ₃ O ₅ (LBO)	c ^E	127.1	126.6	52.0	237.7	57.6	65.5	109.3	86.7	17.8	135
	s ^E	19.64	-8.49	-8.13	9.01	-1.18	22.76	9.15	11.53	56.25	
RbTiOPO ₄	c	163	45	35	165	63	178	58	57	50	136
RTP	s	6.73	-1.54	-0.78	7.36	-2.30	6.59	17.24	17.54	20.0	

TABLE 16 Room-Temperature Elastic Constants for Monoclinic Crystals

Material	c or s	Subscript of stiffness (GPa) or compliance (TPa ⁻¹)							Ref.
		11	12	13	15	22	23	25	
BiB ₃ O ₆	c	159.7	74.2	60.0	-49.7	52.5	13.4	-4.3	137
	s	34.3	-46.8	-0.03	20.1	83.2	-1.27	-27.56	
Material	c or s	33	35	44	46	55	66	—	Ref.
BiB ₃ O ₆	c	205.2	-70.8	23.3	-18.6	74.6	66.9	—	137
	s	7.33	6.87	55.2	15.3	31.7	19.2	—	

TABLE 17 Mechanical Properties of Crystals

Material	Moduli (GPa)				Poisson's Ratio	Flexure Strength (MPa)	Knoop Hardness (kg/mm ²)
	Elastic	Shear	Bulk				
Ag ₃ AsS ₃	30	11	36.8	0.36			
AgBr	24.7	8.8	40.5	0.39 ₉		7.0	
AgCl	22.9	8.1	44.0	0.41	26	9.5	
AgGaS ₂	53	19	67	0.37		320	
AgGaSe ₂	47.8	17.5	60.4	0.37		230	
Al ₂ O ₃	402	163	252	0.23	1200	2250	
ALON	317	128	203	0.24	310	1850	
Ba ₃ B ₆ O ₁₂ , BBO	30	11	60.6	0.41			
BaF ₂	65.9	25.2	57.8	0.31	27	78	
BaTiO ₃	145	53	174	0.36		580	
BeO	395	162	240	0.23	275	1250	
Bi ₁₂ GeO ₂₀ , BGO	82	32	63.3	0.28			
Bi ₁₂ SiO ₂₀ , BSO	84	33	63.1	0.28			
BiB ₃ O ₆ , BIBO	91	37.5	52.5	0.21			
C, diamond	1142	534	442	0.069	2940	9000	
CaCO ₃ , calcite	83	32	73.2	0.31		100	
CaF ₂	110	42.5	85.7	0.29	90	170	
CaMoO ₄	103	40	80	0.29		250	
CaWO ₄	96	37	77	0.29		300	
CdGeAs ₂	74	28	70	0.32		470	
CdS	47.3	17.2	64.0	0.38	28	122	
CdSe	42	15.3	53	0.37	21	65	
CdTe	38.4	14.2	42.4	0.35	26	50	
CsLiB ₆ O ₁₀ , CLBO	70.3	27.5	52.5	0.28			
CsI	18	7.3	12.6	0.26	5.6	12	
CuGaS ₂			95.8			430	
GaAs	116	46.6	75.0	0.24	55	710	
GaN	294	118	195	0.25	70	750	
GaP	139	56.2	88.6	0.24	100	875	
Ge	132	54.8	75.0	0.20 ₆	100	850	

(Continued)

TABLE 17 Mechanical Properties of Crystals (Continued)

Material	Moduli (GPa)			Poisson's Ratio	Flexure Strength (MPa)	Knoop Hardness (kg/mm ²)
	Elastic	Shear	Bulk			
KBr	18	7.2	15.2	0.30	11	6.5
KH ₂ PO ₄ , KDP	38	15.1	27	0.27		
KI	14	5.5	11.9	0.30		5
KNbO ₃	144	53	173	0.36		500
KTaO ₃	316	124	230	0.27		
KTiOPO ₄ , KTP	137	55	88	0.24		700
LiB ₃ O ₅ , LBO	104	41	80	0.28		600
LiCaAlF ₆ , LiCAF	93	36	71	0.28		
LiF	110	45	65.0	0.22 ₅	27	115
LiIO ₃	55	22.4	33.5	0.23		
LiNbO ₃	170	68	112	0.25		630
LiYF ₄ , YLF	85	32	81	0.32	35	300
MgAlO ₄	276	109	198	0.26 ₈	170	1650
MgF ₂	137	54	99.7	0.27	100	500
MgO	310	131	163	0.18	130	675
NaCl	37	14.5	25.3	0.26	9.6	16.5
[NH ₄] ₂ CO, urea	~9	~3	17	0.41		
NH ₄ H ₂ PO ₄ , ADP	29	11	27.9	0.32 ₅		
PbMoO ₄	66	24.6	71	0.34		
PbTe	57.3	22.6	41.9	0.27		
PbTiO ₃	144	56	117	0.20	98	
RbTiOPO ₄ , RTP	140	56.8	87.7	0.23		(500)
Si	165	66.4	97.0	0.221	130	1150
SiO ₂ , α -quartz	95	44	38	0.08		740
SrF ₂	89.2	34.6	70.7	0.29		150
SrMoO ₄	89	35	71	0.29		
SrTiO ₃	283	115	174	0.23		600
TeO ₂	45	17	46	0.33		
TiO ₂	293	115	215	0.27		880
TlBr	24	8.9	22.4	0.32		12
Tl[Br, I], KRS-5	19.6	7.3	20.4	0.34	26	40
Tl[0.7Cl, 0.3Br], KRS-6	24	9.0	32.2	0.33	21	30
Y ₃ Al ₅ O ₁₂ , YAG	280	113	180	0.24	150	1350
Y ₂ O ₃	173	66.4	149.5	0.31	150	700
YVO ₄	133	50	134	0.33		
ZnGeP ₂	[44]	[16]	86	[0.39]		980
β -ZnS	83.5	31.6	77.1	0.32	60	175
α -ZnS	87	33	74	0.30	69	
ZnSe	75.2	28.9	63.1	0.30	55	115
ZnTe	61.1	23.5	51.0	0.30	24	82
ZrO ₂ : 12%Y ₂ O ₃	23.3	88.6	205	0.31	(200)	1150

TABLE 18 Mechanical Properties of Optical and Specialty Glasses and Substrate Materials

Selected Glass Code or Designation	Moduli (GPa)			Poisson's Ratio	Flexure Strength (MPa)	Knoop Hardness (kg/mm ²)
	Elastic	Shear	Bulk			
487704 N-FK5	62	25	39	0.232		520
497816 N-PK52A	71	27	59	0.298		355
508612 N-ZK7	70	29	41	0.214		530
517642 N-BK7	82	34	46	0.206		610
522595 N-K5	71	29	43	0.224		530
523515 N-KF9	66	27	40	0.225		480
548458 LLF1	60	25	34	0.208		450
618634 N-PSK53A	76	30	60	0.288		415
569560 N-BaK4	77	31	49	0.24		550
580537 N-BaLF4	77	31	49	0.245		540
581409 LF5	59	24	36	0.223		450
593355 FF5	[65]	[26]	[41]	[0.238]		500
613443 N-KzFS4	78	31	50	0.241		520
623570 N-SK10	82	32	60	0.273		570
620364 N-F2	82	33	50	0.228		600
658509 N-SSK5	88	34	66	0.278		590
670472 N-BaF10	89	35	65	0.271		620
704394 N-BaSF64	105	42	74	0.264		650
720504 N-LaK10	116	45	90	0.286		780
754523 N-LaF33	111	43	93	0.301		740
743492 NbF1	109	42	95	0.31		790
744447 N-LaF2	94	36	74	0.288		530
805254 N-SF6	93	37	65	0.262		550
883409 N-LaSF31A	126	48	105	0.301		650
Fused silica	72.6	31	36	0.164	110	635
Fused germania	43.1	18	23	0.192		
BS-39B	104	40	83	0.29	90	760
CORTRAN 9754	84.1	33	67	0.290	44	560
IRG 2	95.9	37	73	0.282		481
IRG 9	77.0	30	61	0.288		346
IRG 11	107.5	42	83	0.284		610
IRG 100	21	8	15	0.261		150
HTF-1	64.2	25	49	0.28		320
ZBLAN	60	23	53	0.31		225
Arsenic trisulfide	15.8	6	13	0.295	16.5	180
Arsenic triselenide	18.3	7	14	0.288	16.2	120
AMTIR-1/TI-20	21.9	9	16	0.266	18.6	170
AMTIR-3/TI-1173	21.7	9	15	0.265	17.2	150
Pyrex	62.8	26	35	0.200		
Zerodur	91	37	58	0.24		630
ULE	67.3	29	34	0.17	50	460

TABLE 19 Thermal Properties of Crystals

Material	CC*	Temperature (K)		Heat Capacity (J/g · K)	Thermal Expansion (10 ⁻⁶ /K)	Thermal conductivity (W/m · K)		
		Debye	Melt [†]			@ 250 K	@ 300 K	@ 500 K
Ag ₃ AsS ₃	H		769 m		16 a 12 c		0.092 a 0.110 c	
AgBr	C	145	705 m	0.279	33.8	1.11	0.93	0.57
AgCl	C	162	728 m	0.3544	32.4	1.25	1.19	
AgGaS ₂	T	193	1269 m	0.404	12.7 a -13.2 c		1.5 a 1.4 c	
AgGaSe ₂	T	156	1269 m	0.297	19.8 a -8.1 c		1.1 a 1.0 c	
Al ₂ O ₃	H	1030	2319	0.777	4.15 c 6.65 a 7.15 c	58	46	24.2
ALON	C		2323 m	0.830	5.66		12.6	7.0
Ba ₃ B ₆ O ₁₂ , BBO	H	112	900 p 1368 m	0.490	4.0 a 36.0 c		1.2 a 1.6 c	
BaF ₂	C	283	1553 m	0.4474	18.4	11	7.5	
BaTiO ₃	T	345	267 p 406 p 1898 m	0.439	16.8 a -9.07 c	—	1.3	—
BeO	H	1280	2373 p 2725 m	1.028	5.64 a 7.47 c	420	350	200
Bi ₁₂ GeO ₂₀ (BGO)	C	580	1203	0.242	16.8			
Bi ₁₂ SiO ₂₀ (BSO)	C	350	1274	0.180	15.0			
BiB ₃ O ₆ , BIBO	M		999 m	0.500	-25.6 a 50.4 b 7.7 c			
C, diamond	C	2240	1770 p	0.5169	1.25	2800	2200	1300
CaCO ₃ , calcite	H		323 p 3825 m	0.8820	-3.7 a 25.1 c	5.1 a 6.2 c	4.5 a 5.4 c	(3.4) a (4.2) c
CaF ₂	C	510	1424 p	0.9113	18.9	13	9.7	5.5
CaMoO ₄	T	300	1730 m	0.573	7.6 a 11.8 c		4.0 a 3.8 c	

CaWO ₄	T	245	1855 m	0.396	6.35 a 12.38 c	16	9.5
CdGeAs ₂	T	253	900 p 943 m		8.4 a 0.25 c		
	H	215	1560 m	0.3814	4.6 a 2.5 c	14 a 16 c	
CdSe	H	181	1580 m	0.258	4.40 a 2.45 c	6.2 a 6.9 c	
	C	160	1320 m	0.210	5.0	6.3	
CsLiB ₆ O ₁₀ , CLBO	T		1118 m	0.93	21 a -17 c	2.0 a 2.2 c	
	C	124	898 m	0.2032	48.6	1.05	
CuGaS ₂	T	356	1553 m	0.452	11.2 a 6.9 c	17.9	
	C	344	1511 m	0.318	5.8	54	27
GaN	H	365	1160 d		3.17 a		
	C	460	1740 m	0.435	5.59 c	130 c	
GaP	C	380	1211 m	0.3230	5.3	100	(45)
	C	174	1007 m	0.4400	5.7	59.9	33.8
KH ₂ PO ₄ , KDP	T		123 p 450 p 526 m	0.879	38.5 26.8 a 42.4 c	4.8	2.4
	C	132	954 m	0.3192	40.3	1.3 a 1.2 c	1.7 a 1.3 c
KNbO ₃	O	—	223 p 476 p 1333 m	0.767	(37)	~4	
	C	311	1148 m 1423 d	0.366	5.3	0.17	
KTOPO ₄ , KTP	O	365	1107 p	0.727	11 a 9 b 0 c	2.0 a 3.0 b 3.3 c	
	O	1115 (@300K)	1083 m	0.935	107.1 a -95.4 b 33.7 c	2.7 a 3.1 b 4.5 c	
LiCaAlF ₆ , LiCAF	H		1083 m		22 a 3.6 c	4.58 a 5.14 c	

(Continued)

TABLE 19 Thermal Properties of Crystals (Continued)

Material	CC*	Temperature (K)		Heat Capacity (J/g · K)	Thermal Expansion (10 ⁻⁶ /K)	Thermal conductivity (W/m · K)		
		Debye	Melt†			@ 250 K	@ 300 K	@ 500 K
LiNbO ₃	H	560	1470 c 1530 m	0.648	15.4 a 5.3 c		5.4 a 5.3 c	
LiYF ₄ YLF	T	460	1092 m	0.79	14.3 a 10.1 c	6.3 a 8.8 c	5.3 a 7.2 c	
MgAlO ₄	C	850	2408 m	0.8191	6.97	30	25	
MgF ₂	T	535	1536 m	1.0236	9.4 a 13.6 c		30 a 21 c	
MgO	C	950	3073 m	0.9235	10.6	73	59	32
NaCl	C	321	1074 m	0.8699	41.1	8	6.5	4
[NH ₄] ₂ CO ₂ urea	T	135	408 m	1.553	52.7 a 11.4 c			
NH ₄ H ₂ PO ₄ ADP	T		148 p 463 m	1.236	32 a 4.2 c		1.26 a 0.71 c	
PbMoO ₄	T	190	1338 m	0.326	8.7 a 20.3 c			
PbTe	C	175	1190 m	0.151	19.8	2.5	2.3	1.8
PbTiO ₃	T	337	763 p 1563 m	0.462			4	2.8
RbTiOPO ₄ RTP	O		1213 m 1374 d		10.1 a 13.7 b -4.2 c			
Si	C	645	1680 m	0.7139	2.62	191	140	73.6
SiO ₂ α-quartz	H	271	845 p	0.7400	12.38 a 6.88 c	7.5 a 12.7 c	6.2 a 10.4 c	3.9 a 6.0 c
SrF ₂	C	378	1710 m	0.6200	18.1	11	8.3	
SrMoO ₄	T	260	1763 m	0.619	9.6 a 21.7 c		4.0 a 4.2 c	
SrTiO ₃	C	—	110 p 2358 m	0.536	8.3	12.5	11.2	
TeO ₂	T	~100	1006 m	[0.41]	15.0 a 4.9 c	3.2 a 1.9 c	2.6 a 1.7 c	

TiO ₂	T	760	2128 m	0.6910	6.86 a 8.97 c	8.3 a 11.8 c	7.4 a 10.4 c	(5.5) a (8.0) c
TlBr	C	116	740 m	0.1778	51		0.53	
Tl[Br, I] KRS-5	C	(110)	687 m	(0.16)	58		0.32	
Tl[Cl, Br], KRS-6	C	(120)	697 m	0.201	51		0.50	
Y ₃ Al ₅ O ₁₂ , YAG	C	754	2193 p	0.625	7.7		13.4	
Y ₂ O ₃	C	465	2640 p	0.4567	6.56		13.5	
YVO ₄	T	443	2083	0.56	2.2 a 8.4 c		8.9 a 12.1 c	
ZnGeP ₂	T	428	1225 p 1300 m		7.8 a 5.0 c	60	35 a 36 c	16
β-ZnS	C	340	1293 p	0.4732	6.8		16.7	10
α-ZnS	H	351	2100 m	0.4723	6.54 a 4.59 c			
ZnSe	C	270	1790 m	0.339	7.1		13	8
ZnTe	C	225	1510 m	0.218	8.4		10	
ZrO ₂ :12%Y ₂ O ₃	C	563	3110 m	0.46	8.6		1.8	1.9

*CC = crystal class; C = cubic; H = hexagonal; M = monoclinic; O = orthorhombic; T = tetragonal.

[†]Temperature codes: m = melt temperature; c = Curie temperature; d = decomposition temperature; p = phase change (to different structure) temperature; v = vaporization (sublimation) temperature.

TABLE 20 Thermal Properties of Optical and Specialty Glasses and Substrate Materials

Selected Glass Code	Temperature (K)			Heat Capacity (J/g · K)	Thermal Expansion (10 ⁻⁶ /K)	Thermal Conductivity (W/m · K)
	Glass	Soften	Melt*			
487704 N-FK5	739	945		0.808	9.2	0.925
497816 N-PK52A	740	811		0.67	13.01	0.73
508612 N-ZK7	812	994		0.77	4.5	1.042
517642 N-BK7	830	992		0.858	7.1	1.114
522595 N-K5	819	993		0.783	8.2	0.95
523515 N-KF9	749	913		[0.75]	9.61	1.04
548458 LLF1	704	901		0.65	8.1	[0.960]
618634 N-PSK53A	879	972		0.590	9.56	0.64
569560 N-BaK4	854	998		0.680	6.99	0.88
580537 N-BaLF4	851	934		0.69	6.52	0.827
581409 LF5	692	858		0.657	9.1	0.866
593355 FF5	788	843		[0.80]	8.6	[0.937]
613445 N-KzFS4	820	948		0.76	7.3	0.84
623570 N-SK10	905	1017		0.55	7.0	0.88
620364 N-F2	842	959		0.81	7.84	1.05
658509 N-SSK5	918	1024		0.574	6.8	[0.806]
670472 N-BaF10	933	1063		0.56	6.18	0.78
704394 N-BaSF64	855	985		[0.70]	7.3	[0.90]
720504 N-LaK10	909	987		0.64	5.68	0.86
754523 LaF33	873	946		0.57	5.6	0.8
743492 NbF1	863	898		[0.48]	5.3	[0.845]
744447 N-LaF2	926	1015		0.51	8.06	0.67
805254 N-SF6	862	956		0.69	9.03	0.96
883409 N-LaSF31A	992	1103		0.44	6.74	0.79
Fused silica	1273		1983	0.746	0.51	1.38
Fused germania	800		1388		6.3	
BS-39B		[970]		0.865	8.0	1.23
CORTRAN 9754	1008	1147		0.54	6.2	0.81
IRG 2	975			0.495	8.8	0.91
IRG 9	696			0.695	16.1	0.88
IRG 11	1075			0.749	8.2	1.13
IRG 100	550	624			15.0	0.3
HTF-1	658				16.1	
ZBLAN	543		745	0.520	17.5	0.4
Arsenic trisulfide	460	573		0.473	26.1	0.17
Arsenic triselenide	375	345		0.349	24.6	0.205
AMTIR-1/TI-20	635	678		0.293	12.0	0.25
AMTIR-3/TI-1173	550	570		0.276	14.0	0.22
Zerodur				0.821	0.5	1.46
					(20–300°C)	
Pyrex	560	821		1.05	3.25	1.13
ULE	1000	1490		0.776	±0.03	1.31
					(5–35°C)	

*Or liquidus temperature.

TABLE 21 Summary Optical Properties of Crystals

Material	Transparency (μm)		Refractive Index (n_o)	Thermo-optic coefficient ($10^{-6}/\text{K}$)					Ref.	
	UV	IR		$\lambda(\mu\text{m})$	dn/dT	$\lambda(\mu\text{m})$	dn/dT	$\lambda(\mu\text{m})$		dn/dT
Ag ₃ AsS ₃	0.63(o)	12.5(o)	2.736(o)	150					130	
	0.61(e)	13.3(e)	2.519(e)							
AgBr	0.49	35	2.116			3.39	-61	10.6	-50	139
AgCl	0.42	23	2.002	-61		3.39	-58	10.6	-35	139, 140
AgGaS ₂	0.50(o)	11.4(o)	2.408(o)	258(o)		1.0	176(o)	10.0	153(o)	141
	0.52(e)	12.0(e)	2.354(e)	255(e)			179(e)		155(e)	142
AgGaSe ₂	0.75	17.0	2.617(o)			1.0	75.6(o)	10.0	79.7(o)	143
			2.584(e)				80.5(o)		85.9(e)	
Al ₂ O ₃	0.19(o)	5.0(o)	1.7555(o)	11.7(o)		3.39	11.3(o)	5.0	14.1(o)	144, 145,
			1.7478(e)	12.8(e)		0.5790	-16.4(o)	1.014	-16.8(o)	146, 147
ALON	0.23	4.8	1.771	11.7						148
BBO	0.19	2.6	1.656(o)	-16.6(o)						131
			1.501(e)	-9.8(e)						
BaF ₂	0.14	12.2	1.4663	-16.0		3.39	-15.9	10.6	-14.5	145
				-15.2		3.30	-15.6			149
BaTiO ₃			2.277(o)							
			2.250(e)							
BeO	0.21	3.5	1.710(o)	8.2(o)		0.633	8.2(o)			150
			1.723(e)	13.4(e)			13.4(e)			
BGO	0.50	3.1	2.367	-34.5		0.65	-34.9			151
BSO	0.52		2.397							
BiB ₃ O ₆ , BIBO	0.27	2.7	1.433(x)	69.6(x)		1.0	53.8(x)	3.0	45.5(x)	152
			1.450(y)	63.9(y)			48.1(y)		43.0(y)	
			1.509(z)	85.6(z)			74.0(z)		61.1(z)	
Diamond	0.24	2.7	2.380	10.1				30	9.6	153, 154
	0.24(o)	2.2(o)	1.642(o)	3.6(o)		0.458	3.2(o)	0.633	2.1(o)	2
CaCO ₃	0.21(e)	3.3(e)	1.478(e)	14.4(e)			13.1(e)		11.9(e)	
CaF ₂	0.135	9.4	1.4278	-7.5		0.663	-10.4	3.39	-8.1	149
			0.365	-10.6				3.30	-10.8	155
CaMoO ₄			1.945(o)	-9.6(o)						156
			1.951(e)	-10.0(e)						

(Continued)

TABLE 21 Summary Optical Properties of Crystals (*Continued*)

Material	Transparency (μm)		IR	Refractive Index (n_∞)	Thermo-optic coefficient ($10^{-6}/\text{K}$)					Ref.	
	UV	IR			$\lambda(\mu\text{m})$	dn/dT	$\lambda(\mu\text{m})$	dn/dT	$\lambda(\mu\text{m})$		dn/dT
CaWO_4	(0.2)	5.3		1.884(o) 1.898(e)	0.546	-7.1(o) -10.2(e)				157	
CdGeAs_2	2.5	15		3.522(o) 3.608(e)							
CdS	0.52(o) 0.51(e)	14.8(o) 14.8(e)		2.276(o) 2.293(e)	10.6	58.6(o) 62.4(e)				158	
CdSe	0.75	20		2.448(o) 2.467(e)							
CdTe	0.85	29.9		2.6829	1.15	147	3.39	98.2	10.6	98.0	159
CsI	0.245	62		1.743	0.365	-87.5	0.633	-99.3	30.0	-88.0	160
$\text{CsLiB}_6\text{O}_{10}$	0.18	2.75		1.487(o) 1.435(e)	0.532	-1.9(o) -0.5(e)					161, 162
CuGaS_2				2.493(o) 2.487(e)	0.55	130(o) 173(e)	1.0	59(o) 60(e)	10.0	56(o) 57(e)	163
GaAs	0.90	17.3		3.32	1.15	250	3.39	200	10.6	200	165
GaN				2.31(o) 2.31(e)	1.15	61					
GaP	0.54	10.5		3.014	0.546	200	0.633	160			166
Ge	1.8	15		4.0017	2.5	462	5.0	416	20.0	401	167
$\text{HfO}_2\text{:Y}_2\text{O}_3$	0.35	6.5		2.074	0.365	14.1	0.436	11.0	1.01	5.8	168
KBr	0.200	30.2		1.537	0.458	-39.3	1.15	-41.9	10.6	-41.1	145
KH_2PO_4	0.176	1.42		1.503(o) 1.460(e)	0.624	-39.6(o) -38.2(e)					169
KI	0.250	38.5		1.629	0.458	-41.5	1.15	-44.7	30	-30.8	145
KNbO_3	0.4	5.0		2.103(x) 2.199(y) 2.233(z)	0.436	67(x) -26(y) 125(z)	1.064	23(x) -34(y) 63(z)	3.00	21(x) -23(y) 55(z)	170
KLaO_3				2.14							
KTiOPO_4	0.35	4.5		1.814(x) 1.857(y) 2.143(z)	0.5	10.8(x) 15.5(y) 28.3(z)	1.5	5.3(x) 7.6(y) 12.3(z)			171, 172
LiB_3O_5	0.16	3.6		1.568(x) 1.589(y) 1.608(z)	0.532	-0.9 -13.5 -7.4	1.064	-1.9 -13.0 -8.3			173

LiCaAlF ₆	0.11	1.388(o)	0.546	-4.2(o)						174
LiCAF		1.386(e)		-4.6(e)						
LiF	0.120	1.388	0.458	-16.0	1.15	-16.9	3.39	-14.5		145
			0.365	-15.0	1.53	-15.6				149
LiIO ₃	0.38	1.846(o)	0.4	-74.5(o)	1.0	-84.9(o)				80
		1.708(e)		-63.5(e)		-69.2(e)				
LiNbO ₃	0.35	2.282(o)	0.66	4.4(o)	3.39	0.3(o)				175
		2.006(e)		37.9(e)		28.9(e)				
LiYF ₄	0.18	1.447(o)	0.436	-0.54(o)	0.546	-0.67(o)	0.578	-0.91(o)		176
		1.469(e)		-2.44(e)		-2.30(e)		-2.86(e)		
MgAl ₂ O ₄	0.2	1.701	0.589	9.0						144
MgF ₂	0.13(o)	1.3734(o)	0.633	1.12(o)	1.15	0.88(o)	3.39	1.19(o)		145
	0.13(e)	1.3851(e)		0.58(e)		0.32(e)		0.6(e)		
MgO	0.35	1.720	0.365	19.5	0.546	16.5	0.768	13.6		177, 178
NaCl	0.174	1.526	0.458	-34.2	0.633	-35.4	3.39	-36.3		145
[NH ₄] ₂ CO	0.21	1.477(o)								
		1.586(e)								
NH ₄ H ₂ PO ₄	0.185	1.518(o)	0.624	-47.1(o)						169
		1.471(e)		-4.3(e)						
PbMoO ₄	0.5	2.265(o)	0.588	-75(o)						156
		2.175(e)		-41(e)						
PbTe	4.0	5.57	3.39	-2100	5.0	-1500	10.6	-1200		179
PbTiO ₃		2.52(o)								
		2.52(e)								
RbTiOPO ₄ RTP	0.34	1.763(x)	0.633	56(x)						180
		1.768(y)		91(y)						
		1.8478(z)		66(z)						
Si	1.1	3.4159	2.5	166	5.0	159	10.6	157		167
α-SiO ₂	0.155	1.5352(o)	0.254	-2.9(o)	0.365	-5.4(o)	0.546	-6.2(o)		178
		1.5440(e)		-4.0(e)		-6.2(e)		-7.0(e)		
SrF ₂	0.13	1.4316	0.633	-16.0	1.15	-16.2	10.6	-14.5		145
SrMoO ₄		1.867(o)								
		1.869(e)								
SrTiO ₃	0.5	2.283								

(Continued)

TABLE 21 Summary Optical Properties of Crystals (*Continued*)

Material	Transparency (μm)		Refractive Index (n_o)	Thermo-optic coefficient ($10^{-6}/\text{K}$)					Ref.	
	UV	IR		$\lambda(\mu\text{m})$	dn/dT	$\lambda(\mu\text{m})$	dn/dT	$\lambda(\mu\text{m})$		dn/dT
TeO_2	0.34	4.5	2.177(o) 2.316(e)	0.436	30(o) 25(e)	0.644	9(o) 8(e)		181	
TiO_2	0.42	4.0	2.432(o) 2.683(e)	0.405	4(o) -9(e)				182	
TlBr	0.44	38	2.271							
KRS-5	0.58	42	2.380	0.633 1.014	-250 -237	10.6	-233	30	-195	183 149
KRS-6	0.42	27	2.196							
$\text{Y}_3\text{Al}_5\text{O}_{12}$	0.21	5.2	1.812	0.458	11.9	0.633	9.4	1.06	9.1	184
Y_2O_3	0.29	7.1	1.892	0.663	8.3					148
YVO ₄	0.4	5.0	1.944(o) 2.146(s)	0.488	27.7(o) 22.9(e)	0.633	19.7(o) 12.7(e)	1.34	16.6(o) 10.3(e)	185
ZnGeP ₂	0.8	12.5	3.119(o) 3.156(e)	0.64	359(o) 376(e)	1.0	212(o) 230(e)	10	165(o) 170(e)	164, 186
β -ZnS	0.4	12.5	2.258	0.633	63.5	1.15	49.8	10.6	46.3	159
α -ZnS			2.271(o)							
ZnSe	0.51	19.0	2.275(e)	0.633	91.1	1.15	59.7	10.6	52.0	159
ZnTe	0.6	25	3.15							
$\text{ZrO}_2\text{:Y}_2\text{O}_3$	0.38	6.0	2.113	0.365	16.0	0.458	10.0	0.633	7.9	187

TABLE 22 Summary Optical Properties of Optical and Specialty Glasses

Material	Transparency (μm)		Refractive Index (n_d)	Abbe Number (v_d)	Thermo-optic coeff. ($10^{-6}/\text{K}$)*				Ref.
	UV	IR			$\lambda(\mu\text{m})$	dn/dT	$\lambda(\mu\text{m})$	dn/dT	
N-FK5	0.279	2.78	1.48749	70.41	0.4358	-0.6	1.060	-1.4	188
N-PK52A	0.306		1.49700	81.61	0.4358	-6.0	1.060	-5.7	188
N-ZK7	0.306	2.77	1.50847	61.19	0.4358	7.6	1.060	6.7	188
N-BK7	0.303	2.91	1.51680	64.17	0.4358	3.5	1.060	2.4	188
N-K5	0.313	3.30	1.52249	59.48	0.4358	2.7	1.060	1.4	188
N-KF9	0.346	2.90	1.52346	51.54	0.4358	2.6	1.060	0.9	188
LLF1	0.310	3.30	1.54814	45.75	0.4358	3.9	1.060	2.1	188
N-PSK53A	0.326	2.84	1.61800	63.39	0.4358	-1.8	1.060	-2.9	188
N-BaK4	0.335	3.21	1.56883	55.98	0.4358	4.7	1.060	3.1	188
N-BaLF4	0.340	3.22	1.57956	53.87	0.4358	6.0	1.060	4.2	188
LF5	0.321	3.30	1.58144	40.85	0.4358	3.4	1.060	0.8	188
FF5	0.346		1.59270	35.45	0.6328	0.8			189
N-KzFSN4	0.332	2.65	1.61336	44.49	0.4358	4.7	1.060	2.5	188
N-SK10	0.333	2.74	1.62280	56.90	0.4358	3.5	1.060	2.1	188
N-F2	0.364	3.08	1.62004	36.37	0.4358	5.1	1.060	2.7	188
N-SSK5	0.351	2.93	1.65844	50.88	0.4358	4.2	1.060	2.2	188
N-BaF10	0.356	2.92	1.67003	47.20	0.4358	6.0	1.060	3.8	188
N-BaSF64	0.359	2.98	1.70400	39.38	0.4358	5.9	1.060	2.8	188
N-LaK10	0.346	2.56	1.72000	50.41	0.4358	6.1	1.060	4.2	188
N-LaF33	0.337	2.61	1.78582	44.05	0.4358	10.0	1.060	7.0	188
NbF1	0.313		1.74330	49.22	0.6328	7.9			189
N-LaF2	0.354	2.83	1.74397	44.85	0.4358	2.3	1.060	-0.1	188
N-SF6	0.378	3.21	1.80518	25.43	0.4358	4.8	1.060	-0.8	188
N-LaSF31A	0.344	2.78	1.88300	40.76	0.4358	6.6	1.060	3.3	188
SiO ₂	0.16	3.8	1.45857	67.7	0.5893	10			190
					0.3650	9.3	1.53	8.0	149
GeO ₂	0.30	4.9	$n_D = 1.60832$	41.2					191
BS-39B	0.38	4.9	$n_D = 1.6764$	44.5	0.5893	7.4			192
Corning 9754	0.36	4.8	$n_D = 1.6601$	46.5					193
Schott IRG 2	0.44	5.1	1.8918	30.03	0.436	10.4	3.30	4.4	149, 188
Schott IRG9	0.38	4.1	1.4861	81.02					188
Schott IRG 11	0.44	4.75	1.6809	44.21					188
Schott IRG 100	0.93	13	$n_1 = 2.7235$	—	2,5	103	10.6	56	188
Ohara HTF-1	0.21	6.9	1.44296	92.46					194
ZBLAN	0.25	6.9	$n_D = 1.480$	64	0.6328	-14.5			195, 196
As ₂ S ₃	0.62	11.0	$n_1 = 2.4777$	—	0.6	85	1.0	17	197, 198
As ₂ Se ₃	0.87	17.2	$n_{12} = 2.7728$	—	0.83	55	1.15	33	199, 200
AMTIR-1/TI-20	0.75	(14.5)	$n_1 = 2.6055$	—	1.0	101	10.0	72	200, 201
AMTIR-3/ TI-1173	0.93	16.5	$n_3 = 2.6366$	—	3.0	98	12.0	93	200, 201

*Thermo-optic coefficient in air: $(dn/dT)_{rel}$.

TABLE 23 Room-Temperature Dispersion Formulas for Crystals

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
Ag_3AsS_3	$n_o^2 = 7.483 + \frac{0.474}{\lambda^2 - 0.09} - 0.0019\lambda^2$; $n_e^2 = 6.346 + \frac{0.342}{\lambda^2 - 0.09} - 0.0011\lambda^2$	0.63–4.6(o) 0.59–4.6(e)	202
AgBr	$\frac{n^2 - 1}{n^2 + 2} = 0.453505 + \frac{0.09929\lambda^2}{\lambda^2 - 0.070537} - 0.00150\lambda^2$	0.49–0.67	203
AgCl	$n^2 - 1 = \frac{2.062508\lambda^2}{\lambda^2 - (0.1039054)^2} + \frac{0.9461465\lambda^2}{\lambda^2 - (0.2438691)^2} + \frac{4.300785\lambda^2}{\lambda^2 - (70.85723)^2}$	0.54–21.0	140
AgGaS ₂	$n_o^2 = 5.7975 + \frac{0.2311}{\lambda^2 - 0.0688} - 0.00257\lambda^2$ $n_e^2 = 5.5436 + \frac{0.2230}{\lambda^2 - 0.0946} - 0.00261\lambda^2$	0.58–10.6	142
AgGaSe ₂	$n_o^2 = 6.8507 + \frac{0.4297}{\lambda^2 - 0.1584} - 0.00125\lambda^2$ $n_e^2 = 6.6792 + \frac{0.4597}{\lambda^2 - 0.2122} - 0.00126\lambda^2$	0.73–13.5	143
Al_2O_3 , sapphire	$n_o^2 - 1 = \frac{1.4313493\lambda^2}{\lambda^2 - (0.0726631)^2} + \frac{0.65054713\lambda^2}{\lambda^2 - (0.1193242)^2} + \frac{5.3414021\lambda^2}{\lambda^2 - (18.028251)^2}$ $n_e^2 - 1 = \frac{1.5039759\lambda^2}{\lambda^2 - (0.0740288)^2} + \frac{0.55069141\lambda^2}{\lambda^2 - (0.1216529)^2} + \frac{6.5927379\lambda^2}{\lambda^2 - (20.072248)^2}$	0.2–5.5	204
ALON	$n^2 - 1 = \frac{2.1375\lambda^2}{\lambda^2 - 0.10256^2} + \frac{4.582\lambda^2}{\lambda^2 - 18.868^2}$	0.4–2.3	205
BBO	$n_o^2 = 2.7405 + \frac{0.0184}{\lambda^2 - 0.0179} - 0.0155\lambda^2$ $n_e^2 = 2.3730 + \frac{0.0128}{\lambda^2 - 0.0156} - 0.0044\lambda^2$	0.22–1.06	131
BaF_2	$n^2 - 1 = \frac{0.643356\lambda^2}{\lambda^2 - (0.057789)^2} + \frac{0.506762\lambda^2}{\lambda^2 - (0.10968)^2} + \frac{3.8261\lambda^2}{\lambda^2 - (46.3864)^2}$	0.27–10.3	206
BaTiO_3	$n_o^2 - 1 = \frac{4.187\lambda^2}{\lambda^2 - (0.223)^2}$; $n_e^2 - 1 = \frac{4.064\lambda^2}{\lambda^2 - (0.211)^2}$	0.4–0.7	207
BeO	$n_o^2 - 1 = \frac{1.92274\lambda^2}{\lambda^2 - (0.07908)^2} + \frac{1.24209\lambda^2}{\lambda^2 - (9.7131)^2}$ $n_e^2 - 1 = \frac{1.96939\lambda^2}{\lambda^2 - (0.08590)^2} + \frac{1.67389\lambda^2}{\lambda^2 - (10.4797)^2}$	0.44–7.0	208
BiB_3O_6 , BIBO	$n_x^2 = 3.07403 + \frac{0.03231}{\lambda^2 - 0.03163} - 0.013376\lambda^2$ $n_y^2 = 3.16940 + \frac{0.03717}{\lambda^2 - 0.03483} - 0.01827\lambda^2$ $n_z^2 = 3.6545 + \frac{0.05112}{\lambda^2 - 0.03713} - 0.02261\lambda^2$	0.48–3.1	152
$\text{Bi}_{12}\text{GeO}_{20}$, BGO	$n^2 - 1 = 2.165 + \frac{2.655\lambda^2}{\lambda^2 - 0.07891}$	0.4–0.7	151 209

TABLE 23 Room-Temperature Dispersion Formulas for Crystals (*Continued*)

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
$\text{Bi}_{12}\text{SiO}_{20}$, BSO	$n^2 = 2.72777 + \frac{3.01705\lambda^2}{\lambda^2 - (0.2661)^2}$	0.48–0.7	104 210
C, diamond	$n^2 - 1 = \frac{4.3356\lambda^2}{\lambda^2 - (0.1060)^2} + \frac{0.3306\lambda^2}{\lambda^2 - (0.1750)^2}$	0.225– ∞	211
CaCO_3 , calcite	$n_o - 1 = \frac{0.8559\lambda^2}{\lambda^2 - (0.0588)^2} + \frac{0.8391\lambda^2}{\lambda^2 - (0.141)^2} + \frac{0.0009\lambda^2}{\lambda^2 - (0.197)^2} + \frac{0.6845\lambda^2}{\lambda^2 - (7.005)^2}$ $n_e - 1 = \frac{1.0856\lambda^2}{\lambda^2 - (0.07897)^2} + \frac{0.0988\lambda^2}{\lambda^2 - (0.142)^2} + \frac{0.317\lambda^2}{\lambda^2 - (11.468)^2}$	0.2–2.2	2
CaF_2	$n^2 - 1 = \frac{0.5675888\lambda^2}{\lambda^2 - (0.050263605)^2} + \frac{0.4710914\lambda^2}{\lambda^2 - (0.1003909)^2} + \frac{3.8484923\lambda^2}{\lambda^2 - (34.649040)^2}$ $n^2 - 1 = \frac{0.443749998\lambda^2}{\lambda^2 - (0.00178027854)^2} + \frac{0.444930066\lambda^2}{\lambda^2 - (0.00788536061)^2}$ $+ \frac{0.150133991\lambda^2}{\lambda^2 - (0.0124119491)^2} + \frac{8.85319946\lambda^2}{\lambda^2 - (2752.28175)^2}$	0.23–9.7	155
CaMoO_4	$n_o^2 - 1 = \frac{2.7840\lambda^2}{\lambda^2 - (0.1483)^2} + \frac{1.2425\lambda^2}{\lambda^2 - (11.576)^2}$ $n_e^2 - 1 = \frac{2.8045\lambda^2}{\lambda^2 - (0.1542)^2} + \frac{1.0055\lambda^2}{\lambda^2 - (10.522)^2}$	0.45–3.8	213
CaWO_4	$n_o^2 - 1 = \frac{2.5493\lambda^2}{\lambda^2 - (0.1347)^2} + \frac{0.9200\lambda^2}{\lambda^2 - (10.815)^2}$ $n_e^2 - 1 = \frac{2.6041\lambda^2}{\lambda^2 - (0.1379)^2} + \frac{4.1237\lambda^2}{\lambda^2 - (21.371)^2}$	0.45–4.0	213
CdGeAs_2	$n_o^2 = 10.1064 + \frac{2.2988\lambda^2}{\lambda^2 - 1.0872} + \frac{1.6247\lambda^2}{\lambda^2 - 1370}$ $n_e^2 = 11.8018 + \frac{1.2152\lambda^2}{\lambda^2 - 2.6971} + \frac{1.6922\lambda^2}{\lambda^2 - 1370}$	2.4–11.5	214
CdS	$n_o^2 = 5.1792 + \frac{0.23504}{\lambda^2 - 0.083591} + \frac{0.036927}{\lambda^2 - 0.23504}$ $n_e^2 = 5.2599 + \frac{0.20865}{\lambda^2 - 0.10799} + \frac{0.027527}{\lambda^2 - 0.23305}$	0.51–1.4	215
CdSe	$n_o^2 = 4.2243 + \frac{1.7680\lambda^2}{\lambda^2 - 0.2270} + \frac{3.1200\lambda^2}{\lambda^2 - 3380}$ $n_e^2 = 4.2009 + \frac{1.8875\lambda^2}{\lambda^2 - 0.2171} + \frac{3.6461\lambda^2}{\lambda^2 - 3629}$	1–22	214
CdTe	$n^2 - 1 = \frac{6.1977889\lambda^2}{\lambda^2 - (0.317069)^2} + \frac{3.2243821\lambda^2}{\lambda^2 - (72.0663)^2}$	6–22	216

(Continued)

TABLE 23 Room-Temperature Dispersion Formulas for Crystals (*Continued*)

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
CsLiB ₆ O ₁₀ , CLBO	$n_o^2 = 2.2104 + \frac{0.01018}{\lambda^2 - 0.01424} - 0.01258\lambda^2$ $n_e^2 = 2.0588 + \frac{0.00838}{\lambda^2 - 0.01363} - 0.00607\lambda^2$	0.19–2.75	161
CsI	$n^2 - 1 = \frac{0.34617251\lambda^2}{\lambda^2 - (0.229567)^2} + \frac{1.0080886\lambda^2}{\lambda^2 - (0.1466)^2} + \frac{0.28551800\lambda^2}{\lambda^2 - (0.1810)^2}$ $+ \frac{0.39743178\lambda^2}{\lambda^2 - (0.2120)^2} + \frac{3.3605359\lambda^2}{\lambda^2 - (161.0)^2}$	0.29–50	160
CuGaS ₂	$n_o^2 = 3.9064 + \frac{2.3065\lambda^2}{\lambda^2 - 0.1149} + \frac{1.5479\lambda^2}{\lambda^2 - 738.43}$ $n_e^2 = 4.3165 + \frac{1.8692\lambda^2}{\lambda^2 - 0.1364} + \frac{1.7575\lambda^2}{\lambda^2 - 738.43}$	0.55–11.5	163 164
GaAs	$n^2 = 3.5 + \frac{7.4969\lambda^2}{\lambda^2 - (0.4082)^2} + \frac{1.9347\lambda^2}{\lambda^2 - (37.17)^2}$	1.4–11	217
GaN	$n_o^2 = 3.60 + \frac{1.75\lambda^2}{\lambda^2 - (0.256)^2} + \frac{4.1\lambda^2}{\lambda^2 - (17.86)^2}$ $n_e^2 = 5.35 + \frac{5.08\lambda^2}{\lambda^2 - (18.76)^2}$	<10	218
GaP	$n^2 = 4.1705 + \frac{4.9113\lambda^2}{\lambda^2 - 0.1174} + \frac{1.9928\lambda^2}{\lambda^2 - 756.46}$	0.2–22	219
Ge	$n^2 = 9.28156 + \frac{6.72880\lambda^2}{\lambda^2 - 0.44105} + \frac{0.21307\lambda^2}{\lambda^2 - 3870.1}$	2–12	220
HfO ₂ : 9.8%Y ₂ O ₃	$n^2 - 1 = \frac{1.9558\lambda^2}{\lambda^2 - (0.15494)^2} + \frac{1.345\lambda^2}{\lambda^2 - (0.0634)^2} + \frac{10.41\lambda^2}{\lambda^2 - (27.12)^2}$	0.365–5	168
KBr	$n^2 = 1.39408 + \frac{0.79221\lambda^2}{\lambda^2 - (0.146)^2} + \frac{0.01981\lambda^2}{\lambda^2 - (0.173)^2}$ $+ \frac{0.15587\lambda^2}{\lambda^2 - (0.187)^2} + \frac{0.17673\lambda^2}{\lambda^2 - (60.61)^2} + \frac{2.06217\lambda^2}{\lambda^2 - (87.72)^2}$	0.2–40	221
KH ₂ PO ₄ , KDP	$n_o^2 = 2.259276 + \frac{0.01008956}{\lambda^2 - 0.0129426} + \frac{13.00522\lambda^2}{\lambda^2 - 400}$ $n_e^2 = 2.132668 + \frac{0.008637494}{\lambda^2 - 0.0122810} + \frac{3.2279924\lambda^2}{\lambda^2 - 400}$	0.2–15	70
KI	$n^2 = 1.47285 + \frac{0.16512\lambda^2}{\lambda^2 - (0.129)^2} + \frac{0.41222\lambda^2}{\lambda^2 - (0.175)^2} + \frac{0.44163\lambda^2}{\lambda^2 - (0.187)^2}$ $+ \frac{0.16076\lambda^2}{\lambda^2 - (0.219)^2} + \frac{0.33571\lambda^2}{\lambda^2 - (69.44)^2} + \frac{1.92474\lambda^2}{\lambda^2 - (98.04)^2}$	0.25–50	221
KNbO ₃	$n_x^2 = 4.4222 + \frac{0.09972}{\lambda^2 - 0.05496} - 0.01976\lambda^2$ $n_y^2 = 4.8353 + \frac{0.12808}{\lambda^2 - 0.05674} - 0.02528\lambda^2 + 1.8590 \cdot 10^{-6} \lambda^4 - 1.0689 \cdot 10^{-6} \lambda^6$ $n_z^2 = 4.9856 + \frac{0.15266}{\lambda^2 - 0.06331} - 0.02831\lambda^2 + 2.0754 \cdot 10^{-6} \lambda^4 - 1.2131 \cdot 10^{-6} \lambda^6$	0.40–5.3	222

TABLE 23 Room-Temperature Dispersion Formulas for Crystals (*Continued*)

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
KTaO ₃	$n^2 - 1 = \frac{3.591\lambda^2}{\lambda^2 - (0.193)^2}$	0.4–1.06	223
KTiOPO ₄ KTP	$n_x^2 = 3.29100 + \frac{0.04140}{\lambda^2 - 0.03978} + \frac{9.35522}{\lambda^2 - 31.45571}$ $n_y^2 = 3.45018 + \frac{0.04341}{\lambda^2 - 0.04597} + \frac{16.98825}{\lambda^2 - 39.43799}$ $n_z^2 = 4.59423 + \frac{0.06206}{\lambda^2 - 0.04763} + \frac{110.80672}{\lambda^2 - 86.12171}$	0.43–3.54	172
LiB ₃ O ₅ , LBO	$n_x^2 = 2.45768 + \frac{0.0098877}{\lambda^2 - 0.026095} - 0.013847\lambda^2$ $n_y^2 = 2.52500 + \frac{0.017123}{\lambda^2 - 0.0060517} - 0.0087838\lambda^2$ $n_z^2 = 2.58488 + \frac{0.012737}{\lambda^2 - 0.016293} - 0.016293\lambda^2$	0.29–1.06	224
LiCaAlF ₆ , LiCAF	$n_o^2 = 1.92552 + \frac{0.00492}{\lambda^2 - 0.00569} - 0.00421 \cdot \lambda^2$ $n_e^2 = 1.92155 + \frac{0.00494}{\lambda^2 - 0.00617} - 0.00373 \cdot \lambda^2$	0.4–1.0	174
LiF	$n^2 - 1 = \frac{0.92549\lambda^2}{\lambda^2 - (0.07376)^2} + \frac{6.96747\lambda^2}{\lambda^2 - (32.79)^2}$	0.1–10	221
LiIO ₃	$n_o^2 = 2.03132 + \frac{1.37623\lambda^2}{\lambda^2 - 0.0350823} + \frac{1.06745\lambda^2}{\lambda^2 - 169.0}$ $n_e^2 = 1.83086 + \frac{1.08807\lambda^2}{\lambda^2 - 0.0313810} + \frac{0.554582\lambda^2}{\lambda^2 - 158.76}$	0.5–5	82
LiNbO ₃	$n_o^2 - 1 = \frac{2.9804\lambda^2}{\lambda^2 - 0.01764} + \frac{1.2290\lambda^2}{\lambda^2 - 0.05914} + \frac{12.614\lambda^2}{\lambda^2 - 474.60}$ $n_e^2 - 1 = \frac{2.4272\lambda^2}{\lambda^2 - 0.02047} + \frac{0.5981\lambda^2}{\lambda^2 - 0.06660} + \frac{8.9543\lambda^2}{\lambda^2 - 416.08}$	0.4–5.0	225
LiYF ₄	$n_o^2 = 1.38757 + \frac{0.70757\lambda^2}{\lambda^2 - 0.00931} + \frac{0.18849\lambda^2}{\lambda^2 - 50.99741}$ $n_e^2 = 1.31021 + \frac{0.84903\lambda^2}{\lambda^2 - 0.00876} + \frac{0.53607\lambda^2}{\lambda^2 - 134.9566}$	0.23–2.6	176
MgAl ₂ O ₄	$n^2 - 1 = \frac{1.8938\lambda^2}{\lambda^2 - (0.09942)^2} + \frac{3.0755\lambda^2}{\lambda^2 - (15.826)^2}$	0.35–5.5	226
MgF ₂	$n_o^2 - 1 = \frac{0.48755708\lambda^2}{\lambda^2 - (0.04338408)^2} + \frac{0.39875031\lambda^2}{\lambda^2 - (0.09461442)^2} + \frac{2.3120353\lambda^2}{\lambda^2 - (23.793604)^2}$ $n_e^2 - 1 = \frac{0.41344023\lambda^2}{\lambda^2 - (0.03684262)^2} + \frac{0.50497499\lambda^2}{\lambda^2 - (0.09076162)^2} + \frac{2.4904862\lambda^2}{\lambda^2 - (23.771995)^2}$	0.20–7.04	227
MgO	$n^2 - 1 = \frac{1.111033\lambda^2}{\lambda^2 - (0.0712465)^2} + \frac{0.8460085\lambda^2}{\lambda^2 - (0.1375204)^2} + \frac{7.808527\lambda^2}{\lambda^2 - (26.89302)^2}$	0.36–5.4	177

(Continued)

TABLE 23 Room-Temperature Dispersion Formulas for Crystals (*Continued*)

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
NaCl	$n^2 = 1.00055 + \frac{0.19800\lambda^2}{\lambda^2 - (0.050)^2} + \frac{0.48398\lambda^2}{\lambda^2 - (0.100)^2} + \frac{0.38696\lambda^2}{\lambda^2 - (0.128)^2}$ $+ \frac{0.25998\lambda^2}{\lambda^2 - (0.158)^2} + \frac{0.08796\lambda^2}{\lambda^2 - (40.50)^2} + \frac{3.17064\lambda^2}{\lambda^2 - (60.98)^2} + \frac{0.30038\lambda^2}{\lambda^2 - (120.34)^2}$	0.2–30	221
$[\text{NH}_4]_2\text{CO}_3$ Urea	$n_o^2 = 2.1823 + \frac{0.0125}{\lambda^2 - 0.0300}$ $n_e^2 = 2.51527 + \frac{0.0240}{\lambda^2 - 0.0300} + \frac{0.020(\lambda - 1.52)}{(\lambda - 1.52)^2 + 0.8771}$	0.3–1.06	228
$\text{NH}_4\text{H}_2\text{PO}_4$ ADP	$n_o^2 = 2.302842 + \frac{0.011125165}{\lambda^2 - 0.01325366} + \frac{15.102464\lambda^2}{\lambda^2 - 400}$ $n_e^2 = 2.163510 + \frac{0.009616676}{\lambda^2 - 0.01298912} + \frac{5.919896\lambda^2}{\lambda^2 - 400}$	0.2–1.5	229
PbMoO_4	$n_o^2 - 1 = \frac{3.54642\lambda^2}{\lambda^2 - (0.18518)^2} + \frac{0.58270\lambda^2}{\lambda^2 - (0.33764)^2}$ $n_e^2 - 1 = \frac{3.52555\lambda^2}{\lambda^2 - (0.17950)^2} + \frac{0.20660\lambda^2}{\lambda^2 - (0.32537)^2}$	0.44–1.08	230
PbTe	$n^2 - 1 = \frac{30.046\lambda^2}{\lambda^2 - (1.563)^2}$	4.0–12.5	231
PbTiO_3	$n_o^2 - 1 = \frac{5.363\lambda^2}{\lambda^2 - (0.224)^2}; \quad n_e^2 - 1 = \frac{5.366\lambda^2}{\lambda^2 - (0.217)^2}$	0.45–1.15	232
RbTiOPO_4 , RTP	$n_x^2 = 1.6795 + \frac{1.4281\lambda^2}{\lambda^2 - 0.0325} - 0.0119\lambda^2$ $n_y^2 = 2.0360 + \frac{1.0883\lambda^2}{\lambda^2 - 0.0437} - 0.0090\lambda^2$ $n_z^2 = 2.2864 + \frac{1.1280\lambda^2}{\lambda^2 - 0.0562} - 0.0188\lambda^2$	0.4–1.5	180
Si	$n_2 - 1 = \frac{10.6684293\lambda^2}{\lambda^2 - (0.301516485)^2} + \frac{0.003043475\lambda^2}{\lambda^2 - (1.13475115)^2} + \frac{1.54133408\lambda^2}{\lambda^2 - (1104.0)^2}$	1.36–11	67
α - SiO_2 quartz	$n_o^2 - 1 = \frac{0.663044\lambda^2}{\lambda^2 - (0.060)^2} + \frac{0.517852\lambda^2}{\lambda^2 - (0.106)^2} + \frac{0.175912\lambda^2}{\lambda^2 - (0.119)^2}$ $+ \frac{0.565380\lambda^2}{\lambda^2 - (8.844)^2} + \frac{1.675299\lambda^2}{\lambda^2 - (20.742)^2}$ $n_e^2 - 1 = \frac{0.665721\lambda^2}{\lambda^2 - (0.060)^2} + \frac{0.503511\lambda^2}{\lambda^2 - (0.106)^2} + \frac{0.214792\lambda^2}{\lambda^2 - (0.119)^2}$ $+ \frac{0.539173\lambda^2}{\lambda^2 - (8.792)^2} + \frac{1.807613\lambda^2}{\lambda^2 - (197.70)^2}$	0.18–0.71	178
SrF_2	$n^2 - 1 = \frac{0.67805894\lambda^2}{\lambda^2 - (0.05628989)^2} + \frac{0.37140533\lambda^2}{\lambda^2 - (0.10801027)^2} + \frac{3.8484723\lambda^2}{\lambda^2 - (34.649040)^2}$	0.21–11.5	145

TABLE 23 Room-Temperature Dispersion Formulas for Crystals (*Continued*)

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
SrMoO ₄	$n_o^2 - 1 = \frac{2.4839\lambda^2}{\lambda^2 - (0.1451)^2} + \frac{0.1015\lambda^2}{\lambda^2 - (4.603)^2}$ $n_e^2 - 1 = \frac{2.4923\lambda^2}{\lambda^2 - (0.1488)^2} + \frac{0.1050\lambda^2}{\lambda^2 - (4.544)^2}$	0.45–2.4	213
SrTiO ₃	$n^2 - 1 = \frac{3.042143\lambda^2}{\lambda^2 - (0.1475902)^2} + \frac{1.170065\lambda^2}{\lambda^2 - (0.2953086)^2} + \frac{30.83326\lambda^2}{\lambda^2 - (33.18606)^2}$	0.43–3.8	156
TeO ₂	$n_o^2 - 1 = \frac{2.584\lambda^2}{\lambda^2 - (0.1342)^2} + \frac{1.157\lambda^2}{\lambda^2 - (0.2638)^2}$ $n_e^2 - 1 = \frac{2.823\lambda^2}{\lambda^2 - (0.1342)^2} + \frac{1.542\lambda^2}{\lambda^2 - (0.2631)^2}$	0.4–1.0	181
TiO ₂	$n^2 = 5.913 + \frac{0.2441}{\lambda^2 - 0.0803}; \quad n_e^2 = 7.197 + \frac{0.3322}{\lambda^2 - 0.0843}$	0.43–1.5	182
TlBr	$\frac{n^2 - 1}{n^2 + 2} = 0.48484 + \frac{0.10279\lambda^2}{\lambda^2 - 0.90000} - 0.0047896\lambda^2$	0.54–0.65	203
Tl[Br, I], KRS-5	$n^2 - 1 = \frac{1.8293958\lambda^2}{\lambda^2 - (0.150)^2} + \frac{1.6675593\lambda^2}{\lambda^2 - (0.250)^2} + \frac{1.1210424\lambda^2}{\lambda^2 - (0.350)^2}$ $+ \frac{0.04513366\lambda^2}{\lambda^2 - (0.450)^2} + \frac{12.380234\lambda^2}{\lambda^2 - (164.59)^2}$	0.58–39.4	183
TlCl	$\frac{n^2 - 1}{n^2 + 2} = 0.47856 + \frac{0.07858\lambda^2}{\lambda^2 - 0.08377} - 0.00881\lambda^2$	0.43–0.66	203
Tl[Cl, Br], KRS-6	$n^2 - 1 = \frac{3.821\lambda^2}{\lambda^2 - (0.02234)^2} - 0.000877\lambda^2$	0.6–24	233
Y ₃ Al ₅ O ₁₂ , YAG	$n^2 - 1 = \frac{2.28200\lambda^2}{\lambda^2 - 0.01185} + \frac{3.27644\lambda^2}{\lambda^2 - 282.734}$	0.4–5.5	234
Y ₂ O ₃	$n^2 - 1 = \frac{2.578\lambda^2}{\lambda^2 - (0.1387)^2} + \frac{3.935\lambda^2}{\lambda^2 - (22.936)^2}$	0.2–12	235
YVO ₄	$n_o^2 = 3.778790 + \frac{0.07479}{\lambda^2 - 0.045731} - 0.009701\lambda^2$ $n_e^2 = 4.607200 + \frac{0.108087}{\lambda^2 - 0.052495} - 0.014305\lambda^2$	0.48–1.34	185
ZnGeP ₂	$n_o^2 = 8.0409 + \frac{1.68625\lambda^2}{\lambda^2 - 0.40824} + \frac{1.2880\lambda^2}{\lambda^2 - 611.05}$ $n_e^2 = 8.0929 + \frac{1.8649\lambda^2}{\lambda^2 - 0.41468} + \frac{0.84052\lambda^2}{\lambda^2 - 452.05}$	2.5–9.5	236
β -ZnS	$n^2 - 1 = \frac{0.33904026\lambda^2}{\lambda^2 - (0.31423026)^2} + \frac{3.7606868\lambda^2}{\lambda^2 - (0.1759417)^2} + \frac{2.7312353\lambda^2}{\lambda^2 - (33.886560)^2}$	0.55–10.5	144

(Continued)

TABLE 23 Room-Temperature Dispersion Formulas for Crystals (*Continued*)

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
α -ZnS	$n_o^2 - 1 = 3.4175 + \frac{1.7396\lambda^2}{\lambda^2 - (0.2677)^2}$; $n_e^2 - 1 = 3.4264 + \frac{1.7491\lambda^2}{\lambda^2 - (0.2674)^2}$	0.36–1.4	237
ZnSe	$n^2 - 1 = \frac{4.2980149\lambda^2}{\lambda^2 - (0.1920630)^2} + \frac{0.62776557\lambda^2}{\lambda^2 - 0.37878260^2} + \frac{2.8955633\lambda^2}{\lambda^2 - 46.994595^2}$	0.55–18	145
ZnTe	$n^2 = 9.92 \frac{0.42530}{\lambda^2 - (0.37766)^2} + \frac{2.63580}{\lambda^2 / (56.5)^2 - 1}$	0.55–30	238
ZrO ₂ : 12%Y ₂ O ₃	$n^2 - 1 = \frac{1.347091\lambda^2}{\lambda^2 - (0.166739)^2} + \frac{2.117788\lambda^2}{\lambda^2 - (0.166739)^2} + \frac{9.452943\lambda^2}{\lambda^2 - (24.320570)^2}$	0.36–5.1	186

TABLE 24 Room-Temperature Dispersion Formula for Glasses

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
N-FK5	$n^2 = 1 + \frac{1.07715032\lambda^2}{\lambda^2 - 0.00676601657} + \frac{0.168079109\lambda^2}{\lambda^2 - 0.0230642817} + \frac{0.851889892\lambda^2}{\lambda^2 - 89.0498778}$	0.28–2.4	188
N-PK52A	$n^2 = 1 + \frac{1.029607\lambda^2}{\lambda^2 - 0.00516800155} + \frac{0.1880506\lambda^2}{\lambda^2 - 0.0166658798} + \frac{0.736488165\lambda^2}{\lambda^2 - 138.964129}$	0.3–2.4	188
N-ZK7	$n^2 = 1 + \frac{1.07715032\lambda^2}{\lambda^2 - 0.00676601657} + \frac{0.168079109\lambda^2}{\lambda^2 - 0.0230642817} + \frac{0.851889892\lambda^2}{\lambda^2 - 89.0498778}$	0.3–2.4	188
N-BK7	$n^2 = 1 + \frac{1.03961212\lambda^2}{\lambda^2 - 0.00600069867} + \frac{0.231792344\lambda^2}{\lambda^2 - 0.0200179144} + \frac{1.01046945\lambda^2}{\lambda^2 - 103.560653}$	0.3–2.4	188
N-K5	$n^2 = 1 + \frac{1.08511833\lambda^2}{\lambda^2 - 0.00661099503} + \frac{0.199562005\lambda^2}{\lambda^2 - 0.024110866} + \frac{0.930511663\lambda^2}{\lambda^2 - 111.982777}$	0.3–2.4	188
N-KF9	$n^2 = 1 + \frac{1.19286778\lambda^2}{\lambda^2 - 0.00839154696} + \frac{0.0893346571\lambda^2}{\lambda^2 - 0.0404010786} + \frac{0.920819805\lambda^2}{\lambda^2 - 112.572446}$	0.36–2.4	188
LLF1	$n^2 = 1 + \frac{1.21640125\lambda^2}{\lambda^2 - 0.00857807248} + \frac{0.13366454\lambda^2}{\lambda^2 - 0.0420143003} + \frac{0.883399468\lambda^2}{\lambda^2 - 107.593060}$	0.37–1.01 [†]	188
N-PSK53A	$n^2 = 1 + \frac{1.38121836\lambda^2}{\lambda^2 - 0.00706416337} + \frac{0.196745645\lambda^2}{\lambda^2 - 0.0233251345} + \frac{0.886089205\lambda^2}{\lambda^2 - 97.4847345}$	0.31–2.5	188
N-BaK4	$n^2 = 1 + \frac{1.28834642\lambda^2}{\lambda^2 - 0.00779980626} + \frac{0.132817724\lambda^2}{\lambda^2 - 0.0315631177} + \frac{0.945395373\lambda^2}{\lambda^2 - 105.965875}$	0.334–2.5	188
N-BaLF4	$n^2 = 1 + \frac{1.31004128\lambda^2}{\lambda^2 - 0.0079659645} + \frac{0.142038259\lambda^2}{\lambda^2 - 0.0330672072} + \frac{0.964929351\lambda^2}{\lambda^2 - 109.19732}$	0.36–2.4	188
LF5	$n^2 = 1 + \frac{1.28035628\lambda^2}{\lambda^2 - 0.00929854416} + \frac{0.163505973\lambda^2}{\lambda^2 - 0.0449135769} + \frac{0.893930112\lambda^2}{\lambda^2 - 110.493685}$	0.33–2.4	188
FF5	$n^2 = 1 + \frac{1.3241093\lambda^2}{\lambda^2 - 0.01004} + \frac{0.147104\lambda^2}{\lambda^2 - 0.05705} - 0.0095123\lambda^2$	0.37–1.01	189*
N-KzFSN4	$n^2 = 1 + \frac{1.35055424\lambda^2}{\lambda^2 - 0.0087628207} + \frac{0.197575506\lambda^2}{\lambda^2 - 0.0371767201} + \frac{1.09962992\lambda^2}{\lambda^2 - 90.3866994}$	0.36–2.4	188

TABLE 24 Room-Temperature Dispersion Formula for Glasses (*Continued*)

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
N-SK10	$n^2 = 1 + \frac{1.34972093\lambda^2}{\lambda^2 - 0.00736272269} + \frac{0.238587973\lambda^2}{\lambda^2 - 0.0253765327} + \frac{0.966733600\lambda^2}{\lambda^2 - 103.502909}$	0.334–2.5	188
N-F2	$n^2 = 1 + \frac{1.39757037\lambda^2}{\lambda^2 - 0.00995906143} + \frac{0.159201403\lambda^2}{\lambda^2 - 0.0546931752} + \frac{1.26865430\lambda^2}{\lambda^2 - 119.248346}$	0.33–2.4	188
N-SSK5	$n^2 = 1 + \frac{1.59222659\lambda^2}{\lambda^2 - 0.00920284626} + \frac{0.103520774\lambda^2}{\lambda^2 - 0.0423530072} + \frac{1.05174016\lambda^2}{\lambda^2 - 106.927374}$	0.36–2.4	188
N-BaF10	$n^2 = 1 + \frac{1.58514950\lambda^2}{\lambda^2 - 0.00926681282} + \frac{0.143559385\lambda^2}{\lambda^2 - 0.0424489805} + \frac{1.08521269\lambda^2}{\lambda^2 - 105.613573}$	0.37–1.01	188
N-BaSF64	$n^2 = 1 + \frac{1.65554268\lambda^2}{\lambda^2 - 0.00104485644} + \frac{1.71319770\lambda^2}{\lambda^2 - 0.0499394756} + \frac{1.33664448\lambda^2}{\lambda^2 - 118.961472}$	0.365–2.5	188
N-LaK10	$n^2 = 1 + \frac{1.72878017\lambda^2}{\lambda^2 - 0.00886014635} + \frac{0.169257825\lambda^2}{\lambda^2 - 0.0363416509} + \frac{1.19386956\lambda^2}{\lambda^2 - 82.9009069}$	0.36–2.4	188
N-LaF33	$n^2 = 1 + \frac{1.79653417\lambda^2}{\lambda^2 - 0.00927313493} + \frac{0.311577903\lambda^2}{\lambda^2 - 0.0358201181} + \frac{1.15981863\lambda^2}{\lambda^2 - 87.3448712}$	0.32–2.5	188
NbF1	$n^2 = 1 + \frac{1.9550984\lambda^2}{\lambda^2 - 0.011628} + \frac{0.0165083\lambda^2}{\lambda^2 - 0.064357} - 0.012813\lambda^2$	0.37–1.01	189*
N-LaF2	$n^2 = 1 + \frac{1.80984227\lambda^2}{\lambda^2 - 0.0101711622} + \frac{0.15729555\lambda^2}{\lambda^2 - 0.0442431765} + \frac{1.0930037\lambda^2}{\lambda^2 - 100.687748}$	0.35–2.4	188
N-SF6	$n^2 = 1 + \frac{1.77931763\lambda^2}{\lambda^2 - 0.0133714182} + \frac{0.338149866\lambda^2}{\lambda^2 - 0.0617533621} + \frac{2.08734474\lambda^2}{\lambda^2 - 174.01759}$	0.40–2.4	188
N-LASF31A	$n^2 = 1 + \frac{1.96485075\lambda^2}{\lambda^2 - 0.00982060155} + \frac{0.475231259\lambda^2}{\lambda^2 - 0.0344713438} + \frac{1.48360109\lambda^2}{\lambda^2 - 110.739863}$	0.334–2.5	188
Fused silica	$n^2 - 1 = \frac{0.6961663\lambda^2}{\lambda^2 - (0.0684043)^2} + \frac{0.4079426\lambda^2}{\lambda^2 - (0.1162414)^2} + \frac{0.8974794\lambda^2}{\lambda^2 - (9.896161)^2}$	0.21–3.71	190
Fused germania	$n^2 - 1 = \frac{0.80686642\lambda^2}{\lambda^2 - (0.06897261)^2} + \frac{0.71815848\lambda^2}{\lambda^2 - (0.1539661)^2} + \frac{0.85416831\lambda^2}{\lambda^2 - (11.841931)^2}$	0.43–4.5	191*
BS-93B	$n^2 - 1 = \frac{1.7441\lambda^2}{\lambda^2 - (0.1155)^2} + \frac{1.6465\lambda^2}{\lambda^2 - (14.981)^2}$	0.43–4.5	192
CORTRAN 9754	$n^2 - 1 = \frac{1.66570\lambda^2}{\lambda^2 - (0.10832)^2} + \frac{0.04059\lambda^2}{\lambda^2 - (0.23813)^2} + \frac{1.31792\lambda^2}{\lambda^2 - (13.57622)^2}$	0.4–5.5	193
IRG 2	$n^2 - 1 = \frac{2.07670\lambda^2}{\lambda^2 - (0.11492)^2} + \frac{0.35738\lambda^2}{\lambda^2 - (0.23114)^2} + \frac{2.88166\lambda^2}{\lambda^2 - (17.48306)^2}$	0.365–4.6	188*
IRG 9	$n^2 - 1 = \frac{2.07670\lambda^2}{\lambda^2 - (0.11492)^2} + \frac{0.35738\lambda^2}{\lambda^2 - (0.23114)^2} + \frac{2.88166\lambda^2}{\lambda^2 - (17.48306)^2}$	0.365–4.6	188*

(Continued)

TABLE 24 Room-Temperature Dispersion Formula for Glasses (*Continued*)

Material	Dispersion Formula (Wavelength, λ , in μm)	Range (μm)	Ref.
IRG 11	$n^2 - 1 = \frac{1.7531\lambda^2}{\lambda^2 - (0.1185)^2} + \frac{0.4346\lambda^2}{\lambda^2 - (8.356)^2}$	0.48–3.3	188*
IRG 100	$n^2 = 4.5819 + \frac{2.2693\lambda^2}{\lambda^2 - (0.447)^2} - 0.000928\lambda^2$	1–14	188*
HTF-1	$n^2 - 1 = \frac{1.06375\lambda^2}{\lambda^2 - (0.078958)^2} + \frac{0.80098\lambda^2}{\lambda^2 - (15.1579)^2}$	0.37–5	194*
ZBLAN	$n^2 - 1 = \frac{1.168\lambda^2}{\lambda^2 - (0.0954)^2} + \frac{2.77\lambda^2}{\lambda^2 - (25.0)^2}$	0.50–48	195
AMTIR-1/ TI-20	$n^2 - 1 = \frac{5.298\lambda^2}{\lambda^2 - (0.29007)^2} + \frac{0.6039\lambda^2}{\lambda^2 - (32.022)^2}$	1–14	201*
AMTIR-3/ TI-1173	$n^2 - 1 = \frac{5.8505\lambda^2}{\lambda^2 - (0.29192)^2} + \frac{1.4536\lambda^2}{\lambda^2 - (42.714)^2}$	3–14	201*
	$n^2 - 1 = \frac{5.8357\lambda^2}{\lambda^2 - (0.29952)^2} + \frac{1.064\lambda^2}{\lambda^2 - (38.353)^2}$	0.9–14	239

*Our dispersion equation from referenced data.

†Schott dispersion formula range; data available to 2.3 μm .

Vibrational modes of many optical materials are summarized in Tables 25 through 40 for a number of common optical crystal types. These tables give number and type of zone-center (i.e., the wave vector ≈ 0 , where Γ is the usual symbol denoting the center of the Brillouin zone) optical modes predicted by group theory (and observed in practice) as well as the frequency (in wave number) of the infrared-active and Raman modes. Mulliken notation is used. Table 41 summarizes the available lattice vibration dispersion models for many crystals.

Parameters for ultraviolet and infrared absorption models for crystals and glasses are not included in this edition. Please refer to the previous edition¹ for these parameters.

TABLE 25 Optical Modes of Crystals with Diamond Structure; Space Group: $Fd3m (O_h^h) \#227; \Gamma = F_{2g}(R)$

Material	Raman Mode Location (cm^{-1})	Ref.
C (diamond)	1332.4	240
Si	519.5	241
Ge	300.6	242

TABLE 26 Optical Modes of Crystals with Cesium Chloride Structure; Space Group: $Pm3m (O_h^5) \#221; = F_{1u}(IR)$

Material	Mode location (cm^{-1})			Material	Mode location (cm^{-1})		
	ω_{TO}	ω_{LO}	Ref.		ω_{TO}	ω_{LO}	Ref.
CsBr	75	113	243	TlBr	43	101	245
CsCl	99	160	243, 244	TlCl	63	158	245
CsI	62	88	243	TlI	40	87	246

TABLE 27 Optical Modes of Crystals with Sodium Chloride Structure; Space Group: $Fm\bar{3}m (O_h^1) \#225$; $\Gamma = F_{1u}(\text{IR})$

Material	Mode location (cm^{-1})			Material	Mode location (cm^{-1})		
	ω_{TO}	ω_{LO}	Ref.		ω_{TO}	ω_{LO}	Ref.
AgBr	79	138	245	MgO	401	718	245, 251
AgCl	106	196	245	MgS	237	430	248
BaO	132	425	247	NaBr	134	209	245
BaS	158	230	248	NaCl	164	264	245
CaO	295	557	247, 249	NaF	244	418	245
CaS	256	376	248	NaI	117	176	245
CdO	270	380	250	NiO	401	580	245
KBr	113	165	245	SrO	227	487	247, 249
KCl	142	214	245	SrS	194	284	248
KF	190	326	245	PbS	71	212	252
KI	101	139	245	PbSe	39	116	253
LiF	306	659	245, 251	PbTe	32	112	254

TABLE 28 Optical Modes of Crystals with Zinblende Structure; Space Group: $F\bar{4}3m(T_d^2) \#216$; $\Gamma = F_2(\text{R, IR})$

Material	Mode location (cm^{-1})			Material	Mode location (cm^{-1})		
	ω_{TO}	ω_{LO}	Ref.		ω_{TO}	ω_{LO}	Ref.
AlAs	364	402	255	GaSb	230	240	262
AlSb	319	340	256	InAs	217	239	262, 263
BN	1055	1305	257, 258	InP	304	345	256
BP	799	829	257, 258	InSb	179	191	263, 264
CdTe	141	169	259	-SiC	793	970	242, 265
CuCl	172	210	260	ZnS	282	352	259, 266
GaAs	269	292	256	ZnSe	206	252	259
GaP	367	403	256, 261	ZnTe	178	208	259, 267

TABLE 29 Optical Modes of Crystals with Fluorite Structure; Space Group: $Fm\bar{3}m (O_h^1) \#225$; $\Gamma = F_{1u}(\text{IR}) + F_{2g}(\text{R})$

Material	Mode locations (cm^{-1})			Ref.
	$F_{1u}(\omega_{\text{TO}})$	$F_{1u}(\omega_{\text{LO}})$	F_{2g}	
BaCl ₂			185	268
BaF ₂	184	319	241	269, 270, 271
CaF ₂	258	473	322	269, 270, 271
CdF ₂	202	384	317	271, 272
EuF ₂	194	347	287	268
β -PbF ₂	102	337	256	268, 272
SrCl ₂	147	243	182	268, 273, 274
SrF ₂	217	366	286	269, 270, 271
ThO ₂	281	568		275
UO ₂	281	555	445	275, 276
ZrO ₂	354	680	605	277

TABLE 30 Optical Modes of Crystals with Corundum Structure; Space Group: ($R\bar{3}c(D_{3d}^6)\#167$); $\Gamma = 2A_{1g}(R) + 2A_{1u}(-) + 3A_{2g}(-) + 2A_{2u}(IR, E \parallel c) + 5E_g(R) + 4E_u(IR, E \perp c)$

Material	Mode locations (cm^{-1})												
	Infrared modes (LO in parentheses)						Raman modes						
	E_u	E_u	E_u	E_u	A_{1u}	A_{1u}	E_g	E_g	E_g	E_g	E_g	A_{1g}	A_{1g}
Al_2O_3 (Refs. 35, 278)	385 (388)	422 (480)	569 (625)	635 (900)	400 (512)	583 (871)	378	432	451	578	751	418	645
Cr_2O_3 (Refs. 279, 280)	417 (420)	444 (446)	532 (602)	613 (766)	538 (602)	613 (759)	—	351	397	530	609	303	551
Fe_2O_3 (Refs. 280, 281)	227 (230)	286 (368)	437 (494)	524 (662)	299 (414)	526 (662)	245	293	298	413	612	226	500

TABLE 31 Optical Modes of Crystals with Wurtzite Structure; Space Group: $P6_3mc(C_{6v}^4)\#186$; $\Gamma = A_1(R, IR E \parallel c) + 2B_1(-) + E_1(R, IR E \perp c) + 2E_2(R)$

Material	Mode locations (cm^{-1})				
	IR modes TO & (LO)		Raman modes		Ref.
	E_1	A_1	E_2	E_2	
β -AgI	106 (124)	106 (124)	17	112	282
AlN	673 (916)	614 (893)	252	660	283
BeO	725 (1095)	684 (1085)	340	684	284, 285
CdS	235 (305)	228 (305)	44	252	285, 286
CdSe	172 (210)	166 (211)	34	176	287
GaN	560 (746)	533 (744)	145	568	218, 288
α -SiC	797 970	788 964	149	788	289
ZnO	407 (583)	381 (574)	101	437	285, 290
ZnS	274 (352)	274 (352)	55	280	285

TABLE 32 Optical Modes of Crystals with Trigonal Selenium Structure; Space Group: $P3_121(D_3^4)\#152$; $\Gamma = A_1(R) + A_2(IR, E \parallel c) + 2E(IR, E \perp c, R)$

Material	Mode locations (cm^{-1})				
	Raman	Infrared, $E \perp c$		Infrared, $E \parallel c$	Ref.
	A_1	E	E	A_2	
Se	237	144 (150)	225 (225)	102 (106)	291, 292
Te	120	92 (106)	144 (145)	90 (96)	293

TABLE 33 Optical Modes of Crystals with α -Quartz Structure; Space Group: $P3_221$ (D_3^6)
 $\#154; \Gamma = 4A_1(R) + 4A_2(IR, E \parallel c) + 8E(IR, E \perp c, R)$

Material	Infrared modes, $E \perp c$							
	E	E	E	E	E	E	E	E
SiO ₂ (Refs. 294, 295, 296)	128 (128)	265 (270)	394 (403)	451 (511)	697 (699)	796 (809)	1067 (1159)	1164 (1230)
GeO ₂ (Refs. 297)	121 (121)	166 (166)	326 (372)	385 (456)	492 (512)	583 (595)	857 (919)	961 (972)
Material	Infrared modes, $E \parallel c$				Raman modes			
	A ₂	A ₂	A ₂	A ₂	A ₁	A ₁	A ₁	A ₁
SiO ₂ (Refs. 294, 295, 296)	364 (388)	500 (552)	777 (789)	1080 (1239)	207	356	464	1085
GeO ₂ (Ref. 297)					212	261	440	880

TABLE 34 Optical Modes of Crystals with Rutile Structure; Space Group: $P4_2/mnm$ (D_{2h}^{14}) #136;
 $\Gamma = A_{1g}(R) + A_{2g}(-) + A_{2u}(IR, E \parallel c) + B_{1g}(R) + B_{2g}(R) + 2B_{1u}(-) + E_g(R) + 3E_u(IR, E \perp c)$

Material	Mode locations (cm ⁻¹)							
	Infrared modes (LO in parentheses)				Raman modes			
	E _u	E _u	E _u	A _{2u}	B _{1g}	E _g	A _{1g}	B _{2g}
CoF ₂ (Refs. 298, 299, 230)	190 (234)	270 (276)	405 (529)	345 (506)	68	246	366	494
FeF ₂ (Refs. 299, 300, 301)	173 (231)	244 (248)	405 (530)	307 (487)	73	257	340	496
GeO ₂ (Refs. 297, 303, 304)	300 (345)	370 (470)	635 (815)	455 (755)	97	680	702	870
MgF ₂ (Refs. 301, 302, 305, 306, 307)	247 (303)	410 (415)	450 (617)	399 (625)	92	295	410	515
SnO ₂ (Refs. 308, 309, 310)	243 (273)	284 (368)	605 (757)	465 (703)	123	475	634	776
TiO ₂ (Refs. 32, 36, 301, 311)	189 (367)	382 (444)	508 (831)	172 (796)	143	447	612	826
ZnF ₂ (Refs. 299, 300, 303, 305)	173 (227)	244 (264)	380 (498)	294 (488)	70	253	350	522

TABLE 35 Optical Modes of Crystals with Scheelite Structure; Space Group: $I4_1/a(C_{4h}^6) \#88$; $\Gamma = 3A_g(R) + 4A_u(IR, E||c) + 5E_g(R) + 3E_u(-) + 5E_g(R) + 4E_u(IR, E \perp c)$

Material	Ref.	Infrared modes, $E \perp c$				Infrared modes, $E c$			
		E_u	A_g	E_u	E_u	A_u	A_u	A_u	A_u
CaMoO ₄	312, 313, 314	146 (161)	197 (258)	322 (359)	790 (910)	193 (202)	247 (317)	420 (450)	772 (898)
CaWO ₄	312, 313, 314	142 (153)	200 (248)	313 (364)	786 (906)	177 (181)	237 (323)	420 (450)	776 (896)
SrMoO ₄	313	125	[181]	[327]	[830]	153	[282]	[404]	[830]
SrWO ₄	313	[140]	[168]	[320]	[833]	[150]	[278]	[410]	[833]
PbMoO ₄	313, 314	90 (99)	105 (160)	301 (318)	744 (886)	86 (132)	258 (278)	373 (387)	745 (865)
PbWO ₄	313, 314, 315	73 (101)	104 (137)	288 (314)	756 (869)	58 (109)	251 (278)	384 (393)	764 (866)
LiYF ₄	316, 317	143 (173)	292 (303)	326 (367)	424 (566)	195 (224)	252 (283)	396	490

Material	Ref.	Raman modes												
		A_g	A_g	A_g	B_g	B_g	E_g							
CaMoO ₄	318	205	333	878	110	219	339	393	844	145	189	263	401	797
CaWO ₄	318, 319	218	336	912	84	210	336	401	838	117	195	275	409	797
SrMoO ₄	318	181	327	887	94	157	327	367	842	111	137	231	381	797
SrWO ₄	318, 319	187	334	925	75		334	370	839	101	131	238	378	797
PbMoO ₄	320	164	314	868	64	75	317	348	764	61	100	190	356	744
PbWO ₄	320	178	328	905	54	78	328	358	766	63	78	192	358	753
LiYF ₄	316, 317	[150]	264	425	177	248	329	382	427	153	199	329	368	446

TABLE 36 Optical Modes of Crystals with Spinel Structure; Space Group: $Fd3m (O_h^7) \#227$; $\Gamma = A_{1g}(R) + E_g(R) + F_{1g}(-) + 3F_{2g}(R) + 2A_{2u}(-) + 2E_u(-) + 4F_{1u}(IR) + 2F_{2u}(-)$

Material	Ref.	Mode locations (cm^{-1})								
		Infrared modes				Raman modes				
		F_{1u}	F_{1u}	F_{1u}	F_{1u}	F_{2g}	F_{2g}	F_{2g}	E_g	A_{1g}
MgAl_2O_4	226, 321, 322	305 (311)	428	485 (497)	670 (800)	311	492	611	410	722
CdIn_2S_4	323, 324	68 (69)	171 (172)	215 (270)	307 (311)	93	247	312	185	366
ZnCr_2O_4	325	186 (194)	372 (377)	506 (522)	624 (711)	186	515	610	457	692
ZnCr_2S_4	326, 327, 328	115 (117)	249 (250)	340 (360)	388 (403)	116	[290]	361	249	403

TABLE 37 Optical Modes of Crystals with Cubic Perovskite Structure; Space Group: $\text{Pm}\bar{3}\text{m} (O_h^1) \#221$; $\Gamma = 3F_{1u}(IR) + F_{2u}(-)$

Material	Infrared mode locations (cm^{-1})			Ref.
	F_{1u}^*	F_{1u}	F_{1u}	
KTaO_3	85 (88)	199 (200)	549 (550)	329
SrTiO_3	88 (173)	178 (473)	544 (804)	329, 330, 331
KMgF_3	168 (197)	299 (362)	458 (551)	332, 333
KMnF_3	119 (144)	193 (270)	399 (483)	332, 333

*"Soft" mode with strong temperature dependence.

TABLE 38 Optical Modes of Crystals with Tetragonal Perovskite Structure; Space Group: $\text{P}4_2/\text{mnm} (D_{4h}^{14}) \#136$; $\Gamma = 3A_1(IR, E \parallel c, R) + B_1(R) + 4E(IR, E \perp c, R)$

Material	Mode locations (cm^{-1})								
	Infrared ($E \perp c$)				Raman	Infrared ($E \parallel c$)			
	E	E	E	E	B_1	A_1	A_1	A_1	Ref.
BaTiO_3	34 (180)	181 (306)	306 (465)	482 (706)	305	180 (187)	280 (469)	507 (729)	329, 333, 334
PbTiO_3	89 (128)	221	250 (445)	508 (717)	415	127 (215)	351 (445)	613 (794)	335

TABLE 39 Optical Modes of Crystals with the Chalcopyrite Structure; Space Group: $I\bar{4}2d$ (D_{2d}^{12}) #122; $\Gamma = A_1(R) + 2A_2(-) + 3B_1(R) + 3B_2(IR, E \parallel c, R) + 6E(IR, E \perp c, R)$

Material	Ref.	Mode locations (cm^{-1})						
		Raman modes				Infrared modes ($E \parallel c$)		
		A_1	B_1	B_1	B_1	B_2	B_2	B_2
AgGaS ₂	336 337	295	118	179	334	195 (199)	215 (239)	366 (400)
AgGaSe ₂	113	181	56	160	253	58 (58)	155 (161)	252 (275)
CdGeAs ₂	338	[185]	[105]	[169]	[242]	[108] [(109)]	203 (210)	270 (278)
CuGaS ₂	336	312	138	203	243	259 (284)	339 (369)	371 (402)
ZnGeP ₂	339 340	328	120	247	389	[140]	361 (341)	411 (401)
ZnSiP ₂	341	344	131				352 (362)	511 (535)

Material	Ref.	Infrared mode location ($E \perp c$)					
		E	E	E	E	E	E
AgGaS ₂	336, 337	63 (64)	93 (96)	158 (160)	225 (230)	323 (347)	368 (392)
AgGaSe ₂	336, 342	78 (78)	133 (135)	(112)	160 (163)	208 (213)	247 (274)
CdGeAs ₂	338, 343, 344	95 (95)	[114] [(114)]	159 (161)	200 (206)	255 (258)	272 (280)
CuGaS ₂	336	75 (76)	95 (98)	147 (167)	260 (278)	335 (352)	365 (387)
ZnGeP ₂	339, 340	94 (94)	141 (141)	201 (206)	328 (330)	369 (375)	386 (406)
ZnSiP ₂	341	105 (105)	185 (185)	270 (270)	335 (362)	477 (477)	511 (535)

TABLE 40 Optical Modes of Other Crystals

Material/Space Group	Irreducible Optical Representation and Optical Modes Locations (cm^{-1})
Orpiment, As ₂ S ₃ P2 ₁ /b(C_{2h}^5) #14 (Ref. 345)	$\Gamma = 15A_g(R) + 15B_g(R) + 14A_u(IR, E \parallel b) + 13B_u(IR, E \parallel a, E \parallel c)$ [$7A_1(IR, E \parallel c, R) + 7A_2(R) + 7B_1(IR, E \parallel b, R) + 6B_2(IR, E \parallel a, R)$ for a noninteracting molecular layer structure] $A_g = 136, 154, 204, 311, 355, 382$ $B_g = 140, 159, 198, 278, 311, 354, 383$
Calcite, CaCO ₃ $R\bar{3}c(D_{3d}^6)$ #167 (Refs. 346, 347)	$\Gamma = A_g(R) + 2A_{1u}(-) + 3A_{2u}(IR, E \parallel c) + 4E_g(R) + 5E_u(IR, E \perp c)$ $A_g = 1088$ $E_g = 156, 283, 714, 1432$ $A_{2u} = 92(136), 303(387), 872(890)$ $E_u = 102(123), 223(239), 297(381), 712(715), 1407(1549)$

TABLE 40 Optical Modes of Other Crystals (*Continued*)

Material/Space Group	Irreducible Optical Representation and Optical Modes Locations (cm ⁻¹)
BBO, Ba ₃ [B ₃ O ₆] ₂ R3(C ₃ ²) #146 (Ref. 348)	$\Gamma = 41A(\text{IR}, E \parallel c, R) + 41E(\text{IE}, E \perp c, R)$ (See reference)
BSO, selenite, Bi ₁₂ SiO ₂₀ I23 (T ³) #197 (Ref. 349)	$\Gamma = 8A(\text{R}) + 8E(\text{R}) + 24F(\text{IR})$ A = 92, 149, 171, 282, 331, 546, 785 E = 68, 88, 132, 252, —, 464, 626 F = 44, 51, 59, 89, 99, 106, —, 115, 136, 175, 195, 208, 237, 288, 314, 353, 367, 462, 496, 509, 531, 579, 609, 825
BIBO, BiB ₃ O ₆ C2(C ₂ ²) #5 (Refs. 350, 351)	$\Gamma = 13A + 14B$ (See references)
Iron pyrite, FeS ₂ Pa3 (T ⁶) #205 (Refs. 352, 353)	$\Gamma = A_g(\text{R}) + E_g(\text{R}) + 3F_g(\text{R}) + 2A_u(-) + 2E_u(-) + 5F_u(\text{IR})$ A _g = 379 E _g = 343 F _g = 435, 350, 377 F _u = 293(294), 348(350), 401(411), 412(421), 422(349)
KDP, KH ₂ PO ₄ I42d(D _{2d} ¹²) #122 (Ref. 354)	$\Gamma = 4A_1(\text{R}) + 5S_2(-) + 6B_1(\text{R}) + 6B_2(\text{IR}, E \parallel c, R) + 12E(\text{IR}, E \perp c, R)$ A ₁ = 360, 514, 918, 2700 B ₁ = 156, 479, 570, 1366, 1806, 2390 B ₂ = 80, 174, —, 386, 510, 1350 E = 75, 95, 113, 190, 320, 490, 530, 568, 960, 1145, —, 1325
KTP Pna2 ₁ (C _{2v} ⁹) #33 (Refs. 355, 356)	$\Gamma = 47A_1(\text{IR}, E \parallel c, R) + 48A_2(\text{R}) + 47B_1(\text{IR}, E \parallel a, R) + 47B_2(\text{IR}, E \parallel b, R)$
Lanthanum fluoride, LaF ₃ P3c1(D _{3d} ⁴) #165 (Refs. 357, 358, 359)	$\Gamma = 5A_{1g}(\text{R}) + 12E_g(\text{R}) + 6A_{2u}(\text{IR}, E \parallel c)$ A _{1g} = 120, 231, 283, 305, 390 E _g = 79, 145, 145, 163, 203, 226, 281, 290, 301, 315, 325, 366 A _{2u} = 142(143), 168(176), 194(239), —, 275(296), 323(468) E _u = 100(108), 128(130), 144(145), 168(183), 193(195), 208(222), 245(268), 272(316), 354(364), 356(457)
LBO LiB ₃ O ₅ Pna2 ₁ (C _{2v} ⁹) #33 (Ref. 360)	$\Gamma = 26A_1(\text{IR}, E \parallel z, R) + 27A_2(\text{R}) + 26B_1(\text{IR}, E \parallel x, R) + 26B_2(\text{IR}, E \parallel y, R)$ (See reference)
Lithium iodate, LiIO ₃ P6 ₃ (C ₆ ⁶) #173 (Refs. 361, 362)	$\Gamma = 4A(\text{IR}, E \parallel c, R) + 5B(-) + 4E_1(\text{IR}, E \perp c, R) + 5E_2(\text{R})$ A = 148(148), 238(238), 358(468), 795(817) E ₁ = 180(180), 330(340), 370(460), 764(848) E ₂ = 98, 200, 332, 347, 765
Lithium niobate, LiNbO ₃ R3c(C _{3v} ⁶) #161 (Refs. 363, 364)	$\Gamma = 4A_1(\text{IR}, E \parallel c, R) + 5A_2(+9E(\text{IR}, E \perp c, R))$ A ₁ = 255(275), 276(333), 334(436), 633(876) E = 155(198), 238(243), 265(295), 325(371), 371(428), 431(454), 582(668), 668(739), 743(880)
Potassium niobate, KNbO ₃ Bmm2(C _{2v} ¹⁴) #38 (Ref. 365)	$\Gamma = 4A_1(\text{IR}, E \parallel z, R) + 4B_1(\text{IR}, E \parallel x, R) + 3B_2(\text{IR}, E \parallel y, R) + A_2(\text{R})$ A ₁ = 190(193), 290(296), 299(417), 607(827) B ₁ = 187(190), 243(294), 267(413), 534(842) B ₂ = 56(189), 195(425), 511(838) A ₂ = 283
Tl ₃ AsSe ₅ , TAS R3m(C _{3v} ⁵) #160 (Ref. 366)	$\Gamma = 4A_1(\text{IR}, E \parallel c, R) + 4A_2(-) + 6E(\text{IR}, E \perp c, R)$ A ₁ = 62(66), 93(97), 119(132), 239(247) E = 64.1(67.5), 78.1(82), 92.6(96), 125(133), 240(250)

(Continued)

TABLE 40 Optical Modes of Other Crystals (*Continued*)

Material/Space Group	Irreducible Optical Representation and Optical Modes Locations (cm ⁻¹)
Paratellurite, TeO ₂ P4 ₁ 2 ₁ 2 (D ₄ ^h) #92 (Refs. 367, 368)	$\Gamma = 4A_1(\text{R}) + 4A_2(\text{IR}, E \parallel c) + 5B_1(\text{R}) + 4B_2(\text{R}) + 8E(\text{IR}, E \perp c, \text{R})$ A ₁ = 148, 393, 648 A ₂ = 82(110), 259(263), 315(375), 575(775) B ₁ = 62, 175, 216, 233, 591 B ₂ = 155, 287, 414, 784 E = 121(123), 174(197), 210(237), 297(327), 330(379), 379(415), 643(720), 769(812)
Yttria, Y ₂ O ₃ Ia3 (T _h ^h) #206 (Ref. 369)	$\Gamma = 4E_g(\text{R}) + 4A_g(\text{R}) + 14F_g(\text{R}) + 5E_u(-) + 5A_u(-) + 16F_u(\text{IR})$ E _g = 333, 830, 948 A _g = 1184 F _g = (131), 431, 469, 596 F _u = 120(121), 172(173), 182(183), 241(242), 303(315), 335(359), 371(412), 415(456), 461(486), 490(535), 555(620)
Yttrium Vanadate, YVO ₄ I4 ₁ /amd (D _{4h} ¹⁹) #141 (Refs. 370, 371)	$\Gamma = 2A_{1g}(\text{R}) + 5E_g(\text{R}) + 4B_{1g}(\text{R}) + 1B_{2g}(-) + 3A_{2u}(\text{IR} \parallel c) + 4E_u(\text{IR} \perp c) + 1A_{2g}(-) + 1A_{1u}(-) + 1B_{1u}(-) + 2B_{2u}(-)$ E _u = 195(222), 263(309), 309(311), 780 (930) A _{2u} = 310, 455, 803
Yttrium aluminum garnet (YAG), Y ₃ Al ₂ (AlO ₄) ₃ Ia3d (O _h ¹⁰) #230 (Refs. 372, 373)	$\Gamma = 5A_{1u}(-) + 3A_{1g}(-) + 5A_{2g}(-) + 10E_u(-) + 8E_g(\text{R}) + 14F_{1g}(-) + 17F_{1u}(\text{IR}) + 14F_{2g}(\text{R}) + 16F_{2u}(-)$ A _{1g} = 373, 561, 783 E _g = 162, 310, 340, 403, 531, 537, 714, 758 F _{2g} = 144, 218, 243, 259, 296, 408, 436, 544, 690, 719, 857 F _{1u} = 122(123), 163(172), 177(180), 219(224), 290(296), 330(340), 373(378), 387(388), 395(403), 428(438), 446(472), 472(511), 516(549), 569(585), 692(712), 723(765), 782(841)

TABLE 41 Summary of Available Lattice Vibration Model Parameters

Material	Dispersion model reference		Material	Dispersion model reference	
	Classical	Four-Parameter		Classical	Four-Parameter
AgBr		374	KNbO ₃	365	
AgCl		374	KTaO ₃	329	
AgGaS ₂	336		KTiOPO ₄	355	
AgGaSe ₂	342		LaF ₃		358
Al ₂ O ₃	35	36	La ₂ O ₃	384	
ALON	205		LiF	251	
As ₂ S ₃ (cryst)	345		LiIO ₃	361, 362	
As ₂ S ₃ (glass)	375		LiNbO ₃	385	
As ₂ S ₃ (cryst)	345		YLiF ₄	316	
As ₂ Se ₃ (glass)	375		MgAl ₂ O ₄	226	
BaF ₂	269		MgF ₂	305, 307	302
BaTiO ₃	329, 376	334	MgO	251	
BeO	284		NaF	386	
BN	257		PbF ₂	272	

TABLE 41 Summary of Available Lattice Vibration Model Parameters (*Continued*)

Material	Dispersion model reference		Material	Dispersion model reference	
	Classical	Four-Parameter		Classical	Four-Parameter
CaCO ₃	347		PbSe	387	
CaF ₂	269		PbWO ₄	315	
CaMoO ₄	312		Se	291, 377	
CaWO ₄	312		SiO ₂	296, 388	
CdGeP ₂	343		SrF ₂	269	
CdS	259		SrTiO ₃	329	331
CdSe	287		Te	293	
CdTe	259, 377		TeO ₂	368, 389	389
CsBr	243		TiO ₂	372	32
CsCl	243, 244		TlBr		374
CsI	242		TlCl		374
FeS ₂	378, 379		Y ₃ Al ₅ O ₁₂	390	373
GaAs	380		Y ₂ O ₃	369	
GaN	217, 287		YVO ₄	371	370
GaP	260, 381		ZnGeP ₂	339	
GeO ₂	302		ZnS	259	
HfO ₂ :Y ₂ O ₃	167		ZnSe	259	
KBr	382		ZnTe	259	
KI	383				

2.6 REFERENCES

1. W. J. Tropf, M. E. Thomas, and T. J. Harris, "Properties of Crystals and Glasses," Chap. 33, in M. Bass (ed.), *Handbook of Optics, Second Edition, Vol. 2: Devices, Measurements, and Properties*, McGraw-Hill, New York, 1994.
2. D. E. Gray (ed.), *American Institute of Physics Handbook*, 3d ed., McGraw-Hill, New York, 1972.
3. W. L. Wolfe and G. J. Zissis (eds.), *The Infrared Handbook*, Environmental Research Institute of Michigan, 1985.
4. M. J. Weber (ed.), *Handbook of Laser Science and Technology*, CRC Press, Boca Raton, 1986.
5. E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, Academic Press, Orlando, 1985.
6. E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991.
7. E. D. Palik (ed.), *Handbook of Optical Constants of Solids III*, Academic Press, Orlando, 1997.
8. David N. Nikogosyan, *Nonlinear Optical Crystals: A Complete Survey*, Springer, New York, 2005.
9. N. P. Bansal and R. H. Doremus, *Handbook of Glass Properties*, Academic Press, Orlando, 1986.
10. H. Bach and N. Neuroth, *The Properties of Optical Glass*, Springer-Verlag, Berlin, 1995.
11. B. O. Seraphin and H. E. Bennett, "Optical Constants," in R. K. Willardson and A. C. Beer (eds.), *Semiconductors and Semimetals, Vol. 3: Optical Properties of III-V Compounds*, Academic Press, New York, 1967.
12. F. V. Tooley (ed.), *The Handbook of Glass Manufacture*, Ashlee Publishing Co., New York, 1985.
13. Schott North America, Advanced Optic Home. Contains the optical glass catalog and technical information. http://www.us.schott.com/advanced_optics/english/index.html, accessed Feb 26, 2008.

14. Hoya Corporation USA. <http://www.hoyaoptics.com>, accessed May 21, 2008.
15. Ohara Inc., Products and Services. <http://www.ohara-inc.co.jp/en/product/optical/opticalglass/01000.html>, accessed Sep 11, 2008.
16. D. Krause, "Glasses," in W. Martienssen and H. Warlimont (eds.), *Springer Handbook of Condensed Matter and Materials Data*, Springer, Berlin and Heidelberg, 2005.
17. M. Brinkmann, J. Hayden, M. Letz, et al., "Optical Materials and Their Properties," in F. Trager (ed.), *Springer Handbook of Lasers and Optics*, Springer, New York, 2007.
18. S. Musikant, *Optical Materials: An Introduction to Selection and Application*, Marcel Dekker, New York, 1985.
19. P. Klocek (ed.), *Handbook of Infrared Optical Materials*, Marcel Dekker, New York, 1991.
20. G. W. Fynn and W. J. A. Powell, *Cutting and Polishing Optical and Electronic Materials*, 2d ed., Adam Hilger, Bristol, 1988.
21. W. Zschommler, *Precision Optical Glassworking*, SPIE, Bellingham, WA, 1986.
22. M. E. Thomas, *Optical Propagation in Linear Media: Atmospheric Gases and Particles, Solid-State Components, and Water*, Oxford University Press, USA, 2006.
23. D. C. Harris, *Materials for Infrared Windows and Domes: Properties and Performance*, SPIE Press, Bellingham WA, 1999.
24. J. F. Nye, *Physical Properties of Crystals*, Oxford University Press, Oxford, 1985.
25. P. J. Mohr, B. N. Taylor, and D. B. Newell, "CODATA Recommended Values for the Fundamental Physical Constants: 2006," *J. Phys. Chem. Ref. Data* **37**:1187–1284 (2006).
26. D. Y. Smith, "Dispersion Theory, Sum Rules, and Their Application to the Analysis of Optical Data," in E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, Academic Press, New York, 1985.
27. T. Skettrup, "Urbach's Rule Derived from Thermal Fluctuations in the Band Gap Energy," *Phys. Rev. B* **18**:2622–2631 (1978).
28. F. K. Kneubühl, "Review of the Theory of the Dielectric Dispersion of Insulators," *Infrared Phys.* **29**:925–942 (1989).
29. W. J. Tropf, T. J. Harris, and M. E. Thomas, "Optical Materials: Visible and Infrared," in R. Waynant and W. Ediger (eds.), *Electro-optics Handbook*, McGraw-Hill, New York, 1993.
30. R. H. Lyddane, R. G. Sachs, and E. Teller, "On the Polar Vibrations of Alkali Halides," *Phys. Rev.* **59**:673–676 (1941).
31. W. Cochran and R. A. Cowley, "Dielectric Constants and Lattice Vibrations," *J. Phys. Chem. Solids* **23**:447–450 (1962).
32. F. Gervais and B. Piriou, "Temperature Dependence of Transverse- and Longitudinal-optic Modes in TiO₂ (Rutile)," *Phys. Rev. B* **10**:1642–1654 (1974).
33. J. W. S. Rayleigh, "The Theory of Anomalous Dispersion," *Phil. Mag.* **48**:151–152 (1889).
34. L. Merten and G. Lamprecht, "Directional Dependence of Extraordinary Infrared Oscillator Parameters of Uniaxial Crystals," *Phys. Stat. Sol. (b)* **39**:573–580 (1970). [Also see J. L. Servoin, F. Gervais, A. M. Quittet, and Y. Luspain, "Infrared and Raman Response in Ferroelectric Perovskite Crystals: Apparent Inconsistencies," *Phys. Rev. B* **21**:2038–2041 (1980).]
35. A. S. Barker, "Infrared Lattice Vibrations and Dielectric Dispersion in Corundum," *Phys. Rev.* **132**:1474–1481 (1963).
36. F. Gervais and B. Piriou, "Anharmonicity in Several-Polar-Mode Crystals: Adjusting Phonon Self-Energy of LO and TO Modes in Al₂O₃ and TiO₂ to Fit Infrared Reflectivity," *J. Phys. C* **7**:2374–2386 (1974).
37. R. H. French, "Electronic Band Structure of Al₂O₃ with Comparison to ALON and AlN," *J. Am. Ceram. Soc.* **73**:477–489 (1990).
38. W. Sellmeier, "Zur Erklärung der abnormen Farbenfolge im Spectrum einiger Substanzen," *Ann. Phys. Chem.* **219**(Series 2, **143**):272–282 (1871). [Also see W. Sellmeier, "II. Ueber die durch Aetherschwingungen erregten Mitschwingungen der Körpertheilchen und deren Rückwirkung auf die ersteren, besonders zur Erklärung der Dispersion und ihrer Anomalien," *Ann. Phys. Chem.* **221**:520–549 (1872).]
39. F. Urbach, "The Long-Wavelength Edge of Photographic Sensitivity and the Electronic Absorption of Solids," *Phys. Rev.* **92**:1324 (1953).
40. H. Sumi and A. Sumi, "The Urbach–Martienssen Rule Revisited," *J. Phys. Soc. Japan* **56**:2211–2220 (1987).
41. D. L. Wood and J. Tauc, "Weak Absorption Tails in Amorphous Semiconductors," *Phys. Rev. B* **5**:3144–3151 (1972).

42. H. Mori and T. Izawa, "A New Loss Mechanism in Ultra-Low Loss Optical Fiber Materials," *J. Appl. Phys.* **51**:2270–2271 (1980).
43. M. E. Innocenzi, R. T. Swimm, M. Bass, R. H. French, A. B. Villaverde, and M. R. Kokta, "Room-Temperature Optical Absorption in Undoped α -Al₂O₃," *J. Appl. Phys.* **67**:7542–7546 (1990).
44. S. Adachi, "Model Dielectric Constants of GaP, GaAs, GaSb, InSb, InAs, and InSb," *Phys. Rev. B* **35**:7454–7463 (1987).
45. S. Adachi, "Optical Properties of Al_xGa_{1-x}As Alloys," *Phys. Rev. B* **38**:12345–12352 (1988).
46. S. Adachi, "Optical Properties of In_{1-x}Ga_xAs_yP_{1-y} Alloys," *Phys. Rev. B* **39**:12612–12621 (1989).
47. S. Adachi, "Excitonic Effects in the Optical Spectrum of GaAs," *Phys. Rev. B* **41**:1003–1013 (1990).
48. S. Adachi, "Effects of the Indirect Transitions on Optical Dispersion Relations," *Phys. Rev. B* **41**:3504–3508 (1990).
49. S. Adachi, *Physical Properties of III-V Semiconductor Compounds: InP, InAs, GaAs, GaP, InGaAs, and InGaAsP*, Wiley-Interscience, New York, 1992.
50. S. Adachi, T. Kimura, and N. Suzuki, "Optical Properties of CdTe: Experiment and Modeling," *J. Appl. Phys.* **74**:3435–3441 (1993).
51. K. Sato and S. Adachi, "Optical Properties of ZnTe," *J. Appl. Phys.* **73**:926–931 (1993).
52. D. T. Gillespie, A. L. Olsen, and L. W. Nichols, "Transmittance of Optical Materials at High Temperature," *Appl. Opt.* **4**:1488–1493 (1965).
53. M. A. Ordal, R. J. Bell, R. W. Alexander, L. L. Long, and M. R. Querry, "Optical Properties of Fourteen Metals in the Infrared and Far Infrared: Al, Co, Cu, Au, Fe, Pb, Mo, Ni, Pd, Pt, Ag, Ti, V, and W," *Appl. Opt.* **24**:4493–4499 (1985).
54. T. F. Deutsch, "Absorption Coefficient of Infrared Laser Window Materials," *J. Phys. Chem. Solids* **34**:2091–2104 (1973). [Also see T. F. Deutsch, "Laser Window Materials—An Overview," *J. Electron. Mater.* **4**:663–719 (1975).]
55. D. L. Mills and A. A. Maradudin, "Theory of Infrared Absorption by Crystals in the High Frequency Wing of Their Fundamental Lattice Absorption," *Phys. Rev. B* **8**:1617–1630 (1973). [Also see A. A. Maradudin and D. L. Mills, "Temperature Dependence of the Absorption Coefficient of Alkali Halides in the Multiphonon Regime," *Phys. Rev. Lett.* **31**:718–721 (1973).]
56. H. G. Lipson, B. Bendow, N. E. Massa, and S. S. Mitra, "Multiphonon Infrared Absorption in the Transparent Regime of Alkaline-Earth Fluorides," *Phys. Rev. B* **13**:2614–2619 (1976).
57. M. E. Thomas, R. I. Joseph, and W. J. Tropf, "Infrared Transmission Properties of Sapphire, Spinel, Ytria, and ALON as a Function of Frequency and Temperature," *Appl. Opt.* **27**: 239–245 (1988).
58. M. E. Thomas, W. J. Tropf, and A. Szpak, "Optical Properties of Diamond," *Diamond and Films and Technology* **5**:159–180 (1995).
59. D. Yanga, M. E. Thomas, S. Andersson, and Bob Podgurski, "Infrared Optical Properties of Eight Different Schott Glasses," *Proc. SPIE* **4102**:134–143 (2000).
60. D. V. Hahn, M. E. Thomas, and D. W. Blodgett, "Modeling of the Frequency- and Temperature-Dependent Absorption Coefficient of Long-Wave-Infrared (2–25 μ m) Transmitting Materials," *Appl. Opt.* **44**:6913–6920 (2005).
61. J. A. Harrington, B. L. Bobbs, M. Braunstein, R. K. Kim, R. Stearns, and R. Braunstein, "Ultraviolet-Visible Absorption in Highly Transparent Solids by Laser Calorimetry and Wavelength Modulation Spectroscopy," *Appl. Opt.* **17**:1541–1546 (1978).
62. N. C. Frenelius, R. J. Harris, D. B. O'Quinn, M. E. Gangl, D. V. Dempsey, and W. L. Knecht, "Some Optical Properties of Materials Measured at 1.3 μ m," *Opt. Eng.* **22**:411–418 (1983).
63. J. A. Harrington, D. A. Gregory, and W. F. Ott, "Infrared Absorption in Chemical Laser Window Materials," *Appl. Opt.* **15**:1953–1959 (1976).
64. J. W. Fleming, "Dispersion in GeO₂-SiO₂ Glasses," *Appl. Opt.* **23**:4486–4493 (1984).
65. D. L. Wood, K. Nassau, and T. Y. Kometani, "Refractive Index of Y₂O₃ Stabilized Zirconia: Variation with Composition and Wavelength," *Appl. Opt.* **29**:2485–2488 (1990).

66. N. P. Barnes and M. S. Piltch, "Temperature-Dependent Sellmeier Coefficients and Nonlinear Optics Average Power Limit for Germanium," *J. Opt. Soc. Am.* **69**:178–180 (1979).
67. B. Tatian, "Fitting Refractive-Index Data with the Sellmeier Dispersion Formula," *Appl. Opt.* **23**:4477–4485 (1984).
68. P. G. Nutting, "Dispersion Formulas Applicable to Glass," *J. Opt. Soc. Am.* **2–3**:61–65 (1919).
69. H. H. Li, "Refractive Index of ZnS, ZnSe, and ZnTe and Its Wavelength and Temperature Derivatives," *J. Phys. Chem. Ref. Data* **13**:103–150 (1984).
70. F. Zernike, "Refractive Indices of Ammonium Dihydrogen Phosphate and Potassium Dihydrogen Phosphate between 2000 Å and 1.5, μ ," *J. Opt. Soc. Am.* **54**:1215–1220 (1964). [Errata: *J. Opt. Soc. Am.* **55**:210E (1965)]
71. M. Herzberger, "Colour Correction in Optical Systems and a New Dispersion Formula," *Opt. Acta* **6**:197–215 (1959).
72. M. Herzberger and C. D. Salzberg, "Refractive Indices of Infrared Optical Materials and Color Correction of Infrared Lenses," *J. Opt. Soc. Am.* **52**:420–427 (1962).
73. A. N. Pikhtin and A. D. Yas'kov, "Dispersion of the Refractive Index of Semiconductors with Diamond and Zinc-Blende Structures," *Sov. Phys. Semicond.* **12**:622–626 (1978).
74. J. Riishede, J. Lægsgaard, J. Broeng, and A. Bjarklev, "All-Silica Photonic Bandgap Fibre with Zero Dispersion and Large Mode Area at 730 nm," *J. Optics. A* **6**:667–670 (2004).
75. J. Riishede, N.A. Mortensen, and J. Lægsgaard, "A 'Poor Man's Approach' to Modeling Micro-Structured Optical Fibers," *J. Optics A* **5**:534–538 (2003).
76. J. Lægsgaard, P. J. Roberts, and M. Bache, "Tailoring the Dispersion Properties of Photonic Crystal Fibers," *Opt. Quant. Electron.* **39**:995–1008 (2007).
77. N. A. Mortensen, "Photonic Crystal Fibres: Mapping Maxwell's Equations onto a Schrödinger Equation Eigenvalue Problem," *J. Europ. Opt. Soc. Rapid Public.* **1**:06009 (2006).
78. A. J. Bosman and E. E. Havinga, "Temperature Dependence of Dielectric Constants of Cubic Ionic Compounds," *Phys. Rev.* **129**:1593–1600 (1963).
79. B. Bendow, P. D. Gianino, Y-F. Tsay, and S. S. Mitra, "Pressure and Stress Dependence of the Refractive Index of Transparent Crystals," *Appl. Opt.* **13**:2382–2396 (1974).
80. S. Singh, "Nonlinear Optical Properties," in *Handbook of Laser Science and Technology, Volume III Optical Materials: Part I*, CRC Press, Boca Raton, 1986, pp. 3–228.
81. D. F. Eaton, "Nonlinear Optical Materials," *Science* **253**:281–287 (1991).
82. M. M. Choy and R. L. Byer, "Accurate Second-Order Susceptibility Measurements of Visible and Infrared Nonlinear Crystals," *Phys. Rev. B* **14**:1693–1706 (1976).
83. A. G. Every and A. K. McCurdy, Landolt-Börnstein Numerical Data and Functional Relationships in Science and Technology, New Series, Group III: *Crystal and Solid State Physics, Volume 29A: Low Frequency Properties of Dielectric Crystals*, Springer-Verlag, Berlin, 1993.
84. M. E. Lines, "Scattering Loss in Optic Fiber Materials. I. A New Parametrization," *J. Appl. Phys.* **55**:4052–4057 (1984); "II. Numerical Estimates," *J. Appl. Phys.* **55**:4058–4063 (1984).
85. J. A. Harrington and M. Sparks, "Inverse-Square Wavelength Dependence of Attenuation in Infrared Polycrystalline Fibers," *Opt. Lett.* **8**:223–225 (1983).
86. D. D. Duncan and C. H. Lange, "Imaging Performance of Crystalline and Polycrystalline Oxides," *Proc. SPIE* **1326**:59–70 (1990).
87. J. D. H. Donnay and H. M. Ondik (eds.), *Crystal Data Determination Tables*, 3d ed., U. S. Department of Commerce, 1973.
88. R. W. G. Wyckoff, *Crystal Structures*, John Wiley & Sons, New York, 1963.
89. N. J. Kreidl, "Optical Properties," in F. V. Tooley (ed.), *Handbook of Glass Manufacture*, Books for Industry, New York, 1974, pp. 957–997.
90. H. G. Pfänder, "Geschichte des Glases," Schott Glaslexikon (mvg, Landsberg 1997) pp. 13–23.
91. M. K. Th. Clement, "The Chemical Composition of Optical Glasses and Its Influence on the Optical Properties," in H. Bach, N. Neuroth (eds.), *The Properties of Optical Glass*, Springer, Berlin, Heidelberg, 1998, pp. 58–81.

92. G. F. Brewster, N. J. Kreidl, T. G. Pett, "Lanthanum and Barium in Glass-forming System," *J. Soc. Glass Technol.* **31**:153–169 (1947).
93. W. Jahn, "Mehrstoffsysteme zum Aufbau Optischer Gläser," *Glastechn. Ber.* **43**:107–120 (1961).
94. Schott North America, "Refractive Index and Dispersion," Technical Information 29, (Schott North America Advanced Optics, Duryea, PA 18642, January 2007).
95. S. Wolff and U. Kolberg, "Environmental Friendly Optical Glasses," in H. Bach, N. Neuroth (eds.), *The Properties of Optical Glass*, Springer, Berlin, Heidelberg 1998, pp. 144–148.
96. G. Simmons and H. Wang, *Single Crystal Elastic Constants and Calculated Aggregate Properties: A Handbook*, MIT Press, Cambridge, MA, 1971.
97. R. S. Krishnan, R. Srinivasan, and S. Devanarayanan, *Thermal Expansion of Crystals*, Pergamon Press, Oxford, 1979.
98. Y. S. Touloukian (ed.), *Thermophysical Properties of Matter*, IFI/Plenum, New York, 1970.
99. G. A. Slack, "The Thermal Conductivity of Nonmetallic Crystals," in H. Ehrenreich, F. Seitz, and D. Turnbull (eds.), *Solid State Physics, Vol. 34*, Academic Press, New York, 1979.
100. K. F. Loje and D. E. Schuele, "The Pressure and Temperature Derivative of the Elastic Constants of AgBr and AgCl," *J. Phys. Chem. Solids* **31**:2051–2067 (1970).
101. W. Hildshaw, J. T. Lewis, and C. V. Briscoe, "Elastic Constants of Silver Chloride from 4.2 to 300 K," *Phys. Rev.* **163**:876–881 (1967).
102. Stiffness and compliance of ALON are estimated from the engineering moduli.
103. M. Gospodinov, P. Sveshtarov, N. Petkov, T. Milenov, V. Tassev, and A. Nikolov, "Growth of Large Crystals of Bismuth-Germanium Oxide and Their Physical Properties," *Bulgarian J. Phys.* **16**:520–522 (1989).
104. M. Gospodinov, S. Haussühl, P. Sveshtarov, V. Tassev, and N. Petkov, "Physical Properties of Cubic Bi₁₂SiO₂₀," *Bulgarian J. Phys.* **15**:140–143 (1988).
105. R. D. Greenough and S. B. Palmer, "The Elastic Constants and Thermal Expansion of Single-Crystal CdTe," *J. Phys. D* **6**:587–592 (1973).
106. A. Yoneda, "Pressure Derivatives of Elastic Constants of Single Crystal MgO and MgAl₂O₄," *J. Phys. Earth* **38**:19–55 (1990).
107. A. J. Miller, G. A. Saunders, and Y. K. Yoğurtçu, "Pressure Dependence of the Elastic Constants of PbTe, SnTe and Ge_{0.08}Sn_{0.92}Te," *J. Phys. C* **14**:1569–1584 (1981).
108. J. B. Wachtman, M. L. Wheat, and S. Marzullo, "A Method for Determining the Elastic Constants of a Cubic Crystal from Velocity Measurements in a Single Arbitrary Direction: Application to SrTiO₃," *J. Res. Nat. Bur. Stand.* **67A**:193–204 (1963).
109. Y. K. Yoğurtçu, A. J. Miller, and G. A. Saunders, "Elastic Behavior of YAG Under Pressure," *J. Phys. C* **13**:6585–6597 (1980).
110. J. W. Palko, W. M. Kriven, S. V. Sinogeikin, J. D. Bassn, and A. Sayir, "Elastic Constants of Yttria Y₂O₃ Monocrystals to High Temperatures," *J. Appl. Phys.* **89**:7791–7796 (2001).
111. H. M. Kandil, J. D. Greiner, and J. F. Smith, "Single Crystal Elastic Constants of Yttria-Stabilized Zirconia in the Range 20 to 700°C," *J. Am. Ceram. Soc.* **67**:341–346 (1984).
112. M. H. Grimsditch and G. D. Holah, "Brillouin Scattering and Elastic Moduli of Silver Thiogallate (AgGaS₂)," *Phys. Rev. B* **12**:4377–4382 (1975). [Also see N. S. Orlova, "X-Ray Diffuse Scattering Study of Anisotropy of Elastic Properties in Silver Thiogallate," *Cryst. Res. Technol.* **33**:87–99 (1998).]
113. R. Fourt, P. Derollez, A. Laamyem, B. Hennion, and J. Gonzalez, "Phonons in Silver Gallium Diselenide," *J. Phys.: Condens. Matter* **9**:6579–6589 (1997).
114. Z. Li, S. -K. Chan, M. H. Grimsditch, and E. S. Zouboulis, "The Elastic and Electromechanical Properties of Tetragonal BaTiO₃ Single Crystals," *J. Appl. Phys.* **70**:7327–7332 (1991).
115. T. Hailing, G. A. Saunders, W. A. Lambson, and R. S. Feigelson, "Elastic Behavior of the Chalcopyrite CdGeAs₂," *J. Phys. C* **15**:1399–1418 (1982).
116. I. Martynyuk-Lototska, O. Mys, O. Krupych, V. Adamiv, Ya. Burak, R. Vlokh, and W. Schranz, "Elastic, Piezooptic and Acoustooptic Properties of Borate Crystals (BaB₂O₄, Li₂B₄O₇, and CsLiB₆O₁₀)," *Integrated Ferroelectrics* **63**:99–103 (2004).

117. G. Fischer and J. Zarembowitch, "Elastic Properties of Single-Crystal Urea," *C. R. Acad. Sc. Paris* **270B**:852–855 (1970).
118. Z. Li, M. Grimsditch, C. M. Foster, and S.-K. Chan, "Dielectric and Elastic Properties of Ferroelectric Materials at Elevated Temperature," *J. Phys. Chem. Solids* **57**:1433–1438 (1996).
119. Y. Ohmachi and N. Uchida, "Temperature Dependence of Elastic, Dielectric, and Piezoelectric Constants in TeO₂ Single Crystals," *J. Appl. Phys.* **41**:2307–2311 (1970).
120. I. M. Silvestrova, Y. V. Pisarevskii, P. A. Senyushenkov, A. I. Krupny, R. Voszka, I. F. Földvári, and J. Janszky, "Temperature Dependence of the Elastic Constants of Paratellurite," *Phys. Stat. Sol. (a)* **101**:437–444 (1987).
121. W. Ruju, L. Fengying, W. Xing, and Y. Huaguang, "Ultrasonic Study on Nd:YVO₄ Crystals," *Chinese J. Lasers* **27**:449–454 (2000).
122. A. V. Kopytov and A. S. Poplavnoi, "Crystal Lattice Dynamics of ZnGeP₂ and AgGaS₂ in Hard Ion Model," *Sov. Phys. J.* **23**:353–360 (1980).
123. W. J. Alton and A. J. Barlow, "Acoustic-Wave Propagation in Tetragonal Crystals and Measurements of the Elastic Constants of Calcium Molybdate," *J. Appl. Phys.* **38**:3817–3820 (1967).
124. J. M. Farley and G. A. Saunders, "Ultrasonic Study of the Elastic Behavior of Calcium Tungstate between 1.5 K and 300 K," *J. Phys. C* **5**:3021–3037 (1972). [Also see J. M. Farley and G. A. Saunders, "The Elastic Constants of CaWO₄," *Solid State Commun.* **9**:965–969 (1971); and M. Gluyas, F. D. Hughes, and B. W. James, "The Elastic Constants of Calcium Tungstate, 4.2–300 K," *J. Phys. D: Appl. Phys.* **6**:2025–2037 (1973)]
125. J. M. Farley, G. A. Saunders, and D. Y. Chung, "Elastic Properties of Scheelite Structure Molybdates and Tungstates," *J. Phys. C* **8**:780–786 (1975).
126. J. M. Farley, G. A. Saunders, and D. Y. Chung, "Elastic Properties of Strontium Molybdate," *J. Phys. C* **6**:2010–2019 (1973).
127. P. Blanchfield and G. A. Saunders, "The Elastic Constants and Acoustic Symmetry of LiYF₄," *J. Phys. C* **12**:4673–4689 (1979). [Also see P. Blanchfield, T. Hailing, A. J. Miller, G. A. Saunders, and B. Chapman, "Vibrational Anharmonicity of Oxide and Fluoride Scheelites," *J. Phys. C* **20**:3851–3859 (1983).]
128. V. A. Savastenko and A. U. Sheleg, "Study of the Elastic Properties of Gallium Nitride," *Phys. Stat. Sol. (a)* **48**:K135–K139 (1978).
129. S. Haussühl, "The Propagation of Elastic Waves in Hexagonal Lithium Iodate," *Acustica* **23**:165–169 (1970).
130. I. I. Zubrinov, V. I. Semenov, and D. V. Shelopot, "Elastic and Photoelastic Properties of Proustite," *Sov. Phys. Solid State* **15**:1921–1922 (1974).
131. D. Eimerl, L. Davis, and S. Velsko, "Optical, Mechanical, and Thermal Properties of Barium Borate," *J. Appl. Phys.* **62**:1968–1983 (1987).
132. B. W. Woods, S. A. Payne, J. E. Marion, R. S. Hughes, and L. E. Davis, "Thermomechanical and Thermo-optical Properties of the LiCaAlF₆: Cr³⁺ Laser Material," *J. Opt. Soc. Am.* **B8**:970–977 (1991).
133. A. G. Kalinichev, J. D. Bass, C. S. Zha, P. D. Han, and D. A. Payne, "Elastic Properties of Orthorhombic KNbO₃ Single Crystals by Brillouin Scattering," *J. Appl. Phys.* **74**:6603–6608 (1993).
134. L. Glasser and C. R. A. Catlow, "Modelling Phase Changes in the Potassium Titanyl Phosphate System," *J. Mater. Chem.* **7**:2537–2542 (1997). [Also see I. I. Zubrinov, V. K. Sapozhnikov, E. V. Pestrykov, and V. V. Atuchin, "Elastic and Elastooptic Properties of KTiOPO₄," *Proc. SPIE* **5129**:249–254 (2003).]
135. R. Guo, S.A. Markgraf, Y. Furukawa, M. Sato, and A. S. Bhalla, "Pyroelectric, Dielectric, and Piezoelectric Properties of LiB₃O₅," *J. Appl. Phys.* **78**:7234–7239 (1995).
136. M. Serhane and P. Moch, "Brillouin Scattering in KTiOPO₄, RbTiOPO₄ and TlTiOPO₄: A Study of the Ferroelectric-Paraelectric Phase Transition," *J. Phys.: Condens. Matter* **6**:3821–3830 (1994).
137. S. Haussühl, L. Bohaty, and P. Becker, "Piezoelectric and Elastic Properties of the Nonlinear Optical Material Bismuth Triborate, BiB₃O₆," *Appl. Phys.* **A82**:495–502 (2006).
138. M. V. Hobden, "The Dispersion of the Refractive Indices of Proustite (Ag₃AsS₃)," *Opto-Electronics* **1**:159 (1969).
139. Y. Tsay, B. Bendow, and S. S. Mitra, "Theory of the Temperature Derivative of the Refractive Index in Transparent Crystals," *Phys. Rev. B* **5**:2688–2696 (1972).

140. L. W. Tilton, E. K. Plyler, and R. E. Stephens, "Refractive Index of Silver Chloride for Visible and Infra-Red Radiant Energy," *J. Opt. Soc. Am.* **40**:540–543 (1950).
141. G. C. Bhar, D. K. Ghosh, P. S. Ghosh, and D. Schmitt, "Temperature Effects in AgGaS₂ Nonlinear Devices," *Appl. Opt.* **22**:2492–2494 (1983).
142. E. Takaoka and K. Kato, "Thermo-Optic Dispersion Formula for AgGaS₂," *Appl. Opt.* **38**:4577–4580 (1999).
143. E. Takaoka and K. Kato, "Thermo-Optic Dispersion Formula of AgGaSe₂ and Its Practical Applications," *Appl. Opt.* **37**:561–564 (1998).
144. K. Vedam, J. L. Kirk, and B. N. N. Achar, "Piezo- and Thermo-Optic Behavior of Spinel (MgAl₂O₄)," *J. Solid State Chem.* **12**:213–218 (1975).
145. A. Feldman, D. Horowitz, R. M. Walker, and M. J. Dodge, "Optical Materials Characterization Final Technical Report, February 1, 1978–September 30, 1978," NBS Technical Note 993, February 1979.
146. J. Tapping and M. L. Reilly, "Index of Refraction of Sapphire between 24 and 1060°C for Wavelengths of 633 and 799 nm," *J. Opt. Soc. Am. A* **3**:610–616 (1986).
147. D. Yang, M. E. Thomas, and S. G. Kaplan, "Measurement of the Infrared Refractive Index of Sapphire as Function of Temperature," *Proc. SPIE* **4375**:53–63 (2001).
148. C. H. Lange and D. D. Duncan, "Temperature Coefficient of Refractive Index for Candidate Optical Windows," *SPIE Proc.* **1326**:71–78 (1990).
149. T. Yamamuro, S. Sato, T. Zenno, N. Takeyama, H. Matsuhara, I. Maeda, Y. Matsueda, "Measurement of Refractive Indices of 20 Optical Materials at Low Temperatures," *Optical Engineering* **45**(8):083401 (2006).
150. H. W. Newkirk, D. K. Smith, and J. S. Kahn, "Synthetic Bromellite. III. Some Optical Properties," *Am. Mineralogist* **51**:141–151 (1966).
151. K. Vedam and P. Hennessey, "Piezo- and Thermo-Optical Properties of Bi₁₂GeO₂₀. II. Refractive Index," *J. Opt. Soc. Am.* **65**:442–445 (1975).
152. N. Umemura, K. Miyata, and K. Kato, "New Data on the Optical Properties of BiB₃O₆," *Opt. Matl.* **30**:532–534 (2007).
153. G. N. Ramachandran, "Thermo-Optic Behavior of Solids, I. Diamond," *Proc. Ind. Acad. Sci.* **A25**:266–279 (1947).
154. J. Fontanella, R. L. Johnston, J. H. Colwell, and C. Andeen, "Temperature and Pressure Variation of the Refractive Index of Diamond," *Appl. Opt.* **16**:2949–2951 (1977).
155. I. H. Malitson, "A Redetermination of Some Optical Properties of Calcium Fluoride," *Appl. Opt.* **2**:1103–1107 (1963).
156. M. J. Dodge, "Refractive Index," in M. J. Weber (ed.), *Handbook of Laser Science and Technology, Volume IV, Optical Material: Part 2*, CRC Press, Boca Raton, 1986.
157. T. W. Houston, L. F. Johnson, P. Kisiuk, and D. J. Walsh, "Temperature Dependence of the Refractive Index of Optical Maser Crystals," *J. Opt. Soc. Am.* **53**:1286–1291 (1963).
158. R. Weil and D. Neshmit, "Temperature Coefficient of the Indices of Refraction and the Birefringence in Cadmium Sulfide," *J. Opt. Soc. Am.* **67**:190–195 (1977).
159. R. J. Harris, G. T. Johnson, G. A. Kepple, P. C. Krok, and H. Mukai, "Infrared Thermooptic Coefficient Measurement of Polycrystalline ZnSe, ZnS, CdTe, CaF₂ and BaF₂, Single Crystal KCl, and TI-20 Glass," *Appl. Opt.* **16**:436–438 (1977).
160. W. S. Rodney, "Optical Properties of Cesium Iodide," *J. Opt. Soc. Am.* **45**:987–992 (1955).
161. T. Sasaki, Y. Mori, and M. Yoshimura, "Progress in the Growth of a CsLiB₆O₁₀ Crystal and Its Application to Ultraviolet Light Generation," *Opt. Matl.* **23**:343–351 (2003). [Also see G. C. Bhar, P. Kumbhakar, and A. K. Chaudhary, "Generation of Ultraviolet Radiation with Wide Angular Tolerance in Cesium Lithium Borate Crystal," *PRAMANA J. Phys.* **55**:413–421 (2000).]
162. N. Umemura and K. Kato, "Ultraviolet Generation Tunable to 0.185 μm in CsLiB₆O₁₀," *Appl. Opt.* **36**:6794–6796 (1997).
163. G. D. Boyd, H. Kasper, and J. H. McFee, "Linear and Nonlinear Optical Properties of AgGaS₂, CuGaS₂, and CuInS₂, and Theory of the Wedge Technique for the Measurement of Nonlinear Coefficients," *IEEE J. Quantum Electr.* **7**:563–573 (1971).

164. G. C. Bhar and G. Ghosh, "Temperature-Dependent Sellmeier Coefficients and Coherence Lengths for Some Chalcopyrite Crystals," *J. Opt. Soc. Am.* **69**:730–733 (1979). [Also see G. C. Bhar, "Refractive Index Dispersion of Chalcopyrite Crystals," *J. Phys. D* **13**:455–460 (1980).]
165. M. Bertolotti, V. Bogdanov, A. Ferrari, A. Jascow, N. Nazorova, A. Pikhtin, and L. Schirone, "Temperature Dependence of the Refractive Index in Semiconductors," *J. Opt. Soc. Am. B* **7**:918–922 (1990).
166. D. A. Yas'kov and A. N. Pikhtin, "Optical Properties of Gallium Phosphide Grown by Float Zone. I. Refractive Index and Reflection Coefficient," *Mat. Res. Bull.* **4**:781–788 (1969). [Also see D. A. Yas'kov and A. N. Pikhtin, "Dispersion of the Index of Refraction of Gallium Phosphide," *Sov. Phys. Sol. State* **9**:107–110 (1967).]
167. H. H. Li, "Refractive Index of Silicon and Germanium and Its Wavelength and Temperature Derivatives," *J. Phys. Chem. Ref. Data* **9**:561–658 (1980).
168. D. L. Wood, K. Nassau, T. Y. Kometani, and D. L. Nash, "Optical Properties of Cubic Hafnia Stabilized with Yttria," *Appl. Opt.* **29**:604–607 (1990).
169. C. S. Hoefler, "Thermal Variations of the Refractive Index in Optical Materials," *Proc. SPIE* **681**:135–142 (1986).
170. B. Zysset, I. Biaggio, and P. Günter, "Refractive Indices of Orthorhombic KNbO_3 . I. Dispersion and Temperature Dependence," *J. Opt. Soc. Am. B* **9**:380–386 (1992). [Also see Y. Uematsu, "Nonlinear Optical Properties of KNbO_3 Single Crystal in the Orthorhombic Phase," *Jap. J. Appl. Phys.* **13**:1362–1368 (1974).]
171. J. D. Bierlein and H. Vanherzeele, "Potassium Titanyl Phosphate: Properties and New Applications," *J. Opt. Soc. Am. B* **6**:622–633 (1989).
172. K. Kato and E. Takaoka, "Sellmeier and Thermo-Optic Dispersion Formulas for KTP," *Appl. Opt.* **41**:5040–5044 (2002).
173. S. P. Velsko, M. Webb, L. Davis, and C. Huang, "Phase-Matched Harmonic Generation in Lithium Triborate (LBO)," *IEEE J. Quantum Electron.* **27**:2182–2192 (1991).
174. B. W. Woods, S. A. Payne, J. E. Marion, R. S. Hughes, and L. E. Davis, "Thermomechanical and Thermo-Optical Properties of the $\text{LiCaAlF}_6:\text{Cr}^{3+}$ Laser Material," *J. Opt. Soc. Am.* **B8**:970–977 (1991).
175. D. S. Smith, H. D. Riccius, and R. P. Edwin, "Refractive Indices of Lithium Niobate," *Opt. Commun.* **17**:332–335 (1976).
176. N. P. Barnes and D. J. Gettemy, "Temperature Variation of the Refractive Indices of Yttrium Lithium Fluoride," *J. Opt. Soc. Am.* **70**:1244–1247 (1980).
177. R. E. Stephens and I. H. Malitson, "Index of Refraction of Magnesium Oxide," *J. Nat. Bur. Stand.* **49**:249–252 (1952).
178. T. Radhakrishnan, "Further Studies on the Temperature Variation of the Refractive Index of Crystals," *Proc. Indian Acad. Sci.* **A33**:22–34 (1951).
179. J. N. Zemel, J. D. Jensen, and R. B. Schoolar, "Electrical and Optical Properties of Epitaxial Films of PbS, PbSe, PbTe, and SnTe," *Phys. Rev.* **140A**:330–342 (1965).
180. J. J. Carvajal, P. Segonds, A. Peña, J. Zaccaro, B. Boulanger, F. Diaz, and M. Aguil, "Structural and Optical Properties of $\text{RbTiOPO}_4:\text{Nb}$ Crystals," *J. Phys.: Condens. Matter* **19**:116214 (2007).
181. N. Uchida, "Optical Properties of Single Crystal Paratellurite (TeO_2)," *Phys. Rev. B* **4**:3736–3745 (1971).
182. J. R. DeVore, "Refractive Index of Rutile and Sphalerite," *J. Opt. Soc. Am.* **41**:416–419 (1951).
183. W. S. Rodney and I. H. Malitson, "Refraction and Dispersion of Thallium Bromide Iodide," *J. Opt. Soc. Am.* **46**:956–961 (1956).
184. L. G. DeShazer, S. C. Rand, and B. A. Wechsler, "Laser Crystals," in M. J. Weber (ed.), *Handbook of Laser Science and Technology, Volume V, Optical Material: Part 3*, CRC Press, Boca Raton, 1986.
185. H. S. Shi, G. Zhang, and H. Y. Shen, "Measurement of Refractive Indices and Thermal Refractive Index Coefficients of YVO_4 Crystal," *J. Synthetic Cryst.* **30**:85–88 (2001).
186. G. D. Boyd, E. Beuhler, and F. G. Storz, "Linear and Nonlinear Optical Properties of ZnGeP_2 and CdSe ," *Appl. Phys. Lett.* **18**:301–304 (1971).
187. D. L. Wood and K. Nassau, "Refractive Index of Cubic Zirconia Stabilized with Yttria," *Appl. Opt.* **21**:2978–2981 (1982).

188. Schott Glass Technologies, Duryea, Pa.
189. Hoya Optics, Inc., Fremont, Calif.
190. W. S. Rodney and R. J. Spindler, "Index of Refraction of Fused-Quartz Glass for Ultraviolet, Visible, and Infrared Wavelengths," *J. Res. Nat. Bur. Stand.* **53**:185–189 (1954). [Also see I. H. Malitson, "Interspecimen Comparison of the Refractive Index of Fused Silica," *J. Opt. Soc. Am.* **55**:1205–1209 (1965).]
191. J. W. Fleming, "Dispersion in GeO_2 - SiO_2 Glasses," *Appl. Opt.* **23**:4486–4493 (1984).
192. Barr & Stroud, Ltd., Glasgow, Scotland (UK).
193. Corning, Inc., Corning, N. Y. [Also see W. H. Dumbaugh, "Infrared Transmitting Germanate Glasses," *Proc. SPIE* **297**:80–85 (1981).]
194. Ohara Corporation, Somerville, N. J.
195. R. N. Brown and J. J. Hutta, "Material Dispersion in High Optical Quality Heavy Metal Fluoride Glasses," *Appl. Opt.* **24**:4500–4503 (1985).
196. J. M. Jewell, C. Askins, and I. D. Aggarwal, "Interferometric Method for Concurrent Measurement of Thermo-Optic and Thermal Expansion Coefficients," *Appl. Opt.* **30**:3656–3660 (1991).
197. W. S. Rodney, I. H. Malitson, and T. A. King, "Refractive Index of Arsenic Trisulfide," *J. Opt. Soc. Am.* **48**:633–636 (1958).
198. A. R. Hilton and C. E. Jones, "The Thermal Change in the Nondispersive Infrared Refractive Index of Optical Materials," *Appl. Opt.* **6**:1513–1517 (1967).
199. Y. Ohmachi, "Refractive Index of Vitreous As_2Se_3 ," *J. Opt. Soc. Am.* **63**:630–631 (1973).
200. J. A. Savage, "Optical Properties of Chalcogenide Glasses," *J. Non - Cryst. Solids* **47**:101–116 (1982).
201. Amorphous Materials, Inc., Garland, Tex.
202. K. F. Hulme, O. Jones, P. H. Davies, and M. V. Hobden, "Synthetic Proustite (Ag_3AsS_3): A New Crystal for Optical Mixing," *Appl. Phys. Lett.* **10**:133–135 (1967).
203. H. Schröter, "On the Refractive Indices of Some Heavy-Metal Halides in the Visible and Calculation of Interpolation Formulas for Dispersion," *Z. Phys.* **67**:24–36 (1931) [in German].
204. I. H. Malitson and M. J. Dodge, "Refractive Index and Birefringence of Synthetic Sapphire," *J. Opt. Soc. Am.* **62**:1405A (1972). [Also see M. J. Dodge, "Refractive Index," in *Handbook of Laser Science and Technology, Volume IV, Optical Materials: Part 2*, CRC Press, Boca Raton, 1986, p. 30.]
205. W. J. Tropf and M. E. Thomas, "Aluminum Oxynitride (ALON) Spinel," in E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991, pp. 775–785.
206. I. H. Malitson, "Refractive Properties of Barium Fluoride," *J. Opt. Soc. Am.* **54**:628–632 (1964).
207. S. H. Wemple, M. Didomenico, and I. Camlibel, "Dielectric and Optical Properties of Melt-Grown BaTiO_3 ," *J. Phys. Chem. Solids* **29**:1797–1803 (1968).
208. D. F. Edwards and R. H. White, "Beryllium Oxide," in E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991, pp. 805–814.
209. M. Simon, F. Mersch, C. Kuper, S. Mendricks, S. Wevering, J. Imbrock, and E. Krätzig, "Refractive Indices of Photorefractive Bismuth Titanate, Barium-Calcium Titanate, Bismuth Germanium Oxide, and Lead Germinate," *Phys. Stat. Sol. (a)* **159**:559–562 (1997). [Also see E. Burattini, G. Cappuccio, M. Grandolfo, P. Vecchia, and Sh. M. Efendiev, "Near-Infrared Refractive Index of Bismuth Germanium Oxide ($\text{Bi}_{12}\text{GeO}_{20}$)," *J. Opt. Soc. Am.* **73**:495–497 (1983).]
210. R. E. Aldrich, S. O. Hou, and M. L. Harvill, "Electrical and Optical Properties of $\text{Bi}_{12}\text{SiO}_{20}$," *J. Appl. Phys.* **42**:493–494 (1971).
211. F. Peter, "In Refractive Indices and Absorption Coefficients of Diamond between 644 and 226 Micrometers," *Z. Phys.* **15**:358–368 (1923) [In German].
212. M. Daimon and A. Masumura, "High-Accuracy Measurements of the Refractive Index and Its Temperature Coefficient of Calcium Fluoride in a Wide Wavelength Range from 138 to 2326 nm," *Appl. Opt.* **41**:5275–5281 (2002).
213. W. L. Bond, "Measurement of the Refractive Index of Several Crystals," *J. Appl. Phys.* **36**:1674–1677 (1965). [Our fit to the dispersion data.]

214. G. C. Bhar, "Refractive Index Interpolation in Phase-matching," *Appl. Opt.* **15**:305–307 (1976).
215. M. S. Gomez, J. M. Guerra, and F. Vilches, "Weighted Nonlinear Regression Analysis of a Sellmeier Expansion: Comparison of Several Nonlinear Fits of CdS Dispersion," *Appl. Opt.* **24**:1147–1150 (1985).
216. A. G. DeBell, E. L. Dereniak, J. Harvey, J. Nissley, J. Palmer, A. Selvarajan, and W. L. Wolfe, "Cryogenic Refractive Indices and Temperature Coefficients of Cadmium Telluride from 6 μm to 22 μm ," *Appl. Opt.* **18**:3114–3115 (1979).
217. A. H. Kachare, W. G. Spitzer, and J. E. Fredrickson, "Refractive Index of Ion-Implanted GaAs," *J. Appl. Phys.* **47**:4209–4212 (1976).
218. A. S. Barker and M. Ilegems, "Infrared Lattice Vibrations and Free-Electron Dispersion in GaN," *Phys. Rev. B* **7**:743–750 (1973).
219. F. L. Madaras, J. O. Dimmock, N. Dietz, and K. J. Bachmann, "Sellmeier Parameters for ZnGaP_2 and GaP," *J. Appl. Phys.* **87**:1564–1565 (2000). [Erratum: *J. Appl. Phys.* **87**:7597 (2000).]
220. N. P. Barnes and M. S. Piltch, "Temperature-Dependent Sellmeier Coefficients and Nonlinear Optics Average Power Limit for Germanium," *J. Opt. Soc. Am.* **69**:178–180 (1979).
221. H. H. Li, "Refractive Index of Alkali Halides and Its Wavelength and Temperature Derivatives," *J. Phys. Chem. Ref. Data* **5**:329–528 (1976).
222. N. Umemura, K. Yoshida, and K. Kato, "Phase-Matching Properties of KNbO_3 in the Mid-Infrared," *Appl. Opt.* **38**:991–994 (1999).
223. Y. Fujii and T. Sakudo, "Dielectric and Optical Properties of KTaO_3 ," *J. Phys. Soc. Japan* **41**:888–893 (1976).
224. F. Hanson and D. Dick, "Blue Parametric Generation from Temperature-Tuned LiB_3O_5 ," *Opt. Lett.* **16**:205–207 (1991). [Also see C. Chen, Y. Wu, A. Jiang, B. Wu, G. You, R. Li, and S. Lin, "New Nonlinear-Optical Crystal: LiB_3O_5 ," *J. Opt. Soc. Am. B* **6**:616–621 (1989); K. Kato, "Tunable UV Generation to 0.2325 μm in LiB_3O_5 ," *IEEE J. Quantum Electron.* **QE-26**:1173–1175 (1990); and S. P. Velsko, M. Webb, L. Davis, and C. Huang, "Phase-Matched Harmonic Generation in Lithium Triborate (LBO)," *IEEE J. Quantum Electron.* **QE-27**:2182–2192 (1991).]
225. D. E. Zelmon, D. L. Small, and D. Jundt, "Infrared Corrected Sellmeier Coefficients for Congruently Grown Lithium Niobate and 5 mol% Magnesium Oxide-Doped Lithium Niobate," *J. Opt. Soc. Am.* **B14**:3319–3322 (1997). [Also see D. F. Nelson and R. M. Mikulyak, "Refractive Indices of Congruently Melting Lithium Niobate," *J. Appl. Phys.* **45**:3688–3689 (1974).]
226. W. J. Tropf and M. E. Thomas, "Magnesium Aluminum Spinel (MgAlO_4), in E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991, pp. 881–895. [Improved fit with additional data.]
227. M. J. Dodge, "Refractive Properties of Magnesium Fluoride," *Appl. Opt.* **23**:1980–1985 (1984).
228. M. J. Rosker, K. Cheng, and C. L. Tang, "Practical Urea Optical Parametric Oscillator for Tunable Generation throughout the Visible and Near-infrared," *IEEE J. Quantum Electron.* **QE-21**:1600–1606 (1985).
229. G. E. Jellison, Jr., I. Paulauskas, L. A. Boatner, and D. J. Singh, "Optical Functions of KTaO_3 as Determined by Spectroscopic Ellipsometry and Comparison with Band Structure Calculations," *Phys. Rev.* **B74**:155130 (2006).
230. I. H. Malitson, as quoted by W. L. Wolfe in W. G. Driscoll (ed.), *Handbook of Optics*, 1st ed., McGraw-Hill, New York, 1978. [Our fit to the dispersion data.]
231. F. Weiting and Y. Yixun, "Temperature Effects on the Refractive Index of Lead Telluride and Zinc Sulfide," *Infrared Phys.* **30**:371–373 (1990).
232. S. Singh, J. P. Remeika, and J. R. Potopowicz, "Nonlinear Optical Properties of Ferroelectric Lead Titanate," *Appl. Phys. Lett.* **20**:135–137 (1972).
233. G. Hettner and G. Leisegang, "Dispersion of the Mixed Crystals TlBr-TlI (KRS-5) and TlCl-TlBr (KRS-6) in the Infrared," *Optik* **3**:305–314 (1948) [in German]. [Our fit to the dispersion data.]
234. D. E. Zelmon, D. L. Small, and R. Page, "Refractive-index Measurements of Undoped Yttrium Aluminum Garnet from 0.4 to 5.0 μm ," *Appl. Opt.* **37**:4933–4935 (1998).
235. Y. Nigara, "Measurement of the Optical Constants of Yttrium Oxide," *Jap. J. Appl. Phys.* **7**:404–408 (1968).

236. D. E. Zelmon, E. A. Hanning, and Peter G. Schunemann, "Refractive-Index Measurements and Sellmeier Coefficients for Zinc Germanium Phosphide from 2 to 9 mm with Implications for Phase Matching in Optical Frequency-Conversion Devices," *J. Opt. Soc. Am.* **B18**:1307–1310 (2001).
237. T. M. Bieniewski and S. J. Czyzak, "Refractive Indexes of Single Hexagonal ZnS and CdS Crystals," *J. Opt. Soc. Am.* **53**:496–497 (1963). [Our fit to the dispersion data.]
238. H. H. Li, "Refractive Index of ZnS, ZnSe, and ZnTe and Its Wavelength and Temperature Derivatives," *J. Phys. Chem. Ref. Data* **13**:103–150 (1984).
239. P. Kloeck and L. Colombo, "Index of Refraction, Dispersion, Bandgap and Light Scattering in GeSe and GeSbSe Glasses," *J. Non-Cryst. Solids* **93**:1–16 (1987).
240. S. A. Solin and A. K. Ramdas, "Raman Spectrum of Diamond," *Phys. Rev. B* **1**:1687–1698 (1970). [Also see B. J. Parsons, "Spectroscopic Mode Grüneisen Parameters for Diamond," *Proc. Royal Soc. Lond. A* **352**:397–417 (1977).]
241. B. A. Weinstein and G. J. Piermarini, "Raman Scattering and Phonon Dispersion in Si and GaP at Very High Pressure," *Phys. Rev. B* **12**:1172–1186 (1975).
242. D. Olego and M. Cardona, "Pressure Dependence of Raman Phonons of Ge and 3C-SiC," *Phys. Rev. B* **25**:1151–1160 (1982).
243. P. Vergnat, J. Claudel, A. Hadni, P. Strimer, and F. Vermillard, "Far Infrared Optical Constants of Cesium Halides at Low Temperatures," *J. Phys.* **30**:723–735 (1969) [in French].
244. H. Shimizu, Y. Ohbayashi, K. Yamamoto, K. Abe, M. Midorikawa, and Y. Ishibashi, "Far-Infrared Reflection Spectra of CsCl Single Crystals," *J. Phys. Soc. Japan* **39**:448–450 (1975).
245. S. S. Mitra, "Infrared and Raman Spectra Due to Lattice Vibrations," in *Optical Properties of Solids*, Plenum Press, New York, 1969.
246. R. P. Lowndes and C. H. Perry, "Molecular Structure and Anharmonicity in Thallium Iodide," *J. Chem. Phys.* **58**:271:278 (1973).
247. M. Galtier, A. Montaner, and G. Vidal, "Optical Phonons of CaO, SrO, BaO at the Center of the Brillouin Zone at 300 and 17 K," *J. Phys. Chem. Solids* **33**:2295–2302 (1972) [in French].
248. R. Ramnarine and W. F. Sherman, "The Far-Infrared Investigation of the Mode Frequencies of Some Alkaline Earth Sulfides," *Infrared Phys.* **26**:17–21 (1986).
249. J. L. Jacobson and E. R. Nixon, "Infrared Dielectric Response and Lattice Vibrations of Calcium and Strontium Oxides," *J. Phys. Chem. Solids* **29**:967–976 (1968).
250. Z. V. Popović, G. Stanišić, D. Stojanović, and R. Kostić, "Infrared and Raman Spectral of CdO," *Phys. Stat. Sol. (b)* **165**:K109–K112 (1991).
251. J. R. Jasperse, A. Kahan, J. N. Plendl, and S. S. Mitra, "Temperature Dependence of Infrared Dispersion in Ionic Crystals LiF and MgO," *Phys. Rev.* **146**:526–542 (1966). [Also see A. Kachare, G. Andermann, and L. R. Brantley, "Reliability of Classical Dispersion Analysis of LiF and MgO Reflectance Data," *J. Phys. Chem. Solids* **33**:467–475 (1972).]
252. R. Geick, "Measurement and Analysis of the Fundamental Lattice Vibration Spectrum of PbS," *Phys. Lett.* **10**:51–52 (1964).
253. H. Birkhard, R. Geick, P. Kästner, and K. -H. Unkelback, "Lattice Vibrations and Free Carrier Dispersion in PbSe," *Phys. Stat. Sol. (b)* **63**:89–96 (1974).
254. E. G. Bylander and M. Hass, "Dielectric Constant and Fundamental Lattice Frequency of Lead Telluride," *Solid State Commun.* **4**:51–52 (1966).
255. M. Ilegems and G. L. Pearson, "Infrared Reflection Spectra of Ga_{1-x}Al_xAs Mixed Crystals," *Phys. Rev. B* **1**:1576–1582 (1970).
256. A. Mooradian and G. B. Wright, "First Order Raman Effect in III-V Compounds," *Solid State Commun.* **4**:431–434 (1966). [Also see W. J. Turner and W. E. Reese, "Infrared Lattice Bands in AlSb," *Phys. Rev.* **127**:126–131 (1962).]
257. P. J. Gielisse, S. S. Mitra, J. N. Plendl, R. D. Griffis, L. C. Mansur, R. Marshall, and E. A. Pascoe, "Lattice Infrared Spectra of Boron Nitride and Boron Monophosphide," *Phys. Rev.* **155**:1039–1046 (1967).
258. J. A. Sanjurjo, E. López-Cruz, P. Vogl, and M. Cardona, "Dependence on Volume of the Phonon Frequencies and the IR Effective Charges of Several III-V Semiconductors," *Phys. Rev.* **B28**:4579–4584 (1983).

259. A. Manabe, A. Mitsuishi, and H. Yoshinga, "Infrared Lattice Reflection Spectra of II-VI Compounds," *Jap. J. Appl. Phys.* **6**:593–600 (1967).
260. G. R. Wilkenson, "Raman Spectra of Ionic, Covalent, and Metallic Crystals," in A. Anderson (ed.), *The Raman effect, Vol. 2: Applications*, Marcel Dekker, New York, 1973.
261. A. S. Barker, "Dielectric Dispersion and Phonon Line Shape in Gallium Phosphide," *Phys. Rev.* **165**:917–922 (1968).
262. M. Hass and B. W. Hennis, "Infrared Lattice Reflection Spectra of III-V Compound Semiconductors," *J. Phys. Chem. Solids* **23**:1099–1104 (1962).
263. R. Carles, N. Saint-Cricq, J. B. Renucci, M. A. Renucci, and A. Zwick, "Second-Order Raman Scattering in InAs," *Phys. Rev.* **B22**:4804–4815 (1980).
264. R. B. Sanderson, "Far Infrared Optical Properties of Indium Antimonide," *J. Phys. Chem. Solids* **26**:803–810 (1965).
265. L. Patrick and W. J. Choyke, "Lattice Absorption Bands in SiC," *Phys. Rev.* **123**:813–815 (1965).
266. W. G. Nilsen, "Raman Spectrum of Cubic ZnS," *Phys. Rev.* **182**:838–850 (1969).
267. J. C. Irwin and J. LaCombe, "Raman Scattering in ZnTe," *J. Appl. Phys.* **41**:1444–1450 (1970).
268. R. Shivastava, H. V. Lauer, L. L. Chase, and W. E. Bron, "Raman Frequencies of Fluorite Crystals," *Phys. Lett.* **36A**:333–334 (1971).
269. W. Kaiser, W. G. Spitzer, R. H. Kaiser, and L. E. Howarth, "Infrared Properties of CaF₂, SrF₂, and BaF₂," *Phys. Rev.* **127**:1950–1954 (1962).
270. I. Richman, "Longitudinal Optical Phonons in CaF₂, SrF₂, and BaF₂," *J. Chem. Phys.* **41**:2836–2837 (1966).
271. D. R. Bosomworth, "Far-Infrared Optical Properties of CaF₂, SrF₂, BaF₂, and CdF₂," *Phys. Rev.* **157**:709–715 (1967).
272. J. D. Axe, J. W. Gaglianella, and J. E. Scardefield, "Infrared Dielectric Properties of Cadmium Fluoride and Lead Fluoride," *Phys. Rev.* **139**:A1211–A1215 (1965).
273. R. Droste and R. Geick, "Investigation of the Infrared-Active Lattice Vibration in SrCl₂," *Phys. Stat. Sol. (b)* **62**:511–517 (1974).
274. A. Sadoc, F. Moussa, and G. Pepy, "The Lattice Dynamics of SrCl₂," *J. Phys. Chem. Solids* **37**:197–199 (1976).
275. J. D. Axe and G. D. Pettit, "Infrared Dielectric Dispersion and Lattice Dynamics of Uranium Dioxide and Thorium Dioxide," *Phys. Rev.* **151**:676–679 (1966).
276. P. G. Marlowe and J. P. Russell, "Raman Scattering in Uranium Dioxide," *Phil. Mag.* **14**:409–410 (1966).
277. S. Shin and M. Ishigame, "Defect-Induced Hyper-Raman Spectra in Cubic Zirconia," *Phys. Rev. B* **34**:8875–8882 (1986).
278. S. P. S. Porto and R. S. Krishnan, "Raman Effect of Corundum," *J. Chem. Phys.* **47**:1009–1012 (1967).
279. D. R. Renneke and D. W. Lynch, "Infrared Lattice Vibrations and Dielectric Dispersion in Single-Crystal Cr₂O₃," *Phys. Rev.* **138**:A530–A533 (1965).
280. I. R. Beattie and T. R. Gibson, "The Single-Crystal Raman Spectra of Nearly Opaque Materials. Iron(III) Oxide and Chromium(III) Oxide," *J. Chem. Soc. (A)*, 980–986 (1970). DOI: 10.1039/J19700000980.
281. S. Onari, T. Arai, and K. Kudo, "Infrared Lattice Vibrations and Dielectric Dispersion in α -Fe₂O₃," *Phys. Rev. B* **16**:1717–1721 (1977).
282. G. L. Bottger and C. V. Damsgard, "Raman Scattering in Wurtzite-Type AgI Crystals," *J. Chem. Phys.* **57**:1215–1218 (1972).
283. L. E. McNeil, M. Grimsditch, and R. H. French, "Vibrational Spectroscopy of Aluminum Nitride," *J. Am. Ceram. Soc.* **76**:1132–1136 (1993).
284. E. Loh, "Optical Phonons in BeO Crystals," *Phys. Rev.* **166**:673–678 (1967).
285. C. A. Arguello, D. L. Rousseau, and S. P. S. Porto, "First-Order Raman Effect in Wurtzite-Type Crystals," *Phys. Rev.* **181**:1351–1363 (1969).
286. B. Tell, T. C. Damen, and S. P. S. Porto, "Raman Effect in Cadmium Sulfide," *Phys. Rev.* **144**:771–774 (1966).
287. R. Geick, C. H. Perry, and S. S. Mitra, "Lattice Vibrational Properties of Hexagonal CdSe," *J. Appl. Phys.* **37**:1994–1997 (1966).

288. D. D. Manchon, A. S. Barker, P. J. Dean, and R. B. Zetterstrom, "Optical Studies of the Phonons and Electrons in Gallium Nitride," *Solid State Commun.* **8**:1227–1231 (1970).
289. D. W. Feldman, J. H. Parker, W. J. Choyee, and L. Patrick, "Raman Scattering in 6H SiC," *Phys. Rev.* **170**:698–704 (1968).
290. T. C. Damen, S. P. S. Porto, and B. Tell, "Raman Effect in Zinc Oxide," *Phys. Rev.* **142**:570–574 (1966).
291. R. Geick, U. Schröder, and J. Stuke, "Lattice Vibrational Properties of Trigonal Selenium," *Phys. Stat. Sol.* **24**:99–108 (1967).
292. G. Locovsky, A. Mooradian, W. Taylor, G. B. Wright, and R. C. Keezer, "Identification of the Fundamental Vibrational Modes of Trigonal, α -monoclinic, and Amorphous Selenium," *Solid State Commun.* **5**:113–117 (1967).
293. E. D. Palik, "Tellurium (Te)," in E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991, pp. 709–723.
294. B. D. Saksena, "Analysis of the Raman and Infra-Red Spectra of α -Quartz," *Proc. Ind. Acad. Sci.* **12A**:93–139 (1940).
295. J. F. Scott and S. P. S. Porto, "Longitudinal and Transverse Optical Lattice Vibrations in Quartz," *Phys. Rev.* **161**:903–910 (1967).
296. S. M. Shapiro and J. D. Axe, "Raman Scattering from Polar Phonons," *Phys. Rev. B* **6**:2420–2427 (1972).
297. J. F. Scott, "Raman Spectra of GeO₂," *Phys. Rev. B* **8**:3488–3493 (1970).
298. A. S. Barker and J. A. Detzenberger, "Infrared Lattice Vibrations in CoF₂," *Solid State Commun.* **3**:131–132 (1965).
299. M. Balkanski, P. Moch, and G. Parisot, "Infrared Lattice-Vibration Spectra in NiF₂, CoF₂, and FeF₂," *J. Chem. Phys.* **44**:940–944 (1966).
300. R. M. Macfarlane and S. Ushioda, "Light Scattering from Phonons in CoF₂," *Solid State Commun.* **8**:1081–1083 (1970).
301. S. P. S. Porto, P. A. Fleury, and T. C. Damen, "Raman Spectra of TiO₂, MgF₂, ZnF₂, FeF₂, and MnF₂," *Phys. Rev.* **154**:522–526 (1967).
302. J. Giordano and C. Benoit, "Infrared Spectra of Iron, Zinc, and Magnesium Fluorides: I. Analysis of Results," *J. Phys. C* **21**:2749–2770 (1988). [Also see C. Benoit and J. Giordano, "Dynamical Properties of Crystals of MgF₂, ZnF₂, and FeF₂: II. Lattice Dynamics and Infrared Spectral," *J. Phys. C* **21**:5209–5227 (1988).]
303. A. Kahan, J. W. Goodrum, R. S. Singh, and S. S. Mitra, "Polarized Reflectivity Spectra of Tetragonal GeO₂," *J. Appl. Phys.* **42**:4444–4446 (1971).
304. D. M. Roessler and W. A. Albers, "Infrared Reflectance of Single Crystal Tetragonal GeO₂," *J. Phys. Chem. Solids* **33**:293–296 (1972).
305. A. S. Barker, "Transverse and Longitudinal Optic Mode Study in MgF₂ and ZnF₂," *Phys. Rev.* **136**:A1290–A1295 (1964).
306. R. S. Krishnan and J. P. Russell, "The First-Order Raman Spectrum of Magnesium Fluoride," *Brit. J. Appl. Phys.* **17**:501–503 (1966).
307. J. Giordano, "Temperature Dependence of IR Spectra of Zinc and Magnesium Fluoride," *J. Phys. C* **20**:1547–1562 (1987).
308. R. Summit, "Infrared Absorption in Single-Crystal Stannic Oxide: Optical Lattice-Vibration Modes," *J. Appl. Phys.* **39**:3762–3767 (1967).
309. J. F. Scott, "Raman Spectrum of SnO₂," *J. Chem. Phys.* **53**:852–853 (1970).
310. R. S. Katiyar, P. Dawson, M. M. Hargreave, and G. R. Wilkerson, "Dynamics of the Rutile Structure III. Lattice Dynamics, Infrared and Raman Spectra of SnO₂," *J. Phys. C* **4**:2421–2431 (1971).
311. R. S. Katiyar and R. S. Krishnan, "The Vibrational Spectrum of Rutile," *Phys. Lett.* **25A**:525–526 (1967).
312. A. S. Barker, "Infrared Lattice Vibrations in Calcium Tungstate and Calcium Molybdate," *Phys. Rev.* **135**:A742–A747 (1964).
313. P. Tarte and M. Liegeois-Duyckaerts, "Vibrational Studies of Molybdates, Tungstates and Related Compounds—I. New Infrared Data and Assignments for the Scheelite-Type Compounds X^{II}MoO₄ and X^{II}WO₄," *Spectrochim. Acta* **28A**:2029–2036.

314. V. M. Nagiev, Sh. M. Efendiev, and V. M. Burlakov, "Vibrational Spectra of Crystals with Scheelite Structure and Solid Solutions on Their Basis," *Phys. Stat. Sol. (b)* **125**:467–475 (1984).
315. J. M. Stencel, E. Silberman, and J. Springer, "Temperature-Dependent Reflectivity, Dispersion Parameters, and Optical Constants for PbWO_4 ," *Phys. Rev. B* **12**:5435–5441 (1976).
316. S. A. Miller, H. E. Rast, and H. H. Caspers, "Lattice Vibrations of LiYF_4 ," *J. Chem. Phys.* **53**:4172–4175 (1970).
317. E. Schultheiss, A. Scharmann, and D. Schwabe, "Lattice Vibrations in BiLiF_4 and YLiF_4 ," *Phys. Stat. Sol. (b)* **138**:465–475 (1986).
318. S. P. S. Porto and J. F. Scott, "Raman Spectra of CaWO_4 , SrWO_4 , CaMoO_4 , and SrMoO_4 ," *Phys. Rev.* **157**:716–719 (1967).
319. S. Desgreniers, S. Jandl, and C. Carlone, "Temperature Dependence of the Raman Active Phonons in CaWO_4 , SrWO_4 , and BaWO_4 ," *J. Phys. Chem. Solids* **45**:1105–1109 (1984).
320. R. K. Khanna, W. S. Brower, B. R. Guscott, and E. R. Lippincott, "Laser Induced Raman Spectra of Some Tungstates and Molybdates," *J. Res. Nat. Bur. Std.* **72A**:81–84 (1968).
321. M. P. O'Horo, A. L. Frisillo, and W. B. White, "Lattice Vibrations of MgAl_2O_4 Spinel," *J. Phys. Chem. Solids* **34**:23–28 (1973).
322. M. E. Strifler and S. I. Boldish, "Transverse and Longitudinal Optic Mode Frequencies of Spinel MgAl_2O_4 ," *J. Phys. C* **11**:L237–L241 (1978).
323. K. Yamamoto, T. Murakawa, Y. Ohbayashi, H. Shimizu, and K. Abe, "Lattice Vibrations in CdIn_2S_4 ," *J. Phys. Soc. Japan* **35**:1258 (1973).
324. H. Shimizu, Y. Ohbayashi, K. Yamamoto, and K. Abe, "Lattice Vibrations in Spinel-Type CdIn_2S_4 ," *J. Phys. Soc. Japan* **38**:750–754 (1975).
325. H. D. Lutz, B. Müller, and H. J. Steiner, "Lattice Vibration Spectra. LIX. Single Crystal Infrared and Raman Studies of Spinel Type Oxides," *J. Solid State Chem.* **90**:54–60 (1991).
326. H. D. Lutz, G. Wäschenbach, G. Kliche, and H. Haeuseler, "Lattice Vibrational Spectra, XXXIII: Far-Infrared Reflection Spectra, TO and LO Phonon Frequencies, Optical and Dielectric Constants, and Effective Changes of the Spinel-Type Compounds MCr_2S_4 ($M = \text{Mn, Fe, Co, Zn, Cd, Hg}$), MCr_2Se_4 ($M = \text{Zn, Cd, Hg}$), and MIn_2S_4 ($M = \text{Mn, Fe, Co, Zn, Cd, Hg}$)," *J. Solid State Chem.* **48**:196–208 (1983).
327. K. Wakamura, H. Iwatani, and K. Takarabe, "Vibrational Properties of One- and Two-Mode Behavior in Spinel Type Mixed Systems $\text{Zn}_{1-x}\text{Cd}_x\text{Cr}_2\text{S}_4$," *J. Phys. Chem. Solids* **48**:857–861 (1987).
328. H. C. Gupta, G. Sood, A. Parashar, and B. B. Tripathi, "Long Wavelength Optical Lattice Vibrations in Mixed Chalcogenide Spinels $\text{Zn}_{1-x}\text{Cd}_x\text{Cr}_2\text{S}_4$ and $\text{CdCr}_2(\text{S}_{1-x}\text{Se}_x)_4$," *J. Phys. Chem. Solids* **50**:925–929 (1989).
329. A. S. Barker and J. H. Hopfield, "Coupled-Optical-Phonon-Mode Theory of the Infrared Dispersion in BaTiO_3 , SrTiO_3 , and KTaO_3 ," *Phys. Rev.* **135**:A1732–A1737 (1964).
330. A. S. Barker, "Temperature Dependence of the Transverse and Longitudinal Optic Mode Frequencies and Charges in SrTiO_3 and BaTiO_3 ," *Phys. Rev.* **145**:391–399 (1966).
331. J. L. Servoin, Y. Luspin, and F. Gervais, "Infrared Dispersion in SrTiO_3 at High Temperature," *Phys. Rev. B* **22**:5501–5506 (1980).
332. C. H. Perry and E. F. Young, "Infrared Studies of Some Perovskite Fluorides. I. Fundamental Lattice Vibrations," *J. Appl. Phys.* **38**:4616–4624 (1967).
333. A. Scalabrin, A. S. Chaves, D. S. Shim, and S. P. S. Porto, "Temperature Dependence of the A_1 and E Optical Phonons in BaTiO_3 ," *Phys. Stat. Sol. (b)* **79**:731–742 (1977).
334. J. L. Servoin, F. Gervais, A. M. Quittet, and Y. Luspin, "Infrared and Raman Responses in Ferroelectric Perovskite Crystals," *Phys. Rev. B* **21**:2038–2041 (1980).
335. G. Burns and B. A. Scott, "Lattice Modes in Ferroelectric Perovskites: PbTiO_3 ," *Phys. Rev. B* **7**:3088–3101 (1973).
336. J. P. van der Ziel, A. E. Meixner, H. M. Kasper, and J. A. Ditzenberger, "Lattice Vibrations of AgGaS_2 , AgGaSe_2 , and CuGaS_2 ," *Phys. Rev. B* **9**:4286–4294 (1974).
337. W. H. Koschel and M. Bettini, "Zone-Centered Phonons in $\text{A}^{\text{I}}\text{B}^{\text{III}}\text{S}_2$ Chalcopyrites," *Phys. Stat. Sol. (b)* **72**:729–737 (1975).

338. V. G. Tyuterev and S.I. Skachkov, "Lattice Dynamics, Thermodynamic and Elastic Properties of AgGaS_2 ," *Nuovo Cimento* **14D**:1097–1103 (1992).
339. A. Miller, G. D. Holah, and W. C. Clark, "Infrared Dielectric Dispersion of ZnGeP_2 and CdGeP_2 ," *J. Phys. Chem. Solids* **35**:685–693 (1974).
340. M. Bettini and A. Miller, "Optical Phonons and ZnGeP_2 and CdGeP_2 ," *Phys. Stat. Sol. (b)* **66**:579–586 (1974).
341. I. P. Kaminow, E. Buehler, and J. H. Wernick, "Vibrational Modes in ZnSiP_2 ," *Phys. Rev. B* **2**:960–966 (1970).
342. L. Artus, J. Pascual, A. Goulet, and J. Camassel, "Polarized Infrared Spectra of AgGaSe_2 ," *Solid State Commun.* **69**:753–756 (1989).
343. G. D. Holah, A. Miller, W. D. Dunnett, and G. W. Isler, "Polarised Infrared Reflectivity of CdGeAs_2 ," *Solid State Commun.* **23**:75–78 (1977).
344. E. V. Antropova, A. V. Kopytov, and A. S. Poplavnoi, "Phonon Spectrum and IR Optical Properties of CdGeAs_2 ," *Opt. Spectrosc.* **64**:766–768 (1988).
345. R. Zallen, M. L. Slade, and A. T. Ward, "Lattice Vibrations and Interlayer Interactions in Crystalline As_2S_3 and As_2Se_3 ," *Phys. Rev. B* **3**:4257–4273 (1971).
346. S. P. S. Porto, J. A. Giordmaine, and T. C. Damen, "Depolarization of Raman Scattering in Calcite," *Phys. Rev.* **147**:608–611 (1966).
347. K. H. Hellwege, W. Lesch, M. Plihal, and G. Schaack, "Two Phonon Absorption Spectra and Dispersion of Phonon Branches in Crystals of Calcite Structure," *Z. Physik* **232**:61–86 (1970). [Also see R. K. Vincent, "Emission Polarization Study on Quartz and Calcite," *Appl. Opt.* **11**:1942–1945 (1972).]
348. J. Q. Lu, G. X. Lan, B. Li, Y. Y. Yang, and H. F. Wang, "Raman Scattering Study of the Single Crystal $\beta\text{-BaB}_2\text{O}_4$ under High Pressure," *J. Phys. Chem. Solids* **49**:519–527 (1988).
349. W. Wojdowski, "Vibrational Modes in $\text{Bi}_{12}\text{GeO}_{20}$ and $\text{Bi}_{12}\text{SiO}_{20}$ Crystals," *Phys. Stat. Sol. (b)* **130**:121–130 (1985).
350. X. Hu, J. Wang, B. Teng, C-K. Loong, and M. Grimsditch, "Raman Study of Phonons in Bismuth Triborate BiB_3O_6 Crystal," *J. Appl. Phys.* **97**:033501 (2005).
351. D. Kasprowicz, T. Runka, M. Szybowicz, P. Ziobrowski, A. Majchrowski, E. Michalski, and M. Drozdowski, "Characterization of Bismuth Triborate Single Crystal Using Brillouin and Raman Spectroscopy," *Cryst. Res. Technol.* **40**:459–465 (2005).
352. H. Vogt, T. Chattopadhyay, and H. J. Stolz, "Complete First-Order Raman Spectra of the Pyrite Structure Compounds FeS_2 , MnS_2 , and SiP_2 ," *J. Phys. Chem. Solids* **44**:869–873 (1983).
353. H. D. Lutz, G. Schneider, and G. Kliche, "Far-Infrared Reflection Spectra, TO- and LO-Phonon Frequencies, Coupled and Decoupled Plasmon-Phonon Modes, Dielectric Constants, and Effective Dynamical Charges of Manganese, Iron, and Platinum Group Pyrite Type Compounds," *J. Phys. Chem. Solids* **46**:437–443 (1985).
354. D. K. Agrawal and C. H. Perry, "The Temperature Dependent Raman Spectra of KDP, KD^*P , KDA, and ADP," in M. Balkanski (ed.), *Proceedings of the Second International Conference on Light Scattering in Solids*, Flammarion Sciences, Paris, 1971, pp. 429–435.
355. G. E. Kugel, F. Bréhat, B. Wyncke, M. D. Fontana, G. Marnier, C. Carabatos-Nedelec, and J. Mangin, "The Vibrational Spectrum of KTiOPO_4 Single Crystal Studied by Raman and Infrared Reflective Spectroscopy," *J. Phys. C* **21**:5565–5583 (1988).
356. B. Mohamadou, G. E. Kugel, F. Brehat, B. Wyncke, G. Marnier, P. Simon, "High-Temperature Vibrational Spectra, Relaxation, and Ionic Conductivity Effects in KTiOPO_4 ," *J. Phys. Condens. Matter* **3**:9489–9501 (1991).
357. R. P. Bauman and S. P. S. Porto, "Lattice Vibrations and Structure of Rare-Earth Fluorides," *Phys. Rev.* **161**:842–847 (1967).
358. R. P. Lowndes, J. F. Parrish, and C. H. Perry, "Optical Phonons and Symmetry of Tysonite Lanthanide Fluorides," *Phys. Rev.* **182**:913–922 (1969).
359. E. Liarokapis, E. Anastassakis, and G. A. Kourouklis, "Raman Study of Phonon Anharmonicity in LaF_3 ," *Phys. Rev. B* **32**:8346–8355 (1985).

360. H. R. Xia, L. X. Li, H. Yu, S. M. Dong, J. Y. Wang, Q. M. Lu, C. Q. Ma, and X. N. Wang, "Structure and the Nonlinearity of Lithium Triborate Studied by Raman and Infrared Reflectivity Spectroscopy," *J. Mater. Res.* **16**:3465–3470 (2001).
361. W. Otaguro, E. Weiner-Avnera, C. A. Arguello, and S. P. S. Porto, "Phonons, Polaritons, and Oblique Phonons in LiIO_3 by Raman Scattering and Infrared Reflection," *Phys. Rev. B* **4**:4542–4551 (1971).
362. J. L. Duarte, J. A. Sanjurjo, and R. S. Katiyar, "Off-Normal Infrared Reflectivity in Uniaxial Crystals: α - LiIO_3 and α -quartz," *Phys. Rev. B* **36**:3368–3372 (1987).
363. R. Claus, G. Borstel, E. Wiesendanger, and L. Steffan, "Directional Dispersion and Assignment of Optical Phonons in LiNbO_3 ," *Z. Naturforsch.* **27a**:1187–1192 (1972).
364. X. Yang, G. Lan, B. Li, and H. Wang, "Raman Spectra and Directional Dispersion in LiNbO_3 and LiTaO_3 ," *Phys. Stat. Sol. (b)* **141**:287–300 (1987).
365. D. G. Boziniš and J. P. Hurrell, "Optical Modes and Dielectric Properties of Ferroelectric Orthorhombic KNbO_3 ," *Phys. Rev. B* **13**:3109–3120 (1976).
366. N. N. Syrbu, V. T. Krasovsky, and I. N. Grincheshen, "Infrared Vibrational Modes in Ti_3SbS_3 , Ti_3AsS_3 , and Ti_3AsSe_3 Crystals," *Cryst. Res. Technol.* **29**:1095–1102 (1994).
367. A. S. Pine and G. Dresselhaus, "Raman Scattering in Paratellurite," *Phys. Rev. B* **5**:4087–4093 (1972).
368. D. M. Korn, A. S. Pine, G. Dresselhaus, and T. B. Reed, "Infrared Reflectivity of Paratellurite, TeO_2 ," *Phys. Rev. B* **8**:768–772 (1973).
369. W. J. Tropf and M. E. Thomas, "Yttrium Oxide (Y_2O_3)," in E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991, pp. 1081–1098.
370. C. Z. Bi, J. Y. Ma, J. Yan, X. Fang, D. Z. Yao, B. R. Zhao, and X. G. Qiu, "Far-Infrared Optical Properties of YVO_4 Single Crystal," *Eur. Phys. J.* **B51**:167–171 (2006).
371. S. A. Miller, H. H. Caspers, and H. E. "Rast, Lattice Vibrations of Yttrium Vanadate," *Phys. Rev.* **168**:964–969 (1968). [Also see H. E. Rast, H. H. Caspers, and S. A. Miller, "Infrared Spectral Emittance and Optical Properties of Yttrium Vanadate," *Phys. Rev.* **169**:705–709 (1968).]
372. M. Thirumavalavan, J. Kumar, F. D. Gnanam, and P. Ramasamy, "Vibrational Spectra of $\text{Y}_3\text{Al}_5\text{O}_{12}$ Crystals Grown from Ba- and Pb-Based Flux Systems," *Infrared Phys.* **26**:101–103 (1986).
373. G. A. Gledhill, P. M. Nikolić, A. Hamilton, S. Stojilković, V. Blagojević, P. Mihajlovic, and S. Djurić, "FIR Optical Properties of Single Crystal $\text{Y}_3\text{Al}_5\text{O}_{12}$ (YAG)," *Phys. Stat. Sol. (b)* **163**:K123–K128 (1991).
374. R. P. Lowndes, "Anharmonicity in the Silver and Thallium Halides: Far-Infrared Dielectric Response," *Phys. Rev. B* **6**:1490–1498 (1972).
375. G. Lucovsky, "Optic Modes in Amorphous As_2S_3 and As_2Se_3 ," *Phys. Rev. B* **6**:1480–1489 (1972).
376. W. G. Spitzer, R. C. Miller, D. A. Kleinman, and L. E. Howarth, "Far Infrared Dielectric Dispersion in BaTiO_3 , SrTiO_3 , and TiO_2 ," *Phys. Rev.* **126**:1710–1721 (1962).
377. E. J. Danielewicz and P. D. Coleman, "Far Infrared Optical Properties of Selenium and Cadmium Telluride," *Appl. Opt.* **13**:1164–1170 (1974).
378. J. L. Verble and R. F. Wallis, "Infrared Studies of the Lattice Vibrations in Iron Pyrite," *Phys. Rev.* **182**:783–789 (1969).
379. H. D. Lutz, G. Kliche, and H. Haeuselner, "Lattice Vibrational Spectra XXIV: Far-infrared Reflection Spectra, Optical and Dielectric Constants, and Effective Charges of Pyrite Type Compounds FeS_2 , MnS_2 , MnSe_2 , and MnTe_2 ," *Z. Naturforsch.* **86a**:184–190 (1981).
380. C. J. Johnson, G. H. Sherman, and R. Weil, "Far Infrared Measurement of the Dielectric Properties of GaAs and CdTe at 300 K and 8 K," *Appl. Opt.* **8**:1667–1671 (1969).
381. D. A. Kleinman and W. G. Spitzer, "Infrared Lattice Absorption of GaP" *Phys. Rev.* **118**:110–117 (1960).
382. A. Hadni, J. Claudel, D. Chanal, P. Strimer, and P. Vergnat, "Optical Constants of Potassium Bromide in the Far Infrared," *Phys. Rev.* **163**:836–843 (1967).
383. A. Hadni, J. Claudel, G. Morlot, and P. Strimer, "Transmission and Reflection Spectra of Pure and Doped Potassium Iodide at Low Temperature," *Appl. Opt.* **7**:161–165 (1968) [in French].
384. J. Zarembowitch, J. Gouteron, and A. M. Lejus, "Raman Spectra of Lanthanide Sesquioxide Single Crystals with A-type Structure," *Phys. Stat. Sol. (b)* **94**:249–256 (1979).

385. A. S. Barker and R. Loudon, "Dielectric Properties and Optical Phonons in LiNbO_3 ," *Phys. Rev.* **158**:433–445 (1967).
386. I. F. Chang and S. S. Mitra, "Temperature Dependence of Long-Wavelength Optic Phonons of NaF Single Crystals," *Phys. Rev. B* **5**:4094–4100 (1972).
387. H. Burkhard, R. Geick, P. Kästner, and K. -H. Unkelbach, "Lattice Vibrations and Free Carrier Dispersion in PbSe," *Phys. Stat. Sol. (b)* **63**:89–96 (1974).
388. W. G. Spitzer and D. A. Kleinman, "Infrared Lattice Bands of Quartz," *Phys. Rev.* **121**:1324–1335 (1961).
389. B. Orel and V. Moissenko, "A Vibrational Study of Piezoelectric TeO_2 Crystals," *Phys. Stat. Sol.(b)* **165**: K37–K41 (1991).
390. A. M. Hofmeister and K. R. Campbell, "Infrared Spectroscopy of Yttrium Aluminum, Yttrium Gallium, and Yttrium Iron Garnets," *J. Appl. Phys.* **72**:638–646 (1992).

This page intentionally left blank.

DO NOT DUPLICATE

POLYMERIC OPTICS

John D. Lytle

Advanced Optical Concepts
Santa Cruz, California

3.1 GLOSSARY

$A_{\text{H}_2\text{O}}$	water absorption
K	thermal conductivity
T_s	maximum service temperature
α	thermal expansion coefficient
ρ	density

3.2 INTRODUCTION

A small number of carbon-based polymeric materials possesses some of those qualities which have made glass an attractive optical material. Most of these polymeric materials do exhibit certain physical deficiencies compared to glass. But, despite the fact that “plastic optics” has acquired an image as a low-end technology, it may nonetheless be a better choice, or even the best choice, in certain applications.

Selection Factors

Virtually all of the polymers having useful optical properties are much less dense than any of the optical glasses, making them worthy of consideration in applications where weight-saving is of paramount importance. Many of them exhibit impact resistance properties which exceed those of any silicate glass, rendering them well-suited to military applications (wherein high “g” loads may be encountered), or ideal for some consumer products in which safety may be a critical consideration.

Though the physical properties of the polymers may make them better matched to certain design requirements than glass, by far the most important advantage of polymeric optics is the considerable creative freedom they make available to the optical and mechanical design effort.¹ While the design constraints and guidelines governing glass optics design and fabrication are fairly well defined, the various replication processes which may be put to use in polymer optics fabrication make available unique

opportunities for the creation of novel optical components and systems which would be unthinkable or unworkable in glass. Oftentimes, the differences in the engineering approach, or in the production processes themselves, may make possible very significant cost reductions in high-volume situations.²

3.3 FORMS

Thermoset Resins

Optical polymers fall into basically two categories—the *thermoset* resins and the *thermoplastic* resins. The thermoset resin group consists of chemistries in which the polymerization reaction takes place during the creation of the part, which may be produced by casting, or by transfer replication. The part which has been created at completion of the reaction may then be postprocessed, if desired, by machining. In general, the thermoset resins cannot be melted and re-formed.

The most commonly encountered thermoset optical resin is that used to produce ophthalmic lenses for eyewear.³ The monomer, which is stored in liquid form at reduced temperature, is introduced into a mold, where the polymerization reaction takes place, forming a part which assumes the shape of the cavity containing it. Alternatively, epoxy-based chemistries have been used with some success to form replicated reflecting surface shapes by a transfer process, and to produce aspheric figuring (at relatively modest expense) upon spherical refractive or reflective substrates.

Thermoplastic Resins

With the possible exception of eyewear, most polymeric optics are executed in thermoplastic materials which are supplied in already polymerized form.⁴ These materials are normally purchased in bulk as small pellets. These pellets are heated to a temperature beyond the softening point, so that they flow to become a single viscous mass. This mass is then formed to assume the shape desired in the final part.

Parts may be created by the injection molding process, in which the heated polymer is squirted into a mold at high pressure and allowed to cool in the shape of the desired component. Or the pellets may be directly heated between the two halves of a compression mold, and the mold closed to effect formation of the part. Hybrid molding technologies combining these two processes are recently experiencing increasing popularity in optical molding applications, and have produced optical surface figures of very high quality.

The capability of modern molding technology to produce optics having very good surface-figure quality has made possible the creation of polymeric optical components for a wide variety of applications. Among these are medical disposables, intraocular lenses, a host of consumer products, military optics, and a number of articles in which optical, mechanical, and electrical functions are combined in a single part.⁵

3.4 PHYSICAL PROPERTIES

Density

Optical glass types number in the hundreds (if all manufacturers worldwide are counted). The glass types available from the catalogs cover a wide range of optical, physical, thermal, and chemical properties. The density of these materials varies from about 2.3 g/cm³ to about 6.3 g/cm³. The heaviest optically viable polymer possesses a density of only about 1.4 g/cm³, whereas the lightest of these materials will readily float in water, having a density of 0.83 g/cm³.⁶ All other things being equal, the total element count in an optical system may often be reduced (at modest cost penalty) by the inclusion of nonspherical surfaces. All things considered, then, polymeric optical systems may be made much less massive than their glass counterparts, especially if aspheric technology is applied to the polymer optical trains.

Hardness

Although cosmetic blemishes rarely impact final image quality (except in the cases of field lenses or reticles), optical surfaces are customarily expected to be relatively free of scratches, pits, and

the like. Ordinary usage, especially cleaning procedures, are likely to result in some scratching with the passage of time. Most common optical glasses possess sufficient hardness that they are relatively immune to damage, if some modest amount of care is exercised.

The polymeric optical materials, on the other hand, are often so soft that a determined thumbnail will permanently indent them. The hardness of polymeric optics is difficult to quantify (in comparison to glass), since this parameter is not only material-dependent, but also dependent upon the processing. Suffice it to say that handling procedures which would result in little or no damage to a glass element may produce considerable evidence of abrasion in a polymeric surface, particularly in a thermoplastic. In fact, the compressibility of most thermoplastic polymers is such that the support for hard surface coatings is sufficiently low that protection provides immunity against only superficial abrasion. These deficiencies are of no particular consequence, however, if the questionable surfaces are internal, and thereby inaccessible.

Rigidity

A property closely related to hardness is the elastic modulus, or Young's modulus. This quantity, and the elongation factor at yield, are determinants of the impact resistance, a performance parameter in which the polymers outshine the glasses. These properties are, again, dependent upon the specified polymeric alloy, any additives which may be present, and processing history of the polymer, and cannot be dependably quoted.^{7,8} The reader is referred to any of several comprehensive references listed herein for mechanical properties data. Those properties which create good impact resistance become liabilities if an optical part is subjected to some torsion or compressive stress. Since optical surface profiles must often be maintained to subwavelength accuracy, improper choice of the thickness/diameter ratio, or excessive compression by retaining rings, may produce unacceptable optical figure deformations.

Polymer chemistry is a complex subject probably best avoided in a discussion of polymer optics. Carbon-based polymers have been synthesized to include an extensive variety of chemical subgroups, however. Unfortunately, relatively few of these materials are actually in regular production, and only a handful of those possess useful optical properties for imaging purposes.

Service Temperature

Any decision involving a glass/plastic tradeoff should include some consideration of the anticipated thermal environment. While the optical glasses may exhibit upper service temperature limits of from 400 to 700°C, many of the glass types having the most interesting optical properties are quite fragile, and prone to failure if cooled too quickly. These failures are mostly attributable to cooling-induced shrinkage of the skin layer, which shatters because the insulating properties of the material prevent cooling (and shrinkage) of the bulk material at the same rate.

The polymeric materials, on the other hand, have much lower service temperature limits, in some cases no higher than about 60°C.⁹ The limit may approach 250°C for some of the fluoropolymers. The thermal conductivity of many of these polymers may be as much as an order of magnitude lower than for the glasses and the thermal expansion coefficients characterizing the polymers are often an order of magnitude larger than those associated with optical glasses. Consequently, subjecting any polymeric optical element to a significant thermal transient is likely to create more severe thermal gradients in the material, and result in significant thermally induced optical figure errors.¹⁰ Again, it is suggested that the interested reader consult the plastic handbooks and manufacturer's literature for a complete listing of this behavior, as additives and variation in molecular weight distribution may significantly affect all of these properties. Some of the most important physical properties of the more readily available optical polymers are tabulated in Table 1.

Conductivity (Thermal, Electrical)

Most materials which exhibit poor thermal conductivity are also poor electrical conductors. Since many unfilled polymers are very effective electrical insulators, they acquire static surface charge

TABLE 1 Physical Properties

Material	ρ	α	T_s	K	A_{H_2O}
P-methylmethacrylate	1.18	6.0	85	4–6	0.3
P-styrene	1.05	6.4–6.7	80	2.4–3.3	0.03
NAS	1.13	5.6	85	4.5	0.15
Styrene acrylonitrile (SAN)	1.07	6.4	75	2.8	0.28
P-carbonate	1.25	6.7	120	4.7	0.2–0.3
P-methyl pentene	0.835	11.7	115	4.0	0.01
P-amide (Nylon)	1.185	8.2	80	5.1–5.8	1.5–3.0
P-arylate	1.21	6.3		7.1	0.26
P-sulfone	1.24	2.5	160	2.8	0.1–0.6
P-styrene co-butadiene	1.01	7.8–12			0.08
P-cyclohexyl methacrylate	1.11				
P-allyl diglycol carbonate	1.32		100	4.9	
Cellulose acetate butyrate	1.20			4.0–8.0	
P-ethersulfone	1.37	5.5	200	3.2–4.4	
P-chloro-trifluoroethylene	2.2	4.7	200	6.2	0.003
P-vinylidene fluoride	1.78	7.4–13	150		0.05
P-etherimide	1.27	5.6	170		0.25

fairly easily, and dissipate it very slowly. Not surprisingly, these areas of surface charge quickly attract oppositely charged contaminants, most of which are harder than the plastic. Attempts to clear the accumulated particles from the surfaces by cleaning can, and usually do, result in superficial damage. Application of inorganic coatings to these surfaces may do double duty by providing a more conductive surface (less likely to attract contaminants), while improving the abrasion resistance.

Outgassing

In contrast to glass optical parts, which normally have very low vapor pressure when properly cleaned, most polymers contain lubricants, colorants, stabilizers, and so on, which may outgas throughout the life of the part. This behavior disqualifies most plastic optical elements from serving in space-borne instrumentation, since the gaseous products, once lost, surround the spacecraft, depositing upon solar panels and other critical surfaces. Some, but few, thermoset resins may be clean enough for space applications if their reaction stoichiometry is very carefully controlled in the creation of the part.

Water Absorption

Most polymers, particularly the thermoplastics, are hygroscopic. They absorb and retain water, which must, in most cases, be driven off by heating prior to processing. Following processing, the water will be reabsorbed if the surfaces are not treated to inhibit absorption. Whereas only a very small amount of water will normally attach to the surfaces of a glass optical element, the polymer materials used for optics may absorb from about 0.003 to about 2 percent water by weight. Needless to say, the trapped water may produce dimensional changes, as well as some minor alterations of the spectral transmission. Physical properties of some of the more familiar optical polymers are listed in Table 1. Density = ρ (g/cm^3); thermal expansion coefficient = α ($\text{cm}/\text{cm } ^\circ\text{C} \times 10^{-5}$); max. service temperature = T_s ($^\circ\text{C}$); thermal conductivity = K ($\text{cal}/\text{sec cm } ^\circ\text{C} \times 10^4$); and water absorption (24 h) = A_{H_2O} (%). Values are to be considered approximate, and may vary with supplier and processing variations.

Additives

Polymers are normally available in a variety of “melt flow” grades—each of which possesses viscosity properties best suited to use in parts having specific form factors. A number of additives are commonly present in these materials. Such additives may, or may not, be appropriate in an optical application. Additives for such things as flame retardancy, lubricants, lubrication, and mold release are best avoided if not included to address a specific requirement. Frequently, colorants are added for the purpose of neutralizing the naturally occurring coloration of the material. These additives create an artificial, but “clear,” appearance. The colorants must, of course, absorb energy to accomplish this, resulting in a net reduction in total spectral transmission.

Radiation Resistance

Most of the optical polymers will be seen to exhibit some amount of fluorescence if irradiated by sufficiently intense high-energy radiation.¹¹ High-energy radiation of the ultraviolet and ionizing varieties will, in addition, produce varying amounts of polymer chain crosslinking, depending upon the specific polymer chemistry. Crosslinking typically results in discoloration of the material, and some amount of nonuniform energy absorption. Inhibitors may be added to the polymeric material to retard crosslinking, although, oddly enough, the polymers most susceptible to UV-induced discoloration are generally the least likely to be affected by ionizing radiation, and vice versa.

Documentation

Although polymeric materials suffer some shortcomings in comparison to glass (for optical applications), distinct advantages do exist. The major obstacle to the use of polymers, however, is the spotty and imprecise documentation of many of those properties required for good engineering and design. In general, the resin producers supply these materials in large quantity to markets wherein a knowledge of the optical properties is of little or no importance. With luck, the documentation of optical properties may consist of a statement that the material is “clear.” In the rare case where refractive index is documented, the accuracy may be only two decimal places. In these circumstances, the optical designer or molder is left to investigate these properties independently—a complex task, since the processing itself may affect those properties to a substantial degree.

Unfortunately, optical applications may represent only a small fraction of a percent of the total market for a given resin formulation, and since these materials are sold at prices ranging from less than two dollars to a few dollars per pound, the market opportunity represented by optical applications seems minuscule to most polymer vendors.

3.5 OPTICAL PROPERTIES

Variations

It is only a fortuitous accident that some of the polymers exhibit useful optical behavior, since most all of these materials were originally developed for other end uses. The possible exceptions are the materials used for eyeglass applications (poly-diallylglycol), and the materials for optical information storage (specially formulated polycarbonate). Citation of optical properties for any polymeric material must be done with some caution and qualification, as different melt flow grades (having different molecular weight distribution) may exhibit slightly different refractive index properties. Additives to regulate lubricity, color, and so on can also produce subtle alterations in the spectral transmission properties.

Spectral Transmission

In general, the carbon-based optical polymers are visible-wavelength materials, absorbing fairly strongly in the ultraviolet and throughout the infrared.¹²⁻¹⁴ This is not readily apparent from the absorption spectra published in numerous references, though. Such data are normally generated by spectroscopists for the purpose of identifying chemical structure, and are representative of very thin samples. One can easily develop the impression from this information that the polymers transmit well over a wide spectral range. Parenthetically, most of these polymers, while they have been characterized in the laboratory, are not commercially available. What is needed for optical design purposes is transmission data (for available polymers) taken from samples having sufficient thickness to be useful for imaging purposes.

Some specially formulated variants of poly-methylmethacrylate have useful transmission down to 300 nm.¹⁵ Most optical polymers, however, begin to absorb in the blue portion of the visible spectrum, and have additional absorption regions at about 900 nm, 1150 nm, 1350 nm, finally becoming totally opaque at about 2100 nm. The chemical structure which results in these absorption regions is common to almost all carbon-based polymers, thus the internal transmittance characteristics of these materials are remarkably similar, with the possible exception of the blue and near-UV regions. A scant few polymers do exhibit some spotty narrowband transmission leakage in the far-infrared portion of the spectrum, but in thicknesses suitable only for use in filter applications.

Refractive Index

The chemistry of carbon-based polymers is markedly different from that of silicate glasses and inorganic crystals in common use as optical materials. Consequently, the refractive properties differ significantly. In general, the refractive indices are lower, extending to about 1.73 on the high end, and down to a lower limit of about 1.3. In practice, those materials which are readily available for purchase exhibit a more limited index range—from about 1.42 to 1.65. The Abbe values for these materials vary considerably, though, from about 100 to something less than 20. Refractive index data for a few of these polymers, compiled from a number of sources, is displayed in Table 2. In the chart, PMMA signifies polymethylmethacrylate; P-styr, polystyrene; P-carb, polycarbonate; SAN, styrene acrylonitrile; PEI, polyetherimide; PCHMA, polycyclohexylmethacrylate. The thermo-optic coefficients at room temperature (change in refractive index with temperature) are also listed. Note that these materials, unlike most glasses, experience a reduction in refractive index with increasing temperature. Figure 1, a simplified rendition of the

TABLE 2 Refractive Index of Some Optical Polymers

Line ID	Wavl., nm	PMMA	P-styr	P-carb	SAN	PEI	PCHMA
	1014.0	1.4831	1.5726	1.5672	1.5519		
s	852.1	1.4850	1.5762	1.5710	1.5551		
r	706.5	1.4878	1.5820	1.5768	1.5601		
C	656.3	1.4892	1.5849	1.5799	1.5627		1.502
C'	643.9	1.4896	1.5858	1.5807	1.5634	1.651	
D	589.3	1.4917	1.5903	1.5853	1.5673		
d	587.6	1.4918	1.5905	1.5855	1.5674	1.660	1.505
e	546.1	1.4938	1.5950	1.5901	1.5713	1.668	
F	486.1	1.4978	1.6041	1.5994	1.5790		1.511
F'	480.0	1.4983	1.6052	1.6007	1.5800	1.687	
g	435.8	1.5026	1.6154	1.6115	1.5886		
h	404.7	1.5066	1.6253	1.6224	1.5971		
i	365.0	1.5136	1.6431	1.6432	1.6125		
Abbe number		57.4	30.9	29.9	34.8	18.3	56.1
$dn/dT \times 10^{-4}/^{\circ}\text{C}$		-1.05	-1.4	-1.07	-1.1		

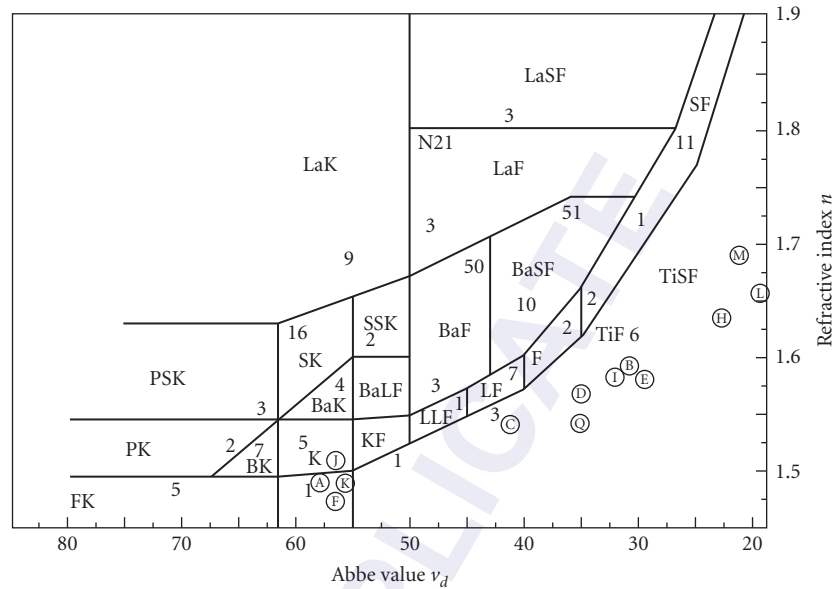


FIGURE 1 Optical glasses and polymers: (a) polymethylmethacrylate; (b) polystyrene; (c) NAS; (d) styrene acrylonitrile; (e) polycarbonate; (f) polymethyl pentene; (g) acrylonitrile-butadiene styrene (ABS); (h) polysulfone; (i) polystyrene co-maleic anhydride; (j) polycyclohexylmethacrylate (PCHMA); (k) polyallyl diglycol carbonate; (l) polyetherimide (PEI); and (m) polyvinyl naphthalene.

familiar glass map (n vs. ν), shows the locations of some of the more familiar polymers. Note that these materials all occupy the lower and right-hand regions of the map. In the Schott classification system, the polymers populate mostly the FK, TiK, and TiF regions of the map.¹⁶

Homogeneity

It must be kept constantly in mind that polymeric optics are molded and not mechanically shaped. The exact optical properties of a piece cannot, therefore, be quantified prior to manufacture of the element. In fact, the precise optical properties of the bulk material in an optical element are virtually certain to be a function of both the material itself, and of the process which produced the part. Some materials, notably styrene and butyrate resins, are crystalline to some degree, and therefore inherently birefringent. Birefringence may develop in amorphous materials, though, if the injection mold and process parameters are not optimized to prevent this occurrence. Likewise, the bulk scatter properties of a molded optical element are a function of the inherent properties of the material, but are also strongly related to the cleanliness of the processing and the heat history of the finished part.

3.6 OPTICAL DESIGN

Design Strategy

Virtually all optical design techniques which have evolved for use with glass materials work well with polymer optics. Ray-tracing formulery, optimization approaches, and fundamental optical construction principles are equally suitable for glass or plastic. The generalized approach to optical design with polymeric materials should be strongly medium-oriented, though. That is, every

effort must be made to capitalize upon the design flexibility which the materials and manufacturing processes afford. Integration of form and function should be relentlessly pursued, since mechanical features may be molded integral with the optics to reduce the metal part count and assembly labor content in many systems.

Aberration Control

The basic optical design task normally entails the simultaneous satisfaction of several first-order constraints, the correction of the monochromatic aberrations, and the control of the chromatic variation of both first-order quantities and higher-order aberrations. It is well known that management of the Petzval sum, while maintaining control of the chromatic defects, may be the most difficult aspect of this effort.^{17,18} It is also widely recognized that the choice of optical materials is key to success. While the available polymer choices cover a wide range of Abbe values, ensuring that achromatization may be accomplished in an all-polymer system, the refractive index values for these materials are not well-positioned on the “glass” map to permit low Petzval sums to be easily achieved.

Material Selection

Simultaneous correction of the Petzval sum and the first-order chromatic aberration may, however, be nicely accomplished if the materials employed possess similar ratios of Abbe number to central refractive index. This implies that the *best* material combinations (involving polymers) should probably include an optical glass. Also implied is the fact that these hybrid material combinations may be inherently superior (in this respect) to all-glass combinations. Ideally, the chosen materials should be well-separated (in Abbe value) on the glass map, so that the component powers required for achromatization do not become unduly high. This condition is satisfied most completely with polymers which lie in the TiF sector of the glass map, coupled with glasses of the LaK, LaF, and LaSF families.

Most lens designers would prefer to utilize high-refractive-index materials almost exclusively in their work. Optical power must be generated in order to form images, and because the combination of optical surface curvature and refractive index creates this refractive power, these two variables may be traded in the lens design process. Since it is well known that curvature generates aberration more readily than does a refractive index discontinuity, one generally prefers to achieve a specified amount of refractive power through the use of low curvature and high-refractive index. From this perspective, the polymers are at a distinct disadvantage, most of them being low-index materials.

Aspheric Surfaces

An offsetting consideration in the use of polymeric optical materials is the freedom to employ nonspherical surfaces. While these may be awkward (and very expensive) to produce in glass, the replication processes which create plastic optical parts do not differentiate between spherical and nonspherical surfaces.

As any lens designer can attest, the flexibility that aspheric surfaces make available is quite remarkable.^{19,20} Spherical surfaces, while convenient to manufacture by grinding and polishing, may generate substantial amounts of high-order aberration if used in any optical geometry which departs significantly from the aplanatic condition. These high-order aberrations are often somewhat insensitive to substantial changes in the optical prescription. Thus, profound configurational alterations may be necessary to effect a reduction in these image defects.

On the other hand, the ability to utilize surface shapes which are more complex than simple spheres permits these high-order aberration components to be moderated at their point of origin, which may in turn reduce the amount of “transferred aberration” imparted to surfaces downstream in the optical train. In a multielement optical system, especially one employing cascaded aspheric surfaces, the required imagery performance may be achieved using fewer total elements. And due to

the fact that the surface aberration contributions are diminished, the sensitivity to positioning errors may also be reduced, with the result that an aspheric optical system may actually be more forgiving to manufacture than its spherical counterpart.

In practice, the use of aspheric surfaces in polymer optical elements appears to more than compensate for the handicap imposed by low-refractive index values. Using aspheric surfaces, it is possible to bend, if not break, many of the rules which limit design with spherical surfaces. Aspherics create extra leverage to deal with the monochromatic aberrations, and with the chromatic variation of these image defects. A designer experienced with aspherics, given a capable set of software tools, can frequently create optical constructions which deliver high performance, despite the fact that they appear odd to those accustomed to the more “classical” spherical surface configurations. Quite often, unfavorable design constraints such as an inconvenient aperture stop location, may be handled with less difficulty using aspherics.

Athermalization

The thermal behavior of the polymers, mentioned previously, may cast a shadow upon some applications where the temperature is expected to vary over a significant range, but the focal surface location must be fixed in space. In such cases, the variation of refractive index usually accounts for the largest share of the variation, with the dimensional changes playing a secondary role. In such situations, the thermally induced excursions of the focal surface may be compensated by modeling these functions and designing mechanical spacers of the proper material to stabilize the detector/image location.

Alternatively, the optical system may be designed to exhibit inherently athermal behavior over the operational temperature range.²¹ Unfortunately, this is not strictly possible using only polymeric materials, as the thermo-optic coefficients display so little variation among themselves that the component powers would be absurdly high.

In combination with one or more glass elements, however, very nicely athermalized design solutions may be obtained with polymer elements.²² Athermal designs may be generated by modeling the optical system in multiconfiguration mode in the lens design software, much as one would develop a zoom lens. The parameters to be “zoomed” in this case are the refractive indices at two or more temperatures within the operating range. The resulting designs frequently concentrate most of the refractive power in the glass elements, with the polymer elements functioning to achieve achromatism and control of the monochromatic aberrations. See also Chap. 8, “Thermal Compensation Techniques,” by Philip J. Rogers and Michael Roberts in Vol. II of this *Handbook*.

Processing Considerations

In much the same manner that optical design with polymer materials is different from optical design with glass, the treatment of the fabrication and assembly issues are also quite different matters. The major issues requiring examination are those related to the materials themselves. While it is possible to characterize the glass for an optical system with complete certainty prior to performing any fabrication operations, with polymers, one’s knowledge of the starting materials is only a rough indication of the properties of the finished optical parts.

When optical properties data are offered by the polymer supplier, it should be realized that these numbers apply *only* to measurement samples which have been predried to specification, have experienced a specified residence time in the extrusion barrel under specific temperature conditions, have been injected into the mold cavity at specific rates and pressures, and so on. Consequently, it is unlikely that the refractive properties of a polymer element will conform closely to catalog values (if such values are indeed supplied). Moreover, homogeneity, bubble content, scatter properties, and so on, are all process-dependent. So while the melt sheets may fix the optical properties of glass materials very precisely, the uncertainty associated with the polymers demands that refractive variations be allocated a significant portion of the fabrication and assembly error budget.

Manufacturing Error Budget

Other constructional parameters, conversely, may be implemented with great precision and repeatability in plastic. The molding process, executed by means of modern equipment, can be exceedingly stable. Vertex thickness, curvature, and wedge may often be maintained to a greater level of precision, with greater economy than is possible with glass fabrication technology. It is not unusual to see part-to-part variations in vertex thickness of less than 0.01 to 0.02 mm over a run of thousands of parts from a single cavity.

Multiple Cavities

The economic appeal of injection molding is the ability to create several parts in one molding cycle. In a multicavity scenario, the parts from different cavities may exhibit some small dimensional differences, depending upon the level of sophistication of the tool design and the quality of its construction. Cavity variations in axial thickness, fortunately, may be permanently minimized by implementing small tooling adjustments after the mold has been exercised. Consequently, part thickness variation rarely consumes a significant fraction of the constructional error budget.

Dimensional Variations

Surface radii, like axial thickness, may be replicated with great repeatability *if* the molding process is adjusted to a stable optimum. Radius errors, if they are present, are usually attributable to incorrect predictions of shrinkage, and may be biased out by correcting the radii of the mold inserts. Thus, the consistency of surface radii achievable with glass may often be equaled in plastic. Thus, radius errors, as well as axial thickness errors, frequently constitute a small portion of the polymer optics manufacturing error budget.

Element wedge, like axial thickness, may be minimized by careful attention to precision in the tool design and construction. It is quite possible to achieve edge-to-edge thickness variations of less than 0.01 mm in molded plastic lenses. With polymer lenses, the azimuthal location of the part gate may be used, if necessary, to define rotational orientation of the element in the optical train. Consequently, rotational alignment of plastic optical parts may be easily indexed.

Optical Figure Variations

Control of optical figure quality is obviously key to the successful execution of a good optical design. In glass, achievement of subfringe of figure conformance is accomplished routinely, albeit at some cost penalty. In polymeric optics, the nonlinear shrinkage, surface tension, and other processing-related effects cause surface figure errors to scale with part size, sometimes at a rate proportional to some exponent of diameter. This limits the practical size range for polymeric optics, although capable optics molders may routinely produce elements in the 10-mm-diameter range to subfringe accuracy.²³

On one hand, it can probably be stated that processing-induced variations in properties, and a dearth of dependable optical data, preclude any serious discussions of such things as apochromatic polymeric optics, or of large polymeric optics operating at the diffraction limit. On the other hand, the consistency with which some dimensional parameters may be reproduced in quantity, and the design freedom and flexibility afforded by molded aspherics, make possible the satisfaction of some design requirements which would be out of range for conventional glass optics.^{24,25}

Specification

Given the fact that the guidelines and restrictions for design and implementation are very different for glass and polymeric optics, it is not surprising that the approach to specification of polymer optical parts and systems should be tailored to the materials and processes of polymer optics. Attempts

to convert a glass optics concept to plastic are frequently unsuccessful if the translation overlooks the fundamental themes of the molding and tooling technologies involved. Much as optimum tube and solid-state electrical circuit topologies should be significantly different, so must the execution of a conceptual optical system, depending upon whether glass or polymer material is the medium.

It follows naturally that manufacturing drawings for polymer optics may contain annotations which seem unfamiliar to those versed in glass optics manufacture. Furthermore, some specifications which are universally present on all glass optics drawings may be conspicuously absent from a polymer optics print.

For example, thermal and cosmetic damage considerations preclude the use of the familiar test glasses in the certification of polymeric optics. Figure conformance, then, need only be specified in “irregularity” or asphericity terms, since the alternative method, use of a noncontacting interferometer, implies that the focus error (*fringe power* in test plate language) will be automatically removed in the adjustment of the test setup.

References to ground surfaces may be omitted from polymer optics drawings, since no such operation takes place. Discussions of “chips” inside the clear aperture, staining, and the like are also superfluous. Beauty defect specifications do apply, although such imperfections are almost always present in every sample from a specific cavity, probably implying the need to rework a master surface.

In general, the lexicon of optics, and that of the molding industry, do not overlap to a great extent. Molding terms like *flash* and *splay* are meaningless to most optical engineers. Those endeavoring to create a sophisticated polymeric optical system, anticipating a successful outcome, are advised to devote some time to the study of molding, and to discussions with the few experts in the arcane field of optics molding, before releasing a drawing package which may be unintelligible to or misunderstood by the vendor.

3.7 PROCESSING

Casting

As mentioned above, polymeric optics may be produced by any of several processes. These include fabrication, transfer replication, casting, compression molding, injection molding, and some combinations of the aforementioned.²⁶ The earliest polymeric optical parts were probably produced by fabrication or precipitation from solution. Large military tank prisms have been made by both processes. In the latter case, the polymer (typically PMMA) was dissolved, and the solvent then evaporated to produce a residue of polymer material in the shape of the mold—a very inefficient technique indeed.

Many of the polymers may be fabricated by cutting, grinding, and polishing, much as one would deal with glass materials. The thermoset resin tradenamed CR-39 (poly-diallylglycol) was formulated specifically to be processed using the same techniques and materials as those used to fabricate glass optics. And this material does indeed produce good results when processed in this manner. It is used extensively in the ophthalmic industry to produce spectacle lenses. The processing, in fact, usually involves casting the thermoset resin to create a lens blank which emerges from the mold with the optical surfaces polished to final form. More conventional fabrication techniques may then be utilized to edge the lens, or perhaps to add a bifocal portion.

Abrasive Forming

Unfortunately, the softness of most of the polymers, coupled with their poor thermal conductivity, complicates the achievement of a truly high quality polish using conventional methods. Even in the case of CR-39, which is relatively hard for a polymer material, some amount of “orange peel” in the polished surface seems unavoidable. Many thermoplastics, most of them softer than CR-39, may be conventionally ground and polished to give the appearance of an acceptable optical surface. Closer

examination, however, reveals surface microstructure which probably does not fall within the standards normally associated with precision optics. Nonetheless, fabrication of optical elements from large slabs of plastic is often the only viable approach to the creation of large, lightweight refractive lenses, especially if cost is an issue.

In general, the harder, more brittle polymers produce better optical surfaces when ground and polished. PMMA and others seem to fare better than, say, polycarbonate, which is quite soft, exhibits considerable elongation at the mechanical yield point, but is in great demand due to its impact resistance.

Single-Point Turning

An alternative approach to fabrication, one that is especially useful for the production of aspheric surfaces, is the computer numerical control (CNC) lathe turning of the bulk material using a carefully shaped and polished tool bit of single-crystal diamond or cubic boron nitride. See also Chap. 10, "Fabrication of Optics by Diamond Turning," by Richard L. Rhorer and Chris J. Evans in Vol. II of this *Handbook*. The lathe required to produce a good result is an exceedingly high precision tool, having vibration isolation, temperature control, hydrostatic or air bearings, and so on. On the best substrate materials (PMMA is again a good candidate), very good microroughness qualities may be achieved. With other materials, a somewhat gummy character (once more, polycarbonate comes to mind) may result in microscopic tearing of the surface, and the expected scatter of the incident radiation.

The diamond-turning process is often applied in conjunction with other techniques in order to speed progress and reduce cost. Parts which would be too large or too thick for economical stand-alone injection molding are frequently produced more efficiently by diamond-turning injection molded, stress-relieved preforms, which require minimal material removal and lathe time for finishing. Postpolishing, asymmetric edging, and other postoperations may be performed as necessary to create the finished part. Optics for illumination and TV projection applications are often produced by some combination of these techniques. Given the fact that the technology in most widespread use for the production of plastic optics involves some form of molding (a front-loaded process, where cost is concerned), diamond-turning is often the preferred production method for short production runs and prototype quantities.

Compression Molding

Most high-volume polymeric optics programs employ a manufacturing technology involving some form of molding to produce the optical surfaces, if not the entire finished part.²⁷ Of the two most widely used approaches, compression molding is best suited to the creation of large parts having a thin cross section. In general, any optical surface possessing relief structure having high spatial frequency is not amenable to injection molding, due to the difficulty of forcing the material through the cavity, and due to the fact that the relief structure in the mold disrupts the flow of the polymer. In addition, the relief structure in the master surfaces may be quite delicate, and prone to damage at the high pressures often present in the mold cavity.

The compression molding process is capable of producing results at considerably lower surface pressure than injection molding, and as long as the amount of material to be formed is small, this molding technology can replicate fine structure and sharp edge contours with amazing fidelity. Since the platens of a compression molding press are normally heated using steam or electrical heaters, most compression molded parts are designed to be executed in polymers having a relatively low temperature softening point, and materials like polyethersulfone are rarely utilized.

Injection Molding

Optical parts having somewhat smaller dimensions may be better suited to production by the injection molding process.²⁸ This is probably the preferred polymer manufacturing technology for optical elements having a diameter smaller than 0.1 m and a thickness not greater than 3 cm. Not only do the economics favor this approach in high production volume, but if properly applied, superior optical surfaces may be produced.²⁹

It should be kept firmly in mind that the basic injection molded process (as it is known to most practitioners) requires a great deal of refinement and enhancement in order to produce credible optical parts.³⁰ Unfortunately, very few molders possess either the molding know-how, or the testing and measurement sophistication to do the job correctly. Given a supply of quality polymer material, the molding machine itself must be properly configured and qualified. Relatively new machinery is a must. The platens to which the mold halves are mounted must be very rigid and properly aligned. And this alignment must be maintainable on a shot-to-shot basis for long periods. The screw and barrel must be kept scrupulously clean, and must be carefully cleaned and purged when switching materials. The shot capacity, in ideal circumstances, should be more carefully matched to the part volume than for non-optical parts. The process control computer must be an inordinately flexible and accurate device, able to profile and servo a number of operational functions that might be of little importance if the molded part were not optical in nature.³¹

Since much of the heating of the injected polymer resin occurs as a result of physical shear and compression (due to a variable pitch screw), the selection of these machine characteristics is critical to success. In addition, the energy supplied to the machine barrel by external electric heaters must be controlled with more care than in standard industrial applications. A failure of a single heater, or a failure of one of the thermal measurement devices which close that servo loop, may result in many defective parts.

Vendor Selection

The injection mold itself requires special attention in both design and execution in order to produce state-of-the-art molded lenses. A number of closely held “trade tricks” normally characterize a mold designed to produce optical parts, and these subtle variations must be implemented with considerably greater accuracy than is normally necessary in ordinary molding. The mold and molding machine are often designed to operate more symbiotically than would be the case in producing non-optical parts. Control of the mold temperature and temperature gradients is extremely critical, as is the control bandwidth of those temperatures and the temperature of the molding room itself. The most important conclusion to be drawn from the preceding paragraphs is that *the molding vendor for polymer optical parts must be selected with great care*. A molding shop, no matter how sophisticated and experienced with medical parts, precision parts for electronics, and so on, will probably consume much time and many dollars before conceding defeat with optical parts.

Although success in molding optical elements is a strong function of equipment, process control, and engineering acumen, *attention to detail in the optical and mechanical design phases will consistently reduce the overall difficulty of manufacturing these items*. An awareness of the basic principles of injection molding procedures and materials is very helpful here, but it is necessary to be aware that, in the optical domain, we are dealing with micrometer-scale deformations in the optical surfaces. Thus, errors or oversights in design and/or molding technique which would totally escape notice in conventional parts can easily create scrap optics.

Geometry Considerations

The lens design effort, for best results, must be guided by an awareness of the basic physics of creating an injection molded part, and of the impact of part cross section, edge configuration, asymmetry, and so on. In general, any lens having refractive power will possess a varying thickness across its diameter. Unfortunately, meniscus-shaped elements may mold best due to the more uniform nature of the heat transfer from the bulk.³² Positive-powered lens elements will naturally shrink toward their center of mass as they cool, and it may be difficult to fill the mold cavity efficiently if the edge cross section is only a small fraction of the center thickness.

Negative lenses, on the other hand, tend to fill in the outer zones more readily, since the thinner portion of the section (the center) tends to obstruct flow directly across the piece from the part gate. In extreme circumstances, it is possible that the outer zones of the lens element will be first to

fill, trapping gases in the center, forming an obvious *sink* in molding terminology. Parts designed with molded-in bores may exhibit the ‘*weld-line*’ phenomenon, which is a visible line in the part where the flow front of the molten plastic is divided by the mold cavity obstruction forming the bore. In the case of both negative and positive lens elements, it is good policy to avoid element forms wherein the center-to-edge thickness ratio exceeds three for positive elements, or is smaller than 0.3 for negative elements.

Shrinkage

Surface-tension effects may play a significant role in the accuracy to which a precision optical surface may be molded.^{33,34} Particularly in areas of the part where the ratio of surface area/volume is locally high (corners, edges), surface tension may create nonuniform shrinkage which propagates inward into the clear aperture, resulting in an edge rollback condition similar to that which is familiar to glass opticians. Surface tension and volumetric shrinkage may, however, actually aid in the production of accurate surfaces. Strongly curved surfaces are frequently easier to mold to interferometric tolerances than those having little or no curvature. These phenomena provide motivation to oversize optical elements, if possible, to a dimension considerably beyond the clear apertures. A buffer region, or an integrally molded flange provides the additional benefit of harmlessly absorbing optical inhomogeneities which typically form near the injection gate. Figure 2 depicts several optical element forms exhibiting favorable (*a–e*) and unfavorable (*f–j*) molding geometries. In some cases, a process combining injection and compression molding may be used to improve optical figure quality. Several variants of this hybrid process are in use worldwide, with some injection molding presses being specifically fitted at the factory to implement this procedure.³⁵

Mechanical Assembly

In order to appreciate fully the design flexibility and cost-saving potential of polymer optics, it is necessary to modify one’s approach to both optical and mechanical design. A fully optimized polymeric optical system not only makes use of aspheric technology and integrally molded features in the optical elements, but embodies an extension of this design philosophy into the lens housing concept and assembly strategy. These issues should ideally be considered in concert from the very

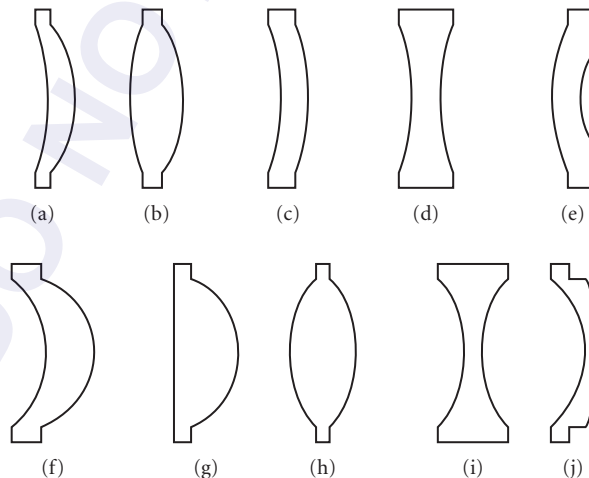


FIGURE 2 Some polymer lens element cross-sectional configurations.

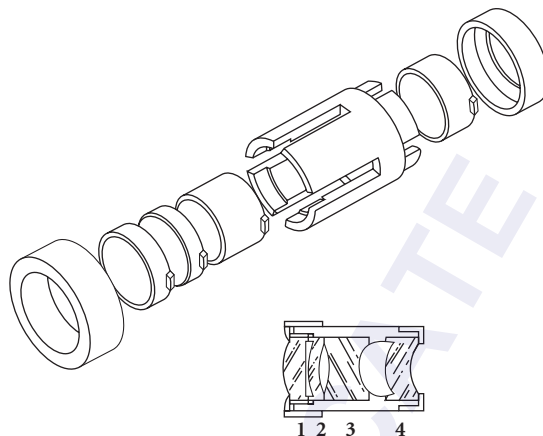


FIGURE 3 Collet-type lens housing. Joining by ultrasound eliminates the possibility of pinching lens elements.

beginning, so that design progress in one aspect does not preclude parallel innovation in other facets of the development process.

It is important to resist the urge to emulate glass-based optomechanical design approaches, since the polymer technology permits design features to be implemented which would be prohibitively expensive (or even impossible) in metal and glass. Spacers required to separate elements may be molded as part of the elements (Fig. 2), reducing the metal part total and simplifying assembly operations. See also Chap. 6, "Mounting Optical Components," by Paul R. Yoder, Jr. in Vol II of this *Handbook*. Housings may be configurations which would be either improbable or unmanufacturable using machine tool technology. The collet-and-cap design shown in Fig. 3 is one such example. Joining might be accomplished by ultrasonic bonding. The clamshell concept shown in Fig. 4 may be designed so that the two halves of the housing are actually the same part, aligned by molded-in locating pins. Joining might be performed by a simple slip-on C ring.

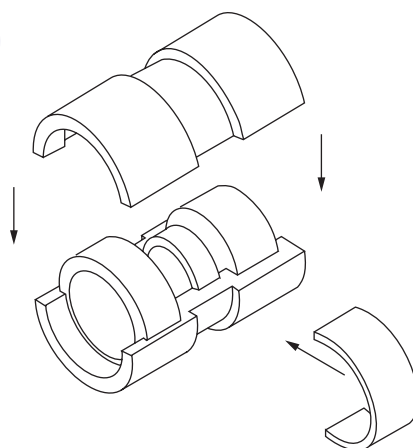


FIGURE 4 Clamshell lens housing. Possibility of lens jamming during assembly is minimized.

Whereas lens assemblies in glass and metal are normally completed by seating threaded retaining rings, their plastic counterparts may be joined by snap-together pieces, ultrasonic bonding, ultraviolet-curing epoxies, expansion C rings, or even solvent bonding.³⁶ Solvent bonding is dangerous, however, since the errant vapors may actually attack the polymer optical surfaces.

Following the basic polymeric optics philosophy, the lens element containment and assembly approach should probably not even consider the disassembly option in the event of a problem. In order to maximize assembly precision and minimize unit cost, the design of the lens cell should evolve alongside that of the optical system, and this cell should be visualized as an extension of a fixture conceived to minimize the labor content of the assembly.

An in-depth treatment of optical mold design and tooling technology is obviously beyond the scope of this discussion. Many of the methods and procedures parallel those in use in the molding industry at large. However, a number of subtle and very important detail differences do exist, and these are not extensively documented in the literature. Issues having to do with metallurgy, heat treatment, chemical passivation, metal polishing, and so on, have little to do with the actual design and engineering of a polymeric optical system. In a modern tool design exercise, though, the flow behavior of the polymer material in the mold, and the thermal behavior of that mold, are carefully modeled in multinode fashion, so that part quality may be maximized, and cycle time minimized.³⁷ A nodding awareness of these methods, and the underlying physics, may be helpful to the person responsible for the engineering of the polymer optical system.

Testing and Qualification

In the process of implementing any optical system design, the matters of testing and certification become key issues. In molded optics, the master surfaces, whose shapes are ultimately transferred to the polymer optical parts, must be measured and documented. A convenient testing procedure for the optical elements replicated from these surfaces must likewise be contrived, in order to optimize the molding process and ensure that the finished assembly will perform to specification. The performance of that assembly must itself be verified, and any disparities from specification diagnosed.

In general, mechanical dimensions of the polymer parts may be verified by common inspection tools and techniques used in the glass optics realm. The possibility of inflicting surface damage, however, dictates that noncontact interferometric techniques be used in lieu of test glasses for optical figure diagnosis. This is a straightforward matter in the case of spherical surfaces, but requires some extra effort in the case of aspherics. See also Chap. 13, "Optical Testing," by Daniel Malacara-Hernández in Vol. II of this *Handbook*.

Obviously, aspheric master surfaces must be scrupulously checked and documented, lest the molder struggle in vain to replicate a contour which is inherently incorrect. The verification of the aspheric masters and their molded counterparts may be accomplished in a variety of ways. Mechanical gauging, if properly implemented, works well, but provides reliable information through only one azimuthal section of the part. Measurement at a sufficient number of points to detect astigmatism is awkward, very time consuming, and expensive. And this is not exactly consistent with the spirit of polymer optics.

Null Optics

An optical *null corrector* permits the aspheric surface to be viewed in its entirety by the interferometer as if it were a simple spherical surface.^{38,39} This is a rapid and convenient procedure. The null optics consist of very accurately manufactured (and precisely aligned) spherical glass elements designed to introduce aberration in an amount equal to, but of opposite sign from, that of the tested aspheric. Thus, interspersing this device permits aspherics to be viewed as if they were spherical. Since there exists no simple independent test of the null compensator, one must depend heavily upon the computed predictions of correction and upon the skill of the fabricator of the corrective optics. See also Chap. 14, "Use of Computer-Generated Holograms in Optical Testing," by Katherine Creath and James C. Wyant in Vol. II of this *Handbook*.

The concept of greatest importance regarding the use of aspheric surfaces is that *successful production of the total system is cast into considerable doubt if a surface is present which is not amenable to convenient testing*. While some aspheric optics may be nulled fairly easily, those which appear in polymer optical systems are frequently strong, exhibiting significant high-order derivatives. If the base curves are strong, especially strongly convex, there may exist no practical geometry in which to create a nulling optical system. And if a favorable geometry does exist, several optical elements may be necessary to effect adequate correction. One can easily approach a practical limit in this situation, since the manufacturing and assembly tolerances of the cascaded spherical elements may themselves (in superposition) exceed the theoretical correction requirement. The bottom line is that one should not proceed with cell design, or any other hardware design and construction, until the aspheric testing issues have been completely resolved.

3.8 COATINGS

Reflective Coatings

Given the fact that optical polymers exhibit specular properties similar to those of glass, it is not surprising that optical coatings are often necessary in polymeric optical systems. The coatings deposited upon polymer substrates fall mostly into four general categories. These include coatings to improve reflectivity, to suppress specular reflection, to improve abrasion resistance, and to retard accumulation of electrostatic charge.

Reflective coatings may be applied by solution plating, or by vacuum-deposition. These are most often metallic coatings, usually aluminum if vacuum-deposited, and normally chromium if applied by plating. The abrasion resistance and general durability of such coatings is rather poor, and susceptibility to oxidation quite high, if no protective coating is applied over the metal film. In some applications, especially involving vacuum deposition, the overcoat may be a thin dielectric layer, deposited during the same process which applies the metal film. If the reflective coating has been applied by plating, the overcoat may be an organic material, perhaps lacquer, and may be deposited separately by spraying or dipping. Not surprisingly, the quality of a surface so treated will be poor by optical standards, and probably suitable only for toy or similar applications.

Antireflection Coatings

Antireflection coatings are frequently utilized on polymer substrates, and may consist of a single layer or a rudimentary multilayer stack yielding better reflection-suppression performance. Due to the stringent requirements for control of the layer thicknesses, such coating formulations may be successfully deposited only in high-vacuum conditions, and only if temperatures in the chamber remain well below the service temperature of the substrate material. Elevated temperatures, necessary for baking the coatings to achieve good adhesion and abrasion resistance, may drive off plasticizing agents, limiting the "hardness" of the chamber vacuum. Such temperatures can ultimately soften the optical elements, so that their optical figure qualities are compromised. Relatively recent developments in the area of ion beam-assisted deposition have made possible improvements in the durability of coatings on polymer materials without having to resort to significantly elevated chamber temperatures.⁴⁰ See also Chap. 7, "Optical Properties of Films and Coatings," by Jerzy A. Dobrowolski in this volume of *Handbook*.

Antiabrasion Coatings

In general, many polymeric optical systems which could benefit from application of coatings are left uncoated. This happens because the expense incurred in cleaning, loading, coating, unloading, and inspecting the optical elements may often exceed that of molding the part itself. Some optics, particularly those intended for ophthalmic applications, are constantly exposed to abuse by abrasion,

and must be protected, cost notwithstanding. Antiabrasion coatings intended to provide immunity to scratching may be of inorganic materials (normally vacuum-deposited), or may be organic formulations.⁴¹

Inorganic antiabrasion coatings may be similar to those used for simple antireflection requirements, except that they may be deposited in thicknesses which amount to several quarter-wave-lengths. The practical thickness is usually limited by internal stress buildup, and by differential thermal expansion between coating and substrate. In general, the inorganic coatings derive their effectiveness by virtue of their hardness, and provide protection only superficially, since sufficient pressure will collapse the underlying substrate, allowing the coating to fracture.

Organic coatings for abrasion resistance normally derive their effectiveness from reduction of the surface frictional coefficient, thereby minimizing the opportunity for a hard contaminant to gain the purchase required to initiate a scratch. These coatings are often applied by dipping, spraying, or spinning. Coatings thus deposited usually destroy the smoothness which is required if the piece is to be qualified as a precision optical element.

Antistatic Coatings

Coatings applied for the purpose of immunization against abrasion, or suppression of specular reflection, often provide a secondary benefit. They may improve the electrical conductivity of the host surface, thus promoting the dissipation of surface static charge, and the accumulation of oppositely charged contaminants. In circumstances where antireflection or antiabrasion coating costs cannot be justified, chemical treatments may be applied which increase conductivity. These materials typically leave a residue sufficiently thin that they are undetectable, even in interferometric testing.

3.9 REFERENCES

1. R. M. Altman and J. D. Lytle, "Optical Design Techniques for Polymer Optics," *S.P.I.E. Proc.* **237**:380–385 (1980).
2. C. Teyssier and C. Tribastone, *Lasers & Optronics* **Dec**:50–53 (1990).
3. PPG Ind., Inc., Tech Bulletin-CR-39.
4. H. Dislich, *Angew. Chem. Int. Ed. Engl.* **18**:49–59 (1979).
5. H. D. Wolpert, *Photonics Spectra* **Feb**:68–71 (1983).
6. Plastics Desk Top Data Bank, pp. 803–837 (1986).
7. *Modern Plastics Encyclopedia—Eng. Data Bank*, McGraw-Hill, New York, 1977, pp. 453–708.
8. *Plastics Technology Manufacturing Hdbk. and Buyer's Guide*, 1986, pp. 358–740.
9. *C.R.C. Hdbk. of Laser Science & Technology*, vol. IV, pp. 85–91.
10. *Encyclopedia of Polymer Science and Technology*, vol. 1, John Wiley & Sons, 1976.
11. *Space Materials Hdbk*, Lockheed Missiles and Space Corp., 1975.
12. J. D. Lytle, G. W. Wilkerson, and J. G. Jaramillo, *Appl. Opt.* **18**:1842–1846 (1979).
13. D. C. Smith, Alpert et al., *N.R.L. Report* 3924, 1951.
14. R. E. Kagarise and L. A. Weinberger, *N.R.L. Report* 4369, 1954.
15. Rohm & Haas Product Bulletin-PL 612d, 1979.
16. *Catalogue of Optical Glasses*, Schott Glass Technologies, 1989.
17. J. D. Lytle, *S.P.I.E. Proc.* **1354**:388–394 (1990).
18. J. Hoogland, *S.P.I.E. Proc.* **237**:216–221 (1980).
19. A. Osawa et al., *S.P.I.E. Proc.* **1354**:337–343 (1990).
20. E. I. Betensky, *S.P.I.E. Proc.* **1354**:663–668 (1990).

21. L. R. Estelle, *S.P.I.E. Proc.* **237**:392–401 (1980).
22. K. Straw, *S.P.I.E. Proc.* **237**:386–391 (1980).
23. M. Muranaka, M. Takagi, and T. Maruyama, *S.P.I.E. Proc.* **896**:123–131 (1988).
24. J. D. Lytle, *S.P.I.E. Proc.* **181**:93–102 (1979).
25. A. L. Palmer, “Practical Design Considerations for Polymer Optical Systems,” *S.P.I.E. Proc.* **306** (1981).
26. D. F. Horne, *Optical Production Technology*, 2d. ed., Adam Hilger, Ltd., 1983 pp. 167–170.
27. J. R. Egger, *S.P.I.E. Proc.* **193**:63–69 (1979).
28. R. Benjamin, *Plastics Design and Processing* **19**(3):39–49 (1979).
29. R. F. Weeks, *Optical Workshop Notebook*, vol. I., O.S.A., 1974–1975, sect. XVII.
30. D. F. Horne, “Lens Mechanism Technology,” Crane, Russack & Co., New York, 1975.
31. J. Sneller, *Modern Plastics Intl.* **11**(6):30–33 (1981).
32. J. D. Lytle, “Workshop on Optical Fabrication and Testing,” *Tech. Digest*, O.S.A., pp. 54–57 (1980).
33. E. C. Bernhardt and G. Bertacchi, *Plastics Technology* **Jan**:81–85 (1986).
34. G. R. Smoluk, *Plastics Engineering* **Jul**:107 (1966).
35. *Plastics News*, Aug. 1989, pp. 8–9.
36. *Plastics World*, July 1979, p. 34.
37. M. H. Naitove, *Plastics Technology* **Apr.** (1984).
38. D. Malacara (ed.), *Optical Shop Testing*, John Wiley & Sons, New York, 1978, chaps. 9, 14.
39. G. W. Hopkins and R. N. Shagam, *Appl. Opt.* **16**(10):2602 (1977).
40. J. D. Rancourt, *Optical Thin Films-User’s Hdbk.*, MacMillan, New York, 1987, pp. 197–199.
41. J. W. Prane, *Polymer News* **6**(4):178–181 (1980).

This page intentionally left blank.

DO NOT DUPLICATE

PROPERTIES OF METALS

Roger A. Paquin

*Advanced Materials Consultant
Tucson, Arizona, and
Optical Sciences Center
University of Arizona
Tucson, Arizona*

4.1 GLOSSARY

a	absorptance, absorptivity
a	plate radius (m)
B	support condition
C_{ij}	elastic stiffness constants (N/m ²)
C_p	specific heat (J/kg K)
CTE	coefficient of thermal expansion (K ⁻¹)
D	flexural rigidity (N m ²)
D	thermal diffusivity (m ² /s)
E	elastic modulus (Young's) (N/m ²)
\mathbf{E}	electromagnetic wave vector (J)
e	electron charge (C)
E_o	amplitude of electromagnetic wave at $x = 0$ (J)
G	load factor (N/kg)
G	shear modulus, modulus of rigidity (N/m ²)
g	acceleration due to gravity (m/s ²)
I	light intensity in medium (W/m ²)
i	$(-1)^{1/2}$
I_o	light intensity at interface (W/m ²)
I_o	section moment of inertia (m ⁴ /m)
K	bulk modulus (N/m ²)
k	extinction coefficient
k	thermal conductivity (W/mK)
L	length (m)
M	materials parameter (kg/Nm)
m	electron mass (kg)
N	complex index of refraction

N	number of dipoles per unit volume (m^{-3})
n	index of refraction
P	plate size (m^4)
q	load (N/m^2)
r	Reflectance, reflectivity
R_f	intensity reflection coefficient
S	structural efficiency (m^{-2})
T	temperature (K)
t	time (s)
t	transmittance, transmissivity
V_0	volume per unit area of surface (m)
x	distance (m)
α	coefficient of thermal expansion (K^{-1})
α	absorption coefficient (m^{-1})
β	deflection coefficient
β	dynamic deflection coefficient
Γ	damping constant
δ	skin depth (nm)
δ	deflection (m)
δ_{DYN}	dynamic deflection (m)
ϵ	emittance, emissivity (W/m^2)
ϵ	complex dielectric constant
ϵ_0	permittivity of free space (F/m)
ϵ_1	real part of dielectric constant
ϵ_2	imaginary part of dielectric constant
θ	angular acceleration (s^2)
λ	wavelength (m)
λ_0	wavelength in vacuum (m)
μ	magnetic susceptibility (H/m)
ν	frequency (s^{-1})
ν	Poisson' ratio
ρ	mass density (kg/m^3)
σ	conductivity (S/m)
ω	radian frequency (s^{-1})

4.2 INTRODUCTION

Metals are commonly used in optical systems in three forms: (1) structures, (2) mirrors, and (3) optical thin films. In this article, properties are given for metal mirror substrate and structural materials used in modern optical systems. Many other materials have not been included due to their limited applicability. Metal film properties are discussed in the context of thick films (claddings) rather than optical thin films that are covered in Chap. 7, "Optical Properties of Films and Coatings." Since mirrors are structural elements, the structural properties are equally important as the optical properties to the designer of an optical system. Therefore, the properties addressed here include physical, mechanical, and thermal properties in addition to optical properties. Mechanical and thermal properties of silicon (Si) and silicon carbide (SiC) are included, but not their optical properties since they are given in Chap. 5, "Optical Properties of Semiconductors."

After brief discussions of optical properties, mirror design, and dimensional stability, curves and tables of properties are presented, as a function of temperature and wavelength, where available. For more complete discussions or listings, the reader should consult the references and/or one of the available databases.¹⁻³ A concise theoretical overview of the physical properties of materials is given by Lines.⁴

Nomenclature

The symbols and units used in this subsection are consistent with usage in other sections of this *Handbook* although there are some unavoidable duplications in the usage of symbols between categories of optical, physical, thermal, and mechanical properties. Definitions of symbols with the appropriate units are contained in the table at the beginning of this chapter.

Optical Properties

The definitions for optical properties given in this section are primarily in the geometric optics realm and do not go into the depth considered in many texts dealing with optical properties of solids.⁵⁻⁸

There is obviously a thickness continuum between thin films and bulk, but for this presentation, bulk is considered to be any thickness of material that has bulk properties. Typically, thin films have lower density, thermal conductivity, and refractive index than bulk; however, current deposition techniques are narrowing the differences. Optical properties of thin films are presented only when bulk properties have not been found in the literature.

The interaction between light and metals takes place between the optical electric field and the conduction band electrons of the metal.⁹ Some of the light energy can be transferred to the lattice by collisions in the form of heat. The optical properties of metals are normally characterized by the two optical constants: index of refraction n and extinction coefficient k that make up the complex refractive index N where:

$$N = n + ik \quad (1)$$

The refractive index is defined as the ratio of phase velocity of light in vacuum to the phase velocity of light in the medium. The extinction coefficient is related to the exponential decay of the wave as it passes through the medium. Note, however, that these “constants” vary with wavelength and temperature. The expression for an electromagnetic wave in an absorbing medium contains both of these parameters:

$$E = E_0 e^{-2\pi kx/\lambda_0} e^{-i(2\pi nx/\lambda_0 - \omega t)} \quad (2)$$

where E_0 is the amplitude of the wave measured at the point $x = 0$ in the medium, E is the instantaneous value of the electric vector measured at a distance x from the first point and at some time t , ω is the angular frequency of the source, and λ_0 is the wavelength in vacuum.

The absorption coefficient α is related to the extinction coefficient by:

$$\alpha = 4k/\lambda_0 \quad (3)$$

and for the general case, the absorption coefficient also appears in the absorption equation:

$$I = I_0 e^{-\alpha x} \quad (4)$$

However, this equation implies that the intensities I and I_0 are measured within the absorbing medium. The complex dielectric constant ϵ for such a material is:

$$\epsilon = \epsilon_1 + i\epsilon_2 \quad (5)$$

where the dielectric constants are related to the optical constants by:

$$\epsilon_1 = n^2 - k^2 \quad (6)$$

$$\epsilon_2 = 2nk \quad (7)$$

Two additional materials properties that influence the light-material interaction are magnetic susceptibility μ and conductivity σ that are further discussed later.

The equations describing the reflection phenomena, including polarization effects for metals, will not be presented here but are explained in detail elsewhere.^{5-8,10,11} After a brief description of Lorentz and Drude theories and their implications for metals, and particularly for absorption, the relationship among reflection, transmission, and absorption is discussed.⁹

The classical theory of absorption in dielectrics is due to H. A. Lorentz¹² and in metals to P. K. L. Drude.¹³ Both models treat the optically active electrons in a material as classical oscillators. In the Lorentz model, the electron is considered to be bound to the nucleus by a harmonic restoring force. In this manner, Lorentz's picture is that of the nonconductive dielectric. Drude considered the electrons to be free, and set the restoring force in the Lorentz model equal to zero. Both models include a damping term in the electron's equation of motion that in more modern terms is recognized as a result of electron-phonon collisions.

These models solve for the electron's motion in the presence of the electromagnetic field as a driving force. From this, it is possible to write an expression for the polarization induced in the medium and from that to derive the dielectric constant. The Lorentz model for dielectrics gives the relative real and imaginary parts of the dielectric constant ϵ_{1R} and ϵ_{2R} in terms of N , the number of dipoles per unit volume; e and m , the electron charge and mass; Γ , the damping constant; ω and ω_0 , the radian frequencies of the field and the harmonically bound electron; and ϵ_0 , the permittivity of free space. These functions are shown in Fig. 1. The range of frequencies where ϵ_1 increases with frequency is referred to as the *range of normal dispersion*, and the region near $\omega = \omega_0$, where it decreases with frequency is called the *range of anomalous dispersion*.

Since the ionic polarizability is much smaller than the electronic polarizability at optical frequencies, only the electronic terms are considered when evaluating optical absorption using the Lorentz model for dielectrics. The Drude model for metals assumes that the electrons are free to move. This means that it is identical to the Lorentz model except that ω_0 is set equal to zero. The real and imaginary parts of the dielectric constant are then given by

$$\epsilon_{1R} = 1 - (Ne^2\epsilon_0 m) \frac{1}{\omega^2 + \Gamma^2} \quad (8)$$

$$\epsilon_{2R} = (Ne^2\epsilon_0 m) \frac{\Gamma}{\omega(\omega^2 + \Gamma^2)} \quad (9)$$

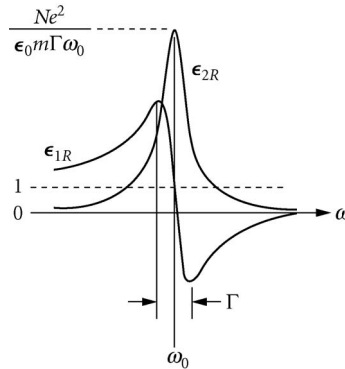


FIGURE 1 Frequency dependences of ϵ_{1R} and ϵ_{2R} .⁹

The quantity Γ is related to the mean time between electron collisions with lattice vibrations, and by considering electronic motion in an electric field \mathbf{E} having radian frequency ω , an expression for the average velocity can be obtained. An expression for the conductivity σ is then obtained and the parts of the dielectric constant can be restated. At electromagnetic field frequencies that are low, it can be shown that $\epsilon_2 \gg \epsilon_1$ and therefore it follows that:

$$\alpha = (\omega\mu\sigma/2)^{1/2} \quad (10)$$

In other words, the optical properties and the conductivity of a perfect metal are related through the fact that each is determined by the motion of free electrons. At high frequencies, transitions involving bound or valence band electrons are possible and there will be a noticeable deviation from this simple result of the Drude model. However, the experimental data reported for most metals are in good agreement with the Drude prediction at wavelengths as short as 1 μm .

From Eq. (10) it is clear that a field propagating in a metal will be attenuated by a factor of $1/e$ when it has traveled a distance:

$$\delta = (2/\omega\mu\sigma)^{1/2} \quad (11)$$

This quantity is called the *skin depth*, and at optical frequencies for most metals it is ~ 50 nm. After a light beam has propagated one skin depth into a metal, its intensity is reduced to 0.135 of its value at the surface.

Another aspect of the absorption of light energy by metals that should be noted is the fact that it increases with temperature. This is important because during laser irradiation the temperature of a metal will increase and so will the absorption. The coupling of energy into the metal is therefore dependent on the temperature dependence of the absorption. For most metals, all the light that gets into the metal is absorbed. If the Fresnel expression for the electric field reflectance is applied to the real and imaginary parts of the complex index for a metal-air interface, the field reflectivity can be obtained. When multiplied by its complex conjugate, the expression for the intensity reflection coefficient is obtained:

$$R_I = 1 - 2\mu\epsilon_0\omega/\sigma \quad (12)$$

Since the conductivity σ decreases with increasing temperature, R_I decreases with increasing temperature, and at higher temperatures more of the incident energy is absorbed.

Since reflection methods are used in determining the optical constants, they are strongly dependent on the characteristics of the metallic surface. These characteristics vary considerably with chemical and mechanical treatment, and these treatments have not always been accurately defined. Not all measurements have been made on freshly polished surfaces but in many cases on freshly deposited thin films. The best available data are presented in the tables and figures, and the reader is advised to consult the appropriate references for specifics.

In this article, an ending of *-ance* denotes a property of a specific sample (i.e., including effects of surface finish), while the ending *-ivity* refers to an intrinsic material property. For most of the discussion, the endings are interchangeable.

Reflectance r is the ratio of radiant flux reflected from a surface to the total incident radiant flux. Since r is a function of the optical constants, it varies with wavelength and temperature. The relationship between reflectance and optical constants is:⁵

$$r = \frac{(n-1)^2 + k^2}{(n+1)^2 + k^2} \quad (13)$$

The reflectance of a good, freshly deposited mirror coating is almost always higher than that of a polished or electroplated surface of the same material. The reflectance is normally less than unity—some transmission and absorption, no matter how small, are always present. The relationship between these three properties is:

$$r + t + a = 1 \quad (14)$$

Transmittance t is the ratio of radiant flux transmitted through a surface to the total incident radiant flux and absorptance a is the ratio of the radiant flux lost by absorption to the total incident radiant flux. Since t and a are functions of the optical constants, they vary with wavelength and temperature. Transmittance is normally very small for metals except in special cases (e.g., beryllium at x-ray wavelengths). Absorptance is affected by surface condition as well as the intrinsic contribution of the material.

The thermal radiative properties are descriptive of a radiant energy-matter interaction that can be described by other properties such as the optical constants and/or complex dielectric constant, each of which is especially convenient for studying various aspects of the interaction. However, the thermal radiative properties are particularly useful since metallic materials are strongly influenced by surface effects, particularly oxide films, and therefore in many cases they are not readily calculated by simple means from the other properties.

For opaque materials, the transmission is near zero, so Eq. (14) becomes:

$$r + a = 1 \quad (15)$$

but since Kirchhoff's law states that absorptance equals emittance, ϵ , this becomes:

$$r + \epsilon = 1 \quad (16)$$

and the thermal radiative properties of an opaque body are fully described by either the reflectance or the emittance. Emittance is the ratio of radiated emitted power (in W/m^2) of a surface to the emissive power of a blackbody at the same temperature. Emittance can therefore be expressed as either *spectral* (emittance as a function of wavelength at constant temperature) or *total* (the integrated emittance over all wavelengths as a function of temperature).

Physical Properties

The physical properties of interest for metals in optical applications include density, electrical conductivity, and electrical resistivity (the reciprocal of conductivity), as well as crystal structure. Chemical composition of alloys is also included with physical properties.

For density, mass density is reported with units of kg/m^3 . Electrical conductivity is related to electrical resistivity, but for some materials, one or the other is normally reported. Both properties vary with temperature.

Crystal structure is extremely important for stability since anisotropy of the elastic, electric, and magnetic properties and thermal expansion depend on the type of structure.¹⁴ Single crystals of cubic metals have completely isotropic coefficient of thermal expansion (CTE), but are anisotropic in elastic properties—modulus and Poisson's ratio. Materials with hexagonal structures have anisotropic expansion and elastic properties. While polycrystalline metals with randomly oriented small grains do not exhibit these anisotropies they can easily have local areas that are inhomogeneous or can have overall oriented crystal structure induced by fabrication methods.

The combined influence of physical, thermal, and mechanical properties on optical system performance is described under "Properties Important in Mirror Design," later in this chapter.

Thermal Properties

Thermal properties of metals that are important in optical systems design include: coefficient of thermal expansion α , referred to in this section as CTE; thermal conductivity k ; and specific heat C_p . All of these properties vary with temperature; usually they tend to decrease with decreasing temperature. Although not strictly a thermal property, the maximum usable temperature is also included as a guide for the optical designer.

Thermal expansion is a generic term for change in length for a specific temperature change, but there are precise terms that describe specific aspects of this material property. ASTM E338 Committee recommends the following nomenclature:¹⁵

Coefficient of linear thermal expansion (CTE or thermal expansivity):

$$\alpha \equiv \frac{1}{L} \frac{\Delta L}{\Delta T} \quad (17)$$

Instantaneous coefficient of linear thermal expansion:

$$\alpha' \equiv \lim_{\Delta T \rightarrow 0} \left(\frac{1}{L} \frac{\Delta L}{\Delta T} \right) \quad (18)$$

Mean coefficient of linear thermal expansion:

$$\bar{\alpha} \equiv \frac{1}{T_2 - T_1} \int_{T_1}^{T_2} \alpha' dT \quad (19)$$

In general, lower thermal expansion is better for optical system performance, as it minimizes the effect of thermal gradients on component dimensional changes. CTE is the prime parameter in materials selection for cooled mirrors.

Thermal conductivity is the quantity of heat transmitted per unit of time through a unit of area per unit of temperature gradient with units of W/mK. Higher thermal conductivity is desirable to minimize temperature gradients when there is a heat source to the optical system.

Specific heat, also called heat capacity per unit mass, is the quantity of heat required to change the temperature of a unit mass of material one degree under conditions of constant pressure. A material with high specific heat requires more heat to cause a temperature change that might cause a distortion. High specific heat also means that more energy is required to force a temperature change (e.g., in cooling an infrared telescope assembly to cryogenic temperatures).

Maximum usable temperature is not a hard number. It is more loosely defined as the temperature at which there is a significant change in the material due to one or more of a number of things, such as significant softening or change in strength, melting, recrystallization, and crystallographic phase change.

Mechanical Properties

Mechanical properties are divided into elastic/plastic properties and strength, and fracture properties. The elastic properties of a metal can be described by a 6×6 matrix of constants called the elastic stiffness constants.¹⁶⁻¹⁸ Because of symmetry considerations, there are a maximum of 21 independent constants that are further reduced for more symmetrical crystal types. For cubic materials there are three constants, C_{11} , C_{12} , and C_{44} , and for hexagonal five constants, C_{11} , C_{12} , C_{13} , C_{33} , and C_{44} . From these, the elastic properties of the material, Young's modulus E (the elastic modulus in tension), bulk modulus K , modulus of rigidity G (also called shear modulus), and Poisson's ratio ν can be calculated. The constants, and consequently the properties, vary as functions of temperature. The properties vary with crystallographic direction in single crystals,¹⁴ but in randomly oriented polycrystalline materials the macroproperties are usually isotropic.

Young's modulus of elasticity E is the measure of stiffness or rigidity of a metal—the ratio of stress, in the completely elastic region, to the corresponding strain. Bulk modulus K is the measure of resistance to change in volume—the ratio of hydrostatic stress to the corresponding change in volume. Shear modulus, or modulus of rigidity, G is the ratio of shear stress to the corresponding shear strain under completely elastic conditions.

Poisson's ratio ν is the ratio of the absolute value of the rate of transverse (lateral) strain to the corresponding axial strain resulting from uniformly distributed axial stress in the elastic deformation region.

For isotropic materials the properties are interrelated by the following equations:¹⁸

$$G = \frac{E}{2(1+\nu)} \quad (20)$$

$$K = \frac{E}{3(1-2\nu)} \quad (21)$$

The mechanical strength and fracture properties are important for the structural aspect of the optical system. The components in the system must be able to support loads with no permanent deformation within limits set by the error budget and certainly with no fracture. For ductile materials such as copper, the yield and/or microyield strength may be the important parameters. On the other hand, for brittle or near-brittle metals such as beryllium, fracture toughness may be more important. For ceramic materials such as silicon carbide, fracture toughness and modulus of rupture are the important fracture criteria. A listing of definitions for each of these terms¹⁹ follows:

creep strength: the stress that will cause a given time-dependent plastic strain in a creep test for a given time.

ductility: the ability of a material to deform plastically before fracture.

fatigue strength: the maximum stress that can be sustained for a specific number of cycles without failure.

fracture toughness: a generic term for measures of resistance to extension of a crack.

hardness: a measure of the resistance of a material to surface indentation.

microcreep strength: the stress that will cause 1 ppm of permanent strain in a given time; usually less than the microyield strength.

microstrain: a deformation of 10^{-6} m/m (1 ppm).

microyield strength: the stress that will cause 1 ppm of permanent strain in a short time; also called precision elastic limit (PEL).

ultimate strength: the maximum stress a material can withstand without fracture.

yield strength: the stress at which a material exhibits a specified deviation from elastic behavior (proportionality of stress and strain), usually 2×10^{-3} m/m (0.2 percent).

Properties Important in Mirror Design

There are many factors that enter into the design of a mirror or mirror system, but the most important requirement is optical performance. Dimensional stability, weight, durability, and cost are some of the factors to be traded off before an effective design can be established.²⁰⁻²³ The loading conditions during fabrication, transportation, and use and the thermal environment play a substantial role in materials selection. To satisfy the end-use requirements, the optical, structural, and thermal performance must be predictable. Each of these factors has a set of parameters and associated material properties that can be used to design an optic to meet performance goals.

For optical performance, the shape or optical figure is the key performance factor followed by the optical properties of reflectance, absorptance, and complex refractive index. The optical properties of a mirror substrate material are only important when the mirror is to be used bare (i.e., with no optical coating).

To design for structural performance goals, deflections due to static (or inertial) and dynamic loads are usually calculated as a first estimate.²⁴ For this purpose, the well-known plate equations²⁵ are invoked. For the static case,

$$\delta = \beta q a^4 / D \quad (22)$$

where δ = deflection

β = deflection coefficient (depends on support condition)

q = normal loading (uniform load example)

a = plate radius (semidiameter)

D = flexural rigidity, defined as:

$$D = EI_0 / (1 - \nu^2) \quad (23)$$

where, in turn E = Young's modulus of elasticity

I_0 = moment of inertia of the section

ν = Poisson's ratio

But

$$q = \rho V_0 G \quad (24)$$

where ρ = material density

V_0 = volume of material per unit area of plate surface

G = load factor (g's)

After substitution and regrouping the terms:

$$\delta = \beta \frac{\rho(1-\nu^2)}{E} \frac{V_0}{I_0} a^4 G \quad (25)$$

or

$$\delta \times B = M \times S \times P \times G \quad (26)$$

where B = support condition

M = materials parameters

S = structural efficiency

P = plate size

G = load factor

This shows five terms, each representing a parameter to be optimized for mirror performance. B , P , and G will be determined from system requirements; S is related to the geometric design of the part; and M is the materials term showing that ρ , ν , and E are the important material properties for optimizing structural performance.

For the dynamic case of deflection due to a local angular acceleration $\ddot{\theta}$ about a diameter (scanning applications), the equation becomes:

$$\delta_{\text{DYN}} = \beta_D \frac{\rho(1-\nu^2)}{E} \frac{V_0}{I_0} a^5 \frac{\ddot{\theta}}{g} \quad (27)$$

The same structural optimization parameters prevail as in the static case. Note that in both cases maximizing the term E/ρ (specific stiffness) minimizes deflection.

The determination of thermal performance^{26,27} is dependent on the thermal environment and thermal properties of the mirror material. For most applications, the most significant properties are the coefficient of thermal expansion CTE or α , and thermal conductivity k . Also important are the specific heat C_p , and thermal diffusivity D , a property related to dissipation of thermal gradients that is a combination of properties and equal to $k/\rho C_p$. There are two important thermal figures of merit, the coefficients of thermal distortion α/k and α/D . The former expresses steady-state distortion per unit of input power, while the latter is related to transient distortions.

Typical room-temperature values for many of the important properties mentioned here are listed for a number of mirror materials in Table 1. It should be clear from the wide range of properties and figures of merit that no one material can satisfy all applications. A selection process is required and a tradeoff study has to be made for each individual application.²⁰

Metal optical components can be designed and fabricated to meet system requirements. However, unless they remain within specifications throughout their intended lifetime, they have failed. The most often noted changes that occur to degrade performance are dimensional instability and/or environment-related optical property degradation. Dimensional instabilities can take many forms with many causes, and there are any number of ways to minimize them. Dimensional instabilities can only be discussed briefly here; for a more complete discussion, consult Refs. 28 to 31. The instabilities most often observed are:

- temporal instability: a change in dimensions with time in a uniform environment (e.g., a mirror stored in a laboratory environment with no applied loads changes figure over a period of time)

TABLE 1 Properties of Selected Mirror Materials

	ρ Density (10^3 kg/m^3)	E Young's Modulus (GN/m ²)	E/ρ Specific Stiffness (arb. units)	CTE Thermal Expansion ($10^{-6}/\text{K}$)	k Thermal Conductivity (W/m K)	C_p Specific Heat (J/kg K)	D Thermal Diffusivity ($10^{-6} \text{ m}^2/\text{s}$)	Distortion coefficient	
								CTE/ k Steady State ($\mu\text{m}/\text{W}$)	CTE/ D Transient ($\text{s}/\text{m}^2 \text{ K}$)
Preferred	small	large	large	small	large	large	large	small	small
Fused silica	2.19	72	33	0.50	1.4	750	0.85	0.36	0.59
Beryllium: 1–70	1.85	287	155	11.3	216	1925	57.2	0.05	0.20
Aluminum: 6061	2.70	68	25	22.5	167	896	69	0.13	0.33
Copper	8.94	117	13	16.5	391	385	115.5	0.53	0.14
304 stainless steel	8.00	193	24	14.7	16.2	500	4.0	0.91	3.68
Invar 36	8.05	141	18	1.0	10.4	515	2.6	0.10	0.38
Silicon	2.33	131	56	2.6	156	710	89.2	0.02	0.03
SiC: RB-30% Si	2.89	330	114	2.6	155	670	81.0	0.02	0.03
SiC: CVD	3.21	465	145	2.4	198	733	82.0	0.01	0.03

- thermal/mechanical cycling or hysteresis instability: a change in dimensions when the environment is changed and then restored, where the measurements are made under the same conditions before and after the exposure (e.g., a mirror with a measured figure is cycled between high and low temperatures and, when remeasured under the original conditions, the figure has changed)
- thermal instability: a change in dimensions when the environment is changed, but completely reversible when the original environment is restored (e.g., a mirror is measured at room temperature, again at low temperature where the figure is different, and finally at the original conditions with the original figure restored)

There are other types of instabilities, but they are less common, particularly in metals. The sources of the dimensional changes cited here can be attributed to one or more of the following:

- externally applied stress
- changes in internal stress
- microstructural changes
- inhomogeneity/anisotropy of properties

In general, temporal and cycling/hysteresis instabilities are primarily caused by changes in internal stress (i.e., stress relaxation). If the temperature is high enough, microstructural changes can take place as in annealing, recrystallization, or second-phase precipitation. Thermal instability is a result of inhomogeneity and/or anisotropy of thermal expansion within the component, is completely reversible, and cannot be eliminated by nondestructive methods.

To eliminate potential instabilities, care must be taken in the selection of materials and fabrication methods to avoid anisotropy and inhomogeneity. Further care is necessary to avoid any undue applied loads that could cause part deformation and subsequent residual stress. The fabrication methods should include stress-relief steps such as thermal annealing, chemical removal of damaged surfaces, and thermal or mechanical cycling. These steps become more critical for larger and more complex component geometries.

Instabilities can also be induced by attachments and amounts. Careful design to minimize induced stresses and selection of dissimilar materials with close thermal expansion matching is essential.³²

4.3 SUMMARY DATA

The properties presented here are representative for the materials and are not a complete presentation. For more complete compilations, the references should be consulted.

Optical Properties

Thin films and their properties are discussed in Chap. 7, "Optical Properties of Films and Coatings," and therefore are not presented here except in the case where bulk (surface) optical properties are not available.

Index of Refraction and Extinction Coefficient The data for the optical constants of metals are substantial, with the most complete listing available in the two volumes of *Optical Constants of Metals*,^{33,34} from which most of the data presented here have been taken. Earlier compilations^{35,36} are also available. While most of the data are for deposited films, the references discuss properties of polished polycrystalline surfaces where available. Table 2 lists room-temperature values for n and k of Al,³⁷ Be,³⁸ Cu,³⁹ Cr,⁴⁰ Au,³⁹ Fe,⁴⁰ Mo,³⁹ Ni,³⁹ Pt,³⁹ Ag,³⁹ W,³⁹ and α -SiC.⁴¹ Figures 2 to 14 graphically show these constants with the absorption edges shown in most cases.

Extensive reviews of the properties of aluminum³⁷ and beryllium³⁸ also discuss the effects of oxide layers on optical constants and reflectance. Oxide layers on aluminum typically reduce the optical constant values by 25 percent in the infrared, 10 to 15 percent in the visible, and very little in the ultraviolet.³⁷ As a result of the high values of n and k for aluminum in the visible and infrared, there are relatively large variations of optical constants with temperature, but they result in only small changes in reflectance.³⁷ The beryllium review³⁸ does not mention any variation of properties with temperature. The optical properties of beryllium and all hexagonal metals vary substantially with crystallographic direction. This variation with crystallography is shown for the dielectric constants of beryllium in Fig. 15.⁴² The optical constants can be obtained from the dielectric constants using the following equations:⁹

$$n = \left\{ \left[(\epsilon_1^2 + \epsilon_2^2)^{1/2} + \epsilon_1 \right] / 2 \right\}^{1/2} \quad (28)$$

$$k = \left\{ \left[(\epsilon_1^2 + \epsilon_2^2)^{1/2} - \epsilon_1 \right] / 2 \right\}^{1/2} \quad (29)$$

This variation in optical properties results in related variations in reflectance and absorptance that may be the main contributors to a phenomenon called *anomalous scatter*, where the measured scatter from polished surfaces does not scale with wavelength when compared to the measured surface roughness.⁴³⁻⁴⁷

The optical constants reported for SiC are for single-crystal hexagonal material.

Reflectance and Absorptance

Reflectance data in the literature are extensive. Summaries have been published for most metals³⁵⁻³⁶ primarily at normal incidence, both as deposited films and polished bulk material. Reflectance as a function of angle is presented for a number of metals in Refs. 48 and 49. Selected data are also included in Ref. 50. Temperature dependence of reflectance is discussed in a number of articles, but little measured data are available. Absorption data summaries are not as readily available, with one summary³⁵ and many articles for specific materials, primarily at laser wavelengths and often as a function of temperature. Table 3 lists values of room-temperature normal-incidence reflectance as a function of wavelength, and Figs. 16 to 26 show r and a calculated from η and k in the range of 0.015 to 10 μm .³⁵ Figure 27³⁵ shows reflectance for polarized radiation as a function of incidence angle for three combinations of n and k , illustrating the tendency toward total external reflectance for angles greater than about 80°.

TABLE 2 n and k of Selected Metals at Room Temperature

Metal	eV	Wavelength Å	μm	n	k
Aluminum ³⁷	300.0	41.3		1.00	0.00
	180.0	68.9		0.99	0.01
	130.0	95.4		0.99	0.02
	110.0	113.0		0.99	0.03
	100.0	124.0		0.99	0.03
	95.0	131.0		1.00	0.04
	80.0	155.0		1.01	0.02
	75.0	165.0		1.01	0.02
	72.0	172.0		1.02	0.00
	50.0	248.0		0.97	0.01
	25.0	496.0		0.81	0.02
	17.0	729.0		0.47	0.04
	12.0	1,033.0	0.10	0.03	0.79
	6.00	2,066.0	0.21	0.13	2.39
	4.00	3,100.0	0.31	0.29	3.74
	3.10	4,000.0	0.40	0.49	4.86
	2.48	5,000.0	0.50	0.77	6.08
	2.07	6,000.0	0.60	1.02	7.26
	1.91	6,500.0	0.65	1.47	7.79
	1.77	7,000.0	0.70	1.83	8.31
1.55	8,000.0	0.80	2.80	8.45	
1.10		1.13	1.20	11.2	
0.827		1.50	1.38	15.4	
0.620		2.00	2.15	20.7	
0.310		4.00	6.43	39.8	
0.177		7.00	14.0	66.2	
0.124		10.0	25.3	89.8	
0.062		20.0	60.7	147.0	
0.039		32.0	103.0	208.0	
Beryllium ³⁸	300.0	41.3		1.00	0.00
	200.0	62.0		0.99	0.00
	150.0	82.7		0.99	0.01
	119.0	104.0		1.00	0.02
	100.0	124.0		0.99	0.00
	50.0	248.0		0.93	0.01
	25.0	496.0		0.71	0.10
	17.0	729.0		0.34	0.42
	12.0	1,033.0	0.10	0.30	1.07
	6.00	2,066.0	0.21	0.85	2.64
	4.00	3,100.0	0.31	2.47	3.08
		4,133.0	0.41	2.95	3.14
		5,166.0	0.52	3.03	3.18
		6,888.0	0.69	3.47	3.23
			1.03	3.26	3.96
			3.10	2.07	12.6
		6.20	3.66	26.7	
		12.0	11.3	50.1	
		21.0	19.9	77.1	
		31.0	37.4	110.0	
		62.0	86.1	157.0	
Copper ³⁹	9,000.0	1.38		1.00	0.00
	4,000.0	3.10		1.00	0.00

TABLE 2 *n* and *k* of Selected Metals at Room Temperature (Continued)

Metal	eV	Wavelength Å	μm	<i>n</i>	<i>k</i>
Copper ³⁹	1,500.0	8.27		1.00	0.00
	1,000.0	12.4		1.00	0.00
	900.0	13.8		1.00	0.00
	500.0	24.8		1.00	0.00
	300.0	41.3		0.99	0.01
	200.0	62.0		0.98	0.02
	150.0	82.7		0.97	0.03
	120.0	103.0		0.97	0.05
	100.0	124.0		0.97	0.07
	50.0	248.0		0.95	0.13
	29.0	428.0		0.85	0.30
	26.0	477.0		0.92	0.40
	24.0	517.0		0.96	0.37
	23.0	539.0		0.94	0.37
	20.0	620.0		0.88	0.46
	15.0	827.0		1.01	0.71
	12.0	1,033.0	0.10	1.09	0.73
	6.50	1,907.0	0.19	0.96	1.37
	5.20	2,384.0	0.24	1.38	1.80
	4.80	2,583.0	0.26	1.53	1.71
	4.30	2,885.0	0.29	1.46	1.64
	2.60	4,768.0	0.48	1.15	2.5
	2.30	5,390.0	0.54	1.04	2.59
	2.10	5,904.0	0.59	0.47	2.81
	1.80	6,888.0	0.69	0.21	4.05
	1.50	8,265.0	0.83	0.26	5.26
	0.950		1.30	0.51	6.92
	0.620		2.00	0.85	10.6
	0.400		3.10	1.59	16.5
	0.200		6.20	5.23	33.0
0.130		9.54	10.8	47.5	
Chromium ⁴⁰	10,000.0	1.24		1.00	0.00
	6,015.0	2.06		1.00	0.00
	5,878.0	2.11		1.00	0.00
	3,008.0	4.12		1.00	0.00
	1,504.0	8.24		1.00	0.00
	992.0	12.5		1.00	0.00
	735.0	16.9		1.00	0.00
	702.0	17.7		1.00	0.00
	686.0	18.1		1.00	0.00
	403.0	30.8		1.00	0.00
	202.0	61.5		0.98	0.00
	100.0	124.0		0.94	0.03
	62.0	200.0		0.88	0.12
	52.0	238.0		0.92	0.18
	29.5	420.0		0.78	0.21
	24.3	510.0		0.67	0.39
	18.0	689.0		0.87	0.70
	14.3	867.0		1.06	0.82
12.8	969.0		1.15	0.75	
11.4	1,088.0	0.109	1.08	0.69	
7.61	1,629.0	0.163	0.66	1.23	

(Continued)

TABLE 2 n and k of Selected Metals at Room Temperature (*Continued*)

Metal	eV	Wavelength Å	μm	n	k
Chromium ⁴⁰	5.75	2,156.0	0.216	0.97	1.74
	4.80	2,583.0	0.258	0.86	2.13
	3.03	4,092.0	0.409	1.54	3.71
	2.42	5,123.0	0.512	2.75	4.46
	1.77	7,005.0	0.700	3.84	4.37
	1.26	9,843.0	0.984	4.50	4.28
	1.12		1.11	4.53	4.30
	0.66		1.88	3.96	5.95
	0.60		2.07	4.01	6.48
	0.34		3.65	2.89	12.0
	0.18		6.89	8.73	25.4
	0.09		13.8	11.8	33.9
	0.06		20.7	21.2	42.0
	0.04		31.0	14.9	65.2
Gold ³⁹	8,266.0	1.50		1.00	0.00
	2,480.0	5.00		1.00	0.00
	2,066.0	6.00		1.00	0.00
	1,012.0	12.25		1.00	0.00
	573.0	21.6		1.00	0.00
	220.0	56.4		0.99	0.01
	150.0	82.7		0.96	0.01
	86.0	144.0		0.89	0.06
	84.5	147.0		0.89	0.07
	84.0	148.0		0.89	0.06
	68.0	182.0		0.86	0.12
	60.0	207.0		0.86	0.16
	34.0	365.0		0.78	0.47
	30.0	413.0		0.89	0.60
	29.0	428.0		0.91	0.60
	27.0	459.0		0.90	0.64
	26.0	480.0		0.85	0.56
	21.8	570.0		1.02	0.85
	19.4	640.0		1.16	0.73
	17.7	700.0		1.08	0.68
	15.8	785.0		1.03	0.74
	12.4	1,000.0	0.10	1.20	0.84
	8.27	1,550.0	0.15	1.45	1.11
	7.29	1,700.0	0.17	1.52	1.07
	6.36	1,950.0	0.20	1.42	1.12
	4.10	3,024.0	0.30	1.81	1.92
	3.90	3,179.0	0.32	1.84	1.90
	3.60	3,444.0	0.34	1.77	1.85
	3.00	4,133.0	0.41	1.64	1.96
2.60	4,769.0	0.48	1.24	1.80	
2.20	5,636.0	0.56	0.31	2.88	
1.80	6,888.0	0.69	0.16	3.80	
1.40	8,856.0	0.89	0.21	5.88	
1.20		1.03	0.27	7.07	
0.82		1.51	0.54	9.58	
0.40		3.10	1.73	19.2	
0.20		6.20	5.42	37.5	
0.125		9.92	12.2	54.7	

TABLE 2 n and k of Selected Metals at Room Temperature (Continued)

Metal	eV	Wavelength Å	μm	n	k
Iron ^{36,40}	10,000.0	1.24		1.00	0.00
	7,071.0	1.75		1.00	0.00
	3,619.0	3.43		1.00	0.00
	1,575.0	7.87		1.00	0.00
	884.0	14.0		1.00	0.00
	825.0	15.0		1.00	0.00
	320.0	38.8		0.99	0.00
	211.0	58.7		0.98	0.01
	153.0	81.2		0.97	0.02
	94.0	132.0		0.94	0.05
	65.0	191.0		0.90	0.12
	56.6	219.0		0.98	0.19
	54.0	230.0		1.11	0.18
	51.6	240.0		0.97	0.05
	30.0	413.0		0.82	0.13
	22.2	559.0		0.71	0.35
	20.5	606.0		0.74	0.42
	18.0	689.0		0.78	0.51
	15.8	785.0		0.77	0.61
	11.5	1,078.0	0.11	0.93	0.84
	11.0	1,127.0	0.11	0.91	0.83
	10.3	1,200.0	0.12	0.87	0.91
	8.00	1,550.0	0.15	0.94	1.18
	5.00	2,480.0	0.25	1.14	1.87
	3.00	4,133.0	0.41	1.88	3.12
	2.30	5,390.0	0.54	2.65	3.34
	2.10	5,903.0	0.59	2.80	3.34
	1.50	8,265.0	0.83	3.05	3.77
	1.24		1.00	3.23	4.35
	0.496		2.50	4.13	8.59
0.248		5.00	4.59	15.4	
0.124		10.0	5.81	30.4	
0.062		20.0	9.87	60.1	
0.037		33.3	22.5	100.0	
0.025		50.0	45.7	141.0	
0.015		80.0	75.2	158.0	
0.010		125.0	120.0	207.0	
0.006		200.0	183.0	260.0	
0.004		287.0	238.0	306.0	
Molybdenum ³⁹	2,000.0	6.19		1.00	0.00
	1,041.0	11.6		1.00	0.00
	396.0	31.3		1.00	0.01
	303.0	40.9		1.00	0.01
	211.0	58.8		0.99	0.00
	100.0	124.0		0.93	0.01
	60.0	207.0		0.90	0.11
	37.5	331.0		0.81	0.29
	35.0	354.0		0.87	0.38
	33.8	367.0		0.91	0.33
	33.0	376.0		0.90	0.33
	31.4	394.0		0.92	0.31

(Continued)

TABLE 2 n and k of Selected Metals at Room Temperature (*Continued*)

Metal	eV	Wavelength Å	μm	n	k
Molybdenum ³⁹	29.2	424.0		0.84	0.26
	23.4	530.0		0.58	0.55
	17.6	704.0		0.94	1.14
	15.6	795.0		1.15	1.01
	15.0	827.0		1.14	0.99
	14.4	861.0		1.13	1.00
	13.2	939.0		1.20	1.03
	12.0	1,033.0	0.10	1.26	0.92
	11.0	1,127.0	0.11	1.05	0.77
	8.80	1,409.0	0.14	0.65	1.41
	6.20	2,000.0	0.20	0.81	2.50
	4.40	2,818.0	0.28	2.39	3.88
	3.30	3,757.0	0.38	3.06	3.18
	3.10	4,000.0	0.40	3.03	3.22
	2.40	5,166.0	0.52	3.59	3.78
	2.30	5,391.0	0.54	3.79	3.61
	2.20	5,636.0	0.56	3.76	3.41
	2.05	6,052.0	0.61	3.68	3.49
	1.90	6,526.0	0.65	3.74	3.58
	1.70	7,293.0	0.73	3.84	3.51
1.50	8,266.0	0.83	3.53	3.30	
1.20		1.03	2.44	4.22	
0.58		2.14	1.34	11.3	
0.24		5.17	3.61	30.0	
0.12		10.3	13.4	58.4	
0.10		12.4	18.5	68.5	
Nickel ³⁹	9,919.0	1.25		1.00	0.00
	4,133.0	3.00		1.00	0.00
	1,771.0	7.00		1.00	0.00
	929.0	13.3		1.00	0.00
	500.0	24.8		1.00	0.00
	300.0	41.3		0.99	0.01
	180.0	68.9		0.98	0.02
	120.0	103.0		0.96	0.05
	84.0	148.0		0.93	0.11
	68.0	182.0		0.98	0.17
	66.0	188.0		1.01	0.16
	64.0	194.0		0.98	0.11
	50.0	248.0		0.93	0.15
	45.0	276.0		0.88	0.13
	35.0	354.0		0.86	0.24
	23.0	539.0		0.92	0.44
	20.5	605.0		0.89	0.49
	13.0	954.0		1.08	0.71
	10.0	1,240.0	0.12	0.95	0.87
	7.20	1,722.0	0.17	1.03	1.27
6.20	2,000.0	0.20	1.00	1.54	
4.80	2,583.0	0.26	1.53	2.11	
4.15	2,988.0	0.30	1.74	2.00	
3.95	3,140.0	0.31	1.72	1.98	
3.15	3,938.0	0.39	1.61	2.33	
2.40	5,166.0	0.52	1.71	3.06	

TABLE 2 *n* and *k* of Selected Metals at Room Temperature (Continued)

Metal	eV	Wavelength Å	μm	<i>n</i>	<i>k</i>
Nickel ³⁹	1.80	6,888.0	0.69	2.14	4.00
	1.20		1.03	2.85	5.10
	0.45		2.76	4.20	10.2
	0.40		3.10	3.84	11.4
	0.28		4.43	4.30	16.0
	0.22		5.64	4.11	20.2
	0.12		10.3	7.11	38.3
	0.10		12.4	9.54	45.8
Platinum ³⁹	2,000.0	6.20		1.00	0.00
	1,016.0	12.2		1.00	0.00
	504.0	24.6		0.99	0.00
	244.0	50.8		0.99	0.01
	121.0	102.0		0.95	0.02
	83.7	150.0		0.88	0.08
	72.9	170.0		0.89	0.10
	53.9	230.0		0.86	0.20
	51.7	240.0		0.88	0.22
	45.6	250.0		0.87	0.16
	32.6	380.0		0.66	0.45
	30.2	410.0		0.70	0.58
	29.5	420.0		0.71	0.57
	28.8	430.0		0.72	0.65
	28.2	440.0		0.72	0.58
	24.8	500.0		0.71	0.72
	20.7	600.0		0.84	0.94
	16.8	740.0		1.05	0.82
	13.2	940.0		1.20	0.93
	12.7	980.0		1.17	0.96
	10.1	1,230.0	0.12	1.36	1.18
	9.05	1,370.0	0.14	1.43	1.14
	8.38	1,480.0	0.15	1.47	1.15
	7.87	1,575.0	0.16	1.46	1.19
	7.29	1,700.0	0.17	1.49	1.22
	6.05	2,050.0	0.20	1.19	1.40
	5.40	2,296.0	0.23	1.36	1.61
	3.00	4,133.0	0.41	1.75	2.92
	2.30	5,390.0	0.54	2.10	3.67
	1.80	6,888.0	0.69	2.51	4.43
1.20		1.03	3.55	5.92	
0.78		1.55	5.38	7.04	
0.70		1.77	5.71	6.83	
0.65		1.91	5.52	6.66	
0.40		3.10	2.81	11.4	
0.20		6.20	5.90	24.0	
0.13		9.54	9.91	36.7	
0.10		12.4	13.2	44.7	
Silver ³⁹	10,000.0	1.24		1.00	0.00
	6,000.0	2.07		1.00	0.00
	3,000.0	4.13		1.00	0.00
	1,500.0	8.26		1.00	0.00
	800.0	15.5		1.00	0.00
	370.0	33.5		1.01	0.01

(Continued)

TABLE 2 n and k of Selected Metals at Room Temperature (*Continued*)

Metal	eV	Wavelength Å	μm	n	k
Silver ³⁹	350.0	35.4		1.00	0.00
	170.0	72.9		0.97	0.00
	110.0	113.0		0.90	0.02
	95.0	131.0		0.86	0.06
	85.0	146.0		0.85	0.11
	64.0	194.0		0.89	0.21
	50.0	248.0		0.88	0.29
	44.0	282.0		0.90	0.33
	35.0	354.0		0.87	0.45
	31.0	400.0		0.93	0.53
	27.5	451.0		0.85	0.62
	22.5	551.0		1.03	0.62
	21.0	590.0		1.11	0.56
	20.0	620.0		1.10	0.55
	15.0	827.0		1.24	0.69
	13.0	954.0		1.32	0.60
	10.9	1,137.0	0.11	1.28	0.56
	10.0	1,240.0	0.12	1.24	0.57
	9.20	1,348.0	0.13	1.18	0.55
	7.60	1,631.0	0.16	0.94	0.83
	4.85	2,556.0	0.26	1.34	1.35
	4.15	2,988.0	0.30	1.52	0.99
	3.90	3,179.0	0.32	0.93	0.50
	3.10	4,000.0	0.40	0.17	1.95
	2.20	5,636.0	0.56	0.12	3.45
1.80	6,888.0	0.69	0.14	4.44	
1.20		1.03	0.23	6.99	
0.62		2.00	0.65	12.2	
0.24		5.17	3.73	31.3	
0.125		9.92	13.1	53.7	
Tungsten ³⁹	2,000.0	6.20		1.00	0.00
	1,016.0	12.2		1.00	0.00
	516.0	24.0		0.99	0.00
	244.0	50.8		0.99	0.02
	100.0	124.0		0.94	0.04
	43.0	288.0		0.74	0.27
	38.5	322.0		0.82	0.33
	35.0	354.0		0.85	0.31
	33.0	376.0		0.82	0.28
	32.0	388.0		0.79	0.30
	30.5	406.0		0.77	0.29
	23.8	521.0		0.48	0.60
	22.9	541.0		0.49	0.69
	22.1	561.0		0.49	0.76
	16.0	775.0		0.98	1.14
	15.5	800.0		0.96	1.12
	14.6	849.0		0.90	1.20
	11.8	1,051.0	0.11	1.18	1.48
	10.8	1,148.0	0.11	1.29	1.39
10.3	1,204.0	0.12	1.22	1.33	
7.80	1,590.0	0.16	0.93	2.06	
5.60	2,214.0	0.22	2.43	3.70	

TABLE 2 n and k of Selected Metals at Room Temperature (*Continued*)

Metal	eV	Wavelength Å	μm	n	k
Tungsten ³⁹	5.00	2,480.0	0.25	3.40	2.85
	4.30	2,883.0	0.29	3.07	2.31
	4.00	3,100.0	0.31	2.95	2.43
	3.45	3,594.0	0.36	3.32	2.70
	3.25	3,815.0	0.38	3.45	2.49
	3.10	4,000.0	0.40	3.39	2.41
	2.80	4,428.0	0.44	3.30	2.49
	1.85	6,702.0	0.67	3.76	2.95
	1.75	7,085.0	0.71	3.85	2.86
	1.60	7,749.0	0.77	3.67	2.68
	1.20		1.03	3.00	3.64
	0.96		1.29	3.15	4.41
	0.92		1.35	3.14	4.45
	0.85		1.46	2.80	4.33
	0.58		2.14	1.18	8.44
	0.40		3.10	1.94	13.2
	0.34		3.65	1.71	15.7
	0.18		6.89	4.72	31.5
	0.12		10.3	10.1	46.4
	0.07		17.7	26.5	73.8
0.05		24.8	46.5	93.7	
Silicon carbide ⁴¹	30.0	413.0		0.74	0.11
	20.5	605.0		0.35	0.53
	13.1	946.0		0.68	1.41
	9.50	1,305.0	0.13	1.46	2.21
	9.00	1,378.0	0.14	1.60	2.15
	7.60	1,631.0	0.16	2.59	2.87
	6.40	1,937.0	0.19	4.05	1.42
	5.00	2,480.0	0.25	3.16	0.26
	3.90	3,179.0	0.32	2.92	0.01
	3.00	4,133.0	0.41	2.75	0.00
	2.50	4,959.0	0.50	2.68	0.00
	1.79	6,911.0	0.69	2.62	—
	1.50	8,266.0	0.83	2.60	—
	0.62		2.00	2.57	0.00
	0.31		4.00	2.52	0.00
	0.12		6.67	2.33	0.02
	0.11		9.80	1.29	0.01
	0.10		10.40	0.09	0.63
	0.10		10.81	0.06	1.57
	0.10		11.9	0.16	4.51
0.09		12.6	8.74	18.4	
0.08		12.7	17.7	6.03	
0.05		13.1	7.35	0.27	
		15.4	4.09	0.02	
		25.0	3.34	—	

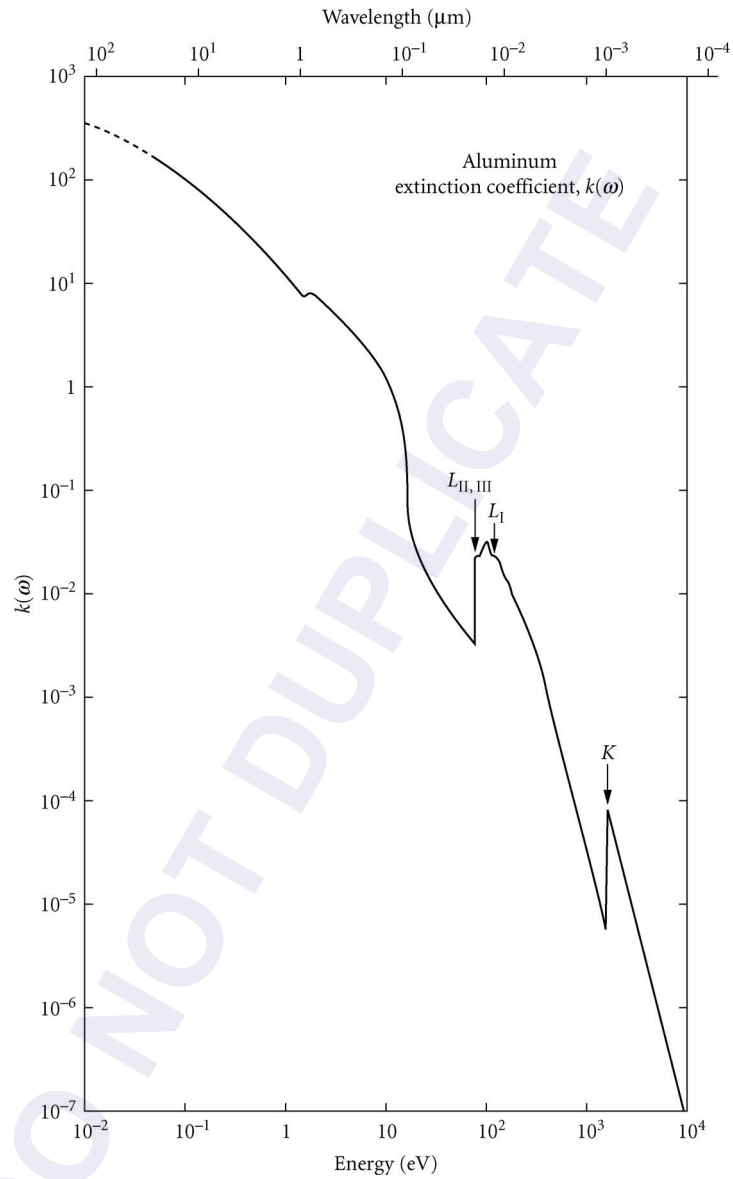


FIGURE 2 k for aluminum vs. photon energy.³⁷

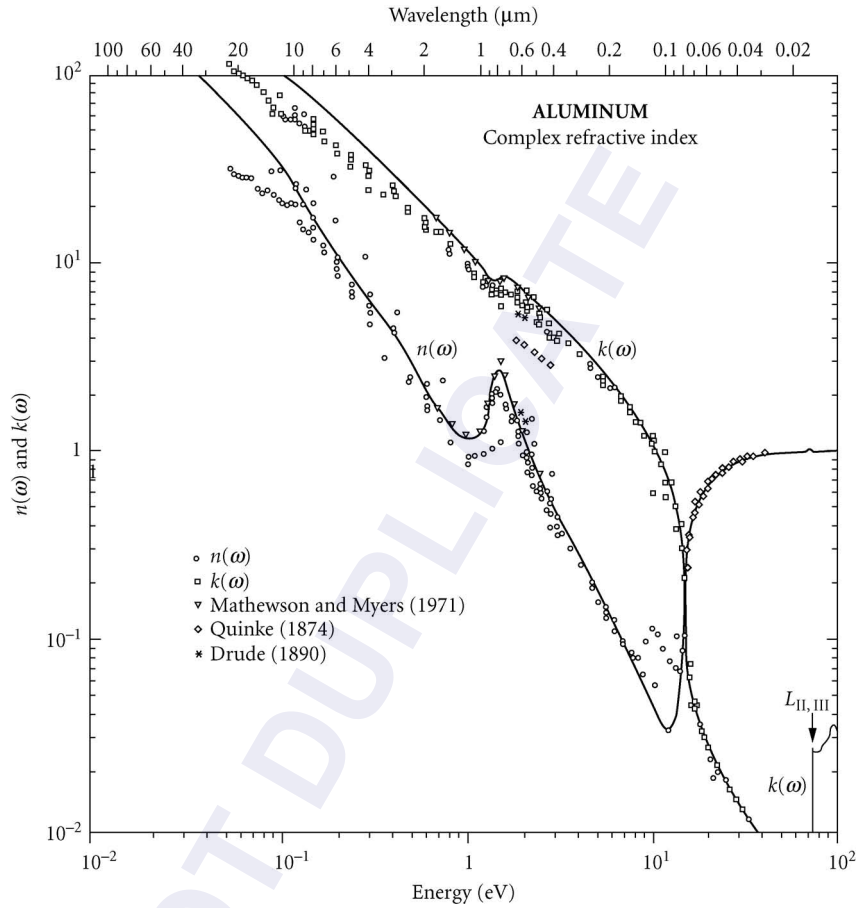


FIGURE 3 n and k for aluminum vs. photon energy.³⁷

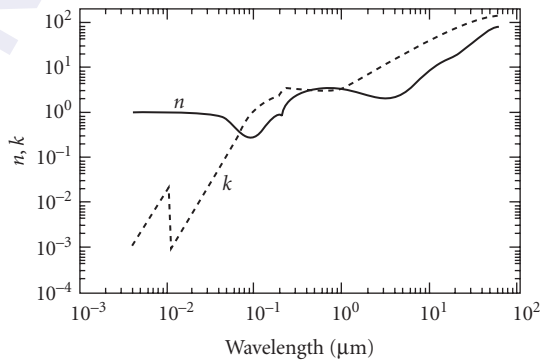


FIGURE 4 n and k for beryllium vs. wavelength.³⁸

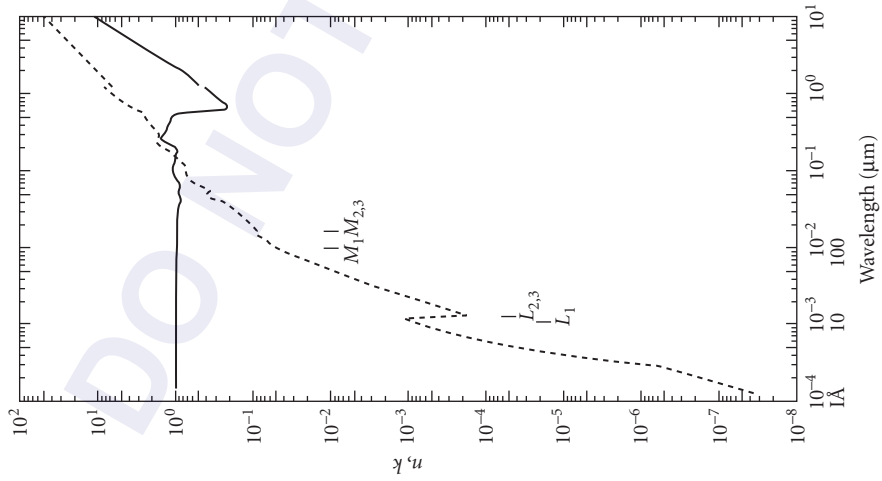


FIGURE 5 n and k for copper vs. wavelength.³⁹

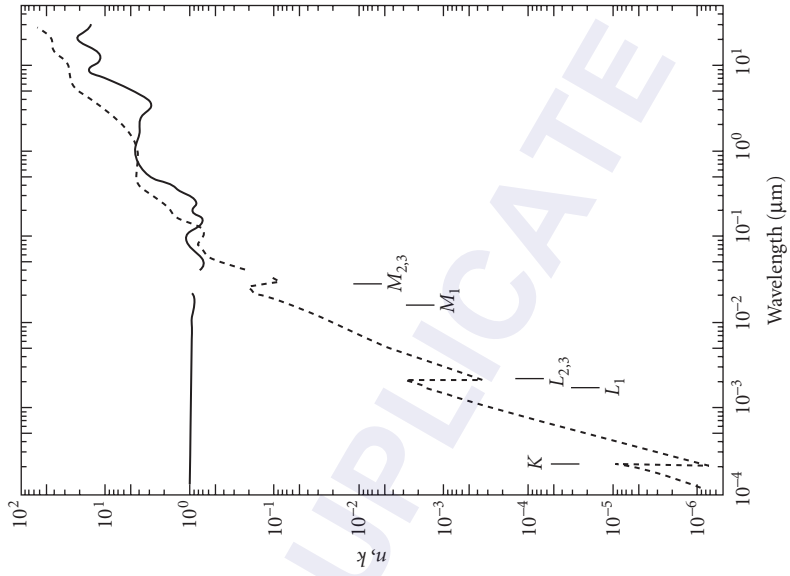


FIGURE 6 n and k for chromium vs. wavelength.⁴⁰

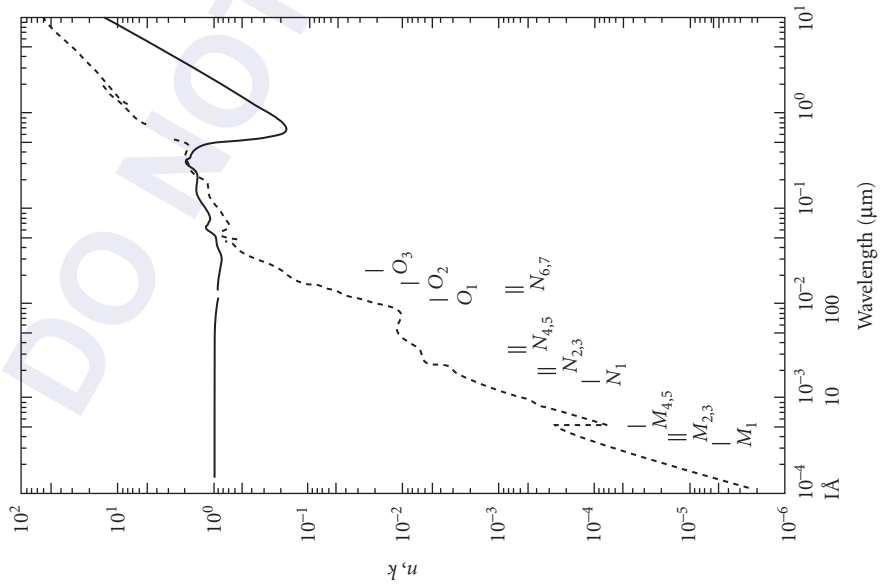


FIGURE 7 n and k for gold vs. wavelength.³⁹

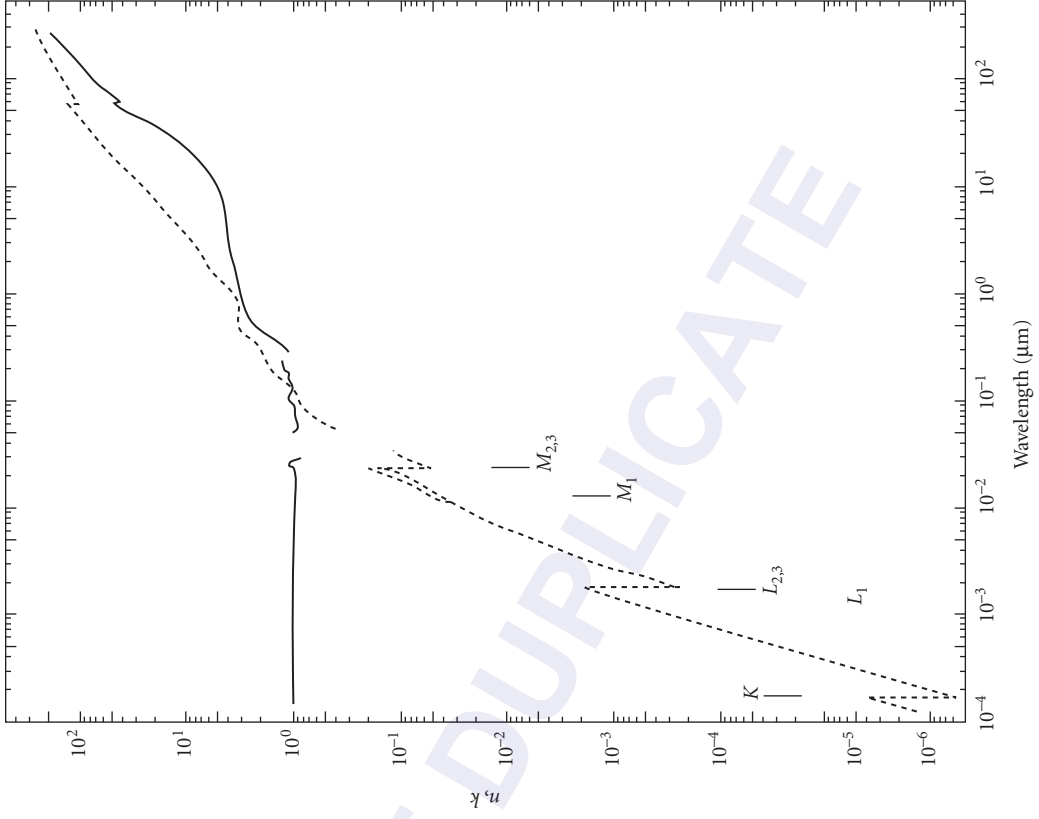


FIGURE 8 n and k for iron vs. wavelength.⁴⁰

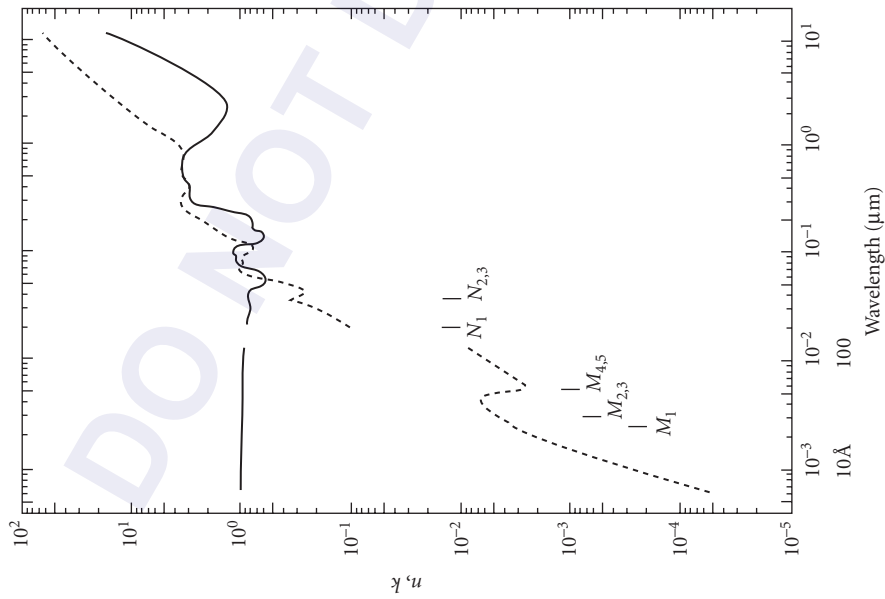


FIGURE 9 n and k for molybdenum vs. wavelength.³⁹

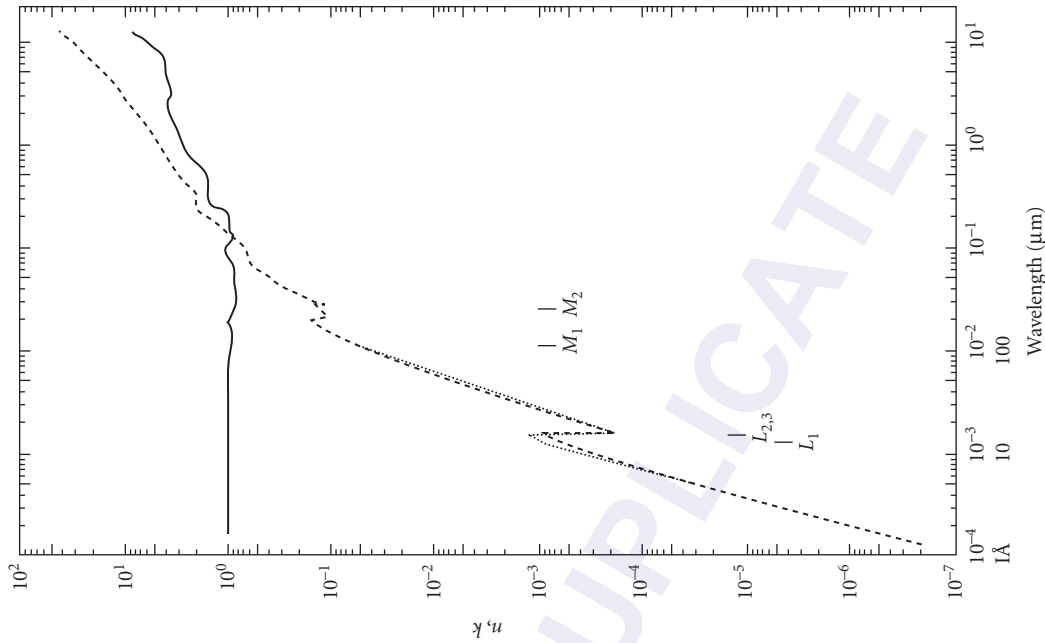


FIGURE 10 n and k for nickel vs. wavelength.³⁹

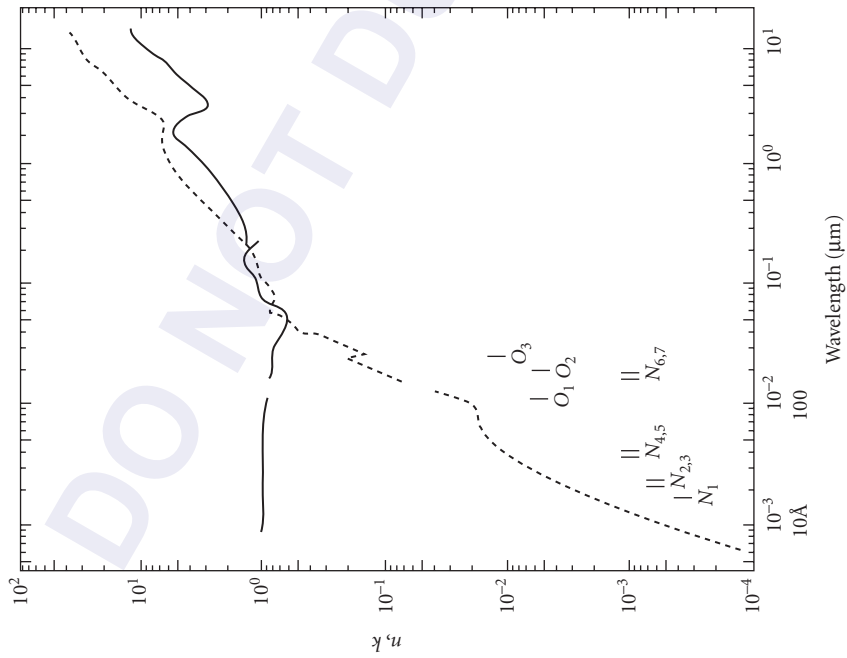


FIGURE 11 n and k for platinum vs. wavelength.³⁹

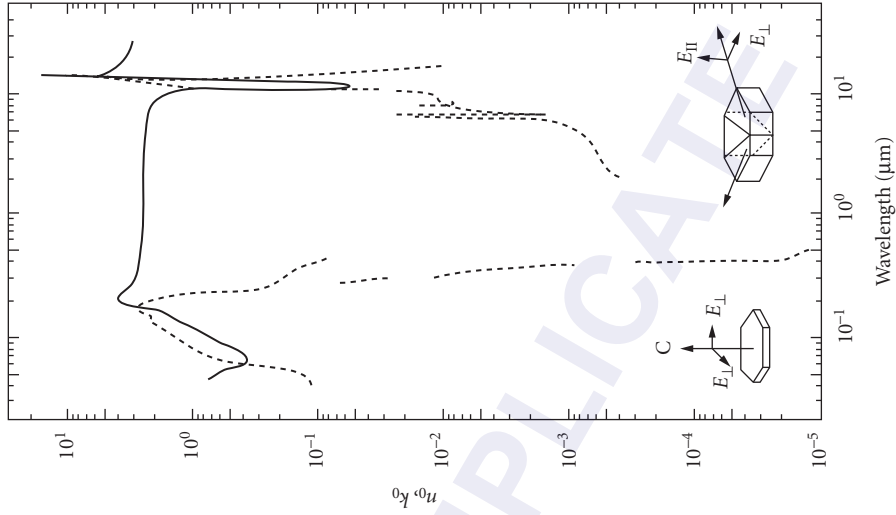


FIGURE 12 n and k for silicon carbide vs. wavelength.⁴¹

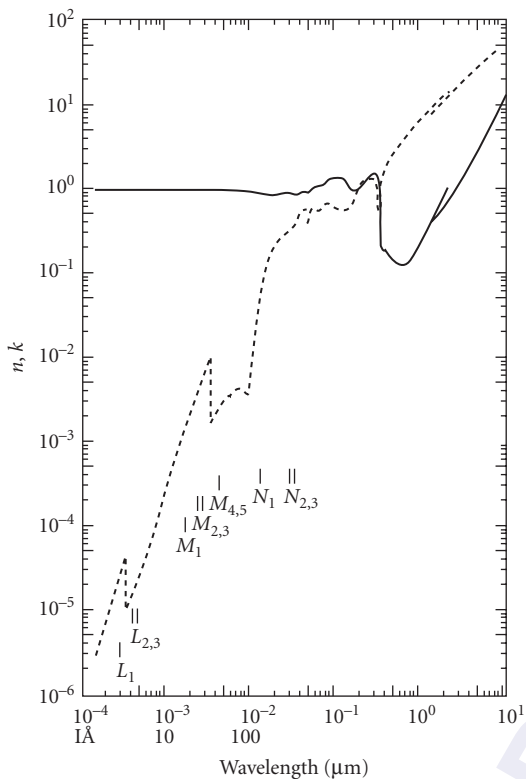


FIGURE 13 n and k for silver as a function of wavelength.³⁹

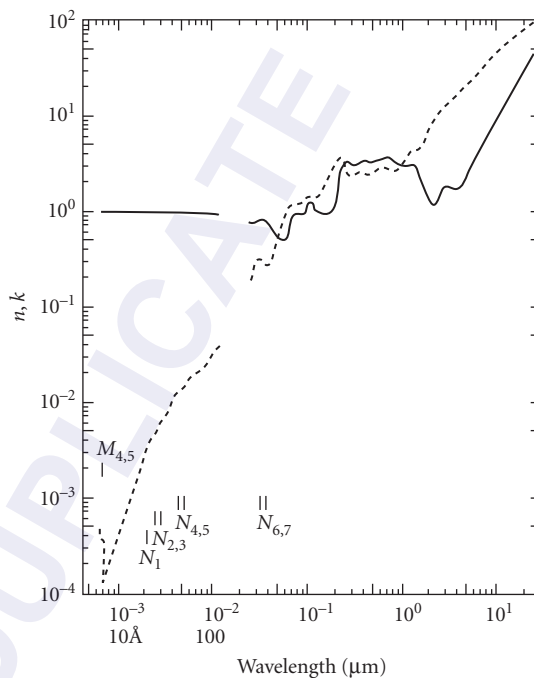


FIGURE 14 n and k for tungsten vs. wavelength.³⁹

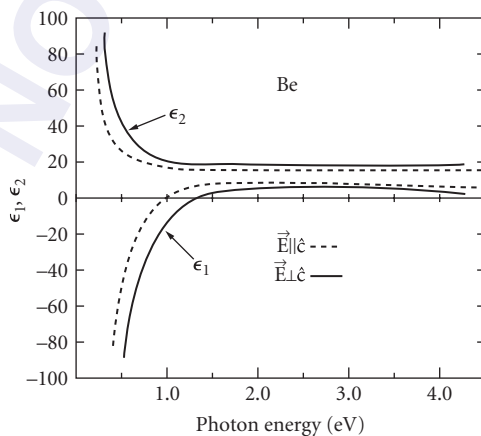


FIGURE 15 Dielectric function for beryllium vs. photon energy showing variation with crystallographic direction.⁴²

TABLE 3 Reflectance of Selected Metals at Normal Incidence

Metal	eV	Wavelength Å	μm	R
Aluminum ³⁶	0.040		31.0	0.9923
	0.050		24.8	0.9915
	0.060		20.7	0.9906
	0.070		17.7	0.9899
	0.080		15.5	0.9895
	0.090		13.8	0.9892
	0.100		12.4	0.9889
	0.125		9.92	0.9884
	0.175		7.08	0.9879
	0.200		6.20	0.9873
	0.250		4.96	0.9858
	0.300		4.13	0.9844
	0.400		3.10	0.9826
	0.600		2.07	0.9806
	0.800		1.55	0.9778
	0.900		1.38	0.9749
	1.00		1.24	0.9697
	1.10		1.13	0.9630
	1.20		1.03	0.9521
	1.30	9,537.0	0.95	0.9318
	1.40	8,856.0	0.89	0.8852
	1.50	8,265.0	0.83	0.8678
	1.60	7,749.0	0.77	0.8794
	1.70	7,293.0	0.73	0.8972
	1.80	6,888.0	0.69	0.9069
	2.00	6,199.0	0.62	0.9148
	2.40	5,166.0	0.52	0.9228
	2.80	4,428.0	0.44	0.9242
	3.20	3,874.0	0.39	0.9243
	3.60	3,444.0	0.34	0.9246
	4.00	3,100.0	0.31	0.9248
	4.60	2,695.0	0.27	0.9249
	5.00	2,497.0	0.25	0.9244
	6.00	2,066.0	0.21	0.9257
	8.00	1,550.0	0.15	0.9269
	10.00	1,240.0	0.12	0.9286
11.00	1,127.0	0.11	0.9298	
13.00	954.0		0.8960	
13.50	918.0		0.8789	
14.00	886.0		0.8486	
14.40	861.0		0.8102	
14.60	849.0		0.7802	
14.80	838.0		0.7202	
15.00	827.0		0.6119	
15.20	816.0		0.4903	
15.40	805.0		0.3881	
15.60	795.0		0.3182	
15.80	785.0		0.2694	
16.00	775.0		0.2326	
16.20	765.0		0.2031	
16.40	756.0		0.1789	
16.75	740.0		0.1460	

(Continued)

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Aluminum ³⁹	17.00	729.0		0.1278
	17.50	708.0		0.1005
	18.00	689.0		0.0809
	19.00	653.0		0.0554
	20.0	620.0		0.0398
	21.0	590.0		0.0296
	22.0	564.0		0.0226
	23.0	539.0		0.0177
	24.0	517.0		0.0140
	25.0	496.0		0.0113
	26.0	477.0		0.0092
	27.0	459.0		0.0076
	28.0	443.0		0.0063
	30.0	413.0		0.0044
	35.0	354.0		0.0020
	40.0	310.0		0.0010
	45.0	276.0		0.0005
	50.0	248.0		0.0003
	55.0	225.0		0.0001
	60.0	206.0		0.0000
	70.0	177.0		0.0000
	72.5	171.0		0.0002
	75.0	165.0		0.0002
	80.0	155.0		0.0002
	85.0	146.0		0.0002
	95.0	131.0		0.0003
	100.0	124.0		0.0002
	120.0	103.0		0.0002
	130.0	95.4		0.0001
150.0	82.7		0.0001	
170.0	72.9		0.0001	
180.0	68.9		0.0000	
200.0	62.0		0.0000	
300.0	41.3		0.0000	
Beryllium ³⁸	0.020		61.99	0.989
	0.040		31.00	0.989
	0.060		20.66	0.988
	0.080		15.50	0.985
	0.100		12.40	0.983
	0.120		10.33	0.982
	0.160		7.75	0.981
	0.200		6.20	0.980
	0.240		5.17	0.978
	0.280		4.43	0.972
	0.320		3.87	0.966
	0.380		3.26	0.955
	0.440		2.82	0.940
	0.500		2.48	0.917
	0.560		2.21	0.887
	0.600		2.07	0.869
0.660		1.88	0.841	
0.720		1.72	0.810	

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Beryllium ³⁸	0.780		1.59	0.775
	0.860		1.44	0.736
	0.940		1.32	0.694
	1.00		1.24	0.667
	1.10		1.13	0.640
	1.20		1.03	0.615
	1.40	8,856.0	0.89	0.575
	1.60	7,749.0	0.77	0.555
	1.90	6,525.0	0.65	0.540
	2.40	5,166.0	0.52	0.538
	2.80	4,428.0	0.44	0.537
	3.00	4,133.0	0.41	0.537
	3.30	3,757.0	0.38	0.536
	3.60	3,444.0	0.34	0.536
	3.80	3,263.0	0.33	0.538
	4.00	3,100.0	0.31	0.541
	4.20	2,952.0	0.30	0.547
	4.40	2,818.0	0.28	0.558
	4.60	2,695.0	0.27	0.575
	Copper ³⁶	0.10		12.4
0.50			2.48	0.979
1.00			1.24	0.976
1.50		8,265.0	0.83	0.965
1.70		7,293.0	0.73	0.958
1.80		6,888.0	0.69	0.952
1.90		6,525.0	0.65	0.943
2.00		6,199.0	0.62	0.910
2.10		5,904.0	0.59	0.814
2.20		5,635.0	0.56	0.673
2.30		5,390.0	0.54	0.618
2.40		5,166.0	0.52	0.602
2.60		4,768.0	0.48	0.577
2.80		4,428.0	0.44	0.545
3.00		4,133.0	0.41	0.509
3.20		3,874.0	0.39	0.468
3.40		3,646.0	0.36	0.434
3.60		3,444.0	0.34	0.407
3.80		3,263.0	0.33	0.387
4.00		3,100.0	0.31	0.364
4.20		2,952.0	0.30	0.336
4.40		2,818.0	0.28	0.329
4.60		2,695.0	0.27	0.334
4.80		2,583.0	0.26	0.345
5.00		2,497.0	0.25	0.366
5.20	2,384.0	0.24	0.380	
5.40	2,296.0	0.23	0.389	
5.60	2,214.0	0.22	0.391	
5.80	2,138.0	0.21	0.389	
6.00	2,066.0	0.21	0.380	
6.50	1,907.0	0.19	0.329	
7.00	1,771.0	0.18	0.271	
7.50	1,653.0	0.17	0.230	

(Continued)

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Copper ³⁶	8.00	1,550.0	0.15	0.206
	8.50	1,459.0	0.15	0.189
	9.00	1,378.0	0.14	0.171
	9.50	1,305.0	0.13	0.154
	10.00	1,240.0	0.12	0.139
	11.00	1,127.0	0.11	0.118
	12.00	1,033.0	0.10	0.111
	13.00	954.0		0.109
	14.00	886.0		0.111
	15.00	827.0		0.111
	16.00	775.0		0.106
	17.00	729.0		0.097
	18.00	689.0		0.084
	19.00	653.0		0.071
	20.00	620.0		0.059
	21.00	590.0		0.048
	22.00	564.0		0.040
	23.00	539.0		0.035
	24.00	517.0		0.035
	25.00	496.0		0.040
	26.00	477.0		0.044
	27.00	459.0		0.043
	28.00	443.0		0.039
	29.00	428.0		0.032
	30.00	413.0		0.025
	32.00	387.0		0.017
	34.00	365.0		0.014
	36.00	344.0		0.012
	38.00	326.0		0.010
	40.00	310.0		0.009
45.00	276.0		0.006	
50.00	248.0		0.005	
55.00	225.0		0.004	
60.00	206.0		0.003	
70.00	177.0		0.002	
90.00	138.0		0.002	
Chromium ³⁶	0.06		20.70	0.962
	0.10		12.40	0.955
	0.14		8.86	0.936
	0.18		6.89	0.953
	0.22		5.64	0.954
	0.26		4.77	0.951
	0.30		4.13	0.943
	0.42		2.95	0.862
	0.54		2.30	0.788
	0.66		1.88	0.736
	0.78		1.59	0.680
	0.90		1.38	0.650
	1.00		1.24	0.639
1.12		1.11	0.631	
1.24	9,998.0	1.00	0.629	
1.36	9,116.0	0.91	0.631	

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Chromium ³⁶	1.46	8,492.0	0.85	0.632
	1.77	7,005.0	0.70	0.639
	2.00	6,199.0	0.62	0.644
	2.20	5,635.0	0.56	0.656
	2.40	5,166.0	0.52	0.677
	2.60	4,768.0	0.48	0.698
	2.80	4,428.0	0.44	0.703
	3.00	4,133.0	0.41	0.695
	4.00	3,100.0	0.31	0.651
	4.40	2,818.0	0.28	0.620
	4.80	2,583.0	0.26	0.572
	5.20	2,384.0	0.24	0.503
	5.60	2,214.0	0.22	0.443
	6.00	2,066.0	0.21	0.444
	7.00	1,771.0	0.18	0.425
	7.60	1,631.0	0.16	0.378
	8.00	1,550.0	0.15	0.315
	8.50	1,459.0	0.15	0.235
	9.00	1,378.0	0.14	0.170
	10.00	1,240.0	0.12	0.120
	11.00	1,127.0	0.11	0.103
	11.50	1,078.0	0.11	0.100
	12.00	1,033.0	0.10	0.101
	13.00	954.0		0.119
	14.00	886.0		0.135
	15.00	827.0		0.143
	16.00	775.0		0.139
	18.00	689.0		0.129
	19.00	653.0		0.131
	20.00	620.0		0.130
22.00	563.0		0.112	
24.00	517.0		0.096	
26.00	477.0		0.063	
28.00	443.0		0.037	
30.00	413.0		0.030	
Gold (electropolished) ³⁶	0.10		12.40	0.995
	0.20		6.20	0.995
	0.40		3.10	0.995
	0.60		2.07	0.994
	0.80		1.55	0.993
	1.00		1.24	0.992
	1.20		1.03	0.991
	1.40	8,856.0	0.89	0.989
	1.60	7,749.0	0.77	0.986
	1.80	6,888.0	0.69	0.979
	2.00	6,199.0	0.62	0.953
	2.10	5,904.0	0.59	0.925
	2.20	5,635.0	0.56	0.880
	2.30	5,390.0	0.54	0.807
2.40	5,166.0	0.52	0.647	
2.50	4,959.0	0.50	0.438	
2.60	4,768.0	0.48	0.331	

(Continued)

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Gold (electropolished) ³⁶	2.70	4,592.0	0.46	0.356
	2.80	4,428.0	0.44	0.368
	2.90	4,275.0	0.43	0.368
	3.00	4,133.0	0.41	0.369
	3.10	3,999.0	0.40	0.371
	3.20	3,874.0	0.39	0.368
	3.40	3,646.0	0.36	0.356
	3.60	3,444.0	0.34	0.346
	3.80	3,263.0	0.33	0.360
	4.00	3,100.0	0.31	0.369
	4.20	2,952.0	0.30	0.367
	4.40	2,818.0	0.28	0.370
	4.60	2,695.0	0.27	0.364
	4.80	2,583.0	0.26	0.344
	5.00	2,497.0	0.25	0.319
	5.40	2,296.0	0.23	0.275
	5.80	2,138.0	0.21	0.236
	6.20	2,000.0	0.20	0.203
	6.60	1,878.0	0.19	0.177
	7.00	1,771.0	0.18	0.162
	7.40	1,675.0	0.17	0.164
	7.80	1,589.0	0.16	0.171
	8.20	1,512.0	0.15	0.155
	8.60	1,442.0	0.14	0.144
	9.00	1,378.0	0.14	0.133
	9.40	1,319.0	0.13	0.122
	9.80	1,265.0	0.13	0.124
	10.20	1,215.0	0.12	0.127
	11.00	1,127.0	0.11	0.116
	12.00	1,033.0	0.10	0.109
14.00	886.0		0.140	
16.00	775.0		0.123	
18.00	689.0		0.109	
20.00	620.0		0.133	
22.00	563.0		0.164	
24.00	517.0		0.125	
26.00	477.0		0.079	
28.00	443.0		0.063	
30.00	413.0		0.064	
Iron ³⁶	0.10		12.40	0.978
	0.15		8.27	0.956
	0.20		6.20	0.958
	0.26		4.77	0.911
	0.30		4.13	0.892
	0.36		3.44	0.867
	0.40		3.10	0.858
	0.50		2.48	0.817
	0.60		2.07	0.783
	0.70		1.77	0.752
	0.80		1.55	0.725
0.90		1.38	0.700	
1.00		1.24	0.678	
1.10		1.13	0.660	

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Iron ³⁶	1.20		1.03	0.641
	1.30	9,537.0	0.95	0.626
	1.40	8,856.0	0.89	0.609
	1.50	8,265.0	0.83	0.601
	1.60	7,749.0	0.77	0.585
	1.70	7,293.0	0.73	0.577
	1.80	6,888.0	0.69	0.573
	1.90	6,525.0	0.65	0.563
	2.00	6,199.0	0.62	0.563
	2.20	5,635.0	0.56	0.563
	2.40	5,166.0	0.52	0.567
	2.60	4,768.0	0.48	0.576
	2.80	4,428.0	0.44	0.580
	3.00	4,133.0	0.41	0.583
	3.20	3,874.0	0.39	0.576
	3.40	3,646.0	0.36	0.565
	3.60	3,444.0	0.34	0.548
	4.00	3,100.0	0.31	0.527
	4.33	2,863.0	0.29	0.494
	4.67	2,655.0	0.27	0.470
	5.00	2,497.0	0.25	0.435
	5.50	2,254.0	0.23	0.401
	6.00	2,066.0	0.21	0.366
	6.50	1,907.0	0.19	0.358
	7.00	1,771.0	0.18	0.333
	7.50	1,653.0	0.17	0.298
	8.00	1,550.0	0.15	0.272
	8.50	1,459.0	0.15	0.251
	9.00	1,378.0	0.14	0.236
	9.50	1,305.0	0.13	0.226
10.00	1,240.0	0.12	0.213	
11.00	1,127.0	0.11	0.162	
11.17	1,110.0	0.11	0.159	
11.33	1,094.0	0.11	0.159	
11.50	1,078.0	0.11	0.160	
12.00	1,033.0	0.10	0.163	
12.50	992.0		0.165	
13.00	954.0		0.162	
13.50	918.0		0.159	
14.00	886.0		0.151	
15.00	827.0		0.135	
16.00	775.0		0.116	
17.00	729.0		0.102	
18.00	689.0		0.091	
20.00	620.0		0.083	
22.00	563.0		0.068	
24.00	517.0		0.045	
26.00	477.0		0.031	
28.00	443.0		0.021	
30.00	413.0		0.014	
Molybdenum ³⁶	0.10		12.40	0.985
	0.20		6.20	0.985
	0.30		4.13	0.983

(Continued)

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Molybdenum ³⁶	0.50		2.70	0.971
	0.70		1.77	0.932
	0.90		1.38	0.859
	1.00		1.24	0.805
	1.10		1.13	0.743
	1.20		1.03	0.671
	1.30	9,537.0	0.95	0.608
	1.40	8,856.0	0.89	0.562
	1.50	8,265.0	0.83	0.550
	1.60	7,749.0	0.77	0.562
	1.70	7,293.0	0.73	0.570
	1.80	6,888.0	0.69	0.576
	2.00	6,199.0	0.62	0.571
	2.20	5,635.0	0.56	0.562
	2.40	5,166.0	0.52	0.594
	2.60	4,768.0	0.48	0.582
	2.80	4,428.0	0.44	0.565
	3.00	4,133.0	0.41	0.550
	3.20	3,874.0	0.39	0.540
	3.40	3,646.0	0.36	0.541
	3.60	3,444.0	0.34	0.546
	3.80	3,263.0	0.33	0.554
	4.00	3,100.0	0.31	0.576
	4.20	2,952.0	0.30	0.610
	4.40	2,818.0	0.28	0.640
	4.60	2,695.0	0.27	0.658
	4.80	2,583.0	0.26	0.678
	5.00	2,497.0	0.25	0.695
	5.20	2,384.0	0.24	0.706
	5.40	2,296.0	0.23	0.706
	5.60	2,214.0	0.22	0.700
	6.00	2,066.0	0.21	0.674
	6.40	1,937.0	0.19	0.641
	6.80	1,823.0	0.18	0.592
	7.20	1,722.0	0.17	0.548
	7.40	1,675.0	0.17	0.542
7.60	1,631.0	0.16	0.552	
7.80	1,589.0	0.16	0.542	
8.00	1,550.0	0.15	0.530	
8.40	1,476.0	0.15	0.495	
8.80	1,409.0	0.14	0.450	
9.20	1,348.0	0.13	0.385	
9.60	1,291.0	0.13	0.320	
10.00	1,240.0	0.12	0.250	
10.40	1,192.0	0.12	0.188	
10.60	1,170.0	0.12	0.138	
11.20	1,107.0	0.11	0.123	
11.60	1,069.0	0.11	0.135	
12.00	1,033.0	0.10	0.154	
12.80	969.0		0.178	
13.60	912.0		0.187	
14.40	861.0		0.182	

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Molybdenum ³⁶	14.80	838.0		0.179
	15.00	827.0		0.179
	16.00	775.0		0.194
	17.00	729.0		0.233
	18.00	689.0		0.270
	19.00	653.0		0.284
	20.00	620.0		0.264
	22.00	563.0		0.207
	24.00	517.0		0.151
	26.00	477.0		0.071
	28.00	443.0		0.036
	30.00	413.0		0.023
	32.00	387.0		0.030
	34.00	365.0		0.034
	36.00	344.0		0.043
	38.00	326.0		0.033
40.00	310.0		0.025	
Nickel ³⁶	0.10		12.40	0.983
	0.15		8.27	0.978
	0.20		6.20	0.969
	0.30		4.13	0.934
	0.40		3.10	0.900
	0.60		2.07	0.835
	0.80		1.55	0.794
	1.00		1.24	0.753
	1.20		1.03	0.721
	1.40	8,856.0	0.89	0.695
	1.60	7,749.0	0.77	0.679
	1.80	6,888.0	0.69	0.670
	2.00	6,199.0	0.62	0.649
	2.40	5,166.0	0.52	0.590
	2.80	4,428.0	0.44	0.525
	3.20	3,874.0	0.39	0.467
	3.60	3,444.0	0.34	0.416
	3.80	3,263.0	0.33	0.397
	4.00	3,100.0	0.31	0.392
	4.20	2,952.0	0.30	0.396
4.60	2,695.0	0.27	0.421	
5.00	2,497.0	0.25	0.449	
5.20	2,384.0	0.24	0.454	
5.40	2,296.0	0.23	0.449	
5.80	2,138.0	0.21	0.417	
6.20	2,000.0	0.20	0.371	
6.60	1,878.0	0.19	0.325	
7.00	1,771.0	0.18	0.291	
8.00	1,550.0	0.15	0.248	
9.00	1,378.0	0.14	0.211	
10.00	1,240.0	0.12	0.166	
11.00	1,127.0	0.11	0.115	
12.00	1,033.0	0.10	0.108	
13.00	954.0		0.105	
14.00	886.0		0.106	

(Continued)

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Nickel ³⁶	15.00	827.0		0.107
	16.00	775.0		0.103
	18.00	689.0		0.092
	20.00	620.0		0.071
	22.00	564.0		0.055
	24.00	517.0		0.051
	27.00	459.0		0.042
	30.00	413.0		0.034
	35.00	354.0		0.022
	40.00	310.0		0.014
	50.00	248.0		0.004
	60.00	206.0		0.002
	65.00	191.0		0.002
	70.00	177.0		0.004
90.00	138.0		0.002	
Platinum ³⁶	0.10		12.40	0.976
	0.15		8.27	0.969
	0.20		6.20	0.962
	0.30		4.13	0.945
	0.40		3.10	0.922
	0.45		2.76	0.882
	0.50		2.50	0.813
	0.55		2.25	0.777
	0.60		2.07	0.753
	0.65		1.91	0.746
	0.70		1.77	0.751
	0.80		1.55	0.762
	0.90		1.38	0.765
	1.00		1.24	0.762
	1.20		1.03	0.746
	1.40	8,856.0	0.89	0.725
	1.60	7,749.0	0.77	0.706
	1.80	6,888.0	0.69	0.686
	2.00	6,199.0	0.62	0.664
	2.50	4,959.0	0.50	0.616
	3.00	4,133.0	0.41	0.565
	4.00	3,100.0	0.31	0.472
	5.00	2,497.0	0.25	0.372
	6.00	2,066.0	0.21	0.276
7.00	1,771.0	0.18	0.230	
8.00	1,550.0	0.15	0.216	
9.00	1,378.0	0.14	0.200	
9.20	1,348.0	0.13	0.198	
9.40	1,319.0	0.13	0.200	
10.20	1,215.0	0.12	0.211	
11.00	1,127.0	0.11	0.199	
12.00	1,033.0	0.10	0.173	
12.80	969.0		0.158	
13.60	912.0		0.155	
14.80	838.0		0.157	
15.20	816.0		0.155	
16.00	775.0		0.146	

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Platinum ³⁶	17.50	708.0		0.135
	18.00	689.0		0.142
	20.00	620.0		0.197
	21.00	590.0		0.226
	22.00	564.0		0.240
	23.00	539.0		0.226
	24.00	517.0		0.201
	26.00	477.0		0.150
	28.00	443.0		0.125
	29.00	428.0		0.118
Silver ³⁶	30.00	413.0		0.124
	0.10		12.40	0.995
	0.20		6.20	0.995
	0.30		4.13	0.994
	0.40		3.10	0.993
	0.50		2.48	0.992
	1.00		1.24	0.987
	1.50	8,265.0	0.83	0.960
	2.00	6,199.0	0.62	0.944
	2.50	4,959.0	0.50	0.914
	3.00	4,133.0	0.41	0.864
	3.25	3,815.0	0.38	0.816
	3.50	3,542.0	0.35	0.756
	3.60	3,444.0	0.34	0.671
	3.70	3,351.0	0.34	0.475
	3.77	3,289.0	0.33	0.154
	3.80	3,263.0	0.33	0.053
	3.90	3,179.0	0.32	0.040
	4.00	3,100.0	0.31	0.103
	4.10	3,024.0	0.30	0.153
	4.20	2,952.0	0.30	0.194
	4.30	2,883.0	0.29	0.208
	4.50	2,755.0	0.28	0.238
	4.75	2,610.0	0.26	0.252
	5.00	2,497.0	0.25	0.257
	5.50	2,254.0	0.23	0.257
	6.00	2,066.0	0.21	0.246
	6.50	1,907.0	0.19	0.225
	7.00	1,771.0	0.18	0.196
7.50	1,653.0	0.17	0.157	
8.00	1,550.0	0.15	0.114	
9.00	1,378.0	0.14	0.074	
10.00	1,240.0	0.12	0.082	
11.00	1,127.0	0.11	0.088	
12.00	1,033.0	0.10	0.100	
13.00	954.0		0.112	
14.00	886.0		0.141	
15.00	827.0		0.156	
16.00	775.0		0.151	
17.00	729.0		0.139	
18.00	689.0		0.124	
19.00	653.0		0.111	

(Continued)

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Silver ³⁶	20.00	620.0		0.103
	21.00	590.0		0.112
	21.50	577.0		0.124
	22.00	564.0		0.141
	22.50	551.0		0.157
	23.00	539.0		0.163
	24.00	517.0		0.165
	25.00	496.0		0.154
	26.00	477.0		0.133
	28.00	443.0		0.090
	30.00	413.0		0.074
	34.00	365.0		0.067
	38.00	326.0		0.043
	42.00	295.0		0.036
	46.00	270.0		0.031
	50.00	248.0		0.027
	56.00	221.0		0.024
	62.00	200.0		0.016
	66.00	188.0		0.016
	70.00	177.0		0.021
76.00	163.0		0.013	
80.00	155.0		0.012	
90.00	138.0		0.009	
100.00	124.0		0.005	
Tungsten ³⁶	0.10		12.40	0.983
	0.20		6.20	0.981
	0.30		4.13	0.979
	0.38		3.26	0.963
	0.46		2.70	0.952
	0.54		2.30	0.948
	0.62		2.00	0.917
	0.70		1.77	0.856
	0.74		1.68	0.810
	0.78		1.59	0.759
	0.82		1.51	0.710
	0.86		1.44	0.661
	0.98		1.27	0.653
	1.10		1.13	0.627
	1.20		1.03	0.590
	1.30	9,537.0	0.95	0.545
	1.40	8,856.0	0.89	0.515
	1.50	8,265.0	0.83	0.500
	1.60	7,749.0	0.77	0.494
	1.70	7,293.0	0.73	0.507
1.80	6,888.0	0.69	0.518	
1.90	6,525.0	0.65	0.518	
2.10	5,904.0	0.59	0.506	
2.50	4,959.0	0.50	0.487	
3.00	4,133.0	0.41	0.459	
3.50	3,542.0	0.35	0.488	
4.00	3,100.0	0.31	0.451	
4.20	2,952.0	0.30	0.440	

TABLE 3 Reflectance of Selected Metals at Normal Incidence (*Continued*)

Metal	eV	Wavelength Å	μm	R
Tungsten ³⁶	4.60	2,695.0	0.27	0.455
	5.00	2,497.0	0.25	0.505
	5.40	2,296.0	0.23	0.586
	5.80	2,138.0	0.21	0.637
	6.20	2,000.0	0.20	0.646
	6.60	1,878.0	0.19	0.631
	7.00	1,771.0	0.18	0.607
	7.60	1,631.0	0.16	0.556
	8.00	1,550.0	0.15	0.505
	8.40	1,476.0	0.15	0.449
	9.00	1,378.0	0.14	0.388
	10.00	1,240.0	0.12	0.287
	10.40	1,192.0	0.12	0.270
	11.00	1,127.0	0.11	0.290
	11.80	1,051.0	0.11	0.318
	12.80	969.0		0.333
	13.60	912.0		0.325
	14.80	838.0		0.276
	15.60	795.0		0.246
	16.00	775.0		0.249
	16.80	738.0		0.273
	17.60	704.0		0.304
	18.80	659.0		0.340
	20.00	620.0		0.354
	21.20	585.0		0.331
	22.40	553.0		0.287
	23.60	525.0		0.252
	24.00	517.0		0.234
	24.80	500.0		0.191
	25.60	484.0		0.150
26.80	463.0		0.105	
28.00	443.0		0.073	
30.00	413.0		0.047	
34.00	365.0		0.032	
36.00	344.0		0.036	
40.00	310.0		0.045	

Figures 28 to 34 show reflectance for polished surfaces and thin films of Al, Be, SiC, and Ni, including effects of oxide films on the surface.⁵¹ One effect of absorption is to limit the penetration depth of incident radiation. Penetration depth is shown in Fig. 35 as a function of wavelength for Al, Be, and Ni.⁵¹

Absorption is a critical parameter for high-energy laser components, and is discussed in hundreds of papers as a function of surface morphology, angle of incidence, polarization state, and temperature. Only a few representative examples of this body of work can be cited here. When absorptance measurements were made of metal mirrors as a function of angle of incidence, polarization state, and wavelength,⁵² it was found that measured values agreed with theory except at high angles of incidence where surface condition plays an undefined role. With the advent of diamond-turning as a mirror-finishing method, many papers have addressed absorptance characteristics of these unique surfaces as a function of surface morphology and angle of incidence, particularly on Ag and Cu mirrors.^{53,54} It has been observed that mirrors have the lowest absorptance when the light is s-polarized and the grooves are oriented parallel to the plane of incidence.⁵⁴ The temperature dependence of optical absorption

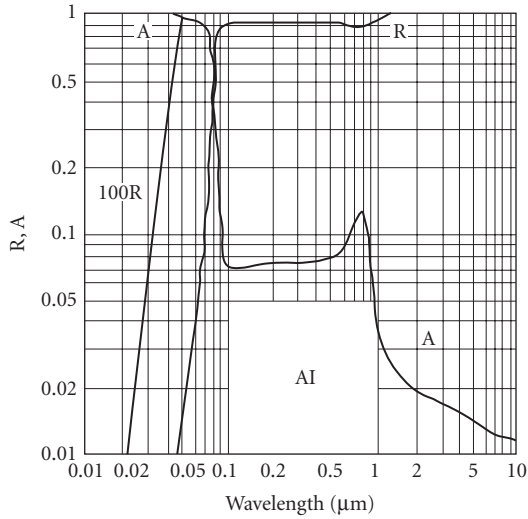


FIGURE 16 Reflectance and absorptance for aluminum vs. wavelength calculated for normal incidence.³⁵ (With permission.)

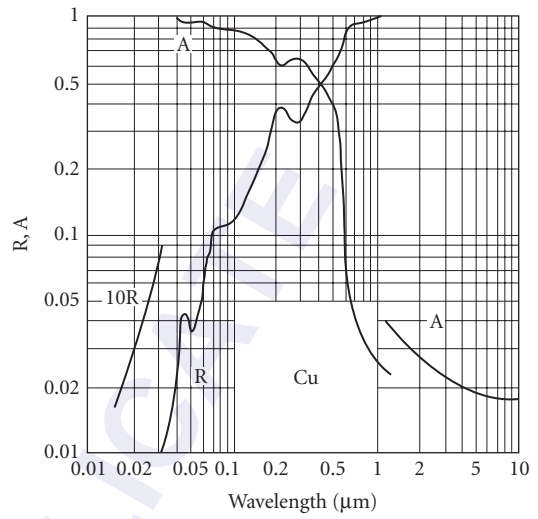


FIGURE 17 Reflectance and absorptance for copper vs. wavelength calculated for normal incidence.³⁵ (With permission.)

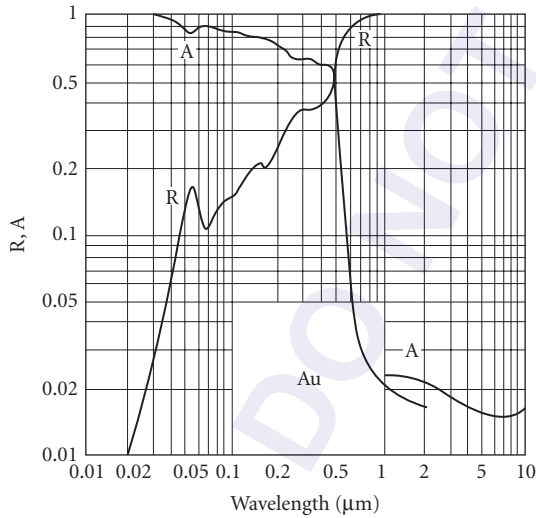


FIGURE 18 Reflectance and absorptance for gold vs. wavelength calculated for normal incidence.³⁵ (With permission.)

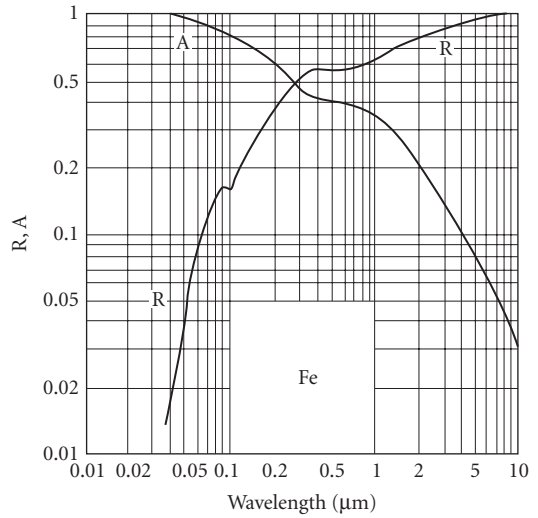


FIGURE 19 Reflectance and absorptance for iron vs. wavelength calculated for normal incidence.³⁵ (With permission.)

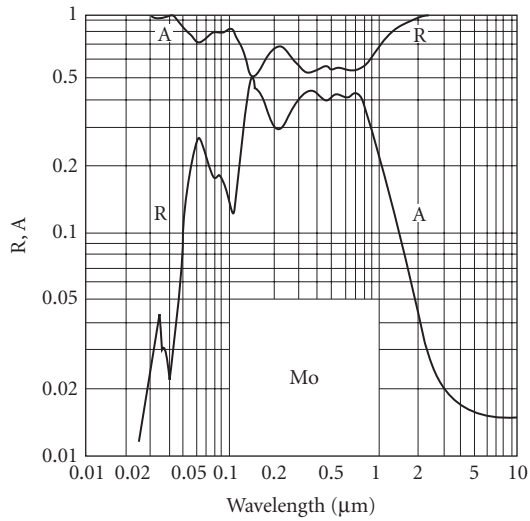


FIGURE 20 Reflectance and absorptance for molybdenum vs. wavelength calculated for normal incidence.³⁵ (With permission.)

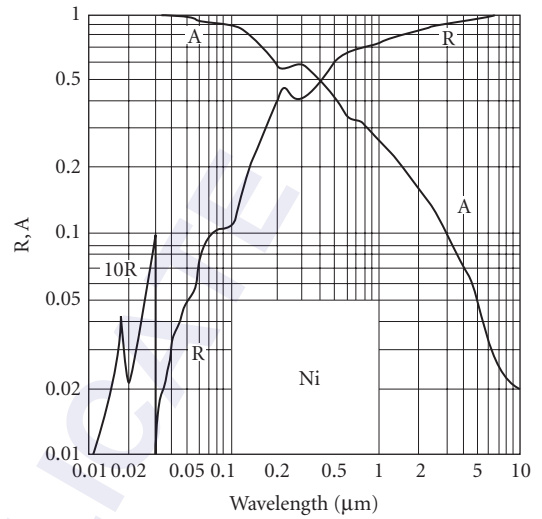


FIGURE 21 Reflectance and absorptance for nickel vs. wavelength calculated for normal incidence.³⁵ (With permission.)

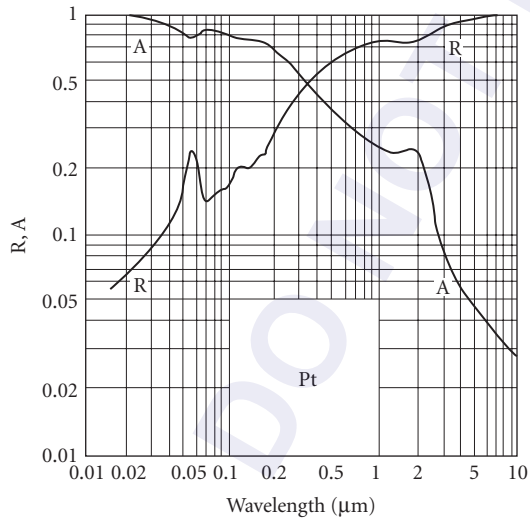


FIGURE 22 Reflectance and absorptance for platinum vs. wavelength calculated for normal incidence.³⁵ (With permission.)

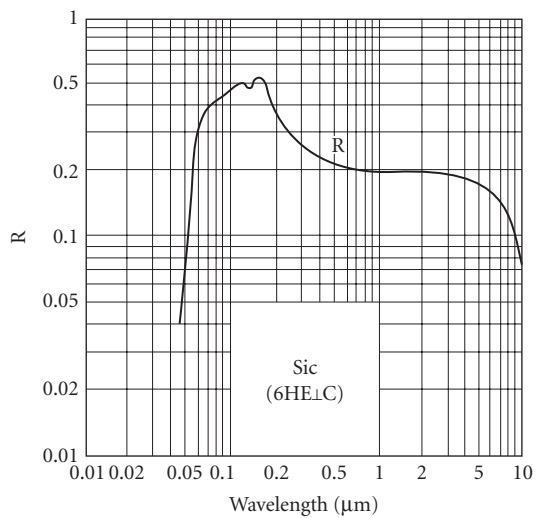


FIGURE 23 Reflectance for the basal plane of hexagonal silicon carbide vs. wavelength calculated for normal incidence.³⁵ (With permission.)

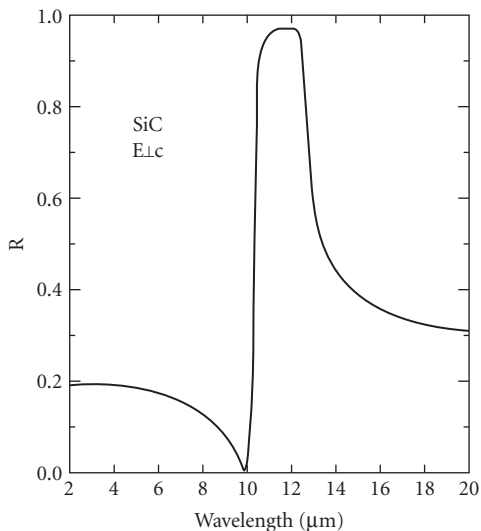


FIGURE 24 Infrared reflectance for the basal plane of hexagonal silicon carbide vs. wavelength calculated for normal incidence.³⁵ (With permission.)

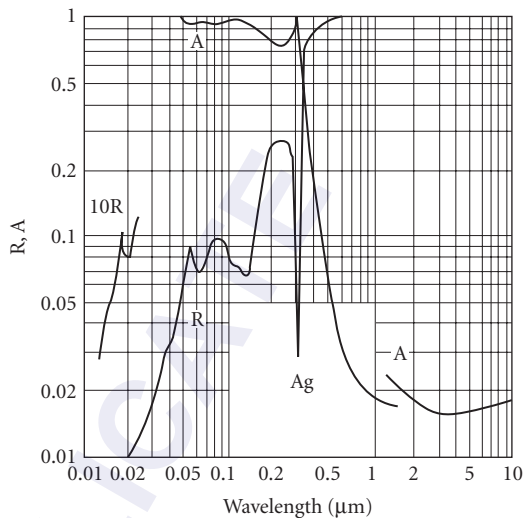


FIGURE 25 Reflectance and absorptance for silver vs. wavelength calculated for normal incidence.³⁵ (With permission.)

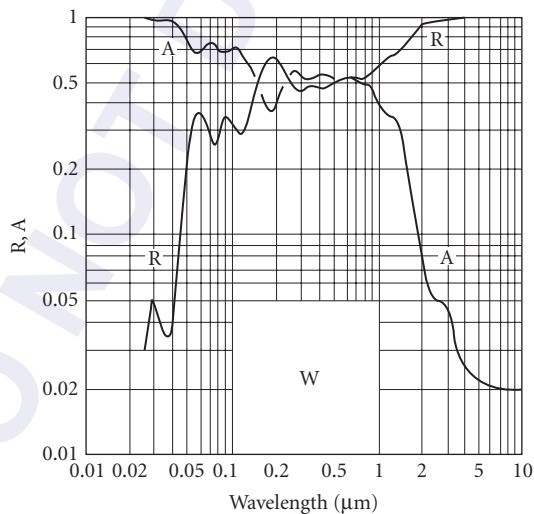


FIGURE 26 Reflectance and absorptance for tungsten vs. wavelength calculated for normal incidence.³⁵ (With permission.)

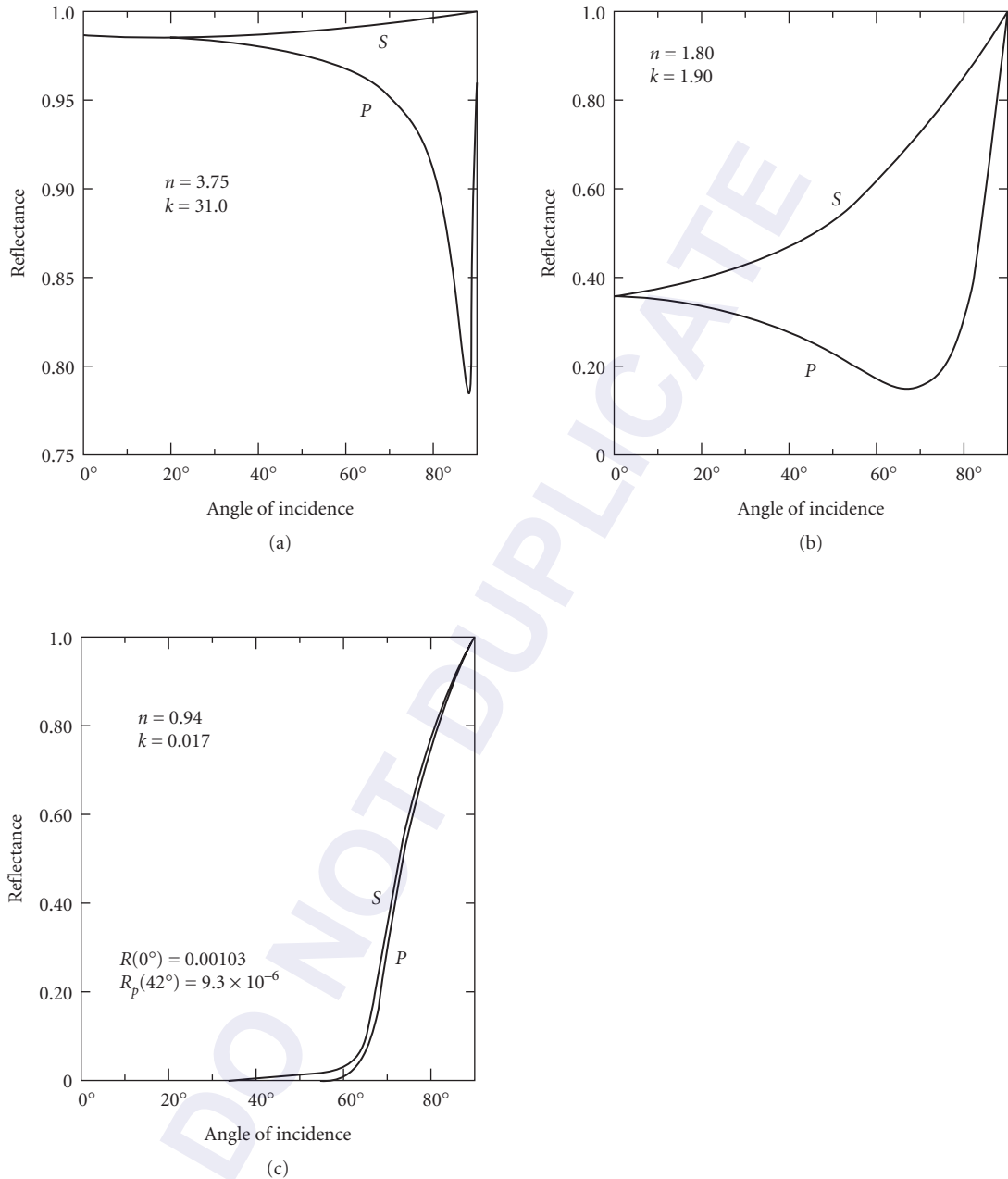


FIGURE 27 Reflectance for polarized radiation vs. angle of incidence for a vacuum-metal interface.³⁵ (With permission.) (a) $n = 3.75$, $k = 31.0$, approximate values for gold at $\lambda = 5 \mu\text{m}$; (b) $n = 1.80$, $k = 1.90$, $\lambda = 0.3 \mu\text{m}$; and (c) $n = 0.94$, $k = 0.017$, approximate values for gold at $\lambda = 0.01 \mu\text{m}$. Note the tendency toward total external reflectance for angle $\geq 80^\circ$.

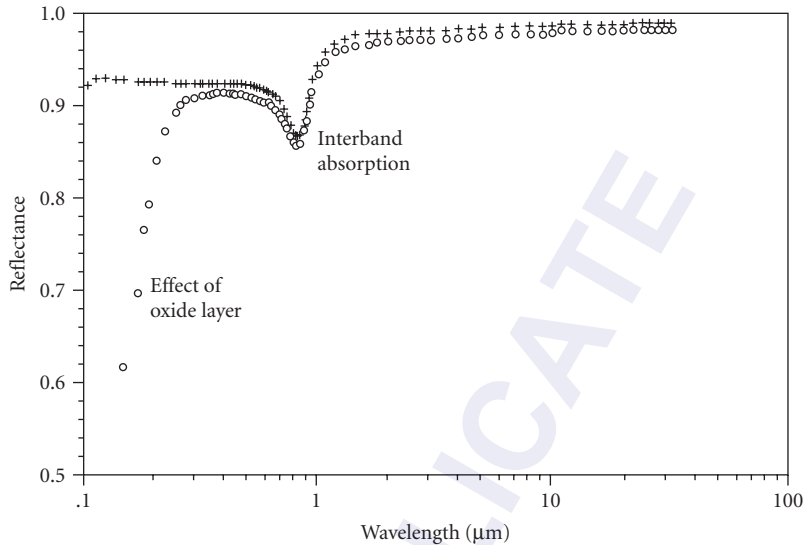


FIGURE 28 Effect of oxide layer on the reflectivity of aluminum vs. wavelength⁵¹ calculated from n and k .³⁷

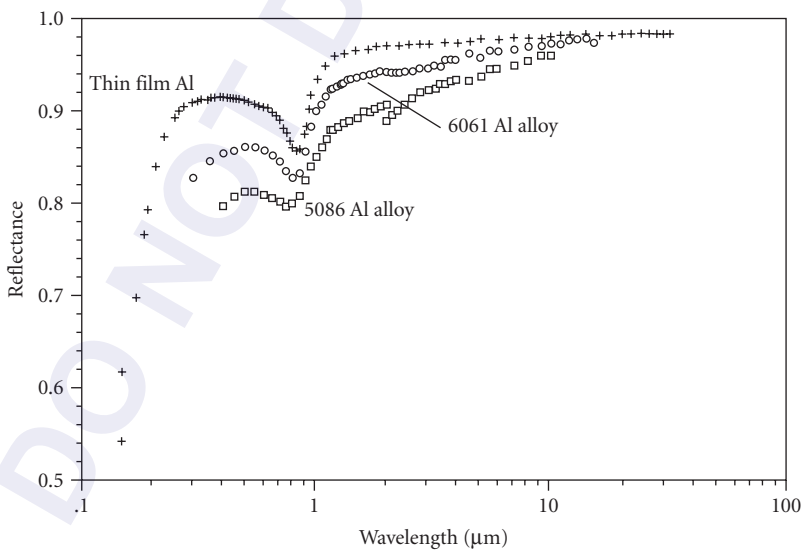


FIGURE 29 Reflectance of optical-grade aluminum alloys vs. wavelength.⁵¹

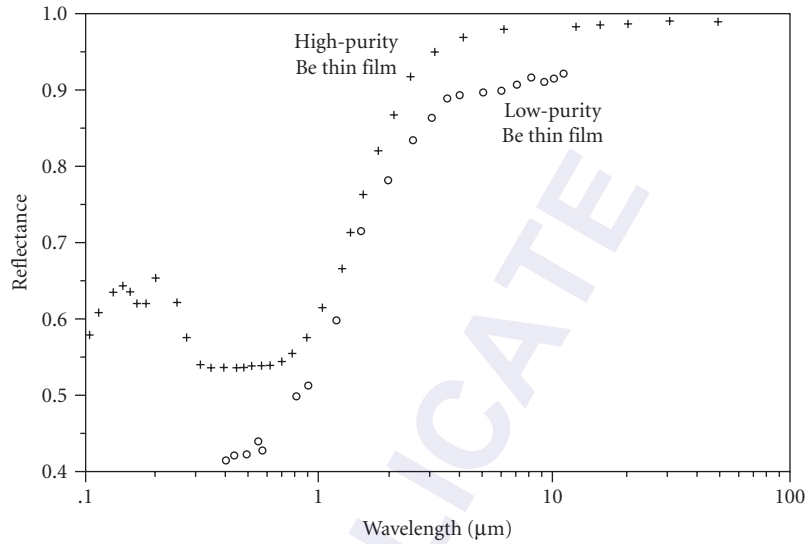


FIGURE 30 Effect of impurities on the reflectance of beryllium thin films vs. wavelength.⁵¹

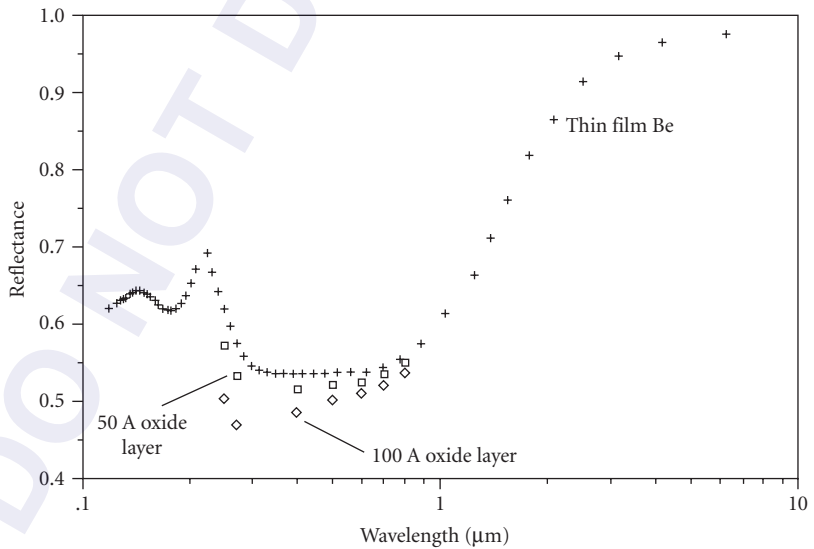


FIGURE 31 Effect of oxide layer thickness on the reflectance of beryllium vs. wavelength⁵¹ calculated from n and k .³⁸

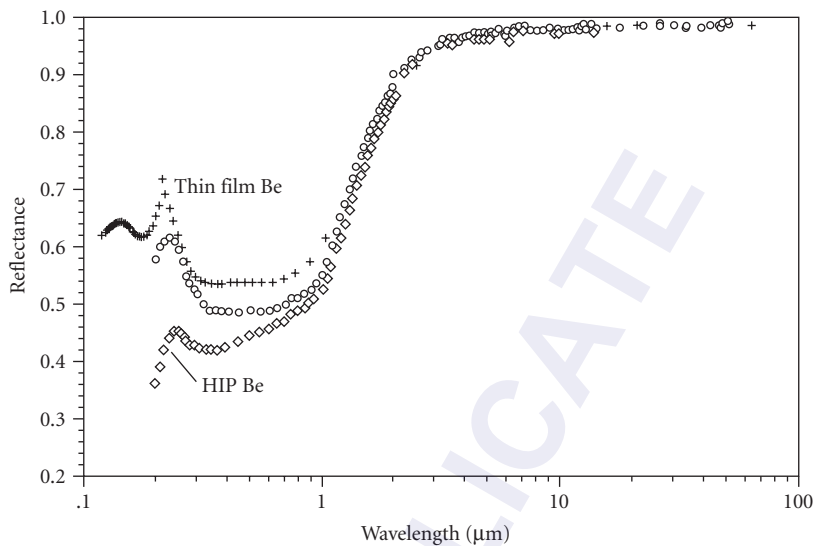


FIGURE 32 Reflectance of polished and evaporated beryllium vs. wavelength;⁵¹ comparison of evaporated high-purity thin film,³⁸ polished high-purity thick film, and polished bulk beryllium (2 percent BeO).

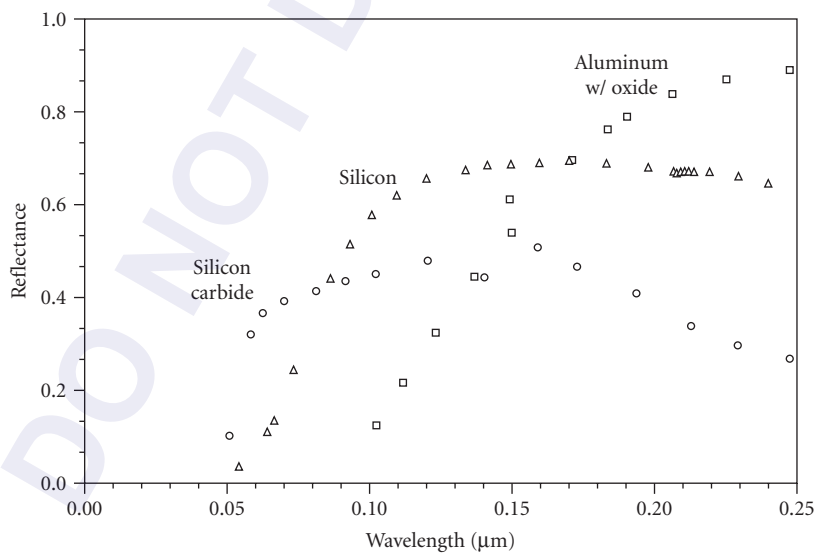


FIGURE 33 Ultraviolet reflectance of aluminum, silicon, and silicon carbide vs. wavelength.⁵¹

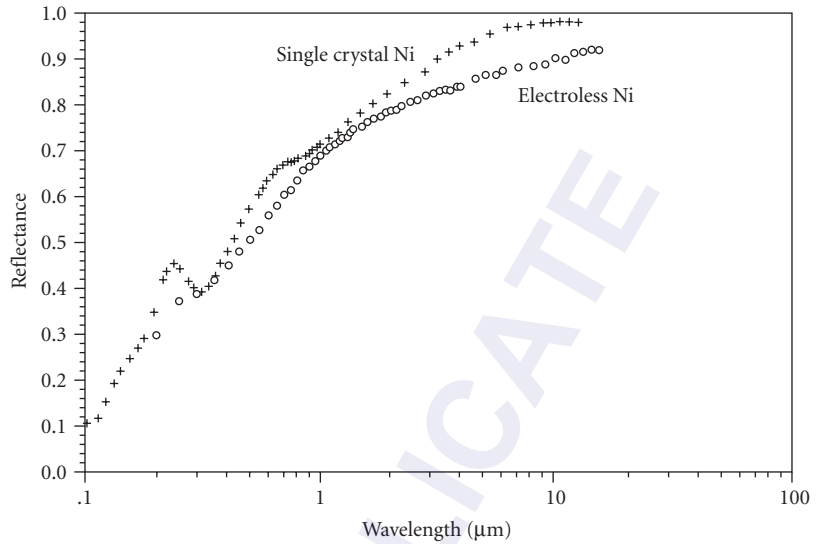


FIGURE 34 Reflectance of pure nickel³⁹ and electroless nickel (Ni-P alloy)⁵¹ vs. wavelength calculated from n and k .

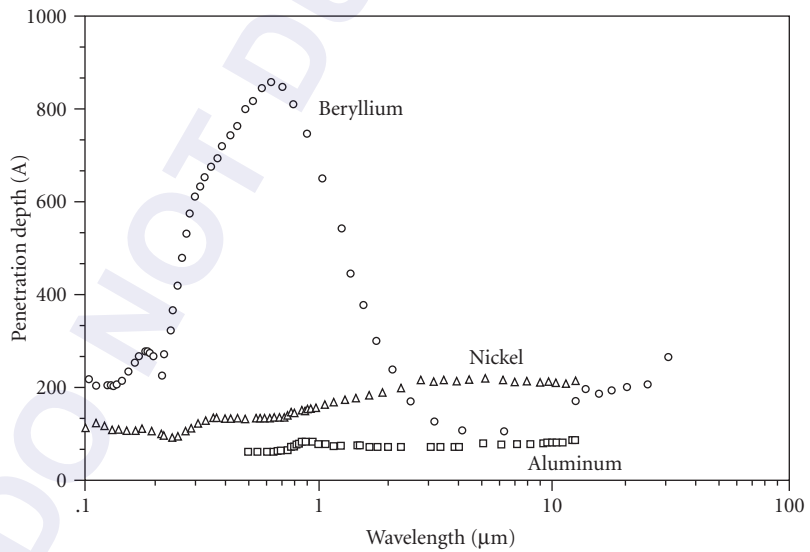


FIGURE 35 Penetration depth in Ångströms vs. wavelength for aluminum, beryllium, and nickel.⁵¹

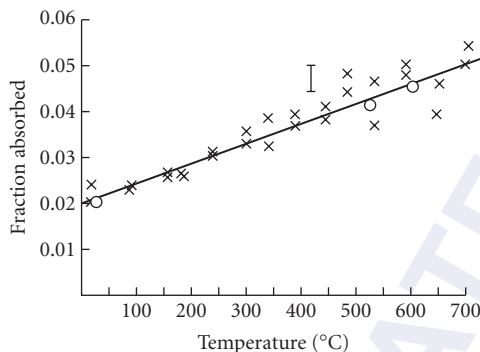


FIGURE 36 The 10.6- μm absorptance of Mo vs. temperature.⁵³ \times is heating and O is cooling. The straight line is a least-squares fit to the data.

has long been known,⁵⁵ but measurements and theory do not always agree, particularly at shorter wavelengths. Figure 36 shows absorptance of Mo as a function of temperature at a wavelength of 10.6 μm .⁵⁵

Mass absorption of energetic photons⁵⁶ follows the same relationship as described in Eq. (4), but with the product mass attenuation coefficient μ and mass density ρ substituted for absorption coefficient α . Table 4 lists mass attenuation coefficients for selected elements at energies between 1 keV (soft x rays) and 1 GeV (hard gamma rays). Units for the coefficient are m^2/kg , so that when multiplied by mass density in kg/m^3 , and depth x in m, the exponent in the equation is dimensionless. To a

TABLE 4 Mass Attenuation Coefficients for Photons⁵⁶

		Mass attenuation coefficient (m^2/kg)						
		Photon energy (MeV)						
Atomic No.		0.001	0.01	0.1	1.0	10.0	100.0	1000.0
Be	4	6.04×10^1	6.47×10^{-2}	1.33×10^{-2}	5.65×10^{-3}	1.63×10^{-3}	9.94×10^{-4}	1.12×10^{-3}
C	6	2.21×10^{-2}	2.37×10^{-1}	1.51×10^{-2}	6.36×10^{-3}	1.96×10^{-3}	1.46×10^{-3}	1.70×10^{-3}
O	8	4.59×10^2	5.95×10^{-1}	1.55×10^{-2}	6.37×10^{-3}	2.09×10^{-3}	1.79×10^{-3}	2.13×10^{-3}
Mg	12	9.22×10^1	2.11	1.69×10^{-2}	6.30×10^{-3}	2.31×10^{-3}	2.42×10^{-3}	2.90×10^{-3}
Al	13	1.19×10^2	2.62	1.70×10^{-2}	6.15×10^{-3}	2.32×10^{-3}	2.52×10^{-3}	3.03×10^{-3}
Si	14	1.57×10^2	3.39	1.84×10^{-2}	6.36×10^{-3}	2.46×10^{-3}	2.76×10^{-3}	3.34×10^{-3}
P	15	1.91×10^2	4.04	1.87×10^{-2}	6.18×10^{-3}	2.45×10^{-3}	2.84×10^{-3}	3.45×10^{-3}
Ti	22	5.87×10^2	1.11×10^1	2.72×10^{-2}	5.89×10^{-3}	2.73×10^{-3}	3.71×10^{-3}	4.56×10^{-3}
Cr	24	7.40×10^2	1.39×10^1	3.17×10^{-2}	5.93×10^{-3}	2.86×10^{-3}	4.01×10^{-3}	4.93×10^{-3}
Fe	26	9.09×10^2	1.71×10^1	3.72×10^{-2}	5.99×10^{-3}	2.99×10^{-3}	4.33×10^{-3}	5.33×10^{-3}
Ni	28	9.86×10^2	2.09×10^1	4.44×10^{-2}	6.16×10^{-3}	3.18×10^{-3}	4.73×10^{-3}	5.81×10^{-3}
Cu	29	1.06×10^3	2.16×10^1	4.58×10^{-2}	5.90×10^{-3}	3.10×10^{-3}	4.66×10^{-3}	5.72×10^{-3}
Zn	30	1.55×10^2	2.33×10^1	4.97×10^{-2}	5.94×10^{-3}	3.18×10^{-3}	4.82×10^{-3}	5.91×10^{-3}
Ge	32	1.89×10^2	3.74	5.55×10^{-2}	5.73×10^{-3}	3.16×10^{-3}	4.89×10^{-3}	6.00×10^{-3}
Mo	42	4.94×10^2	8.58	1.10×10^{-1}	5.84×10^{-3}	3.65×10^{-3}	6.10×10^{-3}	7.51×10^{-3}
Ag	47	7.04×10^2	1.19×10^1	1.47×10^{-1}	5.92×10^{-3}	3.88×10^{-3}	6.67×10^{-3}	8.20×10^{-3}
W	74	3.68×10^2	9.69	4.44×10^{-1}	6.62×10^{-3}	4.75×10^{-3}	8.80×10^{-3}	1.08×10^{-2}
Pt	78	4.43×10^2	1.13×10^1	4.99×10^{-1}	6.86×10^{-3}	4.87×10^{-3}	9.08×10^{-3}	1.12×10^{-2}
Au	79	4.65×10^2	1.18×10^1	5.16×10^{-1}	6.95×10^{-3}	4.93×10^{-3}	9.19×10^{-3}	1.13×10^{-2}

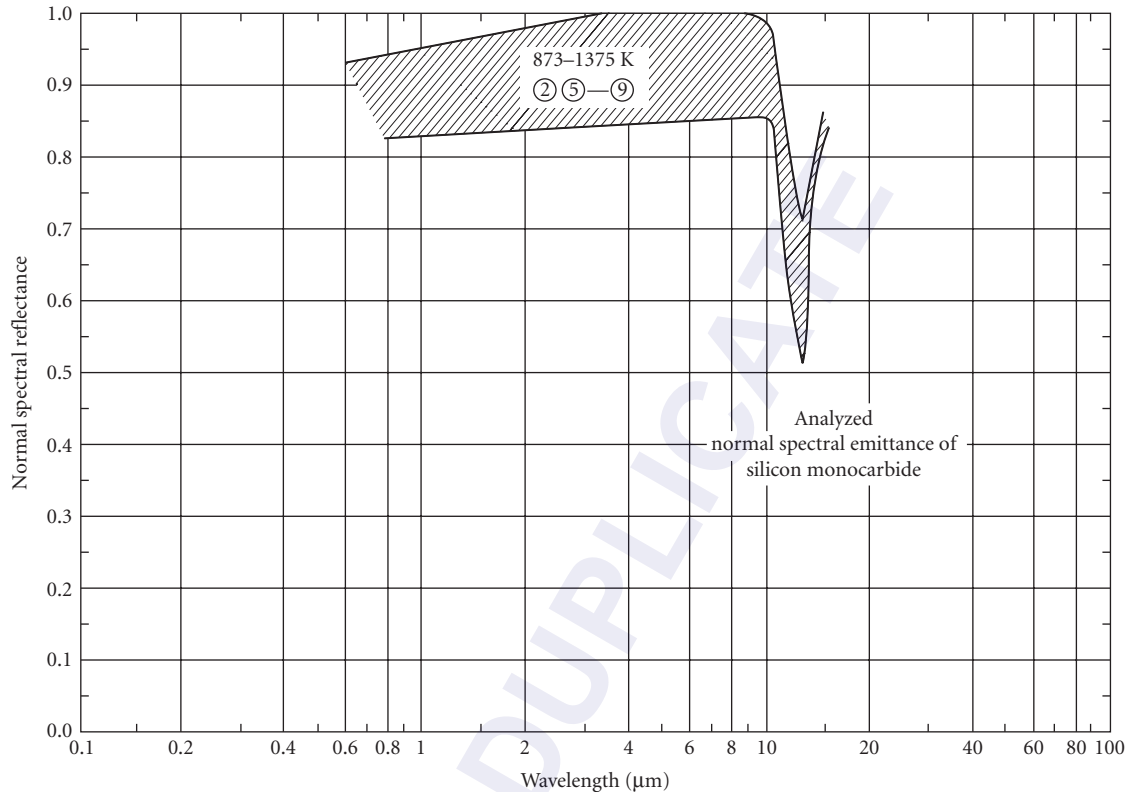


FIGURE 37 Analyzed normal spectral emittance of silicon carbide vs. wavelength.⁵⁷

high approximation, mass attenuation is additive for elements present in a body, independent of the way in which they are bound in chemical compounds. Table 4 is highly abridged; the original⁵⁶ shows all elements and absorption edges.

Emittance Where the transmittance of a material is essentially zero, the absorptance equals the emittance as described above and expressed in Eqs. (15) and (16). Spectral emittance ϵ_s is the emittance as a function of wavelength at constant temperature. These data have been presented as absorptance curves in Figs. 16 to 22, 25, and 26. For SiC, ϵ_s is given in Fig. 37.⁵⁷ Spectral emittance of unoxidized surfaces at a wavelength of 0.65 μm is given for selected materials in Table 5.⁵⁸

Total emittance ϵ_t is the emittance integrated over all wavelengths and usually given as a function of temperature. The total emittance of SiC is given in Fig. 38,⁵⁹ and for selected materials in Table 6.⁶⁰ Numerous papers by groups at the University of New Orleans (Ramanathan et al.⁶¹⁻⁶⁵) and at Cornell University (Sievers et al.^{66,67}) give high- and low-temperature data for the total hemispherical emittance of a number of metals including Ag, Al, Cu, Mo, W, and AISI 304 stainless steel.

Physical Properties

The physical properties at room temperature of a number of metals are listed in Table 7. The crystal form does not appreciably affect the physical properties, but is a factor in the isotropy of thermal and mechanical properties. For most metals, resistivity is directly proportional to temperature and

TABLE 5 Normal Spectral Emittance of Selected Metals ($\lambda = 0.65 \mu\text{m}$)⁵⁸

Metal	Emissivity
Beryllium	0.61
Chromium	0.34
Copper	0.10
Gold	0.14
Iron	0.35
Cast iron	0.37
Molybdenum	0.37
Nickel	0.36
80Ni-20Cr	0.35
Palladium	0.33
Platinum	0.30
Silver	0.07
Steel	0.35
Tantalum	0.49
Titanium	0.63
Tungsten	0.43

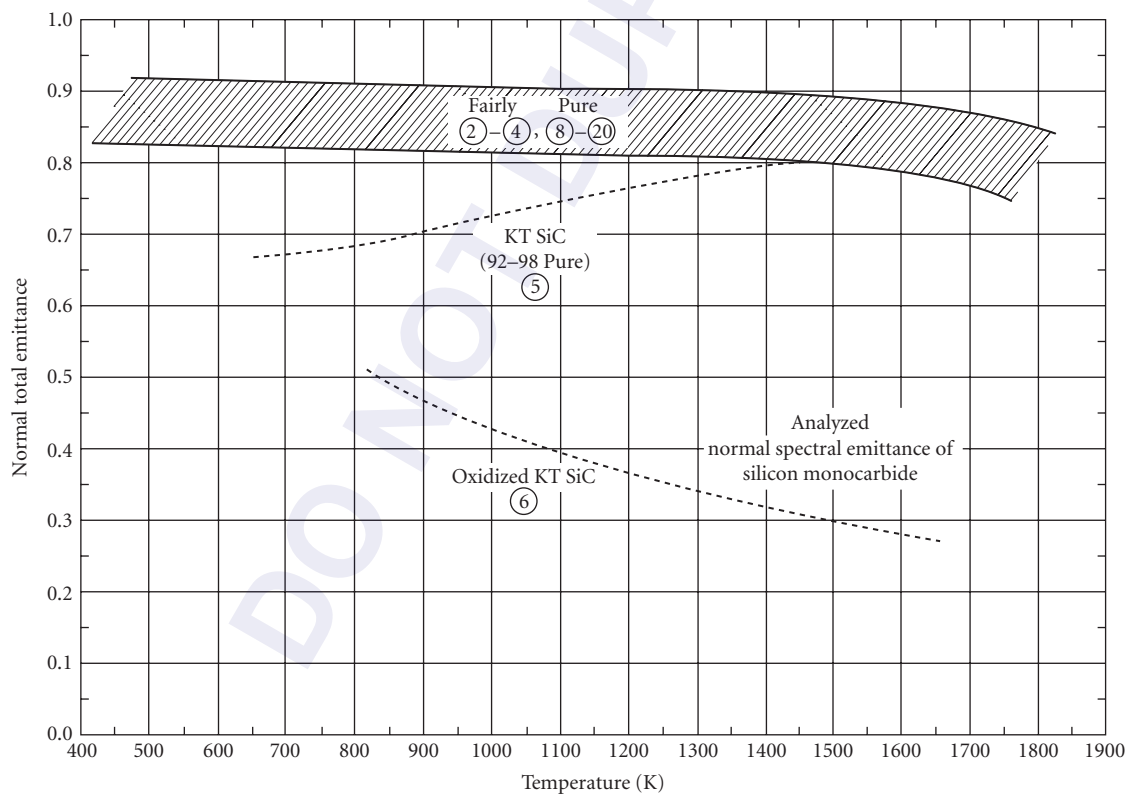
**FIGURE 38** Analyzed normal total emittance of silicon carbide vs. temperature.⁵⁹

TABLE 6 Total Emittance of Selected Materials⁶⁰

Metal	Temperature (°C)	Emissivity
80 Ni-20 Cr	100	0.87
	600	0.87
	1300	0.89
Aluminum	50-500	0.04-0.06
Polished	200	0.11
Oxidized	600	0.19
Chromium	50	0.1
Polished	500-1000	0.28-0.38
Copper Oxidized	50	0.6-0.7
	500	0.88
Polished	50-100	0.02
Unoxidized	100	0.02
Glass	20-100	0.94-0.91
	250-1000	0.87-0.72
	1100-1500	0.7-0.67
Gold		
Carefully polished	200-600	0.02-0.03
Unoxidized	100	0.02
Iron, cast		
Oxidized	200	0.64
	600	0.78
Unoxidized	100	0.21
Molybdenum	600-1000	0.08-0.13
	1500-2200	0.19-0.26
Nickel		
Polished	200-400	0.07-0.09
Unoxidized	25	0.045
	100	0.06
	500	0.12
	1000	0.19
Platinum		
Polished	200-600	0.05-0.1
Unoxidized	25	0.017
	100	0.047
	500	0.096
	1000	0.152
Silver		
Polished	200-600	0.02-0.03
Unoxidized	100	0.02
	500	0.035
Steel		
304 SS	500	0.35
Unoxidized	100	0.08
Tantalum, unoxidized	1500	0.21
	2000	0.26
Tungsten, unoxidized	25	0.024
	100	0.032
	500	0.071
	1000	0.15

TABLE 7 Composition and Physical Properties of Metals

Metal	Mass Density 10 ³ kg/m ³	Electrical Conductivity % IACS ^a	Electrical Resistivity Nohm m ^b	Crystal Form ^c	Chemical Composition Weight %, Typical	Reference
Aluminum: 5086-O	2.66	31	56	fcc	4.0 Mg, 0.4 Mn, 0.15 Cr, bal. Al	84
Aluminum: 6061-T6	2.70	43	40	fcc	1.0 Mg, 0.6 Si, 0.3 Cu, 0.2 Cr, bal. Al	84
Beryllium: I-70-H	1.85	43	40	cph	99.0 Be min., 0.6 BeO, 0.08 Fe, 0.05 C, 0.03 Al, 0.02 Mg	85
Copper: OFC	8.94	101	17	fcc	99.95 Cu min.	84
Gold	19.3	73	24	fcc	99.99 Au min.	84
Invar 36	8.1	<i>d</i>	820	bcc	36.0 Ni, 0.35 Mn, 0.2 Si, 0.02 C, bal. Fe	86
Molybdenum	10.22	34 ^e	52 ^e	bcc	99.9 Mo min., 0.015 C max.	84
Nickel: 200	8.9	18	95	fcc	99.0 Ni min.	84
Nickel: electroless plate	7.75	<i>d</i>	900	fcc	10.5 P, bal. Ni	87
Silicon	2.33	<i>f</i>	<i>f</i>	dia. cubic	99.99 Si	84
Silicon carbide (SiC): CVD	3.21	<i>f</i>	<i>f</i>	cubic	99.99 SiC (beta)	88
SiC: reaction sintered	2.91	<i>f</i>	<i>f</i>	cph + dia. cubic	74.0 SiC (alpha), 26.0 Si	88
Silver	10.49	103	15 ^e	fcc	99.9 Ag min.	84
Stainless steel: 304	8.00	<i>d</i>	720	fcc	19.0 Cr, 9.0 Ni, 1.0 Mn, 0.5 Si, bal. Fe	89
Stainless steel: 416	7.80	<i>d</i>	570	distorted bcc	13.0 Cr, 0.6 Mn, 0.6 Mo, 0.5 Si, bal. Fe	89
Stainless steel: 430	7.80	<i>d</i>	600	bcc	17.0 Cr, 0.5 Mn, 0.5 Si, bal. Fe	89
Titanium: 6Al4V	4.43	<i>d</i>	1710	bcc + cph	6.0 Al, 4.0 V, bal. Ti	90

^aFor equal volume at 293 K.^bAt 293 K.^cfcc = face-centered cubic; cph = close-packed hexagonal; bcc = body-centered cubic.^dNot available.^eAt 273 K.^fDepends on impurity content.

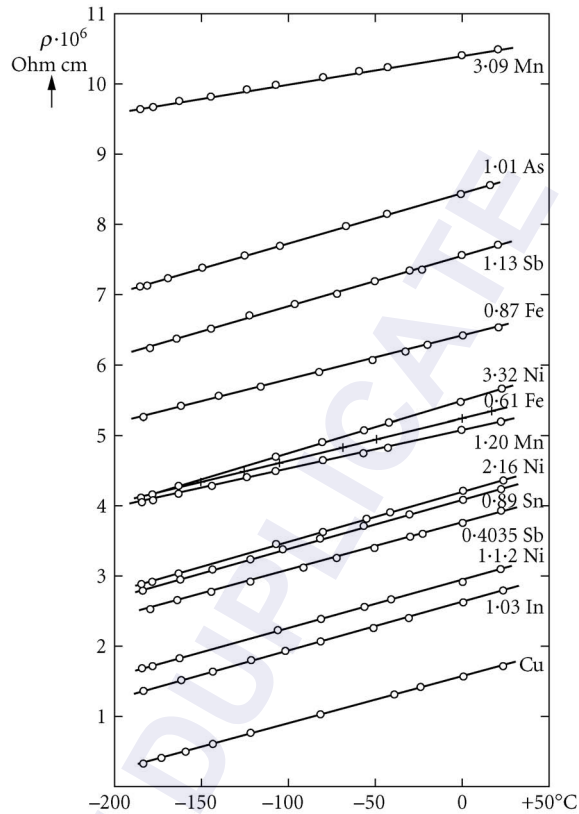


FIGURE 39 Electrical resistance of Cu and Cu alloys vs. temperature;⁶⁸ composition is in atomic percent.

pure metals generally have increased resistivity with increasing amounts of alloying elements. This is shown graphically for copper in Fig. 39.⁶⁸ Resistivity for a number of pure, polycrystalline metals is listed as a function of temperature in Table 8.⁶⁹

Thermal Properties

The thermal properties of materials were documented in 1970 through 1977 in the 13-volume series edited by Touloukian et al.⁷⁰ of the Thermophysical Properties Research Center at Purdue University. The properties database continues to be updated by the Center for Information and Numerical Data Analysis and Synthesis (CINDAS).¹

Selected properties of coefficient of thermal expansion, CTE, thermal conductivity k , and specific heat C_p at room temperature, are listed in Table 9. Maximum usable temperatures are also listed in the table.

The CTE of a material is a measure of length change at a specific temperature, useful for determining dimensional sensitivity to local temperature gradients. The total expansion (contraction) per unit length $\Delta L/L$ for a temperature change ΔT is the area under the CTE vs. T curve between the temperature extremes. Table 10 and Figs. 40 through 42 show recommended^{71,72} CTE vs. T relationships for a number of materials. More recent expansion data have been published for many

TABLE 8 Electrical Resistivity (nohm m) of Pure, Polycrystalline Metals⁶⁹

Temp. (K)	Aluminium	Beryllium	Chromium	Copper	Gold	Iron	Molybdenum	Nickel	Platinum	Silver	Tungsten
1	0.0010	0.332		0.020	0.220	0.225	0.0070	0.032	0.02	0.010	0.0002
10	0.0019	0.332		0.020	0.226	0.238	0.0089	0.057	0.154	0.012	0.0014
20	0.0076	0.336		0.028	0.350	0.287	0.0261	0.140	0.484	0.042	0.012
40	0.181	0.367		0.239	1.41	0.758	0.457	0.68	4.09	0.539	0.544
60	0.959	0.67		0.971	3.08	2.71	2.06	2.42	11.07	1.62	2.66
80	2.45	0.75		2.15	4.81	6.93	4.82	5.45	19.22	2.89	6.06
100	4.42	1.33	16.0	3.48	6.50	12.8	8.58	9.6	27.55	4.18	10.2
150	10.06	5.10	45.0	6.99	10.61	31.5	19.9	22.1	47.6	7.26	20.9
200	15.87	12.9	77.0	10.46	14.62	52.0	31.3	36.7	67.7	10.29	31.8
273	24.17	30.2	118.0	15.43	20.51	85.7	48.5	61.6	96.0	14.67	48.2
293	26.50	35.6	125.0	16.78	22.14	96.1	53.4	69.3	105.0	15.87	52.8
298	27.09	37.0	126.0	17.12	22.55	98.7	54.7	71.2	107.0	16.17	53.9
300	27.33	37.6	127.0	17.25	22.71	99.8	55.2	72.0	108.0	16.29	54.4
400	38.7	67.6	158.0	24.02	31.07	161.0	80.2	118.0	146.0	22.41	78.3
500	49.9	99.0	201.0	30.90	39.70	237.0	106.0	177.0	183.0	28.7	103.0
600	61.3	132.0	247.0	37.92	48.70	329.0	131.0	255.0	219.0	35.3	130.0
700	73.5	165.0	295.0	45.14	58.20	440.0	158.0	321.0	254.0	42.1	157.0
800	87.0	200.0	346.0	52.62	68.10	571.0	184.0	355.0	287.0	49.1	186.0
900	101.8	237.0	399.0	60.41	78.60		212.0	386.0	320.0	56.4	215.0

TABLE 9 Thermal Properties of Metals at Room Temperature

Metal	Coeff. of Thermal Expansion (ppm/K)	Thermal Conductivity (W/m K)	Specific Heat (J/kg K)	Maximum Temperature (K)	Reference
Aluminum: 5086-O	22.6	127	900	475	84
Aluminum: 6061-T6	22.5	167	896	425	84
Beryllium: I-70-H	11.3	216	1925	800	85
Copper: OFC	16.5	391	385	400	84
Gold	14.2	300	130	400	84
Iron	11.8	81	450	900	84
Invar 36	1.0	10	515	475	86
Molybdenum	4.8	142	276	1100	84
Nickel: 200	13.4	70	456	650	84
Nickel: Electroless plate (11% P)	11.0	7	460	425	87
(8% P)	12.8			450	91
Silicon	2.6	156	710	725	84
Silicon Carbide (SiC): CVD	2.2	198	733	1200	88
	2.4	250	700		92
SiC: Reaction sintered	2.6	155	670	1100	92
Silver	19.0	428	235	400	84
Stainless steel: 304	14.7	16	500	700	89
Stainless steel: 416	9.5	25	460	500	89
Stainless steel: 430	10.4	26	460	870	89
Titanium: 6A14V	8.6	7	520	650	84

materials that are too numerous to list here, but those for beryllium⁷³ and beta silicon carbide⁷⁴ are included in Table 10.

Thermal conductivities of many pure polycrystalline materials have been published by the National Bureau of Standards^{75,76} (now National Institute for Science and Technology) as part of the National Standard Reference Data System. Selected portions of these data, along with data from Touloukian et al.,^{77,78} and specific data for beryllium⁷⁹ and beta silicon carbide,⁸⁰ are listed in Table 11 and shown in Figs. 43 through 46.

The specific heat of metals is very well documented.^{73,81-83} Table 12 and Figs. 47 through 49 show the temperature dependence of this property. Table values are cited in J/kg K, numerically equal to W s/kg K.

Mechanical Properties

Mechanical properties are arbitrarily divided between the elastic properties of moduli, Poisson's ratio, and elastic stiffness, and the strength and fracture properties. All of these properties can be anisotropic as described by the elastic stiffness constants, but that level of detail is not included here. In general, cubic materials are isotropic in thermal properties and anisotropic in elastic properties. Materials of any of the other crystalline forms will be anisotropic in both thermal and elastic properties. For an in-depth treatment of this subject see, for example, Ref. 4.

TABLE 10 Temperature Dependence of the Coefficient of Linear Thermal Expansion (ppm/K) of Selected Materials

Temp. (K)	6061 Al	Be	Cu	Au	Fe	304 SS	416 SS	Mo	Ni	Ag	Si	Alpha SiC	Beta SiC
5		0.0003	0.005	0.03	0.01				0.02	0.015		0.01	
10		0.001										0.02	
20		0.005				9.8	4.3	0.3			0		
25		0.009	0.63	2.8	0.2			0.4	0.25	1.9	0	0.03	
50		0.096	3.87	7.7	1.3	10.5	4.9	1	1.5	8.2	-0.2	0.06	
75		0.47							4.3		-0.5	0.09	
100	12.2	1.32	10.3	11.8	5.6	11.4	6	2.8	6.6	14.2	-0.4	0.14	
125	18.7	2.55											
150	19.3	4.01				12.4	7				0.5	0.4	
175	20.3	5.54											
200	20.9	7.00	15.2	13.7	10.1	13.2	7.9	4.6	11.3	17.8	1.5	1.5	
225	21.5	8.32											
250	21.5	9.50				14.1	8.8				2.2	2.8	
293	22.5	11.3	16.5	14.2	11.8	14.7	9.5	4.8	13.4	18.9	2.6	3.3	3.26
300		11.5										3.4	3.29
350	23.8												3.46
400	25.0	13.6	17.6	14.8	13.4	16.3	10.9	4.9	14.5	19.7	3.2	4	3.62
450	26.3												3.77
500	27.5	15.1	18.3	15.4	14.4	17.5	12.1	5.1	15.3	20.6	3.5	4.2	3.92
600	30.1	16.6	18.9	15.9	15.1	18.6	12.9	5.3	15.9	21.5	3.7	4.5	4.19
700		17.8	19.5	16.4	15.7	19.5	13.5	5.5	16.4	22.6	3.9	4.7	4.42
800		19.1	20.3	17	16.2	20.2	13.8	5.7	16.8	23.7	4.1	4.9	4.62
900		20.0	21.3	17.7	16.4		13.9	6	17.1	24.8	4.3	5.1	4.79
1000		20.9	22.4	18.6	16.6	21.1	13.9	6.2	17.4	25.9	4.4	5.3	4.92
Reference	71	73	71	71	71	71	71	71	71	71	72	72	74

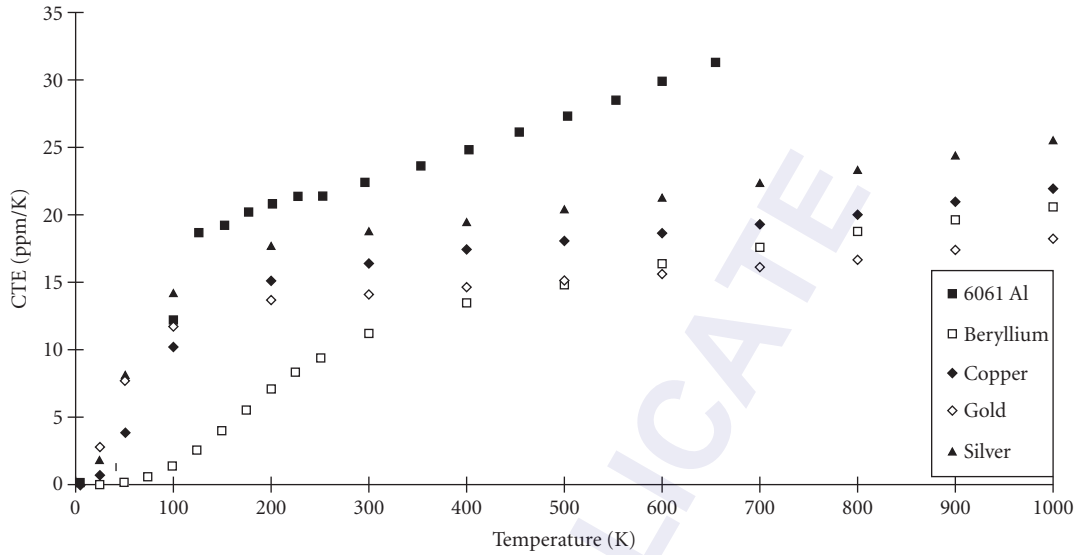


FIGURE 40 Coefficient of linear thermal expansion of 6061 aluminum alloy,⁷¹ beryllium,⁷³ copper,⁷¹ gold,⁷¹ and silver⁷¹ vs. temperature.

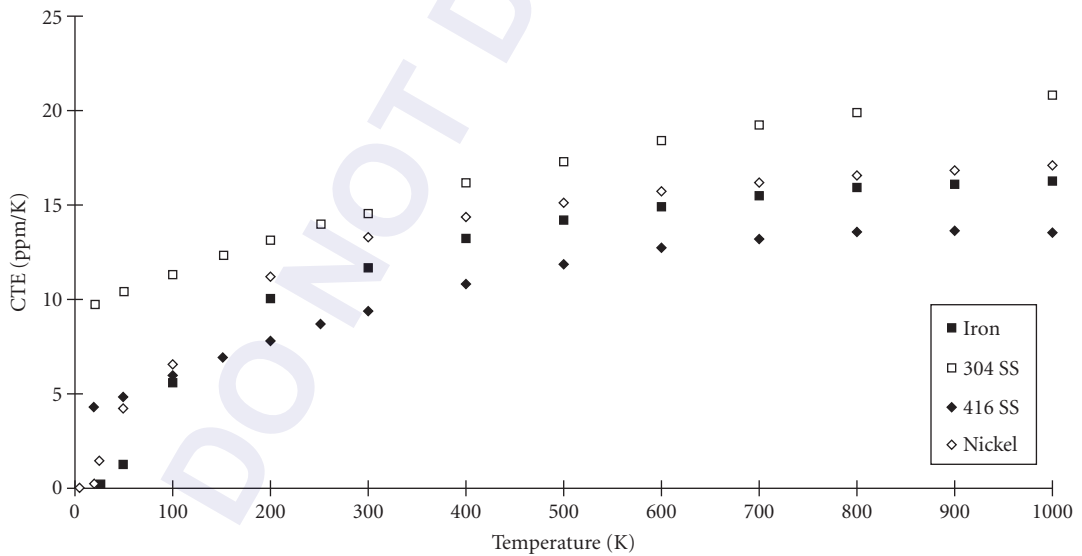


FIGURE 41 Coefficient of linear thermal expansion of iron,⁷¹ stainless steel types 304⁷¹ and 416,⁷¹ and nickel⁷¹ vs. temperature.

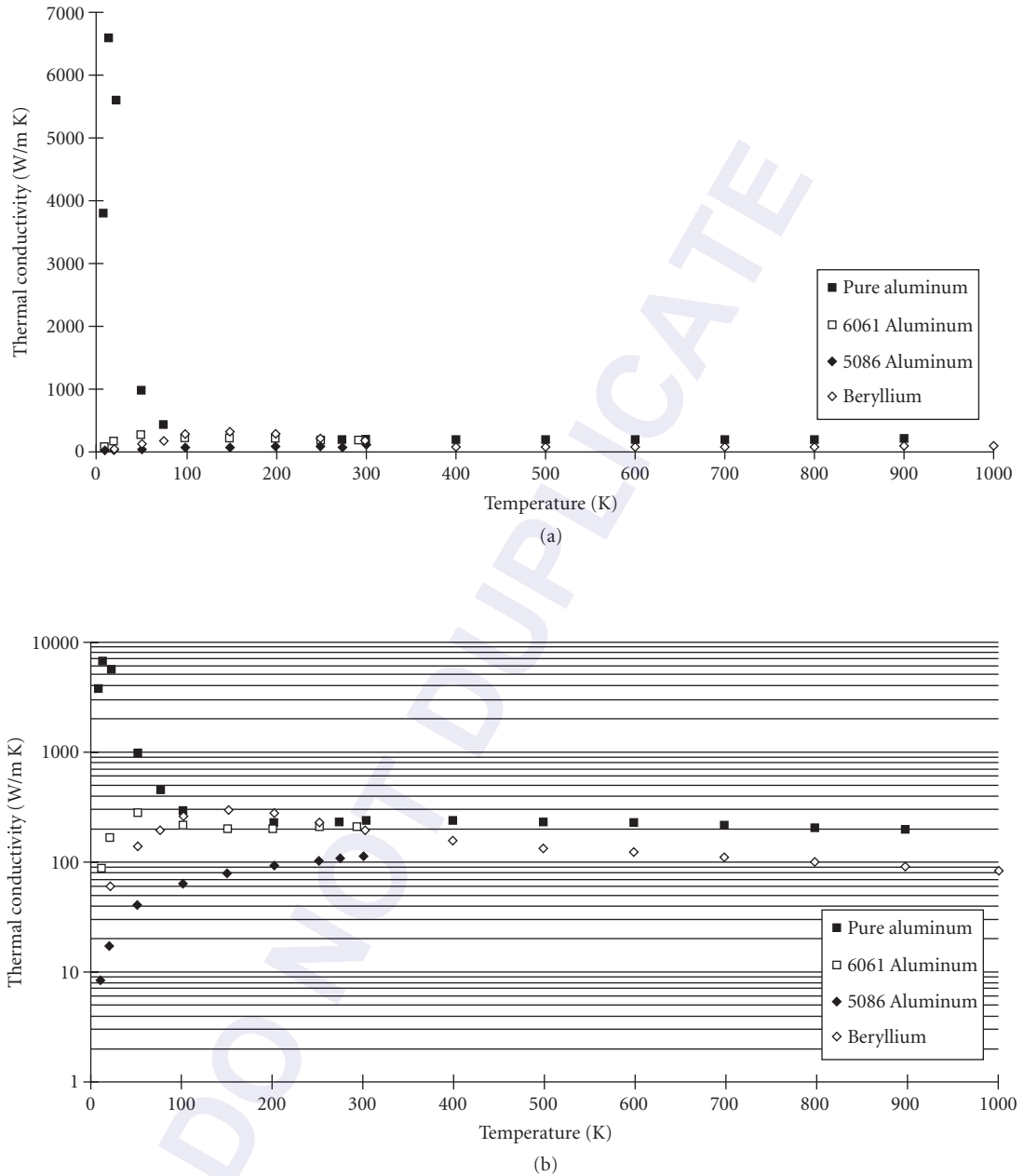


FIGURE 43 Thermal conductivity of three aluminum alloys⁷⁵⁻⁷⁷ and beryllium⁷⁹ vs. temperature.

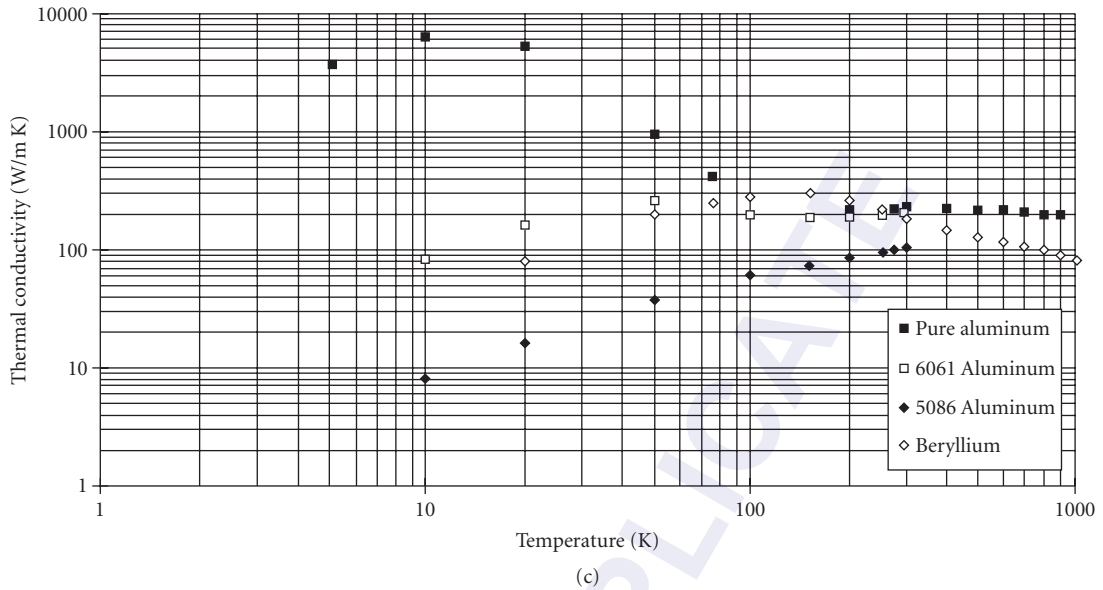


FIGURE 43 (Continued)

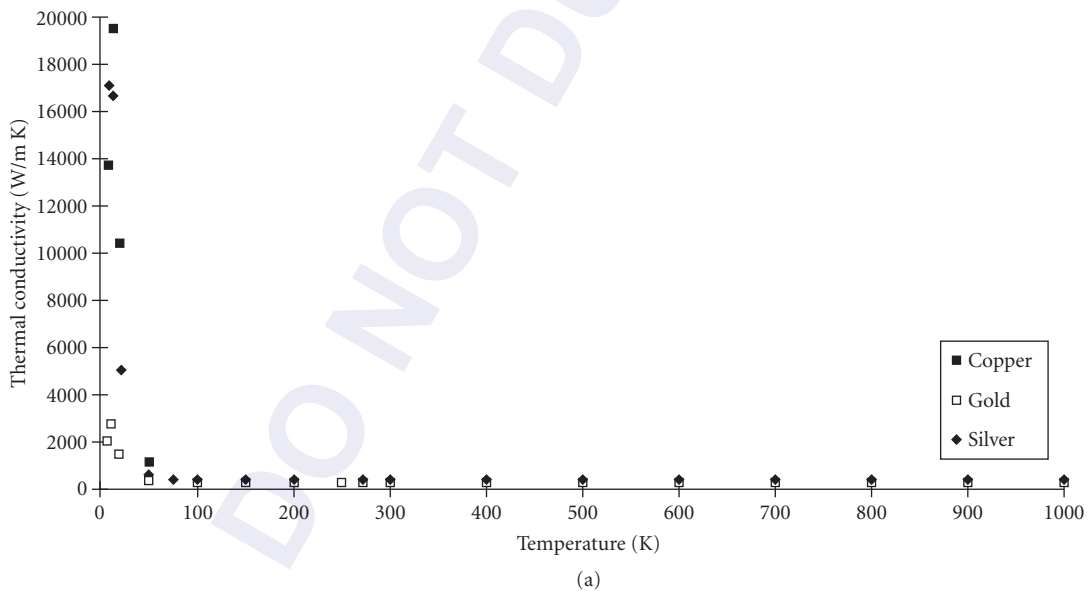


FIGURE 44 Thermal conductivity of copper,⁷⁷ gold,⁷⁷ and silver⁷⁷ vs. temperature.

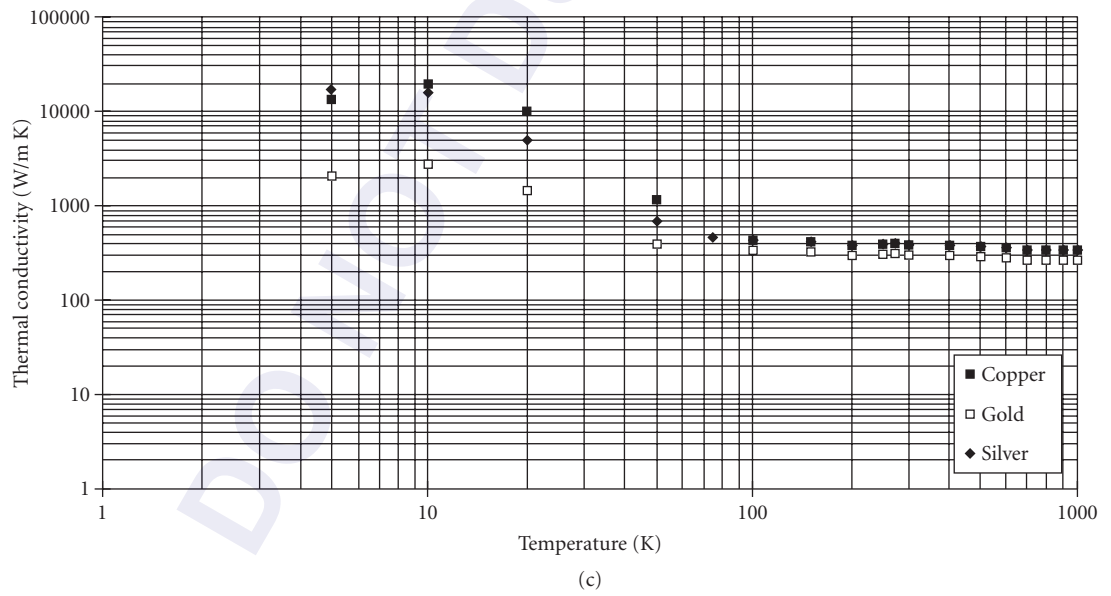
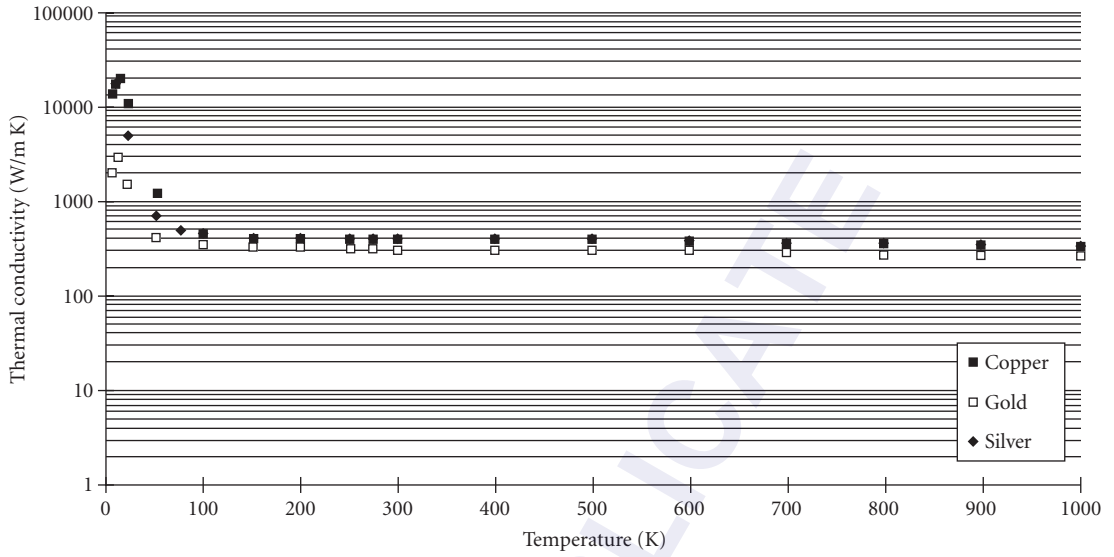


FIGURE 44 (Continued)

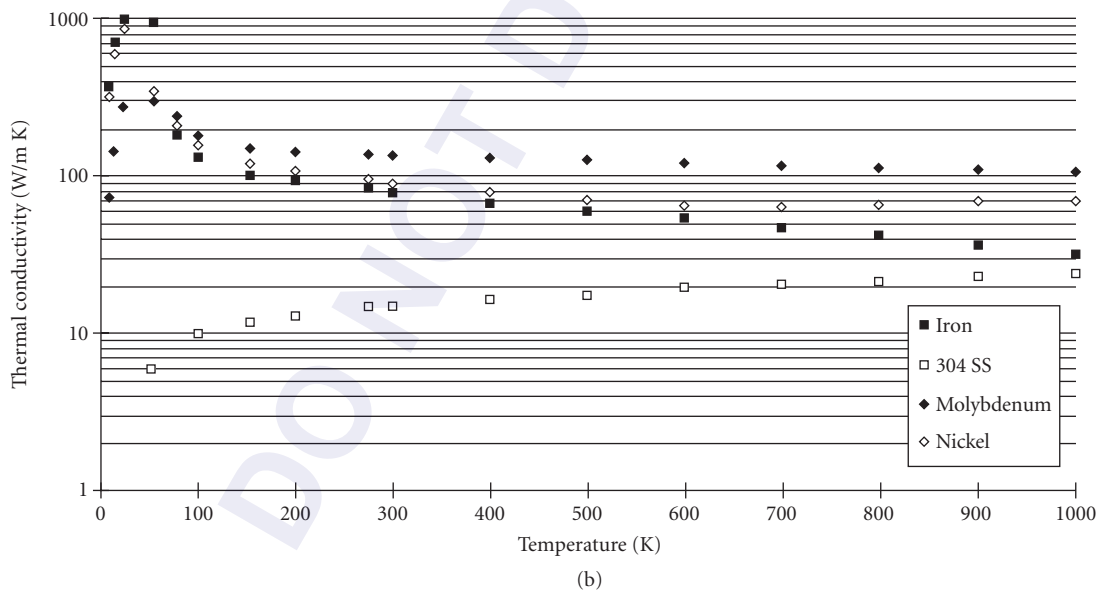
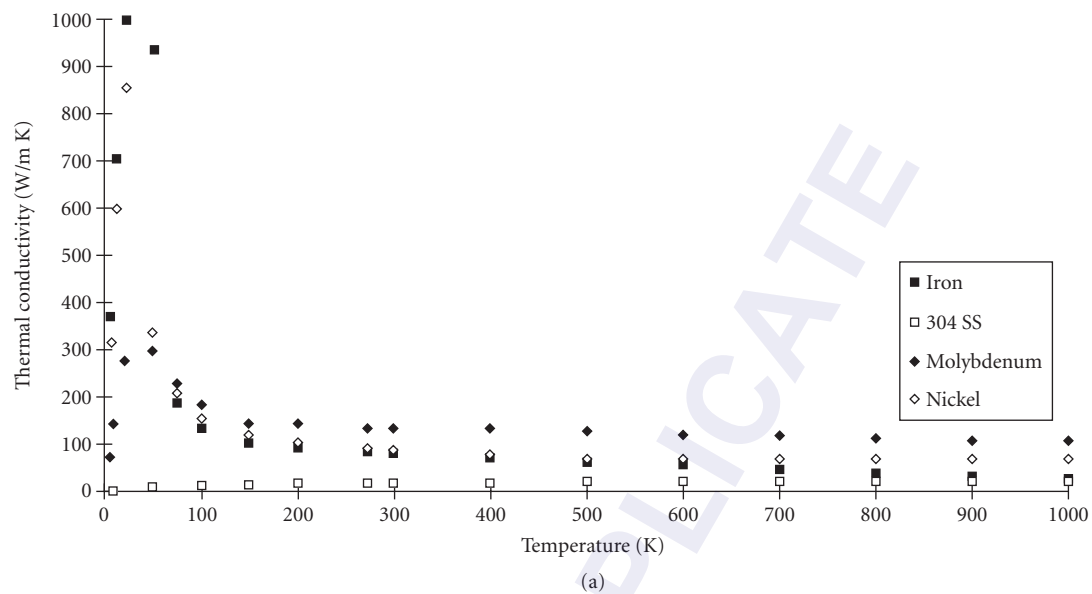


FIGURE 45 Thermal conductivity of iron,⁷⁷ type 304 stainless steel,⁷⁷ molybdenum,⁷⁷ and nickel⁷⁷ vs. temperature.

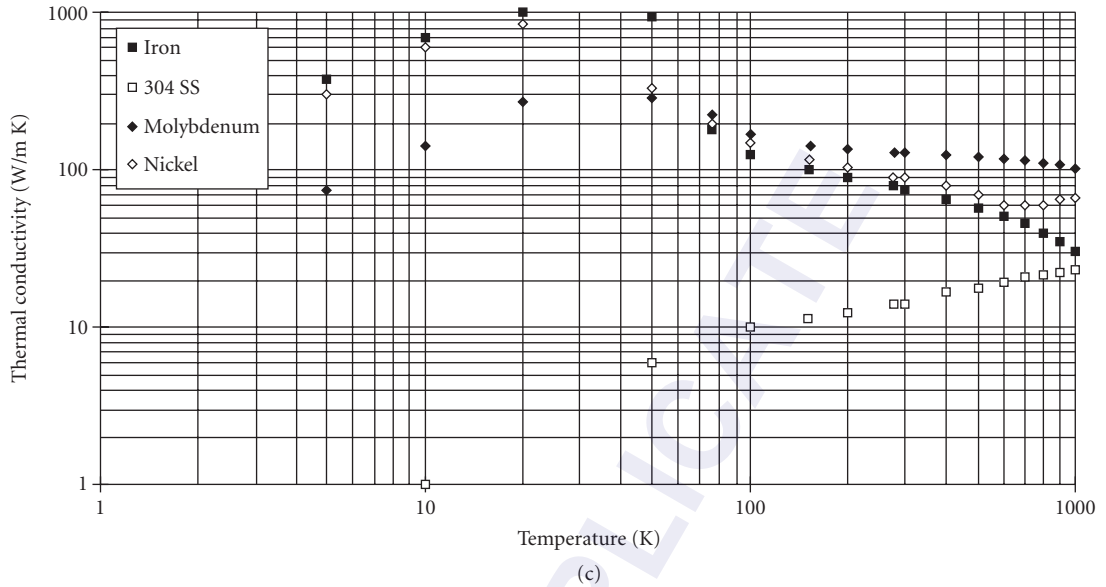


FIGURE 45 (Continued)

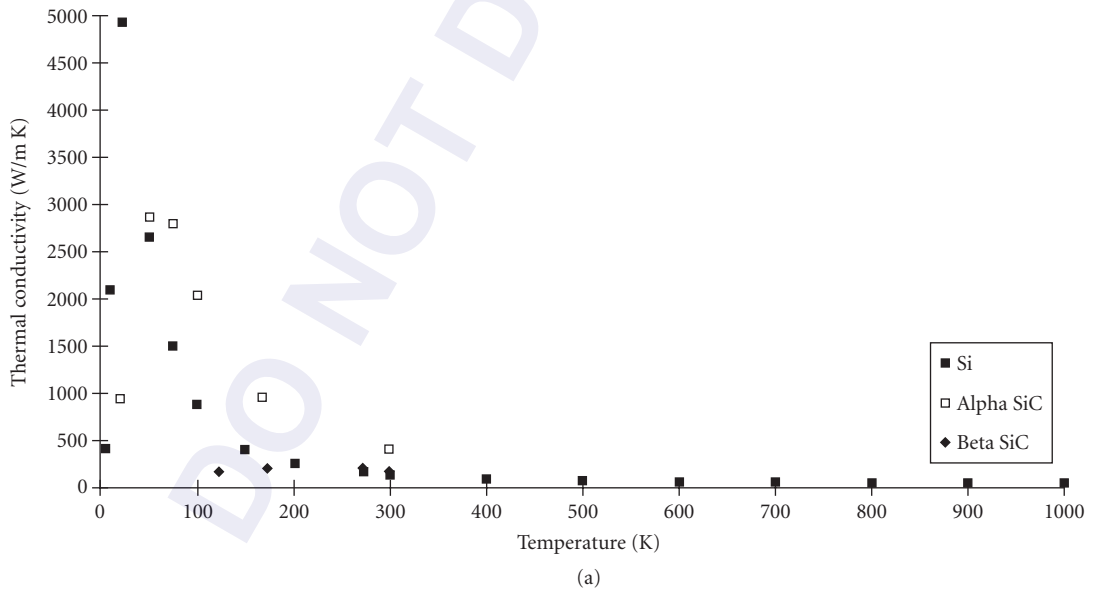
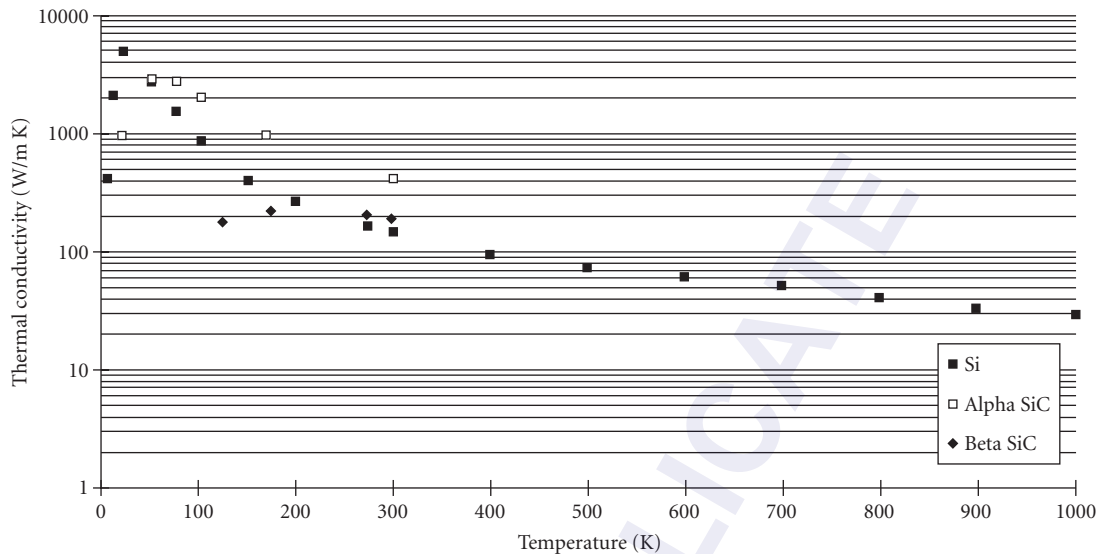
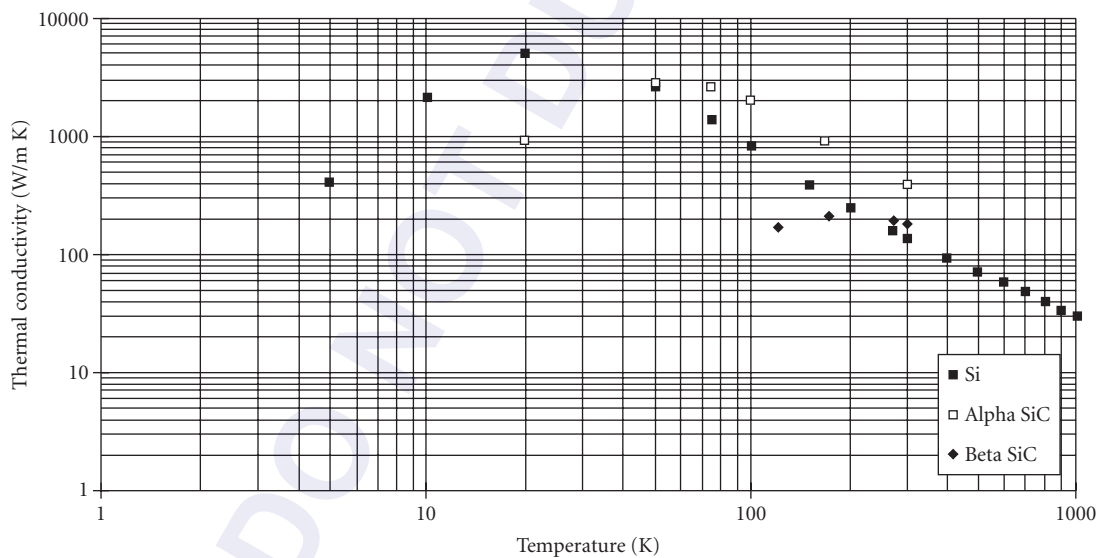


FIGURE 46 Thermal conductivity of silicon⁷⁷ and alpha⁷⁸ and beta⁸⁰ silicon carbide vs. temperature.

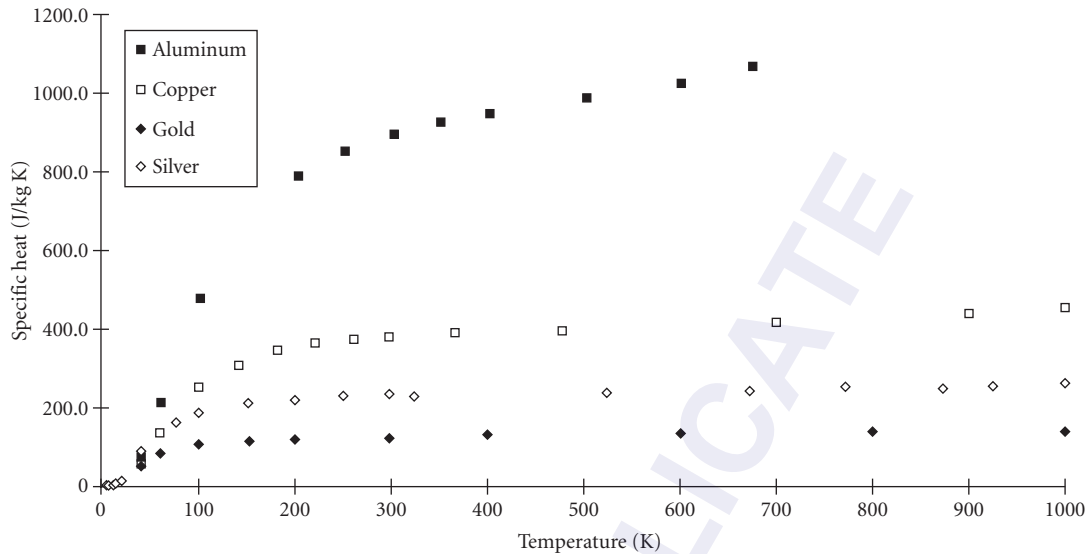


(b)

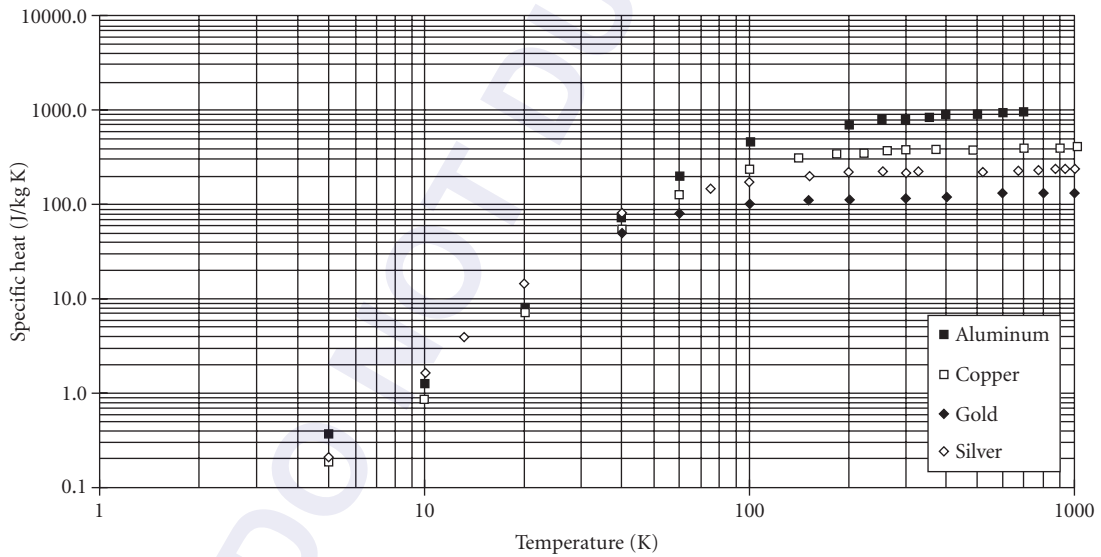


(c)

FIGURE 46 (Continued)

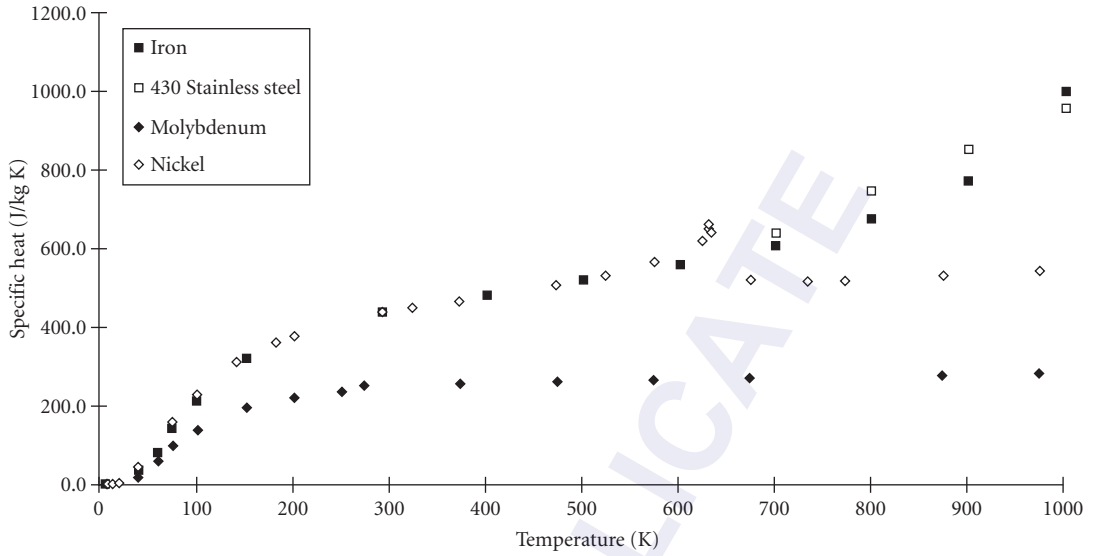


(a)

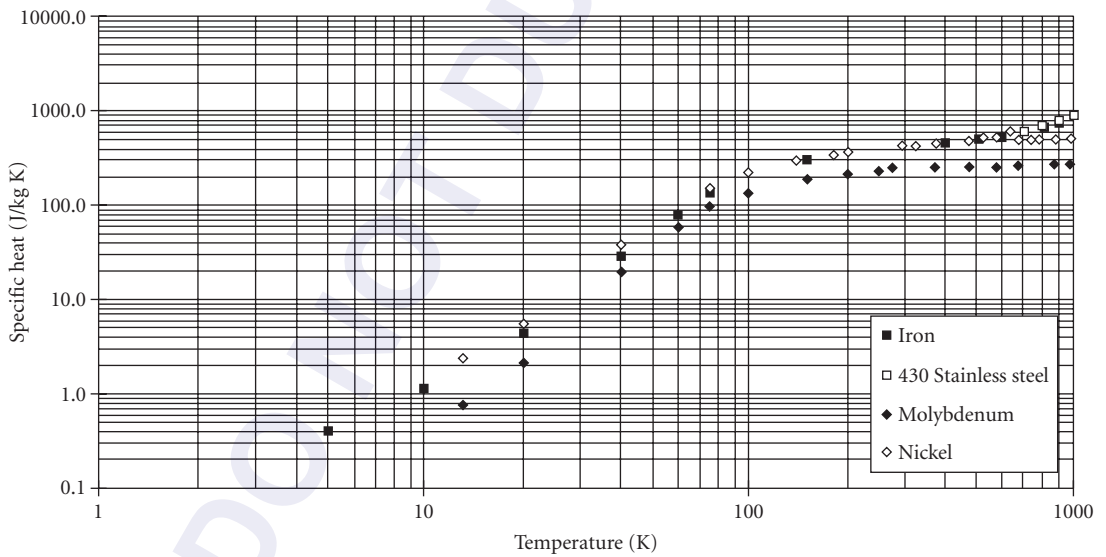


(b)

FIGURE 47 Specific heat of aluminum,^{81,82} copper,⁸³ gold,⁸³ and silver⁸³ vs. temperature.

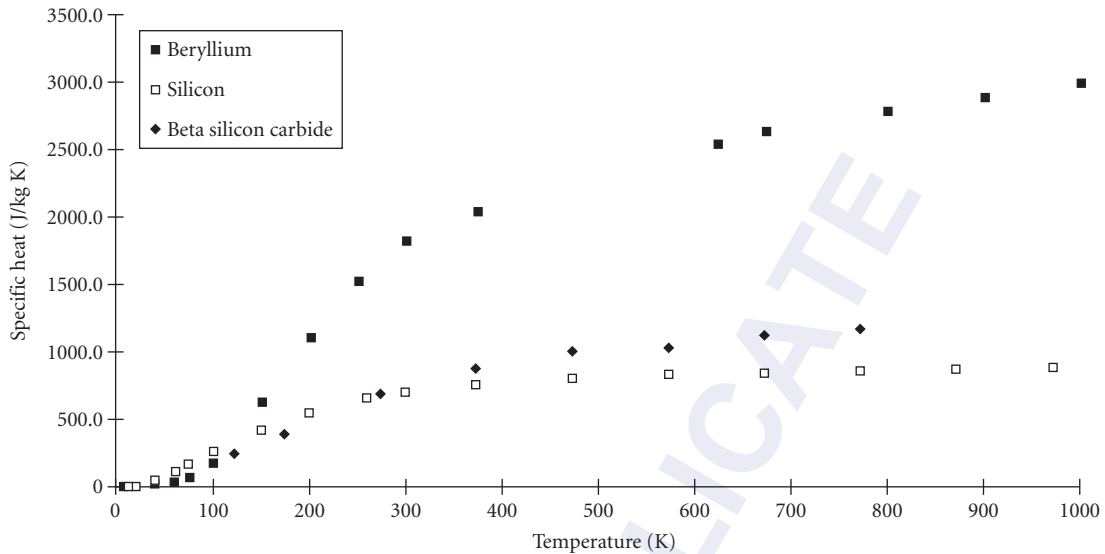


(a)

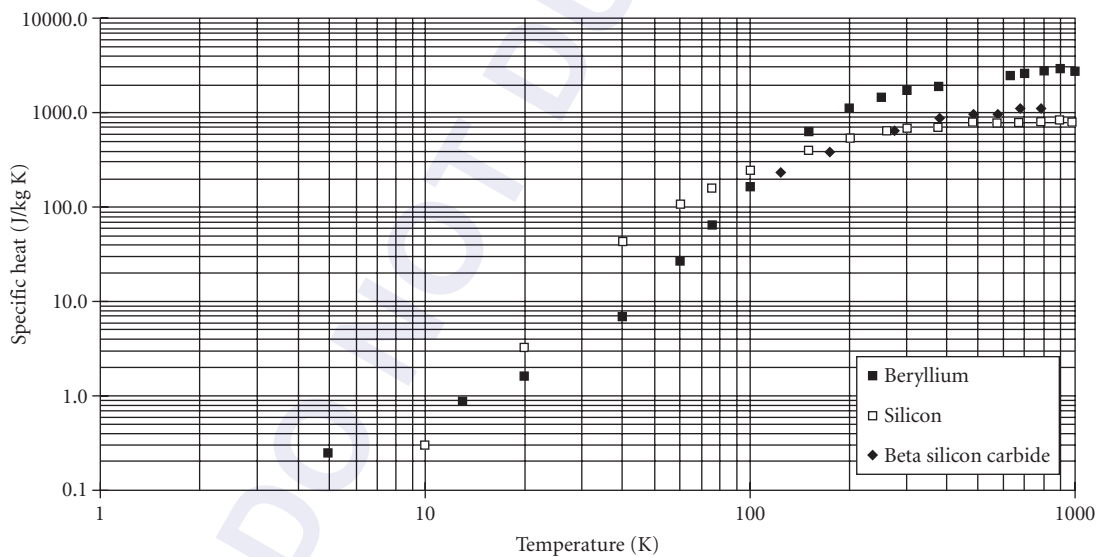


(b)

FIGURE 48 Specific heat of iron,⁸³ type 430 stainless steel,⁸³ molybdenum,⁸³ and nickel⁸³ vs. temperature.



(a)



(b)

FIGURE 49 Specific heat of beryllium,^{73,83} silicon,⁸³ and beta silicon carbide⁸⁰ vs. temperature.

TABLE 13 Elastic Stiffness Constants for Selected Single Crystal Metals

Cubic Metals ⁹³	Elastic stiffness (GN/m ²)				
	C_{11}	C_{44}	C_{12}		
Aluminum	108.0	28.3	62.0		
Chromium	346.0	100.0	66.0		
Copper	169.0	75.3	122.0		
Germanium	129.0	67.1	48.0		
Gold	190.0	42.3	161.0		
Iron	230.0	117.0	135.0		
Molybdenum	459.0	111.0	168.0		
Nickel	247.0	122.0	153.0		
Silicon	165.0	79.2	64.0		
Silicon carbide ⁹⁴	352.0	233.0	140.0		
Silver	123.0	45.3	92.0		
Tantalum	262.0	82.6	156.0		
Tungsten	517.0	157.0	203.0		
Hexagonal Metals	C_{11}	C_{33}	C_{44}	C_{12}	C_{13}
Beryllium ⁹⁵	288.8	354.2	154.9	21.1	4.7
Magnesium ⁹³	22.0	19.7	60.9	-7.8	-5.0
Silicon carbide ⁹⁴	500.0	521.0	168.0	98.0	

Elastic Properties The principal elastic stiffnesses C_{ij} of single crystals of some materials are given in Table 13. The three moduli and Poisson's ratio for polycrystalline materials are given in Table 14. These properties vary little with temperature, increasing temperature causing a gradual decrease in the moduli.

TABLE 14 Elastic Moduli and Poisson's Ratio for Selected Polycrystalline Materials

Materials	Young's Modulus (GN/m ²)	Shear Modulus (GN/m ²)	Bulk Modulus (GN/m ²)	Poisson's Ratio	Reference
Aluminum: 5086-O	71.0	26.4		0.33	84
Aluminum: 6061-T6	68.9	25.9		0.33	84
Beryllium: I-701-H	315.4	148.4	115.0	0.043	96
Copper	129.8	48.3	137.8	0.343	97
Germanium	79.9	29.6		0.32	97
Gold	78.5	26.0	171.0	0.42	97
Invar 36	144.0	57.2	99.4	0.259	97
Iron	211.4	81.6	169.8	0.293	97
Molybdenum	324.8	125.6	261.2	0.293	97
Nickel	199.5	76.0	177.3	0.312	97
Platinum	170.0	60.9	276.0	0.39	97
Silicon	113.0	39.7		0.42	97
Silicon carbide: CVD	461.0			0.21	80
Silicon carbide: reaction sintered	413.0			0.24	88
Silver	82.7	30.3	103.6	0.367	97
Stainless steel: 304	193.0	77.0		0.27	97
Stainless steel: 416	215.0	83.9	166.0	0.283	97
Stainless steel: 430	200.0	80.0		0.27	97
Tantalum	185.7	62.2	196.3	0.342	97
Tungsten	411.0	160.6	311.0	0.28	97

TABLE 15 Strength and Fracture Properties for Selected Materials

Material	Yield Strength (MN/m ²)	Microyield Strength (MN/m ²)	Elongation (in 50 mm) %	Fracture Toughness (MN m ^{-3/2})	Flexural Strength (MN/m ²)	Hardness*	Reference
Aluminum: 5086-O	115.0	40.0	22.0	>25.0	—	55 HRB	84
Aluminum: 6061-T6	276.0	160.0	15.0	<25.0	—	95 HRB	84
Beryllium: I-70-H	276.0	30.0	4.0	12.0	—	80 HRB	84,85
Copper	195.0	12.0	42.0	—	—	10 HRB	84
Germanium	—	—	—	1.0	110.0	800 HK	84
Gold	125.0	—	30.0	—	—	30 HK	84
Invar 36	276.0	37.0	35.0	—	—	70 HRB	86
Molybdenum	600.0	—	40.0	—	—	150 HK	84
Nickel	148.0	—	47.0	—	—	109 HRB	84
Platinum	150.0	—	35.0	—	—	40 HK	84
Silicon	—	—	—	1.0	207.0	1150 HK	84
Silicon carbide: CVD	—	—	—	3.0	595.0	2500 HK	80
Silicon carbide: reaction sintered	—	—	—	2.0	290.0	2326 HK	88
Silver	130.0	—	47.0	—	—	32 HK	84
Stainless steel: 304	241.0	—	60.0	—	—	80 HRB	89
Stainless steel: 416	950.0	—	12.0	—	—	41 HRC	89
Stainless steel: 430	380.0	—	25.0	—	—	86 HRB	89
Tantalum	220.0	—	30.0	—	—	120 HK	84
Tungsten	780.0	—	2.0	—	—	350 HK	84

*HK = Knoop (kg/mm²); HRB = Rockwell B; HRC = Rockwell C.

Strength and Fracture Properties The properties of tensile yield (at 0.2 percent offset), microyield strength, ductility (expressed as percent elongation in 50 mm), fracture toughness, flexural strength, and mechanical hardness are listed in Table 15. Most of these properties vary with temperature: strength and hardness decreasing, and fracture toughness and ductility increasing with temperature.

4.4 REFERENCES

- Center for Information and Numerical Data Analysis and Synthesis (CINDAS), Purdue Univ., 2595 Yeager Rd., W. Lafayette, IN 47906, (800) 428-7675.
- Optical Properties of Solids and Liquids (OPTROP), Sandia National Laboratory, Div. 1824, P.O. Box 5800, Albuquerque, NM 87185, (505) 844-2109.
- H. Wawrousek, J. H. Westbrook, and W. Grattideg (eds.), "Data Sources of Mechanical and Physical Properties of Engineering Materials," *Physik Daten/Physics Data*, No. 30-1, Fachinformationszentrum Karlsruhe, 1989.
- M. E. Lines, "Physical Properties of Materials: Theoretical Overview," in Paul Klocek (ed.), *Handbook of Infrared Optical Materials*, Marcel Dekker, New York, 1991, pp. 1-69.
- N. F. Mott and H. Jones, *The Theory of The Properties of Metals and Alloys*, Dover, New York, 1958, pp. 105-125.
- A. V. Sokolov, *Optical Properties of Metals*, American Elsevier, New York, 1967.
- F. Wooten, *Optical Properties of Solids*, Academic Press, New York, 1972.
- M. Born and E. Wolf, *Principles of Optics*, 5th ed., Pergamon Press, London, 1975, pp. 611-627.

9. This discussion is adapted with permission from M. Bass, "Laser-Materials Interactions," in *Encyclopedia of Physical Science and Technology* **8**, Academic Press, New York, 1992, pp. 415–418.
10. F. Stern in F. Seitz and D. Turnbull (eds.), *Solid State Physics*, vol. 15, Academic Press, New York, 1963, pp. 300–324.
11. L. D. Landau and E. M. Lifshitz, *Electrodynamics of Continuous Media*, Addison-Wesley, Reading, Mass., 1965, pp. 257–284.
12. H. A. Lorentz, *The Theory of Electrons*, Dover, New York, 1952.
13. P. K. L. Drude, *Theory of Optics*, Dover, New York, 1959.
14. C. S. Barrett, *Structure of Metals*, 2d ed., McGraw-Hill, New York, 1952, pp. 521–537.
15. As reported by S. F. Jacobs in "Variable Invariables—Dimensional Instability with Time and Temperature," P. R. Yoder, Jr. (ed.), *Optomechanical Design*, Critical Reviews of Optical Science and Technology, **CR43**, SPIE Optical Engineering Press, Bellingham, Wash., 1992, p. 201.
16. C. Kittel, *Introduction to Solid State Physics*, 3d ed., Wiley, New York, 1966, pp. 111–129.
17. H. Reisman and P. S. Pawlik, *Elasticity*, Wiley, New York, 1980, pp. 128–135.
18. A. Kelly and N. H. Macmillan, *Strong Solids*, 3d ed., Clarendon Press, Oxford, 1986, pp. 382–393.
19. For more complete descriptions see, for example, *Metals Handbook*, 9th ed., **8**, Mechanical Testing, American Society for Metals, Metals Park, OH, 1985, pp. 1–15.
20. R. A. Paquin, "Selection of Materials and Processes for Metal Optics," in *Selected Papers on Optomechanical Design*, *Proc. SPIE*, Milestone Series, **770:27–34** (1987).
21. D. Janeczko, "Metal Mirror Review," in R. Hartmann and W. J. Smith (eds.), *Infrared Optical Design*, Critical Reviews of Optical Science and Technology, **CR38**, SPIE Optical Engineering Press, Bellingham, Wash., 1991, pp. 258–280.
22. M. H. Krim, "Mechanical Design of Optical Systems for Space Operation," in P. R. Yoder, Jr. (ed.), *Optomechanical Design*, Critical Reviews of Optical Science and Technology, **CR43**, SPIE Optical Engineering Press, Bellingham, Wash., 1992, pp. 3–17.
23. P. R. Yoder, Jr. *Opto-Mechanical Systems Design*, 2d ed., Marcel Dekker, New York, 1993, pp. 1–41.
24. This analysis is the same as that used by many structural engineers such as the late G. E. Seibert of Perkin-Elmer and Hughes Danbury Optical Systems.
25. S. Timoshenko and S. Woinowsky-Kreiger, *Theory of Plates and Shells*, 2d ed., McGraw-Hill, New York, 1959, pp. 51–78.
26. P. K. Mehta, "Nonsymmetric Thermal Bowing of Curved Circular Plates," in A. E. Hatheway (ed.), *Structural Mechanics of Optical Systems II*, *Proc. SPIE* **748** (1987).
27. E. Pearson, "Thermo-elastic Analysis of Large Optical Systems," in P. R. Yoder, Jr. (ed.), *Optomechanical Design*, Critical Reviews of Optical Science and Technology, **CR43**, SPIE Optical Engineering Press, Bellingham, Wash., 1992, pp. 123–130.
28. C. W. Marschall and R. E. Maringer, *Dimensional Instability, an Introduction*, Pergamon, New York, 1977.
29. R. A. Paquin (ed.), "Dimensional Stability," *Proc. SPIE* **1335** (1990).
30. R. A. Paquin and D. Vukobratovich (eds.), "Optomechanics and Dimensional Stability," *Proc. SPIE* **1533** (1991).
31. R. A. Paquin, "Dimensional Instability of Materials; How Critical Is It in the Design of Optical Instruments?," in P. R. Yoder (ed.), *Optomechanical Design*, Critical Reviews of Optical Science and Technology, **CR43**, SPIE Optical Engineering Press, Bellingham, Wash., 1992, pp. 160–180.
32. Op. cit. Ref 23, pp. 271–320, 567–584.
33. E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, Academic Press, Orlando, 1985.
34. E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991.
35. D. W. Lynch, "Mirror and Reflector Materials," in M. J. Weber (ed.), *CRC Handbook of Laser Science and Technology*, **IV**, Optical Materials, Part 2: Properties, CRC Press, Boca Raton, Florida, 1986.
36. J. H. Weaver, C. Krafska, D. W. Lynch, and E. E. Koch (eds.), "Optical Properties of Metals," Pts. 1 and 2, *Physik Daten/Physics Data*, Nos. 18-1 and 18-2, Fachinformationszentrum Karlsruhe, 1981.
37. D. Y. Smith, E. Shiles, and Mitio Inokuti, "The Optical Properties of Aluminum," in E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, Academic Press, Orlando, 1985, pp. 369–406.

38. E. T. Arakawa, T. A. Callcott, and Y.-C. Chang, "Beryllium," in E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991, pp. 421–433.
39. D. W. Lynch and W. R. Hunter, in E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, Academic Press, Orlando, 1985, pp. 275–367.
40. D. W. Lynch and W. R. Hunter, in E. D. Palik (ed.), *Handbook of Optical Constants of Solids II*, Academic Press, Orlando, 1991, pp. 341–419.
41. W. J. Choyke and E. D. Palik, "Silicon Carbide," in E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, Academic Press, Orlando, 1985, pp. 587–595.
42. J. H. Weaver, D. W. Lynch, and R. Rossi, "Optical Properties of Single-Crystal Be from 0.12 to 4.5 eV," *Phys. Rev. B* **7**:3537–3541 (1973).
43. J. C. Stover, J. Rifkin, D. R. Cheever, K. H. Kirchner, and T. F. Schiff, "Comparison of Wavelength Scaling Data to Experiment," in R. P. Breault (ed.), *Stray Light and Contamination in Optical Systems, Proc. SPIE* **967**:44–49 (1988).
44. C. L. Vernold, "Application and Verification of Wavelength Scaling for Near Specular Scatter Predictions," in J. C. Stover (ed.), *Scatter from Optical Components, Proc. SPIE* **1165**:18–25 (1989).
45. J. E. Harvey, "Surface Scatter Phenomena: A Linear, Shift-Invariant Process," in J. C. Stover (ed.), *Scatter from Optical Components, Proc. SPIE* **1165**:87–99 (1989).
46. J. C. Stover, M. L. Bernt, D. E. McGary, and J. Rifkin, "An Investigation of Anomalous Scatter from Beryllium Mirrors," in J. C. Stover (ed.), *Scatter from Optical Components, Proc. SPIE* **1165**:100–109 (1989).
47. See also papers in the "Scatter from Be Mirrors" session in J. C. Stover (ed.), *Optical Scatter: Applications, Measurement, and Theory, Proc. SPIE* **1530**:130–230 (1991).
48. Y. S. Touloukian and D. P. DeWitt, "Thermal Radiative Properties, Metallic Elements and Alloys," vol. 7 in Y. S. Touloukian and C. Y. Ho (eds.), *Thermophysical Properties of Matter*, IFI/Plenum, New York, 1970.
49. Y. S. Touloukian and D. P. DeWitt, "Thermal Radiative Properties, Nonmetallic Solids," vol. 8 in Y. S. Touloukian and C. Y. Ho (eds.), *Thermophysical Properties of Matter*, IFI/Plenum, New York, 1971.
50. J. S. Browder, S. J. Ballard, and P. Klocek in Paul Klocek (ed.), *Handbook of Infrared Optical Materials*, Marcel Dekker, New York, 1991, pp. 155–426.
51. C. M. Egert, "Optical Properties of Aluminum, Beryllium, Silicon Carbide (and more)" in *Proc. of Al, Be, and SiC Optics Technologies Seminar*, MODIL, Oak Ridge National Lab., 1993.
52. W. D. Kimura and D. H. Ford, "Absorptance Measurement of Metal Mirrors at Glancing Incidence," *Appl. Optics* **25**:3740–3750 (1986).
53. M. Bass and L. Liou, "Calorimetric Studies of Light Absorption by Diamond Turned Ag and Cu Surfaces and Analyses Including Surface Roughness Contributions," *J. Appl. Phys.* **56**:184–189 (1984).
54. W. D. Kimura and T. T. Saito, "Glancing Incidence Measurements of Diamond Turned Copper Mirrors," *Appl. Optics* **26**:723–728 (1987).
55. M. Bass, D. Gallant, and S. D. Allen, "The Temperature Dependence of the Optical Absorption of Metals," in *Basic Optical Properties of Materials*, NBS SP574, U.S. Govt. Printing Office, Wash. D.C., 1980, pp. 48–50.
56. This discussion and Table 4 are based on the article: M. J. Berger and J. H. Hubbell, "Photon Attenuation Coefficients," in D. R. Lide (editor-in-chief), *CRC Handbook of Chemistry and Physics*, 74th ed., CRC Press, Boca Raton, Fla., 1993, pp. 10-282-10-286.
57. Op. cit., Ref. 49, p. 798.
58. D. R. Lide (editor-in-chief), *CRC Handbook of Chemistry and Physics*, 74th ed., CRC Press, Boca Raton, Fla., 1993, p. 10–299.
59. Op. cit., Ref. 49, p. 792.
60. Op. cit., Ref. 58, p. 10–298.
61. K. O. Ramanathan and S. H. Yen, "High-Temperature Emissivities of Copper, Aluminum, and Silver," *J. Opt. Soc. Am.* **67**:32–38 (1977).
62. E. A. Estalote and K. O. Ramanathan, "Low-Temperature Emissivities of Copper and Aluminum," *J. Opt. Soc. Am.* **67**:39–44 (1977).
63. K. O. Ramanathan, S. H. Yen, and E. A. Estalote, "Total Hemispherical Emissivities of Copper, Aluminum, and Silver," *Appl. Optics* **16**:2810–2817 (1977).

64. D. P. Verret and K. O. Ramanathan, "Total Hemispherical Emissivity of Tungsten," *J. Opt. Soc. Am.* **68**:1167–1172 (1978).
65. C. R. Roger, S. H. Yen, and K. O. Ramanathan, "Temperature Variation of Total Hemispherical Emissivity of Stainless Steel AISI 304," *J. Opt. Soc. Am.* **69**:1384–1390 (1979).
66. R. Smalley and A. J. Sievers, "The Total Hemispherical Emissivity of Copper," *J. Opt. Soc. Am.* **68**:1516–1518 (1978).
67. S. X. Cheng, P. Cebe, L. M. Hanssen, D. M. Riffe, and A. J. Sievers, "Hemispherical Emissivity of V, Nb, Ta, Mo, and W from 300 to 1000 K," *J. Opt. Soc. Am. B* **4**:351–356 (1987).
68. Op. cit. Ref. 5, p. 287.
69. Op. cit. Ref. 58, pp. 12–32–12–33.
70. Y. S. Touloukian et al. (eds.), *Thermophysical Properties of Matter*, **1–13**, IFI/Plenum, New York, 1970–1977.
71. Y. S. Touloukian, R. K. Kirby, R. E. Taylor, and P. D. Desai, "Thermal Expansion, Metallic Elements and Alloys," vol. 12 in Y. S. Touloukian and C. Y. Ho (eds.), *Thermophysical Properties of Matter*, IFI/Plenum, New York, 1975, pp. 23 (Be), 77 (Cu), 125 (Au), 157 (Fe), 208 (Mo), 225 (Ni), 298 (Ag), 1028 (Al), 1138 (SS).
72. Y. S. Touloukian, R. K. Kirby, R. E. Taylor, and T. Y. R. Lee, "Thermal Expansion, Nonmetallic Solids," vol. 13 in Y. S. Touloukian and C. Y. Ho (eds.), *Thermophysical Properties of Matter*, IFI/Plenum, New York, 1977, pp. 154 (Si), 873 (SiC).
73. C. A. Swenson, "HIP Beryllium: Thermal Expansivity from 4 to 300 K and Heat Capacity from 1 to 108 K," *J. Appl. Phys.* **70**(6):3046–3051 (Sep 1991).
74. Z. Li and C. Bradt, "Thermal Expansion of the Cubic (3C) Polytype of SiC," *J. Mater. Sci.* **21**:4366–4368 (1986).
75. C. Y. Ho, R. W. Powell, and P. E. Liley, *Thermal Conductivity of Selected Materials*, NSRDS-NBS-8, National Standard Reference Data System—National Bureau of Standards, Part 1 (1966).
76. C. Y. Ho, R. W. Powell, and P. E. Liley, *Thermal Conductivity of Selected Materials*, NSRDS-NBS-16, National Standard Reference Data System—National Bureau of Standards, part 2 (1968).
77. Y. S. Touloukian, R. W. Powell, C. Y. Ho, and P. G. Klemens, "Thermal Conductivity, Metallic Elements and Alloys," vol. 1 in Y. S. Touloukian and C. Y. Ho (eds.), *Thermophysical Properties of Matter*, IFI/Plenum, New York, 1970.
78. Y. S. Touloukian, R. W. Powell, C. Y. Ho, and P. G. Klemens, "Thermal Conductivity, Nonmetallic Solids," vol. 2 in Y. S. Touloukian and C. Y. Ho (eds.), *Thermophysical Properties of Matter*, IFI/Plenum, New York, 1970.
79. D. H. Killpatrick, private communication, Feb. 1993.
80. "CVD Silicon Carbide," Technical Bulletin #107, Morton International Advanced Materials, 1991.
81. Op. cit. Ref. 58, p. 12–133.
82. E. A. Brandes and G. B. Brook (eds.), *Smithell's Metals Reference Book*, 7th ed., Butterworth Heinmann, Oxford, 1992, pp. 14–3–14–5.
83. Y. S. Touloukian and E. H. Buyco, "Specific Heat, Metallic Elements and Alloys," vol. 4 in Y. S. Touloukian and C. Y. Ho (eds.), *Thermophysical Properties of Matter*, IFI/Plenum, New York, 1970.
84. *Metals Handbook*, **2**, 10th ed., Properties and Selection: Nonferrous Alloys and Special-Purpose Materials, ASM International, Metals Park, OH, 1990, pp. 93–94 and 102–103 (Al), 265 (Cu), 704–705 (Au), 1118–1129 (Fe), 1140–1143 (Mo), 441 (Ni), 1154–1156 (Si), and 699–700 & 1156–1158 (Ag).
85. I-70-H Optical Grade Beryllium Block, Preliminary Material Spec, Brush Wellman Inc., Nov. 1990.
86. *Carpenter Invar "36"*, Technical Data Sheet, Carpenter Technology Corp., Nov. 1980.
87. *Metals Handbook*, **5**, 9th ed., Surface Cleaning, Finishing, and Coating, American Society for Metals, Metals Park, OH, 1982, pp. 223–229.
88. *Engineered Materials Handbook*, **4**, Ceramics and Glasses, ASM International, Metals Park, OH, 1991, pp. 677, 806–808.
89. *Metals Handbook*, **1**, 10th ed., Properties and Selection: Irons, Steels, and High Performance Alloys, ASM International, Metals Park, OH, 1990, p. 871.
90. *Materials Engineering, Materials Selector 1993*, Dec. 1992, p. 104.
91. D. L. Hibbard, "Dimensional Stability of Electroless Nickel Coatings," in R. A. Paquin (ed.), *Dimensional Stability, Proc. SPIE* **1335**:180–185 (1990).

92. G. A. Graves, private communication, Feb. 1993.
93. Op. cit. Ref. 82, pp. 15-5–15-7.
94. Z. Li and R. C. Bradt, “Thermal Expansion and Elastic Anisotropies of SiC as Related to Polytype Structure,” in C. E. Selmer (ed.), *Proceedings of the Silicon Carbide 1987 Symposium 2*, Amer. Ceram. Soc, Westerville, OH, 1989, pp. 313–339.
95. W. D. Rowland and J. S. White, “The Determination of the Elastic Constants of Beryllium in the Temperature Range 25 to 300°C,” *J. Phys. F: Metal Phys.* 2:231–236 (1972).
96. H. Ledbetter, private communication, Oct. 1987.
97. Op. cit., Ref. 82, pp. 15-2–15-3.

DO NOT DUPLICATE

OPTICAL PROPERTIES OF SEMICONDUCTORS

David G. Seiler

*Semiconductor Electronics Division
National Institute of Standards and Technology
Gaithersburg, Maryland*

Stefan Zollner

*Freescale Semiconductor, Inc.
Hopewell Junction, New York*

Alain C. Diebold

*College of Nanoscale Science and Engineering
University at Albany
Albany, New York*

Paul M. Amirtharaj

*Sensors and Electron Devices Directorate
U.S. Army Research Laboratory
Adelphi, Maryland*

5.1 GLOSSARY

- A power absorption
- B magnetic field
- c velocity of light
- D displacement field
- d film thickness
- E applied electric field
- E_c energy, conduction band
- E_{ex} exciton binding energy
- E_g energy band gap
- E_H hydrogen atom ionization energy = 13.6 eV
- E electric field

E_n^\pm	Landau level energy
E_v	energy, valence band
e_i	ionic charge
g^*	effective g -factor
\mathbf{K}	phonon wave vector
k	extinction coefficient
k_B	Boltzmann's constant
\mathbf{k}	electron/hole wave vector
L_\pm	coupled LO phonon — plasmon frequency
m_e^*	electron effective mass
m_h^*	hole effective mass
m_i	ionic mass
m_i'	reduced ionic mass
m_{imp}	impurity ion mass
m_l^*	longitudinal effective mass
m_o	electron rest mass
m_r	electron-hole reduced mass
m_t^*	transverse effective mass
N	volume density
n	refractive index (real part)
$\tilde{n} = (n + ik)$	complex index of refraction
\mathbf{P}	polarization field
\mathbf{q}	photon wave vector
R	power reflection
R_y	effective Rydberg
S	oscillator strength
T	power transmission
T	temperature
V	Verdet coefficient
α	absorption coefficient
α_{AD}	absorption coefficient, allowed-direct transitions
α_{AI}	absorption coefficient, allowed-indirect transitions
δ	skin depth or penetration depth
γ	phenomenological damping parameter
Δ	spin-orbit splitting energy
Γ	Brillouin zone center
ϵ	dielectric function
$\epsilon_{\text{fc}}(\omega)$	free-carrier dielectric function
$\epsilon_{\text{imp}}(\omega)$	impurity dielectric function
$\epsilon_{\text{int}}(\omega)$	intrinsic dielectric function
$\epsilon_{\text{lat}}(\omega)$	lattice dielectric function
$\epsilon(0)$	static dielectric constant
ϵ_0	free-space permittivity
ϵ_1	Real(ϵ)
ϵ_2	Im(ϵ)

ϵ_{∞}	high-frequency limit of dielectric function
η	impurity ion charge
λ	wavelength
λ_c	cut-off wavelength
μ	mobility
μ_B	Bohr magneton
ν	frequency
σ	conductivity
τ	scattering time
ϕ	work function
χ	susceptibility
$\chi^{(n)}$	induced nonlinear susceptibility
Ω	phonon frequency
ω	angular frequency
ω_c	cyclotron resonance frequency
ω_{LO}	longitudinal optical phonon frequency
ω_p	free-carrier plasma frequency
ω_{pv}	valence band plasma frequency
ω_{TO}	transverse optical phonon frequency

5.2 INTRODUCTION

Rapid advances in semiconductor manufacturing and associated technologies have increased the need for optical characterization techniques for materials analysis and in situ monitoring/control applications. Optical measurements have many unique and attractive features for studying and characterizing semiconductor properties: (1) they are contactless, nondestructive, and compatible with any transparent ambient including high-vacuum environments; (2) they are capable of remote sensing, and hence are useful for in situ analysis on growth and processing systems; (3) the high lateral resolution inherent in optical systems may be harnessed to obtain spatial maps of important properties of the semiconductor wafers or devices; (4) combined with the submonolayer sensitivity of a technique such as ellipsometry, optical measurements lead to unsurpassed analytical details; (5) the resolution in time obtainable using short laser pulses allows ultrafast phenomena to be investigated; (6) the use of multichannel detection and high-speed computers can be harnessed for extremely rapid data acquisition and reduction which is crucial for real-time monitoring applications such as in situ sensing; (7) they provide information that complements transport analyses of impurity or defect and electrical behavior; (8) they possess the ability to provide long-range, crystal-like properties and hence support and complement chemical and elemental analyses; and (9) finally, most optical techniques are “table-top” procedures that can be implemented by semiconductor device manufacturers at a reasonable cost. All optical measurements of semiconductors rely on a fundamental understanding of their optical properties. In this chapter, a broad overview of the optical properties of semiconductors is given, along with numerous specific examples.

The optical properties of a semiconductor can be defined as any property that involves the interaction between electromagnetic radiation or light and the semiconductor, including absorption, diffraction, polarization, reflection, refraction, and scattering effects. The electromagnetic spectrum is an important vehicle for giving an overview of the types of measurements and physical processes characteristic of various regions of interest involving the optical properties of semiconductors. The electromagnetic spectrum accessible for studies by optical radiation is depicted in Fig. 1, where both the photon wavelengths and photon energies, as well as the common designations for

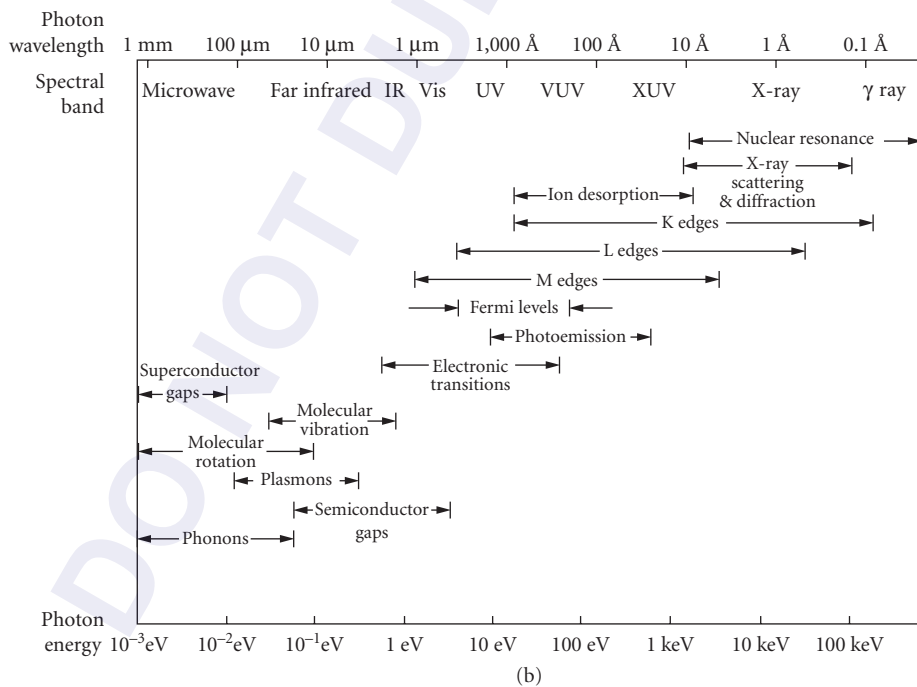
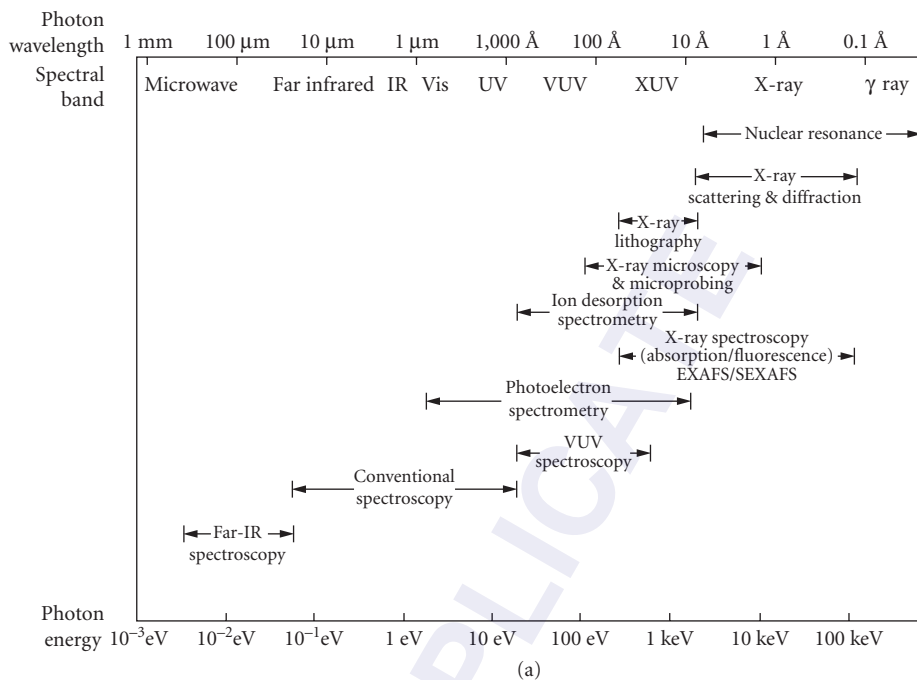


FIGURE 1 The electromagnetic spectrum comprising the optical and adjacent regions of interest: (a) characterization techniques using optical spectroscopy and synchrotron radiation and (b) molecular, atomic, and electronic processes characteristic in various parts of the electromagnetic spectrum plotted as a function of photon energy.¹

the spectral bands, are given.¹ Figure 1a shows the various techniques and spectroscopies and their spectral regions of applicability. Molecular, atomic, and electronic processes characteristic of various parts of the spectrum are shown in Fig. 1b. The high-energy x-ray, photoelectron, and ion desorption processes are important to show because they overlap the region of vacuum ultraviolet (VUV) spectroscopy. The ultraviolet (UV) region of the spectrum has often been divided into three rough regions: (1) the near-UV, between 2000 and 4500 Å; (2) the VUV, 2000 Å down to about 400 Å; and (3) the region below 400 Å covering the range of soft x-rays, 400 to 10 Å.² The spectrum thus covers a broad frequency range which is limited at the high-frequency end by the condition that $\lambda \gg a$, where λ is the wavelength of the light wave in the material and a is the interatomic distance. This limits the optical range to somewhere in the soft x-ray region. Technical difficulties become severe in the ultraviolet region (less than 100 nm wavelength, or greater than 12.3 eV photon energies), and synchrotron radiation produced by accelerators can be utilized effectively for ultraviolet and x-ray spectroscopy without the limitations of conventional laboratory sources. A lower limit of the optical frequency range might correspond to wavelengths of about 1 mm (photon energy of 1.23×10^{-3} eV). This effectively excludes the microwave and radio-frequency ranges from being discussed in a chapter on the optical properties of semiconductors.

From the macroscopic viewpoint, the interaction of matter with electromagnetic radiation is described by Maxwell's equations. The optical properties of matter are introduced into these equations as the constants characterizing the medium such as the dielectric constant, magnetic permeability, and electrical conductivity. (They are not actually "constants" since they vary with frequency.) From our optical viewpoint, we choose to describe the solid by the complex dielectric constant or complex dielectric function $\epsilon(\omega)$. This dielectric constant is a function of the space and time variables and should be considered as a response function or linear integral operator. The complex index of refractive index \tilde{n} is the complex square root of the dielectric function. Its real and imaginary parts are the refractive index n and the extinction coefficient k .

There are a number of methods for determining the optical constants of a semiconductor as a function of wavelength. Some of the most common techniques are as follows:

1. Measure the reflectivity at normal incidence *or* the transmission of a thin slab of known thickness over a wide wavelength range and use a Kramers-Kronig dispersion relation.
2. Measure the transmission of a thin slab of known thickness *and* the absolute reflectivity at normal incidence for a uniform bulk isotropic thin slab of known thickness with smooth surfaces (front and back).
3. Use a polarimetric method like ellipsometry which involves finding the ratio of reflectivities for polarizations perpendicular and parallel to the plane of incidence at a nonnormal incidence together with the difference of phase shifts upon reflection. With recent advancements in the availability of commercial ellipsometers, this is the most common method in use today.
4. Use detailed computer modeling and fitting of reflection, transmission, and/or ellipsometric measurements over a sufficiently large energy range. This method can also be applied in many circumstances to obtain the optical constants of anisotropic materials or thin films.
5. Calculate the band structure of a crystalline material using ab initio methods employing the local density approximation (LDA) with GW corrections (to correct the bandgap error usually found in LDA calculations). Integrate this band structure to determine the joint density of states weighted with the probability of optical transitions by the dipole operator matrix elements. If this technique is used, the Bethe-Salpeter equations need to be solved to take into account many-body excitonic effects (electron-hole interactions). Due to recent advances in computer technology, the accuracy of this method approaches the accuracy of measurements, especially if large samples with high quality cannot be obtained.

These optical constants describe an electromagnetic wave in the medium of propagation; the refractive index n gives the phase shift of the wave, and the extinction coefficient k gives the attenuation of the wave. In practice, one often uses the absorption coefficient α instead of k because of Beer's law formalism describing the absorption.

The field of optical spectroscopy is a very important area of science and technology since most of our knowledge about the structure of atoms, molecules, and solids is based upon spectroscopic investigations. For example, studies of the line spectra of atoms in the late 1800s and early 1900s revolutionized our understanding of the atomic structure by elucidating the nature of their electronic energy levels. Similarly for the case of semiconductors, optical spectroscopy has proven essential to acquiring a systematic and fundamental understanding of the nature of semiconductors. Since the early 1950s, detailed knowledge about the various eigenstates present in semiconductors has emerged including energy bands, excitonic levels, impurity and defect levels, densities of states, energy-level widths (lifetimes), symmetries, and changes in these conditions with temperature, pressure, magnetic field, electric field, etc. One of the purposes of this chapter is to review and summarize the major optical measurement techniques that have been used to investigate the optical properties of semiconductors related to these features. Specific attention is paid to the types of information that can be extracted from such measurements of the optical properties.

Most optical properties of semiconductors are integrally related to the particular nature of their electronic and vibrational structures.³ These electronic and vibrational dispersion relations are in turn related to the type of crystallographic structure, the particular atoms, and their bonding. The full symmetry of the space groups is also essential in determining the structure of the energy bands and vibrational modes. Group theory makes it possible to classify energy eigenstates, determine essential degeneracies, derive selection rules, and reduce the order of the secular determinants which must be diagonalized in order to compute approximate eigenvalues. Often, experimental measurements must be carried out to provide quantitative numbers for these eigenvalues. A full understanding of the optical properties of semiconductors is thus deeply rooted in the foundations of modern solid-state physics. In writing this chapter, the authors have assumed that the readers are familiar with some aspects of solid-state physics such as can be obtained from an advanced undergraduate course.

Most semiconductors have a diamond, zincblende, wurtzite, or rock-salt crystal structure. Elements and binary compounds, which average four valence electrons per atom, preferentially form tetrahedral bonds. A tetrahedral lattice site in a compound AB is one in which each atom A is surrounded symmetrically by four nearest neighboring B atoms. The most important lattices with a tetrahedral arrangement are the diamond, zincblende, and wurtzite lattices. In the diamond structure, all atoms are identical, whereas the zincblende structure contains two different atoms. The wurtzite structure is in the hexagonal crystal class, whereas the diamond and zincblende structures are cubic. Other lattices exist which are distorted forms of these, and many lattices have no relation to the tetrahedral structures.

Band structure calculations show that only the valence band states are important for predicting the following crystal ground-state properties: charge density, Compton profile, compressibility, cohesive energy, lattice parameter, x-ray emission spectra, and hole effective masses. In contrast, both the valence band and conduction band states are important for predicting the following properties: optical dielectric constant or refractive index, optical absorption spectrum, and electron effective masses. Further complexities arise because of the many-body nature of the particle interactions which necessitates understanding excitons, electron-hole droplets, polarons, polaritons, etc.

The optical properties of semiconductors cover a wide range of phenomena which are impossible to do justice to in just one short chapter in this *Handbook*. We have thus chosen to present an extensive, systematic overview of the field, with as many details given as possible. The definitions of the various optical properties, the figures chosen, the tables presented, and the references given all help to orient the reader to appreciate various principles and measurements that form the foundations of the optical properties of semiconductors.

The optical properties of semiconductors are often subdivided into those that are electronic and those that are vibrational (lattice related) in nature. The electronic properties concern processes involving the electronic states of the semiconductor, while the lattice properties involve vibrations of the lattice (absorption and creation of phonons). Lattice properties are of considerable interest for heat dissipation, electronic transport, and lifetimes (broadenings) of electronic states, but it is the electronic properties which receive the most attention in semiconductors

because of the technological importance of their practical applications. Modern-day semiconductor optoelectronic technologies include lasers, light-emitting diodes, photodetectors, optical amplifiers, modulators, switches, etc., all of which exploit specific aspects of the electronic optical properties.

Almost all of the transitions that contribute to the optical properties of semiconductors can be described as one-electron transitions. Most of these transitions conserve the crystal momentum and thus measure the vertical energy differences between the conduction and valence bands. In the one-electron approximation, each valence electron is considered as a single particle, moving in a potential which is the sum of the core potentials and a self-consistent Hartree potential of the other valence electrons.

The phenomena usually studied to obtain information on the optical properties of semiconductors are (1) absorption, (2) reflection or ellipsometry, (3) photoconductivity, (4) emission, (5) light scattering, and (6) modulation techniques. Most of the early information on the optical properties of semiconductors was obtained from measurements of photoconductivity, but these measurements can be complicated by carrier trapping, making interpretation of the results sometimes difficult. Thus, most quantitative measurements are of the type (1), (2), (4), or (5). For example, the most direct way of obtaining information about the energy gaps between band extrema and about impurity levels is by measuring the optical absorption over a wide range of wavelengths. This information can also be obtained by (2) and (4).

The transient nature of the optical properties of semiconductors is important to establish because it gives insight to the various relaxation processes that occur after optical excitation. Because of the basic limitations of semiconductor devices on speed and operational capacity, ultrafast studies have become an extremely important research topic to pursue. The push to extend the technologies in the optoelectronic and telecommunication fields has also led to an explosion in the development and rise of ultrafast laser pulses to probe many of the optical properties of semiconductors: electrons, holes, optical phonons, acoustic phonons, plasmons, magnons, excitons, and the various coupled modes (polaritons, polarons, excitonic molecules, etc.). The time scale for many of these excitations is measured in femtoseconds (10^{-15}) or picoseconds (10^{-12}). Direct time measurements on ultrafast time scales provide basic information on the mechanisms, interactions, and dynamics related to the various optical properties. Some of the processes that have been investigated are the carrier lifetime, the formation time of excitons, the cooling and thermalization rates of hot carriers, the lifetime of phonons, the screening of optical-phonon-carrier interactions, the dynamics of ballistic transport, the mechanism of laser annealing, dephasing processes of electrons and excitons, optical Stark effect, etc. It is not possible in this short review chapter to cover these ultrafast optical properties of semiconductors. We refer the reader to the many fine review articles and books devoted to this field.^{4,5}

The advent of the growth of artificially structured materials by epitaxial methods such as molecular beam epitaxy (MBE) has made possible the development of a new class of materials and heterojunctions with unique electronic and optical properties. Most prominent among these are heterojunction quantum wells and superlattices. The field of microstructural physics has thus been one of the most active areas of research in the past decades. The novel properties of nanostructures fabricated from ultrathin layers of semiconductors of thicknesses $<100 \text{ \AA}$ stem from microscopic quantum mechanical effects. The simplest case to visualize is that of a particle confined in a box which displays distinct quantum energy states, the equivalent of which are electrons and holes confined to a thin layer of a material such as GaAs sandwiched between two thick layers of AlAs. The new energy states produced by the confinement of the charges in the artificially produced potential well can be manipulated, by tailoring the size and shape of the well, to produce a wide variety of effects that are not present in conventional semiconductors. Microstructures formed from alternating thin layers of two semiconductors are called superlattices. They also lead to novel electronic and optical behaviors, most notable of which are large anisotropic properties. The ability to "engineer" the behavior of these nanostructures has led to an explosion of research and applications that is too large to be dealt with in this short review. The reader is referred to several review articles on their optical behavior.^{6,7}

5.3 OPTICAL PROPERTIES

Background

The interaction of the semiconductor with electromagnetic radiation can be described, in the semi-classical regime, using response functions such as ϵ and χ which are defined in the following section. The task of the description is then reduced to that of building a suitable model of χ and ϵ that takes into account the knowledge of the physical characteristics of the semiconductor and the experimentally observed optical behavior. One example of a particularly simple and elegant, yet surprisingly accurate and successful, model of ϵ for most semiconductors is the linear-chain description of lattice vibrations.⁸ This model treats the optical phonons (i.e., the vibrations that have an associated dipole moment) as damped simple harmonic motions. Even though the crystal is made up of $\approx 10^{23}$ atoms, such a description with only a few resonant frequencies and phenomenological terms, such as the damping and the ionic charge, accurately accounts for the optical behavior in the far-infrared region. The details of the model are discussed in the following section “Optical/Dielectric Response.” Such simple models are very useful and illuminating, but they are applicable only in a limited number of cases, and hence such a description is incomplete.

A complete and accurate description will require a self-consistent quantum mechanical approach that accounts for the microscopic details of the interaction of the incident photon with the specimen and a summation over all possible interactions subject to relevant thermodynamical and statistical mechanical constraints. For example, the absorption of light near the fundamental gap can be described by the process of photon absorption resulting in the excitation of a valence band electron to the conduction band. In order to obtain the total absorption at a given energy, a summation has to be performed over all the possible states that can participate, such as from multiple valence bands. Many-body effects (the Coulomb attraction between the photoexcited electron and hole) also need to be considered for an accurate description. Thermodynamic considerations such as the population of the initial and final states have to be taken into consideration in the calculation as well. Hence, a detailed knowledge of the specimen and the photon-specimen interaction can, in principle, lead to a satisfactory description.

Optical/Dielectric Response

Optical Constants and the Dielectric Function In the linear regime for an isotropic solid, the dielectric function ϵ and the susceptibility χ are defined by the following relations:⁹

$$\mathbf{D} = \epsilon_0 \mathbf{E} + \mathbf{P} \quad (1)$$

$$\mathbf{D} = \epsilon_0 (1 + \chi) \mathbf{E} \quad (2)$$

$$\mathbf{D} = \epsilon \mathbf{E} = (\epsilon_1 + i\epsilon_2) \mathbf{E} \quad (3)$$

where \mathbf{E} , \mathbf{D} , and \mathbf{P} are the free-space electric field, the displacement field, and the polarization field inside the semiconductor; ϵ_0 is the permittivity of free space; and ϵ and χ are dimensionless quantities, each of which can completely describe the optical properties of semiconductors. The refractive index \tilde{n} of the material is related to ϵ as shown below:

$$\tilde{n} = \sqrt{\epsilon} = n + ik \quad (4)$$

The real and imaginary parts of the refractive index, n and k , which are also referred to as the optical constants, embody the linear optical property of the material. The presence of k , the imaginary component, denotes absorption of optical energy by the semiconductor. Its relationship to the absorption Coefficient is discussed in the following section. "Reflection, Transmission, and Absorption coefficients." In the spectral regions where absorptive processes are weak or absent, as in the case of the sub-bandgap range, k is very small, whereas in regions of strong absorption, the magnitude of k is large. The optical constants for a large number of semiconductors may be found in Refs. 10 and 11. The variation in the real part n is usually much smaller. For example, in GaAs, at room temperature, in the visible and near-visible region extending from 1.4 to 6 eV, k varies from $<10^{-3}$ at 1.41 eV which is just below the gap, to a maximum of 4.1 at 4.94 eV.¹² In comparison, n remains nearly constant in the near-gap region extending from 3.61 at 1.4 eV to 3.8 at 1.9 eV, with the maximum and minimum values of 1.26 at 6 eV and 5.1 at 2.88 eV, respectively. The real and imaginary components are related by causal relationships that are also discussed in the following sections.

Reflection, Transmission, and Absorption Coefficients The reflection and transmission from a surface at normal incidence are given by

$$\tilde{r} = \frac{(\tilde{n}-1)}{(\tilde{n}+1)} = |\tilde{r}| \cdot \exp(i\theta) \quad (5)$$

$$R = |\tilde{r}|^2 \quad (6)$$

$$T = (1 - R) \quad (7)$$

where \tilde{r} is the complex (Fresnel) reflection coefficient (consisting of a magnitude and a phase θ) and R and T are the reflectance and transmission, describing the ratio of the reflected and transmitted intensity relative to the incident intensity. For a thin slab, in free space, with thickness d and complex refractive index \tilde{n} , the appropriate expressions are:¹³

$$\tilde{r} = \frac{\tilde{r}_1 + \tilde{r}_2 \cdot \exp(i4\pi\tilde{n}d/\lambda)}{1 + \tilde{r}_1 \cdot \tilde{r}_2 \cdot \exp(i4\pi\tilde{n}d/\lambda)} \quad (8)$$

where \tilde{r}_1 and \tilde{r}_2 are the Fresnel reflection coefficients at the first and second interfaces, respectively, and λ is the free-space wavelength.

For most cases of optical absorption, the energy absorbed is proportional to the thickness of the specimen. The variation of intensity inside the absorptive medium is given by the following relationship:

$$I(x) = I(0) \cdot \exp(-\alpha \cdot x) \quad (9)$$

and the absorption coefficient α is related to the optical constants by

$$\alpha = 4\pi k/\lambda \quad (10)$$

Here we note that α (measured in cm^{-1}) describes the attenuation of the radiation intensity rather than that of the electric field.

In spectral regions of intense absorption, all the energy that enters the medium is absorbed. The only part of the incident energy that remains is that which is reflected at the surface. In such a case, it is useful to define a characteristic “skin” thickness that is subject to an appreciable density of optical energy. A convenient form used widely is simply the inverse of α ; that is, $1/\alpha$. This skin depth (or penetration length) is usually denoted by δ :

$$\delta = \frac{1}{\alpha} \quad (11)$$

The skin depths in semiconductors range from >100 nm near the bandgap to <5 nm at the higher energies of ≈ 6 eV.

Kramers-Kronig Relationships A general relationship exists for linear systems between the real and imaginary parts of a response function as shown in the following:

$$\varepsilon_1(\omega) = 1 + \frac{2}{\pi} P \int_0^{\infty} \frac{\omega' \varepsilon_2(\omega')}{\omega'^2 - \omega^2} d\omega' \quad (12)$$

$$\varepsilon_2(\omega) = -\frac{2\omega}{\pi} P \int_0^{\infty} \frac{\varepsilon_1(\omega') d\omega'}{\omega'^2 - \omega^2} + \frac{\sigma_0}{\varepsilon_0 \cdot \omega} \quad (13)$$

where σ_0 is the dc conductivity.

$$n(\omega) = 1 + \frac{2}{\pi} P \int_0^{\infty} \frac{\omega' k(\omega')}{\omega'^2 - \omega^2} d\omega' \quad (14)$$

$$k(\omega) = -\frac{2}{\pi} P \int_0^{\infty} \frac{n(\omega')}{\omega'^2 - \omega^2} d\omega' \quad (15)$$

where P denotes the principal part of the integral. These are referred to as the Kramers-Kronig dispersion relationships.^{14,15} An expression of practical utility is one in which the experimentally measured reflectance R at normal incidence is explicitly displayed as shown:

$$\theta(\omega) = -\frac{\omega}{\pi} P \int_0^{\infty} \frac{\ln(R(\omega'))}{\omega'^2 - \omega^2} d\omega' \quad (16)$$

This is useful since it shows that if R is known for all frequencies, the phase θ of the Fresnel reflection coefficient can be deduced, and hence a complete determination of both n and k can be accomplished. In practice, R can be measured only over a limited energy range, but approximate extrapolations can be made to establish reasonable values of n and k .

The measurement of the reflectivity over a large energy range spanning the infrared to the vacuum ultraviolet, 0.5 to 12 eV range, followed by a Kramers-Kronig analysis, used to be the main method of establishing n and k .¹⁶ However, the advances in spectroscopic ellipsometry in the past 20 years have made this obsolete in all but the highest energy region. A Kramers-Kronig analysis of reflectance data also has experimental difficulties related to surface roughness and native oxides, which are more easily corrected with ellipsometry than with reflectivity measurements. The Kramers-Kronig relations are still important today, since they place consistency conditions on optical constants determined from ellipsometry.

Above 12 eV, the higher energy reflectance is dominated by the valence band plasma edge ω_{pv} and, hence, takes the following forms for $\varepsilon(\omega)$, $\tilde{n}(\omega)$, and $R(\omega)$:

$$\tilde{n} = \sqrt{\varepsilon}(\omega) \approx -\frac{1}{2} \cdot \left(\frac{\omega_{PV}^2}{\omega^2} \right) \quad (17)$$

$$R(\omega) = \frac{(n(\omega)-1)}{(n(\omega)+1)} = \frac{1}{16} \cdot \left(\frac{\omega_{PV}^4}{\omega^4} \right) \quad (18)$$

Sum Rules Having realized the interrelationships between the real and imaginary parts of the response functions, one may extend them further using a knowledge of the physical properties of the semiconductor to arrive at specific equations, commonly referred to as sum rules.^{14,15} These equations are useful in cross-checking calculations and measurements for internal consistency or reducing the computational effort. Some of the often-used relations are shown below:

$$\int_0^{\infty} \omega \varepsilon_2(\omega) d\omega = \frac{\pi}{2} \omega_{PV}^2 \quad (19)$$

$$\int_0^{\infty} \omega \operatorname{Im} \left[\frac{-1}{\varepsilon(\omega)} \right] d\omega = \frac{\pi}{2} \omega_{PV}^2 \quad (20)$$

$$\int_0^{\infty} \omega k(\omega) \cdot d\omega = \frac{\pi}{4} \omega_{PV}^2 \quad (21)$$

$$\int_0^{\infty} [n(\omega)-1] d\omega = 0 \quad (22)$$

where ω_{PV} is the valence band plasma frequency.

The dc static dielectric constant, $\varepsilon(0)$, may be expressed as

$$\varepsilon(0) = 1 + \frac{2}{\pi} \int_0^{\infty} \frac{\varepsilon_2(\omega)}{\omega} \cdot d\omega \quad (23)$$

The reader is referred to Refs. 14 and 15 for more details.

Linear Optical Properties

Overview The optical properties of semiconductors at low enough light levels are often referred to as linear properties in contrast to the nonlinear optical properties described later. There are many physical processes that control the amount of absorption or other optical properties of a semiconductor. In turn, these processes depend upon the wavelength of radiation, the specific properties of the individual semiconductor being studied, and other external parameters such as pressure, temperature, etc. Just as the electrical properties of a semiconductor can be controlled by purposely introducing impurity dopants (both p and n type) or affected by unwanted impurities or defects, so too are the optical properties affected by them. Thus, one can talk about *intrinsic* optical properties of semiconductors that depend upon their perfect crystalline nature and *extrinsic* properties that are introduced by impurities or defects. Many types of defects exist in real solids: point defects, macroscopic structural defects, etc. In this section we review and summarize intrinsic linear optical properties related to lattice effects, interband transitions, and free-carrier or intraband transitions. Impurity- and defect-related extrinsic optical properties are also covered in a separate section and in the discussion of lattice properties affected by them. Figure 2 schematically depicts

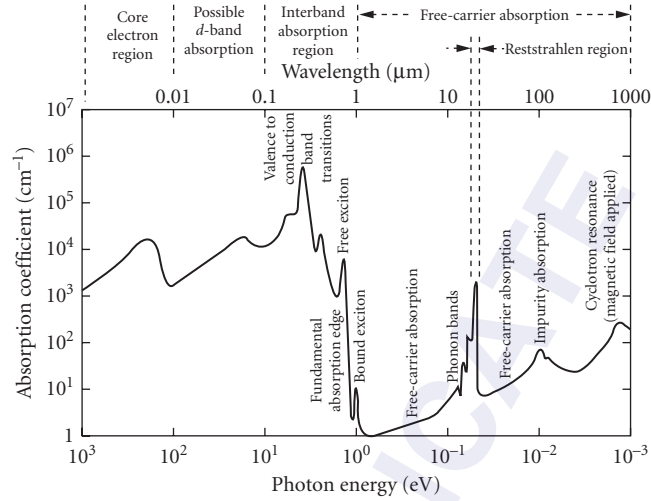


FIGURE 2 Absorption spectrum of a typical semiconductor showing a wide variety of optical processes.

various contributions to the absorption spectrum of a typical semiconductor as functions of wavelength (top axis) and photon energy (bottom axis). Data for a real semiconductor may show more structure than shown here. On the other hand, some of the structure shown may be reduced or not actually present in a particular semiconductor (e.g., impurity absorption, bound excitons, d -band absorption). Table 1 shows the classification of the optical responses of the semiconductor to light in various wavelength regions showing the typical origin of the response and how the measurements are usually carried out. At the longest wavelengths shown in Fig. 2, cyclotron resonance (CR) may occur for a semiconductor in a magnetic field, giving rise to an absorption peak corresponding to a transition of a few meV energy between Landau levels. Shallow impurities may give rise to additional absorption at low temperatures, and here a 10 meV ionization energy has been assumed. If the temperature was high enough so that $k_B T$ was greater than the ionization energy, the absorption peak would be washed out. At wavelengths between 5 and 50 μm , a new set of absorption peaks arises due to the vibrational modes of the lattice. In ionic crystals, the absorption coefficient in the reststrahlen region may reach 10^5 cm^{-1} , whereas in homopolar semiconductors like Si and Ge, only multiphonon features with lower absorption coefficients are present (around 5 to 50 cm^{-1}).

Models of the dielectric function The interaction of light with semiconductors can be completely described by the dielectric function, $\epsilon(\omega)$. The dielectric function $\epsilon(\omega)$ may be divided into independent parts to describe various physical mechanisms so long as the processes do not interact strongly with each other; this is an approximation, referred to as the adiabatic approximation which simplifies the task at hand considerably.¹⁷ The major players that determine the optical behavior of an intrinsic semiconductor are the lattice, particularly in a nonelemental semiconductor; the free carriers (i.e., mobile electrons and holes) and the interband transitions between the energy states available to the electrons. These three mechanisms account for the intrinsic linear properties that lead to a dielectric function as shown:

$$\epsilon_{\text{int}}(\omega) = \epsilon_{\text{lat}}(\omega) + \epsilon_{\text{fc}}(\omega) + \epsilon_{\text{inter}}(\omega) \quad (24)$$

The addition of impurities and dopants that are critical to controlling the electronic properties leads to an additional contribution, and the total dielectric response may then be described as shown:

$$\epsilon(\omega) = \epsilon_{\text{int}}(\omega) + \epsilon_{\text{imp}}(\omega) \quad (25)$$

TABLE 1 Classification by Wavelength of the Optical Responses for Common Semiconductors

Wavelength (nm)	Responses	Physical Origin	Application	Measurement Tech.
$\lambda > \lambda_{TO}$ Far-IR and micro wave region	Microwave <i>R</i> and <i>T</i> Plasma <i>R</i> and <i>T</i>	Free-carrier plasma	Detectors Switches	<i>R</i> , <i>T</i> , and <i>A</i> * Microwave techniques Fourier Transform Spectrometry (FTS), FT-SE
$\lambda_{LO} < \lambda < \lambda_{TO}$ Reststrahlen region	Reststrahlen <i>R</i>	Optical phonons in ionic crystals	Absorbers Filters	<i>R</i> , <i>T</i> , and <i>A</i> FTS & Dispersion spectrometry (DS), FT-SE
$\lambda \sim \lambda_{LO}, \lambda_{TO}, \lambda_P$ Far-IR region	Far-IR <i>A</i>	Optical phonons, impurities (vibrational and electronic), free carriers, intervalence transitions	Absorbers Filters	<i>R</i> , <i>T</i> , and <i>A</i> FTS, DS, and FT-SE
$\lambda_{LO} > \lambda > \lambda_G$ Mid-IR region	Mid-IR <i>T</i> and <i>A</i>	Multiphonon, multiphoton transitions, impurities (vibrational and electronic), intervalence transitions excitons, Urbach tail	Detectors Switches Absorbers Filters	<i>R</i> , <i>T</i> , and <i>A</i> Ellipsometry FTS and DS
$\lambda < \lambda_G$ IR, visible, and UV	<i>R</i> , <i>T</i> , and <i>A</i>	Electronic interband transitions	Reflectors Detectors	Reflection Ellipsometry
$\lambda \sim \lambda_W$ UV, far-UV	Photoemission	Fermi energy to vacuum-level electronic transitions	Photocathodes Detectors	High-vacuum spectroscopy techniques, UV ellipsometry
$\lambda_W > \lambda \geq a$	<i>R</i> , <i>T</i> , and <i>A</i>	Ionic-core transitions	Detectors	Soft x-ray and synchrotron- based analyses
	Diffraction	Photo—ionic-core interactions	X-ray optics and monochromators	X-ray techniques

**R*, *T*, and *A*—Reflection, transmission, and absorption. SE—Spectroscopic ellipsometry.
Note: P—Plasma; G—Energy gap; W—Work functions; *a*—lattice constant.
TO, LO—Transverse and longitudinal optical phonons.

Lattice Absorption

Phonons The dc static response of a semiconductor lattice devoid of free charges to an external electromagnetic field may be described by the single real quantity $\epsilon(0)$. As the frequency of the electromagnetic radiation increases and approaches the characteristic vibrational frequencies associated with the lattice, strong interactions can occur and modify the dielectric function substantially. The main mechanism of the interaction is the coupling between the electromagnetic field with the oscillating dipoles associated with vibrations of an ionic lattice.⁸ The interactions may be described, quite successfully, by treating the solid to be a collection of damped harmonic oscillators with a characteristic vibrational frequency ω_{TO} and damping constant γ . The resultant dielectric function may be written in the widely used CGS units as:

$$\epsilon_{\text{lat}}(\omega) = \epsilon(\infty) + \frac{S\omega_{TO}^2}{(\omega_{TO}^2 - \omega^2 - i\omega\gamma)} \quad (26)$$

where *S* is called the oscillator strength and may be related to the phenomenological ionic charge e_i , reduced mass m'_i , and volume density *N*, through the equation

$$S\omega_{TO}^2 = \frac{4\pi N e_i^2}{m'_i} \quad (27)$$

In the high frequency limit of $\epsilon(\omega)$, for $\omega \gg \omega_{\text{TO}}$, (but well below the onset of interband absorption)

$$\epsilon(\omega) \rightarrow \epsilon_{\infty} \quad (28)$$

The relationship may be easily extended to accommodate more than one characteristic vibrational frequency by summing over all phonon modes i :

$$\epsilon_{\text{lat}}(\omega) = \epsilon_{\infty} + \sum_i \frac{S_i (\omega_{\text{TO}}^i)^2}{[(\omega_{\text{TO}}^i)^2 - \omega^2 - i\omega\gamma^i]} \quad (29)$$

It is worth noting some important physical implications and interrelations of the various parameters in Eq. (26).

For a lattice with no damping, it is obvious that $\epsilon(\omega)$ displays a pole at ω_{TO} and a zero at a well-defined frequency, usually referred to by ω_{LO} . A simple but elegant and useful relationship exists between these parameters as shown by

$$\frac{\epsilon(0)}{\epsilon_{\infty}} = \left(\frac{\omega_{\text{LO}}}{\omega_{\text{TO}}} \right)^2 \quad (30)$$

which is known as the Lyddane-Sachs-Teller relation.¹⁸

The physical significance of ω_{TO} and ω_{LO} is that these are the transverse and longitudinal optical phonon frequencies with zero wave vector, \mathbf{K} , supported by the crystal lattice. The optical vibrations are similar to standing waves on a string. The wave pattern, combined with the ionic charge distribution, leads to oscillating dipoles that can interact with the incident radiation and, hence, the name optical phonons. $\epsilon(\omega)$ is negative for $\omega_{\text{TO}} \geq \omega \geq \omega_{\text{LO}}$, which implies no light propagation inside the crystal and, hence, total reflection of the incident light. The band of frequencies spanned by ω_{TO} and ω_{LO} is referred to as the reststrahlen band.

The reflectivity spectrum of AlSb¹⁹ is shown in Fig. 3. It is representative of the behavior of most semiconductors. Note that the reflectivity is greater than 90 percent at $\approx 31 \mu\text{m}$ in the reststrahlen band spanned by the longitudinal and transverse optical phonons. The two asymptotic limits of the reflection tend to $[(\sqrt{\epsilon(0)}-1)/(\sqrt{\epsilon(0)}+1)]^2$ and $[(\sqrt{\epsilon_{\infty}}-1)/(\sqrt{\epsilon_{\infty}}+1)]^2$, respectively, compare Eqs. (5) and (6). The effects of the phonon damping are illustrated in Fig. 4.²⁰ For the ideal case with zero damping, the reflection in the reststrahlen band is 100 percent. Since absolute reflectance measurements are very difficult to carry out experimentally, the reststrahlen region is nowadays often investigated with infrared ellipsometry. Note that for the elemental semiconductors such as Si and Ge, the lack of a dipole moment associated with the optical vibrations of the lattice leads to the absence of any oscillator strength and hence no interaction with the radiation. The reflection spectrum will, therefore, show no change at or near the optical phonon frequencies. The optical phonon frequencies (ω_{LO} and ω_{TO}) and wavelengths and $\epsilon(0)$ and ϵ_{∞} for the commonly known semiconductors are presented at the end of this chapter in Table 11. For elemental semiconductors, $\epsilon(0) = \epsilon_{\infty}$.

The optical phonons ω_{LO} and ω_{TO} are the frequencies of interest for describing the optical interactions with the lattice. In addition, the lattice is capable of supporting vibrational modes over a wide range of frequencies extending from 0 to $\omega_{\text{LO}}(\Gamma)$, the LO phonon frequency at the center of the Brillouin zone, Γ , as discussed by Cochran in Ref. 21. The vibrational modes can be subdivided into two major categories. The first are optical phonons that possess an oscillating dipole moment and therefore can interact with light. The second group is the acoustic phonons; that is, sound-like vibrations that do not possess a dipole moment and hence are not of primary importance in

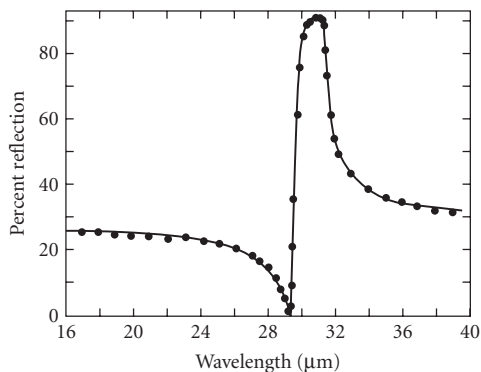


FIGURE 3 Reststrahlen reflection spectrum of AlSb.¹⁹

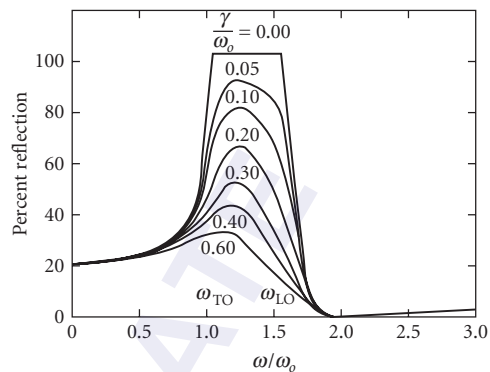


FIGURE 4 Reflection spectra of a damped oscillator for various values of the damping factor.²⁰

determining the optical properties. For the optical band, only the $\mathbf{K} \approx 0$ modes are important since a strong interaction is precluded for other modes due to the large mismatch in the wave vectors associated with the light and vibrational disturbances.

The simple treatment so far, though very useful, is limited to the most obvious and strongest aspect of the interaction of light with the lattice. However, many weaker but important interaction mechanisms have not yet been accounted for. An attempt to produce a complete description starts with the consideration of the total number of atoms N that makes up the crystal. Each atom possesses three degrees of freedom, and hence one obtains $3N$ degrees of freedom for the crystal, a very large number of magnitude, $\approx 10^{24}$. The complexity of the description can be readily reduced when one realizes the severe restriction imposed by the internal symmetry of the crystal. The vibrational characteristics now break up into easily understandable normal modes with well-defined physical characteristics. The translational symmetry associated with the crystal makes it possible to assign a definite wave vector to each lattice mode. In addition, the phonons can be divided into two major classes: the optical vibrations, which we have already discussed, and the acoustic vibrations, which, as the name implies, are sound-like vibrations. The acoustic phonons for $|\mathbf{K}| \rightarrow 0$ are identical to the sound waves. Hence, specifying the type of phonon, energy, wave vector, and polarization uniquely describes each vibrational mode of the crystal. The conventional description of phonons is achieved by graphically displaying these properties in a frequency versus \mathbf{K} plot, as shown in Fig. 5 for GaAs.²²

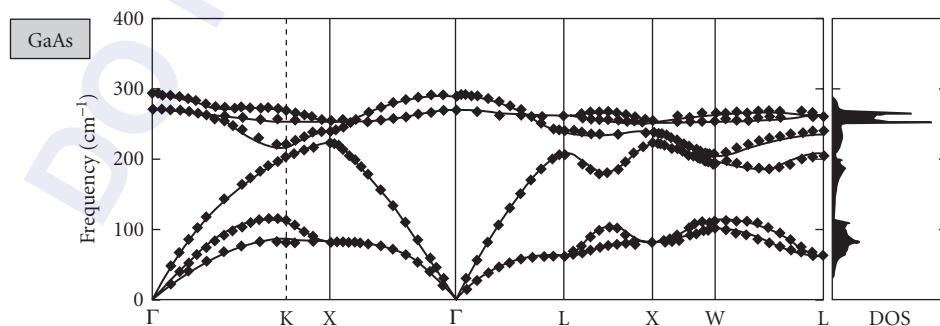


FIGURE 5 Phonon dispersion curves of GaAs. Experimental data are represented by symbols. Solid lines show the results of the calculation of Gianozzi et al.²²

The energies of the phonons are plotted as a function of the wave vector along the high-symmetry directions for each acoustic and optical branch. It can be shown that the number of acoustic branches is always three, two of transverse polarization and one longitudinal, leaving $3p-3$ optical branches, where p is the number of atoms in the primitive cell. The majority of the most important semiconductors fall into the three cubic classes of crystals—namely, the diamond, zincblende, and the rock-salt structures that contain two atoms per primitive cell and hence possess three acoustic and three optical branches each.²³ However, more complicated structures with additional optical phonon branches can be found. The most important group among the second category is the wurtzite structure displayed by CdS, CdSe, GaN, InN, AlN, etc., which contains four atoms in the primitive cell and hence two additional sets of optical phonon branches.²³ The specific symmetry associated with the vibrational characteristics of each mode is used to distinguish them as well as their energies at the high-symmetry points in the phonon dispersion curves.

Multiphonon absorption It has already been pointed out that the interaction of light with phonons is restricted to those with $\mathbf{K} \approx 0$. This is true only when single-phonon interactions are considered in ideal crystals. Higher-order processes, such as multiphonon absorption, can activate phonons with $\mathbf{K} \neq 0$.^{23,24} Symmetry considerations and their implication on the multiphonon absorption are discussed by Birman in Ref. 25. In multiphonon processes, the total momentum of the interacting phonons will be 0, but many modes with $\mathbf{K} \neq 0$ can participate in the interaction.

The energy and momentum conservation conditions may be expressed as follows:

$$\hbar\omega = \sum_i \hbar\Omega_i \quad (31)$$

$$\hbar\mathbf{q} = \sum_i \hbar\mathbf{K}_i \approx 0 \quad (32)$$

where $\hbar\omega$ is the energy of the absorbed photon, $\hbar\Omega_i$ is the energy of the phonons, and $\hbar\mathbf{q}$ and $\hbar\mathbf{K}_i$ are the corresponding momenta. Any number of phonons can participate in the process. However, the strength of interaction between the incident photon and the higher-order processes falls off rapidly with increasing order, making only the two- or three-phonon processes noteworthy in most semiconductors. The well-defined range that spans the phonon energies in most semiconductors, extending from 0 to the LO phonon energy at the center of the zone Γ , restricts the n -phonon process to a maximum energy of $n\hbar\omega_{\text{LO}}(\Gamma)$. Among the participating phonons, those with large values of \mathbf{K} and those in the vicinity of the critical points are the most important owing to their larger populations. These factors are important in understanding the multiphonon absorption behavior, as we now discuss.^{23,24}

Multiphonon processes may be subdivided into two major categories: (1) sum processes where multiple phonons are created, and (2) difference processes in which both phonon creation and annihilation occur with a net absorption in energy. The former process is more probable at higher-incident photon energies and the converse is true for the latter. The reduction of the equilibrium phonon population at low temperatures leads to a lower probability for the difference process, and hence it is highly temperature dependent.

The multiphonon interactions are governed by symmetry selection rules in addition to the energy and momentum conservation laws stated earlier. A list of the possible combination processes is presented in Table 2 for the diamond structure and Table 3 for the zincblende structure. Representative multiphonon absorption spectra are presented in Figs. 6 and 7 for Si²⁶ and GaAs.²⁷ The Si ω_{TO} and ω_{LO} frequency of $\sim 522 \text{ cm}^{-1}$ is indicated in the spectrum for reference. Note that the absence of a dipole moment implies that ω_{TO} and ω_{LO} are degenerate, and no first-order absorption is present. In contrast, for GaAs the very large absorption associated with the one-phonon absorption precludes the possibility of obtaining meaningful multiphonon absorption data in the reststrahlen band that spans the ~ 269 to 295 cm^{-1} (~ 34 – 37 meV) spectral region. The multiphonon spectra obtained from GaAs are displayed in Fig. 7a and b. The absorption associated with the multiphonon processes is much smaller than the single-phonon process. However, the rich structure displayed by the spectra is extremely useful in analyzing the lattice dynamics of the material. In

TABLE 2 Infrared Allowed Processes in the Diamond Structure²⁴

Two-Phonon Processes
TO(X) + L(X)
TO(X) + TA(X)
L(X) + TA(X)
TO(L) + LO(L)
TO(L) + TA(L)
LO(L) + LA(L)
LA(L) + TA(L)
TO(W) + L(W)
TO(W) + TA(W)
L(W) + TA(W)

addition, in applications such as windows for high-power lasers, even the small absorption levels can lead to paths of catastrophic failure.

Impurity-related vibrational optical effects The role of impurities is of primary importance in the control of the electrical characteristics of semiconductors and hence their technological applications. This section outlines the main vibrational features of impurities and the resulting modification of the optical properties of these semiconductors. Impurities can either lead to additional vibrational modes over and beyond that supported by the unperturbed lattice or they can activate normally inactive vibrational modes.²⁸ The perturbation of the lattice by a substitutional impurity is a change in the mass of one of the constituents and a modification of the bonding forces in its vicinity. If the impurity is much lighter than the host atom it replaces, high-frequency vibrational modes above $\omega_{LO}(\Gamma)$, the maximum frequency supported by the unperturbed lattice, are introduced. These vibrational amplitudes are localized in the vicinity of the impurity and hence are known as local vibrational modes (LVM). For heavier impurities, the impurity-related vibrations can occur within the phonon band or in the gap between the acoustic and optical bands. These modes are referred to as resonant modes (RM) or gap modes (GM).²⁸

The qualitative features of an impurity vibrational mode can be understood by considering a simple case of a substitutional impurity atom in a linear chain. The results of a numerical calculation of a 48-atom chain of GaP are presented in Figs. 8 and 9.²⁸ The highly localized character can be seen from Fig. 9. Note that the degree of localization reduces with increasing defect mass for a fixed bonding strength.

TABLE 3 Infrared Allowed Processes in the Zincblende Structure²⁴

Two-Phonon Processes
2LO(Γ), LO(Γ) + TO(Γ), 2TO(Γ)
2TO(X), TO(X) + LO(X), TO(X) + LA(X), TO(X) + TA(X)
LO(X) + LA(X), LO(X) + TA(X)
LA(X) + TA(X)
2TA(X)
2TO(L), TO(L) + LO(L), TO(L) + LA(L), TO(L) + TA(L)
2LO(L), LO(L) + LA(L), LO(L) + TA(L)
2LA(L), LA(L) + TA(L)
2TA(L)
TO ₁ (W) + LO(W), TO ₁ (W) + LA(W)
TO ₂ (W) + LO(W), TO ₂ (W) + LA(W)
LO(W) + LA(W), LO(W) + TA ₁ (W), LO(W) + TA ₂ (W)
LA(W) + TA ₁ (W), LA(W) + TA ₂ (W)

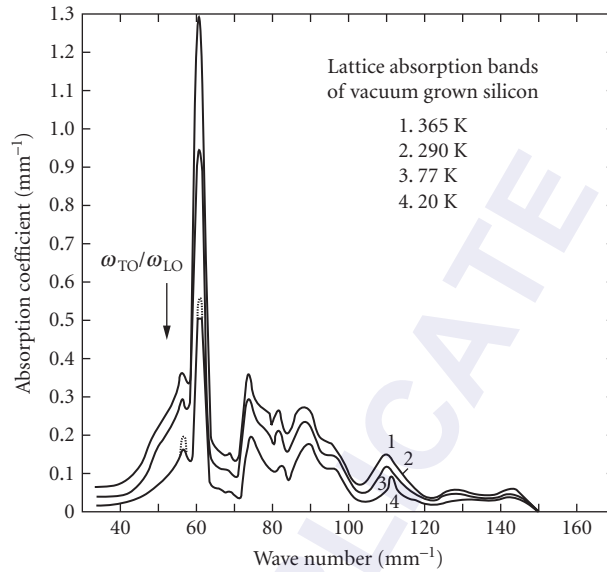


FIGURE 6 Multiphonon absorption of vacuum-grown Si.²⁶

The absorption band produced by the LVM has been successfully used in impurity analyses. Figure 10 shows the IR absorption spectrum associated with interstitial oxygen in Si taken at the National Institute of Standards and Technology (NIST), and Fig. 11 displays the absorption spectrum from a carbon-related LVM in GaAs.^{29,30} Note that the multiple peaks in the high-resolution spectrum are a consequence of the mass perturbations to the local environment resulting from the

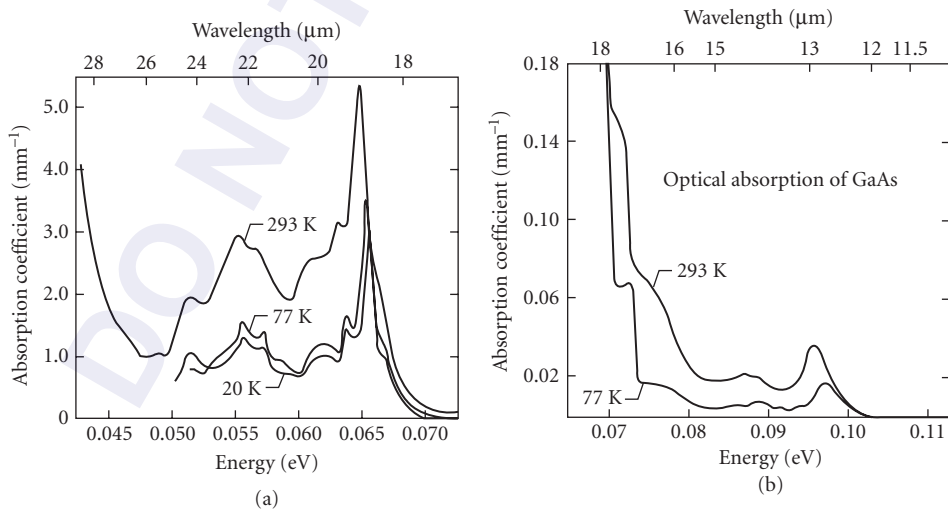


FIGURE 7 Lattice absorption coefficient of high-resistivity *n*-type GaAs vs. wavelength from (a) 18 to 28 μm and (b) 10 to 18 μm at 77 and 293 K.²⁷

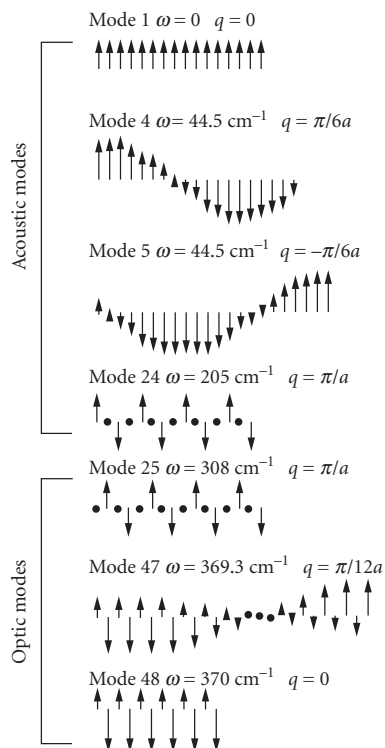


FIGURE 8 Linear-chain model calculations for GaP. A 48-atom chain is considered. Position along the chain is plotted horizontally and ion displacement vertically. Modes 24 and 25 occur on either side of the gap between the acoustic and optic bands.²⁸

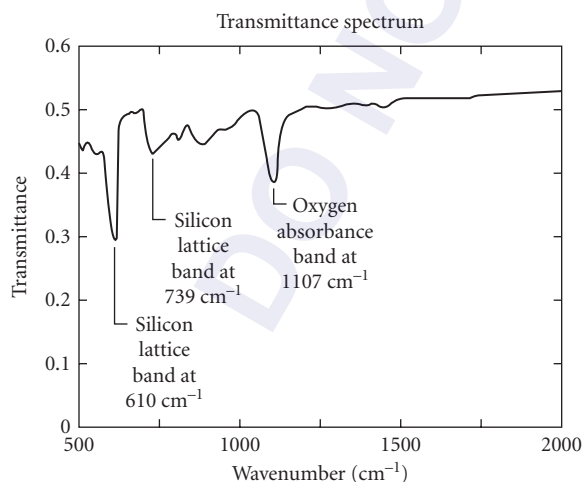


FIGURE 10 IR absorption due to interstitial oxygen in Si.

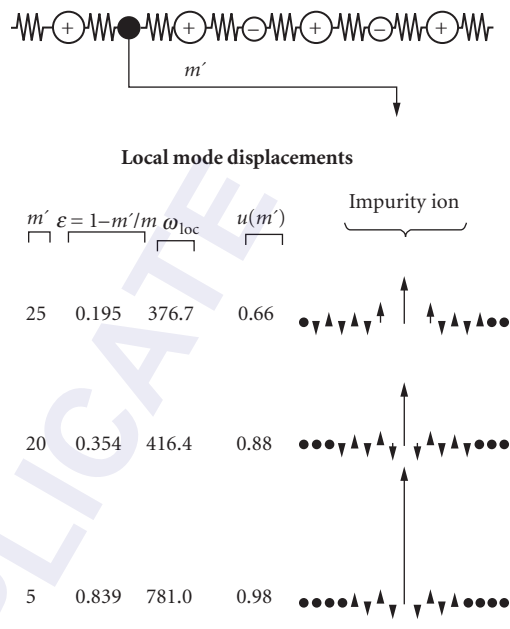


FIGURE 9 Eigenvectors for the highest-frequency (localized) mode for three isotopic substitutions on the $m = 31$ site. Note the extreme localization for a substituent of mass 5.²⁸

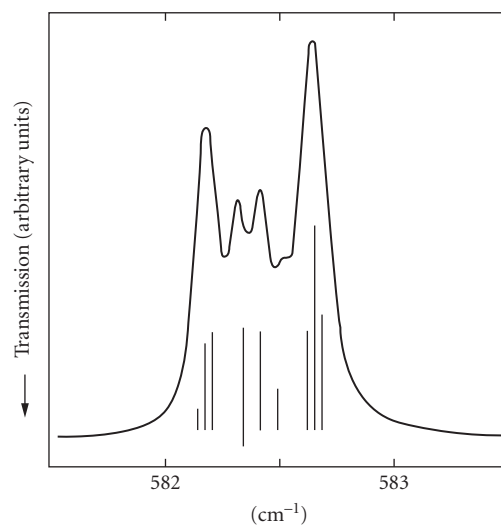


FIGURE 11 ^{12}C local modes in GaAs and the predicted fine structure. The height of each line is proportional to the strength of each mode.³⁰

TABLE 4 Localized Modes in Semiconductors²⁸

Host and Impurity	Mode Frequency (temp. K)	Defect Symmetry; Method of Observation
Diamond N	1340(300)	T_d, A^*
Silicon		
¹⁰ B	644(300), 646(80)	$T_d; A$
¹¹ B	620(300), 622(80)	$T_d; A$
As	366(80) Reson.	$T_d; A$
P	441 Reson, 491(80) Reson.	$T_d; A$
¹⁴ C	570(300), 573(80)	$T_d; A$
¹³ C	586(300), 589(80)	$T_d; A$
¹² C	605(300), 680(80)	$T_d; A$
O	Bands near 30, 500, 1100, 1200	Complex; A, R, T
GaAs		
Al	362(80), \approx 371(4.2)	$T_d; A, R, T$
P	355(80), 353(300), \approx 363(4.2)	$T_d; A, R, T$
Si _{Ga}	384(80)	$T_d; A$
Si _{As}	399(80)	$T_d; A$

*A—absorption; T—transmission; R—reflection; Reson.—resonant mode.

two naturally occurring isotopes of Ga. The fine structure in the spectrum helps identify the site occupied by C as belonging to the As sublattice. When accurate calibration curves are available, the concentration of the impurity can be determined from the integrated intensity of the LVM band. The LVM frequencies of a number of hosts and impurities are presented in Table 4.²⁸ The symmetry of the point group associated with the defect is T_d . The use of LVM in the study of complex defects, particularly those that involve hydrogen, has led to a wealth of microscopic information not easily attainable by any other means.³¹

The quantitative accuracy obtainable from LVM analysis may be illustrated by a simple harmonic model calculation. In such an approximation, the total integrated absorption over the entire band may be expressed as:³²

$$\int \alpha d\omega = \frac{2\pi^2 N \eta^2}{n m_{\text{imp}} c} \quad (33)$$

where N is the volume density of the impurity and η and m_{imp} are the apparent charge and mass of the impurity ion; c is the velocity of light. n is the refractive index of the host crystal. η is an empirically derived parameter that is specific to each center, that is, a specific impurity at a specific lattice location. Once calibration curves are established, measurement of the intensity of absorption can be used to determine N . Figure 12 displays a calibration used to establish the density of interstitial oxygen.³³ Such analyses are routinely used in various segments of the electronic industry for materials characterization.

The effect of the impurities can alter the optical behavior in an indirect fashion as well. The presence of the impurity destroys translational symmetry in its vicinity and hence can lead to relaxation of the wave-vector conservation condition presented earlier in Eq. (32). Hence, the entire acoustic and optical band of phonons can be activated, leading to absorption bands that extend from zero frequency to the maximum $\omega_{\text{LO}}(\Gamma)$. The spectral distribution of the absorption will depend on the phonon density-of-states modulated by the effect of the induced dipole moment.²⁸ The latter is a consequence of the perturbation of the charge distribution by the impurity. The perturbation due to defects may be viewed in a qualitatively similar fashion. For instance, a vacancy may be described as an impurity with zero mass.

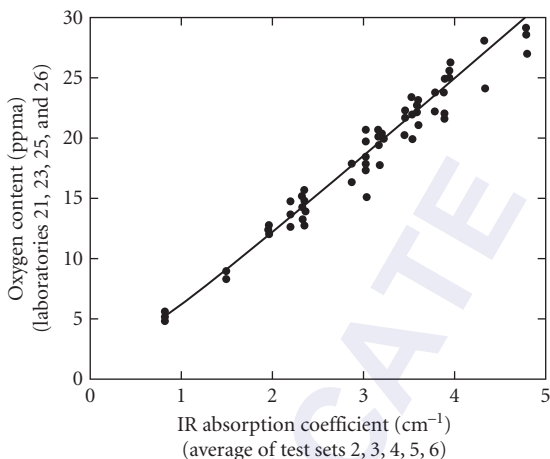


FIGURE 12 The absorption coefficient dependence on the concentration of interstitial oxygen in Si.³³

Interband Absorption

Absorption near the fundamental edge The fundamental absorption edge is one of the most striking features of the absorption spectrum of a semiconductor. Within a small fraction of an electron volt at an energy about equal to the energy gap E_g of the material, the semiconductor changes from being practically transparent to completely opaque—the absorption coefficient changing by a factor of 10^4 or more. This increased absorption is caused by transitions of electrons from the valence band to the conduction band. This characteristic optical property is clearly illustrated in Fig. 13, which shows the transmission versus wavelength for a number of major semiconductors.³⁴ At the lower wavelengths, the transmission approaches zero which defines a cutoff wavelength λ_c for each material. For example, $\lambda_c \approx 7.1 \mu\text{m}$ for InSb; $\lambda_c \approx 4.2 \mu\text{m}$, PbTe; $\lambda_c \approx 3.5 \mu\text{m}$, InAs; $\lambda_c \approx 1.8 \mu\text{m}$, GaSb and Ge; $\lambda_c \approx 1 \mu\text{m}$, Si; and $\lambda_c \approx 0.7 \mu\text{m}$ for CdS. At much longer wavelengths than the edge at λ_c , lattice and free-carrier absorption become appreciable and the transmission drops. Studies of the fundamental absorption edge thus give values for the energy gap and information about the states just above the edge in the conduction band and below it in the valence band. Properties of these states are important to know since they are responsible for electrical conduction. Details of the band structure near the band extrema can be determined from the position and shape of the absorption edge and from its temperature, magnetic field, pressure, impurity concentration, and other parameters' dependence. Finally, this fundamental gap

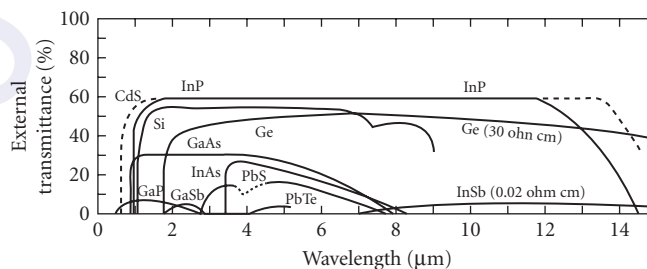


FIGURE 13 The transmission of CdS, InP, Si, Ge, GaAs, GaP, GaSb, InAs, InSb, PbTe, and PbS.³⁴

region is important because usually it is only near the energy gap that phenomena such as excitons (both free and bound), electron-hole drops, donor-acceptor pairs, etc., are seen.

Interband transitions near the fundamental absorption edge are classified as (1) direct or vertical or (2) indirect or nonvertical. The momentum of light ($\hbar q = \hbar n \omega / c$) is negligible compared to the momentum of a \mathbf{k} -vector state at the edge of the Brillouin zone. Thus, because of momentum conservation, electrons with a given wave vector in a band can only make transitions to states in a higher band having essentially the same wave vector. Such transitions are called vertical transitions. A nonvertical transition can take place, but only with the assistance of phonons or other entities which help preserve momentum.

Direct transitions The interband absorption coefficient depends upon the band structure and photon energy $\hbar\omega$. Use of quantum mechanics and, in particular, time-dependent perturbation theory becomes necessary.³⁵ For a nonzero momentum (dipole) matrix element, a simple model gives for allowed direct transitions

$$\alpha_{AD} = C_{AD} (\hbar\omega - E_g)^{1/2} \quad (34)$$

The coefficient C_{AD} involves constants, valence and conduction band effective masses, matrix elements, and only a slight dependence on photon energy. The absorption strengths of direct-gap semiconductors are related to their density of states and the momentum matrix element that couples the bands of interest. Many semiconductors such as AlAs, AlP, GaAs, InSb, CdS, ZnTe, and others have allowed direct transitions; many complex oxides, such as Cu_2O , SiO_2 , and rutile, have forbidden direct absorption.

Figure 14 shows the spectral variation of the absorption coefficient for pure InSb at a temperature of 5 K compared to various theoretical predictions.³⁶ We note the extremely sharp absorption edge which is fit best by the $(\hbar\omega - E_g)^{1/2}$ dependence near the edge. However, a big deviation from the experimental data occurs at higher photon energies. Consideration of two more details allows a better

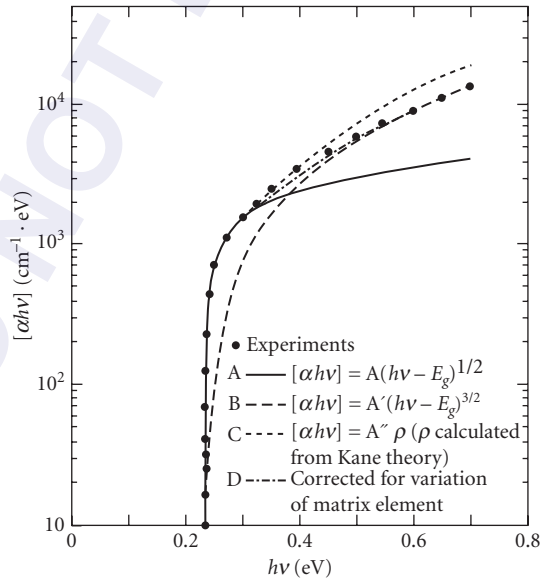


FIGURE 14 Theoretical fit to the experimental absorption edge of InSb at ≈ 5 K.³⁶

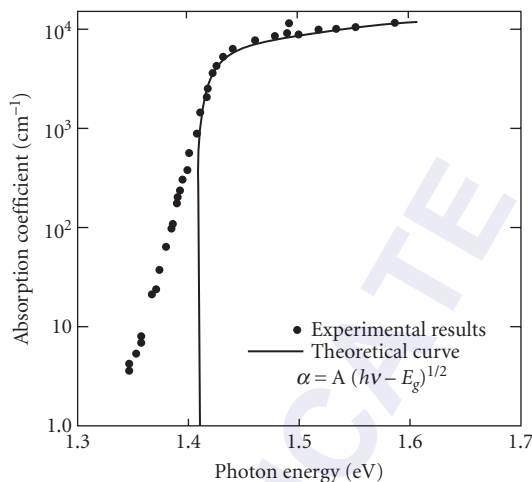


FIGURE 15 Absorption edge of GaAs at room temperature.³⁸

fit: (1) use of a more complicated band model from Kane³⁷ which predicts a more rapidly increasing density of states than for the simple bands, and (2) taking into account a decrease in the optical matrix element at the higher photon energies because of the \mathbf{k} -dependence of the wave functions. The calculated curves in Fig. 14 were arbitrarily shifted so that they look like a better fit than they are. The actual calculated absorption is a factor of about 15 too low at high energies. This discrepancy was attributed to the neglect of exciton effects which can greatly affect the absorption as discussed later. Modern calculations (including the Bethe-Salpeter corrections for many-body effects to the local density approximation) have much better accuracy in describing the optical absorption due to excitons.

Figure 15 shows the absorption behavior of GaAs at room temperature compared with calculations based on Kane's theory.³⁸ Below about 10^3 cm^{-1} , the absorption decreases much more slowly than predicted and absorption is even present for energies below E_g . In practice, there seems to exist an exponentially increasing absorption edge rule (called Urbach's rule) in most direct transition materials which is found to correlate reasonably well with transitions involving band tails. These band tails seem to be related to disorder, such as doping effects, phonon-assisted transitions, and the broadening of electronic states due to electron-phonon interactions.

Indirect transitions Semiconductors such as GaP, Ge, and Si have indirect gaps where the maximum valence band energy and minimum conduction band energy do not occur at the same \mathbf{k} value. In this case, the electron cannot make a direct transition from the top of the valence band to the bottom of the conduction band because this would violate conservation of momentum. Such a transition can still take place but as a two-step process requiring the cooperation of another particle and which can then be described by second-order perturbation theory. The particle most frequently involved is an intervalley phonon of energy $\hbar\Omega_{\mathbf{q}}$ which can be either generated or absorbed in the transition. (In some cases, elastic scattering processes due to impurity atoms or dislocations must be considered; they are less frequent than the phonon interactions in high-purity crystals.) The photon supplies the needed energy, while the phonon supplies the required momentum. The transition probability depends not only on the density of states and the electron-phonon matrix elements as in the direct case, but also on the electron-phonon interaction which is temperature dependent.

Calculations of the indirect-gap absorption coefficient give for the allowed indirect transitions

$$\alpha_{\text{AI}} = C_{\text{AI}}^{(\text{abs})}(\hbar\omega + \hbar\Omega_{\mathbf{q}} - E_g)^2 + C_{\text{AI}}^{(\text{em})}(\hbar\omega - \hbar\Omega_{\mathbf{q}} - E_g)^2 \quad (35)$$

where the superscripts (abs) and (em) refer to phonon absorption and emission, respectively. This expression is only nonzero when the quantities in parentheses are positive, that is, when $\hbar\omega \pm \hbar\Omega q > E_g$. We note that the phonon energies are usually small (≤ 0.05 eV) compared to the photon energy of about 1 eV and thus for the case of allowed indirect transitions with phonon absorption

$$\alpha_{AI} \approx C_{AI}^{(\text{abs})} (\hbar\omega - E_g)^2 \quad (36)$$

Thus the absorption increases as the second power of $(\hbar\omega - E_g)$, much faster than the half-power dependence of the direct transition as seen in Eq. (34).

Figure 16 shows the variation of the absorption coefficient of GaP with photon energy at room temperature near the indirect edge.³⁹ A reasonable fit to the experimental data of Spitzer et al.⁴⁰ is obtained indicating that, for GaP, allowed indirect transitions dominate. Further complications arise because there can be more than one type of phonon emitted or absorbed in the absorption process. Transverse acoustic (TA), longitudinal acoustic (LA), transverse optic (TO), and longitudinal optic (LO) phonons can be involved as shown in the absorption edge data of GaP, as seen in Fig. 17.⁴¹ The phonon energies deduced from these types of experimental absorption edge studies agree with those found from neutron scattering.

Both indirect and direct absorption edge data for Ge are shown in Fig. 18, while the analysis of the 300 K data is plotted in part *b*.⁴² At the lowest energies, α rises due to the onset of indirect absorption as seen by the $\alpha^{1/2}$ dependence on photon energy. At higher energies, a sharper rise is found where direct transitions occur at the zone center and an α^2 dependence on energy is then seen. Note the large shifts of E_g with temperature for both the direct and indirect gaps. Also, the direct-gap absorption is much stronger than that of the indirect-gap absorption.

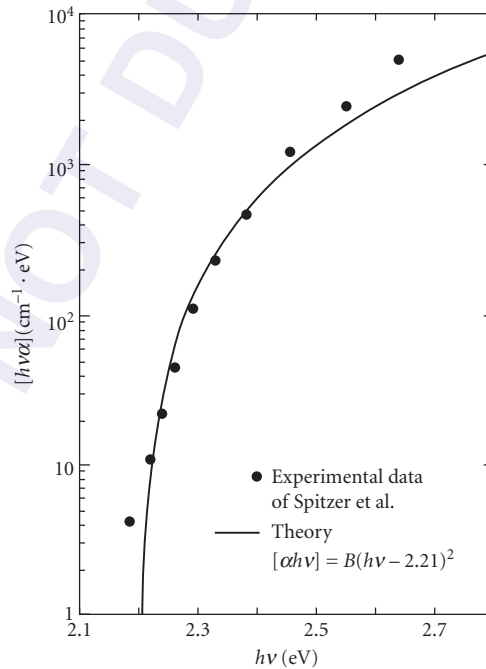


FIGURE 16 Comparison of the experimental data at room temperature for the absorption edge of GaP with the theory for an indirect edge.³⁹

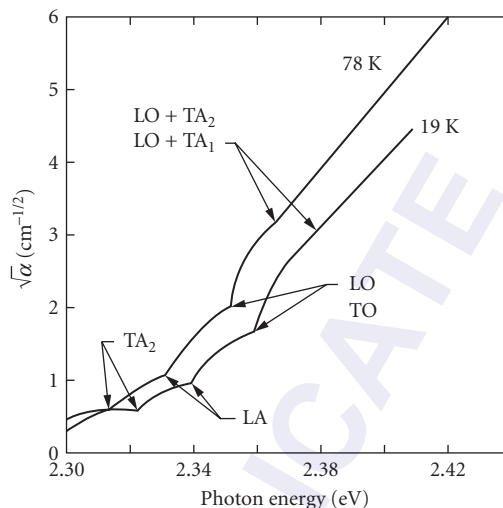


FIGURE 17 Absorption spectra at the edge of GaP, showing thresholds associated with the emission of each of several different phonons.⁴¹

Excitons Among the various optical properties of semiconductors, the subject of excitons has one of the dominant places because of their remarkable and diverse properties. Studies of exciton properties represent one of the most important aspects of scientific research among various solid-state properties. According to Cho,⁴³ there are a number of reasons for this: (1) excitonic phenomena are quite common to all the nonmetallic solids—semiconductors, ionic crystals, rare gas crystals, molecular crystals, etc.; (2) the optical spectra often consist of sharp structure, which allows a detailed theoretical analysis; (3) theories are not so simple as to be understood by a simple application of atomic theory or the Bloch band scheme, but still can be represented by a quasi-hydrogen-like level scheme; (4) sample quality and experimental techniques have continually been improved

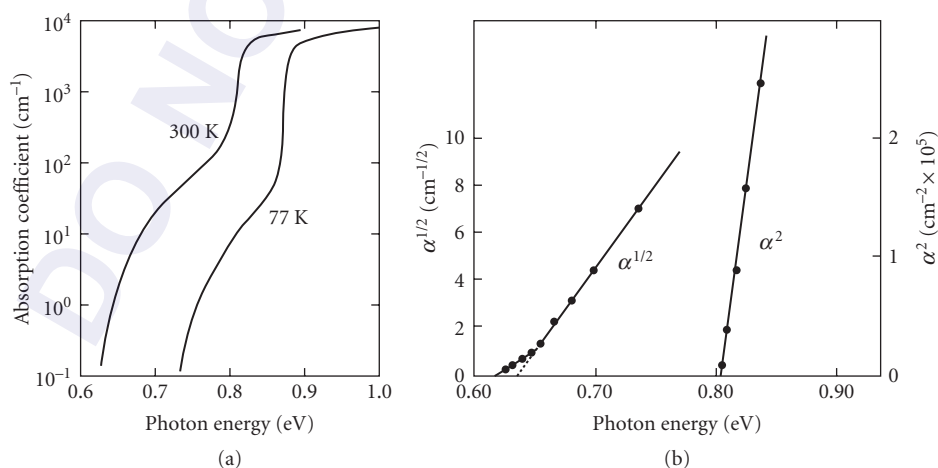


FIGURE 18 (a) α vs. $\hbar\omega$ for Ge and (b) the analysis of the 300 K experimental data.⁴²

TABLE 5 A Glossary of the Main Species of Excitons⁴⁴

Exciton	In essence, an electron and hole moving with a correlated motion as an electron-hole pair
Wannier exciton	Electron and hole both move in extended orbits; energy levels related to hydrogen-atom levels by scaling, using effective masses and dielectric constant; occurs in covalent solids such as silicon
Frenkel exciton	Electron and hole both move in compact orbits, usually essentially localized on adjacent ions; seen in ionic solids, such as KCl, in absorption
Self-trapped exciton	One or both carriers localized by the lattice distortion they cause; observed in ionic solids, such as KCl, in emission
Bound exciton	Only a useful idea when a defect merely prevents translational motion of an exciton and does not otherwise cause significant perturbation
Core exciton	Lowest-energy electronic excitation from a core state, leaving an unoccupied core orbital (e.g., the 1 s level of a heavy atom) and an electron in the conduction band whose motion is correlated with that of the core hole
Excitonic molecule	Complex involving two holes and two electrons
Multiple bound excitons	Complex of many holes and a similar number of electrons, apparently localized near impurities; some controversy exists, but up to six pairs of localized carriers have been suggested
Exciton gas	High concentration of electrons and holes in which each electron remains strongly associated with one of the holes (as insulating phase)
Electron-hole drops	High concentration of electrons and holes in which the motions are plasma-like (a metallic phase), not strongly correlated as in excitons

with subsequent experiments proving existing theories and giving rise to new ones; and (5) the exciton is an elementary excitation of nonmetallic solids, a quantum of electronic polarization. It has a two-particle (electron and hole) nature having many degrees of freedom, and, along with the variety of energy-band structures, this leads to a lot of different properties from material to material or from experiment to experiment. Table 5 gives a definition of the major types of excitons in a glossary obtained from Hayes and Stoneham.⁴⁴ Many examples exist in the literature involving work on excitons in semiconductors to understand their nature and to determine their properties. Besides the references cited in this chapter, the authors refer the reader to the more detailed work presented in Refs. 45 and 46.

An electron, excited from the valence band to a higher energy state in the conduction band, can still be bound by the Coulomb attraction to the hole that the electron leaves in the valence band. This neutral bound-electron hole pair is called an exciton which can move throughout the crystal. Excitons are most easily observed at energies just below E_g using optical absorption or photoluminescence measurements. There are two models used for describing excitons in solids, named after Frenkel and Wannier. In a solid consisting of weakly interacting atoms, Frenkel considered excitons as described by excitations of a single atom or molecule.⁴⁷ An excited electron describes an orbit of atomic dimensions around an atom with a vacant valence state. The empty valence state acts as a mobile hole since the excitation can move from one atom to another. These tightly bound excitons are similar to an ordinary excited state of the atom, except that the excitation can propagate through the solid. The radius of the Frenkel exciton is on the order of the lattice constant. Frenkel excitons are useful to describe optical properties of solids like alkali halides and organic phosphors.

Wannier (or also called Mott-Wannier) excitons are also electrons and holes bound by Coulomb attraction.^{48,49} In contrast to the Frenkel exciton, the electron and hole are separated by many lattice spacings producing a weakly bound exciton which is remarkably similar to a hydrogen-atom-like system. Since the electron and hole are, on the average, several unit cells apart, their Coulomb interaction is screened by the average macroscopic dielectric constant, ϵ_∞ , and electron and hole effective masses can be used. Their potential energy $-e^2/\epsilon_\infty r$ is just that of the hydrogen atom (except for ϵ_∞). The binding energy of the free exciton (relative to a free electron and free hole) is then given by the hydrogen-atom-like discrete energy levels plus a kinetic energy term due to the motion of the exciton:

$$E_{\text{ex}} = \frac{R_y}{n^2} - \frac{\hbar^2 k^2}{2(m_e^* + m_h^*)} \quad (37)$$

$$R_y = \left(\frac{m_r}{\epsilon_\infty^2} \right) \left(\frac{e^4}{2\hbar^2} \right) = \left(\frac{m_r}{m_H \epsilon_\infty^2} \right) E_H \quad (38)$$

where

$$m_r = \left(\frac{1}{m_e^*} + \frac{1}{m_h^*} \right)^{-1} \quad (39)$$

is the reduced mass of the exciton, m_H is the reduced mass of the hydrogen atom, n is the principal quantum number (1, 2, . . . , ∞), E_H is the ionization energy of the hydrogen atom (13.6 eV), and R_y is the effective Rydberg energy. The lowest energy absorption transition of the semiconductor is thus $E_g - E_{\text{ex}}$. As an example, consider CdS where $m_e^* \approx 0.21 m_0$, $m_h^* \approx 0.64 m_0$, and $\epsilon_\infty \approx 8.9$; here $m_r \approx 0.158 m_0$ and $R \approx 27$ meV. The Bohr radius for the $n = 1$ ground state is about 30 Å.

A series of excitonic energy levels thus exists just below the conduction band whose values increase parabolically with \mathbf{k} and whose separation is controlled by n . Excitons are unstable with respect to radiative recombination whereby an electron recombines with a hole in the valence band, with the emission of a photon or phonons. These excitonic levels are shown in Figs. 19 and 20. For

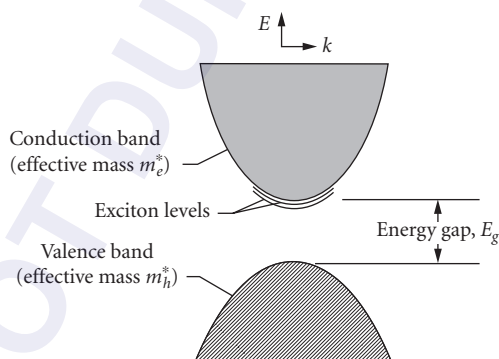


FIGURE 19 Exciton levels in relation to the conduction-band edge, for a simple band structure with both conduction and valence band edges at near $k = 0$.

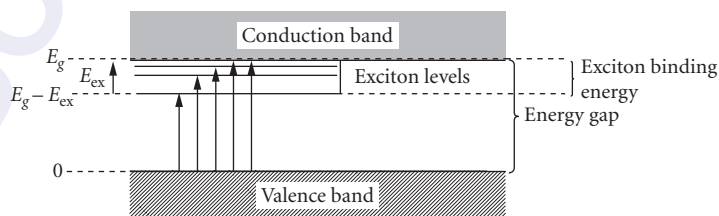


FIGURE 20 Energy levels of a free exciton created in a direct process. Optical transitions from the top of the valence band are shown by the arrows; the longest arrow corresponds to the energy gap.

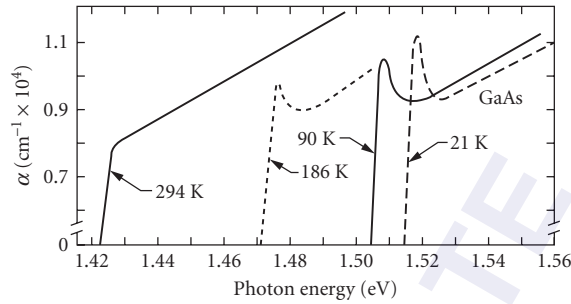


FIGURE 21 Observed exciton absorption spectra in GaAs at various temperatures between 21 and 294 K. Note the decrease in the band edge with increasing temperatures.⁵⁰

many semiconductors only a single peak is observed, as shown in Fig. 21 for GaAs.⁵⁰ However, even though only one line is observed, the exciton states make a sizable contribution to the magnitude of the absorption near and above the edge. At room temperatures, the exciton peak can be completely missing since the binding energy is readily supplied by phonons. In semiconductors with large enough carrier concentrations, no excitons exist because free carriers tend to shield the electron-hole interaction. Neutral impurities can also cause a broadening of the exciton lines and, at large enough concentrations, cause their disappearance.

Extremely sharp exciton states can often be seen as shown in Fig. 22, which shows the absorption spectrum of a very thin, very pure epitaxial crystal of GaAs.⁵¹ The $n = 1, 2,$ and 3 excitons are clearly

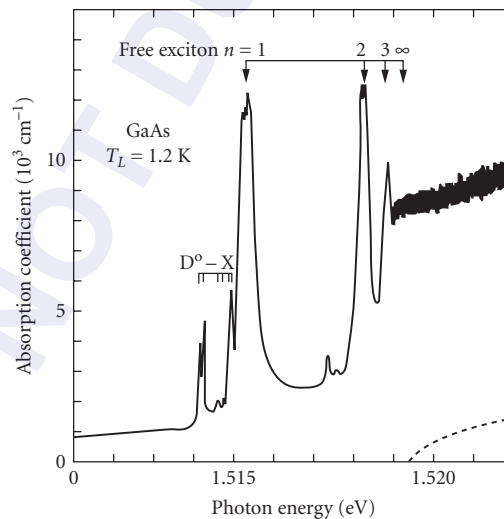


FIGURE 22 Absorption spectrum at 1.2 K of ultra-pure GaAs near the band edge. The $n = 1, 2, 3$ free exciton peaks are shown; also the bandgap E_g , determined by extrapolation to $n = \infty$, and impurity lines (D_v^0X) from the excitons bound to $\approx 10^{15} \text{ cm}^{-3}$ donors. (The rise at high energy is due to substrate absorption.) The dashed line shows the $(E - E_g)$ behavior expected in the absence of electron-hole interaction (the absolute magnitude is chosen to fit the absorption far from the band edge).⁵¹

seen followed by excited states with $n > 3$ leading smoothly into the continuum. The dashed line, calculated neglecting the effects of excitons, illustrates how important exciton effects are in understanding the optical properties of semiconductors and confirms the qualitative picture of exciton absorption.

Excitons in direct-gap semiconductors such as GaAs are called direct excitons. For indirect-gap semiconductors like Si or GaP, the absorption edge is determined by the influence of indirect excitons as revealed by the shape of the absorption. Such indirect-exciton transitions have been observed in several materials including Ge, Si, diamond, GaP, and SiC.

Real semiconductor crystals contain impurities and defects which also can affect the optical properties related to excitonic features, in addition to their causing impurity/defect absorption. A bound exciton (or bound-exciton complex) is formed by binding a free exciton to a chemical impurity atom (ion), complex, or a host lattice defect. The binding energy of the exciton to the defect or impurity is generally weak compared to the free-exciton binding energy. These bound excitons are extrinsic properties of the semiconductor; the centers to which the free excitons are bound can be either neutral donors or acceptors or ionized donors or acceptors. They are observed as sharp-line (width ≈ 0.1 meV) optical transitions in both absorption and photoluminescence spectra. The absorption or emission energies of these bound-exciton transitions always appear below those of the corresponding free-exciton transitions. Bound excitons are very commonly observed because semiconductors contain significant quantities of impurities or defects which produce the required binding. These complexes are also of practical interest because they characterize the impurities often used to control the electrical properties of semiconductors as well as being able to promote radiative recombination near the band gap. Bound excitons exhibit a polarization dependence similar to the free-exciton states from which they originated.

At higher densities of free excitons and low temperatures, they can form an electron-hole droplet by condensing into a “liquid” phase. This condensed phase occurs for electron-hole concentrations of about $2 \times 10^{17} \text{ cm}^{-3}$ and can be thought of as an electron-hole plasma with a binding energy of several meV with respect to the free excitons.⁵²

Polaritons Interesting optical effects arise when one considers explicitly the influence of longitudinal and transverse optical phonons on a transverse electromagnetic wave propagating through the semiconductor. This influence can be taken into account via the dielectric function of the medium. Dispersion curves that arise do not conform either to the photon or to the phonon. The coupling between the photon and phonon becomes so strong that neither can continue to be regarded as an independent elementary excitation, but as a photon-phonon mixture! This mixture can be regarded as a single quantity which can be interpreted as a new elementary excitation, the *polariton*.⁵³ Similar couplings exist between an exciton and the photon. It is an important consideration for interpreting some optical processes involving Raman and luminescence measurements.

High-energy transitions above the fundamental edge The optical properties of most semiconductors have been thoroughly investigated throughout the visible and ultraviolet regions where transitions above the fundamental gap energy give rise to properties strongly dependent upon photon energy. This regime is dominated by optical absorption and reflection of a photon arising from both valence and core electron transitions from the ground state of the system into various bound, auto-ionizing, continuum, or other excited states. The sum of all excitations—both bound (nonionizing) and ionizing—gives the total absorption coefficient and the complex dielectric constant at each photon energy $h\nu$. Photoemission or photoelectron spectroscopy measurements in this high-energy regime provide an alternative to ultraviolet spectroscopy for providing detailed information on the semiconductor. (See the electromagnetic spectrum in Fig. 1 in the Introduction for a reminder that photoemission measurements overlap UV measurements.) Electrons may be ejected from the semiconductor by high-energy photons as shown in Fig. 23.⁵⁴ Their kinetic energies are measured and analyzed to obtain information about the initial electron states. Ionizing excitations involve electron excitations into unbound states above the vacuum level. These excitations result in photoemissions depicted for the two valence bands and the core level shown. For photoemission involving one-electron excitations (usually dominant), binding energies E_b of valence and core levels are given directly by the measured kinetic energies E (from $N(E)$ in the figure and energy conservation):

$$E_b = h\nu - E - \phi \quad (40)$$

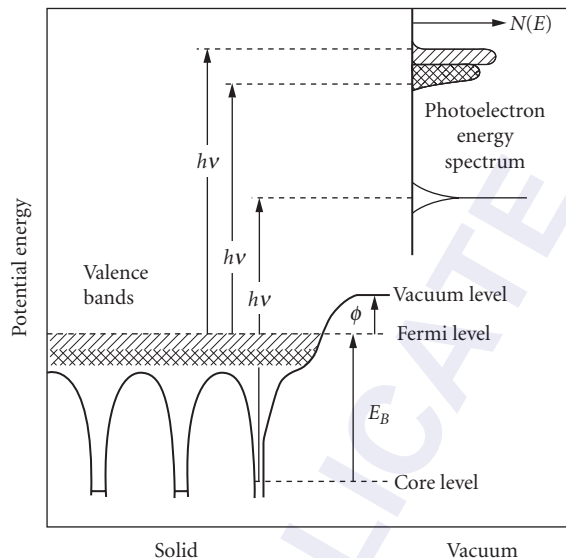


FIGURE 23 A schematic energy-level diagram showing photoemission from the valence bands and a core level in a solid.⁵⁴

Here ϕ is the work function (usually about 2 to 5 eV for most clean solids) which is known or easily measured. X-ray photoemission spectroscopy (XPS) is usually used to study core states and ultraviolet photoemission spectroscopy (UPS) to study valence band states. Photoemission techniques are also used for the study of surface states. Synchrotron radiation provides an intense source of light over a large spectral region. By measuring the angular distribution of the emitted UPS electrons, a direct determination of the E versus k relation for the valence band of GaAs can be made.

There are several regions of importance that must be considered in describing the optical properties of this high-energy region. Figure 24 shows the regions for InSb that are representative of the results for other semiconductors.⁵⁵ InSb is a narrow gap material with a direct bandgap of 0.17 eV at room temperature, so what is shown is at much greater energies than E_g . Sharp structure associated with transitions from the valence band to higher levels in the conduction band characterize the first region that extends to about 8 to 10 eV. This behavior is also characteristic for group-IV and other III-V compound semiconductors, see Fig. 25 which displays the imaginary part of the dielectric function for Si and GaAs. To show how this type of optical spectra can be interpreted in terms of the materials energy band structure, consider Fig. 26 which shows the spectral features of ϵ_2 for Ge in a and the calculated energy bands for Ge in b .^{57,58} Electronic transitions can take place between filled and empty bands subject to conservation of energy and wave vector. The initial and final electron wave vectors are essentially equal, and only vertical transitions between points separated in energy by $\hbar\omega = E_c(k) - E_v(k)$ are allowed.

The intensity of the absorption is proportional to the number of initial and final states and usually peaks when the conduction and valence bands are parallel in k -space. This condition is expressed by

$$\nabla_k [E_c(k) - E_v(k)] = 0 \quad (41)$$

Places in k -space where this is true are called critical points or Van Hove singularities. The experimental peaks can be sharp as shown in Fig. 26a because interband transitions in very pure crystals are not appreciably broadened by damping, and thus the lineshapes are determined primarily by the density of states, especially at low temperature. At room temperature, the broadening is dominated by electron-phonon scattering. Much information is available from the data if a good theory is used.

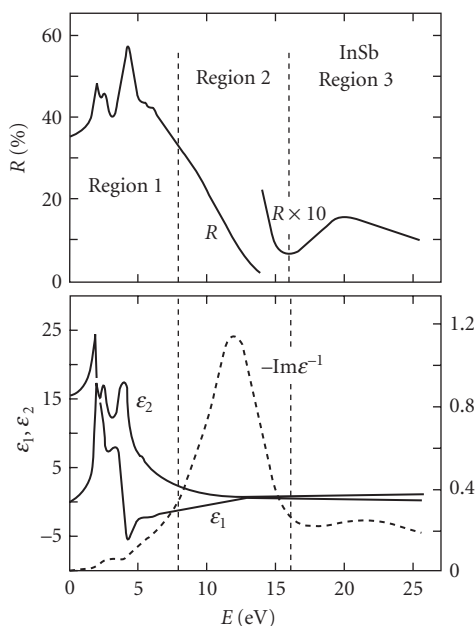


FIGURE 24 The spectral dependence of the reflectance R , the real and imaginary parts of the dielectric constant ϵ_1 and ϵ_2 , and the energy-loss function $-\text{Im}(\epsilon^{-1})$ for InSb.^{55,56}

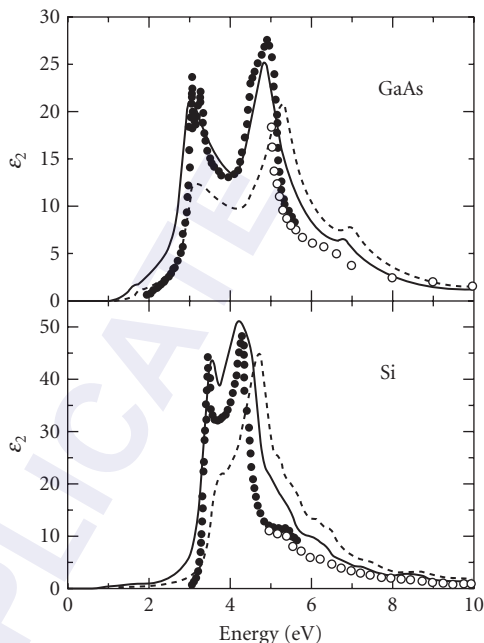


FIGURE 25 Imaginary part of the dielectric function for GaAs and Si. The symbols show experimental results. The dashed line is the result of a calculation based on the local-density approximation with GW corrections. The calculation shown by the solid line also takes into account excitonic effects using the Bethe-Salpeter equations.⁵⁹

Figure 26b shows pseudopotential energy-band calculations that show the special points and special lines in the Brillouin zone that give rise to the data shown in *a*. Band structure calculations using the local density approximation (see Fig. 26) are unable to predict the correct location of the peaks in the dielectric function. In general, the calculated gap is lower than the experimentally observed one (see Fig. 26). Considerable improvement to the accuracy of the excited states is possible when employing GW corrections to the band structure. While these corrections give an accurate band structure (Fig. 25, dashed line), they still ignore the excitonic interactions between the electron and hole participating in the interband transition and thus underestimate the magnitude of the absorption. Quite recently, a number of groups have been successful in considering excitonic Coulomb effects in their calculations using the Bethe-Salpeter equations. Modern calculations of the band structure and the dielectric function are almost as accurate as the best experimental methods, see Fig. 25.⁵⁹

The second region in Fig. 24 extends to about 16 eV and shows a rapid decrease of reflectance due to the excitation of collective plasma oscillations of the valence electrons. The behavior in this second “metallic” region is typical of the behavior of certain metals in the ultraviolet. One can think of the valence electrons as being essentially unbound and able to perform collective oscillations. Sharp maxima in the function $-\text{Im} \epsilon^{-1}$, which describes the energy loss of fast electrons traversing the material, have been frequently associated with the existence of plasma oscillations.

In the third region, the onset of additional optical absorption is indicated by the rise in reflectance. This structure is identified with transitions between filled *d* bands below the valence band and empty conduction band states. As shown in Fig. 27, the structure in region three is present in other III-V compounds, but is absent in Si which does not have a *d*-band transition in this region.⁶⁰

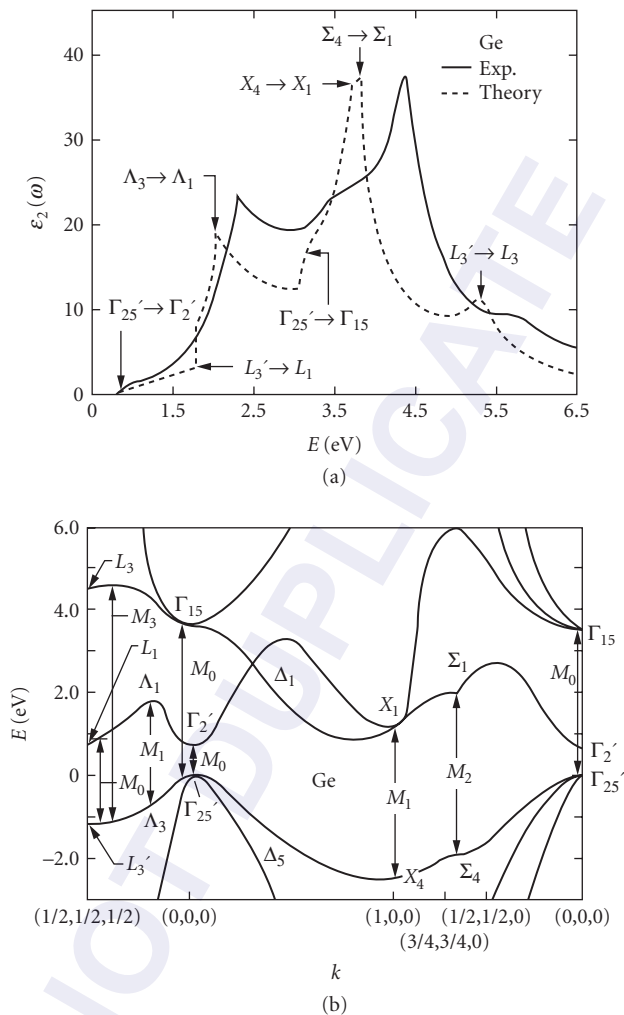


FIGURE 26 (a) Spectral features of ϵ_2 for Ge; (b) the calculated pseudopotential energy bands for Ge along some of the principal axes.^{37,58}

Other structure at higher photon energies is observed for Si as shown in Fig. 28 which shows the imaginary part of the dielectric constant of Si from 1 to 1000 eV.^{56,61} The large peaks on the left are due to excitations of valence electrons, whereas the peaks on the right are caused by excitation of core electrons from L shell states. K shell electrons are excited at energies beyond the right edge of the graph (beyond 1000 eV).

Elemental and compound semiconductors can often be mixed to form mixed-crystal alloys, such as $\text{Al}_{1-x}\text{Ga}_x\text{As}$ or $\text{Si}_{1-x}\text{Ge}_x$. The former alloy is used for optoelectronic applications (e.g., CD players), while the latter is found in high-end microprocessors or cellular telephones. The vibrational spectra of such alloys show multimode behavior (i.e., Si-Si, Ge-Ge, and Si-Ge vibrations in $\text{Si}_{1-x}\text{Ge}_x$ alloys) with phonon mode energies depending on the masses of the atoms. In the visible and UV spectral range, on the other hand, the peaks in the dielectric function given by the energy of the critical

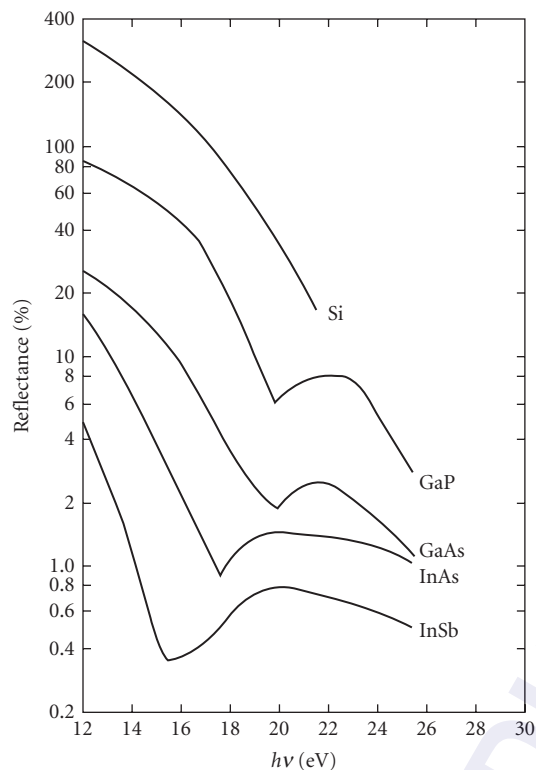


FIGURE 27 Reflectance of several semiconductors at intermediate energies. Starting from 12 eV the reflectance decreases, representing the exhaustion of bonding \rightarrow anti-bonding oscillator strength at energies greater than $2E_g$. The rise in reflectivity in the 15 to 20 eV range in the Ga and In compounds is caused by excitation of electrons from Ga cores (3d states) or In cores (4d states). The ordinate should be multiplied by 2 for InSb, 1 for InAs, 1/2 for GaAs, 1/4 for GaP, and 1/10 for Si.⁶⁰

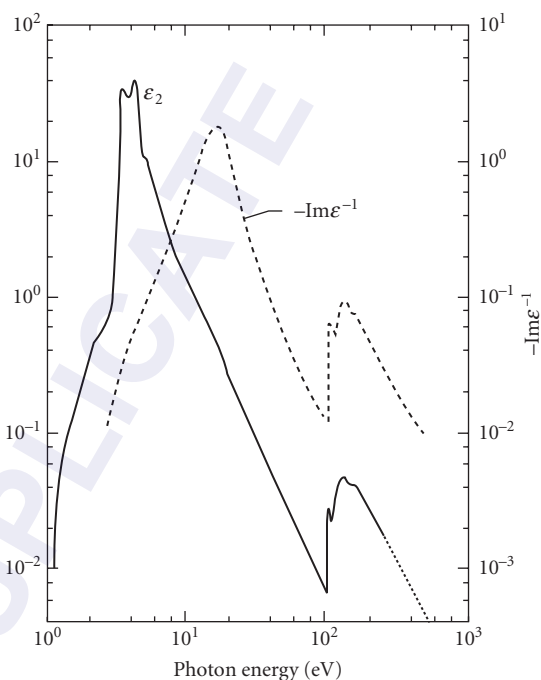


FIGURE 28 Imaginary part of the dielectric function for Si from 1 to 1000 eV.² (From reflectivity measurements of H.R. Philipp and H. Ehrenreich, Ref. 60, out to 20 eV and from transmission measurements of C. Gahwiller and F. C. Brown, Ref. 61, from 40 eV to 200 eV.)

points shift roughly linear with composition as shown in Fig. 29. The dielectric function of such alloys is known with very high accuracy (often also as a function of strain) and used in the production of modern microelectronic circuits to measure thickness and composition of semiconductor alloy films.^{62,63}

Free Carriers

Plasmons Semiconductors, in addition to a crystal lattice that may be ionic, may contain free charges as well.

The free-carrier contribution to the dielectric function is given by Maxwell's equation in CGS (Centimeter, Gram, Second) system of units

$$\epsilon_{fc}(\omega) = -i \frac{4\pi\sigma}{\omega} \quad (42)$$

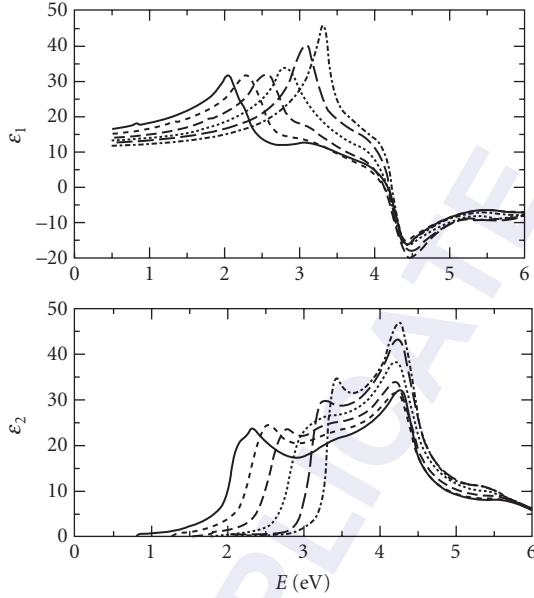


FIGURE 29 Real and imaginary parts of the dielectric function of $\text{Si}_{1-x}\text{Ge}_x$ alloys (bulk polycrystals) for x from 0 to 1 in steps of 0.2. Ge and Si are shown by the solid and long-dashed lines, respectively. The E_1 peak shifts from near 2 eV in Ge to about 3.4 eV in Si. Thin-film effects and strain modify these values in epitaxial films.⁶²

The task of establishing the functional form of $\varepsilon(\omega)$ hence reduces to one of determining the conductivity at the appropriate optical frequencies.

The response of a charge to an externally applied field may be described by classical methods, assuming a damping or resistive force to the charge that is proportional to the velocity of the charge. This simplification is known as the Drude approximation,⁶⁴ and it leads to the following relationship:

$$\sigma = \frac{Ne^2\tau}{m^*} \frac{1}{1 - i\omega\tau} \quad (43)$$

and is related to the dc conductivity by the relationship:

$$\sigma(0) = Ne^2\tau/m^* \quad (44)$$

where N is the free-carrier density, $1/\tau$ is the constant of proportionality for the damping force, and τ is a measure of the electron-electron collision time. Now,

$$\varepsilon_{fc}(\omega) = -\frac{i\omega^2\varepsilon_{p\infty}}{\omega\left(\omega + \frac{i}{\tau}\right)} \quad (45)$$

ω_p is the plasma frequency that describes oscillations of the plasma, that is, the delocalized charge cloud

$$\omega_p^2 = \frac{4\pi N e^2}{m^* \epsilon_\infty} \quad (46)$$

against the fixed crystal lattice.

In an ideal plasma with no damping, the $\epsilon(\omega)$ reduces to

$$\epsilon(\omega) = 1 - \frac{\omega_p^2}{\omega^2} \quad (47)$$

$\epsilon(\omega)$ is negative for $\omega < \omega_p$, which leads to total reflection and, hence, the term plasma reflectivity.

The phenomenon of optical reflection from a plasma reflection and the relationship to the free-carrier density are illustrated in Fig. 30 using the far-infrared reflection spectrum from a series of PbTe samples with hole densities extending from $3.5 \times 10^{18} \text{ cm}^{-3}$ to $4.8 \times 10^{19} \text{ cm}^{-3}$.⁶⁵ The plasma frequency increases with increasing carrier density as described by Eq. (46).

Coupled plasmon-phonon behavior Most semiconductor samples contain free carriers and phonons, and the frequencies of both are comparable. Hence, a complete description of the far-infrared optical properties has to take both into account. This can be achieved readily using Eq. (45) to describe the free carriers. The combined $\epsilon(\omega)$ may then be expressed as

$$\epsilon(\omega) = \epsilon_\infty + \frac{S\omega_{\text{TO}}^2}{\omega_{\text{TO}}^2 - \omega^2 - i\omega\gamma} - \frac{i\omega_p^2\epsilon_\infty}{\omega(\omega + i/\tau)} \quad (48)$$

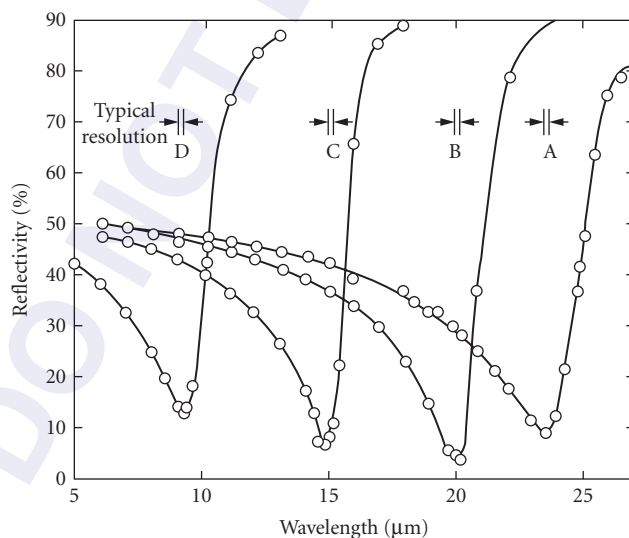


FIGURE 30 Reflectivity at 81 K and normal incidence of variously doped samples of *p*-type PbTe, showing the plasma resonance. Hole concentrations: A, $3.5 \times 10^{18} \text{ cm}^{-3}$; B, $5.7 \times 10^{18} \text{ cm}^{-3}$; C, 1.5×10^{19} ; D, $4.8 \times 10^{19} \text{ cm}^{-3}$.⁶⁵

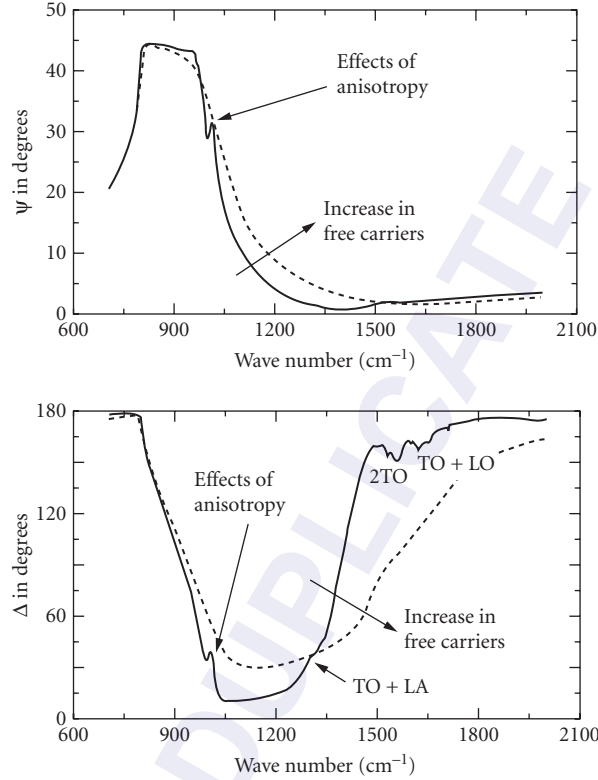


FIGURE 31 Mid-infrared ellipsometry spectra for 4H SiC. The reststrahlen band in the phonon absorption region (near 900 cm^{-1}) of the dielectric function is affected by free carriers and the crystal anisotropy. Two-phonon absorption is also visible.⁶⁶

A good example of the accurate description of the far-infrared behavior of a semiconductor is presented in Fig. 31.⁶⁶ The ellipsometric angles Ψ (related to reflectivity) and Δ (related to the polarization phase shift) of 4H SiC in the vicinity of the reststrahlen band (phonon absorption near 900 cm^{-1}) are affected by the free-carrier effects and the anisotropy of the crystal. All the major features in the complicated spectrum can be well described using the simple oscillatory models described.

The coexistence of phonons and plasmons leads to a coupling between the two participants.⁶⁷ Of particular interest are the coupled plasmon-LO phonon modes denoted by L_+ and L_- that are exhibited as minima in the reflection spectra. As explained earlier, the LO phonon frequencies occur at the zeros of the dielectric function $\epsilon(\omega)$. In the presence of plasmons, the zeros are shifted to the coupled mode frequencies L_+ and L_- . These frequencies can be determined directly for the case of no damping for both the phonon and the plasmon as shown below:

$$L_{\pm} = \frac{1}{2} \{ (\omega_{\text{LO}}^2 + \omega_p^2) \pm [(\omega_{\text{LO}}^2 - \omega_p^2)^2 + 4\omega_{\text{LO}}^2 \omega_p^2 (1 - \epsilon_{\infty} / \epsilon(0))]^{1/2} \} \quad (49)$$

Note that the presence of the plasmon introduces an additional low frequency zero at L_- . The relationship of the L_+ and L_- frequencies with the carrier density is presented in Fig. 32.⁶⁹ The existence of the coupled modes were predicted by Varga⁶⁷ and later observed using Raman scattering^{68,69} and far-infrared reflectivity.⁷⁰

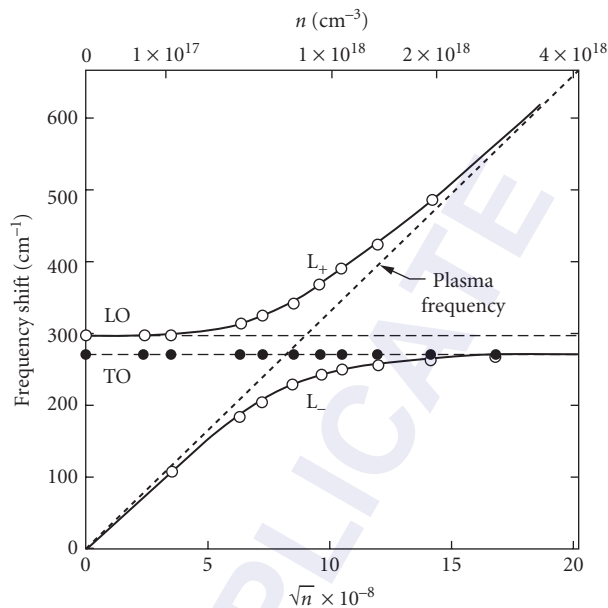


FIGURE 32 The solid curves labeled L_+ and L_- give the calculated frequencies of the coupled longitudinal plasmon-phonon modes, and the measured frequencies are denoted by the open circles.⁶⁹

Impurity and Defect Absorption The extended electronic states, excitons, lattice vibrations, and free carriers discussed thus far are all intrinsic to the pure and perfect crystal. In practice, real-life specimens contain imperfections and impurities. The characteristic optical properties associated with impurities and defects are the subject of discussion in this section. Two representative examples of the most widely observed effects, namely, shallow levels and deep levels in the forbidden gap, are considered in the following discussion.

Some of the effects due to impurities are considered in other parts of this chapter: impurity-related vibrational effects were considered under “Lattice”; excitons bound to impurity states were discussed under “Excitons”; and impurity-related effects in magneto-optical behavior are dealt with under “Magnetic-Optical Properties.” In addition to these effects, optical absorption due to electronic transitions between impurity-related electronic levels may also be observed in semiconductors.

The presence of impurities in a semiconductor matrix leads to both a perturbation of the intrinsic electronic quantum states and the introduction of new states, particularly in the forbidden energy gap. The major classes of electronic levels are the shallow levels that form the acceptor and donor states and lie close to the valence and conduction band extremes, respectively, and those that occur deep in the forbidden gap. The former are well known and are critical in controlling the electrical behavior of the crystal, and the latter are less well known but are, nevertheless, important in determining the sub-bandgap optical behavior.

Direct transitions from the shallow levels to the closest band extrema can be observed in the far-infrared transmission spectra of many semiconductors. An elegant example of this property is illustrated with the spectrum obtained from a high-purity Si wafer⁷¹ as displayed in Fig. 33. Sharp, well-resolved absorption features from electronic transitions due to B, P, As, and Al are present, as are additional features perhaps from unidentified impurities. Note that both acceptors and donor bands are observable due to the highly nonequilibrium state in which the specimen was maintained through the use of intense photoexcitation.

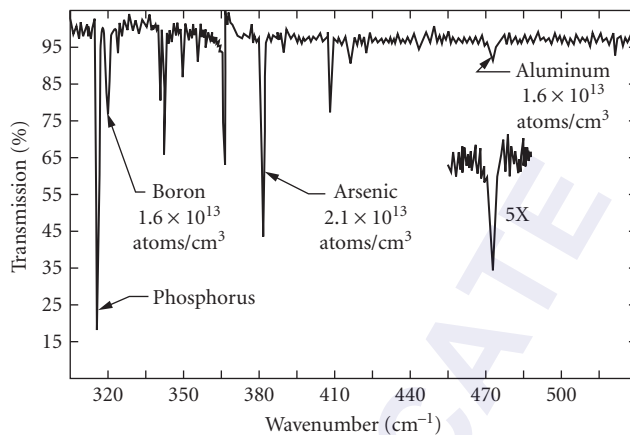


FIGURE 33 Total impurity spectrum of a 265- Ω -cm *n*-type Si sample obtained by the simultaneous illumination method. The input power to the illumination source was about 50 W.⁷¹

Electronic transitions from impurity- and defect-related levels deep in the forbidden gap can significantly alter the sub-bandgap behavior. One of the best known examples of this is the native defect level known as EL2 in GaAs. The level occurs 0.75 eV below the conduction band extremum, and, when present, it can completely dominate the subband gap absorption. The absorption spectrum recorded from a GaAs sample containing EL2 is presented in Fig. 34.⁷² The onset of the absorption at 0.75 eV is due to transitions to the conduction band extremum at the direct gap, and the features at 1.2 and 1.4 eV are due to transitions to higher-lying extrema.⁷² The figure also shows two spectra that exhibit the well-known photoquenching effect associated with EL2. When the specimen is subjected to intense white light radiation, the EL2 absorption is quenched, leaving only the band-to-band transitions with an onset at E_g which occurs at ≈ 1.5 eV at 10 K.

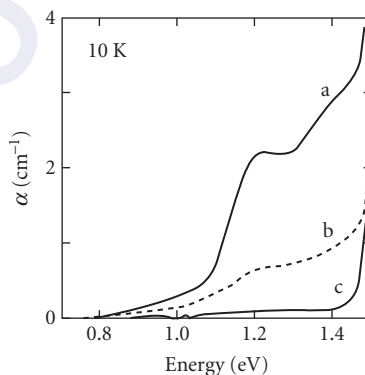


FIGURE 34 EL2 optical absorption spectra recorded at 10 K in the same undoped semi-insulating GaAs material. Curve a: after cooling in the dark; curves b and c: after white light illumination for 1 min and 10 min, respectively.⁷²

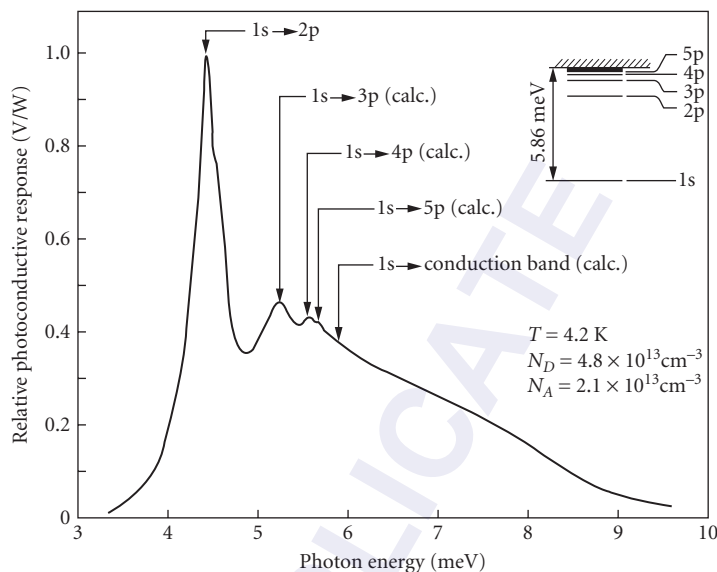


FIGURE 35 Far-infrared photoconductivity spectrum of a high-purity GaAs sample showing the measured transition energies and those calculated from the hydrogenic model using the $(1s \rightarrow 2p)$ transition energy. The hydrogenic energy level diagram is shown in the inset.⁷³

Optical measurements of shallow impurities in semiconductors have been carried out by absorption (transmission) and photoconductivity techniques. The photoconductivity method is a particularly powerful tool for studying the properties of shallow impurity states, especially in samples which are too pure or too thin for precise absorption measurements. In most cases, this type of photoconductivity can only be observed in a specific temperature range, usually at liquid helium temperatures. Figure 35 shows the photoconductivity response of a high-purity GaAs sample with specific transition energies that correspond to hydrogenic-like transitions.⁷³ If the excited states of the impurity were really bound states, electrons in these states could not contribute to the conductivity of the sample, and the excited state absorption would not result in peaks in the photoconductivity spectrum. However, there have been several suggestions as to how the electrons that are excited from the ground state to higher bound excited states can contribute to the sample's conductivity. First, if the excited state is broadened significantly by interactions with neighboring ionized donor and acceptor states, it is essentially unbound or merged with the conduction band. Other mechanisms for impurity excited-state photoconductivity all involve the subsequent transfer of the electron into the conduction band after its excitation to the excited state by the absorption of a photon. Mechanisms that have been considered for this transfer include (1) impact ionization of the electrons in the excited state by energetic free electrons, (2) thermal ionization by the absorption of one or more phonons, (3) photoionization by the absorption of a second photon, and (4) field-induced tunneling from the excited state into the conduction band. All of these mechanisms are difficult to describe theoretically.

Magneto-Optical Properties

Background Phenomena occurring as a result of the interaction of electromagnetic radiation with solids situated in a magnetic field are called magneto-optical (MO) phenomena. Studies of MO phenomena began in 1845 when Michael Faraday observed that the plane of vibration of polarized light rotated as it propagated through a block of glass in a strong magnetic field. By the 1920s, most

MO effects were fairly well understood in terms of the classical dynamics of an electron in a magnetic field. However, when semiconductors were first investigated in the early 1950s, a quantum mechanical interpretation of the MO data in terms of the energy-band structure was found to be necessary. MO spectroscopy was developed in the 1950s and 1960s as a powerful tool for characterizing the fundamental electronic properties of semiconductors. Cyclotron resonance experiments using microwaves, interband studies using broadband visible and infrared sources, and finally lasers were used to characterize the symmetry of band structures, as well as properties such as energy gaps, effective masses, and other band parameters. The two major limitations of the classical theory are that no effects depending on the density of states are predicted and no effects of electron spin are included.

Table 6 presents an overview of the typical types of magneto-optical phenomena observed in semiconductors and the information that can be determined from the experimental measurements. Four classes of MO phenomena can be distinguished: those arising from (1) interband effects, (2) excitonic effects, (3) intraband or free-carrier effects, and (4) impurity magnetoabsorption effects. Further clarification can then be made by determining whether the effect is absorptive or dispersive, resonant or nonresonant, and upon the relative orientation of the magnetic field to the direction of propagation of the electromagnetic radiation and its polarization components. Resonant experiments usually provide more detailed information about the band structure of a semiconductor and often are easier to interpret. There is thus a wide variety of effects as shown which can give different types of information about the crystal's energy-band structure, excitonic properties, and impurity levels. Before summarizing and discussing each of these magneto-optical effects, it is necessary to briefly describe the effects of a magnetic field on the energy-band structure of a semiconductor.

Effect of a Magnetic Field on the Energy Bands Magneto-optical experiments must be analyzed with specific energy-band models in order to extract the related band parameters and to emphasize the underlying physical concepts with a minimum of mathematical complexity. Most often, one deals with only the highest valence bands and the lowest conduction bands near the forbidden energy-gap region. If simple parabolic bands are assumed, a fairly complete analysis of the MO experiments is usually possible, including both the resonant transition frequencies and their line shapes. On the other hand, if more complicated energy bands (e.g., degenerate, non-parabolic) are needed to describe the solid, the detailed analysis of a particular experiment can be complicated.

The effect of a magnetic field on a free electron of mass m^* was determined in 1930 by Landau, who solved the Schroedinger equation. The free electrons experience a transverse Lorentz force which causes them to travel in orbits perpendicular to the magnetic field. The resulting energy eigenvalues corresponding to the transverse components of the wave vector are quantized in terms of harmonic oscillator states of frequency ω_c , while a plane-wave description characterizes the motion along the magnetic field. The allowed energy levels, referred to as Landau levels, are given by

$$E_n^\pm = \left(n + \frac{1}{2} \right) \hbar \omega_c + \frac{\hbar^2 k_z^2}{2m^*} \pm \frac{1}{2} g^* \mu_B B \quad (50)$$

where n is the Landau level number (0, 1, 2, . . .), ω_c (the cyclotron frequency) is equal to eB/m_c^* , and m_c^* is the cyclotron effective mass. The middle term represents the energy of an electron moving along the direction of the \mathbf{B} field in the z direction; it is not quantized. The first term represents the quantized energy of motion in a plane perpendicular to the field. The last term represents the effect of the electron's spin; g^* is the effective spectroscopic g -factor or spin-splitting factor and

$$\mu_B = e\hbar/2m_0 = 5.77 \times 10^{-2} meV/T \quad (51)$$

is the Bohr magneton. The g -factor in a semiconductor has values quite different from the usual value of two found for atomic systems (e.g., for narrow-energy gaps and a strong spin-orbit interaction, g^* can be large and negative).

TABLE 6 Magneto-Optical Phenomena and Typical Information Obtainable

Magneto-Optical Effect	Some Information or Properties Obtainable
Interband effects	
<i>In transmission</i>	
Band-to-band magnetoabsorption.....	Energy gaps, effective masses, g-factors, higher band parameters
Faraday rotation (resonant and nonresonant).....	Energy gaps, effective masses
Faraday ellipticity (nonresonant).....	Relaxation times
Voigt effect (resonant and nonresonant)	Effective masses
Cross-field magnetoabsorption	
<i>In reflection</i>	
Magnetoreflexion	Studies of very deep levels or where absorption is high, similar information as in transmission magnetoabsorption
Kerr rotation	
Kerr ellipticity	
Magneto-excitonic effects	
In transmission, reflection, and photoconductivity	Diamagnetic and Zeeman shifts and splittings; energy gap; effective reduced-mass tensor; effective Rydberg; anisotropy parameter; dielectric tensor components; effective g-factors; effective masses; quality of materials, structures, alloys and interfaces
Intraband or free-carrier effects	
<i>In transmission</i>	
Cyclotron resonance (resonant)	Effective masses, relaxation times, nonparabolicity
Combined resonance.....	Same information as cyclotron resonance plus g-factors
Spin resonance.....	g-factors
Phonon-assisted cyclotron resonance harmonics.....	Same information as cyclotron resonance plus phonon information
Faraday rotation (resonant and nonresonant)	Carrier concentration, effective masses; use for impure materials, flexible, can be used at high temperature
Faraday ellipticity (nonresonant)	
Voigt effect (nonresonant)	
Interference fringe shift (nonresonant)	
Oscillatory variation of the Shubnikov-de Haas type	Carrier concentration
<i>In reflection</i>	
Magnetoplasma reflection	Effective masses, carrier concentration
Magnetoplasma rotation (Kerr effect)	
Magnetoplasma ellipticity	
Impurity magneto-absorption	
Zeeman and diamagnetic effect-type behavior of impurities	Hydrogenic impurity information, binding energy, effective masses, effective Rydberg, central cell corrections, static dielectric constant; both impurity and Landau level information; impurity information as above plus Landau-level information related to effective masses and g-factors
Photoionization behavior (transitions from ground state of impurity to Landau levels)	

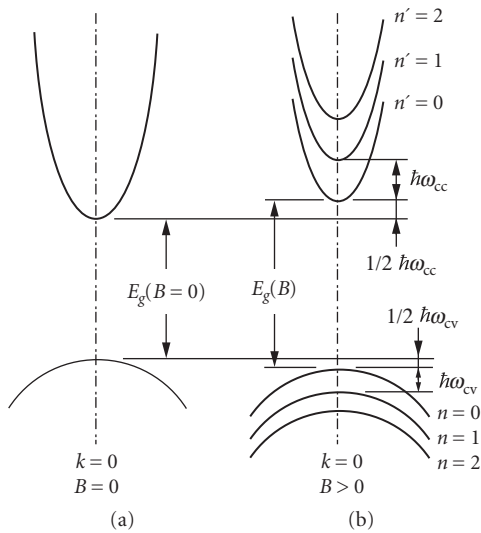
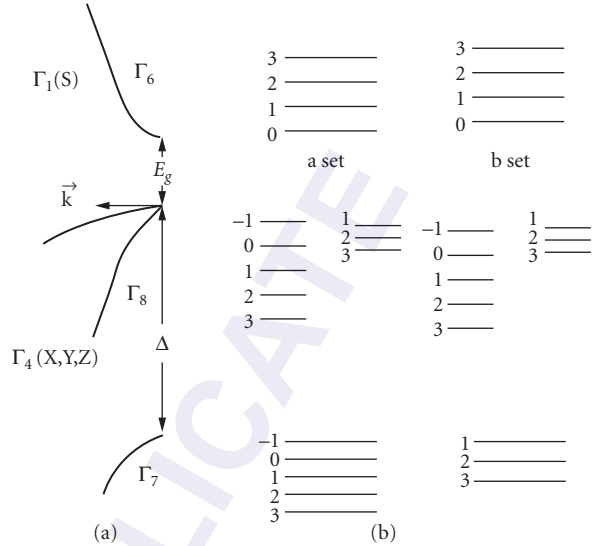

FIGURE 36 Landau levels for simple bands.

FIGURE 37 (a) Zincblende semiconductor energy bands ($H=0$) and (b) in an applied magnetic field H .⁷⁴

Figure 36 shows schematically the effect of a magnetic field on a simple parabolic direct gap extrema at $k=0$. The quasi-continuous parabolic behavior of the nondegenerate conduction and valence bands at $\mathbf{B}=0$ (shown in *a*) is modified by the application of a magnetic field into Landau levels as shown in *b*. Each Landau level is designated by an integer $n=0, 1, 2, 3, \dots$. Finally, in Fig. 37, the Landau effects are shown schematically for zincblende energy bands when both spin and valence band degeneracy are taken into account for both $\mathbf{B}=0$ (*a*) and $\mathbf{B} \neq 0$ (*b*).⁷⁴ With spin included, the valence band splits into the Γ_8 and Γ_7 bands with spin-orbit splitting energy Δ . The quantization of the bands into Landau levels is illustrated in the right-hand side of the figure. The Pidgeon and Brown⁷⁵ model has been successfully used to describe the magnetic field situation in many semiconductors, since it includes both the quantum effects resulting from the partial degeneracy of the p-like bands and the nonparabolic nature of the energy bands. The *a*-set levels are spin-up states, and the *b*-set, the spin-down states. These large changes in the \mathbf{E} versus k relations of the bands when a magnetic field is applied also means large changes in the density of states which become periodic with a series of peaks at energies corresponding to the bottom of each Landau level. This oscillatory variation in the density of states is important for understanding the various oscillatory phenomena in a magnetic field and is a key advantage of using magnetic fields to study semiconductors.

Interband Magneto-Optical Effects Interband transitions in a magnetic field connect Landau-level states in the valence band to corresponding states in the conduction band. Thus, they yield direct information concerning energy gaps, effective masses, effective g-factors, and higher band parameters. The strongest allowed transitions are those that are proportional to the interband matrix element $p = - (i \hbar / m_0) \langle s | p_j | x_j \rangle$, where $j = x, y, z$.

This matrix element directly connects the p-like valence band (x, y, z) which is triply degenerate with the s-like conduction band through the momentum operator p_j . The transition energies can be calculated directly from a knowledge of the selection rules and use of an energy-band model. The selection rules are given by⁷⁶

$$\sigma_L: a(n) \rightarrow a(n-1), b(n) \rightarrow b(n-1)$$

$$\sigma_R: a(n) \rightarrow a(n+1), b(n) \rightarrow b(n+1)$$

$$\pi: a(n) \rightarrow b(n+1), b(n) \rightarrow a(n-1)$$

where σ_L , σ_R , and π are left circular, right circular, and linear ($\mathbf{e} \parallel \mathbf{B}$) polarizations. As discussed earlier, a and b denote the spin-up and spin-down states, and n the Landau-level number. Often, sharp optical transitions between the Landau levels are observed, providing highly accurate information about the fundamental band parameters such as the energy gap E_g , effective masses of the electrons and holes, higher band parameters, etc.

InSb is a material in which interband effects have been studied very extensively; thus, it is a good representative example. Even though it is a narrow-gap semiconductor and thus has a small exciton binding energy, Weiler⁷⁷ has shown that excitonic corrections must be made to properly interpret the magnetic-field-dependent data. This shall be discussed in more detail in the next section. The band models discussed earlier predict that there is an increase in the energy of the absorption edge with magnetic field because of the zero point energy $1/2\hbar\omega_c$ of the lowest conduction band. At larger photon energies, transmission minima (absorption maxima) which are dependent upon magnetic field are observed. By plotting the photon energy positions of the transmission minima against magnetic field, converging, almost linear plots are obtained as shown in Fig. 38.⁷⁵ Extrapolation of the lines to zero field gives an accurate value for the energy gap. Use of a band model and specific transition assignments further allow the determination of other important band parameters, such as effective masses and g -factors.

Magnetoreflexion Besides the changes in absorption brought about by the magnetic field, there are changes in the refractive index. Since the reflectivity depends upon both the real and imaginary

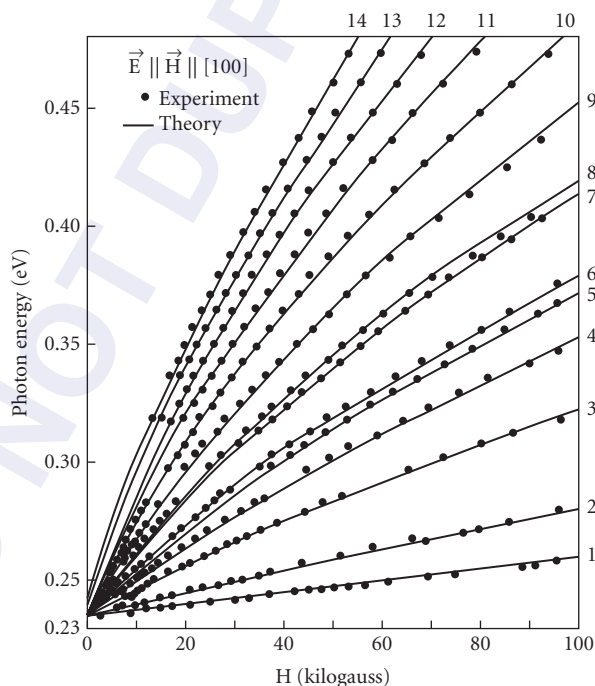


FIGURE 38 Energy values of transmission minima versus magnetic field for electron transitions between Landau levels of valence and conduction bands in InSb. Plot of the photon energy of the principal transmission minima as a function of magnetic field for $\mathbf{E} \parallel \mathbf{H} \parallel [100]$. The solid lines represent the best theoretical fit to the experimental data. The numeral next to each line identifies the quantum assignment.⁷⁵

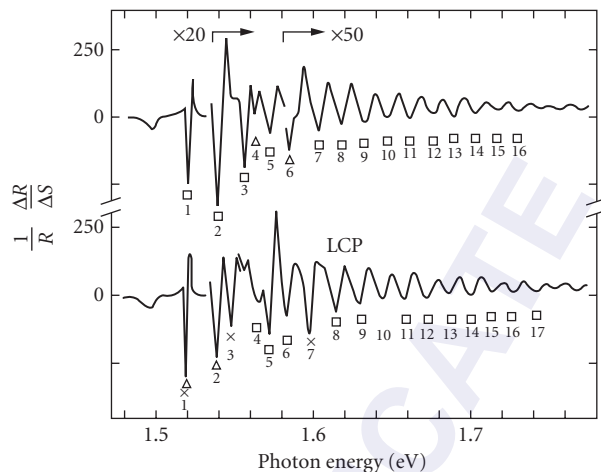


FIGURE 39 Stress-modulated magnetorefectance spectra for the fundamental edge in epitaxially grown high-purity $\langle 211 \rangle$ GaAs at $T \approx 30$ K, with $\Delta S = 5 \times 10^{-5}$, observed in the Faraday configuration with magnetic field $\mathbf{H} \parallel [11\bar{2}]$ and $H = 88.6$ kG. The number directly below each prominent transition refers to the identification of the transitions, Δ , LHA; \times , LHB; \square , HH(AV).⁷⁸

parts of the index, clearly it is affected by the field. Interband transitions are often observed in reflection because of the high absorption coefficients. In addition, modulation spectroscopy techniques, in which a parameter such as stress, electric field, wavelength, magnetic field, etc. is periodically varied and the signal synchronously detected, provide several orders of magnitude enhancement in the sensitivity for observing resonant transitions. This is especially important for the observation of higher energy transitions lying far away from the energy of the fundamental gap. Figure 39⁷⁸ shows the stress-modulated magnetorefectance spectra for the fundamental edge region, the interpretation of which must involve the effect of excitons on the transitions. Quantitative information about the split-off valence to conduction band edge in GaAs is also obtainable from magnetoreflexion.⁷⁹

Faraday rotation A plane-polarized wave can be decomposed into two circularly polarized waves. The rotation of the plane of polarization of light as it propagates through the semiconductor in a direction parallel to an applied magnetic field is called the Faraday effect, or Faraday rotation. The amount of rotation is usually given by the empirical law $\Theta = VB$, where Θ is the angle of rotation, V is the Verdet coefficient, \mathbf{B} is the magnetic field value, and l is the thickness. The Verdet coefficient is temperature, wavelength, and sometimes field dependent. The Faraday effect can then be understood in terms of space anisotropy effects introduced by the magnetic field upon the right and left circularly polarized components. The refractive indices and propagation constants are different for each sense of polarization, and a rotation of the plane of polarization of the linearly polarized wave is observed. The sense of rotation depends on the direction of the magnetic field. Consequently, if the beam is reflected back and forth through the sample, the Faraday rotation is progressively increased. When measurements are thus made, care must be taken to avoid errors caused by multiple reflections. Faraday rotation may also be considered as birefringence of circularly polarized light.

If absorption is present, then the absorption coefficient will also be different for each sense of circular polarization, and the emerging beam will be elliptically polarized. Faraday ellipticity specifies the ratio of the axes of the ellipse.

Faraday rotation can be observed in the Faraday configuration with the light beam propagating longitudinally along the \mathbf{B} -field direction. Light propagating transverse to the field direction is

designated as the Voigt configuration. Two cases must be distinguished—the incident beam may be polarized so that its \mathbf{E} field is either parallel or perpendicular to \mathbf{B} . The Voigt effect is birefringence induced by the magnetic field and arises from the difference between the parallel and perpendicular indices of refraction in the transverse configuration. It is usually observed by inclining the incident plane-polarized radiation with the electric vector at 45° to the direction of \mathbf{B} . The components resolved parallel and perpendicular to \mathbf{B} then have different phase velocities and recombine at the end of the sample to give emerging radiation which is elliptically polarized. Measurements of this ellipticity then determine the Voigt effect.

The interband Faraday effect is a large effect and is therefore useful for characterizing semiconductors. For frequencies smaller than the frequency corresponding to the energy gap, the interband Faraday effect arises from the dispersion associated with the interband magnetoabsorption. In this region, it has been used to determine energy gaps and their pressure and temperature dependence. Since the beam propagates through the crystal in a transparent region of the spectrum, it may be attractive to use in certain applications. For example, in GaAs,⁸⁰ at long wavelengths the Faraday effect has a positive rotation, while near the gap the Verdet coefficient becomes negative. For frequencies equal to or larger than the frequency corresponding to the energy gap, the Faraday rotation is dominated by the nearest magneto-optical transition. Oscillatory behavior, like that seen in the magnetoabsorption, is also often observed.

Diluted magnetic semiconductors (DMS) are a class of materials that have attracted considerable scientific attention. Any known semiconductor with a fraction of its constituent ions replaced by some species of magnetic ions (i.e., ions bearing a net magnetic moment) can be defined as a member of this group. The majority of DMS studied so far have involved Mn^{+2} ions embedded in various II–VI hosts. The optical properties of DMS are controlled by the interaction between the localized magnetic moments of Mn^{+2} and the conduction and/or valence band electrons (referred to as the sp-d interaction), which results in features unique to DMS. The best known (and quite spectacular) of these are the huge Faraday rotations of the visible and near-infrared light in wide-gap DMS. The origin of the large rotations is the sp-d exchange interaction which makes the band structure much more sensitive to the strength of external magnetic fields than in ordinary semiconductors. Figure 40 shows the Faraday effect in $\text{Cd}_{1-x}\text{Mn}_x\text{Se}$ ($x = 0.25$) at $T = 5$ K.⁸¹ Each successive peak represents an additional Faraday rotation of 180° .

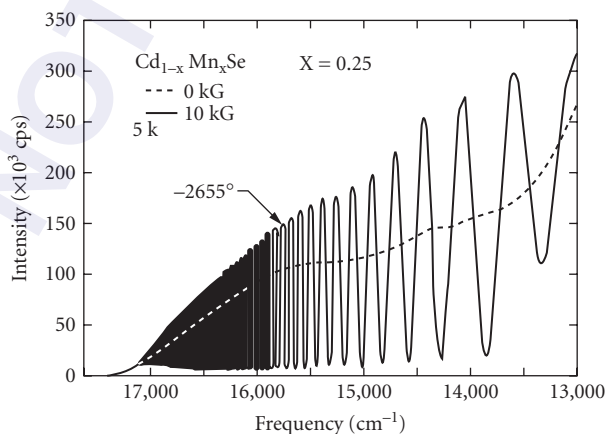


FIGURE 40 Faraday effect in $\text{Cd}_{1-x}\text{Mn}_x\text{Se}$ ($x = 0.25$), $T = 5$ K. The dashed line represents the transmission of light as a function of its frequency for a 3-mm-thick sample located between two polaroids with their axes at 45° and at zero magnetic field. The light propagates along the optic axis, \hat{c} . The oscillations occur when a 10-kG magnetic field is applied along \hat{c} .⁸¹

Excitonic magneto-optical effects As described earlier, a free exciton consists of an electron and hole bound together electrostatically. When the pair has an energy less than that of the energy gap, they orbit around each other. If the orbital radius is large compared with the lattice constant, they can be approximately treated as two point charges having effective masses and being bound together by a Coulomb potential that gives rise to a hydrogen-atom-like behavior. In the presence of a magnetic field, excitons give rise to Zeeman and diamagnetic effects analogous to those in atomic spectra. Fine structure can occur due to motions other than the simple orbiting of an electron and hole—the carriers can have intrinsic motion, motion around an atom, spin motion, and motion of the complete exciton through the lattice. Some of these motions may even be coupled together.

Bound excitons (or bound-exciton complexes) or impurity-exciton complexes are extrinsic properties of materials. Bound excitons are observed as sharp-line optical transitions in both photoluminescence and absorption. The bound exciton is formed by binding a free exciton to a chemical impurity atom (or ion), a complex, or a host lattice defect. The binding energy of the exciton to the impurity or defect is generally smaller than the free-exciton binding energy. The resulting complex is molecular-like (hydrogen-molecule-like), and bound excitons have many spectral properties analogous to those of simple diatomic molecules.

The application of a magnetic field to samples where excitonic features are observed in the absorption spectra results in line splittings, energy shifts, and changes in linewidths. These arise from diamagnetic and Zeeman effects just as in atomic or molecular spectroscopy. The treatment of the problem of an exciton in a magnetic field in zincblende-type structures is difficult due to the complexity of the degenerate valence band. Often a practical solution is adopted that corrects the interband model calculations for exciton binding energies that are different for each Landau level. Elliott and Loudon showed that the absorption spectrum has a peak corresponding to the lowest $N = 0$ hydrogen-like bound state of the free exciton, which occurs below the free interband transition by the exciton binding energy E_B .⁸² Weiler suggests that for transitions to the conduction band Landau level n , the exciton binding energy E_B can be approximated by⁷⁴

$$E_B(n) \approx 1.6R[\gamma_B/(2n+1)]^{1/3} \quad (52)$$

where R is the effective Rydberg,

$$R = R_0 \mu / m_0 \epsilon(0) \quad (53)$$

$R_0 = 13.6$ eV, $\epsilon(0)$ is the static-dielectric constant, μ is the reduced effective mass for the transition, and γ_B is the reduced magnetic field

$$\gamma_B = m_0 S / 2\mu R \quad (54)$$

and

$$S = \hbar e B / m_0 \quad (55)$$

Thus after calculating the interband transition energy for a particular conduction band Landau level n , one subtracts the above binding energy to correct for exciton effects.

Nondegenerate semiconductors, in particular, materials belonging to the wurtzite crystal structure, have been extensively studied both with and without a magnetic field. Exciton states in CdS have been studied by high-resolution two-photon spectroscopy in a magnetic field using a fixed near-infrared beam and a tunable visible dye laser.^{83,84} Figure 41 shows the photoconductive response versus total photon energy near the A-exciton region.⁸³ The two-photon transitions involve P states, and both 2P and 3P states are clearly seen. As the field is increased, both Zeeman splittings and diamagnetic shifts occur. Figure 42 shows both experimental and theoretical transition energies versus \mathbf{B} field.⁸³ Excellent agreement is obtained by using variational calculations that have been successfully used to describe impurity atoms in a magnetic field.

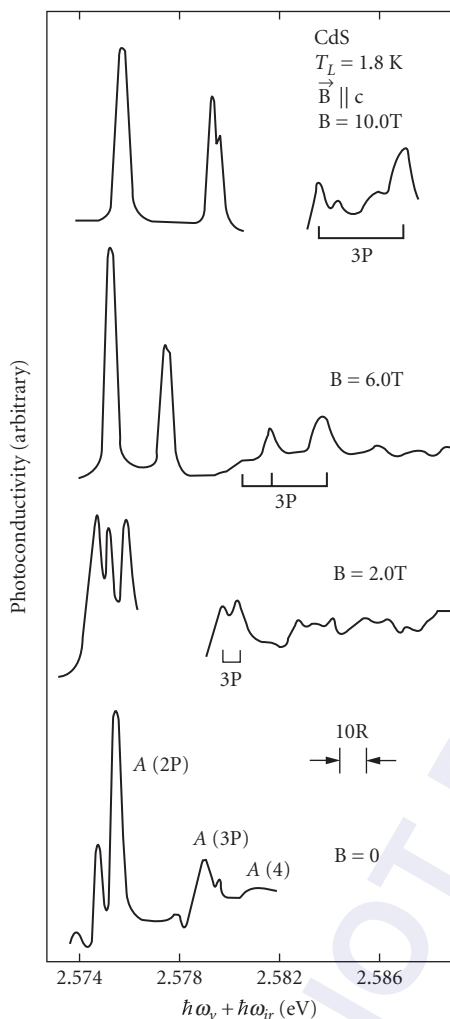


FIGURE 41 Photoconductivity vs. total photon energy $\hbar\omega_v + \hbar\omega_{ir}$ near the A-exciton region in CdS platelets for various magnetic fields. The magnetic field was parallel to the hexagonal c axis in a Voigt configuration with E perpendicular to c for the two photons at a lattice temperature of $T_L = 1.8$ K. The instrumental resolution $R = 0.1$ meV is narrower than the intrinsic linewidths.⁸³

Intraband or Free-Carrier Effects

Cyclotron resonance Cyclotron resonance absorption of free carriers is the simplest and most fundamental magneto-optical effect and provides a direct determination of carrier effective masses. Classically, it is a simple phenomenon—charged particles move in circular orbits (in planes perpendicular to the direction of the magnetic field) whose radii increase as energy is absorbed from the applied electric fields at infrared or microwave frequencies. After a time τ , a collision takes place and

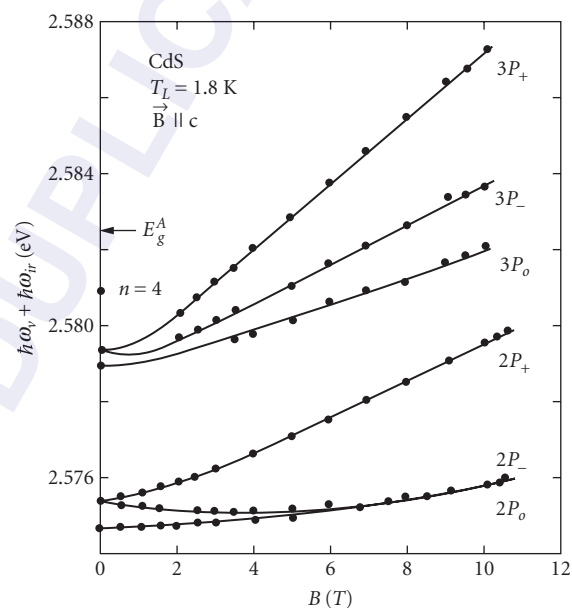


FIGURE 42 Peak positions, in total photon energy $\hbar\omega_v + \hbar\omega_{ir}$, for the 2P and 3P A excitons in CdS platelets as a function of applied \mathbf{B} field. The solid points were determined experimentally, and the solid curves are theoretically obtained from variational calculations of the diamagnetic shifts along with use of the experimental g-factors.⁸³

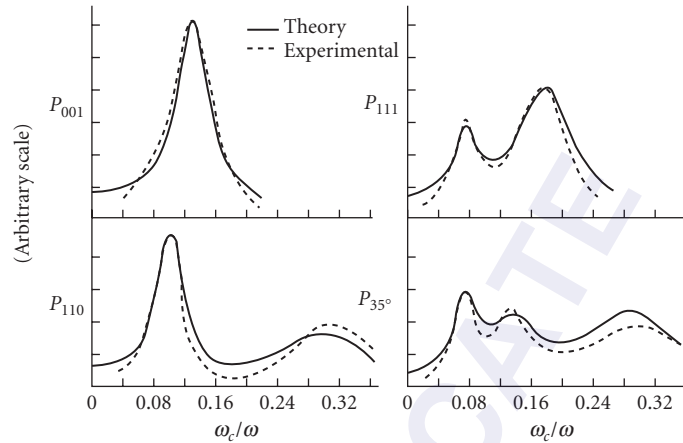


FIGURE 43 Microwave absorption in intrinsic n -type Ge at 4.2 K for four different directions of \mathbf{B} in the (110) plane as a function of the magnetic field. P_{35° represents the absorption with \mathbf{B} at 35° to the [001] axis.⁸⁸

the absorption process begins again. From the resonance relation $\omega_c = e\mathbf{B}/m^*$, extensive and explicit information about the effective masses and the shape of energy surfaces near the band extrema can be obtained. Excellent reviews of cyclotron resonance have been previously published by McCombe and Wagner^{85,86} and Kobori et al.⁸⁷

A classical example of the explicit band structure information that can be obtained from cyclotron resonance experiments is given for Ge. Figure 43 shows the microwave absorption for n -type germanium at 4.2 K for four different \mathbf{B} -field directions, each peak corresponding to a specific electron effective mass.⁸⁸ Figure 44 shows the orientation dependence of the effective masses obtained from the cyclotron resonance experiments which demonstrates that there is a set of crystallographically equivalent ellipsoids oriented along all $\langle 111 \rangle$ directions in the Brillouin zone.⁸⁹ Figure 45 shows the band structure of Ge that these measurements helped to establish⁹⁰; a illustrates the conduction band minima along the $\langle 111 \rangle$ direction at the zone edge and b , the eight half-prolate ellipsoids of revolution or four full ellipsoids. The longitudinal and transverse masses are $m_l^* = 1.6m_0$ and $m_t^* = 0.082m_0$, respectively. Both holes and electrons could be studied by using light to excite extra carriers. This illustrates the use of cyclotron resonance methods to obtain band structure information. Since that time, numerous experiments have been carried out to measure effective mass values for carriers in various materials.

The beginning of modern magneto-optics in which “optical” as opposed to “microwave” techniques were used, began in 1956 with the use of infrared frequencies at high magnetic fields. Far-infrared lasers are extremely important for modern-day measurements of cyclotron resonance as seen in Fig. 46 for n -type InSb.⁹¹ At low temperatures, only the lowest CR transition C_1 (0^+ to 1^+) is seen. Raising the temperature populates the higher-lying Landau levels, and other CR transitions are seen at different fields because of the nonparabolicity of the conduction band which gives rise to an energy-dependent effective mass. At 13 K, a second transition C_2 (0^- to 1^-) is seen and at 92 K, C_3 (1^+ to 2^+). The low field feature denoted by I is called impurity cyclotron resonance because, although it is a neutral donor excitation, its appearance resembles that of the regular CR magnetoabsorption. It results from neutral donors exhibiting a Zeeman transition ($1s \rightarrow 2p^+$), the transition energy of which is much larger than the ionization energy. The I signal gradually disappears as the temperature increases because the donors become ionized.

Cyclotron resonance also serves as a valuable tool for materials characterization through making use of the cyclotron resonance linewidth and intensity. Resonance linewidths can differ considerably from sample to sample. This difference in linewidth is attributed to differences in impurity content, with the higher-purity samples having the narrowest linewidths and largest intensities. Figure 47 shows the electron CR signals for both n -type and p -type GaAs crystals with low compensation.⁸⁷

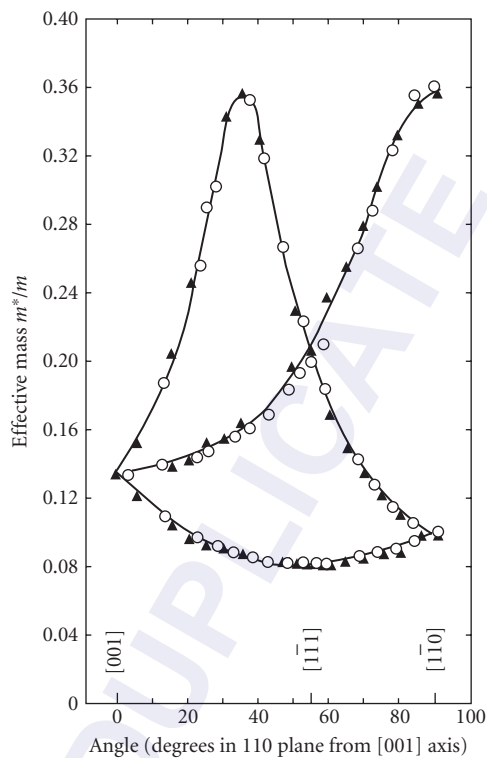


FIGURE 44 Effective mass of electrons in Ge at 4 K for magnetic field directions in a (110) plane.⁸⁹

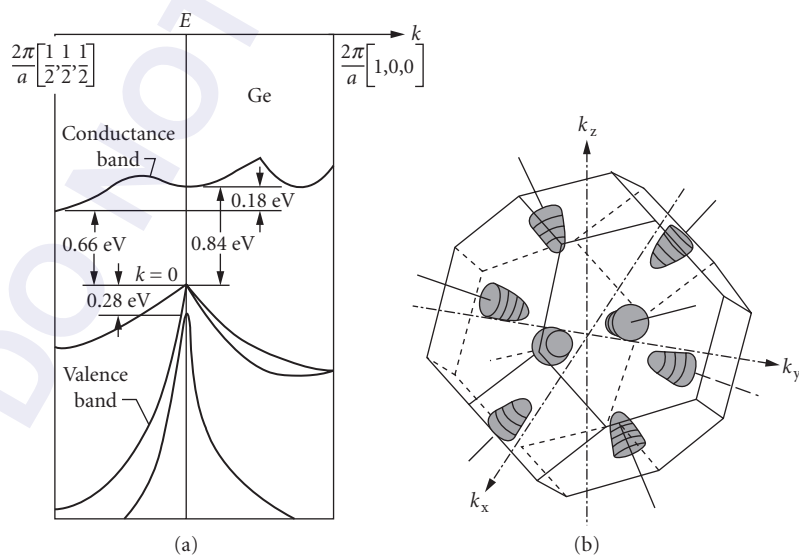


FIGURE 45 (a) Band structure of Ge plotted along the [100] and [111] directions and (b) ellipsoidal energy surface corresponding to primary valleys along the $\langle 111 \rangle$ directions.⁹⁰

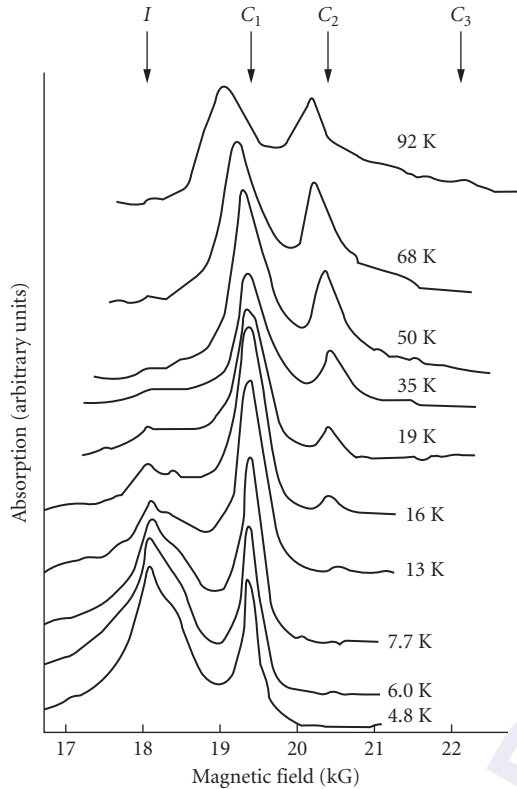


FIGURE 46 Thermal equilibrium resonance traces in n -type InSb at various temperatures. At 4.8 K, only the lowest cyclotron transition $C_1(0^+ \rightarrow 1^+)$ and impurity cyclotron resonance I (ICR) are visible. On raising the temperature, the signal I disappears on complete ionization of donors, while the second and third cyclotron transitions $C_2(0^- \rightarrow 1^-)$ and $C_3(1^+ \rightarrow 2^+)$ start to show up.⁹¹

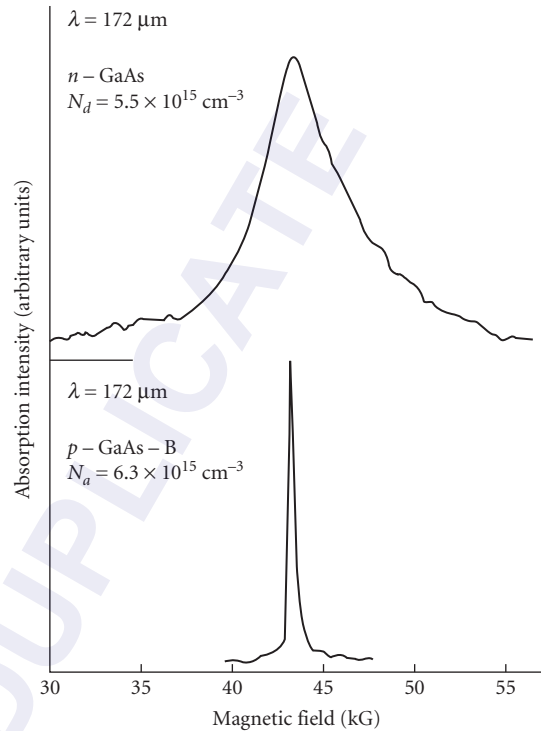


FIGURE 47 Difference in electron cyclotron resonance linewidth between n - and p -type GaAs crystals, having the same order of donor or acceptor concentrations.⁸⁷

The observed large difference in linewidths is primarily considered to reflect the difference in the electron-donor and electron-acceptor scattering rates.

An electron placed in the conduction band of a polar insulator or semiconductor surrounds itself with an induced lattice-polarization charge. The particle called a polaron consists of the electron with its surrounding lattice-polarization charge. The term *magnetopolaron* is also often referred to as a polaron in a magnetic field. Landau-level energies of these magnetopolarons are shifted relative to those predicted for band electrons. These energy shifts give rise to polaron effects that are most clearly evident in optical experiments such as cyclotron resonance. The review by Larsen provides an annotated guide to the literature on polaron effects in cyclotron resonance.⁹²

Free-carrier Faraday rotation and other effects Observation of free-carrier rotation was first reported in 1958. It is best understood as the differential dispersion of the cyclotron resonance absorption, and, as such, it is an accurate method for determining carrier effective masses. It can be measured off resonance and detected under conditions which preclude the actual observations of cyclotron resonance absorption. Cyclotron resonance is a more explicit technique which enables carriers of different mass to be determined by measuring different resonant frequencies. The Faraday effect is an easier and

more flexible technique, but is less explicit, since the dispersion involves the integral of all CR absorption. Also, it is unable to detect any anisotropy of the effective mass in a cubic crystal.

Other free-carrier effects are ellipticity associated with the Faraday rotation, the Voigt effect, and the magnetoplasma effect on the reflectivity minimum.

Impurity magnetoabsorption Impurity states forming shallow levels (i.e., those levels separated from the nearest band by an energy much less than the gap E_g) can be described by the effective mass approximation. This simple model for impurities is that of the hydrogen atom with an electron which has an effective mass m^* and the nuclear charge reduced to e/ϵ_∞ , by the high-frequency dielectric constant of the crystal. This hydrogen-atom-like model leads to a series of energy levels leading up to a photoionization continuum commencing at an energy such that the electron (hole) is excited into the conduction (or valence) band.

The effect of a magnetic field on impurity levels and the related optical transitions is one of the most important tools for the study of the electronic states in semiconductors. The reason for this is that a magnetic field removes all degeneracies including the Kramers' degeneracy due to time reversal. It can produce new quantization rules, and it is a strong perturbation which is not screened like the Coulomb interactions.

For donor states with an isotropic conduction band minimum, as in GaAs or InSb, the hamiltonian acquires terms which are linear (Zeeman term) and quadratic (diamagnetic) in magnetic field. In this case, perturbation theory is appropriate to use until one considers transitions to very high quantum numbers or very high fields. Often complicated variational procedures must be used for accurate solutions. In the high field limit where $\hbar\omega_c/R \gg 1$ or for high energies, the impurity level quantum numbers do not provide a proper classification, and, instead, the states should be related to the continuum Landau levels. The effect of a magnetic field on acceptor impurity states is more complicated because of band degeneracy and of spin-orbit coupling. Also for acceptors, one can be in the high field limit at relatively small values of the field for shallow levels and for excited states.

High-resolution photoconductivity measurements in high-purity GaAs indicate that the behavior of shallow donors in GaAs deviates from that predicted by the simple hydrogenic model. Figure 48

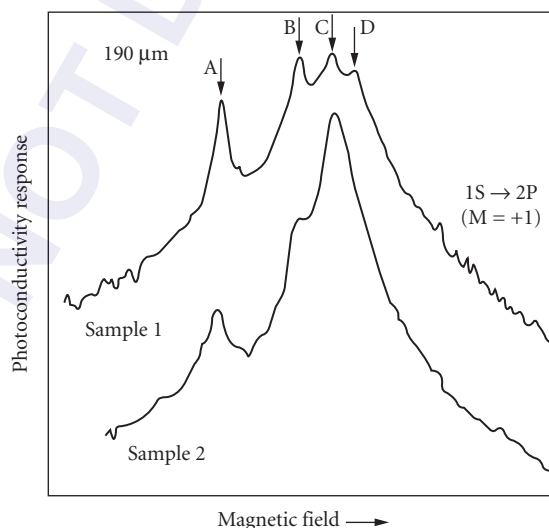


FIGURE 48 The photoconductivity due to the ($1s \rightarrow 2p$, $m = +1$) transition as a function of magnetic field for two different samples when excited by $190 \mu\text{m}$ laser radiation. For sample 1, $N_D = 2.0 \times 10^{14} \text{ cm}^{-3}$, $N_A = 4.0 \times 10^{13} \text{ cm}^{-3}$, and $\mu_{77\text{K}} = 153,000 \text{ cm}^2/\text{V}\cdot\text{s}$. For sample 2, $N_D = 4.3 \times 10^{13} \text{ cm}^{-3}$, $N_A = 2.5 \times 10^{13} \text{ cm}^{-3}$, and $\mu_{77\text{K}} = 180,000 \text{ cm}^2/\text{V}\cdot\text{s}$.⁹³

shows that, under high resolution, the single ($1s \rightarrow 2p$) transition observed at lower resolution actually consists of several different unresolved transitions.⁹³ The photoconductive response is shown for two different high-purity GaAs samples. The magnetic field range shown covers about 1 kG at around 15 kG. The different transitions labeled A, B, C, and D each correspond to a different donor species, and variation of the amplitude of these transitions in different samples results from the different relative concentrations of the particular donor species. Thus, the use of magnetic fields in studying the absorption properties of shallow impurities provides an analytical method of characterizing impurities. Because of the central-cell corrections (chemical shifts), different impurities of the same type (donor or acceptors) can be distinguished and identified by back-doping experiments.

Semiconductor Nanostructures: Low Dimensional Systems Development of growth techniques such as molecular beam epitaxy and metal organic chemical vapor deposition have made possible the precise growth of high-quality, layered, semiconductor heterostructures. Various sequences of thin film layers can be grown to form quantum wells, wires, and dots as well as superlattices and heterojunctions. In these low dimensional systems, the quantization effect of the magnetic field is significantly enhanced, leading to a wealth of new ways to characterize these structures. We refer the readers to recent books that review these systems.^{94,95}

Nonlinear Optical Properties of Semiconductors

Background The study and characterization of nonlinear optical properties of semiconductors are increasingly important topics for research and development, with new effects being discovered and useful nonlinear optical devices being constructed. The underlying cause of these nonlinear optical effects lies in the interaction of electromagnetic radiation with matter. Each nonlinear optical process may be thought to consist of the intense light first inducing a nonlinear response in a medium and then the medium in reacting, modifying the optical fields in a nonlinear way. There is a wide variety of nonlinear optical phenomena in semiconductors, leading to many papers on the subject of nonlinear properties. It is thus not possible here to do justice to this field. We refer the reader to Chap. 10, "Nonlinear Optics."⁹⁶ Here, we give a brief overview of the nonlinear optical processes starting from Maxwell's equations and describe and categorize some important second- and third-order nonlinear processes and properties.

Theoretical Overview of Nonlinear Optical Processes and Properties

Maxwell's equations and polarization power series expansion All electromagnetic phenomena are governed by Maxwell's equations for the electric and magnetic fields $\mathbf{E}(\mathbf{r}, t)$ and $\mathbf{B}(\mathbf{r}, t)$ or by the resulting wave equations

$$\left[\nabla \cdot \nabla + \frac{1}{c^2} \frac{\partial^2}{\partial t^2} \right] \mathbf{E}(\mathbf{r}, t) = \frac{-4\pi}{c^2} \frac{\partial^2}{\partial t^2} \mathbf{P}(\mathbf{r}, t) \quad (56)$$

$$\nabla \cdot \mathbf{E}(\mathbf{r}, t) = -4\pi \nabla \cdot \mathbf{P}(\mathbf{r}, t)$$

Here, \mathbf{P} is the generalized electric polarization which includes not only the electric dipole part but all the multiple contributions. In general, \mathbf{P} is a function of \mathbf{E} which describes fully the response of the medium to the field. It is often known as the constitutive equation since all optical phenomena would be predictable and easily understood from it and its solution for the resulting set of Maxwell's equations with appropriate boundary conditions. Unfortunately, this equation is almost never possible to solve exactly, and physically reasonable approximations must be resorted to for progress to occur.

Most nonlinear optical properties can be described in terms of a power series expansion for the induced polarization. (This assumes that \mathbf{E} is sufficiently weak.) Since lasers are most often used to observe nonlinear optical effects, one usually deals with the interaction of several monochromatic or quasi-monochromatic field components, and \mathbf{E} and \mathbf{P} can be expanded into their Fourier components as

$$\mathbf{E}(\mathbf{r}, t) = \sum_i \mathbf{E}(\mathbf{q}_i, \omega_i), \quad \mathbf{P}(\mathbf{r}, t) = \sum_i \mathbf{P}(\mathbf{q}_i, \omega_i) \quad (57)$$

where

$$\mathbf{E}(\mathbf{q}_i, \omega_i) = \mathbf{E}(\omega_i) \exp(i\mathbf{q}_i \cdot \mathbf{r} - i\omega_i t) + \text{c.c.} \quad (58)$$

The induced polarization is usually written as

$$\begin{aligned} \mathbf{P}(\mathbf{q}_i, \omega_i) = & \chi^{(1)}(\mathbf{q}_i, \omega_i) \cdot \mathbf{E}(\mathbf{q}_i, \omega_i) + \sum_{j,k} \chi^{(2)}(\mathbf{q}_i = \mathbf{q}_j + \mathbf{q}_k, \omega_i = \omega_j + \omega_k) : \mathbf{E}(\mathbf{q}_j, \omega_j) \mathbf{E}(\mathbf{q}_k, \omega_k) \\ & + \sum_{j,k,l} \chi^{(3)}(\mathbf{q}_i = \mathbf{q}_j + \mathbf{q}_k + \mathbf{q}_l, \omega_i = \omega_j + \omega_k + \omega_l) : \mathbf{E}(\mathbf{q}_j, \omega_j) \mathbf{E}(\mathbf{q}_k, \omega_k) \mathbf{E}(\mathbf{q}_l, \omega_l) + \dots \end{aligned} \quad (59)$$

It is, however, sometimes more convenient to use $\mathbf{E}(\mathbf{r}, t)$ and $\mathbf{P}(\mathbf{r}, t)$ directly instead of their Fourier components, especially when dealing with transient nonlinear phenomena.

In the electric dipole approximation, $\chi^{(n)}(\mathbf{r}, t)$ is independent of \mathbf{r} , or $\chi^{(n)}(\mathbf{k}, \omega)$ is independent of \mathbf{q} , and the equations become simpler to write and to work with. These $\chi^{(n)}$ are the susceptibilities, with $\chi^{(1)}$ = linear electric dipole susceptibility and $\chi^{(2)}$ ($\chi^{(3)}$) = nonlinear second-order (third-order) susceptibility tensor. Both absorptive and refractive effects can be described in terms of these complex electric susceptibilities, which have real and imaginary parts for each tensor element. These linear and nonlinear susceptibilities characterize the optical properties of the medium and are related to the microscopic structure of the medium. Knowledge of $\chi^{(n)}$ allows, at least in principle, to predict the n th-order nonlinear optical effects from Maxwell's equations. Consequently, much effort (both experimentally and theoretically) has gone into determining the $\chi^{(n)}$.

The definitions of the nonlinear susceptibilities in the literature vary and have led to some confusion. Shen reviews these definitions and the reasons for the confusion.⁹⁷ In addition to some intrinsic symmetries, the susceptibilities must obey crystallographic symmetry requirements. The spatial symmetry of the nonlinear medium imposes restrictions upon the form of the various $\chi^{(n)}$ tensors. Butcher has determined the structure of the second- and third-order tensors for all crystals.⁹⁸ One important consequence is that, for media with inversion symmetry, $\chi^{(2)} \equiv 0$, and thus $\chi^{(3)}$ represents the lowest-order nonlinearity in the electric-dipole approximation. Of the 12 nonzero elements, only 3 are independent. These susceptibility tensors must transform into themselves under the point group symmetry operations of the medium.

It is often convenient to discuss the various optical processes which might occur in terms of whether they are active or passive. Passive processes involve energy or frequency conservation, and the material medium acts basically as a catalyst. The susceptibilities are predominantly real for passive processes. Of course, as resonances are approached, susceptibilities become complex and may even become totally imaginary. These passive nonlinear optical phenomena are listed in Table 7.⁹⁹ Active nonlinear optical phenomena are listed in Table 8.⁹⁹ In general, energy is exchanged between the radiation and the material only for the active processes. We also note that second-order effects are always passive.

Second-order nonlinear optical properties Most existing nonlinear optical devices are based upon second-order nonlinear optical effects that are quite well understood. Here, we assume the presence of only three quasi-monochromatic fields

$$\mathbf{E} = \mathbf{E}(\omega_1) + \mathbf{E}(\omega_2) + \mathbf{E}(\omega_3) \quad (60)$$

TABLE 7 Passive Nonlinear Optical Phenomena⁹⁹

Frequencies of Incident Fields	Frequencies of Fields Generated by the Polarization of the Medium	Susceptibility	Process (Acronym)
ω_1	No polarization	$0 (\epsilon = 1)$	Vacuum propagation (VP)
ω_1	ω_1	$\chi^{(1)}(\omega_1; \omega_1)$	Linear dispersion (LD)
ω_1, ω_2	$\omega_3 [\omega_3 = \omega_1 + \omega_2]$	$\chi^{(2)}(\omega_3; \omega_1, \omega_2)$	Sum mixing (SM)
ω_1	$\omega_3 [\omega_3 = 2\omega_1]$	$\chi^{(2)}(\omega_3; \omega_1, \omega_1)$	Second-harmonic generation (SHG)
$\omega_1, 0$	ω_1	$\chi^{(2)}(\omega_1; \omega_1, 0)$	Electro-optic linear Kerr effect (EOLKE)
ω_1	$\omega_2, \omega_3 [\omega_1 = \omega_2 + \omega_3]$	$\chi^{(2)}(\omega_2; -\omega_3, \omega_1)$	Difference-frequency mixing (DFM)
ω_1	$\omega_2 [\omega_1 = 2\omega_2]$	$\chi^{(2)}(\omega_2; -\omega_2, \omega_1)$	Degenerate difference-frequency (DDF)
ω_1	0	$\chi^{(2)}(0; -\omega_1, \omega_1)$	Inverse electro-optic effect (IEOE)
$\omega_1, \omega_2, \omega_3$	$\omega_4 [\omega_4 = \omega_1 + \omega_2 + \omega_3]$	$\chi^{(3)}(\omega_4; \omega_1, \omega_1, \omega_1)$	Third-harmonic generation (THG)
ω_1, ω_2	$\omega_3, \omega_4 [\omega_1 + \omega_2 = \omega_3 + \omega_4]$	$\chi^{(3)}(\omega_3; -\omega_4, \omega_1, \omega_2)$ $\chi^{(3)}(\omega_4; -\omega_3, \omega_1, \omega_2)$	Four-wave difference-frequency mixing processes (FWDFMP)
ω_1	$\omega_2, \omega_3, \omega_4 [\omega_1 = \omega_2 + \omega_3 + \omega_4]$	$\chi^{(3)}(\omega_2; -\omega_3, -\omega_4, \omega_1)$	Intensity-dependent refractive index (IDRI)
ω_1	ω_1	$\chi^{(3)}(\omega_1; \omega_1, -\omega_1, \omega_1)$	
$\omega_1, 0$	ω_1	$\chi^{(3)}(\omega_1; 0, 0, \omega_1)$	Quadratic Kerr effect (QKE)

$\omega = 0$ indicates the presence of a uniform electric field.

and

$$\omega_1 = |\omega_2 \mp \omega_3| \quad (61)$$

Thus, Eq. (56) can be decomposed into three sets of equations for each $\mathbf{E}(\omega_i)$. They are then nonlinearly coupled with one another through the polarizations

$$\mathbf{P}(\omega_1) = \chi^{(1)}(\omega_1) \cdot \mathbf{E}(\omega_1) + \chi^{(2)}(\omega_1 = |\omega_2 \pm \omega_3|) : \mathbf{E}(\omega_2)\mathbf{E}(\omega_3) \quad (62)$$

The second-order nonlinear processes are then described by the solutions of the coupled-wave equations with the proper boundary conditions. $\chi^{(2)} = 0$ for materials with a center of inversion. The coefficient $\chi^{(2)}$ is a third-rank tensor. Some second-order processes include sum- and difference-frequency mixing, the electro-optic linear Kerr effect, the inverse electro-optic effect, parametric amplification and oscillation, and second-harmonic generation. The past emphasis has been

TABLE 8 Active Nonlinear Optical Phenomena⁹⁹

Susceptibility	Process
$\chi^{(1)}(\omega_1; \omega_1)$	Linear absorption ($\omega_1 \approx \omega_{10}$)
$\chi^{(3)}(\omega_2; \omega_1, -\omega_1, \omega_2)$	Raman scattering ($\omega_2 \approx \omega_1 \mp \omega_{10}$)
$\chi^{(3)}(\omega_1; \omega_1, -\omega_1, \omega_1)$	Two-photon absorption ($2\omega_1 \approx \omega_{10}$) or Saturable absorption ($\omega_1 \approx \omega_{10}$)
$\chi^{(5)}(\omega_2; \omega_1, \omega_1, -\omega_1, -\omega_1, \omega_2)$	Hyper-Raman scattering ($\omega_2 \approx 2\omega_1 \mp \omega_{10}$)

on finding new nonlinear crystals with a large $\chi^{(2)}$. Semiconductor crystals have received much attention: III-V compounds like GaAs and InSb, II-VI compounds like ZnS and CdSe, I-III-VI compounds like AgGaS₂ and CuInS₂, and II-IV-V compounds like CdSiAs₂ and ZnGeP₂.

In most applications of second-order nonlinear optical effects, it is important to achieve phase-matching conditions

$$\Delta\mathbf{q}=\mathbf{q}_1-\mathbf{q}_2-\mathbf{q}_3=0 \quad (63)$$

where \mathbf{q}_i is the wave vector of $\mathbf{E}(\omega_i)$. This ensures an efficient energy conversion between the pump field(s) and the signal field.

The intensity of the electric field of a laser can produce photons at multiples of the frequency of the probing light signal (ω_0). Second harmonic generation has been widely used to characterize interfaces between silicon and thin dielectric films as well as thin metal films.¹⁰⁰ There are two sources of second harmonic light at the interface between dielectric films and centro-symmetric crystals such as silicon. One is due to the weak electric quadrupole contribution, and the second is due to a surface dipole contribution.^{101,102} The dipole contribution is stronger than the second harmonic signal from the quadrupole one. The detailed crystalline symmetry along the probe beam direction results in a rotational dependence to the quadrupole signal. The surface dipole contribution is strongly affected by chemical changes in the interface making it an ideal means of characterizing subtle process-induced changes in the dielectric film stack used for transistor applications.¹⁰¹ The second harmonic intensity, I , from a p - or s -polarized probe and a p - or s -polarized SHG signal for the Si(100) crystal face are

$$I_{pp}(2\omega)=|a_{pp}^{(0)}+a_{pp}^{(4)}\cos(4\Phi)|^2$$

$$I_{ps}(2\omega)=|a_{ps}^{(4)}\sin(4\Phi)|^2$$

$$I_{sp}(2\omega)=|a_{sp}^{(0)}+a_{sp}^{(4)}\cos(4\Phi)|^2$$

$$I_{ss}(2\omega)=|a_{ss}^{(4)}\sin(4\Phi)|^2$$

Here, the a coefficients are due to either a surface dipole contribution $a^{(0)}$ or a bulk quadrupole $a^{(4)}$ contribution. Φ is the angle between the $[0\bar{1}0]$ crystal axis and the projection of the wave vector of the second harmonic light on the surface. The SHG signal is also strongly influenced by electric fields present in the sample. Changes in the electric field due to trapped charge in hafnium oxide have been characterized using the changes in $a^{(0)}$ due to annealing temperature.¹⁰¹ Both the phase and amplitude of the second harmonic signal from Si(100) surfaces have been measured by SHG.¹⁰²

Third-order nonlinear optical properties In materials with inversion symmetry, third-order processes are the dominant nonlinearity. These processes are described by a fourth-rank nonlinear susceptibility tensor $\chi^{(3)}$ whose contribution to the polarization is given according to Eq. (59) by

$$\mathbf{P}(\omega_i)=\sum_{j,k,l}\omega^{(3)}(\omega_i=\omega_j+\omega_k+\omega_l):\mathbf{E}(\omega_j)\mathbf{E}(\omega_k)\mathbf{E}(\omega_l) \quad (64)$$

In general, this nonlinearity will provide a coupling between four electromagnetic waves. Depending on whether the susceptibility tensor elements are real or imaginary and on whether some of the frequencies are identical or different, a large variety of physical phenomena can be understood and accounted for: third-harmonic generation, two-photon absorption, saturable absorption, intensity-dependent index of refraction, stimulated Raman effect, anti-Stokes generation, stimulated Raleigh scattering, modulation of the index of refraction, and self-focusing of light. We concentrate on discussing only a few of the most pertinent cases of interest.

Third-harmonic generation Here, the output frequency $\omega_4 = \omega_1 + \omega_2 + \omega_3 = 3\omega_1$, since $\omega_1 = \omega_2 = \omega_3$. The polarization at the frequency ω_4 will generate radiation at the third-harmonic frequency. The quantum process responsible for the harmonic generation may be described as a scattering process in which three quanta at the fundamental frequency are annihilated and one quantum at the third-harmonic frequency is created. The system remains in the ground state, although three virtually excited states are involved in the scattering process. Since the phases are important, the process is actually an interference between many four-photon scattering processes.

Two-photon absorption Here, for example, $\omega_1 = -\omega_2 = \omega_3 = \omega_4$ and $\Delta\mathbf{k} = 0$, and the nonlinear polarization has components described by

$$\mathbf{P}(\omega_4) = \chi^{(3)}(\omega_4 = +\omega_1 - \omega_2 + \omega_3) : \mathbf{E}(\omega_1)\mathbf{E}^*(\omega_2)\mathbf{E}(\omega_3) \quad (65)$$

The nonlinear susceptibility $\chi^{(3)}$ is purely imaginary, but positive. One can define an absorption coefficient proportional to the intensity itself. In the important case of resonance, the sum of two frequencies of the exciting field is approximately equal to a transition frequency of the medium, $\omega_{ab} = 2\omega_1$, where $\hbar\omega_{ab}$ is the energy difference between two levels $|a\rangle$ and $|b\rangle$ with the same parity.

The TPA coefficient can be expressed in terms of a third-order nonlinear susceptibility tensor by solving the wave equation using the slowly varying amplitude approximation. The explicit expression for β is related to the imaginary part of the third-order electric dipole susceptibility tensor which depends upon the crystal class and laser electric field direction, for example, in crystals with 43 m symmetry (e.g., Hg_{1-x}Cd_xTe, GaAs, InSb, etc.) and the electric field along the [001] direction,¹⁰³

$$\beta = \frac{32\pi^2\omega}{n^2c^2} \left[3\text{Im} \chi_{1111}^{(3)}(-\omega, \omega, \omega, -\omega) \right] \quad (66)$$

where the convention used is that of Maker and Terhune (1965).¹⁰⁴

Since $\chi^{(3)}$ is a second-rank tensor, there are, in general, nine terms contributing to β with magnitudes which vary with orientation. However, in most systems the symmetry is such that there are relations between some of the terms. For example, in crystals with 43 m symmetry, there are three possible values of χ for each β , and for crystals with m3m symmetry, there are four possible values. These can, in general, be measured by using both linearly and circularly polarized light. It is not always possible to sort out all the different TPA spectra simply by changing the sample orientation or the light polarization because of the competing process of absorption by second-harmonic-produced light.

5.4 MEASUREMENT TECHNIQUES

Overview

The ability to measure the optical response of a semiconductor specimen precisely under well-controlled environmental conditions is of obvious importance in the determination of the optical properties of semiconductors. The significant resolution with which optical spectra can be measured makes it possible to perform precise determinations of the intrinsic properties such as the energy separation between electronic states, lattice vibration frequencies, as well as extrinsic properties due to impurities. The wealth of information that can be obtained and the direct relevance to device-related issues has led to much effort being expended in developing techniques and apparatus to perform a wide range of measurements. This chapter reviews the most widely used procedures and optical components employed in these investigations.

Optical studies, except in rare instances, are contactless and noninvasive. These attractive features have led to widespread use of the techniques in both scientific analyses and, more recently,

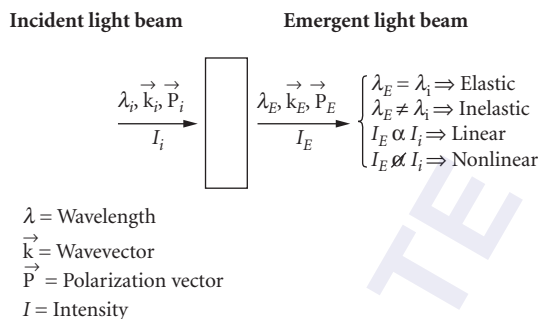


FIGURE 49 Schematic of interaction of light with semiconductors showing the linear, nonlinear, elastic, and inelastic processes.

in manufacturing environments. Specific procedures and experimental apparatus and variations of them are too numerous to be dealt with in this brief chapter. The motivation here is, therefore, to set forth the essentials and dwell on the major aspects of each technique covered, as well as to provide some references which contain more details than could be given here.

The essence of spectroscopic investigations is to determine the interaction of a light beam, with a well-defined wavelength, intensity, polarization, and direction, with a semiconductor specimen, in most cases from the point of view of the light beam. Figure 49 schematically displays this light-specimen interaction. Upon interacting, both the specimen and the light beam will change, and the experimental task is to precisely measure the change in the properties of the light beam. The changes may be classified into the linear and nonlinear regimes based on whether or not the response of the specimen is linear with respect to the incident power. The interactions may be elastic or inelastic; the light beam may also undergo a change in all aspects except its wavelength or photon energy in the former, and the wavelength can also be modified in the latter. These terms arise from the elastic or inelastic interaction of the photon with the specimen where the incident photon energy is preserved or modified in the process. The incident light beam may be reflected, scattered, and transmitted by the specimen. In addition to these processes, a properly excited specimen may emit light as well. The last process, known as luminescence, may also be exploited to gain an insight into the physical behavior of the material.

The linear, elastic regime covers most of the procedures used to elucidate the equilibrium properties related to the optical constants, n and k , introduced previously. The techniques used are comprised of reflection and transmission spectroscopies where the energy reflection R and transmission T of the specimen are studied as function of the wavelength, polarization, and angle of incidence. The net energy absorption A may be determined from R and T as follows: $A = (1 - R - T)$.

As discussed earlier, a complete knowledge of R or T over a large range of energies is required to determine both n and k . Such a task is usually difficult. A more convenient and accurate procedure is to measure the change in polarization properties of an obliquely incident plane-polarized light beam after it interacts with the specimen. This procedure, called ellipsometry, uses the polarization change of the incident beam in the material to extract n and k as well as the thickness of each layer of a multilayer structure.^{105,106} The underlying principles and a more detailed discussion follow.

The nonlinear and inelastic spectroscopic procedures have been popular since the advent of lasers. The very large power densities achievable using lasers over the wide range of energies extending from the far-infrared (FIR) to the UV have driven the rapid developments in this field. At high excitation powers, the absorption of light energy can become superlinear due to the presence of higher-order interactions as discussed previously. Exploitation of these specific interactions as a means of gaining information regarding the specimen is the content of nonlinear spectroscopic techniques. On the other hand, the high degree of wavelength purity and coherence offered by lasers is exploited to perform inelastic spectroscopic analyses such as Raman¹⁰⁷ and Brillouin¹⁰⁷ spectroscopies. The crux of these techniques is to project a high-power highly coherent laser beam onto a

specimen and observe the scattered part of the intensity. The scattered part will be dominated by light with the same wavelength as the incident laser beam; that is, the elastic part, but a small part, usually $<10^{-8}$ of the original intensity, can be observed with well-defined frequency shifts. These additional frequency bands, similar to the sidebands that arise as a consequence of intensity modulation, can be analyzed to provide crucial information regarding the specimen. For instance, optical phonons can interact with the incoming laser photon and energy-shift it by an amount equal to a multiple of the phonon energy. Hence, an analysis of the frequency-shifted bands in the Raman spectrum can be used to establish phonon energies. The major task of Raman spectroscopy is to isolate the very weak frequency-shifted component in the scattered beam.¹⁰⁷

Light emission is important for both spectroscopic analysis and device applications and, hence, has commanded a large amount of attention.^{108,109} Any excited semiconductor will emit light as a means of relaxing to its equilibrium state. Under proper excitation, such as with above-bandgap radiation, the light emission can be made quite intense and can then be easily recorded and subjected to spectroscopic analyses. Within the specimen, the above-bandgap photoexcitation leads to a transition of the electron from a valence band to the conduction band, followed by a rapid process of thermalization whereby the excited electron and hole reach their respective band extrema and recombine from there radiatively; that is, by emitting the potential energy in the form of a photon. The photoexcited free electrons and holes may also form an exciton, or interact with the impurity states in the forbidden band, or both. The final recombination can be mediated by a large number of such intermediate processes, and, hence, the luminescence spectra can display a very rich and complicated structure. The most important aspect of the luminescence spectrum is the fact that nearly all the interactions involve impurity and defect states in the forbidden band. Add to this the rapid thermalization and large self-absorption effects for emission of light with energies greater than the gap, and luminescence is almost entirely a sub-bandgap tool dominated by the impurities and defects present in the semiconductor.

The major categories of spectroscopic procedures are grouped by the energy range of photons used. The commonality of the instrumentation for a given wavelength range accounts for this categorization. The lowest energy region in the FIR spans the phonon energies and a variety of other possible excitations such as plasmons (free-carrier oscillations), magnons, and impurity-related electronic and vibrational excitations. The energy-bandgap of semiconductors ranges from 0 for HgTe to >5 eV¹¹⁰ in diamond and can hence occur anywhere from the far-infrared to the vacuum ultraviolet. For two of the most important electronic materials, namely Si and GaAs, the gaps occur in the near-infrared at 0.8 and 1.5 eV,¹¹⁰ respectively. Hence, the mid- and near-infrared investigations have largely been confined to the study of impurity states within the forbidden gap. The electronic band transitions dominate the higher energies.

Instrumentation

The major components of a spectroscopic system are schematically displayed in Fig. 50. Light from the broadband source enters a monochromator fitted with a dispersion element that separates the various wavelength components and allows a chosen narrow spectral band of light to interact with the specimen. A detector converts the intensity information in the beam to an electrical or digital signal that can then be recorded by a computer or other recording device and then analyzed. Passive components such as lenses and mirrors, filters, polarizers, light pipes, and optical fibers, etc., that are needed to tailor the behavior of the system also form an integral part of the experimental apparatus. A short discussion of each of the major components follows.

Sources

Broadband The ideal broadband source should emit light with sufficient intensity in the wavelength band of interest, possess a stable output and exhibit a minimum amount of noise, and display a slowly varying spectral character; that is, the source should not possess intense spectral features that will interfere with the measurement procedure.¹¹¹ All these characteristics can be

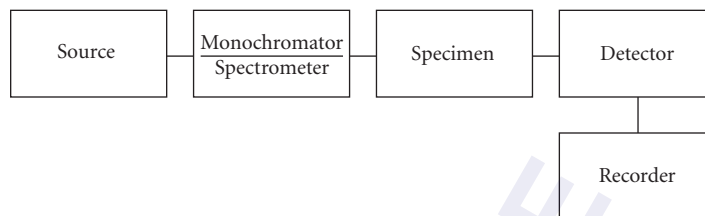


FIGURE 50 Schematic of a spectroscopic measurement apparatus.

satisfied by blackbodies, and, if blackbodies with high enough temperatures can be fabricated, they would be ideal sources for any wavelength region. However, this is not possible since the operating temperatures required to obtain workable energy densities in the ultraviolet are extremely large. Hence, blackbody sources are usually restricted to wavelengths in the red-yellow region starting at ≈ 500 nm or larger and ending at ≈ 2000 μm in the FIR. The incandescent lamp with a hot filament is the best known BB source. Gas emission lamps such as high-pressure arc lamps and low-pressure discharge lamps are useful in the visible and ultraviolet region. Their main feature is the ability to produce a large intensity in the upper energy regions. However, the inherent atomic processes and associated line spectra that are present make using these somewhat complicated. Care should be exercised to avoid wavelength regions where intense spectroscopic features arise from the discharge medium. The unavoidable electronic activity in the discharge media can also be a source of noise.

Laser Laser sources are required in applications where large intensities are essential, as in the studies of nonlinear optical phenomena using Raman scattering or photoluminescence techniques. Lasers are currently available from the UV to the FIR, and often both cw and pulsed operations are possible.¹¹² The argon and krypton ion lasers with emissions in the visible and near-infrared regions, the Nd:YAG with a 1.06 μm emission, and the CO_2 laser with emission in the 9.2 to 10.8 μm range have been the workhorses for a wide variety of semiconductor investigations.

Since efficient laser operation requires a set of excitable electronic or vibrational levels properly arranged to produce population inversion and sufficient amplification, intense laser emission is usually confined to specific wavelengths. However, the use of optically excited dye lasers and tunable solid-state lasers such as the Ti:sapphire lasers can be used to fill the wavelength regions in between most of the visible and near-infrared regions.

The semiconductor lasers that are currently available extend in wavelength from the red to the far-infrared. The III-V alloy-based double heterostructure lasers, fabricated from $\text{Al}_x\text{Ga}_{1-x}\text{As}$ and $\text{In}_{1-x}\text{Ga}_x\text{As}$, are particularly efficient in the near-infrared region. The IV-VI alloy-based lasers, fabricated from $\text{Pb}_{1-x}\text{Sn}_x\text{Te}$, operate at considerably longer wavelengths of ≈ 10 μm ; a small range of emission wavelength tunability has been achieved based on the change of the band gap with the temperature. The intense interest in the development of blue-green laser emission is now having substantial influence.

Spectrometers and Monochromators

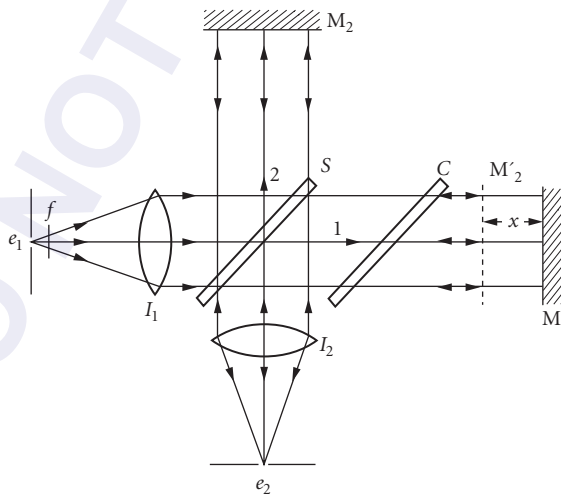
Dispersion spectrometers The monochromator or the spectrometer is the heart of the spectroscopic apparatus, and the dispersion element that analyzes the light is the main component of the monochromator.¹¹³ The simplest and best known dispersion unit is the transparent prism. The physical mechanism that leads to the dispersion is the inherent dispersion in the optical refractive index n of the prism material over the wavelength range of interest. Prism-based monochromators work very well and are still employed. They possess a large throughput; that is, they transmit a large fraction of the incident intensity but suffer from limited resolution since the degree of dispersion is restricted by the characteristics of the prism material.

The most widely used dispersion element is the diffraction grating which is a collection of finely spaced grooves or slits. Diffraction from the multiple slits leads to a dispersive action. The dispersion is given by the following relationship: $n\lambda = d \sin(\theta)$.

The fact that the degree of dispersion can be controlled by the slit spacing reduces the complexity of design and fabrication when compared to the prism. Gratings have a drawback in that they display multiple orders, and hence the throughput at any given order is likely to be very high. Blazing (i.e., control of the shape of the groove) can be used to increase energy in a given band in a particular order to reduce this shortcoming. Both transmission and reflection gratings can be fabricated as well as concave gratings that both disperse and focus the light beam. The major manufacturing flaw in grating fabrication that used mechanical devices was that groove spacing was not well controlled and hence the flawed grating led to the appearance of spectroscopic artifacts that were called "ghosts." The advent of holographic grating fabrication procedures has eliminated these difficulties.

Fourier-transform spectrometers An alternate method for performing spectroscopic measurements is the use of Michelson's interferometer.¹¹⁴ Figure 51 displays the layout of the interferometer. The approach in this procedure is to divide the white-light beam from the source into two wavefronts, introduce a path difference x between the two, recombine them, and record the interference-modulated intensity as a function of x . The recorded intensity variation $B(x)$ is known as the interferogram, and a Fourier transformation of $I(x)$ will yield the spectral distribution of the white light. Experimentally, this is achieved as shown in Fig. 51. The incident beam from the source is collimated by lens I_1 , split into two wavefronts by the beam splitter S . The two wave fronts are directed to a movable mirror, M_1 , and a fixed mirror, M_2 . The reflected beams are combined at the detector where the intensity is recorded as a function of x , the path difference between M_1 and M_2 . The spectrum $B(\omega)$ is related to $I(x)$ by the following relationship:

$$B(\omega) = \int_0^{\infty} [I(x) - I(\infty)] \cos(2\pi\omega x) dx \quad (67)$$



$$B(\omega) = \int_0^{\infty} [I(x) - I(\infty)] \cos(2\pi\omega x) dx$$

FIGURE 51 Schematic of the Fourier-transform interferometer.

This procedure of measuring the spectra is referred to as Fourier-transform spectroscopy and contains two major advantages: the throughput advantage and the multiplex advantage, both of which greatly add to the ultimate signal-to-noise ratios that can be achieved as compared to measurements performed with grating spectrometers under comparable conditions of illumination and recording times. The throughput advantage arises from the fact that improved resolution is not achieved at the expense of reducing slits and reduced throughput, and the multiplex advantage is a consequence of the fact that all wavelength channels are observed all the time as opposed to a one-channel-at-a-time measurement constriction in the dispersion-based instruments.

These significant advantages come with a price. The much larger signal intensities are likely to be seen by the detector and place stringent conditions on the detector performance in its dynamic range and linearity. Less-than-optimum performance may lead to significant distortions in the transformed spectrum that are not intuitively evident. The FT spectrometer was used first in the far-infrared, soon after the advent of high-speed computers that were capable of performing the Fourier transformations. However, the advances in the technology of designing and fabricating complicated optical elements and computer hardware have contributed greatly to the advancement of the field, and FT spectrometers are now available that cover a wide spectral region extending from FIR to the VUV.

Detectors The photomultiplier (PMT)¹¹⁵ is a widely used light detector that is a vacuum-tube-based device that uses a photoemitter followed by a large ($>10^3$) amplification stage so that very low signal levels can be detected. The wavelength band that the detector responds to is determined by the photocathode and window characteristics. Photomultipliers are particularly useful when low-light-level detection is needed as in Raman spectroscopy, but their use is largely confined to the visible and near-infrared as a consequence of the limitation in obtaining photoemitters for lower energies. Use of fluorescent phosphors can extend the upper working region of the PMTs to the VUV and beyond. Since they are vacuum-based devices, they are fragile and have to be handled with care.

Solid-state detectors¹¹⁵ which are almost entirely fabricated from semiconductors have advanced to a state where, in many applications, they are preferable to PMTs. The simplest semiconductor detector is the photoconductor (PC), where absorption of an above-band-gap photon leads to an increase in conductivity. Semiconducting PC detectors are, therefore, sensitive to any radiation with an energy larger than the bandgap. Since the bandgap of semiconductors extends all the way from 0 in HgTe to >5 eV in diamond,¹¹⁰ detectors that function over a very large energy range can be fabricated. Photovoltaic detectors, as the name implies, employ the photovoltaic effect in a p-n junction and can also detect above-bandgap radiation.

The explosive growth in semiconductor technology that has led to large-scale integration has benefited spectroscopic experimenters directly in the field of detectors. Imaging devices such as the CCD (charge-coupled device) array have been incorporated in spectroscopy. The array detectors combined with a dispersion spectrometer can be used for observing multiple channels simultaneously. This has led to an advantage similar to the multiplex advantage in the Fourier spectrometer. In addition, the low noise levels present in these detectors, particularly when they are cooled, have had a large impact on high-sensitivity spectroscopic analysis. Commercially available CCDs are fabricated and, hence, have a lower energy limit of 0.8 eV, the bandgap of Si. However, some linear and 2D arrays, fabricated from InSb and HgCdTe, with longer wavelength response are also available.

Major Optical Techniques

Advances in semiconductor physics have relied on measuring fundamental material and device properties, measuring the quality of the material, and accurately determining the details of thin films, quantum wells, and other microstructures that control or affect device performance. Properties that need to be determined therefore include basic band structure properties such as energy gaps, the presence and concentration of impurities and defects, alloy parameters, layer thicknesses, homogeneity, and uniformity. A very practical review, "Optical Characterization in Microelectronics Manufacturing" describes in detail six techniques: ellipsometry, infrared spectroscopy, microscopy, modulation spectroscopy, photoluminescence, and Raman scattering.¹¹⁶ The

discussion of each technique indicates the basic semiconductor quantities measured, gives the scientific basis of the technique, and indicates how the measurement is made. Illustrative examples from the literature are discussed in detail, showing applications to important semiconductor materials.

Much of the basic physics concerning the optical functions of silicon and related materials has been known since the 1950s. Jellison has reviewed the understanding of the physics of optical measurements in the interest of providing the background required for understanding the many optical techniques.¹¹⁷ The practical value of spectroscopic ellipsometry in gate dielectric metrology for the semiconductor industry is well known and has been documented by Clive Hayzelden.¹¹⁸ Finally, we point out that there are important applications in real-time monitoring of film growth, etching, and surface modification using spectroscopic ellipsometry measurements with monolayer precision.¹¹⁹

Reflection and Transmission/Absorption Measurements of the power reflection, transmission, and absorption, R , T , and A , respectively, are the simplest and most direct methods of spectroscopic analyses of semiconductor materials. The measurements are simple to perform so long as satisfactory spectroscopic apparatus is available in the wavelength region of interest. The major drawback is the less-than-satisfactory absolute accuracies with which the measurements can be performed.

The measurements are usually conducted at near-normal incidence for convenience; normal incidence measurements are difficult to perform and oblique incidence measurements are difficult to analyze. Once the major elements of the spectroscopic system—namely, the source, monochromator, and detector—have been chosen, the R measurements are obtained by directing the light beam on the specimen and measuring the incident and reflected intensities with the detector. The experimental R is determined by ratioing the incident and reflected power, I_0 , and I_R , respectively.

$$R = \frac{I_R}{I_0} \quad (68)$$

The spectral behavior of the measurement apparatus will not affect the final result so long as they remain fixed during the observation of I_0 and I_R . A change in the source intensity or the response of the detector will be reflected in the measurement. Since the measurement conditions cannot be identical in the two measurements, this is an unavoidable source of error.

Several arrangements have been attempted to minimize such an error. The most direct approach is to hold the optical path fixed and instead of measuring I_0 and I_R ; one measures I_{REF} , and I_R , where I_{REF} is the reflected intensity from a well-calibrated reference surface. A ratio of these two measurements leads to the following relationship:

$$\frac{I_R}{I_{\text{REF}}} = \frac{R}{R_{\text{REF}}} \quad (69)$$

A knowledge of R_{REF} , the known reflectivity of the reference surface, may then be used to extract R . This procedure works well and has been employed widely in the infrared region, where reference surfaces from metal mirrors with ≈ 100 percent reflectivity and stable blackbody sources are available. For the visible and higher-energy regions that use gas lamp sources that tend to be not as stable as blackbody sources, the reference method is not satisfactory. The procedure used to overcome the short-term variations in the lamp intensity was to perform the ratioing at each measurement wavelength. Such a “real-time” ratioing arrangement used a dynamical arrangement with a rotating optical element that directed the light beam alternately to the reference and the sample. Many such configurations were used, and a good example of a particularly ingenious arrangement is presented in Fig. 52. The essence of this rotating light pipe reflectometer¹²⁰ was the rotating light pipe that sampled the incident and reflected beam alternately at a frequency of ≈ 70 Hz, a frequency large enough to remove any errors due to variations in the lamp intensity.

The use of reflectometers for the spectroscopic analysis of semiconductors in the visible and near-visible regions has been limited since the advent of the spectroscopic ellipsometer. However, they are still used at higher energies and in applications where their simplicity makes them attractive, for example, to monitor thin films in semiconductor device fabrication.

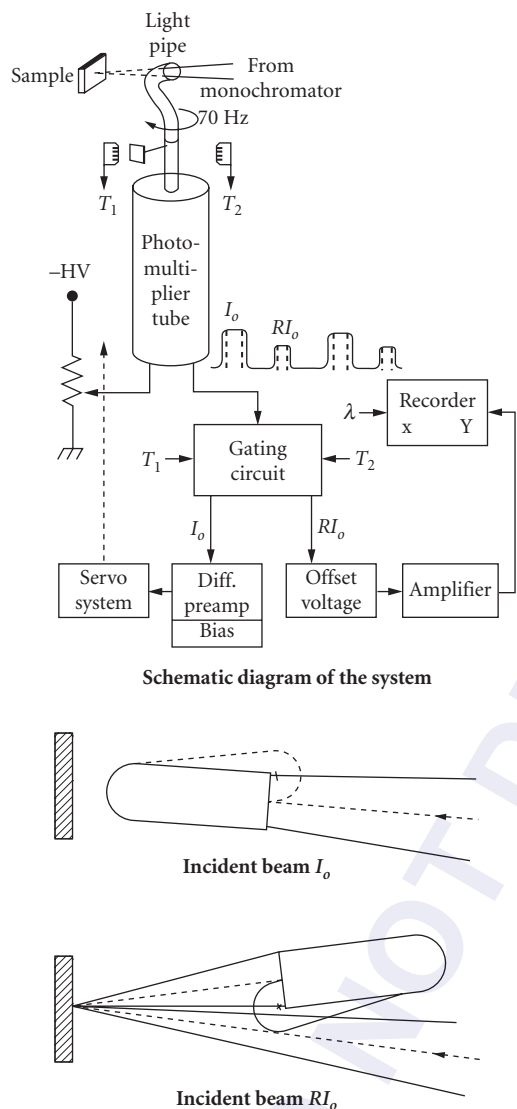


FIGURE 52 Schematic diagram of the rotating light pipe reflectometer and expanded top view showing the geometry of the sample and bent light pipe.¹²⁰

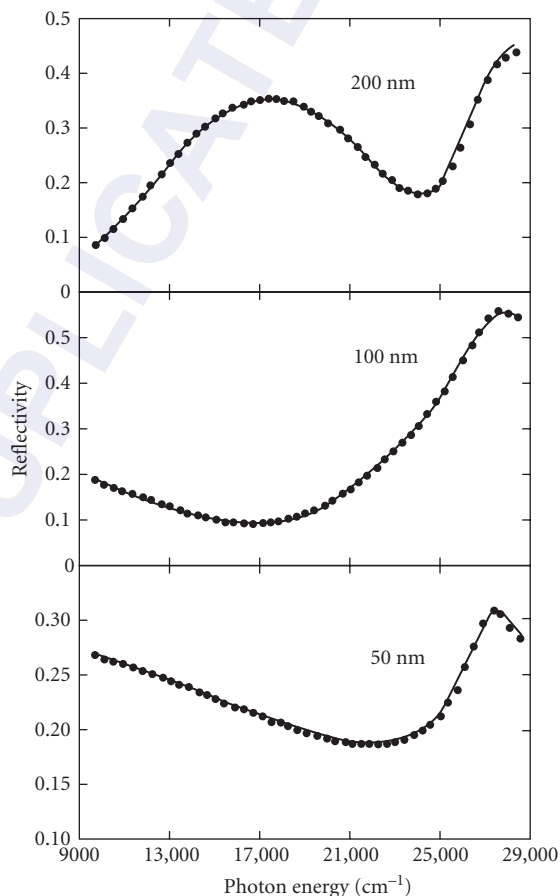


FIGURE 53 Measured (points) and calculated (line) optical reflectivity spectra of silicon dioxide films on silicon.¹²¹

Recent advances in optics have pushed the energy range of commercial reflectometers to just beyond 10 eV (120 nm) and ellipsometers to just beyond 8 eV (150 nm). Reflectometers and ellipsometers are both currently used for measurement of dielectric film thickness during the manufacture of integrated circuits and other semiconducting multilayer film measurements. The reflectivity spectrum obtained from three SiO_2 films on Si is presented in Fig. 53.¹²¹ The results of the computer fit are also displayed. The thicknesses determined from the reflection analyses agree to within ≈ 0.5 percent of the values obtained from ellipsometric values. As illustrated in the figure, the change in the spectral shape from the 200-nm-thick film to 100 nm is much greater than that between the

100- and the 50-nm films which points to reducing sensitivity to film thickness determination as the films get thinner. The ultra-precise measurement of film thickness for dielectric films less than 2 nm in thickness requires use of ellipsometry.

The techniques used to measure the transmission spectra are almost identical to those used for reflection except that nature provides a perfect reference, namely, the absence of the sample in the beam. The sample-in/sample-out reference intensity ratioing works satisfactorily. Since transmission measurements are mainly confined to energies below the forbidden gap and hence lower energies covered by blackbody sources, the difficulties faced by high-energy reflectivity analysis have not been as keenly felt.

The absorption of energy of a specimen may be determined from a knowledge of R and T . The most important absorption mechanism is that associated with the electronic transitions discussed previously. The absorption edge spectra for a number of semiconductors were reproduced in Fig. 13. The sub-bandgap region has been studied extensively using absorption spectroscopy where the spectra are dominated by impurity-related effects. The absorption from the electronic transitions associated with a number of impurities in Si was presented earlier in Fig. 33. The impurity-related vibrational features and the lattice phonon bands can be observed as discussed under the section, "Lattice." The collective charge carrier oscillations and impurity to band-type transitions may also be observed in this region.

Modulation Spectroscopy A very useful variation of the reflection and transmission analysis is the use of modulation techniques that produce a periodic perturbation of a property of the specimen or the light beam and detect in-phase changes in R and T . The schematic diagram of a measurement apparatus is presented in Fig. 54.¹²² The capacity of lock-in amplifiers for very large amplifications of $> 10^5$ with a concomitant reduction in broadband noise is the heart of the procedure. Modulations in R and T , namely, ΔR and ΔT , of 10^{-5} can be easily detected. This means the modulating perturbation can be quite small and accessible in routine experimental systems. In the simplest cases, the modulation response of a specimen to a property x may be expressed as follows:

$$\Delta R = \frac{dR}{dx} \cdot \Delta x \quad (70)$$

where Δx is the intensity of modulation. The modulated property may be internal to the specimen, such as the temperature and pressure, or external, such as the wavelength or polarization of the probe. The most attractive feature, from a measurement point of view, is the fact that modulation

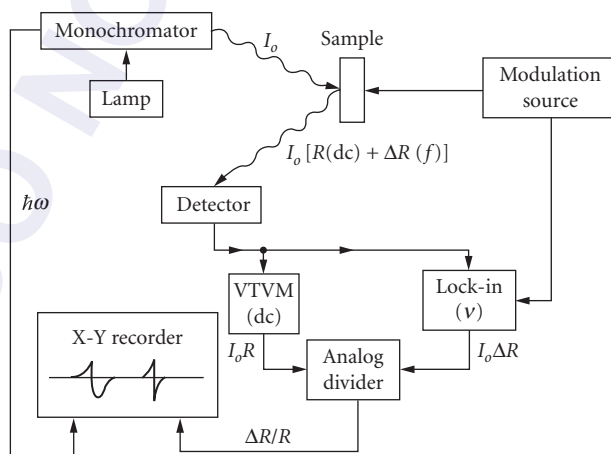


FIGURE 54 Block diagram of a typical modulation spectrometer.¹²²

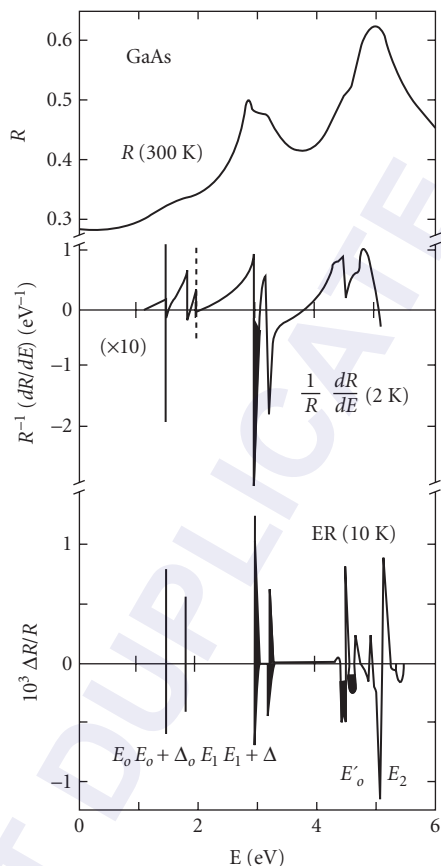


FIGURE 55 A comparison of three types of spectra from 0 to 6 eV for a typical semiconductor, GaAs. *Top*: reflectance R ; *middle*: energy-derivative reflectance; and *bottom*: low-field electroreflectance.¹²³

spectra are derivative-like and hence suppress slowly varying background structure and emphasize features in the vicinity of the critical points in the electronic band structure. Figure 55 presents a comparison of the reflective spectrum of GaAs in the 0 to 6 eV region and a wavelength-modulated spectrum and an electric-field-modulated spectrum.¹²³ Note the narrow linewidths and larger signal-to-noise ratios in both the WMR (wavelength modulation reflectance) and ER (electroreflectance) spectra. For instance, the $E_0 + \Delta_0$, indistinguishable in R , is clearly observable in both WMR and ER. These techniques are crucial in providing accurate values of the interband transition energies, which result in a better understanding of the electronic band structure of a material.

A variety of modulation procedures, underlying mechanisms, measured specimen properties, and other salient features are presented in Table 9.¹²³ The techniques are broadly classified as internal or external modulation to signify if the perturbation is intrinsic to the measurement approach, as in the case of WMR, or is external, as in the case of ER where an externally applied electric-field modulation is needed at the specimen surface. Temperature, stress, and magnetic-field modulation are included for completion. Compositional modulation that compares two nearly identical alloy samples is less frequently used to investigate semiconductors. Spectroscopic ellipsometry, to be discussed later, may

TABLE 9 Characteristics of Some Commonly Used Modulation Techniques¹²³

Technique	Name	Type	Variable	Sample Parameters Affected	Lineshape Type	Principal Parameters Measured	Principal Advantages	Disadvantages
Wavelength modulation: energy derivative reflectance	Internal	Scalar	Wavelength λ ; energy $\hbar\omega$	—	1st derivative	E_g, γ	Universal applicability; minimal sample preparation; fast, convenient	Measurement system can generate intrinsic structure, not easy to eliminate
Spectroscopic ellipsometry	Internal	Scalar	Wavelength λ ; energy $\hbar\omega$	—	Absolute	Dielectric function	As energy derivative reflects but on ϵ_1, ϵ_2 directly	Strongly influenced by surface preparation and thin films
Composition modulation	Internal	Scalar	Sample compared to control sample	—	Complicated	Doping, alloying effects	Obtains differences for parameters impossible to vary cyclically	Two samples/beam motion involved; alignment and surface preparation critical
Thermomodulation	External	Scalar	Temperature T	Threshold E_g ; broadening parameters γ	1st derivative	$E_g, \gamma, dE_g/dT, d\gamma/dT$	Wide applicability; identifies Fermi-level transitions in metals	Slow response (1–40 Hz); broad spectra
Hydrostatic pressure	External	Scalar	Pressure P	Threshold E_g	1st derivative	Deformation potentials		Cannot be modulated; must be used in conjunction with another technique
Light modulation (photoreflectance) (photovoltage)	External	Scalar or tensor	Intensity I of secondary beam	Carrier concentration or surface electric field	Complicated	E_g	Convenient; minimal sample preparation required	Effect on material usually not well-defined
Uniaxial stress	External	Tensor	Stress X	Threshold E_g ; matrix elements	1st derivative	E_g, γ symmetries; deformation potentials	Symmetry determination	Difficult to modulate; limited to high fracture/yield stress materials
Electric field	External	Tensor	Field E	Electron energy $E(k)$ oscillations	3rd derivative (low) Franz-Keldysh symmetries (high)	E_g, γ Effective masses in VUV; impurity concentrations	Very high resolution; only high-resolution technique usable	Requires certain resistivity ranges
Magnetic field	External	Tensor	Field H	Election energy levels	Landau levels	E_g , effective masses	Extremely high resolution	Advantages realized only for lower conduction band minima

also be considered to be a polarization-modulation technique. Modulation techniques are useful in yielding crystal properties such as the electronic transition energies as well as information regarding the perturbation mechanism such as the electro-optic or magneto-optic effects.

A widely used and very useful form of a modulation technique for the study of semiconductors is the electric-field-modulated reflection spectroscopy, referred to as electroreflectance. The basis of the procedure is the electric-field-induced changes in the optical response.

The electric-field-induced perturbations can be treated in detail by considering the effects of the applied potential on the electronic band structure. For relatively weak fields (i.e., for field strengths not large enough to modify the band structure significantly), the major perturbation mechanism may be considered to be the acceleration of the electron to a successive set of momentum states. The perturbation to $\epsilon(E, \mathbf{E})$, for such a simple case may be expressed as:¹²³

$$\Delta\epsilon = \frac{(\hbar\Omega)^3}{3E^2} \frac{\partial^3}{\partial E^3} [E^2 \cdot \epsilon(E)] \quad (71)$$

where E is the energy and

$$(\hbar\Omega)^3 = (e^2 \mathbf{E}^2 \hbar^2) / (8\mu_{\parallel}) \quad (72)$$

e is the electronic charge, and μ_{\parallel} the effective mass parallel to \mathbf{E} .

When the field is large, the electroreflectance spectra display oscillatory structure, known as Franz-Keldysh oscillations, and $\Delta R/R$ has a more complicated functional form that can be found in Ref. 124.

In practical terms, the oscillations may be analyzed to determine $\hbar\Omega$ and, hence, the strength of the electric field that causes the perturbation.

A particularly attractive and widely used form of electroreflectance is the contactless form of electric-field-modulated reflectivity known as photorefectance (PR).¹²⁵ In this procedure, one uses the electric-field modulation produced near the surface of the specimen by a chopped laser beam to modulate the reflectivity of a weak probe beam. PR is observable when there is band bending at the surface of the semiconductor such as that observed in damaged oxide layers on silicon.^{126,127}

The sharp features in the PR spectrum associated with the critical points in the electronic density of states including the direct gap E_g may be employed to determine the transition's energies accurately. The information obtained using ER and PR was instrumental in leading to a detailed understanding of the band structure of semiconductors. The variation of the critical point energies with alloy composition has been used to understand the electronic behavior as well as characterize semiconductor alloys.¹²⁸ Similar studies have been used extensively in the study of microstructures, where PR is particularly useful since it allows the observation of the gap as well as several additional higher-energy transitions as shown in Fig. 56.¹²⁹ Distinct transition from both the well and the barrier region can be observed, and, hence, a complete picture of the microstructure can be obtained. The electric-field-induced Franz-Keldysh oscillations in the ER and PR line shapes have been used to establish the electric field strengths,^{124,130} and, more recently, a similar technique has been used to measure the electric field strength in the surface region of GaAs and the effects of passivation¹³⁰ as shown in Fig. 57. The electric field strength can be determined from the slope of the inset in Fig. 57. Since PR measurements can be performed at room temperature with minimal sample preparation, it is attractive for routine characterization. Commercial PR spectrometers are now available for use in semiconductor fabrication. Several excellent reviews of the applications of PR are available. A compilation of activity may be found in the Proceedings of the 1990 International Conference on Modulation Spectroscopy¹³¹ and the review by Glembocki and Shanabrook,¹³² and the reader is referred to them for more details. A more recent reference is the third International Workshop on Modulation Spectroscopy of Semiconductor Structures, July 3–5, 2008, Wrocław, Poland.

Ellipsometry The reflection and transmission measurements discussed in the previous section consider the ratio of the incident to the reflected or transmitted optical power and hence ignore

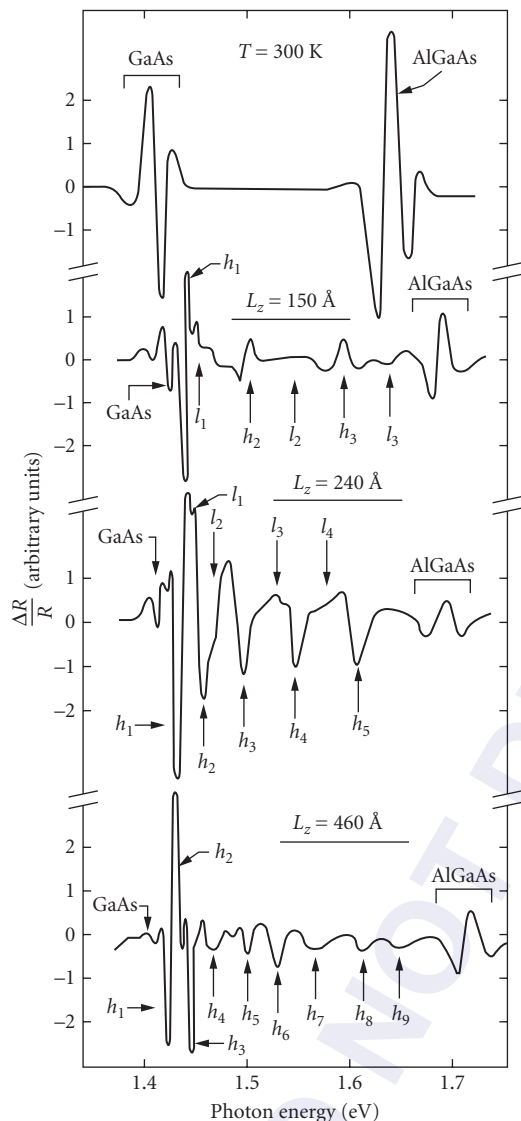


FIGURE 56 Room-temperature PR spectra for an undoped GaAs/AlGa_{1-x}As heterojunction and three multiple quantum-well samples with $x \approx 0.2$, L_z is the well thickness, and arrows labeled h_1, h_2, \dots (l_1, l_2, \dots) correspond to calculated values of interband transitions involving heavy (light) hole valence bands.¹²⁹

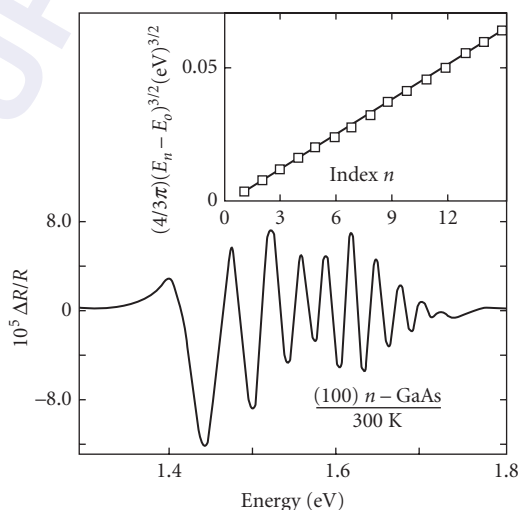


FIGURE 57 Photoreflectance spectrum at room temperature with $I_{\text{pump}} = 3 \mu\text{W}/\text{cm}^2$ and $I_{\text{probe}} = 2 \mu\text{W}/\text{cm}^2$. The inset shows a plot of $(4/3\pi)(E_n - E_0)^{3/2}$ as a function of index n . The slope of the line in the inset yields the electrical field strength.¹³⁰

any information carried by the phase change suffered by the incident beam. A complete description of the reflection or transmission process will have to include the phase information as well. Ellipsometry attempts to obtain part of the phase information by measuring the phase difference, introduced upon reflection, between the normal components of obliquely incident plane-polarized

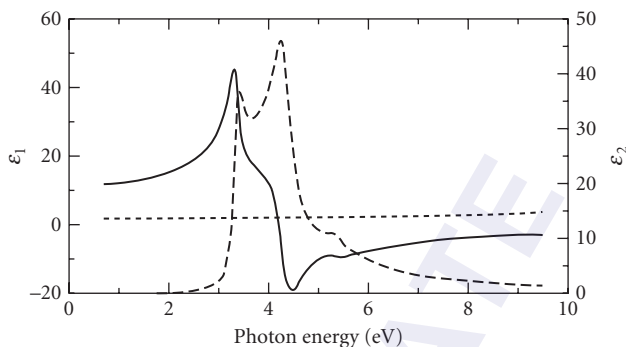


FIGURE 58 Real (solid) and imaginary (dashed) parts of the dielectric function of Si determined using VUV spectroscopic ellipsometry. For comparison, the real part of the dielectric function of SiO_2 is given by the short-dashed line.¹³⁸

light. Even though the absolute phase suffered by each component is not measured, the phase difference may be measured easily as an ellipticity in the polarization state of the reflected light, hence, the name ellipsometry. Recent developments in ellipsometry as well as handbooks covering practical applications can be found in Refs. 133, 134, 135, 136, and 137.

The intrinsic sensitivity of ellipsometry may be illustrated by considering a thin film with thickness d , refractive index n , and measurement wavelength λ . The phase change in the reflected component can be measured to a precision of 0.001° in ellipsometers which translates to a sensitivity of approximately $10^{-4}\lambda$ in thickness. Such extraordinary precision may be used for the study of sub-monolayer films or as a real-time monitor for measuring very small changes in materials' properties. With recent advances in instrumentation, the lower wavelength limit has been extended to below 135 nm (9.5 eV). As an example, Figure 58 shows the dielectric functions of Si and SiO_2 between 0.7 eV and 9.5 eV. Due to the extended spectral limit into the VUV, ellipsometry has become attractive to study wide-bandgap semiconductors, insulators, oxides, and very thin films.¹³⁸

Recent extensions of ellipsometry to long wavelength now permit investigation of lattice vibrations (phonons) and free-charge carriers (plasmons) in complex semiconductor layer structures.¹³⁹ The infrared dielectric function is sensitive to strain, composition, and ordering. Furthermore, concentration, mobility, and effective mass parameters can be determined from far-infrared magneto-optic generalized ellipsometry measurements.^{140–143} Even the sign of the free-charge carrier can be determined, just as in an electrical Hall effect measurement. Finally, magneto-optic generalized ellipsometry has even been extended into the terahertz frequency domain using an intense terahertz synchrotron source at wavelengths from approximately 30 to 270 cm^{-1} (0.9–8.1 THz).¹⁴⁴ Even Landau transitions could be seen in a two-dimensionally confined electron system in a grapheme sample.

Luminescence The process of luminescence, as described earlier, occurs in a suitably excited specimen. It is a mechanism through which the excited specimen relaxes to the equilibrium state.^{108,109} Hence, unlike reflection and transmission spectroscopies, luminescence procedures concentrate on the relaxation of the specimen and often lead to complementary information. For instance, impurities in semiconductors (particularly at low concentrations) are impossible to detect through reflection and more difficult to detect by absorption than by luminescence. Luminescence spectroscopy is thus an important part of the analysis of the optical behavior of semiconductors.¹⁰⁸ Moreover, since one of the main applications of semiconductors is in the arena of light emitters including lasers, the study of luminescence provides direct access to device optoelectronic information.

Luminescence processes may be induced by excitations that produce free electron-hole (e-h) pairs that may recombine across the bandgap or through defect- and impurity-related intermediate steps and emit a photon. The excitations employed most often are (1) an incident intense above-bandgap radiation such as from a laser source, (2) an incident electron beam, (3) electrical injection of electrons and/or holes through an appropriate contact, and (4) thermal excitation. These procedures are known as photoluminescence (PL), cathodoluminescence (CL),¹⁴⁵ electroluminescence (EL),¹⁴⁶ and thermoluminescence (TL), respectively. The most widely used technique for the analysis of semiconductor materials is PL. Cathodoluminescence measurements are usually conducted in a scanning electron microscope (SEM). The SEM electron beam can be focused to a spot size 1000 Å and can be scanned over the area of the sample. Hence, much work has been performed in CL imaging of wafers where one can obtain not only spectroscopic information but also spatial details. Electroluminescence is the most difficult to obtain because of the complexity of producing appropriate contacts. However, in terms of application, EL is the most important since a light emitter has to be able to produce light in an efficient manner under electrical excitation. Thermoluminescence is a technique used with insulators and wide-gap materials and is not widely employed in the analysis of the commercially important materials such as GaAs and Si. The principles of luminescence and the optical information regarding the semiconductor that can be obtained are discussed in the next section using PL.

Photoluminescence Luminescence processes may be excited using an above-bandgap-beam of light that leads to the creation of an electron-hole pair that may recombine across the gap and emit a photon with energy equal to the gap E_g . This process is schematically displayed in Fig. 59.¹⁰⁹ Two

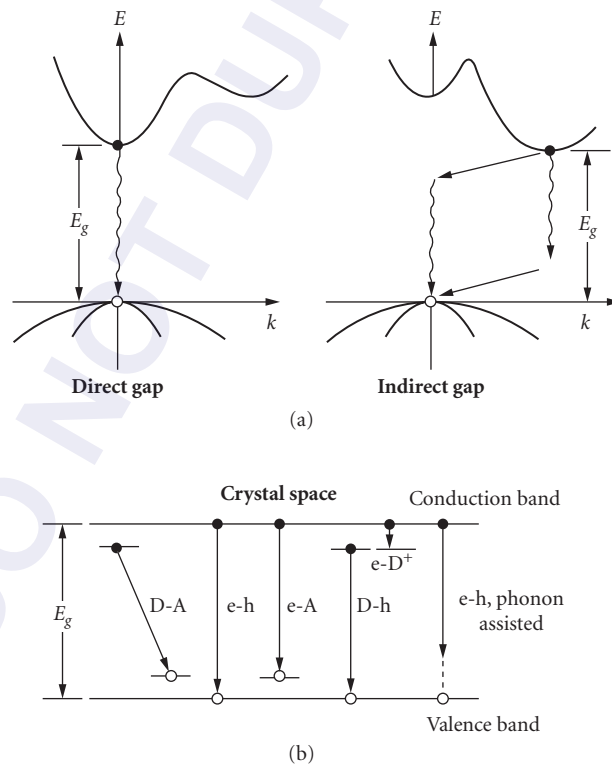


FIGURE 59 Schematic of photoexcitation and relaxation in semiconductors.¹⁰⁹

possibilities are shown: namely, the direct-gap and the indirect-gap semiconductors. The recombination will be direct in the former and will have to involve an additional participant, mostly a phonon, to conserve momentum in the latter. This is the simplest possible recombination mechanism. Several additional routes exist for the relaxation that involves impurities. The photoexcited charges may recombine with ionized acceptors (A) and donors (D) with or without involving a phonon. The process may also be more complicated where, for example, an electron may be trapped by an ionized donor which may subsequently recombine with a hole at a neutral acceptor, leading to a donor-acceptor pair (DAP) transition. In addition, recombination may also occur through the annihilation of excitons considered further in this section.

Excitons, as discussed earlier, are hydrogen-like two-particle electron-hole combinations that are not included in the one-electron energy band description of the solid. The strong Coulomb interaction between the electron and the hole leads to the excitonic coupling.

As the excitation intensity increases, so does the population of the excitons. Higher-order interactions can occur between excitons and entities, such as biexcitons, and may be observed. Under intense excitation, the electrons and holes form a liquid state known as the electron-hole liquid. The presence of impurities and defects can also significantly alter the nature of the luminescence spectra, in particular the excitonic behavior. The electric field in the vicinity of the impurity can trap and localize excitons. Such an interaction leads to bound excitons; the binding energy of the bound exciton E_{BE} will, in addition to E_{ex} given in Eq. (37), contain a localization part as shown:

$$E_{BE} = E_{ex} + \delta \quad (73)$$

The impurity potential that confines an exciton to a given center depends on both the impurity and its local environment. Hence, the impurity-bound excitonic features can be very rich and informative. A detailed knowledge of the PL excitonic spectrum may be used to identify both the chemical species and its environment, as is demonstrated later on. The observed luminescent photon energy at low temperatures may be written as

$$E = E_g - E_{BE} \quad (74)$$

Transitions from a free electron or hole to a neutral acceptor or donor, respectively, will occur at the following energies:

$$E = E_g - E_{ion} \quad (75)$$

where E_{ion} is the relevant ionization energy. A complete description of the free-to-bound transition will have to take into account both the dispersion of the band in the vicinity of the minima and the population distribution. For the donor-to-acceptor transition, an additional electrostatic energy term $e^2/\epsilon(0)r_{sep}$ will have to be accounted for as shown:

$$E = E_g - (E_A + E_D) + \frac{e^2}{\epsilon(0)r_{sep}} \quad (76)$$

where r_{sep} is the distance that separates the two participating centers. The last term is needed to account for the electrostatic energy of the final ionized state of both centers. All of the transitions discussed may occur with phonon participation, and hence the emission energies should be reduced by the quantum of the energy of the phonon; multiple phonons may also be involved in more complex spectra.

The measurement of PL is, in most cases, routine and straightforward, which partly accounts for its popularity as a materials characterization technique. The laser provides the excitation and the

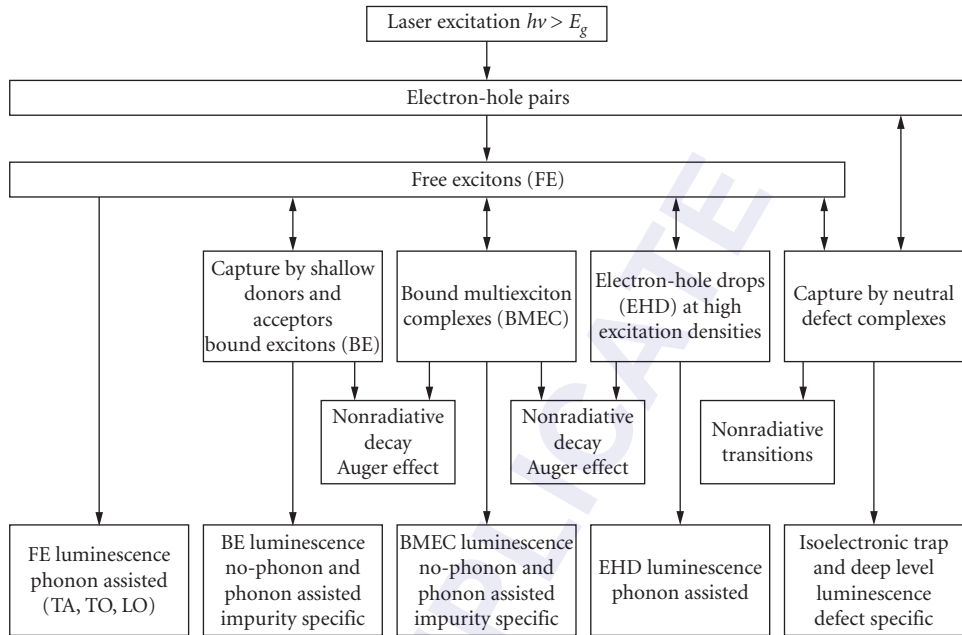


FIGURE 60 Luminescence decay processes stimulated by above-bandgap photoexcitation in high-purity Si.¹⁴⁷

dispersive spectrometer along with a sensitive detector, the detection. The sample should be cooled to ≈ 5 K so that the temperature-induced broadening is kept to a minimum and the population of the processes with a small activation energy, such as the excitons, is sufficiently high to perform accurate measurements. Photoluminescence measurements may also be performed using Fourier transform (FT) spectrometers. The main advantage is the extraordinary resolution that can be achieved with the FT systems.

The extension of PL into a tool with quantitative accuracy, particularly in very high-purity Si, has demonstrated the versatility of the analysis.¹⁴⁷ The paths of luminescence decay in high-purity Si are schematically presented in Fig. 60. The first step in the excitation process is the creation of electron-hole pairs, which subsequently can undergo a wide variety of processes before recombination. Figure 60 represents the possible intermediate states that eventually lead to radiative recombinations. A typical spectrum obtained from a sample containing $1.3 \times 10^{13} \text{ cm}^{-3}$ B, $1.8 \times 10^{12} \text{ cm}^{-3}$ P, and $3 \times 10^{11} \text{ cm}^{-3}$ As is presented in Fig. 61.¹⁴⁷ The spectrum displays both the no-phonon (NP) component (shown as an inset) and TA, TO, and LO phonon-assisted features for the free exciton (FE) and the impurity-related bound exciton features. The measurement was performed at 4.2 K. Since B is an acceptor and P and As are donors in Si, electrical transport analyses are not sufficient to fully analyze the impurities. In contrast, all three impurities can be unambiguously identified using PL. In addition, when dependable calibration curves are available, the concentration of each species can also be established. A representative calibration curve is presented in Fig. 62 for the impurities B, Al, and P.¹⁴⁷ Note that quantitative measurements can be performed down to 10^{12} cm^{-3} for all three impurities. In the case of B, in ultrapure samples, measurements can be performed to levels as low as $\sim 10^{10} \text{ cm}^{-3}$.

Next, we turn our attention to GaAs, which is a direct-bandgap material and hence exhibits efficient luminescence. Figure 63a displays a representative spectrum from a high-quality MOCVD-grown sample doped with C and Zn measured using an FT spectrometer.¹⁴⁸ The spectra shown are intense, sharp, and well-resolved near the band edge, with excitonic features appearing slightly below

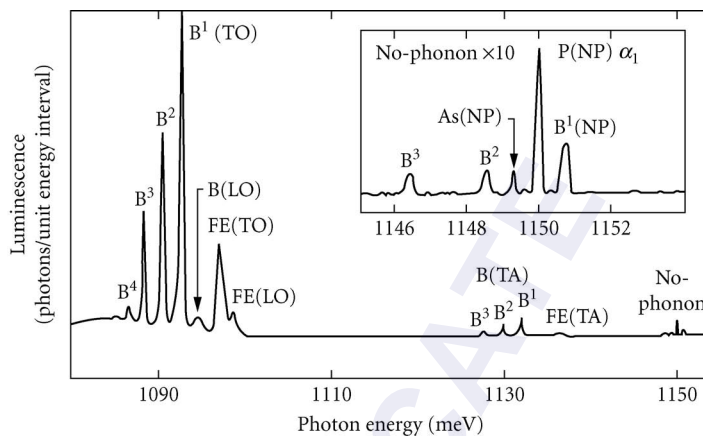


FIGURE 61 Photoluminescence spectrum from a Si sample doped with $1.3 \times 10^{13} \text{ cm}^{-3}$ B, contaminated with $1.8 \times 10^{12} \text{ cm}^{-3}$ P, and $3 \times 10^{11} \text{ cm}^{-3}$ As.¹⁴⁷

the band gap of $\approx 1.519 \text{ eV}$ at 4.2 K. The near-edge features, labeled (A^0, X) and (D^0, X), are due to excitons bound to neutral donors and acceptors, and the deeper-lying features are the free-to-bound transition to the C acceptor, (e, C^0) and the donor-acceptor-pair transitions, (D^0, Zn^0), and (D^0, C^0), to the Zn and C centers, respectively. An expanded high-resolution version of the same spectrum, in the vicinity of (A^0, X) region, is displayed in Fig. 63b.¹⁴⁸ Note the impressive resolution obtainable in the FT system and the clear resolution of the excited states associated with the bound exciton.

The use of PL techniques may be extended deeper in the infrared region as well. Deep levels associated with impurities such as Fe, Cr, and Ag in GaAs are known to produce luminescence bands at ~ 0.35 , ~ 0.6 , and $\sim 1.2 \text{ eV}$.¹⁴⁹ The narrow-gap semiconductors such as InAs, InSb, and HgCdTe luminesce in the 100 to 200 meV region farther into the infrared region. An example of the PL spectra observed from $\text{Hg}_{0.78}\text{Cd}_{0.22}\text{Te}$ narrow-gap alloy¹⁵⁰ is displayed in Fig. 64. The spectra, in comparison to those presented earlier, are rather featureless due to the fact that the effective masses in the

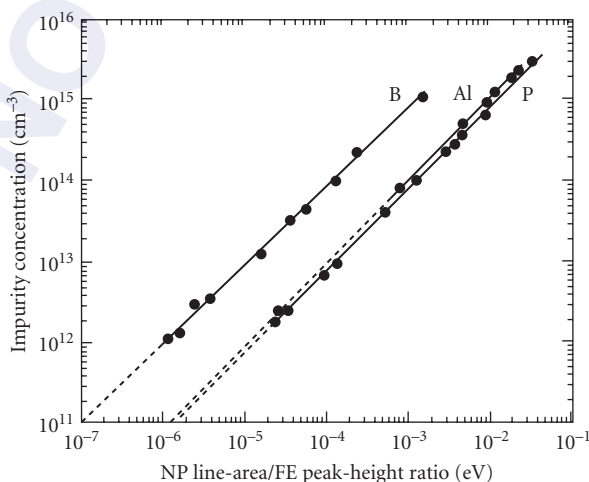


FIGURE 62 Calibration of the photoluminescence technique for measuring B, P, and Al concentrations in Si.¹⁴⁷

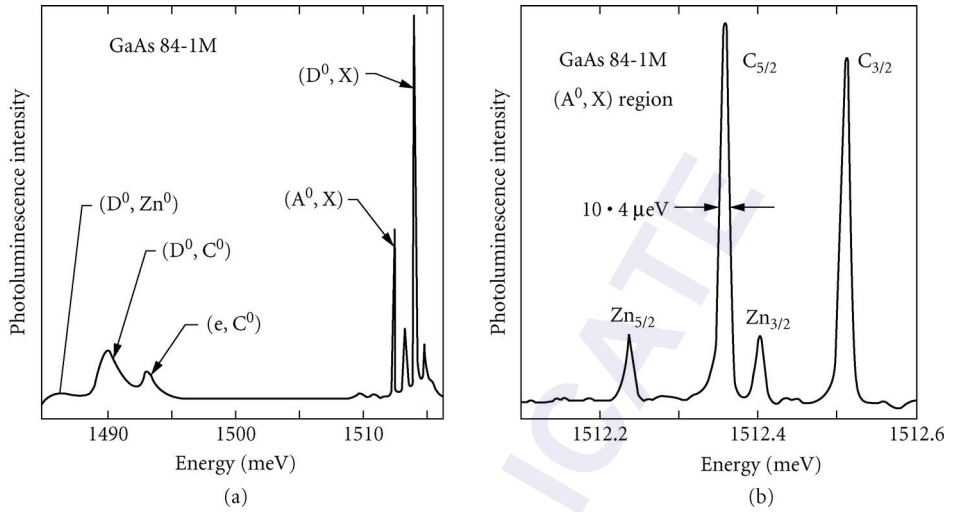


FIGURE 63 (a) Representative PL spectrum from MOCVD-grown GaAs.¹⁴⁸ (b) High-resolution, near-band edge spectrum of the same sample as in (a).¹⁴⁸

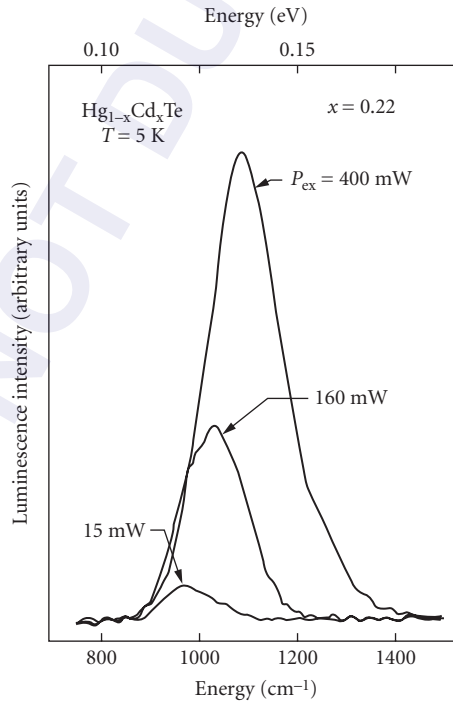


FIGURE 64 Luminescence spectra from $\text{Hg}_{0.78}\text{-Cd}_{0.22}\text{Te}$ measured at 6 K for different 1.06 μm excitation power.¹⁵⁰

narrow-gap semiconductors tend to be small, and hence the excitonic binding energies are small as well. Therefore, sharp, well-resolved spectra are not normally observed. The observation of luminescence is much more difficult in narrow-gap semiconductors.

The luminescence processes described provide information on the relaxation mechanisms and hence are heavily weighted toward transitions that involve only the first excited state and the final ground state of a system. A complete study of a recombination process, for instance the FE, requires information regarding the higher excited states. This can be achieved by a variation of the PL procedure known as photoluminescence excitation (PLE) spectroscopy. The crux of this technique is to concentrate on the PL response with respect to the excitation wavelength and thereby to determine the excitation resonances. Since the resonances occur when the excitation photon energy matches the excited energies, information regarding the excited states may be elucidated. The principle is illustrated with Fig. 65, where an attempt is made to determine the process that the photoexcited electron-hole pair undergoes before the formation of a free exciton in high-purity CdTe.¹⁵¹ The spectrometer was set to the FE energy of 1.596 eV, and the excitation photon energy was scanned from 1.6 to ~ 1.8 eV. The excitation wavelength may be scanned using a dye laser as shown in Fig. 65. Two PLE spectra obtained at 1.8 and 20.6 K are displayed in Fig. 65. The strong oscillatory behavior, with spacing of 21 meV, demonstrates a cascade process through intermediate states separated by the LO phonon energy of 21 meV. This study illustrates the importance of the LO phonon in the electron-hole relaxation process and yields information regarding electronic states that lie above the conduction band extremum, involving the hot electron behavior.

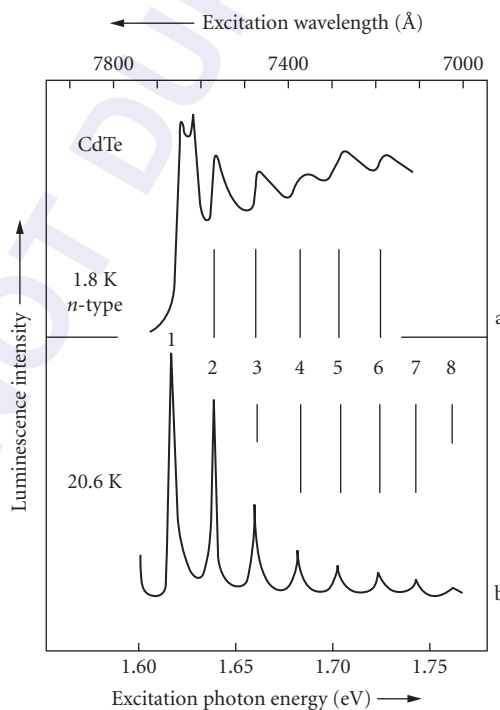


FIGURE 65 Excitation spectra of exciton luminescence of the emission line at 1.596 eV in the energy region higher than the bandgap of CdTe: (a) a pure *n*-type sample at 1.8 K and (b) the sample at 20.6 K.¹⁵¹

Inelastic Light Scattering (Raman and Brillouin) The elastic interaction of the incident photon with the specimen implies a process where the energy or wavelength of the photon is preserved. For instance, reflection and transmission spectroscopies involve the measurement of the fraction of the incident beam that is reflected or transmitted with no change in the incident wavelength. The vast majority of the incident photons undergo only elastic interactions, but a tiny fraction, of the order of $\sim 10^{-8}$, are subjected to inelastic scattering. Inelastic interactions, though very few, are extremely important as probes of the properties of the specimen. The advent of high-powered coherent laser sources and sensitive detectors has made the measurement of the scattered spectra straightforward, leading to an explosive growth in the last 20 years. Several excellent reviews may be found in the five volumes on light scattering, edited by Cardona and Guntherodt.¹⁵² A good introduction to the procedure and theory may be found in Ref. 153.

The experimental procedure involves projecting a high-power laser beam onto a specimen, collecting the back-scattered light, and performing a spectrum analysis and detection. Multiple-dispersion stages are needed for improved resolution and filtering out of the unwanted elastically scattered beam. A typical set of spectra measured from GaAs⁶⁸ is presented in Fig. 66. The axes for the spectra are the detected intensity and the frequency shift suffered by the incident light. The intense peak at 0 cm^{-1} is the elastically scattered peak. The spectrum at the room temperature of $\sim 300 \text{ K}$ displays both the frequency up-shifted and down-shifted components labeled as Stokes and anti-Stokes features corresponding to the TO phonon interaction; these correspond to microscopic processes where the incident photon loses or gains energy due to TO phonon emission or absorption, respectively. At lower temperatures, additional features labeled L_+ and L_- appear near the TO peak as well as a broad feature in the 0 to 100 cm^{-1} region. These features originate from charge carriers and were discussed in earlier sections.

The inelastic scattering techniques are usually divided into two categories, namely, Brillouin and Raman scattering. The former covers low frequencies extending from 0 to $\sim 10 \text{ cm}^{-1}$, while the latter spans the higher frequencies. The low-frequency Brillouin scattering measurements provide access to information on properties of acoustic phonons, spin waves, etc., and involve the use of specialized spectrometers needed to remove the very intense elastic peak.¹⁵⁴ The position of the elastic peak is taken to be the reference zero position, and the scattered spectrum is measured with respect to the zero position; that is, the spectrum is recorded as a function of the frequency shift and not the absolute frequency. Raman scattering is, in general, easier to perform and more informative for the study of semiconductors and hence is emphasized in this section.

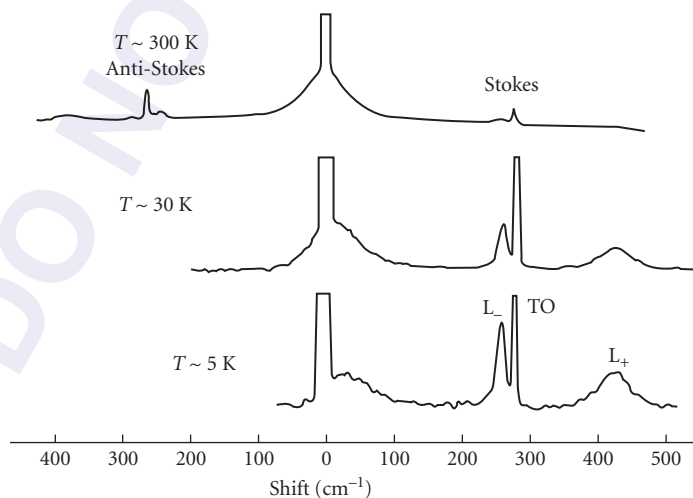


FIGURE 66 Typical Raman spectrum in GaAs, with $n = 1.4 \times 10^{18} \text{ cm}^{-3}$, at 300, 30, and 5 K showing the lineshape change with temperature.⁶⁹

The techniques reviewed so far all involve first-order interactions that provide direct information regarding the electronic, vibrational, and impurity-related behavior of semiconductors. In contrast, inelastic light scattering such as Raman and Brillouin scattering involves the electron-phonon or other quasi-particle interactions and hence can provide additional information regarding these interactions. In addition, since inelastic scattering involves a higher-order interaction, processes that are inactive in first-order may be investigated. For instance, the optical phonons in Si do not possess a dipole moment and therefore do not interact with infrared radiation but can be clearly observed in Raman scattering.

The response of the specimen may be expressed, through the susceptibility χ .¹⁵⁵ However, the presence of the inelastic interactions will give rise to additional contributions. Consider the example of lattice vibrations and their influence on χ in the visible-frequency range. Since the lattice vibrational frequencies correspond to far-infrared light frequencies, no direct contribution is expected. However, the lattice vibrations can influence χ even in the visible-frequency range in an indirect fashion by a small amount. The influence of this interaction may be expressed as follows:

$$\begin{aligned}\chi &= \chi_0 + \sum_i \frac{\partial \chi}{\partial u_i} \cdot \mathbf{u}_i \\ &+ \sum_{i,j} \frac{\partial^2 \chi}{\partial u_i \partial u_j} \mathbf{u}_i \mathbf{u}_j \\ &+ \dots\end{aligned}\quad (77)$$

where \mathbf{u}_i and \mathbf{u}_j are the displacements associated with the normal modes of lattice vibrations or phonons. The first term contains the elastic term, and the second and third terms denote the second- and third-order inelastic interactions with the phonons. Note that the form of χ shown here differs from that used previously in nonlinear optics because it is expressed explicitly in terms of the phonon coordinates. Assuming sinusoidal oscillations for the incident radiation \mathbf{E} and the phonons \mathbf{u} , the polarization induced in the specimen may be expressed as follows:

$$\mathbf{P} = \chi_0 \mathbf{E}_0 \exp(i\omega_0 t) + \sum_i \chi' u_{i0} \exp\{i(\omega_0 \pm \omega_i)t\} + \sum_{ij} \chi'' u_{i0} u_{j0} \exp\{i(\omega_0 \pm \omega_i \pm \omega_j)t\} \quad (78)$$

The source of the inelastically scattered or frequency-shifted components is immediately apparent.

The scattered intensity, assuming only second-order interactions, is usually expressed in terms of a differential cross section; that is, scattered energy per unit time in the solid angle $d\Omega$ as follows:

$$\begin{aligned}\frac{d\sigma}{d\Omega} &= V^2 \left(\frac{\omega_s}{c}\right)^4 \hat{\mathbf{e}}_i \chi' \hat{\mathbf{e}}_s n_j \quad \text{Stokes} \\ &\quad (n_j + 1) \quad \text{Anti-Stokes}\end{aligned}\quad (79)$$

where $\hat{\mathbf{e}}_i$ and $\hat{\mathbf{e}}_s$ denote the polarization state of the scattered and incident light, ω_s the frequency of the scattered beam, and V the scattering volume.

χ' and χ'' are Raman tensors whose symmetry properties may be calculated from a knowledge of the structure of the crystals.¹⁵⁶ The notation used to express the combination of incident and scattered beam directions and polarizations is as follows: $a(b, c)d$, where a and d denote the incident and scattered directions, and b and c denote the respective polarizations. Using the cubic crystal axes x , y , and z , the selection rules may be expressed as follows.

Back scattering from (100) surface:

TO—disallowed for any combination of incident and scattered beam configurations

LO—allowed only for crossed polarizations; that is, $z(x, y) \bar{z}$ and $z(y, x) \bar{z}$, where \bar{z} is the opposite direction to z .

Similar selection rules may be derived for other orientations and crystal structures.

A microscopic description of the scattering event considers quantum mechanical details using perturbation theory. The incident photon with frequency ω_i and wave vector \mathbf{k}_i produces a transition from the initial state i to a virtual state b where the incident photon is annihilated, followed by a transition to the final state f , accompanied by the emission of the “scattered” photon with frequency and wave vector ω_s and \mathbf{q}_s , respectively. The entire process is part of a complete quantum mechanical event and will be governed by momentum, energy, and symmetry conservation rules shown as

$$\begin{aligned}\hbar\omega_i &= \hbar\omega_s \pm \hbar\Omega \\ \hbar\mathbf{q}_i &= \hbar\mathbf{q}_s \pm \hbar\mathbf{K}\end{aligned}\quad (80)$$

where Ω and \mathbf{K} denote the frequency and wave vector of the participating phonon. The exact forms of the scattering cross section will depend on the details of the interactions and are reviewed in Ref. 153.

One of the most powerful aspects of Raman scattering is the ability to perform resonance excitation. When the incident photon energy or the scattered photon energy matches an intrinsic excitation energy, a substantial increase in the scattered intensity is observed. This ability to resonate has been used for both an understanding of the details of the scattering process and applications in which the source of the scattering can be selected. For instance, the study of a particular impurity in a large matrix may be conducted by tuning the resonance to match that of the impurity. Examples of resonance studies are discussed in the following sections.

The power and versatility of Raman scattering to probe many important properties of semiconductors have been exploited for a wide variety of characterizations. Figure 67 schematically displays

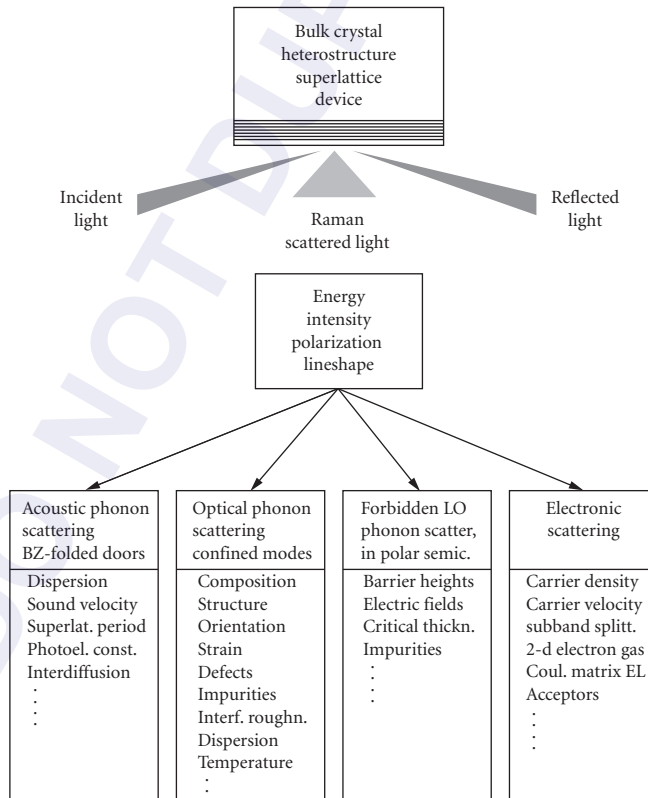


FIGURE 67 Schematics of inelastic light scattering and information which can be extracted from such measurements.¹⁵⁷

the broad areas where Raman scattering has been employed.¹⁵⁷ The study of materials, devices, and microstructural properties include chemical, structural, and electronic properties. The developments in microscopic measurements and analysis have opened up new applications as well. The measurement of the temperature in a spot $\approx 1 \mu\text{m}$ in diameter is one such example that is discussed in the following sections.

The most dominant Raman-scattering mechanism in semiconductors is that due to phonons—in particular, the optical phonons at the center of the Brillouin zone—namely, TO (Γ) and LO (Γ). Even though a semiconductor crystal can support a variety of vibrations, the momentum conservation requirements restrict the interaction of the incident photon to only the TO (Γ) and LO (Γ). The position and shape of the spectra contain important information regarding the structural state of the material: the presence of strain will be reflected as shifts in the line position; degraded crystal quality due to multiple grains and concomitant distributed strain will lead to a broadening of the lineshape. The effects of crystal damage on the Raman spectrum as a result of ion implantation in GaAs is displayed in Fig. 68.¹⁵⁸ Note the large increase in the linewidth of the LO (Γ) feature. In addition, a series of new structures, not present in the undamaged sample, can also be observed. These have been interpreted to be the result of relaxing the momentum conservation laws and hence the activation of phonon-scattering processes not normally allowed in a good crystal. A simple interpretation of the broadening of the LO lineshape was provided using a phonon confinement model, where the LO phonon is described as being confined to small damage-free regions. The shift

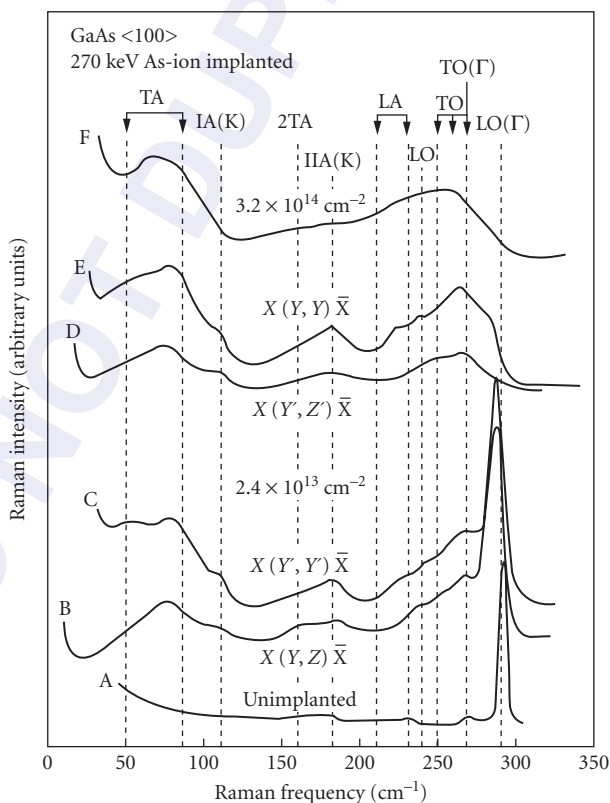


FIGURE 68 Raman spectra of $\langle 100 \rangle$ GaAs before implantation (A), implanted to a fluence of $2.4 \times 10^{13} \text{ cm}^{-2}$ for various polarization configurations (B, C, D, and E) and a fluence of $3.2 \times 10^{14} \text{ cm}^{-2}$ (F).¹⁵⁸

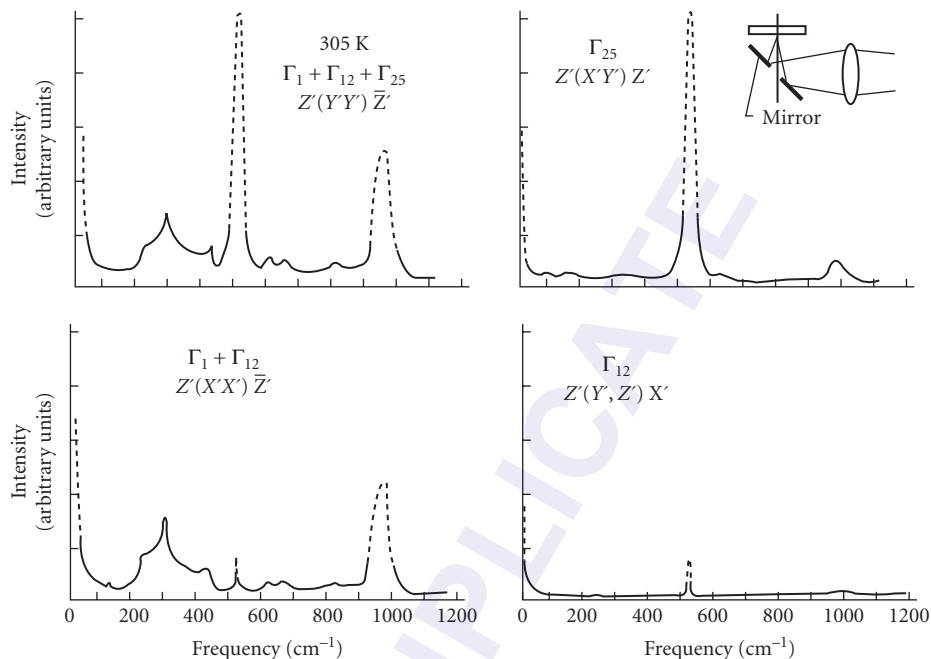


FIGURE 69 The Raman spectra of Si at room temperature showing the multiphonon contributions. The polarization configuration and the representations which were possible contributors to the spectra are shown for each of the four spectra.¹⁵⁹

in the position of the LO peak and the full-width-at-half-maximum were related to the diameter of the undamaged region. The information obtained is useful in characterizing lattice damage, amorphous materials, and degrees of recovery during annealing.

As explained earlier, Raman processes can involve higher-order interactions that involve more than one phonon. Hence, multiphonon scattering can be performed, and information complementary to that discussed in the section entitled “Lattice” on the infrared absorption properties can be obtained. The Raman spectra obtained from Si¹⁵⁹ are displayed in Fig. 69. The spectra were recorded at room temperature where $X' = (100)$, $Y' = (011)$, and $Z' = (0\bar{1}1)$. The irreducible representation of the phonons involved is also noted in the figure. Note the strong peak at $\sim 522 \text{ cm}^{-1}$ which is due to the degenerate $\text{TO}(\Gamma)/\text{LO}(\Gamma)$ phonons that are not observable in infrared absorption measurements. Additional bands present in the 200 to 400 cm^{-1} and the 600 to 1000 cm^{-1} range are due to multiphonon scattering processes. The ability to employ polarization selection rules has been used effectively to isolate the symmetry character of the underlying vibrational features and can be used to eventually identify the source of the various features. Such studies are useful both in understanding the optical behavior of Si and in illuminating the lattice dynamical properties of the material.

In crystals that contain a substantial number of free carriers, the incident photon can scatter off collective charge oscillations, known as plasmons, as discussed earlier. The Raman spectra measured from three GaAs samples¹⁶⁰ with electron densities ranging from 1.95×10^{18} to $6.75 \times 10^{18} \text{ cm}^{-3}$ are presented in Fig. 70. The observed features in addition to the LO phonon peak at $\sim 293 \text{ cm}^{-1}$ are due to the coupled LO phonon-plasmon features, also discussed earlier. The variation of the L_+ mode frequencies with the carrier density may be calculated and compared to measurements, as was shown in Fig. 32. The shape of the L_+ and the L_- features can be used to deduce the mobility of the carriers as well. The uniqueness of the Raman results is that they can be employed to study the behavior of carriers near the

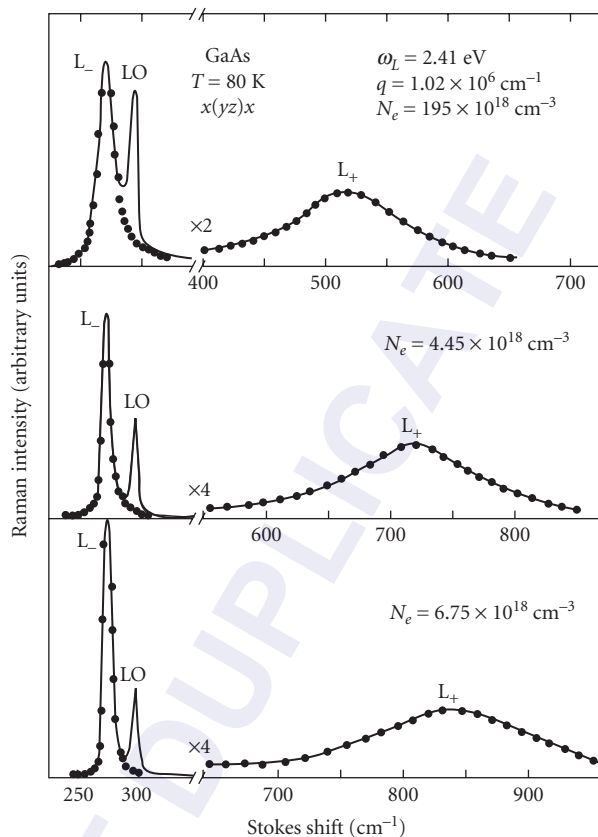


FIGURE 70 Raman spectra of three different *n*-GaAs samples obtained in backscattering geometry from (100) surfaces showing the coupled plasmon-LO phonon modes.¹⁶⁰

surface. Since the incident laser light with a photon energy of ~ 2.5 eV penetrates only 1000 \AA into the sample, the Raman results represent only the behavior of this near-surface region.

The single-particle nature of the free carriers can also be probed by Raman scattering.⁶⁹ The mechanism responsible for the interaction is the scattering of the carriers from inside the Fermi-surface to momentum states that lie outside, the total change in momentum being equal to that imparted by the photons. The integrated effect in the case of the spherical Fermi-surface in GaAs is displayed in Fig. 71. The Fermi wave vector is denoted by \mathbf{p}_F and that of the electron is \mathbf{p} . The wave-vector change as a result of the scattering is \mathbf{q} . Two cases of small and large \mathbf{q} and the resultant single-particle spectrum at 0 K are shown. The net effect in the first case will be a linear increase followed by a rapid fall, and, when \mathbf{q} is large, the spectrum displays a band-like behavior as shown. The single-particle spectrum measured from a sample of *n*-GaAs¹⁶¹ is presented in Fig. 72. The measurements were performed using the 6471 \AA line of the Kr^+ laser at 10 K. The ability to probe the Fermi sphere directly, using a spectroscopic technique, can lead to valuable insights into the electronic distributions that are complementary to those obtained from transport studies that usually provide only information regarding integrated effects of all the carriers.

The effect of resonance enhancement, as mentioned earlier, is one of the most powerful features of Raman scattering. The effect can be illustrated with the variation of the scattering intensity of the

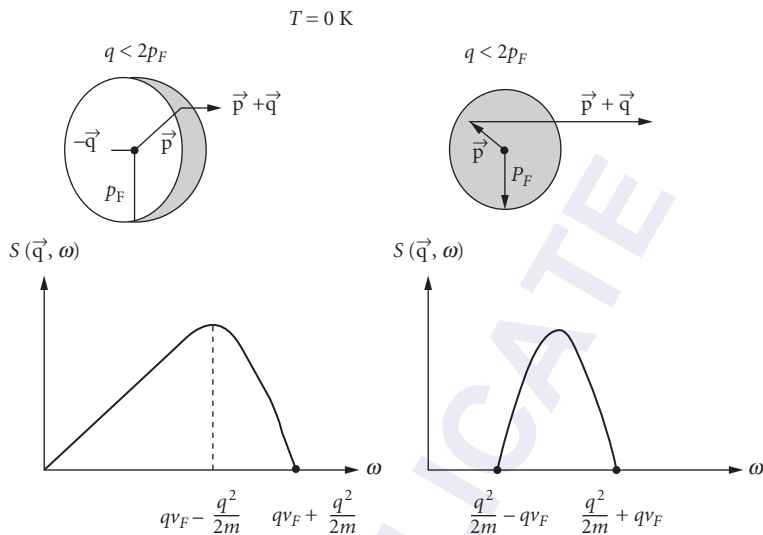


FIGURE 71 Single-particle excitation spectrum at 0 K.⁶⁹

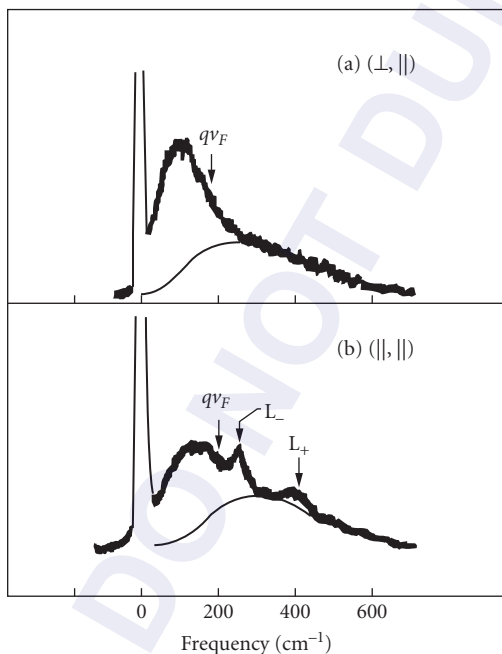


FIGURE 72 Single-particle spectra and coupled LO phonon-plasmon modes (L_{\pm}) for n -GaAs, $n = 1.3 \times 10^{18} \text{ cm}^{-3}$. Temperature: 10 K. Excitation: 6471 Å. The interband transition energy from the split-off valence band to the Fermi level is very close to the laser photon energy. Estimated luminescence background is shown by the dashed line.¹⁶¹

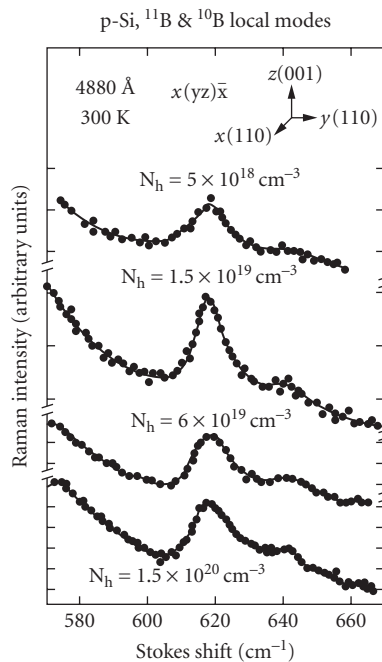


FIGURE 73 Scattering by local modes of boron (isotopes B^{10} and B^{11}) in p -type silicon.¹⁶³

optical phonons in CdS with the incident photon energy as discussed in Ref. 162. An order-of-magnitude enhancement was observed as the excitation energy approached the energy gap at ~ 2.6 eV for all the observed phonon modes. In addition, the TO modes both displayed a reduction before a large enhancement as E_g was approached. The reduction was interpreted as the result of a destructive interference between the resonant and nonresonant terms that contribute to the scattering cross section. Resonance Raman-scattering studies can therefore shed light on the microscopic details of the scattering process.

The application of Raman scattering to the localized vibration due to impurities is illustrated in Fig. 73 which presents the data obtained from B-doped Si.¹⁶³ In most cases, the density of the impurities needs to be quite high to be observable in Raman scattering. However, rapid advances in the measurement procedures may improve the sensitivities. Direct measurement of the electronic transitions related to dopant ions can also be observed in Raman scattering. A good example is the electronic interbound state transitions from three donors in GaP.¹⁶⁴ The normally symmetric, threefold degenerate 1S state of the donor, split into a 1S(A₁) singlet and a 1S(E) doublet due to the interaction with the conduction band valleys in the indirect-gap GaP was clearly observed as well-resolved peaks.¹⁶⁴ The Raman spectra are sensitive to the impurity electronic states and provide a tool to probe them as well.

5.5 ACKNOWLEDGMENTS

The authors express their deep appreciation for the help rendered by Erik Secula in preparing the manuscript.

5.6 SUMMARY AND CONCLUSIONS

An overview of the optical properties of semiconductors has been presented in this chapter. These properties form the foundation for understanding and utilizing the wide variety of optical devices manufactured today. A number of materials can be used together with electronic circuits to generate, detect, and manipulate light signals leading to the field of optoelectronics. Semiconductor materials are becoming increasingly important for use in optoelectronic devices: devices can be made very small, leading to a high degree of compactness and compatibility with other electronic and optical functions; they are robust and highly reliable; they are highly efficient as light-generating sources with internal efficiencies sometimes approaching 100 percent; they are capable of large power dissipation, of very high frequency performance, and can access an enormous range of wavelengths; and performances can be tuned over wavelength, frequency, and efficiency.¹⁶⁵ Table 10 lists some of the most important materials and their applications for optoelectronics.¹⁶⁵

In conclusion, Table 11 presents some of the important parameters for the most common semiconductors that determine the optical behavior of each material. The forbidden-energy gap and higher-energy critical point energies are listed along with optical phonon energies, dielectric constants, refractive index at energies below the energy gap, and free exciton binding energy. Closely related transport parameters such as the charge carrier effective masses and the mobilities are also included for completion. The reader is referred to Palik's compilation in Refs. 10 and 11 for refractive indices for several semiconductors over a wide range of energies. Additional information may also be obtained from many comprehensive collections of physical parameters such as those presented in the Landolt and Bornstein Tables,¹⁶⁶⁻¹⁶⁹ which are now available on the Internet.¹⁷⁰ A more concise semiconductor data handbook was recently published by O. Madelung, and much of the data, figures, references, and more detailed information were shifted to an accompanying CD-ROM.¹⁷¹

TABLE 10 Some Important Semiconductor Materials and Applications for Optoelectronics¹⁶⁵

Material	Type	Substrate	Devices	Wavelength Range (μm)	Applications
Si	IV	Si	Detectors, solar cells	0.5–1	Solar energy conversion, e.g., watches, calculators, heating, cooling, detectors
SiC	IV	SiC	Blue LEDs	0.4	Displays, optical disk memories, etc.
Ge	IV	Ge	Detectors	1–1.8	Spectroscopy
GaAs	III–V	GaAs	LEDs, lasers, detectors, solar cells, imagers, intensifiers, electro-optic modulators, optoisolators	0.85	Remote control TV, etc., video disk players, range-finding, solar energy conversion, optical fiber communication systems (local networks), image intensifiers
AlGaAs	III–V	GaAs	LEDs, lasers, solar cells, imagers	0.67–0.98	
GaInP	III–V	GaAs	Visible lasers, LEDs	0.5–0.7	Displays, control, compact disk players, laser printers/scanners, optical disk memories, laser medical equipment
GaAlInP	III–V	GaAs	Visible lasers, LEDs	0.5–0.7	
GaP	III–V	GaP	Visible LEDs	0.5–0.7	
GaAsP	III–V	GaP	Visible LEDs, optoisolators	0.5–0.7	
InP	III–V	InP	Solar cells	0.9	Space solar cells
InGaAs	III–V	InP	Detectors	1–1.67	Optical fiber communications (long-haul and local loop)
InGaAsP	III–V	InP	Lasers, LEDs	1–1.6	
InAlAs	III–V	InP	Lasers, detectors	1–2.5	
InAlGaAs	III–V	InP	Lasers, detectors	1–2.5	
GaSb/GaAlSb	III–V	GaSb	Detectors, lasers	2–3.5	Imaging
GaAlInN	III–V	Quartz, GaN, AlN	IR to UV source and detector	0.25–0.7	Displays, optical disk memories, bio detectors, solar cells, photochemistry, traffic lights, etc.
CdHgTe	II–VI	CdTe	IR detectors	3–5 and 8–12	Infrared imaging, night vision sights, missile seekers, and many other military applications
ZnSe	II–VI	ZnSe	Short LED, lasers	0.4–0.6	Commercial applications (R&D stage only)
ZnS	II–VI	ZnS	Short LED, lasers	0.4–0.6	
Pb compounds	IV–VI	Pb Compounds	IR lasers, detectors	3–30	Spectroscopy, pollution monitoring

TABLE 11 Material Parameters¹

Material	Type (I/D)	F_g						Higher-Energy Transitions (≈ 300 K)					
		300 K (eV)	77 K (eV)	≈ 0 K (eV)	dE_g/dT (10^{-4} eV/K)	E_0 (eV)	$E_0 + \Delta_0$ (eV)	E_g (eV)	$E_g + \Delta_1$ (eV)	E_0 (eV)	$E_0 + \Delta_0$ (eV)	E_g (eV)	$E_g + \Delta_1$ (eV)
Si	I	1.1242 ^(a-1)	1.169 ^(a-1)	1.170 ^(a-1)	-2.8 ^(a-2)	4.185(4.2 K) ^(a-3)	4.229(4.2 K) ^(a-3)	3.40 ^(a-4)					
Ge	I	0.664 ^(b-1)	0.734 ^(b-1)	0.744 ^(b-2)	-3.7 ^(b-3)	0.888(10 K) ^(b-4)	1.184(10 K) ^(b-4)	2.05 ^(b-5)				2.298 ^(b-6)	
α -Sn	D	0 ^(c-1)	0 ^(c-1)	0 ^(c-1)	0	-0.42(85 K) ^(c-2)	0.8(10 K) ^(c-2)	1.316 ^(c-3)				1.798 ^(c-3)	
GaAs	D	1.424 ^(d-1)	1.5115 ^(d-2)	1.51914 ^(d-3)	-3.9 ^(d-4)	E_g	1.760 ^(d-5)	2.915 ^(d-6)				3.139 ^(d-6)	
AlAs	I	2.153 ^(e-1)	2.223 ^(e-1)	2.229 ^(e-1)	-3.6 ^(e-1)	3.02 ^(e-2)	3.32 ^(e-2)	~ 3.9 ^(e-2)				~ 4.1 ^(e-2)	
InAs	D	0.354 ^(f-1)	0.404 ^(f-2)	0.418 ^(f-3)	-3.5 ^(f-4)	E_g	0.725 ^(f-1)	2.5 ^(f-5)				2.75 ^(f-5)	
InP	D	1.344 ^(g-1)	1.4135 ^(g-2)	1.4236 ^(g-3)	-2.9 ^(g-2)	E_g	1.45 ^(g-4)	3.158 ^(g-5)				3.28 ^(g-6)	
InSb	D	0.18 ^(h-1)	0.23 ^(h-2)	0.2368 ^(h-2)	-2.7 ^(h-2)	E_g	1.16 ^(h-3)	1.88 ^(h-4)				2.38 ^(h-4)	
GaP	I	2.272 ⁽ⁱ⁻¹⁾	2.338 ⁽ⁱ⁻²⁾	2.350 ⁽ⁱ⁻¹⁾	-3.7 ⁽ⁱ⁻²⁾	2.780 ⁽ⁱ⁻³⁾	2.860 ⁽ⁱ⁻³⁾	3.785(10 K) ⁽ⁱ⁻³⁾				3.835(10 K) ⁽ⁱ⁻³⁾	
ZnS (cubic)	D	3.68 ^(j-1)	3.78 ^(j-1)	3.78 ^(j-1)	-4.7 ^(j-2)	E_g	3.752(15 K) ^(j-3)	5.73 ^(j-1)				4.97 ^(k-1)	
ZnSe	D	2.70 ^(k-1)	2.8215 ^(k-2)	2.8215 ^(k-2)	-4.8 ^(k-3)	E_g	3.12 ^(k-1)	4.77 ^(k-1)				4.22 ^(l-3)	
ZnTe	D	2.30 ^(l-1)	2.3941 ^(l-2)	2.3941 ^(l-2)	-4.1 ^(l-1)	E_g	3.18 ^(l-3)	3.64 ^(l-3)				3.87 ^(m-5)	
CdTe	D	1.505 ^(m-1)	1.583 ^(m-2)	1.6063 ^(m-3)	-2.9 ^(m-4)	E_g	2.1 ^(m-5)	3.31 ^(m-6)				2.78 ⁽ⁿ⁻²⁾	
HgTe	D	0 ⁽ⁿ⁻¹⁾	0 ⁽ⁿ⁻¹⁾	0 ⁽ⁿ⁻¹⁾	0	-0.106 ⁽ⁿ⁻¹⁾	1.08 ⁽ⁿ⁻²⁾	2.12 ⁽ⁿ⁻²⁾					
CdS (Hexagonal)	D	2.485 ^(o-1)	2.573 ^(o-1)	2.5825 ^(o-2)	+4.1 ^(o-1)	E_g		1.85 ^(p-4)					
Pbs	D	0.42 ^(p-1)	0.307 ^(p-2)	0.286 ^(p-3)	+5.2 ^(p-3)	E_g		1.24 ^(q-4)					
PbTe	D	0.311 ^(q-1)	0.217 ^(q-2)	0.188 ^(q-1)	+4.5 ^(q-3)	E_g		1.59 ^(r-5)					
PbSe	D	0.278 ^(r-1)	0.176 ^(r-2)	0.1463 ^(r-3)	+4.0 ^(r-4)	E_g							
GaN (hexagonal)	D	3.39 ^(s-1)	3.503 ^(s-2)	3.503 ^(s-2)		E_g							
AlN (hexagonal)	D	6.13 ^(t-1)	6.19 ^(t-1,2)	6.19 ^(t-1,2)		E_g							
InN (hexagonal)	D		1.89 ^(u-1)	1.89 ^(u-1)		E_g							

Material	E_g (FE) (eV)	$\epsilon(0)$	ϵ_∞	$n(\lambda)$ $\lambda > \lambda_c$	dn/dT (10^{-4} /K)	$\hbar\omega_{LO}$ (eV)	$\hbar\omega_{LO}$ (eV)	μ_c^2 (cm ² /V.s)	μ_v^2 (cm ² /V.s)	m_e^* / m_0	m_h^* / m_0
Si	0.014 ^(a-5)	11.9 ^(a-6)	11.9 ^(a-6)	3.4179(10 μ m) ^(a-7)	1.3 ^(a-8)	0.0642 ^(a-9)	0.0642 ^(a-9)	1,500 ^(a-10)	450 ^(a-10)	0.98 ^(a-10)	0.16 ^(a-10)
Ge	0.00415 ^(b-7)	16.2 ^(b-8)	16.2 ^(b-8)	4.00319(10 μ m) ^(b-9)	4.0 ^(b-10)	0.0373 ^(b-11)	0.0373 ^(b-11)	3,900 ^(b-12)	1900 ^(b-12)	0.194 ^(a-10)	0.49 ^(a-10)
α -Sn	0	24 ^(c-4)	24 ^(c-4)	—	—	0.0244 ^(c-5)	0.0244 ^(c-5)	1,400 ^(c-6)	1200 ^(c-6)	1.64 ^(b-12)	0.04 ^(b-12)
GaAs	0.0042 ^(d-7)	13.18 ^(d-8)	10.89 ^(d-8)	3.298(5 μ m) ^(d-9)	1.5 ^(d-10)	0.0333 ^(d-11)	0.0362 ^(d-11)	8,500 ^(d-12)	400 ^(d-12)	0.082 ^(b-13)	0.28 ^(b-12)
AlAs	0.02 ^(e-1)	10.06 ^(e-3)	8.16 ^(e-3)	2.87(2 μ m) ^(e-4)	1.2 ^(e-5)	0.04488 ^(e-6)	0.05009 ^(e-6)	300 ^(e-7)	200 ^(e-8)	0.067 ^(d-12)	0.082 ^(d-12)
InAs	0.0017 ^(f-3)	15.15 ^(f-6)	12.25 ^(f-6)	3.42(10 μ m) ^(f-7)	0.83 ^(g-9)	0.0269 ^(f-8)	0.0296 ^(f-8)	3,300 ^(f-8)	460 ^(f-8)	1.1 ^(g-2)	0.153 ^(e-9)
InP	0.0051 ^(g-3)	12.61 ^(g-7)	9.61 ^(g-7)	3.08(5 μ m) ^(g-8)	0.83 ^(g-9)	0.0377 ^(g-10)	0.0428 ^(g-10)	4,600 ^(g-11)	150 ^(g-11)	0.023 ^(f-8)	0.4 ^(f-8)

(Continued)

Material	E_{∞} (FE) (eV)	$\epsilon(0)$	ϵ_{∞}	$n(\lambda)$ $\lambda > \lambda_c$	dn/dT ($10^{-4}/K$)	$\hbar\omega_{LO}$ (eV)	$\hbar\omega_{LO}$ (eV)	μ_c (cm ² /V·s)	μ_n (cm ² /V·s)	m_e^*/m_0^*	m_h^*/m_0^*
InSb	0.00052 ^(b-5)	16.8 ^(b-6)	15.68 ^(b-7)	3.953(10 μm) ^(b-8)	4.7 ^(b-9)	0.0223 ^(b-10)	0.0237 ^(b-10)	80,000 ^(b-11)	1250 ^(b-11)	0.0145 ^(b-11)	0.40 ^(b-11)
GaP	≈0.02 ^(c-1)	11.1 ^(c-4)	9.11 ^(c-4)	2.90(10 μm) ^(c-5)	—	0.0455 ^(c-6)	0.050 ^(c-6)	110 ^(c-7)	75 ^(c-7)	0.82 ^(c-7)	0.60 ^(c-7)
ZnS (cubic)	0.036 ^(c-4)	5.1 ^(c-5)	2.95 ^(c-6)	2.2014(10 μm) ^(c-7)	0.46 ^(c-8)	0.03397 ^(c-9)	0.04364 ^(c-10)	165 ^(c-10)	5 ^(c-11)	0.28 ^(c-12)	0.49 ^(c-12)
ZnSe	0.018 ^(b-4)	9.6 ^(b-5)	6.3 ^(b-5)	2.410(10 μm) ^(b-6)	0.52 ^(b-7)	0.02542 ^(b-8)	0.03099 ^(b-8)	600 ^(b-9)	100 ^(b-10)	0.142 ^(b-11)	0.57 ^(b-12)
ZnTe	0.0132 ^(c-2)	10.1 ^(c-4)	7.28 ^(c-5)	2.64(5 μm) ^(c-6)	—	0.02194 ^(c-7)	0.02542 ^(c-7)	330 ^(c-8)	—	0.11 ^(c-9)	0.6 ^(c-10)
CdTe	0.0132 ^(m-3)	10.2 ^(m-7)	7.1 ^(m-7)	2.684(5 μm) ^(m-8)	—	0.0174 ^(m-7)	0.0208 ^(m-9)	1,050 ^(m-10)	100 ^(m-10)	0.1 ^(m-11)	0.4 ^(m-12)
HgTe	0	21.0 ^(c-3)	15.2 ^(c-3)	—	—	0.01463 ^(c-4)	0.0171 ^(c-4)	35,000 ^(c-5)	—	0.031 ^(c-6)	0.42 ^(c-7)
CdS (Hexagonal)	0.0274 ^(b-3)	8.7 $(\epsilon_{11}(0))^{(b-4)}$	5.53 $(\epsilon_{1\infty})^{(b-4)}$	2.227(10 μm) $\downarrow c^{(b-5)}$	0.6 ^(b-5)	see ^(c-6)	see ^(c-6)	340 ^(c-7)	50 ^(c-7)	0.21 ^(c-7)	0.8 ^(c-7)
PbS	—	169 ^(p-5)	17.2 ^(p-5)	5.5 $(\epsilon_{\infty})^{(b-4)}$	0.62 ^(b-5)	—	—	—	—	—	—
PbTe	—	41.4 ^(b-5)	33 ^(b-6)	2.245(10 μm) $\downarrow c^{(b-5)}$	—	—	—	—	—	—	—
PbSe	—	210 ^(c-6)	22.9 ^(c-7)	5.66(10 μm) ^(c-7)	—	—	—	—	—	—	—
GaN (Hexagonal)	—	8.9 ^(c-3)	5.35 ^(c-3)	4.75(10 μm) ^(c-8)	—	—	—	—	—	—	—
AlN (Hexagonal)	—	8.5 ^(c-3)	4.77 ^(c-3)	—	—	—	—	—	—	—	—
InN (Hexagonal)	—	15.3 ^(c-3)	8.4 ^(c-3)	—	—	—	—	—	—	—	—

¹Most values quoted are at ≈300 K, except where indicated.

²Highest values reported.

³Longitudinal effective mass.

⁴Transverse effective mass.

⁵Light-hole effective mass.

⁶Heavy-hole effective mass.

⁷Band extrema effective masses are obtained from low-temperature measurements. See quoted references for more details.

a-Si

- (a-1) W. Bludau, A. Onton, and W. Heimke, *J. Appl. Phys.* **45**:1846 (1974).
(a-2) H. D. Barber, *Sol. St. Electronics* **10**:1039 (1967).
(a-3) D. E. Aspnes and A. A. Studna, *Solid State Communications* **11**:1375 (1972).
(a-4) A. Duanois and D. E. Aspnes, *Phys. Rev. B* **18**:1824 (1978).
(a-5) K. L. Shaklee and R. E. Nahory, *Phys. Rev. Lett.* **24**:942 (1970).
(a-6) K. V. Rao and A. Smaekal, *J. Appl. Phys.* **37**:2840 (1966).
(a-7) C. D. Salzberg and J. J. Villa, *J. Opt. Soc. Am.* **47**:244 (1957).
(a-8) M. Cardona, W. Paul, and H. Brooks, *J. Phys. Chem. Solids* **8**:204 (1959).
(a-9) G. Dollinger, *Inelastic Scattering of Neutrons in Solids and Liquids*, vol. II, IAEA, Vienna, 1963, p. 37.
(a-10) S. M. Sze, *Physics of Semiconductor Devices*, 2nd ed., John Wiley, New York, 1981, p. 849.

b-Ge

- (b-1) G. G. Macfarlane, T. P. McLean, J. E. Quarrington, and V. Roberts, *Phys. Rev.* **108**:1377 (1957).
(b-2) S. Zwerdling, B. Lax, L. M. Roth, and K. J. Button, *Phys. Rev.* **114**:880 (1959).
(b-3) T. P. McLean, in *Progress in Semiconductors*, vol. 5, A. F. Gibson (ed.), Heywood, London, 1960.
(b-4) D. E. Aspnes, *Phys. Rev. B* **12**:2297 (1975).

- (b⁵) A. K. Gosh, *Phys. Rev.* **165**:888 (1968).
 (b⁶) L. Vina and M. Cardona, *Physica* **117B** and **118B**:356 (1983).
 (b⁷) V. I. Sidorov and Ya. E. Pokrowski, *Sov. Phys. Semicond.* **6**:2015 (1973).
 (b⁸) E. A. D'Alroy and H. Y. Fan, *Phys. Rev.* **103**:671 (1956).
 (b⁹) R. P. Edwin, M. T. Dudermeil, and M. Lamare, *Appl. Opt.* **21**:878 (1982).
 (b¹⁰) See (a-8).
 (b¹¹) G. Nilsson and G. Nelin, *Phys. Rev. B* **6**:3777 (1972).
 (b¹²) See (a-10).

c- α Sn

- (c¹) J. R. Chelikowsky and M. L. Cohen, *Phys. Rev. B* **14**:556 (1976).
 (c²) S. H. Groves, C. R. Pidgeon, A. W. Ewald, and R. J. Wagner, *J. Phys. Chem. Solids* **31**:2031 (1970).
 (c³) L. Vina, H. F. Hockst, and M. Cardona, *Phys. Rev. B* **31**:958 (1985).
 (c⁴) R. E. Lindquist and A. W. Ewald, *Phys. Rev.* **135**:A191 (1964).
 (c⁵) C. J. Buchenauer, M. Cardona, and F. H. Pollak, *Phys. Rev. B* **3**:1243 (1971).
 (c⁶) See (a-10).

d-GaAs

- (d¹) D. D. Sell, H. C. Casey, and K. W. Wecht, *J. Appl. Phys.* **45**:2650 (1974).
 (d²) M. B. Panish and H. C. Casey, *J. Appl. Phys.* **40**:163 (1969).
 (d³) B. J. Skromme and G. E. Stillman, *Phys. Rev. B* **29**:1982 (1984).
 (d⁴) J. Camassel, D. Auvergne, and H. Mathieu, *J. Appl. Phys.* **46**:2683 (1975).
 (d⁵) C. Allibert, S. Gaillard, M. Erman, and P. M. Frijlink, *J. Phys. Paris* **44**:C10-229 (1983).
 (d⁶) P. Lautenschlager, M. Garriga, S. Logothetidis, and M. Cardona, *Phys. Rev. B* **35**:9174 (1987).
 (d⁷) D. D. Sell, *Phys. Rev. B* **6**:3750 (1972).
 (d⁸) G. A. Samara, *Phys. Rev. B* **27**:3494 (1983).
 (d⁹) K. G. Hambleton, C. Hilsaum, and B. R. Holeman, *Proc. Phys. Soc.* **77**:1147 (1961).
 (d¹⁰) M. Cardona, *Proc. Int. Conf. Phys. Semicond.*, Prague, 1960, Publ. House of the Czech. Acad. of Sciences, Prague, 1960, p. 388.
 (d¹¹) A. Mooradian and G. B. Wright, *Solid State Communications* **4**:431 (1966).
 (d¹²) See (a-10).

e-AlAs

- (e¹) B. Monemar, *Phys. Rev. B* **8**:5711 (1983).
 (e²) S. Adachi, *J. Appl. Phys.* **58**:R1 (1985).
 (e³) R. E. Fern and A. Onton, *J. Appl. Phys.* **42**:3499 (1971).
 (e⁴) M. Hoch and K. S. Hings, *J. Chem. Phys.* **35**:451 (1961).
 (e⁵) H. G. Grimmeiss and B. Monemar, *Phys. Stat. Sol. (a)* **5**:109 (1971).
 (e⁶) O. K. Kim and W. G. Spitzer, *J. Appl. Phys.* **50**:4362 (1979).
 (e⁷) M. Ettenberg, A. G. Sigai, A. Dreeben, and S. L. Gilbert, *J. Electrochem. Soc.* **119**:1355 (1971).
 (e⁸) J. D. Wiley, in *Semiconductors and Semimetals*, vol. 10, R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1975, p. 91.
 (e⁹) M. Huang and W. Y. Ching, *J. Phys. Chem. Solids* **46**:977 (1985).

f-InAs

- (f¹) F. Lukes, *Phys. Stat. Sol. (b)* **84**:K113 (1977).
 (f²) E. Adachi, *J. Phys. Soc. Jpn.* **2**:1178 (1968).
 (f³) A. V. Varfolomeev, R. P. Seisyan, and R. N. Yakimova, *Sov. Phys. Semicond.* **9**:560 (1975).
 (f⁴) F. Matossi and F. Stern, *Phys. Rev.* **111**:472 (1958).
 (f⁵) M. Cardona and G. Harbeke, *J. Appl. Phys.* **34**:813 (1963).
 (f⁶) M. Haas and B. W. Hennis, *J. Phys. Chem. Solids* **23**:1099 (1962).
 (f⁷) O. G. Lorimor and W. G. Spitzer, *J. Appl. Phys.* **36**:1841 (1965).
 (f⁸) R. Carles, N. Saint-Cricq, J. B. Renucci, M. A. Renucci, and A. Zwick, *Phys. Rev. B* **22**:4804 (1980).

g-InP

- (g¹) M. Buganski and W. Lewandowski, *J. Appl. Phys.* **57**:521 (1985).
 (g²) W. J. Turner, W. E. Reese, and G. D. Pettit, *Phys. Rev.* **136**:A1467 (1964).

- (f⁹) H. Mathieu, Y. Chen, J. Camassel, J. Allegre, and D. S. Robertson, *Phys. Rev. B* **32**:4042 (1985).
 (f¹⁰) K. L. Shaklee, M. Cardona, and F. H. Pollak, *Phys. Rev. Lett.* **16**:48 (1966).
 (f¹¹) S. M. Kelso, D. E. Aspnes, M. A. Pollak, and R. E. Nahory, *Phys. Rev. B* **26**:6669 (1982).
 (f¹²) F. Mataragui, A. E. Thompson, and M. Cardona, *Phys. Rev.* **176**:950 (1968).
 (f¹³) See (f-6).
 (f¹⁴) F. Oswald, *Z. Naturforsch* **9a**:181 (1954).
 (f¹⁵) See (d-10).
 (f¹⁶) A. Mooradian and G. B. Wright, *Solid State Communications* **4**:431 (1966).
 (f¹⁷) See (a-10).

h-InSb

- (h¹) F. Lukes and E. Schmidt, *Proc. Int. Conf. Phys. Semicond.*, Exeter, 1962, Inst. of Physics, London, 1962, p. 389.
 (h²) C. L. Littler and D. G. Sella, *Appl. Phys. Lett.* **46**:986 (1985).
 (h³) S. Zwerding, W. H. Kleiner, and J. P. Theriault, *MIT Lincoln Laboratory Report 8C-00M*, 1961.
 (h⁴) M. Cardona, K. L. Shaklee, and F. H. Pollak, *Phys. Rev.* **154**:696 (1967).
 (h⁵) A. Baldereschi and N. O. Lipari, *Phys. Rev. B* **3**:439 (1971).
 (h⁶) J. R. Dixon and J. K. Furdyna, *Solid State Communications* **35**:195 (1980).
 (h⁷) See (f-6).
 (h⁸) T. S. Moss, S. D. Smith, and T. D. F. Hawkins, *Proc. Phys. Soc. B* **70**:776 (1957).
 (h⁹) See (g-9).
 (h¹⁰) W. Kefler, W. Richter, and M. Cardona, *Phys. Rev. B* **12**:2346 (1975).
 (h¹¹) See (a-10).

i-GaP

- (i¹) R. G. Humphreys, U. Rossler, and M. Cardona, *Phys. Rev. B* **18**:5590 (1978).
 (i²) D. Auvergne, P. Merle, and H. Mathieu, *Phys. Rev. B* **12**:1371 (1975).
 (i³) S. E. Stokowski and D. D. Sell, *Phys. Rev. B* **5**:1636 (1972).
 (i⁴) G. A. Samara, *Phys. Rev. B* **27**:3494 (1983).
 (i⁵) H. Welker, *J. Electron* **1**:181 (1955).
 (i⁶) See (g-10).
 (i⁷) See (a-10).

j-ZnS

- (j¹) D. Theis, *Phys. Stat. Sol. (b)* **79**:125 (1977).
 (j²) J. Camassel and D. Auvergne, *Phys. Rev. B* **12**:3258 (1975).
 (j³) B. Segall and D. T. F. Marple, *Physics and Chemistry of II-VI Compounds*, M. Aven and J. S. Prener (eds.), North-Holland, Amsterdam, 1967, p. 318.
 (j⁴) W. Walter and J. L. Birman, *Proc. Int. Conf. on II-VI Semiconducting Compounds*, 1967, D. G. Thomas (ed.), W. A. Benjamin, New York, 1967, p. 89.
 (j⁵) G. Martinez, in *Handbook on Semiconductors*, vol. 2, T. S. Moss (ed.), North-Holland, Amsterdam, 1980, p. 210.
 (j⁶) M. Balkanski and Y. Petroff, *Proc. 7th Intl. Conf. Physics of Semicond.*, Paris, 1964, Dunod, Paris, 1964, p. 245.
 (j⁷) W. L. Wolfe, A. G. DeBell, and J. M. Palmer, *Proc. SPIE* **245**:164 (1980).
 (j⁸) R. J. Harris, G. T. Johnston, G. A. Kepple, P. C. Krock, and M. Mukai, *Appl. Opt.* **16**:436 (1977).
 (j⁹) M. Balkanski, M. Nusimovici, and R. Letoullec, *J. Phys. Paris* **25**:305 (1964).
 (j¹⁰) C. A. Klein and R. N. Donadio, *J. Appl. Phys.* **51**:797 (1980).
 (j¹¹) See (a-10).
 (j¹²) J. C. Miklosz and R. G. Wheeler, *Phys. Rev.* **153**:913 (1967).

k-ZnSe

- (k¹) See (j-1).
 (k²) P. J. Dean, D. C. Herbert, C. J. Werkhoven, B. J. Fitzpatrick, and R. M. Bhargava, *Phys. Rev. B* **23**:4888 (1981).
 (k³) J. Baillou, P. Bugnet, Jac Daumay, C. Auzary, and P. Poindessault, *J. Phys. Chem. Solids* **41**:295 (1980).
 (k⁴) A. K. Ray and F. A. Kroger, *J. Appl. Phys.* **50**:4208 (1979).
 (k⁵) A. Hadni, J. Claudel, and P. Strimmer, *Phys. Stat. Sol.* **26**:241 (1968).
 (k⁶) X. J. Jiang, T. Hisamura, Y. Nossua, and T. Goto, *J. Phys. Soc. Jpn* **52**:4008 (1983).
 (k⁷) See (j-8).

- (k-8) M. Cardona, *J. Phys. Paris* **C8**:29 (1984).
- (k-9) T. Yao, M. Ogura, S. Matsuoka, and T. Morishita, *J. Appl. Phys.* **43**:499 (1983).
- (k-10) G. Jones and J. Woods, *J. Phys. D* **9**:799 (1976).
- (k-11) T. Ohyama, E. Otsuka, T. Yoshida, M. Ishiki, and K. Igaki, *Jpn. J. Appl. Phys.* **23**:L382 (1984).
- (k-12) M. Sondergeld, *Phys. Stat. Sol. (b)* **81**:253 (1977).

I-ZnTe

- (l-1) See (j-1).
- (l-2) M. Venghaus and P. J. Dean, *Phys. Rev. B* **21**:1596 (1980).
- (l-3) See (h-5).
- (l-4) D. Berlincourt, M. Jaffe, and L. R. Shiozawa, *Phys. Rev.* **129**:1009 (1983).
- (l-5) D. T. F. Marple, *J. Appl. Phys.* **35**:539 (1964).
- (l-6) T. L. Chu, S. S. Chu, F. Firszt, and C. Herrington, *J. Appl. Phys.* **59**:1259 (1986).
- (l-7) See (k-8).
- (l-8) A. G. Fisher, J. N. Cardes, and J. Dresner, *Solid State Commun.* **21**:157 (1964).
- (l-9) H. Venghaus, P. J. Dean, P. E. Simmonds, and J. C. Pfister, *Z. Phys. B* **30**:125 (1978).
- (l-10) M. Aven and B. Segall, *Phys. Rev.* **130**:81 (1963).

m-CdTe

- (m-1) P. M. Amirtharaj and D. Chandler-Horowitz, (unpublished).
- (m-2) P. M. Amirtharaj, R. C. Bowman, Jr., and R. L. Alt, *Proc. SPIE* **946**:57 (1988).
- (m-3) N. Nawrocki and A. Twardowski, *Phys. Stat. Sol. (b)* **97**:K61 (1980).
- (m-4) See (j-2).
- (m-5) See (l-3).
- (m-6) A. Moritani, K. Tamiguchi, C. Hamaguchi, and J. Nakai, *J. Phys. Soc. Jpn.* **34**:79 (1973).
- (m-7) T. J. Parker, J. R. Birch, and C. L. Mok, *Solid State Communications* **36**:581 (1980).
- (m-8) L. S. Ladd, *Infrared Phys.* **6**:145 (1966).
- (m-9) J. R. Birch and D. K. Murray, *Infrared Phys.* **18**:283 (1978).
- (m-10) See (a-10).
- (m-11) K. K. Kanazawa and F. C. Brown, *Phys. Rev.* **135**:A1757 (1964).
- (m-12) See (j-3).

n-HgTe

- (n-1) W. Szauskiewicz, *Phys. Stat. Sol. (b)* **81**:K119 (1977).
- (n-2) See (m-6).
- (n-3) J. Baars and F. Sorger, *Solid State Commun.* **10**:875 (1972).
- (n-4) H. Kapa, T. Giebultowicz, B. Buras, B. Lebeck, and K. Clausen, *Physica Scripta* **25**:807 (1982).
- (n-5) T. C. Harmon, in *Physics and Chemistry of II-VI Compounds*, M. Aven and J. S. Prener (eds), North Holland Publishing, Amsterdam, 1967, p. 767.
- (n-6) J. Guldner, C. Rigaux, M. Grynborg, and A. Mycielski, *Phys. Rev. B* **8**:3875 (1973).
- (n-7) K. Shinzui, S. Narita, Y. Nisida, and V. I. Ivanov-Omskii, *Solid State Commun.* (eds.), **32**:327 (1979).

o-Cds

- (o-1) V. V. Sobolev, V. I. Donetskina, and E. F. Zagainov, *Sov. Phys. Semicond.* **12**:646 (1978).
- (o-2) D. G. Seiler, D. Heiman, and B. S. Wherrett, *Phys. Rev. B* **27**:2355 (1983).
- (o-3) A series excitons: D. G. Seiler, D. Heiman, R. Fiegenblatt, R. Aggarwal, and B. Lax, *Phys. Rev. B* **25**:7666 (1982); B series excitons: see (o-2).
- (o-4) A. S. Barker and C. J. Summers, *J. Appl. Phys.* **41**:3552 (1970).
- (o-5) R. Weil and D. Neshmit, *J. Opt. Soc. Am.* **67**:190 (1977).
- (o-6) Complex Phonon Structure with Nine Allowed Optical Modes. See B. Tel, T. C. Damen, and S. P. S. Porto, *Phys. Rev.* **144**:771 (1966).
- (o-7) See (a-10).

p-PbS

- (p-1) R. B. Schoolar and J. R. Dixon, *Phys. Rev. A* **137**:667 (1965).
- (p-2) D. L. Mitchell, E. D. Palik, and J. N. Zemel, *Proc. 7th Int. Conf. Phys. Semicond.*, Paris, 1964, Dunod, Paris, 1964, p. 325.
- (p-3) G. Nimtz and B. Schlichter, *Springer Tracts in Modern Physics*, vol. 98, Springer-Verlag, Berlin, 1983, p. 1.

- (p-1) M. Cardona and D. L. Greenaway, *Phys. Rev. A* **13**:1685 (1964).
 (p-2) R. Davlen, in *Solid State Physics*, vol. 28, H. Ehrenreich, F. Seitz, and D. Turnbull (eds.), Academic, NY, 1973, p. 179.
 (p-3) R. B. Schoolar and J. N. Zemel, *J. Appl. Phys.* **35**:1848 (1964).
 (p-4) M. M. Elcombe, *Proc. Soc. London*, **A300**:210 (1967).
 (p-5) See (a-10).
 (p-6) K. F. Cuff, M. R. Ellet, C. D. Kulgin, and L. R. Williams, in *Proc. 7th Int. Conf. Phys. Semicond.*, Paris, 1964, M. Hulin (ed.), Dunod, Paris, 1964, p. 677.
- q-PbTe**
 (q-1) M. Preier, *Appl. Phys.* **20**:189 (1979).
 (q-2) C. R. Hewes, M. S. Adler, and S. D. Senturia, *Phys. Rev. B* **7**:5195 (1973).
 (q-3) See (p-2).
 (q-4) See (p-4).
 (q-5) W. E. Tennant, *Solid State Communications* **20**:613 (1976).
 (q-6) J. R. Lowney and S. D. Senturia, *J. Appl. Phys.* **47**:1773 (1976).
 (q-7) N. Piccioli, J. B. Beson, and M. Balkanski, *J. Phys. Chem. Solids* **35**:971 (1974).
 (q-8) W. Cochran, R. A. Cowley, G. Dolling, and M. M. Elcombe, *Proc. R. Soc. London*, **A293**:433 (1966).
 (q-9) See (a-10).
 (q-10) See (p-9).
- r-PbSe**
 (r-1) U. Schlögl, Dissertation Technische Universität Berlin, 1970.
 (r-2) D. L. Mitchell, E. D. Palik, and J. N. Zemel, *Proc. 7th Int. Conf. Phys. Semicond.*, Paris, 1964, M. Hulin (ed.), Dunod, Paris, 1964, p. 325.
 (r-3) H. Pasher, G. Bauer, and R. Grisar, *Phys. Rev. B* **38**:3383 (1988).
 (r-4) A. F. Gibson, *Proc. Phys. Soc. (London)* **B65**:378 (1952).
 (r-5) See (p-4).
 (r-6) See (p-5).
 (r-7) J. N. Zemel, J. D. Jensen, and R. B. Schoolar, *Phys. Rev. A* **140**:330 (1965).
 (r-8) H. Burkhard, R. Geick, P. Kastner, and K. H. Unkelbach, *Phys. Stat. Sol. (b)* **63**:89 (1974).
 (r-9) E. Burstein, R. Wheeler, and J. Zemel, *Proc. 7th Int. Conf. Phys. Semicond.*, Paris, 1964, M. Hulin (ed.), Dunod, Paris, 1964, p. 1065.
 (r-10) R. N. Hall and J. H. Racette, *J. Appl. Phys.* **32**:2078 (1961).
 (r-11) J. N. Zemel, J. D. Jensen, and R. B. Schoolar, *Phys. Rev. A* **140**:330 (1965).
 (r-12) U. Schlögl and K. H. Gobrecht, *J. Phys. Chem. Solids* **34**:753 (1973).
 (r-13) See (p-9).
- s-GaN**
 (s-1) H. P. Maruska and J. J. Tietjen, *Appl. Phys. Lett.* **15**: 327 (1969); For a discussion of band parameters in III-Nitrides, see also I. Vurgatman and J. R. Meyer, *J. Appl. Phys.* **94**: 3675 (2003).
 (s-2) O. Madelung, *Semiconductor Data Handbook*, 3rd ed., Springer, New York, 2004, p. 103.
 (s-3) S. N. Mohammad and H. Morfoc, *Prog. Quant. Electron.* **20**: 361 (1996).
 (s-4) U. V. Bhapkar and M. S. Shur, *J. Appl. Phys.* **82**: 1649 (1997).
 (s-5) M. S. Shur and M. Asif Khan, *MRS. Bull.* **22**: 44 (1997).
 (s-6) W. R. L. Lambrecht, K. Kim, S. N. Rashkeev, and B. Segall, Electronic and Optical Properties Group-III Nitrides, their Heterostructures and Alloys, *Mat. Res. Soc. Symp. Proc.* **395**: 455–466 (1996).
- t-AlN**
 (t-1) O. Madelung, *Semiconductor Data Handbook*, 3rd ed., Springer, New York, 2004, p. 88.
 (t-2) W. M. Yim, E. J. Stofo, P. J. Zanuzuchi, J. I. Pankove, M. Ettenburg, and S. L. Gilbert, *J. Appl. Phys.* **44**: 292 (1973).
 (t-3) V. W. L. Chin, T. L. Tansley, and T. Osotchan, *J. Appl. Phys.* **75**: 7365 (1994).
 (t-4) S. K. O'Leary, B. E. Foutz, M. S. Shur, U. V. Bhapkar, and L. F. Eastman, *Solid State Comm.* **105**: 621 (1998).
 (t-5) J. A. Majewski, M. Stadelé, and P. Vogl, Electronic Structure of Biaxially Strained Wurtzite Crystals GaN and AlN, *Mat. Res. Soc. Symp. Proc.* **449**: 887–892 (1997).
- u-InN**
 (u-1) T. L. Tansley and C. P. Foley, *J. Appl. Phys.* **59**: 3241 (1986).
 (u-2) O. Madelung, *Semiconductor Data Handbook*, 3rd ed., Springer, New York, 2004, p. 137.
 (u-3) S. K. O'Leary, B. E. Foutz, M. S. Shur, U. V. Bhapkar, and L. F. Eastman, *J. Appl. Phys.* **83**: 826 (1998).

5.7 REFERENCES

1. D. Attwood, B. Hartline, and R. Johnson, *The Advanced Light Source: Scientific Opportunities*, Lawrence Berkeley Laboratory Publication 5111, 1984, pp.331–389.
2. F. C. Brown, “Ultraviolet Spectroscopy of Solids with the Use of Synchrotron Radiation,” in *Solid State Physics*, vol. 29, H. Ehrenreich, F. Seitz, and D. Turnbull (eds.), Academic Press, New York, 1974, pp. 1–73.
3. P. Y. Yu and M. Cardona, *Fundamentals of Semiconductors*, Springer, Berlin, 1996.
4. A. V. Nurmikko, in *Semiconductors and Semimetals*, vol. 36, D. G. Seiler and C. L. Littler (eds.), Academic Press, New York, 1992, p. 85.
5. R. R. Alfano, (ed.), *Semiconductors Probed by Ultrafast Laser Spectroscopy*, vols. I and II, Academic Press, New York, 1984.
6. G. Bastard, C. Delalande, Y. Guldner, and P. Voisin, in *Advances in Electronics and Electron Physics*, vol. 72, P. W. Hawkes (ed.), Academic Press, New York, 1988, p. 1.
7. C. Weisbuch and B. Vinter, *Quantum Semiconductor Structures, Fundamentals and Applications*, Academic Press, New York, 1991, p. 57.
8. M. Born and K. Huang, *Dynamical Theory of Crystal Lattices*, chap. 2, Oxford University Press, London, 1954, p. 38.
9. W. K. H. Panofsky and M. Phillips, *Classical Electricity and Magnetism*, Addison-Wesley, New York, 1962, p. 29.
10. E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, vol. 1, Academic Press, New York, 1985.
11. E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, vol. 2, Academic Press, New York, 1991.
12. E. D. Palik (ed.), *Handbook of Optical Constants of Solids*, vol. 1, Academic Press, New York, 1985, p. 429.
13. M. Born and E. Wolf, *Principles of Optics*, Pergamon, London, 1970, p. 61.
14. J. S. Toll, “Causality and the Dispersion Relation: Logical Foundations,” *Phys. Rev.* **104**:1760–1770 (1956).
15. D. Y. Smith, “Comments on the Dispersion Relations for the Complex Refractive Index of Circularly and Elliptically Polarized Light,” *J. Opt. Soc. Am.*, **66**(5):454–460 (1976).
16. D. L. Greenaway and G. Harbeke, *Optical Properties and Band Structure of Semiconductors*, Pergamon, London, 1968, p. 9.
17. J. M. Ziman, *Principles of the Theory of Solids*, Cambridge University Press, London, 1972, p. 200.
18. C. Kittel, *Introduction to Solid State Physics*, 4th ed., John Wiley, New York, 1971, p. 184.
19. W. J. Turner and W. E. Reese, “Infrared Lattice Bands in AlSb,” *Phys. Rev.* **127**:126–131 (1962).
20. S. S. Mitra, in *Optical Properties of Solids*, S. Nudelman and S. S. Mitra (eds.), Plenum Press, New York, 1979, p. 333.
21. W. Cochran, *The Dynamics of Atoms in Crystals*, Crane, Rusak and Co., New York, 1973.
22. P. Gianozzi, S. de Gironcoli, P. Pavone, and S. Baroni, “*Ab initio* Calculation of Phonon Dispersions in Semiconductors,” *Phys. Rev. B* **43**:7231–7242 (1991).
23. S. S. Mitra and N. E. Massa, in *Handbook on Semiconductors*, T. S. Moss and W. Paul (eds.), North Holland, Amsterdam, 1982, p. 81.
24. W. G. Spitzer, in *Semiconductors and Semimetals*, vol. 3., R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1967, p. 17.
25. J. L. Birman, *Theory of Crystal Space Groups and Lattice Dynamics*, Springer-Verlag, Berlin, 1974, p. 271.
26. F. A. Johnson, “Lattice Absorption Bands in Silicon,” *Proc. Phys. Soc. (London)* **73**:265 (1959).
27. W. Cochran, S. J. Fray, F. A. Johnson, J. E. Quarrington, and N. Williams, “Lattice Absorption in Gallium Arsenide,” *J. Appl. Phys.* **32**:2102 (1961).
28. A. S. Barker and A. J. Sievers, “Optical Studies of the Vibrational Properties of Disordered Solids,” *Rev. Mod. Phys.* **47**(2):S1–S179 (1975).
29. W. M. Theis, K. K. Bajaj, C. W. Litton, and W. G. Spitzer, “Direct Evidence for the Site of Substitutional Carbon Impurity in GaAs,” *Appl. Phys. Lett.* **41**:70 (1982).
30. R. S. Leigh and R. C. Newman, “Host Isotope Fine Structure of Local Modes: C and Si in GaAs,” *J. Phys. C: Solid State Phys.* **15**:L1045 (1982).

31. M. Stavola and S. J. Pearton, in *Semiconductors and Semimetals*, vol. 34, J. I. Pankove and N. M. Johnson (eds.), Academic Press, New York, 1991, p. 139.
32. R. C. Newman, in *Growth and Characterization of Semiconductors*, R. A. Stradling and P. C. Klipstein (eds.), Adam Hilger, Bristol, 1990, p. 105.
33. A. Baghdadi, W. M. Bullis, M. C. Croarkin, Y. Li, R. I. Scace, R. W. Series, P. Stallhofer, and M. Watanabe, "Interlaboratory Determination of the Calibration Factor for the Measurement of the Interstitial Oxygen Content of Silicon by Infrared Absorption," *J. Electrochem. Soc.* **136**:2015 (1989).
34. W. L. Wolfe, in *The Infrared Handbook*, W. L. Wolfe and G. J. Zeiss (eds.), Environmental Research Institute, Ann Arbor, 1978, pp. 7–39.
35. A. Miller, *Handbook of Optics*, 2nd ed., vol. I, chap. 9, McGraw-Hill, New York, 1994.
36. E. J. Johnson, in *Semiconductors and Semimetals*, vol. 3, R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1967, p. 153.
37. E. O. Kane, "Band Structure of Indium Antimonide," *J. Phys. Chem. Solids* **1**:249 (1957).
38. T. S. Moss and T. D. F. Hawkins, "Infrared Absorption in Gallium Arsenide," *Infrared Phys.* **1**:111 (1961).
39. (op. cit.) E. J. Johnson, p. 191.
40. W. G. Spitzer, M. Gershenzon, C. J. Frosch, and D. F. Gibbs, "Optical Absorption in n-Type Gallium Phosphide," *J. Phys. Chem. Solids* **11**:339 (1959).
41. M. Gershenzon, D. G. Thomas, and R. E. Dietz, *Proc. Int. Conf. Phys. Semicond. Exeter*, Inst. of Physics, London, 1962, p. 752.
42. W. C. Dash and R. Newman, "Intrinsic Optical Absorption in Single-Crystal Germanium and Silicon at 77°K and 300°K," *Phys. Rev.* **99**:1151 (1955). See also G. Burns, *Solid State Physics*, Academic Press, New York, 1985, p. 505.
43. K. Cho, *Excitons*, vol. 14 of *Topics in Current Physics*, K. Cho (ed.), Springer-Verlag, New York, 1979, p. 1.
44. W. Hayes and A. M. Stoneham, *Defects and Defect Processes in Nonmetallic Solids*, John Wiley and Sons, New York, 1985, p. 40.
45. D. C. Reynolds and T. C. Collins, *Excitons: Their Properties and Uses*, Academic Press, New York, 1981, pp. 1–291.
46. E. I. Rashba and M. D. Sturge (eds.), *Excitons*, North Holland, Amsterdam, 1982, pp. 1–865.
47. J. Frenkel, "On the Transformation of Light into Heat in Solids II," *Phys. Rev.* **37**:1276–1294 (1931).
48. G. H. Wannier, "The Structure of Electronic Excitation Levels in Insulating Crystals," *Phys. Rev.* **52**:191–197 (1937).
49. N. F. Mott, "On the Absorption of Light by Crystals," *Proc. Roy. Soc. A* **167**:384 (1938).
50. M. D. Sturge, "Optical Absorption of Gallium Arsenide between 0.6 and 2.75 eV," *Phys. Rev.* **127**:768–773 (1962). For an excellent discussion of theoretical aspects see article by M. Sturge, "Advances in Semiconductor Spectroscopy," in B. DiBartolo (ed.), *Spectroscopy of Laser-Type Materials*, Plenum Press, New York, 1987, p. 267.
51. Adapted by permission from R. G. Ulbrich and C. Weisbuch, Contribution to the study of optical pumping in III-V Semiconductors, These de doctorat d'Etat, Univ. Paris 7, 1977 (Unpublished).
52. (op. cit.) W. Hayes and A. M. Stoneham, p. 51.
53. O. Madelung, *Introduction to Solid-State Theory*, Springer-Verlag, New York, 1978, p. 254.
54. G. Burns, *Solid State Physics*, Academic Press, New York, 1985, p. 969.
55. H. R. Philipp and H. Ehrenreich, in *Semiconductors and Semimetals*, vol. 3, R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1967, p. 93. For a review of interband transitions in Semiconductors see M. L. Cohen and J. R. Chelikowski, *Electronic Structure and Optical Properties of Semiconductors*, Springer-Verlag, Berlin, 1988.
56. H. R. Philipp and H. Ehrenreich, "Optical Properties of Semiconductors," *Phys. Rev.* **129**:1550–1560 (1963).
57. D. Brust, J. C. Phillips, and F. Bassani, "Critical Points and Ultraviolet Reflectivity of Semiconductors," *Phys. Rev. Lett.* **9**:94–97 (1962).
58. D. Brust, "Electronic Spectra of Crystalline Germanium and Silicon," *Phys. Rev.* **134**:A1337–A1353 (1964).
59. M. Rohlfing and S. G. Louie, "Electron-hole Excitations and Optical Spectra from First Principles," *Phys. Rev. B* **62**:4927–4944 (2000).

60. H. R. Philipp and H. Ehrenreich, "Observation of d Bands in 3-5 Semiconductors," *Phys. Rev. Lett.* **8**:92–94 (1962).
61. C. Gahwiller and F. C. Brown, "Photoabsorption near the $L_{II,III}$ Edge of Silicon and Aluminum," *Phys. Rev. B* **2**:1918–1925 (1970).
62. S. Zollner, in *Silicon-Germanium Carbon Alloys: Growth, Properties, and Applications*, S. T. Pantelides and S. Zollner (eds.), Taylor & Francis, New York, 2001, p. 387.
63. J. Humlicek, M. Garriga, M. I. Alonso, and M. Cardona, "Optical Spectra of $\text{Si}_x\text{Ge}_{1-x}$ Alloys," *J. Appl. Phys.* **65**, 2827–2832, 1989.
64. S. Perkowitz, in *Infrared and Millimeter Waves*, vol. 8, K. J. Button (ed.), Academic Press, New York, 1983, p. 71.
65. J. R. Dixon and H. R. Riedl, "Electric-Susceptibility Hole Mass of Lead Telluride," *Phys. Rev.* **138**:A873–A881 (1965).
66. T. E. Tiwald, J. A. Woollam, S. Zollner, J. Christiansen, R. B. Gregory, T. Wetteroth, S. R. Wilson, and A. R. Powell, "Carrier Concentration and Lattice Absorption in Bulk and Epitaxial Silicon Carbide Determined using Infrared Ellipsometry," *Phys. Rev. B* **60**:11464–11474 (1999).
67. B. Varga, "Coupling of Plasmons to Polar Phonons in Degenerate Semiconductors," *Phys. Rev.* **137**:A1896–A1902 (1965).
68. A. Mooradian and G. B. Wright, "Observation of the Interaction of Plasmons with Longitudinal Optical Phonons in GaAs," *Phys. Rev. Lett.* **16**:999–1001 (1966).
69. A. Mooradian, in *Advances in Solid State Physics*, vol. 9, O. Madelung (ed.), Pergamon Press, London, 1969, p. 74.
70. C. G. Olson and D. W. Lynch, "Longitudinal-Optical-Phonon-Plasmon Coupling in GaAs," *Phys. Rev.* **177**:1231–1234 (1969).
71. S. C. Baber, "Net and Total Shallow Impurity Analysis of Silicon by Low Temperature Fourier Transform Infrared Spectroscopy," *Thin Solid Films* **72**:201–210 (1980).
72. G. M. Martin, "Optical Assessment of the Main Electron Trap in Bulk Semi-Insulating GaAs," *Appl. Phys. Lett.* **39**: 747 (1981).
73. C. M. Wolfe and G. E. Stillman, *Gallium Arsenide and Related Compounds*, Inst. Phys., London, 1971, p. 3.
74. M. H. Weiler, *Semiconductors and Semimetals*, vol. 16, R. K. Willardson and A. C. Beer (eds.), Academic Press, New York, 1981, p. 119.
75. C. R. Pidgeon and R. N. Brown, "Interband Magneto-Absorption and Faraday Rotation in InSb," *Phys. Rev.* **146**:575 (1966).
76. M. H. Weiler, R. L. Aggarwal, and B. Lax, "Warping- and Inversion-Asymmetry-Induced Cyclotron-Harmonic Transitions in InSb," *Phys. Rev. B* **17**:3269 (1978).
77. M. H. Weiler, "Conduction and Valence Band Effective Mass Parameters in InSb," *J. Magn. Magn. Mater.* **11**:131–135 (1979).
78. M. Reine, R. L. Aggarwal, and B. Lax, "Stress-Modulated Magnetorefectivity of Gallium Antimonide and Gallium Arsenide," *Phys. Rev. B* **5**:3033 (1972).
79. M. Reine, R. L. Aggarwal, B. Lax, and C. M. Wolfe, "Split-Off Valence-Band Parameters for GaAs from Stress-Modulated Magnetorefectivity," *Phys. Rev. B* **2**:458 (1970).
80. H. Piller, *Proc. 7th Int. Conf. Phys. Semicond.*, Dunod, Paris, 1964, p. 297.
81. E. Oh, D. U. Bartholomew, A. K. Ramdas, J. K. Furdyna, and U. Debska, "Interband Faraday Effect in $\text{Cd}_{1-x}\text{Mn}_x\text{Se}$," *Phys. Rev. B* **38**:13183 (1988).
82. R. J. Elliott and R. Loudon, "Theory of the Absorption Edge in Semiconductors in a High Magnetic Field," *J. Phys. Chem. Solids* **15**:196–207 (1960).
83. D. G. Seiler, D. Heiman, R. Feigenblatt, R. L. Aggarwal, and B. Lax, "Two-Photon Magnetospectroscopy of A -Exciton States in CdS," *Phys. Rev. B* **25**:7666 (1982).
84. D. G. Seiler, D. Heiman, and B. S. Wherrett, "Two-Photon Spectroscopy of B Excitons in CdS," *Phys. Rev. B* **27**:2355 (1983).
85. B. D. McCombe and R. J. Wagner, *Adv. Electron. and Electron. Phys.* **37**:1 (1975).
86. B. D. McCombe and R. J. Wagner, *Adv. Electron. and Electron. Phys.* **38**:1 (1975).

87. H. Kobori, T. Ohyama, and E. Otsuka, "Line-Width of Quantum Limit Cyclotron Resonance. II. Impurity and Carrier-Carrier Scatterings in Ge, InSb and GaAs," *J. Phys. Soc. Jpn.* **59**:2164–2178 (1990).
88. B. Lax, H. J. Zeiger, and R. N. Dexter, "Anisotropy of Cyclotron Resonance in Germanium," *Physica* **20**:818–828 (1954).
89. G. Dresselhaus, A. F. Kip, and C. Kittel, "Cyclotron Resonance of Electrons and Holes in Silicon and Germanium Crystals," *Phys. Rev.* **98**:368 (1955).
90. M. A. Omar, *Elementary Solid State Physics*, Addison-Wesley, Reading, 1975, p. 285.
91. O. Matsuda and E. Otsuka, "Cyclotron Resonance Study of Conduction Electrons in n-Type Indium Antimonide under a Strong Magnetic Field—I: Thermal Equilibrium Case," *J. Phys. Chem. Solids* **40**:809–817 (1979).
92. D. Larsen, in *Landau Level Spectroscopy*, vol. 27.1, chap. 3, G. Landwehr and E. I. Rashba (eds.), North Holland, Amsterdam, 1991, p. 109.
93. H. R. Fetterman, D. M. Larsen, G. E. Stillman, P. E. Tannenwald, and J. Waldman, "Field-Dependent Central-Cell Corrections in GaAs by Laser Spectroscopy," *Phys. Rev. Lett.* **26**:975–978 (1971).
94. G. W. Bryant and G. S. Solomon, *Optics of Quantum Dots and Wires*, Artech House, Boston, 2005.
95. T. N. Muiira, *Physics of Semiconductors in High Magnetic Fields*, Oxford University Press, Oxford, 2008.
96. C. L. Tang, *Handbook of Optics*, 3rd ed., vol. IV, chap. 10, McGraw-Hill, New York, 2009.
97. Y. R. Shen, *The Principles of Nonlinear Optics*, John Wiley and Sons, New York, 1984, p. 38.
98. P. N. Butcher, *Nonlinear Optical Phenomena*, Ohio State University Engineering Publications, Columbus, 1965, p. 1.
99. B. S. Wherrett, *Nonlinear Optics*, P. G. Harper and B. S. Wherrett (eds.), Academic Press, New York, 1977, p. 4.
100. M.C. Downer, D.S. Mendoza, and V.I. Gavrilenko, "Optical Second Harmonic Spectroscopy of Semiconductor Surfaces: Advances in Microscopic Understanding," *Surf. Intef. Anal.* **31**:966–986 (2001).
101. R. Carriles, J. Kwon, Y. Q. An, J. C. Miller, M. C. Downer, J. Price, and A.C. Diebold, "Second-harmonic Generation from Si/SiO₂/Hf_(1-x)Si_xO₂ Structures," *Appl. Phys.Lett.* **88**:161120 (2006).
102. Y. Q. An, R. Carriles, and M. C. Downer, "Absolute Phase and Amplitude of Second-Order Nonlinear Optical Susceptibility Components at Si(001) Interfaces," *Phys. Rev. B* **75**:241307 (2007).
103. J. H. Bechtal and W. L. Smith, "Two-Photon Absorption in Semiconductors with Picosecond Laser Pulses," *Phys. Rev. B* **13**:3515–3522 (1976).
104. P. D. Maker and R. H. Terhune, "Study of Optical Effects Due to an Induced Polarization Third Order in the Electric Field Strength," *Phys. Rev.* **137**:A801–A818 (1965).
105. R. M. A. Azzam and N. M. Bashara, *Ellipsometry and Polarized Light*, North-Holland, Amsterdam, 1987, p. 153. See also Azzam in *Handbook of Optics*, 2nd ed., vol. II, chap. 27, McGraw-Hill, New York, 1994.
106. D. E. Aspnes, "Analysis of Semiconductor Materials and Structures by Spectroellipsometry," *Proc. SPIE* **946**:84 (1988).
107. W. Hayes and R. Loudon, *Scattering of Light by Crystals*, John-Wiley, New York, 1978, p. 53.
108. P. J. Dean, "Photoluminescence as a Diagnostic of Semiconductors," *Prog. Crystal Growth and Characterization* **5**:89 (1982).
109. M. Voos, R. F. Lehney, and J. Shah, in *Handbook of Semiconductors*, vol. 2, T. S. Moss and M. Balkanski (eds.), North-Holland, Amsterdam, 1980, p. 329.
110. See Table 11.
111. A.J. LaRocca, in *Handbook of Optics*, 3rd ed., vol. II, chap. 15, McGraw-McGraw-Hill, New York, 2009.
112. W.T. Silfvast, in *Handbook of Optics*, 3rd ed., vol. II, chap. 16, McGraw-Hill, New York, 2009.
113. B. Henderson in *Handbook of Optics*, 3rd ed., vol. I, chap. 31, McGraw-Hill, New York, 2009.
114. P. Hariharan, in *Handbook of Optics*, 3rd ed., vol. I, chap. 32, McGraw-Hill, New York, 2009.
115. P. R. Norton in *Handbook of Optics*, 3rd ed., vol. II, chap. 24, McGraw-Hill, New York, 2009.
116. S. Perkowitz, D. G. Seiler, and W. M. Duncan, "Optical Characterization in Semiconductor Manufacturing," *J. Res. Natl. Inst. Stand. Technol.* **99**:605 (1994).
117. G. E. Jellison, Jr, "Physics of Optical Metrology of Silicon-Based Semiconductor Devices," in *Handbook of Silicon Semiconductor Metrology*, A. C. Diebold (ed.) Marcel Dekker, Inc., New York, 2001, p. 723.

118. Clive Hayzelden, "Gate Dielectric Metrology," in *Handbook of Silicon Semiconductor Metrology*, A. C. Diebold (ed.), Marcel Dekker, Inc., New York, 2001, p. 17.
119. R. W. Collins, "Automatic Rotating Element Ellipsometers: Calibration, Operation, and Real-Time Applications," *Rev. Sci. Instrum.* **61**:2029 (1990).
120. U. Gerhardt and G. Rubloff, "A Normal Incidence Scanning Reflectometer of High Precision," *Appl. Opt.* **8**:305 (1969).
121. M. I. Bell and D. A. McKeown, "High-Precision Optical Reflectometer for the Study of Semiconductor Materials and Structures," *Rev. Sci. Instrum.* **61**:2542 (1990).
122. D. E. Aspnes and J. E. Fischer, in *Encyclopaedic Dictionary of Physics*, suppl. vol. 5, Thewlis (ed.), Pergamon, Oxford, 1975, p. 176.
123. D. E. Aspnes, in *Handbook on Semiconductors*, vol. 1, T. S. Moss and M. Balkanski (eds.), North-Holland, Amsterdam, 1980, p. 109.
124. D. E. Aspnes and A. A. Studna, "Schottky-Barrier Electroreflectance: Application to GaAs," *Phys. Rev. B* **7**:4605–4625 (1973).
125. R. E. Nahory and J. L. Shay, "Reflectance Modulation by the Surface Field in GaAs," *Phys. Rev. Lett.* **21**:1569–1571 (1968).
126. K. Dev, M. Y. L. Jung, R. Gunawan, R. D. Braatz, and E. G. Seebauer, "Mechanism for Coupling between Properties of Interfaces and Bulk Semiconductors," *Phys. Rev. B* **68**:195311 (2003).
127. K. Dev and E. G. Seebauer, "Band Bending at the Si(1 1 1)–SiO₂ Interface Induced by Low-Energy Ion Bombardment," *Surf. Sci.* **550**:185–191 (2004).
128. F. H. Pollak, "Modulation Spectroscopy as a Technique for Semiconductor Characterization," *Proc. SPIE* **276**:142–156 (1981).
129. O. J. Glembocki, B. V. Shanabrook, N. Bottka, W. T. Beard, and J. Comas, "Photoreflectance Characterization of Interband Transitions in GaAs/AlGaAs Multiple Quantum Wells and Modulation-Doped Heterojunctions," *Appl. Phys. Lett.* **46**:970 (1985).
130. X. Yin, H. M. Chen, F. Pollak, Y. Chan, P. A. Montano, P. D. Kichner, G. D. Pettit, and J. M. Woodall, "Photoreflectance Study of Surface Photovoltage Effects at (100) GaAs Surfaces/Interfaces," *Appl. Phys. Lett.* **58**:260 (1991).
131. F. H. Pollak, M. Cardona, and D. E. Aspnes (eds.), *Proc. Int. Conf. on Modulation Spectroscopy*, Proc. SPIE Bellingham, 1990, p. 1286.
132. O. J. Glembocki and B. V. Shanabrook, in *Semiconductors and Semimetals*, vol. 36, D. G. Seiler and C. L. Littler (eds.), Academic Press, New York, 1992, p. 221.
133. H. G. Tompkins and W. A. McGahan, *Spectroscopic Ellipsometry and Reflectometry*, Wiley, New York, 1999.
134. *The Proc. 2007 Int. Conf. on Spectroscopic Ellipsometry IV*, in *Phys. Stat. Sol. (a)* **205**:4 (2008).
135. H. G. Tompkins and E. A. Irene, *Handbook of Ellipsometry*, Springer, New York, 2005.
136. R. M. A. Azzam, "Ellipsometry," in M. Bass (ed.), *Handbook of Optics*, 3rd ed., vol. V, chap. 2, McGraw-Hill, New York, 2009.
137. R. W. Collins, Ellipsometry, in *The Optics Encyclopedia*, Wiley-VCH, Weinheim, 2004, p. 609.
138. N. V. Edwards, in *Characterization and Metrology for ULSI Technology*, D. G. Seiler A. C. Diebold, Th. J. Shaffner, R. McDonald, S. Zollner, R. P. Khosla, and E. M. Secula (eds.), AIP, Melville, 2003, p. 723.
139. M. Schubert, *Infrared Ellipsometry on Semiconductor Layer Structures*, Springer-Verlag, Berlin, 2004.
140. M. Schubert, T. Hofmann, and C. M. Herzinger, "Generalized Far-Infrared Magneto-Optic Ellipsometry for Semiconductor Layer Structures: Determination of Free-Carrier Effective-Mass, Mobility, and Concentration Parameters in *n*-Type GaAs," *J. Opt. Soc. Am. A* (**20**):347–356 (2003).
141. T. Hofmann, M. Schubert, C. M. Herzinger, and I. Pietzonka, "Far-Infrared-Magneto-Optic Ellipsometry Characterization of Free-Charge-Carrier Properties in Highly Disordered *n*-Type Al_{0.19}Ga_{0.33}In_{0.48}P₅," *Appl. Phys. Lett.* **82**:3463 (2003).
142. T. Hofmann, U. Schade, K. Agarwal, B. Daniel, C. Klingshirn, M. Hetterich, C. Herzinger, and M. Schubert, "Conduction-Band Electron Effective Mass in Zn_{0.87}Mn_{0.13}Se Measured by Terahertz and Far-Infrared Magneto-optic Ellipsometry," *Appl. Phys. Lett.* **88**: 42105 (2006).
143. T. Hofmann, C. M. Herzinger, C. Kraemer, K. Streubel, and M. Schubert, "The Optical Hall Effect," *Phys. Stat. Sol. (a)* **205**(4): 779–783 (2008).

144. T. Hofmann, U. Schade, C. M. Herzinger, P. Esquinazi, and M. Schubert, "Terahertz Magneto-Optic Generalized Ellipsometry Using Synchrotron and Blackbody Radiation," *Rev. Sci. Instr.* **77**:063902 (2006).
145. D. B. Holt and B. G. Yacobi, in *SEM Microcharacterization of Semiconductors*, D. B. Holt and D. C. Joy (eds.), Academic Press, New York, 1989, p. 373.
146. J. I. Pankove, in *Electroluminescence*, J. I. Pankove (ed.), Springer-Verlag, Berlin, 1977, p. 1.
147. E. C. Lightowers, in *Growth and Characterization of Semiconductors*, R. A. Stradling and P. C. Klipstein (eds.), Adam Hilger, Bristol, 1990, p. 135.
148. M. L. W. Thewalt, M. K. Nissen, D. J. S. Beckett, and K. R. Lundgren, "High Performance Photoluminescence Spectroscopy Using Fourier Transform Interferometry," *Mat. Res. Soc. Symp.* **163**:221 (1990).
149. D. E. Aspnes, in *Properties of GaAs*, EMIS Data Reviews, Series #2, p. 229 INSPEC, London, 1990.
150. F. Fuchs, A. Lusson, P. Koidl, and R. Triboulet, "Fourier Transform Infrared Photoluminescence of $\text{Hg}_{1-x}\text{Cd}_x\text{Te}$," *J. Crystal Growth* **101**:722 (1990).
151. P. Hiesinger, S. Suga, F. Willmann, and W. Dreybrodt, "Excitation Spectra of Exciton Luminescence in CdTe," *Phys. Stat. Sol. (b)* **67**:641–652 (1975).
152. *Light Scattering in Solids*, M. Cardona (ed.), vol. 1., Springer-Verlag, Berlin, 1983; *Light Scattering in Solids*, vol. 2, M. Cardona and G. Guntherodt (eds.), Springer-Verlag, Berlin, 1982; *Light Scattering in Solids*, vol. 3, M. Cardona and G. Guntherodt (eds.), Springer-Verlag, Berlin, 1982; *Light Scattering in Solids*, vol. 4, M. Cardona and G. Guntherodt (eds.), Springer-Verlag, Berlin, 1984; *Light Scattering in Solids*, vol. 5, M. Cardona and G. Guntherodt (eds.), Springer-Verlag, Berlin, 1989.
153. (op. cit.) W. Hayes and R. Loudon, pp. 1–360.
154. A. S. Pine, in *Light Scattering of Solids*, vol. 1, M. Cardona (ed.), Springer-Verlag, Berlin, 1982, p. 253.
155. (op. cit.) W. Hayes and R. Loudon, p. 16.
156. (op. cit.) W. Hayes and R. Loudon, p. 44.
157. G. Abstreiter, "Micro-Raman Spectroscopy for Characterization of Semiconductor Devices," *Applied Surface Science* **50**:73–78 (1991).
158. K. K. Tiong, P. M. Amirtharaj, F. H. Pollak, and D. E. Aspnes, "Effects of As^+ Ion Implantation on the Raman Spectra of GaAs: 'Spatial Correlation' Interpretation," *Appl. Phys. Lett.* **44**:122 (1984).
159. P. A. Temple and C. E. Hathaway, "Multiphonon Raman Spectrum of Silicon," *Phys. Rev. B* **7**:3685–3697 (1973).
160. G. Abstreiter, R. Trommer, M. Cardona, and A. Pinczuk, "Coupled Plasmon-LO Phonon Modes and Lindhard-Mermin Dielectric Function of n -GaAs," *Solid State Communications* **30**:703–707 (1979).
161. A. Pinczuk, L. Brillson, E. Burstein, and E. Anastassakis, "Resonant Light Scattering by Single-Particle Electronic Excitations in n -GaAs," *Phys. Rev. Lett.* **27**:317–320 (1971).
162. J. M. Ralston, R. L. Wadsack, and R. K. Chang, "Resonant Cancellation of Raman Scattering from CdS and Si," *Phys. Rev. Lett.* **25**:814–818 (1970).
163. M. Chandrasekhar, H. R. Chandrasekhar, M. Grimsditch, and M. Cardona, "Study of the Localized Vibrations of Boron in Heavily Doped Si," *Phys. Rev. B* **22**:4825–4833 (1980).
164. D. D. Manchon, Jr., and P. J. Dean, *Proc. 10th Int. Conf. on Physics of Semiconductors*, S. P. Keller, J. C. Hensel, and F. Stern (eds.), USAEC, Cambridge, 1970, p. 760.
165. A. W. Nelson, in *Electronic Materials from Silicon to Organics*, L. S. Miller and J. B. Mullin (eds.), Plenum Press, New York, 1991, p. 67.
166. O. Madelung, M. Schulz, and H. Weiss (eds.), *Landolt-Bornstein Numerical Data and Functional Relationships in Science and Technology, Group III—Crystal and Solid State Physics*, vol. 17a, Springer-Verlag, Berlin, 1982.
167. O. Madelung, M. Schulz, and H. Weiss (eds.), *Landolt-Bornstein Numerical Data and Functional Relationships in Science and Technology, Group III—Crystal and Solid State Physics*, vol. 17b, Springer-Verlag, Berlin, 1982.
168. O. Madelung, M. Schulz, and H. Weiss (eds.), *Landolt-Bornstein Numerical Data and Functional Relationships in Science and Technology, Group III—Crystal and Solid State Physics*, vol. 17f, Springer-Verlag, Berlin, 1982.
169. O. Madelung and M. Schulz (eds.), *Landolt-Bornstein Numerical Data and Functional Relationships in Science and Technology, Group III—Crystal and Solid State Physics*, vol. 22a, Springer-Verlag, Berlin, 1987.
170. <http://lb.chemie.uni-hamburg.de/search/index.php>, accessed May 21, 2009.
171. O. Madelung, *Semiconductors: Data Handbook*, 3rd ed., Springer-Verlag, Berlin, 2004.

CHARACTERIZATION AND USE OF BLACK SURFACES FOR OPTICAL SYSTEMS

Stephen M. Pompea

*National Optical Astronomy Observatory
Tucson, Arizona*

Robert P. Breault

*Breault Research Organization, Inc.
Tucson, Arizona*

6.1 INTRODUCTION

Optical instruments and telescopes rely on black baffle and vane surfaces to minimize the effect of stray light on overall system performance. For well-designed and well-baffled systems, the black surfaces chosen for the baffles and vanes can play a significant role in reducing the stray light on the detector.¹⁻⁸ In space-borne systems, a large number of black surfaces play an important role. An excellent and comprehensive review of infrared surfaces is given by Persky.⁹ Additional infrared measurements of materials are given by Miller¹⁰ and Persky and Szczesniak.¹¹ Black surfaces are also used extensively in solar collector applications. Excellent reviews of spectrally selective surfaces for heating and cooling applications are found in Hahn and Seraphim¹² and in Granqvist¹³. In many solar applications, high solar absorptance is desired along with low thermal emittance.¹⁴ In general, surfaces for solar applications will not be addressed in this chapter.

Black coatings are also used in radiometric detectors.¹⁵ Because the surface needed is often small, these surfaces may be even more specialized than the black coatings used for stray light reduction in optical instruments. This chapter will concentrate on the selection and characterization of black surfaces chosen for stray light suppression and suitable for application to relatively large areas of an optical system or optical test bed. Some examples of these uses are seen in Table 1.

The optical system designer has a wide repertoire of baffle surfaces from which to choose. Summaries of optical properties of materials were given by Wolfe,¹⁶ Pompea et al.,⁴ and McCall et al.¹⁷ Reviews of materials by McCall¹⁸ and Smith and Howitt¹⁹ emphasized the ultraviolet/visible and infrared properties, respectively. A number of company databases of scattering data are available, including one using the same instrument (for bidirectional scatter distribution function or BSDF measurements at 0.5145 μm) for approximately 15,000 data runs!²⁰ Large amounts of BSDF data are at Breault Research Organization and a BSDF database format has been proposed by Klicker et al.²¹ An organized effort to create specialized databases of optical properties of surfaces applicable to both ground- and space-based instruments has been undertaken and is proving to have great utility.²²

The choices of optical black surfaces are usually first narrowed by the nature of the application, the substrates available or possible, the wavelength or bandpasses of interest, the angles at which the surfaces must be nonreflective, and a host of system issues and environmental factors.⁵ As the system performance requirements have become more stringent, an array of surfaces has become available to meet these requirements. Many paints (e.g., Chemglaze Z306, SolarChem) (Table 2); anodized surfaces (e.g., Martin Black, Infrablack, Tiodize) (Table 3); etched, electrodeposited, and

6.2 PROPERTIES

TABLE 1 Possible Uses for Black Surfaces in Optical Systems and Test Beds

Apertures and field stops	Baffles	Barrels	Blackbodies	Choppers	Cold shields	Detector housings
Dewar interiors	Enclosures and testing structures	Lens edges	Laser light traps	Radiometers	Radiators	Simulators and targets

TABLE 2 Painted Surfaces

Surface Name or Designation	Manufacturer and/or Distributor (contact person)	Historical Notes	Surface Type	Main Literature References
Aeroglaze L300	Lord Corporation Chemical Products Division Industrial Coatings 2000 West Grandview Boulevard P.O. Box 10038 Erie, PA 16514-0038	Formerly called Chemglaze L300	Paint	24
Chemglaze Z004	Lord Corporation Erie, Pa.		Paint	24
Aeroglaze Z302	Lord Corporation Erie, Pa.	Formerly called Chemglaze Z302	Paint	25, 26, 24
Aeroglaze Z306	Lord Corporation Erie, Pa.	Formerly called Chemglaze Z306	Paint	27, 24, 28, 29, 30, 31, 20, 32, 33, 19, 34, 35, 36, 37, 38, 39, 40, 41, 42
Aeroglaze Z306 with microspheres	Lord Corporation Erie, Pa.	Formerly called Chemglaze Z306 with microspheres	Paint	43, 44, 45, 28, 24
Aeroglaze Z307	Lord Corporation Erie, Pa.	Formerly called Chemglaze Z307 (conductive)	Paint	24
Aeroglaze Z313	Lord Corporation Erie, Pa.	Formerly called Chemglaze Z313	Paint	35, 24
Ames 24E Ames 24E2	NASA Ames Research Center Moffet Field, Calif. S. Smith Sterling Software 1121 San Antonio Rd. Palo Alto, CA 94303		Paint	4, 46, 47, 34
Cardinal 6450	Cardinal Industrial Finishes 1329 Potrero Ave., So. El Monte, CA 91733-3088	Formerly called "Cardinal 6550"	Paint	48, 35
Cornell Black	Prof. J. Houck Department of Astronomy Cornell University Ithaca, NY 14853		Paint	39, 32, 33, 19, 49
DeSoto Flat Black	PRC-Desoto International, Inc. 5454 San Fernando Rd. Glendale, CA 91203		Paint	50, 35
Electrically Conductive Black Optical Paint	Jet Propulsion Laboratory, Caltech, 4800 Oak Grove Dr. Pasadena, CA 91109	Has no tradename	Paint	51, 52, 53

(Continued)

TABLE 2 Painted Surfaces (*Continued*)

Surface Name or Designation	Manufacturer and/or Distributor (contact person)	Historical Notes	Surface Type	Main Literature References
IITRI Bone Black D-111 (IITRI D111)	IIT Research Institute 10 West 35th Street Chicago, IL 60616		Paint	25, 54, 32, 34, 55, 38
LMSC Black	Lockheed Palo Alto Research Lab		Painted multi-player coating	32, 33, 19, 55
MH21-1	IIT Research Institute Chicago, Ill.		Paint	54
MH55	IIT Research Institute Chicago, Ill.		Paint	54
MH2200	IIT Research Institute Chicago, Ill.	Formerly 3M's ECP 2200 paint, but sold to IIT	Paint	26, 27, 56, 54, 57, 33, 19, 35, 58
Solarchem	Eastern Chem Lac Corporation 1080-T Eastern Ave. Malden, MA 02148		Paint	4
463-3-8	Akzo Nobel Coatings, Inc. 434 W. Meats Avenue Orange, CA 92665	Formerly called "Cat-a-lac 463-3-8" diffuse black paint	Paint	27, 56, 59, 60, 30, 58, 42, 55, 41
443-3-8	Akzo Nobel Coatings, Inc. Orange, Calif.	Formerly called "Cat-a-lac 443-3-8"	Paint	61, 25, 56, 59, 60, 58, 55, 41
443-3-17	Akzo Nobel Coatings, Inc. Orange, Calif.	Formerly called "Sikkens 443-3-17" glossy black	Paint	26, 60

Source: Adapted from McCall.²³

TABLE 3 Anodized Surfaces

Surface Name	Manufacturer	Notes	Surface Type	Main Literature References
Infrablack	Lockheed Martin, Denver CO 80201	For Al substrates only	Anodization process	43, 62, 63, 3, 64, 4, 46
Martin Black	Lockheed Martin Denver, CO	For Al substrates only	Anodization process	43, 25, 65, 62, 66, 63, 67, 3, 45, 68, 59, 69, 70, 32, 33, 46, 34, 35, 58, 55, 29, 38, 31, 64, 4, 20, 71, 39, 72, 37, 41, 38
Martin Black, enhanced	Lockheed Martin Denver, Colo.	For Al substrates only	Anodization process	63, 31, 4
Martin Black, posttreated	Lockheed Martin Denver, Colo.	For Al substrates only	Anodization process	63

Source: Adapted from McCall.²³

plasma-sprayed metal surfaces (Table 4) are now available to meet quite specific optical and system performance and environmental requirements.

New classes of surfaces are also being developed to give selected performance at specific angles and wavelengths. Other materials are being developed for applications where "hardened" laser-resistant or radiation-resistant materials are needed. A third area where much development is currently taking place is in the area of materials that are able to withstand severe and unusual forms of environmental exposure for long periods.

TABLE 4 Other Processes

Surface Name	Manufacturer	Notes	Surface Type	Main Literature References
Black chrome-type surfaces	Lockheed Martin Denver, Colo.	For many kinds of metal substrates	Electrodeposition process	63, 4, 41
Black cobalt-type surfaces (Cobalt black)	Lockheed Martin Denver, Colo. And many more companies	For many kinds of metal substrates; references to black copper, black, steel, etc., are covered by black cobalt	Electrodeposition processes, and can be followed by chemical or thermal oxidation	63, 4, 73, 74
Black nickel (NBS black) (Ball black)	Ball Aerospace Electro-Optics Engineering Dept. P.O. Box 1062 Boulder CO 80306 and many others	For many kinds of metal substrate; Ball improved the patent developed by NBS	Deposition and etching process	65, 19, 29, 50, 37, 75, 76
Black Kapton film	DuPont Wilmington, DE 19898		Foil	35
Black Tedlar film	DuPont TEDLAR / Declar, PPD, D-12082 1007 Market Street Wilmington, DE 19898		Foil	35
Boron black	Lockheed Martin Denver, Colo.	For many kinds of metals	Plasma spray deposition process	4, 63
Boron carbide	Lockheed Martin Denver, Colo.	For Ti substrates only	Proprietary process	4
Silicon carbide	Lockheed Martin Denver, Colo. And more companies	For many kinds of metals	Chemical vapor deposition	4
Textured surfaces	NASA Ames Research Center Moffet Field, Calif. (Sheldon Smith) Spire Corporation Patriots Park, Bedford, MA 01730 Optics MODIL, Oak Ridge National Lab Oak Ridge, TN 37831-8039	For many kinds of metal substrates	Either: —sputtered coated —ion beam etched —sputter coated then etched	77, 78, 79, 29, 80, 37, 81, 82, 79, 83
Black optical thin-film interference coatings	National Research Council of Canada Thin Films, Institute for Microstructural Sciences Montreal Road, Building M-36 Ottawa, Ontario, Canada, K1A 0R6	For metal, dielectric, or other substrates	Vacuum deposition techniques: —sputter deposition —ion vapor deposition —resistance-heated source —electron-beam gun deposition	84, 85, 86, 87

Source: Adapted from McCall.²³

The surface morphology of a “black” or diffusely scattering surface is only one determinant of its complex optical properties. However, for many surfaces, surface roughness plays a most significant role. This can be illustrated in Figs. 1 to 8, scanning electron micrographs of some important black or diffuse surfaces. The size and shapes of the surface features provide a valuable indication of how light will be absorbed and scattered or diffracted by such a surface. The creation and design of new

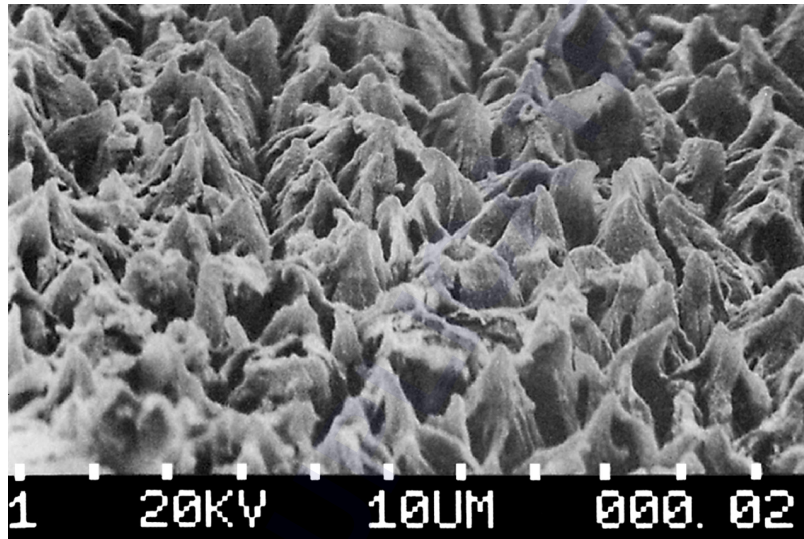


FIGURE 1 Scanning electron micrograph of Martin Optical Black, an anodized aluminum surface for ultraviolet, visible, and infrared use. (Photo courtesy of Don Shepard, Lockheed Martin, Denver.)

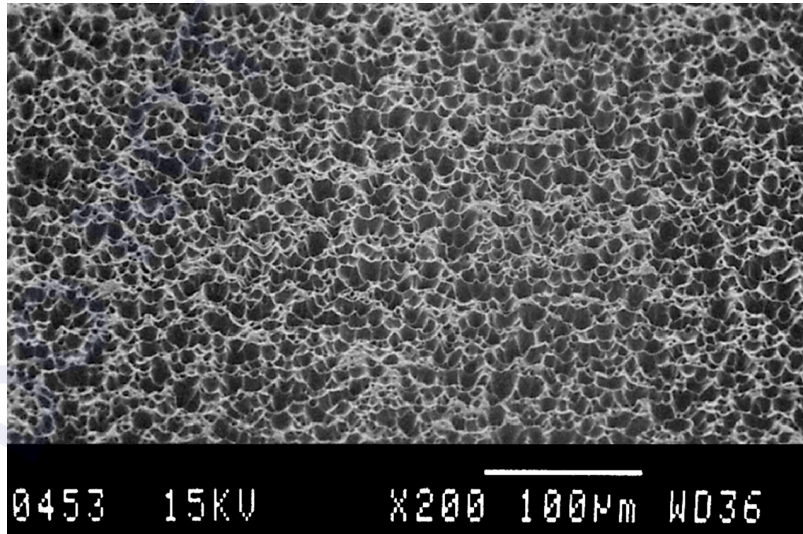


FIGURE 2 Scanning electron micrograph of Ball Black, an etched electroless nickel surface applicable to a variety of substrates. (Photo courtesy of Arthur Olson, Ball Aerospace Systems Group, Boulder, Colorado.)

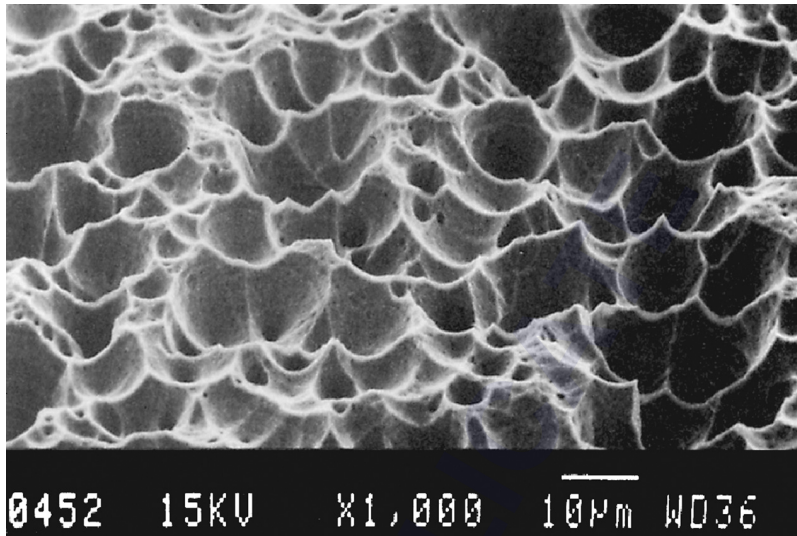


FIGURE 3 Scanning electron micrographs of Ball Black, an etched electroless nickel surface applicable to a variety of substrates. This surface is representative of a class of etched electroless nickel surfaces. (Photo courtesy of Arthur Olson, Ball Aerospace Systems Group, Boulder, Colorado.)

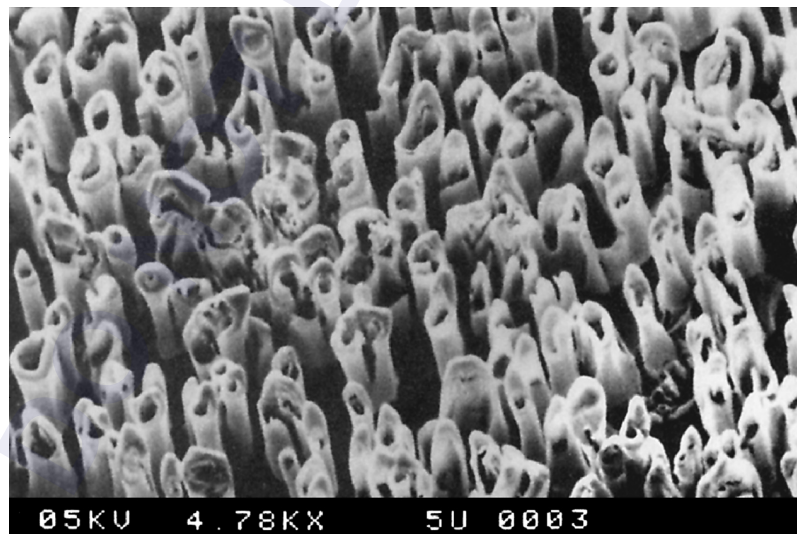


FIGURE 4 Scanning electron micrograph of a sputtered beryllium surface. (Photo courtesy of Roland Seals, Optics MODIL, Oak Ridge National Laboratory.)

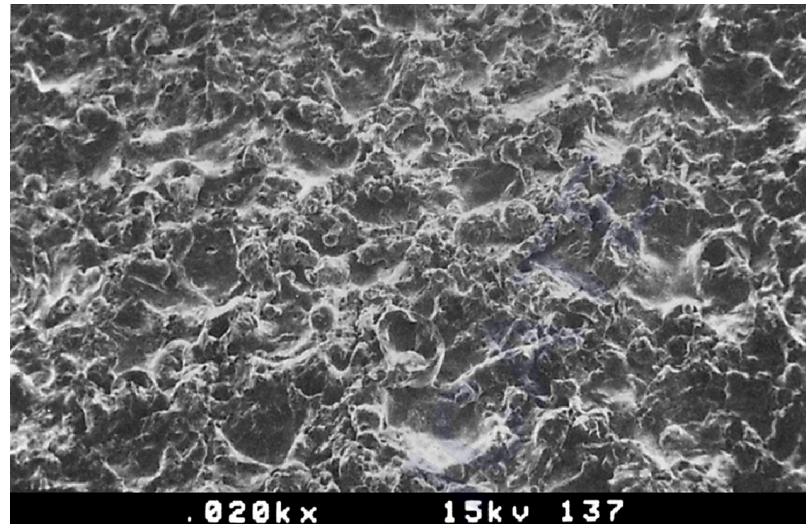


FIGURE 5 Scanning electron micrographs of Ames Perfect Diffuse Reflector (PDR) at 24 magnification. (Photo courtesy of Sheldon Smith, NASA Ames Research Center and Sterling Federal Systems, Palo Alto, Calif.)

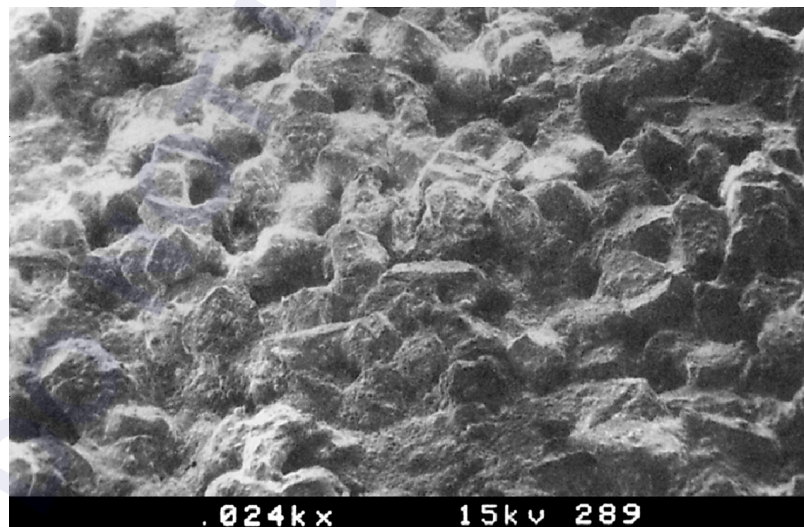


FIGURE 6 Scanning electron micrographs of Ames 24E at 24 magnification. (Photo courtesy of Sheldon Smith, NASA Ames Research Center and Sterling Federal Systems, Palo Alto, Calif.)

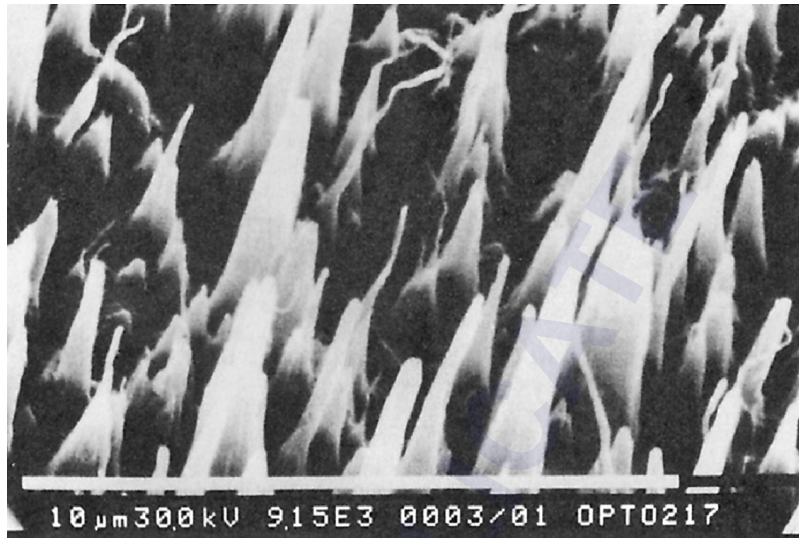


FIGURE 7 Scanning electron micrograph of a textured graphite surface created by bombarding a carbon surface with positive argon ions. (Photo courtesy of Chuck Bowers, Hughes Aircraft, El Segundo, Calif.)

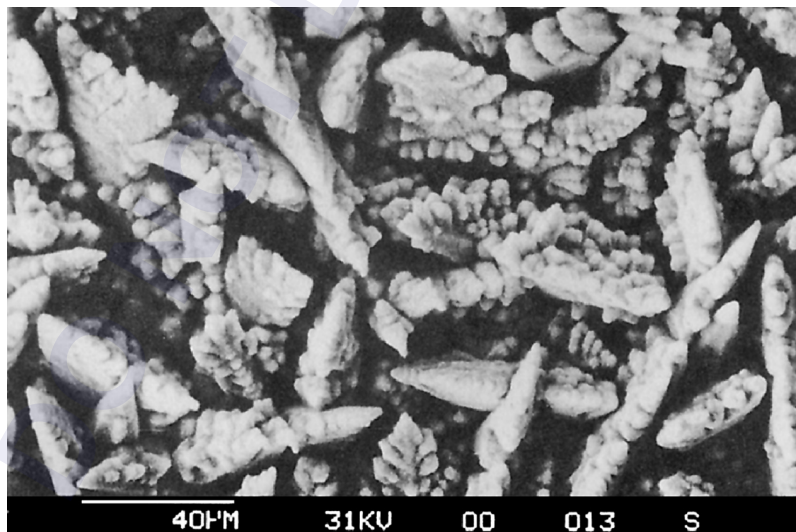


FIGURE 8 Scanning electron micrograph of Orlando Black surface, produced by electrodeposition of copper and subsequent oxidation in a proprietary process. (Photo courtesy of D. Janeczko, Martin Marietta Electronic Systems, Orlando, Fla.)

surfaces will be touched upon later, in the section on design techniques for creating new surfaces for specific applications.

This chapter gives a summary of the materials used in a variety of optical systems (with an emphasis on those that are currently available) and describes their optical and material properties so that the optical designer can begin the material selection process. As such, it is rather an extreme condensation of the data available. However, even the data presented here cannot be considered very definitive. There are a number of reasons for inconsistencies and ambiguities in the data presented. First, many of these surfaces and the processes that create them are evolving, and improving continuously. Even though the actual surface may change, the name may not. Thus, optical measurements of the same named surface that are separated by several years may not be consistent, even if the measurement techniques are consistent.

A second cause of inconsistency among data sets comes from the remeasurement of “archival” samples. When new measurement techniques or improved instruments become available, or the needs of a program demand new measurements, archival samples are retrieved and remeasured. Sometimes these archival samples may not have been stored properly and may not be in pristine or original condition. Other times, these samples may never have been archival in quality. They may have been marketing samples made without specific quality control and distributed widely. These measurements still enter the body of literature with the reader usually unaware of the important circumstances. For robust samples, poor storage may be of little importance. For more exotic materials (e.g., specialized baffles for space applications) that must be handled carefully, lack of proper handling can be of great importance. In this latter case, the measurements made on these degraded surfaces may not be representative.

Wolfe⁸⁸ compares the theory and experiments for bidirectional reflectance distribution function (BRDF) measurements of microrough surfaces. The BRDF of a surface is very useful in understanding the optical performance of a surface and is defined as the ratio of scattered radiance [$W/(cm^2sr)$] to surface irradiance (W/cm^2). Its units are inverse steradians. Radiance is used in the definition of BRDF in order to make the BRDF independent of the parameters of the measuring instrument, such as the detector aperture and distance to the detector. BRDF measurements can be made by a variety of instruments and can be made for any number of wavelengths. For example, BRDF measurements of surfaces have been made in vacuum ultraviolet to far-infrared wavelengths. However, the most common BRDF measurements that are made on black surfaces are probably made at convenient laser wavelengths of 0.6328 and 10.6 μm . It is important to keep in mind that these measurements may be of less predictive value if the system is operating at a substantially different wavelengths than where the measurements were made.

To make a BRDF measurement requires a light source, a sample mounted and illuminated by that source from a variety of angles, a detector to measure the scatter from the sample, and the computer/electronics package to accurately record the detection of the scattered light as a function of angle of incidence and detection angle.

While the optical properties of surfaces can be described or characterized through the use of specular and diffuse reflectance measurements, the use of BRDF measurements are very useful in characterizing both highly specular and highly diffuse surfaces and provide an excellent way to characterize the optical properties of any surface. The BRDF can be used to describe the angle-dependent optical scatter from any surfaces. The angular distribution of scatter can be used, in conjunction with computer modeling to calculate if the scatter from a black surface will be a limiting factor in the resolution or noise level of an optical system. BRDF measurements of surfaces are also used in computer graphics visualization programs to illustrate what a surface may look like under various illumination conditions. For some types of black surfaces, the angular distribution of scatter can be used to calculate a variety of surface parameters (e.g., surface roughness) if certain assumptions about the surface can be made. For opaque black surfaces, appropriate BRDF measurements in the wavelengths, angles, polarizations, and so on of interest have proven to be one of the best overall descriptors of the optical properties of the surface.

Measurements of optical quantities such as the BRDF, that are by definition measurement-device-independent, can also show large variations.^{89,90} No attempt has been made in this chapter to reconcile measurement discrepancies; indeed, no attempt to even identify the areas where conflicting data

occur is made. This would be a herculean task, but more importantly, it is probably unnecessary. The variations in the measurement process ensure that the optical description of black surfaces is still an order-of-magnitude science, or, at best, a half-an-order-of-magnitude science. Optical practitioners should use their own safety factor or better yet, have their own measurements made.^{5,17}

6.2 SELECTION PROCESS FOR BLACK BAFFLE SURFACES IN OPTICAL SYSTEMS

The selection of black surfaces is a systems issue; it must be addressed early and consider all aspects of system design and performance. It must also be examined from a total system performance perspective. Cost and schedule considerations have important effects on the decision and need to be taken seriously. A baffle structure that cannot be built with surfaces that cannot be coated is not a desirable state of affairs!

Similarly, proper financial and schedule support must be present. Black surfaces need to be taken as seriously as any other optical components, such as mirror surfaces or the thin films that act as filters. They deserve serious treatment in their design, fabrication, procurement, and testing. They are not last-minute design decisions or items to be created at the last minute. Since the baffle design and choice of surfaces often sets the final system performance (especially in infrared systems), it makes little sense to expend great energy and money to optimize other more traditional system components while ignoring stray light design and the surfaces so important in the design.

Figures 9 and 10 are a comparison of two processes that could lead to the selection of surfaces for an application. Figure 9 shows a process where design activities are done sequentially. In this process, stray light issues (including black surface issues) are left to the end of the program. In practice, with limited budgets, stray light and other fundamental system engineering issues are often not addressed at all in this process, or are addressed only at stages in the program where the system design is frozen, or cannot be changed without great expense. At this stage, when fundamental problems are found, they often have severe performance, budget, and schedule penalties. That many stray light and black surface problems are rectified (with severe penalties to schedule and budget) at late stages in programs attests to the overall importance of black surfaces to the ultimate success of a system.

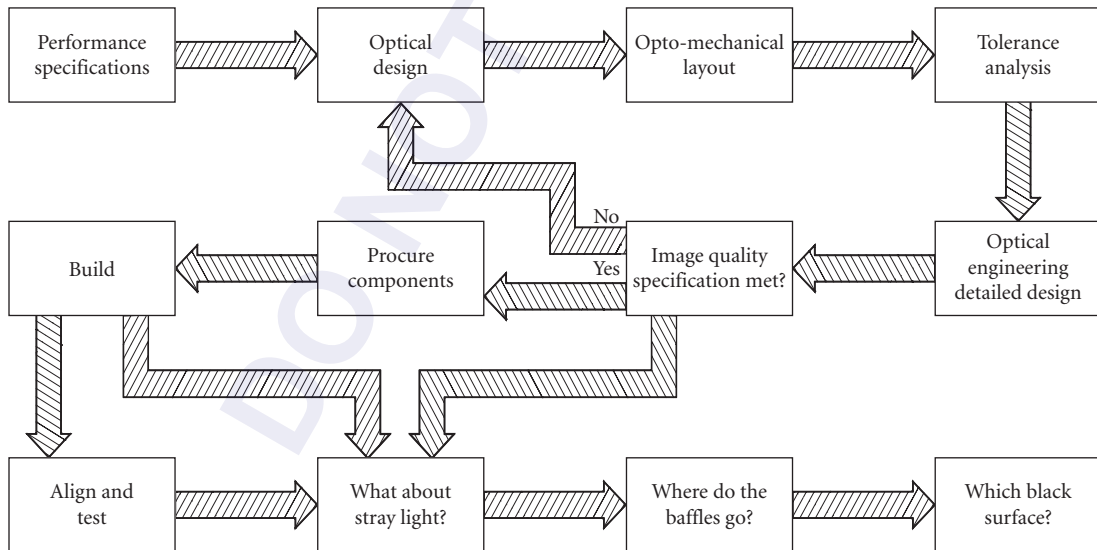


FIGURE 9 When selection of black surfaces is left to the end of a program, serious risks to the program are likely. (From Pompea and McCall.⁵)

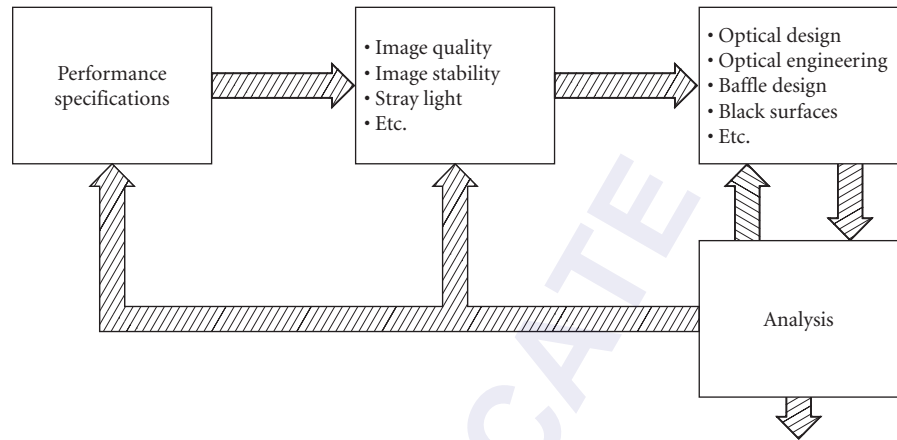


FIGURE 10 The most effective context for black surface selection is near the beginning of the program. This flowchart illustrates the early phase of a system-oriented design process. (From Pompea and McCall.⁵)

In an integrated systems engineering or concurrent engineering environment, illustrated in Fig. 10, the selection process begins early. There are many iterations of the design in the early stages, and close work among designers of different disciplines is essential. The level of detail increases with time, but no fundamental issues are decided without addressing their implication for the system and for each subsystem. This process is far superior to the one in the first flowchart, in which the surface selection process is treated as a trivial one and one waits until the end of the program to evaluate candidate surfaces. A review of general aspects of selection with emphasis on the system level issues inherent in the selection process is given by Pompea and McCall,⁵ McCall et al.,¹⁷ and Pompea.⁹¹

Once it has been decided to use black surfaces in an optical system, a whole host of system-level issues must be addressed to determine an appropriate coating or coatings. The high-level items include (1) purpose and position of the surface in the instrument, (2) wavelength(s) and waveband(s) of interest, (3) general robustness of the surface, (4) nature of the installation process, (5) environment of the optical system, (6) availability, (7) cost, (8) substrate, and (9) other mission or system requirements.

For systems that operate over a wide wavelength range, the selection process may be difficult because data may not be available over the entire wavelength range of interest. For example, the next generation of infrared and submillimeter telescopes are being designed for wavelengths between about 10 and 800 μm . Baffle coatings must be adequate over this entire wavelength range.

The position of the coating in the system and its function at that location is a critical item. For positions where the incident radiation is nearly normal to the surface, the optical considerations become less critical, since the reflectance at near-normal incidence does not vary as significantly at this angle for the most often used surfaces. Also, the BRDF is usually at its lowest value. There is a wealth of data at these near-normal incident angles.⁴

The corollary is that if surfaces must be used close to grazing incidence, the optical properties become extremely important to understand (and quantify for system predictions) just as the amount of data becomes vanishingly small.⁶ Although there is sufficient optical data on specular reflections from specular surfaces, there is relatively little information on specular reflections from diffuse surfaces, let alone BRDFs at large angles of incidence. For these reasons, optical designers try to use lambertian-like surfaces at near to normal incidence light paths as possible.

The robustness of the surface in its environment is a complex subject. Some specific areas that will be addressed here are the atomic oxygen environment of low earth orbit, outgassing in the vacuum of space, and particle generation by surfaces. The practical considerations of expense, schedule, ease of

cleaning, exportability, and availability of specialized surfaces all play important roles. For a unique space instrument, there are different manufacturing considerations than for a large-volume production process. The installation process is often overlooked. It often plays a key role in the successful use of many surfaces.

Exposed spacecraft surfaces often must not contribute to spacecraft charging problems. Some nonconductive surfaces support the buildup of charges that can create large electrical potentials between spacecraft components, or the spacecraft and the plasma and subsequent damage may result. This is particularly important in the region near a detector, which is often very charge-sensitive. Surfaces such as Martin Black do not support charges in excess of 200 V, since the coating can leak charge to the aluminum substrate. Textured metallic coatings, as a general rule, do not have surface charging problems. Some paints as well (e.g., Aeroglaze Z307) are designed to be conductive. Table 5 describes several surfaces used for thermal control in space craft. Some of these surfaces are designed to be conductive, to avoid charging effects. For space applications where charge buildup is not critical, nonconductive surfaces may be used.

Other considerations are the ability of a surface to withstand vibration, acceleration, elevated temperatures, thermal cycling, chemical attack (including atomic oxygen), solar and nuclear radiation, micrometeoroids, and moisture.

For the coating of baffle and (particularly) vane surfaces in ground-based telescopes, a further consideration may often be important. The baffle and vane surfaces in telescopes often have a large view angle or geometrical configuration factor (GCF) with the sky. This is particularly true for telescopes that operate without an enclosure or with a very large slit or opening to the night sky. In this case, the black baffle and vane surfaces must not radiatively couple strongly to the night sky. If there is a substantial view to the sky, and the emissivity of the surface is high, it will cool below the air temperature. This can create temperature-induced “seeing” effects which degrade image quality. The requirement that a surface be strongly absorbing of visible light, yet a poor emitter in the infrared is often a difficult one to meet. Table 6 gives the absorptivity and emissivity of some widely used surfaces.

In many laboratory or instrument applications, the use of one of a few standard black coatings or surfaces (such as the now unavailable 3M Black Velvet 101-C10) was entirely adequate to achieve the necessary performance of the optical system. In these applications, if the stray light performance was inadequate, it usually indicated that the stray light design, not the choice of black surface, was in error. The explanation was that the ignorance of some of the principles of stray light design could produce system performance five orders of magnitude worse than a poor choice of surfaces. The visibility of surfaces from the detector and the paths from those surfaces to the front of the instrument (not the optical properties or blackness of the surface) largely determine the total system performance. In stray light terminology, the GCF (geometrical configuration factor of surfaces) is more important than the BRDF (bidirectional scatter characteristics of surfaces).

TABLE 5 Black Nonconductive and Electrically Conductive Paints and Coatings for Space Use

Name	Type of Coating	Binder	Solar Absorptivity	Thermal Emittance
ML-210-IB	Thermal control	Inorganic	0.98±0.02	0.91±0.02
RM-550-IB	Thermal control	Inorganic	0.97±0.02	0.91±0.02
TMD-560-IB	Thermal control	Inorganic	0.95±0.02	0.91±0.02
AZ-1000-ECB	Thermal control conductive	Inorganic	0.97±0.02	0.89±0.02
MLS-85-SB	Thermal control	Organic	0.98±0.02	0.91±0.02
MLS-85-SB-c	Thermal control conductive	Organic	0.98±0.02	0.91±0.02

Source: Courtesy of J. Zweiner, AZ Technology, Huntsville, Alabama.

TABLE 6 Black Coatings

	$\bar{\alpha}_s$	$\bar{\epsilon}_n$
Anodize Black	0.88	0.88
Carbon Black Paint NS-7	0.96	0.88
Cat-a-lac Black Paint	0.96	0.88
Chemglaze Black Paint Z306	0.96	0.91
Delrin Black Plastic	0.96	0.87
Ebanol C Black	0.97	0.73
Ebanol C Black-384 ESH* UV	0.97	0.75
GSFC Black Silicate MS-94	0.96	0.89
GSFC Black Paint 313-1	0.96	0.86
Hughson Black Paint H322	0.96	0.86
Hughson Black Paint L-300	0.95	0.84
Martin Black Paint N-150-1	0.94	0.94
Martin Black Velvet Paint	0.91	0.94
3M Black Velvet Paint	0.97	0.91
Paladin Black Lacquer	0.95	0.75
Parsons Black Paint	0.98	0.91
Polyethylene Black Plastic	0.93	0.92
Pyramil Black on Beryllium Copper	0.92	0.72
Tedlar Black Plastic	0.94	0.90
Velestat Black Plastic	0.96	0.85

*Note: Solar absorptivity (0.3 to 2.4 μm) and normal emittance (5 to 35 μm) of black coatings.

Source: From Henniger.⁹²

6.3 THE CREATION OF BLACK SURFACES FOR SPECIFIC APPLICATIONS

The optical properties in the visible (e.g., BRDF at 0.6328 μm) at near-normal incidence between painted black surfaces and more specialized high-performance black surfaces usually do not differ by more than a factor of 10. For infrared and ultraviolet wavelengths, the performance of more specialized surfaces is far superior to many paints by a much larger factor. The optical designer must take care not to overspecify the optical performance properties required for a baffle surface. The lowest reflectance surface is not necessary for all applications. On the other hand, the choice of low-performance coatings has important consequences.

Some systems have been designed with an excessive number of vanes on their baffles. In these cases, the baffle system is probably too heavy, costs more than necessary to manufacture, and the excessive number of vane edges may be introducing extra diffraction effects that degrade the performance. The advantage of the high-performance coatings often lies not only in their optical properties, but in the considerable experience of their manufacturers in working on baffle systems, the reliability and consistency in the application process, and the superior documentation of the optical performance characteristics. The environmental resistance of the high-performance coatings is generally much better and much better documented than paints that are used for a wide variety of applications.

Black Surfaces as Optical Design Elements

The systems approach to the choice of black surfaces in optical systems that has been stressed here has led to the rule of thumb in stray light design that the choice of a black surface's optical properties should not be the first design issue considered. If the performance of the preliminary design is inadequate, the designer should look at the design carefully, rather than look for a better black surface to boost the performance. A flaw in the general system design is probably more responsible for the performance than the "blackness" of surfaces.

One reason for this approach is that the choice of a different black surface may not improve the performance, if the surface is used in the system in certain ways (e.g., grazing incidence). An important equalizing factor among many black surfaces is that they become highly specular when used at angles significantly different from normal incidence. This optical performance characteristic makes it critical for the baffle designer to conceptualize the paths and ranges of angles that are required, and to avoid using a surface where its performance is required at grazing incidence or close to it. The writing of scatter specifications is discussed by Stover,⁹³ and Breault² discusses sources of BRDF data and facilities where BRDF measurements are routinely made.

Black surfaces can be contaminated and their absorbing or diffuse structure can be destroyed by improper handling. This can lead to a significant degradation in their optical properties. A margin of safety should always be included in the performance estimates.

Many specific high-performance applications today require that black surfaces must become critical design elements that are used with only small "safety" factors. First, all of the principles of stray light design are used. If the required performance is barely achievable without optimizing the choice of black surfaces, then considerable attention must be paid to the choice of surfaces. Sometimes, designers must often use several black surfaces in an instrument or telescope; each particular black surface is used where its performance is optimal. In the past few years, a variety of black surfaces have become available to supplement the old standbys such as 3M Black Velvet (the mainstay of laboratories 10 years ago), and Chemglaze Z306 (now called Aeroglaze Z306) and Martin Black (the standard bearers for space instruments in the 1980s). Since the applications for these blacks are so specific, we refer to them as *designer blacks*. This term refers not only to new classes of surfaces that have been developed, but also to the creativity of optical and material scientists in modifying existing surfaces to create special properties for specific applications.

This new array of designer optical blacks can still be subdivided into the traditional categories of specular blacks and diffuse (lambertian) blacks. Glossy or specular black surfaces have been used for vane cavities and their use is reviewed by Freniere⁹⁴ and Freniere and Skelton.⁶¹ The approach is to reflect unwanted radiation out of the system or to attenuate it while tightly controlling its reflected direction. The use of specular baffles (reflective black or shiny) requires considerable care to avoid sending specular beams further into the instrument at some particular narrow angles.^{59,95} The increasing sophistication of stray light analysis programs makes this analysis much more complete. The most widely used specular black surfaces are Chemglaze Z302 and a paint formerly called Cat-a-lac Glossy Black, now available from Akzo as 463-3-8.

The lambertian black surface is often the most desirable type of surface to use in baffled telescopes and instruments. However, there is no *true* lambertian black, though several are extremely close for a near-normal angle of incidence. All diffuse blacks have larger BRDFs at large angles of incidence. Furthermore, these black surfaces can be further categorized by the wavelengths where they absorb best, the general categories being the ultraviolet, visible, near-infrared, middle-infrared, and far-infrared. They can also be distinguished by their degree of resistance to the atomic oxygen found in low earth orbit, their resilience under laser illumination, and their performance at narrow laser wavelengths or specific angles of incidence. The method of manufacture and the performance of these specialized surfaces may be proprietary or classified. Even if the general qualities of these surfaces that perform in these areas are known, the full range of test results on them often are not readily available.

Design Techniques for Creating Black Surfaces

Four general tools are available to the black surface designer to tailor the optical properties. Many highly absorbing surfaces use several of these effects to good advantage. The first tool is to use absorbing compounds in the surface. Examples of this are the organic black dye, which is a good visible absorber, in the Martin Black anodized surface and the carbon black particles found in many paints. Similarly, for the infrared paints, the addition of silicon carbide particles provides good absorption near 12 μm . The addition of compounds may create quite different properties at visible and infrared wavelengths; a visual observation of blackness is no guarantee of its infrared properties. The converse is also true. Titanium dioxide is a good infrared absorber, but is very reflective in the visible region. "Black," therefore, means absorptive in the wavelength region being discussed.

The addition of large (relative to the wavelength of interest) cavities, craters, or fissures in the surface can aid in the absorption process by requiring the radiation to make multiple reflections within the material before it leaves the surface. The same effect is achieved by growing large angular projections from a flat substrate. Each reflection allows more of the radiation to be absorbed. For example, a surface with a reflectivity of 15 percent (a rather poor black) becomes a very good black if the light can scatter or reflect three times within the surface before it exits. After three internal reflections, the 15 percent reflectivity black has become a $(0.15)^3$ or 0.34 percent reflective surface. The NBS Black and Ball Black surfaces are good examples of this; they are shiny electroless nickel surfaces until etched. After etching, the surface is filled with a plethora of microscopic, cavities that absorb enough light to create a surface that is black to the eye.

Black appliquéés can also be used to create rough absorbing surfaces.⁹⁶ These surfaces can be "tuned" for different wavelengths or for maximum absorptance at certain angles of incidence. This can be done by angling fibers or surface structure to provide cavities that are oriented normal to the angle of incidence. These materials have low reflectance at a variety of incident angles.⁹⁶⁻⁹⁸

Aerogels can also be used to create very black surfaces. Aerogels are highly porous materials nanometer-scale structures that can be synthesized through the removal of the pore-filling liquid. Of particular interest are carbon aerogels consist of 10-nm-sized particles, which are amorphous, connected in an interconnected network with approximately 10- to 50-nm pores and a pore volume of over 80 percent. Since carbon has inherent absorption in the infrared, this innately rough surface can have extremely low reflectivities over a broad wavelength range (e.g., 2 to 14 μm).⁹⁹

The third phenomenon is the use of scattering from the surface structure, from particles in the surface coating, or from the substrate to diffuse the incoming beam over a hemispherical solid angle. Even without any absorption, this dilution factor can be very important in destroying specular paths, particularly in the far-infrared, where broad wavelength absorption is difficult. For an introduction to surface roughness and scattering phenomena, see Refs. 33, 46, 93, and 100. An in-depth, very readable, treatment of scattering by small particles is given by Bohren and Huffman.¹⁰¹

The fourth phenomenon is based on optical interference of light in thin films. A black optical thin-film multilayer structure can significantly enhance, through optical interference, the amount of absorption of light in a multilayer structure over that of intrinsic absorption alone.⁸⁷ The black layer system technology has been used successfully in a range of applications.^{84-87,102,103} Optical thin films provide tremendous design flexibility. The reflectance, absorptance, and transmittance can be tailored as a function of angle and wavelength. Black multilayers may have a higher cost per unit of surface area, are generally limited to small surfaces, and have a strong angular sensitivity. This last characteristic can be used to great advantage for specialized applications. The films (especially those deposited by ion-beam process) are durable and are suitable for space use.

There is no ideal black surface. For good absorption across a wide wavelength band (e.g., 0.5 to 20 μm), the surface might be excessively thicker than is needed for a narrowband application. Even a "perfect" absorbing surface has disadvantages associated with higher emissivity (e.g., the surface radiates to the detector). In practice, even specular surfaces can be acceptable as black baffle surfaces under the right circumstances, although they are much more difficult to use. By design, they must direct the energy out of the system or in a direction that is not harmful to the system performance.^{95,104,105} Their alignment is critical and there are caveats for their use.⁵⁹ The "blackest" surface is not necessarily the best for a specific application.

In some other applications, a combination of low absorptance and high emittance are desirable. High emittance-low absorptance coatings are possible using a combination of techniques or surfaces. For example, a potassium silicate binder with zinc oxide or zinc orthotitanate particles can be applied to an anodized aluminum surface. In this case, the absorptance is less than 0.16 and the emittance is about 0.92.¹⁰⁶

6.4 ENVIRONMENTAL DEGRADATION OF BLACK SURFACES

Contamination control for terrestrial and spaceborne sensors is a very rapidly developing and important area. A good introduction to the field with extensive references is given in Refs. 2, 107 to 109. A good review of contamination assessment techniques given by Heaney¹¹⁰ and Nahm et al.¹¹¹ discusses scattering from contaminated surfaces. The prevention, detection, and removal of contamination should be important considerations early in the design of the sensor.

Three areas of great interest for spaceborne sensors are the effects of atomic oxygen, outgassing effects, and particle generation by surfaces. A large number of other effects can also influence the performance of surfaces or of the systems where those surfaces are installed over short or long periods. These include adhesion of coatings, radiation effects on coatings, thermal cycling effects, vacuum ultraviolet effects, and electrostatic charging effects, to name a few. It would be helpful if each new coating developed could be tested by a suite of tests relevant to the area of concern, similar to what was done in cryomechanical tests of Ames 24E2 infrared black coating.³⁴ All too often, the tests are done on different formulations at different labs and with different degrees of care. This makes it extremely difficult for the optical designer to have the data necessary to make good decisions.

The surface contamination of a black surface as well as aging effects can affect its optical properties. However, the effect of contaminants from black surfaces on mirror scatter in the optical system is much more pronounced. Williams and Lockie¹¹² use the BRDF as the sensitive parameter to judge the effects of dust, hydrocarbon oil, acrylic, and peelable coating residue on low-scatter mirror surfaces. Young¹¹³ describes the effect of particle scatter on mirrors. Shepard et al.⁷¹ have suggested that scattering properties of black baffle surfaces may be temperature-dependent. This may be a result of chemical changes in the surface or in changes to the surface as absorbed molecules are released at increased temperatures. Further work should be done at the temperatures of interest.

A large number of studies have been done to describe the behavior of black baffle surfaces to intense exposure to high-powered lasers and radiation, where damage to the surface is widespread. Accurate results can often be obtained using other vehicles to deliver energy into the surface. For example, rapid pulsed, low-energy electron beam irradiation of optical surfaces is described by Murray and Johnson¹¹⁴ as a cost effective way to test for radiation hardness. Black baffle surfaces are vulnerable even at low fluences because of their great absorptivity. The damage mechanisms are not well understood. Many of the damage mechanisms are related to the thermal shocks produced by rapid heating of a surface of thin layers on different materials. The destruction of surfaces not only reduces the ability of the surface to absorb stray radiation, but often also creates a contaminating cloud of particles that can greatly affect the optical system.

Atomic Oxygen Effects

Spacecraft surfaces on the early space-shuttle flights showed significant weight loss and aging effects. For example, on STS-1 the forward bulkhead camera blanket was milky yellow after the flight and the white paints on the shuttle exhibited exposure-related degradation effects. On later flights, surfaces such as Kapton showed significant mass losses. These effects were attributed to the interaction of the surfaces with the atomic oxygen present in low earth orbit.¹¹⁵ The atomic oxygen attacks the binder materials in paints, thus removing mass, and causes glossy surfaces to become more lambertian.¹¹⁶ Carbon particles are released when the binder is removed. The weight loss is due both to chemical processes and to erosional (kinetic energy) processes.

A second problem is that many surfaces have a glow associated with them when exposed to the atomic oxygen. Both of these phenomena undermine the effectiveness of baffle surfaces on space instruments. The effect of exposure to atomic oxygen on anodized black surfaces (Martin Black and Enhanced Martin Black) has been measured and simulations of atomic oxygen exposure on Chemglaze (now called Aeroglaze) Z306 paint has been done with a plasma etch chamber.⁶⁴

These tests showed that a carbon-black-based paint could be significantly degraded by atomic oxygen over time while anodized surfaces exhibited only small changes in their surface morphologies or in their visible and near-infrared reflectances. The graying of some smooth painted black surfaces is somewhat compensated for by an etching process that roughens these surfaces. If the exposure times are short, the etching action of the atomic oxygen may even improve the scattering properties of paints. However, removed material may show up on optical surfaces, degrading the optical performance, and it should be obvious that long exposure times would remove too much of the coating.

In general, however, the optical properties of surfaces degrade during longer exposure times and at higher atomic oxygen flux levels. Experiments in space on the long duration exposure facility (LDEF) have indicated how the properties of materials change after nearly six years of exposure to atomic oxygen. For Chemglaze (now called Aeroglaze) Z306 paint, for example, this exposure has led to a loss of the binder material¹¹⁷ and to loss of pigment. Some ion-textured surfaces were also flown as part of this experiment and showed high stability in the space environment.⁷⁹ The effects of space exposure on infrared baffle coatings were described by Blue and Perkowitz¹¹⁸ who showed that long-term exposure decreased the reflectance of these surfaces.

Outgassing

For space applications, the amount of outgassed material and its composition are of great importance. One of the primary dangers is that the outgassed products may form a film with undesirable properties on optical components. This hazard of space exposure is particularly relevant for cooled systems. For cryogenic systems, the continued outgassing of materials as they break down after long exposure to the space environment becomes an additional source of contamination on optical surfaces. Outgassing is also of concern for severe environments on the ground, though the correlation between vacuum testing and a behavior of black surfaces under severe heat in air is not well understood. For the context of this chapter, only vacuum outgassing measures will be considered. In either case, the condensation of volatile materials on optical components can lead to catastrophic system failure.

The two primary measures or tests of outgassing are the *total mass loss* (TML) and the *percent collected volatile condensable materials* (CVCN). An American Society for Testing and Materials (ASTM) Standard Test Method was developed for this purpose and is called E 595-77/84. Several NASA publications give values for TML and CVCN (Ref. 119 gives the most extensive listing; see also Ref. 120). For example, measurements at Goddard Space Flight Center of Chemglaze (now Aeroglaze) Z306 black urethane paint give a TML of 0.92 percent and a CVCN of 0.03 percent. Wood et al. shows the near-infrared absorption due to the outgassing products of Aeroglaze Z306 black paint. For the paint, the total mass loss was 2.07 percent. Ames²⁸ reports outgassing problems with Aeroglaze Z306 if exposed to temperatures greater than 40°C after a room temperature cure. However, many potential space material outgassing problems with this and other surfaces can be avoided completely with proper bakeout procedures. In 1993, a low volatile version of Z306 was being developed.

A third measure of outgassing relates to the water content of the material or surface. The *water vapor regained* (WVR) is a measure of the amount of water readsorbed/reabsorbed in 24 hours, while the sample is exposed to 25°C and 50 percent relative humidity (for extensive listings, see Ref. 119). This determination is done after the vacuum tests for TML and CVCN.

Particle Generation

Materials for space use must be able to withstand launch vibration without generating an excessive number of particles. This is not simply a coating or surface problem—it is a dynamics and substrate problem as well. The details of the baffle design and assembly play an important role in the generation

of particles. If a baffle surface is grossly distorted by the launch vibration, the adhesion of the black surface should not be blamed for the failure; the baffle design failed. A common source of particle generation is abrasion between coated baffle surfaces. For severe vibration environments and for fragile surfaces, special considerations are necessary.⁷⁰ Similarly, surfaces must be able to withstand thermal cycling without producing particles. As painted materials outgas and their composition is altered, particles in the paint can become dislodged.

Materials used in baffles can also get dusty or contaminated by exposure to outgassing products from the spacecraft. For an accessible ground-based black surface, normal cleaning procedures appropriate to the surface may be used, if they do not damage the texture or optical properties of the surface. For a baffle system already in space, the cleaning procedures developed to clean mirrors in situ may apply. These techniques include laser cleaning, plasma or ion cleaning, and jet snow cleaning (see, e.g., Ref. 121, and other papers in Ref. 109).

Baffle Material for Extreme Environments

One important form of degradation for earth-based materials involves the exposure of the surface to the elements of temperature cycling, sunlight, humidity, and other factors of its operating environment. This combination is a severe test of many surfaces. Unfortunately, there is very little data published on the environmental degradation of black surfaces that are exposed to the outdoor environment, or to high temperatures, such as might be found in a closed automobile in the summer. Often, a black surface that is very ordinary in its optical properties is chosen for use in an extreme environment. The reason is that it was the first surface that did not peel or fade significantly.

Space-based surfaces must withstand severe launch vibrations, temperature extremes, collisions with space debris and micrometeoroids, and exposure to ultraviolet radiation. Military missions have also added the requirement that these materials be hardened against nuclear or laser threats. For example, these baffle surfaces and their associated optics must still perform after exposure to high-power lasers.

In the simplest sense, these black surfaces must be able to either selectively reject laser radiation (a difficult requirement) or be able to absorb significant amounts of energy and not degrade. Most of the coatings have been developed for application to substrates, which can survive high temperatures, such as molybdenum, beryllium, tungsten, carbon/carbon, and titanium, because of their ability to withstand high flux densities of laser light. Generally, painted surfaces are completely inadequate. Even anodized surfaces on aluminum may not be adequate at these flux densities.

6.5 OPTICAL CHARACTERIZATION OF BLACK SURFACES

The baffle designer must be concerned about the optical properties of black surfaces over a large range of wavelengths and angles of incidence. Wolfe¹⁶ summarizes the relevant design data. Performance measurements, such as measurements of the specular reflectance and total hemispherical reflectance over a range of wavelengths are now supplemented by a more complete characterization of these optical surfaces: the bidirectional reflectance distribution function (BRDF).¹²² BRDF has units of inverse steradians and is usually measured at several discrete wavelengths, which are easily generated by the most common lasers. In particular, measurements at the HeNe laser wavelengths of 0.6328 μm and at the CO₂ wavelength of 10.6 μm are particularly valuable for visible and infrared applications. There is a notable shortage of BRDF data in the ultraviolet, from 1 to 4 μm , and for far-infrared wavelengths.

BRDF measurements at the same wavelength have the additional advantage that they are directly comparable. As Smith and Wolfe¹²³ point out, the directional reflectivity of a surface is convolved with the instrument function (a weighted mean of the detector solid angle), a unique characteristic of each reflectometer. An approximate deconvolution of the instrument function can be easily

performed by dividing the measured reflectance by the projected measurements, which are instrument independent as mentioned in ASTM Standard E12.09. However, simple specular reflectance measurements and spectra made on different instruments must first be divided by the projected detector solid angle before being compared. Several studies^{89,90} have detailed the variation between BRDF measurements made in different laboratories. With the advent of the ASTM standard, greater congruence among BRDF measurements in the future is probable.

Because of the value of BRDF measurements in comparing the scattering characteristics of surfaces, a large set of BRDF data of black baffle materials that have traditionally been used in optical devices of all types, as well as the scattering functions of some new materials are presented in this chapter. These measurements may be compared to previous measurements made on other instruments.⁶⁷

The specular and total hemispherical reflectance data thus become most valuable in describing some general characteristics of a surface over a wide wavelength range, while the BRDF data at a wavelength of particular interest allow the designer to characterize the angular scattering of the surface for angles of incidence of interest and for view angles from the detector. The total integrated scatter also is extremely useful for surface roughness determination.⁹³

This chapter provides the optical designer a useful summary of data. Sufficient data are presented to allow a determination of which surfaces might be appropriate for a given task. The data generally can be divided into two forms, the first being specular and hemispherical reflectance data over a wavelength range. These data allow the designer to choose a surface with the right basic properties over the wavelength range of interest. BRDF measurement provides a second form of data. The data are usually presented for a few characteristic angles of incidence (5°, 30°, 45°, and 60° are the most common) at each wavelength. These BRDF measurements are particularly useful as input to optical modeling or stray light modeling programs.

In stray light analysis programs such as APART,¹²⁴ ASAP,¹²⁵ and GUERAP,¹²⁶ as well as many other modeling programs, it is the mathematical value of the bidirectional reflectance distribution function that is used by the program models. In these programs, a library of BRDF functions for the most common surfaces is available for use by the analyst. For surfaces that do not explicitly have a BRDF associated with them in the code, the BRDF values can be input.

In most baffle designs, surfaces are used at angles not represented by the available BRDF measurements. Some surfaces are even used at angles close to grazing incidence where few measurements have been made. For these areas, the data available for the models are very limited, and often generous extrapolations are needed. A smooth transition between measurements at different angles or at very large angles of incidence, or from one wavelength to another is often assumed. The more unusual or nonrandom the surface, the more these assumptions or extrapolations should be questioned. Examples of how BRDF measurements are made and used in system designs and analysis are found in, for example, Refs. 3, 43, and 127 to 130.

Black surfaces are also used in radiators used for the thermal control of spacecraft. The hemispherical reflectivity can provide a general indication of the optical properties for a thermal analysis. The normal emittance at room temperature is also an important quantity for spacecraft design. Table 7 gives the hemispherical reflectance and normal emittance of some surfaces. Radiatively cooled surfaces are important for space infrared telescopes and, if designed properly, can provide a reliable means of reducing background radiation. Because of on-orbit degradation, the radiators must be sized for the appropriate performance at the end of the mission.

For these applications, the absorptivity and emissivity of surfaces as a function of wavelength and temperature must be well understood. Thus, room temperature measurements of normal emittance may not be very useful if the surface is to be used at cryogenic temperatures. Stierwalt^{41,132} has contributed greatly to the measurement of far-infrared emissivity of black materials. Clarke and Larkin¹³³ describe measurements on a number of black surfaces suitable for radiometry detectors where a low reflectance is needed, independent of wavelength. For opaque surfaces, the spectral emittance at a given wavelength is equal to the spectral absorptance, which is equal to one minus the spectral reflectance.

These parameters are very sensitive to coating thickness and may change with exposure to ultraviolet radiation from the sun, with damage from solar electron and proton radiation, damage from space debris and micrometeoroids, and with outgassing. In practice, the solar absorptance is measured

TABLE 7 Hemispherical Reflectance and Emittance of Some Common Spacecraft Surfaces

Surface	Thickness	Hemispherical reflectance at 546 nm (incidence angle, degrees)					Normal Emittance (300 K)
		20	40	60	70	80	
Cat-a-lac Black 463-3-8	0.002"	0.054	0.061	0.096	0.141	0.223	0.891
Black Velvet 401-C10	0.025"	0.035	0.036	0.047	0.058	0.085	0.911
Chemglaze Z306	0.005"	0.048	0.052	0.074	0.098	0.144	0.915
Chemglaze Z306 with silica powder	0.031	0.033	0.040	0.050	0.070	0.112	0.912
Chemglaze Z306 with 3M glass microballoons left on 20- μ m screen	0.040	0.048	0.068	0.090	0.130	0.920	0.920
Chemglaze Z306 with 3M glass microballoons left on 44- μ m screen	0.010"	0.036	0.040	0.055	0.070	0.097	0.923
Chemglaze Z306 with 3M glass microballoons left on 63- μ m screen	0.011"	0.037	0.040	0.053	0.064	0.084	0.923
Scotchlite Brand Reflective Sheetting #3285 Black		0.059	0.062	0.095	0.156	0.318	0.909
Scotchlite Brand Reflective Sheetting "C" Black W/A #234		0.146	0.130	0.129	0.135	0.157	0.846
NiS-dyed Anodized Al from Light Metal Coloring Corp.		0.056	0.061	0.097	0.151	0.333	0.912
Sandoz-Bk Organic Black anodized Al from Almag Co.		0.041	0.048	0.064	0.082	0.129	0.884
Black Chrome anodized Al from Goddard SFC		0.028	0.031	0.053	0.083	0.156	0.625
NiS-dyed anodized Al from Langley Research Center		0.050	0.053	0.083	0.138	0.253	0.915
NiS-dyed anodized Al, substrate blasted with glass shot, from Light Metals Coloring Corp.		0.045	0.049	0.073	0.100	0.146	0.920

Source: Data courtesy of J. Heaney, Goddard Space Flight Center.¹³¹

by obtaining the reflectance from 0.3 to 2.4 μm , the region containing 95 percent of the sun's energy. The normal emittance is determined using an infrared spectrophotometer over the wavelength range 5 to 35 μm , where close to 90 percent of the energy of a 300 K blackbody is emitted (see Table 6). While these measurements are useful for design of radiators that operate near room temperature, further measurements are required in order to predict a surface's thermal emittance at lower operating temperature.

6.6 SURFACES FOR ULTRAVIOLET AND FAR-INFRARED APPLICATIONS

Many surfaces are multipurpose and are useful for a variety of wavelengths. Before discussing the properties of these surfaces, it may be instructive to look at several surfaces that are used for ultraviolet and far-infrared applications. In these wavelength regions, the performance of black surfaces is often critical to the performance of the optical instrument.

Black Surfaces for Ultraviolet Applications

For space telescopes or instruments operating in the ultraviolet, surfaces such as Martin Black have been used successfully. However, very few measurements of candidate baffle material have been made at wavelengths less than 0.3 μm (300 nm) and even fewer have been made at shorter ultraviolet wavelengths. Jelinsky and Jelinsky⁴⁰ studied the performance characteristics of a variety of surfaces and surface treatments for use on baffle materials in the Extreme Ultraviolet Explorer. They determined some scattering characteristics of several materials at 30.4, 58.4, and 121.6 nm.

Heaney¹³¹ has made some pioneering bidirectional measurements at 15°, 45°, 75°, and 85° at 58.4 nm (584 Å) of candidate baffle materials for space instrument applications. Interestingly, two of his samples are "bare" aluminum with no applied coating. The natural oxide present on uncoated aluminum samples is strongly absorbing at this wavelength and, by itself, reduces the reflectances to below 5 percent at near-normal incidence. The baffle surface candidates are described here and their optical properties are presented in the figures.

Roughened aluminum: Aluminum 6061 alloy, sheet finish, sandblasted to an unspecified degree of surface roughness. Uncoated (Fig. 11).

Gold iridite: Aluminum with a chromate surface coating produced in a room temperature chemical conversion process (Fig. 12).

Black nickel: Aluminum with a black nickel surface that was produced in an electroplating process with NiS and ZnS (Fig. 13).

Copper black: An aluminum substrate coated with approximately 80 to 100 nanometers of evaporated Cu to provide a conductive layer and then overcoated with an evaporated Cu black deposited at a chamber pressure in the 10^{-3} torr range (Fig. 14).

Black paint: An aluminum substrate painted with a carbon-loaded black silicone paint (NSB69-82) (Fig. 15).

Aluminum, grooved and blazed: An aluminum plate, 6061 alloy, with grooves cut at a blaze angle of about 20°, with a period of 1.1 mm, a depth of 0.4 mm, and otherwise uncoated (Figs. 16, 17, and 18).

Surfaces for Far-Infrared (>30 μm) Use

In the far-infrared wavebands, many surfaces such as paints, which have good visible absorptance, lose their absorbing and scattering characteristics. Although measurements have only been made at wavelengths longer than 15 μm on a few surfaces, the data suggest that most surfaces do not

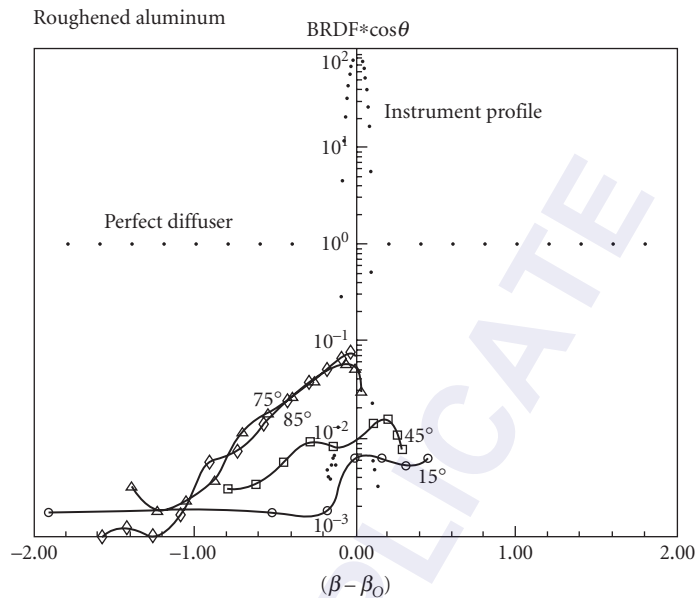


FIGURE 11 BRDF of roughened aluminum. Instrument profile and perfect diffuse reflector shown for reference. Wavelength is 58.4 nm. (From Heaney.¹³¹)

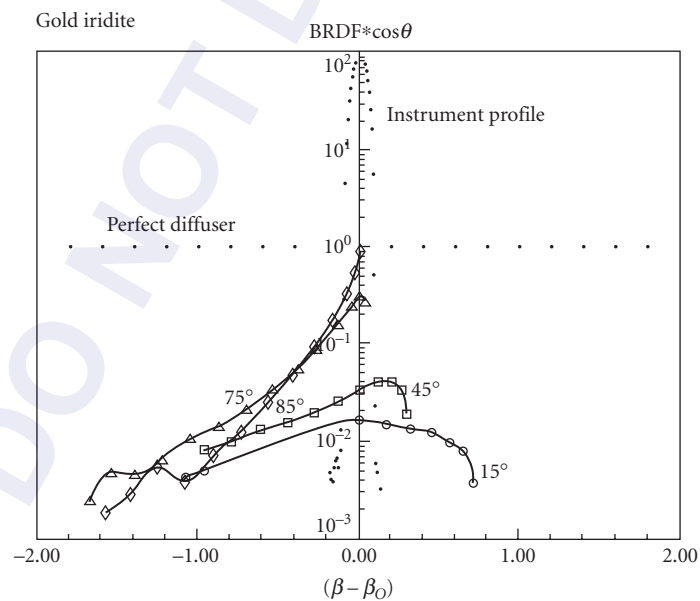


FIGURE 12 BRDF of gold iridite. Instrument profile and perfect diffuse reflector shown for reference. Wavelength is 58.4 nm. (From Heaney.¹³¹)

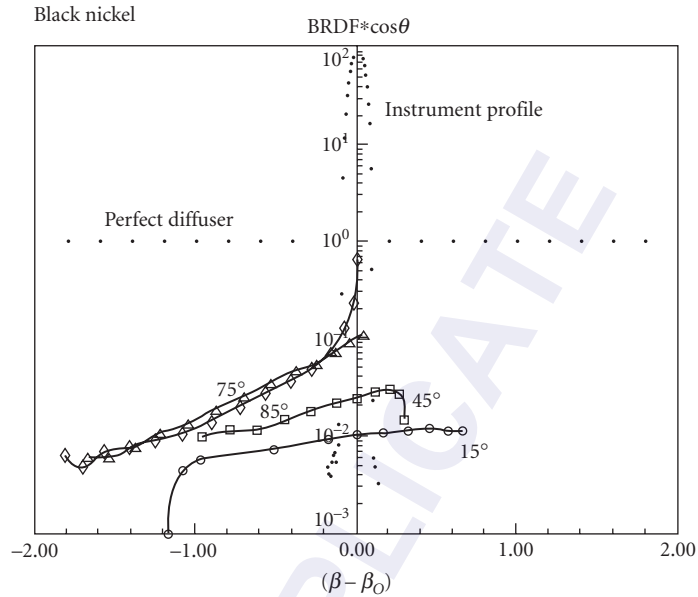


FIGURE 13 BRDF of black nickel. Instrument profile and perfect diffuse reflector shown for reference. Wavelength is 58.4 nm. (From Heaney.¹³¹)

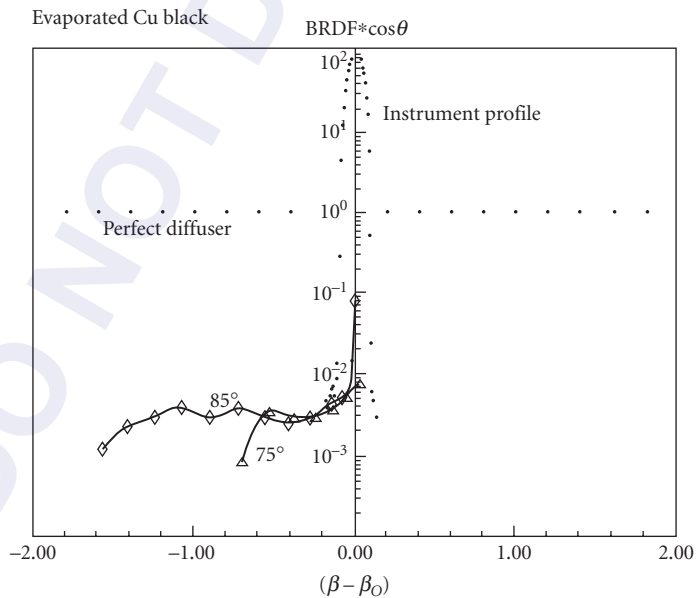


FIGURE 14 BRDF of evaporated Cu black. Instrument profile and perfect diffuse reflector shown for reference. Wavelength is 58.4 nm. (From Heaney.¹³¹)

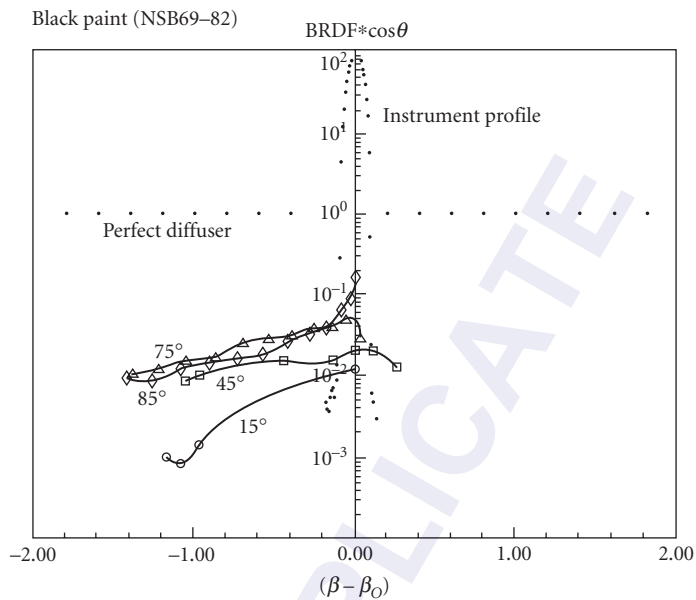


FIGURE 15 BRDF of NSB69-82 black paint. Instrument profile and perfect diffuse reflector shown for reference. Wavelength is 58.4 nm. (From Heaney.¹³¹)

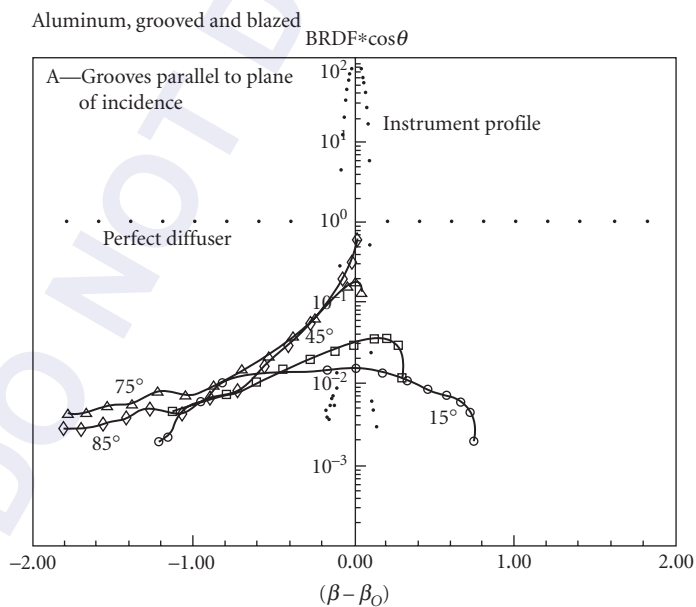


FIGURE 16 BRDF of aluminum, grooved and blazed, with grooves parallel to angle of incidence. Instrument profile and perfect diffuse reflector shown for reference. Wavelength is 58.4 nm. (From Heaney.¹³¹)

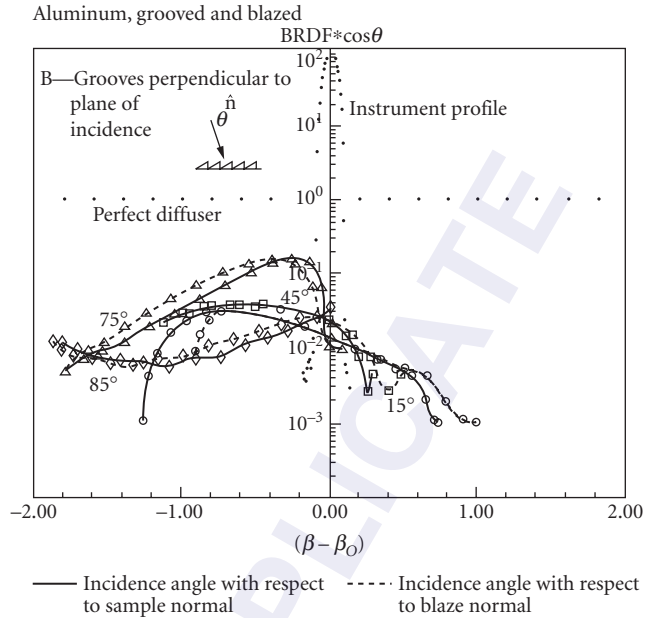


FIGURE 17 BRDF of aluminum, grooved and blazed, with grooves perpendicular to angle of incidence. Instrument profile and perfect diffuse reflector shown for reference. Wavelength is 58.4 nm. (From Heaney.¹³¹)

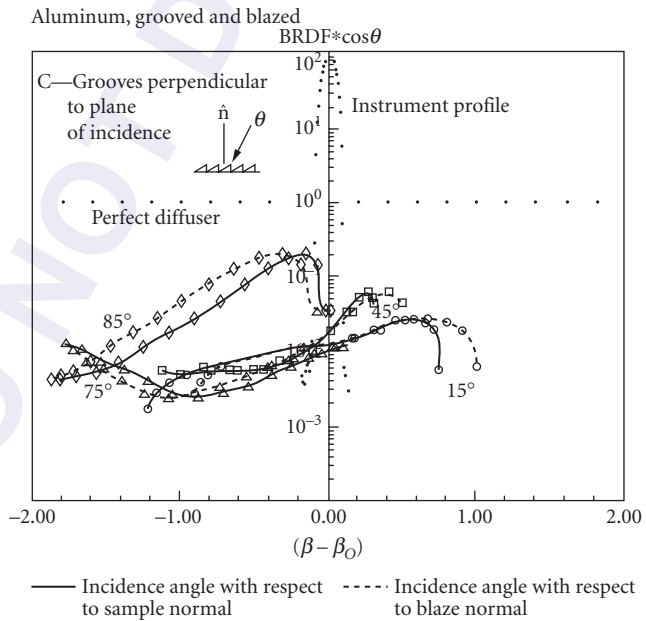


FIGURE 18 BRDF of aluminum, grooved and blazed, with grooves perpendicular to angle of incidence. Instrument profile and perfect diffuse reflector shown for reference. Wavelength is 58.4 nm. (From Heaney.¹³¹)

perform well as “blacks” for long wavelength operation.^{41,153} Smith³² points out that many coatings become transparent at longer wavelength and the substrate roughness enters strongly into the scattering characteristics. An example of this is Infrablack, a rough anodized surface that was previously described. The approach has been carried further in Ames 24E.¹⁵⁴ A number of approaches and surfaces have been developed to create effective surfaces at these wavelengths.

A model has been constructed³³ to describe the effect of coating roughness and thickness on the specular reflectance. Smith and Howitt¹⁹ survey materials for infrared opaque coatings (2 to 700 μm). They highlight the value of silicon carbide and carbon black for far-infrared absorption. Smith^{32,33,39,46,57,72,134-137} provides much of the reference work in this area.

The optical properties of the several infrared coatings and surfaces are given in Figs. 19 to 24. For applications in laboratories and in calibration facilities, surfaces like velvet, commando cloth, and even neoprene may be used. The long-wavelength measurements on these surfaces are given in Figs. 25 to 32. The long-wavelength coatings may be grouped into a few basic categories discussed here. As this area is one of rapid development for space infrared telescopes and space radiators, there will likely be many new approaches in the next few years.

Multiple-Layer Approach For reducing the far-infrared reflectance of baffles, a technique of using antireflection surfaces composed of multiple layers with a different refractive index in each layer was discussed by Grammar et al.⁵⁵ For such a coating to be effective, it must match the optical constants of free space with the large constants of the baffle substrate.

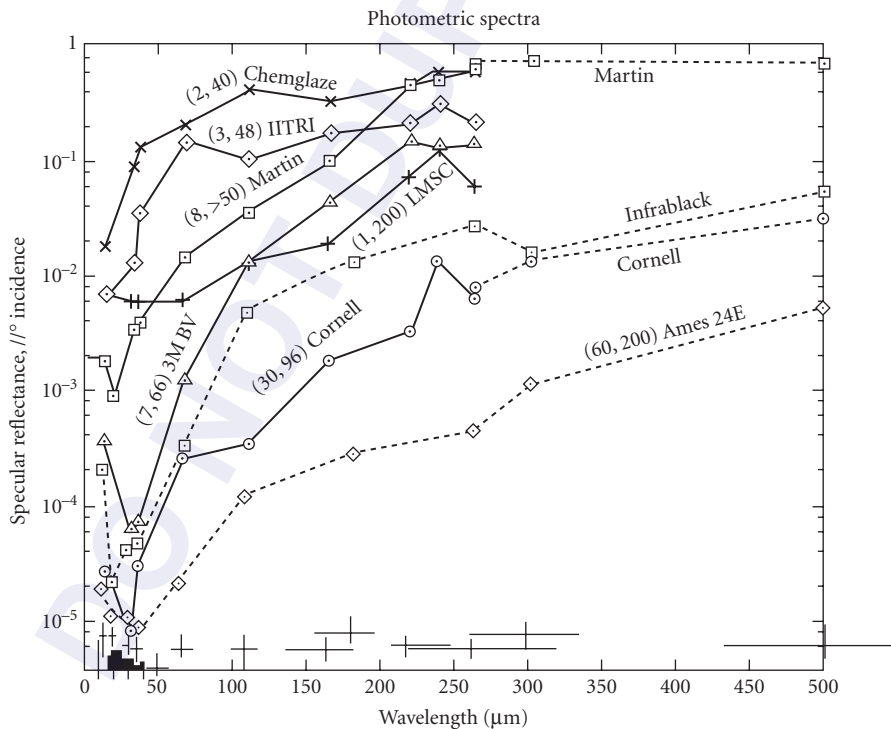


FIGURE 19 Specular reflectance spectra of eight black coatings. Measurements of coating roughness and thickness in micrometers are shown in parentheses before the name of each coating. Filter passbands are shown by horizontal error bars. The absorption bands of amorphous silicates are indicated by the solid histogram in lower left end of the graph. (From Smith.^{32,33,46})

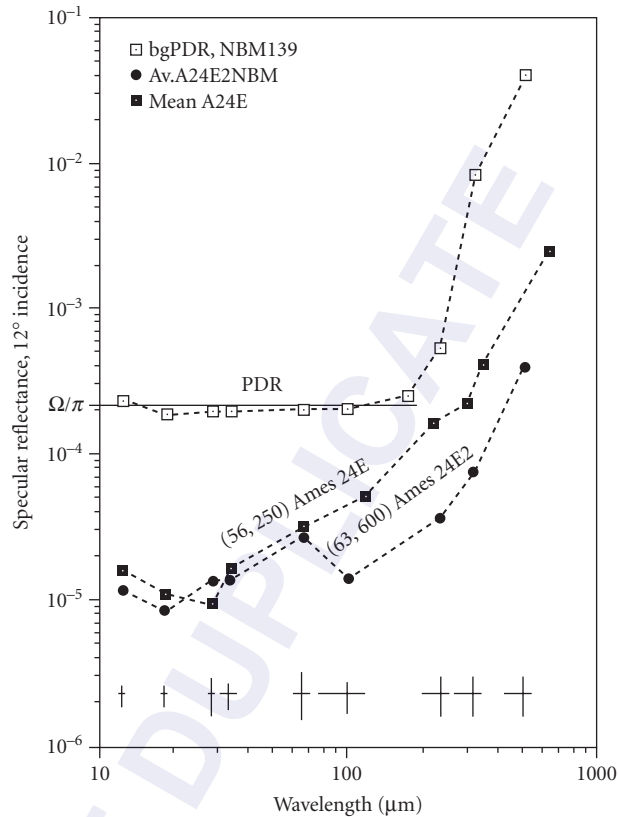


FIGURE 20 Specular reflectance spectra of PDR surface and two infrared black coatings. Dashed lines connect data points of a set, and the solid line indicates the theoretical value of reflectance for a perfect lambertian surface. The coating roughness and thickness are given in parentheses, and filter passbands are indicated by horizontal error bars. (From Smith.^{136,137})

Teflon Overcoat A polytetrafluoroethylene (Teflon) spray-on lubricant was used as an antireflection overcoat for opaque baffle surfaces in the far-infrared and submillimeter by Smith.⁵⁰ A thick Teflon overcoat created by spraying Teflon Wet Lubricant reduced the specular reflectance at millimeter and infrared wavelengths by a factor of two. The refractive index of the coating plays a more significant role than the thickness.

Cornell Black Houck⁴⁹ created a far-infrared black paint based on adding large particles of grit (silicon carbide #80 and #180 grit) to 3M Black Velvet Nextel 101-C10. This paint could be repeatedly cycled to cryogenic temperatures with no flaking or peeling. The far-infrared properties of Cornell Black are described by Smith.^{32,33} Since this 3M paint is no longer available, Cornell Black is not now generally used. The ECP-2200 (now called MH2200) replacement for 3M Black Velvet was studied by Smith and Howitt¹⁹ in an effort to create a paint similar to Cornell Black. Houck⁴⁹ gives complete instructions for making the Cornell Black paint.

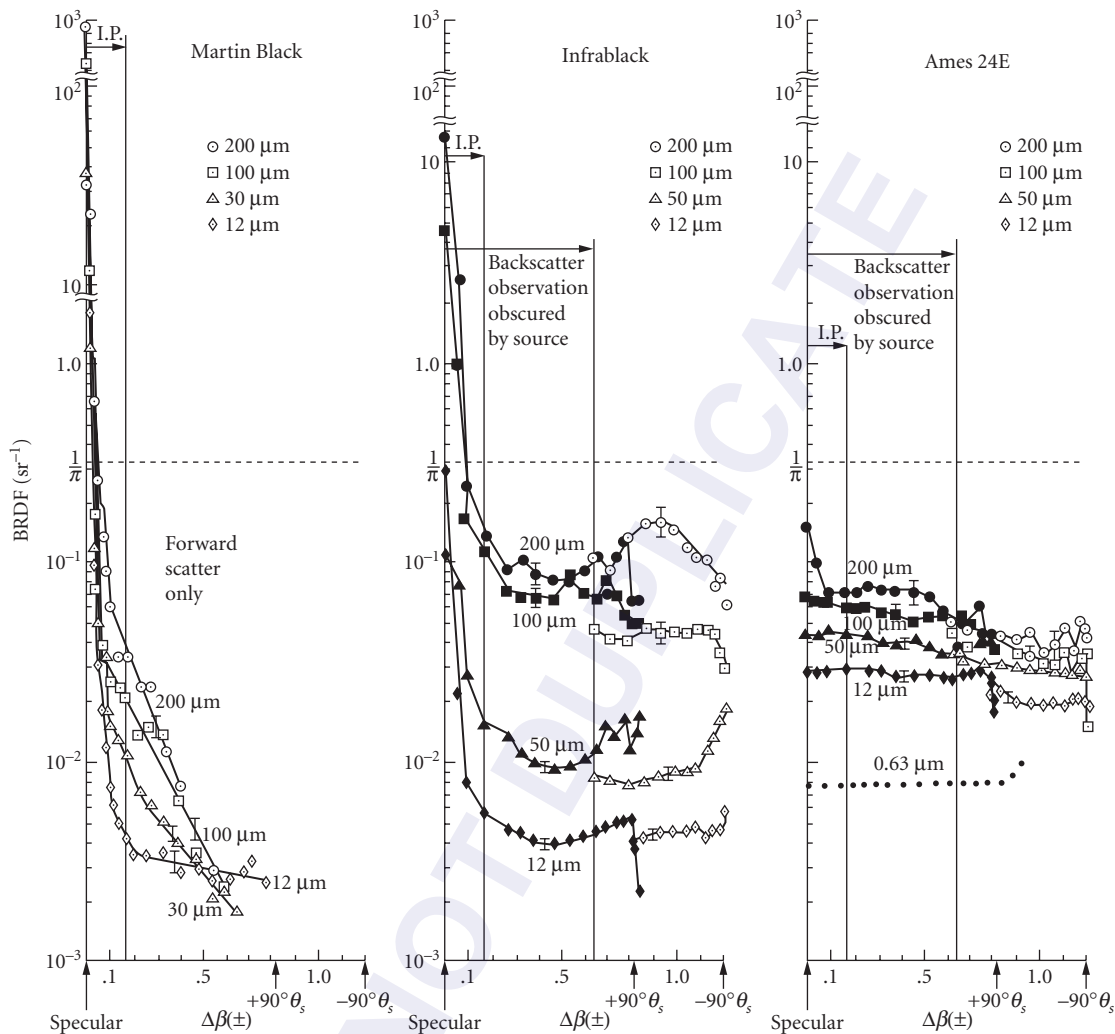


FIGURE 21 BRDFs of Martin Black, Infrablack, and Ames 24E. Note the different breaks in the ordinate scale of the Martin Black measurements used to compress its specular reflectance to fit the scale of the figure. Forward scatter data are solid symbols; backscatter data are open symbols. Forward and backscatter measurements at $0.6328\ \mu\text{m}$ of Ames 24E are by D. Shepard of Martin Marietta (Denver) and are shown schematically by the dotted line in the right hand figure. (From Smith.⁴⁶) Smith¹³⁷ recalibrated this data; the relative comparison of these coatings remained essentially unchanged.

Infrablack Infrablack is a black anodized surface generally applied to 6061 Aluminum.^{31,64,138} It is highly suitable for space radiators, infrared baffles, and for applications where a high emissivity is desired. The specular reflectance of the Infrablack is conservatively one order of magnitude less than that of Martin Black across the spectral region from 12 to 500 μm . The reduced reflectance of Infrablack is attributed to a large increase in the roughness of its substrate prior to anodization. The diffuse reflectance characteristics of the Infrablack are different from Martin Black at the longer wavelengths. At 12 and 50 μm , the BRDF of Infrablack is considerably less than that of Ames 24E, while at 100 and 200 μm they are comparable.

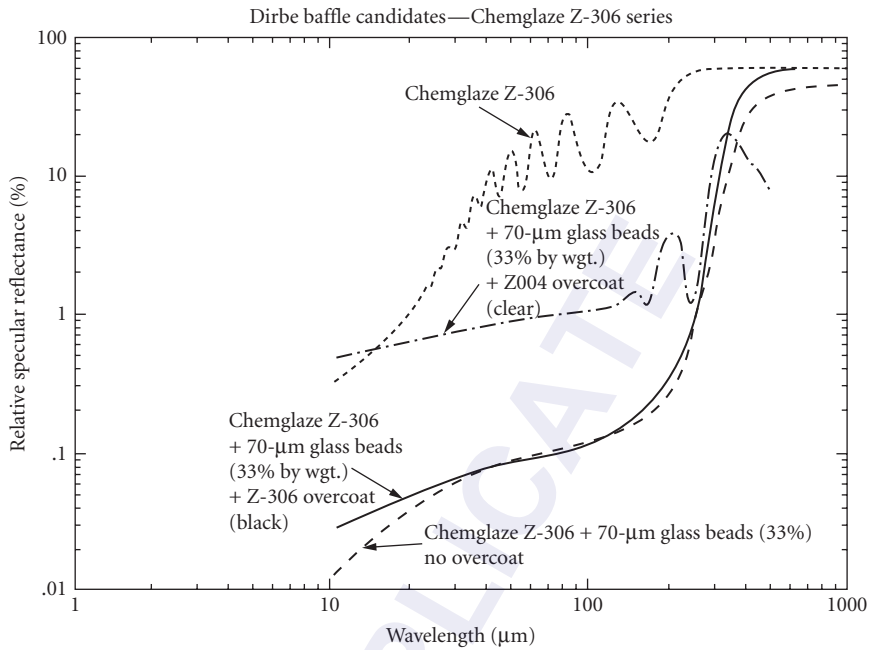


FIGURE 22 Specular reflectances of Diffuse Infrared Background Explorer candidate baffle surfaces based on Chemglaze (now Aeroglaze) Z-306, over the wavelength range 1 to 1000 μm . (From Heaney.¹³¹)

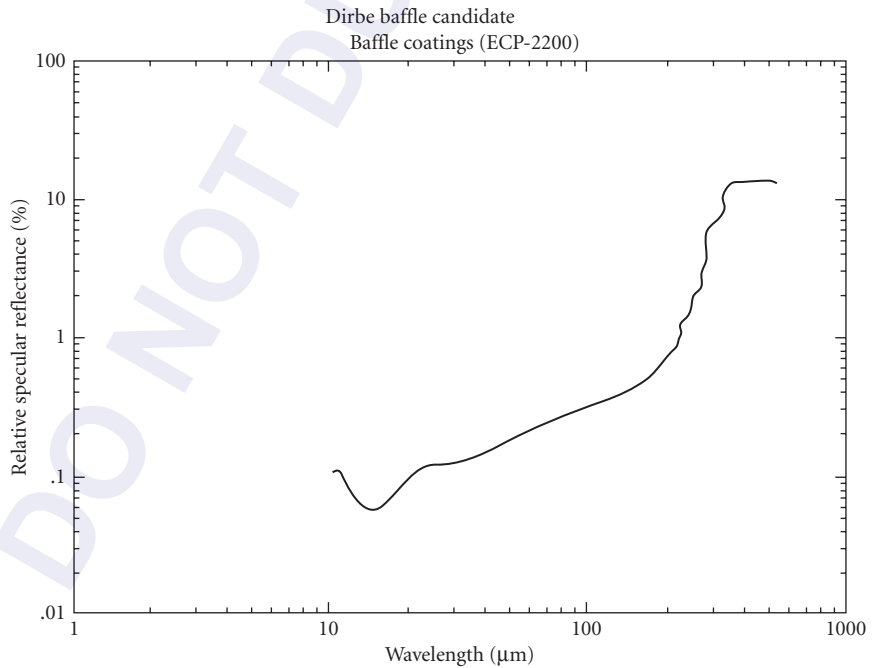


FIGURE 23 Specular reflectance of Diffuse Infrared Background Explorer candidate baffle surface based on ECP-2200 (now MH2200) black paint over the wavelength range 1 to 1000 μm . (From Heaney.¹³¹)

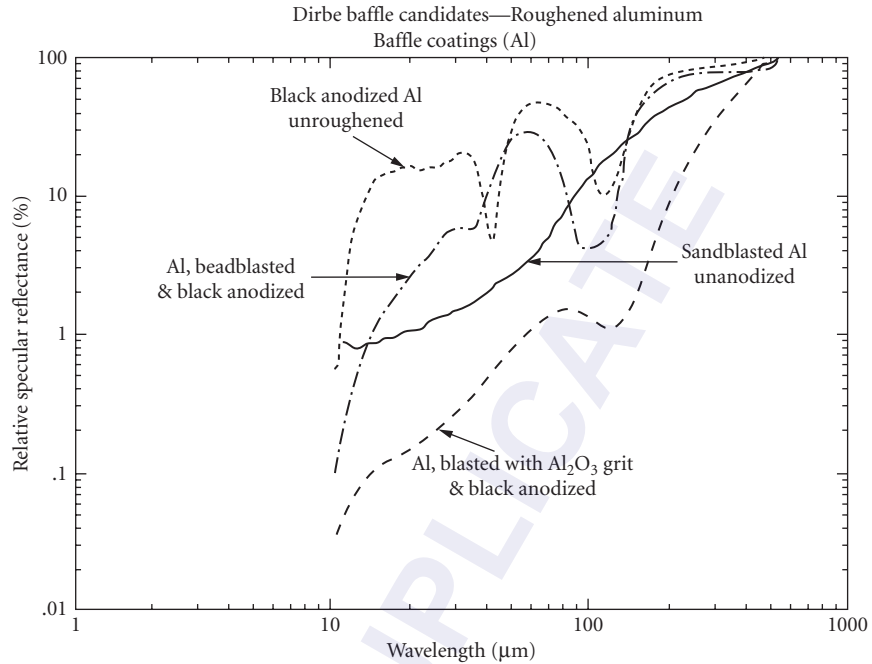


FIGURE 24 Specular reflectance of Diffuse Infrared Background Explorer candidate baffle surface based on roughened aluminum surfaces over the wavelength range 1 to 1000 μm . (From Heaney.¹³¹)

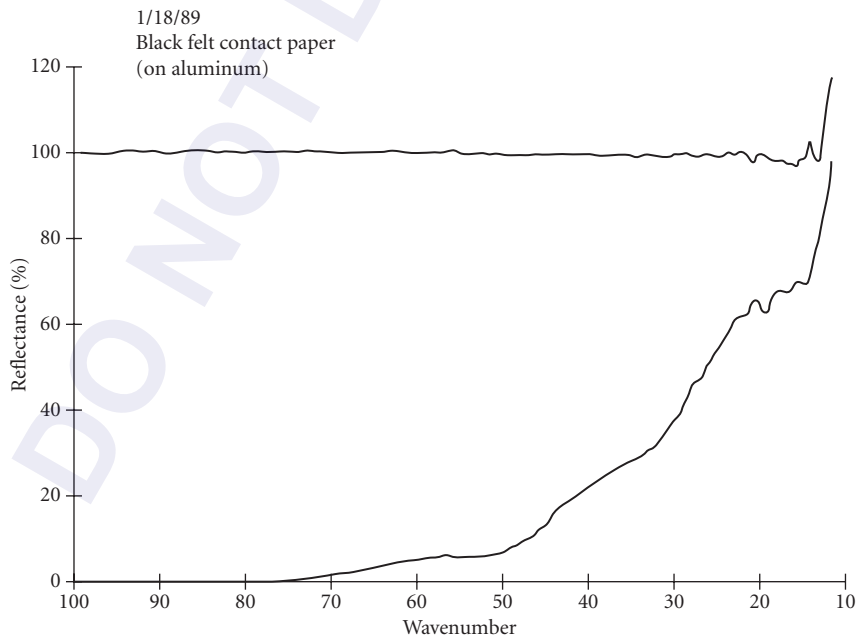


FIGURE 25 Reflectance of black felt contact paper from 100 to 10 wave numbers (100 to 1000 μm). (From Heaney.¹³¹)

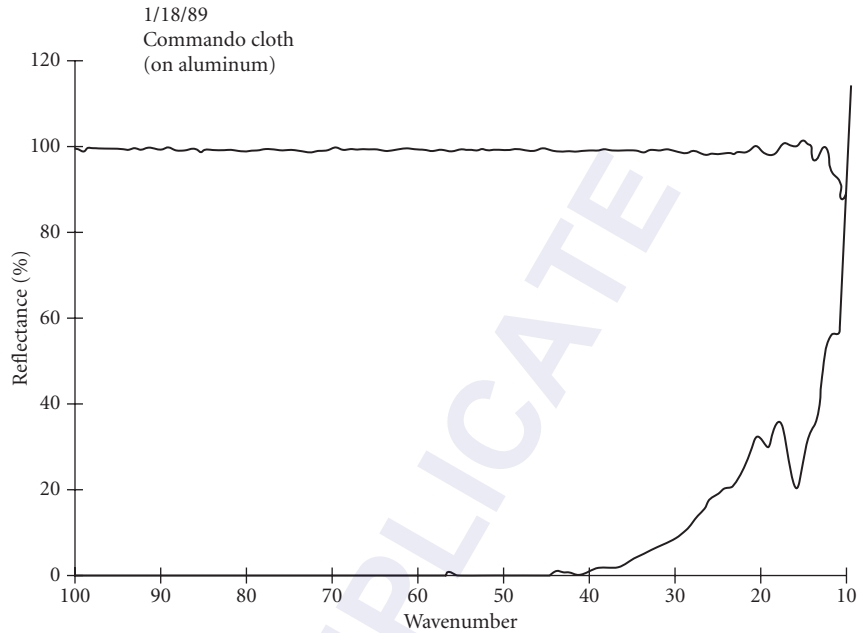


FIGURE 26 Reflectance of commando cloth on aluminum from 100 to 10 wave numbers (100 to 1000 μm). (From Heaney.¹³¹)

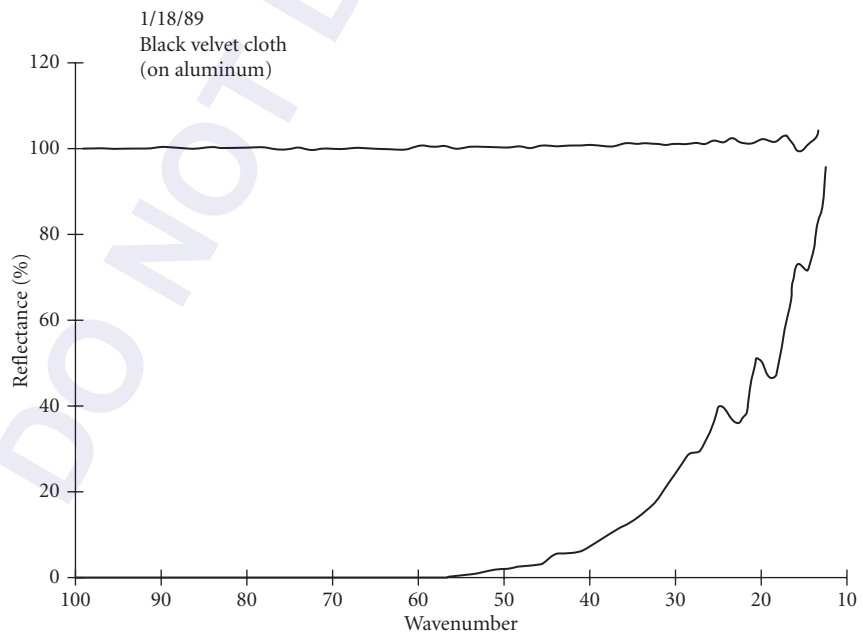


FIGURE 27 Reflectance of black velvet cloth on aluminum from 100 to 10 wave numbers (100 to 1000 μm). (From Heaney.¹³¹)

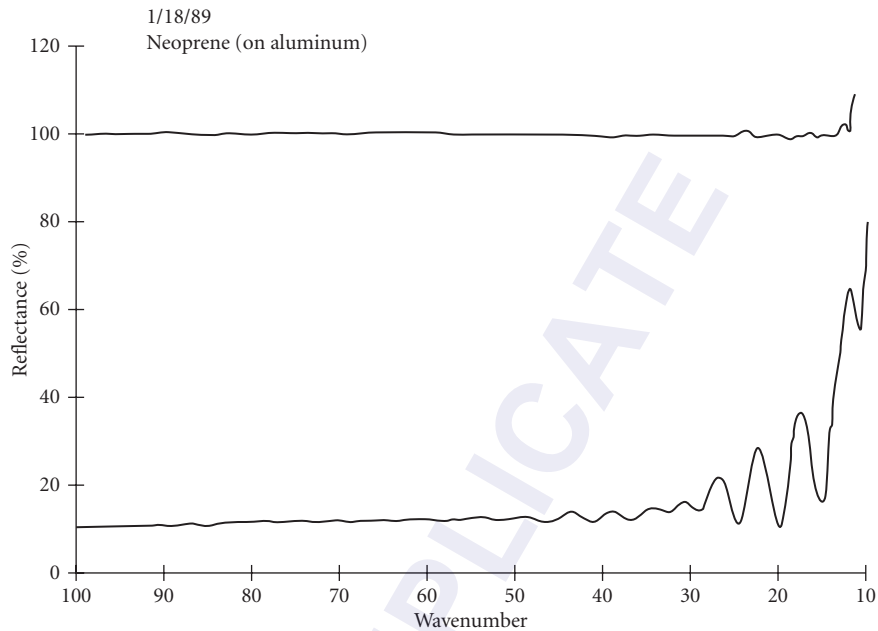


FIGURE 28 Reflectance of neoprene on aluminum from 100 to 10 wave numbers (100 to 1000 μm).
(From Heaney.¹³¹)

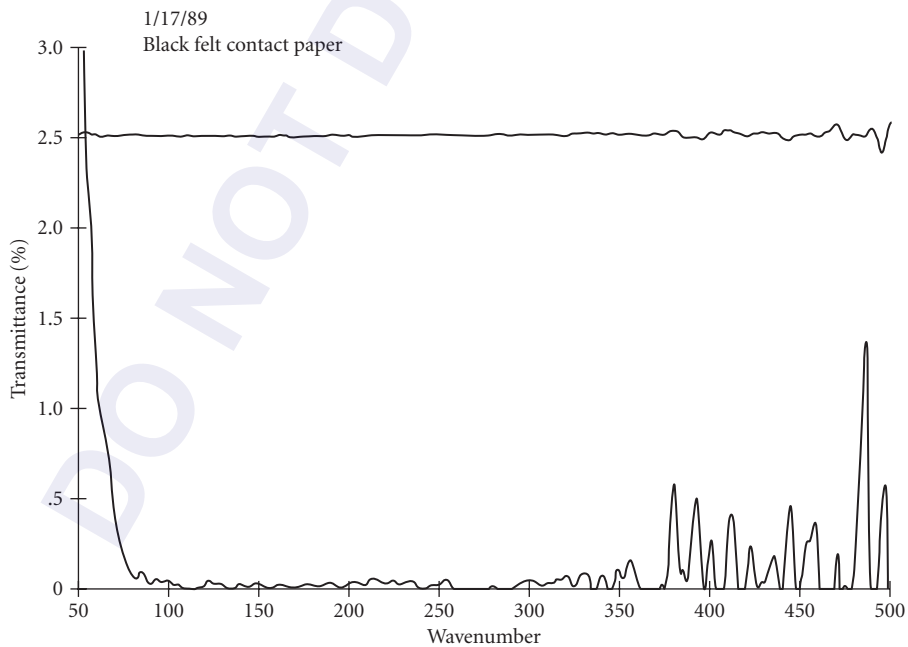


FIGURE 29 Transmittance of black felt contact paper from 50 to 500 wave numbers (200 to 20 μm).
(From Heaney.¹³¹)

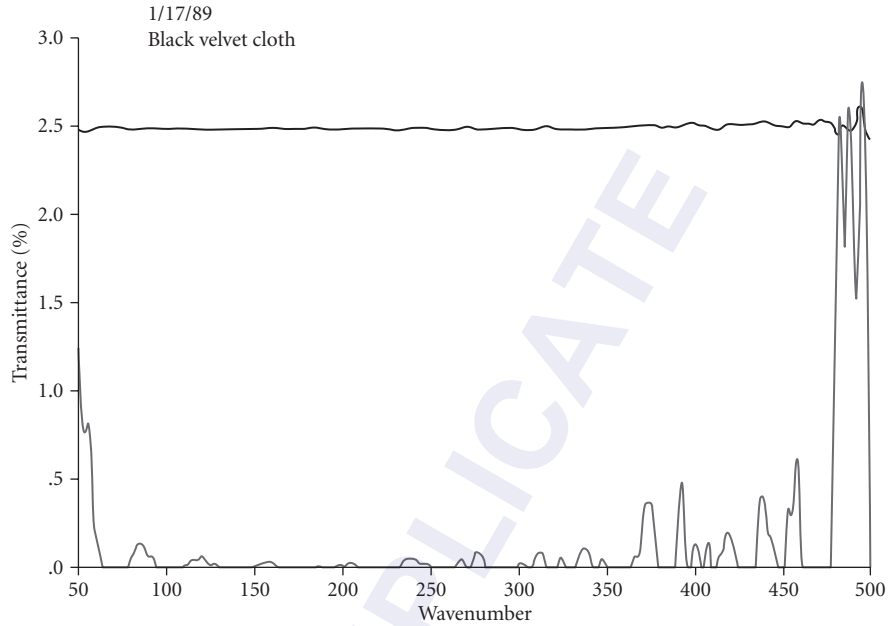


FIGURE 30 Transmittance of black velvet cloth from 50 to 500 wave numbers (200 to 20 μm). (From Heaney.¹³¹)

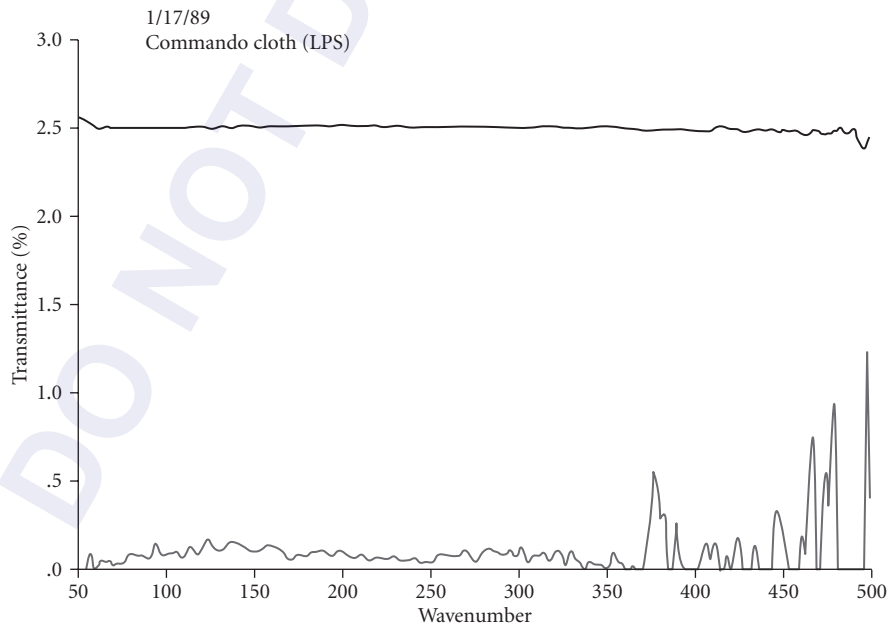


FIGURE 31 Transmittance of commando cloth (lps) from 50 to 500 wave numbers (200 to 20 μm). (From Heaney.¹³¹)

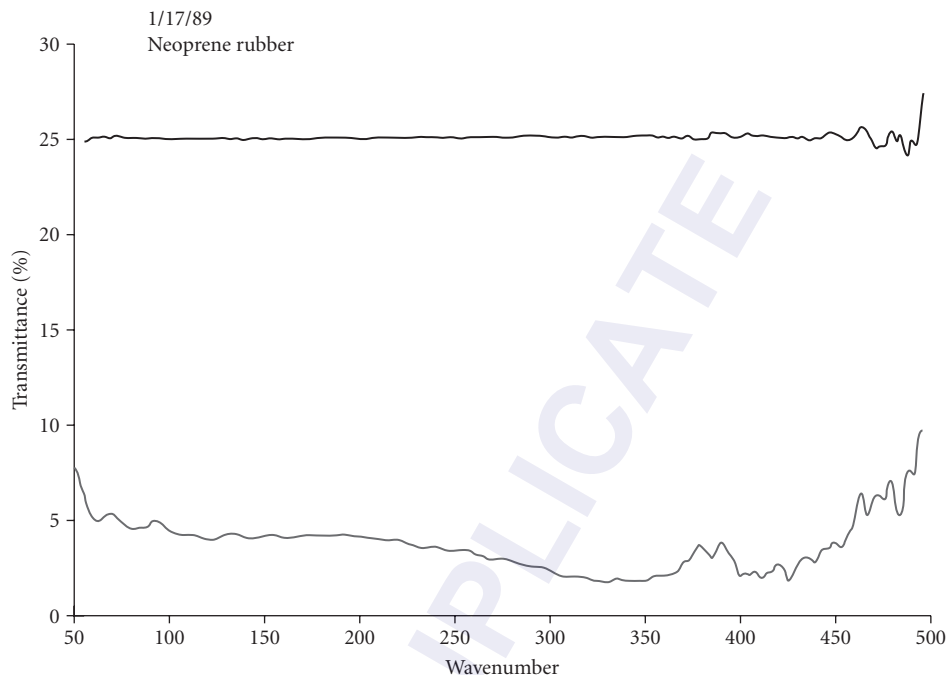


FIGURE 32 Transmittance of neoprene rubber from 50 to 500 wave numbers (200 to 20 μm). (From Heaney.¹³¹)

Ames 24E and Ames 24E2 Ames 24E is a new coating developed at NASA's Ames Research Center.⁴⁶ The coating is a very rough, very thick, highly absorbing coating composed largely of SiC grit and carbon black, and has very low far-infrared reflectance.⁴⁶ BRDF measurements show that this surface is a model lambertian surface at a wavelength of 10.6 μm at near-normal incidence. Ames 24E is nearly lambertian at longer wavelengths, even at 200 μm .

Ames 24E2 is a formulation using ECP-2200, carbon black, silicone resin, and #80 silicon carbide grit. The coating was developed for far-infrared space telescopes and can pass severe cryogenic cycling and vibration tests.³⁴ The coating is considerably less reflective at 100 μm than Ames 24E and has excellent absorption out to about 500 μm . Smith¹³⁷ has also developed a highly lambertian, diffuse reflectance standard, which he calls PDR, for perfect diffuse reflector. It shows lambertian behavior from 10 to about 100 μm .

6.7 SURVEY OF SURFACES WITH OPTICAL DATA

A Note on Substrates and Types of Coatings

Many different substrates are employed for baffle and vane surfaces. In particular, for space or cryogenic use, the 6000 series of alloys of aluminum (particularly 6061) are popular. They represent the baseline by which other materials are judged. Titanium and graphite epoxy are also widely used, and beryllium and carbon fiber substrates are becoming more common. For high-temperature applications or for resistance to lasers, materials with high melting points such as nickel, molybdenum, and carbon/carbon are used.

The distinction between substrate and surface or surface coating is blurry. When paints or other coatings are applied to a substrate, the distinction between the two is very clear. However, many surfaces today involve an extensive alteration of the substrate before the “coating” is applied. They may chemically alter the substrate to create the surface coating, which is an integral part of the substrate. In these cases, the use of the word substrate only serves as a description of the material that must be altered to create the black surface.

The classification of surfaces by substrate and/or type of coating or surface treatment is only one of the many ways available. However, it is a very useful form of classification, and will be used here.

6.8 PAINTS

Paints are easy to apply and their availability and usually low price make them ideal for laboratory experiments and instruments. However, their high outgassing rates and, in some cases, the large number of particles that can be liberated from them compromise their use in space-based applications. However, there are many space-qualified paints that have been flown in space. Notable among them is Aeroglaze Z306, and the derivatives of this basic paint system that have been developed at Goddard Space Flight Center.

The way that paint may be applied can greatly affect the optical properties. The most obvious are the number of coats and the thickness; the substrate preparation also plays a critical role. Brown²⁷ describes the effects of number of coats on optical performance for three diffuse paints. The importance of developing a procedure for painting that incorporates the paint manufacturers' application recommendations cannot be overstated. This procedure must also include process control requirements, cleanliness requirements, application process and surface preparation requirements, safety requirements, quality assurance requirements, and storage shipping requirements, as McCall¹⁸ describes in detail. A short description of a few of the more common paints is given here. Table 2 gives further reference data on these and other paints.

3M Paints and Their Current Derivatives

3M Nextel Black Velvet Several varieties of 3M Black Nextel paint from Minnesota Mining and Manufacturing Company have seen extensive use in ground-based instruments and were considered as standards for black coatings. These paints consisted of pigments with approximately 20 percent carbon black and 80 percent silicon dioxide. However, most of them are presently unavailable. Those that are available are manufactured under license by other companies. The optical properties of the original paint are given in Fig. 33.

The paint labeled 3M 101-C10 historically was one of the most widely used laboratory black paints; however, it is no longer made. Redspot Paint and Varnish (Evansville, Indiana) now makes a similar paint, but uses resin instead of glass microspheres, making it unsuitable for space applications according to Ames.²⁸ Two other paint varieties in the same line were 3M401-C10 and 3M3101 (see Figs. 34 and 35), apparently available from Red Spot Paint and Varnish.

The 3M Nextel Black Velvet was replaced by 3M SCS-2200 (3M Brand ECP-2200) (see Fig. 35). The optical differences between the original and its replacement are discussed by Willey et al.⁴² Illinois Institute of Technology now produces the ECP-2200 under the name MH2200 (see following).

Nextel 2010 A sprayable Nextel 2010 Velvet paint is available from the German company Mankiewicz Gebr. & Co., Hamburg, and has been used on several ground-based astronomical instruments, where cryogenic cycling durability is important. The 2010 surface has an outgassing rate of 1.3×10^{-4} g/cm² at 84°F (29°C) in a 20-hour test at 10^{-5} to 10^{-6} torr. The solar alpha of this surface (from the company literature) is about 0.97 and the normal emittance is about 0.95.

The total emissivity of different black surfaces was investigated with regard to their being used as a total radiation standard in the temperature range -60 to +180°C by Lohrengel.¹³⁹ The investigations showed that the matte varnish Nextel Velvet Coating 2010 black with an emissivity of 0.951 was

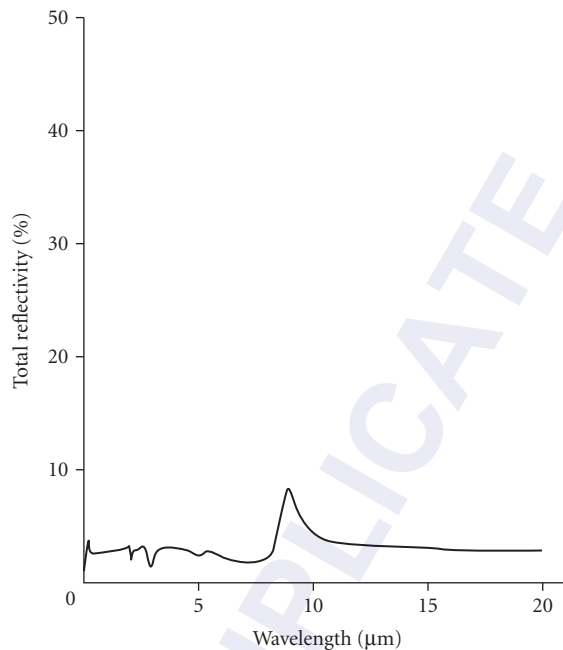


FIGURE 33 Total hemispherical reflectivity of 3M Nextel Black Velvet. (From Willey *et al.*⁴²)

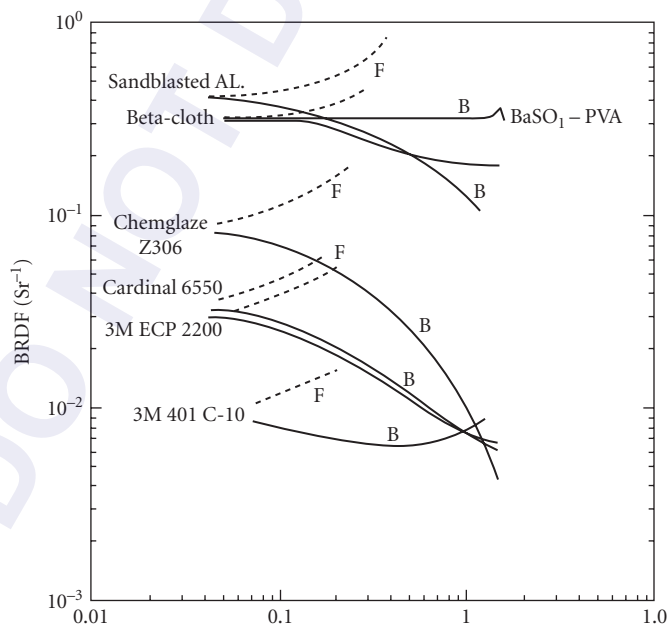


FIGURE 34 Comparison of BRDF profiles at an angle of 45° and 633-nm wavelength. (From Viehmann and Predmore.³⁵)

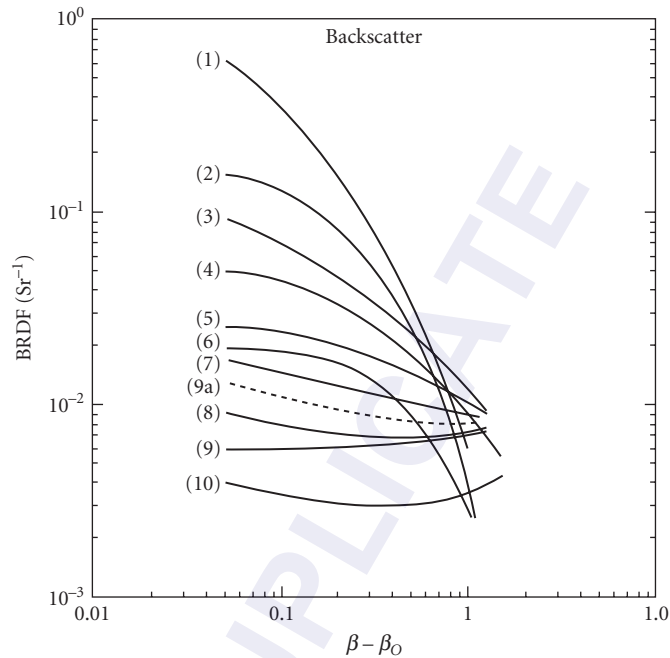


FIGURE 35 Comparison of BRDF profiles at an angle of 45° and 254-nm wavelength: (1) Reynolds Tru Color Diffuse; (2) Chemglaze Z306; (3) Sikkens 443; (4) Sanodal Fast Black GL; (5) Chemglaze Z306 plus microballoons; (6) Sanodal deep black MLW; (7) 3M 401-C10; (8) Chemglaze Z313; (9) DeSoto Flat; (10) DeSoto Flat after EUVSH; (11) Martin Black anodize. (From Viehmann and Predmore.³⁵)

best suited for these purposes. Additional changes to the surface (sprinkling of Cu globules, etching) could be made to increase the total emissivity by a further 0.007 or 0.016, respectively. The present availability of this paint from U.S. manufacturers is unknown.

Nextel Brand Suede Coating Series 3101-C10 C-10 is the black variety of this paint, available from Red Spot (see references). It is a low-gloss paint, with reasonable abrasion resistance (see Fig. 36). It has a soft-cushioned feel and suedelike appearance and was used on office equipment, furniture, and electronic equipment housings. The manufacturer recommends it for interior use only.

MH 2200 MH2200 (formerly ECP-2200, which was derived from the original 3M Nextel Black Velvet) is a one-part, flat black, nonselective solar absorber coating designed for high-temperature service (Fig. 37). When applied to aluminum substrates, the coating integrity and its "solar absorption" of 0.96 are unaffected by temperatures to about 500°C . The coating is 43 percent solids in xylene and is intended for nonwear applications. ECP-2200 is available from IIT Research Institute. It is touted as the substitute for 3M Nextel Black, but does not have the low BRDF at visible wavelengths of its former relative. It has a normal emittance of 0.95 at $10.6\ \mu\text{m}$ and an absorptance of 0.96 over the wavelength range of 0.35 to $2.15\ \mu\text{m}$. Its emittance is 0.86 on an aluminum substrate (normal emittance at 75°F from 2 to $25\ \mu\text{m}$).

Aeroglaze Z Series Aeroglaze (formerly Chemglaze) Z306 diffuse black paints are part of the Z line of single-package, moisture-curing, ASTM Type II, oil-free polyurethanes.²⁴ This paint has excellent chemical, solvent, and (if used with a primer) salt-spray resistance. This is a popular paint for aerospace

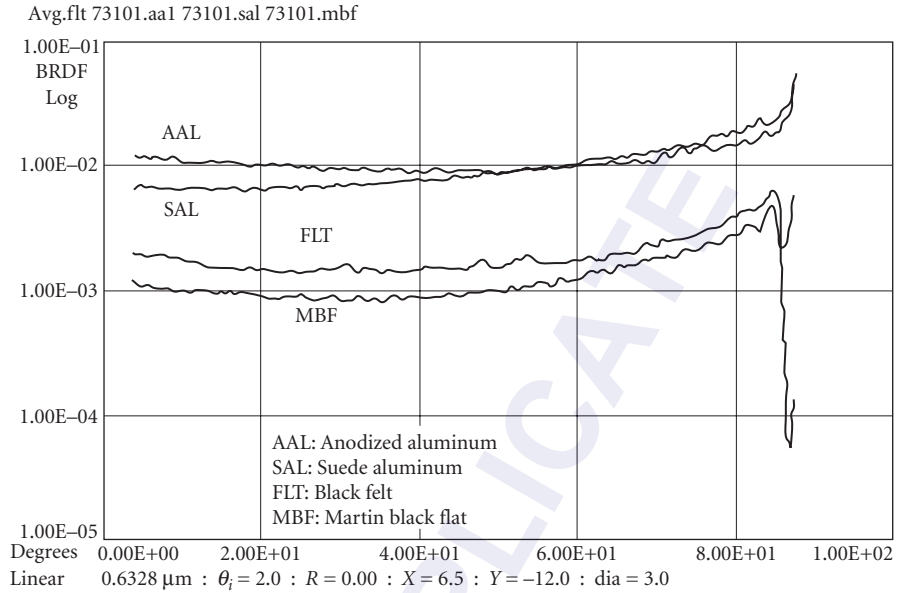


FIGURE 36 BRDF data for anodized aluminum, 3M Nextel Suede paint, black felt, and an aluminum flat with Martin Black wavelength is 633 nm, at near normal incidence. (From Cady *et al.*⁶⁶)

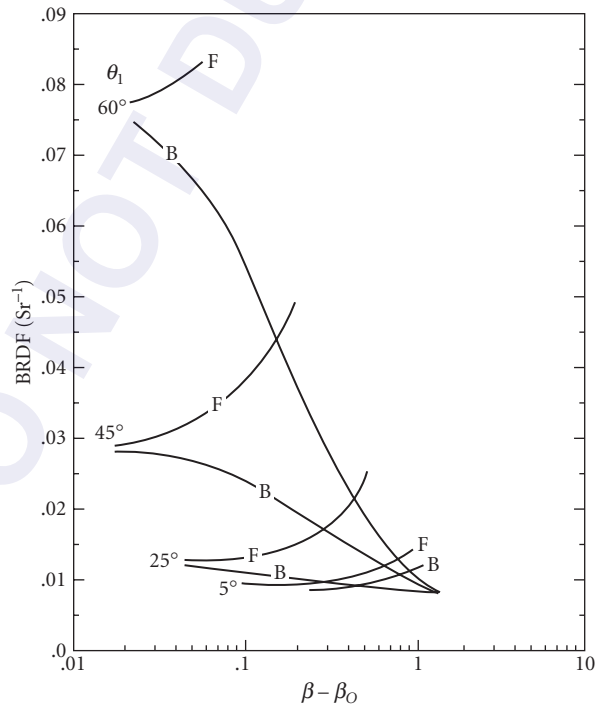


FIGURE 37 BRDF profiles of 3M ECP-2200 black paint 6550 nm. (From Viehmann and Predmore.³⁵)

use, especially after the discontinuation of 3M's 401-C10 Black Velvet Nextel.²⁸ For space use, the Aeroglaze 9924 and Aeroglaze 9929 primers are recommended. The optical properties are given in Figs. 34 and 38 to 43.

A low volatile version of the paint is being developed. In a version developed at Goddard Space Flight Center (based on the Z306)⁴⁴ glass microballoons are added to reduce the specular component, though the particles may be a source of particulate contamination. Other versions were developed at Goddard based on the Z004 paint. Goddard Black is a generic term for several paints of composition and optical properties different from Z306. This name is often incorrectly applied to the modified Z306 described here.

Akzo Nobel Paints Three paints sold formerly under the label Cat-a-lac are now known as Akzo 463-3-8 (formerly Cat-a-lac 463-3-8 diffuse black), 443-3-8 (formerly Cat-a-lac 443-3-8) (see Figs. 34 and 43 to 46) and 443-3-17 (formerly called Sikkens 443-3-17). The Bostic name is also formerly associated with these same paints. The 463-3-8 is the most diffuse of the three blacks, and the other two are considered "glossy." It is an epoxy-based system with low outgassing and low reflectance.

Cardinal Black From Cardinal Industrial Finishes. Optical data are given in Fig. 47.

Cat-a-lac Black (Former Name) This is an epoxy system with low outgassing and low reflectance. It is now manufactured by Sikkens Aerospace Finishes (see Akzo in the tables for address) and is described here.

DeSoto Black The PRC DeSoto Black is a flat black paint most suitable for the range 0.2 to 2.5 μm . It has a reflectance of between 2 and 3 percent over this range. Its emissivity is 0.960 compared to a 300 K blackbody, while its absorptivity over the solar spectrum is 0.924. DeSoto Black can meet outgassing requirements for space applications when vacuum baked. The weight loss is typically below 0.5 percent and the volatile condensable material is below 0.5 percent, similar to 3M Black Velvet.¹⁴⁰ A BRDF curve is given in Fig. 35.

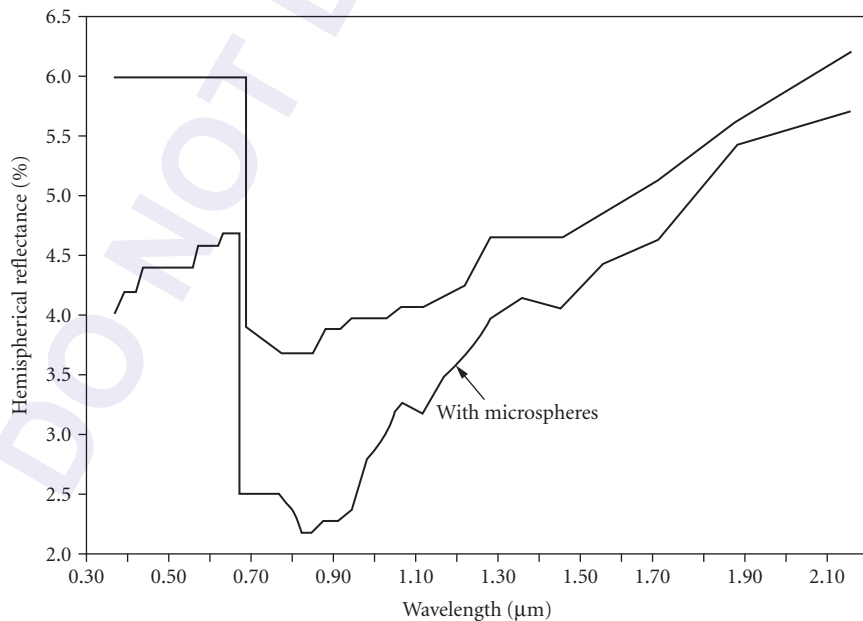


FIGURE 38 Diffuse visible and near-infrared hemispherical reflectance of Chemglaze (now called Aeroglaze) Z306, and Z306 with microspheres. (From Ames.²⁸)

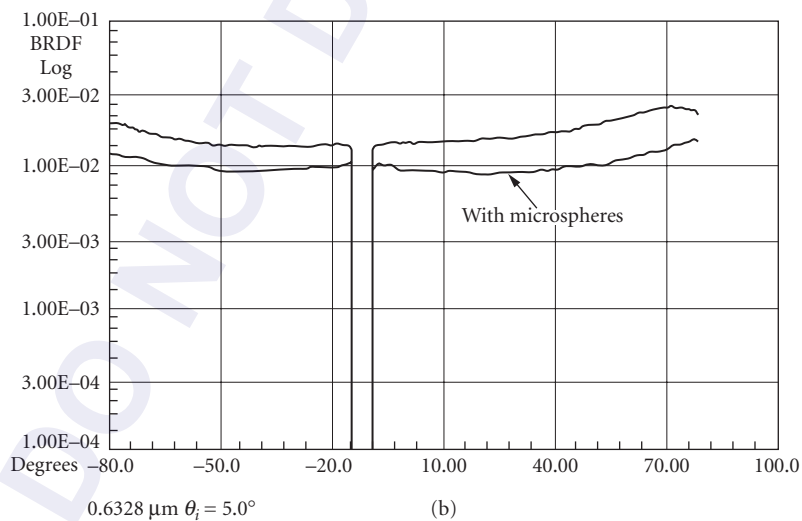
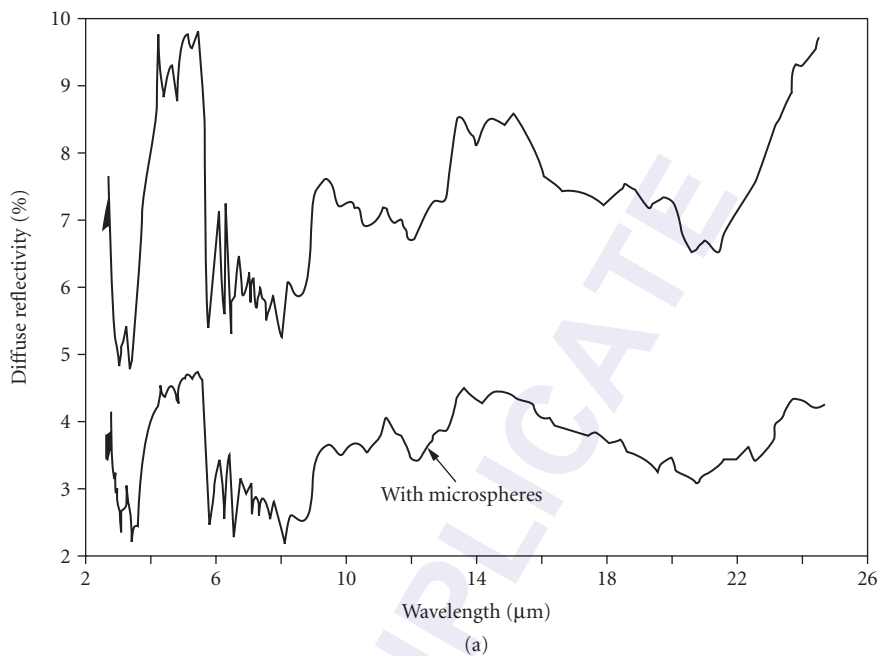


FIGURE 39 (a) Infrared diffuse reflectivity of Chemglaze Z306 (now called Aeroglaze) and Z306 with microspheres. (b) BRDF for Chemglaze (now Aeroglaze) Z306 at near-normal incidence (5°) at a wavelength of $0.6328 \mu\text{m}$. (From Ames.²⁸)

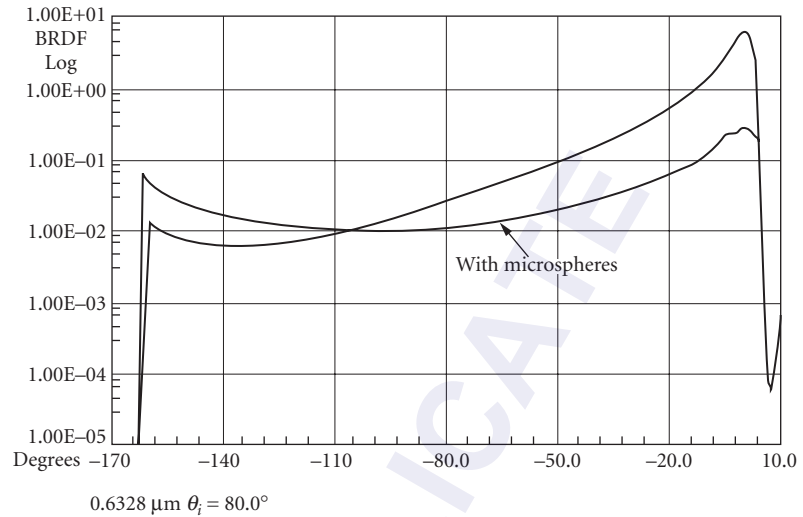


FIGURE 40 BRDF for Chemglaze (now Aeroglaze) Z306 at near grazing incidence (80°) at a wavelength of $0.6328 \mu\text{m}$. (From Ames.²⁸)

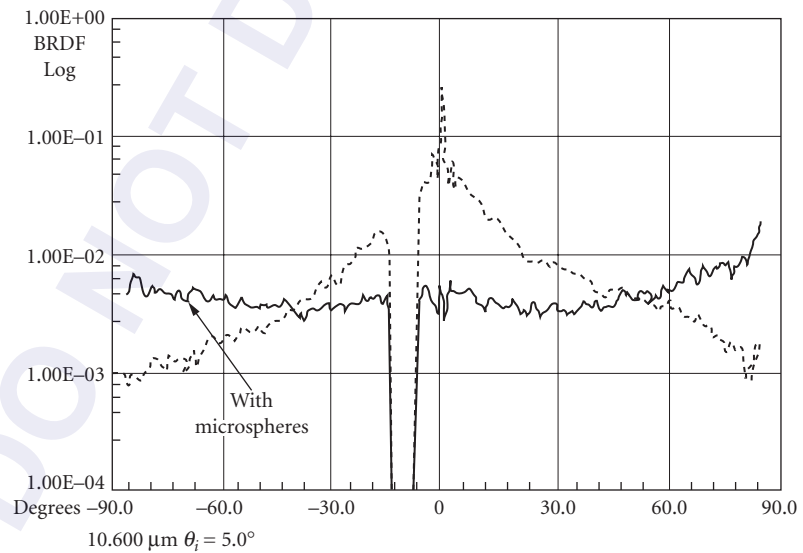


FIGURE 41 BRDF for Chemglaze (now Aeroglaze) Z306 at near-normal incidence (5°) at a wavelength of $10.60 \mu\text{m}$. (From Ames.²⁸)

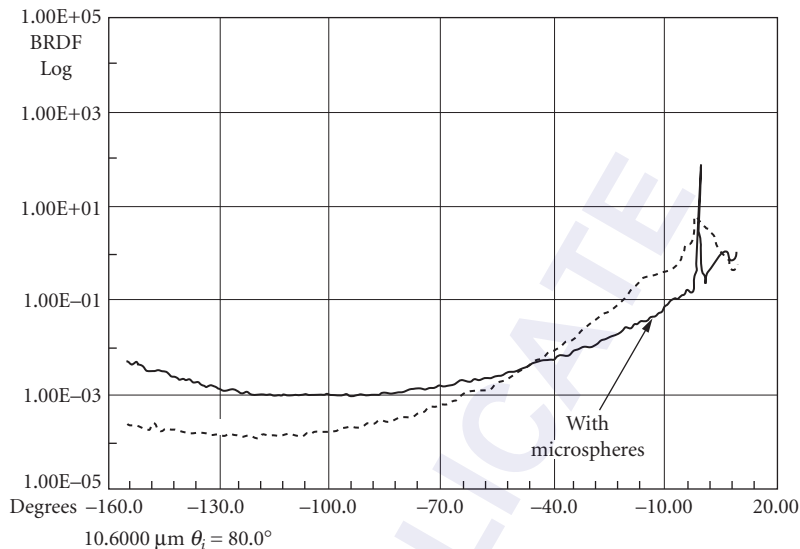


FIGURE 42 BRDF for Chemglaze (now Aeroglaze) Z306 at near grazing incidence (80°) at a wavelength of $10.60 \mu\text{m}$. (From Ames.²⁸)

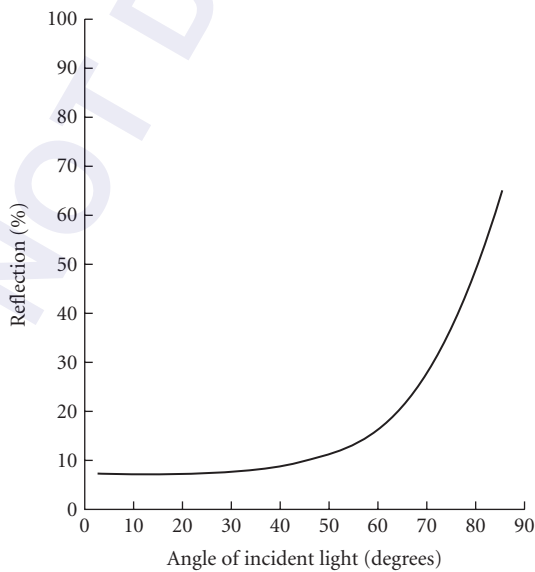


FIGURE 43 Specular reflection measured for Cat-a-lac glossy (now called Akzo 443-3-8) and for Chemglaze glossy (now called Aeroglaze Z302). The wavelength is $0.6328 \mu\text{m}$. The two lines overlap. (From Griner.²⁵)

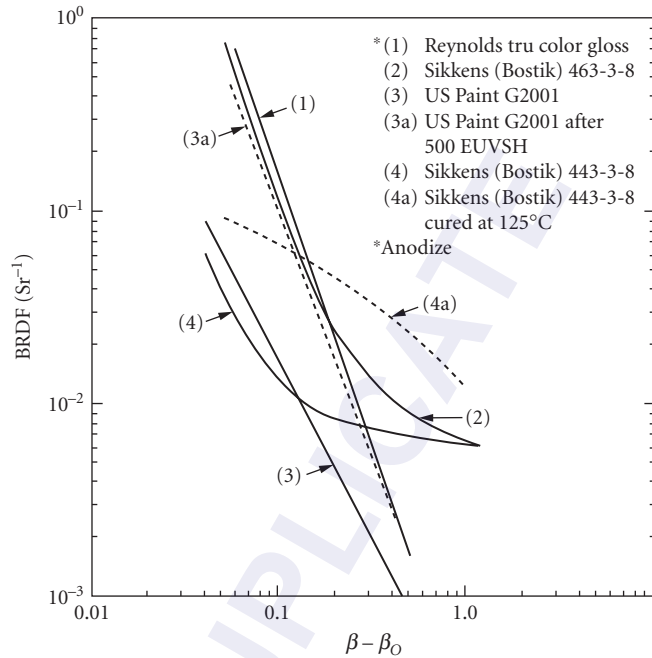


FIGURE 44 BRDFs of glossy surfaces at 0.254 μm and 45° angle of incidence. (From Viehmann and Predmore.³⁵)

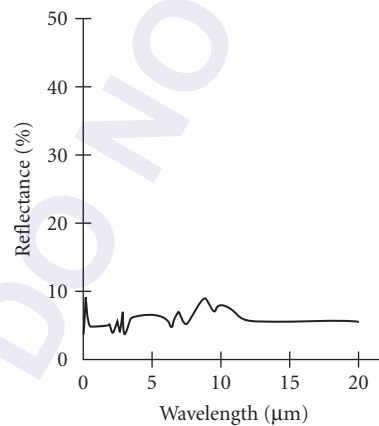


FIGURE 45 Total hemispherical reflectivity of Bostic 463-3-8 (now called Akzo 463-3-8) on bare aluminum. (From Willey et al.⁴²)

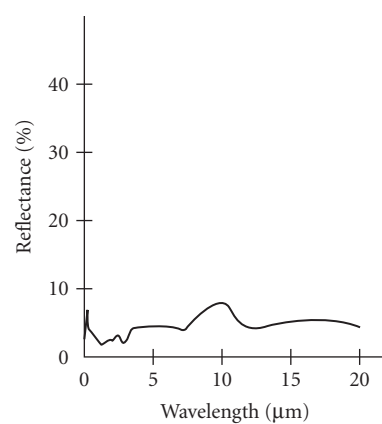


FIGURE 46 Total hemispherical reflectivity of Bostic 463-3-8 (now called Akzo 463-3-8) on ZnCr primer. (From Willey et al.⁴²)

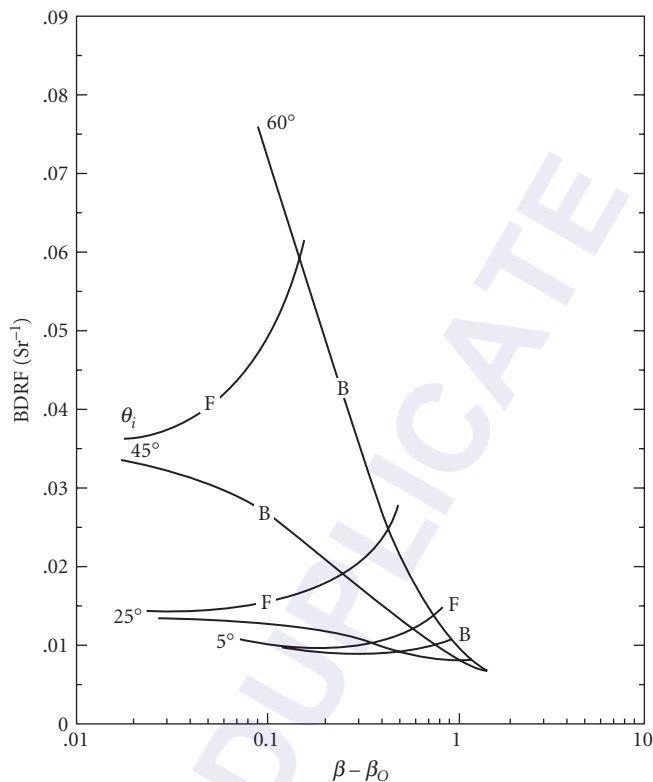


FIGURE 47 BRDF profiles of Cardinal black paint 6550 at 633 nm. (From Viehmann and Predmore.³⁵)

Floquil: Black for Lens Edges Lewis et al.¹⁴¹ discuss experiments to find low-scatter blackening components for refractive elements. With index of refraction matching, gains of several magnitudes in lower BRDF can be achieved over other edge-blackening techniques. The paint Floquil (Polly S Corporation) has been used as an edge-blackening compound and is discussed by Lewis. Smith and Howitt¹⁹ give a specular reflectance curve for Floquil in the infrared from 20 to 200 μm .

Parson's Black Prior to 1970, Parson's Black was a standard black reference.

SolarChem SolarChem is a black paint with excellent visible absorbing properties and has been suggested as a replacement for the unavailable 3M Nextel Black. The visible BRDF of SolarChem has a flat profile at near-normal angles of incidence (Fig. 48), while the 10.6- μm BRDF of SolarChem shows quite a different shape.⁸⁸ The surface is highly specular at 10.6 μm .

Anodized Processes

Anodized processes can produce fine black surfaces. The condition of the substrate plays an important role in the overall scattering properties of the surfaces. A common technique is to sandblast (vapor hone) the surface before anodization. Many varieties of grit may be used and each will produce a different kind of surface treatment. Sometimes, small pieces of grit can become embedded in the surface and this may affect the anodizing process. The BRDF of sandblasted aluminum surfaces are given in Figs. 49 and 50. Table 3 gives reference data on anodized surfaces.

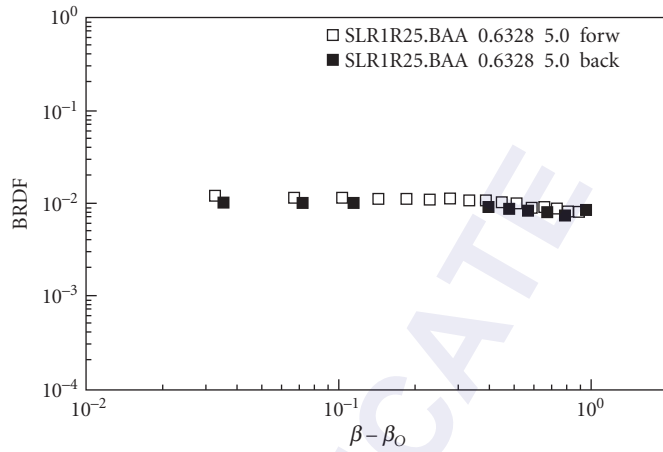


FIGURE 48 Visible (0.6328 μm) BRDF of SolarChem at 5° angle of incidence. (From Pompea et al.⁴)

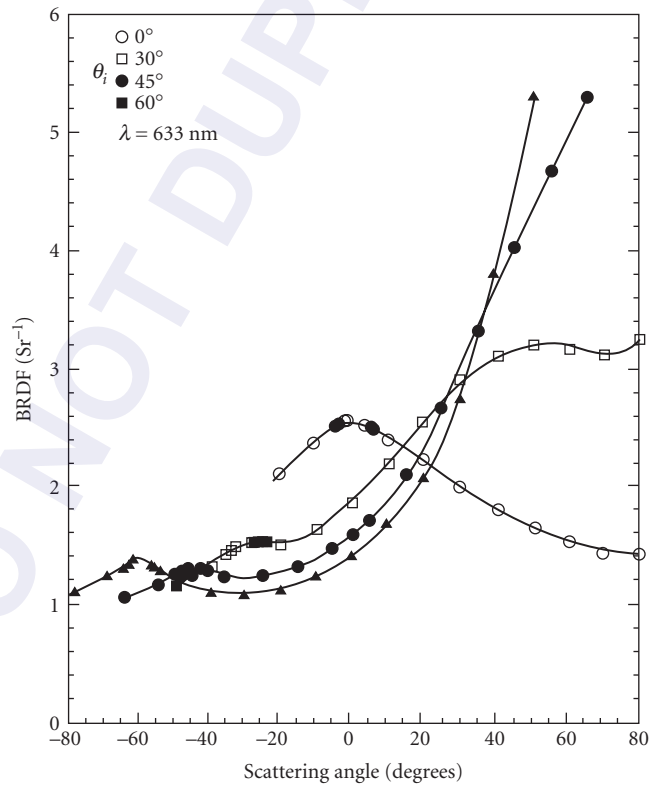


FIGURE 49 BRDF profiles of sandblasted aluminum. Open circle: 0° open square: 30° filled circle: 45° filled square: 60° ; wavelength of 633 nm. (From Viehmann and Predmore.³⁵)

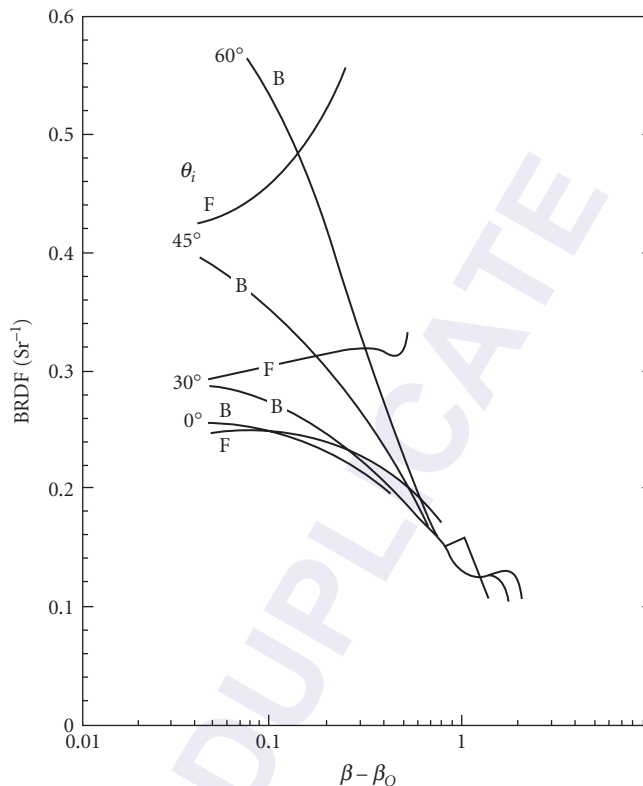


FIGURE 50 Forward (F) and back (B) scatter profiles of sandblasted aluminum. Wavelength of 633 nm. (From Viehmann and Predmore.³⁵)

Martin Black Martin Black is an anodized aluminum surface that is made microrough by a special anodization process developed by Lockheed Martin, Denver.^{142,143} It is made black from the inclusion of an aniline dye that is sealed into the surface. It was developed for the Skylab program and has been used on a wide variety of space instruments operating at vacuum ultraviolet to far-infrared wavelengths. The surface is rough, and scattering at several fundamentally different scale lengths occurs.

Special care in handling surfaces of this class is needed since the surface morphology consists of pyramids extending from the surface, and these features are easily crushed. The surface is still quite black, however, even if some surface damage occurs. The surface morphology works well for space applications as it does not support electrical potentials higher than 200 V, is not affected by temperatures as high as 450°C, by hard vacuum, or by ultraviolet radiation. It also has very low outgassing rates and will not delaminate even under severe environmental conditions.¹⁴⁴ It passes NASA SP-R-0022A specifications for outgassing.

The surface is also highly resistant to chemical attack by atomic oxygen prevalent in low earth orbit.⁶⁴ BRDF measurements of Martin Black show its lambertian character at 0.6328 and 10.6 μm . The reflectance spectrum shows significant reflectances at about 2.5 and 5 μm ; at other wavelengths, the reflectance is low. The optical data are given in Figs. 36, 51, and 52.

Enhanced Martin Black Enhanced Martin Black is similar in its properties to Martin Black, but was created to provide an even more durable surface for long exposure to atomic oxygen in low earth orbit. Experiments aboard the space shuttle confirm this superiority in the shuttle orbit environment.⁶⁴ The

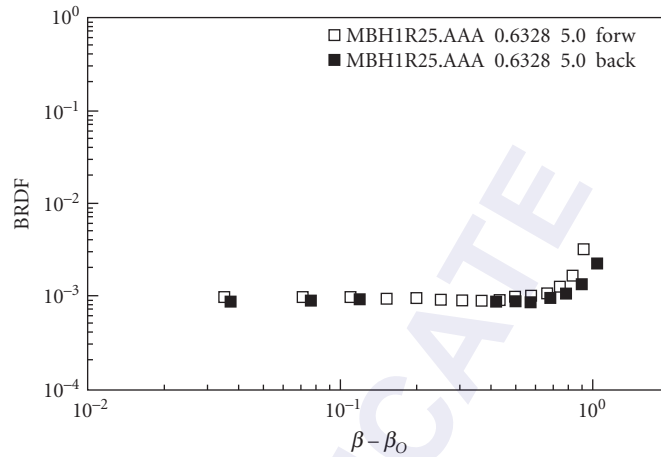


FIGURE 51 Visible ($0.6328 \mu\text{m}$) BRDF of Martin Black at 5° angle of incidence. (From Pompea et al.⁴)

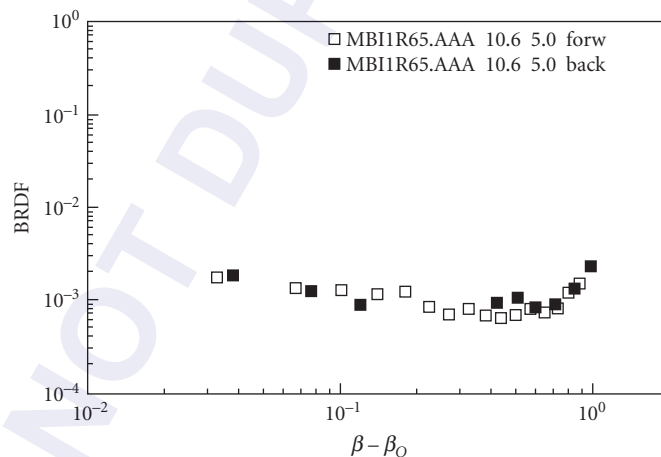


FIGURE 52 Infrared ($10.6 \mu\text{m}$) BRDF of Martin Black at 5° angle of incidence. (From Pompea et al.⁴)

process refinements also reduced the near- and middle-infrared reflectances from about 40 percent at $2.3 \mu\text{m}$ to about 25 percent and from about 15 percent at $5.5 \mu\text{m}$ to about 3.5 percent, while maintaining an absorption of about 99.6 percent in the visible region.

Posttreated Martin Black Another variation of the Martin Black surface is designed for near-infrared applications. Posttreated Martin Black uses hydrogen fluoride to further reduce the near-infrared reflectance peak and eliminate the middle-infrared reflectance peak of Martin Black. The area under the $2.3\text{-}\mu\text{m}$ peak of Martin Black has been reduced by about two-thirds and the $5.8\text{-}\mu\text{m}$ peak of Martin Black has been eliminated.

Infrablack Infrablack is another anodized surface from Martin Marietta, which relies upon a very rough substrate to produce a diffusing and absorbing surface for wavelengths up to $750\ \mu\text{m}$.^{138,145} It was developed for the next generation of far-infrared NASA telescopes. It has also been used in experiments to measure the Stefan-Boltzmann constant where a “blackbody” surface with very high emissivity was needed. This type of surface can also be used in space radiators where high emissivity is beneficial. Like other anodized surfaces in this class, the Infrablack surface is fragile and care in handling is needed to avoid crushing the surface structure. As with Martin Black, even a mistreated or crushed surface will still be significantly more light absorbing than most other black surfaces. Infrablack can be tuned for maximum absorption between 100 and $500\ \mu\text{m}$. The BRDF in the visible and infrared for Infrablack is similar to Martin Black and Enhanced Martin Black; all exhibit relatively lambertian behavior in the visible at near-normal angles of incidence (Figs. 53 and 54).

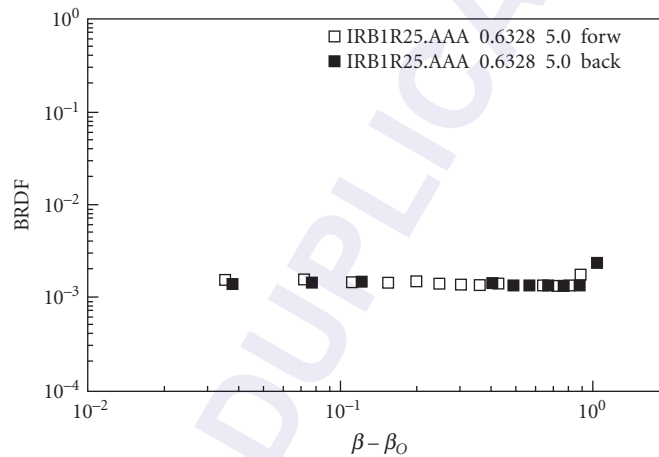


FIGURE 53 Visible ($0.6328\ \mu\text{m}$) BRDF of Infrablack at 5° angle of incidence. (From Pompea et al.⁴)

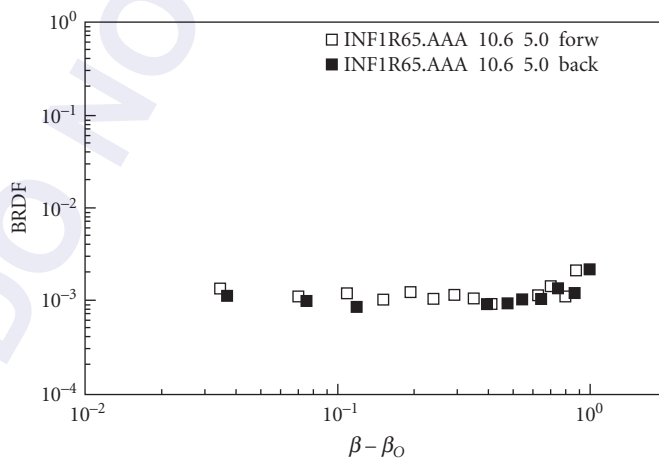


FIGURE 54 Infrared ($10.6\ \mu\text{m}$) BRDF of Infrablack at 5° angle of incidence. (From Pompea et al.⁴)

DEEP SPACE BLACK DEEP SPACE BLACK (N-Science/Advanced Surface Technologies, Arvada, Colorado) shares many if not most characteristics of Martin Optical Black, a similar specialized anodization process. Like Martin Black it is extremely absorptive over a wide range of wavelengths from ultraviolet to the infrared and is best used on 6061 aluminum. The surface has good mechanical stability. The surface can be fine tuned for better performance in the near-infrared (1 to 5 μm) and for maximal durability in the low-earth orbit space environment. The total hemispherical reflectance in the 0.4- to 0.7- μm range is 0.37 and 0.78 percent in the 0.4- to 1.0- μm range. Like Martin Black, the surface does not support an electrical potential greater than 200 V and is robust under vibration and cryogenic cooling. When cryogenically cooled the surface has an emissivity in the far-infrared region (7 to 120 μm) of 0.985 or better.

This surface is used in controlling stray light or as black surface with high emissivity in a variety of optical systems such as detector cold shields and assemblies, detector calibration sources and blackbody cavities, thermal control and radiator surfaces, and spacecraft applications such as star trackers and cameras. DEEP SPACE BLACK is available in three variations:

DEEP SPACE BLACK-VIS: DSB-VIS is optimized for performance in the visible wavelength regime.

DEEP SPACE BLACK-Full Spectrum Basic: DSB-FSB is optimized to provide best performance over the largest wavelength regime possible. **DEEP SPACE BLACK-Full Spectrum Enhanced:** DSB-FSE is a special variant of DSB-FSB that has a more robust finish with the same outstanding optical characteristics of the FSB. DSB-FSE provides a somewhat improved optical performance in the near- and middle-IR. The performance in the visible region of the spectrum is maintained. This product is more suited to spacecraft applications, including those applications concerned with LEO environments.

Tiodize V-E17 Surface The Tiodize process is an electrolytic-conversion hard coating of titanium using an all-alkaline, room-temperature bath. It produces an antigall coating with good optical properties. The process can be used on all forms of titanium and its alloys and has been used in a wide variety of space vehicles and aircraft. The Ultra V-E17 coating is a black organic coating which changes the absorptivity and emissivity of titanium (0.62 and 0.89, respectively) to 0.89 and 0.91, respectively. The emittances were determined from a 25-point integration between 4.8 and 26.2 μm . The absorptance was determined by a 19-point integration between 0.32 and 2.1 μm . The surface has a total mass loss (TML) of 0.91 percent and no detectable volatile condensable material (VCM) in tests run at Ford Aerospace.¹⁴⁶ Its interaction with atomic oxygen is unknown.

TRU-Color Diffuse Black An electrolytic coloring process from Reynolds Aluminum.

Hughes Airborne Optical Adjunct Coating A cleanable specular black for use at 1 to 10 μm .

Etching of Electroless Nickel

The etching of electroless nickel provides a black coating useful for many applications. Some reference data on several electroless nickel processes are presented in Table 4. Kodama et al.¹⁴⁷ developed optical absorbing black films of etched nickel-phosphorus alloy deposited on substrates using an electroless plating process. Results of varying the plating and etching bath components and conditions were studied. They produced an optical absorber with a spectral reflectance of about 0.1 to 0.2 percent in the wavelength range of 488 to 1550 nm. The absorber showed substantial immunity to degradation in two accelerated-aging environments: 2000 hours in a dry-heat environment at 100°C, and 2000 hours in a damp-heat environment at 85°C and 85 percent relative humidity.

NBS Black This is a blackened electroless nickel surface developed for use as a solar collector at the National Bureau of Standards (now NIST) in Gaithersberg, MD.^{75,76,148} It consists of an electroless nickel-phosphorous coating that can be plated onto a wide variety of substrates such as metals, glass, ceramics, and plastics. This coating is subsequently etched using nitric acid, creating conically shaped holes into the surface, which act as light traps.

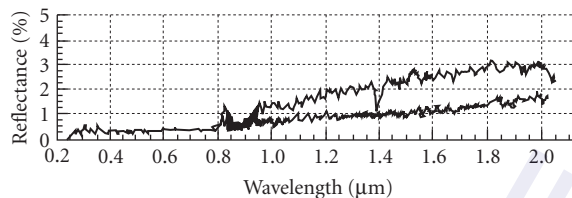


FIGURE 55 Total integrated scatter of NBS Black. The substrates are magnesium AZ31B-F (upper curve) and magnesium ZK60A-T5. Substrate is silica sandblasted. $\alpha = 0.978$. $\epsilon = 0.697$. (From Geikas.⁶⁵)

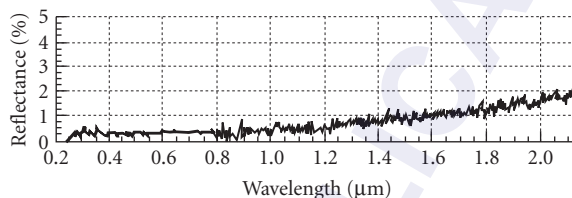


FIGURE 56 Diffuse scatter of NBS Black. The substrates are magnesium AZ31B-F and magnesium ZK60A-T5. The plots are overlaid and are virtually identical. Substrate is silica sandblasted. $\alpha = 0.978$. $\epsilon = 0.697$. (From Geikas.⁶⁵)

The specular reflectance is less than 1 percent over the range of 0.32 to 2.4 μm , making it an extremely black coating (Figs. 55 and 56). As this surface finish was developed for solar collectors, the surface has a high value for absorptivity (0.978), but a low value for room-temperature emissivity (about 0.5). The surface is moderately durable as the surface's relief consists of conical holes into the surface. The NBS Black process has been modified resulting in the next surface, which has greater infrared absorption.

Ball Black Ball Black is an optical black surface produced with modifications to the techniques just described—the selective etching on an electroless nickel surface to produce a multiplicity of conical light traps in the surface.^{29,65} The surface appears intensely black to the eye. The surface, however, is still an unaltered nickel-phosphorous deposit with the same chemical properties as the unetched surface, which was shiny. Ball Black can be plated on aluminum, beryllium, copper, stainless steel, invar, polycarbonate and ABS plastics, titanium, and some magnesium alloys. Since this process involves only deposition and etching of a metallic material, problems with outgassing or volatilization of the surface in a vacuum are minimal.

The surface is able to withstand rapid changes in temperature, with immersion into boiling water after being at liquid nitrogen temperature. The surface is sensitive to handling, though less so than surfaces like Martin Black. Small blemishes in the surface can be restored by re-etching without replating. The solar absorptivity can be made to vary between 0.71 and 0.995 and its room-temperature emissivity can be tailored to values between 0.35 and 0.94 by changes in the etching parameters. The directional reflectance increases by less than a factor of two from 1.5 to 12 μm . It has a BRDF of about 6×10^{-2} at normal incidence at 0.6328 μm . See Figs. 57 and 58 for optical data.

Plasma-Sprayed Surfaces

Boron Black Martin Marietta Boron Black is a surface that is black in the visible (solar absorptivity between 0.89 and 0.97, and emissivity of at least 0.86) and black in the infrared¹⁴⁹ (Fig. 59). Boron

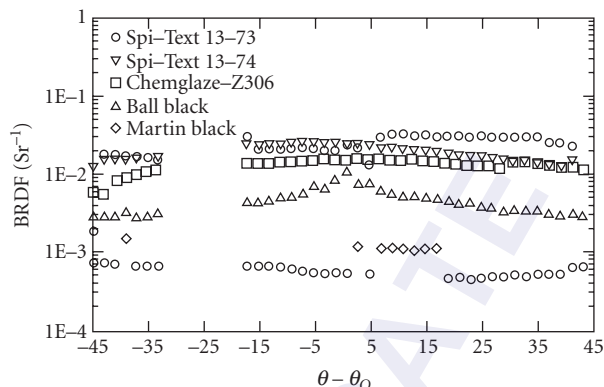


FIGURE 57 BRDF measurements at $0.6328 \mu\text{m}$ of textured metal surfaces from Spire Corporation and Chemglaze Z306 paint, Ball Black electroless nickel coating, and Martin Black anodized coating. (From Lompadó et al.²⁹)

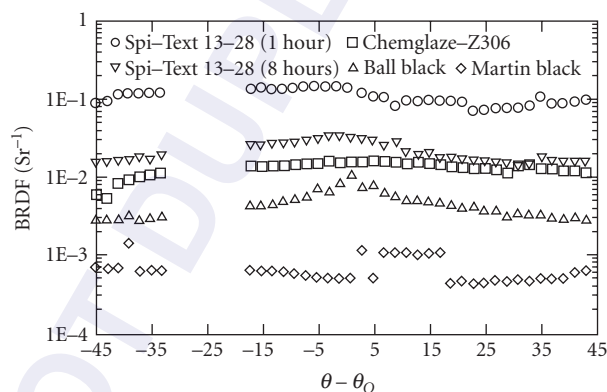


FIGURE 58 BRDF measurements at $0.6328 \mu\text{m}$ of two textured metal surfaces from Spire Corporation and Chemglaze Z306 paint, Ball Black electroless nickel coating, and Martin Black anodized coating. (From Lompadó et al.²⁹)

Black is a plasma-sprayed surface and can probably be applied to most metallic substrates. It has been applied to molybdenum, nickel, and titanium substrates. It offers a low-atomic-number surface with good optical properties and the ease of a plasma-spray deposition. The infrared BRDFs show a profile characteristic of a surface with lambertian profiles at near-normal angles of incidence.

Boron Carbide This is a Martin Marietta proprietary process. The boron carbide surface is applied to a titanium substrate. A 10- to 15- μm layer of B_4C is applied in a proprietary plasma-spraying process. BRDF measurements at 30° angle of incidence at two wavelengths illustrate that this surface is more lambertian in the visible.

Beryllium Surfaces Porous surfaces or those with steep-walled features which can trap radiation have been developed using plasma-sprayed beryllium and sputter-deposited beryllium by workers at Oak Ridge^{82,150,151} and at Spire Corp.¹⁵² who are developing baffles that can operate in severe environments.

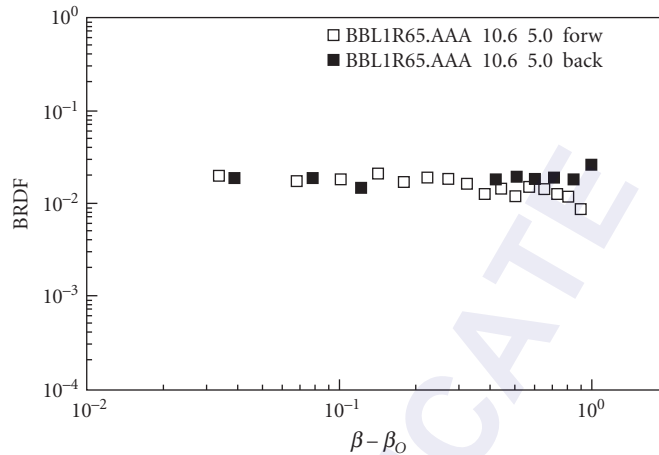


FIGURE 59 Infrared (10.6 μm) BRDF of Boron Black Black at 5° angle of incidence. (From Pompea et al.⁴)

Plasma-sprayed beryllium samples appear visually rough with a matte gray finish. The specular reflectance is approximately constant at about 1.6 percent from 2 to 50 μm , while the BRDF is lambertian in character at near-normal angles of incidence with a value of about $5 \times 10^{-2} \text{ sr}^{-1}$, equivalent to a total hemispherical reflectance of about 15 percent.

Beryllium surfaces with varying thicknesses and columnar grain sizes can also be created by low-temperature magnetron sputtering. The surfaces can then be chemically etched to enhance their absorptive properties. These surfaces exhibit less than 2 percent specular reflectance at near-infrared wavelengths. Very thick surfaces (350 μm) exhibit specular reflectances of less than 0.5 percent up to 50 μm wavelength.

A sputtered coating of Be on a Be surface can be made more absorptive through exposure to an oxygen plasma to form a layer of BeO on the surface of the coating.²⁹ Optical data are presented in Figs. 57, 58, 60, and 61.

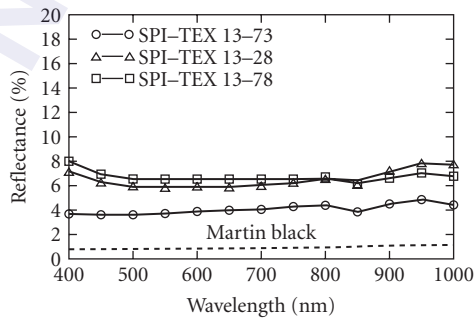


FIGURE 60 Total hemispherical reflectance in the visible bandpass for textured baffle materials from Spire Corp. The 13-73 surface is a textured aluminum surface. (From Spire.⁸¹)

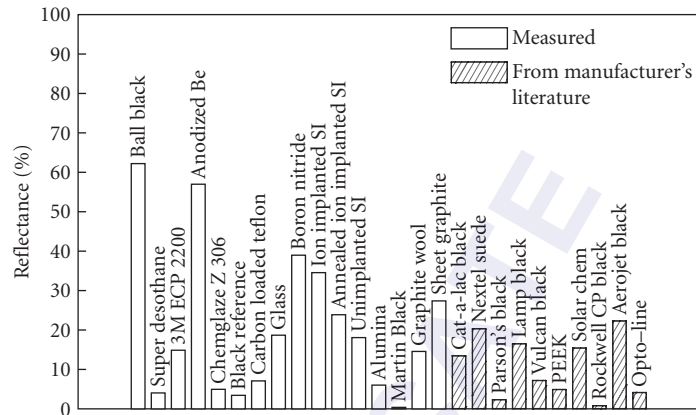


FIGURE 61 Total hemispherical reflectance data for baffle materials at $10.6\ \mu\text{m}$. The shaded measurements are from the manufacturer's literature and have not been independently verified. (From Johnson.¹⁵³)

Ion Beam-Sputtered Surfaces

Ion beam-sputtering processes can roughen surfaces of a large variety of compositions. An excellent review is by Banks in Cuomo et al.⁷⁷ (see also Ref. 78). A variety of sputtering processes can create rough features of various geometries that decrease specular reflections and create a black-appearing surface from a light-appearing substrate. Since the roughness is integral to the surface, these surfaces often have desirable environmental resistance to atomic oxygen. However, these surfaces also have tremendous surface area for trapping of contaminants and thus a potential for large outgassing rates if preventative steps are not taken. These surfaces are expensive to make in large areas, though they are becoming less expensive as coating facilities are built to harness this technology. They do have tremendous potential as diffuse baffle surfaces and as high-emittance surfaces for high-temperature space radiators.^{154,155} They are also cleanable and very durable.

Lompado et al.²⁹ and Blatchley et al.⁸³ describe several textured aluminum surfaces and give plots of their total hemispherical reflectance in the $0.4\text{--}1.0\text{-}\mu\text{m}$ range, as well as BRDF plots at 0.6328 and $10.6\ \mu\text{m}$ (Fig. 60).

Electrodeposited Surfaces

Refer to Table 4 for further reference data on these surfaces.

Black Cobalt Black cobalt is a patented surface developed at Martin Marietta¹⁵⁶ which appears black in the visible, with a solar absorptivity of at least 0.96 and emissivity of at least 0.6. It is applied by an electrodeposition process to a substrate of nickel, molybdenum, or titanium; other substrates are also expected to work. It can be used at higher temperatures than anodized aluminum surfaces such as Martin Black. It is stable against loss of absorptivity at temperatures greater than 450°C for prolonged periods. The surface is not a lambertian surface at $10.6\ \mu\text{m}$ (Fig. 62).

Black Chrome A black chrome surface developed for aerospace use¹⁵⁷ is black in the visible and is created by an electroplating process on a conducting substrate. Many metallic substrates can be used and samples on molybdenum, nickel, and titanium have been made. The preferred substrates are titanium and molybdenum. The coating process involves electrodeposition using chromium and chromium oxides. The visible BRDF is flat (Fig. 63), while the infrared profile shows the surface to

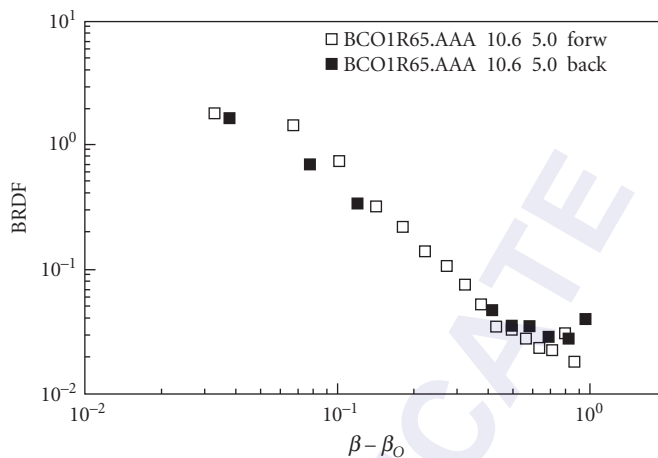


FIGURE 62 Infrared (10.6 μm) BRDF of Black Cobalt at 5° angle of incidence. (From Pompea *et al.*⁴)

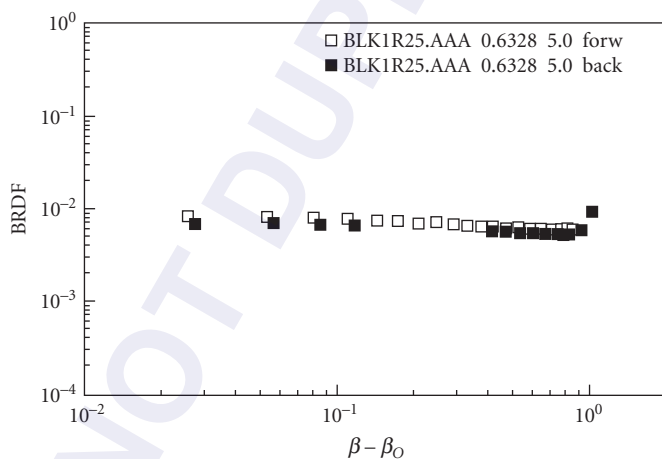


FIGURE 63 Visible (0.6328 μm) BRDF of Black Chrome at 5° angle of incidence. (From Pompea *et al.*⁴)

be largely specular in its scattering profile. The solar absorptivity is at least 0.95 and the emissivity can be tailored in the range of 0.4 to 0.8. The surface has excellent adhesion to the substrate.

Orlando Black Optical Coating A dendritic surface is produced by the electrodeposition of copper in a proprietary copper-plating formulation (Janeczko¹⁵⁸). The dendritic structure is then oxidized to form smaller structures superimposed on the larger ones. The surface has excellent broadband absorption properties from 0.4 to 14 μm and good absorptance at nonnormal angles of incidence (Fig. 64). It may be suitable for use at shorter wavelengths. The surface is usable in vacuum environments as well as a thermal reference in FLIR systems. A wide variety of materials can be coated including ABS plastic, nylon, Ultem, Noryl, fiberglass-reinforced epoxy, glass, stainless steel, copper, brass, nickel, aluminum, gold, and beryllium oxide. Dimensional allowances must be made for the thickness of the finish (approximately 0.001 in).

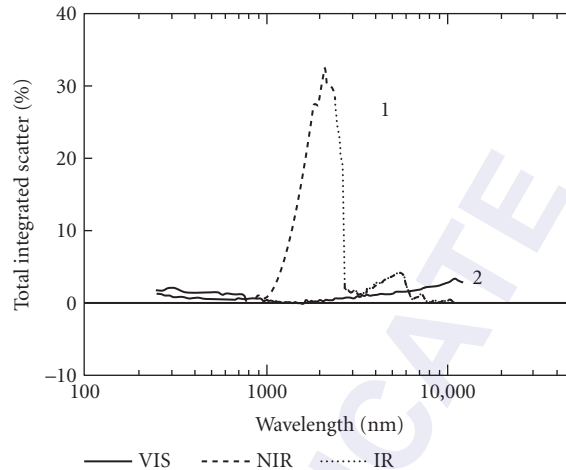


FIGURE 64 Total integrated scatter measurement of (1) and (2) Orlando Black commercially available anodized aluminum. (From Janeczko.¹⁵⁸)

Other Specialty Surfaces

Refer to Table 4 for further data on some of these surfaces.

Acktar Black Coatings (Acktar Ltd., Kiryat-Gat, Israel) Aktar makes a number of black coatings including Nano Black, Magic Black, Vacuum Black, Fractal Black, Ultra Black, and Metal Velvet. These inorganic coatings each have a wide range of working temperatures, very low outgassing (CVCM) in the 0.001 percent range and RML in the 0.167 percent range. Additionally the coatings have excellent adhesion to aluminum, copper, stainless steel, invar, kovar, glass, ceramics, polyimide, nickel, magnesium, titanium and other metals. There are options to make the coatings electrically conductive or nonconductive. The company also makes high emissivity foils.

Carbon Nanotubes and Nanostructured Materials Carbon nanotube arrays can also be used to create ultrablack surfaces. These surfaces are not yet practical for many applications but may have utility in specialized areas. For example, Yang et al.¹⁵⁹ have demonstrated that low-density vertically aligned carbon nanotube arrays can be engineered to have an extremely low index of refraction. The nanoscale surface roughness of the arrays contributes to an ultralow diffuse reflectance of 1×10^{-7} ; an order-of-magnitude lower compared to commercial low-reflectance standard carbon. The nanotube arrays can have an integrated total reflectance of 0.045 percent making it among the darkest materials ever fabricated.

Vorobyev and Guo¹⁶⁰ and Paivasaari et al.¹⁶¹ have also showed that laser ablation and altering of metallic surfaces can lead to greatly enhance optical absorptance.

Nanostructuring alone can enhance the absorptance by a factor of about three or more. The physical mechanism of the increased absorption is due to several effects including nanostructural, microstructural, and macrostructural surface modifications. The value of femtosecond laser ablation techniques using high fluences with a large number of applied pulses in reducing surface reflectivity may be applicable for a wide variety of metallic surfaces. Currently the process is time consuming for large area surfaces.

DURACON (American IMEX Corporation, Monroe, CT) *DURACON* is another high emissivity (>0.98 over the 0.55- to 1.8- μm range; >0.95 over the 2- to 20- μm wavelength range) durable black coating used in applications that cover a variety of wavelengths, especially the middle-infrared. It has

good thermal diffusivity and the interesting property that its conductivity can be tailored from low resistivity (less than $700 \Omega \text{ cm}$) to nearly an insulating surface. Thus, it can be used in environments where electrostatic discharge may be damaging to dielectric and insulating surfaces. It has good adherence and can be used on ceramic, glass, composite, plastic, and metal surfaces. It is durable to handling, unlike some anodized surfaces with surface structures that can be damaged by abrasion or direct compression. It also has very good ability to withstand wind shear at high Mach numbers and to withstand thermal cycling to 1000°F .

The surface can be applied by brush on, spray, or dip coating and any damaged areas can be repaired using a brushed on application of a fresh layer. The product is reported to be ready for normal regular use in 2 hours after this repair.

Electrically Conductive Black Paint Birnbaum et al.⁵¹ developed an electrically conductive, flat black paint for space use in places such as Jupiter's radiation belts where spacecraft charging effects can be important. Its small resistivity prevents the buildup of charge on the surface. It has a visual reflectance of less than 5 percent and other desirable optical properties.

Epner laser black (Brooklyn, New York) This makes a well-known coating used in a variety of application. Laser Black is a multilayer metallic oxide that has a microdendritic structure susceptible to crushing. This inorganic coating does not outgas and when vibrated ultrasonically in acetone shows no flaking or chipping. It has good durability at cryogenic temperatures and can be applied to any metallic substrate and most plastics. It requires a deposition of copper approximately 10 to $12 \mu\text{m}$ thick. Laser Black is used in cold shields for micro-bolometers and in baffles.

Ebanol C This specialty cupric oxide coating can be applied per MIL-C-14550 undercoat and MIL-F-495 processes: immersion in a zincate bath followed by a high cyanide, low copper bath, and copper strike and dip in an Ebanol C solution.⁹

High-Resistivity Coatings Strimer et al.¹⁶² describe a number of black surfaces with high electrical resistivity for visible and infrared applications.

Sputtered and CVD Surfaces Carbon, quartz, and silicon surfaces have been modified by sputtering or by chemical vapor deposition¹⁶³⁻¹⁶⁵ to create surfaces that are black over the wavelength range of 1 to $15 \mu\text{m}$. The textured surfaces are produced by sputtering with a low-energy (e.g., 500 eV) broadbeam ion source while adding impurities to the surface. The results from sputtering with this seeding process are structures in the form of cones, pyramids, and ridges. The exact nature of the surfaces created depends on the substrate temperature, ion and impurity fluxes, and the impurity species. The reflectivity of one modified silicon surface was below 1 percent throughout 1- to $25\text{-}\mu\text{m}$ wavelength range.

Silicon Carbide A silicon carbide surface for high-temperature applications was developed at Martin Marietta for application to a carbon-carbon substrate. It is applied to a hot substrate by a chemical vapor deposition process. The visible BRDF is flat, while the infrared profiles show the surface's specular nature at the longer wavelength.⁴

IBM Black (Tungsten Hexafluoride) This is a surface being produced at Martin Marietta, Orlando, using a process licensed from IBM. Any material that can tolerate 400°C can be used as a substrate for this process which involves vapor deposition of tungsten. The surface has dendritic structures that form light-absorbing traps. The surface looks like a collection of obelisks, does not outgas, and is rugged to the touch. Optical properties and an SEM photograph of the surface are given by Willey et al.⁴² The surface is very black in the 1- to $2\text{-}\mu\text{m}$ wavelength range.

ZO-MOD BLACK ZO-MOD BLACK (ZYP Coatings, Oak Ridge, Tenn.) is a high-emissivity coating applied like paint to ceramic porous and fibrous structures. On heating, a hard abrasion- and chemical-resistant, calcia-stabilized zirconium oxide coating is formed. The coating has high emissivity. The coating has been used inside of ceramic furnaces.

Gold Black Gold blacks are fragile surfaces made by evaporating gold in a low-pressure atmosphere of helium or nitrogen. The surface has good visible light absorption. The preparation and optical properties of gold blacks are given by Harris and McGinnies,¹⁶⁶ Blevin and Brown,¹⁶⁷ and Zaeschmar and Nedoluha.¹⁶⁸

Flame-Sprayed Aluminum Flame-sprayed aluminum is a very durable and very rough surface that has uses as an infrared diffuse reflectance standard. The parameters of the flame-spray process must be tightly controlled to create a similar surface each time.

Black Glass Black glass is available from Schott Glass Technologies, Duryea, Pa., in a glass called UG1, which is an ionically colored glass. Star Instruments, Flagstaff, Ariz., produces a glossy black glass called LOX8. It is a low-expansion copper glass with a coefficient of expansion of about 30 times lower than Pyrex.

Black Kapton The optical properties of Black Kapton are given in Fig. 65.

Other Surfaces Some other data is given for specular metallic anodized baffle surfaces (Fig. 66). This kind of basic treatment is adequate for some applications. The BRDF of Beta cloth is given in Fig. 67. Figure 68 shows the potential of various structured carbon surfaces to be highly absorbing over a wide range of angles. These kind of surfaces have great potential for the future.

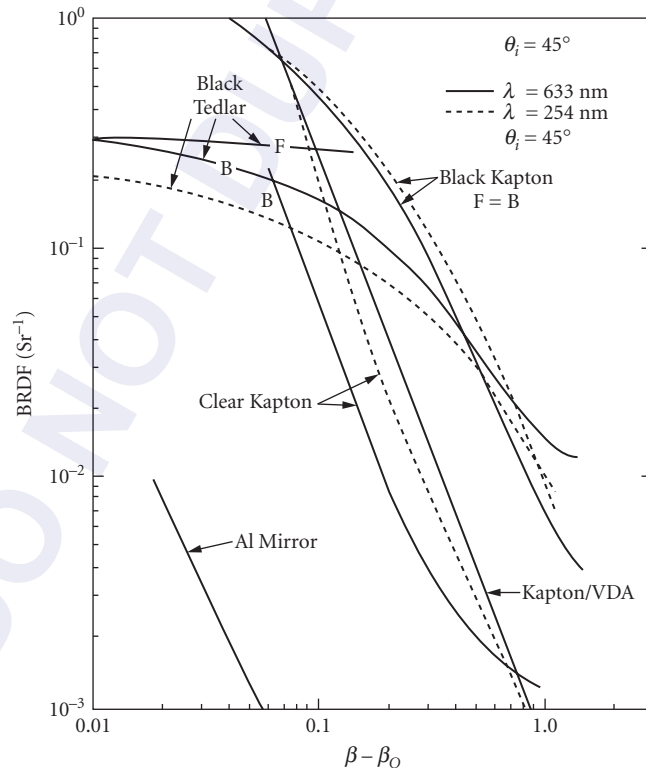


FIGURE 65 BRDFs of black Kapton and black Tedlar films. (From Viehmann and Predmore.³⁵)

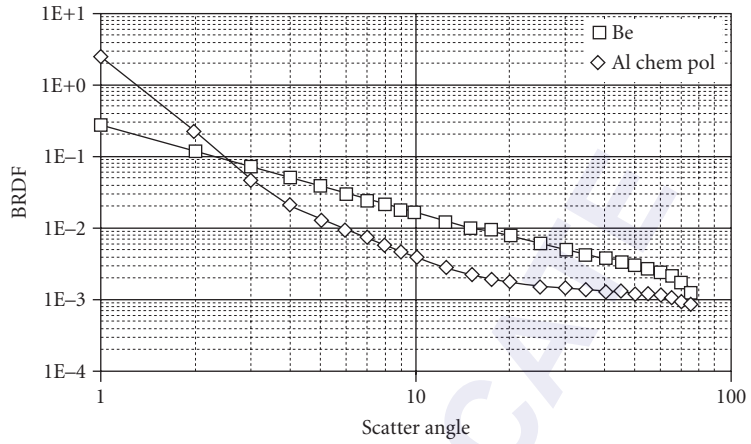


FIGURE 66 BRDF of specular baffle materials using Be and Al substrates that are polished and anodized. The measurements were made at $0.5145 \mu\text{m}$. (From Schaub et al.²⁰)

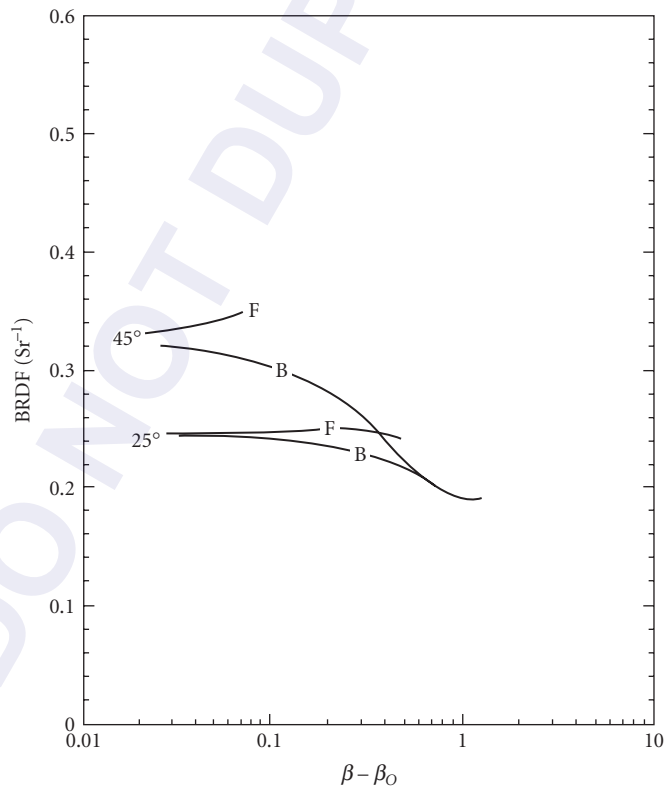


FIGURE 67 BRDF profiles of beta-cloth. Wavelength of 633 nm . (From Viehmann and Predmore.³⁵)

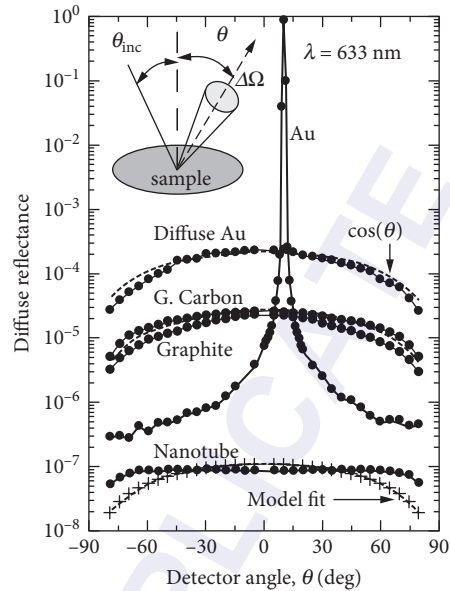


FIGURE 68 Diffuse reflectance of carbon nanotube samples compared to other carbon-based materials. (From Yang et al.¹⁵⁹)

6.9 CONCLUSIONS

Black surfaces are used as optical elements in a variety of ways. There is a great choice of surfaces available, but the selection of the appropriate surface is often problematical. The consequences of choosing the wrong surface are often quite severe.

For the selection process to be most effective, it must be done as early as possible. While the optical data on surfaces is relatively extensive, the data are often inconsistent or are not available for the particular wavelength of interest. The creation of large databases of optical and material properties is a great aid to the optical designer and materials consultant. Of great concern for space materials is the lack of very long term environmental test data. The synergistic effect of a number of degrading influences (e.g., solar ultraviolet, charging effects, atomic oxygen) is unknown for many interesting and potentially effective space materials. The use of specialized surfaces that can be tailored for very specific applications is an emerging trend and is likely to continue.

When black surfaces are used effectively in optical systems, the performance of the system can be greatly enhanced. However, when other approaches to stray light control are not implemented, even the “blackest” possible surfaces cannot establish the desired system performance. This cruel fact is a persuasive basis for an integration of black surface selection into the general stray light design studies, which today are nearly mandatory for the design of high-performance optical devices.

6.10 ACKNOWLEDGMENTS

The National Optical Astronomy Observatory is operated by the Association of Universities for Research in Astronomy (AURA), Inc. under cooperative agreement with the National Science Foundation. Special thanks to J. Heaney (NASA Goddard Space Flight Center), S. H. C. P. McCall (Stellar Optics

Laboratory, Inc.), and S. Smith (NASA Ames Research Center) for the use of their unpublished data. Thanks to R. Bonaminio (Lord Corp.), C. Bowers (Hughes Aircraft Co.), B. Banks (NASA Lewis Research Center), R. Harada (IIT Research Institute), D. Janeczko (Martin Marietta Corp.), A. Olson (Ball Aerospace), R. Seals (Oak Ridge National Laboratory), D. Shepard (Lockheed Martin), and S. Smith (NASA Ames Research Center) for providing data and photographs. Thanks also to J. Martin, L. Bergquist, F. Bartko, R. Culver, and S. Russak for support in the early stages of this work, and to D. Shepard, S. McCall, S. Smith, and W. Wolfe for helpful suggestions on this manuscript.

6.11 REFERENCES

1. W. L. Wolfe, "Scattered Thoughts on Baffling Problems," *Proc. SPIE: Radiation Scattering in Optical Systems*, SPIE, Bellingham, Wash., **257**:2 (1980).
2. R. P. Breault, "Stray Light Technology Overview of the 1980 Decade (And a Peek into the Future)," *Proc. SPIE: Stray Radiation in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1331**:2–11 (1990).
3. D. W. Bergener, S. M. Pompea, D. F. Shepard, and R. P. Breault, "Stray Light Performance of SIRTf: A Comparison," *Proc. SPIE: Stray Radiation IV*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **511**:65–72 (1984).
4. S. M. Pompea, D. F. Shepard, and S. Anderson, "BRDF Measurements at 6328 Angstroms and 10.6 Micrometers of Optical Black Surfaces for Space Telescopes," *Proc. SPIE: Stray Light and Contamination in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **967**:236–248 (1988).
5. S. M. Pompea and S. H. C. P. McCall, "Outline of Selection Process for Black Baffle Surfaces in Optical Systems," *Proc. SPIE: Stray Radiation in Optical Systems II*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1753** (1992a).
6. S. M. Pompea, J. E. Mentzell, and W. A. Siegmund, "A Stray Light Analysis of the 2.5 Meter Telescope for the Sloan Digital Sky Survey," *Proc. SPIE: Stray Radiation in Optical Systems II*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1753** (1992b).
7. S. M. Pompea, R. Pfisterer, and J. Morgan, "A Stray Light Analysis of the Apache Point Observatory 3.5-Meter Telescope System," *Proc. SPIE*, SPIE, Bellingham, Wash., **4842**:128–138, (2003).
8. S. M. Pompea, "The Management of Stray Radiation Issues in Space Optical Systems," *Space Science Reviews* **74**:181–193 (1995).
9. M. J. Persky, "Review of Black Surfaces for Space-Borne Infrared Systems," *Rev. Sci. Instrum.* **70**:2193–2217 (1999).
10. John L. Miller, "Multispectral Infrared BRDF Forward-Scatter Measurements of Common Black Surface Preparations and Materials—or "how black is black in the IR?" *Proc. SPIE*, **5405**:25–35 (2004).
11. M. J. Persky and M. Szczesniak, "Infrared, Spectral, Directional-Hemispherical Reflectance of Fused Silica, Teflon Polytetrafluoroethylene Polymer, Chrome Oxide Ceramic Particle Surface, Pyromark 2500 Paint, Krylon 1602 Paint, and Duraflect Coating," *Appl. Opt.* **47**:1389–1396 (2008).
12. R. E. Hahn and B. O. Seraphin, "Spectrally Selective Surfaces for Photothermal Solar Energy Conversion," in *Physics of Thin Films*, vol. 10, Academic Press, Orlando, Fla., 1978.
13. C. G. Granqvist, *Spectrally Selective Surfaces for Heating and Cooling Applications*, vol. TT 1, SPIE Optical Engineering Text, Bellingham, Wash., 1989.
14. K. M. Yousif, B. E. Smith, and C. Jeynes, "Investigation of Microstructure of Molybdenum-Copper Black Electrodeposited Coatings with Reference to Solar Selectivity" *J. Mater. Sci.* **31**:185–191 (1996).
15. D. B. Betts, F. J. J. Clarke, L. J. Cox, and J. A. Larkin, "Infrared Reflection Properties of Five Types of Black Coating for Radiometric Detectors," *J. Phys. E: (London) Sci. Instrum.* **18**:689–696 (August, 1985).
16. W. L. Wolfe, "Optical Materials," Chap. 7, in *The Infrared Handbook*, W. L. Wolfe and G. J. Zissis (eds.), ERIM, Ann Arbor, Michigan, 1978.
17. S. H. C. P. McCall, S. M. Pompea, R. P. Breault, and N. L. Regens, "Reviews of Properties of Black Surfaces for Ground and Space-Based Optical Systems," *Proc. SPIE: Stray Radiation II*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1753** (1992).
18. S. H. C. P. McCall, "Optical Properties (UV and Visible) of Black Baffle Materials and Processes for Use in Space," *Proprietary Report from Stellar Optics Laboratories* (78 Normark Drive, Thornhill, Ontario, Canada, L3T 3R1) to *Space Astrophysics Laboratory*, ISTS, Ontario, Canada, June 17, 1992.

19. S. M. Smith and R. V. Howitt, "Survey of Materials for an Infrared-Opaque Coating," *Proc. SPIE: Infrared, Adaptive, and Synthetic Aperture Optical Systems*, SPIE, Bellingham, Wash., **643**:53–62 (1986).
20. C. Schaub, M. Davis, G. Inouye, and P. Schaller, "Visible Scatter Measurements of Various Materials," *Proc. SPIE: Stray Radiation in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1331**:293–298 (1990).
21. K. A. Klicker, D. M. Fuhrman, D. R. Bjork, "A BSDF Database," *Proc. SPIE: Stray Radiation in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1331**:270–279 (1990).
22. S. H. C. P. McCall, R. L. Sinclair, S. M. Pompea, and R. P. Breault, "Spectrally Selective Surfaces for Ground and Space-Based Instrumentation: Support for a Resource Base," *Proc. SPIE: Space Astronomical Telescopes and Instruments II*, P. Bely and J. B. Breckinridge (ed.), SPIE, Bellingham, Wash., **1945** (1993).
23. S. H. C. P. McCall, *Black Materials Database of Stellar Optics Laboratories*, 78 Normark Drive, Thornhill, Ontario, Canada, L3T 3R1, unpublished, 1992c.
24. Lord Chemical Products Division, Industrial Coatings Office, 845 Olive Avenue, Novato, CA 94945, 1992.
25. D. B. Griner, "BRDF Measurements of Stray Light Suppression Coatings for the Space Telescope," *Proc. SPIE*, SPIE, Bellingham, Wash., **183**:98 (1979).
26. R. Fernandez, R. G. Seasholtz, L. G. Oberle, and J. R. Kadambi, "Comparison of the Bidirectional Reflectance Distribution Function of Various Surfaces," *Proc. SPIE*, SPIE, Bellingham, Wash., **967**:292 (1988).
27. C. W. Brown and D. R. Smith, "High-Resolution Spectral Reflection Measurements on Selected Optical-Black Baffle Coatings in the 5–20 Micrometer Region," *Proc. SPIE: Stray Radiation in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1331**:210–240 (1990).
28. A. J. Ames, "Z306 Black Paint Measurements," *Proc. SPIE: Stray Radiation in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1331**:299–304 (1990).
29. A. Lompadó, B. W. Murray, J. S. Wollam, and J. F. Meroth, "Characterization of Optical Baffle Materials," *Proc. SPIE: Scatter from Optical Components*, J. Stover (ed.), SPIE, Bellingham, Wash., **1165**:212–226 (1989).
30. J. A. Muscari and T. O'Donnell, "Mass Loss Parameters for Typical Shuttle Materials," *Proc. SPIE*, SPIE, Bellingham, Wash., **287**:20 (1981).
31. S. M. Pompea, D. F. Bergener, D. W. Shepard, S. Russak, and W. L. Wolfe, "Reflectance Measurements on an Improved Optical Black for Stray Light Rejection from 0.3 to 500 Micrometers," *J. Opt. Eng.* **23**:149–152 (1984).
32. S. M. Smith, "Far Infrared Reflectance Spectra of Optical Black Coatings," *Proc. SPIE: Scattering in Optical Materials*, S. Musikant (ed.), Optical Society of America, Washington, DC, **362**:57–59 (1982).
33. S. M. Smith, "Specular Reflectance of Optical-Black Coatings in the Far Infrared," *Appl. Opt.* **23**(14):2311–2326 (1984).
34. S. M. Smith, "Cryo-Mechanical Tests of Ames 24E2 IR-Black Coating," *Proc. SPIE: Stray Radiation in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1331**:241–248 (1990).
35. W. Viehmann and R. E. Predmore, "Ultraviolet and visible BRDF Data on Spacecraft Thermal Control and Baffle Materials," *Proc. SPIE*, SPIE, Bellingham, Wash., **675**:67 (1986).
36. CAL, "Space Systems Process Procedure for the Preparation and Application of Chemglaze Z306 Black Paint for Space Applications," *Document PQ-CAL-SR-10333* from Canadian Astronautics Ltd., 1986.
37. M. M. Nordberg, C. Von Benken, and E. J. Johnson, "Reflectivity Changes in Optical Baffle Materials Following Pulsed Electron Bombardment," *Proc. SPIE*, SPIE, Bellingham, Wash., **1050**:185 (1989).
38. R. J. Noll, R. Harned, R. Breault, and R. Malugin, "Stray Radiation and Infrared Astronomical Satellite (IRAS) Telescope," *Proc. SPIE*, SPIE, Bellingham, Wash., **257**:119 (1980).
39. S. M. Smith, "Far Infrared (FIR) Optical Black Bidirectional Reflectance Distribution Function (BRDF)," *Proc. SPIE: Radiation Scattering in Optical Systems*, SPIE, Bellingham, Wash., **362**:161–168 (1980).
40. P. Jelinsky and S. Jelinsky, "Low Reflectance EUV Materials: A Comparative Study," *Appl. Opt.* **26**(4):613–615 (1987).
41. D. L. Stierwalt, "Infrared Absorption of Optical Blacks," *Opt. Eng.* **18**(2):147–151 (1979).
42. R. R. Willey, R. W. George, J. G. Ohmart, and J. W. Walvoord, "Total Reflectance Properties of Certain Black Coatings from 0.2 to 20 Micrometers," *Proc. SPIE: Generation, Measurement and Control of Stray Radiation III*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **384**:19–26 (1983).
43. D. C. Evans and R. P. Breault, "APART/PADE Analytical Evaluation of the Diffuse Infrared Background Experiment for NASA's Cosmic Background Explorer," *Proc. SPIE: Stray Radiation IV*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **511**:54–64 (1984).

44. T. Heslin, J. Heaney, and M. Harper, "The Effects of Particle Size on the Optical Properties and Surface Roughness of a Glass-Balloon-Filled Black Paint," *NASA Technical Note TND-7643*, Goddard Space Flight Center, Greenbelt, MD., May 1974.
45. D. C. Evans, "Principles of Stray Light Supp and Conceptual Application to the Design of the Diffuse Infrared Background Experiment for NASA's Cosmic Background Explorer," *Proc. SPIE*, SPIE, Bellingham, Wash., **384**:82 (1983).
46. S. M. Smith, "The Reflectance of Ames 24E, Infrablack, and Martin Black," *Proc. SPIE: Stray Radiation V*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **967**:248–254 (1988).
47. S. M. Nee and H. E. Bennett, "Characterization of Optical Blacks by Infrared Ellipsometry and Reflectometry," *Proc. SPIE* **1331**:249 (1990).
48. Cardinal, product information and MSDSs on Velvathane and 6450-01 paints, El Monte, Calif., 1992.
49. J. R. Houck, "New Black Paint for Cryogenic Infrared Applications," *Proc. SPIE*, S. Musikant (ed.), SPIE, Bellingham, Wash., **362**:54–56 (1982).
50. S. M. Smith, "A Simple Antireflection Overcoat for Opaque Coatings in the Submillimeter Region," *Proc. SPIE: Stray Radiation V*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **675**:55–60 (1986).
51. M. B. Birnbaum, E. C. Metzler, and E. L. Cleveland, "Electrically Conductive Black Optical Black Paint," *Proc. SPIE: Scattering in Optical Materials*, S. Musikant (ed.), SPIE, Bellingham, Wash., **362**:60–70 (1982).
52. E. C. Metzler, "Application of Temperature Control Paints," *JPL Doc. No. FS501424 D*, 21 March 1988.
53. R. P. Breault, "Stray Light Technology Overview in 1988," *Proc. SPIE*, SPIE, Bellingham, Wash., **967**:2 (1988).
54. IITRI, product information and MSDSs on paints MH-211, D111, and MH2200, from the Illinois Institute of Technology Research Institute, Chicago, Ill., 1992.
55. J. R. Grammar, L. J. Balin, M. D. Blue, and S. Perkowitz, "Absorbing Coatings for the Far Infrared," *Proc. SPIE: Radiation Scattering in Optical Systems*, G. Hunt (ed.) SPIE, Bellingham, Wash., **257**:192–195 (1980).
56. W. Driscoll, (ed.), and W. Vaughan (assoc. ed.), *Handbook of Optics*, McGraw-Hill, New York, 1978, pp. 7–112.
57. S. M. Smith, "Analysis of 12–700 Micrometer Reflectance Spectra of Three Optical Black Samples," *Proc. SPIE: Generation, Measurement, and Control of Stray Radiation III*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **384**:32–36 (1983).
58. C. L. Wyman, D. B. Griner, G. H. Hunt, and G. B. Shelton, *Opt. Eng.* **14**(6):528 (1975).
59. R. P. Breault, "Specular Vane Cavities," *Proc. SPIE: Generation, Measurement, and Control of Stray Radiation III*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **384**:90–97 (1983).
60. AKZO Coatings, Inc., product information and MSDSs on 443-3-17 and 443-3-8, Orange, Calif., 1992.
61. E. R. Freniere and D. L. Skelton, "Use of Specular Black Coatings in Well-Baffled Optical Systems," *Proc. SPIE: Stray Radiation V*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **675**:126–132 (1986).
62. T. R. Gull, H. Hertzog, F. Osantowski, and A. R. Toft, "Low Orbit Effects on Optical Coatings and Materials as Noted on Early Shuttle Flights," *RAL Workshop on Advanced Technology Reflectors for Space Instrumentation*, M. Grande (ed.), SPIE, Bellingham, Wash., June 16–18, 1986.
63. Martin Marietta, product information from Martin Marietta, 1992.
64. S. M. Pompea, D. W. Bergener, D. F. Shepard, and K. S. Williams, "The Effects of Atomic Oxygen on Martin Black and Infrablack," *Proc. SPIE*, R. P. Breault (ed.), **511**:24–30 (1984).
65. G. Geikas, "Scattering Characteristics of Etched Electroless Nickel Coatings," *Proc. SPIE: Generation, Measurement, and Control of Stray Radiation III*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **384**:10–18 (1983).
66. F. M. Cady, D. R. Cheever, K. A. Klicker, and J. C. Stover, "Comparison of Scatter Data from Various Beam Pumps," *Proc. SPIE*, Spie, Bellingham, Wash., **818**:21 (1987).
67. F. O. Bartell, J. E. Hubbs, M. J. Nofziger, and W. L. Wolfe, *Appl. Opt.* **21**(17):3178 (1982).
68. L. D. Brooks, J. E. Hubbs, F. O. Bartell, and W. L. Wolfe, "Bidirectional Reflectance Distribution Function of the Infrared Astronomical Satellite Solar-Sheild Material," *Appl. Opt.* **21**:2465 (1982).
69. R. P. Breault, "Specification of the Scattering Characteristics of Surfaces and Systems for Use in the Analysis of Stray Light," *Proc. SPIE*, SPIE, Bellingham, Wash., **181**:108 (1979).

70. P. J. Young, R. Noll, L. Andreozzi, and J. Hope, "Particle Contamination from Martin Optical Black," *Proc. SPIE*, SPIE, Bellingham, Wash., **257**:196 (1980).
71. D. F. Shepard, S. M. Pompea, and S. Anderson, "The Effect of Elevated Temperatures on the Scattering Properties of an Optical Black Surface at 0.6328 and 10.6 Micrometers," *Proc. SPIE: Stray Radiation V*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **967**:286–291 (1988).
72. S. M. Smith, "Bidirectional Reflectance Distribution Function (BRDF) Measurements of Sunshield and Baffle Materials for the Infrared Astronomy Satellite (IRAS) Telescope," *Proc. SPIE: Modern Utilization of Infrared Technology VII*, SPIE, Bellingham, Wash., **304**:205–213 (1981).
73. B. Vitt, "Properties of Black Cobalt Coatings," *Proc. SPIE*, SPIE, Bellingham, Wash., **823**:218 (1987).
74. M. G. Hutchins, P. J. Wright, and P. D. Grebenik, "SnO₂: Sb Dip Coated Films on Anodized Aluminium Selective Absorber Plates," *Proc. SPIE*, SPIE, Bellingham, Wash., **653**:188 (1986).
75. C. E. Johnson, "Unique Surface Morphology with Extremely High Light Absorption Capability," *Proc. Electro-less Nickel Conference*, No. 1, Cincinnati, Ohio, November 6–7, 1979.
76. C. E. Johnson, "Black Electroless Nickel Surface Morphologies with Extremely High Light Absorption Capacity," *Metal Finishing*, July, 1980.
77. J. J. Cuomo, S. M. Rossnagel, and H. R. Kaufman, *Handbook of Ion Beam Processing Technology: Principles, Deposition, Film Modification, and Synthesis*, Noyes Publications, Park Ridge, NJ, 1989. (See Chap. 17, Banks, B. A., "Topography Texturing Effects," 1989.)
78. B. A. Banks, "Ion Beam Applications Research—A 1981 Summary of Lewis Research Center Programs," *NASA Technical Memorandum 81721*, 1981.
79. M. J. Mirtich, S. K. Rutledge, N. Stevens, R. Olle, and J. Merrow, "Ion Beam Textured and Coated Surfaces Experiment (IBEX)," Presented at *LDEF First Postretrieval Symposium*, Orlando, Fla., June 3–8, 1991.
80. J. S. Wollam and B. W. Murray, *Proc. SPIE*, SPIE, Bellingham, Wash., **1118**:88 (1989).
81. Spire, product information from Spire Corporation (Appendix A–D of Spire Document FR 10106), Bedford, Mass., (1992).
82. C. M. Egert and D. D. Allred, "Diffuse Absorbing Beryllium Coatings Produced by Magnetron Sputtering," *Proc. SPIE: Stray Radiation in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1331**:170–178 (1990).
83. C. C. Blatchley, E. A. Johnson, Y. K. Pu, and C. Von Benken, "Rugged Dark Materials for Stray Light Suppression by Seeded Ion Beam Texturing," *Proc. SPIE: Stray Radiation II*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1753** (1992).
84. G. Hass, H. H. Schroeder, and A. F. Turner, "Mirror Coatings for Low Visible and High Infrared Reflectance," *J. Opt. Soc. Am.* **46**:31–35 (1956).
85. J. A. Dobrowolski, "Versatile Computer Program for Absorbing Optical Thin Film Systems," *Appl. Opt.* **20** (1981).
86. J. A. Dobrowolski, F. C. Ho, and A. J. Waldorf, "Research on Thin Film Anticounterfeiting Coatings at the National Research Council of Canada," *Appl. Opt.* **28**:2702–2717 (1989).
87. J. A. Dobrowolski, B. T. Sullivan, and R. C. Bajcar, "An Optical Interference, Contrast Enhanced Electroluminescent Device," *Appl. Opt.* **31**:5988 (1992).
88. W. L. Wolfe and Y. Wang, "Comparison of Theory and Experiments for Bidirectional Reflectance Distribution Function (BRDF) of Microrough Surfaces," *Proc. SPIE: Scattering in Optical Materials*, S. Musikant (ed.), SPIE, Bellingham, Wash., **362**:40 (1982).
89. T. A. Leonard, T. A. and M. Pantoliano, "BRDF Round Robin," *Proc. SPIE: Stray Light and Contamination in Optical Systems*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **967**:226–235 (1988).
90. T. A. Leonard, M. Pantoliano, and J. Reilly, "Results of a CO₂ BRDF Round Robin," *Proc. SPIE: Scatter from Optical Components*, J. Stover (ed.), SPIE, Bellingham, Wash., **1165**:444–449 (1989).
91. S. M. Pompea, "Stray Radiation Issues on Adaptive Optics Systems," in *Adaptive Optics for Astronomy*, D. Alloin and J.-M. Mariotti (eds.), Kluwer Academic Publishers, Dordrecht, 1994.
92. J. H. Henninger, "Solar Absorptance and Thermal Emittance of Some Common Spacecraft Thermal Control Coatings," *NASA Reference Publication 1121*, 1984, p. 7.
93. J. Stover, *Optical Scattering: Measurement and Analysis*, 2nd ed., SPIE Press, Bellingham, Wash., 1995.

94. E. R. Freniere, "First-Order Design of Optical Baffles," *Proc. SPIE: Radiation Scattering in Optical Systems*, G. Hunt (ed.), SPIE, Bellingham, Wash., **257**:19–22 (1980).
95. G. Peterson and S. Johnston, "Specular Baffles," *Proc. SPIE: Stray Radiation in Optical Systems II*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **1753** (1992).
96. K. Snail, P. Brown, J. Costantino, W. C. Shemano, C. W. Schmidt, W. F. Lynn, C. L. Seaman, and T. R. Knowles, "Optical Characterization of Black Appliqués" *Proc. SPIE* **2864**:465 (1996).
97. S. R. Meier, "Characterization of Highly Absorbing Black Appliqués in the Infrared," *Appl. Opt.* **40**:2788–2795 (2001).
98. S. R. Meier, "Reflectance and Scattering Properties of Highly Absorbing Black Appliqués over a Broadband Spectral Region," *Appl. Opt.* **40**:6260–6264 (2001).
99. S. R. Meier, M. L. Korwin, and C. I. Merzbacher, "Carbon Aerogel: A New Nonreflective Material for the Infrared," *Appl. Opt.* **39**:3940–3944 (2000).
100. J. M. Bennett and L. Mattsson, *Introduction to Surface Roughness and Scattering*, Optical Society of America, 1989.
101. C. F. Bohren and D. F. Huffman, *Absorption and Scattering of Light by Small Particles*, John Wiley and Sons, New York, 1983.
102. R. Ludwig, "Antireflection Coating on a Surface With High Reflecting Power," U. S. Patent 4,425,022, Jan. 10, 1984.
103. J. A. Dobrowolski, E. H. Hara, B. T. Sullivan, and A. J. Waldorf, "A High Performance Optical Wavelength Multiplexer/Demultiplexer," *Appl. Opt.* **30** (1991).
104. W. I. Linlor, "Reflective Baffle System with Multiple Bounces," *Proc. SPIE: Stray Radiation IV*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **675**:217–239 (1986).
105. A. W. Greynolds and R. K. Melugin, "Analysis of an All-Specular Linlor Baffle Design," *Proc. SPIE: Stray Radiation IV*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **675**:240–248 (1986).
106. H. Babel and H. G. Lee, "High emittance low absorptance coatings," US Patent 5,296,285 (1994).
107. A. P. Glassford, (ed.), *Proc. SPIE: Optical System Contamination: Effects, Measurements, Control*, SPIE, Bellingham, Wash., **777** (1987). (Entire volume)
108. A. P. Glassford, (ed.), *Proc. SPIE: Optical System Contamination: Effects, Measurements, Control II*, SPIE, Bellingham, Wash., **1329** (1990). (Entire volume)
109. A. P. Glassford, (ed.), *Proc. SPIE: Optical System Contamination: Effects, Measurements, Control III*, SPIE, Bellingham, Wash., **1754** (1992). (Entire volume)
110. J. B. Heaney, "A Comparative Review of Optical Surface Contamination Assessment Techniques," *Proc. SPIE: Optical System Contamination: Effects, Measurements, Control*, P. A. Glassford (ed.), SPIE, Bellingham, Wash., **777**:179 (1987).
111. K. Nahm, P. Spyak, and W. Wolfe, "Scattering from Contaminated Surfaces," *Proc. SPIE: Scatter from Optical Components*, J. Stover (ed.), SPIE, Bellingham, Wash., **1165**:294–305 (1989).
112. V. L. Williams and R. T. Lockie, "Optical Contamination Assessment by Bidirectional Reflectance-Distribution Function (BRDF) Measurement," *Opt. Eng.* **18**(2):152–156 (1979).
113. R. P. Young, "Low Scatter Mirror Degradation by Particle Contamination," *Opt. Eng.* **15**(6):516–520 (1976).
114. B. W. Murray and Johnson, E. A., "Pulsed Electron Beam Testing of Optical Surfaces," *Proc. SPIE: Conference on Optical Surfaces Resistant to Severe Environments*, paper no. 1330-01, San Diego, July 11, 1990.
115. M. McCargo, R. E. Dammann, J. C. Robinson, and R. J. Milligan, "Erosion of Diamond Films and Graphite in Oxygen Plasma," *Proceedings of the International Symposium on Environmental and Thermal Control Systems for Space Vehicles*, Joulouse, France, European Space Agency, Noordwijk, The Netherlands, October, 1983, pp. 1–5.
116. A. F. Whittaker, "Atomic Oxygen Effects on Materials," *STS-8 Paint Data Summary*, Marshall Space Flight Center, January 1984.
117. J. L. Golden, "Results of an Examination of the A-276 White and Z-306 Black Thermal Control Paint Disks flown on LDEF," *NASA Conference Publication 10072*, First LDEF Post-Retrieval Symposium Abstracts, June 2–8, 1991.
118. M. D. Blue, and S. Perkowitz, "Space-Exposure Effects on Optical-Baffle Coatings at Far-Infrared Wavelengths," *Appl. Opt.* **31**:4305–4309 (1991).

119. W. A. Campbell and J. J. Scialdone, "Outgassing Data for Selecting Spacecraft Materials," *NASA Reference Publication 1134*, Revision 2, November 1990.
120. R. E. Predmore and E. W. Mielke, *Materials Selection Guide, Revision A*, Goddard Space Flight Center, August 1990.
121. R. V. Peterson, W. Krone-Schmidt, and W. V. Brandt, "Jet-Spray Cleaning of Optics," *Proc. SPIE: Optical System Contamination: Effects, Measurement, Control III*, A. P. Glassford (ed.), SPIE, Bellingham, Wash., **1754** (1992).
122. J. A. Gunderson, "Goniometric Reflection Scattering Measurements and Techniques at 10.6 Micrometers," Thesis, Univ. of Arizona, 1977.
123. S. M. Smith and W. L. Wolfe, "Comparison of Measurements by Different Instruments of the Far-Infrared Reflectance of Rough, Optically Black Coatings," *Proc. SPIE: Scattering of Optical Materials*, S. Musikant (ed.), SPIE, Bellingham, Wash., **362**:46–53 (1983).
124. R. P. Breault, A. Greynold, and S. Lange, "APART/PADE Version 7: A Deterministic Computer Program Used to Calculate Scattered and Diffracted Energy," *Proc. SPIE: Radiation Scattering in Optical Systems*, G. Hunt (ed.), SPIE, Bellingham, Wash., **257** (1980).
125. A. W. Greynolds, *Advanced Systems Analysis Package (ASAP) Users Manual*, Breault Research Organization, Inc., Tucson, Ariz., 1988.
126. S. S. Steadman and B. K. Likeness, "GUERAP III Simulation of Stray Light Phenomena," *Proc. SPIE: Stray Light Problems in Optical Systems*, SPIE, Bellingham, Wash., **107** (1977).
127. R. P. Breault, "Current Technology of Stray Light," *Proc. SPIE: Stray Radiation V*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **675**:4–13 (1986).
128. F. D. Orazio, W. K. Stowell, and R. M. Silva, "Instrumentation of a Variable Angle Scatterometer (VAS)," *Proc. SPIE: Stray Radiation III*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **384**:123–131 (1983).
129. F. O. Bartell, "BRDF Measurement Equipment: Intrinsic Design Considerations," *Proc. SPIE: Stray Radiation IV*, R. P. Breault (ed.) **511**:31–34 (1984).
130. R. M. Silva, F. D. Orazio, and R. B. Sledge, "A New Instrument for Constant (Beta-Beta_g) Scatter Mapping of Contiguous Optical Surfaces of up to 25 Square Inches," *Proc. SPIE: Stray Radiation IV*, R. P. Breault (ed.), SPIE, Bellingham, Wash., **511**:38–43 (1984).
131. J. B. Heaney, Optics Branch, Code 717, Goddard Space Flight Center, Greenbelt, MD., private communication (unpublished data), 1992.
132. D. L. Stierwalt, J. B. Bernstein, and D. D. Kirk, "Measurements of the Infrared Spectral Absorbance of Optical Materials," *Appl. Opt.* **2**(11):1169–1173 (1963).
133. F. J. J. Clarke and J. A. Larkin, "Measurements of Total Reflectance, Transmittance and Emissivity over the Thermal IR Spectrum," *Infrared Physics* **25**:359–367 (1985).
134. S. M. Smith, "Formation of Ames 24E2 IR-Black Coatings," *NASA Tech. Memo 102864*, July (1991).
135. S. M. Smith, "BRDFs of Ames 24E, Ames 24E2, and Ames 47A at Photometric Wavelengths of 220 and 350 micrometers," *Sterling Technical Note TN-91-8441-000-74*, Sterling Software, Moffett Field, Calif., October 1991.
136. S. M. Smith, "An Almost "Perfectly" Diffuse, "Perfect" Reflector for Far-Infrared Reflectance Calibration," *Proc. SPIE: Stray Radiation in Optical Systems II*, Breault (ed.) **1753**:252–261 (1992).
137. S. M. Smith, Sterling Software, NASA Ames Research Center, private communication, April 1992.
138. S. M. Pompea, D. W. Bergener, D. F. Shepard, S. L. Russak, and W. L. Wolfe, "Preliminary Performance Data on an Improved Optical Black for Infrared Use," *Proc. SPIE: New Optical Materials*, SPIE, Bellingham, Wash., **400**:128 (1983).
139. J. Lohrengel, "Total Emissivity of Black Coatings," ("Gesamtemissionsgrad von Schwärzen"), *Wärme - Stoffübertrag* Springer-Verlag, Berlin, **21**(5):311–315 (27 July 1987).
140. V. Freibell, *Colorado, M & P Information Bulletin No. 82.21*, Ball Aerospace, Boulder, 7/27/82.
141. I. T. Lewis, A. R. Telkamp, and A. F. Ledebuhr, "Low Scatter Edge Blackening Compounds for Refractive Optical Elements," *Proc. SPIE: Scatter from Optical Components*, J. Stover (ed.), SPIE, Bellingham, Wash., **1165**:227–236 (1989).
142. J. F. Wade, J. E. Peyton, B. R. Klitzky, and R. E. Groff, "Optically Black Coating and Process for Forming It," U. S. Patent 4,111,762, September 5, 1978.

143. D. F. Shepard, Martin Black and Infrablack Contact at Lockheed Martin, P.O. Box 179, Denver, Colorado 80201.
144. J. F. Wade and W. R. Wilson, *Proc. SPIE*, SPIE, Bellingham, Wash., **67**:59 (1973).
145. S. M. Pompea, D. W. Bergener, and D. F. Shepard, "Reflectance Characteristics of an Infrared Absorbing Surface," *Proceedings 31st National Infrared Information Symposium*, National Bureau of Standards, Boulder, Colorado, 1983, p. 487.
146. W. Friedrichs, product literature from Tiodize Co., Inc., 5858 Engineer Drive, Huntington Beach, CA 92649.
147. S. Kodama, M. Horiuchi, T. Kunii, and K. Kuroda, "Ultra-Black Nickel-Phosphorus Alloy Optical Absorber," *Instrumentation and Measurement, IEEE Transactions* **39**:230–232 (1990).
148. C. Johnson, "Ultra-Black Coating due to Surface Morphology," U. S. Patent 4,361,630, 30, November 1982.
149. D. F. Shepard, R. J. Fenolia, D. C. Nagle, and M. E. Marousek, "High-Temperature, High-Emissivity, Optically Black Boron Surface," U. S. Patent 5,035,949, July 30, 1991, Martin Marietta Corporation.
150. R. D. Seals, C. M. Egert, and D. D. Allred, "Advanced Infrared Optically Black Baffle Materials," *Proc. SPIE: Optical Surfaces Resistant to Severe Environments*, SPIE, Bellingham, Wash., **1330**:164–177 (1990).
151. R. D. Seals, "Advanced Broadband Baffle Materials," *Proc. SPIE: Reflective and Refractive Materials for Earth and Space Applications*, SPIE, Bellingham, Wash., **1485**:78–87 (1991).
152. B. W. Murray and J. S. Wollam, "Space Durable Beryllium Baffle Materials," *SPIE Symposium on Aerospace Sensing*, Orlando, Fla., March 1989.
153. E. Johnson, Spire Corporation, Bedford, Mass., private communication, March 1992.
154. S. K. Rutledge, B. A. Banks, M. J. Mirtich, R. Lebed, J. Brady, D. Hotes, and M. Kussmaul, "High Temperature Radiator Materials for Applications in the Low Earth Orbital Environment," *NASA Technical Memorandum 100190*, 1987.
155. B. A. Banks, S. K. Rutledge, M. J. Mirtich, T. Behrend, D. Hotes, M. Kussmaul, J. Barry, C. Stidham, T. Stueber, and F. DiFillippo, "Arc-Textured Metal Surfaces for High Thermal Emittance Radiators," *NASA Technical Memorandum 100894*, 1988.
156. D. F. Shepard and R. J. Fenolia, "Optical Black Cobalt Surface," U. S. Patent 4,904,353, February 27, 1990, Martin Marietta Corporation.
157. R. J. Fenolia, D. F. Shepard, and S. L. Van Loon, "Optically Black Pliable Foils," U. S. Patent 4,894,125, January 16, 1990, Martin Marietta Corporation.
158. D. J. Janeczko, "Optics and Electro-Optics," Martin Marietta Corporation, Orlando, Fla., personal communication, 1992.
159. Z.-P. Yang, L. Ci, J. A. Bur, S.-Y. Lin, and P. M. Ajayan, "Experimental Observation of an Extremely Dark Material Made by a Low-Density Nanotube Array," *Nano Lett.* **8**: 446–451 (2008).
160. A. Y. Vorobyev and G. Guo, "Effect of Nanostructure-Covered Femtosecond Laser-Induced Periodic Surface Structures on Optical Absorptance of Metals," *Appl. Phys. A* **86**:321–324 (2007).
161. K. Paivasaari, J. J. Kaakkunen, M. Kuittinen, and T. Jaaskelainen, "Enhanced Optical Absorptance of Metals Using Interferometric Femtosecond Ablation," *Optics Express*, **15**:13838–13843 (2007).
162. P. Strimer, X. Gerbaux, A. Hadni, T. Souel, "Black Coatings of Infrared and Visible, with High Electrical Resistivity," *Infrared Physics* **21**:7–39 (1981).
163. R. S. Robinson, and S. M. Rossnagel, "Ion-Beam Induced Topography and Surface Diffusion," *J. Vac. Sci. Technol.* **21**:790 (1982).
164. R. B. Culver, W. A. Solberg, R. S. Robinson, and I. L. Spain, "Optical Absorption of Microtextured Graphite Surfaces in the 1.1–2.4 Micrometer Wavelength Region," *Appl. Opt.* **24**:924 (1984).
165. C. W. Bowers, R. B. Culver, W. A. Solberg, and I. L. Spain, "Optical Absorption of Surfaces Modified by Carbon Filaments," *Appl. Opt.* **26**:4625 (1987).
166. L. Harris, R. T. McGinnies, and B. M. Siegel, "The Preparation and Optical Properties of Gold Black," *J. Opt. Soc. Am.* **38**(7) (1948).
167. W. R. Blevin, and W. J. Brown, "Black Coatings for Absolute Radiometers," *Int. J. Sci. Metrology* **2**(4):139–143 (October 1966).
168. G. Zaeschmar and A. Nedoluha, "Theory of the Optical Properties of Gold Blacks," *J. Opt. Soc. Am.* **62**(3):348–352 (1972).

6.12 FURTHER READINGS

- Anderson, S., S. M. Pompea, D. F. Shepard, and R. Castonguay, "Performance of a Fully Automated Scatterometer for BRDF and BTDF Measurements at Visible and Infrared Wavelengths," *Proc. SPIE: Stray Light and Contamination in Optical Systems* **967**:159 (1988).
- Banks, B. A., *NASA Technical Memorandum 81721*, Lewis Research Center, 1981.
- Bartell, F. O., E. L. Dereniak, and W. L. Wolfe, "The Theory and Measurement of BRDF and BTDF," *Proc. SPIE: Radiation Scattering in Optical Systems* **257**:154 (1980).
- Bennett, H. E. and E. L. Church, "Surface Roughness Measurement in the Submicrometer Range Using Laser Scattering," *Opt. Eng.* **18**:103 (1979).
- Bonnot, A. M., H. Belkhir, D. Pailharey, and P. Mathiez, *Proc. SPIE* **562**:209 (1985).
- Bonnot, A. M., H. Belkhir, and D. Pailharey, *Proc. SPIE: Optical Materials Technology for Energy Efficiency and Solar Energy Conversion V* **653**:215 (1986).
- Brooks, L. D. and W. L. Wolfe, "Microprocessor-Based Instrumentation for Bidirectional Reflectance Distribution Function (BRDF) Measurements from Visible to Far Infrared (FIR)," *Proc. SPIE: Radiation Scattering in Optical Systems* **257**:177 (1980).
- Davis, L. and J. G. Kepros, "Improved Facility for BRDF/BTDF Optical Scatter Measurements," *Proc. SPIE: Stray Radiation V* **675**:24 (1986).
- Duran, M. and C. Gricurt, "Hipparcos Telescope Stray Light Protection," *Proc. SPIE: Stray Radiation V* **675**:134 (1986).
- "Method of Forming Electrodeposited Anti-reflective Surface Coatings, U.S. Patent 5,326,454, Martin Marietta Corporation.
- Freniere, E. R., "Use of Specular Black Coatings in Well-Baffled Optical Systems," *Proc. SPIE: Stray Radiation V* **675**:126 (1986).
- Garrison, J. D., J. C. Haiad, and A. J. Averett, "Progress in the Commercialization of a Carbonaceous Solar Selective Absorber on a Glass Substrate," *Proc. SPIE: Optical Materials Technology for Energy Efficiency and Solar Energy Conversion VI* **823**:225 (1987).
- Hunt, P. J., R. Noll, L. Andreozzi, and J. Hope, "Particle Contamination from Martin Optical Black," *Proc. SPIE: Radiation Scattering in Optical Systems*, G. Hunt (ed.) **257**:196 (1980).
- LDEF, "The Long Duration Exposure Facility (LDEF) Mission I Experiments," *NASA SP-473*, NASA Langley Research Center, Washington, D.C., November 1986.
- Leger, L. J., J. T. Visentine, and J. F. Kuminecz, "Low Earth Orbit Oxygen Effects on Surfaces," presented at *AIAA 22nd Aerospace Sciences Meeting*, Reno, Nevada, January 9–12, 1984.
- Mende, S. B., P. M. Banks, and D. A. Klingensmith III, "Observation of Orbiting Vehicle Induced Luminosities on the STS-8 Mission," *Geophys. Res. Lett.* **11**:527 (1984).
- Ranbar Technology, Inc., Glenshaw, Pa., product information and MSDSs on Ranbar G113 Paint, 1990.
- Redspot Paint and Varnish Co., Inc., product information on Nextel Suede, P.O. Box 418, Evansville, In., 47703, 1992.
- Schiff, T. F., J. C. Stover, D. R. Cheever, and D. R. Bjork, "Maximum and Minimum Limitations Imposed on BSDF Measurements," *Proc. SPIE* **967**:50 (1988).
- Seraphin, B. O., (ed.), "Solar Energy Conversion—Solid State Physics Aspects," in *Topics in Applied Physics*, vol. 31, Springer Verlag, 1979.
- Sikkens Aerospace Finishes Division of Akzo Coatings America, Inc., 434 W. Meats Avenue, Orange, CA 92665.
- Stover, J. C., and C. H. Gillespie, *Proc. SPIE* **362**:172 (1982).
- Stover, J. C., C. H. Gillespie, F. M. Cady, D. R. Cheever, and K. A. Klicker, *Proc. SPIE* **818**:62 (1987).
- Stover, J. C., C. Gillespie, F. M. Cady, D. R. Cheever, and K. A. Klicker, *Proc. SPIE* **818**:68 (1987).
- Visentine, J. T., (ed.), *Atomic Oxygen Effects Measurement for Shuttle Missions STS-8 and 41-G, Volume I–III*, NASA TM-100459, 1988.
- Wood, B. E., W. T. Bertrand, E. L. Kiech, J. D. Holt, and P. M. Falco, *Surface Effects of Satellite Material Outgassing Products*, Final Report AEDC-TR-89-2, Arnold Engineering Development Center, 1989.

This page intentionally left blank.

DO NOT DUPLICATE

OPTICAL PROPERTIES OF FILMS AND COATINGS

Jerzy A. Dobrowolski

*Institute for Microstructural Sciences
National Research Council of Canada
Ottawa, Ontario, Canada*

7.1 INTRODUCTION

Scope of Chapter

In the broadest sense of the term, an optical filter is any device or material which is deliberately used to change the spectral-intensity distribution or the state of polarization of the electromagnetic radiation incident upon it. The change in the spectral intensity distribution may or may not depend on the wavelength. The filter may act in transmission, in reflection, or in both.

Filters can be based on many different physical phenomena, including absorption, refraction, interference, diffraction, scattering, and polarization. For a comprehensive review of this broader topic the interested reader is referred to the chapter entitled “Coatings and Filters” which appeared in the first edition of the *Handbook of Optics*.¹ This chapter deals only with filters that are based on absorption and interference of electromagnetic radiation in thin films. Optical thin-film coatings have numerous applications in many branches of science and technology and there are also many consumer products that use them. The spectral region covered in this chapter extends from about 0.003 to 300 μm (3 to 3×10^5 nm), although the main emphasis is on filters for the visible and adjacent spectral regions.

The discussion in this chapter is largely confined to generic thin-film filters, such as antireflection coatings, cutoff filters, narrowband transmission or rejection filters, reflectors, beam splitters, and so forth. Filters for very specific applications, such as filters for colorimeters and other scientific instruments, color correction filters, and architectural coatings, are, as a rule, not treated. Of the filters described, many are available commercially while others are only research laboratory prototypes. This review does not cover thin-film filters whose properties can be changed by external electric or magnetic fields, temperature, or illumination level.

In this introductory section some general considerations on the use of optical filters are presented. In the following sections, the theory of optical multilayers and the methods for their deposition and characterization are briefly discussed. These sections are useful for gaining a proper understanding of the operation, advantages, and limitations of optical coatings. The remaining sections then describe the properties of various generic thin-film filters.

For further information on this subject the interested reader is particularly encouraged to consult the books by Macleod,² Rancourt,³ and Baumeister.⁴

General Theory of Filters

There are many different ways of describing the performance of optical coatings and filters. For example, transmission and reflection filters intended for visual applications are adequately described by a color name alone, or by reference to one of the several existing color systems (see Chap. 10, Vol. III). There also exist other specialized filter specifications for specific applications. However, the most complete information on the performance of a filter is provided by spectral transmittance, reflectance, absorptance, and optical density curves. This is the method adopted in this chapter.

Referring to Fig. 1, at a wavelength λ the spectral transmittance $T(\lambda)$ of a filter placed between two semi-infinite media is equal to the ratio of the light intensity of that wavelength transmitted $I_T(\lambda)$ by the filter to that incident $I_0(\lambda)$ upon it,

$$T(\lambda) = \frac{I_T(\lambda)}{I_0(\lambda)} \quad (1)$$

The spectral reflectance $R(\lambda)$ of a filter is defined in a similar way,

$$R(\lambda) = \frac{I_R(\lambda)}{I_0(\lambda)} \quad (2)$$

At non-normal incidence the component of the intensity perpendicular to the interface must be used in the above equations.² The relationship between the transmittance $T(\lambda)$ and the density $D(\lambda)$ of a filter is given by

$$D(\lambda) = \log \frac{1}{T(\lambda)} \quad (3)$$

In this chapter, transmittances and reflectances will be plotted using either linear or logarithmic scales. The logarithmic scale is particularly well suited whenever accurate information about the low transmission or reflection region is to be conveyed. However, this is done at the expense of detail at the high end of the scale. A variant, the decibel scale, equal to $-10D(\lambda)$, is frequently used in the telecommunications field. Wavelengths are normally specified in micrometers (μm) because this is the most convenient unit for the whole spectral range covered in this chapter. In the following discussions the dependence of the transmittance, reflectance, and absorptance on wavelength will be implicitly assumed.

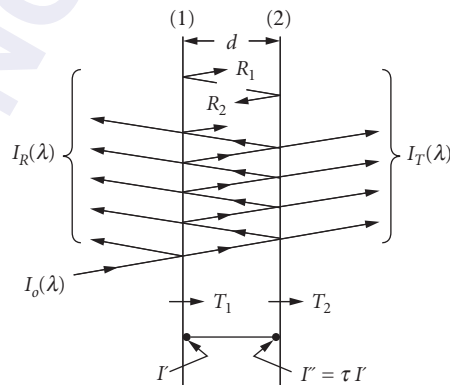


FIGURE 1 Specular transmission and reflection of light by a plane-parallel plate (see Sec. 7.1, subsection "General Theory of Filters").

Transmission and Reflection of Coatings on a Substrate

Many multilayer coatings are deposited onto a transparent or partially transparent substrate. Both the multilayer and the substrate contribute to the overall performance of the filter. For example, absorption in the substrate is frequently used to limit the transmission range of the filter. Reflectances at the filter interfaces need also to be considered. However, they can be reduced by anti-reflection coatings, or by cementing several components together.

In general, a filter can consist of multilayer coatings deposited onto one or both sides of a substrate. The overall transmittance T_{total} of a filter can be expressed in terms of the internal, or intrinsic transmittance τ of the substrate and the transmittances T_1, T_2 and internal reflectances R_1, R_2 of each surface of the substrate (Fig. 1).

The internal transmittance τ of a substrate is defined to be the ratio of the light intensity I'' just before reaching the second interface to the intensity I' just after entering the substrate (Fig. 1):

$$\tau = \frac{I''}{I'} \quad (4)$$

Expressions for the evaluation of the transmittance T and reflectance R of multilayer coatings are given in Sec. 7.2, subsection "Matrix Theory for the Analysis of Multilayer Systems."

Providing that the incident light is not coherent, there will be no interference between the beams reflected from the two surfaces of a substrate, even when the surfaces are plane parallel. A summation of all the partial reflections leads to the following expression for the overall spectral transmittance T_{total} of a filter:

$$T_{\text{total}} = \frac{T_1 T_2 \tau}{1 - R_1 R_2 \tau^2} \quad (5)$$

The reflection coefficients of uncoated interfaces can be calculated from Eq. (80), providing that the complex refractive indices of the substrate and of the medium are known. If all the materials in the filter are nonabsorbing, then

$$T_{\text{total}} = \frac{T_1 T_2}{1 - R_1 R_2} \quad (6)$$

If R_1 is small, an appropriate expression for T_{total} is

$$T_{\text{total}} \approx [1 - R_1(1 - R_2)]T_2 \quad (7)$$

However, this last approximation is not valid in general; some infrared substrate materials have high-reflection coefficients and in such cases Eq. (6) must be used.

Transmission Filters in Series and Parallel

To obtain a desired spectral transmittance it is frequently necessary to combine several filters. One common approach is to place several filters in series (Fig. 2a).

Because of the many different partial reflections that may take place between the various surfaces, precise formulas for the resulting transmittance are complicated.⁵ Accurate calculations are best carried out using matrix methods.⁶

To a first approximation, the resultant transmittance T' of a filter system consisting of k individual filters placed in series is given by

$$T' \approx T'_1 T'_2 T'_3 \cdots T'_k \quad (8)$$

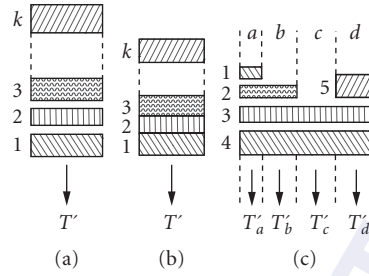


FIGURE 2 Transmission filters arranged in series and in parallel. The filters can be air-spaced (a), (c) or cemented (b).

Here T'_i is the total transmittance of filter i . This expression is valid only if the reflectances of the individual filters are small or if the interference filters are slightly inclined to one another and the optics are arranged in such a way that the detector sees only the direct beam.

Under other circumstances the use of this expression with interference filters can lead to serious errors. Consider two separate filters placed in series and let T'_1 , T'_2 , R'_1 and R'_2 correspond to the transmittances and reflectances of the two filters. If $T'_1 = T'_2 = R'_1 = R'_2 = 0.5$, then according to Eq. (8), the resulting transmittance will be $T' = 0.25$. For this simple case the precise expression can be derived from Eq. (5) and is given by

$$T' = \frac{T'_1 T'_2}{1 - R'_1 R'_2} \quad (9)$$

Evaluating this expression one obtains $T' = 0.33$. This is significantly different from the result obtained from the application of Eq. (8).

Some spectral transmittance curves cannot be easily designed by placing filters in series alone. For certain applications it is quite acceptable to place filters not only in series, but also in parallel.⁷ This introduces areas as additional design parameters. Thus, for example, the effective spectral transmittance T' of the filter shown in Fig. 2c would be given by

$$T' = \left(\frac{a}{A} T'_a + \frac{b}{A} T'_b + \frac{c}{A} T'_c + \frac{d}{A} T'_d \right) \quad (10)$$

where

A = overall area of filter,

a, b, c, d = areas of the four zones,

T'_a, T'_b, T'_c, T'_d = transmittances of four zones

The latter are given by

$$\begin{aligned} T'_a &= T'_1 \cdot T'_2 \cdot T'_3 \cdot T'_4 \\ T'_b &= T'_1 \cdot T'_2 \cdot T'_3 \\ T'_c &= T'_1 \cdot T'_2 \\ T'_d &= T'_1 \cdot T'_2 \cdot T'_3 \cdot T'_5 \end{aligned} \quad (11)$$

Great care must be exercised when using such filters. Because the spectral transmittance of each zone of the filter is different, errors will result unless the incident radiation illuminates the filter uniformly.

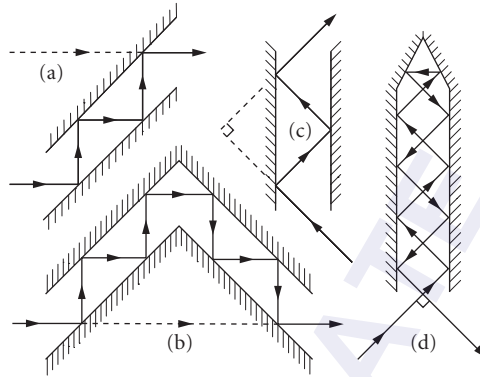


FIGURE 3 Various arrangements (a) to (d) for multiple-reflection filters.

Similar care must be used when employing the filtered radiation. One way proposed to alleviate these problems is to break the filter down into a large number of small, regular elements and to reassemble it in the form of a mosaic.⁸

Reflection Filters in Series

If radiation is reflected from k different filters, the resultant reflectance R' will be given by

$$R' = R'_1 \cdot R'_2 \cdot R'_3 \cdots R'_k \quad (12)$$

which is analogous to Eq. (8) for the resultant transmittance of filters placed in series. Many of the considerations of that section also apply here. For instance, R' will be significant only at those wavelengths at which every one of the reflectors has a significant reflectance. Metal layers (Sec. 7.15, subsection "Reducing Reflection with a Thin Metal Film") and thin-film interference coatings can be used exclusively or in combination.

For the sake of convenience the number of different reflecting surfaces used is normally restricted. The outlines of some possible reflector arrangements given in Fig. 3 are self-explanatory. The arrangement shown in Fig. 3b does not deviate or displace an incident parallel beam. The number of reflections depends in each case on the lengths of the plates and on the angle of incidence of the beam. Other arrangements are possible.

Clearly, reflection filters placed in series require more space and are more complicated to use than transmission filters. But if in a given application these shortcomings can be accepted, multiple-reflection filters offer great advantages, which stem mainly from the nature of the reflectors available for their construction (see Sec. 7.16, subsection "Multiple Reflection Filters").

7.2 THEORY AND DESIGN OF OPTICAL THIN-FILM COATINGS

Design Approaches

A thin-film designer may be asked to design a multilayer coating in which the transmittance, reflectance, and/or absorptance values are specified at a number of wavelengths, angles, and polarizations of the incident light. The designer may be required to provide a coating with many other more

complicated properties, including integral quantities such as CIE color coordinates, solar absorptance, or emissivity.⁹ The parameters that can be used to reach these goals are the number of layers in the multilayer, its overall thickness, the layer thicknesses and the refractive indices and extinction coefficients of the individual layers and of the surrounding media. Clearly, the more demanding the performance specifications, the more complex the resulting system. Many different methods have been developed for the design of multilayer coatings. For a good overview of this topic the interested reader is referred to the books by Macleod,² Knittl,¹⁰ and by Furman and Tikhonravov.¹¹ Here only the most important methods will be mentioned.

Graphical vector methods provide the most understanding of the problem, but the necessary approximations limit them to the solution of problems in which the final reflectance is not too high. Admittance diagrams and similar chart methods do not suffer from this limitation, but they are best applied to problems in which the specifications are relatively simple.¹² Many problems can be solved using the known properties of periodic multilayer systems.¹³ Analytical synthesis methods yield solutions to problems in which quite complex spectral transmittance or reflectance curves are specified.^{14,15} However, solutions frequently obtained in this way call for the use of inhomogeneous layers that are often more difficult to deposit, or for homogeneous layers with optical constants that are outside the range of known materials. Numerical design methods are the most flexible of all because they can be applied to problems with very complex specifications requiring a large number of layers for their solution.^{16,17} They are usually based on the matrix theory of optical multilayer systems and are particularly powerful for the solution of complicated spectral problems when combined with analytical methods.

Matrix Theory for the Analysis of Multilayer Systems

If electromagnetic radiation falls onto a structure consisting of thin films of several different materials, multiple reflections will take place within the structure. Depending on the light source and the layer thicknesses, the reflected beams may be coherent and interfere with one another. This optical interference can be used to design optical multilayer filters with widely varying spectral characteristics. In this section the basic equations for thin-film calculations are presented and some general properties of interference filters are listed. For a thorough discussion of this topic the reader is referred to the work of Macleod^{2,18} and Thelen.¹⁹

Consider the thin-film system consisting of L layers shown in Fig. 4. The construction parameters comprise not only the refractive indices n_j and the thicknesses d_j of the layers $j = 1, 2, \dots, L$, but also

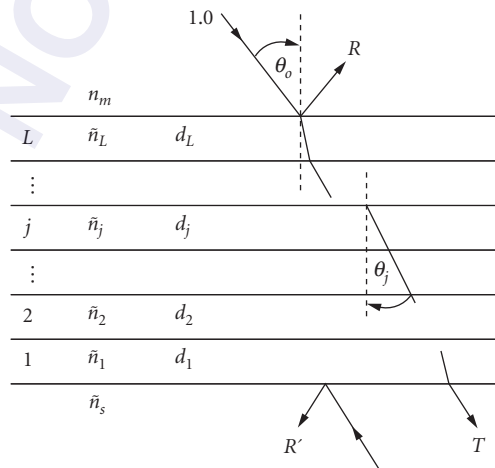


FIGURE 4 Construction parameters of a multilayer.

the refractive indices n_s and n_m of the substrate and the incident medium. The angle of incidence θ , the wavelength λ , and the plane of polarization of the incident radiation are the external variables of the system. It should perhaps be mentioned here that in some application areas it is more customary to use the grazing angle of incidence θ_g , where $\theta_g = 90 - \theta$.

The most general method of calculating the transmittance T and the reflectance R of a multilayer from the above quantities is based on a matrix formulation^{20,21} of the boundary conditions at the film surfaces derived from Maxwell's equations.²¹

It can be shown that the amplitude reflection r and transmission t coefficients of a multilayer coating consisting of L layers bounded by semi-infinite media, for the general case of obliquely incident light, are given by

$$r = \frac{\eta_m E_m - H_m}{\eta_m E_m + H_m} \quad (13)$$

and

$$t = \frac{2\eta_m}{\eta_m E_m + H_m} \quad (14)$$

where

$$\begin{pmatrix} E_m \\ H_m \end{pmatrix} = \mathbf{M} \begin{pmatrix} 1 \\ \eta_s \end{pmatrix} \quad (15)$$

E_m and H_m are the electric and magnetic vectors, respectively, in the incident medium, and \mathbf{M} is a product matrix given by

$$\mathbf{M} = \mathbf{M}_L \mathbf{M}_{L-1} \cdots \mathbf{M}_j \cdots \mathbf{M}_2 \mathbf{M}_1 \quad (16)$$

In the above equation \mathbf{M}_j is a 2×2 matrix that represents the j th film of the system:

$$\mathbf{M}_j = \begin{pmatrix} m_{11} & im_{12} \\ im_{21} & m_{22} \end{pmatrix} = \begin{pmatrix} \cos\delta_j & \frac{i}{\eta_j} \sin\delta_j \\ i\eta_j \sin\delta_j & \cos\delta_j \end{pmatrix} \quad (17)$$

where

$$\delta_j = \frac{2\pi}{\lambda} (n_j d_j \cos\theta_j) \quad (18)$$

the quantity $n_j d_j \cos\theta_j$ is the *effective optical thickness* of the layer j for an angle of refraction θ_j . In Eqs. (13) to (17) η represents the admittance of the medium, substrate, or layer and is given by

$$\eta = \begin{cases} \frac{n}{\cos\theta} & p\text{-polarization} \\ n\cos\theta & s\text{-polarization} \end{cases} \quad (19)$$

depending on whether the incident radiation is polarized parallel (p) or perpendicular (s) to the plane of incidence. Clearly, for normal incidence of light, the value of the admittance is equal to the refractive index. The angle θ_j is related to the angle of incidence θ_0 by Snell's law

$$n_m \sin\theta_0 = n_j \sin\theta_j \quad (20)$$

The intensity transmittance and reflectance are

$$T = \frac{\eta_s}{\eta_m} |t|^2 \quad (21)$$

$$R = |r|^2 \quad (22)$$

and the phase changes on transmission and reflection, ϕ_T and ϕ_R , are given by

$$\phi_T = \arg t \quad (23)$$

$$\phi_R = \arg r \quad (24)$$

More often than not, multilayer coatings are specified in terms of T and R and the phase changes on reflection and transmission are usually not of direct interest to the users. Exceptions are reflectors for interferometers.²² It has also been found that specifying certain values of ϕ_T , ϕ_R was useful in controlling the solution type obtained by some thin-film design methods.²³

If the materials in a multilayer are all nonabsorbing, then $T + R = 1$. Should one or more materials absorb, then in the above equations the refractive indices of these materials must be replaced by their complex refractive indices \tilde{n} , defined by

$$\tilde{n} = n - ik \quad (25)$$

where k is the extinction coefficient of the material. Even though all the elements of the layer matrix for such a material are now complex, its determinant will still be unity. The absorptance of the multilayer is then calculated from

$$A = 1 - T - R \quad (26)$$

Certain important general conclusions about the properties of multilayer filters can be drawn from the above equations.

1. The properties of thin-film systems vary with angle of incidence [Eqs. (18), and (19)]. For some applications this is the major disadvantage of interference filters compared to absorption filters.
2. This variation depends on the polarization of the incident radiation [Eq. (19)]. The following equations define T and R for obliquely incident nonpolarized radiation:

$$T = \frac{1}{2}(T_p + T_s) \quad (27)$$

$$R = \frac{1}{2}(R_p + R_s) \quad (28)$$

The dependence of T and R on polarization has been used for the design of polarizers (see Chap. 13, Vol. I). However, like the angular variation, it is a disadvantage for most other applications. Many researchers have investigated ways of reducing these effects.^{13,24-27}

3. Transmittance curves of nonabsorbing multilayers composed of layers whose optical thicknesses are all multiples of $\lambda/4$ show symmetry about λ_0 when plotted on a relative wavenumber scale λ_0/λ [Eqs. (17) and (18)].
4. A proportional change of all the thicknesses of a nonabsorbing multilayer results merely in a displacement of the transmittance curve on a wave-number scale [Eq. (18)]. Thus a thin-film design can be utilized in any part of the spectrum subject only to the limitations imposed by the dispersion of the optical constants of the materials used.
5. The reflectance and absorptance of a filter containing absorbing layers will depend, in general, on which side of the filter the radiation is incident (Fig. 5). However, the transmittance does not depend on the direction of the incident light.

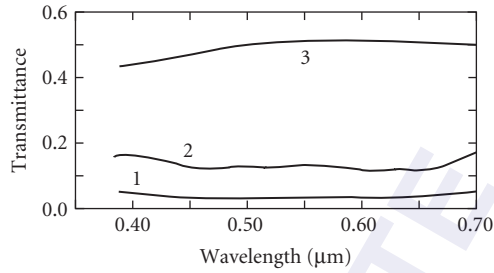


FIGURE 5 Spectral characteristics of a thin chromium alloy film on glass. Curve 1: transmittance; curves 2 and 3: reflectance from the glass and air sides, respectively. (After *Liberty Mirror*.²⁸)

Layers	n_m	Sublayers
		L_{sub}
L	n_L	L_{sub}^{-1}
		\vdots
\vdots		\vdots
		j_{sub}
j	n_j	\vdots
		\vdots
\vdots		\vdots
		\vdots
2	n_2	\vdots
		\vdots
		3
1	n_1	2
		1
	n_s	

FIGURE 6 Subdivision of layers for the evaluation of the electric field within a multilayer (see Sec. 7.2, subsection “Matrix Theory for the Analysis of Multilayer Systems”).

An important parameter is the electric field amplitude squared, $|E|^2$. The susceptibility of a multilayer to high-power laser damage is proportional to the highest value of this quantity within the multilayer. One way to evaluate E is to subdivide each layer of the system into L_{sub} sublayers (Fig. 6) and to evaluate the partial matrix products P_j ,

$$\mathbf{P}_j = \begin{pmatrix} p_{11} & ip_{12} \\ ip_{21} & p_{22} \end{pmatrix} = \prod_{k=1}^j (\mathbf{S}_k) \quad (29)$$

for all L sublayers, where S_k is the matrix of the k th sublayer. Then, the electric field amplitude at a point j in the multilayer is given in terms of the elements of the total and partial product matrices by the following expression:²⁹

$$E_j = \frac{4[p_{11}^2 + (n_s p_{12})^2]}{[m_{11} + (n_s m_{22}/n_m)]^2 + [(m_{21}/n_m) + n_s m_{12}]^2} \quad (30)$$

The analysis of optical thin-film systems [with computer programs based on Eqs. (13) to (30)] is relatively simple. The design of filters with any but the simplest spectral characteristics is a more complicated problem and one of the methods listed in Sec. 7.2, subsection "Design Approaches" must be used.

7.3 THIN-FILM MANUFACTURING CONSIDERATIONS

The optical, mechanical, and environmental properties of multilayer coatings depend on the materials used, on the deposition process and on the surface quality of the substrate.

There are many different methods for the deposition of thin films.³⁰⁻³² Some of the more common processes are reviewed below. The deposition methods and process parameters used affect the microstructure of the resulting layers. The films can be dense with an amorphous or a microcrystalline structure, or they may exhibit a columnar growth with considerable voids. The optical constants clearly depend on this microstructure and films of the same materials may sometimes have very different properties, depending on how they were deposited (see, for example, Ref. 33). The individual films in a multilayer may be under tensile or compressive stress and, unless materials and film thicknesses are selected to compensate for these stresses, the overall stress may be large enough to distort the substrate or cause the multilayer to break up. The mechanical properties of multilayer coatings also critically depend on the microstructure of the films. An excellent discussion of the effects of microstructure on the various properties of optical coatings has been given by Macleod.¹⁸

Optical Coating Materials

Many different materials have been used in the past for the construction of optical multilayer coatings. Some of the compounds used for the deposition of nonabsorbing layers for the ultraviolet, visible, and infrared parts of the spectrum are cryolite (1.35), LiF (1.37), MgF₂ (1.39), ThF₄ (1.52), CeF₃ (1.62), PbF₂ (1.73), ZnS (2.30), ZnSe (2.55), Si (3.5), Ge (4.20), Te (4.80), PbTe (5.50), and the oxides SiO₂ (1.48), Al₂O₃ (1.60), MgO (1.72), Y₂O₃ (1.82), Sc₂O₃ (1.86), SiO (1.95), HfO₂ (1.98), ZrO₂ (2.10), CeO₂ (2.20), Nb₂O₅ (2.20), Ta₂O₅ (2.10), and TiO₂ (2.45). The numbers in parenthesis represent the approximate refractive indices at the midpoints of the material transparency range. Some of the metals used in the same wavelength range for the deposition of reflecting or absorbing layers are Ag(0.12 - i3.45), Al(1.02 - i6.85), Au(0.31 - i2.88), Cu(0.83 - i2.60), Ni(1.80 - i3.33), Cr(3.18 - i4.41), Inconel(2.94 - i2.92), and Rh(2.00 - i5.11). The complex refractive indices of the metals given in the brackets correspond to a wavelength of 0.56 μm. The above values are approximate and are intended as a rough guide only. More extensive listings of coating materials are given, for example, by Macleod² and Costich.³⁴

As already mentioned, the properties of multilayers depend on the materials used for their construction. For example, layers made of oxides are, as a rule, harder than those made of fluorides, sulfides, or semiconductors. They are therefore preferred for use on exposed surfaces. Semiconductor materials should be avoided in filters that are to be used over a wide range of temperatures because their optical constants can change significantly. Some metals are soft and easily damaged while others tarnish when exposed to the atmosphere. Such coatings require further protective coatings, or should be cemented between two transparent plates. Other materials require the precoating with adhesion layers to ensure a good bond to the substrate. For example, frequently a Ni adhesive layer is deposited onto glass before coating with Au.

Evaporation

Conventional (nonreactive) or reactive evaporation from resistance, induction or electron beam gun sources is a low energy process (~ 0.1 eV) and the resulting films frequently have a porous structure. The porosity may vary with the material, the substrate temperature, the residual pressure in the deposition chamber, the deposition rate and angle of incidence of the vapour on the substrate. Values of porosities ranging from 0 to 40 percent have been observed. On exposure to the atmosphere some of the voids in the film may adsorb water vapor. This increases the effective refractive index of the films and results in a shift of the spectral features of the multilayer towards longer wavelengths (“ageing”). This shift is partially reversible—by placing the filter in an inert atmosphere or in a vacuum, or on heating, some of the adsorbed water vapor can be removed. Unless it has been allowed for at the design stage, such ageing can render some filters useless.

The microstructure of the films can be significantly affected by bombarding the substrate during deposition with energetic ions from an auxiliary ion beam source.^{35–37} The additional energy (~ 50 to 100 eV) results in denser films. Hence, coatings produced by *ion-assisted deposition* have higher refractive indices and exhibit less or no ageing on exposure to the atmosphere.

The *ion-plating* process can result in even denser coatings.^{38–40} In this high deposition rate process the starting material must be a good conductor and is usually a metal. Argon and a reactive gas species are introduced into the chamber and, together with the evaporant, are ionized. The ions are then accelerated to the substrate with energies of the order of 10 to 50 eV. Transparent films with near-bulklike densities and low temperature variation of refractive index can be obtained by this process. For most materials the layers are glasslike and the interfaces remain smooth. This results in a lower scatter. *Plasma-ion-assisted deposition* is a process that also produces such good quality coatings, but it has the advantage over the ion-plating process that it does not require a conducting evaporant, so that fluorides and other nonconducting materials can be deposited with it.^{41–43}

Conventionally evaporated thin films can be under compressive or tensile stresses. If not controlled, these stresses can distort the substrate or cause the multilayer to break up. The magnitudes of the stresses depend on the material and on the deposition conditions. It is often possible to select the materials and process parameters so that the stresses of the various layers counteract each other. In contrast, almost all ion-plated layers are under compressive stress. It is therefore more difficult to produce stress-compensated multilayers by this process.

Sputtering

Reactive or nonreactive DC, RF, AC, or pulsed magnetron sputtering are also used to deposit optical multilayer coatings. Very many variants of this process exist.⁴⁴ Most are significantly slower than evaporation and the targets can be quite expensive. Filters produced by magnetron sputtering may therefore also be more expensive. However, the process is stable, provides excellent control over the thicknesses of the layers and can be readily scaled to provide uniform coatings over large areas. Both metal and metal oxide layers can be produced. Sputtering is an energetic process and results in dense, bulklike layers that exhibit virtually no ageing.

In *ion-beam sputtering* an energetic beam of inert ions is aimed at a target made of the material that is to be deposited. Atoms or clusters of atoms of the material are dislodged from the target and land on the substrate with a high energy (about 10 to 200 eV). This is the slowest physical vapor thin-film deposition method described here and it cannot yet be readily scaled for the coating of large components. However, it yields the highest quality coatings. Many of the high reflectance coatings for laser gyroscopes, in which no significant losses can be tolerated, are produced in this way.

Deposition from Solutions

In this procedure the substrate is either dipped in an organo-metallic solution and withdrawn at a very steady rate from it, or the solution is applied from a pipette onto a spinning substrate. The substrate is then placed in an oven to drive off the solvent. The thickness of the film depends on the

concentration of the solvent and on the rate of withdrawal or spinning. Other factors that influence the process are temperature and humidity, as well as the freshness of the solution. Although it yields quite porous films, this method is of interest because many of the layers produced in this way have a high laser damage threshold. The process has also been adapted for the coating of quite large area substrates with multilayer antireflection coatings for picture frame glass and for display windows.

Other Deposition Methods

Great strides have been made in the development of various *chemical vapor deposition methods* for the production of optical coatings.⁴⁵ These processes often take place in the presence of plasmas operating at microwave and/or radio frequencies. The advantage of this approach is that it is relatively easy to produce robust layers with intermediate optical constants through the use of appropriate mixtures of two or more reactant gases. However, most of the reactants used are toxic and their handling requires special care.

Atomic layer deposition is inherently a slow deposition process that is capable of depositing very precise, conformal layers onto irregular-shaped substrates. Recently it has been demonstrated that it can be a commercially viable process if large chambers that contain stacks of closely packed substrate holders are used for the simultaneous coating of hundreds of objects.⁴⁶

Thickness Control during Deposition

The performance of many optical multilayers depends critically on the thicknesses of the individual layers. The control of the layer thicknesses during their formation is therefore very important. Many different methods exist for the monitoring of layer thicknesses. For very steady deposition processes, such as sputtering, simple timing can give good results. However, the most common techniques used are quartz crystal and optical monitoring. The former is very sensitive and can be used for thin and thick films, as well as transparent and opaque films. However, it is an indirect method and requires careful calibration. This is usually not a problem whenever layers of established coating materials are formed using a standard geometry and deposition conditions. Optical monitoring can be performed directly on the substrate, or indirectly on a witness glass. The quantities measured are usually T , R , or the ellipsometric parameters. One advantage of direct optical monitoring is that the parameters measured are usually closely related to the required performance. Furthermore, with optical monitoring, a real-time error determination and compensation is possible after the deposition of each layer, through the reoptimization of the remaining layers of the system.⁴⁷ With this method even quite complicated multilayer structures can be manufactured. This process is now used in many laboratories. Even better results should be possible in the future when the above process is combined with ion-beam etching to adjust the thicknesses of layers that exceed the intended value.⁴⁸

7.4 MEASUREMENTS ON OPTICAL COATINGS

Optical Properties

Transmission, Reflection, and Absorption The most commonly used instrument for the measurement of the optical performance of thin-film coatings is the spectrophotometer. The wavelength dispersion of commercial instruments for the 0.185 to 80.0 μm spectral range is usually provided by prisms or by ruled or holographic gratings. Grazing incidence gratings, crystals, or multilayer coatings are used in the soft x-ray and extreme ultraviolet (XUV) spectral regions. Fourier transform spectrometers are capable of measurements from about 2 to 500 μm . A variety of attachments are available for the measurement of specular and diffuse reflectance. Absolute measurements of T and R that are accurate to within ± 0.1 percent are difficult to make even in the visible part of the electromagnetic spectrum. Measurements at oblique angles of incidence are even more difficult.

Very small absorptions of single layers are normally measured with calorimeters.⁴⁹ The losses of high performance laser reflectors are obtained from measurements of the decay times of Fabry-Perot interferometer cavities formed from these mirrors.⁵⁰

Roughness of the substrate and irregularities occurring within individual films and the layer interfaces give rise to light scattering in all directions.⁵¹⁻⁵³ For many applications it is important to minimize this scatter. Special instruments, called scatterometers, are used to measure the angular variation of the light scatter. Such data provides information about the substrate and multilayer.⁵⁴

The transmittance, reflectance, and absorptance of some optical coatings are affected by exposure to atomic oxygen and by electron, proton, and ultraviolet irradiation. They also depend critically on the cleanliness of the components measured. In space contamination of optical components can also take place.^{55,56}

Optical Constants A reliable knowledge of the optical constants of all the materials used in the construction of optical multilayer coatings is essential. There exist many different methods for their determination.⁵⁷ These include methods that are based on refractometry, photometric, and spectrophotometric measurements of R and/or T ; polarimetry, singlewavelength or spectroscopic ellipsometry, various interferometric methods, attenuated total reflection, or on a combination of two or more of the above methods. Excellent monographs on the various methods will be found in Palik's *Handbook of Optical Constants of Solids*.^{58,59} Some are suitable for measurements on bulk materials and the results are valid only for films produced by the more energetic deposition processes described above. The optical constants of porous films must be measured directly. They will depend on the deposition parameters and on the layer thickness, and may differ significantly from those of bulk materials.⁶⁰ Special methods have to be used for the determination of the optical constants in the x-ray, XUV, and sub-millimeter regions. The accurate measurement of very small, residual extinction coefficients of transparent coating materials is difficult. Generally it involves the use of laser calorimetry or the use of the film as a spacer layer in a bandpass filter. It is also very important to be able to measure the thickness of the film independently.

Laser Damage A measure of the ability of an optical component to withstand high laser irradiations is the laser damage threshold (LDT). There are several ways of defining this quantity. One frequently used definition is based on a plot of the percentage of components that are damaged when they are exposed to different laser fluences. The value of the fluence corresponding to the intersection of a mean curve through the experimental points with the ordinate is defined to be the LDT. It is thus the maximum fluence at which no damage is expected to the component.

Absorption is the main cause for laser damage. The incident radiation that is absorbed in the optical component will be converted into heat. If the thermal conductance of the optical component is too low, the temperature of the local hot spot on the mirror will rise to a value at which damage occurs. The damage will therefore depend on the thermal conduction of the materials of which the mirror blank is made. For example, some high reflectivity mirrors consist of coatings on Si or Cu substrates that are water-cooled during use.

Thermal conductivities of thin films are several orders of magnitude smaller than those of the corresponding bulk materials. (Exceptions to this are some fluoride layers.) This compounds the problem. To increase the LDT, the deposition methods are optimized to obtain thin films with more bulklike thermal conductivities.

Absorption damage usually initiates at defects and other imperfections. In the case of the substrate, it may occur at or below the surface, even when the substrate material is nonabsorbing. It is very important to avoid materials with color centers and subsurface damage and inclusions. The substrate surface must be very smooth and devoid of scratches, digs, and pores, otherwise polishing compounds and other contaminants can be trapped. It is imperative that the substrate be perfectly clean prior to coating. Electric fields associated with the above imperfections increase the absorption by an amount that is proportional to the refractive index of the material.

The coating materials used for the construction of high LDT multilayers must be very pure, with absorption edges far away from the wavelength of interest. As a rule, materials in thin-film form have

extinction coefficients that are orders of magnitude larger than those of the corresponding bulk materials. Currently the processes used to produce high LDT coatings include ion-beam sputtering, sol-gel deposition, and electron beam gun evaporation. To reduce the effects of the residual surface roughness, the thicknesses of the layers are often adjusted to shift the peaks of the electric field away from the layer boundaries.

The form that the damage takes depends to a large extent on the materials. Pitting of the coatings is probably due to the evaporation of the thin-film materials. Delamination may be due to poor adhesion of the layers to the substrate, to undue stresses in the films and/or to a poor match between the expansion coefficients of the layers and the substrate.

The LDT also depends on the laser pulse duration. For very short pulses (<10 ns) thermal conductance does not play a role in the process. For higher values, the LDT is proportional to the square root of the pulse duration. For high repetition rates the LDT depends on the repetition rate. To achieve a long life (10,000 or 100,000 pulses) in an industrial environment, the laser should be operated at a fraction (say, 1/4 or 1/10) of the nominal LDT. In CW lasers it is not the LDT, but the power handling capability that is of essence. Long before damage takes place, the heating can cause a distortion of the surface that, in turn, can result in loss of power and in mode and focussing problems.

The development of high LDT of laser coatings is so important, that conferences have been held in Boulder, Colorado, on this topic every year since 1970. For more detail the reader is referred to the proceedings of these conferences,⁶¹ as well as to an article written by Ristau.⁶² There exists also a draft international standard on this topic. An increasing number of thin-film vendors include in their catalogues LDT information on all or some of their products. Independent LDT test services are provided by several commercial companies and publically funded institutions.

Mechanical Properties

Optical multilayer coatings are frequently required to operate under severe mechanical and environmental conditions. Quite frequently meeting these requirements is more difficult than to achieve the necessary optical performance. A number of standards deal with the substrate and coating quality (MIL-0-13830B), the adhesion of coatings (MIL-M-13508C, MIL C 48497), their abrasion resistance (MIL-C-675A, MIL C 675C, MIL-C-14806-A, MIL C 48497), hardness (MIL-M-13508C), and resistance to humidity (MIL-C-675A, MIL-C-14806-A, MIL-810-C, MIL C 48497), salt solution (MIL-C-675-A) and salt spray (MIL-M-13508-C, MIL-C-14806-A). Most of these standards are reprinted in Rancourt's book.³ Depending on the application, multilayer coatings may be required to meet one or more of the above standards. An overview of the subject of stresses and hardness of thin films on a substrate has been recently published by Baker and Nix.⁶³ Sapieha and her coworkers compared various mechanical properties of three commonly used oxide layers produced by different deposition processes.⁶⁴

Analytical Analysis Methods

In addition to optical and mechanical measurements, multilayer coatings can be subjected to a number of analytical measurements. These include Auger electron spectroscopy, energy dispersive x-ray analysis, Rutherford backscattering, secondary ion mass spectrometry, transmission electron microscopy and x-ray photoelectron spectroscopy.^{63,65} Some of these tests are destructive. However, when a multilayer coating is subjected to these tests, they yield fairly accurate information about the number of layers in a system, and on the composition, thickness, and structure of the individual layers.⁶⁶

7.5 ANTIREFLECTION COATINGS

Effect of Surface Reflections on Performance of Optical Systems

The reflectance of an interface between two nonabsorbing media of refractive indices n_1 and n_2 is given by

$$R = \left[\frac{n_1 - n_2}{n_1 + n_2} \right]^2 \quad (31)$$

An expression for the total transmittance T_0 of a nonabsorbing plane-parallel plate that takes into account the effect of multiple internal reflections within the plate can be obtained from Eq. (9):

$$T_0 = \frac{1-R}{1+R} \quad (32)$$

Of this light only a fraction $(1-R)^2$ passes through the plate without undergoing any reflections.

An expression for the transmission of a number of such plates placed in series is of interest. It helps to explain the effect of multiple reflections between the various plates on the performance of devices such as tripple glazings and photographic objectives. It can be shown⁶ that the total transmittance T_{total} of m plates placed in series is given by

$$T_{\text{total}} = \frac{T_0}{m - T_0(m-1)} \quad (33)$$

The amount of light transmitted directly T_{direct} , is

$$T_{\text{direct}} = (1-R)^{2m} \quad (34)$$

The light T_{stray} that undergoes multiple reflections before transmission is responsible for spurious images and stray light in the image plane. It is given by

$$T_{\text{stray}} = T_{\text{total}} - T_{\text{direct}} \quad (35)$$

The variation of T_{direct} and T_{stray} with R , for several values of m , is shown in Fig. 7. The refractive indices of the plate material plotted against the upper x axis assume that the plates are in air.

It will be seen that even for a relatively small number of low-refractive-index plates the ratio $T_{\text{stray}}/T_{\text{direct}}$ becomes significant. This means that in an image-forming system under unfavorable conditions the stray light can completely obscure the image.⁶⁷ Second, even in nonimaging optical systems the loss of light $(1 - T_{\text{direct}} - T_{\text{stray}})$ can become quite prohibitive.

Both these problems can be overcome by reducing the surface reflection through the application of suitable antireflection coatings to the plate boundaries. Since antireflection coatings with zero reflectance across the whole spectrum cannot be constructed, the spectral reflectance $R(\lambda)$ of antireflection coatings is usually chosen to minimise the integral

$$\int_{\lambda} R(\lambda) I(\lambda) S(\lambda) d\lambda \quad (36)$$

where $I(\lambda)$ and $S(\lambda)$ are the spectral-intensity distribution of the incident radiation and the spectral sensitivity of the detector, respectively. A low reflectance is thus needed only in the spectral region in

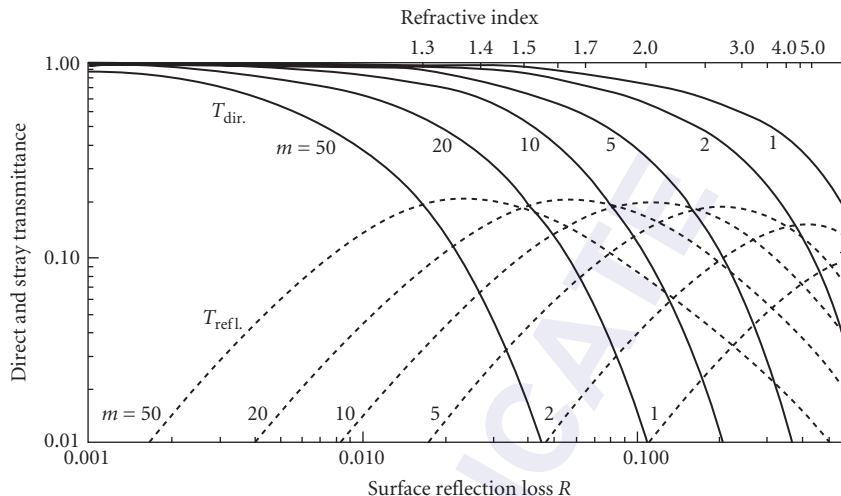


FIGURE 7 Transmittance of a system of m parallel plates of surface reflectances R for directly transmitted radiation T_{direct} and for radiation that suffers multiple reflections before transmission T_{stray} . The upper x axis is calibrated in terms of the refractive index of the plates for the special case when the plates are in air.

which $I(\lambda)S(\lambda)$ is significant. It should be emphasized here that in addition to this requirement, for most applications antireflection coatings must be mechanically very tough, withstand drastic climatic and thermal variations, and stand up to the usual lens-cleaning procedures. Some examples of improvements in the performance of image-forming and non-image-forming optical systems obtained through the use of antireflection coatings are given by Mussett and Thelen⁶⁷ and by Faber et al.⁶⁸

Antireflection coatings can be based on homogeneous layers or on inhomogeneous coatings. One can further classify them into single layer, digital, or structured, and homogeneous multilayer or complex inhomogeneous layer coatings (Fig. 8a to f). Because of their industrial importance, antireflection coatings for the visible and infrared spectral regions have been the subject of much research and development. Two books have been written on this topic^{69,70} and there exists a very extensive literature in scientific and technical journals. For a review of this literature and for a systematic discussion of antireflection coatings, the reader is referred to the excellent review articles by Cox and Hass⁷¹ and by Mussett and Thelen.⁶⁷ A recent report on the present state of the art has been published in the Japanese language. In this section, only a brief summary will be given of the results obtained thus far, intended to aid in the selection of antireflection coatings for particular applications. The calculated data is presented on a logarithmic scale. The relative wavenumber scale facilitates the calculations of the width of the effective region of a coating in different parts of the spectrum.

Antireflection Coatings Made of Homogeneous Layers

The single homogeneous-layer antireflection coating (Fig. 8a) was the first antireflection coating and perhaps is still the most widely used. Theoretically it should be possible to obtain a zero reflectance at one wavelength with single dense films. However, because of a lack of suitable low-index coating materials, this cannot be realized in practice for substrates with indexes less than about 1.9. Nevertheless, even with the available materials a very useful reduction in reflection in a broad spectral region is obtained for all common glass types, the reflectance never rising above that of the uncoated surface (Fig. 9, curve 8). With the sol-gel method it is possible to produce homogeneous porous oxide and fluoride films with very low refractive indices.⁷²⁻⁷⁵ Such films have excellent optical characteristics

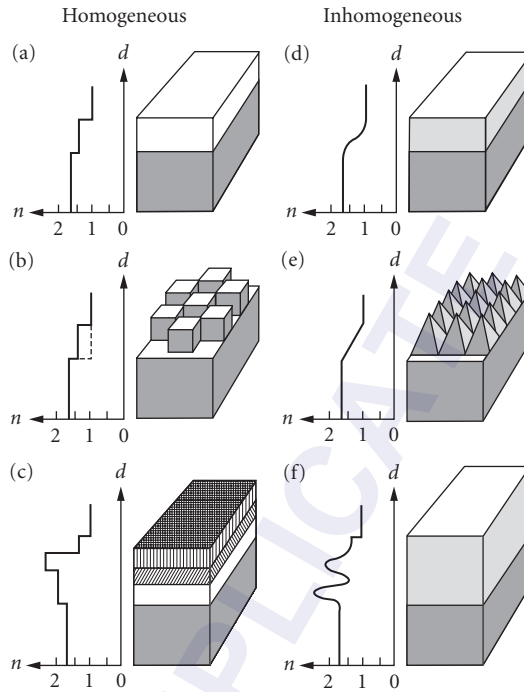


FIGURE 8 Structure and effective refractive index profiles of various types of antireflection coatings. (a), (b), (c) homogeneous single layer, digital and multilayer AR coatings and (d), (e), (f) simple inhomogeneous layer, structured and more complex inhomogeneous layer AR coatings.

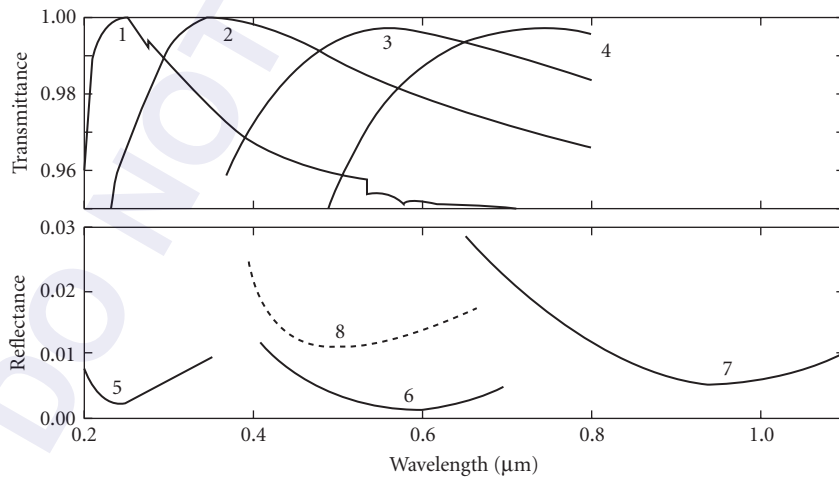


FIGURE 9 Performance of single layer porous homogeneous antireflection coatings on various substrate surfaces. Porous silica on fused silica (1, 3, 5); on glass (6); on KDP (7); SF_8 (4). Porous MgF_2 on fused silica (2). A conventional single layer antireflection coating on glass is shown for comparison (8). (Curve 1 after Wilder;⁷² curve 2 after Thomas;⁷⁴ curves 3 and 4 after Thomas;⁷⁵ curves 5–7 after O'Neill;⁷³ and curve 9 after Balzers.⁷⁶)

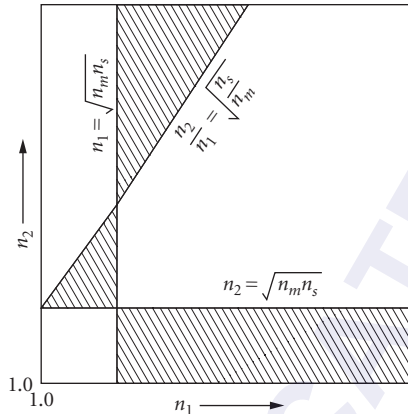


FIGURE 10 Refractive index combinations (*shaded areas*) of two-layer antireflection coatings with which a zero reflectance at one wavelength can be achieved.

(Fig. 9) and have laser damage thresholds that are considerably higher than those produced by conventional means. However, these gains are at the expense of mechanical strength and long-term stability. Low-refractive index coating materials can also be simulated by the deposition or etching of subwavelength structures (Fig. 8*b, e*). The effective refractive index depends on the volume fraction occupied by the structures.⁷⁷

If more than one layer is used, all the degrees of freedom could be used to either (1) obtain a more complete antireflection in one particular spectral region; (2) increase the width of the spectral region over which the reflectance is generally low; or (3) obtain a coating in which the low reflectance is very uniform across the spectrum.^{78,79} Vlasov has shown with the aid of a diagram of the type shown in Fig. 10 that even for a two-layer antireflection coating there exists a large number of refractive index combinations which will yield zero reflectance at one wavelength.⁶⁹ As the number of layers and the overall thickness of the antireflection coating increases, it becomes possible to find solutions for a particular problem that not only fully meet the most important above desiderata and almost satisfy the others but are also based on the use of the mechanically most satisfactory coating materials.

The conditions that are satisfied by the refractive indexes and thicknesses of various types of antireflection coatings are given in Table 1. In a few cases where they are very complicated, reference is made to a paper in which they are set out in full. The calculated transmittance curves of antireflection-coated surfaces of glass (Fig. 11*a to c*), quartz (Fig. 11*d*), germanium (Fig. 11*e and f*), silicon (Fig. 11*g*), and other infrared materials (Fig. 11*h*) utilize refractive indexes that for the most part correspond to real coating materials, and hence the curves represent realistic solutions rather than the theoretically best possible ones. The actual refractive indices used in the calculations (and optical thicknesses where they do not correspond to a multiple of $\lambda/4$) are given in Table 1.

Figure 11*a* shows the performance of several antireflection coatings on glass for applications which require the highest possible efficiency for a limited wavelength region. Of these the coating 2.1 has probably found the greatest acceptance. A typical measured performance curve for such a coating is shown in Fig. 12. Antireflection coatings with a low reflectance in a broad spectral region are shown in Fig. 11*b* and Fig. 13. Coatings of the type 3.4 probably find the widest application in practice. The performance of commercial coatings for quartz substrates are shown in Fig. 14.

Solutions listed in Table 1 presume that coating materials with the required refractive indices exist. This is frequently not the case and the thin-film designer must seek solutions based on available coating materials. For example, Furman⁸⁹ gives a series of practical two-, three-, and four-layer solutions for 10 different values of the substrate index ranging between 1.5 and 4.0, while Stolov provides⁹⁰ seven-layer solutions for 10 values of n_s , $1.46 \leq n_s \leq 1.82$.

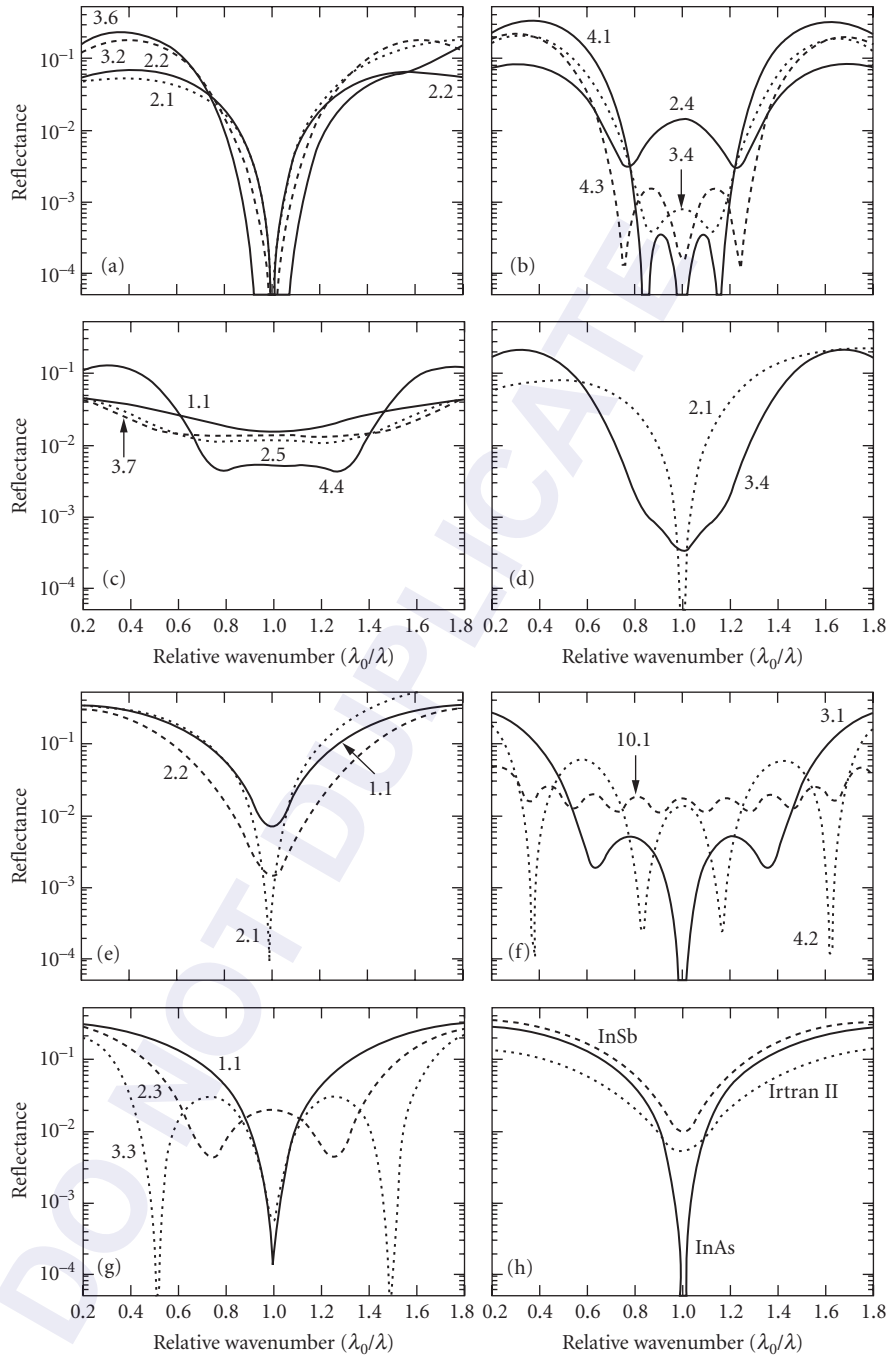


FIGURE 11 Calculated performance of various antireflection coatings. The numbers identifying the individual curves refer to Table 1. (a) High-efficiency antireflection coatings for glass; (b) broadband antireflection coatings for glass; (c) highly achromatic antireflection coatings for glass; and (d) antireflection coatings for quartz. (e) Antireflection coatings for germanium; and (f) broadband antireflection coatings for germanium. (g) Antireflection coatings for silicon and (h) single-layer antireflection coatings for Irtran II, InAs, and InSb.

TABLE 1 Some Antireflection Coatings^{80,81} (*The incident medium in all cases is air.*)

Type	Conditions or Reference	Substrate Material	n_s	n_1 ($n_1 d_1 / \lambda$)	n_2 ($n_2 d_2 / \lambda$)	n_3 ($n_3 d_3 / \lambda$)	n_4 ($n_4 d_4 / \lambda$)	n_m
1.1	$n_1 = \sqrt{n_s n_m}; n_1 d_1 = \frac{\lambda}{4}$	Glass	1.51	1.38				1.00
		Irtran II	2.20	1.59				1.00
		InSb	4.00	2.20				1.00
		InAs	3.40	1.85				1.00
		Ge	4.10	2.20				1.00
		Si	3.50	1.85				1.00
2.1	$\tan^2 \delta_1 = \frac{(n_1^2(n_m - n_s)(n_s n_m - n_1^2))}{(n_s n_2^2 - n_m n_1^2)(n_s n_m - n_1^2)}$	Glass	1.51	2.30 (0.0524)	1.38 (0.3250)			1.00
		Quartz	1.48	2.09 (0.0947)	1.48 (0.3255)			1.00
		Ge	4.10	1.35 (0.0951)	4.10 (0.0586)			1.00
2.2	$n_2^2 n_s = n_1^2 n_m; n_1 d_1 = n_2 d_2 = \frac{\lambda}{4}$	Glass	1.51	1.70	1.38			1.00
		Ge	4.10	3.30	1.57			1.00
2.3	$n_1 n_2 = n_m n_s; n_1 d_1 = n_2 d_2 = \frac{\lambda}{4}$	Si	3.5	2.20	1.35			1.00
2.4	$n_1^2 - \frac{n_1 n_s}{2n_2 n_m}(n_m^2 + n_2^2)(n_1 + n_2) + n_2 n_s^2 = 0$ $\frac{1}{2} n_1 d_1 = n_2 d_2 = \frac{\lambda}{4}$	Glass	1.51	1.70	1.38			1.00
2.5	$n_1 d_1 = n_2 d_2 = \frac{\lambda}{4}$ (see Kard, Ref. 80)	Glass	1.55	1.484	1.32			1.00
3.1	$n_3 n_s = n_2^2 = n_m n_s;$ $n_1 d_1 = n_2 d_2 = n_3 d_3 = \frac{\lambda}{4}$	Ge	4.1	3.30	2.20	1.35		1.00
3.2	$n_1 n_3 = n_2 \sqrt{n_m n_s};$ $n_1 d_1 = n_2 d_2 = n_3 d_3 = \frac{\lambda}{4}$	Glass	1.53	1.80	2.14	1.47		1.00
3.3	$n_2^2 = n_m n_s, n_m n_1^2 = n_2^2 n_s$ $n_1 d_1 = n_2 d_2 = n_3 d_3 = \frac{\lambda}{4}$	Si	3.45	2.56	1.86	1.38		1.00
3.4	$n_2^2 n_s = n_m n_1^2;$ $n_1 d_1 = \frac{1}{2} n_2 d_2 = n_3 d_3 = \frac{\lambda}{4}$	Glass	1.51	1.65	2.10	1.38		1.00
		Quartz	1.48	1.65	2.10	1.38		1.00
3.5	$n_2^2 n_s = n_m n_1^2,$ $\frac{1}{3} n_1 d_1 = \frac{1}{2} n_2 d_2 = n_3 d_3 = \frac{\lambda}{4}$	Glass	1.51	1.659	2.20	1.38		1.00

(Continued)

TABLE 1 Some Antireflection Coatings^{80,81} (*The incident medium in all cases is air.*) (Continued)

Type	Conditions or Reference	Substrate Material	n_s	n_1 ($n_1 d_1 / \lambda$)	n_2 ($n_2 d_2 / \lambda$)	n_3 ($n_3 d_3 / \lambda$)	n_4 ($n_4 d_4 / \lambda$)	n_m
3.6	(see Thetford Ref. 81)	Glass	1.52	1.80	2.20	1.38		1.00
3.7	$n_1 d_1 = n_2 d_2 = n_3 d_3 = \frac{\lambda}{4}$ (see Kard, Ref. 80)	Glass	1.55	1.53	1.454	1.32		1.00
4.1	$n_1 n_4 = n_2 \sqrt{n_m n_s}$; $n_1 d_1 = n_2 d_2 = \frac{1}{2} n_3 d_3 = n_4 d_4 = \frac{\lambda}{4}$	Glass	1.51	1.38	1.548	2.35	1.38	1.00
4.2	$n_1 n_4 = n_2 n_3 = n_m n_s$; $n_1 d_1 = n_2 d_2 = n_3 d_3 = n_4 d_4 = \frac{\lambda}{4}$	Ge	4.0	2.96	2.20	1.82	1.38	1.00
4.3	$n_1 d_1 = n_2 d_2 = n_3 d_3 = n_4 d_4 = \frac{\lambda}{4}$ (see Kard, Ref. 80)	Glass	1.55	1.846	2.289	2.014	1.32	1.00
4.4	(Design derived from Kard, Ref. 82.)	Glass	1.55	1.656 (0.2417)	1.888 (0.2463)	1.832 (0.2390)	1.38 (0.2424)	1.00
10.1	$n_l d_s = n_s - (11-l) \frac{n_s - n_9}{10}$, $l = 1, 2, \dots, 10$ $n_1 d_1 = n_2 d_2 = \dots = n_{10} d_{10} = \frac{\lambda}{4}$ [see Kard, Ref. 80]	Ge	4.00	$n_1 = 3.735$...		$n_9 = 3.735$	$n_{10} = 1.35$	1.00

Wiley has given a useful empirical expression which relates the average reflectance R_{avg} in the visible and near-infrared spectral region $\lambda_{\text{min}} \leq \lambda \leq \lambda_{\text{max}}$ to the overall optical thickness $\Sigma(nd)$ of an antireflection coating composed of layers of refractive indices n_M, n_H and a single, outermost layer of index n_L :⁹¹

$$R_{\text{avg}} = 0.01 \left[\frac{\lambda_{\text{max}}}{\lambda_{\text{min}}} (n_L - 1) \right]^{3.4} \left(\frac{\lambda_{\text{max}}}{\Sigma(nd)} \right)^{0.63} [(1.2 - \Delta n)^2 + 0.42] \quad (37)$$

The above expression is valid for the following parameter values:

$$\begin{aligned} 1.5 \leq \frac{\lambda_{\text{max}}}{\lambda_{\text{min}}} \leq 3.0 \quad 1.17 \leq n_L \leq 1.46 \quad 1.38 \leq n_M, n_H \leq 2.58 \\ 1.0 \leq \frac{\Sigma(nd)}{\lambda_{\text{max}}} \leq 3.0 \quad 0.4 \leq \Delta n = n_H - n_M \leq 1.2 \end{aligned} \quad (38)$$

A method for the design of optimum or near-optimum two-material antireflection coatings for a given substrate, coating materials, and overall thickness has recently been described.⁹²

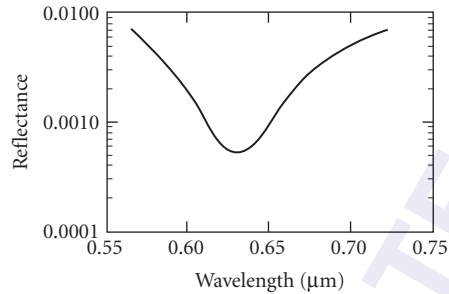


FIGURE 12 Reflectance of a high-efficiency antireflection coating on glass. (After Costich.¹⁷²)

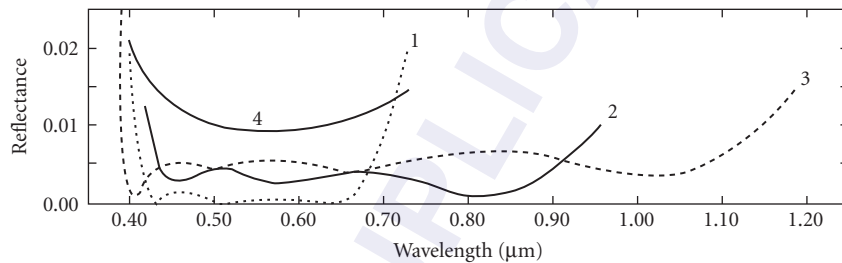


FIGURE 13 Reflectances of three broadband antireflection coatings. (Curve 1 after Turner;¹⁷³ curve 2 after Optical Coating Laboratory;⁸⁴ curve 3 after Thin Film Lab;⁸⁵ and curve 4 after Balzers⁷⁶ corresponds to a single-layer AR coating on glass and is shown for comparison purposes.)

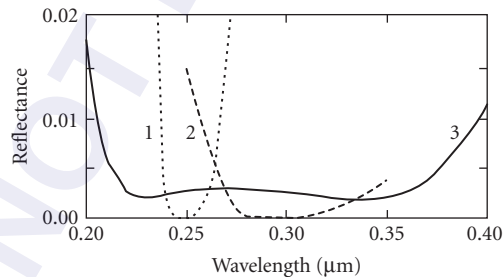


FIGURE 14 Ultraviolet antireflection coatings on fused silica. (Curve 1 after TecOptics;⁸³ curve 2 after Reynard Corp.;⁸⁷ and curve 3 after Spindler & Hoyer.⁸⁸)

Since for some applications the color introduced into an optical system by antireflection coatings is of paramount importance, it has been the subject of many studies.⁹³⁻⁹⁶ One way to avoid the problem is to utilize coatings that are particularly achromatic (Fig. 11c). Thus, for example, 50 surfaces coated with antireflection coating 4.4 would have a transmittance of 78 ± 3 percent across the whole visible spectrum. Nevertheless, of the coatings shown in Fig. 11c only the single-layer antireflection coating is being used extensively. Antireflection coatings can also be used to correct the residual color of lens systems.⁹⁷

Often transmittance rather than reflectance measurements are performed to evaluate antireflection coatings. In general, it is incorrect to assume that $T = 1 - R$. For instance, the transmission of antireflection-coated infrared materials depends not only on the efficiency of the antireflection coating but

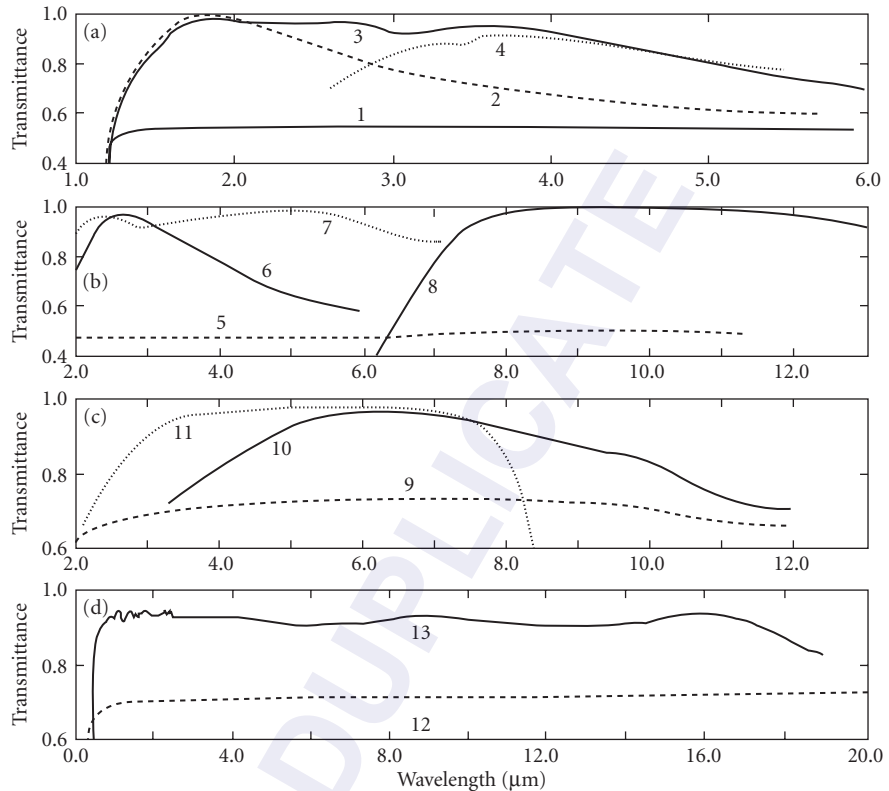


FIGURE 15 Transmittance of plates of infrared materials antireflection-coated on both sides. (a) Silicon. Curve 1: uncoated plate of 1.5 mm thickness; curve 2: single $\lambda_0/4$ layer of SiO_2 ($\lambda_0 = 1.8 \mu\text{m}$); curve 3: $\lambda_0/4$ coatings of MgF_2 , and CeO_2 ($\lambda_0 = 2.2 \mu\text{m}$); and curve 4: hard carbon layer. (b) Germanium. Curve 5: uncoated plate; curve 6: single $\lambda_0/4$ layer of SiO ($\lambda_0 = 2.7 \mu\text{m}$); and curve 7: $\lambda_0/4$ coatings of MgF_2 , CeO_2 , and Si ($\lambda_0 = 3.5 \mu\text{m}$); and curve 8: environmentally stable, high laser damage threshold three-layer design based on ThF_4 and Ge . (c) Irtran II. Curve 9: uncoated plate of 2 mm thickness; curve 10: single $\lambda_0/4$ coating of CeF_3 ; and curve 11: $\lambda_0/4$ layers of MgF_2 and SiO ($\lambda_0 = 4.2 \mu\text{m}$). (d) Zinc selenide. Curve 12: uncoated plate and curve 13: extremely broadband AR coating composed of 398 layers produced by molecular beam deposition. (Curves 1–3, 5–7, 9–11 after Cox and Hass;⁷¹ curve 4 after Balzers;⁷⁶ curve 8 after Oh⁹⁸; and curve 13 after Fisher.⁹⁹)

also on the thickness and temperature of the material. This is because of the finite scatter and absorption in such materials and because of the dependence, in some cases, of the latter on temperature. The measured spectral transmittances of three common antireflection-coated infrared materials at room temperature are shown in Fig. 15. Curves for other materials are given by Cox and Hass.⁷¹

Inhomogeneous and Structured Antireflection Coatings

The interface between two media with refractive indices n_1 and n_2 can be antireflection coated over a very broad spectral region by the application of a transition layer with an index that changes continuously from n_1 to n_2 (Fig. 8d). Many different refractive index profiles have been investigated in the past.^{100–103} Although some of these profiles are more effective than others, all reduce the reflectance to a fraction of a percent over the spectral region in which the coatings are transparent and do not

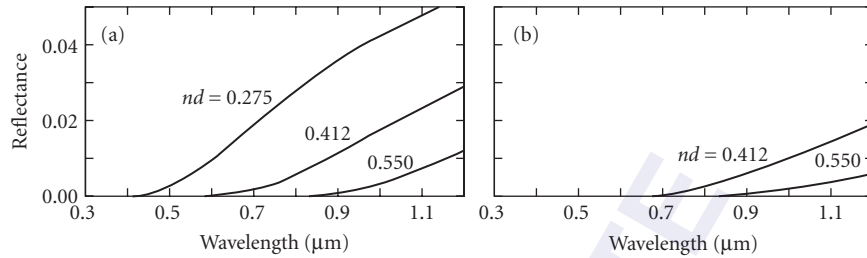


FIGURE 16 Calculated reflectances of the interfaces between two media antireflection-coated with inhomogeneous layers of indicated thicknesses and (complex) refractive indexes that vary smoothly, from the index of one medium to that of the other. (a) Two nonabsorbing media of indexes 1.52 and 2.36. (b) Glass-chromium interface with refractive indexes 1.52 and $2.26 - i0.43$, respectively. (After Anders and Eichinger.¹⁰⁷)

scatter excessively, and for which the optical thickness of the layer is at least one half wavelength (Fig. 16a). A further advantage of inhomogeneous antireflection coatings is that they are not sensitive to the angle of incidence.¹⁰⁴ Processes used for the production of inhomogeneous antireflection coatings include various additive, subtractive, additive/subtractive, and replication methods. Excellent reviews of this topic exist.^{105, 106}

In the additive method relatively dense inhomogeneous layers of varying compositions of two or more compounds are formed on the substrate by physical or chemical deposition processes. Such coatings are mechanically more durable than any other described in this section. However, because of a lack of coating materials with refractive indexes lower than about 1.35, solid inhomogeneous layers are not very suitable for the antireflection coating of air-glass interfaces. The several different inhomogeneous antireflection coatings of this type described in the past do not offer any special advantages over those composed of homogeneous layers.^{70, 108, 109} The situation is different in the case of high-index materials (Fig. 17). An even lower reflectance can be achieved by ending the inhomogeneous layer

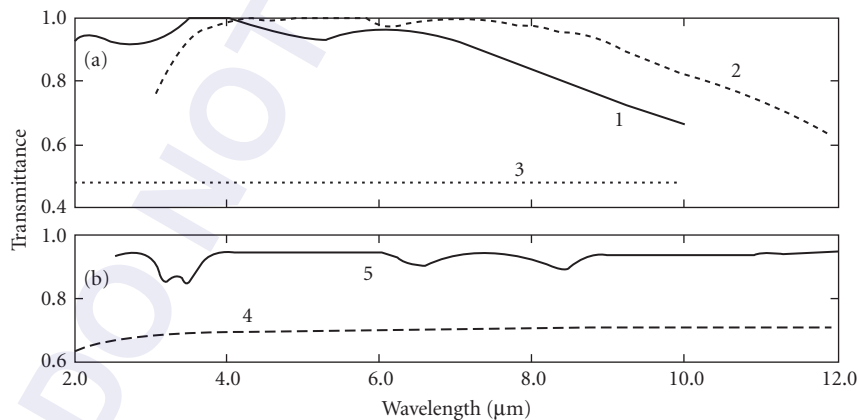


FIGURE 17 (a) Measured transmittance of a germanium plate coated on both sides with inhomogeneous antireflection coatings. Curve 1: 1.2- μm -thick film with an index that changes gradually from that of Ge to that of MgF_2 (after Jacobsson¹¹⁰); curve 2: 1.76- μm -thick film with an index that varies from 4.0 to 1.5, overcoated with a 0.74- μm -thick homogeneous MgF_2 layer (after Scheuerman¹¹¹); and curve 3: transmittance of an uncoated plate. (b) Measured transmittance of a TBI-5 plate. Curve 4: uncoated and curve 5: coated on both sides with ten 2.5- μm -thick homogeneous layers with refractive indexes that vary between 1.3 and 2.29 and that are obtained through the evaporation of suitable NaF-CdTe mixtures. (After Kuznetsov and Perveyev.¹¹²)

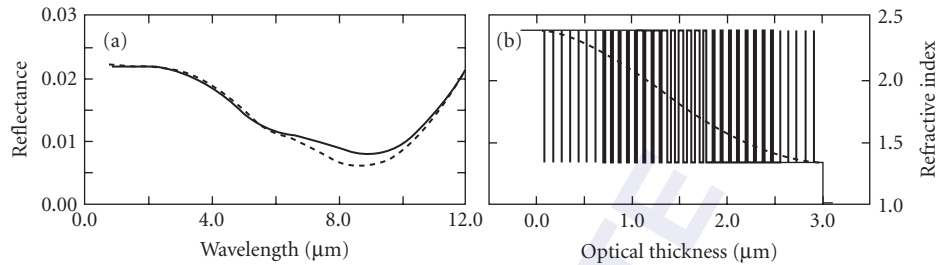


FIGURE 18 Calculated performances (a) and refractive index profiles (b) of an inhomogeneous layer coating and its two-material equivalent. (After Southwell.¹¹⁶)

when its index is equal to the square of that of the lowest-index coating material available. It is then possible to complete the coating by depositing an additional homogeneous quarter-wave-thick layer of that material (Fig. 17a).

A dense inhomogeneous layer can be approximated by a series of homogeneous layers of gradually decreasing refractive indexes (curve 10.1, Fig. 11f).^{113,114} Such layers can be prepared by evaporating a series of appropriate mixtures of two coating materials or, without mixing, by using the Herpin equivalent-index concept¹¹⁵ to simulate intermediate refractive indexes. An even more practical solution is to replace the inhomogeneous layer by a series of thin homogeneous layers of two materials only (Fig. 18).¹¹⁶

In another additive process a refractive index variation down to a value of 1.0 is achieved by depositing onto the substrate microspheres of transparent oxides or fluorides which form pyramid-like clumps (Fig. 8e).¹¹⁷ If losses due to scattering are to be low, the average lateral size of the features of this structure must be a small fraction of the shortest wavelength for which the coating is to be effective. A reflectance of 0.3 percent can be achieved, but the films are fragile.

The subtractive methods are attractive because they do not require expensive deposition equipment. The surface to be antireflection coated is leached and/or etched to form a porous transition layer in which the index varies with thickness. Not all optical materials can be treated in this way. However, special phase-separable glasses have been developed that lend themselves well to this process.^{118–121} Fairly durable antireflection coatings with a reflectance of less than 0.5 percent for the 0.35 to 2.5 μm spectral region have been produced in this way (Fig. 19a). Some other materials, such as Lexan, Mylar, and CR-39 plastic, require an ion implantation pretreatment before the etching can be applied.^{101,122}

In the technologically important additive/subtractive method, a single glasslike film is first deposited by a sol-gel process onto the surface that is to be antireflection coated. The composition of the film is such that, after phase separation, it can be readily leached and/or etched to form a porous microstructure with a controlled refractive index gradient.^{123,126} This eliminates the need for the use of expensive phase separation glass components. Such coatings have reflectances as low as 0.13 percent and a laser damage threshold that is 4 times higher than that of coatings produced by conventional physical vapor deposition techniques (Fig. 19b).¹²⁷ Variants of this process exist.

Microstructured surfaces can also be produced in polymeric and similar materials by a replication process from a suitable cast. Average reflectances in the visible of the order of 0.3 percent have been reported for surfaces treated in this way (Fig. 19c).¹²⁴

Recently there has been a renewed interest in antireflection coatings in which the variable porosities of leached or etched layers are simulated by dense regular shaped structures formed by photochemical or mechanical means. Clapham and coworkers appear to have been the first to demonstrate such devices.¹²⁸ They applied a photoresist to a surface, exposed it to two orthogonal sets of ultraviolet interference fringes, and then developed it to form a regular array of protuberances that could be optionally enlarged by additional ion-beam etching. Such surfaces can reduce the reflectance to less than 0.3 percent in the visible part of the spectrum (Fig. 19d).¹²⁵ The theory of such structures has been investigated by several workers^{129,130} and devices for wavelengths extending into the micrometer

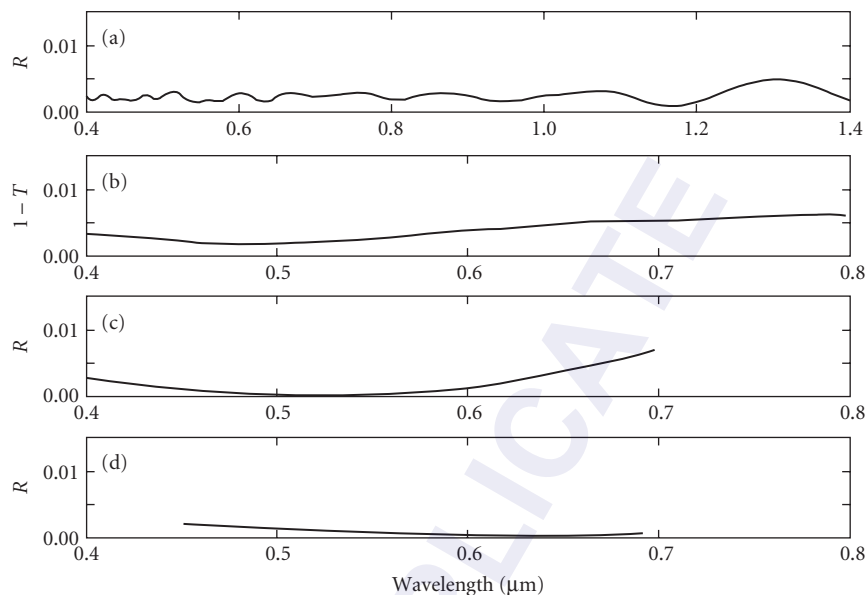


FIGURE 19 Measured performance of various broadband inhomogeneous antireflection coatings. (a) Subtractive process—leached and etched phase separation glass (after Asahara¹²⁰); (b) additive/subtractive process—leached and etched phase-separable film deposited onto fused silica (after Yoldas¹²³); (c) replication process on a cellulose acetate butyrate surface (after Maffitt¹²⁴); and (d) structured antireflection coating in photoresist (after Wilson¹²⁵).

and sub-millimeter region have been fabricated.^{77,121} Efficient laboratory prototype AR coatings produced by the glancing angle deposition process have been demonstrated by Kennedy and Brett.¹³¹

Universal AR Coatings

The idea of a universal AR coating which would reduce the reflectances of a range of substrates with different refractive indices has been first proposed by Vvedenski et al.¹³² Such an AR coating could simplify operations in a coating shop, because all substrates, regardless of the glass that they were made of, could be coated in one deposition run. The prize that one pays for this convenience is that the performance of an AR coating consisting of the same number of layers, but designed for one substrate only will be significantly better.

In Fig. 20 are shown the results obtained for two AR coatings of this type consisting of 11, 17 layers each (rows a, b, respectively). The coatings were designed to simultaneously reduce the reflectances of substrates of index 1.48, 1.55, 1.60, 1.65, 1.73, and 1.75. In column 1 the reflectances of the uncoated and of the coated substrates are shown. Also indicated are the average reflectances of all the coated substrates. The refractive index profiles of the six systems for the different substrate materials shown in a2, b2 differ only in the substrate refractive indices. Increasing the number of layers in the universal AR coatings shown in this figure will result in a diminishing return in the performance. Several papers on this topic have been published by Yeuch-Yeong Liou.¹³⁴

Antireflection Coating of Absorbing and Amplifying Media

Antireflection coatings for glasses and semiconductors in regions of weak absorption, and for present-day laser materials in which k_j , the imaginary part of the complex refractive index [Eq. (25)] is small and negative, differ little from those described above.¹³⁵

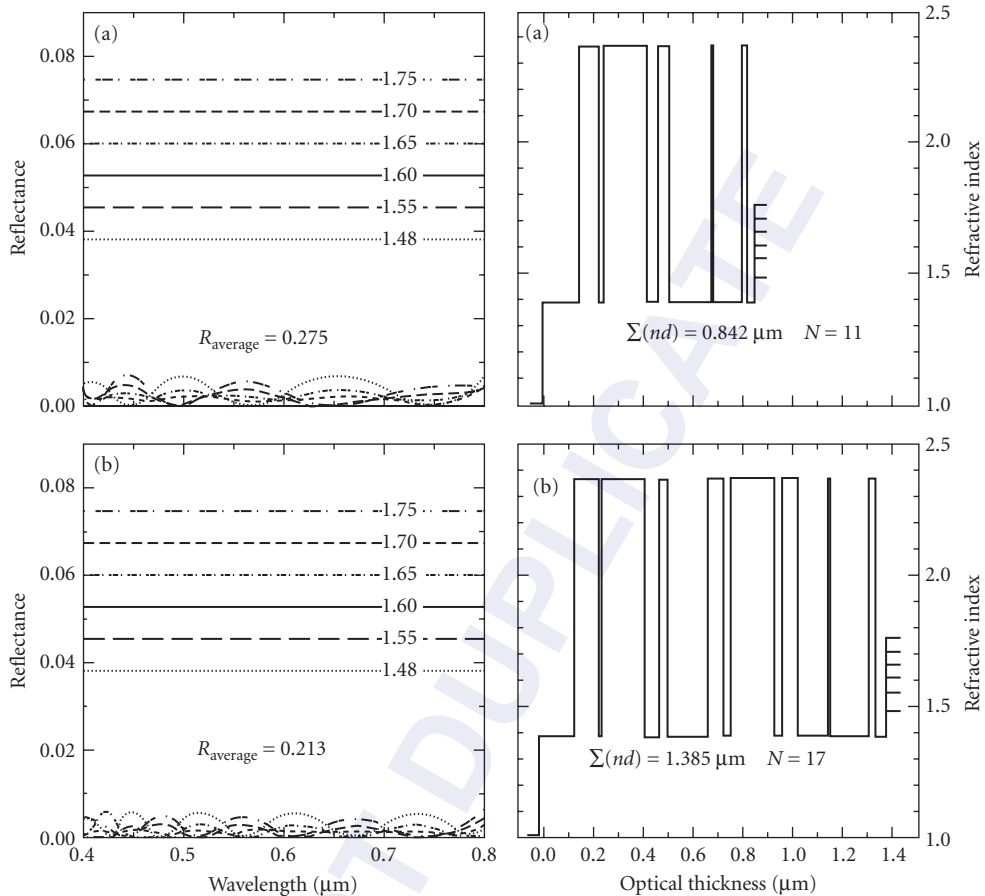


FIGURE 20 Universal AR coatings for substrates with refractive indices 1.48, 1.55, 1.60, 1.65, 1.70, and 1.75. Rows *a*, *b* depict solutions based on 11, 17 layers. Column 1 depicts the reflectances of the substrates before and after the application of the AR coating. In column 2 are shown the refractive index profiles of the solutions. (After Dobrowolski and Sullivan.¹³³)

The reduction of the reflectance of opaque materials for architectural, decorative and technical purposes leads to a corresponding increase in the absorption.^{136,137} This can be utilised to improve the efficiency of radiation detectors and to control the solar-absorptance and thermal-emittance characteristics of surfaces (see, for example, Ref. 138 and 139). A measured example of the reduction in the reflectance of a metal surface attainable with a homogeneous nonabsorbing layer is shown in Fig. 21.

Inhomogeneous transition layers whose refractive index and extinction coefficient change gradually from the values of one of the media to those of the other are also very effective.¹⁰⁷ The calculated reflectance of a glass-chromium interface coated in this way, useful for blackening of prism faces, lens edges, scales, and so on, is shown in Fig. 16*b*. Metal-air interfaces can be treated a little less effectively with single layers because of the lack of suitable low-index materials.

Antireflection Coating of Surfaces Carrying a Thin Film

For some applications it is necessary to deposit a certain film onto a glass surface. The objectionable reflectance that such a film would normally introduce can be avoided by incorporating it into an

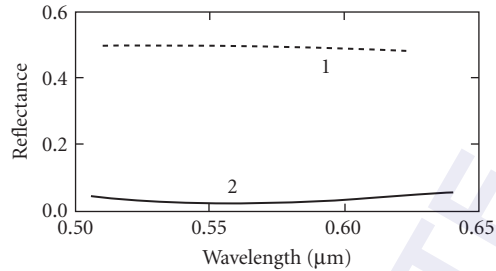


FIGURE 21 Antireflection coating of opaque metals. Curves 1 and 2: reflectance of chromium and of chromium with a ZnS layer. (After Lupashko and Sklyarevskii.¹³⁶)

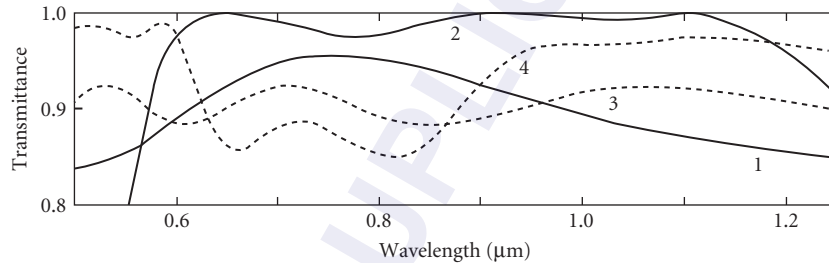


FIGURE 22 Calculated single-surface transmittance of antireflection coatings incorporating a conducting layer. Curves 1 and 2: glass of refractive index 1.755 coated with a 0.566- μm -thick SnO_2 layer before and after three dielectric layers were added to complete the antireflection coating and curves 3 and 4: as above except that the refractive index of the glass and the metric thickness of the SnO_2 film are 1.516 and 0.200 μm , respectively. (After Veremey and Gorbunova.¹⁴²)

antireflection coating. Thus, for example, Fig. 22 shows the calculated transmittances of two conducting layers before and after inclusion in antireflection coatings. The use of homogeneous and of inhomogeneous antireflection coatings with absorbing layers for ophthalmic purposes is described by Katsube et al.¹⁴⁰ and Anders.¹⁴¹

Antireflection Coatings at Nonnormal Angle of Incidence

The calculated performances at 0° , 45° , and 60° incidence of the three commercially most important antireflection coatings and of a 10-layer coating for germanium are shown in Fig. 23. The deterioration with angle of incidence is particularly severe for the narrowband high-efficiency antireflection coating. Figure 23*d* suggests that the closer the design of an antireflection coating approximates an inhomogeneous transition layer (Sec. 7.5, subsection “Inhomogeneous and Structured Antireflection Coatings”) the less angle-dependent is its performance.

To design an antireflection coating for one angle of incidence and one plane of polarization the effective thicknesses and refractive indexes [Eqs. (18) and (19)] of its layers should satisfy the relations set out in Table 1. In practice, small departures from those conditions are required to optimise the performance with good coating materials. Calculated curves for two sets of two- and three-layer coatings designed for use at 45° are shown in Fig. 24.

If the obliquely incident radiation is unpolarized, a compromise is necessary. The effective thicknesses are matched for the design angle, but the refractive-index conditions set out in Table 1

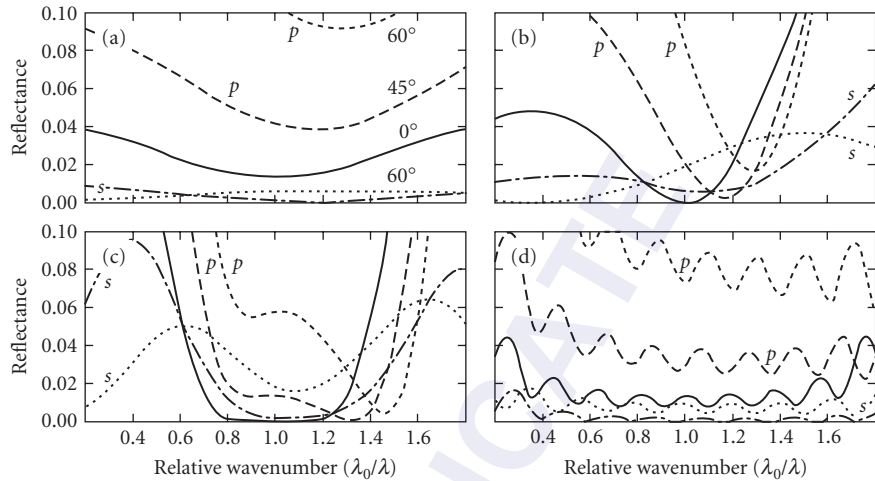


FIGURE 23 Calculated reflectance of (a) one-, (b) two-, (c) three-, and (d) 10-layer antireflection coatings (systems 1.1, 2.1, 3.4, and 10.1 in Table 1) at angles of incidence of 0° (solid curve), 45° (dashed curve), and 60° (dotted curve) for light polarized parallel and perpendicular to the plane of incidence.

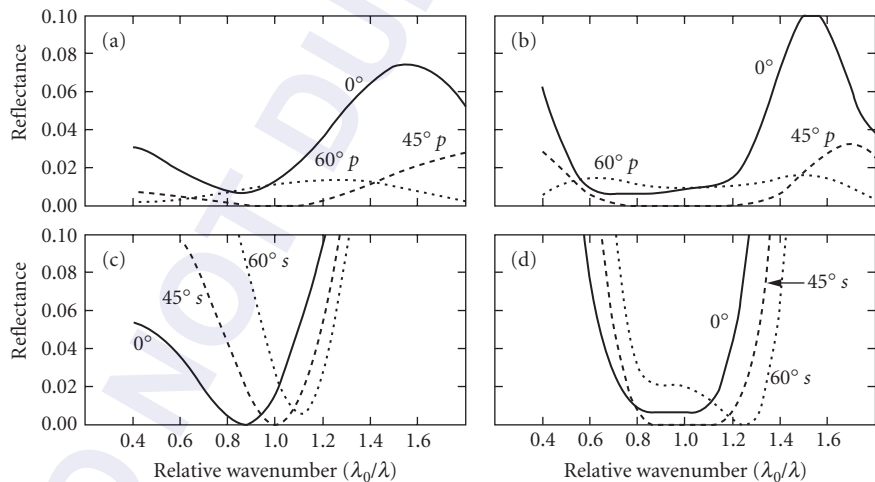


FIGURE 24 Calculated reflectance at 0° (solid curve), 45° (dashed curve), and 60° (dotted curve) of (a) and (c) two-layer and (b) and (d) three-layer antireflection coatings designed for light incident at 45° and polarized parallel or perpendicular to the plane of incidence. (After Turbadar.^{143,144})

are satisfied for normal incidence since they cannot be satisfied for both polarizations at the same time (Fig. 25).

Achromatic antireflection coatings can be designed that are optimized for both polarizations at the same time or that are suitable for use over a range of angles of incidence (Fig. 26b).^{24,145} But the problem becomes more difficult the greater the spectral and angular ranges required, especially if angles of incidence greater than 60° are involved.

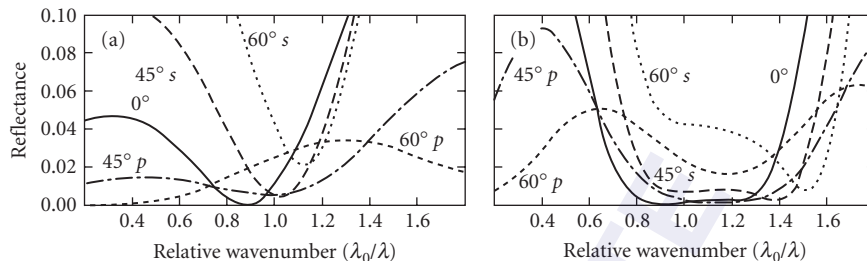


FIGURE 25 Calculated reflectance at 0° (solid curve), 45° (dashed curve), and 60° (dotted curve) of (a) two-layer and (b) three-layer antireflection coatings (systems 2.1 and 3.4 in Table 1) for use with unpolarized light, with the thicknesses of the layers optimized for an angle of incidence of 45° .

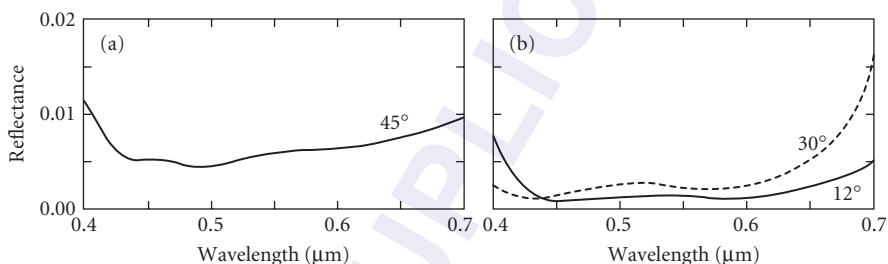


FIGURE 26 Average reflectances for unpolarized light of antireflection coatings designed (a) for 45° incidence and (b) for use with a convergent light cone of semi-angle of 30° . (*Optical Coating Laboratory*.⁸⁴)

In Fig. 27a is shown the variation with angle of incidence of the calculated average reflectance of a six-layer AR coating system for the He-Ne laser wavelength. The reflectance is less than 1 percent for all angles smaller than 60° , but only for a narrow wavelength region around the design wavelength. The performance was limited by the lowest available refractive index at the time.

It has been shown previously (Fig. 19) that, for normal incidence of light, the reflectance from an interface between two media with refractive indices n_s , n_m can be effectively removed over a broad spectral region by an inhomogeneous layer with a refractive index that changes gradually from n_s to n_m , providing that the optical thickness of this inhomogeneous layer is at least one half wavelength thick at the longest wavelength for which the reflectance is to be low. It has been shown numerically¹⁵¹ and theoretically,¹⁵² that this rule continues to hold for up to very high angles of incidence, providing that the effective optical thickness is at least one-half wavelength at the highest angle of interest. This means that broadband inhomogeneous layer AR coatings effective up to very high angles of incidence will be very thick. However, it is possible to transform such solutions to homogeneous layer systems with which, because of the additional thin-film interference effects, similar performances can be obtained with a small number of layers and much smaller overall optical thicknesses. Such solutions still call for the use of at least one layer with a refractive index that is very close to the medium index, n_m . Such solutions can be readily implemented for glass/glass interfaces (Fig. 27c). The angular measurements shown in this diagram were obtained with a laser light source, but the AR region was quite broad. When the medium is air and a maximum angle of incidence is 85° is stipulated, a refractive index of the outermost layer should be of the order of 1.02. In the infrared part of the spectrum there are a number of Reststrahlen materials with a refractive index that falls below unity. The measured average spectral reflectance for a number of angles of incidence of an AR coating based on such a material is shown in Fig. 27b. The antireflection region is narrow because the refractive index of the outermost layer has the required value only in a narrow wavelength band. In the visible part of the spectrum such indices can only be achieved with very porous or structured layers. Ways of producing such layers are

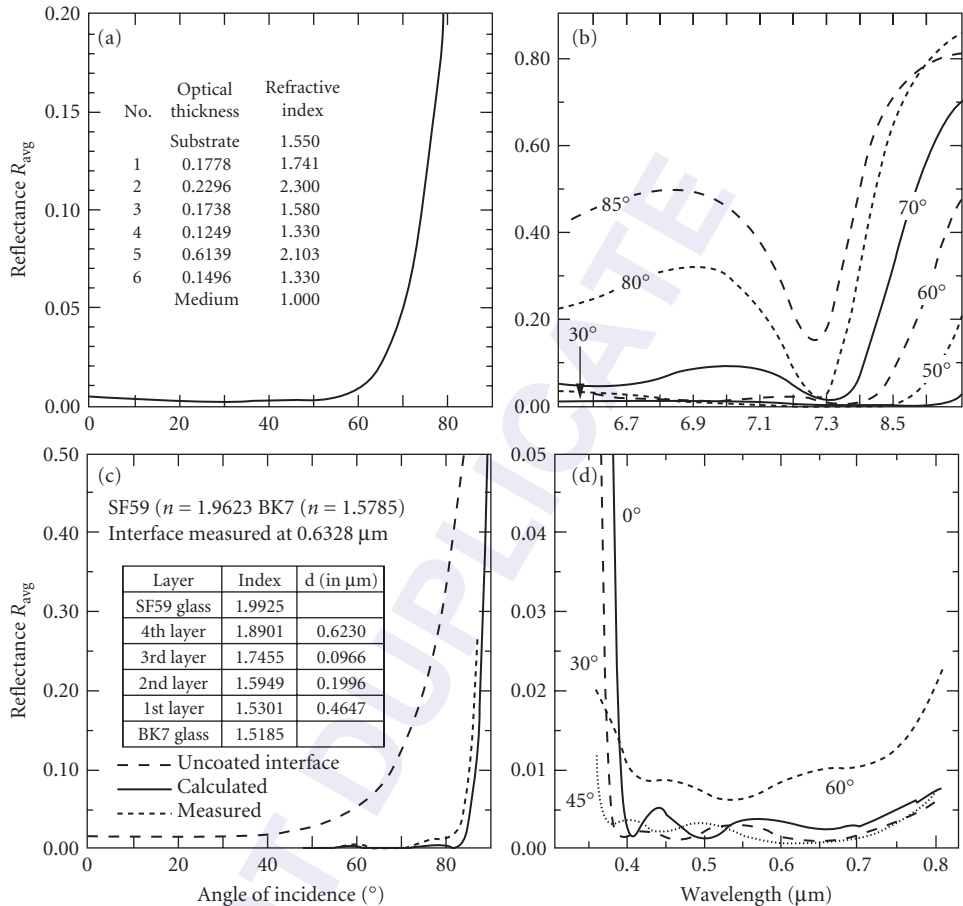


FIGURE 27 Wide angle AR coatings. (a) Calculated average reflectance of a conventional AR coating for the He-Ne laser wavelength for angles of incidence up to 60° (after Dobrowolski *et al.*^{146,147}); (b) measured average reflectance of a Si/air interface for angles of incidence up to 85° with an AR coating based on the use of a Reststrahlen material (after Dobrowolski *et al.*¹⁴⁸); (c) measured average reflectance of a glass/glass interface for angles of incidence up to 85° (after Ma *et al.*¹⁴⁹); and (d) predicted performance of a hybrid wide-band AR coating for angles of incidence up to 60° (after Schulz *et al.*¹⁵⁰).

being investigated at many laboratories. The predicted average reflectance of a hybrid AR coating consisting of seven solid thin films and an outermost structured layer is shown in Fig. 27d. It is expected to be less than 1 percent over a broad spectral region for all angles of incidence smaller than 60°. This will be a huge improvement over the performance of the coating depicted in Fig. 27a.

The design of antireflection coatings for very high angles of incidence has been considered by Monga.^{153,154}

Nonoptical Properties of Antireflection Coatings

Antireflection coatings, in addition to optical specifications, often have to satisfy additional, nonoptical requirements. One such requirement may be to provide scratch resistance to soft, plastic surfaces. In Fig. 28a is shown the measured reflectance of an AR coating produced by Schulz *et al.* for

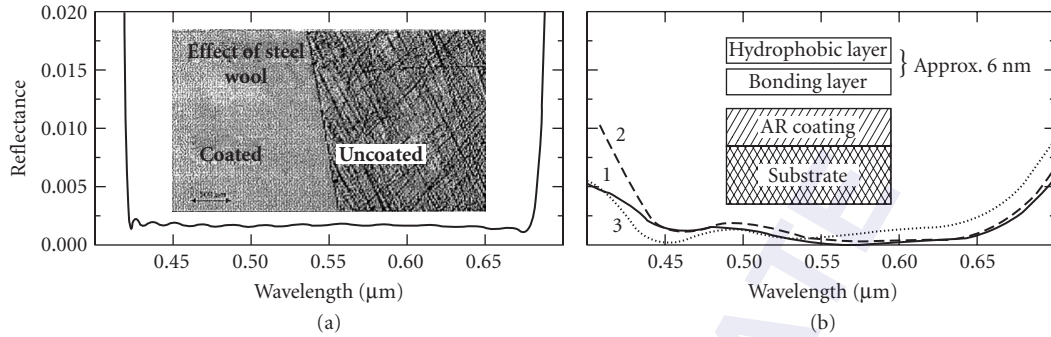


FIGURE 28 Nonoptical properties of antireflection coatings. (a) Measured reflectance of a 27-layer, 3- μm -thick $\text{SiO}_2/\text{Ta}_2\text{O}_5$ AR coating. The insert illustrates the effect of steel wool on the coated and uncoated portions of the substrate (after Schulz *et al.*⁴³). (b) Performance of a water- and dust-repellent AR coating. Curves 1 and 2 are the calculated and measured reflectances. Curve 3 is the performance after cleaning the coating with acetone (after Bruynnghe *et al.*¹⁵⁵).

a PMMA surface.⁴³ It consists of 27 layers of SiO_2 and Ta_2O_5 , with an overall metric thickness of over 3 μm , of which most corresponds to SiO_2 . The inset shows, that the AR coating provides good protection from abrasion. Bruynooghe *et al.* describe an AR coating which repels water, oil, and dust (Fig. 28b).¹⁵⁵ This coating consists of a more or less conventional AR coating onto which is deposited a bilayer coating, the first part of which provides for good adhesion to the AR coating, and the second, the required hydrophobic properties. The total thickness of the bilayer is 0.006 μm .

7.6 TWO-MATERIAL PERIODIC MULTILAYERS THEORY

Nonabsorbing $[AB]^N$ and $[AB]^N A$ Multilayer Types

Let a periodic multilayer be composed of N periods AB , where A, B represent layers of refractive indexes n_A, n_B and optical thicknesses $n_A d_A, n_B d_B$. The most general representation of the complete multilayer is

$$AB \cdot AB \cdot AB \cdots AB = [AB]^N \quad (39)$$

In practice it is customary to write $[HL]^N$ or $[LH]^N$, depending on whether n_A is greater or less than n_B , respectively. The spectral-reflectance curve of the multilayer $[AB]^N$ will lie within a pair of envelopes which, at normal incidence, depend only on $n_A d_A, n_B d_B, n_A, n_B, n_s$, and n_m .^{156,157} For $n_s = n_m$ the lower of the envelopes becomes $R = 0$. The envelopes contain high-reflectance zones within which the reflectance at each wavelength increases monotonically with the number of periods, approaching 1.0 as N tends to infinity (Fig. 29). Outside these high-reflectance zones the curves exhibit subsidiary maxima and minima whose number depends on $n_A d_A / n_B d_B$ and which increases with N . The first-order high-reflectance zone occurs at a wavelength λ_1 given by

$$n_A d_A + n_B d_B = \frac{\lambda_1}{2} \quad (40)$$

and subsequent zones occur at wavelengths λ_q ($\lambda_1 > \lambda_2 > \lambda_3 \dots$) given by

$$N(n_A d_A + n_B d_B) = q \frac{\lambda}{2} \quad q = 2, 3, 4 \dots \quad (41)$$

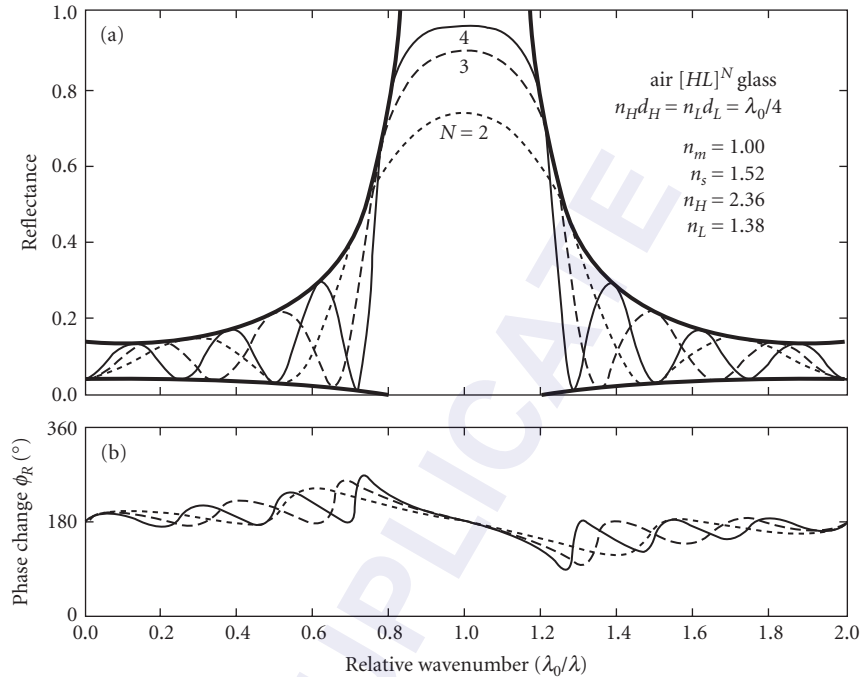


FIGURE 29 (a) Calculated reflectance and (b) phase change on reflection of periodic multilayers of the type $[HLH]^N$. H and L stand for high- and low-refractive-index layers of quarter-wave thickness at $\lambda_0 = 1.0 \mu\text{m}$. The heavy lines are the envelopes of the reflectance curves.

providing that at these wavelengths

$$n_A d_A, n_B d_B \neq p \frac{\lambda_0}{2} \quad p=1,2,3 \dots \quad (42)$$

It is thus possible, by choosing suitable thickness ratios, to arrange for or to suppress high-reflectance zones in several spectral regions at the same time. This is useful in the design of broadband reflectors (Sec. 7.7, "Rejection Filters"), cutoff filters and hot- and cold mirrors (Sec. 7.8), laser reflectors, and so on. Typical curves for thickness ratios 1:1 and 2:1 are given in Fig. 30. Plotted on a λ_0/λ scale, they show symmetry about all wavelengths λ for which $n_A d_A$ and $n_B d_B$ are both equal to some integral multiple of $\lambda/4$.

Maximum Reflectance For a given refractive index ratio (n_A/n_B) and number of periods N the highest reflectance occurs whenever $n_A d_A, n_B d_B$ are each equal to an odd multiple of $\lambda/4$. It is given by

$$R_{\max} = \left[\frac{n_m/n_s - (n_A/n_B)^{2N}}{n_m/n_s + (n_A/n_B)^{2N}} \right]^2 \quad (43)$$

Intermediate reflectances are obtained for the related symmetrical multilayers $[AB]^N A$:

$$R_{\max} = \left[\frac{n_m n_s / n_A^2 - (n_A/n_B)^{2N}}{n_m n_s / n_A^2 + (n_A/n_B)^{2N}} \right]^2 \quad (44)$$

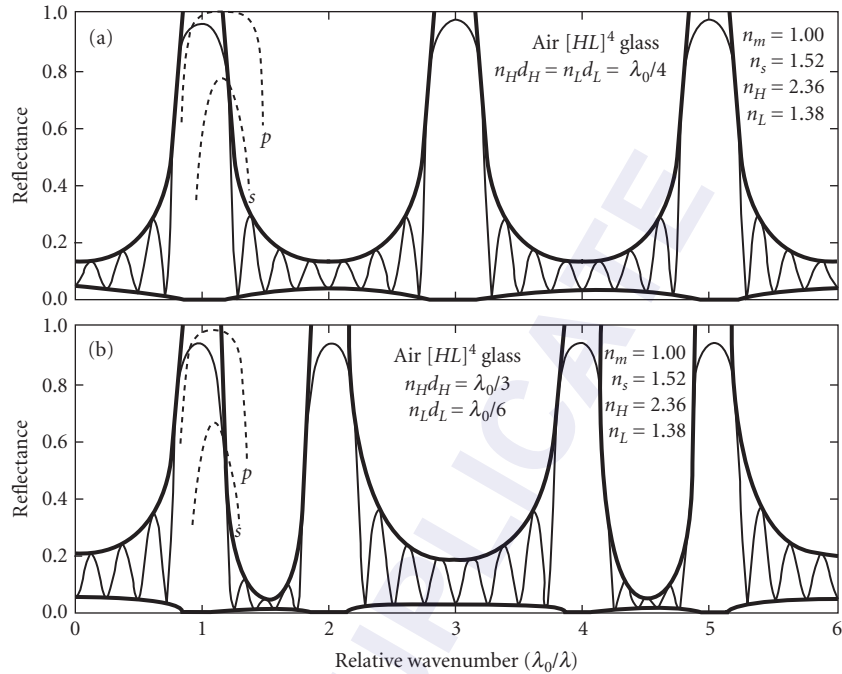


FIGURE 30 Calculated reflectance curves of periodic multilayers with (a) 1:1 and (b) 2:1 thickness ratios. The dotted curves represent the reflectance of polarized radiation incident at an angle of 60° .

Results for a number of such cases are given in Fig. 31. Intermediate reflectances can be obtained by changing the refractive index of any of the layers in the stack.

Explicit expressions for R for other thicknesses are complicated.

Phase Change on Reflection The dispersion of the phase change on reflection from periodic all-dielectric multilayers is much greater than that of metal reflectors. Unless corrected for, it will lead to errors in some metrological and interferometric applications. Like reflectance, it varies rapidly outside the high-reflection zone (Fig. 29b). Within the high-reflection zone it changes almost linearly with wavenumber and is 180° at $\lambda_0/\lambda = 1$. The slope of this portion of the graph increases and approaches a limiting value as N tends to infinity. The limiting value (in degrees per unit wavelength) is given by

$$\left[\frac{d\phi}{d\lambda} \right]_{\lambda} = \left\{ \begin{array}{ll} \frac{180n_m}{\lambda_0(n_H - n_L)} & \text{for } n_H \\ \frac{180n_H n_L}{\lambda_0 n_m (n_H - n_L)} & \text{for } n_L \end{array} \right\} \quad (45)$$

depending on whether the light is first incident on a high- (n_H) or low- (n_L) refractive index layer.¹⁵⁸ The above values should be multiplied by 3, 5, ... in stacks composed of $3\lambda/4$, $5\lambda/4$, ... layers.

Böhme has shown that, by changing the refractive index of one of the layers in a quarter-wave stack, it is possible to obtain a zero value of the phase change on reflection ϕ at λ_0 .¹⁵⁹

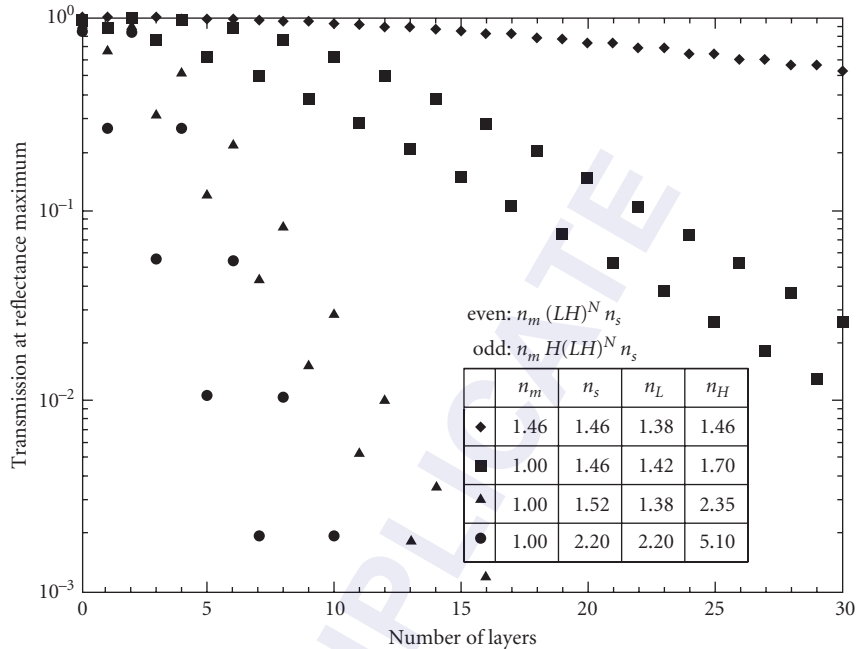


FIGURE 31 Transmittance at the reflection maximum of quarterwave stacks composed of MgF_2 and SiO_2 in the visible region (◆), MgF_2 and MgO in the ultraviolet (■), MgF_2 and ZnS in the visible region (▲), and ZnS and PbTe in the infrared (●).

Periodic Multilayers of the $[(0.5A)B(0.5A)]^N$ Type

The construction of such multilayers differs from that of the type $[AB]^N A$ discussed above by having outermost layers of only half the thickness of the remaining layers.¹¹⁵ The position and width of the high-reflection zones are the same in both cases, but in coatings of the $[(0.5A)B(0.5A)]^N$ type it is possible to reduce the height of the subsidiary maxima on either one or the other side of the main reflectance maximum if $n_A t_A = n_B t_B$. Kard et al. describe the optimum choice of all the construction parameters (N, n_A, n_B, n_s, n_m) of such a coating.¹⁶⁰ If the refractive indices n_s and n_m of the surrounding media must be chosen on some other basis, then for maximum improvement of the short- and long-wavelength transmission the refractive indexes n_A and n_B of the coating materials must satisfy

$$n_s n_m = \frac{n_A^3}{n_B} \quad (46)$$

or

$$n_s n_m = n_A n_B \quad (47)$$

respectively. Estimates of the maximum reflectance for larger values of N can be obtained from Fig. 31. Multilayers of this type are of importance in the design of cutoff filters (Sec. 7.8), and are illustrated in Fig. 32, which should be compared with the curves of Fig. 33.

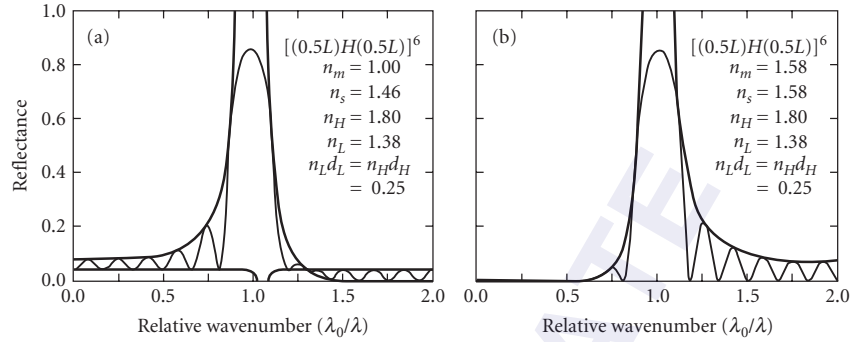


FIGURE 32 Calculated reflectance curves of two periodic multilayers with symmetrical periods of the type $[(0.5L)H(0.5L)]^6$ in which the subsidiary reflectance maxima on (a) the short- and (b) the long-wavelength side of the high-reflectance zone are reduced.

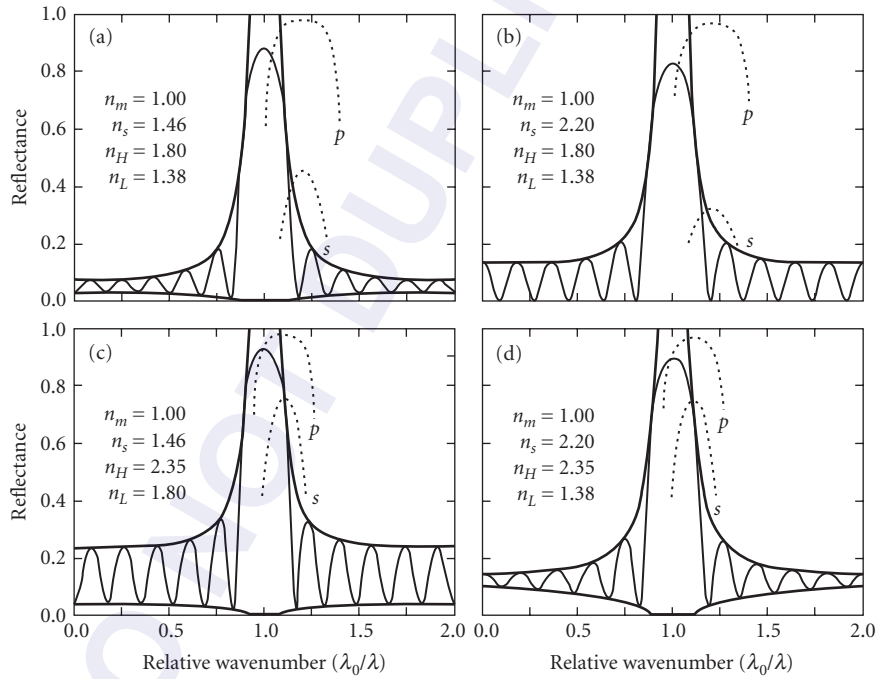


FIGURE 33 Calculated reflectance curves of two quarter-wave stacks of the type $H(LH)^N$ with different values of n_H and n_L but with the same ratio n_H/n_L , (rows 1, 2) deposited onto two different substrate materials (columns 1 and 2). The dotted curves represent the reflectances for light incident at 60° .

Width of the High-Reflectance Zone

For a given value of (n_A/n_B) , the width $\Delta\lambda_R/\lambda$ of the high-reflectance zone is greatest when $n_A d_A = n_B d_B = \lambda/4$, and it is then given by

$$\frac{\Delta\lambda_R}{\lambda} = \frac{4}{\pi} \sin^{-1} \left(\frac{1 - n_A/n_B}{1 + n_A/n_B} \right) \quad (48)$$

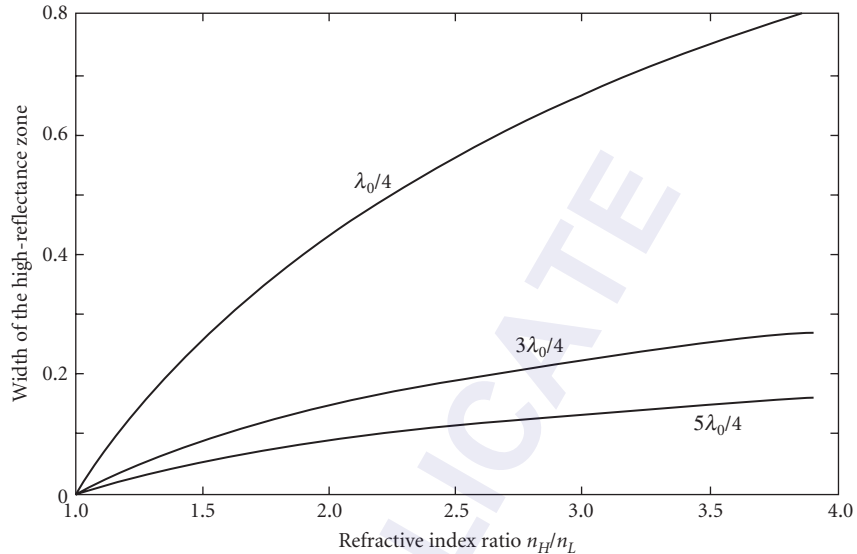


FIGURE 34 Calculated widths of high-reflectance zones of two-material periodic stacks of layers of optical thicknesses $\lambda_0/4$, $3\lambda_0/4$, and $5\lambda_0/4$ for different refractive index ratios.

This width is reduced by a factor of $2N-1$ if N th-order quarter-wavelength layers are used. Graphs of the widths of the high-reflectance zones versus refractive index ratio for $\lambda_0/4$, $3\lambda_0/4$, and $5\lambda_0/4$ layer stacks are given in Fig. 34. Periodic multilayers with equal refractive index ratios have high-reflectance zones of equal widths, but their reflectance curves will not be the same unless the refractive indexes of the surrounding media are also increased by the same ratio (Fig. 33).

Periodic Multilayers of the $[xH \cdot (1-x)L]^N \cdot xH$ Type

As already stated in Sec. 7.6, subsection “Periodic Multilayers of the $[(0.5A)B(0.5A)]^N$ Type” it is not necessary for the optical thicknesses of the layers of a periodic multilayer reflector to be equal to a quarter wave. A multilayer system of the type $[xH \cdot (1-x)L]^N \cdot xH$, where H, L correspond to quarter-wave layers of a high- and low-refractive index and where $0 < x < 1.0$, will have a high reflectance, providing that the number of periods N is high enough.¹⁶¹⁻¹⁶⁴ However, the smaller the value of x , the narrower the width of the high reflectance zone. Given two coating materials, it is thus possible to independently select the reflectance and the width of the rejection region. Figure 35a, and b shows the relationship between R_{\max} , $\Delta\lambda_R/\lambda$ and x , N for values of n_H, n_L , corresponding to the refractive indices of ZnS and polyethylene at about 200 μm .

Angular Sensitivity

For some applications it is important to reduce the angular variation of the reflectance curve. This can be done by using materials with high refractive indexes (Fig. 33) (see also Sec. 7.8, subsection “Angle-of-Incidence Effects”). Another way is to use periods in which the high-index film is thicker than the low-index film (Fig. 30b and Fig. 60a).

Multilayer Reflectors Made of Absorbing Materials

It is possible to achieve a high reflectance with multilayers of the $[AB]^N A$ type even when A, B correspond to absorbing materials. This is of particular interest for spectral regions for which no nonabsorbing coating materials exist. The optical thicknesses of the periods of such systems are still

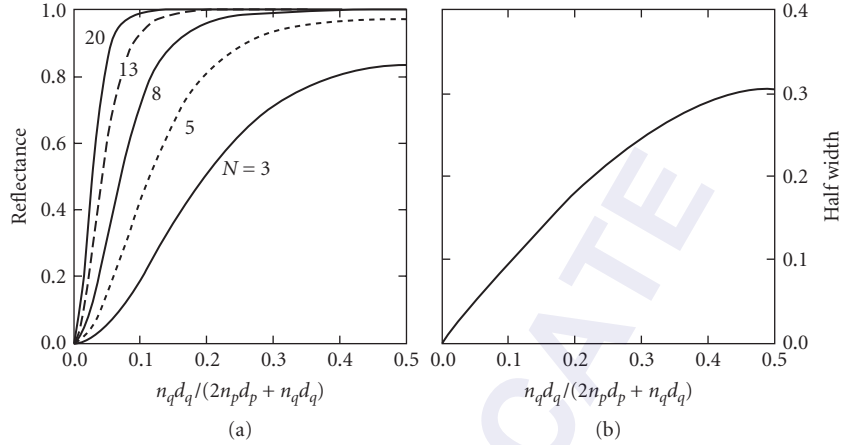


FIGURE 35 Width of the high reflectance zone (a) and maximum reflectance (b) of periodic multilayers of the type $[xH(1-x)L]^N x \cdot H$. The curves were evaluated for refractive indices $n_H = 2.4$ and $n_L = 1.51$. (After Shao.¹⁶⁵)

approximately equal to $\lambda/2$ but, for maximum reflectance, the individual thicknesses d_A , d_B may be quite different, depending on the number of periods N and on the optical constants of the materials used. Reflectors with $k_A > k_B$, have structures that are intermediate to those of quarter-wave stacks (with $k_A = k_B = 0$ and optical thicknesses of $\lambda/4$, in which constructive interference effects are maximized), and those of ideal Bragg crystals (with $k_A \gg k_B$, and in which $d_A \ll d_B$ to minimize absorption losses).

In the extreme ultraviolet (XUV) and the soft x-ray regions ($n - 1$), k is much smaller than 1 for most coating materials.^{165–167} To design a periodic multilayer with a high normal incidence reflectance for a given wavelength, it is first necessary to choose a material ($n_B - ik_B$) with the lowest possible extinction coefficient.^{166,168} Next a second, chemically compatible material ($n_A - ik_A$) is selected with the lowest extinction coefficient that will maximize the normal incidence Fresnel reflection coefficient of the interface between the two materials given by

$$r_{BA} = \frac{(n_B - n_A) - i(k_B - k_A)}{(n_B + n_A) + (k_B + k_A)} \quad (49)$$

Vinogradov and Zeldovich use a factor β_{opt} to relate the metric thicknesses of the layers A and B that yield a maximum reflectance to the overall thickness d_{opt} of the period:¹⁶⁹

$$d_A = \beta_{\text{opt}} \cdot d_{\text{opt}} \quad \text{and} \quad d_B = (1 - \beta_{\text{opt}}) \cdot d_{\text{opt}} \quad (50)$$

where β_{opt} is obtained by the solution of the equation

$$\tan(\pi\beta_{\text{opt}}) = \pi \left[\beta_{\text{opt}} + \frac{n_B k_B}{n_A k_A - n_B k_B} \right] \quad (51)$$

d_{opt} is approximately equal to $\lambda/2$, but Vinogradov and Zeldovich give a more accurate expression for this quantity, as well as for the limiting reflectance R and the number of periods N required to reach that value.¹⁶⁹

7.7 MULTILAYER REFLECTORS—EXPERIMENTAL RESULTS

In the calculations of the data for Figs. 29 to 34, the dispersion of the optical constants of the materials was ignored, and it was assumed that the films were absorption and scatter free and that their thicknesses had precisely the required values. In practice, none of these assumptions is strictly valid, and hence there are departures from the calculated values. In general, the agreement is better within the high-reflectance zone than outside.

Reflectors for Interferometers, Lasers, etc.

The measured reflectances and transmittances of a number of quarter-wave reflectors suitable for use in Fabry-Perot interferometers are shown in Fig. 36. The transmission of a typical commercial laser reflector for $\lambda = 0.6328 \mu\text{m}$ is shown in Fig. 37. The measured reflectances of a number of highly reflecting coatings for the infrared spectral region are shown in Fig. 38.

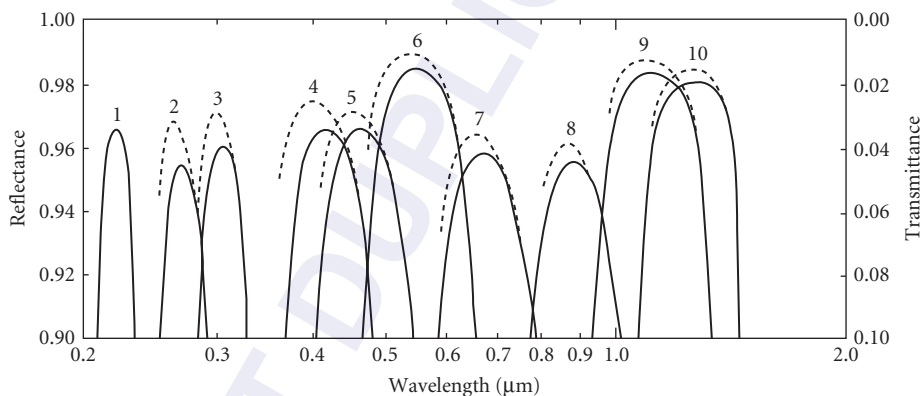


FIGURE 36 Measured spectral reflectance and transmittance curves of periodic all-dielectric reflectors of the type $[HL]^N H$ for the ultraviolet, visible, and near-infrared parts of the spectrum. Curve 1: 27 layers of MgO and MgF_2 ; curves 2 and 3: 11 and 13 layers of PbF_2 and cryolite, respectively; curves 4, 5, 7, and 8: 7 layers of ZnS and cryolite; curves 6, 9-, and 10: 9 layers of ZnS and cryolite. (Curve 1 after Apfel¹⁷⁰ and curves 2 to 10 after Hefft et al.¹⁷¹)

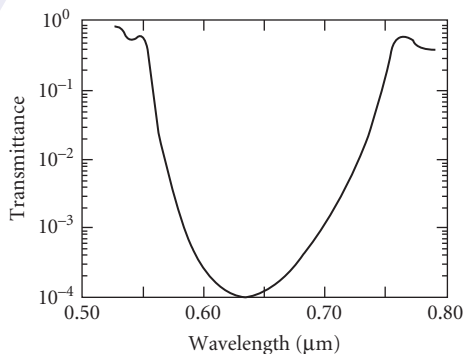


FIGURE 37 Measured spectral transmittance of a laser reflector. (After Costich.¹⁷²)

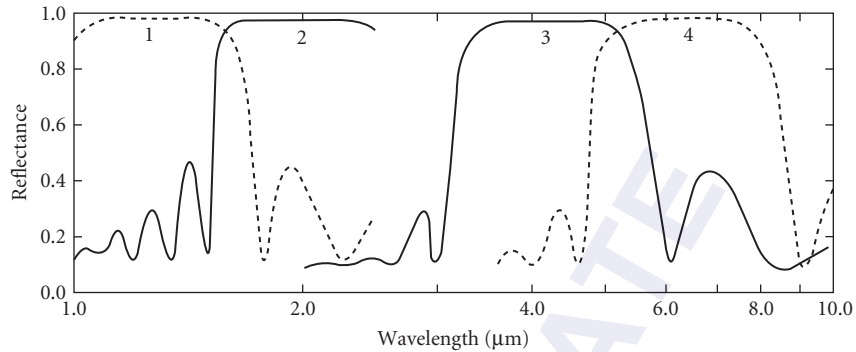


FIGURE 38 Measured spectral-reflectance curves of all-dielectric reflectors for the infrared made of stibnite and chiolite. Curves 1 and 2: multilayers of the type $[HL]^4H$ on glass; curves 3 and 4: type $[(0.5L)H(0.5L)]^4$ on barium fluoride. (After Turner and Baumeister.¹⁷³)

Both “soft” laser coatings that can be dissolved in weak acids and “hard” coatings that can be removed only through polishing are available commercially.

Effects of Imperfections

Thickness Errors Small errors in the thickness of the layers of a quarter-wave stack have only a very small effect on the reflectance and on the phase change on reflection of the multilayer within the high-reflectance zone but they do affect the performance outside the zone.^{174,175} In fact, the effect may be quite serious: thickness variations can give rise to an apparent lack of flatness of the substrate surface.^{176–178}

Dispersion The most noticeable effects of dispersion of the optical constants in quarter-wave stacks of a given multilayer type are the increase in the peak reflectance with decrease in wavelength (see Fig. 36) and the asymmetry of the maxima on either side of the main reflectance maximum.

Absorption The separation between the transmission and reflectance curves in Fig. 36 is mostly due to absorption. These losses can limit the usefulness of the reflectors for some applications. Thus, for example, in interference filters and Fabry-Perot interferometers they lead to a reduction in the peak transmissions and limit the attainable half widths. In optical information-storage devices they set a limit to the highest reflectance attainable. In lasers the losses counteract directly the gain in the laser medium. In addition, absorption within the layers is responsible for damage to laser reflectors (see also Sec. 7.4, subsection “Laser Damage”).

If the two materials used for the construction of periodic quarter-wave stacks have small, but finite extinction coefficients, the resulting absorption at first reduces both the transmission and reflection coefficients.^{156,176,179,180} With an increase in the reflectance of a multilayer of the $[HL]^N H$ type the absorption occurs more and more at the expense of the reflection coefficient, approaching a limiting value of

$$A = -\delta R = \frac{2\pi n_m}{n_H^2 - n_L^2} (k_H + k_L) \quad (52)$$

that is independent of the number of layers.² The corresponding expression for a multilayer $[HL]^N$ in which a low refractive index faces the incident medium is

$$A = -\delta R = \frac{2\pi(n_L^2 k_H + n_H^2 k_L)}{n_m(n_H^2 - n_L^2)} \quad (53)$$

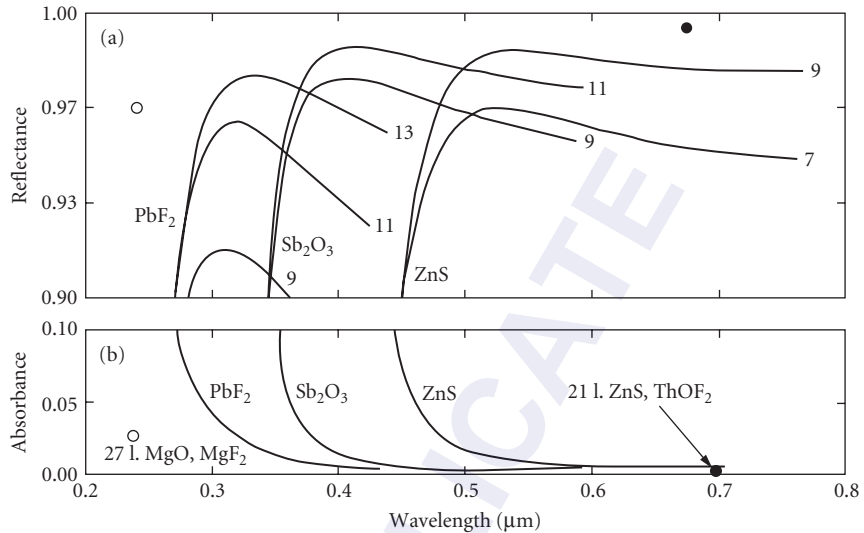


FIGURE 39 Variation with wavelength of the observed peak reflectance of 7-, 9-, 11-, and 13-layer quarter-wave stacks made of PbF₂, Sb₂O₃, or ZnS with cryolite as the low-index material (a). The lower curves indicate the limiting value of the absorbance with these materials (b). (After Honcia and Krebs.¹⁸¹ Results obtained by Apfel¹⁷⁰ (o) and Behrndt and Doughty¹⁸² (●) with other materials are included for comparison.)

Figure 39 shows the spectral variation of experimentally determined maximum reflectances and limiting losses of various quarter-wave multilayer reflectors.

The absorption in periodic multilayers composed of materials having finite extinction coefficients can be reduced below the values given by Eqs. (52) and (53) if a structure of the $[xH \cdot (1-x)L]^N \cdot xH$ type is used.¹⁶³

Surface and Interface Imperfections In thin-film calculations it is usually assumed that substrate surfaces and interfaces between the layers are smooth and that the layers are homogeneous. In practice, substrate surfaces and interfaces have a certain roughness and, at times, thin uniform or inhomogeneous interface layers formed between the boundaries of two layers as a result of oxidation, chemical interactions, or interdiffusion of the two coating materials. The interface layers, as a rule, are different at the *AB* and *BA* boundaries and are typically 0.0003 or 0.0010 μm thick. These and other imperfections of the layer system result in reduced reflection and/or in scatter. If ignored in the model used to represent the multilayer, they add to the discrepancies that are observed between the calculated and experimental data.

Scattering losses are particularly significant at shorter wavelengths. For this reason there have been many theoretical and experimental investigations of scattering of surfaces and thin films. Investigations have shown that when scatter does occur, most of the light is scattered in directions that are close to that of the specularly reflected light.^{156,183} The experimental results for a typical mirror are shown in Fig. 40.

Very low loss reflectors Mirrors with very low losses are required for use in laser cavities and in ring lasers. Mirrors with a combined loss L (= transmission + absorption + scatter) of the order of 5×10^{-5} are commercially available. With special manufacturing techniques practically loss-free reflectors can now be made.

A 41 quarter-wave stack made of Ta₂O₅ and SiO₂ layers with a combined loss of $L = 1.6 \geq 10^{-6}$, corresponding to a reflectance of 0.9999984 at 0.633 μm, has been reported.¹⁸⁵ The absorption and scatter losses were estimated to be of the order of $1.1 \geq 10^{-6}$. The essential starting point for the manufacture of such coatings are superpolished substrates with a surface roughness of 0.5 Å rms or

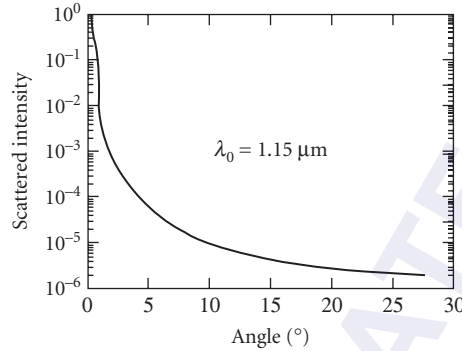


FIGURE 40 Measured intensity of the scattered radiation as a function of the angle away from the direction of specular reflection. (After Blazey.¹⁸⁴)

less. The layers were deposited by reactive ion-beam sputtering from high purity oxide targets in a cryogenically pumped, fully automated deposition system.¹⁸⁶

An even greater challenge will be the requirements of reflecting coatings for Michelson interferometers for the detection of gravitational waves.^{185,187} For this application, not only must the absorption and scattering losses be less than 0.5 and 2 ppm, but in addition, to reduce the thermal noise arising from mechanical loss in the coatings, it will be necessary to match the thermal expansion coefficient and the Young's modulus of the coating materials to those of the substrates.

Multilayers for the soft x-ray and XUV regions The effect of roughness and of interface layers is particularly important in soft x-ray and XUV multilayers because the dimensions of these defects are comparable to the thicknesses of the individual layers and of the wavelength.^{165,166,168} For this reason much attention has been focused on the proper modeling of such coatings. Many workers use several very thin layers and the matrix method outlined in Sec. 7.2, subsection "Matrix Theory for the Analysis of Multilayer Systems" to model the effects of roughness and interface layers.^{188,189} Others make use of the following recursive formula for the amplitude reflectance r_j of the first j layers of the system:

$$r_j = \frac{r_{j-1} + r_{BA} \exp(2i\delta_j)}{1 + r_{j-1} r_{BA} \exp(2i\delta_j)} \quad (54)$$

Here δ_j is the effective optical thickness of the j th layer [Eq. (18)], r_{j-1} is the amplitude reflectance of the first $(j-1)$ layers and r_{BA} is the amplitude reflectance of the interface between the j th and $(j+1)$ th layers. In this approach, if the Fresnel amplitude reflection coefficient r_{BA} is replaced by

$$r_{BA} \exp \left\{ -\frac{1}{2} \left[\frac{4\pi}{\lambda_0} \sigma \Re e(\tilde{n} \cos \tilde{\theta}_0) \right]^2 \right\} \quad (55)$$

the combined effect of roughness and of interface layers σ can be allowed for.¹⁹⁰ In the hard x-ray region, where $n \approx 1$, $k \approx 0$ for all materials, the exponential term in the above expression reduces to the so-called Debye-Waller factor DW,

$$DW = \exp \left[-\frac{1}{2} \left(\frac{4\pi}{\lambda_0} \sigma \cos \theta_j \right)^2 \right] \quad (56)$$

The calculated reflectance of a XUV mirror, with and without the effect of surface imperfections is shown in Fig. 41. These significant reflection losses can be reduced through interface engineering—a

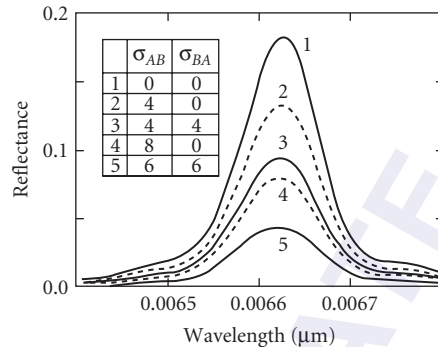


FIGURE 41 Calculated reflectances of 64-period x-ray mirrors with different imperfections σ (in Å) at the ReW-C and C-ReW interfaces. (Afer Spiller.¹⁹¹)

process in which interdiffusion or oxidation is inhibited by the deposition at the interfaces of thin barrier layers of materials such as C or B_4C .¹⁹²

Narrowband Reflection Coatings

Narrowband rejection (or notch, or minus) filters transmit freely all the radiation incident upon them except in one narrow spectral region in which the radiation is either wholly or partially reflected.¹⁹³ Lord Rayleigh observed a corresponding natural phenomenon in potassium chlorate crystals.¹⁹⁴ Subsequent experimentors reported crystals with rejection bands varying between 0.001 and 0.038 μm in width and reflectances between 33 and 99.9 percent.^{195–197} But, unfortunately, at present the size of the crystals that can be grown is insufficient, and the position and width of the rejection region cannot be controlled. The same comments can be made about coextruded polymer films made of the thermoplastic materials polypropylene ($n = 1.49$) and polycarbonate ($n = 1.59$).^{198,199} With thin films the above limitations can be overcome, but until recently it has been difficult to achieve the extremely narrow widths and high rejections observed in the crystals.

A quarter-wave stack of the type $[AB]^N A$ has been suggested as a model for the construction of such filters (see, for example, Refs. 196 and 200). It follows from Figs. 31 and 34 that the closer the refractive index ratio n_A/n_B is to unity, the narrower the width of the reflectance zone and the more layers required to achieve a certain rejection. Shown in Fig. 42a is $(1 - \text{the measured reflectance})$ of a multilayer consisting of 760 layers that was produced by a plasma chemical deposition technique.²⁰¹ Unfortunately, that method could only deposit such coatings on the inside of a tube. To reduce the number of layers, films with higher n_A/n_B ratios could be utilized and the width reduced by using layers with thicknesses that are odd multiples of a quarter wavelength. But this can be done only at the expense of bringing the adjacent higher- and lower-order reflection peaks closer. Resonant reflectors are an extreme example of this. However, recent progress in the manufacture of multilayer systems (see Sec. 7.12) has made it possible to produce notch filters for Raman spectroscopy at any wavelength in the visible or near-infrared that exceed the reported performance of the naturally occurring crystals. Fig. 42b shows the measured transmittance and optical density of a filter that has a 0.020- μm -wide rejection region and a very high transmittance with only a 1 to 2 percent ripple in the wavelength range $0.6 < \lambda < 1.00 \mu\text{m}$. The optical density at $\lambda = 0.75$ is 6.0. The filter consists of 714 layers and it has a metric thickness of 51 μm .²⁰² The measured width of the notch of the filter shown in Fig. 42c is about 0.024 μm and it has an optical density of 5.9 at $\lambda = 1.064 \mu\text{m}$, as measured with a laser. The second surface of this filter was not antireflection coated, and so the transmittance is a little lower, but the high transmittance extends from 0.4 to 1.09 μm . The metric thickness of this filter is 122.2 μm and it consists of 4410 layers.²⁰³

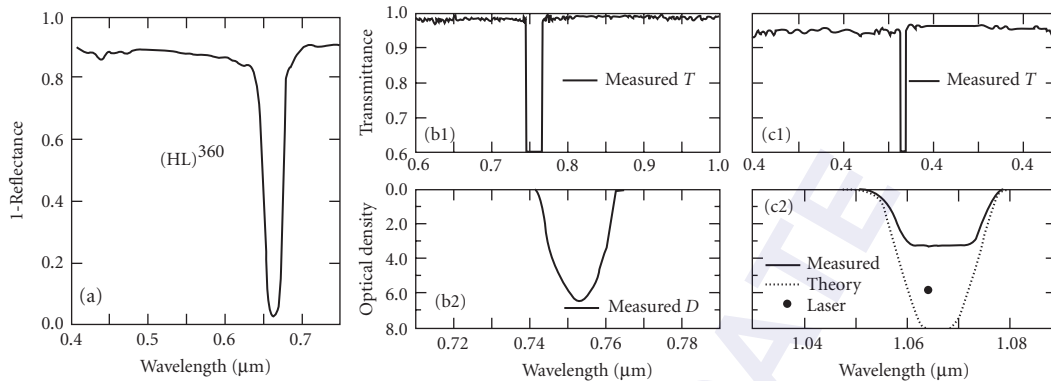


FIGURE 42 Narrow band rejection filters. (a) Measured value of $(1 - \text{reflectance})$ of a 760-layer quarterwave stack produced by chemical vapor deposition using materials with refractive indices of 1.575, 1.585 (after Edmonds *et al.*²⁰¹); (b) measured reflectance and optical density a 714-layer, 51- μm -thick filter (after Iridian Spectral Technologies²⁰²); and (c) measured reflectance and optical density a 4410-layer, 122.2- μm -thick filter (after JDSU²⁰³).

Resonant Reflectors

In the past, even “hard” evaporated coatings could not survive the power densities present in some high-power lasers. To overcome this problem, resonant reflectors were used.^{204–207} They consisted of one or more accurately air-spaced plane-parallel plates of thicknesses of the order of millimeters made of a tough, high-optical-quality material. Because of the long coherence length of the laser radiation incident upon them, interference takes place in the same way as in thin films. Resonant reflectors may be regarded as being quarter-wave reflectors of enormously high order of interference, and all the equations given in Sec. 7.6 apply.

Resonant reflectors made of quartz and sapphire are commercially available. Since the refractive index of quartz is lower than that of sapphire, a larger number of plates is required to attain the same reflectance. On the other hand, quartz is much cheaper and is less temperature sensitive. The calculated reflectance of one-, two-, three-, and four-plate sapphire resonant reflectors are shown in Fig. 43. In another development diffusion-doped quartz plates are used in the construction of resonant

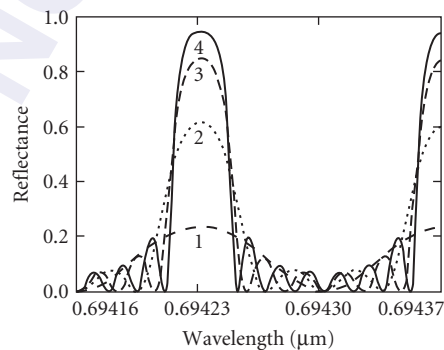


FIGURE 43 Calculated spectral reflectance of resonant reflectors consisting of one, two, three, and four sapphire plates ($n = 1.7$). The optical thicknesses of the plates and of the air spaces between them were assumed to be 1.7 mm.

reflectors.²⁰⁸ The doping process causes the refractive index of the plate to increase smoothly in a 0.5- μm region from that of quartz to about 2.0 at the surface. As a consequence fewer elements are needed to achieve a given reflectance.

All-Dielectric Broadband Reflectors

There are several ways of obtaining a coating with a broad high-reflection region should the width attainable with quarter-wave stacks be inadequate. For a broad region with a very high reflectance one can superimpose two quarter-wave stacks tuned to two different wavelengths. The widest continuous high-reflection region is attained when materials with the highest available refractive-index ratios are used and the thicknesses of the layers are so chosen that the two high-reflection zones are contiguous (see Fig. 44c). For an even broader region further quarter-wave stacks can be superimposed. If high-reflection regions overlap, special precautions must be taken to prevent the appearance of sharp reflectance minima in the high-reflection region.¹⁷³ It is not easy to obtain a very uniform, moderately high reflectance in this way. Another approach is to deposit onto the substrate a series of alternating high- and low-refractive-index films of gradually increasing or decreasing thicknesses (Fig. 45). A broad high-reflection region can also be obtained with a multilayer in which the layers are different multiples of $\lambda/4$ of a selected wavelength (see Fig. 46). Finally, it has been suggested that multilayers with 10:1 high reflectance regions might be produced by depositing hundreds of layers of random thicknesses made of two materials that are nonabsorbing throughout the spectral range of interest.²⁰⁹ As yet experimental data for this approach have not been presented.

If only a relatively small increase in the high-reflection region is required, or if a uniform but only moderately high reflectance is required, the desired result can be achieved by modifying the thicknesses and refractive indexes of a quarter-wave stack using a computer refinement program or by the addition of achromatizing $\lambda/2$ layers. The measured performances of several such reflectors are shown in Figs. 47 and 48. Broadband reflectors with moderate and high reflectances for the ultraviolet spectral region have been reported by Korolev,²¹³ Sokolova,²¹⁴ and Stolov.²¹⁵ A broadband reflector for the XUV spectral region is shown in Fig. 57b.

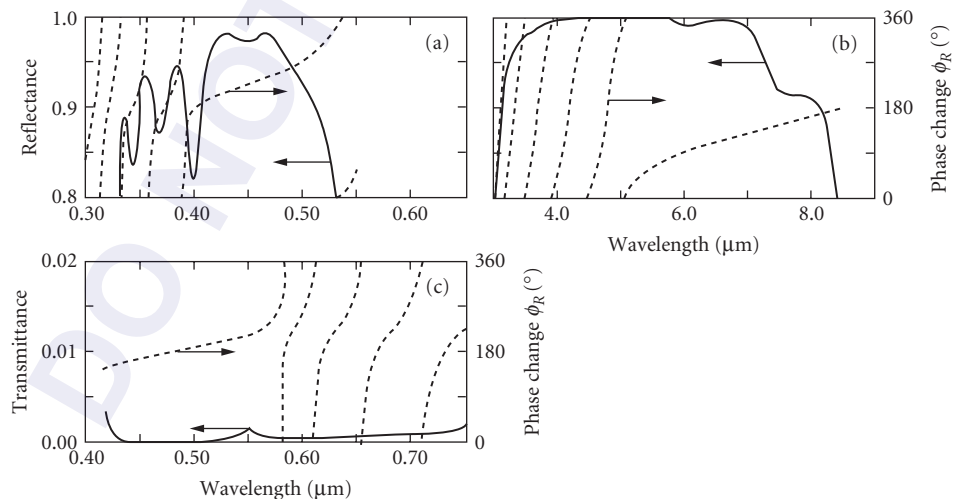


FIGURE 44 Measured reflectance of all-dielectric broadband reflectors consisting of two superimposed quarter-wave stacks with contiguous high-reflection zones. The dotted curves in Figs. 35 to 39 represent the calculated phase changes on reflection. [(a) and (b) after Turner and Baumeister,¹⁷³ (c) after Perry.²¹⁰]

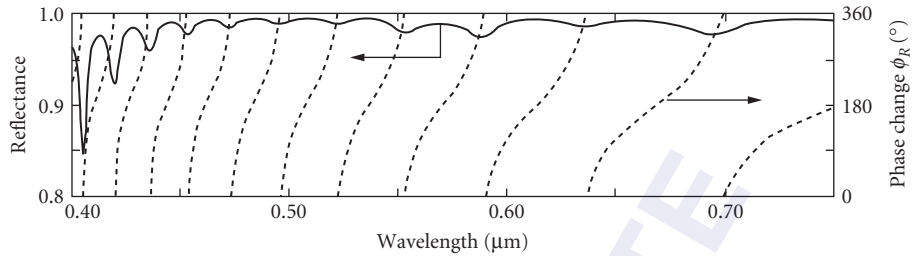


FIGURE 45 Calculated reflectance of an all-dielectric broadband reflector consisting of 35 layers made of a low- and a high-index material and having optical thicknesses that vary in a geometric progression. (After Heavens and Liddell.²¹¹)

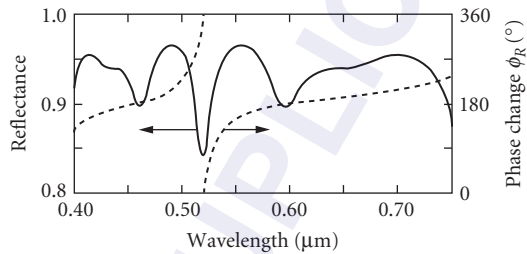


FIGURE 46 Calculated reflectance of an all-dielectric broadband reflector consisting of 11 layers all of which have optical thicknesses that are various multiples of $0.13 \mu\text{m}$. (After Elsner.²¹²)

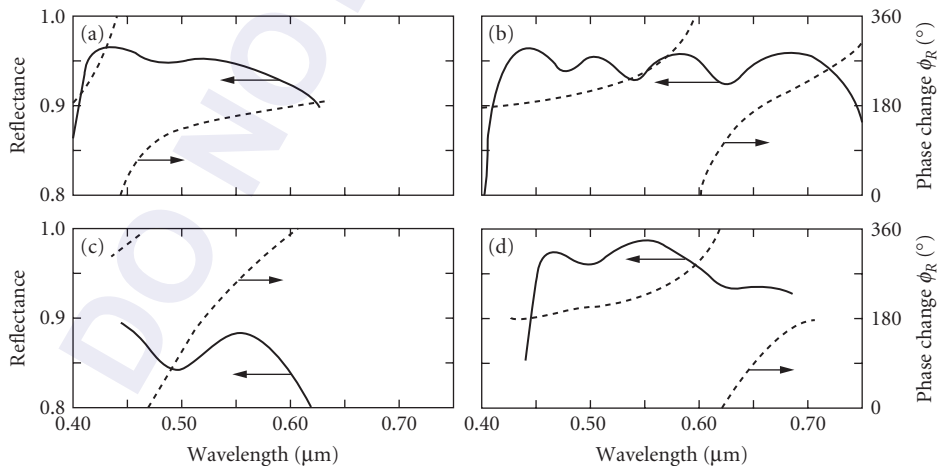


FIGURE 47 Measured reflectances of all-dielectric broadband reflectors designed with refinement programs. [(a) after Penselin and Steudel,²¹⁶ (b) after Baumeister and Stone,²¹⁷ (c) after Ciddor,²¹⁸ and (d) after Ramsay and Ciddor.¹⁷⁸]

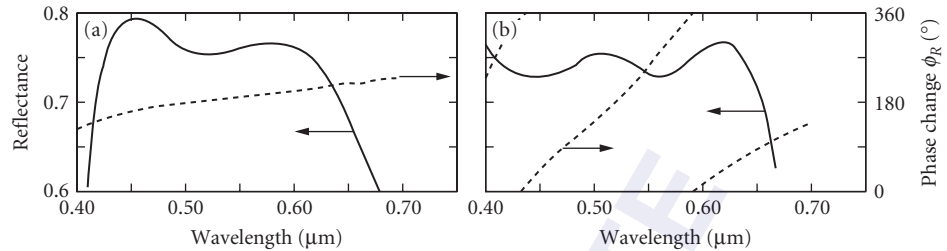


FIGURE 48 Measured reflectances of achromatic all-dielectric reflectors with an intermediate reflectance derived from a quarter-wave stack through (a) the addition of half-wavelength achromatizing layers (after Turner²¹⁹) and (b) designed with a refinement program (after Ciddor²¹⁸).

Broadband Reflectors Effective over a Wide Range of Angles (Perfect Mirrors) At times it is required that the reflectance be high not only over a broad wavelength region, but also for a large range of angles of incidence. It was clear that if this is to be achieved, the superimposed periodic stacks (Fig. 44) must be contiguous at the highest angle required, and that they will therefore overlap at lower angles. This situation has been investigated at some length in a theoretical paper for visible reflectors effective for angles up to 50°. Subsequently a number of workers outside the thin-film field rediscovered this possibility and described similar coatings working for all angles, which they called “perfect mirrors.”

Phase Change on Reflection from Broadband Reflectors The phase change on reflection from broadband reflectors varies even more rapidly with wavelength (Figs. 44 to 47) than that from quarter-wave stacks (Fig. 29). In fact, the phase change can be so large, that it has been proposed for the construction of spacerless transmission interference filters (see Fig. 93h). But a high phase change on reflection can also be a disadvantage in certain metrological applications. Another consequence of this rapid variation is that, in the presence of systematic thickness variations in the layers, it can give an impression of a lack of flatness in the substrate. Figures 47d and 48a represent designs of broadband reflectors in which an effort was made to reduce this effect.

Coatings for Ultrafast Optics

In all sections of this chapter, with the exception of this one, it is assumed that the intensity of the radiation that is falling onto a multilayer does not vary in time and that a steady state exists. However, there are applications in telecommunications, spectroscopy, and in other areas in which the use of very short pulses of light is a great advantage. But, for pulse durations shorter than 100 fs, the transient effects that can be ignored during steady-state operations become significant. In ultrafast optics the incident light can consist of very short Gaussian light pulses with a temporal bandwidth as short as 5 fs. Under these conditions the phase changes ϕ_R , ϕ_T of a multilayer on reflection or transmission, and their first and second derivatives with respect to ω , where $\omega = 2\pi c/\lambda$ and where c is the speed of light in a vacuum, become very important. The first and second derivatives of the phase changes, $-(d\phi/d\omega)$ and $-(d^2\phi/d\omega^2)$, are called the Group Delay (GD) and Group Delay Dispersion (GDD), respectively. Unless they are suitably controlled in a multilayer, the pulse length of the incident beam can be considerably broadened on reflection. In order to avoid this, the thicknesses of layer pairs in a coating are reduced in a systematic way, so that light of shorter wavelengths does not penetrate as deeply into the multilayer structure before a significant part of it is reflected. Such structures resemble those used in broadband reflectors (see, for example, Fig. 45) and are called by some “chirped multilayers.” At times it is useful to use “double chirped” structures in which, in addition to a gradual reduction of the bilayer thickness, the relative thicknesses of the high and low refractive index layers

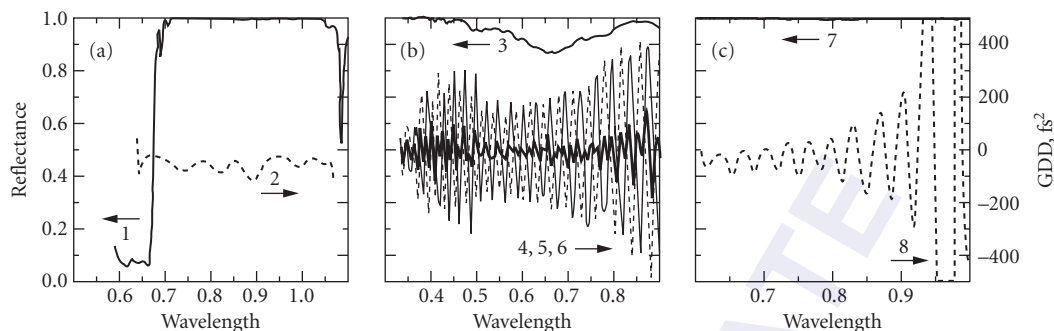


FIGURE 49 Measured reflectances R and the Group Delay Dispersions (GDD) of three different systems designed for femtosecond use. (a) Input coupler.²²⁸ (b) The large fluctuations of the GDD values that occur in mirrors designed for a very wide spectral regions can be reduced by reflecting the light from a pair of mirrors designed to have oscillations that are complementary. In this diagram, the dotted and thin solid curves correspond to the individual GDD curves; the heavy solid curve corresponds to the GDD after reflection from both mirrors.²²⁹ (c) Dispersive mirror designed by the time-domain method.²³⁰

within each pair are also varied. Of course, more degrees of freedom to obtain the required performance are obtained if in the final calculations the thicknesses of all layers of the system are allowed to vary independently.

It has been shown that the temporal and spectral widths of the reflected pulses are related to one another through a Fourier transformation that uses a product of the complex spectrum of the incident pulse with the factor $r_i(\omega)\exp[\phi_i(\omega)]$ for the reflecting coating (see, for example, Refs. 226 and 227). In the numerical design of such coatings, R , GD , and GDD , or even the desired temporal pulse length, may need to be specified and they will depend on the particular problem on hand. A more detailed presentation of the design of such coatings is outside the scope of this chapter and the interested reader is referred to the reviews by MacLeod² and Tempea et al.²²⁶

Some of the multilayer reflector coating types that are required for ultrafast optics include input couplers (with a high R and a negative GDD and a high T in the gain and absorption region of the laser); output couplers (with a certain GDD and a rather low transmittance $0.05 < T < 0.3$ in a wide spectral range of the laser oscillator); chirped mirrors with very large bandwidths. The measured performance of three reflecting coatings for femtosecond applications are shown in Fig. 49.

It should be stated that this is currently a very fast moving field. Tunable dispersion compensators^{202,231} for femtosecond use are currently employed.²³² Compressors and stretchers for nonlinear optics amplifiers with negative and positive GDD as large as of 2000 fs^2 are being developed.²³³ In addition to reflectors, the design of beam splitters²³⁴ and antireflection coatings²³⁵ for ultrafast optics have also been discussed in the literature.

Rejection Filters

Minus Filters Minus filters are, in essence, multilayer reflectors in which the ripples in the transmission regions on either side of the high reflectance zone have been reduced or eliminated. Filters of this kind with various widths and attenuations find applications as correction filters.²³⁶ Narrow minus filters with high attenuations have various scientific and technological uses, including protection of equipment and personnel from harmful laser radiation. Figure 50 shows the measured transmittances of a number of rejection filters of various widths and attenuations. Quite recently spectacular progress has been made in the deposition of narrowband rejection filters that find applications in Raman spectroscopy. The two commercial narrow band rejection filters shown in Fig. 42b and c consist of 713, 4410 layers each, have 0.02- and 0.025- μm -wide rejection regions and measured transmittances of less than 10^{-6} at the center of the notches.^{202,203}

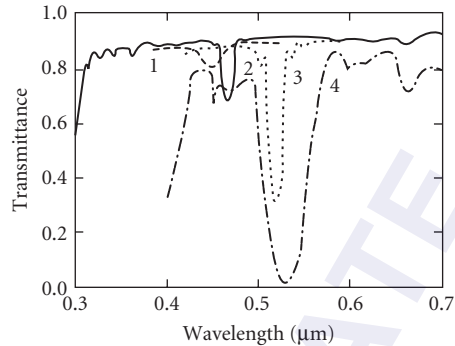


FIGURE 50 Measured spectral transmittances of several narrowband rejection filters. (Curves 1 to 3 after Dobrowolski²³⁶ and curve 4 after Optical Coating Laboratory.²³⁷)

Thelen has shown how to optimize the transmission on both sides of the rejection band of a minus filter simultaneously.²³⁸ If the multilayer is surrounded on both sides by the same medium, it will be symmetrical and can be represented by $C[AB]^NAC$, $DAC[AB]^NACAD$, . . . Here A , B , C , D , . . . are layers of quarter-wave optical thickness at the design wavelength and

$$n_m = n_s = n_A \quad (57)$$

The refractive indices n_C , n_D , . . . depend on n_A and n_B in a more complicated way. The larger the number of different materials used in the construction of these multilayers, the better the performance in the transmission region. Should the use of more than two coating materials not be convenient, it is a simple matter to find a two-material version of this solution. The above points are illustrated in Fig. 51.

Young¹⁹³ described two other design methods for narrowband rejection filters with improved transmission characteristics. The methods are based on analogies with antenna theory, and they yield nonperiodic equi-ripple designs in which all the layers either have equal thicknesses but different indexes or are made of two materials only but have many different thicknesses.

Should the simultaneous rejection of several wavelengths be required, it is possible to achieve this by depositing several minus filters on top of one another (Fig. 52).

Rugate Filters In inhomogeneous layers the refractive index varies continuously in the direction of the thickness of the layer.^{100,240,241} If the refractive index varies in a periodic manner between two extreme values, it is possible to design a minus filter with a high transmission on either side of the rejection band (Fig. 53). Such periodic inhomogeneous layers are sometimes called rugate filters. The rejection wavelength corresponds to that wavelength for which the period of the index variation is equal to a half wavelength. The attenuation depends on the ratio of the highest to lowest refractive index in the design and on the number of periods. As in the multilayer minus filters, the width of the rejection region also depends on the refractive index ratio. Rugate filters do not have the higher-order reflection peaks that are characteristic of periodic multilayers and this is one reason for their attractiveness. Another reason—no sharp interfaces, less scatter, and better mechanical properties. However, they are more difficult to produce. If necessary, they can be approximated by a homogeneous multilayer system consisting of a few (three or four) materials.

As in the case of minus filters, it is possible to reject a number of wavelengths by depositing several rugate filters on top of each other. However, the combined overall thickness of the rugate filters will then be quite high. It is possible to find an inhomogeneous layer solution to this problem in which the refractive index profile is more complicated, but which requires a considerably thinner inhomogeneous layer (Fig. 54).²⁴²

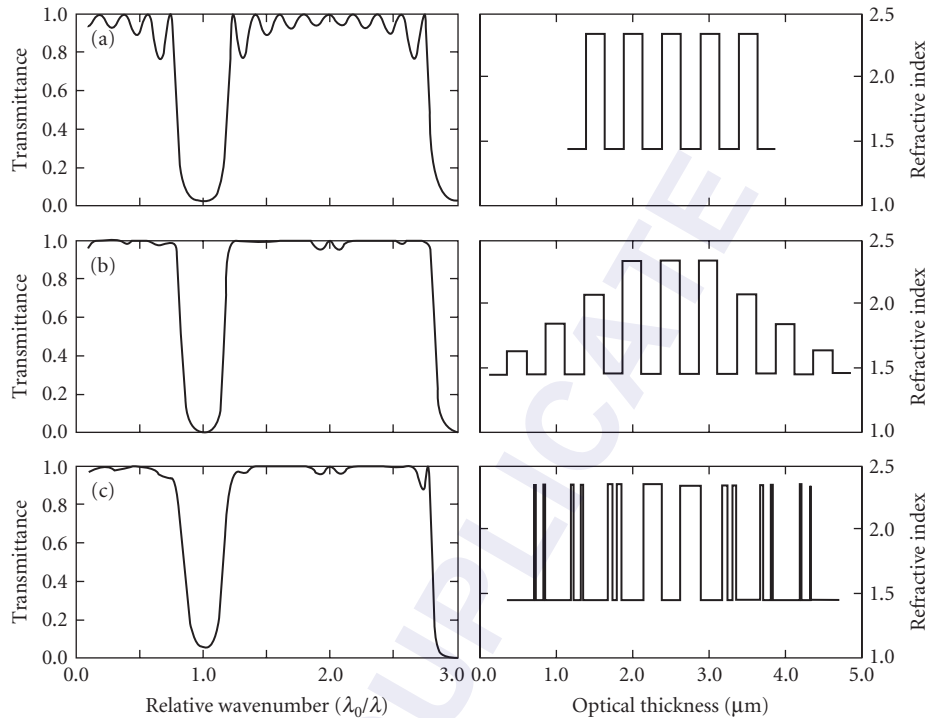


FIGURE 51 Suppression of ripples in the transmission region of rejection filters. Calculated spectral transmittance curves and refractive index profiles of a 9-layer quarter-wave stack (a); a 17-layer, 5 material minus filter (b); and a two material equivalent of the minus filter (c).

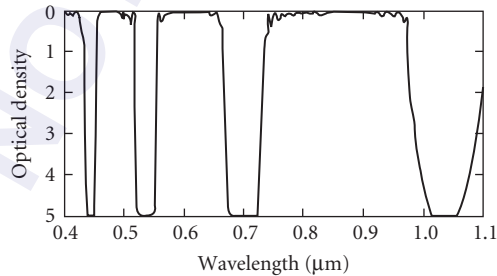


FIGURE 52 Transmission of laser goggles for the rejection of the ruby (0.694 μm) and the NdYag (0.532 and 1.064 μm) laser lines. (After Omega Optical Inc.²³⁹)

In Lippmann-Bragg holographic mirrors the refractive index varies continuously in a direction perpendicular to the plane of the substrate. These devices behave like thin-film systems and have properties similar to those of rugate filters. Holographic edge and narrowband rejection filters are available commercially.^{244,245}

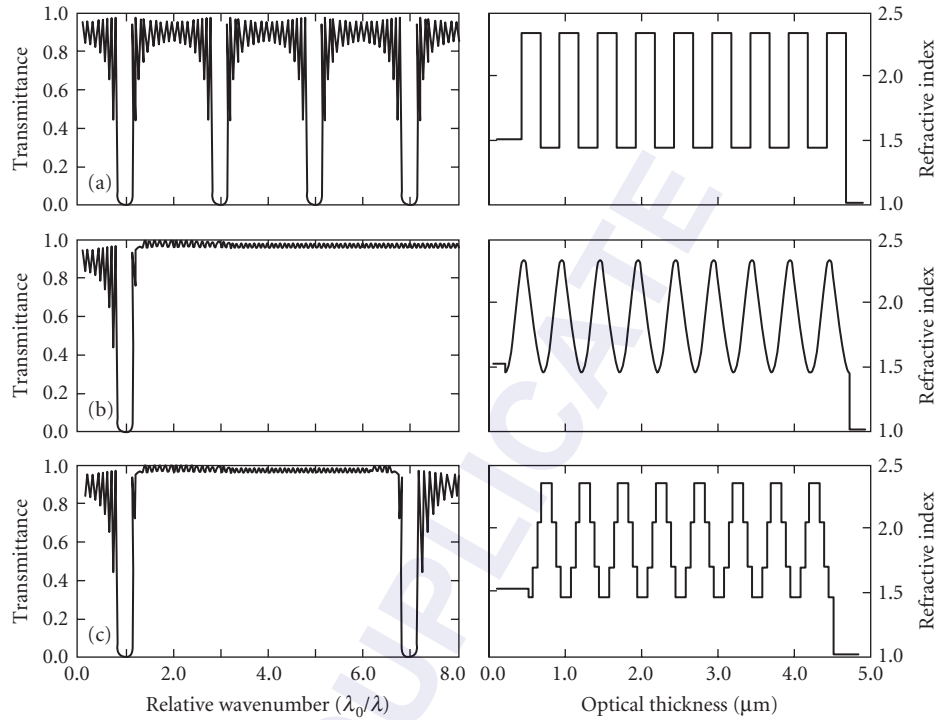


FIGURE 53 Suppression of higher order reflectance peaks in rejection filters. Calculated spectral transmittance curves and refractive index profiles of a 17-layer quarter-wave stack (a); a 9-period rugate filter (b); and a 49-layer four material design (c) by Thelen.¹³

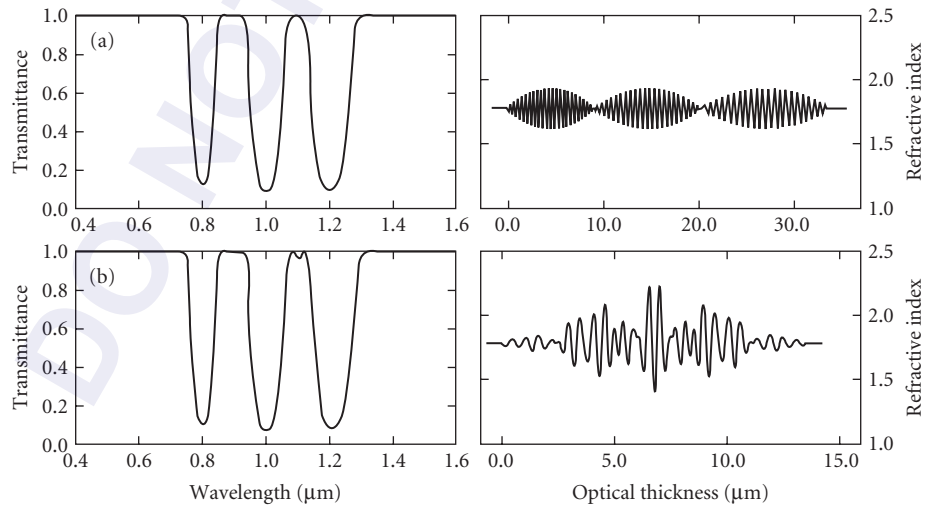


FIGURE 54 Simultaneous suppression of several laser wavelengths. Calculated spectral transmittance curves and refractive index profiles of a series (a) and a parallel (b) solution to the problem (After Verly.²⁴³)

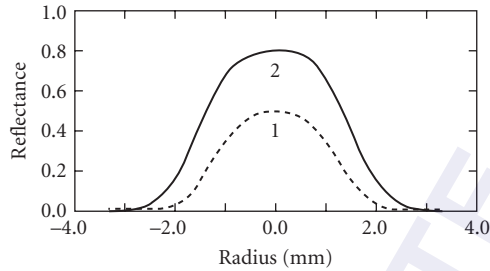


FIGURE 55 Experimental results for two super-Gaussian graded reflectivity mirrors for $\lambda = 1.064 \mu\text{m}$. Curve 1: three layer system with one layer of varying thickness, with $R_{\text{max}} = 0.5$ and $\omega = 1.92 \text{ mm}$ (after Piegari²⁵⁰) and curve 2: fully shaped mirror with $R_{\text{max}} = 0.85$ and $\omega = 1.90 \text{ mm}$ (after Duplain²⁴⁸).

Graded Reflectivity Mirrors

Mirrors in which the absorption or reflection varies radially have been proposed in the past for the control of modes and of edge diffraction effects in laser resonators.^{246,247} In addition to meeting the reflectance specifications, graded reflectivity mirrors must have a sufficiently high laser damage threshold for use with high-power lasers. Graded reflectivity mirrors are produced by depositing thin films through a suitable mask. The substrate and the mask can be stationary, or one or both can rotate. A single shaped layer suffices for a maximum reflectance of intermediate values. For higher values the shaped layer can be inserted between a stack of layers of uniform thickness or, alternatively, all the layers can be deposited through the mask.^{248,249} The experimentally measured reflections of two graded reflectivity mirrors are shown in Fig. 55.

Multilayer Reflectors for the Far-Infrared Region

Thin film-filters cannot be produced by conventional deposition techniques for wavelengths greater than about $80 \mu\text{m}$ because of a lack of low absorption coating materials that can be deposited in the form of thick, stable films. However, a hybrid process in which plastic sheets coated with relatively thin high refractive index materials are heat-bonded can be used to produce self-supporting optical multilayer filters.^{251,252,164} As mentioned before (Sec. 7.6, subsection “Periodic Multilayers of the $[xH \cdot (1-x)L]^N \cdot xH$ Type”), periodic multilayers of unequal optical thickness can have a high reflectance, providing that the number of periods is high enough. The measured reflectance curves of two typical hybrid multilayer filters are shown in Fig. 56.

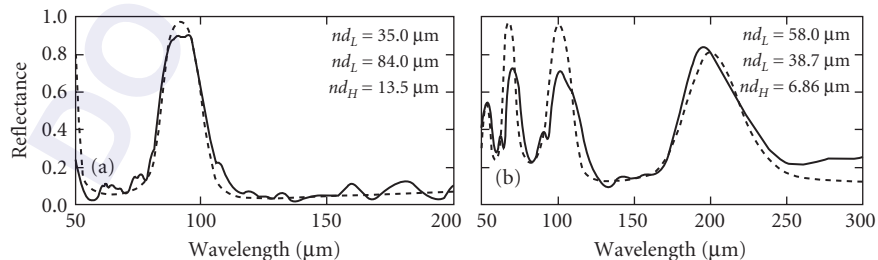


FIGURE 56 Measured and calculated performance of (a) $L'(HL)^{16}$ and (b) $L'(HL)^{16}$ reflectors made of ZnS and polyethylene of thicknesses indicated in the diagram. An instrumental resolution of 3 cm^{-1} was assumed in the calculations. (After Shao.¹⁶⁴)

Multilayer Reflectors for the Soft X-Ray and XUV Regions

There are two main obstacles to obtaining multilayers with a high normal incidence reflectance in the soft x-ray and in the XUV regions. First, at these wavelengths all materials absorb; this limits the number of layers that can contribute to the overall reflectance. Second, roughness and the interdiffusion and alloying of the materials degrade the individual interfaces; this reduces their contribution to the overall reflectance.¹⁸⁹

XUV mirrors are normally produced from two chemically compatible materials by sputtering or by electron beam gun evaporation. As already discussed in Sec. 7.6, subsection “Multilayer Reflectors Made of Absorbing Materials” the XUV multilayers mirrors have a period of optical thickness $\lambda/2$. Within this period, the thickness of the less absorbing material is larger. This material usually has an absorption edge that is close to the design wavelength. Optical constants can vary widely in this region.¹⁶⁷ The second more absorbing material is therefore selected to maximize the Fresnel reflection coefficient of the interface.

Theoretically, the best results are obtained with pure elements. However, sometimes alloys are used because they result in multilayers with better interfaces. At other times, MoSi_2 might be used in place of Mo, or B_4C instead of C or B. Examples of some more commonly used material pairs for near-normal incidence reflectors are (a) Mo/Si for the 130 to 250 Å region; (b) Mo/Y for the 90 to 130 Å region; (c) W/ B_4C , Ru/ B_4C , Mo/ B_4C , etc. for the 70 to 130 Å region; and (d) Co/C, W/C, Re/C, ReW/C, Ni/C, etc. for the 45 to 70 Å. In the above, the second material in each pair has the lower extinction coefficient. A much more complete review of this topic will be found in the review article by Windt et al.²⁵³ Also, an up-to-date list of material pairs and the measured peak reflectances achieved with them will be found on Henke’s website.²⁵⁴

There are many reasons why there are differences between the theoretical performance of a multilayer reflector and the reflectance measured on a synchrotron.^{165,166} The highest near-normal reflectance achieved at 135 Å thus far is $R \approx 0.70$ for an interface-engineered Mo/Si multilayer.²⁵⁵ The measured reflectances of other experimentally produced x-ray and XUV mirrors are shown in Fig. 57. Curve 8 in this diagram shows that, as in the visible part of the spectrum (Figs. 44 to 48), by departing from periodic multilayer systems, it is possible to produce reflectors in this spectral region that have a comparatively wide spectral region with a constant reflectance.

In the XUV and soft x-ray regions, the reflectances that can be achieved with near-normal angles of incidence are not sufficient. For some imaging applications it is necessary to resort to near-grazing incidence mirrors. In Sec. 7.16 the reflectances of some metals, compounds, and two- or three-layer combinations at high angles, as well as periodic and nonperiodic multilayer interference reflectors consisting of tens or even thousands of layers, are presented.

For more information on coatings for the XUV and soft x-ray region, the interested reader is referred to the review article by Yulin.²⁶⁶

7.8 CUTOFF, HEAT-CONTROL, AND SOLAR-CELL COVER FILTERS

Ideal cutoff filters would reject all the radiation below, and transmit all that above a certain wavelength, or vice versa. Real cutoff filters, of course, are not perfect and so, in addition to the cutoff wavelength, the slope of the transition region and the extent and average transmission values of the transmission and rejection regions must be specified. The tolerable departures of these quantities from the ideal values depend greatly on the application. Most all-dielectric cutoff filters are based on periodic multilayers, whose basic properties were described in Sec. 7.6.

Transmission in the Passband Region

The usual way of avoiding the secondary transmission minima in the transmission band of a quarter-wave stack is through the use of eighth-wave layers on both sides of the stack (Sec. 7.6, subsection “Periodic Multilayers of the $[(0.5A) B (0.5A)]^N$ Type”).

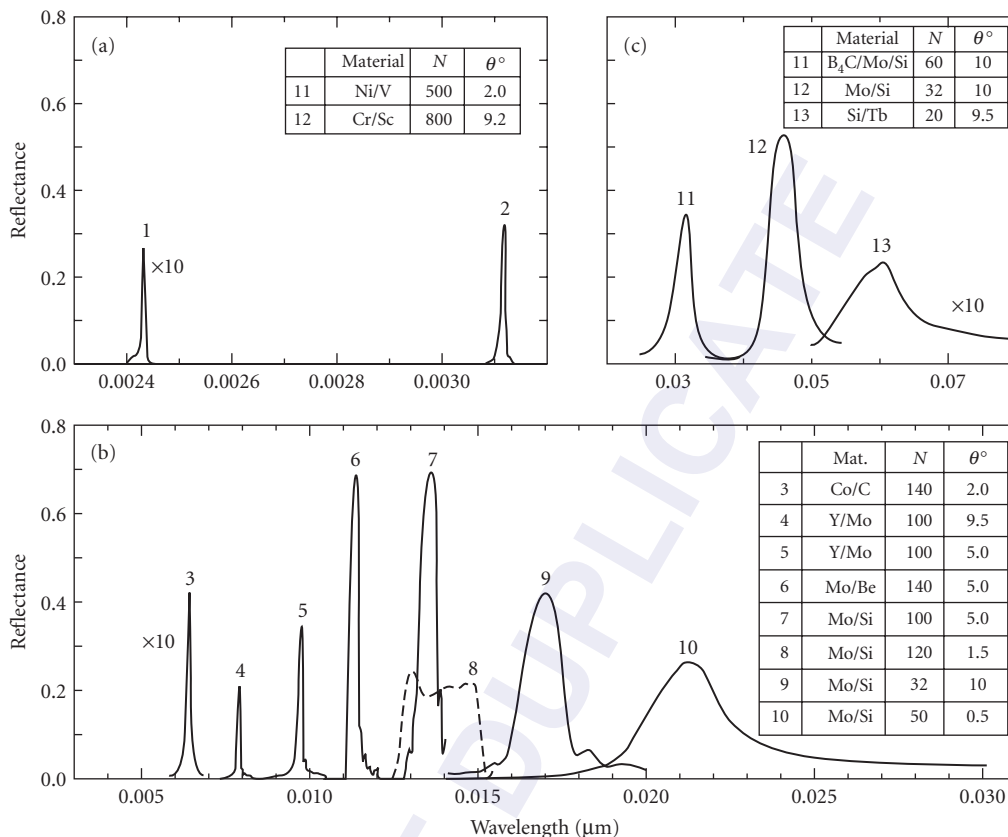


FIGURE 57 Measured reflectance curves of some representative multilayer x-ray mirrors for the (a) 0.0022 to 0.0032 μm ; (b) 0.003 to 0.0029 μm ; and (c) 0.03 to 0.08 μm spectral regions. The materials used, the number of layers N and the angles θ_0 at which the near-normal reflectance was measured are indicated in the tables. (Curve 1 after Eriksson;²⁵⁶ curve 2 after Gullikson;²⁵⁷ curve 3 after Spiller;²⁵⁸ curves 4 and 5 after Montcalm;²⁵⁹ curve 6 after Skulina;²⁶⁰ curve 7 after Bajt;²⁵⁵ curves 8 and 12 after Kaiser;²⁶¹ curve 9 after Ravet;²⁶² curve 10 after Ceglie;²⁶³ curve 11 after Gautier;²⁶⁴ and curve 13 after Kjornrattanawanich.²⁶⁵)

Other, less frequently used methods of smoothing the transmission in the passband are the adjustment of the thicknesses of all the layers of a quarter-wave stack,²⁶⁷ the use of homogeneous²⁶⁸ and inhomogeneous²⁶⁹ layers on either side of the stack; the choice of an optimum set of refractive index values for the substrate and films;²⁷⁰ and the use of an equi-ripple design¹⁹³ in which thicknesses are kept constant but the refractive indices are varied.

The Width of the Transmission Region

For short-wavelength cutoff filters of the type described above, the transmission region is limited only by the transmission characteristics of the materials used for their construction. In long-wavelength cutoff filters the transmission regions can be limited by the appearance of higher-order reflectance maxima (Sec. 7.6, subsection “Nonabsorbing $[AB]^N$ and $[AB]^N A$ multilayer Types”). Should this be a serious limitation, it is possible to suppress a number of adjacent reflectance maxima by using multilayers with periods composed of three or more different materials (Fig. 58).

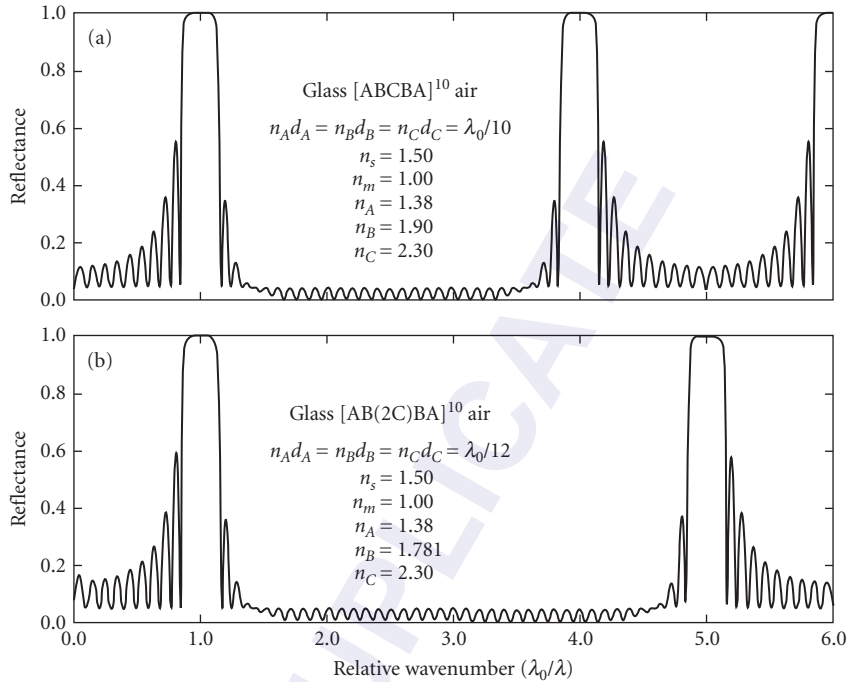


FIGURE 58 Calculated reflectances (a), (b) of two different of three-material periodic multilayers with suppressed higher-order reflectance maxima. (After Thelen.²⁷¹)

By using a period consisting of an inhomogeneous layer with a special refractive-index profile an even larger number of consecutive reflectance maxima can be suppressed (see also Fig. 53).²⁶⁹ Figure 59 shows the measured results for two such experimental coatings in which reflectance maxima are suppressed at three and nine consecutive integer multiples of $1/\lambda_0$, a fact obscured in the case of the second filter by the absorption of the materials used.

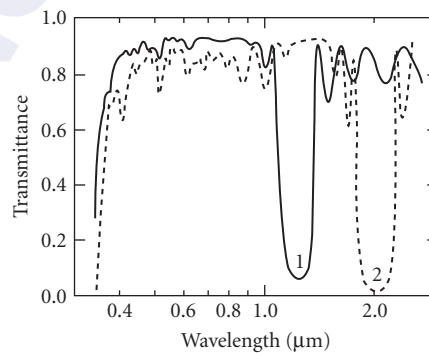


FIGURE 59 Measured transmittance of periodic multilayers in which higher-order reflectance maxima are suppressed through the use of periods that consist of an inhomogeneous layer. (After Scheurman.¹¹¹)

Transmission in the Rejection Region

Figure 31 can be used for an estimate of the number of layers required to achieve a given transmission (Sec. 7.6, subsection “Maximum Reflectance”). Through the use of suitable substrate materials the transmittances throughout the rejection region can typically be below 0.01 and 0.1 percent for short- and long-wavelength cutoff filters, respectively. Rejection filters with higher rejections can be provided or, alternatively, two or more filters in series can be used if they are placed at a small angle to one another.

The Width of the Rejection Region

The width of the high-reflection region of $[(0.5A)B(0.5A)]^N$ coatings can be estimated from Fig. 34. There is no shortage of absorbing materials should it be necessary to extend the rejection region of shortwave cutoff filters. The number of suitable absorbing materials for long-wavelength cutoff filters is more restricted and often it is necessary to use thin films to extend the cutoff region. In addition to the superposition of two or more cutoff filters tuned to different wavelengths (Sec. 7.7, subsection “Rejection Filters”) one can deposit such coatings onto different substrates or onto the opposite sides of the same substrate.²⁷² The resulting transmission will be governed by the considerations of Sec. 7.1, subsection “Transmission Filters in Series and Parallel.”

Slope of the Cutoff

This quantity is defined in a number of ways. One common definition is

$$\left| \frac{\lambda_{0.8} - \lambda_{0.05}}{\lambda_{0.5}} \right| \times 100\% \quad (58)$$

where $\lambda_{0.8}$, $\lambda_{0.5}$ and $\lambda_{0.05}$ refer to the wavelengths at which the transmittances are 0.8, 0.5, and 0.05 of the maximum transmittance of the filter. Explicit formulas for the slope are complicated.¹⁶⁰ The slope increases with the number of periods and with the refractive index ratio. Slopes with values of the order of 5 percent are readily available in practice.

Angle-of-Incidence Effects

The edges of cutoff filters move toward shorter wavelengths as the angle of incidence is increased. The use of higher-index materials reduces the effect. Measured results for a cutoff filter in which the shift was reduced by using high-refractive-index layers that have 3 times the thickness of the low-index layers are shown in Fig. 60a. See also Sec. 7.9, subsection “Nonpolarizing Beam Splitters” on polarization independent color-selective beam splitters.

Experimental Results

The spectral transmittance curves of a number of commercially available short- and long-wavelength cutoff filters are shown in Figs. 61 and 62. High performance long- and short-wavelength cutoff filters are depicted in Figs. 63, 104, and 115. Similar filters for intermediate wavelengths can, of course, be constructed. It is also possible to construct edge filters in which the thicknesses of all the layers vary in proportion around the circumference of the substrate (see also Sec. 7.11, subsection “Linear and Circular Wedge Filters”). A tuning of the cutoff wavelength is thus possible.

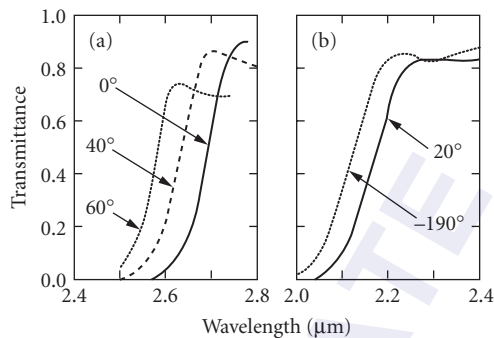


FIGURE 60 Effect of (a) angle of incidence and (b) temperature on the performance of cutoff filters. (*Optical Coating Laboratory*.²⁷³)

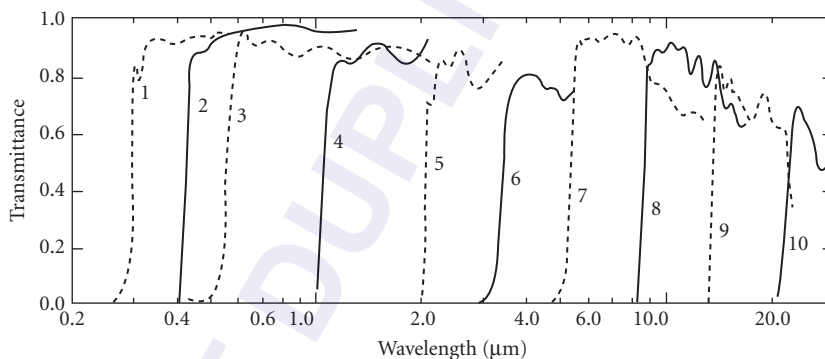


FIGURE 61 A series of commercial short-wavelength cutoff filters. (Curves 1, 5, 8, 9, and 10 after *Optical Coating Laboratory*;^{237,274} curve 2 after *Bausch & Lomb*;²⁷⁵ curves 3 and 6 after *Turner*;⁸³ curve 4 after *Eastman Kodak*;²⁷⁶ and curve 5 after *Infrared Industries*.²⁷⁷)

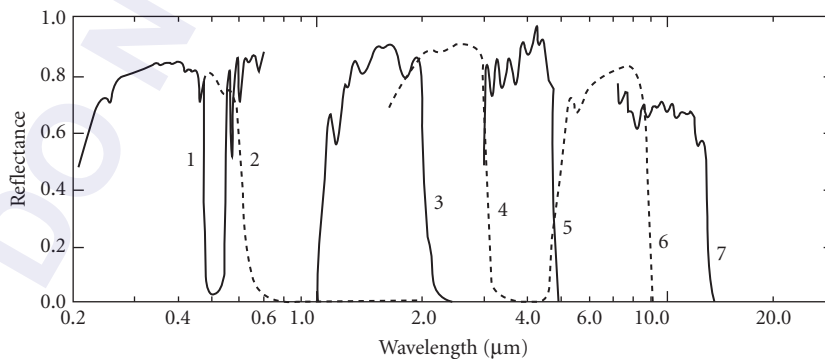


FIGURE 62 A series of commercial long-wavelength cutoff filters. (Curve 1 after *Apfel*;¹⁷⁰ curve 2 after *Eastman Kodak*;²⁷⁶ curves 3, 4, and 7 after *Optical Coating Laboratory*;²⁷⁴ and curves 5 and 6 after *Infrared Industries*.²⁷⁷)

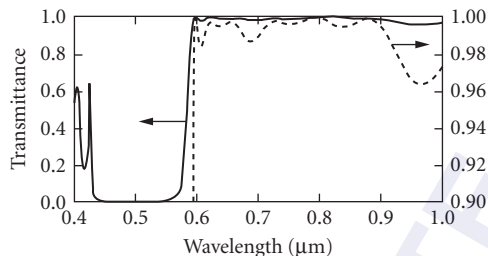


FIGURE 63 Measured transmittance of an unblocked short-wavelength cutoff filter. (After *Thin Film Lab*.⁸⁵)

Heat Reflectors, Cold Mirrors, and Infrared Suppressing Filters

Only 39 percent of the total radiation from carbon arcs and 13 percent from tungsten lamps operated at 3250 K represent visible light. Most of the remaining energy is infrared radiation, which is converted into heat on absorption. The use of heat reflectors and cold mirrors in film projectors,²⁷⁸ in spot lamps for television and film studios,²⁷⁹ and other optical instruments can lead to a very significant reduction of this unwanted heat.

Heat reflectors (also called hot mirrors) are special long-wavelength cutoff filters with a cutoff at 0.7 μm which transmit the visible radiation from 0.4 to 0.7 μm without disturbing the color balance. The width of the rejection region depends on the light source to be used and on whether a heat absorbing glass is also to be used. The spectral-transmittance curves of three typical commercial heat reflectors are shown in Fig. 64a. The measured spectral transmittance and reflectance curves of two heat-reflecting coatings not based on periodic multilayers are shown in Fig. 65.

Cold mirrors reflect as much as possible of the visible light incident upon them and transmit the remaining radiation. The reflectance curves of two commercial cold mirrors are shown in Fig. 64b.

Solar-Cell Covers

Solar-cell covers remove the incident solar energy that does not contribute to the electrical output of the cell and protect it from possible deterioration of its performance through the action of ultraviolet radiation.^{283,284} The spectral transmittance of a blue-red solar-cell cover is shown in Fig. 64b. The earlier blue solar-cell covers (curve 2, Fig. 61) protected the cell only from the adverse effects of ultraviolet radiation.

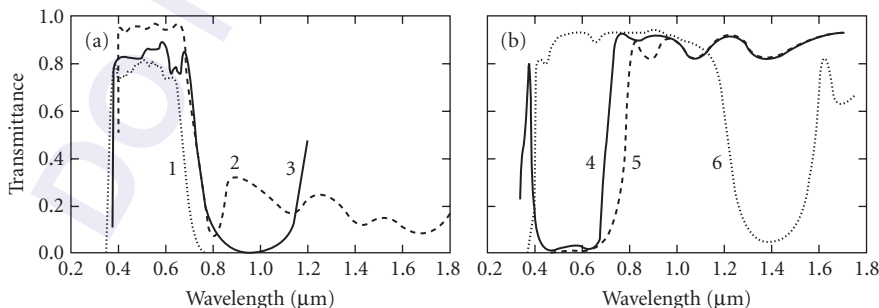


FIGURE 64 Measured performance of commercial multilayer coatings for heat control: (a) heat-reflecting coatings; (b) cold mirrors and blue-red solar-cell cover. (Curve 1 after Corion Corporation,²⁸⁰ curve 2 after Bausch & Lomb;²⁷⁵ curve 3 after Balzers;²⁷⁶ curves 4 and 6 after Optical Coating Laboratory;²⁵⁷ and curve 5 after Heliotek.²⁸¹)

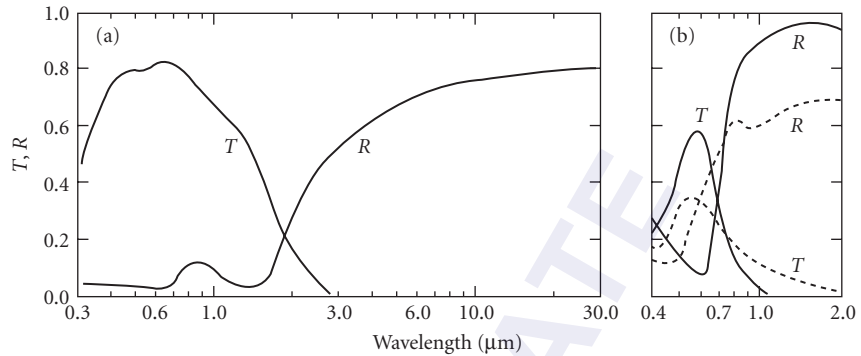


FIGURE 65 Two nonperiodic coatings with heat-reflecting properties: (a) spray-deposited conducting coating of tin oxide and (b) transmittance and reflectance of a gold film (dotted curves) and of a broadband metal dielectric filter (solid curves) with semitransparent gold films of the same total thickness. (After Turner.²⁸²)

Temperature Effects

Refractive indexes of optical materials increase almost linearly with increase in temperature thus causing cutoff edges to move toward longer wavelengths. In actual filters the fractional-wavelength shift varies between 3×10^{-3} and $10^{-4}/^{\circ}\text{C}$. Ion-plated films have a smaller temperature shift than films prepared by conventional e-beam evaporated layers.²⁸⁵ Higher-index materials tend to be more temperature sensitive, thus making it difficult to construct filters that are insensitive both to angle of incidence and temperature changes.²⁸⁶ The measured performance of a cutoff filter at two temperatures is shown in Fig. 60b.

Metal-Dielectric Reflection Cutoff Filters

It is possible to construct metal dielectric cutoff filters that act in reflected light (Fig. 66). Short-wavelength cutoff filters consist of an opaque metal layer and one or more additional layers. The light is removed through absorption within the absorbing layers of the system. The thicknesses of the individual layers are adjusted to maximize the absorption and width of the rejection region. Long-wavelength cutoff filters consist of all-dielectric multilayer reflectors superimposed

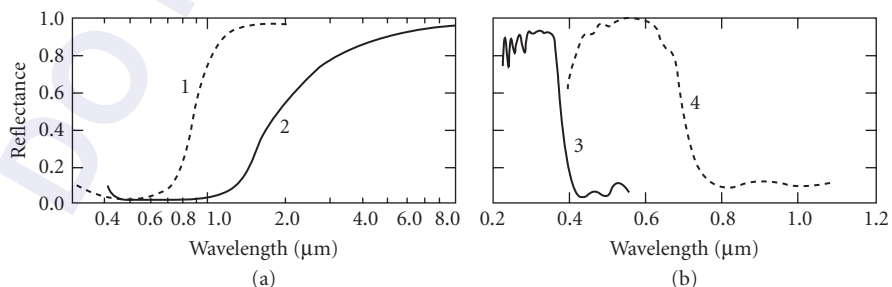


FIGURE 66 Short- and long-wavelength metal-dielectric reflection cutoff filters based on reflection. (a) Curves 1 and 2: three-layer coatings on aluminum (after Drummer and Hass¹³⁸) and (b) curves 4 and 5: multilayer reflecting coatings deposited onto black absorber coatings (after Hoppert²⁸⁷).

onto a black absorbing coating. A high attenuation of the unwanted radiation can be achieved by placing identical filters of either type in a multiple reflection arrangement of the kind depicted in Fig. 3.

Cutoff Filters Based on Absorption

All materials used in multilayer interference coatings possess short- and long-wavelength absorption edges. These can often be used to assist in the blocking of such filters. The admixture of small amounts of absorbing materials to an evaporant or organic coating solution is used at times to tune this absorption edge (Fig. 67*a* and *c*). For example, antireflection coatings containing such ultraviolet-absorbing materials can be used to protect works of art.²⁸⁸

A series of commercial short-wavelength cutoff filters for the infrared spectral region consisting of chemically deposited silver sulfide coatings on silver chloride substrates are also shown in Fig. 67*b*. These filters are quite delicate. When protected with a polystyrene layer, their transmittance is reduced and sharp absorption bands appear. The filters should not be used outdoors unless additionally protected, nor should they be exposed to ultraviolet radiation or temperatures in excess of 110°C.

The greatest advantage of cutoff filters based on absorption in thin films is their much smaller angular dependence.

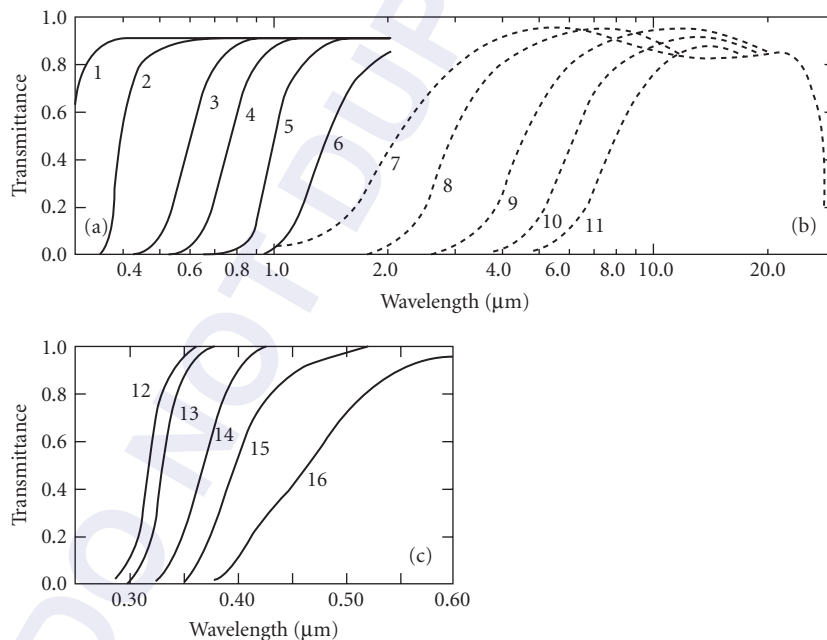


FIGURE 67 Spectral transmittance of absorbing films produced in various ways. (a) Envelopes of the transmission maxima of thick evaporated films of ZnS (curve 2), Ge (curve 6), and various mixtures of ZnS and Ge (curves 3 to 5) on a glass substrate (curve 1) (after Chang²⁸⁹); (b) spectral transmittance of chemically deposited silver sulfide coatings on silver chloride substrates (curves 7 to 11) (Eastman Kodak²⁷⁶); and (c) intrinsic transmission of thin films of titanium dioxide with admixtures of heavy metal oxides, deposited from organic solutions. Curve 12: $\text{TiO}_2 + 1.5\text{SiO}_2$; curve 13: TiO_2 ; curve 14: $\text{TiO}_2 + 0.5\text{PbO}$; curve 15: $\text{TiO}_2 + 0.15\text{Fe}_2\text{O}_3$; curve 16: $\text{TiO}_2 + 5.7\text{UO}_3$ (after Schröder²⁹⁰).

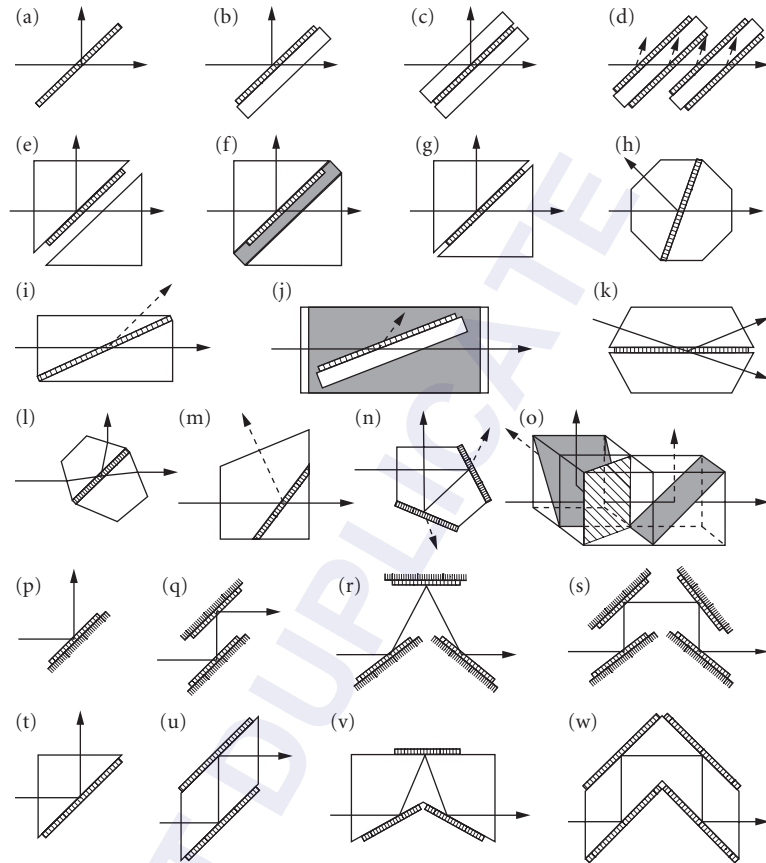


FIGURE 68 Schematic representations (*a* to *w*) of some arrangements for beam splitters, polarizers, phase retarders and multiple reflection devices. Thin films are represented in the diagrams by narrow shaded rectangles. Heavy lines ending in an arrow represent the path of the utilized radiation. Broken lines correspond to beams that are not used. The angles of incidence depend on the application.

7.9 BEAM SPLITTERS AND NEUTRAL FILTERS

Geometrical Considerations

Several different physical forms of beam splitters are illustrated in Fig. 68. The simplest beam splitter (Fig. 68*b*) consists of a coating on a transparent plane-parallel substrate. If the two derived beams are to traverse identical paths, a cemented beam splitter is used (Fig. 68*c*). The lateral displacement of the transmitted beam introduced by the above forms can be avoided with a beam splitting cube (Fig. 68*f*). To reduce the stray reflected light in the system, the free surfaces of the above beam splitters can be antireflection coated. Alternatively the coatings can be deposited onto an approximately 2- μm -thick nitrocellulose pellicle (Fig. 68*a*). The latter is an integral part of the multilayer and may introduce an interference pattern into the spectral reflectance and transmittance characteristics. Pellicle beam splitters are very light and yet quite sturdy.^{111,291,292} They are, however, subject to vibrations

caused by air currents and acoustical waves. The mechanical design of rugged, environmentally stable mounts for the above types of beam splitters have been discussed by Heinrich et al., and Lipshutz.^{293,294} Pellicles made of Mylar have been used at temperatures down to 4 K.²⁹⁵

In general, the transmission and reflection coefficients T and R will depend on the polarization of the incident radiation. The polarization effect can be reduced though the use of more complicated thin-film designs, but usually at the expense of other performance aspects. Achromatic or color-selective beam splitting arrangements have been described in which the two derived beams have intensities that are completely polarization-independent over a very wide spectral region.²⁹⁶ They consist of three identical beam splitters arranged in such a way that each beam undergoes identical reflections and transmissions on passing through the system (Fig. 68*o*).

Achromatic Beam Splitters

These devices are introduced into an incident beam of radiation when it is desired to divide it into two beams of approximately equal relative spectral composition but propagating in two different directions.²⁹⁷ In neutral beam splitters the quantity $0.5(R_p + R_s)_{\theta=45^\circ}$ is always close to the reflectance at normal incidence, even though the individual R_p and R_s values may be quite different. The reflectance of absorbing, uncemented beam splitters depends also on the direction of incidence (Sec. 7.2, subsection "Matrix Theory for the Analysis of Multilayer Systems"). The optimum values of T and R depend on the application. For example, for a binocular eye piece on a nonpolarizing microscope the most important requirement is $T_p + T_s = R_p + R_s$ (Fig. 69*a*). For a vertical illuminator ($R_p T_p + R_s T_s$) should be a maximum (Fig. 69*b*). The condition for maximum fringe contrast in some interferometers requires that $R_{1,p} T_p = R_{2,p} T_p$ and $R_{1,s} T_s = R_{2,s} T_s$ (Fig. 69*c*). This is satisfied automatically by all nonabsorbing and by absorbing cemented beam splitters. The occasional requirement that the phase change on reflection be the same for radiation incident onto the beam splitter from opposite sides is automatically satisfied at the design wavelength by all-dielectric coatings composed of $\lambda/4$ layers, but not by uncemented metal beam splitters.²⁹⁸

For maximum efficiency with unpolarized radiation R_p , R_s , T_p , and T_s should all approach 0.5. However, such a coating will not necessarily exhibit the best ratio of the intensities of the directly transmitted or reflected radiation to that which first undergoes multiple reflections.²⁹⁷ Often beam splitters are required that are uniform over a broad spectral region. Inconel films satisfy this requirement although about one-third of the incident radiation is lost through absorption (see also Sec. 7.9, subsection "Neutral Filters"). The design of achromatic all-dielectric beam splitters has been discussed by many workers.²⁹⁹⁻³⁰⁴ Knittl considered the design of beam splitters in which both the reflectance and the phase change on reflection are achromatized.³⁰⁵ The measured performance of several beam splitters is shown in Figs. 70 and 71.

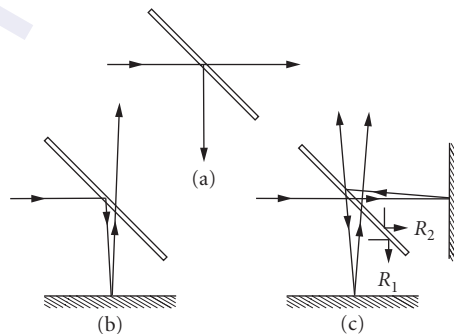


FIGURE 69 Three different ways [(a), (b), and (c)] of using beam splitters.

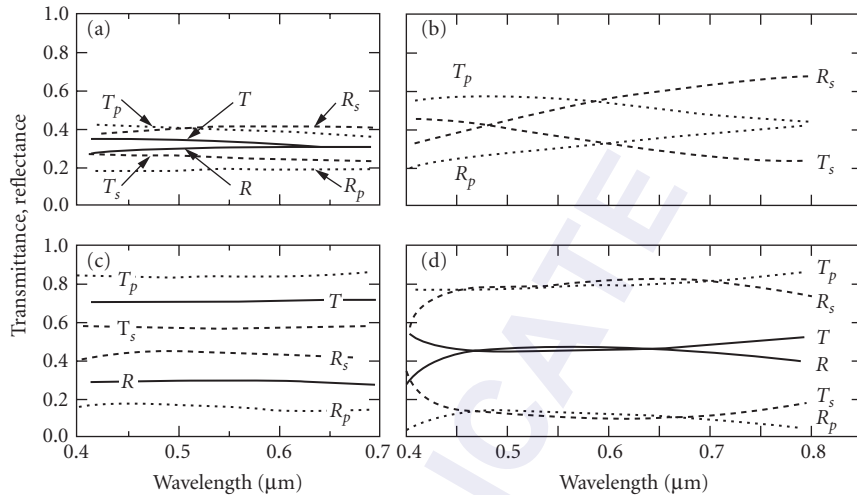


FIGURE 70 Measured spectral transmittance for polarized and unpolarized light of (a) Inconel- and (c) dielectric-coated beam-splitting plates (Oriol³⁰⁶) and of (b) silver- and (d) dielectric-coated beam-splitting cubes (after Anders²⁹⁷).

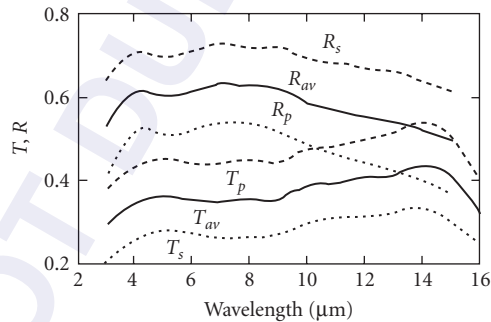


FIGURE 71 Measured performance of a 45° infrared beam splitter consisting of a suitably coated ZnSe plate. (After Pellicori.³⁰⁷)

Beam splitters for the x-ray region described so far operate at close to normal incidence and are effective only over a very narrow range of wavelengths. They consist of multilayer reflecting stacks (see Secs. 7.6, subsection “Multilayer Reflectors Made of Absorbing Materials” and 7.7, subsection “Multilayer Reflectors for the Soft X-Ray and XUV Regions”) deposited onto membranes or onto substrates that are thinned to enhance the transmitted component (Fig. 72).³⁰⁸

Nonpolarizing Beam Splitters For some applications it is important that the beam splitter introduce no polarization effects. Azzam has shown that, with the appropriate single layer on the face of a suitable high-refractive-index prism, it is possible to construct a polarization independent beam splitter.³⁰⁹ This device is quite achromatic and, in addition, by changing the angle of incidence of the beam on the prism, the beam-splitting ratio can be tuned over a wide range of values (Fig. 73a). The principle of frustrated total internal reflection can also be used to design beam splitters that have a very good performance.^{310,311} In these devices radiation is incident at a very oblique angle

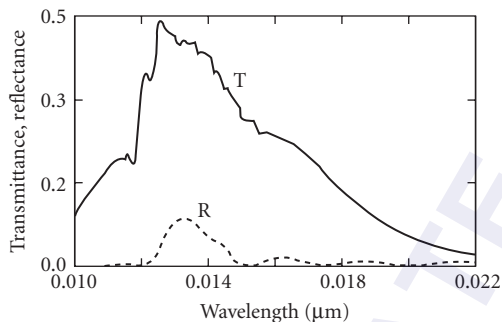


FIGURE 72 Measured performance of an x-ray beam splitter consisting of 11 pairs of Mo and Si layers on a 0.03- μm -thick Si_3N_4 membrane, operating at an angle of incidence of 0.5° . (After Ceglie.²⁶³)

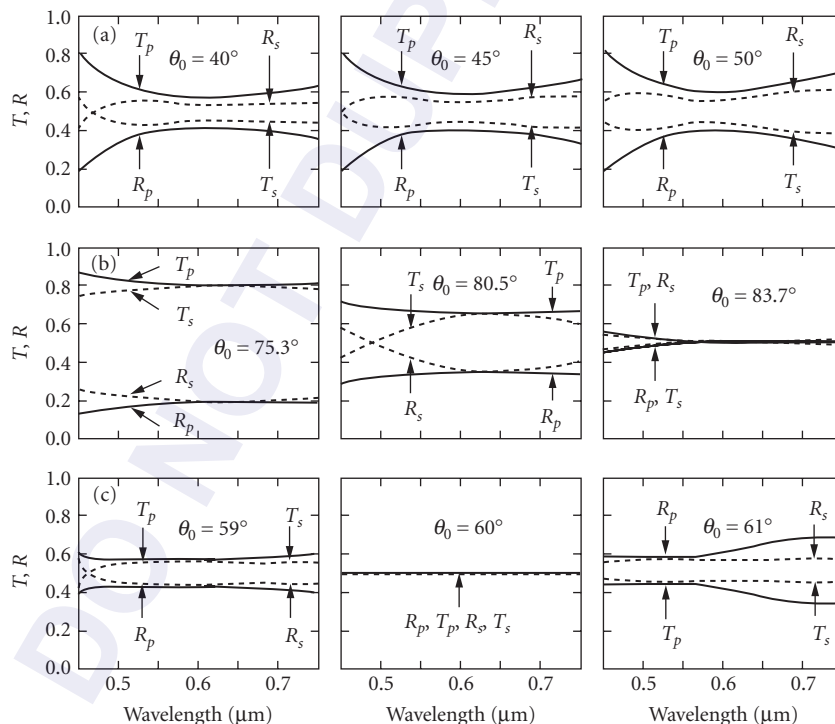


FIGURE 73 Calculated performance at three angles of incidence of beam splitters with an incident medium of air. (a) Beam splitter of the type glass $(HL)^2$ air, where $n_H = 2.35$, $n_L = 1.38$; (b) single layer ($n = 1.533$, $d = 0.1356 \mu\text{m}$) on a prism ($n_s = 2.35$) (after Azzam³⁰⁹); and (c) 15-layer frustrated total internal reflection beam splitter (after Macleod³¹¹). The last two systems are fairly polarization independent and yield different T/R ratios for different angles of incidence.

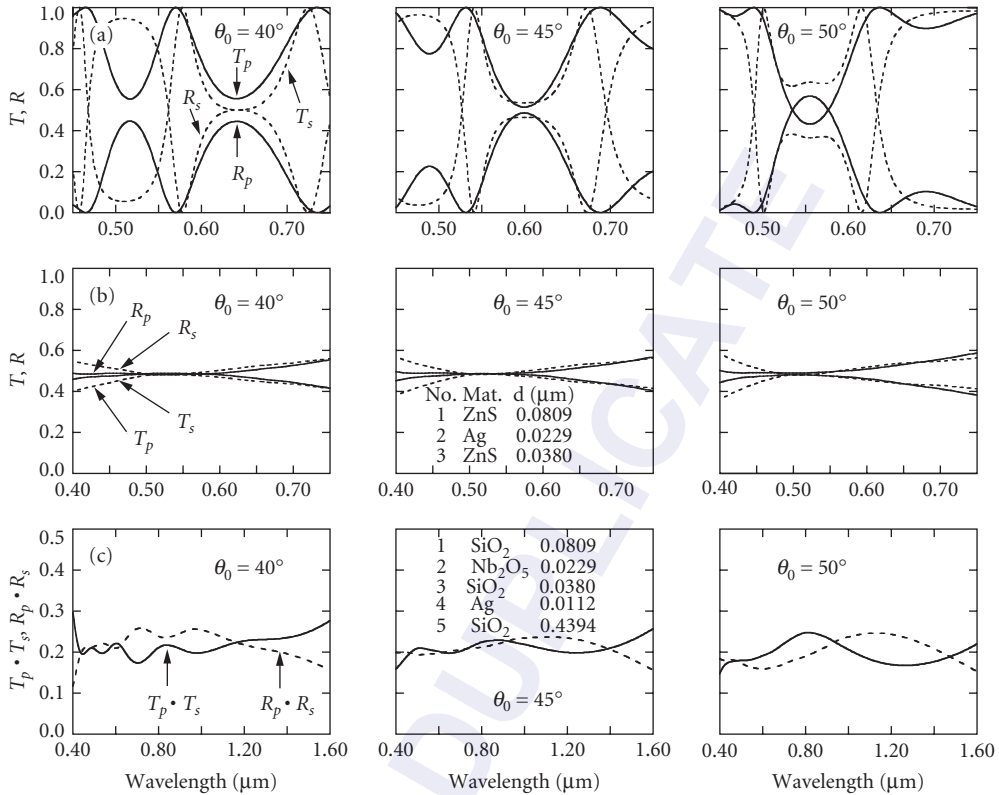


FIGURE 74 Calculated performance of polarization insensitive achromatic beam splitters consisting of layer systems cemented between glass prisms. (a) All-dielectric system of the type $n_s(LMHMHML)^2$, where $n_s = 1.52$, $n_L = 1.38$, $n_M = 1.63$, $n_H = 2.35$; the quarterwave layers are matched for 45° incidence. (after Thelen¹³); (b) three-layer metal/dielectric system (after Chang³²²); and (c) polarization independent beam splitting arrangement of the type of Figure 68o composed of three identical prisms (after Ho²⁹⁶).

onto an air gap or low-refractive-index films at the interface between two prisms (Fig. 68e, and f). Unfortunately, the performance of such systems is very sensitive to the angle of incidence (Fig. 73c).

In many applications it is important that the beam splitter be relatively insensitive to the angle of incidence. One way is to reduce the angle of incidence as much as possible (Fig. 68h).³¹² However, in many cases a 45° angle of incidence is mandatory. Relatively polarization insensitive beam splitters based on dielectric-metal-dielectric layer systems embedded between two prisms have been described.^{313,314} Much work has been done to find solutions based on dielectric layers only.^{24,207,314–321} However, the improvement is frequently at the expense of the width of the spectral region over which the beam splitter is effective. Some typical results are shown in Figs. 74 and 75.

Simple angle- and polarization-insensitive mechanical solutions to the achromatic beamsplitting problem exist if the application can tolerate spacial or temporal beam sharing (Fig. 76).

Color-Selective Beam Splitters

For various technological applications a beam of light must be divided into several components of different color. All-dielectric color selective beam splitters (*dichroics*) are used for this purpose

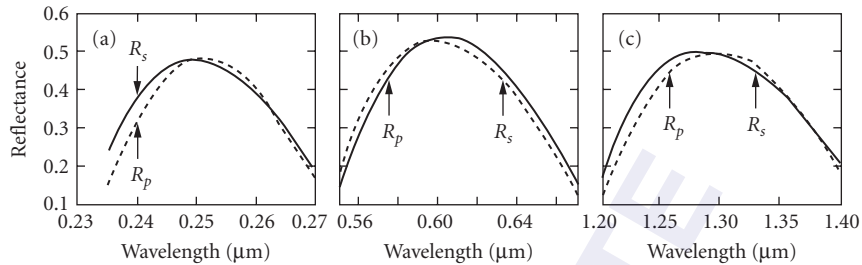


FIGURE 75 Measured performances (a) to (c) of polarization insensitive beam splitters of the type of Fig. 74a produced for three different spectral regions. (After Konoplev.³²³)

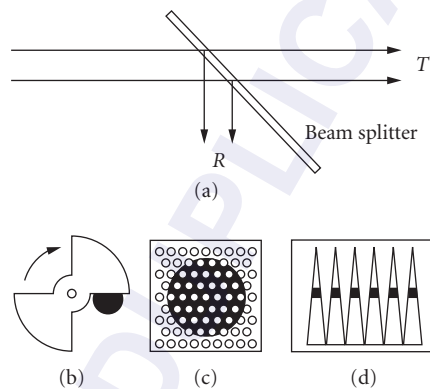


FIGURE 76 Achromatic and angle-insensitive mechanical beam splitting arrangements. (a) and (b) Rotating aluminized blade which alternatively transmits and reflects the incident radiation. (c) Stationary transparent plate with a polka-dot pattern in which either the circles or the background (each of equal total area) are aluminized. (d) Aluminized saw-tooth pattern on a transparent substrate that can be displaced to reflect different fractions of the incident radiation that is collimated into a narrow beam.

because they are practically lossless and because their transition wavelengths can be selected at will. They are essentially cutoff filters (Sec. 7.8) usually designed for use at 45° incidence. Their spectral characteristics normally depend on the polarization of the incident radiation. The effect of this and of the variations in the angle of incidence and thickness of the coatings on the chromaticity coordinates of dichroic beam splitters for television cameras was investigated by Pohlack.³²⁴ If necessary, the polarization of the derived beams can be reduced through the use of auxiliary normal incidence cutoff filters.²¹⁹ Typical transmittance curves of several color-selective beam splitters are shown in Fig. 77.

Nonpolarizing Edge and Bandpass Filters For more exacting applications, such as for use in multiplexers and demultiplexers, or for the separation of emission or absorption lines in atmospheric physics or Raman spectroscopy, it is possible to design and construct short- and long-wavelength color selective beam splitters in which the polarization splitting has been largely eliminated.^{13,325} In the designs described the polarization splitting is usually removed for all angles smaller than the design angle, but the cutoff wavelength still shifts with the angle of incidence (Fig. 78). Some of the designs do not have a wide transmission region.

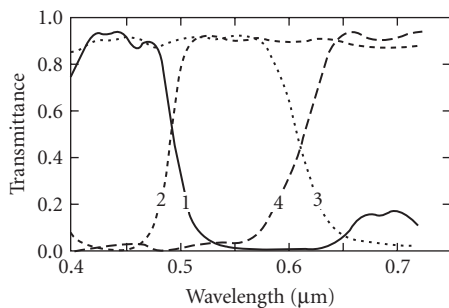


FIGURE 77 Measured spectral transmittance of four color-selective beam splitters. (*Optical Coating Laboratory*.²³⁷)

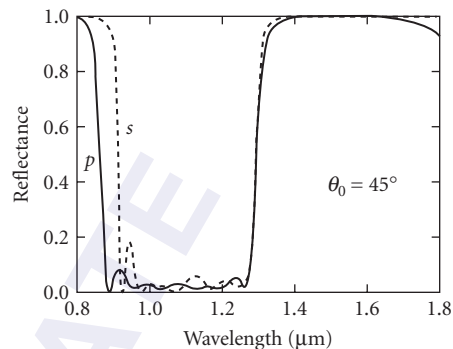


FIGURE 78 Calculated transmittance of a polarization independent color selective beam splitter for several angles of incidence. (*After Thelen*.¹³)

Neutral Filters

These devices are used whenever the intensity of the incident radiation is to be reduced uniformly throughout an extended part of the spectrum. The radiation usually traverses neutral filters at or near-normal incidence. A number of absorbing glasses and gelatin filters are suitable for making neutral-density filters with densities of up to 5.0. However, their spectral transmission curves are not very uniform.

Evaporated films of metals such as aluminum, chromium, palladium, platinum, rhodium, tungsten, and alloys such as chromel, nichrome, and inconel have been used for a long time to produce filters with densities of up to 6.0. A disadvantage of such filters is their high specular reflection. At present inconel is commonly used for high-precision neutral-density filters. Chromium is favored when tough, unprotected coatings are required (Fig. 79). The operating range and neutrality depends on the substrate materials. The spectral-transmittance curves of neutral-density filters on magnesium fluoride, calcium fluoride, quartz, glass, sapphire, and germanium substrates are shown in Fig. 80. Linear and circular metal-film neutral-density wedges and step filters are also available commercially.

At times there may be a need for a neutral attenuation that is not based on absorption. Sets of all-dielectric multilayer coating designs have been published with uniform transmission levels for the ultraviolet,³²⁶ visible,^{215,327} and near-infrared³²⁸ parts of the spectrum (Fig. 81).

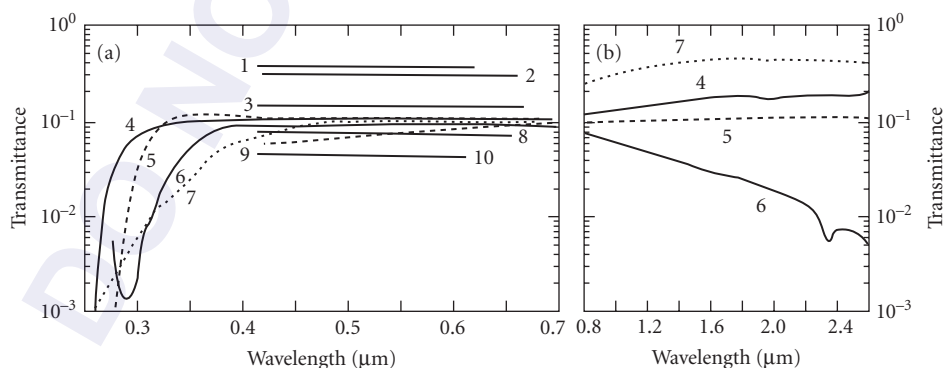


FIGURE 79 (a), (b) Spectral transmittance of various neutral density materials. Curve 1: tungsten film on glass; curves 2 and 3: diffuse and specular transmittance of photographic emulsion; curve 4: M-type carbon suspension in gelatine; curve 5: Inconel film on glass; curve 6: photographic silver density; curve 7: Wratten 96 density filter; curve 8: chromium film; curves 9 and 10: chromel A film on glass evaporated at pressures of 10^{-3} and 10^{-4} Torr, respectively. (Curves 1 to 3, 9, and 10 after Banning;³²⁹ curves 4 to 7 after Eastman Kodak Company;³³⁰ and curve 8 after Optical Coating Laboratory.²³⁷)

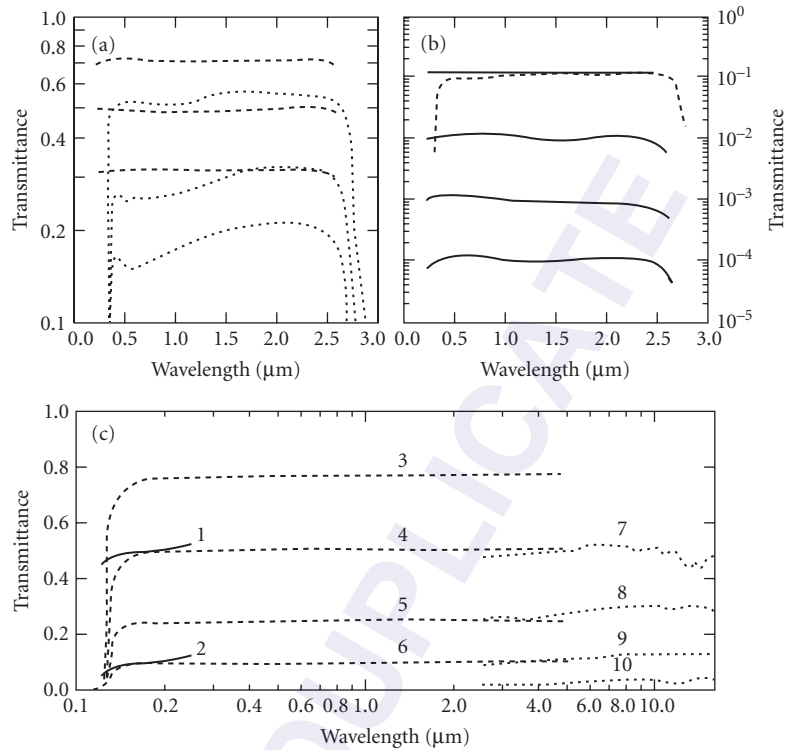


FIGURE 80 Neutral-density attenuators (a), (b) Inconel films on glass substrates (dotted curves, Bausch & Lomb²⁷⁵) and on quartz substrates (full and broken curves, Corion³³¹). (c) Alloy films on MgF₂ (solid curves, after Acton³³²); on CaF₂ (broken curves, after Spindler & Hoyer⁸⁸); and on Ge (dotted curves, after Oriol³⁰⁶)

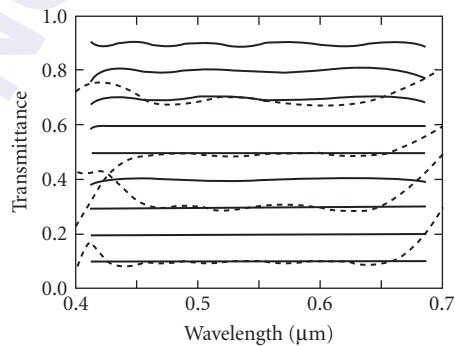


FIGURE 81 Achromatic all-dielectric attenuators. The solid curves are calculated results (after Hodgkinson³²⁷) and the dotted curves represent the performance of commercial coatings (after TechOptics⁸⁶).

7.10 INTERFERENCE POLARIZERS AND POLARIZING BEAM SPLITTERS

The dependence of the optical properties of thin-film systems on the plane of polarization of obliquely incident radiation²⁴ (Sec. 7.2, subsection “Matrix Theory for the Analysis of Multilayers Systems”) can be exploited to design interference polarizers and polarizing beam splitters with properties that augment those attainable by other means (Chap. 13, Vol. 1). The main difference between a polarizer and a polarizing beam splitter is that in the former only one polarized beam is required, whereas in the latter both beams are to be utilized. A polarizing beam splitter can therefore also be used as a polarizer. The performance in transmission or reflection of both devices is usually characterized by their degree of polarization P , or by the extinction ratio ER, both of which are functions of the transmittance T and reflectance R for the desired and undesired polarizations:

$$\left. \begin{array}{l} P = \frac{T_{\text{desired}} - T_{\text{undesired}}}{T_{\text{desired}} + T_{\text{undesired}}} \quad \text{or } P = \frac{R_{\text{desired}} - R_{\text{undesired}}}{R_{\text{desired}} + R_{\text{undesired}}} \\ \text{or by} \\ \text{ER} = \frac{T_{\text{desired}}}{T_{\text{undesired}}} \quad \text{or } \text{ER} = \frac{R_{\text{desired}}}{R_{\text{undesired}}} \end{array} \right\} \quad (59)$$

In a polarizing beam splitter a high degree of polarization is required in both beams and this is more difficult to achieve. Presently efficient interference polarizers and polarizing beam splitters can be constructed for the soft x-ray, ultraviolet, visible, and infrared spectral regions. Interference polarizers and polarizing beam splitters are of particular interest whenever large areas and low losses are required. Schematic representations of the geometries of some of the devices are given in Fig. 68.

Multicomponent Polarizers

It is always possible to find an angle of incidence at which a substrate coated with a film of quarter-wave effective thickness will reflect only s -polarized radiation which is polarized perpendicular to the plane of incidence.³³³ This property can be used to construct efficient transmitting polarizers using far fewer plates than necessary in the conventional pile-of-plates polarizer of equal performance (Fig. 68*d*).³³⁴ The calculated degree of polarization attainable with different numbers of plates is shown in Fig. 82 for a series of film indices and polarizing angles. Experimental results agree closely with the calculations. In polarizers of this type the variation of the degree of polarization is small over a wavelength span of 2:1 and for angular apertures of up to 10°. Yet, because of its bulk, this type of polarizer is not frequently used. An exception is polarizers for the infrared, where it is more difficult to produce multilayers composed of many layers. Because of the high refractive indices available in that spectral region, it is possible to achieve a high degree of polarization even after a single transmission through a plate coated with one layer only.³³⁵

Several geometries of the reflection equivalent of the multiple plate polarizer exist (Fig. 68). The reflectors can consist of one or more dielectric layers deposited onto nonabsorbing parallel plates or prisms (Fig. 68 *t* to *w*). In other devices the substrates are made of metal, or are coated with an opaque metallic film (Fig. 68 *p* to *s*).³³⁶ The angles of incidence on the various mirrors need not be the same. Polarizers of this type are particularly useful in the infrared³³⁷ and vacuum ultraviolet (see subsection “Polarizers for the Extreme Ultraviolet and Soft X-Ray Regions”) spectral regions. Interesting variants are polarizers that are based on total internal reflection or frustrated total internal reflection (Fig. 68 *m* and *n*).

Plate Polarizers

The number of coated plates in the transmission polarizers described above can be reduced without compromising the performance by depositing more than one high-refractive-index layer onto the surface

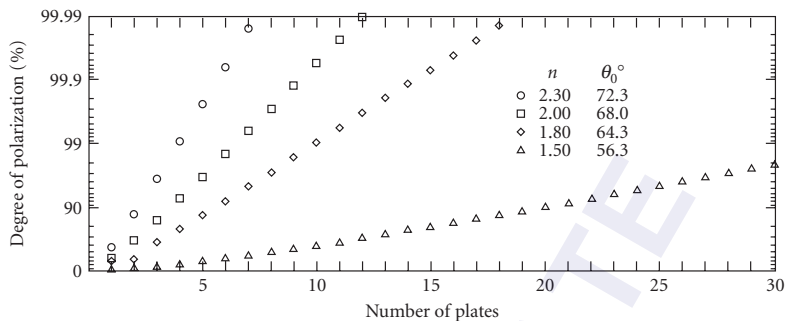


FIGURE 82 Calculated degree of polarization of different numbers of plates of refractive index 1.5 coated on both sides with films of different refractive indexes n . The films have effective optical thicknesses of a quarter wavelength at the appropriate polarizing angle θ .

of a plate and by spacing them with low-index films. In particular, it is possible to minimize the surface scatter, plate absorption, and lateral beam displacement by combining all the layers into one coating.

The most common solution for high-power laser applications are plate polarizers that are based on the polarization splitting that occurs at higher angles, for example, at the edges of quarter-wave stack reflectors (Fig. 68).³³⁸ The plate on which the multilayer is deposited is held at Brewster's angle with respect to the incident light to avoid second surface reflections. Usually the long-wavelength edge of the reflector is used and the design is somewhat modified to remove the ripples in the transmission band of the p -polarized radiation (Fig. 83). The use of other multilayer structures, such as bandpass filters, for the construction of plate polarizers has also been proposed.³³⁹ The wavelength range over which plate polarizers are effective is much smaller than that of polarizers based on a series of coated plates, but this is acceptable for most laser applications.

Other methods for the design of narrowband plate polarizers based on two or three coating materials have been described by Minkov,³⁴⁰ Mahlein,²⁰⁷ and Thelen.³⁴¹ However, these solutions require many layers and sometimes very oblique angles of incidence.

Embedded Polarizers and Polarizing Beam Splitters Based on Two Coating Materials

Polarizers and polarizing beam splitters effective over a wider spectral region are obtained when multilayer coatings of the type $[HL]^N$ or $[(0.5H)L(0.5H)]^N$ are embedded between media of higher refractive index than air.³⁴² The higher the refractive index ratio of the two coating materials used, the fewer the number of layers needed to achieve a certain degree of polarization and the wider the spectral region over which the polarizer will be effective. However, a certain relationship between the refractive indexes of the materials and the angle of incidence must be satisfied.³⁴³ The optical thicknesses of the quarterwave layers should be matched for the angle of incidence.

A particularly convenient polarizer with no lateral beam displacement results when the multilayer is embedded between two right-angled prisms, as shown in Fig. 68i.^{344,345} MacNeille polarizers operate over

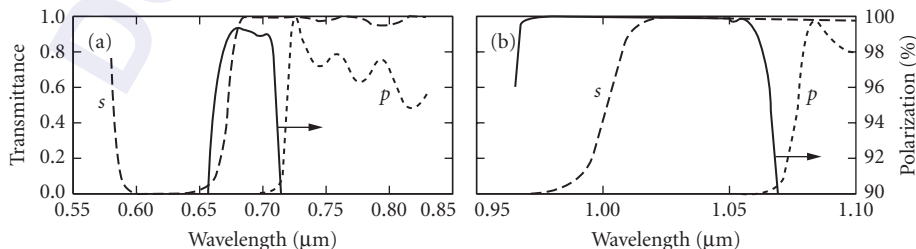


FIGURE 83 Measured performances (a), (b) of two plate polarizers. (After TechOptics.⁸⁶)

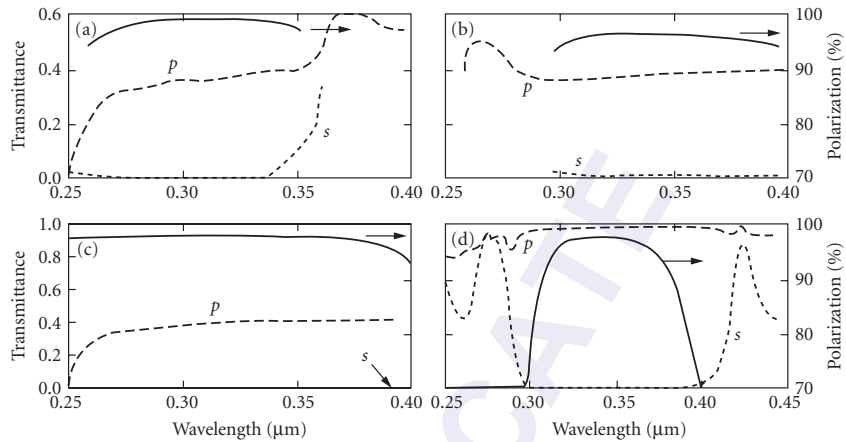


FIGURE 84 The measured degree of polarization P and transmittance of parallel and crossed MacNeille interference polarizers for the ultraviolet spectral region. (c) Represents the results obtained when the cemented polarizers shown in (a) and (b) are placed in series (after Sokolova and Krylova³⁴⁶); and (d) is the measured performance of an optically contacted polarizer with a high laser damage threshold (after Wimperis³⁴⁷).

a very broad range of wavelengths (Fig. 84a and b). For best results, the following relationship between the angle of incidence θ_p and the refractive indices of the prism and the layers should be satisfied:

$$n_p \sin \theta_p = \frac{n_L n_H}{\sqrt{n_L^2 + n_H^2}} \quad (60)$$

This expression is independent of the thicknesses of the layers and, were it not for the dispersion of the optical constants, the transmission for the p polarization would be 1.0 across the entire spectral region at the design angle θ_p . However, the dispersion of the high-index material as well as the limit of the index ratio n_L/n_H will tend to decrease the useful spectral range of the polarizer.³⁴⁸ It is possible to select the V -number of the substrate material in order to decrease this disturbing effect. The rejection of the unwanted polarization will normally not be the same in both beams, although this can be achieved at the expense of the wavelength range.³⁴⁹ If polarizers for an even wider spectral region are required, it is possible to extend the range by placing two polarizers in series (Fig. 84c), by using the technique of superposition of stacks with contiguous high-reflectance zones (Sec. 7.7, subsection “Rejection Filters”) or of two periodic multilayers with thickness ratios of 1:1 and 1:2.¹⁷³ When both beams are used, the useful angular field of the MacNeille polarizer is of the order of $\pm 2^\circ$. It is possible to increase it to $\pm 10^\circ$, but again at the expense of narrowing the spectral range.³⁵⁰

When the MacNeille polarizer is to be used as a polarizing beam splitter, it is possible to design it for normal incidence of the beams onto the prism faces, or for a 90° deflection between the two beams (Fig. 68I).^{349,351} However, it is more convenient to embed the multilayers between two 45° prisms. A higher-refractive-index prism material is required if (60) is to be satisfied³⁵⁰ and the effective wavelength range will once again be reduced (Fig. 85).³⁵²

MacNeille polarizers and polarizing beam splitter cubes suffer from several disadvantages. The diameter of the light beam that can be used with the polarizer is limited by the size, cost, and availability of the prism material. The attainable degree of polarization is also affected by the residual birefringence in the glass prisms. Cemented polarizers cannot be used with high-power lasers because of the absorption of the glue. To overcome these difficulties a liquid prism polarizer has been proposed consisting of a multilayer on a quartz plate that is immersed in distilled water (Figs. 68j and 86).³⁵³ Another way of avoiding the use of a cement is to optically contact the polarizer prisms (Figs. 68g, and 84d).³⁵⁴ More frequently the layers are deposited onto the hypotenuse of one of two air-spaced prisms (Fig. 68e).³⁵⁵ However, this arrangement is similar to a plate polarizer designed for use at 45° and so is its performance.

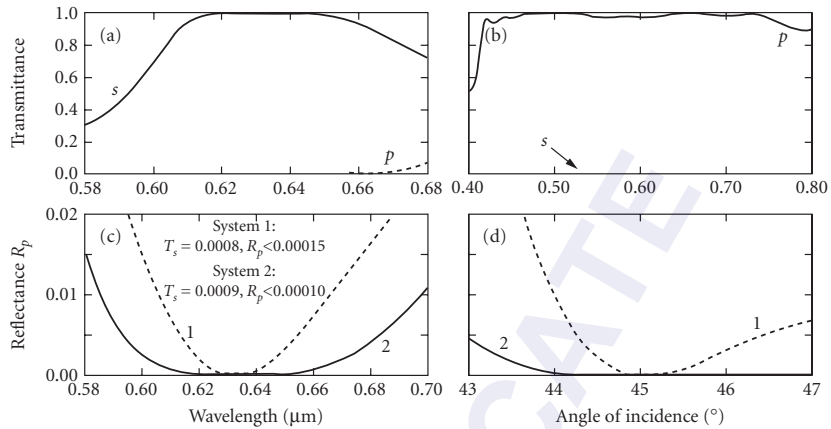


FIGURE 85 Measured performance of polarizing beam splitting cubes. (a) and (b) Measured spectral transmittance curves of two commercial devices (after TechOptics⁸⁶). (c) and (d) Measured spectral and angular performance of two different systems (after Netterfield³⁵²).

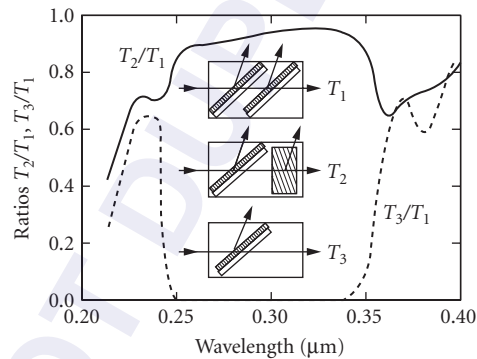


FIGURE 86 Measured spectral performance of two parallel and crossed coated plates forming a liquid prism McNeille polarizer for the ultraviolet part of the spectrum. The multilayers on both plates were identical and consisted of 13 alternate quarterwave layers of HfO_2 and SiO_2 . (After Dobrowolski.³⁵³)

Embedded Polarizing Beam Splitters Based on More than Two Coating Materials

By using coating materials with more than two refractive indices in the design of embedded polarizing beam splitters, it is possible to obtain devices that match the spectral region of a MacNeille polarizer while increasing its angular field. For example, a polarizing beam splitter has been described with a 72-layer 4-material design that has a $\pm 6^\circ$ angular field measured in air that is effective across the 0.4 to 0.7 μm wavelength range.³⁵⁶

Li Li Polarizing Beam Splitter

Li Li described the design of a new type of embedded beam splitter for the visible and infrared spectral regions, using a geometry that is depicted in Fig. 68k.³⁵⁷ What distinguishes this device from the

polarizing beam splitters described in the previous sections is that it operates at angles greater than the critical angle and that it transmits *s*-polarized light and reflects *p*-polarized light, respectively. The calculated performances of these designs in both the transmitted and reflected beams exceed those of any previous devices. For example, one polarizing beam splitter (described in Ref. 357) that consists of 53 layers has a $\pm 7.4^\circ$ angular field measured in air in the 0.4 to 0.7 μm wavelength region.

However, because the Li Li polarizing beam splitters require the use of prisms of high refractive index for which there are no suitable optical cements, the two prisms must be optically contacted. The measured performance of prototypes produced in this manner was in good agreement with the calculated values.³⁵⁸ The need for optical contacting, and the rather large size of the prisms, is the reason why the Li Li polarizing beam splitter is not yet widely used.

Metal-dielectric versions of the Li Li polarizer have been described that are simpler and consist of fewer layers, but they absorb some of the incident light.³⁵⁹

In Fig. 87*a* to *j* the calculated spectral and angular performances of a number of polarizers and polarizing beam splitters for the near-UV, visible, and near-IR spectral regions are compared. A theoretical comparison of the properties of the MacNeille, cube and plate polarizers for one wavelength has been given by Cojocaru.³⁶⁰

Polarizers for the Extreme Ultraviolet and Soft X-Ray Regions

Polarizers based on reflections from two or more single layer coated surfaces have been proposed for the vacuum ultraviolet^{366–368} spectral region. The measured performances in the XUV region of two such polarizers that are based on three reflections are depicted in Fig. 88.

Polarizers for the soft x-ray region can also be based on the fact that the reflectance of an x-ray multilayer mirror at oblique angles is very different for radiation polarized parallel and perpendicular to the plane of incidence. Curve 3, Fig. 91 shows the reflectance of such a periodic Mo/B₁C multilayer with a very high degree of polarization.³⁷⁰ Because the region of high reflection is very narrow, this type of polarizer essentially operates at one wavelength only.³⁷¹ However, if the radiation is reflected from two identical mirrors of this type (Fig. 68*g*), it is possible to construct a device with a reasonable throughput (>0.05) that can be tuned over a wide range of wavelengths without changing the direction of the emerging beam (Fig. 89).³⁷²

More recently multilayer polarizers effective over a broader range of wavelengths have been constructed. A few examples are shown in Fig. 90.^{370,373–375} They are nonperiodic systems made of two suitable materials and they consist of between 60 and 300 layers. The one exception is curve 3 which corresponds to a periodic system and which is shown here for comparison purposes. It will be seen that the effectiveness over a broader range of wavelengths comes at the expense of a lower throughput, but it is also accompanied by a much wider angular field. This is of particular importance when the polarizers are to be used with imaging optics.

7.11 BANDPASS FILTERS

An ideal bandpass filter transmits all the incident radiation in one spectral region and rejects all the other radiation. Such a filter is completely described by the width of the transmission region and the wavelength at which it is centered. Practical filters are not perfect and require more parameters to adequately describe their performance. No uniform terminology has yet been developed for this purpose. Care should be taken when reading and writing specifications since often different terms are used to describe different types of filters, and sometimes quantities bearing the same name are defined differently.

The position of the transmission band is variously specified by the wavelength λ_{max} at which the maximum transmission occurs, the wavelength λ_0 about which the filter passband is symmetrical, or the spectral centre of gravity λ_c of the band. When specifying the tolerance on λ_0 it should be remembered that the peak of interference filters can be moved only toward shorter wavelengths by tilting (see Sec. 7.2, subsection “Matrix Theory for the Analysis of Multilayer System”).

The peak transmittance T_0 may or may not take into account the absorption within the substrate and/or blocking filters used to remove the unwanted transmission of the interference filter away from the principal passband (Fig. 91).

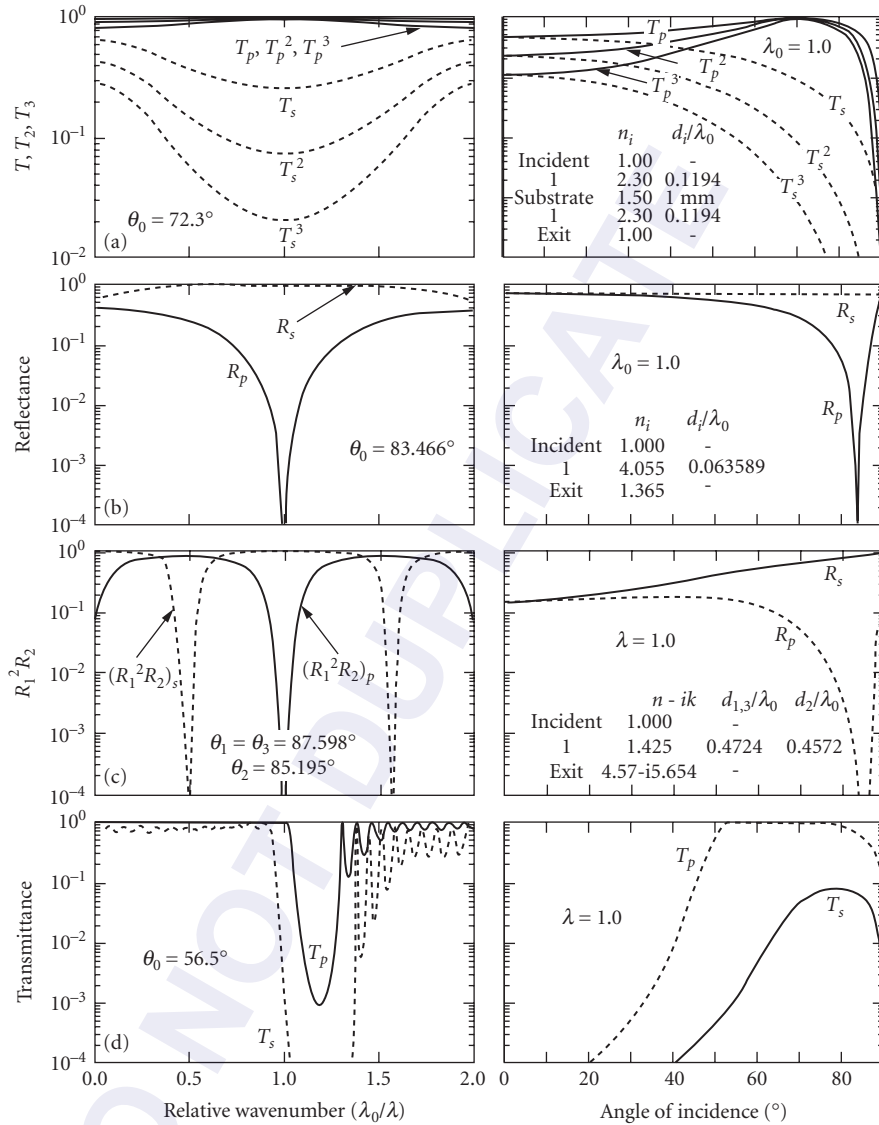


FIGURE 87 Calculated spectral and angular performance of several types of polarizers and polarizing beam splitters with an air incident medium (a to d) and embedded between identical prisms (e to j). Calculations assume that all nonmetals are absorption and dispersion free. (a) Multiple plate polarizer; (b) single reflection polarizer (after Azzam³⁶¹); (c) three reflection polarizer (after Thonn³⁶²); (d) plate polarizer (after Songer³⁶³); (e) MacNeille polarizing beam splitter (after Mouchart³⁵⁰); (f) wide-angle MacNeille polarizing beam splitter (after Mouchart³⁵⁰); (g) frustrated total internal reflection polarizing beam splitter (after Lees and Baumeister³⁶⁴); (h) penta prism polarizer (after Lotem and Rabinovich³⁶⁵); (i) polarizing beam splitter based on more than two materials (after Li Li and Dobrowolski³⁵⁶); and (j) Li Li polarizing beam splitter operating at angles greater than the critical angle (after Li Li and Dobrowolski³⁵⁷).

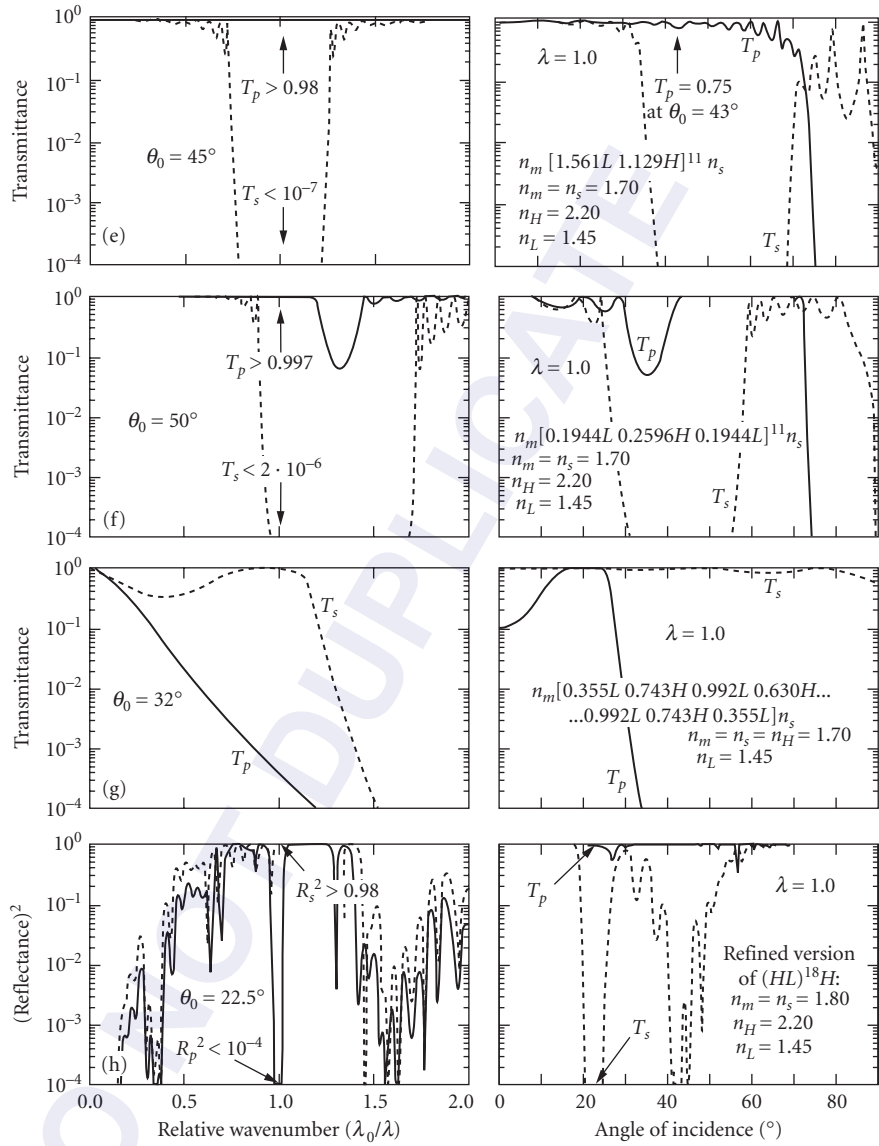


FIGURE 87 (Continued)

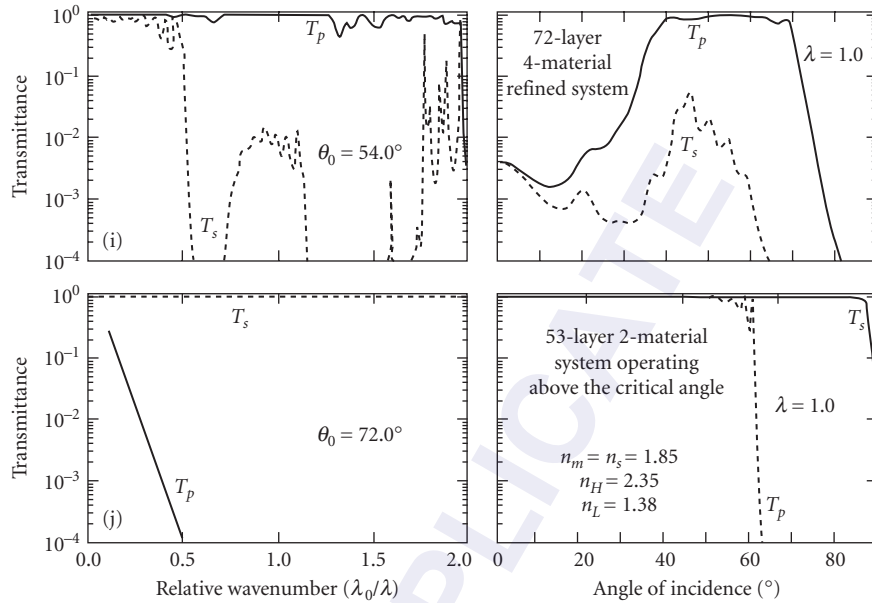


FIGURE 87 (Continued)

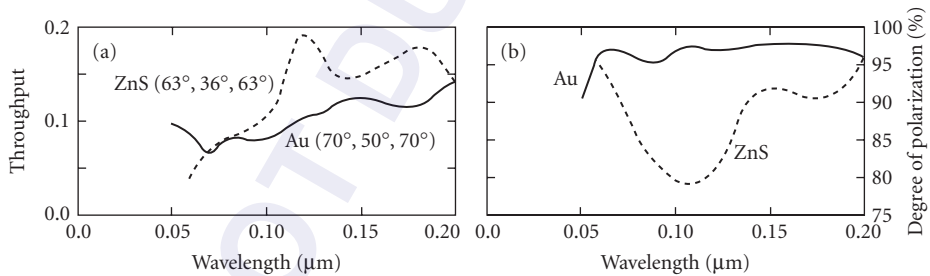


FIGURE 88 Measured throughput (a) and degree of polarization (b) of two extreme ultraviolet polarizers based on three reflections from Au (70°, 50°, 70°) and ZnS (63°, 36°, 63°) surfaces. The angles of incidence on the three mirrors (see Fig. 68r) are given in brackets. (After Remneva et al.³⁶⁹)

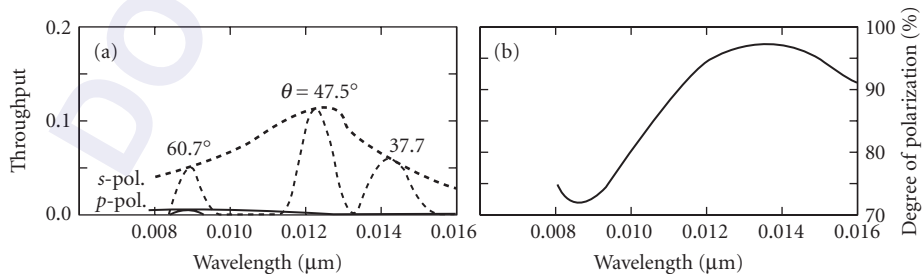


FIGURE 89 Measured throughput (a) and degree of polarization (b) of an x-ray polarizer consisting of two identical 21-layer Ru-C multilayers on Si substrates. The polarizer can be tuned to different wavelengths by changing the angle of incidence θ of the radiation on the two mirrors. (After Yanagihara et al.³⁷²)

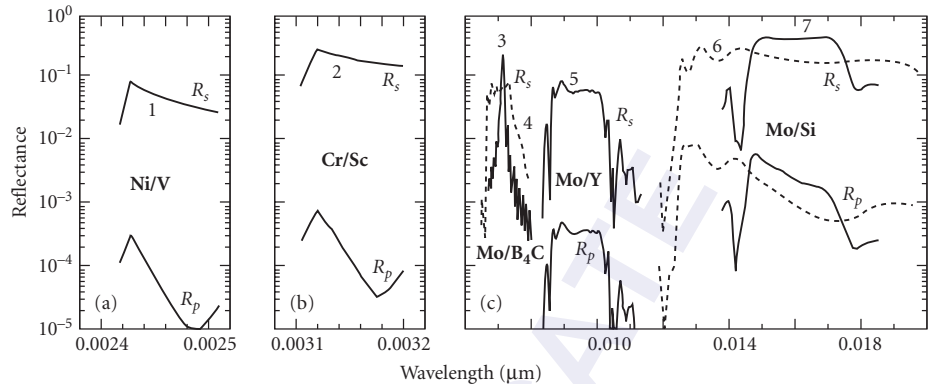


FIGURE 90 Measured spectral performances of two-material multilayer reflection thin-film polarizers for the soft x-ray (a, b) and extreme ultraviolet (c) spectral regions. In the diagrams the reflectances for s - and p -polarized radiation for each polarizer are plotted with the same line type. The materials of which the multilayers are made are indicated in the figures. The measurements were performed at the Brewster angle which, for all the systems, is close to 45° . (Curves 1 and 2 after Eriksson²⁵⁶ and curves 3 to 7 after Zhanshan Wang.^{370,374,375})

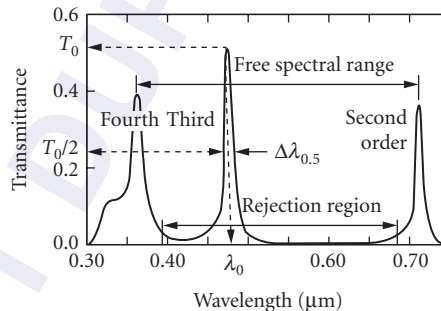


FIGURE 91 Definition of some of the terms used to describe the properties of narrow-bandpass filters (Sec. 7.11). The curve represents the measured transmittance of an unblocked second-order metal-dielectric interference filter of the Fabry-Perot type (Sec. 7.11, subsection “Filters with metallic reflecting coatings”). (After Bausch & Lomb.²⁷⁵)

The half width (HW) $\Delta\lambda_{0.5}$ of the filter is the difference between the wavelengths at which the transmittance is a half of T_0 . This quantity is also sometimes called the full width half maximum (FWHM). It is often expressed as a percentage of λ_0 . The base width (BW) $\Delta\lambda_{0.01}$ is similarly defined. The ratio $\Delta\lambda_{0.01}/\Delta\lambda_{0.5}$, sometimes called the *shape factor*, indicates how “square” the transmission band is. Sometimes widths corresponding to other fractions of the transmittance are used to define it.

The minimum transmittance T_{\min} of the filter does not take into account the effect of blocking filters. The quantity T_{\min}/T_0 is called the *rejection ratio*.

In all-dielectric transmission-band filters the transmittance rises at some distance on either side of the transmission band. The distance over which the transmittance is low is called the *rejection region*. The distance between the two transmission maxima adjacent to the principal transmission band is called the *free spectral range*. Should either of the above two quantities be inadequate, auxiliary blocking filters might have to be provided.

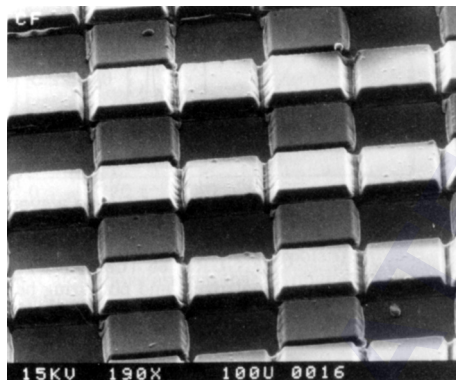


FIGURE 92 Four color checkerboard pattern for use with a 64×64 element focal plane HgCdTe detector array. The dimensions of each element are $100 \times 100 \mu\text{m}$. (Reproduced with permission from Barr Associates.³⁸⁰)

The ultimate measure of the suitability of a bandpass filter with a blocked spectral transmittance $T(\lambda)$ for a particular application is the signal-to-noise ratio SN, defined in terms of the spectral energy distribution $I(\lambda)$ of the source and the spectral detectivity $D(\lambda)$ of the detector,

$$\text{SN} = \frac{\int_{\lambda_1}^{\lambda_2} I(\lambda)T(\lambda)D(\lambda)d\lambda}{\int_0^{\lambda_1} I(\lambda)T(\lambda)D(\lambda)d\lambda + \int_{\lambda_2}^{\infty} I(\lambda)T(\lambda)D(\lambda)d\lambda} \quad (61)$$

where λ_1 , λ_2 are the lower and upper limits of the transmission region of the filter. The SN ratio is sometimes expressed in terms of optical density.

Useful general reviews of bandpass filters exist.^{376,377} Interference filters, and especially bandpass filters, are increasingly deposited in complicated millimeter and sub-millimeter patterns for use with display devices and detectors.^{378,379} Very fine masks or photo-lithographic processes are required to produce such structures (Fig. 92).

Narrow- and Medium-Bandpass Filters (0.1 to 35 percent $h\nu$)

Even though the essential components of Fabry-Perot (FP) interference filters, that is, the spacer and the two reflectors that surround it, can take on many different forms (see Fig. 93), the filters are essentially low-order FP interferometers and hence the theory developed for the latter (see, for example, Born and Wolf¹⁹) applies in full.

The transmittance of a FP type filter, not allowing for absorption and multiple reflections within the substrate, is given by

$$T = \frac{T_R^2}{(1-R)^2 + 4R\sin^2 \delta} \quad (62)$$

where

$$T_R = \sqrt{T_1 T_2} \quad R = \sqrt{R_1 R_2} \quad (63)$$

T_1 , R_1 and T_2 , R_2 are the transmittances and reflectances of the first and second reflectors, respectively, as seen from within the spacer medium, and δ is given by

$$\delta = \frac{2\pi}{\lambda} n d \cos \theta + \phi \quad (64)$$

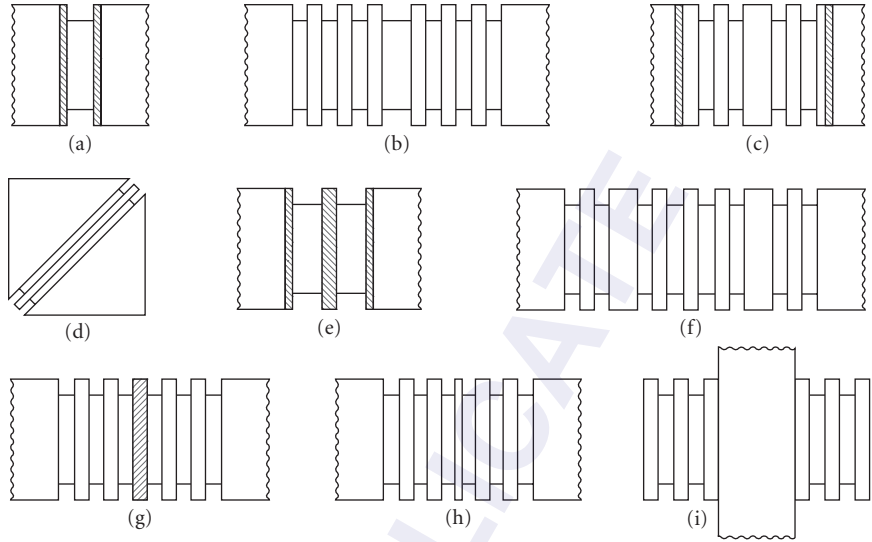


FIGURE 93 Schematic representations of various types of bandpass interference filters: (a) to (c) Fabry-Perot interference filters with metal, dielectric, and metal-dielectric reflectors; (d) frustrated-total-internal-reflection filter; (e) and (f) square-top multicavity filters with metal and dielectric reflectors; (g) induced-transmission filter; (h) phase-dispersion (spacerless) interference filter; and (i) Fabry-Perot filter with a mica or quartz spacer.

$$\phi = \frac{\phi_1 + \phi_2}{2} \quad (65)$$

n , d , and θ are the refractive index, thickness, and angle of refraction of the spacer, respectively, and ϕ_1 and ϕ_2 are the phase changes on reflection from the spacer side of the first and second reflectors at a wavelength λ . Maxima of T occur at wavelengths

$$\lambda_0 = \frac{2nd \cos \theta}{k - \phi/\pi} \quad k=0, 1, 2, \dots \quad (66)$$

and are given by

$$T_0 = \left(\frac{T_R}{1-R} \right)^2 = \left(\frac{1}{1+A/T_R} \right)^2 \quad (67)$$

$A (= 1 - T_R - R)$ is the mean absorption of the reflectors. The minimum transmittance

$$T_{\min} = \left(\frac{T_R}{1+R} \right)^2 \quad (68)$$

occurs at

$$\lambda_{\min} = \frac{2nd \cos \theta}{k - \phi/\pi} \quad k = \frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \dots \quad (69)$$

If T_R and R are essentially the same at λ_0 and λ_{\min} , the rejection ratio is given by

$$\frac{T_{\min}}{T_0} = \left(\frac{1-R}{1+R} \right)^2 \quad (70)$$

For $R > 0.7$ the half width of the transmission band (expressed as a percentage of λ_0) is given by

$$\frac{\Delta\lambda_{0.5}}{\lambda_0} \times 100 \approx \frac{1-R}{\sqrt{R}} \frac{100}{\frac{2\pi n d \cos\theta}{\lambda_0} - \lambda_0 \frac{\partial\phi}{\partial\lambda}} \quad (71)$$

For a given order of interference the half width and rejection ratio cannot be varied independently. The formula

$$T \approx \frac{T_0}{1+4[(\lambda-\lambda_0)/\Delta\lambda_{0.5}]^2} \quad (72)$$

valid for FP filters with small values of $\partial\phi/\partial\lambda$ in the neighbourhood of λ_0 , represents a lorentzian line shape with $\Delta\lambda_{0.1} = 3\Delta\lambda_{0.5}$ and $\Delta\lambda_{0.01} = 10\Delta\lambda_{0.5}$. The shape factor of all FP-type interference filters is therefore of the order of 10.

Southwell has recently shown that a spacer in an interference filter need not consist of a single layer only, and that there are some advantages when it is partitioned into a number of layers.³⁸¹

Fabry-Perot Interference Filters (0.1 to 10 percent hw)

Filters with metallic reflecting coatings This is the first interference bandpass filter ever made, and the simplest.³⁸² It consists of two partially transmitting, highly reflecting metal layers separated by a dielectric film and is symbolically represented by MDM (Fig. 93a). The best metallic reflectors currently available are aluminium and silver for the 0.125 to 0.34 μm and 0.34 to 3.0 μm spectral ranges, respectively. The measured spectral-transmittance curves of a number of typical filters are shown in Fig. 94.

The phase changes on reflection at the spacer-metal-reflector surface are finite and hence they affect the position of the transmission maxima [Eq. (66)]. But the dispersion of the phase change on reflection can be neglected, and so the half width depends only on the reflectance of the metal layers and the order of the spacer. Filters with half widths of 1 to 8 percent are common. Because the ratio A/T is not very small for metals, the maximum transmission is limited [Eq. (67)]. Maximum transmittances of 40 percent are relatively common. For filters with the narrower half widths, or for shorter wavelengths, transmittances of the order of 20 percent have to be accepted. This is much less than the transmittances of all-dielectric filters of comparable half widths. Nevertheless, filters of this type are useful because their transmittances remain low except at wavelengths at which the first- and

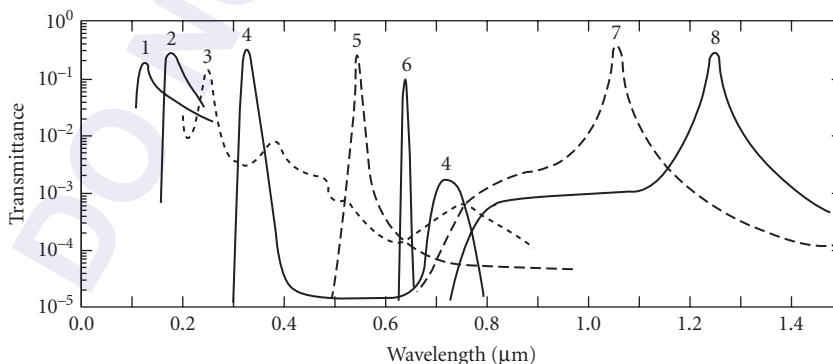


FIGURE 94 Measured transmittance of FP filters with metallic reflecting coatings. Curve 6 represents the transmission of two identical filters cemented together. (Curve 1 after Harrison;³⁸³ curve 2 after Bradley et al.,³⁸⁴ curves 3 and 5 to 7 after Balzers;³⁸⁵ curves 4 and 8 after Schott & Gen.³⁸⁶)

higher-order transmission maxima occur [Eq. (66)]. Blocking is thus easy. In particular, first-order filters usually do not require any blocking on the high-wavelength side—a difficult task at all times. The rejection of the filters is good, though not spectacular. If a better rejection is required and a lower transmittance can be tolerated, two identical filters may be cemented together (curve 6, Fig. 94). This is possible because of the finite absorption in the metal films. Alternatively, metal square-top filters (Sec. 7.11, subsection “Square-Top Multicavity Bandpass Filters”) or filters with even more complicated structures (Sec. XI A 3) can be used.

Filters with all-dielectric reflectors Above 0.2 μm the metallic reflectors can be replaced by all-dielectric quarter-wave stacks (Sec. 7.6).³⁴⁴ The symbolic representation of such a filter (Fig. 93b) is, for example, $[HL]^N 2mH[LH]^N$ or $H[LH]^N 2mL[HL]^N H$, H and L being quarter-wavelength layers of high and low refractive indices, respectively, m is the order of the spacer and N the number of full periods in the reflecting stacks. The phase change on reflection at the boundary between the spacer and such a reflector does not affect the position of λ_0 [Eq. (66)] but the dispersion of the phase change on reflection is finite, depends on the materials used, and for lower-order spacers, contributes very significantly to the reduction of the half width of the transmission band [Eq. (71)]. Expressions for the half width

$$\begin{aligned} \frac{\Delta\lambda_{0.5}}{\lambda_0} \times 100 &= \frac{4n_0 n_L^{2N} (n_H - n_L) \times 100}{m\pi n_H^{2N+1} (n_H - n_L + n_L/m)} \\ &= \frac{4n_0 n_L^{2N-1} (n_H - n_L) \times 100}{m\pi n_H^{2N} (n_H - n_L + n_L/m)} \end{aligned} \quad (73)$$

for high- and low-refractive-index spacers, respectively have been given by Macleod.²

By choosing a suitable combination of the reflectance and order of the spacer almost any half width between 0.1 and 5 percent can be achieved in the visible part of the spectrum while maintaining a useful rejection ratio.

The maximum transmittances of all-dielectric FP filters depart from unity because of the finite absorption, scattering, and errors in the thicknesses and refractive indexes of the films. In the central part of the visible spectrum maximum transmittances of 0.8 are normal for unblocked filters with a half width of 1 percent, although higher transmittances can be achieved. This figure is gradually reduced as λ_0 approaches 0.2 or 20 μm , and filters with narrower half widths become impracticable for lack of adequate transmittance.

The transmittance away from the transmission maximum is low only over the extent of the rejection region of the two materials used for the construction of the reflectors (Fig. 34), and additional blocking is often required on both the long- and the short-wavelength sides. This can result in a considerable lowering of the maximum transmittance of the blocked filter, a 30 to 40 percent loss being not uncommon for filters peaked in the ultraviolet or infrared spectral regions.

For those parts of the visible and infrared for which nonabsorbing mechanically robust coating materials abound, square-top interference filters (Sec. 7.11, subsection “Square-Top Multicavity Bandpass Filters”) are often preferred because of their better shape factor and higher rejection ratio. All-dielectric FP filters are still attractive in the ultraviolet, where there is a lack of such materials and where thickness monitoring is difficult, and also in the far-infrared, where very thick layers are required.

The measured transmittance curves of a number of typical all-dielectric FP filters are shown later in Fig. 98. An additional curve on a smaller scale is given, whenever necessary, to show the transmittance away from the passband.

Filters with metal-dielectric reflectors In these filters the reflectors consist of metal layers whose reflectance has been enhanced through the addition of several dielectric layers (Sec. 7.14, subsection “Fabry-Perot Filters with Solid Spacers”) (Fig. 93c).³⁸⁷ The properties of such filters are intermediate to those described in the two previous sections.

Frustrated total internal reflection filters These are essentially FP filters in which the spacer layer is surrounded by two frustrated-total-reflection surfaces (Fig. 93*d*).^{387,388} They have not found wide applications as bandpass filters because the finite absorption and scattering within the layers have prevented the theoretically expected high transmittance and small half widths from being realised and because the angular variation of the wavelength of the transmission peaks is very high.

Square-Top Multicavity Bandpass Filters (0.1 to 35 percent hw) A filter with a “squarer” shape that does not suffer from some of the disadvantages of the FP filters results when the basic FP structure is repeated two or more times.^{158,389,390} Such filters may be based on metal³⁹¹ or all-dielectric^{392,393} reflectors (Fig. 93*e* and *f*). Thus, for example, $[HL]^N 2H [LH]^N C [HL]^N 2H [LH]^N$ represents an all-dielectric square-top filter in which the FP structure $[HL]^N 2H [LH]^N$ is repeated twice. The quarter-wavelength layer *C* is called a *coupling* or *tie layer*, and the half-wavelength-thick spacer layers $2H$ are often called *cavities*. Small departures from the above model are made at times to improve the transmittance in the pass band or the angular properties of the filter.

The half widths of the narrower multicavity filters of the above type do not differ very significantly from those of the basic FP structure [Eq. (71)]. The shape factors decrease with an increase in the number of the cavities and do not seem to depend on the materials.³⁸⁹ They are approximately 11, 3.5, 2.0, and 1.5 for one, two, three, and four cavities, respectively. The minimum transmittance in the rejection region is roughly that which could be obtained if the filter were composed entirely of $\lambda/4$ layers [Eq. (44)]. Unlike in the FP filter, the half width and rejection ratio can therefore be varied independently. These various points are illustrated in Fig. 95. The peak transmittance of multicavity square-top filters is less affected by the residual absorption in the layers than that of FP type filters. As in the case of their FP counterparts, metal-dielectric square-top filters can be cemented together to enhance the rejection (curve 5, Fig. 96).

The improvement in the performance of square-top bandpass filters over that of the FP type is such that despite their more critical and expensive production, most manufacturers regard them as their standard line of filters. The spectral-transmittance characteristics of typical commercially produced metal-dielectric and all-dielectric bandpass interference filters of different half widths are shown in Figs. 96 and Figs. 97 to 102, respectively. Filters with intermediate half widths and peak wavelengths can readily be obtained.

For very critical or special applications multicavity filters are designed and constructed with properties that exceed those shown in the above figures. For example, for the use in fiber-optic communications systems, multicavity filters are required in which the peak transmittance closely approaches unity. Various procedures for the design of such filters, including some that are based on the use of Chebyshev polynomials, have been described.^{407,408} Special care has to be taken during the manufacture of the coatings to meet this requirement. Typical measured spectral transmittance

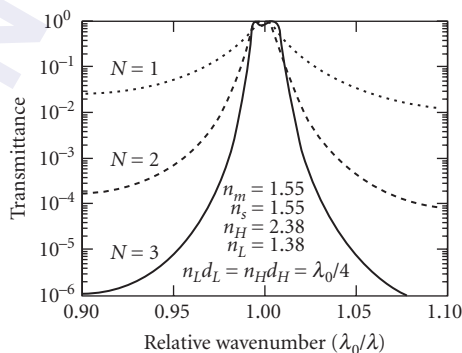


FIGURE 95 Calculated transmittance on a logarithmic scale of the bandpass filters: $air-[(0.5H)L(0.5H)]^3 H [(0.5H)L(0.5H)]^3 N$ -glass, $N = 1, 2$, and 3 .

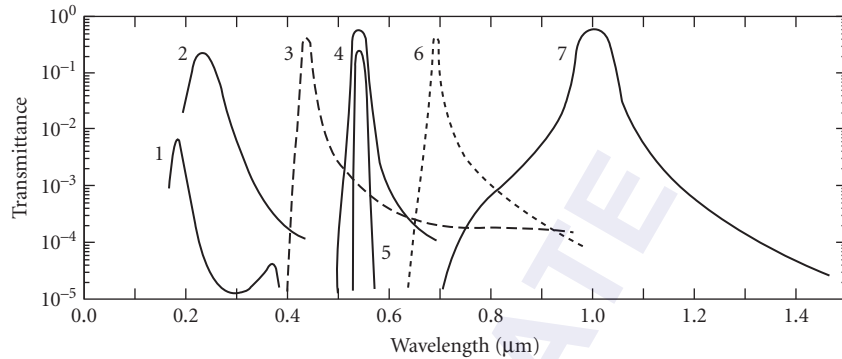


FIGURE 96 Measured transmittance of square-top bandpass filters with metallic reflecting coatings. Curve 5 corresponds to the transmission of two identical filters cemented together. (Curves 1 and 2 after Schröder;³⁹⁴ curves 3 and 6 after Balzers;³⁸⁵ and curves 4, 5, and 7 after Schott & Gen.³⁸⁶)

curves are shown in Fig. 103. For other applications, such as fluorescence or Raman spectroscopy, the peak transmittance is not important, but signal-to-noise ratios of the order of 10^{-8} are required. This necessitates the use of many cavities (Fig. 104).

Induced-Transmission Filters The transmittance of a metal layer can be considerably enhanced by surrounding it with suitable multilayer structures (Fig. 93g). Thus, for example, it is possible to induce a transmittance of 65 percent at $\lambda = 0.25 \mu\text{m}$ in a $0.03\text{-}\mu\text{m}$ -thick aluminum film which, when deposited directly onto a quartz substrate, would transmit only 2.5 percent of the same radiation.⁴¹⁰ The induced transmittance is highly wavelength-sensitive and can be used to construct bandpass filters containing one or more metal layers.⁴¹¹⁻⁴¹⁷ Induced-transmission filters combine the good long-wavelength attenuation properties of the more common types of metal/dielectric filters (Secs. 7.11, subsection "Filters with metallic reflecting coatings" and "Filters with all-dielectric reflectors") with peak transmittances that are closer to those of all-dielectric filters. The performances of some experimentally produced induced-transmission filters are shown in Fig. 105.

Very Narrow Bandpass Filters ($hw = 0.1$ percent)

It follows from (71) that the half widths of interference filters can be reduced by increasing the reflectance of the reflectors, the order of interference of the spacer, or the dispersion of the phase change on reflection. All these approaches have been tried in the past.

Filters with Evaporated Spacers ($hw > 0.03$ percent) In narrowband filters of conventional construction both high reflectance and higher-order spacers are used. The manufacturing process is quite critical, and attention must be paid to details. The films have to be very uniform over the filter area, and they must not absorb or scatter. They must not age, or, alternatively, their aging must be capable of being accelerated or arrested. Monitoring must be precise, so that the peak occurs at or close to the desired wavelength.

Both FP and square-top filters of very narrow bandwidths can be made, the latter being preferable for most applications. The present limit on the half widths of this type of filter seems to be of the order of 0.03 percent.⁴¹⁹ The performance of two commercially produced filters of this type are shown in Fig. 97a and b.

Fabry-Perot Filters with Solid Spacers In practice the half width of interference filters cannot be reduced indefinitely by increasing the optical thickness of an evaporated spacer [Eq. (71)] because

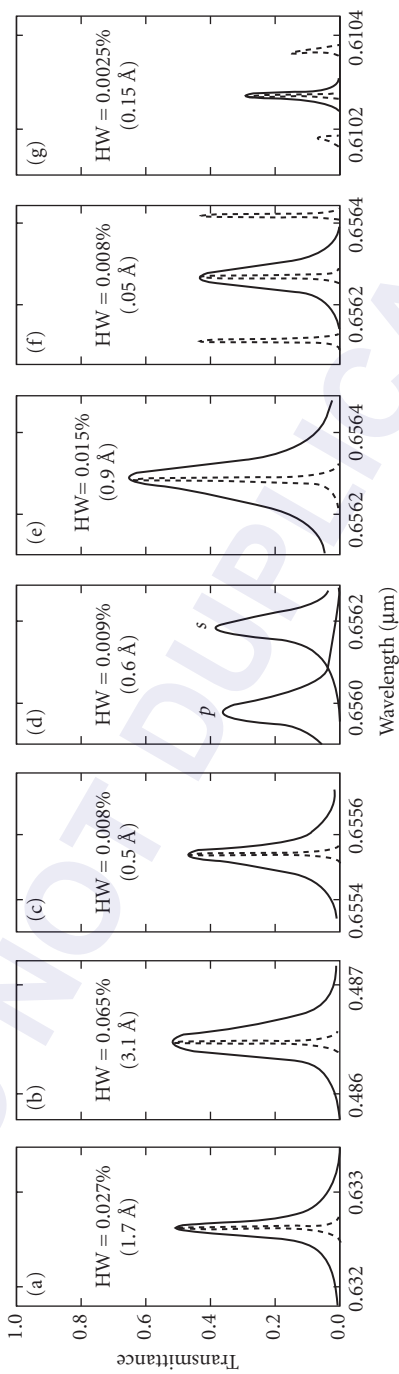


FIGURE 97 Measured transmittance of very narrow bandpass interference filters with half widths less than 0.1 percent. Evaporated spacers: (a) (after Meltzer³⁹⁵); (b) (after Eather and Reasoner³⁹⁶); (c) mica interference filter for H_e^2 ; (d) mica interference filter with transmission bands polarized at right angles to one another (Heliotek²⁸¹); (e) single and (f) and (g) double quartz-spacer interference filters (after Austin³⁹⁷). The dotted curves in Figs. 97 and 98 represent the transmittances of the filters plotted over a ten-times-wider spectral region.

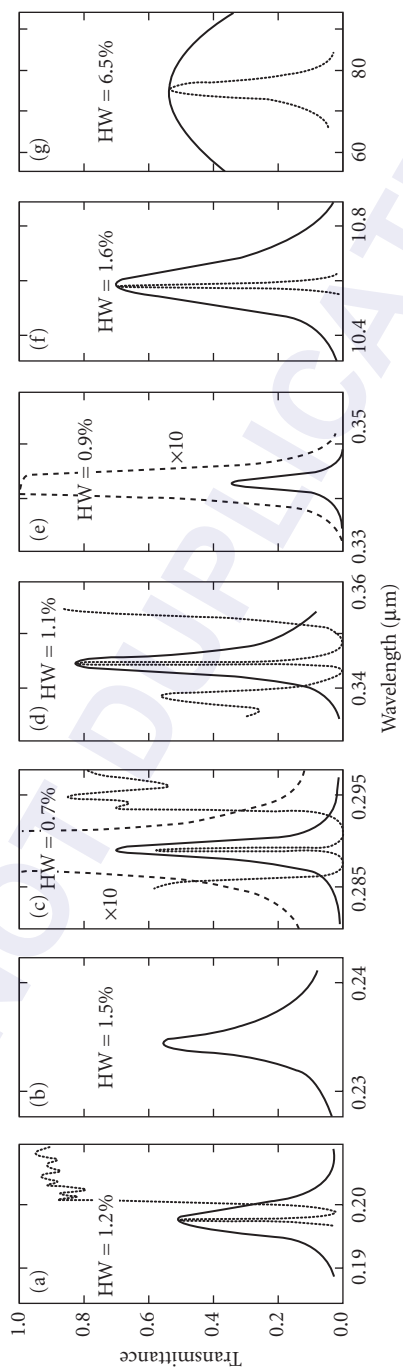


FIGURE 98 Measured transmittance of FP all-dielectric interference filters with narrow half widths. Evaporated spacers: (a) (after Cohendet and Saudreai⁹⁸); (b) and (d) (after Motovilov⁹⁹); (c) and (e) (after Neilson and Ring¹⁰⁰); (f) (after Turner and Walsit¹⁰¹); and (g) (after Smith and Seeley¹⁰²). The dashed curves in Figs. 98 and 99 correspond to a transmittance range of 0.0 to 0.1.

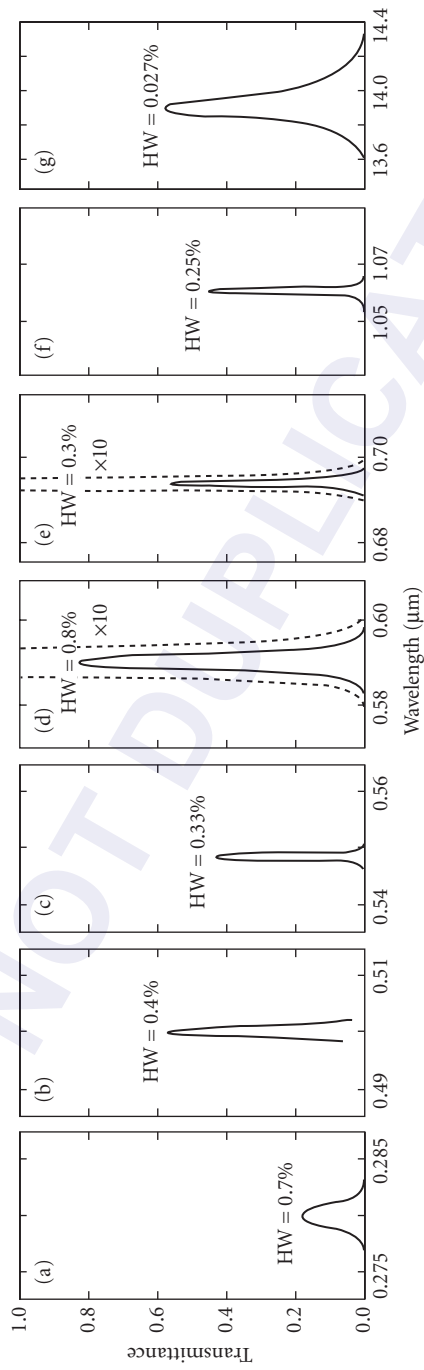


FIGURE 99 Measured transmittances of fully blocked square-top interference filters with half widths between 0.25 and 1.4 percent. (a), (c) and (f) (after Corion²³¹); (b) (after Blifford⁴⁰³); (d) (after Heiotek²⁸¹); (e) (after Baird Atomic⁴⁰⁴); and (g) is the only filter in the series that is not blocked (after Smith and Seeley⁴⁰²).

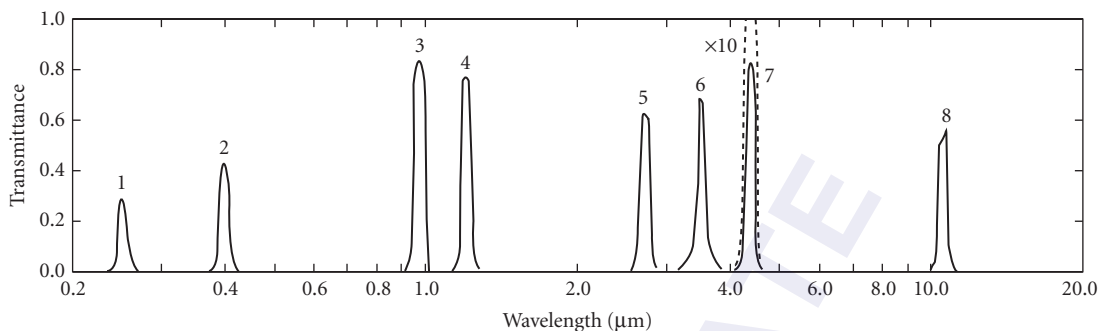


FIGURE 100 Measured transmittances of fully blocked square-top interference filters with half widths of the order of 5 percent. (Curves 1 and 2 after Corion;³³¹ curve 3 after Spectrum Systems;⁴⁰⁵ curves 4 and 6 after Turner;⁸³ curve 5 after Eastman Kodak;²⁷⁶ and curves 7 and 8 after Optical Coating Laboratory.^{237,274})

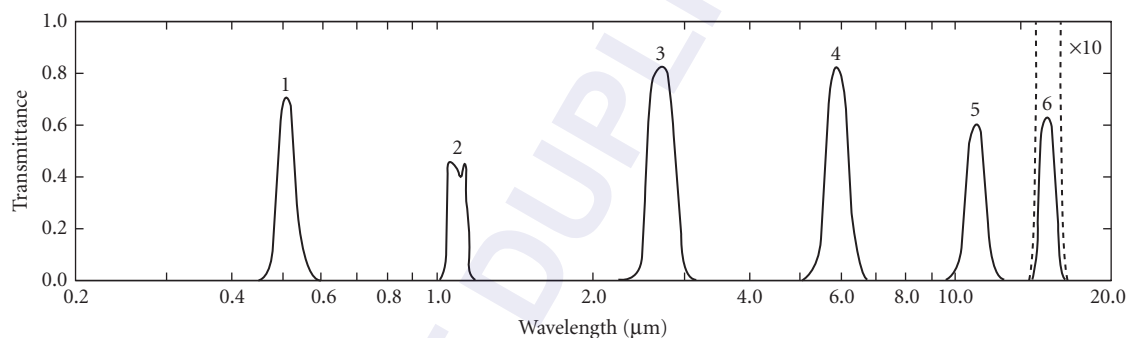


FIGURE 101 Measured transmittances of fully blocked square-top interference filters with half widths of the order of 10 percent. (Curve 1 after Bausch and Lomb;²⁷⁵ curve 2 after Baird Atomic;⁴⁰⁴ curve 3 Infrared Industries;⁴⁰⁶ curve 4 after Turner;⁸³ curve 5 after Eastman Kodak;²⁷⁶ and curve 6 after Optical Coating Laboratory.²⁷⁴)

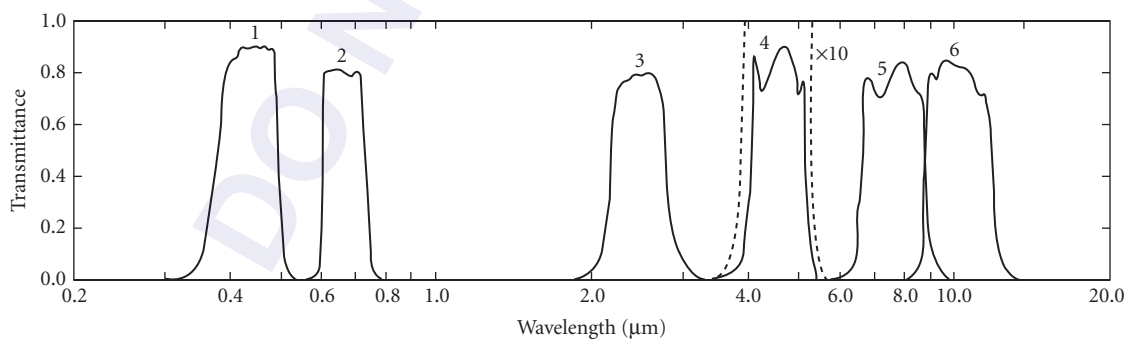


FIGURE 102 Measured transmittances of blocked square-top interference filters with half widths of the order of 25 percent. (Curve 1 after Heliotek;²⁸¹ curves 2 and 4 after Optical Coating Laboratory;²³⁷ curves 3 and 6 after Turner and Walsh⁴⁰¹ and Turner;⁸³ and curve 5 after Infrared Industries.⁴⁰⁶)

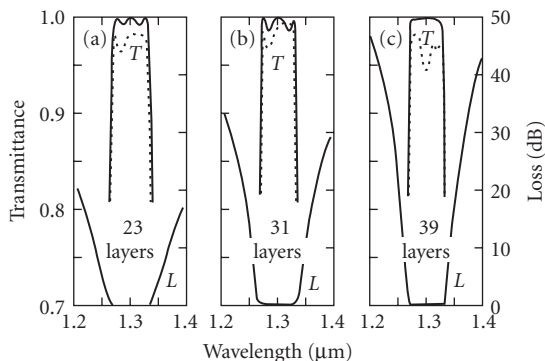


FIGURE 103 Calculated and experimental spectral transmittance and measured attenuation curves of (a) 3-, (b) 4- and (c) 5 cavity bandpass filters. (After Minowa.⁴⁰⁹)

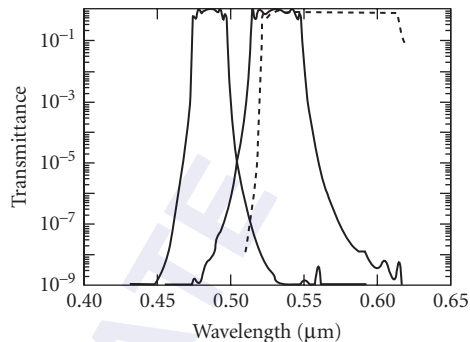


FIGURE 104 Measured spectral transmission characteristics of two bandpass filters for fluorescence applications and of a cutoff filter for Raman spectroscopy. (After Omega Optical Inc.²³⁹)

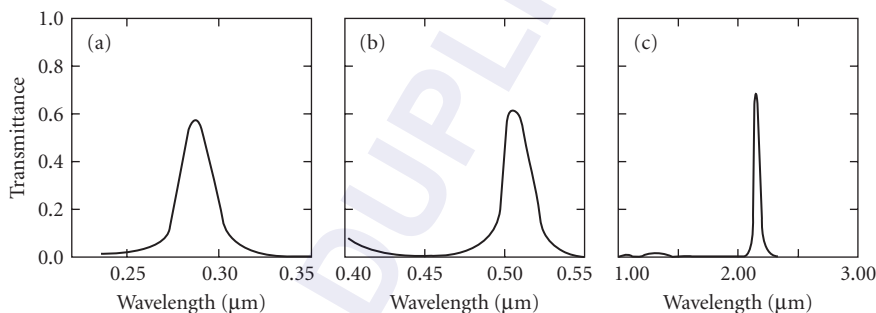


FIGURE 105 Measured transmittance of induced transmission filters for the ultraviolet, visible and infrared spectral regions (a) after Tsypin,⁴¹⁸ (b) after Berning and Turner,⁴¹¹ and (c) after Holloway and Lissberger⁴¹⁴.

when the latter exceeds about two wavelengths, it may become too rough to be useful.³⁹⁶ A high-order filter can, however, be constructed by evaporating reflecting coatings on either side of a thin prefabricated spacer (Fig. 93i).

Mica Spacers ($hw > 0.01$ percent) Transmission bands in silvered mica were probably first observed by Wood⁴²⁰, but the deliberate use of mica to construct filters came much later.^{421–424} The construction of mica interference filters with transmittances of 30 to 80 percent per polarization for half widths of the order of 0.01 to 0.1 percent in the 0.45- to 2.0- μm wavelength region is relatively straightforward.⁴²³ The position of the transmission peak can be located within a fraction of an angstrom, does not change with time, and can be sufficiently uniform over areas of 2 to 5 cm diameter. Because of the very high order of interference (70 to 700 orders) the spectral-free range is quite small, and auxiliary filtering is necessary for most applications. Unless the thickness of the mica is specially selected, the birefringence of mica will result in two mutually perpendicularly polarized sets of transmission bands, a fact that can be used to advantage in some applications. The spectral-transmittance curve of a fully blocked mica interference filter for H_α is shown in Fig. 97c.

Optically polished solid spacers ($hw > 0.002$ percent) It is possible to construct very narrow bandpass filters having thin fused-quartz spacers.^{425–429} A good fused-quartz flat is coated with an all-dielectric reflector, and this coated surface is optically contacted to another flat.^{397,430} The flat is then ground down and polished to form a spacer layer of the required thickness, and the second

reflector layers are applied to complete the filter. As in the mica filters, the position of the transmission band is very stable, and auxiliary blocking filters are needed because of the small spectral-free range. The transmittance of filters with silica spacers is higher than that of corresponding mica filters because fused quartz is highly transparent and is not birefringent. A typical unblocked filter with a clear aperture of 3.5 cm and a half width of 0.007 percent had a transmittance of 45 percent for nonpolarized light. An important advantage of filters with fused-silica spacers is that it has been found possible, by repeating the process described above, to construct square-top filters with rejection ratios of the order of 5×10^{-4} (Fig. 97f) and filters with half widths as low as 0.002 percent (Fig. 97g). However, such filters are very expensive.

Other materials can also be used to produce solid spacers by optical polishing. Germanium was used by Smith and Pidgeon⁴³¹ and by Costich²⁵² to produce very narrow bandpass filters for the infrared. Roche and Title used a substrate made of a combination of yttrium and thorium oxides to produce a filter with a $hw = 0.004$ percent at $3.3 \mu\text{m}$.⁴³²

Plastic spacers ($hw > 0.15$ percent) Mylar has very smooth surfaces and areas can be selected that have a sufficiently uniform thickness to permit the use of this material as a solid spacer. Candille and Saurelle have used this material to produce narrow bandpass filters and obtained half widths of the order of $0.0008 \mu\text{m}$.^{433,434} In some of their designs a second, evaporated narrowband filter deposited onto one side of the solid spacer served to remove unwanted adjacent orders.

Phase-Dispersion Filters ($hw > 0.1$ percent) The dispersion of the phase change on reflection enters into (71) for the half width of FP filters. Typical values of this quantity at $\lambda = 0.5 \mu\text{m}$ for a silver reflector, a nine-layer quarter-wave stack with $n_H/n_L = 1.75$, and for a broadband reflector are -0.5 , -6.8 , and -112.0 , respectively.²²⁴ As a result, the half widths of FP filters constructed with such reflectors should be reduced by factors of about 1.05, 2, and 20. In the last case the contribution of the spacer to the half width is negligible, and a spacerless design is possible (Fig. 93h).^{435,436} Unfortunately, the expected reduction in half width has so far not been fully realized in practice, probably because of errors in the monitoring and lack of uniformity of the layers.⁴³⁷

Tunable filters ($hw > 0.001$ percent) These are usually air-spaced FP interferometers, often provided with elaborate automatic plate-parallelism and spacing control, which are more akin to spectrometers than to filters.^{438,439} The position of the pass band can be tuned quite significantly by changing the separation between the reflector plates. This type of tuning, unlike the tuning of filters by tilting, does not affect the angular field or shape of the transmission band. Ramsey reviews the various problems associated with the construction and use of such instruments.^{440,441}

An electrically tunable 0.005 percent half width interference filter with a lithium niobate spacer sandwiched between two conducting reflecting coatings has also been described (Fig. 106).

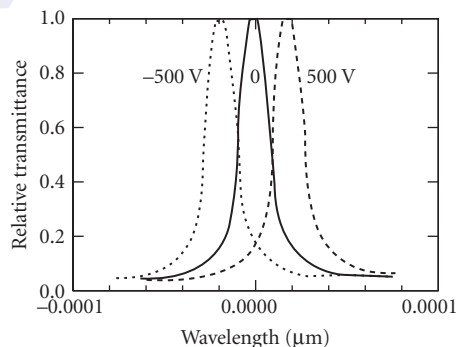


FIGURE 106 Performance of an electrically tunable narrow band filter with a lithium niobate spacer. (After Burton.⁴⁴²)

Wide-Bandpass Filters

Filters with half widths ranging from 10 to 40 percent can be constructed using techniques described in Sec. 7.11, subsection “Square-Top Multicavity Bandpass Filters” (Figs. 101 and 102). Filters with wider transmission bands are usually obtained by combining short- and long-pass filters with cutoffs at the desired wavelengths. The cutoff filters may be all-dielectric (Figs. 61 and 62), glass or gelatine cutoff filters or antireflection-coated infrared materials. Some of the short- or long-pass filters may be regarded as being wide-bandpass filters in their own right. The cutoff filters may be combined into a single filter. Alternatively, it is possible to assemble a number of short- and long-wavelength cutoff filters into sets that make it possible to assemble wide-passband filters of different half widths and peak wavelengths. In this latter arrangement the cutoff positions of all-dielectric short- and long-pass filters can be tuned individually by tilting to coincide with the desired wavelengths. Another way of tuning the edges of the transmission band is to pass the radiation through a pair of circular wedge short- and long-wavelength cutoff filters placed in series.

Filters with very broad transmission bands are also obtained when a multilayer is formed from two suitably displaced long-wavelength cutoff filters separated by an appropriate matching layer (Fig. 107a).⁴⁴³ Automatic computer programs can be used to design high-transmission broadband filters with a high rejection and a shape factor close to unity (Fig. 107b).⁴⁴⁴

Interference Filters with Multiple Peaks

For some applications filters with multiple peaks in one particular spectral region are required. The design of such filters has been considered by Pelletier et al.⁴⁴⁷ Typical results are shown in Fig. 108. Filters with different peak separations, rejections and half widths are possible.

Linear and Circular Wedge Filters

If the thicknesses of all the layers of a bandpass filter vary in proportion across the surface of a substrate, the position of the transmission peak will vary in the same way (Sec. 7.2, subsection “Matrix Theory for the Analysis of Multilayers Systems”). Such wedge filters are as old as the interference filter itself and are available in versions in which the wavelength variation occurs along a straight line or a circle. The latter arrangement is particularly useful because it lends itself well to the construction of low cost, small and light weight rapid-scan monochromators of moderate resolution that are robust and environmentally stable.⁴⁴⁸ Methods for the production of circular variable filters with a linear dependence of wavelength on angle and references to some of the previous work on wedge

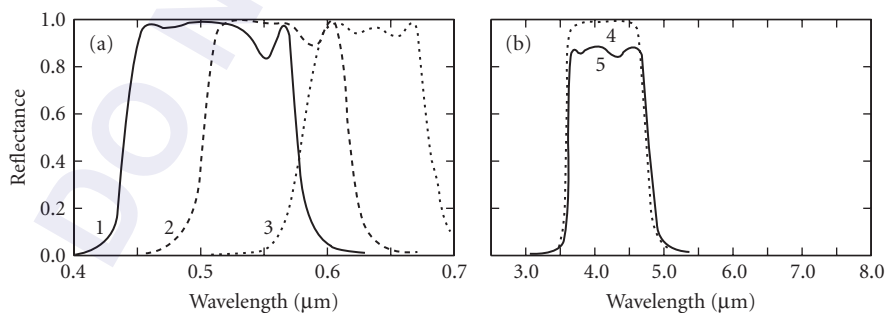


FIGURE 107 Broadband transmission filters: (a) calculated transmittance of three filters consisting essentially of two superimposed suitably tuned long-wavelength cutoff filter structures (after McKenney and Turner⁴⁴⁵) and (b) calculated (curve 4) and measured (curve 5) transmittances of a filter designed with an automatic synthesis program (after Michael⁴⁴⁶).

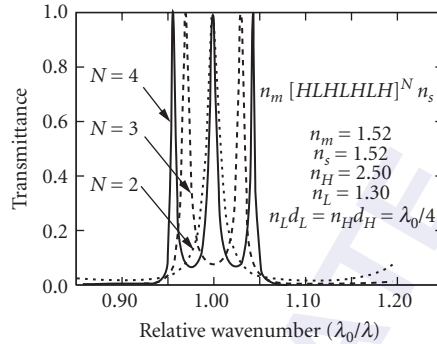


FIGURE 108 Calculated performance of interference filters with two and three closely spaced peaks of the type glass $[HLHLHLH]^N$ glass for values of $N = 3, 4$. The refractive indices of the surrounding media n_m and of the layers n_H, n_L are 1.52 and 2.52, 1.30, respectively. (After Pelletier and Macleod.⁴⁴⁷)

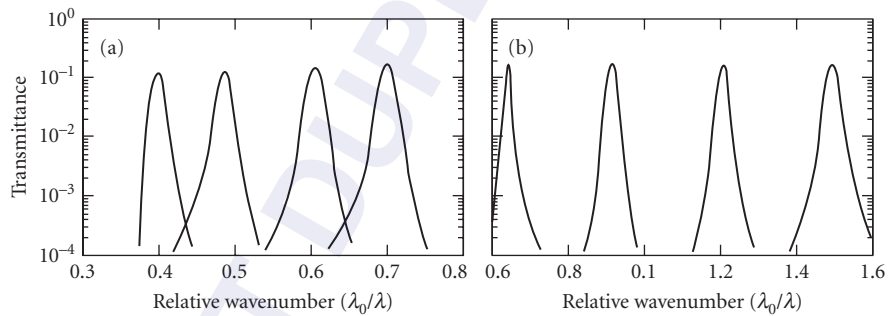


FIGURE 109 The transmission at different angular positions of two circular variable square-top filters for the (a) visible and (b) near-infrared spectral regions. (After Mussett.⁴⁵⁶)

filters are given in several papers.^{449–453} Circular variable square-top filters for the 0.24- to 0.4- and for the 0.4- to 25- μm spectral regions are described by Avilov and by Yen.^{454,455} The maximum transmittances for fully blocked filters vary between 15 and 75 percent, depending on the spectral region and the half width of the filter. Rejection levels of 0.1 or 0.01 percent are possible. Typical transmission curves for several angular positions on two circular variable square-top filters are shown in Fig. 109. The ratio of the maximum to minimum wavelength available in one wheel varies between 1.11 and 16.⁴⁵⁴ The angular width of the slit used in conjunction with a circular variable filter, expressed as a percentage of the angular size of the filter wedge, should not exceed the nominal half width of the filter (expressed in percent) if it is not to cause a marked reduction in the resolution.

Angular Properties of Bandpass Interference Filters

With the gradual increase of the angle of incidence the transmittance maximum of a typical bandpass filter moves toward shorter wavelengths (Fig. 110a). On further increase in the angle of incidence the maximum transmittance and the half width deteriorate; the transmission band becomes asymmetric and eventually splits up into p - and s -polarized components (Fig. 110b). The deterioration is more rapid for nonparallel radiation.

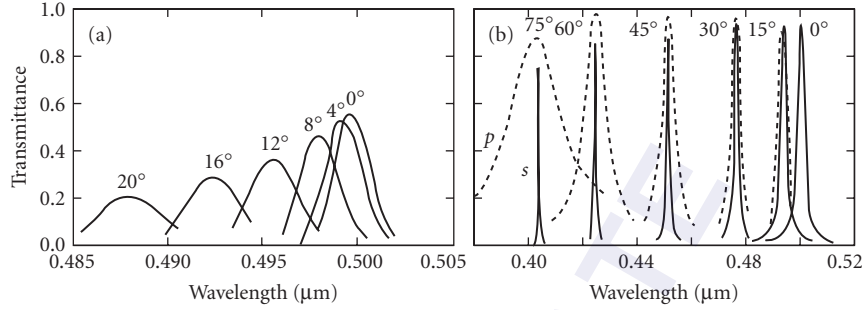


FIGURE 110 Angular properties of all-dielectric interference filters. (a) Measured variation with angle of incidence of the spectral transmittance of a typical commercial interference filter (after Blifford⁴⁰³) and (b) calculated transmittance of a filter in which the peaks of the two polarized transmission bands at nonnormal incidence coincide. The filter is of the type air-[HL]⁴(2A)[LH]⁴-glass, where $n_H t_H = n_L t_L = n_A t_A = \lambda/4$ and $n_s = 1.52$, $n_M = 1.00$, $n_H = 2.30$, $n_L = 1.38$, and $n_A = 1.825$.

Properties of Bandpass Filters for Angles of Incidence Less than 20° The behavior of bandpass filters for angles of incidence $\theta_0 \leq 20^\circ$ can be described quantitatively using the concept of an *effective index* μ^* of the filter. In terms of μ^* the transmittance T in the neighborhood of the transmission peak of any FP filter is given by Lissberger:⁴⁵⁷

$$T \approx \frac{T_0}{1 + \left[\frac{2(\lambda - \lambda_0)}{\Delta\lambda_{0.5}} + \frac{\lambda_0}{\Delta\lambda_{0.5}} \frac{\theta_0^2}{\mu^{*2}} \right]^2} \quad (74)$$

$\Delta\lambda_{0.5}$ and T_0 are the half widths and the maximum transmittance (at λ_0) for normal incidence of the radiation. Formulas for μ^* in terms of the construction parameters have been found for the all-dielectric FP filter and the double-spacer filter^{458,459} and for the metal-dielectric FP and induced-transmission filters.⁴⁶⁰ The change in position of the transmission peak $(\delta\lambda)_\theta$, and the half width $(\Delta\lambda_{0.5})_\theta$ at angle θ are

$$\left(\frac{\delta\lambda}{\lambda_0} \right)_\theta = -\frac{\theta_0^2}{2\mu^{*2}} \quad (75)$$

and

$$\frac{(\Delta\lambda_{0.5})_\theta}{\Delta\lambda_{0.5}} = \left[1 + \left(\frac{\theta_0^2 \lambda_0}{\mu^{*2} \Delta\lambda_{0.5}} \right)^2 \right]^{1/2} \quad (76)$$

For convergent radiation of semiangle α the corresponding expressions are

$$\left(\frac{\delta\lambda}{\lambda_0} \right)_\alpha = -\frac{\alpha^2}{4\mu^{*2}} \quad (77)$$

$$\frac{(\Delta\lambda_{0.5})_\alpha}{\Delta\lambda_{0.5}} = \left[1 + \left(\frac{\alpha^2 \lambda_0}{2\mu^{*2} \Delta\lambda_{0.5}} \right)^2 \right]^{1/2} \quad (78)$$

Linder and Lissberger discuss the requirements and design of filters for this case.^{457,461}

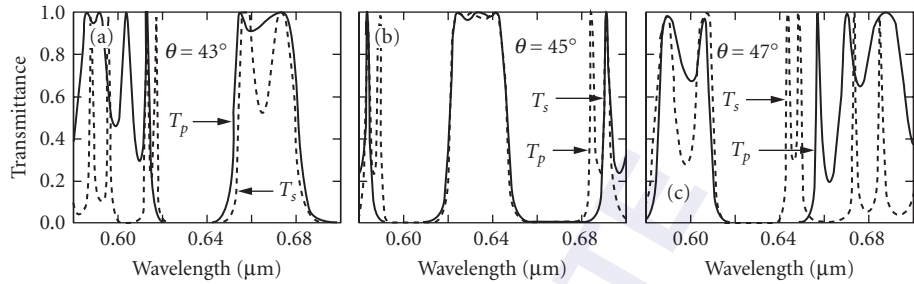


FIGURE 111 Performance at 43° (a), 45° (b), and 47° (c) of a multicity bandpass filter designed to show no polarization splitting for light incident at 45° . (After Baumeister.⁴⁶⁴)

Small tilts are commonly used to tune the peak of a filter to the desired wavelength even though they have an adverse effect on the angular field of the filter.

Bandpass Filters with Little or No Polarization Splitting It has been shown numerically for the phase-dispersion filter,⁴³⁵ for the frustrated-total-reflection filter,⁴⁶² and for the metal-dielectric²⁴ and all-dielectric⁴⁶³ FP filters that it is possible to arrange for the two polarized transmission bands, which may have different widths, to coincide at high angles of incidence (Fig. 110b). The narrow, symmetrical high-transmittance bands that result may be useful for some applications even though the position of the maximum is still displaced with angle. Baumeister has shown how to design multicity filters with no polarization splitting at one angle of incidence (Fig. 111).⁴⁶⁴

Wide-Angle Bandpass Filters Figure 112 shows the variation with effective index μ^* of the angular field of FP filters, defined as being twice the angle of tilt necessary to reduce to $0.8T_0$ the transmittance of the filter for radiation of wavelength λ_0 . To increase the angular field μ^* must be increased. Thus, for example, in an all-dielectric FP filter the expression for μ^* shows that $n_L < \mu^* < n_H$ and that with increasing spacer order μ^* approaches the refractive index of the spacer. The upper limits for μ^* for an all-dielectric FP filter in the ultraviolet, visible, and infrared parts of the spectrum are

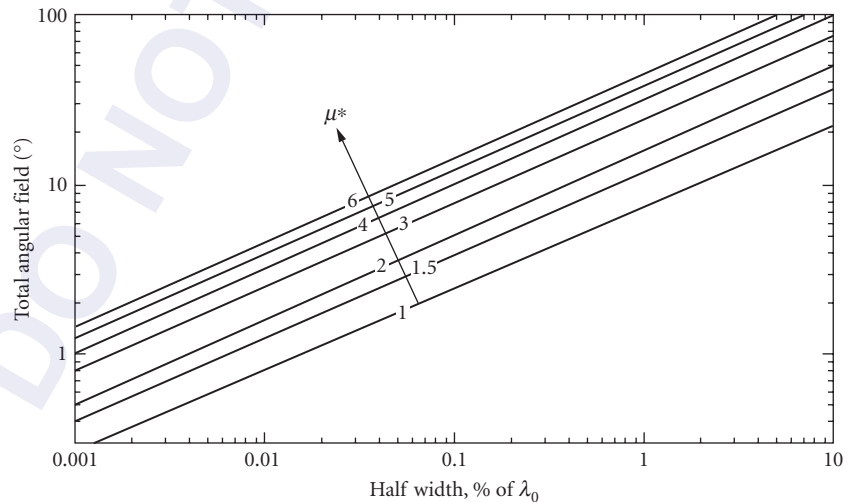


FIGURE 112 Angular field of FP filters as a function of the half width for different effective indexes μ^* .

of the order of 2.0, 2.35, and 5.0, respectively. Little can be done about the angular field of solid spacer filters (Sec. 7.11, subsection “Fabry-Perot Filters with Solid Spacers”).

For metal-dielectric FP and for induced-transmission filters effective indexes μ^* of up to 3.2 and 2.0 have been reported.⁴⁶⁰

Wilmot and Schineller⁴⁶⁵ and Schineller and Flam⁴⁶⁶ have announced a filter consisting of a thin, plane-parallel fiber-optic face plate coated on both sides with all-dielectric mirrors. Since in such a filter the half width is determined only by the thickness of the plate and the reflectivity of the coatings, and since the field of view depends on the ratio of the wavelength to the fiber diameter, the two quantities are independent. The measured transmittance of a 15-Å-half-width, 6-mm-diameter filter composed of 1.5- μm -diameter fibers was 30 percent, and the shift in wavelength with angle of incidence was one-eighth that of a conventional filter.

Stability and Temperature Dependence of Bandpass Filters

The stability of the position of the transmission peak has been studied by many workers.^{389,396,467-475} The observed changes (up to 1 percent of λ_0) seem to depend greatly on the materials and manufacturing conditions. In filters with evaporated spacers both irreversible changes, probably due to changes in the structure of the films, and reversible changes due to the absorption of water vapor, have been observed. Many manufacturers are now able to minimize these effects through the use of more stable materials, improved high energy deposition methods (see Sec. 7.3), or accelerated artificial ageing processes. No changes were observed in solid spacer filters (Sec. 7.11, subsection “Fabry-Perot filters with Solid Spacers”).

Changes in the operating temperature normally do not significantly affect the half widths and peak transmittances of medium and wide bandpass interference filters (see, for example, Refs. 476 to 479). An exception are filters that contain semiconductors that start to absorb significantly on heating (e.g., germanium) or on cooling (PbTe).^{480,481}

The position of the transmission peaks shift linearly toward longer wavelengths with an increase in temperature, the magnitude of the shift depending largely on the spacer material. The temperature coefficient, expressed as a percentage change in λ_0 per degree Celsius change in temperature, lies between 2×10^{-4} and 3×10^{-3} for filters with evaporated spacers for the 0.3- to 1.0- μm spectral region⁴⁰³ and between 2×10^{-3} and 2×10^{-2} for the infrared spectral region.^{389,237} It is of the order of 1×10^{-3} for filters with mica⁴²³ and quartz³⁹⁷ spacers. Unless temperature control is provided, under adverse conditions all the above temperature coefficients could lead to serious shifts of the transmission peaks of very narrow band filters (Sec. 7.11 subsection “Very Narrow Bandpass Filters”). The temperature control can take the form of an external constant-temperature enclosure, or it might be built in right into the filter. Eather and Reasoner³⁹⁶ and Mark et al.⁴⁸² describe arrangements of the latter type in which two sensors embedded in the filters are used to control the current flowing through two transparent conducting coatings that surround the filter.

Deliberate changes in the temperature can be used for a fine tuning of the transmission wavelength without having an adverse effect on the angular field of the filter.

Bandpass Filters for the XUV and X-Ray Regions

The construction of good bandpass filters for the extreme ultraviolet is hampered by the lack of coating materials with suitable optical constants. However, certain metals in thin-film form can be used as rudimentary bandpass filters in the extreme ultraviolet. The primary process in these filters is absorption, although at times interference within the film may have to be considered to explain the spectral transmission characteristics fully.

The measured spectral transmittance of some of these materials are shown in Figs. 113 and 114. By increasing the thicknesses of the films higher rejection ratios could be obtained at the expense of peak transmissions, and vice versa. The transmittance of the most promising material, aluminum, would be higher if it were not for the formation of absorbing oxide layers.

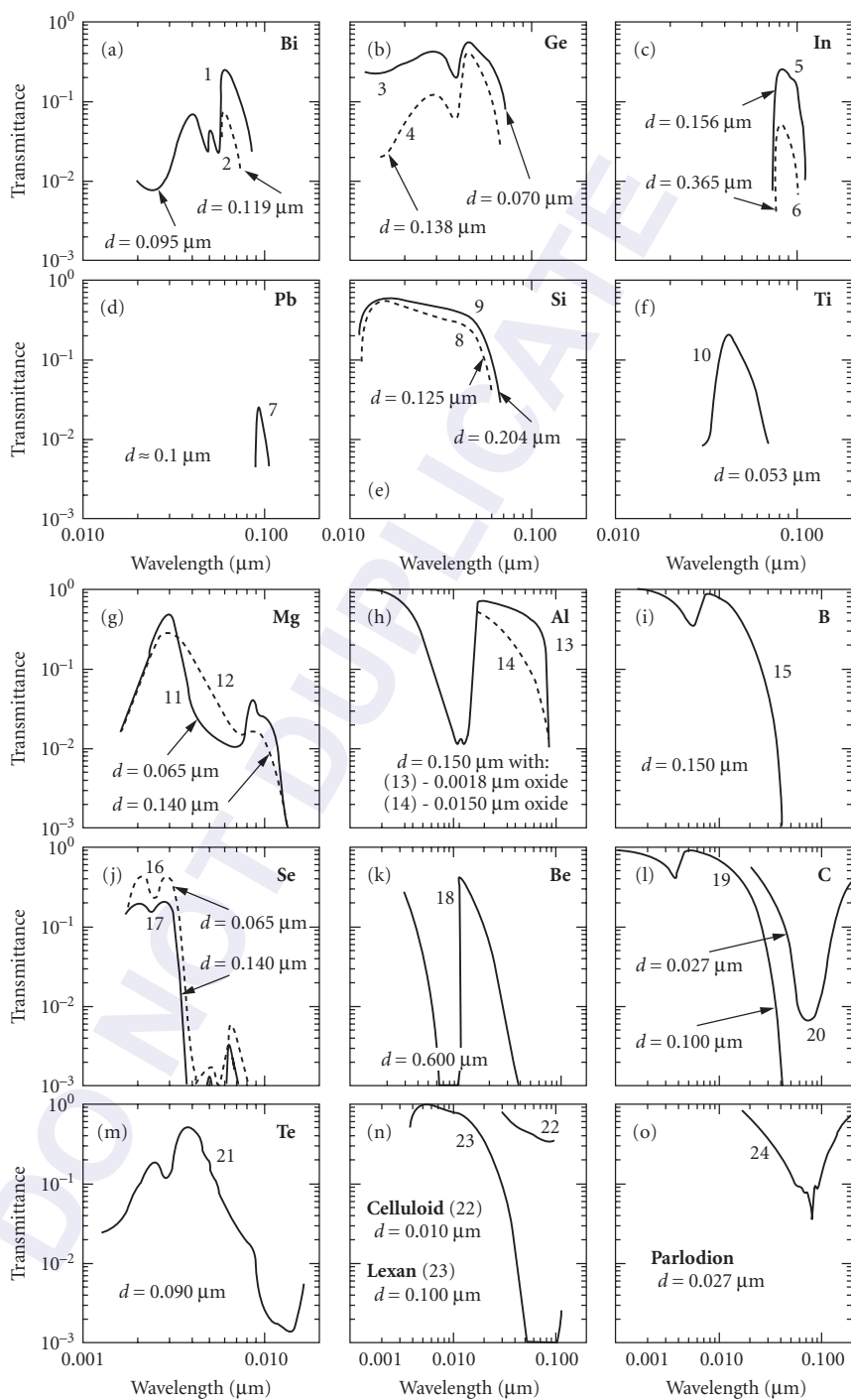


FIGURE 113 Measured extreme-ultraviolet and soft x-ray transmittances (a) to (o) of several self-supporting metal films of indicated thicknesses.

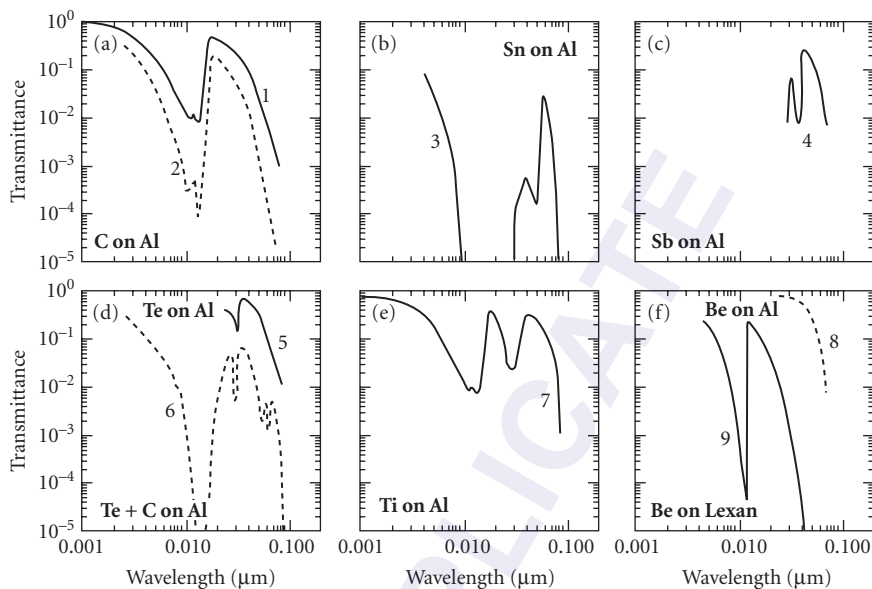


FIGURE 114 Measured extreme-ultraviolet and soft x-ray transmittances (a) to (f) of six metal films of indicated thicknesses deposited onto thin aluminum or plastic films.

Many of the layers are self-supporting (Fig. 113). Others must be deposited onto a suitable transparent substrate (Fig. 114). Thin aluminum films are sometimes used for this purpose. Other materials used in the past are Zapon (cellulose acetate); collodion, Parlodion, and Celluloid (cellulose nitrates); Mylar (polyethylene terephthalate); and Formvar (polyvinyl formal) (Fig. 113n and o). Any residual absorption in the substrate contributes, of course, to the overall-transmission curve. The preparation of self-supporting thin films is described by Novikov and by Sorokin and Blank.^{483,484} Because of their fragility such films are usually mounted on a very fine mesh screen.

Multilayer Fabry-Perot interference filters for the soft x-ray region have also been constructed.⁴⁸⁵⁻⁴⁸⁷ However, thus far the only spectral measurements reported are nonnormal incidence reflection.⁴⁸⁸ The finesse of the filters is low and the modulation of the reflectance curve depends on the thickness of the spacer. The devices are useful for measurement purposes.

7.12 HIGH PERFORMANCE OPTICAL MULTILAYER COATINGS

Advances in the design and the automation of the manufacture of optical multilayer coatings have made possible the routine production of filters with a performance that would have seemed impossible to achieve even fifteen years ago. For the visible and near-infrared parts of the spectrum materials and processes have been developed at a number of commercial and research facilities which permit them to deposit with high accuracy layer systems that consist of more than 4000 layers and with total metric thicknesses exceeding 100 μm . Currently the most common application areas for such coatings are in telecom systems (with filter sizes typically of the order of 1 mm) and in fluorescence and Raman spectroscopy (diameters between 10 and 25 mm). The filter specifications for these filters, as a rule, span such a range of transmittances, or have such sharp gradients, that to depict the results properly, it is usually necessary to present the data in at least two graphs. In this section representative examples will be given to illustrate the type of performance that can be currently achieved.

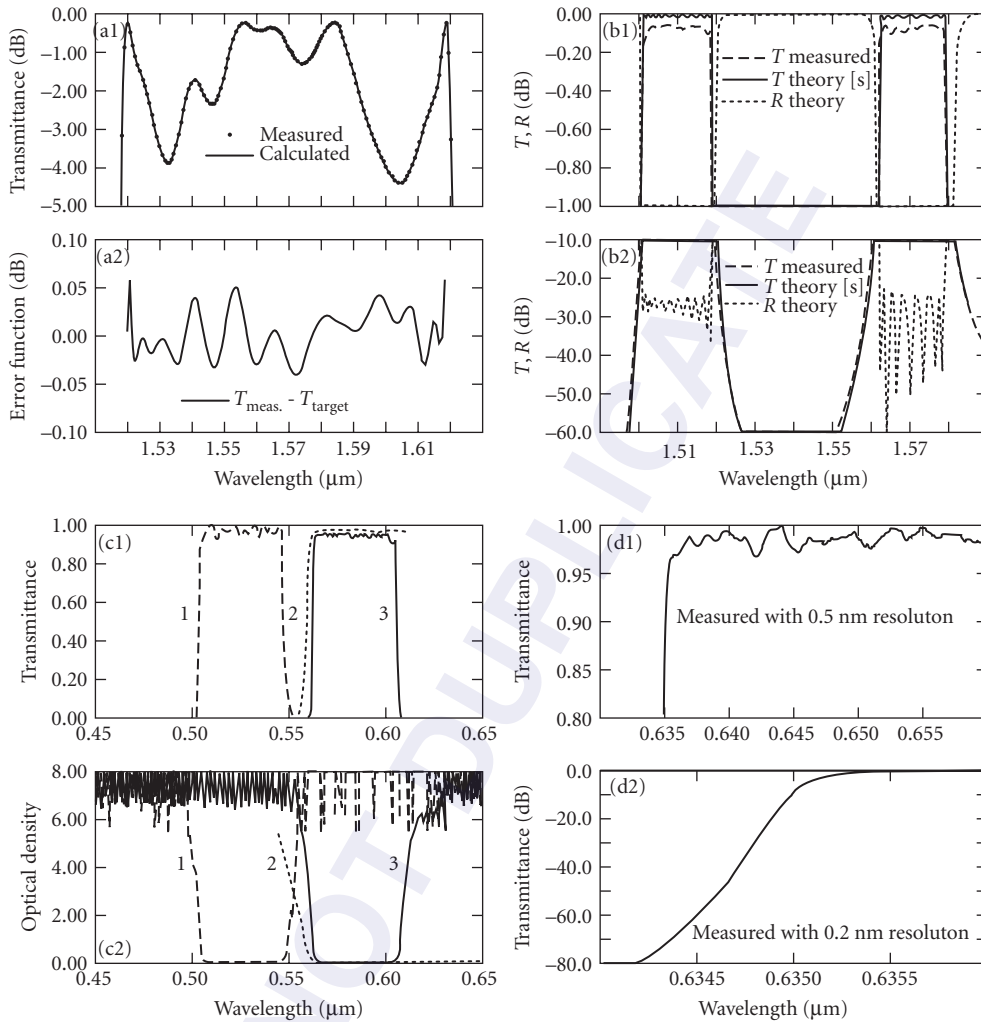


FIGURE 115 Measured data for some high performance optical filters. (a) A precise telecom gain flattening filter with an attenuation that is defined over four decades; (b) dual band filters with very high transmission and reflection regions, for adding or removing signals of certain wavelengths in fiber networks; (c) a set of three filters for fluorescence probe microscopy, having very high transmittances and sharp drop-offs; curves 1, 2, and 3 correspond to the transmittance curves of the excitation, beam splitting and emission filters; (d) Raman edge pass filter with an edge steepness of approx. 86 dB/nm. (After Iridian Spectral Technologies.²⁰²)

In Fig. 115a1 is shown the calculated and measured performance of a 50- μm -thick gain-flattening filter for telecom applications in which the specified transmittance varies in the 1.530- to 1.625- μm spectral region between 0.0 and -4.0 dB. The difference between the target and measured transmittances is shown in Fig. 115a2 and the peak-to-peak variation is <0.1 dB.

Figure 115b1 presents the measured transmittance of a dual band filter of the type that is used for adding or dropping signals of certain wavelengths from a fiber network. Also shown are the theoretical transmittance and reflectance. The main requirements for this application are that the transmittance be as high as possible with low ripple in the transmission bands and have low transmittance

between these bands. The filter shown consisted of more than 300 layers and had a total thickness of 65 μm . Figure 115.b2 shows that the rejection between the two transmission peaks is indeed very high.

In Fig. 115c are shown the measured transmittance curves of the excitation and emission filters and the beam splitter required for fluorescence probe microscopy. For this application the filters must have very high transmittances and sharp drop-offs in order to maximize the signal-to-noise ratio. The emission and excitation filters were about 35 μm thick and consisted of more than 300 layers. The optical density outside the transmission regions are shown in Fig. 115c2. Some of these measurements are noise limited above an optical density 6.

In Fig. 115d is shown the measured transmittance of a steep cutoff filter for Raman spectroscopy. The filter consisted of about 500 layers and had a total thickness of about 40 μm . The ripple in the transmission region is only 4 percent. To appreciate the steepness of the cutoff (about 86 dB/nm) the region around the cutoff is plotted in Fig. 115c2 on an expanded wavelength scale.

Very advanced spectroscopic notch filters are shown in Fig. 42, one of these consisting of as many as 4410 layers. Of course, the processes used to produce the advance coatings described above are used today to produce better coatings with the simpler properties described in the earlier diagrams.

7.13 MULTILAYERS FOR TWO OR THREE SPECTRAL REGIONS

Increasingly there are applications, especially in laser science, in which the spectral transmission and/or reflection has to be controlled at two or more wavelengths.³⁴ The design and construction of such coatings is more difficult than that of systems for one wavelength region only, especially when the ratio of the wavelengths of interest is very large.

Multilayers for Two Spectral Regions

Costich was the first to specify the construction parameters of coatings having all possible combinations of low and high reflection behavior at wavelength ratios of 1.5-, 2.0-, and 3.0: 1.0.⁴⁸⁹ His solutions were based on systems composed of quarter-wave layers or layers with other simple thickness relationships. Experimental results are in good agreement with the calculated values. The calculated performance of two coatings not shown before are given in Figs. 116 and 117 presents the performance of a number of commercially produced coatings of this type. Systems for other wavelength ratios and for combinations of other reflection values are possible.

Sometimes coatings are required in which the reflectance is controlled for wavelength ratios of 10:1 or more.^{490,491} A systematic method for the design of such multilayers with different reflectance characteristics at the two wavelengths has been described.^{492,493} The calculated performance of such coatings for the wavelengths of 0.6328 and 10.6 μm are shown in Fig. 118.

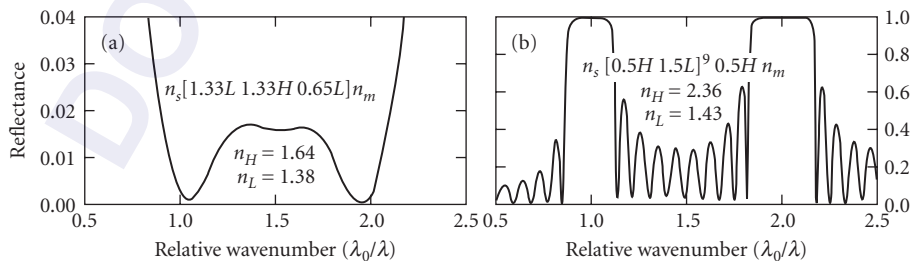


FIGURE 116 Calculated performance of a commercial antireflection coating (a) and a high reflectance coating (b) for wavenumbers 1.0 and 2.0 μm^{-1} . (After Costich.³⁴)

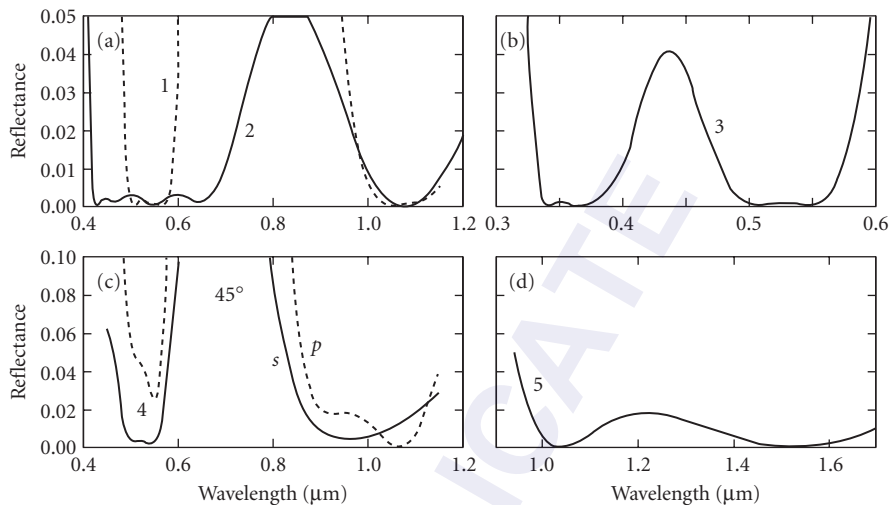


FIGURE 117 Measured reflectances (a) to (d) of some commercial two-wavelength antireflection coatings. (Curves 1, 3, and 4 after TechOptics⁸⁶ and curves 2 and 5 after Thin Film Lab.⁸⁵)

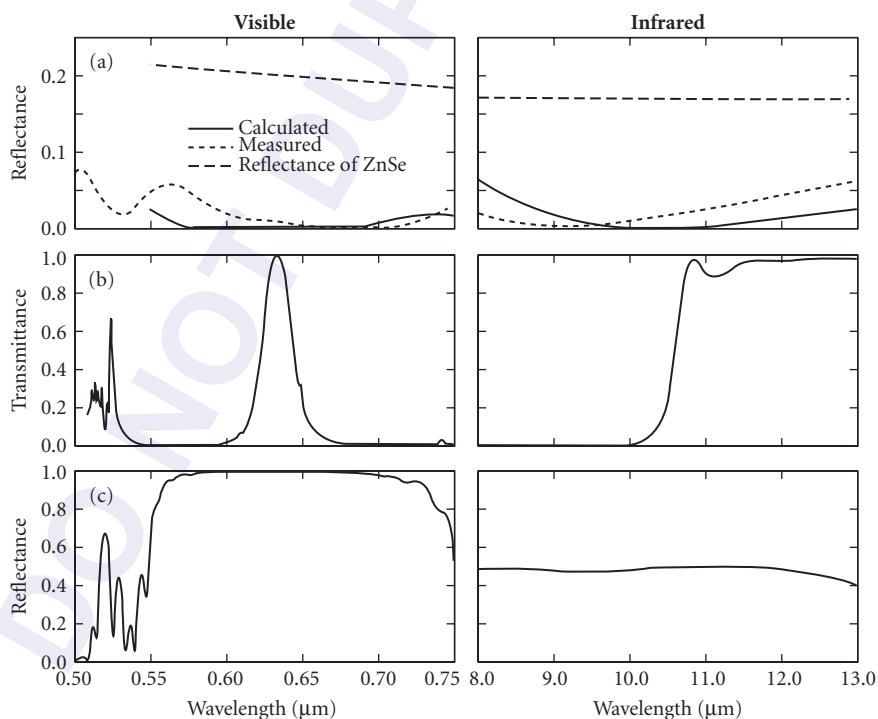


FIGURE 118 Calculated performance of three multilayer coatings (a), (b), and (c) designed for two widely separated spectral regions. Columns 1 and 2 represent the performance of the multilayers in the visible and in the infrared spectral regions, respectively. The experimental measurements for one coating are also shown. (After Li.^{492,493})

Multilayers for Three Spectral Regions

For some laser applications the reflectance or transmittance has to be controlled at three or more wavelengths. Solutions to such problems can also be found. Costich has given designs for all possible combinations of low and high reflection for the important special case of a set of wavenumbers σ_0 , $2\sigma_0$, and $3\sigma_0$.⁴⁸⁹ The designs and performances of his solutions to this problem are shown in Fig. 119.

In principle, the method for the design of coatings for two widely separated spectral regions mentioned above can be extended to the design of coatings for three or more wavelengths. However, the number of layers required increases dramatically as the number of layers required for the longest wavelength region increases. Figure 120 shows the calculated performance of a coating that behaves like a high reflection coating, a beam splitter and an antireflection coating at 0.63, 2.52, and 10.6 μm , respectively.

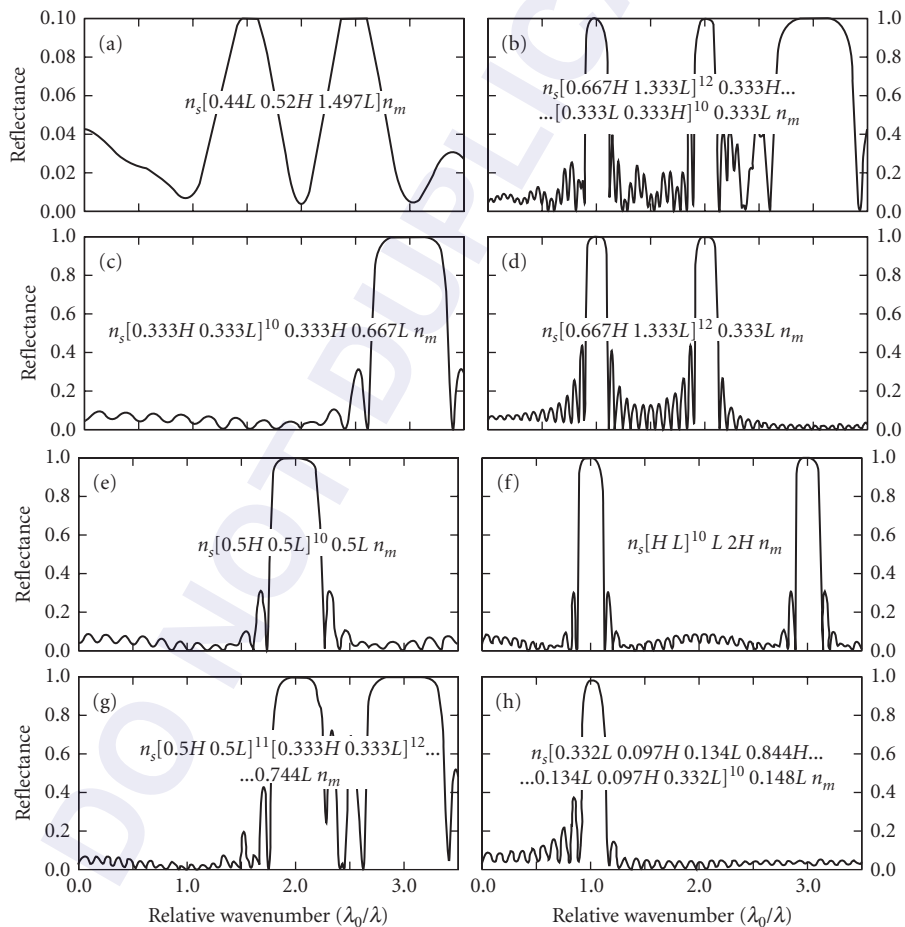


FIGURE 119 Calculated spectral reflectances of multilayer coatings on glass with various combinations of high and low reflectance at relative wavenumbers λ_0/λ 1.0, 2.0, and 3.0. In the designs H , L correspond to quarterwave layers at $\lambda_0 = 1.0 \mu\text{m}$. n_m , n_s , n_{1P} and n_L were assumed to be 1.00, 1.52, 1.95, and 1.43, except in (a) where n_{1P} , n_L were 1.64, 1.38, respectively. (After Costich.³⁴)

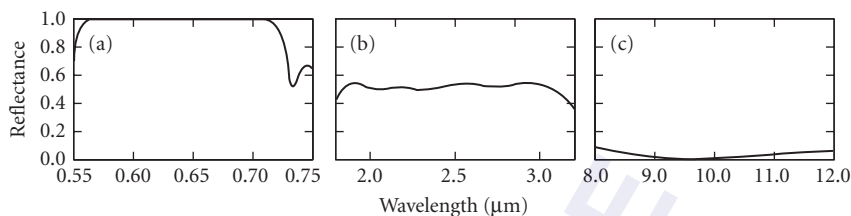


FIGURE 120 Calculated performance of a multilayer with different properties in three spectral regions. (After Dobrowolski.⁴⁹⁴)

7.14 PHASE COATINGS

In some applications, in addition to transmittance or reflectance requirements, special phase relationships have to be satisfied. These may be specific phase changes on reflection ϕ_R or transmission ϕ_T for radiation incident at 0° . At other times it is required to displace or to deflect a beam without affecting its state of polarization. However, most frequently it is necessary to introduce a certain phase difference ($\phi_p - \phi_s$) between p - and s -polarized light. Quarter-wave plates made of birefringent crystals are normally used to provide this phase difference.

Solutions to these and similar problems based on optical interference coatings can also be found. Porous films with an inclined columnar structure, formed in physical vapor deposition processes when the vapor is incident onto the substrate at an oblique angle, can also be birefringent.⁴⁹⁵ Such films have been proposed for the construction of phase retardation plates for use with normal incidence of the radiation.⁴⁹⁶ However, more frequently solutions are based on the difference between the effective indices η_p, η_s of thin films for obliquely incident radiation [Eq. (19)]. Azzam has shown that, when performance at one wavelength only is specified and when oblique angles of incidence are acceptable, elegant solutions to many of the problems can be found that are based on a single layer only.³³⁵

Phase retarding reflectors are commonly designed for use at 45° . Many layers are required when the radiation is incident from the air side. The performances of two multilayers of this type with different phase differences are shown in Fig. 121. The multilayers are optimized to maintain a constant phase difference in the vicinity of the design wavelength. Coatings with other phase differences and

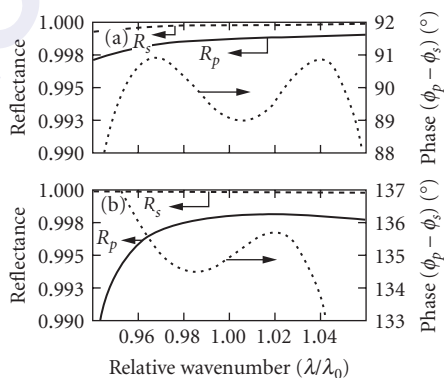


FIGURE 121 Calculated performance of two front surface 45° reflectors with different phase retardations. (a) 20 layers on a Ag substrate (after Southwell⁴⁹⁹) and (b) 22 layers on an Al substrate (after Grishina⁵⁰⁰).

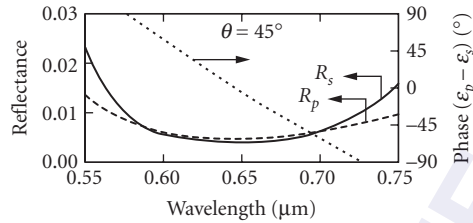


FIGURE 122 Antireflection coating on glass for use at 45° with zero phase retardation at $\lambda = 0.6471 \mu\text{m}$. (After *Thin Film Lab.*⁸⁵)

reflectances can also be constructed.^{497,498} For example, an antireflection coating for 45° incidence in which the phase change is 180° is shown in Fig. 122.

When the radiation is incident on the layers from the substrate side, total internal reflection takes place.⁵⁰¹ The design of thin-film phase retarders based on this approach has also been examined by Apfel⁵⁰² and Azzam.⁵⁰³ Total internal reflection phase retarders operate over broader spectral regions (Fig. 123) but their size is limited by the weight and homogeneity of the prism materials.

More complex phase retardation devices have been constructed in which the radiation is allowed to undergo 2, 3, or even 4 internal reflections.^{342,506,507} The performance of some typical total internal reflection devices that are based on the configurations of Fig. 68*u, v, w* is shown in Fig. 124. For high-power laser beam delivery systems front surface reflectors are usually employed (Fig. 125).

In Fig. 126 are shown the phase changes on reflection of a set of four metal/dielectric interferometer mirrors, all with $R > 0.97$, for which the differences in the phase changes on reflection for adjacent members in the set were approximately 90° over an extended spectral region. Other requirements for phase changes or phase change differences can also be satisfied with thin films.

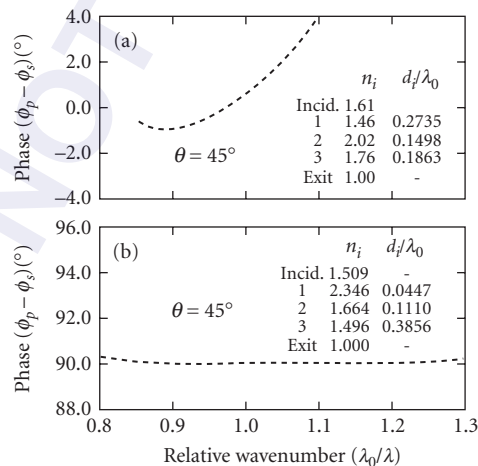


FIGURE 123 Phase retarders based on total internal reflection, consisting of 3 layers on glass and operating at an angle of incidence of 45° . (a) 0° phase retardation (after Cojocar⁵⁰⁴) and (b) 90° phase retardation (after Spiller⁵⁰⁵).

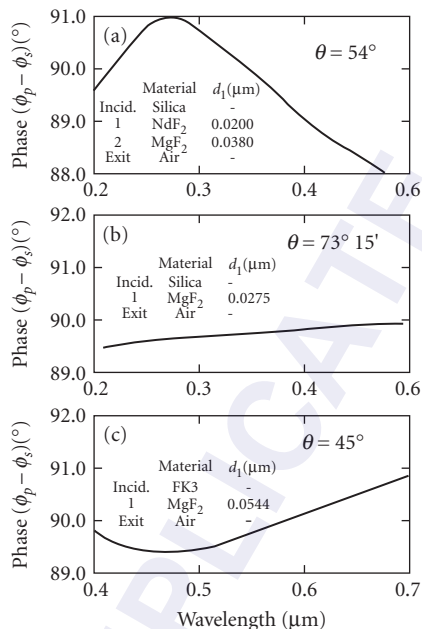


FIGURE 124 90° phase retardation devices based on 2 (a), 3 (b) and 4 (c) total internal reflections (after King,⁵⁰⁸ Clapham,³⁰⁴ and Filinski⁵⁰⁷). The angle of incidence on the first reflecting surface is indicated in the diagrams.

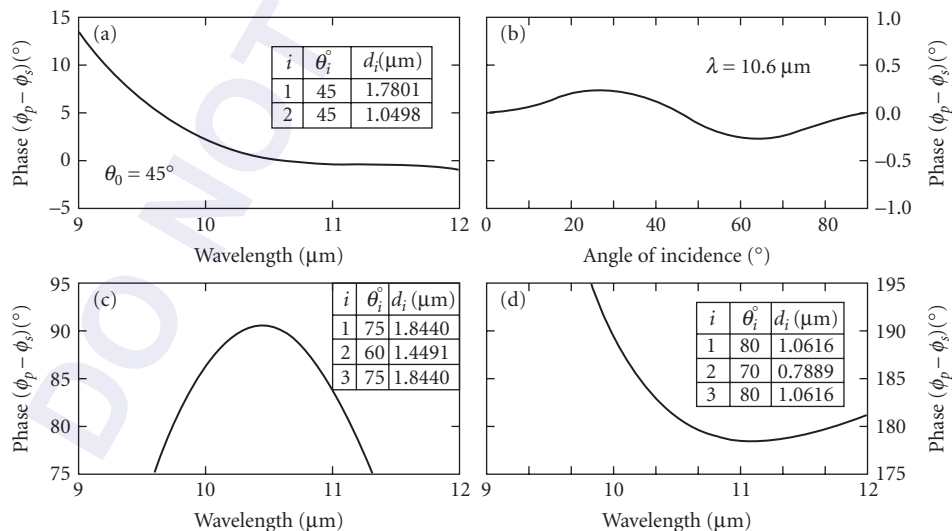


FIGURE 125 Phase retardation devices for $\lambda = 10.6 \mu\text{m}$ based on multiple reflections from surfaces coated with opaque silver films and single layers of ZnS layers having specified thicknesses. (a), (b) 0° phase retardation for all angles of incidence (after Azzam⁵⁰⁹); (c), (d) 90°, 180° phase retardations (after Thonn⁵¹⁰).

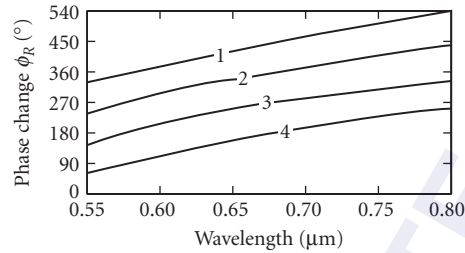


FIGURE 126 Normal incidence phase changes on reflection of a set of four highly reflecting mirrors for Michelson interferometers. (After Piotrowski et al.²²)

7.15 INTERFERENCE FILTERS WITH LOW REFLECTION

Reducing Reflection with a Thin Metal Film

A suitable thin metal film deposited onto a glass surface can act as a very efficient achromatic anti-reflection coating for light incident from the glass side (Fig. 127). The reflectance for light incident from the air side is not reduced and the transmittance suffers as a result of the absorption within the film.⁵¹¹ By combining such films with additional layers, attractive coloured sun glasses or architectural coatings are obtained.

Black Absorbers

Black absorbers efficiently absorb the radiation incident upon them in a specified spectral region. They are used, for example, to control radiant energy,⁵¹² to remove stray light in optical systems, to enhance contrast in display devices⁵¹³ and to increase the signal-to-noise ratio in multiplexers.⁵¹⁴ Black absorber coatings are based on interference in thin film and generally consist of an opaque metal layer and one or more dielectric layers interspersed with partially transparent metal layers (Figs. 21 and 66a). They can be designed for first- and second-surface application (Fig. 128). Coatings of this type can also be designed for the ultraviolet and infrared spectral regions.⁵¹⁵

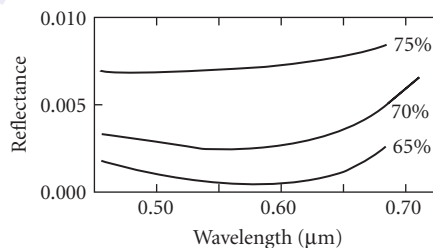


FIGURE 127 Spectral reflectance of thin chromium films on glass for light incident from the substrate side. The transmittance of the layers at $\lambda = 0.565 \mu\text{m}$ is indicated. (After Pohlack.⁵¹¹)

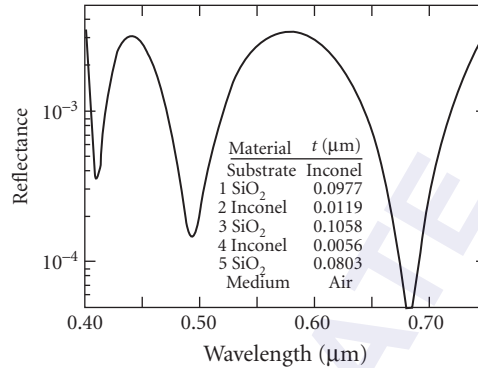


FIGURE 128 Calculated performance of a five layer metal/dielectric black absorber. (After Dobrowolski.⁵¹⁶)

Neutral Attenuators

Conventional metallic film attenuators described in Sec. 7.9, subsection “Neutral Filters” cannot readily be placed in series because multiple reflections between the components may result in unpredictable density values. However, by using metal and dielectric layer combinations of appropriate optical constants and thicknesses, it is possible to reduce the reflection of the metallic film from one or both sides of the substrate.^{82,513–515,517,518} The experimental results for one such attenuator are given in Fig. 129b.

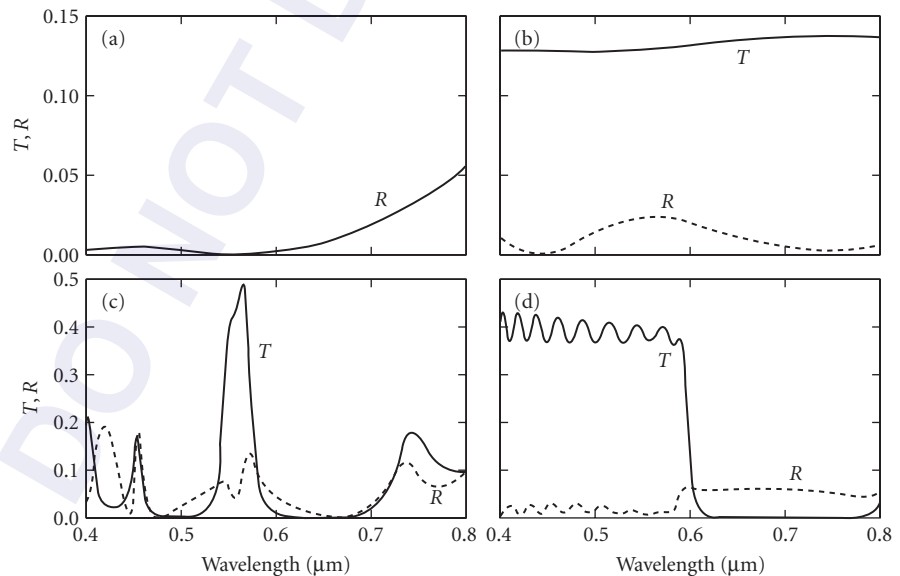


FIGURE 129 Measured performance of interference filters with reduced reflection: (a) black absorber; (b) neutral attenuator; (c) narrow bandpass filter; and (d) long-wavelength cutoff filter. (After Dobrowolski⁵¹⁹ and Sullivan.⁵²⁰)

Other Interference Filters

It is possible, using a similar approach, to reduce the reflection of narrow bandpass filters, cutoff filters, and other filter types. In particular, low reflection narrowband interference filters for welding applications have been described by Jacobsson.⁵²¹ The experimental performance of a bandpass filter and of a long-wavelength cutoff filter are given in Fig. 129c and d. In both cases the luminous reflectance has been reduced by an order of magnitude over that of a conventional design. However, this is at the expense of the transmittance.

7.16 REFLECTION FILTERS AND COATINGS

Metallic Reflectors

The Fresnel reflection coefficient of an interface between two semi-infinite media of complex refractive indices \tilde{n}_m, \tilde{n}_s for polarized radiation incident at nonnormal angle is given by

$$R = \left| \frac{\eta_s - \eta_m}{\eta_s + \eta_m} \right|^2 \quad (79)$$

where η_s, η_m are given by (19). When \tilde{n}_m, \tilde{n}_s correspond to air, and the metal, respectively, and when the angle of incidence is zero, the above expression reduces to

$$R = \frac{(n_s - 1)^2 + k_s^2}{(n_s + 1)^2 + k_s^2} \quad (80)$$

If the substrate is opaque, the above represents the total energy reflected, the remaining energy being absorbed within the material.

Metal reflectors are most commonly made by vacuum deposition of the material onto a suitable glass or quartz substrate. Before deposition aluminum or beryllium mirror surfaces are sometimes first chemically plated with a nickel-phosphorus alloy (Kanigen process). Such deposits have excellent adhesion to the substrate and have a very hard surface that can be optically polished before coating.⁵²²

Visible, infrared, and ultraviolet spectral reflectances of some of the more commonly used metals are shown in Figs. 130 and 131. Using Eq. (80) and the optical constants in Palik's handbook,^{58,59} the

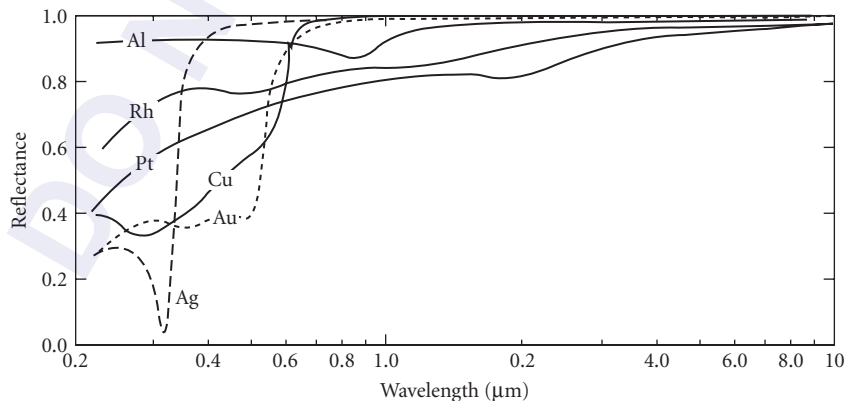


FIGURE 130 Reflectances of some metals. (After Drummeter and Hass.¹³⁸)

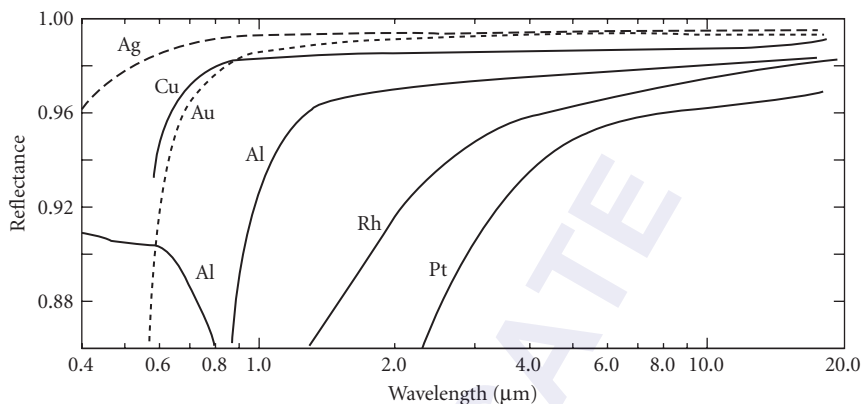


FIGURE 131 Visible and infrared reflectance of certain metals. (Al after Bennett et al.,⁵²⁵ Ag and Au after Bennett and Ashley,⁵²⁶ Cu after Hass and Hadley,⁵²⁷ Rh and Pt after Hass and Fowler.) (See Drummeter and Hass.¹³⁸)

spectral reflectances for many additional metals can be calculated. Silver has the highest visible and infrared reflectance, and hence is used for interferometer mirrors and interference filters. Exposed silver films tarnish readily but they can be protected. Aluminum has the broadest high reflectance region of all metals and is commonly used in front-surface mirrors. It would reflect highly down to 0.1 μm were it not for the absorption below 0.18 μm of the thin oxide layer that starts to form seconds after deposition.^{523,524} Some of the highest known reflectances in the ultraviolet are shown in Fig. 132.

Protective Coatings For many applications the thin aluminum oxide layer on an aluminum surface does not offer sufficient protection against abrasion and chemical attack, and therefore aluminum mirrors are often overcoated with single SiO_2 or of MgF_2 protective layers. Such mirrors can be repeatedly cleaned with water and even withstand boiling in saltwater.⁵³⁴ Single protective layers reduce the reflectance (Fig. 133). If necessary this problem can be overcome by using protective coatings consisting of two or more layers (Fig. 134).

The deterioration of the ultraviolet reflectance of aluminum mirrors due to oxidation can be partially avoided by covering the freshly deposited aluminum layer immediately with a suitable

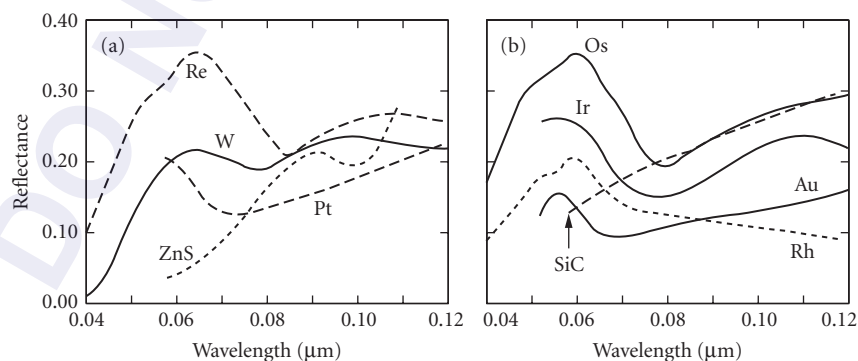


FIGURE 132 Measured ultraviolet reflectance of certain materials (a) and (b). (Pt, Au, and ZnS after Hunter,⁵²⁸ Ir after Hass et al.,⁵²⁹ Rh after Cox et al.,⁵³⁰ Re, and W after Cox et al.,⁵³¹ Os after Cox et al.,⁵³² and SiC after Seely.⁵³³)

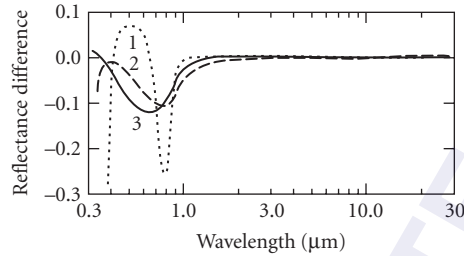


FIGURE 133 Measured difference between the reflectances of protected and unprotected aluminum mirrors. Curve 1: coating 4 of Fig. 134b; curves 2 and 3: 0.1122 ± 0.002 and 0.0752 ± 0.001 μm thick films of MgF_2 and SiO_2 , respectively, on aluminum. (After Bennett.⁵²⁵)

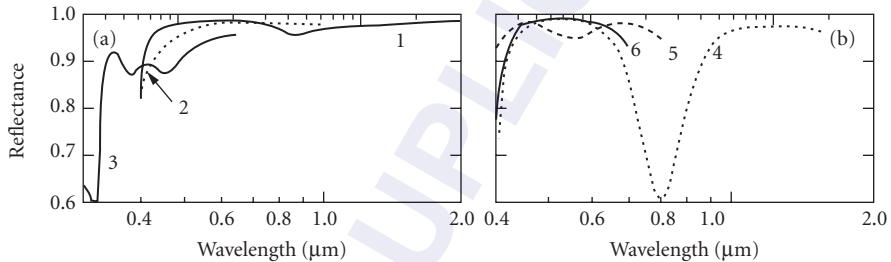


FIGURE 134 Reflectance of very durable overcoated metal mirrors. (a) Silver mirrors: curve 1: protected front-surface mirror (after Denton⁵³⁵); curve 2: enhanced reflection (after Vvedenskiĭ⁵³⁶); curve 3: extended reflection (after Song *et al.*⁵³⁷). (b) Aluminum mirrors: curve 4: with four layers of MgF_2 and CeO_2 (after Hass⁵³⁴); curve 5: with four layers of SiO_2 and TiO_2 (after AIRCO⁵³⁸); and curve 6: with four layers MgF_2 and ZnS (after Furman and Stolov⁵³⁹).

coating of MgF_2 ^{540–542} or LiF .^{540,543,544} The reflectances of such overcoated aluminum reflectors are shown in Fig. 135. Their variation with angle of incidence in the 0.03 to 0.16 μm spectral region is discussed by Hunter.⁵⁶

The reflectance of unprotected and protected silver mirrors has been investigated by Burge *et al.*⁵⁴⁵ Highly adherent and chemically stable mirrors with a reflectance in excess of 0.95 for wavelengths greater than 0.5 μm have been reported.^{546,547}

The reflectance of protected metal mirrors at oblique angles of incidence in the infrared part of the spectrum can be seriously reduced at the short-wavelength side of the Reststrahlen peak of the material used for its protection.^{548–551}

Enhancement of Reflection

By depositing a quarter-wave stack (Sec.7.6 subsections “Nonabsorbing $[AB]^N$ and $[AB]^N A$ Multilayer Types” and “Periodic Multilayers of the $[(0.5A) B (0.5A)]^N$ Type”) onto the metal mirror its reflectance can be enhanced considerably.³⁸⁷ The thickness of the first layer should be adjusted to compensate for the phase change on reflection at the metal surface.⁵⁵² The spectral-reflectance curve dips on either side of the high-reflectance region whose width is governed by the considerations of Sec. 7.6, subsection “Width of the High-Reflectance Zone” and which can be somewhat enhanced by the use of a half-wave outermost layer (Fig. 134). The measured spectral characteristics of three metal-dielectric reflectors for the ultraviolet region are shown in Fig. 136.

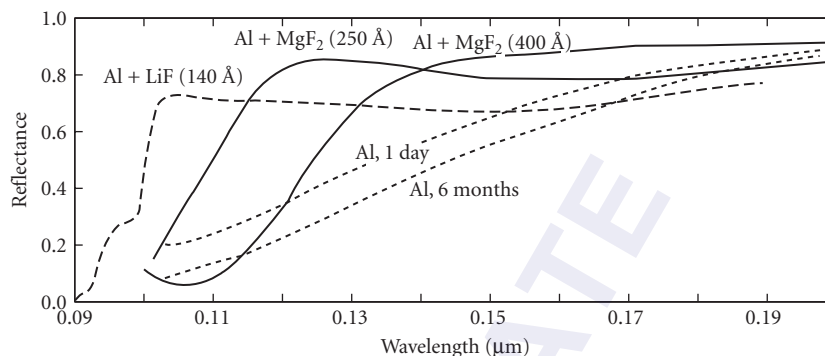


FIGURE 135 Measured spectral-reflectance curves of unprotected aluminum and aluminum overcoated with MgF_2 and LiF films of indicated thicknesses. (*Al + LiF coating after Cox et al.⁵⁴⁴ and all other curves after Canfield et al.⁵⁴¹*)

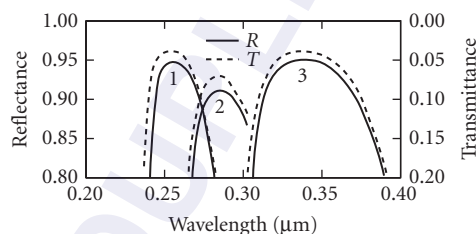


FIGURE 136 Enhanced ultraviolet reflectances of semitransparent aluminum films obtained through the addition of quarter-wave stacks. Curves 1 and 2: eleven layers of PbF_2 and MgF_2 (after *Leš et al.⁵⁵³*); curve 3: nine layers of Sb_2O_3 and MgF_2 (after *Leš and Leš⁵⁵⁴*).

The reflectance of silver, although very high in the visible, falls off rapidly in the near-ultraviolet. Attempts to enhance the reflectance in that part of the spectrum and, at the same time, to protect the silver from tarnish, have been successful (Fig. 134a).⁵³⁷

A different kind of reflection enhancement has been reported for the extreme ultraviolet. By depositing semi-transparent platinum films onto different substrates the opaque-film reflectances of 19.3 and 12.8 percent at 0.0584 and 0.0736 μm were increased by up to 2.8 and 3.8 percent, respectively.⁵⁵⁵ For space applications suitably thick aluminum films on iridium are expected to yield reflectances as high as 40 and 52 percent at the same wavelengths (Fig. 137). Other proposed reflectance-increasing combinations can be found in Madden et al.⁵⁵⁶

Selective Metal-Dielectric Reflectors Several types of coatings that reflect highly in one spectral region, but not in another have been developed in the past for different applications. Hadley and Dennison presented the theory and experimental results of reflection interference filters for the isolation of narrow spectral regions (Fig. 138a).^{558,559} Very narrow reflection filters have been described by Zheng (Fig. 138b).⁵⁶⁰ High-infrared and low-visible reflectance coatings (Fig. 66) are used to control the temperature of satellites.¹³⁸ These coatings could also be used to remove visible light from infrared optical systems. Several reflectors designed to reduce stray visible light in ultraviolet systems are shown in Fig. 139.

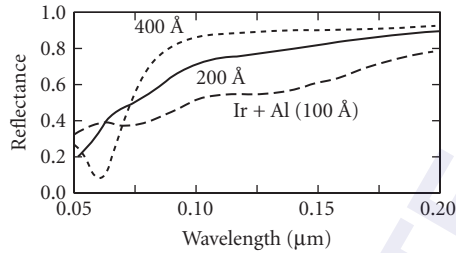


FIGURE 137 Calculated spectral reflectance curves of iridium overcoated with different thicknesses of unoxidized films of aluminum. (After Hass and Hunter.⁵⁵⁷)

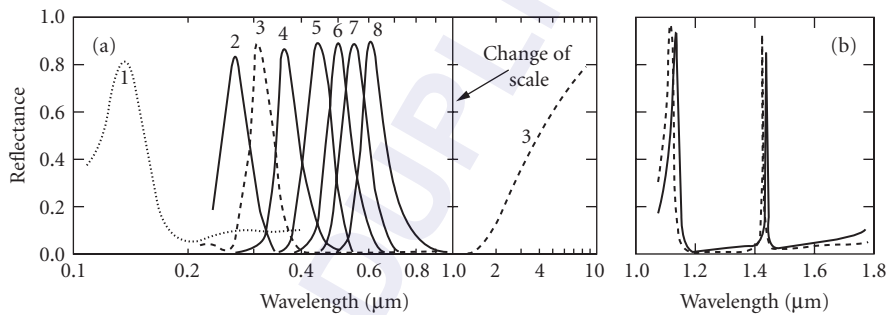


FIGURE 138 Measured performance of selective reflection filters. (a) Reflecting filters for the ultraviolet and visible spectral regions (curve 1 after Stelmack⁵⁶¹). Curves 2 to 8 represent the performance of filters that consist of three half-wave cryolite or MgF_2 spacer layers separated by Inconel films of suitable transmission deposited on an opaque aluminum film (after Turner and Hopkinson⁵⁶²). (b) Calculated and experimental performance of a very narrow band near-infrared reflecting filter (after Gamble⁵⁶³).

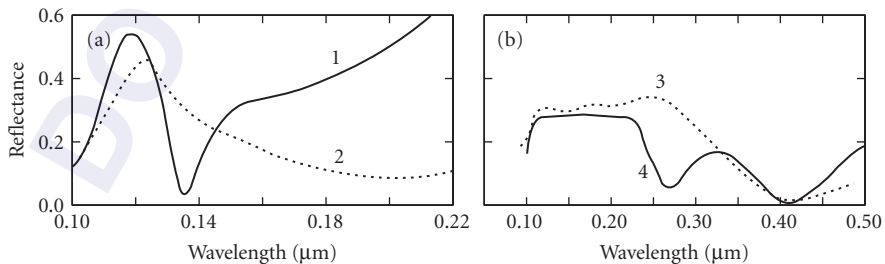


FIGURE 139 Measured performance of four multilayer selective ultraviolet reflectors for the control of stray radiation in the intermediate ultraviolet and visible spectral regions (a) and (b). (Curve 1 after Hunter;⁵²⁸ curve 2 after Berning et al.,⁵⁶⁴ and curves 3 and 4 after Hass and Tousey.⁵⁴²)

Reflection at Angles Close to Grazing Incidence

At wavelengths much shorter than the ultraviolet, all materials have refractive indices that are close to unity and extinction coefficients that are rather small. It follows from Eq. (80) that normal incidence reflectances in that part of the spectrum are small. However, for angles of incidence greater than the critical angle θ_c

$$\theta_c = \cos^{-1}[\sqrt{2(1-n_s)}] \quad (81)$$

total external reflection occurs resulting in high reflectances. Measured values of oblique angle reflection coefficients in the 0.0023 to 0.019 μm spectral region for a number of materials are given by Lukirskii et al.^{565,566} and also on Henke's website.²⁵⁴ Typical spectral reflectance curves of some elements are shown in Fig. 140a and b. An Al reflectance of 0.987 for an angle of incidence of 80° at 0.0584 μm has been reported by Newnam.⁵⁶⁷ It has also been shown that two or three thin films on a suitable substrate can have a better performance than a single metal layer (Fig. 140c).⁵⁶⁸ Higher reflectances can be achieved with periodic and nonperiodic multilayers designed for use at grazing incidence. For use in near-grazing incidence telescope optics for the soft x-ray region it may also be a requirement that, for a given wavelength, the reflectance stays constant over a small range of angles. Nonperiodic systems with such properties are sometimes called "supermirrors" and they can be composed of hundreds of layers. The performances of several such coatings are shown in Fig. 141.

Multiple-Reflection Filters

Metal and Metal-Dielectric Multiple-Reflection Filters Metals such as silver, copper, gold, and metal-dielectric coatings of the type shown in Fig. 64a used in a multiple-reflection arrangement should make low-wavelength cutoff filters with excellent rejection, sharp transition, and a long, unattenuated pass region far superior to those available with transmission filters.

Multiple-Reflection Filters Made of Thin-Film Interference Coatings Interference coatings for use in a multiple-reflection filter need not be deposited onto substrates that transmit well in the spectral region of interest, but they should be used at small angles of incidence if disturbing effects due to polarization are not to occur.⁵⁷⁸ The following are examples of some of the difficult filtering problems that can be easily solved with multiple reflection filters composed of interference coatings, providing that there is space to use a multiple reflection arrangement.

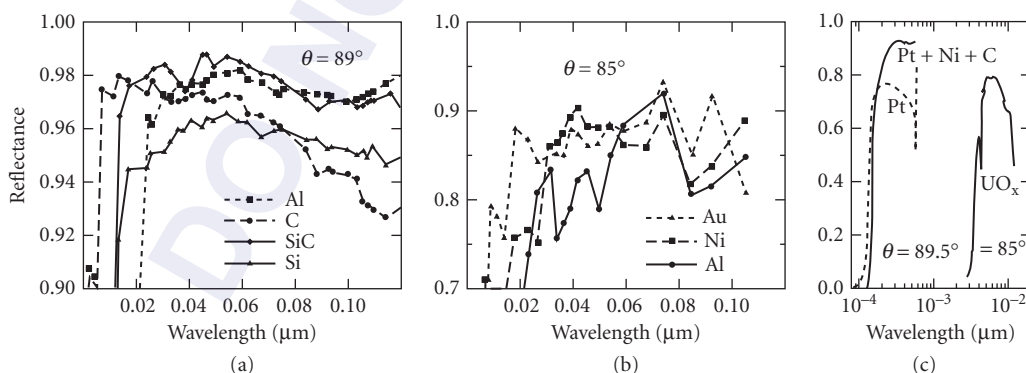


FIGURE 140 Measured XUV and x-ray reflectances of (a) Al, Si, C, and CVD SiC films (after Windt et al.⁵⁶⁹); (b) Ni, Au, and Al films (after Malina and Cash⁵⁷⁰); and (c) uranium oxide (after Sandberg et al.⁵⁷¹) and a single Pt layer and a triple layer of Pt, Ni, and C (after Tamura et al.⁵⁶⁸). The angles of incidence are indicated in the diagrams.

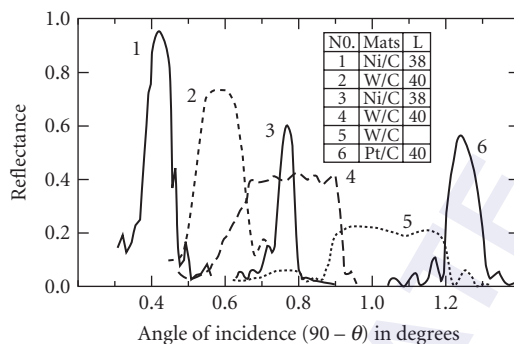


FIGURE 141 Measured reflectances of some mirrors operating at near-grazing angles of incidence for 8 keV ($\lambda = 0.0001543 \mu\text{m}$) radiation. The total number of layers and the materials used in their construction are indicated in the table. (Curve 1 after Spiga;⁵⁷² curve 2 after Protopopov;⁵⁷³ curve 3 after Citterio;⁵⁷⁴ curve 4 after Protopopov;⁵⁷⁵ curve 5 after Wang;⁵⁷⁶ and curve 6 after Yamashita.⁵⁷⁷)

It is difficult to provide adequate blocking with transmission filters for narrow-bandpass filters of the type shown in Figs. 97 to 102 without considerably reducing the peak transmittance. This is done readily with a multiple-reflection filter composed of quarter-wave stacks of the same materials used for the construction of the bandpass filter and centered at the same wavelength (Fig. 142).

It should be possible to construct highly efficient short-pass filters with a very long and low rejection region by using broadband reflectors consisting of several contiguous stacks (Sec. 7.7, subsection "Rejection Filters").

The use of a narrowband transmission filter in a multiple reflection arrangement of the type shown in Fig. 3 results in a high attenuation narrowband rejection filter surrounded by regions of high transmission (Fig. 143). However, such devices must be used with well-collimated light.

The transmittance curves of three multiple-reflection bandpass filters are shown in Fig. 144a to c. By using a number of multiple reflection filters with sharp features it is possible to separate signals

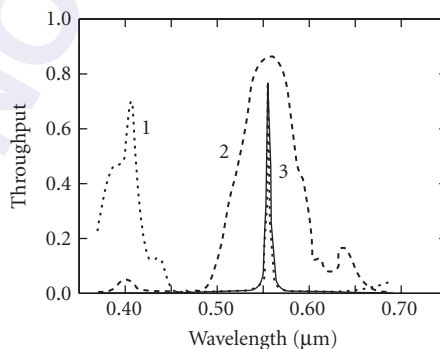


FIGURE 142 Blocking an all-dielectric narrowband interference filter with a multiple-reflection filter. Curve 1: transmittance of interference filter alone; curve 2: spectral reflectance of a quarter-wave stack after fourfold reflection; curve 3: transmission of blocked filter. (After Cohendet and Saudreau.³⁹⁸)

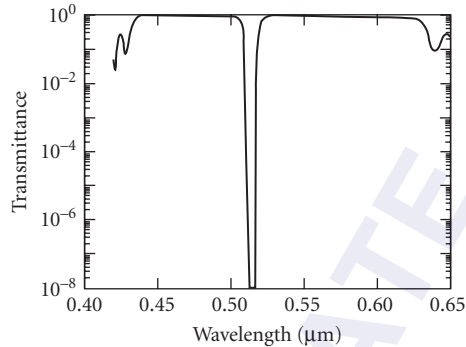


FIGURE 143 Measured transmission after four reflections from identical narrowband interference filters. (After Omega Optical Inc.²³⁹)

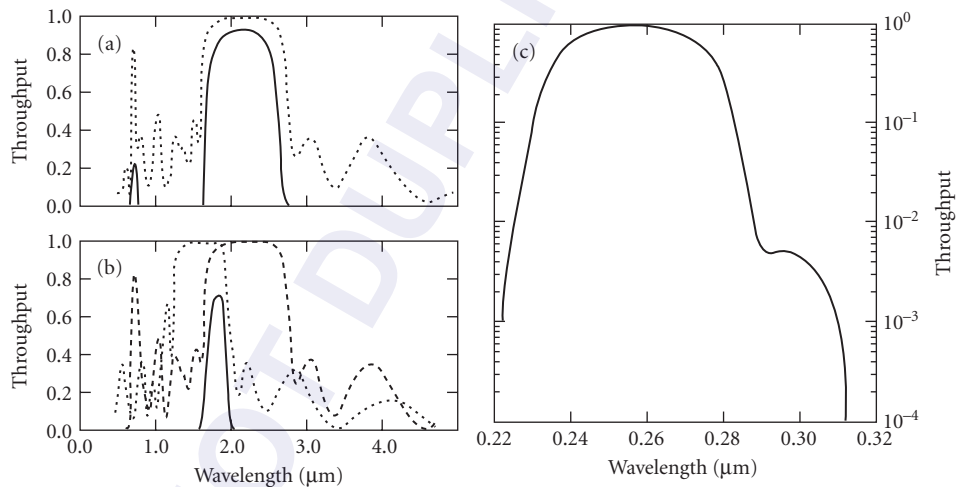


FIGURE 144 Measured spectral characteristics of broad- and narrowband multiple-reflection filters for the infrared and ultraviolet spectral regions. (a) Eight reflections (solid curves) from identical quarter-wave stacks (dotted); (b) six reflections (solid) from each of two quarter-wave stacks (dotted, dashed) tuned to different wavelengths (after Valejev⁵⁷⁹); and (c) commercial multiple-reflection interference filter (Schott & Gen.⁵⁸⁰).

transmitted by radiation of different, closely spaced wavelengths with a very low cross-talk and small insertion loss.⁵¹⁴

Additional information on reflection coatings and filters will be found in the reviews by Hass et al. and by Lynch.^{581,582}

7.17 SPECIAL PURPOSE COATINGS

Space considerations limit the detailed description of coatings and filters for specific applications. Absorbing multilayer coatings on glass for enhancing the visual appearance, thermal and illumination control, and as “one-way-mirrors” find applications in architecture and in the automotive

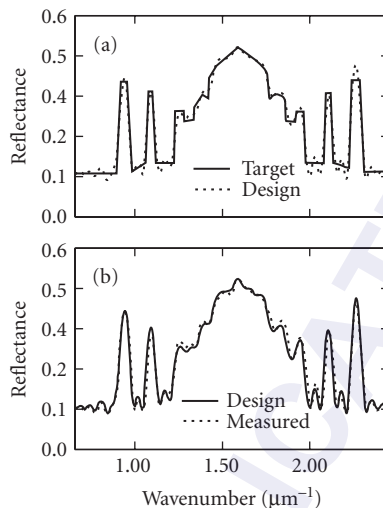


FIGURE 145 Calculated and experimentally measured reflectance of a filter that approximates the silhouette of the Taj Mahal. (After Sullivan and Dobrowolski.⁵⁸³)

industry. They are also used in solar energy conversion and have been proposed for radiative cooling. Thin-film coatings are used in optical recording media and in optical multiplexers/demultiplexers. Bistable Fabry-Perot structures are proposed for use as light switches in optical computers. Special filters and coatings are used in colorimetry, radiometry, detectors, and in high contrast display devices. Consumer-oriented products include various kinds of decorative coatings, as well as coatings for the protection of documents and products from counterfeiting. There is no doubt that in the future even more applications will appear for optical multilayer coatings. Some of these will require very complex spectral characteristics. At the time of writing of the second edition of this *Handbook*, methods for the design and manufacture of such coatings have just been demonstrated (Fig. 145). Now much more complicated correction and gain flattening filters (Fig. 115a) are important components for the telecommunications industry and they are routinely produced by many companies and laboratories.

The author would like to acknowledge the help and encouragement of his colleagues, Brian T. Sullivan, Li Li, Claude Montcalm, Pierre Verly, Daniel Poitras, Jeffrey Wong, and Allan Waldorf.

7. 18 REFERENCES

1. J. A. Dobrowolski, "Coatings and Filters," in *Handbook of Optics* (eds. W. G. Driscoll and W. Vaughan), McGraw-Hill, New York, 1978, pp. 8.1–8.124.
2. H. A. Macleod, *Thin Film Optical Filters*, Institute of Physics Publishing, Bristol, 2001.
3. J. D. Rancourt, *Optical Thin Films Users' Handbook*, Macmillan, New York, 1987.
4. P. Baumeister, *Optical Coating Technology*, SPIE—The International Society for Optical Engineering, Bellingham, Washington, 2004.
5. L. B. Tuckerman, "On the Intensity of the Light Reflected from or Transmitted through a Pile of Plates," *Journal of the Optical Society of America* **37**:818–825 (1947).
6. A. M. Dioffo, "Treatment of General Case of n Dielectric Films with Different Optical Properties," *Revue d'Optique, Theorique et Instrumentale* **47**:117–129 (1968).

7. A. Dresler, "Über eine neuartige Filterkombination zur genauen Angleichung der spektralen Empfindlichkeit von Photozellen an die Augenempfindlichkeitskurve," *Das Licht* 3:41–43 (1933).
8. M. R. Nagel, "A Mosaic Filter Daylight Source for Aerophotographic Sensitometry in the Visible and Infrared Region," *Journal of the Optical Society of America* 44:621–624 (1954).
9. J. A. Dobrowolski, F. C. Ho, A. Belkind, and V. A. Koss, "Merit Functions for More Effective Thin Film Calculations," *Applied Optics* 28:2824–2831 (1989).
10. Z. Knittl, *Optics of Thin Films*, Wiley & Sons, London, 1976.
11. S. A. Furman and A. V. Tikhonravov, *Optics of Multilayer Systems*, Editions Frontieres, Gif-sur-Yvette, 1992.
12. P. H. Berning, "Theory and Calculations of Optical Thin Films," in *Physics of Thin Films* 1, (ed. G. Hass), Academic Press, New York, 1963, pp. 69–121.
13. A. Thelen, *Design of Optical Interference Coatings*, McGraw-Hill Book Company, New York, 1988.
14. P. G. Kard, *Analysis and Synthesis of Multilayer Interference Coatings*, (in Russian) Valrus, Tallin, 1971, pp. 1–236.
15. E. Delano and R. J. Pegis, "Methods of Synthesis for Dielectric Multilayer Filters," in *Progress in Optics* 7, (ed. E. Wolf), North-Holland Publishing Company, Amsterdam, 1969, pp. 68–137.
16. J. A. Dobrowolski and R. A. Kemp, "Refinement of Optical Multilayer Systems with Different Optimization Procedures," *Applied Optics* 29:2876–2893 (1990).
17. L. Li and J. A. Dobrowolski, "Computation Speeds of Different Optical Thin Film Synthesis Methods," *Applied Optics* 31:3790–3799 (1992).
18. H. A. Macleod, "Thin-Film Optical Coatings," in *Applied Optics and Optical Engineering* 10, (eds. R. Kingslake, R. R. Shannon, and J. C. Wyant), Academic Press, 1987, pp. 1–69.
19. M. Born and E. Wolf, *Principles of Optics*, Pergamon Press, New York, 1970.
20. A. Herpin, "Calculation of the Reflecting Power of Any Stratified System," *Comptes Rendus* 225:182–183 (1947).
21. W. Weinstein, "Computations in Thin Film Optics," in *Vacuum* IV, E. T., Heron & Co. Ltd., London, 1954, pp. 3–19.
22. S. H. C. Piotrowski-McCall, J. A. Dobrowolski, and G. G. Shepherd, "Phase Shifting Thin Film Multilayers for Michelson Interferometers," *Applied Optics* 28:2854–2859 (1989).
23. J. A. Dobrowolski and D. Lowe, "Optical Thin Film Synthesis Program Based on the Use of Fourier Transforms," *Applied Optics* 17:3039–3050 (1978).
24. P. Baumeister, "The Transmission and Degree of Polarization of Quarter-Wave Stacks at Non-normal Incidence," *Optica Acta* 8:105–119 (1961).
25. V. R. Costich, "Reduction of Polarization Effects in Interference Coatings," *Applied Optics* 9:866–870 (1970).
26. K. Rabinovitch and A. Pagis, "Polarization Effects in Multilayer dielectric Thin Films," *Optica Acta* 21:963–980 (1974).
27. Z. Knittl, "Control of Polarization Effects by Internal Antireflection," *Applied Optics* 20:105–110 (1981).
28. Liberty Mirror—a Division of Libbey-Owens Ford Glass Co., "Coatings," 851 Third Avenue, Brackensridge, PA 15014, USA, 1966.
29. O. Arnon and P. Baumeister, "Electric Field Distribution and the Reduction of Laser Damage in the Multilayers," *Applied Optics* 19:1853–1855 (1980).
30. J. L. Vossen and W. Kern, *Thin Film Processes*, Academic Press, New York, 1978.
31. J. L. Vossen and W. Kern, *Thin Film Processes II*, Academic Press, Boston, 1991.
32. M. R. Jacobson (ed.), *Selected Papers on Deposition of Optical Coatings*, SPIE Milestone Series vol. MS 63, SPIE Optical Engineering Press, Bellingham, Washington, 1990.
33. J. M. Bennett, E. Pelletier, G. Albrand, J. P. Borgogno, B. Lazarides, C. K. Carniglia, T. H. Allen, T. Tuttle-Hart, K. H. Guenther, and A. Saxer, "Comparison of the Properties of Titanium Dioxide Films Prepared by Various Techniques," *Applied Optics* 28:3303–3317 (1989).
34. V. R. Costich, "Multilayer Dielectric Coatings," in *Handbook of Laser Science and Technology. Optical Materials: Part 3 V*, (ed. M. J. Weber), CRC Press, Inc., Boca Raton, Florida, 1987, pp. 389–430.
35. P. J. Martin, H. A. Maclean, R. P. Netterfield, C. G. Pacey, and W. G. Sainty, "Ion-Assisted Deposition of Thin Films," *Applied Optics* 22:178–184 (1983).

36. J. K. Hirvonen, "Ion Beam Assisted Thin Film Deposition," *Material Science Reports* **6**:215–274 (1991).
37. F. A. Smidt, "Ion-Beam-Assisted Deposition Provides Control over Thin Film Properties," NRL Publication 215-4670, issued May 1992, 1992.
38. H. Pulker, J. Edlinger, and M. Buehler, "Ion Plating Optical Films," in *Proceedings, 6th International Conference on Ion and Plasma Assisted Techniques*, CEP Consultants, Brighton, England, 1987, pp. 371–381.
39. K. H. Guenther, B. Loo, D. Burns, J. Edgell, D. Windham, and K. -H. Muller, "Microstructure Analysis of Thin Films Deposited by Reactive Evaporation and by Reactive Ion Plating," *Journal of Vacuum Science and Technology* **A7**:1435–1445 (1989).
40. A. J. Waldorf, J. A. Dobrowolski, B. T. Sullivan, and L. M. Plante, "Optical Coatings by Reactive Ion Plating," *Applied Optics* **32**:5583–5593 (1993).
41. A. Zöller, B. Beisswenger, R. Götzelmann, and K. Matl, "Plasma Ion-Assisted Deposition: A Novel Technique for the Production of Optical Coatings," *Proc. SPIE* **2253**:394 (1994).
42. A. Zöller, R. Götzelmann, K. Matl, and D. Cushing, "Temperature-Stable Bandpass Filters Deposited with Plasma Ion-Assisted Deposition," *Applied Optics* **35**:5609–5612 (1996).
43. U. Schulz, U. B. Schallenberg, and N. Kaiser, "Antireflection Coating Design for Plastic Optics," *Applied Optics* **41**:3107–3110 (2002).
44. W. D. Westwood, *Sputter Deposition 2*, American Vacuum Society, 2003.
45. L. Martinu and D. Poitras, "Plasma Deposition of Optical Films and Coatings: A Review," *Journal Vacuum Science Technology A* **18**:2619–2645 (2000).
46. S. Sneek, "Atomic Layer Deposition in Mass Production of Optical Coatings," in *Proceedings, 51 Annual Technical Conference*, Society of Vacuum Coaters, Chicago, IL, April 19–24, 2008, pp. 413–416.
47. B. T. Sullivan and J. A. Dobrowolski, "Deposition Error Compensation for Optical Multilayer Coatings: I. Theoretical Description," *Applied Optics* **31**:3821–3835 (1992).
48. D. Poitras, J. A. Dobrowolski, T. Cassidy, and S. Moisa, "Ion-Beam Etching for the Precise Manufacture of Optical Coatings," *Applied Optics* **42**:4037–4044 (2003).
49. P. A. Temple, "Thin-Film Absorptance Measurements Using Laser Calorimetry," in *Handbook of Optical Constants of Solids* (eds. E. D. Palik), Academic Press, Orlando, 1985, pp. 135–153.
50. D. Z. Anderson, J. C. Friesch, and C. S. Masser, "Mirror Reflectometer Based on Optical Cavity Decay Time," *Applied Optics* **23**:1238–1245 (1984).
51. C. Amra, P. Roche, and E. Pelletier, "Interface Roughness Cross-Correlation Laws Deduced from Scattering Diagram Measurements on Optical Multilayers: Effect of the Material Grain Size," *Journal of the Optical Society of America B* **4**:1087–1093 (1987).
52. J. M. Bennett and L. Mattsson, *Introduction to Surface Roughness and Scattering*, Optical Society of America, Washington, DC, 1989.
53. J. M. Bennett, "Recent Developments in Surface Roughness Characterization," *Measurement Science Technology* **3**:1119–1127 (1992).
54. A. Duparré, "Scattering from Surfaces and Thin Films," in *Encyclopedia of Modern Optics* 4, (eds. R. D. Guenther, D. G. Steel and L. Bayvel), Elsevier Academic Press, Amsterdam, 2004, pp. 314–321.
55. G. Hass and W. R. Hunter, "Laboratory Experiments to Study Surface Contamination and Degradation of Optical Coatings and Materials in Simulated Space Environments," *Applied Optics* **9**:2101–2110 (1970).
56. W. R. Hunter, "Optical Contamination: Its Prevention in the XUV Spectrographs Flown by the U. S. Naval Research Laboratory in the Apollo Telescope Mount," *Applied Optics* **16**:909–916 (1977).
57. L. Ward, *The Optical Constants of Bulk Materials and Films*, Adam Hilger, Bristol, 1988.
58. E. D. Palik, *Handbook of Optical Constants of Solids*, Academic Press Inc., Orlando, 1985.
59. E. D. Palik, *Handbook of Optical Constants of Solids II*, Academic Press Inc., Boston, 1991.
60. D. P. Arndt, R. M. A. Azzam, J. M. Bennett, J. P. Borgogno, C. K. Carniglia, W. E. Case, J. A. Dobrowolski, "Multiple Determination of the Optical Constants of Thin-Film Coating Materials," *Applied Optics* **23**:3571–3596 (1984).
61. Boulder Symposium, "Laser-Induced Damage in Optical Materials: Collected Papers, 1969–1998 (CD-ROM)," in *Proceedings, Boulder Laser Damage Symposium*, SPIE, 1999.

62. D. Ristau, "Laser Damage in Thin Film Coatings," in *Encyclopedia of Modern Optics* 3, (eds. R. D. Guenther, D. G. Steel, and L. Bayvel), Elsevier Academic Press, Amsterdam, 2004, pp. 339–349.
63. S. P. Baker and W. D. Nix, "Mechanical Properties of Thin Films on Substrates," *Proceedings of the Society of Photo-Optical Instrumentation Engineers* **1323**:263–276 (1990).
64. J. E. Klemberg-Sapieha, J. Oberste-Berghaus, L. Martinu, R. Blacker, I. Stevenson, G. Sadkhin, D. Morton, et al. "Mechanical Characteristics of Optical Coatings Prepared by Various Techniques: A Comparative Study," *Applied Optics* **43**:2670–2679 (2004).
65. A. W. Czanderna, *Methods of Surface Analysis* 1, Elsevier Scientific Publishing Company, Amsterdam, 1975.
66. J. Bartella, P. H. Berning, B. Bovard, C. K. Carniglia, E. Casparis, V. R. Costich, J. A. Dobrowolski, et al., "Multiple Analysis of an Unknown Optical Multilayer Coating," *Applied Optics* **24**:2625–2646 (1985).
67. A. Mussett and A. Thelen, "Multilayer Antireflection Coatings," in *Progress in Optics* 8, (ed. E. Wolf), 1970, pp. 203–237.
68. W. A. Faber, P. W. Kruse and W. D. Saur, "Improvement in Infrared Detector Performance through Use of Antireflection Film," *Journal of the Optical Society of America* **51**:115 (1961).
69. I. V. Grebenshchikov, *Prosvetlenie Optiki (Antireflection Coating of Optical Surfaces)* State Publishers of Technical and Theoretical Literature, Moscow-Leningrad, (1946).
70. T. Sawaki, "Studies on Anti-Reflection Films," Research Report No. 315, Osaka Industrial Research Institute, issued September, 1960.
71. J. T. Cox and G. Hass, "Antireflection Coatings for Optical and Infrared Optical Materials," in *Physics of Thin Films* 2, (eds. G. Hass and R. E. Thun), Academic Press, New York, 1964, pp. 239–304.
72. J. G. Wilder, "Porous Silica AR Coating for Use at 248nm or 266nm," *Applied Optics* **23**:1448–1449 (1984).
73. F. O'Neill, I. N. Ross, D. Evans, J. U. D. Langridge, B. S. Bilan and S. Bond, "Colloidal Silica Coatings for KrF and Nd: Glass Laser Applications," *Applied Optics* **26**:828–832 (1987).
74. I. M. Thomas, "Porous Fluoride Antireflective Coatings," *Applied Optics* **27**:3356–3358 (1988).
75. I. M. Thomas, "Method for the Preparation of Porous Silica Antireflection Coatings Varying in Refractive Index from 1.22 to 1.44," *Applied Optics* **31**:6145–6149 (1992).
76. Balzers Aktiengesellschaft fur Hochvakuumtechnik und Dünne Schichten, "Kurvenblätter und ihre Bezeichnungen," FL-9496 Balzers, 1962.
77. M. E. Motamedi, W. H. Southwell, and W. J. Gunning, "Antireflection Surfaces in Silicon Using Binary Optics Technology," *Applied Optics* **31**:4371–4376 (1993).
78. P. Baumeister, "Antireflection coatings with Chebyshev or Butterworth Response: Design," *Applied Optics* **25**:4568–4570 (1986).
79. J. Krepelka, "Maximally Flat Antireflection Coatings," *Jemna Mechanika a Optika* 53–56 (1992).
80. P. Kard, "Optical Theory of Anti-Reflection Coatings," *Loodus ja matemaatika* 1 67–85 (1959).
81. A. Thetford, "A Method of Designing Three-Layer Anti-Reflection Coatings," *Optica Acta* **16**:37–43 (1969).
82. J. H. Apfel and R. M. Gelber, "Filter with Neutral Transmitting Multilayer Coating Having Asymmetric Reflectance," U. S. patent 3649359, issued 25 July, 1972.
83. A. F. Turner, Bausch & Lomb, Private communication, 1970.
84. Optical Coating Laboratory Inc., "Multilayer Antireflection Coatings," 2789 North Point Parkway, Santa Rosa, CA 9 5407-7397, USA, 1965.
85. Thin Film Lab, Product information, 501B Basin Rd., West Hurley, NY 12491, USA, 1992.
86. TechOptics Ltd., "Laser Optics and Instrumentation," Second Avenue, Onchan, Isle of Man, British Isles, 1992.
87. Reynard Corporation, "Optical Components; Thin Film Coatings; Optical Instruments," 1020 Calle Sombra, San Clemente, CA 92673-6227, USA, 1992.
88. Spindler & Hoyer Inc., "Precision Optics," 459 Fortune Boulevard, Milford, MA 01757-1757, USA, 1990.
89. S. A. Furman, "Broad-Band Antireflection Coatings," *Soviet Journal of Optical Technology* **33**:559–564 (1966).
90. E. G. Stolov, "New Constructions of Interference Optical Antireflection Coatings," *Soviet Journal of Optical Technology* **58**:175–178 (1991).

91. R. R. Willey, "Broadband Antireflection Coating Design Performance Estimation," in *Proceedings, 34th Annual Technical Conference of the Society of Vacuum Coaters*, Society of Vacuum Coaters, March 17–22, 1991, Philadelphia, pp. 205–209.
92. A. V. Tikhonravov and J. A. Dobrowolski, "A New, Quasi-Optimal Synthesis Method for Antireflection Coatings," *Applied Optics* **32**:4265–4275 (1993).
93. P. T. Scharf, "Transmission Colour in Camera Lenses," *Journal of the SMPTE* **59**:191–194 (1952).
94. A. E. Murray, "Effect of Antireflection Films on Colour in Optical Instruments," *Journal of the Optical Society of America* **46**:790–796 (1956).
95. H. Kubota, "Interference Color," in *Progress in Optics* 1, Elsevier, Amsterdam, 1961, pp. 211–251.
96. V. V. Obodovskiy, V. N. Rozhdstvenskiy, and I. Chernyy, "Amber Coating and Colour Photographic Objectives," *Soviet Journal of Optical Technology* **34**:324–330 (1967).
97. J. A. Dobrowolski and W. Mandler, "Color Correcting Coatings for Photographic Objectives," *Applied Optics* **18**:1879–1880 (1979).
98. T. I. Oh, "Broadband AR Coatings on Germanium Substrates Using Ion-Assisted Deposition," *Applied Optics* **27**:4255–4259 (1988).
99. S. P. Fisher, J. F. Leonard, I. T. Muirhead, G. Buller, and P. Meredith, "The Fabrication of Optical Devices by Molecular Beam Deposition Technology," in *Proceedings, Optical Interference Coatings, Technical Digest Series*, Tucson, AZ, Optical Society of America, Washington, D.C., 1992, pp. 131–133.
100. R. Jacobsson, "Light Reflection from Films of Continuously Varying Refractive Index," in *Progress in Optics* 5, (eds. E. Wolf), North-Holland, Amsterdam, 1966, pp. 247–286.
101. E. Spiller, I. Haller, R. Feder, J. E. E. Baglin, and W. N. Hammer, "Graded-Index AR Surfaces Produced by Ion Implantation on Plastic Materials," *Applied Optics* **19**:3022–3026 (1980).
102. W. H. Southwell, "Gradient-Index Antireflection Coatings," *Optics Letters* **8**:584–586 (1983).
103. B. Sheldon, J. S. Haggerty, and A. G. Emslie, "Exact Computation of the Reflectance of a Surface Layer of Arbitrary Refractive Index Profile and an Approximate Solution of the Inverse Problem," *Journal of the Optical Society of America* **72**:1049–1055 (1982).
104. M. J. Minot, "The Angular Reflectance of Single-Layer Gradient Refractive-Index Films," *Journal of the Optical Society of America* **67**:046–1050 (1977).
105. L. M. Cook and K. -H. Mader, "Integral Antireflective Surfaces on Optical Glass," *Optical Engineering* **21**:SR-008–SR-012 (1982).
106. W. H. Lowdermilk, "Graded-Index Surfaces and Films," in *Handbook of Laser Science and Technology. Optical Materials: Part 3 V*, (ed. M. J. Weber), CRC Press, Inc., Boca Raton, Florida, 1987, pp. 431–458.
107. H. Anders and R. Eichinger, "Die Optische Wirkung und Praktische Bedeutung Inhomogener Schichten," *Applied Optics* **4**:899–905 (1965).
108. G. Lessman, "Optical Properties of Inhomogeneous Thin Films of Varying Modes of Gradation," *Journal of the Optical Society of America* **56**:554 (1966).
109. R. Bertram, M. F. Ouellette, and P. Y. Tse, "Inhomogeneous optical Coatings: An Experimental Study of a New Approach," *Applied Optics* **28**:2935–2939 (1989).
110. R. Jacobsson, "Calculation and Deposition of Inhomogeneous Thin Films," *Journal of the Optical Society of America* **56**:1435 (1966).
111. R. J. Scheuerman, Perkin-Elmer Corporation, Private communication, 1968.
112. A. Y. Kuznetsov and A. F. Pervyev, "Multifilm Anti-Reflection Coatings for Materials with High Refractive Index," *Soviet Journal of Optical Technology* **34**:593–595 (1967).
113. P. H. Berning, "Use of Equivalent Films in the Design of Infrared Multilayer Antireflection Coatings," *Journal of the Optical Society of America* **52**:431–436 (1962).
114. J. A. Dobrowolski and F. C. Ho, "High Performance Step-Down AR Coatings for High Refractive-Index IR Materials," *Applied Optics* **21**:288–292 (1982).
115. L. I. Epstein, "The Design of Optical Filters," *Journal of the Optical Society of America* **42**:806–810 (1952).
116. W. H. Southwell, "Coating Design Using Very Thin High- and Low-Index Layers," *Applied Optics* **24**:457–460 (1985).
117. H. M. Moulton, "Method of Forming a Reflection Reducing Coating," U. S. patent 2536764, issued 2 January, 1951.
118. S. M. Thomsen, "Skeletonizing Glass," U. S. patent 2490662, issued 6 December, 1949.

119. M. J. Minot, "Single Layer, Gradient Refractive Index Antireflection Films—Effective 0.35 to 2.5 μ ," *Journal of the Optical Society of America* **66**:515–519 (1976).
120. Y. Asahara and T. Izumitani, "The Properties of Gradient Index Antireflection Layer on the Phase Separable Glass," *Journal of Non-Crystalline Solids* **42**:269–280 (1980).
121. W. H. Lowdermilk and D. Milam, "Graded-Index Antireflection Surfaces for High-Power Laser Applications," *Applied Physics Letters* **36**:891–893 (1980).
122. S. Fujikawa, Y. Oguri, and E. Arai, "Antireflection Surface Modification of Plastic CR-39 by Means of ^{16}O - and ^{35}Cl -ion Bombardment," *Optik (Stuttgart)* **84**:1–5 (1990).
123. B. E. Yoldas and D. P. Partlow, "Wide Spectrum Antireflective Coating for Fused Silica and Other Glasses," *Applied Optics* **23**:1418–1424 (1984).
124. K. N. Maffitt, H. U. Bruechner, and D. R. Lowrey, "Polymeric Optical Element Having Antireflecting Surface," U. S. patent 4114983, issued 19 September, 1978.
125. S. J. Wilson and M. C. Hutley, "The Optical Properties of 'Moth Eye' Antireflection Surfaces," *Optica Acta* **29**:993–1009 (1982).
126. S. P. Mukherjee, "Gel-Derived Single-Layer Antireflection Films with a Refractive Index Gradient," *Thin Solid Films* **81**:L89–L90 (1981).
127. S. P. Mukherjee and W. H. Lowdermilk, "Gel-Derived Single Layer Antireflection Films," *Journal of Non-Crystalline Solids* **48**:177–184 (1982).
128. P. B. Clapham and M. C. Hutley, "Reduction of Lens Reflexion by the 'Moth Eye' Principle," *Nature* **244**:281–282 (1973).
129. W. H. Southwell, "Pyramid-Array Surface-Relief Structures Producing Antireflection Index Matching on Optical Surfaces," *Journal of the Optical Society of America* **A8**:549–553 (1991).
130. D. H. Raguin and G. M. Morris, "Antireflection Structured Surfaces for the Infrared Spectral Region," *Applied Optics* **32**:1154–1167 (1993).
131. S. R. Kennedy and M. J. Brett, "Porous Broadband Antireflection Coating by Glancing Angle Deposition," *Applied Optics* **42**:4573–4579 (2003).
132. V. D. Vvedenskii, A. A. Metel'nikov and S. A. Furman, "Achromatic Antireflection Coatings for Materials with Refractive Indices of 1.46–1.8," *Soviet Journal of Optical Technology* **47**:157–159 (1980).
133. J. A. Dobrowolski and B. T. Sullivan, "Universal Antireflection Coatings for the Visible Spectral Region," *Applied Optics* **35**:4993–4997 (1996).
134. L. Yeuh Yeong, "Universal Broad-Band Antireflection Coating Designs for Substrates in the Visible Spectral Region," *Japanese Journal of Applied Physics*, Part 1 Regular Papers, Short Notes & Review Papers **41**:6409–6410 (2002).
135. V. N. Smiley, "Conditions for Zero Reflectance of Thin Dielectric Films on Laser Materials," *Journal of the Optical Society of America* **58**:1469–1475 (1968).
136. E. A. Lupashko and I. N. Shkyarevskii, "Multilayer Dielectric Antireflection Coatings," *Optics and Spectroscopy (USSR)* **16**:279–281 (1964).
137. K. C. Park, "The Extreme Values of the Reflectivity and the Conditions for Zero Reflection from Thin Dielectric Films on Metal," *Applied Optics* **3**:877–881 (1964).
138. L. F. Drummeter and G. Hass, "Solar Absorptance and Thermal Emittance of Evaporated Coatings," in *Physics of Thin Films 2*, (eds. G. Hass and R. E. Thun), Academic Press, New York, 1964, pp. 305–361.
139. R. N. Schmidt and K. C. Park, "High-Temperature Space-Stable Selective Solar Absorber Coatings," *Applied Optics* **4**:917–925 (1965).
140. S. Katsube, Y. Katsube, K. Mitome, and S. Furuta, "Non-Metallic Absorbing Films and the Anti-Reflection Coating," *Japanese Journal of Applied Physics* **37**:225–230 (1968).
141. H. Anders, "Inhomogeneous, Absorbing Layers for Sunglasses," *Vakuum-Technik* **15**:123–126 (1966).
142. V. V. Veremey and T. A. Gorbunova, "Achromatic Coating of Glasses Covered by Conducting Films," *Soviet Journal of Optical Technology* **35**:519–521 (1968).
143. T. Turbadar, "Complete Absorption of Plane Polarized Light by Thin Metallic Films," *Optica Acta* **11**:207–210 (1964).
144. T. Turbadar, "Equi-Reflectance Contours of Double-Layer Anti-Reflection Coatings," *Optica Acta* **11**:159–205 (1964).

145. H. Pohlack, "Zum Problem der Reflexionsminderung Optischer Gläser bei nichtsenkrechtem Lichteinfall," in *Jenaer Jahrbuch 1952* (ed. P. Görlich), Fischer, Jena, 1952, pp. 103–118.
146. J. A. Dobrowolski and S. H. C. Piotrowski, "Refractive Index as a Variable in the Numerical Design of Optical Thin Film Systems," *Applied Optics* **21**:1502–1510 (1982).
147. J. A. Dobrowolski, "The Impact of Computers on the Design and Manufacture of Optical Multilayer Coatings During the Past 50 Years," in *Proceedings, 50th Annual Technical Conference*, Society of Vacuum Coaters, Louisville, KY, 2007, pp. 289–301.
148. J. A. Dobrowolski, Y. Guo, T. Tiwald, P. Ma, and D. Poitras, "Toward Perfect Antireflection Coatings. 3. Experimental Results Obtained with the Use of Reststrahlen Materials," *Applied Optics* **45**:1555–1562 (2006).
149. P. Ma, J. A. Dobrowolski, D. Poitras, T. Cassidy, and F. Lin, "Toward the Manufacture of 'Perfect' Antireflection Coatings," in *Proceedings, Proceedings of the 45 Annual Technical Conference*, Society of Vacuum Coaters, Lake Buena Vista, Florida, 2002, pp. 216–218.
150. U. Schulz and N. Kaiser, "Designing Optical Coatings by Using Low-Index Equivalent Layers and Low-Index Effective Media," in *Proceedings, Annual Technical Conference*, Society of Vacuum Coaters, Santa Clara, CA, 2009.
151. J. A. Dobrowolski, D. Poitras, M. Penghui, V. Himanshu, and M. Acree, "Toward Perfect Antireflection Coatings: Numerical Investigation," *Applied Optics* **41**:3075–3083 (2002).
152. D. Poitras and J. A. Dobrowolski, "Toward Perfect Antireflection Coatings. 2. Theory," *Applied Optics* **43**:1286–1295 (2004).
153. J. C. Monga, "Anti-Reflection Coatings for Grazing Incidence Angles," *Journal of Modern Optics* **36**:381–387 (1989).
154. J. C. Monga, "Double-Layer Broadband Antireflection Coatings for Grazing Incidence Angles," *Applied Optics* **31**:546–553 (1992).
155. S. Bruynooghe, S. Spinzig, M. Fliedner, and G. H. Hsu, "Characterization of the Optical Properties of Hydrophobic Coatings and Realization of High Performance AR-Coatings with Dust- and Water-Repellent Properties," *Proc. SPIE* **7101**:71010Q (2008).
156. P. Giacomo, "Les couches Réfléchissantes Multidiélectriques Appliquées a l'interféromètre de Fabry-Perot etude Théorique et Expérimentale Des Couches Réelles," *Revue d'Optique, Théorique et Instrumentale* **35**:318–354 (1956).
157. J. Arndt and P. Baumeister, "Reflectance and Phase Envelopes of an Iterated Multilayer," *Journal of the Optical Society of America* **56**:1760–1762 (1966).
158. J. S. Seeley, "Resolving Power of Multilayer Filters," *Journal of the Optical Society of America* **54**:342–346 (1964).
159. H. Böhme, "Dielektrische Mehrschichtsysteme ohne Dispersion des Phasensprungs," *Optik* **69**:1–7 (1984).
160. P. Kard, E. Nesmelov, and G. Konyukhov, "Theory of a Quarterwave Reflecting Filters," *Eesti NSV Tead. Akad. Toim., Fuus. - Mat.* **17**:314–323 (1968).
161. A. P. Ovcharenko and E. A. Lupashko, "Multilayer Dielectric Coatings of Unequal Thickness," *Optics and Spectroscopy (USSR)* **55**:316–318 (1983).
162. V. A. Smirnova and G. D. Pridatko, "Two-Component Interference Coatings. Optical Properties and Applications," *Optics and Spectroscopy (USSR)* **55**:442–445 (1983).
163. M. Zukic and D. G. Toerr, "Multiple Reflectors as Narrow-Band and Broadband Vacuum Ultraviolet Filters," *Applied Optics* **31**:1588–1596 (1992).
164. J. Shao and J. A. Dobrowolski, "Multilayer Interference Filters for the Far-Infrared and Submillimeter Regions," *Applied Optics* **32**:2361–2370 (1993).
165. T. W. Barbee, "Multilayers for X-Ray Optics," *Optical Engineering* **25**, 898–915 (1986).
166. E. Spiller, *Soft X-Ray Optics* SPIE Optical Engineering Press, 1994.
167. B. L. Henke, "X-Ray Interaction with Matter," http://henke.lbl.gov/optical_constants/ (2008).
168. A. Attwood, *Soft X-Rays and Extreme Ultraviolet Radiation*, Cambridge University Press, New York, 1999.
169. A. V. Vinogradov and B. Y. Zeldovich, "X-Ray and Far UV Multilayer Mirrors: Principles and Possibilities," *Applied Optics* **16**:89–93 (1977).

170. J. H. Apfel, "Multilayer Interference Coatings for the Ultraviolet," *Journal of the Optical Society of America* **59**:553 (1966).
171. K. Hefft, R. Kern, G. Nöldeke, and A. Steudel, "Über Fabry-Perot-Verspiegelungen aus dielektrischen Vielfachschichten für den Spektralbereich von 2350 bis 20000 Å," *Zeitschrift für Physik* **175**:391–404 (1963).
172. V. R. Costich, Spectra-Physics, Private communication, 1968.
173. A. F. Turner and P. W. Baumeister, "Multilayer Mirrors with High Reflectance over an Extended Spectral Region," *Applied Optics* **5**:69–76 (1966).
174. O. S. Heavens, "All-dielectric High-Reflecting Layers," *Journal of the Optical Society of America* **44**:371–373 (1954).
175. D. H. Rank and H. E. Bennett, "Problem of Phase Variation with Wavelength in Dielectric Films. Extension of Interferometric Standards into the Infrared," *Journal of the Optical Society of America* **45**:69–73 (1955).
176. P. Giacomo, "Propriétés Chromatiques Des Couches Réfléchissantes Multi-diélectrique," *Journal de Physique et le Radium* **19**:307–311 (1958).
177. G. Bouwhuis, "A Dispersion Phenomenon Observable on Dielectric Multilayer Mirrors," *Philips Research Reports* **17**:130–132 (1962).
178. J. V. Ramsay and P. E. Ciddor, "Apparent Shape of Broad Band, Multilayer Reflecting Surfaces," *Applied Optics* **6**:2003–2004 (1967).
179. G. Koppelman, "Zur Theorie der Wechschichten aus schwachabsorbierenden Substanzen und ihre Verwendung als Interferometerspiegel," *Annalen der Physik* **7**:387–396 (1960).
180. D. J. Hemingway and P. H. Lissberger, "Properties of Weakly Absorbing Multilayer Systems in Terms of the Concept of Potential Transmittance," *Optica Acta* **20**:85–96 (1973).
181. G. Honcia and K. Krebs, "Highly Reflecting Substances for Dielectric Mirror Systems," *Zeitschrift für Physik* **165**:202–212 (1961).
182. K. H. Behrndt and D. W. Doughty, "High-Reflectance Multilayer Dielectric Mirrors," *Journal of Vacuum Science and Technology* **4**:199–202 (1967).
183. D. Gloge, E. L. Chinnock, and H. E. Earl, "Scattering from Dielectric Mirrors," *Bell System Technical Journal* **48**:511–526 (1969).
184. R. Blazey, "Light Scattering by Laser Mirrors," *Applied Optics* **6**:831–836 (1967).
185. G. Rempe, R. J. Thompson, H. J. Kimble, and R. Lalezari, "Measurement of Ultralow Losses in an Optical Interferometer," *Optics Letters* **17**:363–365 (1993).
186. R. Lalezari, PMS Electro-Optics, 1855 South 57th Court, Boulder, CO80301, USA, Private communication, 1993.
187. G. M. Harry, H. Armandula, E. Black, D. R. M. Crooks, G. Cagnoli, J. Hough, P. Murray, et al., "Thermal Noise from Optical Coatings in Gravitational Wave Detectors," *Applied Optics* **45**:1569–1574 (2006).
188. B. Vidal and P. Vincent, "Metallic Multilayers for X-Rays Using Classical Thin-Film Theory," *Applied Optics* **23**:1794–1801 (1984).
189. D. G. Stearns, "X-Ray Scattering from Interfacial Roughness in Multilayer Structures," *Journal of Applied Physics* **71**:4286–4298 (1992).
190. D. G. Stearns, "The Scattering of X-Rays from Nonideal Multilayer Structures," *Journal of Applied Physics* **65**:491–506 (1989).
191. E. Spiller, "Multilayer Optics for X-Rays," in *Physics, Fabrication, and Applications of Multilayered Structures* (eds. P. D. a. C. Weisbuch), Plenum Publishing Corporation, 1988, pp. 271–309.
192. S. Braun, H. Mai, M. Moss, and R. Scholz, "Microstructure of Mo/Si Multilayers with Barrier Layers," *Proc. SPIE* **4782**:185–195 (2002).
193. L. Young, "Multilayer Interference Filters with Narrow Stop Bands," *Applied Optics* **6**:297–315 (1967).
194. Lord Rayleigh, "On the Remarkable Phenomenon of Crystalline Reflexion Described by Prof. Stokes," *Philosophical Magazine* **26**, **Fifth Series**, 256–265 (1888).
195. R. W. Wood, *Physical Optics* Macmillan, New York, 1940.
196. H. Schröder, "The Light Distribution Functions of Multiple-Layers and Their Applications," *Zeitschrift für Angewandte Physik* **2**:53–66 (1951).
197. J. Strong, "Iridescent KClO₃ Crystals and Infrared Reflection Filters," *Journal of the Optical Society of America* **51**:853–855 (1961).

198. T. J. Alfrey, E. F. Gurnee, and W. J. Schrenk, "Physical Optics of Iridescent Multilayered Plastic Films," *Polymer Engineering and Science* **9**:400–404 (1969).
199. T. J. Alfrey, "Multilayer Thermoplastic Sheets and Films," in *Proceedings, 19th Sagamore Army Materials Research Conference*, (eds. J. J. Burke and V. Weiss), Syracuse University Press, 1973, pp. 195–200.
200. O. S. Heavens, J. Ring, and S. D. Smith, "Interference Filters for the Infra-Red," *Spectrochimica Acta* **10**:179–194 (1957).
201. L. Edmonds, P. Baumeister, M. E. Krisl, and N. Boling, "Spectral Characteristics of a Narrowband Rejection Filter," *Applied Optics* **29**:3203–3204 (1990).
202. A. Badeen, M. Briere, P. Hook, C. Montcalm, R. Rinfret, J. Schneider, and B. T. Sullivan, "Advanced Coatings for Telecom and Spectroscopic Applications," *Proc. SPIE* **7101**:71010H1–H16 (2008).
203. K. D. Hendrix, C. A. Hulse, G. J. Ockenfuss, and R. B. Sarget, "Demonstration of Narrowband Notch and Multi-Notch Filters," *Proc. SPIE* **7067**:7067–7102 (2008).
204. M. Hercher, "Single-Mode Operation of a Q-Switched Ruby Laser," *Applied Physics Letters* **7**:39–41 (1965).
205. G. Magyar, "Simple Giant Pulse Ruby Laser of High Spectral Brightness," *Review of Scientific Instruments* **38**:517–519 (1967).
206. J. K. Watts, "Theory of Multiplate Resonant Reflectors," *Applied Optics* **7**:621–1623 (1968).
207. H. F. Mahlein and G. Schollmeier, "Analysis and Synthesis of Periodic Optical Resonant Reflectors," *Applied Optics* **8**:1197–1202 (1966).
208. Anon, "Quartz Etalon Outperforms Sapphire," *Microwaves* September 1969.
209. K. M. Yoo and R. R. Alfano, "Broad Bandwidth Mirror with Random Layer Thicknesses," *Applied Optics* **28**:2456–2458 (1989).
210. D. L. Perry, "Low-Loss Multilayer Dielectric Mirrors," *Applied Optics* **4**:987–991 (1965).
211. O. S. Heavens and H. M. Liddell, "Staggered Broad-Band Reflecting Multilayers," *Applied Optics* **5**:373–376 (1966).
212. Z. N. Elsner, "On the Calculation of Multilayer Interference Coatings with Given Spectral Characteristics," *Optics and Spectroscopy (USSR)* **17**:238–240 (1964).
213. F. A. Korolev, A. Y. Klementeva, T. F. Meshcheryakova, and I. A. Ramazina, "Wide-Band Reflectors with Multilayer Dielectric Coatings," *Optics and Spectroscopy (USSR)* **28**:416–419 (1970).
214. R. S. Sokolova, "Wide-Band Reflectors for Ultraviolet Radiation," *Soviet Journal of Optical Technology* **38**:295–297 (1971).
215. E. G. Stolov, "Constructions of Neutral Lightsplitters and Wide-Band Mirrors," *Soviet Journal of Optical Technology* **57**:632–634 (1990).
216. S. Penselin and A. Steudel, "Fabry-Perot-Interferometerspiegelungen aus dielektrischen Vielfachschechten," *Zeitschrift für Physik* **142**:21–41 (1955).
217. P. W. Baumeister and J. M. Stone, "Broad-Band Multilayer Film for Fabry-Perot Interferometers," *Journal of the Optical Society of America* **46**:228–229 (1956).
218. P. E. Ciddor, "Minimization of the Apparent Curvature of Multilayer Reflecting Surfaces," *Applied Optics* **7**:2328–2329 (1968).
219. A. F. Turner, "Design Principles for Interference Film Combinations," *Journal of the Optical Society of America* **44**:352 (1954).
220. K. V. Popov, J. A. Dobrowolski, A. V. Tikhonravov, and B. T. Sullivan, "Wide-Band High Reflection Multilayer Coatings at Oblique Angles of Incidence," *Applied Optics* **36**:2139–2151 (1997).
221. J. N. Winn, Y. Fink, S. Fan, and J. D. Joannopoulos, "Omnidirectional Reflection from a One-Dimensional Photonic Crystal," *Optics Letters* **23**:1573–1575 (1998).
222. Y. Fink, J. N. Winn, S. Fan, S. Chen, J. Michel, J. D. Joannopoulos, and E. L. Thomas, "A Dielectric Omnidirectional Reflector," *Science* **282**:1679–1682 (1998).
223. R. Gaughan, "New Coatings Break Reflectivity Barriers," *Photonics Spectra* **33**:29–30 (1999).
224. P. W. Baumeister and F. A. Jenkins, "Dispersion of the Phase Change for Dielectric Multilayers. Application to the Interference Filter," *Journal of the Optical Society of America* **47**:57–61 (1957).
225. P. E. Ciddor, "Phase-Dispersion in Interferometry. Interferometers with Solid Spacers and with Broadband-Reflecting Surfaces," *Optica Acta* **12**:177–183 (1965).
226. G. Tempea, V. Yakovlev, and F. Krausz, "Interference Coatings for Ultrafast Optics," in *Optical Interference Coatings* (eds. N. Kaiser and H. K. Pulker), Springer, Berlin, pp. 393–422, 2003.

227. M. Trubetskov, A. Tikhonravov, and V. Pervak, "Time-Domain Approach for Designing Dispersive Mirrors Based on the Needle Optimization Technique. Theory," *Optics Express* **16**:20637–20647 (2008).
228. N. Matuschek, F. X. Kärtner, and U. Keller, "Theory of Double-Chirped Mirrors," *IEEE Journal Selected Topics in Quantum Electronics* **4**:197–208 (1998).
229. V. Pervak, D. Grupe, M. Trubetskov, A. Tikhonravov, A. Apolonski, and F. Krausz, "1.5-Octave Chirped Dielectric Multilayers for Pulse Compression," *Proc. SPIE* **7101**:71016-1-7 (2008).
230. V. Pervak, I. Ahmad, J. Fulop, M. K. Trubetskov, and A. V. Tikhonravov, "Comparison of Dispersive Mirrors Based on the Time-Domain and Conventional Approaches, for Sub-5-fs Pulses," *Optics Express* **17**:2207–2217 (2009).
231. D. Yang, C. Lin, W. Chen, and G. Barbarossa, "Fiber Dispersion and Dispersion Slope Compensation in a 40-channel 10-Gb/s 3200-km Transmission Experiment Using Cascaded Single-Cavity Gire-Tournois Etalons," *IEEE Photonics Technology Letters* **16**:299–301 (2004).
232. Avanex Corp, "www.avanex.com," 2008.
233. V. Pervak, C. Teisset, A. Sugita, S. Naumov, F. Krausz, and A. Apolonski, "High-Dispersive Mirrors for Femtosecond Lasers," *Optics Express* **16**:10220–10230 (2008).
234. J. Kim, J. R. Birge, V. Sharma, J. G. Fujimoto, F. X. Kärtner, V. Scheuer, and G. Angelow, "Ultrabroadband Beam Splitter with Matched Group-Delay Dispersion," *Optics Letters* **30**:1569–1571 (2005).
235. E. Cojocar, "Second and Third Order Dispersion of Broadband Thin-Film Antireflection Coatings for Ultrafast Lasers," *Optica Applicata* **34**:25–29 (2004).
236. J. A. Dobrowolski, "Optical Interference Filters for the Adjustment of Spectral Response and Spectral Power Distributions," *Applied Optics* **9**:1396–1402 (1970).
237. Optical Coating Laboratory, "Catalog," 2789 Northpoint Parkway, Santa Rosa, CA 95407-7397, USA, 1964.
238. A. Thelen, "Design of Optical Minus Filters," *Journal of the Optical Society of America* **61**:365–369 (1971).
239. Omega Optical Inc., "Bandpass filters," 3 Grove Street, P. O. Box 573, Brattleboro, VT 05302-0573, USA, 1992.
240. R. Jacobsson, "Inhomogeneous and Coevaporated Homogeneous Films for Optical Applications," in *Physics of Thin Films* **8**, (eds. G. Hass, M. H. Francombe and R. W. Hoffman), Academic Press, New York, 1975, pp. 51–98.
241. J. A. Dobrowolski and P. G. Verly (eds.), *Inhomogeneous and Quasi-inhomogeneous Optical Coatings*, Proceedings of the Society of Photo-Optical Instrumentation Engineers vol. 2046, SPIE—The International Society for Optical Engineering, Bellingham, Washington, 1993.
242. W. J. Gunning, R. L. Hall, F. J. Woodberry, W. H. Southwell, and N. S. Gluck, "Codeposition of Continuous Composition Rugate Filters," *Applied Optics* **28**:2945–2948 (1989).
243. P. G. Verly, National Research Council of Canada, Private communication, 1993.
244. Physical Optics Corp., "Holographic Filters Aid Spectroscopy," *Photonics Spectra* **26**:113–114 (1992).
245. Kaiser Optical Systems Inc., "Holographic Notch Filters," 371 Parkland Plaza, P. O. Box 983, Ann Arbor, MI 48106, USA, 1993.
246. S. N. Vlasov and V. I. Talanov, "Selection of Axial Modes in Open Resonators," *Radio Engineering and Electronic Physics* **10**:469–470 (1965).
247. N. G. Vahitov, "Open Resonators with Mirrors Having Variable Reflection Coefficients," *Radio Engineering and Electronic Physics* **10**:1439–1446 (1965).
248. G. Duplain, P. G. Verly, J. A. Dobrowolski, A. Waldorf, and S. Bussière, "Graded Reflectance Mirrors for Beam Quality Control in Laser Resonators," *Applied Optics* **32**:1145–1153 (1993).
249. A. Piegari, "Graded Optical Coatings for Laser Applications," *Proceedings of the Society of Photo-Optical Instrumentation Engineers* **2461**:558–565 (1995).
250. A. Piegari, A. Tirabassi, and G. Emiliani, "Thin Films for Special Laser Mirrors with Radially Variable Reflectance: Production Techniques and Laser Testing," *Proceedings of the Society of Photo-Optical Instrumentation Engineers* **1125**:68–73 (1989).
251. V. G. Vereshchagin and A. D. Zamkovets, "Polymer-Crystalline Multilayer Systems for the Far IR Spectral Region," *Zhurnal Prikladnoi Spektroskopii* **47**:132–135 (1987).
252. V. R. Costich, "Study to Demonstrate a New Process to Produce Infrared Filters," NASA Ames NAS2-12639, issued 1 March 1990.

253. D. L. Windt, S. Donguy, J. Seely, B. Kjornrattanawanich, E. M. Gullikson, L. Golub, and E. DeLuca, "EUV Multilayers for Solar Physics," *Proc. SPIE* **5168**:1–11 (2004).
254. B. L. Henke, "X-Ray Multilayer Results," <http://henke.lbl.gov/multilayer/survey.html> (2008).
255. S. Bajt, J. B. Alameda, T. W. Barbee, W. M. Clift, J. A. Folta, B. Kaufmann, and E. A. Spiller, "Improved Reflectance and Stability of Mo-Si Multilayers," *Optical Engineering* **41**:1797–1804 (2002).
256. F. Eriksson, N. Ghafoor, F. Schäfers, E. M. Gullikson, S. Aouadi, S. Rohde, L. Hultman, and J. Birch, "Atomic Scale Interface Engineering by Modulated Ion-Assisted Deposition Applied to Soft X-Ray Multilayer Optics," *Applied Optics* **47**:4196–4204 (2008).
257. E. M. Gullikson, F. Salmassi, A. L. Aquila, and F. Dollar, "Progress in Short Period Multilayer Coatings for Water Window Applications," in *Proceedings, The 8th International Conference on The Physics of X-Ray Multilayer Structures*, Sapporo, Japan, March 12–16, 2006, abstract S8 O4.
258. E. Spiller and L. Golub, "Fabrication and Testing of Large Area Multilayer Coated X-Ray Optics," *Applied Optics* **28**:2969–2974 (1989).
259. C. Montcalm, B. T. Sullivan, M. Ranger, and H. Pepin, "Ultrahigh Vacuum Deposition Reflectometer System for the In Situ Investigation of Y/Mo Extreme-Ultraviolet Multilayer Mirrors," *Journal of Vacuum Science and Technology a Vacuum Surfaces and Films* **15**:3069–3081 (1997).
260. K. M. Skulina, C. S. Alford, R. M. Bionta, D. M. Makowiecki, E. M. Gullikson, R. Soufli, J. B. Kortright, and J. H. Underwood, "Molybdenum/Beryllium Multilayer Mirrors for Normal Incidence in the Extreme Ultraviolet," *Applied Optics* **34**:3727–3730 (1995).
261. N. Kaiser, S. Yulin, T. Feigl, H. Bernitzki, and H. Lauth, "EUV and Soft X-Ray Multilayer Optics," *Proc. SPIE* **5250**:109–118 (2004).
262. M. F. Ravet, X. Zhang-Song, As. Jerome, F. Delmotte, R. Mercier, M. Bougnet, P. Bouyries, and J. P. Delaboudiniere, "Ion Beam Deposited Mo/Si Multilayers for EUV Imaging Applications in Astrophysics," *Proc. SPIE* **5250**:99–108 (2004).
263. N. M. Ceglio, "Revolution in X-Ray Optics," *Journal of X-Ray Science and Technology* **1**:7–78 (1989).
264. J. Gautier, F. Delmotte, M. Roulliay, F. Bridou, M.-F. Ravet, and A. Jérôme, "Study of Normal Incidence of Three-Component Multilayer Mirrors in the Range 20–40 nm," *Applied Optics* **44**:384–390 (2005).
265. B. Kjornrattanawanich, D. L. Windt, J. F. Seely, and Y. A. Uspenskii, "SiC/Tb and Si/Tb Multilayer Coatings for Extreme Ultraviolet Solar Imaging," *Applied Optics* **45**:1765–1772 (2006).
266. S. Yulin, "Multilayer Coatings for EUV / Soft X-Ray Mirrors," in *Optical Interference Coatings* (eds. N. Kaiser and H. Pulker), Springer, Berlin, 2003, pp. 281–307.
267. P. Baumeister, "Design of Multilayer Filters by Successive Approximations," *Journal of the Optical Society of America* **48**:955–958 (1958).
268. M. A. Gisin, "Cutoff Interference Filters, Transparent to 25 μm ," *Soviet Journal of Optical Technology* **36**:191–194 (1969).
269. R. Jacobsson, "Matching a Multilayer Stack to a High Refractive Index Substrate by Means of an Inhomogeneous Layer," *Journal of the Optical Society of America* **54**:422–423 (1964).
270. E. A. Nesselov and G. P. Konyukhov, "Theory of a Cutoff Interference Filter," *Optics and Spectroscopy (USSR)* **31**:68–70 (1971).
271. A. Thelen, "Multilayer Filters with Wide Transmittance Bands," *Journal of the Optical Society of America* **53**:1266–1270 (1963).
272. M. Ploke, "Berechnung und Anwendung periodischer Mehrfachinterferenzschichten für ein Wärmerreflexionsfilter mit breiter Reflexionsbande," *Zeiss-Mitteilungen* **4**:279–294 (1967).
273. Optical Coating Laboratory Inc., "Effect of the Variation of Angle of Incidence and Temperature on Infrared Filter Characteristics," 2789 Northpoint Parkway, Santa Rosa, CA 95407–7397, USA, 1967.
274. Optical Coating Laboratory Inc., *Infrared Handbook*, Santa Rosa, CA, 1970.
275. Bausch & Lomb, "Bausch and Lomb Multi-Films," 1400 North Goodman Street, P. O. Box 540, Rochester, NY 14692–0450, USA, 1968.
276. Eastman Kodak Company, "Special Filters from Kodak for Technical Applications," 343 State Street, Bldg. 701, Rochester, NY 14650–3512, USA, 1968.
277. Infrared Industries Inc., "Infracron Interference Filters," Thin Film Products Division, 62 Fourth Avenue, Waltham, MA 02154, USA, 1966.

278. H. H. Schroeder and A. F. Turner, "A Commercial Cold Reflector," *Journal of the SMPTE* **69**:351–354 (1960).
279. F. E. Carlson, G. T. Howard, A. F. Turner, and H. H. Schroeder, "Temperature Reduction in Motion-Picture and Television Studios Using Heat-Control Coatings," *Journal of the SMPTE* **65**:136–139 (1956).
280. Corion Corporation, "Optical Filters and Coatings," (1992).
281. Heliotek—a Division of Textron Inc., "Data Sheets," 12500 Gladstane Avenue, Sylmar, CA 91342, USA, 1969.
282. A. F. Turner, "Heat Reflecting Coatings," in *Radiative Heat Transfer from Solid Materials* (eds. H. H. Blau and H. Fischer), Macmillan, New York, 1962.
283. A. Thelen, "The Use of Vacuum Deposited Coatings to Improve the Conversion Efficiency of Silicon Solar Cells in Space," *Progress in Astronautics and Rocketry* **3**:373–383 (1961).
284. D. L. Reynard and A. Andrew, "Improvement of Silicon Solar Cell Performance Through the Use of Thin Film Coatings," *Applied Optics* **5**:23–28 (1966).
285. A. J. Waldorf, National Research Council of Canada, Private communication, 1993.
286. F. Weiting and Z. Pengfei, "Determination of the performance of Edge Filters Containing PbTe on Cooling," *Infrared Physics* **33**:1–7 (1992).
287. B. Hoppert, "Heat Reducing Metal Mirrors," in *Proceedings, Annual Technical Conference*, Society of Vacuum Coaters, 1983, pp. 131–143.
288. A. Piegari and P. Polato, "Wideband Optical Coatings for Artwork Protection from Ultraviolet and Infrared Radiation Damage," *Proc. SPIE* **4829**:64–65 (2003).
289. L. Chang, "Infrared Dispersion of ZnS-Metal Films," Ph. D. thesis, Colorado State University (1964).
290. H. Schröder, "Oxide Layers Deposited from Organic Solutions," in *Physics of Thin Films* **5**, (eds. G. Hass and R. E. Thun), Academic, New York, 1969, pp. 87–141.
291. L. J. Vande Kleff, E. A. Murray, W. M. Frey, and P. W. Yunker, "Construction and Evaluation of Thin Film Beam Dividers," U.S. Army Aberdeen Res. Dev. Cent. Technical Note AD 698025, issued October 1969.
292. K. Hancock, "Membrane Optics," in *Proceedings, Electro-opt. Syst. Des. Conf.*, (editor K. A. Kopetzky), New York, Sept. 16–18, 1969, 1970, Industrial and Scientific Conference Management, Chicago, pp. 231–237.
293. P. L. Heinrich, R. C. Bastien, A. D. Santos, and M. Ostrelich, "Development of an all Dielectric Infrared Beamsplitter Operating in the 5 to 30 Micron Region," Perkin Elmer Technical Report CR-703 to NASA, issued February, 1967.
294. M. L. Lipshutz, "Optomechanical Considerations for Optical Beam Splitters," *Applied Optics* **7**:2326–2328 (1968).
295. H. P. Larson, "Evaluation of Dielectric Film Beamsplitters at Cryogenic Temperature," *Applied Optics* **25**:1917–1921 (1986).
296. F. C. Ho and J. A. Dobrowolski, "Neutral and Color-Selective Beam Splitting Assemblies with Polarization-Independent Intensities," *Applied Optics* **31**: 3813–3820 (1992).
297. H. Anders, *Thin Films in Optics* Focal Press, New York, 1967.
298. W. R. C. Rowley, "Some Aspects of Fringe Counting in Laser Interferometers," *IEEE Transaction Instrumentation Measures* **IM-15**:146–148 (1966).
299. H. Pohlack, "Zur Theorie der absortionsfreien achromatischen Lichtteilungs Spiegel," in *Jenaer Jahrbuch 1956* (eds. P. Görlich), Fischer, Jena, 1956, pp. 79–86.
300. L. A. Catalan and T. Putner, "Study of the Performance of Dielectric Thin Film Beam Dividing Systems," *British Journal of Applied Physics* **12**:499–502 (1961).
301. K. P. Miyake, "Computation of Optical Characteristics of Dielectric Multilayers," *Journal de Physique* **25**:255–257 (1964).
302. R. S. Sokolova, "Beam-Splitter Prisms with Dielectric Layers," *Soviet Journal of Optical Technology* **37**: 318–320 (1970).
303. A. L. Sergejeva, "Achromatic Nonabsorbing Beam Splitters," *Soviet Journal of Optical Technology* **38**:187–188 (1970).
304. P. B. Clapham, "The Design and Preparation of Achromatic Cemented Cube Beam-Splitters," *Optica Acta* **18**:563–575 (1971).
305. Z. Knittl, "Synthesis of Amplitude and Phase Achromatized Dielectric Mirrors," *Le Journal de Physique* **25**:245–249 (1964).

306. Oriol Optics Corporation, "Catalog Section C: Optical Filters," 2789 Northpoint Parkway, Santa Rosa, CA 95407-7397, USA, 1968.
307. S. F. Pellicori, "Beam Splitter and Reflection Reducing Coatings on ZnSe for 3-14 μm ," *Applied Optics* **18**:1966-1968 (1979).
308. C. K. Malek, J. Susini, A. Madouri, M. Ouahabi, R. Rivoira, F. R. Ladan, Y. Lepêtre, and R. Barchewitz, "Semitransparent Soft X-Ray Multilayer Mirrors," *Optical Engineering* **29**:597-602 (1990).
309. R. M. A. Azzam, "Variable-Reflectance Thin-Film Polarization-Independent Beam Splitters for 0.6328- and 10.6- μm Laser Light," *Optics Letters* **10**:110-112 (1985).
310. V. V. Veremei, V. N. Rozhdestvenskii, and A. B. Khazanov, "Radiation Dividers for the Infrared," *Soviet Journal of Optical Technology* **48**:618-619 (1981).
311. H. A. Macleod and Z. Milanovic, "Immersed Beam Splitters—an Old Problem," in *Proceedings, Optical Interference Coatings*, Optical Society of America, Tucson, AZ, 1992, pp. 28-30.
312. J. A. Dobrowolski, F. C. Ho, and A. Waldorf, "Beam Splitter for a Wide Angle Michelson Doppler Imaging Interferometer," *Applied Optics* **24**:1585-1588 (1985).
313. S. J. Refermat and A. F. Turner, "Polarization Free Beam Divider," U. S. patent 3559090, issued 26 January, 1971.
314. S. Itoh and M. Sawamura, "Achromatized Beam Splitter of Low Polarization," U. S. patent 4415233, issued 15 November, 1983.
315. A. A. M. Saleh, "Polarization-Independent, Multilayer Dielectrics at Oblique Incidence," *The Bell System Technical Journal* **54**:1027-1049 (1975).
316. A. Thelen, "Nonpolarizing Interference Films Inside a Glass Cube," *Applied Optics* **15**:2983-2985 (1976).
317. A. R. Henderson, "The Design of Non-Polarizing Beam Splitters," *Thin Solid Films* **51**:339-347 (1978).
318. Z. Knittl and H. Houserikova, "Equivalent Layers in Oblique Incidence: the Problem of Unsplit Admittances and Depolarization of Partial Reflectors," *Applied Optics* **21**:2055-2068. (1982).
319. C. M. d. Sterke, C. J. v. d. Laan, and H. J. Frakena, "Nonpolarizing Beam Splitter Design," *Applied Optics* **22**:595-601 (1983).
320. M. Zukic and K. H. Guenther, "Design of Nonpolarizing Achromatic Beamsplitters with Dielectric Multilayer Coatings," *Optical Engineering* **28**:165-171 (1989).
321. M. Gilo, "Design of a Nonpolarizing Beam Splitter Inside a Glass Cube," *Applied Optics* **31**:5345-5349 (1992).
322. L. Y. Chang and S. H. Mo, "Design of Non-Polarizing Prism Beam Splitter," in *Proceedings, Topical Meeting on Optical Interference Coatings*, Optical Society of America, Tucson, AZ, 1988, pp. 381-384.
323. Y. N. Konoplev, Y. A. Mamaev, V. N. Starostin, and A. A. Turkin, "Nonpolarizing 50% Beam Splitters," *Optics and Spectroscopy (USSR)* **71**:303-305 (1991).
324. H. Pohlack, "Zur theorie der absorptionsfreien achromatischen Lichtteilungs Spiegel," in *Jenaer Jahrbuch* (eds. C. Zeiss), Gustav Fischer Verlag, 1957, pp. 79-86.
325. J. S. Seeley, "Simple Nonpolarizing High-Pass Filter," *Applied Optics* **24**:742-744 (1985).
326. L. N. Kurochkina and N. M. Tulyakova, "Achromatic Lightsplitters for the Ultraviolet Spectrum," *Soviet Journal of Optical Technology* **52**:101-102 (1985).
327. I. J. Hodgkinson and R. G. Stuart, "Achromatic Reflector and Antireflection Coating Designs for Mixed Two-Component Deposition," *Thin Solid Films* **87**:151-158 (1982).
328. J. Mouchart, J. Begel, and S. Chalot, "Dépôts Achromatiques Partiellement Réfléchissants," *Journal of Modern Optics* **37**:875-888 (1990).
329. M. Banning, "Neutral Density Filters of Chromel," *Journal of the Optical Society of America* **37**:686-689 (1947).
330. Eastman Kodak Company, "Kodak Neutral Density Attenuators," 343 State Street, Bldg. 701, Rochester, NY 14650-3512, USA, 1969.
331. Corion Instrument Corporation, "Thin Film Optical Filters," 73 Jeffrey Avenue, Halliston, MA 01746-2082, USA, 1969.
332. Acton Research Corporation, "Optical Filters," P. O. Box 2215, 525 Main Street, Acton, MA 01720, USA, 1992.
333. H. Schröder, "Die Erzeugung von linearpolarisiertem Licht durch dünne dielektrische Schichten," *Optik* **3**:499-503 (1948).

334. R. Messner, "Die theoretischen Grundlagen optischer Interferenzpolarisatoren," *Feinwerktechnik* **57**: 142–147 (1953).
335. R. M. A. Azzam, "Single-Layer-Coated Optical Devices for Polarized Light," *Thin Solid Films* **163**:33–41 (1988).
336. M. Ruiz-Urbieta and E. M. Sparrow, "Reflection Polarization by a Transparent-Film-Absorbing-Substrate System," *Journal of the Optical Society of America* **62**:1188–1094 (1972).
337. J. T. Cox and G. Hass, "Highly Efficient Reflection-Type Polarizers for 10.6- μm CO₂ Laser Radiation Using Aluminum Oxide Coated Aluminum Mirrors," *Applied Optics* **17**:1657–1658 (1978).
338. W. W. Buchman, S. J. Holmes, and F. J. Woodberry, "Single-Wavelength Thin-Film Polarizers," *Journal of the Optical Society of America* **61**:1604–1606 (1971).
339. D. Blanc, P. H. Lissberger, and A. Roy, "The Design, Preparation and Optical Measurement of Thin Film Polarizers," *Thin Solid Films* **57**:191–198 (1979).
340. I. M. Minkov, "Theory of Dielectric Mirrors in Obliquely Incident Light," *Optics and Spectroscopy (USSR)* **33**:175–178 (1973).
341. A. Thelen, "Avoidance or Enhancement of Polarization in Multilayers," *Journal of the Optical Society of America* **70**:118–121 (1980).
342. P. B. Clapham, M. J. Downs, and R. J. King, "Some Applications of Thin Films to Polarization Devices," *Applied Optics* **8**:1965–1974 (1969).
343. P. Kard, "Theory of Achromatic Multilayer Interference Polarizers," *Eesti NSV Tead. Akad. Toim., Fuus. - Mat.* **9**:26–32 (1960).
344. W. Geffcken, "Interference light filter," DB patent 913005, issued June 8, 1954/44.
345. S. M. MacNeille, "Beam splitter," U. S. patent 2403731, issued 6 July, 1946.
346. R. S. Sokolova and T. N. Krylova, "Interference Polarizers for the Ultraviolet Spectral Region," *Optics and Spectroscopy (USSR)* **14**:213–215 (1963).
347. J. Wimperis, Interoptics, A Division of Lumonics, Inc., 14 Capella Court, Nepean (Ottawa) Ontario, Canada K2E 7V6, Private communication, 1993.
348. M. Banning, "Practical Methods of Making and Using Multilayer Filters," *Journal of the Optical Society of America* **37**:792–297 (1947).
349. W. B. Wetherell, "Polarization Matching Mixer in Coherent Optical Communications Systems," *Optical Engineering* **28**:148–156 (1989).
350. J. Mouchart, J. Begel, and E. Duda, "Modified MacNeille Cube Polarizer for a Wide Angular Field," *Applied Optics* **28**:2847–2853 (1989).
351. R. Austin, "Thin Film Polarizing Devices," *Electro-Optical Systems Design* February issue, pp. 30–35, 1974.
352. R. P. Netterfield, "Practical Thin-Film Polarizing Beam-Splitters," *Optica Acta* **24**:69–79 (1977).
353. J. A. Dobrowolski and A. J. Waldorf, "High-Performance Thin Film Polarizer for the UV and Visible Spectral Regions," *Applied Optics* **20**:111–116 (1981).
354. V. P. Sobol', R. A. Petrenko, and A. S. Dimitreev, "Interference Polarizer Employing Optical Contact," *Soviet Journal of Optical Technology* **53**:692 (1986).
355. M. Gilo and K. Rabinovitch, "Design Parameters of Thin-Film Cubic-Type Polarizers for High Power Lasers," *Applied Optics* **26**:2518–2521 (1987).
356. L. Li and J. A. Dobrowolski, "Visible Broadband, Wide-Angle, Thin-Film Multilayer Polarizing Beam Splitter," *Applied Optics* **35**:2221–2225 (1996).
357. L. Li and J. A. Dobrowolski, "High-Performance Thin-Film Polarizing Beam Splitter Operating at Angles Greater Than the Critical Angle," *Applied Optics* **39**:2754–2771 (2000).
358. P. Ma, L. Li, F. Lin, and J. A. Dobrowolski, "Manufacture of High Performance Polarizing Beam Splitter for Projection Display Applications," in *Proceedings, Optical Interference Coatings on CD-ROM*, The Optical Society of America, Washington D. C., 2007, pp. TuC1.
359. L. Li and Z. Pang, "Thin Film Polarizing Device Having Metal-Dielectric Films," US patent 6,317,264, issued November 13, 2001.
360. E. Cojocaru, "Comparison of Theoretical Performances for Different Single-Wavelength Thin-Film Polarizers," *Applied Optics* **31**:4501–4504 (1992).

361. R. M. A. Azzam, "Efficient Infrared Reflection Polarizers Using Transparent High-Index Films on Transparent Low-Index Substrates," *Proc. SPIE* **652**:326–332 (1986).
362. T. F. Thonn and R. M. A. Azzam, "Multiple-Reflection Polarizers Using Dielectric-Coated Metallic Mirrors," *Optical Engineering* **24**:202–206 (1985).
363. L. Songer, "The Design and Fabrication of a Thin Film Polarizer," *Optical Spectra*, October issue, pp. 49–50, 1978.
364. D. Lees and P. Baumeister, "Versatile Frustrated-Total-Reflection Polarizer for the Infrared," *Optics Letters* **4**:66–67 (1979).
365. H. Lotem and K. Rabinovitch, "Penta Prism Laser Polarizer," *Applied Optics* **32**:2017–2020 (1993).
366. H. Winter, H. H. Bukow, and P. H. Heckmann, "A High Transmission Triple-Reflection Polarizer for Lyman- α Radiation," *Optics Communications* **11**:299 (1974).
367. G. Hass and W. R. Hunter, "Reflection Polarizers for the vacuum Ultraviolet Using Al + MgF₂ Mirrors and MgF₂ Plate," *Applied Optics* **17**:76–82 (1978).
368. W. R. Hunter, "Design Criteria for Reflection Polarizers and Analyzers in the Vacuum Ultraviolet," *Applied Optics* **17**:1259–1270 (1978).
369. T. A. Remneva, A. V. Kozhevnikov, and M. M. Nikitin, "Polarizer of Radiation in the Vacuum Ultraviolet," *Journal of Applied Spectroscopy (USSR)* **25**:1587–1590 (1967).
370. Z. Wang, "Non-Periodic Multilayer Coatings in EUV, Soft X-Ray and X-Ray Range," *Proc. SPIE* **7101**:10-1-15 (2008).
371. P. Dhez, "Polarizers and Polarimeters for the X-UV Range," *Nuclear Instrumentation Methods Physics Research* **A261**:66–71 (1987).
372. M. Yanagihara, T. Maehara, H. Normura, M. Yamamoto, T. Namioka, and H. Kimura, "Performance of a wide Bandband Multilayer Polarizer for Soft X-Rays," *Review of Scientific Instruments* **63**:1516–1518 (1992).
373. F. Eriksson, N. Ghafoor, F. Schäfers, E. M. Gullikson, S. Aouadi, S. Rohde, L. Hultman, and J. Birch, "Atomic Scale Interface Engineering by Modulated Ion-Assisted Deposition Applied to Soft X-Ray Multilayer Optics," *Applied Optics* **47**:4196–4204 (2008).
374. Z. Wang, H. Wang, J. Zhu, F. Wang, Z. Gu, L. Chen, A. G. Michette, A. K. Powell, S. J. Pfauntsch, and F. Schäfers, "Broadband Multilayer Polarizers for the Extreme Ultraviolet," *Journal of Applied Physics* **99**:056108 1-3 (2006).
375. Z. Wang, H. Wang, J. Zhu, Y. Xu, S. Zhang, C. Li, and F. Wang, "Extreme Ultraviolet Broadband Mo/Y Multilayer Analyzers," *Applied Physics Letters* **89**:241120-1-3 (2006).
376. Anon, "Metal-Dielectric Interference Filters," in *Physics of Thin Films* 9, (eds. G. Hass, M. H. Francombe, and R. W. Hoffman), Academic Press, New York, (1977), pp. 73–144.
377. E. E. Barr, "The Design and Construction of Evaporated Multilayer Filters for Use in Solar Radiation Technology," in *Advances in Geophysics* 14, (eds. A. J. Drummond), Academic Press, (1970), pp. 391–412.
378. R. Morf and R. E. Kunz, "Dielectric Filter Optimization by Simulated Thermal Annealing," *SPIE—Thin Film Technologies III* **1019**:211 (1988).
379. N. S. Gluck and W. J. Gunning, "Patterned Infrared Spectral Filter Directly Deposited onto Cooled Substrates," *Applied Optics* **28**:5110–5114 (1989).
380. T. A. Mooney, 2 Lyberty Way, Westford, MA 01886, U. S. A. Personal communication, 1993.
381. W. H. Southwell, W. J. Gunning, and R. L. Hall, "Narrow-Bandpass Filter Using Partitioned Cavities," *SPIE - Optical Thin Films II. New Developments* **678**:177–184 (1986).
382. W. Geffcken, "Interference Light Filter," DB patent 716153, issued December 8, 1939.
383. D. H. Harrison, "MDM Bandpass Filters for the Vacuum Ultraviolet," *Applied Optics* **7**:210 (1968).
384. D. J. Bradley, B. Bates, C. O. L. Juulman, and T. Kohno, "Recent Developments in the Application of the Fabry-Perot Interferometer to Space Research," *Journal de Physique* **28**:C2-280–C2-283 (1967).
385. Balzers Aktiengesellschaft für Hochvakuumtechnik und Dünne Schichten, "Interference Filters," FL-9496 Balzers, Principality of Liechtenstein, 1968.
386. Schott Glaswerke, "Monochromatic Interference Filters," Geschäftsbereich Optik, Hattenbergstrasse 10, D-6500 Mainz, Germany, 1965.
387. A. F. Turner, "Some Current Developments in Multilayer Optical Films," *Le Journal de Physique et le Radium* **11**:444–460 (1950).

388. P. W. Baumeister, "Optical Tunneling and Its Applications to Optical Filters," *Applied Optics* **6**:897–905 (1967).
389. S. W. Warren, "Properties and Performance of Basic Designs of Infrared Interference Filters," *Infrared Physics* **8**:65–78 (1968).
390. A. Thelen, "Design of Multilayer Interference Filters," in *Physics of Thin Films* 5, (eds. G. Hass and R. E. Thun), Academic Press, New York & London, 1969, pp. 47–86.
391. W. Geffcken, "The Wave-Band Filter, an Interference Filter of Very High Efficiency," *Zeitschrift für Angewandte Physik* **6**:249–250 (1954).
392. A. F. Turner, "Wide Passband Multilayer Filters," *Journal of the Optical Society of America* **42**:878 (1952).
393. S. D. Smith, "Design of Multilayer Filters by Considering Two Effective Interfaces," *Journal of the Optical Society of America* **48**:43–50 (1958).
394. H. Schröder, "Properties and Applications of Oxide Layers Deposited on Glass from Organic Solutions," *Optica Acta* **9**:249–254 (1962).
395. R. S. Meltzer, Industrial Optics Division, Yardney Razdow Laboratories, Private communications, 1968.
396. R. H. Eather and D. L. Reasoner, "Spectrophotometry of Faint Light Sources with a Tilting-Filter Photometer," *Applied Optics* **8**:227–242 (1969).
397. R. R. Austin, Perkin-Elmer Corporation, Private communication, November 11, 1968.
398. M. A. Cohendet and B. Saudreau, DRME-66-34-443, 1967.
399. O. A. Motovilov, "Narrow-Band Interference Filters for the Ultraviolet Region of the Spectrum," *Optics and Spectroscopy (USSR)* **22**:537–538 (1967).
400. R. G. T. Neilson and J. Ring, "Interference Filters for the Near Ultra-Violet," *Journal de Physique* **28**:C2-270–C2-275 (1967).
401. A. F. Turner and R. Walsh, "Interference Filters for the 8–13 Micron Atmospheric Window," *Bausch & Lomb Technical Report no. 2*, issued 1961.
402. S. D. Smith and J. S. Seely, "Multilayer Filters for the Region 0.8 to 100 Microns," AF61(052)-833, issued May 7, 1968.
403. I. H. Blifford, Jr., "Factors Affecting the Performance of Commercial Interference Filters," *Applied Optics* **5**:105–111 (1966).
404. Baird Atomic, Inc. "Optical Components," 125 Middlesex Turnpike, Bedford, MA 01730–1468, 1968.
405. Spectrum Systems Division—Barnes Engineering Co., "Visible and Near Infrared Filters, Bull. SS-14," 211 Second Avenue, Waltham, MA 02154, USA, 1967.
406. Infrared Industries Inc., "Thin Film Products at Infrared Industries, Inc.," Thin Film Products Division, 62 Fourth Avenue, Waltham, MA 02154, USA, 1967.
407. P. Baumeister, "Use of Microwave Prototype Filters to Design Multilayer Dielectric Bandpass Filters," *Applied Optics* **21**:2965–2967 (1982).
408. A. Zheng, J. S. Seeley, R. Hunneman, and G. J. Hawkins, "Design of Narrowband Filters in the Infrared Region," *Infrared Physics* **31**:237–244 (1991).
409. J. Minowa and Y. Fujii, "High Performance Bandpass Filters for WDM Transmission," *Applied Optics* **23**:193–194 (1984).
410. P. W. Baumeister, V. R. Costich, and S. C. Pieper, "Bandpass Filters for the Ultraviolet," *Applied Optics* **4**:911–914 (1965).
411. P. H. Berning and A. F. Turner, "Induced Transmission in Absorbing Films Applied to Band Pass Filter Design," *Journal of the Optical Society of America* **47**:230–239 (1957).
412. R. L. Maier, "2M Interference Filters for the Ultraviolet," *Thin Solid Films* **1**:31–37 (1967).
413. P. W. Baumeister, "Radiant Power Flow and Absorptance in Thin Films," *Applied Optics* **8**:423–436 (1969).
414. R. J. Holloway and P. H. Lissberger, "The Design and Preparation of Induced Transmission Filters," *Applied Optics* **8**:653–660 (1969).
415. B. V. Landau and P. H. Lissberger, "Theory of Induced-Transmission Filters in Terms of the Concept of Equivalent Layers," *Journal of the Optical Society of America* **62**:1258–1264 (1972).
416. N. P. Matshina, E. A. Nesmelov, I. K. Nagimov, R. M. Validov, and N. N. Soboleva, "Theory of Narrow-Band Filters with Induced Transparency," *Journal of Applied Spectroscopy (USSR)* **55**:1273–1278 (1991).

417. D. Gershenson and L. Sossi, "Conditions for Maximum Transmission of Metal-Dielectric Interference Filters," *Eesti NSV Tead Akad Toim, Fuus. Mat.* **41**:142–149 (1992).
418. V. I. Tsypin and E. A. Sukhanov, "Metal-Dielectric Filters for the 200-350 nm Spectral Regions," *Soviet Journal of Optical Technology* **59**: 438–439 (1992).
419. F. A. Korolev, A. Y. Klement'eva, and T. F. Meshcheryakova, "Interference Filters with a Transmission Band of 1.5 a Half-Width," *Optics and Spectroscopy (USSR)* **9**:341–343 (1960).
420. R. W. Wood, "Some New Cases of Interference and Diffraction," in *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **7**, Taylor and Francis, London, 1904, pp. 376–388.
421. A. I. Kartashev and N. M. Syromyatnikova, "Mica Interference Filters," in *Wavelength of Light as a Standard of Length, Proc. Mendeleev State Sci. Res. Inst.*, **7**, Leningrad, 1947, pp. 86–93.
422. J. Ring, R. Beer, and V. Hewison, "Reflectance Multilayers," *J. Phys. Radium* **19**:321–323 (1958).
423. J. A. Dobrowolski, "Mica Interference Filters with Transmission Bands of Very Narrow Half-Widths," *Journal of the Optical Society of America* **49**:794–806 (1959).
424. G. S. Cheremukhin, V. P. Rozhnov, and G. I. Golubeva, "Narrow-Band Interference Filters Made from Mica," *Soviet Journal of Optical Technology* **43**:312–315 (1976).
425. D. R. Herriott, J. R. Wimperis, and D. L. Perry, "Filters, Wave Plates and Protected Mirrors Made of Thin Polished Supported Layers," *Journal of the Optical Society of America* **4**:546 (abstract) (1965).
426. J. D. Rehnberg, "Design Considerations for an Orbiting Solar Telescope," in *Proceedings, Electro-Opt. Syst. Des. Conf.*, (editor K. A. Kopetzky), New York, Sept 16–18, 1969, 1970, Industrial and Scientific Conference Management, Inc., Chicago, pp. 69–73.
427. A. Title, "Fabry-Perot Interferometers as Narrow-Band Optical Filters: Part one—Theoretical considerations," Harvard College Observatory Technical Report TR-18, issued June, 1970.
428. R. Fisher, "On the Use of a Solid Fabry-Perot Interferometer for Coronal Photography," *Solar Physics* **18**:253–257 (1971).
429. J. F. Markey and R. R. Austin, "High Resolution Solar Observations: The Hydrogen-Alpha Telescopes on Skylab," *Applied Optics* **16**:917–921 (1977).
430. R. R. Austin, "The Use of Solid Etalon Devices as Narrow Band Interference Filters," *Optical Engineering* **11**:65–69 (1972).
431. S. D. Smith and C. R. Pidgeon, "Application of Multiple Beam Interferometric Methods to the Study of CO₂ Emissions at 15 μm ," *Mem. Soc. R. Liege 5ieme serie* **9**:336–349 (1963).
432. A. E. Roche and A. M. Title, "Tilt Tunable Ultra Narrow-Band Filters for High Resolution Infrared Photometry," *Applied Optics* **14**:765–770 (1975).
433. J. -M. Saurel and M. Candille, "Réalisation de filtres multidielectriques sur un substrat biréfringent, déformable et de faible épaisseur (mylar)," *Thin Solid Films* **16**:313–324 (1973).
434. M. Candille and J. M. Saurel, "Réalisation de filtres 'souble onde' à bandes passantes très étroites sur supports en matière plastique (mylar)," *Optica Acta* **21**:947–962 (1974).
435. P. W. Baumeister, F. A. Jenkins, and M. A. Jeppesen, "Characteristics of the Phase-Dispersion Interference Filter," *Journal of the Optical Society of America* **49**:1188–1190 (1959).
436. R. R. Austin, "Narrow Band Interference Light Filter," US patent 3528726 issued 15 September, 1970.
437. P. Giacomo, P. W. Baumeister, and F. A. Jenkins, "On the Limiting Bandwidth of Interference Filters," *Proc. Phys. Soc.* **73**:480–489 (1959).
438. S. D. Smith and O. S. Heavens, "A Tunable Infra-Red Interference Filter," *Journal of Scientific Instruments* **34**:492–496 (1957).
439. P. D. Atherton, N. K. Reay, J. Ring, and T. R. Hicks, "Tunable Fabry-Perot Filters," *Optical Engineering* **20**:806–814 (1981).
440. J. V. Ramsay, "Control of Fabry-Perot Interferometers and Some Unusual Applications," *The Australian Physicist* **5**:87–89 (1968).
441. J. V. Ramsay, H. Kobler, and E. G. V. Mugridge, "A New Tunable Filter with a Very Narrow Pass-Band," *Solar Physics* **12**:492–501 (1970).
442. C. H. Burton, A. J. Leistner, and D. M. Rust, "Electrooptic Fabry-Perot Filter Development for the Study of Solar Oscillations," *Applied Optics* **26**:2637–2642 (1987).

443. D. B. McKenney and P. N. Slater, "Design and Use of Interference Passband Filters with Wide-Angle Lenses for Multispectral Photography," *Applied Optics* **9**:2435–2440 (1970).
444. J. A. Dobrowolski and R. C. Bastien, "Square-Top Transmission Band Interference Filters for the Infra-Red," *Journal of the Optical Society of America* **53**:1332 (1963).
445. D. McKenney and A. F. Turner, "Ultra-Wide Bandpass Filters Developed for Spectral Photography," *Univ. Ariz. Opt. Cent. Newsl.* **1**, (1967).
446. J. Michael and Perkin-Elmer Corporation, Private communication, 1968.
447. E. Pelletier and H. A. Macleod, "Interference Filters with Multiple Peaks," *Journal of the Optical Society of America* **72**:683–687 (1982).
448. W. A. Hovis, Jr., W. A. Kley, and M. G. Strange, "Filter Wedge Spectrometer for Field Use," *Applied Optics* **6**:1057 (1967).
449. A. Thelen, "Circularly Wedged Optical Coatings. I. Theory," *Applied Optics* **4**:977–981 (1965).
450. J. H. Apfel, "Circularly Wedged Optical Coatings. II. Experimental," *Applied Optics* **4**:983–985 (1965).
451. V. A. Martsinovskii and F. K. Safiullin, "Determination of the Conditions of Formation of Annular Coatings with Linear Dependence of Thickness on Rotation Angle," *Journal of Applied Spectroscopy (USSR)* **39**:976–980 (1983).
452. V. P. Avilov and A. Khosilov, "Apparatus for Producing Annular Variable Interference Filters," *Soviet Journal of Optical Technology* **55**:613–615 (1988).
453. I. M. Minkov, "Calculation of Narrow-Band Circular Wedge Filter for 4–12 μm Spectral Region," *Soviet Journal of Optical Technology* **58**:491–492 (1991).
454. V. L. Yen, "Circular Variable Filters," *Optical Spectra* **3**:78–84 (1969).
455. V. P. Avilov, V. I. Tsylin, and E. M. Shipulin, "Circular-Wedge Filter for the 0.24–0.40 μm Region," *Soviet Journal of Optical Technology* **54**:515 (1987).
456. A. Mussett, Optical Coating Laboratory, Private communication, 1969.
457. P. H. Lissberger, "Effective Refractive Index as a Criterion of Performance of Interference Filters," *Journal of the Optical Society of America* **58**:1586–1590 (1968).
458. P. H. Lissberger, "Properties of All-Dielectric Interference Filters. I. A New Method of Calculation," *Journal of the Optical Society of America* **49**:121–125 (1959).
459. C. R. Pidgeon and S. D. Smith, "Resolving Power of Multilayer Filters in Nonparallel Light," *Journal of the Optical Society of America* **54**:1439–1466 (1964).
460. D. J. Hemingway and P. H. Lissberger, "Effective Refractive Indices of Metal-Dielectric Interference Filters," *Applied Optics* **6**:471–476 (1967).
461. S. L. Linder, "Optimization of Narrow Optical Spectral Filters for Nonparallel Monochromatic Radiation," *Applied Optics* **6**:1201–1204 (1967).
462. P. G. Kard, "On Elimination of the Doublet Structure of the Transmission Band in a Total-Reflection Light-Filter," *Optics and Spectroscopy (USSR)* **6**:244–246 (1959).
463. G. P. Konyukhov and E. A. Nesmelov, "On the Theory of Dielectric Narrow-Band Light Filters," *Zhurnal Prikladnoi Spektroskopii* **11**:468– (1969).
464. P. Baumeister, "Bandpass Design—Application to Nonnormal Incidence," *Applied Optics* **31**:504–512 (1992).
465. W. Wilmot and E. R. Schineller, "A Wide-Angle Narrow-Band Optical Filter," *Journal of the Optical Society of America* **56**:549 (1966).
466. E. R. Schineller and R. P. Flam, "Development of a Wide-Angle Narrow-Band Optical Filter," in *Proceedings, Spring Meeting Program—Optical Society of America*, 1968, pp. FF11-1–FF11-4.
467. J. Meaburn, "The Stability of Interference Filters," *Applied Optics* **5**:1757–1759 (1966).
468. J. Schild, A. Steudel, and H. Walter, "The Variation of the Transmission Wavelength of Interference Filters by the Influence of Water Vapor," *Journal de Physique* **28**:C2-276–C2-279 (1967).
469. L. D. Lazareva, "The Effect of Temperature on the Position of the Passband Maximum in Interference Filters," *Optical Technology* **36**:801–802 (1970).
470. S. A. Furman and M. D. Levina, "Effect of Moisture on the Optical Characteristics of Narrow Band Interference Filters," *Optics and Spectroscopy (USSR)* **30**:404–408 (1971).

471. S. A. Furman and M. D. Levina, "Strict Stabilization of the Optical Characteristics of Narrow-Band Interference Filters," *Optical Technology* **38**:374–375 (1971).
472. S. A. Furman and M. D. Levina, "Stabilization of the Location of the Transmission Band of a Narrow-Band Dielectric Interference Filter," *Optical Technology* **38**:272–275 (1971).
473. M. D. Levina and S. A. Furman, "Improving the Stability of the Optical Characteristics of Metal-Dielectric Filters," *Soviet Journal of Optical Technology* **49**:128–129 (1982).
474. D. R. Gibson and P. H. Lissberger, "Optical Properties of Narrowband Spectral Filter Coatings Related to Layer Structure and Preparation," *Applied Optics* **22**:269–281 (1983).
475. A. Brunsting, M. A. Kheiri, D. F. Simonaitis, and A. J. Dosmann, "Environmental Effects on All-Dielectric Bandpass Filters," *Applied Optics* **25**:3235–3241 (1986).
476. S. A. Furman and M. D. Levina, "Investigations of Temperature Dependence of Optical Properties of Dielectric Narrow-Band Filters," *Optics and Spectroscopy (USSR)* **28**:412–416 (1970).
477. S. A. Furman, "Effect of Temperature, Angle of Incidence, and Dispersion of Index of Refraction of Layers on the Position of the Pass Bands of a Dielectric Narrow-Band Filters," *Optics and Spectroscopy (USSR)* **28**:218–222 (1970).
478. A. S. Chaikin and V. V. Pukhonin, "An Investigation of the Temperature and Time Dependence of the Parameters of Narrow-Band Interference Filters," *Journal of Applied Spectroscopy (USSR)* **13**:1513–1515 (1970).
479. G. I. Golubeva, M. S. Abakumova, and A. M. Klochkov, "Results of Industrial Studies," *Optical Technology* **38**:228–229 (1971).
480. A. M. Zheng, J. S. Seeley, R. Hunneman, and G. J. Hawkins, "Ultraviolet Filters with Good Performance when Tilted and Cooled," *Applied Optics* **31**:4336–4338 (1992).
481. W. Feng and Y. Yan, "Shift in Infrared Interference Filters at Cryogenic Temperature," *Applied Optics* **31**:6591–6592 (1992).
482. R. Mark, D. Morand, and S. Waldstein, "Temperature Control of the Bandpass of an Interference Filter," *Applied Optics* **9**:2305–2310 (1970).
483. V. M. Novikov, "Vacuum Installation with a Manipulator," *Soviet Journal of Optical Technology* **35**:121–122 (1968).
484. O. M. Sorokin and V. A. Blank, "Thin-Film Metal and Semiconductor Filters for the Vacuum Ultraviolet," *Optical Technology* **37**:343–346 (1970).
485. T. Barbee and J. H. Underwood, "Solid Fabry-Perot Etalons for X-Rays," *Optics Communications* **48**:161–166 (1983).
486. Y. Lepetre, R. Rivoira, R. Philip, and G. Rasigni, "Fabry-Perot Etalons for X-Rays: Construction and Characterization," *Optics Communications* **51**:127–130 (1984).
487. Y. Lepetre, I. K. Schuller, G. Rasigni, R. Rivoira, and R. Philip, "Novel Characterization of Thin Film Multilayered Structures: Microcleavage Transmission Electron Spectroscopy," *Proceedings of the Society of Photo-Optical Instrumentation Engineers* **563**:258–263 (1985).
488. R. J. Bartlett, W. J. Trela, D. R. Kania, M. P. Hockaday, T. W. Barbee, and P. Lee, "Soft X-Ray Measurements of Solid Fabry-Perot Etalons," *Optics Communications* **55**:229–235 (1985).
489. V. R. Costich, "Coatings for 1, 2, Even 3 Wavelengths," *Laser Focus Magazine* **41**–45 (1969).
490. J. D. Rancourt and W. T. Beauchamp, "Articles Having Improved Reflectance Suppression," U. S. patent 4578527, issued 25 March, 1986.
491. W. C. Herrmann and D. E. Morton, "A Design Technique for Multiband AR Coatings," in *Proceedings, 33rd Annual Technical Conference, Society of Vacuum Coaters, New Orleans, Louisiana, 1990*, pp. 246–249.
492. L. Li, J. A. Dobrowolski, J. D. Sankey, and J. R. Wimperis, "Antireflection Coatings for Both Visible and Far Infrared Spectral Regions," *Applied Optics* **31**:6150–6156 (1992).
493. L. Li and J. A. Dobrowolski, "Design of Optical Coatings for Two Widely Separated Spectral Regions," *Applied Optics* **32**:2969–2975 (1993).
494. J. A. Dobrowolski and L. Li, "Design of Optical Coatings for Three or More Separated Spectral Regions," *Applied Optics* **34**:2934–2940 (1995).
495. H. A. Macleod, "Structure-Related Optical Properties of Thin Films," *Journal of Vacuum Science and Technology* **A4**:418–422 (1986).

496. T. Motohiro and Y. Taga, "Thin Film Retardation Plate by Oblique Deposition," *Applied Optics* **28**:2466–2482 (1989).
497. J. H. Apfel, "Graphical Method to Design Multilayer Phase Retarders," *Applied Optics* **20**:1024–1029 (1981).
498. J. H. Apfel, "Phase Retardance of Periodic Multilayer Mirrors," *Applied Optics* **21**:733–738 (1982).
499. W. H. Southwell, "Multilayer Coating Design Achieving a Broadband 90° Phase Shift," *Applied Optics* **19**:2688–2692 (1980).
500. N. V. Grishina, "Synthesis of Mirrors with Constant Phase Difference Under Oblique Incidence of Light," *Optics and Spectroscopy (USSR)* **69**:262–265 (1990).
501. P. Lostis, "Étude, réalisation et contrôle de lames minces introduisant une différence de marche déterminée entre deux vibrations rectangulaires," *Revue d'Optique, Théorique et Instrumentale* **38**:1–28 (1959).
502. J. H. Apfel, "Graphical Method to Design Internal Reflection Phase Retarders," *Applied Optics* **23**:1178–1183 (1984).
503. R. M. A. Azzam and M. E. R. Khan, "Polarization-Preserving Single-Layer-Coated Beam Displacers and Axicons," *Applied Optics* **21**:3314–3322 (1982).
504. E. Cojocaru, "Polarization-Preserving Totally Reflecting Prisms," *Applied Optics* **31**:4340–4342 (1992).
505. E. Spiller, "Totally Reflecting Thin-Film Phase Retarders," *Applied Optics* **23**:3544–3549 (1984).
506. J. M. Bennett, "A Critical Evaluation of Rhomb-Type Quarterwave Retarders," *Applied Optics* **9**:2123–2129 (1970).
507. I. Filinski and T. Skettrup, "Achromatic Phase Retarders Constructed from Right-Angle Prisms: Design," *Applied Optics* **23**:2747–2751 (1984).
508. R. J. King, "Quarter-Wave Retardation Systems Based on the Fresnel Rhomb Principle," *Journal of Scientific Instruments* **43**:627–622 (1966).
509. R. M. A. Azzam, "ZnS-Ag Film-Substrate Parallel-Mirror Beam Displacers that Maintain Polarization of 10.6 μm Radiation Over a Wide Range of Incidence Angles," *Infrared Physics* **23**:195–197 (1983).
510. T. F. Thonn and R. M. A. Azzam, "Three-Reflection Halfwave and Quarterwave Retarder Using Dielectric-Coated Metallic Mirrors," *Applied Optics* **23**:2752–2759 (1985).
511. H. Pohlack, "Über die reflexionsvermindernde Wirkung dünner Metallschichten auf Glas," in *Jenaer Jahrbuch 1956* (eds. P. Görlich), Fischer, Jena, 1956, pp. 87–93.
512. G. Hass, H. H. Schroeder, and A. F. Turner, "Mirror Coatings for Low Visible and High Infrared Reflectance," *Journal of the Optical Society of America* **46**:31–35 (1956).
513. J. A. Dobrowolski, B. T. Sullivan, and R. C. Bajcar, "Optical Interference, Contrast-Enhanced Electroluminescent Device," *Applied Optics* **31**:5988–5996 (1992).
514. J. A. Dobrowolski, E. H. Hara, B. T. Sullivan, and A. J. Waldorf, "High Performance Optical Wavelength Multiplexer-Demultiplexer," *Applied Optics* **31**:3800–3806 (1992).
515. J. A. Dobrowolski and R. A. Kemp, "Flip-Flop Thin-Film Design Program with Enhanced Capabilities," *Applied Optics* **31**:3807–3812 (1992).
516. J. A. Dobrowolski, "Versatile Computer Program for Absorbing Optical Thin-Film Systems," *Applied Optics* **20**:74–81 (1981).
517. J. H. Apfel and R. M. Gelber, "Multilayer Filter with Metal Dielectric Period," U. S. patent 3679291, issued 14 March, 1972.
518. R. J. Fay and J. R. Cicotta, "Neutral density Filter Element with Reduced Surface Reflection," U. S. patent 3781089, issued December 25, 1973.
519. J. A. Dobrowolski, L. Li, and R. A. Kemp, "Metal/Dielectric Transmission Interference Coatings with Low Reflectance 1. Design," *Applied Optics* **34**: 5673–5683 (1995).
520. B. T. Sullivan and K. L. Byrt, "Metal/Dielectric Transmission Interference Coatings with Low Reflectance 2. Experimental Results," *Applied Optics* **34**: 5684–5694 (1995).
521. J. R. Jacobsson, "Protective Device for Protection Against Radiation During Welding," U. S. patent 4169655, issued October 2, 1979.
522. ASTM, "Symposium on Electroless Nickel Plating," *Spec. Tech. Publ.* **265**, (1959).
523. R. P. Madden, "Preparation and Measurement of Reflecting Coatings for the Vacuum Ultraviolet," in *Physics of Thin Films* 1, (ed. G. Hass), Academic Press, New York, 1963, pp. 123–186.

524. M. L. Scott, P. N. Arendt, B. J. Cameron, J. M. Saber, and B. E. Newnam, "Extreme Ultraviolet Reflectance Degradation of Aluminum and Silicon from Surface Oxidation," *Applied Optics* **27**:1503–1507 (1988).
525. H. E. Bennett, M. Silver, and E. J. Ashley, "Infrared Reflectance of Aluminum Evaporated in Ultra-High Vacuum," *Journal of the Optical Society of America* **53**:1089–1095 (1963).
526. J. M. Bennett and E. J. Ashley, "Infrared Reflectance and Emittance of Silver and Gold Evaporated in Ultrahigh Vacuum," *Applied Optics* **4**:221–224 (1965).
527. G. Hass and L. Hadley, in *American Institute of Physics Handbook* (ed. O. F. Gray), McGraw-Hill, New York, 1972, pp. 6–118.
528. W. R. Hunter, "High Reflectance Coatings for the Extreme Ultra-Violet," *Optica Acta* **9**:255–268 (1962).
529. G. Hass and G. F. Jacobus, "Optical Properties of Evaporated Iridium in the Vacuum Ultraviolet from 500 to 2000 Å," *Journal of the Optical Society of America* **57**:758–762 (1967).
530. J. T. Cox, G. Hass, and W. R. Hunter, "Optical Properties of Evaporated Rhodium Films Deposited at Various Substrate Temperatures in the Vacuum Ultraviolet from 150 to 2000 Å," *Journal of the Optical Society of America* **61**:360–364 (1971).
531. J. T. Cox, G. Hass, and J. B. Ramsey, "Reflectance of Evaporated Rhenium and Tungsten Films in the Vacuum Ultraviolet from 300 to 2000 Å," *Journal of the Optical Society of America* **62**:781–785 (1972).
532. J. T. Cox, G. Hass, and J. B. Ramsey, "Reflectance and Optical Constants of Evaporated Osmium in the Vacuum Ultraviolet from 300 to 2000 Å," *Journal of the Optical Society of America* **63**:435–438 (1973).
533. J. F. Seely, G. E. Holland, W. R. Hunter, R. P. McCoy, K. F. Dymond, and M. Corson, "Effect of Oxygen Atom Bombardment on the Reflectance of Silicon Carbide Mirrors in the Extreme Ultraviolet Region," *Applied Optics* **18**:1805–1810 (1979).
534. G. Hass, "Filmed Surfaces for Reflecting Objects," *Journal of the Optical Society of America* **45**:945–952 (1955).
535. R. Denton and Denton Vacuum Inc., Private communication, 1970.
536. V. D. Vvedenskii, R. Y. Pinskaya, S. A. Furman, and T. V. Shestakova, "Wide-Band Reflectors Based on Silver Films," *Soviet Journal of Optical Technology* **50**:781–782 (1983).
537. D. -Y. Song, R. W. Sprague, H. A. Macleod, and M. R. Jacobson, "Progress in the Development of a Durable Silver-Based High-Reflectance Coating for Astronomical Telescopes," *Applied Optics* **24**:1164–1170 (1985).
538. AIRCO Coating Technology, "Enhanced Front-Surface Aluminum Mirrors," 2700 Maxwell Way, P. O. Box 2529, Fairfield, CA 94533-252, USA, 1990.
539. S. A. Furman and E. G. Stolov, "Synthesis of Achromatic Metal-Dielectric Mirrors," *Soviet Journal of Optical Technology* **44**:359–361 (1977).
540. M. R. Adraens and B. Feuerbacher, "Improved LiF and MgF₂-Overcoated Aluminum Mirrors for Vacuum Ultraviolet Astronomy," *Applied Optics* **10**:958–959 (1971).
541. L. R. Canfield, G. Hass, and J. E. Waylonis, "Further Studies on MgF₂-Overcoated Aluminum Mirrors with Highest Reflectance in the Vacuum Ultraviolet," *Applied Optics* **5**:45–50 (1966).
542. G. Hass and R. Tousey, "Reflecting Coatings for the Extreme Ultraviolet," *Journal of the Optical Society of America* **49**:593–602 (1959).
543. D. W. Angel, W. R. Hunter, and R. Tousey, "Extreme Ultraviolet Reflectance of LiF-Coated Aluminum Mirrors," *Journal of the Optical Society of America* **51**:913–914 (1961).
544. J. T. Cox, G. Hass, and J. E. Waylonis, "Further Studies on LiF-Overcoated Aluminum Mirrors with Highest Reflectance in the Vacuum Ultraviolet," *Applied Optics* **7**:1535–1539 (1968).
545. D. K. Burge, H. E. Bennett, and E. J. Ashley, "Effect of Atmospheric Exposure on the Infrared Reflectance of Silvered Mirrors with and without Protective Coatings," *Applied Optics* **12**:42–47 (1973).
546. G. Hass, J. B. Heaney, J. F. Osantowski, and J. J. Triolo, "Reflectance and Durability of Ag Mirrors Coated with Thin Layers of Al₂O₃ Plus Reactively Deposited Silicon Oxide," *Applied Optics* **14**:2639–2644 (1975).
547. E. A. Volgunova, G. I. Golubeva, and A. M. Klochkov, "Highly Stable Silver Mirrors," *Soviet Journal of Optical Technology* **50**:128–129 (1983).
548. J. T. Cox, G. Hass, and W. R. Hunter, "Infrared Reflectance of Silicon Oxide and Magnesium Fluoride Protected Aluminum Mirrors at Various Angles of Incidence from 8 μm to 12 μm," *Applied Optics* **14**:1247–1250 (1975).
549. J. T. Cox and G. Hass, "Aluminum Mirrors Al₂O₃-Protected with High Reflectance at Normal but Greatly Decreased Reflectance at Higher Angles of Incidence in the 8-12-μm Region," *Applied Optics* **17**:333–334 (1978).

550. W. R. Hunter, J. F. Osantowski, and G. Hass, "Reflectance of Aluminum Overcoated with MgF_2 and LiF in the Wavelength Region from 1600 Å to 300 Å at Various Angles of Incidence," *Applied Optics* **10**:540–544 (1971).
551. S. F. Pellicori, "Infrared Reflectance of a Variety of Mirrors at 45° Incidence," *Applied Optics* **17**:3335–3336 (1978).
552. L. Young, "Multilayer Reflection Coatings on a Metal Mirror," *Applied Optics* **2**:445–447 (1963).
553. M. Z. Leś, F. Leś, and L. Gabla, "Semitransparent Metallic-Dielectric Mirrors with Low Absorption Coefficient in the Ultra-Violet Region of the Spectrum," *Acta Physica Polonica* **23**:211–214 (1963).
554. F. Leś and M. Z. Leś, "Metallic-Dielectric Mirrors with High Reflectivity in the Near Ultra-Violet for the Fabry-Perot Interferometer," *Acta Physica Polonica* **21**:523–528 (1962).
555. G. Hass, J. B. Ramsey, and W. R. Hunter, "Reflectance of Semitransparent Platinum Films on Various Substrates in the Vacuum Ultraviolet," *Applied Optics* **8**:2255–2259 (1969).
556. R. P. Madden, L. R. Canfield, and G. Hass, "On the Vacuum-Ultraviolet Reflectance of Evaporated Aluminum Before and During Oxidation," *Journal of the Optical Society of America* **53**:620–625 (1969).
557. G. Hass and W. R. Hunter, "Calculated Reflectance of Aluminum-Overcoated Iridium in the Vacuum Ultraviolet from 500 Å to 2000 Å," *Applied Optics* **6**:2097–2100 (1967).
558. L. N. Hadley and D. M. Dennison, "Reflection and Transmission Interference Filters. Part I. Theory," *Journal of the Optical Society of America* **37**:451–465 (1947).
559. L. N. Hadley and D. M. Dennison, "Reflection and Transmission Interference Filters. Part II. Experimental, Comparison with Theory, Results," *Journal of the Optical Society of America* **38**:483–496 (1948).
560. S. -y. Zheng and J. W. Y. Lit, "Design of a Narrow-Band Reflection IR Multilayer," *Canadian Journal of Physics* **61**:361–368 (1983).
561. L. A. Stelmack, "Vacuum Ultraviolet Reflectance Filter," U. S. patent 4408825, issued 11 October, 1983.
562. A. F. Turner and H. R. Hopkinson, "Reflection filters for the Visible and Ultraviolet," *Journal of the Optical Society of America* **43**:819 (1953).
563. R. Gamble and P. H. Lissberger, "Reflection Filter Multilayers of Metallic and Dielectric Thin Films," *Applied Optics* **28**:2838–2846 (1989).
564. P. H. Berning, G. Hass, and R. P. Maden, "Reflectance-Increasing Coatings for the Vacuum Ultraviolet and their Applications," *Journal of the Optical Society of America* **50**:586–597 (1960).
565. A. P. Lukirskii, E. P. Savinov, O. A. Ershov, and Y. F. Shepelev, "Reflection Coefficients of Radiation on the Wavelength Range from 23.6 to 113 Å for a Number of Elements and Substances and the Determination of the Refractive Index and Absorption Coefficient," *Optics and Spectroscopy (USSR)* **16**:168–172 (1964).
566. A. P. Lukirskii, E. P. Savinov, O. A. Ershov, I. I. Zhukova, and V. A. Forichev, "Reflection of X-Rays with Wavelengths from 23.6 to 190.3 Å. Some Remarks on the Performance of Diffraction Gratings," *Optics and Spectroscopy (USSR)* **19**:237–241 (1965).
567. J. Hecht, "A Novel approach to High XUV Reflectivity," *Lasers and Optronics* **April**:14 (1989).
568. K. Tamura, C. Sakai, N. Yamada, Y. Ogasaka, and R. Shibata, "Development of X-Ray Mirrors for X-Ray Telescopes," *Proc. SPIE* **5900**:051–059 (2005).
569. D. L. Windt, W. C. Cash, Jr., M. Scott, P. Arendt, B. Newnam, R. F. Fisher, A. B. Swartzlander, P. Z. Takacs, and J. M. Pinneo, "Optical Constants for Thin Films of C. Diamond, Al, Si and CVD SiC from 24 Å to 1216 Å," *Applied Optics* **27**:279–295 (1988).
570. R. F. Malina and W. Cash, "Extreme Ultraviolet Reflection Efficiencies of Diamond-Turned Aluminum, Polished Nickel, and Evaporated Gold Surfaces," *Applied Optics* **17**:3309–3313 (1978).
571. R. L. Sandberg, D. D. Allred, J. E. Johnson, and R. S. Turley, "A Comparison of Uranium Oxide and Nickel as Single-layer Reflectors from 2.7 to 11.6 Nanometers," *Proc. SPIE* **5193**:191–203 (2004).
572. D. Spiga, G. Pareschi, O. Citterio, R. Banham, S. Basso, M. Cassanelli, V. Cotroneo, et al., "Development of Multilayer Coatings (Ni/C-Pt/C) for Hard X-Ray Telescopes by E-Beam Evaporation with Ion Assistance," *Proc. SPIE* **5488**: 813–819 (2004).
573. V. V. Protopopov and V. A. Kalnov, "X-Ray Multilayer Mirrors with an Extended Angular Range," *Optics Communications* **158**:1–6 (1998).
574. O. Citterio, P. Cerutti, F. Mazzoleni, G. Pareschi, E. Poretti, P. Lagana, A. Mengali, C. Misiano, F. Pozzilli, and E. Simonetti, "Multilayer Optics for Hard X-Ray Astronomy by Means of Replication Techniques," *Proc. SPIE* **3766**:310–319 (1999).

575. V. V. Protopopov, A. V. Tikhonravov, A. V. Voronov, M. K. Trubetskov, and G. W. DeBell, "Optimal Design of Graded X-Ray Multilayer Mirrors in the Angular and Spectral Domains," *Proc. SPIE* **3766**:320–326 (1999).
576. Z. Wang, "Multilayers for the EUV and Soft X-Ray Region," *Proc. SPIE* **5963**:0S 1–12 (2005).
577. K. Yamashita, H. Kunieda, Y. Tawara, K. Tamura, Y. Ogasaka, K. Haga, T. Okajima, et al. "New Design Concept of Multilayer Supermirrors for Hard X-Ray Optics," *Proc. SPIE* **3766**:327–335 (1999).
578. E. van Rooyen and E. Theron, "A Multiple Reflection Multilayer Reflection Filter," *Applied Optics* **8**:832–833 (1969).
579. A. S. Valeyev, "Multilayer Reflection-Type Interference Filters," *Soviet Journal of Optical Technology* **34**:317–319 (1967).
580. Schott and Gen., "Interference Reflection Filter UV-R-250," Geschäftsbereich Optik, Hattenbergstrasse 10, D-6500 Mainz, Germany, 1967.
581. G. Hass, J. B. Heaney, and W. R. Hunter, "Reflectance and Preparation of Front Surface Mirrors for Use at Various Angles of Incidence from the Ultraviolet to the Far Infrared," in *Physics of Thin Films* 12, Academic Press, 1982, pp. 1–51.
582. D. W. Lynch, "Mirror and Reflector Materials," in *Handbook of Laser Science and Technology. Optical Materials: Part 2 IV*, (ed. M. J. Weber), CRC Press, Inc., Boca Raton, Florida, 1986, pp. 185–219.
583. B. T. Sullivan and J. A. Dobrowolski, "Deposition of Optical Multilayer Coatings with Automatic Error Compensation: II. Magnetron Sputtering," *Applied Optics* **32**:2351–2360 (1993).

FUNDAMENTAL OPTICAL PROPERTIES OF SOLIDS

Alan Miller

*Scottish Universities Physics Alliance
School of Engineering and Physical Sciences
Heriot-Watt University
Edinburgh, Scotland*

8.1 GLOSSARY

$\underline{\mathbf{A}}$	vector field
$\underline{\mathbf{B}}$	magnetic induction
c	speed of light
$\underline{\mathbf{D}}$	displacement field
\underline{d}_M	penetration depth of material
d_{TIR}	evanescent field depth for total internal reflection
$\underline{\mathbf{E}}$	electric field
\mathcal{E}_g	band gap energy
$\underline{\mathbf{E}}_{\text{loc}}$	local electric field
$-e$	electronic charge
e_i^ρ	components of unit vectors of incident polarization
e_s^ρ	components of unit vectors of scattered polarization
f_j	oscillator strength
$\underline{\mathbf{G}}$	reciprocal lattice vector
$\underline{\mathbf{H}}$	hamiltonian
$\underline{\mathbf{H}}$	magnetic field
I	irradiance
$\underline{\mathbf{I}}$	unit tensor
$\underline{\mathbf{J}}_c$	conduction current density
$\underline{\mathbf{K}}$	dielectric constant or relative permittivity
$\underline{\mathbf{k}}$	electron wavevector
$\underline{\mathbf{M}}$	induced magnetization
m	electron mass
m^*	effective electron mass

N	charge density
N_L	refractive index for left circularly polarized light
N_R	refractive index for right circularly polarized light
n	real part of refractive index
P	Kane momentum matrix element
\underline{P}	induced dipole moment per unit volume
P_D	dielectric amplitude reflection polarization ratio
P_{ji}	momentum matrix element
P_M	metallic amplitude reflection polarization ratio
$P_{p\sigma, ij}$	elasto-optic coefficients
\mathbf{p}	electron momentum
$\tilde{\mathbf{q}}$	complex photon wavevector
$\hat{\mathbf{q}}$	unit photon wavevector
R	reflectance
R_{ex}	exciton Rydberg
\mathbf{P}	Raman scattering coefficient
$r_{\sigma p}$	field reflection amplitude for p-polarization
r_p	field reflection amplitude for s-polarization
r_s	field reflection amplitude for s-polarization
\underline{S}	Poynting vector
s_{ij}	components of lattice strain
S_R	Raman scattering efficiency
S_j	phonon oscillator strength
T	transmittance
t_p	field transmission amplitude for p-polarization
t_s	field transmission amplitude for s-polarization
V	volume of unit cell
$\mathbf{P}(x)$	periodic lattice potential
v	phase velocity
$W(\omega)$	transition rate
Z	number of electrons per atom
α	absorption coefficient
β	two-photon absorption coefficient
Γ	damping constant
Δ_D	dielectric reflection polarization phase difference
Δ_M	metallic reflection polarization phase difference
Δ_S	specific rotary power
Δ_{SO}	spin-orbit splitting
δ	energy walk-off angle
$\underline{\epsilon}$	electrical permittivity
$\tilde{\epsilon}$	complex permittivity
$\epsilon(\infty)$	high frequency permittivity
$\epsilon(0)$	static permittivity
ϵ_1	real part of permittivity
ϵ_2	imaginary part of permittivity
ϵ_0	electrical permittivity of vacuum

$\tilde{\eta}$	complex refractive index
$\bar{\theta}$	principal angle for metallic reflection
θ_{OA}	angle between optic axis and principal axis
θ_B	Brewster's angle
θ_C	total internal reflection critical angle
θ_{TIR}	phase change for total internal reflection
κ	imaginary part of refractive index
$\underline{\underline{\mu}}$	magnetic permeability
μ_o	magnetic permeability of vacuum
ξ	Hooke's law constant
$\hat{\xi}$	unit polarization vector
ρ_f	free charge density
ρ_p	band or polarization charge density
ρ_p	real field amplitude for p-polarization
ρ_s	real field amplitude for s-polarization
ρ_t	total charge density
$\tilde{\sigma}$	complex conductivity
$\underline{\underline{\sigma}}$	electrical conductivity
σ_1	real part of conductivity
σ_2	imaginary part of conductivity
τ	scattering or relaxation time
ϕ	scalar field
ϕ_r	ground state wavefunction
$\tilde{\chi}$	complex susceptibility
χ'	real part of susceptibility
χ''	imaginary part of susceptibility
$\underline{\underline{\chi}}_e$	electrical susceptibility
$\underline{\underline{\chi}}_m$	magnetic susceptibility
Ψ_k	Bloch solution wavefunction
ω	photon frequency
ω_j	resonant oscillation frequency
ω_{LO}	longitudinal optical phonon frequency
ω_p	plasma frequency
ω_s	surface frequency
ω_{TO}	transverse optical phonon frequency

8.2 INTRODUCTION

This chapter describes the fundamental interactions of light with solids. The discussion is restricted to *intrinsic* optical properties. *Extrinsic* optical properties of solids are described in Chap. 10, "Optical Spectroscopy and Spectroscopic Lineshapes," by Brian Henderson in Vol. I. Basic formulas, definitions, and concepts are listed for reference and as a foundation for subsequent chapters of this *Handbook*. More detailed accounts of specific optical properties and particular types of solids are given in later chapters, that is, Chap. 7, "Optical Properties of Films and Coatings,"

by Jerzy A. Dobrowolski in this volume and Chap. 7, “Electro-Optic Modulators” by Georgeanne M. Purvinis and Theresa A. Maldonado in Vol. V. The reader is referred to the many texts which provide more elaborate discussions of the optical properties of solids.¹⁻¹³

Electrical measurements distinguish three general types of solid by their conductivities, i.e., dielectrics, semiconductors, and metals. Optically, these three groups are characterized by different fundamental bandgap energies, \mathcal{E}_g . Although the boundaries are not sharply defined, an approximate distinction is given by metals, $\mathcal{E}_g \leq 0$, semiconductors, $0 < \mathcal{E}_g < 3$ eV, and dielectrics, $\mathcal{E}_g > 3$ eV. Solids may be found in single crystal, polycrystalline, and amorphous forms. Rudimentary theories of the optical properties of condensed matter are based on light interactions with perfect crystal lattices characterized by extended (nonlocal) electronic and vibrational energy states. These eigenstates are determined by the periodicity and symmetry of the lattice and the form of the Coulomb potential which arises from the interatomic bonding.

The principal absorption bands in condensed matter occur at photon energies corresponding to the frequencies of the lattice vibrations (phonons) in the infrared, and electronic transitions in the near infrared, visible, or ultraviolet. A quantum mechanical approach is generally required to describe electronic interactions but classical models often suffice for lattice vibrations. Although the mechanical properties of solids can vary enormously between single crystal and polycrystalline forms, the *fundamental* optical properties are similar, even if the crystallite size is smaller than a wavelength, since the optical interaction is microscopic. However, electronic energy levels and hence optical properties are fundamentally altered when one or more dimensions of a solid are reduced to the scale of the de-Broglie wavelength of the electrons. Modern crystal growth techniques allow fabrication of atomic precision epitaxial layers of different solid materials. Ultrathin layers with dimensions comparable with or smaller than the de-Broglie wavelength of an electron may form quantum wells, quantum wires, and quantum dots in which electronic energy levels are quantized. Amorphous solids have random atomic or molecular orientation on the distance scale of several nearest neighbors, but generally have well-defined bonding and local atomic order which determine the overall optical response.

8.3 PROPAGATION OF LIGHT IN SOLIDS

Dielectrics and semiconductors provide transparent spectral regions at radiation frequencies between the phonon vibration bands and the fundamental (electronic) absorption edge. Maxwell's equations successfully describe the propagation, reflection, refraction, and scattering of harmonic electromagnetic waves.

Maxwell's Equations

Four equations relate the macroscopically averaged electric field \mathbf{E} and magnetic induction \mathbf{B} , to the total charge density, ρ_t , (sum of the bound or polarization charge, ρ_p and the free charge, ρ_f), the conduction current density, \mathbf{J}_c , the induced dipole moment per unit volume, \mathbf{P} , and the induced magnetization of the medium, \mathbf{M} , (expressed in SI units),

$$\nabla \cdot \mathbf{E} = \frac{\rho_t}{\epsilon_0} \quad (1)$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t} \quad (2)$$

$$\nabla \cdot \mathbf{B} = 0 \quad (3)$$

$$\nabla \times \mathbf{B} = \mu_0 \left(\epsilon_0 \frac{\partial \mathbf{E}}{\partial t} + \frac{\partial \mathbf{P}}{\partial t} + \mathbf{J}_c + \nabla \times \mathbf{M} \right) \quad (4)$$

where $\epsilon_0 = 8.854 \times 10^{-12}$ F/m is the permittivity of vacuum, $\mu_0 = 4\pi \times 10^{-7}$ H/m is the permeability of vacuum and c the speed of light in vacuum, $c = (\epsilon_0 \mu_0)^{-1/2}$. By defining a displacement field, \mathbf{D} , and magnetic field, \mathbf{H} , to account for the response of a medium

$$\mathbf{D} = \mathbf{P} + \epsilon_0 \mathbf{E} \quad (5)$$

$$\mathbf{H} = \frac{\mathbf{B}}{\mu_0} - \mathbf{M} \quad (6)$$

and using the relation between polarization and bound charge density,

$$-\nabla \cdot \mathbf{P} = \rho_p \quad (7)$$

Equations (1) and (4) may also be written in the form

$$\nabla \cdot \mathbf{D} = \rho_f \quad (8)$$

and

$$\nabla \times \mathbf{H} = \frac{\partial \mathbf{D}}{\partial t} + \mathbf{J}_c \quad (9)$$

Vector \mathbf{A} and scalar ϕ fields may be defined by

$$\mathbf{B} = \nabla \times \mathbf{A} \quad (10)$$

$$\mathbf{E} = -\nabla \phi - \frac{\partial \mathbf{A}}{\partial t} \quad (11)$$

A convenient choice of gauge for the optical properties of solids is the Coulomb (or transverse) gauge

$$\nabla \cdot \mathbf{A} = 0 \quad (12)$$

which ensures that the vector potential \mathbf{A} is transverse for plane electromagnetic waves, while the scalar potential represents any longitudinal current and satisfies Poisson's equation

$$\nabla^2 \phi = -\frac{\rho_f}{\epsilon_0} \quad (13)$$

Three constitutive relations describe the response of conduction and bound electrons to the electric and magnetic fields

$$\mathbf{J}_c = \underline{\underline{\sigma}} \mathbf{E} \quad (14)$$

$$\mathbf{D} = \mathbf{P} + \epsilon_0 \mathbf{E} = \underline{\underline{\epsilon}} \mathbf{E} \quad (15)$$

$$\mathbf{B} = \mu_0 (\mathbf{H} + \mathbf{M}) = \underline{\underline{\mu}} \mathbf{H} \quad (16)$$

where $\underline{\underline{\sigma}}$ is the electrical conductivity, $\underline{\underline{\epsilon}}$ is the electrical permittivity, and $\underline{\underline{\mu}}$ is the magnetic permeability of the medium and are in general tensor quantities which may depend on field strengths. An alternative relation often used to define a dielectric constant (or relative permittivity), $\underline{\underline{\mathbf{K}}}$ is given by, $\mathbf{D} = \epsilon_0 \underline{\underline{\mathbf{K}}} \mathbf{E}$. In isotropic media using the approximation of linear responses to electric and magnetic fields, $\underline{\underline{\sigma}}$, $\underline{\underline{\epsilon}}$, and $\underline{\underline{\mu}}$ are constant scalar quantities.

Electric, $\underline{\chi}_e$, and magnetic, $\underline{\chi}_m$, susceptibilities may be defined to relate the induced dipole moment, \mathbf{P} , and magnetism, \mathbf{M} to the field strengths \mathbf{E} and \mathbf{H} ,

$$\mathbf{P} = \epsilon_0 \underline{\chi}_e \mathbf{E} \quad (17)$$

$$\mathbf{M} = \epsilon_0 \underline{\chi}_m \mathbf{H} \quad (18)$$

Thus,

$$\underline{\epsilon} = \underline{\epsilon}_0 (\mathbf{I} + \underline{\chi}_e) \quad (19)$$

and

$$\underline{\mu} = \mu_0 (\mathbf{I} + \underline{\chi}_m) \quad (20)$$

where \mathbf{I} is the unit tensor.

Wave Equations and Optical Constants

The general wave equation derived from Maxwell's equations is

$$\nabla^2 \mathbf{E} - \nabla(\nabla \cdot \mathbf{E}) - \epsilon_0 \mu_0 \frac{\partial^2 \mathbf{E}}{\partial t^2} = \mu_0 \left(\frac{\partial^2 \mathbf{P}}{\partial t^2} + \frac{\partial \mathbf{J}_c}{\partial t^2} + \nabla \times \frac{\partial \mathbf{M}}{\partial t} \right) \quad (21)$$

For dielectric, (nonconducting) solids

$$\nabla^2 \mathbf{E} = \mu \epsilon \frac{\partial^2 \mathbf{E}}{\partial t^2} \quad (22)$$

The harmonic plane wave solution of the wave equation for monochromatic light at frequency, ω ,

$$\mathbf{E} = \frac{1}{2} \mathbf{E}_0 \exp i(\mathbf{q} \cdot \mathbf{r} - \omega t) + c.c. \quad (23)$$

in homogeneous ($\nabla \epsilon = 0$), isotropic ($\nabla \cdot \mathbf{E} = 0$, $\mathbf{q} \cdot \mathbf{E} = 0$), nonmagnetic ($\mathbf{M} = 0$) solids results in a complex wavevector, $\tilde{\mathbf{q}}$,

$$\tilde{\mathbf{q}} = \frac{\omega}{c} \sqrt{\frac{\epsilon}{\epsilon_0} + i \frac{\sigma}{\epsilon_0 \omega}} \quad (24)$$

A complex refractive index, $\tilde{\eta}$, may be defined by

$$\tilde{\mathbf{q}} = \frac{\omega}{c} \tilde{\eta} \hat{\mathbf{q}} \quad (25)$$

where $\hat{\mathbf{q}}$ is a unit vector and

$$\tilde{\eta} = n + i\kappa \quad (26)$$

Introducing complex notation for the permittivity, $\tilde{\epsilon} = \epsilon_1 + i\epsilon_2$, conductivity, $\tilde{\sigma} = \sigma_1 + i\sigma_2$, and susceptibility, $\tilde{\chi}_e = \chi'_e + i\chi''_e$, we may relate

$$\epsilon_1 = \epsilon = \epsilon_0 (1 + \chi'_e) = \epsilon_0 (n^2 - \kappa^2) = -\frac{\sigma_2}{\omega} \quad (27)$$

and

$$\varepsilon_2 = \varepsilon_o \chi_e'' = 2\varepsilon_o n \kappa = \frac{\sigma}{\omega} = \frac{\sigma_1}{\omega} \quad (28)$$

Alternatively,

$$n = \left[\frac{1}{2\varepsilon_o} \left(\sqrt{\varepsilon_1^2 + \varepsilon_2^2} + \varepsilon_1 \right) \right]^{1/2} \quad (29)$$

$$\kappa = \left[\frac{1}{2\varepsilon_o} \left(\sqrt{\varepsilon_1^2 + \varepsilon_2^2} - \varepsilon_1 \right) \right]^{1/2} \quad (30)$$

The field will be modified locally by the induced dipoles. If there is no free charge, ρ_f , the local field, \mathbf{E}_{loc} , may be related to the external field \mathbf{E}_i in isotropic solids using the Clausius-Mossotti equation which leads to the relation

$$\mathbf{E}_{\text{loc}} = \frac{n^2 + 2}{3} \mathbf{E}_i \quad (31)$$

Energy Flow

The direction and rate of flow of electromagnetic energy is described by the Poynting vector

$$\mathbf{S} = \frac{1}{\mu_o} \mathbf{E} \times \mathbf{H} \quad (32)$$

The average power per unit area, (irradiance, I), W/m^2 , carried by a uniform plane wave is given by the time averaged magnitude of the Poynting vector

$$I = \langle \mathbf{S} \rangle = \frac{cn |\mathbf{E}_o|^2}{2} \quad (33)$$

The plane wave field in Eq. (23) may be rewritten for absorbing media using Eqs. (25) and (26)

$$\mathbf{E}(\mathbf{r}, t) = \frac{1}{2} \mathbf{E}_o(\mathbf{q}, \omega) \exp\left(-\frac{\omega}{c} \kappa \hat{\mathbf{q}} \cdot \mathbf{r}\right) \exp\left[i\left(\frac{\omega}{c} \mathbf{n} \hat{\mathbf{q}} \cdot \mathbf{r} - \omega t\right)\right] + c.c. \quad (34)$$

The decay of the propagating wave is characterized by the extinction coefficient κ . The attenuation of the wave may also be described by Beer's law

$$I = I_o \exp(-\alpha z) \quad (35)$$

where α is the absorption coefficient describing the attenuation of the irradiance, I , with distance, z . Thus,

$$\alpha = \frac{2\omega\kappa}{c} = \frac{4\pi\kappa}{\lambda} = \frac{\sigma_1}{\varepsilon_o cn} = \frac{\varepsilon_2 \omega}{\varepsilon_o cn} = \frac{\chi_e'' \omega}{cn} \quad (36)$$

The power absorbed per unit volume is given by

$$P_{\text{abs}} = \alpha I = \frac{\omega \chi''}{2} |\mathbf{E}_0|^2 \quad (37)$$

The second exponential in Eq. (34) is oscillatory and represents the phase velocity of the wave, $v = c/n$.

Anisotropic Crystals

Only amorphous solids and crystals possessing cubic symmetry are optically isotropic. In general, the speed of propagation of an electromagnetic wave in a crystal depends both on the direction of propagation and on the polarization of the light wave. The linear electric susceptibility and dielectric constant may be represented by tensors with components of $\underline{\chi}_e$ given by

$$P_i = \epsilon_0 \chi_{ij} E_j \quad (38)$$

where i and j refer to coordinate axes. In an anisotropic crystal, $\mathbf{D} \perp \mathbf{B} \perp \mathbf{q}$ and $\mathbf{E} \perp \mathbf{H} \perp \mathbf{S}$, but \mathbf{E} is not necessarily parallel to \mathbf{D} and the direction of energy flow \mathbf{S} is not necessarily in the same direction as the propagation direction \mathbf{q} .

From energy arguments it can be shown that the susceptibility tensor is symmetric and it therefore follows that there always exists a set of coordinate axes which diagonalize the tensor. This coordinate system defines the principal axes. The number of nonzero elements for the susceptibility (or dielectric constant) is thus reduced to a maximum of three (for any crystal system at a given wavelength). Thus, the dielectric tensor defined by the direction of the electric field vector with respect to the principal axes has the form

$$\begin{bmatrix} \epsilon_1 & 0 & 0 \\ 0 & \epsilon_2 & 0 \\ 0 & 0 & \epsilon_3 \end{bmatrix}$$

The principal indices of refraction are

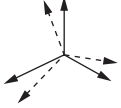
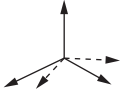

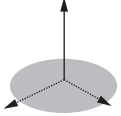
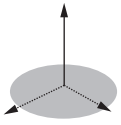
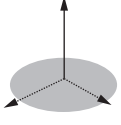
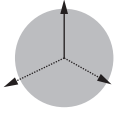
$$n_i = \sqrt{1 + \chi_{ii}} = \sqrt{\epsilon_i} \quad (39)$$

with the \mathbf{E} -vector polarized along any principal axis, i.e., $\mathbf{E} \parallel \mathbf{D}$. This case is designated as an ordinary or \mathbf{o} -ray in which the phase velocity is independent of propagation direction. An extraordinary or \mathbf{e} -ray occurs when both \mathbf{E} and \mathbf{q} lie in a plane containing two principal axes with different n . An *optic axis* is defined by any propagation direction in which the phase velocity of the wave is independent of polarization.

Crystalline solids fall into three classes: (a) optically isotropic, (b) uniaxial, or (c) biaxial (see Table 1). All choices of orthogonal axes are principal axes and $\epsilon_1 = \epsilon_2 = \epsilon_3$ in *isotropic* solids. For a *uniaxial* crystal, $\epsilon_1 = \epsilon_2 \neq \epsilon_3$, a single optic axis exists for propagation in direction 3. In this case, the ordinary refractive index, $n_o = n_1 = n_2$, is independent of the direction of polarization in the 1-2 plane. Any two orthogonal directions in this plane can be chosen as principal axes. For any other propagation direction, the polarization can be divided into an \mathbf{o} -ray component in the 1-2 plane and a perpendicular \mathbf{e} -ray component (see Fig. 1). The dependence of the \mathbf{e} -ray refractive index with propagation direction is given by the ellipsoid,

$$n_j(\theta_i) = \frac{n_i n_j}{(n_i^2 \cos^2 \theta_i + n_j^2 \sin^2 \theta_i)^{1/2}} \quad (40)$$

TABLE 1 Crystallographic Point Groups and Optical Properties

System		Point Group	Symbols	Optical Activity	
		International	Schönflies		
Triclinic	biaxial	1	C ₁	A	
		$\bar{1}$	S ₂	-	
Monoclinic	biaxial	2	C ₂	A	
		m	C _v	-	
		2/m	C _{2h}	-	
Orthorhombic	biaxial	mm	C _{2v}	-	
		222	D ₂	A	
		mmm	D _{2h}	-	
Trigonal	uniaxial	$\frac{3}{2}$	C ₃	A	
		$\bar{3}$	S ₆	-	
		3m	C _{3v}	-	
		$\frac{32}{3}$	D ₃	A	
		$\bar{3}m$	D _{3d}	-	
Tetragonal	uniaxial	$\frac{4}{2}$	C ₄	A	
		$\frac{4}{4}$	S ₄	-	
		4/m	C _{4h}	-	
		4mm	C _{4v}	-	
		$\frac{42m}{2}$	D _{2d}	A	
		42	D ₄	A	
4/mmm	D _{4h}	-			
Hexagonal	uniaxial	$\frac{6}{2}$	C ₆	A	
		$\frac{6}{6}$	C _{3h}	-	
		6/m	C _{6h}	-	
		6mm	C _{6v}	-	
		$\frac{6m2}{2}$	D _{3h}	-	
		62	D ₆	A	
6/mmm	D _{6h}	-			
Cubic	isotropic	23	T	A	
		$\frac{m3}{2}$	T _h	-	
		$\frac{43m}{2}$	T _d	-	
		432	O	A	
		m3m	O _h	-	

where θ_j is defined with respect to optic axis, $i = 3$, and $j = 1$ or 2 . $\theta_j = 90^\circ$ gives the refractive index $n_e = n_3$ when the light is polarized along axis **3**. The difference between n_o and n_e is the birefringence. Figure 1a illustrates the case of positive birefringence, $n_e > n_o$ and Fig. 1b negative birefringence, $n_e < n_o$. The energy walk-off angle, δ , (the angle between **S** and **q** or **D** and **E**) is given by

$$\tan \delta = \frac{n^2(\theta)}{2} \left[\frac{1}{n_3^2} - \frac{1}{n_1^2} \right] \sin 2\theta \quad (41)$$

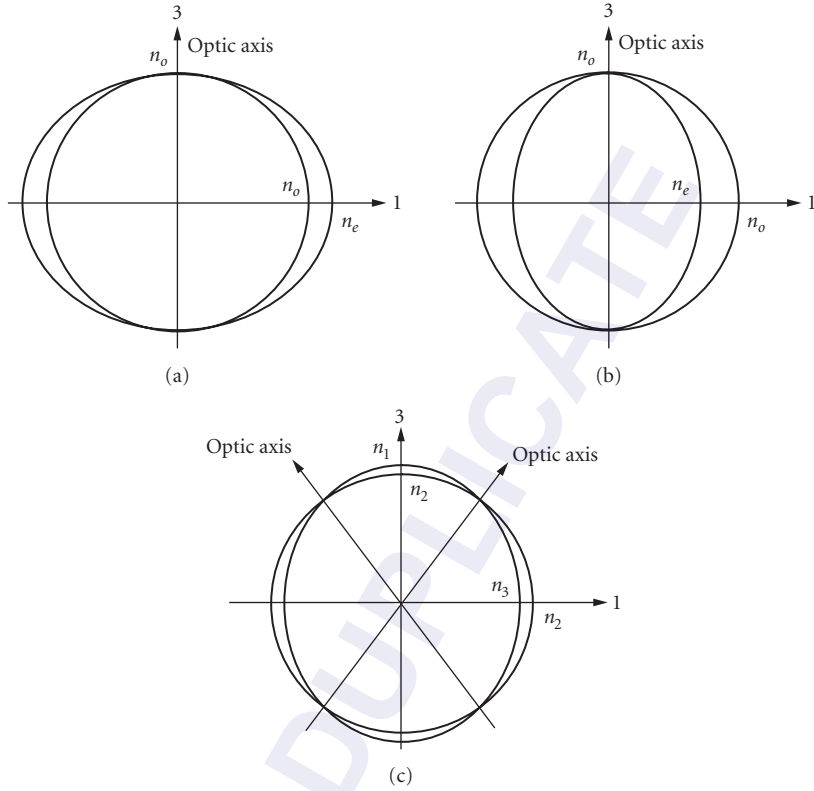


FIGURE 1 Illustration of directional dependence of refractive indices and optic axes in (a) a uniaxial, positive birefringent crystal, (b) a uniaxial, negative birefringent crystal, and (c) biaxial crystal.

In *biaxial* crystals, diagonalization of the dielectric tensor results in three independent coefficients, $\epsilon_1 \neq \epsilon_2 \neq \epsilon_3 \neq \epsilon_1$. For orthorhombic crystals, a single set of orthogonal principal axes is fixed for all wavelengths. However, in monoclinic structures only one principal axis is fixed. The direction of the other two axes rotates in the plane perpendicular to the fixed axis as the wavelength changes (retaining orthogonality). In triclinic crystals there are no fixed axes and orientation of the set of three principal axes varies with wavelength. Equation (40) provides the *e*-ray refractive index within planes containing two principal axes. Biaxial crystals possess two optic axes. Defining principal axes such that $n_1 > n_2 > n_3$, both optic axes lie in the 1-3 plane, at an angle θ_{OA} from axis 1, as illustrated in Fig. 1c, where

$$\sin \theta_{OA} = \pm \frac{n_1}{n_2} \sqrt{\frac{n_2^2 - n_3^2}{n_1^2 - n_3^2}} \quad (42)$$

Crystals with certain point group symmetries (see Table 1) also exhibit optical activity, i.e., the ability to rotate the plane of linearly polarized light. An origin for this phenomenon, is the weak magnetic interaction $\nabla \times \mathbf{M}$ [see Eq. (21)], when it applies in a direction perpendicular to \mathbf{P} (i.e. $\mathbf{M} \parallel \mathbf{P}$). The specific rotary power Δ_S , (angle of rotation of linearly polarized light per unit length) is given by

$$\Delta_S = \frac{\pi}{\lambda} (n_L - n_R) \quad (43)$$

where n_L and n_R are refractive indices for left and right circular polarization. Optical activity is often masked by birefringence; however, polarization rotation can be observed in optically active materials when the propagation is along the optic axis or when the birefringence is coincidentally zero in other directions. In the case of propagation along the optic axis of an optically active uniaxial crystal such as quartz, the susceptibility tensor may be written

$$\begin{bmatrix} \chi_{11} & i\chi_{12} & 0 \\ -i\chi_{12} & \chi_{11} & 0 \\ 0 & 0 & \chi_{33} \end{bmatrix}$$

and the rotary power is proportional to the imaginary part of the magnetic susceptibility, $\chi_m'' = \chi_{12}$,

$$\Delta_S = \frac{\pi\chi_{12}}{n\lambda} \quad (44)$$

Crystals can exist in left- or right-handed versions. Other crystal symmetries, e.g., $\bar{4}2m$, can be optically active for propagation along the 2 and 3 axes, but rotation of the polarization is normally masked by the typically larger birefringence except at accidental degeneracies.

Interfaces

Applying boundary conditions at a plane interface between two media with different indices of refraction leads to the laws of reflection and refraction. Snell's law applies to all o-rays and relates the angle of incidence, θ_A in medium A and the angle of refraction, θ_B in medium B to the respective ordinary refractive indices n_A and n_B ,

$$n_A \sin\theta_A = n_B \sin\theta_B \quad (45)$$

Extraordinary rays do not satisfy Snell's law. The propagation direction for the e-ray can be found graphically by equating the projections of the propagation vectors in the two media along the boundary plane. Double refraction of unpolarized light occurs in anisotropic crystals.

The field amplitude ratios of reflected and transmitted rays to the incident ray (r and t) in isotropic solids (and o-rays in anisotropic crystals) are given by the Fresnel relations. For s- (σ or TE) polarization (\mathbf{E} -vector perpendicular to the plane of incidence) (see Fig. 2a) and p- (π or TM) polarization (\mathbf{E} -vector parallel to the plane of incidence) (see Fig. 2b):

$$r_s = \frac{E_{rs}}{E_{is}} = \frac{n_A \cos\theta_A - \sqrt{n_B^2 - n_A^2 \sin^2\theta_A}}{n_A \cos\theta_A + \sqrt{n_B^2 - n_A^2 \sin^2\theta_A}} \quad (46)$$

$$r_p = \frac{E_{rp}}{E_{ip}} = \frac{n_B^2 \cos\theta_A - n_A \sqrt{n_B^2 - n_A^2 \sin^2\theta_A}}{n_B^2 \cos\theta_A + n_A \sqrt{n_B^2 - n_A^2 \sin^2\theta_A}} \quad (47)$$

$$t_s = \frac{E_{ts}}{E_{is}} = \frac{2n_A \cos\theta_A}{n_A \cos\theta_A + \sqrt{n_B^2 - n_A^2 \sin^2\theta_A}} \quad (48)$$

$$t_p = \frac{E_{tp}}{E_{ip}} = \frac{2n_A n_B \cos\theta_A}{n_B^2 \cos\theta_A + n_A \sqrt{n_B^2 - n_A^2 \sin^2\theta_A}} \quad (49)$$

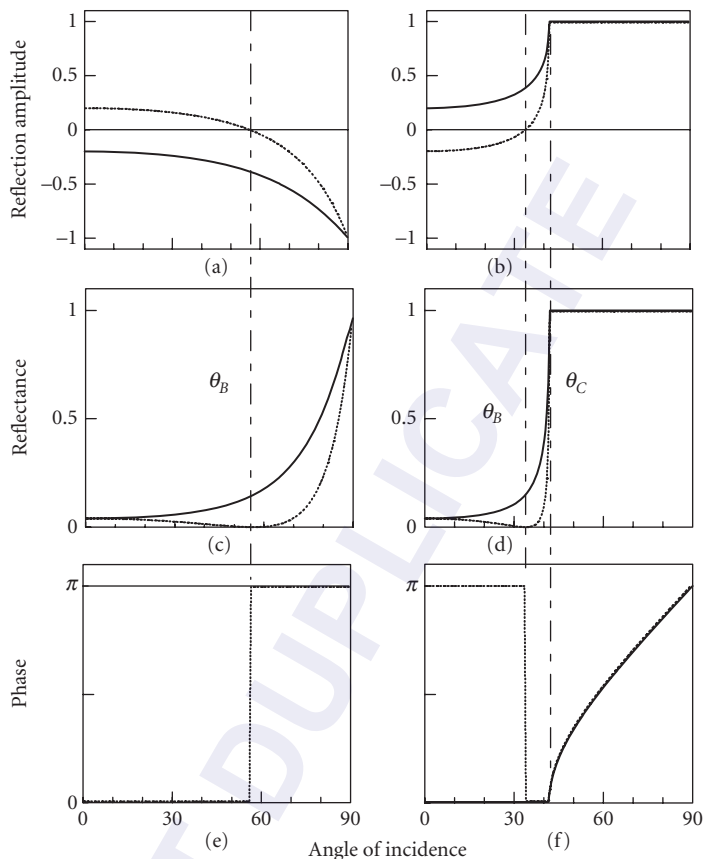


FIGURE 2 The electric field reflection amplitudes (*a, b*), energy reflectance (*c, d*), and phase change (*e, f*) for *s*- (solid lines) and *p*- (dashed lines) polarized light for external (*a, c, e*) and internal (*b, d, f*) reflection in the case $n_A = 1$, $n_B = 1.5$. θ_B is the polarizing or Brewster's angle and θ_C is the critical angle for total internal reflection.²

At normal incidence, the energy reflectance, R , (see Figs. 2*c* and *d*), and transmittance, T , are

$$R = \frac{|E_r|^2}{|E_i|^2} = \left| \frac{n_B - n_A}{n_B + n_A} \right|^2 \quad (50)$$

$$T = \frac{|E_t|^2}{|E_i|^2} = \frac{4n_A^2}{|n_A + n_B|^2} \quad (51)$$

The *p*-polarized reflectivity, Eq. (47), goes to zero at Brewster's angle under the condition

$$\theta_B = \tan^{-1} \left(\frac{n_A}{n_B} \right) \quad (52)$$

If $n_A > n_B$, total internal reflection (TIR) occurs when the angle of incidence exceeds a critical angle,

$$\theta_C = \sin^{-1} \left(\frac{n_B}{n_A} \right) \quad (53)$$

This critical angle may be different for s- and p-polarizations in anisotropic crystals. Under conditions of TIR, the evanescent wave amplitude drops to e^{-1} in a distance

$$d_{\text{TIR}} = \frac{c}{\omega} (n_A^2 \sin^2 \theta_A - n_B^2)^{-1/2} \quad (54)$$

The phase changes on reflection for external and internal reflection from an interface with $n_A = 1.5$, $n_B = 1$ are plotted in Fig. 2e and f. Except under TIR conditions, the phase change is either 0 or π . The complex values predicted by Eqs. (46) and (47) for angles of incidence greater than the critical angle for TIR imply phase changes in the reflected light which are neither 0 nor π . The phase of the reflected light changes by π at Brewster's angle in p-polarization. The ratio of s to p reflectance, P_D , is shown in Fig. 3a and the phase difference, $\Delta_D = \phi_p - \phi_s$ in Fig. 3b. Under conditions of TIR, the phase change on reflection, ϕ_{TIR} , is given by

$$\tan \frac{\phi_{\text{TIR}}}{2} = \frac{\sqrt{\sin^2 \theta_A - \sin^2 \theta_C}}{\cos \theta_A} \quad (55)$$

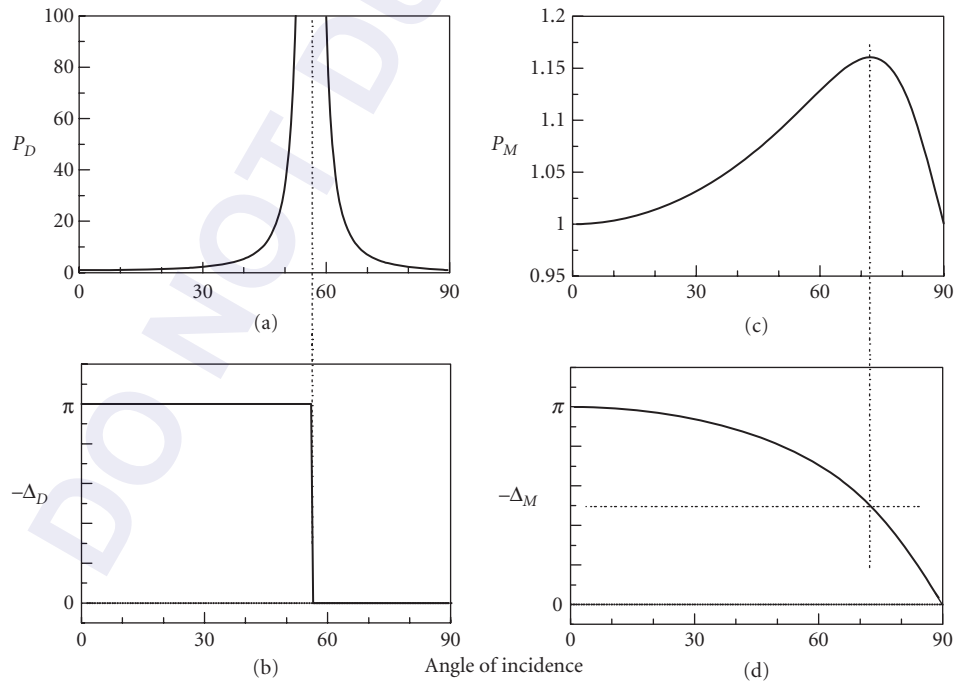


FIGURE 3 Typical polarization ratios, P , (a, c) and phase differences Δ , (b, d) for s- and p-polarizations at dielectric, D, and metallic, M, surfaces.¹

8.4 DISPERSION RELATIONS

For most purposes, a classical approach is found to provide a sufficient description of dispersion of the refractive index within the transmission window of insulators, and for optical interactions with lattice vibrations and free electrons. However, the details of interband transitions in semiconductors and insulators and the effect of d-levels in transition metals requires a quantum model of dispersion close to these resonances.

Classical Model

The Lorentz model for dispersion of the optical constants of solids assumes an optical interaction via the polarization produced by a set of damped harmonic oscillators. The polarization $\underline{\mathbf{P}}$ induced by a displacement \mathbf{r} of bound electrons of density N and charge $-e$ is

$$\mathbf{P} = -N e \mathbf{r} \quad (56)$$

Assuming the electrons to be elastically bound (Hooke's law) with a force constant, ξ ,

$$-e \mathbf{E}_{\text{loc}} = \xi \mathbf{r} \quad (57)$$

the differential equation of motion is

$$m \frac{d^2 \mathbf{r}}{dt^2} + m \Gamma \frac{d \mathbf{r}}{dt} + \xi \mathbf{r} = -e \mathbf{E}_{\text{loc}} \quad (58)$$

where m is the electron mass and Γ is a damping constant. Here the lattice is assumed to have infinite mass and the magnetic interaction has been neglected. Solving the equation of motion for fields of frequency ω gives a relation for the complex refractive index and dielectric constant

$$\tilde{\eta}^2 = \frac{\tilde{\epsilon}}{\epsilon_0} = 1 + \frac{N e^2}{m \epsilon_0} \sum_j \frac{f_j}{(\omega_j^2 - \omega^2 - i \Gamma_j \omega)} \quad (59)$$

We have given the more general result for a number of resonant frequencies

$$\omega_j = \sqrt{\frac{\xi_j}{m}} \quad (60)$$

where the f_j represents the fraction of electrons which contributes to each oscillator with force constant ξ_j . $f_j \cdot \xi_j$ represents oscillator strengths.

A useful semi-empirical relation for refractive index in the transparency region of a crystal known as the Sellmeier formula follows directly from Eq. (59) under the assumption that, far from resonances, the damping constant terms $\Gamma_j \omega$ are negligible compared to $(\omega_j^2 - \omega^2)$

$$n^2 = 1 + \sum_j \frac{A_j \lambda^2}{\lambda^2 - \lambda_j^2} \quad (61)$$

Sum Rules

The definition of oscillator strength results in the sum rule for electronic interactions

$$\sum_j f_j = Z \quad (62)$$

where Z is the number of electrons per atom. The periodicity of the lattice in solids (see “Energy Band Structures” in Sec. 8.7) leads to the modification of this sum rule to

$$\sum_m f_{mn} = \frac{m}{\hbar^2} \frac{\partial^2 \epsilon_n}{\partial k^2} - 1 = \frac{m}{m_n^*} - 1 \quad (63)$$

where m_n^* is an effective mass (see “Energy Band Structures” in Sec. 8.7).

Another sum rule for solids equivalent to Eq. (62) relates the imaginary part of the permittivity or dielectric constant and the plasma frequency, ω_p ,

$$\int_0^\infty \omega \epsilon_2(\omega) d\omega = \frac{1}{2} \pi \omega_p^2 \quad (64)$$

where $\omega_p^2 = Ne^2/\epsilon_0 m$ (see “Drude Model” in Sec. 8.6).

Dispersion relations are integral formulas which relate refractive properties to absorptive process. Kramers-Kronig relations are commonly used dispersion integrals based on the condition of causality which may be related to sum rules. These relations can be expressed in alternative forms. For instance, the reflectivity of a solid is often measured at normal incidence and dispersion relations used to determine the optical properties. Writing the complex reflectivity amplitude as

$$\tilde{r}(\omega) = r_r(\omega) e^{i\theta(\omega)} \quad (65)$$

the phase shift, θ , can be determined by integrating the experimental measurement of the real amplitude r_r

$$\theta(\omega) = -\frac{2\omega}{\pi} \mathcal{P} \int_0^\infty \frac{\ln r_r(\omega')}{\omega'^2 - \omega^2} d\omega' \quad (66)$$

and the optical constants determined from the complex Fresnel relation

$$r_r(\omega) e^{i\theta} = \frac{(n-1+i\kappa)}{(n+1+i\kappa)} \quad (67)$$

Sum rules following from the Kramers-Kronig relations relate the refractive index $n(\omega)$ at a given frequency, ω , to the absorption coefficient, $\alpha(\omega')$, integrated over all frequencies, ω' , according to

$$n(\omega) - 1 = \frac{c}{\omega} \mathcal{P} \int_0^\infty \frac{\alpha(\omega') d\omega'}{\omega'^2 - \omega^2} \quad (68)$$

Similarly, the real and imaginary parts of the dielectric constant, ϵ_1 and ϵ_2 , may be related via the integral relations

$$\epsilon_1(\omega) - 1 = \frac{2}{\pi} \mathcal{P} \int_0^\infty \frac{\omega' \epsilon_2(\omega')}{\omega'^2 - \omega^2} d\omega' \quad (69)$$

$$\epsilon_2(\omega) = -\frac{2\omega}{\pi} \mathcal{P} \int_0^\infty \frac{\epsilon_1(\omega') - 1}{\omega'^2 - \omega^2} d\omega' \quad (70)$$

8.5 LATTICE INTERACTIONS

The adiabatic approximation is the normal starting point for a consideration of the coupling of light with lattice vibrations, i.e., it is assumed that the response of the outer shell electrons of the atoms to an electric field is much faster than the response of the core together with its inner electron shells. Further, the harmonic approximation assumes that for small displacements, the restoring force on the ions will be proportional to the displacement. The solution of the equations of motion for the atoms within a solid under these conditions gives normal modes of vibration whose frequency eigenvalues and displacement eigenvectors depend on the crystal symmetry, atomic separation, and the detailed form of the interatomic forces. The frequency of lattice vibrations in solids is typically in the 100 to 1000 cm^{-1} range (wavelengths between 10 and 100 μm). Longitudinal and doubly degenerate transverse vibrational modes have different natural frequencies due to long range Coulomb interactions. Infrared or Raman activity can be determined for a given crystal symmetry by representing the modes of vibration as irreducible representations of the space group of the crystal lattice.

Infrared Dipole Active Modes

If the displacement of atoms in a normal mode of vibration produces an oscillating dipole moment, then the motion is dipole active. Thus, harmonic vibrations in ionic crystals contribute directly to the dielectric function, whereas higher order contributions are needed in nonpolar crystals. Since photons have small wavevectors compared to the size of the Brillouin zone in solids, only zone center lattice vibrations, (i.e. long wavelength phonons), can couple to the radiation. This conservation of wavevector (or momentum) also implies that only optical phonons interact. In a dipole active, long wavelength optical mode, oppositely charged ions within each primitive cell undergo oppositely directed displacements giving rise to a nonvanishing polarization. Group theory shows that, within the harmonic approximation, the infrared active modes have irreducible representations with the same transformation properties as x , y , or z . The strength of the light-dipole coupling will depend on the degree of charge redistribution between the ions, i.e., the relative ionicity of the solid.

Classical dispersion theory leads to a phenomenological model for the optical interaction with dipole active lattice modes. Because of the transverse nature of electromagnetic radiation, the electric field vector couples with the transverse optical (TO) phonons and the maximum absorption therefore occurs at this resonance. The resonance frequency, ω_{TO} , is inserted into the solution of the equation of motion, Eq. (59). Since electronic transitions typically occur at frequencies 10^2 to 10^3 higher than the frequency of vibration of the ions, the atomic polarizability can be represented by a single high frequency permittivity, $\epsilon(\infty)$. The dispersion relation for a crystal with several zone center TO phonons may be written

$$\tilde{\epsilon}(\omega) = \epsilon(\infty) + \sum_j \frac{S_j}{(\omega_{\text{TO}j}^2 - \omega^2 - i\Gamma_j\omega)} \quad (71)$$

By defining a low frequency permittivity, $\epsilon(0)$, the oscillator strength for a crystal possessing two atoms with opposite charge, Ze , per unit cell of volume, V , is

$$S = \frac{(\epsilon(\infty)/\epsilon_0 + 2)^2 (Ze)^2}{9m_r V} = \omega_{\text{TO}}^2 (\epsilon(0) - \epsilon(\infty)) \quad (72)$$

where Ze represents the “effective charge” of the ions, m_r is the reduced mass of the ions and the local field has been included based on Eq. (31). Figure 4 shows the form of the real and imaginary parts of the dielectric constant, the reflectivity and the polariton dispersion curve. Observing that the real part of the dielectric constant is zero at longitudinal phonon frequencies, ω_{LO} , the

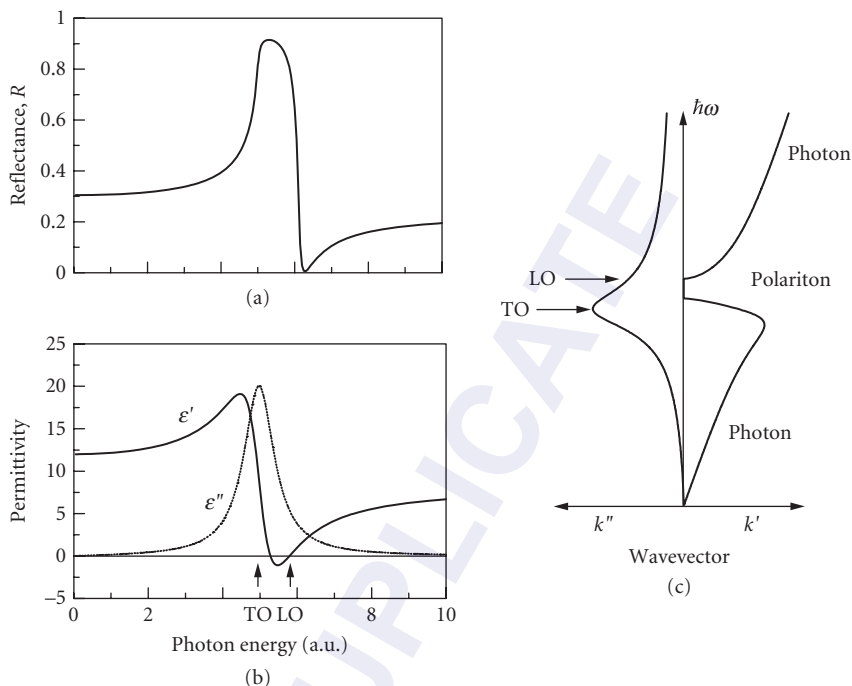


FIGURE 4 (a) Reflectance and (b) real and imaginary parts of the permittivity of a solid with a single infrared active mode. (c) Polariton dispersion curves (real and imaginary parts) showing the frequencies of the longitudinal and transverse optical modes.

Lyddane-Sachs-Teller relation may be derived, which in its general form for a number of dipole active phonons, is given by

$$\frac{\epsilon(0)}{\epsilon(\infty)} = \prod_j \left(\frac{\omega_{Lj}}{\omega_{Tj}} \right)^2 \quad (73)$$

These relations [see Eqs. (71) to (73)] give good fits to measured reflectivities in a wide range of ionically (or partially ionically) bonded solids. The LO-TO splitting and effective charge, Ze , depends on the ionicity of the solid; however, the magnitude of Ze determined from experiments does not necessarily quantify the ionicity since this “rigid ion” model does not account for the change of polarizability due to the distortion of the electron shells during the vibration.

In uniaxial and biaxial crystals, the restoring forces on the ions are anisotropic resulting in different natural frequencies depending on the direction of light propagation as well as the transverse or longitudinal nature of the vibration. Similar to the propagation of light, “ordinary” and “extraordinary” transverse phonons may be defined with respect to the principal axes. For instance, in a uniaxial crystal under the condition that the anisotropy in phonon frequency is smaller than the LO-TO frequency splitting, infrared radiation of frequency ω propagating at an angle θ to the optic axis will couple to TO phonons according to the relation

$$\omega_T^2 = \omega_{T\parallel}^2 \sin^2 \theta + \omega_{T\perp}^2 \cos^2 \theta \quad (74)$$

where $\omega_{T\parallel}$ is a TO phonon propagating with atomic displacements parallel to the optic axis, and $\omega_{T\perp}$ is a TO phonon propagating with atomic displacements perpendicular to the optic axis. The corresponding expression for LO modes is

$$\omega_L^2 = \omega_{T\parallel}^2 \cos^2 \theta + \omega_{T\perp}^2 \sin^2 \theta \quad (75)$$

In Table 2, the irreducible representations of the infrared active normal modes for the different crystal symmetries are labeled x , y , or z .

Brillouin and Raman Scattering

Inelastic scattering of radiation by acoustic phonons is known as Brillouin scattering, while the term Raman scattering is normally reserved for inelastic scattering from optic phonons in solids. In the case of Brillouin scattering, long wavelength acoustic modes produce strain and thereby modulate the dielectric constant of the medium thus producing a frequency shift in scattered light. In a Raman active mode, an incident electric field produces a dipole by polarizing the electron cloud of each atom. If this induced dipole is modulated by a lattice vibrational mode, coupling occurs between the light and the phonon and inelastic scattering results. Each Raman or Brillouin scattering event involves the destruction of an incident photon of frequency, ω_i , the creation of a scattered photon, ω_s , and the creation or destruction of a phonon of frequency, ω_p . The frequency shift, $\omega_i \pm \omega_s = \omega_p$, is typically 100 to 1000 cm^{-1} for Raman scattering but only a few wavenumbers for Brillouin scattering.

Atomic polarizability components have exactly the same transformation properties as the quadratic functions x^2, xy, \dots, z^2 . The Raman activity of the modes of vibration of a crystal with a given point group symmetry can thus be deduced from its group theoretical character table. Polarization selection rules may be deduced from the Raman tensors given in Table 2. The scattering efficiency, S_R , for a mode corresponding to one of the irreducible representations listed is given by

$$S_R = A \left[\sum_{\rho, \sigma} e_i^\sigma R_{\sigma, \rho} e_s^\rho \right]^2 \quad (76)$$

where A is a constant of proportionality, $R_{\sigma, \rho}$ is the Raman coefficient of the representation, and e_i^σ and e_s^ρ are components of the unit vectors of polarization of the incident, i , and scattered, s , radiation along the principal axes, where σ and $\rho = x, y$, and z .

Not all optic modes of zero wavevector are Raman active. Raman activity is often complementary to infrared activity. For instance, since the optic mode in the diamond lattice has even parity, it is Raman active but not infrared active, whereas the zone center mode in sodium chloride is infrared active but not Raman active because the inversion center is on the atom site and so the phonon has odd parity. In piezoelectric crystals, which lack a center of inversion, some modes can be both Raman and infrared active. In this case the Raman scattering can be anomalous due to the long-range electrostatic forces associated with the polar lattice vibrations.

The theory of Brillouin scattering is based on the elastic deformation produced in a crystal by a long wavelength acoustic phonon. The intensity of the scattering depends on the change in refractive index with the strain induced by the vibrational mode. A strain, s_{ij} in the lattice produces a change in the component of permittivity, $\epsilon_{\mu\nu}$, given by

$$\delta \epsilon_{\mu\nu} = - \sum_{\rho, \sigma} \epsilon_{\mu\rho} p_{\rho\sigma, ij} \epsilon_{\sigma\nu} s_{ij} \quad (77)$$

where $p_{\rho\sigma, ij}$ is an elasto-optical coefficient.³ The velocity of the acoustic phonons and their anisotropy can be determined from Brillouin scattering measurements.

TABLE 2 Infrared and Raman-Active Vibrational Symmetries and Raman Tensors¹⁴

Monoclinic		$\begin{bmatrix} a & d \\ d & c \end{bmatrix} \begin{bmatrix} e & f \\ e & f \end{bmatrix}$
2	C_2	$A(y) \quad B(x,z)$
m	C_v	$A'(x,z) \quad A'(y)$
2/m	C_{2h}	$A_g \quad B_g$
Orthorhombic		
mm	C_{2v}	$\begin{bmatrix} a & & \\ & d & \\ & & c \end{bmatrix} \begin{bmatrix} & & e \\ & & e \\ & & f \end{bmatrix}$
222	D_2	$A_1(z) \quad A_2 \quad B_1(x) \quad B_2(y) \quad B_3(x)$
mmm	D_{2h}	$A_g \quad B_1(z) \quad B_2(y) \quad B_3(x) \quad B_{2g} \quad B_{3g}$
Trigonal		
3	C_3	$\begin{bmatrix} a & & \\ & d & e \\ & & c \end{bmatrix} \begin{bmatrix} & & d & -c & -f \\ & & d & -c & -d & e \\ & & e & f & -f & e \end{bmatrix}$
$\bar{3}$	S_6	$A(z) \quad E(x) \quad E(y) \quad E_g$
Tetragonal		
3m	C_{3v}	$\begin{bmatrix} a & & \\ & a & b \\ & & b \end{bmatrix} \begin{bmatrix} & & c & -c & d \\ & & c & -c & d \\ & & d & -d & -d \end{bmatrix}$
32	D_3	$A_1(z) \quad E(y) \quad E(x) \quad E_g$
$\bar{3}m$	D_{3d}	$A_1(z) \quad E(x) \quad E(y) \quad E_g$
4	C_4	$\begin{bmatrix} a & & \\ & a & b \\ & & b \end{bmatrix} \begin{bmatrix} & & c & d \\ & & d & -c \\ & & e & f \\ & & f & -f & e \end{bmatrix}$
$\bar{4}$	S_4	$A(z) \quad B \quad E(x) \quad E(y) \quad E_g$
4/m	C_{4h}	$A \quad A_g \quad B(z) \quad E(x) \quad E(-y) \quad E_g$
4mm	C_{4v}	$\begin{bmatrix} a & & \\ & a & b \\ & & b \end{bmatrix} \begin{bmatrix} & & c & -c \\ & & d & d \\ & & e & e \end{bmatrix}$
$\bar{4}2m$	D_{2d}	$A_1(z) \quad B_1 \quad B_2 \quad E(x) \quad E_g$
42	D_4	$A_1 \quad B_1 \quad B_2(z) \quad E(y) \quad E(x) \quad E_g$
4/mmm	D_{4h}	$A_1(z) \quad A_1 \quad B_1 \quad B_2 \quad E(-y) \quad E_g$
		$A_{1g} \quad B_{1g} \quad B_{2g} \quad E_g$

TABLE 2 Infrared and Raman-Active Vibrational Symmetries and Raman Tensors¹⁴(Continued)

Hexagonal		$\begin{bmatrix} a \\ a \\ b \end{bmatrix}$	$\begin{bmatrix} c \\ c \\ d \end{bmatrix}$	$\begin{bmatrix} -d \\ c \end{bmatrix}$	$\begin{bmatrix} e & f \\ f & -e \end{bmatrix}$	$\begin{bmatrix} f & -e \\ -e & -f \end{bmatrix}$
6	C_6	$A_1(z)$	$E_1(x)$	$E_1(y)$	E_2	E_2
$\bar{6}$	C_3^h	A_1'	$E_1''(x)$	$E_1''(y)$	$E_2'(x)$	$E_2'(y)$
6/m	C_6^h	A_1^g	$E_1^g(x)$	$E_1^g(y)$	$E_2^g(x)$	$E_2^g(y)$
		$\begin{bmatrix} a \\ a \\ b \end{bmatrix}$	$\begin{bmatrix} c \\ c \\ c \end{bmatrix}$	$\begin{bmatrix} -c \\ -c \end{bmatrix}$	$\begin{bmatrix} d \\ d \end{bmatrix}$	$\begin{bmatrix} d \\ -d \end{bmatrix}$
6mm	C_{6v}	$A_1(z)$	$E_1(y)$	$E_1(-x)$	E_2	E_2
$\bar{6}m2$	D_{3h}	A_1	$E_1''(y)$	$E_1''(x)$	$E_2'(x)$	$E_2'(y)$
62	D_6	A_1	$E_1(x)$	$E_1(y)$	E_2	E_2
6/mmm	D_{6h}	A_{1g}	$E_{1g}(x)$	$E_{1g}(y)$	$E_{2g}(x)$	$E_{2g}(y)$
Cubic		$\begin{bmatrix} a \\ a \\ a \end{bmatrix}$	$\begin{bmatrix} b \\ b \\ -2b \end{bmatrix}$	$\begin{bmatrix} -3\frac{1}{2}b \\ -3\frac{1}{2}b \end{bmatrix}$	$\begin{bmatrix} d \\ d \end{bmatrix}$	$\begin{bmatrix} d \\ d \end{bmatrix}$
23	T	A	E	E	$F(x)$	$F(y)$
m3	T_h	A_g	E_g	E_g	$F_g(x)$	$F_g(y)$
$\bar{4}3m$	T_d	A_1	E	E	$F_2(x)$	$F_2(y)$
432	O	A_1	E	E	F_2	F_2
m3m	O_h	A_{1g}	E_g	E_g	$F_{2g}(x)$	$F_{2g}(y)$

8.6 FREE ELECTRON PROPERTIES

Fundamental optical properties of metals and semiconductors with high densities of free carriers are well described using a classical model. Reflectivity is the primary property of interest because of the high absorption.

Drude Model

The Drude model for free electrons is a special condition of the classical Lorentz model (Sec. 8.4) with the Hooke's law force constant, $\xi = 0$, so that the resonant frequency is zero. In this case,

$$\frac{\epsilon_1}{\epsilon_0} = n^2 - \kappa^2 = 1 - \frac{\omega_p^2}{\omega^2 + \tau^{-2}} \quad (78)$$

and

$$\frac{\epsilon_2}{\epsilon_0} = 2n\kappa = \frac{\omega_p^2}{\omega^2 + \tau^{-2}} \left(\frac{1}{\omega\tau} \right) \quad (79)$$

where ω_p is the plasma frequency

$$\omega_p = \sqrt{\frac{Ne^2}{\epsilon_0 m}} = \sqrt{\frac{\mu_0 \sigma c^2}{\tau}} \quad (80)$$

and $\tau (= 1/\Gamma)$ is the scattering or relaxation time for the electrons. In ideal metals ($\sigma \rightarrow \infty$), $n = \kappa$. Figure 5a shows the form of the dispersion in the real and imaginary parts of the dielectric constant for free electrons, while the real and imaginary parts of the refractive index are illustrated in Fig. 5b. The plasma frequency is defined by the point at which the real part changes sign. The reflectivity is plotted in Fig. 5c and shows a magnitude close to 100 percent below the plasma frequency but falls rapidly to a small value above ω_p . The plasma frequency determined solely by the free electron term is typically on the order of 10 eV in metals accounting for their high reflectivity in the visible.

Interband Transitions in Metals

Not all metals are highly reflective below the plasma frequency. The noble metals possess optical properties which combine free electron intraband (Drude) and interband contributions. A typical metal has d-levels at energies a few electron volts below the electron Fermi level. Transitions can be optically induced from these d-states to empty states above the Fermi level. These transitions normally occur in the ultraviolet spectral region but can have a significant influence on the optical properties of metals in the visible spectral region via the real part of the dielectric constant. Describing the interband effects $\delta\epsilon_b$ within a classical (Lorentz) model the combined effects on the dielectric constant may be added.

$$\tilde{\epsilon} = 1 + \delta\tilde{\epsilon}_b + \delta\tilde{\epsilon}_f \quad (81)$$

The interband contribution to the real part of the dielectric constant is positive and shows a resonance near the transition frequency. On the other hand, the free electron contribution is negative below the plasma frequency. The interband contribution can cause a shift to shorter wavelengths of the zero cross-over in ϵ_1 , thus causing a reduction of the reflectivity in the blue. For instance d-states in copper lie only 2 eV below the Fermi level, which results in its characteristic color.

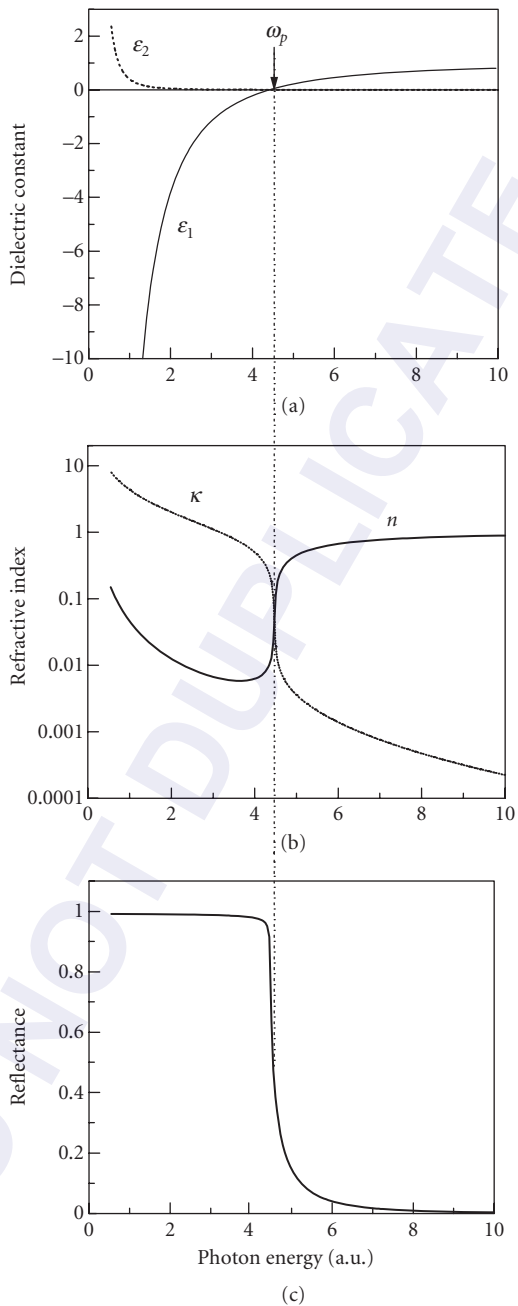


FIGURE 5 Dispersion of: (a) the real and imaginary parts of the dielectric constant; (b) real and imaginary parts of the refractive index; and (c) the reflectance according to the Drude model where ω_p is the plasma frequency.

Reflectivity

Absorption in metals is described by complex optical constants. The reflectivity is accompanied by a phase change and the penetration depth is

$$d_M = \frac{c}{2\omega\kappa} = \frac{\lambda}{4\pi\kappa} \quad (82)$$

At normal incidence at an air-metal interface, the reflectance is given by

$$R = \frac{(n-1)^2 + \kappa^2}{(n+1)^2 + \kappa^2} = 1 - \frac{2}{\kappa} = 1 - 2\sqrt{\frac{2\varepsilon_0}{\varepsilon_2}} \quad (83)$$

By analogy with the law of refraction (Snell's Law) a complex refractive index can be defined by the refraction equation

$$\sin\theta_t = \frac{1}{\tilde{\eta}} \sin\theta_i \quad (84)$$

Since $\tilde{\eta}$ is complex, θ_t is also complex and the phase change on reflection can take values other than 0 and π . For nonnormal incidence, it can be shown that the surfaces of constant amplitude inside the metal are parallel to the surface, while surfaces of constant phase are at an angle to the surface. The electromagnetic wave in a metal is thus inhomogeneous. The real and imaginary parts of the refractive index can be determined by measuring the amplitude and phase of the reflected light. Writing the s and p components of the complex reflected fields in the form

$$E_{rp} = \rho_p e^{i\phi_p}; \quad E_{rs} = \rho_s e^{i\phi_s} \quad (85)$$

and defining the real amplitude ratio and phase differences as

$$P_M = \tan\psi = \frac{\rho_s}{\rho_p}; \quad \Delta_M = \phi_p - \phi_s \quad (86)$$

then the real and imaginary parts of the refractive index are given by

$$n \approx -\frac{\sin\theta_i \tan\theta_i \cos 2\psi}{1 + \sin 2\psi \cos \Delta_M} = -\sin\bar{\theta}_i \tan\bar{\theta}_i \cos 2\bar{\psi} \quad (87)$$

$$\kappa \approx \tan 2\psi \sin \Delta_M = -\tan 2\bar{\psi} \quad (88)$$

$\bar{\theta}_i$ is the principal angle which occurs at the maximum in P_M at the condition $\Delta_M = \pi/2$, (see Fig. 3c and 3d), which is equivalent to Brewster's angle at an interface between two nonabsorbing dielectrics, (see Fig. 3a and 3b).

Plasmons

Plasmons are oscillations of fluctuations in charge density. The condition for these oscillations to occur is the same as the condition for the onset of electromagnetic propagation at the plasma frequency. Volume plasmons are not excited by light at normal incidence since they are purely longitudinal. Oscillations cannot be produced by transverse electromagnetic radiation with zero divergence. However at the surface of a solid, an oscillation in surface charge density is possible. At an interface

between a metal with permittivity, ϵ_m and a dielectric with permittivity, ϵ_d , the condition $\epsilon_m = -\epsilon_d$ such that (neglecting damping and assuming a free electron metal) a surface plasmon can be created with frequency

$$\omega_s = \frac{\omega_p}{(\epsilon_d/\epsilon_o + 1)^{1/2}} \quad (89)$$

By altering the angle of incidence, the component of the electromagnetic wavevector can be made to match the surface plasmon mode.

8.7 BAND STRUCTURES AND INTERBAND TRANSITIONS

Advances in semiconductors for electronic and optoelectronic applications have encouraged the development of highly sophisticated theories of interband absorption in semiconductors. In addition, the development of low dimensional structures (quantum wells, quantum wires, and quantum dots) have provided the means of “engineering” the optical properties of solids. The approach here has been to outline the basic quantum mechanical development for interband transitions in solids.

Quantum Mechanical Model

The quantum theory of absorption considers the probability of an electron being excited from a lower energy level to a higher level. For instance, an isolated atom has a characteristic set of electron levels with associated wavefunctions and energy eigenvalues. The absorption spectrum of the atom thus consists of a series of lines whose frequencies are given by

$$\hbar\omega_{fi} = \mathcal{E}_f - \mathcal{E}_i \quad (\mathcal{E}_f > \mathcal{E}_i) \quad (90)$$

where \mathcal{E}_f and \mathcal{E}_i are a pair of energy eigenvalues. We also know that the spontaneous lifetime, τ , for transitions from any excited state to a lower state sets a natural linewidth of order \hbar/τ based on the uncertainty principle. The Schrödinger equation for the ground state with wavefunction, ϕ_i , in the unperturbed system

$$\mathcal{H}_o \phi_i = \mathcal{E}_i \phi_i \quad (91)$$

is represented by the time-independent hamiltonian, \mathcal{H}_o . The optical interaction can be treated by first order perturbation theory. By introducing a perturbation term based on the classical oscillator

$$\mathcal{H}' = e\mathbf{E} \cdot \mathbf{r} \quad (92)$$

this leads to a similar expression to the Lorentz model, Eq. (59)

$$\tilde{\eta}^2 = \frac{\epsilon}{\epsilon_o} = 1 + \frac{Ne^2}{m\epsilon_o} \sum_m \frac{f_{fi}}{(\omega_{fi}^2 - \omega^2 - i\Gamma_{fi}\omega)} \quad (93)$$

where

$$f_{jj} = \frac{2|\mathbf{p}_{jj}|^2}{m\hbar\omega_{jj}} \quad (94)$$

and \mathbf{p}_{jj} are momentum matrix elements defined by

$$\mathbf{p}_{jj} = \langle \varphi_j | \mathbf{p} | \varphi_j \rangle = \int \varphi_j^* (i\hbar \nabla) \varphi_j d\mathbf{r} \quad (95)$$

Perturbation theory to first order gives the probability per unit time that a perturbation of the form $\mathcal{H}(t) = \mathcal{H}_p \exp(i\omega t)$ induces a transition from the initial to final state,

$$W_{fi} = \frac{2\pi}{\hbar} \left| \langle \varphi_f | \mathcal{H}_p | \varphi_i \rangle \right|^2 \delta(\mathcal{E}_f - \mathcal{E}_i - \hbar\omega) \quad (96)$$

This is known as Fermi's golden rule.

Energy Band Structures

If we imagine N similar atoms brought together to form a crystal, each degenerate energy level of the atoms will spread into a band of N levels. If N is large, these levels can be treated as a continuum of energy states. The wavefunctions and electron energies of these energy bands can be calculated by various approximate methods ranging from nearly free electron to tight binding models. The choice of approach depends on the type of bonding between the atoms.

Within the one electron and adiabatic assumptions, each electron moves in the periodic potential, $V(\mathbf{r})$, of the lattice leading to the Schrödinger equation for a single particle wavefunction

$$\left[\frac{p^2}{2m} + V(\mathbf{r}) \right] \psi(\mathbf{r}) = \mathcal{E} \psi(\mathbf{r}) \quad (97)$$

where the momentum operator is given by $\mathbf{p} = -i\hbar \nabla$. The simple free electron solution of the Schrödinger equation (i.e., for $V(\mathbf{r}) = 0$), is a parabolic relationship between energy and wavevector. The solution including a periodic potential, $V(\mathbf{r})$ has the form

$$\psi_{\mathbf{k}}(\mathbf{r}) = \exp(i\mathbf{k} \cdot \mathbf{r}) \cdot u_{\mathbf{k}}(\mathbf{r}) \quad (98)$$

where \mathbf{k} is the electron wavevector and $u_{\mathbf{k}}(\mathbf{r})$ has the periodicity of the crystal lattice. This is known as the Bloch solution. The allowed values of \mathbf{k} are separated by $2\pi/L$, where L is the length of the crystal. The wavevector is not uniquely defined by the wavefunction, but the energy eigenvalues are a periodic function of \mathbf{k} . For an arbitrarily weak periodic potential

$$\mathcal{E} = \frac{\hbar^2}{2m} |\mathbf{k} + \mathbf{G}|^2 \quad (99)$$

where \mathbf{G} is a reciprocal lattice vector (in one dimension $\mathbf{G} = 2\pi n/a$, where a is the lattice spacing and n is an integer. Thus we need only consider solutions which are restricted to a reduced zone, referred to as the first Brillouin zone, in reciprocal space (between $\mathbf{k} = -\pi/a$ and π/a in one dimension). Higher energy states are folded into the first zone consistent with Eq. (99) to form a series of energy bands. Figure 6 shows the first Brillouin zones for face centered cubic (fcc) crystal lattices and energy levels for a weak lattice potential.

A finite periodic potential, $V(\mathbf{r})$, alters the shape of the free electron bands. The curvature of the bands is described by m^* , an "effective mass," which is defined by the slope of the dispersion curve at a given \mathbf{k} :

$$\frac{1}{m_{\mathbf{k}}^*} = \frac{1}{\hbar^2 k} \frac{d\mathcal{E}}{dk} \quad (100)$$

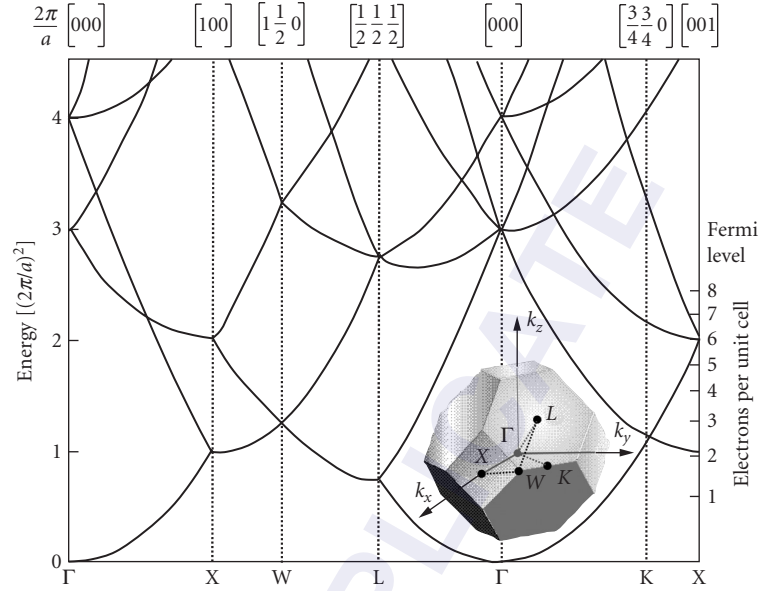


FIGURE 6 Free electron band structure in the reduced Brillouin zone for face-centered-cubic lattices. The insert shows the first Brillouin zone with principal symmetry points labeled. This applies to crystals such as Al, Cu, Ag, Si, Ge, and GaAs.⁵

At the zone center, $\mathbf{k} = 0$, this reduces to the parabolic relationship

$$\mathcal{E} = \frac{\hbar^2 k^2}{2m_o^*} \quad (101)$$

Effective masses can be related to interband momentum matrix elements and energy gaps using perturbation theory. Substituting the Bloch function, Eq. (98), into the Schrödinger equation, Eq. (97), and identifying each band by an index, j , gives

$$\left[\frac{p^2}{2m} + \frac{\hbar}{m} \mathbf{k} \cdot \mathbf{p} + \frac{\hbar^2 k^2}{2m} + V(\mathbf{r}) \right] u_{jk}(\mathbf{r}) = \mathcal{E}_{jk}(\mathbf{r}) u_{jk}(\mathbf{r}) \quad (102)$$

The $\mathbf{k} \cdot \mathbf{p}$ term can be treated as a perturbation about a specific point in \mathbf{k} -space. For any \mathbf{k} , the set of all $u_{jk}(\mathbf{r})$ (corresponding to the N energy levels) forms a complete set, i.e., the wavefunction at any value of \mathbf{k} can be expressed as a linear combination of all wavefunctions at another \mathbf{k} . Second order perturbation theory then predicts an effective mass given by

$$\frac{1}{m^*} = \frac{1}{m} + \frac{2}{m} \sum_j \frac{|\mathbf{p}_{jj}|^2}{\mathcal{E}_j(\mathbf{k}) - \mathcal{E}_j(\mathbf{k})} \quad (103)$$

In principle the summation in Eq. (103) is over all bands; however, this can usually be reduced to a few nearest bands because of the resonant denominator. For example, in diamond- and zinc-blended structured semiconductors, the Kane momentum matrix element, P , defined by

$$P = -\frac{i\hbar}{m} \langle S_{co} | p_x | X_{vo} \rangle \quad (104)$$

successfully characterizes the band structure and optical properties close to zone center. Here S_{CO} is a spherically symmetric s-like atomic wavefunction representing the lowest zone center conduction band state and X_{VO} is a p-like function with x symmetry from the upper valence bands. In this case, including only the three highest valence bands and the lowest conduction band in the summation of Eq. (91), the conduction band effective mass is given by

$$\frac{1}{m_{co}^*} = \frac{1}{m} + \frac{2}{3\hbar^2} \left[\frac{2P^2}{\mathcal{E}_g} + \frac{P^2}{\mathcal{E}_g + \Delta_{SO}} \right] \quad (105)$$

where \mathcal{E}_g is the band gap energy and Δ_{SO} is the spin-orbit splitting. By inverting this expression, the momentum matrix element may be determined from measurements of effective mass and the bandgaps. P is found to have similar magnitudes for a large number of semiconductors. Equation (105) illustrates the general rule that the effective mass of the conduction band is approximately proportional to the band gap energy.

Direct Interband Absorption

In the case of a solid, the first order perturbation of the single electron hamiltonian by electromagnetic radiation is more appropriately described by

$$\mathcal{H}'(t) = \frac{e}{mc} \mathbf{A} \cdot \mathbf{p} \quad (106)$$

rather than Eq. (92). \mathbf{A} is the vector potential,

$$\mathbf{A}(\mathbf{r}, t) = A_0 \hat{\xi} \exp[i(\mathbf{q} \cdot \mathbf{r} - \omega t)] + c.c. \quad (107)$$

\mathbf{q} is the wavevector and $\hat{\xi}$ is the unit polarization vector of the electric field. (Note that this perturbation is of a similar form to the $\mathbf{k} \cdot \mathbf{p}$ perturbation described earlier.) Using Fermi's golden rule, the transition probability per unit time between a pair of bands is given by

$$W_{fi} = \frac{2\pi}{\hbar} \left(\frac{eA_0}{mc} \right)^2 \left| \langle \psi_f | \hat{\xi} \cdot \mathbf{p} | \psi_i \rangle \right|^2 \delta(\mathcal{E}_f(\mathbf{k}) - \mathcal{E}_i(\mathbf{k}) - \hbar\omega) \quad (108)$$

Conservation of momentum requires a change of electron momentum after the transition; however, the photon momentum is very small, so that vertical transitions in \mathbf{k} -space can be assumed in most cases (the electric dipole approximation). The total transition rate per unit volume, $W_T(\omega)$ is obtained by integrating over all possible vertical transitions in the first Brillouin zone taking account of all contributing bands:

$$W_T(\omega) = \frac{2\pi}{\hbar} \left(\frac{2\pi e^2 I}{ncm^2 \omega^2} \right) \sum_f \int \frac{d\mathbf{k}}{(2\pi)^3} \left| \hat{\xi} \cdot \mathbf{p}_f(\mathbf{k}) \right|^2 \delta(\mathcal{E}_f(\mathbf{k}) - \mathcal{E}_i(\mathbf{k}) - \hbar\omega) \quad (109)$$

Here the vector potential has been replaced with the irradiance, I , of the radiation through the relation

$$A_0 = \frac{2\pi c}{n\omega^2} I \quad (110)$$

Note that the momentum matrix element as defined in Eq. (95) determines the oscillator strength for the absorption. \mathbf{p}_{fi} can often be assumed slowly varying in \mathbf{k} so that the zone center matrix element can be employed for interband transitions and the frequency dependence of the absorption coefficient is dominated by the density of states.

Joint Density of States

The delta function in the integration of Eq. (109) represents energy conservation for the transitions between any two bands. If the momentum matrix element can be assumed slowly varying in \mathbf{k} , then the integral can be rewritten in the form

$$J_{fi}(\omega) = \frac{1}{(2\pi)^3} \int d\mathbf{k} \delta(\mathcal{E}_f(\mathbf{k}) - \mathcal{E}_i(\mathbf{k}) - \hbar\omega) = \frac{1}{(2\pi)^3} \int \frac{d\mathbf{S}}{|\nabla_{\mathbf{k}} \mathcal{E}_{fi}(\mathbf{k})|} \quad (111)$$

where $d\mathbf{S}$ is a surface element on the equal energy surface in \mathbf{k} -space defined by $\mathcal{E}_{fi}(\mathbf{k}) = \mathcal{E}_f(\mathbf{k}) - \mathcal{E}_i(\mathbf{k}) = \hbar\omega$. Written in this way, $J(\omega)$ is the joint density of states between the two bands (note the factor of two for spin is excluded in this definition). Points in \mathbf{k} -space for which the condition

$$\nabla_{\mathbf{k}} \mathcal{E}_{fi}(\mathbf{k}) = 0 \quad (112)$$

hold form critical points called van Hove singularities which lead to prominent features in the optical constants. In the neighborhood of a critical point at \mathbf{k}_c , a constant energy surface may be described by the Taylor series

$$\mathcal{E}_{fi}(\mathbf{k}) = \mathcal{E}_c(\mathbf{k}_c) + \sum_{\mu=1}^3 \beta_{\mu} k_{\mu}^2 \quad (113)$$

where μ represents directional coordinates. Minimum, maximum, and saddle points arise depending on the relative signs of the coefficients, β_{μ} . Table 3 gives the frequency dependence of the joint density of states in three-dimensional (3D), two-dimensional (2D), one-dimensional (1D), and zero-dimensional (0D) solids. The absorption coefficient, α , defined by Beer's law may now be related to the transition rate by

$$\alpha(\omega) = I^{-1} \frac{dI}{dz} = \frac{\hbar\omega}{I} W_T \quad (114)$$

Thus, the minimum fundamental absorption edge of semiconductors and insulators (in the absence of excitonic effects) has the general form (Fig. 7a)

$$\alpha = A(\hbar\omega - \mathcal{E}_o)^{1/2} \quad (115)$$

Selection Rules and Forbidden Transitions

Direct interband absorption is allowed when the integral in Eq. (95) is nonzero. This occurs when the wavefunctions of the optically coupled states have opposite parity for single photon transitions. Transitions may be forbidden for other wavefunction symmetries. Although the precise form of the wavefunction may not be known, the selection rules can be determined by group theory from a knowledge of the space group of the crystal and symmetry of the energy band. Commonly, a single photon transition which is not allowed at the zone center because two bands have like parity, will

TABLE 3 Density of States in 3, 2, 1, and 0 Dimensions

	β_1	β_2	β_3	$E < E_c$	$E > E_c$	
3D						
M_0 , min.	+	+	+	0	$C_0(E-E_c)^{1/2}$	
M_1 , saddle	+	+	-	$C_1 - C_1'(E_c - E)^{1/2}$	C_1	
M_2 , saddle	+	-	-	C_2	$C_2 - C_2'(E - E_c)^{1/2}$	
M_3 , max.	-	-	-	$C_3(E_c - E)^{1/2}$	0	
2D						
P_0 , min.	+	+		0	B	
P_1 , saddle	+	-		$-\frac{B}{\pi} \ln \left 1 - \frac{E}{E_c} \right $		
P_2 , max	-	-		B	0	
1D						
Q_0 , min.	+			0	$A(E_c - E)^{-1/2}$	
Q_1 , max	-			$A(E - E_c)^{-1/2}$	0	
0D				$\delta(E - E_c)$		

be allowed at finite \mathbf{k} because wavefunction mixing will give mixed parity states. In this case, the momentum matrix element may have the form

$$p_{fi}(\mathbf{k}) = (\mathbf{k} - \mathbf{k}_0) \cdot \nabla_{\mathbf{k}} [p_{fi}(\mathbf{k})]_{\mathbf{k}=\mathbf{k}_0} \quad (116)$$

that is, the matrix element is proportional to \mathbf{k} . For interband transitions at an M_0 critical point, the frequency dependence of the absorption coefficient can be shown to be (Fig. 7b)

$$\alpha(\omega) = A'(\hbar\omega - \mathcal{E}_o)^{3/2} \quad (117)$$

Indirect Transitions

Interband transitions may also take place with the assistance of a phonon. The typical situation is a semiconductor or insulator which has a lowest conduction band near a Brillouin zone boundary. The phonon provides the required momentum to move the electron to this location but supplies little energy. The phonon may be treated as an additional perturbation and therefore second order

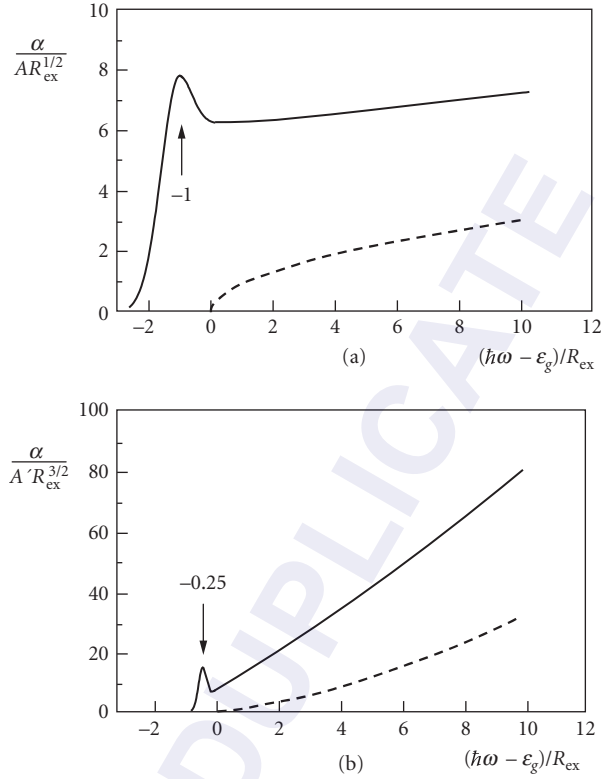


FIGURE 7 Illustration of absorption edge of crystals with: (a) direct allowed transitions and (b) direct forbidden transitions, based on density of state (dashed lines) and excitonic enhanced absorption models (solid lines).

perturbation theory is needed in the analysis of this two step processes. Theory predicts a frequency dependence for the absorption of the form

$$\alpha(\omega) \approx (\hbar\omega \pm \hbar\omega_{\text{ph}} - \mathcal{E}_o)^2 \quad (118)$$

where ω_{ph} is the phonon frequency, absorption or emission being possible. For forbidden indirect transitions, this relationship becomes

$$\alpha(\omega) \approx (\hbar\omega \pm \hbar\omega_{\text{ph}} - \mathcal{E}_o)^3 \quad (119)$$

Multiphoton Absorption

Multiphoton absorption can be treated by higher order perturbation theory. For instance, second order perturbation theory gives a transition rate between two bands

$$W_{\text{fi}}^{(2)} = \frac{2\pi}{\hbar} \left(\frac{eA_o}{mc} \right)^4 \left| \sum_t \frac{\langle \psi_f | \hat{\xi} \cdot \mathbf{p} | \psi_t \rangle \langle \psi_t | \hat{\xi} \cdot \mathbf{p} | \psi_i \rangle}{\mathcal{E}_t - \mathcal{E}_i - \hbar\omega} \right|^2 \delta(\mathcal{E}_f(\mathbf{k}) - \mathcal{E}_i(\mathbf{k}) - 2\hbar\omega) \quad (120)$$

where the summation spans all intermediate states, t . The interaction can be regarded as two successive steps. An electron first makes a transition from the initial state to an intermediate level of the system, t , by absorption of one photon. Energy is not conserved at this stage (momentum is) so that the absorption of a second photon must take the electron to its final state in a time determined by the energy mismatch and the uncertainty principle. In multiphoton absorption, one of the transitions may be an intraband self-transition. Since the probability depends on the arrival rate of the second photon, multiphoton absorption is intensity dependent. The total transition rate is given by

$$W_T^{(2)}(\omega) = \frac{2\pi}{\hbar} \left(\frac{4\pi^2 e^4 I^2}{n^2 c^2 m^4 \omega^4} \right) \sum_f \int \frac{d\mathbf{k}}{(2\pi)^3} \left| \sum_t \frac{\hat{\xi} \cdot \mathbf{p}_{ft}(\mathbf{k}) \hat{\xi} \cdot \mathbf{p}_{ti}(\mathbf{k})}{\mathcal{E}_t - \mathcal{E}_i - \hbar\omega} \right|^2 \delta(\mathcal{E}_f(\mathbf{k}) - \mathcal{E}_i(\mathbf{k}) - 2\hbar\omega) \quad (121)$$

The two photon absorption coefficient is defined by the relation

$$-\frac{dI}{dz} = \alpha I + \beta I^2 \quad (122)$$

so that,

$$\beta(\omega) = \frac{2\hbar\omega}{I^2} W_T^{(2)}(\omega) \quad (123)$$

Excitons

The interband absorption processes discussed earlier do not take into account Coulomb attraction between the excited electron and hole state left behind. This attraction can lead to the formation of a hydrogen-like bound electron-hole state or exciton. The binding energy of free (Wannier) excitons is typically a few meV. If not thermally washed out, excitons may be observed as a series of discrete absorption lines just below the bandgap energy. The energy of formation of an exciton is

$$\mathcal{E}_{\text{ex}} = \mathcal{E}_g + \frac{\hbar^2 |\mathbf{k}|^2}{2(m_e^* + m_h^*)} - \frac{R_{\text{ex}}}{n^2} \quad (124)$$

where R_{ex} is the exciton Rydberg,

$$R_{\text{ex}} = \frac{m_r^* e^4}{2\hbar^2 \epsilon_1^2} \quad (125)$$

m_r^* is the reduced effective mass and n is a quantum number. Optically created electron-hole pairs have equal and opposite momentum which can only be satisfied if $\mathbf{K} = 0$ for the bound pair and results in discrete absorption lines. Coulomb attraction also modifies the absorption above the bandgap energy. The theory of exciton absorption developed by Elliot predicts a modification to Eq. (113) for the direct allowed absorption coefficient above the band edge (Fig. 7a)

$$\alpha = \frac{\pi A R_{\text{ex}}^{1/2} e^{-\pi\gamma}}{\sinh \pi\gamma} \quad (126)$$

where

$$\gamma = \left(\frac{R_{\text{ex}}}{\hbar\omega - \mathcal{E}_g} \right)^{1/2} \quad (127)$$

Excitons associated with direct forbidden interband transitions do not show absorption to the lowest ($n = 1$) state but transitions to excited levels are allowed. Above the band edge for direct forbidden transitions the absorption has the frequency dependence (Fig. 7b)

$$\alpha = \frac{\pi A' R_{\text{ex}}^{3/2} \left(1 + \frac{1}{\gamma^2}\right) e^{\pi\gamma}}{\sinh \pi\gamma} \quad (128)$$

Figure 7 compares the form of the absorption edge based on the density of states function and a discrete exciton absorption line (dashed lines) with the absorption functions based on the Elliot theory and a typically broadened exciton (solid lines). Figure 7a is an illustration of a direct allowed gap, e.g., GaAs with the $n = 1$ exciton visible and Fig. 7b shows a forbidden direct absorption edge, e.g., Cu_2O . In the latter case, optical excitation of the $n = 1$ exciton is forbidden, but the $n = 2$ and higher exciton transitions are allowed.

8.8 REFERENCES

1. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon Press, Oxford, 1980.
2. E. A. Wood, *Crystals and Light: An Introduction to Optical Crystallography*, Van Nostrand, Princeton, 1964.
3. J. F. Nye, *Physical Properties of Crystals*, Oxford University Press, Oxford, 1985.
4. J. N. Hodgson, *Optical Absorption and Dispersion in Solids*, Chapman & Hall, London, 1970.
5. F. Wooten, *Optical Properties of Solids*, North Holland, Amsterdam, 1972.
6. F. Abeles (ed.), *Optical Properties of Solids*, North Holland, Amsterdam, 1972.
7. M. Balkanski (ed.), *Optical Properties of Solids*, North Holland, Amsterdam, 1972.
8. G. R. Fowles, *Introduction to Modern Optics*, rev. 2d ed., Dover, Mineola, 1989.
9. B. O. Seraphin, (ed.), *Optical Properties of Solids: New Developments*, North Holland, Amsterdam, 1976.
10. B. E. A. Saleh and M. C. Teich, *Fundamentals of Photonics*, Wiley, New York, 1991.
11. E. Yariv and P. Yey, *Optical Waves in Crystals*, Wiley, New York, 1983.
12. P. Yey, *Optical Waves in Layered Media*, Wiley, New York, 1988.
13. M. Fox, *Optical Properties of Solids*, Oxford University Press, Oxford, 2001.
14. R. Loudon, "The Raman Effect in Crystals," *Adv. Phys.* **13**:423 (1964).

PHOTONIC BANDGAP MATERIALS

Pierre R. Villeneuve*

*Department of Physics
Massachusetts Institute of Technology
Cambridge, Massachusetts*

9.1 GLOSSARY

a	lattice constant of the periodic structure
c	speed of light in vacuum
E	energy
\mathbf{E}	electric field
f	frequency
f_0	center frequency of the cavity resonance
\mathbf{H}	magnetic field
\mathbf{k}	wave vector
L	cavity length
n	index of refraction
P	power
Q	quality factor
\mathbf{r}	position vector
V_m	modal volume
Δf	frequency width of the cavity resonance
ϵ	macroscopic dielectric function
η	enhancement factor of the spontaneous emission rate
λ	wavelength in vacuum
ω	angular frequency
Θ	differential operator

*Current address: MIT Venture Mentoring Service, Massachusetts Institute of Technology, Cambridge, Massachusetts.

9.2 INTRODUCTION

Electromagnetic waves are known to undergo partial reflection at dielectric interfaces. The magnitude of the reflection is a function of the wave polarization, angle of incidence, and refractive index of the materials at the interface. Inside quarter-wave stacks, electromagnetic waves undergo reflection at multiple interfaces. The multiple reflections can lead to the destructive interference of the waves and the formation of bands of forbidden electromagnetic states. If the frequency of an electromagnetic wave lies inside such a forbidden band, the wave is prevented from propagating inside the stack; instead it is reflected and its amplitude decays exponentially through successive layers.

The operational principle behind fiber Bragg gratings,¹ interference filters,² and distributed feedback (DFB) lasers,³ is also based on multiple reflections that occur inside periodic dielectric materials. The range of frequencies over which waves are reflected (i.e., over which wave propagation is forbidden) defines a stop band, or bandgap, the width of which is proportional to the grating strength (i.e., to the effective index contrast between the different materials). The range is typically less than 1 percent of the midgap frequency, and in some cases much less than 1 percent.

In addition to being forbidden over a small range of frequencies, propagation in dielectric stacks is also forbidden over a small range of angles from normal incidence. This small range of angles defines a cone with its principal axis normal to the surface. Light incident at an angle outside the cone is not reflected, but rather is transmitted through the stack. To increase the angle of the reflection cone, one can increase the index contrast between the different dielectric layers. The cone can be made to extend as far as 90° , allowing light to be reflected off the stack from *any* angle of incidence.⁴ Stacks that reflect light from every direction are referred to as *omnidirectional reflectors*.

The existence of omnidirectional reflectors does not necessarily imply the existence of omnidirectional bandgaps. In fact, omnidirectional reflectors do not have complete three-dimensional (3D) bandgaps. Electromagnetic states exist inside the reflectors at every frequency, but incident light cannot couple to them; the wave vector of the incident light cannot be matched to the wave vector of the electromagnetic states inside the reflector. However, if light were to be generated from *within* the reflector, light could propagate along the dielectric planes—hence the absence of a three-dimensional bandgap.

In order to create a complete three-dimensional bandgap and prevent light from propagating anywhere inside the material, periodic structures must possess a three-dimensional periodicity. The principal feature of three-dimensional (3D) bandgap materials is their ability to eliminate the density of electromagnetic states everywhere inside the materials over a given range of frequencies. Since the rate of spontaneous radiative decay of an atom or molecule scales with the density of allowed states at the transition frequency, photonic bandgap (PBG) materials can be used to greatly affect the radiative dynamics of materials and lead to significant changes in the properties of optical devices.

In addition to affecting the radiative properties of atoms, PBG materials can also be used to control the flow of light by allowing certain states to exist within the bandgap. This feature has triggered the imagination of many researchers as it promises to enable the very large-scale integration of photonic components.

Though three-dimensional PBG materials can completely suppress the density of states, some three-dimensional structures possess partial gaps (i.e., gaps that do not extend along every direction). These pseudogaps can lead to small (but nonzero) densities of states and to significant changes in the radiative properties of materials. Moreover, dielectric stacks, in effect one-dimensional periodic structures, can reduce the density of states by suppressing states with wave vectors normal to the layers but they cannot eliminate every state along every direction.

In this chapter, we discuss the radiative properties of emitters and the control of light flow in PBG materials with pseudogaps and complete gaps. An in-depth review of PBG materials can be found in Ref. 5. Early fabrication efforts of 3D PBG materials are described in Ref. 6; a review of PBG materials at near-infrared frequencies is presented in Ref. 7.

9.3 MAXWELL'S EQUATIONS

Although the word *photon* is used, the appearance of bandgaps arises from a strictly classical treatment of the problem. The properties of PBG materials can be determined from the classical vector wave

equation with a periodic index of refraction. If the fields are expanded in a set of harmonic modes, in the absence of external currents and sources, Maxwell's equations can be written in the following form:

$$\nabla \times \left[\frac{1}{\epsilon(\mathbf{r})} \nabla \times \mathbf{H}(\mathbf{r}) \right] = \frac{\omega^2}{c^2} \mathbf{H}(\mathbf{r}) \quad (1)$$

where $\mathbf{H}(\mathbf{r})$ is the magnetic field, $\epsilon(\mathbf{r})$ is the macroscopic dielectric function equal to the square of the index of refraction, \mathbf{r} is the position vector, ω is the angular frequency, and c is the speed of light in vacuum. The macroscopic dielectric function has a periodic spatial dependence. Equation (1) is an eigenvalue problem; it can be rewritten as

$$\Theta \mathbf{H}_i = \frac{\omega_i^2}{c^2} \mathbf{H}_i \quad (2)$$

where

$$\Theta = \nabla \times \frac{1}{\epsilon(\mathbf{r})} \nabla \times \quad (3)$$

is a periodic Hermitian differential operator and ω_i^2/c^2 is the i th eigenvalue. The solutions \mathbf{H}_i and ω_i are determined entirely from the strength and symmetry properties of $\epsilon(\mathbf{r})$. The solutions are characterized by a wave vector \mathbf{k} and a band number i . The region of all allowed wave vectors is called a *Brillouin zone*, and the collection of all solutions is termed a *band structure*.

Equation (2) closely resembles Schrödinger's equation for the problem of an electronic wave function inside a periodic atomic potential. Since the solutions of Schrödinger's equation lead to band diagrams for allowed and forbidden electronic states in crystalline structures, and since PBG materials have similar effects on electromagnetic waves, PBG materials are often referred to as *photonic crystals*. In this chapter, the terms *PBG material* and *photonic crystal* are used interchangeably.

An interesting feature of Eq. (1) is that there is no fundamental constant with dimensions of length, hence no fundamental length scale other than the assumption that the system is macroscopic. The solution at one length scale determines the solutions at all other length scales, assuming a frequency-independent dielectric function. This simple fact is of considerable practical importance as it allows results to be scaled from one wavelength to another, from the ultraviolet to microwaves and beyond, simply by expanding all distances.

The solutions \mathbf{H}_i and ω_i provide information about the frequency of the allowed electromagnetic modes in a PBG structure and their polarization, symmetry, and field distribution. Although Eq. (1) can be applied to any dielectric structure—the only assumptions made were the absence of external currents and sources—early work in this field focused on the search for a complete bandgap, that is, a range of frequencies with no allowed electromagnetic mode for any wave vector \mathbf{k} inside the Brillouin zone.⁶ A review of three-dimensional photonic crystals follows.

Several numerical methods have been used to solve Maxwell's equations in periodic structures, including the use of a variational approach⁸ where each eigenvalue in Eq. (2) is computed separately by minimizing the functional $\langle \mathbf{H}_i | \Theta | \mathbf{H}_i \rangle$. In this method, fast Fourier transforms are used repeatedly to switch back and forth between real and reciprocal space to avoid storing large matrices.

Other methods include the transfer matrix method⁹ and the finite-difference time-domain (FDTD) method,¹⁰ to name but two. In the former, Maxwell's equations are solved at a fixed frequency by stepping the fields forward in space, one plane at a time, satisfying the continuity conditions at every step. The transfer matrix method is well-suited for transmission and reflection computations in photonic crystals. By imposing Bloch conditions, the transfer matrix method can also be used to compute the band structure. In the case of the FDTD method, Maxwell's equations are discretized on a three-dimensional grid, and the derivatives are approximated at each grid point by a corresponding centered difference. Maxwell's equations are solved everywhere in the computational cell at every time step, allowing the temporal response of the fields to be determined inside photonic crystals.

9.4 THREE-DIMENSIONAL PHOTONIC CRYSTALS

Criteria for 3D Bandgaps

The existence of bandgaps in periodic structures is determined entirely from the symmetry and strength of the periodic dielectric function. Since photonic crystals do not occur naturally, somehow one must arrange dielectric material in a 3D periodic structure, and, as with multilayer dielectric stacks, the length of the repeating unit must be on the order of one-half the wavelength in the material. Most structures exhibiting 3D bandgaps satisfy the following three general criteria: the periodic structure has a spherelike Brillouin zone; the refractive index contrast between the different materials is typically larger than 2; and the high- and low-dielectric materials form connected networks.

Spherelike Brillouin Zone Waves propagating inside a periodic structure sense a periodicity that leads to the formation of stop bands at the edges of the irreducible Brillouin zone. Since the waves sense a different periodicity along the different directions, the wave vectors at the different points on the surface of the Brillouin zone have different magnitudes. Hence, the gaps are likely to be centered at different frequencies. Spherical Brillouin zones (if they were possible) would guarantee the overlap of all the gaps along every direction, since every point on the surface of a sphere is equidistant from the center—but crystal geometries do not allow for spherical Brillouin zones.

Several hundred years of mineralogy and crystallography have led to the classification of the various three-dimensionally periodic lattice geometries. The Brillouin zone of the face-centered-cubic (fcc) lattice is closer to a sphere than any other common crystal geometry. However, despite having the most spherelike Brillouin zone, the farthest point on the surface of the fcc Brillouin zone (i.e., the point with the largest wave vector, the so-called **W** point) lies 29 percent farther from the origin than the closest point, the **L** point. For a gap to open along every direction, the gaps at **W** and **L** must be large enough to overlap.

Large Index Contrast The size of the bandgap at each point on the surface of the Brillouin zone scales with the index contrast between the different materials. For the different gaps to overlap over the entire Brillouin zone the refractive index contrast must be large, typically 2 to 1 or greater. Semiconductor materials such as Si ($n = 3.5$ at $\lambda = 1.5 \mu\text{m}$) and GaAs ($n = 3.4$) in combination with air or low-index oxides are excellent candidates for the fabrication of photonic crystals at infrared wavelengths.

A large index contrast and a spherelike Brillouin zone, however, are not sufficient to guarantee the formation of a bandgap in 3D structures. It is not sufficient to specify the structure in reciprocal space—there are essentially an infinite number of structures with an fcc lattice, since anything can be put inside the fundamental repeating unit. One must also specify the dielectric structure in real space. An example of a successful 3D photonic crystal is shown in the forthcoming section labeled “Examples of 3D Crystals.”

Connected Networks To appreciate the importance of having a connected network, it is useful to consider a one-dimensionally periodic structure such as a multilayer dielectric stack. The energy density of the mode below the stop band is more strongly localized in the high-index layers than the mode above the stop band. The more strongly the energy density of the lower mode is localized in the high-index material and the more strongly the energy density of the upper mode is localized in the low-index material, the larger the bandgap.

In 3D periodic structures, it is generally advantageous for the high-index material to be fully connected to allow the electric field of the mode in the lower band to run through the high-index material as much as possible without having to go through the low-index material. One should be able to connect any point in the high-index material to any other point without having to cross over into the low-index material. The same also holds for the low-index material. Moreover, the low-index material should occupy typically over 50 percent of the total volume. A detailed discussion of the nature of bandgaps in periodic structures is given in Refs. 5 and 11.

The three general criteria just presented should serve only as guidelines. They do not constitute necessary conditions for the creation of 3D bandgaps. For example, though the fcc lattice is the most spherelike, other lattice geometries have been shown to generate 3D bandgaps.

Examples of 3D Crystals

The earliest antecedent to photonic bandgaps is the observation by Sir Lawrence Bragg of narrow stop bands in crystals from x-ray diffraction. The refractive index contrast, however, was very small, typically less than 1.001 to 1, and produced only narrow rings on the surfaces of the Brillouin zone.

The first structure with a full 3D bandgap was discovered by K. M. Ho et al. in 1990 and consisted of a diamond lattice of air spheres (i.e., an fcc lattice with two air spheres per unit cell) inside a high-index material.¹² Since then, there has been considerable effort to develop a process for the manufacturing of diamond (or diamondlike) structures at micrometer wavelengths. One such approach consists of etching a large number of hole triplets at off-vertical angles in a slab;^{13,14} another consists of building an orderly stacking of dielectric rods;¹⁵ yet another consists of etching a series of horizontal grooves into sequentially grown layers and etching vertical holes.¹⁶ These structures are variations of the same diamond lattice grown along either the (1, 1, 1), (0, 0, 1), or (1, 1, 0) directions, respectively.

An example of the structure grown along the (0, 0, 1) direction is shown in Fig. 1. It consists of multiple layers of polycrystalline silicon rods with a stacking sequence that repeats itself every four layers. Within each layer, the rods are parallel to each other; the rods are shifted by half a period every other layer. Only five layers are shown. The structure was fabricated by S. Y. Lin et al. at Sandia National Laboratories in 1998 using a process that involves the repetitive deposition and etching of multiple dielectric films.¹⁷ The width of each rod is roughly 1.2 μm . The bandgap is centered at a wavelength of 10 μm . In addition to fabricating this structure, the researchers also fabricated a structure at shorter wavelengths centered at $\lambda = 1.5 \mu\text{m}$.⁷

An overview of the fabrication of 3D PBG materials at micrometer and submicrometer length scales can be found in Ref. 7.

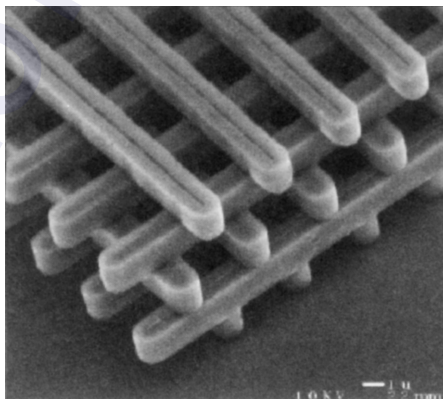


FIGURE 1 Scanning electron micrograph of a three-dimensional photonic crystal built at Sandia National Laboratories. The crystal consists of five layers of polycrystalline silicon rods. The width of the rods is 1.2 μm . The photonic bandgap is centered around a wavelength of 10 μm .

9.5 MICROCAVITIES IN THREE-DIMENSIONAL PHOTONIC CRYSTALS

From Fermi's golden rule, we know that the rate of spontaneous radiative decay of an atom scales with the density of allowed states at the atomic transition frequency. In free space, the density of states scales quadratically with frequency, and the probability of finding an atom in an excited state simply decays exponentially with time.

The introduction of boundaries in the vicinity of the atom has the effect of changing the density of allowed states. For example, in the case of a bounded system with reflecting walls—such as a laser cavity—the density of states is reduced to a spectrally discrete set of peaks, each corresponding to a resonant longitudinal mode of the cavity. When no mode falls within the emission linewidth of the atomic transition, atomic radiative decay is essentially suppressed. However, if the transition frequency overlaps one of the resonant frequencies, the density of available modes for radiative decay becomes very large, which in turn enhances the rate of spontaneous emission. In conventional solid-state lasers, several modes fall within the atomic emission linewidth. The free spectral range of the modes is given by $c/2nL$, where n is the refractive index of the host material and L is the distance between the reflectors. If L was made very small, it would be possible to increase the mode spacing such that only one (or even zero) mode would fall within the emission linewidth.

An example of a small laser cavity is the distributed feedback (DFB) laser consisting of a spatially corrugated waveguide with a quarter-wave phase shift. The phase shift defines a cavity, and the grating on either side acts like a mirror. The length L of the cavity is characterized by the decay length of the evanescent field along the axis of the grating and typically extends over hundreds of wavelengths in the material. The grating creates a stop band along the periodic axis. While the absence of longitudinal modes inside the stop gap reduces the total density of states, the presence of a quarter-wave phase shift generates a resonant mode inside the gap and increases the density of states. The increase is sufficiently large to allow single-mode action of the laser at the resonant frequency. Though DFB lasers have longitudinal stop bands, the total density of states is not zero, since the stop band extends only inside a small cone along one direction. Leaky radiation modes exist along every other direction.

3D PBG materials have the ability to open 3D stop bands that reflect light along every direction in space and that completely eliminate the density of states for a given range of frequencies. In the case where the radiative transition frequency of an atom falls within the frequency gap of the crystal, spontaneous radiative decay is essentially suppressed.

If a small defect (or phase shift) is introduced in the photonic crystal, a mode can be created within the structure at a frequency that lies inside the gap. If the size of the defect is such that it supports a mode, the defect behaves like a microcavity surrounded by reflecting walls. If the radiative transition frequency of the atom matches that of the defect mode, the rate of spontaneous emission can be enhanced.

Figure 2 shows the vector plot of a resonant mode in a 3D photonic crystal similar to the one shown in Fig. 1. The defect is located at the center of the crystal and consists of a broken high-index rib. (The defect could be introduced, for instance, in one of the layers during the growth of the crystal.) The electric field is shown in the vertical plane through the middle of the defect. The mode is strongly localized in all three dimensions, and its amplitude falls off sharply away from the defect. The electric field *jumps* from one edge of the broken rib to the other, while the magnetic field (not shown) has the shape of a torus and runs around the electric field. The frequency of the mode is $f = 0.59c/a$, where a is the lattice constant (i.e., the length of the repeating unit cell) of the crystal. In this particular example, the high-index material has a refractive index of 3.4; the low-dielectric material has an index of 1.0; and the gap extends from $f = 0.52c/a$ to $0.66c/a$.

In contrast to defects in one-dimensional periodic structures (such as DFB lasers), arbitrarily small defects in 3D crystals do not necessarily lead to the creation of localized modes. The volume of the defect must reach a certain threshold to sustain a resonant mode. Furthermore, quarter-wave shifts in DFB lasers lead to resonant modes at the center of the gap. There is no simple equivalent in 3D crystals.

The frequency of the resonant mode changes with the size and shape of the defect. The simple action of adjusting the defect size provides tunability of the resonant mode and affects the localization

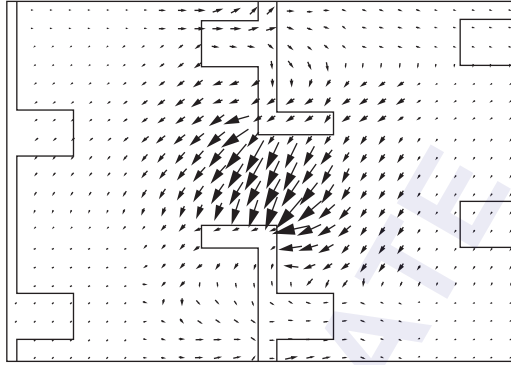


FIGURE 2 Vector plot of the electric field in a 3D PBG with a defect. The overlay indicates the edges of the high-dielectric material. The defect, located at the center of the figure, is fabricated by breaking one of the dielectric ribs. The defect supports a localized resonant mode inside the crystal.

strength. The field attenuation through successive unit cells is stronger for modes lying near the center of the gap than for those lying near the edges.

Although the microcavity in the just-noted example was created by removing part of a high-index rib, a cavity could equally have been created either by adding material between ribs or by changing the shape of one or more ribs. Also, multiple high-order localized modes may appear inside the crystal as the size of the defect is made bigger.

Quality Factor

One important aspect of microcavities in finite-sized crystals is the quality factor Q of the resonator defined as:¹⁸

$$Q = \frac{2\pi f_0 E}{P} = -\frac{2\pi f_0 E}{dE/dt} \quad (4)$$

where f_0 is the resonant frequency, E is the energy stored inside the resonator, and $P = -dE/dt$ is the dissipated power. Hence, a resonator can sustain Q oscillations before its energy decays by a factor of $e^{-2\pi}$ (i.e., a reduction of 99.8 percent) of its initial value. In the specific case where the line-shape of the resonance is a Lorentzian, Eq. (4) reduces to $f_0/\Delta f$, where Δf is the width of the resonance.

Since the quality factor is a measure of the optical energy stored in the microcavity over the cycle-average power radiated out of the cavity, Q is expected to be largest for modes lying near the center of the gap where the field attenuation is strongest. Q is also expected to increase with the size of the crystal, since the reflectivity increases with the number of periods (i.e., the leakage from the edges of the crystal becomes progressively smaller). The quality factor of the mode shown in Fig. 2 is plotted in Fig. 3 as a function of the size of the crystal. The quality factor is computed using the finite-difference time-domain method described in Sec. 9.3. First the resonant mode is excited and the total energy is monitored as a function of time. Then the time required for 99.8 percent of the energy to escape is recorded. Results are shown for crystal sizes of dimension $2N \times 2N \times 2N$. In each case, the defect is surrounded by N unit cells along every direction. Q increases exponentially with the size of the crystal and reaches a value close to 10^4 with as little as four unit cells on either side of the defect. The steepness

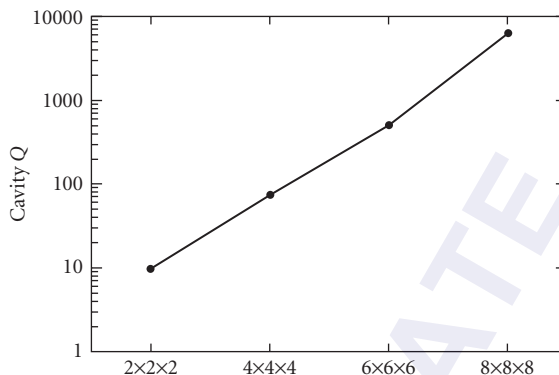


FIGURE 3 Quality factor of the resonant cavity shown in Fig. 2 as a function of the size of the 3D photonic crystal, given in units of cubic lattice constants.

of the slope in Fig. 3 follows directly from the field attenuation through each successive lattice of the crystal. Since the only energy loss in the structure occurs from tunneling through the walls of the finite-sized crystal (i.e., intrinsic losses due to material absorption is not considered), Q does not saturate even for a large number of unit cells. A more detailed description of the properties of resonant modes in photonic crystals can be found in Ref. 19.

Enhancement of Spontaneous Emission

By coupling an optical transition to the microcavity resonance, the spontaneous emission rate can be enhanced by a factor η over the rate without a cavity. The expression for η is given by:²⁰

$$\eta = \frac{2Q}{\pi V_m} \left(\frac{\lambda}{2n} \right)^3 \quad (5)$$

where V_m is the modal volume, n is the refractive index of the medium, and λ is the free-space wavelength of the optical transition. Photonic crystals have the ability to enhance the rate of spontaneous emission by enabling microcavities with large quality factors and small modal volumes. In the case where the modal volume is on the order of a cubic half-wavelength in the material [i.e., $V_m \sim (\lambda/2n)^3$], the enhancement factor is on the order of Q . A detailed example is provided in the following section.

9.6 MICROCAVITIES IN PHOTONIC CRYSTALS WITH TWO-DIMENSIONAL PERIODICITY

Three-dimensional field confinement can be achieved in dielectric structures, in part by the effect of a photonic bandgap and in part by index confinement. An example was given in Sec. 9.5 for the case of a DFB laser (i.e., a structure with a one-dimensional periodicity). One important aspect of structures with dimensional periodicity lower than three is the coupling to radiation modes. By reducing the dimensionality of the periodicity and by resorting to standard index guiding to confine light along the nonperiodic direction(s), one no longer has the ability to contain light completely, and leaves open possible decay pathways through which light can escape.

In this section, we consider a dielectric slab waveguide with a two-dimensional periodic lattice. The periodic lattice is used to confine light in the plane of the waveguide (the xy -plane, say), and the slab keeps the light from escaping along the transverse direction (the z -direction). It is useful to begin with a uniform waveguide, and consider the effect of adding a periodic array of holes. The slab is chosen to have a large refractive index ($n = 3.4$) and, for simplicity, is assumed to lie in air. The thickness of the slab is set equal to $0.5a$, where a is a scaling parameter as defined in the text that follows. The use of a high-index waveguide is twofold: first, the high index provides strong field confinement along the z -direction (i.e., the extent of the guided modes outside the waveguide is small), allowing a large fraction of each mode to interact with the photonic crystal; and second, the high-index contrast between the dielectric material and the holes will increase the likelihood of having a bandgap in the xy -plane.

The waveguide is shown in Fig. 4a. Its corresponding dispersion relation is shown in Fig. 4b. The solid lines correspond to guided modes, and the shaded region corresponds to the continuum of radiation (i.e., nonguided) modes. The guided modes are labeled transverse electric (TE) and transverse magnetic (TM) with respect to the xy -plane of symmetry in the middle of the waveguide. TE (TM) modes are characterized by the absence of electric field components in the z (x and y) direction at the center of the waveguide.

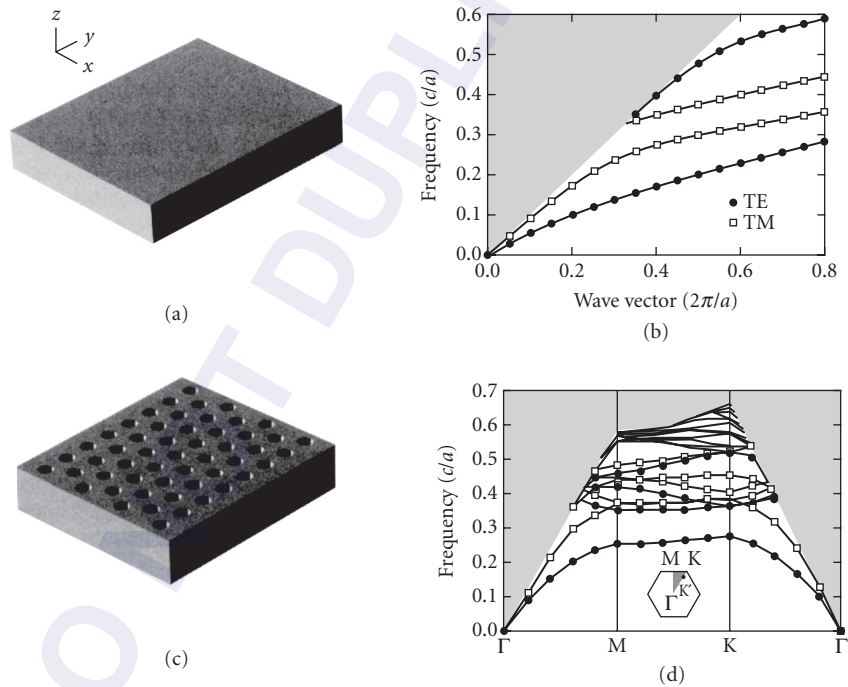


FIGURE 4 (a) Schematic diagram of a dielectric slab waveguide of thickness $0.5a$ and refractive index 3.4. (b) Band diagram of the slab waveguide shown in (a). The solid lines correspond to guided modes; the shaded region corresponds to the continuum of radiation modes. The guided modes are labeled TE or TM with respect to the xy -plane of symmetry in the middle of the slab, (c) Schematic diagram of a slab waveguide with a two-dimensional triangular array of holes with radius $0.3a$, where a is the lattice constant of the periodic array. The parameters of the slab are identical to those in (a). (d) Band diagram for the slab waveguide shown in (c). Only the lowest nine bands are labeled TE and TM. Guided modes do not exist above the cut-off frequency of $0.66c/a$. The inset shows the Brillouin zone and symmetry points for a triangular lattice, with the irreducible zone shaded.

The dispersion relation shown in Fig. 4b extends to the right of the figure; there is no upper bound on the wave vector. The introduction of a periodic array of holes in the waveguide has the effect of folding the dispersion relation into the first Brillouin zone and splitting the guided-mode bands. Figure 4c shows a waveguide with a triangular array of holes. The holes have a radius of $0.30a$, where a is the lattice constant of the array. The associated dispersion relation is shown in Fig. 4d. Again, the shaded region above the light line corresponds to the continuum of radiation modes. The solid lines below the light line correspond to guided modes. These modes remain perfectly guided in spite of the holes and propagate in the waveguide with no loss. A bandgap can be seen between the first and second TE bands. An experimental observation of bandgaps in this type of structure is described in Ref. 21.

The introduction of holes in the waveguide also creates a frequency cutoff for guided modes. Every mode above the frequency $0.66c/a$ is folded into the radiation continuum, and is Bragg-scattered out of the slab. The cutoff frequency is independent of the refractive index of the slab or the size of the holes, and depends only on the lattice geometry of the array of holes.

If a defect is introduced in the PBG structure shown in Fig. 4c, localized modes can be formed in the vicinity of the defect. Since each localized mode has a specific polarization, it is possible to create a TE mode between the first and second TE bands, orthogonal to TM modes. If, for example, light were to originate from a quantum well located at the middle of the waveguide, atomic transitions could be made to couple only to TE modes.

Two competing decay mechanisms contribute to the overall decay rate of the localized mode; horizontal in-plane coupling to guided modes at the edges of the crystal in the unperturbed (i.e., holeless) waveguide, and vertical coupling to radiation modes. For some applications (such as photonic integrated circuits) it may be preferable for the localized mode to decay primarily into guided modes, while for other applications (such as off-chip emission) it may be preferable for the mode to decay primarily into the radiation continuum. These two cases are considered separately in the following text.

The total quality factor of the resonant mode, Q_{tot} , is given by:²²

$$\frac{1}{Q_{\text{tot}}} = \frac{1}{Q_{\text{wg}}} + \frac{1}{Q_{\text{rad}}} \quad (6)$$

where $1/Q_{\text{wg}}$ is a measure of the coupling to waveguide modes and $1/Q_{\text{rad}}$ is a measure of the coupling to radiation modes. The strength of the two competing coupling mechanisms depends on the size of the crystal (i.e., the total number of holes around the defect), the modal volume, and the choice of substrate.

In-Plane Coupling

We present the case of an array of 45 holes with a missing hole at the center (i.e., one hole is filled). The structure supports a localized mode inside the TE bandgap. The total quality factor of the mode is computed using the finite-difference time-domain method described in Sec. 9.3 and is found to be 240. The modal volume, V_m , is defined as:²³

$$V_m = \frac{\int \epsilon(\mathbf{r}) |\mathbf{E}(\mathbf{r})|^2 d^3\mathbf{r}}{(\epsilon(\mathbf{r}) |\mathbf{E}(\mathbf{r})|^2)_{\text{max}}} \quad (7)$$

where $\mathbf{E}(\mathbf{r})$ is the electric field distribution of the mode. The computed modal volume is only three cubic half-wavelengths in the material. The spontaneous emission rate enhancement factor, computed from Eq. (5), is equal to 50.

Since the structure does not have a complete three-dimensional bandgap, Q_{tot} cannot be made arbitrarily large. While the addition of extra holes would reduce the coupling to the guided modes outside the crystal, light could not be prevented from coupling to radiation modes. Any significant

increase in the number of holes would cause the mode to primarily radiate outside of the waveguide. Moreover, coupling to radiation modes would be enhanced if the waveguide was positioned on a substrate. The substrate would provide a favorable pathway for radiation loss. It has been shown, however, that the adverse effects of a substrate could be minimized with the use of a low-index insulating layer between the waveguide and the substrate.^{24,25}

The coupling to radiation modes is also enhanced by reducing the modal volume. The more tightly a mode is confined, the more likely it is to radiate out of the waveguide. Conversely, if the modal volume is made larger, the coupling to radiation modes can be reduced, and, provided the coupling to guided modes remains largely unchanged, Q_{tot} can be increased. To increase the modal volume, one could create a different type of defect in the structure. If, instead of removing a single hole from the two-dimensional array, the radius of seven nearest-neighbor holes was reduced from $0.3a$ to $0.2a$ while otherwise leaving the structure unchanged, the localized mode would become more extended—the modal volume would increase by 20 percent to $3.6(\lambda/2n)^3$ —and Q_{tot} would increase by more than one order of magnitude to 2500. The frequency of the new localized mode would remain unchanged, and the enhancement factor would exceed 400.

Out-of-Plane Coupling

While it may be possible to fabricate high- Q cavities that couple predominantly to guided modes, some applications (such as light-emitting diodes) may require a large fraction of the emitted light to be extracted from the high-index guiding layer. As mentioned previously, the emitted radiation can be made to decay primarily into radiation modes by increasing the total number of holes surrounding the defect. In this case, Q_{wg} would essentially be infinite, and $Q_{\text{tot}} \sim Q_{\text{rad}}$. For simplicity, in this example, we write $Q_{\text{tot}} = Q_{\text{rad}} = Q$.

Light-emitting diodes (LEDs) are widely used as incoherent light sources in applications such as lighting, displays, and short-distance fiber communications. Two important performance characteristics of LEDs are the output efficiency (i.e., the amount of light extracted from the structure for a given injection current) and the modulation rate (i.e., the information emission capacity).

Photonic crystals with two-dimensional periodicity can lead to the enhancement of the rate of spontaneous emission and consequently to higher modulation rates. However, photon reabsorption and nonradiative recombination can affect the performance of LEDs by reducing the extraction efficiency and the modulation rate. High- Q cavities, though seemingly favorable for the enhancement of the rate of spontaneous emission, may cause severe reabsorption in certain material systems, since the likelihood of observing photon reabsorption increases with the photon lifetime inside the cavity.

Display Applications For display applications, it is usually desired to get as much light as possible out of the high-index material over the entire spontaneous emission bandwidth for a constant applied current. If all emitted frequencies fall inside the guided-mode bandgap, all available optical modes can contribute to the output signal. In the ideal case where there are no nonradiative recombination processes, the extraction efficiency is unity; every photon escapes from the high-index waveguide. Even photons reabsorbed by the atomic system, if given enough time, eventually get reemitted and contribute to the output signal. However, when nonradiative processes are present, reabsorbed photons can be lost. In order to achieve high output efficiency, the effective spontaneous emission rate—the spontaneous emission rate reduced by photon reabsorption—has to dominate over the nonradiative recombination rate. The relative rate of the radiative and nonradiative processes can be controlled by modifying the quality factor of the cavity.

Two limit cases are identified: the case where photon reabsorption is negligible, and the case where it is important. The former arises in certain organic emitters, where the energy levels of the molecules are such that absorption and spontaneous emission are spectrally separated. The latter arises in most semiconductor systems, where both absorption and emission processes occur between the conduction and valence bands.

In the case of low reabsorption, if the cavity linewidth is larger than the emission linewidth, an increase of the cavity Q can result in an increase of the effective spontaneous emission rate and of

the output efficiency. However, a reduction of the cavity linewidth beyond the material emission linewidth does not further enhance the spontaneous emission rate or output efficiency. In the case of large reabsorption, the rate of spontaneous emission and the output efficiency reach a maximum when the cavity linewidth is comparable to the material linewidth, but fall to zero when the cavity linewidth becomes much smaller than the material linewidth. A more detailed description of these conditions can be found in Ref. 26.

Communications Applications For communications applications, it is advantageous to reduce the emission linewidth below the material emission linewidth to improve the temporal coherence of the emitted light. It is also advantageous to increase the modulation speed to improve the information emission capacity. If a time-varying current is applied to the LED, the response time of the electron-photon system will be determined by the slowest of the different relaxation processes.

While electronic recombination lifetimes are typically on the order of a few nanoseconds in both semiconductors and organic dyes, the photon lifetime in a cavity depends on the cavity Q and, in the case where, say, $Q = 1000$, is on the order of several picoseconds. Since the modulation speed is limited by the slower of the two processes, the electronic recombination rate, which is a sum of the effective spontaneous emission rate and the nonradiative recombination rate, constitutes the limiting factor. To achieve high modulation speeds, it is therefore necessary to increase the spontaneous emission rate.

In the case where photon reabsorption is small, such as in organic dyes, the rate of spontaneous emission and the modulation speed increase with the cavity Q . Conversely, when photon reabsorption is large, such as in semiconductors, the maximum rate of spontaneous emission and the maximum modulation rate are achieved when the cavity linewidth is comparable to the material linewidth. These conclusions are similar to those found for display applications.

Examples of Low- Q and High- Q Cavities Low- Q cavities can be fabricated in high-index dielectric waveguides by introducing an array of holes with no defects. The absence of defects ensures that photons inside the waveguide—emitted from a quantum well, say—spend as little time as possible inside the waveguide and minimize the risk of being reabsorbed. The cavity is defined by the waveguide itself, which provides vertical field confinement. To avoid removing active material, the holes can be made to extend only partly into the guiding layer so as to not penetrate into the quantum well. Bandgaps for guided modes can be generated even when the holes do not extend through the entire thickness of the waveguide. The waveguide can also be positioned on a dielectric or metallic mirror to ensure that the output radiation escapes through the top surface.

High- Q cavities can be fabricated in structures similar to those used for low- Q cavities except that, in the case of high- Q cavities, defects are introduced in the periodic array. The introduction of defects creates highly confined modes in the area of the defects, hence only a small fraction of the quantum well overlaps with the resonant modes (i.e., only a fraction of the electron-hole pairs contributes to the emitted signal). To eliminate this problem, high- Q cavities could be generated by placing the dielectric layer between two vertical Bragg mirrors—in analogy to resonant-cavity LEDs—and by getting rid of the defects. The entire active region would then overlap with the resonant cavity mode.

Experimental results of quantum-well emitters and dyes in photonic crystals with two-dimensional periodicity can be found in Refs. 27 through 31. A detailed analysis of the output efficiency and modulation rate of LEDs can be found in Refs. 26 and 32.

9.7 WAVEGUIDES

While three-dimensional field confinement can be achieved by introducing local-point defects in photonic crystals, two-dimensional field confinement can be achieved by introducing extended line defects. Both point defects and line defects can generate localized modes with

frequencies that lie inside the bandgap. However, unlike point defects, line defects can generate modes that propagate along the lines with nonzero group velocity. Line defects can be made, for example, by *carving* channels in photonic crystals or by creating line dislocations. Electromagnetic waves propagating along the lines are guided not from total internal reflection but from the bandgap effect; they are prevented from leaking into the crystal since their frequencies lie inside the bandgap.

The absence of radiation modes in three-dimensional photonic crystals suggests that it may also be possible to create waveguides with very sharp bends. Since electromagnetic waves are prevented from propagating inside photonic crystals, the waves would only either propagate through the bend or be reflected back. It will be shown in the subsection labeled “Waveguide Bends” that, for certain frequencies, reflection may be eliminated altogether, leading to complete transmission. In three-dimensional crystals, waveguide bends could extend along any direction and could be used for the implementation of interconnected integrated optical circuits on multiple planes. In this chapter, however, we focus only on line defects in photonic crystals with two-dimensional periodicity.

Waveguides in Photonic Crystals with Two-Dimensional Periodicity

As we saw in Sec. 9.6, photonic crystals with two-dimensional periodicity rely on the existence of bandgaps to control propagation in the plane and on index guiding to confine electromagnetic fields along the third dimension. An example of a photonic crystal with two-dimensional periodicity was shown in Fig. 4c; its corresponding dispersion relation was shown in Fig. 4d. In the bandgap, no *guided* mode existed for TE polarization.

In this section, a line defect is introduced in the photonic crystal shown in Fig. 4c by increasing the radius of a line of nearest-neighbor holes along the Γ -K direction from $0.30a$ to $0.45a$. The resulting dispersion relation is shown in Fig. 5. The dispersion relation is computed using the plane-wave expansion method described in Sec. 9.3. The wave vector along the line defect is plotted on the abscissa.

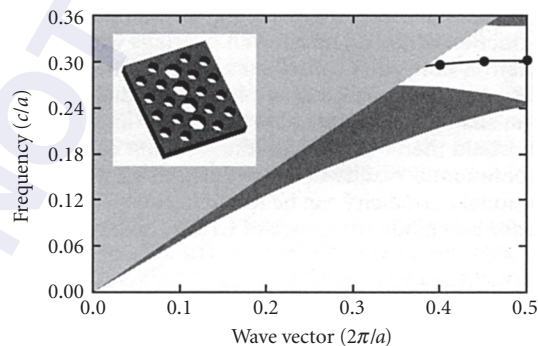


FIGURE 5 Projected dispersion relation of the TE modes in the waveguide structure shown in the inset. The dispersion relation is projected along the axis of the waveguide (i.e., along the line defined by the series of larger holes). The light gray region corresponds to the continuum of radiation modes, and the dark gray regions correspond to modes inside the bulk PBG dielectric slab. The thickness of the slab is $0.5a$, the refractive index is 3.4, the radius of the small holes is $0.3a$, and the radius of the large holes is $0.45a$. The figure is to be compared with Fig. 4d along the Γ -K direction.

In this structure, it is necessary to distinguish between the modes which are guided inside the dielectric slab (the so-called bulk crystal modes that correspond to the different bands in Fig. 4d) and the modes which are guided along the line defect. The dispersion relation is obtained from Fig. 4d by projecting the wave vector of every mode along the Γ -K direction; the dark gray regions correspond to the continuum of bulk crystal modes and the light gray region corresponds to the continuum of radiation modes. The bulk crystal modes and radiation modes are depicted with a uniform shading despite the nonuniform density of states in these regions. Since the structure retains an inherent periodicity along the line defect, the wave vector has an upper limit. However, although the line defect extends along the Γ -K direction, K is not the point at the edge of the dispersion relation. The boundary is located at the projected M point along the Γ -K direction, labeled K' as shown in the inset of Fig. 4d.

Only modes lying outside the shaded regions are truly guided along the line defect. A single guided mode appears inside the bandgap. Since the line defect consists of a series of larger holes, the effective index of the waveguide is lower than that of the surrounding photonic crystal. Hence, the mode is not index-guided in the plane; it is constrained horizontally by the bandgap. The effective index, however, is higher in the waveguide than in the regions above and below the slab, allowing the mode to be guided vertically by index confinement. The electric field of the guided mode is mostly concentrated in the dielectric material. The fraction of electric-field energy inside the high-dielectric material at K', for example, is close to 75 percent.

Alternatively, a line defect could have been created by reducing the radius of a series of holes, or by creating lattice dislocations. Also, instead of using a high-index slab with holes, one could have used an array of high-index posts. High-index posts can generate dispersion relations similar to the one shown in Fig. 4c except that the open (solid) circles would now correspond to TE (TM) polarization. A more detailed analysis of these and other similar structures can be found in Refs. 33 and 34.

Waveguide Bends

If a sharp bend is introduced in a PBG waveguide—with a radius of curvature on the order of a few lattice constants—it may be possible to obtain high transmission through the bend for a wide range of frequencies. To obtain high transmission, the waveguide must support a single mode at the frequency of interest, and the radiation losses must be small, since coupling to high-order guided modes and to radiation modes reduces the transmission and increases the reflection.

While it may be possible to obtain 100 percent transmission in photonic crystals with two-dimensional periodicity, we choose to consider waveguide bends in purely two-dimensional crystals. Two-dimensional crystals can be viewed either as flat structures in a two-dimensional Cartesian space or as structures of infinite thickness with no field variation along the vertical direction. Since there is no index confinement along the vertical direction, there are no radiation modes and no light cone. The bandgap in a 2D structure is analogous to a three-dimensional bandgap in that there are truly no modes inside the bandgap.

For simplicity, we consider a 2D photonic crystal of dielectric columns on a square lattice, surrounded by air. The refractive index of the rods is chosen to be 3.4 and the radius $0.20a$, where a is the lattice constant of the array. A large bandgap appears in this structure for TM polarization (electric field parallel to the axis of the columns). A line defect is created inside the crystal by removing a row of rods. The line defect introduces a single guided TM mode inside the gap, similar to the one shown in Fig. 5. The main difference between the dispersion relation for this 2D crystal and the one shown in Fig. 5 is the absence of radiation modes in the 2D crystal. The bandgap extends over the entire range of wave vectors. If a bend is introduced in the waveguide, light will either travel through the bend or be reflected back, since there are no radiation modes to which light can couple. Only back reflection can hinder perfect transmission.

The transmission and reflection can be studied using the finite-difference time-domain method described in Sec. 9.3. In this method, a dipole located at the entrance of the waveguide creates a pulse with a Gaussian envelope in time. The field amplitude is monitored inside the waveguide at two points, one before the bend and one after the bend. The pulses are then Fourier-transformed to

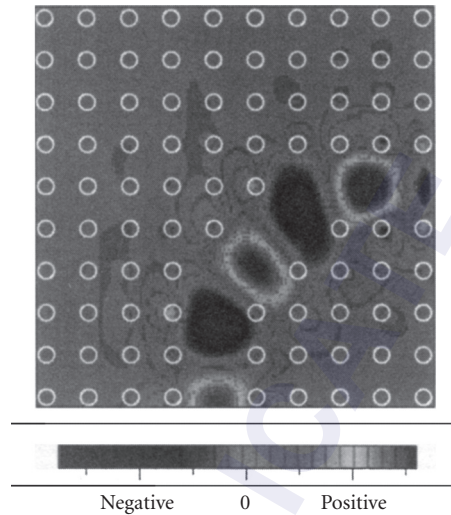


FIGURE 6 Electric field pattern of a guided mode in a photonic crystal in the vicinity of a bend. The white circles indicate the position of the high-dielectric columns. The electric field is polarized along the axis of the columns. The mode is strongly confined inside the guide and is completely transmitted through the bend. The radius of curvature of the bend is on the order of the wavelength of the guided mode.

obtain the reflection and transmission coefficients for each frequency. A detailed description of this method and computational results are presented in Ref. 35.

The electric field pattern of a mode propagating through the bend is shown in Fig. 6. The mode is strongly guided inside the photonic crystal. One hundred percent of the light travels through the bend despite a radius of curvature on the order of one wavelength.

The transmission through the bend can be modeled as a simple one-dimensional scattering process. The bend can be broken down into three separate waveguide sections: the input waveguide in the (01) direction; the output waveguide in the (10) direction; and a short waveguide section in the (11) direction, connecting the input and output waveguides. Each section supports a single guided mode with wave vector $\mathbf{k}_1(f)$ for propagation along the (01) or (10) direction, and $\mathbf{k}_2(f)$ for propagation along (11). These wave vectors are given by dispersion relations similar to the one shown in Fig. 5. The mode propagating along the (01) direction is scattered into the mode propagating along (11), then into the mode propagating along (10). At the interfaces, the fields and their derivatives must be continuous. By complete analogy with the one-dimensional Schrödinger equation for a square potential well, the transmission through the sharp bend can be mapped onto that of a wave propagating in a square *dielectric potential*. This potential consists of three constant pieces corresponding to the (01), (11), and (10) directions, respectively. The model differs from the standard one-dimensional scattering problem in that the depth of the well, determined by the difference $|\mathbf{k}_1(f)|^2 - |\mathbf{k}_2(f)|^2$, now depends on the frequency of the traveling wave. The scattering model correctly predicts the general quantitative features of the transmission spectrum obtained from the FDTD method, as well as the frequencies where the reflection coefficient vanishes.³⁵

The results have been experimentally confirmed using a structure consisting of a square array of tall circular rods.³⁶ The rods were made of alumina with a refractive index of 3.0 and a radius of 0.25 mm. The lattice constant was chosen to be 1.27 mm and the rods were close to 10 cm in length. The large aspect ratio between the length and the lattice constant provided a good approximation

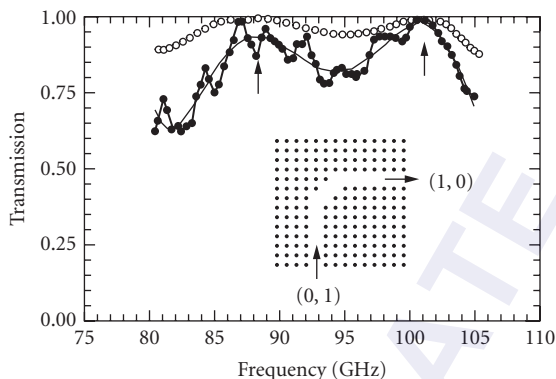


FIGURE 7 Normalized transmission spectrum for the PBG structure shown in the inset. The solid circles correspond to experimental data; the open circles are computed from the one-dimensional scattering model. Near-perfect transmission is observed through the bend near 87 and 101 GHz. The arrows indicate the positions of the reflection nodes from theory. The experimental data is fitted with a polynomial curve.

of a two-dimensional system. Because of the absence of vertical confinement, the waveguides were made to extend over less than 100 lattice constants to minimize loss in the vertical direction. The bandgap extended from 76 to 105 GHz. The experiment was carried out at millimeter-wave frequencies to facilitate the fabrication of structures with a large aspect ratio.

To test the PBG structure, millimeter-wave transmitters and receivers were placed next to the entrance and exit of the PBG waveguide. This coupling scheme closely resembled the setup used in the computational simulations. The transmitted signal is shown in Fig. 7. The signal is normalized to the transmitted signal of a straight waveguide. The PBG bend exhibits near-perfect transmission around 87 and 101 GHz. The two arrows indicate the expected positions of the reflection nodes computed from the one-dimensional scattering model. The positions of the nodes confirm a subtle and important point about PBG waveguides: The detection of light at the end of a straight waveguide would not be a sufficient condition, in itself, to confirm PBG guiding. It is the existence of transmission peaks around the sharp bend, along with the specific position of these peaks, that confirms PBG guiding.

Waveguide Intersections

In addition to sharp bends, photonic crystals can be used to fabricate waveguide intersections with low crosstalk. If two waveguides intersect each other on the same plane, light traveling along one waveguide typically leaks into the second waveguide, causing signal loss and crosstalk. The insertion of a microcavity at the center of the intersection of two PBG waveguides can reduce the crosstalk and increase the throughput. If the resonant mode inside the cavity is such that it can couple only to one waveguide, the crosstalk can be essentially eliminated. In this case, the problem reduces to the well-known phenomenon of resonant tunneling through a cavity.

Figure 8a shows two intersecting waveguides in a two-dimensional photonic crystal identical to the one shown in Fig. 7. At the center of the intersection, a microcavity is created by adding rods inside the waveguides and by increasing the radius of one rod by 60 percent. The cavity is outlined by a dashed box. The cavity supports two degenerate modes with opposite symmetry at a frequency lying inside the bandgap. From symmetry, each resonant mode can couple to only one waveguide, as shown schematically in Fig. 8b. Therefore, under the approximation that the waveguides couple

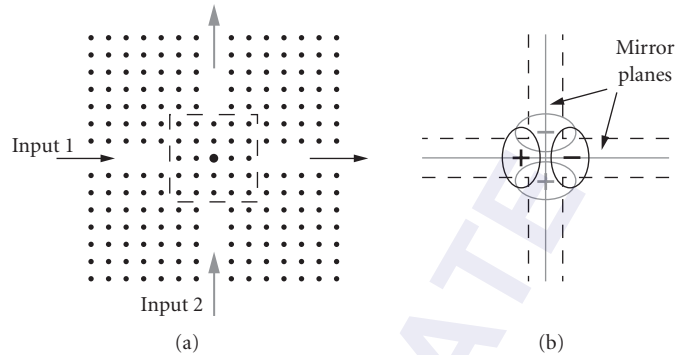


FIGURE 8 (a) Diagram of two intersecting waveguides inside a photonic crystal. The two waveguides are aligned along the (10) and (01) directions. A microcavity—outlined by the dashed line—is created at the center of the intersection by adding columns inside the waveguides and by increasing the size of the dielectric column at the center. The microcavity supports two degenerate modes with opposite symmetry. The mode contours are shown schematically in (b). By symmetry, the modes corresponding to the black contour lines cannot couple to even modes in the waveguide along the (01) direction, and the modes corresponding to the gray contour lines cannot couple to even modes in the waveguide along the (10) direction.

to one another only through the resonant cavity, crosstalk is prohibited. The throughput in each waveguide is described by resonant tunneling; the throughput spectrum is a Lorentzian function with 100 percent transmission at resonance. The width of the resonance is given by the inverse of the quality factor of the microcavity.

In general, large throughput and low crosstalk can be achieved if each waveguide has a single guided mode in the frequency range of interest, and if the microcavity supports two resonant modes, each mode having even symmetry with respect to the mirror plane along one waveguide and odd symmetry with respect to the other mirror plane. The presence of radiation loss would reduce the throughput and increase the crosstalk. A detailed description of PBG waveguide intersections can be found in Ref. 37.

9.8 CONCLUSION

The routing and interconnection of optical signals through narrow channels and around sharp bends are important for large-scale all-optical circuit applications. In addition to sharp bends and low-crosstalk intersections, photonic bandgap materials can also be used for narrowband filters, add/drop filters, light emitters, low-threshold lasers, modulators, attenuators, and dispersion compensators. PBG materials may enable the high-density integration of optical components on a single chip.

While this chapter focuses mostly on applications for high-density optical circuits, many other applications have been proposed for PBG materials. One such application is the PBG fiber.³⁸ While photonic crystals can guide light along a periodic plane (as shown in Sec. 9.7), they can also guide light along the direction perpendicular to the plane of periodicity. A PBG fiber is a two-dimensional periodic structure that essentially extends to infinity along the nonperiodic direction. Light is confined inside the fiber by a defect located at the center. PBG fibers may have interesting features such as single-mode operation over a large bandwidth and preferred dispersion compensation properties. Other applications can be found in Ref. 39.

9.9 REFERENCES

1. Alan Michette, "Zone and Phase Plates, Bragg-Fresnel Optics," in *Handbook of Optics*, vol. III, McGraw-Hill, New York, 2000.
2. J. A. Dobrowolski, "Optical Properties of Films and Coatings," in *Handbook of Optics*, vol. I, McGraw-Hill, New York, 1978.
3. T. L. Koch, F. J. Leonberger, and P. G. Suchoski, "Integrated Optics," in *Handbook of Optics*, vol. II, McGraw-Hill, New York, 1995.
4. J. N. Winn, Y. Fink, S. Fan, and J. D. Joannopoulos, "Omnidirectional Reflection from a One-Dimensional Photonic Crystal," *Opt. Lett.* **23**:1573–1575 (1998).
5. J. D. Joannopoulos, R. D. Meade, and J. N. Winn, *Photonic Crystals*, Princeton Press, Princeton, New Jersey, 1995.
6. E. Yablonovitch, "Photonic Band-Gap Structures," *J. Opt. Soc. Am. B* **10**:283–295 (1993).
7. B. Goss Levi, "Visible Progress Made in Three-Dimensional Photonic 'Crystals,'" *Phys. Today*, January, 17–19 (1999).
8. R. D. Meade, A. M. Rappe, K. D. Brommer, and J. D. Joannopoulos, "Accurate Theoretical Analysis of Photonic Band-Gap Materials," *Phys. Rev. B* **48**:8434–8437 (1993). Erratum: S. G. Johnson, *Phys. Rev. B* **55**:15942 (1997).
9. J. B. Pendry, "Photonic Band Structures," *J. Mod. Optics* **41**:209–229 (1994).
10. K. S. Kunz and R. J. Luebbers, *The Finite-Difference Time-Domain Method for Electronics*, CRC Press, Boca Raton, Florida, 1993.
11. R. D. Meade, K. D. Brommer, A. M. Rappe, and J. D. Joannopoulos, "Nature of the Photonic Band Gap: Some Insights from a Field Analysis," *J. Opt. Soc. Am. B* **10**:328–332 (1993).
12. K. M. Ho, C. T. Chan, and C. M. Soukoulis, "Existence of a Photonic Gap in Periodic Dielectric Structures," *Phys. Rev. Lett.* **65**:3152–3155 (1990).
13. E. Yablonovitch, T. J. Gmitter, and K. M. Leung, "Photonic Band Structure: The Face-Centered-Cubic Case Employing Nonspherical Atoms," *Phys. Rev. Lett.* **67**:2295–2298 (1991).
14. C. C. Cheng and A. Scherer, "Fabrication of Photonic Band-Gap Crystals," *J. Vac. Sci. Technol. B* **13**:2696–2700 (1995).
15. E. Ozbay, A. Abeyta, G. Tuttle, M. Tringides, R. Biswas, C. T. Chan, C. M. Soukoulis, and K. M. Ho, "Measurement of a Three-Dimensional Photonic Band Gap in a Crystal Structure Made of Dielectric Rods," *Phys. Rev. B* **50**:1945–1948 (1994).
16. S. Fan, P. R. Villeneuve, R. D. Meade, and J. D. Joannopoulos, "Design of Three-Dimensional Photonic Crystals at Submicron Lengthscales," *Appl. Phys. Lett.* **65**:1466–1468 (1994).
17. S. Y. Lin, J. G. Fleming, D. L. Hetherington, B. K. Smith, R. Biswas, K. M. Ho, M. M. Sigalas, W. Zubrzycki, S. R. Kurtz, and J. Bur, "A Three-Dimensional Photonic Crystal Operating at Infrared Wavelengths," *Nature* **394**:251–253 (1998).
18. A. Yariv, *Optical Electronics*, Saunders, Philadelphia, Pennsylvania, 1991.
19. P. R. Villeneuve, S. Fan, and J. D. Joannopoulos, "Microcavities in Photonic Crystals: Mode Symmetry, Tunability, and Coupling Efficiency," *Phys. Rev. B* **54**:7837–7842 (1996).
20. H. Yokoyama and S. D. Brorson, "Rate Equation Analysis of Microcavity Lasers," *J. Appl. Phys.* **66**:4801–4805 (1989).
21. T. F. Krauss, R. M. De La Rue, and S. Band, "Two-Dimensional Photonic Bandgap Structures Operating at Near-Infrared Wavelengths," *Nature* **383**:699–702 (1996).
22. H. A. Haus, *Waves and Fields in Optoelectronics*, Prentice Hall, Englewood Cliffs, New Jersey, 1984.
23. R. Coccioli, M. Boroditsky, K. W. Kim, Y. Rahmat-Samii, and E. Yablonovitch, "Smallest Possible Electromagnetic Mode Volume in a Dielectric Cavity," *IEE Proc.-Optoelectron.* **145**:391–397 (1998).
24. P. R. Villeneuve, S. Fan, S. G. Johnson, and J. D. Joannopoulos, "Three-Dimensional Photon Confinement in Photonic Crystals of Low-Dimensional Periodicity," *IEE Proc.-Optoelectron.* **145**:384–390 (1998).
25. J. S. Foresi, P. R. Villeneuve, J. Ferrera, E. R. Thoen, G. Steinmeyer, S. Fan, J. D. Joannopoulos, L. C. Kimerling, Henry I. Smith, and E. P. Ippen, "Photonic-Bandgap Microcavities in Optical Waveguides," *Nature* **390**:143–145 (1997).

26. S. Fan, P. R. Villeneuve, and J. D. Joannopoulos, "Rate-Equation Analysis of Output Efficiency and Modulation Rate of Photonic-Crystal Light Emitting Diodes," *IEEE J. Quantum Electron* **36**: October (2000).
27. R. K. Lee, O. J. Painter, B. D'Urso, A. Scherer, and A. Yariv, "Measurement of Spontaneous Emission from a Two-Dimensional Photonic Band Gap Defined Microcavity at Near-Infrared Wave lengths," *Appl. Phys. Lett.* **74**:1522–1524 (1999).
28. M. Meier, A. Mekis, A. Dodabalapur, A. Timko, R. E. Slusher, J. D. Joannopoulos, and O. Nalamasu, "Laser Action from Two-Dimensional Feedback in Photonic Crystals," *Appl. Phys. Lett.* **74**:7–9 (1999).
29. K. Inoue, M. Sasada, J. Kawamata, K. Sakoda, and J. Haus, "A Two-Dimensional Photonic Crystal Laser," *Jpn. J. Appl. Phys.* **38**:L157–L159 (1999).
30. T. Baba and T. Matsuzaki, "Fabrication and Photoluminescence Studies of GaInAsP/InP 2-Dimensional Photonic Crystals," *Jpn. J. Appl. Phys.* **35**:1348–1352 (1996).
31. P. L. Gourley, J. R. Wendt, G. A. Vawter, T. M. Brennan, and B. E. Hammons, "Optical Properties of Two-Dimensional Photonic Lattices Fabricated as Honeycomb Nanostructures in Compound Semiconductors," *Appl. Phys. Lett.* **64**:687–689 (1994).
32. M. Boroditsky, R. Vrijen, T. F. Krauss, R. Coccioli, R. Bhat, and E. Yablonovitch, "Spontaneous Emission Extraction and Purcell Enhancement from Thin-Film 2-d Photonic Crystals," *J. Lightwave Technol.* **17**:2096–2112 (1999).
33. S. G. Johnson, S. Fan, P. R. Villeneuve, and J. D. Joannopoulos, "Guided Modes in Photonic Crystal Slabs," *Phys. Rev. B* **60**:5751–5758 (1999).
34. S. G. Johnson, P. R. Villeneuve, S. Fan, and J. D. Joannopoulos, "Linear Waveguides in Photonic-Crystal Slabs," *Phys. Rev. B* **62**: September (2000).
35. A. Mekis, J. C. Chen, I. Kurland, S. Fan, P. R. Villeneuve, and J. D. Joannopoulos, "High Transmission Through Sharp Bends in Photonic Crystal Waveguides," *Phys. Rev. Lett.* **77**:3787–3790 (1996).
36. S. Y. Lin, E. Chow, V. Hietala, P. R. Villeneuve, and J. D. Joannopoulos, "Experimental Demonstration of Guiding and Bending of Electromagnetic Waves in a Photonic Crystal," *Science* **282**:274–276 (1998).
37. S. G. Johnson, C. Manolatu, S. Fan, P. R. Villeneuve, J. D. Joannopoulos, and H. A. Haus, "Elimination of Cross Talk in Waveguide Intersections," *Opt. Lett.* **23**:1855–1857 (1998).
38. J. C. Knight, J. Broeng, T. A. Birks, and P. St. J. Russel, "Photonic Band Gap Guidance in Optical Fibers," *Science* **282**:1476–1478 (1998).
39. C. M. Soukoulis, ed., *Photonic Band Gap Materials*, NATO ASI Series E: Applied Sciences, Kluwer Academic, Dordrecht, 1996.

This page intentionally left blank.

DO NOT DUPLICATE

PART

2

NONLINEAR
OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

Chung L. Tang

*School of Electrical and Computer Engineering
Cornell University
Ithaca, New York*

10.1 GLOSSARY

c	velocity of light in free space
\mathbf{D}	displacement vector
d_{mn}	Kleinman's \mathbf{d} -coefficient
\mathbf{E}	electric field in lightwave
$\tilde{\mathbf{E}}$	complex amplitude of electric field
e	electronic charge
f	oscillator strength
\hbar	Planck's constant
I	intensity of lightwave
\mathbf{k}	propagation vector
m	mass of electron
N	number of equivalent harmonic or anharmonic oscillators per volume
$n_{1,2}$	index of refraction at the fundamental and second-harmonic frequencies, respectively
\mathbf{P}	macroscopic polarization
$\mathbf{P}^{(n)}$	n th-order macroscopic polarization
$P_{0,2,+,-}$	power of lightwave at the fundamental, second-harmonic, sum-, and difference-frequencies, respectively
$\tilde{\mathbf{P}}$	complex amplitude of macroscopic polarization
\mathbf{Q}	amplitude of vibrational wave or optic phonons
\mathbf{S}	strain of acoustic wave or acoustic phonons
T_{mn}	relaxation time of the density matrix element ρ_{mn}
Γ_j	damping constant of j th optical transition mode
δ	Miller's coefficient
$\epsilon(\mathbf{E})$	field-dependent optical dielectric tensor
ϵ_n	n th-order optic dielectric tensor
ϵ_0	optical dielectric constant of free space
η	amplitude of plasma wave or plasmons
λ	wavelength
ρ_{mn}	density matrix element

$\chi(\mathbf{E})$	field-dependent optic susceptibility tensor
χ_1 or $\chi^{(1)}$	linear optic susceptibility tensor
χ_n or $\chi^{(n)}$	n th-order optic susceptibility tensor
ω_p	plasma frequency
$\langle a \mathbf{p} b \rangle$	dipole moment between states a and b

10.2 INTRODUCTION

For linear optical materials, the macroscopic polarization induced by light propagating in the medium is proportional to the electric field:

$$\mathbf{P} = \epsilon_0 \chi_1 \cdot \mathbf{E} \quad (1)$$

where the linear optical susceptibility χ_1 and the corresponding linear dielectric constant $\epsilon_1 = \epsilon_0(1 + \chi_1)$ are field-independent constants of the medium.

With the advent of the laser, light intensities orders of magnitude brighter than what could be produced by any conventional sources are now possible. When the corresponding field strength reaches a level on the order of, say, 100 KV/m or more, materials that are normally “linear” at lower light-intensity levels may become “nonlinear” in the sense that the optical “constants” are no longer “constants” independent of the light intensity. As a consequence, when the field is not weak, the optical susceptibility χ and the corresponding dielectric constant ϵ of the medium can become functions of the electric field $\chi(\mathbf{E})$ and $\epsilon(\mathbf{E})$, respectively. Such a field-dependence in the optical parameters of the material can lead to a wide range of nonlinear optical phenomena and can be made use of for a great variety of new applications.

Since the first experimental observation of optical second-harmonic generation by Franken¹ and the formulation of the basic principles of nonlinear optics by Bloembergen and coworkers² shortly afterward, the field of nonlinear optics has blossomed into a wide-ranging and rapidly developing branch of optics. There is now a vast literature on this subject including numerous review articles and books.³⁻⁶ It is not possible to give a full review of such a rich subject in a short introductory chapter in this *Handbook*; only the basic principles underlying the lowest order, the second-order, nonlinear optical processes and some illustrative examples of related applications will be discussed here. The reader is referred to the original literature for a more complete account of the full scope of this field.

If the light intensity is not so weak that the field dependence can be neglected and yet not too strong, the optical susceptibility and the corresponding dielectric constant can be expanded in a Taylor series:

$$\chi(\mathbf{E}) = \chi_1 + \chi_2 \cdot \mathbf{E} + \chi_3 : \mathbf{E}\mathbf{E} + \dots \quad (2)$$

or

$$\epsilon(\mathbf{E}) = \epsilon_1 + \epsilon_2 \cdot \mathbf{E} + \epsilon_3 : \mathbf{E}\mathbf{E} + \dots \quad (3)$$

where

$$\epsilon_1 = \epsilon_0(1 + \chi_1) \quad (4)$$

$$\epsilon_n = \epsilon_0 \chi_n \quad \text{for } n \geq 2 \quad (5)$$

and ϵ_0 is the dielectric constant of free space. When these field-dependent terms in the optical susceptibility are not negligible, the induced macroscopic polarization in the medium contains terms that are proportional nonlinearly to the field:

$$\begin{aligned}\mathbf{P} &= \epsilon_0 \chi_1 \cdot \mathbf{E} + \epsilon_0 \chi_2 : \mathbf{E}\mathbf{E} + \epsilon_0 \chi_3 : \mathbf{E}\mathbf{E}\mathbf{E} + \dots \\ &= \mathbf{P}^{(1)} + \mathbf{P}^{(2)} + \mathbf{P}^{(3)} + \dots\end{aligned}\quad (6)$$

As the field intensity increases, these nonlinear polarization terms $\mathbf{P}^{(n>1)}$ become more and more important, and will lead to a large variety of nonlinear optical effects.

The more widely studied of these nonlinear optical effects are, of course, those associated with the lower-order terms in Eq. (6). The second-order nonlinear effects will be discussed in some detail in this chapter. Many of the higher-order nonlinear terms have been observed and are the bases of a variety of useful nonlinear optical devices. Examples of the third-order effects are third-harmonic generation^{7,8} associated with $|\chi^{(3)}(3\omega = \omega + \omega + \omega)|^2$, two-photon absorption⁹ associated with $\text{Im } \chi^{(3)}(\omega_1 = \omega_1 + \omega_2 - \omega_2)$, self-focusing^{10,11} and light-induced index-of-refraction¹² change associated with $\text{Re } \chi^{(3)}(\omega = \omega + \omega - \omega)$, four-wave mixing¹³ $|\chi^{(3)}(\omega_4 = \omega_1 + \omega_2 - \omega_3)|^2$, degenerate four-wave mixing or phase-conjugation^{14,15} $|\chi^{(3)}(\omega = \omega + \omega - \omega)|^2$, optical Kerr effect¹⁶ $\text{Re } \chi^{(3)}(\omega = 0 + 0 + \omega)$, and many others.

There is also a large variety of dynamic nonlinear optical effects such as photon echo,¹⁷ optical nutation¹⁸ (or optical Rabi effect¹⁹), self-induced transparency,²⁰ picosecond²¹ and femtosecond²² quantum beats, and others.

In addition to the nonlinear optical processes involving only photons that are related to the nonlinear dependence on the E-field as shown in Eq. (6), the medium can become nonlinear indirectly through other types of excitations as well. For example, the optical susceptibility can be a function of the molecular vibrational amplitude Q in the medium, or the stress associated with an acoustic wave S in the medium, or the amplitude η of any space-charge or plasma wave, or even a combination of these excitations as in a polariton, in the medium:

$$\begin{aligned}\mathbf{P} &= \epsilon_0 [\chi_1 + \chi_2 : \mathbf{E} + \chi_3 : \mathbf{E}\mathbf{E} + \dots] \mathbf{E} \\ &+ \epsilon_0 [\chi_q : \mathbf{Q} + \chi_a : \mathbf{S} + \chi_\eta : \boldsymbol{\eta} + \dots] \mathbf{E}\end{aligned}\quad (7)$$

giving rise to the interaction of optical and molecular vibrational waves, or optical and acoustic phonons, etc. Nonlinear optical processes involving interaction of laser light and molecular vibrations in gases or liquids or optical phonon in solids can lead to stimulated Raman^{23–25} processes. Those involving laser light and acoustic waves or acoustic phonons lead to stimulated Brillouin^{26–28} processes. Those involving laser light and mixed excitations of photons and phonons lead to stimulated polariton²⁹ processes. Again, there is a great variety of such general nonlinear optical processes in which excitations other than photons in the medium may play a role. It is not possible to include all such nonlinear optical processes in the discussions here. Extensive reviews of the subject can be found in the literature.^{3–5}

10.3 BASIC CONCEPTS

Microscopic Origin of Optical Nonlinearity

Classical Harmonic Oscillator Model of Linear Optical Media The linear optical properties, including dispersion and single-photon absorption, of optical materials can be understood phenomenologically on the basis of the classical harmonic oscillator model (or Drude model). In this simple model, the optical medium is represented by a collection of independent identical harmonic oscillators

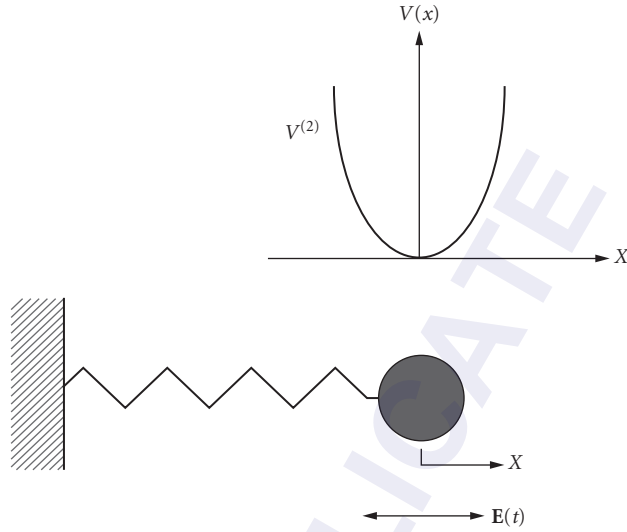


FIGURE 1 Harmonic oscillator model of linear optical media.

embedded in a host medium. The harmonic oscillator is characterized by four parameters: a spring constant k , a damping constant Γ , a mass m , and a charge $-e\sqrt{f}$ as shown schematically in Fig. 1. f is also known as the oscillator-strength and $-e$ is the charge of an electron. The resonance frequency ω_0 of the oscillator is then equal to $[k/m]^{1/2}$.

In the presence of, for example, a monochromatic wave:

$$\mathbf{E} = \frac{1}{2} [\tilde{\mathbf{E}} e^{-i\omega t} + \tilde{\mathbf{E}}^* e^{i\omega t}] \quad (8)$$

the response of the medium is determined by the equation of motion of the oscillator in the presence of the field:

$$\frac{\partial^2 X^{(1)}(t)}{\partial t^2} + \Gamma \frac{\partial X^{(1)}(t)}{\partial t} + \omega_0^2 X^{(1)}(t) = \frac{-e\sqrt{f}}{2m} [\tilde{\mathbf{E}} e^{-i\omega t} + \text{c.c.}] \cdot x \quad (9)$$

where $X^{(1)}(t)$ is the deviation of the harmonic oscillator from its equilibrium position in the absence of the field. The corresponding linear polarization in the steady state and linear complex susceptibility are from Eqs. (8) and (9):

$$\begin{aligned} \mathbf{P}^{(1)} &= -NeX^{(1)}(t)x = \frac{1}{2} [\tilde{\mathbf{P}}^{(1)} e^{-i\omega t} + \tilde{\mathbf{P}}^{(1)*} e^{i\omega t}] \\ &= \frac{Ne^2 f \tilde{\mathbf{E}}}{2mD(\omega)} e^{-i\omega t} + \text{c.c.} \end{aligned} \quad (10)$$

and

$$\epsilon_0 \chi^{(1)} = \frac{|\tilde{\mathbf{P}}|}{|\tilde{\mathbf{E}}|} = \frac{Ne^2 f}{mD(\omega)} \quad (11)$$

where N is the volume density of the oscillators and $D(\omega) = \omega_0^2 - \omega^2 - i\omega\Gamma$. The corresponding real and imaginary parts of the corresponding linear complex dielectric constant of the medium ϵ_1

and $\text{Im } \epsilon_1$, respectively, describe then the dispersion and absorption properties of the linear optical medium. To represent a real medium, the results must be summed over all the effective oscillators (j):

$$\text{Re } \epsilon_1 = \epsilon_0 + \sum_j \frac{\omega_{pj}^2 f_j (\omega_{0j}^2 - \omega^2)}{(\omega_{0j}^2 - \omega^2)^2 + \omega^2 \Gamma_j^2} \quad (12)$$

and

$$\text{Im } \epsilon_1 = \sum_j \frac{\omega_{pj}^2 f_j \omega \Gamma_j}{(\omega_{0j}^2 - \omega^2)^2 + \omega^2 \Gamma_j^2} \rightarrow \frac{\omega_{pj}^2 f_j}{2\omega} \frac{\Gamma_j / 2}{(\omega - \omega_{0j})^2 - (\Gamma_j / 2)^2} \quad \text{for } \omega \approx \omega_{0j} \quad (13)$$

where $\omega_{pj}^2 = 4\pi N_j e^2 / m$ is the plasma frequency for the j th specie of oscillators. Each specie of oscillators is characterized by four parameters: the plasma frequency ω_{pj} , the oscillator-strength f_j , the resonance frequency ω_{0j} , and the damping constant Γ_j . These results show the well-known anomalous dispersion and lorentzian absorption lineshape near the transition or resonance frequencies.

The difference between the results derived using the classic harmonic oscillator or the Drude model and those derived quantum mechanically from first principles is that, in the latter case, the oscillator strengths and the resonance frequencies can be obtained directly from the transition frequencies and induced dipole moments of the transitions between the relevant quantum states in the medium. For an understanding of the macroscopic linear optical properties of the medium, extended versions of Eqs. (12) and (13), including the tensor nature of the complex linear susceptibility, are quite adequate.

Anharmonic Oscillator Model of the Second-Order Nonlinear Optical Susceptibility An extension of the Drude model with the inclusion of suitable anharmonicities in the oscillator serves as a useful starting point in understanding the microscopic origin of the optical nonlinearity classically. Suppose the spring constant of the oscillator representing the optical medium is not quite linear in the sense that the potential energy of the oscillator is not quite a quadratic function of the deviation from the equilibrium position, as shown schematically in Fig. 2. In this case, the response of the oscillator to a harmonic force is asymmetric. The deviation (solid line) from the equilibrium position is larger and smaller on alternate half-cycles than that in the case of the harmonic oscillator. This means that there must be a second-harmonic component (dark shaded curve) in the response of the oscillator as shown schematically in Fig. 3. It is clear, then, that the larger the anharmonicity and the corresponding asymmetry in the oscillator potential, the larger the second-harmonic in the response. Extending this kind of consideration to a three-dimensional model, it implies that to have second-harmonic generation, the material must not have inversion symmetry and, therefore, must be crystalline. It is also clear that for the third and higher odd harmonics, the anharmonicity in the oscillator potential should be symmetric. Even harmonics will always require the absence of inversion symmetry. Beyond that, obviously, the larger the anharmonicities, the larger the nonlinear effects.

Consider first the second-harmonic case. The corresponding anharmonic oscillator equation is:

$$\frac{\partial^2 X(t)}{\partial t^2} + \Gamma \frac{\partial X(t)}{\partial t} + \omega_0^2 X(t) + \nu X(t)^2 = \frac{-e\sqrt{f}}{2m} [\tilde{\mathbf{E}} e^{-i\omega t} + \text{c.c.}] \cdot x \quad (14)$$

Solving this equation by perturbation expansion in powers of the \mathbf{E} -field:

$$X(t) = X^{(1)}(t) + X^{(2)}(t) + X^{(3)}(t) + \dots \quad (15)$$

leads to the second-order nonlinear optical susceptibility

$$\epsilon_0 \chi^{(2)} = \frac{|\tilde{\mathbf{P}}^{(2)}|}{|\tilde{\mathbf{E}}^2|} = \frac{Ne^3 f \nu}{2m^2 D^2(\omega) D(2\omega)} \quad (16)$$

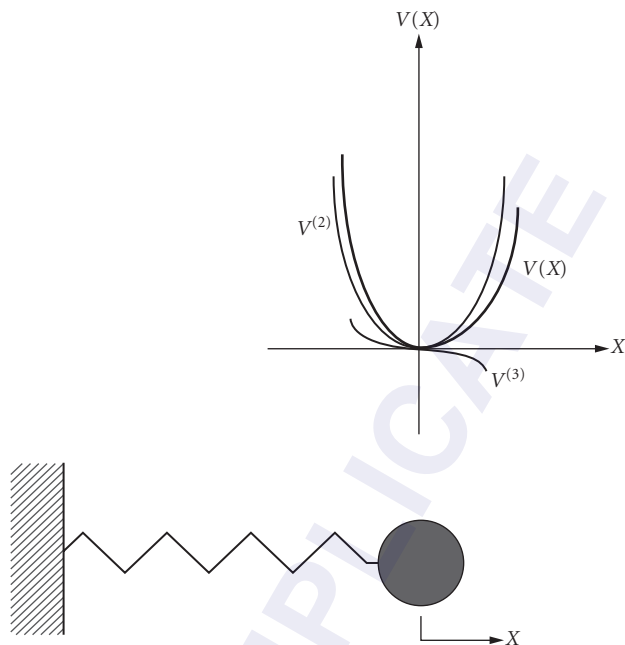


FIGURE 2 Anharmonic oscillator model of nonlinear optical media.

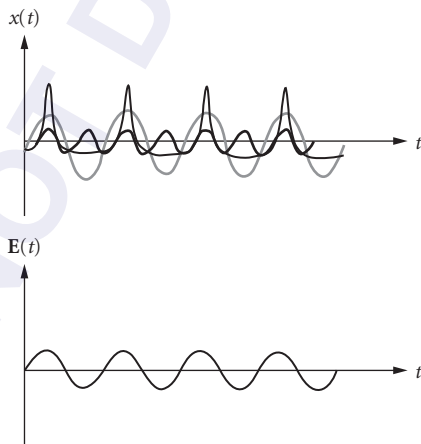


FIGURE 3 Response $[x(t)]$ of anharmonic oscillator to sinusoidal driving field $[E(t)]$.

Unlike in the linear case, a more exact expression of the nonlinear susceptibility-derived quantum mechanically will, in general, have a more complicated form and will involve the excitation energies of, and dipole matrix elements between, all the states. Nevertheless, an expression like Eq. (16) obtained on the basis of the classical anharmonic oscillator model is very useful in discussing qualitatively the second-order nonlinear optical properties of materials. Equation (16) is particularly useful in understanding the dispersion properties of the second nonlinearity.

It is also the basis for understanding the so-called Miller's rule³⁰ which gives a very rough estimate of the order of magnitude of the nonlinear coefficient. We note that the strong frequency dependence in the denominator of the $\chi^{(2)}$ involves factors that are of the same form as those that appeared in $\chi^{(1)}$. Suppose we divide out these factors and define a parameter which is called Miller's coefficient:

$$\delta = \chi^{(2)}(2\omega) / [\chi^{(1)}(\omega)]^2 \chi^{(1)}(2\omega) \epsilon_0^2 = mv / 2e^3 f^{1/2} N^2 \epsilon_0^2 \quad (17)$$

from Eqs. (16) and (11). For many inorganic second-order nonlinear optical crystals, it was first suggested by R. C. Miller that δ was approximately a constant for all materials, and its value was found empirically to be on the order of $2-3 \times 10^{-6}$ esu. If this were true, to find materials with large nonlinear coefficients, one should simply look for materials with large values of $\chi^{(1)}(\omega)$ and $\chi^{(1)}(2\omega)$. This empirical rule was known as Miller's rule. It played an important historical role in the search for new nonlinear optical crystals and in explaining the order of magnitude of nonlinear coefficients for many classes of nonlinear optical materials including such well-known materials as the ADP-isomorphs—for example, KH_2PO_4 (KDP), NH_4PO_4 (ADP), etc.—and the ABO_3 type of ferroelectrics—for example, LiIO_3 , LiNbO_3 , etc.—or III-V and II-VI compound semiconductors in the early days of nonlinear optics.

On a very crude basis, a value of δ can be estimated from Eq. (17) by assuming that the anharmonic potential term in Eq. (14) becomes comparable to the harmonic term when the deviation X is on the order of one lattice spacing in a typical solid, or on the order of an Angstrom. Thus, using standard numbers, Eq. (17) predicts that, in a typical solid, δ is on the order of 4×10^{-6} esu in the visible. It is now known that there are many classes of materials that do not fit this rule at all. For example, there are organic crystals with Miller's coefficients thousands of times larger than this value.

A more rigorous theory for the nonlinear optical susceptibility will clearly have to come from appropriate calculations based upon the principles of quantum mechanics.

Quantum Theory of Nonlinear Optical Susceptibility Quantum mechanically, the nonlinearities in the optical susceptibility originate from the higher-order terms in the perturbation solutions of the appropriate Schrödinger's equation or the density-matrix equation.

According to the density-matrix formalism, the induced macroscopic polarization \mathbf{P} of the medium is specified completely in terms of the density matrix:

$$\mathbf{P} = N \text{Trace}[\mathbf{p}\rho] \quad (18)$$

where \mathbf{p} is the dipole moment operator of the essentially noninteracting individual polarizable units, or "atoms" or molecules or unit cells in a solid, as the case may be, and N is the volume density of such units.

The density-matrix satisfies the quantum mechanical Boltzmann equation or the density-matrix equation:

$$\frac{\partial \rho_{mn}}{\partial t} + i\omega_{mn} \rho_{mn} + \frac{\rho_{mn} - \bar{\rho}_{mn}}{T_{mn}} = \frac{i}{\hbar} \sum_k [\rho_{mk} V_{kn} - V_{mk} \rho_{kn}] \quad (19)$$

where $\bar{\rho}_{mn}$ is the equilibrium density matrix in the absence of the perturbation V and T_{mn} is the relaxation time of the density-matrix element ρ_{mn} . The n th-order perturbation solution of Eq. (19) in the steady state is:

$$\rho_{mn}^{(n)}(t) = \frac{i}{\hbar} \sum_k \int_{-\infty}^t [\rho_{mk}^{(n-1)}(t') V_{kn}(t') - V_{mk}(t') \rho_{kn}^{(n-1)}(t')] \exp\left[\left(i\omega_{mn} + \frac{1}{T_{mn}}\right)(t'-t)\right] dt' \quad (20)$$

The zeroth-order solution is clearly that in the absence of any perturbation or:

$$\rho_{mn}^{(0)} = \bar{\rho}_{mn} \quad (21)$$

In principle, once the zeroth-order solution is known, one can generate the solution to any order corresponding to all the nonlinear optical processes. While such solutions are formally complete and correct, they are generally not very useful, because it is difficult to know all the excitation energies and transition moments of all the states needed to calculate $\chi^{(n\geq 2)}$. For numerical evaluations of $\chi^{(n)}$, various simplifying approximations must be made.

To gain some qualitative insight into the microscopic origin of the nonlinearity, it can be shown on the basis of a simple two-level system that the second-order solution of Eq. (20) leads to the approximate result:

$$\chi^{(2)} \propto [(\omega_{ge} - \omega)(\omega_{ge} - 2\omega)]^{-1} |\langle g|\mathbf{p}|e\rangle|^2 [\langle e|\mathbf{p}|e\rangle - \langle g|\mathbf{p}|g\rangle] \quad (22)$$

It shows that for such a two-level system at least, there are three important factors: the resonance denominator, the transition-moment squared, and the change in the dipole moment of the molecule going from the ground state to the excited state. Thus, to get a large second-order optical nonlinearity, it is preferable to be near a transition with a large oscillator strength and there should be a large change in the dipole moment in going from the ground state to that particular excited state. It is known, for example, that substituted benzenes with a donor and an acceptor group have strong charge-transfer bands where the transfer of charges from the donor to the acceptor leads to a large change in the dipole moment in going from the ground state to the excited state. The transfer of the charges is mediated by the delocalized π electrons along the benzene ring. Thus, there was a great deal of interest in organic crystals of benzene derivatives. This led to the discovery of many organic nonlinear materials. In fact, it was the analogy between the benzene ring structure and the boroxal ring structure that led to the discovery of some of the best known recently discovered inorganic nonlinear crystals such as β -BaB₂O₄ (BBO)³¹ and LiB₃O₅ (LBO).³²

In general, however, there are few rules that can guide the search for new nonlinear optical crystals. It must be emphasized, however, that the usefulness of a material is not determined by its nonlinearity alone. Many other equally important criteria must be satisfied for the nonlinear material to be useful, for example, the transparency, the phase-matching property, the optical damage threshold, the mechanical strength, chemical stability, etc. Most important is that it must be possible to grow single crystals of this material of good optical quality for second-order nonlinear optical applications in bulk crystals. In fact, optical nonlinearity is often the easiest property to come by. It is these other equally important properties that are often harder to predict and control.

Form of the Second-Order Nonlinear Optical Susceptibility Tensor

The simple anharmonic oscillator model shows that to have second-order optical nonlinearity, there must be asymmetry in the crystal potential in some direction. Thus, the crystal must not have inversion symmetry. This is just a special example of how the spatial symmetry of the crystal affects the form of the optical susceptibility. In this case, if the crystal contains inversion symmetry, all the elements of the susceptibility tensor must be zero. In a more general way, the form of the optical susceptibility tensor is dictated by the spatial symmetry of the crystal structure.³³

For second-order nonlinear susceptibilities in the cartesian coordinate system:

$$P_i^{(2)} = \sum_{j,k} \epsilon_0 \chi_{ijk}^{(2)} E_j E_k \quad (23)$$

$\chi_{ijk}^{(2)}$ in general has 27 independent coefficients before any symmetry conditions are taken into account. Taking into account the permutation symmetry condition, namely, the order E_j and E_k appearing in Eq. (23) is not important, or

$$\chi_{ijk}^{(2)} = \chi_{ikj}^{(2)} \quad (24)$$

the number of independent coefficients reduces down to 18. With 18 coefficients, it is sometimes more convenient to define a two-dimensional 3×6 tensor, commonly known as the Kleinman \mathbf{d} -tensor:³⁴

$$\begin{pmatrix} P_x \\ P_y \\ P_z \end{pmatrix} = \epsilon_0 \begin{pmatrix} d_{11} & d_{12} & d_{13} & d_{14} & d_{15} & d_{16} \\ d_{21} & d_{22} & d_{23} & d_{24} & d_{25} & d_{26} \\ d_{31} & d_{32} & d_{33} & d_{34} & d_{35} & d_{36} \end{pmatrix} \begin{pmatrix} E_x^2 \\ E_y^2 \\ E_z^2 \\ 2E_y E_z \\ 2E_x E_z \\ 2E_x E_y \end{pmatrix} \quad (25)$$

rather than the three-dimensional tensor $\chi_{ijk}^{(2)} = \chi_{ikj}^{(2)}$. One obvious advantage of the \mathbf{d}_{im} = tensor form is that the full tensor can be written in the two-dimensional matrix form, whereas it would be difficult to exhibit on paper any three-dimensional matrix.

An additional important point about the \mathbf{d} -tensor is that it is defined in terms of the complex amplitudes of the E-field and the induced polarization with the 1/2 factor explicitly separated out in the front as shown in Eq. (8). In contrast, the definition of $\chi_{ijk}^{(2)}$ may be ambiguous in the literature because not all the authors define the complex amplitude with a 1/2 factor in the front. For linear processes, it makes no difference, because the 1/2 factors in the induced polarization and the E-field cancel out. In nonlinear processes, the 1/2 factors do not cancel and the numerical value of the complex susceptibility will depend on how the complex amplitudes of the E-field and polarization are defined.

For crystalline materials, the remaining 18 coefficients are, in general, not all independent of each other. Spatial symmetry requires, in addition, that they must satisfy the characteristic equation:

$$\chi_{ijk}^{(2)} = \sum_{\alpha\beta\gamma} \chi_{\alpha\beta\gamma}^{(2)} R_{\alpha i} R_{\beta j} R_{\gamma k} \quad (26)$$

where $R_{\alpha i}$, etc., represent the symmetry operations contained in the space group for the particular crystal structure and Eq. (26) must be satisfied for all the Rs in the group. For example, if a crystal has inversion symmetry, or $R_{\alpha i, \beta j, \gamma k} = (-1)\delta_{\alpha i, \beta j, \gamma k}$, Eq. (26) implies that $\chi_{ijk}^{(2)} = (-1)\chi_{ikj}^{(2)} = 0$ as expected. From the known symmetry elements of all 32 crystallographic point groups, the forms of the corresponding second-order nonlinear susceptibility tensors can be worked out and are tabulated. Equation (26) can in fact be generalized³³ to an arbitrarily high order n :

$$\chi_{ijk\dots}^{(n)} = \sum_{\alpha\beta\gamma\dots} \chi_{\alpha\beta\gamma\dots}^{(n)} \dots R_{\alpha i} R_{\beta j} R_{\gamma k} \dots \quad (27)$$

for all the Rs in the group. Thus, the forms of any nonlinear optical susceptibility tensors can in principle be worked out once the symmetry group of the optical medium is known.

The \mathbf{d} -tensors for the second-order nonlinear optical process for all 32-point groups derived from Eq. (26) are shown in, for example, Ref. 34. Similar tensors can in principle be derived from Eq. (27) for the nonlinear optical susceptibilities to any order for any point group.

Phase-Matching Condition (or Conservation of Linear Photon Momentum) in Second-Order Nonlinear Optical Processes

On a microscopic scale, the nonlinear optical effect is usually rather small even at relatively high light-intensity levels. In the case of the second-order effects, the ratio of the second-order term to

the first-order term in Eq. (2), for example, is very roughly the ratio of the applied E-field strength to the “atomic E-field” in the material or:

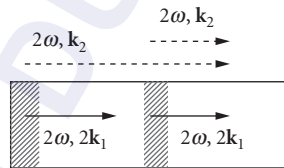
$$\frac{\chi_2 E}{\chi_1} \approx \frac{E}{E_{\text{atomic}}} \tag{28}$$

which is on the order of 10^{-4} even at an intensity level of 1 MW/cm^2 . The same ratio holds very roughly in each successively higher order. To see such a small effect, it is important that the waves generated through the nonlinear optical process add coherently on a macroscopic scale. That is, the new waves generated over different parts of the optical medium add coherently on a macroscopic scale. This requires that the phase velocities of the generated wave and the incident fundamental wave be “matched.”³⁵

Because of the inevitable material dispersion, in general the phases are not matched because the freely propagating second-harmonic wave will propagate at the phase velocity corresponding to the second-harmonic while the source polarization at the second-harmonic will propagate at the phase velocity of the fundamental. Phase matching requires that the propagation constant of the source polarization $2\mathbf{k}_1$ be equal to the propagation constant \mathbf{k}_2 of the second-harmonic or:

$$2\mathbf{k}_1 = \mathbf{k}_2 \tag{29}$$

Multiplying Eq. (29) by \hbar implies that the linear momentum of the photons must be conserved. As shown in the schematic diagram in Fig. 4, in a normally dispersive region of an optical medium, \mathbf{k}_2 is always too long and must be reduced to achieve proper phase matching.



Phase-matching condition: $2k_1 = k_2$

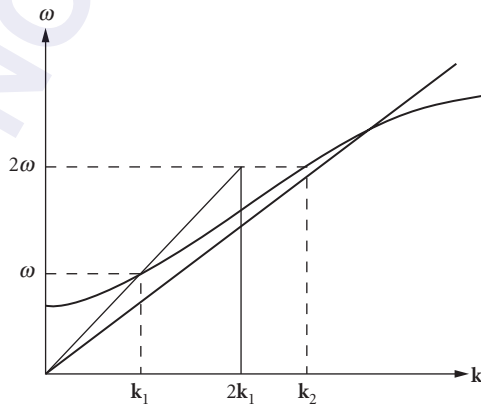


FIGURE 4 Phase-matching requirement and the effect of materials dispersion on momentum mismatch in second-harmonic process.

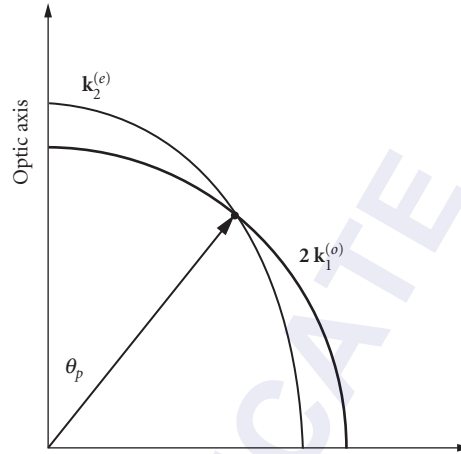


FIGURE 5 Phase matching using birefringence to compensate material dispersion in second-harmonic generation.

In bulk crystals, the most effective and commonly used method is to use birefringence to compensate for material dispersion, as shown schematically in Fig. 5. In this scheme, the \mathbf{k} -vector of the extraordinary wave in the anisotropic crystal is used to shorten \mathbf{k}_2 or lengthen $2\mathbf{k}_1$ as needed. For example, in a negative uniaxial crystal, the fundamental wave is sent into the crystal as an ordinary wave and the second-harmonic wave is generated as an extraordinary wave in a so-called Type I phase-matching condition:

$$2\mathbf{k}_1^{(o)} = \mathbf{k}_2^{(e)} \quad (30)$$

or the fundamental wave is sent in both as an ordinary wave and an extraordinary wave while the second-harmonic is generated as an extraordinary wave in the so-called Type II phase-matching condition:

$$\mathbf{k}_1^{(o)} + \mathbf{k}_1^{(e)} = \mathbf{k}_2^{(e)} \quad (31)$$

In a positive uniaxial crystal, $\mathbf{k}_2^{(e)}$ in Eqs. (30) and (31) should be replaced by $\mathbf{k}_2^{(o)}$ and $\mathbf{k}_1^{(o)}$ in Eq. (30) should be replaced by $\mathbf{k}_1^{(e)}$. Crystals with isotropic linear optical properties clearly lack birefringence and cannot use this scheme for phase matching. Semiconductors of zinc-blende structure, such as the III-V and some of the II-VI compounds, have very large second-order optical nonlinearity but are nevertheless not very useful in the bulk crystal form for second-order nonlinear optical processes because they are cubic and lack birefringence and, hence, difficult to phase match. Phase matching can also be achieved by using waveguide dispersion to compensate for material dispersion. This scheme is often used in the case of III-V and II-VI compounds of zinc-blende structure. Other phase-matching schemes include the use of the dispersion of the spatial harmonics of artificial period structures to compensate for material dispersion.

These phase-matching conditions for the second-harmonic processes can clearly be generalized to other second-order nonlinear optical processes such as the sum- and difference-frequency processes in which two photons of different frequencies and momenta \mathbf{k}_1 and \mathbf{k}_2 either add or subtract to create a third photon of momentum \mathbf{k}_3 . The corresponding phase-matching conditions are:

$$\mathbf{k}_1 \pm \mathbf{k}_2 = \mathbf{k}_3 \quad (32)$$

The practical phase-matching schemes for these processes are completely analogous to those for the second-harmonic process. For example, one can use the birefringence in a bulk optical crystal or the waveguide dispersion to compensate for the material dispersion in a sum- or difference-frequency process.

Conversion Efficiencies for the Second-Harmonic and Sum- and Difference-Frequency Processes

With phase matching, the waves generated through the nonlinear optical process can coherently accumulate spatially. The spatial variation of the complex amplitude of the generated wave follows from the wave equation:

$$\frac{\partial^2}{\partial z^2} E_i(z, t) - \frac{1}{c^2} \frac{\partial^2}{\partial t^2} E_i(z, t) = \frac{1}{c^2 \epsilon_0} \frac{\partial^2}{\partial t^2} P_i(z, t) \quad (33)$$

where

$$E_i(z, t) = \frac{1}{2} [\tilde{E}_{0,i} e^{ik_0 z - i\omega_0 t} + \text{c.c.}] + \frac{1}{2} [\tilde{E}_{2,i}(z) e^{-ik_2 z - i\omega_2 t} + \text{c.c.}] \quad (34)$$

and

$$P_i(z, t) = P_i^{(\omega_0)}(z, t) + [P_{\text{source},i}^{(2\omega_0)}(z, t) + P_i^{(2\omega_0)}(z, t)] \quad (35)$$

$$P_i^{(2\omega_0)}(z, t) = \frac{1}{2} \left[\sum_j \epsilon_0 \chi_{ij}^{(1)}(2\omega_0) \tilde{E}_{2,j}(z) e^{ik_2 z - i\omega_2 t} + \text{c.c.} \right] \quad (36)$$

$$P_{\text{source},i}^{(2\omega_0)}(z, t) = \frac{1}{2} [\tilde{P}_{s,i}^{(2\omega_0)}(z) e^{-i2k_0 z - i\omega_2 t} + \text{c.c.}] \quad (37)$$

$$\tilde{P}_{s,i}^{(2\omega_0)}(z, t) = \frac{1}{2} \sum_{jk} \epsilon_0 \chi_{ijk}^{(2)}(2\omega_0) \tilde{E}_{0,j} \tilde{E}_{0,k} \quad (38)$$

The spatial variation in the complex amplitude of the fundamental wave $\tilde{E}_{0,i}$ in Eq. (34) is assumed negligible and, in fact, we assume it to be that of the incident wave in the absence of any nonlinear conversion in the medium. It is, therefore, implied that the nonlinear conversion efficiency is not so large that the fundamental intensity is appreciably depleted. In other words, the small-signal approximation is implied. Solving Eq. (33) with the boundary conditions that there is no second-harmonic at the input and no reflection at the output end of the crystal, one finds the second-harmonic at the output end of the crystal $z = L$ to be:

$$I^{2\omega_0}(z = L) = \frac{2d^2 I_0^2}{c \epsilon_0 n_2 (n_1 - n_2)^2} \sin^2 \left(\frac{L \omega_0}{c} \right) (n_1 - n_2) \quad (39)$$

where d is the appropriate Kleinman \mathbf{d} -coefficient and the intensities refer to those inside the medium. When the phases of the fundamental and second-harmonic waves are not matched, or $n_1 \neq n_2$, it is clear from Eq. (39) that the second-harmonic intensity is an oscillating function of the crystal length. The maximum intensity is reached at a crystal length of:

$$L_{\text{max}} = \frac{\lambda}{4|n_2 - n_1|} \quad (40)$$

which is also known as the coherence length for the second-harmonic process. The maximum intensity that can be reached is:

$$I_{\max}^{2\omega_0}(z = L_{\max}) = \frac{2d^2}{c\epsilon_0 n_2 (n_1 - n_2)^2} I_0^2 \quad (41)$$

regardless of the crystal length as long as it is greater than the coherence length. The coherence length for many nonlinear optical materials could be on the order of a few microns. Therefore, without phase matching, the second-harmonic intensity in such crystals corresponds to what is generated within a few microns of the output surface of the nonlinear crystal. A much more interesting or important case is clearly when there is phase matching or $n_1 = n_2$.

The second-harmonic intensity under the phase-matched condition is, from Eq. (39):

$$I_2 = \left(\frac{8\pi^2}{c\epsilon_0 n_1^2 n_2} \right) \left(\frac{L}{\lambda_1} \right)^2 d_{\text{eff}}^2 I_0^2 \quad (42)$$

where d_{eff} is the effective \mathbf{d} -coefficient which takes into account the projections of the \mathbf{E} -field and the second-harmonic polarization along the crystallographic axes and the form of the proper \mathbf{d} -tensor for the particular crystal structure. The intensities in this equation refer to the intensities inside the nonlinear medium and the wavelength refers to the free-space wavelength. Equation (42) shows that the second-harmonic intensity under phase-matched conditions is proportional to the square of the length of the crystal measured in the wavelength, as expected for coherent processes. The second-harmonic intensity is also proportional to the effective \mathbf{d} -coefficient squared and the fundamental intensity squared, as expected.

One might be tempted to think that, to increase the second-harmonic power conversion efficiency indefinitely, all one has to do is to focus the beam very tight since the left-hand side is inversely proportional to the beam cross section while the right-hand side is inversely proportional to the cross section squared. Because of diffraction, however, as the fundamental beam is focused tighter and tighter, the effective focal region becomes shorter and shorter. Optimum focusing is achieved when the Rayleigh range of the focal region becomes the limiting interaction length rather than the crystal length. A rough estimate assumes that a beam of square cross section doubles in width (w) due to diffraction in an "optimum focusing length," $L_{\text{opt}} \sim w^2/\lambda_2$, and that this optimum focusing length is equal to the crystal length L . Under such a nominally optimum focusing condition, the maximum second-harmonic power that can be generated in practice is, therefore, approximately:

$$P_2^{(\text{opt})} = \left(\frac{2\pi^2}{c\epsilon_0 n_1^2 n_2} \right) \left(\frac{L}{\lambda_2^3} \right) d_{\text{eff}}^2 P_0^2 \quad (43)$$

Note that this maximum power is linearly proportional to the crystal length. It must be emphasized, however, that this linear dependence is not an indication of incoherent optical process. It is because the beam spot size (area) under the optimum focusing condition is linearly proportional to the crystal length. Numerically, for example, approximately 3 W of second-harmonic power could be generated under optimum focusing in a 1-cm-long LiIO_3 crystal with 30 W of incident fundamental power at 1 μm .

Equation (43) can, in fact, be generalized to other three-photon processes such as the sum-frequency and difference-frequency processes:

$$P_+^{(\text{opt})} = \left(\frac{2\pi^2}{c\epsilon_0 n_1 n_2 n_+} \right) \left(\frac{L}{\lambda_+^3} \right) d_{\text{eff}}^2 (\omega_+ = \omega_1 + \omega_2) P_1 P_2 \quad (44)$$

and

$$P_{-}^{(\text{opt})} = \left(\frac{2\pi^2}{c\epsilon_0 n_1 n_2 n_{-}} \right) \left(\frac{L}{\lambda_{-}^3} \right) d_{\text{eff}}^2 (\omega_{-} = \omega_1 - \omega_2) P_1 P_2 \quad (45)$$

In using Eqs. (44) and (45), one must be especially careful in relating the numerical values of the d_{eff} coefficients for the sum- and difference-frequency processes to that measured in the second-harmonic process because the two low-frequency photons degenerate in frequency in the latter process.

The Optical Parametric Process

A somewhat different, but rather important, second-order nonlinear optical process is the optical parametric process.^{36,37} Optical parametric amplifiers and oscillators powerful solid-state sources of broadly tunable coherent radiation capable of covering the entire spectral range from the near-UV to the mid-IR and can operate down to the femtosecond time domain. The basic principles of optical parametric process were known even before the invention of the laser, dating back to the days of the masers. The practical development of the optical parametric oscillator had been impeded, however, due to the lack of suitable nonlinear optical materials. As a result of recent advances³⁸ in nonlinear optical materials research, these oscillators are now practical devices with broad potential applications in research and industry. The basic physics of the optical parametric process and recent developments in practical optical parametric oscillators are reviewed in this section as an example of wavelength-shifting nonlinear optical devices.

Studies of the optical parameters of materials clearly have always been a powerful tool to gain access to the atomic and molecular structures of optical materials and have played a key role in the formulation of the basic principles of quantum mechanics and, indeed, modern physics. Much of the information obtained through linear optics and linear optical spectroscopy came basically from just the first term in the expansion of the complex susceptibility, Eq. (2). The possibility of studying the higher-order terms in the complex susceptibility through nonlinear optical techniques greatly expands the power of such studies to gain access to the basic building blocks of materials on the atomic or molecular level. Of equal importance, however, are the numerous practical applications of nonlinear optics. Although there are now thousands of known laser transitions in all kinds of laser media, the practically useful ones are still relatively few compared to the needs. Thus, there is always a need to shift the laser wavelengths from where they are available to where they are needed. Nonlinear optical processes are the way to accomplish this. Until recently, the most commonly used wavelength-shifting processes were harmonic generation, sum-, and difference-frequency generation processes. In all these processes, the generated frequencies are always uniquely related to the frequencies of the incident waves. The parametric process is different. In this process, there is the possibility of generating a continuous range of frequencies from a single-frequency input.

For harmonic, sum-, and difference-frequency generation, the basic devices are nothing more than suitably chosen nonlinear optical crystals that are oriented and cut according to the basic principles already discussed in the previous sections and there is a vast literature on all aspects of such devices. The spontaneous optical parametric process can be viewed as the inverse of the sum-frequency process and the stimulated parametric process, or the parametric amplification process, can be viewed as a repeated difference-frequency process.

Spontaneous Parametric Process The spontaneous parametric process, also known as the parametric luminescence or parametric fluorescence process, is described by a simple Feynman diagram as shown in Fig. 6. It describes the process in which an incident photon, called a pump photon, propagating in a nonlinear optical medium breaks down spontaneously into two photons of lower frequencies, called signal and idler photons using a terminology borrowed from earlier microwave

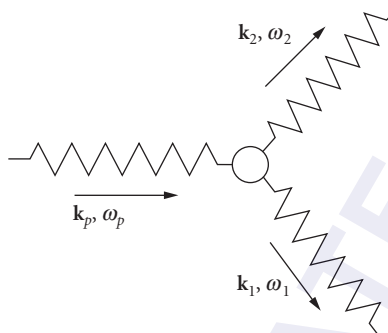


FIGURE 6 Spontaneous breakdown of a pump photon into a signal and an idler photon.

parametric amplifier work, with the energy and momentum conserved:

$$\omega_p = \omega_s + \omega_i \quad (46)$$

$$\mathbf{k}_p = \mathbf{k}_s + \mathbf{k}_i \quad (47)$$

The important point about this second-order nonlinear optical process is that the frequency condition Eq. (46) does not predict a unique pair of signal and idler frequencies for each fixed pump frequency ω_p . Neglecting the dispersion in the optical material, there is a continuous range of frequencies that can satisfy this condition. Taking into account the dispersion in real optical materials, the frequency and momentum matching conditions Eqs. (46) and (47), in general, cannot be satisfied simultaneously. In analogy with the second-harmonic or the sum- or difference-frequency processes, one can use the birefringence in the material to compensate for the material dispersion for a set of photons propagating in the nonlinear crystal. By rotating the crystals, the birefringence in the direction of propagation can be tuned, thereby leading to tuning of the signal and idler frequencies. This tunability gives rise to the possibility of generating photons over a continuous range of frequencies from incident pump photons at one particular frequency, which means the possibility of constructing a continuously tunable amplifier or oscillator by making use of the parametric process.

A complete theory for the spontaneous parametric emission is beyond the scope of this introductory chapter because, as all spontaneous processes, it requires the quantization of the electromagnetic waves. Detailed descriptions of the process can be found in the literature.⁴

Stimulated Parametric Process, or the Parametric Amplification Process With only the pump photons present in the initial state, spontaneous emission occurs at the signal and idler frequencies under phase-matched conditions. With signal and pump photons present in the initial state, stimulated parametric emission occurs in the same way as in a laser medium, except here the pump photons are converted directly into the signal and the corresponding idler photons through the second-order nonlinear optical process and no exchange of energy with the medium is involved. The stimulated parametric process can also be viewed as a repeated difference-frequency process in which the signal and idler photons repeatedly mix with the pump photons in the medium, generating more and more signal and idler photons under the phase-matched condition.

The spatial dependencies of the signal and idler waves can be found from the appropriate coupled-wave equations under the condition when the pump depletion can be neglected. The corresponding complex amplitude of the signal wave at the output $\mathbf{E}_s(L)$ is proportional to that at the input $\mathbf{E}_s(0)$, as in any amplification process:³⁹

$$\mathbf{E}_s(L) = \mathbf{E}_s(0) \cosh gL \quad (48)$$

where

$$g = \frac{d_{\text{eff}} |E_p| \sqrt{k_s k_i}}{2n_s n_i} \quad (49)$$

is the spatial gain coefficient of the parametric amplification process. d_{eff} is the effective Kleinman d -coefficient for the parametric process. k_s and k_i are the phase-matched propagation constants of the signal and idler waves, respectively; n_s and n_i are the corresponding indices of refraction.

Optical Parametric Oscillator Given the parametric amplification process, a parametric oscillator can be constructed by simply adding a pair of Fabry-Perot mirrors, as in a laser, to provide the needed optical feedback of the stimulated emission. The optical parametric oscillator has the unique characteristic of being continuously tunable over a very broad spectral range. This is perhaps one of the most important applications of second-order nonlinear optics.

The basic configuration of an optical parametric oscillator (OPO) is extremely simple. It is shown schematically in Fig. 7. Typically, it consists of a suitable nonlinear optical crystal in a Fabry-Perot cavity with dichroic cavity mirrors which transmit at the pump frequency and reflect at the signal frequency or at the signal and idler frequencies. In the former case, the OPO is a singly resonant OPO (SRO) and, in the latter case, it is a doubly resonant OPO (DRO). The threshold for the SRO is much higher than that for the DRO. The trade-off is that the DRO tends to be highly unstable and, thus, not as useful.

Tuning of the oscillator can be achieved by simply rotating the crystal relative to the direction of propagation of the pump beam or the axis of the Fabry-Perot cavity. As an example of the spectral range that can be covered by the OPO, Fig. 8 shows the tuning curve of a β -barium borate OPO pumped by the third-harmonic output at 355 μm and the fourth-harmonic at 266 μm of a Nd:YAG laser. Also shown are the corresponding spontaneous parametric emissions. The symbols correspond to the experimental data and the solid curves are calculated.³⁸ With a single set of mirrors to resonate the signal wave in the visible, the entire spectral range from about 400 nm to the IR absorption edge of the β -barium borate crystal can be covered. With KTiO_2PO_4 (KTP) or the more recently developed $\text{KTiO}_2\text{AsO}_4$ (KTA) crystals, the tuning range can be extended well into the mid-IR range to the 3- to 5- μm range. With AgGaSe_2 , the potential tuning range could be extended to the 18- to 20- μm range.

The efficiency of the SRO that can be achieved in practice is relatively high, typically over 30 percent on a pulsed basis. Since the OPO is scalable, the output energy is only limited by the pump energy available and can be in the multijoule range.

A serious limitation at the early stage of development is the oscillator linewidth that can be achieved. Without rather complicated and special arrangements, the oscillator linewidth is typically a few Angstroms or more, which is not useful for high-resolution spectroscopic applications. The linewidth problem is, however, not a basic limitation inherent in the parametric process. It is primarily due to the finite pulse length of the pump sources, which limits the cavity length that can be used so that the number of passes by the signal through the nonlinear crystal is not too small. As more suitable pump sources are developed, various line-narrowing schemes⁴⁰ typically used in tunable lasers can be adapted for use in OPOs as well.

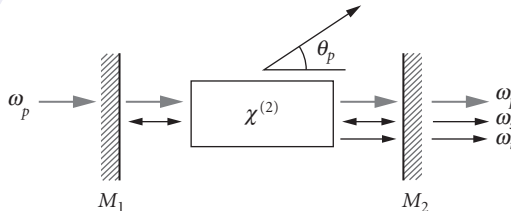


FIGURE 7 Schematic of singly resonant optical parametric oscillator.

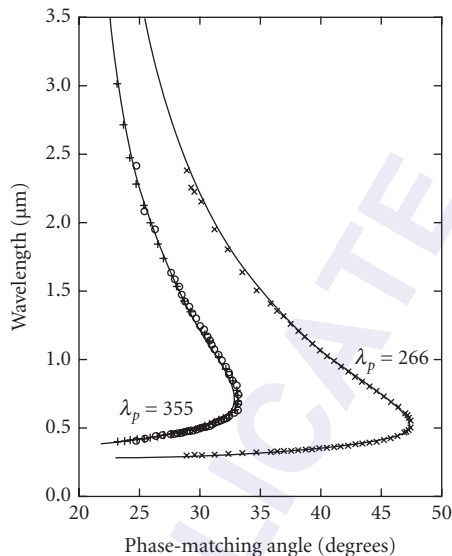


FIGURE 8 Tuning characteristics of BBO spontaneous parametric emission (\times and $+$) and OPO (circles) pumped at the third (355 nm) and fourth (266 nm) harmonics of Nd-YAG laser output. Solid curves are calculated.

The OPO holds promise to become a truly continuously tunable powerful solid-state source of coherent radiation with broad applications as a research tool and in industry.

10.4 MATERIAL CONSIDERATIONS

The second-order nonlinearity is the lowest-order nonlinearity and the first to be observed as the intensity increases. As the discussion following Eq. (26) indicates, only materials without inversion symmetry can have second-order nonlinearity, which means that these must be crystalline materials. The lowest-order nonlinearity in a centrosymmetric system is the third-order nonlinearity.

To observe and to make use of the second-order nonlinear optical effects in a nonlinear crystal, an effective \mathbf{d} -coefficient on the order of 10^{-13} m/V or larger is typically needed. In the case of the third-order nonlinearity, the effect becomes nonnegligible or useful in most applications when it is on the order of 10^{-21} MKS units or more.

Ever since the first observation of the nonlinear optical effect¹ shortly after the advent of the laser, there has been a constant search for new efficient nonlinear materials. To be useful, a large nonlinearity is, however, hardly enough. Minimum requirements in other properties must also be satisfied, such as transparency window, phase-matching condition, optical damage threshold, mechanical hardness, thermal and chemical stability, etc. Above all, it must be possible to grow large single crystals of good optical quality for second-order effects. The perfection of the growth technology for each crystal can, however, be a time-consuming process. All these difficulties tend to conspire to make good nonlinear optical materials difficult to come by.

The most commonly used second-order nonlinear optical crystals in the bulk form tend to be inorganic crystals such as the ADP-isomorphs $\text{NH}_4\text{H}_2\text{PO}_4$ (ADP), KH_2PO_4 (KDP), $\text{NH}_4\text{H}_2\text{AsO}_4$ (ADA), CsH_2AsO_4 (CDA), etc. and the corresponding deuterated version; the ABO_3 type of ferroelectrics such as LiIO_3 , LiNbO_3 , KNbO_3 , etc.; and the borates such as $\beta\text{-BaB}_2\text{O}_4$, LiB_3O_5 , etc.

Although the III-V and II-VI compounds such as GaAs, InSb, GaP, ZnTe, etc. generally have large d -coefficients, because their structures are cubic, there is no birefringence that can be used to compensate for material dispersion. Therefore, they cannot be phase-matched in the bulk and are useful only in waveguide forms. Organic crystals hold promise because of the large variety of such materials and the potential to synthesize molecules according to some design principles. As a result, there have been extensive efforts at developing such materials for applications in nonlinear optics, but very few useful second-order organic crystals have been identified so far. Nevertheless, organic materials, especially for third-order processes, continue to attract a great deal of interest and remain a promising class of nonlinear materials.

To illustrate the important points in considering materials for nonlinear optical applications, a few examples of second-order nonlinear crystals with their key properties are tabulated in Tables 1 through 3. It must be emphasized, however, that because some of the materials are relatively new, some of the numbers listed are subject to confirmation and possibly revision. Discussions of other inorganic and organic nonlinear optical crystals can be found in the literature.⁴¹

As nonlinear crystals and devices become more commercialized, the issues of standardization of nomenclature and conventions and quantitative accuracy are becoming increasingly important. Some of these issues are being addressed⁴² but much work remains to be done.

TABLE 1 Properties of Some Nonlinear Optical Crystals*

Crystal	LiB ₃ O ₅	β -BaB ₂ O ₄ ^f
Point group	mm ^{2a}	3 m
Birefringence	$n_{x=a} = 1.5656^b$ $n_{y=c} = 1.5905$ $n_{z=b} = 1.6055$	$n_e = 1.54254$ $n_o = 1.65510$
Nonlinearity [pm/V]	$d_{32} = 1.16^b$	$d_{22} = 16$ $d_{31} = 0.08$
Transparency [μ m]	0.16–2.6 ^c	0.19–2.5
Γ_{\max} [GW/cm ²]	$\sim 25^b$	$\sim 5^e$
SHG cutoff [nm]	555 ^d	411
$\ell\Delta T$ [°C · cm]	3.9 ^e	55
$\ell\Delta\Theta$ [mrad · cm], CPM	31.3 ^e	0.52
$\ell^{1/2}\Delta\Theta$ [mrad (cm) ^{1/2}]	71.9 ^e NCPM @ 148.0°C	Not available
$\ell\Delta\lambda$ [Å · cm]	Not available	21.1
Δv_g^{-1} @ 630 nm [fs/mm]	240 ^d	360
OPO tuning range [nm]	~ 415 – 2500^d ($\lambda_p = 355$)	~ 410 – 2500 ($\lambda_p = 355$)
Boule size	20 × 20 × 15 mm ^{3e}	$\varnothing 84$ mm × 18 mm
Growth	TSSG ^e @ $\sim 810^\circ\text{C}$	TSSG from Na ₂ O @ $\sim 900^\circ\text{C}$
Predominant growth defects	Flux ^e inclusions	Flux and bubble inclusions
Chemical properties	Nonhygroscopic ^e (m.p. $\sim 834^\circ\text{C}$)	Slightly hygroscopic ($\beta \rightarrow \alpha \sim 925^\circ\text{C}$)

*Data shown is at 1.064 μ m unless otherwise indicated. Γ_{\max} —surface damage threshold; $\ell\Delta T$ —temperature-tuning bandwidth; $\ell\Delta\Theta$, CPM—critical phase-matching acceptance angle; $\ell^{1/2}\Delta\Theta$ —noncritical phase-matching acceptance angle; $\ell\Delta\lambda$ —SHG bandwidth; Δv_g^{-1} —group-velocity dispersion for SHG at 630 nm.

^aVon H. Konig and A. Hoppe, *Z. Anorg. Allg. Chem.* **439**:71 (1978); M. Ihara, M. Yuge, and J. Krogh-Moe, *Yogyo-Kyokai-Shi* **88**:179 (1980); Z. Shuqing, H. Chaoen, and Z. Hongwu, *J. Cryst. Growth* **99**:805 (1990).

^bC. Chen, Y. Wu, A. Jiang, B. Wu, G. You, R. Li, and S. Lin, *J. Opt. Soc. Am.* **B6**:616 (1989); S. Liu, Z. Sun, B. Wu, and C. Chen, *J. App. Phys.* **67**:634 (1989). On the basis of $d_{32} = 2.69 \times d_{36}$ (KDP) and using the value d_{36} (KDP) = 0.39 pm/V according to R. C. Ekart et al., *J. Quan. Elec.* **26**:922 (May 1990).

^c0.16–2.6 μ m: C. Chen, Y. Wu, A. Jiang, B. Wu, G. You, R. Li, and S. Lin, *J. Opt. Soc. Am.* **B6**:616 (1989). 0.165–3.2 μ m; S. Zhao, C. Huang, and H. Zhang, *J. Cryst. Growth* **99**:805 (1990).

^dCalculated by using Sellmeier equations reported in reference; B. Wu, N. Chen, C. Chen, D. Deng, and Z. Xu, *Opt. Lett.* **14**:1080 (1989).

^eT. Ukachi and R. J. Lane, measurements carried out on Cornell LBO crystals grown by self-flux method.

^fReference sources given in: "Growth and Characterization of Nonlinear Optical Crystals Suitable for Frequency Conversion," by L. K. Cheng, W. R. Bosenberg, and C. L. Tang, review article in *Progress in Crystal Growth and Characterization* **20**:9–57 (Pergamon Press, 1990), unless indicated otherwise.

^gEstimated surface damage threshold scaled from detailed bulk damage results reported by H. Nakatani et al., *Appl. Phys. Lett.* **53**:2587 (26 December, 1988).

TABLE 2 Properties of Several Visible Near-IR Nonlinear Optical Crystals*

Characteristics	KNbO ₃ [†]	LiNbO ₃ [‡]	Ba ₂ NaNb ₅ O ₁₅
Point group	mm ²	3 m	mm ²
Transparency [μm]	0.4–5.5	0.4–5.0	0.37–5.0
Birefringence	negative biaxial $n_{x=c} = 2.2574$ $n_{y=a} = 2.2200$ $n_{z=b} = 2.1196$	negative uniaxial $n^o = 2.2325$ $n^e = 2.1560$	negative biaxial $n_{x=b} = 2.2580$ $n_{y=a} = 2.2567$ $n_{z=c} = 2.1700$
Second-order nonlinearity [pm/V]	$d_{32} = 12.9, d_{31} = -11.3$ $d_{24} = 11.9, d_{15} = -12.4$ $d_{33} = -19.6$	$d_{33} = -29.7$ $d_{31} = -4.8$ $d_{32} = 2.3$	$d_{32} = -12.8, d_{31} = -12.8$ $d_{24} = 12.8, d_{15} = -12.8$ $d_{33} = -17.6$
$\partial(n^o - n^{2o})/\partial T$ [°C ⁻¹]	1.6×10^{-4}	-5.9×10^{-5}	1.05×10^{-4}
T _{pm} [°C]	181, d_{32}	-8, d_{31}	89, d_{32} 101, d_{31}
$\ell\Delta T$ [°C-cm]	0.3	0.8	0.5
λ_{SHG} (cutoff) [μm] @ 25°C	0.860	~1.08	1.01
Γ_{max} [MW/cm ²]	Not available	~120	40
Phase transition temperature (°C)	225 and 435	~1000	300
Growth technique	TSSG from K ₂ O @ ~ 1050°C	Czochralski @ ~ 1200°C	Czochralski @ ~ 1440°C
Predominant growth problems	Cracks, blue coloration, multidomains	Temp, induced compositional striations	Striations, microtwinning, multidomains
Postgrowth processing	Poling	Poling	Poling and detwinning
Crystal size	20 × 20 × 20 mm ³ (single domain)	Ø 100 mm × 200 mm (as grown boule)	Ø 20 mm × 50 mm (with striations)

*Unless otherwise specified, data are for $\lambda = 1.064 \mu\text{m}$. (Data taken from: $a, e-i$; $a, b-c$; and a, d , respectively.)

[†]There is a disagreement on the sign of the nonlinear coefficients of KNbO₃ in the literature. Data used here are taken from Ref. *e* with the appropriate correction for the IRE convention.

[‡]Data are for congruent melting LiNbO₃. Five-percent MgO doped crystals gives photorefractive damage threshold about 10–100 times higher.^{*k,l*} The phase-matching properties for these crystals may differ due to the resulting changes in the lattice constants.^{*j*}

^aS. Singh in *CRC Handbook of Laser Science and Technology*, vol. 4, *Optical Materials*, part I, M. J. Weber (ed.), CRC Press, 1986, pp. 3–228.

^bR. L. Byer, J. F. Young, and R. S. Feigelson, *J. Appl. Phys.* **41**:2320 (1970).

^cR. L. Byer in *Quantum Electronics: A Treatise*, H. Rabin and C. L. Tang (eds), vol. 1, part A, Academic Press, 1975.

^dS. Singh, D. A. Draeger, and J. E. Geusic, *Phys. Rev. B* **2**:2709 (1970).

^eY. Uematsu, *Jap. J. Appl. Phys.* **13**: 1362 (1974).

^fP. Gunter, *Appl. Phys. Lett.* **34**:650 (1979).

^gW. Xing, H. Looser, H. Wuest, and H. Arend, *J. Crystal Growth* **78**:431 (1986).

^hD. Shen, *Mat. Res. Bull.* **21**: 1375 (1986).

ⁱT. Fukuda and Y. Uematsu, *Jap. J. Appl. Phys.* **11**:163 (1972).

^jB. C. Grabmaier and F. Otto, *J. Crystal Growth* **79**: 682 (1986).

^kD. A. Bryan, R. Gerson, and H. E. Tomaschke, *Appl. Phys. Lett.* **44**:847 (1984).

^lG. Zhong, J. Jian, and Z. Wu, *11th International Quantum Electronics Conference*, IEEE Cat. No. 80 CH 1561-0, June 1980, p. 631.

10.5 APPENDIX

The results in this article are given in the rationalized MKS systems. Unfortunately, many of the pioneering papers on nonlinear optics were written in the cgs gaussian system. In addition, different conventions and definitions of the nonlinear optical coefficients are used in the literature by different authors. These choices have led to a great deal of confusion. In this Appendix, we give a few key results to facilitate comparison of the results using different definitions and units.

First, in the MKS system, the displacement vector **D** is related to the E-field and the induced polarization **P** in the medium as follows:

$$\begin{aligned} \mathbf{D} &= \epsilon_0 \mathbf{E} + \mathbf{P} \\ &= \epsilon_0 \mathbf{E} + \mathbf{P}^{(1)} + \mathbf{P}^{(2)} + \dots \end{aligned} \quad (\text{A-1})$$

TABLE 3 Properties of Several UV, Visible, and Near-IR Crystals*

Crystal	KDP	KTP (II) [†]
Point group	42 m	mm ²
Birefringence	$n_e = 1.4599$ $n_o = 1.4938$	$n_{x=a} = 1.7367$ $n_{y=b} = 1.7395$ $n_{z=c} = 1.8305$
Nonlinearity [pm/V]	$d_{36} = 0.39$	$d_{32} = 5.0, d_{31} = 6.5$ $d_{24} = 7.6, d_{15} = 6.1$ $d_{33} = 13.7$
Transparency [μm]	0.2–1.4	0.35–4.4
Γ_{max} [GW/cm ²]	~3.5	~15.0
SHG cutoff [nm]	487	~990
$\ell\Delta T$ [$^{\circ}\text{C}\cdot\text{cm}$]	7	22
$\ell\Delta\theta$ [mrad-cm]	1.2	15.7
$\ell\Delta\lambda$ [$\text{\AA}\cdot\text{cm}$]	208 [‡]	4.5
Δv^{-1} @ 630 nm [fs/mm]	185	Not applicable
OPPO tuning range [nm]	~430–700	~610–4200
[nm]	($\lambda_p = 266$)	($\lambda_p = 532$)
ΔT_F [$^{\circ}\text{C}$]	12	Not available
Boule size	40 × 40 × 100 cm ³	~20 × 20 × 20 mm ³
Growth technique	Solution growth from H ₂ O	TSSG from 2KPO ₃ -K ₄ P ₂ O ₇ @ ~1000 $^{\circ}\text{C}$
Predominant growth defects	Organic impurities	Flux inclusions
Chemical properties	Hygroscopic (m.p. ~253 $^{\circ}\text{C}$)	Nonhygroscopic (m.p. ~1172 $^{\circ}\text{C}$)

*Unless otherwise stated, all data for 1064 nm. (Data taken from c, e, a, b, f, m ; and $d, g-i$, respectively.)

[†]KTP Type I interaction gives $d_{\text{eff}} \sim d_{36}$ (KDP) or less for most processes.^m The d_{ij} values^d are for crystals grown by the hydrothermal technique.^{f-l} Significantly lower damage thresholds were reported for hydrothermally grown crystals.

[‡]The anomalously large spectral bandwidth is a manifestation of the λ -noncritical phase matching.ⁿ This is equivalent to a very good group-velocity matching ($\Delta v_g^{-1} \sim 8$ fs/mm) for this interaction in KDP.

^dD. Eimerl, *J. Quant. Elect.* **QE-23**:575 (1987).

^eD. Eimerl, L. Davis, S. Velsko, E. K. Graham, and A. Zalkin, *J. Appl. Phys.* **62**:1968 (1987).

^fD. Eimerl, *Ferroelectrics* **72**:95 (1987).

^gY. S. Liu, L. Drafall, D. Dentz, and R. Belt, *G. E. Technical Information Series Report*, 82CRD016, Feb. 1982.

^hY. Nishida, A. Yokotani, T. Sasaki, K. Yoshida, T. Yamanaka, and C. Yamanaka, *Appl. Phys. Lett.* **52**:420 (1988).

ⁱA. Jiang, F. Cheng, Q. Lin, Z. Cheng, and Y. Zheng, *J. Crystal Growth* **79**:963 (1986).

^jP. Bordui, in *Crystal Growth of KTiOPO₄ from High Temperature Solution*, Ph.D. thesis, Massachusetts Institute of Technology, 1987.

^k*Information Sheet on KTiOPO₄*, Ferroxcube, Division of Ampere Electronic Corp., Saugerties, New York, 1987.

^lP. Bordui, J. C. Jacco, G. M. Loiacono, R. A. Stolzenberger, and J. J. Zola, *J. Crystal Growth* **84**:403 (1987).

^mF. C. Zumsteg, J. D. Bierlein, and T. E. Gier, *J. Appl. Phys.* **47**:4980 (1976).

ⁿR. A. Laudis, R. J. Cava, and A. J. Caporaso, *J. Crystal Growth* **74**:275 (1986).

^oS. Jia, P. Jiang, H. Niu, D. Li, and X. Fan, *J. Crystal Growth* **79**:970 (1986).

^pL. K. Cheng, unpublished.

^qJ. Zyss and D. S. Chemla, in *Nonlinear Optical Properties of Organic Molecules and Crystals*, vol. 1, D. S. Chemla and J. Zyss (eds), Academic Press, 1987, pp. 146–159.

The corresponding wave equation is given in Eq. (33). For the second-order polarization and the corresponding Kleinman \mathbf{d} -coefficients, two definitions are in use. A more popular definition in the current literature is as follows:

$$\mathbf{P}^{(2)} = \epsilon_0 \mathbf{d}_2 : \mathbf{E}\mathbf{E} \quad (\text{A-2})$$

In an earlier widely used reference,³⁴ Yariv defined his \mathbf{d} -coefficient as follows:

$$\mathbf{P}^{(2)} = \mathbf{d}_2^{(\text{Yariv})} : \mathbf{E}\mathbf{E} \quad (\text{A-3})$$

The numerical values of $\mathbf{d}_2^{(\text{Yariv})}$ in this reference (e.g., Table 16.2)³⁴ are given in $(1/9) \times 10^{-22}$ MKS units. The numerical value of ϵ_0 in the MKS system is $10^7 \times (1/4\pi c^2)$ in MKS units. Thus, for example, a tabulated value of $\mathbf{d}_2^{(\text{Yariv})} = 0.5 \times (1/9) \times 10^{-22}$ MKS units in Ref. 34 converts to a numerical value of $d_2 = 0.628$ pm/V in MKS units.

In the cgs gaussian system, the displacement vector \mathbf{D} is related to the \mathbf{E} -field and the induced polarization \mathbf{P} in the medium as follows:

$$\begin{aligned} \mathbf{D} &= \epsilon_0 \mathbf{E} + 4\pi \mathbf{P} \\ &= \epsilon_0 \mathbf{E} + 4\pi \mathbf{P}^{(1)} + 4\pi \mathbf{P}^{(2)} + \dots \end{aligned} \quad (\text{A-4})$$

The corresponding wave equation is:

$$\frac{\partial^2}{\partial z^2} E_i(z, t) - \frac{1}{c^2} \frac{\partial^2}{\partial t^2} E_i(z, t) = \frac{4\pi \partial^2}{c^2 \partial t^2} P_i(z, t) \quad (\text{A-5})$$

The conventional definition of \mathbf{d}_2 is as follows:

$$\mathbf{P}^{(2)} = \mathbf{d}_2 : \mathbf{E}\mathbf{E} \quad (\text{A-6})$$

The numerical value of \mathbf{d}_2 in cgs gaussian units is, therefore, equal to $(3 \times 10^4/4\pi)$ times the numerical value of \mathbf{d}_2 in rationalized MKS units. Thus, continuing with the numerical example given in the preceding paragraph, $\mathbf{d}_2 = 0.628$ pm/V is equal to 1.5×10^{-9} cm/Stat-Volt or 1.5×10^{-9} esu.

As a final check, the expression Eq. (42) for the second-harmonic intensity in the MKS system becomes, in the cgs gaussian system:

$$I_2 = \left(\frac{512\pi^5}{cn_1^2 n_2} \right) \left(\frac{L}{\lambda_1} \right)^2 d_{\text{eff}}^2 I_0^2 \quad (\text{A-7})$$

All the intensities refer to those inside the medium, and the wavelength is the free-space wavelength.

10.6 REFERENCES

1. P. A. Franken, A. E. Hill, C. W. Peters, and G. Weinreich, *Phys. Rev. Lett.* **7**:118 (1961).
2. J. A. Armstrong, N. Bloembergen, J. Ducuing, and P. S. Pershan, *Phys. Rev.* **127**:1918 (1962); N. Bloembergen and Y. R. Shen, *Phys. Rev.* **133**:A37 (1964).
3. N. Bloembergen, *Nonlinear Optics*, Benjamin, New York, 1965.
4. See, for example, H. Rabin and C. L. Tang (eds.), *Quantum Electronics: A Treatise*, vol. 1A and B *Nonlinear Optics*, Academic Press, New York, 1975, and the references therein.
5. See, for example, Y. R. Shen, *The Principles of Nonlinear Optics*, J. W. Wiley Interscience, New York, 1984.
6. M. D. Levenson and S. S. Kano, *Introduction to Nonlinear Laser Spectroscopy*, Academic Press, New York, 1988, and the references therein.
7. P. D. Maker and R. W. Terhune, *Phys. Rev. A* **137**:801 (1965).
8. See, for example, secs. 7.3 and 7.4 of Ref. 5.
9. H. Mahr, "Two-Photon Absorption Spectroscopy," in Ref. 4.
10. G. A. Askar'yan, *Sov. Phys. JETP* **15**:1088, 1161 (1962); M. Hercher, *J. Opt. Soc. Am.* **54**:563 (1964); R. Y. Chiao, E. Garmire, and C. H. Townes, *Phys. Rev. Lett.* **13**:479 (1964) [Erratum, **14**:1056 (1965)].
11. See, for example, Y. R. Shen, "Self-Focusing," chap. 17 in Ref. 5.

12. See, for example, R. W. Boyd, *Nonlinear Optics*, chap. 4, Academic Press, 1992.
13. See, for example, chap. 15 in Ref. 5.
14. Y. B. Zeldovich, V. I. Popoviecher, V. V. Ragul'skii, and F. S. Faizullov, *JETP Letters* **15**:109 (1972).
15. R. W. Hellwarth, *J. Opt. Soc. Am.* **68**:1050 (1978); A. Yariv, *IEEE J. Quant. Elect.* **QE-14**:650 (1978).
16. See, for example, A. Yariv and P. Yeh, *Optical Waves in Crystals*, Wiley, New York, 1984, p. 221.
17. N. A. Kurnit, I. D. Abella, and S. R. Hartmann, *Phys. Rev. Lett.* **13**:567 (1964); S. Hartmann, in R. Glauber (ed.), *Proc. of the Int. School of Phys. Enrico Fermi Course XLII*, Academic Press, New York, 1969, p. 532.
18. C. L. Tang and B. D. Silverman, "Physics of Quantum Electronics," P. Kelley, B. Lax, and P. E. Tannenwald (eds.), McGraw-Hill, 1966, p. 280. G. B. Hocker and C. L. Tang, *Phys. Rev. Lett.* **21**:591 (1969); *Phys. Rev.* **184**:356 (1969).
19. R. G. Brewer, *Phys. Today*, May 1977.
20. S. L. McCall and E. L. Hahn, *Phys. Rev. Lett.* **18**:908 (1967); *Phys. Rev.* **183**:457 (1969).
21. N. Bloembergen and A. H. Zewail, *J. Phys. Chem.* **88**:5459 (1984).
22. M. J. Rosker, F. W. Wise, and C. L. Tang, *Phys. Rev. Lett.* **57**:321 (1986); *J. Chem. Phys.* **86**:2827 (1987).
23. E. J. Woodbury and W. K. Ng, *Proc. IRE* **50**:2347 (1962); R. W. Hellwarth, *Phys. Rev.* **130**:1850 (1963).
24. E. Garmire, E. Pandarese, and C. H. Townes, *Phys. Rev. Lett.* **11**:160 (1963).
25. C. S. Wang, "The Stimulated Raman Process," chap. 7 in Ref. 4.
26. R. Y. Chiao, C. H. Townes, and B. P. Stoicheff, *Phys. Rev. Lett.* **12**:592 (1964); E. Garmire and C. H. Townes, *App. Phys. Lett.* **5**:84 (1964).
27. C. L. Tang, *J. App. Phys.* **37**:2945 (1966).
28. I. L. Fabelinskii, "Stimulated Mandelstam-Brillouin Process," chap. 5 in Ref. 4.
29. See, for example, sec. 10.7 in Ref. 5.
30. R. C. Miller, *App. Phys. Lett.* **5**:17 (1964).
31. C. Chen, B. Wu, A. Jiang, and G. You, *Sci. Sin. Ser. B* **28**:235 (1985).
32. C. Chen, Y. Wu, A. Jiang, B. Wu, G. You, R. Li, and S. Lin, *J. Opt. Soc. Am.* **B6**:616 (1989).
33. P. A. Franken and J. F. Ward, *Rev. of Mod. Phys.* **35**:23 (1963).
34. A. Yariv, *Quantum Electronics*, John Wiley, New York, 1975, pp. 410–411.
35. J. A. Giordmaine, *Phys. Rev. Lett.* **8**:19 (1962); P. D. Maker, R. W. Terhune, M. Nisenhoff, and C. M. Savage, *Phys. Rev. Lett.* **8**:21 (1962).
36. W. H. Louisell, *Coupled Mode and Parametric Electronics*, John Wiley, New York, 1960.
37. N. Kroll, *Phys. Rev.* **127**:1207 (1962).
38. See, for example, C. L. Tang, *Proc. IEEE* **80**:365 (March 1992).
39. Ref. 4, p. 428.
40. See, for example, L. F. Mollenauer and J. C. White (eds.), *Tunable Lasers*, Springer-Verlag, Berlin, 1987.
41. See, for example, S. K. Kurtz, J. Jerphagnon, and M. M. Choy, in *Landolt-Boerstein Numerical Data and Functional Relationships in Science and Technology*, New Series, K. H. Wellwege (ed.), Group III, vol. 11, Springer-Verlag, Berlin, 1979; *Nonlinear Optical Properties of Organic and Polymeric Materials*, D. Williams (ed.), Am. Chem. Soc., Wash., D.C., 1983; *Nonlinear Optical Properties of Organic Molecules*, D. Chemla and J. Zyss (eds.), Academic Press, New York, 1987.
42. D. A. Roberts, *IEEE J. Quant. Elect.* **28**:2057 (1992).

COHERENT OPTICAL TRANSIENTS

Paul R. Berman and Duncan G. Steel

*Physics Department
University of Michigan
Ann Arbor, Michigan*

11.1 GLOSSARY

$\mathbf{E}(\mathbf{R}, t)$	electric field vector
$E(t)$	electric field amplitude
ω	field frequency
$\phi(t)$	field phase
ω_0	atomic transition frequency
$\bar{\omega}_0$	average atomic transition frequency
δ	atom-field detuning
δ_0	average atom-field detuning
$\Omega_0(t)$	Rabi frequency
$\mathbf{\Omega}(t)$	pseudofield vector
$\Omega(t)$	generalized Rabi frequency
$\mathbf{U}(t)$	Bloch vector
$\rho_{ij}(Z, t)$	density matrix element in a field interaction representation
$\rho_{ij}^T(\mathbf{R}, t)$	density matrix element in Schrödinger representation
(u, v, w)	elements of Bloch vector
γ	transverse relaxation rate
γ_2	excited-state decay rate or longitudinal relaxation rate
$\mathbf{P}(\mathbf{R}, t)$	polarization vector
\mathbf{k}	field propagation vector
$E_S(Z, t)$	complex signal electric field amplitude
$P(Z, t)$	complex polarization field amplitude
\mathcal{N}	atomic density

This chapter is dedicated to Richard G. Brewer, a pioneer in coherent optical transients, a mentor and a friend.

L	sample length
μ	dipole moment matrix element
τ, τ_{ref}	pulse durations
θ	pulse area
Δ	difference between local and average transition frequency in a solid
$W_f(\Delta)$	distribution of frequencies in a solid
σ_w	width of $W_f(\Delta)$
\mathbf{v}	atomic velocity
u	most probable atomic speed
$W_0(\mathbf{v})$	atomic velocity distribution
$I(L, t)$	signal intensity exiting the sample
T_{21}, T	time interval between pulses
Γ_t	transit time decay rate
$\Gamma_{2,0}, \Gamma_{2,1}$	branching decay rates of the excited state
$\Omega_i^{(s)} t$	two-photon Rabi frequency
ω_k	recoil frequency
\mathbf{P}	center-of-mass momentum
$E_b(t_1, t_2)$	backscattered electric field amplitude
$J_N(x)$	Bessel function
ξ	one-half of the frequency chirp rate
ω_{RD}	frequency offset between reference and data pulses

11.2 INTRODUCTION

Optical spectroscopy is a traditional method for determining transition frequencies in atoms and molecules. One can classify optical spectroscopy into two broad categories: *continuous-wave* (CW) or stationary spectroscopy and *time-dependent* or transient spectroscopy. In CW spectroscopy, one measures absorption or emission line shapes as a function of the incident frequency of a probe field. The absorption or emission maximum determines the transition frequency, while the width of the line is a measure of relaxation processes affecting the atoms or molecules. It is necessary to model the atom-field interaction to obtain predictions for the line shapes, but, once this is done, it is possible to extract the relevant transition frequencies and relaxation rates from the line shapes. In transient spectroscopy, one can also determine relaxation rates and transition frequencies, but the methodology is quite different. Atomic state populations or coherences between atomic states are excited by pulsed optical fields. Following the excitation, the time-evolution of the atoms is monitored, from which transition frequencies and relaxation rates can be obtained. In certain cases the transient response is studied as a function of incident field frequency or intensity. Whether or not transient or CW spectroscopy offers distinct advantages depends on a number of factors, such as signal to noise and the reliability of lineshape formulas.¹

In this chapter, we present basic concepts of coherent optical transient spectroscopy,²⁻¹² along with applications involving atomic vapors or condensed matter systems. Experimental techniques are discussed in Sec. 11.11. As in the case of CW spectroscopy, it will prove useful to consider both linear and nonlinear interactions of the atoms with the fields. The examples chosen to illustrate the concepts are relatively simple, but it is important to note that sophisticated coherent transient techniques can now be used to probe complex structures, such as liquids and semiconductors. Although we consider ensembles of atoms interacting with the applied fields, current technology allows one to study the transient response of single atoms or molecules.¹³

11.3 OPTICAL BLOCH EQUATIONS

Many of the important features of coherent optical transients can be illustrated by considering the interaction of a radiation field with a two-level atom, with lower state $|1\rangle$ having energy $-\hbar\omega_0/2$ and upper state $|2\rangle$ having energy $\hbar\omega_0/2$. For the moment, the atom is assumed to be fixed at $\mathbf{R} = 0$ and all relaxation processes are neglected. The incident electric field is

$$\mathbf{E}(\mathbf{R}=0, t) = \frac{1}{2}\boldsymbol{\epsilon}\{E(t)\exp[-i[\omega t - \phi(t)]] + E(t)\exp[i[\omega t - \phi(t)]]\} \quad (1)$$

where $E(t)$ is the field amplitude, $\boldsymbol{\epsilon}$ is the field polarization, $\phi(t)$ is the field phase, and ω is the carrier frequency. The time dependence of $E(t)$ allows one to consider pulses having arbitrary shape while the time dependence of $\phi(t)$ allows one to consider arbitrary frequency *chirps*. It is convenient to expand the atomic state wave function in a field interaction representation as

$$|\psi(t)\rangle = c_1(t)\exp\left\{\frac{i[\omega t - \phi(t)]}{2}\right\}|1\rangle + c_2(t)\exp\left\{-\frac{i[\omega t - \phi(t)]}{2}\right\}|2\rangle \quad (2)$$

The atom-field interaction potential is

$$V(\mathbf{R}, t) = -\boldsymbol{\mu} \cdot \mathbf{E}(\mathbf{R}, t) \quad (3)$$

where $\boldsymbol{\mu}$ is a dipole moment operator. Substituting the state vector Eq. (2) into Schrödinger's equation and neglecting rapidly varying terms (rotating-wave approximation), one finds that the state amplitudes evolve according to

$$i\hbar \frac{d\mathbf{c}}{dt} = \tilde{\mathbf{H}}\mathbf{c} \quad \tilde{\mathbf{H}} = \left(\frac{\hbar}{2}\right) \begin{pmatrix} -\delta(t) & \Omega_0(t) \\ \Omega_0(t) & \delta(t) \end{pmatrix} \quad (4)$$

where \mathbf{c} is a vector having components (c_1, c_2) ,

$$\delta = \omega_0 - \omega \quad (5)$$

is an atom-field detuning,

$$\Omega_0(t) = -\frac{\mu E(t)}{\hbar} = -\frac{\mu}{\hbar} \sqrt{\frac{2S(t)}{\epsilon_0 c}} \quad (6)$$

is a Rabi frequency, $\mu = \langle 1|\boldsymbol{\mu} \cdot \boldsymbol{\epsilon}|2\rangle = \langle 2|\boldsymbol{\mu} \cdot \boldsymbol{\epsilon}|1\rangle$ is a dipole moment matrix element, ϵ_0 is the permittivity of free space, $\mathbf{S}(t)$ is the time-averaged Poynting vector of the field, and

$$\delta(t) = \delta + \frac{d\phi(t)}{dt} \quad (7)$$

is a generalized atom-field detuning. Equation (4) can be solved numerically for arbitrary pulse envelope and phase factors.

Expectation values of physical observables are conveniently expressed in terms of density matrix elements defined by

$$\rho_{ij} = c_i c_j^* \quad (8)$$

which obey equations of motion

$$\begin{aligned}
 \dot{\rho}_{11} &= -\frac{i\Omega_0(t)}{2}(\rho_{21} - \rho_{12}) \\
 \dot{\rho}_{22} &= \frac{i\Omega_0(t)}{2}(\rho_{21} - \rho_{12}) \\
 \dot{\rho}_{12} &= -\frac{i\Omega_0(t)}{2}(\rho_{22} - \rho_{11}) + i\delta(t)\rho_{12} \\
 \dot{\rho}_{21} &= \frac{i\Omega_0(t)}{2}(\rho_{22} - \rho_{11}) - i\delta(t)\rho_{21}
 \end{aligned} \tag{9}$$

An alternative set of equations in terms of real variables can be obtained if one defines new parameters

$$\begin{aligned}
 u &= \rho_{12} + \rho_{21} & \rho_{12} &= \frac{u + iv}{2} \\
 v &= i(\rho_{21} - \rho_{12}) & \rho_{21} &= \frac{u - iv}{2} \\
 w &= \rho_{22} - \rho_{11} & \rho_{22} &= \frac{m + w}{2} \\
 m &= \rho_{11} + \rho_{22} & \rho_{11} &= \frac{m - w}{2}
 \end{aligned} \tag{10}$$

which obey the equations of motion

$$\begin{aligned}
 \dot{u} &= -\delta(t)v \\
 \dot{v} &= \delta(t)u - \Omega_0(t)w \\
 \dot{w} &= \Omega_0(t)v \\
 \dot{m} &= 0
 \end{aligned} \tag{11}$$

The last of these equations reflects the fact that $\rho_{11} + \rho_{22} = 1$, while the first three can be rewritten as

$$\dot{\mathbf{U}} = \mathbf{\Omega}(t) \times \mathbf{U} \tag{12}$$

where the *Bloch vector* \mathbf{U} has components (u, v, w) and the *pseudofield vector* $\mathbf{\Omega}(t)$ has components $[\Omega_0(t), 0, \delta(t)]$. An important feature of a density matrix description is that relaxation can be incorporated easily into density matrix equations, but *not* into amplitude equations.

Equations (9) or (11) constitute the *optical Bloch equations* without decay.^{8-10,14} The vector \mathbf{U} has unit magnitude and precesses about the pseudofield vector with an instantaneous rate

$$\Omega(t) = \sqrt{[\Omega_0(t)]^2 + [\delta(t)]^2} \tag{13}$$

that is referred to as the *generalized Rabi frequency*. The tip of the Bloch vector traces out a path on the *Bloch sphere*. The component w is the population difference of the two atomic states, while u and v are related to the quadrature components of the atomic polarization (see following discussion).

It is possible to generalize Eqs. (9) and (11) to include relaxation. In the most general situation, each density matrix element can be coupled to all density matrix elements via relaxation. For optical transitions,

however, it is often the case that the energy separation of levels 1 and 2 is sufficiently large to preclude any relaxational transfer of population from level 1 to level 2, although state $|2\rangle$ can decay to state $|1\rangle$ via spontaneous emission. For the present, we also assume that $\rho_{11} + \rho_{22} = 1$; there is no relaxation outside the two-level subspace. In this limit, relaxation can be included in Eq. (9) by modifying the equations as

$$\dot{\rho}_{11} = -\frac{i\Omega_0(t)}{2}(\rho_{21} - \rho_{12}) + \gamma_2\rho_{22} \quad (14a)$$

$$\dot{\rho}_{22} = \frac{i\Omega_0(t)}{2}(\rho_{21} - \rho_{12}) - \gamma_2\rho_{22} \quad (14b)$$

$$\dot{\rho}_{12} = -i\frac{\Omega_0(t)}{2}(\rho_{22} - \rho_{11}) - [\gamma - i\delta(t)]\rho_{12} \quad (14c)$$

$$\dot{\rho}_{21} = i\frac{\Omega_0(t)}{2}(\rho_{22} - \rho_{11}) - [\gamma + i\delta(t)]\rho_{21} \quad (14d)$$

where γ_2 is the spontaneous decay rate of level 2, γ is the real part of the decay rate of the coherence ρ_{12} , and the detuning

$$\delta(t) = \delta + \frac{d\phi(t)}{dt} - s \quad (15)$$

is modified to include the imaginary part s of the decay rate of the coherence ρ_{12} . The corresponding equations for the Bloch vector are

$$\begin{aligned} \dot{u} &= -\delta(t)v - \gamma u \\ \dot{v} &= \delta(t)u - \Omega_0(t)w - \gamma v \\ \dot{w} &= \Omega_0(t)v - \gamma_2(w+1) \\ \dot{m} &= 0 \end{aligned} \quad (16)$$

With the addition of decay, the length of the Bloch vector is no longer conserved.

One can write

$$\gamma = \frac{\gamma_2}{2} + \text{Re}(\Gamma_{12}) \quad s = \text{Im}(\Gamma_{12}) \quad (17)$$

where $\gamma_2/2$ is a radiative component and Γ_{12} is a complex decay parameter that could arise, for example, as a result of phase-interrupting collisions with a background gas. The quantity γ_2 is referred to as the *longitudinal relaxation rate*. Moreover, one usually refers to $T_1 = \gamma_2^{-1}$ as the *longitudinal relaxation time* and $T_2 = \gamma^{-1}$ as the *transverse relaxation time*. In the case of purely radiative broadening, $\gamma_2 = 2\gamma$ and $T_1 = T_2/2$.

The optical Bloch equations are easily generalized to include additional levels and additional fields. In particular, for an ensemble of three-level atoms interacting with two radiation fields, new

classes of coherent optical transient effects can appear. Moreover, the sensitivity of the coherent transients to the polarization of the applied fields offers an additional degree of selectivity in the detection of the transient signals. Some examples of coherent transient phenomena in three-level and multilevel systems can be found in the references.^{11,15–20} Chapter 14, “Electromagnetically Induced Transparency,” contains some interesting phenomena associated with such multilevel systems.

11.4 MAXWELL-BLOCH EQUATIONS

The optical Bloch equations must be coupled to Maxwell’s equations to determine in a self-consistent way the modification of the atoms by the fields and the fields by the atoms. To accomplish this task, we start with Maxwell’s equations, setting $\mathbf{D} = \epsilon_0 \mathbf{E} + \mathbf{P}$. The wave equation derived from Maxwell’s equations is

$$\nabla^2 \mathbf{E} - \nabla(\nabla \cdot \mathbf{E}) = \frac{1}{c^2} \frac{\partial^2 \mathbf{E}}{\partial t^2} + \frac{1}{\epsilon_0 c^2} \frac{\partial^2 \mathbf{P}}{\partial t^2} \quad (18)$$

As a result of atom-field interactions, it is assumed that a polarization is created in the medium of the form

$$\mathbf{P}(\mathbf{R}, t) = \frac{1}{2} \epsilon [P(\mathbf{R}, t) \exp i(kZ - \omega t) + P^*(\mathbf{R}, t) \exp -i(kZ - \omega t)] \quad (19)$$

which gives rise to a signal electric field of the form

$$E_s(\mathbf{R}, t) = \frac{1}{2} \epsilon [E_s(\mathbf{R}, t) \exp i(kZ - \omega t) + E_s^*(\mathbf{R}, t) \exp -i(kZ - \omega t)] \quad (20)$$

The Z axis has been chosen in the direction of \mathbf{k} . It is assumed that the complex field amplitudes $P(\mathbf{R}, t)$ and $E_s(\mathbf{R}, t)$ vary slowly in space compared with $\exp(ikZ)$ and slowly in time compared with $\exp(i\omega t)$. To simplify matters further, transverse effects such as self-trapping, self-focusing, diffraction, and ring formation^{21–25} are neglected. In other words, we take $P(\mathbf{R}, t)$ and $E_s(\mathbf{R}, t)$ to be functions of Z and t only, choose $\epsilon \cdot \mathbf{k} = 0$, and drop the $\nabla(\nabla \cdot \mathbf{E})$ term in the wave equation. When Eqs. (19) and (20) are substituted into the wave equation and terms of order $\partial^2 E_s(Z, t)/\partial t^2$, $\partial^2 E_s(Z, t)/\partial Z^2$, $\partial^2 P(Z, t)/\partial t^2$, and $\partial P(Z, t)/\partial t$ are neglected, one finds

$$2ik \left(\frac{\partial}{\partial Z} + \frac{\omega}{kc^2} \frac{\partial}{\partial t} \right) E_s(Z, t) - \left(k^2 - \frac{\omega^2}{c^2} \right) E_s(Z, t) = -\frac{\omega^2}{\epsilon_0 c^2} P(Z, t) \quad (21)$$

It is important to note that the polarization field acts as the source of the signal field. Consequently, the signal field does *not* satisfy Maxwell’s equations in vacuum, implying that there can be cases when $k \neq \omega/c$. For the moment, however, we assume that this *phase-matching* condition is met. Moreover, it is assumed that a quasi-steady state has been reached in which one can neglect the $\partial E_s(Z, t)/\partial t$ in Eq. (21). With these assumptions, Eq. (21) reduces to

$$\frac{\partial E_s(Z, t)}{\partial Z} = \frac{ik}{2\epsilon_0} P(Z, t) \quad (22)$$

Additional equations would be needed if the applied fields giving rise to the polarization of the medium are themselves modified to any extent by the signal field.

The polarization $P(Z, t)$ is the link between Maxwell's equations and the optical Bloch equations. The polarization is defined as the average dipole moment per unit volume, or

$$\begin{aligned} \mathbf{P}(\mathbf{R}, t) &= \text{Tr}(\rho^T \boldsymbol{\mu}) = \frac{1}{V} \text{Tr} \left[\sum_j \int d\mathbf{R}_j \rho^{T(j)}(\mathbf{R}_j, t) \boldsymbol{\mu} \delta(\mathbf{R} - \mathbf{R}_j) \right] \\ &= \mathcal{N} \boldsymbol{\mu} [\langle \rho_{12}^T(\mathbf{R}, t) \rangle + \langle \rho_{21}^T(\mathbf{R}, t) \rangle] \end{aligned} \quad (23)$$

where $\rho^{T(j)}(\mathbf{R}_j, t)$ and $\rho_{12}^T(\mathbf{R}, t)$ are single-particle density matrix elements, V is the sample volume, and \mathcal{N} is the atomic density. The superscript T indicates that these are "total" density matrix elements written in the Schrödinger representation rather than the field interaction representation. The relationship between the two is

$$\rho_{12}^T(\mathbf{R}, t) = \rho_{12}(Z, t) \exp[-i(kZ - \omega t)] \quad (24)$$

The angle brackets in Eq. (23) indicate that there may be additional averages that must be carried out. For example, in a vapor, there is a distribution of atomic velocities that must be summed over, while in a solid, there may be an inhomogeneous distribution of atomic frequencies owing to local strains in the media. By combining Eqs. (19), (23), and (24), one arrives at

$$\frac{\partial E_s(Z, t)}{\partial Z} = \frac{ik\mathcal{N}\boldsymbol{\mu}}{\epsilon_0} \langle \rho_{21}(Z, t) \rangle = \frac{ik\mathcal{N}\boldsymbol{\mu}}{2\epsilon_0} \langle u(z, t) - iv(z, t) \rangle \quad (25)$$

which, together with Eqs. (14) or (16), constitute the Maxwell-Bloch equations.

11.5 FREE POLARIZATION DECAY

As a first application of the Maxwell-Bloch equations, we consider free polarization decay (FPD),²⁶⁻³⁰ which is the analog of free induction decay (FID) in nuclear magnetic resonance (NMR).^{31,32} The basic idea behind FPD is very simple. An external field is applied to an ensemble of atoms and then removed. The field creates a phased array of atomic dipoles that radiate coherently in the direction of the incident applied field. The decay of the FPD signal provides information about the transverse relaxation times. We will discuss several possible scenarios for observing FPD. For the present, we assume that there is no inhomogeneous broadening of the sample (all atoms have the same frequency). Moreover, in this and all future examples, it is assumed that one can neglect any changes in the applied fields' amplitudes or phases as the fields propagate in the medium; the Rabi frequencies of the applied fields are functions of t only.

A short pulse is applied at $t = 0$, *short* meaning that

$$|\delta(t)| \tau, \gamma\tau, \gamma_2\tau \ll 1 \quad (26)$$

where τ is the pulse duration. The inequality Eq. (26) allows one to neglect any effects of detuning or relaxation during the pulse's action. Before the pulse arrives, the atom is in its ground state, implying that the components of the Bloch vector are $u = v = 0$, $w = -1$; that is, the Bloch vector points down (see Fig. 1a). *During* the pulse, the pseudofield vector can be approximated as $\boldsymbol{\Omega}(t) = [\boldsymbol{\Omega}_0(t), 0, 0]$, owing to Eq. (26). The Bloch vector precesses in the wv plane and reaches a final value following the pulse given by

$$u(0^+) = 0 \quad v(0^+) = \sin\theta \quad w(0^+) = -\cos\theta \quad (27)$$

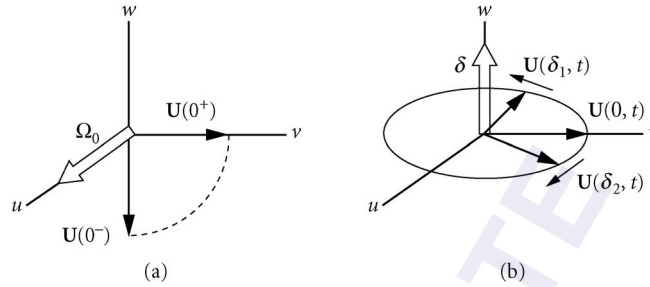


FIGURE 1 Evolution of the Bloch vector in free polarization decay. (a) A radiation pulse brings the Bloch vector from its initial position along the $-w$ axis to the uv plane. (b) With the field off, the Bloch vector precesses in the uv plane. In an inhomogeneously broadened sample, atoms having different detunings δ precess at different rates. The decay of the Bloch vector is not indicated in the figure.

where

$$\theta = \int_{-\infty}^{\infty} \Omega_0(t) dt \quad (28)$$

is a pulse area and 0^+ is a time just after the pulse. Following the pulse, the pseudofield vector is $\Omega = (0, 0, \delta)$, and the Bloch vector precesses about the w axis as it decays (see Fig. 1b). Explicitly, one finds

$$\begin{aligned} u(t) &= -\sin\theta \sin(\delta t) \exp(-\gamma t) \\ v(t) &= \sin\theta \cos(\delta t) \exp(-\gamma t) \\ w(t) &= -1 + (1 - \cos\theta) \exp(-\gamma_2 t) \end{aligned} \quad (29)$$

From this result, we can draw two conclusions. First, since the applied field is off when the atoms radiate, radiation is emitted at the natural frequency ω_0 . This conclusion follows formally from Eqs. (29), (10), (23), and (24), which can be used to show that $P(Z, t)$ varies as $\exp(-i\delta t)$, and both $\mathbf{P}(\mathbf{R}, t)$ and $\mathbf{E}(\mathbf{R}, t)$ oscillate at frequency $\omega + \delta = \omega_0$. Second, one can use Eqs. (29) and (25) to obtain

$$\frac{\partial E_s(Z, t)}{\partial Z} = \frac{kN\mu}{2\epsilon_0} \sin\theta \exp[-(\gamma + i\delta)t] \quad (30)$$

If the sample under consideration is optically thin, the power exiting a sample of length L in the \hat{Z} direction is proportional to

$$I(L, t) = |E_s(L, t)|^2 = \left(\frac{kN\mu L}{2\epsilon_0} \sin\theta \right)^2 \exp(-2\gamma t) \quad (31)$$

The signal is maximal for a pulse area of $\pi/2$. A measure of the output power as a function of time following the excitation pulse enables one to obtain a value for the transverse relaxation rate. For a pencil-like sample, the neglect of cooperative effects such as superradiance is based on the assumption that $NL/k^2 \ll 1$.^{25,33-35}

An alternative means for observing an FPD signal is to use an atomic beam that traverses a field interaction zone. The atom “sees” a radiation pulse in the atomic rest frame. If the atoms all have

the same longitudinal velocity u_0 , then the FPD signal measured at a distance L from the field interaction zone arises from atoms which were excited at time $t_e = t - L/u_0$. This implies that the phase factor $\exp[-(\gamma + i\delta)t]$ in Eq. (30) should be replaced by $\exp[-(\gamma + i\delta)(t - t_e)] = \exp[-(\gamma + i\delta)L/c]$. The emitted field is radiated at the incident field frequency ω rather than the atomic frequency ω_0 . If the intensity is monitored as a function of L , one can obtain the transverse relaxation rate.

Often the atoms or molecules are characterized by an inhomogeneous distribution of frequencies. In a solid, this can occur as a result of different strains in the host medium. In a vapor, the velocity distribution of the atoms is equivalent to a distribution of atomic transition frequencies, when viewed in the laboratory frame. To discuss both solids and vapors in the same context, we define

$$\delta \equiv \delta(\omega_0, \mathbf{v}) = \omega_0 - \omega + \mathbf{k} \cdot \mathbf{v} = \delta_0 + \Delta + \mathbf{k} \cdot \mathbf{v} \quad (32)$$

where

$$\delta_0 = \bar{\omega}_0 - \omega \quad \Delta = \omega_0 - \bar{\omega}_0 \quad (33)$$

and $\bar{\omega}_0$ is the average transition frequency. In a solid, $\mathbf{v} = 0$, but there is an inhomogeneous distribution of frequencies given by

$$W_f(\Delta) = \frac{1}{\sqrt{\pi}\sigma_w} \exp\left(-\frac{\Delta}{\sigma_w}\right)^2 \quad (34)$$

where σ_w characterizes the width of the inhomogeneous distribution. In a vapor, $\Delta = \sigma_w = 0$, but there is a Maxwellian velocity distribution

$$W_0(\mathbf{v}) = \frac{1}{(\pi u^2)^{3/2}} \exp[-(v/u)^2] \quad (35)$$

where u is the most probable atomic speed. The net effect is that Eq. (31) must be replaced by

$$\begin{aligned} I(L, t) &= \left(\frac{kN\mu L}{\epsilon_0}\right)^2 |\langle \rho_{21}(t) \rangle|^2 \\ &= \left(\frac{kN\mu L}{2\epsilon_0} \sin\theta\right)^2 \left| \int d\mathbf{v} W_0(\mathbf{v}) \int d\Delta W_f(\Delta) \exp(-[\gamma - i(\bar{\omega}_0 + \Delta + \mathbf{k} \cdot \mathbf{v})]t) \right|^2 \\ &= \left(\frac{kN\mu L}{2\epsilon_0} \sin\theta\right)^2 \exp(-2\gamma t) \exp\left[\frac{-(\sigma_w t)^2}{2}\right] \exp\left[\frac{-(kut)^2}{2}\right] \end{aligned} \quad (36)$$

If $\sigma_w \gg \gamma$ (solids) or $ku \gg \gamma$ (vapors), the signal decays mainly owing to inhomogeneous broadening. Bloch vectors corresponding to different frequencies precess about the w axis at different rates, implying that the optical dipoles created by the pulse lose their relative phase in a time of order $T_2^* = (2\sigma_w)^{-1}$ or $(2ku)^{-1}$ (see Fig. 1b). The FPD signal can be used to measure T_2^* , which can be viewed as an inhomogeneous, transverse relaxation time. At room temperature, ku/γ is typically on the order of 100. In a solid, σ_w/γ can be orders of magnitude larger. An experimental FPD signal is shown in Fig. 2.

It is also possible to produce an FPD signal by preparing the atoms with a CW laser field and suddenly turning off the field. This method was used by Brewer and coworkers in a series of experiments on coherent optical transients in which Stark fields were used to tune molecules in a vapor into and out of resonance.²⁶ The CW field modifies the velocity distribution for the molecules. In

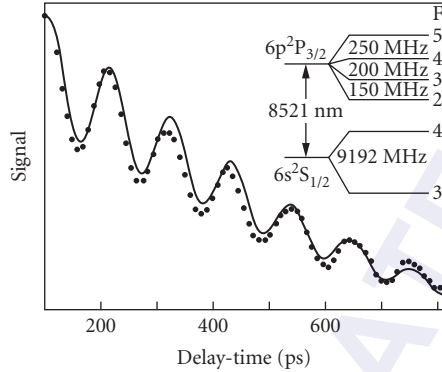


FIGURE 2 An FPD signal obtained on the D_2 transition in cesium. The excitation pulse has a duration of 20 ps. The decay time of the signal is determined by the inhomogeneous transverse relaxation rate $T_2^* = 1.4$ ns. Oscillations in the signal originate from the ground state hyperfine splitting. (From Ref. 28. Reprinted with permission.)

the linear field regime, this again leads to an FPD signal that decays on a time scale of order $(ku)^{-1}$. It is fairly easy to obtain this result. To first order in the field, the steady-state solution of Eq. (14d), generalized to include the Doppler shift $\mathbf{k} \cdot \mathbf{v}$ and initial velocity distribution $W_0(\mathbf{v})$, is

$$\rho_{21}(\mathbf{v}) = -i \left(\frac{\Omega_0}{2} \right) [\gamma - i(\delta_0 + \mathbf{k} \cdot \mathbf{v})]^{-1} W_0(\mathbf{v}) \quad (37)$$

With this initial condition, it follows from Eq. (14d) that, for $t > 0$,

$$\begin{aligned} \langle \rho_{21}(\mathbf{v}, t) \rangle &= \int d\mathbf{v} \frac{-i(\Omega_0/2) \exp[-\gamma t - i(\delta_0 + \mathbf{k} \cdot \mathbf{v})t]}{\gamma + i(\delta_0 + \mathbf{k} \cdot \mathbf{v})} W_0(\mathbf{v}) \\ &= -i \frac{\sqrt{\pi} \Omega_0}{2ku} \exp \left[\left(\frac{\gamma + i\delta_0}{ku} \right)^2 \right] \left[1 - \Phi \left(\frac{\gamma + i\delta_0}{ku} + \frac{kut}{2} \right) \right] \end{aligned} \quad (38)$$

where Φ is the error function. For $kut > 1$ and $|\gamma + i\delta_0|/ku \ll 1$, $\langle \rho_{21}(\mathbf{v}) \rangle \sim -i[\Omega_0/(k^2 u^2 t)] \exp(-k^2 u^2 t^2/4)$.

When one considers *nonlinear* interactions with the field, the situation changes. The CW field excites only those atoms having $\mathbf{k} \cdot \mathbf{v} = -\delta \pm \gamma'$, where γ' is a power-broadened homogeneous width. These velocity-selected atoms are no longer subject to inhomogeneous broadening and give rise to a contribution to the FPD signal that decays with rate $T_2^{-1} = (\gamma + \gamma')$.²⁶ Thus, by using *nonlinear* atom-field interactions, one can extract *homogeneous* decay rates in situations where there is large inhomogeneous broadening. The price one pays is that only a small percentage of the atoms (those that have velocities for which the applied field is resonant) contribute to the signal. An example of an FPD signal of this type is shown in Fig. 3.

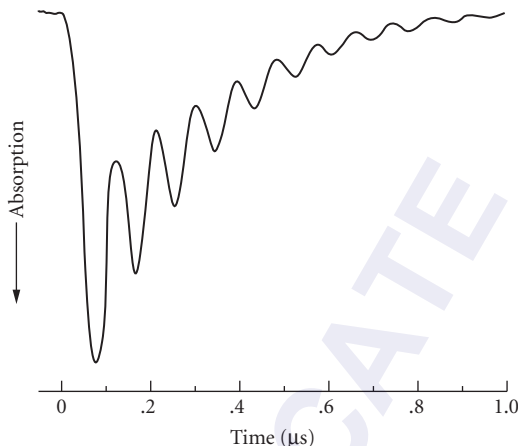


FIGURE 3 An FPD signal obtained in NH_2D at $10.6 \mu\text{m}$ using the method of Stark switching of a molecular transition frequency. Molecules that interact resonantly with a CW field are suddenly tuned out of resonance by the Stark switching. The oscillations result from the heterodyne detection method that is used, while the slowly varying increase in the signal is the result of optical nutation of molecules switched *into* resonance by the Stark pulse. The FPD signal manifests itself as a *reduction of the amplitude of the oscillation* with time. This amplitude decays with the (power-broadened) *homogeneous* decay rate T_2 , owing to the fact that the nonlinear interaction with a CW field results in the excitation of only a small fraction of the Doppler profile of the molecules. (From Ref. 26. Reprinted with permission.)

11.6 PHOTON ECHO

We have seen that it is possible to measure the transverse relaxation rate T_2 in an inhomogeneously broadened sample using FPD, but only a relatively small percentage of the atoms participate. The question arises as to whether other techniques would allow for full participation of the atoms. The *photon echo* is one such method.^{8–10,36–38} The photon echo has very little to do with either photons or echoes, but the name has a nice ring to it. The photon echo, the optical analog of the spin echo,³⁹ was first observed in ruby by Kurnit et al.³⁶ A pulse having propagation vector \mathbf{k}_1 is applied at $t = 0$, a second pulse having propagation vector $\mathbf{k}_2 = \mathbf{k}_1$ is applied at $t = T_{21}$, and an echo is radiated at time $t = 2T_{21}$ in a direction $\mathbf{k} = \mathbf{k}_1$. There are many ways to explain echo formation, and some of these are indicated in the following discussion.

In the Bloch vector picture, a $\pi/2$ pulse excites the optical dipoles at $t = 0$, bringing the Bloch vector along the v axis (Fig. 4a). The Bloch vector then begins to precess in the uv plane at a rate equal to the atom-field detuning. In an inhomogeneously broadened medium, different atoms have different resonant frequencies. As a consequence, the Bloch vectors associated with different atoms precess at different rates and dephase relative to each other in a time T_2^* (Fig. 4b). The dipole coherence is not lost, however. If at time T_{21} a π pulse is applied, the net effect of the pulse is to cause a reflection about the uv plane (Fig. 4c). As the atoms continue to precess at different rates (Fig. 4d), the rates are such that the Bloch vectors for *all* the atoms will become aligned with the $-v$ axis at time $t = 2T_{21}$ and an *echo signal* is emitted (Fig. 4e). From time $t = 0$ to $t = 2T_{21}$ the dipoles decay with the homogeneous decay rate γ .

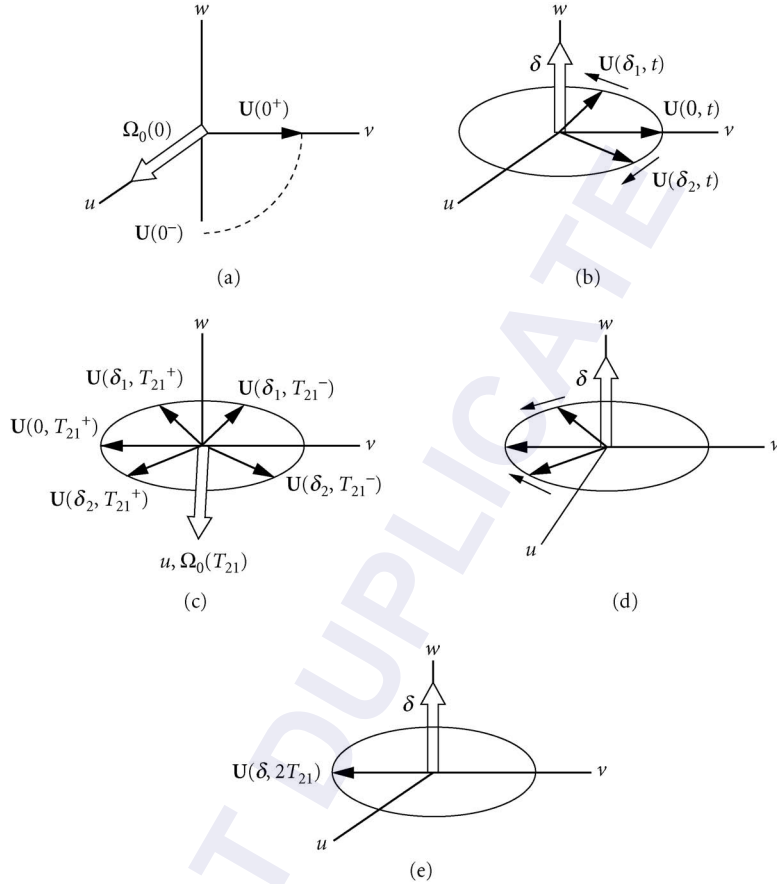


FIGURE 4 Bloch vector picture of echo formation. (a) An initial $\pi/2$ pulse brings the Bloch vector to the uv plane. (b) In a field-free region, the Bloch vector precesses in the uv plane; atoms with different detunings dephase relative to one another in a time equal to the inhomogeneous, transverse relaxation rate T_2^* . (c) At time T_{21} , a second pulse, this time a π pulse, reflects the Bloch vectors with respect to the uw plane. (d) In a field-free region, the Bloch vectors continue to precess. (e) At time $t = 2T_{21}$, all the vectors are aligned along the $-v$ axis (the optical dipoles have rephased), and an echo signal is emitted. The decay of the Bloch vector is not indicated in the figure.

By measuring the echo signal as a function of delay time T_{21} between the pulses, one can obtain the transverse relaxation time $T_2 = \gamma^{-1}$.

It is not necessary that the pulse areas be equal to $\pi/2$ and π , although these areas lead to a maximal signal. What is necessary is that the second pulse produce at least a partial reflection about the uw plane. This reflection takes the Bloch vector components $u + iv$ into $u - iv$, or, equivalently, takes density matrix element ρ_{12} into ρ_{21} . Since ρ_{12} and ρ_{21} are related to the real and imaginary parts of the average dipole moment operator, the second pulse must couple these real and imaginary parts. *Such coupling is impossible for a linear atom-field interaction.* Thus, by its very nature, the photon echo can occur only when nonlinear atom-field interactions are present.

An alternative way to picture echo formation is to use double-sided Feynman diagrams⁴⁰ that keep track of the relative phase of the different dipoles. Diagrams similar to those indicated in Fig. 5 were

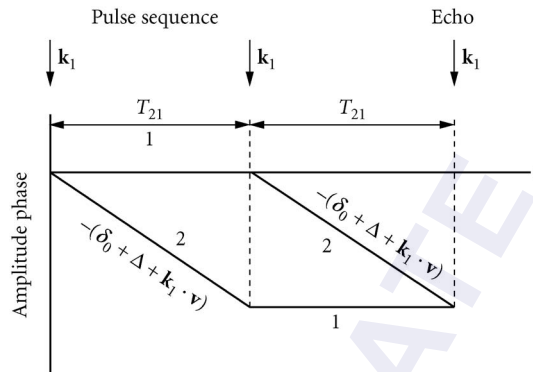


FIGURE 5 A phase diagram that can be used to analyze coherent transient phenomena. Each line corresponds to a state amplitude, and density matrix elements are obtained by multiplying a top line by the conjugate of a bottom line *at the same time*. The relative phase of a given density matrix element is given by the vertical separation of the two lines. The slopes of various line segments are shown on the graph. Lines corresponding to atoms having different atom-field detunings Δ (solid) or different velocities \mathbf{v} (vapor) would have different slopes. The pulse sequence is shown at the top of the figure along with the temporal position and direction of the echo. This diagram corresponds to the two-pulse photon echo. Regardless of the atom-field detuning or atomic velocity, all lines cross at $t = 2T_{21}$, indicating that all the optical dipoles are in phase at this point.

introduced by Hartmann and Friedberg in the context of a billiard ball echo model⁴¹ and have been used extensively in theories of atom interferometry.^{42–44} Each line represents a field amplitude. The abscissa is time, and the ordinate is the *phase* associated with the amplitude. In the absence of any interactions, it follows from Eq. (4) that the phase associated with the state $|1\rangle$ amplitude is $\delta t/2$ and that associated with the state $|2\rangle$ amplitude is $-\delta t/2$ so that the state $|1\rangle$ amplitude evolves without any phase change and the state $|2\rangle$ amplitude evolves with a phase equal to $-\delta t$.

We start with the atom in state $|1\rangle$ at $t = 0$. The atom-field interaction takes state $|1\rangle$ to $|2\rangle$ on absorption with a phase factor $\exp(i\mathbf{k}_1 \cdot \mathbf{R})$ and $|2\rangle$ to $|1\rangle$ on emission with a phase factor $\exp(-i\mathbf{k}_1 \cdot \mathbf{R})$. A vertical cut establishes the density matrix element of interest and the vertical distance between the two amplitudes is a measure of the *relative phase* of the amplitudes. For example, between $t = 0$ and $t = T_{21}$ in Fig. 5, the density matrix element ρ_{12} has been created with relative phase $(\delta_0 + \Delta + \mathbf{k}_1 \cdot \mathbf{v})t$ which grows with increasing t . One finds significant contributions to the dipole coherence at a given time *only* when the relative phase is the *same* for all the optical dipoles at that time. In a solid $\mathbf{v} = 0$, but $\Delta = \omega_0 - \bar{\omega}_0$ is different for different atoms owing to variations in ω_0 ; in a vapor $\Delta = 0$, but $\mathbf{k}_1 \cdot \mathbf{v}$ is different for different velocity subclasses of atoms. Thus, the slopes of the lines in Fig. 5 would differ for different atoms in both solids and vapors. On averaging over an inhomogeneous frequency distribution, the average dipole coherence would vanish, except at times near crossings of the state amplitudes, where the relative phase of all the dipoles is nearly equal to zero. Between $t = 0$ and $t = T_{21}$, this occurs only near $t = 0$, where an FPD signal is emitted. The application of a second pulse at $t = T_{21}$, however, converts ρ_{12} into ρ_{21} and begins a rephasing process for the dipoles. The state amplitudes in Fig. 5 intersect and the dipoles are rephased at $t = 2T_{21}$, *independent of the value of \mathbf{v} or Δ* . The echo signal is radiated for times $t \approx 2T_{21}$.

Analytical calculations of the signal intensity can be carried out using Eqs. (14) and (16). One simply pieces together periods in which the pulses act with periods of free evolution. The results are rather complicated, in general. However, if $\sigma_w T_{21} \gg 1$ (solid) or $k_1 u T_{21} \gg 1$ (vapor), only terms in the density matrix sequence indicated schematically in Fig. 5 survive the average over the inhomogeneous frequency distribution in the vicinity of the echo at $t \approx 2T_{21}$. At these times, one finds a total averaged density matrix element

$$\langle \rho_{21}^T(\mathbf{R}, t) \rangle = \langle \rho_{21}(t) \rangle \exp(i[kZ - \omega t])$$

where

$$\mathbf{k} = -\mathbf{k}_1 + 2\mathbf{k}_1 = \mathbf{k}_1$$

$$\begin{aligned} \langle \rho_{21}(t) \rangle &= \frac{1}{2} \sin \theta_1 \sin^2 \left(\frac{\theta_2}{2} \right) \exp(-2\gamma T_{21}) \exp[-i\delta_0(t - 2T_{21})] \\ &\times \exp \left[\frac{-\sigma_\omega^2(t - 2T_{21})^2}{4} \right] \exp \left[\frac{-k_1^2 u^2(t - 2T_{21})^2}{4} \right] \end{aligned} \quad (39)$$

and θ_i is the pulse area of pulse i . The corresponding echo intensity is

$$\begin{aligned} I(L, t) &= \left(\frac{kN\mu L}{\epsilon_0} \right)^2 |\langle \rho_{21}(t) \rangle|^2 = \left[\frac{kN\mu L}{2\epsilon_0} \sin \theta_1 \sin^2 \left(\frac{\theta_2}{2} \right) \right]^2 \\ &\times \exp(-4\gamma T_{21}) \exp \left[\frac{-\sigma_\omega^2(t - 2T_{21})^2}{2} \right] \exp \left[\frac{-k_1^2 u^2(t - 2T_{21})^2}{2} \right] \end{aligned} \quad (40)$$

It is interesting to note that the echo intensity near $t = 2T_{21}$ mirrors the FPD intensity immediately following the first pulse.

For experimental reasons it is often convenient to use a different propagation vector for the second pulse. Let \mathbf{k}_1 and \mathbf{k}_2 be the propagation vectors of the first and second pulses, which have identical carrier frequencies ω . In this case, one must modify the definition (24) of the field interaction representation to account for the different \mathbf{k} vectors. The final result for the total averaged density matrix element in the vicinity of the echo is

$$\langle \rho_{21}^T(\mathbf{R}, t) \rangle = \langle \rho_{21}(t) \rangle \exp(i[kZ - \omega t])$$

where

$$\mathbf{k} = 2\mathbf{k}_2 - \mathbf{k}_1$$

and

$$\begin{aligned} \langle \rho_{21}(t) \rangle &= \frac{i}{2} \sin \theta_1 \sin^2 \left(\frac{\theta_2}{2} \right) \exp(-2\gamma T_{21}) \exp[-i\delta_0(t - 2T_{21})] \\ &\exp \left[\frac{-\sigma_\omega^2(t - 2T_{21})^2}{4} \right] \int d\mathbf{v} W_0(\mathbf{v}) \exp\{i\mathbf{k}_1 \cdot \mathbf{v} T_{21} - \mathbf{k}_2 \cdot \mathbf{v}(t - T_{21})\} \end{aligned} \quad (41)$$

Recall that $\sigma_\omega = 0$ for a vapor and $W_0(\mathbf{v}) = \delta_D(\mathbf{v})$ for a solid, where δ_D is the Dirac delta function. In these equations there are three things to note. First, the signal is emitted in a direction different from that of the applied fields, a desirable feature from an experimental point of view. Second, the phase

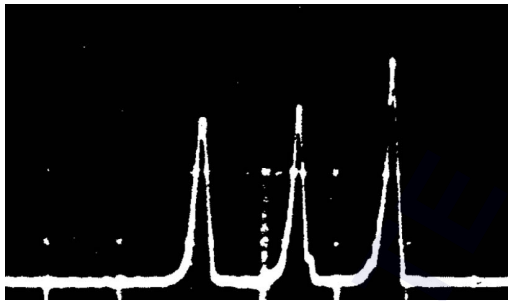


FIGURE 6 A photon echo signal from ruby. Time increases to the right with a scale of 100 ns/division. The pulse on the right is the echo signal, while the first two pulses are the (attenuated) input pulses. The echo appears at $t = 2T_{21}$ where T_{21} is the separation of the input pulses. (From Ref. 36. Reprinted with permission.)

matching condition $k = \omega/c$ is no longer satisfied since $k = |2\mathbf{k}_2 - \mathbf{k}_1|$ and $k_1 = k_2 = \omega/c$; however, if the fields are nearly collinear such that $(k^2 - \omega^2/c^2)L^2 \ll 1$, the effects of phase mismatch are negligible. Third, there is now a qualitative difference between the solid and vapor case. Owing to the fact that the detuning depends on the propagation vectors for the vapor, it is not possible to exactly rephase all the dipoles in the vapor when $\mathbf{k}_1 \neq \mathbf{k}_2$. If $|\mathbf{k}_2 - \mathbf{k}_1|/k_1 \ll 1$, however, nearly complete rephasing of the dipoles occurs for $t \approx 2T_{21}$. The photon echo signal observed by Kurnit et al.³⁶ is shown in Fig. 6.

As can be deduced from Fig. 5, the signal is sensitive only to off-diagonal density matrix elements in the entire time interval of interest. Thus, any disturbance of the off-diagonal density matrix elements or optical coherence will be reflected as a decrease in the echo intensity. As such, echo signals can serve as a probe of all contributions to transverse relaxation. Transverse relaxation generally falls into two broad categories which can lead to qualitatively different modifications of the coherent transient signals. First, there are *dephasing processes*, which produce an exponential damping of the coherences and contribute to γ . Second, there are *spectral diffusion* (solid)^{45–48} or *velocity-changing collisions* (vapor),^{47–49} which change the frequency associated with the optical coherences. Such terms enter the optical Bloch equations as *integral* terms, transforming the equations into differentio-integral equations. In a solid the change in frequency can be produced by fluctuating fields acting at each atomic site. Spectral diffusion of coherences in solids is difficult to detect, since phase-interrupting processes often dominate the signals. It has been measured in FPD using impurity ions in a crystalline host.⁵⁰ The situation in vapors is a bit more subtle. The phase-changing and velocity-changing aspects of collisions are entangled and cannot be separated, in general.⁴⁹ If the collisional interaction is *state independent*, however, as it is for some molecular transitions, then collisions are purely velocity changing in nature, leading to an echo that decays exponentially as T_{21}^3 for early times and T_{21} for later times.⁵¹ For electronic transitions, collisions are mainly phase changing in nature, but there is a velocity-changing contribution that persists in the forward diffractive-scattering cone. This diffractive scattering has been observed for Na–,⁵² Li–,⁵³ and Yb-rare gas collisions⁵⁴ using photon echo techniques.

11.7 STIMULATED PHOTON ECHO

Up to this point, we have considered pulse sequences that are useful for measuring transverse relaxation times. Now we examine *stimulated photon echoes*,^{8–10} which can be used to *simultaneously* measure both transverse and longitudinal relaxation times. Stimulated photon echoes have become an important diagnostic probe of relaxation in condensed-matter systems. The pulse sequence consists of three pulses, having

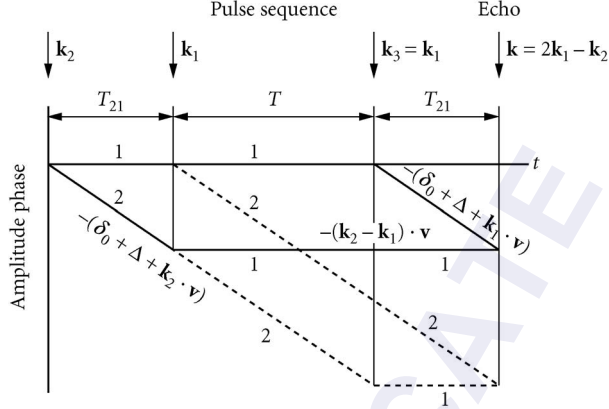


FIGURE 7 A phase diagram for the stimulated photon echo in which field 2 acts first and the echo is emitted in the $(-\mathbf{k}_2 + \mathbf{k}_1 + \mathbf{k}_3) = (2\mathbf{k}_1 - \mathbf{k}_2)$ direction. The solid lines involve the intermediate state population ρ_{11} while the dashed lines involve the intermediate state population ρ_{22} . All dipoles are in phase at $t = 2T_{21} + T$.

areas $\theta_1, \theta_2, \theta_3$ and propagation vectors $\mathbf{k}_1, \mathbf{k}_2, \mathbf{k}_3$, with $k_i = \omega/c$. In this section we take $\mathbf{k}_3 = \mathbf{k}_1$, and in the next section we will set $\mathbf{k}_3 = -\mathbf{k}_1$. The time interval between the first two pulses is T_{21} , and pulse 1 can precede or follow pulse 2. Pulse 3 occurs at time $t = T_{21} + T$ (Fig. 7). Signals can be generated in many different directions. For the sake of definiteness, we consider only the signal radiated in the $-\mathbf{k}_2 + \mathbf{k}_1 + \mathbf{k}_3$ direction. The phase diagram giving rise to this signal is shown in Fig. 7. For radiation to be emitted in the $-\mathbf{k}_2 + \mathbf{k}_1 + \mathbf{k}_3$ direction when $\mathbf{k}_3 = \mathbf{k}_1$ pulse 2 must precede pulse 1. (Of course, there are diagrams with pulse 1 preceding pulse 2, but these give rise to radiation in the $-\mathbf{k}_1 + \mathbf{k}_2 + \mathbf{k}_3$ direction.) It is assumed that T_{21} is greater than the inhomogeneous relaxation time T_2^* . The signal contains contributions from the optical coherence (off-diagonal density matrix elements) in the time intervals $(0, T_{21})$, $(T + T_{21}, t)$ and contributions from atomic state populations (diagonal density matrix elements) in the time interval T between the second and the third pulses. The echo appears when $t - (T_{21} + T) = T_{21}$.

The calculation of the echo signal is straightforward. Just before the second pulse, the density matrix $\rho_{12}(T_{21})$ varies as $\exp[-(\gamma - i\delta_2)T_{21}]$, where $\delta_1 = \omega_0 - \omega + \mathbf{k}_1 \cdot \mathbf{v}$. In the time interval T , the population difference $\omega = \rho_{22} - \rho_{11}$ decays at rate γ_2 and oscillates at frequency $\delta_2 - \delta_1 = (\mathbf{k}_2 - \mathbf{k}_1) \cdot \mathbf{v}$. In the final time interval, $\rho_{21}(t)$ varies as $\exp[-(\gamma + i\delta_3)(t - T - T_{21})]$, where $\delta_3 = \delta_1$ since $\mathbf{k}_3 = \mathbf{k}_1$. Combining the various field interaction and free propagation zones, one finds that the total averaged density matrix element in the vicinity of the echo $t \approx T + 2T_{21}$ is

$$\langle \rho_{21}^T(\mathbf{R}, t) \rangle = \langle \rho_{21}(t) \rangle \exp[i(kZ - \omega t)]$$

with

$$\mathbf{k} = -\mathbf{k}_2 + \mathbf{k}_1 + \mathbf{k}_3 = 2\mathbf{k}_1 - \mathbf{k}_2$$

and

$$\begin{aligned} \langle \rho_{21}(t) \rangle = & \left(\frac{i}{8} \right) \sin\theta_1 \sin\theta_2 \sin\theta_3 \sin(-\gamma_2 T) \exp(-2\gamma T_{21}) \exp\{-i\delta_0[t - T_{21} - T] - T_{21}\} \\ & \times \exp\left\{ \frac{-\sigma_\omega^2 [(t - T_{21} - T) - T_{21}]^2}{4} \right\} \int d\mathbf{v} W_0(\mathbf{v}) \\ & \times \exp\{i[\mathbf{k}_2 \cdot \mathbf{v} T_{21} + (\mathbf{k}_2 - \mathbf{k}_1) \cdot \mathbf{v} T - (2\mathbf{k}_1 - \mathbf{k}_2) \cdot \mathbf{v} (t - T_{21} - T)]\} \end{aligned} \quad (42)$$

The optimal pulse sequence consists of three $\pi/2$ pulses. Phase matching can be achieved only for $|\mathbf{k}_1 - \mathbf{k}_2|L \ll 1$. In a solid, the integral in Eq. (42) is equal to unity, and the echo signal is maximal for $t = T + 2T_{21}$. In a vapor, the echo signal is degraded if $\mathbf{k}_1 \neq \mathbf{k}_2$; however, if $\mathbf{k}_1 \approx \mathbf{k}_2$, then, at $t = T + 2T_{21}$, the echo amplitude varies as

$$\exp(-\gamma_2 T) \exp(-2\gamma T_{21}) \exp\left[-\frac{|\mathbf{k}_1 - \mathbf{k}_2|^2 u^2 (T + 2T_{21})^2}{4}\right]$$

By varying the angle between \mathbf{k}_1 and \mathbf{k}_2 , one can determine the Doppler width ku . By monitoring the echo signal as a function of $T_{21}(T)$, one obtains information on the transverse (longitudinal) relaxation.

Relaxation other than spontaneous emission can occur for the populations in the time interval T . The inhomogeneous phase $\delta_1 T_{21}$, acquired in the time interval T_{21} , is canceled by the phase $\delta_3 T_{21}$, acquired in the interval T_{21} following the third pulse. If, between the second and third pulses, the frequency (solid) or velocity (vapor) has changed owing to spectral diffusion (solid)^{45–48} or velocity-changing collisions (vapor),^{47–49,55,56} the phase cancellation will not be complete. Thus, the echo signal as a function of T provides information on these relaxation processes. The rate of spectral diffusion or velocity-changing collisions must be of order or greater than γ_2 to be observable. One would have a longer time to observe such effects if it were the ground-state lifetime rather than the excited-state lifetime that was the relevant time scale, but, in a closed two-level system, such is not the case.

The situation changes if a three-level system, such as the one shown in Fig. 8, is used. The fields drive the 1–2 transition, but level 2 decays to both level 1 and level 0. The total population of the 1–2 state subsystem is no longer conserved, requiring an *additional* decay rate to account for relaxation. Let us suppose that all states decay with rate Γ_t as a result of their finite time in the laser beams. Moreover, let $\Gamma_{2,1}$ and $\Gamma_{2,0}$ be the decay rates of level 2 to levels 1 and 0, respectively, such that $\gamma_1 = \gamma_0 = \Gamma_t$ and $\gamma_2 = \Gamma_{2,1} + \Gamma_{2,0} + \Gamma_t$. For simplicity, let us also take $\mathbf{k}_1 \approx \mathbf{k}_2$. In the interval T , the decay dynamics resulting from spontaneous emission and transit time effects is

$$\begin{aligned} \dot{\rho}_{22} &= -\gamma_2 \rho_{22} = -(\Gamma_{2,1} + \Gamma_{2,0} + \Gamma_t) \rho_{22} \\ \dot{\rho}_{11} &= -\Gamma_t \rho_{11} + \Gamma_{2,1} \rho_{22} \end{aligned}$$

Assuming that $\gamma_2 T \gg 1$, one finds that at time $t = T_{21} + T$, $\rho_{22}(T_{21} + T) \sim 0$, and

$$\rho_{11}(T_{21} + T) \sim \frac{\Gamma_{2,0}}{\Gamma_{2,1} + \Gamma_{2,0}} \exp(-\Gamma_t T)$$

which then replaces the factor $\exp(-\gamma_2 T)$ in Eq. (42). If $\Gamma_{2,0} \neq 0$, there is a long-lived component in the ground-state population. One can exploit this feature of *open systems* to study spectral diffusion or velocity-changing collisions with very high sensitivity.^{57,58} In Fig. 9, stimulated echo data is shown

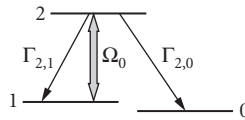


FIGURE 8 Open atomic-level scheme that can be used to observe transient signals limited only by some effective ground state lifetime. The field couples only states $|1\rangle$ and $|2\rangle$, but level 2 decays to both levels 1 and 0.

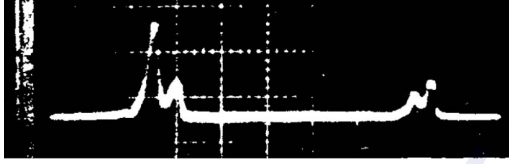


FIGURE 9 Stimulated photon echo observed on the D_1 transition in sodium. This is an “open” system for which a stimulated echo signal can be produced for separations T between the second and third pulses much larger than the excited-state lifetime. In this diagram T is 17 times the 16-ns lifetime of the excited state. The first three pulses are scattered light from the three input pulses and the fourth pulse is the echo. The echo appears at $t = T + 2T_{21}$, where T_{21} is the separation of the first two input pulses. (From Ref. 58. Reprinted with permission.)

that was used to measure a cross section for collisions between ground-state Na and He atoms.⁵⁸ The echo occurs for time separations T much greater than the excited-state lifetime.

Open systems also offer interesting possibilities as storage devices. Since the effective ground-state lifetime can be as long as days in certain solids, one can write interferometric information into the sample by replacing one of the first two pulses by a signal pulse and reading it out at a later time with the third pulse.^{59–62} In the case of vapors, it is also possible to replace some of the incident pulses by standing-wave fields.^{63–71} In this manner, modulated ground-state populations with associated Doppler phases of order kuT can be created and rephased, providing sensitivity to velocity-changing collisions as small as a few centimeters per second.⁷²

Before leaving this section, it is perhaps useful to make a slight digression on *homogeneously* broadened systems. A diagram similar to that shown in Fig. 10, *in which field 1 acts first*, also leads to a signal in the $\mathbf{k} = \mathbf{k}_1 - \mathbf{k}_2 + \mathbf{k}_3 = 2\mathbf{k}_1 - \mathbf{k}_2$, provided $\mathbf{k}_3 = \mathbf{k}_1$ (in the figure, $\mathbf{k}_3 = -\mathbf{k}_1$). We have not considered this contribution for inhomogeneously broadened systems since such a diagram leads to an overall phase of $\phi_d = -2(\delta_0 + \Delta + 2\mathbf{k}_1 \cdot \mathbf{v})T_{21}$ at time $t = T + 2T_{21}$. On averaging over either Δ or \mathbf{v} in an inhomogeneously broadened sample, this contribution would vanish. In a homogeneously broadened sample, however, $\Delta = 0$ and $\mathbf{v} = 0$, giving an *identical* relative phase ϕ_d to all the atoms. For $t > T + T_{21}$, the corresponding density matrix element associated with this diagram is

$$\langle \rho_{21}(t) \rangle = \left(\frac{i}{8} \right) \sin\theta_1 \sin\theta_2 \sin\theta_3 \exp(-\gamma_2 T) \exp[-\gamma(t - T)] \exp[-i\delta_0(t - T)] \quad (43)$$

Equation (43) does not constitute an echo in the usual sense, since there is no dephasing-rephasing cycle. The signal appears promptly (it is actually an FPD signal) following the third pulse. If one measures the *time-integrated* intensity in the signal following the third pulse, however (as is often the case with ultrafast pulses in which time resolution of the echo is not possible), it is impossible to tell directly whether an echo has occurred or not. For such measurements, a signal emitted in the $\mathbf{k} = 2\mathbf{k}_1 - \mathbf{k}_2$ direction when pulse 1 acts first is a clear signature of a homogeneously broadened system, since such a signal vanishes for inhomogeneously broadened samples.

The time-integrated signal is proportional to

$$\int_{T+T_{21}}^{\infty} |\langle \rho_{21}(t) \rangle|^2 dt$$

When field 2 acts first, the time integrated, inhomogeneously broadened signal varies as $\exp(-4\gamma T_{21})$, while the homogeneously broadened signal varies as $\exp(-2\gamma T_{21})$. When field 1 acts

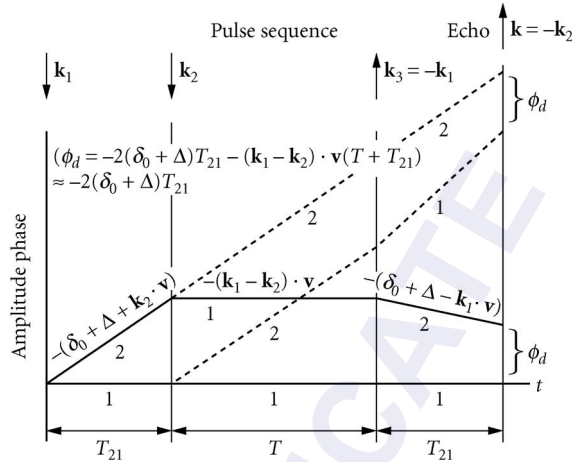


FIGURE 10 A phase diagram for the stimulated photon echo in which field 1 acts first. The direction of field 3 is opposite to that of field 1, and the echo is emitted in the $(\mathbf{k}_1 - \mathbf{k}_2 + \mathbf{k}_3) = -\mathbf{k}_2$ direction. We have taken $\mathbf{k}_2 \approx \mathbf{k}_1$. The solid lines involve the intermediate state population ρ_{11} , while the dashed lines involve the intermediate state population ρ_{22} . At time $t = 2T_{21} + T$, the relative phase is $\phi_d \approx -2(\delta_0 + \Delta)T_{21}$. For solids, the average over Δ washes out the signal. For vapors, $\Delta = 0$, the phase $\phi_d \approx -2\delta_0 T_{21}$ is the same for all the atoms, and an echo is emitted.

first, the signal is vanishingly small {varying as $\exp[-(\sigma_w^2 + k^2 u^2)T_{21}^2 / 2]$ } for inhomogeneously broadened samples, while the homogeneously broadened signal strength is essentially unchanged from that when pulse 2 acts first. This time-ordering asymmetry can be used to distinguish between homogeneously and inhomogeneously broadened samples.⁷³

11.8 PHASE CONJUGATE GEOMETRY AND OPTICAL RAMSEY FRINGES

A qualitative difference between stimulated photon echo signals arises when \mathbf{k}_3 is in the $-\mathbf{k}_1$ direction rather than the \mathbf{k}_1 direction.^{61,74,75} In this case, it is possible to generate a phase-matched signal in the

$$\mathbf{k} = \mathbf{k}_1 - \mathbf{k}_2 + \mathbf{k}_3 = -\mathbf{k}_2$$

direction for *both* time orderings of fields 1 and 2. Moreover, in weak fields, the amplitude of the signal field is proportional to the conjugate of input field 2. As a consequence, the signal is referred to as a *phase conjugate* signal for this geometry.^{76,77} To simplify matters, we will set $\mathbf{k}_1 \approx \mathbf{k}_2 = -\mathbf{k}$ and neglect terms of order $|\mathbf{k}_1 - \mathbf{k}_2|u(T + 2T_{21})$.

The appropriate phase diagrams are shown in Fig. 10 when field 1 acts before field 2 and Fig. 11 when field 2 acts before field 1. There is a qualitative difference between the phase diagrams of Fig. 10 and Fig. 7. At time $t = 2T_{21} + T$, the lines representing the state amplitudes *do not* cross in Fig. 10. Rather, they are separated by a phase difference of $\phi_d = -2(\delta_0 + \Delta)T_{21}$. The phase shift resulting from *Doppler shifts* cancels at $t = 2T_{21} + T$, but *not* the phase shift resulting from the atom-field detuning.

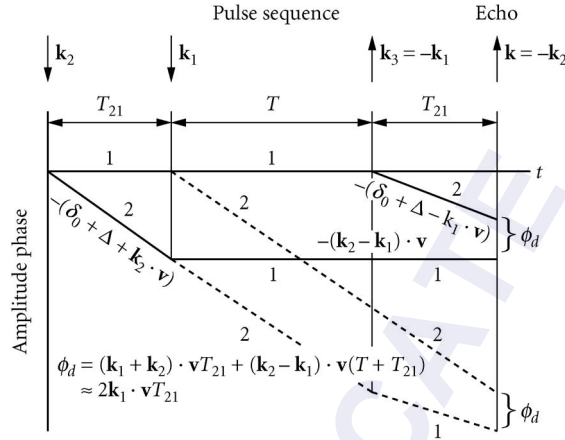


FIGURE 11 A phase diagram for the stimulated photon echo in which field 2 acts first and $\mathbf{k}_3 = -\mathbf{k}_1$. The echo is emitted in the $(-\mathbf{k}_2 + \mathbf{k}_1 + \mathbf{k}_3) = -\mathbf{k}_2$ direction, and we have taken $\mathbf{k}_2 \approx \mathbf{k}_1$. The solid lines involve the intermediate state population ρ_{11} , while the dashed lines involve the intermediate state population ρ_{22} . At time $t = 2T_{21} + T$, the relative phase is $\phi_d \approx 2\mathbf{k}_1 \cdot \mathbf{v}T_{21}$. For vapors, the average over the velocity distribution washes out the signal. For solids, $\mathbf{v} = 0$, the relative phase is zero for all the atoms, and an echo is emitted.

The significance of these results will become apparent immediately. The averaged density matrix element in the vicinity of the echo is

$$\begin{aligned} \langle \rho_{21}(t) \rangle &= \left(\frac{i}{8} \right) \sin\theta_1 \sin\theta_2 \sin\theta_3 \exp(-\gamma_2 T) \exp(-2\gamma T_{21}) \exp\{-i\delta_0(t - T_{21} - T) + T_{21}\} \\ &\times \exp\left[\frac{-\sigma_\omega^2 [(t - T_{21} - T) + T_{21}]^2}{4} \right] \exp\left[\frac{-k^2 u^2 [(t - T_{21} - T) - T_{21}]^2}{4} \right] \end{aligned} \quad (44)$$

In a solid, the signal vanishes near $t = 2T_{21} + T$, since $\sigma_\omega T_{21} \gg 1$.

In a vapor, an echo is formed at time $t = T + 2T_{21}$. At this time, the averaged density matrix element varies as $\exp(-2i\delta_0 T_{21})$, a factor which was absent for the nearly collinear geometry. This phase factor is the optical analog⁷⁸⁻⁸⁰ of the phase factor that is responsible for the generation of Ramsey fringes.⁸¹ One can measure the phase factor directly by heterodyning the signal field with a reference field, or by converting the off-diagonal density matrix element into a population by the addition of a *fourth pulse* in the \mathbf{k}_3 direction at time $t = T + 2T_{21}$.^{*} In either case, the signal varies as $\cos(2\delta_0 T_{21})$. In itself, this dependence is useless for determining the optical frequency since one cannot identify the fringe corresponding to $\delta_0 = 0$. To accomplish this identification, there are two possibilities. If the experiment is carried out using an atomic beam rather than atoms in a cell, $T_{21} = L/u_0$ will be different for atoms having different u_0 ($L =$ spatial separation of the first two pulses and u_0 is the

^{*}In the case of atoms moving through spatially separated fields with different longitudinal velocities, one must first average Eq. (44) over *longitudinal* velocities before taking the absolute square to calculate the radiated field. As a result of this averaging, the radiated signal intensity is maximum for $\delta_0 = 0$. Consequently, for spatially separated fields, heterodyne detection or a fourth field is not necessary since the radiated field intensity as a function of δ_0 allows one to determine the line center.

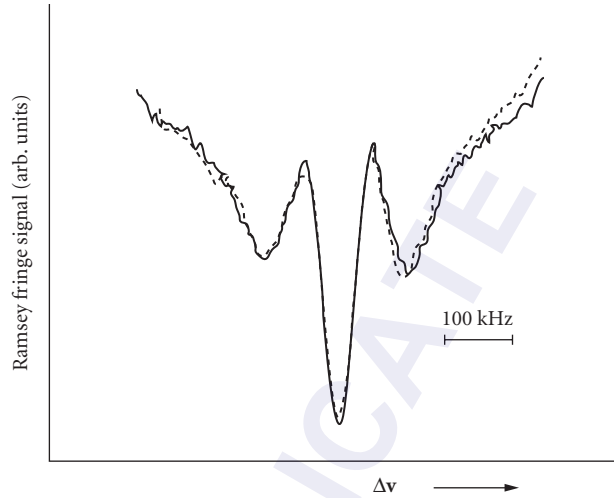


FIGURE 12 An optical Ramsey fringe signal on the 657-nm intercombination line in calcium. Four field zones were used. The most probable value of T_{21} was about 10^{-5} s for an effusive beam having a most probable longitudinal speed equal to 800 m/s, giving a central fringe width on the order of 60 kHz. The dashed and dotted lines represent runs with the directions of the laser field reversed, to investigate phase errors in the signals. (From Ref. 82. Copyright © 1994; reprinted with permission from Elsevier Science.)

longitudinal velocity of the atoms). When a distribution of u_0 is averaged over, the fringe having $\delta_0 = 0$ will have the maximum amplitude. Experiments of this type allow one to measure optical frequencies with accuracy of order T_{21}^{-1} (see Fig. 12). For experiments carried out using temporally separated pulses acting on atoms in a cell, it is necessary to take data as a function of δ_0 for several values of T_{21} , and then average the data over T_{21} ; in this manner the central fringe can be identified. The optical Ramsey fringe geometry has been reinterpreted as an atom interferometer.^{42,83} Atom interferometers are discussed in Sec. 11.9.

We now move to Fig. 11, in which field 2 acts before field 1. At time $t = 2T_{21} + T$, the lines representing the state amplitudes are separated by a phase difference of $\phi_d = -2\mathbf{k} \cdot \mathbf{v}T_{21}$. The phase shift resulting from the atom-field detuning cancels at $t = 2T_{21} + T$, but *not* the phase shift resulting from the Doppler effect. The corresponding density matrix element is

$$\begin{aligned} \langle \rho_{21}(t) \rangle &= \frac{i}{8} \sin\theta_1 \sin\theta_2 \sin\theta_3 \exp(-\gamma_2 T) \exp(-2\gamma T_{21}) \exp\{-i\delta_0[(t - T_{21} - T) - T_{21}]\} \\ &\times \exp\left[\frac{-\sigma_\omega^2[(t - T_{21} - T) - T_{21}]^2}{4}\right] \exp\left[\frac{-k^2 u^2[(t - T_{21} - T) + T_{21}]^2}{4}\right] \end{aligned} \quad (45)$$

Near $t = T + 2T_{21}$, the signal vanishes for a vapor since $kuT_{21} \gg 1$, but gives rise to a phase conjugate signal in solids ($u = 0$). There are no Ramsey fringes in this geometry; optical Ramsey fringes cannot be generated in an inhomogeneously broadened solid.

When the first two fields are *identical*, there is no way to distinguish which field acts first, and Eqs. (44) and (45) must be added before taking the absolute square to determine the radiated electric field. There is no interference between the two terms, however, since one of the terms is approximately equal to zero in the vicinity of the echo for either the solid or the vapor.

11.9 TWO-PHOTON TRANSITIONS AND ATOM INTERFEROMETRY

In the previous section, we have already alluded to the fact that optical Ramsey fringes can serve as the basis of an atom interferometer.⁴⁴ There is some disagreement in the literature as to exactly what constitutes an atom interferometer. Ramsey fringes and optical Ramsey fringes were developed without any reference to quantization of the atoms' center-of-mass motion. As such, optical Ramsey fringes can be observed in situations where quantization of the center-of-mass motion is irrelevant. The interference observed in these interferometers is based on an *internal state* coherence of the atoms. Matter-wave effects (that is, effects related to quantization of the center-of-mass motion) may play a role under certain circumstances, but they are not critical to the basic operating principle associated with optical Ramsey fringes.

In this section, we consider a time-domain, matter-wave atom interferometer⁸⁴ which relies on the wave nature of the center-of-mass motion for its operation. Moreover, the interferometer illustrates some interesting features of coherent optical transients not found in NMR. We return to an ensemble of two-level atoms, which have been cooled in a magneto-optical trap. See Chap. 20, "Laser Cooling and Trapping of Atoms," for a more detailed description of trapping of atoms. The atoms are subjected to two *standing-wave* optical pulses separated in time by T . The electric field amplitude of pulse i ($i = 1, 2$) is given by $\mathbf{E}_i(\mathbf{Z}, t) = \mathcal{E}E_i(t) \cos(kZ) \cos(\omega t)$. Either off-resonant^{84,85} or resonant⁸⁶ pulses can be used. For resonant pulses, grating echoes can be observed in situations where a classical description of the center-of-mass motion is valid.⁶⁶⁻⁶⁸ We consider only off-resonant pulses in this discussion, for which echoes can occur only when quantized motion of the atoms is included. For an atom-field detuning $|\delta| \gg \Omega_0, \gamma, \gamma_2, ku$, it is possible to adiabatically eliminate the excited state amplitude and arrive at an effective hamiltonian for the ground state atoms given by

$$H = \frac{P^2}{2M} - \sum_{i=1,2} \hbar \Omega_i^{(s)}(t) \cos(2kZ)$$

where \mathbf{P} is the center-of-mass momentum operator, M is the atomic mass, and

$$\Omega_i^{(s)}(t) = \frac{\mu^2 E_i^2(t)}{8\hbar^2 \delta}$$

is a two-photon Rabi frequency. A spatially homogeneous term has been dropped from the hamiltonian.

The net effect of the field is to produce a spatially modulated, AC Stark or light shift of the ground state energy. Let us assume that the pulse duration τ is sufficiently short to ensure that $\tau^{-1} \gg \omega_{2k}, \gamma, \gamma_2, ku, \sqrt{\omega_{2k} \Omega_i^{(s)}(t)}$, where

$$\omega_{2k} = \frac{\hbar(2k)^2}{2M}$$

is a two-photon recoil frequency whose importance will become apparent. In this limit, any motion of the atoms during the pulses can be neglected. The net effect of pulse i is to produce a ground state amplitude that varies as $\exp[i\theta_i^{(s)} \cos(2kZ)]$, where $\theta_i^{(s)} = \int \Omega_i^{(s)}(t) dt$ is a pulse area. In other words, the standing-wave field acts as a *phase grating* for the atoms. One can think of the two traveling-wave

components of the standing-wave field exchanging momentum via the atoms. All even integral multiples of $2\hbar k$ can be exchanged by the fields, imparting impulsive momenta of $2n\hbar k$ (n is a positive or negative integer) to the atoms. The frequency change associated with this momentum change for an atom having momentum \mathbf{P} is

$$\frac{E_{\mathbf{P}, \mathbf{P} \pm 2n\hbar k \hat{z}}}{\hbar} = \frac{\left(\frac{p^2}{2M} \frac{|\mathbf{P} \pm 2n\hbar k \hat{z}|^2}{2M} \right)}{\hbar} = \mp \frac{2n\mathbf{P}_z k}{M} - \omega_{2nk}$$

and consists of two parts. The first part is independent of \hbar and represents a classical Doppler shift, while the second part is proportional to \hbar and represents a quantum matter-wave effect. The quantum contribution will become important for times of order ω_{2nk}^{-1} . In other words, following the first pulse, the atomic density will remain approximately constant for times $t < \omega_{2nk}^{-1}$ (the maximum value of n is the larger of $\theta^{(s)}$ and unity). For times larger than this, the quantum evolution of the center-of-mass motion can transform the phase grating into an amplitude grating which can be deposited on a substrate or probed with optical fields. In contrast to closed two-level systems, the signals can persist here for arbitrarily long times. The recoil associated with absorption and emission “opens” the system and allows for long-lived transients.⁸⁷

The evolution of the system can be followed using phase diagrams in a manner similar to that used in Secs. 11.6 to 11.8. The situation is more complex, however, since a standing-wave field generates an infinity of different phase shifts $\mp 2nkvt$, where $v = P_z/M$. Details of the calculation can be found in the article by Cahn et al.⁸⁴ Here we sketch the general idea. The first pulse creates all even spatial harmonics of the fields, with weighting functions that are Bessel functions of the pulse area. The atomic density remains constant until a time $t \sim \omega_{2nk}^{-1}$. At this time one would expect to find a spatially modulated atomic density; however, if $ku \gg \omega_{2k}$ as is assumed, by the time the spatial modulation is established, the modulation is totally destroyed as a result of Doppler dephasing. As in the photon echo experiment, the Doppler dephasing can be reversed by the second pulse at time $t = T$. Since standing waves are used, there is an infinity of echo positions possible, corresponding to different dephasing-rephasing conditions for the various momentum components created by the fields.

Of the many echoes that can be produced, we consider only those echoes that are formed at times $t_N = (N+1)T$, $N = 1, 2, \dots$. Moreover, in an expansion of the atomic density in harmonics of the field, we keep only the second harmonic, since it can be probed by sending in a traveling-wave field and observing a *backscattered* signal. Phase matching is automatically guaranteed for the backscattered signal. For times $t = t_N + t_d$, with $t_d \gtrsim \frac{1}{2}ku \ll T$, the backscattered electric field amplitude varies as⁸⁴

$$E_b(t_d, NT) = \exp\left[\frac{-(2k)^2 u^2 (t_d)^2}{4}\right] J_N[2\theta_1 \sin(\omega_{2k} t_d)] J_{N+1}[2\theta_2 \sin(N\omega_{2k} T + \omega_{2k} t_d)] \quad (46)$$

where the J s are Bessel functions. The electric field can be measured using a heterodyne technique. The experimental data is shown as a function of t_d and T for $N = 1, 2$ in Figs. 13 and 14. One sees that the signal vanishes identically at the echo times, but not in the immediate vicinity of the echo points. The uniform atomic density at the echo point mirrors the uniform atomic density immediately after excitation by the first pulse. As a function of T , the signal is periodic with period $\pi/N\omega_{2k}$ for N odd and $2\pi/N\omega_{2k}$ for N even. By measuring the period, one can obtain values for \hbar/M . The interferometer can also be used to measure inertial effects such as the acceleration of the atoms owing to gravity. The advantage of this interferometer is that large interaction times are possible—one is limited only by the time it takes for the atoms to leave the atom-field interaction zone.

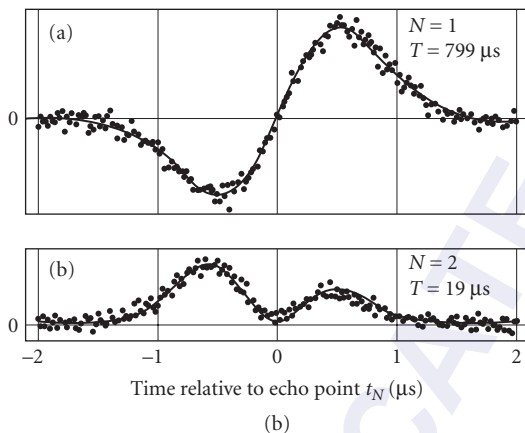


FIGURE 13 A time-domain atom interferometer. Two off-resonant, standing-wave optical pulses separated in time by T are applied to rubidium atoms in a magneto-optical trap and a probe field is applied near (a) $t = 2T$ or (b) $t = 3T$, giving rise to backscattered electric field signals. The electric field amplitude $E_b(t_p, NT)$ is recorded in the graphs as a function of t_p , the time from the echo position. The solid line is theory and the dots are experimental points. Note that the time delay between pulses is $799 \mu\text{s}$ in (a), indicating that these ground-state transients are limited only by some effective ground-state lifetime. (From Ref. 84. Reprinted with permission.)

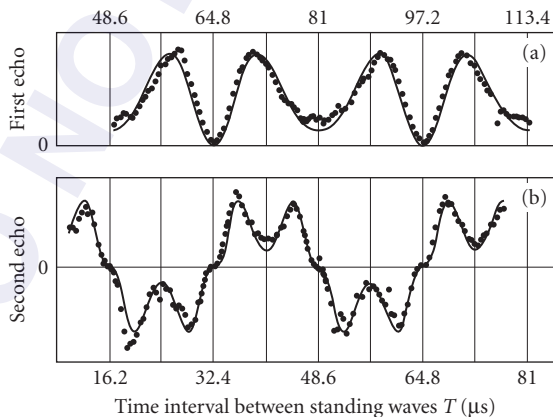


FIGURE 14 Same as in Fig. 13, but $E_b(t_p, NT)$ is now recorded as a function of T , the time separation of the pulses. The period of the signals is $\pi/\omega_{2k} = 32.39 \mu\text{s}$. (From Ref. 84. Reprinted with permission.)

11.10 CHIRPED PULSE EXCITATION

The discussion in this chapter has focused on transform-limited optical pulses, possessing smooth Fourier transforms centered about the carrier frequency. Alternative pulse shapes offer new and interesting possibilities. If one uses *stochastic* pulse envelope functions,⁸⁸⁻⁹⁰ the correlation time τ_c associated with the pulse is much smaller than the pulse duration. In some sense, *each* pulse can be viewed as a sequence of pulses having duration of order τ_c . As a consequence, stimulated photon echoes using stochastic pulses can be used to measure relaxation times as short as τ_c rather than the pulse duration.

The idea of using a pulse whose effective coherence time is shorter than the pulse duration has been exploited by others⁹¹⁻⁹³ in schemes involving chirped pulses. It was shown both theoretically⁹¹ and experimentally^{92,93} that, by sweeping the pulse frequency, one can write and read data encoded in solids, using the equivalent of stimulated echoes, a process Mossberg⁹¹ refers to as *swept-carrier time-domain optical memory*. Without going into the mathematical details of the calculations needed to arrive at expressions for the signals, we present the underlying physical concepts pertinent to this excitation scheme.

There are three pulses, as in a traditional stimulated photon echo, but the pulse characteristics differ markedly from those discussed in Sec. 11.7. The first pulse is a reference pulse having electric field

$$\mathbf{E}_1(\mathbf{R}, t) = \boldsymbol{\epsilon} E_1(t) \cos(\mathbf{k}_1 \cdot \mathbf{R} - \bar{\omega}_0 t - \xi t^2)$$

where the amplitude $E_1(t)$ is a smooth function of t centered at $t = 0$ having temporal width τ_{ref} and $\phi(t) = -\xi t^2$ is the pulse phase. The central frequency of the pulse coincides with the optical frequency, but the frequency is chirped at rate 2ξ , giving an instantaneous frequency $\omega(t) = \omega - \dot{\phi} = \omega + 2\xi t$ and an atom-field detuning $\delta(t) = \Delta - 2\xi t$. The frequency shifts $\pm \xi \tau_{\text{ref}}$ are assumed to be less than the inhomogeneous width σ_w . As the frequency is scanned, different atoms in the inhomogeneous distribution in the sample come into resonance with the field at different times. If $\sqrt{\xi} \gg \gamma, \tau_{\text{ref}}^{-1}$, for atoms having detuning $\Delta = \omega_0 - \omega$, the field comes into resonance at time $t_1 = \Delta/2\xi$ for a duration of order $\xi^{-1/2} \ll \tau_{\text{ref}}$. Thus, for each frequency group of atoms, the field acts as a *pulse having temporal width much less than the width of the pulse*. In calculating ρ_{12} resulting from this pulse one finds a phase factor of the form $\exp[i(-\mathbf{k}_1 \cdot \mathbf{R} + \Delta t_1 - \xi t_1^2)] = \exp[-i(\mathbf{k}_1 \cdot \mathbf{R} - \Delta^2/4\xi)]$.

The second or data pulse

$$\mathbf{E}_2(\mathbf{R}, t) = \boldsymbol{\epsilon} E_2(t) \cos[\mathbf{k}_2 \cdot \mathbf{R} - (\bar{\omega}_0 - \omega_{\text{RD}})t - \xi t^2]$$

is similar to the first except that it is offset from the first by frequency ω_{RD} . Moreover, the field amplitude $E_2(t)$ is now assumed to consist of a sequence of input data, such as a number of individual pulses contained in the overall pulse envelope. If the Fourier spectrum of $E_2(t)$ contains frequency components ω_f , then the second pulse will come into resonance with atoms having detuning Δ at time $t_2 = (\Delta + \omega_{\text{RD}} - \omega_f)/2\xi$. To have pulse 2 act on the same atoms at a time greater than t_1 , one must restrict the maximum value of ω_f to be less than ω_{RD} . We shall neglect ω_f in what follows. For each frequency group, the first two pulses act as a *sequence of short pulses*, separated in time by

$$T_{21} = \omega_{\text{RD}}/2\xi$$

As in the normal stimulated photon echo, the second pulse converts the density matrix element ρ_{12} created by the first pulse into population. One finds a population difference $(\rho_{22} - \rho_{11})$ that varies as $\exp[-i(\mathbf{k}_1 \cdot \mathbf{R} - \Delta^2/4\xi)] \exp(-\gamma T_{21}) \exp\{i[\mathbf{k}_2 \cdot \mathbf{R} - (\Delta + \omega_{\text{RD}})^2/4\xi]\} + \text{c.c.}$ Although the pulse durations are of order τ_{ref} , homogeneous decay occurs only on a time scale $T_{21} = \omega_{\text{RD}}/2\xi \ll \tau_{\text{ref}}$ which is

the effective pulse separation for a given frequency group. The phase factor $\exp\{i[(\mathbf{k}_2 - \mathbf{k}_1) \cdot \mathbf{R} - \Delta T_{21}]\}$ associated with $(\rho_{22} - \rho_{11})$ is identical in form to that found in the normal stimulated photon echo. It is not surprising then that a third pulse, identical to the first but propagating in the \mathbf{k}_3 direction and displaced in time from the first by $T > \tau_{\text{ref}}$ leads to a reconstruction of the data pulse propagating in the $(-\mathbf{k}_1 + \mathbf{k}_2 + \mathbf{k}_3)$ direction. In other words, on averaging over the inhomogeneous frequency distribution, one finds contributions to the signal only for those times that correspond to the time sequence of pulse 2, displaced by time T .

11.11 EXPERIMENTAL CONSIDERATIONS

Early coherent optical transient experiments were the optical analogs of NMR experiments. These experiments firmly established coherent transient spectroscopy as a viable technique in the optical domain. The relatively simple atomic and molecular vapors or rare earth-doped crystals that were used for these studies were chosen for practical reasons—their transition frequencies coincided with available laser frequencies and relaxation times (T_1 is typically on the order of 10 ns or longer) were longer than the lasers' pulse widths. The methodology is well described in numerous earlier reviews (see, for example, Levenson⁹). These relaxation time scales are relatively long by today's standards. Hence, the technology that was available based on photomultipliers, high-speed diodes, and fast oscilloscopes made it possible to observe the coherent transient phenomena that were created by fast laser-frequency or Stark switching. With high-speed detectors, photon echoes were readily observed with single-shot Q-switched lasers.

The continued advances in the development of ultrafast lasers have now reduced the pulse widths by nearly 6 orders of magnitude compared to the early Q-switched lasers. In addition, modern laser systems are characterized by high repetition rates (100 MHz if no amplification is used) compared to the relatively slow repetition rates of older laser systems (1 to 10 Hz). Ultrafast lasers have opened a host of possibilities for studying complex molecules, fluids, and solids, including semiconductors. This new capability was accompanied by new challenges, since standard detectors and electronics were not capable of time-resolving the emitted signals. In some cases, it was not even known whether the materials being investigated were homogeneously or inhomogeneously broadened. For these cases, it is important to check for the asymmetry predicted as a function of time delay in the stimulated photon echo, as described in Sec. 11.7. For the very shortest pulses (typically < 100 fs), even the simplest laboratory operations of reflection from a mirror or transmission through a cryostat window or beam splitter become an issue. The transform-limited pulse bandwidth is so large that the linear dispersion in these systems leads to a chirp in the pulse, which can give rise to artifacts in the data if the chirp is not compensated at the sample by incorporation of grating or prism pairs in the system. High-repetition-rate systems also can give rise to thermal heating of the sample, leading to gratings which give signals that easily dominate the electronic signals of interest. Discrimination against these signals is critical. Sometimes, it is possible to perform the measurements of interest using orthogonally polarized fields, giving rise to a signal that is sensitive to spatially modulated magnetic-state coherence, but not thermal gratings. An alternative approach is to amplitude modulate one of the optical fields at a high frequency and use phase-sensitive detection. If the modulation frequency chosen is sufficiently high, then the modulation of the thermal grating is weak, and the electronic term dominates.

In the case of photon echo spectroscopy, determining the time origin $t = 0$ is important for accurate analysis of the signals. Although this does not pose any serious technical problems when dealing with nanosecond time-scale resolution, for femtosecond laser pulses, where distance scales can be as small as a few microns, the problem is not trivial. One solution, discussed in a recent review of photon echo spectroscopy,⁷³ is to measure the time-integrated signal as a function of pulse separation for the ordinary signal (in the $2\mathbf{k}_1 - \mathbf{k}_2$ direction for self-diffracted four-wave mixing) and the complimentary echo (in the $2\mathbf{k}_1 - \mathbf{k}_3$ direction). The intersection of the two superimposed mirror images allows one to determine the time origin. In some cases, it may be necessary to time-resolve the emission either to confirm absolutely that the signal is indeed an echo or to determine the

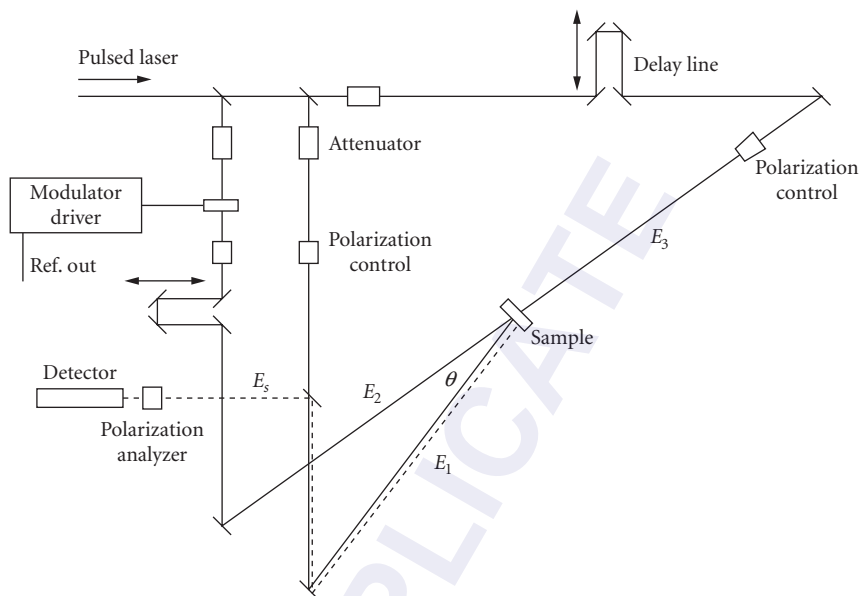


FIGURE 15 A typical experimental configuration, employing a high-repetition short-pulse laser system, that can be used for stimulated photon echo or transient four-wave mixing studies of an arbitrary sample.

inhomogeneous broadening. The usual approach to achieve this goal is to mix some of the original laser beam with the signal beam in a second harmonic crystal and detect the upconverted signal. The signal is time resolved by measuring the upconverted signal as a function of delay between the signal and reference fields.^{94,95}

An experimental configuration incorporating many of these features is shown in Fig. 15. Measurements of the dephasing rate are made by monitoring the signal amplitude as a function of time delay between fields E_1 and E_2 , while measurements of the energy relaxation rate are made as a function of the time delay between fields E_2 and E_3 . In the latter case, measurement of the relaxation rate as a function of the angle between fields E_1 and E_2 allows one to measure the grating (e.g., spatially modulated population) relaxation rate, due, for example, to diffusion.⁹⁶ Control of the fields' polarizations is essential in these experiments since the allowed electronic excitation depends critically on the selection rules. The modulator is used to amplitude modulate one of the optical fields to allow phase-sensitive detection. Because of scattering, it may also be necessary in some cases to modulate field E_1 at a different frequency and detect the sum or difference frequency in the lock-in amplifier. The forward three-pulse geometry is most prevalent in the literature; however, the more recently developed phase conjugate geometry shown in the figure is more desirable, owing to reduced incoherent scattering into the detector. In addition, this system is much easier to align, since the signal is exactly counterpropagating with respect to field E_2 . Usually, the feedback into the oscillator for spectroscopy applications is small because of the presence of the attenuators. However, for high-intensity studies it may be necessary to add an optical isolator to eliminate feedback, or, as is more typical, simply arrange for a slight misalignment of fields E_1 and E_3 . This results in a slight deviation of the signal field from the $-\mathbf{k}_2$ direction, but this poses no problem. It should also be noted that for ultrafast laser systems, the narrow pulses may experience group velocity dispersion (GVD) in propagating through various optical components and in reflection, requiring GVD compensation to avoid artifacts due to frequency chirps. In some cases, additional information can be obtained by time-resolving the optical phase⁹⁷ and/or polarization⁹⁸ of the signal, both of which may change during the time interval in which the signal is emitted.

Atoms undergoing collisions in a vapor and rare earth-doped solids provided a convenient laboratory environment both to study coherent transients and to examine new physical phenomena associated with reservoir interactions. Moreover, such studies allowed one to gain theoretical and experimental insights upon which the new technology could build. It has now been possible to elevate the utility of these spectroscopic tools for application to more complex systems. In particular, coherent optical transients have yielded important data in both chemical physics^{99,100} and condensed-matter physics.¹⁰¹ In many cases, the theoretical framework described here, with appropriate refinements of the unperturbed hamiltonian and reservoir interactions, can adequately model these systems. For example, local field corrections and excitation-induced scattering can be accounted for by a relatively simple modification of the optical Bloch equations. As a result, additional phase diagrams, with time orderings different from those shown in Figs. 10 and 11, can contribute to echo signals.¹⁰² However, in other cases, a major revision of the theoretical picture is needed to account for more complex many-body interactions.^{103–105} An approach based on modified optical Bloch equation may not be appropriate in this limit.

11.12 CONCLUSION

Information on level structure and relaxation processes in vapors, solids, and liquids can be extracted from coherent optical transient signals. The most basic coherent optical transients have been reviewed in this chapter. Coherent optical transient spectroscopy is still an evolving field, as new techniques are being added to established ones. It is likely that one will see increased use of both temporal and spatial masks for coherent control of atomic state coherence. Although many of the coherent optical transients are direct analogs of similar effects in NMR, others are unique to the optical domain. The velocity selectivity associated with the Doppler effect offers unique possibilities. Velocity diffusion has been studied extensively using coherent optical transients, and coherent optical transients are being rediscovered as an important probe of cold atoms and Bose condensates.^{106,107} Beyond the gas phase, developments in the area of coherent optical transient spectroscopy are moving very rapidly, as the power of this methodology is seen as a key that can help unlock the decay dynamics of complex molecules and semiconductor systems. One is also examining whether coherent optical transient methods similar to those employed in multidimensional NMR^{108,109} can be used to probe electronic and molecular structure.

11.13 REFERENCES

1. H. Metcalf and W. D. Phillips, "Time-Resolved Subnatural-Width Spectroscopy," *Opt. Lett.* **5**:540 (1980); W. D. Phillips and H. J. Metcalf, "Time Resolved Sub-Natural Width Spectroscopy," in *Precision Measurements and Fundamental Constants II*, B. N. Taylor and W. D. Phillips (eds.), U.S. National Bureau of Standards Special Publication 617 (1984).
2. R. G. Brewer, "Coherent Optical Spectroscopy," in *Frontiers in Laser Spectroscopy* (Les Houches, Session 27), R. Balian, S. Haroche, and S. Liberman (eds.) (North-Holland, Amsterdam, 1977), Vol. 1, pp. 341–396.
3. R. G. Brewer, "Coherent Optical Spectroscopy," in *Nonlinear Spectroscopy* (Proc. Int. School Phys., Enrico Fermi, Course 64), N. Bloembergen (ed.) (North-Holland, Amsterdam, 1977).
4. R. G. Brewer, "Coherent Optical Transients," *Physics Today* **30**(5):50 (1977).
5. R. L. Shoemaker, "Coherent Transient Infrared Spectroscopy," in *Laser and Coherence Spectroscopy*, J. I. Steinfeld (ed.) (Plenum Press, New York, 1978), pp. 197–371.
6. R. G. Brewer and E. L. Hahn, "Optical Memory," *Scientific American* **271**(6):50 (1984).
7. Special issue on coherent optical spectroscopy, *J. Opt. Soc. Amer.* **3**(4) (1986).
8. L. Allen and J. Eberly, *Optical Resonance and Two-Level Atoms* (Wiley, New York, 1975), Chaps. 2, 4, and 9.
9. M. D. Levenson and S. S. Kano, *Introduction to Nonlinear Laser Spectroscopy* (Academic Press, Boston, 1988), Chap. 6.

10. P. Meystre and M. Sargent III, *Elements of Quantum Optics* (Springer-Verlag, Berlin, 1991), Chap. 11.
11. T. W. Mossberg, R. Kachru, S. R. Hartmann, and A. M. Flusberg, "Echoes in Gaseous Media: A Generalized Theory of Rephasing Phenomena," *Phys. Rev. A* **20**:1976 (1979).
12. W. Zinth and W. Kaiser "Ultrafast Coherent Spectroscopy," in *Ultrashort Laser Pulses and Applications*, W. Kaiser (ed.), Vol. 60, Springer Topics in Applied Physics series (Springer-Verlag, Berlin, 1988), pp. 235–277.
13. Special issue on single-molecule spectroscopy, *Science* **283**:1167–1694 (1999).
14. R. P. Feynman, F. L. Vernon, Jr., and R. W. Hellwarth, "Geometrical Representation of the Schrödinger Equation for Solving Maser Problems," *J. Appl. Phys.* **28**:49 (1957).
15. R. G. Brewer and E. L. Hahn, "Coherent Raman Beats," *Phys. Rev. A* **8**:464 (1973).
16. P. Hu, S. Geschwind, and T. M. Jedju, "Spin-Flip Raman Echo in *n*-Type CdS," *Phys. Rev. Lett.* **37**:1357 (1976).
17. T. Mossberg, A. Flusberg, R. Kachru, and S. R. Hartmann, "Tri-Level Echoes," *Phys. Rev. Lett.* **39**:1523 (1977).
18. M. Ducloy, J. R. R. Leite, and M. S. Feld, "Laser Saturation Spectroscopy in the Time-Delayed Mode: Theory of Optical Free Induction Decay in Coupled Doppler-Broadened Systems," *Phys. Rev. A* **17**:623 (1978).
19. N. Lu and P. R. Berman, "Photon Echoes Using Double-Resonance Optical Pulses," *J. Opt. Soc. Amer. B* **2**:1883 (1985).
20. I. V. Yevseyev, V. M. Yermachenko, and V. A. Reshtov, "The Stimulated Photon Echo as a New Method for Measuring Population, Orientation, and Alignment Relaxation Times," *J. Phys. B* **19**:185 (1986).
21. See, for example, R. Boyd, *Nonlinear Optics* (Academic Press, Boston, 1992), Chap. 6, and references therein.
22. See, for example, Y. R. Shen, *The Principles of Nonlinear Optics* (Wiley, New York, 1984), Chap. 17, and references therein.
23. See, for example, J. F. Valley, G. Khitrova, H. M. Gibbs, J. W. Grantham, and X. Jiajin, "CW Conical Emission: First Comparison and Agreement between Theory and Experiment," *Phys. Rev. Lett.* **64**:2362 (1990), and references therein.
24. See, for example, B. D. Paul, M. L. Dowell, A. Gallagher, and J. Cooper, "Observation of Conical Emission from a Single Self-Trapped Beam," *Phys. Rev. A* **59**:4784 (1999), and references therein.
25. A. I. Lvovsky and S. R. Hartmann, "Superradiant Self-Diffraction," *Phys. Rev. A* **59**:4052 (1999).
26. R. G. Brewer and R. L. Shoemaker, "Optical Free Induction Decay," *Phys. Rev. A* **6**:2001 (1972).
27. R. G. Brewer and A. Z. Genack, "Optical Coherent Transients by Laser Frequency Switching," *Phys. Rev. Lett.* **36**:959 (1976).
28. H. Lehmitz and H. Harde, "Measurement of First-Order Free-Induction Decay," in *Methods of Laser Spectroscopy*, Y. Prior, A. Ben-Reuven, and M. Rosenbluh (eds.) (Plenum Press, New York, 1986), pp. 109–112.
29. P. Dubé, M. D. Levenson, and J. L. Hall, "Free-Induction Decay in Molecular Iodine Measured with an Extended-Cavity Diode Laser," *Opt. Lett.* **22**:184 (1997).
30. R. M. Macfarlane and M. Zhu, "Observation of Coherent Transients by Use of Current Switching of a Semiconductor Diode Laser," *Opt. Lett.* **22**:248 (1997).
31. A. Abragam, *The Principles of Nuclear Magnetism* (Oxford University Press, New York, 1961).
32. C. P. Slichter, *Principles of Magnetic Resonance* (Harper & Row, New York, 1963).
33. R. H. Dicke, "Coherence in Spontaneous Radiation Processes," *Phys. Rev.* **93**:99 (1954).
34. See, for example, I. P. Herman, J. C. MacGillivray, N. Skribanowitz, and M. S. Feld, "Self-Induced Emission in Optically Pumped HF Gas: The Rise and Fall of the Superradiant State," in *Laser Spectroscopy*, R. G. Brewer and A. Mooradian (eds.) (Plenum Press, New York, 1974), pp. 379–412.
35. D. Polder, M. F. H. Schuurmans, and Q. H. F. Vreken, "Superfluorescence: Quantum-Mechanical Derivation of Maxwell-Bloch Description with Fluctuating Field Source," *Phys. Rev. A* **19**:1192 (1979).
36. N. A. Kurnit, I. D. Abella, and S. R. Hartmann, "Observation of a Photon Echo," *Phys. Rev. Lett.* **13**:567 (1964).
37. R. G. Brewer and R. L. Shoemaker, "Photon Echo and Optical Nutation in Molecules," *Phys. Rev. Lett.* **27**:631 (1971).
38. L. S. Vasilenko and N. N. Rubtsova, "Photon Echo in Molecular Gases: I. Spatial, Temporal, Polarization, and Spectral Properties, II. Investigation of Collisional Relaxation," *Laser Phys.* **6**:821 (1996); **7**:903 (1997), and references therein.

39. E. L. Hahn, "Spin Echoes," *Phys. Rev.* **80**:580 (1950).
40. P. R. Berman, "Theory of Collision Effects on Line Shapes Using a Quantum-Mechanical Description of the Atomic Center-of-Mass Motion—Application to Lasers. I," *Phys. Rev. A* **2**:2435 (1970).
41. R. Beach, S. R. Hartmann, and R. Friedberg, "Billiard Ball Echo Model," *Phys. Rev. A* **25**:2658 (1982); R. Friedberg and S. R. Hartmann, "Billiard Balls and Matter-Wave Interferometry," *Phys. Rev. A* **48**:1446 (1993); "Echoes and Billiard Balls," *Laser Phys.* **3**:1128 (1993). See also T. W. Mossberg and S. R. Hartmann, "Diagrammatic Representation of Photon Echoes and Other Laser-Induced Ordering Processes in Gases," *Phys. Rev. A* **23**:1271 (1981).
42. C. J. Bordé, "Atomic Interferometry with Internal State Labeling," *Phys. Lett. A* **140**:10 (1989).
43. R. Friedberg and S. R. Hartmann, "Relaxation and Interferometry via Multiple Order Coherent Scattering in Atomic Vapors," *Laser Phys.* **5**:526 (1993).
44. P. R. Berman (ed.), *Atom Interferometry* (Academic Press, San Diego, 1997).
45. P. W. Anderson, B. I. Halperin, and C. M. Varma, "Anomalous Low Temperature Thermal Properties of Glasses and Spin Glasses," *Phil. Mag.* **25**:1 (1972).
46. W. M. Yen and P. M. Silzer (eds.), *Laser Spectroscopy of Solids*, Vol. 49, Springer Topics in Applied Physics series (Springer-Verlag, Berlin, 1986).
47. P. R. Berman, "Validity Conditions for the Optical Bloch Equations," *J. Opt. Soc. Amer. B* **3**:564 (1986), and references therein.
48. P. R. Berman, "Markovian Relaxation Processes for Atoms in Vapors and in Solids: Calculation of Free-Induction Decay in the Weak External-Field Limit," *J. Opt. Soc. Amer. B* **3**:572 (1986), and references therein.
49. P. R. Berman, in *New Trends in Atomic Physics* (Les Houches, Session 38), G. Grynberg and R. Stora (eds.) (North-Holland, Amsterdam, 1984), Vol. 1, pp. 451–514.
50. R. G. Devoe and R. G. Brewer, "Experimental Test of the Optical Bloch Equations for Solids," *Phys. Rev. Lett.* **50**:1269 (1983).
51. P. R. Berman, J. M. Levy, and R. G. Brewer, "Coherent Optical Transient Study of Molecular Collisions: Theory and Observations," *Phys. Rev. A* **11**:1668 (1975).
52. T. W. Mossberg, R. Kachru, and S. R. Hartmann, "Observation of Collisional Velocity Changes Associated with Atoms in a Superposition of Dissimilar Electronic States," *Phys. Rev. Lett.* **44**:73 (1980).
53. R. Kachru, T. J. Chen, S. R. Hartmann, T. W. Mossberg, and P. R. Berman, "Measurement of a Total Atomic-Radiator-Perturber Scattering Cross Section," *Phys. Rev. Lett.* **47**:902 (1981).
54. R. A. Forber, L. Spinelli, J. E. Thomas, and M. S. Feld, "Observation of Quantum Diffractive Velocity-Changing Collisions by Use of Two-Level Heavy Optical Radiators," *Phys. Rev. Lett.* **50**:331 (1982).
55. J. C. Keller and J. L. LeGouët, "Stimulated Photon Echo for Collisional Study in Yb Vapor," *Phys. Rev. Lett.* **52**:2034 (1984).
56. A. G. Yodh, J. Golub, and T. W. Mossberg, "Collisional Relaxation of Excited State Zeeman Coherences in Atomic Ytterbium Vapor," *Phys. Rev. A* **32**:844 (1985).
57. R. Kachru, T. W. Mossberg, and S. R. Hartmann, "Stimulated Photon Echo Study of Na($3^2S_{1/2}$)-CO Velocity-Changing Collisions," *Opt. Comm.* **30**:57 (1979).
58. T. Mossberg, A. Flusberg, R. Kachru, and S. R. Hartmann, "Total Scattering Cross Section for Na on He Measured by Stimulated Photon Echoes," *Phys. Rev. Lett.* **42**:1665 (1979).
59. T. W. Mossberg, "Time-Domain Frequency-Selective Optical Storage Data," *Opt. Lett.* **7**:77 (1982).
60. N. W. Carlson, W. R. Babbitt, and T. W. Mossberg, "Storage and Phase Conjugation of Light Pulses Using Stimulated Photon Echoes," *Opt. Lett.* **8**:623 (1983).
61. M. K. Kim and R. Kachru, "Long Term Image Storage and Phase Conjugation by a Backward-Stimulated Echo in $\text{Pr}^{3+}\text{LaF}_3$," *J. Opt. Soc. Amer. B* **4**:305 (1987).
62. X. A. Shen and R. Kachru, "High Speed Recognition by Using Stimulated Echoes," *Opt. Lett.* **17**:520 (1992).
63. L. S. Vasilenko, N. M. Dyuba, and M. N. Skvortsov, "Coherent Emission in Time-Separated Fields," *Sov. J. Quant. Electron.* **8**:980 (1978).
64. Y. V. Baklonov, B. Y. Dubetsky, and V. P. Chebotaev, "Non-linear Ramsey Resonance in the Optical Region," *Appl. Phys.* **9**:171 (1976).

65. E. V. Baklanov, B. Y. Dubetsky and V. M. Semibalamut, "Theory of Stimulated Coherent Emission from Atoms in Spatially Separated Optical Fields," *Sov. Phys. JETP* **49**:244 (1979).
66. T. W. Mossberg, R. Kachru, E. Whittaker, and S. R. Hartmann, "Temporally Recurrent Spatial Ordering of Atomic Population in Gases: Grating Echoes," *Phys. Rev. Lett.* **43**:851 (1979).
67. J. L. LeGouët and P. R. Berman, "Photon Echoes in Standing Wave Fields: Time Separation of Spatial Harmonics," *Phys. Rev. A* **20**:1105 (1979).
68. R. Kachru, T. W. Mossberg, E. Whittaker, and S. R. Hartmann, "Optical Echoes Generated by Standing Wave Fields: Observations in Atomic Vapors," *Opt. Comm.* **31**:223 (1979).
69. M. V. Belyayev, V. P. Chebotaev, M. N. Skvortsov, and L. S. Vasilenko, "Resonant Coherent Transients in a Gas in the Standing Wave Field," *Appl. Phys. B* **26**:67 (1981).
70. L. S. Vasilenko, I. D. Matveyenko, and N. N. Rubtsova, "Study of Narrow Resonances of Coherent Radiation in Time Separated Fields in SF₆," *Opt. Comm.* **53**:371 (1985).
71. B. Dubetsky, P. R. Berman, and T. Sleator, "Grating Stimulated Echo," *Phys. Rev. A* **46**:2213 (1992).
72. P. R. Berman, "Collisional Decay and Revival of the Grating Stimulated Echo," *Phys. Rev. A* **49**:2922 (1994).
73. A. M. Weiner, S. De Silvestri, and E. P. Ippen, "Three-Pulse Scattering for Femtosecond Dephasing Studies: Theory and Experiments," *J. Opt. Soc. Amer. B* **2**:624 (1985).
74. M. Fujita, H. Nakatsuka, H. Nakanishi, and M. Matsuoka, "Backward Echo in Two-Level Systems," *Phys. Rev. Lett.* **42**:974 (1979).
75. A. I. Alekseev and V. N. Beloborodov, "Forward and Backward Photon Echoes in Gases and Solids," *Opt. Spectrosc.* **57**:277 (1984).
76. R. A. Fisher (ed.), *Optical Phase Conjugation* (Academic Press, New York, 1983).
77. M. Gower and D. Proch (eds.), *Optical Phase Conjugation* (Springer-Verlag, Berlin, 1994).
78. J. C. Bergquist, S. A. Lee, and J. L. Hall, "Saturated Absorption with Spatially Separated Laser Fields: Observation of Optical 'Ramsey' Fringes," *Phys. Rev. Lett.* **38**:159 (1977), and references therein.
79. R. L. Barger, "Influence of Second-Order Doppler Effect on Optical Ramsey Fringe Profiles," *Opt. Lett.* **6**:145 (1981).
80. C. J. Bordé, C. Salomon, S. Avrillier, A. Van Lerberghe, C. Bréant, D. Bassi, and G. Scoles, "Optical Ramsey Fringes with Traveling Waves," *Phys. Rev. A* **30**:1836 (1984).
81. N. Ramsey, "A Molecular Beam Resonance Method with Separated Oscillating Fields," *Phys. Rev.* **78**:695 (1950).
82. N. Ito, J. Ishikawa, and A. Morinaga, "Evaluation of the Optical Phase Shift in a Ca Ramsey Fringe Stabilized Optical Frequency Standard by Means of Laser Beam Reversal," *Opt. Comm.* **109**:414 (1994).
83. B. Y. Dubetsky, A. P. Kazantsev, V. P. Chebotaev, and V. P. Yakolev, "Interference of Atoms in Separated Optical Fields," *Sov. Phys. JETP* **62**:685 (1985).
84. S. B. Cahn, A. Kumarakrishnan, U. Shim, T. Sleator, P. R. Berman, and B. Dubetsky, "Time-Domain de Broglie Wave Interferometry," *Phys. Rev. Lett.* **79**:784 (1997).
85. E. M. Rasel, M. K. Oberthaler, H. Batelaan, J. Schmiedmayer, and A. Zeilinger, "Atom Wave Interferometry with Diffraction Gratings of Light," *Phys. Rev. Lett.* **75**:2633 (1995).
86. B. Dubetsky and P. R. Berman, "Matter-Wave Interference Using Two-Level Atoms and Resonant Optical Fields," *Phys. Rev. A* **59**:2269 (1999).
87. J. Guo, P. R. Berman, D. Dubetsky, and G. Grynberg, "Recoil-Induced Resonances in Nonlinear Spectroscopy," *Phys. Rev. A* **46**:1426 (1992).
88. R. Beach and S. R. Hartmann, "Incoherent Photon Echoes," *Phys. Rev. Lett.* **53**:663 (1984).
89. M. Mitsunaga, "CW Photon Echo: Theory and Observations," *Phys. Rev. A* **42**:1617 (1990).
90. See, for example, B. Do, J. Cha, D. S. Elliott, and S. J. Smith, "Phase Conjugate Four Wave Mixing with Partially-Coherent Laser Fields," *Phys. Rev. A* **60**:508 (1999). This paper contains an extensive bibliography citing earlier work in this field.
91. T. W. Mossberg, "Swept-Carrier Time-Domain Optical Memory," *Opt. Lett.* **17**:535 (1992).
92. H. Lin, T. Wang, G. A. Wilson, and T. W. Mossberg, "Experimental Demonstration of Swept-Carrier Time-Domain Optical Memory," *Opt. Lett.* **20**:91 (1995).

93. K. D. Merkel and W. R. Babbitt, "Chirped-Pulse Programming of Optical Coherent Transient True-Time Delay," *Opt. Lett.* **23**:528 (1998).
94. L. Schultheis, M. D. Sturge, and J. Hegarty, "Photon Echoes from Two-Dimensional Excitons in GaAs-AlGaAs Quantum Wells," *Appl. Phys. Lett.* **47**:995 (1985).
95. M. D. Webb, S. T. Cundiff, and D. G. Steel, "Observation of Time-Resolved Picosecond Stimulated Photon Echoes and Free Polarization Decay in GaAs/AlGaAs Multiple Quantum Wells," *Phys. Rev. Lett.* **66**:934 (1991).
96. H. J. Eichler, P. Günter, and D. W. Pohl, *Laser-Induced Dynamic Gratings* (Springer-Verlag, Berlin, 1986).
97. J. Y. Bigot, M. A. Mycek, S. Weiss, R. G. Ulbrich, and D. S. Chemla, "Instantaneous Frequency Dynamics of Coherent Wave Mixing in Semiconductor Quantum-Wells," *Phys. Rev. Lett.* **70**:3307 (1993).
98. A. L. Smirl, "Coherent Exciton Dynamics: Time-Resolved Polarimetry," in *Semiconductor Quantum Optoelectronics: From Quantum Physics to Smart Devices*, A. Miller and D. M. Finlayson (eds.) (Institute of Physics, London, 1999).
99. C. J. Bardeen, W. Wang, and C. V. Shank, "Femtosecond Chirped Pulse Excitation of Vibrational Wave Packets in LD690 and Bacteriorhodopsin," *J. Phys. Chem. A* **102**:2759 (1998).
100. W. P. de Boeij, M. S. Pshenichnikov, and D. A. Wiersma, "Ultrafast Solvation Dynamics Explored by Femtosecond Photon Echo Spectroscopies," *Ann. Rev. Phys. Chem.* **49**:99 (1998).
101. J. Shah, *Ultrafast Processes in Semiconductors and Semiconductor Nanostructures*, Vol. 115 in Springer Solid-State Sciences series (Springer, Berlin, 1996).
102. See, for example, H. Wang, K. B. Ferrio, D. G. Steel, P. R. Berman, Y. Z. Hu, R. Binder, and S. W. Koch, "Transient Four-Wave Mixing Line Shapes: Effects of Excitation Induced Dephasing," *Phys. Rev. A* **49**:1551 (1994).
103. H. Haug and S. W. Koch, *Quantum Theory of the Optical and Electronic Properties of Semiconductors* (World Scientific, Singapore, 1993).
104. C. Sieh, T. Meier, F. Jahnke, A. Knorr, S. W. Koch, P. Brick, M. Hübner, C. Ell, J. Prineas, G. Khitrova, and H. M. Gibbs, "Coulomb Memory Signatures in the Excitonic Optical Stark Effect," *Phys. Rev. Lett.* **82**:3112 (1999).
105. D. S. Chemla, "Ultrafast Transient Nonlinear Optical Processes in Semiconductors," *Semiconductors and Semimetals* **58**:175 (1999).
106. J. Stenger, S. Inouye, A. P. Chikkatur, D. M. Stamper-Kurn, D. E. Pritchard, and W. Ketterle, "Bragg Spectroscopy of a Bose-Einstein Condensate," *Phys. Rev. Lett.* **82**:4569 (1999).
107. Y. B. Ovchinnikov, J. H. Müller, M. R. Doery, E. I. D. Vredenburg, K. Helmerson, S. L. Rolston, and W. D. Phillips, "Diffraction of a Released Bose-Einstein Condensate by a Pulsed Standing Light Wave," *Phys. Rev. Lett.* **83**:284 (1999).
108. S. Mukamel, A. Piryatinski, and V. Chernyak, "Two-Dimensional Raman Echoes: Femtosecond View of Molecular Structure and Vibrational Coherence" *Acts. Chem. Res.* **32**:145 (1999).
109. P. Hamm, M. Lim, W. F. DeGrado, and R. M. Hochstrasser, "The Two-Dimensional IR Nonlinear Spectroscopy of a Cyclic Penta-peptide in Relation to its Three-Dimensional Structure," *Proc. Natl. Acad. Sci. USA* **96**:2036 (1999).

PHOTOREFRACTIVE MATERIALS AND DEVICES

Mark Cronin-Golomb

*Department of Biomedical Engineering
Tufts University
Medford, Massachusetts*

Marvin Klein

*Intelligent Optical Systems, Inc.
Torrance, California*

12.1 INTRODUCTION

The photorefractive effect is a real time holographic optical nonlinearity that is effective for low-power lasers over a wide range of wavelengths. It is relatively easy to use even in modestly equipped laboratories: all that is needed to get started is a sample of photorefractive material, almost any laser operating in the visible or near-infrared, and a few simple optics such as lenses and beam splitters. The diffraction efficiency of photorefractive holograms is roughly independent of the intensity of the writing beams and in many materials, the diffraction efficiency of these holograms can approach 100 percent, so that sophisticated detectors are not required. Its simplicity of use has been largely responsible for its widespread popularity. As is often the case, however, attractive features such as these come only with associated disadvantages. For the photorefractive effect, the main disadvantage is one of speed. The nonlinearities come to steady state at a rate which is approximately inversely proportional to intensity. The fastest high-diffraction efficiency materials have response times of the order of 1 ms at 1 W/cm². Even so, in certain applications, the characteristic slowness is not a disadvantage. In the first part of this chapter, we explain the basic mechanisms of the photorefractive effect. The second part deals with material selection considerations and the third part describes some typical applications. For the reader in need of an extensive overview of the photorefractive effect and its applications, we recommend a two-volume set edited by Gunter and Huignard and the three-volume updated version two.¹

Grating Formation

The photorefractive effect is observed in materials which

1. Exhibit a linear electro-optic effect
2. Are photoconductive
3. Have a low dark conductivity

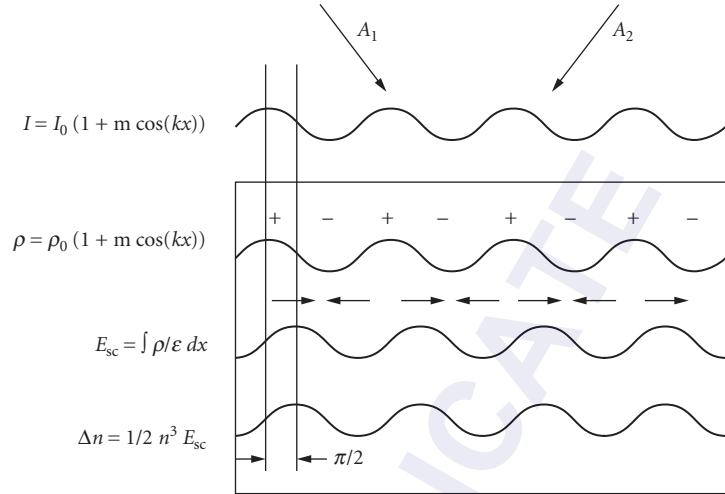


FIGURE 1 The photorefractive mechanism. Two laser beams intersect, forming an interference pattern $I(x)$. Charge is excited where the intensity is large and migrates to regions of low intensity. The electric field E_{sc} associated with the resultant space charge ρ_{sc} operates through the linear electro-optic coefficients to produce a refractive index grating Δn .

Two laser beams record a photorefractive hologram when their interference pattern is incident on a photorefractive crystal (Fig. 1). Charge carriers are preferentially excited in the bright fringes and are then free to drift and diffuse until they recombine with traps, most likely in the darker regions of the interference pattern. In this way, a space charge builds up inside the crystal in phase with the interference pattern. The electric field of this space charge acts through the linear electro-optic effect to form a volume holographic refractive index grating in real time. Real time means that the development occurs with a time constant of the order of the response time of the photorefractive crystal. The writing beams then diffract from the hologram into each other. Even though the process depends on the linear electro-optic effect, which is second-order, the whole process acts effectively as a third-order nonlinearity as far as the writing beams are concerned, and third-order coupled-wave equations can be written for their amplitudes. It is important to notice that the electric field is spatially shifted by 90° with respect to the interference pattern because of the Gauss's law integration that links the space charge to the electric field. The refractive index gradient is also shifted by $\pm 90^\circ$. This shift is possible because the refractive index perturbation depends on the direction of the electric field, not just on its magnitude. The direction of the grating shift is determined by the sign of the electro-optic coefficient and crystal orientation. This dependence on crystal orientation is due to the lack of inversion symmetry associated with the linear electro-optic coefficient: if the crystal is inverted through its origin, the sign of the phase shift changes. In uniaxial crystals, inversion corresponds simply to reversing the direction of the optic axis. These symmetry effects are intimately related to the origin of photorefractive beam amplification.

The reason that use of the linear electro-optic effect is important is that a Bragg-matched volume hologram should have the same period as the optical interference pattern that wrote it. The refractive index change should be directly proportional to the space charge electric field. This is only possible if the material displays the linear electro-optic effect leading to a refractive index distribution $\Delta n \propto r_{\text{eff}} E_{sc}$, where r_{eff} is an effective electro-optic coefficient and E_{sc} is the space charge field. The lack of inversion symmetry needed by the linear electro-optic effect may be found in ferroelectric

materials such as barium titanate, optically active materials such as bismuth silicon oxide and cubic compound semiconductors such as GaAs and InP.

When the charge transport is purely diffusive, the magnitude of the spatial phase shift is 90° . However, in certain circumstances the phase shift can depart from 90° . This occurs if electric fields are applied to the crystal or if the crystal exhibits the photovoltaic effect² so that drift mechanisms come into play, or, if the writing beams have different frequencies, so that the interference pattern moves in the crystal with the index grating lagging behind it.

The second requirement implies that the material should contain photoexcitable impurities. Direct band-to-band photoconductivity is usually not useful since it limits the optical interaction distances to rather small absorption depths.

The requirement for low dark conductivity ensures that the space charge can support itself against decay by leakage through background conduction.

The Standard Rate Equation Model

The simplest model, as formulated by Vinetskii and Kukhtarev,³ involves optical excitation of charge carriers. For the purposes of this introduction, we will assume that the carriers are electrons which can be excited from a donor species such as Fe^{2+} and which can recombine into an acceptor such as Fe^{3+} . Let n be the local number density of mobile excited electrons and $N_D = N + N^+$ be the number density of the impurities or defects responsible for the photorefractive effect, where N^+ is the number density of acceptor dopants (e.g., Fe^{3+}) and N is the number density of donor dopants (e.g., Fe^{2+}). Let N_A be the number density of negative ions that compensate the excess positive charge of acceptor dopants when the charge is uniformly distributed in the dark. Neglecting the photovoltaic effect, we may write the following rate equations of generation and recombination, continuity, electric field (Poisson equation), and total drift and diffusion current:

$$\begin{aligned}\frac{\partial N^+}{\partial t} &= (sI + \beta)N - \gamma n N^+ \\ \frac{\partial N^+}{\partial t} &= \frac{\partial n}{\partial t} - \frac{1}{e} \nabla \cdot \mathbf{j} \\ \nabla \cdot \mathbf{E} &= (N^+ - N_A - n)e/\epsilon \\ \mathbf{j} &= \mu_e n \mathbf{E} + k_B T \mu \nabla n\end{aligned}\quad (1)$$

where s is proportional to the photoionization cross section ($s = \sigma/h\nu$), β is the dark generation rate, γ is a recombination coefficient, \mathbf{j} is the electric current density. At steady state, the space charge is determined by a balance between charge diffusion away from bright fringes and electrostatic repulsion from charge concentrations. Extensions of these equations to include the effects of electron-hole competition,⁴ multiple dopants,^{5,6} photovoltaic effects,^{7,8} and short pulse excitation⁹ have been developed over the past few years. Nevertheless, the most important features of the photorefractive effect may be well-modeled by the simple equations shown here. Assuming that the number density of charge carriers is much less than the optically induced donor density perturbation ($N^+ - N_A$), a situation that almost always holds, the rate equations can be linearized to give the following solution for the fundamental spatial Fourier component of the space charge field E_{sc} induced by a sinusoidal optical fringe pattern of wave number k_g and fringe visibility m when a dc field E_0 is applied to the crystal:

$$E_{sc} = \frac{m/2}{1 + \beta/sI_0} \frac{E_q(iE_0 - E_d)}{E_0 + i(E_d + E_q)} \quad (2)$$

where I_0 is the total average intensity of the interacting beams, E_q and E_d are characteristic fields of maximum space charge and diffusion, respectively $E_d = k_B T k_g / e$, $E_q = e N_A / \epsilon k_g$. The response time τ is given by

$$\tau = \frac{N_A}{SN_D(I_0 + \beta/s)} \frac{E_0 + i(E_d + E_\mu)}{E_0 + i(E_d + E_q)} \quad (3)$$

where E_μ is the characteristic mobility field $E_\mu = \gamma N_A / \mu k_g$. These results show that the steady-state space charge field is approximately independent of total intensity and that the response time is inversely proportional to total intensity if the intensity exceeds the saturation intensity β/s .

In the absence of a dc applied field, the fundamental Fourier component of the space charge field is purely imaginary, indicating a 90° spatial phase shift between the interference pattern and the index grating. The effect of an applied dc field is to increase the magnitude of the space charge field and to move the spatial phase shift from 90° . The 90° phase shift is optimal for two-beam coupling amplifiers, as will be indicated, and can be restored by inducing a compensating phase difference through detuning the writing beams from each other to cause the grating to lag behind the interference pattern.¹⁰ Alternatively, an ac applied field may be used to enhance the magnitude of the photorefractive grating and maintain the 90° phase shift.^{11,12} For externally pumped four-wave mixing, the 90° phase shift is not optimum, so the applied field-induced deviation of the phase shift is advantageous.¹³

Wave Interactions

Two-Beam Coupling Consider two laser beams writing a grating in a photorefractive medium as depicted in Fig. 2. The effect of the photorefractive nonlinearity on the writing beams can be described quite accurately by conventional coupled-wave theory. Coupled-wave equations for two-beam coupling can be found by taking the slowly varying envelope approximation and substituting the optical electric fields into the scalar wave equation.¹⁴

$$\nabla^2 E + k^2 E = 0 \quad (4)$$

where

$$k = \frac{\omega n}{c} = \frac{\omega(n_0 + \Delta n)}{c} \quad (5)$$

with Δn being the optically induced refractive index change.

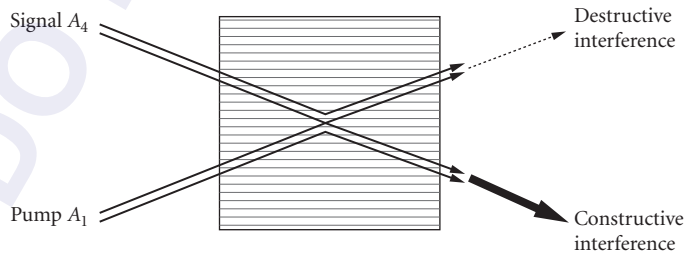


FIGURE 2 Two-beam coupling amplification. Beams 1 and 4 write a diffraction grating. Beam 1 diffracts from the grating to constructively interfere with and amplify beam 4. Beam 4 diffracts from the grating to destructively interfere with and attenuate beam 1.

The resultant coupled-wave equations read:

$$\begin{aligned}\frac{dA_1}{dz} &= -\gamma \frac{A_1 A_4^*}{I_0} A_4 \\ \frac{dA_4^*}{dz} &= \gamma \frac{A_1 A_4^*}{I_0} A_1^*\end{aligned}\quad (6)$$

where I_0 is sum of the intensities I_1 and I_4 of the interacting beams. The coupling constant γ is related to the space charge field by

$$\gamma = \frac{i\omega}{c \cos\theta} \frac{r_{\text{eff}} n^3 E_{\text{sc}}/m}{2} \quad (7)$$

It is proportional to the effective electro-optic coefficient r_{eff} ¹⁵ and is real when the index grating is 90° out of phase with the interference pattern (as in the case without an applied dc electric field, or with high-frequency ac applied fields). In that case, the coupled-wave equations show that beam 4 is amplified at the expense of beam 1 if γ is positive. The amplification effect can be explained in physical terms: the light diffracted by the grating from beam 1 interferes constructively with beam 4, so beam 4 is amplified. The light diffracted by the grating from beam 4 interferes destructively with beam 1, which consequently loses energy. At the same time, the phases of the interacting beams are preserved. If the spatial phase shift departs from 90° then the energy transfer effect becomes less, and phase coupling begins to appear.

The photorefractive beam coupling gain can be substantial in materials with large electro-optic coefficients. The intensity gain coefficient Γ which characterizes the transfer of energy between two beams [$\Gamma = \gamma + \gamma^*$; see Eq. (8)] can exceed 60 cm⁻¹ in BaTiO₃¹⁶ and in SBN.¹⁷ Such high gain makes possible the construction of devices such as photorefractive parametric oscillators and recursive image processors. In the high gain case, signal-to-noise issues become important. Generally, defects in photorefractive crystals act as scattering centers. In a process similar to amplified spontaneous emission, this scattered light can be very strongly amplified into a broad fan of light. This fanning effect¹⁸⁻²⁰ is a significant source of noise for photorefractive image amplifiers, and considerable effort has been devoted to lessening the effect by growing more uniform crystals and making device design modifications.^{21, 22} On the other hand, the fanning effect can be a useful source of seed beams for various oscillators^{23, 24} or as the basis for optical limiters.²⁵

In the undepleted pump approximation ($I_1 \gg I_4$), theoretical analysis is extremely simple: the gain is exponential, with amplitude gain coefficient simply γ . But even in the pump depletion case, analysis is quite straightforward because, as in the case of second-harmonic generation, the nonlinear coupled-wave equations can be solved exactly.¹⁴

$$\begin{aligned}I_1(z) &= \frac{I_0}{1 + (I_4(0)/I_1(0)) \exp(\Gamma z)} \\ I_4(z) &= \frac{I_0}{1 + (I_1(0)/I_4(0)) \exp(-\Gamma z)} \\ \psi_1(z) &= \psi_1(0) - \Gamma' z + \frac{\Gamma'}{\Gamma} \ln(I_4(z)/I_4(0)) \\ \psi_4(z) &= \psi_4(0) - \Gamma' z - \frac{\Gamma'}{\Gamma} \ln(I_1(z)/I_1(0))\end{aligned}\quad (8)$$

where $\Gamma = \gamma + \gamma^*$ and $\Gamma' = (\gamma - \gamma^*)/2i$. The physical implications of these equations are clear: for $I_4(0) \exp(\Gamma z) \ll I_1(0)$, beam 4 is exponentially amplified; for $I_4(0) \exp(\Gamma z) \gg I_1(0)$, beam 4

receives all of the intensity of both input beams, and beam 1 is completely depleted. There is phase transfer only if γ has an imaginary component. Linear absorption is accounted for simply by multiplying each of the intensity equations by $\exp(-\alpha z)$, where α is the intensity absorption coefficient.

The plane-wave transfer function of a photorefractive two-beam coupling amplifier may be used to determine thresholds (given in terms of $\gamma\ell$, where ℓ is the interaction length), oscillation intensities, and frequency pulling effects in unidirectional ring oscillators based on two-beam coupling.^{26–29}

Four-Wave Mixing Four-wave mixing and optical phase conjugation may also be modeled by plane-wave coupled-wave theories for four interacting beams¹³ (Fig. 3). In general, when all four beams are mutually coherent, they couple through four sets of gratings: (1) the transmission grating driven by the interference term $(A_1A_4^* + A_2^*A_3)$, (2) the reflection grating driven by $(A_2A_4^* + A_1^*A_3)$, (3) the counterpropagating pump grating $(A_1A_2^*)$, and (4) the signal/phase conjugate grating $(A_3A_4^*)$. The theories can be considerably simplified by modeling cases in which the transmission grating only or the reflection grating only is important. The transmission grating case may be experimentally realized by making beams 1 and 4 mutually coherent but incoherent with beam 2 (and hence beam 3, which is derived directly from beam 2). The four-wave mixing coupled-wave equations can be solved analytically in several useful cases by taking advantage of conservation relations inherent in the four-wave mixing process,^{30,31} and by using group theoretic arguments.³² Such

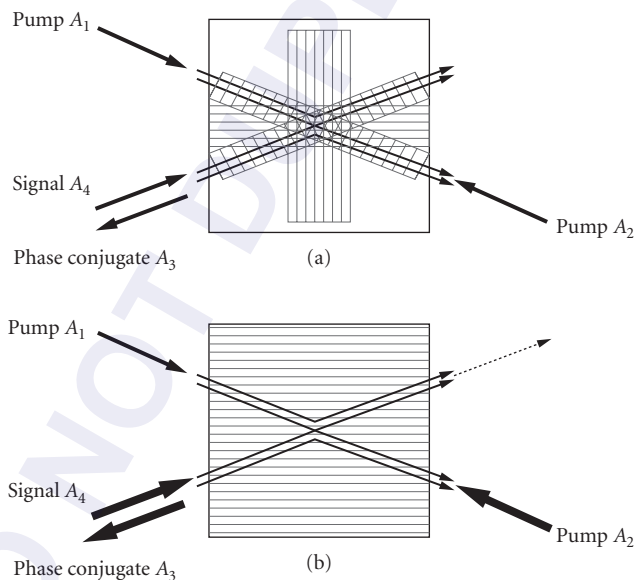


FIGURE 3 Four-wave mixing phase conjugation. (a) Beams 1 and 2 are pump beams, beam 4 is the signal, and beam 3 is the phase conjugate. All four interaction gratings are shown: reflection, transmission, counterpropagating pump, and signal/phase conjugate. Beam 1 interferes with beam 4 and beam 2 interferes with beam 3 to write the transmission grating. Beam 2 interferes with beam 4 and beam 1 interferes with beam 3 to write the reflection grating. Beam 1 interferes with beam 2 to write a counterpropagating beam grating, as do beams 3 and 4. (b) If beams 1 and 4 are mutually incoherent, but incoherent with beam 2, only the transmission grating will be written. This interaction is the basis for many self-pumped phase conjugate mirrors.

analytic solutions are of considerable assistance in the design and understanding of various four-wave mixing devices. But as in other types of four-wave mixing, the coupled-wave equations can be linearized by assuming the undepleted pump approximation. Considerable insight can be obtained from these linearized solutions.¹³ For example, they predict phase conjugation with gain and self-oscillation. The minimum threshold for phase conjugation with gain in the transmission grating case is $\gamma\ell = 2\ln(1 + \sqrt{2}) \approx 1.76$, whereas the minimum threshold for the usual $\chi^{(3)}$ nonlinearities is $\gamma\ell = i\pi/4 \approx 0.79i$.

Self-oscillation (when the phase conjugate reflectivity tends to infinity) can only be achieved if the coupling constant is complex. In the $\chi^{(3)}$ case this is the normal state of affairs, but in the photorefractive case, as mentioned, it requires some additional effect such as provided by applied electric fields, photovoltaic effect, or frequency detuning. Photorefractive four-wave mixing thresholds are usually higher than the corresponding $\chi^{(3)}$ thresholds because photorefractive symmetry implies that either the signal or the phase conjugate tend to be deamplified in the interaction. In the $\chi^{(3)}$ case, both signals and phase conjugate can be amplified. This effect results in the fact that while the optimum pump intensity ratio is unity in the $\chi^{(3)}$ case, the optimum beam intensities are asymmetric in the photorefractive case. A consequence of the interplay between self-oscillation and coupling constant phase is that high-gain photorefractive phase conjugate mirrors tend to be unstable. Even if the crystal used is purely diffusive, running gratings can be induced which cause the coupling constant to become complex, giving rise to self-oscillation instabilities.^{33,34}

Anisotropic Scattering Because of the tensor nature of the electro-optic effect, it is possible to observe interactions with changing beam polarization. One of the most commonly seen examples is the anisotropic diffraction ring of ordinary polarization that appears on the opposite side of the amplified beam fan when a single incident beam of extraordinary polarization propagates approximately perpendicular to the crystal optic axis.^{35–37} Since the refractive indices for ordinary and extraordinary waves differ from each other, phase matching for such an interaction can only be satisfied along specific directions, leading to the appearance of phase-matched rings. There are several other types of anisotropic scattering, such as broad fans of scattered light due to the circular photovoltaic effect,³⁸ and rings that appear when ordinary and extraordinary polarized beams intersect in a crystal.³⁹

Oscillators with Photorefractive Gain and Self-Pumped Phase Conjugate Mirrors

Photorefractive beam amplification makes possible several four-wave mixing oscillators, including the unidirectional ring resonator and self-pumped phase conjugate mirrors (SPPCMs). The simplest of these is the linear self-pumped phase conjugate mirror.⁴⁰ Two-beam coupling photorefractive gains supports oscillation in a linear cavity (Fig. 4a). The counterpropagating oscillation beams pump the crystal as a self-pumped phase conjugate mirror for the incident beam. The phase conjugate reflectivity of such a device can theoretically approach 100 percent, with commonly available crystals. In practice, the reflectivity is limited by parasitic fanning loss. Other types of self-pumped phase conjugate mirrors include the following. The semilinear mirror, consisting of a linear mirror with one of its cavity mirrors removed⁴⁰ (Fig. 4b). The ring mirror (transmission grating⁴¹ and reflection grating⁴² types). The transmission grating type is shown in Fig. 4c. The double phase conjugate mirror⁴³ (Fig. 4d). This device is also sometimes known as a mutually pumped phase conjugator (MPPC). Referring to Fig. 4b, it can be seen that the double phase conjugate mirror is part of the semilinear mirror. Several variants involving combinations of the transmission grating ring mirror and double phase conjugate mirror: the cat mirror²³ (Fig. 4e) so named after its first subject, frogs legs⁴⁴ (Fig. 4f), bird-wing⁴⁵ (Fig. 4g), bridge⁴⁶ (Fig. 4h), and mutually incoherent beam coupler⁴⁷ (Fig. 4i). The properties of these devices are sometimes influenced by the additional simultaneous presence of reflection gratings and gratings written between the various pairs of counterpropagating beams. The double phase conjugate mirror can be physically understood as being supported by a special sort of photorefractive self-oscillation in which beams 2 and 4 of Fig. 3 are

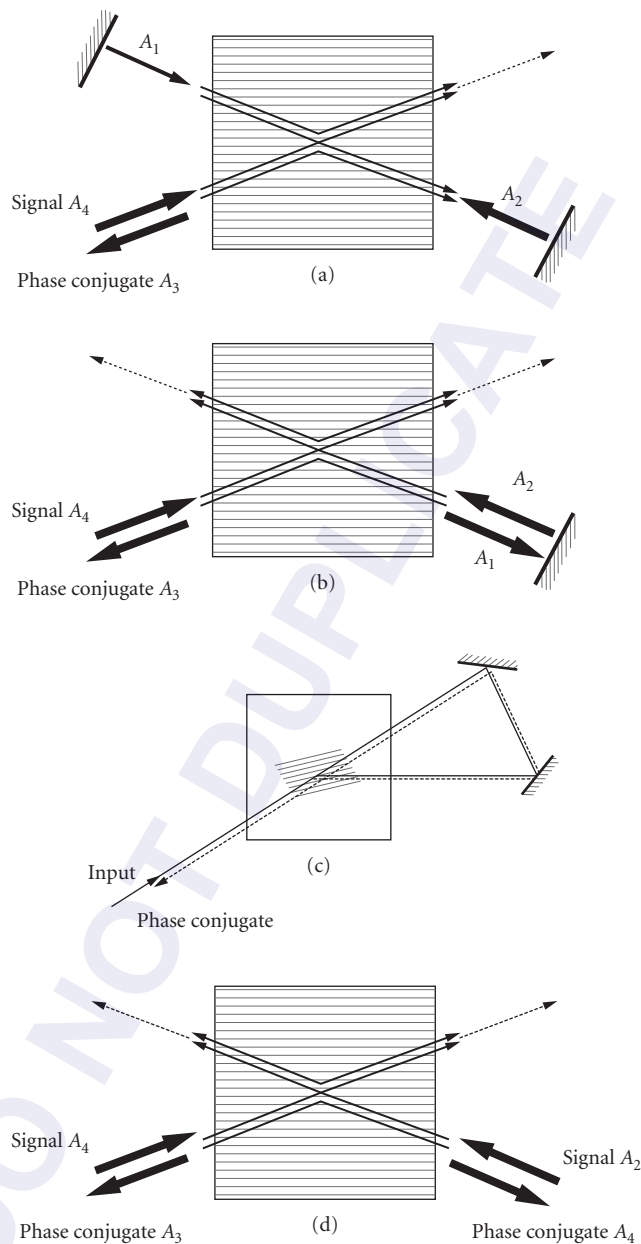


FIGURE 4 Self-pumped phase conjugate mirrors: those with external feedback: (a) linear; (b) semilinear; (c) ring; those self-contained in a single crystal with feedback (when needed) provided by total internal reflection; (d) double phase conjugate mirror; (e) cat; (f) frogs' legs; (g) bird-wing; (h) bridge; and (i) mutually incoherent beam coupler.

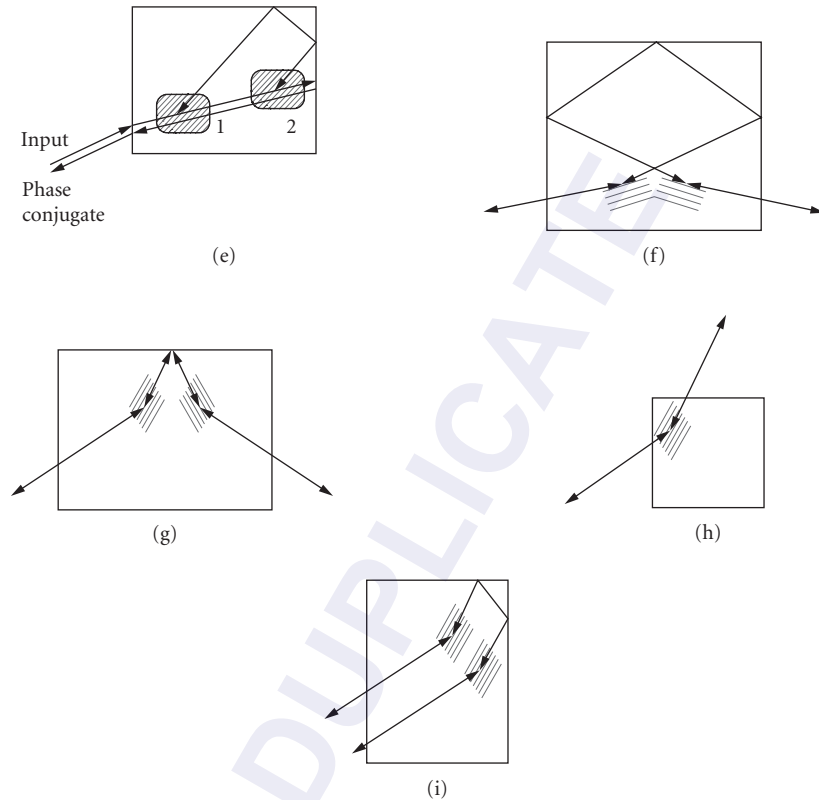


FIGURE 4 (Continued)

taken as strong and depleted. The self-oscillation threshold for the appearance of beams 1 and 3 is $\gamma\ell = 2.0$ with purely diffusive coupling.⁴⁸ It can be shown that a transmission grating ring mirror with coupling constant $\gamma\ell$ is equivalent in the plane-wave theory to a double phase conjugate mirror with coupling constant $2\gamma\ell$.⁴⁹

The matter of whether the various devices represent the results of absolute instabilities or convective instabilities has been the subject of some debate.^{50, 51}

Stimulated Photorefractive Scattering

It is natural to ask whether there is a photorefractive analogue to stimulated Brillouin scattering (SBS, a convective instability). In SBS, an intense laser beam stimulates a sound wave whose phase fronts are the same as those of the incident radiation. The Stokes wave reflected from the sound grating has the same phase fronts as the incident beam, and is traveling in the backward phase conjugate sense. The backward wave experiences gain because the sonic grating is 90° spatially out of phase with the incident beam. In the photorefractive case, the required 90° phase shift is provided automatically, so that stimulated photorefractive scattering (SPS) can be observed without any Stokes frequency shift.^{52, 53} However, the fidelity of SPS tends to be worse than that of SBS because the intensity gain discrimination mechanism is not as strong, as can be seen by examining the coupled-wave equations in each case.⁵⁴

Time-Dependent Effects

Time-dependent coupled-wave theory is important for studying the temporal response and temporal stability of photorefractive devices. Such a theory can be developed by including a differential equation for the temporal evolution of the grating. The spatiotemporal two-beam coupled-wave equations can be written, for example, as

$$\begin{aligned}\frac{\partial A_1}{\partial z} &= -GA_4 \\ \frac{\partial A_4^*}{\partial z} &= GA_1^* \\ \frac{\partial G}{\partial t} + G/\tau &= \frac{\gamma}{\tau} \frac{A_1 A_4^*}{I_0}\end{aligned}\quad (9)$$

where τ is an intensity-dependent possibly complex response time determined from photorefractive charge transport models. Models like this can be used to show that the response time for beam amplification and deamplification is increased by a factor of $\gamma\ell$ over the basic photorefractive response time τ .⁵⁵ The potential for photorefractive bistability can be studied by examining the stability of the multiple solutions of the steady-state four-wave mixing equations. In general, there is no reversible plane-wave bistability except when other nonlinearities are included in photorefractive oscillator cavities.^{56,57} (See the section on “Thresholding.”) Temporal instabilities generate interesting effects such as deterministic chaos, found both experimentally⁵⁸ and theoretically^{59,60} in high-gain photorefractive devices.

Influence of the Nonlinearity on Beam Spatial Profiles

In modeling the changes in the transverse cross section of beams as they interact, it is necessary to go beyond simple one-dimensional plane-wave theory. Such extensions are useful for analyzing fidelity of phase conjugation and image amplification, and for treating transverse mode structure in photorefractive oscillators. Several different methods have been used to approach the transverse profile problem.⁶¹ In the quasi-plane-wave method, one assumes that each beam can be described by a single plane wave whose amplitude varies perpendicular to its direction of propagation. The resulting two-dimensional coupled partial differential equations give good results when propagative diffraction effects can be neglected.⁶² Generalization of the one-dimensional coupled-wave equations to multi-coupled wave theory also works quite well,⁶³ as does the further generalization to coupling of continuous distributions of plane waves summed by integration.²⁰ This latter method has been used quite successfully to model beam fanning.²⁰ The split-step beam propagation method has also been used to include diffractive beam propagation effects as well as nonlinear models of the optically induced grating formation.⁶⁴

Spatiotemporal instabilities have been studied in phase conjugate resonators.^{65,66} These instabilities often involve optical vortices.⁶⁷

12.2 MATERIALS

Introduction

Photorefractive materials have been used in a wide variety of applications, as will be discussed later. These materials have several features which make them particularly attractive.

- The characteristic phase shift between the writing intensity pattern and the induced space charge field leads to energy exchange between the two writing beams, amplified scattered light (beam fanning), and self-pumped oscillators and conjugators.

- Photorefractive materials can be highly efficient at power levels obtained using CW lasers. Image amplification with a gain of 4000⁶⁸ and degenerate four-wave mixing with a reflectivity of 2000 percent⁶⁹ have been demonstrated.
- In optimized bulk photorefractive materials, the required energy to write a grating can approach that of photographic emulsion (50 $\mu\text{J}/\text{cm}^2$), with even lower values of write energy measured in photorefractive multiple quantum wells.
- The response time of most bulk photorefractive materials varies inversely with intensity. Gratings can be written with submillisecond response times at CW power levels and with nanosecond response times using nanosecond pulsed lasers. Most materials have a useful response with picosecond lasers.
- The high dark resistivity of oxide photorefractive materials allows the storage of holograms for time periods up to a year in the dark.

In spite of the great appeal of photorefractive materials, they have specific limitations which have restricted their use in practical devices. For example, oxide ferroelectric materials are very efficient, but are rather insensitive. Conversely, the bulk compound semiconductors are extremely sensitive, but suffer from low efficiency in the absence of an applied field. In this section we will first review the figures of merit used to characterize photorefractive materials, and then discuss the properties of the different classes of materials.

Figures of Merit

The figures of merit for photorefractive materials can be conveniently divided into those which characterize the steady-state response, and those which characterize the early portion of the transient response⁷⁰⁻⁷². Most applications fall into one or the other of these regimes, although some may be useful in either regime. For example, aberration correction, optical limiting, and laser coupling are applications which generally require operation in the steady state. On the other hand, certain optical processing applications require a response only to a given level of index change or efficiency, and are thus better characterized by the initial recording slope of a photorefractive grating.

Steady-State Performance The steady-state change in the refractive index is related to the space charge electric field by

$$\Delta n_{ss} = \frac{1}{2} n_b^3 r_{\text{eff}} E_{sc} \quad (10)$$

where n_b is the background refractive index, r_{eff} is the effective electro-optic coefficient (which accounts for the specific propagation direction and optical polarization in the sample), and E_{sc} is the space charge electric field. For large grating periods where diffusion limits the space charge field, the magnitude of the field is independent of material parameters and thus $\Delta n_{ss} \propto n_b^3 r_{\text{eff}}$. In this case (which is typical for many applications), the ferroelectric oxides are favored, because of their large electro-optic coefficients. For short grating periods or for very large applied fields, the space charge field is trap-limited and $\Delta n_{ss} \propto n_b^3 r_{\text{eff}} / \epsilon_r$, where ϵ_r is the relative dielectric constant.

The temporal behavior of the local space charge field in a given material depends on the details of the energy levels which contribute to the photorefractive effect. In many cases, the buildup or decay of the field is exponential. The fundamental parameter which characterizes this transient response is the write or erase energy W_{sat} . In many materials the response time τ at an average intensity I_0 is simply given by

$$\tau = W_{\text{sat}} / I_0 \quad (11)$$

This clearly points out the dependence of the response time on the intensity. As long as enough energy is provided, photorefractive gratings can be written with beams ranging in intensity from mW/cm^2 to MW/cm^2 . The corresponding response time can be calculated simply from Eq. (11).

In the absence of an applied or photovoltaic field (assuming only a single charge carrier), the response time can be written as⁷⁰

$$\tau = \tau_{\text{di}}(1 + 4\pi^2 L_d^2 / \Lambda_g^2) \quad (12)$$

where τ_{di} is the dielectric relaxation time, L_d is the diffusion length, and Λ_g is the grating period. The diffusion length is given by

$$L_d = (\mu \tau_r k_B T / e)^{1/2} \quad (13)$$

where μ is the mobility, τ_r is the recombination time, k_B is Boltzmann's constant, T is the temperature, and e is the charge of an electron. The dielectric relaxation time is given by

$$\tau_{\text{di}} = \epsilon_r \epsilon_0 / (\sigma_d + \sigma_p) \quad (14)$$

where σ_d is the dark conductivity and σ_p is the photoconductivity, given by

$$\sigma_p = \alpha e \mu \tau_r I_0 / h\nu \quad (15)$$

where α is the absorption coefficient.

In as-grown oxide ferroelectric materials, the diffusion length is usually much less than the grating period. In this case, the response time is given by

$$\tau = \tau_{\text{di}} \quad (16)$$

In addition, the contribution from dark conductivity in Eq. (14) can be neglected for intensities greater than $\sim \text{mW}/\text{cm}^2$. In this regime, materials with large values of absorption coefficient and photoconductivity ($\mu \tau_r$) are favored.

In bulk semiconductors, the diffusion length is usually much larger than the grating period. In this case, the response time is given by

$$\tau \approx \tau_{\text{di}}(4\pi^2 L_d^2 / \Lambda_g^2) \quad (17)$$

If we again neglect the dark conductivity in Eq. (14), then

$$\tau \approx \pi \epsilon k_B T / e^2 \alpha I_0 \Lambda_g^2 \quad (18)$$

In this regime, materials with small values of dielectric constant and large values of absorption coefficient are favored. For a typical bulk semiconductor with $\epsilon_r \approx 12$, $\alpha = 1 \text{ cm}^{-1}$, and $\Lambda_g = 1 \mu\text{m}$, we find $W_{\text{sat}} \approx 100 \mu\text{J}/\text{cm}^2$. This saturation energy is comparable to that required to expose high resolution photographic emulsion. In photorefractive multiple quantum wells (with $\alpha \approx 10^{13} \text{ cm}^{-1}$), the saturation energy can be much smaller. Note finally that an applied field leads to an increase in the write energy in the bulk semiconductors.⁷²

Transient Performance In the transient regime, we are typically concerned with the time or energy required to achieve a design value of index change or diffraction efficiency. This generally can be obtained from the initial *recording slope* of a photorefractive grating. One common figure of merit which characterizes the recording slope is the sensitivity,⁷⁰⁻⁷² defined as the index change per absorbed energy per unit volume:

$$S = \Delta n / \alpha I_0 \tau = \Delta n / \alpha W_{\text{sat}} \quad (19)$$

TABLE 1 Materials Parameters for BaTiO₃, BSO, and GaAs

Material Class	Ferroelectric Oxide	Ferroelectric Nonoxide	Sillenite	Compound Semiconductor
Material	BaTiO ₃	Sn ₂ P ₂ S ₆	BSO	GaAs
Wavelength range (μm)	0.4–1.1	0.65–1.3	0.45–0.65	0.9–1.3
Electro-optic coefficient r_{eff} (pm/V)	100 (r_{33}) 1640 (r_{42})	174(r_{11}) 92(r_{21}) 140(r_{31}) –25(r_{51})	4(r_{41})	1.4(r_{41})
Dielectric constant	135(ϵ_{33}) 3700(ϵ_{11})	230–300 (ϵ_{11})	56	13.2
$n_b^3 r_{\text{eff}}^3 / \epsilon$ (pm ³ /V)	10(r_{33}) 6(r_{42})	18(r_{11})	1.4	3.3
Mobility μ (cm ² /V-s)	0.01		0.1	6000
Recombination time τ_r (s)	10 ^{–8}		10^{–6}	3 × 10 ^{–8}
Diffusion length L_d (μm)	0.01	0.55	0.5	20
Photoconductivity $\mu \tau_r$ (cm ² /V)	10 ^{–10}	1.6 × 10 ^{–7}	10 ^{–7}	1.8 × 10 ^{–4}

(The parameters in bold type are particularly distinctive for that material.)

In the absence of an applied field and for large diffusion lengths, we find the limiting value of the sensitivity:

$$S \approx 1/4\pi(n_b^3 r_{\text{eff}}^3 / h\nu\epsilon)(me/\epsilon_0)\Lambda_g \quad (20)$$

where m is the modulation index.

We will see that the only material dependence in the limiting sensitivity is through the figure of merit $n_b^3 r_{\text{eff}}^3 / \epsilon$. This quantity varies little from material to material. Using typical values of the materials parameters and assuming $\Lambda_g = 1 \mu\text{m}$, we find the limiting value $S = 400 \text{ cm}^3/\text{kJ}$. Values in this range are routinely observed in the bulk semiconductors. In the as-grown ferroelectric oxides, in which the diffusion length is generally less than the grating period, the sensitivity values are typically two to three orders of magnitude smaller.

Comparison of Materials In the following sections, we will briefly review the specific properties of photorefractive materials, organized by crystalline structure. To introduce this discussion, we have listed relevant materials parameters in Table 1. This table allows the direct comparison among BaTiO₃, Sn₂P₂S₆ (SPS), Bi₁₂SiO₂₀ (BSO), and GaAs, which are representative of the four most common classes of photorefractive materials.

The distinguishing feature of BaTiO₃ is the magnitude of its electro-optic coefficients, leading to large values of steady-state index change. The sillenites are distinguished by their large value of recombination time, leading to a larger photoconductivity and diffusion length. The compound semiconductors are distinguished by their large values of mobility, leading to very large values of photoconductivity and diffusion length. Note also the different spectral regions covered by these four materials.

Ferroelectric Materials

The photorefractive effect was first observed in ferroelectric oxides that were of interest for electro-optic modulators and second-harmonic generation.⁷³ Initially, the effect was regarded as “optical

damage” that degraded device performance.⁷⁴ Soon, however, it became apparent that refractive index gratings could be written and stored in these materials.⁷⁵ Since that time, extensive research on material properties and device applications has been undertaken.

The photorefractive ferroelectric oxides can be divided into three structural classes: ilmenites (LiNbO_3 , LiTaO_3), perovskites [BaTiO_3 , KNbO_3 , $\text{KTa}_{1-x}\text{Nb}_x\text{O}_3$ (KTN)], and tungsten bronzes [$\text{Sr}_{1-x}\text{Ba}_x\text{Nb}_2\text{O}_6$ (SBN), $\text{Ba}_2\text{NaNb}_5\text{O}_{15}$ (BNN) and related compounds]. In spite of their varying crystal structure, these materials have several features in common. They are transparent from the bandgap (~ 350 nm) to the intrinsic IR absorption edge near $4 \mu\text{m}$. Their wavelength range of sensitivity is also much broader than that of other photorefractive materials. For example, useful photorefractive properties have been measured in BaTiO_3 from 442 nm ⁷⁶ to $1.09 \mu\text{m}$,⁷⁷ a range of a factor of $2\frac{1}{2}$. Ferroelectric oxides are hard, nonhygroscopic materials—properties which are advantageous for the preparation of high-quality surfaces. Their linear and nonlinear dielectric properties are inherently temperature-dependent, because of their ferroelectric nature. As these materials are cooled below their melting point, they undergo a structural phase transition to a ferroelectric phase. Additional transitions may occur in the ferroelectric phase on further lowering of the temperature. In general, samples in the ferroelectric phase contain regions of differing polarization orientation called domains, leading to a reduction in the net polar properties of the sample. To make use of the electro-optic and nonlinear optic properties of the ferroelectric oxides, these domains must be aligned to a single domain state. This process, called poling, can take place during the growth process, or more commonly, after polydomain samples have been cut from an as-grown boule.

Growth of large single crystals of ferroelectric oxides has been greatly stimulated by the intense interest in photorefractive and nonlinear optic applications. Currently, most materials of interest are commercially available. However, considerably more materials development is required before optimized samples for specific applications can be purchased.

Lithium Niobate and Lithium Tantalate LiNbO_3 was the first material in which photorefractive “damage” was observed.⁷³ This material has been developed extensively for frequency conversion and integrated optics applications. It is available in large samples with high optical quality. For photorefractive applications, iron-doped samples are generally used. The commonly observed valence states are Fe^{2+} and Fe^{3+} . The relative populations of these valence states can be controlled by annealing in an atmosphere with a controlled oxygen partial pressure. In a reducing atmosphere (low oxygen partial pressure), Fe^{2+} is favored, while Fe^{3+} is favored in an oxidizing atmosphere. The relative $\text{Fe}^{2+}/\text{Fe}^{3+}$ population ratio will determine the relative contributions of electrons and holes to the photoconductivity.⁷⁸ When Fe^{2+} is favored, the dominant photocarriers are electrons; when Fe^{3+} is favored, the dominant photocarriers are holes. In most oxides, electrons have higher mobilities, so that electron-dominated samples yield faster photorefractive response times. Even in heavily reduced LiNbO_3 , the write energy is rarely lower than 10 J/cm^2 , so this material has not found use for real-time applications. The properties of LiTaO_3 are essentially the same as those of LiNbO_3 .

Currently, the most promising application of LiNbO_3 is for holographic storage. LiNbO_3 is notable for its very large value of dark resistivity, leading to very long storage times in the dark. In addition, the relatively large write or erase energy of LiNbO_3 makes this material relatively unsusceptible to erasure during readout of stored holograms.

Improved retention of stored holograms can be obtained by fixing techniques. The most common fixing approach makes use of complementary gratings produced in an ionic species which is not photoactive.^{79, 80} Typically, one or more holograms are written into the sample by conventional means. The sample is then heated to 150°C , where it is annealed for a few hours. At this temperature, a separate optically inactive ionic species is thermally activated and drifts in the presence of the photorefractive space charge field until it compensates this field. The sample is then cooled to room temperature to “freeze” the compensating ion grating. Finally, uniform illumination washes out the photorefractive grating and “reveals” the permanent ion grating.

Another important feature of LiNbO_3 for storage applications is the large values of diffraction efficiency (approaching 100 percent for a single grating) which can be obtained. These large efficiencies arise primarily from the large values of space charge field, which, in turn, result from the very large value of photovoltaic field.

Barium Titanate BaTiO₃ was one of the first ferroelectric materials to be discovered, and also one of the first to be recognized as photorefractive.⁸¹ The particular advantage of BaTiO₃ for photorefractive applications is the very large value of the electro-optic coefficient r_{42} (see Table 1), which, in turn, leads to large values of grating efficiency, beam-coupling gain, and conjugate reflectivity. For example, four-wave mixing reflectivities as large as 20 have been observed,⁶⁹ as well as an image intensity gain of 4000.⁶⁸

After the first observation of the photorefractive effect in BaTiO₃ in 1970, little further research was performed until 1980 when Feinberg et al.⁸² and Krätzig et al.⁸³ pointed out the favorable features of this material for real-time applications. Since that time, BaTiO₃ has been widely used in a large number of experiments in the areas of optical processing, laser power combining, spatial light modulation, optical limiting, and neural networks.

The photorefractive properties of BaTiO₃ have been reviewed in Ref. 84. Crystals of this material are grown by top-seeded solution growth in a solution containing excess TiO₂.⁸⁵ Crystal growth occurs while cooling the melt from 1400 to 1330°C. At the growth temperature, BaTiO₃ has the cubic perovskite structure, but on cooling through $T_c = 132^\circ\text{C}$, the crystal undergoes a transition to the tetragonal ferroelectric phase. Several approaches to poling have been successfully demonstrated. In general, the simplest approach is to heat the sample to just below or just above the Curie temperature, apply an electric field, and cool the sample with the field present.

Considerable efforts have been expended to identify the photorefractive species in as-grown BaTiO₃. Early efforts suggested that transition metal impurities (most likely iron) were responsible.⁸⁶ In later experiments, samples grown from ultrapure starting materials were still observed to be photorefractive.⁸⁷ In this case, barium vacancies have been proposed as the dominant species.⁸⁸ Since that time, samples have been grown with a variety of transition metal dopants. All dopants produce useful photorefractive properties, but cobalt-doped samples⁸⁹ and rhodium-doped samples⁹⁰ appear particularly promising, because of their reproducible high gain in the visible and enhanced sensitivity in the infrared.⁹⁰

One particular complication in developing a full understanding of the photorefractive properties of BaTiO₃ is the presence of shallow levels in the bandgap, in addition to the deeper levels typically associated with transition-metal impurities or dopants. The shallow levels are manifested in several ways. Perhaps the most prominent of these is the observation that the response time (and photoconductivity) of as-grown samples does not scale inversely with intensity [see Eq. (11)], but rather has a dependence of the type $\tau \sim (I_0)^{-x}$ is observed, where $x = 0.6\text{--}1.0$.^{81-84,91} Several models relating to the sublinear behavior of the response time and photoconductivity to shallow levels have been reported.⁹²⁻⁹⁵

The characteristic sublinear variation of response time with intensity implies that the write or erase energy W_{sat} increases with intensity, which is a clear disadvantage for high-power, short-pulse operation. Nevertheless, useful gratings have been written in BaTiO₃ using nanosecond pulses^{96,97} and picosecond pulses.⁹⁸ Another manifestation of the presence of shallow levels is intensity-dependent absorption.⁹⁹ The shallow levels have been attributed to oxygen vacancies or barium vacancies, but no unambiguous identification has been made to date.

The major limitation of BaTiO₃ for many applications is the relatively slow response time of this material at typical CW intensity levels. In as-grown crystals, typical values of response time are 0.1 to 1 s at 1 W/cm². These values are approximately three orders of magnitude longer than theoretical values determined from the band transport model (see "Steady-State Performance"), or from more fundamental arguments.^{100,101}

Two different approaches have been studied to improve the response time of BaTiO₃. In the first, as-grown samples can be operated at an elevated temperature (but below the Curie temperature). In a typical experiment (see Fig. 5), an improvement in response time by two orders of magnitude was observed¹⁰² when different samples were operated at 120°C. In some of these samples, the magnitude of the peak beam-coupling gain did not vary significantly with temperature. In these cases, the improvement in response time translates directly to an equivalent improvement in sensitivity.

While operation at an elevated temperature may not be practical for many experiments, the importance of the preceding experiment is that it demonstrates the *capability for improvement in response time*, based on continuing materials research.

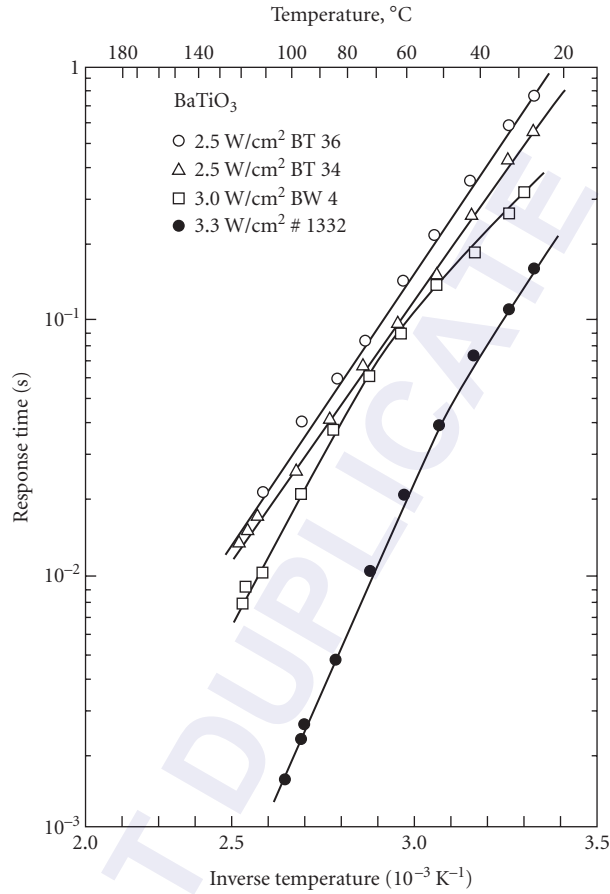


FIGURE 5 Measured response time as a function of temperature for four samples of BaTiO₃. The measurement wavelength was 515 nm and the grating period was 0.79 μm.

Other materials research efforts concentrated on studies of new dopants, as well as heat treatments in reducing atmospheres.^{88,89,103,104} The purpose of the reducing treatments is to control the valence states of the dopants to produce beneficial changes of the trap density and the sign of the dominant photocarrier. While some success has been achieved,^{89,104} a considerably better understanding of the energy levels in the BaTiO₃ bandgap is required before substantial further progress can be made.

Potassium Niobate KNbO₃ is another important photorefractive material with the perovskite structure. It undergoes the same sequence of phase transitions as BaTiO₃, but at higher transition temperatures. At room temperature it is orthorhombic, with large values of the electro-optic coefficients r_{42} and r_{51} .

KNbO₃ has been under active development for frequency conversion and photorefractive experiments since 1977 (Ref. 70). Unlike BaTiO₃, undoped samples of KNbO₃ have weak photorefractive properties. Iron doping has been widely used for photorefractive applications,⁷⁰ but other transition metals have also been studied.

Response times in as-grown KNbO₃ at 1 W/cm² are somewhat faster than those of BaTiO₃, but are still several orders of magnitude longer than the limiting value. The most common approach

to improving the response time of KNbO_3 is electrochemical reduction. In one experiment at 488 nm,¹⁰¹ a photorefractive response time of 100 μs at 1 W/cm^2 was measured in a reduced sample. This response time is very close to the limiting value, which indicates again the promise for faster performance in all the ferroelectric oxides.

Strontium Barium Niobate and Related Compounds $\text{Sr}_{1-x}\text{Ba}_x\text{Nb}_2\text{O}_6$ (SBN) is a member of the tungsten bronze family,¹⁰⁵ which includes materials such as $\text{Ba}_2\text{NaNb}_5\text{O}_{15}$ (BNN) and $\text{Ba}_{1-x}\text{Sr}_x\text{K}_{1-y}\text{Na}_y\text{Nb}_5\text{O}_{15}$ (BSKNN). SBN is a mixed composition material with a phase transition temperature which varies from 60 to 200°C as x varies from 0.75 to 0.25. Of particular interest is the composition SBN-60, which melts congruently,¹⁰⁶ and is thus easier to grow with high quality. SBN is notable for the very large values of the electro-optic tensor component r_{33} . In other materials such as BSKNN, the largest tensor component is r_{42} . In this sense it resembles BaTiO_3 .

In general, the tungsten bronze system contains a large number of mixed composition materials, thus offering a rich variety of choices for photorefractive applications. In general, the crystalline structure is quite open, with only partial occupancy of all lattice sites. This offers greater possibilities for doping, but also leads to unusual properties at the phase transition, due to its diffuse nature.¹⁰⁵

The photorefractive properties of SBN were first reported in 1969,¹⁰⁷ very soon after gratings were first recorded in LiNbO_3 . Since that time, there has been considerable interest in determining the optimum dopant for this material. The most common dopant has been cerium.^{108–110} Cerium-doped samples can be grown with high optical quality and large values of photorefractive gain.^{111,112} Another promising dopant is rhodium, which also yields high values of gain coefficient.¹¹³

As with BaTiO_3 and KNbO_3 , as-grown samples of SBN and other tungsten bronzes are relatively slow at an intensity of 1 W/cm^2 .¹¹⁴ Doping and codoping has produced some improvement. In addition, the use of an applied dc electric field has led to improvement in the response time.¹¹⁵

The photorefractive effect has also been observed in fibers of SBN.¹¹⁶ The fiber geometry has promise in holographic storage architectures.

Tin Hypothiodiphosphate $\text{Sn}_2\text{P}_2\text{S}_6$ has been known as a ferroelectric since 1974,¹¹⁷ but has been investigated as a photorefractive material only since 1991.¹¹⁸ Its Curie point is 337 K, only a few tens of degrees above room temperature, so its electro-optic coefficients are expected to be high (Table 1). It is distinguished from the ferroelectric oxide photorefractives by its useful wavelength range and speed of response. Its bandgap is narrower than that of typical photorefractive oxides, so its wavelength range is pushed deeper into the infrared, and will operate with high gain from 0.65 to 1.3 μm ^{119,120} and at least as far as 1.55 μm for tellurium doped crystals.¹²¹ There are several variants of nominally undoped material, known by their color (type I yellow, type II yellow, and modified brown), as well as doped crystals, each with their different characteristics. The properties of type I yellow crystals depend on their history of illumination, and are characterized by the existence of a photoinduced fast grating mediated by positive charge carriers and a thermally induced slow grating mediated by electrons. By virtue of the fact that these are due to oppositely charged carriers, they are 180° out of phase with each other, and tend to cancel each other out. In type II yellow crystals, the slow grating is suppressed, thus improving the steady-state gain.¹²² The response time of the fast grating at 1.06 μm is 300 ms at 1 W/cm^2 and is inversely proportional to intensity while the response time of the slow grating is of the order of 100 s, and is approximately independent of intensity. Brown crystals are produced by modifying the vapor transport crystal growth method in such a way as to increase the concentration of intrinsic defects. Typical photorefractive properties for type II, brown, and Te-doped crystals are shown in Table 2.

Cubic Oxides (Sillenites)

The cubic oxides are notable for their high photoconductivity, leading to early applications for spatial light modulation¹²³ and real-time holography using the photorefractive effect.¹²⁴ The commonly used sillenites are $\text{Bi}_{12}\text{SiO}_{20}$ (BSO), $\text{Bi}_{12}\text{GeO}_{20}$ (BGO), and $\text{Bi}_{12}\text{TiO}_{20}$ (BTO). Some relevant properties of these materials are listed in Table 3.

The sillenites are cubic and noncentrosymmetric, with one nonzero electro-optic tensor component r_{41} . The magnitude of r_{41} in the sillenites is small, ranging from approximately 4 to 6 pm/V in the

TABLE 2 Typical Photorefractive Parameters of Various $\text{Sn}_2\text{P}_2\text{S}_6$ Crystals at Two Light Wavelengths

$\text{Sn}_2\text{P}_2\text{S}_6$ Sample	λ (nm)	α_x (cm^{-1})	Γ_{max} (cm^{-1})	τ (ms)	$N_{\text{eff}} 10^{16} \text{cm}^{-1}$
Yellow type II	633	0.5	4–7	10–50	0.7
	780	0.2	2–5	100	0.2
Brown	633	5.7	38	4	2.5
	780	1.0	18	10	0.7
Te-doped (1%)	633	1.0	10	2.5	0.9
	780	0.4	6	7	1.0

λ , without pre-illumination; α_x , absorption coefficient for x -polarized light; Γ_{max} , maximal two-wave mixing gain; τ , faster response time at a grating spacing of 1 μm and scaled to a light intensity of 1 W/cm^2 ; N_{eff} , effective trap density. (After Grabar et al.¹¹⁷)

TABLE 3 Material Properties of BSO, BGO, and BTO

Material	BSO	BGO	BTO
Wavelength range (μm)	0.5–0.65	0.5–0.65	0.6–0.75
Electro-optic coefficient r_{41} (pm/V)	4.5	3.4	5.7
$n_b^3 r_{41}$ (pm/V)	81	56	89
Dielectric constant	56	47	48
$n_b^3 r_{41} / \epsilon$ (pm/V)	1.4	1.2	1.9
Optical activity at 633 nm (degrees/mm)	21	21	6

visible. In addition, the sillenites are optically active, with a rotatory power (at 633 nm) of 21°/mm in BSO and BGO, and 6°/mm in BTO. These values increase sharply at shorter wavelengths. The optical activity of the sillenites tends to reduce the effective gain or diffraction efficiency of samples with normal thickness, but in certain experiments it also allows the use of an output analyzer to reduce noise.

The energy levels due to defects and impurities tend to be similar in each of the sillenites. In spite of many years of research, the identity of the photorefractive species is still not known. It is likely that intrinsic defects such as metal ion vacancies play an important role. With only one metal ion for each 12 bismuth ions and 20 oxygen ions, small deviations in metal ion stoichiometry can lead to large populations of intrinsic defects. In each of the sillenites, the effect of the energy levels in the bandgap is to shift the fundamental absorption edge approximately 100 nm to the red.

BSO and BGO melt congruently and can be grown from stoichiometric melts by the Czochralski technique. On the other hand, BTO melts incongruently and is commonly grown by the top-seeded solution growth technique, using excess Bi_2O_3 . BTO is particularly interesting for photorefractive applications (compared with BSO and BGO), because of its lower optical activity at 633 nm and its slightly larger electro-optic coefficient (5.7 pm/V).¹²⁵ It has been studied extensively at the Ioffe Institute in Russia, where both material properties and device applications have been examined.^{126–128}

In the sillenites it is very common to apply large dc or ac electric fields to enhance the photorefractive space charge field, and thus provide useful values of gain or diffraction efficiency. A dc field will increase the amplitude of the space charge field, but the spatial phase will decrease from the value of 90° which optimizes the gain. In order to restore the ideal 90° phase shift, moving grating techniques are typically used.^{129,130} By contrast, an ac field can enhance the amplitude of the space charge field, while maintaining the spatial phase at the optimum value of 90°.¹² In this case, the best performance is obtained when a square waveform is used, and when the period is long compared with the recombination time, and short compared with the grating formation time. Both dc and ac field techniques have produced large gain enhancements in sillenites and semiconductors, but only when the signal beam is very weak, i.e., when the pump/signal intensity ratio is large. As the amplitude of the signal beam increases, the gain decreases sharply, by an amount which cannot be explained by pump depletion. This effect is significant for applications such as self-pumped phase conjugation, in which the buildup of the signal wave will reduce the effective gain and limit the device performance.

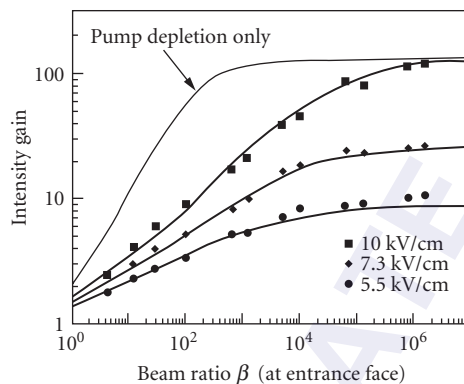


FIGURE 6 Measured two-wave mixing gain as a function of input pump-to-signal beam ratio in BTO. The measurement wavelength was 633 nm, the applied field was a 60-Hz ac square wave, and the grating period was 5.5 μm . The individual points are experimental data; the bold curves are fits using a large signal model. The thin curve is the standard pump depletion theory for the 10-kV/cm case.

A typical plot of intensity gain in BTO as a function of beam ratio for several values of ac square-wave voltage amplitude is given in Fig. 6.¹²⁷ Note that the highest gain is observed only for a beam ratio on the order of 10^5 (small signal limit). The simplest physical description of this nonlinearity is that the internal space charge field is clamped to the magnitude of the applied field; this condition only impacts performance for decreasing beam ratios (large signal limit). In a carefully established experiment using a very large beam ratio (10^5), gain coefficients approaching 35 cm^{-1} have been measured using an ac square wave field with an amplitude of 10 kV/cm (see Fig. 7).¹³¹

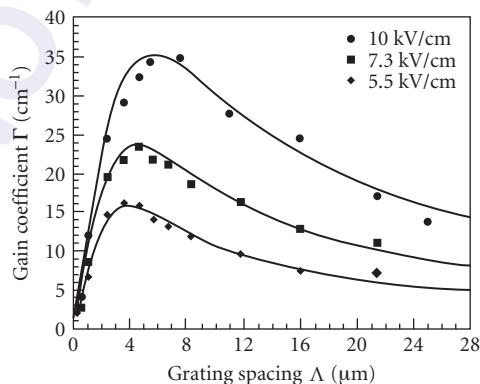


FIGURE 7 Measured gain coefficient as a function of grating spacing in BTO. The measurement wavelength was 633 nm, the applied field was a 60-Hz ac square wave, and the beam ratio was 10^5 . The individual points are experimental data; the solid curves are fits using the basic band transport model.

Bulk Compound Semiconductors

The third class of commonly used photorefractive materials consists of the compound semiconductors (Si and Ge are cubic centrosymmetric materials, and thus have no linear electro-optic effect). Gratings have been written in CdS,¹³² GaAs:Cr,¹³³ GaAs:EL2,¹³⁴ InP:Fe,¹³³ CdTe,¹³⁵ GaP,¹³⁶ and ZnTe.¹³⁷ These materials have several attractive features for photorefractive applications (see Table 4). First, many of these semiconductors are readily available in large sizes and high optical quality, for use as electronic device substrates. These substrates are generally required to be semi-insulating; the deep levels provided for this purpose are generally photoactive, with favorable photorefractive properties. Second, the semiconductors have peak sensitivity for wavelengths in the red and near-infrared. The range of wavelengths extends from 633 nm in GaP,¹³⁶ CdS,¹³⁹ and ZnTe¹³⁷ to 1.52 μm in CdTe:V.¹³⁵ Third, the mobilities of the semiconductors are several orders of magnitude larger than those in the oxides. There are several important consequences of these large mobilities. Most importantly, the resulting large diffusion lengths lead to fast response times [see Eqs. (12) and (18)]. The corresponding values of write/erase energies (10 to 100 $\mu\text{J}/\text{cm}^2$) are very near the limiting values. These low values of write/erase energy have been observed not only at the infrared wavelengths used for experiments in InP and GaAs, but also at 633 nm in ZnTe.¹³⁷

The large mobilities of the compound semiconductors also yield large values of dark conductivity (compared with the oxides), so the storage times in the dark are normally less than 1 s. Thus, these materials are not suited for long-term storage, but may still be useful for short-term memory applications. Finally, the short diffusion times in the semiconductors yield useful photorefractive performance with picosecond pulses.⁹⁶

The electro-optic coefficients for the compound semiconductors are quite small (see Table 4), leading to low values of beam-coupling gain and diffraction efficiency in the absence of an applied electric field. As in the sillenites, both dc and ac field techniques have been used to enhance the space charge field. Early experiments with applied fields produced enhancements in the gain or diffraction efficiency which were considerably below the calculated values.^{140–142} The causes of these discrepancies are now fairly well understood. First, space charge screening can significantly reduce the magnitude of the applied field inside the sample. This effect is reduced by using an ac field, but even in this case the required frequencies to overcome all screening effects are quite high.¹⁴² Second, the mobility-lifetime product is known to reduce at high values of electric field due to scattering of electrons into other conduction bands and cascade recombination. This effect is particularly prominent in GaAs.¹⁴³ Third, large signal effects act to reduce the gain when large fields are used.¹⁴⁴ As in the sillenites, the highest gains are only measured when weak signal beams (large pump/signal beam ratios) are used. Finally, when ac square-wave fields are used, the theoretical gain value is only obtained for sharp transitions in the waveform.^{145,146}

Another form of electric field enhancement has been demonstrated in iron-doped InP.^{147,148} This material is unique among the semiconductors in that the operating temperature and incident intensity can be chosen so that the photoconductivity (dominated by holes) exactly equals the dark conductivity (dominated by electrons). In this case, an applied field will enhance the amplitude of the space charge field, while maintaining the ideal spatial phase of 90°. Gain coefficients as high as 11 cm^{-1} have been reported in InP:Fe at 1.06 μm using this technique.¹⁴⁸

TABLE 4 Relevant Materials Properties of Photorefractive Compound Semiconductors

Material	GaAs	InP	GaP	CdTe	ZnTe
Wavelength range (μm)	0.92–1.3	0.96–1.3	0.63	1.06–1.5	0.63–1.3
EO coeff. r_{41} (pm/V)	1.2	1.45	1.1	6.8	4.5
$n_b^3 r_{41}$ (pm/V)	43	52	44	152	133
$n_b^3 r_{41}/\epsilon$	3.3	4.1	3.7	16	13
Dielectric constant	13.2	12.6	12	9.4	10.1

(Most of the values are taken from Ref. 138.)

In early research on the photorefractive semiconductors, the wavelengths of operation were determined by available laser sources. Thus, all early experiments were performed at 1.06 μm (Nd:YAG), $\sim 1.3 \mu\text{m}$ (Nd:YAG or laser diode), and $\sim 1.5 \mu\text{m}$ (laser diode). Later, the He-Ne laser (633 nm) was used to study wider bandgap materials. As the Ti:sapphire laser became available, interest turned to investigating the wavelength variation of photorefractive properties, especially near the band edge. Near the band edge, a new nonlinear mechanism contributes to the refractive index change: the Franz-Keldysh effect.¹⁴⁹ In this case, the internal space charge field develops as before. This field slightly shifts the band edge, leading to characteristic electroabsorption and electrorefraction. These effects can be quite large at wavelengths near the band edge, where the background absorption is also high. However, the peak of the electrorefraction spectrum is shifted slightly to longer wavelengths, where the background absorption is smaller. This is generally the wavelength region where these effects are studied.

The electrorefractive photorefractive (ERPR) effect has different symmetry properties than the conventional electro-optic photorefractive effect. It is thus possible to arrange an experiment so that only the ERPR effect contributes, or both effects contribute to the gain. In addition, the ERPR effect is quadratic in applied electric field. Thus, energy transfer between two writing beams only occurs when a dc field is present. The direction of energy transfer is determined by the sign of the electric field; this allows switching energy between two output beams via switching of the sign of the applied field.

In the first report of the band-edge photorefractive effect,¹⁴⁹ a gain coefficient of 7.6 cm^{-1} was measured in GaAs:EL2 at 922 nm, for a field of 10 kV/cm. In this case, both nonlinear mechanisms contributed to the gain. When a moving grating was used to optimize the spatial phase of the grating, the gain coefficient increased to 16.3 cm^{-1} .

In InP the temperature/intensity resonance can be used to optimize the spatial phase, thus eliminating the need for a moving grating. In the first experiment using band-edge resonance and temperature stabilization, gain coefficients approaching 20 cm^{-1} were measured in InP:Fe (see Fig. 8).¹⁵⁰ Later experiments on a thin sample using a beam ratio of 10^6 resulted in a measured gain coefficient of 31 cm^{-1} .¹⁵¹

The photorefractive effect can also be used to measure basic materials properties of electro-optic semiconductors, without the need for electrical contacts.^{152–154} Quantities which can be measured include the populations of filled and empty traps and the mobility-recombination time product. One particular feature of the photorefractive technique is the ability to map properties across a wafer.¹⁵⁴

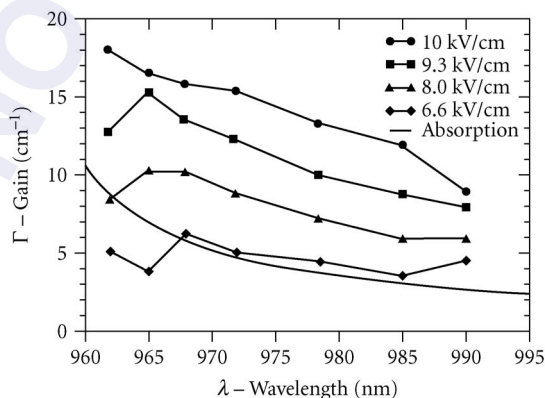


FIGURE 8 Measured gain coefficient as a function of wavelength in InP:Fe, for a grating period of 5 μm and four values of applied dc field. The beam ratio was 1000, and the intensity was adjusted at each point to produce the maximum gain. The background absorption coefficient is also plotted.

Multiple Quantum Wells

While the enhancement of the electro-optic effect near the band edge of bulk semiconductors is significant, much larger nonlinearities are obtained at wavelengths near prominent band-edge exciton features in multiple quantum wells (MQWs). In addition, the large absorption in these structures yields much faster response times than those in bulk semiconductors. Finally, the small device thickness of typical MQW structures (typically 1 to 2 μm) provides improved performance of Fourier plane processors such as optical correlators.^{155,156} One disadvantage of the small device thickness is that diffraction from gratings in these devices is in the Raman-Nath regime, yielding multiple diffraction orders.

In their early stages of development, MQWs were not optimized for photorefractive applications because of the absence of deep traps and the large background conductivity within the plane of the structure. It was later recognized that defects resulting from ion implantation can provide the required traps and increase the resistivity of the structure.

The first photorefractive MQWs were GaAs/AlGaAs structures which were proton-implanted for high resistivity ($\rho = 10^9/\text{ohm-cm}$).^{157,158} Two device geometries were considered (see Fig. 9), but only devices with applied fields parallel to the layers were studied. The principles of operation are initially the same for both device geometries. When two incident waves interfere in the sample, the spatially modulated intensity screens the applied field in direct proportion to the intensity, leading to a spatially modulated internal field. This spatially modulated field induces changes in both the refractive index and the absorption coefficient. The mechanism for these changes¹⁵⁸ is field ionization of excitons (Franz-Keldysh effect) in the parallel geometry and the quantum-confined Stark effect in the perpendicular geometry.

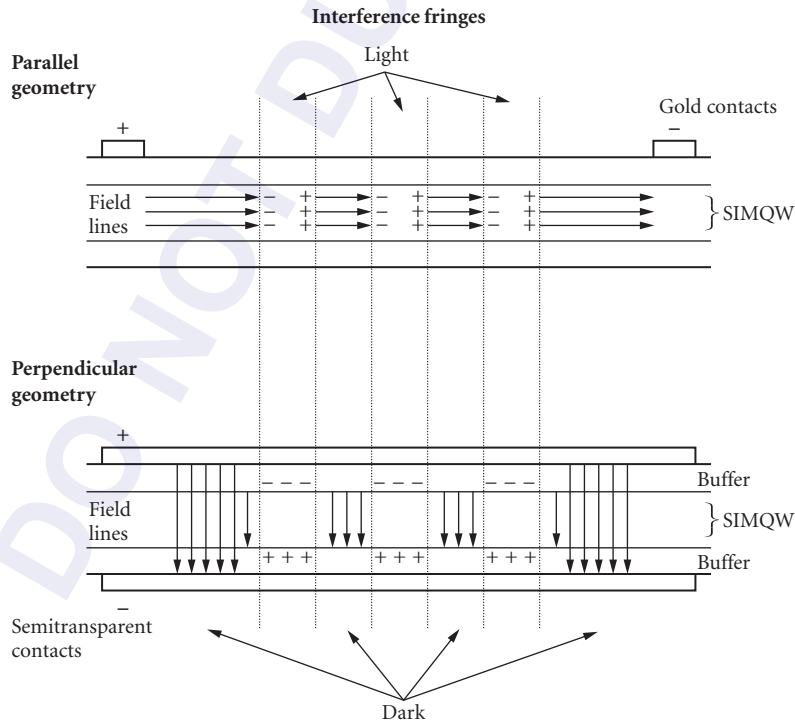


FIGURE 9 Device geometries for photorefractive MQWs.

The magnitudes of the change in refractive index and absorption coefficient are strongly dependent on wavelength near the characteristic exciton peak. Both the index and the absorption grating contribute to the diffraction efficiency through the relationship

$$\eta = (2\pi \Delta n L / \lambda)^2 + (\Delta \alpha L / 2)^2 \quad (21)$$

where Δn and $\Delta \alpha$ are the amplitudes of the index and absorption gratings, respectively, and L is the device thickness. Although the device thickness of typical MQWs is much less than the thickness of a typical bulk sample (by 3 to 4 orders of magnitude), the values of Δn can be made larger, leading to practical values of diffraction efficiency (see following).

The first III-V MQWs using the parallel geometry^{157,158} had rather small values of diffraction efficiency (10^{-5}). In later experiments, a diffraction efficiency of 3×10^{-4} and a gain coefficient of 1000 cm^{-1} were observed.¹⁵⁹ Still higher values of diffraction efficiency (on the order of 1.3 percent) were obtained using the perpendicular geometry in CdZnTe/ZnTe MQWs.¹⁶⁰ These II-VI MQWs have the added feature of allowing operation at wavelengths in the visible spectral region, in this case 596 nm.

In recent work on GaAs/AlGaAs MQWs in the perpendicular geometry, several device improvements were introduced.¹⁶¹ First, Cr-doping was used to make the structure semi-insulating, thus eliminating the added implantation procedure and allowing separate control of each layer. Second, the barrier thickness and Al ratio were adjusted to give a reduced carrier escape time, leading to a larger diffraction efficiency. In these samples, a diffraction efficiency of 3 percent was observed at 850 nm for an applied voltage of 20 V across a 2- μm -thick device. The response time was 2 μs at an intensity of 0.28 W/cm^2 , corresponding to a very low write energy of $0.56 \mu\text{J/cm}^2$. The diffraction efficiency cited here was obtained at a grating period of 30 μm . For smaller values of grating period, the diffraction efficiency was smaller, due to charge smearing effects. The fast response time and small thickness of these structures make them ideal candidates for Fourier plane processors such as optical correlators. Competing bulk semiconductors or spatial light modulators have frame rates which are 2 to 3 orders of magnitude below the potential frame rate of $\sim 10^6 \text{ s}^{-1}$ which is available from photorefractive MQWs.

Future work on photorefractive MQWs would include efforts to grow thicker devices (so as to reduce the diffraction into higher grating orders) and to improve the diffraction efficiency at high spatial frequencies.

Organic Crystals and Polymer Films

Organic materials are increasingly providing a viable alternative to their inorganic counterparts. Examples include organic crystals for frequency conversion applications and polymer films for electro-optic waveguide devices. These materials are, in general, simpler to produce than their inorganic counterparts. In addition, the second-order nonlinear coefficients in these materials can be quite large, with values comparable to those of the well-known inorganic material LiNbO_3 . In most organic materials the electronic nonlinearity results from an extended system of π electrons produced by electron donor and acceptor groups. For a purely electronic nonlinearity, the dielectric constant ϵ is just the square of the refractive index, leading to much smaller values of ϵ than those in inorganic crystals. Thus, the electro-optic figure of merit n_b^3/ϵ is enhanced in organic materials. This enhancement makes organic materials very appealing for photorefractive applications.

The first experiments reported on the photorefractive effect in organic crystals were those of Sutter et al.^{162,163} on 2-cyclooctylamino-5-nitropyridine (COANP) doped with the electron acceptor 7,7,8,8-tetracyanoquinodimethane (TCNQ). Pure COANP crystals (used for frequency doubling) are yellow, whereas the TCNQ-doped samples are green, due to a prominent extrinsic absorption band between 600 and 700 nm. In experiments at 676 nm with a grating period of 1.2 μm , both absorption and refractive index gratings were observed. Typical diffraction efficiencies were 0.1 percent, with a corresponding refractive index grating amplitude of 10^{-6} . The recorded buildup times of the index gratings were on the order of 30 to 50 min at 3.2 W/cm^2 . Following this initial demonstration of the photorefractive effect in organic crystals there has been very little progress since most attempts at doping these

crystals for photoexcitation and charge transport have simply resulted in expulsion of the dopants from the crystal structure. Only one additional crystal has been reported.¹⁶⁴

The situation is much more promising for the photorefractive effect in composite polymer films. These materials were first reported in 1991^{165,166} and from a materials science point of view, they are much easier to prepare than organic single crystals. In addition, there is greater flexibility in modifying the films to optimize photorefractive performance. In this respect, there are four requirements for an efficient photorefractive material:

1. A linear electro-optic effect, or a quadratic electrooptic effect with a linear component induced by a dc bias electric field
2. A source of photoionizable charges
3. A means of transporting these charges
4. A means of trapping the charges

In an organic polymer composite, each of these functions can be separately optimized. The space charge field due to the holographic interference pattern not only perturbs the refractive index via the electro-optic effect, but can also reorient molecules in the medium giving rise to a refractive index variation via polarizability anisotropy, an effect known as orientational enhancement.¹⁶⁷

In spite of the large EO coefficients and the great flexibility in the materials engineering of organic polymers, there are also some practical problems which need to be addressed. First, polymers are most easily prepared as thin films. If propagation in the plane of the films is desired, then extremely high optical quality and low absorption are required. If propagation through the film is desired, then the grating diffraction efficiency is reduced. As the film thickness is made larger to enhance efficiency, the quality of the films is harder to maintain.

One important requirement of polymer films is that they must be poled to induce a linear electro-optic effect. If the poling voltage is applied normal to the film plane, then there is no electro-optic effect for light diffracted from gratings written with their wave-vectors in the plane of the film. In a practical sense, this means that the writing beams must enter the film at large angles to its normal. In addition, it becomes more difficult to provide large poling fields as the film thickness increases.

The first polymeric photorefractive material^{165,166} was composed of the epoxy polymer bisphenol-A-diglycidylether 4-nitro-1,2-phenylenediamine (bisA-NPDA) made photoconductive by doping with the hole transport agent diethylamino-benzaldehyde diphenyl hydrazone (DEH). In this case, the polymer provided the nonlinearity leading to the electro-optic effect, as well as a mechanism for charge generation. The dopant provided a means for charge transport, while trapping was provided by intrinsic defects.

Films of this material with thicknesses between 200 and 500 μm were prepared. The material was not cross-linked, so a large field was required at room temperature to maintain the polarization of the sample. For an applied field of 120 kV/cm, the measured value of the electro-optic figure of merit $n_i^3 \gamma / \epsilon$ was 1.4 pm/V. Using interference fringes with a spacing of 1.6 μm oriented 25° from the film plane, the measured grating efficiency at 647 nm was 2×10^{-5} . The grating buildup time was on the order of 100 s at an intensity of 25 W/cm². Analysis of the data showed that the photorefractive trap density had the relatively small value of 2×10^{15} cm⁻³. In spite of the low value of trap density, relatively large values of space charge electric field were obtained, due to the low value of dielectric constant.

Subsequent research has led to general design principles for photorefractive polymer composites. Photoexcitation of charge is often accomplished by using donor-acceptor charge transfer complexes. In this way, the absorption spectrum can be tailored to the wavelengths of interest. Carbazole is often used as an electron donor entity, coupled with electron acceptors 2,4,7-trinitro-9-fluorenone (TNF),¹⁶⁸ (2,4,7-trinitro-9-fluorenylidene) malononitrile (TNFM),^{168,169} or C₆₀.^{170,171} As in the case of the formation of inorganic photorefractive gratings in response to illumination by the optical interference pattern of two intersecting laser beams, the photoexcited charges should be free to move away from the site of excitation and be preferentially retrapped in the darker regions of the interference pattern to form a spatially varying charge distribution following the interference pattern. However, in contrast to the case of inorganic crystals, diffusion is not effective in driving charge separation, so electric fields

have to be applied to force charge separation by drift. These fields are applied by sandwiching a several micrometer thick layer of the polymer system between transparent electrodes. The electrodes have to be tilted with respect to the grating wave vector so as to provide a component of the bias field parallel to the grating vector. Another way in which organics differ from inorganics is that both the photogeneration rate and mobility depend on the electric fields in the material. Hole mobility is usually much greater than electron mobility. This has the effect of allowing the hole grating to dominate the electron grating. If the hole and electron gratings were of similar magnitude, their electric fields would tend to cancel each other out and weaken the photorefractive grating.

The holes migrate via hopping along a network of oxidizable charge transport agents. This network can be provided by the donor entity carbazole itself, or by hydrazones such as DEH or arylamines such as tri-tolylamine (TTA) or *N,N*-bis(4-methylphenyl)-*N,N*-bis-(phenyl)-benzidine (TPD). These can be added as dopants, or attached to the polymer backbone of the host polymer, as is the case in PVK. As in the case of inorganic photorefractives, a population of empty shallow traps is required to enable the nonuniform space charge distribution of the photorefractive grating. This is often achieved by providing a population of deep traps for some of the shallow traps to empty into, and it has been shown that the nonlinear optical chromophores can serve this purpose in PVK-based materials. These chromophores serve double duty as the moieties providing the optical nonlinearity. They should have large hyperpolarizability β for electro-optic susceptibility and/or large polarizability anisotropy $\Delta\alpha = \alpha_{\parallel} - \alpha_{\perp}$, where parallel and perpendicular refer to the molecular axis. They should also have a large ground state dipole moment μ_g to enable the molecular orientation effect. These parameters can be combined into a single figure of merit (FOM) defined as

$$\text{FOM} = \frac{1}{M} \left[9\mu_g \beta + \frac{2\mu_g^2 \Delta\alpha}{k_B T} \right] \quad (22)$$

where M is the molar mass, k_B is Boltzmann's constant, and T is the temperature.

The final component in the composite is a plasticizer to lower the glass transition temperature T_g so as to better enable the orientational orientation effect.

An example of a complete composite comprising hole transporter-electron donor/nonlinear chromophore-deep trap/plasticizer/sensitizer-electron acceptor is PVK/AODCST/BBP/ C_{60} in the ratio 49.5:35:15:0.5 percent, where AODCST 2-[[4-[bis(2-methoxyethyl)amino]phenyl]methylene]-malononitrile and BBP is butyl benzyl phthalate. It showed a gain coefficient of 235 cm^{-1} with a response time of 5 ms at 1 Wcm^{-2} for 647-nm light.^{170,172} There are many variations on this theme for the design of photorefractive polymer systems, including the use of sol-gel processing,¹⁷¹ and the use of alternative sensitizers such as gold nanoparticles,¹⁷³ transition metal complexes,^{174–177} and quantum dots. Quantum dots have been investigated as sensitizers;¹⁷⁸ this is attractive since the spectral sensitivity of the system could be tuned through selection of the size of the quantum dots. It is tempting to try to increase the nonlinearity by increasing the proportion of nonlinear chromophore; however, this can lead to phase separation in a composite polymer. This drawback can be overcome by using an organic amorphous glass as photoconductor and NLO molecule simultaneously, or by using fully functionalized polymers in which the charge generator, charge transporter, and NLO components are incorporated as side chains. Liquid crystals have large orientational nonlinearity, and they have been successfully made photorefractive via the addition of small amounts of sensitizer.¹⁷⁹ They have also been combined with photoconductive polymers as polymer-dispersed or polymer-dissolved liquid crystals.^{180,181} Another approach is to replace the transparent electrodes that bias the liquid crystal with thin plates of inorganic photorefractive material such as cerium-doped strontium barium niobate.¹⁸² In this way, the large photorefractive space charge generated in the inorganic plates can extend into the liquid crystal layer and generate a large orientational nonlinearity. This removes the need to tilt the liquid crystal cell and resulted in gain coefficients as large as 1600 cm^{-1} and grating periods as small as 300 nm.

Table 5, reprinted from a review article by Ostroverkhova and Moerner,¹⁸³ shows the characteristics of several organic photorefractive systems. That review provides many further details on modeling, design, and characterization of organic photorefractive materials.

TABLE 5 PR Properties of High Performance Organic Materials in the Visible Part of the Spectrum^a

Composite (conc of Constituents, wt %)	T_g , °C	α , cm ⁻¹	d , μm	λ , nm	Γ , cm ⁻¹ (E, V/ μm)	τ_g^{-1} , s ⁻¹ (J, W/cm ²)	$\eta_{\text{max}}^{\%}$ (E, V/ μm)	τ_{FWM}^{-1} , s ⁻¹ (J, W/cm ²)	Δn , 10 ⁻³ (E, V/ μm)	Refs.
Polymer composites										
PVK/AODCST/BBP/C ₆₀ (49.5/35/15/0.5)		9	80	647	235 (100)	200 (1)				170, 172
PVK/DCDHF-6/BBP/C ₆₀ (49.5/30/20/0.5)		15	80	647	400 (100)	6 (0.1) 4 (0.1)				172
PVK/BDMPAB/TNF (55/44/1)	43		100	633	195 (85)	~1 (0.004)	40 ^{int} (70)		4.2 (92)	184
PVK/6OCB/C ₆₀ (49.8/50/0.2)	47.1		70		210 (65)					185
PSX/DB-IP-DC/TNF (69/30/1)	27.5	60	100	633	390 (100)	30 (0.04)	92 ^{int} (30)		3 (30)	186, 187
PSX/DMNPAA/TNF (53/46/1)	25		670	670	221 (80)			0.2(1.2)	5.8 (80)	188
PSX/stilbene A/TNF (51/48/1)	25		40	670	53 (100)		100 ^{int} , 60 ^{ext} (70)	0.017 (1.2)	10.5 (100)	188
DBOP-PPV/DMNPAA/ MNPAA/DPP/PCBM (52/20/20/5/3)	14.4	34	105	633			90 ^{int} (62)	1.7 (0.305)	2.6 (62)	189
<i>p</i> -PMEH-PPV/DO3/DPP/C ₆₀ (74/5/20/1)	45			633	403 (0 ^b)	0.003 (0.28)				190
PPT-Cz/DDCST/C ₆₀ (64.5/35/0.5)	-7	36.6	100	633	250 (60)		93 ^{int} (100)	0.37 (0.034)	1 (50)	191
PTCB/DHADC-MPN/DIP/ TNFM (49.7/37.6/12.5/0.18)		22.6	105	633	225 (50)		71 ^{ext} (28)	0.07 (0.78)	8.5 (50)	192
Amorphous glasses										
2BNCM/PMMA/TNF (90/9/7/0.3)	22	4	150	676	69 (40)		80 (40)	0.012 (1)	10 (40)	193
DCDHF-6/C ₆₀ (99.5/0.5)	19	12.7	70	676	240 (30)	0.6 (0.1)		0.41 (0.8)		194
DCDHF-6-CF3/C ₆₀ (99.5/0.5)	17	19.9	70	676	255 (40)	0.116 (0.1)		0.21 (0.8)		194
EHCN/TNF (99/1)	25	41	100	633	84 (40)		90 ^{int} (30)	0.67 (0.121)	1.3 (30)	195
Cz-C6-THDC/ECZ/TNF (89/10/1) ^c	33		50				65 (70)		4.5 (70)	196
Methine A	6	1.64	130	633	118 (89)		74 ^{int} (53)		5.6 (53)	197

Fully functional polymers						
Ru-FPP	130	102	690	380 (0 ^b)	0.0014	174
Polymer-dispersed liquid crystals and liquid crystals						
PMMA/TL202/ECZ/TNEM (42/40/17/1)		99 ^d	105	633	136 (10)	100 ^{int} (8) 3.2 (22)
PMMA/TL202/ECZ/ICdS (42/40/16/2)		7.5	129	514.5	30 (31)	56 ^{ext} (22) 72 ^{ext} , 90 ^{int} (50)
SCLP/E7/C ₆₀ (50/49.95/0.05) E7 on PVK/TNF (83/17) ^e		<20 20	50 10	633 514.5	640 (0.7) 3700 ^f (0.9)	~0.1 (3) 44 ^{ext} , f (0.9) ~100
Hybrid organic-inorganic composites, glasses, and sol-gels						
PIBM/AZPON (40/60)	113		55	633	350 (35)	80 ^{int} (13)
BEPON	24	49	40	633	750 (100)	40 ^{int} (28) ^c
PVK/DCVDEA/TNF/Au (70/28.6/1.4/1)	40		70	633	240 (130)	43 (90)
Sol-gel DMHINAB-urethane- SiO1.5/ SiO1.5OH/ECZ/TNF (1:1:1:0.2:0.002) ^g		29	30	633	444 (0 ^b)	25.6 ^{int} , h (0 ^b) 0.0056
Sol-gel SG-Cz/SG-MN/PEG/ TNF (45/45/9/1)	2		75	633	55 (94)	82.4 ^{ext} (94) 0.59 (0.14)

^aColumns represent: (1) composition (concentration of the constituents in wt %, unless stated otherwise); (2) glass transition temperature T_g ; (3) absorption coefficient α ; (4) sample thickness d ; (5) wavelength of the PR experiments λ ; (6) 2BC gain coefficient, measured with p -polarized writing beams, Γ (electric field E , at which the indicated Γ was obtained); (7) PR response time τ^{-1} , obtained from fits to 2BC dynamics (total light intensity of writing beams, at which the indicated value of τ^{-1} was obtained); (8) maximal diffraction efficiency η_{max} , measured with p -polarized probe and s -polarized writing beams. External (η_{ext}) or internal (η_{int}) diffraction efficiency is indicated, where applicable (electric field E , at which the indicated η was obtained); (9) PR speed τ_{FWM} , obtained from fits to either formation or erasure of the PR grating measured in the FWM experiment (total light intensity of either writing beams or erasing beam); (10) refractive index modulation Δn , calculated from the diffraction efficiency (electric field E , at which the indicated refractive index modulation was obtained); (11) reference to work from which the data was taken. All data reported were obtained in Bragg regime (volume grating) at room temperature, unless stated otherwise. ^bPrepooled material. ^cTemperature of the measurements T_m , 30°C. ^dIncludes scattering losses. ^eAll the measurements were done in the Raman-Nath regime. See discussion about the relevance of the gain coefficient in the text. ^fMaximal diffraction efficiency possible in the Raman-Nath regime is ~34 percent. The authors attributed their high diffraction efficiency by nonsinusoidal space-charge field. ^gMolar concentrations. ^hThis value was obtained with p -polarized writing beams and p -polarized readout.

12.3 DEVICES

Real-Time Holography

The real-time phase holograms produced by the photorefractive effect can be used to perform any of the functions of regular holograms. In fact, many of the first applications of photorefractive nonlinear optics were replications of experiments and ideas first introduced in the 1960s in terms of the then new static holography. These include distortion correction by phase conjugation and one-way imaging through distortions, holographic interferometry for nondestructive testing, vibration mode visualization, and pattern recognition by matched filtering. A concise overview of these applications may be found in Goodman's classic text.²⁰² Their photorefractive realizations are well described in Huignard and Gunter.¹ Some of the advantages of photorefractive holography over conventional holography are

- Photorefractive holograms are volume phase holograms. This results in high diffraction efficiencies.
- Photorefractive holograms are self-developing.
- Photorefractive holograms diffract light from the writing beams into each other during the writing process. This gives rise to a dynamic feedback process in which the grating and writing beams influence each other.
- Photorefractive holograms adapt to changing optical fields.

It is the last two features that most strongly differentiate the photorefractive effect from regular holography.


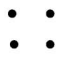
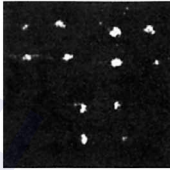


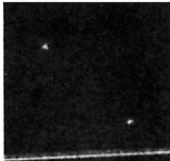

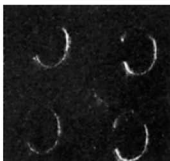
Pattern Recognition As an example of the transfer of applications of conventional holography to photorefractive nonlinear optics we consider the case of pattern recognition by matched filtering.

A matched filter for pattern recognition is simply a Fourier transform hologram of the desired impulse response. It can be made in real time in a photorefractive crystal, a lens being used to produce the Fourier transform. The best-known system is a nonlinear optical triple processor based on four-wave mixing²⁰³ (Fig. 10). For a phase conjugate mirror, the coupled-wave equation for beam 3 which is the phase conjugate of beam 4 pumped by beams 1 and 2 is

$$\begin{aligned} \frac{dA_3}{dz} &= \gamma \frac{(A_1 A_4^* + A_2^* A_3)}{I_0} A_2 \\ &= \frac{\gamma}{I_0} I_2 A_3 + \frac{\gamma}{I_0} A_1 A_2 A_4^* \end{aligned} \quad (23)$$

where the A_j and I_j are the amplitudes and intensities of beams j , respectively, and I_0 is the total intensity of the interacting beams. The first term on the right-hand side simply corresponds to amplification of beam 3. The second term is the source for beam 3. Thus, the amplitude of beam 3 is proportional to the product $A_1 A_2 A_4^*/I_0$. As depicted in Fig. 10, A_1 , A_2 , and A_4 are the Fourier transforms of the spatially varying input fields a_1 , a_2 , and a_4 . The output a_3 is proportional to the inverse Fourier transform of the product of the three Fourier transforms $A_1 A_2 A_4^*/I_0$. If I_0 is spatially constant the output is beam 1 convolved with beam 2 correlated with beam 4 ($a_3 \propto a_1 \otimes a_2 * a_4$) where $*$ represents the spatial correlation operation and \otimes represents the convolution operation, all produced in real time. Modified filters such as phase-only filters can be produced by taking advantage of energy transfer during the filter writing process,^{204–206} or by saturation induced by the presence of the possibly spatially varying intensity denominator I_0 .^{207–209}

Sometimes, applications directly transferred from static holography inspire further developments made possible by exploitation of the physical mechanisms involved in the photorefractive effect. For example, acoustic signals can be temporally correlated with optical signals. An acoustic

U_1	U_2	U_4	U_3
	DELTA FUNCTION		
	DELTA FUNCTION	E	
C	DELTA FUNCTION	CAL TECH	
C		DELTA FUNCTION	

(a)

FIGURE 10 (a) Results demonstrating real-time spatial convolution and correlation of two-dimensional images. The input fields are labeled E_1 , E_2 , and E_p ; the output is labeled E_c . (After White and Yariv.²⁰³)

signal applied to a photorefractive crystal induces piezoelectric fields. If the crystal is illuminated by a temporally varying optical signal, then the photorefractive space charge generated by the photocurrent will be proportional to the product of the time-varying photoconductivity and the time-varying piezoelectric fields. Such correlators can be used to make photorefractive tapped delay lines with tap weights proportional to correlation values.^{210–212} It is also possible to make acoustic filters which detect Bragg-matched retroreflection of acoustic waves from a photorefractive grating in a crystal such as lithium niobate that has low acoustic loss.²¹³

Applications of Photorefractive Gain in Two-Beam Coupling

Coherent Image Amplification Coherent image amplification is especially important for coherent optical processors. Without it, the losses introduced by successive filtering operations would soon become intolerable. Practical considerations include maintenance of signal-to-noise ratio and

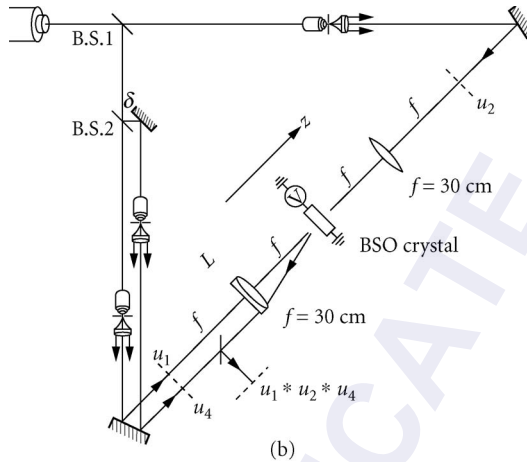


FIGURE 10 (b) Experimental apparatus for performing spatial convolution and correlation using four-wave mixing in photorefractive bismuth silicon oxide. Input and output planes are shown by dashed lines. (After White and Yariv.²⁰³)

amplification fidelity. The main contribution to noise introduced by photorefractive amplifiers is from the fanning effect. Although it can be reduced in a given crystal in a variety of different ways, such as by crystal rotation²¹ and multiwavelength recording,²² by far the best approach to this problem would be to undertake research to grow cleaner crystals. Amplification fidelity is determined by gain uniformity. In the spatial frequency domain, it is limited because photorefractive gain depends on the grating period. Images with a high spatial bandwidth write holograms with a wide range of spatial frequencies and grating periods. Optimal uniformity is obtained for reflection gratings, in which the image beam counterpropagates with respect to the pump. In that case, the change in grating period depends only to second order on the image spatial frequency. Any remaining first-order nonuniformity is due to the angular dependence of the effective electro-optic coefficient. In the space domain, gain uniformity is limited by pump depletion so that the most accurate results will be obtained if the pump beam is strong enough that it is not significantly depleted in the interaction.

Two-beam coupling amplification can also be used for beam cleanup:²¹⁴ a badly distorted beam, say from a laser diode, can be converted to a gaussian beam. A small sample of the beam is split off, spatially filtered and amplified in two-beam coupling by the remaining bulk of the distorted beam. The efficiency of the method can be quite high: fidelity limitations due to spatial variations in gain are usually quite small and can be removed by a second round of spatial filtering. A unidirectional ring resonator with an intracavity spatial filter can also be used for beam cleanup.²¹⁵

Laser Power Combining An application related to two-beam coupling image amplification is coherent power combining, which would be especially useful for semiconductor lasers. The output of a single laser gain stripe is limited to values of the order of a few hundred milliwatts. Some applications require diffraction-limited beams containing many watts produced at high efficiency. Such a source can in principle be made by using two-beam coupling amplification of a diffraction-limited seed by the mutually injection locked outputs of many diode stripes. The injection locking can be achieved by evanescent coupling between laser gain stripes^{216,217} or by retroreflecting a portion of the amplified beam with a partially transmitting mirror.²¹⁸ Another possibility involves forming a double phase conjugate mirror^{219,220} or ring self-pumped phase conjugate mirror with a master laser providing one input, and the light from the gain elements loosely focused into the crystal providing the other inputs. Such a system will be self-aligning and will correct intracavity distortions by phase conjugation.

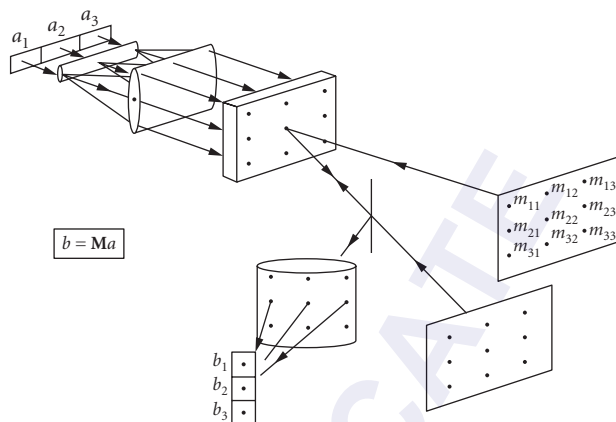


FIGURE 11 Schematic drawing of the basic principle of optical matrix vector multiplication through four-wave mixing in nonlinear media. Light from a linear source array a is fanned out by a cylindrical lens, where it diffracts from an image plane hologram of the matrix \mathbf{M} to produce a set of beams bearing the required products $\mathbf{M}_{ij}a_j$. A second cylindrical lens sums the diffracted beams to form $b = \mathbf{M}a$. (After Yeh and Chiou.²²⁹)

Phase conjugate master oscillator/power amplifiers have also been used with some success.²²¹ Practical problems include the need to control the spectral effects of the associated multiple coupled cavities. Reference²²² gives an excellent exposition of these problems. Also, while self-pumped phase conjugate reflectivities and two-beam coupling efficiencies can theoretically approach 100 percent, in practice these efficiencies rarely exceed 80 percent. Among oxide ferroelectrics, barium titanate exhibits high gain at GaAs laser wavelengths. However, while some bulk semiconductors such as InP:Fe and CdTe:V are sensitive in the 1.3- to 1.5- μm wavelength range of interest for optical fiber communications, high gain requires the application of high electric fields.

Optical Interconnects Use of the double phase conjugate mirror for laser locking suggests another application. The beam-coupling crystal can be viewed as a device that provides optical interconnection of the laser gain elements to each other.²²³ The basic idea exists in the realm of static holography in terms of computer interconnection by holographic optical elements^{224,225} (HOEs). The use of photorefractive crystals should enable the construction of reprogrammable interconnects that would be self-aligning if phase conjugation were used:²²⁶ the laser gain elements in the power combining case can be imagined as the input/output ports of an electronic chip.²²⁷ Another way to go about the interconnection problem is to design in terms of an optical crossbar switch, or matrix vector multiplier.^{228,229} The vector is an array of laser diode sources and the matrix describes connection patterns of the sources to a vector array of detectors (Fig. 11). The interconnection matrix is realized as a photorefractive hologram.

Applications of Photorefractive Loss in Two-Beam Coupling

If the sign of the coupling constant is reversed (for example, by rotation of the uniaxial crystal by 180° so that the direction of the optic axis is reversed) the pump interferes destructively with the signal so that the output is reduced to a very low level. This resulting photorefractive loss can be used in a number of applications such as for optical limiters, optical bistability,²³⁰ and novelty filters and achieved using a variety of different devices such as ring resonators and phenomena such as beam fanning.

Optical Limiters The process of beam depletion in photorefractive materials forms the basis of their use as optical limiters. For this application, another important property of photorefractive materials is their ability to respond selectively to coherent optical inputs. Any portion of the input which is temporally incoherent is transmitted through a photorefractive material in a linear manner. Thus, photorefractive limiters will also selectively attenuate (or *excise*) a coherent beam while transmitting an incoherent beam; these devices have thus also been referred to as excisors.

The first studies of photorefractive limiters were published in 1985.²⁵ A number of device architectures involving resonators and self-pumped phase conjugate mirrors were discussed, but the primary emphasis was on beam fanning. Several features of limiters or excisors using beam fanning were pointed out: (1) the device design is very simple, with a single input beam and no separate external paths required; (2) the limiting mechanism is due to scattering (and not absorption), so that added heating is not present; (3) the device will operate at one or more wavelengths within the bandwidth of its photorefractive response, which (for ferroelectric oxides) extends over the entire visible band; and (4) the device will respond to sources with a relatively small coherence length, including mode-locked lasers producing picosecond pulses.

Experiments at 488 nm using BaTiO₃ in the beam fanning geometry²⁵ produced a steady-state device transmission of 2.5 percent, and a response time of 1.1 s at 1 W/cm². Using the measured response time and intensity, we note that 1.1 J/cm² will pass through the device before it fully activates. In later measurements in the beam fanning geometry, attenuation values exceeding 30 dB and device activation energies as low as 1 to 10 mJ/cm² have been measured.

Two-beam coupling amplification of a second beam produced by beam splitters^{231,232} or gratings in contact with the crystal^{233,234} has also been studied as a mechanism for optical limiting. This mechanism is closely related to fanning, with the only difference that the second beam in a fanning device is produced internally by scattering.

Novelty Filters Novelty filters are devices whose output consists of only the changing part of the input. The photorefractive effect can be used to realize the novelty filter operation in several different ways. The simplest way is to use two-beam coupling for image deamplification as was used in the fanning and two-beam coupling optical limiters. In that case, the pump interferes destructively with the signal so that the output is reduced to a very low level. Now if the signal suddenly changes, the output will be the difference between the new input signal and the reconstruction of the old signal by diffraction of the pump from the old grating. Thus, the output will show the changed parts of the scene until the grating adapts during the photorefractive response time to the new scene.^{235,236} Such interferometers have been used, for example, to map turbulent flow,^{237,238} to make photothermal measurements,²³⁹ and to build acoustic spectrum analyzers.²⁴⁰

Phase Conjugate Interferometry

Another way to produce a novelty filter is to use a phase conjugate interferometer.^{241,242} This is an interferometer in which some or all of the conventional mirrors are replaced by phase conjugate mirrors, thus achieving the benefit of self-alignment. The effects of phase objects inserted in the interferometer are canceled out by phase conjugation. One of the most common realizations is a phase conjugate Michelson interferometer (Figs. 12 and 13). If the phase conjugate mirrors have common pumping beams (this can be achieved by illuminating the same self-pumped phase conjugator with the beams in both arms), the phase of reflection for both beams will be the same and there will be a null at the output from the second port of the device. If a phase object inserted in one of the arms suddenly changes, the null will be disturbed, and the nonzero output will represent the changing parts of the phase object. The nonzero output persists until the gratings in the phase conjugate mirror adapt to the new input fields.²⁴³ If two amplitude objects are inserted in the interferometer, one in each arm, the intensity of the output at the nulling port is the square of the difference between the squared moduli of the objects. This architecture thus gives rise to image subtraction.^{244,245} A slightly modified version can be used to measure thin-film properties (refractive index, absorption coefficient, and thickness): the film under test on its substrate is used as the interferometer beam splitter.²⁴⁶

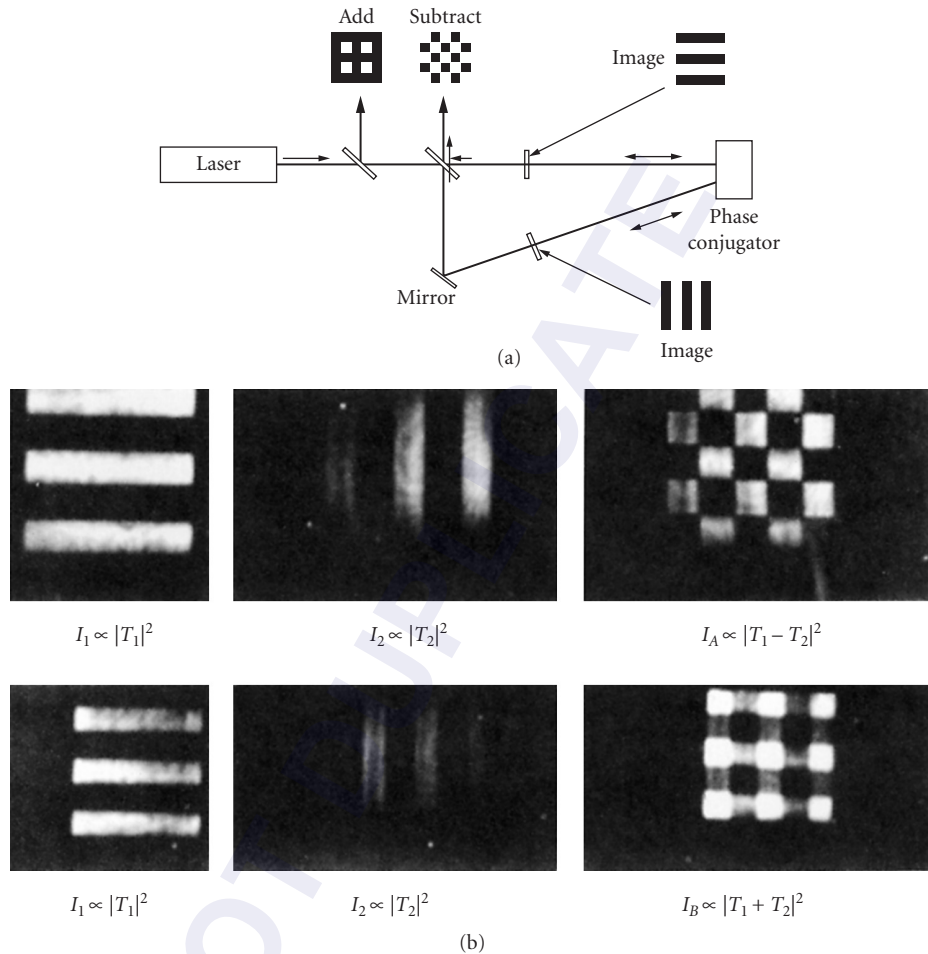


FIGURE 12 Phase conjugate interferometer for image subtraction: (a) amplitude images are placed in the interferometer arms. Their difference appears at the nulling output, the sum appears at the retroreflection output. (b) Real-time image subtraction and addition of images using above apparatus. (After Chiou and Yeh,²⁴⁵)

Associative Memories and Neural Networks

There are a number of adaptive processors that can be designed using photorefractive beam coupling. In addition to novelty filters, these include associative memories, neural network models, and other recursive image processors.

Photorefractive phase conjugate mirrors provide an elegant way to realize linear associative memories in which a fragment of an image can be used to recall the entire image from a bank of multiplexed holograms stored in a long-term storage photorefractive crystal such as lithium niobate.^{247,248} Neural networks, on the other hand, are nets of interconnected signals with nonlinear feedback. They have been extensively investigated in the artificial intelligence community.²⁴⁹ Typical applications are modeling of neural and cognitive systems and the construction of classification machines. As we have seen, there are a number of different ways to realize optical interconnections

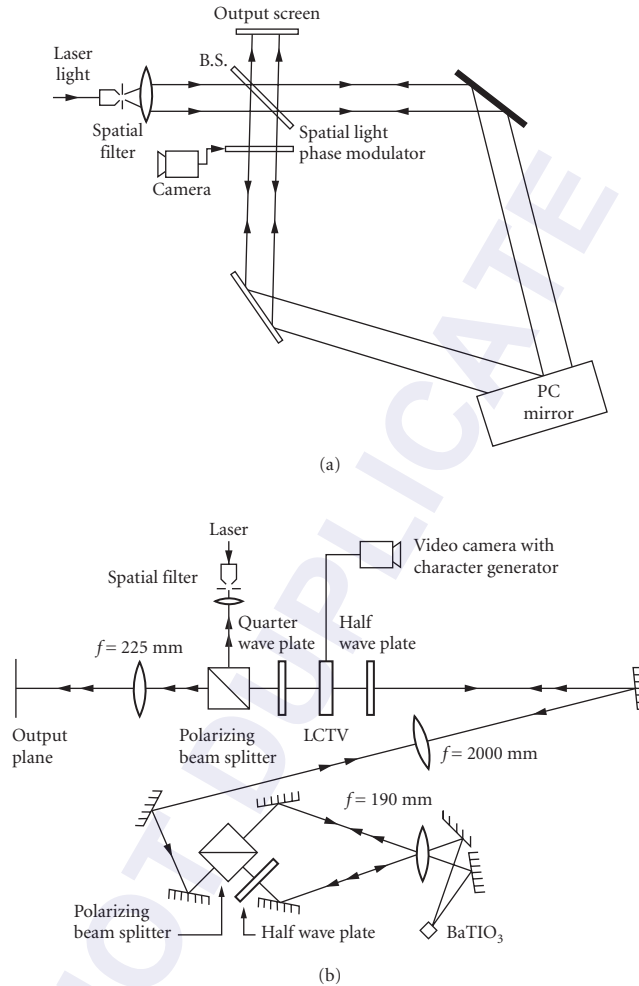


FIGURE 13 Phase conjugate interferometer as a novelty filter: any phase change in the object arm disturbs the null at the output until the phase conjugate mirror adjusts to the change. (a) Optical tracking novelty filter incorporating a spatial phase modulator. BS, beam splitter; PC, phase conjugate mirror. (b) Modification of preceding device to enable the use of a polarization modulating liquid crystal television (LCTV).

using the photorefractive effect, and nonlinear feedback can be introduced by using the nonlinear transfer properties of pumped photorefractive crystals. A pattern classifier can be built by recording a hologram for each image class (e.g., represented by a clear fingerprint) in a photorefractive crystal with a long time constant. Many holograms can be superposed if they are recorded with spatially orthogonal reference beams. As in the case of the linear associative memory, a smudged fingerprint introduced to the system will partially reconstruct each of the reference beams. The brightest reconstruction will be the reference associated with the fingerprint most like the smudged input. An oscillator with internal saturable absorption is built to provide competitive feedback of the reference reconstructions to themselves. The oscillator mode should be the mode associated with the proper fingerprint. That

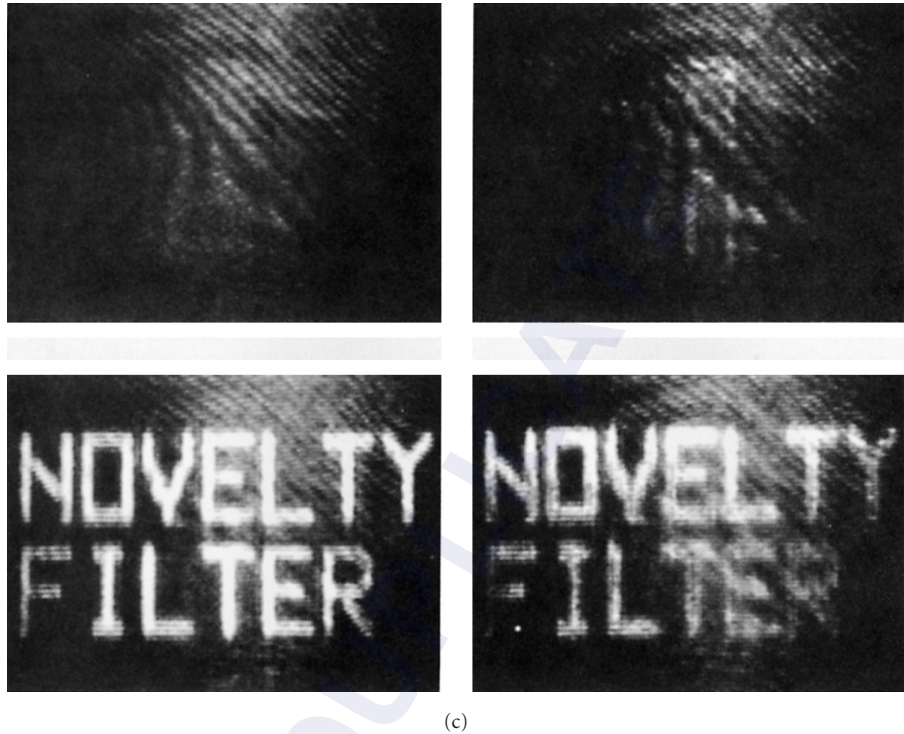


FIGURE 13 (c) Photograph of the output of the tracking novelty filter shown in (b). Input to the LCTV is taken from a character generator driving a video camera: (1) the character generator is off; the interferometer is essentially dark; (2) the character-generator display, showing the phrase NOVELTY FILTER is activated; (3) the filter adapts to the new scene and becomes nearly blank, as in (1). Some letters are visible; (4) the character-generator display is deactivated. The previous phrase appears at the output of the interferometer. Shortly thereafter it fades to (1). (After Anderson *et al.*²⁴³)

mode then reconstructs the clear fingerprint. The output of the device is the stored fingerprint which is most like the smudged input. We have described just one optical neural network model, but just as there are many theoretical neural network models, there are also many optical neural network models.^{250–254} Each of them has its own practical difficulties, including those of reliability, suitability of available threshold functions, and stability.

In addition, the optical gain of photorefractive oscillators makes possible the design of other recursive image processors, for example, to realize Gerschberg-Saxton-type algorithms in phase conjugate resonators.²⁵⁵

Thresholding

In optical data processing it is often necessary to determine if one or more elements of an optical pattern has an intensity above (or below) a set threshold value. For example, in optical associative memory applications, it is necessary to select the stronger modes among many in an optical resonator. In optical correlator applications, the output information plane may be thresholded to determine whether a correlation has been obtained, and to determine the location of the correlation peak(s). A closely related operation which is also useful for these applications is the Max or winner-take-all operation.

The ideal thresholding device should have the following properties: (1) the capability to process complex images with high resolution, (2) high sensitivity, (3) low crosstalk between pixels, (4) a sharp threshold, with constant output for intensities above the threshold value, (5) large signal-to-noise, and (6) the ability to control the threshold level by external means.

A large number of thresholding schemes using the photorefractive effect have been proposed or demonstrated, although only a few experimental demonstrations of thresholding of spatial patterns or images have been reported. Techniques used in early investigations include (1) uniform incoherent erasure of self-pumped phase conjugate mirrors²⁵⁶ and of photorefractive end mirrors in phase conjugate resonators^{257,258} and (2) pump depletion in externally pumped phase conjugate mirrors and double phase conjugate mirrors.⁵⁵ A number of thresholding devices using ring resonators have also been demonstrated.^{55,259} One way of increasing the sharpness of the threshold is to insert additional nonlinear media in the photorefractive resonator cavities.

Ingold et al.²⁶⁰ demonstrated winner-take-all behavior in a nonresonant cavity containing a nematic liquid crystal and a photorefractive crystal. In later experiments,²⁶¹ a thresholding phase conjugate resonator containing a saturable absorber consisting of a thin film of fluorescein-doped boric acid glass demonstrated several performance improvements. In this architecture (shown schematically in Fig. 14) the input image was amplitude-encoded as a two-dimensional array of pixels on an incoherent control beam that was incident on the saturable absorber (a single pixel is shown in Fig. 14). If the intensity at a given pixel was above a threshold intensity, the saturable absorber bleached locally by an amount sufficient to switch on the phase conjugate resonator. The phase conjugate resonator continued to oscillate at these pixel locations even when the control beam was removed. The output was thus bistable and latching.

Photorefractive Holographic Storage

The neural networks described above rely on optical information storage in a material whose gratings are long-lived. One of the earliest potential applications for the photorefractive effect was for

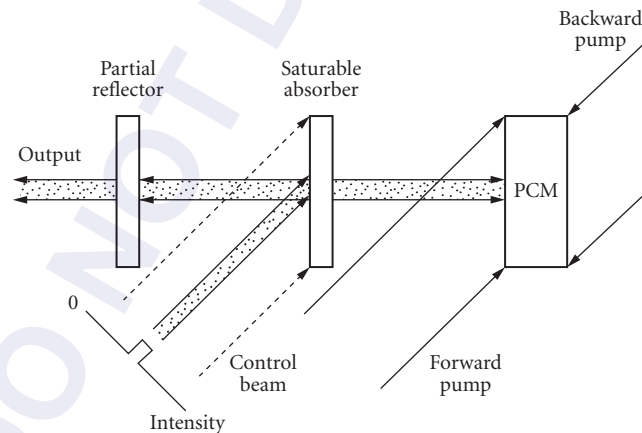


FIGURE 14 Schematic diagram of linear phase conjugate resonator containing an intracavity saturable absorber. Information is read into the resonator by means of a separate control beam incident on the saturable absorber. The control beam can be brought in through a beam splitter, or at an angle to the optic axis (as shown). The fluorescein-doped boric acid glass saturable absorber had a saturation intensity ($I_s = 20 \text{ mW/cm}^2$) at its absorption peak (450 nm). The BaTiO_3 phase conjugate mirror was pumped at the neighboring Ar laser wavelength of 458 nm. The control beam was generated from the same laser, and cross-polarized to a void feedback into the resonator. (After McCahon et al.²⁶¹)

holographic information storage. One of the best materials for the purpose was and still is lithium niobate; it has storage times which can be as long as years. The principal concerns for holographic storage are information density, crosstalk minimization, fixing efficiency, signal-to-noise ratio, and development of practical readout/writing architectures. There has been a recent renewed interest in optical memory design, now that materials and computer technology have improved and schemes for rapid addressing have matured. The two main classes of addressing scheme use spatially orthogonal reference beams^{262–264} and temporally orthogonal reference beams²⁶⁵ (spatial multiplexing vs. frequency multiplexing). Spatial multiplexing has the advantage that relatively simple optical sources can be used. Frequency multiplexing has the advantage that spatial crosstalk is reduced compared to that associated with spatial multiplexing. However, it has the drawback of requiring frequency-tunable laser sources. An important consideration is the need to pack as much information as possible into each individual hologram. Here the frequency multiplexing approach is superior, because the information density can reach much higher values before crosstalk sets in. Crosstalk can also be reduced by storing information in photorefractive fiber bundles instead of bulk, so that the information is more localized.^{266,267}

Holographic Data Storage

A photorefractive holographic digital storage system demonstrated in 1994 had a storage capacity of 163 kB using lithium niobate as the storage medium.²⁶⁸ It used angular multiplexing to record data pages as separate holograms and distributed consecutive bits over multiple pages to reduce the probability of burst errors. The raw bit error rate was between 10^{-3} and 10^{-4} but was improved to 10^{-6} by use of a Hamming error correcting code. Shortly after that, the capacity was increased to 5 MB using Reed-Solomon error correcting codes,²⁶⁹ and in 2000 a 1-GB lithium niobate system was demonstrated with 50 μ s seek time.²⁷⁰ Lithium niobate suffers from a low recording sensitivity, of the order of 0.1 cm/J and it has been largely replaced as a holographic recording medium by photopolymers, whose sensitivity can be several orders of magnitude larger (1000 cm/J).²⁶⁹

Photorefractive Waveguides

There are three essentially different ways to prepare electro-optic (possibly photorefractive) waveguides. One is to produce local alterations in the chemistry of an electro-optic substrate, for example, by titanium in-diffusion into LiNbO₃ or ion implantation in BaTiO₃.²⁷¹ Another way is to grow waveguides in layers by techniques such as RF sputtering, liquid phase epitaxy, laser ablation,²⁷² and metalorganic chemical vapor deposition (MOCVD).²⁷³ A third way is to polish bulk material down to a thin wafer.²⁷⁴

The performance of electro-optic waveguide devices such as couplers and switches can be seriously compromised by the refractive index changes induced by the photorefractive effect in the host electro-optic materials such as lithium niobate. Therefore, one of the main motivations for understanding the photorefractive effect in waveguides is to develop ways to minimize its effects. MgO doping of lithium niobate is commonly used in attempts to reduce the photorefractive effect.²⁷⁵

Some researchers have taken advantage of waveguide photorefractivity. Optical confinement in waveguides enhances the effectiveness of optical nonlinearities. With conventional $\chi^{(3)}$ materials, waveguiding confinement increases the coupling constant-length product by maintaining high intensity over longer distances than would be possible in bulk interactions. In the photorefractive case, optical confinement reduces the response time.

An excellent review of earlier work on photorefractivity in waveguides may be found in Ref. 276. More recently Eason and coworkers have measured a response time improvement by a factor of 100 in an ion-implanted BaTiO₃ waveguide.²⁷¹ A bridge mutually pumped phase conjugator was also demonstrated.²⁷⁷

Photorefractive Solitons

In 1992, the possibility that photorefractive crystal might be able to support spatial solitons was proposed,²⁷⁸ and subsequently demonstrated, first as a transient effect,²⁷⁹ and then in the steady state.²⁸⁰ Optical spatial solitons are beams of light in which the normal tendency toward diffractive spreading is counterbalanced by a nonlinear optical self-focusing effect. Solitons that are stable in 2 transverse directions and the 1 propagation direction (2 + 1 solitons) are made possible by the fact that the photorefractive nonlinearity is saturable, in contrast to the situation with Kerr nonlinearities where self-focusing leads to catastrophic collapse.²⁸¹ Solitons can form in photorefractive crystals, where drift is made to dominate diffusion through the application of a DC electric field (screening solitons) or by use of the photovoltaic effect (photovoltaic solitons) and where the degree of saturation of the nonlinearity is controlled through provision of background illumination.²⁸² These solitons can be described by the nonlinear Schrodinger equation with saturable nonlinearity:²⁸³

$$i \frac{\partial u}{\partial z} + \frac{1}{2} \frac{\partial^2 u}{\partial x^2} - \frac{u}{|u|^2 + 1} = 0 \quad (24)$$

where u is the amplitude of the soliton normalized by the square root of the sum of the background and dark intensities, x is the transverse coordinate, and z is the propagation direction. More exact versions of this equation have been used, for example, including the diffusion component of the photorefractive effect,²⁸⁴ but Eq. (24) does describe the basic features.

Since the initial demonstration of basic screening solitons, many different types of solitons have been found, providing a very active research area both in fundamental physics and potential applications. These include bright and dark solitons,²⁸⁵ vortex solitons,²⁸⁶ vector solitons,²⁸⁷ incoherent,²⁸⁸ and white light solitons.²⁸⁹ One of the main research interests has been in the study of soliton collisions and interactions.²⁹⁰ Applications have included optically induced waveguides for optical routing,²⁹¹ and for beam confinement for nonlinear frequency conversion.²⁹²

12.4 REFERENCES

1. Gunter, P. and J.-P. Huignard (eds), *Photorefractive Materials and Their Applications 1, 2, and 3*. Springer-Verlag, New York, 2005, 2006, 2007.
2. A. M. Glass, D. von der Linde, and T. J. Negran, *Appl. Phys. Lett.* **25**:233 (1974).
3. V. L. Vinetskii and N. V. Khukhtarev, *Sov. Phys. Solid State* **16**:2414 (1975).
4. F. P. Strohkendl, J. M. C. Jonathan, and R. W. Hellwarth, *Opt. Lett.* **11**:312 (1986).
5. N. V. Kukhtarev, G. E. Dovgalenko, and V. N. Starkov, *Appl. Phys. A* **33**:227 (1984).
6. G. C. Valley, *Appl. Opt.* **22**:3160 (1983).
7. N. V. Khukhtarev, V. B. Markov, S. G. Odoulov, M. S. Soskin, and V. L. Vinetskii, *Ferroelectrics* **22**:949, 961 (1979).
8. B. Belinicher and B. Sturman, *Usp. Fiz. Nauk.* **130**:415 (1980).
9. G. C. Valley, *IEEE J. Quantum Electron.* **QE19**:1637 (1983).
10. J. P. Huignard and J. P. Herriau, *Appl. Opt.* **24**:4285 (1985).
11. S. I. Stepanov and M. P. Petrov, *Sov. Tech. Phys. Lett.* **10**:572 (1984).
12. M. Ziari, W. H. Steier, P. M. Ranon, M. B. Klein, and S. Trivedi, *J. Opt. Soc. Am. B* **9**:1461 (1992).
13. B. Fischer, M. Cronin-Golomb, J. O. White, and A. Yariv, *Opt. Lett.* **6**:519 (1981).
14. D. W. Vahey, *J. Appl. Phys.* **46**:3510 (1975).
15. J. Feinberg and K. R. MacDonald, "Phase Conjugate Mirrors and Resonators with Photorefractive Materials," in P. Gunter and J. P. Huignard (eds.), *Photorefractive Materials and Their Applications II, Survey of Applications*, Springer Verlag, Berlin, 1989.

16. M. D. Ewbank, R. A. Vasquez, R. S. Cudney, G. D. Bacher, and J. Feinberg, Paper FS1, *Technical Digest, 1990 OSA Annual Meeting*, Boston, MA, Nov. 4–9 (1990).
17. R. A. Vasquez, R. R. Neurgaonkar, and M. D. Ewbank, *J. Opt. Soc. Am. B* **89**:1416 (1992).
18. V. V. Voronov, I. R. Dorosh, Yu. S. Kuz'minov, and N. V. Tkachenko, *Sov. J. Quantum Electron.* **10**:1346 (1980).
19. J. Feinberg, *J. Opt. Soc. Am.* **72**:46 (1982).
20. M. Segev, Y. Ophir, and B. Fischer, *Opt. Commun.* **77**:265 (1990).
21. H. Rajbenbach, A. Delboulbe, and J. P. Huignard, *Opt. Lett.* **14**:1275 (1989).
22. W. S. Rabinovich, B. J. Feldman, and G. C. Gilbreath, *Opt. Lett.* **16**:1147 (1991).
23. J. Feinberg, *Opt. Lett.* **7**:486 (1982).
24. A. Zozulya, *Sov. J. Quantum Electron.* **22**:677 (1992).
25. M. Cronin-Golomb and A. Yariv, *J. Appl. Phys.* **57**:4906 (1985).
26. J. O. White, M. Cronin-Golomb, B. Fischer, and A. Yariv, *Appl. Phys. Lett.* **40**:450 (1982).
27. S. K. Kwong, A. Yariv, M. Cronin-Golomb, and I. Ury, *Appl. Phys. Lett.* **47**:460 (1985).
28. M. D. Ewbank and P. Yeh, *Opt. Lett.* **10**:496 (1985).
29. D. Anderson and R. Saxena, *J. Opt. Soc. Am. B* **4**:164 (1987).
30. M. Cronin-Golomb, B. Fischer, J. O. White, and A. Yariv, *IEEE J. Quantum Electron.* **QE20**:12 (1984).
31. A. A. Zozulya and V. T. Tikonchuk, *Sov. J. Quantum Electron.* **18**:981 (1988).
32. D. A. Fish, A. K. Powell, and T. J. Hall, *Opt. Commun.* **88**:281 (1992).
33. W. Krolikowski, K. D. Shaw, and M. Cronin-Golomb, *J. Opt. Soc. Am. B* **6**:1828 (1989).
34. W. Krolikowski, M. Belic, M. Cronin-Golomb, and A. Bledowski, *J. Opt. Soc. Am. B* **7**:1204 (1990).
35. R. A. Rupp and F. W. Drees, *Appl. Phys. B* **39**:223 (1986).
36. D. A. Temple and C. Warde, *J. Opt. Soc. Am. B* **3**:337 (1986).
37. M. D. Ewbank, P. Yeh, and J. Feinberg, *Opt. Commun.* **59**:423 (1986).
38. E. M. Avakyan, K. G. Belabaev, V. Kh. Sarkisov, and K. M. Tumanyan, *Sov. Phys. Solid State* **25**:1887 (1983).
39. L. Holtmann, E. Kratzig, and S. Odoulov, *Appl. Phys. B* **53**:1 (1991).
40. M. Cronin-Golomb, B. Fischer, J. O. White, and A. Yariv, *Appl. Phys. Lett.* **41**:689 (1982).
41. M. Cronin-Golomb, B. Fischer, J. O. White, and A. Yariv, *Appl. Phys. Lett.* **42**:919 (1983).
42. V. A. D'yakov, S. A. Korol'kov, A. Mamaev, V. V. Shkunov, and A. A. Zozulya, *Opt. Lett.* **16**:1614 (1991).
43. S. Weiss, S. Sternklar, and B. Fischer, *Opt. Lett.* **12**:114 (1987).
44. M. D. Ewbank, R. A. Vasquez, R. R. Neurgaonkar, and J. Feinberg, *J. Opt. Soc. Am. B* **7**:2306 (1990).
45. M. D. Ewbank, *Opt. Lett.* **13**:47 (1988).
46. D. Wang, Z. Zhang, Y. Zhu, S. Zhang, and P. Ye, *Opt. Commun.* **73**:495 (1989).
47. A. M. C. Smout and R. W. Eason, *Opt. Lett.* **12**:498 (1987).
48. B. Fischer, S. Sternklar, and S. Weiss, *IEEE J. Quantum Electron.* **QE25**:550 (1989).
49. M. Cronin-Golomb, *Opt. Lett.* **15**:897 (1990).
50. A. A. Zozulya, *Opt. Lett.* **16**:2042 (1991).
51. K. D. Shaw, *Opt. Commun.* **94**:458 (1992).
52. T. Y. Chang and R. W. Hellwarth, *Opt. Lett.* **10**:408 (1985).
53. R. A. Mullen, D. J. Vickers, L. West, and D. M. Pepper, *J. Opt. Soc. Am. B* **9**:1726 (1992).
54. G. C. Valley, *J. Opt. Soc. Am. B* **9**:1440 (1992).
55. M. Horowitz, D. Kliger, and B. Fischer, *J. Opt. Soc. Am. B* **8**:2204 (1991).
56. W. Krolikowski, K. D. Shaw, M. Cronin-Golomb, and A. Bledowski, *J. Opt. Soc. Am. B* **6**:1828 (1989).
57. W. Krolikowski and M. Cronin-Golomb, *Appl. Phys. B* **52**:150 (1991).
58. K. D. Shaw, *Opt. Commun.* **97**:148 (1993).
59. D. J. Gauthier, P. Narum, and R. W. Boyd, *Phys. Rev. Lett.* **58**:1640 (1987).
60. W. Krolikowski, M. Belic, M. Cronin-Golomb, and A. Bledowski, *J. Opt. Soc. Am. B* **7**:1204 (1990).

61. F. Vacchs and P. Yeh, *J. Opt. Soc. Am. B* **6**:1834 (1989).
62. W. Krolikowski and M. Cronin-Golomb, *Opt. Commun.* **89**:88 (1992).
63. J. Hong, A. E. Chiou, and P. Yeh, *Appl. Opt.* **29**:3027 (1990).
64. M. Cronin-Golomb, *Opt. Commun.* **89**:276 (1992).
65. G. C. Valley and G. J. Dunning, *Opt. Lett.* **9**:513 (1984).
66. R. Blumrich, T. Kobialka, and T. Tschudi, *J. Opt. Soc. Am. B* **7**:2299 (1990).
67. S. R. Liu and G. Indebetouw, *J. Opt. Soc. Am. B* **9**:1507 (1992).
68. F. Laeri, T. Tschudi, and J. Libers, *Opt. Commun.* **47**:387 (1983).
69. J. Feinberg and R. W. Hellwarth, *Opt. Lett.* **5**:519 (1980).
70. P. Gunter, *Phys. Reports* **93**:200 (1983).
71. D. von der Linde and A. M. Glass, *Appl. Phys.* **8**:85 (1975).
72. G. C. Valley and M. B. Klein, *Opt. Eng.* **22**:704 (1983).
73. A. Ashkin, G. D. Boyd, J. M. Dziedzic, R. G. Smith, A. A. Ballman, J. J. Levinstein, and K. Nassau, *Appl. Phys. Lett.* **9**:72 (1966).
74. F. S. Chen, *J. Appl. Phys.* **38**:3418 (1967).
75. F. S. Chen, J. T. LaMacchia, and D. B. Fraser, *Appl. Phys. Lett.* **3**:213 (1968).
76. M. B. Klein and G. C. Valley, *J. Appl. Phys.* **57**:4901 (1985).
77. M. Cronin-Golomb, K. Y. Lau, and A. Yariv, *Appl. Phys. Lett.* **47**:567 (1985).
78. R. Orlowski and E. Kratzig, *Solid St. Comm.* **27**:1351 (1978).
79. J. J. Amodei and D. L. Staebler, *Appl. Phys. Lett.* **18**:540 (1971).
80. R. L. Townsend and J. T. LaMacchia, *J. Appl. Phys.* **41**:5188 (1970).
81. R. L. Townsend and J. T. LaMacchia, *J. Appl. Phys.* **41**:5188 (1970).
82. J. Feinberg, D. Heiman, A. R. Tanguay, Jr., and R. W. Hellwarth, *J. Appl. Phys.* **51**:1297 (1980).
83. E. Kratzig, F. Welz, R. Orlowski, V. Doorman, and M. Rosenkranz, *Solid St. Comm.* **34**:817 (1980).
84. M. B. Klein, "Photorefractive Properties of BaTiO₃," in P. Gunter and J.-P. Huignard (eds.), *Photorefractive Materials and Their Applications*, Springer-Verlag, Berlin, 1988.
85. V. Belruss, J. Kalnajs, A. Linz, and R. C. Folweiler, *Mater. Res. Bull.* **6**:899 (1971).
86. M. B. Klein and R. N. Schwartz, *J. Opt. Soc. Am. B* **3**:293 (1986).
87. P. G. Schunemann, T. M. Pollak, Y. Yang, Y. Y. Teng, and C. Wong, *J. Opt. Soc. Am. B* **5**:1702 (1988).
88. B. A. Wechsler and M. B. Klein, *J. Opt. Soc. Am. B* **5**:1713 (1988).
89. D. Rytz, R. R. Stephens, B. A. Wechsler, M. S. Keirstad, and T. M. Baer, *Opt. Lett.* **15**:1279 (1990).
90. G. W. Ross, P. Hribek, R. W. Eason, M. H. Garrett, and D. Rytz, *Opt. Commun.* **101**:60 (1993); B. A. Wechsler, M. B. Klein, C. C. Nelson, and R. N. Schwartz, *Opt. Lett.* **19**, April 15 (1994).
91. S. Ducharme and J. Feinberg, *J. Appl. Phys.* **56**:839 (1984).
92. R. A. Rupp, A. Maillard, and J. Walter, *Appl. Phys.* **A49**:259 (1989).
93. L. Holtmann, *Phys. Status Solidi A* **113**:K89 (1989).
94. D. Mahgerefteh and J. Feinberg, *Phys. Rev. Lett.* **64**:2195 (1990).
95. G. A. Brost and R. A. Motes, *Opt. Lett.* **15**:1194 (1990).
96. L. K. Lam, T. Y. Chang, J. Feinberg, and R. W. Hellwarth, *Opt. Lett.* **6**:475 (1981).
97. N. Barry and M. J. Damzen, *J. Opt. Soc. B* **9**:1488 (1992).
98. A. L. Smirl, G. C. Valley, R. A. Mullen, K. Bohnert, C. D. Mire, and T. F. Boggess, *Opt. Lett.* **12**:501 (1987).
99. G. A. Brost, R. A. Motes, and J. R. Rotge, *J. Opt. Soc. Am. B* **5**:1879 (1988).
100. P. Yeh, *Appl. Opt.* **26**:602 (1987); A. M. Glass, M. B. Klein, and G. C. Valley, *Appl. Opt.* **26**:3189 (1987).
101. E. Voit, M. Z. Zha, P. Amrein, and P. Gunter, *Appl. Phys. Lett.* **51**:2079 (1987).
102. D. Rytz, M. B. Klein, R. A. Mullen, R. N. Schwartz, G. C. Valley, and B. A. Wechsler, *Appl. Phys. Lett.* **52**:1759 (1988).

103. S. Ducharme and J. Feinberg, *J. Opt. Soc. Am. B* **3**:283 (1986).
104. M. H. Garrett, J. Y. Chang, H. P. Jenssen, and C. Warde, *Opt. Lett.* **17**:103 (1992).
105. M. E. Lines and A. M. Glass, *Principles And Applications of Ferroelectrics And Related Materials*, Clarendon Press, Oxford, 1977, pp. 280–292.
106. K. Megumi, N. Nagatsuma, Y. Kashiwada, and Y. Furuhashi, *J. Matls. Sci.* **11**:1583 (1976).
107. J. B. Thaxter, *Appl. Phys. Lett.* **15**:210 (1969).
108. K. Megumi, H. Kozuka, M. Kobayashi, and Y. Furuhashi, *Appl. Phys. Lett.* **30**:631 (1977).
109. V. V. Voronov, I. R. Dorosh, Y. S. Kuzminov, and N. V. Tkachenko, *Sov. J. Quantum Electron.* **10**:1346 (1980).
110. R. R. Neurgaonkar, W. K. Cory, J. R. Oliver, M. D. Ewbank, and W. F. Hall, *Opt. Eng.* **26**:392 (1987).
111. G. L. Wood and R. R. Neurgaonkar, *Opt. Lett.* **17**:94 (1992).
112. R. A. Vasquez, F. R. Vachss, R. R. Neurgaonkar, and M. D. Ewbank, *J. Opt. Soc. Am. B* **8**:1932 (1991).
113. R. A. Vasquez, R. R. Neurgaonkar, and M. D. Ewbank, *J. Opt. Soc. Am. B* **9**:1416 (1992).
114. M. D. Ewbank, R. R. Neurgaonkar, W. K. Cory, and J. Feinberg, *J. Appl. Phys.* **62**:374 (1987).
115. K. Sayano, A. Yariv, and R. R. Neurgaonkar, *Opt. Lett.* **15**:9 (1990).
116. L. Hesselink and S. Redfield, *Opt. Lett.* **13**:877 (1988).
117. A. A. Grabar, M. Jazbinsek, A. N. Shumelyuk, Y. M. Vysochanskii, G. Montemezzani, and P. Gunter, in *Photorefractive Materials and Their Applications 2: Materials*, P. Gunter and J. -P. Huignard, (eds.) Springer, New York, 2006, Chap. 10, pp. 327–362.
118. A. A. Grabar, R. I. Muzhikash, A. D. Kostyuk, and Y. M. Vysochanskii, *Fizika Tverdogo Tela* **33**:2335 (1991).
119. M. Jazbinsek, D. Haertle, G. Montemezzani, P. Gunter, A. A. Grabar, I. M. Stoika, and Y. M. Vysochanskii, *J. Opt. Soc. Am. B Opt. Phys.* **22**:2459 (2005).
120. A. Shumelyuk, S. Odoulov, O. Oleynik, G. Brost, and A. Grabar, *Appl. Phys. B-Lasers and Optics* **88**:79 (2007).
121. R. Mosimann, P. Marty, T. Bach, F. Juvalta, M. Jazbinsek, P. Gunter, and A. A. Grabar, *Opt. Lett.* **32**:3230 (2007).
122. M. Jazbinsek, D. Haertle, G. Montemezzani, P. Gunter, A. A. Grabar, I. M. Stoika, and Y. M. Vysochanskii, *J. Opt. Soc. Am. B-Opt. Phys.* **22**:2459 (2005).
123. B. A. Horowitz and F. J. Corbitt, *Opt. Eng.* **17**:353 (1978).
124. J. P. Huignard and F. Micheron, *Appl. Phys. Lett.* **29**:591 (1976).
125. J. P. Wilde, L. Hesselink, S. W. McCahon, M. B. Klein, D. Rytz, and B. A. Wechsler, *J. Appl. Phys.* **67**:2245 (1990).
126. S. I. Stepanov and M. P. Petrov, *Opt. Commun.* **52**:292 (1985).
127. S. L. Sochava, S. I. Stepanov, and M. P. Petrov, *Sov. Tech. Phys. Lett.* **13**:274 (1987).
128. M. P. Petrov, S. L. Sochava, and M. P. Petrov, *Opt. Lett.* **14**:284 (1989).
129. S. I. Stepanov, V. V. Kulikov, and M. P. Petrov, *Opt. Commun.* **44**:19 (1982).
130. B. Imbert, H. Rajbenbach, S. Mallick, J. P. Herriau, and J.-P. Huignard, *Opt. Lett.* **13**:327 (1988).
131. J. E. Millerd, E. M. Garmire, M. B. Klein, B. A. Wechsler, F. P. Strohkendl, and G. A. Brost, *J. Opt. Soc. Am. B* **9**:1449 (1992).
132. R. Baltrameyunas, Yu. Vaitkus, D. Veletskas, and I. Kapturauskas, *Sov. Tech. Phys. Lett.* **7**:155 (1981).
133. A. M. Glass, A. M. Johnson, D. H. Olson, W. Simpson, and A. A. Ballman, *Appl. Phys. Lett.* **44**:948 (1984).
134. M. B. Klein, *Opt. Lett.* **9**:350 (1984).
135. A. Partovi, J. Millerd, E. M. Garmire, M. Ziari, W. H. Steier, S. B. Trivedi, and M. B. Klein, *Appl. Phys. Lett.* **57**:846 (1990).
136. K. Kuroda, Y. Okazaki, T. Shimura, H. Okimura, M. Chihara, M. Itoh, and I. Ogura, *Opt. Lett.* **15**:1197 (1990).
137. M. Ziari, W. H. Steier, P. M. Ranon, S. Trivedi, and M. B. Klein, *Appl. Phys. Lett.* **60**:1052 (1992).
138. A. M. Glass and J. Strait, "The Photorefractive Effect in Semiconductors," in P. Gunter and J.-P. Huignard (eds.), *Photorefractive Materials and Their Applications I*, Springer-Verlag, Berlin, 1989, vol. 61, pp. 237–262.
139. P. Tayebati, J. Kumar, and S. Scott, *Appl. Phys. Lett.* **59**:3366 (1991).
140. J. Kumar, G. Albanese, and W. H. Steier, *J. Opt. Soc. Am. B* **4**:1079 (1987).
141. B. Imbert, H. Rajbenbach, S. Mallick, J. P. Herriau, and J.-P. Huignard, *Opt. Lett.* **13**:327 (1988).

142. M. B. Klein, S. W. McCahon, T. F. Boggess, and G. C. Valley, *J. Opt. Soc. Am. B* **5**:2467 (1988).
143. G. C. Valley, H. Rajbenbach, and H. J. von Bardeleben, *Appl. Phys. Lett.* **56**:364 (1990).
144. Ph. Refregier, L. Solymar, H. Rajbenbach, and J.-P. Huignard, *J. Appl. Phys.* **58**:45 (1985).
145. K. Walsh, A. K. Powell, C. Stace, and T. J. Hall, *J. Opt. Soc. Am. B* **7**:288 (1990).
146. M. Ziari, W. H. Steier, P. M. Ranon, M. B. Klein, and S. Trivedi, *J. Opt. Soc. Am. B* **9**:1461 (1992).
147. P. Gravey, G. Picoli, and J. Y. Labandibar, *Opt. Commun.* **70**:190 (1989).
148. G. Picoli, P. Gravey, C. Ozkul, and V. Vieux, *J. Appl. Phys.* **66**:3798 (1989).
149. A. Partovi, A. Kost, E. M. Garmire, G. C. Valley, and M. B. Klein, *Appl. Phys. Lett.* **56**:1089 (1990).
150. J. E. Millerd, S. D. Koehler, E. M. Garmire, A. Partovi, A. M. Glass, and M. B. Klein, *Appl. Phys. Lett.* **57**:2776 (1990).
151. J. E. Millerd, E. M. Garmire, and M. B. Klein, *Opt. Lett.* **17**:100 (1992).
152. G. C. Valley, S. W. McCahon, and M. B. Klein, *J. Appl. Phys.* **64**:6684 (1988).
153. A. Partovi, E. M. Garmire, G. C. Valley, and M. B. Klein, *Appl. Phys. Lett.* **55**:2701 (1989).
154. R. B. Bylisma, D. H. Olson, and A. M. Glass, *Appl. Phys. Lett.* **52**:1083 (1988).
155. D. M. Pepper, J. AuYeung, D. Fekete, and A. Yariv, *Opt. Lett.* **3**:7 (1978).
156. L. Pichon and J.-P. Huignard, *Opt. Commun.* **36**:277 (1981).
157. A. M. Glass, D. D. Nolte, D. H. Olson, G. E. Doran, D. S. Chemla, and W. H. Knox, *Opt. Lett.* **15**:264 (1990).
158. D. D. Nolte, D. H. Olson, G. E. Doran, W. H. Knox, and A. M. Glass, *J. Opt. Soc. Am. B* **7**:2217 (1990).
159. Q. N. Wang, R. M. Brubaker, D. D. Nolte, and M. R. Melloch, *J. Opt. Soc. Am. B* **9**:1626 (1992).
160. A. Partovi, A. M. Glass, D. H. Olson, G. J. Zydzik, K. T. Short, R. D. Feldman, and R. F. Austin, *Opt. Lett.* **17**:655 (1992).
161. A. Partovi, A. M. Glass, D. H. Olson, G. J. Zydzik, H. M. O'Bryan, T. H. Chiu, and W. H. Knox, *Appl. Phys. Lett.* **62**:464 (1993).
162. K. Sutter and P. Gunter, *J. Opt. Soc. Am. B-Opt. Phys.* **7**:2274 (1990).
163. K. Sutter, J. Hulliger, and P. Gunter, *Solid St. Comm.* **74**:867 (1990).
164. K. Sutter, J. Hulliger, R. Schlessler, and P. Gunter, *Opt. Lett.* **18**:778 (1993).
165. S. Ducharme, J. C. Scott, R. J. Twieg, and W. E. Moerner, *Phys. Rev. Lett.* **66**:1846 (1991).
166. W. E. Moerner, C. Walsh, J. C. Scott, S. Ducharme, D. M. Burland, G. C. Bjorklund, and R. J. Twieg, *Nonlinear Optical Properties of Organic Materials IV* **1560**:278 (1991).
167. W. E. Moerner, S. M. Silence, F. Hache, and G. C. Bjorklund, *J. Opt. Soc. Am. B-Opt. Phys.* **11**:320 (1994).
168. B. Kippelen, S. R. Marder, E. Hendrickx, J. L. Maldonado, G. Guillemet, B. L. Volodin, D. D. Steele, et al., *Science* **279**, 54 (1998).
169. F. Wurthner, S. Yao, J. Schilling, R. Wortmann, M. Redi-Abshiro, E. Mecher, F. Gallego-Gomez, and K. Meerholz, *J. Am. Chem. Soc.* **123**:2810 (2001).
170. D. Wright, M. A. Diaz-Garcia, J. D. Casperson, M. DeClue, W. E. Moerner, and R. J. Twieg, *Appl. Phys. Lett.* **73**:1490 (1998).
171. J. D. Wright and N. A. J. M. Sommerdijk, *Sol-Gel Materials: Chemistry and Applications*, Gordon and Breach, Amsterdam, 2001.
172. D. Wright, U. Gubler, Y. Roh, W. E. Moerner, M. He, and R. J. Twieg, *Appl. Phys. Lett.* **79**:4274 (2001).
173. F. Wang, Z. J. Chen, B. Zhang, Q. H. Gong, K. W. Wu, X. S. Wang, B. W. Zhang, and F. Q. Tang, *Appl. Phys. Lett.* **75**:3243 (1999).
174. Z. H. Peng, A. R. Gharavi, and L. P. Yu, *J. Am. Chem. Soc.* **119**:4622 (1997).
175. Q. Wang, L. M. Wang, J. J. Yu, and L. P. Yu, *Adv. Mater.* **12**:974 (2000).
176. W. You, L. M. Wang, Q. Wang, and L. P. Yu, *Macromolecules* **35**:4636 (2002).
177. I. Aiello, D. Dattilo, M. Ghedini, A. Bruno, R. Termine, and A. Golemme, *Adv. Mater.* **14**:1233 (2002).
178. J. G. Winiarz and P. N. Prasad, *Opt. Lett.* **27**:1330 (2002).
179. G. P. Wiederrecht, B. A. Yoon, and M. R. Wasielewski, *Science* **270**:1794 (1995).
180. A. Golemme, B. Kippelen, and N. Peyghambarian, *Appl. Phys. Lett.* **73**:2408 (1998).

181. R. Termine and A. Golemme, *Opt. Lett.* **26**:1001 (2001).
182. G. Cook, C. A. Wyres, M. J. Deer, and D. C. Jones, *Proc. SPIE* **5213**:63 (2003).
183. O. Ostroverkhova and W. E. Moerner, *Chem. Rev.* **104**:3267 (2004).
184. Z. J. Chen, F. Wang, Z. W. Huang, Q. H. Gong, Y. W. Chen, Z. J. Zhang, and H. Y. Chen, *J. Phys. D-Appl. Phys.* **31**:2245 (1998).
185. Y. W. Bai, X. F. Chen, X. H. Wan, Q. F. Zhou, H. Liu, B. Zhang, and Q. H. Gong, *Appl. Phys. Lett.* **80**:10 (2002).
186. H. Chun, I. K. Moon, D. H. Shin, S. Song, and N. Kim, *J. Mater. Chem.* **12**:858 (2002).
187. W. J. Joo, N. J. Kim, H. Chun, I. K. Moon, and N. Kim, *Polymer* **42**:9863 (2001).
188. S. Schlöter, U. Hofmann, P. Strohrriegl, H. W. Schmidt, and D. Haarer, *J. Opt. Soc. Am. B-Opt. Phys.* **15**:2473 (1998).
189. E. H. Mecher, C. Brauchle, H. H. Horhold, J. C. Hummelen, and K. Meerholz, *Phys. Chem. Chem. Phys.* **1**:1749 (1999).
190. D. J. Suh, O. O. Park, T. Ahn, and H. K. Shim, *Jpn. J. Appl. Phys. Part 2-Letters* **41**:L428 (2002).
191. O. P. Kwon, S. H. Lee, G. Montemezzani, and P. Gunter, *Adv. Function. Mater.* **13**:434 (2003).
192. E. Hendrickx, J. Herlocker, J. L. Maldonado, S. R. Marder, B. Kippelen, A. Persoons, and N. Peyghambarian, *Appl. Phys. Lett.* **72**:1679 (1998).
193. P. M. Lundquist, R. Wortmann, C. Geletneky, R. J. Twieg, M. Jurich, V. Y. Lee, C. R. Moylan, and D. M. Burland, *Science* **274**:1182 (1996).
194. O. Ostroverkhova, D. Wright, U. Gubler, W. E. Moerner, M. He, A. Sastre-Santos, and R. J. Twieg, *Adv. Function. Mater.* **12**:621 (2002).
195. J. Sohn, J. Hwang, S. Y. Park, and G. J. Lee, *Jpn. J. Appl. Phys. Part 1-Regular Papers Short Notes & Review Papers* **40**:3301 (2001).
196. H. Chun, N. J. Kim, W. J. Joo, J. W. Han, C. H. Oh, and N. Kim, *Synth. Met.* **129**:281 (2002).
197. L. M. Wang, M. K. Ng, and L. P. Yu, *Appl. Phys. Lett.* **78**:700 (2001).
198. H. Ono and N. Kawatsuki, *J. Appl. Phys.* **85**:2482 (1999).
199. S. Bartkiewicz, K. Matczyszyn, A. Miniewicz, and F. Kajzar, *Opt. Commun.* **187**:257 (2001).
200. P. Cheben, F. del Monte, D. J. Worsfold, D. J. Carlsson, C. P. Grover, and J. D. Mackenzie, *Nature* **408**:64 (2000).
201. D. H. Choi, H. T. Hong, W. G. Jun, and K. Y. Oh, *Opt. Mater.* **21**:373 (2003).
202. J. W. Goodman, *Introduction to Fourier Optics*, 3rd ed, Roberts, Greenwood Village, CO.
203. J. O. White and A. Yariv, *Appl. Phys. Lett.* **37**:5 (1980).
204. J. Joseph, K. Singh, and P. K. C. Pillai, *Opt. Commun.* **85**:389 (1991).
205. T. Y. Chang, J. H. Hong, S. Campbell, and P. Yeh, *Opt. Lett.* **17**:1694 (1992).
206. J. Khoury, T. C. Fu, M. Cronin-Golomb, and C. Woods, *J. Opt. Soc. Am. B* **11**:1960 (1994).
207. J. P. Huignard and J. P. Herriau, *Appl. Opt.* **17**:2671 (1978).
208. J. Feinberg, *Opt. Lett.* **5**:330 (1980).
209. E. Ochoa, J. W. Goodman, and L. Hesselink, *Opt. Lett.* **10**:430 (1985).
210. J. J. Berg, J. N. Lee, M. W. Casseday, and B. J. Udelson, *Opt. Eng.* **19**:359 (1980).
211. H. Lee and D. Psaltis, *Opt. Lett.* **12**:459 (1987).
212. R. M. Montgomery and M. R. Lange, *Appl. Opt.* **30**:2844 (1991).
213. D. E. Oates, P. G. Gottschalk, and P. B. Wright, *Appl. Phys. Lett.* **46**:1125 (1985).
214. A. E. T. Chiou and P. Yeh, *Opt. Lett.* **10**:621 (1985).
215. S. K. Kwong and A. Yariv, *Appl. Phys. Lett.* **48**:564 (1986).
216. S. MacCormack and R. W. Eason, *Opt. Lett.* **15**:1212 (1990).
217. S. MacCormack and R. W. Eason, *J. Appl. Phys.* **67**:7160 (1990).
218. W. R. Christian, P. H. Beckwith, and I. McMichael, *Opt. Lett.* **14**:81 (1989).

219. M. Segev, S. Weiss, and B. Fischer, *Appl. Phys. Lett.* **50**:1397 (1987).
220. S. Weiss, M. Segev, and B. Fischer, *IEEE J. Quantum Electron.* **QE24**:706 (1988).
221. R. R. Stephens, R. C. Lind, and C. R. Guiliano, *Appl. Phys. Lett.* **50**:647 (1987).
222. P. D. Hillman and M. Marciniak, *J. Appl. Phys.* **66**:5731 (1989).
223. S. Weiss, M. Segev, S. Sternklar, and B. Fischer, *Appl. Opt.* **27**:3422 (1988).
224. J. W. Goodman, F. I. Leonberger, S. Y. Kung, and R. A. Athale, *Proc. IEEE* **72**:850 (1984).
225. J. W. Goodman, in *Optical Processing and Computing*, H. H. Arsenault, T. Szoplik, and B. Macukow (eds.), Academic Press, San Diego, 1989, chap. 1.
226. K. Wagner and D. Psaltis, *Appl. Opt.* **26**:5061 (1987).
227. M. Cronin-Golomb, *Appl. Phys. Lett.* **23**:2189 (1989).
228. A. E. T. Chiou and P. Yeh, *Appl. Opt.* **31**:5536 (1992).
229. P. Yeh and A. E. T. Chiou, *Opt. Lett.* **12**:138 (1987).
230. D. M. Lininger, P. J. Martin, and D. Z. Anderson, *Opt. Lett.* **14**:697 (1989).
231. S. W. McCahon and M. B. Klein, *Proc. SPIE* **1105**:119 (1989).
232. G. L. Wood, W. W. Clark III, G. J. Salamo, A. Mott, and E. J. Sharp, *J. Appl. Phys.* **71**:37 (1992).
233. M. B. Klein and G. J. Dunning, *Proc. SPIE* **1692**:73 (1992).
234. J. L. Schultz, G. J. Salamo, E. J. Sharp, G. L. Wood, R. J. Anderson, and R. R. Neurgaonkar, *Proc. SPIE* **1692**:78 (1992).
235. J. E. Ford, Y. Fainman, and S. H. Lee, *Opt. Lett.* **13**:856 (1988).
236. M. Cronin-Golomb, A. M. Biernacki, C. Lin, and H. Kong, *Opt. Lett.* **12**:1029 (1987).
237. G. F. Albrecht, H. F. Robey, and T. R. Moore, *Appl. Phys. Lett.* **57**:864 (1990).
238. H. F. Robey, *Phys. Rev. Lett.* **65**:1360 (1990).
239. S. D. Kalaskar and S. E. Bialkowski, *Anal. Chem.* **64**:1824 (1992).
240. G. Zhou, L. Bintz, and D. Z. Anderson, *Appl. Opt.* **31**:1740 (1992).
241. J. Feinberg, *Opt. Lett.* **8**:569 (1983).
242. M. D. Ewbank, P. Yeh, M. Khoshnevisan, and J. Feinberg, *Opt. Lett.* **10**:282 (1985).
243. D. Z. Anderson, D. M. Lininger, and J. Feinberg, *Opt. Lett.* **12**:123 (1987).
244. S. K. Kwong, G. A. Rakuljic, and A. Yariv, *Appl. Phys. Lett.* **48**:201 (1986).
245. A. E. Chiou and P. Yeh, *Opt. Lett.* **11**:306 (1986); J. Feinberg and K. R. MacDonald, in *Photorefractive Materials and their Applications II*, Springer Verlag, Berlin (1989).
246. E. Parshall and M. Cronin-Golomb, *Appl. Opt.* **30**:5090 (1991).
247. G. J. Dunning, E. Marom, Y. Owechko, and B. H. Softer, *Opt. Lett.* **12**:346 (1987).
248. Y. Owechko, *IEEE J. Quantum Electron.* **QE25**:619 (1989).
249. T. Kohonen, *Self-Organization and Associative Memory*, Springer-Verlag, New York, 1984.
250. D. Z. Anderson, *Opt. Lett.* **11**:56 (1986).
251. J. Hong, S. Campbell, and P. Yeh, *Appl. Opt.* **29**:3019 (1990).
252. D. Psaltis, D. Brady, and K. Wagner, *Appl. Opt.* **27**:1752 (1988).
253. A. V. Huynh, J. F. Walkup, and T. F. Krile, *Opt. Eng.* **31**:979 (1992).
254. E. G. Paek, P. F. Liao, and H. Gharavi, *Opt. Eng.* **31**:986 (1992).
255. K. P. Lo and G. Indebetouw, *Appl. Opt.* **31**:1745 (1992).
256. M. Cronin-Golomb and A. Yariv, *Proc. SPIE* **700**:301 (1986).
257. M. B. Klein, G. J. Dunning, G. C. Valley, R. C. Lind, and T. R. O'Meara, *Opt. Lett.* **11**:575 (1986).
258. S. W. McCahon, G. J. Dunning, K. W. Kirby, G. C. Valley, and M. B. Klein, *Opt. Lett.* **17**:517 (1992).
259. D. M. Lininger, P. J. Martin, and D. Z. Anderson, *Opt. Lett.* **14**:697 (1989).
260. M. Ingold, P. Gunter, and M. Schadt, *J. Opt. Soc. Am. B* **7**:2380 (1990).
261. S. W. McCahon, G. J. Dunning, K. W. Kirby, G. C. Valley, and M. B. Klein, *Opt. Lett.* **17**:517 (1992).

262. C. Denz, G. Paulat, G. Roosen, and T. Tschudi, *Opt. Commun.* **85**:171 (1991).
263. F. Mok, M. C. Tackitt, and H. M. Stoll, *Opt. Lett.* **16**:605 (1991).
264. Y. Taketomi, J. E. Ford, H. Sasaki, J. Ma, Y. Fainman, and S. H. Lee, *Opt. Lett.* **16**:1774 (1991).
265. G. A. Rakuljic, V. Levya, and A. Yariv, *Opt. Lett.* **17**:1471 (1992).
266. L. Hesselink and S. Redfield, *Opt. Lett.* **13**:877 (1988).
267. F. Ito, K. Kitayama, and H. Oguri, *J. Opt. Soc. Am. B* **9**:1432 (1992).
268. J. F. Heanue, M. C. Bashaw, and L. Hesselink, *Science* **265**:749 (1994).
269. L. Hesselink, S. S. Orlov, and M. C. Bashaw, *Proc IEEE* **92**:1231 (2004).
270. J. A. Ma, T. Chang, S. Choi, and J. Hong, *Opt. Quant. Electron.* **32**:383 (2000).
271. K. E. Youden, S. W. James, R. W. Eason, P. J. Chandler, L. Zhang, and P. D. Townsend, *Opt. Lett.* **17**:1509 (1992).
272. K. E. Youden, R. W. Eason, M. C. Gower, and N. A. Vainos, *Appl. Phys. Lett.* **59**:1929 (1991).
273. Y. Nagao, H. Sakata, and Y. Mimura, *Appl. Opt.* **31**:3966 (1992).
274. B. Fischer and M. Segev, *Appl. Phys. Lett.* **54**:684 (1989).
275. J. L. Jackel, *Electron. Lett.* **21**:509 (1985).
276. V. E. Wood, P. J. Cressman, R. J. Holman, and C. M. Verber, in P. Gunter and J. P. Huignard (eds), *Photorefractive Materials and their Applications*, Springer-Verlag, Berlin, 1989.
277. S. W. James, K. E. Youden, P. M. Jeffrey, R. W. Eason, P. J. Chandler, L. Zhang, and P. D. Townsend, *Appl. Opt.* **32**:5299 (1993).
278. M. Segev, B. Crosignani, A. Yariv, and B. Fischer, *Phys. Rev. Lett.* **68**:923 (1992).
279. G. C. Duree, J. L. Shultz, G. J. Salamo, M. Segev, A. Yariv, B. Crosignani, P. Diporto, E. J. Sharp, and R. R. Neurgaonkar, *Phys. Rev. Lett.* **71**:533 (1993).
280. M. D. I. Castillo, P. A. M. Aguilar, J. J. Sanchezmondragon, S. Stepanov, and V. Vysloukh, *Appl. Phys. Lett.* **64**:408 (1994).
281. P. L. Kelley, *Phys. Rev. Lett.* **15**:1005 (1965).
282. K. Kos, H. X. Meng, G. Salamo, M. F. Shih, M. Segev, and G. C. Valley, *Phys. Rev. E* **53**:R4330 (1996).
283. W. Krolkowski, B. Luther-Davies, and C. Denz., *IEEE J. Quant. Electron.* **39**:3 (2003).
284. A. A. Zozulya and D. Z. Anderson, *Phys. Rev. A* **51**:1520 (1995).
285. Z. G. Chen, M. Mitchell, M. F. Shih, M. Segev, M. H. Garrett, and G. C. Valley, *Opt. Lett.* **21**:629 (1996).
286. Z. G. Che, M. Segev, D. W. Wilson, R. E. Muller, and P. D. Maker, *Phys. Rev. Lett.* **78**:2948 (1997).
287. M. I. Carvalho, S. R. Singh, D. N. Christodoulides, and R. I. Joseph, *Phys. Rev. E* **53**:R53 (1996).
288. M. Mitchell, M. Segev, T. H. Coskun, and D. N. Christodoulides, *Phys. Rev. Lett.* **79**:4990 (1997).
289. M. Mitchell and M. Segev, *Nature* **387**:880 (1997).
290. G. I. Stegeman and M. Segev, *Science* **286**:1518 (1999).
291. M. Tiemann, J. Schmidt, V. M. Petrov, J. Petter, and T. Tschudi, *Appl. Opt.* **46**:2683 (2007).
292. S. Lan, M. F. Shih, G. Mizell, J. A. Giordmaine, Z. G. Chen, C. Anastassiou, J. Martin, and M. Segev, *Opt. Lett.* **24**:1145 (1999).

12.5 FURTHER READING

Books and Review Articles on Photorefractive Materials, Effects, and Devices

- Gunter, P., "Holography, Coherent Light Amplification and Phase Conjugation in Photorefractive Materials," *Phys. Rep.* **93**:199 (1982).
- Gunter, P., "Photorefractive Materials," in *CRC Handbook of Laser Science and Technology*, vol. IV, part 2, CRC Press Boca Raton, 1986.

- Gunter, P. and J.-P. Huignard (eds.), *Photorefractive Materials and Their Applications I and II*, Springer Science, Berlin, 1988, 1989.
- Gunter, P. and J.-P. Huignard (eds), *Photorefractive Materials and Their Applications 1, 2, and 3*. Springer-Verlag, New York, 2005, 2006, 2007.
- Klein, M. B. and G. C. Valley, "Optimal Properties of Photorefractive Materials for Optical Data Processing," *Opt. Engin.* **22**:704 (1983).
- Odoulov, S., M. Soskin, and A. Khizniak, *Optical Oscillators with Degenerate Four-Wave Mixing (Dynamic Grating Lasers)*, Harwood Academic Publishers, London, 1989.
- Petrov, M. P., S. I. Stepanov, and A. V. Khomenko, *Photorefractive Crystals in Coherent Optical Systems*, Springer-Verlag, Berlin, 1991.
- Solymar, L., D. J. Webb, and A. Grunnet-Jepsen, *The Physics and Applications of Photorefractive Materials* (Oxford Series in Optical and Imaging Sciences) 1996.
- Yeh, P., *Introduction to Photorefractive Nonlinear Optics*, Wiley, New York, 1993.
- Yu, F. and S. Yin (eds), *Photorefractive Optics*, Academic Press, San Diego, 2000.

OPTICAL LIMITING

David J. Hagan

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

13.1 INTRODUCTION

As the name implies, an optical limiter is a device designed to keep the power, irradiance, energy, or fluence transmitted by an optical system below some specified maximum value, regardless of the magnitude of the input. It must do this while maintaining high transmittance at low input powers. The most important application of such a device is the protection of sensitive optical sensors and components from laser damage. There are many other potential applications for such devices, including laser power regulation or stabilization, or restoration of signal levels in optical data transmission or logic systems, but this chapter will primarily concentrate on devices for sensor protection.

Perhaps the most obvious way of achieving optical limiting is by active control, where input light levels are monitored by a sensor, which through some processor activates a modulator or shutter that in turn limits the transmitted light. The best-known examples of these are the iris and blink response of the eye. However, these are limited in speed to about 0.1 second, so that any intense pulse of light shorter than this can get past these defense mechanisms and damage the retina before they can respond. Speed is an issue with most active control systems for optical limiting. To protect sensors, the transmittance must be reduced in a time much shorter than the width of the potentially damaging pulse. Even very fast electro-optic shutters are limited to rise times on the order of 1 ns, which may be insufficient to adequately block Q-switched pulses shorter than 50 ns or so in duration. Even for protection against longer pulses, cost and complexity are disadvantages of active optical limiting systems. Another direct way to protect sensors against high-power lasers is to use narrow-line spectral filters. These can work well when the laser wavelength is known, such as in laboratory laser safety goggles, but would be largely ineffective against tunable lasers.

Passive systems, on the other hand, use a nonlinear optical material that functions as a combined sensor, processor, and modulator. This offers the potential for high speed, simplicity, compactness, and low cost. However, passive systems place severe requirements on the nonlinear medium.¹⁻³ While many materials exhibit the type of effects that produce optical limiting, usually these effects prove not to be large enough. Because of this, all prototype passive systems demonstrated to date place the nonlinear optical component in or near a focal plane. In a focal plane, the energy density of a beam from a distant laser source is $10^5 \sim 10^8$ times greater than in a pupil plane. Even in this focused geometry, material nonlinearities are barely large enough, and systems that adequately protect eyes and other common sensors over a broad wavelength band have yet to be demonstrated, at least in the

visible and near-infrared. The idea of using a nonlinear material in a pupil plane (e.g., a coating on the surface of goggles) is therefore far from reality. For mid-infrared (3 to 12 μm wavelength range), optical nonlinearities are typically much larger than in the visible and results have been more promising than in the visible.⁴⁻⁶ Still, however, limiting elements must be placed near a focal plane. Hence, research to date on optical limiting has predominantly focused on the search for new or modified materials with stronger nonlinearities, and on how to optimally use the best available nonlinear materials. This chapter will concentrate on passive devices.

The response of an *ideal* optical limiter is shown in Fig. 1, along with some typical responses of passive limiters. Clearly, an optical sensor requires high linear transmittance, T_L , at low input light levels for the transmission of images. Meanwhile, for higher inputs the limiter must clamp the transmitted energy below some maximum value, E_{max} , up to the maximum energy the limiter can withstand, E_D . This is usually the energy damage threshold for the limiting material itself. Usually the minimum transmittance of the device, T_{min} , occurs at this energy. Often, the performance of a limiting system or device is characterized by some type of figure of merit (FOM).^{7,8} One of the most commonly used is $\text{FOM} = T_L/T_{\text{min}}$, which states that a large linear transmittance combined with a low minimum transmittance is desirable.⁸ A slightly different way of expressing this is in terms of the optical density (OD), defined as $\text{OD} = -\log_{10}(T)$, so that the FOM may be reexpressed as the change in OD, or $\Delta\text{OD} = \log_{10}(T_L/T_{\text{min}})$. For the ideal limiter shown in Fig. 1, the FOM is equivalent to the *dynamic range*, which is defined as $\text{D.R.} = E_D/E_L$.⁷ However, although such merit figures are of some use, it is usually necessary to separately specify parameters such as linear transmittance, maximum transmitted energy, and damage energy. For example, in some applications, a linear transmittance of >50 percent could be an absolute requirement that cannot be offset by an improvement in

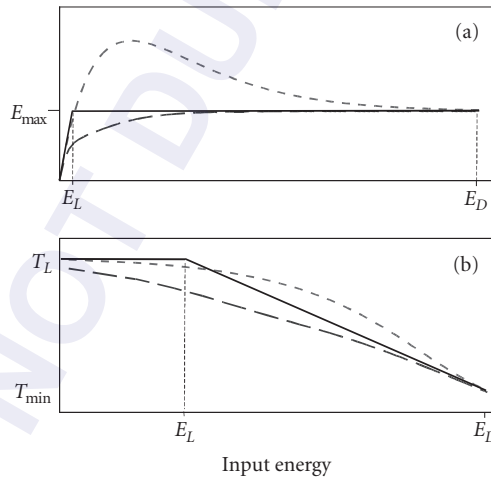


FIGURE 1 Typical limiting curves, drawn as (a) transmitted energy versus input energy, and (b) transmittance versus input energy, on a log-log scale. The solid line is the *ideal* optical limiter response, while the short- and long-dashed lines are typical of real systems. E_D is the energy at which the limiter undergoes irreversible laser damage. E_{max} is the maximum transmitted energy, here measured at the maximum input energy E_D . E_L is the limiting threshold for the ideal limiter, T_L is the linear transmittance, and T_{min} is the minimum transmittance, usually obtained for the maximum input energy, E_D .

protection. In addition, for nonideal limiting responses, the FOM may not give a clear indication of whether a device provides adequate protection at all input energies. For example, the dashed curves in Fig. 1b look very similar when plotted as transmittance versus input energy. Although they have the same FOM, when plotted as transmitted versus input energy, it is clear that the long-dashed curve provides considerably better protection than the short-dashed curve.

The maximum permissible transmitted energy is highly dependent on the threat laser wavelength and pulse width, on the sensor to be protected, and on the f -number (defined as the ratio of focal length to lens diameter) of the imaging system. Most practical imaging systems use an f -number of 5 or less. However, one very common sensor for which we can specify maximum safe exposure levels is the human eye.⁹ Mostly, we are concerned with retinal damage. Visible or near-infrared radiation is not absorbed in the cornea or lens of the eye, and the focusing of light onto the retina produces an optical gain on the order of 10^4 [i.e., the fluence (incident energy/unit area) at the retina is $\sim 10^4$ times that incident on the cornea]. Hence for visible or near-infrared radiation, damage will always occur first at the retina. For wavelengths in the UV and further into the infrared, light does not reach the retina, as it is absorbed in the lens or cornea where the fluence is much lower. If necessary, damage may be avoided by use of optical elements that simply block those wavelengths by reflection or absorption, as the eye cannot detect those wavelengths anyway. As shown in Fig. 2, the ANSI standard for the maximum safe energy entering the eye for pulse lengths less than $17\ \mu\text{s}$ is $0.2\ \mu\text{J}$.¹⁰ However, larger energies may be tolerated where there is a finite probability of a retinal lesion but little chance of retinal hemorrhaging. For example, the ED-50 exposure level, for which there is a 50 percent chance of a retinal lesion but little chance of permanent damage, corresponds to $\sim 1\ \mu\text{J}$ in the visible. Therefore, ideally one would desire $E_{\text{max}} \ll 1\ \mu\text{J}$ for an eye-protection limiter. However, since no practical prototype limiter so far has come close to this value, a more common target value for E_{max} in recent literature has been $\sim 1\ \mu\text{J}$.^{11,12} As will be subsequently described, the total energy entering the eye is not a complete measure of performance, as many nonlinear optical mechanisms that give rise to limiting strongly distort a laser beam as well as controlling its total transmitted energy. Therefore, a better measure of limiting performance is the focusable component of the energy entering the eye, E_{foc} . E_{foc} is defined as the energy falling within a 1.5-mrad-diameter circle in the retinal plane.¹³ The accepted value for the minimum resolution of the eye is 1.5 mrad. Should the limiter defocus the beam enough that the focused beam significantly exceeds 1.5 mrad, an E_{max} of $> 1\ \mu\text{J}$ may be tolerable as long as $E_{\text{foc}} < 1\ \mu\text{J}$.

In the remainder of this chapter, we describe some of the fundamental principles of passive optical limiting, including nonlinear mechanisms and optimization of geometry. We also present a few examples of experimental demonstrations of some types of optical limiters, although this is by no means intended

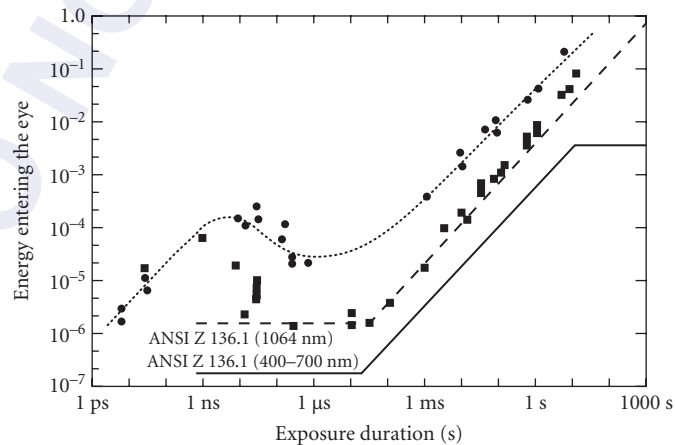


FIGURE 2 Damage thresholds for the eye. (After Ref. 10.)

to be a comprehensive review. Rather, this chapter is intended as a starting point for newcomers to the field of optical limiting. For more detailed surveys of published research in this field, the reader is referred to the review papers,^{1,3,12} journal special issues,¹⁴ and conference proceedings^{15–22} on the subject.

13.2 BASIC PRINCIPLES OF PASSIVE OPTICAL LIMITING

By way of an introduction to passive optical limiting, we briefly describe five of the most common nonlinear mechanisms used. As summarized in Fig. 3, these are (a) nonlinear absorption, (b) nonlinear refraction, (c) nonlinear scattering, (d) photorefraction, and (e) optically induced phase transitions. There have been many other schemes proposed for passive optical limiting, but those mentioned here form the basis for the vast majority of practical limiting devices that have appeared in the literature. A common theme to all schemes is that they each require the nonlinear optical material to be placed in or near a focal plane. Here we concentrate on how each nonlinear optical property results in limiting and we avoid detailed descriptions of the nonlinear mechanisms. For a more complete description of nonlinear optical phenomena and mechanisms, the reader is referred to Chap. 16, “Third-Order Optical Nonlinearities,” of this volume.

Limiting via Nonlinear Absorption

Perhaps the most obvious and direct way to produce passive optical limiting is via nonlinear absorption (NLA), where we require the absorption to increase with increasing incident pulse fluence or irradiance. This can occur in a number of ways, as illustrated in Fig. 4, which shows some of the possible optical transitions for a generic material system. These could represent electronic transitions in many different material types (for example, an organic molecule or a semiconductor crystal).

Two-photon absorption (2PA) is a third-order nonlinear optical process that involves the simultaneous absorption of two photons.²³ For 2PA, the absorption increases in proportion to the incident irradiance, I . Another possibility is a two-step absorption process, where linear absorption populates excited states, from which a second absorption to a higher-lying energy state is possible.^{24,25} If the

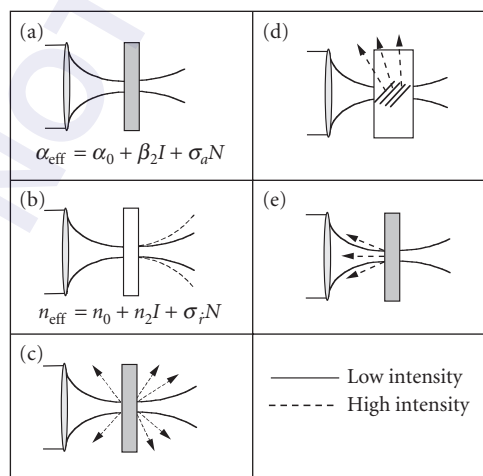


FIGURE 3 Fundamental mechanisms for passive optical limiting: (a) nonlinear absorption; (b) nonlinear refraction; (c) nonlinear scattering; (d) photorefraction; and (e) optically induced phase transitions.

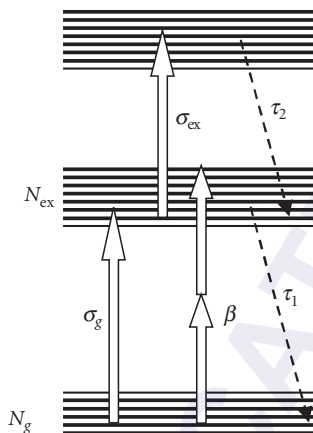


FIGURE 4 Nonlinear absorption processes.

excited state cross section, σ_{ex} , exceeds the ground state cross section, σ_g , then the absorption will increase with increasing excited state population, N_{ex} , and hence with increasing incident fluence. This is usually referred to as *excited state absorption* (ESA) or *reverse-saturable absorption* (RSA). The latter nomenclature grew out of the more commonly observed *saturable absorption*, where $\sigma_{\text{ex}} < \sigma_g$ and the absorption decreases with increasing fluence. In materials with suitable energy levels, it is also possible to populate the excited state by two-photon absorption, and still produce excited state absorption at the excitation wavelength.^{23,26} For any of these cases, we may write an approximate effective absorption coefficient, α_{eff}

$$\alpha_{\text{eff}} = \alpha + \beta I + \sigma_{\text{ex}} N_{\text{ex}} = \sigma_g N_g + \beta I + \sigma_{\text{ex}} N_{\text{ex}} \quad (1)$$

where β is the 2PA coefficient and N_g and N_{ex} are the ground state and excited state absorption cross sections, respectively. For a laser pulse shorter than the excited state lifetime, τ_1 , and for low excitation levels, N_{ex} is directly proportional to the incident fluence. However, an optical limiter is required to work under high levels of excitation, so usually rate equations must be solved to determine the overall transmittance. Nevertheless, it is clear that a large ratio of $\sigma_{\text{ex}}/\sigma_g$ is desirable, as we want large ESA, but small linear absorption. However, σ_g cannot be too small, as N_{ex} must become large enough to produce a strong limiting effect. The minimum achievable transmittance should occur when all molecules have been promoted to the first excited state, so that the transmittance is $T_{\text{min}} = \exp(-\sigma_{\text{ex}}NL)$, where N is the total molecular density and L is the material thickness. Hence the maximum achievable FOM is $T_{\text{min}}/T_L = \exp[-(\sigma_{\text{ex}} - \sigma_g)NL]$.²⁷ However, it is not practical to expect such a physical situation. Before such a high excitation is reached, other effects, including ionization, heating, and so forth will occur.¹¹

In theory, 2PA is ideal for optical limiting, as the linear absorption is zero. Moreover, it can be shown that there is an absolute upper limit to the irradiance that can propagate through a two-photon absorber, given by $I_{\text{max}} = (\beta L)^{-1}$. 2PA also may populate excited states without the inconvenience of linear absorption, so that one can obtain limiting due to both 2PA and ESA. However, it is difficult to find materials with sufficiently large and broadband 2PA coefficients to work well with nanosecond or longer pulses.^{26,28,29}

Regardless of the excitation mechanism, it is desirable to have an excited state lifetime, τ_1 , longer than the laser pulse width, so that each electron or molecule need only be excited one time per pulse. A short upper excited state lifetime, τ_2 , is required to reduce saturation of the excited state absorption, which detracts from limiting performance.

Optimization of NLA Limiters While NLA gives the best optical limiting by placing the nonlinear material in a focal plane, this is also where the damage energy threshold of the nonlinear material is lowest. Very often, this gives an unacceptably low FOM. The damage threshold may be greatly increased by placing another nonlinear absorber in front of the one at focus, hence protecting it from damage. In the case of 2PA, this can be achieved by using a *thick* 2PA material, as illustrated in Fig. 5a.^{7,28,30} Here, the term *thick* means that the material thickness is much greater than the depth of focus of the beam. The front surface is far from focus, and the 2PA away from focus protects the material near focus. Theoretically, this can be done with no reduction in T_L , as there need be no linear absorption. In reality, parasitic linear transmittance and scatter may reduce T_L . For RSA materials, the intrinsic linear absorption does not permit us to use an arbitrarily thick medium. Instead, discrete elements can be used. This geometry is usually referred to a *cascaded* or *tandem* limiter.³¹ In the simplest case, this can consist of two or more elements in tandem, as illustrated in Fig. 5b. The elements can be positioned so that the damage fluence, F_d , is reached simultaneously by all elements. This gives the maximum damage energy.^{8,11,27} It has been shown that the total FOM of the limiter is given by the product of the FOMs of each individual element. However, because of beam distortion due to NLA or to any NLR that may be present in the material, the FOM of an individual element in a cascaded geometry does not usually approach the FOM of the same element when used alone. Miles¹¹ pointed out that this geometry helps keep the fluence high through the length of the limiter, by balancing the decrease in fluence due to absorption with the increase in fluence due to focusing. This is illustrated in Fig. 5b, which shows a sketch of on-axis fluence versus distance for a four-element tandem limiter. This can be understood by considering the example of a limiter with $T_{\min} = 10^{-4}$. If such a limiter were to have only a single element, the fluence on the front surface of the cell would have to exceed that on the rear surface by 10^4 . If damage to the front of the cell were to be avoided, the fluence on the rear surface would be so low that the molecules near the rear surface could not contribute significantly to the limiting. Therefore, these molecules only serve to reduce the linear transmittance. However, for a 4-cell tandem limiter, $T_{\min} = 0.1$ for each element, and the net value of T_{\min} for the tandem limiter is 10^{-4} . This is much easier to achieve. The greater the number of cells, the larger the average fluence in each cell, and the separation of the cells is proportional to the square of the distance from focus, Z .⁸ This can be extended to the limiting case of a single element with a graded molecular density, $N(z) \approx 1/\sigma_{\text{ex}}|Z|$.^{11,32} In this case, the on-axis fluence

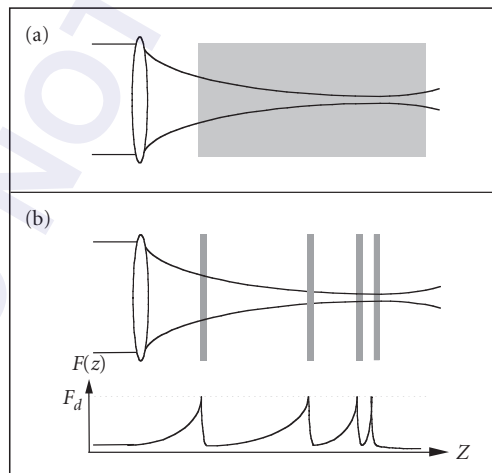


FIGURE 5 Optimization of limiter geometry: (a) for a 2PA limiter and (b) for an RSA limiter. The graph in (b) is a sketch of the on-axis fluence through the limiter near the maximum operating energy.

would remain constant through the RSA material at some designed value of the input energy, usually just below the damage threshold. To avoid problems of generating the exact molecular density distribution, approximating the distribution with a steplike series of adjacent cells of different thickness and density has also been proposed. Like the tandem devices, these designs must be modified to account for beam distortion.²⁷

Limiting via Nonlinear Refraction

From Kramers-Krönig relations, we know that all materials exhibiting nonlinear absorption must also exhibit nonlinear refraction.³³ A consequence of this is that each process that gives rise to optically induced changes in absorption must also result in changes to the refractive index. This can usually be expressed as

$$n_{\text{eff}} = n_0 + n_2 I + \sigma_r N_{\text{ex}} \quad (2)$$

where n_2 describes instantaneous index changes proportional to incident irradiance and σ_r describes the change in index due to population of excited states. As described in Chap. 16, “Third-Order Optical Nonlinearities,” of this volume and in Refs. 33 and 34, n_2 is related to the 2PA coefficient, β , by Kramers-Krönig relations. σ_r and σ_{ex} are related in a similar manner. Such index changes can occur even at wavelengths where there is no change in absorption.

As a focused beam has a spatially varying irradiance, then the induced index change varies across the beam profile, causing the beam to be strongly distorted upon propagation. Near focus, the beam is usually brightest in the center, so for a negative index change (where the index decreases with increasing irradiance or fluence), the nonlinear material will behave like a lens with negative focal length, and the beam is defocused. This process is referred to as *self defocusing*. If the sign of the index change is positive, *self-focusing* results. Both of these effects can cause the beam to spread in the far field and hence limit the energy density in the far field, although the geometrical alignment may be different for optimal limiting in each case.³⁵ This means that E_{foc} may be strongly limited without necessarily limiting E_{max} . The presence of pupil-plane apertures in the imaging system combined with the beam distortion may also result in limiting of E_{max} . An advantage of this method over nonlinear absorption devices is that there is no need to absorb large amounts of energy in the nonlinear material, which could cause thermal damage problems. A potential problem is that inadvertent refocusing of the eye could reduce the defocusing effect of the limiter. However, the nonlinear refraction usually aberrates the beam sufficiently that this is not a concern.

In terms of NLR, a *thick* limiter is defined as one where the propagation path in the nonlinear material is long enough that the index changes cause the beam to change its size inside the material. This process is sometimes referred to as *internal self-action*.³⁶ In this situation, the limiting behavior differs considerably between positive and negative index changes. Self-focusing causes the irradiance to be increased, causing more self-focusing, which becomes a catastrophic effect once a critical input power is reached.³⁷ This results in breakdown of the medium, which can strongly scatter the light and hence effectively limit the transmitted energy. It also causes damage to the material, but this is not a problem if the nonlinear medium is a liquid. Self-defocusing, on the other hand, reduces the irradiance, so that the limiting occurs more gradually as the input energy is increased.

Thermally induced index changes are also important in optical limiters. In liquids, where the thermal expansion is large, the index change results from the change in density upon heating due to laser absorption.³⁸ Hence the index decreases with temperature, giving a selfdefocusing effect. The turn-on time for defocusing is dictated by the time taken for the liquid to expand across the width of the laser beam, which is roughly given by the beam radius divided by the speed of sound.³⁹ The turn-off time is dictated by thermal diffusion. In solids, thermal expansion is much smaller, but other effects, such as temperature dependence of the absorption edge, can cause thermally induced index changes. These usually result in an increase in index with temperature. As this is a local effect, the turn-on time is usually very fast, but turn-off times depend again on thermal diffusion. Thermal

self-focusing in solid-state limiters can be a problem, leading to optical damage. Although thermal defocusing in liquids can be used to produce limiting, and some of the first passive optical limiters were based on this effect,³⁸ it usually degrades the performance of limiters based on NLA. This is one reason to use solid polymer host matrices for RSA dyes.⁴⁰

Limiting via Nonlinear Scattering

Like absorption, scattering is also capable of strongly attenuating a transmitted beam. Nonlinear scattering is possible by laser-induced creation of new scatter centers or by laser-induced changes in the refractive index difference between existing scatter centers and their surroundings. In the latter case, glass scatterers (such as small particles, a rough surface, highly porous glass, or a regular array of holes in glass) are index-matched by immersion in a liquid.^{41–43} In this state, the composite material is clear and highly transparent. A small amount of absorber dye is dissolved in the liquid, so that when illuminated by a strong laser pulse, the solvent is heated and the index matching is lost, resulting in strong scattering.

In the former case, new scatter centers can be created when small particles, such as carbon black, molecular clusters, or other absorbing particles, are exposed to intense laser radiation.^{44,45} In their normal state, such particles are very small, and although they may have a very strong optical absorption coefficient, due to their small size they neither absorb nor scatter much radiation. Upon absorption of radiation, the particles rapidly heat and ionize. This can cause the formation of microplasmas, which grow rapidly and strongly scatter the laser radiation. If the particles are suspended in a liquid, the heating can cause subsequent formation of microbubbles, which also strongly scatter light.^{46,47} In either event, the scattering produces strong optical limiting. As the size of the scattering particles approaches the wavelength, the scattering is predominantly in the forward direction, which could reduce performance for low f -number imaging systems. Another problem is that the limiting process destroys the particles, making this mechanism unsuited to protection against high repetition-rate lasers. This might be overcome by flowing the suspensions. It has been shown that nonlinear scattering may also be an important yet unintentional mechanism in the operation of RSA limiters based on organic molecules.⁴⁷ It is likely that that this is due to incomplete dissolution of the organic material, which leaves small clusters of undissolved material in suspension in the solvent.

Other Mechanisms

While the majority of the results reported on passive limiting employ one of the three previously noted mechanisms, many other schemes have been proposed and demonstrated. Most of these may involve some sort of change in refractive index or absorption, but they may use the change in a manner different from those just described. Although such schemes are too numerous to fully document here, the two following examples are worthy of mention.

Photorefractive Photorefractives change their index when exposed to light, and they do so in such a way that the index change is in proportion to the gradient of light intensity. This is achieved by a complex process involving photoexcitation of charge carriers and diffusion of those carriers that results in a space charge field.⁴⁸ This field in turn causes an index change via the electro-optic effect. Due to the dependence on the gradient of intensity, the photo-refractive effect is usually employed in situations where a periodic modulation of the irradiance induces a phase grating (a periodic modulation of the index). For optical-limiting applications, the interference is obtained by picking off a portion of the input laser beam with a beam splitter and overlapping it with the original beam in a photorefractive crystal.⁴⁹ Alternatively, a reflection from the rear surface of the crystal is used to interfere with the forward-going beam.⁵⁰ In both cases, a grating is produced that, via *two-beam coupling*, strongly scatters the incoming laser beam, limiting the transmitted energy.⁵¹ An interesting side effect of this mechanism is that the limiting is coherence-dependent as well as

intensity-dependent. Due to the requirement for charge diffusion, the turn on time is relatively slow, so that this type of limiter is usually only suitable for pulses of millisecond or longer duration.

Optically Induced Phase Change A number of materials show a reversible, thermally induced semiconductor-metal phase transition upon illumination with strong laser radiation.⁵² These materials are transparent to infrared radiation in their semiconducting state but highly reflective in their metallic state. Hence, in the infrared these materials may be transparent for low powers, while at high powers, weak optical absorption and subsequent heating may induce the strongly reflecting metallic phase, blocking the transmitted light. Some examples of materials of this type are Ag_2S , $\text{TmSe}_x\text{Te}_{1-x}$, and vanadium oxides, $\text{V}_x\text{O}_{2x-1}$. To be effective, such a material must be stable in its transparent state, have a small latent heat associated with the phase transition, and require a reasonably small temperature change (~ 100 K) to induce the phase transition. Vanadium oxides with compositions close to VO_2 or V_2O_3 comprise the most-studied class of these materials for optical limiting, having a phase transition temperature at around 70°C . In thin-film form and with appropriate antireflection coatings, these materials have high broadband transmission through the infrared (3 to $12\ \mu\text{m}$) which drops to around 1 percent in the metallic phase.⁵²

13.3 EXAMPLES OF PASSIVE OPTICAL LIMITING IN SPECIFIC MATERIALS

Unavoidably, most of the materials used for passive optical limiting exhibit a combination of the nonlinear properties previously described. This usually complicates matters, but in some cases the different nonlinearities may be used to complement each other. In this section, we briefly present a few specific examples of limiting devices that illustrate how some of the principles just described may be applied in practice.

Semiconductors

Semiconductors exhibit a variety of strong nonlinear absorption and refraction effects.⁵³ Due to their broad absorption bands, they are capable of producing 2PA over a broad wavelength range where the linear absorption is low. Moreover, the carriers excited by 2PA produce very strong absorption and negative nonlinear refraction. The wavelength range of operation for 2PA, avoiding linear interband absorption, is $E_g/2 < h\nu < E_g$. E_g is the bandgap of the semiconductor and $h\nu$ is the photon energy. For example, in ZnSe, this corresponds to an operating wavelength range of about 480 to 900 nm, while in InSb the range is 7 to $14\ \mu\text{m}$ at room temperature. Over this range, the combined effects of 2PA and free carrier absorption and refraction are more or less constant for a given semiconductor. However, the nonlinearities scale very strongly with bandgap. If we keep the ratio of photon energy to bandgap energy fixed, 2PA scales as E_g^{-3} and free carrier absorption and refraction scale approximately as E_g^{-2} .^{6,23} Hence, semiconductors work significantly better for the infrared than for the visible. A problem with semiconductors is that they tend to have low damage thresholds. It was shown by Van Stryland et al.²⁸ that this can be overcome by the use of the thick-limiter geometry. In this case, a thick sample of the large-gap semiconductor ZnSe was used to demonstrate limiting of 30-ps pulses at a wavelength of 532 nm, as shown in Fig. 6. Due to a combination of 2PA and free-carrier absorption and refraction that occurs prior to focus, it was not possible to damage the ZnSe in the bulk. Hence, these devices were labeled *self-protecting* limiters. The FOM was measured as 2×10^4 . The linear transmittance was 40 percent, probably due to a combination of scatter and parasitic absorption. Despite this good performance with ps pulses, the self-protection does not prevail for nanosecond pulses. This is thought to be due to the effects of carrier recombination and diffusion, which reduce the carrier defocusing effect, allowing positive thermal index changes to dominate.

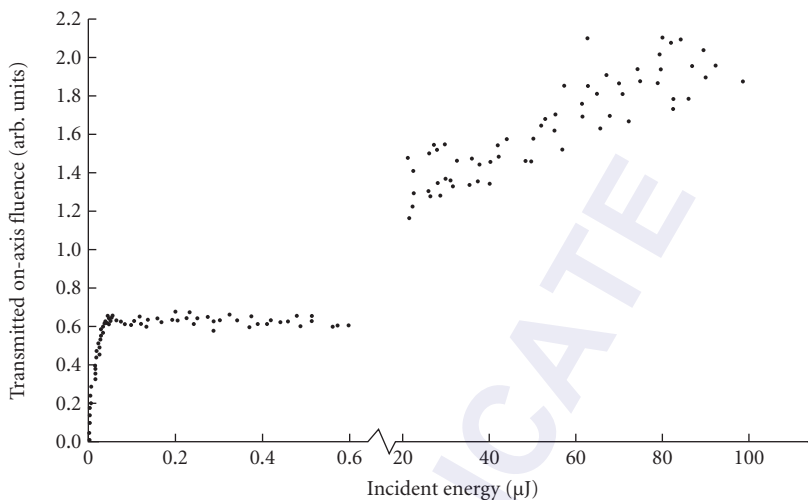


FIGURE 6 Limiting curve for ZnSe thick limiter for 30 ps, 532 nm. (After Ref. 28.)

Organics

Organic molecules have attracted interest for optical limiting for their attractive NLA properties in the visible and near-infrared. While some materials have shown promise for their combined 2PA and excited-state absorption properties,^{26,29} so far these effects have not proven large enough over a sufficiently broad wavelength band to be practical. By far the most attention has been paid to organics for their RSA properties.^{25,54} Generally, the mechanism is as shown in Fig. 4, but often the excited states may also relax into long-lived triplet states. This effectively increases the lifetime of the highly absorbing excited state. The transfer to triplet states may be enhanced by the addition of heavy atoms or paramagnetic groups to the molecule or solvent.⁵⁵ Several groups have demonstrated promising results with optimized limiting devices using phthalocyanines or other similar RSA molecules.^{40,56} One example of an optimized limiter based on this type of material is that of Perry et al.⁴⁰ Here an indium phthalocyanine chloride (InClPc) was incorporated into a PMMA polymer host to make a solid-state limiting material. Slices of this material were used to make a three-element tandem limiter. The device, which had a linear transmittance of about 55 percent (70 percent internal transmittance), was designed to operate with a maximum fluence of 3 J/cm² in an $f/5$ focusing geometry. The combination of solid-state host and low design fluence helps minimize the detrimental effects of thermally induced refractive index changes. The device had a minimum transmittance of 0.185 percent at the maximum input energy of 6.5 mJ, corresponding to a maximum output energy of 12 μ J. This was a factor of four greater than predicted by simple design models, which assume a constant beam shape. This discrepancy is small compared with similar liquid-based limiters, which suffer from much greater thermal refraction. To properly design limiters of this type, propagation codes have been developed that account for all NLA and NLR mechanisms, including thermal refraction, and that are capable of modeling internal self-action.⁵⁷

Carbon-Black Suspensions

It was shown by Mansour et al.⁴⁴ that dilute ink exhibits very strong, broadband optical limiting properties for nanosecond pulses. Ink is a liquid suspension of amorphous carbon particles. Figure 7 shows an example of limiting of 14-ns, 532-nm pulses using a carbon-black suspension (CBS) in a 50:50 water/ethanol mixture with $T_l = 70$ percent. Similar results are obtained with a

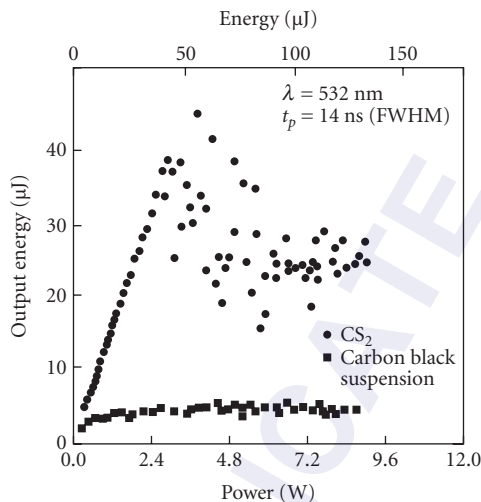


FIGURE 7 Optical limiting of 14-ns pulses at 532-nm wavelength in 1-cm thick, 70 percent linear transmittance CBS. Also shown are results of limiting in CS_2 , for comparison. (After Ref. 44.)

1064-nm laser wavelength, indicating an extremely broadband limiting response. By observing scattered light intensity as a function of input irradiance, it was clearly shown that the incident light becomes very strongly scattered as the incident energy increases. Measurements of the angular distribution of scattered light show a Mie scattering pattern typical of scattering particles a factor of ≈ 3 larger than the original carbon-black particles. Similar results were found for a layer of small carbon particles deposited on a glass surface. It was concluded from this that the limiting is a result of scattering and absorption by microplasmas formed after thermionic emission from the laser-heated carbon particles. Later studies have indicated that the nonlinear scattering may result from microbubbles formed in the solvent by heating of the carbon. There is clear evidence that for longer pulses, the limiting is dependent on the volatility of the solvent,⁴⁶ and imaging techniques have shown that bubbles may persist in the focal volume 100 ns after the pulse.⁴⁷ It is quite feasible that microplasmas may be responsible for limiting on shorter (< 10 ns) time scales, while bubbles play a more important role for longer (~ 100 ns) pulses.

Photorefractives

Although it has been shown that the photorefractive effect can occur over a vast range of time scales, effects large enough for practical devices typically have millisecond and longer response times. Cook et al.⁵⁰ exploited the large photorefractive two-beam coupling gain in $\text{Fe}:\text{LiNbO}_3$ to demonstrate strong optical limiting for millisecond pulses or c.w. lasers. A weak reflection of the input beam from the rear surface of the lithium niobate crystal was sufficient to initiate the two-beam coupling. As shown in Fig. 8, this produced a change in OD of up to 4 with a response time of about 2 ms. The linear transmittance for these samples is typically in the range of 30 to 60 percent. The limiting effect is not strongly sensitive to wavelength, operating between 420 and 670 nm. It is found that the optimum f -number for limiting was about $f/20$, and limiting performance drops off rapidly as the f -number is decreased. This results from the trade-off between the high irradiance produced by small focused spot sizes and the long interaction length produced by large spot sizes.

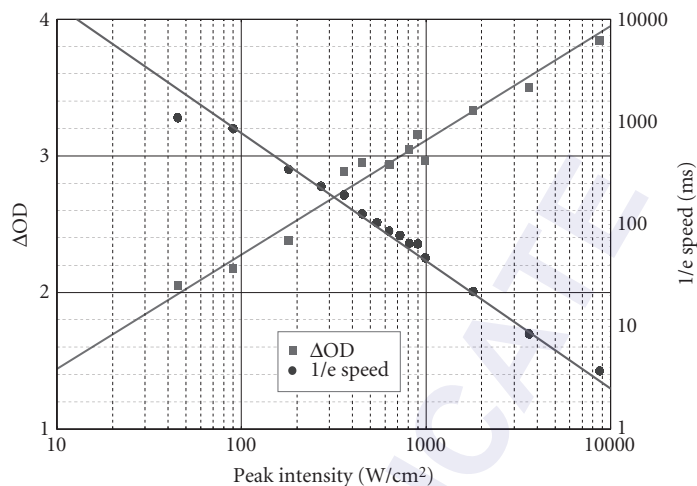


FIGURE 8 Far-field optical limiting and response time in a 3.4-mm path crystal of 0.03 percent Fe + 0.05 percent Tb:LiNbO₃ at 523.5 nm. (After Ref. 50.)

Liquid Crystals

Liquid crystals are composed of highly anisotropic molecules. By illuminating a nematic liquid crystal with a linearly polarized laser beam, the molecules can align with the electric field in the beam, inducing an irradiance-dependent birefringence in the bulk liquid. These effects are large, but typically take milliseconds to seconds for the realignment process. Recently, it has been shown that doping a liquid crystal with certain dyes can induce molecular reorientation at extremely low powers.⁵⁸ Khoo et al.⁵⁹ have shown that this effect may be used to achieve extremely low-power c.w. limiting. Using a twisted nematic 5CB film with 1 percent methyl red doping, between crossed polarizers, the maximum transmitted energy of a c.w. argon ion laser beam was kept below 13 μ W for inputs up to 140 mW, with $T_L = 10$ percent, including Fresnel reflections and losses at the polarizers. Figure 9 shows a photograph of the transmitted beam for low and high powers, with an image that was simultaneously transmitted by the system. An advantage of liquid crystals is that they can be highly transparent across the entire visible spectrum.

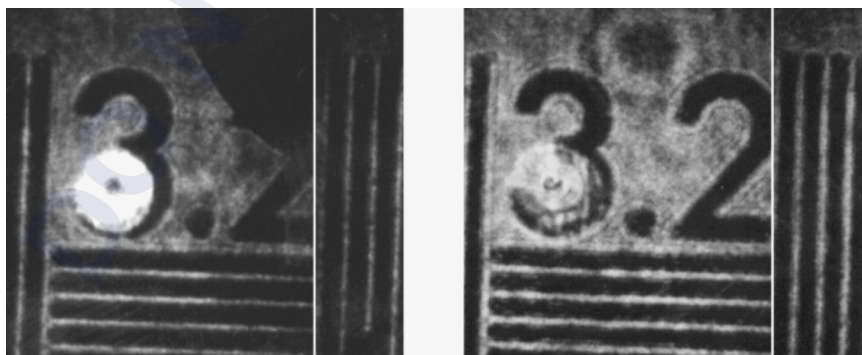


FIGURE 9 Suppression of a c.w. visible laser beam while transmitting the image of a resolution test chart, using a dye-doped nematic liquid crystal passive limiter. (a) Low laser intensity, showing no limiting effect and (b) above the limiting threshold. (After Ref. 59.)

13.4 REFERENCES

1. L. W. Tutt and T. F. Boggess, "A Review of Optical Limiting Mechanisms and Devices Using Organics, Fullerenes, Semiconductors and Other Materials," *Prog. Quantum. Electron.* **17**:299–338 (1993).
2. G. L. Wood, W. W. Clark, M. J. Miller, G. J. Salamo, and E. J. Sharp, "Evaluation of Passive Optical Limiters and Switches," *Proceedings of SPIE Conference on Materials for Optical Switches, Isolators and Limiters*, **1105**:154–181 (1989).
3. E. W. Van Stryland, M. J. Soileau, S. Ross, and D. J. Hagan, "Passive Optical Limiting: Where Are We?" *Nonlinear Optics* **21**:29–38 (1998).
4. M. Sheik-Bahae, D. C. Hutchings, D. J. Hagan, and E. W. Van Stryland, "Dispersion of Bound Electronic Nonlinear Refraction in Solids," *IEEE J. Quantum Electron.* **27**:1296–1309 (1991).
5. A. A. Said, M. Sheik-Bahae, D. J. Hagan, T. H. Wei, J. Wang, J. Young, and E. W. V. Stryland, "Determination of Bound and Free-Carrier Nonlinearities in ZnSe, GaAs, CdTe and ZnTe," *J. Opt. Soc. Am. B* **9**:409 (1991).
6. B. S. Wherrett, "Scaling Rules for Multiphoton Interband Absorption in Semiconductors," *J. Opt. Soc. Am. B-Optical Physics* **1**:67–72 (1984).
7. D. J. Hagan, E. W. Van Stryland, M. J. Soileau, Y. Y. Wu, and S. Guha, "Self-Protecting Semiconductor Optical Limiters," *Opt. Lett.* **13**:315–317 (1988).
8. D. J. Hagan, T. Xia, A. A. Said, T. H. Wei, and E. W. V. Stryland, "High Dynamic Range Passive Optical Limiters," *Int. J. Nonlin. Opt. Phys.* **2**:483–501 (1993).
9. D. H. Sliney and M. L. Wolbarsht, *Safety with Lasers and Other Optical Sources*, Plenum, New York, 1980.
10. D. H. Sliney, "Retinal Damage from Laser Radiation," *Nonlinear Optics* **21**:1–18 (1998).
11. P. A. Miles, "Bottleneck Optical Limiters: The Optimal Use of Excited-State Absorbers," *Appl. Opt.* **33**:6965 (1994).
12. R. C. Hollins, "Optical Limiters, Spatial, Temporal and Spectral Effects," *Nonlinear Optics* **21**:39–48 (1998).
13. J. A. Zuclich, D. J. Lund, P. R. Edsall, R. C. Hollins, P. A. Smith, B. E. Stuck, and L. N. McLin, "Experimental Study of the Variation of Laser-Induced Retinal Damage Threshold with Retinal Image Size," *Nonlinear Optics* **21**:19–28 (1998).
14. J. A. Hermann and J. Staromlynska, "Special Issue on Optical Limiters, Switches and Discriminators," *J. Nonlinear Optical Phys. and Materials* **2**:271–337(1993).
15. F. Kajzar, "Proceedings of First International Workshop on Optical Power Limiting, Cannes, France," in *Nonlinear Optics, Principles, Materials, Phenomena, Devices*, Gordon & Breach, 1998.
16. C. M. Lawson, "Nonlinear Optical Liquids," *Proceedings of SPIE*, SPIE, **2853**, Denver, 1996.
17. C. M. Lawson, "Nonlinear Optical Liquids for Power Limiting and Imaging," *Proceedings of SPIE*, SPIE, **3472**, San Diego, 1998.
18. C. M. Lawson, "Power Limiting Materials and Devices," *Proceedings of SPIE*, SPIE, **3798**, Denver, 1999.
19. M. J. Soileau, "Materials for Optical Switches, Isolators and Limiters," *Proceedings of SPIE*, SPIE, **1105**, Orlando, 1989.
20. M. J. Soileau, *Electro-Optical Materials for Switches, Coatings, Sensor Optics and Detectors*, SPIE, **1307**, Orlando, 1990.
21. Crane, K. Lewis, E. W. Van Stryland, and M. Khoshnevisan, *Materials for Optical Limiting*, **374**, MRS, Boston, 1994.
22. R. Sutherland, R. Pachter, P. Hood, D. J. Hagan, K. Lewis, and J. Perry, *Materials for Optical Limiting II*, **479**, MRS, San Francisco, 1997.
23. E. W. Van Stryland, H. Vanherzeele, M. A. Woodall, M. J. Soileau, A. L. Smirl, S. Guha, and T. F. Boggess, "2 Photon-Absorption; Nonlinear Refraction; and Optical Limiting in Semiconductors," *Opt. Eng.*, **24**:613–623 (1985).
24. C. R. Guliano and L. D. Hess, "Nonlinear Absorption of Light: Optical Saturation of Electronic Transitions in Organic Molecules with High Intensity Laser Radiation," *IEEE J. Quantum Electron.* **QE3**:358–367 (1967).
25. T. H. Wei, D. J. Hagan, M. J. Sence, E. W. V. Stryland, J. W. Perry, and D. R. Coulter, "Direct Measurements of Nonlinear Absorption and Refraction in Solutions of Phthalocyanines," *Appl. Phys.* **B54**:46 (1992).
26. A. A. Said, C. Wamsley, D. J. Hagan, E. W. V. Stryland, B. A. Reinhardt, P. Roderer, and A. G. Dillard, "Third and Fifth Order Optical Nonlinearities in Organic Materials," *Chem. Phys. Lett.* **228**:646–650 (1994).

27. T. Xia, D. J. Hagan, A. Dogariu, A. A. Said, and E. W. V. Stryland, "Optimization of Optical Limiting Devices Based on Excited-State Absorption," *Appl. Opt.* **36**:4110–4122 (1997).
28. E. W. Van Stryland, Y. Y. Wu, D. J. Hagan, M. J. Soileau, and K. Mansour, "Optical Limiting with Semiconductors," *J. Opt. Soc. Am. B* **5**:1980–1989 (1988).
29. M. Albota, D. Beljonne, J. L. Brédas, J. Ehrlich, J. Fu, A. Heikal, T. Kogej, M. Levin, S. R. Marder, D. McCord-Maughon, J. W. Perry, H. Röckel, M. Rumi, G. Subramaniam, W. W. Webb, X. L. Wu, and C. Xu, "Design of Organic Molecules with Large Two-Photon Absorption Cross Sections," *Science* **281**:1653–1656 (1998).
30. M. Sheik-Bahae, A. A. Said, D. J. Hagan, M. J. Soileau, and E. W. Van Stryland, "Nonlinear Refraction and Optical Limiting in Thick Media," *Opt. Eng.* **30**:1228–1235 (1991).
31. A. A. Said, R. DeSalvo, M. Sheik-Bahae, D. J. Hagan, and E. W. V. Stryland, "Self Protecting Optical Limiters Using Cascading Geometries," *Proceedings of SPIE*, SPIE, Orlando, Florida, 1992.
32. S. W. McCahon and L. W. Tutt, "Optical Limiter Including Optical Congruence and Absorbing Body with Inhomogeneous Distribution of Reverse Saturable Absorption," U.S. Patent 5,080,469, 1992.
33. D. C. Hutchings, M. Sheik-Bahae, D. J. Hagan, and E. W. Van Stryland, "Kramers-Kronig Relations in Nonlinear Optics," *Opt. and Quantum Electron.* **24**:1–30 (1992).
34. M. Sheik-Bahae, D. J. Hagan, and E. W. V. Stryland, "Dispersion of Bound Electronic Nonlinear Refraction in Solids," *Phys. Rev. Lett.* **65**:96–99 (1990).
35. M. Sheik-Bahae, A. A. Said, and E. W. V. Stryland, "High-Sensitivity, Single-Beam n_2 Measurements," *Opt. Lett.* **14**:955–957 (1989).
36. A. E. Kaplan, "External Self-Focusing of Light by a Nonlinear Layer," *Radiophys. Quantum Electron.* **12**:692–696 (1969).
37. M. J. Soileau, W. E. Williams, and E. W. Van Stryland, "Optical Power Limiter with Picosecond Response Time," *IEEE J. Quantum Electron.* **QE-19**:731–735 (1983).
38. R. C. Leite, S. P. S. Porto, and T. C. Damen, *Appl. Phys. Lett.* **10**:100 (1967).
39. D. I. Kovsh, D. J. Hagan, and E. W. V. Stryland, "Numerical Modeling of Thermal Refraction in Liquids in the Transient Regime," *Opt. Exp.* **4** (1999).
40. J. W. Perry, K. Mansour, I.-Y. S. Lee, X.-L. Wu, P. V. Bedworth, C.-T. Chen, D. Ng, S. R. Marder, P. Miles, T. Wada, M. Tian, and H. Sasabe, "Organic Optical Limiter with a Strong Nonlinear Absorptive Response," *Science* **273**:1533 (1996).
41. B. L. Justus, A. J. Campillo, and A. L. Huston, "Thermal-Defocusing/Scattering Optical Limiter," *Opt. Lett.* **19**:673 (1994).
42. H.-B. Lin, R. J. Tonucci, and A. J. Campillo, "Two-Dimensional Photonic Bandgap Optical Limiter in the Visible," *Opt. Lett.* **23**:94–96 (1998).
43. D. J. Hagan, S. S. Yang, C. Basanez, E. W. Van Stryland, W. Moreshead, and J. L. Nogueas, "Optical Limiting via Nonlinear Scattering with Solgel Hosts," *Proceedings of SPIE Conf. on Power Limiting Materials and Devices*, SPIE, **3798**, Denver, Colorado, 1999.
44. K. Mansour, E. W. V. Stryland, and M. J. Soileau, "Nonlinear Properties of Carbon Black Suspensions," *J. Opt. Soc. Am. B* **9**:1100 (1992).
45. T. Xia, A. Dogariu, K. Mansour, D. J. Hagan, A. A. Said, E. W. V. Stryland, and S. Shi, "Nonlinear Optical Properties of Inorganic Metallic Clusters," *J. Opt. Soc. Am. B* **15**:1497 (1998).
46. K. J. McEwan, P. K. Milsom, and D. B. James, "Nonlinear Optical Effects in Carbon Suspensions," *Proceedings of Nonlinear Optical Liquids for Power Limiting and Imaging*, SPIE, **3472**, San Diego, California, 1998.
47. R. V. Goedert, R. J. Becker, A. Clements, and T. A. Whitaker III, "Overview of the Shadowgraphic Imaging Technique and Results for Various Materials," *Proceedings of Nonlinear Optical Liquids for Power Limiting and Imaging*, SPIE, **3472**, San Diego, California, 1998.
48. M. Cronin-Golumb and A. Yariv, "Optical Limiters Using Photorefractive Limiters," *J. Appl. Phys.* **57**:4906 (1985).
49. S. W. McCahon and M. V. Klein, "Coherent Beam Excisors Using the Photorefractive Effect in BaTiO₃," *Proceedings of SPIE Conference on Materials for Optical Switches, Isolators and Limiters*, SPIE, **1105**, Orlando, Florida, 1989.
50. G. Cook, D. C. Jones, C. J. Finnan, L. L. Taylor, and T. W. Vere, "Optical Limiting with Lithium Niobate," *Proceedings of Power-Limiting Materials and Devices*, SPIE, **3798**, Denver, Colorado, 1999.

51. R. W. Boyd, *Nonlinear Optics*, Academic Press, Boston, 1992.
52. K. L. Lewis, A. M. Pitt, T. Wyatt-Davies, and J. R. Milward, "Thin Film Thermochromic Devices for Non-Linear Optical Devices," *Proceedings of MRS Symposium on Materials for Optical Limiting*, MRS, Boston, 1994.
53. N. Peyghambarian, S. W. Koch, and A. Mysyrowicz, *Introduction to Semiconductor Optics*, Prentice Hall, Englewood Cliffs, New Jersey, 1993.
54. J. S. Shirk, R. G. S. Pong, F. J. Bartoli, and A. W. Snow, "Optical Limiter Using a Lead Phthalocyanine," *Appl. Phys. Lett.* **63**:1880 (1993).
55. J. W. Perry, "Organic and Metal-Containing Reverse Saturable Absorbers for Optical Limiters," in *NonLinear Optics of Organic Molecules and Polymers*, H. S. Nalwa and S. Miyata, (ed.), CRC Press, New York, 1997.
56. A. A. Said, T. Xia, D. J. Hagan, A. Wajsgrus, S. Yang, D. Kovsh, and E. W. V. Stryland, "Liquid-Based Multicell Optical Limiter," *Proceedings of SPIE Conference on Nonlinear Optical Liquids*, SPIE, **2853**, Denver, Colorado, 1996.
57. D. Kovsh, S. Yang, D. J. Hagan, and E. W. Van Stryland, "Nonlinear Optical Beam Propagation for Optical Limiting," *App. Opt.* **38**:5168–5180 (1999).
58. L. Marrucci, D. Paparo, G. Abbate, E. Santamo, M. Kreuzer, P. Lehnert, and T. Vogeler, "Enhanced Optical Nonlinearity by Photoinduced Molecular Orientation in Absorbing Liquids," *Phys. Rev. A* **58**:4926 (1998).
59. I. C. Khoo, M. V. Wood, M. Y. Shih, and P. H. Chen, "Extremely Nonlinear Photosensitive Liquid Crystals for Image Sensing and Sensor Protection," *Opt. Exp.* **4**:432–442 (1999).

This page intentionally left blank.

DO NOT DUPLICATE

ELECTROMAGNETICALLY INDUCED TRANSPARENCY

Jonathan P. Marangos

*Quantum Optics and Laser Science Group
Blackett Laboratory, Imperial College
London, United Kingdom*

Thomas Halfmann

*Institute of Applied Physics
Technical University of Darmstadt
Darmstadt, Germany*

14.1 GLOSSARY

Terms and Acronyms

EIT	electromagnetically induced transparency
CPT	coherent population trapping
STIRAP	stimulated Raman adiabatic passage
RAP	rapid adiabatic passage
SCRAP	Stark chirped rapid adiabatic passage
CPR	coherent population return
LWI	lasing without inversion
lambda (Λ) scheme	three coupled atomic levels with the initial and final states at lower energy than the intermediate state
ladder scheme	three coupled atomic levels with the energy of the initial state below the intermediate state, and the energy of the intermediate state below the final state
vee (V) scheme	three coupled atomic levels with the initial and final states at higher energy than the intermediate state
cw	continuous wave

Symbols

$ a\rangle$	quantum state of atom
E_a	energy of quantum state $ a\rangle$
ω_{ab}	angular transition frequency between states $ a\rangle$ and $ b\rangle$ (rad s^{-1})
Δ_{ab}	detuning of a light field from an atomic transition at ω_{ab} (rad s^{-1})
μ_{ab}	transition dipole moment between two states $ a\rangle$ and $ b\rangle$ (Cm)

$\text{Re } \chi^{(1)}$	real part of linear susceptibility (dispersion)
$\text{Im } \chi^{(1)}$	imaginary part of linear susceptibility (absorption)
$\chi^{(3)}$	nonlinear susceptibility of third order ($\text{m}^2 \text{V}^{-2}$)
N_A	number of atoms
n_A	atomic density (cm^{-3})
E	electric field strength (Vm^{-1})
ω	angular frequency of a radiation field (rad s^{-1})
λ	wavelength of a radiation field (nm)
P	macroscopic polarization (Cm^{-2})
Ω	Rabi frequency (rad s^{-1})
ρ_{ab}	density matrix element
Γ_a	decay rate of state $ a\rangle$ (rad s^{-1})
γ_{ab}	dephasing rate of the coherence ρ_{ab} (rad s^{-1})
γ_{Doppler}	Doppler (inhomogeneous) linewidth (rad s^{-1})
γ_{laser}	laser bandwidth (rad s^{-1})

14.2 INTRODUCTION

Electromagnetically induced transparency (EIT) is a quantum interference phenomenon that arises when coherent optical fields couple to the states of a material quantum system. In EIT the interference occurs between alternative transition pathways, driven by radiation fields within the internal states of the quantum system. Interference effects arise, because in quantum mechanics the probability amplitudes (which may be positive or negative in sign), rather than probabilities, must be summed and squared to obtain the total transition probability between the relevant quantum states. Interference between the amplitudes may lead to either an enhancement (constructive interference) or a complete cancellation (destructive interference) in the total transition probability. As a consequence, interference effects can lead to profound modification of the optical and nonlinear optical properties of a medium. Thus, control of optical or nonlinear optical properties and processes becomes possible.

Interference effects of this kind are well known in physics. These occur naturally if there are two transition pathways available to the same final state. Fano interferences exhibit an example of interference between two transition pathways. In this case the two pathways are direct photoionization from a quantum state to the ionization continuum and photoionization from the quantum state via an intermediate autoionizing state.^{1,2} The interference between these two pathways leads to asymmetric resonances in the photoionization spectrum. The photoionization cross section vanishes at certain excitation frequencies, that is, complete destructive interference occurs. This process is well known for radiative transitions to autoionizing states in atoms or to predissociating states in molecules. It was also predicted to occur in semiconductor quantum wells.³

Interfering transition pathways can also be deliberately induced by application of resonant laser fields to multilevel atomic systems. Perhaps the most striking example for this type of interference is the cancellation of absorption for a probe radiation field, tuned in resonance to an atomic transition. Usually the resonant excitation would lead to strong absorption of the probe field. However, if the atoms are prepared by EIT,⁴⁻⁶ the absorption is essentially switched off. EIT exhibits a laser-induced interference effect between the quantum coherences in the atom, which renders an initially highly opaque medium into an almost transparent medium. Similarly the refractive properties of the medium may be greatly modified.^{5,7,8} For instance the usual correlation of high refractive index with high absorption can be broken, leading to the creation of media with unique optical properties.

There has been a considerable research effort devoted to EIT and related topics over the last few years. This has been motivated by the recognition of a number of new potential applications, for example, lasers without inversion, highly efficient nonlinear optical processes, storage of light pulses

and quantum information, lossless propagation of laser beams through optically thick media, and highly efficient and selective population transfer via coherent adiabatic processes. EIT is one of an interrelated group of processes, for example, including coherent population trapping (CPT) and coherent adiabatic population transfer, that result from externally induced quantum mechanical coherence and interference. In contrast, the earlier ideas associated with CPT (first observed in 1976)⁹ had found application mostly as a tool of high-resolution spectroscopy, rather than as a new direction in nonlinear optics. Therefore, the concept of EIT has contributed a distinctive new thrust to work on atomic coherence and its applications—a thrust, which is of direct interest to optical scientists and engineers.

To explain the basic idea and applications of EIT, an equivalent picture to interfering transition pathways is provided by the concept of laser-dressed states. In terms of quantum mechanics, the *dressed states* are the eigenstates of the quantum system, including strong interaction with driving radiation fields. The dressed states are coherent superpositions of the *bare states*, that is, the eigenstates of the quantum system without external interaction. The coherent superpositions have well-defined amplitudes and phases that describe the relationship between the atomic states in the superposition. The reader is referred to *The Theory of Coherent Atomic Excitation* by Shore¹⁰ for a complete account of these ideas. An important feature of EIT is the preparation of large populations of these coherently driven, uniformly phased atoms. Such media are termed *phasesonium* by Scully,⁷ to convey the basic idea. The (both linear and nonlinear) optical properties of the coherent medium are very different from those of a normal, incoherently driven medium. In the dressed medium the terms of linear and nonlinear susceptibilities can be retained only to the extent that it is recognized that all these resonant processes are highly nonperturbative. As we will discuss later, even the so-called “linear” processes now involve the coupling of atoms with many photons. An important consequence of this is that the magnitudes of linear and nonlinear susceptibilities can reach equality in a phase-coherent medium. This is in marked contrast to the normal situation. Usually the nonlinear susceptibility would give rise to nonlinear optical processes, which are many orders of magnitude weaker than those arising from the linear susceptibility.

The exceptionally high efficiencies of nonlinear optical processes in gas phase media, prepared by EIT, therefore constitute an important feature of EIT. The large conversion efficiencies in gas phase media, driven to EIT, become comparable to nonlinear optical processes in optical crystals. Thus a renewed interest in gas phase nonlinear optical devices possessing unique capabilities [e.g., high conversion efficiencies into the spectral regions of vacuum-ultraviolet (VUV) and far-infrared (IR) radiation] has occurred. There have been a number of notable recent demonstrations of EIT, applied to frequency conversion. A near-unity frequency conversion into the far-ultraviolet spectral region was reported for a four-wave mixing scheme in lead vapor. The lead atoms were prepared in the state of *maximal coherence*, that is, maximal polarization, by EIT. The uniquely high conversion efficiency arises since the nonlinear terms become equal in magnitude to the linear terms. Besides applications in dense gas phase media with thermal velocity distribution, large optical nonlinearities, induced by EIT were also studied in laser-cooled atoms. In such media, a successful implementation of EIT requires only quite weak laser couplings. Thus, at maximum transparency there is an extremely steep dispersion as a function of the driving laser frequency, that is, the detuning from the atomic resonance. The consequence of this steep slope is a very slow group velocity for optical pulses propagating through the medium.¹¹ Massive optical nonlinearities accompany the steep dispersion. These are manifested as very large nonlinear refractive indices that are many orders of magnitude larger than any previously observed. These huge nonlinearities are the subject of current research activity. They offer the likelihood of efficient nonlinear optical processes at the few photon level.

In addition to the applications, described above, the coherence and interference effects related to EIT may also permit new possibilities to build short wavelength lasers, that is, lasers in the x-ray spectral range. As the Einstein coefficient for spontaneous emission scales with the cube of the transition frequency ω ,³ it is usually hard to achieve population inversion by optical pumping in an x-ray laser. In contrast, the laser concept based on EIT does not rely on population inversion in the atomic laser medium anymore.^{12,13} Lasing without inversion (LWI) has been demonstrated in sodium atoms and rubidium atoms in the visible range.^{14,15} The prospects that this might lead to the construction of lasers which are able to circumvent the usual constraints of achieving inversion in short wavelength

lasers has been much discussed.^{16,17} Related effects on LWI and EIT in semiconductor quantum wells have also been theoretically explored using laser-induced processes¹⁸ or bandgap engineering¹⁹ to create the necessary coherences. Moreover, atomic and molecular coherences, driven by EIT, were also exploited for efficient frequency conversion and the generation of ultrashort laser pulses.^{20–27,28–39}

It is the aim of this chapter to provide an accessible summary of EIT and to present some of the main results obtained in recent research. It is not possible in the space available to cover all work in this field, and we apologize to authors whose important contributions are not directly mentioned. An extensive review of all theoretical and experimental work on related atomic coherence phenomena is also beyond the scope of the present article. The reader is advised to look at a number of reviews on LWI,^{16,17} coherent population trapping,^{10,40–43} coherent adiabatic population transfer,⁴⁴ efficient frequency conversion in coherently prepared media⁴⁵ and laser-induced continuum structure^{46–50} to find these topics presented in detail. The theory pertinent to EIT is sketched in the text, but again the reader is referred to the more detailed treatments published in the literature.⁵ References to relevant literature will be given as they arise. The purpose of the chapter is (1) to communicate the underlying physical principles of EIT and related effects, (2) to describe the manifestations of EIT and to summarize the conditions required to create EIT, and (3) to introduce some potential applications in optical technology.

14.3 COHERENCE IN TWO- AND THREE-LEVEL ATOMIC SYSTEMS

The first experimental work on laser-induced atomic coherence was carried out in the 1970s. Earlier relevant work includes the investigation of dressing two-level systems by strong microwave fields. This led to the observation of splittings between dressed states, that is, Autler-Townes splittings,⁵¹ and photon echoes in two-level systems.⁵² Mollow^{53,54} reported novel features, subsequently termed the Mollow triplet, of resonance fluorescence in a two-level system driven by a strong resonant laser. Much work on two-level systems has been carried out since.^{10,44,55,56} Although two-level systems remain a subject of considerable interest, our concern here is primarily with three-level systems (and in some cases four-level systems).

Atomic coherence and interference in three-level systems was first observed experimentally by Alzetta, Arimondo et al.,^{9,57} and by Gray et al.⁵⁸ The first group of authors performed experiments that established coherence between the Zeeman split lower levels of a sodium atom using a multimode laser. By employing a spatially varying magnetic field Alzetta, Arimondo et al. observed a series of spatially separated dark lines. These resonances correspond to the locations in the magnetic field where the Zeeman splitting matched the frequency difference between modes of the coupling laser. This situation corresponds to a two-photon resonant lambda-type level scheme. Thus, this was the first experimental observation of CPT. The experiments of Gray et al.⁵⁸ involved the preparation of coherence between the hyperfine lower levels of sodium atoms. In these experiments two coincident laser fields are coupled to a three-level lambda scheme of states to create superpositions of the two lower states $|1\rangle$ and $|2\rangle$ (see Fig. 1). One of the superpositions, that is, the coupled state or *bright state* $|C\rangle$ interacts with the fields [see Eq. (1)]. For the other superposition, that is, the noncoupled or *dark state* $|NC\rangle$, interference causes cancellation of the two driven dipoles. Once the coherent states are formed, the population in the system will all be “optically pumped” into the dark state through spontaneous emission from the intermediate state. The optical pumping process occurs on the timescale of a few times the radiative decay time. Once in the dark state there is no process to remove the population. Thus the population is trapped in the dark state.

The basic idea of CPT has been extended to systems, driven by time-varying optical fields to yield very efficient excitation of atomic and molecular states.^{44,59–61} In stimulated Raman adiabatic passage (STIRAP), the noncoupled state $|NC\rangle$ exhibits a specific evolution in time. Initially $|NC\rangle$ is prepared such that it is composed purely of the lowest state $|1\rangle$. For intermediate times $|NC\rangle$ evolves as a superposition of the two lower states $|1\rangle$ and $|2\rangle$, with no contribution from the intermediate state $|3\rangle$. Finally $|NC\rangle$ is composed purely of state $|2\rangle$. Thus, population is transferred completely

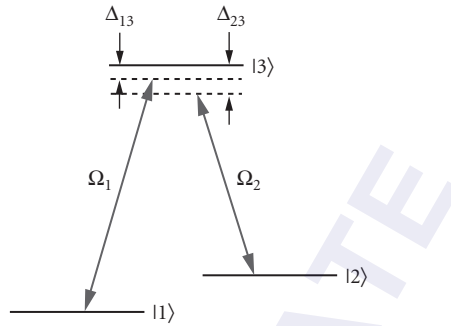


FIGURE 1 CPT in a three-level lambda configuration. Two radiation fields, that is, couplings Ω_1 and Ω_2 , are applied with frequencies close to the single-photon resonances at ω_{13} and ω_{23} .

between the lower states $|1\rangle$ and $|2\rangle$, with no intermediate storage in state $|3\rangle$. The bright state is not populated during the process. This specific preparation of a dark state in STIRAP is achieved by employing counterintuitive pulse sequences, that is, the couplings (or the laser pulses) Ω_1 and Ω_2 still partly overlap in time, but Ω_2 reaches its peak value *prior* to the peak of Ω_1 . The timescale for STIRAP is determined by the evolution of the laser pulses. This contrasts to conventional CPT, that is, driven with coincident radiation fields, when the time evolution is determined by spontaneous emission. We also note that recently extended work has been performed to utilize CPT also in laser cooling and manipulation of trapped atoms. The technique, used in these experiments was termed *velocity-selective coherent population trapping* (VSCPT).^{62–64}

In experiments, dedicated to CPT the primary concern is focused on the manipulation of state populations, essentially of individual atoms. In contrast, for EIT the main interest is the optical response, rather than simply the populations, of the entire medium. The optical response is determined by the coherences rather than the populations. In terms of density matrix elements, in CPT the pertinent quantities are the on-diagonal density matrix elements, that is, the populations; while in EIT they are the off-diagonal density matrix elements, that is, the coherences. Most important, in the limit of a strong coupling field Ω_2 and population initially in the ground state, the coherences, driven by EIT are almost instantaneously established. The timescale of the evolution in EIT is determined by $1/\Omega_2$. For strong excitation, even driven by quite long nanosecond (ns) pulses, the timescale of $1/\Omega_2$ easily reaches the regime of picoseconds, that is, well below the pulse duration. For a successful implementation of population trapping a timescale of several radiative lifetimes is required, that is, typically in the regime of many nanoseconds. From these considerations we see, that though EIT and CPT are closely related, some of their important features as well as their aim are very different.

14.4 THE BASIC PHYSICAL CONCEPT OF ELECTROMAGNETICALLY INDUCED TRANSPARENCY

As discussed in the preceding section, there is a close link between EIT and other phenomena, relying on atomic coherence, that is, adiabatic population transfer processes.^{9,43,44,57–61,65–67} In all these processes, three-level atomic systems are involved—that is, systems that can be adequately reduced to three levels when interaction with the pertinent electromagnetic fields are considered. The atomic dipole selection rules require that two pairs of levels are dipole-coupled, while the transition between the third pair is dipole-forbidden. In Fig. 2, we show the basic three-level schemes. All of the level schemes, discussed in this paper can be reduced to one or other of these schemes. Following

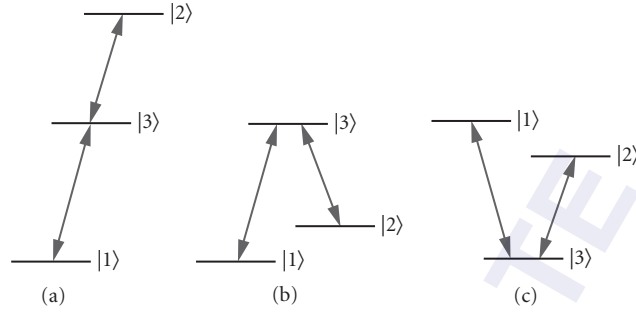


FIGURE 2 Basic three-level schemes: (a) ladder (or cascade) scheme with $E_1 < E_3 < E_2$; (b) lambda (Λ) scheme with $E_1 < E_2 < E_3$; and (c) vee scheme with $E_3 < E_1$ and $E_3 < E_2$.

the nomenclature of Harris et al,⁶ we label the levels $|1\rangle$, $|2\rangle$, and $|3\rangle$. The dipole-allowed transitions are between states $|1\rangle$ and $|3\rangle$ and between states $|2\rangle$ and $|3\rangle$. Classification of the schemes then depends upon the relative energies of the three states: (1) ladder (or cascade) scheme with $E_1 < E_3 < E_2$ (2) lambda (Λ) scheme with $E_1 < E_2 < E_3$, and (3) vee (V) scheme with $E_3 < E_1$ and $E_3 < E_2$. EIT has been extensively studied in all three of these configurations. In a lambda or ladder scheme state $|1\rangle$ is normally the ground state of the atom. This is, where initially the majority of the population resides. In EIT there is no need for significant population transfer. Thus states $|2\rangle$ and $|3\rangle$ remain essentially unpopulated throughout the process. It should be noted that the lambda scheme has a special importance due to the metastability of state $|2\rangle$. As a consequence, very long-lived coherences can be established between states $|1\rangle$ and $|2\rangle$. This leads to near-perfect conditions for EIT.

To understand in more detail, how laser fields interact with a three-level atom to create coherent superpositions of the atomic bare states, we will consider now CPT in a lambda scheme. A three-level lambda system (see Fig. 1) is coupled by two near-resonant laser fields. The interaction strength is defined by the Rabi frequencies $\Omega_1 = \mu_{13} E_1 / \hbar V$ at frequency ω_1 and $\Omega_2 = \mu_{23} E_2 / \hbar V$ at frequency ω_2 with the dipole transition moments μ_{13} and μ_{23} , and the electric fields E_1 and E_2 . The transition frequencies are defined as ω_{12} and ω_{23} . The one-photon detunings of the radiation fields from the atomic resonances are $\Delta_{13} = \omega_{13} - \omega_1$, $\Delta_{23} = \omega_{23} - \omega_2$. The two-photon (Raman) detuning is $\Delta = [(\omega_{13} - \omega_{23}) - (\omega_1 - \omega_2)]$. The Hamiltonian of the bare atom H_0 must be modified to include the interactions due to the two couplings. Thus $H = H_0 + V_1 + V_2$, with the interactions $V_j = \hbar \Omega_j$. The eigenstates of this new Hamiltonian will be linear superpositions of the bare atomic states $|1\rangle$, $|2\rangle$, and $|3\rangle$ (see Refs. 10, 42, 68). For exact two-photon resonance and, that is, $\Delta = 0$, two of the three eigenstates of the total Hamiltonian H turn out to be symmetric and antisymmetric coherent superpositions of the two lower bare states. These superpositions read

$$|C\rangle = \frac{1}{\Omega'} (\Omega_1 |1\rangle + \Omega_2 |2\rangle) \quad (1a)$$

$$|NC\rangle = \frac{1}{\Omega'} (\Omega_2 |1\rangle - \Omega_1 |2\rangle) \quad (1b)$$

where $\Omega' = [\Omega_1^2 + \Omega_2^2]^{1/2}$. It is important to note that no component of the bare state $|3\rangle$ appears in these superpositions. The superposition state $|C\rangle$ is coupled to the intermediate state $|3\rangle$ via electric dipole interaction, that is, $|C\rangle$ is a bright state. In contrast, the other state $|NC\rangle$ is not coupled to state $|3\rangle$, that is, $|NC\rangle$ is a *dark state* or *trapped state*. This is obvious from the total dipole moment for a transition from state $|NC\rangle$ to the bare state $|3\rangle$. If the magnitudes of the coupling fields Ω_1 and Ω_2 are appropriately balanced, the negative sign in the superposition of $|1\rangle$ and $|2\rangle$, which forms the state $|NC\rangle$, causes the transition moment $\langle NC | \mu | 3 \rangle$ to vanish. In effect, the two terms

that are summed to give the transition amplitude between $|NC\rangle$ and $|3\rangle$ are of equal and opposite magnitude, and hence the total amplitude will vanish. In a classical picture, this corresponds to the electron driven by two fields, both of which may be strong, but which exert forces of exactly equal magnitude and opposite directions. This interaction leads to a zero net force and hence the electron stays at rest. In conventional CPT and assuming steady-state conditions, the superposition state $|NC\rangle$ will acquire all of the population of the system through optical pumping. Thus spontaneous emission from state $|3\rangle$ populates state $|NC\rangle$, but absorption losses from state $|NC\rangle$ back to state $|3\rangle$ are not possible.

The noncoupled state $|NC\rangle$ also serves as the key component for coherent population transfer by STIRAP,⁴⁴ which was already briefly discussed above. We assume, that initially all population is in state $|1\rangle$. State $|2\rangle$ is assumed to be metastable, for example, if the level scheme is of lambda-type configuration. If at early times the Rabi frequencies (i.e., the corresponding radiation fields) are applied such that $\Omega_2 \gg \Omega_1$, state $|NC\rangle$ is equal to state $|1\rangle$ [see (Eq. 1)]. All population of the system is prepared in the dark state $|NC\rangle$. No population is in the bright state $|C\rangle$. If at the end of the interaction $\Omega_1 \gg \Omega_2$, state $|NC\rangle$ aligns now parallel to the target state $|2\rangle$. As $\Omega_1 \gg \Omega_2$, the contribution of state $|1\rangle$ is negligible, that is, all population is transferred to the target state $|2\rangle$. The sequence of an initially strong coupling between the states $|2\rangle$ and $|3\rangle$ and a finally strong coupling between the states $|1\rangle$ and $|2\rangle$ exhibits a counterintuitive laser pulse order, which is the typical feature of STIRAP.

We note that in the previous description we ignored a fast time oscillation of the bare states $|1\rangle$ and $|2\rangle$ in the superpositions in Eq. (1). The oscillation occurs at frequencies ϵ_1/\hbar and ϵ_2/\hbar , with the energy of the bare states ϵ_i . In fact these terms disappear when the dipole moments are formed. In typical implementations of CPT, for example, STIRAP, the couplings are of comparable strength, that is, $\Omega_1 \approx \Omega_2$, and the two-photon transition is strongly driven. The transition is “saturated” if the terminology of incoherent excitation is used. We note an interesting feature of CPT with respect to the coherences for the case of the laser frequencies tuned to exact two-photon resonance, but with large single-photon detunings. In this case, state $|3\rangle$ can be adiabatically eliminated from the level scheme.¹⁰ Thus state $|3\rangle$ does not enter into the consideration of the coupling between atoms and fields any more. However, also in this case, that is, even far detuned from the single-photon resonances, the two-photon resonance condition alone is sufficient to drive large coherences between states $|1\rangle$ and $|2\rangle$. We also stress the point that in general in CPT the initial population may be distributed between the lower states $|1\rangle$ and $|2\rangle$. This is usually the case if the lower states are provided by sublevels of an atomic ground state, for example, for Zeeman or hyperfine split states. As the energy difference of the lower states $|1\rangle$ and $|2\rangle$ is very small in this case, the initial thermal populations will be almost the same. In this case, a careful analysis of states $|NC\rangle$ and dark state $|C\rangle$ is required to determine the population dynamics for the specific CPT process under consideration, for example, STIRAP. However, also in the case of an initial population in both lower states, the state $|NC\rangle$ is still a dark state. In contrast, in implementations of EIT, the population is initially and for all times completely stored in state $|1\rangle$.

In CPT, interference effects arise from both coupling fields, since they are of comparable strength. If only one of the fields is strong, that is, $\Omega_1 \ll \Omega_2$, only interference effects due to processes driven by Ω_2 will be important. This is the situation in many implementations of EIT and is discussed by a number of authors (see, e.g., Ref. 5 and references therein). In such EIT experiments, Ω_2 is usually called the coupling field and labeled Ω_C , and Ω_1 is a weaker probe field, labeled Ω_p . In the following we will use the notations Ω_C and Ω_p , whenever appropriate. Based on the considerations, presented above, we will now discuss some simple and straightforward alternative approaches to understand the basic concept of EIT.

First, let us consider the basis, formed by the coupled state $|C\rangle$ and the noncoupled state $|NC\rangle$ [see Eq. (1)]. We can write the bare state $|1\rangle$ in this basis:

$$|1\rangle = \frac{1}{\Omega} (\Omega_C |NC\rangle + \Omega_p |C\rangle) \quad (2)$$

Very obviously, for the case $\Omega_p \ll \Omega_C$ state $|1\rangle$ is almost equivalent to $|NC\rangle$, that is, the dark state. Thus absorption to state $|3\rangle$ vanishes. The population remains in the ground state $|1\rangle$ throughout

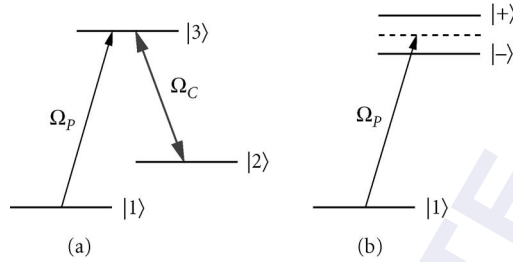


FIGURE 3 EIT in a lambda scheme (compare Fig. 4) viewed in terms of (a) bare atomic states, driven by a weak probe field Ω_p and a strong coupling field Ω_C and (b) dressed states, generated by the strong coupling Ω_C . The probe field is still at the bare state resonance frequency ω_{13} .

the interaction with the two radiation fields. The probe laser propagates through the medium without any absorption losses, that is, the medium is driven to EIT.

Alternatively, as $\Omega_p \ll \Omega_C$, we can treat the three-level system of states $|1\rangle$, $|2\rangle$, and $|3\rangle$ as a composition of a strongly coupled two-level subsystem of states $|2\rangle$ and $|3\rangle$, with the weakly coupled state $|1\rangle$ attached to the subsystem. Thus it is straightforward to describe the subsystem in terms of the dressed states, which arise due to the strong interaction² (see Fig. 3). For a strong resonant coupling at the single-photon resonance $\Delta_{23} = 0$ the dressed states of the subsystem are

$$|+\rangle = \frac{1}{\sqrt{2}}(|2\rangle + |3\rangle) \quad (3a)$$

$$|-\rangle = \frac{1}{\sqrt{2}}(|2\rangle - |3\rangle) \quad (3b)$$

The transition amplitude at the (undressed) resonant frequency ω_{13} from the ground state $|1\rangle$ to the dressed states will be the sum of the contributions to the dressed states $|+\rangle$ and $|-\rangle$, that is, $[\langle 1 | \mu | + \rangle + \langle 1 | \mu | - \rangle] \sim [\langle 1 | \mu | 2 \rangle + \langle 1 | \mu | 3 \rangle + \langle 1 | \mu | 2 \rangle - \langle 1 | \mu | 3 \rangle] = \mu_{12} + \mu_{13} + \mu_{12} - \mu_{13}$, with the transition moments μ_{ij} . As we assumed in the definition of our three-level system, the transition between states $|1\rangle$ and $|2\rangle$ is forbidden, that is, $\mu_{12} = 0$. The transition moments μ_{13} enter the sum with opposite signs. Thus the transition amplitude, that is, the absorption, reduces exactly to zero. The system is driven now to EIT.

Finally, another approach to an understanding of EIT is based on the concept of quantum interferences.^{13,69-71} Interference, associated with EIT, arises because the transition amplitude between states $|1\rangle$ and $|3\rangle$ includes different pathways from the ground state $|1\rangle$ to the excited state $|3\rangle$: One term, which is due to excitation by the resonant field Ω_p only, that is, a direct path from state $|1\rangle$ to state $|3\rangle$; an additional term, which is due to the presence of the second field Ω_C (see Fig. 4), that is, a indirect path from state $|1\rangle$ to $|3\rangle$ further on to state $|2\rangle$ and back to $|3\rangle$. The additional term and similar higher-order terms have a negative sign with respect to the direct path. Hence the higher-order terms cancel completely the direct path. This situation is closely related to interferences, mediated at Fano-type resonances,¹ for example, via autoionizing states, or to laser-induced continuum structure.⁴⁶⁻⁵⁰

Equivalently within the picture of EIT in terms of the interfering pathways between the bare atomic states, the coherences are the quantities pertinent to the interference. Coherences can be thought of, in a semiclassical picture, as associated with the oscillating electric dipoles driven by the coupling fields applied between pairs of quantum states of the system, for example, between states $|i\rangle$ and $|j\rangle$. Strong excitation of these dipoles occurs whenever electromagnetic fields are applied close to resonance with an electric dipole transition between two states. If there are several pathways

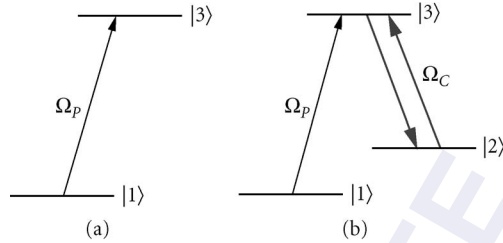


FIGURE 4 EIT in a lambda scheme, viewed in terms of Fano-type interference (compare Fig. 3), (a) shows the direct channel for the excitation $|1\rangle \rightarrow |3\rangle$ by the probe field Ω_p , (b) shows the lowest-order multiphoton channel induced by the coupling field Ω_c that is, the sequence $|1\rangle \rightarrow |3\rangle \rightarrow |2\rangle \rightarrow |3\rangle$. Interference between pathway (a) and (b) (also including higher-order terms) occurs.

to excite the oscillating dipole at frequency ω_{ij} , then interference arises between the various contributions to this dipole. These contributions must be summed to give the total amplitude to the electric dipole oscillation (see Fig. 4).

Mathematically, coherences are identified with the off-diagonal density matrix elements ρ_{ij} . These are formed by taking bilinear combinations of probability amplitudes of two quantum states, that is, by the weighting factors associated with the outer products such as $|i\rangle\langle j|$.⁷² Off-diagonal elements in the density matrix play a critical role in the evolution of an atom coupled to electro-magnetic fields.⁷³ Many calculations of atomic coherence effects and of EIT, as well as general nonlinear optics and laser action, in three-level systems are therefore developed in terms of the density matrix. The magnitudes of the relevant density matrix elements are computed from the basic coupled evolution equations, that is, the Liouville equation,^{10,72} and are found to depend upon experimental parameters (e.g., detunings and laser intensities). This approach also naturally lends itself to the inclusion of dampings that cause the decay of populations and coherences (e.g., radiative decay and collisions).⁶⁷

EIT serves to control the absorption of a probe laser on the transition between the states $|1\rangle$ and $|3\rangle$ in the three-level systems, defined and discussed above. Thus EIT will manifest itself in the density matrix element ρ_{13} . The real and imaginary parts of ρ_{13} both vanish at zero detuning, that is, the coherence is cancelled by interference of the excitation pathways. A set of coupled equations connecting the density matrix elements (e.g., ρ_{12} , ρ_{23} , and ρ_{13}) and their temporal derivatives is deduced from the Liouville equation. These coupled equations are then solved for various sets of conditions by either analytical or numerical means. Interference that leads to EIT arises from the coherences ρ_{23} and ρ_{12} , which are coupled to ρ_{13} . The coherence ρ_{12} between the ground state $|1\rangle$ and state $|2\rangle$ is present only due to the additional laser coupling. The contribution to the coherence ρ_{13} from the coherences ρ_{23} and ρ_{12} (driven by both laser fields) cancels with the direct contribution (driven by the probe laser field alone).

Although this use of density matrix elements is convenient, it is by no means essential, and many theoretical treatments that give clear physical insight have been performed in terms of probability amplitudes. Additional physical insights have been obtained by adopting alternative approaches, for instance by a careful consideration of the Feynman diagrams representing the various processes involved that lead to interference,^{69–71} or by applying a quantum jump approach.⁷⁴ In all cases the predictions are essentially identical.

Analytical solutions are generally only possible for steady-state conditions [corresponding to continuous wave (cw) laser fields]. A time-dependent calculation of the density matrix is appropriate, if laser pulses rather than cw radiation drives the atom. A time-dependent calculation is also vital to account for transient effects or pulse propagation. Some analytical solutions also for time-dependent calculations of the density matrix have been obtained, but unless restrictive simplifying assumptions are applied,⁷⁵ time-dependent calculations must be performed numerically. In many

cases the results of a full time-dependent treatment will be comparable to the results, obtained in a steady-state approach. This holds true at least in so far, as qualitative trends are concerned. To calculate the propagation of laser pulses through an extended medium the time-dependent equations for the density matrix elements must be coupled to Maxwell's equations.¹⁰ This is necessary, for example, to compute the propagation of matched pulses,^{76–78} the efficiency of frequency conversion processes in coherently prepared media,^{20–27,28–39} to account for losses in the driving laser fields or to model pulse shape modifications.⁷⁹

14.5 MANIPULATION OF OPTICAL PROPERTIES BY ELECTROMAGNETICALLY INDUCED TRANSPARENCY

Any optical process in a medium, driven by radiation pulses, is determined by the polarization, that is induced by the light fields. The macroscopic polarization P at the transition frequency ω_{13} can be related to the microscopic coherence ρ_{13} via the expression

$$P_{13} = n_A \mu_{13} \rho_{13} \quad (4)$$

where n_A is the number of atoms per unit volume in the ground state within the medium, and μ_{13} is the dipole matrix element associated with the (undressed) transition.⁷³ In this way imaginary and real parts of the linear susceptibility at frequency ω can be directly related to ρ_{13} via the macroscopic polarization.⁵⁵ The latter is defined as

$$P_{13}(\omega) = \epsilon_0 \chi(\omega) E \quad (5)$$

introducing the susceptibility $\chi(\omega)$. In this paper, the microscopic coherences are treated quantum mechanically, while the electromagnetic fields are treated classically (i.e., using Maxwell's equations and the susceptibilities). This semiclassical approach is not essential, and fully quantum mechanical treatments for CPT (see, Ref. 41) and EIT^{78–80} have been developed. These fully quantum approaches are appropriate for cases such as the coupling of atoms to modes in cavities,^{80,81} or when the statistical properties of the light fields are of interest. The latter is the case, for example, in proposals to generate squeezed light using EIT.⁷⁹ For relatively strong light fields, present in most laser experiments, a semiclassical treatment, with spontaneous decay added as a phenomenological damping process, proves adequate.

The real and imaginary parts of the (dressed) linear susceptibility, associated with the dispersion and absorption of the medium respectively, are given by^{5,6}

$$\text{Re} \chi_D^{(1)}(-\omega_p, \omega_p) = \frac{|\mu_{13}|^2 n_A}{\epsilon_0 \hbar} \left[\frac{-4\Delta_{21}(|\Omega_C|^2 - 4\Delta_{21}\Delta_{32}) + 4\Delta_{31}\Gamma_2^2}{(4\Delta_{31}\Delta_{21} - \Gamma_3\Gamma_2 - |\Omega_C|^2)^2 + 4(\Gamma_3\Delta_{21} + \Gamma_2\Delta_{31})^2} \right] \quad (6a)$$

$$\text{Im} \chi_D^{(1)}(-\omega_p, \omega_p) = \frac{|\mu_{13}|^2 n_A}{\epsilon_0 \hbar} \left[\frac{8\Delta_{21}^2\Gamma_3 + 2\Gamma_2(|\Omega_C|^2 + \Gamma_3\Gamma_2)}{(4\Delta_{31}\Delta_{21} - \Gamma_3\Gamma_2 - |\Omega_C|^2)^2 + 4(\Gamma_3\Delta_{21} + \Gamma_2\Delta_{31})^2} \right] \quad (6b)$$

To deduce the linear susceptibility, monochromatic fields and negligible collisional or Doppler broadening were assumed. The real and imaginary parts of the linear susceptibility, along with the nonlinear susceptibility, are plotted in Fig. 5 (see figure caption for explanation of the labeling) as a function of the detuning Δ_{13} , at $\Delta_{23} = 0$, that is, for resonant excitation by the coupling field Ω_C . A striking result, the absorption for the probe field vanishes at exact resonance, if the coupling field Ω_C is switched on and state $|2\rangle$ is perfectly metastable (i.e., $\Gamma_2 = 0$). Simultaneously the dispersion is significantly modified. The dispersion is still zero at line center, that is, the same value as in the case of the coupling field Ω_C switched off. However, the group velocity (dependent upon the slope of $[\text{Re} \chi^{(1)}]$) becomes anomalously low^{5,82} when absorption has vanished. This offers the possibility of

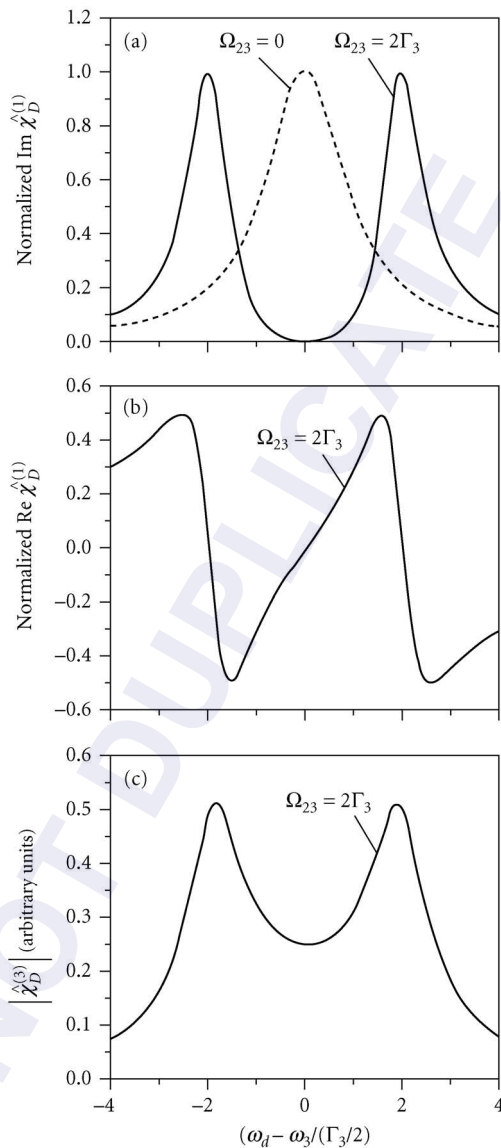


FIGURE 5 The dressed susceptibilities in terms of the normalized detuning $(\omega_d - \omega_3)/(\Gamma_3/2)$ [in our notation this corresponds to the detuning of the probe laser, scaled to the decay rate, i.e., $\Delta_{13}/(\Gamma_3/2)$] for a value of the coupling Rabi field $\Omega_{23} = 2\Gamma_3$ [in our notation $\Omega_C = 2\Gamma_3$]. (a) $[\text{Im } \chi^{(1)}]$, i.e. absorption; (b) $[\text{Re } \chi^{(1)}]$, i.e., dispersion; and (c) $|\chi^{(3)}|$, i.e., nonlinear response. [Reprinted figure with permission from Ref 6. Copyright (1990) by the American Physical Society.]

slowing down the speed of light by EIT, that is, a most prominent example for the striking features of EIT, which is discussed elsewhere in this book. Inclusion of finite laser linewidths, collisional and Doppler broadening in these deductions is straightforward.⁶ Including these effect, the dressed susceptibilities still retain the key features, provided some limits for the experimental parameters are kept in view.

The medium that would in the absence of the coupling field be optically thick is now rendered transparent. The reduction in absorption is not merely that caused by the effective detuning induced by the Autler-Townes splitting of the dressed state absorption peaks (see Fig. 3), that is, the absorption that would be measured if the probe field were interrogating the absorption coefficient of the medium in the wing of the absorption lines of the dressed states. Additionally, there is destructive interference at the transition frequency ω_{13} that leads to *complete* cancellation of all absorption, provided there are no additional dephasing processes in the system. Even if the transition dipole moment μ_{12} is not zero (i.e., if the spontaneous decay rate is $\Gamma_2 \neq 0$), the absorption due to EIT will be reduced compared to the weak field absorption by a factor of Γ_2/Γ_3 .⁶

In the preceding considerations it was implicitly assumed that the probability amplitudes of state $|3\rangle$ remain close to zero (i.e., the probe field is very weak). If there is an incoherent population pump process into the upper states $|2\rangle$ or $|3\rangle$, such that these populations no longer remain negligible, then gain on the transition between the states $|3\rangle$ and $|1\rangle$ can result (see Sec. 14.8). The remarkable feature of this gain is that under the circumstances in which EIT occurs, that is, when absorption is cancelled, the gain can be present without the requirement of population inversion in the bare atomic states. This is an example of amplification without inversion. The process has successfully been implemented in a vee-type scheme in Rubidium atoms¹⁵ and a lambda scheme in sodium atoms.⁸³ Much theoretical work^{12-14, 84-86} has been reported on LWI. In early work, LWI was predicted by Arkhipkin and Heller,⁸⁵ and then further elucidated by Harris,⁸⁶ Kocharovskaya and Khanin,¹² and Scully et al.⁸⁴ A long-term objective in this work is the prospect of overcoming the familiar difficulties of constructing short wavelength lasers, that is, the requirement of very high energy pump fields in order to compensate for the small transition moment at far-infrared wavelengths.

In addition to gain without population inversion, any incoherent pumping of population into the upper states also modifies the dispersion in the medium. In particular, it is then possible to obtain a very large refractive index for specific wavelength regimes. The refractive index can reach values comparable to those normally encountered very close to an absorption line, while here the absorption now vanishes.⁸ The prospects for engineering the refractive properties of media to give novel combinations of absorption, gain, and dispersion have been explored in a number of theoretical⁸⁷⁻⁹³ and experimental studies (see Sec. 14.9).

The successful implementation of EIT depends upon a number of critical parameters, both inherent to the quantum system and the experimental setup, for example, the driving laser pulses. A correct choice of the atomic energy level configuration is essential. The configuration must satisfy the conditions, already discussed above, that is, dipole allowed transitions $|1\rangle \rightarrow |3\rangle$ and $|2\rangle \rightarrow |3\rangle$, while the transition $|1\rangle \rightarrow |2\rangle$ is dipole forbidden. Radiative couplings to other energy levels outside of these states, that lead to an open three-level system, and additional level substructure must also be considered. Collisions with other species in the medium or photoionization must be minimized in order to prohibit perturbing decay or dephasing of the coherence ρ_{12} , which is essential to EIT. The couplings may be either driven by cw or pulsed lasers. In both cases the couplings must be sufficiently strong to overcome the inhomogeneous broadening. Moreover, sufficiently monochromatic or radiation with transform-limited bandwidth in the pulsed case is required, in order not to dephase the coherence ρ_{12} . These critical parameters are summarized in Table 1. In the following section we discuss in more detail the most crucial parameters for a successful implementation of EIT.

Intrinsic Dephasing of Atomic Coherence in Gas Phase Media

For processes of laser-induced atomic coherence in a realistic medium the maintenance of the phase of the coherence during the interaction is essential for effective quantum interference. Any dephasing of the coherence will wash out and eventually nullify the interference effects. Dephasing can

TABLE 1 Summary of Critical Experimental Parameters for a Successful Implementation of EIT

Physical Parameter	Constraint	Typical Values
Radiative decay rate Γ_3 of state $ 3\rangle$	–	1–100 MHz
Radiative decay rate Γ_2 of state $ 2\rangle$	$\Gamma_2 \ll \Gamma_3$	<1 MHz
Photoionization rate Γ_{ion}	$\Gamma_{\text{ion}} \gg \Gamma_3$	Depends upon the laser intensity
Coherence dephasing γ_{ij}	$\gamma_{12} < \gamma_{13}, \gamma_{23}$	0–1 MHz (gases), 0–100 GHz (solids)
Laser linewidth γ_{laser}	Transform-limited	<1 MHz (cw) or $1/\tau_{\text{pulse}}$ (pulsed)
Doppler linewidth γ_{Doppler}	$\gamma_{\text{Doppler}} < \Omega_C$	1 GHz
Rabi frequency Ω_C	$\Omega_C > \gamma_{12}$	$\Omega_C = \mu_{23} E_C / \hbar \nabla$
Laser pulse energy E_C [see Eq. (8)]	$E_C > f_{13} / f_{23} \hbar \omega N_A$	

arise from a variety of different sources, for example, the excitation of a multitude of closely spaced hyperfine or Zeeman components (see, e.g., Refs. 10,15,94), radiative decay of state $|2\rangle$, photoionization channels,⁹⁵ and collisions.^{6,96} Following these arguments, it is obvious, that in a ladder scheme perfect EIT is usually not possible, because state $|2\rangle$ is not metastable, but undergoes spontaneous emission to state $|3\rangle$ —as well as to other states outside the three-state system. Collisional broadening with atoms of the same and other species is also critical and places strict limits on sample purity, otherwise leading to foreign gas broadening. Possibly also limits are imposed on the sample density. Moreover, at large atomic density the local field effects due to dipole-dipole couplings between the atoms may be important.^{96,97} In this case the simple relationship between macroscopic polarization and the coherences in Eq. (5) may break down.

Dephasing due to Phase Fluctuations in the Laser Fields

If cw lasers with narrow bandwidth $\Delta\nu < 1$ MHz (i.e., negligible phase fluctuations) are used, the quality of an implementation of EIT usually approaches the expectations, based on steady-state conditions and strictly monochromatic excitation. Due to phase-diffusion processes^{92–102} laser line broadening gives rise to linewidths above the limit of the allowed radiative decay rate. This will destroy EIT. EIT may also be implemented with pulsed laser, provided the laser pulses exhibits transform limited bandwidth. Though these transform-limited bandwidth is inevitably larger than those of cw lasers, EIT is not reduced at all. A single-mode transform-limited laser pulse (i.e., without excess phase fluctuation) will introduce insufficient dephasing during the interaction time (i.e., the pulse duration τ_{pulse}) to disturb the phases of the atomic coherences. It should also be appreciated that hyperfine sublevels will in general cause dephasing of coherences on a timescale, which is given by the inverse of their frequency separation $\Delta\omega_{\text{HFS}}$.¹⁰ In a pulsed excitation the dephasing due to hyperfine splitting within the laser bandwidth will therefore be negligible, provided $\tau_{\text{pulse}} < 1/\Delta\omega_{\text{HFS}}$ (i.e., if the hyperfine splittings are sufficiently small).

Dephasing Processes in Solids

A few experiments on EIT have also been performed in solid-state media.^{103–115} In contrast to implementations of EIT in gaseous media, coherent interactions in solids suffer from additional dephasings. These dephasings are induced by quantum processes in the crystal lattice (e.g., by photon-phonon interactions). The dephasings lead to significant additional homogeneous broadening in spectral lines of optical transitions. At room temperature the broadening may reach the regime of many gigahertz. In the gas phase, homogeneous broadening (e.g., as observed in the natural linewidth of a transitions) is usually much smaller than inhomogeneous broadening [e.g., induced by Doppler shifts (see below)]. In contrast, in solids it is usually the homogeneous broadening which dominates. This homogeneous broadening (i.e., the dephasing) will wash out any quantum interferences. In principle, there are two ways to deal with the fast dephasing processes in solids: (1) to “freeze” phonon processes, that is, to cool

the medium to cryogenic temperatures such that phonon interactions may be neglected; (2) to drive the medium with ultrafast laser pulses, that is, pulse durations below the timescale of the dephasing processes. However, ultrashort laser pulses usually do not provide sufficient pulse energy to saturate the transitions in EIT. Thus, in most experiments on EIT in solids, cryogenic cooling of the medium is a necessary requirement.

Inhomogeneous Broadening

In many experiments (inhomogeneous) Doppler broadening presents a serious limit since it introduces a randomization in the effective laser detunings over the ensemble of atoms in the sample.^{6,116,117} Various methods have been employed to eliminate this effect, for example, Doppler-free excitation^{118,119} or using cooled atoms in a magneto-optical trap.^{120,121} Alternatively, at coupling Rabi frequencies larger than the Doppler width the influence of inhomogeneous broadening can be overcome.⁶ In spectroscopic terms, the laser-induced power broadening beats the Doppler broadening in this case. This concept requires Rabi frequencies exceeding the Doppler width, that is, typically in the regime of 1 GHz. Such Rabi frequencies can be provided by a cw laser, with a typical power of 1 W or less, only under conditions of tight focusing. This may lead to undesirable effects such as defocusing due to the interplay between the dressed refractive index and the transverse intensity variation across the beam waist in the region of the focus.¹¹⁹ For pulsed lasers, with intrinsically high peak power, it is usually not necessary to focus the laser to reach the required intensity. Thus defocusing effects are negligible and the interaction volume (i.e., the atom number) will be much larger than under conditions of focusing. Moreover, for short- or ultrashort laser pulses the laser bandwidth already exceeds the Doppler width. However, short- or ultrashort laser pulses are usually no good choice to drive EIT (see below).

Inhomogeneous broadening also plays a major role for applications of EIT in solid-state systems. Such inhomogeneous broadenings are induced, for example, when doped atoms in a crystal experience different electric fields in the background of the host crystal. This leads to an inhomogeneous distribution of transition frequencies. In general, the implementation of EIT in inhomogeneously broadened solids requires large laser intensities. This increases the risk of optical damage to the sample. An exception to this is the work on solids under conditions of cryogenic cooling. In such media, exceptionally small inhomogeneous widths are encountered.¹²² Moreover, inhomogeneous broadening in doped solids can be overcome by appropriate optical preparation (e.g., spectral hole burning). In this case, a specific ensemble of atoms in the inhomogeneously broadened medium is prepared to exhibit spectrally narrow transitions.

Coupling Laser Power

In addition to the conditions, described above, there are some more constraints for the coupling laser Rabi frequency: (1) The coupling laser Rabi frequency must be sufficiently large to induce a transparency with a spectral width exceeding the linewidth of the Raman (two-photon) transition. For Raman transitions involving large detunings from the single-photon resonances, this may require large laser powers. (2) Harris and Luo⁷⁹ derived a condition for the laser pulse energy, which demands a sufficient number of photons in the laser pulse to match the number of atoms, weighted by the transition oscillator strength, in the laser path [see Eq. (8)]. Essentially, the number of photons in the coupling laser pulse must exceed the number of atoms in the medium. (3) In the adiabatic limit, the pulse durations must exceed the time evolution of the transparency, which is in the order of $1/\Omega_C$. Thus, the product $\Omega \cdot \tau_{\text{pulse}}$ must be large, that is, $\Omega \cdot \tau_{\text{pulse}} \gg 1$. In terms of incoherent excitation, the laser must “saturate” the transition. This adiabaticity criterion can be derived in a very similar form also for other adiabatic processes.⁴⁴ As $\Omega \cdot \tau_{\text{pulse}}$ is a combination of the electric field and the pulse duration, laser pulses with large intensity and/or large interaction time are a good choice to fulfill the adiabaticity criterion. An analysis of typical laser systems with specific pulse duration reveals that laser pulses with medium pulse duration [i.e., in the regime of short nanosecond (ns)

or long picosecond (ps) pulses] yield the largest product $\Omega \cdot \tau_{\text{pulse}}$, if one-photon transitions are driven. Shorter laser pulses usually cannot compensate for the reduced pulse duration by a sufficient increase in the electric field. On the other hand, long pulses or cw radiation permits for long interaction time, but the electric field is weak. Thus, laser pulses with intermediate pulse duration and pulse energies in the regime of mJ are the best choice. However, if multiphoton excitations or specific atomic systems with large transition moments are considered, also ultrashort laser pulses may drive adiabatic interactions.

14.6 ELECTROMAGNETICALLY INDUCED TRANSPARENCY, DRIVEN BY PULSED LASERS

In the earliest work on EIT, driven by pulsed lasers, the *linear* optical response of an extended ensemble of atoms in the gas phase was investigated. In these experiments the transmission of a weak probe laser pulse, propagating through an otherwise optically dense medium, was measured. The medium was rendered transparent in the presence of a strong coupling laser pulse. In pulsed laser experiments there is usually no difficulty to induce Rabi frequencies, which exceed the inhomogeneous bandwidth of the medium (i.e., to drive the complete medium in EIT). It is essential, however, that the laser pulse exhibits transform-limited bandwidth. Thus, for example, single-mode transform-limited nanosecond lasers are an appropriate choice (see Sec. 14.5). Such laser pulses are, for example, provided by injection seeding an amplifier or optical parametric oscillator (OPO) with narrow-bandwidth cw radiation of appropriate frequency.

The first demonstration of EIT, driven with pulsed lasers, was performed by Harris et al. in strontium vapor¹²³ and lead vapor.¹²⁴ In both experiments laser pulses with transform-limited radiation were used. In the experiment on strontium (see Fig. 6) the atoms are initially optically pumped into the excited state $5s\ 5p\ ^1P_1$. The transition between state $|1\rangle = 5s\ 5p\ ^1P_1$ and the autoionizing state $|3\rangle = 4d\ 5d\ ^1D_2$ at a wavelength of $\lambda_p = 337.1\ \text{nm}$ is rendered transparent. A coupling laser, derived from a single-mode Littman dye laser, at wavelength $\lambda_c = 570.3\ \text{nm}$ drove the transition between the metastable state $|2\rangle = 4d\ 5p\ ^1D_2$ and $|3\rangle = 4d\ 5d\ ^1D_2$. As in the prototypical scheme, the probe field excited the system to a state $|3\rangle$ with a large decay rate. In the absence of the coupling laser the probe laser experienced strong absorption. Thus the strontium vapor was completely opaque at resonance, with an inferred transmission of $\exp(-20 \pm 1)$. When the coupling laser was applied, the transmission at line center increased dramatically to $\exp(-1 \pm 0.1)$. It was pointed out by the authors that for this large transparency the interference effect is essential. The detuning from the (dressed) absorption lines, separated by the Autler-Townes splitting alone would only account for an increase in transmission to a value of $\exp(-7.0)$.

The experiment in lead vapor¹²⁴ demonstrated EIT within the bound states of a medium, experiencing significant collisional broadening. Here EIT was implemented in a ladder configuration with the probe laser, driving the transition between the ground state $|1\rangle = 6s^2\ 6p^2\ ^3P_0$ and the excited state $|3\rangle = 6s^2\ 6p\ 7s\ ^3P_1^0$. The coupling laser drove the transition between states $|2\rangle = 6s^2\ 6p\ 7p\ ^3D_1$ and $|3\rangle = 6s^2\ 6p\ 7s\ ^3P_1^0$. The reduction in opacity, driven by the transform-limited coupling laser, reached a factor of $\exp(-10)$. The particular coupling scheme in lead was chosen because of approximate coincidence between the frequency of an injection seeded Nd:YAG laser at $\lambda_c = 1064\ \text{nm}$ and the transition frequency ω_{23} . The detuning was $\Delta_{23} = 6\ \text{cm}^{-1}$. An important feature of this experiment was the role of resonance broadening, which was the dominant broadening channel for state $|3\rangle$, that is, about 40 times larger than the natural linewidth. Due to the destructive interference between the contributions to state $|3\rangle$ in the two dressed states (see Sec. 14.4), these collisions have no effect on transparency. In contrast, the collisions, which dephase state $|2\rangle$, also affect the degree of EIT. As these are no resonance collisions, their strength is small.

Both experiments served to demonstrate the principle of EIT in a three-level system. They show, that EIT also works in systems including autoionization or collisional broadening. In both cases the coupling laser exhibited near transform-limited bandwidth. No special requirements were imposed on the probe laser, although the probe laser bandwidth must be less than the spectral width of the EIT.

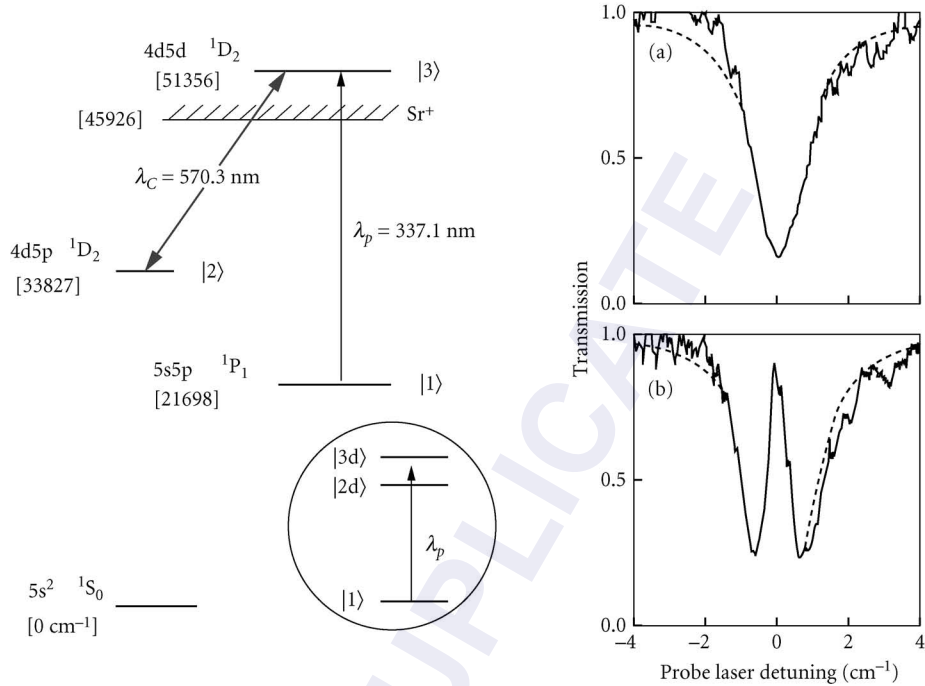


FIGURE 6 Coupling scheme for EIT in strontium atoms and transmission of a probe laser pulse versus the probe laser detuning. Inset in the coupling scheme: dressed states [in our notation state $|-\rangle$ corresponds to $|2d\rangle$ and state $|+\rangle$ corresponds to $|3d\rangle$; compare Fig. 3]. (a) When the coupling laser is switched off, the probe laser experiences absorption in the line center. (b) When the coupling laser is switched on, the probe laser absorption is dramatically reduced in the line center. [Reprinted figure with permission from Ref. 123. Copyright (1991) by the American Physical Society.]

Experiments using pulsed lasers continue to be important, most especially in the context of nonlinear optics and matched pulse propagation. A related resonant EIT scheme in lead has been explored by Kasapi as a technique for enhanced isotope discrimination.¹²⁵ This method utilized the resonant opacity of a low abundance lead isotope ^{207}Pb at EIT for the most common lead isotope ^{208}Pb in their sample. Additional work has also recently illustrated how EIT can be established in the lead isotope ^{207}Pb despite the presence of hyperfine structure. This was done by adjusting the laser frequencies to coincide with the center of gravity of the hyperfine transitions.^{126,127} Under this condition, interference of the manifold of hyperfine states yields EIT.

14.7 STEADY STATE ELECTROMAGNETICALLY INDUCED TRANSPARENCY, DRIVEN BY CW LASERS

Continuous wave (cw) lasers permit the implementation of EIT under steady-state conditions. Such experiments provide an excellent case to test the theoretical concept of EIT against experimental data from excitations under near-ideal conditions. Investigations of new effects in the cw regime permit straightforward comparison with theoretical predictions. A monochromatic laser is required, with a linewidth significantly less than the radiative decay rate Γ_3 (i.e., in the range of 10 kHz to 10 MHz). Such radiation is typically provided by either dye or titanium sapphire ring lasers or more cheaply,

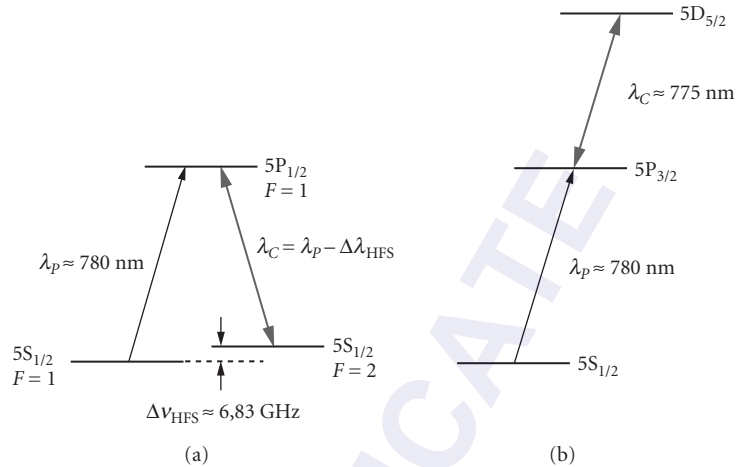


FIGURE 7 Coupling schemes for EIT in rubidium atoms. (a) A lambda scheme involving the hyperfine sublevels of the ground state [which correspond to states |1] and |2] in our notation] of ^{87}Rb (or ^{85}Rb) and the excited states $5P_{1/2}$ or $5P_{3/2}$ [which correspond to state |3] in our notation] and (b) a ladder-type scheme involving the ground state $5S_{1/2}$ [i.e., state |1] in our notation], and the excited states $5P_{3/2}$ [i.e., state |3] in our notation] and $5D_{5/2}$ [i.e., state |2] in our notation].

but also with limited tunability and power, by external cavity stabilized laser diodes. However, in contrast to pulsed laser experiments, in cw EIT experiments it is more difficult to reach sufficient coupling strengths Ω_C in order to exceed the inhomogeneous broadening. This has required the employment of both Doppler-free techniques and the reduction of the Doppler width by atom cooling methods.

Much work has been carried out in rubidium vapor. This is due to suitability of the rubidium atom's energy level configuration for EIT (see Fig. 7), the possibility of near-complete elimination of Doppler shifts in certain configurations, and the ease of handling the vapor. Xiao et al. as well as Moseley et al. performed important demonstrations of steady state EIT in rubidium atoms. A near-ideal lambda scheme is formed in rubidium (see Fig. 7a) between the ground states $5S_{1/2}$ ($F=1$) and ($F=2$) and the excited state $5P_{1/2}$ ($F=1$) state. The transitions in this case are around 780 nm, separated by the ground state hyperfine splitting of $\Delta\nu_{\text{HFS}} \approx 6.83$ GHz. In a detailed theoretical treatment of this system it is necessary to include all hyperfine sublevels of the three states. However, in essence the behavior is that of a simple three-level system. Likewise, also a ladder scheme is formed in this atom between the states $5S_{1/2}$, $5P_{3/2}$, and $5D_{5/2}$ (see Fig. 7b), with the transition wavelengths at $\lambda_p = 780$ nm and $\lambda_C = 775$ nm. Excitation of the transition at these wavelengths is easily possible using either cw titanium sapphire ring lasers or grating stabilized laser diodes. These schemes in rubidium also exhibit the additional advantage of Doppler-free excitation with the two closely spaced wavelengths. In such experimental configuration, Xiao et al. and Moseley et al. studied transparency,^{116,128,129} refractive index modification,¹¹⁸ and propagation effects.^{119,130}

Observation of EIT in the rubidium ladder scheme showed good agreement with a steady-state calculation, involving residual effects of inhomogeneous (Doppler) broadening.¹¹⁶ Due to the near-frequency coincidence between the transitions at 780 and 775 nm in this scheme, the effect of inhomogeneous broadening on the experiment was almost eliminated by using counterpropagating beams. This elimination of inhomogeneous broadening permitted the application of grating-stabilized laser diodes, operating with relatively low power ($P < 10$ mW), to provide both probe and coupling fields in these experiments.

Results of experiments on the rubidium lambda scheme also proved to be consistent with theory. EIT was observed at the probe transition line center with a linewidth and depth in reasonable

agreement with the steady-state calculation.¹²⁸ In the case of the lambda scheme, copropagating beams lead to a Doppler-free excitation. Again this was possible due to the closely spaced probe and coupling laser wavelengths. As mentioned before, also here only low laser powers were required for EIT. This elimination of Doppler broadening was exploited by the same authors to perform an experiment that stresses the quantum interference nature of the EIT effect. By employing a coupling laser strength $\Omega_C < \Gamma_3$ (i.e., the Rabi splitting is too small to give rise, on its own, to any significant absorption reduction) a well-developed transparency, with a depth limited only by the laser linewidths, was reported.¹²⁹ This experiment illustrates clearly how the additional coherence due to the coupling laser causes interference that cancels the effect of probe absorption.

The limits for a successful implementation of steady-state EIT with respect to the probe laser power has also been examined in the rubidium ladder scheme, as discussed above.¹³¹ When the coupling laser was strong, but the probe relatively weak, EIT was induced as usual. In the case of a strong probe with a power comparable to the coupling laser ($\Omega_p \sim \Omega_C$), the EIT was destroyed and replaced by enhanced absorption (i.e., electromagnetically induced absorption). This is explained by the opening of additional pathways in the absorption process (due to higher-order interactions with the probe field) leading to constructive interference in absorption. This result has implications in certain nonlinear frequency mixing schemes where a strong field should be generated at the probe frequency. The results suggest that there may be a limitation to the strength of the generated fields (see Sec. 14.12).

Besides operation in Doppler-free excitation schemes, also laser cooling of the atoms can be used to fully eliminate Doppler broadening. This also enables steady-state EIT in systems where the laser frequencies ω_c and ω_p differ significantly. Recently work has been reported on EIT (and CPT) in lambda systems in cold rubidium and cesium atoms, confined in a magneto-optical trap (MOT).^{120,121} In rubidium a coupling scheme involving Zeeman splittings has also been studied.¹²⁹ These systems are close to ideal as the trapped atoms are very cool, that is, Doppler broadening is almost absent. Moreover the system exhibits very low density and hence may be considered as collisionless. If required, larger densities of trapped atoms can be provided by using dark-spot trap techniques.¹³² Work has been carried out that exploits the characteristics of cold, confined rubidium atoms to study nonlinear absorption and dispersion.¹³³ Also temporal evolution of EIT in the transient regime has been investigated in a MOT with rubidium atoms.¹³⁴ In cold cesium atoms in a MOT the nonlinear sum rule¹³⁵ for EIT-type situations has been experimentally verified.¹³⁶

14.8 GAIN WITHOUT INVERSION AND LASING WITHOUT INVERSION

We consider now the case of population, transferred via an incoherent pump process into the excited states of a three-level system, driven in EIT. Thus a small population in the upper state of the probe transition results in inversionless gain of the probe field. The reader is referred to the extensive literature on this subject for further discussion (see, Ref. 16). In studying gain without inversion, precautions should be taken to confirm that the system is indeed noninverted. This is in practice rather difficult to confirm. Experimenters must verify that there is clear evidence for gain on the probe transition and truly no population inversion on this transition.

Gao et al.¹³⁷ performed an early experiment that stimulated discussion and subsequent work. A four-level system of the hyperfine ground states $3S_{1/2}$ ($F = 0, 1$) and the excited states $3P_{3/2}$ and $3P_{1/2}$ in sodium atoms, driven in a Raman process, was investigated. A laser pulse excited the transition between the ground states $3S_{1/2}$ ($F = 0, 1$) and state $3P_{1/2}$. This created a coherence between the two hyperfine ground states. A cw probe laser was tuned close to the transition between the ground states $3S_{1/2}$ ($F = 0, 1$) and state $3P_{3/2}$. The probe laser experienced amplification, when the probe laser frequency was appropriately tuned—and provided a small amount of population was pumped into the $3P_{3/2}$ state. The incoherent pumping process was driven by a DC gas discharge in the sodium vapor. The authors claimed, that the amplification process was inversionless. Initial criticism arose, as there was no independent monitoring of the excited state populations in this experiment.

Subsequent work on this issue, involving measurements of absorption, provided firmer evidence for inversionless conditions.¹³⁸ A similar excitation and amplification scheme in potassium vapor was reported by Kleinfeld et al.¹³⁹ The authors performed a careful numerical analysis that supports their claims for amplification without inversion. Further evidence for amplification without inversion has been found in several systems. In a lambda scheme in sodium atoms, with additional incoherent pumping into the upper state, Fry et al.¹⁴⁰ observed effects, based on atomic coherence, leading to amplification without inversion. The role of atomic coherence was confirmed by switching the coupling field on or off. Amplification was only observed when the field was present, and when it was absent the large population always present in the lowest state led to absorption of the probe. In another demonstration, picosecond pulses were used to excite atomic coherence among the Zeeman sublevels of the ground state in sodium atoms. Amplification without inversion was monitored and unambiguously confirmed by Nottelman et al.¹⁴¹ Amplification without inversion was also demonstrated in a transient scheme in cadmium vapor through the formation of a linear superposition of coherently populated Zeeman sublevels by van de Veer et al.¹⁴² In this experiment nanosecond laser pulses were used. The coherent nature of the process was proven (1) by the dependence of the gain on the time delay between the coherence preparation and probe pulse; and (2) by the dependence of the gain on the magnitude of the Zeeman splitting, which controlled the period for coherent transfer of population in the atom. Recently a double lambda scheme in helium atoms, driven by infrared radiation at 877.9 nm radiation in a helium-neon gas discharge, was used to observe amplification at both wavelengths 1079.8 and 611 nm (i.e., the latter in an up-conversion process). Here the evidence for amplification without inversion rests on comparison to calculation.¹⁴³

To demonstrate LWI, the gain medium must be placed within an optical cavity. In two experiments^{15,83} on this subject, amplification of a probe laser in the inversionless medium was demonstrated. Then the cavity was set up and lasing was observed even under conditions where no inversion was possible at all. The first of these experiments by Zibrov et al.¹⁵ was implemented in a vee-type scheme, formed on the D1 and D2 lines of rubidium, with incoherent pumping from the $F = 2$ hyperfine level into the upper state of the D1 transition. The latter provided the lasing transition. Laser diodes were used to derive all driving fields. The incoherent pump was generated by injecting white noise into an acousto-optic modulator (AOM). The AOM modulated one of the diodes. This work was also the first experiment to demonstrate amplification without inversion using laser diodes. An important conceptual advantage of the vee scheme is that there is no possibility for “hidden” inversion in a dressed basis at all. Thus the vee scheme serves as a very appropriate basis to demonstrate inversionless gain. The coupling scheme in the experiment could be considered in a simplified form as a four-level system, that is, three levels coherently coupled and a fourth coupled via the incoherent field. There are, however, 32 hyperfine sublevels in the particular experiment in rubidium, which must be considered in a detailed analysis. This complete analysis was carried out by the authors to yield predictions in good agreement with their experiment. Subsequently work was reported by Padmabandu et al.⁸³ demonstrating LWI in the same lambda scheme in sodium atoms, as also considered by the authors of Ref. 140.

14.9 MANIPULATION OF THE INDEX OF REFRACTION IN DRESSED ATOMS

Besides the manipulation of absorption or transmission, EIT also permits control of the refractive index of a laser-driven medium. Figure 5 shows the dependence of the real part of the linear susceptibility $[\text{Re } \chi^{(1)}]$, which determines the refractive index, on the probe laser detuning. The dispersion is most significantly modified between the absorption peaks, that is, for probe laser detunings in the range $\Delta_{13} = \pm\Omega_C/2$. In the absence of a coupling field the usual form of $[\text{Re } \chi^{(1)}]$, for an atom leads to anomalous dispersion, that is, a negative slope in $[\text{Re } \chi^{(1)}]$, in the vicinity of the resonance (i.e., for detunings $\Delta_{13} = \pm\Gamma_3$). The dispersion $[\text{Re } \chi^{(1)}]$ vanishes at exact resonance. This usual behavior of $[\text{Re } \chi^{(1)}]$ is not significant, since the medium is highly opaque in the frequency range close to the resonance. However, in EIT the absorption is nearly zero close to resonance. Thus the modified dispersion can have a large effect on the refractive properties of the medium.

In a medium, driven to EIT, the dispersion $[\text{Re } \chi^{(1)}] = 0$ at resonance. Thus the refractive index will attain the vacuum value ($n = 1$), while the medium is fully transparent. As the assumption of a (closed) three-level atomic system is only an approximation, there will always be a contributions to the total refractive index due to all other states of the atom. These other levels, however, are typically far from resonant with the driving laser fields. Thus they lead to relatively small contributions to the dispersion. In nonlinear frequency mixing the vanishing dispersion results in near-perfect phase-matching, which is essential for efficient frequency up-conversion (see Sec. 14.12). A second important modification to the dispersion in the frequency range $\Delta_{13} = \pm\Omega_C/2$ is that the medium shows *normal* dispersion, that is, a positive slope in $[\text{Re } \chi^{(1)}]$. The value of the dispersion in this frequency range will depend on the shape of the curve and hence on the coupling laser intensity, which determines Ω_C . In Sec. 14.11, we will examine the situation, when Ω_C is small. In this case the dispersion profile can be very steep, and consequently very low group velocities result. The intensity dependence of $[\text{Re } \chi^{(1)}]$ leads also to the strong spatial dependence of the refractive index across the intensity profile of a focused laser beam.

Investigations of the modification of dispersion (i.e., the refractive index) as induced by EIT have been carried out for steady-state excitation with cw lasers.^{118,119,130} Direct measurements of the modified refractive index confirm the theoretical predictions. The dispersive properties of rubidium atoms, driven by EIT, was investigated using a Mach-Zehnder interferometer.¹¹⁸ The dispersion measured at the center frequency was inferred to be equivalent to a small group velocity of $v_g = c/13.2$. In the last years, many experiments have been implemented in order to slow down the speed of light in media, driven to EIT. In this article we will only briefly mention some of the experiments, while a detailed discussion of slow light and storage of photons is subject to another article in this book.

A number of observations on electromagnetically induced focusing and defocusing, based on spatial variations in the index of refraction, have also been performed. A wavelength-dependent-induced focusing or defocusing was reported by Moseley et al.^{119,130} employing a coupling scheme in rubidium. Constraints are introduced to the tightness of focusing in strongly driven media. However, these experiments also indicate possibilities to control the spatial properties of a beam at frequency ω_p by a beam at another (i.e., perhaps very different, frequency ω_C).

In a scheme where a small amount of population is incoherently pumped into the excited state, gain was predicted⁸ at exact resonance (i.e., $\Delta_{13} = 0$). Thus, the imaginary part of the linear susceptibility $[\text{Im } \chi^{(1)}]$ becomes negative at resonance. When absorption vanishes at nearby frequencies, that is, $[\text{Im } \chi^{(1)}] = 0$, the value of $[\text{Re } \chi^{(1)}]$ can be very large. This situation is termed *enhanced refraction*. Enhanced refraction has been observed in a lambda scheme in rubidium, provided there is an additional pumping field to result in a small (noninverted) population in the upper state of the system.¹⁴⁴ In this experiment an enhanced index of refraction was found at frequencies where the absorption was zero. A proposed application for refractive index modifications of this kind is high-sensitivity magnetometry.¹⁴⁵ The large dispersion at the point of vanishing absorption could be used to detect magnetic level shifts via optical phase measurements in a Mach-Zehnder interferometer with high accuracy.

14.10 PULSE PROPAGATION EFFECTS

Propagation of pulses is significantly modified in the presence of EIT. Figure 5b shows the changes to $[\text{Re } \chi^{(1)}]$ in media, driven in EIT. An analysis of the refractive changes has been provided by Harris et al.⁸² who expanded the susceptibilities (both real and imaginary parts) of the dressed atom in a series around the resonance frequency to determine various terms in $[\text{Re } \chi^{(1)}]$. The first term of the series (zero order) $[\text{Re } \chi^{(1)}(\omega_{13})] = 0$ corresponds to the vanishing dispersion at resonance. The next term $\partial/\partial\omega [\text{Re } \chi^{(1)}]$ gives the slope of the dispersion curve. At the transition frequency ω_{13} this slope yields

$$\frac{\partial}{\partial\omega}[\text{Re } \chi^{(1)}] = \mu_{13}^2 \frac{4n_A}{\hbar\epsilon_0} \frac{(\Omega_C^2 - \Gamma_2^2)}{(\Omega_C^2 + \Gamma_2\Gamma_3)^2} \quad (7)$$

This expression shows the dependence of the slope of $[\text{Re } \chi^{(1)}]$ on Ω_C . The latter parameter permits controls of the group velocity of a laser pulse at frequency ω_{13} propagating through the medium. Higher-order terms in the expansion lead to pulse distortions (group velocity dispersion), but at resonance the lowest nonvanishing term is of third order.

The first experimental studies of pulse propagation in media, driven to EIT system, were conducted by Kasapi et al.¹⁴⁶ In a lambda scheme in lead vapor, Kasapi et al. measured the delay of a probe pulse for various coupling laser strengths. Whilst the transmission through the medium was still large (i.e., 55%), pulses were found to propagate with velocities as low as $v_g = c/165$. The authors showed that the delay time τ_{delay} for the pulse in the medium compared to a propagation through the vacuum was correlated to the attenuation of the transmitted pulse and the residual decay rate γ_{12} of the (dipole forbidden) transition $|1\rangle \rightarrow |2\rangle$ —transition via the relation $\ln\{E_{\text{out}}/E_{\text{in}}\} = -\gamma_{12} \cdot \tau_{\text{delay}}$, with E_{out} and E_{in} as the energies of the probe pulse leaving and entering the medium. This idea was subsequently demonstrated as a method for measuring Lorentzian linewidths.¹⁴⁷ It was also demonstrated¹⁴⁶ that the presence of the coupling pulse leads to a near-diffraction-limited transmitted beam quality for a strong probe field under conditions where severe spatial distortion was present in the absence of the coupling laser. The situation where probe and coupling pulses were both strong was studied in a subsequent experiment that investigated further the elimination of optical self-focusing and filamentation, which afflict a strong probe field.¹⁴⁸

The prospects of controlling the refractive index using strong off-resonant pulses was examined theoretically for a lambda scheme.⁹³ In this treatment it was shown that the off-resonant bound and continuum states lead to Stark shifts of states $|1\rangle$ and $|2\rangle$, which can be compensated by detuning the lasers from the exact Raman resonance between the bare states. If this is done correctly, the additional coherence ρ_{12} will lead to EIT-like modification of the refractive index experienced by both pulses. This extends to the situation for which the probe pulse is also strong. With both strong laser pulses, off-resonant CPT is possible. Formation of an off-resonance trapped state is an important aspect of an experiment investigating the elimination of optical-self focusing,¹⁴⁸ in which a nonlinear refractive index would otherwise lead to self-focusing, filamentation, and beam breakup of the strong probe field. Off-resonant CPT is also important in nonlinear frequency mixing, as demonstrated in experiments in lead vapor¹⁴⁹ and solid hydrogen.¹²⁶ In this case the condition $\Omega' > [\Delta \cdot \gamma_{12}]$ must be met, where Δ is the detuning of the fields from the intermediate resonances. This leads to the requirement of high peak powers, if the detunings Δ are large.

Propagation of two coupled strong pulses in a lambda-type system was discussed by Harris^{76,77} and Eberly et al.^{42,78} The discussed excitations scheme was equivalent to EIT. However, if both fields Ω_p and Ω_C are strong, the dressed atomic system reacts back on the field modes. This results in lossless propagation through the medium for both fields. For laser pulses with matched intensity envelopes (i.e., an identical form of temporal variation) any losses are minimal. If two initially matched pulses are simultaneously launched into a medium comprising atoms in the ground state, the system will self-organize so as to preserve the matched pulses and generate states, showing CPT. As the atoms are initially in the ground state $|1\rangle$, the probe pulse will initially experience loss and will have a lower group velocity than the coupling pulse. It will then start to lag the coupling pulse and so the pulse pair will automatically satisfy the condition for adiabatic preparation of trapped states (i.e., a counterintuitive pulse sequence). So, following the initial loss of probe pulse energy, the medium is set up for lossless transmission.

A proper insight into this process is best obtained in terms of CPT. The laser fields cause the formation of the superposition states $|C\rangle$ and $|NC\rangle$ [see Eq. (1)]. However, the atoms in the phase-coherent medium that is formed are also responsible for driving the two fields, and this means that even intensity envelopes which are initially different will evolve into matched pulses. This process results in the self-consistent formation of stable normal modes of the driving fields, one of which is uncoupled from the “uncoupled” atomic state and the other of which is “uncoupled” from the coupled atomic state. These new field modes result in the lossless propagation of pulses through a normally lossy medium, once a certain preparation energy has been extracted from the laser fields.

In the adiabatic limit the pulses are sufficiently intense, such that the timescale to establish EIT (i.e., $1/\Omega_C$) in the whole medium is fast with respect to the envelope evolution. Indeed, EIT can be much faster than the timescale required to establish population-trapped states in an individual

atom. This is due to the fact, that the latter process requires transfer of population, that is, irreversible exchange of energy between the field and the medium. In fact, the preparation time for EIT corresponds to a certain necessary pulse energy. The minimum energy $E_C^{(\min)}$ required to prepare the medium in EIT is essentially given by the photon energy of the coupling laser multiplied by the oscillator strength-weighted number of atoms:⁷⁹

$$E_C^{(\min)} = \frac{f_{13}}{f_{23}} \hbar \omega N_A \quad (8)$$

with the number of atoms N_A in the interaction volume and the oscillator strengths of the transitions at ω_{13} and ω_{23} . If the oscillator strength (or dipole moments) of the two transitions are comparable, Eq. (8) simply demands photon numbers, which exceed the number of atoms. Once the coupling laser pulse fulfills this condition (e.g., is long enough to provide a sufficient number of photons), the medium will be rendered transparent for all subsequent times. Application of this effect to the propagation of strong picosecond (ps) and femtosecond (fs) pulses a straightforward consideration. The preparation energy is not transferred irreversibly to the medium but is stored reversibly in the coherent excitation of state |2).

14.11 ULTRASLOW LIGHT PULSES

As already briefly discussed above (see Sec. 14.9), the steep, positive slope of the dispersion [$\text{Re } \chi^{(1)}(\omega)$] leads to a reduced group velocity v_g . The group velocity depends upon the slope, that is, the derivative of $\chi^{(1)}(\omega)$, as follows:

$$\frac{1}{v_g} = \frac{1}{c} + \frac{\pi}{\lambda} \left(\frac{\partial}{\partial \omega} [\text{Re } \chi^{(1)}] \right) \quad (9)$$

From the expression for the derivative $\partial \chi^{(1)}/\partial \omega$ [Eq. (7)], we see that this slope is steepest (i.e., v_g reaches a minimum) for $\Omega_C \gg \Gamma_2$ and $\Omega_C^2 \gg \Gamma_2 \Gamma_3$, but when Ω_C is still small compared to Γ_3 , hence $\partial \chi^{(1)}/\partial \omega \sim 1/\Omega_C^2$. In the limit of small coupling Rabi frequency Ω_C the group velocity v_g therefore yields

$$v_g = \frac{\hbar c \epsilon_0 \Omega_C^2}{2 n_A \omega_p \mu_{13}^2} \quad (10)$$

In an inhomogeneously broadened system the requirement for EIT is $\Omega_C > \gamma_{\text{Doppler}}$ (see Sec. 14.5). For typical atomic systems this usually also means $\Omega_C > \Gamma_3$. This condition constrains the group velocity reduction to relatively modest values. Thus in experiments in lead vapor¹⁴⁶ (see also Sec. 14.13), a group velocity reduction was clearly demonstrated, but the absolute reduction was only $v_g = c/165$. Further reductions in v_g are possible in media with small Doppler broadening γ_{Doppler} . In these media lower values for Ω_C (i.e., $\Omega_C \ll \Gamma_3$) still permit implementation of EIT. An elimination of Doppler broadening of spectral transitions is, for example, possible in a coupling scheme with transition frequencies $\omega_{13} \approx \omega_{23}$ and employing a Doppler-free geometry for the laser beams, that is, copropagating beams in a vee or lambda scheme and counterpropagating beam in a ladder scheme. Alternatively, laser-cooled and -trapped atomic samples can be used to reduce the Doppler width such that $\gamma_{\text{Doppler}} \ll \Gamma_3$. Cooled atoms also permit an excitation at frequencies $\omega_{13} \neq \omega_{23}$ and no restrictions are set on the beam geometry.

The spectral bandwidth of a light pulse determines a limitation for the possible reduction in the group velocity in a medium, or vice versa. The limit is set by the spectral linewidth of the EIT and by increased group velocity dispersion at frequencies detuned from exact resonance $\omega_p = \omega_{13}$. These effects

limit the permitted bandwidth of a light pulse with reduced group velocity to values significantly less than Ω_C . Thus only light pulses with sufficiently long duration (i.e., small spectral bandwidth without significant contribution from frequency components outside the bandwidth of EIT) are subject to an efficient reduction of the group velocity.

A number of experimental observations of ultraslow group velocity have been reported. In a Doppler-free excitation scheme, closely related to CPT, in cesium vapor the dispersion profile near resonance was determined using a Mach-Zehnder interferometer.¹⁵⁰ From the measured steep dispersion, a group velocity of $v_g = c/3000$ was inferred. A Doppler-free experiment in rubidium vapor was carried out, which directly measured a very large group velocity reduction. A probe beam with fast modulated amplitude was passed through the sample. The phase lag of the modulation with reference to the input signal was used to measure the group velocity. Very low values of $v_g = 90 \text{ ms}^{-1}$ were determined, even though the temperature of the sample was as high as 360 K.¹⁵¹ A lambda-type coupling scheme with Zeeman sublevels in rubidium was recently investigated and a magnetic field dependent group velocity as low as $v_g = 8 \text{ ms}^{-1}$ was measured.¹⁵²

A very dramatic demonstration of ultraslow light propagation was provided by Hau et al.¹¹ They employed a laser-cooled sodium sample trapped in a novel magnetic trap. Evaporative cooling of this sample was used to prepare a Bose-Einstein condensed (BEC) state. Experiments were performed (see Fig. 8) to investigate the delay in propagation of a probe pulse with pulse duration $\tau_p = 2.5 \mu\text{s}$, when a coupling laser was applied to the sample in the perpendicular direction. The reduced group velocity was studied both above and below the BEC transition temperature T_C . The lower group velocities found below T_C were due to the increased density of the sample when in this state. The lowest value of v_g that was measured was $v_g = 17 \text{ ms}^{-1}$. Other experiments on EIT also permitted the reduction of the speed of light in solids. We will discuss these experiments below (see Sec. 14.15).

In this section we gave only a very short introduction to the field of ultraslow light in coherently prepared media. The subject of ultraslow light as well as the related topic of storing light in media, prepared by EIT, attracted a huge attention in the last years. For a detailed review on ultraslow light and storage of light pulses, we like to draw the attention of the reader to a related article on this subject elsewhere in this book.

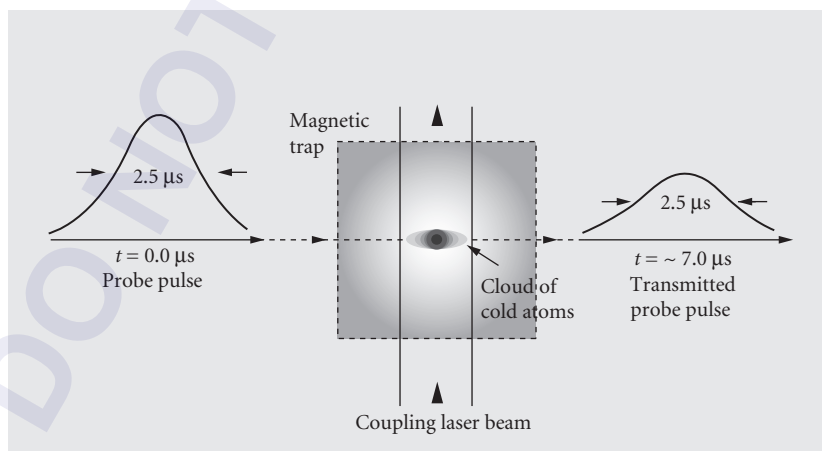


FIGURE 8 Schematic experimental setup for the observation of ultraslow light in a cloud of cold atoms, as reported by Hau et al.¹¹ A circularly polarized probe pulse passes through a cloud of cold atoms (length $L = 0.2 \text{ mm}$) with a transit time of $T = 7 \mu\text{s}$. This corresponds to a group velocity of 1/10 millionth of the speed of light in vacuum. [Reprinted by permission from Ref. 153. Copyright (1999).]

14.12 NONLINEAR OPTICAL FREQUENCY CONVERSION

In the previous sections, we discussed the manipulation of linear optical properties (i.e., absorption and dispersion) by EIT. In addition to these effects, also the nonlinear optical response of a medium is significantly modified by EIT. This has been shown to lead to frequency mixing with greatly enhanced efficiencies and to large Kerr-type nonlinearities. In addition to the dressed linear susceptibility [see Eq. (6)], there is also a dressed nonlinear susceptibility. This nonlinear susceptibility describes the coupling of the atom to fields at frequencies other than the frequencies of the probe laser and the coupling laser, that is, to frequencies which are generated by nonlinear interaction between the driving radiation fields. Consider, for example, a four-wave mixing process. Three fields at frequencies ω_A , ω_B , and ω_C couple via the nonlinear response of the medium to generate a new field at frequency ω_p with $\omega_p = \omega_A \pm \omega_B \pm \omega_C$. In general, the nonlinear susceptibility reaches maxima close to atomic resonances. Thus laser fields with frequencies tuned close to atomic resonances will yield enhanced frequency conversion. Figure 9 shows some basic coupling schemes for resonantly enhanced frequency conversion. Figure 9a depicts a situation, which we are already familiar with: Two radiation fields at frequencies ω_A and ω_B are tuned close to single-photon resonances at transitions $|1\rangle \rightarrow |4\rangle$ and $|2\rangle \rightarrow |4\rangle$ in a lambda-type level scheme. A third radiation field at frequency ω_C (i.e., tuned to resonance at the transition $|2\rangle \rightarrow |3\rangle$) mixes with the radiation fields at ω_A and ω_B to produce a wave at $\omega_p = \omega_A - \omega_B + \omega_C$. It is obvious from the coupling scheme, that the radiation fields ω_A and ω_B , as well as the fields ω_C and ω_p , will be automatically tuned close to two-photon resonance between the two lower states $|1\rangle$ and $|2\rangle$. The same holds true for the slightly modified coupling scheme in Fig. 9b. Here the lasers at ω_A and ω_B are far detuned from any single-photon resonance. Still, the two-photon resonance between the lower states $|1\rangle$ and $|2\rangle$ is maintained. If the laser frequencies ω_A and ω_B are degenerate (i.e., $\omega_A = \omega_B$) and the lambda system is modified to a ladder-type system, as depicted in Fig. 9c, the probe field is generated at frequency $\omega_p = 2\omega_A + \omega_C$. The latter is a typical scheme for frequency conversion to the regime of short-wavelength radiation [e.g., to the regime of vacuum-ultraviolet (VUV) or extreme-ultraviolet (XUV) radiation].

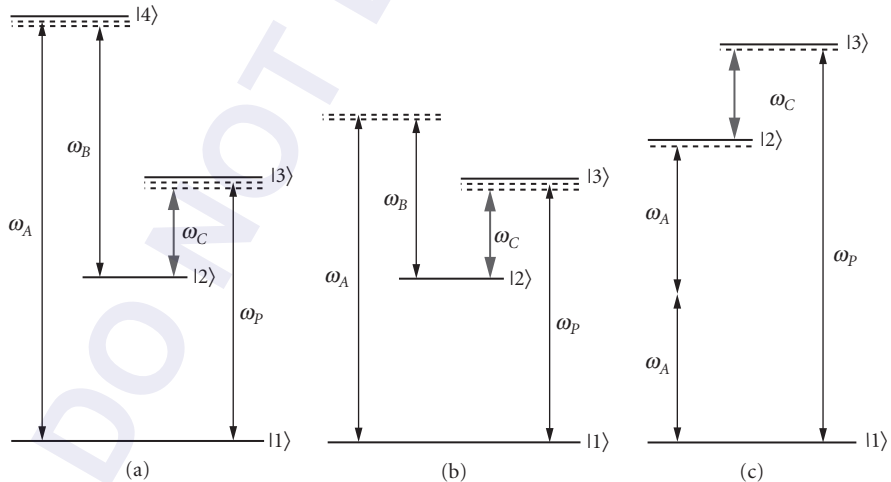


FIGURE 9 Three basic coupling schemes for four-wave mixing in coherently driven media. (a) Double-lambda system, driven on single-photon transitions by four radiation fields ω_A , ω_B , and ω_C generating a probe radiation field at frequency $\omega_p = \omega_A - \omega_B + \omega_C$. (b) Modified lambda system with the fields at ω_A and ω_B still tuned close to two-photon, but far detuned from any single-photon resonance. (c) Ladder system with degenerated fields at $\omega_A = \omega_B$, driving a four-wave mixing process to generate the probe field at $\omega_p = 2\omega_A + \omega_C$.

In conventional ‘‘incoherent’’ frequency mixing, the laser-driven resonances will enhance the efficiency of the frequency conversion process. There is also a significant drawback of this resonant excitation scheme: As also the generated wave at ω_p is resonant with the transition $|1\rangle \rightarrow |3\rangle$, it will be reabsorbed in the medium. However, if coherent interaction by EIT is considered, reabsorption is suppressed. To understand this feature of four-wave mixing, supported by EIT, let us consider the excitation processes in Fig. 9a in a slightly different way. The excitation scheme is equivalent to two coupled lambda systems. One lambda system is defined by the states $|1\rangle$, $|2\rangle$, and $|4\rangle$. The other lambda system is defined by the states $|1\rangle$, $|2\rangle$, and $|3\rangle$. If the laser at frequency ω_c is strong, it will drive EIT in the second lambda system of states $|1\rangle$, $|2\rangle$, and $|3\rangle$. In our previous notation, the generated wave at ω_p corresponds to the probe laser, and the laser at ω_c corresponds to the coupling laser. As the medium is prepared in EIT now for the probe laser, the generated field at frequency ω_p will experience no reabsorption any more, that is, the dressed linear susceptibility at the transition $|1\rangle \rightarrow |3\rangle$ is zero. In contrast, the dressed nonlinear susceptibility is not reduced to zero (see below and Fig. 5c). Thus the susceptibility provides a nonlinear coupling between the radiation fields, while absorption is suppressed.

The nonlinear susceptibility is calculated from Eqs. (4) and (5) in the same way as the linear susceptibilities in Eq. (6), however, now terms in ρ_{13} involving the fields at frequencies ω_A and ω_B are also included. The nonlinear susceptibility for the four-wave mixing process reads⁵

$$\begin{aligned} \chi_D^{(3)}(-\omega_p, \omega_A, \omega_B, \omega_C) &= \frac{2\mu_{23}\mu_{31}n_A}{3\epsilon_0\hbar^3} \frac{1}{\Omega_C^2 + (\Gamma_3 + 2i\Delta_{13})(\Gamma_2 + 2i[\Delta_{13} - \Delta_{23}])} \\ &\times \sum_j \mu_{1j}\mu_{j2} \left(\frac{1}{\omega_{j1} - \omega_A} + \frac{1}{\omega_{j1} - \omega_B} \right) \end{aligned} \quad (11)$$

where the sum on the right-hand side of Eq. (11) represents the contributions to the nonlinear susceptibility by all the states $|j\rangle$ of the atom, the fields at frequencies ω_A and ω_B are close to the two-photon resonance between the states $|1\rangle$ and $|2\rangle$. The modulus of this nonlinear susceptibility is plotted in Fig. 5c, although superficially similar to $[\text{Im} \chi_D^{(1)}]$ (see Fig. 5a) in the sense that both display the familiar Autler-Townes splitting. However, there is a fundamental difference. Paying attention to the center of the profile at detuning $\Delta_{13} = 0$, we see that, instead of destructive interference, there is constructive interference at this point. The value of the nonlinear susceptibility is indeed larger than the incoherent sum of the contributions from the dressed states.

Constructive interference in the nonlinear optical response of the atom is one of the most significant consequences of EIT. It accompanies the destructive interference which causes transparency at the same frequency. Already in 1990 Harris et al. recognized how this could be used to greatly improve four-wave mixing efficiencies.⁶ Enhancement of frequency mixing by atomic coherence and interference has been examined theoretically by a number of authors (see, e.g., Refs. 5, 6, 154, 155). The interference effects lead to improved frequency mixing in resonant systems because of three connected effects: (1) The resonant reabsorption in the medium is reduced due to the creation of transparency. (2) Phase matching is optimized due to vanishing dispersion near resonance. (3) Although the coupling field causes some reduction in the absolute value of the nonlinear susceptibility at resonance, the susceptibility is subject to constructive rather than destructive interference. For a medium with a large product of density and length ($n_A L$) (i.e., the uncoupled system is optically thick) there will be a large enhancement in the four-wave mixing conversion efficiency. This persists in a Doppler-broadened medium, providing the Rabi frequency Ω_C exceeds the Doppler width.

In a finite medium a calculation including absorption and phase-matching must be performed to predict the enhancement and the frequency dependence of the nonlinear response.^{5,6,101,102} In the limit of an optically deep medium and considering plane wave fields, the figure of merit for the conversion efficiency is given by the ratio $|\chi_D^{(3)}/\chi_D^{(1)}|$. This is physically reasonable, since $\chi_D^{(1)}$ characterizes the reabsorption and wave-vector mismatch in the medium, while $\chi_D^{(3)}$ determines the strength of the nonlinear coupling, which generates the new field. Thus we can see the importance of the destructive interference in the value of $\chi_D^{(1)}$ in the same frequency range, where $\chi_D^{(3)}$ experiences constructive interference. This enhances the generation efficiency in a Doppler-broadened

medium by many orders of magnitude, provided the Rabi frequency Ω_C exceeds the Doppler width. Modifications of $\chi_D^{(1)}$ result in reduced absorption and the essentially perfect phase matching for all resonant fields. Residual phase mismatch then appears only because of the dispersion caused by the remaining off-resonant states.

Electromagnetically induced phase matching was observed in an off-resonance four-wave mixing scheme in lead vapor.¹⁵⁶ In this experiment, the small detuning from resonance $\Delta = 6 \text{ cm}^{-1}$ indicated that transparency played little role in the enhancement (since the sample was optically thin). However, once the coupling laser strength exceeded a critical value, that is, $\Omega_C > [\gamma_{\text{doppler}}^2 + \Delta^2]^{1/2}$, the dispersion became effectively zero. This led to perfect phase matching and enhanced conversion efficiency.

In the first experiment to demonstrate enhancement by EIT, Hakuta, Stoicheff et al. investigated three-wave mixing (normally forbidden) in hydrogen when a DC electric field was applied to the sample.^{157,158} This experiment can also be viewed in the context of EIT. The DC field may be considered as an electromagnetic field at zero frequency. This field dresses the atom, creating two Stark states between which absorption cancellation occurs through interference. Second harmonic generation of a field at wavelength 243 nm is then resonantly enhanced at the two-photon transition between the states 1s and 2s. Due to EIT, the generated radiation at the Lyman- α wavelength 121.6 nm is no more reabsorbed. Further experiments by the same group then demonstrated EIT enhancement of a resonant four-wave mixing scheme in hydrogen atoms.¹⁵⁹ This was achieved in a scheme involving the states 1s, 2s, and 3p (corresponding to states |1>, |2> and |3> in our notation). The coupling laser at $\lambda_C = 656 \text{ nm}$ drove the transition between states 2s and 3p. The coupling laser mixed with a field at $\lambda_A = 243 \text{ nm}$, driving the two-photon transition between states 1s and 2s. This nonlinear interaction generated a field at $\lambda_p = 103 \text{ nm}$ (see Fig. 10, compare also Fig. 9c). Both lasers at λ_C and λ_A were derived from single-mode pulsed dye lasers. In this scheme photoionization of state 2s was an order of magnitude smaller than the ionization loss from state 3p. This is an important feature, since the coherence between the states 1s and 2s is essential for EIT effect to survive. Thus decay due to photoionization is a critical effect. However, so long as the decay rate of state 2s is significantly smaller than the rate of state 3p, the destructive interference in the absorption will persist and still lead to enhancement of the frequency mixing efficiency.

High conversion efficiencies require that the product of density and path length ($n_A L$) is large enough to ensure that the medium is initially strongly opaque. An atomic hydrogen beam was specially designed in the experiments, discussed above, in order to increase the beam density and path length. Thus EIT could be studied for the regime of an optically thin up to an optically thick medium.^{95,160} Working with products ($n_A L$) up to 10^{16} cm^{-2} , conversion efficiencies of $>10^{-3}$ were found. This is an extraordinarily large value for frequency conversion in a gaseous medium. There was, however, a limit to the obtainable conversion efficiency in this scheme, due to the large Doppler width in hydrogen. This is a consequence of the low mass of hydrogen and the elevated temperatures in the discharge used to produce the atoms. The coupling Rabi frequency Ω_C had to exceed the large Doppler width, leading to large Autler-Townes splittings and reduction in the magnitude of $\chi^{(3)}$ available for mixing. To overcome this, four-wave mixing schemes in atoms with smaller Doppler widths (e.g., krypton) have been studied. Experiments with a four-wave mixing scheme at a large product ($n_A L \approx 5 \times 10^{16} \text{ cm}^{-2}$ in krypton at room temperature have recently demonstrated a conversion efficiency of 10^{-2} for generation of a field at 123.6 nm.¹⁶¹ Quantum interference effects arising from the generated field itself have been reported to a limit in the optimum density for resonantly enhanced four-wave mixing. In a conversion scheme in rubidium¹⁶² at higher density the generated field itself became strong enough to cause a significant perturbation to the coherences in the system. This latter work, however, did not employ a single-mode coupling laser. Thus, the low limit in the rubidium density ($n_A < 10^{15} \text{ cm}^{-3}$), measured in these experiments, may not be reflected in situations where EIT is present.

Finally, we mention some alternative frequency conversion processes in lambda-type coupling schemes. In such schemes, the metastable state |2> often exhibits very long lifetimes. This leads to intrinsically low decay rates for the coherence ρ_{12} and near-perfect transparency is possible. Also in some selected systems (e.g., rubidium or sodium atoms) where the two lower states are hyperfine sublevels, Doppler-free configurations can be employed. These features make lambda systems a

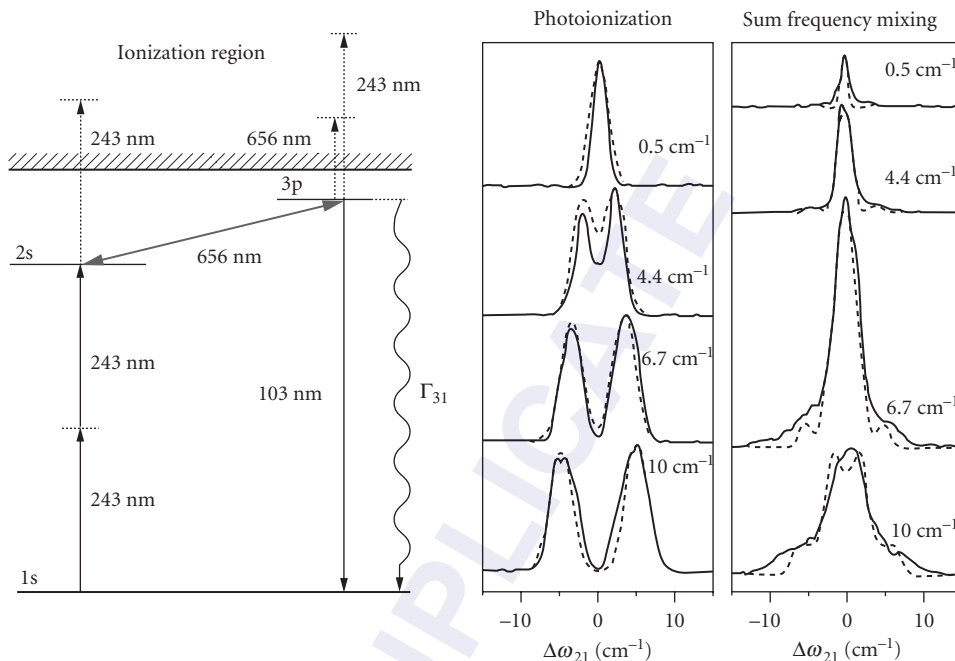


FIGURE 10 Four-wave mixing scheme in hydrogen atoms (compare also Fig. 9c). EIT is driven by the coupling laser at wavelength $\lambda_C = 656$ nm, tuned to the single-photon transition between states 2s and 3p [corresponding to states |2) and |3) in our notation]. A laser at $\lambda_A = 243$ nm drove the two-photon transition between states 1s and 2s [corresponding to states |1) and |2) in our notation]. A probe field was generated at wavelength $\lambda_p = 103$ nm, i.e., on the transition between the states 1s and 3p [corresponding to states |1) and |3) in our notation]. The photoionization yield and the sum-frequency mixing efficiency were measured versus the detuning $\Delta\omega_{21}$ [in our notation Δ_{21}] for different values of the coupling Rabi frequency Ω_C (given in wavenumbers). The photoionization spectra show that the Autler-Townes splitting increases with the Rabi frequency Ω_C . The conversion efficiency increases faster than the Autler-Townes splitting. A maximum is reached for a Rabi frequency of $\Omega_C = 6.7$ cm^{-1} . The medium is driven in EIT. For larger Rabi frequency there is no further enhancement in the conversion efficiency, as the re-absorption is already completely cancelled. However, the Autler-Townes splitting still increases with $\Omega_C = 6.7$ cm^{-1} . Thus the nonlinear coupling (i.e., the value of the nonlinear susceptibility) is reduced for further increasing Rabi frequency (compare the separation of the dressed states in Fig. 5c). [Reprinted figure with permission Ref 159. Copyright (1993) by the American Physical Society.]

very favorable choice for frequency conversion processes. Let us consider again the lambda scheme in Fig. 9a. In contrast to the above discussion, we will permit now for alternative combinations of the applied laser frequencies. For example, consider the coupling field at ω_C applied so as to create EIT. With two additional fields, a large number of different frequency mixing processes may arise. Radiation at frequencies ω_p and ω_A (or ω_B) in combination with the coupling field at ω_C drive the generation of new fields via four-wave mixing. The sign of the detuning from resonance for the applied or generated fields ω_A or ω_B may be positive or negative. The additional state |4) may be present at small detuning (i.e., a double-lambda scheme is prepared) or the detuning may be very large (i.e., the lasers drive essentially a three-level system) (compare Fig. 9b). Also if only two fields (e.g., ω_C and ω_p) are applied, they will drive a coherence ρ_{12} that can give rise to four-wave mixing processes. The coherence ρ_{12} can mix again with either of the fields ω_C and ω_p to generate, via a stimulated Raman process, Stokes and anti-Stokes fields. Moreover, if strong fields at ω_C and ω_A and a weak field at ω_B are applied, four-wave mixing yields a field at frequency ω_p , which is phase-conjugated to ω_B .

Nondegenerate four-wave mixing (NDFWM) based on EIT in a lambda scheme has been experimentally studied. In rubidium¹⁶³ a coupling field ω_c was applied resonantly, while a second field ω_A and a weak probe at ω_B were both applied with a detuning of 450 MHz from the resonance to state $|3\rangle$. A phase-conjugate field was generated at frequency ω_B . Due to EIT, absorption of this generated field vanished, but the nonlinearity remained resonantly enhanced. The susceptibilities $[\text{Im } \chi^{(1)}]$ and $\chi^{(3)}$ were measured independently under conditions of an optically thin medium. The data confirmed that $\chi^{(3)}$ was indeed enhanced by constructive interference. If an optically dense medium was used, a significant enhancement in NDFWM was observed. High phase-conjugate gain was also recently observed, applying very low laser powers. The effect arose from the presence of population trapping in a double-lambda scheme in sodium atoms.^{105,164} Resonant four-wave mixing,¹⁶⁵ driven by cw lasers, and frequency up-conversion¹⁶⁶ have also been observed in an experiment investigating a double-lambda scheme in sodium dimers.

14.13 NONLINEAR OPTICS AT MAXIMAL ATOMIC COHERENCE

When two laser pulses (i.e., a probe and a coupling laser pulse) propagate as matched pulses through a medium (see Sec. 14.10), large amounts of populations are prepared in trapped states. This is possible under conditions of strong excitation, that is, when the laser electric field strengths are large enough to drive adiabatic evolution for all atoms in the beam path. In this case the coherence ρ_{12} will reach a maximum magnitude (i.e. $|\rho_{12}| = |c_1^* c_2| = 1/2$). Thus, all the atoms are prepared in the trapped state $|\text{NC}\rangle$, that is, in a coherent superposition of the ground and the excited state with equal amplitudes $|c_1| = |c_2| = \sqrt{1/2}$. As the polarization of an atom is directly proportional to the coherence [compare Eq. (4)], also the polarization reaches a maximum at maximal atomic coherence. Under these circumstances, mixing of additional fields with the atom will become extremely efficient. We must not fail to mention, that the preparation of maximal coherences is also possible by adiabatic passage processes, other than EIT. Thus, STIRAP (see Sec. 14.3), coherent population return (CPR), rapid adiabatic passage (RAP), or Stark chirped rapid adiabatic passage (SCRAP) serve as alternative tools.⁴⁵

We consider now again a frequency conversion process under conditions of EIT in the coupling scheme, depicted in Fig. 9a and b. If the lasers at frequencies ω_p and ω_c induce a maximal coherence, any frequency conversion process, driven by the fields ω_A or ω_B occurs with high efficiency. This is due to the fact that the linear susceptibility, which governs the amount of absorption and dispersion, and the nonlinear susceptibility, which govern the frequency conversion efficiency, both have only a single nonresonant denominator with respect to the detuning. Thus the strengths for the nonlinear and the linear susceptibility for the fields at frequencies ω_A or ω_B are of the same order of magnitude. As a consequence, efficient energy exchange between the electric fields can occur in a distance of the order of the optical coherence length. The latter is determined by the real part of the susceptibility. This is equivalent to near-vacuum conditions for the dispersion and absorption of the medium while the nonlinearity is large.

The preparation of trapped states leading to the formation of a large atomic coherence has been applied by the Harris et al. to drive very efficient nonlinear frequency conversion.¹⁴⁹ The adiabatic frequency conversion process was implemented in a lambda-type scheme in lead vapor, driven by two strong lasers, that is, the coupling laser at $\lambda_c = 406$ nm and the probe laser at $\lambda_p = 283$ nm. A large, near-maximal coherence was created between the states $|1\rangle$ and $|2\rangle$ (see Fig. 11). The phase-coherent atoms acted as a local oscillator that mixed with a third laser field at 425 nm, detuned by 1112 cm^{-1} from the resonance frequency ω_3 . By a four-wave mixing process, involving the three fields, a new field at 293 nm was generated. The conversion efficiency was exceptionally large, reaching ~40 percent. The high conversion efficiency occurs since in this system the preparation of the optimal coherence ρ_{12} yielded a large nonlinear susceptibility, which was of the same size as the linear susceptibility. The same scheme was used, this time mixing a field at 233 nm to generate a field in the far-ultraviolet spectral region at 186 nm. In this case, the nonresonant detuning from the

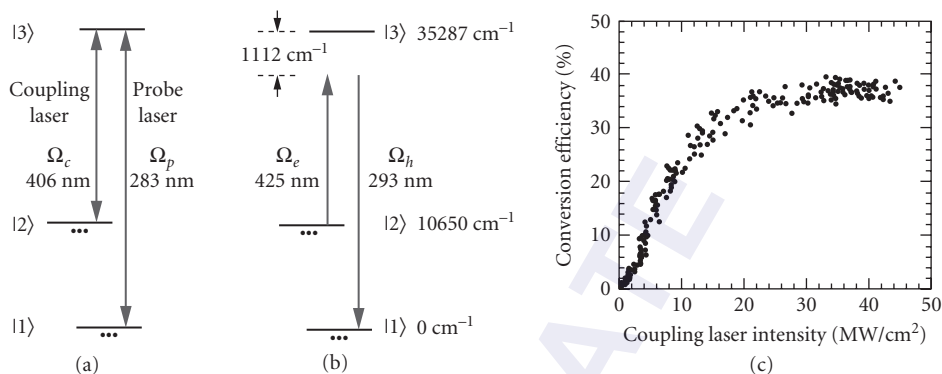


FIGURE 11 Efficient frequency conversion at maximal atomic coherence in lead vapor. (a) A large atomic coherence ρ_{12} was prepared by the probe and the coupling laser at wavelength $\lambda_p = 283$ nm and $\lambda_c = 406$ nm. (b) A laser at $\lambda_c = 425$ nm [in our notation λ_A] mixes with this coherence to generate a strong sum frequency mixing signal at $\lambda_h = 293$ nm [in our notation λ_B]. (c) Conversion efficiency versus coupling laser intensity. The efficiency increases linearly, till it reaches a plateau of 40 percent. This exhibits an extraordinary large value for a conversion process in a gaseous medium. [Reprinted figure with permission from Ref. 149. Copyright (1996) by the American Physical Society.]

third state in lead vapor was only ~ 40 cm^{-1} . A near-unity conversion efficiency was demonstrated in this case.²⁰

The coupling scheme in lead vapor, as discussed above, was recently the subject of a proposal for a broadband, high-efficiency, optical parametric oscillator (OPO).^{21,22} In this case, a maximal coherence ρ_{12} between the states $|1\rangle$ and $|2\rangle$ was created in the same fashion as discussed above. The coherence ρ_{12} then acts as a local oscillator in an optical parametric down-conversion process generating signal and idler waves in the infrared and far-infrared spectral region. In this system the nonlinear and linear responses of the medium were calculated to be of the same order and high conversion efficiencies up to 10 percent were predicted for the center of the OPO tuning range. Furthermore the device was predicted to cover the entire spectrum from the infrared to very long wavelength (i.e., essentially almost to the regime of DC fields).

The concept of maximal atomic coherence has been proposed to eliminate phase-mismatch in Raman scattering (i.e., the generation of Stokes and anti-Stokes radiation).¹²² In this scheme, the vibrational states $v = 0$ and $v = 1$ in the electronic ground state of a molecule form the lower states $|1\rangle$ and $|2\rangle$ of a Raman-type excitation scheme. The Raman scheme is equivalent to a lambda-type excitation scheme with large single-photon detunings Δ_{12} and Δ_{13} , while still the laser frequencies are tuned close to two-photon resonance ($\Delta_{12} - \Delta_{13}$). In such an excitation scheme, the dephasing rate γ_{12} is very small. Thus interference can occur and cause the dispersion to become negligible. Because of the removal of the usual phase mismatch, efficient operation of Raman scattering over a broad range of frequencies, that is, from the infrared to the vacuum-ultraviolet is possible. Another important prediction concerned efficient generation of broadband coherent spectra, associated with strong-field refractive index control under conditions of maximal coherence.²³ In this case a pair of laser pulses are slightly detuned from exact two-photon (Raman) resonance. Also in this case, which is closely related to EIT, a maximal coherence is generated.

To understand this feature of adiabatic excitation, we consider a two-level system of states $|1\rangle$ and $|2\rangle$, driven by a single laser pulse (see Fig. 12a). The dressed (adiabatic) eigenstates of the system read

$$|+\rangle = \sin \vartheta(t) |1\rangle + \cos \vartheta(t) |2\rangle \quad (12a)$$

$$|-\rangle = \cos \vartheta(t) |1\rangle - \sin \vartheta(t) |2\rangle \quad (12b)$$

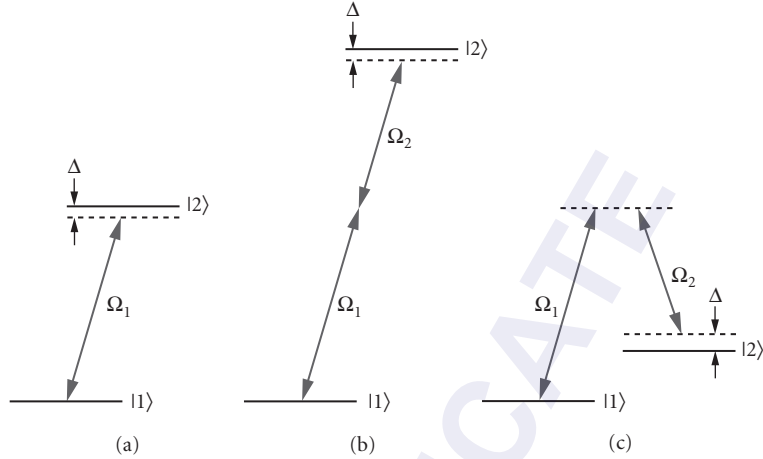


FIGURE 12 Coupling schemes for preparation of a maximal coherence by CPR. (a) Two-level system, driven by a single laser pulse on a single-photon transition. (b) Effective two-level system, driven by two laser pulses on a two-photon transition. (c) Effective two-level system, driven by two laser pulses on a Raman-type transition.

with the mixing angle $\vartheta(t)$, defined by $\vartheta(t) = 1/2 \arctan [\Omega_1(t)/\Delta]$. As we will consider now two strong laser pulses, we return to our previous designation of Ω_1 and Ω_2 (rather than Ω_p and Ω_c). We note that for resonant excitation at $\Delta = 0$ Eq. (12) yields the resonant form of the dressed states, as already introduced in Eq. (3). Let us assume now the Rabi frequency $\Omega_1(t)$ to be negligibly small outside a finite time interval $t_0 < t < t_1$, that is, outside the pulse duration $\tau_{\text{pulse}} = t_1 - t_0$. Consider now the case of the laser frequency detuned from exact resonance (i.e., $|\Delta| > 1/\tau_{\text{pulse}}$). If at the beginning of the interaction all the population is in the ground state, the state vector $\Psi(t)$ of the system at time $t = -\infty$ is aligned parallel to the adiabatic state $|-\rangle$. If the evolution of the system is adiabatic, the state vector $\Psi(t)$ remains always aligned with the adiabatic state $|-\rangle$ [see Eq. (12b) and the definition of the mixing angle ϑ]. Thus, during the excitation process (i.e., at intermediate times $t_0 < t < t_1$) the state vector $\Psi(t)$ is a coherent superposition of the bare states $|1\rangle$ and $|2\rangle$. Therefore, population is transiently excited to the upper state. However, at the end of the interaction at $t = +\infty$, the state vector $\Psi(t)$ becomes once again aligned with the initial state $|1\rangle$ [see Eq. (12b) and the definition of the mixing angle ϑ]. The population transferred during the process from the ground state $|1\rangle$ to the excited state $|2\rangle$ returns completely to the ground state after the excitation process. No population resides permanently in the excited state, no matter how large the transient intensity of the laser pulse may be. This phenomenon is called coherent population return (CPR), and is also known from coherent spectroscopy.^{167–169} Figure 13 shows the bare state population dynamics in CPR. As expected from the analytical considerations, population flows from the ground state to the excited state during the excitation and returns back to the ground state after the excitation process. If the peak Rabi frequency $\Omega_1^{(\text{max})}$ is sufficiently large (i.e., $\Omega_1^{(\text{max})} \gg \Delta$), the mixing angle during the process becomes $\vartheta = \pi/4$. Thus the coherence is $|\rho_{12}| = |c_1^* c_2| = \sin \vartheta \cos \vartheta = 1/2$ (i.e., a maximal coherence is prepared during the process).

The dynamics, discussed above, are not restricted to excitations of single-photon transitions, as depicted in Fig. 12a. Also multiphoton transitions [e.g., a two-photon transition (see Fig. 12b)], are driven in CPR, provided the laser frequencies are slightly detuned from the multiphoton transition frequency. This holds true also for Raman-type excitation schemes (see Fig. 12c), when the laser frequencies are tuned such that large detunings from any single-photon resonance as well as a small detuning from exact two-photon resonance occur. Thus, also in the case of Raman transitions between molecular vibrational states, a maximal coherence can be established and used for efficient frequency conversion processes. Moreover, a detailed theoretical analysis shows, that the adiabatic

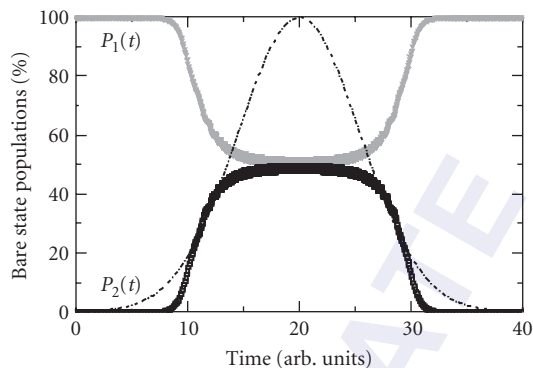


FIGURE 13 Numerical simulation of CPR.¹⁶⁷ When a transition in a two-level system is coherently driven, slightly detuned from the exact transition frequency (compare Fig. 12), population flows from the ground state $|1\rangle$ to the excited state $|2\rangle$ and back again [the probabilities $P_i(t)$ are indicated by the solid data points]. A transient maximal coherence is established during the process. The dashed line shows the temporal profile of the driving laser pulse.

excitation also permits control of the refractive index, that is, phase matching or phase mismatch plays only a minor role in the conversion process.

Harris et al. proposed, demonstrated, and applied the efficient generation of Raman sidebands in gaseous media under the conditions of maximal coherence.^{22–24,26–28,33,35} The experiment was performed in deuterium molecules in the gas phase. The deuterium molecules were cooled to 77 K in order to increase the population in the lowest rotational state and to reduce dephasing by collisions. Two nanosecond radiation pulses at frequencies ω_1 (pump) and ω_2 (Stokes) excited a Raman transition between the vibrational ground state $v = 0$ and the first excited state $v = 1$ (compare Fig. 12c). If the laser frequencies are slightly detuned from the Raman resonance, a maximum coherence is established. The medium acts now like a molecular modulator, oscillating at the frequency $\Delta\omega$ (i.e., the difference frequency between the vibrational states $v = 0$ and $v = 1$). The generation of Raman sidebands can be viewed as subsequent mixing process of the molecular modulator and the laser fields: Interaction of the field at ω_2 with the coherence at $\Delta\omega$ produces the first anti-Stokes sideband at $\omega_2 + \Delta\omega$ with large efficiency. Interaction of the first sideband with the coherence generates the second anti-Stokes sideband at $\omega_2 + 2\Delta\omega$. The process proceeds to higher-order sidebands. In their experiment, Harris et al. demonstrated conversion of the two driving radiation fields into Raman sidebands, covering a spectral region from 2.94 μm to 195 nm (i.e., from the far-infrared to the far-ultraviolet).²⁴ The largest conversion efficiency occurs, when the two radiation fields are slightly detuned from exact two-photon resonance. This, on the first glance surprising feature, confirms the theoretical expectations, as discussed above: Adiabaticity is maintained and a maximal coherence is prepared for detunings $|\Delta| > 1/\tau_{\text{pulse}}$.

Harris, Sokolov et al. applied the scheme for efficient generation of Raman sideband, as discussed above, to generate intense ultrashort radiation pulses.^{27,28,33} In an impressive experiment, Raman sidebands, with pulse durations in the nanosecond regime, were overlapped and combined. When their relative phases were appropriately adjusted, the combination of the phase-locked frequency components yielded a train of ultrashort radiation pulses with pulse duration in the regime of down to 1.6 fs (see Fig. 14).³³

We must not fail to mention that Hakuta et al. performed some of the early investigations on Raman sideband generation in coherently prepared media.^{25,31,32,34} These experiments were conducted in solid hydrogen. Other experiments on Raman sideband generation in media, coherently driven by nanosecond radiation pulses, were also performed by Marangos et al.^{29,30} In these experiments

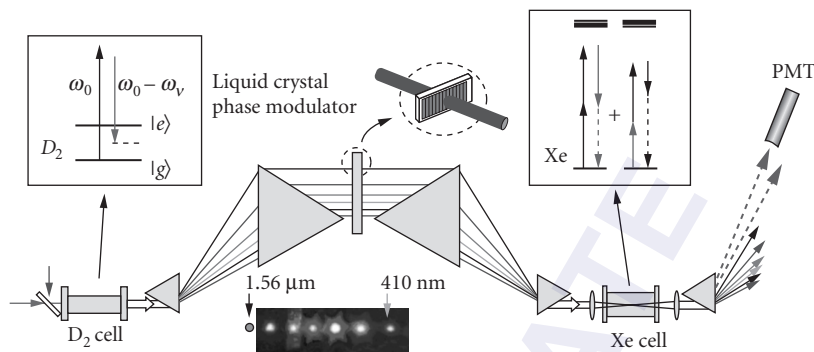


FIGURE 14 Experimental setup for temporal synthesis and characterization of single cycle optical pulses, generated by combination of Raman sidebands. The sidebands are generated in a cell with deuterium molecules, dispersed and their phases are independently varied by a liquid crystal modulator. The sidebands are recombined and focused into a target cell with xenon gas. The figure also shows the spectrum of the seven sidebands in the wavelength regime from 1.56 μm to 410 nm, used for generation of the ultrashort pulse. The duration of the ultrashort single-cycle pulses was 1.6 fs with a peak power of 1 MW. [Reprinted figure with permission from Ref. 33. Copyright (2005) by the American Physical Society.]

the coherence was probed by long (nanosecond) as well as ultrashort (femtosecond) radiation pulses.

Other work extended the concept of maximal coherence to the regime of excitations by ultrashort radiation pulses.^{170,171} As discussed above (see Sec. 14.5), ultrashort pulses are usually not favourable choice to drive adiabatic excitations on single-photon transitions. Thus, usually a maximal coherence cannot be prepared by ultrashort radiation pulses. However, in selected molecular media and with sufficient pulse energies, already a significantly enhanced molecular coherence (though not maximal) may enable efficient frequency conversion. Thus the efficient generation of Raman sidebands with pulse durations in the femtosecond time domain was demonstrated in hydrogen and methane molecules.^{170,171} The total conversion efficiency approached 10 percent.

Other extensions of the concepts, discussed above, utilize the excitation of Raman transitions by lasers of ultrabroad bandwidth.^{36–39} If the bandwidth γ_{laser} of a single ultrashort laser pulse covers the spacing $\Delta\omega$ between the vibrational ground state and the first excited state, the Raman transition is driven with frequency components, deduced from the single radiation pulse. Using SF_6 molecules as the Raman-active medium, an ultrashort pump pulse with wavelength $\lambda_p = 800$ nm and an additional ultrashort probe pulse with wavelength $\lambda_{pr} = 400$ nm, the molecular modulation technique has led to the generation of pulse trains and isolated pulses with a duration of a few femtoseconds.^{37–39}

14.14 NONLINEAR OPTICS AT THE FEW PHOTON LEVEL

As we have noted, one of the most remarkable features of EIT is that the nonlinear susceptibility undergoes constructive interference, while the linear susceptibility undergoes destructive interference. In a system without inhomogeneous broadening, perfect transparency can be induced for a coupling Rabi frequency $\Omega_C \ll \Gamma_3$. In this case the medium becomes transparent. Due to the small size of Ω_C with respect to the decay rate Γ_3 the nonlinear susceptibility will reach a value essentially identical to the value in the bare atomic system. Since atomic resonances are generally much narrower than those of a solid-state system, the magnitudes of the nonlinearity in this case can be many orders of magnitude larger than that of any solid. For instance, a very large nonlinear refractive

index was reported by Hau et al.¹¹ for EIT in a ultracold sodium vapor. They measured a value for the nonlinear refractive index that was $0.18 \text{ cm}^2 \text{ W}^{-1}$. This was 10^6 times larger than the nonlinear refractive index measured for cold cesium atoms in a system, which was not driven to EIT.¹⁷² In fact, the refractive index in cesium was itself much larger than the refractive index in a solid. For comparison, the largest nonlinear refractive indices of solid-state materials are on the order of less than $10^{-12} \text{ cm}^2 \text{ W}^{-1}$.¹⁷³

An important feature in the regime, discussed above, is the Kerr-type nonlinearity. This can be understood from examination of Fig. 9a. If fields at frequencies ω_p , ω_C , and ω_B are considered, the Kerr-type nonlinearity describes the coupling between the fields at ω_p and ω_B . The detuning between of field ω_A from state $|4\rangle$ is usually smaller than the detuning from state $|3\rangle$. Via the susceptibility $[\text{Re } \chi^{(3)}]$, the field at ω_A causes a cross-modulation with the field at ω_p . This effect was recently predicted to lead to a giant Kerr nonlinearity if the field at ω_C is of moderate strength. This is viable even when extremely low-power laser fields at ω_p and ω_B are applied. In a theoretical analysis,¹⁷⁴ values of $[\text{Re } \chi^{(3)}] \approx 3 \times 10^{-6} \text{ m}^2 \text{ V}^{-2}$ are found. These are large enough to permit measurable cross-phase modulations between fields containing only a few photons.

In the experiments of Hau et al.¹¹ under the same conditions where ultraslow light pulse propagation was measured, a large nonlinearity was deduced. This was extrapolated from the measurement of the AC Stark shift of the transparency minimum due to the interaction between the field at ω_C and state $|4\rangle$. This interaction is equivalent to the Kerr nonlinearity discussed in Ref. 174, when the field at ω_C also plays the role of the additional field (previously ω_B). The connection between the large nonlinearity and the ultraslow light propagation has been studied by several authors.^{151,175} In the case of the experiments discussed here, the connection is seen clearly to arise through the AC Stark shift of the steep dispersion curve responsible for ultraslow group velocities. For even very modest fields at ω_C , large changes in the absolute value of the dispersion are obtained, which result in significant phaseshifts. The latter are imparted on the probe field by itself. Similarly, in an experiment on ultraslow light, driven in a Doppler-free excitation scheme in rubidium,¹⁵¹ very large nonlinearities are evidenced by the highly efficient four-wave mixing process that was observed. Therefore, the authors highlight the strong connection between ultraslow group velocities and ultralarge nonlinearity.

A number of investigations showed, that the huge nonlinearities available in the excitation schemes, discussed above, are large enough to mediate significant interaction between pairs of photons. There are important implications if a single photon can cause a measurable modification to another single photon (e.g., by a strong cross-phase modulation). For instance, applications are proposed for quantum nondemolition measurements and for quantum information processing. In the latter case strong mutual interactions can be utilized to generate entangled pairs of photons that form the basis of quantum logic gates. Discussions of this topic include an analysis of frequency mixing and nonlinear phase shifts at the few photon level,¹⁷⁵ techniques for generating squeezed light at very low input powers,¹⁷⁶ photon switching,¹⁷⁷ and photon blockade within an atom-cavity system.¹⁷⁸

14.15 ELECTROMAGNETICALLY INDUCED TRANSPARENCY IN SOLIDS

While adiabatic processes (e.g., EIT) have been very extensively studied in the gas phase, only a few experiments were performed in solids.^{105–107,108–115} In contrast to gaseous media, solids offer large density, scalability, compactness, and convenience in their preparation and usage. Moreover, atomic diffusion is absent in solids. All of these features are of significant interest for applications in optical data storage and quantum information processing. However, large homogeneous and inhomogeneous broadenings exhibit an obstacle for the implementation of EIT in solids (see Sec. 14.5). Thus, only very specific systems in the solid phase permit appropriate conditions for EIT.

Some larger classes of solids, which are appropriate candidates for the implementation of adiabatic processes, combine the advantages of gaseous media and solids. Such media are (e.g., quantum

wells) electron-hole pairs (excitons), doped solids, and specific color centers. Quantum wells show a quite defined level structure, which is very similar to atoms in the gas phase. However, quantum wells are also subject to tremendous inhomogeneous broadening, because it is very hard to control the size, shape, and depths of a larger ensemble during the manufacturing process. Thus, usually only coherent excitation of a very small ensemble (or even single) quantum wells is possible—which is a significant drawback for applications, as one major advantage of a solid state medium (i.e., a large density) is lost. Electron-hole pairs have also been discussed in the scientific community as appropriate quantum systems for the implementation of adiabatic interactions. Like quantum dots, the electron-hole pairs also exhibit relatively well-defined level structure. As the dephasing times in quantum wells and electron-hole pairs are in the order of picoseconds or below, experimental implementations require ultrashort laser pulses. On the other hand, ultrashort laser pulses are usually not a good choice for driving adiabatic excitation of single-photon transitions (see Sec. 14.5). However, if the transition moments and the laser pulse energies are large enough, quantum wells and electron-hole pairs may serve as appropriate media for the implementation of EIT. Thus, coherent excitations in quantum wells or related nanostructures attracted significant attention in the recent years and already lead to first successful implementations of EIT in quantum wells.^{108–111}

A class of solid media, which combines the properties of the gas phase and the solid phase in a very advantageous way, are doped solids. In such media, atoms or ions are doped in a host crystal. The dopants in the host crystal still exhibit the level structure of the free atom. At room temperature, interaction with the host lattice leads to large homogeneous broadening. Moreover, as the dopants experience spatially varying electric fields from the host lattice, the atomic transitions are subject to substantial inhomogeneous broadening. At cryogenic temperatures, the homogeneous linewidth can be reduced well below the width of the inhomogeneous broadening, that is, selective optical excitation at zero-phonon transitions becomes possible. The still remaining inhomogeneous broadening might be considered as an advantage, as it permits the selective excitation of different ensembles of dopants. This was already utilized two decades ago, when spectral hole burning¹⁷⁹ was discovered as an appropriate way to address specific ensembles in an inhomogeneously broadened medium (e.g., for applications in frequency-selective, multichannel storage of information). Thus, when the dopants in a cryogenically cooled host crystal are appropriately prepared by spectral hole burning or related techniques, a specific subset of the atoms exhibits a quantum system with very narrow spectral linewidth.

Some basic adiabatic interactions have been demonstrated in the recent years in rare-earth doped solids. Thus, EIT,^{105–107,112,113} RAP,^{180,183} and recently STIRAP^{182,183} were successfully implemented in praseodymium ions Pr^{3+} , doped in a cryogenically cooled Y_2SiO_5 host crystal. $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$ is a well-established and commercially available material for optical data storage (e.g., by spectral hole burning). The spin coherence, that is, the dephasing time in the medium is pretty large (i.e., exceeding the timescale of 100 μs).^{114,179} This enables the implementation and observation of EIT with long interaction times, driven, for example, by modulated radiation, derived from cw lasers.

In the first successful implementation of EIT in a doped solid, Hemmer et al.¹⁰⁵ used a cw laser, operating at 606 nm, split in three beams, which are frequency modulated in acousto-optical modulators. Two of the three beams form the probe and the coupling laser. A third beam (i.e., the repump laser) is necessary to refill the spectral holes, burned by the probe and coupling laser. The $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$ crystal was cooled to a temperature of approximately 2 K to reduce the otherwise huge homogeneous broadening. Figure 15 shows the coupling scheme and the transmission of a probe laser in the medium. The praseodymium ions show a quite complicated level structure, as the ground state $^3\text{H}_4$ and the excited state $^1\text{D}_2$ interact with the electric field in the background of the host crystal and split in Stark components. The local variations in the electric field cause substantial inhomogeneous broadening in the crystal. However, by appropriate tuning and matching of the laser frequencies, a specific ensemble of ions is selected and driven in EIT. When the intensity of the coupling laser is large enough, the transmission of the probe laser is significantly enhanced in the center of the bare state resonance (see Fig. 15). Similar experiments on EIT in $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$ were also conducted by Ichimura et al.¹⁸³

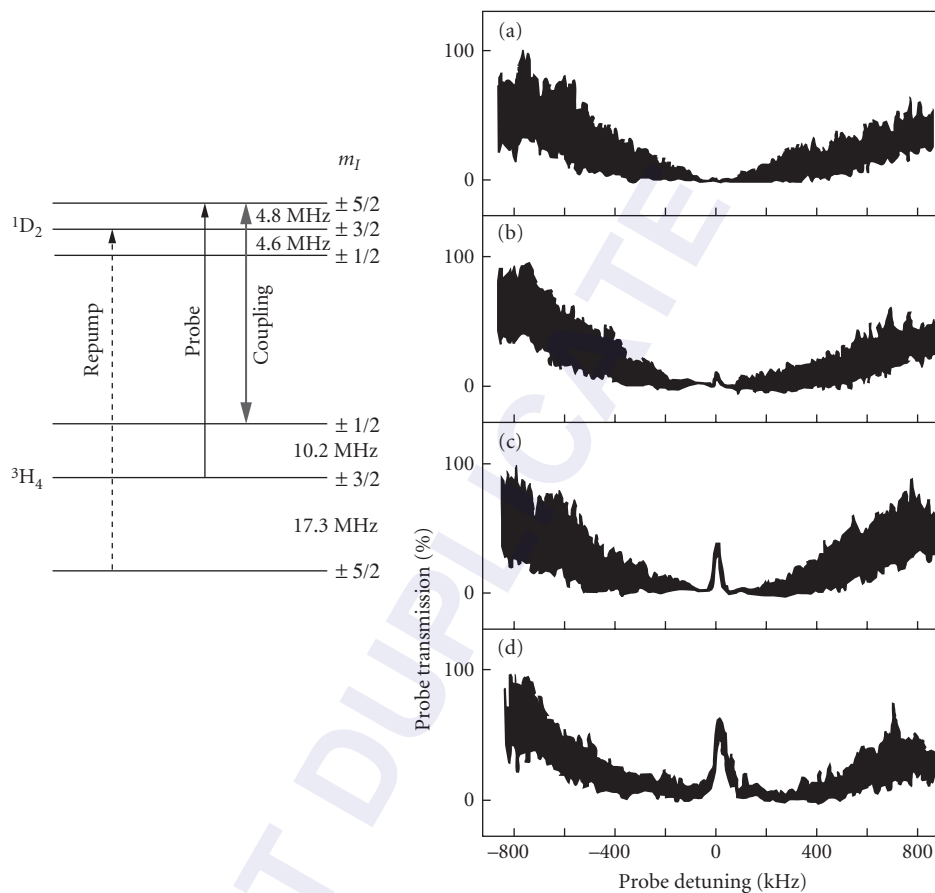


FIGURE 15 EIT in $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$. Coupling scheme and transmission of the probe laser versus detuning. The intensity of the coupling laser increases from graph (a) to (d). When the coupling laser is strong, the medium is driven to EIT and the probe laser, tuned to the bare state resonance, passes the medium with largely reduced losses. [The figure was published from Ref. 105. Copyright Elsevier (1997).]

In projects, dealing with the manipulation of optical processes by EIT, Hemmer et al. also reported efficient phase conjugation as well as ultraslow propagation of light pulses in $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$.¹¹³ The speed of light was reduced down to values of $v_g = 45$ m/s. Moreover, in these experiments also stopping of a light pulse in the crystal was observed. Recently, Longdell et al. reported stopping of light with storage times larger than 1 second in $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$, prepared by EIT (see Fig. 16).¹¹⁵

In addition to rare-earth doped solids, also special color centers are appropriate media for adiabatic interactions. Thus, Hemmer et al. also demonstrated EIT in nitrogen vacancy (N-V) color centers in diamond.¹⁸⁴ The transition moments in these color centers are much larger than in $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$. Therefore lower laser intensities are required to drive the medium. N-V color centers are interesting candidates for applications in quantum computation.¹⁸⁵ Thus the experimental investigations in $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$ and N-V diamonds, as well as theoretical proposals for EIT and other adiabatic interactions in nanostructures (see, e.g. Refs. 186 to 188 and references therein), may pave the way toward the implementation of solid state quantum information processors.

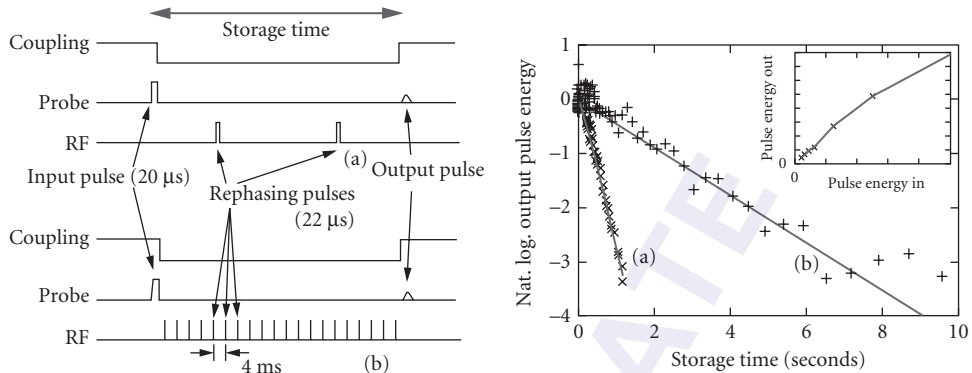


FIGURE 16 Demonstration of storage of light in $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$, prepared by EIT. Pulse sequences and output pulse energy storage time, for two different techniques (a, b) to control decoherence by an additional radio frequency (RF) field. The coupling laser pulse prepares the medium in EIT for the probe laser pulse. When the coupling pulse is switched off, the probe pulse is stored in the medium and recovered, when the coupling pulse is switched on again. The storage time exceeds several seconds. [Reprinted figure with permission from Ref. 115. Copyright by the American Physical Society (2005).]

14.16 CONCLUSION

In this chapter, we reviewed the basic concepts, implementations, and applications of EIT. We discussed the principles of EIT in terms of quantum interferences, atomic coherences, and dressed atomic states in a system of three quantum states, coherently driven by two radiation fields, that is, the probe laser field and the coupling laser field. The concept of EIT is closely related to other coherent, adiabatic interactions, for example, RAP, SCRAP, STIRAP, CPT, and CPR. However, while the latter techniques are focused to a large extent on the manipulation of population distributions, the ideas of EIT are mainly connected with atomic coherences and the manipulation of optical properties.

From some very basic considerations, we derived the main features of EIT, that is, vanishing absorption and dispersion at atomic resonances (due to destructive quantum interference), a steep slope of the dispersion at the atomic resonance, and enhanced nonlinear optical response of the medium (by constructive quantum interference). Thus the optical properties of quantum systems, manipulated by EIT, are very different from the properties of common optical media. These features lead to a rich variety of applications in the field of optical physics, for example, transmission through otherwise optically dense media, gain in quantum systems without population inversion, propagation of light with ultraslow velocity, storage of light pulses, and the significant enhancement of nonlinear optical processes. In the latter field, EIT permits frequency conversion in gaseous media with extraordinary large efficiencies.

We reviewed convincing experiments and striking applications of EIT. The majority of these experiments were implemented in atoms or molecules in the gas phase, but some also in selected solids. The investigations, as conducted so far, already lead to tremendous impact in the fields of optical information processing and quantum computation, applied optics, laser physics and physics on ultrashort timescales. However, we believe, there is a huge potential of EIT and related concepts, which is still to be explored in the future.

14.17 FURTHER READING

Arimondo, E., "Coherent Population Trapping in Laser Spectroscopy," *Progr. Opt.* **35**:259 (1996).

Fleischhauer, M., Imamoglu, A., and Marangos, J. P., "Electromagnetically Induced Transparency: Optics in Coherent Media," *Rev. Mod. Phys.* **77**:633 (2005).

- Harris, S. E., "Electromagnetically Induced Transparency," *Phys. Today* **50**:36 (1997).
- Kocharovskaya, O., "Amplification and Lasing without Inversion," *Phys. Rep.* **219**:175 (1992).
- Marangos, J. P., "Electromagnetically Induced Transparency," *J. Mod. Opt.* **45**:471 (1998).
- Sargent III, M., Scully, M. O., and Lamb, W. E., Jr., *Laser Physics*, Addison-Wesley, Reading (1974).
- Scully, M. O., "From Lasers and Masers to Phaseonium and Phasers," *Phys. Rep.* **219**:191 (1992).
- Scully, M. O., and Zubairy, M. S., *Quantum Optics*, Cambridge University Press, Cambridge (1997).
- Shore, B. W., "The Theory of Coherent Atomic Excitation," John Wiley & Sons, New York (1990).
- Shore, B. W., "Counter Intuitive Physics," *Contemp. Phys.* **36**:15 (1995).
- Vitanov, N. V., Halfmann, T., Shore, B. W., and Bergmann, K., "Laser-Induced Population Transfer by Adiabatic Passage Techniques," *Annu. Rev. Phys. Chem.* **52**:763 (2001).

14.18 REFERENCES

1. U. Fano, "Effects of Configuration Interaction on Intensities and Phase-Shifts," *Phys. Rev.* **124**:1866 (1961).
2. R. B. Madden and K. Codling, "Two-Electron Excitation States in Helium," *Astrophys. J.* **141**:364 (1965).
3. K. Maschke, P. Thomas, and E. O. Gobel, "Fano Interference in Type-11 Semiconductor Quantum-Well Structures," *Phys. Rev. Lett.* **67**:2646 (1991).
4. S. E. Harris, "Electromagnetically Induced Transparency," *Phys. Today* **50**:36 (1997).
5. M. Fleischhauer, A. Imamoglu, and J. P. Marangos, "Electromagnetically Induced Transparency: Optics in Coherent Media," *Rev. Mod. Phys.* **77**:633 (2005).
6. S. E. Harris, J. E. Field, and A. Imamoglu, "Non-Linear Optical Processes Using Electromagnetically Induced Transparency," *Phys. Rev. Lett.* **64**:1107 (1990).
7. M. O. Scully, "From Lasers and Masers to Phaseonium and Phasers," *Phys. Rep.* **219**:191 (1992).
8. M. O. Scully, "Enhancement of the Index of Refraction via Quantum Coherence," *Phys. Rev. Lett.* **67**:1855 (1991).
9. E. Arimondo and G. Orriols, "Nonabsorbing Atomic Coherences by Coherent Two-Photon Transitions in a Three-Level Optical Pumping," *Lett. Nuovo Cimento* **17**:333 (1976).
10. B. W. Shore, "The Theory of Coherent Atomic Excitation," John Wiley & Sons., New York (1990).
11. L. V. Hau, S. E. Harris, Z. Dunon, and C. H. Behroozi, "Light Speed Reduction to 17 Metres per Second in an Ultracold Atomic Gas," *Nature* **397**:594 (1999).
12. O. Kocharovskaya and Ya. I. Khanin, "Coherent Amplification of an Ultrashort Pulse in a Three Level Medium without Population Inversion," *JETP Lett.* **48**:630 (1988).
13. A. Imamoglu and S. E. Harris, "Lasers without Inversion: Interference of Dressed Lifetime Broadened States," *Opt. Lett.* **14**:1344 (1989).
14. E. E. Fill, M. O. Scully, and S.-Y. Zhu, "Lasing without Inversion via the Lambda Quantum-Beat Laser in the Collision Dominated Regime," *Opt. Comm.* **77**:36 (1990).
15. A. S. Zibrov, M. D. Lukin, D. E. Nikonov, L. Hollberg, M. O. Scully, V. L. Velichansky, and H. G. Robinson, "Experimental Observation of Laser Oscillation without Population Inversion via Quantum Interference in Rb," *Phys. Rev. Lett.* **75**:1499 (1995).
16. O. Kocharovskaya, "Amplification and Lasing without Inversion," *Phys. Rep.* **219**:175 (1992).
17. P. Mandel, "Lasing without Inversion: A Useful Concept?," *Contemp. Phys.* **34**:235 (1993).
18. D. S. Lee and K. J. Malloy, "Analysis of Reduced Interband Absorption Mechanisms in Semiconductor Quantum Wells," *IEEE J. Quant. Electr.* **30**:85 (1994).
19. A. Imamoglu and R. J. Ram, "Semiconductor Lasers without Population Inversion," *Opt. Lett.* **19**:1744 (1994).
20. A. J. Merriam, S. J. Sharpe, H. Xia, D. Manuszak, G. Y. Yin, and S. E. Harris, "Efficient Gas-Phase Generation of Coherent Vacuum Ultra-Violet Radiation," *Opt. Lett.* **24**:625 (1999).
21. S. E. Harris and M. Jain, "Optical Parametric Oscillators Pumped by Population-Trapped Atoms," *Opt. Lett.* **22**:636 (1997).

22. S. E. Harris, G. Y. Yin, M. Jain, H. Xia, and A. J. Merriam, "Nonlinear Optics at Maximum Coherence," *Phil. Trans. Roy. Soc. London (Ser A. Mathematical, Physical and Engineering Sciences)* **355**:2291 (1997).
23. S. E. Harris and A. V. Sokolov, "Broadband Spectral Generation with Refractive Index Control," *Phys. Rev. A* **55**:R4019 (1997).
24. S. E. Harris and A. V. Sokolov, "Raman Generation by Phased and Antiphased Molecular States," *Phys. Rev. Lett.* **85**:562 (2000).
25. F. E. Kien, J. Q. Liang, M. Katsuragawa, K. Ohtsuki, K. Hakuta, and A. V. Sokolov, "Sub-Femtosecond Pulse Generation with Molecular Coherence Control in Stimulated Raman Scattering," *Phys. Rev. A* **60**:1562 (1999).
26. D. D. Yavuz, A. V. Sokolov, and S. E. Harris, "Eigenvectors of a Raman Medium," *Phys. Rev. Lett.* **84**:000075 (2000).
27. S. E. Harris and A. V. Sokolov, "Subfemtosecond Pulse Generation by Molecular Modulation," *Phys. Rev. Lett.* **81**:002894 (1998).
28. A. V. Sokolov, D. R. Walker, D. D. Yavuz, G. Y. Yin, and S. E. Harris, "Femtosecond Light Source for Phase-Controlled Multiphoton Ionization," *Phys. Rev. Lett.* **87**:033402 (2001).
29. S. Gundry, M. P. Anscombe, A. M. Abdulla, S. D. Hogan, E. Sali, J. W. G. Tisch, and J. P. Marangos, "Off-Resonant Preparation of a Vibrational Coherence for Enhanced Stimulated Raman Scattering," *Phys. Rev. A* **72**:033824 (2005).
30. S. Gundry, M. P. Anscombe, A. M. Abdulla, E. Sali, J. W. G. Tisch, P. Kinsler, G. H. C. New, and J. P. Marangos, "Ultrashort-Pulse Modulation in Adiabatically Prepared Raman Media," *Opt. Lett.* **30**:180 (2005).
31. Fam Le Kien, Nguyen Hong Shon, and K. Hakuta, "Generation of Subfemtosecond Pulses by Beating a Femtosecond Pulse with a Raman Coherence Adiabatically Prepared in Solid Hydrogen," *Phys. Rev. A* **64**:051803(R) (2001).
32. Fam Le Kien, J. Q. Liang, M. Katsuragawa, K. Ohtsuki, and K. Hakuta, "Subfemtosecond Pulse Generation with Molecular Coherence Control in Stimulated Raman Scattering," *Phys. Rev. A* **60**:1562 (1999).
33. M. Y. Shverdin, D. R. Walker, D. D. Yavuz, G. Y. Yin, and S. E. Harris, "Generation of a Single-Cycle Optical Pulse," *Phys. Rev. Lett.* **94**:033904 (2005).
34. J. Q. Liang, M. Katsuragawa, Fam Le Kien, and K. Hakuta, "Sideband Generation Using Strongly Driven Raman Coherence in Solid Hydrogen," *Phys. Rev. Lett.* **85**:2474 (2000).
35. Fam Le Kein, K. Hakuta, and A. V. Sokolov, "Pulse Compression by Parametric Beating with a Prepared Raman Coherence," *Phys. Rev. A* **66**:023813 (2002).
36. B. R. Washburn, S. R. Ralph, R. S. Windeler, "Ultrashort Pulse Propagation in Air-Silica Microstructure Fiber," *Opt. Expr.* **10**:575 (2002).
37. M. Wittmann, A. Nazarkin, and G. Korn, "Fs-Pulse Synthesis Using Phase Modulation by Impulsively Excited Molecular Vibrations," *Phys. Rev. Lett.* **84**:55508 (2000).
38. M. Wittmann, A. Nazarkin, and G. Korn, "Synthesis of Periodic Femtosecond Pulse Trains in the Ultraviolet by Phase-Locked Raman Sideband Generation," *Opt. Lett.* **26**:298 (2001).
39. N. Zhavoronkov and G. Korn, "Generation of Single Intense Short Optical Pulses by Ultrafast Molecular Phase Modulation," *Phys. Rev. Lett.* **88**:203901 (2002).
40. B. J. Dalton and P. L. Knight, *Laser Physics, Lecture Notes in Physics*, Vol. 182, I. D. Harvey and D. F. Walls (eds.), Springer-Verlag, Berlin (1983).
41. H. I. Yoo and J. H. Eberly, "Dynamical Theory of an Atom with Two or Three Levels Interacting with Quantized Cavity Fields," *Phys. Rep.* **118**:239 (1985).
42. J. H. Eberly, "Transmission of Dressed Fields in Three-Level Media," *Quant. Semiclass. Opt.* **7**:373 (1995).
43. E. Arimondo, "Coherent Population Trapping in Laser Spectroscopy," *Progr. Opt.* **35**:259 (1996).
44. N. V. Vitanov, T. Halfmann, B. W. Shore, and K. Bergmann, "Laser-Induced Population Transfer by Adiabatic Passage Techniques," *Annu. Rev. Phys. Chem.* **52**:763 (2001).
45. T. Halfmann, "Enhanced Frequency Conversion in Coherently Prepared Media," in *Recent Research Activities in Chemical Physics: From Atomic Scale to Macroscale*, edited by A. F. Terzis and E. Paspalakis, Research Signpost, India (2007).
46. P. L. Knight, "Laser Induced Continuum Structure," *Comm. Ac. Mol. Phys.* **15**:193 (1984).
47. P. L. Knight, M. A. Lauder, and B. J. Dalton, "Laser Induced Continuum Structure," *Phys. Rep.* **190**:1 (1990).

48. K. Böhmer, T. Halfmann, L. P. Yatsenko, D. Charalambidis, A. Horsmans, and K. Bergmann, "Laser-Induced Continuum Structure in the Two Ionization Continua of Xenon," *Phys. Rev. A* **66**:013406 (2002).
49. T. Halfmann, L. P. Yatsenko, M. Shapiro, B. W. Shore, and K. Bergmann, "Population Trapping and Laser-Induced Continuum Structure in Helium: Experiment and Theory," *Phys. Rev. A* **58**:R46 (1998).
50. L. P. Yatsenko, T. Halfmann, B. W. Shore, and K. Bergmann, "Photoionization Suppression by Continuum Coherence: Experiment and Theory," *Phys. Rev.* **59**:2926 (1999).
51. S. H. Autler and C. H. Townes, "Stark Effect in Rapidly Varying Fields," *Phys. Rev.* **100**:703 (1955).
52. I. D. Abella, N. A. Kunitz, and S. R. Hanmann, "Photon Echoes," *Phys. Rev.* **141**:391 (1966).
53. B. R. Mollow, "Power Spectrum of Light Scattered by 2-Level Systems," *Phys. Rev.* **188**:1969 (1969).
54. B. R. Mollow, "Stimulated Emission and Absorption near Resonance for Driven Systems," *Phys. Rev. A* **5**:2217 (1972).
55. P. L. Knight and P. W. Miloni, "The Rabi Frequency in Optical Spectra," *Phys. Rep.* **66**:21 (1980).
56. L. Alien and J. H. Eberly, *Optical Resonance and 2-Level Atoms*, Dover, New York (1987).
57. G. Alzetta, L. Moi, and G. Orriols, "Nonabsorption Hyperfine Resonances in a Sodium Vapour Irradiated with a Multimode Dye Laser," *Nuovo Cimento B* **52**:209 (1979).
58. H. R. Gray, R. M. Whitley, and C. R. Stroud, "Coherent Trapping of Atomic Populations," *Opt. Lett.* **3**:218 (1978).
59. J. Oreg, F. T. Hioe, and J. H. Eberly, "Adiabatic Following in Multilevel Systems," *Phys. Rev. A* **29**:690 (1984).
60. U. Gaubatz, P. Rudecki, S. Schiemann, and K. Bergmann, "Population Transfer between Molecular Vibrational Levels by Stimulated Raman-Scattering with Partially Overlapping Laser Fields: A New Concept and Experimental Results," *J. Chem. Phys.* **92**:5363 (1990).
61. N. Dam, L. Oudejans, and J. Reuss, "Relaxation Rates of Ethylene Obtained from Their Effect on Coherent Transients," *Chem. Phys.* **140**:217 (1990).
62. A. Aspect, E. Arimondo, R. Kaiser, N. Vansteenkiste, and C. Cohen-Tannoudji, "Laser Cooling below the One-Photon Recoil Energy by Velocity-Selective Coherent Population Trapping," *Phys. Rev. Lett.* **61**:826 (1988).
63. F. Bardou, B. Saubamea, J. Lawall, K. Shimizu, O. Emile, C. Westbrook, A. Aspect, and C. Cohen-Tannoudji, "Sub-Recoil Laser Cooling with Pre-Cooled Atoms," *Comptes Rendus de Acad. Sci. Paris, Ser. II*, **318**:877 (1994).
64. M. R. Doery, M. T. Widmer, M. J. Bellanca, W. F. Buell, T. H. Bergemann, H. Metcalf, and E. J. D. Vrenbregt, "Population Accumulation in Dark States and Sub-Recoil Laser Cooling," *Phys. Rev. A* **52**:2295 (1995).
65. R. M. Whitley and C. R. Stroud, Jr, "Double Optical Resonance," *Phys. Rev. A* **14**:1498 (1976).
66. P. M. Radmore and P. L. Knight, "Population Trapping and Dispersion in a Three-Level System," *J. Phys. B* **15**:561 (1982).
67. B. W. Shore, "Counter Intuitive Physics," *Contemp. Phys.* **36**:15 (1995).
68. M. O. Scully and M. Suhail Zubairy, *Quantum Optics*, Cambridge University Press, Cambridge (1997).
69. B. Lounis and C. Cohen-Tannoudji, "Coherent Population Trapping and Fano Profiles," *J. Phys.* **11**(2):579 (1992).
70. S. E. Harris and J. J. Macklin, "Lasers without Inversion: Single-Atom Transient Response," *Phys. Rev. A* **40**:4135 (1989).
71. J. L. Cohen and P. R. Berman, "Amplification without Inversion: Understanding Probability Amplitudes, Quantum Interferences and Feynman Rules in a Strongly Driven System," *Phys. Rev. A* **55**:3900 (1997).
72. D. ter Haar, "Theory and Application of the Density Matrix," *Rep. Prog. Phys.* **24**:304 (1961).
73. M. Sargent III, M. O. Scully, and W. E. Lamb, Jr, *Laser Physics*, Addison-Wesley, Reading (1974).
74. C. Cohen-Tannoudji, B. Zambon, and E. Arimondo, "Quantum-Jump Approach to Dissipative Processes: Application to Amplification without Inversion," *J. Opt. Soc. Am. B* **10**:2107 (1993).
75. Y.-Q. Li and M. Xiao, "Transient Properties of an Electromagnetically Induced Transparency in Three-Level Atoms," *Opt. Lett.* **20**:1489 (1995).
76. S. E. Harris, "Electromagnetically Induced Transparency with Matched Pulses," *Phys. Rev. Lett.* **72**:552 (1993).
77. S. E. Harris, "Normal-Modes for Electromagnetically Induced Transparency," *Phys. Rev. Lett.* **72**:52 (1994).
78. J. H. Eberly, M. L. Pons, and H. R. Haq, "Dressed-Field Pulses in an Absorbing Medium," *Phys. Rev. Lett.* **72**:56 (1994).

79. S. E. Harris and Z.-F. Luo, "Preparation Energy for Electromagnetically Induced Transparency," *Phys. Rev. A* **52**:928 (1995).
80. J. E. Field, "Vacuum-Rabi-Splitting-Induced Transparency," *Phys. Rev. A* **47**:5064 (1993).
81. J. E. Field and A. Imamoglu, "Spontaneous Emission into an Electromagnetically Induced Transparency," *Phys. Rev. A* **48**:2486 (1993).
82. S. E. Harris, J. E. Field, and A. Kasapi, "Dispersive Properties of Electromagnetically Induced Transparency," *Phys. Rev. A* **46**:R29 (1992).
83. G. G. Padmabandu, G. R. Welch, I. N. Shubin, E. S. Fry, D. E. Nikonov, M. D. Lukin, and M. O. Scully, "Laser Oscillation without Population Inversion in a Sodium Atomic Beam," *Phys. Rev. Lett.* **76**:2053 (1996).
84. M. O. Scully, S.-Y. Zhu, and A. Gavrielides, "Degenerate Quantum-Beat Laser: Lasing without Inversion and Inversion without Lasing," *Phys. Rev. Lett.* **62**:2813 (1989).
85. V. G. Arkhipkin and Yu. I. Heller, "Radiation Amplification without Population Inversion at Transitions to Autoionising States," *Phys. Lett. A* **98**:12 (1983).
86. S. E. Harris, "Lasers without Inversion: Interference of a Lifetime Broadened Resonance," *Phys. Rev. Lett.* **62**:1033 (1989).
87. M. O. Scully and S.-Y. Zhu, "Ultra-Large Index of Refraction via Quantum Interference," *Opt. Comm.* **87**, 134 (1992).
88. M. Fleischhauer, C. H. Keitel, M. O. Scully, and C. Su, "Lasing without Inversion and Enhancement of the Index of Refraction via Interference of Incoherent Pump Processes," *Opt. Comm.* **87**:109 (1992).
89. A. D. Wilson-Gordon and H. Friedman, "Enhanced Index of Refraction: A Comparison between Two- and Three-Level Systems," *Opt. Comm.* **94**:238 (1992).
90. M. Fleischhauer, C. H. Keitel, M. O. Scully, C. Su, B. T. Ulrich, and S.-Y. Zhu, "Resonantly Enhanced Refractive Index without Absorption via Atomic Coherence," *Phys. Rev. A* **46**:1468 (1992).
91. U. Rathe, M. Fleischhauer, S.-Y. Zhu, T. W. Hänsch, and M. O. Scully, "Non-Linear Theory of Index Enhancement via Quantum Coherence and Interference," *Phys. Rev. A* **47**:4994 (1993).
92. A. S. Manka, J. P. Dowling, C. M. Bowden, and M. Fleischhauer, "Piezophotonic Switching due to Local-Field Effects in a Coherently Prepared Medium of 3-Level Atoms," *Phys. Rev. Lett.* **73**:1789 (1994).
93. S. E. Harris, "Refractive-Index Control with Strong Fields," *Opt. Lett.* **19**:2018 (1994).
94. D. J. Fulton, R. R. Moseley, S. Shepherd, B. D. Sinclair, and M. H. Dunn, "Effects of Zeeman Splitting on Electromagnetically Induced Transparency," *Opt. Comm.* **116**:231 (1995).
95. G. Z. Zhang, M. Kasuragawa, K. Hakuta, R. I. Thompson, and B. P. Stoicheff, "Sum-Frequency Generation Using Strong-Field Coupling and Induced Transparency in Atomic Hydrogen," *Phys. Rev. A* **52**:1584 (1995).
96. R. R. Moseley, B. D. Sinclair, and M. H. Dunn, "Local Field Effects in the Three-Level Atom," *Opt. Comm.* **108**:247 (1994).
97. J. P. Dowling and C. M. Bowden, "Near Dipole-Dipole Effects in Lasing without Inversion: An Enhancement of Gain and Absorptionless Index of Refraction," *Phys. Rev. Lett.* **70**:1421 (1993).
98. G. S. Agarwal, "Exact Solution for the Influence of Laser Temporal Fluctuations on Resonance Fluorescence," *Phys. Rev. Lett.* **37**:1383 (1976).
99. S. Swain, "Theory of Atomic Processes in Strong Resonant Electromagnetic Fields," *Adv. At. Mol. Phys.* **16**:159 (1980).
100. B. J. Dalton and P. L. Knight, "The Effect of Laser Field Fluctuations on Coherent Population Trapping," *J. Phys. B* **15**:3997 (1982).
101. P. A. Lakshmi and S. Swain, "Effects of Laser Linewidths on Nonlinear Optical Processes Exploiting Electromagnetically Induced Transparency," *J. Mod. Opt.* **38**:2031 (1991).
102. J. C. Petch, C. H. Keitel, P. L. Knight, and J. P. Marangos, "Role of Electromagnetically Induced Transparency in Resonant Four-Wave Mixing Schemes," *Phys. Rev. A* **53**:543 (1996).
103. Y. Zhao, C. Wu, B. S. Ham, M. K. Kirn, and E. Awad, "Microwave Induced Transparency in Ruby," *Phys. Rev. Lett.* **79**:64 (1997).
104. G. B. Serapiglia, E. Paspalakis, C. Sirtori, K. L. Vodopyanov, and C. C. Phillips, "Laser-Induced Quantum Coherence in a Semiconductor Quantum Well," *Phys. Rev. Lett.* **84**:1019 (2000).
105. B. S. Ham, P. R. Hemmer, and M. S. Shahar, "Efficient Electromagnetically Induced Transparency in Rare-Earth Doped Crystals," *Opt. Comm.* **144**:227 (1997).

106. B. S. Ham, S. M. Shahriar, and P. R. Hemmer, "Electromagnetically Induced Transparency over Spectral Hole-Burning Temperature in a Rare-Earth-Doped Solid," *J. Opt. Soc. Am. B* **16**:801 (1999).
107. P. R. Hemmer, D. P. Katz, J. Donoghue, M. Cronin-Golomb, M. S. Shahriar, and P. Kumar, "Efficient Low-Intensity Optical Phase Conjugation Based on Coherent Population Trapping in Sodium," *Opt. Lett.* **20**:982 (1995).
108. M. Philipps and H. Wang, "Spin Coherence and Electromagnetically Induced Transparency via Exciton Correlations," *Phys. Rev. Lett.* **84**:186401 (2002).
109. M. Philipps and H. Wang, "Electromagnetically Induced Transparency due to Intervalence Band Coherence in a GaAs Quantum Well," *Opt. Lett.* **28**:831 (2002).
110. M. C. Phillips, H. Wang, I. Romyantsev, N. H. Kwong, R. Takayama, and R. Binder, "Electromagnetically Induced Transparency in Semiconductors via Biexciton Coherence," *Phys. Rev. Lett.* **91**:183602 (2003).
111. W. W. Chow, H. C. Schneider, and M. C. Phillips, "Theory of Quantum-Coherence Phenomena in Semiconductor Quantum Dots," *Phys. Rev. A* **68**:053802 (2003).
112. K. Ichimura, K. Yamamoto, and N. Gemma, "Evidence for Electromagnetically Induced Transparency in a Solid Medium," *Phys. Rev. A* **58**:4116 (1998).
113. A. V. Turukhin, V. S. Sudarshanam, M. S. Shahriar, J. A. Musser, B. S. Ham, and P. R. Hemmer, "Observation of Ultraslow and Stored Light Pulses in a Solid," *Phys. Rev. Lett.* **88**:023602 (2002).
114. E. Kuznetsova, O. Kocharovskaya, P. Hemmer, and M. O. Scully, "Atomic Interference Phenomena in Solids with a Long-Lived Spin Coherence," *Phys. Rev. A* **66**:063802 (2002).
115. J. J. Longdell, E. Fraval, M. J. Sellars, and N. B. Manson, "Stopped Light with Storage Times Greater than One Second Using Electromagnetically Induced Transparency in a Solid," *Phys. Rev. Lett.* **95**:063601 (2005).
116. J. Y. Gea-Banacloche, Y.-Q. Li, S.-Z. Jin, and M. Xiao, "Electromagnetically Induced Transparency in Ladder-Type Inhomogeneously Broadened Media: Theory and Experiment," *Phys. Rev. A* **51**:576 (1995).
117. S. Shepherd, D. J. Fulton, and M. H. Dunn, "Wavelength Dependence of Coherently Induced Transparency in a Doppler-Broadened Cascade Medium," *Phys. Rev. A* **54**:5394 (1996).
118. M. Xiao, Y.-Q. Li, and J. Gea-Banacloche, "Measurements of Dispersive Properties of Electromagnetically Induced Transparency in Rubidium Atoms," *Phys. Rev. Lett.* **74**:666 (1995).
119. R. R. Moseley, S. Shepherd, D. J. Fulton, B. D. Sinclair, and M. H. Dunn, "Spatial Consequences of Electromagnetically Induced Transparency: Observation of Electromagnetically Induced Focusing," *Phys. Rev. Lett.* **74**:670 (1995).
120. M. Misunga, T. Mukai, K. Watanabe, and T. Mukai, "Dressed-Atom Spectroscopy of Cold Cs Atoms," *J. Opt. Soc. B* **13**:2696 (1996).
121. S. A. Hopkins, E. Usadi, H. X. Chen, and A. V. Durrant, "Electromagnetically Induced Transparency of Laser-Cooled Rubidium Atoms in Three-Level Λ Type Systems," *Opt. Comm.* **134**:185 (1997).
122. K. Hakuta, M. Suzuki, M. Katsuragawa, and J. Z. Li, "Self-Induced Phase Matching in Parametric Anti-Stokes Stimulated Raman Scattering," *Phys. Rev. Lett.* **79**:209 (1997).
123. K. J. Boller, A. Imamoglu, and S. E. Harris, "Observation of Electromagnetically Induced Transparency," *Phys. Rev. Lett.* **66**:2593 (1991).
124. J. E. Field, K. H. Hahn, and S. E. Harris, "Observation of Electromagnetically Induced Transparency in Collisionally Broadened Lead Vapor," *Phys. Rev. Lett.* **67**:3062 (1991).
125. A. Kasapi, "Enhanced Isotope Discrimination Using Electromagnetically Induced Transparency," *Phys. Rev. Lett.* **77**:1035 (1996).
126. Hui Xia, A. J. Merriam, S. J. Sharpe, G. Y. Yin, and S. E. Harris, "Electromagnetically Induced Transparency with Spectator Momenta," *Phys. Rev. A* **59**:R3190 (1999).
127. Hui Xia, S. J. Sharpe, A. J. Merriam, and S. E. Harris, "Electromagnetically Induced Transparency in Atoms with Hyperfine Structure," *Phys. Rev. A* **56**:R3362 (1998).
128. Y.-Q. Li and M. Xiao, "Electromagnetically Induced Transparency in a Three-Level System in Rubidium Atoms," *Phys. Rev. A* **51**:R2703 (1995).
129. Y.-Q. Li and M. Xiao, "Observation of Quantum Interference between the Dressed States in an Electromagnetically Induced Transparency," *Phys. Rev. A* **51**:4959 (1995).
130. R. R. Moseley, S. Shepherd, D. J. Fulton, B. D. Sinclair, and M. H. Dunn, "Electromagnetically Induced Focusing," *Phys. Rev. A* **53**:408 (1996).

131. S. Wielandy and Alexander L. Gaeta, "Investigation of Electromagnetically Induced Transparency in the Strong Probe Regime," *Phys. Rev. A* **58**:2500 (1998).
132. W. Ketterle, K. B. Davis, M. A. Joffe, A. Martin, and D. E. Pritchard, "High Densities of Cold Atoms in a Dark Spontaneous-Force Optical Trap," *Phys. Rev. Lett.* **77**:2253 (1993).
133. T. van der Veldt, J.-F. Roch, P. Grelu, and P. Grangier, "Nonlinear Absorption and Dispersion of Cold Rb-87 Atoms," *Opt. Comm.* **137**:420 (1997).
134. H. X. Chen, A. V. Durran, J. P. Marangos, and J. A. Vaccaro, "Observation of Transient Electromagnetically Induced Transparency in a Rubidium Lambda System," *Phys. Rev. A* **58**:1545 (1998).
135. S. Scandolo and F. Bassani, "Non-Linear Sum Rules: The Three-Level and Anharmonic-Oscillator Modes," *Phys. Rev. B* **45**:13257 (1992).
136. F. S. Cataliotti, C. Fort, T. W. Hansch, M. Inguscio, and M. Prevedelli, "Electromagnetically Induced Transparency in Cold Free Atoms: Test of a Sum Rule for Non-Linear Optics," *Phys. Rev. A* **56**:2221 (1997).
137. J. Gao, C. Guo, X. Guo, G. Jin, P. Wang, J. Zhao, H. Zhang, Y. Jiang, D. Wang, and D. Jiang, "Observation of Light Amplification without Population Inversion," *Opt. Comm.* **93**:323 (1992).
138. J. Gao, H. Z. Zhang, H. F. Cui, X. Z. Guo, Y. Jiang, Q. W. Wang, G. X. Jin, and J. S. Li, "Inversionless Amplification in Sodium," *Opt. Comm.* **110**:590 (1994).
139. J. A. Kleinfeld and A. D. Streater, "Observation of Gain due to Coherence Effects in Potassium-Helium Mixture," *Phys. Rev. A* **49**:R4301 (1994).
140. E. S. Fry, X. Li, D. Nikonov, G. G. Padmabandu, M. O. Scully, A. V. Smith, F. K. Tittel, C. Wang, S. R. Wilkinson, and S.-Y. Zhu, "Atomic Coherence Effects within the Sodium D1 Line: Lasing without Inversion via Population Trapping," *Phys. Rev. Lett.* **70**:3235 (1993).
141. A. Nottelman, C. Peters, and W. Lange, "Inversionless Amplification of Picosecond Pulses due to Zeeman Coherence," *Phys. Rev. Lett.* **70**:1783 (1993).
142. W. E. van der Veer, R. J. J. van Dienst, A. Donselmann, and H. B. van Linden van den Huevcll, "Experimental Demonstration of Light Amplification without Population Inversion," *Phys. Rev. Lett.* **70**:3243 (1993).
143. C. Peters and W. Lang, "Laser Action Below Threshold Inversion due to Coherent Population Trapping," *Appl. Phys. B* **62**:221 (1996).
144. A. S. Zibrov, M. D. Lukin, L. Hollberg, D. E. Nikonov, M. O. Scully, H. G. Robinson, and V. L. Velichansky, "Experimental Demonstration of Enhanced Index of Refraction via Quantum Coherence in Rb," *Phys. Rev. Lett.* **76**:3935 (1996).
145. M. Fleischhauer and M. O. Scully, "Quantum Sensitivity Limits of an Optical Magnetometer Based on Phase Coherence," *Phys. Rev. A* **49**:1973 (1994).
146. A. Kasapi, M. Jain, G. Y. Yin, and S. E. Harris, "Electromagnetically Induced Transparency: Propagation Dynamics," *Phys. Rev. Lett.* **74**:2447 (1995).
147. A. Kasapi, G. Y. Yin, M. Jain, and S. E. Harris, "Measurement of Lorentzian Linewidth by Pulse-Propagation Delay," *Phys. Rev. A* **53**:4547 (1996).
148. M. Jain, A. J. Merriam, A. Kasapi, G. Y. Yin, and S. E. Harris, "Elimination of Optical Self-Focusing by Population Trapping," *Phys. Rev. Lett.* **75**:4385 (1995).
149. M. Jain, H. Xia, G. Y. Yin, A. J. Merriam, and S. E. Harris, "Efficient Non-Linear Frequency Conversion with Maximal Atomic Coherence," *Phys. Rev. Lett.* **77**:4326 (1996).
150. O. Schmidt, R. Wynands, Z. Hussein, and D. Meschede, "Steep Dispersion and Group Velocity Dispersion below $c/3000$ in Coherent Population Trapping," *Phys. Rev. A* **53**:R27 (1996).
151. M. M. Kash, V. A. Sautenkov, A. S. Zibrov, L. Hollberg, G. R. Welch, M. D. Lukin, Y. Roslovtsev, E. S. Fry, and M. O. Scully, "Ultraslow Group Velocity and Enhanced Non-Linear Optical Effects in a Coherently Driven Hot Atomic Gas," *Phys. Rev. Lett.* **82**:5229 (1999).
152. D. Budker, D. F. Kimball, S. M. Rochester, and V. V. Yaschuk, "Nonlinear Magneto-Optics and Reduced Velocity of Light in Atomic Vapour with Slow Ground State Relaxation," *Phys. Rev. Lett.* **83**:1767 (1999).
153. J. P. Marangos, "Slow Light in Cool Atoms," *Nature* **397**:559 (1999).
154. S. P. Tewari and G. S. Agarwal, "Control of Phase-Matching and Nonlinear Generation in Dense Media by Resonant Fields," *Phys. Rev. Lett.* **56**:1811 (1986).
155. G. S. Agarwal and S. P. Tewari, "Large Enhancements in Non-Linear Generation by External Electromagnetic Fields," *Phys. Rev. Lett.* **70**:1417 (1993).

156. M. Jain, G. Y. Yin, J. E. Field, and S. E. Harris, "Observation of Electromagnetically Induced Phase Matching," *Opt. Lett.* **18**:998 (1993).
157. K. Hakuta, L. Marmet, and B. P. Stoicheff, "Electric-Field-Induced Second-Harmonic Generation with Reduced Absorption," *Phys. Rev. Lett.* **66**:596 (1991).
158. K. Hakuta, L. Marmet, and B. P. Stoicheff, "Nonlinear Optical Generation with Reduced Absorption Using Electric-Field Coupling in Atomic Hydrogen," *Phys. Rev. A* **45**:5152 (1992).
159. G. Z. Zhang, K. Hakuta, and B. P. Stoicheff, "Nonlinear Optical Generation Using Electromagnetically Induced Transparency in Hydrogen," *Phys. Rev. Lett.* **71**:3099 (1993).
160. R. S. D. Sihombing, M. Katsuragawa, G. Z. Zhang, and K. Hakuta, "Quantum Interference in Resonant Multi-Photon Ionisation Processes for Strongly Coupled Atomic Systems," *Phys. Rev. A* **54**: 1551(1996).
161. C. Dorman, I. Kucukcara, and J. P. Marangos, "Measurement of High Conversion Efficiency to 123.6 nm Radiation in a Four-Wave Mixing Scheme Enhanced by Electromagnetically Induced Transparency," *Phys. Rev. A* **61**:013802-1 (2000).
162. L. Deng, W. R. Garret, M. G. Payne, and D. Z. Lee, "Observation of a Critical Concentration in Laser-Induced Transparency and Multi-Photon Excitation and Ionisation in Rubidium," *Opt. Lett.* **21**:928 (1996).
163. Y.-Q. Li and M. Xiao, "Enhancement of Non-Degenerate Four-Wave-Mixing Based on Electromagnetically Induced Transparency in Rubidium Atoms," *Opt. Lett.* **21**:1064 (1996).
164. T. T. Grove, M. S. Shahriar, P. R. Hemmer, P. Kumar, V. S. Sudarshanam, and M. Cronin-Golomb, "Distortion-Free Gain and Noise Correlation in Sodium Vapor with Four-Wave Mixing and Coherent Population Trapping," *Opt. Lett.* **22**:769 (1997).
165. S. Babin, U. Hinze, E. Tiemann, and B. Wellegehausen, "Continuous Resonant Four-Wave-Mixing in Double-A Level Configurations of Na," *Opt. Lett.* **21**:1186 (1996).
166. A. Apolonskii, S. Balischev, U. Hinze, E. Tiemann, and B. Wellegehausen, "Continuous Frequency Up-Conversion in a Double-A Scheme in Na," *Appl. Phys. B* **64**:435 (1997).
167. A. Peralta Conde, L. Brandt, and T. Halfmann, "Trace Isotope Detection Enhanced by Coherent Elimination of Power Broadening," *Phys. Rev. Lett.* **97**:243004 (2006).
168. T. Halfmann, T. Rickes, N. Vitanov, and K. Bergmann, "Lineshapes in Coherent Two-Photon Excitation," *Opt. Comm.* **220**:353 (2003).
169. N. V. Vitanov, B. W. Shore, L. P. Yatsenko, T. Halfmann, T. Rickes, K. Böhmer, and K. Bergmann, "Power Broadening Revisited: Theory and Experiment," *Opt. Comm.* **199**:117 (2001).
170. E. Sali, P. Kinsler, G. H. New, K. J. Mendham, T. Halfmann, J. W. G. Tisch, and J. P. Marangos, "Behavior of High-Order Stimulated Raman Scattering in a Highly Transient Regime," *Phys. Rev. A* **72**:013813 (2005).
171. E. Sali, K. J. Mendham, J. W. G. Tisch, T. Halfmann, and J. P. Marangos, "High-Order Stimulated Raman Scattering in a Highly Transient Regime Driven by a Pair of Ultrashort Pulses," *Opt. Lett.* **29**:495 (2004).
172. A. Lambrecht, J. M. Courty, S. Reynaud, and E. Giacobino, "Cold Atoms: A New Medium for Quantum Optics," *Appl. Phys. B* **60**:129 (1995).
173. D. N. Nikogosyan, *Properties of Optical and Laser-Related Materials: A Handbook*, John Wiley & Sons, New York (1997).
174. H. Schmidt and A. Imamoglu, "Giant Kerr Nonlinearities Obtained by Electromagnetically Induced Transparency," *Opt. Lett.* **21**:1936 (1996).
175. S. E. Harris and L. V. Hau, "Non-Linear Optics at Low Light Levels," *Phys. Rev. Lett.* **82**:4611 (1999).
176. M. D. Lukin, A. B. Matsko, M. Fleischhauer, and M. O. Scully, "Quantum Noise and Correlations in Resonantly Enhanced Wave-Mixing Based on Atomic Coherence," *Phys. Rev. Lett.* **82**:1847 (1999).
177. S. E. Harris and Y. Yamamoto, "Photon Switching by Quantum Interference," *Phys. Rev. Lett.* **81**:3611 (1998).
178. A. Imamoglu, H. Schmidt, G. Woods, and M. Deutsch, "Strongly Interacting Photons in a Non-Linear Cavity," *Phys. Rev. Lett.* **79**:1467 (1997).
179. W. E. Moerner (ed.), *Persistent Spectral Hole Burning: Science and Applications*, Springer, Berlin (1988).
180. J. Klein, F. Beil, and T. Halfmann, "Rapid Adiabatic Passage in a $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$ Crystal," *J. Phys. B*, in press (2007).
181. L. Rippe, M. Nilsson, and S. Kröll, "Experimental Demonstration of Efficient and Selective Population Transfer and Qubit Distillation in a Rare-Earth-Metal-Ion-Doped Crystal," *Phys. Rev. A* **71**:062328 (2005).

182. J. Klein, F. Beil, and T. Halfmann, "Stimulated Raman Adiabatic Passage in a $\text{Pr}^{3+}:\text{Y}_2\text{SiO}_5$ Crystal," submitted (2007).
183. H. Goto and K. Ichimura, "Population Transfer via Stimulated Raman Adiabatic Passage in a Solid," *Phys. Rev. A* **74**:049902 (2006).
184. P. R. Hemmer, A. V. Turukhin, M. S. Shahriar, and J. A. Musser, "Raman-Excited Spin Coherences in Nitrogen-Vacancy Color Centers in Diamond," *Opt. Lett.* **26**:361 (2001).
185. L. Childress, M. V. Gurudev-Dutt, J. M. Taylor, A. S. Zibrov, F. Jelezko, J. Wrachtrup, P. R. Hemmer, and M. D. Lukin, "Coherent Dynamics of Coupled Electron and Nuclear Spin Qubits in Diamond," *Science* **314**:281 (2006).
186. U. Hohenester, J. Fabian, and F. Troiani, "Adiabatic Passage Schemes in Coupled Semiconductor Nanostructures," *Opt. Comm.* **264**:426 (2006).
187. J. Siewert, T. Brandes, and G. Falci, "Adiabatic Passage with Superconducting Nanocircuits," *Opt. Comm.* **264**:435 (2006).
188. W. W. Chow, H. C. Schneider, and M. C. Philipps, "Theory of Quantum Coherence Phenomena in Semiconductor Quantum Dots," *Phys. Rev. A* **68**:053802 (2003).

STIMULATED RAMAN AND BRILLOUIN SCATTERING

John Reintjes and Mark Bashkansky

*Optical Sciences Division
Naval Research Laboratory
Washington, D.C.*

15.1 INTRODUCTION

Raman and Brillouin scattering are inelastic scattering processes in which the wavelength of the scattered radiation is different from that of the incident light and a change in the internal energy of the scattering medium occurs. The main distinction between Raman and Brillouin scattering is the type of internal mode involved. Raman scattering involves nonpropagating collective modes in the material. Examples include electronic excitations and molecular vibrations and rotations in liquids and gases, electronic excitations and optical phonons in solids, and electron-plasma oscillations in plasmas. Brillouin scattering involves low-frequency propagating waves, for example acoustic waves in solids, liquids, and gases and ion-acoustic waves in plasmas. The two processes exhibit a range of similarities and differences in the properties of the scattering process as well as in the materials that are involved. General descriptions of Raman and Brillouin scattering are given in Refs. 1–8. Collections of papers on specific aspects of Raman and Brillouin scattering are contained in Refs. 9–11. Review articles are given in Refs. 12 and 13.

15.2 RAMAN SCATTERING

Raman scattering occurs in a wide variety of solids, liquids, gases, and plasmas. The most common form of Raman scattering is one in which the incident light, termed the *pump*, is scattered into light at a longer wavelength, termed the *Stokes wave*, with the energy difference between the incident and scattered photons being taken up in excitation of the appropriate mode of the material. The difference between the incident and scattered photon energy is termed the *Stokes shift*. The identification of the scattered wave as the Stokes wave is made in analogy with the Stokes shift to longer wavelengths in fluorescence, but the dynamics of the two processes are different except for interactions that are resonant with an allowed single-photon resonant transition. Raman scattering in which the incident light is scattered to a light wave at a shorter wavelength, accompanied by a deexcitation of an internal mode of the medium, is termed *anti-Stokes scattering*, and the scattered wave is termed the *anti-Stokes wave*. The difference between the anti-Stokes and pump photon energies is termed the

anti-Stokes shift. Again, the analogy is with the corresponding process in fluorescence. Raman shifts can range from tens to tens of thousands of wavenumbers, and are determined entirely by the material and the mode involved.

Raman Interactions

Raman scattering can be viewed in the semiclassical model as a two-photon interaction in which the material makes a real transition from an initial to a final state and a pump photon is destroyed while a Stokes or anti-Stokes photon is created. Several types of Raman interactions are possible. These are illustrated in the level diagrams of Fig. 1, which show Stokes scattering, anti-Stokes scattering, anti-Stokes scattering with four-wave mixing, multiple Stokes scattering, and hyper-Raman scattering. Of these, the most common is Stokes scattering (Fig. 1a), in which the pump photon at frequency ω_L is scattered into a longer-wavelength Stokes photon ω_S , accompanied by the excitation of an internal mode of the medium at frequency ω_o .

Anti-Stokes scattering (Fig. 1b), in which the pump photon is scattered into a shorter-wavelength anti-Stokes photon ω_{AS} accompanied by the deexcitation of an internal mode of the medium, requires initial excitation into upper levels of the medium. The anti-Stokes interaction illustrated in Fig. 1b is less common than the Stokes interaction, occurring most often in spontaneous Raman scattering when the levels are excited thermally. In stimulated processes, this interaction incurs exponential loss unless a population inversion is created between the initial and final states.

Anti-Stokes Raman scattering involving a four-wave mixing interaction, as illustrated in Fig. 1c, is much more common in Raman scattering. In this interaction, two pump photons are scattered into a Stokes and anti-Stokes interaction with no net excitation or deexcitation of the medium. This interaction requires perfect or approximate phase matching, depending on the conditions of the scattering interaction. Multiple Raman scattering (Fig. 1d) occurs when the Stokes wave becomes

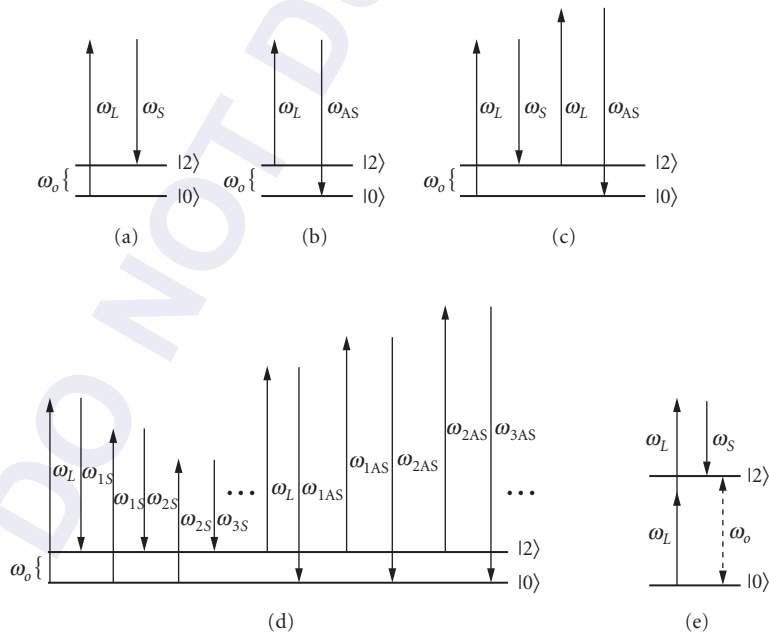


FIGURE 1 Level diagrams showing (a) stimulated Raman Stokes scattering; (b) stimulated Raman anti-Stokes scattering; (c) coherent anti-Stokes four-wave mixing; (d) multiple Stokes and anti-Stokes scattering; and (e) hyper-Raman scattering.

TABLE 1 Frequency and k -Vector Relationships in Various Raman Interactions

Stokes scattering	$\omega_s = \omega_L - \omega_o$ $k_o = k_L - k_s$
Anti-Stokes scattering	$\omega_{AS} = \omega_L + \omega_o$ $k_o = k_{AS} - k_L$
Coherent anti-Stokes scattering	$\omega_s + \omega_{AS} = 2\omega_L$ $k_{AS} = 2k_L - k_s$
Multiple Stokes scattering	$\omega_{ns} = \omega_{(n-1)s} - \omega_o = \omega_L - n\omega_o$ $k_o = k_{(n-1)s} - k_{ns}$
Hyper-Raman scattering	$\omega_s = 2\omega_L - \omega_o$ $k_o = 2k_L - k_s$

powerful enough to drive its own Raman interaction. This generally occurs when the pump wave is significantly above the stimulated Raman threshold. Under these conditions, multiple Stokes waves are generated, each one shifted from its effective pump wave by the frequency of the internal mode of the medium. Multiple Stokes and anti-Stokes waves can also be created through four-wave mixing processes involving one or more of the pump- or frequency-shifted waves. Hyper-Raman scattering (Fig. 1e) involves multiphoton interactions in which two or more pump photons are scattered into a single Stokes photon with excitation of an appropriate mode of the material.

In all of these interactions, energy is conserved among the incident and scattered photons and internal energy of the medium. The appropriate relations are summarized in Table 1.

Not all levels in a material can be involved in Raman scattering. In general, Raman scattering follows the rules for two-photon dipole transitions. In materials with inversion symmetry, the initial and final states must have the same parity, and therefore are mutually exclusive with absorptive transitions. In materials without inversion symmetry, internal levels can be both Raman and optically active.

Regimes of Raman Scattering

Raman scattering can occur in the spontaneous and stimulated regimes. In the spontaneous regime, the power of the Stokes and anti-Stokes waves is proportional to the power of the pump wave. The entire manifold of Raman-active internal modes is present in the scattered spectrum, with relative intensities of the Stokes components being determined by the relative Raman scattering cross sections for the various modes. Anti-Stokes scattering arises in spontaneous Raman scattering through thermal excitation of the internal modes. Therefore, the intensity of anti-Stokes modes is reduced from that of the Stokes modes for the same internal level by the thermal excitation factor $e^{-h\omega_o/kT}$. Anti-Stokes scattering in the spontaneous regime is generally less common than Stokes scattering, except when low-lying rotational levels of molecules or low-frequency phonons in solids are involved, because thermal excitation of the internal mode is required. Anti-Stokes scattering can be more prominent in stimulated or coherent scattering processes, as will be described in later sections. Spontaneous Raman scattering is used primarily for spectroscopic studies, especially for modes that are forbidden in single-photon absorption or emission measurements.

Stimulated Raman Scattering (SRS) occurs when the intensity of the incident pump wave is strong enough to initiate a positive feedback effect in the medium, resulting in exponential growth of the scattered wave. Stimulated Raman scattering is used for wavelength shifting of coherent light, amplification, improved optical beam properties, pulse compression, phase conjugation, and beam combining. Coherent Raman interactions are used for spectroscopy.

Stimulated Raman Scattering

In stimulated Raman scattering, the internal mode of the medium is driven by the interference of the pump and Stokes waves while the Stokes wave is driven by the modulation of the pump wave

by the material oscillation. Thus, the growth rate of the Stokes wave, which is determined by the strength of the internal excitation, increases as the Stokes intensity increases. This is exactly the condition needed for exponential growth of the Stokes wave. Stimulated Raman scattering was observed soon after the development of Q-switched lasers.¹⁴

Stimulated Raman Geometries Stimulated Raman scattering can occur in several geometries, each with its own applications. Stimulated Raman generators are illustrated in Fig. 2*a* and *b*. In this arrangement, only a pump wave is incident on the medium, and the Stokes wave grows from quantum noise. Forward scattering is the most common form of Raman generator (Fig. 2*a*), but backward wave generation can occur under some conditions (Fig. 2*b*). Single-pass Raman generators involve amplified spontaneous emission (ASE), and the coherence and divergence properties are characteristic of ASE devices. If the interaction is strong enough to involve pump depletion, the coherence of the Stokes wave can approach that of the pump wave. Raman generators are typically used for frequency conversion from the pump to the Stokes wavelength or for creation of a Stokes wave for other applications. Conversion to multiple Stokes or anti-Stokes waves can also be done with Raman generators. Mirrors can be used with the Raman generator to create a Raman oscillator as in Fig. 2*c*. This arrangement is used with low-intensity continuous-wave (cw) or pulsed pump lasers to reduce the Raman threshold.¹³

Raman amplification can be achieved with the arrangements shown in Fig. 2*d* and *e*, in which a Stokes wave is supplied along with the pump wave. The Stokes wave is amplified at the expense of the pump wave. The Stokes wave is usually created in a separate low-power Raman generator and may be spatially filtered for control of the spatial divergence. Raman amplification is used when a high-quality Stokes beam is desired or when highly efficient conversion of the pump energy is desired without creation of multiple Stokes components. Backward amplification (Fig. 2*e*) is used for Stokes amplification or for pulse compression. In Raman generators, usually only a single Stokes frequency is generated corresponding to the material mode with the highest gain. Amplification of any of the Raman modes is possible if a suitable input Stokes signal is provided.

Coherent anti-Stokes Raman scattering (CARS) is illustrated in Fig. 2*f*. Here the pump and Stokes waves are supplied, often having comparable intensities and propagating at the appropriate phase-matching angle. The anti-Stokes wave is created at the appropriate angle for phase matching. CARS is used for spectroscopy and other diagnostic applications.

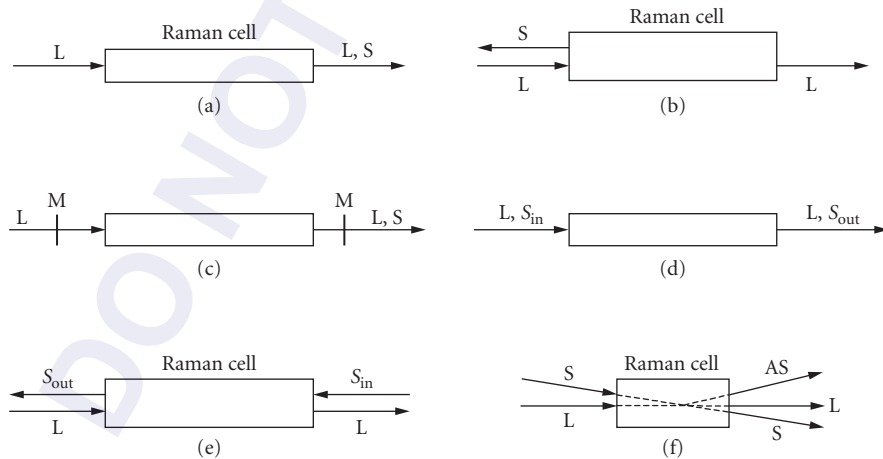


FIGURE 2 Schematic diagrams of various Raman interaction geometries. (a) Forward Raman generator. (b) Backward Raman generator. (c) Raman oscillator. (d) Forward Raman amplifier. (e) Backward Raman amplifier. (f) Coherent anti-Stokes four-wave mixing. AS, generated anti-Stokes mixing; L, pump laser; M, mirrors; S, Stokes wave; S_{in} , Stokes wave from Raman generator; S_{out} , amplified Stokes wave.

Raman Susceptibilities Classically, the Raman effect occurs because of modulation of the polarization in the medium induced by the pump at the difference frequency between the pump and material mode. This arises from the hyperpolarizability of the medium. The polarization of the medium is given by^{6,7}

$$P = \epsilon_0 \left(\mu_0 + \mu E + \frac{\partial \alpha}{\partial Q} Q E + \dots \right) \quad (1)$$

where μ_0 is the permanent dipole moment of the material, μ is the first-order dipole moment, $(\partial \alpha / \partial Q)$ is the hyperpolarizability, and Q is the normal mode coordinate of the material oscillation. Spontaneous Raman scattering is given classically by the relation

$$P_s = NA \frac{\partial \sigma}{\partial \Omega} d\Omega P_L L \quad (2a)$$

where $d\sigma/d\Omega$ is the differential cross section of the Raman transition, L is the length of the scattering medium, A is the cross-sectional area, and $d\Omega$ is its solid angle. The Raman cross section is related to the hyperpolarizability by

$$\frac{\partial \sigma}{\partial \Omega} = \frac{\pi h n_s v_s^4}{4 n_L m \omega_0 c^2} \left(\frac{\partial \alpha}{\partial Q} \right)^2 \quad (2b)$$

The Raman effect can also be analyzed quantum-mechanically through the use of the density matrix¹⁵ using the third-order elements ρ_{0i} and ρ_{i2} and the second-order element ρ_{02} , where the subscripts 0 and 2 designate the initial and final Raman levels and the subscript i designates intermediate levels. The classical parameter Q is related to the off-diagonal density matrix element ρ_{02} , while the population density is related to the element ρ_{22} . Thus the normal mode coordinate of the material oscillation is associated with quantum mechanical transition probabilities and not directly with the population of the excited state of the material.

The Raman cross section can be related to the susceptibilities of nonlinear optics¹⁵ through density matrix calculations involving the third-order elements ρ_{0i} and ρ_{i2} and the second-order element ρ_{02} . The nonlinear Raman polarization amplitude is

$$P(\omega_s) = \frac{3}{2} \epsilon_0 \chi^{(3)}(-\omega_s, \omega_L, -\omega_L, \omega_s) |A_L|^2 A_s \quad (3a)$$

and the Raman susceptibility per atom or molecule is

$$\chi_{\text{Raman}}^{(3)} = \frac{1}{6\hbar^3 \epsilon_0} \frac{1}{\omega_{20} - (\omega_L - \omega_s) + i/T_2} \left| \sum_i \mu_{0i} \mu_{i2} \left(\frac{1}{\omega_{i0} - \omega_L} + \frac{1}{\omega_{i0} + \omega_s} \right) \right|^2 \quad (3b)$$

where T_2 is the homogeneous Raman dephasing time and M_{ij} is the transition dipole moment between states i and j .

The Raman susceptibility has real and imaginary parts. The imaginary part, which is negative and proportional to T_2 , is responsible for spontaneous and stimulated Raman scattering. The real part contributes to the nonlinear refractive index through Raman-type interactions. The Raman susceptibility is related to the hyperpolarizability as

$$\chi'' = -\frac{1}{12\Gamma m \omega_0} \left(\frac{\partial \alpha}{\partial Q} \right)^2 \quad (4)$$

where χ'' is the imaginary part of the Raman susceptibility, Γ is the inverse of the dephasing time, m is the effective mass of the material oscillator, and ω_0 is the frequency of the material transition. The Raman susceptibility is defined only for steady-state interactions. In transient interactions, the Raman polarization must be solved for as a dynamic variable along with the optical fields.

SRS Equations Stimulated Raman scattering is described in the most general form by the equations:²

$$\nabla_{\perp}^2 A_S + 2ik_S \left(\frac{\partial A_S}{\partial z} \pm \frac{1}{v_{gS}} \frac{\partial A_S}{\partial t} \right) = -\frac{\Delta N}{2} \frac{\omega_S^2}{c^2} \left(\frac{\partial \alpha}{\partial Q} \right)_S Q^* A_L \quad (5a)$$

$$\nabla_{\perp}^2 A_L + 2ik_L \left(\frac{\partial A_L}{\partial z} + \frac{1}{v_{gL}} \frac{\partial A_L}{\partial t} \right) = -\frac{\Delta N}{2} \frac{\omega_L^2}{c^2} \left[\left(\frac{\partial \alpha}{\partial Q} \right)_S Q A_S + \left(\frac{\partial \alpha}{\partial Q} \right)_{AS} Q^* A_{AS} e^{i\Delta kz} \right] \quad (5b)$$

$$\nabla_{\perp}^2 A_{AS} + 2ik_{AS} \left(\frac{\partial A_{AS}}{\partial z} \pm \frac{1}{v_{gAS}} \frac{\partial A_{AS}}{\partial t} \right) = -\frac{\Delta N}{2} \frac{\omega_{AS}^2}{c^2} \left(\frac{\partial \alpha}{\partial Q} \right)_{AS} Q A_L e^{-i\Delta kz} \quad (5c)$$

$$\frac{\partial Q^*}{\partial t} + \left[i \left(\frac{\omega_0^2 - \Omega^2}{2\Omega} \right) + \Gamma \right] Q^* = -\frac{i}{4\omega_0 m} \left[\left(\frac{\partial \alpha}{\partial Q} \right)_S A_L^* A_S + \left(\frac{\partial \alpha}{\partial Q} \right)_{AS} A_{AS}^* A_L e^{-i\Delta kz} \right] \quad (5d)$$

$$\frac{\partial \Delta N}{\partial t} + \frac{1}{T_1} (\Delta N - N) = -i \frac{\Delta N \epsilon_0}{4\hbar} \left[\left(\frac{\partial \alpha}{\partial Q} \right)_S (Q^* A_L A_S^* - Q A_L^* A_S) + \left(\frac{\partial \alpha}{\partial Q} \right)_{AS} (Q^* A_L^* A_{AS} e^{i\Delta kz} - Q A_L A_{AS}^* e^{-i\Delta kz}) \right] \quad (5e)$$

where A_S , A_L , and A_{AS} are the slowly varying optical field amplitudes of the Stokes, laser (pump), and anti-Stokes waves given by

$$E_{S,L,AS}(x, y, z, t) = \frac{1}{2} \left[A_{S,L,AS}(x, y, z, t) e^{-i(\omega_{S,L,AS} t - k_{S,L,AS} z)} + c.c. \right] \quad (6)$$

where $\omega_{S,L,AS}$ are the optical frequencies, $k_{S,L,AS}$ are the k vectors, $v_{S,L,AS}$ are the group velocities, ΔN is the difference in the population between the lower and the upper Raman transition levels, $\Delta N = N_0 - N_2$, N is the total population density, and Q is the amplitude of the normal mode coordinate of the material excitation defined by

$$Q(x, y, z, t) = \frac{1}{2} \left[Q(x, y, z, t) e^{-i(\Omega t - k_0 z)} + c.c. \right] \quad (7)$$

$(\partial \alpha / \partial Q)_S$ and $(\partial \alpha / \partial Q)_{AS}$ are the hyperpolarizabilities for the Stokes and anti-Stokes waves, respectively, ω_0 is the Raman transition frequency, and k_0 is the nonpropagating k vector of the material excitation. Γ is the half-width at half-maximum of the Raman linewidth given by $\Gamma = 1/T_2$, Δk is the phase mismatch of the anti-Stokes scattering, m is the effective reduced mass of the material oscillation, and ϵ_0 is the permittivity of the free space.

The relation between various frequencies and k vectors is

$$\omega_s = \omega_L - \Omega \quad (8a)$$

$$\omega_{AS} = \omega_L + \Omega \quad (8b)$$

$$k_o = k_L - k_s \quad (8c)$$

$$\Delta k = k_s + k_{AS} - 2k_L \quad (8d)$$

The intensity of the various light waves is given by

$$I_{S,L,AS} = \frac{1}{2} (cn\epsilon_o) |A_{S,L,AS}|^2 \quad (9)$$

Phase matching, which is of central importance in many nonlinear interactions, is automatically satisfied in stimulated Raman scattering because k_o automatically adjusts to satisfy Eq. (8c). Phase matching in anti-Stokes scattering through four-wave mixing is not automatic, and phase-matching conditions determine many of the properties of the four-wave mixing anti-Stokes scattering interaction.

These equations describe most of the effects that are commonly encountered in Raman scattering, including amplification, diffraction, growth from noise, multiple Stokes generation, forward and backward interactions, coherent anti-Stokes interactions, and transient and steady-state interactions. Each of these will be discussed individually in subsequent parts of this chapter. In most of these effects the population of the ground state is undisturbed and the approximation $\Delta N = N$ is valid. In some situations involving high-power lasers or resonant interactions, ΔN will vary and Eq. (5e) must be used.

Steady-State Stokes Scattering Many of the effects associated with stimulated Raman scattering can be illustrated by considering plane-wave steady-state forward Raman amplification. In this interaction, an incident Stokes wave at frequency ω_s is amplified by a pump wave at frequency ω_L in the geometry of Fig. 2d. The plane-wave steady-state Raman amplification equations are obtained from Eq. (5a and b) by assuming that the time variation of the wave envelopes is slow compared to the dephasing time T_2 , allowing the time derivatives to be neglected; by neglecting the transverse spatial derivative; and by assuming that the ground state population is unchanged, allowing Eq. (5e) to be neglected and the ground state population $N_o = N$ to be treated as a constant. The resulting equations are

$$\frac{dA_s}{dz} = i\kappa_2 Q^* A_L \quad (10a)$$

$$\frac{dA_L}{dz} = i \frac{\omega_L n_s}{\omega_s n_L} \kappa_2 A_s Q \quad (10b)$$

$$Q^* = -i \frac{1}{4\omega_o m} \left(\frac{\partial \alpha}{\partial Q} \right) \frac{1}{\Gamma + i(\omega_o - \Omega)} A_L^* A_s \quad (10c)$$

where

$$\kappa_2 = \frac{N\omega_s}{4n_s c} \left(\frac{\partial \alpha}{\partial Q} \right) \quad (11)$$

Equation (10c) results from solution of Eq. (5d) neglecting the time derivative and the anti-Stokes term.

Steady-state gain If the pump wave intensity is taken as a constant, then Eq. (10a) can be solved for the Stokes intensity:

$$I_S(L) = I_S(0)e^{g_{ss}L} \quad (12)$$

where L is the length of the interaction region and g_{ss} is the steady-state Raman gain coefficient, given by

$$g_{ss} = \frac{N\omega_s}{4\Gamma n_s n_L m \omega_o c^2 \epsilon_o} \left(\frac{\partial \alpha}{\partial Q} \right)^2 \quad (13)$$

This result is also written in the literature in the forms

$$I_S(L) = I_S(0)e^{G_{ss}} \quad (14a)$$

$$I_S(L) = I_S(0)e^{N\sigma_R I_L L} \quad (14b)$$

where

$$G_{ss} = g_{ss} I_L L \quad (15a)$$

is the total steady-state Raman gain and

$$\sigma_R = \frac{\omega_s}{4\Gamma n_s n_L m \omega_o c^2 \epsilon_o} \left(\frac{\partial \alpha}{\partial Q} \right)^2 \quad (15b)$$

is the stimulated Raman cross section. The gain coefficient is given in terms of the Raman susceptibility as

$$g_{ss} = -\frac{3N\omega_s}{c^2 n_s n_L \epsilon_o} \chi'' \quad (16)$$

The classical Raman scattering cross section is related to the Raman gain, given in Eq. (13), by

$$g_{ss} = \frac{2Nc^2}{\pi n_s^2 h v_s^3 \Delta v_R} \frac{\partial \sigma}{\partial \Omega} \quad (17)$$

where Δv_R is the Raman linewidth [full width at half-maximum (FWHM)] related to the dephasing time as

$$\Delta v_R = \frac{1}{\pi T_2} \quad (18)$$

In transparent regions of the spectrum, the Raman cross section scales as v_s^4 and the Raman gain coefficient scales as v_s . In dispersive regions of the spectrum, as for example when the pump wavelength is in the ultraviolet, additional frequency variation in the Raman gain arises from the resonant term $1/(\omega_o - \omega_L)$ in the susceptibility.

Initially, stimulated Raman gain coefficients were calculated from measurements of the spontaneous Raman scattering cross section or estimated from measurements of the Raman threshold. The most accurate determinations of Raman gain coefficients are now made with steady-state amplification measurements in the low-gain regime.

The conditions for steady-state Raman amplification are encountered when both the pump and incident Stokes radiation are either cw or narrowband pulses with pulse duration longer than a steady-state time given by¹⁶⁻¹⁹

$$t_{ss} = G_{ss} T_2 \quad (19)$$

and contain no rapid internal temporal variation, requiring that the linewidths $\Delta\nu_L$ and $\Delta\nu_S$ satisfy the condition

$$\Delta\nu_{L(S)} \ll \Delta\nu_R \quad (20)$$

Raman linewidths The steady-state Raman gain scales as the medium density and inversely with the Stokes wavelength and the material linewidth. In materials such as some molecular gases, the linewidth is pressure dependent due to pressure broadening with a variation of the form:²⁰

$$\Delta\nu_R = \Delta\nu_o + \beta\rho \quad (21)$$

where ρ is the density of the material and $\Delta\nu_o$ is the Raman linewidth at low pressures. In materials that exhibit this behavior, the gain increases at low pressures but levels off at higher pressures, becoming independent of pressure in the limit of high pressures. In hydrogen, for example, the steady-state gain for the Q(1) vibrational mode is effectively independent of pressure for pressures above about 10 atmospheres. The rotational lines in some gases such as hydrogen reach their pressure-broadened limit at lower pressures than do the vibrational lines. As a result, rotational Raman scattering in these materials is more prominent at low pressures, while the vibrational scattering is dominant at higher pressures.

The linewidth for forward scattering in some gases such as hydrogen exhibits Dicke narrowing,²¹ in which inelastic collisions cause a narrowing of the linewidth at small but nonzero pressures before the material enters a pressure-broadened regime. The linewidth exhibits a variation in a region above some cutoff pressure of the form:²²

$$\Delta\nu_R = \frac{A}{\rho} + B\delta\rho \quad (22)$$

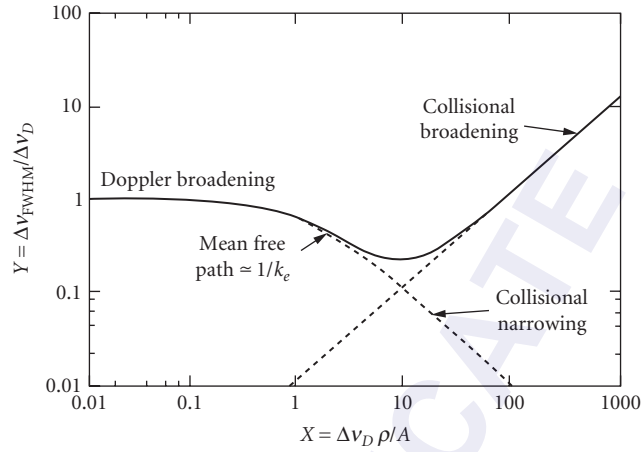
This behavior is shown for hydrogen in Fig. 3. Raman frequency shifts for a number of materials are shown in Table 2. Linewidths and gain coefficients for selected materials are given in Table 3. The temperature dependence of the frequency shift parameters for hydrogen is given in Table 4. Temperature dependence of line-broadening parameters for hydrogen and nitrogen is given in Table 5. Formulas for quantities appropriate for Raman scattering in molecular gases are given in Table 6.

Pump depletion If the Stokes intensity becomes large enough, the depletion of the pump wave must be taken into account and Eq. (10a and b) must be solved together.

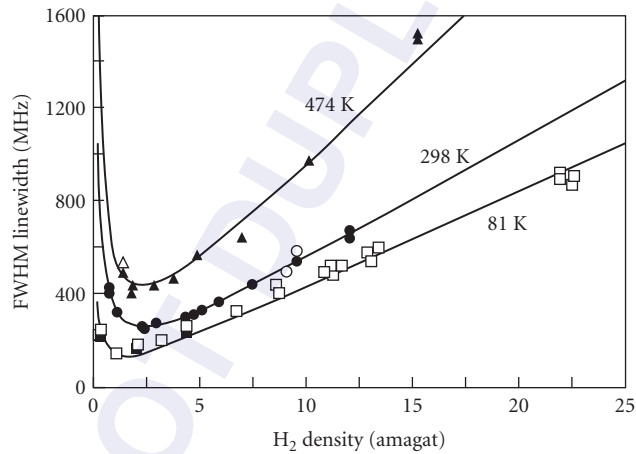
The solutions are

$$I_S(z) = \frac{\left[I_S(0) + \frac{\omega_S}{\omega_L} I_L(0) \right] \frac{\omega_L}{\omega_S} \frac{I_S(0)}{I_L(0)} e^{\left[I_L(0) + \frac{\omega_L}{\omega_S} I_S(0) \right] gz}}{1 + \frac{\omega_L}{\omega_S} \frac{I_S(0)}{I_L(0)} e^{\left[I_L(0) + \frac{\omega_L}{\omega_S} I_S(0) \right] gz}} \quad (23a)$$

$$I_L(z) = \frac{\left[1 + \frac{\omega_S}{\omega_L} \frac{I_L(0)}{I_S(0)} \right] I_L(0) e^{-\left[I_L(0) + \frac{\omega_L}{\omega_S} I_S(0) \right] gz}}{1 + \frac{\omega_S}{\omega_L} \frac{I_L(0)}{I_S(0)} e^{-\left[I_L(0) + \frac{\omega_L}{\omega_S} I_S(0) \right] gz}} \quad (23b)$$



(a)



(b)

FIGURE 3 Calculated (a) and experimental (b) variation with pressure of the linewidth of the Q(1) transition in hydrogen gas showing Dicke narrowing and pressure broadening. (From Ref. 22; copyright 1986 by the American Physical Society.)

When $I_s(0) \ll I_L(0)$, the pump depletion solutions simplify to

$$I_s(z) = \frac{I_s(0)e^{I_L(0)gz}}{1 + \frac{\omega_L I_s(0)}{\omega_S I_L(0)} e^{I_L(0)gz}} \quad (24a)$$

$$I_L(z) = \frac{\left[1 + \frac{\omega_S I_L(0)}{\omega_L I_s(0)}\right] I_L(0) e^{-I_L(0)gz}}{1 + \frac{\omega_S I_L(0)}{\omega_L I_s(0)} e^{-I_L(0)gz}} \quad (24b)$$

TABLE 2 Raman Transition Frequencies of Selected Materials

Liquids		
Substance	ν_R (cm ⁻¹)	Reference
Bromoform	222	23
Tetrachloroethylene	448	24
Carbon tetrachloride ^a	460	25
Ethyl iodide	497	26
Hexafluorobenzene ^a	515	25
Bromoform	539	23
Chlorine	552	2
Methylene bromide	580	25
Trichloroethylene	640	23
Carbon disulfide	655	27
Ethylene bromide	660	28
Chloroform	667	23
<i>o</i> -Xylene	730	29
FC104 ^b	757	30
Sulfur hexafluoride	775	31
α -Dimethylphenethylamine	836	32
Dioxane	836	23
Morpholine ^a	841	25
Thiophenol ^a	916	25
Nitromethane ^a	927	25
Deuterated benzene	944	14
Potassium dihydrogen phosphate	980	33
Cumene ^a	990	25
Pyridine	991	14
1,3-Dibromobenzene	992	24
Benzene	992	14
Aniline	997	34
Styrene	998	35
<i>m</i> -Toluidine ^a	999	25
Acetophenone	999	36
Bromobenzene	1000	34
Chlorobenzene ^a	1001	25
<i>tert</i> -Butylbenzene	1000	24
Benz aldehyde ^c	1001	24
Ethylbenzoate	1001	36
Benzonitrile	1002	34
Ethylbenzene	1002	29
Toluene	1004	14
Fluorobenzene	1012	37
γ -Picoline	1016	25
<i>m</i> -Cresol ^a	1029	25
<i>m</i> -Dichlorobenzene ^a	1034	25
1-Fluoro-2-chlorobenzene ^d	1034	24
Iodo-Benzene ^a	1070	25
Benzoyl chloride ^a	1086	25
Benzaldehyde ^a	1086	25
Anisole ^a	1097	25
Pyrrrole ^a	1178	25
Furan ^a	1180	25
Nitrous oxide	1289	31
Styrene	1315	35

(Continued)

TABLE 2 Raman Transition Frequencies of Selected Materials (*Continued*)

Liquids		
Substance	ν_R (cm ⁻¹)	Reference
Nitrobenzene	1344	14
1-Bromonaphthalene	1363	14
1-Chloronaphthalene	1374	38
2-Ethyl-naphthalene	1382	24
<i>m</i> -Nitrotoluene ^a	1389	25
Carbon dioxide	1392	31
Quinoline ^a	1427	25
Bromocyclohexane	1438	26
Furan ^a	1522	25
Methyl salicylate ^a	1612	25
Cinnamaldehyde	1624	38
Styrene	1631	35
3-Methylbutadiene	1638	39
Pentadiene	1655	39
Isoprene	1792	32
1-Hexyne	2116	24
Dimethyl sulfoxide ^c	2128	40
<i>o</i> -Dichlorobenzene ^a	2202	25
Benzonitrile	2229	38
Acetonitrile	2250	26
1,2-Dimethylaniline ^a	2292	25
Nitrogen	2327	41
Hydrobromic acid	2493	30
Hydrochloric acid	2814	42
Methylcyclohexane ^a	2817	25
Methanol	2831	23
<i>cis,trans</i> -1,3-Dimethylcyclohexane	2844	24
Tetrahydrofuran	2849	38
Cyclohexane	2852	14
<i>cis</i> -1,2-Dimethylcyclohexane	2853	24
α -Dimethylphenethylamine	2856	32
Dioxane	2856	23
Decahydronaphthalene	2860	30
Cyclohexane	2863	23
Cyclohexanone	2863	29
<i>cis,trans</i> -1,3-Dimethylcyclohexane	2866	24
<i>cis</i> ,1,4-Dimethylcyclohexane	2866	24
Cyclohexane	2884	23
Dichloromethane ^a	2902	25
Dimethyl sulfoxide	2916	40
Morpholine ^a	2902	25
Cargille 5610 ^f	2908	30
2,3-Dimethyl-1,5-hexadiene	2910	24
Limonene	2910	32
<i>o</i> -Xylene	2913	29
1-Hexyne	2916	24
<i>cis</i> -2-Heptene	2916	24
2-Octene	2918	24
Acetonitrile	2920	30
Mesitylene	2920	32
2-Bromopropane	2920	24

TABLE 2 Raman Transition Frequencies of Selected Materials (*Continued*)

Liquids		
Substance	ν_R (cm ⁻¹)	Reference
Acetone	2921	29
Ethanol	2921	23
<i>cis</i> -1,2-Dimethylcyclohexane	2921	24
Carvone	2922	32
2-Chloro-2-methylbutane	2927	24
Dimethylformamide	2930	23
<i>Cis,trans</i> -1,3-Dimethylcyclohexane	2926	24
<i>m</i> -Xylene	2933	29
1,2-Diethyl tartrate	2933	32
<i>o</i> -Xylene	2933	29
Piperidine	2933	29
1,2-Diethylbenzene	2934	24
1-Bromopropane	2935	24
Piperidine	2936	29
Tetrahydrofuran	2939	38
Decahydronaphthalene	2940	30
Piperidine	2940	29
Cyclohexanone	2945	29
2-Nitropropane	2945	24
1,2 Diethyl carbonate ^a	2955	25
1,2 Dichloroethane ^a	2956	25
<i>trans</i> -Dichloroethylene	2956	23
Methyl fluoride	2960	31
1-Bromopropane	2962	24
2-Chloro-2-methylbutane	2962	24
α -Dimethylphenethylamine	2967	32
Dioxane	2967	23
Methyl chloride	2970	31
Cyclohexanol ^a	2982	25
Cyclopentane ^a	2982	25
Cyclopentanol ^a	2982	25
Bromocyclopentane ^a	2982	25
<i>o</i> -Dichlorobenzene ^a	2982	25
<i>p</i> -Chloro toluene ^a	2982	25
α -Picoline ^a	2982	25
<i>p</i> -Xylene	2988	29
<i>o</i> -Xylene	2992	29
Dibutyl-phthalate ^a	2992	25
1,1,1-Trichloroethane	3018	23
Ethylene chlorohydrin ^a	3022	25
Isophorone ^a	3022	25
Nitrosodimethylamine ^a	3022	25
Propylene glycol ^a	3022	25
Cyclohexane ^a	3038	25
Styrene	3056	35
Pyridine	3058	24
Benzene	3064	14
<i>tert</i> -Butylbenzene	3065	24
1-Fluoro-2-chlorobenzene	3082	24
Turpentine ^a	3090	25
Pseudocumene ^a	3093	25

(Continued)

TABLE 2 Raman Transition Frequencies of Selected Materials (*Continued*)

Liquids		
Substance	ν_R (cm ⁻¹)	Reference
Acetic acid ^a	3162	25
Acetylacetone ^a	3162	25
Methyl methacrylate ^a	3162	25
γ -Picoline ^a	3182	25
Aniline	3300	34
Water ^a	3651	25
Solids		
Substance	ν_R (cm ⁻¹)	Reference
Quartz	128	43
Lithium niobate	152	44
α -Sulfur	216	45
Lithium niobate	248	44
Bromine	295	45
Bromine	303	45
Quartz	466	43
α -Sulfur	470	45
Chlorine	543	46
Lithium niobate	628	44
Carbon disulfide	656	46
Potassium iodate	746	47
Potassium bromate ^a	780	47
Potassium periodate	790	47
Potassium bromate	798	47
Potassium chromate	844	47
Sodium molybdate ^a	884	47
Potassium dichromate	906	47
Calcium tungstate	911	38
Potassium dihydrogen phosphate	915	33
Ammonium vanadate	915	47
Sodium tungstate	915	47
Potassium perchlorate ^a	936	47
Potassium chlorate ^a	938	47
Ammonium sulfate ^a	975	47
Potassium sulfate ^a	985	47
Stilbene	997	48
Polystyrene	1001	23
Calcium nitrate	1050	47
Calcium nitrate tetrahydrate ^a	1052	47
Potassium nitrate	1060	47
Magnesium nitrate dehydrate ^a	1060	47
Ammonium nitrate ^a	1062	47
Magnesium nitrate hexahydrate ^a	1063	47
Sodium nitrate	1075	47
77 K (not observed at 293 K) ^a		
Calcite	1086	45
Diamond	1332	45
Naphthalene	1380	38
Anthracene	1403	49
Stilbene	1591	48
Potassium thiocyanate	1040	47

TABLE 2 Raman Transition Frequencies of Selected Materials (*Continued*)

Solids		
Substance	ν_R (cm ⁻¹)	Reference
Potassium ferricyanide ^a	2100	47
Triglycine sulfate	2422	25
Triglycine sulfate	2702	25
Triglycine sulfate	3022	25
Polystyrene	3054	23
Gases		
Substance	ν_R (cm ⁻¹)	Reference
Barium vapor ^g	IR ^h	50
Cesium vapor ^g	IR ^h	51,52
Hydrogen fluoride	FIR ^h	53
Potassium vapor ^g	IR ^h	51,52
Rubidium vapor ^g	IR ^h	54
ρ -H ₂	354	55,56
Carbon tetrafluoride	980	57
Oxygen	1552	29
Nitrogen	2331	58
Potassium vapor	2721	59
Methane	2916	60
Deuterium	2991	60
Hydrogen deuteride	3628	61
Hydrogen	4155	60

^aObserved at low resolution.

^bProduct of 3M Co., St. Paul, Minnesota.

^c1:1 mixture with tetrachloroethylene.

^dVery weak and diffuse.

^eDeuterated.

^fProduct of Cargille Laboratories, Cedar Falls, N.J.

^gStimulated electronic Raman scattering (SERS).

^hGenerally tunable transitions in the infrared (IR) and far infrared (FIR).

Reprinted with permission from M. J. Weber (ed.), *CRC Handbook of Laser Science and Technology: Optical Materials*, part 1, vol. III. Copyright CRC Press, Boca Raton, FL, 1988.

The conditions for exponential Raman amplification can now be identified. In order for the approximate solution of Eq. (12) to be valid, it is necessary for the condition

$$I_S(0)e^{g_{SS}I_L z} \ll I_L(0) \quad (25)$$

to be satisfied. When the condition in Eq. (25) is not satisfied, the full solution of Eq. (23) [or Eq. (24), if $I_S(0) \ll I_L(0)$] must be used. The calculated behavior of the pump and Stokes waves in the depletion regime is shown in Fig. 4.

In the limit of large pump depletion, the Stokes intensity becomes

$$I_S(z) = I_S(0) + \frac{\omega_S}{\omega_L} I_L(0) \quad (26)$$

Equation (26) indicates that the maximum energy that can be added to the Stokes wave is reduced from the energy of the pump wave by the fraction (ω_S/ω_L) . This ratio is termed the *Manly Rowe fraction*. The difference between the energy given up by the pump wave and that gained by the Stokes wave is taken up by energy deposited in the material excitation. Although complete energy

TABLE 3 Linewidths and Gain Coefficients for Selected Materials

Raman transition frequencies, linewidths, and gains of selected materials at room temperature ^a						
Material	Mode	ν_0 (cm^{-1})	$\Delta\nu$ (MHz)	λ_L (nm)	g (cm^2/GW)	Reference
H ₂ gas (20 atm)	Q(1)	4155	309/ ρ + 5.22 <i>p</i>	532	2.5	22, 62, 63
H ₂ gas (high)	S(1)	587	119 <i>p</i>	350	1.2	64
D ₂ gas (60 atm)	Q(2)	2987	101/ ρ + 120 <i>p</i>	532	0.45	63, 65, 66, 67
D ₂ gas	S(2)	414	124 <i>p</i>	350		64
HD	Q(1)	3628	693 <i>p</i>	532		68
HD	S(1)	443	760 <i>p</i>	350	0.098	64
CH ₄ (1115 atm)	ν_1	2913	8220 + 384 <i>p</i>	532	1.26	63, 65
N ₂	Q	2327	4500 ($\rho < 10$)	248	0.003 <i>p</i>	69
N ₂	S(6)	60	3570 <i>p</i>	566	0.0063	70
O ₂	Q	1552	54000	248	0.012 <i>p</i>	69
SiH ₄	Q	2186	15000 (est.)	248	0.19 <i>p</i>	69
GeH ₄	ν_1	2111	15000 (est.)	248	0.27 <i>p</i>	69
CF ₄	ν_1	908	21000 (est.)	248	0.008 <i>p</i>	69
SF ₆	ν_1	775	30000 (est.)	248	0.014 <i>p</i>	69
Liquid N ₂		2326.5	0.067	694	16 ± 5	2
Liquid O ₂		1552	0.117	694	14.5 ± 4	2
H ₂ O		3290	200	694	0.14	1, 2, 71
Benzene		992	2.15		2.8	2
CS ₂		655.6	0.50	694	24	2
Nitrobenzene		1345	6.6		2.1	2
Bromobenzene		1000	1.9		1.5	2
Chlorobenzene		1002	1.6		1.2	2
Toluene		1003	1.94		1.2	2
LiNbO ₃		637	20	694	9.4	2
Ba ₂ Nb ₅ O ₁₅		650		694	6.7	2
LiTaO ₃		201	22	694	4.4	2
SiO ₂		467		694	0.8	2
Ba(NO ₃) ₂		1047	1.5 cm^{-1}	532	47	72–76

Calculated^b and measured gain (g_s) factor for stimulated Raman transitions in liquids (l), gases (g), and solids (s)

Substance	Pump		Frequency (cm^{-1})	Linewidth $\Delta\nu_R$ (cm^{-1})	Gain (g_s) $\times 10^9$ (cm/W)	g_s calc. $\times 10^9$ (cm/W)	g_s at 532 nm ^c	g_s Relative to C_6H_6 (l) ^d	Reference
	Wavelength (nm)	Wavelength (nm)							
Benzene (l)	532	532	992	2.15	5.5 ^e		5.5	1.0	132
	532	532	992	2.15	5.5 ^f		5.5		133
	532	532	992	2.15	4.3 \pm 0.9 ^g		4.3 \pm 0.9		134
Oxygen (l)	694.3	694.3	992	2.15		2.8	5.9 ^{h,i}		135
	694.3	694.3	1552	0.117	16.0 \pm 0.5	14.5 \pm 0.4	21.5	3.9	135
Nitrogen (l)	532	532	2327	0.067	30.0	24.0 \pm 7.0	30.0	5.4	136
	694.3	694.3	2327	0.067	16.0 \pm 0.55	17.0 \pm 0.5	21.5	3.9	135
Carbon disulfide (l)	1060	1060	2327	0.067	10.0	9.0 \pm 3.0	23.2	4.2	137
	1315	1315	2327	0.067	5.0 \pm 2.0	6.0 \pm 2.0	15.6 \pm 6.2	2.9	138
Methanol (l)	694.3	694.3	655.6	0.50	24.0		32.2	5.9	135
	694.3	694.3	2837	18	0.4		0.53	0.10	139
Carbon tetrachloride (l)	597.6	597.6	458		1.3		1.5	0.27	140
	532	532						0.12 ^j	141
Acetone (l)	532	532						0.20 ^j	141
Cyclohexane (l)	532	532						0.25 ^j	141
	694.3	694.3	4155			1.5 ⁱ	2.14		142
Hydrogen ^k (g)	694.3	694.3	4155		1.9 \pm 0.3 ⁱ	1.5 ⁱ	2.7		142
	694.3	694.3	4155		1.5 ⁱ		2.1		143
Nitrogen (g , 1 atm)	353	353	4155		5.0 ⁱ	5.7 ⁱ			131
	308	308	4155			6.7 ⁱ			131
Hydrogen deuteride (g)	10600	10600	354			0.07	0.09		144
	353	353	3628		0.2				131
Nitrogen (g , 1 atm)	694.3	694.3	2330.7		0.0022		0.0030		128
	694.3	694.3	2330.7	0.075		0.0027 ^m	0.0038		145
Carbon tetrafluoride (g , 500 atm)	532	532	980	0.17	24.6 \pm 2 ⁿ		2.46 \pm 2 ⁿ		127
Calcite (s)	532	532	1086	1.2	5.5		5.5		146

(Continued)

TABLE 3 Linewidths and Gain Coefficients for Selected Materials (*Continued*)

Substance	Pump		Frequency (cm^{-1})	Linewidth $\Delta\nu_R$ (cm^{-1})	Gain (g_s) $\times 10^9$ (cm/W)	g_s calc. $\times 10^9$ (cm/W)	g_s at 532 nm ^c	g_s Relative to $C_6H_6(l)^d$	Reference
	Wavelength (nm)	Power (mW)							
Fused quartz (s)	495.4	420	420	18	0.017	0.016	0.016		147
Potassium dihydrogen phosphate (s)	532	918	918	18	0.21 ± 0.05^s	0.21 ± 0.05	0.21 ± 0.05		134

Calculated^b and measured gain (g_s) factor for stimulated Raman transitions in liquids (l), gases (g), and solids (s)

^a ρ is measured in amagats (1 amagat = $2.68 \times 10^{-19} \text{ cm}^{-3}$).

^bOnly those supporting measured values of special interest are given.

^cExcept where noted, $g/532 = (\omega/532)/((\omega_s/\text{meas}) g/\text{meas})$.

^dFor qualitative use only; see previous text.

^eEstimated from simulated threshold.

^fEstimated from stimulated conversion.

^gDirect measurements with single-frequency lasers.

^hExtrapolated using $I \propto \omega_s^3(\omega_a - \omega_s)^{-2}$; $\omega_a = 39,000 \text{ cm}^{-1}$.

ⁱSee Ref. 77 for discussion of peak cross section measurements.

^jRelative threshold measurement.

^kFor detailed analysis of H_2 cross-section pump frequency dependence, see Ref. 63.

^lPressure independent gain; Q(1) transition.

^mCorrected using linewidth of Ref. 58; Q(6) transition.

ⁿTransient gain.

Reprinted with permission from M. J. Weber (ed.), *CRC Handbook of Laser Science and Technology: Optical Materials*, suppl. 2, copyright CRC Press, Boca Raton, FL, 1988; R. W. Boyd, *Nonlinear Optics*, Academic Press, New York, 1992; M. J. Weber (ed.), *CRC Handbook of Laser Science and Technology: Optical Materials*, part 1, vol. III, copyright CRC Press, Boca Raton, FL, 1988.

TABLE 4 Temperature Dependence of Line Shift Parameters for H₂

T (K)	Transition	$\nu_R(0)$ (cm ⁻¹)	C (MHz/amagat)	D (MHz/amagat ²)	Reference
474	Q(l)	4155	9.5 ± 0.9	0.51 ± 0.07	22
298	Q(l)		-96 ± 1		22
298	Q(0)		-64 ± 5		22
81	Q(l)		-305 ± 10		22
81	Q(0)		-250 ± 10		22
81	Q(0)		-336 ± 10		22
<i>p</i> -H ₂					
295	S(0)		6.5 ± 3		78
	S(l)		5.9 ± 1		78
80	S(0)		-22.3 ± 0.6		78
	S(l)		-23 ± 0.6		78

Line shift parameters: $\nu_R = \nu_R(0) + Cp + Dp^2$

Reprinted with permission from M. J. Weber (ed.), *CRC Handbook of Laser Science and Technology: Optical Materials*, suppl. 2. Copyright CRC Press, Boca Raton, FL, 1988.

TABLE 5 Temperature Dependence of Line Broadening Parameters for H₂ and N₂

		H ₂						
T (K)	Transition	A (MHz – amagat)	B (MHz/amagat)	ρ_c (amagat)	Reference			
474	Q(1)	508 ± 29	94.0 ± 2.0	1.23	22			
298	Q(l)	309 ± 11	52.2 ± 0.5	0.92	22			
298	Q(0)	257 ± 12	76.6 ± 0.8	0.79	22			
81	Q(1)	107 ± 20	41.5 ± 0.6	0.63	22			
81	Q(0)	189 ± 40	29 ± 1	1.1	22			
81	Q(0)	76 ± 6	45.4 ± 0.8	0.45	22			
<i>p</i> -H ₂								
295	S(l)	6.15	114 ± 5	0.134	78			
	S(0)	1.87	77 ± 2	0.068	78			
80	S(l)	2.1	110 ± 3	0.088	78			
	S(0)	1.37	67 ± 2	0.095	78			
		N ₂						
T (K)	Transition	ν_R (cm ⁻¹)	λ_s (nm)	B_0 (MHz/amagat)	γ	$\Delta N/N$	g (cm/GW)	Reference
298	S(6)	60	568	3560	0.26 ± .04	0.0282	0.0036	70, 78
	S(8)	76	568	3270	0.33 ± .004	0.0337	0.0046	70, 78
	S(10)	92	568	3060	0.35 ± .002	0.0335	0.0048	70, 78
	S(12)	108	568	2870	0.39 ± .03	0.0289	0.0043	70, 78
195	S(6)		568	3070		0.0485	0.0072	70, 78
	S(8)		568	2860		0.0491	0.0076	70, 78
	S(10)		568	2660		0.0399	0.0065	70, 78
	S(12)		568	2340		0.0271	0.0049	70, 78
80	S(6)		568	2520		0.0897	0.0161	70, 78
	S(8)		568	2120		0.0452	0.0093	70, 78
	S(10)		568	1940		0.0156	0.0034	70, 78
	S(12)		568	1690		0.00378	0.00091	70, 78

Broadening coefficients for H₂: $\Delta\nu_R = A/\rho + B\rho$

$A = 309/\rho (T/298)^{0.92}$ $B = [51.8 + 0.152(T - 298) + 4.85 \times 10^{-4}(T - 298)] \rho$

Broadening coefficients for N₂: $\Delta\nu_R = B_0 \rho (T/295)^7$

Reprinted with permission from M. J. Weber (ed.), *CRC Handbook of Laser Science and Technology: Optical Materials*, suppl. 2. Copyright CRC Press, Boca Raton, FL, 1988.

TABLE 6 Formula Summary for SRS in Molecules 1

Quantity	Formula	Reference
$\partial\sigma/\partial\Omega$ (vibrational, averaged over rotational levels)	$(2\pi\nu_s)^4(\alpha_{01}^2 + 7/45\sum_j f_j b_{jj}\gamma_{01}^2)$	62
$\partial\sigma/\partial\Omega$ (vibrational, single rotational levels)	$(2\pi\nu_s)^4(\alpha_{01}^2 + 4/45b_{jj}\gamma_{01}^2)$	62
	$b_{jj} = J(J+1)(2J-1)(2J+1)$ f_j -relative population of level J For hydrogen: $\gamma_{01} = 1.11 \times 10^{-25} \text{ cm}^{-3}$ $\gamma_{01} = 0.951 \times 10^{-25}$ for $\lambda_L = 500 \text{ nm}$	
$\partial\sigma/\partial\Omega$ (vibrational, wavelength dependence for hydrogen)	$A\nu_s^4/(\nu_i^2 - \nu_L^2)^2$	62
ν_R (Raman shift in H2, pressure dependence)	$\nu_R(0) + C\rho + \Delta\rho^2$ (ρ in amagats)	62
$\Delta\nu_R$ (linewidth of H2, pressure and temperature dependence)	$(309/\rho)(T/298)^{0.92} + [51.8 + 0.152(T - 298) + 4.85 \times 10^{-4}(T - 298)^2]\rho$ (ρ in amagats)	62
$\partial\sigma/\partial\Omega$ (rotational scattering in linear molecules)	$(2/15)(2\pi\nu_s/c)^4\gamma^2[(J+1)(J+2)]/[(2J+1)(2J+3)]$	70
	J is lower-state rotational number $\gamma(\nu) = \alpha_{s\parallel}[v_{\parallel}^2/(v_{\parallel}^2 - \nu^2)] - \alpha_{s\perp}[v_{\perp}^2/(v_{\perp}^2 - \nu^2)]$ For nitrogen: $v_{\parallel} = 1.26 \times 10^5 \text{ cm}^{-1}$ $v_{\perp} = 1.323 \times 10^5 \text{ cm}^{-1}$ $\alpha_{s\parallel} = 2.2 \times 10^{-24} \text{ cm}^{-3}$ $\alpha_{s\perp} = 1.507 \times 10^{-24} \text{ cm}^{-3}$	

Reprinted with permission from M. J. Weber (ed.), *CRC Handbook of Laser Science and Technology, Optical Materials*, suppl. 2. Copyright CRC Press, Boca Raton, FL, 1988.

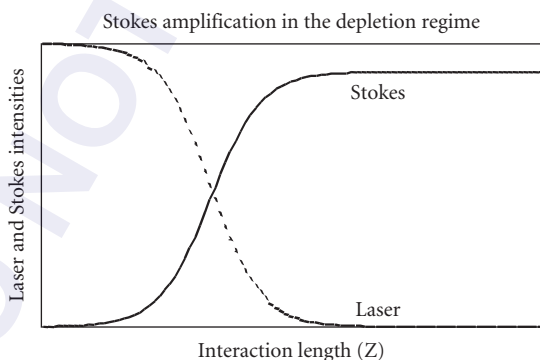


FIGURE 4 Calculated pump and Stokes dependence on interaction length in pump depletion regime.

conversion is not possible with Raman scattering, in principle 100 percent of the pump photons can be converted to the Stokes wave. In practice, high photon conversion efficiency to a single Stokes wave is limited in Raman generators by multiple Stokes or anti-Stokes scattering. High conversion efficiencies are generally obtained in practice in Raman amplifiers or in Raman generators where multiple Stokes generation is unfavorable because of energy level structure. Photon conversion efficiencies in excess of 90 percent have been reported with pulsed lasers.

Gain narrowing Gain narrowing can occur when the Raman gain varies in one of the parameters of the interaction—for example, time, because of pulsed radiation; space, because of a special mode profile; or spectrum, because of use of broadband light; or when a Raman generator is used.

Spectral gain narrowing Spectral gain narrowing is commonly encountered in Raman generators. The spontaneous Stokes emission has the spectral variation of the appropriate Raman transition. For a lorentzian line shape, this variation is given by

$$P_s(\Delta\omega) = P_s(0) \frac{1}{\Delta\omega^2 T^2 + 1} \quad (27)$$

where $\Delta\omega = \omega_o - (\omega_L - \omega_s)$.

This variation with detuning also appears in the Raman gain of Eq. (13) or (15). The Stokes light generated from such a signal can be represented approximately as

$$I(\Delta\omega, L) = I(0) \frac{1}{\Delta\omega^2 T^2 + 1} \exp\left[\frac{G(0)}{\Delta\omega^2 T^2 + 1}\right] \quad (28)$$

The reduced gain at nonzero detunings from the Raman resonance leads to narrowing of the amplified Stokes spectrum relative to the spontaneous spectrum.²² If we use the quantity $\Delta\omega(G)$ to describe the spectral width (FWHM) of the generated or amplified Stokes light, then the spectral width of the amplified Stokes light is

$$\Delta\omega(G) = \Delta\omega(0) \sqrt{\frac{\ln 2/G(0)}{1 - \ln 2/G(0)}} \quad (29)$$

where $G(0)$ is the gain at the center of the line. If we consider the range of values of G that can be achieved without driving the interaction into saturation to be $23 < G_{\max} < 40$, then the maximum gain narrowing that can be experienced without driving the interaction into saturation is on the order of 5.5 to 7.5.

Photon description The steady-state Raman gain can be recast in the formulation of photon interactions. In this situation, the pump and Stokes intensities are given by

$$I_s = h\nu_s n_s \quad (30a)$$

$$I_L = h\nu_L n_L \quad (30b)$$

where n_s and n_L are the photon flux densities (photons per square centimeter per second), which are related to the corresponding photon densities (photons per cubic centimeter) without dispersion as

$$n_s = N_s c \quad (31a)$$

$$n_L = N_L c \quad (31b)$$

Equations for the growth of pump and Stokes photon densities can be derived from quantum mechanical rate equations as⁷

$$\frac{dN_s}{dt} = K(N_s + 1)N_L \quad (32a)$$

$$\frac{dN_L}{dt} = -K(N_L + 1)N_s \quad (32b)$$

where K is the appropriate transition rate for the Raman interaction.

These can be converted to the spatial gain equations for the photon flux density through the relations

$$\frac{dN_s}{dt} = \frac{dn_s}{dz} \quad (33a)$$

$$\frac{dN_L}{dt} = \frac{dn_L}{dz} \quad (33b)$$

giving

$$\frac{dn_s}{dz} = K(n_s + 1)n_L \quad (34a)$$

$$\frac{dn_L}{dz} = -K(n_L + 1)n_s \quad (34b)$$

The 1 in Eq. (34a and b) arises from quantum mechanical commutators and represents spontaneous Raman scattering. When the number of Stokes photons per mode is small compared to unity, the spontaneous Raman emission result is obtained:

$$n_s = Kn_L L \quad (35)$$

where L is the length of the scattering medium. When $n_s \gg 1$, the exponential gain of stimulated Raman scattering is obtained:

$$n_s(L) = n_s(0)e^{Kn_L z} \quad (36)$$

where $K = \hbar \nu g_{SS}$. The relation of spontaneous and stimulated Raman scattering will be discussed again later in the chapter.

Transient Effects When the temporal structure of the pump or Stokes wave varies on a time scale comparable to or shorter than the steady-state time given in Eq. (19), transient effects must be taken into account. In this situation, the integrating effects of the material excitation affect the properties of the Raman interaction. Depending on the situation, this can result in reduced Raman gain for pulses of a given intensity or alteration of the coherence properties of the Stokes radiation. Transient effects can occur for short pump pulses or broadband pump and/or Stokes waves.

Pulsed transient effects Transient Raman scattering with pulsed Raman radiation is described by the equations

$$\frac{\partial A_s(z,t)}{\partial z} + \frac{n_s}{c} \frac{\partial A_s(z,t)}{\partial t} = i\kappa_2 Q^*(z,t)A_L(z,t) \quad (37a)$$

$$\frac{\partial A_L(z,t)}{\partial z} + \frac{n_L}{c} \frac{\partial A_L(z,t)}{\partial t} = i\frac{\omega_L n_s}{\omega_s n_L} \kappa_2 Q(z,t)A_s(z,t) \quad (37b)$$

$$\frac{\partial Q^*(z,t)}{\partial t} + \Gamma Q^* = -i\kappa_1 A_s(z,t)A_L^*(z,t) \quad (37c)$$

where

$$\kappa_1 = \frac{1}{4\omega_0 m} \left(\frac{\partial \alpha}{\partial Q} \right) \quad (38)$$

and κ_2 is given in Eq. (11).

In the most general form, when pump depletion is involved, these equations must be solved numerically. When pump depletion and dispersion can be neglected, the pump field can be taken as a prescribed function of z and t , and the equations take on the form

$$\frac{\partial A_S(z', \tau)}{\partial z'} = i\kappa_2 Q^*(z', \tau) A_L(\tau) \quad (39a)$$

$$\frac{\partial Q^*(z', \tau)}{\partial \tau} + \Gamma Q^*(z', \tau) = -i\kappa_1 A_S(z', \tau) A_L^*(\tau) \quad (39b)$$

where z' and τ are coordinates moving with the common velocity of the pump and Stokes pulses:

$$z' = z \quad (40a)$$

$$\tau = t - z/c \quad (40b)$$

An integral solution for these equations can be written as¹⁶⁻¹⁹

$$A_S(z, \tau) = A_S(0, \tau) + \sqrt{\kappa_1 \kappa_2} z A_L(\tau) \int_{-\infty}^{\tau} \frac{e^{-\Gamma(\tau-\tau')} A_L^*(\tau') A_S(0, \tau') I_1 \left(\sqrt{4\kappa_1 \kappa_2} z [p(\tau) - p(\tau')] \right)}{\sqrt{p(\tau) - p(\tau')}} d\tau' \quad (41a)$$

where I_0 and I_1 are modified Bessel functions and

$$p(\tau) = \int_{-\infty}^{\tau} |A_L(\tau')|^2 d\tau' \quad (42)$$

is proportional to the total pump energy integrated to time τ .

In the extreme transient regime, when $t_p \ll T_2$, and when the incident Stokes pulse has the same functional form as the pump, an analytic solution for the Stokes intensity can be found:

$$I_S(z, \tau) = I_S(0, \tau) I_0^2(u(z, \tau)) \quad (43)$$

where

$$u(z, \tau) = \sqrt{4\kappa_1 \kappa_2} z p(\tau) = \sqrt{2g_{ss} \Gamma z} \int_{-\infty}^{\tau} I_L(\tau') d\tau' \quad (44)$$

and g_{ss} is the steady-state gain coefficient as given in Eq. (13).

When the conditions for Eq. (43) are valid, the solution can be approximated for values of u greater than about 3, corresponding to intensity amplifications of about 24, with the first term in the asymptotic expansion:

$$I_S(z, \tau) = I_S(0, \tau) \frac{e^{2u}}{2\pi u} \quad (45)$$

Quantities other than the instantaneous intensity, such as power density, energy density, or total pulse energy, are useful for characterizing transient Raman measurements done with short pulses. Expressions corresponding to Eq. (45) and for the first and second terms of the expansion of Eq. (43) in u are given in Table 7. The variation of these quantities with the transient gain parameter

TABLE 7 Analytic Expressions and High-Gain Limiting Forms for Small-Signal Transient Raman Amplification of Various Quantities^{19,a}

	Intensity, I_S^b	Energy Density, \mathcal{W}_S^c	Power, \mathcal{P}_S^d	Energy, \mathcal{E}_S^e
Approximation	$I_S(r, z, t) = \frac{I_S(r, z, t)}{I_S(r, 0, t)}$	$\mathcal{W}_S(r, z) = \frac{\int_{-\infty}^{\infty} I_S(r, z, t') dt'}{\int_{-\infty}^{\infty} I_S(r, 0, t') dt'}$	$\mathcal{P}_S(r, z) = \frac{\int_0^{\infty} I_S(r, z, t) dr^2}{\int_0^{\infty} I_S(r, 0, t) dr^2}$	$\mathcal{E}_S(r, z) = \frac{\int_0^{\infty} \int_{-\infty}^{\infty} I_S(r, z, t') dt' dr^2}{\int_0^{\infty} \int_{-\infty}^{\infty} I_S(r, 0, t') dt' dr^2}$
Exact	$I_0^2(u)$	$I_0^2(u) - I_1^2(u)$	$I_0^2(u) - I_1^2(u)$	$I_0^2(u) - 2I_1^2(u) + I_0(u)I_2(u)$
High-gain limit (1st term)	$e^{2u}/2\pi u$	$e^{2u}/2\pi u^2$	$e^{2u}/2\pi u^2$	$e^{2u}/2\pi u^3$
High-gain limit (1st 2 terms)	$(e^{2u}/2\pi u)(1 + 1/4u)$	$(e^{2u}/2\pi u^2)(1 + 1/4u)$	$(e^{2u}/2\pi u^2)(1 + 1/4u)$	$(e^{2u}/2\pi u^3)(1 + 3/4u)$

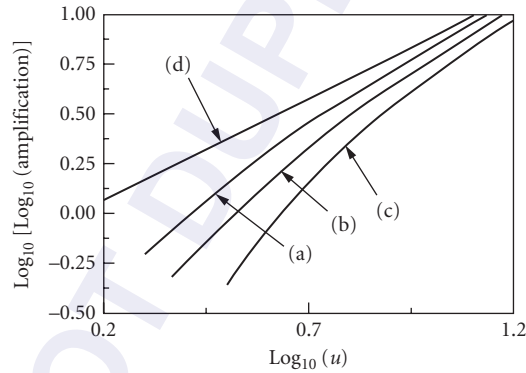
$$^a u = u_m(r, z, t) = 2[\kappa_1 \kappa_2 \int_{-\infty}^t |A_L(r, t')|^2 dt']^{1/2}$$

$$^b u = u_m(r, z, t)$$

$$^c u = u_m(r, z, \infty)$$

$$^d u = u_m(0, z, t)$$

$$^e u = u_m(0, z, \infty)$$

**FIGURE 5** Theoretical dependence of the small-signal amplification on the transient gain parameter u for (a) the intensity, (b) the energy density or power, and (c) the energy as predicted by the exact Bessel function solution. The exponential form $\exp(2u)$ is shown for comparison in (d). (From Ref. 19.)

u is shown in Fig. 5. These solutions are useful for modeling the approximate properties of transient interactions. However, in most practical situations the incident pump and Stokes pulses do not have the same functional form, and the more general integral solution of Eq. (41a) must be used for accurate results, taking into account the specific variation of the Stokes amplitude and phase and the relative timing between the incident Stokes and pump pulses. An approximate analytic expression has been given in the limit of high conversion in which the Stokes pulse evolves to an approximate constant form and the pump pulse is described by regular Bessel functions.⁷⁹

In the extreme transient regime $t_p \ll T_2$, growth of the Stokes pulse in the leading part of the pulse is dominated by the Bessel function and the Stokes pulse rises more slowly than the pump. In the trailing part of the pulse, the integral contribution remains fairly constant and the Stokes amplitude

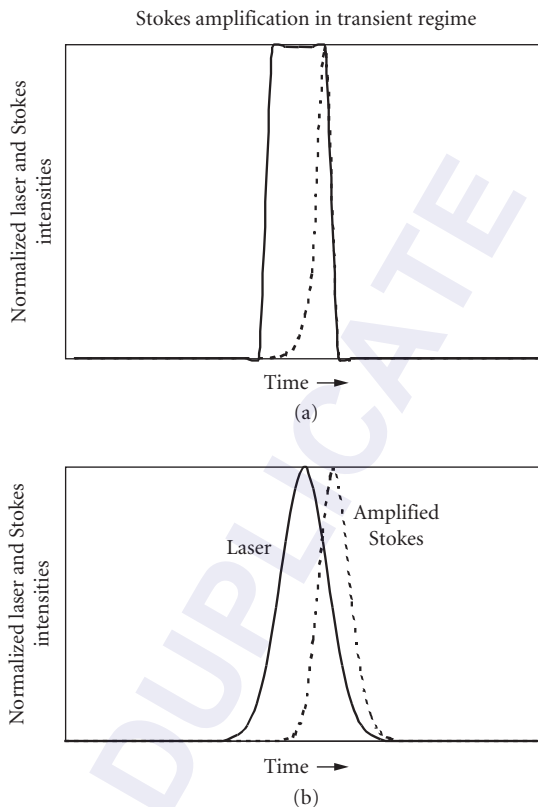


FIGURE 6 Calculated Stokes intensity (dotted line) and pump (solid line) for square pump pulses (a) and gaussian pump pulses (b). The input Stokes signal was assumed to be constant in both cases.

follows the pump amplitude.¹⁷ The result is that the Stokes pulse is delayed relative to the pump and is shorter than the pump. Ideal square pulses can be shortened by an arbitrarily large factor, while more realistic gaussian pulses can be shortened by about a factor of 2.^{17,80} The transient response for square and gaussian pulses is illustrated in Fig. 6. Optimal amplification of the Stokes pulse has been shown to require that the incident Stokes pulse arrive earlier than the pump by about half the pulse duration.¹⁹ An example of experimentally measured transient Raman amplification is given in Fig. 7 along with a theoretical comparison. The theoretical curve was obtained by integrating the square magnitude of Eq. (41a) over space and time using the experimentally measured amplitude and phase variations of the incident Stokes pulse. Agreement between experiment and theory was obtained over approximately eight orders of magnitude using one adjustable parameter in the region below pump depletion.

The peak intensity amplification in the extreme transient limit is reduced from the steady-state amplification for pump pulses with the same peak intensity. The intensity amplification grows as the pulse length increases, approaching the steady-state value when $t_p \gg G_{ss} T_2$. For pulses of a given energy, the total integrated energy amplification increases steadily as the pulse duration decreases, reaching a maximum in the extreme transient regime.

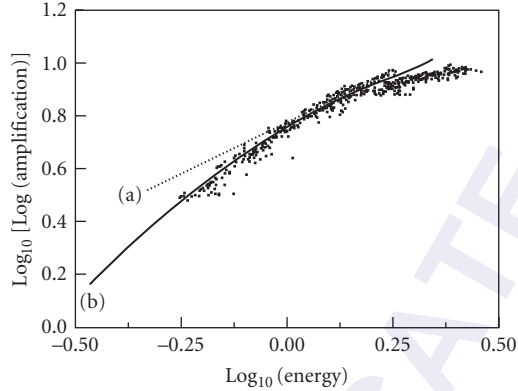


FIGURE 7 Experimental measurement of transient Stokes amplification in hydrogen gas. Solid line includes measured phase and amplitude structure of incident Stokes pulse. Dotted line was calculated for the same conditions, neglecting the phase structure of the incident Stokes pulse. The deviation of the solid curve from the measurements at the high end is due to pump depletion. (From Ref. 19.)

Phase pulling There is a tendency for the Stokes phase to become locked to the pump phase in the transient regime. This can be seen from the equations for the phases of the pump and Stokes pulses:¹⁹

$$\frac{\partial \phi_s}{\partial z} = -\kappa_2 \frac{A_L}{A_S} Q \sin \Phi \quad (46a)$$

$$\frac{\partial \phi_L}{\partial z} = -\kappa_2 \frac{n_s \omega_L}{n_L \omega_S} \frac{A_S}{A_L} Q \sin \Phi \quad (46b)$$

$$\frac{\partial \phi_Q}{\partial z} = -\kappa_1 \frac{A_L A_S}{Q} \sin \Phi \quad (46c)$$

where $\Phi = \phi_s + \phi_Q - \phi_L$.

The phase driving terms vanish for the condition $\Phi = 0$ or π , which are the conditions for optimum power transfer from the pump beam and to the pump beam, respectively. Phase pulling occurs only when $\Phi \neq 0$ or π and when the time derivatives are important. When the time derivatives are not important, ϕ_Q automatically adjusts itself so that the condition $\Phi = 0$ is satisfied. When the transient response is important, phase pulling occurs whenever ϕ_s and ϕ_L have different time dependencies. If the amplitude of the initial Stokes wave is small compared to that of the pump, the phase of the material excitation will be established at some constant value early in the pulse after the material excitation has been amplified above the noise level. At later times the phase of the Stokes wave will be driven according to Eq. (46a) so as to establish the condition $\Phi = 0$. In the early stages of the amplification, the corresponding driving term for the pump phase is much smaller, so that the Stokes phase becomes effectively locked to that of the incident pump, with a possible constant offset due to the phase of the material excitation. If the amplitude of the incident Stokes wave is comparable to that of the pump, the phases of all of the waves can evolve during the interaction, and the

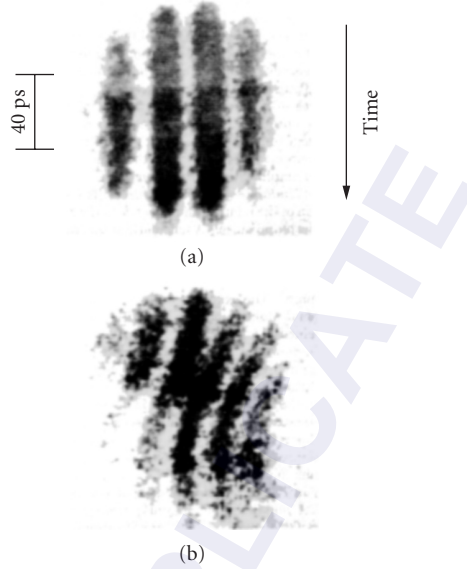


FIGURE 8 Experimental demonstration of phase pulling in transient Raman amplification. The pump beam and incident Stokes beam both carry phase modulation due to the nonlinear index in the laser. The measurement shows time-dependent interference between the incident and amplified Stokes waves. When the incident Stokes phase structure is aligned with the pump phase structure, no phase pulling is observed (a). When the incident Stokes pulse is advanced to misalign the pump and Stokes phase structure, phase pulling is observed (b). (From Ref. 19.)

transfer of power from the pump to the Stokes wave can be affected or greatly reduced. An example of phase pulling in short-pulse Raman interactions is shown in Fig. 8.

Solitons Soliton and other self-similar solutions to the transient Raman equations can also be derived. Using Eq. (37) in the form

$$\frac{\partial A_S}{\partial \xi} = i\kappa_2 A_L Q^* \quad (47a)$$

$$\frac{\partial A_L}{\partial \xi} = i\frac{\omega_L}{\omega_S} \kappa_2 A_S Q \quad (47b)$$

$$\frac{\partial Q}{\partial \tau} + \Gamma Q = i\kappa_1 A_L A_S^* \quad (47c)$$

where κ_1 and κ_2 are given in Eqs. (38) and (11), respectively, and ξ and τ are moving coordinates related to the laboratory coordinates by

$$\xi = z \quad (48a)$$

$$\tau = t - zn/c \quad (48b)$$

In the extreme transient limit, the term with Γ in Eq. (47c) can be neglected and the equations can be expressed in normalized quantities as

$$\frac{\partial A_1}{\partial \chi} = -A_2 X \quad (49a)$$

$$\frac{\partial A_2}{\partial \chi} = A_1 X^* \quad (49b)$$

$$\frac{\partial X}{\partial \tau'} = A_1 A_2^* \quad (49c)$$

where $A_1 = \sqrt{\omega_s n_L \kappa_1 / \omega_L n_s \kappa_2} A_L$, $A_2 = \sqrt{\kappa_1 / \kappa_2} A_S$, and $X = -i \sqrt{\omega_L n_s / \omega_s n_L} Q$, $\chi = \kappa_2 \xi$, $\tau' = (\omega_L n_s / \omega_s n_L) \kappa_2 \tau$.

The soliton solutions are:⁸¹

$$A_1(\chi, \tau) = K(\tau) \left(1 - \frac{a^2}{4\alpha^2}\right)^{1/2} \exp(i\alpha\chi/2) \operatorname{sech} \left\{ \left(\alpha^2 - \frac{a^2}{4}\right)^{1/2} \left[\chi - \frac{I(\tau)}{\alpha^2}\right] \right\} \quad (50a)$$

$$A_2(\chi, \tau) = K(\tau) \left(1 - \frac{a^2}{4\alpha^2}\right)^{1/2} \tanh \left\{ \left(\alpha^2 - \frac{a^2}{4}\right)^{1/2} \left[\chi - \frac{I(\tau)}{\alpha^2}\right] \right\} \quad (50b)$$

$$X(\chi, \tau) = \left(\alpha^2 - \frac{a^2}{4}\right)^{1/2} \exp(i\alpha\chi/2) \operatorname{sech} \left\{ \left(\alpha^2 - \frac{a^2}{4}\right)^{1/2} \left[\chi - \frac{I(\tau)}{\alpha^2}\right] \right\} \quad (50c)$$

where $K^2(\tau) = |A_1(\chi, \tau)|^2 + |A_2(\chi, \tau)|^2$, $I(\tau) = \int_{-\infty}^{\tau} K^2(\tau') d\tau'$ and a and α are arbitrary parameters such that $a > \alpha/2$.

These are transient solutions in which the Stokes pulse is large everywhere except for a dip centered on the position $\chi = I(\tau)/\alpha^2$. The Stokes pulse also has a phase shift at its center. When $a = 0$, the phase shift is π , and when $a \neq 0$, it is smaller. The pump intensity is small except for a narrow region centered about $\chi = I(\tau)/\alpha^2$, and the behavior of the material excitation is similar to that of the pump.

The form of these solutions is opposite to those normally encountered in Raman interactions. Once they are established, they propagate in self-similar form at a speed that is slower than that of the optical pulses, gradually "walking" backward in the optical pulse and eventually disappearing. However, the pulse shapes necessary for soliton formation do not occur in all experimental situations. In normal transient experiments, with the incident Stokes wave weak and the pump strong, the Stokes phase will lock to the pump phase, and solitons cannot develop. They have been produced with a phase shift introduced onto the Stokes pulse externally.⁸² In this situation, the soliton pulse of Eq. (50a), which is significantly shorter than T_2 , evolves from longer pulses and damping plays a central role in its formation. An experimental demonstration of this behavior has been done by Druhl et al.⁸² and their results are reproduced in Fig. 9. Conditions for soliton formation can also be encountered in growth from noise using narrowband pump pulses due to phase fluctuations in the zero point starting signal of the Stokes field. Other self-similar solutions are possible as well. Forms with damped oscillations, termed *accordion solutions*, have been described by Menyuk.⁸³

Broadband effects Transient effects also appear in the conversion of broadband radiation, which is produced, for example, in many types of pulsed lasers, when the overall pulse duration is longer than T_2 but the pump linewidth is wider than the Raman linewidth:⁸⁴⁻⁹⁹

$$t_p \gg T_2 \quad (51a)$$

$$\Delta\nu_L > \frac{1}{\pi T_2} \quad (51b)$$

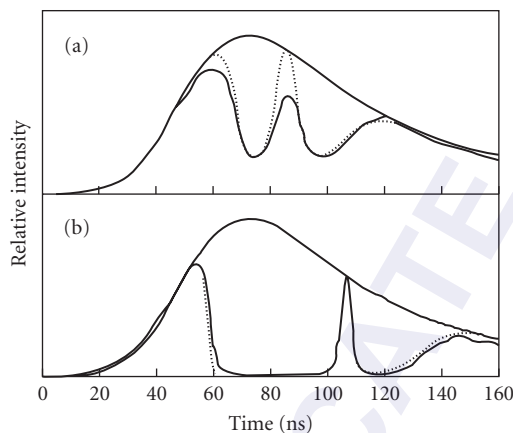


FIGURE 9 Experimental (solid curves) and theoretical (dashed curves) behavior of Raman soliton formation in hydrogen gas showing the input (upper curves) and output (lower curves) pump pulses. The soliton was initiated by introducing a phase shift on the incident Stokes pulse. The overall pulse duration was about 70 ns. The curves in (b) were obtained with higher pump power than the curves in (a). (From Ref. 82; copyright 1983 by American Physical Society.)

The equations for the Stokes and material excitation are as given in Eq. (37). The interaction can be modeled in the time domain, or in the frequency domain using a mode model for the laser radiation.

In the mode picture, the pump and Stokes radiation is modeled as being made up of a combination of randomly phased modes separated by an amount Δ , shown in Fig. 10:

$$E_S = \frac{1}{2} \sum_m A_{S,m} e^{-i(m\Delta t + \phi_{S,m})} e^{-i(\omega_S t - k_S z)} \quad (52a)$$

$$E_L = \frac{1}{2} \sum_m A_{L,m} e^{-i(m\Delta t + \phi_{L,m})} e^{-i(\omega_L t - k_L z)} \quad (52b)$$

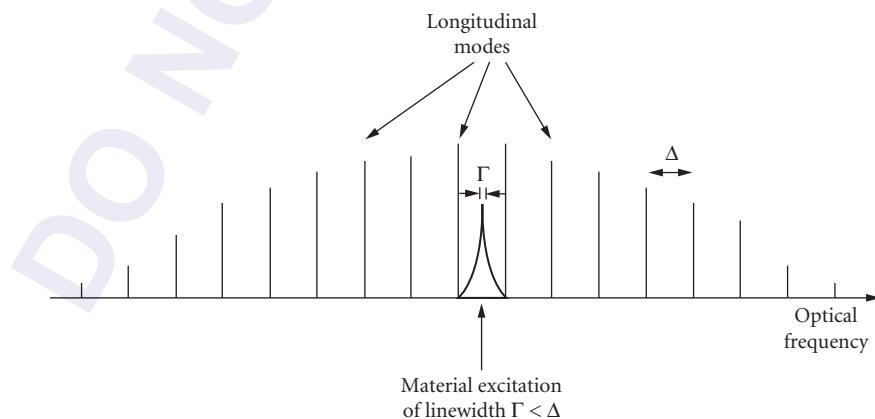


FIGURE 10 Mode structure used to model the broadband Raman scattering. The laser and Stokes modes are spaced by frequency Δ , and, for approximations of Eq. (56), the Raman linewidth is narrower than the mode spacing.

In the absence of dispersion, Eq. (37c) can be solved in the frequency domain:

$$Q_k^* = i\kappa_2 \sum_m \frac{A_{L,m+k}^* A_{S,m}}{\Gamma - im\Delta} \quad (53)$$

The average Stokes intensity is given by⁹³

$$\bar{I}_S(z) = \bar{I}_S(0) + \bar{I}_S(0) \sum_k \left[\frac{\left| \sum_m A_{L,m+k} A_{S,m}^*(0) \right|^2}{\sum_m A_{L,m} A_{L,m}^* \sum_m A_{S,m} A_{S,m}^*(0)} \right] \left\{ \exp \left[\frac{\bar{I}_L g z}{1 + (k\Delta/\Gamma)^2} \right] - 1 \right\} \quad (54)$$

where the bar symbol denotes the average intensity given by

$$\bar{I}_S = \frac{1}{2} c n_s \epsilon_0 \sum_m A_{S,m}^* A_{S,m} \quad (55)$$

If the longitudinal mode spacing is much wider than the Raman linewidth and the pump linewidth is broad enough, the material excitation Q will not be able to follow the temporal variations of the light due to the mode structure. This is equivalent to the steady-state approximation. Only the $k = 0$ term survives in Eq. (54), and the average intensity is given by

$$\bar{I}_S(z) = \bar{I}_S(0) [1 + R(e^{g\bar{I}_L z} - 1)] \quad (56)$$

where R is the field cross-correlation function given by

$$R = \frac{\left| \sum_n A_{S,n}(0) A_{P,n}^*(0) \right|^2}{\sum_n |A_{S,n}|^2 \sum_n |A_{L,n}|^2} \quad (57)$$

In this approximation the double summation in the product of the Stokes and pump waves has collapsed to a single sum. Each Stokes mode m interacts only with the corresponding pump mode m .

When the Stokes field is correlated with the pump, the correlation function is unity and the broadband Stokes light has the same exponential gain as in a narrowband interaction with the same average pump intensity. When the incident Stokes wave is not fully correlated with the pump, the component of the Stokes wave that is correlated with the pump has the largest gain, while Stokes components that are not correlated with the pump ($R = 0$) do not receive amplification.^{86–88} The Stokes wave thus becomes more correlated with the pump wave as it is amplified. This behavior is equivalent to the phase-pulling effects discussed earlier. The effective Stokes input signal to the amplifier is reduced by the factor $1/M$, where M is the number of spectral modes of the Stokes wave. However, when growth from noise is considered, the total input Stokes signal increases in proportion to the Stokes bandwidth and the effective Stokes input is one photon per mode independent of the Stokes bandwidth. The threshold for growth from noise is therefore the same for broadband and narrowband interactions.

When dispersion is taken into account,^{84,88,96} the exponential gain becomes:

$$G = G_0 - \frac{4\tau_w^2}{G_0 \tau_c^2} \quad (58)$$

where G_0 is the gain without dispersion, τ_c is the correlation time of the laser radiation, defined by $\tau_c = [1/\delta\omega_L^2]^{1/2}$ where $\delta\omega_L^2$ is the variance of the laser spectrum, and τ_w is the beam walkoff time, given by

$$\tau_w = z\Delta(1/\nu) \quad (59)$$

where $\Delta(1/\nu) = 1/\nu_{gL} - 1/\nu_{gS}$.

Solutions also exist for pump depletion when the mode spacing is large and the dispersion can be neglected:⁹⁸

$$I_S(z) = I_S(0) \frac{1 + \alpha(1 + \beta)\gamma + (1 - \beta)e^{-\beta \left[I_L(0) + \frac{\omega_L}{\omega_S} I_S(0) \right] g z}}{2 \left[\gamma + e^{-\beta \left[I_L(0) + \frac{\omega_L}{\omega_S} I_S(0) \right] g z} \right]} \quad (60)$$

where

$$\alpha = \frac{\omega_S I_L(0)}{\omega_L I_S(0)} \quad (61)$$

$$\beta = \frac{\left[\left(I_S(0) - \frac{\omega_S}{\omega_L} I_L(0) \right)^2 + 4 \frac{\omega_S}{\omega_L} I_L(0) I_S(0) R \right]^{1/2}}{\left[I_S(0) - \frac{\omega_S}{\omega_L} I_L(0) \right]} \quad (62)$$

$$\gamma = \frac{(1 + \beta) - \alpha(1 - \beta)}{\alpha(1 + \beta) - (1 - \beta)} \quad (63)$$

Broadband Raman scattering has also been analyzed within the time domain.⁹⁹ Here the starting point is Eq. (39a and b). The average of quantities is calculated as

$$\langle f(t) \rangle = \frac{1}{T} \int_t^{t+T} f(t') dt' \quad (64)$$

where the interval T is chosen to be large enough to provide a stationary average of the temporal structure in the signal. Generally speaking,

$$T \gg \frac{1}{\Delta \nu_L} \quad (65)$$

If $\Delta \nu_L \gg \Delta \nu_R$, again the material excitation cannot follow the time variations of the optical signals and Eq. (39b) can be solved as

$$\langle Q^*(t) \rangle = -\frac{i\kappa_1}{\Gamma} \langle A_S(t) A_L^*(t) \rangle \quad (66)$$

$\langle Q^*(t) \rangle$ is a slowly varying quantity even though $A_S(t)$ and $A_L(t)$ individually have rapid time variations. The equation for the average Stokes intensity is

$$\frac{\partial}{\partial z} \langle A_S(t) A_S^*(t) \rangle = i\kappa_2 \langle A_S^*(t) A_L(t) \rangle \langle Q^*(t) \rangle - A_S(t) A_L^*(t) \langle Q(t) \rangle \quad (67)$$

Making use of the fact that Q is not correlated with A_L or A_S and using Eq. (66) gives

$$\frac{\partial}{\partial z} \langle A_S(t) A_S^*(t) \rangle = i\kappa_2 \left[\langle A_S^*(t) A_L(t) \rangle \langle Q^*(t) \rangle - \langle A_S(t) A_L^*(t) \rangle \langle Q(t) \rangle \right] \quad (68)$$

$$\frac{\partial}{\partial z} \langle A_S(t) A_S^*(t) \rangle = \frac{\kappa_1 \kappa_2}{\Gamma} \left[\langle A_S^*(t) A_L(t) \rangle \langle A_S(t) A_L^*(t) \rangle \right] \quad (69)$$

The Stokes intensity is given by

$$\langle I_S(z, t) \rangle = \langle I_S(0, t) \rangle [1 + R(e^{g \langle I_L \rangle L} - 1)] \quad (70a)$$

where R is the normalized Stokes pump cross-correlation function at the input:

$$R = \frac{\langle A_S(0, t) A_L^*(0, t) \rangle \langle A_S^*(0, t) A_L(0, t) \rangle}{\langle |A_S(0, t)|^2 \rangle \langle |A_L(0, t)|^2 \rangle} \quad (70b)$$

The result of Eq. (70a) has the same form as that of Eq. (56). Akhmanov et al.⁹⁹ have discussed the statistical properties of stimulated Raman scattering with broadband radiation under a number of other conditions.

Spectral properties As described, the Stokes radiation produced in a Raman generator in the steady-state regime is expected to be a gain-narrowed version of the spontaneous Stokes emission. Druhl et al.⁸² have shown that when narrowband pump radiation is used, the linewidth of the generated Stokes radiation with single pulses varies randomly from the same width as the pump radiation to a value several times greater than the spontaneous Raman linewidth. Only in the ensemble average does the linewidth of the generated Stokes radiation coincide with the gain-narrowed spontaneous line. Individual pulses exhibit considerable spectral structure. This behavior is traceable to the stochastic nature of the damping process, by which the Raman coherence has decreased to $1/e$ of its initial value on a statistical basis.

When broadband radiation is used, the Stokes wave has a tendency to be pulled into correlation with the pump and the Stokes spectral variation matches that of the pump. Duncan et al.¹⁰⁰ have shown that the spectrum of transient spontaneous Raman scattering matches that of the pump.

Anti-Stokes Raman Scattering Anti-Stokes scattering produces a scattered wave at a shorter wavelength than the pump with frequency

$$\omega_{AS} = \omega_L + \omega_o \quad (71)$$

Anti-Stokes scattering can occur either as a two-photon transition between an upper and lower state, as illustrated in Fig. 1b, or as a resonant four-wave mixing process, as illustrated in Fig. 1c. The first interaction is directly analogous with the transitions involved in stimulated Stokes Raman scattering that have been discussed in previous sections. For a normal thermal distribution of population, the stimulated version of the anti-Stokes interaction incurs exponential loss. Anti-Stokes components are produced in the spontaneous Raman spectrum. When a population inversion is created between the upper and lower states, the anti-Stokes process has exponential gain, with properties similar to those of normal stimulated Stokes Raman scattering. This interaction, termed as *anti-Stokes Raman laser*, has been described by several authors.^{101–104}

Anti-Stokes radiation can also be produced through a four-wave mixing process, illustrated in Fig. 1c (Refs. 1, 71, 105–107, 107a, 107b). In this interaction, two pump wave photons are converted to one Stokes and one anti-Stokes photon with the relation

$$2\omega_L = \omega_S + \omega_{AS} \quad (72a)$$

The process is sensitive to the phase mismatch given by

$$\Delta k = k_{AS} - 2k_L + k_S \quad (72b)$$

Usually, four-wave mixing processes are optimized when the phase mismatch is zero. For materials with normal dispersion, this occurs when the Stokes and anti-Stokes waves propagate at angles to the pump light, as shown in Fig. 11. The angles θ_S and θ_{AS} are given in the small angle and small Δk approximation by

$$\theta_S = \sqrt{\frac{2(k_S + k_{AS} - 2k_L)}{k_S(1 + k_S/k_{AS})}} \quad (73a)$$

$$\theta_{AS} = \sqrt{\frac{2(k_S + k_{AS} - 2k_L)}{k_{AS}(1 + k_{AS}/k_S)}} \quad (73b)$$

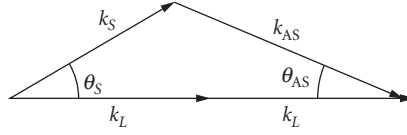


FIGURE 11 k vector diagram for coherent anti-Stokes Raman scattering showing laser and Stokes and anti-Stokes propagation directions for a medium with positive dispersion.

When the dispersion of the material is small, so that the refractive indexes at the various wavelengths can be approximated as $n_s = n_L - \delta$ and $n_{AS} = n_L + \delta$, the phase-matching angles are given by

$$\theta_s = \sqrt{\frac{2(\lambda_s - \lambda_{AS})\delta}{n_L \lambda_{AS} (1 + \lambda_{AS} / \lambda_s)}} \quad (74a)$$

$$\theta_{AS} = \sqrt{\frac{2(\lambda_s - \lambda_{AS})\delta}{n_L \lambda_s (1 + \lambda_s / \lambda_{AS})}} \quad (74b)$$

The plane-wave steady-state equations describing anti-Stokes generation with four-wave mixing are

$$\frac{dA_s^*}{dz} = K_3 |A_L|^2 A_s^* + K_2 \left(\frac{\omega_s n_{AS}}{\omega_{AS} n_s} \right) A_L^{*2} A_{AS} e^{i\Delta k z} \quad (75a)$$

$$\frac{dA_{AS}}{dz} = -K_1 |A_L|^2 A_{AS} - K_2 A_L^2 A_s^* e^{-i\Delta k z} \quad (75b)$$

where $K_1 = -iK_{AS} \chi_{AS}$; $K_2 = -iK_{AS} \sqrt{\chi_{AS} \chi_s^*}$; $K_3 = -iK_s \chi_s^*$; χ_s , and χ_{AS} are nonlinear susceptibilities for stimulated growth of the Stokes and anti-Stokes waves, respectively; $K_{S(AS)} = N \omega_{S(AS)} / n_{S(AS)} c$; N is the number density, and $\Delta k = k_s + k_{AS} - 2k_L$ is the phase mismatch.

General solutions have been discussed by Bloembergen and Shen.¹⁰⁵ These have shown that the Stokes and anti-Stokes waves grow as part of a mixed mode with amplitudes

$$A_{\text{Raman}} = \begin{pmatrix} A_s \\ A_{AS} \end{pmatrix} \quad (76)$$

One mode is primarily anti-Stokes in character and has exponential loss. The other mode is primarily Stokes in character and has exponential gain given by

$$A_{\text{Raman}} = \begin{pmatrix} A_s \\ A_{AS} \end{pmatrix} e^{gz} \quad (77)$$

where g is given by

$$g = \text{Re}\{(1/2)(K_3 - K_1)|A_L|^2 - (i/2)[\Delta k^2 + 2i\Delta k(K_3 + K_1)|A_L|^2 - (K_1 - K_3)^2|A_L|^4]^{1/2}\} \quad (78)$$

The exponential gain for the coupled mode is zero for exact phase matching, $\Delta k = 0$, and increases for nonzero Δk until it reaches its full decoupled value for $\Delta k > 2g_{ss}$, where g_{ss} is the steady-state gain

coefficient. For nonzero Δk , the maximum gain occurs at small detunings from exact resonance. The ratio of anti-Stokes to Stokes intensities is given by

$$\frac{|A_{AS}|^2}{|A_S|^2} = \left(\frac{\omega_S^2}{c^2 k_{Sz}} \right)^2 |\chi_S|^2 |A_L|^4 \Delta k^{-2} \quad (79)$$

Solutions for phase-matched conditions have been discussed by Duncan et al.¹⁰⁶ They have the form

$$A_S(z) = A_S(0) \left\{ \frac{K_1}{K_1 - K_3} - \left(\frac{K_3}{K_1 - K_3} \right) \exp[-(K_1 - K_3)|A_L|^2 z] \right\} \\ + \left(\frac{\omega_S n_{AS}}{\omega_{AS} n_S} \right) \left(\frac{K_2}{K_1 - K_3} \right) A_{AS}(0) \{1 - \exp[-(K_1 - K_3)|A_L|^2 z]\} \quad (80)$$

$$A_{AS}(z) = A_{AS}(0) \left\{ \left(\frac{K_1}{K_1 - K_3} \right) \exp[-(K_1 - K_3)|A_L|^2 z] - \left(\frac{K_3}{K_1 - K_3} \right) \right\} \\ - \left(\frac{K_2}{K_1 - K_3} \right) A_S(0) \{1 - \exp[-(K_1 - K_3)|A_L|^2 z]\} \quad (81)$$

Initially the Stokes and anti-Stokes amplitudes grow linearly in z with opposite phases. The growth slows down as z increases and the condition

$$\frac{A_{AS}(z)}{A_S(z)} = -\frac{K_2}{K_1} = \sqrt{\frac{\chi_S^*}{\chi_{AS}}} \quad (82)$$

is approached asymptotically in the limit of large z . This ratio is approximately equal to unity except when there is strong resonant enhancement of the anti-Stokes susceptibility. The maximum value of the anti-Stokes amplitude is

$$A_{AS,\max}(z) = -A_{AS}(0) \frac{K_3}{K_1 - K_3} - A_S(0) \frac{K_2}{K_1 - K_3} \quad (83a)$$

$$\approx -A_{AS}(0) \frac{\omega_S}{\omega_{AS} - \omega_S} - A_S(0) \frac{\omega_{AS}}{\omega_{AS} - \omega_S} \quad (83b)$$

where the second of these relations is approximately valid when $\chi_S^* \approx \chi_{AS}$ and $n_{AS} \approx n_S$.

If the initial Stokes amplitude is small, as for example in a Raman generator, the limiting value of the anti-Stokes amplitude will be small and the predominant anti-Stokes generation will occur at small but finite phase mismatches. For Raman generators, the anti-Stokes radiation is produced in a cone about the phase-matching angle with a dark ring at exact phase matching. Experimental limitations can make the phase-matching minimum difficult to observe, but measurements of the dark ring in the phase-matching cone have been reported, as shown in Fig. 12. If the incident Stokes intensity is comparable to the pump intensity, considerable conversion can be made to the anti-Stokes wave at exact phase matching.¹⁰⁷

In spectroscopic applications of CARS,¹⁰⁸ the usual input condition is for approximately equal pump and Stokes amplitudes with no anti-Stokes input. These experiments are usually performed under conditions of low exponential gain well below the limiting conditions of Eq. (82). Under these conditions, the anti-Stokes generation is maximized at the phase-matching conditions.

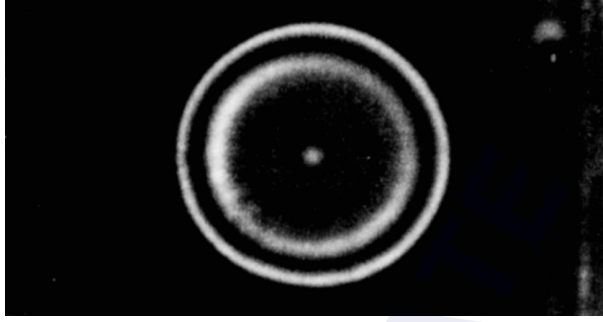


FIGURE 12 Photograph of the far field of the anti-Stokes emission pattern in hydrogen gas at a pressure of 14 atm. The anti-Stokes radiation is emitted in a cone about the phase-matching angle. The dark ring in the center of the cone is due to parametric gain suppression. (From Ref. 106.)

Growth from Noise The most common configuration for Stokes Raman interactions is a Raman generator in which only a pump signal is provided at the input, as shown in Fig. 1a. The Stokes wave is generated in the interaction. The Stokes generation process can be viewed as one in which the effective Stokes noise at the beginning of the cell is amplified in the stimulated Raman interaction as described in the previous sections.

Classically, the Stokes noise is considered as arising from the spontaneous Raman scattering that is produced at the beginning of the cell. If we consider the effective spontaneous Stokes radiation that serves as a source for amplification to be that generated in the first e -folding length of the Raman generator, the Stokes intensity at the output is

$$I_s = \frac{\partial \sigma}{\partial \Omega} d\Omega N A (e^{N\sigma I_L} - 1) \quad (84)$$

where N is the number density of the medium, A is the cross-sectional area, and $d\Omega$, is the solid angle of the gain column.

The growth from noise has been modeled more rigorously in terms of quantum fluctuations of the Stokes field amplitude and material excitation.^{97,100,109-112} In this treatment, the Stokes and material oscillators are described by quantum mechanical creation and annihilation operators while the pump field is treated classically. The equations for the Stokes and material oscillators are:⁹⁷

$$\frac{\partial}{\partial z} \hat{A}_s^{(-)}(z, \tau) = i\kappa_2 A_L(t) \hat{Q}^{(+)}(z, \tau) \quad (85a)$$

$$\frac{\partial}{\partial \tau} \hat{Q}^{(+)}(\tau) + \Gamma \hat{Q}^{(+)}(\tau) = -i\kappa_1 A_L^*(\tau) \hat{A}_s^{(-)}(z, \tau) + \hat{F}^{(+)}(z, \tau) \quad (85b)$$

where the symbol $\hat{}$ indicates a quantum mechanical operator, the symbols $(-)$ and $(+)$ indicate creation and annihilation operators, respectively, and \hat{F} is a Langevin operator that ensures the correct longtime behavior of Q .

The initial fluctuations of the Stokes and material oscillators satisfy the conditions:

$$\langle \hat{A}_s^{(+)}(0, t) \hat{A}_s^{(-)}(0, t') \rangle = \frac{2\hbar\omega_s}{cn_s \epsilon_0 a} \delta(t-t') \quad (86a)$$

$$\langle \hat{A}_s^{(-)}(0, t) \hat{A}_s^{(+)}(0, t') \rangle = 0 \quad (86b)$$

$$\langle \hat{Q}^{(+)}(z, 0) \hat{Q}^{(-)}(z', 0) \rangle = \frac{1}{\rho} \delta(z - z') \quad (86c)$$

$$\langle \hat{Q}^{(-)}(z, 0) \hat{Q}^{(+)}(z', 0) \rangle = 0 \quad (86d)$$

$$\langle \hat{F}^{(+)}(z, t) \hat{F}^{(-)}(z', t') \rangle = \frac{2\Gamma}{\rho} \delta(z - z') \delta(t - t') \quad (86e)$$

$$\langle \hat{F}^{(-)}(z, t) \hat{F}^{(+)}(z', t') \rangle = 0 \quad (86f)$$

where a is the cross-sectional area of the beam. The average Stokes intensity is given by the expectation of the normally ordered number operator:

$$I_s(z, \tau) = \frac{1}{2} c n_s \epsilon_o \langle A_s^{(-)}(z, \tau) A_s^{(+)}(z, \tau) \rangle \quad (87)$$

The formal solution is⁹⁷

$$\begin{aligned} \hat{A}_s^{(-)}(z, \tau) = & \hat{A}_s^{(-)}(0, \tau) \\ & + (\kappa_1 \kappa_2 z)^{1/2} A_L(\tau) \int_{-\infty}^{\tau} d\tau' \hat{A}_s^{(-)}(0, \tau') A_L^*(\tau') e^{-\Gamma(\tau-\tau')} \frac{I_1(\{4\kappa_1 \kappa_2 z [p(\tau) - p(\tau')] \}^{1/2})}{[p(\tau) - p(\tau')]^{1/2}} \\ & - i \kappa_2 A_L(\tau) e^{-\Gamma\tau} \int_0^z dz' \hat{Q}^+(z', 0) I_0(\{4\kappa_1 \kappa_2 (z - z') p(\tau)\}^{1/2}) \\ & - i \kappa_2 A_L(\tau) \int_{-\infty}^{\tau} d\tau' \int_0^z dz' \hat{F}(z', \tau') e^{-\Gamma(\tau-\tau')} I_0(\{4\kappa_1 \kappa_2 (z - z') [p(\tau) - p(\tau')] \}^{1/2}) \end{aligned} \quad (88)$$

The Stokes intensity obtained from use of Eq. (87) is

$$\begin{aligned} I_s(z, \tau) = & \frac{1}{2} c n_s \epsilon_o |\kappa_2 A_L(\tau)|^2 z \{ e^{-2\Gamma\tau} (I_0^2[4\kappa_1 \kappa_2 z p(\tau)]^{1/2}) - I_1^2([4\kappa_1 \kappa_2 z p(\tau)]^{1/2}) \\ & + 2\Gamma \int_{-\infty}^{\tau} e^{-2\Gamma(\tau-\tau')} (I_0^2(4\kappa_1 \kappa_2 z [p(\tau) - p(\tau')]^{1/2}) - I_1^2(\{4\kappa_1 \kappa_2 z [p(\tau) - p(\tau')] \}^{1/2})) d\tau' \} \end{aligned} \quad (89)$$

In this formulation, only the third and fourth terms of Eq. (88) survive because the expectation values on the right side are taken over the initial state, which contains no quanta in either the Stokes or molecular fields. The first and second terms involve a Stokes annihilation operator acting on the Stokes ground state and return zero. The third term returns a nonzero result because it involves a creation operator acting on the molecular ground state, as discussed by Raymer.⁹⁷ In this treatment the Stokes light is generated entirely from fluctuations in the material oscillators, while the material excitation is generated from zero-point fluctuations in the material oscillators.

In the extreme transient regime, this result reduces to

$$I_s(z, \tau) = \frac{1}{2} \Gamma g_{ss} I_L(\tau) z \{ I_0^2((2g_{ss} z I_L \Gamma \tau)^{1/2}) - I_1^2((2g_{ss} z I_L \Gamma \tau)^{1/2}) \} \quad (90)$$

Comparison of this result with that of Eq. (43) shows a different functional dependence on the modified Bessel functions, which reflects the effects of buildup of the signal from the distributed noise source.

In the steady state, Eq. (89) reduces to

$$I_s(z, \tau = \infty) = \frac{1}{2} \Gamma g_{ss} I_L z [I_0(g_{ss} I_L z / 2) - I_1(g_{ss} I_L z / 2)] e^{g_{ss} I_L z / 2} \quad (91)$$

An alternative analysis for transient scattering has been presented using antinormal ordering of the creation and annihilation operators for the intensity.¹⁰⁹ In this formalism, the zero-point term must be subtracted explicitly. The intensity is given by

$$I_S(z, \tau) = \frac{1}{2} c n_s \varepsilon_o \left\{ \langle \hat{A}_S^{(+)}(z, \tau) \hat{A}_S^{(-)}(z, \tau) \rangle - \langle \hat{A}_S^{(-)}(0, \tau) \hat{A}_S^{(+)}(0, \tau) \rangle \right\} \quad (92)$$

The Stokes intensity is given by

$$\begin{aligned} I_S(z, \tau) = & \frac{1}{2} c n_s \varepsilon_o \left\{ \kappa_2^2 |A_L(\tau)|^2 \iint dz' dz'' I_o(\sqrt{4\kappa_1 \kappa_2 (z-z') p(\tau)}) I_o(\sqrt{4\kappa_1 \kappa_2 (z-z'') p(\tau)}) \right. \\ & \times \langle \hat{Q}^{(-)}(z', 0) \hat{Q}^{(+)}(z'', 0) \rangle + \kappa_1 \kappa_2 z |A_L(\tau)|^2 \iint I_1(\sqrt{4\kappa_1 \kappa_2 z [p(\tau) - p(\tau')]})) \\ & \times I_1(\sqrt{4\kappa_1 \kappa_2 z [p(\tau) - p(\tau'')]})) / \sqrt{[p(\tau) - p(\tau')] [p(\tau) - p(\tau'')]} \\ & \times A_L(\tau') A_L^*(\tau'') \langle \hat{A}_S^{(+)}(0, \tau') \hat{A}_S^{(-)}(0, \tau'') \rangle d\tau' d\tau'' + \sqrt{\kappa_1 \kappa_2 z} \left\{ A_L^*(\tau) - \int_{-\infty}^{\tau} A_L(\tau') \langle \hat{A}_S^{(+)}(0, \tau') \hat{A}_S^{(-)}(0, \tau) \rangle \right. \\ & \times I_1(\sqrt{4\kappa_1 \kappa_2 z [p(\tau) - p(\tau')]})) / \sqrt{p(\tau) - p(\tau')} d\tau' + A_L(\tau) \int_{-\infty}^{\tau} A_L^*(\tau') \\ & \left. \times \langle \hat{A}_S^{(+)}(0, \tau) \hat{A}_S^{(-)}(0, \tau') \rangle I_1(\sqrt{4\kappa_1 \kappa_2 z [p(\tau) - p(\tau')]})) / \sqrt{p(\tau) - p(\tau')} d\tau' \right\} \end{aligned} \quad (93)$$

Here the first term in Q is 0 because it represents an annihilation operator operating on the ground state. Further analysis shows that this result is identical to the one in Eq. (90). The second term gives the transient stimulated Raman signal, and the last two terms in brackets describe spontaneous Raman scattering. In this formalism, the Stokes wave is started by its own zero-point motion and does not involve the zero-point motion of the molecular oscillators. The zero-point motion of the molecular oscillators is responsible for the initiation of the molecular excitation. Further analysis has shown that the Stokes signal can be viewed as arising from quantum fluctuations of the Stokes radiation,⁹⁷ the material oscillator,¹⁰⁹ or a combination of both.¹⁰⁰

The effective Stokes noise amplitude corresponds to one Stokes photon at the input of the generator. This result is expected for this model, which assumes plane wave propagation, effectively assuming a single mode in the amplifier. In a more general case, the effective Stokes noise level will be one Stokes photon for each temporal and spatial mode of the amplifier. The number of spatial modes is given by the square of the effective Fresnel number of the amplifier:¹¹³

$$N_{\text{spatial modes}} = F^2 \quad (94)$$

where

$$F = \frac{A}{\lambda L} \quad (95)$$

where A is the effective area of the gain region and L is the interaction length. Because of spatial gain narrowing, the effective area of the generator can be significantly less than the nominal diameter of the pump beam and can depend on the gain level. The number of temporal modes depends on the relation of the pump pulse duration to the dephasing time T_2 . A single temporal mode will be present when $t_p \ll T_2$. For longer pulses, the number of temporal modes has been modeled as¹¹⁰

$$N_{\text{temporal modes}} = 1, \quad t_p < t_{ss} \quad (96a)$$

$$N_{\text{temporal modes}} = \frac{t_p}{G_{ss} T_2}, \quad t_p > t_{ss} \quad (96b)$$

or alternatively as¹¹⁴

$$N_{\text{temporal modes}} = \sqrt{\frac{\ln 2/G_{ss}}{1 - \ln 2/G_{ss}}} \frac{t_p}{\pi T_2} \quad (97)$$

The first of these derives the number of temporal modes from the steady-state time of Eq. (19) and the second from the gain-narrowing formula of Eq. (29).

Raman threshold Generation of first Stokes radiation from noise in a single-pass generator passes smoothly from exponential amplification of noise to depletion of the pump radiation without a true threshold. Thresholdlike behavior has been reported in some situations but has been due to multimode structure of the radiation or secondary reflections. A Raman threshold is, however, commonly associated with a single-pass gain. This is done by assigning the threshold to a pump value (intensity or energy) at which the Stokes signal from a Raman generator is an arbitrary fraction of the incident pump (typically of the order of 1 percent). At this level, pump saturation is generally not important, but for higher pump powers the process rapidly transitions to saturation. Thus, the concept of Raman threshold in Raman generators is reasonably practical, if not technically precise. The gain that is required to reach threshold depends on the number of noise modes present in the generator. For typical geometries of long, narrow interaction lengths, the Raman gain at threshold is of the order of e^{23} to e^{40} . Raman amplifiers are typically operated at gain levels below threshold. Stable amplifier gains of the order of e^{19} are achievable.

Quantum fluctuations Macroscopic manifestations of the stochastic nature of the Raman initiating fluctuations have been reported in spatial fluctuations of the output Stokes intensity profile, in the pointing of the Stokes beam, and in the spectral and temporal structure of the Stokes signal generated by narrowband radiation.^{110,111,115–120} The stochastic nature of the starting signal is manifest in the pulse energy statistics of Raman generators operated below threshold. In this regime, the statistical distribution of the Stokes pulse energy is of the form $p(W_s) = \exp\{-W_s/\langle W_s \rangle\}$, where W_s is the energy of a Stokes pulse and $\langle W_s \rangle$ is the average energy over an ensemble. When the Raman generator is operated below threshold, energy of the output pulses fluctuates in accordance with this distribution. As the generator approaches pump depletion, the Stokes pulse energy distribution approaches one that is peaked about the average value. Statistical distributions of pulse radiation in short pulse experiments show exponential behavior characteristic of stochastic input for gains below threshold, and a gradual evolution to coherent behavior as saturation is approached.

Competition of the quantum noise with real Stokes signals in Raman amplifiers at the quantum level has been reported by Duncan et al.¹¹³ Their results show experimentally that the effective initiating signal is consistent with a noise level of one photon per mode of the amplifier. When the incident Stokes signal exceeds the noise level of one photon per mode by a sufficient amount, the fluctuations are effectively suppressed in both the spatial profile and the pulse energy statistics. An example of the evolution of the amplified Stokes signal from one dominated by quantum noise to one dominated by the coherent Stokes input signal in the image of bars in a resolution chart is shown in Fig. 13.

Multiple Stokes Generation Once a significant signal is produced in the first Stokes wave, it can serve as a pump wave for a second Raman process, generating a second Stokes wave at

$$\omega_{2S} = \omega_S - \omega_o = \omega_L - 2\omega_o \quad (98)$$

This usually occurs at a pump power such that significant Stokes conversion occurs within the first half of the Raman cell, allowing generation of the second Stokes wave in the last part of the cell. Once the first Stokes wave is generated, the second Stokes radiation can also arise from a four-wave mixing interaction of the form

$$\omega_{2S} = 2\omega_S - \omega_L \quad (99a)$$

which has properties similar to the anti-Stokes four-wave mixing interaction.

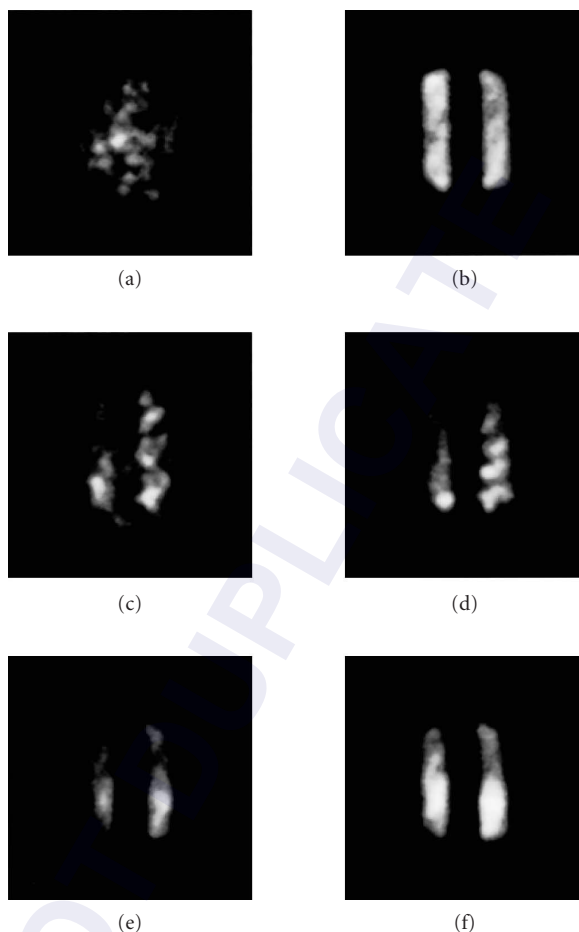


FIGURE 13 Images of bars in a resolution chart from a Stokes amplifier with amplification of 1.4×10^4 for different levels of incident Stokes energy, showing effects of competition between incident Stokes energy and quantum noise, (a) Incident Stokes wave blocked. (b) Incident Stokes wave. (c) 210 seed Stokes photons, camera sensitivity 1. (d) 800 seed Stokes photons, camera sensitivity 0.3. (e) 3.2×10^4 seed Stokes photons, camera sensitivity 7.4×10^{-3} . (f) 1.5×10^{10} seed Stokes photons, camera sensitivity 3.2×10^{-8} . (From Ref. 121.)

The four-wave mixing interaction is a coherent one and will produce second Stokes radiation with coherent statistics when the first Stokes radiation has saturated the pump. The second Stokes radiation generated from stimulated scattering will exhibit the stochastic behavior characteristic of growth from quantum noise. The relative importance of the two sources of second Stokes radiation is affected by the phase mismatch for the four-wave mixing interaction and depends on the density of the material.

Still higher pump powers can result in conversion to third or higher Stokes orders. Each of the orders involves a frequency shift due to the same material transition, rather than higher excitation of the material system. In most materials, conditions for multiple Stokes generation are sufficient

to produce significant amounts of anti-Stokes energy through the four-wave mixing interaction. Multiple-order anti-Stokes energy can also be generated using various combinations of Stokes and anti-Stokes orders of the form

$$\omega_{nAS} = \omega_{mS} - \omega_{(m+1)S} + \omega_{(n-1)AS} \quad (99b)$$

where n and m are orders of Stokes and anti-Stokes radiation and $\omega_{0S} = \omega_{0AS} = \omega_L$. Multiple-order Stokes radiation can also be produced through four-wave mixing involving similar terms.

An ideal progression of Raman scattering through multiple Stokes modes is shown in Fig. 14a. Such a progression is seldom seen in practice because of the onset of anti-Stokes and four-wave mixing interactions. Higher-order Stokes energy can be suppressed through choice of resonant structure in the material or through use of high pressures that suppress four-wave mixing through disruption of phase matching. An example of second Stokes generation in hydrogen is shown in Fig. 14b, in which initiation through four-wave mixing occurs at low powers and initiation through a stimulated process, evidenced by the wide scatter of points, occurs at higher powers.¹²²

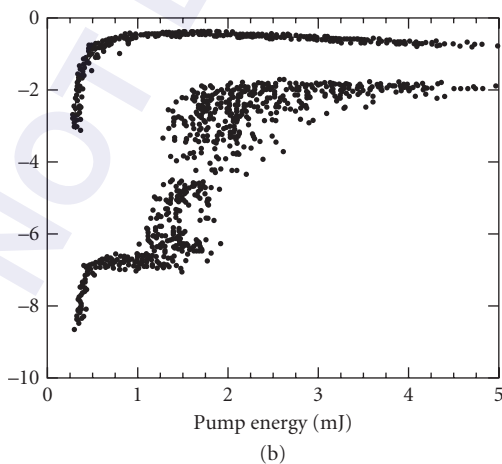
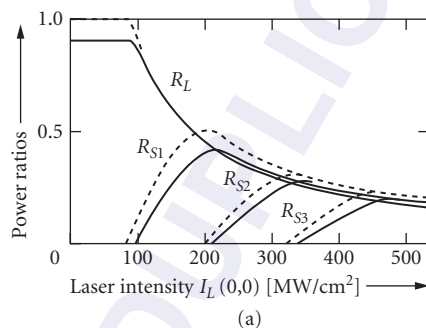


FIGURE 14 (a) Theoretical calculation of multiple Stokes generation for gaussian pulses. (From Ref. 2.) (b) Second Stokes generation (lower curve) in hydrogen at 1600 psi showing growth from four-wave mixing at low pump energies and the transition to stimulated emission from quantum fluctuations at higher pump energies. The first Stokes wave is shown in the upper curve. (From Ref. 122.)

Focused Beams The effects of focusing are described by the spatial derivative in Eq. (5a and b). When the pump intensity varies with z , the gain must be integrated over the interaction length. For steady-state interactions, the Stokes intensity takes the form

$$I_s(r, z) = I_s(r, 0) \exp \left[g \int_0^L I_L dz \right] \quad (100)$$

The profile of a gaussian beam is described by a $1/e$ field radius w , given by

$$w(z) = w_o \sqrt{1 + \xi^2} \quad (101)$$

where w_o is the radius of the beam waist, $\xi = 2z/b$, and $b = 2\pi w_o^2/\lambda$ is the confocal parameter. When $b \gg L$, the pump beam is collimated over the interaction length and the primary effect of the gaussian profile is to produce gain narrowing, effectively confining the Stokes intensity near the beam axis. When the beam is tightly focused, so that $b \ll L$, the integrated gain is independent of the interaction length and depends only on the total pump power:

$$G = g \int I_L dz = g4P/b\lambda \quad (102)$$

where P is the total pump power. Amplification or generation with focused beams can result in changes of the Stokes beam divergence and displacements of the apparent source point for divergence of the Stokes beam.

Backward Raman Scattering Stimulated Raman scattering also occurs in the backward direction, for which the Stokes and pump waves travel in opposite directions.^{2,69} This type of interaction is much more dependent on the geometry and the laser linewidth than the forward interaction. Backward interaction can involve growth from noise (Fig. 2b) or amplification (Fig. 2d). Generally, it occurs for sufficiently narrow line widths that the coherence length of the radiation is longer than the interaction length. One of the characteristic differences of the backward interaction is that the growing Stokes wave continually interacts with undepleted pump as it propagates. Thus, for conditions in which the Stokes wave can grow to saturate the Raman interaction, the Stokes wave continually experiences the full pump intensity rather than a continually decreasing intensity as in the forward direction. One consequence of this property is that the intensity of the backward Stokes pulse can grow to be much higher than the initial pump intensity, while the pulse duration decreases, producing pulse compression, one of the common applications for backward Raman scattering.⁶⁹ For this application, the optimal duration for the pump pulse is twice the length of the Raman cell. The Stokes wave grows to depletion level in the first half of the pump pulse and depletes the pump in the second half. Factors affecting pulse compression have been discussed in Ref. 69, where designs of systems to give a factor-of-10 shortening were described.

Backward scattering also results in Stokes pulses of higher spatial quality than are obtained in the forward direction. If the backward Stokes wave depletes the pump pulse too much before it travels an entire cell length, the initiating signal for Stokes radiation in subsequent parts of the pulse can be suppressed and oscillations can result.¹²³

Polarization Dependence The Raman gain for various polarizations depends on the symmetry of the Raman transition and is governed by the depolarization ratio. For many materials and transitions, the maximum gain occurs for pump and Stokes polarizations that are linear and parallel. For other types of transitions—for example, rotational transitions in diatomic molecules—circular polarization is preferred with the pump and Stokes waves polarized in the opposite sense. The relative gains for various polarization combinations are shown in Table 8.¹²⁴

The earliest and most common application of stimulated Raman scattering is the production of coherent sources at wavelengths different from those of the pump laser. Prior to the introduction of tunable lasers, this was one of the few methods available for obtaining coherent radiation at any but a small number of wavelengths at which fixed-frequency lasers existed. This application remains an important one today for extending the versatility of tunable lasers, and for generation of radiation

TABLE 8 Polarization Dependence of Relative Gain for SRS

	Laser Polarization	Stokes Polarization	Relative Gain
Rotational	Linear	Linear, parallel	1
scattering,	Linear	Linear, perpendicular	0.75
linear	Circular	Circular, same sense	0.25
molecules	Circular	Circular, opposite sense	1.5

at wavelengths required for specific applications such as the eye-safe region around 1.5 Mm. Other applications of Raman scattering include Raman beam cleanup, pulse compression, time-gated imaging, and coherent spectroscopy. Two of these are described in the following text.

Coherent Spectroscopy Several Raman interactions are used for coherent spectroscopy.^{108,125} These interactions have the advantage of producing a stronger signal than incoherent spectroscopy under some conditions. These advantages occur primarily when prominent Raman modes are studied in pure materials or the dominant constituent of a mixture. A summary of interactions used in coherent spectroscopy along with names given to them is given in Table 9. The most common of these is coherent anti-Stokes Raman scattering (CARS). A typical CARS spectroscopic interaction is shown in Fig. 11. Incident light at both the pump and Stokes wavelengths is supplied. The anti-Stokes signal is measured as the wavelength of the Stokes light is varied. The spectral structure is recorded as the frequency difference $\omega_L - \omega_S$ is tuned through Raman resonances of the material.

The anti-Stokes intensity is given in a steady-state plane-wave approximation as

$$I_{AS} = |\chi_R + \chi_{NR}|^2 I_L^2 I_S \quad (103)$$

where χ_R is the Raman susceptibility that carries the resonance and χ_{NR} is the nonresonant nonlinear susceptibility due to electronic or other transitions in the medium. The Raman susceptibility exhibits resonant behavior as $\omega_L - \omega_S$ is tuned through the Raman resonance ω_o , while the nonresonant susceptibility is generally constant through the resonance region. The wavelength variation of the anti-Stokes intensity reflects the influence of the interference between the resonant and nonresonant susceptibilities. In particular, the anti-Stokes intensity goes through zero when $\chi_R = -\chi_{NR}$, which always occurs on the high-frequency side of the resonance. Depending on the situation, the anti-Stokes intensity can exhibit a resonant-like peak if χ_{NR} is small compared to χ_R at the resonance or a dispersive-like behavior if χ_{NR} is larger than χ_R . Intermediate behavior between these two extremes is also possible depending on the relative magnitudes of χ_R and χ_{NR} . Dispersive-like behavior can occur for relatively weak Raman transitions, or for Raman transitions in materials that are minor constituents in mixtures. An example of CARS spectral behavior for various conditions is shown in Fig. 15.

Time-Gated Imaging Raman amplification has also been applied to time-gated imaging to suppress background light in highly scattering materials, for example biological tissue or materials such as ceramics.^{127–129} Time gating with Raman scattering can be done with either short pulses, in which case the gate time is comparable to the pulse duration, or with broadband radiation, in which case the time gate is determined by the coherence time of the pulse. In pulsed gating, only that part of the signal that overlaps with the pump pulse in time is amplified, and light that is delayed through multiple scattering is suppressed. In coherence gating, only the Stokes components that are coherent with the pump are amplified, and components that

TABLE 9 Raman Interactions Used for Coherent Spectroscopy

Interaction	Frequency Relations	Measured Quantity
Coherent anti-Stokes Raman scattering (CARS)	$\omega_{AS} = 2\omega_L - \omega_S$	Anti-Stokes power
Coherent Stokes Raman spectroscopy (CSRS)	$\omega_S = \omega_L - \omega_o$	Amplified Stokes power
Raman-induced Kerr effect spectroscopy (RIKES)	$\omega_S = \omega_L - \omega_o$	Stokes power in orthogonal polarization component

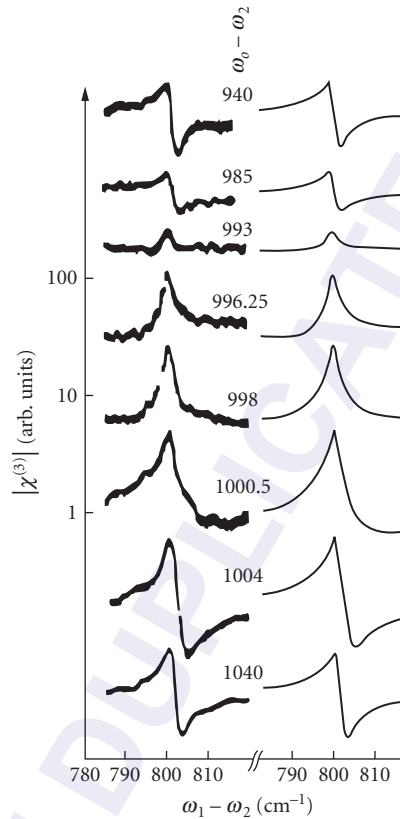


FIGURE 15 Signal from coherent anti-Stokes Raman scattering showing variation of line shape due to interference of Raman and nonresonant susceptibilities. (From Ref. 126; copyright 1976 by American Physical Society.)

are decorrelated because of multiple-path scattering are suppressed. Gate contrasts of the order of 10^9 are possible with pulse gating, while contrasts of the order of 10^6 can be obtained with correlation gating. An example of the use of pulsed gating for imaging in a scattering solution is shown in Fig. 16.

15.3 STIMULATED BRILLOUIN SCATTERING

Brillouin scattering involves low-frequency propagating waves—for example, acoustic waves in solids, liquids, and gases and ion-acoustic waves in plasmas. Again, scattering can be to a longer or shorter wavelength than the incident radiation, with the long-wavelength scattered wave being termed the *Stokes wave* and the short-wavelength scattered wave termed the *anti-Stokes wave*. The difference between the incident and scattered frequencies is again termed the *Stokes shift* or *anti-Stokes shift* as appropriate. For Brillouin scattering, the energies of the modes are much lower than for Raman scattering, and anti-Stokes radiation is much more common. Common Brillouin shifts are typically on the order of 0.1 to 100 GHz, and depend on the excitation wavelength and interaction geometry as well as on material properties. Brillouin scattering is used most commonly for phase conjugation and pulse compression. It is also prominent as a limiting process for intensity in fiber-optic systems.

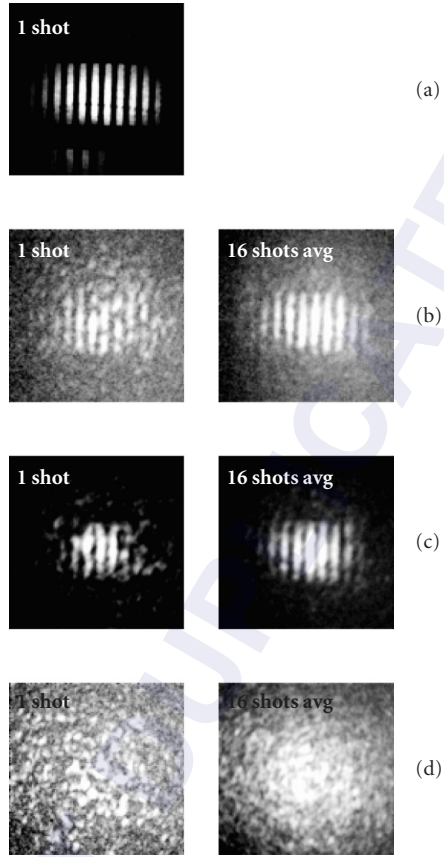


FIGURE 16 Use of Raman amplification with short pulses for time gating to photograph objects through scattering materials. (a) Image of bar chart, no scatterer. The relative timing of the pump and Stokes pulses were (b) -16 ps, (c) 0 ps, and (d) $+24$ ps. (From Ref. 128.)

The equations for Stokes Brillouin scattering for forward ($+v_g$) and backward ($-v_g$) waves are^{2,6,7}

$$\frac{\partial A_S}{\partial z} \pm \frac{1}{v_g} \frac{\partial A_S}{\partial t} = i \frac{\omega_s}{4n_s c} \gamma \frac{\Delta \rho^*}{\rho_o} A_L \quad (104a)$$

$$\begin{aligned} -2i\Omega \frac{\partial \Delta \rho^*}{\partial t} + (v^2 q^2 - \Omega^2 + iq^2 \Gamma' \Omega) \Delta \rho^* + 2iqv^2 \frac{\partial \Delta \rho^*}{\partial z} - \frac{v^2 s q^2 \beta_p \rho_o}{\gamma} \Delta T^* \\ + 2iq \frac{v_s^2 \beta_p \rho_o}{\gamma} \frac{\partial \Delta T^*}{\partial z} = q^2 \gamma A_L^* A_S \end{aligned} \quad (104b)$$

$$\rho_o C_v \frac{\partial \Delta T^*}{\partial t} - i\Omega \rho_o C_v \Delta T^* - \frac{C_v (\gamma - 1)}{\beta_p} \frac{\partial \Delta \rho}{\partial t} + i\Omega \frac{C_v (\gamma - 1)}{\beta_p} \Delta \rho + \kappa q^2 (\Delta T) = cn \epsilon_o \alpha A_L^* A_S \quad (104c)$$

where v_s is the sound velocity and

$$\Gamma' = \frac{1}{\rho} \left[\frac{4}{3} \eta_s + \eta_b + \frac{\kappa}{C_p} \left(\frac{C_p}{C_v} - 1 \right) \right] \quad (105)$$

and η_s is the coefficient of shear viscosity, η_b is the coefficient of bulk viscosity, κ is the thermal conductivity, C_p is the specific heat at constant pressure, and C_v is the specific heat at constant volume.

The fields are described by the equations

$$E_L(z, t) = \frac{1}{2} [A_L(z, t) e^{-i(\omega_L t - k_L z)} + \text{c.c.}] \quad (106a)$$

$$E_S(\xi, t) = \frac{1}{2} [A_S(z, t) e^{-i(\omega_S t - k_S \xi)} + \text{c.c.}] \quad (106b)$$

$$\rho(\zeta, t) = \rho_o + \frac{1}{2} [\Delta \rho e^{-i(\Omega t - q \zeta)} + \text{c.c.}] \quad (106c)$$

where A_L , A_S , and ρ are the amplitudes of the laser, scattered optical wave, and sound wave density, respectively, and ω_L , ω_S , Ω , and k_L , k_S , and q are the frequencies and k vectors of the various waves. The scattered optical wave propagates along the direction ξ and the sound wave propagates along direction ζ , neither of which is required to coincide with z . The frequencies and k vectors obey the following relations:

$$\omega_L - \omega_S = \Omega \quad (107a)$$

$$\vec{k}_L - \vec{k}_S = \vec{q} \quad (107b)$$

$$q = \frac{2\pi}{\Lambda_s} \quad (107c)$$

where Λ_s is the sound wavelength. The k vectors of the various waves are arranged according to the diagram in Fig. 17.

Because the sound frequency is much less than the optical frequency, $|k_L| \approx |k_S|$ and

$$q \approx 2k_L \sin(\theta/2) \quad (108)$$

The Brillouin frequency shift can then be written as

$$\Omega_B = 2n_L \omega_L \frac{v_s}{c} \sin(\theta/2) \quad (109)$$

Unlike the Raman frequency shift, the Brillouin frequency shift depends on the laser frequency and the interaction geometry. It has its maximum value for backward scattering ($\theta = 180^\circ$).

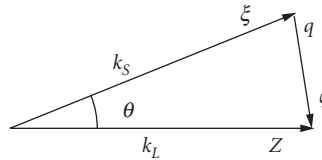


FIGURE 17 k vector diagram for stimulated Brillouin scattering showing k vectors for laser k_L Stokes k_S and sound wave q .

The Brillouin interaction involves both electrostrictive and thermal effects. Equation (104c) was obtained from the equation for the low-frequency acoustic and thermal response of liquids and gases to optical radiation:^{2,6,7}

$$-\frac{\partial^2 \Delta \rho}{\partial t^2} + \frac{v_s^2}{\gamma} \nabla^2 (\Delta \rho) + \frac{2\eta_s + \eta_b}{\rho_o} \nabla^2 \frac{\partial \Delta \rho}{\partial t} + \frac{v_s^2 \beta_p \rho_o}{\gamma} \nabla^2 (\Delta T) = \frac{1}{2} c n \epsilon_o \gamma_e \nabla^2 \langle E_{\text{tot}} \rangle^2 \quad (110a)$$

$$\rho_o C_v \frac{\partial \Delta T}{\partial t} - \frac{C_v (\gamma - 1)}{\beta_p} \frac{\partial \Delta \rho}{\partial t} - \kappa \nabla^2 (\Delta T) = c n \epsilon_o \alpha \langle E_{\text{tot}} \rangle^2 \quad (110b)$$

where $\gamma = C_p/C_v$ is the ratio of specific heats at constant pressure and volume, η_s is the coefficient of shear viscosity, η_b is the coefficient of bulk viscosity, β_p is the thermal expansion coefficient at constant pressure, $\gamma_e = \rho(\partial \epsilon / \partial \rho)$ is the electrostrictive constant, α is the absorption coefficient, κ is the thermal conductivity, $E_{\text{tot}} = E_L + E_S$ is the total optical field, and ΔT is the change in temperature. The following form for the temperature

$$T = T_o + \frac{1}{2} [\Delta T e^{-i(\Omega t - qz)} + \text{c.c.}] \quad (111)$$

can be used to reduce Eq. (110a and b) to first order, giving Eq. (104b and c).

The Brillouin equations can be solved under various approximations. The propagating terms in $\partial \Delta \rho / \partial z$ and $\partial \Delta T / \partial z$ are usually neglected because the sound waves are strongly damped. The steady-state solution for the Stokes intensity is

$$I_S(L) = I_S(0) e^{gL} \quad (112)$$

where the gain coefficient g is given by $g = g^e + g^a$, where g^e is the electrostrictive gain coefficient given by

$$g^e = \frac{\omega_s \gamma_e^2 \Omega_B}{n_s n_L c^2 \epsilon_o \rho_o \Gamma_B v_s^2} \frac{1}{1 + \left(\frac{2\delta\Omega}{\Gamma_B} \right)^2} \quad (113a)$$

and g^a is the absorptive Brillouin gain coefficient given by

$$g^a = \frac{\omega_s \gamma_e \gamma_a \Omega_B}{2 n_s n_L c^2 \epsilon_o \rho_o \Gamma_B v_s^2} \frac{4\delta\Omega / \Gamma_B}{1 + \left(\frac{2\delta\Omega}{\Gamma_B} \right)^2} \quad (113b)$$

where $\gamma_a = 2\alpha n v^2 \beta_p / C_p \Omega_B$,

$$\Omega_B = v_s q \quad (114)$$

is the Brillouin frequency shift, $\delta\Omega = \Omega - \Omega_B$, and

$$\Gamma_B = q^2 \left[\frac{1}{\rho} (2\eta_s + \eta_b) + \frac{\kappa}{C_p} \left(\frac{C_p}{C_v} - 1 \right) \right] \quad (115a)$$

$$= q^2 \Gamma'' \quad (115b)$$

is the Brillouin damping constant.

The gain due to the electrostrictive interaction is peaked about the Brillouin frequency shift $\Omega = \Omega_B$ with a linewidth (FWHM) of $\Delta\nu_B = \Gamma_B/2\pi$. The acoustic energy damping time is given by

$$\tau_B = 1/2\pi\Delta\nu_B \quad (116)$$

The electrostrictive gain coefficient g can also be written as

$$g^e = \frac{2\omega_S \gamma_e^2 \omega_L}{n_s c^3 \epsilon_o \rho_o \Gamma_B v_s} \sin(\theta/2) \frac{1}{1 + \left(\frac{\delta\Omega}{\Gamma_B/2}\right)^2} \quad (117a)$$

$$= \frac{\omega_S \gamma_e^2}{2\omega_L n_s n_L^2 c \epsilon_o \rho_o \Gamma_B v_s} \frac{1}{\sin(\theta/2)} \frac{1}{1 + \left(\frac{\delta\Omega}{\Gamma_B/2}\right)^2} \quad (117b)$$

The maximum steady-state gain is

$$g_{B,\max}^e = \frac{\omega_S \omega_L \gamma_e^2}{\pi n_s c^3 \epsilon_o \rho_o v_s \Delta\nu_B} \quad (118)$$

and occurs for backward scattering ($\theta = 180^\circ$). The frequency shift for backward scattering is

$$\Omega_B(180^\circ) = 2n\omega_L \frac{v_s}{c} \quad (119)$$

Since the Brillouin frequency shift is small, $\omega_S \approx \omega_L$ and the maximum Brillouin gain apparently scales as ω_L^2 . However, the linewidth also scales as ω_L^2 [see Eq. (115b)], leaving the maximum Brillouin gain independent of wavelength as indicated in Eq. (117b).

The dependence of the linewidth on q^2 causes the steady-state gain in the forward direction to go to infinity [Eq. (117b)] rather than to zero [Eq. (117a)]. However, the forward interaction is always transient because the steady-state time for forward scattering also goes to infinity. As a result, the forward gain is zero. Brillouin scattering at 90° is important in propagation of high-power laser radiation through large-diameter optics.

Transient Brillouin scattering has been described by Kroll¹⁶ and Faris et al.¹³⁰ The transient solutions to the stimulated Brillouin scattering are formally similar to those for Raman scattering and can be written as

$$A_S(z, \tau) = A_S(0, \tau) + \sqrt{\Gamma_B g_B z / cn_L \epsilon_o} A_L(\tau) \int_{-\infty}^{\tau} \frac{e^{-(\tau-\tau')/T_2} A_L^*(\tau') A_S(0, \tau') I_1(\sqrt{\Gamma_B g_B z [p(\tau) - p(\tau')] / cn_L \epsilon_o})}{\sqrt{p(\tau) - p(\tau')}} d\tau \quad (120)$$

Again, the transient gain depends on the time integral of the laser intensity, and the scattered intensity grows as a Bessel function. The steady-state solution of Eq. (112) applies when the pulse duration is greater than the steady-state time given by

$$t_{ss} = \frac{G_{ss}}{\Gamma_B} \quad (121)$$

where $G_{ss} = G_L L$.

Brillouin scattering in solids has been discussed in Refs. 131 to 134. The formal equations for electrostrictive Brillouin scattering are similar to those for liquids and gases. However, the gain coefficient depends on the polarization of the laser and scattered light, and Brillouin resonances exist for both longitudinal and shear acoustic waves. The Brillouin gain is given by

$$g^e = \frac{\omega_s \Omega_B n_s^3 n_L^4 p'^2}{n_s n_L c^2 \epsilon_0 \rho_o \Gamma_B v_s^2} \frac{1}{1 + \left(\frac{2\delta\Omega}{\Gamma_B} \right)} \quad (122)$$

where p' is the photoelastic constant appropriate for the specific combination of polarizations for the optical and acoustic waves. For longitudinal acoustic waves in isotropic materials¹³⁴

$$p' = p_{12}(\hat{\epsilon}_L \cdot \hat{\epsilon}_s) + p_{44}[(\hat{\epsilon}_L \cdot \hat{\kappa}_a)(\hat{\epsilon}_L \cdot \hat{\epsilon}_s) + (\hat{\epsilon}_L \cdot \hat{\epsilon}_s)(\hat{\epsilon}_s \cdot \hat{\kappa}_a)] \quad (123a)$$

while for shear acoustic waves

$$p' = p_{44}[(\hat{\epsilon}_L \cdot \hat{\kappa}_a)(\hat{\epsilon}_s \cdot \hat{\epsilon}_a) + (\hat{\epsilon}_L \cdot \hat{\epsilon}_a)(\hat{\epsilon}_s \cdot \hat{\kappa}_a)] \quad (123b)$$

where $\hat{\epsilon}_{L(S)(a)}$ is the unit polarization vector of the pump (Stokes) (acoustic) wave and $\hat{\kappa}_a$ is the unit propagation vector of the acoustic wave.

Buildup of noise in a generator with Brillouin scattering is similar in principle to buildup with Raman scattering. However, the primary noise source for Brillouin scattering is the thermally excited acoustic phonons. As a result, the equivalent noise level for Brillouin scattering can be several orders of magnitude larger than the Raman noise level.¹³⁵⁻¹³⁷ The noise power into a given solid angle $d\Omega = \pi\theta^2/4$ is given by¹³⁵

$$p_{\text{noise}} = \frac{\omega_s}{\Omega} k_B T \Delta v_B \left(\frac{\theta}{2\theta_D} \right)^2 \quad (124)$$

where $\theta_D = \lambda/2\pi D$ is the diffraction angle that can be resolved by a gain column of diameter D . Brillouin scattering with broadband radiation is discussed in Refs. 138 and 139. Multiline Brillouin scattering is treated in Ref. 140. As with backward Raman scattering, the gain decreases when the coherence length becomes comparable to or shorter than the interaction length.

Thermal Brillouin scattering is driven by nonuniform heating due to absorption and is described by g^a in Eq. (113b). The gain shows a dispersive behavior with frequency, with zero gain at the Brillouin frequency, gain for frequencies less than Ω_B , and loss for frequencies greater than Ω_B . Brillouin scattering parameters for various materials are listed in Table 10.

Brillouin Phase Conjugation

One of the major uses of stimulated Brillouin scattering (SBS) is phase conjugation.^{140a} Phase conjugation is also possible with stimulated Raman scattering, but SRS is not used for this purpose as much as SBS. Phase conjugation is used to correct distortions on optical beams that arise from propagation through nonideal optical media, such as the atmosphere or low-quality optical components. It can also be used for correcting or stabilizing aiming errors that arise from motion of components in an optical train, for improving the beam quality of laser oscillators and oscillator amplifier systems, and for beam combining.

A typical arrangement for phase conjugation is shown in Fig. 18. The distorting medium impresses a transverse phase variation on the optical beam propagating from left to right that can result in increased divergence, intensity structure, or reduced focal plane intensity if focused by a subsequent lens or mirror. If the beam entering the distorting material carried an initial distortion that would be undone by the material, then the beam would emerge from the medium undistorted.

TABLE 10 Brillouin Parameters for Various Materials

Liquids									
	Laser Wavelength (nm)	Frequency Shift (GHz)	ΔV (MHz)	τ_B (ns)	g_B (cm/GW)	n	Density (g/cm ³)	Reference	
Acetone	1059	2.987	119 ± 5	1.34	15.8	1.355	0.791	141	
	532	5.93	361	0.44	12.9	1.359(Na-D)		142	
	532	6.0	320	0.497	20			138	
Benzene	1059	4.124	228	0.7	9.6	1.4837	0.879	141	
	532	8.33	515	0.31	12.3	1.501(Na-D)		142	
Benzyl alcohol	532	9.38	2120	0.08	5.75	1.54(Na-D)	1.045	142	
Butyl acetate	532	6.23	575	0.28	9.13	1.394(Na-D)	0.882	142	
CS ₂	1060	3.761	50	3.2	68	1.595	1.262	141	
	532	7.7	120	1.9	130			138	
CCl ₄	1060	2.772	528	0.3	3.8	1.452	1.595	141	
	532	5.72	890	0.18	8.77	1.4595	1.594	142	
Chloroform	532	5.75	635	0.25	11.7	1.446(Na-D)	1.492	142	
Cyclohexane	532	7.19	1440	0.11	5.8	1.426(Na-D)	0.779	142	
N,N-Dimethyl formamide	532	7.93	615	0.26	7.8	1.431(Na-D)	0.944	142	
Dichloromethane	532	5.92	255	0.62	16.8	1.424	1.325	142	
o-Dichlorobenzene	532	8.03	1340	0.12	4.7	1.551	1.306	142	
Ethanol	532	5.91	546	0.29		1.36	0.785	142	
Ethylene glycol	532	10.2	3630	0.04	0.85	1.431	1.113	142	
Freon 113	532	3.72	865	0.18	5.5	1.3578	1.575	142	
n-Hexane	532	5.64	580	0.27	8.8	1.379	0.67	142	
Nitrobenzene	1060	4.255	396	0.4	7.2	1.5297	1.206	141	
Methanol	532	5.47	325	0.49	10.6	1.329	0.791	142	
	530	5.6	210	0.334	13			138	
Pyridine	532	8.92	746	0.21	14	1.51	0.978	142	
Tin tetrachloride	1064	2.21 ± 0.02	182 ± 12	0.874	11.2 ± 0.5	1.36	2.226	143	
	532	4.71	357	0.45					
Titanium tetrachloride	1060	3.070	216	0.735	14.2	1.577	1.73	141	
Toluene	532	7.72	1314	0.12	8.4	1.496	0.867	142	
Trichloroethylene	532	5.94	765	0.21	12	1.4755	1.464	142	
Water	1060	3.703	170	0.935	3.8	1.324	1	141	
	532	7.4	607	0.26	2.94	1.333	1	142	
Xylenes	532	7.74	1211	0.13	9.3	1.497	0.86	142	

(Continued)

TABLE 10 Brillouin Parameters for Various Materials (Continued)

Gases									
	Laser Wavelength (nm)	Frequency Shift (GHz)	$\Delta\nu$ (MHz)	τ_b (ns)	g_b (cm/GW)	n	Density (g/cm ³)	Reference	
Xenon									
7599 Torr	532	0.654 ± 0.024	98.1 ± 8.9		1.38 ± 0.19	1.0069	0.05767	134	
6840 Torr	532	0.627 ± 0.030	107.4 ± 16.9			1.0062	0.05159	134	
CClF ₃				$0.65 \lambda_p^2 P$				144	
3310 kPa (liquid)	1060		305					145	
3860	1060		155					145	
3950	1060		200					145	
32 atm (liquid)				6.2 ± 0.4				146	
Gas				$(4.78 \times 10^9 / \rho \lambda_L^2 + 3.25 \times 10^5)^{-1}$				144	
SF ₆				$(5.9 \times 10^9 / \rho \lambda_L^2 + 1.6 \times 10^5)^{-1}$				144	

P in atmospheres, ρ in kg/m^3

Solids

Substance	Polarization		Laser Wavelength (nm)	Frequency Shift (GHz)	$\Delta\nu$ (MHz)	τ_b (ns)	g_b (cm/GW)	n	Density (g/cm ³)	Reference
	P	k								
d-LAP			1053							
	x	$y = b = Y$		28.9 ($\theta = 180^\circ$)		5.5 ($\theta = 180^\circ$)	18 ($\theta = 180^\circ$)			147
				14.4 ($\theta = 90^\circ$)		11 ($\theta = 90^\circ$)	26 ($\theta = 90^\circ$)			147
	x	z	532	84.1 ± 3.5			20.7 ± 2.99	1.5090	1.600	134
	x	x	532	25.207 ± 0.042			10.99 ± 1.88	1.5090	1.600	134
	y	x	532	19.590 ± 0.087			27.96 ± 2.85	1.5764	1.600	134
	y	z	532	26.415 ± 0.009			12.25 ± 0.93	1.5764	1.600	134
	z	x	532	19.644 ± 0.048			29.85 ± 2.4	1.5847	1.600	134
	z	y	532	26.709 ± 0.039			24.33 ± 3.26	1.5847	1.600	134
		THG	532	94.8 ± 8.4			17.45 ± 3.41			134
			351	261 ($\theta = 180^\circ$)		0.61 ($\theta = 180^\circ$)	22 ($\theta = 180^\circ$)			147
				132 ($\theta = 90^\circ$)		1.2 ($\theta = 90^\circ$)	31 ($\theta = 90^\circ$)			147
Fused silica			1053							
			1053	40.8 ($\theta = 180^\circ$)		3.9 ($\theta = 180^\circ$)	4.8 ($\theta = 180^\circ$)		2.202	147
			532	20.4 ($\theta = 90^\circ$)		7.8 ($\theta = 90^\circ$)	6.8 ($\theta = 90^\circ$)			147
			532	163 ± 7.6			2.9 ± 0.015	1.4607		134
			532	32.65 ± 0.054			2.69 ± 0.22			134
			351	32.62						147
			351	177.6 ± 13.5						147
			351	185 ($\theta = 180^\circ$)		0.43 ($\theta = 180^\circ$)	5.4 ($\theta = 180^\circ$)			147
			532	101.5 ± 7.5		0.86 ($\theta = 90^\circ$)	7.6 ($\theta = 90^\circ$)			147
KD*P	$x = c = Z$	$z = Y$	532	29.763 ± 0.06			3.53 ± 0.31	1.4683	2.355	134
	z	x	532	27.627 ± 0.087			4.57 ± 0.38	1.5073	2.355	134
	z	Z	532	30.525 ± 0.156			5.09 ± 0.40	1.5073	2.355	134

KDP, THG	532				28.554 ± 0.036	72.9 ± 5.7											
LAA	532	b	y		20.892 ± 0.009	100.4 ± 7.6								6.5 ± 0.95	1.5073	2.355	134
BK3	532				31.383 ± 0.036	198.6 ± 6.6								24.9 ± 3.75	1.5008	2.37	134
LHG-8	532				27.786 ± 0.024	219.0 ± 6.2								1.78 ± 0.13	1.5316	2.83	134
BK7	532				34.65 ± 0.039	165.0 ± 8.6								2.74 ± 0.23	1.5195	2.51	134
CaF2	532				37.164 ± 1.185	45.6 ± 8.8								2.15 ± 0.21	1.4354	3.179	134
Plexiglas	532				15.687 ± 0.036	253.7 ± 12.6								4.11 ± 0.65	1.4938	1.19	134
GGG	532				26.283 ± 0.005	12.5 ± 6.9								1.02 ± 0.5	1.9788	7.09	134

Glasses

	Laser Wavelength (nm)	Frequency Shift (GHz)	ΔV (MHz)	v_s (m/s)	g_B (cm/GW)	n	Density (g/cm ³)	Reference
SiO ₂	488	35.6	156	5944.2	4.482	1.462	2.203	132
ZBL	488	25.0	213.6	3979.0	2.832	1.530	4.672	132
ZBLA	488	25.2	98.7	3968.4	1.713	1.548	4.579	132
ZBLAN	488	26.6	96.0	4270	3.608	1.521	4.301	132
HBL	488	22.4	151.4	3608.9	1.127	1.514	5.78	132
HBLA	488	22.1	162.3	3470	0.96	1.554	5.83	132
HBLAPC	488	25.2	179.5	4035	1.023	1.524	5.1	132
BeF ₂	488	24.3	52.5	4634.1	16.06	1.28	2.01	132
95BeF ₂ -5THF ₄	488	24.9	74.8	4638.5	11.54	1.31	2.1	132
4.97 Li ₂ O 94.27 B ₂ O ₃ 0.13 Al ₂ O ₃	488	26.9	100	4364	12.88	1.5056	1.9434	132
4.97 Li ₂ O 94.27 B ₂ O ₃ 0.13 Al ₂ O ₃	488	26.8	116	4351	14.29	1.5053	1.9439	132
4.97 Li ₂ O 94.27 B ₂ O ₃ 0.13 Al ₂ O ₃	488	24.8	113	4050	11.74	1.4946	1.9025	132
4.97 Li ₂ O 94.27 B ₂ O ₃ 0.13 Al ₂ O ₃	488	32.1	124	5109	10.93	1.5324	2.0753	132
4.97 Li ₂ O 94.27 B ₂ O ₃ 0.13 Al ₂ O ₃	488	29.0	138	4662	12.66	1.5189	1.9823	132
4.97 Li ₂ O 94.27 B ₂ O ₃ 0.13 Al ₂ O ₃	488	33.0	104	5205	3.441	1.5455	2.0541	132
4.97 Li ₂ O 94.27 B ₂ O ₃ 0.13 Al ₂ O ₃	488	33.1	184	5458	3.905	1.4781	2.2416	132
6 K ₂ O-94SiO ₂ , 673°C 78.56 h	488	33.1	186	5465	2.783	1.4781	2.2519	132
6 K ₂ O-94SiO ₂ , 604.8°C 231.8 h	488	32.4	175	5330	3.587	1.4845	2.2703	132
8 K ₂ O-92SiO ₂ , 547°C	488	32.4	32.5	5335	3.587	1.4844	2.2725	132
8 K ₂ O-92SiO ₂ , 603.5°C 103.42 h	488	32.4	177	5332	5.181	1.4844	2.2716	132
8 K ₂ O-92SiO ₂ , 589.5°C 141.79 h	488	32.4	170	5321	3.704	1.4844	2.2696	132
8 K ₂ O-92SiO ₂ , 574.3°C 345.68 h	488	32.5	187	5339	4.038	1.4846	2.2702	132
8 K ₂ O-92SiO ₂ , 573.4°C 84.66 h	488	30.8	208	5039	3.122	1.4923	2.2991	132
10 K ₂ O-90SiO ₂	488	25.8	104	4190	8.169	1.5037	2.0281	132
10 K ₂ O-90SiO ₂	488	30.2	134	4839	4.675	1.5231	2.2699	132

Reprinted with permission from M. X. Weber (ed), *CRC Handbook of Laser Science and Technology: Optical Materials*, suppl. 2. Copyright CRC Press, Boca Raton, FL, 1988.

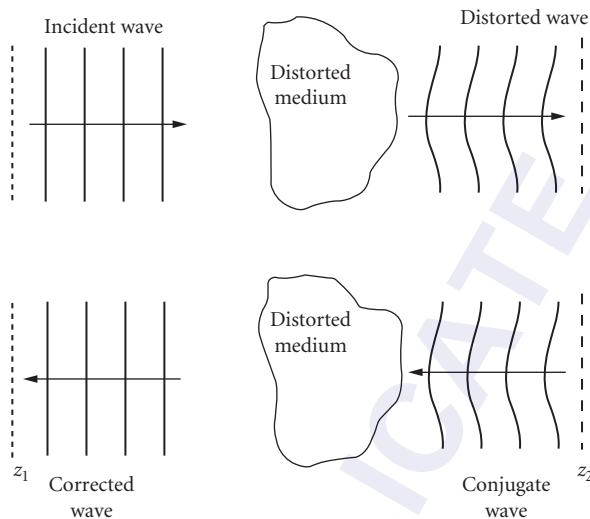


FIGURE 18 Schematic diagram for correction of phase distortions caused by distorting medium with double-pass propagation.

The general concept is that the initial distortion required for this compensation is the inverse of the distortion imposed by the medium. However, determining precisely what that distortion is and imposing it on the initial beam is the heart of the matter in phase compensation.

The arrangement shown in Fig. 19 illustrates how this is done with nonlinear optical phase conjugation. The beam propagating from left to right is initially undistorted. It emerges from the medium at point B with a complete record of the distortion imposed by the medium. In the nonlinear medium, the beam undergoes a backward nonlinear optical interaction in which a second beam is generated that travels in the opposite direction and has phase variations that are reversed from the original distorted beam. When this beam propagates through the distorting medium again from right to left, the medium again impresses a phase variation on it, but this time the process simply undoes the initial phase distortion instead of creating one. The beam emerging from the distorting medium at point A is undistorted.

Several types of nonlinear interactions have been used for phase conjugation, including degenerate and nondegenerate four-wave mixing, stimulated Brillouin scattering, Brillouin-enhanced four-wave mixing, and stimulated Raman scattering. Of these, degenerate four-wave mixing is used most often for low-power interactions and stimulated Brillouin scattering is used for high-power applications. Brillouin-enhanced four-wave mixing provides a high-gain Brillouin system.

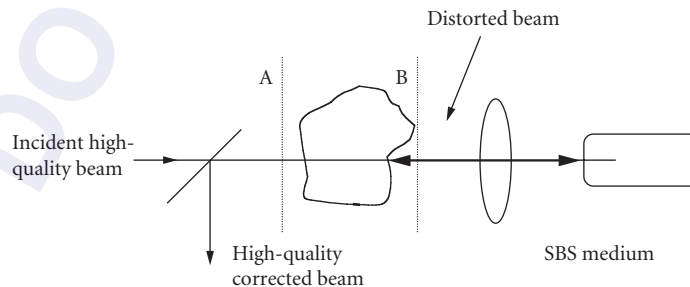


FIGURE 19 Schematic diagram of use of stimulated Brillouin scattering for phase distortion correction with phase conjugation.

Phase conjugation with SBS is the result of mode competition in which the mode that corresponds to the conjugate of the pump has higher gain than other possible modes. The gain for the backward-generated beam involves a sum over all possible spatial modes of the backward Stokes wave. The mode that matches the pump wave has the highest gain because the high- and low-intensity points coincide, providing reinforcement of that mode over the entire interaction length. The exponential gain for this mode is twice that for the noncorrelated modes when the pump beam carries a large number of spatial modes.

In order for the conjugate mode to dominate over the other modes, it is necessary that it experiences preferential gain over an extended region. Essentially, it is required that the Stokes beam diffract across the entire pump beam in the interaction region before it grows to a saturation level. This has been accomplished by using a light guide to confine the beams. In free-focused interactions, it requires that the divergence of the pump beam be sufficient to provide a wide focal area relative to its depth. An example of phase conjugate correction in a two-pass Q-switched Nd:YAG amplifier system is shown in Fig. 20 using liquid Freon-113 as the conjugating medium.¹⁴⁸ Properties of phase-conjugated beams, quality of correction, and efficiency of conversion are given in Ref. 10.

Brillouin-Enhanced Four-Wave Mixing Brillouin-enhanced four-wave mixing (BEFWM) is a nearly degenerate four-wave mixing process where the four waves are coupled by the Brillouin nonlinearity as shown schematically in the Fig. 21. BEFWM was first observed and explained by Basov et al.¹⁴⁹ In addition to the beam to be conjugated, two oppositely propagating pump beams are

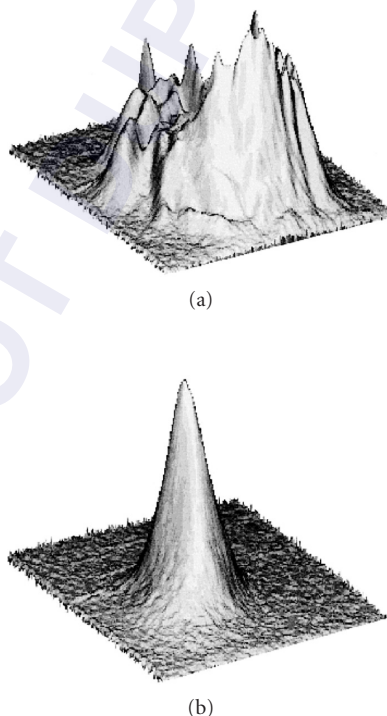


FIGURE 20 Spatial beam profiles of a Q-switched Nd:YAG oscillator-amplifier system using a two-pass amplifier configuration with a conventional mirror in the two-pass amplifier (a) and an SBS mirror (b). The Brillouin mirror was a cell of liquid Freon-113. (From Ref. 148.)



FIGURE 21 k vector diagram for Brillouin-enhanced four-wave mixing.

provided at frequency ω_L . In general, the two pump beams do not need to be at the same frequency, but in practice they usually are. The signal beam k_s at a frequency $\omega_s = \omega_p - \Delta\omega_B$ interacts with the laser beam k_1 at a frequency ω_p to produce an acoustic wave that moves in the direction of the laser beam. The second laser beam k_2 scatters off the moving grating to generate the conjugated anti-Stokes beam k_{AS} at a frequency $\omega_{AS} = \omega_p + \Delta\omega_B$. The phase-conjugated beam, traveling in the opposite direction of the signal beam, can be amplified when the pump beams do not interact with each other. Therefore, BEFWM represents a mirror with reflectivity greater than 1. The pump beams are usually decoupled by frequency or polarization control. A reflectivity as high as $\sim 7 \times 10^5$ has been observed.¹⁵⁰ The process can also be done with the anti-Stokes wave as the signal beam. The signal Stokes beam is usually generated externally via a separate SBS process. The angles between the beams can be adjusted to make the phase matching possible. Various aspects of BEFWM are described in detail in Ref. 151.

15.4 REFERENCES

1. N. Bloembergen, *Nonlinear Optics*, Benjamin, New York, 1965.
2. W. Kaiser and M. Maier, "Stimulated Rayleigh, Brillouin and Raman Spectroscopy," in F. T. Arecchi and E. O. Schulz-DuBois (eds.), *Laser Handbook*, vol. 2, North-Holland, Amsterdam, New York, 1972, pp. 1077–1150.
3. C. S. Wang, "The Stimulated Raman Process," in Herbert Rabin and C. L. Tang (eds.), *Quantum Electronics: A Treatise*, vol. 1, part A, Academic Press, New York, 1975.
4. A. Penzkofer, A. Laubereau, and W. Kaiser, "High Intensity Raman Interactions," *Prog. Quant. Electron.* **6**:55–140 (1979).
5. F. Milanovich, "Stimulated Raman Scattering," in M. J. Weber (ed.), *Handbook of Laser Science and Technology: Optical Materials*, vol. III, CRC, Boca Raton, FL, 1986.
6. Y. R. Shen, *Principles of Nonlinear Optics*, Wiley, New York, 1984.
7. R. W. Boyd, *Nonlinear Optics*, Academic, New York, 1992.
8. J. F. Reintjes, "Stimulated Raman and Brillouin Scattering," in M. J. Weber (ed.), *Handbook of Laser Science and Technology: Optical Materials*, suppl. 2, CRC, Boca Raton, FL, 1995, pp. 334–364.
9. "Stimulated Raman and Brillouin Scattering for Laser Beam Control," special issue of *JOSA B* **3** (October 1986).
10. J. Reintjes (ed.), "Laser Wavefront Control," *SPIE* **1000** (1988).
11. D. Levi, C. R. Menyuk, and P. Winternitz (eds.), *Self-Similarity in Stimulated Raman Scattering*, Les Publications CRM, Montreal, Canada, 1994.
12. N. Bloembergen, "The Stimulated Raman Effect," *Am. J. Phys.* **35**:989 (1967).
13. A. Z. Grasyuk, "Raman Lasers (Review)," *Sov. J. Quant. Electron.* **4**:269 (1974).
14. G. Eckardt, R. W. Hellwarth, F. J. McClung, S. E. Schwarz, and D. Weiner, "Stimulated Raman Scattering from Organic Liquids," *Phys. Rev. Lett.* **9**:455 (1962).
15. Y. R. Shen and N. Bloembergen, "Theory of Stimulated Raman and Brillouin Scattering," *Phys. Rev.* **137**:A1787 (1965).

16. N. M. Kroll, "Excitation of Hypersonic Vibrations by Means of Photoelastic Coupling of High Intensity Light Waves to Elastic Waves," *J. Appl. Phys.* **36**:34 (1965).
17. R. L. Carman, F. Shimizu, C. S. Wang, and N. Bloembergen, "Theory of Stokes Pulse Shapes in Transient Stimulated Raman Scattering," *Phys. Rev. A* **2**:60 (1970).
18. C. S. Wang, "Theory of Stimulated Raman Scattering," *Phys. Rev.* **182**:482 (1969).
19. M. D. Duncan, R. Mahon, L. L. Tankersley, and J. Reintjes, "Transient Stimulated Raman Amplification in Hydrogen," *JOSA B* **5**:37 (1988).
20. J. R. Murray and A. Javan, "Effects of Collisions on Raman Line Profiles of Hydrogen and Deuterium Gas," *J. Mol. Spectrosc.* **42**:1-26 (1972).
21. R. H. Dicke, *Phys. Rev.* **89**:472 (1953).
22. W. K. Bischel and M. J. Dyer, "Temperature Dependence of the Raman Linewidth and Lineshift for the Q(1) and Q(0) Transitions in Normal and Para H₂," *Phys. Rev. A* **33**:3113 (1986).
23. S. Kern and B. Feldman, *Stimulated Raman Emission*, vol. 3, Massachusetts Institute of Technology, Lincoln Laboratory, Bedford, MA, 1974, p. 18.
24. J. J. Barrett and M. C. Tobin, "Stimulated Raman Emission Frequencies in 21 Organic Liquids," *J. Opt. Soc. Am.* **56**:129 (1966).
25. M. D. Martin and E. L. Thomas, "Infrared Difference Frequency Generation," *IEEE J. Quant. Electron.* **QE-2**:196 (1966).
26. M. A. El-Sayed, F. M. Johnson, and J. Duardo, "A Comparative Study of the Coherent Raman Processes Using the Ruby and the Second Harmonic Neodymium Giant-Pulsed Lasers," *J. Chem. Phys.* **1**:227 (1967).
27. J. A. Giordmaine and J. A. Howe, "Intensity-Induced Optical Absorption Cross Section in CS₂," *Phys. Rev. Lett.* **11**:207 (1963).
28. T. A. Prasada Rao and N. Seetharaman, "Amplification of Stimulated Raman Scattering by a Dye," *Ind. J. Pure Appl. Phys.* **13**:207 (1975).
29. M. Geller, D. P. Bortfeld, and W. R. Sooy, "New Woodbury-Raman Laser Materials," *Appl. Phys. Lett.* **3**:36 (1963).
30. W. L. Smith and F. P. Milanovich, Lawrence Livermore National Laboratory, Livermore, CA, Private Communication, 1983.
31. J. R. Maple and J. T. Knudtson, "Transient Stimulated Vibrational Raman Scattering in Small Molecule Liquids," *Chem. Phys. Lett.* **56**:241 (1978).
32. J. K. Wright, C. H. H. Carmichael, and B. J. Brown, "Narrow Linewidth Output from a Q-switched Nd 3-/Glass Laser," *Phys. Lett.* **16**:264 (1965).
33. M. K. Srivastava and R. W. Crow, "Raman Susceptibility Measurements and Stimulated Raman Effect in KDP," *Opt. Commun.* **8**:82 (1973).
34. P. D. Maker and R. W. Terhune, "Study of Optical Effects due to an Induced Polarization Third Order in the Electric Field Strength," *Phys. Rev.* **137**:A801 (1965).
35. D. P. Bortfeld, M. Geiller, and G. Eckhardt, "Combination Lines in the Stimulated Raman Spectrum of Styrene," *J. Chem. Phys.* **40**:1770 (1964).
36. V. A. Orlovich, "Measurement of the Coefficient of Stimulated Raman Scattering in Organic Liquids with the Aid of an Amplifier with Transverse Pumping," *Zh. Prikl. Spektrosk.* **23**:224 (1975).
37. J. A. Calvieilo and Z. H. Heller, "Raman Laser Action in Mixed Liquids," *Appl. Phys. Lett.* **5**:112 (1964).
38. G. Eckhardt, "Selection of Raman Laser Materials," *IEEE J. Quant. Electron.* **QE-2**:1 (1966).
39. V. A. Subov, M. M. Sushchinskii, and I. K. Shuvalton, "Investigation of the Excitation Threshold of Induced Raman Scattering," *J. Exp. Theor. Phys. USSR* **47**:784 (1964).
40. C. D. Decker, "High-Efficiency Stimulated Raman Scattering/Dye Radiation Source," *Appl. Phys. Lett.* **33**:323 (1978).
41. B. P. Stoicheff, "Characteristics of Stimulated Raman Radiation Generated by Coherent Light," *Phys. Lett.* **7**:186 (1963).
42. D. Cotter, D. C. Hanna, and R. Wyatt, "A High Power, Widely Tunable Infrared Source Based on Stimulated Electronic Raman Scattering in Caesium Vapour," *Opt. Commun.* **16**:256 (1976).
43. P. E. Tannenwaid and J. B. Thaxter, "Stimulated Brillouin and Raman Scattering in Quartz at 2.1 to 293 Kelvin," *Science* **134**:1319 (1966).

44. J. Gelbwachs, R. H. Pantell, H. E. Puthoff, and J. M. Yarborough, "A Tunable Stimulated Raman Oscillator," *Appl. Phys. Lett.* **14** (1969).
45. G. Eckhardt, D. P. Bortfeld, and M. Geller, "Stimulated Emission of Stokes and Anti-Stokes Raman Lines from Diamond, Calcite and α -Sulfur Single Crystals," *Appl. Phys. Lett.* **3**:137 (1963).
46. I. V. Aleksandrov, Y. S. Bobovitch, A. V. Bortkevich, and M. Y. Tsenter, "Raman Scattering in Crystalline Chlorine and Bromine," *Opt. Spektrosk.* **36**:150 (1974).
47. I. I. Kondilenko, P. A. Korotkov, and V. I. Maly, "Temperature Dependence of SRS Thresholds for Some Salts of Inorganic Acids in the Crystalline Phase," *Opt. Commun.* **10**:50 (1974).
48. D. L. Weinberg, "Stimulated Raman Emission in Crystals and Organic Liquids," *Mass. Inst. Technol. Lincoln Lab. Solid-State Res. Rep.* **3**:31 (1965).
49. O. S. Avanesyan, V. A. Benderskii, V. K. H. Brikshtein, V. L. Broude, A. G. Lavrushko, I. I. Tartakovskii, and P. V. Filippov, "Characteristics of Stimulated Light Generation and Stimulated Raman Scattering in Anthracene Crystals," *Sov. J. Quant. Electron.* **7**:403 (1977).
50. J. L. Carlsten and P. C. Dunn, "Stimulated Stokes Emission with a Dye Laser: Intense Tunable Radiation in the Infrared," *Opt. Commun.* **14**:8 (1975).
51. D. Cotter, D. C. Hanna, P. A. Kärkkäinen, and R. Wyatt, "Stimulated Electronic Raman Scattering as a Tunable Infrared Source," *Opt. Commun.* **15**:143 (1975).
52. P. P. Sorokin, J. J. Wynne, and J. R. Landkard, "Tunable Coherent IR Source Based upon Four-Wave Parametric Conversion in Alkali Metal Vapors," *Appl. Phys. Lett.* **22**:342 (1973).
53. A. DeMartino, R. Frey, and F. Pradere, "Tunable Far Infrared Generation in Hydrogen Fluoride," *Opt. Commun.* **27**:262 (1978).
54. P. May, P. Bernage, and H. Boequet, "Stimulated Electronic Raman Scattering in Rubidium Vapour," *Opt. Commun.* **29**:369 (1979).
55. R. L. Byer and W. R. Trutna, "16- μ m Generation by CO₂-Pumped Rotational Raman Scattering in H₂," *Opt. Lett.* **3**:144 (1978).
56. P. Rabinowitz, A. Stein, R. Brickman, and A. Kaldor, "Efficient Tunable H₂ Raman Laser," *Appl. Phys. Lett.* **35**:739 (1979).
57. E. Pochon, "Determination of the Spontaneous Raman Linewidth of CF₄ by Measurements of Stimulated Raman Scattering in Both Transient and Steady States," *Chem. Phys. Lett.* **77**:500 (1981).
58. B. E. Kinkaid and J. R. Fontana, "Raman Cross-Section Determination by Direction Stimulated Raman Gain Measurements," *Appl. Phys. Lett.* **28**:12 (1975).
59. M. Rokni and S. Yatsiv, "Resonance Raman Effects in Free Atoms of Potassium," *Phys. Lett.* **24**:277 (1967).
60. R. W. Minek, R. W. Terhune, and W. G. Rado, "Laser-Stimulated Raman Effect and Resonant Four-Photon Interactions in Gases H₂, D₂, and CH₄," *Appl. Phys. Lett.* **3**:181 (1963).
61. H. Komine, Northrop Corp., Palos Verdes, Calif., Private Communication, 1983.
62. W. K. Bischel and M. J. Dyer, "Wavelength Dependence of the Absolute Raman Gain Coefficient for the Q(1) Transition in H₂," *JOSA* **B5**:677 (1986).
63. W. K. Bischel and G. Black, "Wavelength Dependence of Raman Scattering Cross Sections from 200–600 nm," in C. K. Rhodes, H. Egger, and H. Pummer (eds.), *Excimer Lasers—1983*, AIP, New York, 1983.
64. W. K. Bischel, "Stimulated Raman Gain Processes in H₂, HD and D₂," unpublished data.
65. J. J. Ottusch and D. A. Rockwell, "Measurement of Raman Gain Coefficients of Hydrogen, Deuterium and Methane," *IEEE J. Quant. Electron.* **QE-24**:2076 (1988).
66. D. A. Russel and W. B. Roh, "High Resolution CARS Measurements of Raman Linewidths of Deuterium," *J. Mol. Spect.* **24**:240 (1987).
67. K. C. Smyth, G. J. Rosasco, and W. S. Hurst, "Measurements and Rate Law Analysis of D2 Q-Branch Line Broadening Coefficients for Collisions with D₂, He, Ar, H₂ and CH₄," *J. Chem. Phys.* **87**:1001 (1987).
68. D. Haner and I. S. McDermid, "Stimulated Raman Shifting of the Nd:YAG Fourth Harmonic (266 nm) in H₂, HD and D₂," *IEEE J. Quant. Electron.* **QE-26**:1292 (1990).
69. J. R. Murray, J. Goldhar, D. Eimerl, and A. Szoke, "Raman Pulse Compression of Excimer Lasers for Application to Laser Fusion," *IEEE J. Quant. Electron.* **QE-15**:342 (1979).
70. G. C. Herring, M. X Dyer, and W. K. Bischel, "Temperature and Wavelength Dependence of the Rotational Raman Gain Coefficient in N₂," *Opt. Lett.* **11**:348 (1986).

71. A. P. Hickman, J. A. Paisner, and W. K. Bischel, *Phys. Rev. A* **33**:1788 (1986).
72. H. Chuan and T. H. Chyba, "Solid-State Barium Nitrate Raman Laser in the Visible Region," *Opt. Commun.* **135**:273 (1997).
73. P. G. Zverev, T. T. Basiev, V. V. Osiko, A. M. Kulkov, V. N. Voitsekhovskii, and V. E. Yakobsen, "Physical, Chemical and Optical Properties of Barium Nitrate Raman Crystal," *Opt. Mat.* **11**:315 (1999).
74. P. G. Zverev and T. T. Basiev, "Investigation of the Line Broadening of an SRS-Active Vibration in a Barium Nitrate Crystal by Two-Photon Raman Amplification Spectroscopy," *Quant. Electron.* **25**:1204 (1995).
75. V. N. Voitsekhovskii, S. N. Karpukhin, and V. E. Yakobson, "Single-Crystal Barium Nitrate and Sodium Nitrate as Efficient Materials for Laser-Radiation Frequency Conversion Based on Stimulated Raman Scattering," *J. Opt. Tech.* **62**:770 (1995).
76. P. G. Zverev, J. T. Murray, R. C. Powell, R. J. Reeves, and T. T. Basiev, "Stimulated Raman Scattering of Picosecond Pulses in Barium Nitrate Crystals," *Opt. Commun.* **97**:59 (1993).
77. A. Owyong and P. S. Percy, "Precise Characterization of the Raman Nonlinearity in Benzene Using Nonlinear Interferometry," *J. Appl. Phys.* **48**:674 (1977).
78. G. Herring, M. J. Dyer, and W. K. Bischel, "Temperature and Density Dependence of the Linewidths and Lineshifts of the Rotational Raman Lines in N_2 and H_2 ," *Phys. Rev. A* **34**:1944 (1986).
79. G. Hilfer and C. R. Menyuk, "Stimulated Raman Scattering in the Transient Limit," *J. Opt. Soc. Am. B* **7**:739–749 (1990).
80. M. J. Colles, *Appl Phys. Lett.* **19**:23 (1971).
81. C. R. Menyuk, in D. Levi, C. R. Menyuk, and P. Winternitz (eds.), *Self-Similarity in Stimulated Raman Scattering*, Les Publications CRM, Montreal, Canada, 1994.
82. K. J. Druhl, R. G. Wenzel, and J. L. Carlsten, *Phys. Rev. Lett.* **51**:1171 (1983).
83. C. R. Menyuk, *Phys. Rev. A* **47**:2235 (1993).
84. W. R. Trutna, Y. K. Park, and R. L. Byer, *IEEE J. Quant. Electron.* **QE-15**:648 (1979).
85. M. G. Raymer, J. Mostowski, and J. L. Carlsten, "Theory of Stimulated Raman Scattering with Broadband Lasers," *Phys. Rev. A* **19**:2304 (1979).
86. J. Eggleston and R. L. Byer, *IEEE J. Quant. Electron.* **QE-16**:850 (1980).
87. E. A. Stappaerts, W. H. Long Jr., and H. Komine, *Opt. Lett.* **5**:4 (1980).
88. A. Flusberg, D. Korff, and C. Duzy, "The Effect of Weak Dispersion on Stimulated Raman Scattering," *IEEE J. Quant. Electron.* **QE-21**:232 (1985).
89. K. A. Druhl, "Coherence Properties of Stokes Beams for Incoherent Broadband Pumps," *JOSA B* **3**:1363 (1986).
90. D. Korff, E. Mazur, C. Duzy, and A. Flusberg, "Raman Conversion Using Crossed Broadband Pump Beams and Bisecting Stokes," *JOSA B* **3**:1333 (1986).
91. G. C. Lombardi and H. Injeyan, "Phase Correlation Effects in a Raman Amplifier," *JOSA B* **3**:1461 (1986).
92. M. Bashkansky and J. Reintjes, "Correlation Effects in Pump-Depleted Broad-Band Stimulated Raman Amplification," *Opt. Commun.* **83**:103 (1991).
93. M. Bashkansky and J. Reintjes, "Incoherent Multimode Raman Amplification Theory," *JOSA B* **8**: 1843 (1991).
94. M. D. Skeldon and R. Bahr, "Stimulated Rotational Raman Scattering in Air with a High Power Broadband Laser," *Opt. Lett.* **16**:366 (1991).
95. M. Trippenbach, K. Rzazewski, and M. G. Raymer, "Stimulated Raman Scattering of Colored Light," *JOSA B* **1**:671 (1984).
96. B. W. Shore, S. Lowder, and M. A. Johnson, "Some General Properties of Stimulated Raman Propagation with Pump Depletion, Transiency and Dispersion," LLNL Report No. UCRL-ID-1-7967 R1, 1991.
97. M. G. Raymer and J. Mostowski, "Stimulated Raman Scattering: Unified Treatment of Spontaneous Initiation and Spatial Propagation," *Phys. Rev. A* **24**:1980 (1981).
98. G. P. Dzhotyan, Y. E. D'yakov, I. G. Zubarev, A. B. Mironov, and S. I. Mikhailov, "Amplification during Stimulated Raman Scattering in a Nonmonochromatic Pump Field," *Sov. Phys. JETP* **46**:431 (1978).
99. S. A. Akhmanov, Y. E. D'yakov, and L. I. Pavlov, "Statistical Phenomena in Raman Scattering Stimulated by a Broad-Band Pump," *Sov. Phys. JETP* **39**:249 (1975).

100. M. D. Duncan, R. Mahon, L. L. Tankersley, and J. Reintjes, "Spectral and Temporal Characteristics of Spontaneous Raman Scattering in the Transient Regime," *JOSA B* **8**:300 (1991).
101. J. C. White and D. Henderson, "Anti-Stokes Raman Lasers," *AIP Conf. Proc.* **90**:117–127 (1982).
102. J. C. White and D. Henderson, "TI: Threshold and Dispersion Effects in the Anti-Stokes Raman Laser," *Opt. Lett.* **7**:517 (1982).
103. K. Ludewigt, R. Dierking, W. Pflingsten, and B. Wellegehausen, "Vacuum Ultraviolet Anti-Stokes Raman Lasers," *IEEE J. Quant. Electron.* **QE 22**:1967–1974 (1986).
104. A. Z. Grasyuk, L. L. Losev, A. P. Lutsenko, and S. N. Sazonov, "Parametric Raman Anti-Stokes Laser," *Sov. J. Quant. Electron.* **20**:1153–1155 (1990).
105. N. Bloembergen and Y. R. Shen, "Coupling between Vibrations and Light Waves in Raman Laser Media," *Phys. Rev. Lett.* **12**:504 (1964).
106. M. D. Duncan, R. Mahon, L. L. Tankersley, and J. Reintjes, "Parametric Raman Gain Suppression in D₂ and H₂," *Opt. Lett.* **11**:803 (1986).
107. C. Reiser, T. D. Raymond, R. B. Michie, and A. P. Hickman, "Efficient Anti-Stokes Raman Conversion in Collimated Beams," *JOSA B* **6**:1859 (1989).
- 107a. K. Leung, M. Oron, D. Klimek, R. Holmes, and A. Flusberg, "Observation of Parametric Gain Suppression in Rotational Raman Transitions in N₂ and H₂," *Opt. Lett.* **13**:33 (1988).
- 107b. B. N. Perry, P. Raninowitz, and O. S. Bomse, *Opt. Lett.* **10**:146 (1985).
108. M. D. Levenson, "Coherent Raman Spectroscopy," (see also included references), *Phys. Today* (May 1977).
109. C. M. Bowden and J. C. Englund, "Macroscopic Manifestation of Quantum Noise," *Opt. Commun.* **67**:71 (1988).
110. M. G. Raymer, I. A. Walmsley, X. Mostowski, and B. Sobolewska, "Quantum Theory of Spatial and Temporal Coherence Properties of Stimulated Raman Scattering," *Phys. Rev. A* **32**:332 (1985).
111. M. G. Raymer, Z. W. Li, and I. A. Walmsley, "Temporal Quantum Fluctuations in Stimulated Raman Scattering," *Phys. Rev. Lett.* **63**:1586 (1989).
112. M. G. Raymer and L. A. Westling, "Quantum Theory of Stokes Generation with a Multimode Laser," *JOSA B* **2**:1417 (1985).
113. M. D. Duncan, R. Mahon, L. L. Tankersley, and J. Reintjes, "Low-Light Level, Quantum-Noise-Limited Amplification in a Stimulated Raman Amplifier," *JOSA B* **9**:2107 (1992).
114. R. C. Swanson, P. R. Battle, and J. Carlsten, "Interferometric Measurement of Quantum Noise in a Raman Amplifier," *Phys. Rev. Lett.* **67**:38 (1991).
115. N. Fabricius, K. Nattermann, and D. von der Linde, "Macroscopic Manifestation of Quantum Fluctuations in Transient Stimulated Raman Scattering," *Phys. Rev. Lett.* **52**:113 (1984).
116. K. Nattermann, N. Fabricius, and D. von der Linde, "Observation of Transverse Effects on Quantum Fluctuations in Stimulated Raman Scattering," *Opt. Commun.* **57**:212 (1986).
117. M. G. Raymer, K. Rzazewski, and J. Mostowski, "Pulse Energy Statistics in Stimulated Raman Scattering," *Opt. Lett.* **7**:71 (1982).
118. K. Rzazewski, M. Lewenstein, and M. G. Raymer, "Statistics of Stimulated Stokes Pulse Energies in the Steady State Regime," *Opt. Commun.* **43**:451 (1982).
119. D. C. MacPherson, R. C. Swanson, and J. L. Carlsten, "Quantum Fluctuations in the Stimulated Scattering Linewidth," *Phys. Rev. Lett.* **61**:66 (1988).
120. M. G. Raymer and I. A. Walmsley, "Quantum Statistics of Stimulated Raman Scattering," in L. Mandel (ed.), *Coherence and Quantum Optics*, vol. V, Plenum, New York, 1983, p. 63.
121. M. D. Duncan, R. Mahon, L. L. Tankersley, and J. Reintjes, "Imaging through a Low Light Level Amplifier," *SPIE* **1409**:127–134 (1991).
122. M. D. Duncan, R. Mahon, L. L. Tankersley, and J. Reintjes, "Second Stokes Generation in Deuterium and Hydrogen," *Opt. Commun.* **86**:538–546 (1991).
123. G. I. Kachen and W. H. Lowdermilk, "Subnanosecond Pulsations in Forward and Backward Stimulated Raman Scattering," *Opt. Commun.* **18**:112 (1976).
124. M. Rokni and A. Flusberg, "Stimulated Rotational Raman Scattering in the Atmosphere," *IEEE J. Quant. Electron.* **QE-22**:1102 (1986).

125. M. D. Levenson, *Introduction to Nonlinear Laser Spectroscopy*, Academic Press, New York, 1982.
126. H. Lotem and R. T. Lynch Jr., *Phys. Rev. Lett.* **37**:334 (1976).
127. M. D. Duncan, R. Mahon, L. L. Tankersley, and J. Reintjes, "Time-Gated Imaging through Scattering Media Using Stimulated Raman Amplification," *Opt. Lett.* **16**:1868–1870 (1991).
128. R. Mahon, M. D. Duncan, L. L. Tankersley, and J. Reintjes, "Time-Gated Imaging through Dense Scatterers with a Raman Amplifier," *Appl. Opt.* **32**:7425–7433 (1993).
129. M. Bashkansky, P. R. Battle, R. Mahon, and J. Reintjes, "Subsurface Defect Detection in Ceramic Materials Using Ultrafast Optical Techniques," *22nd Ann. Rev. of Prog. in Quantitative Nondestructive Evaluation*, Seattle, WA, 1995.
130. G. W. Faris, M. J. Dyer, and A. Peet Hickman, "Transient Effects on Stimulated Brillouin Scattering," *Opt. Lett.* **17**:1049 (1992).
131. L. G. Hwa, J. Schroeder, and X.-S. Zhao, "Intrinsic Brillouin Linewidths and Stimulated Brillouin Coefficients in Glasses Studied by Inelastic Light Scattering," *JOSA B* **6**:833 (1989).
132. J. Schroeder, L. G. Hwa, G. Kendall, C. S. Dumais, M. C. Shyong, and D. A. Thompson, "Inelastic Light Scattering in Halide and Oxide Glasses: Intrinsic Brillouin Linewidths and Stimulated Brillouin Gain," *J. Noncryst. Solids* **102**:240 (1988).
133. G. W. Faris, L. E. Jusinski, M. J. Dyer, W. K. Bischel, and A. Peet Hickman, "High Resolution Brillouin Gain Spectroscopy in Fused Silica," *Opt. Lett.* **15**:703 (1990).
134. G. W. Faris, L. E. Jusinski, and A. Peet Hickman, "High Resolution Stimulated Brillouin Gain Spectroscopy in Glasses and Crystals," *JOSA B* **10**:587 (1990).
135. A. M. Scott, D. E. Watkins, and P. Tapster, "Gain and Noise Characteristics of Brillouin Amplifier and Their Dependence on the Spatial Structure of the Pump Beam," *JOSA B* **7**:929 (1990).
136. R. W. F. Gross, S. T. Amimoto, and L. Garman-DuVall, "Gain and Phase Conjugation Fidelity of a Four Wave Brillouin Mirror Based on Methane," *Opt. Lett.* **16**:94 (1991).
137. V. I. Bespalov, O. V. Kulagin, A. I. Makarov, G. A. Pasmanik, A. K. Potjomkin, P. B. Potlov, and A. A. Shilov, "High-Sensitivity Optical System with Laser Amplifiers and Phase Conjugating Mirrors," *Opt. Acoust. Rev.* **1**:71 (1989).
138. P. Narum, M. D. Skeldon, and R. W. Boyd, "Effect of Laser Mode Structure on Stimulated Brillouin Scattering," *IEEE J. Quant. Electron.* **QE-22**:2161 (1986).
139. G. C. Valley, "A Review of Stimulated Brillouin Scattering Excited with a Broad Pump Laser," *IEEE J. Quant. Electron.* **QE-22**:704 (1986).
140. W. T. Whitney, M. T. Duignan, and B. J. Feldman, "Stimulated Brillouin Scattering and Phase Conjugation of Multiline Hydrogen Fluoride Laser Radiation," *JOSA B* **7**:2160 (1990).
- 140a. Robert A. Fisher (ed.), *Optical Phase Conjugation*, Academic Press, New York, 1983.
141. A. I. Erokhin, V. I. Kovalev, and F. S. Faizullov, "Determination of the Parameters of a Nonlinear Response of Liquids in an Acoustic Resonance Region by the Method of Nondegenerate Four Wave Interaction," *Sov. J. Quant. Electron.* **16**:872 (1986).
142. M. J. Dyer and W. K. Bischel, unpublished data.
143. S. T. Amimoto, R. W. F. Gross, L. Garman-DuVall, T. W. Good, and J. D. Piranian, "Stimulated Brillouin Scattering Properties of SnCl₄," *Opt. Lett.* **16**:1382 (1991).
144. M. J. Damzen, M. H. R. Hutchinson, and W. A. Schroeder, "Direct Measurement of the Acoustic Decay Times of Hypersonic Waves Generated by SBS," *IEEE J. Quant. Electron.* **QE-23**:328 (1987).
145. F. E. Hovis and J. D. Kelley, "Phase Conjugation by Stimulated Brillouin Scattering in CClF₃ near the Gas-Liquid Critical Temperature," *JOSA B* **6**:840 (1989).
146. S. Y. Tang, C. Y. She, and S. A. Lee, "Continuous Wave Rayleigh-Brillouin Gain Spectroscopy in SF₆," *Opt. Lett.* **12** (1987).
147. D. Milam, LLNL Report No. 90-011/6330K.
148. P. Klovekom and J. Munch, "Variable Stimulated Brillouin Scattering Pulse Compressor for Nonlinear Optical Measurements," *Appl. Opt.* **36**:5913 (1997).
149. N. G. Basov, I. G. Zubarev, A. V. Kotov, S. I. Mikhailov, and M. G. Smirnov, "Small-Signal Wavefront Reversal in Non-Threshold Reflection from a Brillouin Mirror," *Sov. J. Quant. Electron.* **9**:237–239 (1979).

150. N. F. Andreev, V. I. Bespalov, A. M. Kiselev, A. Z. Matreev, G. A. Pasmanik, and A. A. Shilov, "Wavefront Inversion of Weak Optical Signals with a Large Reflection Coefficient," *JETP Lett.* **32**:625–629 (1980).
151. A. M. Scott and K. D. Ridley, "A Review of Brillouin-Enhanced Four-Wave Mixing," *IEEE J. Quant. Electron.* **25**:438–459 (1989).

15.5 ADDITIONAL REFERENCES

- S. A. Akmanov, B. V. Zhdanov, A. I. Kovrigin, and S. A. Pershin, "Effective Stimulated Scattering in the Ultraviolet and Dispersion of Gain in the 1.06–0.26 μ Band," *JETP Lett.* **15**:185 (1972).
- F. Ausenegg and V. Deserno, "Stimulated Raman Scattering Excited by Light of 5300 \AA ," *Opt. Commun.* **2**:295 (1970).
- G. Bisson and G. Mayer, "Effets Raman stimulés dans la calcite," *Crit. Acad. Set Paris* **265**:397 (1967).
- N. Bloembergen, B. P. Lallemand, A. Pine, and P. Simova, "Controlled Stimulated Raman Amplification and Oscillation in Hydrogen Gas," *IEEE J. Quant. Electron.* **QE-3**:197 (1967).
- R. L. Byer, "A 16- μm Source for Laser Isotope Enrichment," *IEEE J. Quant. Electron.* **QE-12**:732 (1976).
- M. J. Colles, "Efficient Stimulated Raman Scattering from Picosecond Pulses," *Opt. Commun.* **1**:169 (1969).
- J. Gazengel, N. P. Xuan, and G. Rivoire, "Stimulated Raman Scattering Thresholds for Ultra-Short Excitation," *Opt. Acta* **26**:1245 (1979).
- H. Gorner, M. Maier, and W. Kaiser, "Raman Gain in Liquid Core Fibers," *J. Raman Spectrosc.* **2**:363 (1974).
- A. Z. Grasyuk, V. F. Erinkov, I. G. Zubarev, V. I. Mishin, and V. G. Smirnov, "Laser Based on Raman Scattering in Liquid Nitrogen," *JETP Lett.* **8**:291 (1968).
- J. B. Grun, A. K. McQuillan, and B. P. Stoicheff, "Intensity and Gain Measurements on the Stimulated Raman Emission in Liquid O_2 and N_2 ," *Phys. Rev.* **181**:61 (1969).
- E. E. Hagenlocker, R. W. Minck, and W. G. Rado, "Effects of Phonon Lifetime on Stimulated Optical Scattering in Gases," *Phys. Rev.* **154**:226 (1967).
- E. P. Ippen, "Low-Power Quasi-cw Raman Oscillator," *Appl. Phys. Lett.* **16**:303 (1970).
- R. W. Minck, E. E. Hagenlocker, and W. G. Rado, "Consideration and Evaluation of Factors Influencing the Stimulated Optical Scattering in Gases," Scientific Laboratory, Ford Motor Company, Dearborn, MI, SC66-24, 1966.
- I. Reinhold and M. Maier, "Gain Measurements of Stimulated Raman Scattering Using a Tunable Dye Laser," *Opt. Commun.* **5**:31 (1972).
- W. L. Smith, F. P. Milanovich, and M. Henesian, Lawrence Livermore National Laboratory, Private Communication, 1983.
- M. B. Vakhonev, V. N. Volkov, A. Z. Grasyuk, and A. N. Kirkin, "Determination of the Gain in Stimulated Raman Scattering under Spatially Inhomogeneous Pumping Conditions," *Sov. J. Quant. Electron.* **6**:1369 (1976).

THIRD-ORDER OPTICAL NONLINEARITIES

Mansoor Sheik-Bahae and Michael P. Hasselbeck

*Department of Physics and Astronomy
University of New Mexico
Albuquerque, New Mexico*

16.1 INTRODUCTION

The subject of this chapter could well fill a textbook, and indeed the topic comprises a significant portion of the many books on nonlinear optics. A large (but by no means exhaustive or complete) list of texts that provide extensive discussion of high-order nonlinearities appears in Refs. 1 through 30. We have not attempted to write a review chapter nor mentioned or even listed every known third-order nonlinear optical phenomenon. Our aim is to illustrate important and representative third-order effects, emphasizing qualitative descriptions. Details can be found in the references. An exception is our discussion of the Kramers-Kronig relations in nonlinear optics. A fundamental premise of this transformation is the causal link between nonlinear refraction and nonlinear absorption, which is a key aspect of the third-order susceptibility. It has not been treated in most texts; some of the important mathematical steps are given here. Our treatment of third-order nonlinear optics assumes that the reader is familiar with electromagnetic theory, physical optics, and quantum mechanical energy level diagrams.

Any real, physical oscillating system will exhibit a nonlinear response when it is overdriven. In an optical system, a nonlinear response can occur when there is sufficiently intense illumination. The nonlinearity is exhibited in the polarization ($\bar{\mathcal{P}}$) of the material, which is often represented by a power series expansion of the total applied optical field ($\bar{\mathcal{E}}$):

$$\bar{(\mathcal{P})} = \epsilon_0 \chi^{(1)} \bar{\mathcal{E}} + \epsilon_0 \chi^{(2)} \bar{\mathcal{E}}^2 + \epsilon_0 \chi^{(3)} \bar{\mathcal{E}}^3 + \dots \quad (1)$$

Here $\chi^{(1)}$ is the linear susceptibility representing the linear response (i.e., linear absorption and the refractive index) of the material. The two lowest-order nonlinear responses are accounted for by the second- and third-order nonlinear susceptibilities $\chi^{(2)}$ and $\chi^{(3)}$. The subject of this chapter is third-order effects. Processes arising from the second-order response (including second-harmonic generation and optical parametric processes) are discussed elsewhere (see Chap. 10, “Nonlinear Optics” and Chap. 17, “Continuous-Wave Optical Parametric Oscillators”). We will, however, briefly consider the cascading of second-order nonlinearities that appear as an *effective* third-order process in Sec. 16.10.

Third-order optical nonlinearities cover a vast and diverse area in nonlinear optics. A simple illustration of this point is the reported range of magnitudes and response times for $\chi^{(3)}$ in various materials, which span 15 orders of magnitude! This has led to unavoidable inconsistency and confusion in the definition and interpretation of the nonlinear susceptibility. We will not be immune from such inconsistencies and errors. In that spirit, we note that the simple power series representation of the nonlinear optical response described by Eq. (1) is not rigorously correct because it assumes the response is instantaneous. In the case of the bound electronic nonlinearity, for example, this assumption is excellent because the response is exceedingly fast. The response is not *infinitely* fast, however. Response times can vary by orders of magnitude depending on the physical mechanism and resonance conditions involved. Furthermore, Eq. (1) assumes locality, which implies that the nonlinear polarization at a given point in space depends on the magnitude of the electric field only at that point. This condition is not always satisfied. The electrostrictive nonlinearity, for example, is the result of physical displacement of charged particles in a material subject to the ponderomotive force due to the gradient of light irradiance. It is therefore nonlocal. It is nevertheless instructive to apply Eq. (1) to describe various third-order effects that are local and (single photon) nonresonant.

The nonlinear polarization represented in Eq. (1) excludes *effective* third-order nonlinear processes involving linear absorption ($\chi^{(1)}$ process) of one of the excitation beams. An example is the thermal nonlinearity resulting from linear absorption and heating that causes a change of refractive index. Although this is effectively a third-order nonlinear response, we group this and similar phenomena in Sec. 16.9 on cascaded $\chi^{(1)}\chi^{(1)}$ effects.

The term involving $\bar{\mathcal{E}}^3$ in Eq. (1) implies that three optical fields interact to produce a fourth field. The $\chi^{(3)}$ interaction is thus a four-photon process. This is a consequence of the quantum mechanical picture of the nonlinear susceptibility. Conservation of photon energy is always required to complete the interaction process. Assuming the applied optical fields are monochromatic plane waves, we write the total input electric field $\bar{\mathcal{E}}$ as

$$\bar{\mathcal{E}}(r, t) = \bar{\mathcal{E}}_1(r, t) + \bar{\mathcal{E}}_2(r, t) + \bar{\mathcal{E}}_3(r, t) \quad (2)$$

In general, each beam has a different frequency (ω) and wave vector (\vec{k}), represented in complex notation:

$$\bar{\mathcal{E}}_j(r, t) = \frac{\bar{E}_j}{2} \exp(i\omega_j t - i\vec{k}_j \cdot r) + \text{c.c.} \quad \text{for } j=1, 2, 3 \quad (3)$$

where c.c. stands for complex conjugate and \bar{E}_j is a complex vector describing the amplitude, phase, and polarization of each beam. It is important to realize that there can be up to three different input laser frequencies, but there can also be as few as one. Ignoring the $\chi^{(1)}$ and $\chi^{(2)}$ components in Eq. (1), the nonlinear polarization resulting from the $\bar{\mathcal{E}}^3$ interaction leads to a total of 108 terms involving all possible permutations of the fields at three frequencies. The nonlinear polarization occurs at frequencies given by

$$\omega_4 = \pm\omega_i \pm \omega_j \pm \omega_k \quad \text{for } i, j, k=1, 2, 3 \quad (4)$$

The existence of 108 terms does not mean there are as many distinct mechanisms involved. For instance, three terms give $\omega_4 = 3\omega_j$ for $j=1, 2, 3$, describing exactly the same process of third-harmonic generation (THG). Furthermore, THG is a special case of sum frequency generation (SFG) involving one, two, or three different frequencies giving $\omega_4 = \omega_i + \omega_j + \omega_k$, $i, j, k=1, 2, 3$ accounting for 27 terms.

One realizes 108 permutations with different time ordering of three different laser beams distinguished by frequency, and/or wave vector, and/or polarization. If only two distinguishable laser beams are available, the number of permutations decreases to 48. When the system is driven by a single beam, the third-order response involves only four terms in three fields. In general, the $\chi^{(3)}$ coefficients associated with each term will be different due to the ever-present dispersion (i.e., frequency dependence) of the susceptibilities. The frequency dependence is a direct consequence of the finite response time of the interaction. We will expand on this subject in the discussion of the Kramers-Kronig dispersion relations in Sec. 16.4.

Another important property of nonlinear susceptibilities is their tensor nature. Because of the molecular or lattice structure of materials, the nonlinear response will depend on the state of

polarization of the optical fields. For the sake of brevity, we neglect the tensor properties of $\chi^{(3)}$ and treat all the susceptibilities and electric fields as scalar quantities. The reader may consult textbooks on nonlinear optics for detailed discussions of this subject.

Propagation of interacting beams is also an important consideration, and one must account for wave vector summation (i.e., conservation of momentum) that results from the $\bar{\epsilon}^3$ operation. It is useful to invoke the four photon picture, recalling that the momentum of each photon is given by $\hbar k_j$. Taking the resultant nonlinear polarization to be a plane wave with a wave vector k_4 , momentum conservation requires that:

$$\bar{k}_4 = \pm \bar{k}_i \pm \bar{k}_j \pm \bar{k}_k \quad (5)$$

where $|\bar{k}_j| = n(\omega_j)\omega_j/c$. This phase-matching requirement is not necessarily satisfied in every interaction due to dispersion of the linear refractive index in the material. Phase matching can be a serious obstacle in interactions leading to new-frequency generation, that is, when $\omega_4 \neq \omega_1, \omega_2$ and ω_3 (e.g., Sec. 16.6 on THG). When the nonlinear polarization is at one of the driving frequencies, $\omega_4 = \omega_i$ for example, conservation of energy [Eq. (4)] implies that $\omega_j = -\omega_k$. In this case, Eq. (5) reduces to a vector-matching condition that depends only on the geometry (i.e., direction) of the beams (Secs. 16.11 and 16.12).

The frequency terms arising from the third-order nonlinear polarization described by Eq. (1) are collected in Table 1. In the following section, we discuss the physical mechanisms and important features of these processes.

TABLE 1 Frequency Terms Arising from Third-Order Nonlinear Polarization

Nonlinear Process	$8P(\omega_4)/\epsilon_0\chi^{(3)}$	ω_4
• Third-harmonic generation (THG)	$E_j^3, j=1, \dots, 3$	$3\omega_1, 3\omega_2, 3\omega_3$
• Sum frequency generation (SFG)	$3E_i E_j^2, i, j=1, 2, 3, i \neq j$	$2\omega_1 + \omega_2, 2\omega_1 + \omega_3,$ $2\omega_2 + \omega_3, 2\omega_2 + \omega_1,$ $2\omega_3 + \omega_1, 2\omega_3 + \omega_2,$
• Frequency mixing	$6E_1 E_2 E_3$	$\omega_1 + \omega_2 + \omega_3$
• Parametric amplification	$3E_i^* E_j^2, i, j=1, 2, 3, i \neq j$	$2\omega_1 - \omega_2, 2\omega_1 - \omega_3,$ $2\omega_2 - \omega_3, 2\omega_2 - \omega_1,$ $2\omega_3 - \omega_1, 2\omega_3 - \omega_2,$
• Coherent Stokes and anti-Stokes Raman scattering (CSRS and CARS)	$6E_i^* E_j E_k, i, j, k=1, 2, 3, i \neq j \neq k$	$\omega_1 + \omega_2 - \omega_3,$ $\omega_1 - \omega_2 + \omega_3,$ $-\omega_1 + \omega_2 + \omega_3$
	$6E_i^* E_j^* E_k, i, j, k=1, 2, 3, i \neq j \neq k$	$\omega_1 - \omega_2 - \omega_3,$ $-\omega_1 - \omega_2 + \omega_3,$ $-\omega_1 + \omega_2 - \omega_3$
• Bound electronic optical Kerr effect	$3E_i^2 E_i^*, i=1, 2, 3$	$\omega_1, \omega_2, \omega_3$
• Raman-induced Kerr effect (RIKE)		
• Molecular orientational Kerr effect	$6E_i E_j E_j^*, i, j=1, 2, 3, i \neq j$	
• Two-photon absorption (2PA)		
• ac Stark effect		
• Stimulated Raman scattering (SRS)		
• Stimulated Rayleigh-Wing Scattering		

16.2 QUANTUM MECHANICAL PICTURE

The conservation of energy shown in the frequency summation of Eq. (4) contains both positive and negative signs from each of the input beams. A positive sign represents the annihilation (absorption) of a photon, while the negative sign is interpreted as the generation (gain) of a photon. Both annihilation and generation of photons involve atomic and/or molecular transitions from one state to another. It is instructive to use diagrams to keep track of transitions participating in the nonlinear interactions. Let us take the most general case, where four distinguishable beams (i.e., three input photons at $\omega_1, \omega_2, \omega_3$, and the final photon at ω_4) are involved. The third-order nonlinear interaction follows a path corresponding to one of 14 sign-ordering possibilities, assuming emission at the photon frequency ω_4 . All possible time-ordering sequences are illustrated in Fig. 1. In addition, the interacting photons are in general distinguishable; to preserve the clarity of presentation we have not shown this in Fig. 1. Because the photons are (in general) distinct, we must allow for permutation of frequencies in the diagrams. Assuming emission at ω_4 (i.e., we can only assign ω_4 to a downward-pointing arrow), we count up the various time-ordering permutations for each interaction path shown in Fig. 1. This gives a total of 168 terms! Little would be gained by a tedious analysis of all these terms, and such a task is far beyond the scope of this chapter. Instead, we illustrate some important third-order mechanisms and the role of resonances in Fig. 2, where we have labeled the three most important diagrams in Fig. 1. The energy level $|g\rangle$ is the ground state, while $|a\rangle, |u\rangle$, and $|b\rangle$ are intermediate states of the system in a sequence of transitions involving photons with frequencies $\omega_1, \omega_2, \omega_3$, and ω_4 ($i, j, k, l = 1, 2, 3, 4$) such that $\pm\omega_i \pm \omega_j \pm \omega_k \pm \omega_l = 0$. The three time-ordering processes shown in the figure are:

Figure 2a. Consecutive absorption of three photons followed by the generation of the final photon, partly describing sum frequency generation and third-harmonic generation. The reverse process is third-order parametric amplification, which is the absorption of a photon together with emission of three photons.

Figure 2b. An absorption-emission-absorption-emission sequence. Difference frequency generation and frequency mixing are examples of this type of interaction. Coherent anti-Stokes Raman spectroscopy (CARS) is also represented by this transition sequence.

Figure 2c. Absorption of two photons followed by emission of the two photons. As can be seen in the third section of Table 1, a variety of physical mechanisms fall under this general description. Note that the essential difference between (b) and (c) is the time ordering of the transitions. This is extremely important in resonant cases: a Raman-type resonance occurs in (b), and a two-photon resonance exists in (c).

Energy conservation is strictly obeyed upon the completion of the interaction [as dictated by Eq. (4)] but may be violated in the time frame of intermediate state transitions. This is allowed by Heisenberg's Uncertainty Principle. In many cases, an intermediate state is a virtual state, which is a convenient way of stating that a real, intermediate state of the system does not exist to support the transition of a photon at the selected wavelength. The virtual, intermediate state allows for energy bookkeeping in transition diagrams, but a physical description of the optical interaction using quantum mechanics involves only real eigenstates of the system. In particular, there must be a dipole-allowed transition between the initial state $|g\rangle$ and a real state *associated* with the virtual state. The time scale and strength of the interaction is partly determined by the energy mismatch between the virtual, intermediate state and an associated real, electronic state. This means a system can absorb a photon of energy $\hbar\omega_i$ and make a transition from the ground state $|g\rangle$ to a real intermediate state $|a\rangle$ even though there is insufficient photon energy to bridge the gap (i.e., there is an energy mismatch $\Delta E = |\hbar\omega_i - E_a + E_g| > 0$). This is possible, provided the interaction occurs in a time faster than the observation time $\Delta t \sim \hbar/\Delta E$ permitted by the Uncertainty Principle. Transitions of this type are called *virtual* transitions, as opposed to *real* transitions, where energy is conserved. In the former case, Δt is known as the virtual lifetime of the transition.

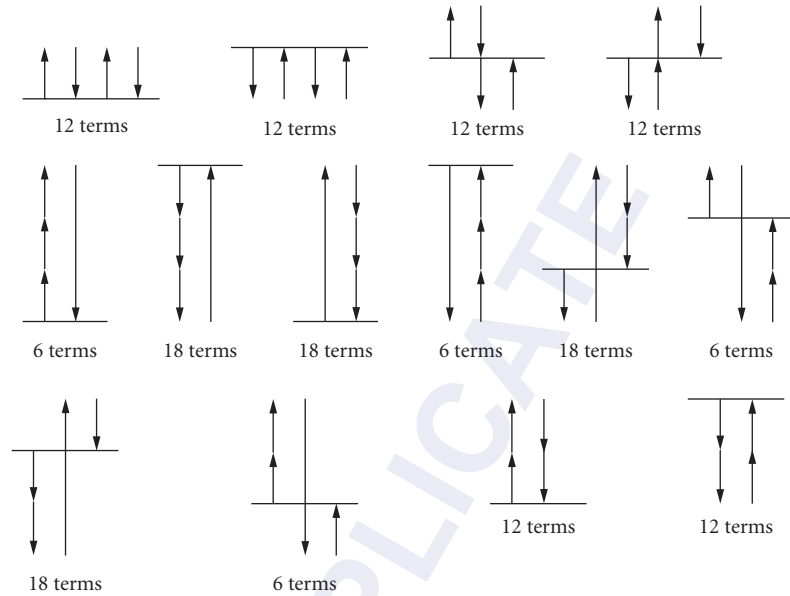


FIGURE 1 Time-ordering sequence illustrating all possible third-order paths. Arrows depict photons. Note that in general, the size of the arrows can be different provided their vector sum is zero. The number of terms is obtained assuming emission of a photon at ω_4 . For clarity, arrows are not marked.

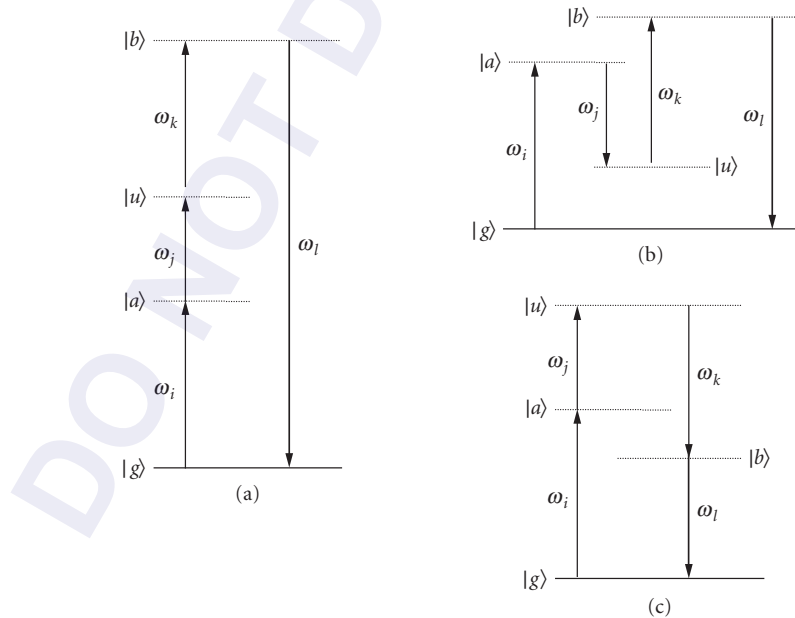


FIGURE 2 Energy level diagrams for some important third-order nonlinear optical processes: (a) third-harmonic generation (THG); (b) coherent anti-Stokes Raman scattering (CARS); and (c) two-photon absorption (2PA).

If the entire sequence of transitions comprising the third-order interaction is not completed within the virtual lifetime, the intermediate state collapses back to the ground state, and no nonlinear interaction occurs. In other words, all the required particles must be present in the system during the virtual lifetime. The longer the virtual lifetime, the greater the probability that the required photons will appear, allowing the multiparticle interaction to run to completion. A longer virtual lifetime translates to a larger third-order nonlinear susceptibility $\chi^{(3)}$. The closer an input photon moves to a dipole-allowed system resonance, the longer the virtual lifetime and the stronger the resulting $\chi^{(3)}$ will be.

These quantum mechanical issues are manifest in the mathematical formulation of $\chi^{(3)}$ derived from perturbation theory:^{1,3,17}

$$\chi^{(3)}(\pm\omega_1, \pm\omega_2, \pm\omega_3) = \frac{N}{\hbar^3} \sum_{i,j,k,l} \sum_{a,u,b} \mu_{ga} \frac{\mu_{au}}{(\omega_{ag} \mp \omega_i)} \frac{\mu_{ub}}{(\omega_{ug} \mp \omega_j \mp \omega_i)} \frac{\mu_{bg}}{\underbrace{(\omega_{bg} \mp \omega_i \mp \omega_j \mp \omega_k)}_{\mp\omega_l}} \quad (6)$$

In Eq. (6), N is the total population in the ground state $|g\rangle$, and μ 's are dipole-moment matrix elements associated with each of the transitions. The first sum describes the frequency permutations: i, j, k , and l can take any integer value 1, 2, 3, 4, provided energy conservation ($\pm\omega_i \pm \omega_j \pm \omega_k \pm \omega_l = 0$) is obeyed. The second sum is over all possible real, intermediate quantum eigenstates of the system. This complicated-looking equation is nothing more than the sequence of optical transitions weighted by the appropriate virtual lifetime. The first coefficient represents the virtual transition initiated by a photon of energy $\hbar\omega_i$ from the ground state to the intermediate state $|a\rangle$ with strength given by the matrix element μ_{ga} . The next three matrix elements are weighted by the virtual lifetimes of their initial state. The virtual lifetimes are represented by energy (i.e., frequency) denominators; as the photon frequency approaches a system resonance, the virtual lifetime and magnitude of $\chi^{(3)}$ grow accordingly.

The \pm signs in front of the frequency arguments in Eq. (6) indicate there is a physical significance to the time ordering of the participating photons. This representation distinguishes the various components of the third-order susceptibility. In many textbooks, a permutation of all the frequencies (including the signs) is already incorporated in the final calculation of $\chi^{(3)}$.^{3,4} In that case, for a given ω_i , one obtains the total contribution to $\chi^{(3)}$ with the order of frequency arguments having no particular physical relevance.

We also point out that the nonlinear susceptibility described by Eq. (6) and shown in our example is *real*. A resonance condition occurs when any one of the energy/frequency denominators approaches zero. This not only enhances $\chi^{(3)}$ but also makes it a complex quantity (i.e., a resonance condition introduces an imaginary component to $\chi^{(3)}$). This is better understood by making the following substitution:

$$\frac{1}{\Delta\omega} \rightarrow \frac{1}{\Delta\omega + i\Gamma} \quad (7)$$

where Γ represents a phenomenological broadening of the particular transition. This complex damping term accounts for the physical impossibility of the nonlinear susceptibility becoming infinite in a resonance condition. Even in the case of vanishing damping, a basic theorem of complex variables can be applied to Eq. (7):

$$\lim_{\Gamma \rightarrow 0} \frac{1}{\Delta\omega + i\Gamma} = \mathcal{P} \left(\frac{1}{\Delta\omega} \right) + i\pi\delta(\Delta\omega) \quad (8)$$

where \mathcal{P} stands for the principle value and δ is the Dirac delta function. The important message is that in general, the nonlinear susceptibility $\chi^{(3)}$ is a complex quantity that will be dominated by its imaginary component when photon frequencies move into resonance with real eigenstates of the system.

The resonance conditions leading to a strong imaginary $\chi^{(3)}$ are associated with one or more of the following processes: three-photon ($\omega_i + \omega_j + \omega_k = -\omega_l \approx \omega_{bg}$), two-photon ($\omega_i + \omega_j = -\omega_k - \omega_l \approx \omega_{ug}$), Raman-type ($\omega_i - \omega_j = \omega_l - \omega_k \approx \omega_{ug}$), and/or single-photon ($\omega_i \approx \omega_{ag}$) resonances. The latter cases (i.e., those having linear resonance) will be discussed in Sec. 16.9, which deals with cascaded $\chi^{(1)}$ - $\chi^{(1)}$ nonlinearities. A special case of linear resonance can occur in Raman-type transitions, where $\omega_{ug} = 0$ (i.e., when the second intermediate state is degenerate with the ground state). This corresponds to the optical Stark effect (ac Stark effect). The three-photon resonance that gives rise to an imaginary $\chi^{(3)}$ in third-harmonic generation does not have a significant physical implication. It only influences the phase of the interacting fields, similar to the case of second-order effects (e.g., second-harmonic generation).¹ The remaining two processes involving two-photon and Raman resonances are of significant interest and will be discussed in detail.

16.3 NONLINEAR ABSORPTION AND NONLINEAR REFRACTION

Just as the real and imaginary components of the linear susceptibility $\chi^{(1)}$ are associated with refraction and absorption, the real and imaginary parts of $\chi^{(3)}$ describe nonlinear refraction (NLR) and nonlinear absorption (NLA) or gain. This can be understood by considering situations in which the nonlinear polarization is at one of the driving frequencies. These are particular cases of Fig. 2b and c, with corresponding polarization terms given in the third section of Table 1. Taking the interacting photons to have frequencies ω_a and ω_b , the total polarization (linear and third order) at ω_a can be written as

$$P(\omega_a) = \epsilon_0 \left\{ \frac{1}{2} \chi^{(1)}(\omega_a) E_a + \frac{3}{8} \chi^{(3)}(\omega_a, \omega_a, -\omega_a) E_a^2 E_a^* + \frac{6}{8} \chi^{(3)}(\omega_a, \omega_b, -\omega_b) E_a E_b E_b^* \right\} \quad (9)$$

For the sake of brevity, we ignore time ordering in the frequency arguments of $\chi^{(3)}$. This means that the $\chi^{(3)}$ component in Eq. (9) is assumed to contain the various permutations of frequencies including, for example, two-photon as well as Raman transitions shown in Fig. 2b and c. From Eq. (9), we introduce an effective susceptibility χ_{eff} :

$$\chi_{\text{eff}}(\omega_a) = \chi^{(1)}(\omega_a) + \frac{3}{4} \chi^{(3)}(\omega_a, \omega_a, -\omega_a) |E_a|^2 + \frac{6}{4} \chi^{(3)}(\omega_a, \omega_b, -\omega_b) |E_b|^2 \quad (10)$$

Deriving the coefficients of nonlinear absorption and refraction from Eq. (10) is now straightforward. The complex refractive index is defined as

$$n + i\kappa = (1 + \chi_{\text{eff}})^{1/2} \quad (11)$$

Making the very realistic assumption that the nonlinear terms in Eq. (10) are small compared to the linear terms, we use the binomial expansion to simplify Eq. (11):

$$n + i\kappa \cong n_0 + i \frac{c}{2\omega_a} \alpha_0 + \Delta n + i \frac{c}{2\omega_a} \Delta \alpha \quad (12)$$

where $n_0 = (1 + \Re\{\chi^{(1)}\})^{1/2}$. We also assume the background linear absorption coefficient is small, that is, $\alpha_0 \propto \Im\{\chi^{(1)}\} \ll \Re\{\chi^{(1)}\}$. We define the irradiance as $I_i = (1/2)c\epsilon_0 n_0(\omega_i) |E_i|^2$ ($i = a, b$) and the nonlinear refraction coefficient n_2 and the nonlinear absorption coefficient α_2 as follows:

$$n_2(\omega_a; \omega_b) = \frac{3}{4\epsilon_0 n_0(\omega_a) n_0(\omega_b) c} \Re\{\chi^{(3)}(\omega_a, -\omega_b, \omega_b)\} \quad (13)$$

$$\alpha_2(\omega_a; \omega_b) = \frac{3\omega_a}{4\epsilon_0 n_0(\omega_a) n_0(\omega_b) c^2} \Im\{\chi^{(3)}(\omega_a, -\omega_b, \omega_b)\} \quad (14)$$

The change of refractive index due to the presence of fields E_a and E_b is

$$\Delta n(\omega_a) = n_2(\omega_a; \omega_a) I_a + 2n_2(\omega_a; \omega_b) I_b \quad (15)$$

and the change of absorption is

$$\Delta \alpha(\omega_a) = \alpha_2(\omega_a; \omega_a) I_a + 2\alpha_2(\omega_a; \omega_b) I_b \quad (16)$$

where I_a and I_b are the irradiances of the two beams. Note that without loss of generality we assume the measurement is performed on the laser beam corresponding to field E_a , while the field E_b acts as an excitation source only. The first terms on the right-hand side of the just-noted equations correspond to self-action (i.e., commonly performed single-beam experiments). The second terms correspond to the case of an excite-probe experiment where the two beams are distinguishable either by frequency and/or wave vector. The factor of 2 in front of the second terms in Eqs. (15) and (16) arises from the larger number of permutations in this component of the nonlinear susceptibility.^{3,4,31} The stronger nondegenerate response (i.e., distinguishable beams) is sometimes referred to as *weak-wave retardation*.³² While most reported measurements and applications involve degenerate self-action processes (i.e., a single laser beam), the theoretical treatment presented in this chapter considers the more general nondegenerate case. One must keep in mind that degenerate third-order coefficients are only the limit of the nondegenerate case, where $\omega_a = \omega_b$. The need for generality in the theoretical approach is very important for correct implementation of the Kramers-Kronig dispersion relations in nonlinear optics. This allows us to establish a rigorous mathematical relation between NLR and NLA, discussed in the next section.

Another commonly used coefficient for describing the nonlinear index is \tilde{n}_2 defined as

$$n = n_0 + \tilde{n}_2(\omega_a; \omega_a) \frac{|E_a|^2}{2} + 2\tilde{n}_2(\omega_a; \omega_b) \frac{|E_b|^2}{2} \quad (17)$$

where \tilde{n}_2 is usually given in Gaussian units (esu). \tilde{n}_2 is related to n_2 by

$$\tilde{n}_2(\text{esu}) = \frac{cn}{40\pi} n_2(\text{SI}) \quad (18)$$

where the right-hand side is in SI/MKS units. The reader is cautioned that in the literature various symbols and definitions different from those given here are often used to describe the nonlinear refractive index. The symbol β is commonly used in place of α_2 to denote two-photon absorption (2PA).

The propagation of electromagnetic waves E_a and E_b through a nonlinear medium, ignoring the effect of diffraction and dispersion (i.e., pulse distortion), is governed by the following equations for the irradiance and phase of the probe beam (E_a):

$$\frac{dI_a}{dz} = -\alpha_0(\omega_a)I_a - \alpha_2(\omega_a; \omega_a)I_a^2 - 2\alpha_2(\omega_a; \omega_b)I_aI_b \quad (19)$$

and

$$\frac{d\phi_a}{dz} = \frac{\omega_a}{c} [n_0(\omega_a) + n_2(\omega_a; \omega_a)I_a + 2n_2(\omega_a; \omega_b)I_b] \quad (20)$$

The coefficient n_2 is often used to describe the nonlinear index change due to mechanisms such as thermally induced material changes, molecular orientation effects, saturation of absorption, and ultra-fast $\chi^{(3)}$ processes. Here, consistent with our definition of $\chi^{(3)}$, we designate the n_2 notation for local and linearly nonresonant nonlinearities only. Processes that appear as an *effective* n_2 are treated separately as cascaded $\chi^{(1)}; \chi^{(1)}$ or $\chi^{(2)}; \chi^{(2)}$ phenomena.

As a consequence of the principle of causality, the real and imaginary parts of the linear susceptibility are connected through the Kramers-Kronig relations of linear optics. Equations (19) and (20) suggest a similar relation in nonlinear optics. We discuss the Kramers-Kronig relations of nonlinear optics and their underlying physics next.

16.4 KRAMERS-KRONIG DISPERSION RELATIONS

The complex response function of any linear, causal system obeys a dispersion relation linking its real and imaginary parts as Hilbert transform pairs. In linear optics, causality is manifest in the Kramers-Kronig (K-K) dispersion relations (see Chap. 5, “Optical Properties of Semiconductors”) that tie the frequency-dependent refractive index, $n(\omega)$, to the absorption coefficient $\alpha(\omega)$ and vice versa:

$$n(\omega) - 1 = \frac{c}{\pi} \mathcal{P} \int_0^\infty \frac{\alpha(\omega')}{\omega'^2 - \omega^2} d\omega' \quad (21)$$

where \mathcal{P} denotes the Cauchy principal value. The principal value is really just a warning to be careful when integrating near the singularity in the denominator of the integrand. We drop the \mathcal{P} notation for simplicity, although it is always implied. There is an equivalent relation for the real and imaginary parts of the linear susceptibility:

$$\Re\{\chi^{(1)}(\omega)\} = \frac{1}{\pi} \int_{-\infty}^\infty \frac{\Im\{\chi^{(1)}(\omega')\}}{\omega' - \omega} d\omega' \quad (22)$$

The K-K relation is the mathematical expression of causality, and a simple, intuitive derivation of these relations can be made.^{33,34}

Causality clearly holds for any real, linear system. Real, nonlinear systems must also be causal—does that imply there are dispersion relations as well? If so, what form do they take? The Kramers-Kronig relations of linear optics are derived from linear dispersion theory, suggesting this procedure is completely inappropriate for a nonlinear system. Fortunately, this is not the case, and since the early days of nonlinear optics, many authors have addressed the K-K relations in the nonlinear regime.³⁵⁻³⁹

The usefulness of these relations was not fully appreciated until recently, however.^{31,39,40} The key insight is that one can linearize the system; we view it as the material *plus* a strong perturbing light beam. This new linear system, which is different from the system in the presence of weak light, has a modified absorption spectrum. The linear Kramers-Kronig relation is applied in the presence of and in the absence of a high field perturbation, and we study the difference between these two regimes. It is important to appreciate the fact that our new system is causal even in the presence of an external perturbation. This allows us to write down a modified form of the Kramers-Kronig relation linking the index of refraction to the absorption^{31,39}:

$$[n(\omega) + \Delta n(\omega; \zeta)] - 1 = \frac{c}{\pi} \int_{-\infty}^{\infty} \frac{\alpha(\omega') + \Delta\alpha(\omega'; \zeta)}{\omega'^2 - \omega^2} d\omega' \quad (23)$$

After subtracting the linear terms n and α , we are left with a relationship between the changes in refractive index and change of absorption:

$$\Delta n(\omega; \zeta) = \frac{c}{\pi} \int_0^{\infty} \frac{\Delta\alpha(\omega'; \zeta)}{\omega'^2 - \omega^2} d\omega' \quad (24)$$

where ζ denotes the perturbation. An equivalent relation also exists that allows calculation of the change in absorption coefficient, given the change in the refractive index. This relation is rarely used in practice for reasons described momentarily. In evaluating Eq. (24), it is essential that the perturbation be independent of the frequency of observation (ω'). In other words, the excitation ζ must remain constant as ω' is varied.

It is an interesting fact that calculation of the refractive index change from data obtained in nonlinear optics experiments is often easier than extracting the absolute refractive index from the K-K transform in linear optics! The reason is that absorption changes in nonlinear optics are usually strong only in a very limited frequency range; the integration range in Eq. (24) need only consider this spectrum. In contrast, evaluation of the linear index based on the linear absorption spectrum normally involves a much larger amount of data. One must take full account of the entire linear absorption data to obtain quantitative agreement with experiments that measure the refractive index. In the same way, the reverse transformation in nonlinear optics (obtaining $\Delta\alpha$ from Δn) is not as accommodating as the transformation of Eq. (24). Experiments show that irradiance-dependent changes of the refractive index occur over a relatively broad frequency spectrum. A large and impractical amount of nonlinear dispersion data must be collected and incorporated into a K-K calculation of nonlinear absorption. The reverse transformation is thus difficult to accomplish in practice.

Equation (24) has been used to determine refractive changes due to *real* excitations (i.e., $\chi^{(1)}; \chi^{(1)}$ cascaded processes) such as thermal and free-carrier nonlinearities in semiconductors.^{41,42} In these examples, ζ denotes either change of temperature or change of free-carrier density, respectively. This K-K methodology has also been applied with great success to the situation where the perturbation is *virtual* or nonresonant. This work unified the bound electronic Kerr effect in bulk semiconductors (i.e., the dispersive nonlinearity resulting from anharmonic motion of bound, valence electrons) to its absorptive counterparts: two-photon absorption, the electronic Raman effect, and the ac Stark effect.^{31,39,43-45} The dispersion relation between α_2 and n_2 is given by

$$n_2(\omega_a; \omega_b) = \frac{c}{\pi} \int_0^{\infty} \frac{\alpha_2(\omega'; \omega_b)}{\omega'^2 - \omega_a^2} d\omega' \quad (25)$$

Note that in the general case we are dealing with two frequencies: ω_a and ω_b . Even in the degenerate situation (i.e., $n_2 = n_2(\omega_a; \omega_a)$) where we desire the nonlinear index coefficient at a single frequency ω_a , we are still required to provide the nondegenerate absorption spectrum $\alpha_2(\omega'; \omega_b)$ at

all frequencies ω' as input to the calculation. We also point out that Eq. (25) can be used to derive relations linking the real and imaginary parts of the nonlinear susceptibility via Eq. (13) [the inverse transformation is obtained with Eq. (14)]:

$$\Re\{\chi^{(3)}(\omega_a, \omega_b, -\omega_b)\} = \frac{2}{\pi} \int_0^{\infty} \frac{\omega' \Im\{\chi^{(3)}(\omega', \omega_b, \omega_b)\}}{\omega'^2 - \omega_a^2} d\omega' \quad (26)$$

Using the symmetry properties of $\chi^{(3)}$, an equivalent representation is

$$\Re\{\chi^{(3)}(\omega_a, \omega_b, -\omega_b)\} = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{\Im\{\chi^{(3)}(\omega', \omega_b, \omega_b)\}}{\omega' - \omega_b} d\omega' \quad (27)$$

Equation (27) can also be derived in a very general way from a first-principles approach that applies the causality condition directly in the temporal nonlinear response. In this way, one can obtain the dispersion relations for the n th order optical susceptibility:³⁹

$$\chi^{(n)}(\omega_1, \omega_2, \dots, \omega_j, \dots, \omega_n) = \frac{-i}{\pi} \int_{-\infty}^{\infty} \frac{\chi^{(n)}(\omega_1, \omega_2, \dots, \omega', \dots, \omega_n)}{\omega_j - \omega'} d\omega' \quad (28)$$

where $\chi^{(n)}$ is a complex quantity. By separating the real and imaginary parts of this equation, we get the generalized Kramers-Kronig relation for a nondegenerate, n th order nonlinear susceptibility:

$$\Re\chi^{(n)}(\omega_1, \omega_2, \dots, \omega_j, \dots, \omega_n) = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{\Im\chi^{(n)}(\omega_1, \omega_2, \dots, \omega', \dots, \omega_n)}{\omega' - \omega_j} d\omega' \quad (29)$$

Note that for $n = 3$ with the substitutions $\omega_1 = \omega_a$, $\omega_2 = \omega_b$, and $\omega_3 = -\omega_b$, Eq. (29) becomes identical to Eq. (27). When the susceptibilities for generating frequency harmonics are included, Eq. (29) can be further generalized.³⁹ If we consider the $\chi^{(3)}$ associated with third-harmonic generation, it can be shown that:

$$\Re\{\chi^{(3)}(+\omega, +\omega, +\omega)\} = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{\Im\{\chi^{(3)}(+\omega', +\omega', +\omega')\}}{\omega' - \omega} d\omega' \quad (30)$$

Such relationships have been utilized in calculations of the total $\chi^{(3)}$ (THG) in semiconductors. It is computationally more convenient to first calculate the imaginary part because of the presence of δ functions in its frequency domain. The real part is then calculated using the K-K dispersion relations.⁴⁶

16.5 OPTICAL KERR EFFECT

The $\chi^{(3)}$ process leading to an intensity-dependent refractive index is known as the *optical Kerr effect* (OKE). Experimental observation is relatively straightforward, usually requiring just a single laser beam. The OKE is described by Eq. (15) where, for the sake of brevity, we ignore cross-modulation terms and drop the frequency terms to write

$$n = n_0 + n_2 I \quad (31)$$

There are a variety of physical mechanisms that submit to this mathematical representation, many different ways to observe the effect experimentally, and an assortment of practical devices that can be built.

Optical transitions giving rise to a nonlinear susceptibility $\chi^{(3)}$ (see Figs. 1 and 2, for example) are intimately related to the energy eigenstates of the system. These eigenstates can be associated with bound electronic motion, molecular vibrations, or molecular rotations of the system. Electronic transitions involve the largest energy separation and rotational transitions the smallest. In a given material (gas, liquid, or solid), one or more of these excitations may contribute to the optical Kerr effect. In general, the various contributions differ in their response time, magnitude, and frequency dependence. Referring to our earlier discussion of Eq. (6), the time response of an optical nonlinearity is governed by the virtual lifetime of the relevant transitions. This implies that the nonlinearities associated with the electronic transitions give the fastest response time because they possess large energy denominators. Experiments have shown that electronic nonlinearities are usually faster than the time resolution provided by the shortest optical pulses available at the time of this writing (<10 fs). For practical purposes, then, the nonlinearity associated with the motion of bound electrons can be regarded as instantaneous. At the other extreme, the nonlinearity associated with rotational motion of molecules is relatively sluggish—response times in the picosecond regime have been measured. In the middle range are nonlinearities arising from molecular vibrations. For a Raman-type transition as shown in Fig. 2*b* this effect is manifest as the Raman-induced Kerr effect (RIKE). We discuss two important cases of NLR: the bound electronic Kerr effect in solids (n_2), and the rotational (or orientational) Kerr effect in liquids.

Bound Electronic Optical Kerr Effect in Solids

The optical Kerr effect in solids has been extensively studied in materials ranging from large-gap dielectrics and glasses to narrow-gap semiconductors.^{31,44,47–50} The fundamental energy gap E_g turns out to be a parameter of critical importance. Because of direct, linear absorption of the incident laser light, we are interested in the transparency regime where the photon frequency is less than the bandgap energy (i.e., $\omega < \omega_g = E_g / \hbar$). We estimate the response time (τ_r) of the nonlinearity using the virtual lifetime of the transition: $\tau_r \approx 1/|\omega - \omega_g|$. Far below the bandgap, where $\omega \ll \omega_g$, the response time can be very fast ($\ll 10$ fs) and for most applications can be assumed as instantaneous. This ultrafast response has been exploited in soliton propagation in glass fibers^{5,51} (Sec. 16.11) and in the generation of femtosecond pulses in solid-state lasers.⁵² The optical Kerr effect also causes self-focusing, sometimes resulting in beam distortion and damage in transparent media (Sec. 16.11). Another significant application is the development of ultrafast all-optical-switching devices.^{53,54} Although much progress has been made in this area of research, development of a practical switch has been hindered by the relatively small magnitude of bound electronic nonlinearities.

The bound electronic optical Kerr effect in optical solids has been analyzed using semiempirical methods⁵⁵ and, more recently, by simple two-band models appropriate for semiconductors.^{31,44} The latter treatment provides information about the dispersion, bandgap scaling, and the relationship between the NLA (e.g., two-photon absorption) and NLR through Kramers-Kronig transformation. The resulting simple formula allows one to predict the nonlinear refraction coefficient n_2 knowing only the photon frequency (ω), energy bandgap (E_g), and linear index (n_0):

$$n_2(m^2/W) = \frac{A}{n_0^2 E_g^4} \bar{G}_2 \left(\frac{\hbar\omega}{E_g} \right) \quad (32)$$

where $A \approx 3 \times 10^{-35}$ (MKS); E_g is in joules. The function \bar{G}_2 describes the normalized dispersion of the coefficient n_2 and is depicted in Fig. 3, along with the normalized two-photon absorption spectrum for bulk semiconductors (Sec. 16.8).

Recall that NLA and NLR correspond to the imaginary and real parts of the third-order susceptibility, respectively. The derivation of Eq. (32) by a Kramers-Kronig transformation required knowledge of the NLA spectrum. Three different mechanisms of NLA were employed in the analysis, corresponding to the three relevant time-ordering sequences depicted in Fig. 2: two-photon absorption, electronic

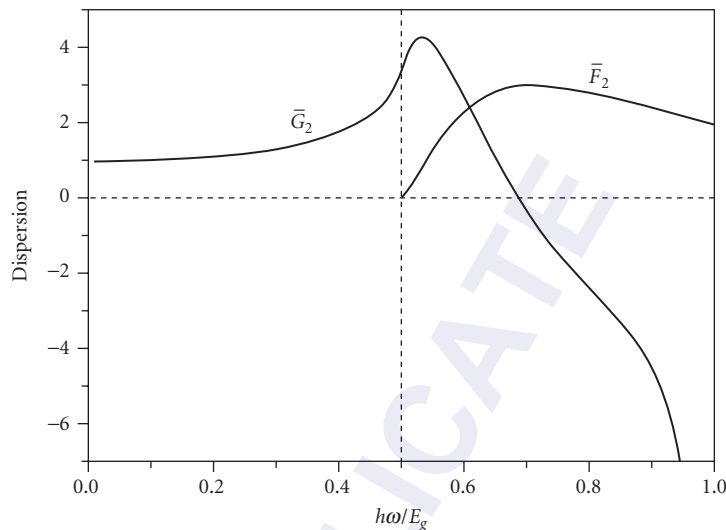


FIGURE 3 Calculated dispersion of nonlinear absorption. Also shown is the two-photon absorption spectrum.

Raman scattering, and the optical (ac) Stark shift of electronic states. The optical Stark effect is a change of the fundamental energy gap that occurs when the oscillating laser field becomes comparable to the electric field binding valence electrons to the positively charge nuclei. In experiments where only a single beam is used (i.e., all the input frequencies are the same), the only observable NLA effect is two-photon absorption, which is discussed in Sec. 16.8.

The plot of the dispersion function \bar{G}_2 (Fig. 3) is consistent with our intuitive arguments about resonance enhancement presented in Sec. 16.2. At long wavelengths (i.e., for $\hbar\omega \ll E_g$), we are far from resonance so $\bar{G}_2 \approx 1$ and is nearly frequency independent (i.e., there is very weak dispersion). Approaching the two-photon resonance (i.e., as $2\hbar\omega$ gets close to E_g), there is an approximately fivefold enhancement of n_2 . With increasing photon energy, the sign of n_2 reverses; this is a direct consequence of the sign change in the energy denominator associated with two-photon absorption. With further photon energy increase, there is enhancement due to resonances ascribed to two-photon absorption and the optical Stark effect. The inverse fourth-power bandgap scaling (i.e., E_g^{-4}) and the dispersion predicted by this simple expression display remarkable agreement with data obtained with many different semiconductors and dielectrics.⁴⁴

Reorientational Kerr Effect in Liquids

The reorientational Kerr effect involves transitions between rotational energy levels of a molecule. It is nonresonant and therefore associated with the real part of $\chi^{(3)}$. The absorptive nonlinearity associated with these rotational levels gives an imaginary component to $\chi^{(3)}$; this is a Raman-type transition known as *Rayleigh-wing scattering*.^{3,17} For simplicity, we make a classical description of this phenomenon. Consider a carbon disulfide (CS_2) molecule as shown in Fig. 4. This is a cigar-shaped molecule with different polarizabilities along its principal axes (here we show $\alpha_3 > \alpha_1$). As discussed in Sec. 16.7, the polarizability describes the propensity for an external field to produce a dipole in a molecule. In the first step of the interaction, the optical field polarizes this molecule (i.e., induces a dipole moment). The induced dipole interacts with the applied field and aligns itself along the direction of polarization. This molecular reorientation (rotation) causes a birefringence in an isotropic

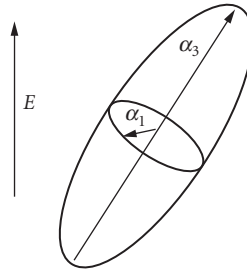


FIGURE 4 Simple picture of CS₂ molecule.

solution; initially, the molecules were randomly oriented and there was no birefringence. The response time of the molecule depends on its mass: the heavier molecule, the slower the response. As an example, CS₂ has a reorientational $n_2 \approx 3.4 \times 10^{-18} \text{ m}^2/\text{W}$ with a relaxation time $\tau \approx 1$ to 2 ps.⁵⁶

16.6 THIRD-HARMONIC GENERATION

In crystals where there is no inversion symmetry, $\chi^{(2)}$ vanishes, making sum and difference frequency mixing impossible. The possibility of third-harmonic generation always exists in principle, although it usually suffers from practical drawbacks. Typical values of $\chi^{(3)}$ are orders of magnitude smaller than $\chi^{(2)}$ coefficients found in popular frequency conversion crystals. This means the laser irradiance must be increased to compensate, often leading to material damage. Moreover, third-harmonic generation in crystals is usually difficult to phase match. Because of these obstacles, cascading of two second-order effects (second-harmonic generation, followed by sum-frequency generation) in two separate crystals is usually the preferred method of obtaining high multiples of the pump laser frequency (see Sec. 16.10).

Gases do not have the damage limitations of crystals. Third-harmonic generation was extensively studied in many different gases around the time that high-power, Q-switched lasers became widely available, and conversion efficiencies as high as several percent were obtained. Studies of sodium vapor have been helpful in elucidating the resonant enhancement that occurs near ω , 2ω , and 3ω .

16.7 STIMULATED SCATTERING

Useful spectroscopic information can be obtained when light is scattered from material, often at frequencies far removed from absorption and emission resonances. Spontaneous scattering is a linear process, in which the material is unmodified by the probing light beam. The various forms of spontaneous scattering (Raman, Brillouin, and Rayleigh) have been known for the better part of a century. In the presence of a sufficiently intense laser beam, however, these scattering processes can be strongly amplified by a nonlinear interaction of the excitation beam with the material, resulting in stimulated scattering.

Raman scattering is most commonly described as the interaction of light with vibrational waves in a material. (Electronic and magnetic excitations can also be measured in Raman experiments.) These vibrational frequencies are typically in the infrared, meaning that Raman-scattered light can have a substantial spectral shift with respect to near-infrared or visible excitation light. Brillouin scattering involves the interaction of light with acoustic waves—waves associated with the propagation of pressure in the medium, leading to periodic density fluctuations. Acoustic waves occur at frequencies that are orders of magnitude smaller than material vibrations; Brillouin scattered light is therefore frequency-shifted from the incident light by a much smaller amount. Rayleigh scattering

results from the interaction of light with stationary density variations—variations much smaller than the wavelength of the incident light. Scattering takes place without any frequency shift relative to the incident light. Rayleigh scattering is of interest because of its strong wavelength dependence and polarization properties. Spontaneous scattering processes scale linearly with input irradiance.

Stimulated Raman Scattering

In Raman scattering (Fig. 5), a photon is absorbed by a material that makes a quantum-mechanical transition from a low energy state $|1\rangle$ to a high energy state $|2\rangle$. At some short time later (i.e., not instantaneously), the material relaxes to a lower energy state $|3\rangle$ different from the original state, giving up its energy in the form of a photon of different energy than the excitation photon. If the lower state $|3\rangle$ is at a higher energy than state $|1\rangle$, the emitted photon will be at a longer wavelength than the excitation light. This is called *Stokes shifted Raman scattering*. If the terminal state $|3\rangle$ is at a lower energy than state $|1\rangle$, the emitted photon will be shorter in wavelength than the incident light, leading to *anti-Stokes shifted Raman scattering*. The difference between the incident and emitted light thus provides information about the relative positions of the different energy levels. Maintaining the same nomenclature, there is also Stokes and anti-Stokes shifted Brillouin scattering. Note that when state $|1\rangle$ and $|3\rangle$ are the same, there is no frequency shift and we have Rayleigh scattering.

The intermediate state can be a real state corresponding to a quantum mechanical energy level of the system; this is known as *resonant Raman scattering*. In the theme of this chapter, resonant Raman scattering is an example of a cascaded linear process leading to an effective $\chi^{(3)}$. More often, the intermediate level is not resonant with the photon, and the transition from $|1\rangle$ to $|2\rangle$ is virtual (illustrated by a horizontal dotted line in Fig. 5). The distinction between the resonant and nonresonant processes can be confusing because both are referred to as Raman scattering. To maintain consistency with standard nomenclature, we briefly depart from the logical organization of this chapter and discuss both resonant and nonresonant Raman scattering in this section.

The essential physics of Raman scattering can be understood from the classical picture of a diatomic molecule of identical atoms vibrating back and forth at frequency ω_L . The diatomic molecule is an illustrative example; in principle all Raman-active and some normal modes of vibration of a solid, liquid, or gas can be probed with Raman techniques.⁵⁷ We assume that the electronic charge distribution on the molecule is perfectly symmetric, hence there is no permanent dipole or a dipole moment modulated by the vibration. This normal mode is therefore not dipole active, that is, it cannot absorb electromagnetic radiation (see Chap. 8, “Fundamental Optical Properties of Solids”).

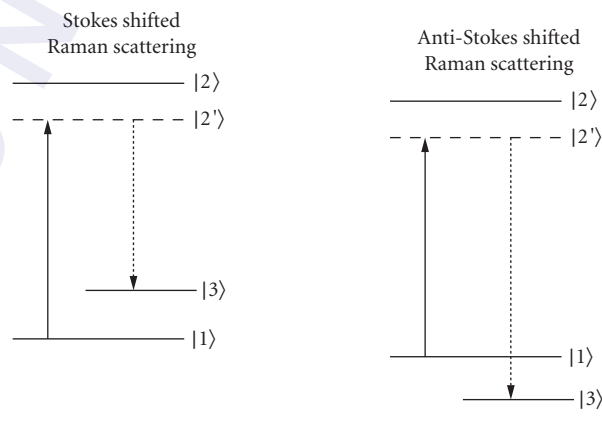


FIGURE 5 Raman scattering.

When an external electric field is applied, the situation changes. The field in an electromagnetic wave polarizes the charge distribution on the molecule and it acquires a dipole. If the induced dipole is also modulated by a normal mode of vibration, the mode is said to be Raman active. The extent to which an external field can polarize the molecule is quantified by the following equation:

$$\mathbf{p}(r, t) = \alpha \mathbf{E}(r, t) \quad (33)$$

where $\mathbf{p}(r, t)$ is the induced dipole moment, α is the polarizability, $\mathbf{E}(r, t)$ is the time- and spatially-varying electric field, and bold type denotes vector quantities. The polarizability is not constant, however, but rather is a function of the molecular separation distance q . Writing the first two terms of a Taylor series expansion of $\alpha(q)$ we have:

$$\alpha(q) = \alpha_0 + \left(\frac{\partial \alpha}{\partial q} \right)_{q_0} q \quad (34)$$

where α_0 is a constant representing the polarizability at the equilibrium position of the molecule (q_0). The molecule vibrates at a frequency $\pm \omega_L$, which is the energy difference between the states $|1\rangle$ and $|3\rangle$ in Fig. 5, hence:

$$q = q_0 \exp(\pm i \omega_L t) \quad (35)$$

Inserting Eqs. (34) and (35) into Eq. (33), and realizing that the electromagnetic field varies sinusoidally at the optical frequency ω , we find that the second term in the polarizability expansion is responsible for the appearance of induced dipoles oscillating at a frequency offset from the incident electromagnetic wave by $\pm \omega_L$:

$$p(r, t)_{\text{Raman}} = E_0(r) \left(\frac{\partial \alpha}{\partial q} \right)_{q_0} q_0 \exp(i \omega t \pm i \omega_L t) \quad (36)$$

These dipoles can radiate and are the origin of spontaneous Raman scattering. There is also an oscillating dipole unaffected by the vibration corresponding to the term α_0 . This dipole oscillation is exactly at the frequency of the incident light and corresponds to spontaneous Rayleigh scattering:

$$p(r, t)_{\text{Rayleigh}} = E_0(r) \alpha_0 \exp(i \omega t) \quad (37)$$

In stimulated scattering, we have to consider the force exerted on the vibrating molecule by the external field as a consequence of its polarizability. This force involves only the second term in Eq. (34):

$$F = \frac{1}{2} \left(\frac{\partial \alpha}{\partial q} \right)_{q_0} \langle E^2(r, t) \rangle \quad (38)$$

where the angular brackets represent a time average over an optical period. In a dipole-active interaction, the lowest order forcing term is proportional to E , resulting in linear absorption of light. In the case of a Raman-active mode, Eq. (38) shows the force scales as E^2 ; therefore, the force is nonlinear in the field. The forcing term is negligible at low light intensities, but becomes important when large electromagnetic field levels generated by lasers are encountered.

Because the Raman active mode of the molecule is subject to a force proportional to E^2 , there must be two input photons driving the interaction. If the two photons are at different frequencies,

the molecule will experience a force at the beat frequency of the two photons. If the wavelengths of the two photons are chosen so that their beat frequency equals that of the molecular vibration ω_L , strong amplification of all three waves (two input electromagnetic waves and the molecular vibration) can occur, resulting in stimulated scattering. The molecular polarizability thus acts as a nonlinear mixing term to reinforce and amplify the interacting waves. It is important to realize this is of practical consequence only when the input electromagnetic fields are sufficiently high.

The nonlinear polarizability impresses sidebands on the pump light, resulting in three distinct electromagnetic waves (laser beam, Stokes shifted Raman, and anti-Stokes shifted Raman) propagating in the medium. The same nonlinear mixing process that led to the generation of the Raman sidebands in the first place can cause coherent excitation of additional molecules due to their polarizability. In this way, a coherent vibrational wave builds up, which in turn feeds more energy into the Raman-shifted components, thus amplifying them. In stimulated scattering, the fluctuations of the optical medium (vibrations, density variations, etc.) are *induced* and *amplified* by the external electromagnetic radiation. In contrast, spontaneous scattering originates from the naturally occurring (thermally driven, for example) fluctuations of the material. Because the linear optical properties of the medium are modified by the presence of an exciting laser beam (specifically its irradiance), the various stimulated scattering mechanisms are classified as third-order nonlinear optical processes.

In stimulated Raman scattering (SRS), one is most often looking for a new frequency generation at the wavelength corresponding to the energy difference of levels $|2'\rangle$ (or $|2\rangle$) and $|3\rangle$ shown in Fig. 5. Stimulated Raman gain and loss applied to an input beam at this frequency can be obtained as well. Polarization effects occurring in the nonlinear wave mixing process can also be studied in what is known as *Raman-induced Kerr effect spectroscopy* (RIKES, see Chap. 5, “Optical Properties of Semiconductors”). We also mention two other classes of third-order nonlinear spectroscopy: coherent anti-Stokes Raman spectroscopy (CARS) and coherent Stokes Raman spectroscopy (CSRS). In these interactions, illustrated in Fig. 6, two external laser fields at frequencies ω_1 and ω_2 are supplied. There must be a third-order nonlinear polarization present as in SRS, leading to frequency mixing and new wavelengths.

The somewhat subtle differences distinguishing SRS, CARS, and CSRS are the number and location of intermediate levels (designated by dashed lines in Fig. 6). Consider two excitation frequencies

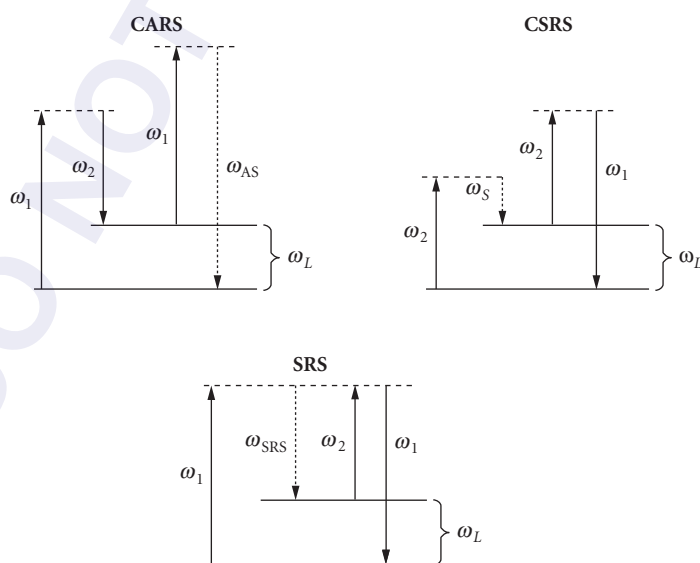


FIGURE 6 Stimulated Raman processes.

with $\omega_1 > \omega_2$ probing a given material system with real energy levels separated by frequency ω_L in Fig. 6. In CARS, short-wavelength photons are detected at $\omega_{as} = 2\omega_1 - \omega_2$. For the CSRS arrangement, the excited intermediate states are at lower energy and a longer-wavelength photon at frequency $\omega_s = 2\omega_2 - \omega_1$ is detected. Note that SRS is obtained when the intermediate levels are degenerate; SRS is thus a special case of the nonlinear interaction. In linear, spontaneous Raman scattering, a single exciting electromagnetic wave is required. In SRS, two input fields are involved; they just happen to be at the same frequency and invariably are supplied by a single laser source. Unlike spontaneous Raman scattering, however, SRS is a function of the third-order nonlinear susceptibility and hence depends nonlinearly on the irradiance of the excitation laser.

It is important to emphasize that all three of these stimulated Raman processes (SRS, CARS, and CSRS) are essentially a mixing of three waves to produce a fourth wave via the third-order nonlinear polarization. Analysis of the problem is made using second-order perturbation theory in quantum mechanics. The tensor nature of the third-order susceptibility and the multitude of ways the interacting waves of various polarization states can mix lead to complicated expressions. One finds resonance denominators quantifying the efficiency of the wave mixing process. The scattering efficiency is governed by the proximity of photon energies to real energy eigen states in the system. In principle, SRS, CSRS, and CARS can all take place in an experiment; the various generated beams can be distinguished by substantially different angles of propagation when leaving the irradiated sample. These angles are readily determined by phase-matching conditions for the nonlinear interaction. The wave vectors of the interacting photons can be arranged for maximum output signal of the desired Raman process. The phase-matching condition is obtained automatically in SRS, but careful orientation of the interacting beams can lead to very narrow linewidths and extremely accurate spectroscopic measurements. Some applications of stimulated Raman scattering include high-resolution spectroscopy of gases,^{11,17,58,59} spin-flip Raman scattering, stimulated polariton (the quanta of photon-phonon coupling) scattering, and ultrafast time-resolved measurements.^{60,61} The reader should also be aware that Raman spectroscopy beyond the third-order nonlinear susceptibility has been demonstrated. Further information can be obtained in texts on nonlinear laser spectroscopy.^{3,6,10,16,62–65}

Stimulated Brillouin Scattering

Stimulated Brillouin scattering (SBS) is an important third-order nonlinear optical effect that has been widely used for efficient phase conjugate reflection of high-power lasers.²⁶ An incident laser beam can scatter with the periodic refractive index variations associated with a propagating acoustic wave. The scattered light, depending on the propagation direction of the acoustic wave, will be Stokes or anti-Stokes shifted by the frequency of the acoustic wave. The process is stimulated because the interference of the incident and scattered wave can lead to an amplification of the acoustic wave, which then tends to pump more energy into the scattered wave. This positive feedback process can cause an exponential growth of the SBS beam and very high efficiencies in the right circumstances. Optical feedback to the medium is accomplished in one of two ways: (1) *electrostriction* is local compression of the material in response to the strength of the electromagnetic field with a commensurate refractive index change; and (2) *linear optical absorption* by the laser field leads to local heating, expansion, density fluctuations, and thus periodic modulation of the refractive index. The latter effect is an example of a cascaded $\chi^{(1)}:\chi^{(1)}$ process, which is the subject of Sec. 16.9. Electrostriction is usually associated with SBS, and we discuss it here.

Consider again the diatomic molecule that was used to illustrate the Raman effect. In the presence of an external electric field, it acquires a polarizability described by Eq. (34). As we have seen, the induced dipole can interact with the field. Electrostriction accounts for the ability of the electric field to do work on the polarized molecule—pulling and pushing it by electrostatic forces. The molecules will move and tend to pile up in regions of high field, increasing the local density. Associated with these density changes will be a change of refractive index. Density fluctuations can also be generated by the change in pressure that accompanies a propagating acoustic wave: Pressure nodes will exist in the peaks and valleys of the acoustic wave. Electrostriction therefore provides a coupling mechanism between acoustic waves and electromagnetic waves.

It is important to emphasize that the periodic modulations in an electrostrictive medium are propagating spatial fluctuations modulated at the frequency of traveling acoustic waves. When the density fluctuations are stationary, we can have stimulated Rayleigh scattering. A thorough, detailed discussion of the many (often intricate) issues in stimulated Brillouin and Rayleigh scattering can be found in textbooks on nonlinear optics.^{3,11,16,17,28}

16.8 TWO-PHOTON ABSORPTION

Two-photon absorption (2PA) is the process by which the energy gap between two real states is bridged by the simultaneous (in the context of the Uncertainty Principle discussed in Sec. 16.2) absorption of two photons, not necessarily at the same frequency. Both photons have insufficient energy to complete the transition alone; 2PA is thus observed in the spectral range where the material is normally transparent. When the two photons are present together for a fleeting instant of time determined by the Uncertainty Principle, an optical transition can take place.

Quantum mechanically, we can think of the first photon making a virtual transition to a non-existent state between the upper and lower levels (Fig. 2c). If the second photon appears within the virtual lifetime of that state, the absorption sequence to the upper state can be completed. If not, the virtual transition collapses back to the ground state, and no absorption takes place. To have an appreciable rate of 2PA, photons must be supplied at a rate high enough that there is a reasonable probability two photons will both be present during the virtual lifetime. Because the virtual lifetime is so short, photon fluxes must be high, and therefore power levels from laser beams are required.

The efficiency of 2PA is affected by the proximity of the input photons to a real state of the system. It is important to note that there must be an allowed optical transition linking the initial state and this real state. The closer one of the input photons coincides with a real state, the stronger the 2PA. When the intermediate state of 2PA is also a system resonance, the situation is commonly referred to as *excited state absorption* (ESA)—a sequence of two linear absorption processes that leads to an effective third-order nonlinearity. Excited state absorption is thus a cascaded $\chi^{(1)}:\chi^{(1)}$ effect, giving rise to an effective third-order nonlinearity (Sec. 16.9). It has implications for optical power limiting and is discussed in Chap. 13, “Optical Limiting.”

In stimulated scattering, the *difference* frequency of two input electromagnetic fields $\hbar\omega_i - \hbar\omega_j$ equals a characteristic energy resonance of the material system. In 2PA, an energy resonance exists at the *sum* of the two input fields: $\hbar\omega_i + \hbar\omega_j$. In Secs. 16.2 and 16.3, 2PA was shown to be associated with the imaginary part of $\chi^{(3)}$. This is because it is an absorption process (i.e., it is exactly resonant with two eigenstates of the system). It is the only NLA process (i.e., a process associated with the imaginary part of $\chi^{(3)}$) that can be simply studied with a single photon frequency.

Two-photon absorption in semiconductors is one of the most thoroughly studied subjects in the entire field of nonlinear optics.⁶⁶ The 2PA coefficient (often written β or α_2^{2PA}) of bulk semiconductors has been calculated using models involving only two parabolic bands and also with more complex band structure.⁶⁷ It is defined by the rate of electron-hole pair excitation: $dN/dt = \beta I^2 / (2\hbar\omega)$. The two-parabolic band model gives a comparatively simple yet general and accurate description of 2PA for a large class of semiconductors. The theoretical result for single frequency excitation can be expressed as

$$\beta(m/W) \approx \frac{B}{n_0^2 E_g^3} \bar{F}_2 \left(\frac{\hbar\omega}{E_g} \right) \quad (39)$$

where $\bar{F}_2(x) = 2(2x-1)^{3/2}/x^5$ and $B = 5.67 \times 10^{-66}$ ($x = \hbar\omega/E_g$ and the energy bandgap E_g is in joules). The best empirical fit to experimental data is obtained with B adjusted to a slightly higher value of 9.06×10^{-66} . The function \bar{F}_2^2 describes the dispersion of 2PA and is plotted in Fig. 3. The intimate relation between NLA and NLR and the role of 2PA in semiconductors is explored in Sec. 16.4 and in Refs. 31 and 67 to 69. There are important practical implications for 2PA in

semiconductors and dielectrics. It can enhance or degrade optical switching performance in semiconductor devices and lead to optical damage in laser window materials. 2PA is also the basis of Doppler-free spectroscopy of gases.^{57,59}

16.9 EFFECTIVE THIRD-ORDER NONLINEARITIES; CASCADED $\chi^{(1)}:\chi^{(1)}$ PROCESSES

Effective third-order nonlinearities occur when one of the transitions in our four-photon interaction picture is resonant, providing a path of linear absorption. Linear absorption is a mechanism to directly couple laser light into the system—with sufficiently intense laser light, the linear optical properties of a material can be modified. An effective third-order nonlinearity occurs when linear absorption affects the refractive index. We give some examples here.

Optically Generated Plasmas

Optical generation of stable plasmas is readily obtained in semiconductors and most studies of this subject have been made with these materials (see Chap. 8, “Fundamental Optical Properties of Solids” and Chap. 5, “Optical Properties of Semiconductors”). For cascaded linear processes that we discuss here, the formation of a free electron-hole pair occurs by direct bandgap excitation by an incident photon. The optically produced carriers augment the background electron-hole density, and the plasma remains electrostatically neutral. If the generation of an excess amount of plasma exceeds the rate of loss (by recombination or diffusion) on the time scale of interest, the plasma will modify the linear optical properties of the semiconducting material. The simplest way to see this is via the classical Drude model, where the refractive index of a metal or semiconductor is^{70,71}

$$n = n_0 \sqrt{1 - \frac{\omega_p^2}{\omega^2}} \quad (40)$$

In this equation, ω is the angular frequency of the light, n_0 is the linear index in the absence of significant free carrier density, and ω_p is the density-dependent plasma frequency:

$$\omega_p = \sqrt{\frac{Ne^2}{m\epsilon}} \quad (41)$$

where N is the electron-hole pair density, e is the electronic charge, m is the reduced mass of the positive and negative charge carriers (electrons, holes, or ions), and ϵ is an appropriate background dielectric constant. Note that as the carrier density increases, the refractive index decreases. The material is usually excited by a laser beam with a nonuniform spatial profile such as a Gaussian, giving rise to negative-lensing and self-defocusing, assuming $\omega < \omega_p$.

The situation becomes complicated at densities where many-body effects become important or when the Drude model ceases to be valid. We also note that optical generation of plasmas can also occur as the result of nonlinear mechanisms in the presence of high laser fields, and we aren't, of course, restricted only to solid-state plasmas. Examples of nonlinear plasma production are multiphoton absorption, laser-induced impact ionization, and tunneling. Because plasma generation and the concomitant refractive index modification are caused by a nonlinear optical process, the order of the nonlinearity is higher than three. Although high-order nonlinearities are a rich subset of the field, they are not dealt with in this chapter. A number of review articles and textbooks on laser-induced change to the refractive index due to plasma generation via linear absorption are available.^{31,42,72–74}

Absorption Saturation

Absorption saturation is a well-known example of a cascaded linear process. Consider a homogeneously broadened system of two-level atoms (i.e., a system of identical two-level atoms) with the energy diagram shown in Fig. 7. The lower and upper states are resonant with a photon depicted by the vertical arrow, and there is linear absorption of incident light. Associated with the absorption is a spectral linewidth (with a Lorentzian shape for a homogeneous system), illustrated on the right side of the diagram in Fig. 7. An induced dipole or polarization is set up between the two states upon excitation by a photon. The states are coherently coupled, and this coherence is in phase with the exciting electromagnetic field. In a real system, this coherence will be quickly destroyed by collisions with other atoms. The rate at which the coherence is washed out determines the spectral width of the absorption profile and hence the frequency response of the imaginary component of the linear susceptibility. The real part of the linear susceptibility is obtained by the Kramers-Kronig transformation, giving rise to what is traditionally known as anomalous dispersion of the refractive index shown on the left side of Fig. 7. Note that the refractive index is positive for frequencies below resonance and negative at frequencies above it.

We have assumed that the rate at which photons are supplied to the system produces a negligible change of population in the upper and lower states. This means that the rate of population relaxation (recombination and diffusion, for example) is much faster than excitation. When the incident irradiance is sufficiently high, however, this may no longer be the case. The upper level can become appreciably occupied, reducing the availability of terminal states for optical transitions. The absorption thus decreases or bleaches, indicated by the dashed lines in Fig. 7. Associated with the change of absorption is a change of refractive index. The relationship between absorption and refraction can again be handled with the Kramers-Kronig transformation *provided* we consider the system as being composed of both the atoms *and* the input light beam, specifically the nonequilibrium change of population created by the input light. This is exactly the same mathematical formulation used in the description of nonresonant third-order nonlinearities introduced in Sec. 16.4. The difference here is that the absorption of photons is a resonant, linear process. The linear absorption of light then affects the linear optical properties of absorption and dispersion. Because the reduction of absorption depends directly on laser irradiance, it behaves like a third-order nonlinearity. We point out that the resonant nature of absorption saturation (and the associated nonlinear refraction) can lead to an unacceptable deviation from a third-order susceptibility description when the input irradiance goes even higher. In the high-power regime, a nonperturbative approach (i.e., not represented by a power series expansion) must be used.⁷⁵

Absorption saturation spectroscopy (with emphasis on gases) is discussed at length in Ref. 59. It is also a principle nonlinear effect in bulk and quantum-confined semiconductors (where the simple two-level picture previously outlined must be substantially modified), having implications for optical switching and bistability.^{42,53,72,76}

We briefly mention the density matrix (see Chap. 10,⁷⁶ “Nonlinear Optics”). This is a powerful method of analysis for the resonant interaction of light with a two-level system. In addition to absorption saturation, the nonlinear optical effects described by the density matrix include Rabi oscillations, photon echoes, optical nutation, superradiance, self-induced transparency, and optical-free induction

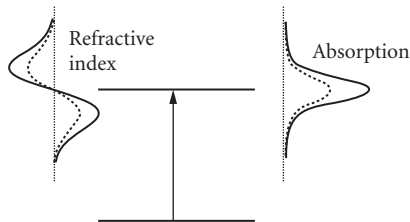


FIGURE 7 Absorption saturation and refraction of a two-level atomic system.

decay, which are not immediately associated with a third-order nonlinearity derived from a perturbation expansion of the polarization. Experimental manifestations of these phenomena, however, can often be represented by an effective third-order nonlinearity. The nonlinear optics of the two-level system and the associated optical Bloch equations derived from the density matrix formulation of the problem are discussed in Chap. 11, “Coherent Optical Transients.” The reader is also referred to many excellent textbooks and monographs.^{3,11,24,77,78}

Thermal Effects

Linear absorption of light must result in energy deposition in the irradiated material. If the rate of energy deposition significantly exceeds its rate of removal, heating can take place. As the energy of a collection of atoms and molecules increases, it's easy to understand that their macroscopic optical properties will be altered. When there is a linear relationship between laser irradiance and the refractive index, an effective third-order nonlinearity will result. The thermo-optic effect has a relatively straightforward physical interpretation and is arguably one of the most important optical nonlinearities. It is often the power-limiting mechanism of a high-power solid-state laser—where excessive circulating optical flux can cause *thermal blooming* in the laser crystal. This heat-induced lensing effect can destroy the beam quality. Thermal-index coefficients have been extensively tabulated, and the reader is referred to Refs. 18 and 42.

Photorefraction

Photorefraction results from the spatial redistribution of electrons and/or holes in a solid (see Chap. 12, “Photorefractive Materials and Devices”). Electron-hole pairs are generated by linear absorption of laser light. If the excitation geometry produces a spatial modulation of irradiance such as a two-beam interference pattern, the electrons and holes will arrange themselves in accordance with the spatial irradiance profile. Often the linear absorption involves impurity levels. Mobile carriers tend to diffuse from the bright regions, leaving fixed charges behind. If a sinusoidally modulated, two-beam interference fringe pattern is written in a doped photorefractive crystal, for example, fixed charges of ionized states will be prevalent in the high irradiance regions while mobile charge carriers will tend to accumulate in regions with low light levels. A modulated space charge must exist, and therefore a modulated electric field pattern must be present as well. This field alters the refractive index via the linear electro-optic effect (i.e., Pockel's effect). The photorefractive nonlinearity is clearly nonlocal, as it requires a spatial modulation of charge density. Manipulation of the carriers can also be obtained with static electric fields, and the response time tends to be of the order of seconds.^{3,18,79,80}

16.10 EFFECTIVE THIRD-ORDER NONLINEARITIES; CASCADED $\chi^{(2)}:\chi^{(2)}$ PROCESSES

Materials lacking a center of inversion symmetry have nonzero $\chi^{(2)}$ and exhibit a second-order nonlinear polarization. This is the second term in the polarization power series expansion of Eq. (1) that is responsible for the most well-known effects in nonlinear optics, including second-harmonic generation (sum and difference frequency generation, optical rectification), and optical parametric processes (see Chap. 17, “Continuous-Wave Optical Parametric Oscillators”). It is also possible to cascade two $\chi^{(2)}$ processes to produce an effect that mimics a $\chi^{(3)}$ process. The most common and efficient way of producing THG, for example, is by a $\chi^{(2)}$ cascade process. In this interaction, an input source at ω generates SHG via the second-order susceptibility $\chi^{(2)}$ ($2\omega = \omega + \omega$); the second harmonic and fundamental then mix in a second (or the same) nonlinear crystal to produce the third harmonic by sum frequency generation $\chi^{(2)}$ ($3\omega = 2\omega + \omega$). This type of nonlinearity is nonlocal because the two processes (SHG and SFG) take place in spatially separate regions.

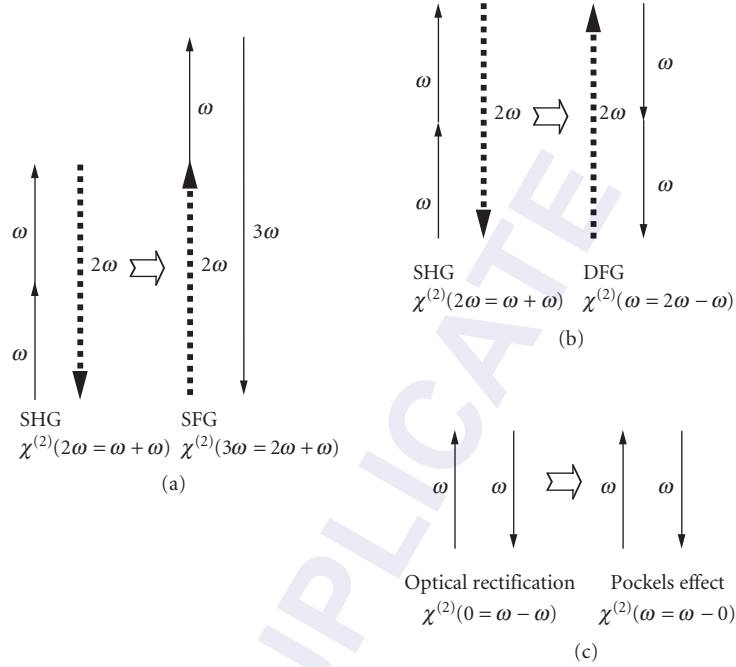


FIGURE 8 Cascaded $\chi^{(2)}:\chi^{(2)}$ effective third-order nonlinearities.

In addition to THG, all other $\chi^{(3)}$ effects have an analogous process in $\chi^{(2)}:\chi^{(2)}$ cascading. Consider the case where we have one frequency ω . The corresponding cascaded $\chi^{(2)}$ effects are depicted in Fig. 8. Recall that the intrinsic $\chi^{(3)}$ is manifest as: (1) THG, (2) 2PA, and (3) the ac Stark effect. The horizontal block arrows indicate the point of cascade (i.e., the propagation of the beams between the two $\chi^{(2)}$ interactions). There are always two distinguishable $\chi^{(2)}$ interaction regions; hence $\chi^{(2)}:\chi^{(2)}$ is clearly nonlocal. We also point out that in general the photon frequencies can be different.

Although the analogy is limited, cascaded $\chi^{(2)}$ effects exhibit NLA and NLR mimicking $\chi^{(3)}$. In the case of SHG, for example, the manifestation of $\chi^{(3)}$ NLA is depletion of the pump beam. The utility of cascaded $\chi^{(2)}:\chi^{(2)}$ for producing large nonlinear phase shifts has been realized only recently.^{81–83} The analysis of $\chi^{(2)}$ is relatively straightforward, involving the coupled amplitude equations governing the propagation of the interacting beams. For example, the nonlinear phase shift imposed on a fundamental beam (ω) as it propagates through an SHG crystal of length L with a phase mismatch $\Delta k = k(2\omega) - 2k(\omega)$ and assuming small depletion is⁸²

$$\Delta\phi \approx \tan^{-1}\left(\frac{\Delta k L \tan(\beta L)}{2} - \frac{\Delta k L}{2}\right) \quad (42)$$

where $\beta = \sqrt{(\Delta k L / 2)^2 + \Gamma^2}$, $\Gamma = \omega \chi^{(2)} |E_0| / 2c \sqrt{n(2\omega)n(\omega)}$, and E_0 denotes the electric field of the fundamental beam. An effective $\chi^{(3)}$ can be obtained if we expand $\Delta\phi$ to lowest order in Γ and use $\Delta\phi = \omega L n_2^{\text{eff}} I / c$ to give⁸⁴

$$n_2^{\text{eff}} = \frac{9}{c^2 4\pi\epsilon_0} \frac{\omega d_{\text{eff}}^2 L}{n^2(\omega)n(2\omega)} \left[\frac{\pi}{\Delta k L} \left(1 - \frac{\sin(\Delta k L)}{\Delta k L} \right) \right] \quad (43)$$

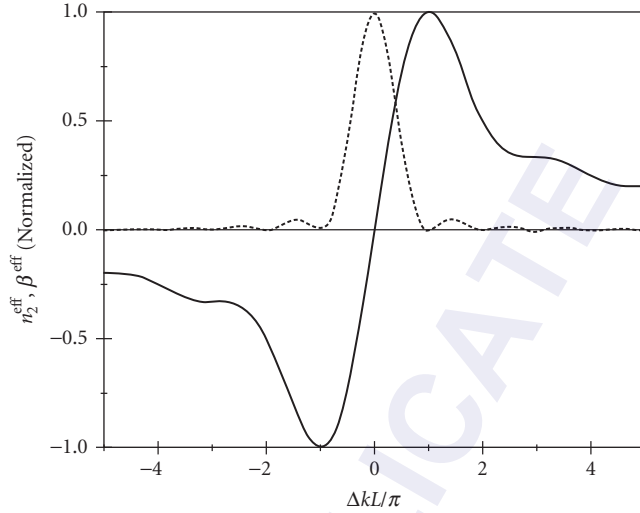


FIGURE 9 Calculated phase shifts in cascaded $\chi^{(2)}:\chi^{(2)}$.

Here $d_{\text{eff}} = \chi^{(2)}(2\omega = \omega + \omega)/2$ is the effective tensor component of the second-order nonlinear susceptibility for a given experimental geometry.

Phase mismatch is represented by the bracketed term in Eq. (43), plotted in Fig. 9. Also shown in Fig. 9 is the depletion of the fundamental beam which, on the same level of approximation, can be regarded as an effective two-photon absorption coefficient that scales as $\beta^{\text{eff}} \sim \text{sinc}^2(\Delta kL)$. Cascaded $\chi^{(2)}$ leads to an effective refractive index modulation (n_2^{eff}) and two-photon absorption (β^{eff}), but one should not conclude that the material's index of refraction is altered or energy is deposited in the material. These coefficients describe only nonlinear phase shifts and the conversion of the fundamental beam to second-harmonic beams. The nonlinear phase shift from the $\chi^{(2)}:\chi^{(2)}$ process has been used to demonstrate nonlinear effects analogous to those observed previously with the intrinsic optical Kerr effect. These include self-focusing and self-defocusing,⁸¹ all-optical switching,^{85,86} soliton propagation,^{87,88} and laser mode-locking.⁸⁹

16.11 PROPAGATION EFFECTS

When the nonlinear optical polarization $P^{\text{NL}}(t)$ is known, the propagation of optical fields in a nonlinear medium can be analyzed with the aid of Maxwell's equations:

$$\nabla \times \nabla \times E + \frac{1}{c^2} \frac{\partial^2 E}{\partial t^2} = -\mu \frac{\partial^2 P}{\partial t^2} \quad (44)$$

where $P(t)$ is the total polarization including the linear and nonlinear terms. The slowly varying envelope approximation is then usually made to reduce the above equation to a system of four coupled nonlinear differential equations for the four interacting fields. For *thin* nonlinear media, where there is no significant distortion of the spatial beam and temporal shape upon propagation, the problem simplifies greatly. Equations (19) and (20) were obtained with this approximation. For *thick* nonlinear media, however, linear and nonlinear diffraction as well as dispersion cannot be ignored. In this section, we discuss two important propagation phenomena: *self-focusing* and *soliton formation*.

Self-Focusing

Self-focusing occurs in materials with a positive intensity-dependent refractive index coefficient ($n_2 > 0$). Self-focusing (or *Kerr-lensing*) causes spatial collapse of the laser beam when it propagates through transparent optical materials, often leading to optical damage. It is a consequence of the nonuniform spatial profile of the laser beam.^{2,3,17,18} For a *thin* nonlinear material, one makes the so-called parabolic approximation for the nonlinear phase shift to obtain an approximate Kerr-lens focal length, assuming a Gaussian beam of radius w ($1/e$ of the electric field profile):⁹⁰

$$f_{\text{NL}} \approx \frac{aw^2}{4Ln_2I} \quad (45)$$

where L is the thickness of the medium, I is the irradiance, and $6 > a > 4$ is a correction term. Note that when n_2 is negative, Eq. (45) shows there will be a negative focal length and thus self-defocusing of the incident beam.

Equation (45) is valid for $f_{\text{NL}} \gg L$ and $Z_0 \gg L$, where Z_0 is the diffraction length (Rayleigh range) of the incident beam. This is the so-called *external self-action* regime.⁹¹ This approximation fails for thick nonlinear media and/or at high irradiance (i.e., *internal self-action*). Equation (44) must then be solved numerically. Analysis shows that self-lensing of a Gaussian beam overcomes diffraction at a distinct threshold power (i.e., the self-focusing threshold), given by the approximate formula:^{90,92}

$$P_{\text{cr}} \approx \frac{a\lambda^2}{8\pi n_2 n_0} \quad (46)$$

Note that for sufficiently thick media, the self-focusing threshold occurs at a critical power, not at a threshold irradiance (i.e., the power at which the self-focusing overcomes diffraction). Self-focusing and diffraction both scale with beam area, thus canceling out the spot size dependence in Eq. (46). Self-focusing and self-defocusing (collectively called *self-action effects*) are often employed in optical limiting applications. Self-action is also the essential mechanism for mode-locking cw solid-state lasers, commonly known as *Kerr-lens mode-locking*.⁹³

Solitons

Soliton waves are realized in many different physical circumstances, ranging from mechanical motion to light propagation. In general, they are robust disturbances that can propagate distortion-free for relatively long distances. The robustness of optical solitons can be manifest in the time domain (temporal solitons), transverse space (spatial solitons), or both (light bullets). Temporal solitons have been extensively studied in optical fibers because of their tremendous utility in long-distance optical communication. They exist as a consequence of a balance between the competing effects of linear refractive index dispersion and nonlinear phase modulation.

Assume a single beam propagating in a long nonlinear waveguide characterized by an instantaneous nonlinear index coefficient n_2 and a linear refractive index $n(\omega)$. Ignoring spatial effects (i.e., diffraction), we write the electric field as $\mathcal{E}(z, t) = A_0 u(z, t - z/v_g) \exp(i\omega_0 t - ik_0 z) + \text{c.c.}$ From Eq. (44), one derives a differential equation describing the evolution of the soliton field envelope $u(z, t)$:^{3,51,94}

$$-i \frac{\partial u}{\partial z} + \frac{k_2}{2} \frac{\partial^2 u}{\partial \tau^2} = \Delta k_{\text{NL}} |u|^2 u \quad (47)$$

Here $v_g = d\omega/dk|_{\omega=\omega_0}$ is the soliton pulse group velocity, τ is a retarded time $\tau = t - z/v_g$, $k_2 = d^2k/d\omega^2|_{\omega=\omega_0}$ gives the group velocity dispersion (GVD), and $\Delta k_{\text{NL}} = n_2 I_0 \omega_0 / c$ is the irradiance-dependent change of the propagation wave vector. In MKS units, the peak intensity of the soliton pulse is $I_0 = n_0 \epsilon_0 c |A_0|^2 / 2$. Equation (47) is called the *nonlinear Schrödinger equation* (NLSE) and can be solved exactly. One solution gives the fundamental soliton pulse:

$$u(z, \tau) = \text{sech}(\tau/\tau_0) e^{ikz} \quad (48)$$

where $\tau_0^2 = -k_2/\Delta k_{\text{NL}}$ is the soliton pulsewidth and $\kappa = -k_2/2\tau_0^2$. Note that the modulus of the soliton pulse envelope $|u|$ remains unperturbed upon propagation. For this solution to exist, the GVD (k_2) and the nonlinear refraction (Δk_{NL}) must have opposite signs. For transparent optical solids, including silica glass optical fibers, the nonlinear index coefficient n_2 is almost always positive, which means that a negative GVD is required. In fused silica fibers, the point of balance is attained at a wavelength of $\lambda \approx 1.55 \mu\text{m}$. This is also a spectral region with very low absorption loss. This wavelength has become the standard for the telecommunications industry. Optical solitons in fibers were first reported by Mollenauer et al.⁵¹

Spatial solitons refer to the propagation of an optical beam without any change or distortion to its spatial irradiance distribution. In this type of soliton, a point of stability is achieved between linear diffraction (causing the beam to diverge) and nonlinear self-focusing. In the presence of a $\chi^{(3)}$ nonlinearity, a stable solution to the NLSE can be found in one spatial dimension only.^{95,96} Using a cascaded $\chi^{(2)};\chi^{(2)}$ nonlinearity, however, two-dimensional spatial solitons (or solitary waves) have been demonstrated.^{87,88} Two-dimensional spatial solitons can be realized for a cascaded $\chi^{(2)};\chi^{(2)}$ process because of its different behavior compared to $\chi^{(3)}$ at large nonlinear phase shifts. Specifically, cascaded $\chi^{(2)};\chi^{(2)}$ exhibits a saturation of the nonlinear phase-shift that is a direct consequence of depletion of the fundamental beam.

16.12 COMMON EXPERIMENTAL TECHNIQUES AND APPLICATIONS

There are a variety of experimental methods for determining the characteristics (magnitude, response time, spectrum, etc.) of $\chi^{(3)}$ (or $\chi_{\text{eff}}^{(3)}$). The merit of a technique depends on the nature of the nonlinearity and/or the specific property that we wish to measure. An example is obtaining short-time resolution at the expense of sensitivity. Nonlinear optical coefficients can be determined absolutely or relative to a reference material. In the former case, accuracy is determined by the ability to precisely characterize the incident beams. There are many potential sources of error and misinterpretation in nonlinear optical measurements. In *thick* samples, for example, the phase shift that occurs during beam propagation can lead to varying irradiance at different points within the sample. This can be quite difficult to account for and properly model. It is usually best to work in the *external self-action* regime (i.e., thin-sample limit so that beam propagation effects can be ignored)⁹¹ (see also Sec. 16.3). This greatly simplifies data analysis, since the equation describing nonlinear absorption can be separated from nonlinear refraction. Even if the thin-sample approximation is satisfied, nonlinear refraction can deflect light so strongly *after* the sample that the detector does not collect all the transmitted energy. This will lead to an overestimation of the nonlinear loss. Particular care must be exercised when using ultrashort pulses. Pulse broadening effects due to group velocity dispersion, for example, may cast ambiguity on the magnitude as well as response time associated with a nonlinearity.

We briefly discuss a few of the commonly used experimental methods: four-wave mixing, excite-probe techniques, interferometry, and Z-scan. It is practically impossible for any single technique to unambiguously separate the different nonlinear responses. Experiments are generally sensitive to several different nonlinearities at once. Different measurements are usually required to unravel the underlying physics, by varying parameters such as irradiance and pulse width. Near-instantaneous nonlinearities such as two-photon absorption and the optical Kerr effect should be independent of pulse width. Slower nonlinear responses will change as the pulse width approaches the response time. Ultrafast and cumulative nonlinearities are often present simultaneously in experiments (e.g., semiconductors), thus hindering their experimental isolation.

Time-Resolved Excite-Probe Techniques

Pump-probe (excite-probe) measurements allow the study of temporal dynamics in nonlinear absorption.^{27,60} In the usual implementation, a relatively strong pump pulse excites the sample and

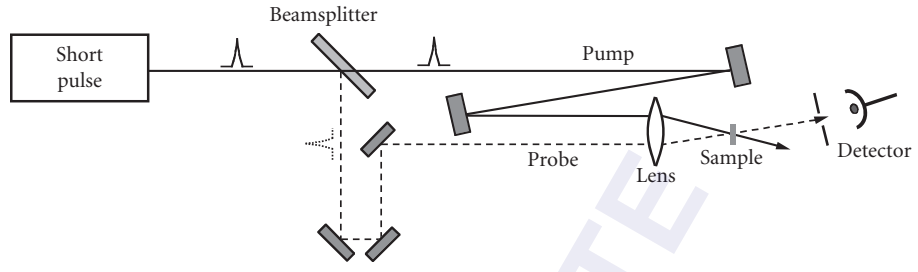


FIGURE 10 Pump-probe experiment.

changes its optical properties (see Fig. 10). A weaker probe pulse interrogates the excitation region and detects changes. By varying the relative time separation of the two pulses (i.e., by appropriately advancing and delaying the probe pulse), the temporal response can be mapped out. Specifically, relatively slow and fast nonlinear responses can be identified. Often, but not always, the probe is derived from the excitation beam. In degenerate pump-probe (identical frequencies), the probe beam is isolated from the pump beam in a noncollinear geometry (i.e., spatial separation as in Fig. 10) or by orienting the probe with a different polarization.

Nondegenerate nonlinear absorption spectra can also be measured; one approach is to use a fixed-frequency laser pump and continuum (white light) probe such as the output of a flash-lamp.⁹⁷ The temporal width of a flashlamp source is usually much longer than the laser pulse, which causes a convolution of the fast two-photon response with much longer-lived cumulative nonlinearities in the probe spectrum. The availability of femtosecond white-light continuum sources has allowed nondegenerate spectra to be obtained on short time scales where the ultrafast response dominates.^{27,98}

Interpretation of the nonlinear response is complicated by the fact that pump-probe experimental methods are sensitive to any and all induced changes in transmission (or reflection); pump-induced phase shifts on the probe are not readily detected. A time-resolved technique that is very sensitive to index changes is the optical Kerr-gate. This is a form of the pump-probe experiment where induced anisotropy in a time-gated crystal leads to polarization changes.⁹⁹ Beyond the two-beam pump-probe, three-beam interactions can produce a fourth beam through NLA and/or NLR. This is known as *four-wave mixing* and is discussed in the next section. The development of high-irradiance, femtosecond-pulsed laser systems has led to the evolution of pump-probe measurements that automatically yield the nondegenerate nonlinear absorption spectrum. The femtosecond pulse is split into two beams: one beam is used for sample excitation while the other beam is focused into an appropriate material to produce a white-light continuum for probing. This white-light continuum is used to measure the response over a range of frequencies ω' ; this data is suitable for numerical evaluation via the K-K integral in Eq. (24). For a sufficiently broad spectrum of data, the K-K integration yields the nondegenerate n_2 coefficient.¹⁰⁰

Four-Wave Mixing

The most general case of third-order interaction has all four interacting waves (three input and one scattered) at different frequencies. Generating and phase-matching three different laser wavelengths in an experiment is a formidable task; the benefit is often an improved signal-to-noise ratio.

The other extreme is when all four waves have identical frequency, a situation known as *degenerate four-wave mixing* (DFWM), although it is commonly (and less precisely) referred to as *four-wave mixing* (see Chap. 12, "Photorefractive Materials and Devices" and Chap. 5, "Optical Properties of Semiconductors"). DFWM is readily implemented in the laboratory, since only a single laser source is needed. There are two general cases: *nonresonant* and *resonant* DFWM. In transparent media (i.e., nonresonant) the index of refraction is usually a linear function of laser irradiance, and nonresonant DFWM (wavelength far from an absorption resonance) leads to optical phase conjugation. Phase

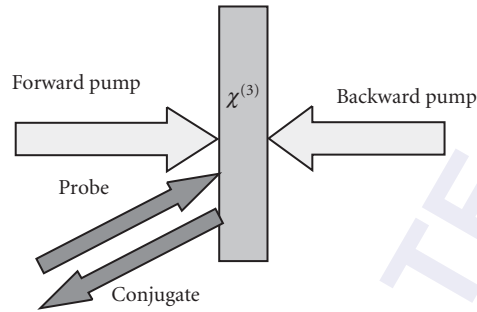


FIGURE 11 Schematic diagram of a four-wave mixing experiment.

conjugation by the optical Kerr effect (Sec. 16.5) is one of the most important applications involving third-order nonlinearities.²⁶ Nonresonant DFWM leads to the formation of a phase grating due to the spatial modulation of the refractive index. Two of the beams write the phase grating while a third reads or probes the grating by diffracting from it, thereby generating a fourth beam (see Fig. 11). The diffracted beam can either be transmitted or reflected (i.e., a phase-conjugate beam) from the material in a direction determined by the wave vectors of the interacting photons. In some experiments, the writing beams also serve to read the grating. One of the difficulties in interpreting DFWM data for third-order nonlinearities is that the signal is proportional to $|\chi^{(3)}|^2 = |\Re\{\chi^{(3)}\} + \Im\{\chi^{(3)}\}|^2$ (i.e., NLA and NLR both contribute). Separating the effects is difficult without performing additional experiments. Techniques that study different polarizations can provide information on the symmetry properties of $\chi^{(3)}$.

In resonant DFWM, there is the added complication of optical absorption at the frequency of the interacting light beams.¹⁰¹ This is an example of a cascaded $\chi^{(1)}$; $\chi^{(1)}$ effective third-order nonlinearity discussed in Sec. 16.9, where absorption causes population in excited states, resulting in a spatial grating due to the spatial modulation of population. In principle, both phase and absorption gratings are present in resonant DFWM. In practice, it is usually the intensity-dependent changes of population (i.e., effective $\chi^{(3)}$) that dominate the nonlinear polarization, although this is not always the case.⁷⁶ For example, photocarrier generation in a semiconductor can alter the bulk plasma frequency and thus modulate the refractive index, leading to a strong phase grating (see Sec. 16.9).

The diffracted beam contains a wealth of information about the system under study. In nonresonant DFWM, the absolute magnitude and spectral width of the Kerr-effect nonlinearity (n_2) can be obtained. Even more can be deduced from time-resolved DFWM, where the interacting beams are short laser pulses. If the pulses (two and sometimes three separate pulses) are delayed with respect to each other, the dynamic response of the nonlinear polarization can be measured. In resonant DFWM, the diffracted beam measures the coherent response of the optically coupled eigenstates of the system. The linewidth of the diffracted beam indicates the rate at which various physical processes broaden the transition. The nonlinear polarization can be washed out by mechanisms such as population relaxation and diffusion and scattering events associated with optically coupled states. Because of selection rules linking resonant states, various polarization geometries can be employed to study specific transitions. This can be very useful in studies of a complex system such as a semiconductor. Time-resolved experiments with short pulses provide information that complements and elucidates spectral linewidth data obtained from measurements with long duration or continuous laser beams.^{26,60,102}

Interferometry

Interferometric methods can be used to measure nonlinearly induced phase distortion.^{103,104} One implementation of this approach places a sample in one path (e.g., arm) of an interferometer,

and the interference fringes are monitored as a function of irradiance. If the interferometer is set up to give a series of straight-line interference fringes for low-input irradiance (linear regime), the fringes become curved near the region of high irradiance, such as the center of a Gaussian beam. The addition of a streak camera can add time resolution.²⁷ Alternatively, a third beam can be added to the experiment. The sample is in the path of one weak beam and the strong third beam. The fringe pattern of the two weak beams is monitored as a function of sample irradiance provided by the strong beam. The relative fringe shift observed when the strong beam is present and blocked gives the optical path length change. The nonlinear phase shift can thus be determined. Interferometric experiments require excellent stability and precise alignment. When these conditions are met, sensitivities of better than $\lambda/10^4$ induced optical path length change can be measured.

Z-Scan

The Z-scan was developed to measure the magnitude and sign of nonlinear refraction (NLR). It is also useful for characterizing nonlinear absorption (NLA) and for separating the effects of NLR from NLA.^{105,106} The essential geometry is shown in Fig. 12. Using a single, focused laser beam, one measures the transmittance of a sample through a partially obscuring circular aperture (Z-scan) or around a partially obscuring disk (EZ-scan¹⁰⁷) placed in the far field. The transmittance is determined as a function of the sample position (Z) measured with respect to the focal plane. Employing a Gaussian spatial profile beam simplifies the analysis.

We illustrate how Z-scan (or EZ-scan) data is related to the NLR of a sample. Assume, for example, a material with a positive nonlinear refractive index. We start the Z-scan far from the focus at a large value of negative Z (i.e., close to the lens). The beam irradiance is low and negligible NLR occurs; the transmittance remains relatively constant near this sample position. The transmittance is normalized to unity in this linear regime as depicted in Fig. 13. As the sample is brought closer to focus, the beam irradiance increases, leading to self-focusing. This positive NLR moves the focal point closer to the lens, causing greater beam divergence in the far field. Transmittance through the aperture is reduced. As the sample is moved past the focus, self-focusing increasingly collimates the beam, resulting in enhanced transmittance through the aperture. Translating the sample farther toward the detector reduces the irradiance to the linear regime and returns the normalized transmittance to unity. Reading the data right to left, a valley followed by peak is indicative of positive NLR. In negative NLR, one finds exactly the opposite: a peak followed by a valley. This is due to laser-induced self-defocusing. Characteristic curves for both types of NLR are shown in Fig. 12. The EZ-scan reverses the peak and valley in both cases. In the far field, the largest fractional changes in irradiance occur in the wings of a Gaussian beam. For this reason, the EZ-scan can be more than an order of magnitude more sensitive than the Z-scan.

We define an easily measurable quantity ΔT_{pv} as the difference between the normalized peak and valley transmittance: $T_p - T_v$. Analysis shows that variation of ΔT_{pv} is linearly dependent on the

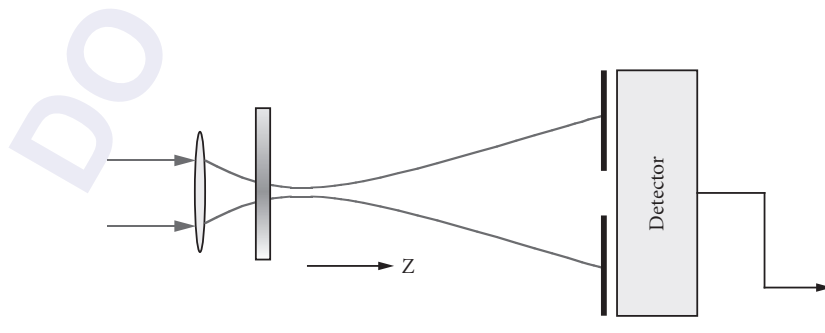


FIGURE 12 Z-scan experimental arrangement.

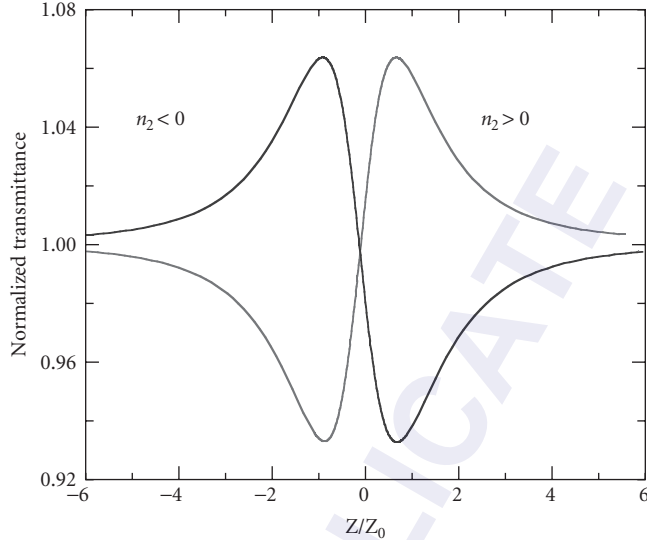


FIGURE 13 Representative curves depicting nonlinear refraction (both positive and negative) as measured by the Z-scan.

temporally averaged induced phase distortion, defined here as $\Delta\Phi_0$. If the Z-scan aperture is closed to allow linear transmission of less than 10 percent, and $\Delta T_{pv} < 1$:^{105,108}

$$\Delta T_{pv} \cong 0.41 |\Delta\Phi_0| \quad (49)$$

assuming cw illumination. If the experiment is capable of resolving transmission changes $\Delta T_{pv} \cong 1\%$, the Z-scan will be sensitive to wavefront distortion of less than $\lambda/250$ (i.e., $\Delta\Phi_0 = 2\pi/250$). The Z-scan has demonstrated sensitivity to a nonlinearly induced optical path length change of nearly $\lambda/10^3$, while the EZ-scan has shown a sensitivity of $\lambda/10^4$, including temporal averaging over the pulsewidth.

To this point in the discussion, we have assumed a purely refractive nonlinearity with no NLA. It has been shown that two-photon absorption will suppress the peak and enhance the valley. If NLA and NLR are present simultaneously, a numerical fitting procedure can extract both the nonlinear refractive and absorptive coefficients. Alternatively, a second Z-scan with the aperture removed (all the transmitted light collected) can independently determine the NLA. Considering 2PA only and a Gaussian input beam, the Z-scan traces out a symmetric Lorentzian shape. The so-called *open aperture* Z-scan is sensitive *only* to NLA. One can then divide the partially obscuring Z-scan data by the open aperture data to give a curve that shows only nonlinear refraction. By performing these two types of Z-scans, we can isolate NLR and NLA without the need for a complicated numerical analysis of a single data set obtained with an aperture.

All-Optical Switching and Optical Bistability

Since the early 1980s, there has been substantial interest in third-order nonlinear optical behavior in materials because of the potential for performing high-speed switching operations—gate speeds many orders of magnitude faster than conventional electronics have been demonstrated. The possibility of increasing data rates on information networks provides the obvious motivation for this

research. A bistable optical switch has two stable output states for a given input (i.e., a specific optical power level). This has implications for applications such as optical data storage and power limiting. Both nonresonant $\chi^{(3)}$ and resonant effective $\chi^{(3)}$ processes have been extensively studied, primarily in semiconductors. The early work looked at bulk semiconductor behavior, but as the technology matured the emphasis shifted to specially designed optical waveguides made from suitable material. At the time of this writing, both resonant and nonresonant approaches have encountered problems that have limited practical use. The bound electronic nonlinearity responds on a femtosecond time-scale but is inherently weak. The laser irradiance must be increased to compensate, but this leads to the unwelcome presence of 2PA and associated losses. Resonant nonlinearities must involve the generation of carriers (electrons and holes). While such nonlinearities can be exceptionally strong, the speed of an optical switch depends crucially on the ability to manipulate the carriers. Generation of electron-hole pairs, for example, may dramatically affect the refractive index of a semiconductor and its ability to modulate light, but if the carriers have a long recombination lifetime the switch recovery time will be relatively slow. Other significant issues that must be weighed when comparing optical switching schemes to the all-electronic approach (i.e., transistors and integrated circuits) include device packaging density and heat removal.^{1,8,42,53,72,109,110}

16.13 REFERENCES

1. N. Bloembergen, *Nonlinear Optics*, Addison-Wesley, Redwood City, California, 1992.
2. S. A. Akhmanov and R. V. Khokhlov, *Problems of Nonlinear Optics: Electromagnetic Waves in Nonlinear Dispersive Media*, Gordon and Breach Science Publishers, New York, 1972.
3. R. W. Boyd, *Nonlinear Optics*, Academic Press, Boston, 1992.
4. P. N. Butcher and D. Cotter, *The Elements of Nonlinear Optics*, Cambridge University Press, Cambridge and New York, 1990.
5. G. P. Agrawal, *Nonlinear Fiber Optics*, 2nd ed., Academic Press, San Diego, California, 1995.
6. W. Demtröder, M. Inguscio, and North Atlantic Treaty Organization, Scientific Affairs Division, *Applied Laser Spectroscopy*, Plenum Press, New York, 1990.
7. D. C. Hanna, M. A. Yuratich, and D. Cotter, *Nonlinear Optics of Free Atoms and Molecules*, Springer-Verlag, Berlin and New York, 1979.
8. F. A. Hopf and G. I. Stegeman, *Applied Classical Electrodynamics*, Reprint ed., Krieger Pub. Co., Malabar, Florida, 1992.
9. I.-C. Khoo, J.-F. Lam, and F. Simoni, *Nonlinear Optics and Optical Physics*, World Scientific, Singapore and River Edge, New Jersey, 1994.
10. V. S. Letokhov and V. P. Chebotayev, *Nonlinear Laser Spectroscopy*, Springer-Verlag, Berlin and New York, 1977.
11. M. D. Levenson and S. Kano, *Introduction to Nonlinear Laser Spectroscopy*, Rev. ed., Academic Press, Boston, 1988.
12. S. R. Marder, J. E. Sohn, G. D. Stucky, American Chemical Society, Division of Organic Chemistry, American Chemical Society, Division of Inorganic Chemistry, and American Chemical Society, Meeting, *Materials for Nonlinear Optics: Chemical Perspectives*, American Chemical Society, Washington, D.C., 1991.
13. D. L. Mills, *Nonlinear Optics: Basic Concepts*, 2nd enl. ed., Springer, New York, 1998.
14. A. C. Newell and J. V. Moloney, *Nonlinear Optics*, Addison-Wesley, Redwood City, California, 1992.
15. E. G. Sauter, *Nonlinear Optics*, John Wiley & Sons, New York, 1996.
16. M. Schubert and B. Wilhelmi, *Nonlinear Optics and Quantum Electronics*, John Wiley & Sons, New York, 1986.
17. Y. R. Shen, *The Principles of Nonlinear Optics*, John Wiley & Sons, New York, 1984.
18. R. L. Sutherland, *Handbook of Nonlinear Optics*, Marcel Dekker, New York, 1996.
19. P. Yeh, *Nonlinear Optics and Applications*, vol. 613, Society of Photo-Optical Instrumentation Engineers, Bellingham, Washington, 1986.
20. P. Yeh, *Introduction to Photorefractive Nonlinear Optics*, John Wiley & Sons, New York, 1993.

21. F. Zernike and J. E. Midwinter, *Applied Nonlinear Optics*, John Wiley & Sons, New York, 1973.
22. B. E. A. Saleh and M. C. Teich, *Fundamentals of Photonics*, John Wiley & Sons, New York, 1991.
23. A. Yariv, *Optical Electronics in Modern Communications*, 5th ed., Oxford University Press, New York, 1997.
24. A. Yariv, *Quantum Electronics*, 3rd ed., John Wiley & Sons, New York, 1989.
25. A. Yariv and P. Yeh, *Optical Waves in Crystals: Propagation and Control of Laser Radiation*, John Wiley & Sons, New York, 1984.
26. R. A. Fisher, *Optical Phase Conjugation*, Academic Press, New York, 1983.
27. J.-C. Diels and W. Rudolph, *Ultrashort Laser Pulse Phenomena: Fundamentals, Techniques, and Applications on a Femtosecond Time Scale*, Academic Press, San Diego, 1996.
28. B. I. A. Zel'dovich, N. F. Pilipetskii, and V. V. Shkunov, *Principles of Phase Conjugation*, Springer-Verlag, Berlin and New York, 1985.
29. B. S. Wherrett and P. G. Harper, *Nonlinear Optics*, Academic Press, London, 1977.
30. E. Garmire and A. Kost, "Nonlinear Optics in Semiconductors I & II," in *Semiconductors and Semimetals*, R. K. Wilardson and E. R. Webber, eds., Academic Press, 1999.
31. M. Sheik-Bahae and E. W. Van Stryland, "Optical Nonlinearities in the Transparency Region of Bulk Semiconductors," in *Nonlinear Optics in Semiconductors I*, vol. 58, *Semiconductors and Semimetals*, E. Garmire and A. Kost, eds., Academic Press, 1999.
32. E. W. Van Stryland, A. L. Smirl, T. F. Boggess, M. J. Soileau, B. S. Wherrett, and F. Hopf, "Weak-Wave Retardation and Phase-Conjugate Self-Defocusing in Si," in *Picosecond Phenomena III*, R. M. H. K. B. Eisenthal, W. Kaiser, and A. Laubereau, eds., Springer-Verlag, Berlin, 1982.
33. J. S. Toll, "Casuality and Dispersion Relation: Logical Foundation," *Phys. Rev.* **104**:1760–1770 (1956).
34. H. M. Nussenzveig, *Causality and Dispersion Relations*, Academic Press, New York, 1972.
35. P. J. Caspers, "Dispersion Relations for Nonlinear Response," *Phys. Rev. A.* **133**:1249 (1964).
36. S. M. Kogan, "On the Electromagnetics of Weakly Nonlinear Media," *Sov. Phys. JETP* **16**:217 (1963).
37. P. J. Price, "Theory of Quadratic Response Functions," *Phys. Rev.* **130**:1792 (1964).
38. F. L. J. Ridener and R. H. J. Good, "Dispersion Relations for Nonlinear Systems or Arbitrary Degree," *Phys. Rev. B* **11**:2768 (1975).
39. D. C. Hutchings, M. Sheik-Bahae, D. J. Hagan, and E. W. Van Stryland, "Kramers-Kronig Relations in Nonlinear Optics," *Optical and Quantum Electronics* **24**:1–30 (1992).
40. F. Bassani and S. Scandolo, "Dispersion-Relations and Sum-Rules in Nonlinear Optics," *Phys. Rev. B-Condensed Matter* **44**:8446–8453 (1991).
41. D. A. B. Miller, C. T. Seaton, M. E. Prise, and S. D. Smith, "Band-Gap-Resonant Non-Linear Refraction in Iii-V Semiconductors," *Phys. Rev. Lett.* **47**:197–200 (1981).
42. B. S. Wherrett, A. C. Walker, and F. A. P. Tooley, "Nonlinear Refraction for CW Optical Bistability," in *Optical Nonlinearities and Instabilities in Semiconductors*, H. Haug, ed., Academic Press, Boston, 1988.
43. M. Sheik-Bahae, D. J. Hagan, and E. W. Van Stryland, "Dispersion and Band-Gap Scaling of the Electronic Kerr Effect in Solids Associated with 2-Photon Absorption," *Phys. Rev. Lett.* **65**:96–99 (1990).
44. M. Sheik-Bahae, D. C. Hutchings, D. J. Hagan, and E. W. Van Stryland, "Dispersion of Bound Electronic Nonlinear Refraction in Solids," *IEEE J. Quantum Electron.* **27**:1296–1309 (1991).
45. M. Sheik-Bahae, "Nonlinear Optics of Bound Electrons in Solids," in *Nonlinear Optical Materials*, J. V. Moloney, ed., Springer, New York, 1998.
46. D. J. Moss, E. Ghahramani, J. E. Sipe, and H. M. Van Driel, "Band-Structure Calculation of Dispersion and Anisotropy in $\chi^{(3)}$ For 3rd-Harmonic Generation in Si, Ge, and GaAs," *Phys. Rev. B* **41**:1542–1560 (1990).
47. R. Adair, L. L. Chase, and S. A. Payne, "Nonlinear Refractive-Index Measurements of Glasses Using 3-Wave Frequency Mixing," *Journal of the Optical Society of America B-Optical Physics* **4**:875–881 (1987).
48. R. Adair, L. L. Chase, and S. A. Payne, "Nonlinear Refractive-Index of Optical-Crystals," *Phys. Rev. B-Condensed Matter* **39**:3337–3350 (1989).
49. S. R. Friberg, A. M. Weiner, Y. Silberberg, B. G. Sfez, and P. S. Smith, "Femtosecond Switching in a Dual-Core-Fiber Nonlinear Coupler," *Opt. Lett.* **13**:904–906 (1988).
50. M. J. Weber, D. Milam, and W. L. Smith, "Non-Linear Refractive-Index of Glasses and Crystals," *Optical Engineering* **17**:463–469 (1978).

51. L. F. Mollenauer, R. H. Stolen, and J. P. Gordon, *Phys. Rev. Lett.* **45**:1095 (1980).
52. D. E. Spence, P. N. Kean, and W. Sibbett, "60-Fsec Pulse Generation from a Self-Mode-Locked Ti-Sapphire Laser," *Opt. Lett.* **16**:42–44 (1991).
53. H. M. Gibbs, *Optical Bistability: Controlling Light with Light*, Academic Press, Orlando, Florida, 1985.
54. M. N. Islam, *Ultrafast Fiber Switching Devices and Systems*, Cambridge University Press, Cambridge and New York, 1992.
55. N. L. Boling, A. J. Glass, and A. Owyong, "Empirical Relationships for Predicting Nonlinear Refractive Index Changes in Optical Solids," *IEEE J. Quantum Electron.* **QE-14**:601 (1978).
56. P. P. Ho and R. R. Alfano, "Optical Kerr Effect in Liquids," *Phys. Rev. A* **20**:2170 (1979).
57. F. A. Hopf and G. I. Stegeman, *Applied Classical Electrodynamics*, John Wiley & Sons, New York, 1985.
58. J. F. Reintjes, *Nonlinear Optical Parametric Processes in Liquids and Gases*, Academic Press, New York, 1984.
59. M. D. Levenson, *Introduction to Nonlinear Laser Spectroscopy*, Academic Press, New York, 1982.
60. W. Kaiser and D. H. Auston, *Ultrashort Laser Pulses: Generation and Applications*, 2nd ed., Springer, Berlin and New York, 1993.
61. M. Cardona and G. Güntherodt, *Light Scattering in Solids VI: Recent Results, Including High-Tc Superconductivity*, Springer-Verlag, Berlin and New York, 1991.
62. J. J. Valentini, "Laser Raman Techniques," in *Optical Engineering: Laser Spectroscopy and Its Applications*, vol. 11, L. J. Radziemski, R. W. Solarz, and J. A. Paisner, eds., Dekker, New York, 1987.
63. D. A. Long, *Raman Spectroscopy*, McGraw-Hill, New York, 1977.
64. G. L. Eesley, *Cohernet Raman Spectroscopy*, Pergamon, Oxford, 1981.
65. L. A. Woodward, *Raman Spectroscopy: Theory and Practice*, H. A. Szymanski, ed., Plenum, New York, 1987.
66. E. W. Van Stryland and L. Chase, "Two-Photon Absorption: Inorganic Materials," in *Handbook of Laser Science and Technology, Supplement 2, Optical Materials*, M. Weber, ed., CRC Press, 1994.
67. E. W. Van Stryland, H. Vanherzeele, M. A. Woodall, M. J. Soileau, A. L. Smirl, S. Guha, and T. F. Boggess, "2 Photon-Absorption; Nonlinear Refraction; and Optical Limiting in Semiconductors," *Optical Engineering*, **24**:613–623 (1985).
68. E. W. Van Stryland, M. A. Woodall, H. Vanherzeele, and M. J. Soileau, "Energy Band-Gap Dependence of 2-Photon Absorption," *Opt. Lett.* **10**:490–492 (1985).
69. B. S. Wherrett, "Scaling Rules for Multiphoton Interband Absorption in Semiconductors," *Journal of the Optical Society of America B-Optical Physics* **1**:67–72 (1984).
70. J. I. Pankove, *Optical Processes in Semiconductors*, Dover, New York, 1971.
71. C. Kittel, *Introduction to Solid State Physics*, 7th ed., John Wiley & Sons, New York, 1996.
72. A. Miller, D. A. B. Miller, and S. D. Smith, "Dynamic Non-Linear Optical Processes in Semiconductors," *Advances in Physics* **30**:697–800 (1981).
73. H. Haug and S. W. Koch, *Quantum Theory of the Optical and Electronic Properties of Semiconductors*, 3rd ed., World Scientific, Singapore, 1994.
74. W. W. Chow, S. W. Koch, and M. Sargent, *Semiconductor-Laser Physics*, corrected printing ed., Springer-Verlag, Berlin and New York, 1997.
75. L. Banyai and S. W. Koch, "A Simple Theory for the Effects of Plasma Screening on the Optical Spectra of Highly Excited Semiconductors," *Zeitschrift Fur Physik B-Condensed Matter* **63**:283–291 (1986).
76. R. K. Jain and M. B. Klein, "Degenerate Four-Wave Mixing in Semiconductors," in *Optical Phase Conjugation*, R. A. Fisher, ed., Academic Press, New York, 1983.
77. J. I. Steinfeld, *Laser and Coherence Spectroscopy*, Plenum Press, New York, 1978.
78. M. Sargent, M. O. Scully, and W. E. Lamb, *Laser Physics*, Addison-Wesley, Reading, Massachusetts and London, 1977.
79. J. Feinberg, "Optical Phase Conjugation in Photorefractive Materials," in *Optical Phase Conjugation*, R. A. Fisher, ed., Academic Press, New York, 1983.
80. J. E. Millard, M. Ziari, and A. Partovi, "Photorefractivity in Semiconductors," in *Nonlinear Optics in Semiconductors I*, vol. 58, *Semiconductors and Semimetals*, E. Garmire and A. Kost, eds., Academic Press, 1999.

81. R. DeSalvo, D. J. Hagan, M. Sheik-Bahae, G. Stegeman, E. W. Van Stryland, and H. Vanherzeele, "Self-Focusing and Self-Defocusing by Cascaded 2nd-Order Effects in Ktp," *Opt. Lett.* **17**:28–30 (1992).
82. N. R. Belashenkov, S. V. Gagarskii, and M. V. Inochkin, "Nonlinear Refraction of Light on Second-Harmonic Generation," *Opt. Spektrosk.* **66**:1383–1386 (1989).
83. G. I. Stegeman, M. Sheik-Bahae, E. Van Stryland, and G. Assanto, "Large Nonlinear Phase-Shifts in 2nd-Order Nonlinear-Optical Processes," *Opt. Lett.* **18**:13–15 (1993).
84. E. W. Van Stryland, "Third-Order and Cascaded Nonlinearities," presented at Laser Sources and Applications, 1996.
85. G. Assanto, G. Stegeman, M. Sheik-Bahae, and E. Van Stryland, "All-Optical Switching Devices Based on Large Nonlinear Phase-Shifts from 2nd Harmonic-Generation," *Appl. Phys. Lett.* **62**:1323–1325 (1993).
86. G. Assanto, G. I. Stegeman, M. Sheik-Bahae, and E. Van Stryland, "Coherent Interactions for All-Optical Signal-Processing via Quadratic Nonlinearities," *IEEE J. Quantum Electron.* **31**:673–681 (1995).
87. W. Torruellas, B. Lawrence, and G. I. Stegeman, "Self-Focusing and 2d Spatial Solutions in Pts," *Electronics Letters* **32**:2092–2094 (1996).
88. W. E. Torruellas, Z. Wang, D. J. Hagan, E. W. Vanstryland, G. I. Stegeman, L. Torner, and C. R. Menyuk, "Observation of 2-Dimensional Spatial Solitary Waves in a Quadratic Medium," *Phys. Rev. Lett.* **74**:5036–5039, 1995.
89. L. J. Qian, X. Liu, and F. W. Wise, "Femtosecond Kerr-Lens Mode Locking with Negative Nonlinear Phase Shifts," *Opt. Lett.* **24**:166–168, 1999.
90. M. Sheik-Bahae, A. A. Said, D. J. Hagan, M. J. Soileau, and E. W. Van Stryland, "Nonlinear Refraction and Optical Limiting in Thick Media," *Optical Engineering* **30**:1228–1235 (1991).
91. A. E. Kaplan, "External Self-Focusing of Light by a Nonlinear Layer," *Radiophys. Quantum Electron.* **12**:692–696 (1969).
92. J. H. Marburger, in *Progress in Quantum Electronics*, J. H. Sanders and S. Stenholm, eds., Pergamon Press, New York, 1977.
93. H. A. Haus, J. G. Fujimoto, and E. P. Ippen, "Analytic Theory of Additive Pulse and Kerr Lens Mode-Locking," *IEEE J. Quantum Electron.* **28**:2086–2096 (1992).
94. G. P. Agrawal, *Nonlinear Fiber Optics*. Academic Press, Boston, 1989.
95. A. Barthelemy, S. Maneuf, and C. Froehly, "Soliton Propagation and Self-Confinement of Laser-Beams by Kerr Optical Non-Linearity," *Optics Communications* **55**:201–206 (1985).
96. S. Maneuf, A. Barthelemy, and C. Froehly, "Soliton Beam Propagation: Space-Time Behavior and Spectral Features," *Journal of Optics* **17**:139–145 (1986).
97. J. J. Hopfield, J. M. Worlock, and K. Park, "Two-Quantum Absorption Spectrum of KI," *Phys. Rev. Lett.* **11**:414 (1963).
98. J. A. Bolger, A. K. Kar, B. S. Wherrett, R. Desalvo, D. C. Hutchings, and D. J. Hagan, "Nondegenerate 2-Photon Absorption-Spectra of Znse; Zns and Zno," *Optics Communications* **97**:203–209 (1993).
99. P. Maker, R. Terhune, and C. Savage, "Intensity-Dependent Changes in the Refractive Index of Liquids," *Phys. Rev. Lett.* **12**:507 (1964).
100. D. J. Hagan, E. Miesak, R. Negres, S. Ross, J. Lim, and E. W. Van Stryland, "Nonlinear Spectrometry of Chromophores for Optical Limiting," *SPIE Proceedings* **3472**:80–90 (1998).
101. R. L. Abrams, J. F. Lam, R. C. Lind, D. G. Steel, and P. F. Liao, "Phase Conjugation and High Resolution Spectroscopy by Resonant Degenerate Four-Wave Mixing," in *Optical Phase Conjugation*, R. A. Fisher, ed., Academic Press, New York, 1983.
102. J. Shah, *Ultrafast Spectroscopy of Semiconductors and Semiconductor Nanostructures*, Springer, Berlin and New York, 1996.
103. M. J. Weber, D. Milam, and W. L. Smith, "Nonlinear Refractive Index of Glasses and Crystals," *Opt. Eng.* **463** (1978).
104. M. J. Moran, C. Y. She, and R. L. Carmen, "Interferometric Measurements of the Nonlinear Refractive Index Coefficient Relative to CS₂ in the Laser System Related Materials," *IEEE J. Quantum Electron.* **11**:259 (1975).
105. M. Sheik-Bahae, A. A. Said, T. H. Wei, D. J. Hagan, and E. W. Van Stryland, "Sensitive Measurement of Optical Nonlinearities Using a Single Beam," *IEEE J. Quantum Electron.* **26**:760–769 (1990).

106. E. W. Van Stryland and M. Sheik-Bahae, "Z-Scan," in *Characterization Techniques and Tabulations for Organic Nonlinear Optical Materials*, M. G. Kuzyk and C. W. Dirk, eds., Marcel Dekker, New York, 1998.
107. T. Xia, D. J. Hagan, M. Sheik-Bahae, and E. W. Van Stryland, "Eclipsing Z-Scan Measurement of $\Lambda/10(4)$ Wave-Front Distortion," *Opt. Lett.* **19**:317–319 (1994).
108. M. Sheik-Bahae, A. A. Said, and E. W. Van Stryland, "High-Sensitivity; Single-Beam N2 Measurements," *Opt. Lett.* **14**:955–957 (1989).
109. G. Stegeman and E. Wright, "All-Optical Waveguide Switching," *Opt. and Quantum Electron.* **22**:95 (1990).
110. G. I. Stegeman, A. Villeneuve, J. Kang, J. S. Aitchison, C. N. Ironside, K. Alhemyari, C. C. Yang, C. H. Lin, H. H. Lin, G. T. Kennedy, R. S. Grant, and W. Sibbett, "AlGaAs Below Half Bandgap: The Silicon of Nonlinear-Optical Materials," *International Journal of Nonlinear Optical Physics* **3**:347–371 (1994).

This page intentionally left blank.

DO NOT DUPLICATE

CONTINUOUS-WAVE OPTICAL PARAMETRIC OSCILLATORS

Majid Ebrahim-Zadeh

*ICFO—Institut de Ciències Fòniques
Mediterranean Technology Park
Barcelona, Spain, and
Institutio Catalana de Recerca i Estudis Avancats (ICREA)
Passeig Lluís Companys
Barcelona, Spain*

17.1 INTRODUCTION

Since the publication of an earlier review on optical parametric oscillators (OPOs) in 2000,¹ there has been remarkable progress in the technological development and applications of OPO devices. Once considered an impractical approach for the generation of coherent radiation, OPOs have now been finally transformed into truly viable, state-of-the-art light sources capable of accessing difficult spectral regions and addressing real applications beyond the reach of conventional lasers. While the first experimental demonstration of an OPO was reported in 1965,² for nearly two decades thereafter there was little or no progress in the practical development of OPO devices, owing to the absence of suitable nonlinear materials and laser pump sources. With the advent of a new generation of birefringent nonlinear crystals, most notably β -BaB₂O₄ (BBO), LiB₃O₅ (LBO), and KTiOPO₄ (KTP), but also KTiOAsO₄ (KTA) and RbTiOAsO₄ (RTA) in the mid-1980s, and advances in solid-state laser technology, there began a resurgence of interest in OPOs as potential alternatives to conventional lasers for the generation of coherent radiation in new spectral regions. The high optical damage threshold, moderate optical nonlinearity, and favorable phase-matching properties of the newfound materials led to important breakthroughs in OPO technology. In the years to follow, tremendous progress was achieved in the development of OPO devices, particularly in the pulsed regime, and a variety of OPO systems from the nanosecond to the ultrafast picosecond and femtosecond timescales, and operating from the near-ultraviolet (near-UV) to the infrared (IR) were rapidly developed. These developments led to the availability of a wide range of practical OPO devices and their deployment in new applications, with some systems finding their way to the commercial market. A decade later, in the mid-1990s, the emergence of quasi-phase-matched (QPM) ferroelectric nonlinear crystals, particularly periodically poled LiNbO₃ (PPLN) stimulated new impetus for the advancement of continuous-wave (cw) OPO devices, traditionally the most challenging regime for OPO operation due to almost negligible nonlinear gains available under cw pumping. The flexibility offered by grating-engineered QPM materials, allowing access to the highest nonlinear tensor coefficients, combined with noncritical phase matching (NCPM) and long interaction lengths (>50 mm in PPLN), enabled the low available nonlinear gains to be overcome, hence permitting the development of practical cw OPOs in a variety of resonance configurations. As such, the advent of QPM

materials, most notably PPLN, but also periodically poled KTP (PPKTP), RbTiOAsO₄ (PPRTA), and LiTaO₃ (PPLT), has had an unparalleled impact on cw OPO technology. Combined with advances in novel high-power solid-state crystalline, semiconductor, and fiber pump sources over the past decade, these developments have led to the practical realization of a new class of cw OPOs with previously unattainable performance capabilities with regard to wavelength coverage, output power and efficiency, frequency and power stability, spectral and spatial coherence, and fine frequency tuning.

With their exceptional spectral coverage and tuning versatility, temporal flexibility from the cw to femtosecond timescales, practical performance parameters, and compact solid-state design, OPO devices have now been firmly established as truly competitive alternatives to conventional lasers and other technologies for the generation of widely tunable coherent radiation in difficult spectral and temporal domains. In the current state of technology, OPO devices can provide spectral access from ~400 nm in the ultraviolet (UV) to ~12 μm in the mid-infrared (mid-IR), as well as the terahertz (THz) spectral region. They can also provide temporal output from the cw and long-pulse microsecond regime to nanosecond, picosecond, and ultrafast sub-20 fs timescales. Many of the developed OPO systems are now routinely deployed in a variety of applications including spectroscopy, optical microscopy, environmental trace gas detection and monitoring, life sciences, biomedicine, optical frequency metrology and synthesis, and imaging.

The aim of this chapter is to provide an overview of the advances in OPO device technology and applications since the publication of the earlier review in 2000.¹ The chapter is concerned only with the developments after 2000, since many of the important advances in this area prior to that date can already be found in the previous treatment¹ as well as other reviews on the subject.^{3–10} Because of limited scope, and given that most of the important advances over the last decade have been in the CW operating regime, the chapter is focused only on a discussion of cw OPOs. Reviews on pulsed and ultrafast OPOs can be found elsewhere.^{3,4,6–10} This chapter also does not include a description of the fundamental concepts in nonlinear and crystal optics, parametric generation, amplification and gain, or a comprehensive description of the design criteria and operating principles of OPO devices, which have been the subject of several earlier treatments.^{11–16}

17.2 CONTINUOUS-WAVE OPTICAL PARAMETRIC OSCILLATORS

Of the different types of OPO devices developed to date, advancement of practical OPOs in the cw operating regime has been traditionally most difficult, since the substantially lower nonlinear gains available under cw pumping necessitate the use of high-power cw pump laser or the deployment of multiple-resonant cavities to reach operation threshold. As in a conventional laser oscillator, the OPO is characterised by a threshold condition, defined by the pumping intensity at which the growth of the parametric waves in one round-trip of the optical cavity just balances the total loss in that round-trip. Once threshold has been surpassed, coherent light at macroscopic levels can be extracted from the oscillator. In order to provide feedback in an OPO, a variety of resonance schemes may be deployed by suitable choice of mirrors forming the optical cavity, as illustrated in Fig. 1a to e. The mirrors may be highly reflecting at only one of the parametric waves (*signal* or *idler*, but not both), as in Fig. 1a, in which case the device is known as a *singly resonant oscillator* (SRO). This configuration is characterised by the highest cw operation threshold. In order to reduce threshold, alternative resonator schemes may be employed where additional optical waves are resonated in the optical cavity. These include the *doubly resonant oscillator* (DRO), Fig. 1b, in which both the signal and idler waves are resonant in the optical cavity, and the *pump-resonant* or *pump-enhanced* SRO, Fig. 1c, where the *pump* as well as one of the generated waves (signal or idler) is resonated. In an alternative scheme, Fig. 1d, the pump may be resonated together with both parametric waves, in which case the device is known as a *triply resonant oscillator* (TRO). Such schemes can bring about substantial reductions in threshold from the cw SRO configuration, with the TRO offering the lowest operation threshold. In an alternative scheme, the external pump power threshold for a cw SRO may also be substantially reduced by deploying internal pumping, where the OPO is placed inside

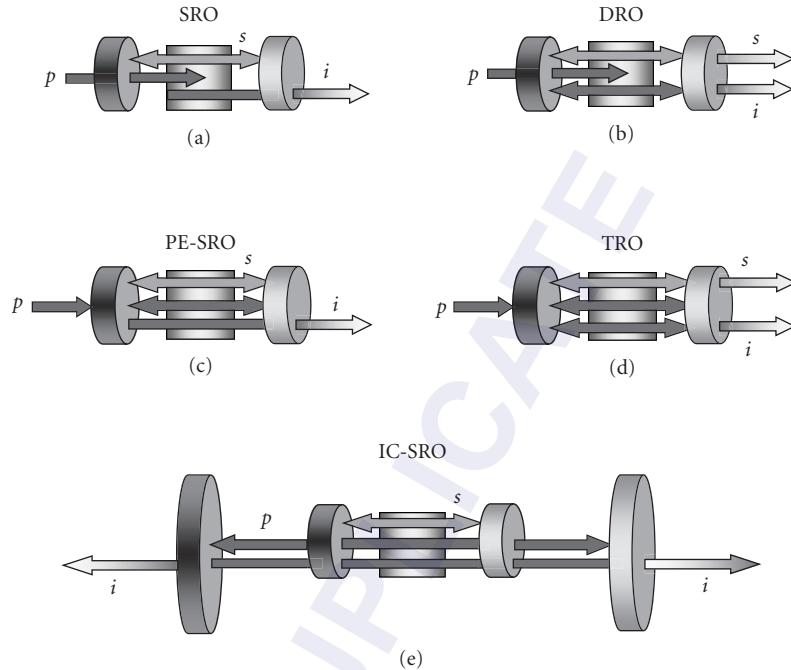


FIGURE 1 Cavity resonance configurations for cw OPOs. The symbols p , s , and i denote pump, signal, and idler, respectively.

a minimally output-coupled pump laser. A schematic of such an *intracavity SRO (IC-SRO)* is illustrated in Fig. 1e.

The comparison of steady-state threshold for conventional externally pumped cw OPOs under different resonance schemes is shown in Fig. 2, where the calculated external pump power threshold is plotted as a function of the effective nonlinear coefficient of several materials including LBO, KTA, KTP, KNbO_3 , PPLN, and PPRTA. From the plot, it is clear that for the majority of birefringent materials the attainment of cw SRO threshold requires pump powers on the order of tens of watts, well outside the range of the most widely available cw laser sources. However, in the case of PPLN, the cw SRO threshold is substantially reduced to acceptable levels below ~ 1 W, bringing operation of cw SROs within the convenient range of widespread cw solid-state pump lasers. With the cw PE-SRO, considerably lower thresholds can be achieved, from a few hundred milliwatts to ~ 1 W for birefringent materials and below ~ 100 mW for PPLN. In the case of cw DRO, still lower thresholds of the order of 100 mW are attainable with birefringent materials, with only a few milliwatts for PPLN, whereas with the cw TRO, thresholds from below 1 mW to a few milliwatts can be obtained in birefringent materials.

It is thus clear that practical operation of cw OPOs in SRO configurations is generally beyond the reach of birefringent materials, but requires DRO, TRO, and PE-SRO cavities. On the other hand, implementation of cw SROs necessitates the use of PPLN or similar QPM materials, offering enhanced optical nonlinearities, and long interaction lengths under NCPM. However, the threshold reduction from SRO to PE-SRO, DRO, and TRO cavity configurations is often achieved at the expense of increased spectral and power instability in the OPO output arising from the difficulty in maintaining resonance for more than one optical wave in a single optical cavity. For this reason, the cw SRO offers the most direct route to the attainment of high output stability and spectral control without stringent demands on the frequency stability of the laser pump source. On the other hand,

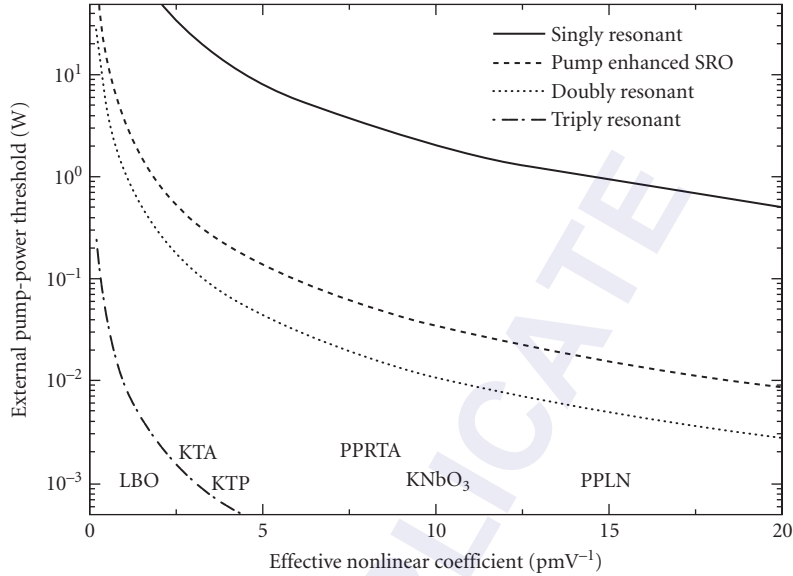


FIGURE 2 Calculated minimum thresholds for different OPO resonance configurations versus the effective nonlinear coefficients in various nonlinear materials. The calculation assumes confocal focusing and loss values that are typically encountered in experimental cw OPOs, the finesse representing round-trip power losses of approximately 2.0 percent. The plots correspond to a pump wavelength of 800 nm, degenerate operation, a pump refractive index of ~ 1.7 , a crystal length of 20 mm, signal and idler cavity finesse of ~ 300 , and a pump enhancement factor of ~ 30 . In the case of PE-SRO and TRO, the enhancement factor of 30 represents the maximum enhancement attainable with parasitic losses of ~ 3 percent at the pump.¹⁷

practical implementation of cw PE-SRO, DRO, and TRO requires active stabilization techniques to control output power and frequency stability, with the PE-SRO offering the most robust configuration for active stabilization and TRO representing the most difficult in practice. In addition, practical operation of OPOs in multiple resonant cavities can only be achieved using stable, single-frequency pump lasers and such devices also require more complex protocols for frequency tuning and control than the cw SRO. More detailed description of the different resonance and pumping schemes for OPOs and analytical treatment of tuning mechanisms, spectral behavior, frequency control, and stabilisation can be found in an earlier review.¹

Singly Resonant Oscillators

By deploying the intracavity pumping scheme using a Ti:sapphire laser in combination with a 20-mm PPKTP crystal, Edwards et al.¹⁸ reported a cw IC-SRO capable of providing up to 455 mW of non-resonant infrared idler power at a down-conversion efficiency of 87 percent. Using a combination of pump tuning at room temperature and crystal temperature tuning, idler (signal) coverage in the 2.23 to 2.73 μm (1.14 to 1.27 μm) spectral ranges was demonstrated. By configuring the Ti:sapphire pump laser and the SRO in ring cavity geometries and using intracavity etalons, 115 mW of unidirectional, single-frequency idler power was generated at 2.35 μm with mode-hop-free operating time intervals of about 10 s under free-running conditions. The resonant signal was measured to have a linewidth < 15 MHz for a pump linewidth < 25 MHz.

The advent of PPLN with large effective nonlinearity ($d_{\text{eff}} \sim 15$ pm/V) and long interaction lengths (currently up to 80 mm) under NCPM has enabled the development of cw SROs in conventional

external pumping configurations using more commonly available, moderate- to high-power solid-state pump sources. By deploying a fixed-frequency, cw, single-mode Nd:YVO₄ pump laser at 1.064 μm , Bisson et al.¹⁹ developed a portable source for mid-IR photoacoustic spectroscopy based on a PPLN cw SRO by using discrete mode-hop tuning of the idler. The SRO, based on a 50-mm PPLN crystal with fanned grating ($\Lambda = 29.3$ to $30.1 \mu\text{m}$), was configured in a ring cavity, and frequency selection and fine tuning was implemented using solid or air-spaced intracavity etalons. With an uncoated, 400- μm -thick, solid Nd:YAG etalon, a total mode-hop-tuning range of $\sim 4 \text{ cm}^{-1}$ for the idler in discrete steps of 0.02 to 0.1 cm^{-1} was achieved by rotation of the etalon. The SRO could deliver a maximum idler power of $\sim 120 \text{ mW}$ at a pump depletion of 40 to 50 percent for 6 W of pump power. Using the mode-hop-tuned idler output near 3.3 μm , photoacoustic spectroscopy of the methane Q branch was performed at atmospheric pressure by simultaneous tuning of the PPLN crystal combined with etalon rotation. A total of four etalon scans covering $\sim 10 \text{ cm}^{-1}$ was necessary to trace the Q branch spectrum. In an effort to achieve a constant tuning rate as well as minimize insertion loss due to etalon rotation, which in turn leads to mode hops arising from variable heating of the PPLN crystal due to the changes in intracavity power, an alternative air-spaced fused silica etalon with ~ 0.5 to 1.5 mm spacing and ~ 5 percent reflectivity at the signal ($\sim 1.57 \mu\text{m}$) was also employed in the present device. While resulting in a higher oscillation threshold ($\sim 4 \text{ W}$) and lower idler output ($\sim 80 \text{ mW}$), the combination of PPLN tuning and piezoelectric scan of the etalon over a distance of 3 μm (at 1.5 mm separation) yielded a total mode-hop tuning range of $\sim 14 \text{ cm}^{-1}$ for the idler at a constant tuning rate and in discrete steps of 0.1 cm^{-1} , providing sufficient resolution for atmospheric sensing and pressure-broadened spectroscopy. The measured linewidth of the idler was $< 10 \text{ MHz}$ with a passive stability of $\sim 50 \text{ MHz}$ over 30 s.

By using a 10-W cw single-frequency diode-pumped Nd:YAG laser at 1.064 μm , Van Herpen et al.²⁰ demonstrated a cw SRO based on PPLN with a mid-IR idler tuning range of 3.0 to 3.8 μm . The SRO, configured in a ring cavity and using a crystal with fanned grating ($L = 50 \text{ mm}$, $\Lambda = 29.3$ to $30.1 \mu\text{m}$) exhibited a pump power threshold of $\sim 3 \text{ W}$ and could provide a maximum idler output power of 1.5 W at 3.3 μm for 9 W of pump power. The combination of the single-mode pump laser, a ring cavity for the SRO, and the inclusion of an intracavity air-spaced etalon enabled mode-hop-free tuning of the idler over 12 GHz by tuning the pump frequency over 24 GHz, with the idler mode-hop tuning range limited by mode hopping in the pump laser. Under this condition, 700 mW of single-frequency, smoothly tunable idler power could be provided by the SRO. In a later experiment,²¹ using the same PPLN crystal and identical cavity design for the SRO, the authors were able to improve the idler output power in the 3.0 to 3.8 μm range by increasing the available Nd:YAG pump power to 15 W and by optimizing pump focusing and the SRO cavity length. The SRO similarly exhibited a cw power threshold of $\sim 3 \text{ W}$, but could provide 2.2 W of idler power for 10.5 W of input pump power. The coarse and fine tuning properties of this SRO were similar to the earlier device. For fine tuning, an intracavity air-spaced etalon with variable spacing of 0.2 to 3 mm (FSR = 50 to 750 GHz) was used. Continuous scanning of etalon spacing resulted in discrete mode-hop tuning of the idler over 100 GHz. With a 400- μm uncoated solid YAG etalon (FSR = 207 GHz), an idler mode-hop tuning range of 10 cm^{-1} in steps of 0.02 to 0.1 cm^{-1} (0.6 to 3 GHz) could be obtained by rotation of the etalon. Subsequently, using the same pump laser, the authors reported a cw SRO based on a multigrating PPLN crystal ($\Lambda = 25.9$ to $28.7 \mu\text{m}$) and providing extended idler coverage into the 3.7 to 4.7 μm spectral range in the mid-IR.²² The ring-cavity SRO exhibited an oscillation threshold of between 5 and 7.5 W over this spectral range and for an input pump power of 11 W could provide a maximum idler output of 1.2 W at 3.9 μm , decreasing to 120 mW at 4.7 μm . The increase in SRO threshold and corresponding decrease in output power were attributed to the increasing idler absorption in PPLN at longer wavelengths toward 5 μm . With the inclusion of the same 400- μm uncoated YAG etalon to stabilize the resonant signal frequency, continuous mode-hop free tuning of the idler was achieved by tuning the pump frequency over 24 GHz, but with a reduction in idler power by as much as 50 percent. Discontinuous mode-hop tuning of the idler output could also be obtained through rotation of the intracavity etalon. In a later report, the use of a tunable high-power ($> 20 \text{ W}$) diode-pumped Yb:YAG laser in combination with two PPLN crystals with fanned gratings ($\Lambda = 28.5$ to $29.9 \mu\text{m}$) and two sets of OPO mirrors enabled the generation of widely tunable idler radiation with a total tuning range of 2.6 to 4.66 μm , and at increased

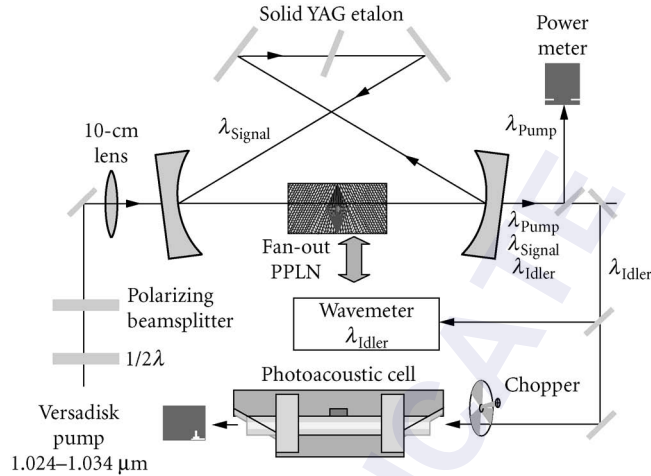


FIGURE 3 Experimental setup of the cw SRO. The pump wavelength varies from 1024 to 1034 nm and the idler wavelength from 2.6 to 4.7 μm . A pump rejecter mirror separates the pump light from the idler and signal beams, after which the idler beam is reflected toward the wavemeter and photoacoustic cell. The signal wavelength can be measured with the same wavemeter by replacing the idler reflector with a signal reflector.²³

cw power levels up to 3 W.²³ For frequency stability, a 400- μm uncoated YAG etalon (FSR = 207 GHz) was similarly used internal to the SRO cavity (Fig. 3). The SRO had a threshold of 8 W and, with nonoptimized mirror and crystal coatings, could provide 3.0 W of mid-IR idler output at 2.954 μm for 18 W of pump power. The SRO could provide an idler mode-hop tuning range of 25 GHz in steps of 100 MHz (FSR of the pump laser cavity) by tuning the intracavity pump etalon. Combined with the tuning of the Lyot filter within the pump laser, a total mode-hop tuning range of 190 GHz could be scanned, limited by a mode hop in signal frequency of 207 GHz corresponding to the FSR of the YAG etalon within the SRO cavity. By recording the photoacoustic signal in ethane, the authors characterized the frequency stability of the SRO. Due to unoptimized coatings, the idler exhibited frequency instabilities of 90 MHz/s, while temperature fluctuations in the PPLN crystal resulted in an idler frequency drift of 250 MHz over 200 s. In the same report, the authors demonstrated extension of the idler wavelength to 3.3 to 4.66 μm using the broad tuning of the pump laser (1.024 to 1.034 μm) in combination with grating tuning of the PPLN crystal, providing 200 mW of idler power at 4.235 μm , corresponding to the strongest CO_2 absorption line.

Subsequently, Ngai et al.²⁴ reported a cw SRO with automatic tuning control based on a multigrating MgO:PPLN ($L = 50$ mm, $\Lambda = 29.0$ to 31.5 μm). A schematic of the experimental setup is shown in Fig. 4. The SRO was pumped by a master oscillator-power amplifier (MOPA) laser at 1064 nm, providing 11.5 W of single-frequency output with a linewidth of 5 kHz (over 1 ms), frequency stability of 50 MHz/h, and continuous tuning over 48 GHz. The combination of temperature and grating tuning in the MgO:PPLN crystal provided coarse coverage over 2.75 to 3.83 μm in the idler and 1.47 to 1.73 μm in the signal, with a maximum idler power of 2.75 W. By using a ring SRO cavity containing a 400- μm -thick uncoated solid YAG etalon (FSR = 207 GHz), a short-term frequency stability of 4.5 MHz over 1 s was attainable in the absence of active stabilization. Fine wavelength scanning of the idler output was achieved through a combination of pump tuning, etalon rotation, and temperature tuning using an automated process with computer control. First, by continuous tuning of the pump frequency over 48 GHz at a fixed etalon angle, the idler could be tuned over 12 GHz before the occurrence of a mode hop in the pump laser (Fig. 5). The total idler tuning range attainable in this way was 207 GHz, limited by an etalon mode hop. Then, by rotation of the

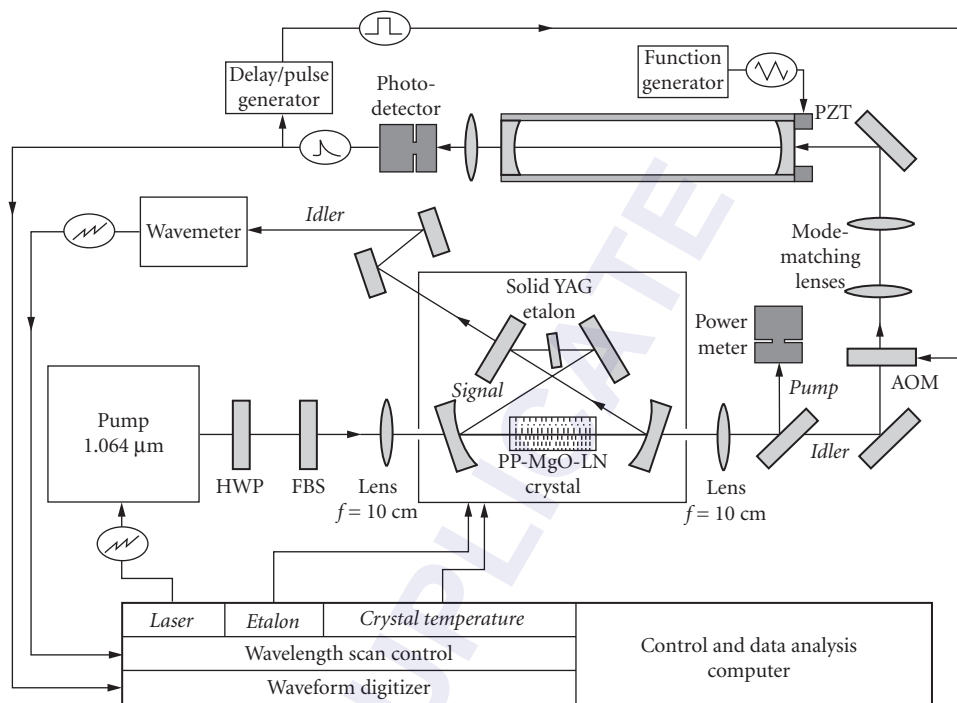


FIGURE 4 Experimental setup of automatically tunable cw SRO combined with continuous-wave cavity leak-out spectroscopy. The OPO cavity is resonant for the signal wavelength. The idler beam is sent to a cw leak-out cavity and to a wavemeter.²⁴

etalon to a new angle, the pump was again scanned until a new total tuning range of 207 GHz was covered, and process was repeated. Finally, changing the crystal temperature by 2 to 5°C, and repeating the entire process, wavelength scans of up to 450 cm^{-1} with a resolution of $<5 \times 10^{-4} \text{ cm}^{-1}$ could be obtained with a single grating period. Using this automated tuning process, the utility of the cw SRO for sensitive detection of CO_2 , methane, and ethane was demonstrated with photoacoustic and cavity leak-out spectroscopy, and analysis of human breath was performed by recording the absorption spectra of methane, ethane, and water in two test persons using photoacoustic spectroscopy.

With the continued advances in pump laser technology, the development of cw SROs based on high-power diode-pumped fiber lasers and amplifiers has also become a reality. Fiber lasers are attractive alternatives as pump sources for cw SROs, because they combine the high-power properties of crystalline solid-state laser materials with significant wavelength tuning and excellent spatial beam quality in compact and portable design. The pump tuning capability allows rapid and wide tuning of the SRO output without recourse to temperature or grating period variation, while the high available powers and excellent beam quality allow access to SRO threshold and enable the generation of practical output powers. The use of fiber pump lasers can thus provide a versatile class of cw SROs for the mid-IR that offer the advantages of simplicity, compact all-solid-state design, portability, reduced cost, improved functionality, and high output power and efficiency. Operation of a cw SRO pumped by a fiber laser was first reported by Gross et al.²⁵ using a tunable Yb-doped fiber laser. The laser delivered more than 8 W of cw output power in excellent spatial beam quality and was tunable over the wavelength range of 1031 to 1100 nm. With the use of a 40-mm-long multi-grating PPLN crystal and a ring cavity for the SRO, a cw idler output power of 1.9 W was generated at a wavelength of 3.2 μm in the mid-IR for 8.3 W of fiber pump power, with a corresponding SRO power threshold of 3.5 W. Idler wavelength tuning over 3.057 to 3.574 μm could be accomplished

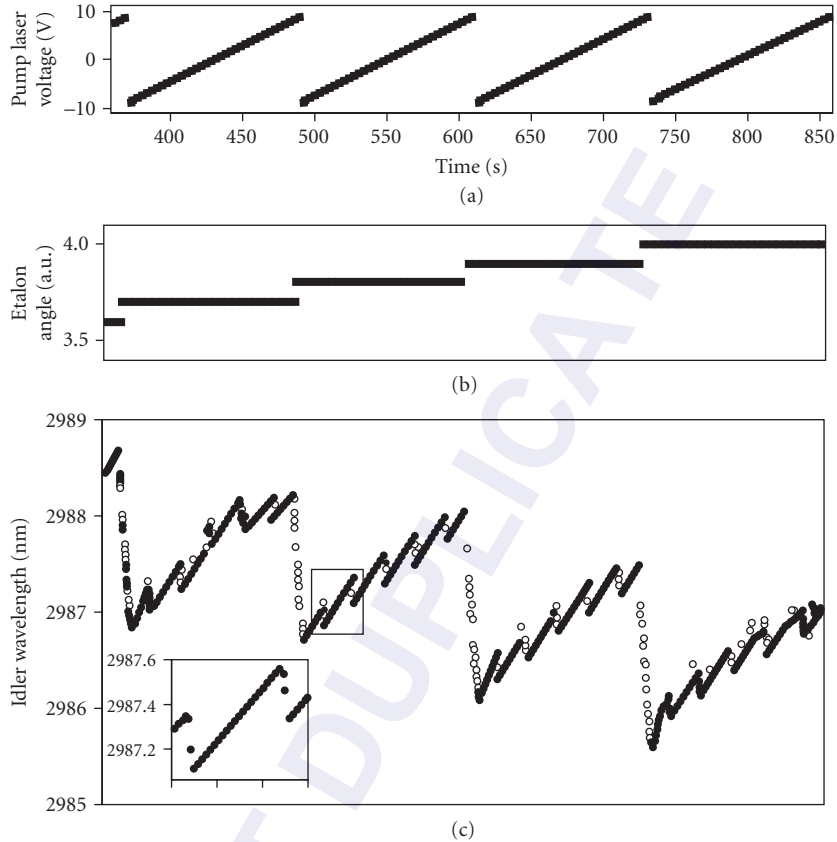


FIGURE 5 Combined pump-etalon scan. By scanning the pump laser (a) and stepping the etalon angle after each pump laser scan (b), a continuous wavelength coverage over 207 GHz can be realized (c). The resolution of the idler frequency is limited by the resolution of the wavemeter [inset in (c)].²⁴

by varying the crystal temperature or changing the grating period. However, wider and more convenient wavelength tuning was also available by exploiting the tuning capability of the fiber pump laser, where an idler tuning range of more than 700 nm over 2.980 to 3.700 μm was obtained by varying the pump wavelength between 1.032 and 1.095 μm . In a subsequent experiment, Klein et al.²⁶ demonstrated rapid wavelength tuning of a similar cw SRO by using electronic wavelength control of the Yb-doped fiber pump laser with an acousto-optic tunable filter. The SRO, based on a 40-mm-long single-grating PPLN crystal, was arranged in a similar ring cavity and, at a fixed crystal temperature and grating period, could be rapidly tuned over 3.160 to 3.500 μm in the idler wavelength by electronically tuning the fiber pump laser from 1060 to 1094 nm. The 340-nm idler tuning could be achieved within a time interval of 330 μs , representing a frequency tuning rate of 28 THz/ms. The overall electronic tuning range of the fiber pump laser over 1.057 to 1.100 μm resulted in an SRO idler tuning range of 437 nm in the mid-IR, from 3.132 to 3.569 μm . For the maximum fiber pump power of 6.6 W at 1.074 μm , the SRO generated an idler output power of 1.13 W at 3.200 μm .

More recently, operation of a low-threshold mid-IR cw SRO was reported by Henderson and Stafford²⁷ using MgO:PPLN and an all-fiber laser pump source. A schematic of the experimental setup is shown in Fig. 6. The cw single-frequency pump at 1083 nm was configured in a MOPA

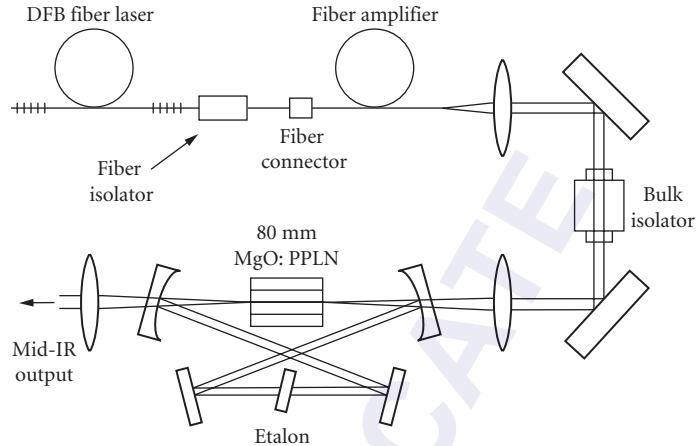


FIGURE 6 Schematic of experimental configuration for the fiber-pumped cw SRO.²⁷

arrangement using a 20-mW distributed feedback (DFB) fiber laser with 50-kHz linewidth as the seed and a polarization-maintaining fiber as the amplifier. The use of fiber connection between the two stages ensured an all-fiber configuration with no free-space components, alignment-free injection, and minimum long-term cavity misalignment. The MOPA could provide up to 3.5 W of amplified single-mode pump power for 20 mW of input seed power. Using multigrating and fanned crystals ($\Lambda = 31.3$ to $32.5 \mu\text{m}$) of 80-mm interaction length and operating the SRO just above room temperature (30°C), oscillation thresholds as low as 780 mW were obtained, with up to 750 mW of idler power generated for 2.8 W of fiber pump power. The idler output was tunable over 2650 to 3200 nm with a near-diffraction-limited spatial mode up to 500 mW and beam quality factor $M^2 = 1.04$. By exploiting the tunability of the pump laser through application of a voltage to the piezoelectric transducer attached to the fiber (rapid) and temperature variation of the seed source (slow), continuous mode-hop-free tuning of the idler over more than 120 GHz was demonstrated (Fig. 7). Using a Fabry-Perot interferometer, the idler linewidth was measured to be 1.1 MHz at $3.17 \mu\text{m}$. The narrow linewidth, broad coarse wavelength coverage, and rapid mode-hop-free tuning of the idler through piezoelectric tuning of the pump enabled high-resolution spectroscopy in a variety of mid-IR gases including water vapor, CO_2 , and methane.

The development of PPLN has also led to substantial reductions in cw SRO power threshold, compatible with the direct use of semiconductor diode lasers as pumps for cw SROs. In addition to a compact design, an important advantage of this approach is the tunability of diode laser, which allows rapid and continuous tuning of SRO output at a fixed temperature and grating period through pump tuning. However, to provide the sufficiently high cw pump powers (typically >1 W) and the highest beam quality to attain SRO threshold, it has been necessary to boost the available power from single-mode diode lasers using amplification schemes. By employing a grating stabilized, extended-cavity single-stripe InGaAs semiconductor diode laser at 924 nm as a master oscillator and a single-pass tapered amplifier, Klein et al.²⁸ demonstrated operation of a cw SRO based on a 38-mm-long PPLN crystal with a pump power threshold of 1.9 W. For 2.25 W of diode pump power, 200 mW of single-frequency idler radiation was generated at $2.11 \mu\text{m}$. Wavelength tuning was achieved by electronic control of the master oscillator cavity, providing continuous mode-hop-free tuning of the diode pump radiation over 60 GHz from the power amplifier with a corresponding linewidth of <4 MHz. By using an intracavity etalon to fix the resonant signal frequency, a continuous mode-hop-free idler tuning of 56 GHz was obtained at $2.11 \mu\text{m}$ by tuning the pump wavelength. In an alternative scheme, using a distributed Bragg reflector (DBR) diode laser at 1082 nm, which was amplified in an Yb-doped fiber, Lindsay et al.²⁹ achieved rapid mode-hop-free tuning of a mid-IR cw SRO. A schematic of the experimental configuration is shown in Fig. 8, and the SRO idler

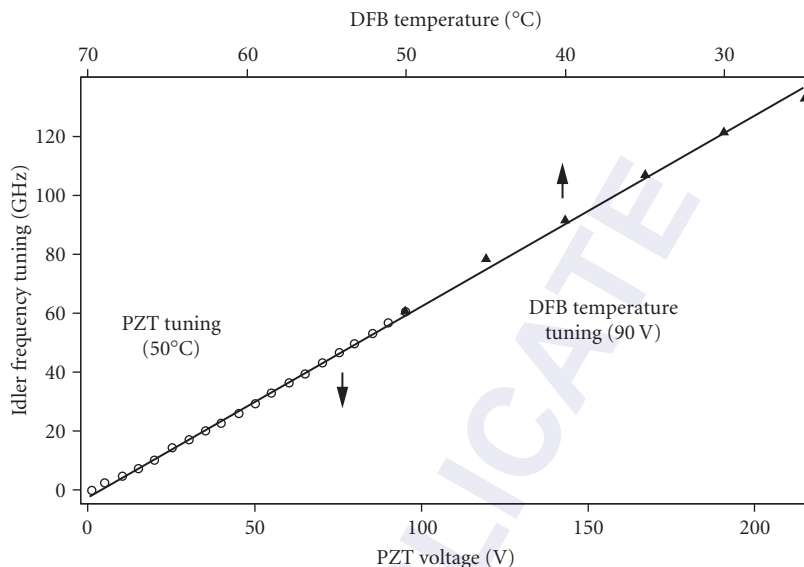


FIGURE 7 Fine tuning of the OPO idler frequency measured as a function of pump tuning parameter, performed by PZT voltage and fiber temperature variation.²⁷

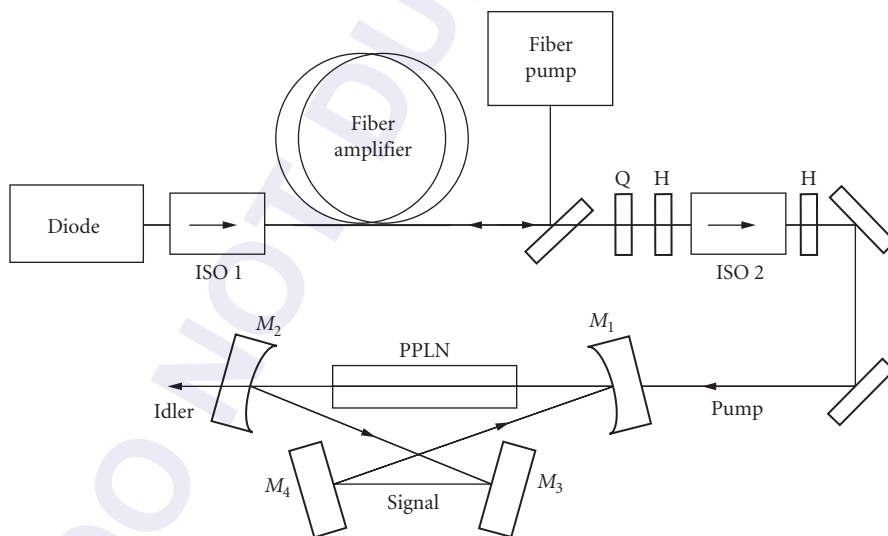


FIGURE 8 Schematic of experimental arrangement for the cw SRO pumped by a fiber-amplified DBR diode laser.²⁹

output power and tuning range are shown in Fig. 9. The SRO was based on a 40-mm PPLN crystal and could provide rapid continuous tuning over 110 GHz in 29 ms. Coarse and discontinuous wavelength tuning of the idler wave was also obtained over 20 nm by tuning the DBR diode laser, and more than 1 W of idler output power was generated across the 3.405 to 3425 μm range for 6.9 W of input pump power. An overall idler tuning range of 300 nm in the 3 to 3.5 μm band in the mid-IR was also available by varying the temperature of the PPLN crystal.

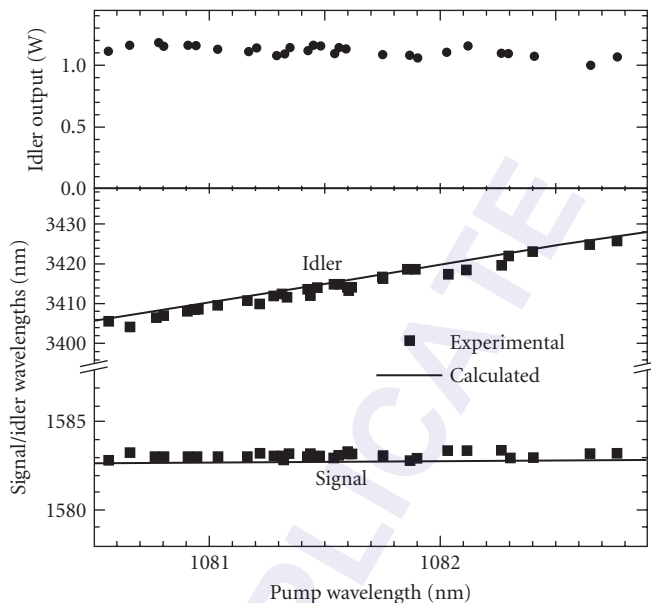


FIGURE 9 Variation of OPO output wavelengths (lower plot), and corresponding idler output power (upper plot), during pump tuning by seed laser DBR section alone. PPLN grating period was $29.75 \mu\text{m}$ and temperature was 180.5°C . Solid lines are calculated tuning range.²⁹

The availability of increasingly powerful pump sources such as cw fiber lasers, together with the high nonlinear coefficient ($d_{\text{eff}} \sim 17 \text{ pm/V}$) and long interaction lengths (50 to 80 mm) in PPLN, can now readily permit practical operation of cw SROs many times above operation threshold, providing multiwatt idler output powers. At the same time, the presence of high optical powers can lead to additional linear and nonlinear optical effects which can modify SRO output characteristics. These include thermal loading of the crystal due to linear absorption, which can result in thermal lensing, thermal phase mismatching and output beam quality degradation, or spectral generation and broadening due to higher-order nonlinear optical effects. As such, optimum performance of cw SROs at high pump powers requires strategies to combat such effects in order to achieve maximum conversion efficiency and output extraction at full pump power, while maintaining the highest spectral and spatial beam quality, power, and frequency stability. The performance characteristics of cw SROs at high pump powers many times threshold have been studied by Henderson and Stafford.³⁰ By deploying a 15-W cw single-frequency Yb fiber laser at 1064 nm as the pump and 50-mm MgO:PPLN crystals with multiple ($\Lambda = 31.5$ to $32.1 \mu\text{m}$) and fanned ($\Lambda = 30.8$ to $31.65 \mu\text{m}$) gratings, they investigated the effects of pump power on crystal heating, wavelength tuning, beam quality, and optimum output power and extraction efficiency. With the high beam quality of the fiber laser ($M^2 \sim 1.06$), and using a ring cavity with mirrors of highest reflectivity at the signal ($R \sim 99.9$ percent) and optimum mode-matching, they achieved a threshold as low as 1.0 W, enabling SRO operation at up to 15 times threshold. With the multi-grating crystal, the SRO reached a pump depletion of 91 percent at 2.5 times threshold, remaining constant to within ~ 10 percent up to the maximum pump power at 15 times above threshold. The idler output, measured at 2610 nm, exhibited a linear increase with input pump power, reaching 4.5 W at 15 W of pump, with a corresponding external photon conversion efficiency of ~ 74 percent. However, operation of SRO at increasing levels of pump power was found to result in a passive increase in crystal temperature and thus a shift in the output wavelength. At the highest pump power the rise in crystal temperature was as much as 23°C , leading to a significant shift in

signal (26 nm) and idler (57 nm) wavelengths compared to operation at low pump power. Given the minimal absorption of the MgO:PPLN crystal at idler wavelengths of 2 to 3 μm in this SRO, the self-heating effect was attributed to the finite absorption of the intracavity signal power. To confirm this, the authors deployed output coupling of the signal by replacing one of the high reflectors with a 4.2 percent output coupler. By operating the SRO at an ambient temperature of 26°C, they observed a 22°C rise in crystal temperature to 47°C under minimum output coupling at the maximum pump power. However, when using the 4.2 percent output coupler, the corresponding temperature rise was only 2.5°C, from 26°C to 28.5°C. Using measurements of signal output power, they estimated the circulating signal power to be as high as ~ 1.4 kW at the maximum pump power under minimum output coupling, decreasing to ~ 100 W with the 4.2 percent output coupler. By estimating the total absorption in the 50-mm crystal as 0.4 percent (0.08 percent/cm), they were able to conclude that an absorbed signal power of 5 W was responsible for the 22°C rise in crystal temperature. These measurements clearly confirmed the role of the intracavity signal power in heating of the MgO:PPLN crystal and its influence on spectral shifting of SRO output. The rise in crystal temperature was also observed to have a significant influence on the degradation of spatial quality of the idler beam by inducing thermal lensing effects within the crystal. From measurement of idler beam quality at the same output power level of 3.2 W, they were able to deduce a quality factor of $M^2 \sim 1.35$ under minimum output coupling compared to $M^2 \sim 1.0$ when using the 4.2 percent output coupler, hence confirming the deleterious effects of high circulating signal power on SRO beam quality and thus the need for optimization of output coupling at a given pump power to achieve the highest beam quality while maintaining maximum extraction efficiency. To this end, the authors also investigated the optimization of SRO output power and extraction efficiency at the maximum pump power by using variable output coupling (0 to 5 percent) for the signal across a limited tuning range. Using the fanned crystal, they found the optimum output coupling value to be 3.0 percent, resulting in the simultaneous extraction of 3.0 W of idler and 4.2 W of signal at an overall extraction efficiency of 48 percent. Under this condition, the pump depletion was 78 percent and SRO threshold was 5.8 W, corresponding to the optimum pumping ratio of ~ 2.5 for maximum power extraction. The effect of use of signal output coupling as a means of optimizing the performance of cw SROs was also later investigated in a separate experiment by Samanta and Ebrahim-Zadeh.³¹ Using a cw SRO based on MgO:sPPLT pumped at 532 nm, the authors demonstrated improvements of 1.08 W in total output power, 10 percent in total extraction efficiency, and a 130-nm extension in the useful tuning range, while maintaining pump depletions of 70 percent, idler output powers of 2.59 W, and a minimal increase in oscillation threshold of 24 percent. The output-coupled cw SRO could deliver a total power of up to 3.6 W at 40 percent extraction efficiency across 848 to 1427 nm. The single-frequency resonant signal also exhibited a higher spectral purity than the nonresonant idler output.

The high nonlinear gain coefficient of PPLN combined with the large optical powers present in cw SROs has also been observed to give rise to higher-order nonlinear effects in addition to the second-order parametric process. In a recent example of such an effect,³² operation of a cw SRO based on MgO:PPLN was reported together with simultaneous Raman action driven by the high intracavity signal intensity. The SRO, based on a multigrating MgO:PPLN crystal ($L = 50$ mm, $\Lambda = 28.5$ to 31.5 μm), was configured in a linear standing-wave cavity and pumped by a 10-W Yb fiber laser at 1070 nm. Two sets of cavity mirrors were used for the SRO, providing different reflectivities for the signal over 1500 to 1700 nm. With the low-Q cavity ($R = 98.2$ to 99 percent; $Q \sim 10^8$), normal cw SRO operation with the expected signal and idler spectra was achieved with a 3.3-W threshold, and 1.6 W of idler power was generated at 3620 nm for 8 W of pump at an optical efficiency of 20 percent and slope efficiency of 35 percent. With the high-Q SRO cavity ($R = 99.4$ to 99.8 percent; $Q \sim 10^9$), stimulated Raman action with characteristic spectra was simultaneously observed in the vicinity of signal spectrum, driven by the tenfold increase in intracavity signal power to ~ 100 W. The cw SRO threshold in this case was reduced to 0.5 W, with a corresponding reduction in optical efficiency to 16 percent and slope efficiency to 15 percent. The pump power threshold for Raman conversion was 1.9 W. While stimulated by intracavity signal power, Raman action was present only for grating periods and mirror reflectivities with lowest loss at the corresponding wavelengths, confirming the resonant nature of the observed effect. It was also observed that the presence of Raman oscillation with the high-Q SRO cavity resulted in improved idler RMS power stability of 1.46 percent compared to a 4.1 percent variation with the low-Q cavity, suggesting power limiting of intracavity signal by the Raman conversion.

In a subsequent experiment, Henderson and Stafford³³ also observed stimulated Raman oscillation in a high-power cw SRO based on MgO:PPLN. Using a 14.5-W cw single-frequency Yb fiber laser at 1064 nm and the same SRO arrangement as in Ref. 30, they observed Raman conversion of the intracavity signal under minimum output coupling and at pump powers more than 2 times above threshold, corresponding to circulating signal powers in excess of 230 W. Because of the increasing loss of SRO cavity across an extended tuning range, only two components of the Raman spectrum could be observed. However, under conditions of output coupling no Raman generation was observed up to the maximum available pump corresponding to 170 W of intracavity signal power. In the same cw SRO, the authors also observed spectral broadening of the resonant signal wave at high pump powers. Using highly reflecting mirrors to minimize threshold to 1.5 W, they were able to investigate the evolution of signal spectrum with pump power above threshold. It was observed that while at pump powers up to 3 times above threshold, the signal spectrum remained single-frequency, at pumping ratios between 3 to 4.7 the spectrum exhibited broadening with a symmetric pattern of side modes. The side modes were separated by between 0.2 and 0.5 nm, with their number and intensity increasing with pump power. Above a pumping ratio of 4.7, the signal spectrum was observed to become continuous with a FWHM bandwidth of ~2 nm. These observations, which were found to be in qualitative agreement with predicted theory, confirm that the operation of cw SROs at high pump powers and under the conditions of minimum signal coupling must be limited below a critical pumping ratio of ~4.5, if single-frequency oscillation is to be maintained. Since the maximum conversion in the same experiments was found to be attainable at a pumping ratio of 2.5, by choosing an optimum output coupling of 3.0 percent, the authors increased the SRO threshold to 5.1 W and so were able to maintain single-frequency operation up to the full available pump power of 14.5 W by remaining above the optimum pumping ratio (~2.5) for optimum conversion, but below the critical ratio (~4.5) for spectral broadening and multimode operation. Under this condition, 5.1 W of single-mode signal and 3.5 W of single-mode idler were simultaneously generated for 14.5 W of pump at an overall extraction efficiency of nearly 60 percent, with a measured idler bandwidth of 30 kHz over 500 μs .

The advent of QPM nonlinear materials has had a profound impact on cw SROs, with the vast majority of devices developed to date based on PPLN as the nonlinear material. When pumped near ~1 μm by solid-state, amplified semiconductor, or fiber lasers, they can provide potential coverage from above ~1.3 μm up to the absorption edge of the material near ~5 μm . For wavelength generation below ~1.3 μm , the use of PPLN is generally precluded by photorefractive damage induced by visible pump or signal radiation. As such, the development of practical cw OPOs for visible and near-IR at wavelength below ~1.3 μm has remained difficult, particularly in high-power SRO configuration where strong visible pump and signal radiation are present. This has thus necessitated the use of additional frequency conversion schemes or deployment of alternative QPM materials such as PPKTP and, more recently, MgO-doped periodically poled stoichiometric LiTaO₃ (MgO:sPPLT). To extend the tunable range of cw SROs to the visible range, Strossner et al.³⁴ used an approach based on second harmonic generation (SHG) of the idler output from a cw SRO in an external enhancement cavity. By deploying a 10-W, single-frequency, cw pump laser at 532 nm in combination with multigrating PPKTP ($L = 24$ mm, $\Lambda = 8.96$ to 12.194 μm) and PPLN ($L = 25$ mm, $\Lambda = 6.51$ to 9.59 μm) crystals as the OPO gain medium, and PPLN ($L = 43$ mm, $\Lambda = 6.51$ to 20.93 μm) for SHG, a visible green-to-red tuning range of 550 to 770 nm in the frequency-doubled idler was demonstrated. Together with direct signal (656 to 1035 nm) and idler (1096 to 2830 nm) tuning, this resulted in a total system tuning range of 550 to 2830 nm, with a tuning gap of ~60 nm over 1035 to 1096 nm. The output power, limited by photorefractive damage to the PPLN crystal, and optical damage to the PPKTP crystal and coatings induced by input pump, was 60 mW (signal), 800 mW (idler), and 70 mW (visible frequency-doubled idler) for up to 3.3 W of pump. The output signal from the free-running SRO exhibited a short-term linewidth of 20 kHz over 50 μs , with a jitter of 300 kHz over 5 ms, and 5 MHz over 1 s. By frequency locking the SRO to a monolithic Nd:YAG laser, a jitter-free linewidth of 20 kHz was measured at a signal wavelength of 946 nm. In the absence of pump tuning, mode-hop-free tuning of SRO output was obtained by adjustment of the cavity length using piezo control and synchronous rotation of the etalon using a feedback loop, resulting in 38 GHz of fine tuning in the signal for PPKTP and 5 to 16 GHz for PPLN, limited by photorefractive effects. The

free-running SRO exhibited a mode hop over a free-spectral range (680 MHz) every 10 min, but locking the SRO cavity to the 532-nm pump ensured a long-term frequency drift of <50 MHz/h in the signal and idler, accompanied by a reduction in output power by ~10 percent.

More recently, the development of MgO:sPPLT has brought about new opportunities for the advancement of practical cw SROs for the visible and near-IR at wavelength below ~1.3 μm with the direct use of high-power solid-state laser sources in the green. By deploying a 10-W, single-frequency, cw, frequency-doubled Nd:YVO₄ pump laser at 532 nm, MgO:sPPLT ($L = 30.14$ mm, $\Lambda = 7.97$ μm) as the nonlinear crystal and temperature tuning, Samanta et al.³⁵ demonstrated a cw SRO with a tunable range of 848 to 1430 nm. Using a linear standing-wave cavity and double-pass pumping, the cw SRO had an oscillation threshold of 2.88 W, and could provide >1.51 W of single-pass idler power for 6 W of pump at an extraction efficiency of >25 percent and photon conversion efficiency of >56 percent. The maximum idler power and conversion efficiency in this SRO was limited by thermal lensing effects, attributed to the finite liner absorption of the green pump light in the MgO:sPPLT crystal. Despite this, the SRO could deliver >500 mW of single-pass power across the entire idler tuning range of 1104 to 1430 nm, and in a Gaussian profile, confirming the absence of photorefractive damage as is present in PPLN. With a standing-wave SRO cavity and in the absence of intracavity frequency selection, the output frequency in both signal and idler was characterized by mode hops. Soon after, by deploying a compact ring cavity with a 500- μm intracavity etalon, the authors demonstrated single-frequency operation of the cw SRO.³⁶ Using the same pump laser and MgO:sPPLT crystal in a single-pass pumping arrangement, the SRO had a pump power threshold of 2.84 W and could deliver 1.59 W of single-mode idler power over 1140 to 1417 nm for 7.8 W of pump at >20 percent extraction efficiency. The total SRO tuning range was 852 to 1417 nm, obtained for a variation in crystal temperature from 61 to 236°C. Under free-running conditions, the idler had an instantaneous linewidth of ~7 MHz and exhibited a peak-to-peak power stability of 16 percent over 5 hours. Measurements of idler power at different crystal temperatures revealed stronger thermal lensing at higher temperatures. In a separate experiment, operation of a similar cw SRO based on MgO:sPPLT and pumped by a Nd:YVO₄ laser at 532 nm was reported by Melkonian et al.,³⁷ providing tunable signal over 619 to 640 nm in the red. Using a ring cavity for the SRO containing a 30-mm multigrating crystal ($\Lambda = 11.55$ to 12.95 μm) and a 2-mm-thick intracavity silica etalon, the SRO could provide ~100 mW of nonresonant idler power. The resonant signal was extracted using a 1.7 percent output coupler, providing 100 mW of single-frequency red radiation for 10 W of input pump power. The cw SRO threshold varied from 3.6 W in the absence of the intracavity etalon up to 6.6 W with signal output coupling, and rising to 6.8 W depending on the exact signal wavelength. The maximum pump depletion was 15 percent, limited by thermal effects attributed to pump and signal absorption. The output signal frequency could be mode-hop-tuned over a total range of 27 GHz by rotation of the intracavity etalon, in steps of 255 MHz corresponding to free-spectral-range of the SRO cavity. With active stabilization of SRO cavity length, a frequency stability of 20 MHz over 3 min was obtained for the signal.

The development of practical, high-power, single-frequency cw SROs based on MgO:sPPLT pumped in the green and operating below 1 μm ^{35–37} has also provided new motivation for spectral extension to shorter wavelengths. By using internal SHG of the resonant near-IR signal in a cw SRO based on MgO:sPPLT, Samanta and Ebrahim-Zadeh³⁸ demonstrated the first cw SRO tunable in the blue. A schematic of the SRO configuration is shown in Fig. 10. The device was based on similar experimental design as in Ref. 36, except for the exclusion of the intracavity etalon and inclusion of a 5-mm BiB₃O₆ (BIBO) crystal at the secondary waist of the bow-tie SRO ring resonator to frequency double the circulating signal radiation in a single direction. By varying the temperature of the MgO:sPPLT crystal to tune the signal over 850 to 978 nm, and simultaneous rotation of the BIBO phase-matching angle, a wavelength range of 425 to 489 nm in the blue was accessed. The generated blue power varied from 45 to 448 mW across the tuning range, with the variation arising from the non-optimum reflectivity of the blue coupling mirror over the signal wavelength range. The output power behavior and pump depletion of the SRO with pump power is shown in Fig. 11. The frequency-doubled SRO had a threshold of 4 W (2.4 W without the BIBO crystal), and exhibited a pump depletion of up to ~73 percent under blue generation. In addition to the blue, the device could provide in excess of 100 mW of signal and as much as 2.6 W of idler output power. Without an intracavity etalon, the

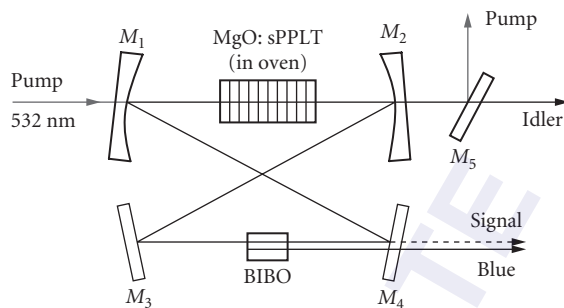


FIGURE 10 Schematic of the intracavity frequency-doubled MgO:sPPL cw SRO for blue generation.³⁸

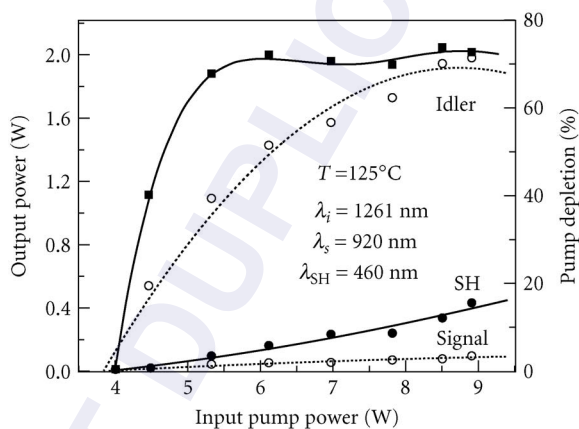


FIGURE 11 Single-frequency blue power, signal power, idler power, and pump depletion as functions of input pump power to the frequency-doubled cw SRO. Solid and dotted lines are guide for the eye.³⁸

single-mode nature of the pump and resonant signal resulted in single-frequency blue generation and a measured instantaneous linewidth of ~ 8.5 MHz in the absence of active stabilization. The blue output beam also exhibited a gaussian spatial profile. In the meantime, operation of an intracavity frequency-doubled cw SRO based on MgO:sPPL was also reported by My et al.,³⁹ providing tunable output in the orange-red. By resonating the idler wave in the 1170 to 1355 nm range in a ring resonator and employing a 10-mm intracavity β -BaB₂O₄ (BBO) crystal for doubling, tuning output over 585 to 678 nm was generated. With a 30-mm MgO:sPPL crystal ($\Lambda = 7.97\ \mu\text{m}$), up to 485 mW of visible radiation was internally generated for 7.6 W of pump, with 170 mW extracted as useful output. The device could also provide up to 3 W of nonresonant infrared signal power. The power threshold for the cw SRO was 4.5 W (4 W without the BBO crystal) and pump depletions of ~ 80 percent were measured for input powers > 6 W. Without active stabilization, the visible SHG output was single mode with a frequency stability of 12 MHz over 12 min, and mode-hop-free operation could be maintained over several minutes.

In a departure from conventional cw OPOs based on bulk materials, the use of guided-wave nonlinear structures can also in principle offer an attractive approach to the realization of OPO sources in miniature integrated formats. The tight confinement of optical waves in a waveguide can

provide substantial enhancement in nonlinear gain per input pump power compared with the bulk materials, but a major drawback of the approach is the unacceptably high input and output coupling losses in the waveguide, hindering OPO operation. The problem is further exacerbated in the cw regime, and in the SRO configuration, which is characterized by the highest oscillation threshold. In an effort to overcome this difficulty, Langrock et al.⁴⁰ deployed the technique of fiber-feedback to achieve operation of a cw SRO based on a reverse-proton-exchanged (RPE) PPLN waveguide. In this approach, the SRO cavity was formed in a ring using a single-mode optical fiber pigtailed to both ends of the waveguide, providing feedback at the resonant signal wave. The pump was similarly coupled into the PPLN waveguide using a separate fiber and, together with the nonresonant idler, exited the waveguide in a single pass. The configuration resulted in minimum coupling losses of 0.7 dB (signal input) and 0.6 dB (signal output), and alignment-free operation. Using a 67-mm RPE PPLN waveguide containing a 49-mm grating ($\Lambda = 16.1 \mu\text{m}$), and a tunable external cavity diode laser at 779 nm as the pump, cw SRO threshold was reached at ~ 200 mW of coupled pump power. The waveguide cw SRO exhibited gain bandwidths in excess of 60 nm. Ultimately, however, practical realization of such waveguide cw SROs offering significant output will require further optimization of waveguide fabrication process to minimize propagation losses ($\alpha = 0.2$ dB/cm) and loop losses (1.5 dB) in the present device, as well as further reductions in the waveguide-to-fiber coupling losses.

Multiple-Resonant Oscillators

Because the high pump power requirement for cw SROs can be prohibitive for many practical applications, extensive efforts have been directed to the development of cw OPOs in alternative resonance configurations, from the traditional DRO to the more recently devised PE-SRO and TRO, with the goal minimizing the pump power thresholds. These efforts have brought the operation of cw OPOs within the reach of more commonly available low- to moderate-power cw laser sources, albeit at the expense of added system complexity arising from the need for more elaborate cavity designs, more complex tuning protocols and the imperative requirement for active stabilization and control. In particular, the use of DRO and PE-SRO resonance schemes in combination with novel cavity designs have led to the practical generation of cw mid-IR radiation with the highest degree of frequency stability, and continuous mode-hop-free tuning capability over extended frequency spans at practical powers. These efforts have led to the realization of novel cw OPO systems in PE-SRO, DRO, and TRO configurations pumped by a variety of laser sources. These sources offer practical cw output powers in the mW to 100s mW range, high frequency stability, significant mode-hop-free tuning capability and extended wavelength coverage in the 1 to 5 μm spectral range.

Doubly Resonant Oscillators The use of DRO configurations in combination with PPLN has permitted substantial reductions in cw pump power threshold in cw OPOs, to levels compatible with the direct use single-mode semiconductor diode lasers, without the need for power amplification. The smooth wavelength tuning capability of the diode laser pump can then be similarly exploited to achieve continuous mode-hop-free tuning of the DRO output. In an example of such an approach, Henderson et al.⁴¹ demonstrated a PPLN cw DRO pumped directly with a 150-mW, single-mode, single-stripe, DBR diode laser at 852 nm. Configured in a single, linear, standing-wave resonator and using a 19-mm-long multigrating crystal ($\Lambda = 23.0$ to 23.45 μm), the DRO exhibited a pump power threshold of ~ 17 mW, with thresholds as low as 5 mW under optimum alignment and minimum output coupling. Using three DRO mirror sets, a signal (idler) wavelength range of 1.1 to 1.4 μm (2.2 to 3.7 μm) was accessed by temperature tuning the PPLN crystal. The DRO generated a total signal and idler power of 18 mW at 1.3 and 2.3 μm , respectively, for 89 mW of input diode pump power, with 4 mW of output in the idler beam. Continuous mode-hop-free tuning of the signal (idler) at 1.3 μm (2.3 μm) could be obtained over 12 GHz (7 GHz) by smooth tuning the frequency diode laser using temperature variation and over 17 GHz (10 GHz) using current control in combination with active servo control of the DRO cavity length to follow the pump frequency scan. The continuous mode-hop-free tuning ranges were limited by the restrictions on DRO cavity length variation imposed by the piezoelectric transducer. With a pump linewidth of <3 MHz, the signal wave was measured to have

a linewidth of <7 MHz and a free-running frequency stability of ~ 20 MHz over 10 s without active stabilization. To demonstrate the utility of the diode-pumped cw DRO, continuous tuning of the idler was used in single-pass absorption spectroscopy of R6 line in CO gas at $2.3 \mu\text{m}$.

In an effort to extend the wavelength coverage of cw OPOs to the visible spectrum, Petelski et al.⁴² demonstrated a DRO tunable in the yellow, using a cascaded frequency scheme comprising three resonators including two external enhancement cavities, and a cw single-frequency monolithic ring Nd:YAG laser as the primary pump source. The pump light at $1.064 \mu\text{m}$ was frequency doubled in an external enhancement cavity using $\text{MgO}:\text{LiNbO}_3$ as the nonlinear crystal, to provide the input radiation for the cw DRO based on an identical crystal. With a threshold of 15 mW, the DRO delivered 95 mW of idler power for 450 mW of input at 532 nm, with the idler tunable over 1.130 to $1.190 \mu\text{m}$ by changing the temperature of the $\text{MgO}:\text{LiNbO}_3$ crystal. The generated idler output was then frequency doubled in an external enhancement cavity based on a 16-mm PPLN crystal ($\Lambda = 8.0$ to $8.6 \mu\text{m}$) to provide tunable radiation across the 565 to 590 nm range. For 1.05 W of primary pump power at $1.064 \mu\text{m}$, an output power of 3.8 mW was obtained at 580 nm, with active stabilization of the DRO and enhancement cavities providing single-mode output with a 3 percent intensity noise and stable operation over 10 hours. Fine tuning of the yellow output was obtained by smooth tuning of the Nd:YAG pump source, as well as mode-hop tuning of the DRO by scanning the cavity length. By tuning the pump frequency over 10 GHz, the yellow output at 580 nm was tuned continuously over 18 GHz, while mode-hop tuning could provide 160 GHz of step tuning across 20 mode pairs.

Pump-Enhanced Singly Resonant Oscillators In separate experiments, using a single-mode, single-stripe, grating-stabilized AlGaAs diode lasers at ~ 810 nm in both solitary and external-cavity configurations, Lindsay et al.⁴³ demonstrated a PPLN cw OPO in pump-enhanced configuration. Using a single, linear, standing-wave cavity for the OPO and a 50-mm, multigrating ($\Lambda = 21.0$ to $22.4 \mu\text{m}$) crystal, a typical pump power threshold of 25 to 30 mW over a signal (idler) tuning range of 1.06 to $1.19 \mu\text{m}$ (2.58 to $3.44 \mu\text{m}$) was demonstrated, with ~ 4 mW of single-mode idler power available for 62 mW of pump. Wavelength tuning could be achieved by the variation of pump wavelength, PPLN crystal temperature, or grating period. In the external cavity configuration, locking of the OPO cavity length to the pump laser frequency using the Pound-Drever technique enabled stabilized single-mode operation with a signal (idler) intensity fluctuation of ± 3.5 percent (2.6 percent) at 0.5 mW (4 mW) power level over 1 hour. Over this period, pump frequency fluctuations of ± 125 MHz resulted in ± 100 MHz variation in signal frequency. However, mode-hop-free operation was maintained throughout the entire period. By continuous tuning of the external cavity diode laser over 510 MHz, a mode-hop-free tuning range of 377 MHz for the signal and 133 MHz for the idler was obtained, limited by the relative tuning rate of parametric gain curve and resonant signal modes with frequency tuning of the pump in the common-cavity PE-SRO, and so could be extended using a dual-cavity arrangement for the OPO.

By deploying a single-frequency cw Ti:sapphire pump laser, Turnbull et al.⁴⁴ demonstrated a cw PE-SRO based on a multigrating PPLN crystal ($L = 19$ mm, $\Lambda = 21.0$ to $22.4 \mu\text{m}$) with an extended mid-IR idler coverage to $5.26 \mu\text{m}$, well into the strong absorption region of the material beyond $\sim 4.5 \mu\text{m}$. The PE-SRO could provide idler tuning in two regions from 2.71 to $3.26 \mu\text{m}$ and from 4.07 to $5.26 \mu\text{m}$, using a combination of temperature tuning and grating period variation. The oscillator exhibited a typical cw pump power threshold of 100 mW, increasing to more than 500 mW near $5 \mu\text{m}$. For 750 mW of input pump power, the PE-SRO could typically provide a maximum one-way cw idler power of 16 mW. By employing a dual-cavity arrangement for the PE-SRO and a solid etalon in one arm of the resonant signal cavity, the authors were able to demonstrate mode-hop-free tuning of the idler over 10.8 GHz by smoothly tuning the pump laser over 12.3 GHz. In a further experiment, Stothard et al.⁴⁵ extended operation of the same PE-SRO to a single-cavity, traveling-wave ring geometry with the aim of extending the total fine-tuning performance of the oscillator and improving the idler output power. By using a low-finesse intracavity etalon to control the resonant signal frequency, the authors were able to achieve discontinuous mode-hop tuning of the idler frequency across the entire free spectral range of the etalon, corresponding to 83 GHz. The PE-SRO exhibited an oscillation threshold of 250 mW and could generate 35 mW of cw single-frequency mid-IR

idler radiation in the 2.8 to 3- μm spectral range for 600 mW of input Ti:sapphire pump power. In a separate experiment, Lindsay et al.⁴⁶ reported the operation of a cw PE-SRO based on PPRTA ($L = 20$ mm, $\Lambda = 39.6$ μm) as the nonlinear crystal. The PE-SRO was pumped by a diode-pumped cw single-frequency Nd:YVO₄ laser at 1.064 μm and was configured in a linear standing-wave resonator with the nonresonant idler double-passed through the cavity. With the 20-mm crystal length, the oscillator had an external pump power threshold of 250 mW and could deliver 87 mW of cw mid-IR idler output for 900 mW of input pump power. Coarse tuning of the idler over the wavelength range of 3.245 to 3.520 μm was achieved by temperature tuning the single-grating PPRTA crystal. Continuous mode-hop-free tuning of the PE-SRO over 0.7 GHz was also demonstrated by fine tuning the Nd:YVO₄ pump laser frequency.

In subsequent experiments, Muller et al.⁴⁷ investigated long-term frequency stability and linewidth properties of cw PE-SROs in common-cavity and dual-cavity configurations. Using multi-grating PPLN crystals ($L = 50$ mm, $\Lambda = 28.64$ to 30.16 μm), a 2.5 W cw Nd:YAG pump laser (linewidth ~ 1 kHz/100 ms, frequency drift ~ 1 MHz/min) at 1064 nm, and linear standing-wave cavity arrangements for both configurations, they studied frequency stability and linewidth of the idler output tunable in the 3.1 to 3.9 μm spectral range. Schematics of the experimental setups for both cavity configurations are shown in Figs. 12 and 13. In the common-cavity PE-SRO (Fig. 12) with the pump and signal resonant in the same cavity, the OPO was locked to the pump laser using the Pound-Drever-Hall (PDH) technique, and the oscillator exhibited a threshold of 280 mW. In dual-cavity arrangement (Fig. 13), the pump and signal were resonated in separate linear cavities, with an intracavity 500- μm -thick Nd:YAG etalon inserted into the signal cavity for improved frequency control. The pump cavity was locked to the pump laser using the PDH method, while the signal cavity was locked to the point of maximum idler power using an intensity lock without an external reference. The use of dual-cavity configuration increased the threshold to 380 mW, but resulted in stable, mode-hop-free operation over 30 min. While an intrinsic advantage of the common-cavity PE-SRO is direct stabilization of the signal frequency to pump, it is more sensitive to mechanical perturbations, leading to mode hops. Moreover, reliable mode-hop-free operation and continuous frequency tuning by tuning the pump laser are more difficult due to the simultaneous resonance of two different wavelengths within a single cavity. On the other hand, the dual-cavity approach can overcome spontaneous mode hops and continuous tuning limitations of common-cavity PE-SROs, and can also offer several tuning methods combining etalon, signal cavity and pump frequency tuning. The

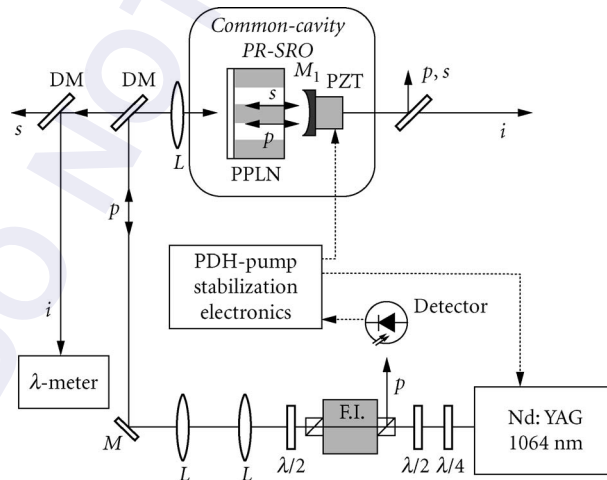


FIGURE 12 Common-cavity cw PE-SRO setup; M = mirror, DM = dichroic mirror, FI = Faraday isolator, L = lens, p = pump, s = signal, i = idler.⁴⁷

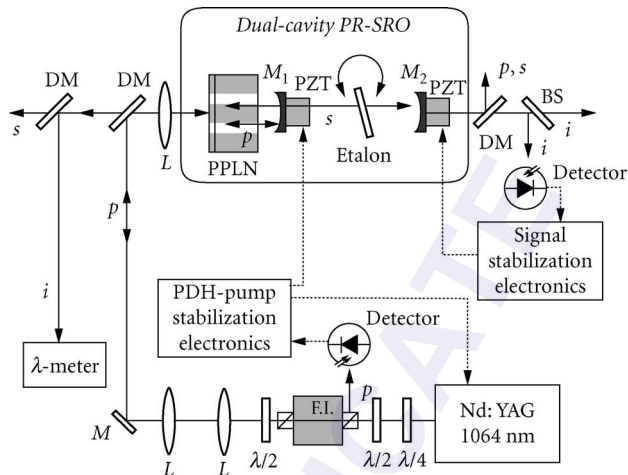


FIGURE 13 Dual-cavity cw PE-SRO setup; M = mirror, DM = dichroic mirror, FI = Faraday isolator, L = lens, p = pump, s = signal, i = idler.⁴⁷

results of the investigations revealed long-term frequency stability better than ± 30 MHz over more than 30 min for both configurations (Figs. 14 and 15), limited by the resolution of the wavemeter. The short-term frequency jitter was 56 kHz over 1.8 s for the common-cavity PE-SRO and 13.5 MHz over 1.5 s for the dual-cavity PE-SRO. The short-term linewidths, measured using the cavity leak-out technique in external high-finesse cavities, were (9 ± 2) kHz for the common-cavity and (6 ± 1) kHz for the dual-cavity over 20 μ s. The difference in frequency stability and linewidth of the two configurations is a result of the stabilization methods used. In the common-cavity PE-SRO, direct locking of the cavity to the pump laser provides a strong stabilization of the signal to pump frequency. In the

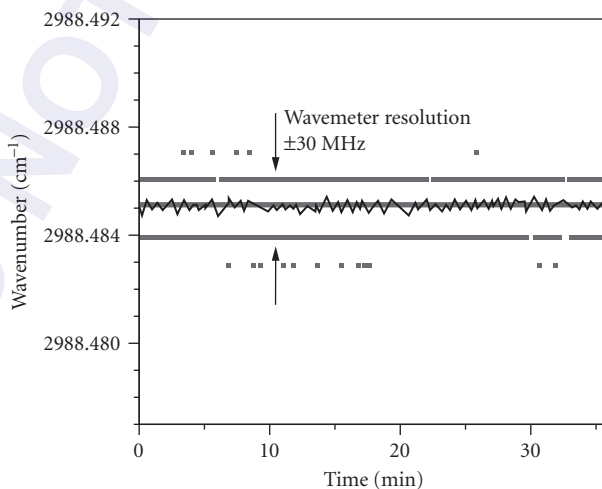


FIGURE 14 Long-term frequency stability (digital wavemeter read-out in 30 MHz steps and running average over 20 points) of the common-cavity cw PE-SRO.⁴⁷

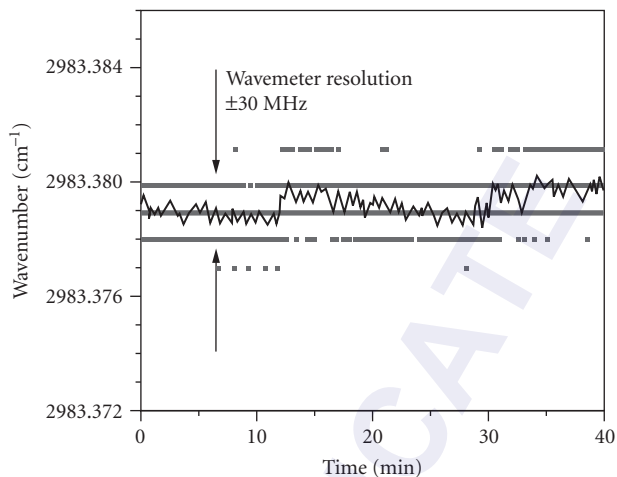


FIGURE 15 Long-term frequency stability (digital wavemeter read-out in 30 MHz steps and running average over 20 points) of the dual-cavity cw PE-SRO.⁴⁷

dual-cavity PE-SRO, the signal cavity is locked to the maximum idler output power and so the idler frequency is coupled to the maximum of the phase-matching gain curve. Any shifts in the gain curve in time will also result in corresponding variations in the signal and idler frequencies. However, the reported results confirm that despite this limitation, the dual-cavity PE-SRO provides the required frequency stability and linewidths as well as continuous mode-hop-free tuning necessary for high-resolution spectroscopy.

Triply Resonant Oscillators Characterized by the lowest oscillation threshold, the TRO represents the least demanding configuration for cw OPOs with regard to pump power. However, practical operation of such oscillators requires active cavity length control to maintain simultaneous resonance of the pump, signal and idler within the OPO resonator. In an example of such a device, Gross et al.⁴⁸ reported a TRO based on a 58-mm PPLN crystal using only 14 mW of pump power from a grating-stabilized, single-frequency, extended-cavity diode laser at 805 nm. By deploying a linear two-mirror cavity, highly reflecting mirrors, optimum mode-matching and active stabilization, they achieved a threshold pump power as low as 600 μ W at resonance, with a maximum total (signal and idler) output power of 47 μ W generated for 13.5 mW of pump. With mirrors of higher transmission (4 percent) at the signal and idler, a total output power of 2.1 mW was generated for 13.8 mW of pump, but at the expense of an increase in pump power threshold to 4.5 mW. Using a segmented design for the PPLN crystal consisting of two 19-mm outer sections poled with multiple grating of identical periods ($\Lambda = 20.2$ to 20.8 μ m) at 50 percent duty cycle and a single-domain section of length 20 mm at the centre, they demonstrated electro-optic tuning of the TRO output wavelengths at a fixed temperature, grating period, and pump wavelength. By applying an electric voltage of up to +1230 V across the single-domain section, wavelength tuning of the signal and idler over 1560 to 1660 nm could be obtained through modification of the phase-matching gain spectrum induced by the electro-optic effect. By applying a voltage modulation of amplitude 513 V at 0.11 Hz, signal (idler) tuning over 9.7 nm (10.8 nm) was demonstrated over 4.6 s, limited by the bandwidth of servo electronics.

Because of their phase coherent properties and the ability to generate exactly correlated frequencies, cw OPOs are also uniquely versatile sources for applications in optical frequency synthesis and metrology. Combined with compact solid-state semiconductor, or fiber pump lasers, they can offer practical tools for precision frequency generation, measurement, and control across extended regions spanning

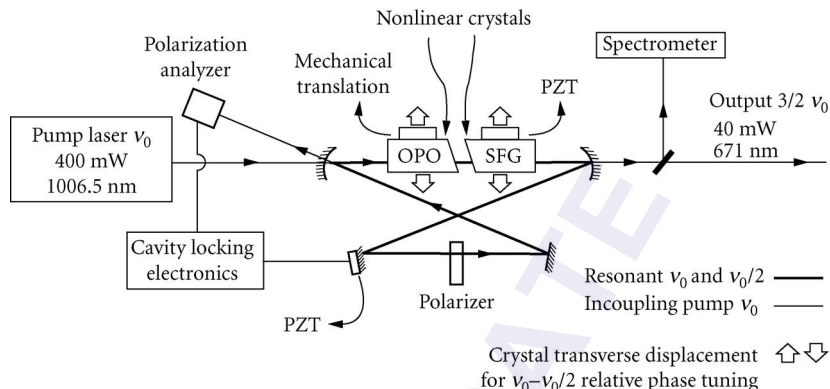


FIGURE 16 Schematic of the cw TRO for 3/2 frequency multiplier, converting a cw single-frequency radiation at 1006.5 nm into 671 nm. The wedged surfaces of the crystals are cut at an angle of 100 mrad with respect to the crystal axis. The input (output) facet of the OPO (SFG) crystal is at normal incidence. The two inclined surfaces facing each other are parallel. The transverse displacement of the nonlinear crystals provides an independent control over the cavity dispersion, ensuring simultaneous resonance of the two infrared fields.⁴⁹

the optical spectrum. In a report, Ferrari⁴⁹ demonstrated a particular architecture for optical frequency synthesis based on a cw OPO in TRO configuration. By taking advantage of the lowest oscillation threshold offered by the TRO, the author developed a 3/2 pump frequency multiplier from the near-IR to visible using a 400-mW, single-mode semiconductor MOPA device at 1006.5 nm as the pump source (Fig. 16). The TRO, based on PPKTP as the nonlinear crystal ($L = 20 \text{ mm}$, $\Lambda = 38 \mu\text{m}$), was operated at degeneracy to provide identical signal and idler frequencies at half the pump frequency. A second intracavity PPKTP crystal ($L = 20 \text{ mm}$, $\Lambda = 19.8 \mu\text{m}$) was then used to sum the pump with the degenerate frequency, resulting in 3/2 frequency multiplication of the pump, corresponding to a wavelength of 671 nm in the red. To provide stable operation, the TRO cavity length was locked using the pump resonance, and independent control of the degenerate signal and idler oscillating modes was obtained by using wedged crystals. Fine tuning of the TRO cavity modes could be obtained while maintaining pump resonance by lateral translation of the wedged crystals to alter in the optical path lengths within the crystals, thus enabling single-frequency operation to be achieved at degeneracy. The TRO had a threshold of $<50 \text{ mW}$ and could provide 40 mW of single-mode output at degeneracy in a near-gaussian spatial mode with an RMS amplitude noise of 1.4 percent (50 kHz bandwidth) over 3 min at full power.

17.3 APPLICATIONS

The important advances in cw OPOs over the last decade have led to the realization of a new generation of practical coherent light sources in new spectral regions from the visible to the near- and mid-IR offering unprecedented optical powers, high frequency stability and narrow linewidth, excellent beam quality, and extended fine and coarse tuning. These capabilities have paved the way for the deployment of cw OPOs in new application areas, in particular spectroscopy. Most notably, cw OPOs based on PPLN have found important applications in sensitive detection and analysis of trace gases in mid-IR, where a variety of important molecular fingerprints exist. A wide range of experiments, from simple single-pass absorption to high-resolution Doppler-free and cavity leak-out spectroscopy have been successfully performed by deploying cw OPOs based mainly on PPLN and operating in different resonance configurations of SRO, PE-SRO, or DRO. The higher cw output powers available to SROs and PE-SROs have also enabled detection of trace gases in the mid-IR with unprecedented sensitivity using photoacoustic spectroscopy.

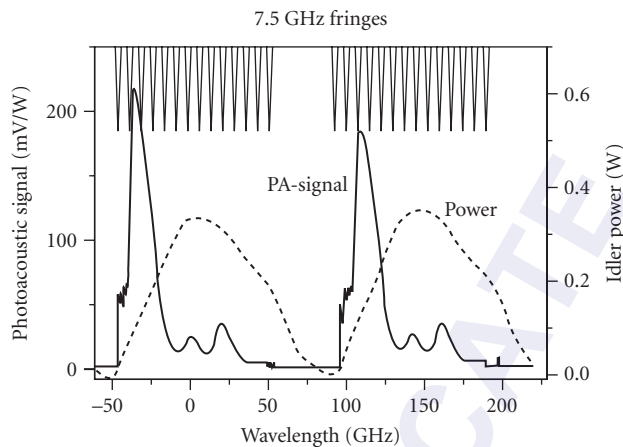


FIGURE 17 100-GHz-wide mode-hop scan of 10 ppm ethane in nitrogen was made around 2996.9 cm^{-1} . The solid line at the bottom shows the photoacoustic signal, the dashed line the idler power and the solid line at the top the fringes from a 7.5-GHz external Fabry-Perot etalon.²¹

By exploiting fine frequency tuning in a diode-pumped PPLN cw SRO, Klein et al.²⁸ performed single-pass absorption spectroscopy of rovibrational transitions in N_2O gas near $2.1\text{ }\mu\text{m}$. The wide and continuous mode-hop-free tuning of the idler over 56 GHz enabled monitoring of three molecular lines separated by $\sim 20\text{ GHz}$ within a single frequency scan.

Using a cw SRO based on PPLN, providing 700 mW of idler power and a total single-frequency fine tuning range of 24 GHz, Van Herpen et al.²⁰ recorded absorption line of ethane in nitrogen using the photoacoustic spectroscopy technique. In a later experiment,²¹ by deploying a more powerful pump laser and a similar SRO configuration, photoacoustic spectroscopy of ethane in nitrogen was demonstrated with a detection sensitivity of 10 parts per trillion (ppt) (Figs. 17 and 18). The

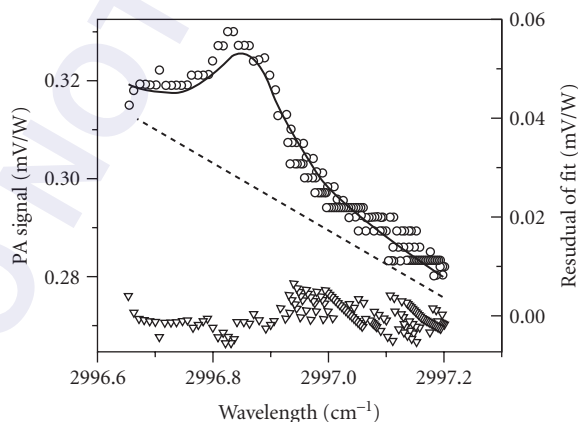


FIGURE 18 Pump-scan around 2996.9 cm^{-1} of 0.4-ppb ethane in nitrogen at atmospheric pressure. A Lorentzian fit with linearly decreasing background has been plotted through the data (dashed line). The linearly decreasing background is also shown separately with a dotted line. The bottom of the picture shows the residual if the fit is subtracted from the data.²¹

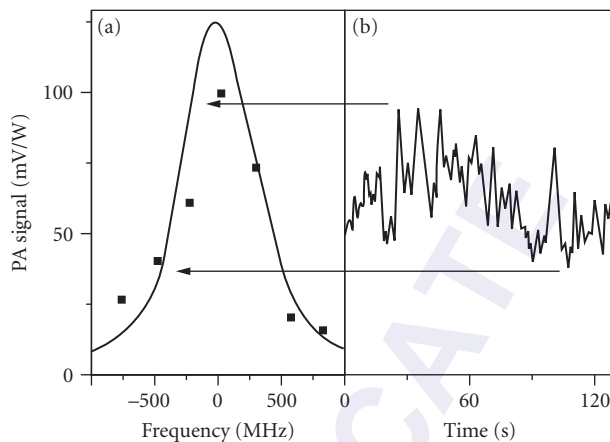


FIGURE 19 The high resolution performance and the wavelength stability of the SRO are demonstrated by recording the photoacoustic signal from the half maximum of a 77 mbar pressure broadened absorption line of 20 ppm of ethane in nitrogen at 2996.9 cm^{-1} . When not tuning the pump frequency, the photoacoustic signal shows random oscillations at a rather high frequency (90 MHz/s), combined with an oscillation at a lower frequency (250 MHz in 200 s).²³

extension of operation of the SRO to the 3.7 to 4.7- μm spectral range enabled photoacoustic detection of the CO_2 absorption line in nitrogen near 4.235 μm using 24-GHz continuous scan of the idler frequency.²² By using a high-power cw SRO based on PPLN operating in the 2.6 to 4.7 μm in the mid-IR, photoacoustic spectroscopy of a mixture of 20-ppm ethane in nitrogen was performed near 3.33 μm .²³ The SRO was pumped by a 20-W, cw single-frequency Yb:YAG laser (tunable over 1.024 to 1.034 μm) and could deliver 3 W of idler power at 2.954 μm (Fig. 3). By step-tuning the pump frequency using a combination of intracavity etalon and a Lyot filter, mode-hop-tuning of the idler over 190 GHz could be achieved, enabling coverage of several absorption lines of ethane. With 2.15 W of idler power available at 3.33 μm , a photoacoustic detection sensitivity of 0.005 ppb was deduced. By monitoring the photoacoustic signal corresponding to the strongest absorption peak in the same 20-ppm mixture of ethane in nitrogen at low pressure (77 mbar) as a function of time (Fig. 19), a slow idler frequency drift of 250 MHz over 200 s was measured, which was attributed to the fluctuations in the PPLN crystal temperature and lack of thermal isolation of the SRO from the environment. The idler frequency also exhibited fast frequency fluctuations of 90 MHz/s, attributed to the nonoptimized coatings of the PPLN crystal and cavity mirrors, which resulted in unwanted etalon and resonance effects in the SRO cavity. In a later report, Ngai et al.²⁴ performed photoacoustic and cavity leak-out spectroscopy of several trace gases including CO_2 and multicomponent gas mixtures of methane, ethane, and water in human breath using an automatically tunable cw SRO based on a multigrating MgO:PPLN crystal (Figs. 4 and 5). By deploying a ring SRO cavity with an uncoated intracavity YAG etalon and using a combination of pump tuning and etalon rotation, step-tuning of the idler over 207 GHz could be achieved. Adjustment of the crystal temperature could further be used to repeat this process, to provide wavelength scans of up to 450 cm^{-1} with a single grating period at high resolution ($<5 \times 10^{-4}\text{ cm}^{-1}$). By translating the crystal to other grating periods, an extended idler wavelength range of 2.75 to 3.83 μm could be accessed. Using the wide wavelength scanning capability, extended wavelength coverage, and automatic tuning capability of this source, photoacoustic spectroscopy of strong CO_2 combination bands in laboratory air (460 ppmv) extending over 14 nm (from 2788 to 2802 nm) was recorded by a combination of pump tuning, etalon rotation and crystal temperature tuning with a spectral resolution of 0.01 nm and a recording time of 1 hour. The recorded spectra were corrected for the water-vapor contribution to reveal the true

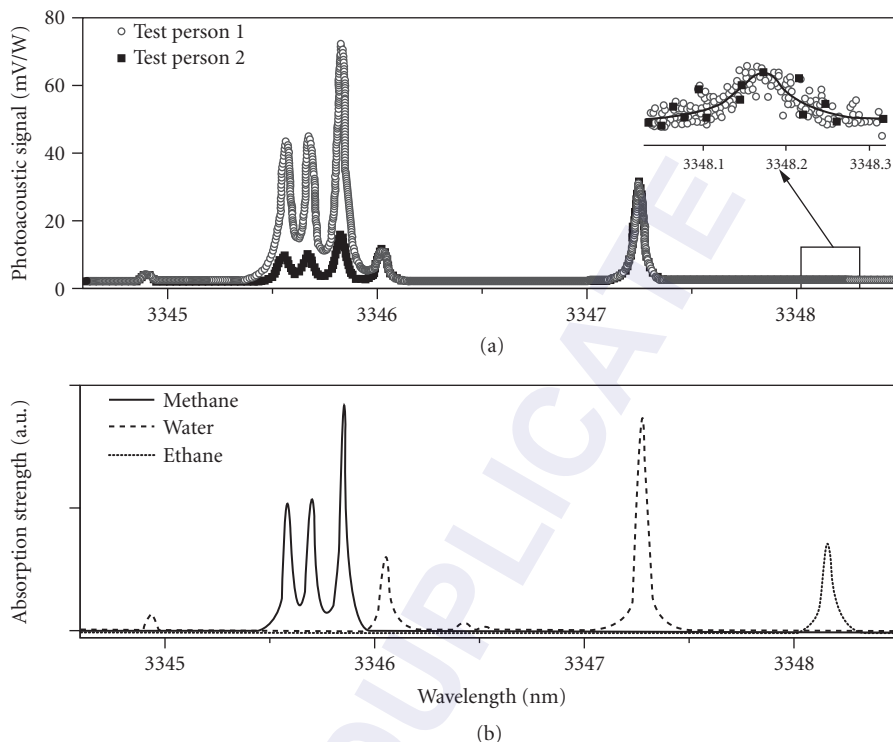


FIGURE 20 (a) Photoacoustic spectra measured from the breath of two different test persons. The recorded spectra reveal a higher methane concentration in the breath of person 1, indicating the presence of methanogenic flora. (b) Calculated absorption spectra based on the HITRAN database for methane, ethane, and water.²⁴

CO₂ spectra, with the results in excellent agreement with calculations. Similarly, using the same automatic tuning protocol, photoacoustic spectroscopy of methane, ethane, and water vapor in human breath were simultaneously recorded by scanning the idler wavelength in the 3344 to 3349 spectral range (Fig. 20). Concentrations of 21 ppmv and 13 ppbv were calculated for methane and ethane, respectively, in the first sample, and 4.5 ppmv and 13 ppbv in the second sample. Using the same cw SRO, cavity leak-out spectroscopy was performed for gas mixtures containing methane (at 3.221 μm) and ethane (at 3.337 μm) in N₂ at 100 mbar pressure. In the case of methane, the absorption peak was scanned over 200 s with a spectral resolution of 0.001 cm^{-1} , resulting in a background methane concentration of 31 ppbv in N₂, a noise-equivalent detection limit of 0.16 ppbv, and a minimum detectable absorption coefficient of $2.0 \times 10^{-9} \text{ cm}^{-1}$. For ethane, the idler was scanned over 60 s with a lower resolution of 0.01 cm^{-1} , resulting in concentrations from 5 to 100 ppbv, a noise-equivalent detection limit of 0.07 ppbv, and a minimum detectable absorption coefficient of $1.4 \times 10^{-9} \text{ cm}^{-1}$. The cw SRO was also used to record the spectrum of ¹²CH₄ and ¹³CH₄ isotopes of methane in laboratory air with cavity leak-out spectroscopy. Scanning the idler wavelength from 3210 to 3211.5 nm revealed the absorption features corresponding to the two isotopes as well as water vapor, with the measured data in good agreement with calculated spectra.

The rapid and continuous mode-hop-free idler tuning over 110 GHz in a PPLN cw SRO pumped by a fiber-amplified DBR diode laser²⁹ has been used to perform real-time single-pass absorption spectroscopy of methane near 3.39 μm with a refresh period of 29 ms and a corresponding refresh rate of 34 Hz, demonstrating the suitability of the system for rapid spectroscopic measurements (Fig. 21). In a subsequent experiment,⁵⁰ by taking advantage of frequency modulation capabilities

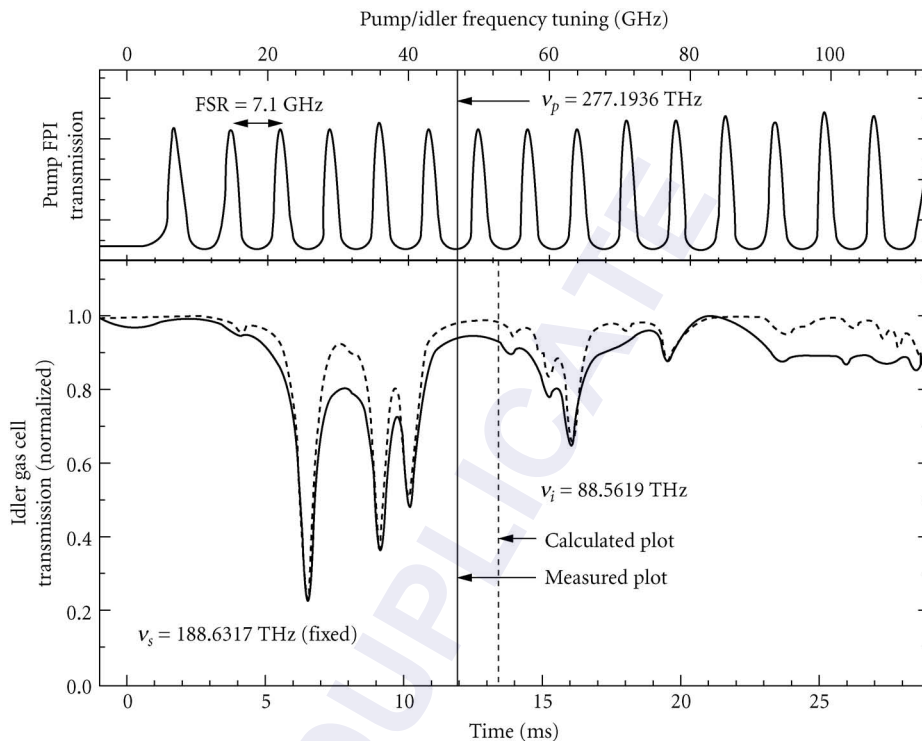


FIGURE 21 Rapid mode-hop-free tuning of a cw SRO pumped by a fiber-amplified DBR diode laser, and application absorption spectroscopy of CH_4 .²⁹

of the idler output from a similar fiber-amplified DBR diode-pumped cw SRO, sensitive detection of multicomponent trace gas mixtures was performed using quartz-enhanced photoacoustic spectroscopy (QEPAS). The SRO, based on a multigrating MgO:PPLN crystal ($L = 50 \text{ mm}$, $\Lambda = 28.5$ to $31.5 \mu\text{m}$), was similarly configured in a ring cavity with a $400\text{-}\mu\text{m}$ -thick uncoated intracavity YAG etalon to provide enhanced frequency selection and stability (Fig. 22). The pump laser delivered 7.9 W of single-frequency power at 1082 nm with a linewidth of $\sim 100 \text{ MHz}$. The idler wavelength was tuned coarsely over 3 to $4 \mu\text{m}$ by a combination of grating and temperature tuning, in discrete steps of 1 to 4 cm^{-1} due to the combination of SRO cavity modes, etalon mode selection and phase-matching bandwidth of the crystal. Pump tuning was then used to provide fine wavelength control and mode-hop-free tuning of the idler over 5.2 cm^{-1} , which could be shifted within a total range of 16.5 cm^{-1} by current control of the DBR diode laser. By monitoring the idler wavelength with a wavemeter, and using a computer-controlled feedback loop to adjust the phase and current to the DBR, locking of idler frequency at 2990.076 cm^{-1} could be achieved with a stability of $1.7 \times 10^{-3} \text{ cm}^{-1}$ over 30 min . With the programmed tuning, idler wavenumbers from 2987 to 2994 cm^{-1} could be accessed in 1-cm^{-1} steps within the pump tuning range, with any desired wavelength reached within $\sim 20 \text{ s}$ of tuning the pump. Using this source, QEPAS spectral data of 2-ppmv ethane in nitrogen were successfully recorded at 2990.1 cm^{-1} , as well as isolated 1.2 percent water vapor at 2994.4 cm^{-1} (Fig. 23). By locking the pump laser frequency to the maximum ethane absorption peak at 2990.08 cm^{-1} and measuring concentrations from 20 ppmv down to 100 ppbv , a QEPAS detection limit of 25 ppbv was deduced. By keeping the idler wavelength fixed at the same absorption peak with a 2 ppmv methane concentration, a linear dependence of QEPAS signal on input power was verified, confirming the advantage of higher powers for improved detection sensitivity. To further demonstrate the

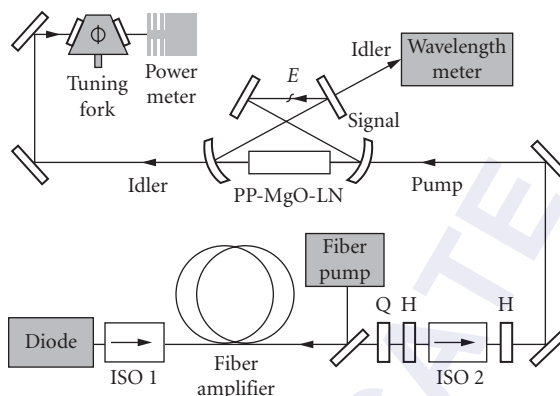


FIGURE 22 Schematic of the cw SRO experimental setup for quartz-enhanced photoacoustic spectroscopy.⁵⁰

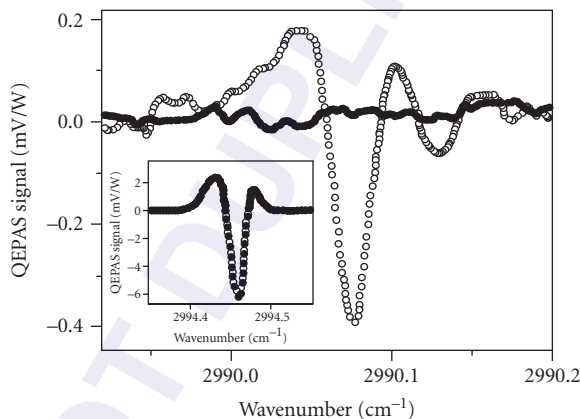


FIGURE 23 Example of a QEPAS scan over a 2 ppmv ethane peak at 2990.08 cm^{-1} (open circles) in 200 mbar nitrogen gas. QEPAS signal detected from pure nitrogen in the same spectral range is also shown, revealing a flat background signal (filled circles). Another example of a QEPAS scan is given in the inset for 1.2 percent water at 2994.4 cm^{-1} .⁵⁰

unique advantages offered by the wide tuning coverage of the cw SRO, QEPAS data corresponding to a multiple gas mixture containing 2.2 ppmv ethane, 1.1 percent water, and 1.5 ppmv methane were successfully recorded by pump-tuning the idler from 2979.4 to 2994.6 cm^{-1} , with the obtained spectra in good agreement with calculations. From the obtained data, it was also possible to deduce an improved QEPAS detection limit of 13 ppbv for the stronger methane absorption line at 2983.3 cm^{-1} and 6.8 ppmv for water at 2994.4 cm^{-1} .

By deploying an all-fiber-pumped PPLN cw SRO,²⁷ single-pass absorption spectroscopy of several gases over the wavelength range of 2700 to 3200 nm was performed by continuous scanning of the idler frequency with a linewidth of $\sim 1\text{ MHz}$. Using piezoelectric tuning of the fiber pump laser allowed mode-hop-free tuning of the idler frequency up to 60 GHz, enabling high-resolution spectroscopy of various absorption features in water vapor (2709 nm), carbon dioxide (2810 nm), nitrous oxide (2879 nm), ammonia (2897 nm), and methane (3167 nm).

In a separate experiment, using a PPLN cw DRO pumped directly by a 150-mW DBR diode laser and delivering 4 mW of idler output with a mode-hop-free tuning range of 10 GHz, simple single-pass absorption spectroscopy of CO molecule at 2.3 μm was demonstrated.⁴¹ Spectroscopic applications of cw OPOs in the visible spectral range has also been performed by external frequency doubling of the idler output from a cw DRO in an external enhancement cavity to generate yellow radiation in the 565 to 590 nm region.⁴² The DRO was pumped by a cw, monolithic ring Nd:YAG laser at 1064 nm and by tuning the pump frequency over 10 GHz, the yellow output at 580 nm could be tuned smoothly over 18 GHz, while mode-hop tuning could provide 160 GHz of step-tuning across 20 mode pairs. This tuning capability enabled spectroscopy of the ${}^5D_0 \rightarrow {}^7F_0$ transition in $\text{Eu}^{3+}:\text{Y}_2\text{SiO}_5$ at 4 K, which is of interest because it exhibits the lowest known homogeneous linewidth for an optical transition in a solid. The step-tuning capability enabled the full spectrum of the transition containing two absorption peaks at 580.070 and 580.224 nm to be scanned, while the fine mode-hop-free tuning enabled continuous scan of each inhomogeneously broadened absorption peak, resulting in linewidths of $\Delta\nu = 3$ GHz and $\Delta\nu = 2.5$ GHz, respectively. Persistent spectral hole-burning was also observed by continuous tuning of the pump, and linewidths <1 MHz were recorded after 40 min of burning. By monitoring the hole spectrum every few hours, it was possible to follow the hole decay as long as 15 hours, demonstrating the repeatability of frequency tuning and stability of the cw DRO over many hours.

By deploying a cw PE-SRO based on a 19-mm-long PPLN crystal, Kovalchuk et al.⁵¹ performed high-resolution Doppler-free spectroscopy of rovibrational transition of methane molecule at 3.39 μm . The PE-SRO, configured in semi-monolithic cavity design and containing an intracavity etalon, was pumped by a cw single-mode Nd:YAG ring laser with a linewidth of ~ 5 kHz. It could provide single-frequency idler output with a linewidth of ~ 100 kHz and a mode-hop-free tuning range of 1 GHz, obtained by smooth tuning of the pump laser. The PE-SRO had a minimum oscillation threshold of 305 mW and generated a total two-way idler power of 58 mW at 3.39 μm for 808 mW of input pump power. With this set-up, Doppler-free resonances with a resolution of ~ 500 kHz were recorded in methane, limited by the frequency jitter of the PE-SRO idler output and pressure broadening. Subsequently, Muller et al.⁵² reported a transportable, all-solid-state photoacoustic spectrometer based on a cw PE-SRO using an alternative dual-cavity semi-monolithic cavity design. The PE-SRO used a 19-mm-long multigrating PPLN crystal and was pumped by a 2.5-W cw single-frequency Nd:YAG laser at 1.064 μm . The oscillator was characterized by a pump power threshold of 380 mW and could generate a total two-way idler power of 200 mW. Using active signal cavity length stabilization, an idler frequency stability better than ± 30 MHz was obtained. Coarse tuning of the idler wavelength over 3.1 to 3.9 μm was performed by a combination of temperature and grating period variation, while the use of an etalon inside the signal cavity enabled continuous mode-hop-free tuning of the idler over 1.5 GHz by tuning the pump laser. Mode-hop tuning of the idler frequency in steps of 450 MHz could also be achieved over 52 GHz by rotation of the intracavity etalon. Using this mode-hop tuning method and with 70 mW of available idler power, photoacoustic spectroscopy of ethane was performed with a detection sensitivity of 110 ppt. By deploying the same dual-cavity cw PE-SRO in combination with cavity leak-out spectroscopy, substantial improvement in detection sensitivity of ethane down to 0.5 ppt was demonstrated near 3 μm ,⁵³ and by exploiting the frequency tuning capability of the PE-SRO, simultaneous monitoring of ethane, methane, and water vapor in human breath could be performed, without significant interference from other gases.

In a further application, Stothard et al.⁵⁴ demonstrated the use of a cw PE-SRO for hyperspectral imaging of gases in the mid-IR. Based on a 20-mm-long, single-grating crystal of PPRTA and pumped by a 1-W diode-pumped Nd:YVO₄ laser, the PE-SRO could provide ~ 50 mW of idler output power and coverage in the 3.18 to 3.50- μm spectral range by temperature tuning the crystal. The use of a three-mirror standing-wave PE-SRO cavity and an intracavity etalon at the signal wavelength-enabled mode-hop tuning of the idler frequency over ~ 30 GHz, in steps of 1.4 to 2 GHz, by rotation of the etalon. This frequency resolution was sufficient compared to the typical linewidth of ~ 5 GHz for pressure-broadened transitions in gases under atmospheric pressure, thus enabling the deployment of the cw PE-SRO for imaging of methane gas in the atmosphere. By tuning the cw PE-SRO idler frequency to the strong methane absorption lines near 3.27 and 3.35 μm , gas concentrations of

the order of 30 ppm-m could be detected and significant target areas ($\sim 4 \text{ m}^2$ at 3 m) could be effectively imaged with the subsecond acquisition times.

As well as versatile spectroscopic tools, cw OPOs offer unique sources of correlated twin beams and nonclassical states of light for applications in quantum optics and quantum information processing. When pumped by cw all-solid-state lasers, they can provide compact twin-beam optical sources for quantum cryptography and sub-shot-noise measurement. In one configuration of such a device, a cw TRO based on a 10-mm KTP crystal and using a linear semi-monolithic resonator to separate the pump and parametric wave cavities, was reported by Hayasaka et al.⁵⁵ The TRO was pumped at 540 nm by the second harmonic of an extended-cavity single-stripe cw diode laser and was operated close to degeneracy. The use of this pump wavelength permitted type II NCPM in KTP, resulting in a TRO pump power threshold as low as 2.5 mW near degeneracy and providing 5.1 mW of cw output power for 16 mW of pump power. By recording the noise spectrum of the twin-beam intensities, 4.3 dB of intensity-difference squeezing was observed at ~ 3 MHz. In a subsequent experiment, Su et al.⁵⁶ demonstrated quantum entanglement between the signal and idler twin beams in a nondegenerate cw DRO above threshold. The DRO, based on a 10-mm PPKTP crystal cut for type II nondegenerate phase matching, was configured in a linear semi-monolithic cavity and was pumped by the stabilized cw single-mode output of a frequency-doubled Nd:YAP laser at 540 nm. Using a pair of unbalanced Mach-Zender interferometers with unequal arm lengths, the amplitude and phase noise of the signal and idler beams above threshold were recorded at 20 MHz, enabling quantum correlations to be deduced from the noise levels of the intensity difference and phase sum of the photocurrents measured by the unbalanced interferometers. Using this method, the authors were able to deduce correlations of amplitude and phase quadratures of signal and idler below the shot-noise-limit, amounting to ~ 2.58 and ~ 1.05 dB, respectively, and from the sum of the amplitude and phase correlation variances, demonstrate quantum entanglement of the twin beams below the shot-noise-limit. For 230 mW of input pump power, nearly twice the 120-mW DRO threshold, the output power in the correlated twin beams was 22 mW. The nondegenerate signal and idler twin beams were at wavelengths of 1079.130 and 1080.215 nm, respectively, separated by 1.085 nm. In another experiment, Tanimura et al.⁵⁷ deployed a cw DRO below threshold to generate squeezed vacuum on resonance with the rubidium *D* line at 795 nm. The DRO, based on a 10-mm crystal of PPKTP, was configured in a ring cavity and pumped at 397.5 nm by the second harmonic of a cw single-frequency Ti:sapphire laser to provide near-degenerate signal and idler frequencies at 795 nm (Fig. 24). Operating the cw DRO below threshold and using homodyne detection, the authors were able to measure strongly squeezed vacuum at the OPO output. With the DRO operated at 61 mW of input pump power, below a calculated threshold of 150 mW, a squeezing level of -2.75 dB below the shot-noise-limit and anti-squeezing level of $+7.00$ dB above the shot-noise-limit was observed (Fig. 25). Such a system could find useful applications for ultraprecise measurements of atomic spins as well as quantum information processing by mapping squeezed vacuum onto an atomic ensemble.

Because of their phase coherent properties and the ability to lock correlated frequencies to stable optical references, cw OPOs also represent highly promising light sources for applications in optical frequency synthesis and metrology. In an example of such application, a novel approach based on the combination of a cw OPO with a femtosecond Ti:sapphire frequency comb was used to provide a phase-coherent bridge from the visible to mid-IR spectral regions.⁵⁸ The cw PE-SRO, based on a 19-mm PPLN crystal and deploying a similar configuration to that in Ref. 51 was pumped by a cw single-frequency Nd:YAG laser. The oscillator provided idler emission in the 2.4 to 3.7- μm spectral range, with 50 mW of idler power and an instantaneous linewidth of ~ 10 kHz at 3.39 μm . The technique takes advantage of the fact that in a PPLN cw PE-SRO, in addition to the phase-matched signal and idler, there are also non-phase-matched frequencies generated by mixing of the resonant pump and signal waves. This process provides a range of visible frequencies within the emission bandwidth of a femtosecond Ti:sapphire laser, which can be used for frequency comparison with the nearest comb lines. By forming suitable differences of the heterodyne beat frequencies between the visible frequency components from the cw OPO and adjacent comb lines in the femtosecond laser, mutual phase locking of OPO optical frequencies, Ti:sapphire repetition frequency and carrier-envelope offset frequency could be obtained. Using this method, the authors performed direct frequency comparison between an iodine-stabilized Nd:YAG laser at 1.064 μm and a mid-IR methane

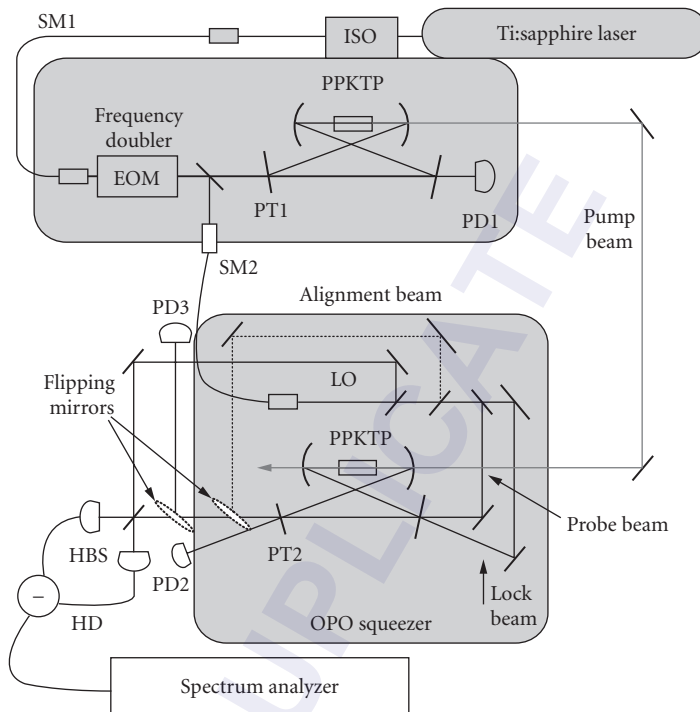


FIGURE 24 Experimental setup for PPKTP cw DRO for the generation of squeezed vacuum. ISO, optical isolator; EOM, electro-optic modulator; OPO, subthreshold degenerate optical parametric oscillator; HBS, half-beam splitter; PT, partial transmittance mirror; HD, balanced homodyne detector; PD, photodiode; SM, single-mode fiber.⁵⁷

optical frequency standard at $3.39 \mu\text{m}$. Subsequent demonstrations, taking advantage of the unique coherence properties of cw OPOs for frequency synthesis, include the development of a precise $3/2$ frequency multiplier from a near-IR pump to the visible using a cw TRO.⁴⁹

17.4 SUMMARY

This chapter has provided an overview of the latest advances in cw OPOs and their applications over the past decade. The advent of QPM nonlinear materials has had an unprecedented impact on cw OPO technology which, combined with major advances in solid-state laser technology and innovative design architectures, has led to the realization of a new generation of truly practical cw OPOs with performance capabilities surpassing conventional lasers. The mature technology and ready availability of PPLN has enabled the development of near- and mid-IR cw OPOs for the 1 to $5 \mu\text{m}$ spectral range in various resonance configurations, from TROs offering minimal oscillation threshold and lowest output power to cw SROs with highest threshold and watt-level output power, previously unattainable with birefringent nonlinear materials, and using a variety of pump sources from miniature semiconductor diode lasers to high-power solid-state and fiber lasers. The application of novel cavity designs and resonance schemes including pump enhancement, dual cavities and

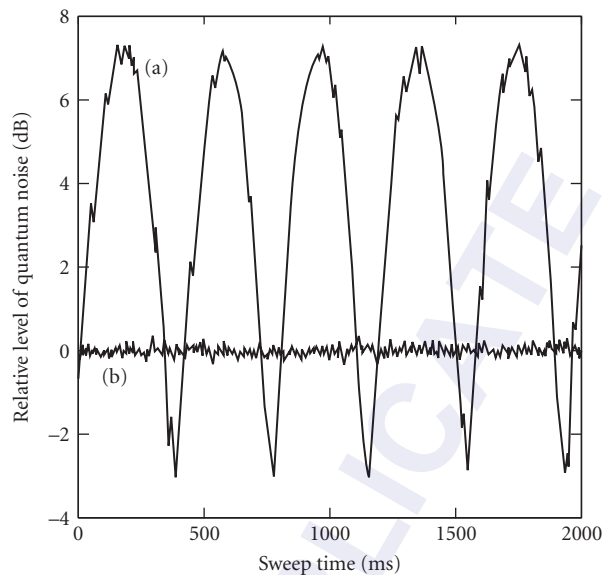


FIGURE 25 Measured quantum noise levels in cw DRO base on PPKTP. (a) Local oscillator beam phase was scanned. (b) Shot-noise level. Noise levels are displayed as the relative power level to the shotnoise level (0 dB). The settings of the spectrum analyzer were zero-span mode at 1 MHz, resolution bandwidth = 100 kHz, and video bandwidth = 30 Hz.⁵⁷

intracavity pumping, together with innovative tuning and stabilization techniques have enabled the generation of coherent output with excellent power and frequency stability and short-term linewidths down to a few kHz, long-term stability of a few MHz, and smooth mode-hop-free tuning of more than 100 GHz from cw OPOs. With the continuing advances in QPM material technology, reliable fabrication of MgO-doped PPLN has led to significant reductions in photorefractive damage, enabling its use at lower temperatures and higher powers with reduced output beam degradation. By deploying high-power fiber pump lasers in combination with MgO:PPLN, output powers in excess of 10 W have now been realized in cw SROs and power scaling to several tens of watts appears a clear possibility.

While photorefractive damage has placed limitations on the use of (MgO:)PPLN under visible pumping and confined operation of cw OPOs to near- and mid-IR above $\sim 1 \mu\text{m}$, advances in QPM material technology have enabled reliable fabrication of alternative QPM crystals with short grating periods, long interaction lengths, and immunity to photorefractive. This has paved the way for spectral extension of cw OPOs to the wavelengths below $1 \mu\text{m}$ by exploiting QPM materials including MgO:sPPLT and PPKTP and using direct pumping in the visible followed by additional frequency upconversion steps internal or external to the OPO cavity. By exploiting such techniques, spectral regions in orange and red have been successfully accessed and wavelengths down to 425 nm in the blue have been generated at several hundred milliwatts of output power in single-frequency spectrum and high beam quality. These developments have opened up new opportunities for the realization of practical solid-state sources with wide tunability across the visible and into the UV, where there is a severe shortage of conventional laser sources or other solid-state technologies.

The rapid advances in cw OPOs over the past decade have paved the way for the routine deployment of practical devices in a wide range of applications, from spectroscopy and imaging to frequency metrology and quantum optics. The unique capabilities of cw OPOs with regard to spectral versatility, output power, frequency and power stability, and compact solid-state design

open up new avenues for the ultimate development of these devices in real applications including portable gas detectors, mobile breath analyzers, handheld imaging systems or transportable optical frequency meters. At the same time, significant challenges and opportunities remain for further advancement of cw OPOs. In particular, wavelength extension of these devices into the mid-IR wavelength regions beyond 5 μm remains difficult because of the absorption of oxide-based QPM nonlinear materials, the key building blocks for the successful development of cw OPOs. The 5 to 12- μm spectral range is a particularly important and interesting region because of the presence of several atmospheric windows as well as many molecular absorption finger prints, including volatile substances in exhaled breath. The development of cw OPOs for this wavelength range will thus be of significant interest for spectroscopy, gas sensing, biomedicine, and atmospheric transmission. Progress toward the development of cw OPOs in this spectral range will require the deployment of alternative mid-IR birefringent nonlinear materials such as ZnGeP_2 , or the more recently developed nonlinear crystals such as CdSiP_2 and orientation-patterned GaAs, together with suitable near- to mid-IR laser pump sources and cascaded two-step pumping schemes. On the other hand, advancement of cw OPOs toward shorter UV wavelengths below the current 425-nm limit will be feasible with the use of recently demonstrated frequency upconversion techniques internal to cw SROs by exploiting existing birefringent crystals such as BBO, BIBO, and LBO. With the rapid advances in cw fiber laser technology and the potential for wavelength extension toward the extremes of the optical spectrum, the realization of compact, high-power and practical coherent solid-state light sources across the entire 300 to 12000 nm range based on cw OPOs appears a clear possibility in not too distant a future.

17.5 REFERENCES

1. M. Ebrahim-Zadeh and M. H. Dunn, "Optical Parametric Oscillators," in *OSA Handbook of Optics*, Vol. 4, (McGraw-Hill, New York, 2000) pp. 22.1–22.72.
2. J. A. Giordmaine and R. C. Miller, "Tunable Coherent Parametric Oscillation in LiNbO_3 at Optical Frequencies," *Phys. Rev. Lett.* **14**:973–976 (1965).
3. R. L. Byer and A. Piskarskas, "Optical Parametric Oscillation and Amplification," Special Issue, *J. Opt. Soc. Am. B* **10**:1656–1791; 2148–2243 (1993).
4. W. R. Bosenberg and R. C. Eckardt, "Optical Parametric Devices," Special issue, *J. Opt. Soc. Am. B* **12**: 2084–2322 (1995).
5. S. Schiller and J. Mlynek, "Continuous-Wave Optical Parametric Oscillators," Special Issue, *Appl. Phys. B* **66**: 661–764 (1998).
6. M. Ebrahim-Zadeh, R. C. Eckardt, and M. H. Dunn, "Optical Parametric Devices and Processes," Special Issue, *J. Opt. Soc. Am. B* **16**:1477–1602 (1999).
7. M. Ebrahim-Zadeh, "Mid-Infrared Ultrafast and Continuous-Wave Optical Parametric Oscillators," in *Solid-State Mid-Infrared Laser Sources* (Springer, Berlin, Heidelberg, 2003) pp. 179–218.
8. M. Ebrahim-Zadeh, "Optical Parametric Devices," in *Handbook of Laser Technology and Applications* (Institute of Physics Publishing, London, 2003) pp. 1347–1392.
9. K. L. Vodopyanov, "Pulsed Mid-Infrared Optical Parametric Oscillators," in *Solid-State Mid-Infrared Laser Sources*, I. T. Sorokina and K. L. Vodopyanov, eds. (Springer, Berlin, Heidelberg, 2003) pp. 141–178.
10. M. Ebrahim-Zadeh, "Mid-Infrared Optical Parametric Oscillators and Applications," in *Mid-Infrared Coherent Sources and Applications*, M. Ebrahim-Zadeh and I. T. Sorokins, eds. (Springer, Berlin, Heidelberg, 2007) pp. 347–375.
11. R. L. Sutherland, *Handbook of Nonlinear Optics* (Marcel Dekker, New York, 1996).
12. M. H. Dunn and M. Ebrahim-Zadeh, "Parametric Generation of Tunable Light from Continuous-Wave to Femtosecond Pulses," *Science* **286**:1513–1517 (1999).
13. M. Ebrahim-Zadeh, "Parametric Light Generation," *Phil. Trans. Roy. Soc. London A* **263**:2731–2750 (2003).
14. S. E. Harris, "Tunable Optical Parametric Oscillators," *Proc. IEEE* **57**:2096–2113 (1969).

15. R. L. Byer, "Optical Parametric Oscillators," in *Treatise in Quantum Electronics* (Academic Press, New York, 1973) pp. 587–702.
16. R. G. Smith, "Optical Parametric Oscillators," in *Lasers* (Marcel Dekker, New York, 1976) pp. 189–307.
17. I. D. Lindsay, "High Spatial and Spectral Quality Diode-Laser-Based Pump Sources for Solid-State Lasers and Optical Parametric Oscillators," Ph.D. Thesis, University of St Andrews (1999).
18. T. J. Edwards, G. A. Turnbull, M. H. Dunn, and M. Ebrahim-Zadeh, "Continuous-Wave, Singly-Resonant, Optical Parametric Oscillator Based on Periodically Poled KTiOPO_4 ," *Opt. Exp.* **16**:58–63 (2000).
19. S. E. Bisson, K. M. Armstrong, T. J. Kulp, and M. Hartings, "Broadly Tunable, Mode-Hop-Tuned CW Optical Parametric Oscillator Based on Periodically Poled Lithium Niobate," *Appl. Phys. B* **40**:6049–6055 (2001).
20. M. Van Herpen, S. te Lintel Hekkert, S. E. Bisson, and F. J. M. Harren, "Wide Single-Mode Tuning of a 3.0–3.8- μm , 700-mW, Continuous-Wave Nd:YAG-Pumped Optical Parametric Oscillator Based on Periodically-Poled Lithium Niobate," *Opt. Lett.* **27**:640–642 (2002).
21. M. M. J. W. Van Herpen, S. Li, S. E. Bisson, S. te Lintel Hekkert, and F. J. M. Harren, "Tuning and Stability of Continuous-Wave Mid-Infrared High-Power Single Resonant Optical Parametric Oscillator," *Appl. Phys. B* **75**:329–333 (2002).
22. M. M. J. W. Van Herpen, S. E. Bisson, and F. J. M. Harren, "Continuous-Wave Operation of a Single-Frequency Optical Parametric Oscillator at 4–5 μm Based on Periodically Poled LiNbO_3 ," *Opt. Lett.* **28**:2497–2499 (2003).
23. M. M. J. W. Van Herpen, S. E. Bisson, A. K. Y. Ngai, and F. J. M. Harren, "Combined Wide Pump Tuning and High Power of a Continuous-Wave, Singly Resonant Optical Parametric Oscillator," *Appl. Phys. B* **78**:281–286 (2004).
24. A. K. Y. Ngai, S. T. Persijn, G. Von Basum, and F. J. M. Harren, "Automatically Tunable Continuous-Wave Optical Parametric Oscillator for High-Resolution Spectroscopy and Sensitive Trace-Gas Detection," *Appl. Phys. B* **85**:173–180 (2006).
25. P. Gross, M. E. Klein, T. Walde, K.-J. Boller, M. Auerbach, P. Wessels, and C. Fallnich, "Fiber-Laser-Pumped Continuous-Wave Singly Resonant Optical Parametric Oscillator," *Opt. Lett.* **27**:418–420 (2002).
26. M. E. Klein, P. Gross, K.-J. Boller, M. Auerbach, P. Wessels, and C. Fallnich, "Rapidly Tunable Continuous-Wave Optical Parametric Oscillator Pumped by a Fiber Laser," *Opt. Lett.* **28**:920–922 (2003).
27. A. Henderson and R. Stafford, "Low Threshold, Singly Resonant CW OPO Pumped by an All-Fiber Pump Source," *Opt. Exp.* **14**:767–772 (2006).
28. M. E. Klein, C. K. Laue, D.-H. Lee, K.-J. Boller, and R. Wallenstein, "Diode-Pumped Singly Resonant Continuous-Wave Optical Parametric Oscillator with Wide Continuous Tuning of the Near-Infrared Idler Wave," *Opt. Lett.* **25**:490–492 (2000).
29. I. D. Lindsay, B. Adhimoolum, P. Gross, M. E. Klein, and K.-J. Boller, "110 GHz Rapid, Continuous Tuning from an Optical Parametric Oscillator Pumped by a Fiber-Amplified DBR Diode Laser," *Opt. Exp.* **13**:1234–1239 (2005).
30. A. Henderson and R. Stafford, "Intracavity Power Effects in Singly Resonant CW OPOs," *Appl. Phys. B* **85**:181–184 (2006).
31. G. K. Samanta and M. Ebrahim-Zadeh, "Continuous-Wave Singly Resonant Optical Parametric Oscillator with Resonant Wave Coupling," *Opt. Exp.* **16**:6883–6888 (2008).
32. A. V. Okishev and J. D. Zuegel, "Intracavity-Pumped Raman Laser Action in a Mid-IR, Continuous-Wave (CW) MgO:PPLN Optical Parametric Oscillator," *Opt. Exp.* **14**:12169–12173 (2006).
33. A. Henderson and R. Stafford, "Spectral Broadening and Stimulated Raman Conversion in a Continuous-Wave Optical Parametric Oscillator," *Opt. Lett.* **32**:1281–1283 (2007).
34. U. Strossner, J. P. Meyn, R. Wallenstein, P. Urenski, A. Arie, G. Rosenman, J. Mlynek, S. Schiller, and A. Peters, "Single-Frequency Continuous-Wave Optical Parametric Oscillator System with an Ultrawide Tuning Range of 550–2830 nm," *J. Opt. Soc. Am. B* **19**:1419–1424 (2002).
35. G. K. Samanta, G. R. Fayaz, Z. Sun, and M. Ebrahim-Zadeh, "High-Power, Continuous-Wave, Singly Resonant Optical Parametric Oscillator Based on MgO:sPPLT ," *Opt. Lett.* **32**:400–402 (2007).
36. G. K. Samanta, G. R. Fayaz, and M. Ebrahim-Zadeh, "1.59- μm , Single-Frequency, Continuous-Wave Optical Parametric Oscillator Based on MgO:sPPLT ," *Opt. Lett.* **32**:2623–2625 (2007).
37. J.-M. Melkonian, T.-H. My, F. Bretenaker, and C. Drag, "High Spectral Purity and Tunable Operation of a Continuous Singly Resonant Optical Parametric Oscillator Emitting in the Red," *Opt. Lett.* **32**:518–520 (2007).

38. G. K. Samanta and M. Ebrahim-Zadeh, "Continuous-Wave, Single-Frequency, Solid-State Blue Source for the 425–489 nm Spectral Range," *Opt. Lett.* **33**:1228–1230 (2008).
39. T.-H. My, C. Drag, and F. Bretenaker, "Single-Frequency and Tunable Operation of a Continuous Intracavity Frequency Doubled Singly Resonant Optical Parametric Oscillator," *Opt. Lett.* **33**:1455–1457 (2008).
40. C. Langrock and M. M. Fejer, "Fiber-Feedback Continuous-Wave and Synchronously Pumped Singly Resonant Ring Optical Parametric Oscillators Using Reverse-Proton-Exchanged Periodically Poled Lithium Niobate Waveguides," *Opt. Lett.* **32**:2263–2265 (2007).
41. A. J. Henderson, P. M. Roper, L. A. Borschowa, and R. D. Mead, "Stable, Continuously Tunable Operation of a Diode-Pumped Doubly Resonant Optical Parametric Oscillator," *Opt. Lett.* **25**:1264–1266 (2000).
42. T. Petelski, R. S. Conroy, K. Benecheikh, J. Mlynek, and S. Schiller, "All-Solid-State, Tunable, Single-Frequency Source of Yellow Light for High-Resolution Spectroscopy," *Opt. Lett.* **26**:1013–1015 (2001).
43. I. D. Lindsay, C. Petridis, M. H. Dunn, and M. Ebrahim-Zadeh, "Continuous-Wave Pump-Enhanced Singly Resonant Optical Parametric Oscillator Pumped by an Extended-Cavity Diode Laser," *Appl. Phys. Lett.* **78**:871–873 (2001).
44. G. A. Turnbull, D. McGloin, I. D. Lindsay, M. Ebrahim-Zadeh, and M. H. Dunn, "Extended Mode-Hop-Free Tuning by Use of Dual-Cavity, Pump-Enhanced Optical Parametric Oscillator," *Opt. Lett.* **25**:341–343 (2000).
45. D. J. M. Stothard, I. D. Lindsay, M. H. Dunn, "Continuous-Wave Pump-Enhanced Optical Parametric Oscillator with Ring Resonator for Wide and Continuous Tuning of Single-Frequency Radiation," *Opt. Exp.* **12**:502–511 (2004).
46. I. D. Lindsay, D. J. M. Stothard, C. F. Rae, and M. H. Dunn, "Continuous-Wave Pump-Enhanced Optical Parametric Oscillator Based on Periodically Poled RbTiOAsO₄," *Opt. Exp.* **11**:134–140 (2003).
47. F. Muller, G. Von Basum, A. Pop, D. Halmer, P. Hering, M. Murtz, F. Kunnermann, and S. Schiller, "Long-Term Frequency Stability and Linewidth Properties of Continuous-Wave Pump-Resonant Optical Parametric Oscillators," *Appl. Phys. B* **80**:307–313 (2005).
48. P. Gross, M. E. Klein, H. Ridderbusch, D.-H. Lee, J.-P. Meyn, R. Wallenstein, and K.-J. Boller, "Wide Wavelength Tuning of an Optical Parametric Oscillator Through Electro-Optic Shaping of the Gain Spectrum," *Opt. Lett.* **27**:1433–1435 (2002).
49. G. Ferrari, "Generating Green to Red Light with Semiconductor Lasers," *Opt. Exp.* **15**:1672–1678 (2007).
50. A. K. Y. Ngai, S. T. Persijn, I. D. Lindsay, A. A. Kosterev, P. Gross, C. J. Lee, S. M. Cristescu, F. K. Tittel, K.-J. Boller, and F. J. M. Harren, "Continuous Wave Optical Parametric Oscillator for Quartz-Enhanced Photoacoustic Trace Gas Sensing," *Appl. Phys. B* **89**:123–128 (2007).
51. E. V. Kovalchuk, D. Dekorsy, A. I. Lvovsky, C. Braxmaier, J. Mlynek, A. Peters, and S. Schiller, "High-Resolution Doppler-Free Molecular Spectroscopy with a Continuous-Wave Optical Parametric Oscillator," *Opt. Lett.* **26**:1430–1432 (2001).
52. F. Muller, A. Popp, F. Kuhnemann, and S. Schiller, "Transportable, Highly Sensitive Photoacoustic Spectrometer Based on Continuous-Wave Dual-Cavity Optical Parametric Oscillator," *Opt. Exp.* **11**:2820–2825 (2003).
53. G. von Basum, D. Halmer, P. Hering, M. Murtz, S. Schiller, F. Muller, A. Popp, and F. Kuhnemann, "Parts per Trillion Sensitivity for Ethane in Air with an Optical Parametric Oscillator Cavity Leak-Out Spectrometer," *Opt. Lett.* **29**:797–799 (2004).
54. D. J. M. Stothard, M. H. Dunn, and C. F. Rae, "Hyperspectral Imaging of Gases with a Continuous-Wave Pump-Enhanced Optical Parametric Oscillator," *Opt. Exp.* **12**:947–955 (2004).
55. K. Hayasaka, Y. Zhang, and K. Kasai, "Generation of Twin Beams from an Optical Parametric Oscillator Pumped by a Frequency-Doubled Diode Laser," *Opt. Lett.* **29**:1665–1667 (2004).
56. X. Su, A. Tan, X. Jia, Q. Pan, C. Xie, and K. Peng, "Experimental Demonstration of Quantum Entanglement between Frequency Non-Degenerate Optical Twin Beams," *Opt. Lett.* **31**:1133–1135 (2006).
57. T. Tanimura, D. Akamatsu, Y. Yokoi, A. Furusawa, and M. Kozuma, "Generation of a Squeezed Vacuum Resonant on a Rubidium *D₁* Line with Periodically Poled KTiOPO₄," *Opt. Lett.* **31**:2344–2364 (2006).
58. E. V. Kovalchuk, T. Schuldt, and A. Peters, "Combination of a Continuous-Wave Optical Parametric Oscillator and a Femtosecond Frequency Comb for Optical Frequency Metrology," *Opt. Lett.* **30**:3141–3143 (2005).

This page intentionally left blank.

DO NOT DUPLICATE

NONLINEAR OPTICAL PROCESSES FOR ULTRASHORT PULSE GENERATION

Uwe Siegner and Ursula Keller

*Institute of Quantum Electronics
Physics Department
Swiss Federal Institute of Technology (ETH)
Zurich, Switzerland*

18.1 GLOSSARY

- A_A laser beam cross section on the absorber
 A_L laser beam cross section in the laser gain material. If the cavity mode is not constant inside the gain medium, this area corresponds to an effective averaged value.
 c light velocity in vacuum
 D_2 total group delay dispersion inside the laser cavity per cavity round-trip (i.e., second-order dispersion)
 d sample thickness
 DR differential reflectivity
 DT differential transmission
 E_p intracavity pulse energy
 $E_{\text{sat},A}$ saturation energy of the saturable absorber
 $E_{\text{sat},L}$ saturation energy of laser material
 $F_{p,A}$ pulse fluence on the absorber

$$F_{p,A} = \frac{E_p}{A_A} = \int I_A(t) dt$$

- $F_{\text{sat},A}$ saturation fluence of the absorber

$$F_{\text{sat},A} = \frac{E_{\text{sat},A}}{A_A}$$

- $F_{\text{sat},L}$ saturation fluence of laser material

$$F_{\text{sat},L} = \frac{h\nu}{2\sigma_L} = \frac{E_{\text{sat},L}}{A_L}$$

The factor of 2 is used in case of a linear resonator with a standing wave.

g	saturated amplitude gain cross section
g_0	small signal amplitude gain coefficient
$I(t)$	time-dependent intensity
$I_A(t)$	time-dependent intensity on absorber
$I_L(t)$	time-dependent intensity inside the laser gain material
$I_{\text{sat},A}$	saturation intensity of the absorber

$$I_{\text{sat},A} = \frac{F_{\text{sat},A}}{\tau_A}$$

$I_{\text{sat},L}$	saturation intensity of laser material
--------------------	--

$$I_{\text{sat},L} = \frac{F_{\text{sat},L}}{\tau_L}$$

k	wave number in vacuum $k = \omega/c = 2\pi/\lambda$
k_n	wave number in a dielectric with refractive index n : $k_n = kn$
\mathbf{k}_1	wave vector of the pump pulse in a pump-probe experiment or of the first pulse applied in an FWM experiment for positive time delay Δt
\mathbf{k}_2	wave vector of the probe pulse in a pump-probe experiment or of the second pulse applied in an FWM experiment for positive time delay Δt
L_L	length of the laser gain medium
n_2	nonlinear refractive index
$P(t)$	time-dependent power: $E_p = \int P(t)dt$
$q(t)$	saturable amplitude loss coefficient (does not include any nonsaturable losses)
q_0	unsaturated amplitude loss coefficient, also corresponds to the maximal loss coefficient
α	intensity absorption constant
$\Delta\alpha$	nonlinear change of the intensity absorption constant
ΔR	maximum modulation depth of a saturable absorber integrated within a mirror structure (i.e., the maximum nonlinear intensity reflectivity change)
ΔR_{ns}	nonsaturable loss of a saturable absorber integrated within a mirror structure
Δt	time delay between the excitation pulses in a pump-probe or FWM experiment
ν	frequency
ν_0	center frequency
$\Delta\nu_g$	FWHM (full width half maximum) gain bandwidth
λ	wavelength in vacuum
λ_n	wavelength in a dielectric medium with a refractive index n
λ_0	center wavelength
σ_L	gain cross section: $F_{\text{sat},L} = h\nu/\sigma_L$, where h is the Planck's constant
σ_A	absorber cross section: $F_{\text{sat},A} = h\nu/\sigma_A$
τ_p	FWHM (full width half maximum) of pulse intensity
τ_A	recovery time of the saturable absorber
τ_L	upper state lifetime of laser medium

18.2 ABBREVIATIONS

A-FPSA	antiresonant Fabry-Perot saturable absorber
APM	additive pulse modelocking
As_{Ga}	arsenic antisite
As_{Ga}^0	neutral arsenic antisite
As_{Ga}^+	ionized arsenic antisite
CB	conduction band
CPM	colliding pulse modelocking
FWHM	full width half maximum
FWM	four-wave mixing
KLM	Kerr lens modelocking
LT	low temperature
MBE	molecular beam epitaxy
MOCVD	metal-organic chemical vapor deposition
SAM	self-amplitude modulation
SBR	saturable Bragg reflector
SESAM	semiconductor saturable absorber mirror
SPM	self-phase modulation
V_{Ga}	Ga vacancy
VB	valence band

18.3 INTRODUCTION

Since 1990 we have observed tremendous progress in ultrafast laser sources—a development that was triggered by the invention of the Ti:sapphire laser.¹ The strong interest in all-solid-state ultrafast laser technology was the driving force and formed the basis for many new inventions and discoveries. Solid-state lasers provide some of the best laser qualities in terms of high-quality spatial modes, high output power, energy storage and large pulse energy when Q-switched, and large optical bandwidth necessary for ultrashort pulse generation. Today, the most important nonlinear optical processes that support short and ultrashort passive pulse generation are based on Kerr and/or semiconductor nonlinearities. Saturable absorbers based on either the Kerr effect or on semiconductors play a major role in ultrashort pulse generation. However, independent of the specific saturable absorber material or mechanism, we can define a few macroscopic saturable absorber parameters that will determine the pulse formation process. These parameters are defined in Sec. 18.4 and discussed in more detail for the specific cases of a slow and a fast saturable absorber. The material properties then can be modified over an even larger range if the absorber is integrated within a device structure. This will be discussed at the end of Sec. 18.4. Stable pulse generation can then be obtained when these macroscopic saturable absorber parameters are designed correctly. The crucial role of the optical Kerr effect in ultrashort pulse generation will be discussed in Sec. 18.5. Optimization of semiconductor saturable absorber parameters normally requires a better understanding of the underlying physics, that is, the microscopic properties (Sec. 18.6). Ultrafast nonlinear optical processes can be analyzed in various ways to extract the information about the material system involved in the process. Numerous studies with ultrafast laser pulses have been performed in atomic, molecular, and condensed matter systems. Thus, some of the most common experimental techniques in ultrafast spectroscopy will also be reviewed in Sec. 18.6. The results of such ultrafast spectroscopy measurements are then summarized for semiconductor materials because of their importance in ultrashort pulse generation.

Today ultrafast lasers demonstrate unsurpassed performances: pulse duration in the two-cycle regime,^{2,3} compact and reliable picosecond and femtosecond all-solid-state lasers,^{4,5} pulse repetition rates in the 100-GHz regime,^{6,7} average powers well above 10 W,⁸ and novel Q-switching performances that bridge the gap between modelocking and Q-switching both in terms of pulse durations and pulse repetition rates.⁹ A more recent review is given in Ref. 10 with detailed tables summarizing the different modelocking results for many solid-state lasers. Further average power scaling towards the 1-kW regime looks very promising with SESAM modelocked Yb-doped thin disk lasers¹¹ and this concept could be transferred to optically pumped semiconductor lasers.^{12,13}

Optical pulses in the two-cycle regime have been produced by a variety of methods: direct generation in a modelocked laser oscillator,^{2,3} continuum generation together with parametric optical amplification,¹⁴ and external pulse compression.^{15,16} In Fig. 1, interferometric autocorrelation measurements of the two-cycle pulses are shown. Although this is a relatively crude characterization technique, it is the only common measurement that has been performed on these sources; thus it allows for a direct comparison. These pulse generation techniques rely essentially on three identical ingredients:¹⁷ (1) an ultrabroadband amplifying process, (2) precise control of dispersion, and (3) the nonlinear optical Kerr effect.

Ultrafast all-solid-state lasers are based on diode-pumped solid-state lasers. Semiconductor saturable absorbers were the first intracavity saturable absorbers that reliably started and sustained passive modelocked diode-pumped solid-state lasers.^{18,19} Any previous attempts to passively modelock such lasers with intracavity saturable absorbers resulted in Q-switching instabilities or Q-switched modelocking (see Fig. 2). The precise control of optical nonlinearities was necessary to resolve

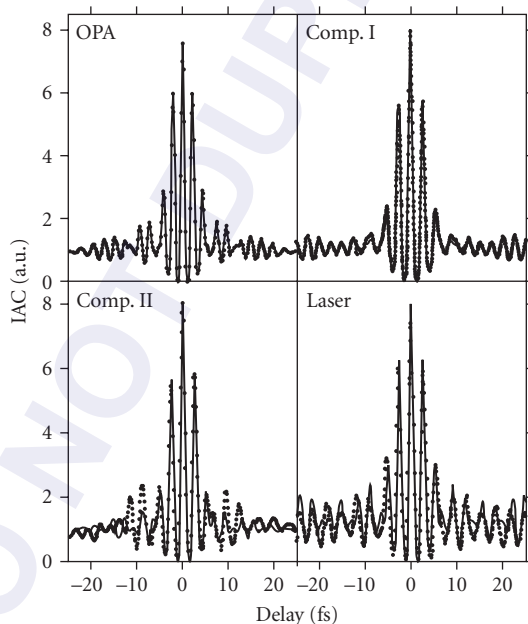


FIGURE 1 Interferometric autocorrelation measurements of different sources in the 5-fs range: OPA: optical parametric amplification.¹⁴ Comp. I: compression of Ti:sapphire cavity-dumped pulses in silica fiber.¹⁵ Comp. II: compression of amplified μJ pulses in a hollow fiber filled with Krypton.¹⁶ Laser: pulses directly obtained from a Ti:sapphire laser without any external pulse compression.² Dots are measured data and lines correspond to fits that were used to estimate pulse duration from autocorrelation.

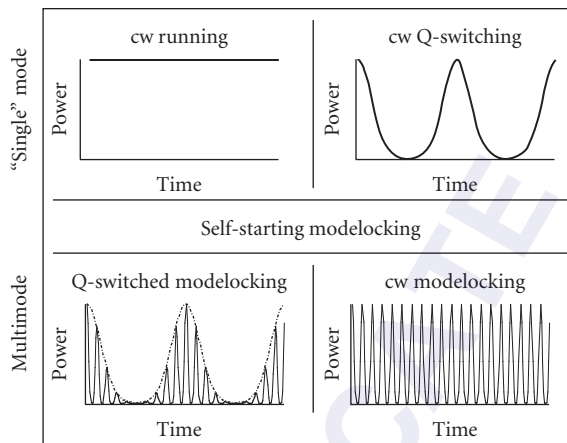


FIGURE 2 Different modes of operation of a laser with a saturable absorber. Continuous wave (cw) Q-switching typically occurs with much longer pulses and lower pulse repetition rates than cw mode-locking.

the Q-switching problem. Such a control can be achieved with epitaxially grown semiconductor saturable absorbers, which makes them very attractive for use as saturable absorbers in solid-state lasers. Semiconductor saturable absorbers provide a variety of bandgaps, ranging from the visible to the infrared, and can be easily integrated into different device structures (such as a mirror, for example), which allows for the control of the absorber parameters over an even larger range. This is actually required to obtain stable pulse generation with many solid-state lasers. The specific saturable absorber nonlinearities required for stable cw modelocking or Q-switching have recently been discussed in much more detail in Ref. 4 and will be briefly summarized in Sec. 18.4.

In this section we will put the emphasis on the nonlinearities used for ultrashort pulse generation. It is not our goal to provide a tutorial for pulse generation techniques such as passive modelocking or Q-switching. A more tutorial-type overview of the physics of ultrashort pulse generation was given in Refs. 5, 10, 20, and 21. We would like to refer the interested readers to those chapters and the extensive references therein.

18.4 SATURABLE ABSORBERS: MACROSCOPIC DESCRIPTION

Saturable Absorber: Self-Amplitude Modulation

Saturable absorbers have been used to passively Q-switch and modelock many different lasers. Different saturable absorbers, such as organic dyes, colored filter glasses, dye-doped crystals, and semiconductors have been used. Independent of the specific saturable absorber material, we can define a few macroscopic absorber parameters that will determine the pulse generation process. The macroscopic properties of a saturable absorber are the modulation depth, the nonsaturable loss, the saturation fluence, the saturation intensity, and the impulse response or recovery times. These parameters determine the operation of a passively mode locked or Q-switched laser. In our notation we assume that the saturable absorber is integrated within a mirror structure. Thus we are interested in the nonlinear reflectivity change or the *differential reflectivity* $DR(t)$ as a function of time or in the reflectivity $R(F_{p,A})$ as a function of the incident pulse energy fluence on the saturable absorber. If the saturable absorber is used in transmission, we simply characterize the absorber by nonlinear transmission measurements. Both the saturation fluence $F_{\text{sat},A}$ and the absorber recovery time τ_A are determined

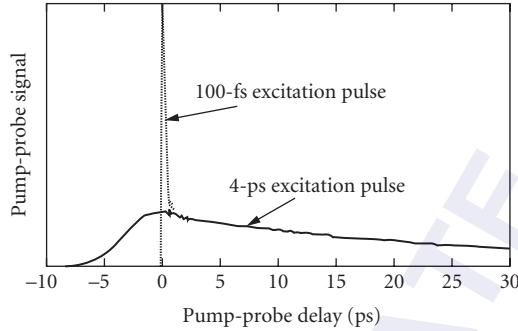


FIGURE 3 Standard pump-probe techniques determine the impulse response $DR(t)$. (Here we assume that the saturable absorber is integrated within a mirror structure.) $DR(t)$ for the same saturable absorber is different for different excitation pulse durations. For excitation with a picosecond pulse, the pump-probe trace clearly shows no significant modulation depth with a fast time constant.

experimentally without any need to determine the microscopic properties of the nonlinearities. Thus, the saturation fluence of the absorber is not only dependent on material properties but also on the specific device structure the absorber is integrated in.

Standard pump-probe techniques determine the impulse response $DR(t)$ and therefore τ_A (see Fig. 3). In the picosecond regime we normally only have to consider one recovery time, because much faster femtosecond nonlinearities in the saturable absorber result in a negligible modulation depth. This is shown in Fig. 3, where the impulse response $DR(t)$ was measured for two different excitation pulse durations. For excitation with a picosecond pulse, the pump-probe trace clearly shows no significant modulation depth with a fast time constant. In the femtosecond pulse regime we normally have to consider more than one absorber recovery time. In this case the slow component normally helps to start the initial pulse formation process. The modulation depth of the fast component then determines the pulse duration at steady state. Further improvements of the saturable absorber normally require some better understanding of the underlying physics, which will be discussed in more detail in Sec. 18.6.

The saturation fluence $F_{\text{sat}, A}$ is determined and defined by the measurement of the nonlinear change in reflectivity $R(F_{p, A})$ as a function of increased incident pulse fluence (see Fig. 4). The common traveling wave rate equations²² in the slow absorber approximation normally give a very good fit and determine the saturation fluence $F_{\text{sat}, A}$, modulation depth ΔR , and nonsaturable losses ΔR_{ns} of the absorber. The modulation depth is typically small to prevent Q-switching instabilities in passively modelocked solid-state lasers.²³ Thus it is reasonable to make the following approximation:

$$\Delta R = 1 - e^{-2q_0} \approx 2q_0, q_0 \ll 1 \quad (1)$$

where q_0 is the unsaturated amplitude loss coefficient. Here, it is assumed that the nonsaturable losses are negligible.

The saturation of an absorber can be described with the following differential equation:²²

$$\frac{dq(t)}{dt} = -\frac{q(t) - q_0}{\tau_A} - \frac{q(t)P(t)}{E_{\text{sat}, A}} \quad (2)$$

where $q(t)$ is the saturable amplitude loss coefficient that does not include any nonsaturable losses. At any time t the reflected (or transmitted) intensity $I_{\text{out}}(t)$ from the saturable absorber is given by

$$I_{\text{out}}(t) = R(t)I_{\text{in}}(t) = e^{-2q(t)}I_{\text{in}}(t) \quad (3)$$

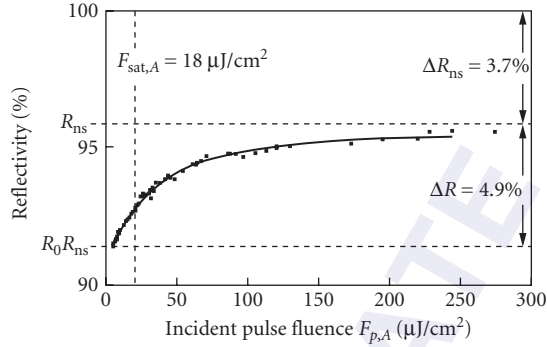


FIGURE 4 Nonlinear Reflectivity as function of the incident pulse fluence $R(F_{p,A})$. The measurement was made with a SESAM supporting 34-fs pulses in a Ti:sapphire laser. The measurements are fitted with common traveling wave rate equations²² in the slow absorber approximation which normally give a very good fit and determine the saturation fluence $F_{sat,A}$, modulation depth ΔR , and nonsaturable losses ΔR_{ns} of the absorber. Detailed guidelines how to measure the macroscopic SESAM parameters are given in Refs. 24 and 25.

Then the total net reflectivity is given by

$$R_{tot} = \frac{\int I_{out}(t) dt}{\int I_{in}(t) dt} = \frac{F_{out}}{F_{in}} = 1 - \frac{2}{F_{in}} \int q(t) I_{in}(t) dt \quad (4)$$

This determines the total absorber loss coefficient q_p , which results from the fact that part of the excitation pulse needs to be absorbed to saturate the absorber:

$$R_{tot} = e^{-2q_p} \approx 1 - 2q_p \quad (5)$$

From Eqs. (4) and (5) it then follows that for $q_p \ll 1$

$$q_p = \frac{1}{F_{in}} \int q(t) I_{in}(t) dt = \int q(t) f(t) dt \quad (6)$$

where

$$f(t) \equiv \frac{I_{in}(t)}{F_{in}} = \frac{P_{in}(t)}{E_{p,in}}, \quad \text{with} \quad \int f(t) dt = \frac{1}{F_{in}} \int I_{in}(t) dt = 1 \quad (7)$$

We then distinguish between two typical cases: a slow and a fast saturable absorber.

Slow Saturable Absorber

In the case of a slow saturable absorber, we assume that the excitation pulse duration is much shorter than the recovery time of the absorber (i.e., $\tau_p \ll \tau_A$). Thus, we can neglect the recovery of the absorber during pulse excitation, and Eq. (2) reduces to:

$$\frac{dq(t)}{dt} \approx -\frac{q(t)P(t)}{E_{sat,A}} \quad (8)$$

This differential equation can be solved, and we obtain for the self-amplitude modulation (SAM):

$$q(t) = q_0 \exp \left[-\frac{E_p}{E_{\text{sat},A}} \int_0^t f(t') dt' \right] \quad (9)$$

Equation (6) then determines the total absorber loss coefficient for a given incident pulse fluence $F_{p,A}$:

$$q_p(F_{p,A}) = \int q(t) f(t) dt = q_0 \frac{F_{\text{sat},A}}{F_{p,A}} (1 - e^{-F_{p,A}/F_{\text{sat},A}}) \quad (10)$$

It is not surprising that q_p does not depend on any specific pulse form, because $\tau_p \ll \tau_A$.

A slow saturable absorber has been successfully used to passively modelock dye and semiconductor lasers. In this case, dynamic gain saturation was supporting the pulse formation process (Fig. 5a), and much shorter pulses than the recovery time of the saturable absorber were obtained.^{26,27} Dynamic gain saturation means that the gain experiences a fast pulse-induced saturation that then recovers again between consecutive pulses. Therefore, an ultrashort net-gain window can be formed by the combined saturation of absorber and gain for which the absorber has to saturate and recover faster than the gain, while the recovery time of the saturable absorber can be much longer than the pulse duration. Haus's master equation formalism can describe this passive modelocking technique very well.²⁸ This formalism is based on linearized differential operators that describe the temporal evolution of a pulse envelope inside the laser cavity. Generally, the linearized master equations of Haus describe modelocking very well, as long as they only have to deal with small nonlinearities and loss modulation. This is true for most lasers; otherwise a strong tendency for instabilities is observed. A review of Haus's modelocking formalism is presented in more recent articles and book chapters.^{5,10,21,29} Assuming passive modelocking according to Fig. 5a, Haus predicts a mode locked pulse duration for a fully saturated absorber

$$\tau \approx \frac{4}{\pi} \frac{1}{\Delta \nu_g} \quad (11)$$

$$\tau_p \approx 1.76 \times \frac{4}{\pi} \frac{1}{\Delta \nu_g}$$

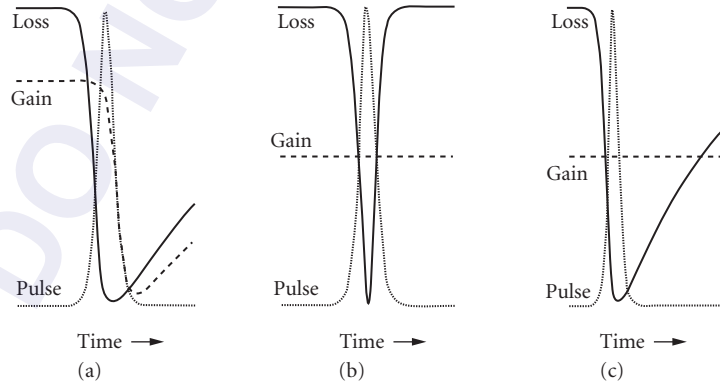


FIGURE 5 Passive modelocking mechanisms explained by three fundamental models: (a) slow saturable absorber modelocking with dynamic gain saturation, (b) fast saturable absorber modelocking, and (c) soliton modelocking.

with a predicted pulse form of

$$I(t) = I_0 \operatorname{sech}^2(t/\tau), \quad \tau_p = 1.76 \cdot \tau \quad (12)$$

where τ_p is the FWHM pulse width of the pulse intensity. For dye lasers using Rhodamin 6G with a gain bandwidth $\Delta\nu \approx 4 \cdot 10^{13}$ Hz and DODCI as the saturable absorber with an absorber cross section $\sigma_A = 0.52 \cdot 10^{-16}$ cm², we would then predict a pulse duration τ_p of 56 fs at a center wavelength of 620 nm [Eqs. (11) and (12)]. Slightly shorter pulses of 27 fs duration were demonstrated.^{30,31} The interplay of self-phase modulation and negative group velocity dispersion can result in pulses that are shorter than would be predicted by the SAM alone [Eq. (10)]. Martinez estimated the additional pulse shortening to be about a factor of 2.^{32,33} This will be discussed in Sec. 18.5.

For solid-state lasers we cannot apply slow saturable absorber modelocking as shown in Fig. 5a, because no significant dynamic gain saturation is taking place, due to the small gain cross section and the long upper state lifetime of the laser. The gain cross section of ion-doped solid-state lasers is typically 10^{-19} cm² and smaller. This is at least 1000 times smaller than dye, semiconductor, or color center lasers. In addition, the upper-state lifetime of ion-doped solid-state lasers is typically in the μ s to ms regime—much longer than the pulse repetition period that is typically in the ns regime. We therefore do not observe any significant dynamic gain saturation, and the gain is only saturated to a constant level by the average intracavity intensity (Fig. 5b and c).

Fast Saturable Absorber

In the case of a fast saturable absorber, the absorber recovery time is much faster than the pulse duration (i.e., $\tau_p \gg \tau_A$). Thus, we can assume that the absorption instantaneously follows the absorption of a certain power $P(t)$, and Eq. (2) reduces to

$$0 = -\frac{q(t) - q_0}{\tau_A} - \frac{q(t)P(t)}{E_{\text{sat},A}} \quad (13)$$

The saturation of the fast absorber then follows directly from Eq. (13):

$$q(t) = \frac{q_0}{1 + \frac{I_A(t)}{I_{\text{sat},A}}} \quad (14)$$

where we used the fact that $P_{\text{sat},A} = E_{\text{sat},A}/\tau_A$ and $P(t)/P_{\text{sat},A} = I_A(t)/I_{\text{sat},A}$. In the linear regime we can make the following approximation in Eq. (14):

$$q(t) \approx q_0 - \gamma I_A(t), \quad \text{with} \quad \gamma \equiv \frac{q_0}{I_{\text{sat},A}} \quad (15)$$

The total absorber loss coefficient q_p [Eqs. (10) to (12)] now depends on the pulse form, and for a sech^2 -pulse shape we obtain for an incident pulse fluence $F_{p,A}$ and the linear approximation of $q(t)$ [Eq. (15)]:

$$q_p(F_{p,A}) = \frac{1}{F_{p,A}} \int q(t) I_A(t) dt = q_0 \left(1 - \frac{1}{3} \frac{F_{p,A}}{\tau I_{\text{sat},A}} \right) \quad (16)$$

We only obtain an analytic solution for fast saturable absorber modelocking if we assume an ideal fast absorber that saturates linearly with pulse intensity over the full modulation depth [Eq. (15)].

For a maximum modulation depth, we then can assume that $q_0 = \mathcal{I}_{0,A}^s$, where $I_{0,A}$ is the peak intensity on the saturable absorber. We then obtain with Eq. (16) a residual saturable absorber loss of $q_0/3$, which the pulse experiences to fully saturate the ideal fast saturable absorber.

A fast saturable absorber has been successfully used to passively modelock solid-state lasers (Fig. 5b). An analytical solution with Haus's master formalism³⁴ is only obtained if we assume an ideal fast saturable absorber that produces a decreased loss directly proportional to intensity of the incident laser pulse [Eq. (15)]. Then again a sech²-pulse shape [Eq. (12)] with the following pulse duration is predicted:

$$\tau_p = 1.76 \frac{4D_g}{\gamma F_{p,A}}, \quad \text{with} \quad D_g = \frac{g}{\pi^2 (\Delta\nu_g)^2} \quad (17)$$

where D_g is the gain dispersion. Unfortunately, fast saturable absorbers with femtosecond recovery times are often not sufficient for reliable self-starting of the modelocking process. In passive modelocking, pulse formation should start from normal noise fluctuations in a laser. One noise spike is strong enough to start saturating the absorber and thereby lowers the loss. This noise spike begins to grow in amplitude and becomes shorter until a stable pulse duration is obtained. Initially, these noise spike durations are on the order of the cavity round-trip time, introducing only very small loss modulations in a fast saturable absorber (see Fig. 3). A combination of a fast and a slow saturable absorber can help to solve this problem. Thus, semiconductor saturable absorbers are very interesting because they typically have a bitemporal impulse response which even can be modified.

For picosecond solid-state lasers, the self-amplitude modulation of a fast saturable absorber with a picosecond recovery time is sufficient for stable pulse generation. A picosecond recovery time can be achieved with low-temperature-grown semiconductor saturable absorbers where midgap defect states form very efficient traps for the photoexcited electrons in the conduction band. A more detailed description of the microscopic nonlinearities is given in Sec. 18.6. In the picosecond regime, we developed a very simple stability criteria for stable passive modelocking without Q-switching instabilities:²³

$$E_p^2 > E_{p,c}^2 = E_{\text{sat},L} E_{\text{sat},A} \Delta R \quad (18)$$

The critical intracavity pulse energy $E_{p,c}$ is the minimum intracavity pulse energy that is required to obtain stable cw modelocking; that is, for $E_p > E_{p,c}$ we obtain stable cw modelocking and for $E_p < E_{p,c}$ we obtain Q-switched modelocking (see Fig. 2). For good stability of a mode locked laser against unwanted fluctuations of pulse energy, operation close to the stability limit is not recommended. Thus, a large modulation depth supports shorter pulses [Eqs. (1), (15), and (17)], but an upper limit is given by the onset of self-Q-switching instabilities [Eq. (18)].

Semiconductor Saturable Absorber Mirrors

Semiconductor saturable absorbers were used as early as 1974 in CO₂ lasers and as early as 1980 for semiconductor diode lasers. A color center laser was the first solid-state laser that was cw modelocked with an intracavity semiconductor saturable absorber.³⁵ However, for both the diode and color center laser, dynamic gain saturation supported pulse formation, and a slow saturable absorber was sufficient for pulse generation (Fig. 5a). In addition, because of the much larger gain cross section and therefore smaller saturation energy $E_{\text{sat},L}$ [Eq. (18)] (i.e., typically 1000 to 10,000 times smaller than ion-doped solid-state lasers), Q-switching instabilities were not a problem. Thus, semiconductor saturable absorber parameters (see Figs. 3 and 4) have to be chosen much more carefully for stable cw modelocking.

We typically integrate the semiconductor saturable absorber into a mirror structure that results in a device whose reflectivity increases as the incident optical intensity increases. This general class of device is called a *semiconductor saturable absorber mirror* (SESAM).^{18,36} A detailed description and guidelines on how to design a SESAM for either passive modelocking or Q-switching for different laser parameters is

given in Ref. 4. Such a SESAM device structure can be a simple Bragg mirror, where at least one quarter-wave-layer contains an absorber layer (saturable Bragg reflector -SBR).^{37–39} A larger parameter range for the saturation fluence and modulation depth can be achieved if the saturable absorber is integrated inside a Fabry-Perot structure that is operated at antiresonance (A-FPSA).^{18,40} The antiresonant Fabry-Perot structure is broadband and can be designed to have no significant bandwidth limitations even in the sub-10-fs pulse-width regime.⁴¹ The top reflector of the A-FPSA provides an adjustable parameter that determines the intensity entering the semiconductor saturable absorber and therefore the saturation fluence of the saturable absorber device. A $\text{SiO}_2/\text{TiO}_2$ dielectric top reflector has an additional advantage in that the damage threshold for this SESAM design is significantly higher compared to a SESAM design with the same saturation fluence but based on semiconductor materials alone. We also distinguish between resonant and antiresonant SESAM designs. There are different trade-offs between these design regimes, but a compromise can be found to obtain both low saturation fluence and sufficiently small group delay dispersion.⁴²

More recently, semiconductor-doped dielectric films have been demonstrated for saturable absorber applications.⁴³ Semiconductor-doped glasses have been used to modelock lasers as early as 1990.⁴⁴ However, the recently developed InAs-doped thin-film rf-sputtering technology offers similar advantages as SESAMs, which allows for the integration of the absorber into a device structure. At this point, however, the saturation fluence of $\approx 10 \text{ mJ/cm}^2$ is still rather high for stable solid-state laser modelocking [Eq. (18)]. In comparison, MBE or MOCVD-grown SESAMs have typically a saturation fluence of $\approx 10 \mu\text{J/cm}^2$, even though they can be modified from the $\mu\text{J/cm}^2$ to the mJ/cm^2 -range depending on the specific device structure.^{4,36}

A more detailed discussion of microscopic semiconductor nonlinearities will be presented in Sec. 18.6. In principle, knowledge of the macroscopic absorber parameters is sufficient to understand pulse generation. However, further improvements of saturable absorbers will require a more detailed understanding of microscopic optical nonlinearities.

18.5 KERR EFFECT

Longitudinal and Transverse Kerr Effect

The extremely rapid response and broad bandwidth of the Kerr nonlinearity are very attractive for a modelocking process. For high intensities, polarization inside a dielectric medium does not proportionally follow the electric field anymore. This gives rise to an index change that is proportional to intensity. Off-resonance, this nonlinear optical effect is extremely fast, with estimated response times in the few-femtosecond range. The transverse and longitudinal effects resulting from the intensity dependence are shown schematically in Fig. 6. The transverse Kerr effect retards the central and most intense part of a plane wavefront, thus acting as a focusing lens, referred to as the Kerr lens. Later we will indicate concepts that make use of the Kerr lens to produce short pulses. Along the axis of propagation, the longitudinal Kerr effect retards the center of an optical pulse, producing a red shift of the leading part of the pulse and a blue shift in the trailing part. Consequently, the longitudinal Kerr effect has been named *self-phase modulation* (SPM).

Longitudinal Kerr Effect for External Pulse Compression

SPM generates extra bandwidth; in other words, it spectrally broadens the pulse. SPM alone does not modify the pulse width—but a much shorter pulse can be generated with the extra bandwidth and proper dispersion compensation.⁴⁵ To create a short pulse, the blue spectral components have to be advanced relative to the red ones, exactly counteracting the phase delays induced by the SPM (assuming $n_2 > 0$). To do so an effect opposite to the normal material dispersion is needed. This type of dispersion is called *anomalous* or *negative* dispersion. A careful balance between a nonlinear spectral broadening process and negative dispersion is needed for efficient compression of a pulse.

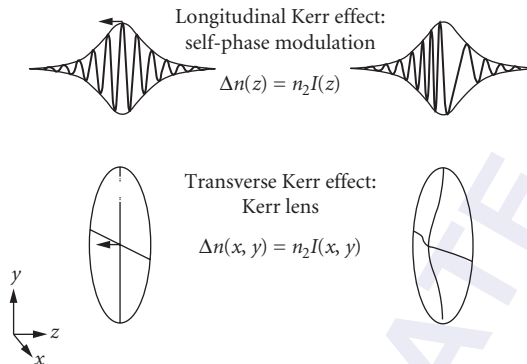


FIGURE 6 The Kerr effect produces a nonlinear intensity-dependent refractive index change. The longitudinal Kerr effect (top) causes self-phase modulation (SPM). The transverse Kerr effect (bottom) causes a nonlinear focusing lens (Kerr lens).

Typically, self-phase modulation in a single-mode fiber is used to chirp the pulse, which is then compressed with a grating pair compressor.⁴⁶

Ultimately, compression schemes are limited by uncompensated higher-order dispersion and higher-order nonlinearities. For pulses shorter than 100 fs, compression is typically limited to factors of less than 10. Compression of amplified CPM dye laser pulses with 50 fs duration produced the long-standing world record of 6 fs for short pulses.⁴⁷ Similar concepts have been recently used for external pulse compression of 13-fs pulses from a cavity dumped Ti:sapphire laser¹⁵ and of 20-fs pulses from a Ti:sapphire laser amplifier¹⁶ resulting, in both cases, in approximately 4.5-fs pulses (see Fig. 1). In the latter case, the use of a noble-gas-filled hollow fiber resulted in unsurpassed pulse energies of about 0.5 mJ, with 5.2-fs pulses and a peak power of 0.1 TW.⁴⁸

Longitudinal Kerr Effect for Broadband Parametric Amplification

The extra bandwidth obtained with SPM can be extremely large, producing a white-light continuum⁴⁹ that can be used as a seed for broadband parametric amplification. Parametric processes can provide amplification with even broader bandwidth than can typically be achieved in laser amplifiers. Noncollinear phase-matching at a crossing angle of 3.8° in Barium beta borate (BBO) provides more than 150 THz amplification bandwidth.⁵⁰ With this type of setup, parametric amplification has been successfully demonstrated with pulse durations of less than 5 fs.¹⁴

Longitudinal Kerr Effect for Passive Modelocking: Soliton Modelocking

It was recognized early on that the longitudinal Kerr effect or SPM together with negative dispersion results in soliton formation and further reduces pulse duration by about a factor of 2 in dye lasers.^{32,33} However, at that time an analytic solution for the pulse-shortening effect was not presented. Using soliton perturbation theory, an analytic solution has been derived that describes how to make use of much more significant soliton pulse shortening in solid-state lasers.^{51,52} This modelocking model, referred to as *soliton* modelocking, is fundamentally different from any previous models, because it treats soliton pulse shaping as the dominant pulse formation process and the saturable absorber as a perturbation to the soliton. This strongly relaxes the requirements of the saturable absorber as compared to pure saturable absorber modelocking. However, the saturable

absorber is still required to start the modelocking process and to stabilize the soliton pulses against continuous wave breakthrough.

The soliton modelocking model was experimentally confirmed by the production of 300-fs-long soliton pulses with a saturable absorber response time of only 10 ps.⁵³ Pulses of 13 fs were achieved with a saturable absorber response time of about 60 fs.⁵² Further improvements of extremely broad-band saturable absorbers⁴¹ with a higher modulation depth would be necessary to obtain even shorter pulses based on this modelocking mechanism.

In soliton modelocking the pulse duration is given by the soliton condition:

$$\tau_p = 1.76 \frac{|D_2|}{kn_2 L_L F_{p,L}} \quad (19)$$

where D_2 is the total group delay dispersion inside the laser cavity per cavity round-trip. Here we assume that the dominant SPM is produced in the laser material (i.e., n_2 is the nonlinear refractive index of the laser material, L_L is the length of the laser material, and $F_{p,L}$ the pulse fluence inside the laser material). In other cases we have to add all other contributions as well. The pulse duration scales linearly with the negative intracavity dispersion. Reducing the intracavity dispersion results in shorter transform-limited pulses.⁵³ However, there is a limit. The soliton loses energy due to gain dispersion and losses in the cavity. This lost energy, called *continuum* in soliton perturbation theory,⁵⁴ is initially contained in a low-intensity background pulse, which experiences negligible SPM, but spreads in time due to group velocity dispersion. This continuum experiences a higher gain compared to the soliton pulse, because it only sees the gain at line center (while the soliton sees an effectively lower average gain due to its larger bandwidth). After a sufficient build-up time, the continuum would actually grow until it reaches lasing threshold, destabilizing the soliton. However, we can stabilize the soliton by introducing a “slow” saturable absorber into the cavity. This “slow” absorber has to be fast enough to add sufficient additional loss for the growing continuum that spreads in time so that it no longer reaches lasing threshold. For a given recovery time of the saturable absorber, the continuum pulse will not broaden fast enough when the negative dispersion becomes too small. Therefore there is a minimum pulse duration that can be achieved for a given set of absorber and laser parameters:

$$\tau_{p,\min} = 1.7627 \left(\frac{1}{\sqrt{6\pi\Delta\nu_g}} \right)^{3/4} \phi_s^{-1/8} \left(\frac{\tau_A g^{3/2}}{q_0} \right)^{1/4} \quad (20)$$

where ϕ_s is the phase shift of the soliton per cavity round trip (assuming that the dominant SPM occurs in the laser gain medium) and is given by

$$\phi_s = \phi_s(z=2L_L) = kn_2 L_L I_{0,L} \quad (21)$$

Here we assume a fully saturated slow absorber with a linear approximation for the exponential decay of the slow saturable absorber, and a slow saturable absorber.

In the femtosecond regime, we observe a significant reduction of the tendency of Q-switching instabilities compared to pure saturable absorber modelocked picosecond lasers [Eq. (18)]. This can be explained as follows: If the energy of an ultrashort pulse rises slightly due to relaxation oscillations, SPM and/or SAM broadens the pulse spectrum. A broader spectrum, however, reduces the effective gain due to the finite gain bandwidth, which provides some negative feedback, thus decreasing the critical pulse energy which is necessary for stable cw modelocking. The simple stability requirement of Eq. (18) then has to be modified as follows:²³

$$E_{\text{sat},L} g K^2 E_p^2 + E_p^2 > E_{\text{sat},L} E_{\text{sat},A} \Delta R \quad (22)$$

where K is given by

$$K \equiv \frac{0.315}{1.76} \frac{4\pi n_2 L_L}{D_2 A_L \lambda_0 \Delta v_g} \quad (23)$$

Here we assume that the dominant SPM is produced in the laser medium. In other cases we have to add all other contributions as well.

Longitudinal Kerr Effect for Passive Modelocking with a Coupled Cavity: Soliton Laser, Additive Pulse Modelocking

The longitudinal Kerr effect can also be used to produce the same effect as a fast saturable absorber. To do this, the phase nonlinearity provided by the longitudinal Kerr effect has to be converted into an effective amplitude nonlinearity. The earliest modelocking schemes (based only on SPM) used a coupled cavity to convert SPM into SAM. In the soliton laser,⁵⁵ pulses compressed by SPM and anomalous dispersion in the coupled cavity are directly coupled back into the main laser cavity. This provides more gain for the center of the pulse. Pulses as short as 19 fs have been demonstrated with color center lasers.⁵⁶ Later, the SPM-to-SAM conversion with a coupled cavity was demonstrated for a case when the pulses inside the coupled cavity were broadened due to positive group velocity dispersion.⁵⁷ In this case, no compressed pulse was fed back into the main cavity. An effective SAM was obtained because SPM inside the coupled cavity generates a phase modulation on the pulse that adds constructively at the peak of the pulse in the main cavity and destructively in the wings, thus shortening the pulse duration inside the main cavity. This was also referred to as additive pulse modelocking (APM).⁵⁸ Although very powerful in principle, these coupled-cavity schemes have the severe disadvantage that the auxiliary cavity has to be stabilized interferometrically. An alternative method for converting the reactive Kerr nonlinearity into an effective saturable absorber has been developed: Kerr-lens modelocking (KLM).⁵⁹

Transverse Kerr Effect for Passive Modelocking: Kerr Lens Modelocking

The discovery of Kerr lens modelocking has been a breakthrough in ultrashort pulse generation.⁵⁹ Initially the modelocking mechanism was not understood and was somewhat of a mystery. But within a short time after the initial discovery it became clear that the transverse Kerr effect provides a fast saturable absorber. In KLM, the transverse Kerr effect produces a Kerr lens (see Fig. 6) that focuses the high intensity part of the beam more strongly than the low intensity part. Thus, combined with an intracavity aperture the Kerr lens produces less loss for high intensity and forms an effective fast saturable absorber.^{60–62} A similar modelocking effect can be obtained without a hard aperture when the Kerr lens produces an increased overlap of the laser mode with the pump profile in the gain medium.⁶³ The Kerr lens provides the strongest advantage for the pulsed operation when the cavity is operated close to the stability limit. Optimization guidelines for SAM produced by the Kerr lens in different cavities can be found in Ref. 64. Unfortunately, the transverse Kerr effect couples the modelocking process with the laser cavity mode. In contrast, the use of only the longitudinal Kerr effect in modelocking decouples the modelocking process from the laser mode. This allows optimum cavity design for scaling the laser to higher powers and to higher pulse repetition rates without being constrained by the Kerr lens.

KLM is well described by the fast absorber modelocking model discussed above⁶⁵ even though it is not so easy to determine the exact saturable absorber parameters such as the effective saturation fluence. However, the linearized model does not describe the pulse generation with Ti:sapphire lasers in the sub-10-fs regime very well. Pulse-shaping processes in these lasers are more complex.^{66,67} Under the influence of the different linear and nonlinear pulse shaping mechanism, the pulse is significantly broadened and recompressed, giving rise to a “breathing” of the pulse width. The order of

the pulse shaping elements in the laser cavity becomes relevant and the spectrum of the modelocked pulses becomes more complex. In this case, an analytical solution can no longer be obtained. As a rough approximation, the pulses still behave like solitons and consequently these lasers are also called solitary lasers.⁶⁶

Longitudinal Kerr Effect for Passive Modelocking with Nonlinear Polarization Rotation

A Kerr effect-induced nonlinear polarization rotation in a weakly birefringent fiber has been used as a “pulse cleaner” to reduce the low-intensity pulse pedestals.^{68,69} The same effect can also be used to form an effective fast saturable absorber.⁷⁰ Pulses as short as 38 fs have been generated with Nd-doped fiber lasers.

18.6 SEMICONDUCTOR ULTRAFAST NONLINEARITIES: MICROSCOPIC PROCESSES

The discussion of saturable absorbers in Sec. 18.4 has shown that semiconductors are well-suited absorber materials for ultrashort pulse generation. In contrast to saturable absorber mechanisms based on the Kerr effect (Sec. 18.5), ultrafast semiconductor nonlinearities can be studied outside the laser. Such studies give insight into the microscopic processes that determine the nonlinear optical properties of semiconductors on ultrashort time scales. This insight has substantially contributed to our understanding of the physics of semiconductors. Moreover, the information obtained from ultrafast semiconductor spectroscopy provides the basis for further improvement of ultrashort pulse generation with semiconductor saturable absorbers.

In this section, we will first give an overview of ultrafast semiconductor dynamics, which will be followed by a description of the most important experimental techniques for the study of ultrafast processes and nonlinearities. Then we summarize some of the results that have been obtained with these experimental techniques. This summary emphasizes aspects that are particularly relevant for saturable absorber applications, but goes beyond the saturable absorber issue if this is helpful to illustrate general concepts.

In ultrafast semiconductor spectroscopy, it is often convenient to distinguish between excitonic excitations (i.e., Coulomb-bound electron-hole pairs at the band edge⁷¹) and unbound electron-hole pairs in the continuum of the spectrum. Laser pulses with a temporal width well below 100 fs have a spectral bandwidth that is much larger than the spectral width of the exciton resonance and the exciton binding energy in most semiconductors. This is illustrated in Fig. 7 for the example of

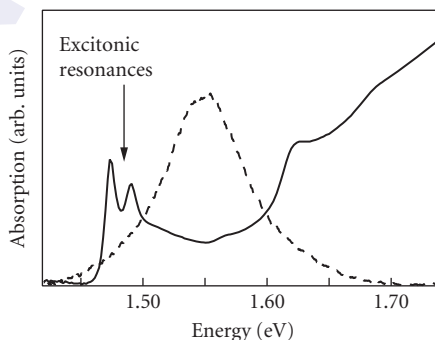


FIGURE 7 Room temperature linear absorption spectrum (solid) of a GaAs/AlGaAs semiconductor quantum well and power spectrum of a 16-fs laser pulse (dashed).

a 16-fs pulse and the absorption spectrum of a GaAs/AlGaAs quantum well⁷² at room temperature. Therefore, saturable absorber applications with sub-100-fs pulses very often involve broadband continuum excitations. For this reason, we will focus on ultrafast continuum nonlinearities and dynamics. Exciton dynamics will be discussed only to outline some general concepts of ultrafast semiconductor spectroscopy. For a comprehensive, in-depth review of ultrafast semiconductor spectroscopy the interested reader is referred to Ref. 73.

Overview

Semiconductors are characterized by closely spaced electronic eigenstates in energy space. This electronic structure gives rise to strong interaction among optical excitations on ultrafast time scales and very complex dynamics. Despite the complexity of the dynamics, different time regimes can be distinguished in the evolution of optical excitations in semiconductors.^{73,74} These different time regimes are schematically illustrated in Fig. 8, which shows the energy dispersion diagram of a 2-band bulk semiconductor. Optical excitation with an ultrafast laser pulse prepares the semiconductor in the coherent regime (time regime I in Fig. 8). In this regime, a well-defined phase relation exists between the optical excitations and the electric field of the laser pulse and among the optical excitations themselves. The coherence among the excitations in the semiconductor gives rise to a macroscopic polarization (dipole moment density). Since the macroscopic polarization enters as a source term in Maxwell's equations, it leads to an electric field which is experimentally accessible. The magnitude and decay of the polarization provide information on the properties of the semiconductor in the coherent regime. The irreversible decay of the polarization is due to scattering processes and is usually described by the so-called dephasing or transversal relaxation time. For a mathematical definition of this time constant the reader is referred to Refs. 73, 75, 76, and 77. Some more details about dephasing and coherent dynamics in semiconductors will be given later.

After the loss of coherence, ultrafast spectroscopy of semiconductors is solely concerned with the dynamics of the population (i.e., electron and hole distributions). In this incoherent regime, the time regimes II–IV can be distinguished, as described in the text that follows. The initial electron and hole distributions are nonthermal in most cases (i.e., they cannot be described by Fermi-Dirac statistics with a well-defined temperature). Scattering among charge carriers is mainly responsible for the redistribution of energy within the carrier distributions and for the formation of thermal

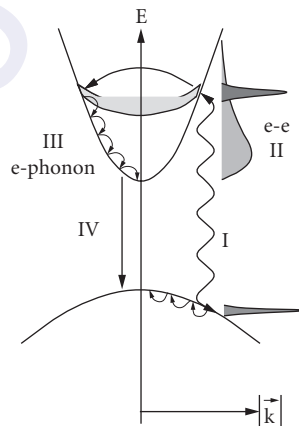


FIGURE 8 Schematic dispersion diagram of a two-band bulk semiconductor showing the time regimes I–IV after optical excitation; see text for more details. e–e: electron–electron scattering, e–phonon: electron–phonon scattering.

distributions. This thermalization is shown as time regime II in Fig. 8, for the example of a thermalizing electron distribution where thermalization occurs through scattering among the electrons. For excitation of the continuum, thermalization usually occurs on a time scale of 100 fs under most experimental conditions. More details about the dynamics in the thermalization regime will be presented later.

In general, the carriers have a temperature different from the lattice temperature after thermalization has been completed. In Fig. 8 it is assumed that the carriers have a higher temperature than the lattice. For this case, Fig. 8 schematically shows the cooling of carriers by the emission of phonons (i.e., energy transfer to the lattice). Cooling defines the time regime III. Typical time constants are in the picosecond and tens of picosecond range.

Finally, the optically excited semiconductor returns to thermodynamic equilibrium by the recombination of electron-hole pairs. Recombination is shown as time regime IV in Fig. 8. In a perfect semiconductor crystal, recombination proceeds via the emission of photons or Auger processes at high carrier densities. These recombination processes take place on time scales of tens of picoseconds and longer. These slow recombination processes as well as the relatively slow carrier cooling will not be discussed in more detail in this chapter. An excellent review can be found in Ref. 73.

Another ultrafast process is encountered if large densities of deep-level traps are incorporated in a semiconductor. Trapping of carriers into deep levels can proceed on subpicosecond time scales (not shown in Fig. 8). Since carrier trapping is important in many saturable absorber applications, it is discussed at the end of this section.

We note that the different time regimes temporally overlap. For example, a scattering process may destroy the coherence and contribute to thermalization. Nevertheless, it is very useful to distinguish between the different time regimes because they are a convenient means for the description of the complex semiconductor dynamics. The schematic picture of the different time regimes also demonstrates that two or more time constants are usually required to describe the temporal response of a semiconductor absorber. For example, we recall that thermalization typically takes place on the 100-fs time scale, while carrier trapping proceeds on times scales from a few hundreds of femtoseconds to picoseconds.

Experimental Techniques

In the following, we describe two very common experimental techniques for the study of ultrafast processes and nonlinearities. The discussion focuses on semiconductors. However, the experimental techniques have also been intensively used for the study of other condensed matter systems or molecules.

Transient Four-Wave Mixing Transient four-wave mixing (FWM) is an experimental technique for the study of the coherent regime. A schematic diagram of a transient FWM experiment is shown in Fig. 9. Two excitation pulses with wave vectors \mathbf{k}_1 and \mathbf{k}_2 excite the sample. In most FWM experiments the two pulses have the same spectrum (degenerate FWM). We assume here that the case of degenerate FWM is realized. An optical delay line is used to introduce a time delay Δt between the pulses where positive Δt refers to pulse \mathbf{k}_1 arriving before pulse \mathbf{k}_2 at the sample.

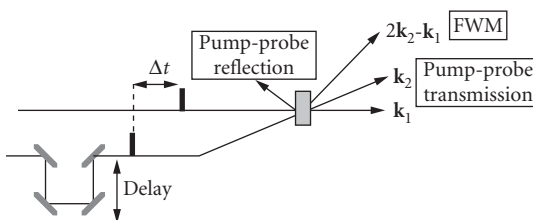


FIGURE 9 Schematic experimental setup for four-wave-mixing (FWM) and pump-probe spectroscopy.

Pulse \mathbf{k}_1 generates a coherent polarization in direction \mathbf{k}_1 in the sample. If the time delay Δt is smaller than the dephasing time of the coherent polarization, in a second step, this polarization interacts with the electric field of pulse \mathbf{k}_2 to set up an interference grating. In a third step, pulse \mathbf{k}_2 is self-diffracted from this grating and an electric field is emitted in the phase-matching direction $2\mathbf{k}_2 - \mathbf{k}_1$. The diffracted field constitutes the FWM signal. A more detailed analysis shows that the FWM signal is due to a nonlinear polarization in direction $2\mathbf{k}_2 - \mathbf{k}_1$.^{73,78}

The FWM emission can be analyzed in different ways. First, the time integral over the FWM intensity can be measured versus the time delay. Such measurements usually provide information about the dephasing time since the strength of the FWM emission is determined by the coherent polarization in direction \mathbf{k}_1 that is left when pulse \mathbf{k}_2 is applied.^{73,78} Spectral information is obtained if the FWM emission is analyzed with a spectrometer at fixed time delays. If more than one resonance is excited, the FWM spectrum shows the magnitude of the optical nonlinearity of the different resonances. Moreover, for homogeneously broadened transitions, the dephasing time can be obtained from the spectral width of the FWM emission. For a fixed time delay, the decay of the FWM signal can be measured in real time by optical gating with a reference pulse. Usually, sum frequency generation in a nonlinear crystal is used for this purpose. This is essentially a cross correlation measurement between the FWM signal pulse and the reference pulse, in which the time resolution is determined by the duration of the reference pulse. In such a correlation measurement, a slow photodetector can be used to detect the sum frequency signal. Real-time detection of FWM signals shows whether the excited optical transitions are homogeneously or inhomogeneously broadened. For homogeneous broadening, the FWM signal in direction $2\mathbf{k}_2 - \mathbf{k}_1$ immediately sets in when pulse \mathbf{k}_2 is applied (for positive time delays). For inhomogeneous broadening, the FWM signal is emitted time-delayed with respect to pulse \mathbf{k}_2 as a so-called photon echo.⁷⁵ More details as well as a mathematical analysis of degenerate two-pulse FWM can be found in Refs. 73 and 78. Transient FWM with three pulses is treated in Ref. 79.

For excitation of interband transitions in semiconductors, occupation of valence and conduction band states contributes to the optical nonlinearity exploited in the FWM process. Occupation effects are a source of nonlinearity in many electronic systems, as shown by the analysis of simple two-level systems.^{73,78} Other sources of nonlinearity will be briefly mentioned when we present some results of coherent semiconductor spectroscopy.

Pump-Probe Spectroscopy Pump-probe spectroscopy is the most widely used technique for the study of ultrafast optical nonlinearities. As shown in Fig. 9, the sample is excited by the pump pulse \mathbf{k}_1 . The nonlinear changes of the transmission or reflectivity of the sample are detected by the time-delayed probe pulse \mathbf{k}_2 . The time delay is defined as positive if the pump pulse precedes the probe. Often, the pump pulse train is amplitude-modulated and a lock-in amplifier is used for detection of the nonlinear transmission or reflectivity changes at the modulation frequency. It is important to note that the pump-probe experiment is a correlation experiment, in which a slow photodetector can be used.

With this technique, nonlinear reflectivity changes are measured in saturable absorbers which are integrated within a mirror structure, as mentioned in Sec. 18.4. In studies of microscopic processes, very often the nonlinear changes of the probe transmission are measured, which is referred to as *differential transmission* (DT) spectroscopy. In the simplest arrangement, the DT signal is spectrally integrated over the spectrum of the probe pulse. If the probe pulse has a large bandwidth, the DT signal at various photon energies can be obtained from measurements of the spectrum of the transmitted probe pulse in the presence and absence of the pump. The difference between those spectra represents the DT spectrum. Spectral information can also be obtained if a tunable narrowband probe pulse is scanned over a spectral window.

The differential transmission signal contains information on the coherent and the incoherent regime. In the coherent regime, the DT signal is determined by the nonlinear polarization in the direction of the probe pulse. The DT signal in the coherent regime is analyzed in more detail in Refs. 80, 81, 82, and 83.

Coherent effects can be neglected on time scales much longer than the dephasing time of the polarization. The DT signal is then determined by the population (i.e., the electron and hole distributions) generated by the pump pulse. In the incoherent regime, one can roughly distinguish between

two different sorts of effects that determine the DT signal of semiconductors: (1) occupation effects and (2) many-body Coulomb effects. Occupation effects are based on the fermionic nature of optical excitations and can be understood in the following way. A nonequilibrium carrier population in excited states reduces the optical transition rate into these states due to the reduced density of empty final states. This manifests itself by the reduction of the absorption at certain photon energies. Occupation effects are important in many electronic systems. Many-body effects are particularly important in semiconductors. They include the renormalization of the bandgap and screening of the Coulomb interaction. Details can be found in Refs. 84 through 87. Of course, the many-body effects are related to the occupation effects since they depend on the distribution of carriers.

If many-body effects can be neglected and if the optical matrix elements do not depend on the pump excitation, the nonlinear change $\Delta\alpha$ of the absorption constant α is directly proportional to the magnitude of the population (i.e., the electron and hole densities). In general, both refractive index changes and absorption changes $\Delta\alpha$ contribute to the differential transmission signal.⁸⁸ A considerable simplification is obtained if refractive index changes can be neglected and if the relation $\Delta\alpha d \ll 1$ holds (d sample thickness). Then the DT signal is proportional to $\Delta\alpha$, and the decay of the DT signal can be identified with the decay of the electron and hole densities. We note that in studies of carrier trapping, pump-probe data are very often analyzed under the just-described assumptions.

Results

In this subsection, examples are given that show how ultrafast spectroscopy has contributed to the understanding of ultrafast processes and nonlinearities in semiconductors. We will highlight those aspects that are particularly relevant for broadband saturable absorbers. The material is organized according to the different time regimes which have been identified at the beginning of Sec. 18.6.

Coherent Regime: Excitonic Excitations Studies of excitons in the coherent regime are an illustrative example for the issues addressed in coherent ultrafast spectroscopy of semiconductors. Therefore, we briefly discuss some aspects of coherent exciton dynamics even though excitonic excitations are less important for many ultrafast saturable absorber applications. This discussion will illustrate the main issues addressed in ultrafast semiconductor spectroscopy in the coherent regime: (1) dephasing times and underlying scattering mechanisms, (2) coupling and interference between optical excitations, and (3) the nature of the optical nonlinearity.

Four-wave mixing has been intensively used to study the dephasing of excitons and the underlying scattering mechanisms. Varying the excitation intensity, the effects of exciton-exciton and exciton-electron scattering on the dephasing of excitons have been investigated in bulk semiconductors⁸⁹ and quantum wells.⁹⁰ Varying the temperature, dephasing of excitons due to interaction with lattice vibrations (phonons) has been studied in Ref. 91. Summarizing these results, exciton dephasing times can be in the picosecond range at moderate exciton and free-carrier densities and Helium temperatures. Increasing the temperature decreases the dephasing time due to enhanced phonon scattering. Likewise, enhanced carrier or exciton densities cause faster dephasing. A detailed analysis of exciton dephasing in three- and two-dimensional semiconductors is presented in Ref. 92. We note that exciton dephasing can also give insight into the structural properties of a semiconductor, such as interface roughness in quantum wells and alloy disorder in mixed crystals.^{93,94}

More recent four-wave mixing investigations of coherent exciton dynamics focus on quantum-mechanical coupling and interference between different exciton transitions. Beating phenomena between two exciton transitions have been observed in various systems^{95–98} as well as beats involving the whole excitonic Rydberg series.⁹⁹ The study of beating phenomena allows for the extraction of level splittings and yields information on the quantum-mechanical coupling in multilevel systems (i.e., their internal structure^{100,101}).

The nature of the optical nonlinearity in the coherent regime is another subject that has been intensively studied in four-wave-mixing experiments on excitons.^{87,102–110} It has been shown that many-body Coulomb interaction substantially contributes to the coherent optical nonlinearity. With respect to the essence of this result we note that the polarization of a certain interband excitation

gives rise to an electric field which has to be added to the external laser field in order to determine the coherent dynamics of other excitations. A rigorous theoretical treatment can be found in Refs. 85 and 87.

Coherent Regime: Continuum Excitations Four-wave-mixing studies of continuum dephasing have been performed with broadband sub-10-fs pulses in three-dimensional bulk semiconductors¹¹¹ and quasi two-dimensional quantum wells.¹¹² These studies have shown that the decay of the coherent polarization of the semiconductor continuum is extremely fast. Decay times are only about 10 fs at high carrier densities, which are likely to be obtained in a saturable absorber in a laser cavity. The ultrafast dephasing is mainly due to carrier-carrier scattering with a density dependence that reflects the dimensionality of the semiconductor.¹¹² More recent work¹¹³ shows that at reduced carrier densities interaction with the lattice also needs to be considered in continuum dephasing experiments. The interaction with the lattice has been identified as electron-LO-phonon scattering,¹¹³ which has a time constant of about 200 fs.^{114–116}

The internal structure of the continuum has been experimentally investigated in Ref. 117. The observation of a photon echo has demonstrated that the semiconductor continuum can be treated as an ensemble of uncoupled excitations at higher carrier densities.¹¹⁷ Interaction between continuum and exciton transitions in coherent nonlinear optics has been studied both for degenerate^{118–120} and nondegenerate excitons and continua.^{121–124} This work has demonstrated the importance of many-body coupling effects between different interband transitions at lower carrier densities.

We note that it is not yet clear how coherence in a semiconductor saturable absorber affects the formation of broadband ultrashort pulses in a laser cavity. Some theoretical predictions can be found in Refs. 125 and 126. Experimental results about this issue are missing so far. Given the ultrafast dephasing times in semiconductors under the conditions in a laser cavity, coherence effects are most likely to be important for the generation of sub-10-fs pulses.

Thermalization Regime Here we will focus on the excitation of semiconductor continuum states. Studies of the thermalization of free electron and hole distributions are an instructive example for the usefulness of the differential transmission (DT) technique. In particular, measurements of differential transmission spectra have yielded considerable insight into the complex processes that determine the dynamics in the thermalization regime. For excitation well above the bandgap, nonthermal carrier distributions can be observed in bulk semiconductors, such as GaAs,¹²⁷ and quantum wells.¹²⁸ These nonthermal distributions manifest themselves as positive signal in the DT spectrum with a shape that is approximately given by the spectrum of the pump pulse. Often, this signature in the DT spectrum is referred to as a *spectral hole*. The decay of the spectral hole and the thermalization of the carrier distributions lead to a substantial change of the shape of the DT spectrum. Thermalization occurs on the 100-fs time scale in undoped semiconductors and is determined by carrier-carrier scattering in many experiments. The exact thermalization time strongly depends on the carrier density, the excess photon energy with respect to the band edge, and the type of carrier.^{73,127–129} Thermalization of optically excited carriers in the presence of cold electron or hole plasmas has also been investigated by the differential transmission technique.^{73,130} These experiments have been performed in modulation-doped quantum wells⁷² and show that thermalization can occur in less than 10 fs.

Besides carrier-carrier scattering, intervalley scattering is another process which can affect carrier dynamics in the thermalization regime. Semiconductors such as GaAs possess several conduction band minima at different points of the Brillouin zone. If electrons are created near the center of the Brillouin zone with large enough excess energy, they can scatter to the side valleys. This process has been investigated in DT experiments, and intervalley scattering times in the sub-100-fs range have been deduced.^{131,132}

Spectrally resolved DT measurements have also revealed interesting many-body effects. A closer inspection of DT spectra has shown that the spectral hole is redshifted with respect to the pump spectrum.¹³³ Moreover, at the high-energy edge of the pump spectrum a negative DT signal is observed. These signatures have been interpreted in terms of Fermi edge singularities at the upper and lower edge of the nonthermal electron distribution generated by the pump pulse.^{133,134} The

work on Fermi edge singularities is an instructive example for the interplay between occupation and many-body effects in differential transmission experiments in semiconductors.

Thermalization of carrier distributions and many-body effects also change the spectrally integrated DT signal or the impulse response of a semiconductor saturable absorber. In particular, thermalization contributes to the fast decay seen in the impulse response at early times (cp. Fig. 3). It is important to note that this fast decay is the result of the complex redistribution of nonlinear transmission or reflectivity changes in frequency space. As a consequence, no general concept has evolved so far for the engineering of this fast decay.

The fast decay of a spectrally integrated DT curve in the thermalization regime also depends on the temporal structure of the pump and probe pulses themselves.^{135,136} This point has been recently demonstrated in experiments in which a so-called frequency chirp was imposed on 20-fs pulses. In a frequency-chirped laser pulse the frequency varies over the temporal profile of the pulse. Such chirped pulses have been used in DT experiments on continuum transitions in bulk semiconductors. The results show that the fast decay of spectrally integrated DT curves can be enhanced by an appropriate frequency chirp.^{135,136} The manipulation of ultrafast nonlinearities by the chirp of the pulses can be viewed as an example of a much more general concept known as *coherent control*.¹³⁷ In coherent control experiments, the temporal shape or spectral content of ultrafast laser pulses are adjusted¹³⁸ to reach preset goals. The chirp control experiments in Refs. 135 and 136 demonstrate that the detailed temporal and spectral structure of the laser pulses should be included in the optimization of semiconductor saturable absorbers. That probably also means that the position of the SESAM within the laser cavity may play a role in sub-10-fs pulse generation.

Carrier Trapping Processes that remove electrons and holes from the bands of a semiconductor lead to a decay of the nonlinear transmission or to reflectivity changes resulting from the interband transitions. As discussed in Sec. 18.4, semiconductor saturable absorber applications in ultrashort pulse generation often require picosecond or subpicosecond absorber recovery times. The simplest way to obtain such short absorber recovery times would be to remove the optically excited carriers from the bands a few hundreds of femtoseconds after they have been created. Ultrafast depletion of the semiconductor band states is also important in all-optical switching devices^{139,140} and optoelectronics.¹⁴¹ However, intrinsic recombination processes are usually too slow to deplete the band states of a semiconductor on picosecond or subpicosecond time scales. Therefore, one generates defect states in the bandgap which give rise to fast carrier trapping, thereby depleting the bands. The trapping time is determined by the density and the type of the traps. Higher trap densities give rise to faster trapping.

Standard methods for the controlled incorporation of defect and trap states are ion implantation¹⁴² and low-temperature (LT) molecular beam epitaxy.¹⁴³ In ion-implanted semiconductors, the trap density and the type of defect are determined by the implantation dose. The growth temperature controls the defect density in LT semiconductors, where larger defect densities are incorporated at lower temperatures.¹⁴⁴ GaAs is the best understood LT-grown III-V semiconductor. Low-temperature growth of GaAs is performed at temperatures of 200 to 300°C, as compared to about 600°C in standard molecular beam epitaxy. During LT growth of GaAs, excess arsenic is incorporated in the form of arsenic antisites (As on Ga lattice site: As_{Ga}) at densities as large as 10^{20} cm^{-3} .^{144,145} In undoped LT GaAs, more than 90 percent of the antisites are neutral, while the rest is singly ionized due to presence of Ga vacancies (V_{Ga}) which are the native acceptors in the material (see Fig. 10a).^{144,146} The ionized arsenic antisites have been identified as electron traps.¹⁴⁷ Annealing at higher temperatures (typically 600°C and higher) converts the arsenic antisite point defects into arsenic clusters, so-called As precipitates (see Fig. 10b).¹⁴⁸ A detailed review of the properties of LT GaAs can be found in Refs. 149 and 150.

The carrier trapping times in as-grown LT GaAs can be in the subpicosecond regime and show the expected decrease with decreasing growth temperature.^{151,152} Subpicosecond recovery times of nonlinear transmission or reflectivity changes are also found in annealed LT GaAs, indicating that arsenic precipitates efficiently deplete the band states.^{153,154} For more details about carrier trapping in LT semiconductors the reader is referred to Refs. 155 through 162.

Picosecond and subpicosecond carrier trapping times have also been found in semiconductors implanted with various ion species.^{163–168} A decrease of the trapping time with increasing ion dose

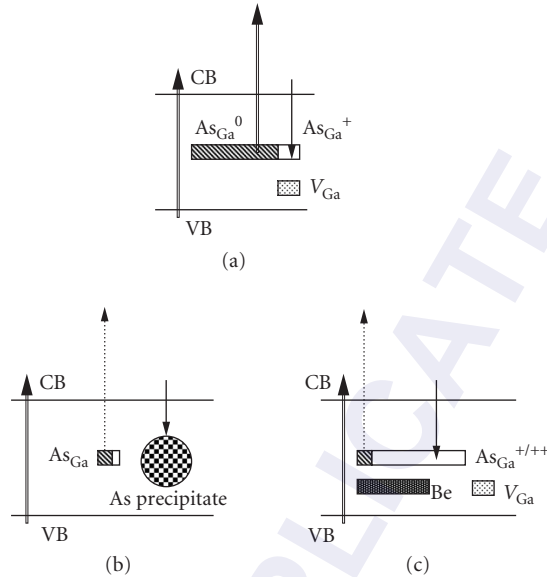


FIGURE 10 Electronic structure of undoped as-grown (a), undoped annealed (b), and Beryllium doped as-grown LT GaAs (c). The double arrows mark strong optical absorption transitions. Weak optical absorption transitions are indicated by dotted arrows. Trapping processes are shown by arrows that point downward.

was observed at lower doses.^{163,164,167} At higher ion doses, the trapping time can increase with the dose.¹⁶⁸ The correlation of trapping times with structural properties of ion implanted semiconductors has given more insight into this unexpected dose dependence of the trapping time.¹⁶⁹ This work indicates that not only the defect density but also the type of defect depends on the ion dose.¹⁶⁹ Both the density and the type of defect affect carrier trapping, leading to longer trapping times if less effective traps are generated at higher ion doses.¹⁶⁹

Besides an ultrafast carrier trapping and absorber recovery time, other important saturable absorber parameters are the modulation depth and the nonsaturable losses which remain even at the highest pump energy fluences (see Sec. 18.4). Optimized materials combine an ultrafast recovery time with high modulation and small nonsaturable losses. This material optimization issue has been addressed in recent publications.^{154,168,170} In these studies, the nonlinearity of continuum transitions was investigated in different modifications of GaAs. The preparation of the semiconductor layers ensured that the modulation depth and the nonsaturable losses were determined by nonlinear absorption changes.

It has been shown that standard as-grown LT GaAs with an ultrafast carrier trapping time suffers from a small absorption modulation and high nonsaturable absorption losses.^{154,170} Note that large nonsaturable absorption decreases the modulation depth and causes large nonsaturable losses when the semiconductor absorber is integrated within a mirror structure. The high nonsaturable absorption mainly results from the strong defect absorption from the neutral As antisites to the conduction band (see Fig. 10a) whose saturation fluence has been shown to be extremely high.¹⁵⁴ Therefore, the goals of material optimization are (1) to reduce the nonsaturable absorption by the reduction of the density of neutral As antisites and (2) to maintain a fast trapping and absorber recovery time.

We have demonstrated two different ways to reach those goals. Annealing of LT GaAs strongly reduces the density of neutral As antisites and the nonsaturable absorption.¹⁵⁴ The simultaneous reduction of the density of useful ionized As antisite electron traps does not substantially increase the

absorber recovery time due to the presence of the As precipitates (see Fig. 10b). Alternatively, doping with acceptors, such as Beryllium, can be used to reduce the density of neutral As antisites.¹⁷¹ The simultaneous increase of the ionized As antisite density results in ultrafast carrier trapping and absorber recovery times (see Fig. 10c).¹⁷⁰ Annealed LT GaAs and Be doped LT GaAs combine ultrafast recovery times with high modulation depth and small nonsaturable losses. These materials are well suited for saturable absorber devices in laser physics and for all-optical switching applications.

Studies of the modulation depth ΔR , the nonsaturable losses ΔR_{ns} , and the recovery time τ_A in ion implanted GaAs have shown that ΔR decreases and ΔR_{ns} increases with decreasing recovery time.¹⁶⁸ Nevertheless, if the ion species, ion dose, and annealing conditions are properly chosen, combinations of ΔR , ΔR_{ns} , and τ_A can be obtained which are appropriate for saturable absorber applications. Ion-implanted GaAs is an alternative to annealed or Be doped LT GaAs as a material for saturable absorber devices.¹⁶⁸

18.7 REFERENCES

1. P. F. Moulton, "Spectroscopic and Laser Characteristics of Ti:Al₂O₃," *J. Opt. Soc. Am. B.* **3**:125–132 (1986).
2. D. H. Sutter, G. Steinmeyer, L. Gallmann, N. Matuschek, F. Morier-Genoud, U. Keller, V. Scheuer, G. Angelow, and T. Tschudi, "Semiconductor Saturable-Absorber Mirror-assisted Kerr-Lens Mode-Locked Ti:Sapphire Laser Producing Pulses in the Two-Cycle Regime," *Optics Lett.* **24**:631–633 (1999).
3. U. Morgner, F. X. Kärtner, S. H. Cho, Y. Chen, H. A. Haus, J. G. Fujimoto, E. P. Ippen, V. Scheuer, G. Angelow, and T. Tschudi, "Sub-Two-Cycle Pulses from a Kerr-Lens Mode-Locked Ti:Sapphire Laser: Addenda," *Optics Lett.* **24**:920 (1999).
4. U. Keller, "Semiconductor Nonlinearities for Solid-State Laser Modelocking and Q-Switching," in Chapter 4 in *Nonlinear Optics in Semiconductors*, E. Garmire and A. Kost, (eds.), Academic Press, Inc., Boston, 1999, vol. 59, pp. 211–286.
5. U. Keller, "Ultrafast Solid-State Lasers," *Progress in Optics*, Elsevier, 2004, vol. 46, pp. 1–115.
6. L. Krainer, R. Paschotta, S. Lecomte, M. Moser, K. J. Weingarten, and U. Keller, "Compact Nd:YVO₄ Lasers with Pulse Repetition Rates up to 160 GHz," *IEEE J. Quantum Electron.* **38**:1331–1338 (2002).
7. A. E. H. Oehler, T. Südmeyer, K. J. Weingarten, and U. Keller, "100-GHz Passively Modelocked Er:Yb:Glass Laser at 1.5 μ m with 1.6-ps Pulses" *Optics Express* **16**(26): 21930–21935 (2008).
8. R. Paschotta, J. Aus der Au, G. J. Spühler, F. Morier-Genoud, R. Hövel, M. Moser, S. Erhard, M. Karszewski, A. Giesen, and U. Keller, "Diode-Pumped Passively Modelocked Lasers with High Average Power," *Appl. Phys. B.* **70**:S25–S32 (2000).
9. G. J. Spühler, R. Paschotta, R. Fluck, B. Braun, M. Moser, G. Zhang, E. Gini, and U. Keller, "Experimentally Confirmed Design Guidelines for Passively Q-Switched Microchip Lasers Using Semiconductor Saturable Absorbers," *J. Opt. Soc. Am. B.* **16**:376–388 (1999).
10. U. Keller, *Ultrafast Solid-State Lasers*, Landolt-Börnstein, Group VIII/1B1, Laser Physics and Applications. Subvolume B: Laser Systems. Part 1. G. Herziger, H. Weber, R. Proprawe (eds.), Springer-Verlag, Berlin, Heidelberg, New York, pp. 33–167, October, 2007.
11. T. Südmeyer, S. V. Marchese, S. Hashimoto, C. R. E. Baer, G. Gingras, B. Witzel, and U. Keller, "Femtosecond Laser Oscillators for High-Field Science," *Nature Photonics* **2**:599–604 (2008).
12. U. Keller and A. C. Tropper, "Passively Modelocked Surface-Emitting Semiconductor Lasers," *Physics Reports* **429**(2):67–120 (2006).
13. A. R. Bellancourt, D. J. H. C. Maas, B. Rudin, M. Golling, T. Südmeyer, and U. Keller, "Modelocked Integrated External-Cavity Surface Emitting Laser (MIXSEL)," *IET Optoelectronics* **3**:61–72 (2009).
14. A. Shirakawa, I. Sakane, M. Takasaka, and T. Kobayashi, "Sub-5-fs Visible Pulse Generation by Pulse-Front-Matched Noncollinear Optical Parametric Amplification," *Appl. Phys. Lett.* **74**:2268–2270 (1999).
15. A. Baltuska, Z. Wei, M. S. Pshenichnikov, D. A. Wiersma, and R. Szipöcs, "All-Solid-State Cavity Dumped Sub-5-fs laser," *Appl. Phys. B.* **65**:175–188 (1997).
16. M. Nisoli, S. Stagira, S. D. Silvestri, O. Svelto, S. Sartania, Z. Cheng, M. Lenzner, C. Spielmann, and F. Krausz, "A Novel High-Energy Pulse Compression System: Generation of Multigigawatt Sub-5-fs Pulses," *Appl. Phys. B.* **65**:189–196 (1997).

17. G. Steinmeyer, D. H. Sutter, L. Gallmann, N. Matuschek, and U. Keller, "New Frontiers in Ultra-short Pulse Generation: Pushing the Limits in Linear and Nonlinear Optics," *Science* **286**:1507–1512 (1999).
18. U. Keller, D. A. B. Miller, G. D. Boyd, T. H. Chiu, I. F. Ferguson, and M. T. Asom, "Solid-State Low-Loss Intracavity Saturable Absorber for Nd:YLF Lasers: An Antiresonant Semiconductor Fabry-Perot Saturable Absorber," *Optics Lett.* **17**:505–507 (1992).
19. K. J. Weingarten, U. Keller, T. H. Chiu, and J. F. Ferguson, "Passively Mode-Locked Diode-Pumped Solid-State Lasers Using an Antiresonant Fabry-Perot Saturable Absorber," *Optics Lett.* **18**:640–642 (1993).
20. X. M. Zhao and J.-C. Diels, "Ultrashort Laser Sources," in *Handbook of Optics*, M. Bass, E. W. Stryland, D. R. Williams, and W. L. Wolfe, (eds.) McGraw-Hill, Inc., New York, 1995, vol. 1, pp. 14.1–14.29.
21. H. A. Haus, "Short Pulse Generation," in *Compact Sources of Ultrashort Pulses I*. I. N. Duling (ed.) Cambridge University Press, New York, 1995, pp. 1–56.
22. G. P. Agrawal and N. A. Olsson, "Self-Phase Modulation and Spectral Broadening of Optical Pulses in Semiconductor Laser Amplifiers," *IEEE J. Quantum Electron.* **25**:2297–2306 (1989).
23. C. Hönninger, R. Paschotta, F. Morier-Genoud, M. Moser, and U. Keller, "Q-Switching Stability Limits of Continuous-Wave Passive Mode Locking," *J. Opt. Soc. Am. B.* **16**:46–56 (1999).
24. M. Haiml, R. Grange, and U. Keller "Optical Characterization of Semiconductor Saturable Absorbers," *Appl. Phys. B.* **79**:331–339 (2004).
25. D. J. H. Maas, B. Rudin, A.-R. Bellancourt, D. Iwaniuk, S. V. Marchese, T. Südmeier, and U. Keller, "High Precision Optical Characterization of Semiconductor Saturable Absorber Mirrors," *Optics Express* **16**:7571–7579 (2008).
26. G. H. C. New, "Modelocking of Quasi-Continuous Lasers," *Opt. Commun.* **6**:188–192 (1972).
27. G. H. C. New, "Pulse Evolution in Mode-Locked Quasi-Continuous Lasers," *IEEE J. Quantum Electron.* **10**:115–124 (1974).
28. H. A. Haus, "Theory of Mode Locking with a Slow Saturable Absorber," *IEEE J. Quantum Electron.* **11**:736–746 (1975).
29. U. Keller, "Recent Developments in Compact Ultrafast Lasers," *Nature* **424**:831–838 (2003).
30. J. A. Valdmanis, R. L. Fork, and J. P. Gordon, "Generation of Optical Pulses as Short as 27 fs Directly from a Laser Balancing Self-Phase Modulation, Group-Velocity Dispersion, Saturable Absorption, and Saturable Gain," *Optics Lett.* **10**:131–133 (1985).
31. J. A. Valdmanis and R. L. Fork, "Design Considerations for a Femtosecond Pulse Laser Balancing Self-Phase Modulation, Group Velocity Dispersion, Saturable Absorption, and Saturable Gain," *IEEE J. Quantum Electron.* **22**:112–118 (1986).
32. O. E. Martinez, R. L. Fork, and J. P. Gordon, "Theory of Passively Modelocked Lasers Including Self-Phase Modulation and Group-Velocity Dispersion," *Optics Lett.* **9**:156–158 (1984).
33. O. E. Martinez, R. L. Fork, and J. P. Gordon, "Theory of Passively Modelocked Lasers for the Case of a Nonlinear Complex Propagation Coefficient," *J. Opt. Soc. Am. B.* **2**:753 (1985).
34. H. A. Haus, "Theory of Modelocking with a Fast Saturable Absorber," *J. Appl. Phys.* **46**:3049–3058 (1975).
35. K. N. Islam, E. R. Sunderman, C. E. Soccolich, I. Bar-Joseph, N. Sauer, T. Y. Chang, and B. I. Miller, "Color Center Lasers Passively Mode Locked by Quantum Wells," *IEEE J. Quantum Electron.* **25**:2454–2463 (1989).
36. U. Keller, K. J. Weingarten, F. X. Kärtner, D. Kopf, B. Braun, I. D. Jung, R. Fluck, C. Hönninger, N. Matuschek, and J. Aus der Au, "Semiconductor Saturable Absorber Mirrors (SESAMs) for Femtosecond to Nanosecond Pulse Generation in Solid-State Lasers," *IEEE J. Sel. Top. Quantum Electron.* **2**:435–453 (1996).
37. B. G. Kim, E. Garmire, S. G. Hummel, and P. D. Dapkus, "Nonlinear Bragg Reflector Based on Saturable Absorption," *Appl. Phys. Lett.* **54**:1095–1097 (1989).
38. L. R. Brovelli, I. D. Jung, D. Kopf, M. Kamp, M. Moser, F. X. Kärtner, and U. Keller, "Self-Starting Soliton Modelocked Ti:sapphire Laser Using a Thin Semiconductor Saturable Absorber," *Electronics Lett.* **31**:287–289 (1995).
39. S. Tsuda, W. H. Knox, E. A. D. Souza, W. Y. Jan, and J. E. Cunningham, "Low-Loss Intracavity AlAs/AlGaAs Saturable Bragg Reflector for Femtosecond Mode Locking in Solid-State Lasers," *Optics Lett.* **20**:1406–1408 (1995).
40. L. R. Brovelli, U. Keller, and T. H. Chiu, "Design and Operation of Antiresonant Fabry-Perot Saturable Semiconductor Absorbers for Mode-Locked Solid-State Lasers," *J. Opt. Soc. Am. B.* **12**:311–322 (1995).

41. I. D. Jung, F. X. Kärtner, N. Matuschek, D. H. Sutter, F. Morier-Genoud, Z. Shi, V. Scheuer, M. Tilsch, T. Tschudi, and U. Keller, "Semiconductor Saturable Absorber Mirrors Supporting Sub-10 fs Pulses," *Applied Physics B: Special Issue on Ultrashort Pulse Generation* **65**:137–150 (1997).
42. G. J. Spühler, K. J. Weingarten, R. Grange, L. Krainer, M. Haiml, V. Liverini, M. Golling, S. Schön, and U. Keller, "Semiconductor Saturable Absorber Mirror Structure with Low Saturation Fluence," *Appl. Phys. B* **81**(1): 27–32 (2005).
43. I. P. Bilinsky, J. G. Fujimoto, J. N. Walpole, and L. J. Missaggia, "InAs-Doped Silica Films for Saturable Absorber Applications," *Appl. Phys. Lett.* **74**:2411–2413 (1999).
44. N. Sarukura, Y. Ishida, T. Yanagawa, and H. Nakano, "All Solid-State cw Passively Modelocked Ti:Sapphire Laser Using a Colored Glass Filter," *Appl. Phys. Lett.* **57**:229–230 (1990).
45. L. F. Mollenauer, R. H. Stolen, and J. P. Gordon, "Experimental Observation of Picosecond Pulse Narrowing and Solitons in Optical Fibers," *Phys. Rev. Lett.* **45**:1095–1098 (1980).
46. W. J. Tomlinson, R. H. Stolen, and C. V. Shank, "Compression of Optical Pulses Chirped by Self-Phase Modulation in Fibers," *J. Opt. Soc. Am. B* **1**:139–149 (1984).
47. R. L. Fork, C. H. B. Cruz, P. C. Becker, and C. V. Shank, "Compression of Optical Pulses to Six Femtoseconds by Using Cubic Phase Compensation," *Optics Lett.* **12**:483–485 (1987).
48. S. Sartania, Z. Cheng, M. Lenzner, G. Tempea, C. Spielmann, F. Krausz, and K. Ferencz, "Generation of 0.1-TW 5-fs Optical Pulses at a 1-kHz Repetition Rate," *Optics Lett.* **22**:1562–1564 (1997).
49. R. L. Fork, C. V. Shank, C. Hirlimann, R. Yen, and W. J. Tomlinson, "Femtosecond White-Light Continuum Pulses," *Optics Lett.* **8**:1–3 (1983).
50. G. M. Gale, M. Cavallari, T. J. Driscoll, and F. Hache, "Sub-20 fs Tunable Pulses in the Visible from an 82 MHz Optical Parametric Oscillator," *Optics Lett.* **20**:1562–1564 (1995).
51. F. X. Kärtner and U. Keller, "Stabilization of Soliton-Like Pulses with a Slow Saturable Absorber," *Optics Lett.* **20**:16–18 (1995).
52. F. X. Kärtner, I. D. Jung, and U. Keller, "Soliton Modelocking with Saturable Absorbers," Special Issue on Ultrafast Electronics, Photonics and Optoelectronics, *IEEE J. Sel. Topics in Quantum Electronics (JSTQE)* **2**:540–556 (1996).
53. I. D. Jung, F. X. Kärtner, L. R. Brovelli, M. Kamp, and U. Keller, "Experimental Verification of Soliton Modelocking Using Only a Slow Saturable Absorber," *Optics Lett.* **20**:1892–1894 (1995).
54. D. J. Kaup, "Perturbation Theory for Solitons in Optical Fibers," *Phys. Rev. A* **42**:5689–5694 (1990).
55. L. F. Mollenauer and R. H. Stolen, "The Soliton Laser," *Optics Lett.* **9**:13–15 (1984).
56. F. M. Mitschke and L. F. Mollenauer, "Ultrashort Pulses from the Soliton Laser," *Optics Lett.* **12**:407–409 (1987).
57. P. N. Kean, X. Zhu, D. W. Crust, R. S. Grant, N. Landford, and W. Sibbett, "Enhanced Modelocking of Color Center Lasers," *Optics Lett.* **14**:39–41 (1989).
58. E. P. Ippen, H. A. Haus, and L. Y. Liu, "Additive Pulse Modelocking," *J. Opt. Soc. Am. B* **6**:1736–1745 (1989).
59. D. E. Spence, P. N. Kean, and W. Sibbett, "60-fsec Pulse Generation from a Self-Mode-Locked Ti:Sapphire Laser," *Optics Lett.* **16**:42–44 (1991).
60. U. Keller, G. W. 'tHooft, W. H. Knox, and J. E. Cunningham, "Femtosecond Pulses from a Continuously Self-Starting Passively Mode-Locked Ti:Sapphire Laser," *Optics Lett.* **16**:1022–1024 (1991).
61. F. Salin, J. Squier, and M. Piché, "Modelocking of Ti:Sapphire Lasers and Self-Focusing: A Gaussian Approximation," *Optics Lett.* **16**:1674–1676 (1991).
62. D. K. Negus, L. Spinelli, N. Goldblatt, and G. Feugnet, "Sub-100 Femtosecond Pulse Generation by Kerr Lens Modelocking in Ti:Sapphire," in *Advanced Solid-State Lasers*, G. Dubé and L. Chase (eds.) Optical Society of America, Washington, D.C., 1991, vol. 10, pp. 120–124.
63. M. Piché and F. Salin, "Self-Mode Locking of Solid-State Lasers without Apertures," *Optics Lett.* **18**:1041–1043 (1993).
64. V. Magni, G. Cerullo, S. D. Silvestri, and A. Monguzzi, "Astigmatism in Gaussian-Beam Self-Focusing and in Resonators for Kerr-Lens Mode-Locking," *JOSA B* **12**:476–485 (1995).
65. H. A. Haus, J. G. Fujimoto, and E. P. Ippen, "Analytic Theory of Additive Pulse and Kerr Lens Mode Locking," *IEEE J. Quantum Electron.* **28**:2086–2096 (1992).
66. T. Brabec, C. Spielmann, and F. Krausz, "Mode Locking in Solitary Lasers," *Optics Lett.* **16**:1961–1963 (1991).

67. F. Krausz, M. E. Fermann, T. Brabec, P. F. Curley, M. Hofer, M. H. Ober, C. Spielmann, E. Wintner, and A. J. Schmidt, "Femtosecond Solid-State Lasers," *IEEE J. Quantum Electron.* **28**:2097–2122 (1992).
68. J.-L. Tapié, and G. Mourou, "Shaping of Clean, Femtosecond Pulses at 1.053 μm for Chirped-Pulse Amplification," *Optics Lett.* **17**:136–138 (1992).
69. Y. Beaudoin, C. Y. Chien, J. S. Coe, J. L. Tapié, and G. Mourou, "Ultra-high-Contrast Ti:Sapphire/Nd:Glass Terawatt Laser System," *Optics Lett.* **17**:865–867 (1992).
70. M. Hofer, M. H. Ober, F. Haberl, and M. E. Fermann, "Characterization of Ultrashort Pulse Formation in Passively Mode-Locked Fiber Lasers," *IEEE J. Quantum Electron.* **28**:720–728 (1992).
71. R. J. Elliot, in *Polarons and Excitons*, C. G. Kuper and G. D. Whitefield (eds.), Plenum, New York, 1963.
72. C. Weisbuch and B. Vinter, *Quantum Semiconductor Structures: Fundamentals and Applications*, Academic, San Diego, 1991.
73. J. Shah, *Ultrafast Spectroscopy of Semiconductors and Semiconductor Nanostructures*, Springer-Verlag, Berlin, 1996.
74. E. O. Göbel, "Ultrafast Spectroscopy of Semiconductors," in *Advances in Solid State Physics*, U. Rössler (ed.) Friedrich Vieweg & Sohn, Braunschweig/Wiesbaden, 1990, vol. 30, pp. 269–294.
75. L. Allen and J. H. Eberly, *Optical Resonance and Two-Level Atoms*, Dover, New York, 1975.
76. Y. R. Shen, *The Principles of Nonlinear Optics*, John Wiley & Sons, New York, 1984.
77. K. Shimoda, *Introduction to Laser Physics* (Springer Series in Optical Science) Springer, Heidelberg, 1984.
78. T. Yajima, and Y. Taira, "Spatial Optical Parametric Coupling of Picosecond Light Pulses and Transverse Relaxation Effects in Resonant Media," *Journal of the Physical Society of Japan* **47**:1620–1626 (1979).
79. A. M. Weiner, S. De Silvestri, and E. P. Ippen, "Three-Pulse Scattering for Femtosecond Dephasing Studies: Theory and Experiment," *J. Opt. Soc. Am. B.* **2**:654–661 (1985).
80. C. H. Brito Cruz, J. P. Gordon, P. C. Becker, R. L. Fork, and C. V. Shank, "Dynamics of Spectral Hole Burning," *IEEE J. Quantum Electron.* **24**:261–266 (1988).
81. B. Fluegel, N. Peyghambarian, G. Olbright, M. Lindberg, S. W. Koch, M. Joffe, D. Hulin, A. Migus, and A. Antonetti, "Femtosecond Studies of Coherent Transients in Semiconductors," *Phys. Rev. Lett.* **59**:2588–2591 (1987).
82. M. Lindberg and S. W. Koch, "Transient Oscillations and Dynamics Stark Effect in Semiconductors," *Phys. Rev. B.* **38**:7607–7614 (1988).
83. J. P. Sokoloff, M. Joffe, B. Fluegel, D. J. Hulin, M. Lindberg, S. W. Koch, A. Migus, A. Antonetti, and N. Peyghambarian, "Transient Oscillations in the Vicinity of Excitons and in the Band of Semiconductors," *Phys. Rev. B.* **38**:7615–7621 (1988).
84. H. Haug and S. Schmitt-Rink, "Basic Mechanisms of the Optical Nonlinearities of Semiconductors near the Band Edge," *J. Opt. Soc. Am. B.* **2**:1135–1142 (1985).
85. H. Haug and S. W. Koch, *Quantum Theory of the Optical and Electronic Properties of Semiconductors*, World Scientific, Singapore, 1993.
86. S. Schmitt-Rink, D. S. Chemla, and D. A. B. Miller, "Linear and Nonlinear Optical Properties of Semiconductor Quantum Wells," *Advances in Physics*, **38**:89–188 (1989).
87. D. S. Chemla, "Ultrafast Transient Nonlinear Optical Processes in Semiconductors," in *Semiconductors and Semimetals*, Academic Press, Boston, 1999, vol. 58, pp. 175–256.
88. S. Hunsche, H. Heesel, A. Ewertz, H. Kurz, and J. H. Collet, "Spectral-Hole Burning and Carrier Thermalization in GaAs at Room Temperature," *Phys. Rev. B.* **48**:17818–17826 (1993).
89. L. Schultheis, J. Kuhl, A. Honold, and C. W. Tu, "Ultrafast Phase Relaxation of Excitons via Exciton-Exciton and Exciton-Electron Collisions," *Phys. Rev. Lett.* **57**:1635–1638 (1986).
90. A. Honold, L. Schultheis, J. Kuhl, and C. W. Tu, "Collision Broadening of Two-Dimensional Excitons in a GaAs Single Quantum Well," *Phys. Rev. B.* **40**:6442–6445 (1989).
91. L. Schultheis, A. Honold, J. Kuhl, K. Kohler, and C. W. Tu, "Optical Dephasing of Homogeneously Broadened Two-Dimensional Exciton Transitions in GaAs Quantum Wells," *Phys. Rev. B.* **34**:9027–9030 (1986).
92. J. Kuhl, A. Honold, L. Schultheis, and C. W. Tu, "Optical Dephasing and Orientational Relaxation of Wannier Excitons and Free Carriers in GaAs and GaAs/AlGaAs Quantum Wells," in *Advances in Solid State Physics*, U. Rössler (ed.) Friedrich Vieweg & Sohn, Braunschweig/Wiesbaden, 1989, vol. 29, pp. 157–181.

93. U. Siegner, D. Weber, E. O. Göbel, D. Bennhardt, V. Heuckeroth, R. Saleh, S. D. Baranovskii, P. Thomas, H. Schwab, C. Klingshirn, J. M. Hvam, and V. G. Lyssenko, "Optical Dephasing in Semiconductor Mixed Crystals," *Phys. Rev. B.* **46**:4564–4581 (1992).
94. S. T. Cundiff and D. G. Steel, "Coherent Transient Spectroscopy of Excitons in GaAs-AlGaAs Quantum Wells," *IEEE J. Quantum Electron.* **28**:2423 (1992).
95. E. O. Göbel, K. Leo, T. C. Damen, J. Shah, S. Schmitt-Rink, W. Schäfer, J. F. Müller, and K. Köhler, "Quantum Beats of Excitons in Quantum Wells," *Phys. Rev. Lett.* **64**:1801–1804 (1990).
96. K. Leo, T. C. Damen, J. Shah, E. O. Göbel, and K. Köhler, "Quantum Beats of Light and Heavy Hole Excitons in Quantum Wells," *Appl. Phys. Lett.* **57**:19–21 (1990).
97. B. F. Feuerbacher, J. Kuhl, R. Eccleston, and K. Ploog, "Quantum Beats Between the Light and Heavy Hole Excitons in a Quantum Well," *Solid State Communications* **74**:1279–1283 (1990).
98. K. Leo, T. C. Damen, J. Shah, and K. Köhler, "Quantum Beats of Free and Bound Excitons in GaAs/Al_xGa_{1-x}As Quantum Wells," *Phys. Rev. B.* **42**:11359 (1990).
99. J. Feldmann, T. Meier, G. V. Plessen, M. Koch, E. O. Göbel, P. Thomas, G. Bacher, C. Hartmann, H. Schweizer, W. Schäfer, and N. Nickel, "Coherent Dynamics of Excitonic Wave Packets," *Phys. Rev. Lett.* **70**:3027–3030 (1993).
100. M. Koch, J. Feldmann, G. V. Plessen, E. O. Göbel, P. Thomas, and K. Köhler, "Quantum Beats versus Polarization Interference: An Experimental Distinction," *Phys. Rev. Lett.* **69**:3631–3634 (1992).
101. V. G. Lyssenko, J. Erland, I. Balslev, K.-H. Pantke, B. S. Razbirin, and J. M. Hvam, "Nature of Nonlinear Four-Wave-Mixing Beats in Semiconductors," *Phys. Rev. B.* **48**:5720–5723 (1993).
102. K. Leo, M. Wegener, J. Shah, D. S. Chemla, E. O. Göbel, T. C. Damen, S. Schmitt-Rink, and W. Schäfer, "Effects of Coherent Polarization Interaction on Time-Resolved Degenerate Four-Wave Mixing," *Phys. Rev. Lett.* **65**:1340–1343 (1990).
103. M. Wegener, D. S. Chemla, S. Schmitt-Rink, and W. Schäfer, "Line Shape of Time-Resolved Four-Wave Mixing," *Phys. Rev. A.* **42**:5675–5683 (1990).
104. D. S. Kim, J. Shah, T. C. Damen, W. Schäfer, F. Jahnke, S. Schmitt-Rink, and K. Köhler, "Unusual Slow Temporal Evolution of Femtosecond Four-Wave-Mixing Signals in Intrinsic GaAs Quantum Wells: Direct Evidence for the Dominance of Interaction Effects," *Phys. Rev. Lett.* **69**:2725 (1992).
105. S. Weiss, M.-A. Mycek, J.-Y. Bigot, S. Schmitt-Rink, and D. S. Chemla, "Collective Effects in Excitonic Free Induction Decay: Do Semiconductors and Atoms Emit Coherent Light in Different Ways?" *Phys. Rev. Lett.* **69**:2685–2688 (1992).
106. H. Wang, K. Ferrio, D. G. Steel, Y. Z. Hu, R. Binder, and S. W. Koch, "Transient Nonlinear Optical Response from Excitation Induced Dephasing in GaAs," *Phys. Rev. Lett.* **71**:1261–1264 (1993).
107. P. Kner, S. Bar-Ad, M. V. Marquezini, D. S. Chemla, and W. Schäfer, "Magnetically Enhanced Exciton-Exciton Correlations in Semiconductors," *Phys. Rev. Lett.* **78**:1319–1322 (1997).
108. P. Kner, W. Schäfer, R. Lövenich, and D. S. Chemla, "Coherence of Four-Particle Correlations in Semiconductors," *Phys. Rev. Lett.* **81**:5386–5389 (1998).
109. G. Bartels, A. Stahl, V. M. Axt, B. Haase, U. Neukirch, and J. Gutowski, "Identification of Higher-Order Electronic Coherences in Semiconductors by their Signature in Four-Wave-Mixing Signals," *Phys. Rev. Lett.* **81**:5880–5883 (1998).
110. P. Kner, S. Bar-Ad, M. V. Marquezini, D. S. Chemla, R. Lövenich, and W. Schäfer, "Effect of Magnetoexciton Correlations on the Coherent Emission of Semiconductors," *Phys. Rev. B.* **60**:4731 (1999).
111. P. C. Becker, H. L. Fragnito, C. H. Brito Cruz, R. L. Fork, J. E. Cunningham, J. E. Henry, and C. V. Shank, "Femtosecond Photon Echoes from Band-to-Band Transitions in GaAs," *Phys. Rev. Lett.* **61**:1647–1649 (1988).
112. J.-Y. Bigot, M. T. Portella, R. W. Schoenlein, J. E. Cunningham, and C. V. Shank, "Two-Dimensional Carrier-Carrier Screening in a Quantum Well," *Phys. Rev. Lett.* **67**:636–639 (1991).
113. S. Arlt, U. Siegner, J. Kunde, F. Morier-Genoud, and U. Keller, "Ultrafast Dephasing of Continuum Transitions in Bulk Semiconductors," *Phys. Rev. B.* **59**:14860–14863 (1999).
114. J. A. Kash, J. C. Tsang, and J. M. Hvam, "Subpicosecond Time-Resolved Raman Spectroscopy of LO Phonons in GaAs," *Phys. Rev. Lett.* **54**:2151–2154 (1985).
115. G. Fasol, W. Hackenburg, H. P. Hughes, K. Ploog, E. Bauser, and H. Kano, "Continuous-Wave Spectroscopy of Femtosecond Carrier Scattering in GaAs," *Phys. Rev. B.* **41**:1461–1478 (1990).

116. J. H. Collet, "Screening and Exchange in the Theory of the Femtosecond Kinetics of the Electron-Hole Plasma," *Phys. Rev. B.* **47**:10279–10291 (1993).
117. A. Lohner, K. Rick, P. Leisching, A. Leitenstorfer, T. Elsaesser, T. Kuhn, F. Rossi, and W. Stolz, "Coherent Optical Polarization of Bulk GaAs Studied by Femtosecond Photon-Echo Spectroscopy," *Phys. Rev. Lett.* **71**:77–80 (1993).
118. U. Siegner, M.-A. Mycek, S. Glutsch, and D. S. Chemla, "Ultrafast Coherent Dynamics of Fano Resonances in Semiconductors," *Phys. Rev. Lett.* **74**:470–473 (1995).
119. U. Siegner, M.-A. Mycek, S. Glutsch, and D. S. Chemla, "Quantum Interference in the System of Lorentzian and Fano Magnetoexciton Resonances in GaAs," *Phys. Rev. B.* **51**:4953–4961 (1995).
120. S. Arlt, U. Siegner, F. Morier-Genoud, and U. Keller, "Ultrafast Coherent Dynamics in Semiconductor Quantum Wells for Multi-Subband Excitation in Different Density Regimes," *Phys. Rev. B.* **58**:13073–13080 (1998).
121. T. Rappen, U. Peter, M. Wegener, and W. Schäfer, "Coherent Dynamics of Continuum and Exciton States Studied by Spectrally Resolved fs Four-Wave Mixing," *Phys. Rev. B.* **48**:4879–4882 (1993).
122. M. U. Wehner, D. Steinbach, and M. Wegener, "Ultrafast Coherent Transients due to Exciton-Continuum Scattering in Bulk GaAs," *Phys. Rev. B.* **54**:R5211–R5214 (1996).
123. D. Birkedal, V. G. Lyssenko, J. M. Hvam, and K. ElSayed, "Continuum Contribution to Excitonic Four-Wave Mixing Due to Interaction-Induced Nonlinearities," *Phys. Rev. B.* **54**:R14250–R14253 (1996).
124. S. T. Cundiff, M. Koch, W. H. Knox, J. Shah, and W. Stolz, "Optical Coherence in Semiconductors: Strong Emission Mediated by Nondegenerate Interactions," *Phys. Rev. Lett.* **77**:1107–1110 (1996).
125. V. P. Kalosha, M. Müller, and J. Herrmann, "Coherent-Absorber Mode Locking of Solid-State Lasers," *Optics Lett.* **23**:117–119 (1998).
126. V. P. Kalosha, M. Müller, and J. Herrmann, "Theory of Solid-State Laser Mode Locking by Coherent Semiconductor Quantum-Well Absorbers," *J. Opt. Soc. Am. B.* **16**:323 (1999).
127. J. L. Oudar, D. Hulin, A. Migus, A. Antonetti, and F. Alexandre, "Subpicosecond Spectral Hole Burning due to Nonthermalized Photoexcited Carriers in GaAs," *Phys. Rev. Lett.* **55**:2074–2077 (1985).
128. W. H. Knox, C. Hirlimann, D. A. B. Miller, J. Shah, D. S. Chemla, and C. V. Shank, "Femtosecond Excitation of Nonthermal Carrier Populations in GaAs Quantum Wells," *Phys. Rev. Lett.* **56**:1191–1193 (1986).
129. R. W. Schoenlein, W. Z. Lin, E. P. Ippen, and J. G. Fujimoto, "Femtosecond Hot-Carrier Energy Relaxation in GaAs," *Appl. Phys. Lett.* **51**:1442–1444 (1987).
130. W. H. Knox, D. S. Chemla, G. Livescu, J. E. Cunningham, and J. E. Henry, "Femtosecond Carrier Thermalization in Dense Fermi Seas," *Phys. Rev. Lett.* **61**:1290–1293 (1988).
131. J.-Y. Bigot, M. T. Portella, R. W. Schoenlein, J. E. Cunningham, and C. V. Shank, "Resonant Intervalley Scattering in GaAs," *Phys. Rev. Lett.* **65**:3429–3432 (1990).
132. P. C. Becker, H. L. Fragnito, C. H. Brito Cruz, J. Shah, R. L. Fork, J. E. Cunningham, J. E. Henry, and C. V. Shank, "Femtosecond Intervalley Scattering in GaAs," *Appl. Phys. Lett.* **53**:2089 (1988).
133. J.-P. Foing, D. Hulin, M. Joffre, M. K. Jackson, J. L. Oudar, C. Tanguy, and M. Combescot, "Absorption Edge Singularities in Highly Excited Semiconductors," *Phys. Rev. Lett.* **68**:110–113 (1992).
134. C. Tanguy and M. Combescot, "X-Ray-Like Singularities for Nonequilibrium Fermi Sea," *Phys. Rev. Lett.* **68**:1935–1938 (1992).
135. J. Kunde, U. Siegner, S. Arlt, F. Morier-Genoud, and U. Keller, "Chirp-Controlled Ultrafast Optical Nonlinearities in Semiconductors," *Appl. Phys. Lett.* **73**:3025–3027 (1998).
136. J. Kunde, U. Siegner, S. Arlt, G. Steinmeyer, F. Morier-Genoud, and U. Keller, "Potential of Femtosecond Chirp Control of Ultrabroadband Semiconductor Continuum Nonlinearities," *J. Opt. Soc. Am. B.* **16**:2285–2294 (1999).
137. W. S. Warren, H. Rabitz and M. Dahleh, "Coherent Control of Quantum Dynamics: The Dream Is Alive," *Science* **259**:1581–1589 (1993).
138. A. M. Weiner, "Femtosecond Pulse Shaping Using Spatial Light Modulators," *Rev. Sci. Instruments* **71**:1929–1960 (2000).
139. R. Takahashi, Y. Kawamura, and H. Iwamura, "Ultrafast 1.55 μm All-Optical Switching Using Low-Temperature-Grown Multiple Quantum Wells," *Appl. Phys. Lett.* **68**:153–155 (1996).
140. H. S. Loka and P. W. E. Smith, "Ultrafast All-Optical Switching in an Asymmetric Fabry-Perot Device Using Low-Temperature-Grown GaAs," *IEEE Phot. Tech. Lett.* **10**:1733–1735 (1998).

141. J. F. Whitaker, "Optoelectronic Applications of LTMBE III-V Materials," *Mater. Sci. Eng.* **B22**:61–67 (1993).
142. J. F. Ziegler, J. P. Biersack, and U. Littmark, *The Stopping and Range of Ions in Solids*, Pergamon, New York, 1989.
143. F. W. Smith, A. R. Calawa, C.-L. Chen, M. J. Manfra, and L. J. Mahoney, "New MBE Buffer Used to Eliminate Backgating in GaAs MESFETs," *IEEE Electron. Device Lett.* **9**:77–80 (1988).
144. X. Liu, A. Prasad, W. M. Chen, A. Kurpiewski, A. Stoschek, Z. Liliental-Weber, and E. R. Weber, "Mechanism Responsible for the Semi-Insulating Properties of Low-Temperature-Grown GaAs," *Appl. Phys. Lett.* **65**:3002–3004 (1994).
145. X. Liu, A. Prasad, J. Nishio, E. R. Weber, Z. Liliental-Weber, and W. Walukiewicz, "Native Point Defects in Low-Temperature Grown GaAs," *Appl. Phys. Lett.* **67**:279 (1995).
146. M. Luysberg, H. Sohn, A. Prasad, P. Specht, Z. Liliental-Weber, E. R. Weber, J. Gebauer, and R. Krause-Rehberg, "Effects of the Growth Temperature and As/Ga Flux Ratio on the Incorporation of Excess As into Low Temperature Grown GaAs," *J. Appl. Phys.* **83**:561–6, 1998.
147. Z. Liliental-Weber, J. Ager, D. Look, X. W. Lin, X. Liu, J. Nishio, K. Nichols, W. Schaff, W. Swider, K. Wang, J. Wasburn, E. R. Weber, and J. Whitaker, in *Proceedings of the Eighth Conference on Semi-insulating III-V Materials*, edited by M. Godlewski (World Scientific, Singapore, 1994), p. 305.
148. M. R. Melloch, N. Otsuka, J. M. Woodall, A. C. Warren, and J. L. Freeouf, "Formation of Arsenic Precipitates in GaAs Buffer Layers Grown by Molecular Beam Epitaxy at Low Substrate Temperatures," *Appl. Phys. Lett.* **57**:1531 (1990).
149. G. L. Witt, R. Calawa, U. Mishra, E. Weber (eds.), *Low Temperature (LT) GaAs and Related Materials*, Mat. Res. Soc. Symp. Proceedings, vol. 241 Pittsburgh, 1992.
150. G. L. Witt, "LTMBE GaAs: Present Status and Perspectives," *Mater. Sci. Eng.* **B22**:9 (1993).
151. S. Gupta, M. Y. Frankel, J. A. Valdmanis, J. F. Whitaker, G. A. Mourou, F. W. Smith, and A. R. Calawa, "Subpicosecond Carrier Lifetime in GaAs Grown by Molecular Beam Epitaxy at Low Temperatures," *Appl. Phys. Lett.* **59**:3276–3278 (1991).
152. S. Gupta, J. F. Whitaker, and G. A. Mourou, "Ultrafast Carrier Dynamics in III-V Semiconductors Grown by Molecular-Beam Epitaxy at Very Low Substrate Temperatures," *IEEE J. Quantum Electron.* **28**:2464–2472 (1992).
153. K. A. McIntosh, K. B. Nichols, S. Vergheese, and E. R. Brown, "Investigation of Ultrashort Photo-carrier Relaxation Times in Low-Temperature-Grown GaAs," *Appl. Phys. Lett.* **70**:354–356 (1997).
154. M. Haiml, U. Siegner, F. Morier-Genoud, U. Keller, M. Luysberg, R. C. Lutz, P. Specht, and E. R. Weber, "Optical Nonlinearity in Low-Temperature Grown GaAs: Microscopic Limitations and Optimization Strategies," *Appl. Phys. Lett.* **74**:3134 (1999).
155. E. S. Harmon, M. R. Melloch, J. M. Woodall, D. D. Nolte, N. Otsuka, and C. L. Chang, "Carrier Life time versus Anneal in Low Temperature Growth GaAs," *Appl. Phys. Lett.* **63**:2248–2250 (1993).
156. Y. Kostoulas, L. J. Waxer, I. A. Walmsley, G. W. Wicks, and P. M. Fauchet, "Femtosecond Carrier Dynamics in Low-Temperature-Grown Indium Phosphide," *Appl. Phys. Lett.* **66**:1821–1823 (1995).
157. Y. Kostoulas, K. B. Ucer, G. W. Wicks, and P. M. Fauchet, "Femtosecond Carrier Dynamics in Low-Temperature Grown $\text{Ga}_{0.51}\text{In}_{0.49}\text{P}$," *Appl. Phys. Lett.* **67**:3756–3758 (1995).
158. U. Siegner, R. Fluck, G. Zhang, and U. Keller, "Ultrafast High-Intensity Nonlinear Absorption Dynamics in Low-Temperature Grown Gallium Arsenide," *Appl. Phys. Lett.* **69**:2566–2568 (1996).
159. A. J. Lochtefeld, M. R. Melloch, J. C. P. Chang, and E. S. Harmon, "The Role of Point Defects and Arsenic Precipitates in Carrier Trapping and Recombination in Low-Temperature Grown GaAs," *Appl. Phys. Lett.* **69**:1465–1467 (1996).
160. G. Segsneider, T. Dekorsky, H. Kurz, R. Hey, and K. Ploog, "Energy Resolved Ultrafast Relaxation Dynamics Close to the Band Edge of Low-Temperature Grown GaAs," *Appl. Phys. Lett.* **71**:2779–2781 (1997).
161. T. S. Sosnowski, T. B. Norris, H. H. Wang, P. Grenier, and J. F. Whitaker, "High-Carrier-Density Electron Dynamics in Low-Temperature-Grown GaAs," *Appl. Phys. Lett.* **70**:3245–3247 (1997).
162. H. S. Loka, S. D. Benjamin, and P. W. E. Smith, "Optical Characterization of Low-Temperature-Grown GaAs for Ultrafast All-Optical Switching Devices," *IEEE J. Quantum Electron.* **34**:1426–1437 (1998).
163. M. B. Johnson, T. C. McGill, and N. G. Paulter, "Carrier Lifetimes in Ion-Damaged GaAs," *Appl. Phys. Lett.* **54**:2424–2426 (1989).

164. M. Lambsdorff, J. Kuhl, J. Rosenzweig, A. Axmann, and J. Schneider, "Subpicosecond Carrier Lifetimes in Radiation-Damaged GaAs," *Appl. Phys. Lett.* **58**:1881–1883 (1991).
165. A. Krotkus, S. Marcinkevicius, J. Jasinski, M. Kaminska, H. H. Tan, and C. Jagadish, "Picosecond Carrier Lifetime in GaAs Implanted with High Doses of As Ions: An Alternative Material to Low-Temperature GaAs for Optoelectronic Applications," *Appl. Phys. Lett.* **66**:3304–3306 (1995).
166. F. Ganikhanov, G.-R. Lin, W.-C. Chen, C.-S. Chang, and C.-L. Pan, "Subpicosecond Carrier Lifetimes in Arsenic-Implanted GaAs," *Appl. Phys. Lett.* **67**:3465–3467 (1995).
167. C. Jagadish, H. H. Tan, A. Krotkus, S. Marcinkevicius, K. P. Korona, and M. Kaminska, "Ultrafast Carrier Trapping in High Energy Ion Implanted Gallium Arsenide," *Appl. Phys. Lett.* **68**:2225–2227 (1996).
168. M. J. Lederer, B. Luther-Davies, H. H. Tan, C. Jagadish, M. Haiml, U. Siegner, and U. Keller, "Non linear Optical Absorption and Temporal Response of Arsenic- and Oxygen-Implanted GaAs," *Appl. Phys. Lett.* **74**:1993–1995 (1999).
169. H. H. Tan, C. Jagadish, M. J. Lederer, B. Luther-Davies, J. Zou, D. J. H. Cockayne, M. Haiml, U. Siegner, and U. Keller, "Role of Implantation-Induced Defects on the Response Time of Semiconductor Saturable Absorbers," *Appl. Phys. Lett.* **75**:1437–1439 (1999).
170. M. Haiml, U. Siegner, F. Morier-Genoud, U. Keller, M. Luysberg, P. Specht, and E. R. Weber, "Femtosecond Response Times and High Optical Nonlinearity in Beryllium Doped Low-Temperature Grown GaAs," *Appl. Phys. Lett.* **74**:1269–1271 (1999).
171. P. Specht, S. Jeong, H. Sohn, M. Luysberg, A. Prasad, J. Gebauer, R. Krause-Rehberg, and E. R. Weber, "Defect Control in As-rich GaAs," *Mater. Sci. Forum* **951**:258–263 (1997).

LASER-INDUCED DAMAGE TO OPTICAL MATERIALS

Marion J. Soileau

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

19.1 INTRODUCTION

For the general public the word “laser” brings to mind images of space ships zapping evil invaders, or perhaps earthbound military planners engineering speed-of-light weapons for destroying incoming ICBMs. However, the reality is that the things most often “zapped” or damaged by lasers are the various optical components used to make or direct laser beams.

This chapter deals with laser-induced damage (LID) in optical materials. LID refers to permanent damage produced by melting, ablation, cracking, plasma formation (spark), and so on in or on an optical material as a result of exposure to laser radiation. The LID threshold refers to the fluence or irradiance which causes such damage. In early literature on this topic LID thresholds were defined to be the average between the highest fluence or irradiance levels for which no damage was observed and the lowest levels for which damage was observed. Most recent literature defines the LID threshold as the fluence or irradiance level for which probability for damage goes to zero (a much more useful definition for use in designing laser systems). For more defect-free samples these two definitions give approximately the same values for LID thresholds.

The problem of LID is as old as the first pulsed laser. The reason for this is that the power density, or fluence, is most often largest inside the laser resonator itself. The problem can, in principle, be avoided by lowering the power density or fluence inside the laser system by appropriate expansion of the beam. However, such expansions can greatly increase the cost, weight, and volume of laser systems. In many cases, larger size components are simply not available.

The net result is that systems are constrained by cost, weight, and volume to be as compact as possible. In practical terms, this means that the LID thresholds of critical components form the limit of laser output power, or energy, for many pulsed laser systems.

Thus, the critical nature of LID has led to many efforts to understand the mechanisms of LID and to improve the damage resistance of optical components. As will be shown later, the LID processes are often highly nonlinear as well as complex. The range of operating wavelengths of lasers (infrared to the ultraviolet), pulse widths (continuous output to tens of femtoseconds), and variable repetition rates further complicate the situation.

There is vast literature on the subject of LID. A good start for data on LID of optical materials can be found in the proceedings of the annual Conference on Laser-Induced Optical Materials, also known as Boulder Damage Symposium, held in Boulder, Colorado, since 1969.¹ These proceedings

are available from the SPIE in hardcopy, online, or a set of CD-ROMs. The CDs and online versions are searchable by topics, key words, and authorship. The proceedings' of the 2008 meeting contain review articles on LID to surfaces and mirrors,² materials and measurements,³ thin films,⁴ and fundamental mechanisms,⁵ for the previous 40 years of research on these topics. Much of the material in this chapter on fundamental mechanism is taken from Ref. 5.

19.2 PRACTICAL ESTIMATES

With few exceptions LID practical limits are determined at the surfaces and interfaces of optical materials. The reason for this is that these are the locations with the most defects and impurities. LID at such locations can occur at fluences that are orders of magnitude smaller than that in the bulk of materials, and below that expected from measurements of average absorption of light, or required for failure through operation of some fundamental mechanisms such as multiphoton absorption.

Given the large parameter space of laser operations many authors have tried to develop scaling rules for LID to aid the design of laser systems.⁶ These models try to incorporate parametric dependence in various mechanisms but the parameter space can be huge. For example, in one of the earlier papers on LID⁷ E.L. Bliss noted that LID depends on laser frequency, pulse duration, beam diameter, temperature, beam focusing, and the details of the materials growth and preparation. Even after more than 40 years of study, the parametric dependence of the various mechanisms of LID are not understood well enough to allow accurate scaling models that can ensure adequate design of laser systems.

However, there is a "rule of thumb," which gives a good starting point in estimating limits imposed by LID:⁵

$$E_d = (10 \text{ J/cm}^2)(t_p/1 \text{ ns})^{1/2} \quad (1)$$

where E_d = damage threshold fluence, J/cm^2
 t_p = laser pulse width, ns

Equation (1) gives an *approximation* of the threshold fluence to within plus or minus an order of magnitude for optical materials (transparent dielectrics) over the wavelength range from the UV to the infrared. This is admittedly a gross over simplification. The proper use of this equation is simply to get an idea as to whether or not one should be concerned about the possibility of LID. Laser systems designs that anticipate fluences on the high side of this equation will likely be dominated by LID considerations. For fluences below this range, one can reasonably assume that LID will not be a major factor limiting system performance.

Note that Eq. (1) has little theoretical basis except for the case where damage is caused by a thin absorbing film on a surface or interface. For that case, damage fluence should be scaled as the square root of the pulse width due to one-dimensional heat dissipation into the surface.

If one must know the damage threshold to better than the approximation in Eq. (1), then one should make a careful measurement of the threshold for conditions (t_p , wavelength, etc.) similar to those expected in use.

19.3 SURFACE DAMAGE

A good optical surface can be polished and cleaned to have optical absorption of order 10^{-4} of the incident radiation,⁸ which by itself should not result in optical damage. However, surfaces are subject to contamination and will have some amount of subsurface defects and imbedded impurities⁹ (such as polishing particles). While minor surface contamination and point defects may not substantially contribute to the average linear absorption of light by a surface such factors can, and usually do, dominate LID of optical surfaces. Examination of surface damage sites indicates a morphology consistent with defect-initiated damage.

Surface damage thresholds can be raised by

- Using grinding and polishing techniques that produce minimum subsurface damage¹⁰
- Postpolishing etching or annealing¹¹
- Proper surface cleaning¹²

A note of caution: Care must be used in cleaning optical surfaces as to not scratch the surface and not leave behind residue from the cleaning process. It is best not to get surfaces dirty! Holding a laser component on its edge by an ungloved hand will result in the diffusion of body oils across the surface sufficient to cause LID at low fluences. The references cited above are just a sampling. The Boulder proceedings¹ contain over 100 papers that relate to this one topic.

Standard optical polishing techniques can lead to substantial subsurface damage and contamination that in turn leads to LID. Typical procedures for preparing optical surfaces consist of a series of polishing or grinding steps with each subsequent step using finer grinding or polishing agents. A process sometimes referred to as “controlled grinding”⁸ can raise LID thresholds. In such processes, each step first removes damaged layer (approximately 3 times the size of the previous grinding or polishing grit), then proceeds with the next step in surface finishing.

An effective method of surface annealing is pre-irradiation of surfaces with fluence levels below the single shot LID threshold of a surface. Two methods are used: the so-called N on 1 method, where multiple subthreshold irradiations are undertaken,¹³ and the S on 1 process,¹⁴ in which the surface is irradiated at levels substantially below the single-shot damage threshold and then ramped up to the expected operational level. Substantial increases in LID fluences have been observed by using these techniques.

Entrance versus Exit Surface Damage

Consider an optical component illuminated by laser radiation at normal incidence. An interesting damage phenomenon is that for a given power or fluence incident upon an optical component, the exit (or rear) surface will damage before the entrance (or front) surface. This is at first counterintuitive, since Fresnel reflection loss at the entrance surfaces results in less energy reaching the rear surface.

The puzzling fact that the exit surface damages first provided an important clue as to the more fundamental aspects of LID. Several exotic explanations were proposed, but the explanation from linear optics is simply that the electric field associated with the laser is larger at the exit surface than at the entrance (even though the beam is partially depleted by Fresnel reflection at the front surface).¹⁵

This is understood by considering that the reflected wave at the entrance surface is 180° out of phase with the incident field (Fresnel reflection for propagating from a low-index medium to a high-index medium). At the exit surface, propagation is from the high-index medium (say, glass) to the low (say, air). In that case, the reflected field is in phase with the incident field. Therefore, the exit surface damages first simply because the electric field is higher at the rear surface. The fluence and irradiance are proportional to the electric field squared thus for even linear damage mechanisms, for example, simple absorption, the LID thresholds (expressed in terms of incident energy or power) for the exit surface will be lower than that of the entrance surface.

This field dependence of LID has played an important role in the design of thin film coatings (“move” the E-field away from coating interfaces), understanding the effects of cracks, voids, and other surface defects, and in formulation of various models for fundamental mechanisms of LID.

LID in Optical Coatings

Section 19.3 emphasizes the dominance of surface damage in most damage-limited laser systems. In fact, most optical surfaces in high-performance laser systems are coated with thin films (usually multiple layers) to control or manage reflectivity. Examples include high-reflection mirrors, antireflection coatings on transmitting optics, partial reflecting beam splitters, and wavelength selectivity elements.

LID to thin films is the most likely factor limiting the performance of high-power laser systems. (See Ref. 4 and references contained therein.) Multilayer (ML) coatings are affected by the defects, contamination, and subsurface damage of the surfaces upon which they are applied. In addition, each coating interface contributes an additional surface to the system and thus additional opportunities for surface defects to effect damage. Among other thin-film problems affecting LID are: adhesion; differential expansion; defects in the films themselves; porosity, which allows diffusion of defects; and enhancement of E-fields within the thin film stack.

A few general statements and guidelines to minimize thin-film damage are

1. Proper surface preparation and surface cleaning is essential to maximize LID thresholds in thin films.
2. The problem is more severe for transmitting thin films. Among the issues is the simple fact that the film-substrate interface cannot be avoided and this is a common initiator of LID. In general, antireflection (AR) coatings exhibit lower thresholds than high-reflection (HR) coatings.
3. Stress in thin films can reduce LID thresholds so that coating materials choices and deposition conditions should be selected to minimize internal stress in the thin-film stack.
4. In general, lower index coatings will have higher damage thresholds than higher index films. However, multilayer systems must have both high- and low-index components and the smaller the index difference, the more layers are needed to achieve the same desired reflectivity or antireflectivity.
5. Amorphous coatings have less problems with impurity diffusion than do single or multicrystalline materials and so are better choices when available.
6. Inherent in thin-film design is arranging thicknesses and indices so that multiple reflections either enhance the E-field or reduce the E-field. One important design consideration is to ensure that regions of enhanced E-field are moved away from surfaces and interfaces. The simple addition of a half-wave layer can move the high-field region away from coatings interfaces.¹⁶
7. Modern thin-film design codes allow easy design of extremely complex multilayer systems for managing reflection and transmission of multiple laser systems and wavelengths. However, designers should be mindful that larger multilayer stacks and more complex designs are likely to significantly lower the LID thresholds of these optical elements and thus limit system performance. The general design rule is “keep it simple!”

Multishot LID

LID for sites experiencing multiple irradiations are different than those of sites irradiated once. Several multishot effects have been observed. Irradiation below the single-shot damage threshold can produce LID thresholds substantially larger than the single-shot values for surfaces, thin films, and even in the bulk of a material. This is sometimes referred to as laser conditioning or laser annealing. The N on 1¹³ and S on 1¹⁴ previously discussed are examples of this phenomena.

N on 1 or S on 1 conditioning are typically used in very large, high-power laser systems such as those used in laser fusion experiments. Such systems are typically dominated by LID limits so extra efforts in laser conditioning are essential for reliably achieving maximum output power. In some cases the LID threshold seems to be reduced by repetitive pulses. In such cases, small damage sites, which initially go undetected, grow in size until the damage reaches catastrophic portions. There has been much work done on this critical question. In some cases, damaged sites can be effectively annealed by multiple shots below the damage threshold, thus extending the useful life of a component.

19.4 PACKAGE-INDUCED DAMAGE

The cliché “cleanliness is next to godliness” is usually associated with LID. Care must be taken to avoid dust or other particulates on optical surfaces and samples should be handled with gloved hands to avoid diffusion of body oils onto surfaces. Exposure to water and other solvents can

substantially lower LID thresholds, and indeed many procedures meant to clean surfaces can leave behind residue that can substantially reduce the LID threshold.

The problem of contamination of optics has led to construction of laser systems, in which the laser cavity is backfilled with clean, inert gas (such as nitrogen). This approach is particularly attractive for laser systems that are expected to operate for extended periods without possibility of access or repair, for example, a laser rangefinder in orbit around Mars to map the planet's surface.

A nonintuitive consequence of such packaging is a substantial reduction in system's lifetime due to LID, even for systems conservatively designed with fluence levels safely below the expected LID thresholds. What is observed in repetitively pulsed laser systems is that any residual hydrocarbons inside the sealed system can be dissociated due to complex interactions with the packaging gas (which acts like a buffer), the residual hydrocarbons, and the laser pulses. The "freed" carbon in the system can then be deposited, with the assistance of the laser light, on the various optical surfaces. The carbon is deposited on the surfaces in the beam path. LID will be initiated at such sites if these carbon deposits become sufficiently dense.¹⁷

The problem of package-induced LID can be avoided by back filling the laser system with oxygen (or just air.) In this case any free carbon from laser-assisted decomposition of hydrocarbons forms carbon dioxide gas, which is transparent for lasers operated from the UV to the infrared. Since any free carbon is combined with oxygen, laser-assisted carbon deposition is avoided.

19.5 NONLINEAR OPTICAL EFFECTS

Since the time of James Clark Maxwell, it has been known that light is an electromagnetic wave. However, for "natural," incoherent light there are no easily observed effects of the E-field (other than that of the "carrier" of the light energy in the Poynting vector). The field strength in some cases results in substantial nonlinear optical (NLO) effects, some of which can cause LID. Examples include electrostriction (density changes due to an impressed electric field) and E-field-dependent electronic polarizability of optical materials.¹⁸ These effects result in an increase in the index of refraction in regions of high fields. This in turn produces lensing within a material, which further increases the E-field, resulting in further increases in self-lensing until the threshold for damage is reached. This process, called self-focusing and is reviewed in the classic paper by Marburger.¹⁸

An important feature of self-focusing is that one can calculate a critical *power* at which it takes place. This is the beam power at which nonlinear refractive effects (self-focusing) overcome linear refractive effects (diffraction) leading to beam collapse. For a beam focused into the bulk of the sample this collapse can result in LID.¹⁹

Among the many effects of self-focusing is to confuse efforts to identify fundamental mechanisms of LID and the dependence of LID thresholds on various parameters, for example, pulse width, spot size, and wavelength dependence. Since self-focusing depends on the state of polarization of the light, linear or circular, a simple test for its contribution to LID in isotropic media is to measure the polarization dependence of the damage threshold.^{20,21} If a polarization dependence is found the measured threshold is due to self-focusing and not the mechanism of LID.

19.6 AVOIDANCE OF DAMAGE

Given the NLO effects that influence damage, it is essential to avoid situation where field enhancement occurs. A simple example is diffraction associated with an aperture or edge. This diffraction produces regions of enhanced (coherent sum of fields) electric fields. This enhancement of fields can also cause regions of self-focusing along the beam, which in turn cause damage.

Such effects are minimized through the use of single-mode beams, avoidance of linear diffraction by apertures or other hard edges, and special filters to remove high spatial frequencies within the beam. These techniques are particularly useful in oscillator-amplifier systems. The amplifier(s)

will further enhance the interference patterns produced by diffraction, and thus increase the likelihood of self-focusing causing damage.

19.7 FUNDAMENTAL MECHANISMS

A subset of LID is laser-induced breakdown (LIB). LIB is the damage observed with pulsed lasers for highly transparent materials. Early observations of LIB include the following:

1. Abrupt truncation of the beam transmitted through the sample, that is, threshold like behavior. (See e.g., the work of Anthes and Bass²² that showed that when damage occurs in a highly transparent material it goes from transparent to highly absorbing within less than 2 ps.)²²
2. Bright, hot plasma accompanying damage.
3. Damage occurs at the peak of the pulse when at the threshold irradiance.
4. LIB in transparent dielectrics is the most complicated since it depends on materials, properties (most of which are not properly understood), laser parameters (wavelength and pulse width), beam spatial modes, and focusing conditions.

These observations plus the lack of any reasonable alternative for energy coupling in otherwise transparent materials with bandgaps 5 to 20 times the photon energy of the laser that produced damage led early investigators to conclude that LID was due to electron avalanche breakdown.

As early as in 1965 Yasojima showed²³ that

1. Laser damage in NaCl at 10.6 μm showed LIB, that is, avalanche like behavior.
2. LIB could be substantially reduced by “seeding” the damage with a flash of blue light, thus producing “starter” electrons for the LIB.
3. The author speculated that this avalanche was seeded by nonlinear processes.

This was a remarkable bit of work and went unnoticed in work reported in Western and Soviet literature on this subject.

Based upon his seeding experiments, Yasojima speculated that LIB was initiated by either multiphoton absorption or tunneling to produce free carrier that could be rapidly accelerated by the several megavolt E-fields associated with focused pulsed laser beams.

In the 43 years, since Yasojima’s work, most work done by others has produced further evidence supporting the ideas of this unrecognized scholar!

Great progress in understanding multiphonon absorption—new materials, new methods, and design considerations—has been made over the past 40 years.¹ Among the many important fundamental results reported are multiphonon limits of linear absorption, very good models for LID to metals, the role of stress in high-power CW applications and the effects of absorption on beam quality. Many of these advances were made possible by progress in precision measurement of small levels of absorption, surface roughness, and scattering.

Linear absorption by defects, inclusions, and contamination remains practical concern, but the fundamentals of damage due to defects are reasonably well understood. LIB, the catastrophic, threshold like process in highly transparent materials has been the subject of much controversy and is where understanding still lags. The quest for understanding fundamental limits runs head on with the many varied and complex nonlinear processes involved. However, there is general agreement that for damage to occur there must be absorption of laser radiation by the material.

Beyond that there is little agreement, or perhaps more properly stated, no detailed understanding as to how a material goes from highly transparent to totally absorbing within just a few picoseconds.²² The discussion of LIB mechanism is a discussion of the process(s) by which this sort of “change of phase” occurs.

One reason to seek understanding of fundamental mechanisms is to be able to do reasonable scaling with respect to various parameters such as pulse width, wavelength, beam properties, and

materials' properties. Another reason to seek this understanding is that if the fundamental mechanisms of failure are understood, ways to improve materials could be found.

Self-focusing complicates efforts to understand LIB. Self-focusing is not in itself a mechanism for LIB. Rather, it is a mechanism for concentrating the beam so that local intensities reach critical values at which other nonlinear processes occur.

Among the major advances in improving LID limits was understanding of how to pick and sometimes design materials to minimize the nonlinear refractive index, n_2 , that governs self-focusing in a medium. Probably the most practical advance was the understanding that this problem could be managed by careful control of the mode of laser oscillators and properly accounting for self-focusing in the propagation of beams in large oscillator-amplifier systems. The key is to avoid and remove high spatial frequencies. In practical terms, this means oversizing the optics so that the beam is not clipped and using spatial filters to remove any high spatial frequencies that may occur during beam propagation. As a result, self-focusing in gain media in large laser systems such as those used in laser fusion experiments can be effectively managed. However, the problem often reappears as laser makers try to maximize energy output from gain media of limited size by multimode operation, a problem made worse by trying to compact multimode systems in limited space. This prevents the loss of high spatial frequencies and in turn results in catastrophic damage due to small-scale self-focusing.

The influence of self-focusing in damage experiments is a bit more problematic and has been a major issue in understanding the parametric dependence of LIB in transparent solids. Examples include the fact that the critical power for self-focusing scales as the square of the wavelength and as the inverse of the pulse width. Multiple mechanisms exist for self-focusing and these can scale as the spot size and pulse width together (electrostriction) and self-focusing can exhibit large dispersion. Thus self-focusing can be confused with the wavelength, pulse width, and spot size dependences of LIB.

The problem is trying to describe such a nonlinear process in enough detail such that materials may be improved, and predictions of the dependence on pulse width, wavelength, and spot size may be made with confidence. This situation was made more complex by the simple fact that no precursor to LIB has been found and experiments are hard to replicate since the sample(s) studied are consumed in the experiment and new samples often produced different results.

Early work concluded²⁴ that LIB was an intrinsic property of materials. This claim was based on the following factors:

- The observations of breakdown-like behavior previously mentioned
- The observed frequency dependence (really very little frequency dependence) of LIB in NaCl from 0.35 to 10.6 μm

This work assumed that the breakdown irradiance was an intrinsic property of a material, thus any sample-to-sample variations measured by others were due to poor quality samples and any spot-size dependence was purely due to self-focusing.

Other work produced different results^{25–28} such as

- There are just not enough free carriers present in good insulators (highly transparent materials) to initiate damage for very small spot sizes, that is, it takes at least a few free electrons to be present to initiate the avalanche.
- Better materials with higher thresholds than the so-called “intrinsic” thresholds became available, some with a factor of 20 higher threshold than the previously reported “intrinsic” results.
- Subsequent frequency dependence measurements showed a decrease in threshold at 0.5 μm , inconsistent with a pure avalanche model.
- Evidence was produced that indicated that not all spot-size dependence was due to self-focusing. Feldman's²¹ work on the polarization dependence of self-focusing was used to rule out self-focusing in some spot-size dependence measurements.

Where does this leave us? Figure 1 is an illustration of the problem. Here the spot size is constant and the pulse width and wavelength dependence is given for NaCl and SiO₂. Note that in this case

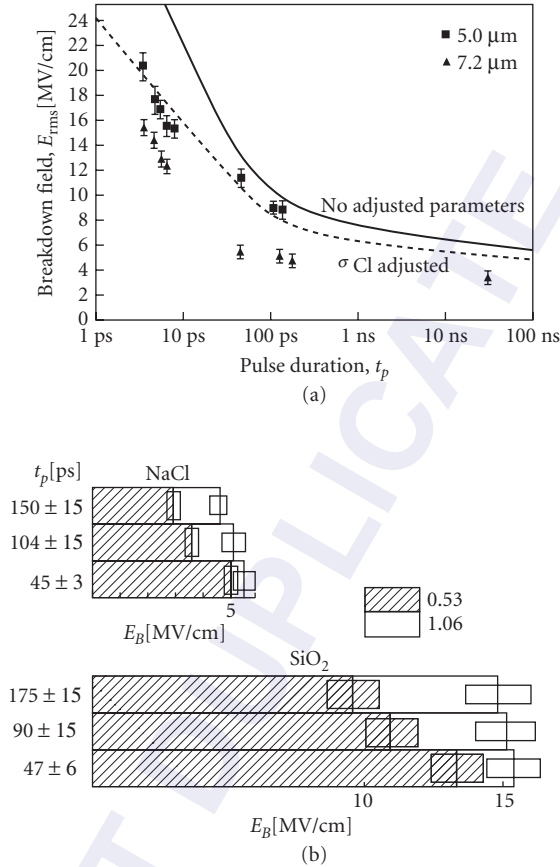


FIGURE 1 (a) The RMS breakdown field data for NaCl (78-NC-6) at 1.06 μm are plotted as a function of pulse duration, t_p . The solid line and dotted line were obtained from the theory developed by Sparks *et al.*³⁰ for NaCl at room temperature. The dotted line uses a different value for the absorption cross section for Cl ions in the theory than the solid line. (b) Wavelength dependence of the breakdown field E_B for NaCl and SiO₂ for a variety of laser pulse widths. All the above data was taken on the same sample of NaCl and the same sample of SiO₂. The 1.06 μm thresholds are taken from Ref. 2 and are interpolated from measurements made at spot size 6.1 and 10.3 μm . (From Ref. 29.)

the spot size was the same at both wavelengths, and lack of self-focusing was confirmed from measurement of lack of polarization dependence in the LIB thresholds.

The decrease in LIB field with pulse width in Fig. 1 is consistent with an avalanche process for these pulses in the picosecond regime. However, the wavelength dependence is not consistent with a purely avalanche model, and not strong enough to suggest a multiphoton model for LIB.

The pulse width dependence shown in Fig. 1 is consistent with the avalanche ionization model. Note that this model is very much dependent on the very nonlinear ionization rates, which in fact

are estimated from DC measurements! These results are *not inconsistent* with a model that assumes an avalanche initiated by multiphoton or by tunneling.

After more than 40 years of “*not inconsistent*” is the best one can do. We do not have quantitative information about avalanche ionization dependence of laser fields in solids. We have no accurate parameters for n photon absorption (n PA) (for $n > 3$) or tunneling models and we do not have definitive information that confirms that all the critical parameters have been considered.

19.8 PROGRESS IN MEASUREMENTS OF CRITICAL NLO PARAMETERS

A difficulty in the past has been lack of accurate measurements of NLO parameters. A great breakthrough was the invention of the so-called Z-scan^{31,32} technique for measuring the sign and magnitude of the nonlinear refractive index as a function of wavelength and pulse width. This technique also allows one to measure the n PA coefficients and free carrier cross section. Shown in Fig. 2^{31,32} is the Z-scan technique for measuring nonlinear refraction (with sign as well as magnitude) and nonlinear absorption.

Note, in Fig. 3 the difference in signal for negative nonlinear refraction compared to positive nonlinear refraction. The sign of nonlinear refraction changes from 1 to 0.5 μm ! By fully opening the aperture in front of the detector one can measure nonlinear absorption as is shown in Fig. 4.

The bottom curve shows the results of increasing the irradiance by about a factor of 5, and the fit is excellent, with no additional adjustment of parameters, indicating that the essential physics [linear and nonlinear absorption, positive and negative nonlinear refraction, and two photon absorption (2PA) and excited state absorption] is accounted for in this measurement. ZnSe is a 2PA material at 0.5 μm and a 4PA material at 1 μm . The nonlinearity shown is for picosecond pulses and is due to electronic effects. For nanosecond pulses electrostriction can also play a role for small spot sizes. An in-depth understanding of fundamental mechanisms LIB awaits such detailed, reproducible, non-destructive measurements of processes leading to damage.

What about many or n photon absorption (n PA)? Can this process be the mechanism leading to LIB? Figure 5³² shows results for 2, 3, and 4PA materials at 1 μm for 49 ps pulses. The net result is that for 2 or 3PA the generated carriers prevent damage by defocusing the beam. Self-defocusing

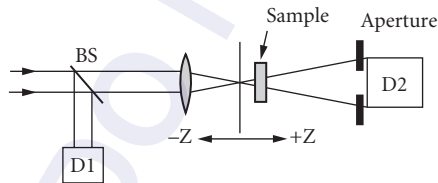


FIGURE 2 Note that the sample is scanned through the beam focus. The ratio of the readings of D1 and D2 is measured. The fit to the plot of this ratio versus Z (position relative to the focus) gives the nonlinear index of refraction. If the aperture D2 is fully open and captures all the transmitted beam then the nonlinear absorption is measured. (From Ref. 31.)

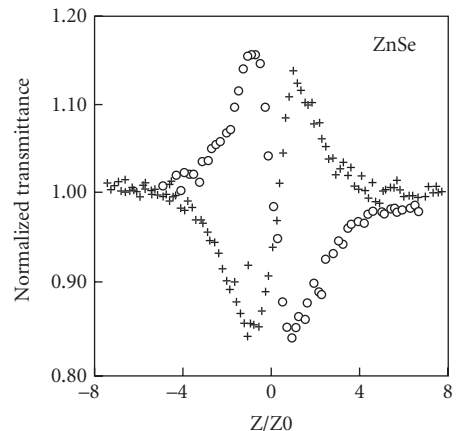


FIGURE 3 Measurement of n_2 (circles) at 1.06 and 0.53 μm (crosses) showing the sign change for ZnSe (positive n_2 and negative n_2 , respectively). (From Ref. 32.)

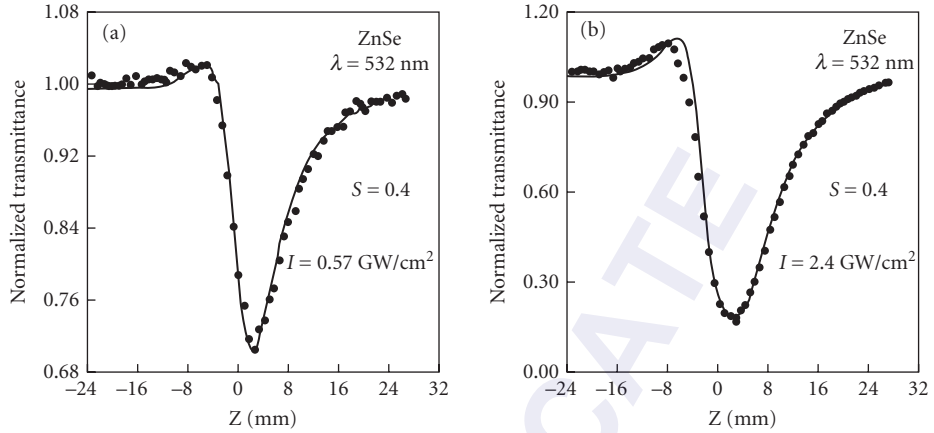


FIGURE 4 This figure shows the power of the Z-scan technique to help understand relevant nonlinear interactions, including two photon absorption, excited state absorption (free carrier absorption), and nonlinear refraction for ZnSe at 0.53 μm and 30 ps pulses. (From Ref. 31.)

of the 2 to 3PA-generated carriers prevents bulk and rear surface damage. What happens when the process is 4PA? The carrier generation rate by 4PA is not sufficient to cause defocusing and thus allows the irradiance at the exit surface to increase to the LID threshold. Recent measurements using femtosecond pulses and interferometric techniques for measuring carrier generation at focus, and so on offer hope of helping understand LIB in highly transparent materials.³³ One must be able to

Nonlinear spectroscopy

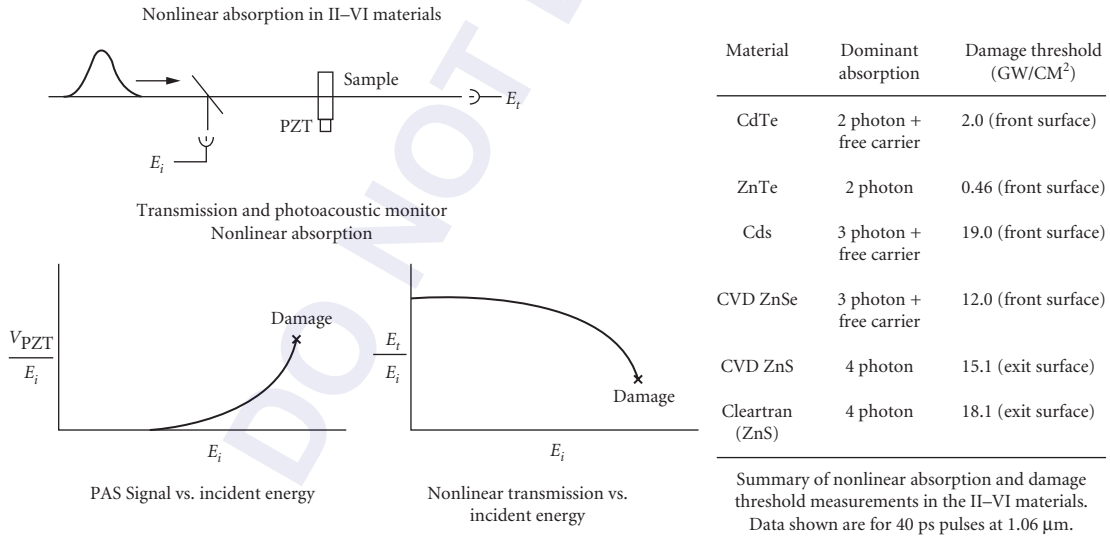


FIGURE 5 These data demonstrate different responses for nPA. The dynamics of 2 and 3PA are such that the rate of generation of free carriers is such that self-defocusing dominates propagation, preventing damage in the bulk of the material or the rear surface. For bandgaps beyond 3PA the LID of the rear surface occurs prior to field strengths that would cause significant 4PA. (From Ref. 32.)

measure the dynamical parameters with sufficient accuracy and time resolution to sort out different nonlinear parameters critical to the damage process. Only then will a more complete description and rigorous models for LID be possible.

19.9 REFERENCES

1. Proceedings of Laser-Induced Damage to Optical Materials, available from the SPIE, The International Society for Optical Engineering, Bellingham, WA, in hard-copy, via the web (<http://www.spiedigitallibrary.org>) or in searchable CD-ROM format.
2. S. Papernov and A. W. Schmid, "Laser-Induced Surface Damage of Optical Materials: Absorption Sources, Initiation, Growth, and Mitigation," in G. J. Exarhos, D. Ristau, M. J. Soileau, and C. J. Stolz (eds.), *Proceedings of the SPIE*, Vol. 7132 (2008).
3. M. Commandre, "The Transparent-Media Characterization Dedicated to Laser Damage Studies: A Key Task, Multi-Faceted and Always Renewed," in G. J. Exarhos, D. Ristau, M. J. Soileau, and C. J. Stolz (eds.), *Proceedings of the SPIE*, Vol. 7132 (2008).
4. B. E. Newnam, "Laser-Induced Damage in Thin-Film Optical Materials: 1970–2008," in G. J. Exarhos, D. Ristau, M. J. Soileau, and C. J. Stolz (eds.), to be published in the *Proceedings of the Symposium on Laser-Induced Damage to Optical Materials*, Sept. 2008, SPIE.
5. M. J. Soileau, "40 Year Retrospective of Fundamental Mechanisms" in G. J. Exarhos, D. Ristau, M. J. Soileau, and C. J. Stolz (eds.), *Proceedings of the SPIE*, Vol. 7132 (2008).
6. J. R. Bettis, R. A. House II, "Damage in Laser Materials: 1976," in A. J. Glass and A. H. Guenther (eds.), *Nat. Bur. Stand. (U.S.) Spec. Publ.* 462 (1976), p. 338., available from the SPIE, see Ref. 1.
7. E. L. Bliss, "Damage in Laser Glass, ASTM," in A. J. Glass and A. H. Guenther (eds.), *Spec. Tech. Publ.* 469 (1969), ASTM, Philadelphia, PA, available from the SPIE, see Ref. 1.
8. M. J. Soileau, H. E. Bennett, J. M. Bethke, and J. Shaffer, "Damage in Laser Materials: 1975," in A. J. Glass and A. H. Guenther (eds.), *Nat. Bur. Stand. (U.S.) Spec. Publ.* 435 (1975), available from the SPIE, see Ref. 1.
9. F. Cooke, N. Brown, and E. Prochnow, "Annular Lapping of Precision Optical Flatware," *Opt. Eng.* 15:407 (1976).
10. M. R. Kozlowski, J. Carr, I. Hutcheon, R. Torres, L. Sheehan, D. Camp, M. Yan, "Laser-Induced Damage in Optical Materials: 1997," in G. J. Exarhos, A. H. Guenther, M. R. Kozlowski, and M. J. Soileau (eds.), *Proceedings of SPIE*, Vol. 3244 (1998), p. 365.
11. V. Wang, J. E. Rudisill, C. R. Giuliano, M. Braunstein, and A. Braunstein, "Damage in Laser Materials: 1975," in A. J. Glass and A. H. Guenther (eds.), *Nat. Bur. Stand. (U.S.) Spec. Publ.* 435 (1975), available from the SPIE, see Ref. 1.
12. J. O. Porteus, P. C. Archibald, J. W. Bethke, J. H. Dancy, W. N. Faith, J. B. Franck, and P. A. Tmepel, "Influence of Cleaning Solvents, Sunlight, Humidity, and HF Gas on Pulsed Damage and Optical Characteristics of 3.8- μm Multilayer Coatings," in H. E. Bennett, A. H. Guenther, D. Milam, and B. E. Newnam (eds.), *NBS Special Publication #638, 397* (1983), available from the SPIE, see Ref. 1.
13. See for example, K. Yoshida, H. Yoshida, and S. Nakai, "Laser-Induced Damage in Optical Materials: 1991," in H. E. Bennett, L. L. Chase, A. H. Guenther, B. E. Newnam, and M. J. Soileau (eds.), *Proceedings of SPIE*, Vol. 1624 (1992).
14. For examples of S-on-1 damage see: F. Y. Ganin, K. Michlitsch, J. FLU-I, M. R. Korlouski, and P. Krulevitch, "Laser-Induced Damage in Optical Materials: 1996," in H. E. Bennett, A. H. Guenther, M. R. Kozlowski, B. E. Newman, and M. J. Soileau (eds.), *Proceedings of SPIE*, Vol. 2966 (1997).
15. M. D. Crisp, "Damage in Laser Materials: 1973," in A. J. Glass and A. H. Guenther (eds.), *Nat. Bur. Stand. (U.S.) Spec. Publ.* 387 (1973), available from the SPIE, see Ref. 1.
16. B. E. Newnam, D. H. Gill, and G. Faulkner, "Damage in Laser Materials: 1975," in A. J. Glass and A. H. Guenther (eds.), *Nat. Bur. Stand. (U.S.) Spec. Publ.* 435 (1975). Also see Ref. 4 of this chapter and references contained therein.
17. J. A. Sharps, "Laser-Induced Damage in Optical Materials: 1995," in H. E. Bennett, A. H. Guenther, M. R. Kozlowski, B. E. Newnam, and M. J. Soileau (eds.), *Proceedings of SPIE*, Vol. 2714 (1996), and the results of a Mini-Symposium on semiconductor laser lifetimes which is related to this topic is published in "Laser-Induced Damage to Optical Materials: 1997," *Proceedings of the SPIE*, Vol. 3244.

18. J. Marburger, "Theory of Self-Focusing for Fast Nonlinear Response," *NBS Special Publication #356*, SPIE, 1971, p. 51.
19. J. Marburger, "Self-Focusing: Theory," *Prog. Quant. Electr.*, No. 4, Pergamon Press, 1975, pp. 35–110, and references contained therein.
20. M. J. Soileau, M. Bass, and E. W. Van Stryland, "Frequency Dependence of Breakdown Fields in Single-Crystal NaCl and KCl," in A. J. Glass and A. H. Guenther (eds.), *NBS Special Publication #541*, SPIE, 1978, p. 309.
21. A. Feldman, D. Horowitz, and R. M. Waxler, "Mechanisms for Self-Focusing in Optical Glasses," *IEEE J. Quant. Elec.* **QE-9**:1054 (1973).
22. J. Anthes and M. Bass, "Direct Observation of the Dynamics of Picosecond-Pulse Optical Breakdown," *Appl. Phys. Lett.* **31**:412 (1977).
23. Y. Yasojima, M. Takeda, and Y. Inuishi, "Laser-Induced Breakdown in Ionic Crystals and a Polymer," *Jpn. J. Appl. Phys.* **7**(5):552 (1968).
24. N. Bloembergen, "Laser-Induced Electric Breakdown in Solids," *IEEE J. Quant. Elec.* **10**(3):375 (1974). M. J. Weber, D. Milam, and W. L. Smith, "Non-linear Refractive Index of Glasses and Crystals," *Opt. Engineering* **17**(5):489 (1978).
25. M. Sparks, "Current Status of Electron-Avalanche-Breakdown Theories," in A. J. Glass and A. H. Guenther (eds.), *NBS Special Publication #435*, SPIE, 1975, p. 331.
26. A. A. Manenkov, "New Results on Avalanche Ionization as a Laser Damage Mechanism in Transparent Solids," in A. J. Glass and A. H. Guenther (eds.), *NBS Special Publication #509*, SPIE, 1977, p. 455.
27. M. J. Soileau, M. Bass, and E. W. Van Stryland, "Frequency Dependence of Breakdown Fields in Single-Crystal NaCl and KCl," in A. J. Glass and A. H. Guenther (eds.), *NBS Special Publication #541*, SPIE, 1978, p. 309.
28. M. Soileau and M. Bass, "Optical Breakdown in NaCl and KCl from 0.53 to 10.6 Microns," *Appl. Phys. Lett.* **35**:370 (1979).
29. M. J. Soileau, W. E. Williams, E. W. Van Stryland, T. F. Boggess, and A. L. Smirl, "Temporal Dependence of Laser-Induced Breakdown in NaCl and SiO₂," in H. E. Bennett, A. H. Guenther, D. Milam, and B. E. Newnam (eds.), *NBS Special Publication #669*, SPIE, 1982, p. 387.
30. M. Sparks, T. Holstein, R. Warren, D. L. Mills, A. A. Maradudin, L. J. Sham, E. Loh, Jr., and F. King, "Theory of Electron-Avalanche Breakdown in Solids," in H. E. Bennett, A. J. Glass, A. H. Guenther, and B. E. Newnam (ed.), *NBS Special Publication #568*, SPIE, 1979, p. 467.
31. E. W. Van Stryland, M. Sheik Bahae, A. A. Said, D. J. Hagan, and M. J. Soileau, "Laser-Induced Damage in Optical Materials: 1993," in H. E. Bennett, L. L. Chase, A. H. Guenther, B. E. Newnam, and M. J. Soileau (eds.), *Proceedings of SPIE*, Vol. 2114 (1994); available from the SPIE, see Ref. 1.
32. E. W. Van Stryland, M. A. Woodall, W. E. Williams, and M. J. Soileau, "Damage in Laser Materials: 1981," in H. E. Bennett, A. H. Guenther, D. Milam, and B. E. Newnam (eds.), *Nat. Bur. Stand. (U.S.) Spec. Publ.* 638 (1983), available from the SPIE, see Ref. 1.
33. S. Garnov, G. J. Exarhos, A. H. Guenther, K. L. Lewis, D. Ristau, M. J. Soileau, C. J. Stolz (eds.), *Proceedings of the 40 Anniversary Boulder Damage Symposium*, 2008. Available from the SPIE, see Ref. 1.

PART

3

QUANTUM AND
MOLECULAR
OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

LASER COOLING AND TRAPPING OF ATOMS

Harold J. Metcalf

*Department of Physics
State University of New York
Stony Brook, New York*

Peter van der Straten

*Debye Institute
Department of Atomic and Interface Physics
Utrecht University
Utrecht, The Netherlands*

20.1 INTRODUCTION

This chapter begins with some of the general ideas about laser cooling. One of the characteristics of optical control of atomic motion is that the speed of atoms can be considerably reduced. Since the spread of velocities of a sample of atoms is directly related to its temperature, the field has been dubbed *laser cooling*, and this name has persisted throughout the years.

In Sec. 20.2 we introduce the general idea of optical forces and how they can act on atoms. We show how such forces can be velocity dependent, and thus nonconservative, which makes it possible to use optical forces for cooling. The section concludes with the discussion of a few special temperatures. Section 20.3 presents a quantum mechanical description of the origin of the force resulting from the atomic response to both stimulated and spontaneous emission processes. This is quite different from the familiar quantum mechanical calculations using state vectors to describe the state of the system, since spontaneous emission causes the state of the system to evolve from a pure state into a mixed state. Since spontaneous emission is an essential ingredient for the dissipative nature of the optical forces, the density matrix is introduced to describe it. The evolution of the density matrix is given by the optical Bloch equations (OBE), and the optical force is calculated from them. It is through the OBE that the dissipative aspects of laser cooling are introduced to the otherwise conservative quantum mechanics. The velocity dependence is treated as an extension of the force on an atom at rest.

In Sec. 20.4 the first modern laser cooling experiments are described. Atoms in beams were slowed down from thermal velocity to a few m/s, and the dominant problem was the change in Doppler shift arising from such a large change in velocity. Some typical values of parameters are discussed and tabulated. Section 20.5 introduces true cooling by optical forces to the μK regime. Such experiments require at least two laser beams, and are called *optical molasses* because the resulting viscous force can slow atoms to extremely slow velocities, and hence compress the width of the velocity distribution. The limits of such laser cooling are discussed, as well as the extension from experiments in 1D to 3D. Here the velocity dependence of the force is built into the description via the Doppler shift instead of being added in as an extension of the treatment. In 1988 some experiments reported temperatures below the limit calculated for optical molasses, and Sec. 20.6 presents the new description of laser cooling that emerged from this surprise. For the first time, the force

resulting from spontaneous emission in combination with the multiple level structure of real atoms were embodied in the discussion. Here the new limits of laser cooling are presented.

The discussion up to this point has been on atomic velocities, and thus can be described in terms of a velocity space. Laser cooling thus collects atoms near the origin of velocity space. It is also possible to collect atoms into a small region of ordinary configuration space, and such trapping is discussed in Sec. 20.7. Neutral atom traps can employ magnetic fields, optical fields, and both working together. However, such traps are always very shallow, and so only atoms that have been cooled to the few mK domain can be captured. The combination of laser cooling and atom trapping has produced astounding new tools for atomic physicists, and Sec. 20.8 describes some of the applications and uses of these wonderful new capabilities.

20.2 GENERAL PROPERTIES CONCERNING LASER COOLING

These experiments almost always involve atomic absorption of nearly resonant light. The energy of the light $\hbar\omega$ raises the internal energy of the atom, and the angular momentum \hbar changes the internal angular momentum ℓ of the electron, as described by the well-known selection rule $\Delta\ell = \pm 1$. By contrast, the linear momentum of the light $p = E/c = h/\lambda(\vec{p} = \hbar\vec{k})$ cannot be absorbed by internal atomic degrees of freedom, and therefore must change the motion of the atoms in the laboratory frame. The force resulting from this momentum exchange between the light field and the atoms can be used in many ways to control atomic motion, and is the subject of this chapter.

Absorption of light populates the atomic excited state, and the return to the ground state can be either by spontaneous or by stimulated emission. The nature of the optical force that arises from these two different processes is quite different, and will be described separately. Such atomic transitions (i.e., the motion of the atomic electrons) must be described quantum mechanically in the well-known form of the Schrödinger equation. By contrast, the center-of-mass motion of the atoms can usually be described classically, but there are many cases where even this is not possible so it must also involve quantum mechanics.

In the simplest possible case, the absorption of well-directed light from a laser beam, the momentum exchange between the light field and the atoms results in a force

$$\vec{F} = d\vec{p}/dt = \hbar\vec{k}\gamma_p \quad (1)$$

where γ_p is the excitation rate of the atoms. The absorption leaves the atoms in their excited state, and if the light intensity is low enough so that they are much more likely to return to the ground state by spontaneous emission than by stimulated emission, the resulting fluorescent light carries off momentum $\hbar k$ in a random direction. The momentum exchange from the fluorescence averages zero, so the net total force is given by Eq. (1).

The scattering rate γ_p depends on the laser detuning from atomic resonance $\delta = \omega_\ell - \omega_a$, where ω_ℓ is the laser frequency and ω_a is the atomic resonance frequency. This detuning is measured in the atomic reference frame, and it is necessary that the Doppler-shifted laser frequency in the moving atoms' reference frame be used to calculate the absorption and scattering rate. Then γ_p is given by the Lorentzian

$$\gamma_p = \frac{s_0\gamma/2}{1 + s_0 + [2(\delta + \omega_D)/\gamma]^2} \quad (2)$$

where $\gamma \equiv 1/\tau$ is the angular frequency corresponding to the decay rate of the excited state. Here $s_0 = I/I_s$ is the ratio of the light intensity I to the saturation intensity $I_s \equiv \pi\hbar c/3\lambda^3\tau$, which is a few mW/cm² for typical atomic transitions (λ is the optical wavelength). The Doppler shift seen by the moving atoms is $\omega_D = -\vec{k} \cdot \vec{v}$ (note that \vec{k} opposite to \vec{v} produces a positive Doppler shift). The force is thus velocity-dependent, and the experimenter's task is to exploit this dependence to the desired goal, for example, optical friction for laser cooling.

The spontaneous emission events produce unpredictable changes in atomic momenta so the discussion of atomic motion must also include a “random walk” component. This can be described as a diffusion of the atomic momenta in momentum space, similar to Brownian motion in real space. The evolution of the momentum in such circumstances is described by the Fokker-Planck equation, and it can be used for a more formal treatment of the laser cooling process. Solutions of the Fokker-Planck equation in limiting cases can ultimately be used to relate the velocity distribution of the atoms with their temperature.

The idea of “temperature” in laser cooling requires some careful discussion and disclaimers. In thermodynamics, temperature is carefully defined as a parameter of the state of a closed system in thermal equilibrium with its surroundings. This, of course, requires that there be thermal contact (i.e., heat exchange) with the environment. In laser cooling this is clearly not the case because a sample of atoms is always absorbing and scattering light. Furthermore, there is essentially no heat exchange (the light cannot be considered as heat even though it is indeed a form of energy). Thus the system may very well be in a steady-state situation, but certainly not in thermal equilibrium, so that the assignment of a thermodynamic “temperature” is completely inappropriate.

Nevertheless, it is convenient to use the label of temperature to describe an atomic sample whose average kinetic energy $\langle E_k \rangle$ in one dimension has been reduced by the laser light, and this is written simply as $k_B T/2 = \langle E_k \rangle$, where k_B is Boltzmann’s constant. It must be remembered that this temperature assignment is absolutely inadequate for atomic samples that do not have a Maxwell-Boltzmann velocity distribution, whether or not they are in thermal equilibrium: there are infinitely many velocity distributions that have the same value of $\langle E_k \rangle$ but are so different from one another that characterizing them by the same “temperature” is a severe error.

With these ideas in mind, it is useful to define a few rather special values of temperatures associated with laser cooling. The highest of these temperatures corresponds to the energy associated with atoms whose speed and concomitant Doppler shift puts them just at the boundary of absorption of light. This velocity is $v_c \equiv \gamma/k \sim \text{few m/s}$, and the corresponding temperature is $k_B T_c \equiv M\gamma^2/k^2$, and is typically several mK. (Here M is the atomic mass.)

The next characteristic temperature corresponds to the energy associated with the natural width of atomic transitions, and is called the *Doppler temperature*. It is given by $k_B T_D \equiv \hbar\gamma/2$. Because it corresponds to the limit of certain laser cooling processes, it is often called the *Doppler limit*, and is typically several hundred μK . Associated with this temperature is the one-dimensional velocity $v_D = \sqrt{k_B T_D/M} \sim 30 \text{ cm/s}$.

The last of these three characteristic temperatures corresponds to the energy associated with a single photon recoil. In the absorption or emission process of a single photon, the atoms obtain a recoil velocity $v_r \equiv \hbar k/M$. The corresponding energy change can be related to a temperature, the *recoil limit*, defined as $k_B T_r \equiv \hbar^2 k^2/M$, and is generally regarded as the lower limit for optical cooling processes (although there are a few clever schemes that cool below it). It is typically a few μK , and corresponds to speeds of $v_r \sim 1 \text{ cm/s}$.

These three temperatures are related to one another through a single dimensionless parameter $\varepsilon \equiv \omega/\gamma$ that is ubiquitous in describing laser cooling. It is the ratio of the recoil frequency $\omega_r \equiv \hbar k^2/2M$ to the natural width γ , and as such embodies most of the important information that characterize laser cooling on a particular atomic transition. Typically $\varepsilon \sim 10^{-3} - 10^{-2}$, and clearly $T_r = 4\varepsilon T_D = 4\varepsilon^2 T_c$.

In laser cooling and related aspects of optical control of atomic motion, the forces arise because of the exchange of momentum between the atoms and the laser field. Since the energy and momentum exchange is necessarily in discrete quanta rather than continuous, the interaction is characterized by finite momentum *kicks*. This is often described in terms of *steps* in a fictitious space whose axes are momentum rather than position. These steps in momentum space are of size $\hbar k$ and thus are generally small compared to the magnitude of the atomic momenta at thermal velocities \bar{v} . This is easily seen by comparing $\hbar k$ with $M\bar{v}$,

$$\frac{\hbar k}{M\bar{v}} = \sqrt{\frac{T_r}{T}} \ll 1 \quad (3)$$

Thus the scattering of a single photon has a negligibly small effect on the motion of thermal atoms, but repeated cycles of absorption and emission can cause a large change of the atomic momenta and velocities.

20.3 THEORETICAL DESCRIPTION

Force on a Two-Level Atom

We begin the calculation of the optical force on atoms by considering the simplest schemes, namely, a single-frequency light field interacting with a two-level atom confined to one dimension. It is based on the interaction of two-level atoms with a laser field as discussed in many textbooks.¹

The philosophy of the correspondence principle requires a smooth transition between quantum and classical mechanics. Thus the force F on an atom is defined as the expectation value of the quantum mechanical force operator \mathcal{F} , as defined by $F = \langle \mathcal{F} \rangle = d\langle p \rangle / dt$. The time evolution of the expectation value of a time-independent quantum mechanical operator \mathcal{A} is given by

$$\frac{d}{dt} \langle \mathcal{A} \rangle = \frac{i}{\hbar} \langle [\mathcal{H}, \mathcal{A}] \rangle \quad (4)$$

The commutator of \mathcal{H} and P is given by $[\mathcal{H}, p] = i\hbar(\partial\mathcal{H}/\partial z)$, where the operator p has been replaced by $-i\hbar(\partial/\partial z)$. The force on an atom is then given by

$$F = - \left\langle \frac{\partial\mathcal{H}}{\partial z} \right\rangle \quad (5)$$

This relation is a specific example of the Ehrenfest theorem and forms the quantum mechanical analog of the classical expression that the force is the negative gradient of the potential.

Discussion of the force on atoms caused by light fields begins with that part of the Hamiltonian that describes the electric dipole interaction between the atom and the light field. The electric field of the light is written as $\mathcal{E}(\vec{r}, t) = E_0 \hat{e} \cos(kz - \omega t)$ and the interaction Hamiltonian is $\mathcal{H}' = e\mathcal{E}(\vec{r}, t) \cdot \vec{r}$ where \vec{r} is the electron coordinate. It has only off-diagonal matrix elements given by $\mathcal{H}'_{eg} = -eE_0 \hat{e} \cdot \langle e | \vec{r} | g \rangle$ where e and g represent the excited and ground states respectively. The force depends on the atomic state as determined by its interaction with the light, and is calculated from the expectation value $\langle \mathcal{A} \rangle = \text{Tr}(\rho \mathcal{A})$ as in Eq. (4), where ρ is the density matrix found by solving the optical Bloch equations (OBE).¹ Then

$$F = \hbar \left(\frac{\partial\Omega}{\partial z} \rho_{eg}^* + \frac{\partial\Omega^*}{\partial z} \rho_{eg} \right) \quad (6)$$

where the Rabi frequency is defined as $\hbar\Omega \equiv \mathcal{H}'_{eg}$. Note that the force depends on the state of the atom, and in particular, on the optical coherence between the ground and excited states, ρ_{eg} .

Although it may seem a bit artificial, it is instructive to split $\partial\Omega/\partial z$ into its real and imaginary parts (the matrix element that defines Ω can certainly be complex):

$$\frac{\partial\Omega}{\partial z} = (q_r + iq_i)\Omega \quad (7)$$

Here $q_r + iq_i$ is the logarithmic derivative of Ω . In general, for a field $E(z) = E_0(z) \exp(i\phi(z)) + \text{c.c.}$, the real part of the logarithmic derivative corresponds to a gradient of the amplitude $E_0(z)$ and the imaginary part to a gradient of the phase $\phi(z)$. Then the expression for the force becomes

$$F = \hbar q_r (\Omega \rho_{eg}^* + \Omega^* \rho_{eg}) + i\hbar q_i (\Omega \rho_{eg}^* - \Omega^* \rho_{eg}) \quad (8)$$

Equation (8) is a very general result that can be used to find the force for any particular situation as long as the OBE for ρ_{eg} can be solved. In spite of the chosen complex expression for Ω , it is important to note that the force itself is real, and that first term of the force is proportional to the real part of $\Omega\rho_{eg}^*$, whereas the second term is proportional to the imaginary part.

A Two-Level Atom at Rest

There are two important special optical arrangements to consider. The first one is a traveling wave whose electric field is $E(z) = (E_0/2)(e^{i(kz-\omega t)} + c.c.)$. In calculating the Rabi frequency from this, the rotating wave approximation (RWA) causes the positive frequency component of $E(z)$ to drop out. Then the gradient of the Rabi frequency becomes proportional to the gradient of the surviving negative frequency component, so that $q_r = 0$ and $q_i = k$. For such a traveling wave the amplitude is constant but the phase is not, and this leads to the nonzero value of q_i .

This is in direct contrast to the case of a standing wave, composed of two counterpropagating traveling waves so its amplitude is twice as large, for which the electric field is given by $E(z) = E_0 \cos(kz)(e^{-i\omega t} + c.c.)$, so that $q_r = -k \tan(kz)$ and $q_i = 0$. Again, only the negative frequency part survives the RWA, but the gradient does not depend on it. Thus a standing wave has an amplitude gradient, but not a phase gradient.

The steady-state solutions of the OBE for a two-level atom at rest provide simple expressions for ρ .¹ Substituting the solution for ρ_{eg} into Eq. (8) gives

$$F = \frac{\hbar s}{1+s} \left(-\delta q_r + \frac{1}{2} \gamma q_i \right) \quad (9)$$

where $s \equiv s_0/[1 + (2\delta/\gamma)^2]$ is the off-resonance saturation parameter. Note that the first term is proportional to the detuning δ , whereas the second term is proportional to the decay rate γ . For zero detuning, the force for a traveling wave becomes $F = (\hbar k \gamma / 2)[s_0 / (s_0 + 1)]$, a very satisfying result because it is simply the momentum per photon $\hbar k$, times the scattering rate γ_p at resonance of Eq. (2).

It is instructive to identify the origin of both terms in Eq. (9). Absorption of light leads to the transfer of momentum from the optical field to the atoms. If the atoms decay by spontaneous emission, the recoil associated with the spontaneous fluorescence is in a random direction, so its average over many emission events results in zero net effect on the atomic momentum. Thus the force from absorption followed by spontaneous emission can be written as $F_{sp} = \hbar k \gamma \rho_{ee}$, where $\hbar k$ is the momentum transfer for each photon, γ is the rate for the process, and ρ_{ee} is the probability for the atoms to be in the excited state. Using Eq. (2), the force resulting from absorption followed by spontaneous emission becomes

$$F_{sp} = \frac{\hbar k s_0 \gamma / 2}{1 + s_0 + (2\delta/\gamma)^2} \quad (10)$$

which saturates at large intensity as a result of the term s_0 in the denominator. Increasing the rate of absorption by increasing the intensity does not increase the force without limit, since that would only increase the rate of stimulated emission, where the transfer of momentum is opposite in direction compared to the absorption. Thus the force saturates to a maximum value of $\hbar k \gamma / 2$, because ρ_{ee} has a maximum value of $1/2$.

Examination of Eq. (10) shows that it clearly corresponds to the second term of Eq. (8). This term is called the *light pressure force*, *radiation pressure force*, *scattering force*, or *dissipative force*, since it relies on the scattering of light out of the laser beam. It vanishes for an atom at rest in a standing wave where $q_i = 0$, and this can be understood because atoms can absorb light from either of the two counterpropagating beams that make up the standing wave, and the average momentum transfer then vanishes. This force is dissipative because the reverse of spontaneous emission is not possible, and therefore the action of the force cannot be reversed. It plays a very important role in the slowing and cooling of atoms.

By contrast, the first term in Eq. (8) derives from the light shifts of the ground and excited states that depend on the strength of the optical electric field. A standing wave is composed of two counter-propagating laser beams, and their interference produces an amplitude gradient that is not present in a traveling wave. The force is proportional to the gradient of the light shift, and the ground-state light shift $\Delta E_g = \hbar\Omega^2/4\delta$ can be used to find the force on ground-state atoms in low-intensity light:

$$F_{\text{dip}} = -\frac{\partial(\Delta E_g)}{\partial z} = \frac{\hbar\Omega}{2\delta} \frac{\partial\Omega}{\partial z} \quad (11)$$

For an amplitude-gradient light field such as a standing wave, $\partial\Omega/\partial z = q_r\Omega$, and this force corresponds to the first term in Eq. (8) in the limit of low saturation ($s \ll 1$).

For the case of a standing wave Eq. (9) becomes

$$F_{\text{dip}} = \frac{2\hbar k\delta s_0 \sin 2kz}{1 + 4s_0 \cos^2 kz + (2\delta/\gamma)^2} \quad (12)$$

where s_0 is the saturation parameter of each of the two beams that form the standing wave. For $\delta < 0$ the force drives the atoms to positions where the intensity has a maximum, whereas for $\delta > 0$ the atoms are attracted to the intensity minima. The force is conservative and therefore cannot be used for cooling. This is called the *dipole force*, *reactive force*, *gradient force*, or *redistribution force*. It has the same origin as the force of an inhomogeneous DC electric field on a classical dipole, but relies on the redistribution of photons from one laser beam to the other.

It needs to be emphasized that the forces of Eqs. (10) and (12) are two fundamentally different kinds of forces. For an atom at rest, the scattering force vanishes for a standing wave, whereas the dipole force vanishes for a traveling wave. The scattering force is dissipative, and can be used to cool, whereas the dipole force is conservative, and can be used to trap. Dipole forces can be made large by using high-intensity light because they do not saturate. However, since the forces are conservative, they cannot be used to cool a sample of atoms. Nevertheless, they can be combined with the dissipative scattering force to enhance cooling in several different ways. By contrast, scattering forces are always limited by the rate of spontaneous emission γ and cannot be made arbitrarily strong, but they are dissipative and are required for cooling.

Atoms in Motion

Laser cooling requires dissipative or velocity-dependent forces that cannot be conservative. The procedure followed here is to treat the velocity of the atoms as a small perturbation, and make first-order corrections to the solutions of the OBE obtained for atoms at rest.² It begins by adding drift terms in the expressions for the relevant quantities. Thus the Rabi frequency satisfies

$$\frac{d\Omega}{dt} = \frac{\partial\Omega}{\partial t} + v \frac{\partial\Omega}{\partial z} = \frac{\partial\Omega}{\partial t} + v(q_r + iq_i)\Omega \quad (13)$$

where Eq. (7) has been used to separate the gradient of Ω into real and imaginary parts. Differentiating the steady-state density matrix elements found by solving the OBE¹ leads to

$$\frac{dw}{dt} = \frac{\partial w}{\partial t} + v \frac{\partial w}{\partial z} = \frac{\partial w}{\partial t} - \frac{2vq_r s}{(1+s)^2} \quad (14)$$

since $s_0 = 2|\Omega|^2/\gamma^2$ and Ω depends on z . Here $w \equiv \rho_{gg} - \rho_{ee}$. Similarly,

$$\frac{d\rho_{eg}}{dt} = \frac{\partial\rho_{eg}}{\partial t} + v \frac{\partial\rho_{eg}}{\partial z} = \frac{\partial\rho_{eg}}{\partial t} - \frac{iv\Omega}{2(\gamma/2 - i\delta)(1+s)} \left[q_r \left(\frac{1-s}{1+s} \right) + iq_i \right] \quad (15)$$

Since neither w nor ρ_{eg} is explicitly time-dependent, both $\partial w/\partial t$ and $\partial \rho_{eg}/\partial t$ vanish. Equations (14) and (15) are still difficult to solve analytically for a general optical field, and the results are not very instructive. However, the solution for the two special cases of the standing and traveling waves provide considerable insight.

For a traveling wave $q_r = 0$, and the velocity-dependent force can be found by combining Eqs. (14) and (15) with the OBE to eliminate the time derivatives. The resulting coupled equations can be separated and substituted into Eq. (8) for the force to find, after considerable algebra,

$$F = \hbar q_i \frac{s\gamma/2}{1+s} \left(1 + \frac{2\delta v q_i}{(1+s)(\delta^2 + \gamma^2/4)} \right) \equiv F_0 - \beta v \quad (16)$$

The first term is the velocity-independent force F_0 for an atom at rest given by Eq. (9). The second term is velocity-dependent and can lead to compression of the velocity distribution. For a traveling wave $q_i = k$ and thus the damping coefficient β is given by

$$\beta = -\hbar k^2 \frac{4s_0(\delta/\gamma)}{(1+s_0 + (2\delta/\gamma)^2)^2} \quad (17)$$

Such a force can compress the velocity distribution of an atomic sample for negative values of δ (i.e., for red detuned light). For small detuning and low intensity the damping coefficient β is linear in both parameters. However, for detunings much larger than γ and intensities much larger than I_s , β saturates and even decreases as a result of the dominance of δ in the denominator of Eq. (17). This behavior can be seen in Fig. 1, where the damping coefficient β has been plotted as a function of detuning for different saturation parameters. The decrease of β for large detunings and intensities is caused by saturation of the transition, in which case the absorption rate becomes only weakly dependent on the velocity. The maximum value of β is obtained for $\delta = -\gamma/2$ and $s_0 = 2$, and is given by $\beta_{\max} = \hbar k^2/4$. The damping rate Γ is given by $\Gamma = \beta/M$, and its maximum value is $\Gamma_{\max} = \omega_r/2$, where ω_r is the recoil frequency. For the alkalis this rate is of the order of 10^4 – 10^3 s $^{-1}$, indicating that atomic velocity distributions can be compressed in about 10 to 100 μ s. Furthermore, F_0 in Eq. (16) is always present and so the atoms are *not* damped toward any constant velocity.

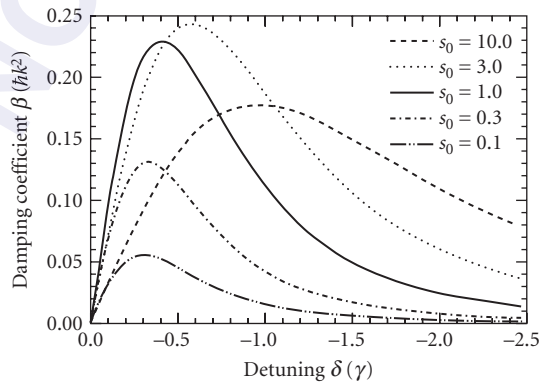


FIGURE 1 The damping coefficient β for an atom in a traveling wave as a function of the detuning for different values of the saturation parameter s_0 . The damping coefficient is maximum for intermediate detunings and intensities.

For a standing wave $q_i = 0$, and just as above, the velocity-dependent force can be found by combining Eqs. (14) and (15) with the OBE to eliminate the time derivatives. The resulting coupled equations can again be separated and substituted into Eq. (8) for the force to find

$$F = -\hbar q_r \frac{s\delta}{1+s} \left(1 - v q_r \frac{(1-s)\gamma^2 - 2s^2(\delta^2 + \gamma^2/4)}{(\delta^2 + \gamma^2/4)(1+s)^2 \gamma} \right) \quad (18)$$

where $q_r = -k \tan(kz)$. In the limit of $s \ll 1$, this force is

$$F = \hbar k \frac{s_0 \delta \gamma^2}{2(\delta^2 + \gamma^2/4)} \left(\sin 2kz + kv \frac{\gamma}{(\delta^2 + \gamma^2/4)} (1 - \cos 2kz) \right) \quad (19)$$

Here s_0 is the saturation parameter of each of the two beams that compose the standing wave. The first term is the velocity-independent part of Eq. (9) and is sinusoidal in space, with a period of $\lambda/2$. Thus its spatial average vanishes. The force remaining after such averaging is $F_{av} = -\beta v$, where the damping coefficient β is given by

$$\beta = -\hbar k^2 \frac{8s_0(\delta/\gamma)}{(1+(2\delta/\gamma)^2)^2} \quad (20)$$

In contrast to the traveling-wave case, this is a true damping force because there is no F_0 , so atoms are slowed toward $v = 0$ independent of their initial velocities. Note that this expression for β is valid only for $s \ll 1$ because it depends on spontaneous emission to return atoms to their ground state.

There is an appealing description of the mechanism for this kind of cooling in a standing wave. With light detuned below resonance, atoms traveling toward one laser beam see it Doppler shifted upward, closer to resonance. Since such atoms are traveling away from the other laser beam, they see its light Doppler shifted further downward, hence further out of resonance. Atoms therefore scatter more light from the beam counterpropagating to their velocity, so their velocity is reduced. This damping mechanism is called *optical molasses*, and is one of the most important tools of laser cooling.

Needless to say, such a pure damping force would reduce the atomic velocities, and hence the absolute temperature, to zero. Since this violates thermodynamics, there must be something left out of the description. It is the discreteness of the momentum changes in each case, $\Delta p = \hbar k$, that results in a minimum velocity change. The consequences of this discreteness can be described as a diffusion of the atomic momenta in momentum space by finite steps as discussed earlier.

The Fokker-Planck Equation

The random walk in momentum space associated with spontaneous emission is similar to Brownian motion in coordinate space. There is an analogous momentum diffusion constant D , and so the atomic motion in momentum space can be described by the Fokker-Planck equation

$$\frac{\partial W(p, t)}{\partial t} = -\frac{\partial [F(p, t)W(p, t)]}{\partial p} + \frac{\partial^2 [D(p, t)W(p, t)]}{\partial p^2} \quad (21)$$

where $W(p, t)$ is the momentum distribution of the atoms. For the special case when both the force and the diffusion are independent of time, the formal stationary solution is

$$\bar{W}(p) = \frac{C}{D(p)} \exp \left(\int_0^p \frac{F(p')}{D(p')} dp' \right) \quad (22)$$

where C is an integration constant. Once the force and diffusion are known, the stationary solution of the Fokker-Planck equation emerges easily.

In the simplest and most common case in laser cooling the force is proportional to the velocity and the diffusion is independent of velocity:

$$F(v) = -\beta v \quad \text{and} \quad D(v) = D_0 \quad (23)$$

Then the stationary solution of Eq. (21) for $\bar{W}(v)$ is

$$\bar{W}(p) \propto e^{-\beta p^2 / 2MD_0} \quad (24)$$

This is indeed a Maxwell-Boltzmann distribution. For low intensity where spontaneous emission dominates, $D_0 = s\gamma(\hbar k)^2/2$, so the steady-state temperature is given by $k_B T = D_0/\beta = \hbar\gamma/2$ for $\delta = -\gamma/2$, its optimum value.¹ This is called the *Doppler temperature* because the velocity dependence of the cooling mechanism derives from the Doppler shift. The fact that the conditions of Eq. (23) for the force and diffusion are often approximately correct explains why the notion of temperature often appears as a description of a laser-cooled sample.

One of the most important properties of laser cooling is its ability to change the phase space density of an atomic sample. Changing the phase space density provides a most important distinction between light optics and atom optics. The Hamiltonian description of geometrical optics leads to the brightness theorem, that can be found in many optics books. Thus bundles of light rays obey a similar phase space density conservation. But there is a fundamental difference between light and atom optics. In the first case, the “forces” that determine the behavior of bundles of rays are “conservative” and phase space density is conserved. For instance, a lens can be used to focus a light beam to a small spot; however, at the same time the divergence of the beam must be increased, thus conserving phase space density. By contrast, in atom optics dissipative forces that are velocity-dependent can be used, and thus phase space density is no longer conserved. Optical elements corresponding to such forces cannot exist for light, but in addition to the atom optic elements of lenses, collimators, and others, phase space compressors can also be built. Such compression is essential in a large number of cases, such as atomic beam brightening for collision studies or cooling for the achievement of Bose-Einstein condensation.

20.4 SLOWING ATOMIC BEAMS

Among the earliest laser cooling experiments was deceleration of atoms in a beam.³ The authors exploited the Doppler shift to make the momentum exchange (hence the force) velocity-dependent. It worked by directing a laser beam opposite to an atomic beam so the atoms could absorb light, and hence momentum $\hbar k$, very many times along their paths through the apparatus as shown in Fig. 2.^{3,4} Of course, excited-state atoms cannot absorb light efficiently from the laser that excited them, so between absorptions they must return to the ground state by spontaneous decay, accompanied by emission of fluorescent light. The spatial symmetry of the emitted fluorescence results in an average of zero net momentum transfer from many such fluorescence events. Thus the net force on the atoms is in the direction of the laser beam, and the maximum deceleration is limited by the spontaneous emission rate γ .

The maximum attainable deceleration is obtained for very high light intensities, and is limited because the atom must then divide its time equally between ground and excited states. High-intensity light can produce faster absorption, but it also causes equally fast stimulated emission; the combination produces neither deceleration nor cooling because the momentum transfer to the atom in emission is then in the opposite direction to what it was in absorption. The force is limited to $F = \hbar k \gamma_p$, and so the deceleration therefore saturates at a value $\bar{a}_{\max} = \hbar k \gamma / 2M$ [see Eq. (2)]. Since the maximum deceleration \bar{a}_{\max} is fixed by atomic parameters, it is straightforward to calculate the minimum stopping length L_{\min} and time t_{\min} for the rms velocity of atoms $\bar{v} = 2\sqrt{k_B T / M}$ at the chosen temperature. The result is $L_{\min} = \bar{v}^2 / 2\bar{a}_{\max}$ and $t_{\min} = \bar{v} / \bar{a}_{\max}$. In Table 1 are some of the parameters for slowing a few atomic species of interest from the peak of the thermal velocity distribution.

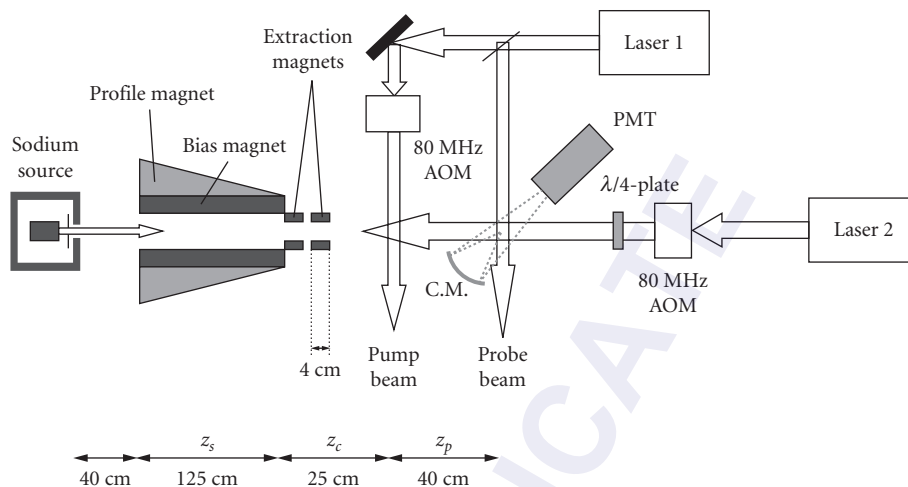


FIGURE 2 Schematic diagram of the apparatus for beam slowing. The tapered magnetic field is produced by layers of varying length on the solenoid, and the bias field is produced by full-length windings. The TOF aspect is implemented with Laser 1 whose beams optically pump the atoms between the hfs states.

TABLE 1 Parameters of Interest for Slowing Various Atoms

Atom	T_{oven} (K)	\bar{v} (m/s)	L_{min} (m)	t_{min} (ms)
H	1000	5000	0.012	0.005
He*	4	158	0.03	0.34
He*	650	2013	4.4	4.4
Li	1017	2051	1.15	1.12
Na	712	876	0.42	0.96
K	617	626	0.77	2.45
Rb	568	402	0.75	3.72
Cs	544	319	0.93	5.82

The stopping length L_{min} and time t_{min} are minimum values. The oven temperature T_{oven} that determines the peak velocity is chosen to give a vapor pressure of 1 Torr. Special cases are H at 1000 K and He in the metastable triplet state, for which two rows are shown: one for a 4-K source and another for the typical discharge temperature.

Maximizing the scattering rate γ_p requires $\delta = -\omega_D$ in Eq. (2). If δ is chosen for a particular atomic velocity in the beam, then as the atoms slow down, their changing Doppler shift will take them out of resonance. They will eventually cease deceleration after their Doppler shift has been decreased by a few times the power-broadened width $\gamma' = \gamma\sqrt{1+s_0}$, corresponding to Δv of a few times γ'/k . Although this Δv of a few m/s is considerably larger than the typical atomic recoil velocity $\hbar k/M$ of a few cm/s, it is still only a small fraction of the atoms' average thermal velocity \bar{v} , so that significant further cooling or deceleration cannot be accomplished.

In order to achieve deceleration that changes the atomic speeds by hundreds of m/s, it is necessary to maintain $(\delta + \omega_D) \ll \gamma$ by compensating such large changes of the Doppler shift. This can be done by changing ω_D , or δ via either ω_l or ω_a . The two most common methods for maintaining this resonance are sweeping the laser frequency ω_l along with the changing ω_D of the decelerating atoms,⁵⁻⁷ or by spatially varying the atomic resonance frequency with an inhomogeneous DC magnetic field to keep the decelerating atoms in resonance with the fixed frequency laser.^{1,3,8}

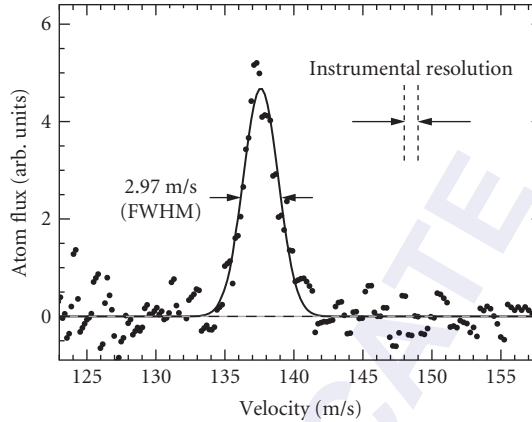


FIGURE 3 The velocity distribution measured with the time-of-flight (TOF) method. The experimental width of approximately $\frac{1}{6}(\gamma/k)$ is shown by the dashed vertical lines between the arrows. The Gaussian fit through the data yields a FWHM of 2.97 m/s. (Figure from Ref. 78.)

The use of a spatially varying magnetic field to tune the atomic levels along the beam path was the first method to succeed in slowing atoms.³ It works as long as the Zeeman shifts of the ground and excited states are different so that the resonant frequency is shifted. The field can be tailored to provide the appropriate Doppler shift along the moving atom's path. For uniform deceleration $a \equiv \eta a_{\max}$ from initial velocity v_0 , the appropriate field profile is $B(z) = B_0 \sqrt{1 - z/z_0}$, where $z_0 \equiv Mv_0^2 / \eta \hbar k \gamma$ is the length of the magnet, $B_0 = \hbar k v_0 / \mu'$, $\mu' \equiv (g_e M_e - g_g M_g) \mu_B$, subscripts g and e refer to ground and excited states, $g_{g,e}$ is the Landé g -factor, μ_B is the Bohr magneton, and $M_{g,e}$ is the magnetic quantum number. The design parameter $\eta < 1$ determines the length of the magnet z_0 . A solenoid that can produce such a spatially varying field has layers of decreasing lengths as shown schematically in Fig. 2. The technical problem of extracting the beam of slow atoms from the end of the solenoid can be simplified by reversing the field gradient and choosing a transition whose frequency decreases with increasing field.⁹

For alkali atoms such as Na, a time-of-flight (TOF) method can be used to measure the velocity distribution of atoms in the beam. It employs two additional beams labeled *pump* and *probe* from laser 1 as shown in Fig. 2. Because these beams cross the atomic beam at 90° , $\omega_D = -\mathbf{k} \cdot \mathbf{v} = 0$ and they excite atoms at all velocities. The pump beam is tuned to excite and empty a selected ground hyperfine state (hfs), and it transfers more than 98 percent of the population as the atoms pass through its 0.5 mm width. To measure the velocity distribution of atoms in the selected hfs, this pump laser beam is interrupted for a period $\Delta t = 10$ to $50 \mu\text{s}$ with an acoustic optical modulator (AOM). A pulse of atoms in the selected hfs passes the pump region and travels to the probe beam. The time dependence of the fluorescence induced by the probe laser, tuned to excite the selected hfs, gives the time of arrival, and this signal is readily converted to a velocity distribution. Figure 3 shows the measured velocity distribution of the atoms slowed by laser 2.

20.5 OPTICAL MOLASSES

Doppler Cooling

In Sec. 20.3 there was a discussion of the radiative force on atoms moving in a standing wave (counterpropagating laser beams). The slowing force is proportional to velocity for small enough velocities, resulting in viscous damping^{10,11} that gives this technique the name *optical molasses* (OM). By

using three intersecting orthogonal pairs of oppositely directed beams, the movement of atoms in the intersection region can be severely restricted in all three dimensions, and many atoms can thereby be collected and cooled in a small volume. OM has been demonstrated at several laboratories,¹² often with the use of low-cost diode lasers.¹³

It is straightforward to estimate the force on atoms in OM from Eq. (10). The discussion here is limited to the case where the light intensity is low enough so that stimulated emission is not important. In this low-intensity case the forces from the two light beams are simply added to give $\vec{F}_{\text{OM}} = \vec{F}_+ + \vec{F}_-$, where

$$\vec{F}_{\pm} = \pm \frac{\hbar k \gamma}{2} \frac{s_0}{1 + s_0 + [2(\delta \mp |\omega_D|)/\gamma]^2} \quad (25)$$

Then the sum of the two forces is

$$\vec{F}_{\text{OM}} \cong \frac{8\hbar k^2 \delta s_0 \vec{v}}{\gamma(1 + s_0 + (2\delta/\gamma)^2)} \equiv -\beta \vec{v} \quad (26)$$

where terms of order $(kv/\gamma)^4$ and higher have been neglected [see Eq. (20)].

These forces are plotted in Fig. 4. For $\delta < 0$, this force opposes the velocity and therefore viscously damps the atomic motion. \vec{F}_{OM} has maxima near $v \approx \pm \gamma'/2k$ and decreases rapidly for larger velocities.

If there were no other influence on the atomic motion, all atoms would quickly decelerate to $v = 0$, and the sample would reach $T = 0$, a clearly unphysical result. There is also some heating caused by the light beams that must be considered, and it derives from the discrete size of the momentum steps the atoms undergo with each emission or absorption as previously discussed for Brownian motion (see Sec. 20.3). Since the atomic momentum changes by $\hbar k$, their kinetic energy changes on the average by at least the recoil energy $E_r = \hbar^2 k^2 / 2M = \hbar \omega_r$. This means that the average frequency of each absorption is $\omega_{\text{abs}} = \omega_a + \omega_r$ and the average frequency of each emission is $\omega_{\text{emit}} = \omega_a + \omega_r$. Thus the light field loses an average energy of $\hbar(\omega_{\text{abs}} - \omega_{\text{emit}}) = 2\hbar \omega_r$ for each scattering. This loss occurs at a rate $2\gamma_p$ (two beams), and the energy is converted to atomic kinetic energy because the atoms recoil from each event. The atomic sample is thereby heated because these recoils are in random directions.

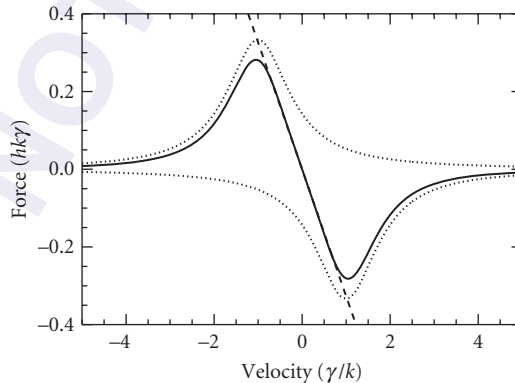


FIGURE 4 Velocity dependence of the optical damping forces for one-dimensional optical molasses. The two dotted traces show the force from each beam, and the solid curve is their sum. The straight line shows how this force mimics a pure damping force over a restricted velocity range. These are calculated for $s_0 = 2$ and $\delta = -\gamma$ so there is some power-broadening evident.

The competition between this heating with the damping force of Eq. (26) results in a nonzero kinetic energy in steady state where the rates of heating and cooling are equal. Equating the cooling rate, $\vec{F} \cdot \vec{v}$, to the heating rate, $4\hbar\omega_r\gamma_p$, the steady-state kinetic energy is $(\hbar\gamma/8)(2|\delta|/\gamma + \gamma/2|\delta|)$. This result is dependent on $|\delta|$, and it has a minimum at $2|\delta|/\gamma = 1$, whence $\delta = -\gamma/2$. The temperature found from the kinetic energy is then $T_D = \hbar\gamma/2k_B$, where T_D is called the *Doppler temperature* or the *Doppler cooling limit*. For ordinary atomic transitions, T_D is typically below 1 mK.

Another instructive way to determine T_D is to note that the average momentum transfer of many spontaneous emissions is zero, but the rms scatter of these about zero is finite. One can imagine these decays as causing a random walk in momentum space with step size $\hbar k$ and step frequency $2\gamma_p$, where the factor of 2 arises because of the two beams. The random walk results in diffusion in momentum space with diffusion coefficient $D_0 \equiv 2(\Delta p)^2 / \Delta t = 4\gamma_p(\hbar k)^2$ as discussed in Sec. 20.3. Then Brownian motion theory gives the steady-state temperature in terms of the damping coefficient β to be $k_B T = D_0/\beta$. This turns out to be $\hbar\gamma/2$ as above for the case $s_0 \ll 1$ when $\delta = -\gamma/2$. This remarkable result predicts that the final temperature of atoms in OM is independent of the optical wavelength, atomic mass, and laser intensity (as long as it is not too large).

Atomic Beam Collimation—One-Dimensional Optical Molasses

When an atomic beam crosses a one-dimensional OM as shown in Fig. 5, the transverse motion of the atoms is quickly damped while the longitudinal component is essentially unchanged. This transverse cooling of an atomic beam is an example of a method that can actually increase its brightness (atoms/sec-sr-cm²) because such active collimation uses dissipative forces to compress the phase space volume occupied by the atoms. By contrast, the usual realm of beam focusing or collimation techniques for light beams and most particle beams is restricted to selection by apertures or conservative forces that preserve the phase space density of atoms in the beam.

This velocity compression at low intensity in one dimension can be simply estimated for two-level atoms in 1D to be about $v_c/v_D = \sqrt{\gamma/\omega_r} = \sqrt{l/\epsilon}$. For Rb, $v_D = 12$ cm/s, $v_c = \gamma/k \approx 4.6$ m/s, $\omega_r \approx 2\pi \times 3.8$ kHz, and $l/\epsilon \approx 1600$. Including two transverse directions along with the longitudinal slowing and cooling just discussed, the decrease in phase space volume from the momentum contribution alone for laser cooling of a Rb atomic beam can exceed 10^6 .

Clearly optical techniques can create atomic beams enormously more times intense than ordinary thermal beams, and also many orders of magnitude brighter. Furthermore, this number could be increased by several orders of magnitude if the transverse cooling could produce temperatures

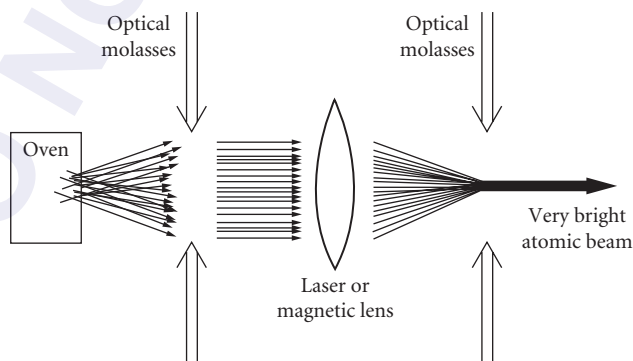


FIGURE 5 Scheme for optical brightening of an atomic beam. First the transverse velocity components of the atoms are damped out by an optical molasses, then the atoms are focused to a spot, and finally the atoms are recollimated in a second optical molasses. (Figure from Ref. 79.)

below the Doppler temperature. For atoms cooled to the recoil temperature $T_r = \hbar \omega_r / k_B$, where $\Delta p = \hbar k$ and $\Delta x = \lambda / \pi$, the brightness increase could be 10^{17} .

Experiments in Three-Dimensional Optical Molasses

Optical molasses experiments can also work in three dimensions at the intersection of three mutually orthogonal pairs of opposing laser beams (see Fig. 6). Even though atoms can be collected and cooled in the intersection region, it is important to stress again that this is *not* a trap. That is, atoms that wander away from the center experience no force directing them back. They are allowed to diffuse freely and even escape, as long as there is enough time for their very slow diffusive movement to allow them to reach the edge of the region of the intersection of the laser beams. Because the atomic velocities are randomized during the damping time $M/\beta = 2/\omega_r$, atoms execute a random walk with a step size of $2v_D/\omega_r = \lambda/\pi\sqrt{2\varepsilon} \cong \text{few } \mu\text{m}$. To diffuse a distance of 1 cm requires about 10^7 steps or about 30 s.^{14,15}

Three-dimensional OM was first observed in 1985.¹¹ Preliminary measurements of the average kinetic energy of the atoms were done by blinking off the laser beams for a fixed interval. Comparison of the brightness of the fluorescence before and after the turnoff was used to calculate the fraction of atoms that left the region while it was in the dark. The dependence of this fraction on the duration of the dark interval was used to estimate the velocity distribution and hence the temperature. The result was not inconsistent with the two-level atom theory previously described.

A few months later a more sensitive ballistic technique was devised at the National Institute of Standards and Technology (NIST) that showed the astounding result that the temperature of the atoms in OM was very much lower than T_D .¹⁶ These experiments also found that OM was less sensitive to perturbations and more tolerant of alignment errors than was predicted by the 1D, two-level atom theory. For example, if the intensities of the two counter-propagating laser beams forming an OM were unequal, then the force on atoms at rest would not vanish, but the force on atoms with some nonzero drift velocity *would* vanish. This drift velocity can be easily calculated by using Eq. (25) with unequal intensities s_{0+} and s_{0-} , and following the derivation of Eq. (26). Thus atoms would drift out of an OM, and the calculated rate would be much faster than observed by deliberately unbalancing the beams in the experiments.¹²

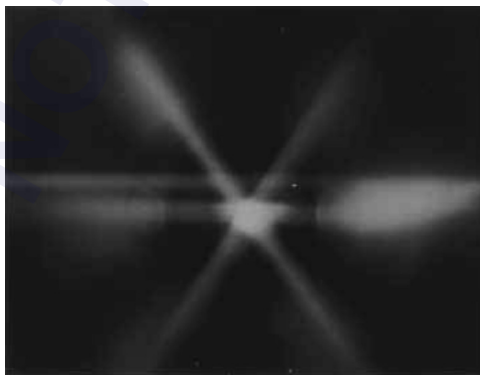


FIGURE 6 Photograph of optical molasses in Na taken under ordinary snapshot conditions in the lab at NIST. The upper horizontal streak is from the slowing laser while the three beams that cross at the center are on mutually orthogonal axes viewed from the (111) direction. Atoms in the optical molasses glow brightly at the center. (Figure from Ref. 19.)

It was an enormous surprise to observe that the ballistically measured temperature of the Na atoms was as much as 10 times *lower* than $T_D = 240 \mu\text{K}$,¹⁶ the temperature minimum calculated from the theory. This breaching of the Doppler limit forced the development of an entirely new picture of OM that accounts for the fact that in three dimensions, a two-level picture of atomic structure is inadequate. The multilevel structure of atomic states, and optical pumping among these sublevels, must be considered in the description of 3D OM, as discussed in the text that follows.

20.6 COOLING BELOW THE DOPPLER LIMIT

Introduction

In response to the surprising measurements of temperatures below T_D , two groups developed a model of laser cooling that could explain the lower temperatures.^{17,18} The key feature of this model that distinguishes it from the earlier picture was the inclusion of the multiplicity of sublevels that make up an atomic state (e.g., Zeeman and hfs). The dynamics of optically pumping atoms among these sublevels provides the new mechanism for producing the ultralow temperatures.¹⁹

The dominant feature of these models is the nonadiabatic response of moving atoms to the light field. Atoms at rest in a steady state have ground-state orientations caused by optical pumping processes that distribute the populations over the different ground-state sublevels. In the presence of polarization gradients, these orientations reflect the local light field. In the low-light-intensity regime, the orientation of stationary atoms is completely determined by the ground-state distribution: The optical coherences and the excited-state population follow the ground-state distribution adiabatically.

For atoms moving in a light field that varies in space, optical pumping acts to adjust the atomic orientation to the changing conditions of the light field. In a weak pumping process, the orientation of moving atoms always lags behind the orientation that would exist for stationary atoms. It is this phenomenon of nonadiabatic following that is the essential feature of the new cooling process.

Production of spatially dependent optical pumping processes can be achieved in several different ways. As an example, consider two counterpropagating laser beams that have orthogonal polarizations (as will be discussed shortly). The superposition of the two beams results in a light field having a polarization that varies on the wavelength scale along the direction of the laser beams. Laser cooling by such a light field is called *polarization gradient cooling*. In a three-dimensional optical molasses, the transverse wave character of light requires that the light field always has polarization gradients.

Linear \perp Linear Polarization Gradient Cooling

One of the most instructive models for discussion of sub-Doppler laser cooling was introduced in Ref. 17 and very well described in Ref. 19. If the polarizations of two counterpropagating laser beams are identical, the two beams interfere and produce a standing wave. When the two beams have orthogonal linear polarizations (same frequency ω_l) with their \hat{e} vectors perpendicular (e.g., \hat{x} and \hat{y}), the configuration is called *lin \perp lin* or *lin-perp-lin*. Then the total field is the sum of the two counterpropagating beams given by

$$\begin{aligned}\vec{E} &= E_0 \hat{x} \cos(\omega_l t - kz) + E_0 \hat{y} \cos(\omega_l t + kz) \\ &= E_0 [(\hat{x} + \hat{y}) \cos \omega_l t \cos kz + (\hat{x} - \hat{y}) \sin \omega_l t \sin kz]\end{aligned}\quad (27)$$

At the origin, where $z = 0$, this becomes

$$\vec{E} = E_0 (\hat{x} + \hat{y}) \cos \omega_l t \quad (28)$$

which corresponds to linearly polarized light at an angle $+\pi/4$ to the x -axis. The amplitude of this field is $\sqrt{2}E_0$. Similarly, for $z = \lambda/4$, where $kz = \pi/2$, the field is also linearly polarized but at an angle $-\pi/4$ to the x -axis.

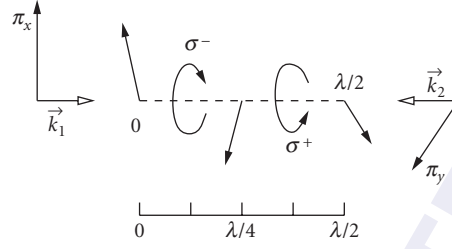


FIGURE 7 Polarization gradient field for the lin \perp lin configuration.

Between these two points, at $z = \lambda/8$, where $kz = \pi/4$, the total field is

$$\vec{E} = E_0 [\hat{x} \sin(\omega t + \pi/4) - \hat{y} \cos(\omega t + \pi/4)] \quad (29)$$

Since the \hat{x} and \hat{y} components have sine and cosine temporal dependence, they are $\pi/2$ out of phase, and so Eq. (29) represents circularly polarized light rotating about the z -axis in the negative sense. Similarly, at $z = 3\lambda/8$ where $kz = 3\pi/4$, the polarization is circular but in the positive sense. Thus in this lin \perp lin scheme the polarization cycles from linear to circular to orthogonal linear to opposite circular in the space of only half a wavelength of light, as shown in Fig. 7. It truly has a very strong polarization gradient.

Since the coupling of the different states of multilevel atoms to the light field depends on its polarization, atoms moving in a polarization gradient will be coupled differently at different positions, and this will have important consequences for laser cooling. For the $J_g = 1/2 \rightarrow J_e = 3/2$ transition (the simplest transition that shows sub-Doppler cooling), the optical pumping process in purely σ^+ light drives the ground-state population to the $M_g = +1/2$ sublevel. This optical pumping occurs because absorption always produces $\Delta M = +1$ transitions, whereas the subsequent spontaneous emission produces $\Delta M = \pm 1, 0$. Thus the average $\Delta M \geq 0$ for each scattering event. For σ^- -light the population is pumped toward the $M_g = -1/2$ sublevel. Thus atoms traveling through only a half wavelength in the light field need to readjust their population completely from $M_g = +1/2$ to $M_g = -1/2$ and back again.

The interaction between nearly resonant light and atoms not only drives transitions between atomic energy levels, but also shifts their energies. This light shift of the atomic energy levels plays a crucial role in this scheme of sub-Doppler cooling, and the changing polarization has a strong influence on the light shifts. In the low-intensity limit of two laser beams, each of intensity $s_0 I_s$, the light shifts ΔE_g of the ground magnetic substates are given by¹

$$\Delta E_g = \frac{\hbar \delta s_0 C_{ge}^2}{1 + (2\delta/\gamma)^2} \quad (30)$$

where C_{ge} is the Clebsch-Gordan coefficient that describes the coupling between the atom and the light field.

In the present case of orthogonal linear polarizations and $J = 1/2 \rightarrow 3/2$, the light shift for the magnetic substate $M_g = 1/2$ is three times larger than that of the $M_g = -1/2$ substate when the light field is completely σ^+ . On the other hand, when an atom moves to a place where the light field is σ^- , the shift of $M_g = -1/2$ is three times larger. So in this case the optical pumping previously discussed causes there to be a larger population in the state with the larger light shift. This is generally true for any transition J_g to $J_e = J_g + 1$. A schematic diagram showing the populations and light shifts for this particular case of negative detuning is shown in Fig. 8.

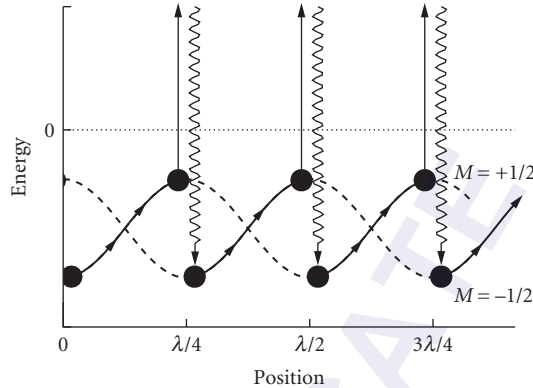


FIGURE 8 The spatial dependence of the light shifts of the ground-state sublevels of the $J = 1/2 \leftrightarrow 3/2$ transition for the case of the lin \perp lin polarization configuration. The arrows show the path followed by atoms being cooled in this arrangement. Atoms starting at $z = 0$ in the $M_g = +1/2$ sublevel must climb the potential hill as they approach the $z = \lambda/4$ point where the light becomes σ^- polarized, and there they are optically pumped to the $M_g = -1/2$ sublevel. Then they must begin climbing another hill toward the $z = \lambda/2$ point where the light is σ^+ polarized and they are optically pumped back to the $M_g = +1/2$ sublevel. The process repeats until the atomic kinetic energy is too small to climb the next hill. Each optical pumping event results in absorption of light at a lower frequency than emission, thus dissipating energy to the radiation field.

Origin of the Damping Force

To discuss the origin of the cooling process in this polarization gradient scheme, consider atoms with a velocity v at a position where the light is σ^+ -polarized, as shown at the lower left of Fig. 8. The light optically pumps such atoms to the strongly negative light-shifted $M_g = +1/2$ state. In moving through the light field, atoms must increase their potential energy (climb a hill) because the polarization of the light is changing and the state $M_g = 1/2$ becomes less strongly coupled to the light field. After traveling a distance $\lambda/4$, atoms arrive at a position where the light field is σ^- -polarized, and are optically pumped to $M_g = -1/2$, which is now lower than the $M_g = 1/2$ state. Again the moving atoms are at the bottom of a hill and start to climb. In climbing the hills, the kinetic energy is converted to potential energy, and in the optical pumping process, the potential energy is radiated away because the spontaneous emission is at a higher frequency than the absorption (see Fig. 8). Thus atoms seem to be always climbing hills and losing energy in the process. This process brings to mind a Greek myth, and is thus called *Sisyphus laser cooling*.

The cooling process just described is effective over a limited range of atomic velocities. The force is maximum for atoms that undergo one optical pumping process while traveling over a distance $\lambda/4$. Slower atoms will not reach the hilltop before the pumping process occurs, and faster atoms will already be descending the hill before being pumped toward the other sublevel. In both cases the energy loss is smaller and therefore the cooling process less efficient. Nevertheless, the damping constant β for this process is much larger than for Doppler cooling, and therefore the final steady-state temperature is lower.^{17,19}

In the experiments of Ref. 20, the temperature was measured in a 3D molasses under various configurations of the polarization. Temperatures were measured by a ballistic technique, where

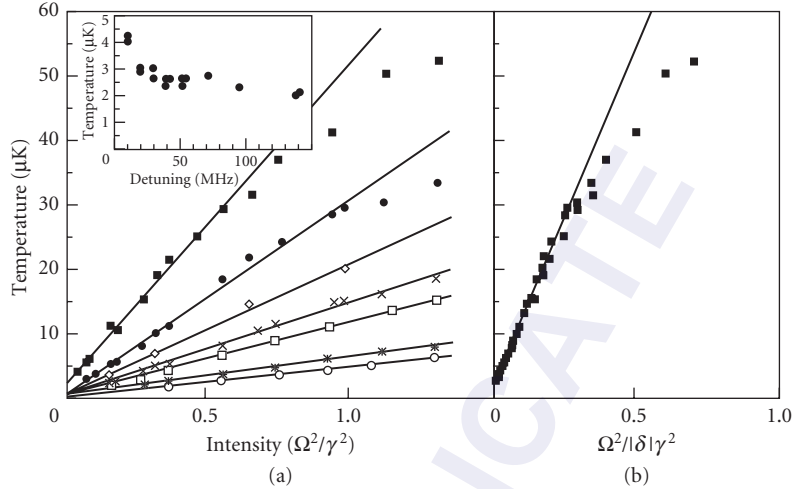


FIGURE 9 Temperature as a function of laser intensity and detuning for Cs atoms in an optical molasses from Ref. 20. (a) Temperature as a function of the detuning for various intensities. (b) Temperature as a function of the light shift. All the data points are on a universal straight line.

the flight time of the released atoms was measured as they fell through a probe a few cm below the molasses region. Results of their measurements are shown in Fig. 9a, where the measured temperature is plotted for different detunings as a function of the intensity. For each detuning, the data lie on a straight line through the origin. The lowest temperature obtained is 3 μK , which is a factor 40 below the Doppler temperature and a factor 15 above the recoil temperature of Cs. If the temperature is plotted as a function of the light shift (see Fig. 9b), all the data are on a single universal straight line.

The Limits of Laser Cooling

The lower limit to Doppler laser cooling of two-level atoms arises from the competition with heating. This cooling limit is described as a random walk in momentum space whose steps are of size $\hbar k$ and whose rate is the scattering rate, $\gamma_p = s_0 \gamma / 2$ for zero detuning and $s_0 \ll 1$. As long as the force can be accurately described as a damping force, then the Fokker-Planck equation is applicable, and the outcome is a lower limit to the temperature of laser cooling given by the Doppler temperature $k_B T_D \equiv \hbar \gamma / 2$.

The extension of this kind of thinking to the sub-Doppler processes described in Sec. 20.5 must be done with some care, because a naive application of the consequences of the Fokker-Planck equation would lead to an arbitrarily low final temperature. In the derivation of the Fokker-Planck equation it is explicitly assumed that each scattering event changes the atomic momentum p by an amount that is a small fraction p as embodied in Eq. (3), and this clearly fails when the velocity is reduced to the region of $v_r \equiv \hbar k / M$.

This limitation of the minimum steady-state value of the average kinetic energy to a few times $2E_r \equiv k_B T_r = M v_r^2$ is intuitively comforting for two reasons. First, one might expect that the last spontaneous emission in a cooling process would leave atoms with a residual momentum of the order of $\hbar k$, since there is no control over its direction. Thus the randomness associated with this would put a lower limit on such cooling of $v_{\min} \sim v_r$. Second, the polarization gradient cooling mechanism just described requires that atoms be localizable within the scale of $\sim \lambda / 2\pi$ in order to

be subject to only a single polarization in the spatially inhomogeneous light field. The uncertainty principle then requires that these atoms have a momentum spread of at least $\hbar k$.

The recoil limit discussed here has been surpassed by evaporative cooling of trapped atoms²¹ and two different optical cooling methods, neither of which can be based in simple notions. One of these uses optical pumping into a velocity-selective dark state and is described in Ref. 1. The other one uses carefully chosen, counterpropagating laser pulses to induce velocity-selective Raman transitions, and is called *Raman cooling*.²²

20.7 TRAPPING OF NEUTRAL ATOMS

Introduction

Although ion trapping, laser cooling of trapped ions, and trapped ion spectroscopy were known for many years,²³ it was only in 1985 that neutral atoms were first trapped.²⁴ Confinement of neutral atoms depends on the interaction between an inhomogeneous electromagnetic field and an atomic multipole moment. Unperturbed atoms do not have electric dipole moments because of their inversion symmetry, and therefore electric (e.g., optical) traps require induced dipole moments. This is often done with nearly resonant optical fields, thus producing the optical traps that will be discussed shortly. On the other hand, many atoms have ground- or metastable-state magnetic dipole moments that may be used for trapping them magnetically.

In order to confine any object, it is necessary to exchange kinetic for potential energy in the trapping field, and in neutral atom traps the potential energy must be stored as internal atomic energy. There are two immediate and extremely important consequences of this requirement. First, the atomic energy levels will necessarily shift as the atoms move in the trap, and these shifts will affect the precision of spectroscopic measurements, perhaps severely. Second, practical traps for ground-state neutral atoms are necessarily very shallow compared with thermal energy because the energy level shifts that result from convenient size fields are typically considerably smaller than $k_B T$ for $T = 1$ K. Neutral atom trapping therefore depends on substantial cooling of a thermal atomic sample, and is often connected with the cooling process.

The small depth of neutral atom traps also dictates stringent vacuum requirements, because an atom cannot remain trapped after a collision with a thermal energy background gas molecule. Since these atoms are vulnerable targets for thermal energy background gas, the mean free time between collisions must exceed the desired trapping time. The cross section for destructive collisions is quite large because even a gentle collision (i.e., large impact parameter) can impart enough energy to eject an atom from a trap. At pressure P sufficiently low to be of practical interest, the trapping time is $\sim(10^{-8}/P)$ s, where P is in Torr.

Magnetic Traps

An atom with a magnetic moment $\vec{\mu}$ can be confined by an inhomogeneous magnetic field because of an interaction between the moment and the field. This produces a force given by $\vec{F} = \nabla(\vec{\mu} \cdot \vec{B})$. Several different magnetic traps with varying geometries that exploit this force have been studied in some detail, and their general features have been presented.²⁵ The simplest magnetic trap is a quadrupole comprised of two identical coils carrying opposite currents (see Fig. 10) that has a single center where the field is zero. When the coils are separated by 1.25 times their radius, such a trap has equal depth in the radial (x - y plane) and longitudinal (z -axis) directions.²⁵ Its experimental simplicity makes it most attractive, both because of ease of construction and of optical access to the interior. Such a trap was used in the first neutral atom trapping experiments at NIST.

The magnitude of the field is zero at the center of this trap, and increases in all directions as $B = A\sqrt{\rho^2 + 4z^2}$, where $\rho^2 \equiv x^2 + y^2$, and the field gradient A is constant. The field gradient is fixed along any line through the origin, but has different values in different polar directions. Therefore the

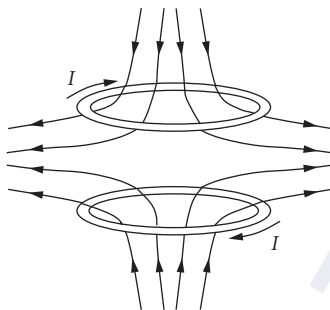


FIGURE 10 Schematic diagram of the coil configuration used in the quadrupole trap and the resultant magnetic field lines. Because the currents in the two coils are in opposite directions, there is a $|\vec{B}|=0$ point at the center.

force that confines the atoms in the trap is neither harmonic nor central, and angular momentum is not conserved. There are several motivations for studying the motion of atoms in a magnetic trap. Knowing their positions may be important for trapped atom spectroscopy. Moreover, simply studying the motion for its own sake has turned out to be an interesting problem because the distorted conical potential of the quadrupole trap does not have analytic solutions, and its bound states are not well known. For the two-coil quadrupole magnetic trap of Fig. 10, stable circular orbits can be found classically.¹ The fastest trappable atoms in circular orbits have $v_{\max} \sim 1$ m/s, so the orbital frequency becomes $\omega_r/2\pi \sim 20$ Hz. Because of the anharmonicity of the potential, the orbital frequencies depend on the orbit size, but in general, atoms in lower energy orbits have higher frequencies.

Because of the dependence of the trapping force on the angle between the field and the atomic moment the orientation of the magnetic moment with respect to the field must be preserved as the atoms move about in the trap. This requires velocities low enough to ensure that the interaction between the atomic moment $\vec{\mu}$ and the field \vec{B} is adiabatic, especially when the atom's path passes through a region where the field magnitude is small. This is especially critical at the low temperatures of the Bose condensation experiments. Therefore energy considerations that focus only on the trap depth are not sufficient to determine the stability of a neutral atom trap: orbit and/or quantum state calculations and their consequences must also be considered.

The condition for adiabatic motion can be written as $\omega_z \gg |dB/dt|/B$, where $\omega_z = \mu B/\hbar$ is the Larmor precession rate in the field. The orbital frequency for circular motion is $\omega_r = v/\rho$, and since $v/\rho = |dB/dt|/B$ for a uniform field gradient, the adiabaticity condition is $\omega_z \gg \omega_r$. For the two-coil quadrupole trap, the adiabaticity condition can be easily calculated.¹ A practical trap ($A \sim 1$ T/m) requires $\rho \gg 1$ μm as well as $v \gg 1$ cm/s. Note that violation of these conditions results in the onset of quantum dynamics for the motion (deBroglie wavelength \approx orbit size). Since the nonadiabatic region of the trap is so small (less than 10^{-18} m³ compared with typical sizes of ~ 2 cm corresponding to 10^{-5} m³), nearly all the orbits of most atoms are restricted to regions where they are adiabatic.

Modern techniques of laser and evaporative cooling have the capability to cool atoms to energies where their deBroglie wavelengths are on the micron scale. Such cold atoms may be readily confined to micron-size regions in magnetic traps with easily achievable field gradients, and in such cases, the notion of classical orbits is inappropriate. The motional dynamics must be described in terms of quantum mechanical variables and suitable wave functions. Furthermore, the distribution of atoms confined in various quantum states of motion in quadrupole as well as other magnetic traps is critical for interpreting the measurements on Bose condensates.

Studying the behavior of extremely slow (cold) atoms in the two-coil quadrupole trap begins with a heuristic quantization of the orbital angular momentum using $Mr^2\omega_r = n\hbar$ for circular orbits.¹ For velocities of optically cooled atoms of a few cm/s, $n \sim 10 - 100$. By contrast, evaporative cooling²¹ can produce velocities ~ 1 mm/s resulting in $n \sim 1$. It is readily found that $\omega_z = n\omega_r$ so that

the adiabatic condition is satisfied only for $n \gg 1$. The separation of the rapid precession from the slower orbital motion is reminiscent of the Born-Oppenheimer approximation for molecules, and three-dimensional quantum calculations have also been described.¹

Optical Traps

Optical trapping of neutral atoms by electrical interaction must proceed by inducing a dipole moment. For dipole optical traps, the oscillating electric field of a laser induces an oscillating atomic electric dipole moment that interacts with the laser field. If the laser field is spatially inhomogeneous, the interaction and associated energy level shift of the atoms varies in space and therefore produces a potential. When the laser frequency is tuned below atomic resonance ($\delta < 0$), the sign of the interaction is such that atoms are attracted to the maximum of laser field intensity, whereas if $\delta > 0$, the attraction is to the minimum of field intensity.

The simplest imaginable trap (see Fig. 11) consists of a single, strongly focused Gaussian laser beam^{26,27} whose intensity at the focus varies transversely with r as $I(r) = I_0 e^{-r^2/w_0^2}$, where w_0 is the beam waist size. Such a trap has a well-studied and important macroscopic classical analog in a phenomenon called *optical tweezers*.^{28–30} With the laser light tuned below resonance ($\delta < 0$), the ground-state light shift is everywhere negative, but is largest at the center of the Gaussian beam waist. Ground-state atoms therefore experience a force attracting them toward this center given by the gradient of the light shift. In the longitudinal direction there is also an attractive force that depends on the details of the focusing. Thus this trap produces an attractive force on atoms in three dimensions.

The first optical trap was demonstrated in Na with light detuned below the D-lines.²⁷ With 220 mW of dye laser light tuned about 650 GHz below the Na transition and focused to a $\sim 10 \mu\text{m}$ waist, the trap depth was about $15\hbar\gamma$ corresponding to 7 mK. Single-beam dipole force traps can be made with the light detuned by a significant fraction of its frequency from the atomic transition. Such a far-off-resonance trap (FORT) has been developed for Rb atoms using light detuned by nearly 10 percent to the red of the D_1 transition at $\lambda = 795 \text{ nm}$.³¹ Between 0.5 and 1 W of power was focused to a spot about $10 \mu\text{m}$ in size, resulting in a trap 6 mK deep where the light scattering rate was only a few hundred/s. The trap lifetime was more than half a second.

The dipole force for blue light repels atoms from the high-intensity region, and offers the advantage that trapped atoms will be confined where the perturbations of the light field are minimized.¹ On the other hand, it is not as easy to produce hollow light beams compared with Gaussian beams, and special optical techniques need to be employed.

In a standing wave the light intensity varies from zero at a node to a maximum at an antinode in a distance of $\lambda/4$. Since the light shift, and thus the optical potential, vary on this same scale, it is possible to confine atoms in wavelength-size regions of space. Of course, such tiny traps are usually very shallow, so loading them requires cooling to the μK regime. The momentum of such cold atoms is then so small that their deBroglie wavelengths are comparable to the optical wavelength, and hence to the trap size. In fact, the deBroglie wavelength equals the size of the optical traps ($\lambda/2$) when the momentum is $2\hbar k$, corresponding to a kinetic energy of a few μK . Thus the atomic motion in the trapping volume is not classical, but must be described quantum mechanically. Even atoms whose energy exceeds the trap depth must be described as quantum mechanical particles moving in a periodic potential that display energy band structure.³²

Atoms trapped in wavelength-sized spaces occupy vibrational levels similar to those of molecules. The optical spectrum can show Raman-like sidebands that result from transitions among the

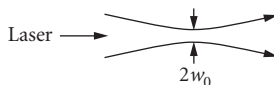


FIGURE 11 A single focused laser beam produces the simplest type of optical trap.

quantized vibrational levels^{33,34} as shown in Fig. 19. These quantum states of atomic motion can also be observed by spontaneous or stimulated emission.^{33,35} Considerably more detail about atoms in such optical lattices is to be found in Ref. 34.

Magneto-Optical Traps

The most widely used trap for neutral atoms is a hybrid, employing both optical and magnetic fields. The resultant *magneto-optical trap* (MOT) was first demonstrated in 1987.³⁶ The operation of an MOT depends on both inhomogeneous magnetic fields and radiative selection rules to exploit both optical pumping and the strong radiative force.^{1,36} The radiative interaction provides cooling that helps in loading the trap, and enables very easy operation. The MOT is a very robust trap that does not depend on precise balancing of the counter-propagating laser beams or on a very high degree of polarization. The magnetic field gradients are modest and can readily be achieved with simple, air-cooled coils. The trap is easy to construct because it can be operated with a room-temperature cell where alkali atoms are captured from the vapor. Furthermore, low-cost diode lasers can be used to produce the light appropriate for all the alkalis except Na, so the MOT has become one of the least expensive ways to produce atomic samples with temperatures below 1 mK. For these and other reasons it has become the workhorse of cold atom physics, and has also appeared in dozens of undergraduate laboratories.

Trapping in an MOT works by optical pumping of slowly moving atoms in a linearly inhomogeneous magnetic field $B = B(z) \equiv Az$, such as that formed by a magnetic quadrupole field. Atomic transitions with the simple scheme of $J_g = 0 \rightarrow J_e = 1$ have three Zeeman components in a magnetic field, excited by each of three polarizations, whose frequencies tune with field (and therefore with position) as shown in Fig. 12 for one dimension. Two counter-propagating laser beams of opposite circular polarization, each detuned below the zero field atomic resonance by δ , are incident as shown.

Because of the Zeeman shift, the excited state $M_e = +1$ is shifted up for $B > 0$, whereas the state with $M_e = -1$ is shifted down. At position z' in Fig. 12 the magnetic field therefore tunes the

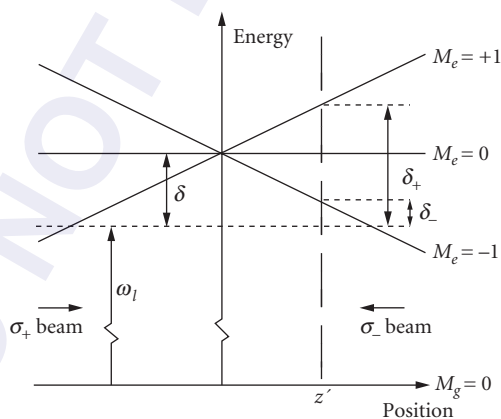


FIGURE 12 Arrangement for a magneto-optical trap (MOT) in 1D. The horizontal dashed line represents the laser frequency seen by an atom at rest in the center of the trap. Because of the Zeeman shifts of the atomic transition frequencies in the inhomogeneous magnetic field, atoms at $z = z'$ are closer to resonance with the σ^- laser beam than with the σ^+ beam, and are therefore driven toward the center of the trap.

$\Delta M = -1$ transition closer to resonance and the $\Delta M = +1$ transition further out of resonance. If the polarization of the laser beam incident from the right is chosen to be σ^- and correspondingly σ^+ for the other beam, then more light is scattered from the σ^- beam than from the σ^+ beam. Thus the atoms are driven toward the center of the trap, where the magnetic field is zero. On the other side of the center of the trap, the roles of the $M_e = \pm 1$ states are reversed, and now more light is scattered from the σ^+ beam, again driving the atoms toward the center.

The situation is analogous to the velocity damping in an optical molasses from the Doppler effect as previously discussed, but here the effect operates in position space, whereas for molasses it operates in velocity space. Since the laser light is detuned below the atomic resonance in both cases, compression and cooling of the atoms is obtained simultaneously in an MOT.

For a description of the motion of the atoms in an MOT, consider the radiative force in the low-intensity limit [see Eq. (10)]. The total force on the atoms is given by $\vec{F} = \vec{F}_+ + \vec{F}_-$, where

$$\vec{F}_{\pm} = \pm \frac{\hbar \vec{k} \gamma}{2} \frac{s_0}{1 + s_0 + (2\delta_{\pm}/\gamma)^2} \quad (31)$$

and the detuning δ_{\pm} for each laser beam is given by

$$\delta_{\pm} = \delta \mp \vec{k} \cdot \vec{v} \pm \mu' B / \hbar \quad (32)$$

Here $\mu' \equiv (g_e M_e - g_g M_g) \mu_B$ is the effective magnetic moment for the transition used. Note that the Doppler shift $\omega_D \equiv -\vec{k} \cdot \vec{v}$ and the Zeeman shift $\omega_z = \mu' B / \hbar$ both have opposite signs for opposite beams.

When both the Doppler and Zeeman shifts are small compared to the detuning δ , the denominator of the force can be expanded and the result becomes $\vec{F} = -\beta \vec{v} - \kappa \vec{r}$, where β is the damping coefficient. The spring constant κ arises from the similar dependence of \vec{F} on the Doppler and Zeeman shifts, and is given by $\kappa = \mu' A \beta / \hbar k$. This force leads to damped harmonic motion of the atoms, where the damping rate is given by $\Gamma_{\text{MOT}} = \beta / M$ and the oscillation frequency $\omega_{\text{MOT}} = \sqrt{\kappa / M}$. For magnetic field gradients $A \approx 10$ G/cm, the oscillation frequency is typically a few kHz, and this is much smaller than the damping rate that is typically a few hundred kHz. Thus the motion is overdamped, with a characteristic restoring time to the center of the trap of $2\Gamma_{\text{MOT}} / \omega_{\text{MOT}}^2 \sim$ several ms for typical values of the detuning and intensity of the lasers.

Since the MOT constants β and κ are proportional, the size of the atomic cloud can easily be deduced from the temperature of the sample. The equipartition of the energy of the system over the degrees of freedom requires that the velocity spread and the position spread are related by $k_B T = m v_{\text{rms}}^2 = \kappa z_{\text{rms}}^2$. For a temperature in the range of the Doppler temperature, the size of the MOT should be of the order of a few tenths of a mm, which is generally the case in experiments.

So far the discussion has been limited to the motion of atoms in one dimension. However, the MOT scheme can easily be extended to 3D by using six instead of two laser beams (see Fig. 13). Furthermore, even though very few atomic species have transitions as simple as $J_g = 0 \rightarrow J_e = 1$, the scheme works for any $J_g \rightarrow J_e = J_g + 1$ transition. Atoms that scatter mainly from the σ^+ laser beam will be optically pumped toward the $M_g = +J_g$ substate, which forms a closed system with the $M_e = +J_e$ substate.

The atomic density in an MOT cannot increase without limit as more atoms are added. The density is limited to $\sim 10^{11}/\text{cm}^3$ because the fluorescent light emitted by some trapped atoms is absorbed by others, and this diffusion of radiation presents a repulsive force between the atoms.^{37,38} Another limitation lies in the collisions between the atoms, and the collision rate for excited atoms is much larger than for ground-state atoms. Adding atoms to an MOT thus increases the density up to some point, but adding more atoms then expands the volume of the trapped sample.

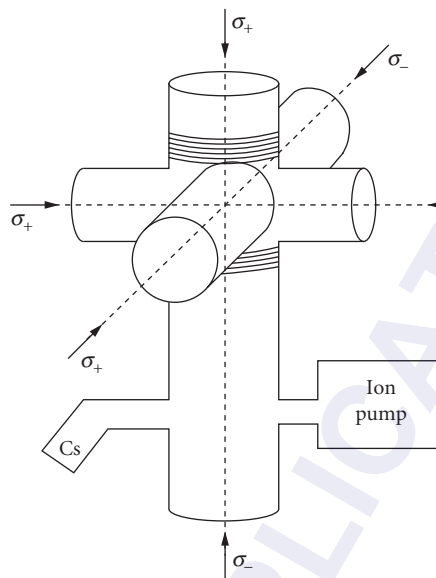


FIGURE 13 The schematic diagram of an MOT shows the coils and the directions of polarization of the six light beams. It has an axial symmetry and various rotational symmetries, so some exchanges would still result in a trap that works, but not all configurations are possible. Atoms are trapped from the background vapor of Cs that arises from a piece of solid Cs in one of the arms of the setup.

20.8 APPLICATIONS

Introduction

The techniques of laser cooling and trapping as described in the previous sections have been used to manipulate the positions and velocities of atoms with unprecedented variety and precision.¹ These techniques are currently used in the laboratories to design new, highly sensitive experiments that move experimental atomic physics research to completely new regimes. In this section only a few of these topics will be discussed. One of the most straightforward of these is the use of laser cooling to increase the brightness of atomic beams, which can subsequently be used for different types of experiments. Since laser cooling produces atoms at very low temperatures, the interaction between these atoms also takes place at such very low energies. The study of these interactions, called *ultra-cold collisions*, has been a very fruitful area of research in the last decade.

The atom-laser interaction not only produces a viscous environment for cooling the atoms down to very low velocities, but also provides a trapping field for the atoms. In the case of interfering laser beams, the size of such traps can be of the order of a wavelength, thus providing microscopic atomic traps with a periodic structure. These optical lattices described later in this section provide a versatile playground to study the effects of a periodic potential on the motion of atoms and thus simulate the physics of condensed matter. Another topic of considerable interest that will be discussed exists only because laser cooling has paved the way to the observation of Bose-Einstein condensation. This was predicted theoretically more than 80 years ago, but was observed in a dilute gas for the first time in 1996. Finally, the physics of dark states is also discussed in this section. These show a rich variety of effects caused by the coupling of internal and external coordinates of atoms.

Atomic Beam Brightening

In considering the utility of atomic beams for the purposes of lithography, collision studies, or a host of other applications, maximizing the beam intensity may not be the best option. Laser cooling can be used for increasing the phase space density, and this notion applies to both atomic traps and atomic beams. In the case of atomic beams, other quantities than phase space density have been defined as well, but these are not always consistently used. The geometrical solid angle occupied by atoms in a beam is $\Delta\Omega = (\Delta v_{\perp} \sqrt{v})^2$, where \bar{v} is some measure of the longitudinal velocity of atoms in the beam and Δv_{\perp} is a measure of the width of the transverse velocity distribution of the atoms. The total current or flux of the beam is Φ , and the flux density or intensity is $\Phi/\pi(\Delta x)^2$ where Δx is a measure of the beam's radius. Then the beam brightness or radiance R is given by $R = \Phi/\pi(\lambda x_{\perp})^2 \Delta\Omega$. Optical beams are often characterized by their frequency spread, and, because of the deBroglie relation $\lambda = h/p$, the appropriate analogy for atomic beams is the longitudinal velocity spread. Thus the spectral brightness or brilliance B is given by $B = R\bar{v}/\Delta v_z$. Note that both R and B have the same dimensions as flux density, and this is often a source of confusion. Finally, B is simply related to the 6D phase space density. Recently a summary of these beam properties has been presented in the context of phase space (see Fig. 14).

One of the first beam-brightening experiments was performed by Nellesen et al.^{39,40} where a thermal beam of Na was slowed with the chirp technique.¹ Then the slow atoms were deflected out of the main atomic beam and transversely cooled. In a later experiment⁴¹ this beam was fed into a two-dimensional MOT where the atoms were cooled and compressed in the transverse direction by an optical molasses of $\sigma^+ - \sigma^-$ polarized light. Another approach was used by Riis et al. who directed a slowed atomic beam into a hairpin-shaped coil that they called an *atomic funnel*.⁴² The wires of this coil generated a two-dimensional quadrupole field that was used as a two-dimensional MOT as described before.

These approaches yield intense beams when the number of atoms in the uncooled beam is already high. However, if the density in the beam is initially low, for example in the case of metastable noble gases or radioactive isotopes, one has to capture more atoms from the source in order

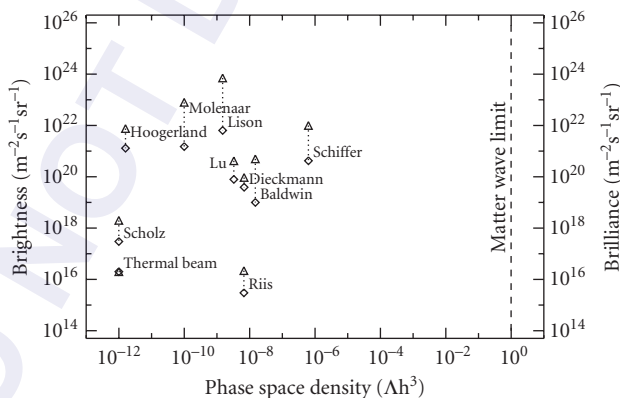


FIGURE 14 Plot of brightness (diamonds) and brilliance (triangles) versus phase space density for various atomic beams cited in the literature. The lower-left point is for a normal thermal beam, and the progression toward the top and right has been steady since the advent of laser cooling. The experimental results are from Riis et al.,⁴² Scholz et al.,⁸⁰ Hoogerland et al.,⁴⁴ Lu et al.,⁸¹ Baldwin et al.,⁸² Molenaar et al.,⁷⁸ Schiffer et al.,⁸³ Lison et al.,⁸⁴ and Dieckmann et al.⁸⁵ The quantum boundary for Bose-Einstein condensation, where the phase space density is unity, is shown by the dashed line on the right. (Figure adapted from Ref. 84.)

to obtain an intense beam. Aspect et al.⁴³ have used a quasi-standing wave of converging laser beams whose incidence angle varied from 87° to 90° to the atomic beam direction, so that a larger solid angle of the source could be captured. In this case they used a few mW of laser light over a distance of 75 mm. One of the most sophisticated approaches to this problem has been developed for metastable Ne by Hoogerland et al.⁴⁴ They used a three-stage process to provide a large solid angle capture range and produce a high brightness beam.

Applications to Atomic Clocks

Perhaps one of the most important practical applications of laser cooling is the improvement of atomic clocks. The limitation to both the accuracy and precision of such clocks is imposed by the thermal motion of the atoms, so a sample of laser-cooled atoms could provide a substantial improvement in clocks and in spectroscopic resolution.

The first experiments intended to provide slower atoms for better precision or clocks were attempts at an atomic fountain by Zacharias in the 1950s.^{1,45} This failed because collisions depleted the slow atom population, but the advent of laser cooling enabled an atomic fountain because the slow atoms far outnumber the faster ones. The first rf spectroscopy experiments in such a fountain using laser-cooled atoms were reported in 1989 and 1991,^{46,47} and soon after that some other laboratories also reported successes.

Some of the early best results were reported by Gibble and Chu.^{48,49} They used an MOT with laser beams 6 cm in diameter to capture Cs atoms from a vapor at room temperature. These atoms were launched upward at 2.5 m/s by varying the frequencies of the MOT lasers to form a moving optical molasses as described in Sec. 20.5, and subsequently cooled to below 3 μ K. The atoms were optically pumped into one hfs sublevel, then passed through a 9.2-GHz microwave cavity on their way up and again later on their way down. The number of atoms that were driven to change their hfs state by the microwaves was measured versus microwave frequency, and the signal showed the familiar Ramsey oscillations. See Chap. 11, “Coherent Optical Transients,” for a discussion of Ramsey fringes. The width of the central feature was 1.4 Hz and the S/N was over 50. Thus the ultimate precision was 1.5 mHz corresponding to $\delta\nu/\nu \cong 10^{-12}/\tau^{1/2}$, where τ is the number of seconds for averaging.

The ultimate limitation to the accuracy of this experiment as an atomic clock was collisions between Cs atoms in the beam. Because of the extremely low relative velocities of the atoms, the cross sections are very large (see the next subsection) and there is a measurable frequency shift.⁵⁰ By varying the density of Cs atoms in the fountain, the authors found frequency shifts of the order of a few mHz for atomic density of $10^9/\text{cm}^3$, depending on the magnetic sublevels connected by the microwaves. Extrapolation of the data to zero density provided a frequency determination of $\delta\nu/\nu \cong 4 \times 10^{-14}$. More recently the frequency shift has been used to determine a scattering length of $-400a_0$ ⁵¹ so that the expected frequency shift is 10^4 times larger than other limitations to the clock at an atomic density of $n = 10^9/\text{cm}^3$. Thus the authors suggest possible improvements to atomic timekeeping of a factor of 1000 in the near future. Even more promising are cold atom clocks in orbit (microgravity) where the interaction time can be very much longer than 1 s.⁵²

Ultracold Collisions

Laser-cooling techniques were developed in the early 1980s for a variety of reasons, such as high-resolution spectroscopy.¹ During the development of the techniques to cool and trap atoms, it became apparent that collisions between cold atoms in optical traps was one of the limiting factors in the achievement of high-density samples. Trap loss experiments revealed that the main loss mechanisms were caused by laser-induced collisions. Further cooling and compression could only be achieved by techniques not exploiting laser light, such as evaporative cooling in magnetic traps. Elastic collisions between atoms in the ground state are essential in that case for the rethermalization of the sample, whereas inelastic collisions lead to destruction of the sample. Knowledge about collision physics at these low energies is therefore essential for the development of high-density samples of atoms using either laser or evaporative cooling techniques.

Ground-state collisions play an important role in evaporative cooling. Such elastic collisions are necessary to obtain a thermalization of the gas after the trap depth has been lowered, and a large elastic cross section is essential to obtain a rapid thermalization. Inelastic collisions, on the other hand, can release enough energy to accelerate the atoms to energies too high to remain trapped. Ground-state collisions for evaporative cooling can be described by one parameter, the scattering length a . At temperatures below T_D , these collisions are in the s-wave scattering regime where only the phase shift δ_0 of the lowest partial wave $\ell=0$ is important. Moreover, for sufficiently low energies, such collisions are governed by the Wigner threshold laws where the phase shift δ_0 is inversely proportional to the wavevector k of the particle motion. Taking the limit for low energy gives the proportionality constant, defined as the scattering length $a = -\lim_{k \rightarrow 0} (\delta_0/k)$. The scattering length plays an important role not only in ultracold collisions, but also in the formation of Bose-Einstein condensates. In the Wigner threshold regime the cross section approaches a constant, $\sigma = 8\pi a^2$.⁵³

Although ground-state collisions are important for evaporative cooling and BEC, they do *not* provide a very versatile research field from a collision physics point of view. The situation is completely different for the excited-state collisions. For typical temperatures in optical traps, the velocity of the atoms is sufficiently low that atoms excited at long range by laser light decay before the collision takes place. Laser excitation for low-energy collisions has to take place during the collision. By tuning the laser frequency, the collision dynamics can be altered and information on the states formed in the molecular system can be obtained. This is the basis of the new technique of photo-associative spectroscopy, which for the first time has identified purely long-range states in diatomic molecules.^{1,54}

For atoms colliding in laser light closely tuned to the S-P transition, the potential is a C_3/R^3 dipole-dipole interaction when one of the atoms is excited. Absorption takes place at the Condon point R_C given by $\hbar\delta = -C_3/R_C^3$ or $R_C = (C_3/\hbar|\delta|)^{1/3}$. Note that the light has to be tuned below resonance, which is mostly the case for laser cooling. The Condon point for laser light detuned a few γ below resonance is typical 1000 to 2000 a_0 .

Once the molecular complex becomes excited, it can evolve to smaller internuclear distances before emission takes place. Two particular cases are important for trap loss: (1) the emission of the molecular complex takes place at much smaller internuclear distance, and the energy gained between absorption and emission of the photon is converted into kinetic energy, or (2) the complex undergoes a transition to another state and the potential energy difference between the two states is converted into kinetic energy. In both cases the energy gain can be sufficient to eject one or both atoms out of the trap. In the case of the alkalis, the second reaction can take place because of the different fine-structure states and the reaction is denoted as a fine-structure changing collision. The first reaction is referred to as *radiative escape*.

Trap loss collisions in MOTs have been studied to great extent, but results of these studies have to be considered with care. In most cases, trap loss is studied by changing either the frequency or the intensity of the trapping laser, which also changes the conditions of the trap. The collision rate is not only changed because of a change in the collision cross section, but also because of changes in both the density and temperature of the atoms in the trap. Since these parameters cannot be determined with high accuracy in a high-density trap, where effects like radiation trapping can play an important role, obtaining accurate results this way is very difficult.

The first description of such processes was given by Gallagher and Pritchard.⁵⁵ In their semiclassical model (the GP-model), the laser light is assumed to be weak enough that the excitation rate can be described by a quasi-static excitation probability. Atoms in the excited state are accelerated toward one another by the C_3/R^3 potential. In order to calculate the survival of the atoms in the excited state, the elapsed time between excitation and arrival is calculated. The total number of collisions is then given by the number of atoms at a certain distance, the fraction of atoms in the excited state, and the survival rate, integrated over all distances. For small detunings, corresponding to large internuclear distances, the excitation rate is appreciable over a very large range of internuclear distances. However the excitation occurs at large internuclear distances, so the survival rate of the excited atoms is small. For large detunings the excitation is located in a small region at small internuclear distances, so the total excitation rate is small, but the survival rate is large. As a result of this competition, the collision rate peaks at intermediate detunings.

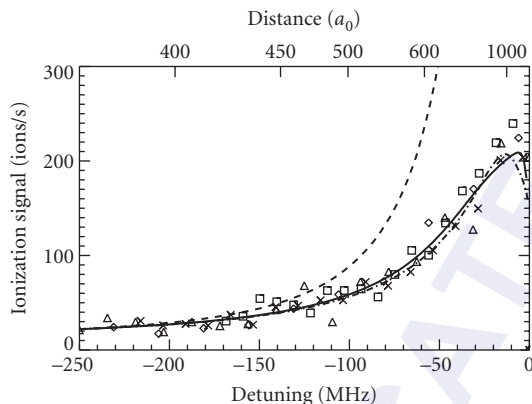


FIGURE 15 The frequency dependence for the associative ionization rate of cold He^* collisions. The experimental results (symbols) are compared with the semiclassical model (solid line), JV-model (dashed line), and modified JV-model (dashed-dotted line). The axis on top of the plot shows the Condon point, where the excitation takes place.

Another description of optical collisions is given by Julienne and Vigue.⁵⁶ Their description of optical collisions (JV model) is quantum mechanical for the collision process, where they make a partial wave expansion of the incoming wavefunction. The authors describe the excitation process in the same way as it was done in the GP model. Thus the excitation is localized around the Condon point with a probability given by the quasi-static Lorentz formula.

In still another approach, a completely semiclassical description of optical collisions has been given by Mastwijk et al.⁵⁷ These authors start from the GP model, but make several important modifications. First, the Lorentz formula is replaced by the Landau-Zener formula. Second, the authors consider the motion of the atoms in the collision plane. At the Condon point, where the excitation takes place, the trajectory of the atom in the excited state is calculated by integration of the equation of motion. The results for their model are shown in Fig. 15, and are compared with experiment and the JV model. The agreement between the theory and experiment is rather good. For the JV model two curves are shown. The first curve shows the situation for the original JV model. The second curve shows the result of a modified JV model, where the quasi-static excitation rate is replaced by the Landau-Zener formula. The large discrepancies between the results for these two models indicates that it is important to use the correct model for the excitation. The agreement between the modified JV model and the semiclassical model is good, indicating that the dynamics of optical collisions can be described correctly quantum mechanically or semiclassically. Since the number of partial waves in the case of He^* is in the order of 10, this is to be expected.

The previous description of optical collisions applies to the situation that the quasi-molecule can be excited for each frequency of the laser light. However, the quasi-molecule has well-defined vibrational and rotational states and the excitation frequency has to match the transition frequency between the ground and excited rovibrational states. Far from the dissociation limit, the rovibrational states are well-resolved and many resonances are observed. This has been the basis of the method of photo-associative spectroscopy (PAS) for alkali-metal atoms, where detailed information on molecular states of alkali dimers have been obtained recently. Here *photo-association* refers to the process where a photon is absorbed to transfer the system from the ground to the excited state where the two atoms are bound by their mutual attraction.

The process of PAS is depicted graphically in Fig. 16. When two atoms collide in the ground state, they can be excited at a certain internuclear distance to the excited molecular state and the two atoms may remain bound after the excitation and form a molecule. This transient molecule lives as long as the system remains excited. The number of rotational states that can contribute to the spectrum is

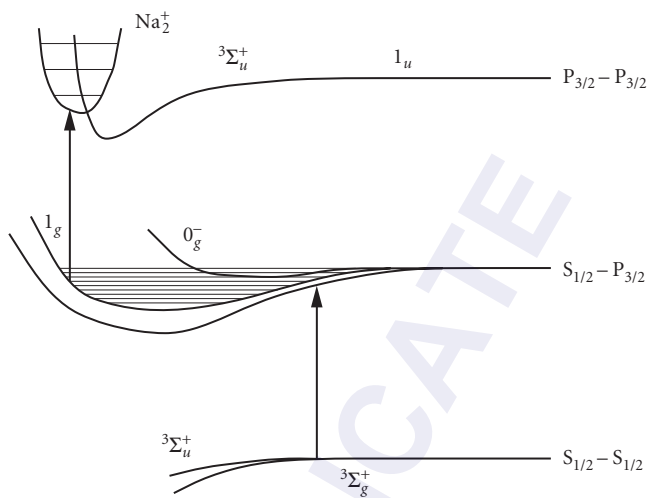


FIGURE 16 Photoassociation spectroscopy of Na. By tuning the laser below atomic resonance, molecular systems can be excited to the first excited state in which they are bound. By absorption of a second photon the system can be ionized, providing a high detection efficiency.

small for low temperature. The resolution is limited only by the linewidth of the transition, which is comparable to the natural linewidth of the atomic transition. With PAS, molecular states can be detected with a resolution of ≈ 10 MHz, which is many orders of magnitude better than traditional molecular spectroscopy. The formation of the molecules is probed by absorption of a second photon of the same color, which can ionize the molecule.

PAS has also been discussed in the literature as a technique to produce cold molecules. The methods discussed employ a double resonance technique, where the first color is used to create a well-defined rovibrational state of the molecule and a second color causes stimulated emission of the system to a well-defined vibrational level in the ground state. Although such a technique has not yet been shown to work experimentally, cold molecules have been produced in PAS recently using a simpler method.⁵⁸ The 0_g^- state in Cs₂ has a double-well structure, where the top of the barrier is accidentally close to the asymptotic limit. Thus atoms created in the outer well by PAS can tunnel through the barrier to the inner well, where there is a large overlap of the wavefunction with the vibrational levels in the ground state. These molecules are then stabilized against spontaneous decay and can be observed. The temperature of the cold molecules has been detected and is close to the temperature of the atoms. This technique and similar techniques will be very important for the production and study of cold molecules.

Optical Lattices

In 1968, Letokhov⁵⁹ suggested that it is possible to confine atoms in the wavelength-size regions of a standing wave by means of the dipole force that arises from the light shift. This was first accomplished in 1987 in one dimension with an atomic beam traversing an intense standing wave.⁶⁰ Since then, the study of atoms confined in wavelength-size potential wells has become an important topic in optical control of atomic motion because it opens up configurations previously accessible only in condensed matter physics using crystals.

The basic ideas of the quantum mechanical motion of particles in a periodic potential were laid out in the 1930s with the Kronig-Penney model and Bloch's theorem, and optical lattices offer important opportunities for their study. For example, these lattices can be made essentially free of

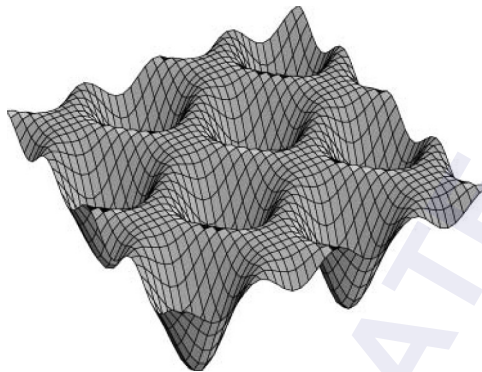


FIGURE 17 The “egg-crate” potential of an optical lattice shown in two dimensions. The potential wells are separated by $\lambda/2$.

defects with only moderate care in spatially filtering the laser beams to assure a single transverse mode structure. Furthermore, the shape of the potential is exactly known, and doesn’t depend on the effect of the crystal field or the ionic energy level scheme. Finally, the laser parameters can be varied to modify the depth of the potential wells without changing the lattice vectors, and the lattice vectors can be changed independently by redirecting the laser beams. The simplest optical lattice to consider is a one-dimensional pair of counter-propagating beams of the same polarization, as was used in the first experiment.⁶⁰

Because of the transverse nature of light, any mixture of beams with different \vec{k} -vectors necessarily produces a spatially periodic, inhomogeneous light field. The importance of the “egg-crate” array of potential wells arises because the associated atomic light shifts can easily be comparable to the very low average atomic kinetic energy of laser-cooled atoms. A typical example projected against two dimensions is shown in Fig. 17.

The name *optical lattice* is used rather than *optical crystal* because the filling fraction of the lattice sites is typically only a few percent (as of 1999). The limit arises because the loading of atoms into the lattice is typically done from a sample of trapped and cooled atoms, such as an MOT for atom collection, followed by an optical molasses for laser cooling. The atomic density in such experiments is limited to a few times $10^{11}/\text{cm}^3$ by collisions and multiple light scattering. Since the density of lattice sites of size $\lambda/2$ is a few times $10^{13}/\text{cm}^3$, the filling fraction is necessarily small.

At first thought it would seem that a rectangular 2D or 3D optical lattice could be readily constructed from two or three mutually perpendicular standing waves.^{61,62} However, a sub-wavelength movement of a mirror caused by a small vibration could change the relative phase of the standing waves. In 1993 a very clever scheme was described.⁶³ It was realized that an N -dimensional lattice could be created by only $n + 1$ traveling waves rather than $2n$. Instead of producing optical wells in 2D with four beams (two standing waves), these authors used only three. The k -vectors of the coplanar beams were separated by $2\pi/3$, and they were all linearly polarized in their common plane (not parallel to one another). The same immunity to vibrations was established for a 3D optical lattice by using only four beams arranged in a quasi-tetrahedral configuration. The three linearly polarized beams of the 2D arrangement just described were directed out of the plane toward a common vertex, and a fourth circularly polarized beam was added. All four beams were polarized in the same plane.⁶³ The authors showed that such a configuration produced the desired potential wells in 3D.

The NIST group studied atoms loaded into an optical lattice using Bragg diffraction of laser light from the spatially ordered array.⁶⁴ They cut off the laser beams that formed the lattice, and before the atoms had time to move away from their positions, they pulsed on a probe laser beam at the Bragg angle appropriate for one of the sets of lattice planes. The Bragg diffraction not only enhanced the reflection of the probe beam by a factor of 10^5 , but by varying the time between the shut-off of the lattice and turn-on of the probe, they could measure the “temperature” of the atoms

in the lattice. The reduction of the amplitude of the Bragg scattered beam with time provided some measure of the diffusion of the atoms away from the lattice sites, much like the Debye-Waller factor in X-ray diffraction.

Laser cooling has brought the study of the motion of atoms into an entirely new domain where the quantum mechanical nature of their center-of-mass motion must be considered.¹ Such exotic behavior for the motion of whole atoms, as opposed to electrons in the atoms, has not been considered before the advent of laser cooling simply because it is too far out of the range of ordinary experiments. A series of experiments in the early 1990s provided dramatic evidence for these new quantum states of motion of neutral atoms, and led to the debut of de Broglie wave atom optics.

The limits of laser cooling discussed in Sec. 20.6 suggest that atomic momenta can be reduced to a “few” times $\hbar k$. This means that their de Broglie wavelengths are equal to the optical wavelengths divided by a “few.” If the depth of the optical potential wells is high enough to contain such very slow atoms, then their motion in potential wells of size $\lambda/2$ must be described quantum mechanically, since they are confined to a space of size comparable to their de Broglie wavelengths. Thus they do not oscillate in the sinusoidal wells as classical localizable particles, but instead occupy discrete, quantum-mechanical bound states, as shown in the lower part of Fig. 18.

The group at NIST also developed a new method that superposed a weak probe beam of light directly from the laser upon some of the fluorescent light from the atoms in a 3D optical molasses, and directed the light from these combined sources onto a fast photodetector.⁶⁵ The resulting beat signal carried information about the Doppler shifts of the atoms in the optical lattices.³⁴ These Doppler shifts were expected to be in the sub-MHz range for atoms with the previously measured $50 \mu\text{K}$ temperatures. The observed features confirmed the quantum nature of the motion of atoms in the wavelength-size potential wells (see Fig. 19).¹⁶

In the 1930s Bloch realized that applying a uniform force to a particle in a periodic potential would not accelerate it beyond a certain speed, but instead would result in Bragg reflection when its

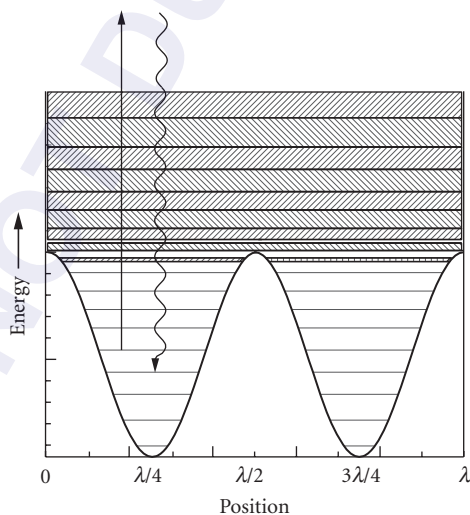


FIGURE 18 Energy levels of atoms moving in the periodic potential of the light shift in a standing wave. There are discrete bound states deep in the wells that broaden at higher energy, and become bands separated by forbidden energies above the tops of the wells. Under conditions appropriate to laser cooling, optical pumping among these states favors populating the lowest ones as indicated schematically by the arrows.

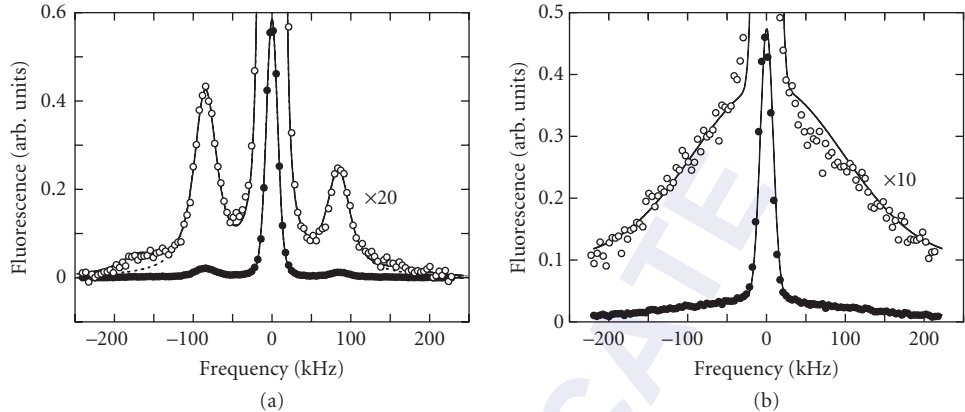


FIGURE 19 (a) Fluorescence spectrum in a 1D lin \perp lin optical molasses. Atoms are first captured and cooled in an MOT, then the MOT light beams are switched off leaving a pair of lin \perp lin beams. Then the measurements are made with $\delta = -4\gamma$ at low intensity. (b) Same as (a) except the 1D molasses is $\sigma^+ - \sigma^-$ which has no spatially dependent light shift and hence no vibrational motion. (Figure from Ref. 34.)

de Broglie wavelength became equal to the lattice period. Thus an electric field applied to a conductor could not accelerate electrons to a speed faster than that corresponding to the edge of a Brillouin zone, and that at longer times the particles would execute oscillatory motion. Ever since then, experimentalists have tried to observe these Bloch oscillations in increasingly pure and/or defect-free crystals.

Atoms moving in optical lattices are ideally suited for such an experiment, as was beautifully demonstrated in 1996.⁶⁶ The authors loaded a one-dimensional lattice with atoms from a 3D molasses, further narrowed the velocity distribution, and then instead of applying a constant force, simply changed the frequency of one of the beams of the 1D lattice with respect to the other in a controlled way, thereby creating an accelerating lattice. Seen from the atomic reference frame, this was the equivalent of a constant force trying to accelerate them. After a variable time t_a the 1D lattice beams were shut off and the measured atomic velocity distribution showed beautiful Bloch oscillations as a function of t_a . The centroid of the very narrow velocity distribution was seen to shift in velocity space at a constant rate until it reached $v_r = \hbar k/M$, and then it vanished and reappeared at $-v_r$ as shown in Fig. 20. The shape of the “dispersion curve” allowed measurement of the “effective mass” of the atoms bound in the lattice.

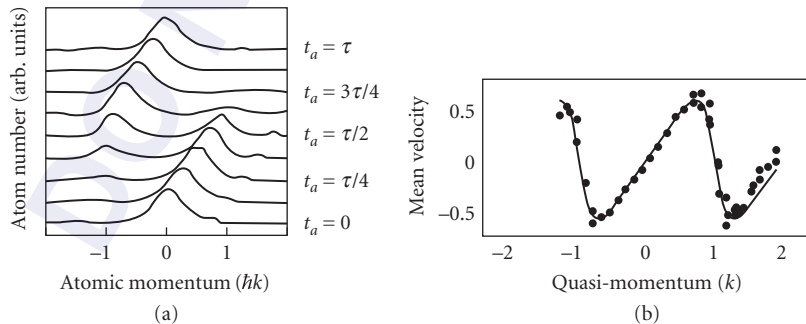


FIGURE 20 Plot of the measured velocity distribution versus time in the accelerated 1D lattice. The atoms accelerate only to the edge of the Brillouin zone where the velocity is $+v_r$, and then the velocity distribution appears at $-v_r$. (Figure from Ref. 66.)

Bose-Einstein Condensation

In 1924 S. Bose found the correct way to evaluate the distribution of identical entities, such as Planck's radiation quanta, that allowed him to calculate the Planck spectrum using the methods of statistical mechanics. Within a year Einstein had seized upon this idea, and generalized it to identical particles with discrete energies. This distribution is

$$N(E) = \frac{1}{e^{\beta(E-\mu)} - 1} \quad (33)$$

where $\beta \equiv 1/k_B T$ and μ is the chemical potential that vanishes for photons: Eq. (33) with $\mu = 0$ is exactly the Planck distribution. Einstein observed that this distribution has the peculiar property that for sufficiently low average energy (i.e., low temperature), the total energy could be minimized by having a discontinuity in the distribution for the population of the lowest allowed state.

The condition for this Bose-Einstein condensation (BEC) in a gas can be expressed in terms of the de Broglie wavelength λ_{dB} associated with the thermal motion of the atoms as $n\lambda_{\text{dB}}^3 \geq 2.612\dots$, where n is the spatial density of the atoms. In essence, this means that the atomic wave functions must overlap one another.

The most familiar elementary textbook description of BEC focuses on noninteracting particles. However, particles *do* interact, and the lowest order approximation that is widely used to account for the interaction takes the form of a mean-field repulsive force. It is inserted into the Hamiltonian for the motion of each atom in the trap (*n.b.*, not for the internal structure of the atom) as a term V_{int} proportional to the local density of atoms. Since this local density is itself $|\Psi|^2$, it makes the Schrödinger equation for the atomic motion nonlinear, and the result bears the name *Gross-Pitaevski equation*. For N atoms in the condensate it is written

$$\left[-\frac{\hbar^2}{2M} \nabla_{\vec{R}}^2 + V_{\text{trap}}(\vec{R}) + NV_{\text{int}} |\Psi(\vec{R})|^2 \right] \Psi(\vec{R}) = E_N \Psi(\vec{R}) \quad (34)$$

where \vec{R} is the coordinate of the atom in the trap, $V_{\text{trap}}(\vec{R})$ is the potential associated with the trap that confines the atoms in the BEC, and $V_{\text{int}} \equiv 4\pi\hbar^2 a/M$ is the coefficient associated with strength of the mean field interaction between the atoms. Here a is the scattering length, and M is the atomic mass.

For $a > 0$ the interaction is repulsive so that a BEC would tend to disperse. This is manifest for a BEC confined in a harmonic trap by having its wavefunction somewhat more spread out and flatter than a Gaussian. By contrast, for $a < 0$ the interaction is attractive and the BEC eventually collapses. However, it has been shown that there is metastability for a sufficiently small number of particles with $a < 0$ in a harmonic trap, and that a BEC can be observed in vapors of atoms with such negative scattering length as ^7Li .⁶⁷⁻⁶⁹ This was initially somewhat controversial.

Solutions to this highly nonlinear Eq. (34), and the ramifications of those solutions, form a major part of the theoretical research into BEC. Note that the condensate atoms all have exactly the same wave function, which means that adding atoms to the condensate does not increase its volume, just like the increase of atoms to the liquid phase of a liquid-gas mixture makes only an infinitesimal volume increase of the sample. The consequences of this predicted condensation are indeed profound. For example, in a harmonic trap, the lowest state's wavefunction is a Gaussian. With so many atoms having *exactly* the same wave function they form a new state of matter, unlike anything in the familiar experience.

Achieving the conditions required for BEC in a low-density atomic vapor requires a long and difficult series of cooling steps. First, note that an atomic sample cooled to the recoil limit T_r would need to have a density of a few times 10^{13} atoms/cm³ in order to satisfy BEC. However, atoms cannot be optically cooled at this density because the resulting vapor would have an

absorption length for on-resonance radiation approximately equal to the optical wavelength. Furthermore, collisions between ground- and excited-state atoms have such a large cross section that at this density the optical cooling would be extremely ineffective. In fact, the practical upper limit to the atomic density for laser cooling in a 3D optical molasses (see Sec. 20.6) or MOT (see Sec. 20.7) corresponds to $n \sim 10^{10}$ atoms/cm³. Thus it is clear that the final stage of cooling toward a BEC must be done in the dark. The process typically begins with an MOT for efficient capture of atoms from a slowed beam or from the low-velocity tail of a Maxwell-Boltzmann distribution of atoms at room temperature. Then a polarization gradient optical molasses stage is initiated that cools the atomic sample from the mK temperatures of the MOT to a few times T_r . For the final cooling stage, the cold atoms are confined in the dark in a purely magnetic trap and a forced evaporative cooling process is used to cool.¹

The observation of BEC in trapped alkali atoms in 1995 has been the largest impetus to research in this exciting field. As of this writing (1999), the only atoms that have been condensed are Rb,⁷⁰ Na,⁷¹ Li,⁷² and H.⁷³ The case of Cs is special because, although BEC is certainly possible, the presence of a near-zero energy resonance severely hampers its evaporative cooling rate.

The first observations of BEC were in Rb,⁷⁰ Li,⁷² and Na,⁷¹ and the observation was done using ballistic techniques. The results from one of the first experiments are shown in Fig. 21. The three panels show the spatial distribution of atoms some time after release from the trap. From the ballistic parameters, the size of the BEC sample as well as its shape and the velocity distribution of its atoms could be inferred. For temperatures too high for BEC, the velocity distribution is Gaussian but asymmetrical. For temperatures below the transition to BEC, the distribution is also not symmetrical, but now shows the distinct peak of a disproportionate number of very slow atoms corresponding to the ground state of the trap from which they were released. As the temperature is lowered further, the number of atoms in the narrow feature increases very rapidly, a sure signature that this is truly a BEC and not just very efficient cooling.

The study of this “new form” of matter has spawned innumerable subtopics and has attracted enormous interest. Both theorists and experimentalists are addressing the questions of its behavior

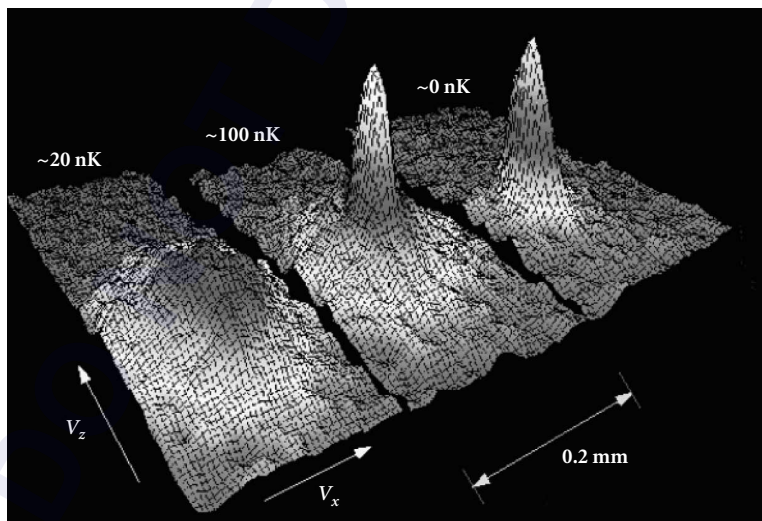


FIGURE 21 Three panels showing the spatial distribution of atoms after release from the magnetostatic trap following various degrees of evaporative cooling. In the first one, the atoms were cooled to just before the condition for BEC was met, in the second one, to just after this condition, and in the third one to the lowest accessible temperature consistent with leaving some atoms still in the trap. (Figure taken from the JILA Web page.)

in terms of rigidity, acoustics, coherence, and a host of other properties. Extraction of a coherent beam of atoms from a BEC has been labeled an “atom laser” and will surely open the way for new developments in atom optics.¹

Dark States

The BEC discussed in the previous subsection is an example of the importance of quantum effects on atomic motion. It occurs when the atomic de Broglie wavelength λ_{dB} and the interatomic distances are comparable. Other fascinating quantum effects occur when atoms are in the light and λ_{dB} is comparable to the optical wavelength. Some topics connected with optical lattices have already been discussed, and the *dark states* described here are another important example. These are atomic states that cannot be excited by the light field.

The quantum description of atomic motion requires that the energy of such motion be included in the Hamiltonian. The total Hamiltonian for atoms moving in a light field would then be given by

$$\mathcal{H} = \mathcal{H}_{\text{atom}} + \mathcal{H}_{\text{rad}} + \mathcal{H}_{\text{int}} + \mathcal{H}_{\text{kin}} \quad (35)$$

where $\mathcal{H}_{\text{atom}}$ describes the motion of the atomic electrons and gives the internal atomic energy levels, \mathcal{H}_{rad} is the energy of the radiation field and is of no concern here because the field is not quantized, \mathcal{H}_{int} describes the excitation of atoms by the light field and the concomitant light shifts, and \mathcal{H}_{kin} is the kinetic energy E_k of the motion of the atoms' center of mass. This Hamiltonian has eigenstates of not only the internal energy levels and the atom-laser interaction that connects them, but also of the kinetic energy operator $\mathcal{H}_{\text{kin}} \equiv \mathcal{P}^2/2M$. These eigenstates will therefore be labeled by quantum numbers of the atomic states as well as the center of mass momentum p . For example, an atom in the ground state, $|g; p\rangle$, has energy $E_g + p^2/2M$ which can take on a continuous range of values.

To see how the quantization of the motion of a two-level atom in a monochromatic field allows the existence of a velocity-selective dark state, consider the states of a two-level atom with single internal ground and excited levels, $|g; p\rangle$ and $|e; p'\rangle$. Two ground eigenstates $|g; p\rangle$ and $|g; p''\rangle$ are generally not coupled to one another by an optical field except in certain cases. For example, in oppositely propagating light beams (1D) there can be absorption-stimulated emission cycles that connect $|g; p\rangle$ to itself or to $|g; p \pm 2\rangle$ (in this section, momentum is measured in units of $\hbar k$). The initial and final E_k of the atom differ by $\pm 2(p \pm 1)/M$ so energy conservation requires $p = \mp 1$ and is therefore velocity-selective (the energy of the light field is unchanged by the interaction since all the photons in the field have energy $\hbar\omega_l$).

The coupling of these two degenerate states by the light field produces off-diagonal matrix elements of the total Hamiltonian \mathcal{H} of Eq. (35), and subsequent diagonalization of it results in the new ground eigenstates of \mathcal{H} given by (see Fig. 22) $|\pm\rangle \equiv (|g; -1\rangle \pm |g; +1\rangle)/\sqrt{2}$. The excitation rate of these eigenstates $|\pm\rangle$ to $|e; 0\rangle$ is proportional to the square of the electric dipole matrix element $\bar{\mu}$ given by

$$|\langle e; 0 | \bar{\mu} | \pm \rangle|^2 = |\langle e; 0 | \bar{\mu} | g; -1 \rangle \pm \langle e; 0 | \bar{\mu} | g; +1 \rangle|^2 / 2 \quad (36)$$

This vanishes for $|-\rangle$ because the two terms on the right-hand side of Eq. (36) are equal since $\bar{\mu}$ does not operate on the external momentum of the atom (dotted line of Fig. 22). Excitation of $|\pm\rangle$ to $|e; \pm 2\rangle$ is much weaker since it's off resonance because its energy is higher by $4\hbar\omega_l = 2\hbar^2 k^2 / M$, so that the required frequency is higher than to $|e; 0\rangle$. The resultant detuning is $4\omega_l = 8\varepsilon(\gamma/2)$, and for $\varepsilon \sim 0.5$, this is large enough so that the excitation rate is small, making $|-\rangle$ quite dark. Excitation to any state other than $|e; \pm 2\rangle$ or $|e; 0\rangle$ is forbidden by momentum conservation. Atoms are therefore optically pumped into the dark state $|-\rangle$ where they stay trapped, and since their momentum components are fixed, the result is velocity-selective coherent population trapping (VSCPT).

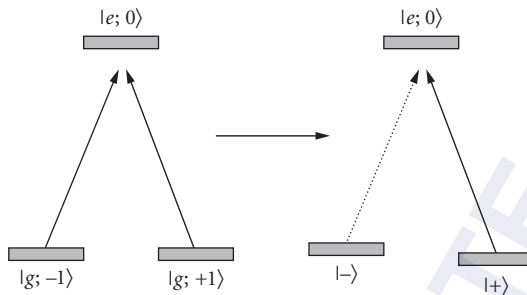


FIGURE 22 Schematic diagram of the transformation of the eigenfunctions from the internal atomic states $|g; p\rangle$ to the eigenstates $|\pm\rangle$. The coupling between the two states $|g; p\rangle$ and $|g; p'\rangle$ by Raman transitions mixes them, and since they are degenerate, the eigenstates of \mathcal{H} are the nondegenerate states $|\pm\rangle$.

A useful view of this dark state can be obtained by considering that its components $|g; \pm 1\rangle$ have well-defined momenta, and are therefore completely delocalized. Thus they can be viewed as waves traveling in opposite directions but having the same frequency, and therefore they form a standing de Broglie wave. The fixed spatial phase of this standing wave relative to the optical standing wave formed by the counterpropagating light beams results in the vanishing of the spatial integral of the dipole transition matrix element so that the state cannot be excited. This view can also help to explain the consequences of p not exactly equal ± 1 , where the de Broglie wave would be slowly drifting in space. It is common to label the average of the momenta of the coupled states as the *family momentum*, \mathcal{P} , and to say that these states form a *closed family*, having family momentum $\mathcal{P} = 0$.^{74,75}

In the usual case of laser cooling, atoms are subject to both a damping force and to random impulses arising from the discrete photon momenta $\hbar k$ of the absorbed and emitted light. These can be combined to make a force versus velocity curve as shown in Fig. 23a. Atoms with $\mathcal{P} \neq 0$ are always subject to the light field that optically pumps them into the dark state and thus produces random impulses as shown in Fig. 23b. There is no damping force in the most commonly studied case of a real atom, the $J = 1 \rightarrow 1$ transition in He^* , because the Doppler and polarization gradient cooling cancel one another as a result of a numerical “accident” for this particular case.

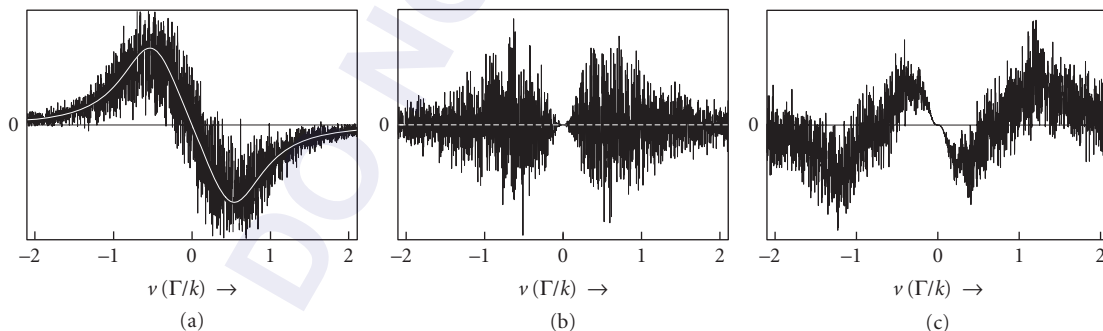


FIGURE 23 Calculated force versus velocity curves for different laser configurations showing both the average force and a typical set of simulated fluctuations. Part (a) shows the usual Doppler cooling scheme that produces an atomic sample in steady state whose energy width is $\hbar\gamma/2$. Part (b) shows VSCPT as originally studied in Ref. 74 with no damping force. Note that the fluctuations vanish for $\mathcal{P} = 0$ because the atoms are in the dark state. Part (c) shows the presence of both a damping force and VSCPT. The fluctuations vanish for $\mathcal{P} = 0$, and both damping *and* fluctuations are present at $\mathcal{P} \neq 0$.

Figures 23a and b should be compared to show the velocity dependence of the sum of the damping and random forces for the two cases of ordinary laser cooling and VSCPT. Note that for VSCPT the momentum diffusion vanishes when the atoms are in the dark state at $\mathcal{P} = 0$, so they can collect there. In the best of both worlds, a damping force would be combined with VSCPT as shown in Fig. 23c. Such a force was predicted in Ref. 76 and was first observed in 1996.⁷⁷

20.9 REFERENCES

1. H. J. Metcalf and P. van der Straten, *Laser Cooling and Trapping*, Springer-Verlag, New York, 1999.
2. J. Gordon and A. Ashkin, "Motion of Atoms in a Radiation Trap," *Phys. Rev. A* **21**:1606 (1980).
3. W. Phillips and H. Metcalf, "Laser Deceleration of an Atomic Beam," *Phys. Rev. Lett.* **48**:596 (1982).
4. J. Prodan, W. Phillips, and H. Metcalf. "Laser Production of a Very Slow Monoenergetic Atomic Beam," *Phys. Rev. Lett.* **49**:1149 (1982).
5. J. Prodan and W. Phillips, "Chirping the Light Fantastic—Recent NBS Atom Cooling Experiments," *Prog. Quant. Elect.* **8**:231 (1984).
6. W. Ertmer, R. Blatt, J. L. Hall, and M. Zhu, "Laser Manipulation of Atomic Beam Velocities: Demonstration of Stopped Atoms and Velocity Reversal," *Phys. Rev. Lett.* **54**:996 (1985).
7. R. Watts and C. Wieman, "Manipulating Atomic Velocities Using Diode Lasers," *Opt. Lett.* **11**:291 (1986).
8. V. Bagnato, G. Lafyatis, A. Martin, E. Raab, R. Ahmad-Bitar, and D. Pritchard, "Continuous Stopping and Trapping of Neutral Atoms," *Phys. Rev. Lett.* **58**:2194 (1987).
9. T. E. Barrett, S. W. Dapore-Schwartz, M. D. Ray, and G. P. Lafyatis, "Slowing Atoms with (σ^-)-Polarized Light," *Phys. Rev. Lett.* **67**:3483–3487 (1991).
10. J. Dalibard and W. Phillips, "Stability and Damping of Radiation Pressure Traps," *Bull. Am. Phys. Soc.* **30**:748 (1985).
11. S. Chu, L. Hollberg, J. Bjorkholm, A. Cable, and A. Ashkin, "Three-Dimensional Viscous Confinement and Cooling of Atoms by Resonance Radiation Pressure," *Phys. Rev. Lett.* **55**:48 (1985).
12. P. D. Lett, R. N. Watts, C. E. Tanner, S. L. Rolston, W. D. Phillips, and C. I. Westbrook, "Optical Molasses," *J. Opt. Soc. Am. B* **6**:2084–2107 (1989).
13. D. Sesko, C. Fan, and C. Wieman. "Production of a Cold Atomic Vapor Using Diode-Laser Cooling," *J. Opt. Soc. Am. B* **5**:1225 (1988).
14. P. Gould, P. Lett, and W. D. Phillips, "New Measurement with Optical Molasses," in *Laser Spectroscopy VIII*, W. Persson and S. Svanberg, (eds.) Springer, Berlin, 1987.
15. T. Hodapp, C. Gerz, C. Westbrook, C. Furtlehner, and W. Phillips, "Diffusion in Optical Molasses," *Bull. Am. Phys. Soc.* **37**:1139 (1992).
16. P. Lett, R. Watts, C. Westbrook, W. Phillips, P. Gould, and H. Metcalf, "Observation of Atoms Laser Cooled Below the Doppler Limit," *Phys. Rev. Lett.* **61**:169 (1988).
17. J. Dalibard and C. Cohen-Tannoudji, "Laser Cooling below the Doppler Limit by Polarization Gradients—Simple Theoretical-Models," *J. Opt. Soc. Am. B* **6**:2023–2045 (1989).
18. P. J. Ungar, D. S. Weiss, S. Chu, and E. Riis, "Optical Molasses and Multilevel Atoms—Theory," *J. Opt. Soc. Am. B* **6**:2058–2071 (1989).
19. C. Cohen-Tannoudji and W. D. Phillips, "New Mechanisms for Laser Cooling," *Phys. Today* **43**: October, 33–40 (1990).
20. C. Salomon, X Dalibard, W. D. Phillips, A. Clairon, and S. Guellati, "Laser Cooling of Cesium Atoms below 3 μ K," *Europhys. Lett.* **12**:683–688 (1990).
21. W Ketterle and N. J. van Druten, "Evaporative Cooling of Trapped Atoms," *Adv. Atom. Mol. Opt. Phys.* **37**:181 (1996).
22. M. Kasevich and S. Chu, "Laser Cooling below a Photon Recoil with 3-Level Atoms," *Phys. Rev. Lett.* **69**:1741–1744 (1992).

23. D. Wineland, W. Itano, J. Bergquist, and J. Bollinger, "Trapped Ions and Laser Cooling," Technical Report 1086, N.I.S.T (1985).
24. A. Migdall, J. Prodan, W. Phillips, T. Bergeman, and H. Metcalf, "First Observation of Magnetically Trapped Neutral Atoms," *Phys. Rev. Lett.* **54**:2596 (1985).
25. T. Bergeman, G. Erez, and H. Metcalf, "Magnetostatic Trapping Fields for Neutral Atoms," *Phys. Rev. A* **35**:1535 (1987).
26. A. Ashkin, "Acceleration and Trapping of Particles by Radiation Pressure," *Phys. Rev. Lett.* **24**:156 (1970).
27. S. Chu, J. Bjorkholm, A. Ashkin, and A. Cable, "Experimental Observation of Optically Trapped Atoms," *Phys. Rev. Lett.* **57**:314 (1986).
28. A. Ashkin, "Application of Laser Radiation Pressure," *Science* **210**:1081–1088 (1980).
29. A. Ashkin and J. M. Dziedzic, "Observation of Radiation-Pressure Trapping of Particles by Alternating Light Beams," *Phys. Rev. Lett.* **54**:1245 (1985).
30. A. Ashkin and J. M. Dziedzic, "Optical Trapping and Manipulation of Viruses and Bacteria," *Science* **235**:1517 (1987).
31. J. D. Miller, R. A. Cline, and D. J. Heinzen, "Far-Off-Resonance Optical Trapping of Atoms," *Phys. Rev. A* **47**:R4567–R4570 (1993).
32. Y. Castin and J. Dalibard, "Quantization of Atomic Motion in Optical Molasses," *Europhys. Lett.* **14**:761–766 (1991).
33. P. Verkerk, B. Lounis, C. Salomon, C. Cohen-Tannoudji, J. Y. Courtois, and G. Grynberg, "Dynamics and Spatial Order of Cold Cesium Atoms in a Periodic Optical-Potential," *Phys. Rev. Lett.* **68**:3861–3864 (1992).
34. P. S. Kessen, C. Gerz, P. D. Lett, W. D. Phillips, S. L. Rolston, R. J. C. Spreeuw, and C. I. Westbrook, "Observation of Quantized Motion of Rb Atoms in an Optical-Field," *Phys. Rev. Lett.* **69**:49–52 (1992).
35. B. Lounis, P. Verkerk, J. Y. Courtois, C. Salomon, and G. Grynberg, "Quantized Atomic Motion in 1D Cesium Molasses with Magnetic-Field," *Europhys. Lett.* **21**:13–17 (1993).
36. E. Raab, M. Prentiss, A. Cable, S. Chu, and D. Pritchard, "Trapping of Neutral-Sodium Atoms with Radiation Pressure," *Phys. Rev. Lett.* **59**:2631 (1987).
37. T. Walker, D. Sesko, and C. Wieman, "Collective Behavior of Optically Trapped Neutral Atoms," *Phys. Rev. Lett.* **64**:408–411 (1990).
38. D. W. Sesko, T. G. Walker, and C. E. Wieman, "Behavior of Neutral Atoms in a Spontaneous Force Trap," *J. Opt. Soc. Am. B* **8**:946–958 (1991).
39. J. Nellessen, J. H. Muller, K. Sengstock, and W. Ertmer, "Laser Preparation of a Monoenergetic Sodium Beam," *Europhys. Lett.* **9**:133–138 (1989).
40. J. Nellessen, J. H. Muller, K. Sengstock, and W. Ertmer, "Large-Angle Beam Deflection of a Laser-Cooled Sodium Beam," *J. Opt. Soc. Am. B* **6**:2149–2154 (1989).
41. J. Nellessen, J. Werner, and W. Ertmer, "Magneto-optical Compression of a Monoenergetic Sodium Atomic-Beam," *Opt. Commun.* **78**:300–308 (1990).
42. E. Riis, D. S. Weiss, K. A. Moler, and S. Chu, "Atom Funnel for the Production of a Slow, High-Density Atomic-Beam," *Phys. Rev. Lett.* **64**:1658–1661 (1990).
43. A. Aspect, N. Vansteenkiste, R. Kaiser, H. Haberland, and M. Karrais, "Preparation of a Pure Intense Beam of Metastable Helium by Laser Cooling," *Chem. Phys.* **145**:307–315 (1990).
44. M. D. Hoogerland, J. P. J. Driessen, E. J. D. Vredendregt, H. J. L. Megens, M. P. Schuwer, H. C. W. Beijerinck, and K. A. H. van Leeuwen, "Bright Thermal Atomic-Beams by Laser Cooling—A 1400-Fold Gain in-Beam Flux," *App. Phys. B* **62**, 323–327 (1996).
45. R. A. Nauman and H. Henry Stroke, "Apparatus Upended: A Short History of the Fountain A-Clock," *Phys. Today* **89** (May 1996).
46. M. A. Kasevich, E. Riis, S. Chu, and R. G. Devoe, "RF Spectroscopy in an Atomic Fountain," *Phys. Rev. Lett.* **63**:612–616 (1989).
47. A. Clairon, C. Salomon, S. Guellati, and W. D. Phillips, "Ramsey Resonance in a Zacharias Fountain," *Europhys. Lett.* **16**:165–170 (1991).

48. K. Gibble and S. Chu, "Future Slow-Atom Frequency Standards," *Metrologia* **29**:201–212 (1992).
49. K. Gibble and S. Chu, "Laser-Cooled Cs Frequency Standard and a Measurement of the Frequency-Shift Due to Ultracold Collisions," *Phys. Rev. Lett.* **70**:1771–1774 (1993).
50. K. Gibble and B. Verhaar, "Eliminating Cold-Collision Frequency Shifts," *Phys. Rev. A* **52**:3370 (1995).
51. R. Legere and K. Gibble, "Quantum Scattering in a Juggling Atomic Fountain," *Phys. Rev. Lett.* **81**:5780 (1998).
52. Ph. Laurent, P. Lemonde, E. Simon, G. Santorelli, A. Clairon, N. Dimarcq, P. Petit, C. Audoin, and C. Salomon, "A Cold Atom Clock in the Absence of Gravity," *Eur. Phys. J. D* **3**:201 (1998).
53. P.S. Julienne and E.H. Mies, "Collisions of Ultracold Trapped Atoms," *J. Opt. Soc. Am.* **B6**:2257–2269 (1989).
54. P. D. Lett, P. S. Julienne, and W. D. Phillips, "Photoassociative Spectroscopy of Laser-Cooled Atoms," *Annual Rev. Phys. Chem.* **46**:423 (1995).
55. A. Gallagher and D. E. Pritchard, "Exoergic Collisions of Cold Na*–Na," *Phys. Rev. Lett.* **63**:957–960 (1989).
56. P. S. Julienne and J. Vigue, "Cold Collisions of Ground-State and Excited-State Alkali-Metal Atoms," *Phys. Rev. A* **44**:4464–4485 (1991).
57. H. Mastwijk, J. Thomsen, P. van der Straten, and A. Niehaus, "Optical Collisions of Cold, Metastable Helium Atoms," *Phys. Rev. Lett.* **80**:5516–5519 (1998).
58. A. Fioretti, D. Comparat, A. Crubellier, O. Dulieu, F. Masnou-Seeuws, and P. Pillet, "Formation of Cold Cs₂ Molecules through Photoassociation," *Phys. Rev. Lett.* **80**:4402–4405 (1998).
59. V. S. Lethokov, "Narrowing of the Doppler Width in a Standing Light Wave," *JETP Lett.* **7**:272 (1968).
60. C. Salomon, J. Dalibard, A. Aspect, H. Metcalf, and C. Cohen-Tannoudji, "Channeling Atoms in a Laser Standing Wave," *Phys. Rev. Lett.* **59**:1659 (1987).
61. K. I. Petsas, A. B. Coates, and G. Grynberg, "Crystallography of Optical Lattices," *Phys. Rev. A* **50**:5173–5189 (1994).
62. P. S. Jessen and I. H. Deutsch, "Optical Lattices," *Adv. Atom. Mol. Opt. Phys.* **37**:95–138 (1996).
63. G. Grynberg, B. Lounis, P. Verkerk, J. Y. Courtois, and C. Salomon, "Quantized Motion of Cold Cesium Atoms in 2-Dimensional and 3-Dimensional Optical Potentials," *Phys. Rev. Lett.* **70**:2249–2252 (1993).
64. G. Birkl, M. Gatzke, I. H. Deutsch, S. L. Rolston, and W. D. Phillips, "Bragg Scattering from Atoms in Optical Lattices," *Phys. Rev. Lett.* **75**:2823–2826 (1995).
65. C. I. Westbrook, R. N. Watts, C. E. Tanner, S. L. Rolston, W. D. Phillips, P. D. Lett, and P. L. Gould, "Localization of Atoms in a 3-Dimensional Standing Wave," *Phys. Rev. Lett.* **65**:33–36 (1990).
66. M. Dahan, E. Peik, J. Reichel, Y. Castin, and C. Salomon, "Bloch Oscillations of Atoms in an Optical Potential," *Phys. Rev. Lett.* **76**:4508 (1996).
67. H. T. C. Stoof, "Atomic Bose-Gas with a Negative Scattering Length," *Phys. Rev. A* **49**:3824–3830 (1994).
68. T. Bergeman, "Hartree-Fock Calculations of Bose-Einstein Condensation of ⁷Li Atoms in a Harmonic Trap for $T > 0$," *Phys. Rev. A* **55**:3658 (1997).
69. T. Bergeman, "Erratum: Hartree-Fock Calculations of Bose-Einstein Condensation of ⁷Li Atoms in a Harmonic Trap for $T > 0$," *Phys. Rev. A* **56**:3310 (1997).
70. M. H. Anderson, J. R. Ensher, M. R. Matthews, C. E. Wieman, and E. A. Cornell, "Observation of Bose-Einstein Condensation in a Dilute Atomic Vapor," *Science* **269**:198–201 (1995).
71. K. Davis, M.-O. Mewes, M. Andrews, M. van Druten, D. Durfee, D. Kurn, and W. Ketterle, "Bose-Einstein Condensation in a Gas of Sodium Atoms," *Phys. Rev. Lett.* **75**:3969 (1995).
72. C. C. Bradley, C. A. Sackett, J. J. Tollett, and R. G. Hulet, "Evidence of Bose-Einstein Condensation in an Atomic Gas with Attractive Interactions," *Phys. Rev. Lett.* **75**:1687–1690 (1995).
73. D. Fried, T. Killian, L. Willmann, D. Landhuis, S. Moss, D. Kleppner, and T. Greytak, "Bose-Einstein Condensation of Atomic Hydrogen," *Phys. Rev. Lett.* **81**:3811 (1998).
74. A. Aspect, E. Arimondo, R. Kaiser, N. Vansteenkiste, and C. Cohen-Tannoudji, "Laser Cooling below the One-Photon Recoil Energy by Velocity-Selective Coherent Population Trapping," *Phys. Rev. Lett.* **61**:826 (1988).

75. A. Aspect, C. Cohen-Tannoudji, E. Arimondo, N. Vansteenkiste, and R. Kaiser, "Laser Cooling Below the One-Photon Recoil Energy by Velocity-Selective Coherent Population Trapping—Theoretical-Analysis," *J. Opt. Soc. Am. B* **6**:2112–2124 (1989).
76. M. S. Shahriar, P. R. Hemmer, M. G. Prentiss, P. Marte, J. Mervis, D. P. Katz, N. P. Bigelow, and T. Cai, "Continuous Polarization-Gradient Precooling-Assisted Velocity-Selective Coherent Population Trapping," *Phys. Rev. A* **48**:R4035–R4038 (1993).
77. M. Widmer, M. J. Bellanca, W. Buell, H. Metcalf, M. Doery, and E. Vredenburg, "Measurement of Force-Assisted Population Accumulation in Dark States," *Opt. Lett.* **21**:606–608 (1996).
78. P. A. Molenaar, P. van der Straten, H. G. M. Heideman, and H. Metcalf, "Diagnostic-Technique for Zeeman-Compensated Atomic-Beam Slowing—Technique and Results. *Phys. Rev. A* **55**:605–614 (1997).
79. B. Sheehy, S. Q. Shang, P. van der Straten, and H. Metcalf, "Collimation of a Rubidium Beam Below the Doppler Limit," *Chem. Phys.* **145**:317–325 (1990).
80. A. Scholz, M. Christ, D. Doll, J. Ludwig, and W. Ertmer, "Magneto-optical Preparation of a Slow, Cold and Bright Ne* Atomic-Beam," *Opt. Commun.* **111**:155–162 (1994).
81. Z. T. Lu, K. L. Corwin, M. J. Renn, M. H. Anderson, E. A. Cornell, and C. E. Wieman, "Low-Velocity Intense Source of Atoms from a Magneto-optical Trap," *Phys. Rev. Lett.* **77**:3331–3334 (1996).
82. K. G. H. Baldwin, private communication.
83. M. Schiffer, M. Christ, G. Wokurka, and W. Ertmer, "Temperatures Near the Recoil Limit in an Atomic Funnel," *Opt. Commun.* **134**:423–430 (1997).
84. F. Lison, P. Schuh, D. Haubrich, and D. Meschede, "High Brilliance Zeeman Slowed Cesium Atomic Beam," *Phys. Rev. A* **61**:013405 (2000).
85. K. Dieckmann, R. J. C. Spreeuw, M. Weidemuller, and J. T. M. Walraven, "Two-Dimensional Magneto-Optical Trap as a Source of Slow Atoms," *Phys. Rev. A* **58**:3891 (1998).

Todd Ditmire

*Texas Center for High Intensity Laser Science
Department of Physics
The University of Texas at Austin
Austin, Texas*

21.1 GLOSSARY

a_0	normalized peak vector potential of the intense laser pulse
a_{Bohr}	Bohr radius
b	confocal parameter
c	speed of light
e	charge of the electron
f_p	relativistic ponderomotive force
E_a	atomic field strength
E_{cr}	Schwinger critical electric field
E_0	peak electric field amplitude of the intense laser pulse
f_x	laser fractional absorption for x process
k_0	laser wavenumber
k_e	electron wavenumber
$\text{KE}_{e-,ATI}$	kinetic energy of electrons
I	laser intensity
I_H	ionization potential of hydrogen
I_p	ionization potential
ℓ	plasma scale length
L_c	coherence length
m_e, m_i	mass of the electron, ion
n	refractive index
n_2	nonlinear refractive index
n_{crit}	electron critical density
n_e, n_i	electron and ion density
p	effective harmonic nonlinear order

$P_{ ,\perp}$	polarizability tensor of a molecule parallel and perpendicular to molecular axis
P_c	critical power for self-focusing
q	harmonic order
Q_{clust}	charge on cluster from outer ionization
R_0	initial radius of a cluster
R_c	critical ionization distance in molecules
u_{ion}	ion velocity
U_p	ponderomotive potential energy
v_g	laser group velocity
v_{osc}	electron oscillation velocity
v_D	electron drift velocity
w	$1/e^2$ focal spot radius of a focused Gaussian laser beam
W_x	ionization rate for x process
z_R	Rayleigh range of focused laser
Z	charge state of ions
α	fine structure constant
β_{Rot}	molecular rotation constant
Δk	phase mismatch
ϵ_{CE}	Coulomb explosion energy
ϵ_D	dielectric function
γ	relativistic Lorentz factor for electrons
γ_{osc}	cycle-averaged relativistic Lorentz factor for electrons in strong field
γ_{SRS}	stimulated Raman scattering growth rate
Λ	Coulomb logarithm
λ_{Debye}	plasma Debye length
λ_p	optical scale length in an overdense plasma
ν_{ei}	electron-ion collision frequency
σ_N	generalized N -photon cross section for multiphoton ionization
σ_T	Thomson scattering cross section
τ_p	laser pulse duration
ω_a	atomic unit of frequency
ω_{BG}	Bohm-Gross frequency
ω_0	angular oscillation frequency of the intense laser pulse
ω_p	plasma frequency
$\omega_{\text{S,A}}$	Stokes, anti-Stokes frequency

21.2 INTRODUCTION AND HISTORY

Strong field physics (or “high field physics” in much of the literature) refers to the phenomena that occur during the interaction of intense electromagnetic waves with matter of various forms. It is characterized by interactions that are often highly nonlinear. While such interactions have been accessed with microwave radiation,¹ traditionally, strong field physics has been studied with intense optical and near-infrared (IR) pulses generated by high-intensity lasers. These interactions occur in a regime in which the electric field of the optical wave dominates the motion and dynamics of electrons subject to these fields. At the highest intensities that are now accessible, the motion of

electrons can become relativistic during each optical cycle, and the magnetic field of the light pulse starts to become important in affecting the motion of electrons in the field.

While it is possible to access strong field effects with what are presently considered rather modest intensities in certain situations, it is customary to consider strong field physics as the regime in which the light intensity is high enough that the peak electric field of the wave $E_0 = \sqrt{8\pi I/c}$ becomes comparable to the atomic unit of electric field $E_a = e/a_{\text{Bohr}}^2 = 5.1 \times 10^9$ V/cm, the field felt by an electron in a hydrogen atom. Light acquires this electric field at an intensity of 3.5×10^{16} W/cm², though there are many strong field physics effects which manifest themselves at fields of about 10 percent of E_a (at intensity $\sim 10^{14}$ W/cm²). At these intensities light interaction with atoms can no longer be described by standard perturbation theory, and light interactions with electrons in a plasma dominate the thermal motion of the free electrons. At higher intensities, beyond 10^{18} W/cm², the field becomes high enough that an electron in an optical frequency wave can be accelerated to relativistic velocity in less than one optical cycle. Such intensities are also characterized by high magnetic fields and optical forces. For example, in a pulse with intensity of 10^{18} W/cm², an intensity quite modest by modern standards, the peak electric field is 3×10^{10} V/cm and the optical magnetic field is 100 MG. The light pressure, I/c , is ~ 0.3 Gbar. The highest intensity lasers can now reach intensity approaching 10^{22} W/cm².

The theoretical study of high field physics can be said to have started in earnest with a classic paper by L. V. Keldysh in 1964.² In this paper, the rate of ionization of an atom or ion in a strong laser field was first derived with a nonperturbative theory. The first real experimental observation of nonperturbative high-field effects occurred in the ground breaking experiment of Agostini et al. in 1979.³ Their experiment observed, for the first time, truly nonperturbative multiphoton effects in laser-atom interactions by examining photo-electron production from intense 6-photon ionization of Xe atoms at intensity up to 4×10^{13} W/cm². They found that electrons were ejected during ionization with energy higher than that expected from absorption of the minimum number of photons needed for ionization, an effect that came to be known as above threshold ionization.⁴ This observation sparked a long campaign of experiments and theoretical work on strong laser field ionization of atoms and ions that continues to this day. These early experiments in strong field multiphoton ionization were followed by the first observation of nonperturbative nonlinear optical phenomena in high order harmonic generation by Rhodes et al. in 1987⁵ in which highly nonlinear interactions of an intense laser pulse with a gas of atoms led to emission of a range of high harmonics of the laser frequency. The initial observations of high harmonics were striking in that a range of harmonics extended to very high orders with almost constant intensity, completely at odds with lowest order perturbation theory. In fact, very high nonlinear orders, >100 , have been reported in studies of this effect,^{6,7} resulting in the production of coherent light into the soft x-ray region. High harmonic generation with intense laser pulses continues to be studied actively and has led to a revolution in the production of electromagnetic pulses with durations of a few hundred attoseconds.^{8,9}

These early nonlinear multiphoton discoveries were followed by the realization that much of these effects could be understood by treating the field classically and the interaction with electrons semiclassically. This simplification in describing strong-field interactions occurred nearly simultaneously by Corkum et al.¹⁰ and Kulander.¹¹ The semiclassical treatment is now the basis for much of our understanding of strong field ionization, above threshold ionization and high harmonic generation.

While the study of strong field physics has its origins in the study of atomic ionization, it was also realized early on that strong field interactions with plasmas would manifest unique effects, not only through the ionization of atoms and ions in the plasma but in the collective motion of the plasma electrons driven by the strong forces of an intense laser pulse. The study of intense laser interactions with plasmas has been a very important aspect of strong field physics leading to numerous breakthroughs, such as the development of compact x-ray lasers^{12,13} and plasma accelerators.¹⁴ For example, it was realized as early as 1979, in a classic paper by Tajima and Dawson,¹⁵ that an intense laser pulse could be used to drive a plasma wave that could, in turn, accelerate electrons with very high gradient, far above that of traditional accelerators.¹⁶ This has led to a steady advance in understanding intense laser light propagation in underdense plasma. Recent years have seen other important advances in strong field laser interactions. The latest developments in laser technology now enable experiments in plasmas at intensities in which the free electron velocity becomes relativistic.

This has led to new range of nonlinear phenomena created by the relativistic mass increase of the electron in the strong laser field. For example, absorption of light in plasmas becomes much more complicated in the high-field regime, with collective effects playing a much bigger role. This leads to plasma interactions that exhibit large “anomalous absorption” deviating from simple linear kinetic theories of light interactions with plasmas.

This chapter is intended to introduce many of the fundamental concepts underlying the modern strong field physics research. These concepts span descriptions of intense light interactions with single electrons, individual atoms, ensembles of atoms in molecules and clusters, and many charged particles in plasmas. This chapter does not represent a comprehensive review of modern strong field physics research and is not a survey of recent results in the field. No attempt is made to discuss specific experimental results that confirm the phenomena presented (though citations to such work are given). Instead, the basic phenomena underlying the more complex effects observed in strong field physics will be discussed, and the basic equations needed to describe these high-field effects will be presented. (Equations here are presented without proof; the reader is encouraged to seek detailed derivations from the references provided.) If a more detailed review of the various aspects of strong field physics is desired, there have been a number of excellent review articles published in recent years (including a number of older articles which are still relevant). A listing of some of these review articles appears in Sec. 21.12 in Refs. 17 to 41 for the interested reader. (In the remainder of this chapter all units are CGS unless otherwise stated.)

21.3 LASER TECHNOLOGY USED IN STRONG FIELD PHYSICS

Before discussing strong field phenomena, it is important to note that advances in this area of physics have been driven by many important leaps in laser technology over the last 20 years. The enabling technology advancement for creating the increasingly higher intensities needed to access strong fields was the invention of chirped pulse amplification (CPA) lasers.⁴² The CPA technique, first demonstrated by Strickland and Mourou in 1985,⁴² is illustrated in Fig. 1. The goal of CPA is to amplify picosecond to femtosecond duration pulses to high energy in laser gain media, thereby

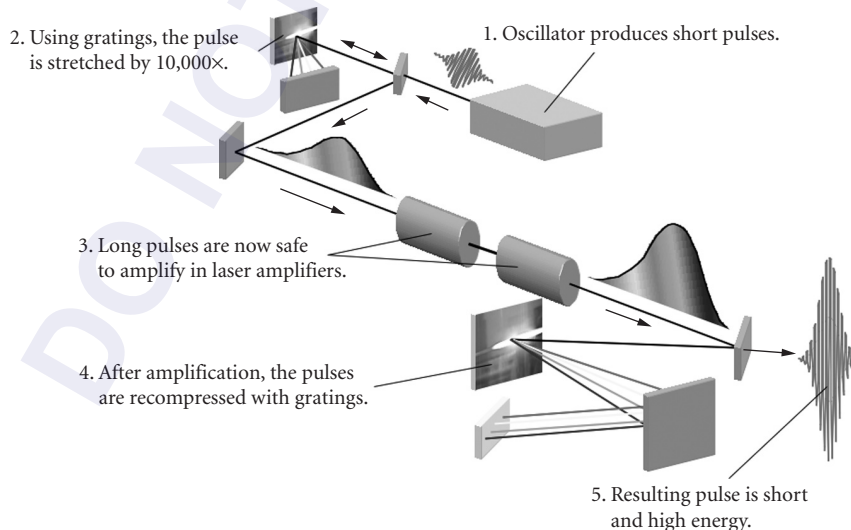


FIGURE 1 Architecture used for chirped pulse amplification (CPA) lasers.

reaching the terawatts to petawatts of peak power needed to access strong field phenomena. As illustrated in Fig. 1 a broad bandwidth, mode-locked laser produces a low power, ultrafast pulse of light, usually with duration of 20 to 500 fs. This short pulse is first stretched in time by a factor of around 10 thousand from its original duration using diffraction gratings. This allows the pulses, now of much lower peak power, to be safely amplified in the laser, avoiding the deleterious nonlinear effects which would occur if the pulses had higher peak power.⁴³ These amplified pulses are, finally, recompressed in time, (again using gratings) in a manner that preserves the phase relationship between the component frequencies in the pulse. The CPA laser pulse output has a duration near that of the original pulse but with an energy greater by the amplification factor of the laser chain. In high-energy CPA systems ($> \sim 1$ J), severe nonlinearities occurring when the pulse propagates in air can be a major problem, so the pulse must be recompressed in an evacuated chamber. The state of the art in CPA now enables focused intensity of up to 10^{21} W/cm² (Ref. 44) with peak intensity up to 1 PW (10^{15} W).⁴⁵ Table top CPA lasers can usually access intensity of $\sim 10^{19}$ W/cm² and high repetition rate (~ 1 kHz) lasers usually operate with peak intensity of $\sim 10^{16}$ W/cm².

The first generation of CPA lasers was based primarily on flashlamp pumped Nd:glass amplifiers.^{42,46–48} These glass-based lasers, operating at a wavelength near 1 μ m, are usually limited to pulse duration of greater than about 400 fs because of gain narrowing in the amplifiers.⁴⁹ The most significant scaling of this approach to CPA was demonstrated by the petawatt laser at Lawrence Livermore National Laboratory in the late 1990s.⁵⁰ This laser demonstrated the production of 500 J of energy per pulse with duration of under 500 fs, yielding over 10^{15} W of peak power. Since this demonstration, a number of petawatt laser projects have been undertaken around the world.⁵¹

The second common approach to CPA uses Ti:sapphire as the amplifier material. This material permits amplification of 800 nm wavelength pulses with much shorter pulse durations, often down to ~ 30 fs. However, the short excited state lifetime of Ti:sapphire (3 μ s) requires that the material be pumped by a second laser (usually a frequency doubled Nd:YAG or Nd:glass laser). The inherent inefficiencies of this two-step pumping usually limit the output energy of such a laser to under a few joules of energy per pulse. A number of multiterawatt lasers based on Ti:sapphire now operate in many high-intensity laser labs worldwide.^{52–55} These laser typically yield energy of 1 mJ to 1 J (though higher energy examples with energy ~ 10 J do exist), and they typically exhibit repetition rate of 1 kHz at the 1 mJ level or ~ 10 Hz at the 0.1 to 1 J level.⁵⁶ To date, the largest scaling of Ti:sapphire technology has been to power levels approaching 1 PW.^{54,57}

The third major technology now commonly used in CPA lasers is based on a technique known as optical parametric chirped pulse amplification (OPCPA).⁵⁸ In this approach, amplification of the stretched pulses occurs not with an energy storage medium like Nd:glass or Ti:sapphire but via parametric interactions in a nonlinear crystal. This approach is quite attractive because of the very high gain per stage possible (often in excess of 10^4 per pass) and the very broad gain bandwidth possible, in principal. To date, a number of CPA lasers based on OPCPA have been demonstrated.^{59–61}

21.4 STRONG FIELD INTERACTIONS WITH SINGLE ELECTRONS

We begin our discussion of the various strong field phenomena by considering strong field interactions with individual, free electrons (which encompasses the situation in which the electrons are not affected by the electrostatic forces of a collective electron plasma). This discussion will be followed in the next sections by overviews of strong laser field interactions with atoms, molecules, clusters, and then plasmas.

The Ponderomotive Force

When a strong laser field interacts with a free electron, the field can almost always be treated classically and the trajectory of the electron can be found using classical mechanics. If the intensity is

not too high, below about 10^{18} W/cm² for optical and near-infrared frequencies, then motion of the electron can be treated nonrelativistically and the magnetic field of the laser can be ignored. In that case, the electron oscillates at the laser frequency in the direction of the laser polarization. Solution of Newton's equation yields for electron velocity $v(t) = (eE_0/m_e\omega_0) \sin(\omega_0 t)$, where E_0 is the laser's peak electric field, ω_0 is the frequency of the laser light, and $v_{\text{osc}} = (eE_0/m_e\omega_0)$ is the classical electron oscillation velocity amplitude. While complications arise for very short laser pulses (with envelope comparable to the wavelength) or tightly focused pulses, in a plane wave it is useful to consider the cycle-averaged kinetic energy of this oscillating electron. This energy, called the ponderomotive potential, is

$$U_p = \frac{e^2 E_0^2}{4m_e \omega_0^2} \quad (1)$$

In practical units, the ponderomotive energy is equal to $9.33 \times 10^{-14} I$ (W/cm²) λ^2 (μm) in eV and is, for example, roughly 10 keV in a Nd:glass laser field at 1.054 μm focused to intensity of 10^{17} W/cm². This ponderomotive energy usually sets the energy scale for most strong field interactions. In a focused laser beam a force $-\nabla U_p$, called the ponderomotive force, will act on an electron. The ponderomotive force will tend to accelerate electrons transversely out of the focus from high to low intensity (increasing the electron's energy by $\sim U_p$). In the absence of a strong transverse intensity gradient, a free electron illuminated by a strong laser field will begin to oscillate as the field amplitude increases, but will then come back completely to rest as the intensity falls back to zero, acquiring no net energy from the laser field.

Relativistic Effects in Strong Field Interactions with Free Electrons

The dynamics of electrons in field strengths at which relativistic effects become important are considerably more complicated. These effects become important when the classical, *nonrelativistic* quantity v_{osc} becomes comparable to or greater than c . At optical and near-IR wavelengths relativistic effects become important at intensity approaching 10^{18} W/cm² [where Eq. (1) predicts that the ponderomotive energy exceeds 100 keV in near-IR fields and is, therefore, a large fraction of the 511 keV electron rest mass]. The extent to which the field interaction with the electron is relativistic can be quantified by the dimensionless *normalized vector potential*, a_0 , which is v_{osc}/c , and given by

$$a_0 = \frac{eE_0}{m_e c \omega_0} \quad (2)$$

When a_0 approaches 1 (which occurs for 1 μm light at intensity equal to 1.4×10^{18} W/cm²), Eq. (1) breaks down. In this regime, the electron motion in the strong field becomes significantly affected by the laser's magnetic field, and, while complex, has been solved analytically by a number of authors.⁶²⁻⁶⁵

The electron's motion deviates from the simple harmonic oscillation described above because the $\mathbf{v} \times \mathbf{B}$ force drives the electron forward. The electron now acquires a significant velocity component in the laser's k direction and, as a result, no longer experiences linearly varying phase and a perfectly sinusoidal oscillation of the electric field. The electron's motion becomes highly anharmonic. Figure 2 illustrates the trajectory of an electron irradiated by a laser field of 1 μm wavelength at intensity of 10^{19} W/cm², ($a_0 = 2.7$). The longitudinal momentum, p_z can be easily derived from the relativistic equations of motion for the electron in an EM wave. For an electron initially at rest, the field will yield forward momentum such that

$$p_z = \frac{p_x^2}{2m_e c} \quad (3)$$

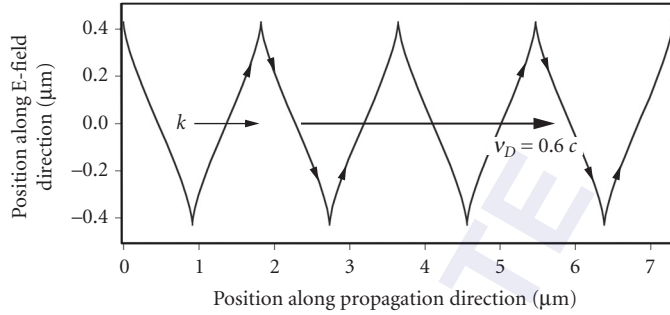


FIGURE 2 Trajectory of an electron driven by a relativistic laser field with $1\ \mu\text{m}$ wavelength at an intensity of $10^{19}\ \text{W}/\text{cm}^2$, ($a_0 = 2.7$ and $\gamma^{\text{osc}} = 2.1$). This illustrates that the electron acquires significant drift velocity along the laser propagation direction and oscillates anharmonically.

if p_x is the transverse, E-field driven momentum.^{66–68} For an electron in a plane wave, a cycle-averaged forward drift velocity v_D will be acquired in the lab frame, as pictured in Fig. 2. This *average* forward drift velocity is

$$\frac{v_D}{c} = \frac{a_0^2}{4 + a_0^2} \quad (4)$$

This equation indicates that the electron will drift in the forward direction at nearly the speed of light when a_0 is roughly 10, corresponding to a near-IR intensity of about $10^{20}\ \text{W}/\text{cm}^2$.

Finding the transverse oscillation velocity is more complex, but a useful result can be found for weak relativistic fields (when a_0 is between ~ 0.3 and 3 or intensity is in the 10^{17} to $10^{19}\ \text{W}/\text{cm}^2$ range for near-IR light). In this regime, one can neglect the longitudinal velocity in the relativistic equations of motion to find an approximate result for the transverse oscillation velocity

$$\frac{v_x}{c} = \frac{a_0}{1 + a_0} \quad (5)$$

which reduces to v_{osc} when $a_0 \ll 1$.

It is now possible to talk in terms of a relativistic ponderomotive energy which can be written as²²

$$U_p^{\text{rel}} = (\gamma_{\text{osc}} - 1)mc^2 \quad (6)$$

where $\gamma_{\text{osc}} = \sqrt{1 + a_0^2/2}$ is the effective cycle-averaged relativistic Lorentz factor (ignoring the contribution from the slower drift of the electron) for linearly polarized light. (Note that in circularly polarized light $\gamma_{\text{osc}}^{\text{circ}} = \sqrt{1 + a_0^2}$.) Equation (6) reduces to Eq. (1) when $a_0 \ll 1$. Deriving a “ponderomotive force” in a relativistic light beam is a less concrete concept, as the idea of a cycle-averaged force in a weak intensity gradient is of limited utility when, at strongly relativistic intensity, the electron surfs along, almost in phase with the light field at c . However, a heuristic treatment of the electron dynamics yields a relativistic ponderomotive force that can be written as⁶⁹

$$\mathbf{f}_p = -m_e c^2 \nabla \gamma_{\text{osc}} \quad (7)$$

Finally, we note that this relativistic motion will also eject electrons from the focus of an intense laser. However, unlike the nonrelativistic case in which the ponderomotive force drives the electrons out of the focus at 90° to the laser propagation, the magnetic force causes ejection along an axis folded forward toward the k direction. This ejection angle can be simply calculated by considering the relativistic kinematics of the absorption of many photons of momentum $\hbar k$ and the relativistic

relationship between transverse and longitudinal momentum in the EM field, Eq. (3). This yields an ejection angle, θ , of electrons with Lorentz factor outside the laser focus of γ

$$\theta = \arctan \left[\frac{2}{\gamma - 1} \right] \quad (8)$$

It should be noted that this equation is only valid for purely plane wave interactions. The field near a real focus will contain longitudinal components which alter the ejection angle of the electrons somewhat.⁷⁰ Nonetheless, Eq. (8) correctly indicates that in a strongly relativistic focus, electrons will be ejected with high Lorentz factor and will come out in a cone around the propagation axis of the laser, having surfed along with the field through the laser focus. If electrons are ejected into the field at the peak of an oscillation, a situation that occurs when highly charged ions are ionized by the field (see the next section) the relativistic electrons propagating with small ejection angle θ can pick up substantial energy from the laser field. A useful estimate for *maximum* ejected electron energy derived from this relativistic free-wave acceleration mechanism is⁷¹

$$\gamma_{\max} \approx \frac{eE_0 \theta z_R}{2m_e c^2} \quad (9)$$

where $z_R = \pi w^2 / \lambda$ is the Rayleigh range of the Gaussian focus. Equation (9) indicates that the ejected electron in the relativistic intensity regime acquire energy proportional to the laser's electric field ($\sim I^{1/2}$) and not its square, as they do in the nonrelativistic case [see Eq. (1)]. When $a_0 \gg 1$, θ will be small, and electrons will gain substantial energy from the field. For example, calculations show that an electron produced by ionization at intensity of 5×10^{21} W/cm² in a 1- μ m wavelength laser field focused to a 5- μ m Gaussian spot can acquire energy up to ~ 1 GeV and are ejected at an angle of $\sim 3^\circ$.⁷¹

Nonlinear Thomson Scattering by Electrons in an Intense Laser Field

The acceleration associated with oscillation of a free electron in an electromagnetic wave gives rise to emitted radiation, a process known as Thomson scattering, which has the well-known scattering cross section $\sigma_T = 8\pi e^2 / 3m_e c^2$.⁷² The *anharmonic* motion of an electron in a strong laser field, as illustrated in Fig. 2, alters the Thomson scattering in an important manner. First, the light scattered from the electron has an emission pattern folded forward toward the laser propagation axis, and second, the radiated light from the electron will contain higher harmonic components.^{64,73–77} In a strongly relativistic field ($a_0 \gg 1$), the radiated emission will be forward folded by the effective Lorentz boost, into an angle $\theta \approx 3/a_0$. In addition, the anharmonic motion of the electron induced by the magnetic field results in scattered light at even and odd harmonics of the incident light field. It can be shown that the total integrated scattered power into the first three harmonics of the laser field are⁷⁴

$$P_1 \cong \frac{e^2 \omega_0^2 c}{3} a_0^2 \quad P_2 \cong \frac{7e^2 \omega_0^2 c}{20} a_0^4 \quad P_3 \cong \frac{621e^2 \omega_0^2 c}{1792} a_0^6 \quad (10)$$

illustrating the nonlinear intensity dependence of the second and third harmonic. These equations also illustrate that as $a_0 \rightarrow 1$ the power scattered into harmonics (P_2, P_3) will be comparable to the power scattered by linear Thomson scattering, P_1 . The angular distribution of the scattered radiation for linearly polarized light is analytically complex. Figure 3 illustrates the polar distribution of light from the second and third harmonics of nonlinear Thomson scattering at $a_0 = 1$. A simple formula for the angular distribution, $D(\theta)$, of scattered light from a circularly polarized beam as a function of polar angle θ with respect to the laser propagation direction can be derived,⁷⁴ where

$$D(a_0, \theta) = \frac{1}{\left(1 + \frac{a_0^2}{2} \sin^2 \frac{\theta}{2} \right)^4} \quad (11)$$

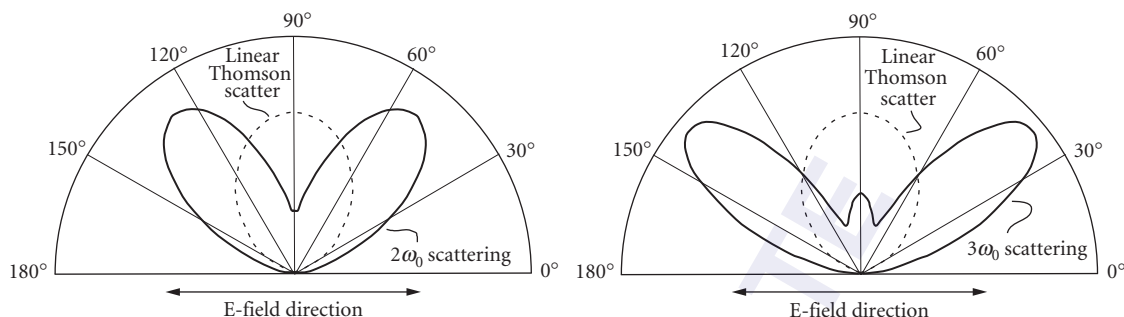


FIGURE 3 Azimuthal angular distribution (in the plane perpendicular to the laser propagation direction) of nonlinear relativistic Thomson scattering by an electron in a light field with $a_0 = 1$. The dipole distribution from linear Thomson scattering from an electron is shown for comparison. (These plots were adapted from Ref. 75.)

which, of course, reduces to the isotropic polar emission of Thomson scattered light in a weak circularly polarized field.⁷²

High-Field Interactions with Relativistic Electron Beams

The previous discussion considered the interaction of intense light radiation with electrons at rest or nearly at rest initially in the laboratory frame. If the laser collides with electrons that are already relativistic, which occurs when an intense laser pulse interacts with a beam of electrons from a high energy accelerator or synchrotron, the scattered radiation is altered by the fact that the electron sees a laser photon whose energy is upshifted by a factor γ , the Lorentz factor of the relativistic electron beam. If the scattering is linear Thomson scattering, the scattered photon will acquire another factor of γ in its energy when transformed back into the laboratory frame. This γ^2 upshift in photon energy can be exploited to produce femtosecond pulses in the x-ray regime by colliding an ultrashort pulse with a relativistic electron beam.⁷⁸ The scattered photon will have photon energy given by

$$\hbar\omega_{\text{scat}} = 2\gamma^2\hbar\omega_0 \frac{1 - \cos\phi}{1 + \gamma^2\theta^2} \quad (12)$$

where ϕ is the angle between the laser and the electron beam and θ is the angle of the scattered photon with respect to the electron propagation direction. This indicates that if a laser is scattered from the electron beam at a 180° angle, the photons scattered can be upshifted by as much as $4\gamma^2$. Furthermore, Eq. (12) indicates that the scattered photons will be emitted in a directed cone with angle $\sim 1/\gamma$, with an angular dependence on the upshifted photon energy. This picture must be amended somewhat if the electron beam Lorentz factor is high enough that the laser photon in the electron frame is upshifted in the electron rest frame such that $\gamma\hbar\omega_0/m_e c^2 \sim 1$. In this case, the situation becomes that of inverse Compton scattering and the kinematics of the electron recoil from the photon scattering must be considered. This process is, strictly speaking, a linear process; however, practical experimental implementation of this technique has usually been in the high intensity laser regime because of the low scattering cross section of free electrons ($\sigma_T = 6.6 \times 10^{-25} \text{ cm}^2$).

The situation becomes more complex when the laser is intense enough to cause multiphoton Compton scattering, whose scattering efficiency will then scale as a_0^{2n} , where n is the multiphoton order.⁶⁷ Accessing this regime in the lab is difficult because of the extremely low cross section but is made easier with a very high energy electron beam because an intense laser will have its intensity boosted in the electron frame through relativistic time compression.⁷⁹ At relativistic laser intensity

($a_0 \geq 1$) Eq. (12) must be amended, and the maximum scattered photon energy for head on collision and direct backscatter becomes

$$\hbar\omega_{\text{scat}} = \frac{4n\hbar\omega_0\gamma^2}{1 + \frac{4n\hbar\omega_0\gamma^2}{m_e c^2} + a_0^2} \quad (13)$$

The factor of a_0^2 in the denominator arises from the mass shift of the electron in the strong laser field.

21.5 STRONG FIELD INTERACTIONS WITH ATOMS

Keldysh Parameter and Transition from the Multiphoton to the Quasi-Classical Regime

Perhaps the most fundamental process that occurs when a single atom or ion is immersed in a strong laser field is the ionization of the most weakly bound electron. With the exception of recollision double ionization, discussed below, this ionization process is almost always a single electron process, involving the removal of the outermost bound electron by the light field.¹⁸ (This approach to understanding ionization and nonlinear optical dynamics in strong field atomic interactions is termed the “single active electron approximation,” and it underlies most of the theory presented in this section.) High-field interactions with single atoms essentially split into two regimes. The first occurs when the field can be treated quantum mechanically as an ensemble of photons, and the second arises when so many photons participate that the light can be treated as a classical field. In the second case, which occurs at sufficiently long wavelength or high intensity, the motion of electrons in the field can be treated classically. Generally, the second situation arises if the free electron wavepacket is much smaller than its classical oscillation amplitude. In this case a free electron wavepacket is localized to the extent that it can be considered a point particle; in strong field physics this is called the quasi-classical regime. The uncertainty principle implies that this occurs for a free electron when $U_p/\hbar\omega \gg 1$.

In the context of atomic ionization, these two pictures of the liberated free electron naturally lead to two regimes of ionization, the multiphoton regime and the semiclassical tunneling regime. These two regimes of ionization can be quantitatively differentiated by the Keldysh parameter²

$$\gamma_K = \sqrt{\frac{I_p}{2U_p}} \quad (14)$$

where I_p is the ionization potential of the atom or ion to be ionized. This quantity can be physically thought of as the ratio of the time it takes for an electron wavepacket to tunnel through the potential barrier of an ion immersed in a uniform electric field (see Fig. 5) to the period of the light oscillation. The Keldysh parameter delineates the barrier between the multiphoton ionization regime, which occurs when $\gamma_K \gg 1$, and the tunneling regime, which is the predominant ionization mechanism when $\gamma_K \ll 1$. Note that, given the scaling of U_p from Eq. (1), the latter situation predominates at higher intensity and longer wavelength. In practice, γ_K is a “soft” parameter in which tunneling ionization (described below) can be considered to be a very good approximation, even when γ_K is only slightly less than 1.

Multiphoton Ionization

When the Keldysh parameter is larger than 1, the ionization of an atom or ion in a strong laser field takes on a predominantly multiphoton character. The atom or ion can be thought of as absorbing a number of discrete photons from the field simultaneously, so that a bound electron acquires enough

energy to be promoted to the continuum.^{80–86} This view tends to predominate practically at modest intensities ($\sim 10^{13}$ W/cm²)^{82,85} or with light at wavelengths shorter than optical (UV).^{81,84} The most straightforward picture for this process is to say that the ionization rate for N -photon ionization is given by lowest order perturbation theory^{4,87,88} so that the ionization rate can be written as

$$W_N = \sigma_N I^N / (\hbar \omega_0)^N \quad (15)$$

where σ_N is the generalized multiphoton ionization (MPI) cross section. Equation (15) indicates that the ionization rate can be extremely nonlinear with laser intensity. For example, $N = 22$ when helium is ionized by a 1- μm wavelength laser,⁸⁹ though this multiphoton picture turns out to be accurate only for lower ionization potential atoms irradiated at shorter wavelengths. The practical difficulty with this model is in the calculation of σ_N . Lowest order perturbation theory is accurate for a very limited range of intensities and higher order processes soon become important in the calculation of σ_N as intensity is increased. Figure 4 plots the order of magnitude of the generalized cross section as a function of multiphoton order, N . This plot shows, for example, that the multiphoton ionization cross section for 11 photon ionization of Xe by a 1- μm laser field is about 10^{-350} cm²² – s¹⁰. This implies a saturation intensity for multiphoton ionization of atomic Xe by 100 fs pulses of about 2×10^{14} W/cm².

Calculation in the lowest order perturbation theory framework of σ_N is complicated by other factors as well. Resonances with intermediate states complicate the calculation and drastically affect the multiphoton ionization rate. Furthermore, when the field's ponderomotive potential becomes a significant fraction of the ionization potential, the bound state levels can no longer be thought of as unperturbed ion eigenstates; these levels move in energy via the AC Stark shift.⁹⁰ Also, as the ponderomotive potential increases, the minimum number of photons needed for ionization may change. Since the ionized electron is “born” into the continuum which has an oscillating field, it must acquire energy equal to U_p to enter the continuum. As a result, the number of photons needed for ionization is that which overcomes the ionization potential of the unperturbed atom or ion plus the ponderomotive energy, such that $N\hbar\omega_0 = I_p + U_p$.

A great number of nonperturbative approaches have been developed to derive more accurate ionization rates in a strong laser field in the multiphoton regime. One of the best known such approaches derives from the original work of Keldysh.² This approach has been extensively

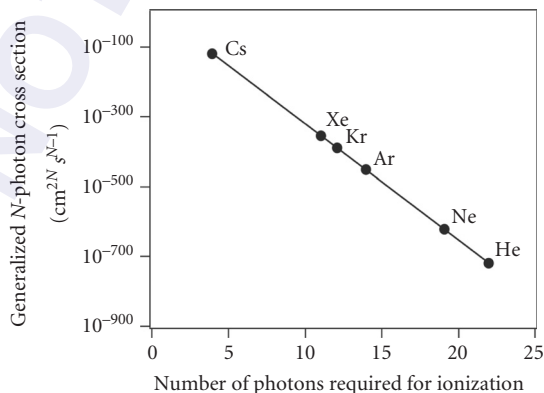


FIGURE 4 Generalized multiphoton cross section for atoms as a function of the multiphoton order, N . The location of the cross section for ionization of various neutral atoms by a laser at a wavelength of 1 μm are plotted to illustrate the increasing nonlinearity of MPI for more tightly bound atoms. (This plot was adapted from Refs. 80 and 21.)

developed by F. Faisal⁹¹ and H. Reiss;⁹² ionization calculations based on this approach are usually termed KFR theories or the strong field approximation (SFA). In this approximation, the ionization rate is determined by calculating a quantum mechanical transition probability directly between the ground state of the atom or ion and a continuum state. The ground state is assumed to be the field-free eigenstate unperturbed by the field (accurate for a tightly bound electron) and the final continuum state is assumed to be that of a free electron wavefunction in a plane electromagnetic wave (termed a Volkov state), thereby ignoring the effect of the ion Coulomb field on the outgoing ionized electron. A general equation valid for ionization rate for any Keldysh parameter can be derived in this way. In the multiphoton regime ($\gamma_K \gg 1$) the ionization rate predicted by this theory is²

$$W_K = A \omega_0 \left(\frac{I_p}{\hbar \omega_0} \right)^{3/2} \exp \left[2N_{\text{eff}} - \frac{I_p^{\text{eff}}}{\hbar \omega} \left(1 + \frac{2U_p}{I_p} \right) \right] \left(\frac{U_p}{2I_p} \right)^{N_{\text{eff}}} \Phi \left[\left(2N_{\text{eff}} - \frac{2I_p^{\text{eff}}}{\hbar \omega_0} \right)^{1/2} \right] \quad (16)$$

where $I_p^{\text{eff}} = I_p + U_p$ is the effective ionization potential of the atom dressed by the light's ponderomotive potential, N_{eff} is the minimum number of photons required to ionize the ion with this I_p^{eff} , $\Phi[z] = \int \exp[y^2 - z^2] dy$ is the probability integral and A is a numerical cofactor of the order of unity that accounts for the weak dependence on the details of the atom. Experiments have illustrated, for example, in 580 nm light that $A = 24$ for the first ionization of Ar, $A = 18$ for Kr and $A = 4$ for Xe.⁸³ While Eq. (16) is not particularly accurate for most ions and has a rather limited range of applicability, it is useful for estimating the order of magnitude of the ionization rate when $\gamma_K > 1$.

Tunnel Ionization

When the laser field is strong enough and the laser frequency is not too high, γ_K becomes less than 1 and a different picture of strong field ionization emerges. In this regime it is accurate to think of the bound electron wavepacket as evolving in a binding potential that is distorted by the strong light field, a situation known as tunnel ionization, illustrated in Fig. 5. The laser field represents a slowly varying deformation of the ion's confining Coulomb potential which oscillates back and forth. Near the peak of the light field oscillation, the electron can tunnel through the confining potential (pictured at the right in Fig. 5) freeing it and releasing it into the continuum. Because of the exponential nature of the quantum mechanical tunneling rate through a potential barrier, this

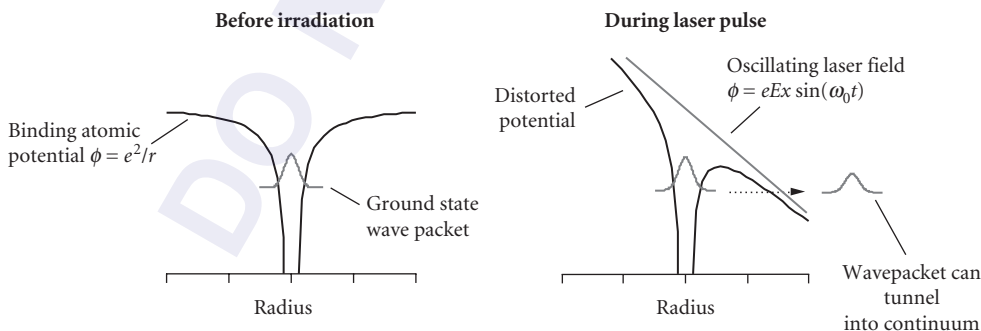


FIGURE 5 Illustration of the potential of an ion distorted by the application of a strong, adiabatically varying electric field. A strong enough field allows tunneling of the bound electron into the continuum thereby ionizing the atom/ion.

method of ionization is strongly nonlinear with increasing electric field. In a Coulomb potential this confining barrier will have a width roughly $\delta r \approx I_p/eE_0$ so the electron will tunnel through this barrier with a time

$$\tau_{\text{tun}} \approx \delta r / v \approx \frac{\sqrt{2I_p m_e}}{eE_0} \quad (17)$$

When this time is faster than a laser oscillation cycle, the tunneling picture is valid (equivalent to $\gamma_K < 1$).

The instantaneous ionization rate from tunneling by an electron from a hydrogenlike ion in a quasi-static field is given by⁹³

$$W_{\text{H-like}} = 4\omega_a \left(\frac{I_p}{I_H} \right)^{5/2} \frac{E_a}{E(t)} \exp \left[-\frac{2}{3} \left(\frac{I_p}{I_H} \right)^{3/2} \frac{E_a}{E(t)} \right] \quad (18)$$

where I_H is the ionization potential of hydrogen (13.6 eV), $\omega_a = 4.13 \times 10^{16} \text{ s}^{-1}$ is the atomic unit of frequency, $E_a = 5.14 \times 10^9 \text{ V/cm}$ is the atomic unit of electric field and $E(t)$ is the instantaneous applied electric field strength. The total ionization rate can be found by integrating Eq. (18) over the entire optical cycle.

There have been many published improvements on this simple tunneling formula.^{34,94,95} In fact, the general equation derived by Keldysh retrieves a tunneling rate when it is taken in the limit that $\gamma_K \ll 1$.² All tunneling models result in an ionization rate with the exponential field dependence $\sim \exp[-2a/3E(t)]$, with “ a ” depending on the model. The most sophisticated and most widely accepted model for the tunnel ionization of a complex atom or ion (i.e., multielectron nonhydrogenic ion) was developed by Ammosov, Delone, and Krainov, usually termed the ADK ionization rate.⁹⁵ This model predicts that the tunnel ionization rate averaged over one full optical cycle is

$$W_{\text{ADK}} = \omega_a C_{nl} f(l, m) \frac{3^{1/2}}{2\pi^{1/2}} \left(\frac{I_p}{I_H} \right)^{1/4} \left(\frac{E_0}{E_a} \right)^{1/2} \left[2 \frac{E_a}{E_0} \left(\frac{I_p}{I_H} \right)^{3/2} \right]^{2n-|m|-1} \exp \left[-\frac{2}{3} \left(\frac{I_p}{I_H} \right)^{3/2} \frac{E_a}{E_0} \right] \quad (19a)$$

This formula accounts for atomic structure through the principal orbital quantum number n , the orbital angular momentum l , and individual magnetic quantum number m . The atom-dependent cofactors are

$$C_{nl} = (2 \exp[1/n])^n (2\pi n)^{-1/2} \quad (19b)$$

$$f(l, m) = \frac{(2l+1)(l+|m|)!}{2^{|m|} (|m|)! (l-|m|)!} \quad (19c)$$

The usual method of deriving a net ionization rate is to sum over all m states of the ionizing ion. The ADK ionization formula has been found experimentally to be quite accurate over a wide range of intensities and ionic species.⁹⁶

This cycle averaged tunnel ionization rate is very nonlinear. It usually exhibits an intensity dependence that varies as I^6 to I^9 . As a result, strong field ionization tends to exhibit a thresholdlike behavior that quickly saturates once the intensity rises slightly above a threshold value. (Saturation in this context means that $W \tau_p \approx 1$, where τ_p is the laser pulse duration.) This threshold intensity can be easily estimated by determining when the field becomes high enough to suppress completely the confining Coulomb potential and the electron can freely escape the ion during the peak of the field

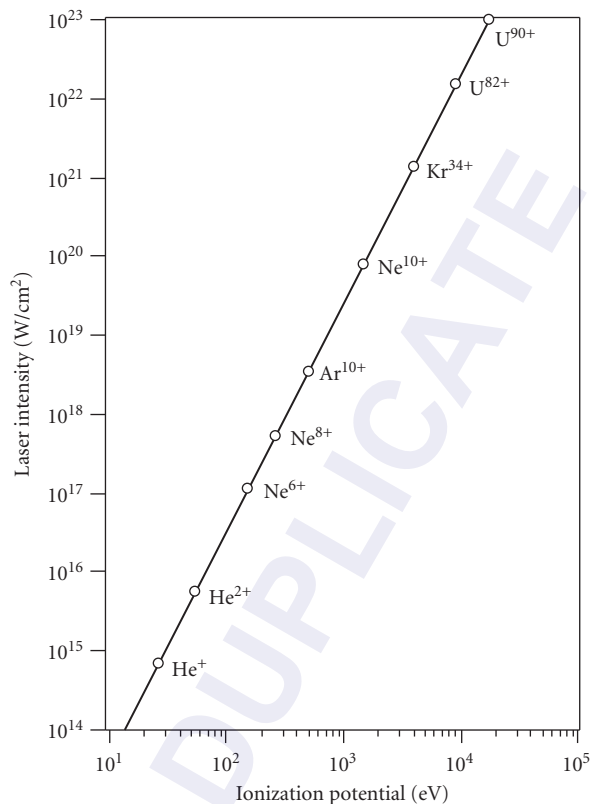


FIGURE 6 Plot of the Barrier Suppression Ionization (BSI) intensity for various ions.

cycle. This simple model, known as the Barrier Suppression Ionization (BSI) model,⁹⁶ indicates that the ionization threshold intensity occurs at

$$I_{\text{BSI}} = \frac{cI_p^4}{128\pi Z^2 e^6} \quad (20)$$

where Z is the charge state of the ion being created by the ionization event. In practical units the BSI intensity is $I_{\text{BSI}} [\text{W}/\text{cm}^2] = 4.0 \times 10^9 I_p^4 [eV] Z^{-2}$. Equation (20) turns out to be remarkably accurate despite its simplicity in predicting the intensity of significant ionization in a strong field. The predicted BSI threshold intensity for a variety of ion species is plotted as a function of ionization potential in Fig. 6, illustrating the extent to which very high charge states can be produced by tunneling with modern high intensity lasers.

Above Threshold Ionization

When an atom or ion is subject to a strong ionizing laser field an electron acquires some kinetic energy upon ionization.^{3,17,90,97-102} In the weak field regime, the energy acquired by an electron in the light field upon multiphoton ionization can be thought of as a simple extension of the photoelectric effect, in which the electron gains an energy $\text{KE}_{e^-} = N\hbar\omega_0 - I_p$, where again, N is the minimum number of photons needed for ionization.¹⁷ However, as the intensity increases and becomes nonperturbative, the

electron can absorb more than the minimum number of photons needed for ionization, a consequence of the electron remaining in the vicinity of the nucleus long enough during ionization to absorb more than N photons. The electron now leaves the ion with energy

$$\text{KE}_{e^-} = (N+s)\hbar\omega_0 - I_p^{\text{eff}} \quad (21)$$

where the ionization potential is dressed (in other words, shifted in a time-averaged sense) by the ponderomotive potential [as in Eq. (16)], and s is the number of additional absorbed laser photons. This effect is termed above threshold ionization (ATI) and leads to electron kinetic energy distributions that look qualitatively like those in Fig. 7. These electrons are typically ejected from the focus along the laser polarization direction (at least in the nonrelativistic limit).¹⁰³ The absorption of s additional photons leads to electron energy peaks separated by $\hbar\omega_0$, shifted down by an energy I_p^{eff} (though ponderomotive acceleration in the focus of a long pulse laser can shift these electron energies back to the undressed energy).

At modest field strengths in the multiphoton regime ($\sim 10^{13}$ W/cm² at IR wavelengths), the yield of electrons in each peak drops off exponentially, as predicted by lowest order perturbation theory. At nonperturbative intensities, however, the yield in each electron ATI peak flattens out and the ATI spectrum develops a plateau over the first few orders, out to an energy of $\sim U_p$, (illustrated in Fig. 7).¹⁰⁰ There have been many theories published to explain the quantitative behavior of ATI in

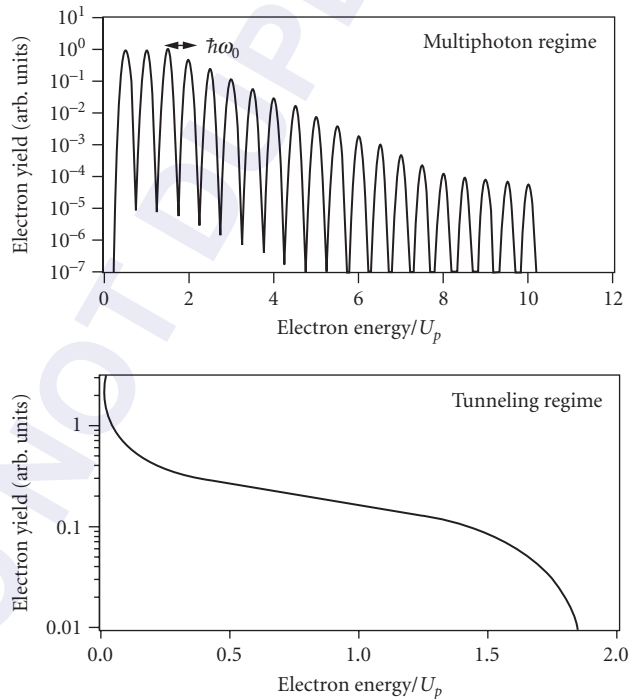


FIGURE 7 Illustration of the general nature of electron energy spectra resulting from intense ionization of atoms, manifesting the phenomenon known as ATI. At the top the spectrum common for ionization in the multiphoton regime is shown. Here the spectrum is composed of peaks, each corresponding to electrons that have absorbed s photons above that needed for ionization. On the bottom, a spectrum characteristic of ATI in the tunneling regime is illustrated showing that the spectrum is smooth and extends out to $2 U_p$.

the multiphoton regime. In fact, the strong field approximation of the KFR theories predicts the presence of these multiphoton ATI peaks. Resonances often play an important role in ATI in the multiphoton regime,¹⁰⁴ leading to a rich variety of effects that appear in the ATI energy spectra of optical and near-IR pulses with intensity 10^{13} to 10^{15} W/cm².

As the laser field increases, the character of the electron ATI spectrum takes on a qualitatively different character. At intensities entering the quasi-classical tunneling regime ($\gamma_K < 1$), the spectrum loses its multiphoton character composed of distinct electron energy peaks and becomes a smooth, monotonically decreasing energy distribution, illustrated in Fig. 7. The electron energies have kinetic energy predominantly between 0 and $2U_p$, though there is a small component of electrons with energies that reach up to as much as $10U_p$ through a rescattering process described below.¹⁰² This loss of distinct peaks occurs when the field loses its quantum character and drives the ionization as if it were a classical field.^{97,98}

Quasi-Classical ATI When ions are ionized in the tunneling regime, which tends to occur in optical and near-IR laser intensities of 10^{15} to 10^{18} W/cm², the shape and energy of electron ATI can be determined by a simple model. As with tunnel and BSI ionization, the field is treated classically and the electron is treated as a compact wavepacket that propagates in the continuum subsequent to tunnel ionization by the classical equations of motion.¹⁰⁵ This so-called quasi-classical model is useful in explaining not only ATI in the strong field, long wavelength regime (i.e., when $\gamma_K < 1$) but also aids in explaining a large number of other strong field phenomena.¹⁰ As such, the quasi-classical model has become one of the major building blocks for understanding modern strong field physics.

In the quasi-classical ATI model, the electron is considered to tunnel into the continuum at a well-defined phase in the oscillating field (see Fig. 8) and propagates in the field as a pointlike charge. When the electron is “born” in to the continuum in this way, it not only oscillates but also acquires some directed drift velocity in the direction of the laser’s polarization, with an energy that is a function of the phase in the field at which it was ionized. (This is true for linear polarized fields; in circularly polarized fields, the electron acquires constant drift velocity.) Solution of

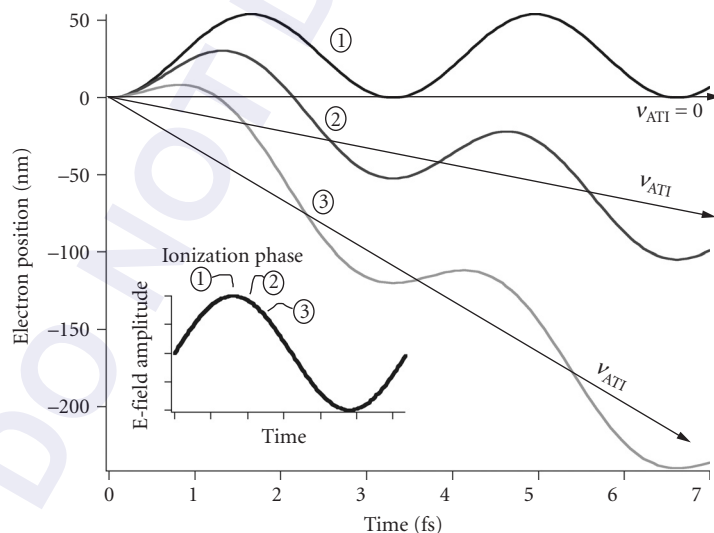


FIGURE 8 Plot of a calculation showing the classical trajectory of electrons born at three different phases in the laser’s field. These three phase locations are illustrated in the inset. When born at the peak of the field (trajectory 1), the electron oscillates but acquires no net drift velocity. However, when born off the peak, the electron acquires nonzero drift velocity on top of its quiver velocity.

the nonrelativistic equations of motion indicates that the drift component of the kinetic energy acquired by an electron is

$$\text{KE}_{\text{ATI}} = 2U_p \sin^2 \Delta\phi_0 \quad (22)$$

where $\Delta\phi_0$ is the phase at which the electron is “born” with respect to the peak of the field. Equation (22) indicates that if the electron is born at the peak of the electric field, where the tunnel ionization rate is highest, it will acquire no ATI drift energy. The electron can acquire up to $2U_p$ of kinetic ATI energy if it is ionized well off of the field peak. If Eq. (22) is combined with a model for static tunnel ionization, an electron energy distribution can be derived. Use of the simple H-like tunnel formula of Eq. (18) leads to an ATI energy spectrum from a linearly polarized field of¹⁰⁵

$$f(\varepsilon) = \frac{a}{(1 - \varepsilon/2U_p)(\varepsilon/2U_p)^{1/2}} \exp \left[-\frac{2}{3} \left(\frac{I_p}{I_H} \right)^{3/2} \left(\frac{E_a}{E_0} \right) \left(1 - \frac{\varepsilon}{2U_p} \right)^{-1/2} \right] \quad (23)$$

where a is just a normalization constant. The shape of this spectrum for tunnel ionization is illustrated in Fig. 7. This distribution is peaked near zero energy (because the ionization probability is greatest at the peak of the field where $\Delta\phi_0 = 0$ and $\text{KE}_{\text{ATI}} = 0$) but stretches out to a maximum electron energy of $2U_p$. In linearly polarized light, as a rule of thumb, the average electron energy is roughly 10 percent of U_p at the intensity where the ionization rate saturates.¹⁰⁶

Rescattering Effects

The quasi-classical picture of strong field ionization can be extended to explain other observed effects. In this picture once an electron is liberated from its binding potential by tunneling, its motion in the field can be described classically. Once freed into the continuum, the electron can recollide with its parent ion if it is ionized within a certain range of phase of the field,^{21,107–116} a process illustrated in Fig. 9. That this happens can be seen by examining the trajectories shown in Fig. 8, in which the electrons return to the $x = 0$ position in all three cases after ionization at $t = 0$ (though with different return energy in each case).

If the field is linearly polarized and the electron is born by tunnel ionization at some time t_0 and $x = 0$ (the location of the parent ion) in an oscillating electric field $E_0 \cos\omega_0 t$, its position as a function of time is

$$x(t) = -\frac{eE_0}{m\omega_0^2} [\cos\omega_0 t + \omega_0 t \sin\omega_0 t_0 - \cos\omega_0 t_0 - \omega_0 t_0 \sin\omega_0 t_0] \quad (24)$$

This classical trajectory will result in the recollision of the tunnel ionized electron with its parent nuclear core for ionization phases of $\omega_0 t_0 = 0^\circ$ to 90° and 180° to 270° . The highest energy recollisions occur at phases of 17° and 197° . The free electron will have a kinetic energy after it is born by ionization given by

$$\text{KE}(t) = \frac{e^2 E_0^2}{2m_e \omega_0^2} [\sin\omega_0 t - \sin\omega_0 t_0]^2 \quad (25)$$

The energy upon return to the nucleus can be found with Eqs. (24) and (25). The maximum energy an electron can have upon returning to its parent nucleus (corresponding to ionization at 17° and 197°) is $3.17U_p$. The recollision energy and tunnel ionization probability as a function of phase in the laser field is illustrated in Fig. 9. This strong field driven electron recollision can manifest itself in a number of important effects. The generation of short wavelength harmonic radiation in this way is discussed in Sec. 21.7, but there are two other consequences of this effect:

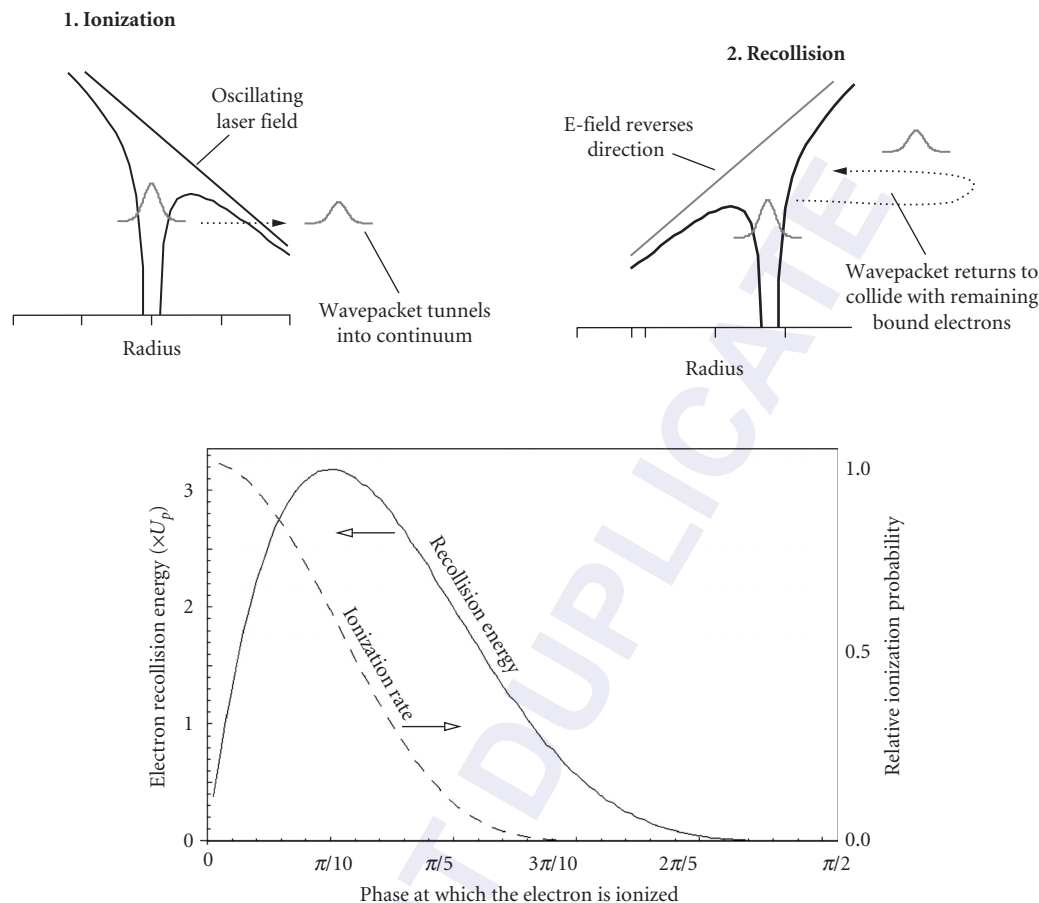


FIGURE 9 The upper illustrations show the process by which a tunnel ionized electron can return to its parent ion and interact with the bound electrons remaining in the ion. The bottom plot shows the energy that such a returning electron will have as a function of the phase in the field at which it was ionized. This shows that the maximum return energy can be $3.17U_p$. The relative ionization probability as a function of the field phase is also shown to illustrate how many electrons will return with the energies indicated.

Strong Field Double Ionization It has been observed for some time that strong field ionization exhibits a small but significant signature of simultaneous double ionization of ions where two electrons are apparently liberated at the same time in the field.^{112,115,117} This effect manifests itself experimentally in the yield of a given ion charge as a function of laser intensity, illustrated in Fig. 10. The “knee” that is seen in the ionization yield of a next higher charge state when the ionization of a lower charge state saturates is clear evidence of double ionization. In addition to rescattering, a number of ideas have been forwarded to explain this strong field behavior.²⁶ These include:

- *Shake off model*, in which the sudden removal of one electron by tunneling leads to a quantum mechanical relaxation of the second electron into a new set of eigenstates, some of which rest in the continuum and, therefore, lead to ionization.^{117,118} This effect is well known and documented in single photon photo-ionization.¹¹⁹
- *Collective tunneling*, in which two electrons tunnel out simultaneously from the ion.¹²⁰

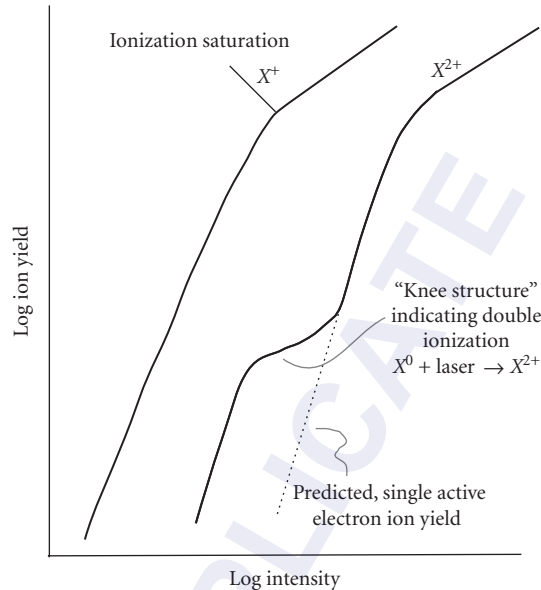


FIGURE 10 Generalized illustration of the ion yield typically measured for ions produced by tunnel ionization, showing the strongly nonlinear increase of ion yield with intensity followed by a roll over at the saturation intensity, where most ions become depleted. The “knee” structure seen in the second ion is the usual experimental signature for nonsequential double ionization.

- *Rescattering*, in which the recolliding electron, described above, collisionally ionizes a second electron. The recolliding electron can acquire sufficient kinetic energy in the field to ionize a second electron through collisional ionization on its return. This is presently the best accepted explanation for strong field double ionization. For example, the double ionization yield is seen to drop dramatically when circularly polarized light is employed.¹²¹ This indicates that the recollision mechanism (at least in the tunneling regime) is likely the dominant mechanism because an electron ionized in circularly polarized light will propagate in such a way that it will not return to the parent ion and cannot ionize a second electron.

ATI Plateau Extension A second consequence of the strong field driven recollision is manifested in the electron ATI energy spectrum. As discussed above, ATI electrons in the quasi-classical regime can acquire up to $2U_p$ of drift energy. However, a laser field-driven recollision can lead to a small number of scattered electrons with energy up to $10U_p$.¹⁰² This rescattering leads to the production of an electron ATI spectrum with a large predominance of electrons with energy below $2U_p$ but with a small fraction of electrons with an energy spectrum plateau that extends out to $10U_p$.

Relativistic Effects

When $a_0 \rightarrow 1$ strong field ionization dynamics are altered in a number of ways.⁴⁰ The most significant effects arise from the forward force exerted by the laser magnetic field, though the relativistic mass increase of the electron does play a role in certain effects.

Relativistic Tunnel Ionization It turns out that tunnel ionization rates at relativistic intensity do not deviate significantly from the nonrelativistic rates.¹²² The most significant relativistic effects in tunneling occur when the bound state energy I_p becomes comparable to the rest energy $m_e c^2$; it is the Coulomb correction to the mass of the electron in the ground state that is the principal effect. Therefore, relativistic effects will be important for ions of charge $Z > \sim 50$. A relativistic generalization of the Keldysh theory,³⁴ indicates that the relativistically corrected tunnel ionization rate, W_{Rel} will be higher than the nonrelativistic rate, $W_{\text{non-Rel}}$ by a factor

$$\frac{W_{\text{Rel}}}{W_{\text{non-Rel}}} \approx \exp\left[-\frac{1}{36}(Z\alpha)^5 \frac{E_{\text{cr}}}{E_0}\right] \quad (26)$$

where α is the fine structure constant (1/137) and E_{cr} is the Schwinger critical field from quantum electro-dynamics theory (1.3×10^{16} V/cm). Using the barrier suppression ionization model to estimate the appropriate Z , Eq. (26) suggests that the nonrelativistic tunneling rates should be good up to an intensity of $\sim 10^{26}$ W/cm².

Relativistic Electron ATI The ejection of free electrons following tunnel ionization in a relativistic field can be influenced by the laser's magnetic field. At nonrelativistic intensities, electrons are ejected in a rather narrow distribution along the laser polarization axis. As a_0 approaches 1, the magnetic field pushes the electron distribution toward the k direction of the light propagation, a consequence predicted by Eq. (8).⁶⁸ This distribution shift is illustrated in Fig. 11. At highly relativistic intensities (i.e., when $a_0 \gg 1$, an intensity of $> 10^{21}$ W/cm² at near-IR wavelengths) tunnel ionized electrons are quickly bent toward the laser propagation axis by the magnetic field. The electron, which will have a velocity near c , will then “surf” along with the laser pulse acquiring energy from the laser field. Such electrons will be ejected from the laser focus in a narrow cone along the laser direction and will acquire many MeV or even GeV of energy.^{71,123} An energy versus angle distribution of electrons ejected in this regime is illustrated in Fig. 11.

Relativistic Suppression of Rescattering Another consequence of the forward ejection of electrons in a relativistic light pulse is that the nonsequential double ionization that normally accompanies rescattering of the electrons on their return after ionization is suppressed.^{40,124–126} The forward motion of the electron imparted by the magnetic field (see Fig. 2) forces the electron away from the core and prevents the collisional ionization of a second electron. The fall-off of nonsequential ionization yield occurs at intensity as low as 10^{17} W/cm² ($a_0 \approx 0.3$). The rescattering plateau in ATI spectra associated with electrons with energy in the $2U_p$ to $10U_p$ range also decreases in magnitude because of the rescattering suppression. Furthermore, this phenomenon leads to suppression of single atom high harmonic generation at high intensity.

Ionization Stabilization

While experimental evidence is scant, there is a strong theoretical basis for believing that, in certain situations, the ionization rate of an atom in a strong field actually declines with increasing intensity.²⁷ This phenomena has come to be called ionization stabilization. There are usually two manifestations of this stabilization discussed in the literature.

Adiabatic Stabilization Quantum mechanical calculations of ionization rates have shown that the ionization rate can be stabilized in unionized atoms at field strengths well above one atomic unit. A simple picture to explain this can be constructed if one considers that the electron wavefunction in the ground state of the atom is considerably altered by the strong oscillating field. The wave function is thought to evolve into a time averaged structure with peaks away

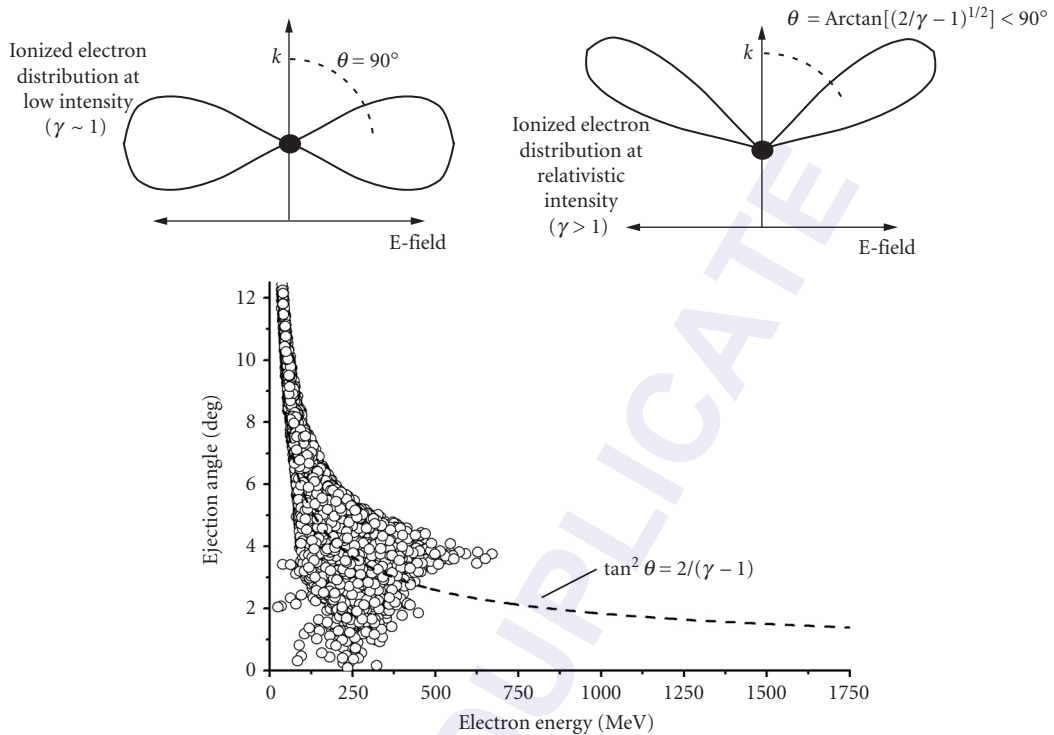


FIGURE 11 The upper illustrations are a generalized illustration of the angular distribution of electrons produced during ionization with respect to the laser field and propagation directions. In the medium intensity regime, electrons are ejected by tunnel ionization along the E-field direction (left); however, at relativistic intensity, the magnetic field pushes the distribution toward the direction of laser propagation. The latter effect can be understood as the conservation of momentum following the absorption of an extremely large number ($\sim 10^6$) of photons. The bottom plot (*adapted from Ref. 71*) shows the results of a Monte Carlo simulation yielding the ejection angle and energy of electrons produced by irradiation of an Ar ion at intensity of 5×10^{21} W/cm². Here the dashed line is the prediction of Eq. (8); the deviation of the simulation from this equation is a consequence of longitudinal fields at the focus.

from the nuclear center. In a time averaged sense, the bound electron sees *two centers* for the atomic potential.^{27,127–129} This deformation of the electron wave function away from the nucleus lowers the ionization rate at higher intensity. Figure 12 illustrates the calculated ionization rates of an electron in an excited state of hydrogen as a function of intensity illustrating the fall of ionization rate at high I .

Dynamic Stabilization This mechanism of ionization stabilization, often termed interference stabilization, arises most prominently in calculations of strong field ionization of Rydberg atoms.^{129–131} It is thought to arise from quantum destructive interference of pathways into the continuum. There has been some experimental evidence for this form of stabilization in Rydberg atoms¹³² but has yet to be demonstrated in atoms in the ground state.

Numerical simulations have suggested that the magnitude of ionization stabilization decreases at relativistic intensity due to the effects of the magnetic field and the Lorentz force on the electron.¹³³

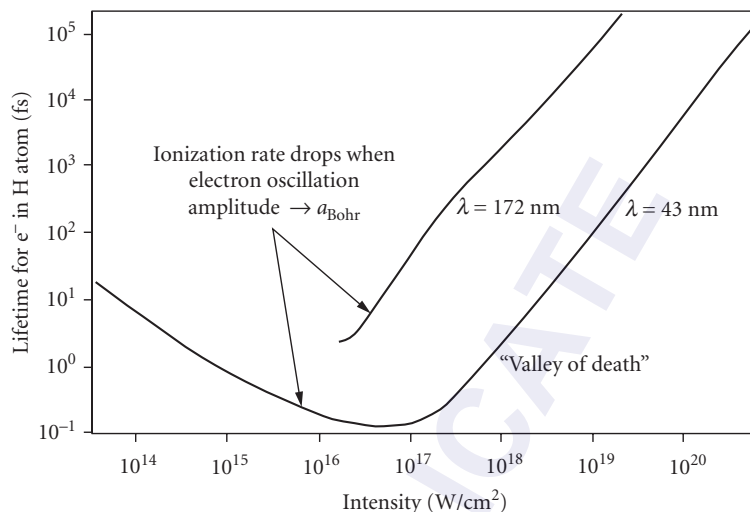


FIGURE 12 Calculation of the ionization rate of a hydrogen atom in an intense, short wavelength field. (Adapted from Ref. 27.) This calculation shows that at around the intensity at which the electron quiver amplitude is comparable to the size of the hydrogen ground state (~ 1 Bohr radius) the ionization rate actually decreases (lifetime increases) as the intensity increases from the delocalization of the electron wavefunction.

21.6 STRONG FIELD INTERACTIONS WITH MOLECULES

The ionization of small molecules (of less than ~ 10 atoms) by an intense laser field is qualitatively similar to the ionization of single atoms. The two limiting regimes for ionization (multiphoton and tunneling) as discussed in the previous section are still relevant for molecules and the tunnel ionization rate formulas tend to work reasonably well in predicting the ionization rate of electrons in a molecule if appropriately chosen ionization potentials are utilized.¹³⁴ Unlike single atoms, however, the motion of the molecule's nuclei during the interaction with the laser pulse can affect the dynamics of the electron ionization and energy gain from the field. Fragmentation of the molecule is one significant consequence of irradiation at high intensity, and the motion of the molecular nuclei has a dynamic impact on the structure of the molecule during its interaction with the intense laser pulse.

Nuclear Motion and Molecular Alignment in Strong Fields

Because small molecules tend to fragment rapidly in an intense light field, the regime of strong field laser interactions with these molecules is usually limited to rather modest intensity, below about 10^{15} W/cm². Higher intensity pulses tend to destroy a small molecule well before the high intensity can be reached. Much of current research has concentrated on diatomic molecules.

At modest intensity a small molecule will experience a force by the light field which will tend to align it.^{135–137} In the absence of the light field, the molecules are randomly oriented and exist in a range of molecular rotation states with energy eigenvalues of $\epsilon_{\text{Rot}} = \beta_{\text{Rot}} J(J+1)$ and with rotational quantum number J . Some values of β_{Rot} are tabulated in Table 1.¹³⁸ When a moderately strong light

TABLE 1 Rotational Constant, and Alignment Well Depth for Three Example Molecules Irradiated at 10^{15} W/cm²

Molecule	β_{Rot} (meV)	Max Well Depth (meV)
H ₂	3.89	21.9
N ₂	0.25	96.8
CO ₂	0.048	212

Source: Table Adapted from Ref. 138.

field is applied (below the intensity at which the molecule ionizes), the induced dipole causes the molecule to see a cycle-averaged potential given by¹³⁸

$$V(\theta) = -\frac{I(t)}{2c} [(P_{\parallel} - P_{\perp}) \cos^2(\theta) + P_{\perp}] \quad (27)$$

where the P terms are the parallel and perpendicular components of the polarizability tensor, θ is the angle between the molecular axis and the polarization axis and $I(t)$ is the time dependent intensity. The maximum well depth for some molecules in a field of intensity 10^{15} W/cm² is tabulated in Table 1.

The potential of Eq. (27) will tend to align a linear molecule, such as a diatomic along the polarization axis of a linearly polarized field. Because molecules will usually feel a lower intensity field as the pulse of an intense laser ramps up in time, it is usually a good approximation to say that the nuclei of a linear molecule will be (partially) aligned along the laser electric field at subsequent higher intensity. A molecule with nonzero rotational momentum will evolve, classically speaking, in a pendulumlike motion around the laser electric field axis; these states are often called “pendular” states.

Even at modest intensity, the molecule will begin to dissociate via a process known as bond softening.^{35,139,140} In a small diatomic molecule, such as H₂ the first electron will be ionized at the equilibrium distance of the two atoms via multiphoton or tunnel ionization. If the molecule is indeed aligned along the light electric field, the molecular nuclei will then begin to separate by bond softening. Molecular dissociation begins to occur when the potential that binds the nuclear wavepackets couples to photons in the strong laser field. This coupling leads to a ladder of potential curves, each shifted by one photon in energy (described by what is known as Floquet theory^{141,142}). If the field-dressed states are treated as if they are quasi-static, the Hamiltonian of the molecule can then be diagonalized, distorting the bound potential curves, in a manner illustrated in Fig. 13 (these distorted potentials are termed adiabatic potentials). As can be seen in this figure, the distorted curves allow molecules in excited vibrational states to dissociate, (and some in lower states can tunnel through the distorted potential barrier). This potential distortion is termed bond softening and is the main mechanism by which a molecule begins to fragment as an intense laser is ramped up in intensity. In H₂ this bond softening occurs, for example, at intensity of $\sim 10^{13}$ W/cm².³⁵

There are other means by which molecules can dissociate in midstrength fields ($I < 10^{14}$ W/cm²). In the picture in which the field is treated as a classical, periodically varying alteration to the Hamiltonian of the molecule, the molecule sees states that are shifted down (“dressed”) by an energy equal to the photon energy. The nuclear wavepacket can couple from a bound state to a dissociating state and begin to separate. For example, in H₂⁺ immersed in a visible light field of intensity around 10^{13} W/cm², this occurs by coupling first to an unbound state shifted by $3\hbar\nu$ and then, after some expansion of the internuclear separation, coupling to the continuum of the bound state dressed by two photons. This multiphoton process is often termed above threshold dissociation^{143,144} because it results in the absorption of more photons than are energetically required to dissociate the molecule (much like above threshold ionization is absorption of more photons by an electron than needed to ionize).

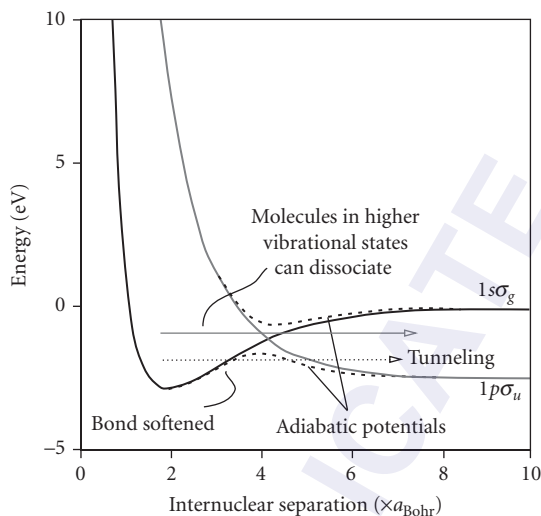


FIGURE 13 Nuclear potential energy curves for the first bound and first unbound state of H_2^+ . The coupling of these states to a strong field leads to mixing and new distorted potentials called “adiabatic states,” which are pictured as dashed curves. The bond softening resulting from these distorted curves permits molecules left in high lying vibrational states to dissociate immediately and those in lower states to tunnel out on a longer time scale.

Coulomb Explosion

At higher intensity (10^{14} to 10^{15} W/cm²), multiple ionization of the molecular nuclei will lead to a Coulomb repulsion of the ions, a process known as a Coulomb explosion. This explosion happens very quickly and subsequent multiple ionization of the nuclei while in proximity to each other occurs only if the rise time of the laser is comparable or faster than this Coulomb explosion separation time. The characteristic time for explosion of a diatomic is roughly

$$\tau_{\text{CE}} \approx \frac{3}{4} \frac{\mu^{1/2} R_0^{3/2}}{(q_1 q_2)^{1/2} e} \quad (28)$$

where $\mu = m_1 m_2 / (m_1 + m_2)$ is the reduced mass of the repelling system, q_1 and q_2 are the charge states of the ionized nuclei, and R_0 is the initial separation of the exploding nuclei. Equation (28) indicates that H_2^{2+} Coulomb explodes in under 1 fs, while N_2^{2+} explodes in about 5 fs.

Upon Coulomb explosion, the ejected ions acquire a kinetic energy just given by their initial Coulomb potential energy. For a diatomic molecule with ions charged by field ionization to q_1 and q_2 , the Coulomb explosion energy is

$$\mathcal{E}_{\text{CE}} = \frac{q_1 q_2 e^2}{R} \quad (29)$$

where R is the separation of the two ions at the point of ionization. In almost all cases, the observed Coulomb explosion energy of an exploding diatomic or other small molecule is less than that expected of an explosion from the molecules equilibrium distance (which is, for example, $1.4 a_B$ for H_2).¹⁴⁵ This arises because of pre-expansion of the nuclear separation from bond softening and Coulomb repulsion of lower charge states before final ionization. A good rule of thumb is that \mathcal{E}_{CE}

upon strong field ionization will be about 45 to 55 percent of the energy calculated from equilibrium distances.¹³⁸ For example, Coulomb explosion of H_2 usually leads to the ejection of protons with energy between 3 and 5 eV.¹⁴⁶ ϵ_{CE} is largely independent of pulse duration and wavelength.¹⁴⁷

Molecular Tunnel Ionization and the Critical Ionization Distance

The ionization rate for the first electrons of diatomic molecules in a strong laser field for low charge states in the tunneling regime (i.e., $\gamma \ll 1$) frequently follows standard atomic tunnel ionization rates, such as the ADK rate of Eq. (19), with the appropriate molecular ionization potential inserted into the formula.¹³⁴ Nonsequential recollision double ionization of diatomics also occurs in a manner similar to that of atoms.^{148,149} There are, however, some subtleties which manifest in the tunnel ionization rate of certain diatomic species; these deviations from standard tunneling formulas are thought to arise from destructive interference between electrons liberated from the two nuclear centers.¹⁵⁰

The production of higher charge states in a molecule under higher intensity irradiation and the subsequent Coulomb explosion of the highly charged fragments deviates significantly from the predictions of tunnel ionization of single atoms. In particular, charge states from field ionization of molecular nuclei tend to occur at an intensity much lower than that of the same charge state in an isolated ion. This occurs because of the presence of what is known as the critical ionization radius in the molecule.^{151–155} This effect can be explained simply in the quasi-classical tunneling/barrier suppression ionization model described above. It arises because the field-induced tunneling of an electron from the molecule can be aided by the presence of the second charged nucleus of the molecule near the first. This effect is illustrated for a diatomic molecule in Fig. 14, in which a diatomic molecule

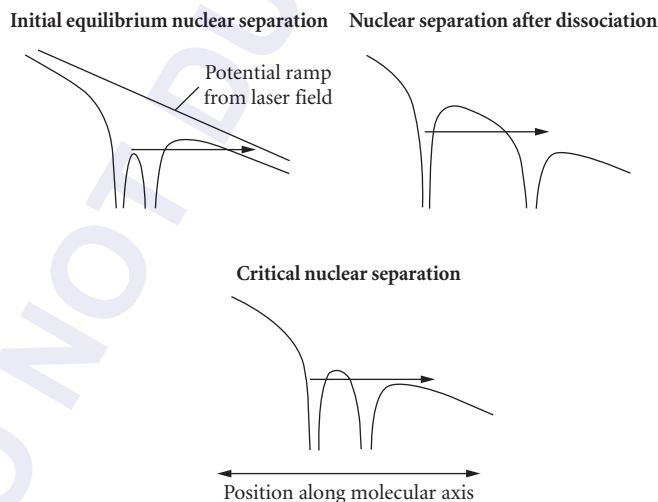


FIGURE 14 Drawing of the electron potential energy curves for a diatomic molecule immersed in a strong slowly varying field for three nuclear spacings. At the initial equilibrium position, pictured in the upper left, ionization by tunneling occurs when the electron tunnels out of the combined Coulomb well into the continuum. On the other hand, if the molecular nuclear have become greatly separated by dissociation, ionization must occur by tunneling from the potential of an isolated ion. However, at an intermediate separation, the electrons from the left-most ion need only tunnel through a narrow barrier formed from the combination of both Coulomb fields. This results in a greatly enhanced ionization rate at this nuclear separation.

is aligned along the axis of the laser electric field. When the nuclei are close together, the electrons are confined in the combined potential of the two nuclei; an applied field then requires an electron to tunnel from the combined potential. At very large nuclear separation, the electron is effectively localized at one of the nuclei and the tunneling rate is just that of an electron tunneling from an isolated Coulomb potential well. However, as the spacing of the diatomic nuclei begins to expand from the close initial position, the outer barrier (on the right in the figure) begins to drop and the central barrier begins to rise. There will occur an optimum distance in which an electron confined to the “uphill” potential need only tunnel through the relatively thin central potential barrier (see Fig. 14 bottom). This separation is the critical ionization distance. A simple analysis of the shape of the potential surrounding two ions each with charge q , indicates that this critical ionization distance will occur roughly at

$$R_c \approx \frac{4qe^2}{I_p} \quad (30)$$

If we make the approximation that the binding potential of charged ions in a diatomic are roughly given by the unionized molecular ionization potential divided by q , we find that the critical distance is independent of charge state. For a diatomic such as H_2 , the ionization potential of 15.4 eV suggests a critical ionization distance of about 3.7 Å (about 7 atomic units or 5 times the equilibrium separation). A similar calculation for I_2 suggests a critical ionization distance of 10 a.u. for all of the iodine charge states.¹⁴⁵

This ionization enhancement at the critical distance explains the appearance of higher charge states in molecular strong field ionization at lower intensity than might be expected.³⁵ The multiple ionization dynamics follow a multistep process. After some initial ionization of a high Z molecule, the ions begin to separate. If they come apart to the critical distance on a time scale faster than the laser pulse passes [predicted by Eq. (28) for most small molecules], the nuclei will be rapidly ionized to higher charge states at the critical distance, leading to an energetic Coulomb explosion of the higher charged ions from R_c separation. This process is described in Fig. 15¹⁵⁶ in which the barrier suppression ionization thresholds are plotted for various charge states of iodine as a function of I_2 separation along with the trajectory of an iodine molecule in a strong field.

There is some theoretical evidence for the presence of a second, more closely spaced critical ionization distance in molecule strong field ionization.¹⁵⁷ This effect, often called charge-resonance enhanced ionization (CREI), is a quantum mechanical effect and results from a localization of part of the electron wavepacket in the upper well at an internuclear separation smaller than R_c described above. There is presently no experimental evidence for this second critical distance.

Triatomic and Larger Molecules in Strong Fields

The fragmentation of molecules larger than a diatomic quickly becomes more difficult to describe, though the principles described in the last two subsections apply, particularly if the molecule is linear. However, the ionization dynamics are complicated and many fragmentation channels are usually observed at intensity $>10^{14}$ W/cm².¹³⁸ For example, even in a molecule as simple as CO_2 multiple channels such as $O^+ + CO^+ + 4.7$ eV and $O^{2+} + CO^+ + 8.5$ eV are observed with about equal probability when irradiated in near-IR pulses of intensity $\sim 10^{15}$ W/cm².¹³⁸ Many fragmentation channels occur not as a sequence of two independent bond fragmentations but occur as one, nonsequential rupture of both bonds.

Tunnel ionization rates for multinuclear molecules which are linear can often be found by treating the molecule as a single elongated well potential from which the electron can tunnel in the laser field.¹⁵⁸ These so-called structural tunnel ionization models have been employed successfully to explain ionization of more complex molecules, such as benzene,¹⁵⁹ though multielectron effects play a much larger role, and the single active electron approximation implicit in tunnel ionization theory is inaccurate.

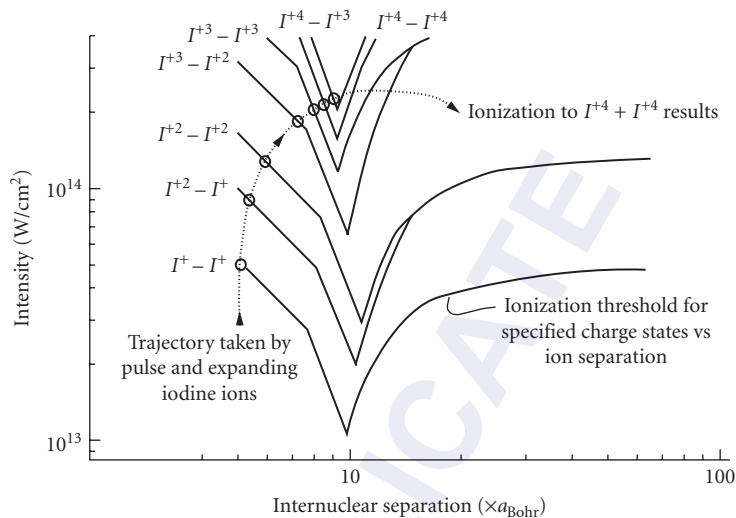


FIGURE 15 The solid lines in this plot show the BSI ionization threshold for various iodine ion pairs aligned along the laser field as a function of nuclear separation. The large dip in ionization threshold at around $10a_{\text{Bohr}}$ for all charge state pairs is a consequence of the critical ionization distance described in Fig. 14. The dashed line shows the trajectory in terms of laser intensity and nuclear separation that an iodine molecule undergoes when it is irradiated by a 150-fs pulse at an intensity of $2.5 \times 10^{14} \text{ W/cm}^2$ showing the various charge states created by BSI ionization in the pulse as the I_2 molecule Coulomb explodes. The presence of the critical ionization dip results in higher charge state (+4 for both ions) than would be achieved at this intensity by irradiation of a single iodine atom. (Plot adapted from that of Ref. 156.)

21.7 STRONG FIELD NONLINEAR OPTICS IN GASES

High Order Harmonic Generation

When a strong field pulse of light propagates through an extended gas a number of new phenomena occur. These occur when the gas density is high enough that coherent optical and nonlinear optical effects become important, a point that occurs in a practical sense in gases with density above about 10^{17} cm^{-3} . The single atom or single molecule interactions described above come into play and often are the seeds to the initiation of other physical effects, the most common of which is strong field ionization leading to plasma formation. However, another important strong field effect which manifests itself at intensity just at or slightly below the intensity threshold for ionization is the generation of radiation at harmonics of the laser field.^{7,24,38,160,161} This process, termed high order harmonic generation (HHG), is described in greater detail in other chapters of this volume so only a rudimentary overview is given here.

Harmonic generation in gases subject to strong laser fields arises from the nonlinear oscillations of the bound electrons in the gas's atoms. This is generally performed in gases of rare gas atoms, though HHG from molecular¹⁶² and clustering¹⁶³ gases is also possible. At modest intensity ($<10^{14} \text{ W/cm}^2$), HHG can be thought of as a multiphoton process, illustrated in the upper left of Fig. 16. In this picture an atom simultaneously absorbs q photons from the intense light field and then re-emits a photon with energy $q\hbar\omega_0$. Since free atoms in a gas exhibit a centro-symmetric potential, angular momentum must be conserved. This constraint, combined with the fact that the absorption and emission of a

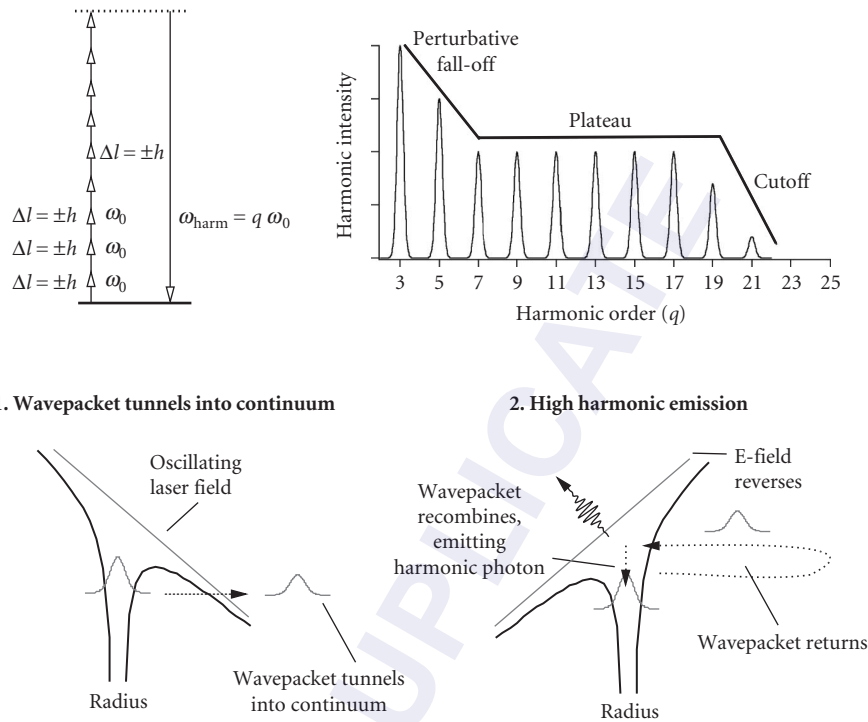


FIGURE 16 Simple view of the various aspects of high order harmonic generation physics. In the upper left, the multiphoton view of HHG is illustrated showing why only odd harmonics are possible given the need to conserve momentum by the interaction with atoms in a gas. The plot in the upper right is a generalized illustration of the character of a usual HHG spectrum at nonperturbative intensities ($>10^{14}$ W/cm 2) showing the presence of an initial fall-off of yield with harmonic order at low orders followed by a long plateau of harmonics at roughly constant yield terminated by a cutoff. The bottom pictures illustrate the quasi-classical picture of HHG in which the generation of harmonic photons can be thought of as a result of a tunnel ionized electron recombining with its parent ion upon return to the nucleus in the field.

photon is accompanied by the change in angular momentum of $\pm\hbar$, indicates that only odd harmonics of the light field can be emitted. (It also indicates that harmonics cannot be produced in a gas with circularly or elliptically polarized light.^{164,165}) Because all of the atoms in a gas are driven in this way harmonically, their emitted radiation will add coherently and the emitted harmonic radiation will propagate in the same direction as the drive laser and will retain many of the temporal and coherence properties¹⁶⁶ of the driving laser field. More accurately, HHG can be thought of as arising from the nonlinear polarization induced in a medium of density n_g such that $\mathbf{P}(t) = n_g \mathbf{d}(t)$, where $\mathbf{d}(t)$ is the induced atomic dipole. In the single active electron approximation, $\mathbf{d}(t) = \langle \psi(t) | \mathbf{er} | \psi(t) \rangle$, where $\psi(t)$ is the time-dependent wave function of the laser-driven atom.

This dipole, when driven very nonlinearly by a strong light field, will have Fourier components, $\mathbf{d}(q\omega_0)$, at frequencies $q\omega_0$ out to rather high harmonic orders (where q is an odd integer). In the weak field regime, where perturbation theory is appropriate, $\mathbf{d}(q\omega_0)$ varies as the q th power of the electric field, resulting in a harmonic yields which increase as the q th power of the incident intensity. Such behavior is indeed observed at intensity below about 10^{13} W/cm 2 in gases of noble gas atoms at moderate harmonic number (say $q = 3$ to 9). However, in a strong field, at intensity of 10^{14} to 10^{15} W/cm 2 the nonperturbative behavior of the laser driven atom leads to an atomic dipole of the q th harmonic

that will vary with a power law, I^p , that is usually well below the harmonic order ($p < q$).¹⁶⁷ In most strong field interactions in near-IR pulses at orders of $q = 11$ to 101, p will typically rest in the range of 5 to 8.⁷

Harmonic generation in this nonperturbative regime leads to a dramatic spectrum of harmonics that differs markedly from that expected from simple, multiphoton arguments. The emitted harmonic spectrum in this nonperturbative regime usually exhibits behavior illustrated in Fig. 16. The yield of the low order harmonics, out to perhaps the 5th or 7th order, falls exponentially, as expected under perturbation theory. However, the yield of harmonics at higher orders will then remain roughly constant, out to rather high order, a feature usually termed the HHG plateau. This plateau is then followed by an abrupt cutoff in harmonic production. Extremely high orders of harmonics from near-IR lasers have been demonstrated, with $q > 101$ achievable with sub-100 fs laser pulses. When extremely short (<20 fs) pulses are used, the gas atoms can survive to an even higher intensity resulting in extremely high order harmonic production. Orders well over $q \sim 201$ can be produced with such pulses in high ionization potential atoms such as helium,^{6,168} though the spectrum loses its distinct harmonic character at these very high orders. Obviously, conversion of a near-IR laser to such high-order results in light with wavelengths in the XUV and soft x-ray region (2 to 30 nm), and as such represents an attractive means for generation coherent pulses in this soft wavelength region.

Quasi-Classical Model of High Harmonic Generation In the strong field regime, it is possible to describe HHG with a quasi-classical model, much as strong field ionization and ATI can be described in this regime.^{10,11} Again, such a model is appropriate when the light field meets the condition that $\gamma_K < 1$. The mechanism for producing harmonic radiation by this model is illustrated in Fig. 16 at the bottom. The strong laser field can induce tunneling of the bound electron wavepacket; this freed electron wavepacket oscillates in the laser electric field. Electrons freed at certain phases of the laser oscillation can return to the vicinity of the nucleus and recombine back down to its initial ground state, emitting a photon of energy $I_p + \epsilon_{\text{osc}}$. The periodic return of many wavepackets leads to emission of radiation at well-defined harmonics of the light field.¹⁶⁹

The energy of the emitted photon depends on the electron's energy upon return to the nucleus. As Eq. (24) illustrates, if the electron is born by tunneling at a phase ϕ_0 , then the electron can return to the nucleus to produce a harmonic photon with energy

$$\hbar\omega_{\text{Harm}} = \frac{1}{2}m_e v_{\text{osc}}^2 (\sin\phi - \sin\phi_0)^2 + I_p \quad (31)$$

As mentioned in Sec. 21.5, analysis of Eqs. (24) and (31) indicates that the electron will return with maximum energy at $\phi_0 = 17^\circ$ and 197° and does so, in that case, with energy $3.17U_p$. This analysis leads to a well-known formula for the *maximum* photon energy of the high harmonic spectrum, known commonly as the cutoff harmonic

$$q_{\text{cutoff}} \hbar\omega_0 \approx I_p + 3.2U_p \quad (32)$$

In most experimental situations, U_p is usually evaluated at the intensity at which ionization starts to saturate. Above this intensity, there are no further atoms to participate in the harmonic generation process so further increase in intensity does not lead to a higher order cutoff harmonic. This cutoff formula was first discovered via numerical simulations¹⁷⁰ and has been well confirmed in many experiments.^{6,7,168}

This quasi-classical model for HHG predicts some important features of HHG. Because some time must elapse between the freeing of the electron by tunneling and its return to the nucleus to emit a photon, the emitted harmonic will pick up a phase shift with respect to the laser field.^{169,171} This model indicates that this phase shift is intensity dependent and, therefore, has a consequence on the macroscopic phase matching of the harmonics.¹⁷² Also note that Eqs. (24) and (31) indicate that most return energies can be generated by two different electron trajectories, resulting from liberation at two distinct phases in the laser field. The two trajectories, a short and a long trajectory, can

destructively interfere in the harmonic generation because of their differing phases resulting from their different times spent in the continuum before recombining with the parent nucleus.^{173,174}

Harmonic Yield and Phase-Matching Considerations Because HHG is a parametric process that results from the coherent addition of radiation produced from the nonlinear oscillations of electrons from many atoms in the gas, the harmonic intensity is strongly affected not only by the single atom physics of the atomic dipole, but also by phase-matching processes. As in standard nonlinear optics and usual harmonic generation, the harmonic field will continue to grow along with the propagating laser as long as the HHG field is in phase and newly generated harmonic light can add coherently.¹⁷⁵ If this condition is preserved, the harmonic yield will increase as the square of propagation distance. However, phase shifts associated with differences in refractive index between the fundamental and the harmonic as well as phases intrinsic to a focused laser beam will lead to mismatches that will clamp the harmonic generation.

While there are a number of ways to generate harmonics, one of the most common ways is to focus a Gaussian-shaped beam into a gas medium of length L with density n_0 . Calculating the harmonic yield is a very difficult problem, however, estimates can be simply made with a few assumptions. A simple model for the harmonic conversion yield can be delineated for a focused Gaussian beam with confocal parameter b (defined by usual Gaussian optics as $b = 2\pi w_0^2/\lambda$, where w_0 is the usual $1/e^2$ focal radius) and with the induced single atom dipole at the q th harmonic that varies as $|d(q\omega)| = \xi_p |E_p|$. The integrated energy yield of the q th harmonic generated by a square top pulse of duration t_p then is given by¹⁶⁷

$$\text{Energy}_q = \frac{\pi^2 q^2 \omega_0 n_0^2 b \tau_p |d(q\omega)|^2 L^2}{p} \left[\frac{\sin\{(\Delta k + 2q/b - 2q/pb)L/2\}}{\{(\Delta k + 2q/b - 2q/pb)L/2\}} \right]^2 \quad (33)$$

where Δk is a phase mismatch induced by the medium itself. The final sinc² factor is a strong field generalization of the phase-matching factor found in standard nonlinear optics and determines the coherence length

$$L_c = [\pi(\Delta k/2 + q/b - q/pb)]^{-1} \quad (34)$$

over which harmonic generation can build up. A medium length, longer than this coherence length will result in destructive interference and a clamp on the harmonic conversion.

The harmonic yield is determined by a complicated interplay of the laser propagation, the medium density, the atomic response, and loss of media through ionization. As a result, the harmonic conversion efficiency of high harmonics in the plateaus varies broadly and can range from 10^{-8} times the input laser energy for very high harmonics ($q \sim 31$ to 101) in tightly bound atoms (such as He or Ne)¹⁷⁶ up to as high as $\sim 10^{-5}$ for moderate order harmonics ($q \sim 23$ to 31).¹⁷⁷ The phase mismatch Δk usually arises at low intensity from the intrinsic dispersion of the gas itself between the fundamental and the harmonic. At higher intensity, as ionization becomes important, Δk will be determined by the dispersion of the plasma itself. In this regime when a plasma density of n_e is created by ionization

$$\Delta k_{\text{plasma}} \equiv \frac{2\pi n_e e^2}{qc\omega_0 m_e} (q^2 - 1) \quad (35)$$

At higher intensity in weakly focused beams (i.e., when $b \gg L$) this plasma-induced phase mismatch dominates the high harmonic production process. For example, at a plasma density of 10^{18} cm^{-3} , the coherence length of the 31st harmonic implied by Eq. (35) is only $\sim 10 \mu\text{m}$. Therefore, harmonics are produced only over this length, even if the medium is substantially longer (as is usually the case).

Attosecond Pulse Generation

The harmonic spectrum schematically illustrated in Fig. 16 in fact has a very broad bandwidth if the entire spectrum is considered as having a coherent phase relationship over the entire spectral window. This implies that such a broad spectrum, when Fourier transformed, results in a pulse, or a train of pulses, with duration well under 1 fs, that is, in the attosecond regime. This situation can indeed be achieved experimentally leading to the production of attosecond pulses with duration approaching 100 as.^{9,31,178–180} The physical origin of this can be easily seen within the context of the quasi-classical model. The return of an electron during its HHG generating recollision can be thought of as producing a short burst of bremsstrahlung radiation with duration comparable to the return encounter of the electron. If a laser pulse is short enough, such a bright burst can be made to occur for only one laser cycle and, therefore, produce a single isolated burst of attosecond radiation. This process is described at length in another chapter in this volume.

21.8 STRONG FIELD INTERACTIONS WITH CLUSTERS

When a strong laser field interacts with a cluster of atoms, collective effects not present in the interaction of strong field pulses with ions or small molecules come into play.^{19,39,181–183} Here clusters refer to assemblages of greater than ~ 100 atoms on one hand, but assemblages whose spatial dimension is still well below the laser wavelength, that is, particles with diameter < 100 nm or $< 10^6$ atoms. Such clusters are usually van der Waals bonded assemblies of atoms or small molecules. Large, easily polarizable atoms such as Xe or Kr form clusters most easily, though even small atoms or molecules, such as H_2 , can be made to cluster under the right circumstances.

Ionization Mechanisms in Clusters

Clusters will ionize in strong laser fields at intensities in which single atoms begin to ionize. In the cluster, however, there are two aspects of the ionization that must be considered, the inner ionization of the constituent atoms and ions and the outer ionization, resulting from the removal of free electrons within the cluster out, away from the cluster.³⁹ Ionization of the cluster is shaped by the fact that the laser field penetrates completely through the cluster, even when many free electrons are retained in the cluster and act as a plasma. This can be seen by noting that the plasma skin depth $\lambda_p = c/(\omega_p^2 - \omega_0^2)^{1/2}$ is almost always much smaller than the cluster diameter, even at solid electron density n_e (note that $\lambda_p \sim 20$ nm at solid density). Here $\omega_p = (4\pi e^2 n_e / m_e)^{1/2}$ is the plasma frequency (the resonant frequency at which electron waves in a plasma oscillate).¹⁸⁴

Inner Ionization Processes in the Cluster When a strong field light pulse begins to interact with the initially unionized cluster the constituent atoms in the cluster undergo ionization. This process is usually termed inner ionization because the electrons liberated from the atoms within the cluster do not necessarily exit the cluster as a whole. If these electrons linger in the cluster, they produce a nanoplasma which has unique collective properties.¹⁸¹ Inner ionization in the early stages of the laser interaction is usually dominated by quasi-classical tunnel ionization. The tunnel ionization of atoms within the cluster can be enhanced in a manner similar to the enhancement of ionization in diatomic molecules resulting from the suppression of the binding potential from neighboring ions in the molecule (discussed in Sec. 21.6).¹⁸⁵ The presence of nearby ions in the cluster can succeed in lowering the Coulomb binding potential (as illustrated in Fig. 17), and tunnel ionization can be greatly enhanced. This process has been termed “ionization ignition.”¹⁸⁶ A second process can also increase the rate of inner ionization in the cluster. This ionization occurs by collisional ionization by free electrons in the cluster, driven in an oscillatory motion by

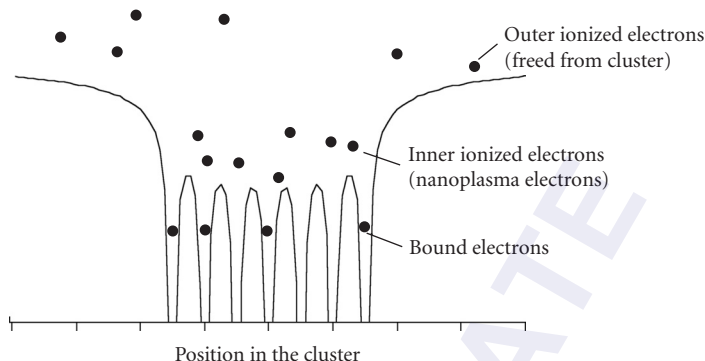


FIGURE 17 Plot of the Coulomb potential of ions arrayed in a cluster. The close proximity of many ions can suppress the barrier of the binding potentials, enhancing the tunnel ionization rate in the cluster. Some ionized electrons, however, will be confined to the cluster potential as a whole, resulting in the formation of a nanoplasma.

the strong laser field. The rate of this laser-driven collisional ionization can be calculated by using the well-known Lotz formula for ionization¹⁸⁷

$$W_{\text{las}} = n_e \frac{a_i Q_i}{\pi I_p} \sqrt{\frac{1}{m_e U_p}} \int_{I_p}^{2U_p} \frac{\ln(K_e/I_p)}{2K_e} \frac{1}{\sqrt{1-K_e/2U_p}} dK_e \quad (36)$$

Here $a_i \approx 1.1 \times 10^{-37} \text{ cm}^2 - \text{erg}^2$, and Q_i is the number of electrons in the outer shell of the ion. This laser driven collisional ionization rate tends to dominate the ionization in the cluster, usually being higher than the tunneling rate for most solid density cluster plasmas.¹⁸¹ This laser driven collisional ionization in the cluster leads to the production of charges states well above those which are observed in strong field ionization of single atoms at similar intensity.¹⁸⁸⁻¹⁹⁰ For example, Eq. (36) predicts that the ionization rate for $\text{Ar}^{+15} \rightarrow \text{Ar}^{+16}$ is about 0.001 fs^{-1} at an intensity of $5 \times 10^{17} \text{ W/cm}^2$ in a 1- μm laser pulse. This would yield approximately 10 percent ionization to He-like Ar during a 100-fs pulse. On the other hand, the BSI theory indicates that an intensity of $2 \times 10^{21} \text{ W/cm}^2$ would be needed to achieve significant ionization to Ar^{+16} by direct field ionization. This laser driven collisional ionization can lead to the production of very high charge states in clusters, even at modest laser intensity.^{189,190}

Outer Ionization of the Cluster If space charge forces retain the inner ionized electrons in the body of the cluster, the removal of electrons from the cluster sphere, outer ionization, may take much longer than the initial inner ionization (see Fig. 17). There are two regimes for outer ionization in the strong field irradiation of the cluster. First, the laser field itself can extract the electrons from the cluster. In a very strong field (i.e., $E_0 \gg Q_{\text{clust}} e/R_0^2$, where Q_{clust} is the number of electrons removed by outer ionization from the cluster, and R_0 is the radius of the cluster), outer ionization occurs almost instantaneously and the cluster enters the Coulomb explosion regime. For laser fields comparable to or smaller than the binding field of the laser, it can be shown that the number of electrons outer ionized by the laser field is¹⁹¹

$$Q_{\text{clust}} \approx 12\pi \frac{n_e R_0^2 e E_0}{m_e \omega_p^2} \quad (37)$$

which is proportional to the square root of intensity. Alternatively, outer ionization can occur by electrons that have been heated sufficiently to escape the binding potential of the cluster. In this case, the rate of outer ionization by “free streaming” can be estimated assuming a Maxwellian electron energy distribution with temperature, T_e ¹⁸¹

$$W_{\text{fs}} = n_e \frac{2\sqrt{2\pi}}{m_e^{1/2}(kT_e)^{1/2}} (K_{\text{esc}} + kT_e) \exp\left[-\frac{K_{\text{esc}}}{kT_e}\right] \times \begin{cases} \frac{\lambda_e}{4r} (12r^2 - \lambda_e^2) & \lambda_e < 2r \\ 4r^2 & \lambda_e > 2r \end{cases} \quad (38)$$

where $K_{\text{esc}} = Q_{\text{clust}} e^2 / R_0$ is the energy needed to escape, $\lambda_e = (k_B T_e)^2 / 4\pi n_e (Z+1) e^4 \ln \Lambda$ is the electron mean free path in the cluster and $\ln \Lambda$ is the well-known plasma Coulomb logarithm.

Coulomb Explosion of Small Clusters

If a cluster is stripped of most of its electrons via the outer ionization process very quickly, much faster than the cluster can expand, the cluster will disassemble by a Coulomb explosion. For the cluster to evolve in this limit, two conditions must be met (1) the intensity of the light field must be high enough and (2) the laser pulse must ramp up to the intensity needed for complete outer ionization must faster than the cluster expands.

The first condition is difficult to determine as the intensity needed for complete outer ionization is not easily calculated analytically; it depends on the dynamics of the driven electron cloud in the cluster. This can be estimated by requiring that the laser ponderomotive energy dominates the electron dynamics over the confining potential energy of the cluster,

$$U_p > \frac{Q_{\text{clust}} e^2}{R_0} \quad (39)$$

Here Q_{clust} is the total charge on the cluster sphere from electrons that have previously exited the cluster. The second condition mandates that the rise time of the laser pulse to the stripping ponderomotive potential be faster than the explosion time of the cluster. This characteristic explosion time can be estimated by calculating the time required for a charged cluster to expand from its initial radius, a , to twice its initial radius. Integration of the motion of a charged deuterium cluster yields for this characteristic explosion time:

$$\tau_{\text{Coul}} \approx \sqrt{\frac{m_i}{n_i Z^2 e^2}} \quad (40)$$

where m_i is the mass of the ions and n_i is the density of atoms in the cluster. Note that τ_{Coul} is independent of cluster radius and equals about 15 fs for fully stripped hydrogen clusters. Equation (40) indicates that the Coulomb explosion limit can usually be accessed in clusters only with intense, sub-100-fs laser pulses.¹⁹²

If a Coulomb explosion is driven and it can be assumed that all electrons are removed prior to any ion movement, the ion energy spectrum from a single exploding cluster, denoted $f_{\text{sc}}(\mathcal{E})$, can be stated as¹⁹³

$$\begin{aligned} f_{\text{sc}}(\mathcal{E}) &= \frac{3}{2} \mathcal{E}_{\text{max}}^{-3/2} \mathcal{E}^{1/2} & \mathcal{E} \leq \mathcal{E}_{\text{max}} \\ &= 0 & \mathcal{E} > \mathcal{E}_{\text{max}} \end{aligned} \quad (41)$$

where \mathcal{E}_{max} is the maximum energy of ions ejected and corresponds to those ions at the surface of the cluster, $\mathcal{E}_{\text{max}} = 4\pi e^2 n_i R_0^2 / 3$. In hydrogen clusters, for example, \mathcal{E}_{max} is about 2.5 keV for 5 nm

clusters. This ion spectrum is peaked near ϵ_{\max} . In most experiments, however, the clusters irradiated are composed of a broad size distribution. This tends to broaden the ion energy distribution observed.

Nanoplasma Description of Large Clusters

In the limit of a Coulomb explosion, the principal absorption mechanism for laser light is in the deposition of the energy needed to expel the electrons from the charged cluster. The other limit of strong field laser cluster interactions occurs when there is little or no outer ionization subsequent to significant inner ionization in the cluster. This tends to occur in larger clusters and clusters composed of higher Z atoms which can become more highly charged. When outer ionization lags inner ionization, a nanoplasma is formed in the cluster.¹⁸¹ In this case, the laser interactions with the cluster can be dominated by collective oscillations of the electron cloud.

Cluster Electron Heating Because the electrons are confined to the cluster by space charge forces, they can acquire energy from the intense laser field. Because of the collective oscillation of the electron cloud, the electric field inside the cluster will be⁷²

$$E = \frac{3}{|\epsilon_D + 2|} E_0 \quad (42)$$

where the dielectric constant can usually be taken to be that from a simple Drude model for a plasma $\epsilon_D = 1 - \omega_p^2 / \omega_0(\omega_0 + i\nu)$. ν is the electron-ion collision frequency (discussed below). Equation (42) indicates that the field is enhanced in the cluster when the laser frequency is $3^{1/2}$ times the plasma frequency, that is, when $n_e/n_{\text{crit}} = 3$. ($n_{\text{crit}} = m_e \omega_0^2 / 4\pi e^2$ is the plasma critical density, at which the laser oscillates at the plasma resonance frequency.) This resonance condition corresponds to a state in which the laser frequency matches the natural collective oscillation frequency of the electron cloud in the spherical cluster. This resonance condition is accompanied not only by an increase of the field in the cluster but also an increase in the electron heating rate and cluster absorption of energy from the laser. The heating rate of electrons in the cluster is then¹⁹⁴

$$\frac{\partial U}{\partial t} = \frac{9\omega_0^2 \omega_p^2 \nu}{8\pi} \frac{1}{9\omega_0^2(\omega_0^2 + \nu^2) + \omega_p^2(\omega_p^2 - 6\omega_0^2)} |E_0|^2 \quad (43)$$

The choice of collision frequency in this equation is complicated by various aspects of strongly coupled plasma physics, but for most interactions in which the electrons in the cluster are dominated by the driven motion of the laser electric field, we can say that the collision frequency most relevant for strong field cluster interactions is¹⁹⁵

$$\nu_E = \frac{16Z^2 e n_i m_e \omega_0^3}{E_0^3} \left(\ln \left[\frac{eE_0}{2m_e \omega_0 \nu_e} \right] + 1 \right) \ln \Lambda \quad (44)$$

where Z is the average charge state of ions in the cluster and $\nu_e = (k_B T_e / m_e)^{1/2}$. The principal consequence of this resonance occurs after the cluster has been initially inner ionized to an electron density near that of a solid ($\sim 10^{23} \text{ cm}^{-3}$). As the cluster expands, its resonance frequency falls and, if the laser pulse is long enough, will come into resonance with the laser. This is accompanied by a violent driving of the cluster nanoplasma cloud with rapid energy deposition in the cluster. Electron temperatures of many tens of keV are possible, even with modest ($< 10^{16} \text{ W/cm}^2$) drive intensities. When such cluster nanoplasmas are irradiated in a dense gas jet, very high (near 100 percent) absorption of the laser occurs¹⁹⁶ and bright x-ray emission, high ion temperatures, and even nuclear fusion can be observed.^{182,197–200} Because of this dynamic plasma resonance, there is often an optimum pulse

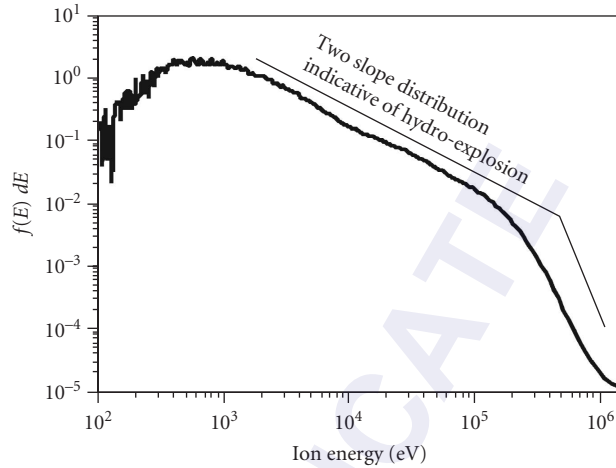


FIGURE 18 Example of the ion spectrum resulting from the irradiation of Xe clusters in the nanoplasma regime. These ~ 2500 atom Xe clusters were irradiated by a 100-fs pulse at an intensity of $\sim 2 \times 10^{16}$ W/cm 2 . This illustrates how ions with energies (~ 1 MeV in these data) much higher than U_p (~ 1.2 keV here) can be generated during the interaction with clusters. (Figure adapted from Ref. 189.)

duration for heating clusters in this manner.²⁰¹ For example, small Ar clusters can expand into resonance on a 100-fs timescale, while large (>10 nm) Xe clusters can take 1 to 5 ps to reach resonance with a near-IR laser.²⁰²

Cluster Expansion Mechanisms In the nanoplasma regime, the cluster will expand, though not by Coulomb repulsion forces between ions, as it does in the Coulomb explosion regime. Instead, it will expand by the ambipolar potential created by the thermal pressure of the hot electrons in the nanoplasma, $p_e = n_e k T_e$. This pressure will lead the cluster to explode on a time scale of $t_{\text{expl}} \approx R_0 (m_e / Z k T_e)^{1/2}$. The resulting ion spectrum will be characteristic of that from a hydrodynamically exploding plasma, as illustrated in Fig. 18. The hot tail that results from such an expansion can lead to the production of ions with hundreds of keV to MeV of energy, even in laser pulses with ponderomotive energies of only ~ 1 keV (an intensity of about 10^{16} W/cm 2 in a near-IR field).

It should be noted that the simple model used here to describe the clusters in the nanoplasma regime rely on the assumption that the cluster remains more or less uniform density. During the expansion of the cluster, it will devolve from a uniform density to a plasma with a radial plasma gradient. In this case, the absorption is likely dominated by resonance absorption (described in Sec. 21.10) at the critical density surface around the cluster circumference.²⁰³ In addition to the absorption processes described, various plasma processes can occur, such as electron ion equilibration and electron recombination, affecting the cluster dynamics.

Intense Laser Pulse Interactions with Clusters in the Nonneutral Regime

The description of the dynamics of a cluster in a strong laser field in a regime intermediate to the Coulomb explosion and the nanoplasma regimes, when outer ionization is only partial, is complicated. Simple electrostatic theory tells us that if some of the electrons have been outer ionized, then

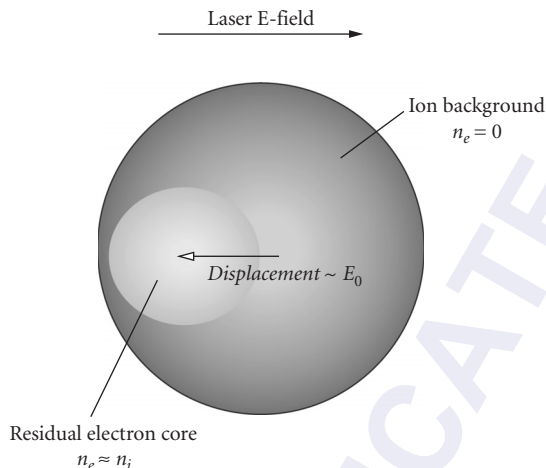


FIGURE 19 Schematic showing the geometry of the electron cloud within the ion sphere of a partially charged cluster. The electron cloud contracts to create a quasi-neutral core which can be driven from side to side by the strong laser field. When the oscillations are large enough that some of the electron cloud is pulled out away from the ion sphere, outer ionization of these electrons occurs.

the cluster, no longer quasi-neutral will evolve so that the remaining electrons will collapse into an inner, neutral cloud within the cluster, a situation illustrated in Fig. 19.¹⁹¹ The laser field will then pull the electron cloud over by a distance $d_0 = 3eE_0/m_e\omega_p^2$. This extracts electrons through laser field acceleration.

Electron Stochastic Heating In the intermediate regime of cluster ionization, there will be a population of electrons that will be driven by the laser in the vacuum surrounding the cluster sphere. These electrons can pass in and out of the cluster a number of times, picking up energy from the laser at each pass, in a manner similar to the vacuum heating described in Sec. 21.10. This heating is called stochastic heating and leads to the generation of a population of very high energy electrons. It can be shown that the maximum energy that can be reached by electrons in this manner is $\epsilon_{\max} \sim m_e R_0^2 \omega^2$, an energy which can be well in excess of the ponderomotive energy if the cluster is much larger than a quiver amplitude.¹⁹¹

21.9 STRONG FIELD PHYSICS IN UNDERDENSE PLASMAS

In this section, we will discuss the physics involved in the interaction of strong field electromagnetic pulses with underdense plasma, namely, plasma in which the plasma frequency $\omega_p = \sqrt{4\pi e^2 n_e / m_e}$ is smaller than the laser frequency, ω_0 . In this situation, the refractive index of the plasma $n = \sqrt{1 - \omega_p^2 / \omega_0^2}$, is real and the laser field can propagate through the plasma. At high intensity, various field driven coupling mechanisms occur which deposit energy into the plasma fluid through interactions with charged particles individually or through interactions with plasma waves. These interactions tend to heat the plasma electrons as a whole or they directly couple laser energy into a small population of fast electrons.

Strong Field Inverse Bremsstrahlung Heating

Perhaps the most basic mechanism for an intense light pulse to heat underdense plasma is through inverse bremsstrahlung: collisional heating. In the strong field regime with IR or optical wavelength pulses, the process can be accurately described classically. The nature of the heating is illustrated in Fig. 20. When an oscillating electron in the laser field collides with an ion and scatters from the Coulomb field of the ion, its adiabaticity is broken, acquiring some random, thermal energy from the laser field. The heating rate of electrons is then some appropriate electron-ion collision frequency times the amount of energy gained per collision, which can usually be taken to be the ponderomotive energy. So the heating rate per electron is $dU/dt|_{IB} \approx v_{ei} U_p$.¹⁸⁴

At low intensities, the electron-ion collision frequency can be taken to be the usual temperature-dependent plasma electron-ion collision rate.²⁰⁴ In the strong field regime the picture is a little different. When $U_p > k_B T_e$ the electron motion is dominated not by thermal motion but instead by the ponderomotive motion of the laser, which means that the standard collision frequencies utilized in normal plasma physics cannot be used. There have been various models published to describe the electron heating and dynamics in this strong field limit.^{195,205} A good model for heating of electrons in a strong field leads to a heating rate of²⁰⁶

$$\frac{\partial U}{\partial t}\bigg|_{\text{las}} = \frac{16Z^2 n_i e^3 m_e \omega_0}{3E_0} \ln \left[\frac{eE_0}{m_e^{1/2} \omega_0 (kT_e)^{1/2}} \right] \ln \Lambda \quad (45)$$

Again, $\ln \Lambda$ is the usual plasma Coulomb logarithm, which can be taken for the underdense plasmas treated here to be $\Lambda = (kT_e)^{3/2} / 4\pi^{1/2} Z e^3 n_i^{1/2}$.

Notice that in this high-intensity regime, the heating rate actually decreases with increasing intensity (as $I^{-1/2}$), which results from the strong decrease of the Coulomb scattering cross section with increasing electron velocity. (The heating rate also decreases with increasing wavelength for the same reason.) It is interesting to note that, when the full kinetic evolution of the electron energy distribution function is solved in the high-field limit, the electron distribution approaches that naturally of a Maxwellian, independent of any electron equilibration.²⁰⁶

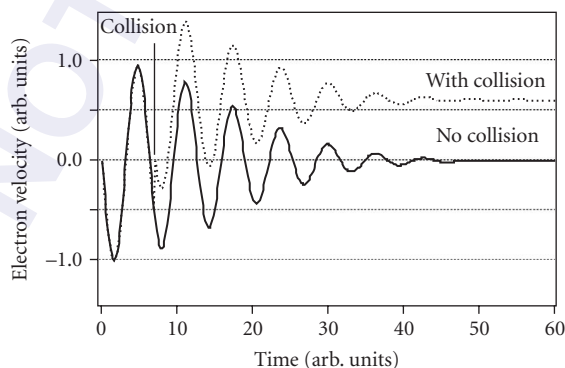


FIGURE 20 Calculation of the classical trajectory of an electron in a laser pulse which decays in time adiabatically to zero. The electron velocity falls to zero with the laser field and acquires no net energy when no collision takes place. If, however, an instantaneous 90° collision occurs, this breaks the adiabaticity of the electron's oscillation leaving it with some residual velocity after the pulse field has fallen to zero. This is the origin of inverse bremsstrahlung heating.

Plasma Instabilities Driven by Intense Laser Pulses

A salient feature of high-intensity laser interactions with plasmas is that plasmas can support collective motion, including coherent waves, which couple to the electromagnetic field of the laser. These plasma waves manifest themselves in various ways, such as in ion acoustic waves, (which are essentially sound waves in the plasma gas¹⁸⁴). Because a plasma is composed of a positively charged ion fluid and a negatively charged, light electron fluid, it can also support electrostatic and propagating electromagnetic waves, often called Langmuir waves, composed of an oscillating electron density fluctuation.

Because the laser field is an electromagnetic wave, it can couple energy to these plasma waves as the pulse propagates through the plasma. For example, if the pulse drives ion acoustic waves, the process is termed Brillouin scattering. This process is usually not significant in ultrashort pulse, high-intensity interactions because the mass of the ions are so large that the growth rate of such instabilities is too slow for the laser pulse duration. Of greater importance to intense ultrashort laser pulses (with $I > \sim 10^{16}$ W/cm²) is the coupling of laser energy to electron plasma waves. This coupling can occur through a panoply of mechanisms, such as the so-called $2\omega_p$ mechanism,²⁰⁷ but the most important process in the strong field regime is electron Raman scattering.²⁰⁸ The energetics of this process are described in Fig. 21. In Raman scattering, one laser photon is coupled to an electron plasma wave (which naturally oscillates at the plasma frequency) and results in the production (or destruction) of one quanta of the plasma wave with the simultaneous production of a photon at a shifted wavelength. The scattered light has frequency downshifted if energy added to the plasma wave (the Stokes process) or upshifted if energy is absorbed from the plasma (the anti-Stokes process).

In an intense laser pulse this process can undergo positive feedback, resulting in an exponentially growing instability, coupling a significant amount of light energy into the electron plasma wave (which is usually a longitudinal electron density fluctuation). How this happens can be seen if one considers what happens if a Stokes photon is produced in a Raman scattering event at a frequency of $\omega_s = \omega_0 - \omega_p$. This new field can add to the fundamental laser oscillation creating a beat frequency of $\omega_b = \omega_0 - \omega_s = \omega_p$, which can then, in turn, drive the plasma wave resonantly, increasing its amplitude and further driving the production of more Stokes photons. This feedback process leads to exponential growth of the Stokes field and the plasma wave, and is usually referred to as stimulated Raman scattering (or SRS).

The growth rate, defined by noting that the scattered wave amplitude and scattered light field grow as $\sim e^{\Gamma t}$, can be found from the coupled nonlinear equations describing the Stokes field (which is the one which will be resonant, growing the fastest) and the plasma wave. In the nonrelativistic limit, this growth rate is¹⁸⁴

$$\Gamma_{\text{SRS}} = \frac{k_e}{2} \sqrt{\frac{U_p}{m_e} \frac{\omega_p^2}{\omega_{\text{BG}}(\omega_0 - \omega_{\text{BG}})}} \quad (46)$$

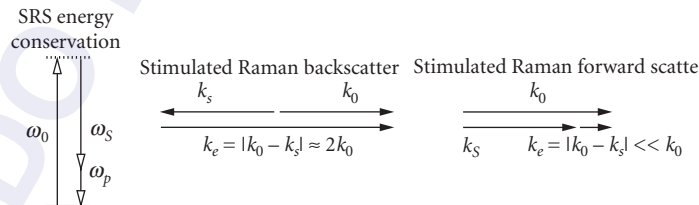


FIGURE 21 On the left, the energetics of Raman scattering in a plasma are demonstrated for the Stokes process in which a laser photon is scattered by a plasma wave creating a lower energy Stokes photon and depositing $\hbar\omega_p$ of energy into the plasma wave. The right-hand drawing illustrates the conservation of momentum for SRS backscatter and forward scatter, showing that the electron wave k -vector, and hence the growth rate, are maximum for the backscatter case.

here ω_{BG} is the Bohm-Gross plasma frequency $\omega_{\text{BG}} = \sqrt{\omega_p^2 + 3k_e^2(k_B T_e)^2/m_e}$, which is simply equal to the plasma frequency for a cold ($T_e = 0$) plasma, and $k_e = |\mathbf{k}_e|$ is the magnitude of the wave-number of the plasma wave, defined by conservation of momentum by $\mathbf{k}_e = \mathbf{k}_0 + \mathbf{k}_s$. Equation (46) illustrates that the growth rate of this instability increases as the square root of the light intensity. In underdense plasma such that $\omega_0 \gg \omega_p$, ω_{BG} the growth rate also scales with the square root of laser drive wavelength. Consequently, high intensity, long wavelength light pulses are much more susceptible to the SRS instability. In practice, an intense light pulse will cause the instability to grow from density fluctuations arising from thermal noise present in any plasma. When the laser pulse has passed, the plasma waves will tend to damp either by collisional processes in the plasma or by noncollisional processes such as Landau damping.¹⁸⁴

The growth of SRS has two significant consequences: (1) it acts to absorb energy from the intense laser pulse coupling its energy to the plasma electrons and (2) it can drive plasma waves to large amplitude, which, in turn generate high energy, nonthermal electrons. Often the hot electrons generated can have energy many times that of the ponderomotive energy, which means electrons of tens to hundreds of keV even at modest intensity (say $\sim 10^{16}$ to 10^{17} W/cm²). There are two important regimes of SRS.

Stimulated Raman Backscatter In a low density plasma ($\omega_0 \gg \omega_p$) $|\mathbf{k}_0| \approx |\mathbf{k}_s|$. From Eq. (46) it can be seen that for a given plasma density and laser intensity the growth rate is maximum when k_e is maximum. As Fig. 21 illustrates this occurs when the generated Stokes light is directly backscattered. In this case, when the plasma is cold, the growth rate in the nonrelativistic limit is $\Gamma_{\text{SRS-bs}} \approx \sqrt{\omega_p \omega_0 U_p / m_e c^2}$. The growth time in this case is just the temporal pulse duration of the main pulse, since the Stokes field propagates backward through the main pulse. If the laser pulse spatial length is comparable to the length of the medium of interaction, a condition satisfied, for example, when an intense picosecond pulse propagates through a gas jet of a few mm in length, SRS backscatter will be the dominant plasma instability.²⁰⁹ Practically speaking, SRS backscatter will be important—resulting in a significant amount of laser energy backscattered (>1 percent) and significant plasma heating—when $\Gamma_{\text{SRS}} \tau_{\text{pulse}} \sim 6$ to 7. Equation (46) indicates that this will occur in a 100-fs pulse with wavelength near 1 μm propagating through a plasma of $n_e \sim 10^{18}$ W/cm² when the pulse intensity exceeds $\sim 2 \times 10^{17}$ W/cm².

Stimulated Raman Forward Scatter When the laser pulse is very intense, and the duration is short, so that SRS backscatter has little time to grow, the situation is different. In this case, SRS forward scattering is possible, with geometry illustrated in Fig. 21.²¹⁰ Now, though k_e is small the Stokes or anti-Stokes radiation copropagate with the intense drive laser, (and, in fact, both Stokes and anti-Stokes can be nearly resonant and grow at about the same rate²¹¹). The scattered radiation growth factor is now determined by the length of the pulse propagation in the plasma. This effect tends to be particularly important for pulses at relativistic intensity (i.e., $>10^{18}$ W/cm²) so relativistic effects must be included in deriving a growth rate. In this case, the growth rate, in terms of the normalized vector potential a_0 , is⁶⁶

$$\Gamma_{\text{SRS-fs}} = \frac{\omega_p^2 a_0}{\sqrt{8\omega_0(1+a_0^2/2)}} \quad (47)$$

SRS forward scattering becomes important when $\Gamma_{\text{SRS-fs}} L/c \sim 6$, where L is the propagation length through the plasma. Equation (47) yields the surprising result that at strongly relativistic intensity $\Gamma \sim 1/a_0 \sim 1/I^{1/2}$ so the growth rate actually decreases with increasing intensity. This is a consequence of the increase of the effective mass of the electrons when they are driven at relativistic oscillation velocities. Equation (47) indicates that a 1- μm wavelength pulse traversing 5 mm of plasma at density 10^{18} W/cm² will exhibit significant forward scattering when $a_0 \approx 0.8$, an intensity of $\sim 1 \times 10^{19}$ W/cm².

Wakefield Generation and Electron Acceleration

The plasma instabilities discussed in the section “Plasma Instabilities Driven by Intense Laser Pulses” essentially arise from the fact that electron plasma oscillations can be driven to large amplitude by

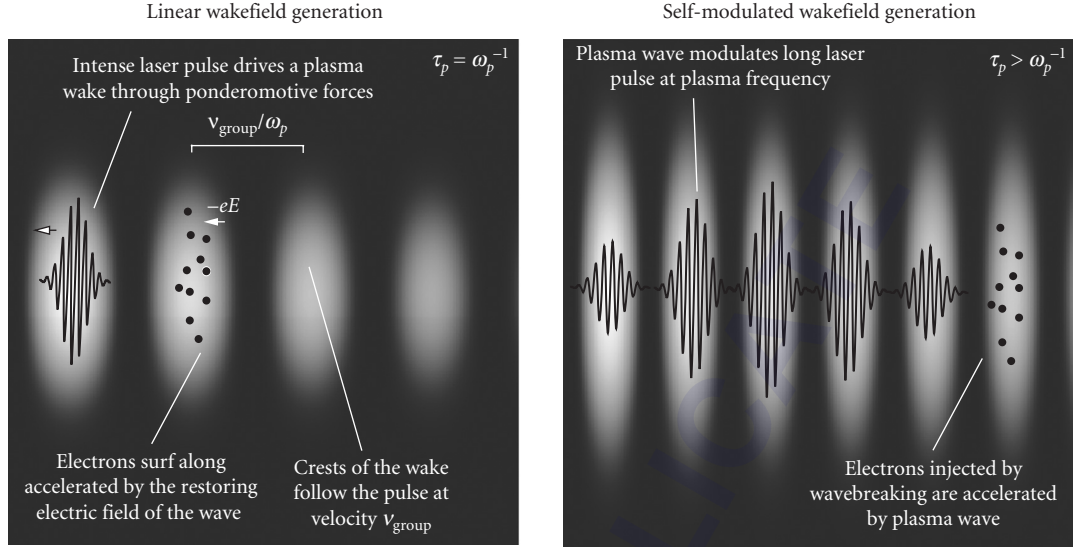


FIGURE 22 Schematic illustration of how wakefield acceleration works in two situations. On the left, the plasma wave is excited by the ponderomotive force of a laser pulse with duration comparable to the plasma oscillation time. On the right is illustrated the case in which the driving laser is longer than a plasma period, but through Raman forward scattering has become modulated at the plasma frequency.

the oscillating ponderomotive force of the strong light field. These plasma oscillations are longitudinal waves, with longitudinal electric fields acting as restoring forces on the oscillating electron density fluctuations. This situation is shown schematically in Fig. 22. The presence of these strong, oscillating electric fields has been investigated for a number of years as a means to accelerate electrons.^{14,15,32,212–217} Electrons to be accelerated are either injected externally or get accelerated from the free electrons in the plasma itself, often through a process known as wave breaking (discussed below). Since the acceleration is accomplished through the creation of a plasma wave left in the wake of the ponderomotive force of the laser, this idea is termed plasma wakefield acceleration.

In a general sense, the idea is to expel electrons from a region in the essentially immobile background ions through the ponderomotive force of the laser. The phase velocity of this plasma wave, then, is just that of the group velocity of the laser pulse traveling through the plasma, $v_g = c\sqrt{1-\omega_p^2/\omega_0^2}$. Electrons injected at the right phase of the plasma wave, if traveling at a velocity near this group velocity (i.e. $\sim c$) can surf along with the wave, acquiring energy from it. The laser intensity required is set by the scheme in which the plasma wave is produced (detailed below) but in general the ponderomotive force, $-\nabla U_p$ at the peak of the pulse should be sufficient to expel a sizable number of electrons from their background neutralizing ions. Practically this means that the required intensity is 10^{16} to 10^{18} W/cm² for ~ 1 - μm wavelength pulses in moderate density ($\sim 10^{18}$ cm⁻³) plasmas. The longitudinal electric field that is produced in such an oscillating plasma wake, E_{WF} can be easily estimated using Poisson's equation and the fact that the wave moves at nearly c with an oscillation frequency of ω_p ³²

$$E_{\text{WF}} = \eta \sqrt{4\pi m_e c^2 n_e} \quad (48)$$

where η is the fractional plasma wave amplitude of the wakefield ($\Delta n/n$). Equation (48) indicates that very large acceleration gradients are possible with plasmas. In a plasma of density 10^{18} cm⁻³ nearly 1 GeV/cm acceleration can be achieved if the laser pulse is intense enough to approach $\eta \sim 1$.

The accelerating field is limited by the fact that, as one pushes to higher fields, eventually the plasma wave becomes nonlinear. This means that the electron fluid velocity begins to exceed the wave phase velocity and the wave loses its coherence, a situation called wave breaking.¹⁸⁴ This nonlinear behavior can serve to inject electrons out of the coherent wave oscillation into the accelerating gradient of the wakefield.²¹³ It does, however, set a limit on the maximum electric field attainable. In general this maximum field is very difficult to predict, but a rough estimate can be made using the “cold wave-breaking” limit, which says that the maximum wakefield strength occurs when $\eta \sim 1$, that is when $E_{\max} \approx m_e c \omega_p / e$.⁶⁶

The limit of acceleration with such laser driven plasma waves, at least in a single plasma stage, arises from the fact that there is a slight mismatch between the electron velocity and phase velocity of the plasma wave, so that electrons eventually move out of the accelerating phase of the wave and catch up to the electric field in the other direction, decelerating the electrons. This dephasing length can be found by assuming that the electrons are relativistic with velocity very near c , but that the plasma wave has velocity given by that of the laser pulse group velocity. Then the maximum length over which an electron can be accelerated before dephasing is $L_{\text{dephase}} \approx 2\pi\omega_0^2 c / \omega_p^3$. For the example just cited, this implies a maximum per-stage acceleration length of 3 cm. Because the acceleration field scales as $n_e^{1/2}$, while the dephasing length scales as $n_e^{-3/2}$ there tends to be an advantage in using lower plasma density (subject to the constraint that it becomes more difficult to propagate an intense focused pulse over longer distances). In general, there are three methods of optical wakefield generation.

1. *Linear wakefield generation.* This method relies on nearly resonantly driving a wakefield in the plasma by matching the laser pulse duration to the plasma oscillation period, such that $\tau_{\text{pulse}} \approx 2\pi/\omega_p$.^{14,212} This situation is schematically illustrated in Fig. 22 on the left. For the conditions cited in the above examples, the optimum plasma wave amplitude is achieved with a pulse of ~ 100 -fs duration.
2. *Self-modulated wakefield generation.* When the laser pulse is much longer than the plasma oscillation period, it cannot resonantly drive the plasma wave. However, as described in the section “Plasma Instabilities Driven by Intense Laser Pulses,” plasma instabilities can lead to the creation of radiation with frequency shifted by the plasma frequency.²¹⁰ This copropagating light, stimulated Raman forward scattering, modulates the laser pulse at the plasma period (illustrated in Fig. 22, right). Consequently, the self-modulated wakefield requires some period of propagation for instability growth and, therefore, typically requires higher intensities than the linear wakefield. In both linear and self-modulated wakefield experiments, the fast electrons that emerge from the plasma typically have a very broad energy spectrum, essentially arising from the fact that if electrons are not externally injected, wavebreaking will tend to inject electrons at ALL phases of the wakefield, leading to a range of accelerated energies.^{213,214,216}
3. *Plasma beat wave acceleration.* The third approach to wakefield generation involves copropagating two pulses at slightly different frequencies such that their difference “beat” frequency is resonant with the plasma frequency. This approach usually requires difficult laser technology and is not commonly pursued in the strong field laser regime.²¹⁵

Bubble Acceleration A somewhat different regime exists in the plasma when the driving laser pulse is very intense ($>10^{19}$ W/cm²), the pulse duration is shorter than the characteristic plasma oscillation time (about half an oscillation time has been found optimal in simulations²⁸), and the pulse is focused to a spot comparable to a plasma oscillation wavelength. In this regime the plasma electron density is driven strongly nonlinearly to an amplitude in which wave breaking occurs over a fraction of a plasma cycle.^{218–221} The plasma structure changes and cannot be described as a simple harmonic plasma oscillation like that pictured in Fig. 22. Instead, the plasma electrons are completely expelled by the three-dimensional ponderomotive force of the pulse, leading to a cavity depleted of electrons in the region immediately behind the laser pulse. This “bubble” of electrons filled with the ion background, depicted in Fig. 23, can then accelerate electrons which get trapped in a region just ahead of the closing rear wall of the bubble. The interesting aspect of such a structure is not that the accelerating

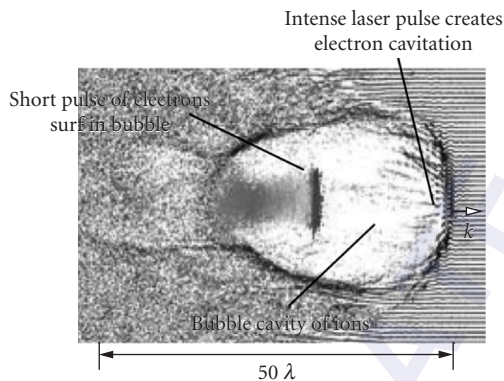


FIGURE 23 Illustration of how a plasma responds and electrons are accelerated when irradiated in the bubble regime. (This figure was adapted from Ref. 28.)

gradient is substantially different than the strongly driven linear wakefield case, but that self-injected electrons tend to get accelerated at one spot in the bubble, leading to electrons accelerated with gradients >1 GeV/cm with a quasi-monoenergetic spectrum, often with energy spread of only a few percent.²¹⁹

Ponderomotive Channel Formation

The strong ponderomotive forces which lead to the bubble formation described above can also lead to a Coulomb explosion of the ions due to their space charge repulsion while the electron expulsion exists.²²² If the laser pulse is substantially longer than a plasma period (a situation different than the bubble regime just discussed) the laser pulse will ponderomotively expel the plasma electrons in a more or less adiabatic manner from the focal region. This electron cavitation is nearly 100 percent if the following cavitation condition is met: $a_0 > w\omega_p/c$,⁶⁶ where w is the spot size radius of the focused laser beam. This ponderomotive cavitation, held during the duration of the laser pulse can then result in an ion radial explosion which will persist even after the pulse has departed because of the inertia imparted to the ions. This ponderomotive channel formation has been observed in experiments with picosecond duration laser pulses (which are long enough to initiate the ion expansion), leading to a channel formed on the laser axis.²²³ The radial Coulomb explosion of ions in the channel can lead to the production of radially directed ions with many MeV of kinetic energy.

Direct Laser Acceleration and Betatron Resonance

The expulsion of electrons by the strong ponderomotive force of a relativistic pulse propagating in an underdense plasma can have other consequences. On a timescale prior to the Coulomb explosion of the heavy ions (often taking many picosecond) the ions do not move significantly, producing a transient radial electric field that will confine a small number of electrons to the pondermotively cleared core, acting as potential well in which electrons can oscillate. As illustrated in Sec. 21.4, electrons in a relativistic intensity laser pulse are actually driven in a forward motion by the combination of the electric and magnetic fields of the laser; this laser driven current in turn creates a toroidal magnetic field. Trapped electrons in the ponderomotive channel will oscillate in this potential well. A relativistic electron trapped in the channel will oscillate at the betatron frequency, which is roughly $\omega_B = \omega_p/2\gamma^{1/2}$. If the electron is propagating at relativistic speed along with the laser, the betatron frequency can be in resonance with the oscillations of the laser field. This leads to significant

acceleration of the electron while it remains in resonance with the laser field, a situation termed “direct laser acceleration” (DLA).²²⁴

This situation is, in a sense, the inverse to the free electron laser situation. In an FEL energy is coupled from relativistic electrons to an electromagnetic wave in an undulator; in DLA the reverse occurs with the betatron oscillations of the ion channel playing the role of the undulator. The oscillating relativistic electrons will also produce soft x-ray radiation through their interactions with the quasi-static toroidal magnetic field,²⁵ a situation similar to electrons in a synchrotron.

Ionization-Induced Defocusing

As an intense focused laser pulse propagates into a gas it will ionize the atoms by tunnel and multi-photon ionization, initiating the creation of a plasma. As this ionization occurs, the electron density of the plasma increases, affecting the propagation of the laser in a process known as plasma-induced defocusing.^{225,226} Because the spatial profile of any focused laser pulse is inevitably peaked with higher intensity near the center of the propagation axis, ionization tends to occur in the center of the beam earlier in the pulse, resulting in a higher density plasma on axis than at the edges of the beam. The situation that arises is illustrated in Fig. 24. A plasma has a refractive index less than 1, meaning that the denser plasma yields a *faster* phase velocity. Because of this, the pulse’s phase fronts will advance in the center of the beam, where the ionized plasma has higher electron density, causing a net *defocusing* of the beam. This has the practical experimental consequence of clamping the maximum intensity that can be achieved by a focused laser in a gas target.²²⁶

The magnitude and specifics of the defocusing effects are difficult to quantify, as they depend on the kind of gas irradiated, the intensity of the laser and the spatial profile of the light beam near the focus (not to mention that the electron density radial profile is transient evolving during the laser pulse itself). However, an estimate for the maximum electron density that can be reached before defocusing clamps the laser intensity can be made by noting that the refractive index near the focus acquires a radially (and time) dependent profile:

$$n(r, t) = \sqrt{1 - \frac{n_e(r, t)}{n_{\text{crit}}}} \approx 1 - \frac{1}{2} \frac{n_e(r)}{n_{\text{crit}}} \quad (49)$$

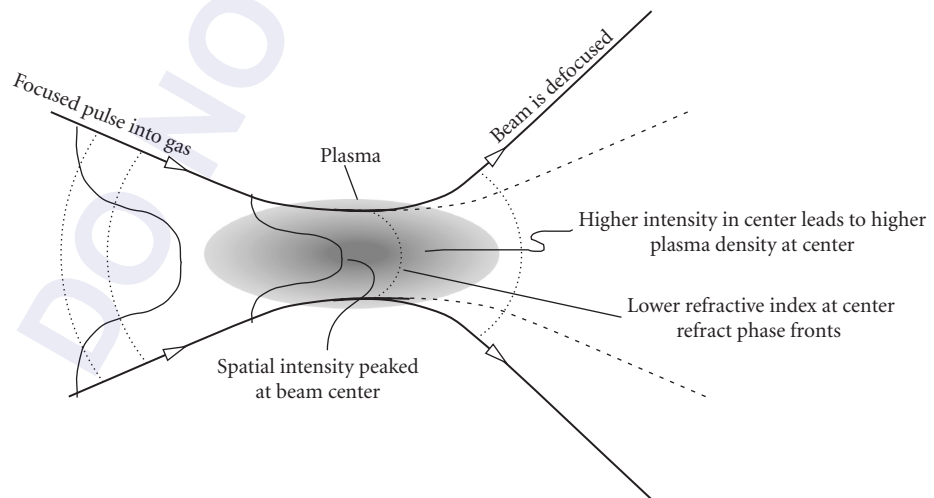


FIGURE 24 Description of how plasma formation at the focus of a laser can induce defocusing.

(Assuming that the electron density is well below the critical electron density, $n_{\text{crit}} = me\omega_0^2/4\pi e^2$ which is about 10^{21} cm^{-3} for 1- μm wavelength light.) Noting that the higher electron density acts as a negative lens, we can find the focal length of this negative lens by assuming that the laser propagates through the plasma over a length equal to its focused Rayleigh length $z_R = \pi w^2/\lambda$. In this case, ionization-induced defocusing will dominate the propagation of the focused pulse when

$$n_e(r=0) \approx \frac{\lambda}{\pi z_R} n_{\text{crit}} \quad (50)$$

When the on-axis electron density reaches this level, the pulse focusing will be clamped and will begin to defocus. For example, a 1- μm wavelength laser focused to a 10- μm spot will defocus when the on-axis electron density reaches about 10^{18} cm^{-3} . Creating any greater electron density will be difficult unless the gas geometry is fashioned so that the beam is focused to its tightest spot outside of the gas.

Relativistic Interactions in Plasma: Self-Channeling and Self-Phase Modulation

When the light pulse intensity is relativistic ($a_0 > 1$) the electrons oscillating in the strong field of the laser, acquire relativistic velocity during a single cycle. This means that the electron mass changes during the course of the light period affecting the optical properties of the plasma.²⁵ This can be seen by noting that the refractive index of a plasma, $n = \sqrt{1 - \omega_p^2/\omega_0^2}$, depends on the plasma frequency, which in turn depends on the square root of the electron mass. An increase in the effective mass of the electrons in the plasma through their relativistic oscillation will decrease the instantaneous value of the plasma frequency and, therefore, *increase* the value of the refractive index. Because the relativistic mass increase in the field depends on field amplitude (and hence intensity), the refractive index of the plasma in the relativistic regime is now intensity dependent.

Recalling the field-induced oscillatory gamma $\gamma_{\text{osc}} = \sqrt{1 + a_0^2/2}$ the refractive index in the plasma is now

$$n = \sqrt{1 - \frac{\omega_p^2}{\gamma_{\text{osc}} \omega_0^2}} \quad (51)$$

Many of the results of standard nonlinear optics can now be applied to this relativistic underdense plasma. In low-density plasma ($n_e \ll n_{\text{crit}}$) at modest a_0 (< 1) Taylor expansion of refractive index allows us to retrieve the refractive index in a standard nonlinear optics format in which $n = n_0 - n_2 I$,¹⁷⁵ where n_0 is the usual, nonrelativistic linear refractive index of the plasma, $n = \sqrt{1 - \omega_p^2/\omega_0^2}$ and the nonlinear refractive index is

$$n_2 = \frac{\pi e^2}{m_e^2 \omega_0^2 c^3} \frac{n_e}{n_c} \quad (52)$$

The intensity dependence of the plasma refractive index has two significant consequences, both with analogues in traditional nonlinear optics.

Relativistic Self-Focusing and Self-Channeling When a relativistically intense pulse, with a radially dependent intensity profile propagates through the plasma, relativistic effects can cause an increase in the refractive index on the center of the propagation axis, as illustrated in Fig. 25. This causes the phase fronts near the axis to be retarded with respect to those at the pulse's outer edges, resulting in self-focusing of the pulse as it passes through the plasma.^{20,227–230} Focusing will occur if the focusing power of the nonlinear refractive index is stronger than the natural defocusing of a pulse from diffraction of the wave (or any ionization-induced defocusing). Because both of these competing

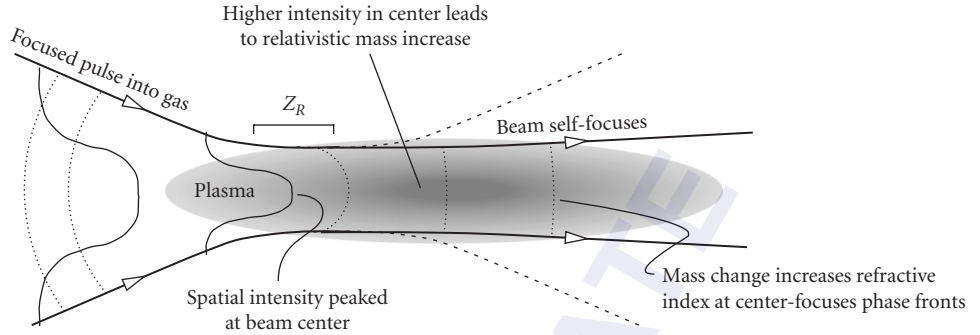


FIGURE 25 Description of relativistic induced self-focusing and channeling in an underdense plasma.

effects depend on the spot size (a larger spot means lower intensity and less relativistic self-focusing but, at the same time, a larger spot means less defocusing from diffraction). It can be shown that the threshold for self-focusing to occur is independent of intensity and depends only on the power of the light pulse. This critical power is²¹⁷

$$P_c = \frac{2c^5 m_e^2}{e^2} \left(\frac{\omega_0}{\omega_p} \right)^2 \quad (53)$$

or in practical units, (and in a form often quoted in the literature) $P_c = 17 (\omega_0/\omega_p)^2$ in GW. This means that if a 1- μm wavelength pulse propagates through a plasma at an electron density of 10^{18} cm^{-3} , the critical power for the pulse to undergo relativistic self-focusing is $\sim 20 \text{ TW}$.

If this power is achieved in the plasma, the pulse can potentially (though not necessarily) undergo relativistic channeling. The power of the laser must exceed P_c , but, in addition, the self-focal length must be fast enough to force a collapse of the pulse spatially within the experimental extent of the plasma (see Fig. 25). When these conditions are fulfilled, self-channeling occurs, permitting the intense pulse to propagate through the plasma at its focused diameter over a length which is many times that which would be allowed by simple diffraction (i.e., $\sim z_R$). Self-channeling in underdense plasma of distances up to 100 Rayleigh ranges has been observed.²³⁰ Using these criteria, and Eq. (52) the intensity required for self-channeling in a plasma of length l_p is $I > w^2/n_e I_p^2$. This relationship implies that a 1- μm wavelength pulse, focused to a radius of 10 μm in a plasma of $n_e \sim 10^{18} \text{ cm}^{-3}$ with length of 1 mm, will require a peak intensity of about 10^{18} W/cm^2 to undergo self-channeling.

Relativistic Self-Phase Modulation The second nonlinear optical consequence of the relativistic refractive index results from the temporal variation of the intensity of the pulse. This temporal variation leads to a time-dependent refractive index which, in turn, induces a time-varying phase on the pulse, a process known as relativistic self-phase modulation (SPM).^{36,231} As in the analogue of SPM in traditional nonlinear optical materials, the relativistic SPM can significantly broaden the spectrum of the pulse. The effects of SPM on a pulse are well treated in the optics literature and the relativistic case is essentially similar to these treatments utilizing Eq. (52). The spectrum of the pulse will be broadened by the SPM through the frequency shift acquired over a propagation distance dx

$$d\omega = -\frac{\omega_0}{c} n_2 \frac{\partial I}{\partial t} dx \quad (54)$$

This relationship shows that a strongly phase modulated pulse (initially transform limited) in a plasma of length, L will have its initial spectrum broadened significantly when $I > c/\omega_0 n_2 L$. This will

be a significant effect, for example, on a 1- μm wavelength pulse in a 1-mm-long plasma of 10^{18} cm^{-3} at an intensity of $> 2 \times 10^{18}\text{ W/cm}^2$ (independent of pulse duration).

21.10 STRONG FIELD PHYSICS AT SURFACES OF OVERDENSE PLASMAS

Structure of a Solid Target Plasma Irradiated at High Intensity

The irradiation of a solid target with a strong field laser pulse leads to an enormous plethora of effects, far too numerous to detail in this chapter. However, a common aspect of such interactions is that the intense laser pulse deposits energy very quickly into the solid, creating an overdense plasma ($n_e > n_{\text{crit}}$). This plasma expands and, depending on the pulse duration and the temporal structure of the laser pulse, can create an underdense region through which the intense pulse must propagate before reaching the critical density surface (where $n_e = n_{\text{crit}}$) and being reflected. The morphology of such interactions is illustrated in Fig. 26. Solid target interactions are accompanied by a range of mechanisms which serve to absorb laser energy, many of which result in the production of hot electrons. These hot electrons can often have energies in the relativistic range if $I \lambda^2$ of the laser is greater than $\sim 10^{17}\text{ W/cm}^2\text{-}\mu\text{m}^2$ (an a_0 of about 0.3).

It is usually convenient to think of strong field laser interactions with solid target plasmas in two regimes.

1. When the laser pulse has very little “prepulse” (i.e., laser energy which precedes the primary, typically Gaussian-shaped pulse, by many picoseconds or nanoseconds), and is very fast, the solid target retains a sharp density profile with electron density rising from vacuum to an over critical value within much less than a laser wavelength. This condition can be achieved if the plasma expands much less than a wavelength within the laser pulse duration, a condition which can be estimated by noting that the plasma expands like an ideal gas so a sharp plasma interface is retained if $\tau_p (k_B T_e / m_e)^{1/2} \ll \lambda$.
2. If the intense laser pulse is longer or is preceded by significant prepulse (intense enough to preionize the solid before the main pulse arrives), the solid plasma will expand, creating a lower density plasma scale length (as pictured in Fig. 26). In this case, the incoming laser must interact with some region of underdense plasma many wavelengths long, in which effects such as SRS or

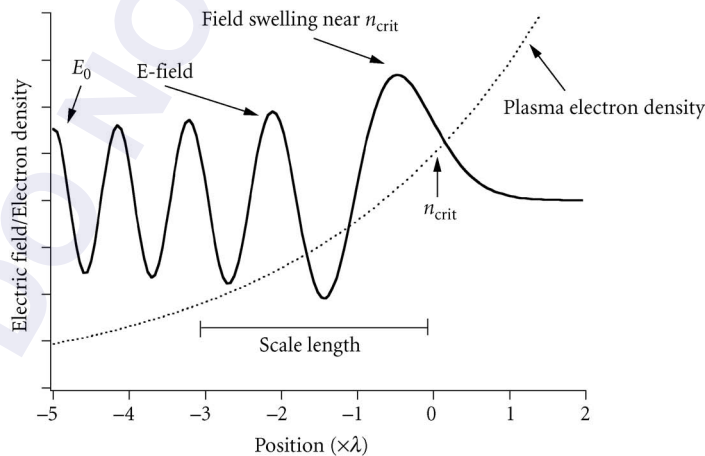


FIGURE 26 Illustration of the laser field as it propagates into a plasma density scale length up to the critical density at the surface of an overdense plasma.

self-channeling as discussed in Sec. 21.9 can take place. The specific profile of the plasma, and the shape of the incoming light field will be complicated, but it is useful to examine the solution when the plasma density profile is linear with position characterized by density scale length ℓ , (the length over which the electron density ramps up from 0 to n_{crit}). In this case, the incoming electric field of the laser is of the form¹⁸⁴

$$E(z) = 2\pi^{1/2}(\omega_0 \ell / c)^{1/6} E_0 \text{Ai}[(\omega_0^2 / c^2 \ell)^{1/3}(z - \ell)] \quad (55)$$

where E_0 is the strength of the incident field in vacuum and $\text{Ai}[x]$ is the well-known Airy function. The field swells near the critical density reflection point, which has a consequence on many strong field aspects of longer pulse interactions. The extent of this swelling is approximately given as $E_{\text{MAX}} \approx 1.9(\omega \ell / c)^{1/6} E_0$.

Resonance Absorption and Vacuum Heating

When an intense laser creates conditions like those shown in Fig. 26, various absorption mechanisms, almost all of them coupling laser light into plasma electron energy, rise up. Typically, all of these absorption mechanisms will conspire to create a situation in which the absorption of the laser by the solid is in the vicinity of 20 to 60 percent, depending on laser intensity and focal geometry.^{232,233} One particularly important absorption mechanism arises when the laser is obliquely incident on the target surface with p -polarization, a situation with a nonzero component of the laser field oscillating in and out of the target surface, as illustrated in Fig. 27. The component of the laser field perpendicular to the surface of the target E_{\perp} leads to electrons being driven across a density gradient, depositing energy into those electrons near the critical density surface. The nature of the laser coupling of energy to these electrons depends on whether the plasma interface has a long scale length or is sharp.

Resonance Absorption When a plasma scale length greater than the incident laser wavelength exists, E_{\perp} of the laser can drive resonant plasma oscillations at the critical density surface (where the

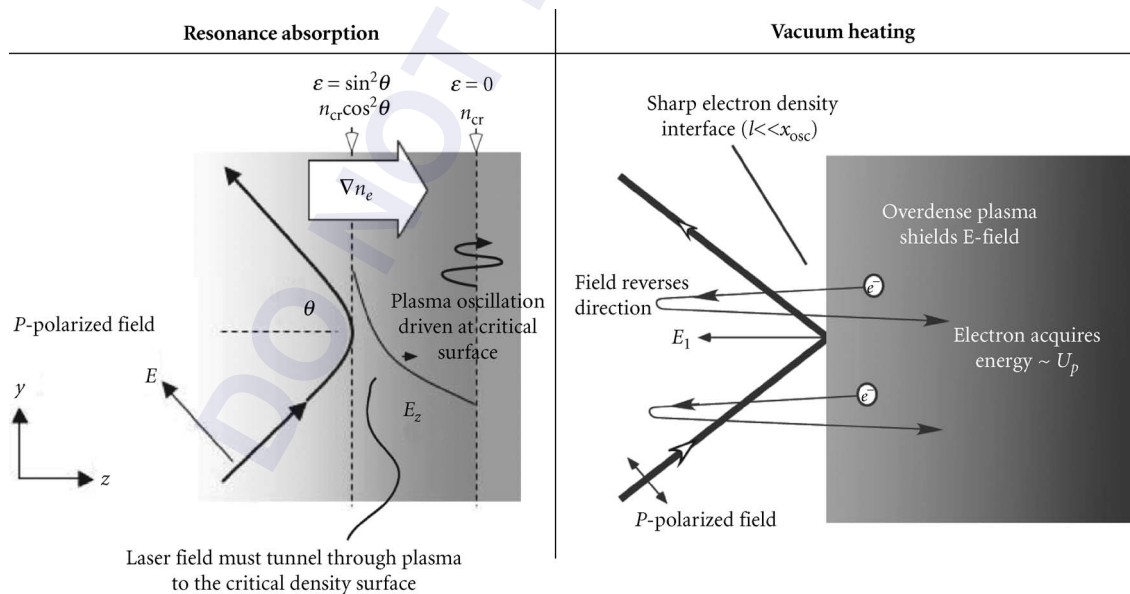


FIGURE 27 Illustration of resonance absorption on the left and vacuum heating on the right.

laser frequency is equal to the plasma oscillation frequency) along a direction parallel to the plasma gradient. These driven plasma oscillations lead to what is known as resonance absorption.^{184,234} This process, illustrated on the left in Fig. 26, leads to deposition of laser energy into an electron wave which can be damped by collisions or by other collisionless effects, such as Landau damping or wave breaking, in turn leading to the production of supra-thermal electrons.²³⁵ Because the driving field E_{\perp} determined two competing effects with incidence angle there is an optimum incidence angle for maximum resonant absorption. Greater incidence angle θ produces greater perpendicular component of E-field, however, the laser field reflects further away from the critical density surface at higher incidence angle¹⁸⁴ requiring the field to tunnel further to the resonant electron density at n_{crit} . Consequently, the efficiency of resonance absorption peaks at an angle given approximately by

$$\theta_{\text{MAX}} \approx \sin^{-1} \left[\left(\frac{c}{2\omega_0 \ell} \right)^{1/3} \right] \quad (56)$$

Energy deposited by resonance absorption per unit time in the plasma is linear with laser intensity, and the fraction of laser power absorbed in this manner, f_{RA} , can be shown to be equal to²² $f_{\text{RA}} \cong 2.6q^2 e^{4q^{7/3}}$, where $q = (\omega \ell / c)^{1/3} \sin \theta$. This absorption mechanism tends to produce a component of hot electrons which travel into the target in a direction normal to the target surface. The electrons usually exhibit an energy distribution much like that of a Maxwellian. Particle-in-cell (PIC) simulations and experiments have shown that, in the strong field, ultrashort pulse regime, temperature of this hot Maxwellian from resonance absorption is roughly, (in terms of indicated units), $T_{\text{hot}} [\text{keV}] \approx 2 \times 10^{-4} (T_0 [\text{keV}] I [\text{W}/\text{cm}^2] \lambda [\mu\text{m}]^2)^{1/3}$,²³⁶ where T_0 is the background plasma temperature. So a 1- μm laser at an intensity of $10^{18} \text{ W}/\text{cm}^2$ incident on a $1 \times \text{keV}$ plasma typically produces a distribution of hot electrons with temperature in the vicinity of 200 keV.

Vacuum Heating When the plasma density gradient is much smaller than a wavelength (or more accurately, when the plasma density gradient is smaller than the oscillation amplitude of electrons in E_{\perp}) a different description of the resonant heating is required. Instead, a situation like that shown on the right side of Fig. 27 exists. Here, an obliquely incident intense laser can pull electrons from the overdense plasma, accelerate them in vacuum and, once the oscillating electric field reverses direction, accelerate them back toward the overdense plasma interface. Once the electrons are propelled back into the overdense plasma, the laser field is shielded and the electrons retain a large fraction of the energy they acquired in their vacuum excursion. This process is commonly called vacuum heating or “Brunel absorption.”^{237–239} Like resonance absorption, this process requires p -polarized light so that $E_{\perp} \neq 0$, and results in fast electrons being accelerated in a direction perpendicular to the target surface. The electrons’ average energy will, naturally be roughly equal to U_p since this is about the energy they gain in their vacuum half-cycle oscillation. When the laser intensity is nonrelativistic, the absorbed power varies as the 3/2 power of intensity (unlike resonance absorption where the power deposition is linear with intensity). The fraction of absorbed laser power in vacuum heating in this case is

$$f_{\text{VH}} = \frac{4eE_0}{\pi m_e \omega_0 c} \left(\frac{\sin^3 \theta}{\cos \theta} \right) \quad (57)$$

where, again, θ is the incidence angle of the laser. This absorption is maximum at the highest angles and increases with the square root of the intensity. Note, however, that the *hottest electrons* are created when E_{\perp} is maximum, a condition which occurs at 55° if we assume nearly 100 percent reflection from the plasma surface. Note that the absorption fraction predicted in Eq. (57) can be >1 at high-intensity and incidence angle. This is because 100 percent reflection was assumed and no relativistic effects are included. Therefore, Eq. (57) is only appropriate for the low-absorption regime. In the strongly relativistic limit (when $a_0 \gg 1$) the dynamics of vacuum heating become

more complicated. It can be shown that in this regime the absorption approaches 100 percent at an optimum angle of 73° and that the fractional absorption is⁶⁶

$$f_{\text{VH}}^{\text{rel}} \approx \frac{4\pi(\sin^2\theta/\cos\theta)}{(\pi + \sin^2\theta/\cos\theta)^2} \quad (58)$$

Other Collisionless Absorption Mechanisms: $\mathbf{j} \times \mathbf{B}$ Heating and Anomalous Skin Effect

Resonance absorption or vacuum heating are usually the most important absorption mechanisms for an intense laser pulse incident on a solid-density plasma (though at near normal incidence, or with s -polarization, simple collisional absorption is most important when $\ell \gg \lambda$). However, there are other mechanisms which can lead to absorption of an intense light pulse on the surface of an overdense plasma (and the potential for producing hot electrons). In almost all cases these mechanisms rely on the presence of a steep density gradient, just as one finds in the vacuum heating regime discussed above. At modest intensity, defined as the regime in which $U_p < k_B T_e$, a phenomena known as the anomalous skin effect takes place.²⁴⁰ Absorption through this effect essentially results when electrons in the plasma have a large mean free path when compared to the skin depth (c/ω_p) in the plasma. Electrons migrate via their thermal motion into this skin depth region and acquire energy from the electric field's evanescent wave in the plasma. This absorption effect is found to play a minor role in short-pulse laser-plasma interactions (with absorption typically less than 5 percent) and, because of the requirement of low U_p , is usually not important for laser intensity above about 10^{15} W/cm². Unlike RA or vacuum heating, this absorption mechanism does not require oblique incidence, and therefore, can play a nonnegligible role for moderate intensity, normal incidence situations.

A much more important absorption mechanism at a sharp dense plasma interface occurs in the relativistic intensity regime. In this regime, the magnetic field of the light pulse begins to play a significant role in the trajectory of the driven plasma electrons. Plasma electrons at the surface of the plasma, driven transversely by the E-field with velocity approaching c , feel an oscillating force in the direction of the laser's \mathbf{k} vector from the now significant $\mathbf{V}_{\text{osc}} \times \mathbf{B}$ force.^{241,242} The electrons acquire energy in a manner similar to vacuum heating except now the oscillatory driving force is the $\mathbf{j} \times \mathbf{B}$ term, oscillating at a frequency of $2\omega_0$. Oblique incidence is not necessary for this $\mathbf{j} \times \mathbf{B}$ heating to occur; the heating, in fact, is maximum at normal incidence (where E-field driven vacuum heating vanishes). As a consequence, at relativistic intensity, the relative importance of vacuum heating and $\mathbf{j} \times \mathbf{B}$ heating depends on incidence angle, polarization (s - or p -) and intensity. In general, $\mathbf{j} \times \mathbf{B}$ heating will dominate when $v_{\text{osc}} B > 2E_0 \sin \theta$, even with p -polarized light, which is equivalent to $a_0 > 2^{1/2}$ or an intensity greater than $\sim 4 \times 10^{18}$ W/cm² in a 1- μm wavelength field at $\theta = 45^\circ$. There are situations when both vacuum heating and $\mathbf{j} \times \mathbf{B}$ occur, with vacuum heating producing hot electrons normal to the target and $\mathbf{j} \times \mathbf{B}$ producing electrons parallel to the laser direction. Again, the energy of the accelerated electrons from $\mathbf{j} \times \mathbf{B}$ heating will roughly equal the relativistic ponderomotive potential, $(\gamma_{\text{osc}} - 1) m_e c^2$.²²

Ponderomotive Steepening and Hole Boring

An intense pulse incident on a solid density plasma surface will exert a force on that plasma. Because the pulse will have some temporal rise, there is naturally a gradient in the ponderomotive energy $-\nabla U_p$, which will drive the plasma inward, as illustrated in Fig. 28.²⁴² This ponderomotive gradient is really nothing other than the manifestation of the pressure associated with the light pulse I/c , which can be very high. The light pressure of a pulse at an intensity of 10^{18} W/cm² is ~ 330 Mbar, equivalent to the thermal pressure of a plasma at solid density of $n_e \sim 10^{23}$ cm⁻³ when the electron temperature is 2 keV, well above the thermal temperatures of most laser plasmas. At modest intensities, the ponderomotive gradient will simply slow the expansion of the plasma outward. However, when $I/c > n_e k_B T_e$, the laser pulse will push the plasma back toward higher density, decreasing the plasma density

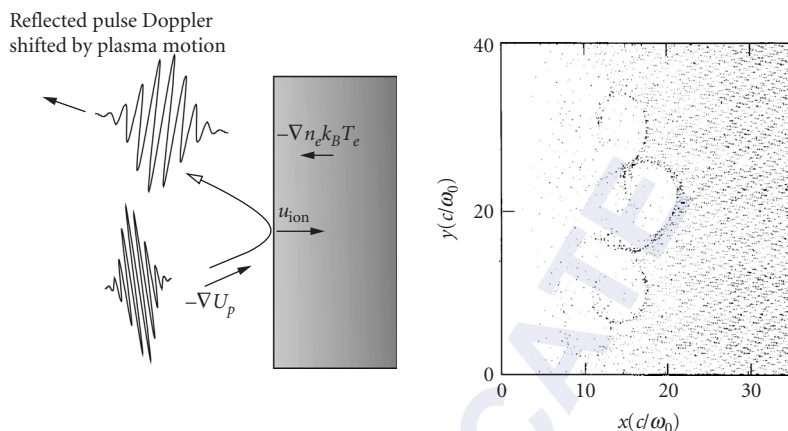


FIGURE 28 On the left is an illustration of how ponderomotive steepening leads to hole boring when intense pulse is focused on a solid target plasma. On the right, a PIC simulation from Ref. 242 is reproduced showing the bubblelike structure that develops in the plasma ions at the surface of a plasma as a result of hole boring by a picosecond duration intense laser.

gradient, a process referred to as ponderomotive steepening. This reversal of the plasma velocity will manifest itself as a red shift in the spectrum of the reflected laser pulse or its harmonics.^{243,244}

At very high intensities, such as 10^{19} W/cm² or greater, the light pressure can be sufficient to push the plasma, even over the short timescales of laser pulses usually employed, back into the target. The high-intensity limit of ponderomotive steepening is commonly called hole boring, a process pictured in the PIC simulations of Fig. 28. Hole boring drives ions into the plasma at high velocity. This ion velocity can be estimated by balancing the inward ion momentum flux with the light pressure. Accounting for some incidence angle θ and some potential fractional absorption, f_{abs} , the inward ion velocity is

$$u_{\text{ion}} = \sqrt{\frac{(2 - f_{\text{abs}}) I \cos \theta}{2n_i m_i c}} \quad (59)$$

Equation (59) predicts, for example, that a solid aluminum plasma irradiated at 10^{19} W/cm² can acquire ion velocities during hole boring of $\sim 4 \times 10^7$ cm/s. It is possible that the high energy ions driven inward at ultrahigh intensities ($> 10^{20}$ W/cm²) can produce a collisionless shock in the underlying overdense plasma.²⁴⁵ The amount of laser energy absorbed from the laser into these fast ions can be easily estimated and is (for small absorption) $f_{\text{hole boring}} \approx 2u_{\text{ion}}/c$. This relation suggests that a few percent of the incident laser energy can be transferred to ions through ponderomotive acceleration.

High Harmonic Generation from Solid Plasmas

The light reflected from the surface of an overdense plasma at modest intensity will typically retain with good fidelity the spectrum of the incident laser light, though at higher intensity this spectrum can be Doppler shifted from ponderomotive steepening or hole boring. At these high intensities, however, nonlinear interactions can occur at the reflection surface, in such a way that harmonics of the light field are generated.^{246–249} At very high intensity ($> \sim 10^{19}$ W/cm²) quite high orders can be generated at the plasma surface, up to orders of $q \sim 50$ to 100 .²⁴⁸

The origin of these harmonics can be seen by noting that the ponderomotive pressure at the surface of the plasma is $\sim \nabla a^2/2$, which not only has a time averaged term (the force which gives

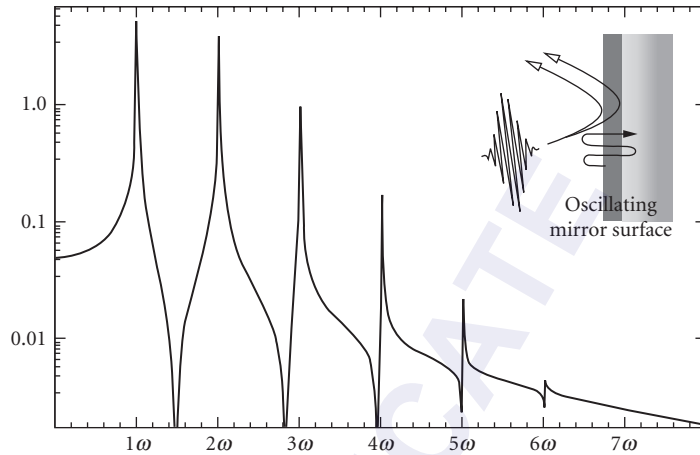


FIGURE 29 Plot of the Fourier transform of Eq. (61) with $kA_{\text{mirror}} = 1$, showing how harmonics can arise from the oscillating mirror model. The inset shows the geometry of the oscillating mirror.

rise to the hole boring discussed above) but also has a fast time varying component which varies as $p_{\text{osc}} \sim a_0^2 \sin 2\omega_0 t$. This rapidly oscillating force will drive the plasma electrons collectively in and out of the plasma. A similar situation occurs for p -polarized light incident on a solid plasma at an angle, though now the plasma electrons are driven at the laser frequency by the laser's electric field (see the inset to Fig. 29). If the intense light pulse is incident on a sharp-density gradient these driving forces effectively act to create an oscillating critical density surface, or, in other words, an oscillating mirror.^{246,247} It is the oscillation of the reflection point in the plasma which yields the nonlinearity leading to high harmonic generation at the plasma surface.

The theory describing this situation is complex, particularly in the relativistic limit, however, a simple analysis of the oscillating mirror amplitude, A_{mirror} , in the nonrelativistic limit, shows that⁶⁶

$$A_{\text{mirror}} \cong \frac{\lambda_0}{\pi} \left(\frac{\omega_0}{\omega_p} \right)^3 a_0^2 \quad (60)$$

This equation illustrates that the efficiency of creating these surface high harmonics increases with increasing intensity and with lower surface plasma density (because a lower density plasma is easier to drive by the ponderomotive forces of the laser). Simulations have shown that high harmonic generation from dense plasma surfaces is maximized with the target density in the vicinity of $4n_{\text{crit}}$,²⁴⁷ a consequence of an interaction between the second harmonic of the laser and the target surface, leading to violent resonant driving of the plasma surface. Equation (60) also shows that the mirror amplitude begins to approach the scale of one laser wavelength when $a_0 \rightarrow (n/n_{\text{crit}})^{3/4}$. In the nonrelativistic limit when the mirror oscillates with the laser frequency (the case for obliquely incident p -polarized light), the reflected light will have an electric field at the plasma surface ($z = 0$) such that

$$E_{\text{refl}}(t) \cong \frac{\omega_p}{2\omega_0} a_0 \sin(\omega_0 t) + kA_{\text{mirror}} \sin(\omega_0 t) \quad (61)$$

A Fourier transform of this field when $kA_{\text{mirror}} = 1$ is illustrated in Fig. 29, showing how the nonlinearity introduced by the oscillating mirror amplitude gives rise to numerous harmonics (at both even and odd harmonics of the laser in this case).

In the strongly relativistic limit, the theory becomes much more difficult. PIC simulations have shown, however, that the efficiency of the high orders ($q \gg 1$) can be estimated with the empirical relation $\eta_q \sim I \lambda^2 (1/q)^5$.²⁴⁹ As this indicates the high harmonic spectrum from solids does not exhibit a plateau and cutoff, but instead, harmonics are created in a spectrum which simply falls off with order as a power law. The yield of the harmonic becomes nearly linear with intensity.

Relativistic Effects and Induced Transparency

As in underdense plasmas, relativistic effects can affect the optical properties of overdense plasmas as well. Self-induced relativistic transparency is one such effect.²⁵⁰ At relativistic intensity, the mass change of the electrons leads to an effective shift in the plasma frequency, such that

$$\omega_p^{\text{rel}} = \omega_p / \gamma^{1/2} = \sqrt{\frac{4\pi e^2}{m_e}} \frac{1}{(1+a_0^2/2)^{1/4}} \quad (62)$$

Since ω_p^{rel} drops with the square root of intensity, an overdense, reflecting plasma can be shifted by the relativistic factor in Eq. (62) to one which becomes effectively underdense, with $\omega_p^{\text{rel}} < \omega_0$, allowing the high-intensity light to propagate through the plasma. The critical intensity for this to occur, in the strong relativistic limit ($a_0 \gg 1$) is

$$a_0^{\text{trans}} = \sqrt{2} \frac{n_e}{n_{\text{crit}}} \quad (63)$$

Consequently, a solid density plasma, with $n_e \sim 10^{23} \text{ cm}^{-3}$, will become transparent to 1- μm laser light at an a_0 of ~ 140 . This is an intensity near $3 \times 10^{22} \text{ W/cm}^2$, just within reach of the highest intensity lasers now available. That this relativistically induced transparency can occur in a plasma was mentioned in the literature many years ago, however, to date there has yet to be a direct experimental observation of this effect.

21.11 APPLICATIONS OF STRONG FIELD INTERACTIONS WITH PLASMAS

Modern strong field research now focuses on a range of effects that arise from the basic phenomena that we have described in this chapter. It is impossible to summarize all of these research avenues here, but we conclude in this section with a very brief sampling of some of the applications of high-field research that are among the most active at this writing.

Femtosecond X-Ray Production

Many of the strong field effects discussed in previous sections result in the production of energetic electrons (through, e.g., plasma or free wave acceleration via wakefield, SRS, resonant absorption, or vacuum and $\mathbf{j} \times \mathbf{B}$ heating). When electrons are accelerated by these mechanisms at the surface of a cold, solid target, the penetration of these hot electrons will lead to x-ray production. If the laser generating these fast electrons is short, it is possible to produce bright sources of x-rays with femtosecond time duration, which, can, in turn, be used for various pump-probe applications.³⁸ There are generally three principal ways in which x-rays are generated in these experiments: (1) through the direct interaction with the laser field or other strong fields in and around a plasma; (2) through the ejection of inner shell electrons in the underlying cold solid leading to K_α line emission; or

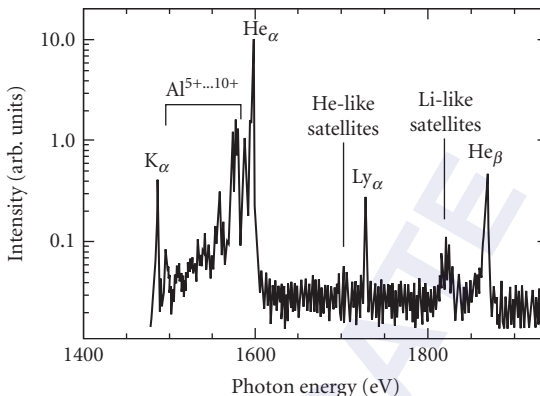


FIGURE 30 An example of an x-ray spectrum produced by irradiation of a solid Al target at intensity of $\sim 10^{18}$ W/cm². This spectrum shows that not only are H-like and He-like lines emitted from the hot plasma, but also that K_α x-rays are emitted when the hot electrons produced by the laser penetrate the cold target and knock out K-shell holes in the unionized Al atoms. This x-ray pulse is usually very fast, <1 ps in duration. (*Adapted from Ref. 252.*)

(3) through free-free transitions of the electrons as they pass by highly charged nuclei of the solid, yielding bremsstrahlung radiation. The later process has been shown to yield x-rays with energies well above 1 MeV when intensities of $>10^{18}$ W/cm² are used to irradiate a solid target.²⁵¹ A characteristic x-ray spectrum resulting from an example experiment is shown in Fig. 30.

Fusion Neutron Production

The production of fast ions in strong field interactions, through processes such as ponderomotive acceleration of ions at the surface of a target or through the ejection of ions from the Coulomb explosion of irradiated clusters, can be harnessed to drive nuclear fusion. The production of bursts of 2.45 MeV neutrons when various targets contain deuterium (solid or cluster targets) has been well studied by experiment.^{199,253}

High Magnetic Field Production

In addition to the very strong magnetic fields associated with high intensity EM waves, the intense irradiation of solid target plasmas can produce strong transient DC magnetic fields. These fields essentially arise from two sources: (1) Thermo-electric magnetic field generation. In the high temperatures and steep density gradients found in intense ultrafast laser irradiation of solids, magnetic fields can be generated by thermal transport effects. In this case the magnetic field is proportional to $\nabla T_e \times \nabla n_e$ which favors high B-field production in the steep density gradients found in femtosecond laser plasma interactions.²⁵⁴ (2) Fast electron magnetic field generation. The fast electrons produced by the panoply of mechanisms discussed in previous sections, frequently lead to the production of high peak currents, perhaps exceeding a MA. This leads to enormous B-fields. Fields of upto many 100s of MG are possible with this mechanism.²⁵⁵

MeV Proton Acceleration

When fast electrons are generated at the surface of a solid target plasma by an intense laser and the target is sufficiently thin such that the electron mean free path is greater than the target thickness, the fast electrons can exit the back surface and set up a strong ambipolar field. This situation is illustrated in Fig. 31. These electrons will set up an electric field which ionizes ions on the back surface (often covered with water and hydrocarbon contamination in vacuum) and accelerate them to high velocity.^{256–258} The field created will be of the order of $eE_{\text{sheath}} \approx k_B T_{\text{hot}} / \lambda_{\text{Debye}}^{\text{hot}}$, where $\lambda_{\text{Debye}}^{\text{hot}} = \sqrt{k_B T_{\text{hot}} / 4\pi n_{\text{hot}} e^2}$ is the hot electron Debye radius and T_{hot} and n_{hot} are the temperature and density of the hot electrons. The sheath field can often be of the order of a few MeV/ μm , so protons can easily acquire many MeV of energy during their acceleration in the direction of the target back surface normal. This ion acceleration mechanism is termed “target normal sheath acceleration” (TNSA). Ion temperatures of a few MeV are often observed when sub-picosecond (200 to 1000 fs) pulses are focused at intensity in the vicinity of 10^{20} W/cm². It has been found that this TNSA mechanism is most effective for higher energy pulses (>1 J) in the >500 fs range since such pulses tend to produce a greater number of hot electrons than sub-100 fs pulses at similar intensities.

Fast Ignition

The production of fast electrons by an intense short pulse laser has been proposed as a means to ignite a fusion capsule imploded by a large scale laser or Z-pinch, in the so-called inertial confinement approach to fusion (ICF).²⁵⁹ One variant of ICF, known as fast ignition, is described in Fig. 32. The basic idea is that a separate, high energy driver will assemble at high density (>200 g/cm³) a deuterium/tritium fusion fuel, and the short pulse laser will be focused from the side, onto the edge of the fusion fuel, generating a beam of hot electrons with temperature in the range of 1 to 10 MeV which penetrate

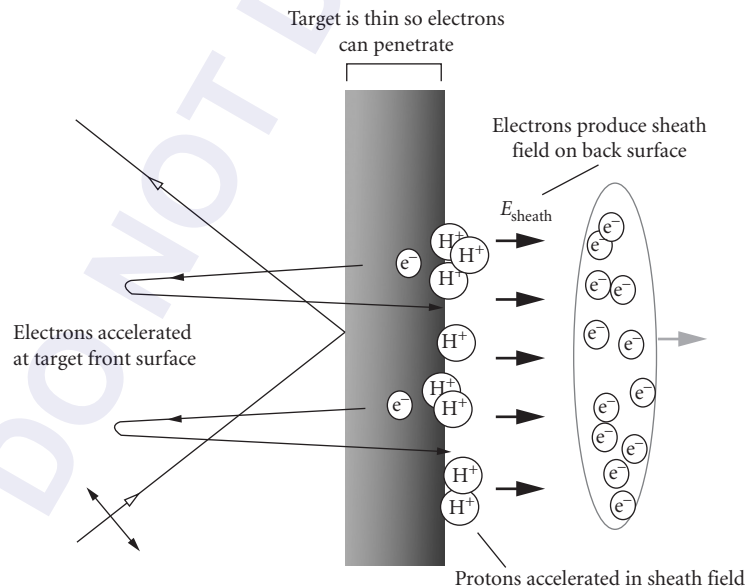


FIGURE 31 Illustration of how protons can be accelerated from the back surface of an intensely irradiated target by the target normal sheath acceleration phenomenon.

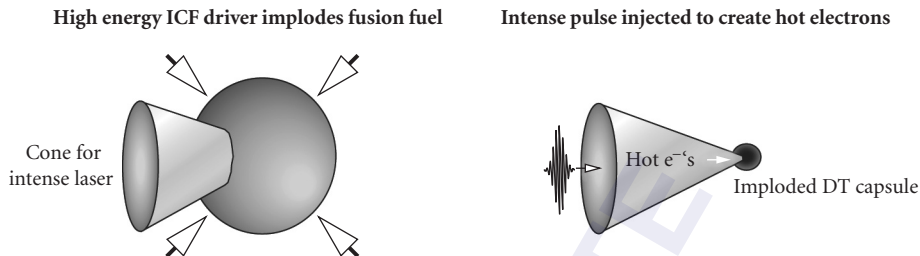


FIGURE 32 Illustration of how fast ignition would utilize an intense pulse to ignite an imploded inertial confinement fusion capsule.

the fuel and heat it to fusion ignition temperature. At present, the most promising method of doing this is to embed a cone in the fusion fuel prior to its compression (see Fig. 32) to permit effective injection of the intense laser at the peak of the fuel's compression. There remain many technical challenges to the ultimate implementation of this idea, with numerous physics issues remaining unsolved (such as understanding how such high peak currents of electrons transport into such a dense plasma or what fraction of short pulse laser energy converts to hot electrons). Nonetheless, there have been experimental studies conducted in Japan in recent years which have yielded promising results on the fast ignition approach at modest laser compression and short pulse laser drive energy.²⁶⁰

Raman Amplification

It has been recently demonstrated that it is possible to amplify ultrashort pulses in a plasma by counterpropagating a long pulse laser with an intense short pulse laser of lower frequency.²⁶¹ This superradiant amplification essentially occurs through a parametric process in the plasma closely akin to stimulated Raman scattering. With this high-intensity technique, it may be possible using kJ-class nanosecond lasers to amplify kJ-energy sub-100 fs pulses in a plasma.

21.12 REFERENCES

1. T. F. Gallagher, "Rydberg Atoms in Strong Microwave Fields," in *Atoms in Intense Laser Fields*, M. Gavrila, (ed.), Academic Press, San Diego, 1992, pP. 67–108.
2. L. V. Keldysh, "Ionization in the Field of a Strong Electromagnetic Wave," *Soviet Physics JETP* **20**:1307–1314 (1965).
3. P. Agostini, F. Fabre, G. Mainfray, G. Petite, and N. K. Rahman, "Free-Free Transitions Following 6-Photon Ionization of Xenon Atoms," *Physical Review Letters* **42**:1127–1130 (1979).
4. F. Fabre, G. Petite, P. Agostini, and M. Clement, "Multi-Photon above-Threshold Ionization of Xenon at 0.53 and 1.06- μm ," *Journal of Physics B-Atomic Molecular and Optical Physics* **15**:1353–1369 (1982).
5. A. McPherson, G. Gibson, H. Jara, U. Johann, T. S. Luk, I. A. McIntyre, K. Boyer, and C. K. Rhodes, "Studies of Multiphoton Production of Vacuum Ultraviolet-Radiation in the Rare-Gases," *Journal of the Optical Society of America B-Optical Physics* **4**:595–601 (1987).
6. Z. H. Chang, A. Rundquist, H. W. Wang, M. M. Murnane, and H. C. Kapteyn, "Generation of Coherent Soft X Rays at 2.7 nm Using High Harmonics," *Physical Review Letters* **79**:2967–2970 (1997).
7. A. L'Huillier and P. Balcou, "High-Order Harmonic-Generation in Rare-Gases with a 1-Ps 1053-Nm Laser," *Physical Review Letters* **70**:774–777 (1993).
8. F. Krausz and P. Corkum, "Research Supports Observation of Attosecond Pulses," *Laser Focus World* **38**:7–7 (2002).
9. M. Hentschel, R. Kienberger, C. Spielmann, G. A. Reider, N. Milosevic, T. Brabec, P. Corkum, U. Heinzmann, M. Drescher, and F. Krausz, "Attosecond Metrology," *Nature* **414**:509–513 (2001).

10. P. B. Corkum, "Plasma Perspective on Strong-Field Multiphoton Ionization," *Physical Review Letters* **71**:1994–1997 (1993).
11. K. C. Kulander, K. J. Schafer, and J. L. Krause, in *Super-Intense Laser Atom Physics*, NATO ASI Ser. , A. L. H. B. Piraux, and K. Rzaewski, (eds.), Plenum Press, New York, 1993.
12. H. Daido, "Review of Soft X-Ray Laser Researches and Developments," *Reports on Progress in Physics* **65**:1513–1576 (2002).
13. J. J. Rocca, "Table-Top Soft X-Ray Lasers," *Review of Scientific Instruments* **70**:3799–3827 (1999).
14. V. Malka, J. Faure, Y. A. Gauduel, E. Lefebvre, A. Rousse, and K. T. Phuoc, "Principles and Applications of Compact Laser-Plasma Accelerators," *Nature Physics* **4**:447–453 (2008).
15. T. Tajima and J. M. Dawson, "Laser Electron-Accelerator," *Physical Review Letters* **43**:267–270 (1979).
16. S. Tazzari and M. Ferrario, "Trends in High Energy Particle Accelerators," *Reports on Progress in Physics* **66**:1045–1094 (2003).
17. J. H. Eberly, J. Javanainen, and K. Rzaewski, "Above-Threshold Ionization," *Physics Reports-Review Section of Physics Letters* **204**:331–383 (1991).
18. K. Burnett, V. C. Reed, and P. L. Knight, "Atoms in Ultra-Intense Laser Fields," *Journal of Physics B-Atomic Molecular and Optical Physics* **26**:561–598 (1993).
19. T. Ditmire, "Atomic Clusters in Ultrahigh Intensity Light Fields," *Contemporary Physics* **38**:315–328 (1997).
20. E. Esarey, P. Sprangle, J. Krall, and A. Ting, "Self-Focusing and Guiding of Short Laser Pulses in Ionizing Gases and Plasmas," *IEEE Journal of Quantum Electronics* **33**:1879–1914 (1997).
21. M. Protopapas, C. H. Keitel, and P. L. Knight, "Atomic Physics with Super-High Intensity Lasers," *Reports on Progress in Physics* **60**:389–486 (1997).
22. S. C. Wilks and W. L. Kruer, "Absorption of Ultrashort, Ultra-Intense Laser Light by Solids and Overdense Plasmas," *IEEE Journal of Quantum Electronics* **33**:1954–1968 (1997).
23. P. Lambropoulos, P. Maragakis, and J. Zhang, "Two-Electron Atoms in Strong Fields," *Physics Reports-Review Section of Physics Letters* **305**:203–293 (1998).
24. T. Brabec and F. Krausz, "Intense Few-Cycle Laser Fields:Frontiers of Nonlinear Optics," *Reviews of Modern Physics* **72**:545–591 (2000).
25. D. Umstadter, "Review of Physics and Applications of Relativistic Plasmas Driven by Ultra-Intense Lasers," *Physics of Plasmas* **8**:1774–1785 (2001).
26. R. Dorner, T. Weber, M. Weckenbrock, A. Staudte, M. Hattass, H. Schmidt-Bocking, R. Moshhammer, and J. Ullrich, "Multiple Ionization in Strong Laser Fields," *Advances in Atomic, Molecular, and Optical Physics* **48**:1–34 (2002).
27. M. Gavrilá, "Atomic Stabilization in Superintense Laser Fields," *Journal of Physics B-Atomic Molecular and Optical Physics* **35**:R147–R193 (2002).
28. A. Pukhov, "Strong Field Interaction of Laser Radiation," *Reports on Progress in Physics* **66**:47–101 (2003).
29. J. Ullrich, R. Moshhammer, A. Dorn, R. Dorner, L. P. H. Schmidt, and H. Schmiidt-Bocking, "Recoil-Ion and Electron Momentum Spectroscopy:Reaction-Microscopes," *Reports on Progress in Physics* **66**:1463–1545 (2003).
30. D. Umstadter, "Relativistic Laser-Plasma Interactions," *Journal of Physics D-Applied Physics* **36**:R151–R165 (2003).
31. P. Agostini and L. F. DiMauro, "The Physics of Attosecond Light Pulses," *Reports on Progress in Physics* **67**:813–855 (2004).
32. R. Bingham, J. T. Mendonca, and P. K. Shukla, "Plasma Based Charged-Particle Accelerators," *Plasma Physics and Controlled Fusion* **46**:R1–R23 (2004).
33. T. Ditmire, S. Bless, G. Dyer, A. Edens, W. Grigsby, G. Hays, K. Madison, et al., "Overview of Future Directions in High Energy-Density and High-Field Science Using Ultra-Intense Lasers," *Radiation Physics and Chemistry* **70**:535–552 (2004).
34. V. S. Popov, "Tunnel and Multiphoton Ionization of Atoms and Ions in a Strong Laser Field (Keldysh Theory)," *Physics-Uspokhi* **47**:855–885 (2004).
35. J. H. Posthumus, "The Dynamics of Small Molecules in Intense Laser Fields," *Reports on Progress in Physics* **67**:623–665 (2004).
36. D. Umstadter, S. Sepker, and S. Y. Chen, "Relativistic Nonlinear Optics," *Advance in Atomic, Molecular, and Optical Physics* **52**:331–389 (2005).

37. M. Borghesi, J. Fuchs, S. V. Bulanov, A. J. Mackinnon, P. K. Patel, and M. Roth, "Fast Ion Generation by High-Intensity Laser Irradiation of Solid Targets and Applications," *Fusion Science and Technology* **49**:412–439 (2006).
38. T. Pfeifer, C. Spielmann, and G. Gerber, "Femtosecond X-Ray Science," *Reports on Progress in Physics* **69**:443–505 (2006).
39. U. Saalmann, C. Siedschlag, and J. M. Rost, "Mechanisms of Cluster Ionization in Strong Laser Pulses," *Journal of Physics B-Atomic Molecular and Optical Physics* **39**:R39–R77 (2006).
40. Y. I. Salamin, S. X. Hu, K. Z. Hatsagortsyan, and C. H. Keitel, "Relativistic High-Power Laser-Matter Interactions," *Physics Reports-Review Section of Physics Letters* **427**:41–155 (2006).
41. R. A. Ganeev, "High-Order Harmonic Generation in a Laser Plasma: A Review of Recent Achievements," *Journal of Physics B-Atomic Molecular and Optical Physics* **40**:R213–R253 (2007).
42. D. Strickland and G. Mourou, "Compression of Amplified Chirped Optical Pulses," *Optics Communication* **56**:219–221 (1985).
43. W. Koehner, *Solid-State Laser Engineering*, Springer Verlag, Berlin, 1996.
44. S. W. Bahk, P. Rousseau, T. A. Planchon, V. Chvykov, G. Kalintchenko, A. Maksimchuk, G. A. Mourou, and V. Yanovsky, "Characterization of Focal Field Formed by a Large Numerical Aperture Paraboloidal Mirror and Generation of Ultra-High Intensity (10(22)W/cm(2)) (vol 80, p. 823, 2005)," *Applied Physics B-Lasers and Optics* **81**:727–727 (2005).
45. P. A. Norreys, K. M. Krushelnick, and M. Zepf, "PW Lasers: Matter in Extreme Laser Fields," *Plasma Physics and Controlled Fusion* **46**:B13–B21 (2004).
46. C. N. Danson, J. Collier, D. Neely, L. J. Barzanti, A. Damerell, C. B. Edwards, M. H. R. Hutchinson, et al., "Well Characterized 10(19) W cm(2) Operation of VULCAN—An Ultra-High Power Nd: Glass Laser," *Journal of Modern Physics* **45**:1653–1669 (1998).
47. Y. Kitagawa, Y. Sentoku, S. Akamatsu, M. Mori, Y. Tohyama, R. Kodama, K. A. Tanaka, et al., "Progress of Fast Ignitor Studies and Petawatt Laser Construction at Osaka University," *Physics of Plasmas* **9**:2202–2207 (2002).
48. F. G. Patterson and M. D. Perry, "Design and Performance of a Multiterawatt, Subpicosecond Neodymium—Glass-Laser," *Journal of the Optical Society of America B-Optical Physics* **8**:2384–2391 (1991).
49. N. Blanchot, C. Rouyer, C. Sauteret, and A. Migus, "Amplification of Sub-100 TW Femtosecond Pulses by Shifted Amplifying Nd:Glass Amplifiers: Theory and Experiment," *Optics Letter* **20**:395–397 (1995).
50. M. D. Perry, D. Pennington, B. C. Stuart, G. Tietbohl, J. A. Britten, C. Brown, S. Herman, B. et al., "Petawatt Laser Pulses," *Optics Letter* **24**:160–162 (1999).
51. R. F. Service, "Laser Labs Race for the Petawatt," *Science* **301**:154–156 (2003).
52. K. Yamakawa, S. Matsuoka, M. Aoyama, T. Kase, Y. Akahane, H. Takuma, and C. P. J. Barty, "Design and Performance of a 100 TW, Sub-20 fs Ti: Sapphire Laser System," in *X-Ray Lasers 1998*. 1999. pp. 645–648.
53. B. C. Walker, C. Toth, D. N. Fittinghoff, T. Guo, D. E. Kim, C. Rose-Petrucci, J. A. Squier, K. Yamakawa, K. R. Wilson, and C. P. J. Barty, "A 50-EW/cm(2) Ti: Sapphire Laser System for Studying Relativistic Light-Matter Interactions," *Optics Express* **5**:196–202 (1999).
54. F. G. Patterson, J. Bonlie, D. Price, B. W. Hite, and P. Springer, *LLNL Internal Report UCRL-JC-134912*: (1999).
55. V. Yanovsky, V. Chvykov, G. Kalinchenko, P. Rousseau, T. Planchon, T. Matsuoka, A. Maksimchuk, et al., "Ultra-High Intensity-300-TW Laser at 0.1 Hz Repetition Rate," *Optics Express* **16**:2109–2114 (2008).
56. S. Backus, C. G. Durfee, M. M. Murnane, and H. C. Kapteyn, "High Power Ultrafast Lasers," *Review of Scientific Instruments* **69**:1207–1223 (1998).
57. K. Yamakawa and C. P. J. Barty, "Ultrafast, Ultrahigh-Peak, and High-Average Power Ti: Sapphire Laser System and Its Applications," *IEEE Journal of Selected Topics in Quantum Electronics* **6**:658–675 (2000).
58. I. N. Ross, P. Matousik, M. Towrie, A. J. Langley, and J. L. Collier, "The Prospects for Ultrashort Pulse Duration and Ultrahigh Intensity Using Optical Parametric Chirped Pulse Amplifiers," *Optics Communication* **144**:125–133 (1997).
59. I. N. Ross, J. L. Collier, P. Matousek, C. N. Danson, D. Neely, R. M. Allot, D. A. Pepler, C. Hernandez-Gomez, and K. Osvay, "Generation of Terawatt Pulses by Use of Optical Parametric Chirped Pulse Amplification," *Applied Optics* **39**:2422–2427 (2000).
60. I. Jovanovic, C. A. Ebberts, and C. P. J. Barty, "Hybrid Chirped-Pulse Amplification," *Optics Letter* **27**:1622–1624 (2002).

61. Y. X. Leng, L. H. Lin, X. D. Yang, H. H. Lu, Z. Q. Zhang, and Z. Z. Xu, "Regenerative Amplifier with Continuously Variable Pulse Duration Used in an Optical Parametric Chirped-Pulse Amplification Laser System for Synchronous Pumping," *Optical Engineering* **42**:862–866 (2003).
62. J. E. Gunn and J. P. Ostriker, "On the Motion and Radiation of Charged Particles in Strong Electromagnetic Wave. I. Motion in Plane and Spherical Waves," *The Astrophysical Journal* **165**:523–541 (1971).
63. J. H. Eberly and A. Sleeper, "Trajectory and Mass Shift of a Classical Electron in a Radiation Pulse," *Physical Review* **176**:1570–1573 (1968).
64. F. V. Hartemann, "High-Intensity Scattering Processes of Relativistic Electrons in Vacuum," *Physics of Plasmas* **5**:2037–2047 (1998).
65. J. N. Bardsley, B. M. Penetrante, and M. H. Mittleman, "Relativistic Dynamics of Electrons in Intense Laser Fields," *Physical Review A* **40**:3823–3836 (1989).
66. P. Gibbon, *Short Pulse Laser Interaction with Matter: An Introduction*, Imperial College Press, London, 2005.
67. D. D. Meyerhofer, "High-Intensity-Laser-Electron Scattering," *IEEE Journal of Quantum Electronics* **33**:1935–1941 (1997).
68. C. I. Moore, J. P. Knauer, and D. D. Meyerhofer, "Observation of the Transition from Thomson to Compton-Scattering in Multiphoton Interactions with Low-Energy Electrons," *Physical Review Letters* **74**:2439–2442 (1995).
69. E. A. Startsev and C. J. McKinstrie, "Multiple Scale Derivation of the Relativistic Ponderomotive Force," *Physical Review E* **55**:7527–7535 (1997).
70. B. Quesnel and P. Mora, "Theory and Simulation of the Interaction of Ultraintense Laser Pulses with Electrons in Vacuum," *Physical Review E* **58**:3719–3732 (1998).
71. A. Maltsev and T. Ditmire, "Above Threshold Ionization in Tightly Focused, Strongly Relativistic Laser Fields," *Physical Review Letters* **90**:053002 (2003).
72. J. D. Jackson, *Classical Electrodynamics*, Wiley, New York, 1975.
73. L. S. Brown and T. W. B. Kibble, "Interaction of Intense Laser Beams with Electrons," *Physical Review* **133**:A705–A719 (1964).
74. E. S. Sarachik and G. T. Schappert, "Classical Theory of the Scattering of Intense Laser Radiation by Free Electrons," *Physical Review D* **1**:2738–2753 (1970).
75. S. Y. Chen, A. Maksimchuk, and D. Umstadter, "Experimental Observation of Relativistic Nonlinear Thomson Scattering," *Nature* **396**:653–655 (1998).
76. E. Esarey, S. K. Ride, and P. Sprangle, "Nonlinear Thomson Scattering of Intense Laser-Pulses from Beams and Plasmas," *Physical Review E* **48**:3003–3021 (1993).
77. C. I. Castilloherrera and T. W. Johnston, "Incoherent Harmonic Emission from Strong Electromagnetic-Waves in Plasmas," *IEEE Transactions on Plasma Science* **21**:125–135 (1993).
78. R. W. Schoenlein, W. P. Leemans, A. H. Chin, P. Volfbeyn, T. E. Glover, P. Balling, M. Zolotarev, K. J. Kim, S. Chattopadhyay, and C. V. Shank, "Femtosecond X-Ray Pulses at 0.4 Angstrom Generated by 90 Degrees Thomson Scattering: A Tool for Probing the Structural Dynamics of Materials," *Science* **274**:236–238 (1996).
79. C. Bula, K. T. McDonald, E. J. Prebys, C. Bamber, S. Boege, T. Kotseroglou, A. C. Melissinos, et al., "Observation of Nonlinear Effects in Compton Scattering," *Physical Review Letters* **76**:3116–3119 (1996).
80. G. Mainfray and C. Manus, "Multiphoton Ionization of Atoms," *Reports on Progress in Physics* **54**:1333–1372 (1991).
81. M. D. Perry and O. L. Landen, "Resonantly Enhanced Multiphoton Ionization of Krypton and Xenon with Intense Ultraviolet-Laser Radiation," *Physical Review A* **38**:2815–2829 (1988).
82. M. D. Perry, A. Szoke, O. L. Landen, and E. M. Campbell, "Nonresonant Multiphoton Ionization of Noble-Gases—Theory and Experiment," *Physical Review Letters* **60**:1270–1273 (1988).
83. M. D. Perry, O. L. Landen, A. Szoke, and E. M. Campbell, "Multiphoton Ionization of the Noble-Gases by an Intense 10(14)-W/Cm² Dye-Laser," *Physical Review A* **37**:747–760 (1988).
84. T. S. Luk, H. Pummer, K. Boyer, M. Shahidi, H. Egger, and C. K. Rhodes, "Anomalous Collision-Free Multiple Ionization of Atoms with Intense Picosecond Ultraviolet-Radiation," *Physical Review Letters* **51**:110–113 (1983).
85. A. Lhuillier, L. A. Lompre, G. Mainfray, and C. Manus, "Multiply Charged Ions Induced by Multiphoton Absorption in Rare-Gases at 0.53 μm ," *Physical Review A* **27**:2503–2512 (1983).
86. A. Lhuillier, L. A. Lompre, G. Mainfray, and C. Manus, "Multiply Charged Ions Formed by Multi-Photon Absorption Processes in the Continuum," *Physical Review Letters* **48**:1814–1817 (1982).

87. G. Petite, F. Fabre, P. Agostini, M. Crance, and M. Aymar, "Nonresonant Multiphoton Ionization of Cesium in Strong Fields—Angular-Distributions and above-Threshold Ionization," *Physical Review A* **29**:2677–2689 (1984).
88. P. Lambropoulos, "Topics on Multiphoton Processes in Atoms," *Advances in Atomic and Molecular Physics* **12**:87–164 (1976).
89. L. A. Lompre, G. Mainfray, C. Manus, and J. Thebault, "Multiphoton Ionization of Rare-Gases by a Tunable-Wavelength 30-Psec Laser-Pulse at 1.06 μm ," *Physical Review A* **15**:1604–1612 (1977).
90. R. R. Freeman, P. H. Bucksbaum, H. Milchberg, S. Darack, D. Schumacher, and M. E. Geusic, "Above-Threshold Ionization with Subpicosecond Laser-Pulses," *Physical Review Letters* **59**:1092–1095 (1987).
91. F. Faisal, *Theory of Multiphoton Processes*, Plenum, New York, 1987.
92. H. R. Reiss, "Effect of an Intense Electromagnetic-Field on a Weakly Bound System," *Physical Review A* **22**:1786–1813 (1980).
93. L. D. Landau and E. M. Lifshitz, *Quantum Mechanics*, Pergamon, New York, 1965, sec 73.
94. A. M. Perelomov, V. S. Popov, and M. V. Terent'ev, "Ionization of Atoms in an Alternating Electric Field," *Soviet Physics JETP* **23**:924–934 (1966).
95. M. V. Ammosov, N. B. Delone, and V. P. Krainov, "Tunnel Ionization of Complex Atoms and Atomic Ions in a Varying Electromagnetic-Field," *Zhurnal Eksperimentalnoi i Teoreticheskoi Fiziki* **91**:2008–2013 (1986).
96. S. Augst, D. D. Meyerhofer, D. Strickland, and S. L. Chin, "Laser Ionization of Noble-Gases by Coulomb-Barrier Suppression," *Journal of the Optical Society of America B-Optical Physics* **8**:858–867 (1991).
97. U. Mohideen, M. H. Sher, H. W. K. Tom, G. D. Aumiller, O. R. Wood, R. R. Freeman, J. Bokor, and P. H. Bucksbaum, "High-Intensity above-Threshold Ionization of He," *Physical Review Letters* **71**:509–512 (1993).
98. P. B. Corkum, N. H. Burnett, and F. Brunel, "Above-Threshold Ionization in the Long-Wavelength Limit," *Physical Review Letters* **62**:1259–1262 (1989).
99. J. Javanainen, J. H. Eberly, and Q. C. Su, "Numerical Simulations of Multiphoton Ionization and above-Threshold Electron-Spectra," *Physical Review A* **38**:3430–3446 (1988).
100. P. Kruit, J. Kimman, H. G. Muller, and M. J. Vanderwiel, "Electron-Spectra from Multiphoton Ionization of Xenon at 1064, 532, and 355 Nm," *Physical Review A* **28**:248–255 (1983).
101. J. L. Chaloupka, J. Rudati, R. Lafon, P. Agostini, K. C. Kulander, and L. F. DiMauro, "Observation of a Transition in the Dynamics of Strong-Field Double Ionization," *Physical Review Letters* **90**:033002 (2003).
102. B. R. Yang, K. J. Schafer, B. Walker, K. C. Kulander, P. Agostini, and L. F. DiMauro, "Intensity-Dependent Scattering Rings in High-Order above-Threshold Ionization," *Physical Review Letters* **71**:3770–3773 (1993).
103. R. R. Freeman, T. J. McIlrath, P. H. Bucksbaum, and M. Bashkansky, "Ponderomotive Effects on Angular-Distributions of Photoelectrons," *Physical Review Letters* **57**:3156–3159 (1986).
104. G. N. Gibson, R. R. Freeman, and T. J. McIlrath, "Verification of the Dominant Role of Resonant Enhancement in Short-Pulse Multiphoton Ionization," *Physical Review Letters* **69**:1904–1907 (1992).
105. N. H. Burnett and P. B. Corkum, "Cold-Plasma Production for Recombination Extreme-Ultraviolet Lasers by Optical-Field-Induced Ionization," *Journal of the Optical Society of America B-Optical Physics* **6**:1195–1199 (1989).
106. T. Ditmire, "Simulations of Heating and Electron Energy Distributions in Optical Field Ionized Plasmas," *Physical Review E* **54**:6735–6740 (1996).
107. X. Liu, H. Rottke, E. Eremina, W. Sandner, E. Goulielmakis, K. O. Keeffe, M. Lezius, et al., "Nonsequential Double Ionization at the Single-Optical-Cycle Limit," *Physical Review Letters* **93**:263001 (2004).
108. B. Feuerstein, R. Moshhammer, D. Fischer, A. Dorn, C. D. Schroter, J. Deipenwisch, J. R. C. Lopez-Urrutia, et al., "Separation of Recollision Mechanisms in Nonsequential Strong Field Double Ionization of Ar: The Role of Excitation Tunneling," *Physical Review Letters* **87**:043003 (2001).
109. M. Lein, E. K. U. Gross, and V. Engel, "Intense-Field Double Ionization of Helium:Identifying the Mechanism," *Physical Review Letters* **85**:4707–4710 (2000).
110. R. Kopold, W. Becker, H. Rottke, and W. Sandner, "Routes to Nonsequential Double Ionization," *Physical Review Letters* **85**:3781–3784 (2000).
111. H. W. van der Hart, and K. Burnett, "Recollision Model for Double Ionization of Atoms in Strong Laser Fields," *Physical Review A* **62**:013407 (2000).

112. R. Moshhammer, B. Feuerstein, W. Schmitt, A. Dorn, C. D. Schroter, J. Ullrich, H. Rottke, et al., "Momentum Distributions of Nen^+ Ions Created by an Intense Ultrashort Laser Pulse," *Physical Review Letters* **84**:447–450 (2000).
113. B. Sheehy, R. Lafon, M. Widmer, B. Walker, L. F. DiMauro, P. A. Agostini, and K. C. Kulander, "Single- and Multiple-Electron Dynamics in the Strong-Field Tunneling Limit," *Physical Review A* **58**:3942–3952 (1998).
114. B. Walker, B. Sheehy, K. C. Kulander, and L. F. DiMauro, "Elastic Rescattering in the Strong Field Tunneling Limit," *Physical Review Letters* **77**:5031–5034 (1996).
115. B. Walker, B. Sheehy, L. F. DiMauro, P. Agostini, K. J. Schafer, and K. C. Kulander, "Precision-Measurement of Strong-Field Double-Ionization of Helium," *Physical Review Letters* **73**:1227–1230 (1994).
116. W. Becker, A. Lohr, and M. Kleber, "Effects of Rescattering on above-Threshold Ionization," *Journal of Physics B-Atomic Molecular and Optical Physics* **27**:L325–L332 (1994).
117. D. N. Fittinghoff, P. R. Bolton, B. Chang, and K. C. Kulander, "Observation of Nonsequential Double Ionization of Helium with Optical Tunneling," *Physical Review Letters* **69**:2642–2645 (1992).
118. R. Dorn, H. Brauning, J. M. Feagin, V. Mergel, O. Jagutzki, L. Spielberger, T. Vogt, et al., "Photo-Double-Ionization of He: Fully Differential and Absolute Electronic and Ionic Momentum Distributions," *Physical Review A* **57**:1074–1090 (1998).
119. R. Wehlitz, F. Heiser, O. Hemmers, B. Langer, A. Menzel, and U. Becker, "Electron-Energy and Electron-Angular Distributions in the Double Photoionization of Helium," *Physical Review Letters* **67**:3764–3767 (1991).
120. U. Eichmann, M. Dorr, H. Maeda, W. Becker, and W. Sandner, "Collective Multielectron Tunneling Ionization in Strong Fields," *Physical Review Letters* **84**:3550–3553 (2000).
121. D. N. Fittinghoff, P. R. Bolton, B. Chang, and K. C. Kulander, "Polarization Dependence of Tunneling Ionization of Helium and Neon by 120-fs Pulses at 614 nm," *Physical Review A* **49**:2174–2177 (1994).
122. N. Milosevic, V. P. Krainov, and T. Brabec, "Relativistic Theory of Tunnel Ionization," *Journal of Physics B-Atomic Molecular and Optical Physics* **35**:3515–3529 (2002).
123. S. X. Hu and A. F. Starace, "GeV Electrons from Ultraintense Laser Interaction with Highly Charged Ions," *Physical Review Letters* **88**:245003 (2002).
124. S. Palaniyappan, I. Ghebregziabher, A. DiChiara, J. MacDonald, and B. C. Walker, "Emergence from non-relativistic strong-field rescattering to ultrastrong-field laser-atom physics: A semiclassical analysis," *Physical Review A* **74**:033403 (2006).
125. I. Ghebregziabher, S. Palaniyappan, J. MacDonald, and B. C. Walker, "Impact of the Laser Magnetic Field on Recombination and Bremsstrahlung Radiation from Atomic Ionization Rescattering in Ultraintense Fields," *Physical Review A* **73**:033419 (2006).
126. S. Palaniyappan, A. DiChiara, E. Chowdhury, A. Falkowski, G. Ongadi, E. L. Huskins, and B. C. Walker, "Ultrastrong Field Ionization of Nen^+ ($n \leq 8$): Rescattering and the Role of the Magnetic Field," *Physical Review Letters* **94**:243003 (2005).
127. M. Pont and M. Gavrila, "Stabilization of Atomic-Hydrogen in Superintense, High-Frequency Laser Fields of Circular-Polarization," *Physical Review Letters* **65**:2362–2365 (1990).
128. Q. Su, J. H. Eberly, and J. Javanainen, "Dynamics of Atomic Ionization Suppression and Electron Localization in an Intense High-Frequency Radiation-Field," *Physical Review Letters* **64**:862–865 (1990).
129. K. Burnett, P. L. Knight, B. R. M. Piraux, and V. C. Reed, "Suppression of Ionization in Strong Laser Fields," *Physical Review Letters* **66**:301–304 (1991).
130. M. V. Fedorov and O. V. Tikhonova, "Strong-Field Short-Pulse Photoionization of Rydberg Atoms: Interference Stabilization and Distribution of the Photoelectron Density in Space and Time," *Physical Review A* **58**:1322–1334 (1998).
131. O. V. Tikhonova, E. A. Volkova, A. M. Popov, and M. V. Fedorov, "Interference Stabilization of Rydberg Atoms: Analytical Investigation and Numerical Simulations," *Laser Physics* **8**:85–92 (1998).
132. M. P. DeBoer, J. H. Hoogenraad, R. B. Vrijen, L. D. Noordam, and H. G. Muller, "Indications of High-Intensity Adiabatic Stabilization in Neon," *Physical Review Letters* **71**:3263–3266 (1993).
133. L. N. Gaier and C. H. Keitel, "Relativistic Classical Monte Carlo Simulations of Stabilization of Hydrogenlike Ions in Intense Laser Pulses," *Physical Review A* **65**:023406 (2002).
134. A. Talebpour, J. Yang, and S. L. Chin, "Semi-Empirical Model for the Rate of Tunnel Ionization of $N-2$ and $O-2$ Molecule in an Intense Ti: Sapphire Laser Pulse," *Optics Communications* **163**:29–32 (1999).

135. J. H. Posthumus, J. Plumridge, L. J. Frasinski, K. Codling, A. J. Langley, and P. F. Taday, "Double-Pulse Measurements of Laser-Induced Alignment of Molecules," *Journal of Physics B-Atomic Molecular and Optical Physics* **31**:L985–L993 (1998).
136. J. H. Posthumus, J. Plumridge, M. K. Thomas, K. Codling, L. J. Frasinski, A. J. Langley, and P. F. Taday, "Dynamic and Geometric Laser-Induced Alignment of Molecules in Intense Laser Fields," *Journal of Physics B-Atomic Molecular and Optical Physics* **31**:L553–L562 (1998).
137. P. Hering and C. Cornaggia, "Coulomb Explosion of N-2 and CO2 Using Linearly and Circularly Polarized Femtosecond Laser Pulses," *Physical Review A* **59**:2836–2843 (1999).
138. C. Cornaggia, "Small Polyatomic Molecules in Intense Laser Fields", in *Molecules and Clusters in Intense Laser Fields*, J. Posthumus, (ed.), Cambridge University Press, Cambridge, 2001, pP. 84–113.
139. P. H. Bucksbaum, A. Zavriyev, H. G. Muller, and D. W. Schumacher, "Softening of the H₂⁺ Molecular-Bond in Intense Laser Fields," *Physical Review Letters* **64**:1883–1886 (1990).
140. B. Sheehy and L. F. DiMauro, "Atomic and Molecular Dynamics in Intense Optical Fields," *Annual Review of Physical Chemistry* **47**:463–494 (1996).
141. J. H. Shirley, "Solution of the Schrödinger Equation with a Hamiltonian Periodic in Time, *Physical Review* **138**:B979–B987 (1965).
142. A. D. Bandrauk and M. L. Sink, "Photo-Dissociation in Intense Laser Fields—Predissociation Analogy," *Journal of Chemical Physics* **74**:1110–1117 (1981).
143. A. Giustisuzor, X. He, O. Atabek, and F. H. Mies, "Above-Threshold Dissociation of H₂⁺ in Intense Laser Fields," *Physical Review Letters* **64**:515–518 (1990).
144. G. Jolicard and O. Atabek, "Above-Threshold-Dissociation Dynamics of H₂⁽⁺⁾ with Short Intense Laser-Pulses," *Physical Review A* **46**:5845–5855 (1992).
145. J. H. Posthumus and J. F. McCann, "Diatomic Molecules in Intense Laser Fields," in *Molecules and Clusters in Intense Laser Fields*, J. Posthumus, (ed.), Cambridge University Press, Cambridge, 2001, pP. 27–83.
146. M. R. Thompson, M. K. Thomas, P. F. Taday, J. H. Posthumus, A. J. Langley, L. J. Frasinski, and K. Codling, "One and Two-Colour Studies of the Dissociative Ionization and Coulomb Explosion of H₂ with intense Ti: Sapphire Laser Pulses," *Journal of Physics B-Atomic Molecular and Optical Physics* **30**:5755–5772 (1997).
147. C. Cornaggia, J. Lavancier, D. Normand, J. Morellec, P. Agostini, J. P. Chambaret, and A. Antonetti, "Multielectron Dissociative Ionization of Diatomic-Molecules in an Intense Femtosecond Laser Field," *Physical Review A* **44**:4499–4505 (1991).
148. C. Cornaggia and P. Hering, "Nonsequential Double Ionization of Small Molecules Induced by a Femtosecond Laser Field," *Physical Review A* **62**:023403 (2000).
149. A. Talebpour, S. Larochelle, and S. L. Chin, "Non-Sequential and Sequential Double Ionization of No in an Intense Femtosecond Ti:Sapphire Laser Pulse," *Journal of Physics B-Atomic Molecular and Optical Physics* **30**:L245–L250 (1997).
150. C. Guo, M. Li, J. P. Nibarger, and G. N. Gibson, "Single and Double Ionization of Diatomic Molecules in Strong Laser Fields," *Physical Review A* **58**:R4271–R4274 (1998).
151. S. Chelkowski and A. D. Bandrauk, "Two-Step Coulomb Explosions of Diatoms in Intense Laser Fields," *Journal of Physics B-Atomic Molecular and Optical Physics* **28**:L723–L731 (1995).
152. H. T. Yu, T. Zuo, and A. D. Bandrauk, "Intense Field Ionization of Molecules with Ultra-Short Laser Pulses-Enhanced Ionization and Barrier-Suppression Effects," *Journal of Physics B-Atomic Molecular and Optical Physics* **31**:1533–1551 (1998).
153. M. Schmidt, D. Normand, and C. Cornaggia, "Laser-Induced Trapping of Chlorine Molecules with Picosecond and Femtosecond Pulses," *Physical Review A* **50**:5037–5045 (1994).
154. J. H. Posthumus, L. J. Frasinski, A. J. Giles, and K. Codling, "Dissociative Ionization of Molecules in Intense Laser Fields—a Method of Predicting Ion Kinetic Energies and Appearance Intensities," *Journal of Physics B-Atomic Molecular and Optical Physics* **28**:L349–L353 (1995).
155. E. Constant, H. Stapelfeldt, and P. B. Corkum, "Observation of Enhanced Ionization of Molecular Ions in Intense Laser Fields," *Physical Review Letters* **76**:4140–4143 (1996).
156. J. F. McCann and J. H. Posthumus, "Molecular Dynamics in Intense Laser Fields," *Philosophical Transactions of the Royal Society of London Series a-Mathematical Physical and Engineering Sciences* **357**:1309–1329 (1999).

157. T. Zuo and A. D. Bandrauk, "Charge-Resonance-Enhanced Ionization of Diatomic Molecular-Ions by Intense Lasers," *Physical Review A* **52**:R2511–R2514 (1995).
158. R. J. Levis and M. J. DeWitt, "Photoexcitation, Ionization, and Dissociation of Molecules Using Intense Near-Infrared Radiation of Femtosecond Duration," *Journal of Physical Chemistry A* **103**:6493–6507 (1999).
159. M. J. DeWitt and R. J. Levis, "Concerning the Ionization of Large Polyatomic Molecules with Intense Ultrafast Lasers," *Journal of Chemical Physics* **110**:11368–11375 (1999).
160. J. J. Macklin, J. D. Kmetec, and C. L. Gordon, "High-Order Harmonic-Generation Using Intense Femtosecond Pulses," *Physical Review Letters* **70**:766–769 (1993).
161. M. D. Perry and J. K. Crane, "High-Order Harmonic Emission from Mixed Fields," *Physical Review A* **48**:R4051–R4054 (1993).
162. M. Y. Ivanov and P. B. Corkum, "Generation of High-Order Harmonics from Inertially Confined Molecular-Ions," *Physical Review A* **48**:580–590 (1993).
163. T. D. Donnelly, T. Ditmire, K. Neuman, M. D. Perry and R. W. Falcone, "High-Order Harmonic Generation in Atom Clusters," *Physical Review Letters* **76**:2472–2475 (1996).
164. K. S. Budil, P. Salieres, A. Lhuillier, T. Ditmire, and M. D. Perry, "Influence of Ellipticity on Harmonic-Generation," *Physical Review A* **48**:R3437–R3440 (1993).
165. P. Antoine, A. Lhuillier, M. Lewenstein, P. Salieres, and B. Carre, "Theory of High-Order Harmonic Generation by an Elliptically Polarized Laser Field," *Physical Review A* **53**:1725–1745 (1996).
166. T. Ditmire, E. T. Gumbrell, R. A. Smith, J. W. G. Tisch, D. D. Meyerhofer, and M. H. R. Hutchinson, "Spatial Coherence Measurement of Soft X-Ray Radiation Produced by High Order Harmonic Generation," *Physical Review Letters* **77**:4756–4759 (1996).
167. A. L'Huillier, L. -A. Lompre, G. Mainfray, and C. Manus, "High-Order Harmonic Generation in Rare Gases," in *Atoms in Intense Laser Fields*, M. Gavrila, (ed.), Academic Press, San Diego, 1992, pP. 139–206.
168. C. Spielmann, N. H. Burnett, S. Sartania, R. Koppitsch, M. Schnurer, C. Kan, M. Lenzner, P. Wobrowschek, and F. Krausz, "Generation of Coherent X-Rays in the Water Window Using 5-Femtosecond Laser Pulses," *Science* **278**:661–664 (1997).
169. M. Lewenstein, P. Balcou, M. Y. Ivanov, A. Lhuillier, and P. B. Corkum, "Theory of High-Harmonic Generation by Low-Frequency Laser Fields," *Physical Review A* **49**:2117–2132 (1994).
170. J. L. Krause, K. J. Schafer, and K. C. Kulander, "High-Order Harmonic-Generation from Atoms and Ions in the High-Intensity Regime," *Physical Review Letters* **68**:3535–3538 (1992).
171. M. Lewenstein, P. Salieres, and A. Lhuillier, "Phase of the Atomic Polarization in High-Order Harmonic-Generation," *Physical Review A* **52**:4747–4754 (1995).
172. P. Salieres, A. Lhuillier, and M. Lewenstein, "Coherence Control of High-Order Harmonics," *Physical Review Letters* **74**:3776–3779 (1995).
173. P. Balcou, A. S. Dederichs, M. B. Gaarde, and A. L'Huillier, "Quantum-Path Analysis and Phase Matching of High-Order Harmonic Generation and High-Order Frequency Mixing Processes in Strong Laser Fields," *Journal of Physics B-Atomic Molecular and Optical Physics* **32**:2973–2989 (1999).
174. M. B. Gaarde, F. Salin, E. Constant, P. Balcou, K. J. Schafer, K. C. Kulander, and A. L'Huillier, "Spatiotemporal Separation of High Harmonic Radiation into Two Quantum Path Components," *Physical Review A* **59**:1367–1373 (1999).
175. R. W. Boyd, *Nonlinear Optics*, Academic Press, Boston, 1992.
176. T. Ditmire, J. K. Crane, H. Nguyen, L. B. Dasilva, and M. D. Perry, "Energy-Yield and Conversion-Efficiency Measurements of High-Order Harmonic Radiation," *Physical Review A* **51**:R902–R905 (1995).
177. E. Takahashi, Y. Nabekawa, M. Nurhuda, and K. Midorikawa, "Generation of High-Energy High-Order Harmonics by Use of a Long Interaction Medium," *Journal of the Optical Society of America B-Optical Physics* **20**:158–165 (2003).
178. M. Drescher, M. Hentschel, R. Kienberger, M. Uiberacker, V. Yakovlev, A. Scrinizi, T. Westerwalbesloh, U. Kleineberg, U. Heinzmann, and F. Krausz, "Time-Resolved Atomic Inner-Shell Spectroscopy," *Nature* **419**:803–807 (2002).
179. A. Baltuska, T. Udem, M. Uiberacker, M. Hentschel, E. Goulielmakis, C. Gohle, R. Holzwarth, et al., "Attosecond Control of Electronic Processes by Intense Light Fields," *Nature* **421**:611–615 (2003).

180. P. Antoine, A. Lhuillier, and M. Lewenstein, "Attosecond Pulse Trains Using High-Order Harmonics," *Physical Review Letters* **77**:1234–1237 (1996).
181. T. Ditmire, T. Donnelly, A. M. Rubenchik, R. W. Falcone, and M. D. Perry, "Interaction of Intense Laser Pulses with Atomic Clusters," *Physical Review A* **53**:3379–3402 (1996).
182. A. McPherson, B. D. Thompson, A. B. Borisov, K. Boye, and C. K. Rhodes, "Multiphoton-Induced X-Ray Emission at 4–5 keV from Xe Atoms with Multiple Core Vacancies," *Nature* **370**:631–634 (1994).
183. V. P. Krainov and M. B. Smirnov, "Cluster Beams in the Super-Intense Femtosecond Laser Pulse," *Physics Reports-Review Section of Physics Letters* **370**:237–331 (2002).
184. W. L. Kruer, *The Physics of Laser Plasma Interactions*, Addison-Wesley, Redwood City, 1988.
185. C. Siedschlag and J. M. Rost, "Enhanced Ionization in Small Rare-Gas Clusters," *Physical Review A* **67**:013404 (2003).
186. C. RosePetrucci, K. J. Schafer, K. R. Wilson, and C. P. J. Barty, "Ultrafast Electron Dynamics and Inner-Shell Ionization in Laser Driven Clusters," *Physical Review A* **55**:1182–1190 (1997).
187. W. Lotz, "An Empirical Formula for the Electron-Impact Ionization Cross-Section," *Z. fur Physik* **206**:205–211 (1967).
188. E. M. Snyder, S. A. Buzza, and A. W. Castleman, "Intense Field-Matter Interactions: Multiple Ionization of Clusters," *Physical Review Letters* **77**:3347–3350 (1996).
189. T. Ditmire, J. W. G. Tisch, E. Springate, M. B. Mason, N. Hay, R. A. Smith, J. Marangos, and M. H. R. Hutchinson, "High-Energy Ions Produced in Explosions of Superheated Atomic Clusters," *Nature* **386**:54–56 (1997).
190. M. Lezius, S. Dobosz, D. Normand, and M. Schmidt, "Explosion Dynamics of Rare Gas Clusters in Strong Laser Fields," *Physical Review Letters* **80**:261–264 (1998).
191. B. N. Breizman, A. V. Arefiev, and M. V. Fomyts'kiy, "Nonlinear Physics of Laser-Irradiated Microclusters," *Physics of Plasmas* **12**:056706 (2005).
192. K. W. Madison, P. K. Patel, M. Allen, D. Price, R. Fitzpatrick, and T. Ditmire, "Role of Laser-Pulse Duration in the Neutron Yield of Deuterium Cluster Targets," *Physical Review A* **70**:053201 (2004).
193. K. W. Madison, P. K. Patel, D. Price, A. Edens, M. Allen, T. E. Cowan, J. Zweiback, and T. Ditmire, "Fusion Neutron and Ion Emission from Deuterium and Deuterated Methane Cluster Plasmas," *Physics of Plasmas* **11**:270–277 (2004).
194. L. D. Landau and E. M. Lifshitz, *Electrodynamics of Continuous Media*, Pergamon Press, Oxford, 1984, pp. 272–273.
195. V. P. Silin, "Nonlinear High-Frequency Plasma Conductivity," *Soviet Physics, JETP* **20**:1510 (1965).
196. T. Ditmire, R. A. Smith, J. W. G. Tisch, and M. H. R. Hutchinson, "High Intensity Laser Absorption by Gases of Atomic Clusters," *Physical Review Letters* **78**:3121–3124 (1997).
197. T. Ditmire, T. Donnelly, R. W. Falcone, and M. D. Perry, "Strong X-Ray Emission from High-Temperature Plasmas Produced by Intense Irradiation of Clusters," *Physical Review Letters* **75**:3122–3125 (1995).
198. T. Ditmire, R. A. Smith, R. S. Marjoribanks, G. Kulcsa, and M. H. R. Hutchinson, "X-Ray Yields from Xe Clusters Heated by Short Pulse High Intensity Lasers," *Applied Physics Letters* **71**:166–168 (1997).
199. T. Ditmire, J. Zweiback, V. P. Yanovsky, T. E. Cowan, G. Hays, and K. B. Wharton, "Nuclear Fusion from Explosions of Femtosecond Laser-Heated Deuterium Clusters," *Nature* **398**:489–492 (1999).
200. G. C. Junkel-Vives, J. Abdallah, F. Blasco, F. Dorchies, T. Caillaud, C. Bonte, C. Stenz, et al., "Evidence of Supercritical Density in 45-fs-Laser-Irradiated Ar-Cluster Plasmas," *Physical Review A* **66**:033204 (2002).
201. U. Saalmann and J. M. Rost, "Ionization of Clusters in Intense Laser Pulses Through Collective Electron Dynamics," *Physical Review Letters* **91**:223401 (2003).
202. J. Zweiback, T. Ditmire, and M. D. Perry, "Femtosecond Time-Resolved Studies of the Dynamics of Noble-Gas Cluster Explosions," *Physical Review A* **59**:R3166–R3169 (1999).
203. H. M. Milchberg, S. J. McNaught, and E. Parra, "Plasma Hydrodynamics of the Intense Laser-Cluster Interaction," *Physical Review E* **64**:056402 (2001).
204. A. B. Langdon, "Non-Linear Inverse Bremsstrahlung and Heated-Electron Distributions," *Physical Review Letters* **44**:575–579 (1980).
205. G. J. Pert, "Inverse Bremsstrahlung Absorption In Large Radiation Fields During Binary Collisions—Classical Theory," *Journal of Physics A: General Physics* **5**:506–520 (1972).

206. R. D. Jones and K. Lee, "Kinetic-Theory, Transport, and Hydrodynamics of a High-Z Plasma in the Presence of an Intense Laser Field," *Physics of Fluids* **25**:2307–2323 (1982).
207. C. S. Liu and M. N. Rosenbluth, "Parametric Decay of Electromagnetic-Waves into 2 Plasmons and Its Consequences," *Physics of Fluids* **19**:967–971 (1976).
208. D. W. Forslund, J. M. Kindel, and E. L. Lindman, "Theory of Stimulated Scattering Processes in Laser-Irradiated Plasmas," *Physics of Fluids* **18**:1002–1016 (1975).
209. C. B. Darrow, C. Coverdale, M. D. Perry, W. B. Mori, C. Clayton, K. Marsh, and C. Joshi, "Strongly Coupled Stimulated Raman Backscatter from Subpicosecond Laser-Plasma Interactions," *Physical Review Letters* **69**:442–445 (1992).
210. W. B. Mori, C. D. Decker, D. E. Hinkel, and T. Katsouleas, "Raman Forward Scattering of Short-Pulse High-Intensity Lasers," *Physical Review Letters* **72**:1482–1485 (1994).
211. S. Guerin, G. Laval, P. Mora, J. C. Adam, A. Heron, and A. Bendib, "Modulational and Raman Instabilities in the Relativistic Regime," *Physics of Plasmas* **2**:2807–2814 (1995).
212. V. Malka, S. Fritzler, E. Lefebvre, M. M. Alonard, F. Burgy, J. P. Chambaret, J. F. Chemin, et al., "Electron Acceleration by a Wake Field Forced by an Intense Ultrashort Laser Pulse," *Science* **298**:1596–1600 (2002).
213. A. Modena, Z. Najmudin, A. E. Dangor, C. E. Clayton, K. A. Marsh, C. Joshi, V. Malka, et al., "Electron Acceleration from the Breaking of Relativistic Plasma-Waves," *Nature* **377**:606–608 (1995).
214. K. Nakajima, D. Fisher, T. Kawakubo, H. Nakanishi, A. Ogata, Y. Kato, Y. Kitagawa, et al., "Observation of Ultrahigh Gradient Electron Acceleration by a Self-Modulated Intense Short Laser-Pulse," *Physical Review Letters* **74**:4428–4431 (1995).
215. M. Everett, A. Lal, D. Gordon, C. E. Clayton, K. A. Marsh, and C. Joshi, "Trapped Electron Acceleration by a Laser-Driven Relativistic Plasma-Wave," *Nature* **368**:527–529 (1994).
216. J. Krall, A. Ting, E. Esarey, and P. Sprangle, "Enhanced Acceleration in a Self-Modulated-Laser Wake-Field Accelerator," *Physical Review E* **48**:2157–2161 (1993).
217. P. Sprangle, E. Esarey, A. Ting, and G. Joyce, "Laser Wakefield Acceleration and Relativistic Optical Guiding," *Applied Physics Letters* **53**:2146–2148 (1988).
218. S. P. D. Mangles, C. D. Murphy, Z. Najmudin, A. G. R. Thomas, J. L. Collier, A. E. Dangor, E. J. Divall, et al., "Monoenergetic Beams of Relativistic Electrons from Intense Laser-Plasma Interactions," *Nature* **431**:535–538 (2004).
219. C. G. R. Geddes, C. Toth, J. van Tilborg, E. Esarey, C. B. Schroeder, D. Bruhwiler, C. Nieter, J. Cary, and W. P. Leemans, "High-Quality Electron Beams from a Laser Wakefield Accelerator Using Plasma-Channel Guiding," *Nature* **431**:538–541 (2004).
220. J. Faure, Y. Glinec, A. Pukhov, S. Kiselev, S. Gordienko, E. Lefebvre, J. P. Rousseau, F. Burgy, and V. Malka, "A Laser-Plasma Accelerator Producing Monoenergetic Electron Beams," *Nature* **431**:541–544 (2004).
221. A. Pukhov and J. Meyer-ter-Vehn, "Laser Wake Field Acceleration: The Highly Non-Linear Broken-Wave Regime," *Applied Physics B-Lasers and Optics* **74**:355–361 (2002).
222. P. Sprangle, E. Esarey, J. Krall, and G. Joyce, "Propagation and Guiding of Intense Laser-Pulses in Plasmas," *Physical Review Letters* **69**:2200–2203 (1992).
223. K. Krushelnick, E. L. Clark, Z. Najmudin, M. Salvati, M. I. K. Santala, M. Tatarakis, A. E. Dangor, et al., "Multi-MeV Ion Production from High-Intensity Laser Interactions with Underdense Plasmas," *Physical Review Letters* **83**:737–740 (1999).
224. A. Pukhov, Z. M. Sheng, and J. Meyer-ter-Vehn, "Particle Acceleration in Relativistic Laser Channels," *Physics of Plasmas* **6**:2847–2854 (1999).
225. S. C. Rae, "Ionization-Induced Defocusing of Intense Laser-Pulses in High-Pressure Gases," *Optics Communications* **97**:25–28 (1993).
226. P. Monot, T. Auguste, L. A. Lompre, G. Mainfray, and C. Manus, "Focusing Limits of a Terawatt Laser in an Underdense Plasma," *Journal of the Optical Society of America B-Optical Physics* **9**:1579–1584 (1992).
227. M. Borghesi, A. J. MacKinnon, L. Barringer, R. Gaillard, L. A. Gizzi, C. Meyer, O. Willi, A. Pukhov, and J. Meyer-ter-Vehn, "Relativistic Channeling of a Picosecond Laser Pulse in a Near-Critical Preformed Plasma," *Physical Review Letters* **78**:879–882 (1997).
228. P. Monot, T. Auguste, P. Gibbon, F. Jakober, and G. Mainfray, "Collimation of an Intense Laser Beam by a Weakly Relativistic Plasma," *Physical Review E* **52**:R5780–R5783 (1995).

229. P. Monot, T. Auguste, P. Gibbon, F. Jakober, G. Mainfray, A. Dulieu, M. Louisjacquet, G. Malka, and J. L. Miquel, "Experimental Demonstration of Relativistic Self-Channelling of a Multiterawatt Laser-Pulse in an Underdense Plasma," *Physical Review Letters* **74**:2953–2956 (1995).
230. A. B. Borisov, A. V. Borovskiy, O. B. Shiryayev, V. V. Korobkin, A. M. Prokhorov, J. C. Solem, T. S. Luk, K. Boyer, and C. K. Rhodes, "Relativistic and Charge-Displacement Self-Channelling of Intense Ultrashort Laser-Pulses in Plasmas," *Physical Review A* **45**:5830–5845 (1992).
231. I. Watts, M. Zepf, E. L. Clark, M. Tatarakis, K. Krushelnick, A. E. Dangor, R. Allott, R. J. Clarke, D. Neely, and P. A. Norreys, "Measurements of Relativistic Self-Phase-Modulation in Plasma," *Physical Review E* **66**:036409 (2002).
232. U. Teubner, I. Uschmann, P. Gibbon, D. Altenbernd, E. Forster, T. Feurer, W. Theobald, et al., "Absorption and Hot Electron Production by High Intensity Femtosecond UV-Laser Pulses In Solid Targets," *Physical Review E* **54**:4167–4177 (1996).
233. D. F. Price, R. M. More, R. S. Walling, G. Guethlein, R. L. Shepherd, R. E. Stewart, and W. E. White, "Absorption of Ultrashort Laser-Pulses by Solid Targets Heated Rapidly to Temperatures 1–1000 eV," *Physical Review Letters* **75**:252–255 (1995).
234. V. L. Ginzberg, *The Properties of Electromagnetic Waves in Plasma*, Pergamon, New York, 1964.
235. D. W. Forslund, J. M. Kindel, and K. Lee, "Theory of Hot-Electron Spectra at High Laser Intensity," *Physical Review Letters* **39**:284–288 (1977).
236. F. N. Beg, A. R. Bell, A. E. Dangor, C. N. Danson, A. P. Fews, M. E. Glinsky, B. A. Hammel, P. Lee, P. A. Norreys, and M. Tatarakis, "A Study of Picosecond Laser-Solid Interactions up to $10(19) \text{ W cm}^{-2}$," *Physics of Plasmas* **4**:447–457 (1997).
237. M. K. Grimes, A. R. Rundquist, Y. S. Lee, and M. C. Downer, "Experimental Identification of 'Vacuum Heating' at Femtosecond-Laser-Irradiated Metal Surfaces," *Physical Review Letters* **82**:4010–4013 (1999).
238. P. Gibbon and A. R. Bell, "Collisionless Absorption in Sharp-Edged Plasmas," *Physical Review Letters* **68**:1535–1538 (1992).
239. F. Brunel, "Not-So-Resonant, Resonant Absorption," *Physical Review Letters* **59**:52–55 (1987).
240. W. Rozmus and V. T. Tikhonchuk, "Skin Effect and Interaction of Short Laser-Pulses with Dense-Plasmas," *Physical Review A* **42**:7401–7412 (1990).
241. W. L. Kruer and K. Estabrook, "X-ray Heating by Very Intense Laser-Light," *Physics of Fluids* **28**:430–432 (1985).
242. S. C. Wilks, W. L. Kruer, M. Tabak, and A. B. Langdon, "Absorption of Ultra-Intense Laser-Pulses," *Physical Review Letters* **69**:1383–1386 (1992).
243. M. Zepf, M. CastroColin, D. Chambers, S. G. Preston, J. S. Wark, J. Zhang, C. N. Danson, et al., "Measurements of the Hole Boring Velocity from Doppler Shifted Harmonic Emission from Solid Targets," *Physics of Plasmas* **3**:3242–3244 (1996).
244. X. Liu and D. Umstadter, "Competition between Ponderomotive and Thermal Forces in Short-Scale-Length Laser Plasmas," *Physical Review Letters* **69**:1935–1938 (1992).
245. L. O. Silva, M. Marti, J. R. Davies, R. A. Fonseca, C. Ren, F. S. Tsung, and W. B. Mori, "Proton Shock Acceleration in Laser-Plasma Interactions," *Physical Review Letters* **92**:015002 (2004).
246. D. vonderLinde and K. Rzaewski, "High-Order Optical Harmonic Generation from Solid Surfaces," *Applied Physics B-Lasers and Optics* **63**:499–506 (1996).
247. R. Lichters, J. MeyerterVehn, and A. Pukhov, "Short-Pulse Laser Harmonics from Oscillating Plasma Surfaces Driven at Relativistic Intensity," *Physics of Plasmas* **3**:3425–3437 (1996).
248. P. A. Norreys, M. Zepf, S. Moustazis, A. P. Fews, J. Zhang, P. Lee, M. Bakarezos, et al., "Efficient Extreme UV Harmonics Generated from Picosecond Laser Pulse Interactions with Solid Targets," *Physical Review Letters* **76**:1832–1835 (1996).
249. P. Gibbon, "Harmonic Generation by Femtosecond Laser-Solid Interaction: A Coherent 'water-window' Light Source?," *Physical Review Letters* **76**:50–53 (1996).
250. P. Kaw and J. Dawson, "Relativistic Nonlinear Propagation of Laser Beams in Cold Overdense Plasma," *Physics of Fluids* **13**:472–481 (1970).
251. J. D. Kmetec, C. L. Gordon, J. J. Macklin, B. E. Lemoff, G. S. Brown, and S. E. Harris, "MeV X-Ray Generation with a Femtosecond Laser," *Physical Review Letters* **68**:1527–1530 (1992).

252. A. Saemann, K. Eidmann, I. E. Golovkin, R. C. Mancini, E. Andersson, E. Forster, and K. Witte, "Isochoric Heating of Solid Aluminum by Ultrashort Laser Pulses Focused on a Tamped Target," *Physical Review Letters* **82**:4843–4846 (1999).
253. D. Hilscher, O. Berndt, M. Enke, U. Jahnke, P. V. Nickles, H. Ruhl, and W. Sandner, "Neutron Energy Spectra from the Laser-Induced D(d,n)He-3 Reaction," *Physical Review E* **64**:016414 (2001).
254. A. R. Bell, F. N. Beg, Z. Chang, A. E. Dangor, C. N. Danson, C. B. Edwards, A. P. Fews, et al., "Observation of Plasma-Confinement in Picosecond Laser-Plasma Interactions," *Physical Review E* **48**:2087–2093 (1993).
255. M. Tatarakis, A. Gopal, I. Watts, F. N. Beg, A. E. Dangor, K. Krushelnick, U. Wagner, et al., "Measurements of Ultrastrong Magnetic Fields During Relativistic Laser-Plasma Interactions," *Physics of Plasmas* **9**:2244–2250 (2002).
256. R. A. Snavely, M. H. Key, S. P. Hatchett, T. E. Cowan, M. Roth, T. W. Phillips, M. A. Stoyer, et al., "Intense High-Energy Proton Beams from Petawatt-Laser Irradiation of Solids," *Physical Review Letters* **85**:2945–2948 (2000).
257. S. C. Wilks, A. B. Langdon, T. E. Cowan, M. Roth, M. Singh, S. Hatchett, M. H. Key, D. Pennington, A. MacKinnon, and R. A. Snavely, "Energetic Proton Generation in Ultra-Intense Laser-Solid Interactions," *Physics of Plasmas* **8**:542–549 (2001).
258. S. P. Hatchett, C. G. Brown, T. E. Cowan, E. A. Henry, J. S. Johnson, M. H. Key, J. A. Koch, et al., "Electron, Photon, and Ion Beams from the Relativistic Interaction of Petawatt Laser Pulses with Solid Targets," *Physics of Plasmas* **7**:2076–2082 (2000).
259. M. Tabak, J. Hammer, M. E. Glinsky, W. L. Kruer, S. C. Wilks, J. Woodworth, E. M. Campbell, M. D. Perry, and R. J. Mason, "Ignition and High-Gain with Ultrapowerful Lasers," *Physics of Plasmas* **1**:1626–1634 (1994).
260. R. Kodama, P. A. Norreys, K. Mima, A. E. Dangor, R. G. Evans, H. Fujita, Y. Kitagawa, et al., "Fast Heating of Ultrahigh-Density Plasma as a Step Towards Laser Fusion Ignition," *Nature* **412**:798–802 (2001).
261. J. Ren, W. F. Cheng, S. L. Li, and S. Suckewer, "A New Method for Generating Ultraintense and Ultrashort Laser Pulses," *Nature Physics* **3**:732–736 (2007).

SLOW LIGHT PROPAGATION IN ATOMIC AND PHOTONIC MEDIA

Jacob B. Khurgin

*Department of Electrical and Computer Engineering
Johns Hopkins University
Baltimore, Maryland*

22.1 GLOSSARY

ω	field frequency
ω_{12}	resonant frequency of an atomic transition
f_{12}	oscillator strength of an atomic transition
γ_{12}	damping coefficient (dephasing rate) of the transition
$\epsilon(\omega)$	dielectric constant
$n(\omega)$	refractive index
e	electron charge
m	free electron mass
Ω_p	plasma frequency
Ω	Rabi frequency
c	speed of light in vacuum
η_0	impedance of vacuum
$\alpha(\omega)$	absorption coefficient
$k(\omega)$	wave vector
λ	wavelength
v_g	group velocity
$S(\omega)$	slow down factor
$\beta_n(\omega)$	n th order dispersion
B	bit rate
N_{st}	number of bits stored (delayed) in an optical buffer
Λ	period of Bragg grating or other periodic photonic structure
δn	refractive index modulation
ω_B	Bragg frequency
$\Delta\omega_{gap}$	Photonic bandgap width

κ	coupling coefficient
τ	one half of the resonator round trip time
Δv_{pass}	width of passband of the photonic structure
$g(\omega)$	gain per unit length of optical amplifier
ΔT_d	delay time

22.2 INTRODUCTION

While the subject of slow light had become immensely popular in the last decade, the physics of the phenomenon of light propagation in media and structures with reduced group velocity, for which the term “slow light” had been coined, can be traced to the 19th century when the classical theory of dispersion of electromagnetic waves had been first formulated in works of Lorentz¹ and others. Slow wave propagation has also been observed and widely used in the microwave range since as early as the 1940s.² Building on this history and following the pioneering works of Refs. 3 and 4, the science and technology of slow light had been transformed from a scientific curiosity to a rapidly evolving field with many potential applications. Slow light propagation had been observed in a wide variety of media and structures ranging from Bose-Einstein condensates and low-pressure metal vapors on one hand, to optical fibers and photonic bandgap structures on the other hand. This makes slow light a truly interdisciplinary field that is not easy to describe in a short review. There exists a body of review work on slow light, starting in 2002 with reviews by Boyd and Gauthier⁵ and Milloni⁶ who also published a book in 2005 including chapters on slow, fast, and left-handed light.⁷ The first comprehensive treatment of diverse slow light, schemes has been given in the recent book,⁸ which contains contributions from 18 groups that have been actively involved in the slow-light field and have all made significant contributions in recent years.

With all the seeming diversity of slow light schemes, they can be all characterized by a single common feature—the existence of sharp single or multiple resonances. The resonance can be defined by a simple atomic transition, by a Bragg grating or other resonant photonic structure, or by an external laser as in the schemes involving various nonlinear processes—resonant scattering, spectral hole burning, or four-wave mixing. In this chapter, we try to emphasize the commonality of all slow light approaches as well as their distinctive features.

22.3 ATOMIC RESONANCE

Let us consider a simple atomic resonance characterized by the resonant frequency ω_{12} and the oscillator strength f_{12} as shown in Fig. 1a. The dielectric constant in the vicinity of the resonance has a familiar expression

$$\varepsilon(\omega) = \bar{\varepsilon} + \frac{\Omega_p^2}{\omega_{12}^2 - \omega^2 - j\omega\gamma_{12}} \quad (1)$$

where $\bar{\varepsilon} = \bar{n}^2$ is the nonresonant part or “background” part of the dielectric constant, and γ_{12} is the dephasing rate of the polarization. The expression in the numerator Ω_p is the “plasma frequency” that can be found as

$$\Omega_p^2 = \frac{N_a e^2}{\varepsilon_0 m_0} f_{12} \quad (2)$$

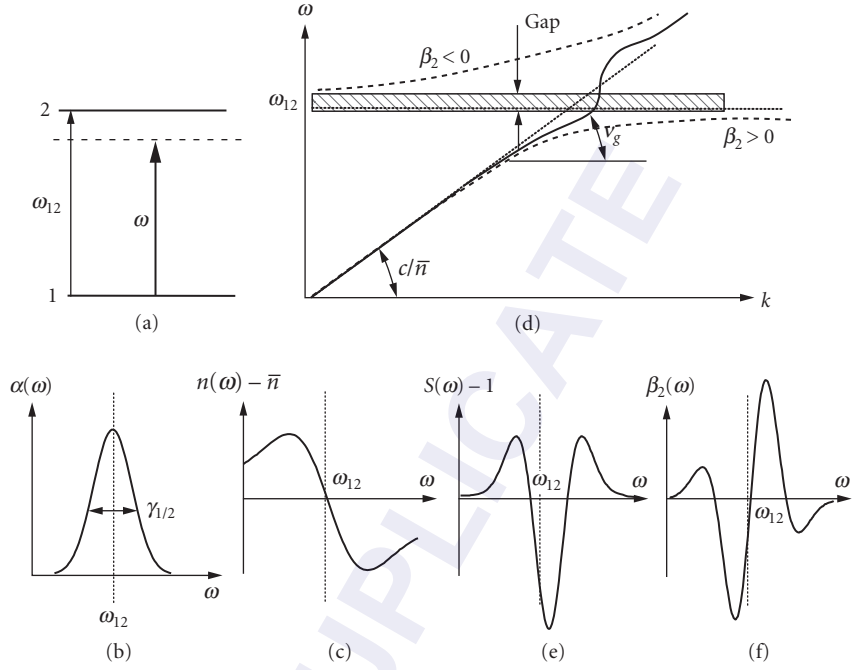


FIGURE 1 (a) Two-level atomic resonance; (b) spectrum of absorption coefficient; (c) spectrum of refractive index; (d) dispersion diagram and group velocity in the vicinity of a resonance; (e) spectrum of slowdown factor; and (f) spectrum of group velocity dispersion.

where e is an electron charge, ϵ_0 is the dielectric permittivity of vacuum, m_0 is a free electron mass, and N_a is the concentration of active atoms. From real and imaginary parts of Eq. (1) one immediately obtains the expressions for the absorption coefficient shown in Fig. 1b

$$\alpha(\omega) = \frac{2\omega}{c} \text{Im}(\epsilon^{1/2}) \approx \frac{1}{4\bar{n}c} \frac{\Omega_p^2 \gamma_{12}}{(\omega_{12} - \omega)^2 + \gamma_{12}^2/4} \quad (3)$$

and refractive index (Fig. 1c)

$$n(\omega) = \text{Re}(\epsilon^{1/2}) \approx \bar{n} + \frac{1}{4\bar{n}\omega} \frac{\Omega_p^2 (\omega_{12} - \omega)}{(\omega_{12} - \omega)^2 + \gamma_{12}^2/4} = \bar{n} + \frac{c}{\omega} \alpha(\omega) \frac{\omega_{12} - \omega}{\gamma_{12}} \quad (4)$$

the last relation being a particular form of more general Kramer-Kronig relation. In the vicinity of the resonance the refractive index first increases, then rapidly decreases to the point where $n(\omega) < \bar{n}$, and then increases once again to the value of background index. Introducing the wave vector as

$$k(\omega) = n(\omega) \cdot \omega / c \quad (5)$$

One can obtain the dispersion relation between the frequency and wave vector shown in Fig. 1d (solid line) and then introduce group velocity as

$$v_g = \partial\omega / \partial k \quad (6)$$

and slow down factor defined as a relative reduction of group velocity due to resonance

$$S(\omega) = \frac{c}{\bar{n}v_g(\omega)} = 1 + \frac{\Omega_p^2}{4\bar{E}} \frac{(\omega_{12} - \omega)^2 - \gamma_{12}^2/4}{[(\omega_{12} - \omega)^2 + \gamma_{12}^2/4]^2} = 1 + \frac{c\alpha(\omega)}{\bar{n}\gamma_{12}} \frac{(\omega_{12} - \omega)^2 - \gamma_{12}^2/4}{(\omega_{12} - \omega)^2 + \gamma_{12}^2/4} \quad (7)$$

shown in Fig. 1e. One can now see that there are two regions in which $S > 0$, that is, slow light regime, and one region where $S < 1$, which can be either the fast light regime ($0 < S < 1$) or negative group velocity regime ($S < 0$) indicating that the light gets reflected. One can get better intuitive picture of the physical phenomena governing slow and fast light propagation by plotting the dispersion curve in the absence of loss (i.e., $\gamma_{12} = 0$) as dashed line in Fig. 1d. This curve corresponds to a well-known coupled modes model, also known as polariton dispersion curve in solid state physics. The first mode is a photon which in the absence of atomic transition is described by linear dispersion curve $\omega_p = ck/\bar{n}$. The second mode is the atomic polarization characterized by a resonance frequency ω_{12} . The dispersion curve of atomic polarization is a horizontal line $\omega_a = \omega_{12}$ indicating the obvious fact that it has zero group velocity as atoms do not move, at least not on the scale of the speed of light. In the vicinity of the resonance, two modes couple into each other and the modified dispersion curve is split into two branches separated by the gap in which the light cannot propagate. Notice that for each wave vector there are two coupled solutions characterized by two different group velocities. The one further away from the resonance has higher group velocity and is usually referred to as a “photonlike,” while the one closer to the resonance has lower group velocity and is usually referred to as “atomlike.” One can then interpret the slow light propagation phenomenon in the following way: The energy gets constantly coupled from the electromagnetic field to the atomic polarization and back. The longer is the time the energy spends in the form of atomic excitation, the slower the coupled mode propagates. Thus the slow light propagation in atomic system can be understood as constant excitation and de-excitation of atoms in which coherence is preserved.

One important implication is related to the strength of the electric field in the electromagnetic wave propagating in a slow light regime caused by an atomic resonance. Since the group velocity is the velocity with which the energy propagates, the local energy density of light beam with power density P is

$$U = P/v_g = (P\bar{n}/c)S \quad (8)$$

and it gets enhanced by a slowdown factor in the slow light medium. But since the energy density is related to electric field as

$$U = \frac{1}{4}\epsilon_0 \frac{\partial(\omega\mathcal{E})}{\partial\omega} E^2 + \frac{1}{4}\mu_0 H^2 = S \frac{1}{2}\epsilon_0 \bar{E} E^2 \quad (9)$$

where the first term corresponds to electric energy and the second term to the magnetic energy. It follows from Eqs. (8) and (9) that $E^2 = 2\eta_0 P/\bar{n}$, where $\eta_0 = \sqrt{\mu_0/\epsilon_0} = 377\Omega$ is the vacuum impedance. Thus the electric field does not get enhanced in the atomic slow light medium. In our intuitive picture this simply means that all the additional energy *compressed* into the medium gets stored in the atomic polarization. This fact has important implications in the nonlinear optics. Also, one can see that without absorption the region of fast light ($0 < S < 1$) is absent in Fig. 1d—and it is important that fast light is associated with the absorption peak, while the slow light is associated only with the off-resonant absorption—hence significant delays of the signal in the slow light are quite possible, while significant advances are difficult to observe due to absorption.

When it comes to practical applications of the slow light, it is important to achieve large delays over significant bandwidth and in a compact device. The obstacles on the way to this goal include the aforementioned loss and the dispersion of group velocity and dispersion of absorption.^{9–11} To ascertain the importance of the group velocity dispersion (GVD), one only needs to use the Taylor

series expansion of the dispersion relation $k(\omega)$ near the signal frequency ω_0 to obtain the expression for the group velocity

$$v_g^{-1}(\omega) = v_g^{-1}(\omega_0) + \beta_2(\omega_0) \cdot (\omega - \omega_0) + \frac{1}{2} \beta_3(\omega_0) \cdot (\omega - \omega_0)^2 + \dots \quad (10)$$

where $\beta_n(\omega) = \partial^n k(\omega) / \partial \omega^n$. One can estimate the limitations imposed by GVD by noticing that the time delay of the signal can be written as

$$T_d(\omega) = v_g^{-1}(\omega)L = v_g^{-1}(\omega_0)L + \beta_2(\omega_0)L + \frac{1}{2} \beta_3(\omega_0)L^2 + \dots \quad (11)$$

and then introducing a criterion for the on-off keyed gaussian signal of bit rate B the difference between the delay times of spectral components at the edges of the signal bandwidth, that is, at $\omega_0 \pm \Delta\omega_{\text{sig},0}/2$, should be less than one-half of the bit interval, that is

$$\Delta T_d(L) \approx \beta_2 \Delta\omega_{\text{sig}} L < 1/2B \quad (12)$$

which leads to the condition relating the maximum allowable length and bit rate

$$|\beta_2| B^2 L < \frac{1}{16 \ln 2} \quad (13)$$

This condition shows the severe limitations imposed by the second order GVD term β_2 whose spectrum is plotted in Fig. 1f. In addition to GVD, the signal bandwidth is also limited by the dispersion of absorption—the fact that different frequency components of the signal experience different attenuation causes signal distortion, but it is usually GVD that is the main reason that limits slow light scheme performance. Hence, while slow light had been observed in single-atomic resonance schemes as early as 1960s and 1970s^{12–16} it was only in 1990s when a double resonance schemes with complete cancellation of β_2 had been discovered¹⁷ that the slow light research had really taken off.

Double Atomic Resonance

Although the first practical slow light results were achieved using phenomenon of electromagnetic transparency (EIT),^{17,18} where the double resonance is created by strong pump beam, the main features of the double-resonant atomic schemes can be understood easier if one considers closely spaced narrow resonances do occur naturally in metal vapors, such as, for instance, in Rb⁸⁵ Ref. 19 where two D² resonances near 780 nm separated by $\nu_{32} = 3$ GHz have been used most successfully to this date in SL experiments in atomic medium. The rationale for using double resonance can be seen from Fig. 1d and f, which show that the lowest order GVD term β_2 has opposite signs below and above resonance. Then, if one can combine two resonances as in Fig. 2a, only the third-order GVD β_3 will be a factor for signals centered at frequency ω_0 in the middle between two transitions where the absorption (Fig. 2b) is low, as evident from the dispersion curves Fig. 2c and d.

A closer look at the dispersion curve Fig. 2c reveals the basic trade-off inherent in every slow light scheme—the dispersion curve is split into three branches of which the central one, “squeezed” between two atomic resonances, is the one with a slow group velocity. Clearly, the group velocity, being a slope of the curve is inversely proportional to the splitting between two resonances ν_{32} , that is, to the maximum theoretical bandwidth of the scheme. In reality, the

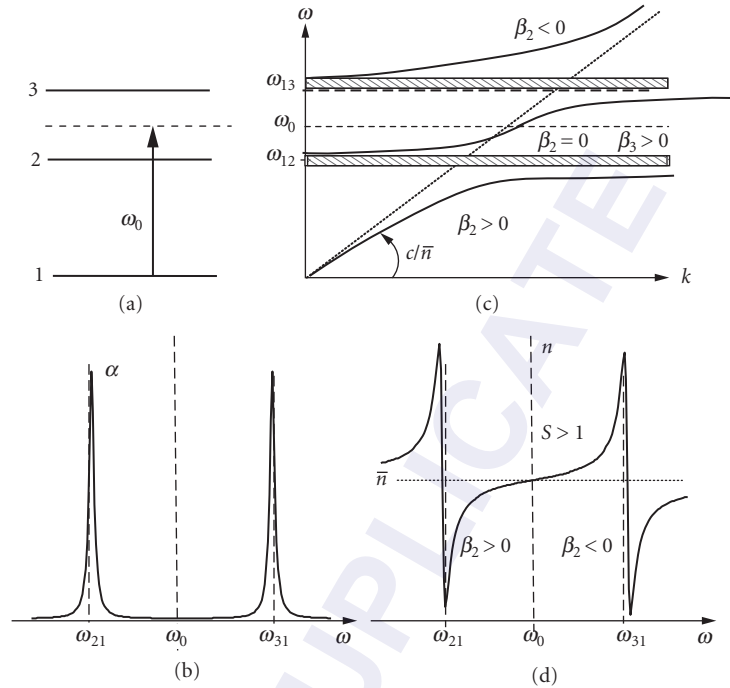


FIGURE 2 (a) Double atomic resonance and its (b) dispersion and group velocity near it; (c) absorption spectrum; and (d) refractive index spectrum.

practical bandwidth (or bit rate for digital signals) is even smaller than that and is mostly limited by the third-order dispersion^{11,20}

$$|\beta_3|B^3L < 1/16(\ln 2)^2 \quad (14)$$

In the works¹⁹ group delays in excess of 100 ns with 2-ns pulses at 780 nm had been obtained—that correspond to about 50-bit tunable optical buffers. The delay could be tuned by changing concentration of Rb vapor as shown in Fig. 3.

Tunable Double Resonance—Electromagnetically-Induced Transparency

The double atomic resonant scheme described cannot be adapted to variable bandwidth because the width of passband cannot be changed. To change the passband width, one can consider an alternative of spectral hole burning in the inhomogeneously broadened transition.^{21,22} As shown in Fig. 4, a strong pump pulse creates a situation at which the absorption in the frequency range $\Delta\omega$ becomes depleted. The profile of the absorption spectrum shown in Fig. 4b looks remarkably like the double-resonant profile of Fig. 2. With the refractive index profile shown in Fig. 4c, one can see that a strong reduction of group velocity can be expected near the center of the spectral hole.

By changing the spectrum of the pump, for instance, using intensity or frequency modulation, one can change $\Delta\omega$ to achieve the maximum delay without distortion for a given bit rate. Delays of 2-bit intervals were achieved in Ref. 23 for a moderate bandwidth of 100 MHz but only in a

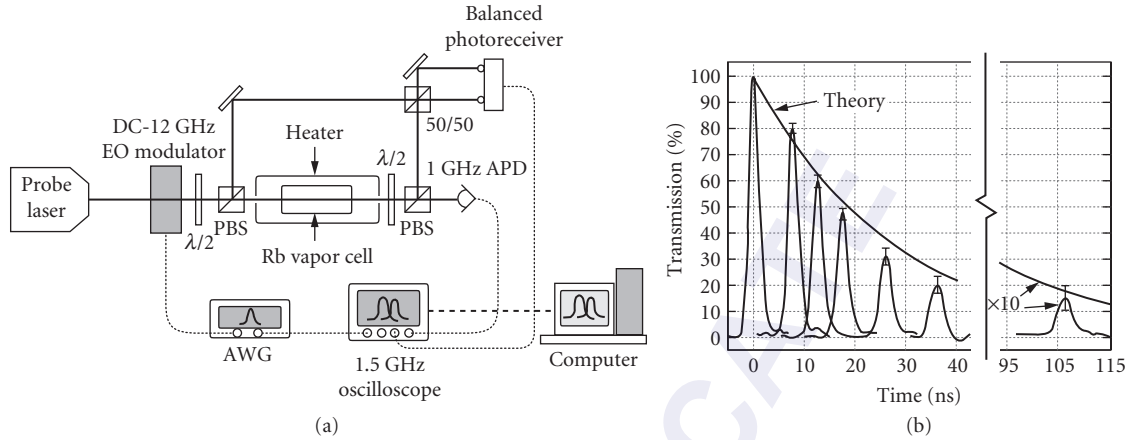


FIGURE 3 Experimental apparatus (a) used in Ref. 19 to obtain variable pulse delays at various optical depths using double resonance in Rb vapor (b).

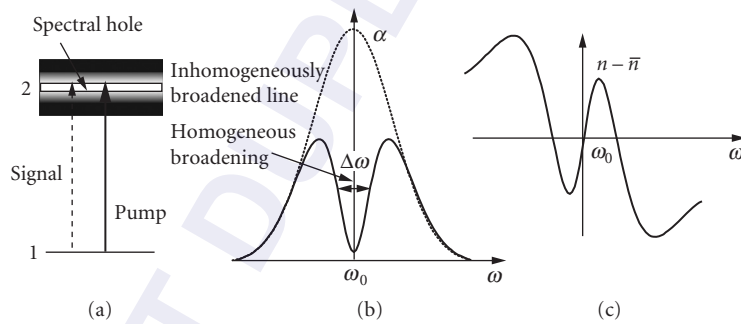


FIGURE 4 (a) Slow light scheme based on spectral hole burning; (b) absorption spectrum; and (c) refractive index spectrum.

40-cm-long Rb vapor delay line. Since the background absorption in the hole burning is always high, it is the dispersion of loss α_2 that causes the signal distortion and is in fact a limitation in this scheme. The scheme also suffers from the large energy dissipation as the pump gets absorbed.

To avoid large background absorption and to achieve wide passband tunability one uses an entirely different SL scheme based on electromagnetic transparency (EIT), first considered by Harris.^{24–26} Without trying to explain all the intricacies of EIT, one can understand the rationale of using it. Since finding two closely spaced atomic resonances is not trivial, one should consider the means for their creation artificially.

Now, the atomic oscillator in the absence of external modulation has just one resonant frequency ω_1 in its response spectrum, just like any harmonic wave whose spectrum contains just one frequency component ω_1 . But if the wave is amplitude modulated with some frequency Ω , there will appear two sidebands at frequencies $\omega_1 \pm \Omega$ in its spectrum, and when the modulation depth reaches 100 percent the carrier frequency ω_0 would get entirely suppressed and the spectrum shows just two sidebands separated by 2Ω . Now, according to this analogy, if one strongly modulates the strength of atomic oscillator with some external frequency Ω , one should expect the absorption spectrum to behave in fashion similar to the spectrum of amplitude modulated wave, that is, it should show two absorption lines separated by 2Ω . The material should then become more transparent at the resonance frequency ω_0 —hence the term “EIT.”

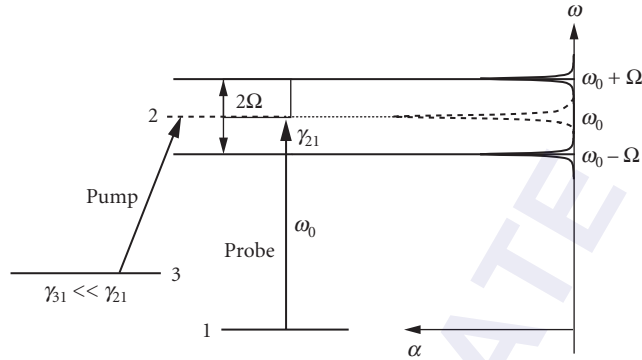


FIGURE 5 Principle of electromagnetic transparency in atomic “ Λ ” scheme.

To accomplish the EIT transmission modulation there exist numerous schemes, but we shall consider only one—the most widely used three-level “ Λ ” scheme¹ shown in Fig. 5 in which the ground-to-excited state transition ω_{12} is resonant with the frequency of the optical signal carrier ω_0 and has a dephasing rate of γ_{21} . In the absence of pump, the absorption spectrum (dashed line in Fig. 5) is a normal Lorentzian line. There also exists a strong transition that couples the excited level 2 with the third level 3, but notably, the transition between levels 1 and 3 is forbidden. When a strong resonant pump with intensity I_{pump} at frequency ω_{23} is turned on, the mixing of states 2 and 3 causes modulation of the absorption of signal. As expected, the Lorentzian peak in the absorption spectrum splits into two smaller peaks at frequencies $\omega_0 \pm \Omega$, where the Rabi frequency is

$$\Omega = \left(f_{23} \frac{12\pi\alpha_f I_{\text{pump}}}{m\bar{n}\omega_{23}} \right)^{\frac{1}{2}} \quad (15)$$

where $\alpha_f = 1/137$ is a fine structure constant. Thus changing pump intensity allows one to achieve full tunability of the group velocity and to achieve very small group velocities.^{3,4} But aside from being a convenient technique of implementing the tunable double resonant scheme, EIT has a significant advantage over other schemes because the residual absorption rate at the resonant frequency

$$\alpha(\omega_0) = \frac{1}{nc} \frac{\Omega_p^2}{8\Omega^2} \gamma_{31} \quad (16)$$

is proportional to the dephasing rate of the intra-atomic excitation 31, which is not coupled to outside world. Thus typically $\gamma_{31} \ll \gamma_{21}$ and the residual absorption is much weaker in the EIT than in the case of two independent resonances. It indicates that EIT is a coherent effect and the reduction of absorption occurs because of the destructive interference of the absorption by two sidebands. This effect has a rather simple physical interpretation. While in the simple single, or double-resonant slow light scheme, the energy transferred from the electromagnetic wave to the atomic excitation and back in the EIT scheme the process involves more steps. First, as the signal photon propagates in the EIT medium, it transfers its energy to the excitation of atomic transition between levels 1 and 2. Due to the presence of strong pump wave coupling between levels 2 and 3, the excitation is almost instantly transferred to the long-lived excitation between the levels 1 and 3 and then the process occurs in reverse until the energy is transferred back into the photon. Then the process repeats itself. Overall, most of the time the energy gets stored in the form of 1-3 excitations and thus it propagates with a very slow group velocity. Furthermore, the actual absorption event occurs only when

the excitation 1-3 loses coherence and the energy cannot get back to the photon. Naturally, it is the dephasing rate of this excitation, that is, $\gamma_{31} \ll \gamma_{21}$, which determines the residual absorption loss in Eq. (16). We stress here once again that since the energy gets stored in the form of atomic excitation, one cannot expect enhancement in the strength of the optical field.

22.4 BANDWIDTH LIMITATIONS IN ATOMIC SCHEMES

The most spectacular results were achieved in slow light experiments of EIT in which the light velocity was slowed down to the pedestrian speed³ and then even stopped⁴ first in metal vapors and then in solid state medium containing rare earth ions.²⁷ These achievements are of great importance for physics when it comes to manipulating single photons²⁸ and coherent control.²⁹ Also important are imaging applications of slow light where indeed impressive results were observed by a number of groups,^{30,31} extra high resolution interferometers,^{32,33} and rotation sensors.³⁴ There has also been significant progress on using electrically pumped semiconductor medium to achieve slow light in EIT configuration or using coherent population oscillations.^{35–37} But overall the best results in term of the delay-bit rate product were achieved in metal vapors^{19,24} at relatively narrow bandwidths. The delays were usually limited by the third-order dispersion. These results follow the basic properties of the lorentzian dispersion—according to Eq. (7) the slowdown factor S is proportional to $(\omega_{12} - \omega)^{-2}$ —hence the total delay can be very large in the narrow frequency band near the resonance, but then it changes and becomes much smaller. A number of publications have been devoted to the limits of delay-bit rate product.^{9–11,38,39} In fact, it was shown in Ref. 11 and then in Ref. 39 that the most relevant figure of merit for slow light delay line is the minimum length of the delay line L required to store a number of bits N_{st} at a given bit rate B . The required length was found in Ref. 11 to be

$$L \sim cBN_{st}^2[\Omega_p/2\pi]^{-2} \quad (17)$$

indicating that the performance suffers at higher bit rate and also that the required length increases nonlinearly with storage capacity. Indeed, most of the demonstrated and proposed slow light schemes based on atomic transitions do not show spectacular results at bit rates above 1 Gbit/s and for this reason a different class of resonances shall be considered.

22.5 PHOTONIC RESONANCE

As we have already mentioned, in case of atomic resonance the apparent slow down of light is caused by the resonant energy transfer to and from the excitation of atomic polarization. An entirely different resonance is the photonic resonance in which the energy is resonantly transferred between two or more modes of electromagnetic radiation. When the transfer takes place between forward and backward propagating wave slow light, the group velocity gets reduced and we are once again faced with the slow light phenomenon, albeit of entirely different nature from the slow light in atomic medium.

The most common photonic resonance is the Bragg grating (Fig. 6a)—a structure in which the refractive index is periodically modulated with period Λ ,

$$n = \bar{n} + \delta n \cos\left(\frac{2\pi}{\Lambda}z\right) \quad (18)$$

As a result, a photonic bandgap opens in the vicinity of Bragg frequency⁴⁰

$$\omega_B = \frac{\pi c}{\Lambda \bar{n}} \quad (19)$$

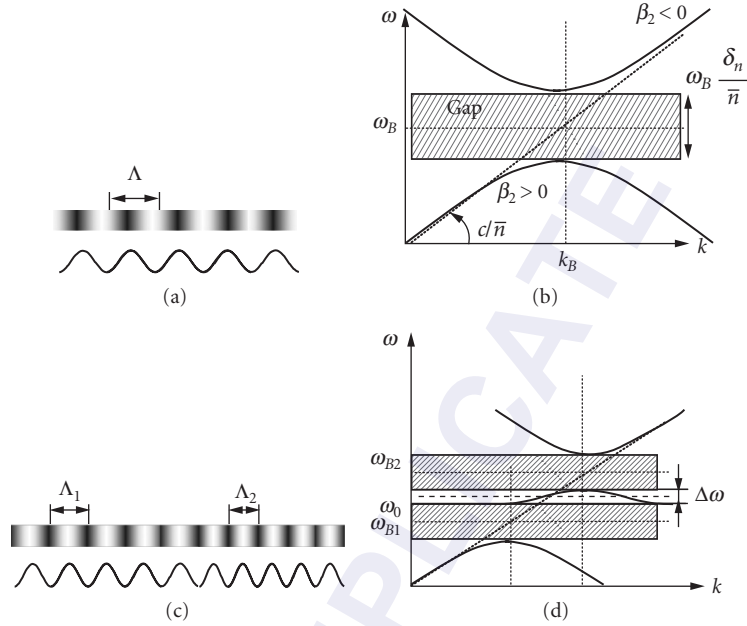


FIGURE 6 (a) Bragg grating and its index profile; (b) dispersion of Bragg grating; (c) cascaded Bragg grating; and (d) dispersion curve of cascaded grating.

and the dispersion law becomes modified as

$$\frac{k - k_B}{k_B} = \sqrt{\left(\frac{\omega - \omega_B}{\omega_B}\right)^2 - \left(\frac{\delta n}{2\bar{n}}\right)^2} \quad (20)$$

This dispersion law is plotted in Fig. 6b. One can see the similarities between it and the dispersion of a single atomic resonance (Fig. 1d), especially the presence of the gap near the resonance frequency indicating that the light gets reflected from the grating just as it gets reflected from the atomic medium at resonance. Close to the gap, the group velocity indeed gets reduced with the slowdown factor being

$$S = \frac{\frac{|\omega - \omega_B|}{\omega_B}}{\sqrt{\left(\frac{\omega - \omega_B}{\omega_B}\right)^2 - \left(\frac{\delta n}{2\bar{n}}\right)^2}} \quad (21)$$

Furthermore, the width of forbidden gap is $\Delta\omega_{\text{gap}} = \omega_B(\delta n/\bar{n})$, and hence the index contrast $\delta n/\bar{n}$ can be called the “strength” of the grating. This “grating strength” plays a role equivalent to that played by the oscillator strength of the atomic resonance. But the physics is quite different—the slow down effect in a photonic structure is the result of the transfer of energy between the forward and backward propagating waves—no energy is transferred to the medium. Hence the strength of the electric field in photonic SL structures gets greatly enhanced with important implications for non-linear optics. One can also use simple photonic crystals^{41,42} whose dispersion curves are similar to

Bragg gratings but the problem of structures with a single photonic resonance is identical to the single atomic resonance—strong second-order dispersion β_2 . Clearly, in order to compensate it one should consider structures with more than one photonic resonance.

Double Resonant Photonic Structures

Since any atomic SL scheme requires operation near a particular narrow linewidth absorption resonance, finding such a resonance near a particular wavelength is not an easy task, and, in fact, only a very few absorption lines have been employed in practice, Rb vapors being a “workhorse.” Finding two closely spaced narrow lines is even more difficult, and even if such two lines can be found, the splitting between them ν_{32} is fixed—hence the SL delay will be optimized for one particular combination of storage capacity and bit rate.

In contrast, the photonic double resonant can be easily implemented by simply combining two Bragg gratings with slightly different periods Λ_1 and Λ_2 as shown in Fig. 6c. Such combination was first suggested for dispersion compensation^{43,44} and then considered for application in electro-optic modulators.⁴⁵ As long as one deals with linear devices, such as delay lines, one can simply cascade two Bragg gratings sequentially and the resulting dispersion curve will be simply the mean of the individual dispersion curves. For the nonlinear and electro-optic devices one can alternate the short segments of Bragg gratings with periods Λ_1 and Λ_2 . The dispersion curve of Fig. 6d is remarkably similar to the dispersion curve of the atomic double resonance in Fig. 2b. Two gratings engender two photonic bandgaps, centered at $\omega_{B,i} = \pi c / \Lambda_i \bar{n}$ of almost equal widths $\Delta\omega_{\text{gap},i} = \omega_{B,i} (\delta n / \bar{n}) \approx \omega_0 (\delta n / \bar{n})$ with a narrow passband $\Delta\omega$ between them. By choosing the periods Λ_1 and Λ_2 for a given index modulation δn one can design $\Delta\omega$ to be arbitrarily narrow or wide. This fact gives the designer true flexibility. In fact, one can show the slowdown factor in this scheme is

$$S(\omega_0) \approx \left[1 + \frac{\Delta\omega_{\text{gap}}}{2\Delta\omega} \right]^{1/2} \quad (22)$$

Thus, as one can see, the dispersion curve gets *squeezed* between two gaps and the group velocity decreases with the passband, but the dependence is not as strong as in the case of atomic resonance, thus the photonic structures in general should have far superior performance at large bandwidth when compared to atomic medium. A similar approach can be cascaded photonic crystal waveguides as demonstrated in Refs. 46 to 50.

However, cascaded Bragg gratings have also a number of disadvantages, the first of which is a relatively small index contrast available, and the second is difficulty in fabricating two gratings with a prescribed value of the frequency offset. Furthermore, the cascaded geometry is applicable only to the linear devices that do not incorporate any nonlinear or electro-optic component. For this reason it is preferable to use alternating short segments of gratings with different periods. But a periodic sequence of short Bragg grating segments can be considered a new grating with periodically modulated properties, or Moiré grating (Fig. 7a). In a Moiré grating the segments are not independent but interact coherently—hence its properties are somewhat different from the cascaded grating as shown in Ref. 51, with the main distinction being the fact that the dispersion curve of Moiré grating with period d is also periodic in wave-vector space with a period $2\pi/d$. The ability of a Moiré grating to slow down the light was first predicted in Ref. 51 and demonstrated in Ref. 52. It was also noted that the Moiré grating is only example of periodically structured photonic media in which slow light can be observed. In the periodically structured media, the light energy density is distributed periodically and the periodically spaced regions of high intensity can be thought of as the resonators coupled to each other. Thus we shall refer to them as coupled resonator structures (CRS).⁵³ Aside from using Moiré gratings, CRS can be fabricated by coupling Fabry-Perot resonators (Fig. 7b), ring resonators (Fig. 7c),⁵⁴ or so-called “defect modes” in the photonic crystal (Fig. 7d).^{55,56} These and other CRS implementations are discussed at length in the literature^{41,57–61} and here we shall only give a short description of their properties and compare them with the EIT-like photonic structures.

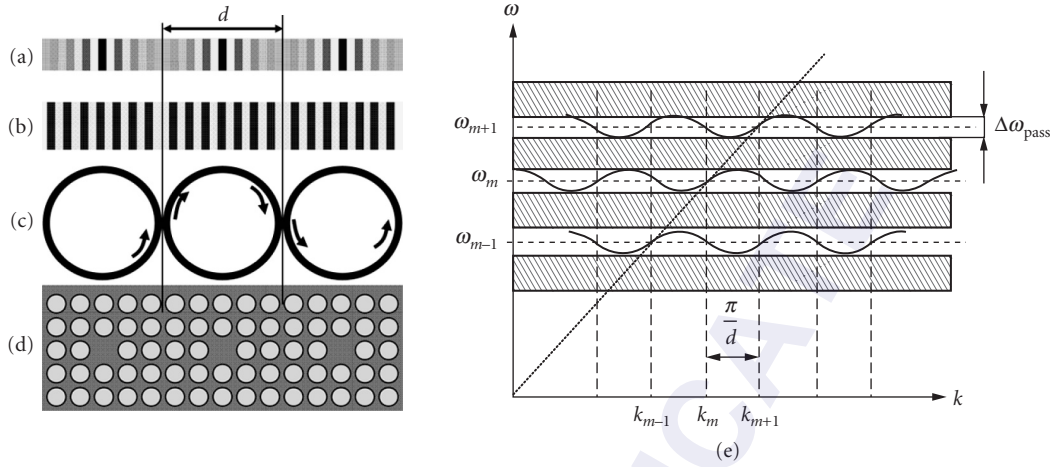


FIGURE 7 Photonic SL structures based on coupled resonators (CRS): (a) Moiré grating; (b) coupled Fabry-Perot resonators; (c) coupled ring resonators; (d) coupled defect modes in photonic crystal; and (e) dispersion in a typical CRS.

A periodic chain of coupled resonators is characterized by three parameters: period d , the time of one way pass through each resonator τ , and the coupling (or transmission) coefficient κ . The dispersion relation in this chain can be written as

$$\sin \omega \tau = \kappa \sin kd \quad (23)$$

The dispersion curves are shown in Fig. 7e and consist of the series of passbands around resonant frequencies $\nu_m = m/2\tau$ separated by the wide gaps. Presence of multiple resonances indicates that the light propagating through CRS can be considered a superposition of more than one forward and more than one backward propagating wave. Alternatively, one can also think about the resonators as “photonic atoms” analogous to real atoms in EIT. At any rate, the dispersion curve is quite similar to the EIT dispersion curve in sense that it gets *squeezed* into a narrow passband of width

$$\Delta \nu_{\text{pass}} = (\pi \tau)^{-1} \sin^{-1}(\kappa) \quad (24)$$

The three parameters d , τ , and κ are not independent of each other. First of all d and τ are obviously related to each other. This relation can be obtained from taking a limit of Eq. (23) at $\kappa = 1$

$$\frac{d}{\tau} = \frac{\omega}{k} = \frac{c}{\bar{n}} \quad (25)$$

which simply indicates that with 100 percent coupling the light simply propagates through medium without reflections. Also related are the size of resonator, that is, d and the coupling coefficient κ ; to achieve small κ , one needs to confine the light tightly within the resonators which requires large spacing between them. If the index contrast $\delta n/\bar{n}$ is large, a high degree of confinement can be achieved within a relatively small resonator; otherwise the light will leak from one resonator to another. This issue has been addressed in detail in Ref. 11, but here we simply assume that one uses the smallest resonator size that can be fabricated using a technology with a given index contrast.

Using Taylor expansion of the dispersion relation [Eq. (24)], one obtains the group velocity

$$v_g^{-1} = \frac{\tau}{d\kappa} = \frac{\bar{n}}{c} \kappa^{-1} \quad (26)$$

Hence the slowdown factor is simply

$$S = \kappa^{-1} \quad (27)$$

At the same time the third-order dispersion is

$$\beta_3 = v_g^{-1} \left(\frac{\tau}{\kappa} \right)^2 (1 - \kappa^2) \quad (28)$$

indicating that one can always optimize the performance of CRS delay line by choosing proper coupling coefficient to strike the balance between sufficient delay and low distortion caused by dispersion. As a result one can achieve the following relation between the storage capacity N_{st} and the required delay line length

$$L \sim c N_{\text{st}}^{3/2} \left(\omega \frac{\delta n}{n} \right)^{-1} \quad (29)$$

Comparing Eq. (29) with the results for atomic delay line Eq. (17), one can first notice that the required length of the buffer (essentially the number of coupled resonators) does not increase with the bit rate. Thus the CRS should have much better performance at high bit rates compared to EIT delay lines. This result (and the fact that required length increase only as a power 3/2 of storage capacity) follows the fact that the slowdown factor is not as strong function of passband width in CRS as in EIT. Furthermore, the expression inside the parenthesis can become commensurate with the optical frequency when large index contrast is used. This contrast can be as high as a factor of 2 in Si on SOI ring resonators^{61,62} or photonic crystal,⁴¹ which means that one can ultimately store N_{st} bits of optically encoded information in roughly $L \sim N_{\text{st}}^{3/2} \lambda$ length.⁶³

The limitations imposed by high order dispersion can be mitigated by using a combination of CRS with β_3 's of different sign as shown in.⁶⁴ Alternatively, one can use the so-called dynamic slow light scheme⁶⁵⁻⁶⁷ in which the properties of the CRS are adiabatically tuned when the light enters it allowing the compression of spectrum and thus reducing the deleterious effect of dispersion. Similar results can be in principle achieved using semiconductor devices^{68,69} combining Bragg grating with quantum wells that are resonant near the Bragg frequency. It should be noted, however, that even dynamic structures have limited storage capacity due to dispersion occurring before the spectral compression, that is, when the light enters the delay line.^{70,71}

22.6 SLOW LIGHT IN OPTICAL FIBERS

In the prior subsections, we have described the slow light in the vicinity of atomic and photonic resonances that are always associated with loss. In case of atomic resonance, the loss is inherent and can be traced to dephasing of atomic polarization. In case of photonic resonance, the loss occurs simply because the light spends more time inside the delay line (or, in one case, it takes longer effective path by bouncing back and forth or making circles). But it is well known that strong dispersion also takes place in the spectral vicinity of the resonant gain. The advantages of using gain are manifold—first of all one is not faced with attenuation, and second the gain can be changed at will by changing the pump strength and spectrum—hence the delay and the bandwidth can be made tunable.

Consider a gain profile $g(\omega)$ shown in Fig. 8a and characterized by its peak value $g_0 = g(\omega_0)$ and the FWHM linewidth $\delta\omega_{1/2}$. Usually the peak gain is inversely proportional to the linewidth and therefore it makes sense to introduce the integrated optical gain that is proportional to the total pump power $\gamma \sim g_0 \delta\omega_{1/2}$. Applying Kramers-Kronig transform to the gain profile, one can obtain

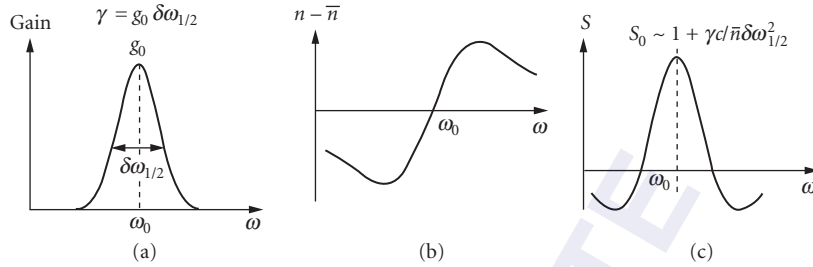


FIGURE 8 Slow light in an optical amplifier: (a) gain; (b) index; and (c) slowdown factor.

the refractive index spectrum shown in Fig. 8*b*. There is a positive slope at the center of the gain and this leads to slowdown factor of

$$S(\omega_0) \sim 1 + g_0 c / \bar{n} \Delta\omega_{1/2} \sim 1 + \gamma c / \bar{n} \Delta\omega_{1/2}^2 \quad (30)$$

shown in Fig. 8*c*. To achieve tunable gain, most often one uses stimulated Brillouin or stimulated Raman scattering. One unique advantage of these processes is that it can be observed in many media including the ubiquitous silica optical fiber which allows the slow light devices to be easily integrated into communication systems.

In the stimulated Brillouin scattering (SBS) process, a high-frequency acoustic wave is induced in the material via electrostriction for which the density of a material increases in regions of high optical intensity. The process of SBS can be described classically as a nonlinear interaction between the pump (at angular frequency ω_p) and a probe (or Stokes) field ω_s through the induced acoustic wave of frequency Ω_A .⁷² The acoustic wave in turn modulates the refractive index of the medium and scatters pump light into the probe wave when its frequency is downshifted by the acoustic frequency. This process leads to a strong coupling between the three waves when this resonance condition is satisfied, which results in exponential amplification (absorption) of the probe wave. Efficient SBS occurs when both energy and momentum are conserved, which is satisfied when the pump and probe waves counterpropagate. This leads to atypically very narrow gain bandwidth of typically about 30 MHz. With the peak values of Brillouin gain of about $g_0 \sim 0.02 \text{ m}^{-1}$ one can then expect the slowdown factor to be on the order of only a few percent different from unity—hence the value of this scheme is not in the absolute delay time but in adjustable additional delay

$$\Delta T_d \sim g_0 L / \Delta\omega_{1/2} \quad (31)$$

Since the large value of total gain $g_0 L$ can lead to Brillouin lasing, it is usually limited to about 10 to 15 meaning that maximum pulse delay of about 2 pulse lengths.⁷³ It is important that the main limitation for the achievable delay in amplifiers is the dispersion of gain rather than dispersion of group velocity.⁷³ While this delay of a few pulse lengths is relatively small, it is tunable and as such may be sufficient for many applications.

Tunable slow-light delay via SBS in an optical fiber was first demonstrated independently by Song et al.⁷⁴ and Okawachi et al.⁷⁵ In the experiment performed in Ref. 74 delay could be tuned continuously by as much as 25 ns by adjusting the intensity of the pump field and the technique can be applied to pulses as short as 15 ns. A fractional slow-light delay of 1.3 was achieved for the 15-ns-long input pulse with a pulse broadening of 1.4.

Following the first demonstrations of SBS slow light in optical fibers, there has been considerable interest in exploiting the method for telecommunication applications. One line of research

has focused on reducing pulse distortion by reducing the distortion caused by the gain dispersion performance^{76–81} by essentially trying to shape the Brillouin gain with multiple pumps. In numerous experimental demonstrations of multiple closely spaced SBS gain lines generated by a multi-frequency pump, significant increases in slow-light pulse delay were achieved as compared with the optimum single-SBS-line delay.

Another line of research has focused on broadband SBS slow light.^{82–86} The width of the resonance which enables the slow-light effect limits the minimum duration of the optical pulse that can be effectively delayed without much distortion, and therefore limits the maximum data rate of the optical system. Herraes et al.⁸² were the first to increase the SBS slow-light bandwidth and achieved a bandwidth of about 325 MHz by broadening the spectrum of the SBS pump field. Zhu et al. extended this work to achieve an SBS slow-light bandwidth as large as 12.6 GHz, thereby supporting a data rate of over 10 Gb/s.⁸³ The latest results on expanding bandwidth are summarized in the review article.⁸⁷

An alternative way to achieve tunable delays in optical fibers is via stimulated Raman scattering (SRS) can also be achieved in optical fibers. The SRS arises from exciting vibrations in individual molecules also known as optical phonons—as opposed to exciting sound waves (acoustic phonons) as in the SBS process. The optical phonons, unlike acoustic phonons, are localized and have very short lifetimes, measured in picoseconds or fractions of picoseconds. Furthermore, in amorphous materials, such as glass, the frequencies of optical phonons are spread over the large interval (measured in THz)—hence the Raman gain is much broader than Brillouin gain, but is also much smaller in absolute value. Integrated gain $\gamma \sim g_0 \delta\omega_{1/2}$ is similar for SBS and SRS but in SRS it is spread out over much wider $\delta\omega_{1/2}$.

Sharping et al. demonstrated an ultrafast all-optical controllable delay in a fiber Raman amplifier.⁸⁸ In this experiment a 430-fs pulse is delayed by 85 percent of its pulse width using SRS in 1-km-long high nonlinearity fiber. The ability to accommodate the bandwidth of pulses shorter than 1 ps in a fiber-based system makes SRS slow light useful for producing controllable delays in ultrahigh bandwidth telecommunication systems.

In addition to optical fibers, SRS slow light has also been demonstrated in a silicon-on-insulator planar waveguide.⁸⁹ Since silicon is a single-crystalline material, the Raman gain is concentrated into the narrower bandwidth than in glass, but this bandwidth is still sufficient for delaying short optical pulses of 3 to 4 ps in a very short (8 mm) waveguide. This scheme represents an important step in the development of chip-scale photonics devices for telecommunication and optical signal processing.

Slow light propagation was also demonstrated in erbium-doped fibers^{90,91} using coherent population oscillations, but the bandwidth, related to the relaxation time in erbium ion was only on the order of kHz. Other methods included using the parametric gain⁹² in the optical fiber as well as taking advantage of EIT in hollow optical fibers.⁹³

22.7 CONCLUSION

In this chapter, we have only given a very brief overview of the physics involved in the phenomenon of slow (and fast) light propagation. We have shown that while the underlying mechanisms are diverse, they all are characterized by a relatively narrow resonance, whether it is due to an atomic transition, resonance in the photonic structure, or is determined by a pump in an optical amplifier. As a result, the delay bandwidth and physical dimensions of all slow-light media are interdependent and must be optimized for a given set of objectives that depends on application. The applications of slow light, in both linear and nonlinear optics, are so diverse that it is impossible to cover them all even briefly in this chapter. They involve optical communications, optical signal processing, microwave photonics, precise interferometric instruments, and many others, and many new ones appear every year. The multidisciplinary field of slow light experiences rapid growth and this brief survey of the current state of affairs is nothing more than a snapshot intended to provide sufficient background to the newcomers into this field and to suggest a number of sources from which a deeper knowledge can be gained.

22.8 REFERENCES

1. H. A. Lorentz, "Über die Beziehung zwischen der Fortpflanzung des Lichtes," *Wiedem. Ann.* **9**:641 (1880).
2. J. R. Pierce, "Traveling-Wave Tubes," *Bell Syst. Tech. J.* **29**:1 (1950).
3. L. V. Hau, S. E. Harris, Z. Dutton, and C. H. Behroozi, "Light Speed Reduction to 17 Metres per Second in an Ultracold Atomic Gas," *Nature* **397**:594–596 (1999).
4. D. F. Phillips, A. Fleischhauer, A. Mair, R. L. Walsworth, and M. D. Lukin, "Storage of Light in Atomic Vapor," *Phys. Rev. Lett.* **86**:783–786 (2001).
5. R. W. Boyd and D. J. Gauthier, *Progress in Optics*, Elsevier, Radarweg, Amsterdam, 2002, pp. 497–530.
6. P. W. Milonni, "Controlling the speed of light pulses," *J. Phys. B-At. Mol. Opt. Phys.* **35**:R31 (2002).
7. P. W. Milonni, *Fast Light, Slow Light, and Left-Handed Light*, Institute of Physics, Bristol and Philadelphia, 2005.
8. J. B. Khurgin and R. S. Taylor, *Slow Light: Science and Applications*, Taylor and Francis, New York, 2009.
9. R. W. Boyd, D. J. Gauthier, A. L. Gaeta, and A. E. Willner, "Maximum Time Delay Achievable on Propagation through a Slow-Light Medium," *Phys. Rev. A* **71**:023801 (2005).
10. A. B. Matsko, D. V. Strekalov, and L. Maleki, "On the Dynamic Range of Optical Delay Lines Based on Coherent Atomic Media," *Opt. Exp.* **13**:2210–2223 (2005).
11. J. B. Khurgin, "Optical Buffers Based on Slow Light in EIT Media and Coupled Resonator Structures—Comparative Analysis," *J. Opt. Soc. Am. B*, **22**:1062–1074 (2005).
12. S. L. McCall and E. L. Hahn, "Self-Induced Transparency by Pulsed Coherent Light," *Phys. Rev. Lett.* **18**:908 (1967).
13. D. Grischkowsky, "Adiabatic Following and Slow Optical Pulse Propagation in Rubidium Vapor," *Phys. Rev. A* **7**:2096 (1973).
14. S. Chu and S. Wong, "Linear Pulse Propagation in an Absorbing Medium," *Phys. Rev. Lett.* **48**:738 (1982).
15. S. Chu and S. Wong, "Linear Pulse Propagation in an Absorbing Medium—Response," *Phys. Rev. Lett.* **49**:1293 (1982).
16. B. Segard and B. Macke, "Observation of Negative Velocity Pulse Propagation," *Phys. Lett. A* **109**:213 (1985).
17. S. E. Harris, J. E. Field, and A. Imamoglu, "Nonlinear Optics Using Electromagnetically Induced Transparency," *Phys. Rev. Lett.* **64**:1107 (1990).
18. M. M. Kash, V. A. Sautenkov, A. S. Zibrov, L. Hollberg, G. R. Welch, M. D. Lukin, Y. Rostovtsev, E. S. Fry, and M. O. Scully, "Ultraslow Group Velocity and Enhanced Nonlinear Optical Effects in a Coherently Driven Hot Atomic Gas," *Phys. Rev. Lett.* **82**:5229 (1999).
19. R. M. Camacho, M. V. Pack, and J. C. Howell, "Low-Distortion Slow Light Using Two Absorption Resonances," *Phys. Rev. A* **73**:063812 (2006).
20. R. M. Camacho, M. V. Pack, J. C. Howell, A. Schweinsberg, and R. W. Boyd, "Wide-Bandwidth, Tunable, Multiple-Pulse-Width Optical Delays Using Slow Light in Cesium Vapor," *Phys. Rev. Lett.* **98**:153601 (2007).
21. M. S. Bigelow, N. N. Lepeshkin, and R. W. Boyd, "Observation of Ultraslow Light Propagation in a Ruby Crystal at Room Temperature," *Phys. Rev. Lett.* **88**:023602 (2002).
22. A. V. Turukhin, V. S. Sudarshanam, M. S. Shahriar, and P. R. Hemmer, "Observation of Ultraslow and Stored Light Pulses in a Solid," *Phys. Rev. Lett.* **88**:023602 (2002).
23. R. M. Camacho, M. V. Pack, and J. C. Howell, "Slow Light with Large Fractional Delays by Spectral Hole-Burning in Rubidium Vapor," *Phys. Rev. A* **74**:033801 (2006).
24. K-J Boller, A. Imamoglu, and S. E. Harris, "Observation of EIT," *Phys. Rev. Lett.* **66**:2593–2596 (1991).
25. S. E. Harris, J. E. Field, and A. Kasapi, "Dispersive properties of EIT," *Phys. Rev. A* **46**:R39–R32 (1992).
26. A. Kasapi, M. Jain, G. Y. Jin, and S. E. Harris, "EIT: Propagation Dynamics," *Phys. Rev. Lett.* **74**:2447–2450 (1995).
27. A. V. Turukhin, V. S. Sudarshanam, M. S. Shahriar, J. A. Musser, B. S. Ham, and P. R. Hemmer, "Observation of Ultraslow and Stored Light Pulses in a Solid," *Phys. Rev. Lett.* **88**(2):023602 (2002).
28. M. D. Lukin, A. Imamoglu, "Nonlinear Optics and Quantum Entanglement of Ultraslow Single Photons," *Phys. Rev. Lett.* **84**:1419 (2000).

29. F. E. Zimmer, A. Andre, M. D. Lukin, and M. Fleischhauer, "Coherent Control of Stationary Light Pulses," *Opt. Commun.* **264**:441 (2006).
30. R. M. Camacho, C. J. Broadbent, I. Ali-Khan, and J. C. Howell, "All-Optical Delay of Images Using Slow Light," *Phys. Rev. Lett.* **98**:043902 (2007).
31. O. Firstenberg, O. M. Shuker, M. N. Davidson, N. A. Ron, "Elimination of the Diffraction of Arbitrary Images Imprinted on Slow Light," *Phys. Rev. Lett.* **102**:043601 (2009).
32. Z. Shi and R. W. Boyd, "Slow-Light Interferometry: Practical Limitations to Spectroscopic Performance," *J. Opt. Soc. Am.* **25**:136 (2008).
33. Z. Shi, R. W. Boyd, D. J. Gauthier, and C. C. Dudley, "Enhancing the Spectral Sensitivity of Interferometers Using Slow-Light Media," *Opt. Lett.* **32**:915 (2007).
34. M. S. Shakhriar, G. S. Pati, R. Tripathi, V. Gopal, M. Messall, and K. Salit, "Ultrahigh Enhancement in Absolute and Relative Rotation Sensing Using Fast and Slow Light," *Phys. Rev. A* **75**:53807 (2007).
35. P. K. Kondratko, S. W. Chang, H. Su, and S. L. Chuang, "Slow Light with Tunable Bandwidth in p-doped and Intrinsic Quantum Dot Electro-Absorbers," *Appl. Phys. Lett.* **90**:251108 (2007).
36. F. Ohman, K. Yivind, and J. Mork, "Voltage-Controlled Slow Light in an Integrated Semiconductor Structure with Net Gain," *Opt. Express* **14**:9955 (2006).
37. F. G. Sedgwick, B. Pesala, J. Y. Lin, W. S. Ko, X. X. Zhao, and C. J. Chang-Hasnain, "THz-Bandwidth Tunable Slow Light in Semiconductor Optical Amplifiers," *Opt. Express* **15**:747 (2007).
38. R. S. Tucker, P. C. Ku, and C. J. Chang-Hasnain, "Slow-Light Optical Buffers: Capabilities and Fundamental Limitations," *J. Lightwave Technol.* **23**:4046 (2005).
39. D. A. B. Miller, "Fundamental Limit to Linear One-Dimensional Slow Light Structures," *Phys. Rev. Lett.* **99**:203903 (2007).
40. T. Erdogan, "Fiber Grating Spectra," *J. Lightwave Technol.* **15**(8):1277–1294 (1997).
41. T. Baba, "Slow Light in Photonic Crystals," *Nature Photon.* **2**:465 (2008).
42. T. F. Krauss, "Slow Light in Photonic Crystal Waveguides," *J. Phys. D. Appl. Phys.* **40**:2666 (2007).
43. N. M. Litchinitser, B. J. Eggleton, and G. P. Agrawal, "Dispersion of Cascaded Fiber Gratings in WDM Lightwave Systems," *J. Lightwave Technol.* **16**:1523–1529 (1999).
44. S. Wang, H. Erlig, H. R. Fetterman, E. Yablonovitch, V. Grubsky, D. S. Starodubov, and J. Feinberg, "Group Velocity Dispersion Cancellation and Additive Group Delays by Cascaded Fiber Bragg Gratings in Transmission," *Microwave and Guided Wave Letters* **8**:327–329 (1998).
45. Khurgin JB, Kang JU, and Ding YJ, "Ultrabroad-Bandwidth Electro-Optic Modulator Based on a Cascaded Bragg Grating," *Opt. Lett.* **25**:70–72 (2000).
46. D. Mori and T. Baba, "Wideband and Low Dispersion Slow Light by Chirped Photonic Crystal Coupled Waveguide," *Opt. Express* **13**:9398–9408 (2005).
47. M. L. Povinelli, S. G. Johnson, and J. D. Joannopoulos, "Slow-Light, Band-Edge Waveguides for Tunable Time Delays," *Opt. Express* **13**:7145–7159 (2005).
48. D. Mori, S. Kubo, H. Sasaki, and T. Baba, "Experimental Demonstration of Wideband Dispersion-Compensated Slow Light by a Chirped Photonic Crystal Directional Coupler," *Opt. Express* **15**:5264–5270 (2007).
49. S. C. Huang, M. Kato, E. Kuramochi, C. P. Lee, and M. Notomi, "Time-Domain and Spectral-Domain Investigation of Inflection-Point Slow-Light Modes in Photonic Crystal Coupled Waveguides," *Opt. Express* **15**:3543–3549 (2007).
50. T. Kawasaki, D. Mori, and T. Baba, "Experimental Observation of Slow Light in Photonic Crystal Coupled Waveguides," *Opt. Express* **15**:10274–10281 (2007).
51. J. B. Khurgin, "Light Slowing Down in Moire Fiber Gratings and Its Implications for Nonlinear Optics," *Phys. Rev. A* **62**:3821–3824 (2000).
52. S. Longhi, D. Janner, G. Galzerano, G. Della Valle, D. Gatti, and P. Laporta, "Optical Buffering in Phase-Shifted Fibre Gratings," *Electron. Lett.* **41** (2005).
53. A. Yariv, Y. Xu, R. K Lee, and A. Scherer, "Coupled-Resonator Optical Waveguide: A Proposal and Analysis," *Opt. Lett.* **24**:711–713 (1999).
54. C. K. Madsen and G. Lenz, "Optical All-Pass Filters for Phase Response Design with Applications for Dispersion Compensation," *IEEE Photon. Tech. Lett.* **10**:994–996 (1998).

55. A. Melloni, F. Morichetti, and M. Martelli, "Linear and Nonlinear Pulse Propagation in Coupled Resonator Slow-Wave Optical Structures," *Opt. Quantum Electron.* **35**:365–378 (2003).
56. Z. Wang and S. Fan, "Compact All-Pass Filters in Photonic Crystals as the Building Block for High-Capacity Optical Delay Lines," *Phys. Rev. E* **68**:066616–06661623 (2003).
57. D. D. Smith, C. Hongrok, K. A. Fuller, A. T. Rosenberger, and R.W. Boyd, "Coupled-Resonator-Induced Transparency," *Phys. Rev. A-At. Mol. Opt. Phys.* **69**:63804 (2004).
58. J. E. Heebner and R.W. Boyd, "'Slow' and 'Fast' Light in Resonator-Coupled Waveguides," *J. Mod. Opt.* **49** (2002).
59. S. Mookhejea, "Dispersion Characteristics of Coupled-Resonator Optical Waveguides," *Opt. Lett.* **30**:2406 (2005).
60. M. Notomi, K. Yamada, A. Shinya, J. Takahashi, C. Takahashi, and I. Yokohama, "Extremely Large Group-Velocity Dispersion of line-defect waveguides in Photonic Crystal Slabs," *Phys. Rev. Lett.* **87**(25):253902 (2001).
61. F. Xia, L. Sekaric, and Y. Vlasov, "Ultracompact Optical Buffers on a Silicon Chip," *Nature Photon.* **1**:65 (2007).
62. A. B. Matsko, A. A. Savchenlov, and L. Maleki, "Vertically Coupled Whispering-Gallery-Mode Resonator Waveguide," *Opt. Lett.* 3066 (2005).
63. J. B. Khurgin, "Dispersion and Loss Limitations on the Performance of Optical Delay Lines Based on Coupled Resonant Structures," *Opt. Lett.* **32**:163–165 (2007).
64. J. B. Khurgin, "Expanding the Bandwidth of Slow-Light Photonic Devices Based On Coupled Resonators," *Opt. Lett.* **30**:513 (2005).
65. M. F. Yanik and S. Fan, "Stopping Light All Optically," *Phys. Rev. Lett.* **92**:083901 (2004).
66. M. F. Yanik and S. Fan, "Stopping Light in a Waveguide with an All-Optical Analog of Electromagnetically Induced Transparency," *Phys. Rev. Lett.* **93**:233903 (2004).
67. M. F. Yanik and S. Fan, "Dynamic Photonic Structures: Stopping, Storage, and Time-Reversal of Light," *Stud. Appl. Math.* **115**:233 (2005).
68. Z. S. Yang, N. H. Kwong, R. Binder, and A. L. Smirl, "Distortionless Light Pulse Delay in Quantum-Well Bragg Structures," *Opt. Lett.* **30**:2790–2792 (2005).
69. Z. S. Yang, N. H. Kwong, R. Binder, and A. L. Smirl, "Stopping, Storing, and Releasing Light in Quantum-Well Bragg Structures," *J. Opt. Soc. Am B* **22**:2144 (2005).
70. J. B. Khurgin, "Adiabatically tunable optical delay lines and Their Performance Limitations," *Opt. Lett.* **30**:2778 (2005).
71. R. S. Tucker, "The Role of Optics as Deelectronics in High-Capacity Routers," *J. Lightwave Technol.* **24**:4655 (2006).
72. R. W. Boyd, *Nonlinear Optics*, 2nd ed. Academic Press, San Diego, 2003.
73. J. B. Khurgin, "Performance Limits of Delay Lines Based on Optical Amplifiers," *Opt. Lett.* **31**(7):948–950 (2006).
74. K. Y. Song, M. G. Herr´aez, and L. Th´evenaz, "Observation of Pulse Delaying and Advancement in Optical Fibers Using Stimulated Brillouin Scattering," *Opt. Express* **13**:82–88 (2005).
75. Y. Okawachi, M. S. Bigelow, J. E. Sharping, Z. Zhu, A. Schweinsberg, D. J. Gauthier, R. W. Boyd, and A. L. Gaeta, "Tunable All-Optical Delays via Brillouin Slow Light in an Optical Fiber," *Phys. Rev. Lett.* **94**:153, 902 (2005).
76. M. D. Stenner, M. A. Neifeld, Z. Zhu, A. M. C. Dawes, and D. J. Gauthier, "Distortion Management in Slow-Light Pulse Delay," *Opt. Express* **13**:9995–10,002 (2005).
77. A. Minardo, R. Bernini, and L. Zeni, "Low Distortion Brillouin Slow Light in Optical Fibers Using AM Modulation," *Opt. Express* **14**:5866–5876 (2006).
78. Z. Shi, R. Pant, Z. Zhu, M. D. Stenner, M. A. Neifeld, D. J. Gauthier, and R. W. Boyd, "Design of a Tunable Time-Delay Element Using Multiple Gain Lines for Increased Fractional Delay with High Data Fidelity," *Opt. Lett.* **32**:1986–1988 (2007).
79. K. Y. Song, K. S. Abedin, K. Hotate, M. G. Herraiez, and L. Thevenaz, "Highly Efficient Brillouin Slow And Fast Light Using As₂Se₃ Chalcogenide Fiber," *Opt. Express* **14**:5860–5865 (2006).
80. A. Zadok, A. Eyal, and M. Tur, "Extended Delay of Broadband Signals in Stimulated Brillouin Scattering Slow Light Using Synthesized Pump Chirp," *Opt. Express* **14**:8498–8505 (2006).

81. T. Schneider, R. Henker, K. U. Lauterbach, and M. Junker, "Comparison of Delay Enhancement Mechanisms for SBS-Based Slow Light Systems," *Opt. Express* **15**:9606–9613 (2007).
82. M. G. Herraiez, K. Y. Song, and L. Thevenaz, "Arbitrary-Bandwidth Brillouin Slow Light in Optical Fibers," *Opt. Express* **14**:1395–1400 (2006).
83. Z. Zhu, A. M. C. Dawes, D. J. Gauthier, L. Zhang, and A. E. Willner, "Broadband SBS Slow Light in an Optical Fiber," *J. Lightwave Technol.* **25**:201–206 (2007).
84. T. Schneider, M. Junker, and K.-U. Lauterbach, "Potential Ultra Wide Slow-Light Bandwidth Enhancement," *Opt. Express* **14**:11082–11087 (2006).
85. K. Y. Song and K. Hotate, "25 GHz Bandwidth Brillouin Slow Light in Optical Fibers," *Opt. Lett.* **32**:217–219 (2007).
86. L. Yi, L. Zhan, W. Hu, and Y. Xia, "Delay of Broadband Signals Using Slow Light In Stimulated Brillouin Scattering with Phase-Modulated Pump," *IEEE Photon. Technol. Lett.* **19**:619–621 (2007).
87. L. Thevenaz, "Slow Light in Optical Fibers," *Nature Photonics* **2**:474–481 (2008).
88. J. E. Sharping, Y. Okawachi, and A. L. Gaeta, "Wide Bandwidth Slow Light Using a Raman Fiber Amplifier," *Opt. Express* **13**:6092–6098 (2005).
89. Y. Okawachi, M. A. Foster, J. E. Sharping, A. L. Gaeta, Q. Xu, and M. Lipson, "All-Optical Slow-Light on a Photonic Chip," *Opt. Express* **14**:2317–2322 (2006).
90. A. Schweinsberg, N. N. Lepeshkin, M. S. Bigelow, R. W. Boyd, and S. Jarabo, "Observation of Superluminal and Slow Light Propagation in Erbium-Doped Optical Fiber," *Europhys. Lett.* **73**:218–224 (2006).
91. H. Shin, A. Schweinsberg, G. Gehring, K. Schwertz, H. J. Chang, R. W. Boyd, Q.-H. Park, and D. J. Gauthier, "Reducing Pulse Distortion in Fast-Light Pulse Propagation Through an Erbium-Doped Fiber Amplifier," *Opt. Lett.* **32**:906–908 (2007).
92. D. Dahan and G. Eisenstein, "Tunable All Optical Delay via Slow and Fast Light Propagation in a Raman Assisted Fiber Optical Parametric Amplifier: A Route to All Optical Buffering," *Opt. Express* **13**:6234–6249 (2005).
93. S. Ghosh, A. R. Bhagwat, C. K. Renshaw, S. Goh, A. L. Gaeta, and B. J. Kirby, "Low-Lightlevel Optical Interactions with Rubidium Vapor in a Photonic Band-Gap Fiber," *Phys. Rev. Lett.* **97**:023, 603 (2006).

This page intentionally left blank.

DO NOT DUPLICATE

QUANTUM ENTANGLEMENT IN OPTICAL INTERFEROMETRY

Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver, and
Jonathan P. Dowling

*Hearne Institute for Theoretical Physics
Department of Physics and Astronomy
Louisiana State University
Baton Rouge, Louisiana*

23.1 INTRODUCTION

The newly emergent fields of quantum sensing and imaging utilize quantum entanglement—the same subtle effects exploited in quantum information processing—to push the capability of precision measurements and image construction using interferometers to the ultimate quantum limit of resolution.^{1,2} Alan Migdall at the U.S. National Institute of Standards and Technology, for example, has proposed and implemented a quantum optical technique for calibrating the efficiency of photodetectors using the temporal correlations of entangled photon pairs.³ It was one of the first practical applications of quantum optics to optical metrology, and has produced a technique to calibrate detectors without the need for an absolute standard.

These quantum effects can also be applied to increase the signal-to-noise ratio in an array of sensors from Laser Interferometer Gravitational Wave Observatory (LIGO) to Laser Light Detection and Ranging (LIDAR) systems, and to synchronized atomic clocks. Quantum imaging exploits similar quantum ideas to beat the Rayleigh diffraction limit in resolution of an imaging system, such as used in optical lithography. We present an introduction to these exciting fields and their recent development.

Entanglement is the most profound property of quantum mechanical systems.^{4,5} First we need to define entanglement. For simplicity, let us consider a system of two modes A and B only. Modes A and B may describe the two spatial paths of a Mach-Zehnder interferometer or two different polarization modes in an optical cavity. We can put a photon in either of the modes or let them remain empty. Let us suppose a general state of mode A which is a superposition of the two possible states, therefore we obtain $\alpha|0\rangle_A + \beta|1\rangle_A$, where α and β may be complex and $|\alpha|^2 + |\beta|^2 = 1$ is required for proper normalization of the state. We also write a superposition for a general state for mode B, that is, $\gamma|0\rangle_B + \delta|1\rangle_B$ and $|\gamma|^2 + |\delta|^2 = 1$. Now consider the combined two-mode state $|\Psi\rangle = (|1\rangle_A |0\rangle_B + |0\rangle_A |1\rangle_B) / \sqrt{2}$, where either a photon is in mode A or B. It is easy to see that this state cannot be decomposed into a product state for mode A and mode B only, that is, we *cannot* find any coefficients α , β , γ , δ such that the equation $(\alpha|0\rangle_A + \beta|1\rangle_A)(\gamma|0\rangle_B + \delta|1\rangle_B) = (|1\rangle_A |0\rangle_B + |0\rangle_A |1\rangle_B) / \sqrt{2}$ is satisfied. The state $|\Psi\rangle$ is an example for a nonseparable state. In general, any nonseparable state of two or more systems is called entangled. Erwin Schrödinger was the first who coined the term entanglement,⁶ although he is by far more prominent for his Schrödinger cat. Now we have a proper definition for entangled

pure states. But what is the definition for entangled mixed states? Let us suppose two systems A and B that can be in n different mixed states ρ_i^A and ρ_i^B with $i=1, \dots, n$. A separable mixed state may be written as $\rho = \sum_i p_i \rho_i^A \otimes \rho_i^B$, where the p_i are probabilities. Mathematicians call this particular sum a convex combination of product states. Now we can define entanglement for mixed states in excluding just these separable states. We say, if the bipartite system *cannot* be written in the above way, we call it entangled.⁷

Entanglement is not necessarily a useful physical quantity. Entanglement is usually discussed together with nonlocal correlations. These nonlocal correlations can be verified by Bell experiments.⁸ For a nice review of the current status of Bell experiments see, for example, Ref. 9. The violation of a Bell inequality^{10,11} by a specific quantum state is an indication that the state is able to exhibit nonlocal correlations. It is known that for any entangled *pure* state of any number of quantum systems one may violate a generalized Bell inequality.^{12,13} An extension of this statement for mixed entangled states has not been found. Furthermore Werner provided in 1989 an example of nonseparable mixed states that do not violate a Bell inequality.⁷ This demonstrated that the class of entangled mixed states decomposes into states that are entangled but do not show nonlocal correlations and those that are entangled and are nonlocal. The state $|\Psi\rangle$ is clearly entangled but only until very recently it has been proven in several experiments that this state violates a Bell inequality.¹⁴⁻¹⁷ Generalizations of this state where the single photon is replaced by N photons will play an important role in applications described later in this chapter. That also this generalization for N photons in either of the modes violates a Bell inequality has been proven in Ref. 18 very recently. This shows that the connection of entanglement and nonlocal correlations is still a very hot research topic. Now let us come back to a very practical application, an optical interferometer.

In order to see the role of quantum mechanics in optical interferometers, we first consider a prototype Mach-Zehnder interferometer (MZI).¹⁹ The task is to measure the path-length difference between the two arms of the interferometer. In the standard approach, as shown in Fig. 1, an input laser beam is launched into the first 50-50 beam splitter (BS) on the left in port A, bounced off of the two mirrors in the middle, and recombined at the second 50-50 BS on the right. Light then emerges from the top and bottom ports, C and D, of the second BS, and is then made incident on two photodetectors \mathcal{C} and \mathcal{D} , as shown. Typically, the intensities in each port, I_C and I_D , are measured at each detector and the result is combined to yield the difference intensity, $I = I_D - I_C$, which

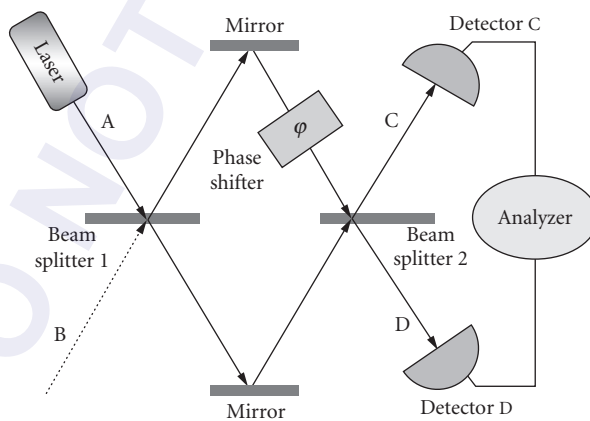


FIGURE 1 Schematic of a Mach-Zehnder interferometer. Laser light in port A is split by the first 50-50 beam splitter, acquires a phase difference, and recombines at the second beam splitter and emerges in ports C and D. We take a convention such that, for a balanced interferometer, port C is the dark port. Hence, any light emergent here is indicative of an arm displacement and can be detected by the two detectors and the analyzer.

we shall call the signal. To indicate the phase induced by the path difference between the upper and lower arms, we place an icon for a phase shift φ , which in this example has the value $\varphi = kx$. Here x is the path-length difference between the two arms. The wave number, $k = 2\pi/\lambda$, is a predetermined constant, given the optical wavelength λ .

We adopt the convention that the light field always picks up a $\pi/2$ phase shift upon reflection off of a mirror or off of a BS, and also no phase shift upon transmission through a BS. Then, the two light fields emerging from the second BS out the upper port C are precisely π out of phase with each other, and hence completely cancel out due to destructive interference (the dark port). Consequently, the two light fields recombine completely in phase as they emerge from the lower port D and add up due to constructive interference (the bright port). Hence for a balanced MZI all of the energy that comes in port A emerges out of port D and none out of port C. Clearly, any change in the path difference x away from the $x = 0$ balanced condition will cause light to appear in the formerly dark port, and in this way we can measure x by simply measuring intensities at the detectors.

The question is: How precise a measurement of the path difference x can we make? If the light intensity incident on port A is I_A , then in terms of the phase shift φ the output-port intensities can be written as

$$I_C = I_A \sin^2(\varphi/2) \quad (1a)$$

$$I_D = I_A \cos^2(\varphi/2) \quad (1b)$$

It is typical for the analyzer in Fig. 1 to compute the difference intensity $M = I_D - I_C$ (where M stands for “minus”) such that

$$M(\varphi) \equiv I_D - I_C = I_A \cos(\varphi) \quad (2)$$

Since $\varphi/2 = kx/2 = \pi x/\lambda$, we have that $I_C = 0$ and $I_D = I_A$ whenever $x/\lambda = 0, 1, 2, 3, \dots$. Hence, our ruler is the light wave itself and the tick marks are spaced the wavelength λ apart. We may start with a balanced interferometer with equal arm lengths, $x = 0$ (and $M = I_A$), and then slowly move the upper mirror upward increasing x . As we do we will break the balance and begin to see light emerging from the formally dark port C (M decreases in the plot).

At the point $\varphi = \pi/2$, when $I_C = I_D$, then $M = 0$. Eventually, we will see port C attain maximum brightness and port D will go dark ($M = -I_A$). As we continue the mirror displacement, this process will reverse, as sine and cosine are periodic, and finally port C will go dark again (M is maximum again with $M = I_A$). At this point, we can stop moving the upper mirror and we are assured that now the path difference x has gone from 0 to λ . If we take $\lambda = 1.0 \mu\text{m}$, then it would seem we have a machine capable of measuring distances to an accuracy of about $\lambda = 1.0 \mu\text{m}$. This is consistent with the Rayleigh diffraction limit, typically invoked in classical optics.

Let us now balance the interferometer such that we start at the point $\varphi = \pi/2$, when $I_C = I_D$, and hence $M = 0$ in Fig. 2. Note this is where the curve crosses the horizontal axis and the slope of the M -curve is steepest. If we call the horizontal displacement change $\Delta\varphi$, then we can see this is related to the vertical intensity change ΔM . For small changes, we may approximate this relation using differentials, that is, $\Delta M = I_A \Delta\varphi$ or

$$\frac{\Delta M}{\Delta\varphi} = \frac{\partial M}{\partial\varphi} = I_A \sin(\varphi) \quad (3)$$

which may be written as

$$\Delta\varphi = \frac{\Delta M}{\partial M/\partial\varphi} = \frac{\Delta M}{I_A \sin(\varphi)} \quad (4)$$

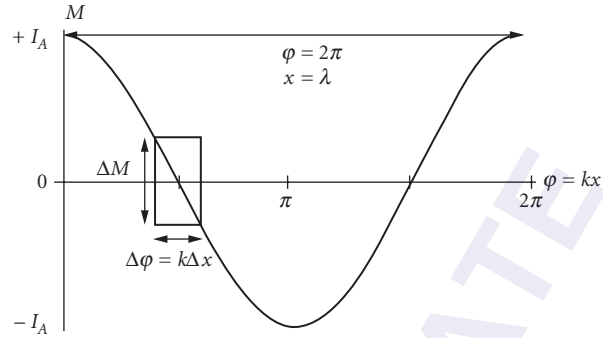


FIGURE 2 Typical Mach-Zehnder analyzer. The difference intensity M is plotted as a function of the phase shift $\varphi = kx$, where x is the arm displacement to be measured. The minimal detectable displacement, Δx , is limited by the fluctuations in the optical intensity, ΔM . These fluctuations are quantum mechanical in nature.

The quantity $\partial M/\partial\varphi$ is the slope of the curve, which is largest at the crossing point, implying our minimum detectable phase $\Delta\varphi$ is smallest there, via Eq. (4). At the crossing point $\varphi = \pi/2$ and $\sin(\pi/2) = 1$, and so this relation would seem to indicate that *if* we can measure the intensity displacement ΔM with infinite precision ($\Delta M = 0$), we can measure the phase (and hence distance) with infinite precision ($\Delta\varphi = 0$). In practice, various technical imperfections tend to set the limit for the finite precision ΔM .²⁰ However, as far as classical electromagnetic waves are concerned, nothing “fundamentally” prevents ΔM being zero. Hence, it would appear that *if* we try hard, we could detect any amount of phase shift no matter how small it is.

23.2 SHOT-NOISE LIMIT

The problem is that the simple classical arguments we used above do not take into account the effects of quantum mechanics. Specifically it does not take into account the fact that the intensity of the light field is not a constant, which can be measured with infinite precision, but that it fluctuates about some average value, and those fluctuations have their origin in the vacuum fluctuations of the quantized electromagnetic field.^{21,22} According to quantum mechanics, optical intensity can never be measured with infinite precision. Hence the uncertainty, in the curve of Fig. 2, always has some finite value, indicated by the box of height ΔM . The intensity displacement M can never be measured with infinite precision and has a fundamental uncertainty ΔM , and therefore the consequent phase φ will always have its related uncertainty $\Delta\varphi$, which is the width of the box. These fundamental quantum intensity fluctuations suggest that there is a Heisenberg uncertainty principle at work, which in our example implies that the intensity I and the phase φ cannot both simultaneously be measured with infinite precision.

For a quantum analysis of this phenomenon, we introduce the mean number of photons in the laser field as the dimensionless quantity n , and note that the intensity I is then proportional to n for a steady-state system. If we denote the fluctuation in the phase as $\Delta\varphi$ and that in the intensity as Δn , we can then write down the Heisenberg number-phase uncertainty relation as^{23–25}

$$\Delta n \Delta\varphi \geq 1 \quad (5)$$

This is closely related to the better known energy-time uncertainty principle $\Delta E \Delta t \geq \hbar$, where ΔE is the uncertainty in the energy, Δt is the uncertainty in the time, and \hbar is Dirac’s constant (Planck’s constant

divided by 2π). For a standing, monochromatic, electromagnetic wave we have $E = \hbar n \omega$, where ω is the frequency. This is just the energy per photon multiplied by the average number of photons. Since there is no propagation for a standing wave we have $\varphi = \omega t$ as the accumulated phase at any point. Approximating both of these expressions with differentials gives $\Delta E = \hbar \Delta n \omega$ and $\Delta \varphi = \omega \Delta t$. Inserting these two expressions into the energy-time uncertainty relation yields the number-phase relation, Eq. (5).

For a laser beam, the quantum light field is well approximated by a coherent state, denoted as $|\alpha\rangle$, where the complex number $\alpha = |\alpha|e^{i\varphi}$ is proportional to the electric field amplitude E such that $|\alpha|^2 = n$, the latter of which we recall is the dimensionless field intensity.²¹ This is the dimensionless quantum version of the classical relation $|E|^2 = I = I_0 n$. The full dimensional form is $E = E_0 \sqrt{n}$, where $I_0 = |E_0|^2 = \hbar \omega / (\epsilon_0 V)$, which in SI units, \hbar is Dirac's constant, ϵ_0 is the free-space permittivity, and V is the mode volume for the electromagnetic field. Hence I_0 is the intensity of a single photon. The fluctuations are typically represented in a phasor diagram, as shown in Fig. 3. Here the phase is the polar angle φ measured counterclockwise off the horizontal axis. The radius from the origin to the center of the coherent-state disk is $R = |\alpha|^2 = n$. The diameter of the disk d is on the order of $d = \Delta n = \sqrt{n}$. From simple geometry, we can then approximate $d = R \Delta \varphi$, where $\Delta \varphi$ is the uncertainty or fluctuation in the angular φ direction.

Combining all this we arrive at the fundamental relationships between number (intensity) and phase uncertainty for a coherent-state laser beam,

$$\Delta n \Delta \varphi = 1 \quad (6a)$$

$$\Delta n = \sqrt{n} \quad (6b)$$

$$\Delta \varphi_{\text{SNL}} = \frac{1}{\Delta n} = \frac{1}{\sqrt{n}} \quad (6c)$$

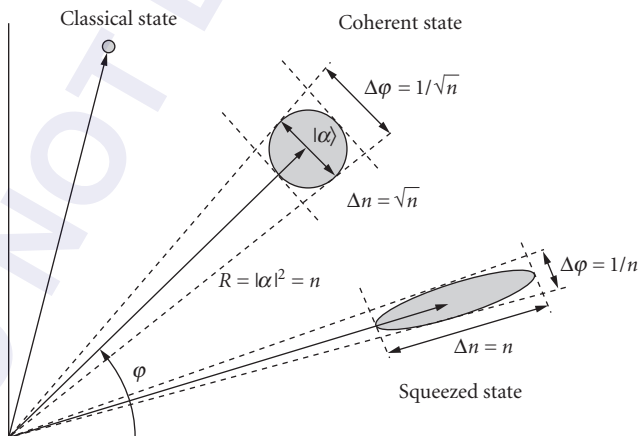


FIGURE 3 Phase-space diagram showing quantum fluctuations. Fluctuations in the radial direction correspond to intensity and those in the angular direction phase. A coherent state is a disk and has fluctuations equal in intensity and phase (a “true” classical state is a point and has no fluctuations). Also shown is a phase squeezed state, which has fluctuations decreased in the angular (phase) direction, at the expense of increase fluctuation in the radial (intensity) direction. Such a phase-squeezed state can be used to beat the shot-noise limit.

The first relation, Eq. (6a), tells us that we have equality in Eq. (5); that is a coherent state in a minimum uncertainty state (MUS). Such a state saturates the Heisenberg number-phase uncertainty relation with equality. This is the best you can do according to the laws of quantum mechanics. The second relation, Eq. (6b), describes the fact that the number fluctuations are poissonian with a mean of n and a deviation of $\Delta n = \sqrt{n}$, a well-known property of the Poisson distribution and the consequent number statistics for coherent-state laser beams.²² Putting back the dimensions we arrive at

$$\Delta\varphi_{\text{SNL}} = \sqrt{\frac{I_0}{I}} \quad (7)$$

which is called the shot-noise limit (SNL). The term “shot noise” comes from the notion that the photon-number fluctuations arise from the scatter in arrival times of the photons at the beam splitter, much like buckshot from a shotgun ricocheting off a metal plate. We can also import the SNL into our classical analysis above. Consider Eq. (4), where we now take $I_A = I_0 n$, $\Delta M = \sqrt{I_0 n}$, and $\varphi = \pi/2$. We again recover Eq. (6c) for the phase uncertainty. Hence quantum mechanics puts a quantitative limit on the uncertainty of the optical intensity, and that intensity reflects itself in a consequent quantitative uncertainty of the phase measurement.

In classical electromagnetism, we can also represent a monochromatic plane wave on the phasor diagram of Fig. 3—but instead of a disk the classical field is depicted as a point. The radial vector to the point is proportional to the electric field amplitude E and the phase angle corresponds to the classical phase of the field. The phase-space point represents the idea that, classically, we can measure number and phase simultaneously and with infinite precision. As we have seen above, quantum mechanically this is not so. The Heisenberg uncertainty principle (HUP) of Eq. (5) tells us that both phase and intensity cannot be measured simultaneously with infinite precision. For a minimum uncertainty state (MUS), such as a coherent state $|\alpha\rangle$, we have equality in the HUP, as given in Eq. (6a). Then, combined with the poissonian-statistical distribution of photon number for a coherent state, Eq. (6b), we arrive at the shot-noise limit.

23.3 HEISENBERG LIMIT

In 1981, Carlton Caves first proposed the idea of using nonclassical states of light—the so-called squeezed states—to improve the sensitivity of optical interferometers to below even the shot-noise limit.²⁶ This notion came as somewhat of a surprise to the interferometer community, as it was thought at the time that the shot-noise limit was the ultimate limit on sensitivity as imposed by quantum mechanics. However, there are other minimum uncertainty states besides the coherent state. The easiest way to see this is to look again at the representation of the coherent state as a disk in phase space (Fig. 3). The fact that it is a disk indicates that the fluctuations are the same in all directions, and that the area of the disk is a constant \mathcal{A} . The pictogram and the HUP then tells us that any quantum state must have an area greater or equal to \mathcal{A} , and that the MUS has an area equal to \mathcal{A} . This is, for a coherent state, equivalent to stating the three conditions of Eq. (6). However, we can relax Eq. (6b) and (6c), while still maintaining the HUP of Eq. (6a).

That is, we can decrease $\Delta\varphi$, at the expense of increasing Δn at the same time, so that the product $\Delta\varphi\Delta n = 1$ remains constant and the area of the disk remains the same value \mathcal{A} . Pictorially this amounts to squeezing the coherent-state disk in the angular direction, while allowing it to expand in the radial direction, as shown in Fig. 3. The important point is that the area \mathcal{A} of the ellipse remains unchanged so that the HUP is obeyed. However, we can decrease phase uncertainty at the expense of increasing the number uncertainty. Furthermore, it is possible to produce such squeezed states of light in the laboratory, using nonlinear optical devices and ordinary lasers.^{27–30}

Now the question is: What is the most uncertainty we can produce in photon number, given that the mean photon number n is a fixed constant, and that we still want to maintain the MUS condition—that the area of the ellipse remains a constant \mathcal{A} . Intuitively one cannot easily imagine a scenario where

the fluctuations in the energy, $\Delta E = \hbar \omega \Delta n$, exceeds the total energy of the laser beam, $E = \hbar \omega n$. Hence the best we can hope to achieve is $\Delta E = E$ or, canceling out some constants, $\Delta n = n$. Inserting this expression in the HUP of Eq. (6a), we obtain what is called the Heisenberg limit:

$$\Delta \phi_{\text{HL}} = \frac{1}{n} \quad (8)$$

Putting back the dimensions we get

$$\Delta \phi_{\text{HL}} = \frac{I_0}{I} \quad (9)$$

This is exactly the limit one gets with a rigorous derivation using squeezed light in the limit of infinite squeezing.^{31,32} It is called the Heisenberg limit as it saturates the number-phase HUP, and also because it can be proven that this is the best you can do in a passive interferometer with finite average photon number n . Converting to minimum detectable displacement we get

$$\Delta x = \frac{\lambda}{n} = \lambda \frac{I_0}{I} \quad (10)$$

where I_0 is the single photon intensity, defined above.

So far, we have considered the situation that we send light in port A and analyzed what came out of ports C and D for the MZI shown in Fig. 1. What about input port B? Classically there is no light coming in port B, and hence it is irrelevant. But, it is not so. In his 1981 paper, Caves showed that no matter what state of the photon field you put in port A, so long as you put nothing (quantum vacuum) in port B, you will always recover the SNL. In quantum electrodynamics, even an interferometer mode with no photons in it experiences electric field fluctuations in that mode.

In the MZI these vacuum fluctuations have another important effect; at the first BS they enter through port B and mix with whatever is coming in port A to give the SNL in overall sensitivity. It becomes clear then, from this result, that the next thing to try would be to plug that unused port B with something besides vacuum. It was Caves' idea to plug the unused port B with squeezed light (squeezed vacuum to be exact). That, with coherent laser light in port A as before—and in the limit of infinite squeezing—then the SNL rolls over into the HL.

In the laboratory, however, infinite squeezing is awfully hard to come by. With current technology,^{33–35} the expected situation is to sit somewhere between SNL and HL but a lot closer to the former than the latter. Recent analyses by a Caltech group, on exploiting squeezed light in LIGO, indicates a potential for about a one-order-of-magnitude improvement in a future LIGO upgrade.³⁶ Not the twelve orders of magnitude that was advertised above, but enough to allow the observatory to sample about 80 times the original volume of space for gravitational-wave sources. That, for LIGO, is a big deal.

23.4 “DIGITAL” APPROACHES

Squeezing is an “analog” approach to quantum optical interferometry, in that the average photon number and the degree of squeezing are continuous variables. There is another approach, exploiting discrete photon number and path-entangled optical states, where the photon number is fixed. There is a large body of literature on using entangled particles or photons in a Mach-Zehnder interferometer in order to beat the shot-noise limit. The first such proposal was by Bernard Yurke in the context of neutron interferometry.³⁷

In 1986, the use of photon-number eigenstates, the so-called Fock states, to improve the interferometer phase sensitivity by Yurke et al.,³⁸ and independently by Yuen.³⁹ They showed that with a

suitably correlated input state the phase sensitivity can be improved to the one proportional to $1/N$. The correlated input state—which we call the Yurke state—is given by

$$|\psi\rangle_Y = \left| \frac{N}{2}, \frac{N}{2} \right\rangle + \left| \frac{N}{2} + 1, \frac{N}{2} - 1 \right\rangle \quad (11)$$

where the normalization constant of $1/\sqrt{2}$ has been dropped for convenience. The two numbers inside the brackets represent the number of photons entering the input ports A and B, respectively. As is the case of the squeezed state, the Fock-state approaches also require the light field incident upon *both* input ports of the MZI. If N photons entered into each input port of the interferometer in nearly equal numbers, then it is possible to obtain the asymptotic phase sensitivity scaling of order $1/N$ for large N , the Heisenberg-limit.⁴⁰

In 1993, Hillery and Mlodinow showed that the $SU(2)$ -squeezed minimum-uncertainty states can also be used for the input state of the Mach-Zehnder interferometer to achieve Heisenberg-limited sensitivity.⁴¹ They called such a state the “intelligent” state. The intelligent state has a fixed number of photons. However, it consists of all the possible combination of the photon number distribution in the input modes such that

$$|\psi\rangle_{INT} = \sum_{k=0}^N C_{N,k} |N-k, k\rangle \quad (12)$$

where the probability amplitude $C_{N,k}$ varies with N and k as well as the degree of squeezing.⁴² Although it contains a fixed number of photons, the exact form of the intelligent state depends on the degree of squeezing. Hence, we may put this approach in between the analog and the digital ones.

In 1993, on the other hand, Holland and Burnett proposed the use of the so-called dual-Fock states for the two input ports of the MZI in order to achieve Heisenberg-limited sensitivity.⁴³ The dual-Fock state has the same number of photons at each input mode such that the quantum state can be written as

$$|\psi\rangle_{DF} = \left| \frac{N}{2}, \frac{N}{2} \right\rangle \quad (13)$$

We see that the dual-Fock state has a much simpler form, compared to the Yurke state of Eq. (11). Why is it that this simple looking dual-Fock state had been overlooked in the earlier investigation? Using the dual-Fock state input, the measurement of the difference intensity does not yield any signal. In other words, the difference between the two output modes is always zero regardless of the amount of phase shift. This prevents the use of dual-Fock state for the interferometer input when the difference signal is measured, and because of that it was discarded, not overlooked. What Holland and Burnett suggested is to construct a probability distribution of the estimated phase, conditioned on the number of photons at each output port. Then, they showed, after sufficient number of trials, the probability distribution becomes narrow such that the phase uncertainty approaches to the $1/N$ Heisenberg limit. Heisenberg-limited phase measurements have also been proposed for Ramsey-type atom interferometry, where the degenerate Bose-Einstein condensates play the role of atomic dual-Fock states.^{44,45}

The conditional probability distributions are generally used for parameter estimation. The phase is not a quantum mechanical observable (that is represented by a self-adjoint operator),⁴⁶ and it has to be estimated. Here, instead of direct inversion for the phase shift as performed with the difference-intensity measurement in Eq. (2), the probability distribution is constructed to infer the phase shift given the measurement results (photon counts at C and D ports separately).⁴⁷ The error in the phase estimation is then given by the variance of that distribution.

Such phase estimation protocols with conditional probabilities are applied to look for optimal quantum states of the input light field as well as the optimal output measurement schemes.^{48–52} Given that the output detection does not have to be the conventional difference intensity measurement,

Sanders and Milburn proposed an ideal canonical measurement based on phase state projection.⁴⁸ Using such an ideal measurement strategy, Berry and Wiseman suggested an optimal state that minimizes the so-called Holevo phase variance.⁴⁹ They further developed an adaptive measurement scheme that approximates the ideal measurement scheme by Sanders and Milburn. With other measurement strategies such as coincident measurement that exploits the fourth-order correlation, it is also possible to obtain the phase sensitivity beyond the shot-noise limit with the dual-Fock state.^{53–57}

23.5 NOON STATE

Let us now consider a two-mode, path-entangled, photon-number state, commonly called the NOON state. The idea is that we have a fixed finite number of photons N that are either all in the upper mode A or all in the lower mode B, but we cannot tell—even in principle—which is which. The state of all up and none down is written $|\text{up}\rangle = |N\rangle_A |0\rangle_B$ and the state of all down and none up is similarly $|\text{down}\rangle = |0\rangle_A |N\rangle_B$. The notation indicates a product state of N photons either in A or B (but not both) such that

$$|\text{N00N}\rangle \equiv |\text{up}\rangle + |\text{down}\rangle = |N, 0\rangle + |0, N\rangle \quad (14)$$

where a normalization constant of $1/\sqrt{2}$ has again been dropped for convenience. The NOON state was first discussed in 1989 by Barry Sanders, who was particularly interested in the Schrödinger-cat aspect and how that affected quantum decoherence.⁵⁸ It was rediscovered in the context of quantum imaging—particularly for quantum lithography⁵⁹—to circumvent the Rayleigh diffraction limit. The NOON state has the interesting property that it is quantum entangled between the two modes and rigorously violates what is known as a Bell inequality for nonclassical correlations.¹⁸

To see why a NOON state has the magical properties—super sensitivity and super resolution, in particular—we take a brief look at the difference in behavior between a number state $|N\rangle$ and a coherent state $|\alpha\rangle$ in an MZI. When a coherent state passes through a phase shifter φ , such as depicted in Fig. 1, it picks up a phase of φ . This is a property of a classical monochromatic light beam that coherent states inherit quantum mechanically. However, number states are already highly nonclassical states to begin with. Their behavior in the phase shifter is radically different.

When a monochromatic beam of number states passes through a phase shifter, the phase shift is directly proportional to N , the number of photons. There is no n -dependence in the phase shift for the coherent state (n is the average number of photons). In terms of a unitary evolution of the state, the evolution for any photon state passing through a phase shifter φ is governed by

$$\hat{U}(\varphi) \equiv \exp(i\varphi\hat{n}) \quad (15)$$

where \hat{n} is the photon number operator. The phase shift operator can be shown to have the following two different effects on coherent versus number states,²²

$$\hat{U}_\varphi |\alpha\rangle = |e^{i\varphi}\alpha\rangle \quad (16a)$$

$$\hat{U}_\varphi |N\rangle = e^{iN\varphi} |N\rangle \quad (16b)$$

Notice that the phase shift for the coherent state is independent of number, but that there is an N dependence in the exponential for the number state. The number state then evolves in phase N times more rapidly than the coherent state. After the phase shifter the NOON state evolves into

$$|N, 0\rangle + |0, N\rangle \rightarrow e^{iN\varphi} |N, 0\rangle + |0, N\rangle \quad (17)$$

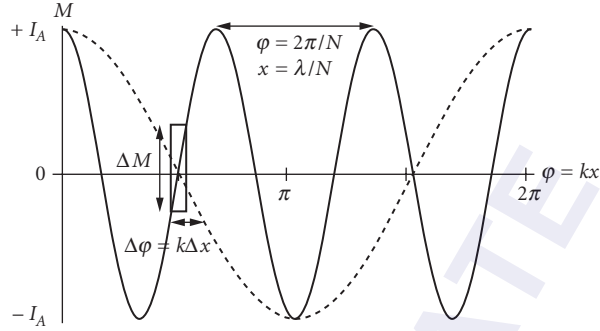


FIGURE 4 Comparison of the detection signal of a coherent state (dotted) and a N00N state (solid). The N00N state signal oscillates N times as fast as the coherent state and has maximum slope that is N times as steep ($N = 3$). The effect is as if the N00N state was composed of photons with an effective wavelength of λ/N instead of λ .

which is the origin of the quantum improvement phase sensitivity. If we now carry out an N -photon detecting analyzer (still different from the conventional difference intensity measurement), we obtain

$$M_{\text{N00N}}(\phi) = I_A \cos(N\phi) \quad (18)$$

which is the solid curve in Fig. 4. The N00N signal (solid) oscillates N times as fast as the coherent state (dotted). Two things are immediately clear. The distance between peaks goes from $\lambda \rightarrow \lambda/N$, which is the quantum lithography effect—we now beat the Rayleigh diffraction limit of λ by a factor of N . This sub-Rayleigh-diffraction-limit effect is now commonly called “super resolution.” The slope of the curve at the horizontal axis crossing point gets larger as well, also by a factor of N . Now the minimal detectable phase, given by Eq. (4), consequently goes down. However, the signal M for this N00N state is not the same as for the coherent state scheme, as we are now counting photons N at a time. And it turns out then that $\Delta M_{\text{N00N}} = 1$ for the new scheme, and then Eq. (4) gives

$$\Delta\phi_{\text{N00N}} = 1/N \quad (19)$$

which is precisely the Heisenberg limit of Eq. (8). This Heisenberg limit, or the beating of the shot-noise limit, is now commonly called “super sensitivity.”

For $N = 1$ and $N = 2$ (low N00N) it is fairly straightforward to make such states with nonclassical sources of photon number states of either the form $|1\rangle_A|0\rangle_B$ or $|1\rangle_A|1\rangle_B$, that is one photon in mode A and none in B, or one photon in each of modes A and B. The standard approach utilizes spontaneous parametric down-conversion (SPDC), where an ultraviolet (UV) photon is down converted into a pair of number states. The effect of a simple beam splitter transformation on these states²² is to convert them to low-N00N states, as follows:

$$|1\rangle_A|0\rangle_B \xrightarrow{\text{BS}} |1\rangle_{A'}|0\rangle_{B'} + |0\rangle_{A'}|1\rangle_{B'} \quad (20a)$$

$$|1\rangle_A|1\rangle_B \xrightarrow{\text{BS}} |2\rangle_{A'}|0\rangle_{B'} + |0\rangle_{A'}|2\rangle_{B'} \quad (20b)$$

where Eq. (20a) shows that a single photon cannot be split in two, and Eq. (20b) is illustrative of the more subtle Hong-Ou-Mandel effect—if two single photons are incident on a 50-50 beam splitter they will “stick” and both photons will go one way or both will go the other way, but you never get one photon out each port.⁶⁰

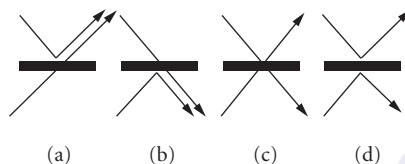


FIGURE 5 Four possibilities of the output, when sending a $|1\rangle_A|1\rangle_B$ state through a beam splitter. The diagrams (c) and (d) lead to the same final state, interfering destructively; (c) no phase acquired for transmission-transmission; and (d) a total of π phase shift for reflection-reflection.

As depicted in Fig. 5, it is the probability amplitude for the transition $|1\rangle_A|1\rangle_B \xrightarrow{\text{BS}} |1\rangle_{A'}|1\rangle_{B'}$ that completely cancels out due to destructive interference. On the other hand, the probability amplitude for the transition indicated by Eq. (20b) adds up, due to constructive interference. So it is relatively easy, once you have a source of single photons, to create low-N00N states. The challenge is then, how to go to high N00N?

One of the first proposals for making high-N00N states was introduced in 2002 by Gerry and Campos, motivated by their applications in lithography and metrology.⁶¹ The idea is to make a kind of quantum computing gate (a Fredkin gate) so that a single photon in the upper MZI-1 controls the phase shift in the lower MZI-2, as shown in Fig. 6.⁶² A nonlinear optical material called a cross-Kerr phase shifter couples the two MZIs.

The strength of the cross-Kerr effect, however, is so tiny at the single photon level⁶³⁻⁶⁵ that is it comes up 20 orders of magnitude too small for making a good N00N-state generator. There are two

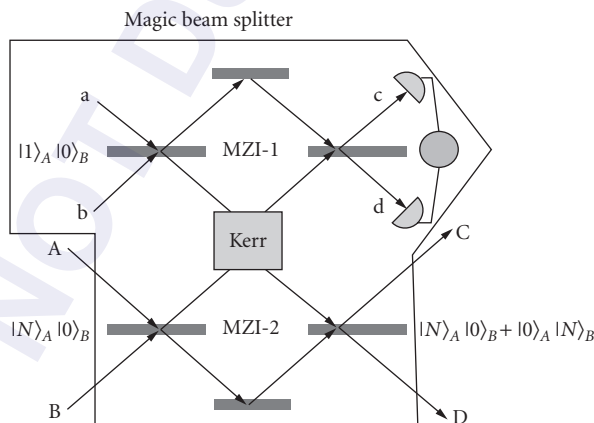


FIGURE 6 N00N-state generator by Gerry and Campos.⁶¹ In the lower interferometer if N photons enter upper port A, they will always emerge in lower port D for a balanced Mach-Zehnder. Now if the lower interferometer is coupled to the upper one, via a strong cross-Kerr nonlinearity, a single photon in the lower branch of the upper interferometer causes a π phase shift, directing all N photons to emerge out port C. If the upper device is also an interferometer, one can arrange a superposition of zero and one photons in the Kerr box, giving rise to a superposition of a 0 and π phase shift. This Kerr superposition results in the N00N state, the number state superposition in the output modes C and D of the lower interferometer.

well-known ways to boost the Kerr effect: put in an optical micro-cavity around the atoms in the Kerr material,⁶⁶ or coherently lock the atoms together in an approach known as electromagnetically induced transparency.^{67–69} Both roads have their complication and technical challenges. There is, however, a third path.

Here is where the universe of quantum-optical-metrology took some hints from the world of quantum-optical computing. The Gerry-Campos idea is based on an optical Fredkin gate, a sort of a single photon transistor. If such a device could be made easily, it would be a quick and easy road to the all-optical quantum computer. In 2001, Knill, LaFlamme, and Milburn (KLM) proposed an all-optical, digital, scheme for quantum computing that exploited discrete entangled photon states distributed over the modes of an optical interferometer.⁷⁰ This discovery ignited a huge international collaborative research and development program on the optical quantum computer. In this scheme, the Kerr nonlinearity is replaced with additional ancillary, mirrors, phase shifters, beam splitters and—most importantly—detectors. The KLM idea is that the detection process in the ancillary devices induces an effective Kerr nonlinearity.⁷¹

While still not perfect, the effective Kerr produced this way can be much stronger than the usual Kerr materials. So instead of working one time in 10^{20} , our new device works one time in 10, which is a 19-order-of-magnitude improvement. Our group, then at the NASA Jet Propulsion Laboratory, proposed the first high-N00N generation scheme based on all linear-optical devices.⁷² The idea is to make the upper and lower mirrors into additional beam splitters and to put detectors just past them. Such a N00N-generating scheme can be stacked to produce N00N states with N that is arbitrarily high.^{73,74} However, at least initially, all such schemes produced N00N states with a probability that scaled exponentially badly as N increased. The larger and larger you make N , the far more likely it is that you get only non-N00N states. Recently there has been a great deal of work on the development of schemes that produce N00N “efficiently.”^{75–77} So we are well along the digital quantum-computer-paved road to super resolution and super sensitivity—at least in theory.

The demonstration of the doubling of the effective wavelength was made in 1999,⁷⁸ and a proof-of-principle demonstration of the super-resolving application to sub-Rayleigh lithography was carried out by the group of Yanhua Shih in 2001.⁷⁹ Then in 2004 the $N = 2$ barrier was breached; the group of Steinberg demonstrated super resolution for $N = 3$, and, for $N = 4$, the group of Zeilinger demonstrated it as well.^{80–82} These 2004 experiments exploited the *effective* Kerr nonlinearities extracted from the photodetection process. While not “optical” metrology, there has been a lot of work on precision frequency measurements with trapped ions. In particular, the so-called maximally entangled states have been exploited to show super resolution and super sensitivity.^{83–89}

Then, in 2007, the group of Andrew White demonstrated $N = 6$ super resolution in a process that used classical photons with a nonlinear N -photon detection scheme.⁸⁶ While interesting, such an approach can never achieve super sensitivity due to its semiclassical nature, as the authors pointed out. The first demonstration of both super sensitivity and super resolution in a single experiment also appeared in late 2007 in a collaborative Japanese and UK experiment.⁸⁷ This was the first experiment to beat the shot-noise limit, using N00N states, with $N > 2$. However, the results still were not quite at the Heisenberg limit. It lies between the shot-noise limit and the Heisenberg limit.

Another Australian collaboration, of the groups of Geoffrey Pryde and Howard Wiseman, carried out an experiment that produced N single-photon states of the form of Eq. (20a), and recycled them through the MZI in a feedback-loop implementation of a quantum computing protocol known as the Kitaev phase estimation algorithm^{88–90} with the effective $N = 378$. The idea is that instead of making such a large N00N state, they make 378 passes through the interferometer, with feedback, using the low $N = 1$ N00N state. Since all N00N states are equally entangled and violate a Bell inequality for any nonzero value of N , the trade-off is that of a complicated N00N-state generating scheme with a less complicated single-photon detection scheme with some electronic feedback.¹⁸ Such a protocol is, surprisingly, easier to implement than the high-N00N approach, and arrives at the same performance in sensitivity scaling, the Heisenberg limit. This Australian experiment, once again, illustrates the close connection between quantum optical computing and quantum metrology, as it achieves super sensitivity by running a quantum computational algorithm. A quantum computer is nothing more than a (complicated) quantum sensor, and hence one can design a quantum sensor by exploiting ideas from quantum computing.

23.6 QUANTUM IMAGING

Quantum imaging is a new subfield of quantum optics that exploits quantum correlations, such as quantum entanglement of the electromagnetic field, in order to image objects with a resolution (or other imaging criteria) that is beyond what is possible in classical optics. Examples of quantum imaging are quantum ghost imaging, quantum lithography, and sub-Rayleigh imaging.^{91,92} In 2000, it was pointed out that N00N states had the capability to beat the Rayleigh diffraction limit by a factor of N . This super-resolution feature is due to the high-frequency oscillations of the N00N state in the interferometer, as illustrated in Fig. 4. For the quantum lithography application, the idea is to realize that if one has an N -photon absorbing material, used as a lithographic resist, then these high-frequency oscillations are written onto the material in real space and are not just a trace on an oscilloscope. Mathematically, the N -photon absorption and the N -photon detection process have a similar structure, that is,

$$\langle N00N | (\hat{a}^\dagger)^N (\hat{a})^N | N00N \rangle = 1 + \cos(N\varphi) \quad (21)$$

where \hat{a} and \hat{a}^\dagger are the mode annihilation and creation operators. From Fig. 4, we see in the solid curve this oscillates N times faster than if we were using single photons, or coherent light, as in the dotted curve. Recall that, for our MZI, we have $\varphi = kx = 2\pi x / \lambda$, where x is the displacement between the two arms. For lithography x is also the distance measured on the photographic plate or lithographic resist. If we compare the classical resolution to the N00N resolution we may write, $\varphi_{N00N} = N\varphi_{\text{classical}}$, which we can solve for

$$\lambda_{N00N} = \frac{\lambda_{\text{classical}}}{N} \quad (22)$$

Written this way, we can say the effective wavelength of the N photons bundled together N at a time into the N00N state is N times smaller than the classical wavelength. This is another way to understand the super-resolution effect. The N entangled photons conspire to behave as a single classical photon of a wavelength smaller by a factor of N .⁹³ Since the Rayleigh diffraction limit for lithography is couched in terms of the minimal resolvable distance $\Delta x = \lambda_{\text{classical}}$, then we have $\Delta x_{N00N} = \lambda_{N00N} = \lambda_{\text{classical}}/N$.

Another interesting application is so-called “ghost imaging.” This effect exploits the temporal and spatial correlations of photon pairs, also from spontaneous parametric down conversion, to image an object in one branch of the interferometer by looking at correlations in the coincidence counts of the photons.⁹⁴ There is no image in the single-photon counts in either arm, but only in the double photon counts in both arms. The image is in a sense stored nonlocally.

A particular application of this more general idea of quantum imaging has been seen in quantum coherence tomography.⁹⁵ In this experiment, they image a phase object placed in one arm of the interferometer, using entangled photons in an $N = 2$ N00N state. They see not only the factor of two improvement in resolving power, predicted by Eq. (22), but also as a bonus they get a dispersion cancellation in the imaging system due to frequency entanglement between the photons.

Current experiments on N00N states have used rather dim sources of entangled photons, from UV pumped $\chi^{(2)}$ crystals in a spontaneous parametric down conversion (SPDC) setup.⁹⁶ For bright sources of N00N states, one can turn to optical parametric amplifiers (OPA), which is the same setup as SPDC, but in which we crank up the pump power.⁹⁷ In this regime of high gain, the creation of entangled photon pairs of the form of Eq. (13) occurs, but we have many, many pairs and the output can be written

$$|\text{OPA}\rangle = \sum_{n=0}^{\infty} a_n |n\rangle_A |n\rangle_B \quad (23)$$

where the probability of a large twin-number state $|N\rangle_A |N\rangle_B$ is given by $|a_N|^2$, which can be quite large in the limit of high pump powers. Passing the OPA state through a 50-50 beam splitter gives the generalized Hong-Ou-Mandel effect term by term, so that we get,

$$|\text{OPA}\rangle \xrightarrow{\text{BS}} \sum_{n=0}^{\infty} \sum_{m=0}^n c_{nm} |2n-2m\rangle |2m\rangle \quad (24)$$

where again the coefficients c_{nm} can be quite large for high pump powers. Taking the term $n = 1$ we immediately get the $N = 2$ N00N state from the regular Hong-Ou-Mandel effect. For larger $n \geq 1$, we find that there is always a large N00N component along with the non-N00N. For an $N = 2$ absorber, the visibility of the $N = 2$ N00N oscillations was predicted to saturate at a visibility of 20 percent.^{98,99} This 20 percent visibility is more than enough to exploit for lithography and imaging, and has recently been measured in a recent experiment in the group of DeMartini in Rome,¹⁰⁰ in collaboration with our activity at Louisiana State University.

Current commercial lithography exploits extreme ultraviolet light of around 100 nm and plans are to go to x-ray in the future. The problem is that the lithography system for x-ray cannot use the same lenses, mirrors, and other imaging devices as did the optical system and so each reduction in wavelength involves a huge cost in technology and hardware investment. But what if we could etch 50-nm-sized features using 500-nm wavelength photons by exploiting quantum entanglement? This is the promise of quantum lithography. However, no real demonstration of quantum lithography has been had, so far, due to the N-photon resist problem.¹⁰¹

23.7 TOWARD QUANTUM REMOTE SENSING

Improvements in optical metrology and imaging have a natural application in the realm of optical remote sensing, such as in coherent optical laser interferometric radar (LIDAR) or in sensor miniaturization.^{102,103} On the other hand, the group of Lloyd at MIT proposed a quantum optical clock synchronization protocol that eliminates the timing jitter of optical pulses that are transmitted through a fluctuating atmosphere.¹⁰⁴ These atmospheric fluctuations are currently the limiting source of noise in the global positioning system. Of course, when the N00N states propagating over distances of kilometers through the atmosphere, then photon scattering and loss and other issues such as atmospheric turbulence become an issue that are not apparent in a table top quantum interferometry demonstration.

The primary issue associated with photon loss is that the visibility of the interference pattern decreases, and that of the N00N state pattern decreases more rapidly than that of the single photon or coherent state interferometer. Hence, when the loss is sufficiently high, the slope of the N00N oscillations in Fig. 4 decreases to the point that, as far as super sensitivity is concerned, we do worse with N00N states than with either single photons or coherent states.^{105,106} There are, however, Fock states that offer super resolution and sensitivity better than shot noise that degrade less quickly under loss than N00N states.¹⁰⁷ What we've termed M&M states, of the form

$$|m, m'\rangle \equiv |m, m'\rangle + |m', m\rangle \quad (25)$$

(where a normalization constant of $1/\sqrt{2}$ has again been dropped for convenience) contain a trade-off of being more resilient to environmental decoherence, but having a greater minimum phase uncertainty of $1/(m-m')$.

In quantum optics, the photon loss (as well as the detector inefficiency) is modeled by a series of beam splitters.²⁵ In doing so, we first need to enlarge the Hilbert space to include the modes that represent the scattered photons and then, after the scattering, trace out those modes.¹⁰⁷ Here, instead of going into details, we use a simple (not necessarily correct) phase-shifter model and give a brief estimation of the loss effect. Consider Eq. (16), describing how coherent states and number states

behave upon passing through a phase shifter. The photon loss might be incorporated in the phase shifter by making the substitution $\varphi \rightarrow \varphi + i\gamma$, where γ is the rate at which photons are absorbed. We see then the effect of this loss in Eq. (16) is to produce an exponential loss factor that depends on N , for number states

$$|\alpha\rangle \xrightarrow{\varphi+i\gamma} |e^{-\gamma} e^{i\varphi} \alpha\rangle \quad (26a)$$

$$|N\rangle \xrightarrow{\varphi+i\gamma} e^{-N\gamma} e^{iN\varphi} |N\rangle \quad (26b)$$

Typically, we have $\gamma = gL$, where g is the loss per unit length and L the distance traveled through the lossy medium. The exponential dependence of the loss in the coherent (classical) state of Eq. (26a) is called Beer's law for optical absorption. We see that for N -photon number states, Eq. (26b), we have a super-exponential behavior, or what we call super-Beer's law. It implies that the $N00N$ states are much more fragile in a lossy environment than a classical coherent state.

Of course, for the number states, such a simple phase-shifter model fails to describe the quantum states of the lesser number of photons, and we need to carry out the detailed analysis mentioned above. Will we then be able to overcome such fragility of the entangled states and achieve super sensitivity in remote sensing? We may have to seek an answer from the field of quantum computing—the tools to fight against decoherence such as quantum error correction or decoherence free subspace.^{108,109}

Due to its “digital” nature, the entangled-number-state approach has recently benefited tremendously from an influx of ideas and experimental techniques originally developed in the context of all-optical “digital” quantum computing.¹¹⁰ The idea is that an optical quantum computer is a giant optical quantum interferometer, where the quantum entanglement between photons is exploited to carry out mathematical calculations, which are impossible on any classical computer. However, in the proposed quantum interferometers the entanglement is exploited to make ultimate precise measurements not possible with any classical optical machine. The optical quantum computer can be turned into an optical quantum interferometric measuring device, and vice versa. Theoretical and experimental tricks, devised for the former, can be exploited in the latter. Since, for over the past 10 years, a lot of efforts have gone into the development of quantum computers, we are now able to leverage this research for quantum optical metrology, imaging, and sensing.

23.8 REFERENCES

1. H. Lee, P. Kok, and J. P. Dowling, “A Quantum Rosetta Stone for Interferometry,” *J. Mod. Opt.* **49**:2325 (2002).
2. V. Giovannetti, S. Lloyd, and L. Maccone, “Quantum-Enhanced Measurements: Beating the Standard Quantum Limit,” *Science* **306**:1330 (2004).
3. A. L. Migdall, “Absolute Quantum Efficiency Measurements Using Correlated Photons: Toward a Measurement Protocol,” *IEEE Trans. Instrum. Meas.* **50**:478 (2001).
4. A. Einstein, B. Podolsky, and N. Rosen, “Can Quantum-Mechanical Description of Physical Reality Be Considered Complete?” *Phys. Rev.* **47**:777 (1935).
5. N. Bohr, “Can Quantum-Mechanical Description of Physical Reality Be Considered Complete?” *Phys. Rev.* **48**:696 (1935).
6. E. Schrödinger, “Die gegenwärtige Situation in der Quantenmechanik,” *Die Naturwissenschaften* **2**:827 (1935).
7. R. F. Werner, “Quantum States with Einstein-Podolsky-Rosen Correlations Admitting a Hidden-Variable Model,” *Phys. Rev. A* **40**:4277 (1989).
8. J. S. Bell, “On the Einstein-Podolsky-Rosen Paradox,” *Physics* **1**:195 (1964).
9. M. Genovese, “Research on Hidden Variable Theories: A Review of Recent Progresses,” *Phys. Rep.* **413**:319 (2005).

10. J. S. Bell, "On the Problem of Hidden Variables in Quantum Mechanics," *Rev. Mod. Phys.* **38**:447 (1966).
11. J. F. Clauser, M. A. Horne, A. Shimony, and R. A. Holt, "Proposed Experiment to Test Local Hidden-Variable Theories," *Phys. Rev. Lett.* **23**:880 (1969).
12. N. Gisin and A. Peres, "Maximal Violation of Bell's Inequality for Arbitrary Large Spin," *Phys. Lett. A* **162**:15 (1992).
13. S. Popescu and D. Rohrlich, "Generic Quantum Nonlocality," *Phys. Lett. A* **166**:293 (1992).
14. S. J. van Enk, "Single-Particle Entanglement," *Phys. Rev. A* **72**:064306 (2005).
15. B. Hessmo, P. Usachev, H. Heydari, and G. Björk, "Experimental Demonstration of Single Photon Nonlocality," *Phys. Rev. Lett.* **92**:180401 (2004).
16. S. A. Babichev, J. Appel, and A. I. Lvovsky, "Homodyne Tomography Characterization and Nonlocality of a Dual-Mode Optical Qubit," *Phys. Rev. Lett.* **92**:193601 (2004).
17. M. D'Angelo, A. Zavatta, V. Parigi, and M. Bellini, "Tomographic Test of Bell's Inequality for a Time-Delocalized Single Photon," *Phys. Rev. A* **74**:052114 (2006).
18. C. F. Wildfeuer, A. P. Lund, and J. P. Dowling, "Strong Violations of Bell-Type Inequalities for Path-Entangled Number States," *Phys. Rev. A* **76**:052101 (2007).
19. M. Born and E. Wolf, *Principles of Optics*, 7th ed., Cambridge University Press, Cambridge (2002).
20. J. Fraden, *Handbook of Modern Sensors*, Springer-Verlag, Berlin (1996).
21. M. O. Scully and M. S. Zubairy, *Quantum Optics*, Cambridge University Press, Cambridge (1997).
22. C. C. Gerry and P. L. Knight, *Introductory Quantum Optics*, Cambridge University Press, Cambridge (2005).
23. P. A. M. Dirac, "The Quantum Theory of the Emission and Absorption of Radiation," *Proc. Roy. Soc. A* **114**:243 (1927).
24. W. Heitler, *The Quantum Theory of Radiation*, Oxford University Press, Oxford (1954).
25. R. Loudon, *The Quantum Theory of Light*, Oxford University Press, Oxford (2000).
26. C. M. Caves, "Quantum Mechanical Noise in an Interferometer," *Phys. Rev. D* **23**:1693 (1981).
27. D. Walls, "Squeezed States of Light," *Nature* **306**:141 (1983).
28. M. Xiao, L.-A. Wu, and H. J. Kimble, "Precision Measurement beyond the Shot-Noise Limit," *Phys. Rev. Lett.* **59**:278 (1987).
29. P. Grangier, R. E. Slusher, B. Yurke, and A. La Porta, "Squeezed-Light-Enhanced Polarization Interferometer," *Phys. Rev. Lett.* **59**:2153 (1987).
30. R. Loudon and P. L. Knight, "Squeezed Light," *J. Mod. Opt.* **34**:709 (1987).
31. A. S. Lane, S. L. Braunstein, and C. M. Caves, "Maximum-Likelihood Statistics of Multiple Quantum Phase Measurements," *Phys. Rev. A* **47**:1667 (1993).
32. Z. Y. Ou, "Complementarity and Fundamental Limit in Precision Phase Measurement," *Phys. Rev. Lett.* **77**:2352 (1996).
33. V. V. Dodonov, "Nonclassical States in Quantum Optics: A Squeezed Review of the First 75 Years," *J. Opt. B* **4**:R1 (2002).
34. T. Corbitt and N. Mavalvala, "Quantum Noise in Gravitational-Wave Interferometers," *J. Opt. B* **6**:S675 (2004).
35. H.-A. Bachor and T. C. Ralph, *A Guide to Experiments in Quantum Optics*, 2nd ed., Wiley-VCH, Berlin (2004).
36. H. J. Kimble, Y. Levin, A. B. Matsko, K. S. Thorne, and S. P. Vyatchanin, "Conversion of Conventional Gravitational-Wave Interferometers into Quantum Nondemolition Interferometers by Modifying Their Input and/or Output Optics," *Phys. Rev. D* **65**:022002 (2001).
37. B. Yurke, "Input States for Enhancement of Fermion Interferometer Sensitivity," *Phys. Rev. Lett.* **56**:1515 (1986).
38. B. Yurke, S. L. McCall, and J. R. Klauder, "SU(2) and SU(1,1) Interferometers," *Phys. Rev. A* **33**:4033 (1986).
39. H. P. Yuen, "Generation, Detection, and Application of High-Intensity Photon-Number-Eigenstate Fields," *Phys. Rev. Lett.* **56**:2176 (1986).
40. J. P. Dowling, "Correlated Input-Port, Matter-Wave Interferometer: Quantum-Noise Limits to the Atom-Laser Gyroscope," *Phys. Rev. A* **57**:4736 (1998).

41. M. Hillery and L. Mlodinow, "Interferometers and Minimum-Uncertainty States" *Phys. Rev. A* **48**:1548 (1993).
42. C. Brif and A. Mann, "Nonclassical Interferometry with Intelligent Light," *Phys. Rev. A* **54**:4505 (1996).
43. M. J. Holland and K. Burnett, "Interferometric Detection of Optical Phase Shifts at the Heisenberg Limit," *Phys. Rev. Lett.* **71**:1355 (1993).
44. P. Bouyer and M. A. Kasevich, "Heisenberg-Limited Spectroscopy with Degenerate Bose-Einstein Gases," *Phys. Rev. A* **56**:R1083 (1997).
45. J. A. Dunningham, K. Burnett, and S. M. Barnett, "Interferometry below the Standard Quantum Limit with Bose-Einstein Condensates," *Phys. Rev. Lett.* **89**:150401 (2002).
46. S. M. Barnett and D. T. Pegg, "On Hermitian Optical Phase Operator," *J. Mod. Opt.* **36**:7 (1989).
47. Z. Hradil, "Estimation of Counted Quantum Phase," *Phys. Rev. A* **51**:1870 (1995).
48. B. C. Sanders and G. J. Milburn, "Optimal Quantum Measurements for Phase Estimation," *Phys. Rev. Lett.* **75**:2944 (1995).
49. D. W. Berry and H. M. Wiseman, "Optimal States and Almost Optimal Adaptive Measurements for Quantum Interferometry," *Phys. Rev. Lett.* **85**:5098 (2000).
50. L. Pezze and A. Smerzi, "Phase Sensitivity of a Mach-Zehnder Interferometer," *Phys. Rev. A* **73**:011801(R) (2006).
51. H. Uys and P. Meystre, "Quantum States for Heisenberg-Limited Interferometry," *Phys. Rev. A* **76**:013804 (2007).
52. G. A. Durkin and J. P. Dowling, "Local and Global Distinguishability in Quantum Interferometry," *Phys. Rev. Lett.* **99**:070801 (2007).
53. T. Kim, O. Pfister, M. J. Holland, J. Noh, and O. L. Hall, "Quantum States for Heisenberg-Limited Interferometry," *Phys. Rev. A* **57**:4004 (1998).
54. A. Kuzmich and L. Mandel, "Sub-Shot-Noise Interferometric Measurements with Two-Photon States," *Quantum Semicl. Opt.* **10**:493 (1998).
55. R. A. Campos, C. C. Gerry, and A. Benmoussa, "Optical Interferometry at the Heisenberg Limit with Twin Fock States and Parity Measurements," *Phys. Rev. A* **68**:023810 (2003).
56. R. C. Pooser and O. Pfister, "Particle-Number Scaling of the Phase Sensitivity in Realistic Bayesian Twin-Mode Heisenberg-Limited Interferometry," *Phys. Rev. A* **69**:043616 (2004).
57. F. W. Sun, B. H. Liu, Y. X. Gong, Y. F. Huang, Z. Y. Ou, and G. C. Guo, "Experimental Demonstration of Phase Measurement Precision Beating Standard Quantum Limit by Projection Measurement," *Europhys. Lett.* **82**:24001 (2008).
58. B. C. Sanders, "Quantum Dynamics of the Nonlinear Rotator and the Effects of Continual Spin Measurement," *Phys. Rev. A* **40**:2417 (1989).
59. A. N. Boto, P. Kok, D. S. Abrams, S. L. Braunstein, C. P. Williams, and J. P. Dowling, "Quantum Interferometric Optical Lithography: Exploiting Entanglement to Beat the Diffraction Limit," *Phys. Rev. Lett.* **85**:2733 (2000).
60. C. K. Hong, Z. Y. Ou, and L. Mandel, "Measurement of Subpicosecond Time Intervals between Two Photons by Interference," *Phys. Rev. Lett.* **59**:2044 (1987).
61. C. C. Gerry and R. A. Campos, "Generation of Maximally Entangled Photonic States with a Quantum-Optical Fredkin Gate," *Phys. Rev. A* **64**:063814 (2001).
62. G. J. Milburn, "Quantum Optical Fredkin Gate," *Phys. Rev. Lett.* **62**:2124 (1989).
63. N. Imoto, H. A. Haus, and Y. Yamamoto, "Quantum Nondemolition Measurement of the Photon Number via Optical Kerr Effect," *Phys. Rev. A* **32**:2287 (1985).
64. M. Kitagawa and Y. Yamamoto, "Number-Phase Minimum-Uncertainty States with Reduced Number Uncertainty in a Kerr Nonlinear Interferometer," *Phys. Rev. A* **34**:3974 (1986).
65. P. Kok, H. Lee, and J. P. Dowling, "Single-Photon Quantum Nondemolition Detectors Constructed with Linear Optics and Projective Measurements," *Phys. Rev. A* **66**:063814 (2002).
66. Q. A. Turchette, C. J. Hood, W. Lange, H. Mabuchi, and H. J. Kimble, "Measurement of Conditional Phase-Shifts for Quantum Logic," *Phys. Rev. Lett.* **75**:4710 (1995).
67. H. Schmidt and A. Imamoglu, *Opt. Lett.* **21**:1936 (1996).

68. M. D. Lukin, "Trapping and Manipulating Photon States in Atomic Ensembles," *Rev. Mod. Phys.* **75**:457 (2003).
69. M. Fleischhauer, A. Imamoglu, and J. Marangos, "Electromagnetically Induced Transparency: Optics in Coherent Media," *Rev. Mod. Phys.* **77**:633 (2005).
70. E. Knill, R. Laflamme, and G. J. Milburn, "A Scheme for Efficient Quantum Computation with Linear Optics," *Nature* **409**:46 (2001).
71. H. Lee, P. Kok, and J. P. Dowling, "From Linear Optical Quantum Computing to Heisenberg-Limited Interferometry," *J. Opt. B* **6**:S769 (2004).
72. H. Lee, P. Kok, N. J. Cerf, and J. P. Dowling, "Linear Optics and Projective Measurements Alone Suffice to Create Large-Photon-Number Path Entanglement," *Phys. Rev. A* **65**:030101(R) (2002).
73. P. Kok, H. Lee, and J. P. Dowling, "Creation of Large-Photon-Number Path Entanglement Conditioned on Photodetection," *Phys. Rev. A* **65**:052104 (2002).
74. J. Fiurasek, "Conditional Generation of n -Photon Entangled States of Light," *Phys. Rev. A* **65**:053818 (2002).
75. H. Cable and J. P. Dowling, "Efficient Generation of Large Number-Path Entanglement Using Only Linear Optics and Feed-Forward," *Phys. Rev. Lett.* **99**:163604 (2007).
76. K. T. Kapale and J. P. Dowling, "Bootstrapping Approach for Generating Maximally Path-Entangled Photon States," *Phys. Rev. Lett.* **99**:053602 (2007).
77. N. M. VanMeter, P. Lougorski, D. B. Uskor, K. Kiding, J. Eisert, and J. P. Dowling, "General Linear-Optical Quantum State Generation Scheme: Applications to Maximally Path-Entangled States," *Phys. Rev. A* **76**:063808 (2007).
78. E. J. S. Fonseca, C. H. Monken, and S. Padua, "Measurement of the de Broglie Wavelength of a Multiphoton Wave Packet," *Phys. Rev. Lett.* **82**:2868 (1999).
79. M. D'Angelo, M. V. Chekhova, and Y. Shih, "Two-Photon Diffraction and Quantum Lithography," *Phys. Rev. Lett.* **87**:013602 (2001).
80. D. Bouwmeester, "High $N00N$ for Photons," *Nature* **429**:139 (2004).
81. P. Walther, J. W. Pan, M. Aspelmeyer, R. Ursin, S. Gasparoni, and A. Zeilinger, "De Broglie Wavelength of a Non-Local Four-Photon State," *Nature* **429**:158 (2004).
82. M. W. Mitchell, J. S. Lundeen, and A. M. Steinberg, "Super-Resolving Phase Measurements with a Multiphoton Entangled State," *Nature* **429**:161 (2004).
83. J. J. Bollinger, W. M. Itano, and D. J. Wineland, "Optimal Frequency Measurements with Maximally Correlated States," *Phys. Rev. A* **54**:R4649 (1996).
84. D. Leibfried, M. D. Barrett, T. Schaertz, J. Britton, J. Chiaverini, W. M. Itano, J. D. Jost, C. Langer, and D. J. Wineland, "Toward Heisenberg-Limited Spectroscopy with Multiparticle Entangled States," *Science* **304**:1476 (2004).
85. D. Leibfried, E. Knill, S. Seidelin, J. Britton, R. B. Blakestad, J. Chiaverini, D. B. Hume, et al., "Creation of a Six-Atom Schrödinger Cat State," *Nature* **438**:639 (2005).
86. K. J. Resch, K. L. Pregnell, R. Prevedel, A. Glichrist, G. J. Pryde, J. L. O'Brien and A. G. White, "Time-Reversal and Super-Resolving Phase Measurements," *Phys. Rev. Lett.* **98**:223601 (2007).
87. T. Nagata, R. Okamoto, J. L. O'Brien, K. Sasaki, and S. Takeuchi, "Beating the Standard Quantum Limit with Four-Entangled Photons," *Science* **316**:726 (2007).
88. J. L. O'Brien, "Precision without Entanglement," *Science* **318**:1393 (2007).
89. J. P. Dowling, "Kittens Catch Phase," *Nature* **450**:362 (2007).
90. B. L. Higgins, D. W. Berry, J. D. Bartlett, H. M. Wiseman, and G. J. Pryde, "Entanglement-Free Heisenberg-Limited Phase Estimation," *Nature* **450**:393 (2007).
91. L. A. Lugiato, A. Gatti, E. Brambilla, "Quantum Imaging," *J. Opt. B* **4**:S176 (2002).
92. J. P. Dowling, A. Gatti, and A. Sergienko, "Special Issue: Quantum Imaging," *J. Mod. Opt.* **53**:(5–6) (2006).
93. J. Jacobson, G. Bjork, I. Chuang, and Y. Yamamoto, "Photonic de Broglie Waves," *Phys. Rev. Lett.* **74**:4835 (1995).
94. Y. Shih, "Quantum Imaging," *IEEE J. Sel. Top. Quantum Electron.* **13**:1016 (2007).
95. M. B. Nasr, B. E. A. Saleh, A. V. Sergienko, and M. C. Teich, "Demonstration of Dispersion-Canceled Quantum-Optical Coherence Tomography," *Phys. Rev. Lett.* **91**:083601 (2003).

96. J. G. Rarity, P. R. Tapster, E. Jakeman, T. Larchuk, R. A. Campos, M. C. Teich, and B. E. A. Saleh, "Two-Photon Interference in a Mach-Zehnder Interferometer," *Phys. Rev. Lett.* **65**:1348 (1990).
97. R. W. Boyd, *Nonlinear Optics*, 2nd ed., Academic Press, San Diego (2003).
98. E. M. Nagasako, S. J. Bently, R. W. Boyd, G. S. Agarwal, "Nonclassical Two-Photon Interferometry and Lithography with High-Gain Parametric Amplifiers," *Phys. Rev. A* **64**:043802 (2001).
99. G. S. Agarwal, K. W. Chan, R. W. Boyd, H. Cable, and J. P. Dowling, "Quantum States of Light Produced by a High-Gain Optical Parametric Amplifier for Use in Quantum Lithography," *J. Opt. Soc. Am. B* **24**:270 (2007).
100. F. Scarrino, C. Vitelli, F. DeMartini, R. Glasser, H. Cable, and J. P. Dowling, "Experimental Sub-Rayleigh Resolution by an Unseeded High-Gain Optical Parametric Amplifier for Quantum Lithography," *Phys. Rev. A* **77**:012324 (2008).
101. H. J. Chang, H. Shin, M. N. O'Sullivan-Hale, and R. W. Boyd, "Implementation of Sub-Rayleigh-Resolution Lithography Using an n -Photon Absorber," *J. Mod. Opt.* **53**:2271 (2006).
102. K. T. Kapale, L. D. DiDomenico, H. Lee, P. Kok and J. P. Dowling, "Quantum Interferometric Sensors," *Concepts of Physics* **2**:225 (2005).
103. G. Gilbert, M. Hamrick, and Y. S. Weinstein, "Quantum Sensor Miniaturization," *IEEE Photon. Tech. Lett.* **19**:1798 (2007).
104. V. Giovannetti, S. Lloyd, L. Maccone, and F. N. C. Wong, "Clock Synchronization with Dispersion Cancellation," *Phys. Rev. Lett.* **87**:117902 (2001).
105. G. Gilbert, M. Hamrick, and Y. S. Weinstein, "On the Use of Photonic N00N States for Practical Quantum Interferometry," *J. Opt. Soc. Am. B* **25**:1336 (2008).
106. M. A. Rubin and S. Kaushik, "Loss-Induced Limits to Phase Measurement Precision with Maximally Entangled States," *Phys. Rev. A* **75**:053805 (2007).
107. S. D. Huver, C. F. Wildfeuer, and J. P. Dowling, "Entangled Fock States for Robust Quantum Optical Sensors," *Phys. Rev. A* **78**:063828 (2008).
108. P. G. Kwiat, A. J. Berglund, J. B. Altepeter, and A. G. White, "Experimental Verification of Decoherence-Free Subspace," *Science* **290**:498 (2000).
109. M. A. Nielsen and I. L. Chuang, *Quantum Computation and Quantum Information*, Cambridge University Press, Cambridge (2000).
110. P. Kok, W. J. Munro, K. Memoto, T. C. Ralph, J. P. Dowling and G. J. Milburn, "Linear Optical Quantum Computing with Photonic Qubits," *Rev. Mod. Phys.* **79**:135 (2007).

This page intentionally left blank.

DO NOT DUPLICATE

INDEX

Index note: The *f* after a page number refers to a figure, the *n* to a note, and the *t* to a table.

- Abbe number, 2.23, 2.28, 2.29*f*
Aberration control, 3.8
Above threshold ionization (ATI), 21.14–21.17
 plateau extension of, 21.19
 quasi-classical, 21.16–21.17, 21.16*f*
 relativistic electron, 21.20, 21.21*f*
 strong field interactions with atoms,
 21.14–21.17, 21.15*f*, 21.16*f*
Abrasion resistance, of antireflection coatings,
 7.31–7.32, 7.32*f*
Abrasive forming, polymer, 3.11–3.12
Absorbance, water, 1.10
Absorbers:
 black, 7.104
 saturable, 18.5–18.11
 fast, 18.9–18.10
 self-amplitude modulation, 18.5–18.7,
 18.6*f*, 18.7*f*
 semiconductor saturable absorber mirrors,
 18.3, 18.10–18.11
 slow, 18.7–18.9, 18.8*f*
Absorbing compounds (in black surfaces), 6.15
Absorbing media:
 antireflection coatings for, 7.26, 7.27
 in multilayer reflectors, 7.37–7.38, 7.38*f*
Absorptance:
 of metals, 4.6, 4.39, 4.40*f*–4.42*f*, 4.48, 4.49
 and emittance, 4.49, 4.49*f*, 4.50*t*, 4.51*t*
 and mass attenuation coefficients for
 photons, 4.48*t*
 of optical coatings, 7.12–7.13
 of water, 1.5*t*, 1.9
Absorption:
 bio-optical models for, 1.27–1.29,
 1.27*f*, 1.28*t*
 cutoff filters based on, 7.60, 7.60*f*
 defect-related, 5.37–5.39, 5.38*f*, 5.39*f*
 by detritus in water, 1.26–1.27
 in dielectrics, 4.4
 Absorption (*Cont.*):
 direct interband, 8.27–8.28
 by dissolved organic matter, 1.22–1.23, 1.23*t*
 fundamental edge of (*see* Fundamental
 absorption edge)
 impurity-related, 5.37–5.39, 5.38*f*, 5.39*f*,
 5.51–5.52, 5.51*f*
 interband (*see* Interband absorption)
 lattice, 5.13–5.20, 5.15*f*, 5.17*t*, 5.18*f*–5.19*f*,
 5.20*t*, 5.21*f*
 of light in laser cooling, 20.4
 magnetoabsorption, 5.51–5.52, 5.51*f*
 measurement of, 5.64
 in multilayer reflectors, 7.40–7.41, 7.41*f*
 multiphonon, 5.16–5.17, 5.17*t*, 5.18*f*
 nonlinear, 16.7–16.9
 by organic detritus, 1.25–1.27, 1.25*t*, 1.26*f*
 in overdense plasmas, 21.47–21.48, 21.47*f*
 phonon, 5.13–5.16, 5.15*f*
 by phytoplankton, 1.23–1.25, 1.24*f*–1.25*f*
 resonance, 21.47–21.48, 21.47*f*
 by sea water, 1.21, 1.22*t*
 in solids, 8.27–8.28
 spectral, 1.26–1.27, 1.26*f*
 superlinear, 5.57
 two-photon, 5.56
 by water, 1.20–1.29
 Absorption coefficient(s), 1.23*t*
 for optical constants, 5.9–5.10
 spectral
 of natural waters, 1.27*f*
 for phytoplankton, 1.24, 1.24*f*,
 1.25*f*, 1.25*t*
 of water, 1.5*t*, 1.7*f*, 1.10, 1.17, 1.18*t*
 natural water, 1.20–1.21
 sea water, 1.21, 1.22*t*
 Absorption coefficient dependence, 5.20, 5.21*f*
 Absorption saturation, 16.21–16.22, 16.21*f*
 Absorption spectrum, 5.12, 22.6, 22.6*f*, 22.7*f*

- Absorption transitions:
 direct interband, 5.22–5.23, 5.22f–5.23f, 5.25f
 indirect, 5.22–5.24, 5.24f–5.25f
- ac Stark effect, 16.3t, 16.7
- Acceleration:
 bubble, 21.41–21.42, 21.42f
 direct laser, 21.43
 electron, 21.39–21.42, 21.40f
 MeV proton, 21.54, 21.54f
 plasma beat wave, 21.41
 target normal sheath, 21.54
- Accordion solutions, 15.28
- Achromatic antireflection coatings, 7.29
- Achromatic beam splitters, 7.62–7.65, 7.62f–7.65f
- Acktar black coatings, 6.55
- Acoustic optical modulators (AOMs), 14.19, 20.13
- Acoustic phonons, 5.14–5.16
- Active nonlinear optical phenomena, 5.54, 5.54t
- Active optical limiting, 13.1–13.3
- Adachi dispersion model, 2.15, 2.22
- Additive pulse modelocking (APM), 18.3, 18.14
- Additives, polymers as, 3.5
- Adiabatic ionization stabilization, 21.20, 21.21, 21.22f
- Adiabatic potentials, 21.23, 21.24f
- Aerogels, 6.15
- Aeroglaze Z series, 6.36f, 6.37, 6.37f, 6.39, 6.39f–6.42f
- Aeroglaze Z302, 6.14
- Aeroglaze Z306, 6.14, 6.17, 6.29f, 6.35
- Akzo (company), 6.14
- Akzo Nobel paints, 6.39, 6.42f, 6.43f
- All-dielectric color selective (dichroic) beam splitters, 7.65–7.66, 7.67f
- All-dielectric reflectors:
 broadband, 7.39, 7.40f, 7.45–7.47, 7.45f–7.47f
 interference filters with, 7.81
- Altaite (PbTe), 2.40t, 2.44t, 2.48t, 2.52t, 2.57t, 2.64t, 2.69t
- Alumina (Al₂O₃), 2.38t, 2.46t, 2.47t, 2.50t, 2.55t, 2.60t, 2.70t, 2.76t
- Aluminum:
 absorptance of, 4.40f, 4.48t, 4.51t
 anodized, 6.5f, 6.38f, 6.58f
 black velvet cloth on, 6.31f
 commando cloth on, 6.31f
 flame-sprayed, 6.57
 grooved and blazed, 6.21, 6.24f, 6.25f
 optical constants for, 4.11
 optical properties of, 4.12t, 4.20f, 4.21f
- Aluminum (Cont.):
 penetration depth of, 4.47f
 physical properties of, 4.52t, 4.54t
 reflectance of, 4.27t–4.28t, 4.40f, 4.44f, 4.46f
 roughened, 6.21, 6.22f, 6.30f
 sandblasted, 6.45f
 thermal properties of
 coefficient of linear thermal expansion, 4.56t, 4.57f
 elastic properties, 4.69t
 at room temperature, 4.55t
 specific heat, 4.65t, 4.66f
 strength and fracture properties, 4.70t
 thermal conductivity, 4.58t, 4.59f–4.60f
 and UV light, 6.21
- Aluminum alloys:
 reflectance of, 4.44f
 thermal conductivity of, 4.59f–4.60f
- Aluminum mirrors, 7.106–7.108, 7.106f–7.108f
- Aluminum oxynitride (Al₂₃O₂₇N₅) (ALON), 2.38t, 2.44t, 2.47t, 2.50t, 2.55t, 2.60t, 2.76t
- American Society for Testing and Materials (ASTM), 6.17
- Ames 24E, 6.7f, 6.26, 6.26f, 6.27f, 6.34
- Ames 24E2, 6.27f, 6.28f, 6.34
- Ames Perfect Diffuse Reflector, 6.7f
- Ames Research Center, 6.34
- Ammonium phosphate (NH₄H₂PO₄, ADP), 2.40t, 2.45t, 2.48t, 2.52t, 2.57t, 2.64t
- Ammosoc-Delone-Krainov (ADK) ionization rate, 21.13
- Amorphous materials, 2.5, 12.26t [*See also* Glass(es)]
- Ampere law, 2.6
- Amplifiers:
 flashlamp pumped Nd:glass, 21.5
 lock-in, 5.64
 optical parametric, 23.13–23.14
 Raman, 15.4, 15.4f, 22.15
 Ti:sapphire, 21.5
- Amplifying media, antireflection coatings for, 7.26, 7.27
- AMTIR-1 glass, 2.43t, 2.49t, 2.54t, 2.59t, 2.68t
- AMTIR-3 glass, 2.43t, 2.49t, 2.54t, 2.59t, 2.68t
- ance suffix, 4.5
- Angle-of-incidence effects, for cutoff filters, 7.56
- Angular distribution, 1.12, 21.8–21.9, 21.9f
- Angular sensitivity, in two-material periodic multilayers theory, 7.37

- Anharmonic oscillator model, of second-order nonlinear optical susceptibility, **10.7–10.9, 10.8f**
- Anisotropic crystals, propagation of light and, **8.8–8.11, 8.9t, 8.10f**
- Anisotropic scattering, **12.7**
- Annealing:
of glass, **2.28**
of optical surfaces, **19.3**
- Anodized aluminum surface, **6.5f, 6.38f, 6.58f**
- Anodized surface treatments, **6.44–6.49**
(*See also specific anodized treatments, e.g.: Martin Black*)
- Anodized surfaces, **6.3t**
- Anomalous (negative) dispersion, **4.4, 18.11**
- Anomalous skin effect, **21.49**
- Antimony flint glass, **2.43t**
- Antireflection (AR) coatings, **7.15–7.32**
of absorbing and amplifying media, **7.26, 7.27**
homogeneous-layer, **7.16–7.23, 7.17f–7.19f, 7.20t–7.21t, 7.22f–7.23f**
inhomogeneous and structured, **7.23–7.26, 7.24f, 7.26f**
at nonnormal angle of incidence, **7.28–7.31, 7.28f–7.31f**
nonoptical properties of, **7.31–7.32, 7.32f**
surface reflections and optical performance, **7.15–7.16, 7.16f**
of surfaces carrying thin films, **7.27–7.28, 7.28f**
universal, **7.26, 7.27f**
- Antiresonant Fabry-Perot saturable absorber (A-FPSA), **18.3, 18.11**
- Anti-Stokes four-wave mixing, coherent, **15.2t, 15.3t, 15.4, 15.4f**
- Anti-Stokes scattering, **15.1–15.3, 15.3t**
coherent Raman, **15.4, 15.4f, 15.34, 15.42, 15.42t, 15.43f**
multiple, **15.2f**
Raman, **15.32–15.34, 15.33f, 15.35f**
shifted Raman, **16.15, 16.15f**
stimulated Raman, **15.2t**
- Anti-Stokes shift, **15.2, 15.43**
- Anti-Stokes wave, **15.1, 15.43**
- APART (stray light analysis program), **6.19**
- Apparent optical properties (AOPs), of water, **1.4, 1.5t–1.6t, 1.12–1.13**
- Appliqués, black, **6.15**
- Arsenic antite (As_{Ga}), **18.3**
- Arsenic triselenide glass, **2.43t, 2.49t, 2.54t, 2.59t**
- Arsenic trisulfide glass, **2.43t, 2.49t, 2.54t, 2.59t**
- ASP (stray light analysis program), **6.19**
- Aspheric surfaces, **3.8–3.9**
- Assembly of polymers, mechanical, **3.14–3.16, 3.14f, 3.15f**
- Associative memory, optical, **12.34**
- Athermal glasses, **2.26**
- Athermalization, **3.9**
- Atom interferometry, **11.22–11.23, 11.24f**
- Atomic beams:
brightening of, **20.27–20.28, 20.27f**
collimation of, **20.15–20.16, 20.15f**
slowing of, **20.11–20.13, 20.12f, 20.12t, 20.13f**
- Atomic clocks, **20.28**
- Atomic coherence, maximal, **14.28–14.32, 14.29f–14.32f**
- Atomic funnels, **20.27**
- Atomic ionization, **21.3**
- Atomic layer deposition, **7.11**
- Atomic oxygen, black surfaces and, **6.16–6.17**
- Atomic resonance, **22.2–22.9, 22.3f**
about, **22.2–22.5, 22.3f**
double, **22.5–22.9**
features of, **22.5–22.6, 22.6f**
tunable double resonance, **22.6–22.9, 22.7f, 22.8f**
- Atomic systems, coherence in, **14.4–14.5, 14.4f**
- Atom-laser interactions:
applications of, **20.26–20.39**
atomic beam brightening, **20.27–20.28, 20.27f**
atomic clocks, **20.28**
Bose-Einstein condensation, **20.35–20.37, 20.36f**
dark states, **20.37–20.39, 20.38f**
optical lattices, **20.31–20.34, 20.32f–20.34f**
ultracold collisions, **20.28–20.31, 20.30f, 20.31f**
cooling atoms with, **20.3–20.21**
below Doppler limit, **20.17–20.21, 20.18f–20.20f**
history of, **20.3–20.4**
optical molasses, **20.13–20.17, 20.14f–20.16f**
properties of lasers, **20.4–20.6**
slowing atomic beams, **20.11–20.13, 20.12f, 20.12t, 20.13f**
theoretical description, **20.6–20.11, 20.9f**
trapping atoms with, **20.21–20.39**
magnetic traps, **20.21–20.23, 20.22f**
magneto-optical traps, **20.24–20.25, 20.24f, 20.26f**
optical traps, **20.23–20.24, 20.23f**

- Atoms (generally):
 in motion, 20.8–20.10, 20.9f
 trapping of neutral (*see* Trapping atoms)
 two-level, 20.6–20.8
- Atoms, strong field interactions with, 21.10–21.21
 above threshold ionization, 21.14–21.17,
 21.15f, 21.16f
 ionization stabilization, 21.20–21.21, 21.22f
 Keldysh parameter, 21.10
 multiphoton and quasi-classical regimes,
 21.10
 multiphoton ionization, 21.10–21.12, 21.11f
 relativistic effects, 21.19–21.20, 21.21f
 rescattering effects, 21.18–21.19, 21.18f, 21.19f
 tunnel ionization, 21.12–21.14, 21.12f, 21.14f
- Attenuation, in water:
 beam, 1.40–1.41, 1.41f, 1.42f
 diffuse and Jerlov water types, 1.42–1.46,
 1.43t–1.45t, 1.44f, 1.45f
- Attenuation functions, of water, 1.13
- Attenuators, neutral, 7.105, 7.105f
- Attosecond pulse generation, 21.31
- Axial thickness, of polymers, 3.10
- Backward Raman amplifiers, 15.4, 15.4f
 Backward Raman generators, 15.4, 15.4f
 Backward Raman scattering, 15.41, 21.38f,
 21.39
- Bacteria, in water, 1.14
- Baffle surfaces:
 design of, 6.19
 in extreme environments, 6.18
 selection process for black, 6.10–6.12,
 6.10f–6.11f, 6.12t–6.13t
- Ball Black, 6.5f, 6.6f, 6.50, 6.51f, 6.53f
- Band structures:
 defined, 9.3
 of solids, 8.24–8.27, 8.26f
- Bandpass filters, 7.73–7.96
 about, 7.73, 7.77–7.78, 7.77f–7.78f
 angular properties of, 7.91–7.94, 7.92f, 7.93f
 with multiple peaks, 7.90, 7.91f
 narrow- and medium-, 7.78–7.83, 7.79f,
 7.80f, 7.82f–7.88f
 nonpolarizing, 7.66, 7.67f
 square-top multicavity, 7.82–7.83, 7.82f–7.88f
 stability and temperature dependence of, 7.94
 very narrow, 7.83, 7.88–7.89, 7.89f
 wedge filters, 7.90, 7.91, 7.91f
 wide, 7.90, 7.90f
- Bandpass filters (*Cont.*):
 wide-angle, 7.93–7.94, 7.93f
 for XUV and x-ray regions, 7.94–7.96,
 7.95f–7.96f
- Bare states, 14.3
- Barium beta borate (BBO), 2.38t, 2.46t, 2.47t,
 2.50t, 2.55t, 2.60t, 2.75t, 17.1, 18.12
- Barium beta borate (BBO) optical parametric
 oscillators, 10.18, 10.19f
- Barium crown glass, 2.41t
- Barium dense flint glass, 2.42t
- Barium flint glass, 2.42t
- Barium fluoride (BaF₂), 2.38t, 2.44t, 2.47t,
 2.50t, 2.55t, 2.60t, 2.69t, 2.76t
- Barium titanate (BaTiO₃), 2.38t, 2.44t, 2.47t,
 2.55t, 2.76t, 12.13t, 12.14–12.16, 12.16f
- Barrier suppression ionization (BSI), 21.14, 21.14f
- Barrier suppression ionization thresholds,
 21.26, 21.27f
- Beam attenuation, in water, 1.13, 1.40–1.41,
 1.41f, 1.42f
- Beam attenuation coefficient, 1.7f, 1.10
- Beam cleanup, 12.30
- Beam propagation, split-step, 12.10
- Beam splitters (BSs), 7.61–7.67, 7.62f–7.68f,
 23.2, 23.2f, 23.14
- Beam walkoff time, 15.30
- Beating (phenomena), 18.19
- Beer's law, 8.7, 8.28
- Beryllium:
 absorptance of, 4.48t, 4.50t
 optical constants for, 4.11
 optical properties of, 4.12t, 4.21f, 4.26f
 penetration depth of, 4.47f
 physical properties of, 4.52t, 4.54t
 reflectance of, 4.28t–4.29t, 4.45f, 4.46f
 thermal properties of
 coefficient of linear thermal expansion,
 4.56t, 4.57f
 elastic properties, 4.69t
 at room temperature, 4.55t
 specific heat, 4.65t, 4.68f
 strength and fracture properties, 4.70t
 thermal conductivity, 4.58t, 4.59f–4.60f
- Beryllium surfaces, 6.51, 6.52, 6.53f, 6.58f
- Beta cloth, 6.57, 6.58f
- Betatron resonance, 21.42–21.43
- Biaxial crystals, 8.8, 8.9t, 8.10, 8.10f
- Bidirectional scatter distribution function
 (BSDF), 6.1, 6.9–6.10, 6.18–6.19

- Bio-optical models, of absorption, 1.27–1.29, 1.27f, 1.28t
- Bird-wing mirror, 12.7, 12.8f
- Birefringence, 8.9, 17.1
- Bismuth germanium oxide ($\text{Bi}_{12}\text{GeO}_{20}$) (BGO), 2.38t, 2.44t, 2.47t, 2.50t, 2.55t, 2.60t
- Bismuth triborate (BiB_3O_6) (BIBO), 2.38t, 2.46t, 2.47t, 2.50t, 2.55t, 2.60t, 2.75t, 17.14
- Bistable optical switches, 16.31
- Black absorbers, 7.104, 7.105f
- Black chrome, 6.53, 6.54, 6.54f
- Black coatings, 6.13t
- Black cobalt, 6.53, 6.54f
- Black dye, 6.15
- Black felt contact paper, 6.30f, 6.32f
- Black glass, 6.57
- Black Kapton, 6.57, 6.57f
- Black layer system, 6.15
- Black nickel, 6.21, 6.23f
- Black paint, 6.21, 6.24f
- Black surfaces, 6.1–6.59
 creation of, 6.13–6.16
 environmental degradation of, 6.16–6.18
 atomic oxygen effects, 6.16–6.17
 extreme environments, 6.18
 outgassing, 6.17
 particle generation, 6.17–6.18
 for far-infrared applications, 6.21, 6.26–6.34, 6.28f–6.34f
 Ames 24E and 24E2, 6.26f, 6.28f, 6.34
 Cornell Black, 6.26f, 6.27
 Infrablack, 6.26f, 6.28, 6.28f
 multiple-layer approach, 6.26, 6.26f–6.27f
 Teflon overcoat, 6.27
 optical characterization of, 6.18–6.21, 6.20t
 paints and surface treatments, 6.35–6.58, 6.37f, 6.43f, 6.53f
 Acktar black coatings, 6.55
 Aeroglaze Z series, 6.36f, 6.37, 6.37f, 6.39, 6.39f–6.42f
 Akzo Nobel paints, 6.39, 6.42f, 6.43f
 anodized processes, 6.44–6.49, 6.47f, 6.48f, 6.51f, 6.53f
 black glass, 6.57
 Black Kapton, 6.57, 6.57f
 carbon nanotubes and nanostructured materials, 6.55, 6.59f
 Cardinal Black, 6.36f, 6.39, 6.44f
 Cat-a-lac Black, 6.39, 6.42f, 6.53f
 DeSoto Black, 6.37f, 6.39
- Black surfaces, paints and surface treatments (*Cont.*):
 DURACON, 6.55–6.56
 electrically conductive black paint, 6.56
 electrodeposited surfaces, 6.53–6.54, 6.54f, 6.55f
 etching of electroless nickel, 6.49–6.50, 6.50f, 6.51f, 6.53f
 flame-sprayed aluminum, 6.57
 Floquil, 6.44
 gold blacks, 6.57
 high-resistivity coatings, 6.56
 IBM Black (tungsten hexafluoride), 6.56
 ion beam-sputtered surfaces, 6.53
 Parson's Black, 6.44, 6.53f
 plasma-sprayed surfaces, 6.50–6.52, 6.51f–6.53f
 silicon carbide, 6.56
 SolarChem, 6.44, 6.48f, 6.53f
 specular metallic anodized surfaces, 6.57, 6.58f
 sputtered and CVD surfaces, 6.56
 3M paints and derivatives, 6.35–6.37, 6.36f, 6.38f, 6.53f
 ZO-MOD BLACK, 6.56
 selection process for, 6.10–6.12, 6.10f–6.11f, 6.12t–6.13t
 and substrates, 6.34–6.35
 types and morphologies of, 6.1–6.10, 6.2t–6.4t, 6.5f–6.8f
 for ultraviolet applications, 6.21, 6.22f–6.25f
- Black Tedlar, 6.57f
- Black velvet cloth, 6.31f, 6.33f
- Blazing, gratings and, 5.60
- Blink response (of eye), 13.1
- Bloch equations, optical, 11.3–11.6
- Bloch solution, 8.25
- Bloch sphere, 11.4
- Blooming, thermal, 16.22
- Bombardment, of carbon surface, 6.8f
- Boron, 5.82f, 5.83
- Boron Black, 6.50, 6.51, 6.52f
- Boron carbide, 6.51
- Borosilicate crown glass, 2.41t
- Bose-Einstein condensation (BEC), 14.22, 20.26, 20.35–20.37, 20.36f
- Bose-Einstein distribution of states, 2.16
- Boulder Damage Symposium, 19.1–19.2
- Bound electronic optical Kerr effect, 16.12–16.13, 16.13f

- Bound excitons, 5.26*t*, 5.29, 5.46
 Bound-electronic optical Kerr effect, 16.3*t*
 Bragg grating, 22.9–22.11, 22.10*f*
 Bragg reflection, 20.33
 Breault Research Organization, 6.1
 Bremsstrahlung heating, inverse, 21.37, 21.37*f*
 Brewster's angle, 8.12, 8.23
 Bridge mirror, 12.7, 12.8*f*
 Bright (coupled) state, 14.4
 Brightening, atomic beam, 20.27–20.28, 20.27*f*
 Brillouin gain, 15.48
 Brillouin scattering, 21.38
 for crystals and glasses, 2.27
 defined, 2.27
 in measurement, 5.76, 5.77
 in nonlinear optics, 16.14, 16.18–16.19
 Raman vs., 15.1
 in solids, 8.18
 stimulated [*see* Stimulated Brillouin scattering (SBS)]
 in strong-field physics, 21.38
 Brillouin spectroscopy, 5.57–5.58
 Brillouin zone, 8.25, 8.26*f*, 8.27, 8.29, 9.3, 9.4
 Brillouin-enhanced four-wave mixing (BEFWM), 15.53, 15.54, 15.54*f*
 Broadband light sources, 5.58–5.59
 Broadband parametric amplification, 18.12
 Broadband reflectors, all-dielectric, 7.45–7.47, 7.45*f*–7.47*f*
 Broadband SBS slow light, 22.15
 Broadband transient Raman scattering, 15.28–15.32, 15.29*f*
 Broadening, in spectral lines, 14.13, 14.14
 Bromellite (BeO), 2.38*t*, 2.46*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.70*t*
 Bromyrite (AgBr), 2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.69*t*, 2.76*t*
 BS-39B glass, 2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
 Bubble acceleration, 21.41–21.42, 21.42*f*
 Bulk (term), 4.3
 Bulk compound semiconductors, 12.20–12.21, 12.20*t*, 12.21*f*
 Bulk modulus, for metals, 4.69*t*
 Cadmium germanium diarsenide (CdGeAs₂), 2.38*t*, 2.45*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.61*t*, 2.74*t*
 Cadmium selenide (CdSe), 2.38*t*, 2.46*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.61*t*, 2.70*t*, 2.77*t*
 Cadmium telluride (CdTe), 2.39*t*, 2.44*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.69*t*, 2.77*t*
 Calcite (CaCO₃), 2.38*t*, 2.46*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.61*t*, 2.74*t*, 2.77*t*
 Calcium molybdate (powellite) (CaMoO₄), 2.38*t*, 2.45*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.61*t*, 2.72*t*, 2.77*t*
 Carbamide ((NH₄)₂CO), 2.40*t*
 Carbon black particles, 6.15
 Carbon deposition, laser-assisted, 19.5
 Carbon disulfide (CS₂), 16.13–16.14, 16.14*f*
 Carbon nanotubes, 6.55, 6.59*f*
 Carbon surface bombarded with positive argon ions, 6.8*f*
 Carbon-black suspensions (CBSs), 13.10–13.11, 13.11*f*
 Cardinal Black, 6.36*f*, 6.39, 6.44*f*
 Carrier trapping, 18.21–18.23, 18.22*f*
 Carrier-carrier scattering, 18.20
 Cascaded limiters, 13.6
 Cascaded $x^{(1)}:x^{(1)}$ processes, of third-order optical nonlinearities, 16.20–16.22, 16.21*f*
 Cascaded $x^{(2)}:x^{(2)}$ processes, of third-order optical nonlinearities, 16.22–16.24, 16.23*f*, 16.24*f*
 Casting, of polymers, 3.11
 Cat mirror, 12.7, 12.8*f*
 Cat-a-lac Black, 6.39, 6.42*f*, 6.53*f*
 Cauchy dispersion formula, 2.21
 Causality, principle of, 2.8
 Cavities, surface absorption and, 6.15
 Cavity resonance, for cw optical parametric oscillators, 17.2–17.4, 17.3*f*, 17.4*f*
 Cellulose acetate butyrate, 3.4*t*
 Center-of-mass motion of atoms, 20.4
 Ceragyrite (AgCl), 2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.69*t*, 2.76*t*
 Cesium chloride (CsCl), 2.68*t*
 Cesium iodide (CsI), 2.39*t*, 2.44*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.62*t*, 2.68*t*, 2.77*t*
 Cesium lithium borate (CsLiB₆O₁₀) (CLBO), 2.39*t*, 2.45*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.62*t*
 Chalcopyrite, 2.74*t*
 Charge-coupled devices (CCDs), 5.61
 Charge-resonance enhanced ionization (CREI), 21.26
 Chemglaze Z series (*see* Aeroglaze Z series)
 Chemical vapor deposition (CVD), 6.56, 7.11
 Chirped multilayers, 7.47, 7.48
 Chirped pulse amplification (CPA) lasers, 21.4–21.5, 21.4*f*
 Chirped pulse excitation, 11.25–11.26
 Chirps, in spectroscopy, 11.3

- Chlorophyll:
 and absorption by phytoplankton, 1.23,
 1.24, 1.28
 and beam attenuation, 1.41
 and diffuse attenuation, 1.44
 fluorescence by, 1.49
 and remote sensing, 1.46
- Chloroplasts, 1.23
- Chromatic aberration control, 3.8
- Chromium:
 absorptance of, 4.48*t*, 4.50*t*
 optical properties of, 4.13*t*–4.14*t*, 4.22*f*
 physical properties of, 4.54*t*
 reflectance in, 4.30*t*–4.31*t*
 thermal properties of, 4.69*t*
- Clamshell housings, 3.15, 3.15*f*
- Classical electronic polarization theory, 2.14
- Classical harmonic oscillator model, 10.5–10.7,
 10.6*f*
- Clausius-Mossotti equation, 8.7
- Clausius-Mossotti relationships, 2.24
- Cleaning and cleanliness, of optical surfaces,
 19.3–19.5
- Clocks, atomic, 20.28
- Closed family, 20.38
- Cluster electron heating, 21.34, 21.35
- Cluster expansion, 21.35, 21.35*f*
- Clusters, strong field interactions with, 21.31–21.36
 Coulomb explosion, 21.33–21.34
 intense laser pulse interactions, 21.35–21.36,
 21.36*f*
 ionization mechanisms in, 21.31–21.33, 21.32*f*
 nanoplasma description, 21.34–21.35, 21.35*f*
- Coatings:
 antireflection, 7.15–7.32
 of absorbing and amplifying media, 7.26, 7.27
 homogeneous-layer, 7.16–7.23, 7.17*f*–
 7.19*f*, 7.20*t*–7.21*t*, 7.22*f*–7.23*f*
 inhomogeneous and structured, 7.23–7.26,
 7.24*f*, 7.26*f*
 at nonnormal angle of incidence, 7.28–7.31,
 7.28*f*–7.31*f*
 nonoptical properties of, 7.31–7.32, 7.32*f*
 surface reflections and optical perfor-
 mance, 7.15–7.16, 7.16*f*
 of surfaces carrying thin films, 7.27–7.28,
 7.28*f*
 universal, 7.26, 7.27*f*
 filters with metallic reflecting, 7.80–7.81
 high performance optical multilayer, 7.96–7.98
- Coatings (*Cont.*):
 and interference polarizers, 7.70–7.72,
 7.70*f*–7.72*f*
 laser-induced damage in, 19.3–19.4
 measurements on, 7.12–7.14
 narrowband reflection, 7.43, 7.44*f*
 phase, 7.101, 7.101*f*–7.104*f*, 7.102
 reflection, 7.106–7.113, 7.106*f*–7.113*f*
 thin-film
 and antireflection coatings, 7.27–7.28,
 7.28*f*
 manufacturing of, 7.10–7.12
 of metal, 7.104, 7.104*f*
 theory and design of, 7.5–7.10, 7.6*f*, 7.9*f*
 transmission and reflection of, 7.3
 types of, 3.17–3.18, 6.13*t*
 for ultrafast optics, 7.47–7.48, 7.48*f*
 (*See also specific coatings, e.g.: Ebanol C coating*)
- Coefficient of linear thermal expansion, 4.7,
 4.56*t*, 4.57*f*–4.58*f*
- Coefficient of thermal expansion (CTE), of
 metals, 4.6–4.7, 4.10*t*, 4.53, 4.55*t*
- Coherence:
 in atomic systems, 14.4–14.5, 14.4*f*,
 14.28–14.32, 14.29*f*–14.32*f*
 maximal, 14.3, 14.28–14.32, 14.29*f*–14.32*f*
- Coherence length, 10.15
- Coherent anti-Stokes four-wave mixing, 15.2*t*,
 15.3*t*, 15.4, 15.4*f*
- Coherent anti-Stokes Raman scattering
 (CARS), 15.4, 15.4*f*, 15.34, 15.42, 15.42*t*,
 15.43*f*, 16.3*t*, 16.4, 16.5*f*, 16.17–16.18,
 16.17*f*
- Coherent control, 18.21
- Coherent excitons, 18.19–18.20
- Coherent image amplification, 12.29, 12.30
- Coherent optical transients, 11.1–11.28
 chirped pulse excitation, 11.25–11.26
 and cw spectroscopy, 11.2
 experimental considerations, 11.26–11.28,
 11.27*f*
 free polarization decay, 11.7–11.11, 11.8*f*,
 11.10*f*, 11.11*f*
- Maxwell-Bloch equations, 11.6–11.7
 optical Bloch equations, 11.3–11.6
 phase conjugate geometry and optical
 Ramsey fringes, 11.19–11.22, 11.20*f*,
 11.21*f*
- photon echo, 11.11–11.15, 11.12*f*, 11.13*f*,
 11.15*f*

- Coherent optical transients (*Cont.*):
 stimulated photon echo, 11.15–11.19,
 11.16*f*–11.19*f*
 two-photon transitions and atom
 interferometry, 11.22–11.23, 11.24*f*
- Coherent population return (CPR), 14.1, 14.30,
 14.30*f*, 14.31*f*
- Coherent population trapping (CPT), 14.1,
 14.3–14.5, 14.7, 20.37
- Coherent Raman scattering, 15.3
- Coherent Stokes Raman scattering (CSRS),
 16.3*t*, 16.17–16.18, 16.17*f*
- Cold mirrors, 7.58
- Collected volatile condensable materials
 (CVCMM), 6.17
- Collective tunneling, 21.18
- Collet-and-cap housings, 3.15, 3.15*f*
- Colliding pulse modelocking (CPM), 18.3
- Collimation, of atomic beams, 20.15–20.16, 20.15*f*
- Collisional broadening, spectral-line, 14.13
- Collisional heating, 21.37, 21.37*f*
- Collisional ionization, 21.31, 21.32
- Collisions, 20.28–20.31, 20.30*f*, 20.31*f*
 excited-state, 20.29
 ground-state, 20.29
 trap loss, 20.29
- Colloids, in water, 1.14
- Color, ocean, 1.46
- Color center lasers, 18.10
- Color-selective beam splitters, 7.65–7.66, 7.66*f*,
 7.67*f*
- Commando cloth, 6.31*f*, 6.33*f*
- Common glasses, 2.3
- Communications, out-of-plane coupling for, 9.12
- Complex Fresnel relation, 8.15
- Complex refractive index, 1.16–1.17, 2.8
- Compliance tensors, 2.30, 2.31*t*
- Composites, 12.26*t*, 12.27*t*
- Compositional modulation, 5.65, 5.66*t*, 5.67
- Compression molding, of polymers, 3.12
- Computer numerical control (CNC) lathe
 turning, 3.12
- Condon point, for laser light, 20.29
- Conduction band (CB), 18.3
- Conductivity:
 of metals, 4.6
 of paints, 6.12, 6.12*t*
 of polymers, 3.3–3.4
 of solids, 8.4
 of water, 1.16
- Conjugate mirrors, phase (*see* Phase conjugate
 mirrors)
- Connected networks, 3D photonic crystals
 and, 9.4–9.5
- Constringence, 2.23 (*See also* Abbe number)
- Contamination, of optical surfaces, 19.4–19.5
- Contamination control, 6.16
- Continuous-wave (cw) lasers, 7.14, 14.16–14.18,
 14.17*f*
- Continuous-wave optical parametric oscillators
 (cw OPOs), 17.1–17.31
 cavity resonance configurations for, 17.2–17.4,
 17.3*f*, 17.4*f*
 for correlated twin beams of light, 17.28,
 17.29*f*, 17.30*f*
 for hyperspectral imaging, 17.27–17.28
 limitations of, 17.30
 for metrology and optical frequency synthesis,
 17.28, 17.29
 multiple-resonant oscillators, 17.16–17.21
 doubly resonant, 17.16–17.17
 pump-enhanced singly resonant,
 17.17–17.20, 17.18*f*–17.20*f*
 triply resonant, 17.20–17.21, 17.21*f*
 singly resonant oscillators, 17.4–17.16
 guided-wave nonlinear structures,
 17.15–17.16
 MgO:sPPLT in, 17.14–17.15, 17.15*f*
 PPLN crystals in, 17.4–17.13, 17.6*f*–17.11*f*
 QPM nonlinear materials, 17.13–17.14
 in spectroscopy, 17.21–17.27
 high-resolution Doppler-free, 17.27
 photoacoustic, 17.22–17.24, 17.22*f*–17.24*f*
 single-pass absorption, 17.24–17.27,
 17.25*f*, 17.26*f*
 technological advances in, 17.1–17.2,
 17.30–17.31
- Continuous-wave (cw) Q-switched modelocking,
 18.5*f*
- Continuous-wave (cw) spectroscopy, 11.2
- Continuum excitations, 18.20
- Continuum pulse generation, 18.4
- Controlled grinding, 19.3
- Conventional evaporation, 7.11
- Cooling:
 of atoms with atom-laser interactions,
 20.3–20.21
 below Doppler limit, 20.17–20.21,
 20.18*f*–20.20*f*
 history of, 20.3–20.4

- Cooling, of atoms with atom-laser interactions
(*Cont.*):
 optical molasses, 20.13–20.17,
 20.14f–20.16f
 properties of lasers, 20.4–20.6
 slowing atomic beams, 20.11–20.13,
 20.12f, 20.12t, 20.13f
 theoretical description, 20.6–20.11, 20.9f
 Doppler, 20.13–20.15, 20.14f
 laser (*see* Laser cooling)
 polarization gradient, 20.17
 Raman, 20.21
 Cooling rate, for glasses, 2.5n
 Copper:
 absorptance of, 4.40f, 4.48t, 4.50t
 optical properties of, 4.12t–4.13t, 4.22f
 physical properties of, 4.52t–4.54t
 reflectance of, 4.29t–4.30t, 4.40f
 thermal properties of
 coefficient of linear thermal expansion,
 4.56t, 4.57f
 elastic properties, 4.69t
 at room temperature, 4.55t
 specific heat, 4.65t, 4.66f
 strength and fracture properties, 4.70t
 thermal conductivity, 4.58t, 4.60f–4.61f
 Copper black, 6.21, 6.23f
 Copper gallium sulfide (CuGaS₂), 2.39t, 2.45t,
 2.47t, 2.51t, 2.56t, 2.62t, 2.74t
 Core excitons, 5.26t
 Cornell Black, 6.26f, 6.27
 Cornu equation, for refraction index, 2.22
 Correlated twin beams of light, 17.28, 17.29f,
 17.30f
 CORTAN glass, 2.43t, 2.49t, 2.54t, 2.59t, 2.67t
 Corundum, 2.70t
 Coulomb attraction, 8.31
 Coulomb explosions:
 cluster, 21.33–21.34
 molecular, 21.24–21.25
 Coulomb gauge, for solids, 8.5
 Coulomb potentials, 21.31, 21.32f
 Coupled plasmon-phonon behavior, 5.35, 5.36,
 5.36f, 5.37f
 Coupled resonator structures (CRS), 22.11–22.13,
 22.12f
 Coupling:
 in-plane, 9.10–9.11
 out-of-plane, 9.11–9.12
 two-beam, 12.4–12.6, 12.4f, 13.8–13.9
 Coupling laser power, 14.14–14.15
 Coupling (tie) layer, of bandpass filters, 7.82
 Coupling schemes, 14.30f (*See also specific
 coupling schemes, e.g.:* Lambda coupling)
 CR-39 resin (poly-diallylglycol), 3.11
 Craters, surface absorption and, 6.15
 Creep strength, of metals, 4.8
 Cross-Kerr phase shifter, 23.11
 Crown glasses, 2.28, 2.41t, 2.42t
 Crystals:
 anisotropic, 8.8–8.11, 8.9t, 8.10f
 biaxial, 8.8, 8.9t, 8.10, 8.10f
 and dielectric tensor and optical indicatrix,
 2.17–2.19, 2.19f
 and dispersion formulas for refractive index,
 2.21–2.22
 and glasses, 2.1–2.77
 lattice vibration model parameters for,
 2.76t–2.77t
 material properties of, 2.27–2.36
 characteristic temperatures, 2.32, 2.33
 combinations of, 2.36
 correlations of, 2.36
 elastic properties, 2.30–2.31, 2.31t
 hardness and strength, 2.31–2.32, 2.32f,
 2.32t
 heat capacity and Debye temperature,
 2.33–2.34
 material designation and composition,
 2.27, 2.29f
 naming of, 2.27
 thermal conductivity, 2.35–2.36, 2.35f
 thermal expansion, 2.34–2.35, 2.34f
 unit cell parameters, molecular weight,
 and density, 2.30
 mechanical properties of, 2.47t–2.48t
 nonlinear optical, 10.19–10.20, 10.20t–10.22t
 and nonlinear optical coefficients, 2.26–2.27,
 2.27t
 optical activity of, 8.10–8.11
 optical applications of, 2.17–2.27
 as optical materials, 2.4–2.5
 optical modes of, 2.68t–2.76t
 with cesium chloride structure, 2.68t
 with chalcopyrite structure, 2.74t
 with corundum structure, 2.70t
 with cubic perovskite structure, 2.73t
 with diamond structure, 2.68t
 with fluorite structure, 2.69t
 other structures, 2.74t–2.76t

- Crystals, optical modes of (*Cont.*):
 with α -quartz structure, 2.71*t*
 with rutile structure, 2.71*t*
 with scheelite structure, 2.72*t*
 with sodium chloride structure, 2.69*t*
 with spinel structure, 2.73*t*
 with tetragonal perovskite structure, 2.73*t*
 with trigonal selenium structure, 2.70*t*
 with wurtzite structure, 2.70*t*
 with zinblende structure, 2.69*t*
- optical properties of, 2.6, 2.8–2.9
 origin and models of, 2.9–2.17, 2.10*f*
 absorption in transparent region, 2.17
 electronic transitions, 2.12–2.15, 2.13*f*
 lattice vibrations, 2.11–2.12
 multiphoton absorption and refraction,
 2.15–2.17, 2.16*f*, 2.17*f*
- and photoelastic coefficients, 2.24
 physical properties of, 2.37, 2.38*t*–2.43*t*
 classes and symmetry properties, 2.7*t*
 composition, structure, and density,
 2.38*t*–2.41*t*
 physical constants, 2.8*t*
 symmetry properties, 2.5, 2.6*t*–2.8*t*
- Raman scattering in, 8.19*t*–8.20*t*
 room-temperature dispersion formulas,
 2.60*t*–2.66*t*
- room-temperature elastic constants of
 cubic crystals, 2.44*t*–2.49*t*
 hexagonal crystals, 2.46*t*
 monoclinic crystals, 2.47*t*
 orthorhombic crystals, 2.46*t*
 tetragonal crystals, 2.44*t*–2.45*t*
- and scatter, 2.27
 thermal properties, 2.50*t*–2.53*t*
 and thermo-optic coefficients, 2.24–2.26
 and total power law, 2.19–2.20, 2.20*f*
 uniaxial, 8.8, 8.9*t*, 8.10*f*
 (*See also Solids*)
- Cubic crystals, 8.9*t*, 8.20*t*
 dielectric constants of, 2.18
 room-temperature elastic constants of,
 2.44*t*–2.49*t*
 symmetries of, 2.7*t*
- Cubic oxides (sillenites), 12.17–12.19, 12.18*t*,
 12.19*f*
- Cubic perovskite structure, of crystals and
 glasses, 2.73*t*
- Cubic zirconia ($\text{ZrO}_2\cdot 0.12\text{Y}_2\text{O}_3$), 2.41*t*, 2.44*t*,
 2.48*t*, 2.69*t*
- Cumulative size distribution, of particles in
 water, 1.15
- Curie temperature, of crystals and glasses, 2.33
- Cutoff, slope of, 7.56
- Cutoff filters, 7.53–7.60, 7.54*f*, 7.55*f*, 7.57*f*,
 7.59*f*–7.61*f*
- Cw modelocking, 18.5*f*
- Cyclotron resonance (CR), 5.12, 5.12*f*, 5.40,
 5.47–5.50, 5.48*f*–5.50*f*
- Damage, laser-induced [*see Laser-induced
 damage (LID)*]
- Damping:
 in laser cooling, 20.19–20.20, 20.20*f*
 phonon, 5.14
- Dark (noncoupled) states, 14.4, 14.6–14.7,
 20.37–20.39, 20.38*f*
- Data storage, photorefractive holographic, 12.37
- De-Broglie wavelengths, 8.4
- Debye molar heat capacity, 2.33–2.34
- Debye temperature, for crystals and glasses,
 2.33–2.34
- Decay, homogeneous, 11.10, 11.11*f*
- DEEP SPACE BLACK, 6.49
- Defect modes (in photonic crystals), 22.11,
 22.12*f*
- Defect-related absorption, 5.37–5.39, 5.38*f*,
 5.39*f*
- Defocusing:
 ionization-induced, 21.43–21.44, 21.43*f*
 self, 13.7, 13.8, 19.9–19.11, 19.10*f*
 thermal, 13.8
- Degenerate four-wave mixing (DFWM),
 16.27–16.28, 16.28*f*, 18.17
- Dense crown flint glass, 2.42*t*
- Dense flint glass, 2.43*t*
- Dense phosphate crown glass, 2.41*t*
- Density matrix, 16.21, 16.22
- Dephasing, 11.15, 14.12–14.14
- Deposition method, of manufacturing
 thin-films, 7.11–7.12
- Designer blacks, 6.14
- DeSoto Black, 6.37*f*, 6.39
- Detectors, semiconductor, 5.61
- Detritus, organic:
 absorption by, 1.25–1.27, 1.25*t*, 1.26*f*
 in water, 1.14
- Detuning, 14.11*f*
- Diamond (crystal), 2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*,
 2.55*t*, 2.61*t*, 2.68*t*

- Diamond lattice structure, 5.6, 5.16, 5.17*t*
of air spheres, 9.5
of crystals and glasses, 2.68*t*
- Diamond-turning process, 3.12
- Dichroic beam splitters, 7.65–7.66, 7.67*f*
- Dicke narrowing, 15.9
- Dielectric broadband reflectors, 7.39, 7.40*f*,
7.45–7.47, 7.45*f*–7.47*f*
- Dielectric color selective (dichroic) beam
splitters, 7.65–7.66, 7.67*f*
- Dielectric constant, 2.14, 2.17–2.18
of crystals and glasses, 2.6*t*
dispersion of, 8.22*f*
for solids, 8.15
- Dielectric function, 4.26*f*, 5.8–5.9, 5.12
- Dielectric multiple reflection cutoff filters,
metal, 7.59–7.60
- Dielectric multiple reflection filters, metal, 7.111
- Dielectric potential, 9.15
- Dielectric reflectors:
interference filters with, 7.81
metal, 7.81–7.82, 7.108, 7.109, 7.110*f*
- Dielectric tensor, of crystals and glasses, 2.17–2.18
- Differential reflectivity (DR), 18.1, 18.5, 18.6, 18.6*f*
- Differential transmission (DT), 18.1
- Differential transmission (DT) spectroscopy,
18.18–18.19
- Diffraction grating, 5.59–5.60
- Diffuse attenuation:
and Jerlov water types, 1.42–1.46, 1.43*t*–1.45*t*,
1.44*f*, 1.45*f*
in water, 1.13
- Diffuse attenuation coefficient, 1.12
- Diffuse Infrared Background Explorer, 6.29*f*
- Diffusely scattering surfaces (*see* Black surfaces)
- Diluted magnetic semiconductors (DMSs), 5.45
- Dipole active modes, infrared, 8.16–8.18, 8.17*f*
- Dipole force, 20.8
- Direct excitons, 5.29
- Direct interband absorption, in solids, 8.27–8.28
- Direct (vertical) interband absorption transitions,
5.22–5.23, 5.22*f*–5.23*f*, 5.25*f*
- Direct laser acceleration (DLA), 21.43
- Direct pulse generation, 18.4
- Dispersion, 18.11
Cauchy, 2.21
for crystals and glasses, 2.21–2.23, 2.60*t*–2.66*t*,
2.66*t*–2.68*t*
Drude, 2.21–2.22
Maxwell-Helmholtz-Drude, 2.12, 2.21–2.22
- Dispersion (*Cont.*):
in multilayer reflectors, 7.40
normal vs. anomalous, 4.4
principle of, 2.23
range of anomalous, 4.4
room-temperature, 2.60*t*–2.66*t*, 2.66*t*–2.68*t*
in solids, 8.14–8.16, 8.22*f*
Zernike, 2.22
- Dispersion spectrometers, 5.59–5.60
- Displays, out-of-plane coupling and, 9.11–9.12
- Dissipative force, 20.7
- Dissolved substances, in water, 1.13
- Distributed Bragg reflector (DBR) diode laser,
17.9, 17.10*f*, 17.11*f*
- Distributed feedback (DFB) lasers, 9.6
- Distribution functions, of water, 1.6*t*, 1.12
- Documentation, for polymers, 3.5
- Donor-acceptor pair (DAP) transition, 5.71
- Doppler cooling, 20.13–20.15, 20.14*f*
- Doppler cooling limit, 20.15
- Doppler limit, 20.5
- Doppler shift, 11.19, 20.4
- Doppler temperature, 20.5, 20.11, 20.15
- Double atomic resonance, 22.5–22.9, 22.6*f*–22.8*f*
- Double ionization, strong field, 21.18–21.19,
21.19*f*
- Double phase conjugate mirrors, 12.7, 12.8*f*
- Double photonic resonance, 22.11–22.13, 22.12*f*
- Double-lambda coupling, 14.24*f*
- Double-sided Feynman diagrams, 11.12–11.13,
11.13*f*
- Doubly resonant optical parametric oscillators
(DR OPOs), 10.18
- Doubly resonant oscillators (DROs), 17.2–17.4,
17.3*f*, 17.4*f*, 17.16–17.17
- Down-conversion, spontaneous parametric,
23.10, 23.13
- Downwelling average cosine, of water, 1.6*t*,
1.7*f*, 1.12
- Downwelling diffuse attenuation coefficient,
for sea water, 1.44*t*–1.45*t*
- Downwelling irradiance, of water, 1.5*t*, 1.7*f*, 1.8
- Downwelling irradiance diffuse attenuation
coefficients, for sea water, 1.43*t*
- Downward scalar irradiance, of water, 1.5*t*, 1.7*f*, 1.8
- Dressed states, 14.3
- Drift velocity, 21.7
- Drude approximation, 5.34
- Drude dispersion formula, for crystals and
glasses, 2.21–2.22

- Drude model, 4.4–4.5, 8.15, 8.21, 8.22f, 16.20
d-tensor, 10.11
 Dual-cavity PE-SROs, 17.18–17.20, 17.19f, 17.20f
 Dual-Fock states, 23.8–23.9
 Ductility, of metals, 4.8, 4.70
 DURACON, 6.55–6.56
 Dynamic ionization stabilization, 21.21
- Ebanol C coating, 6.56
 Echo signal, 11.11–11.12, 11.12f
 ECP-2200, 6.27, 6.29f
 ECP-2200 coating (*see* MH 2200 coating)
 Edge filters, nonpolarizing, 7.66, 7.67f
 Effective mass, 8.25–8.26, 8.26f
 E-field-dependent electronic polarizability, 19.5
 “Egg-crate” array, of potential wells, 20.32
 Ehrenfest theorem, 20.6
 Einstein-Smoluchowski theory of scattering, 1.30
 Elastic constants, for crystals:
 cubic crystals, 2.44t–2.49t
 hexagonal crystals, 2.46t
 monoclinic crystals, 2.47t
 orthorhombic crystals, 2.46t
 Elastic properties, of crystals and glasses,
 2.30–2.31, 2.31t
 Elastic stiffness, of metals, 4.69, 4.69t
 Elasto-optic coefficients, for crystals and
 glasses, 2.21
 Electric field amplitude, in multilayer systems,
 7.9–7.10
 Electric field reflectance, Fresnel expression
 for, 4.5
 Electrically conductive black paint, 6.56
 Electric-field-modulated reflection spectroscopy,
 5.66t, 5.67
 Electrodeposited surfaces, 6.8f, 6.53–6.54,
 6.54f, 6.55f
 Electroless nickel, etched, 6.5f, 6.6f, 6.49–6.50,
 6.50f, 6.51f, 6.53f
 Electromagnetic spectrum, semiconductor
 interactions with, 5.3–5.6, 5.4f
 Electromagnetically induced transparency
 (EIT), 14.1–14.36, 22.5–22.9, 22.7f, 22.8f
 coherence in two- and three-level atomic
 systems, 14.4–14.5, 14.4f
 and cw lasers, 14.16–14.18, 14.17f
 at few photon level, 14.32–14.33
 gain and lasing without inversion,
 14.18–14.19
 as interference effect, 14.2–14.4
- Electromagnetically induced transparency
 (EIT) (*Cont.*):
 manipulation of optical properties by,
 14.10–14.15, 14.11f, 14.13t
 coupling laser power, 14.14–14.15
 dephasing and fluctuations in laser fields,
 14.13
 dephasing in gas phase media,
 14.12–14.13
 dephasing in solids, 14.13–14.14
 inhomogeneous broadening, 14.14
 and maximal atomic coherence,
 14.28–14.32, 14.29f–14.32f
 nonlinear optical frequency conversion,
 14.24–14.28, 14.24f, 14.27f
 physical concept of, 14.5–14.10, 14.6f, 14.8f,
 14.9f
 pulse propagation effects, 14.20–14.22
 and pulsed lasers, 14.15–14.16, 14.16f
 and refraction index in dressed atoms,
 14.19–14.20
 in solids, 14.33–14.36, 14.35f, 14.36f
 ultraslow light pulses, 14.22–14.23, 14.23f
- Electron(s):
 relativistic, above threshold ionization,
 21.20, 21.21f
 strong field interactions with single, 21.5–21.10,
 21.7f, 21.9f
- Electron acceleration, 21.39–21.42, 21.40f
 Electron beams, strong field interactions with
 relativistic, 21.9–21.10
 Electron stochastic heating, 21.36
 Electron-hole drops, 5.26t
 Electron-hole pairs, 14.34, 16.31
 Electro-optic coefficients, for crystals and
 glasses, 2.21
 Electrorreflectance, 5.66t, 5.67
 Electrorefractive photorefractive (ERPR) effect,
 12.21
 Electrostriction, 16.18–16.19, 19.5
 Element wedge, of polymers, 3.10
 Ellipsometers and ellipsometry, 5.5, 5.57, 5.63,
 5.66t, 5.67–5.69, 5.68f, 5.69f
 Elongation, of metals, 4.70, 4.70t
 Embedded polarizers, 7.70–7.71, 7.71f, 7.72f
Emiliana huxleyi, 1.15
 Emittance:
 for crystals and glasses, 2.20
 of metals, 4.6, 4.49, 4.49f, 4.50t, 4.51t
 and surface coatings, 6.19, 6.20t

- Energy:
 flow of, in solids, 8.7–8.8
 Landau levels of, 5.40, 5.42f
- Energy bands:
 magnetic field effects on, 5.40
 for solids, 8.25–8.27, 8.26f
- Energy walk-off angle, 8.9
- Energy-time uncertainty principle, 23.4
- Engineering moduli, for crystals and glasses, 2.37
- Enhanced Martin Black, 6.46, 6.47
- Enhanced refraction, 14.20
- Entanglement, quantum (*see* Quantum entanglement, in optical interferometry)
- Entrance damage, laser-induced, 19.3
- Environmental degradation, of black surfaces, 6.16–6.18
- Environmentally responsible glass, 2.29–2.30
- Epner Laser Black, 6.56
- Etched electroless nickel surface, 6.5f, 6.6f, 6.49–6.50, 6.50f, 6.51f, 6.53f
- Etching, surface, 6.15
- Euphotic zone, of water, 1.46, 1.46t
- Evaporated spacers, 7.83, 7.84f, 7.85f
- Evaporation method, of manufacturing thin-films, 7.11
- Excitance, total integrated, 2.19
- Excitation(s):
 chirped pulse, 11.25–11.26
 continuum, 18.20
 excitonic, 18.19–18.20
 photoexcitation, 5.70f
 single-particle, 5.81, 5.82f
- Excited state absorption (ESA), 13.5, 16.19
- Excited state collisions, 20.29
- Exciton(s):
 free, luminescence, 5.72, 5.73f
 in semiconductors, 5.25–5.29, 5.26t, 5.27f–5.28f, 5.46
 and solids, 8.31–8.32
- Exciton gases, 5.26t
- Exciton Rydberg, 8.31
- Excitonic excitations, 18.19–18.20
- Excitonic magneto-optical effects, 5.46, 5.47f
- Excitonic molecules, 5.26t
- Exit damage, laser-induced, 19.3
- External pulse compression, 18.4, 18.11–18.12
- External self-action, 16.25
- Extinction coefficient, of metals, 4.3, 4.11, 4.12t–4.19t, 4.20f–4.26f
- Extreme environments, black surface degradation in, 6.18
- Extreme Ultraviolet Explorer, 6.21
- Extreme ultraviolet (XUV) light:
 bandpass filters for, 7.94–7.96, 7.95f–7.96f
 interference polarizers for, 7.73, 7.76f–7.77f
 multilayer reflectors for, 7.42–7.43, 7.53
- Extrinsic optical properties:
 of semiconductors, 5.11
 of solids, 8.3
- Fabry-Perot interference filters, 7.78–7.82, 7.79f, 7.80f, 7.92–7.94, 7.93f, 7.96
- Fabry-Perot interferometers, 7.13, 7.39, 7.39f, 7.40, 7.89
- Fabry-Perot resonators, 22.11, 22.12f
- Family momentum, 20.38
- Fano interferences, 14.2
- Faraday effect, 5.50–5.51
- Faraday rotation:
 free-carrier, 5.50–5.51
 interband, 5.44–5.45, 5.45f
- Faraday's law, 2.6
- Far-infrared (FIR) lasers, 5.48
- Far-infrared (FIR) radiation:
 and black surfaces, 6.21, 6.26–6.34, 6.28f–6.34f
 Ames 24E and 24E2, 6.26f, 6.28f, 6.34
 Cornell Black, 6.26f, 6.27
 Infrablack, 6.26f, 6.28, 6.28f
 multiple-layer approach, 6.26, 6.26f–6.27f
 Teflon overcoat, 6.27
 and EIT, 14.3
- Far-infrared region, multilayer reflectors for, 7.52, 7.52f
- Far-infrared (FIR) telescopes, 6.48
- Far-off-resonance traps (FORTs), 20.23
- Fast ignition, 21.54–21.55, 21.55f
- Fast saturable absorbers, 18.8f, 18.9–18.10
- Fatigue strength, of metals, 4.8
- Femtosecond, 5.7
- Femtosecond x-ray production, 21.52–21.53, 21.53f
- Fermi level, 8.21
- Fermi's golden rule, 8.25
- Ferroelectric oxides, 12.13–12.14, 12.13t
- Ferroelectric photorefractive materials, 12.13–12.17, 12.13t
 barium titanate, 12.15–12.16, 12.16f
 lithium niobate and lithium tantalate, 12.14
 potassium niobate, 12.16–12.17

- Ferroelectric photorefractive materials (*Cont.*):
 strontium barium niobate and related compounds, 12.17
 tin hypothydiphosphate, 12.17, 12.18*t*
- FF5 glass (593355), 2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- Fiber feedback, 17.16
- Fiber pump lasers, 17.7–17.11, 17.9*f*–17.11*f*
- Fiber Raman amplifiers, 22.15
- Fibers:
 photonic bandgap, 2.23
 slow light propagation in, 22.13–22.15, 22.14*f*
- Films:
 and black surfaces, 6.15
 polymer, 12.23–12.25, 12.26*t*–12.27*t*
 semiconductor-doped dielectric, 18.11
 thin (*see* Thin-film coatings)
- Filters and filtering, 7.1–7.114
- antireflection coatings, 7.15–7.32
 of absorbing and amplifying media, 7.26, 7.27
 homogeneous-layer, 7.16–7.23, 7.17*f*–7.19*f*, 7.20*t*–7.21*t*, 7.22*f*–7.23*f*
 inhomogeneous and structured, 7.23–7.26, 7.24*f*, 7.26*f*
 at nonnormal angle of incidence, 7.28–7.31, 7.28*f*–7.31*f*
 nonoptical properties of, 7.31–7.32, 7.32*f*
 surface reflections and optical performance, 7.15–7.16, 7.16*f*
 of surfaces carrying thin films, 7.27–7.28, 7.28*f*
 universal, 7.26, 7.27*f*
- bandpass, 7.73–7.96
 about, 7.73, 7.77–7.78, 7.77*f*–7.78*f*
 angular properties of, 7.91–7.94, 7.92*f*, 7.93*f*
 with multiple peaks, 7.90, 7.91*f*
 narrow- and medium-, 7.78–7.83, 7.79*f*, 7.80*f*, 7.82*f*–7.88*f*
 stability and temperature dependence of, 7.94
 very narrow, 7.83, 7.88–7.89, 7.89*f*
 wedge filters, 7.90, 7.91, 7.91*f*
 wide, 7.90, 7.90*f*
 for XUV and x-ray regions, 7.94–7.96, 7.95*f*–7.96*f*
- beam splitters, 7.61–7.67, 7.62*f*–7.68*f*
 achromatic beam splitters, 7.62–7.65, 7.62*f*–7.65*f*
 color-selective beam splitters, 7.65–7.66, 7.66*f*, 7.67*f*
 geometrical considerations for, 7.61–7.62
- Filters and filtering (*Cont.*):
 with coatings
 measurements on, 7.12–7.14
 transmission and reflection of, 7.3
 cutoff, heat-control, and solar-cell cover, 7.53–7.60
 cutoff filters, 7.53–7.60, 7.54*f*, 7.55*f*, 7.57*f*, 7.59*f*–7.61*f*
 heat reflectors, 7.58
 solar-cell cover filters, 7.58
 high performance optical multilayer coatings, 7.96–7.98, 7.97*f*
 interference polarizers and polarizing beam splitters, 7.69–7.73, 7.70*f*–7.72*f*, 7.76*f*–7.77*f*
 with low reflection, 7.104–7.106, 7.104*f*–7.105*f*
 matched, 12.28–12.29, 12.29*f*, 12.30*f*
 multilayer reflectors, 7.39–7.53
 all-dielectric broadband reflectors, 7.39, 7.40*f*, 7.45–7.47, 7.45*f*–7.47*f*
 coatings for ultrafast optics, 7.47–7.48, 7.48*f*
 for far-infrared region, 7.52, 7.52*f*
 graded reflectivity mirrors, 7.52
 imperfections in, 7.40–7.43, 7.41*f*–7.43*f*
 for interferometers and lasers, 7.39–7.40, 7.39*f*–7.40*f*
 narrowband reflection coatings, 7.43, 7.44*f*
 rejection filters, 7.48–7.50, 7.49*f*–7.51*f*
 resonant reflectors, 7.43–7.45, 7.44*f*
 for soft x-ray and XUV regions, 7.53
 neutral filters, 7.67, 7.67*f*–7.68*f*
 novelty, 12.32, 12.33*f*–12.35*f*
 phase coatings, 7.101, 7.101*f*–7.104*f*, 7.102
 reflection, 7.5, 7.5*f*
 reflection coatings and, 7.106–7.113, 7.106*f*–7.113*f*
 special purpose coatings, 7.113–7.114, 7.114*f*
 theory of, 7.1*f*, 7.2
 thin-film coatings
 and antireflection coatings, 7.27–7.28, 7.28*f*
 manufacturing of, 7.10–7.12
 of metal, 7.104, 7.104*f*
 theory and design of, 7.5–7.10, 7.6*f*, 7.9*f*
 transmission, 7.3–7.5, 7.4*f*
 for two or three spectral regions, 7.98–7.100, 7.98*f*–7.101*f*
 two-material periodic multilayers theory for, 7.32–7.38, 7.33*f*–7.38*f*

- Filtrate absorption, 1.21
- Finite-difference time-domain (FDTD) solution
(to Maxwell's equations), 9.3
- Fissures, and surface absorption, 6.15
- Flame-sprayed aluminum, 6.57
- Flashlamp pumped Nd:glass amplifiers, 21.5
- Flexural strength, of metals, 4.70, 4.70*t*
- Flint glasses, 2.28, 2.41*t*–2.43*t*
- Floquet theory, 21.23
- Floquil, 6.44
- Fluor crown glass, 2.41*t*
- Fluorescence, chlorophyll, 1.49
- Fluorite (CaF₂), 2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*,
2.61*t*, 2.69*t*, 2.77*t*
- Fluoro flint glass, 2.42*t*
- Fock states, 23.8–23.9, 23.14
- Focused beams, 15.41
- Fokker-Planck equation, 20.10–20.11
- Forbidden bands, 9.2
- Forward drift velocity, 21.7
- Forward Raman amplifiers, 15.4, 15.4*f*
- Forward Raman generators, 15.4, 15.4*f*
- Forward-looking infrared (FLIR)
systems, 6.54
- Fourier-transform spectrometers, 5.60–5.61,
5.60*f*, 5.72, 7.12
- Four-photon absorption (4PA), 19.10
- Four-wave mixing (FWM), 18.3
coherent anti-Stokes, 15.2*t*, 15.3*t*, 15.4, 15.4*f*
degenerate, 18.17
and EIT, 14.24–14.26, 14.24*f*, 14.27*f*
resonant, 14.28
and third-order optical nonlinearities,
16.27–16.28, 16.28*f*
transient, 18.17–18.18, 18.17*f*
- Four-wave mixing phase conjugation, 12.6–12.7,
12.6*f*
- Fracture toughness:
of crystals and glasses, 2.32, 2.32*t*
of metals, 4.8, 4.70, 4.70*t*
- Franz-Keldysh effect, 12.22
- Fredkin gate, 23.11
- Free carriers and free-carrier effects, 5.47–5.52
in crystals, 2.15
cyclotron resonance, 5.47–5.50, 5.48*f*–5.50*f*
free-carrier Faraday rotation, 5.50–5.51
impurity magnetoabsorption, 5.51–5.52,
5.51*f*
in semiconductors, 5.33–5.36, 5.35*f*–5.37*f*,
5.81, 5.82*f*
- Free electron properties, of solids, 8.21–8.24
Drude model, 8.21, 8.22*f*
interband transitions in metals, 8.21
plasmons, 8.23–8.24
reflectivity, 8.23
- Free polarization decay (FPD), 11.7–11.11,
11.8*f*, 11.10*f*, 11.11*f*
- Free spectral range, of bandpass filters, 7.77
- Free-carrier Faraday rotation, 5.50–5.51
- Free-exciton (FE) luminescence, 5.72, 5.73*f*
- Frenkel excitons, 5.26, 5.26*t*
- Frequency conversion, nonlinear optical,
14.24–14.28, 14.24*f*, 14.27*f*
- Frequency mixing, 16.3*t*
- Fresnel reflection coefficient, for metallic
reflectors, 7.106
- Fresnel relations, 8.11, 8.15
- Fringe power, of polymers, 3.11
- Frogs legs mirror, 12.7, 12.8*f*
- Full width half maximum (FWHM), 18.3
- Fully functional polymers, 12.27*t*
- Fundamental absorption edge, 5.21
absorption near, 5.21–5.22, 5.21*f*
high-energy transitions above, 5.29–5.33,
5.30*f*–5.34*f*
- Fused glass, 2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- Fusion neutron production, 21.53
- Gain narrowing, in steady-state Stokes
scattering, 15.21
- Gain saturation, dynamic, 18.8, 18.8*f*
- Gain without inversion, and EIT, 14.18–14.19
- Gallagher-Pritchard (GP) model (of trap-loss
collisions), 20.29
- Gallium arsenide (GaAs):
composition, structure, and density of, 2.39*t*
dispersion formulas for, 2.62*t*
elastic constants of, 2.44*t*
lattice vibration model parameters for, 2.77*t*
linear-chain model calculations for, 5.19*f*
local vibrational modes for, 5.19*f*
luminescence in, 5.72, 5.73, 5.74*f*
mechanical properties of, 2.47*t*
multiphonon absorption of vacuum-grown,
5.18*f*
optical modes of, with zincblende
structure, 2.69*t*
optical properties of, 2.56*t*
Raman scattering of, 5.80, 5.81, 5.81*f*
thermal properties of, 2.51*t*

- Gallium nitride (GaN), 2.39t, 2.46t, 2.47t, 2.51t, 2.56t, 2.62t, 2.70t, 2.77t
- Gallium phosphide (GaP), 2.39t, 2.44t, 2.47t, 2.51t, 2.56t, 2.62t, 2.69t, 2.77t
- Gap modes (GMs), 5.17
- Gas phase media, dephasing in, 14.12–14.13
- Gases, strong field nonlinear optics in, 21.27–21.31, 21.28f
- Gauss law, 2.6
- Gelbstoff, in water, 1.13 (*See also* Yellow matter)
- Generalized Rabi frequency, 11.4
- Geometrical configuration factor (GCF), 6.12
- Germania glass, fused, 2.43t, 2.49t, 2.54t, 2.59t, 2.67t
- Germanium:
 - absorptance of, 4.48t
 - in crystal form, 2.39t, 2.44t, 2.47t, 2.51t, 2.56t, 2.62t, 2.68t
 - thermal properties of
 - elastic stiffness, 4.69t
 - moduli and Poisson's ratio, 4.69t
 - strength and fracture properties, 4.70t
- Ghost imaging, 23.13
- "Ghosts" (in gratings), 5.60
- Gilvin, 1.13 (*See also* Yellow matter)
- Glass(es):
 - amorphous, 12.26t
 - antireflection coatings for, 7.26–7.28, 7.28f
 - common, 2.3
 - and crystals, 2.1–2.77
 - defined, 2.33
 - hybrid organic-inorganic, 12.27t
 - material properties of, 2.27–2.36
 - characteristic temperatures, 2.32, 2.33
 - combinations of, 2.36
 - correlations of, 2.36
 - elastic properties, 2.30–2.31, 2.31t
 - hardness and strength, 2.31–2.32, 2.32f, 2.32t
 - heat capacity and Debye temperature, 2.33–2.34
 - material designation and composition, 2.28–2.30, 2.29f
 - naming of, 2.27
 - thermal conductivity, 2.35–2.36, 2.35f
 - thermal expansion, 2.34–2.35, 2.34f
 - unit cell parameters, molecular weight, and density, 2.30
 - mechanical properties of, 2.49t
- Glass(es) (*Cont.*):
 - optical applications of, 2.17–2.27
 - dielectric tensor and optical indicatrix, 2.17–2.19, 2.19f
 - dispersion formulas, 2.21–2.23
 - nonlinear optical coefficients, 2.26–2.27, 2.27t
 - scatter, 2.27
 - thermo-optic coefficients, 2.24–2.26, 2.24f
 - total power law, 2.19–2.20, 2.20f
 - as optical materials, 2.4–2.5
 - optical properties of, 2.6, 2.8–2.9
 - lattice vibration model parameters, 2.76t–2.77t
 - origin and models, 2.9–2.17, 2.10f, 2.13f, 2.16f, 2.17f
 - room-temperature dispersion formulas, 2.66t–2.68t
 - summary table, 2.59t
 - physical properties of, 2.37, 2.38t–2.43t
 - optical glass reference table, 2.41t–2.43t
 - physical constants, 2.8t
 - specialty, and substrate materials, 2.43t
 - symmetry properties, 2.5, 2.6t, 2.8t
 - thermal properties of, 2.54t
- Glass scatterers, 13.8
- Glass-based lasers, 21.5
- Glass-ceramics, 2.33
- Goddard Space Flight Center, 6.35, 6.39
- Gold:
 - absorptance of, 4.40f, 4.48t, 4.50t, 4.51t
 - optical properties of, 4.14t, 4.23f
 - physical properties of, 4.52t, 4.54t
 - reflectance of, 4.31t–4.32t, 4.40f
 - thermal properties of
 - coefficient of linear thermal expansion, 4.56t, 4.57f
 - elastic stiffness, 4.69t
 - moduli and Poisson's ratio, 4.69t
 - at room temperature, 4.55t
 - specific heat, 4.65t, 4.66f
 - strength and fracture properties, 4.70t
 - thermal conductivity, 4.58t, 4.60f–4.61f
- Gold black surfaces, 6.57
- Gold iridite, 6.21, 6.22f
- Gordon inequality, 1.21
- Graded reflectivity mirrors, 7.52
- Gradient force, 20.8
- Grain boundaries, of crystals, 2.4

- Gratings:
 Bragg, 22.9–22.11, 22.10f
 diffraction, 5.59–5.60
 formation of, 12.1–12.3, 12.2f
 Moiré, 22.11, 22.12f
 transmission, 12.7, 12.8f
- Greenockite (CdS), 2.38t, 2.46t, 2.47t, 2.51t, 2.56t, 2.61t, 2.70t, 2.77t
- Grinding, controlled, 19.3
- Gross-Pitaevski equation, 20.34
- Ground-based telescopes, 6.12
- Ground-state collisions, 20.29
- Group Delay Dispersion (GDD), 7.47–7.48, 7.48f
- Group Delay (GD) phase changes, 7.47, 7.48
- Group theory, 5.6
- Group velocity, 22.3–22.5
- Group velocity dispersion (GVD), 11.27, 22.4–22.5
- Grüneisen relationship, 2.34, 2.35
- GUERAP (stray light analysis program), 6.19
- Guided-wave nonlinear structures, 17.15–17.16
- Hafnium dioxide-yttrium oxide ($\text{HfO}_2\text{:Y}_2\text{O}_3$), 2.56t, 2.62t, 2.77t
- Halite (NaCl), 2.40t, 2.44t, 2.48t, 2.52t, 2.57t, 2.64t, 2.69t
- Hardness:
 of crystals and glasses, 2.31, 2.32f
 of metals, 4.8
 of polymers, 3.2–3.3
- Harmonic generation:
 high, 21.27–21.30
 harmonic yield and phase matching, 21.30
 quasi-classical model, 21.28f, 21.29–21.30
 from solid plasmas, 21.50–21.52, 21.51f
 third-order, 16.2, 16.3t
 in crystals, 16.14
 energy level diagrams for, 16.5f
 and semiconductors, 5.56
- Harmonic yield, 21.30
- Hartmann equation, for refraction index, 2.22
- Heat capacity (specific heat):
 of crystals and glasses, 2.6t, 2.33
 of metals, 4.7, 4.10t, 4.53, 4.55, 4.55t, 4.65t, 4.66f–4.69f
- Heat reflectors, 7.58
- Heat-induced lensing effect, 16.22
- Heating:
 cluster electron, 21.34, 21.35
 inverse Bremsstrahlung, 21.37, 21.37f
 $j \times B$, 21.49
- Heating (*Cont.*):
 self-heating, 17.12
 vacuum, 21.47f, 21.48–21.49
- Heisenberg limit, 23.6–23.7
- Heisenberg number-phase uncertainty relation, 23.4
- Heisenberg uncertainty principle (HUP), 23.6
- Herzberger equation, 2.22
- Hexagonal crystals:
 anisotropic, 8.9t
 dielectric constants of, 2.18
 room-temperature elastic constants, 2.46t
 symmetries of, 2.7t, 8.20t
- High emittance-low absorptance coatings, 6.16
- High harmonic generation (HHG), 21.27–21.30
 harmonic yield and phase matching, 21.30
 quasi-classical model, 21.28f, 21.29–21.30
 from solid plasmas, 21.50–21.52, 21.51f
- High magnetic field production, 21.53
- High performance optical multilayer coatings, 7.96–7.98, 7.97f
- High-energy transitions above fundamental edge, 5.29–5.33, 5.30f–5.34f
- High-Q cavities, 9.12
- High-reflectance zones, of multilayers, 7.36–7.37, 7.37f
- High-repetition short-pulse lasers, 11.26–11.27, 11.27f
- High-resistivity coatings, 6.56
- High-resolution Doppler-free spectroscopy, 17.27
- Hilbert transforms, 2.8–2.9
- Hole boring, 21.50, 21.50f
- Holographic optical elements (HOEs), 12.31
- Holographic storage:
 data, 12.37
 photorefractive, 12.36–12.37
- Holography, real-time, 12.28–12.29, 12.29f, 12.30f
- Homogeneity, in polymeric optics, 3.7
- Homogeneous broadening, spectral-line, 14.13
- Homogeneous decay, 11.10, 11.11f
- Homogeneous-layer antireflection coatings, 7.16–7.23, 7.17f–7.19f, 7.20t–7.21t, 7.22f, 7.23f
- Homogeneously broadened systems, 11.18–11.19, 11.19f
- Hong-Ou-Mandel effect, 23.10, 23.14
- Hooke's law, 8.14, 8.21
- Hot mirrors, 7.58
- Housings, lens, 3.15–3.16, 3.15f

- HTF-1 glass, 2.43t, 2.49t, 2.54t, 2.59t, 2.68t
 Hughes Airborne Optical Adjunct Coating, 6.49
 Human eye:
 active optical limiting by, 13.1
 damage thresholds for, 13.3, 13.3f
 Hybrid organic-inorganic composites, glasses,
 and sol-gels, 12.27t
 Hydrologic optics, 1.3 (*See also* Water)
 Hydrostatic pressure, 5.66t
 Hyperbolic cumulative size distribution, 1.15
 Hyper-Raman scattering, 15.2t, 15.3, 15.3t
 Hyperspectral imaging, 17.27–17.28
 Hysteresis instability, of metals, 4.10
- IBM Black, 6.56
 Illinois Institute of Technology, 6.35
 Imaging:
 ghost, 23.13
 hyperspectral, 17.27–17.28
 quantum, 23.13–23.14
 time-gated, 15.42, 15.43, 15.44f
 Impurity magnetoabsorption, 5.51–5.52, 5.51f
 Impurity-related absorption, 5.37–5.39,
 5.38f, 5.39f
 Impurity-related vibrational optic effects,
 5.17–5.20, 5.18f–5.19f, 5.20t, 5.21f
 Index contrast, 3D photonic crystals and, 9.4
 Index ellipsoid, of crystals and glasses, 2.18–2.19,
 2.19f
 Index grating, 12.1–12.3, 12.2f
 Index of absorption, 2.8
 Index of refraction (*see* Refractive index)
 Indirect (nonvertical) absorption transitions,
 5.22–5.24, 5.24f–5.25f
 Indirect excitons, 5.29
 Indirect interband transitions, 8.29–8.30, 8.29f
 Indium phosphide (InP), 12.21
 Induced transparency, 21.52
 Induced-transmission filters, 7.83, 7.88f
 Inelastic scattering:
 of light, 5.76–5.83, 5.76f, 5.78f–5.82f
 and polarization, 1.47–1.49, 1.48f, 1.49f
 Inertial confinement fusion (ICF), 21.54, 21.55f
 Infrablack, 6.26, 6.26f, 6.28, 6.28f, 6.48, 6.48f
 Infrared (IR) dipole active modes, 8.16–8.18,
 8.17f
 Infrared (IR) region:
 absorption in, 5.19f
 all-dielectric reflectors for, 7.39, 7.40f
 multilayer reflectors for far-, 7.52, 7.52f
 Infrared (IR) suppressing filters, 7.58, 7.58f
 Inherent optical properties (IOPs), of water,
 1.4, 1.5t, 1.9–1.12, 1.10f
 Inhomogeneous antireflection coatings,
 7.23–7.26, 7.24f–7.26f
 Inhomogeneous broadening, spectral-line,
 14.14
 Injection molding, of polymers, 3.2, 3.12–3.13
 Inner ionization, of cluster, 21.31–21.32, 21.32f
 Inorganic particles, in water, 1.14–1.15
 In-plane coupling, 9.10–9.11
 Instantaneous coefficient of linear thermal
 expansion, 4.7
 Instrumentation, spectroscopic, 5.58–5.61, 5.59f
 detectors, 5.61
 and light sources, 5.58–5.59
 broadband, 5.58–5.59
 laser, 5.59
 spectrometers and monochromators,
 5.59–5.61
 dispersion spectrometers, 5.59–5.60
 Fourier-transform spectrometers,
 5.60–5.61, 5.60f, 5.72
 “Intelligent” state, 23.8
 Intense laser pulses:
 cluster interactions with, 21.35–21.36, 21.36f
 plasma instabilities driven by, 21.38–21.39,
 21.38f
 Interacting beams, propagation of, 16.3
 Interband absorption, of semiconductors,
 5.21–5.33
 absorption near fundamental edge,
 5.21–5.22, 5.21f
 direct transitions, 5.22–5.23, 5.22f–5.23f
 excitons, 5.25–5.29, 5.26t, 5.27f–5.28f
 high-energy transitions above fundamental
 edge, 5.29–5.33, 5.30f–5.34f
 indirect transitions, 5.23–5.24, 5.24f–5.25f
 near fundamental edge, 5.21–5.22, 5.21f
 polaritons, 5.29
 Interband magneto-optical effects, 5.42–5.46, 5.43f
 excitonic, 5.46, 5.47f
 Faraday rotation, 5.44–5.45, 5.45f
 magnetoreflexion, 5.43, 5.44, 5.44f
 Interband transitions, of solids, 8.27–8.32
 direct interband absorption, 8.27–8.28
 excitons, 8.31–8.32
 indirect transitions, 8.29–8.30, 8.29f
 joint density of states, 8.28, 8.29f
 in metals, 8.21

- Interband transitions, of solids (*Cont.*):
 multiphoton absorption, 8.30–8.31
 selection rules and forbidden transitions,
 8.28, 8.29
- Interference effect, EIT as, 14.2–14.4
- Interference filters, Fabry-Perot, 7.78–7.82,
 7.79f, 7.80f, 7.92–7.94, 7.93f, 7.96
- Interference polarizers, 7.69–7.73, 7.70f–7.72f,
 7.76f–7.77f
- Interfering transition pathways, 14.2–14.3
- Interferometers and interferometry:
 atom, 11.22–11.23, 11.24f
 Fabry-Perot, 7.13, 7.39, 7.39f, 7.40, 7.89
 Mach-Zehnder, 23.2–23.4, 23.2f
 Michelson, 7.42, 7.104f
 multilayer reflectors for, 7.39, 7.39f
 phase conjugate, 12.32, 12.33f, 12.34f
 quantum entanglement in, 23.1–23.15
 concepts and equations for, 23.1–23.4,
 23.2f, 23.4f
 digital approaches, 23.7–23.9
 Heisenberg limit, 23.6–23.7
 NOON state, 23.9–23.12, 23.10f, 23.11f
 and quantum imaging, 23.13–23.14
 and remote sensing, 23.14–23.15
 shot-noise limit, 23.4–23.6, 23.5f
 and third-order optical nonlinearities,
 16.28–16.29
 time-domain atom, 11.22–11.24, 11.24f
- Internal self-action, 13.7, 16.25
- International Association for Physical Sciences
 of the Ocean (IAPSO), 1.4, 1.5t–1.6t
- Intervally scattering, 18.20
- Intraband magneto-optical effects, 5.47–5.52
 cyclotron resonance, 5.47–5.50, 5.48f–5.50f
 free-carrier Faraday rotation, 5.50–5.51
 impurity magnetoabsorption, 5.51–5.52,
 5.51f
- Intracavity singly resonant oscillators
 (IC-SROs), 17.3f
- Intrinsic optical properties:
 of semiconductors, 5.11
 of solids, 8.3
- Invar 36, 4.10t, 4.52t, 4.55t, 4.69t, 4.70t
- Inverse Bremsstrahlung heating, 21.37, 21.37f
- Inverse dielectric tensor, 2.6t, 2.19
- Ion beam-sputtered surfaces, 6.53
- Ion-assisted deposition, 7.11
- Ion-beam sputtering, 7.11, 7.14
- Ion-implanted semiconductors, 18.21
- Ionization:
 above-threshold (*see* Above threshold
 ionization)
 atomic, 21.3
 barrier suppression, 21.14, 21.14f
 in clusters, 21.31–21.33, 21.32f
 collisional, 21.31, 21.32
 double, 21.18–21.19, 21.19f
 inner, 21.31–21.32, 21.32f
 molecular tunnel (*see* Molecular tunnel
 ionization)
 multiphoton, 21.10–21.12, 21.11f
 outer, 21.32–21.33
 stabilization of, 21.20–21.21, 21.22f
 threshold (*see* Threshold ionization)
 tunnel (*see* Tunnel ionization)
- Ionization distance, 21.25–21.26, 21.25f, 21.27f
- Ionization ignition, 21.31
- Ionization rate, ADK, 21.13
- Ionization stabilization, 21.20–21.21, 21.22f
- Ionization-induced defocusing, 21.43–21.44,
 21.43f
- Ionized arsenic antisite (As_{Ga}^+), 18.3
- Ion-plating process, 7.11
- IRG 2 glass, 2.43t, 2.49t, 2.54t, 2.59t, 2.67t
- IRG 9 glass, 2.43t, 2.49t, 2.54t, 2.59t, 2.67t
- IRG 11 glass, 2.43t, 2.49t, 2.59t, 2.68t
- IRG 100 glass, 2.43t, 2.49t, 2.54t, 2.59t, 2.68t
- Iris (in eye), 13.1
- Iron:
 absorptance of, 4.40f, 4.48t, 4.50t
 optical properties of, 4.15t, 4.23f
 physical properties of, 4.54t
 reflectance of, 4.32t–4.33t, 4.40f
 thermal properties of
 coefficient of linear thermal expansion,
 4.56t, 4.57f
 elastic stiffness, 4.69t
 moduli and Poisson's ratio, 4.69t
 at room temperature, 4.55t
 specific heat, 4.65t, 4.67f
 thermal conductivity, 4.58t, 4.62f–4.63f
- Irradiance, of water, 1.5t, 1.8–1.9
- Irradiance reflectance (irradiance ratio), of
 water, 1.6t, 1.7f, 1.12, 1.46–1.47, 1.47f
- Irradiated plasma, 21.46–21.47, 21.46f
- Isotropic crystals:
 dielectric constants of, 2.18
 symmetries of, 2.7t
- Isotropic solids, 8.8, 8.9t

- Iturriaga R., 1.25
 -ivity suffix, 4.5
- $\mathbf{j} \times \mathbf{B}$ heating, 21.49
- Jerlov water types, 1.42–1.46
- Joint density of states, 8.28, 8.29*f*
- Jumps, electric-field, 9.6
- Junge (hyperbolic) cumulative size distribution, 1.15
- JV model (of trap-loss collisions), 20.30
- K (optical constant of water), 1.17, 1.17*f*
- Kane momentum matrix elements, 8.26
- Keldysh parameter, 21.10
- Keldysh-Faisal-Reiss (KFR) theories, 21.12
- Kerr effect, 18.11–18.15
 longitudinal, 18.11–18.15, 18.12*f*
 optical, 16.11–16.14
 Raman-induced, 16.3*t*, 16.12, 16.17
 transverse, 18.11, 18.14–18.15
- Kerr lens modelocking (KLM), 16.25, 18.3, 18.14–18.15
- Kerr-lensing (*see* Self-focusing)
- Kerr-type nonlinearity, 14.33
- Kirchhoff's law, 4.6
- Kitaev phase estimation algorithm, 23.12
- Kleinman \mathbf{d} -tensor, 10.11
- Knill-LaFlamme-Milburn (KLM) scheme, 23.12
- Knoop test, 2.31, 2.32*f*
- Kopelevich model of absorption, 1.28
- Kramers-Kronig (K-K) relations:
 dielectric-constant, 2.9, 2.12, 2.22
 dispersion, 5.10–5.11, 16.1, 16.9–16.11
 for solids, 8.15
- KRS-5 crystals, 2.40*t*, 2.44*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.65*t*
- KRS-6 crystals, 2.40*t*, 2.44*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.65*t*
- Ladder coupling, 14.1, 14.6*f*, 14.24*f*
- Lambda coupling, 14.1, 14.6*f*, 14.8*f*, 14.9*f*, 14.26, 14.27
- Lambertian black surface, 6.14
- Landau levels (of energy), 5.40, 5.42*f*
- Langmuir waves, 21.38
- Lanthanum glass, 2.42*t*, 2.43*t*
- Laser:
 continuous-wave, 7.14, 14.16–14.18, 14.17*f*
 multilayer reflectors for, 7.39–7.40, 7.39*f*
 theoretical description of, 20.6–20.11
 atoms in motion, 20.8–20.10, 20.9*f*
 Fokker-Planck equation, 20.10–20.11
 force on two-level atom, 20.6–20.7
- Laser Black, 6.56
- Laser conditioning, 19.4
- Laser cooling, 20.3–20.21, 20.26–20.39
 about, 20.3–20.4
 in atomic beam brightening, 20.27–20.28, 20.27*f*
 in atomic clocks, 20.28
 below Doppler limit, 20.17–20.21, 20.18*f*–20.20*f*
 in Bose-Einstein condensation, 20.35–20.37, 20.36*f*
 in dark states, 20.37–20.39, 20.38*f*
 defined, 20.3
 history of, 20.3–20.4
 in optical lattices, 20.31–20.34, 20.32*f*–20.34*f*
 by optical molasses, 20.13–20.17, 20.14*f*–20.16*f*
 properties of, 20.4–20.6
 Sisyphus, 20.19
 by slowing of atomic beams, 20.11–20.13, 20.12*f*, 20.12*t*, 20.13*f*
 in ultracold collisions, 20.28–20.31, 20.30*f*, 20.31*f*
- Laser damage threshold (LDT), of coatings, 7.13–7.14, 7.18
- Laser fields, 14.13
- Laser Interferometer Gravitational Wave Observatory (LIGO), 23.1, 23.7
- Laser Light Detection and Ranging (LIDAR) systems, 23.1
- Laser light sources, 5.59
- Laser power combining, 12.30, 12.31
- Laser pulses, 21.35–21.36, 21.36*f*, 21.38–21.39, 21.38*f*
- Laser technology, for strong field interactions, 21.4–21.5, 21.4*f*
- Laser-induced breakdown (LIB), 19.6–19.9, 19.8*f*
- Laser-induced damage (LID), 19.1–19.11
 avoidance of, 19.5–19.6
 and critical NLO parameters, 19.9–19.11, 19.9*f*, 19.10*f*
 estimates of, 19.2
 mechanisms of, 19.6–19.9, 19.8*f*
 and nonlinear optical effects, 19.5
 package-induced, 19.4–19.5
 surface damage, 19.2–19.4
- Lasing without inversion (LWI), 14.1, 14.3–14.4, 14.18–14.19

- Lattice absorption, semiconductor, 5.13–5.20
 impurity-related vibrational optic effects,
 5.17–5.20, 5.18f–5.19f, 5.20t, 5.21f
 multiphonon absorption, 5.16–5.17, 5.17t
 phonons, 5.13–5.16, 5.15f
- Lattice vibrations:
 of crystals and glasses, 2.11–2.12, 2.76t–2.77t
 linear-chain model of, 5.8
 optical, 20.31–20.34, 20.32f–20.34f
 in solids, 8.16–8.18, 8.17f, 8.19t–8.20t
- Lead titanate (PbTiO_3), 2.40t, 2.45t, 2.48t,
 2.52t, 2.57t, 2.64t
- Lead vapor, 14.15, 14.16
- Lens housings, 3.15–3.16, 3.15f
- Lenses:
 negative, 3.13
 positive-powered, 3.13
- Lensing, 16.22, 19.5
- LF5 glass (581409), 2.49t, 2.54t, 2.59t, 2.66t
- Li Li polarizing beam splitter, 7.72–7.73, 7.76f
- Light flint glass, 2.42t
- Light modulation, 5.66t
- Light pipe reflectometers, rotating, 5.62, 5.63f
- Light pressure force, 20.7
- Light propagation, in solids, 8.4–8.13
 anisotropic crystals, 8.8–8.11, 8.9t, 8.10f
 energy flow, 8.7–8.8
 interfaces, 8.11–8.13, 8.12f, 8.13f
 Maxwell's equations, 8.4–8.6
 wave equations and optical constants, 8.6–8.7
- Light sources, 5.58–5.59
- Limiters:
 cascaded, 13.6
 optical, 12.32
 self-protecting, 13.9
 tandem, 13.6, 13.6f
- Limiting, optical, 13.2 (*See also* Passive optical limiting)
- Line defects, 9.12–9.13
- Linear electro-optic effect, 12.2–12.3
- Linear optical absorption, 16.18
- Linear optical properties (of semiconductors),
 5.11–5.39, 5.12f, 5.13t
 free carriers, 5.33–5.36, 5.35f–5.37f
 impurity and defect absorption, 5.37–5.39,
 5.38f, 5.39f
 interband absorption, 5.21–5.33
 absorption near fundamental edge,
 5.21–5.22, 5.21f
 direct transitions, 5.22–5.23, 5.22f–5.23f
- Linear optical properties (of semiconductors),
 interband absorption (*Cont.*):
 excitons, 5.25–5.29, 5.26t, 5.27f–5.28f
 high-energy transitions above fundamental
 edge, 5.29–5.33, 5.30f–5.34f
 indirect transitions, 5.23–5.24, 5.24f–5.25f
 polaritons, 5.29
 lattice absorption, 5.13–5.20
 impurity-related vibrational optic effects,
 5.17–5.20, 5.18f–5.19f, 5.20t, 5.21f
 multiphonon absorption, 5.16–5.17,
 5.17t
 phonons, 5.13–5.16, 5.15f
 and models of dielectric function, 5.12
- Linear-chain model, of lattice vibrations, 5.8
- Lin-perp-lin polarization gradient cooling,
 20.17–20.18, 20.18f, 20.19f
- Lippmann-Bragg holographic mirrors, 7.50
- Liquid crystals, 13.12, 13.12f
- Lithium fluoride (LiF), 2.39t, 2.44t, 2.48t,
 2.57t, 2.69t
- Lithium iodate ($\alpha\text{-LiIO}_3$), 2.39t, 2.46t, 2.48t,
 2.57t, 2.75t
- Lithium niobate (LiNbO_3), 2.39t, 2.46t, 2.48t,
 2.52t, 2.57t, 2.63t, 2.75t, 12.14, 17.1,
 17.4–17.13, 17.6f–17.11f
- Lithium tantalate (LiTaO_3), 12.14, 17.14–17.15,
 17.15f
- Lithium triborate (LiB_3O_5) (LBO), 2.39t, 2.46t,
 2.51t, 2.56t, 2.63t, 2.75t, 17.1
- Lithium-calcium-aluminum fluoride
 (LiCaAlF_6) (LiCAF), 2.39t, 2.46t, 2.51t,
 2.57t, 2.63t
- LLF1 glass (548458), 2.49t, 2.54t, 2.59t, 2.66t
- Local density approximation (LDA), 5.5
- Local vibrational modes (LVM), 5.17, 5.18,
 5.19f, 5.20, 5.20f
- Localized vibration, 5.82f, 5.83
- Lockheed Martin, 6.46
- Lock-in amplifiers, 5.64
- Long duration exposure facility (LDEF), 6.17
- Longitudinal acoustic (LA) phonons, 5.24,
 5.25f
- Longitudinal Kerr effect, 18.11–18.15, 18.12f
 (*See also* Self-phase modulation)
- Longitudinal optic (LO) phonons, 5.24, 5.25f,
 5.79, 5.79f, 5.80
- Longitudinal relaxation rate and time, 11.5
- Longitudinal-mode (LO) frequencies, for
 crystals and glasses, 2.11, 2.12

- Lorentz model:
 of absorption, 4.4
 of dispersion, 8.14, 8.21
- Low temperature (LT), 18.3
- Low-*Q* cavities, 9.12
- Low-temperature (LT) molecular beam epitaxy, 18.21
- LOX8 glass, 6.57
- Luminescence spectroscopy, 5.69–5.75, 5.70*f*, 5.72*f*–5.75*f*
- Lyddane-Sachs-Teller (LST) relation, 2.11, 5.14, 8.17
- Mach-Zehnder interferometer (MZI), 23.2–23.4, 23.2*f*
- MacNeille polarizers, 7.70–7.72, 7.71*f*, 7.72*f*, 7.75*f*
- Magnesium-oxide doped stoichiometric lithium tantalite (MgO:sPPLT), 17.14–17.15, 17.15*f*
- Magnetic field modulation, 5.66*t*
- Magnetic permeability, of water, 1.16
- Magnetic traps, 20.21–20.23, 20.22*f*
- Magnetoabsorption, impurity, 5.51–5.52, 5.51*f*
- Magneto-optical (MO) properties
 (of semiconductors), 5.39–5.52, 5.41*t*
 effect of magnetic field on energy bands, 5.40, 5.42, 5.42*f*
 interband effects, 5.42–5.46, 5.43*f*–5.45*f*, 5.47*f*
 intraband or free-carrier effects, 5.47–5.52, 5.48*f*–5.51*f*
 semiconductor nanostructures, 5.52
- Magneto-optical traps (MOTs), 14.18, 20.24–20.25, 20.24*f*, 20.26*f*
- Magnetopolarons, 5.50
- Magnetroreflection, 5.43, 5.44, 5.44*f*
- Magnetron sputtering, 7.11
- Mankiewicz Gebr. & Co., 6.35
- Manly Rowe fraction, 15.15
- Manufacturing error budget, for polymeric optics, 3.10
- “Marine snow,” 1.14, 1.29
- Martin Black, 6.3*t*, 6.12, 6.14, 6.15, 6.26*f*, 6.28*f*, 6.46, 6.47*f*, 6.51*f*, 6.53*f*
- Martin Marietta, 6.48, 6.53
- Martin Optical Black, 6.5*f*
- Mass attenuation coefficients, for photons
 and metals, 4.48, 4.48*t*, 4.49
- Mass density, of metals, 4.6
- Master oscillator-power amplifier (MOPA) laser, 17.6
- Matched filtering, 12.28–12.29, 12.29*f*, 12.30*f*
- Material designation, of crystals and glasses, 2.27–2.30, 2.29*f*
- Matrix theory for multilayer systems, 7.6–7.10, 7.9*f*
- Maximal atomic coherence, 14.28–14.32, 14.29*f*–14.32*f*
- Maximal coherence, 14.3
- Maximum usable temperature, of metals, 4.7, 4.55*t*
- Maxwell wave equation, 2.12
- Maxwell-Bloch equations, 11.6–11.7
- Maxwell-Boltzmann distribution, 20.11
- Maxwell-Helmholtz-Drude dispersion formula, 2.12, 2.21–2.22
- Maxwell’s equations:
 for electric and magnetic fields, 1.16, 1.17, 5.52
 methods for solving, 9.2–9.3
 for optical fields, 16.24
 and refractive index, 2.6
 for solids, 8.4–8.6
- Mean coefficient of linear thermal expansion, 4.7
- Measurement techniques, for semiconductors, 5.7, 5.56–5.83
 ellipsometry, 5.67–5.69, 5.68*f*, 5.69*f*
 inelastic light scattering, 5.76–5.83, 5.76*f*, 5.78*f*–5.82*f*
 instrumentation, 5.58–5.61, 5.59*f*, 5.60*f*, 5.72
 luminescence, 5.69–5.75, 5.70*f*, 5.72*f*–5.75*f*
 modulation spectroscopy, 5.64–5.67, 5.64*f*–5.65*f*, 5.66*t*, 5.68*f*
 reflection and transmission/absorption, 5.62–5.64, 5.63*f*
 spectroscopic procedures, 5.56–5.58
- Measurements:
 on coatings, 7.12–7.14
 with surfaces and processes, 6.9–6.10
- Mechanical assembly, of polymers, 3.14–3.16, 3.14*f*, 3.15*f*
- Mechanical cycling, of metals, 4.10
- Medium-bandpass filters, 7.78–7.83, 7.79*f*, 7.80*f*, 7.82*f*–7.88*f*
- Melt data sheets, for glass, 2.29
- Memory, optical, 11.25, 12.34
- Meniscus-shaped elements, 3.13
- Mercury cadmium telluride ($\text{Hg}_{0.78}\text{Cd}_{0.22}\text{Te}$) narrow-gap alloy, 5.73, 5.74*f*
- Metal-dielectric multiple reflection cutoff filters, 7.59–7.60

- Metal-dielectric multiple reflection filters, 7.111
- Metal-dielectric reflectors, 7.81–7.82, 7.108–7.109, 7.109f, 7.110f
- Metallic reflecting coatings, 7.80–7.81
- Metallic reflectors, 7.106–7.108, 7.106f–7.109f
- Metal-organic chemical vapor deposition (MOCVD), 18.3
- Metals, 4.1–4.70
- absorptance of, 4.39, 4.40f–4.42f, 4.48, 4.49
 - and emittance, 4.49, 4.49f, 4.50t, 4.51t
 - and mass attenuation coefficients for photons, 4.48t
 - aluminum and aluminum alloys
 - absorptance, 4.40f, 4.48t, 4.51t
 - optical properties, 4.12t, 4.20f, 4.21f
 - penetration depth, 4.47f
 - physical properties, 4.52t, 4.54t
 - reflectance, 4.27t–4.28t, 4.40f, 4.44f, 4.46f
 - thermal properties, 4.55t, 4.56t, 4.57f, 4.58t, 4.59f–4.60f, 4.65t, 4.66f, 4.69t, 4.70t
 - beryllium
 - absorptance, 4.48t, 4.50t
 - optical properties, 4.12t, 4.21f, 4.26f
 - penetration depth, 4.47f
 - physical properties, 4.52t, 4.54t
 - reflectance, 4.28t–4.29t, 4.45f, 4.46f
 - thermal properties, 4.55t, 4.56t, 4.57f, 4.58t, 4.59f–4.60f, 4.65t, 4.68f, 4.69t, 4.70t
 - chromium
 - absorptance, 4.48t, 4.50t
 - optical properties, 4.13t–4.14t, 4.22f
 - physical properties, 4.54t
 - reflectance, 4.30t–4.31t
 - thermal properties, 4.69t
 - copper
 - absorptance, 4.40f, 4.48t, 4.50t
 - optical properties, 4.12t–4.13t, 4.22f
 - physical properties, 4.52t–4.54t
 - reflectance, 4.29t–4.30t, 4.40f
 - thermal properties, 4.55t, 4.56t, 4.57f, 4.58t, 4.60f–4.61f, 4.65t, 4.66f, 4.69t, 4.70t
 - germanium
 - absorptance, 4.48t
 - thermal properties, 4.69t, 4.70t
 - gold
 - absorptance, 4.40f, 4.48t, 4.50t, 4.51t
 - optical properties, 4.14t, 4.23f
 - physical properties, 4.52t, 4.54t
 - Metals, gold (*Cont.*):
 - reflectance, 4.31t–4.32t, 4.40f
 - thermal properties, 4.55t, 4.56t, 4.57f, 4.58t, 4.60f–4.61f, 4.65t, 4.66f, 4.69t, 4.70t
 - interband transitions in, 8.21
 - Invar 36, 4.10t, 4.52t, 4.55t, 4.69t, 4.70t
 - iron
 - absorptance, 4.40f, 4.48t, 4.50t
 - optical properties, 4.15t, 4.23f
 - physical properties, 4.54t
 - reflectance, 4.32t–4.33t, 4.40f
 - thermal properties, 4.55t, 4.56t, 4.57f, 4.58t, 4.62f–4.63f, 4.65t, 4.67f, 4.69t
 - mechanical properties of, 4.7–4.8
 - for mirror design, 4.8–4.10, 4.10t
 - molybdenum
 - absorptance, 4.41f, 4.48f, 4.48t, 4.50t, 4.51t
 - optical properties, 4.15t–4.16t, 4.24f
 - physical properties, 4.52t, 4.54t
 - reflectance, 4.33t–4.35t, 4.41f
 - thermal properties, 4.55t, 4.56t, 4.58f, 4.58t, 4.62f–4.63f, 4.65t, 4.67f, 4.69t, 4.70t
 - nickel and nickel alloys
 - absorptance, 4.41f, 4.48t, 4.50t, 4.51t
 - optical properties, 4.16t–4.17t, 4.24f
 - penetration depth, 4.47f
 - physical properties, 4.52t, 4.54t
 - reflectance, 4.35t–4.36t, 4.41f, 4.47f
 - thermal properties, 4.55t, 4.56t, 4.57f, 4.58f, 4.62f–4.63f, 4.65t, 4.67f, 4.69t, 4.70t
 - nomenclature for, 4.3
 - optical properties of, 4.3–4.6, 4.4f, 4.11
 - dielectric function, 4.26f
 - extinction coefficient, 4.11, 4.12t–4.19t, 4.20f–4.26f
 - refraction index, 4.11, 4.12t–4.19t, 4.21f–4.26f
 - penetration depth of, 4.47f
 - physical properties of, 4.6, 4.49, 4.52t–4.54t
 - platinum
 - absorptance, 4.41f, 4.48t, 4.50t, 4.51t
 - optical properties, 4.17t, 4.25f
 - physical properties, 4.54t
 - reflectance, 4.36t–4.37t, 4.41f
 - thermal properties, 4.69t, 4.70t
 - reflectance of, 4.11, 4.27t–4.39t, 4.40f–4.47f
 - semiconductors and dielectrics vs., 8.4
 - silicon
 - absorptance, 4.48t
 - physical properties, 4.52t

- Metals, silicon (*Cont.*):
 reflectance, 4.46f
 thermal properties, 4.55t, 4.56t, 4.58f, 4.58t,
 4.63f–4.64f, 4.65t, 4.68f, 4.69t, 4.70t
- silicon carbide
 absorptance, 4.49f, 4.50f
 optical properties, 4.19t, 4.25f
 physical properties, 4.52t
 reflectance, 4.41f, 4.42f, 4.46f
 thermal properties, 4.55t, 4.56t, 4.58f, 4.58t,
 4.63f–4.64f, 4.65t, 4.68f, 4.69t, 4.70t
- silver
 absorptance, 4.42f, 4.48t, 4.50t, 4.51t
 optical properties, 4.17t–4.18t, 4.26f
 physical properties, 4.52t, 4.54t
 reflectance, 4.37t–4.38t, 4.42f
 thermal properties, 4.55t, 4.56t, 4.57f, 4.58t,
 4.60f–4.61f, 4.65t, 4.66f, 4.69t, 4.70t
- stainless steel
 physical properties, 4.52t
 thermal properties, 4.55t, 4.56t, 4.57f, 4.58t,
 4.62f–4.63f, 4.65t, 4.67f, 4.69t, 4.70t
- steel, 4.50t, 4.51t
- tantalum, 4.50t, 4.69t, 4.70t
- thermal properties of, 4.6–4.7, 4.53, 4.55
 coefficient of linear thermal expansion,
 4.56t, 4.57f, 4.58f
 elastic properties, 4.69, 4.69t
 at room temperature, 4.55t
 specific heat, 4.65t, 4.66f–4.69f
 strength and fracture properties, 4.70, 4.70t
 thermal conductivity, 4.58t, 4.59f–4.63f
- titanium, 4.48t, 4.50t, 4.52t, 4.55t
- tungsten
 absorptance, 4.42f, 4.48t, 4.50t, 4.51t
 optical properties, 4.18t–4.19t, 4.26f
 physical properties, 4.54t
 reflectance, 4.38t–4.39t, 4.42f
 thermal properties, 4.69t, 4.70t
- zinc, 4.48t
- Metrology, cw OPOs for, 17.28, 17.29
- MeV proton acceleration, 21.54, 21.54f
- MH 2200 coating, 6.37, 6.38f
- Mica spacers, 7.84f, 7.88–7.89
- Michelson interferometers, 5.60, 7.42, 7.104f
- Microcavities, in 3D photonic crystals, 9.6–9.12,
 9.7f, 9.8f
- Microcreep strength, of metals, 4.8
- Microstrain, of metals, 4.8
- Microyield strength, of metals, 4.8, 4.70, 4.70t
- Mie scattering, 1.15, 1.32, 2.27
- Mie theory, 1.40
- Miller's rule, 10.9
- Minus filters, 7.43, 7.48, 7.49, 7.50f
- Mirrors:
 aluminum, 7.106–7.108, 7.106f–7.108f
 graded reflectivity, 7.52
 hot and cold, 7.58
 Lippmann-Bragg holographic, 7.50
 metals for, 4.8–4.10, 4.10t
 perfect, 7.47
 phase conjugate, 12.7, 12.8f–12.9f, 12.33–12.35
 self-pumped phase conjugate, 12.7,
 12.8f–12.9f
 semiconductor saturable absorber, 18.3,
 18.10–18.11
 semilinear, 12.7, 12.8f
 silver, 7.106f–7.108f, 7.107–7.109
 supermirrors, 7.111
 (*See also* Reflectors)
- M&M states, 23.14
- Modelocking:
 additive pulse, 18.3, 18.14
 colliding pulse, 18.3
 cw, 18.5f
 cw Q-switched, 18.5f
 Kerr lens, 18.3, 18.14–18.15
 passive, 18.8–18.9, 18.8f, 18.12–18.15
 Q-switched, 18.4, 18.5, 18.5f
 soliton, 18.8f, 18.12–18.14
- Modified lambda coupling, 14.24f
- Modulation spectroscopy, 5.64–5.67,
 5.64f–5.65f, 5.66t, 5.68f
- Mohs scale, 2.31, 2.32f
- Moiré grating, 22.11, 22.12f
- Molasses, optical [*see* Optical molasses (OM)]
- Molding, of polymers, 3.2, 3.12–3.13
- Molecular alignment, in strong fields,
 21.22–21.23, 21.23t, 21.24f
- Molecular beam epitaxy (MBE), 5.7, 18.3
- Molecular dissociation, 21.23
- Molecular orientational Kerr effect, 16.3t
- Molecular tunnel ionization, 21.25–21.26,
 21.25f, 21.27f
- Molecular weight, of crystals and glasses, 2.30
- Molecules, strong field interactions with,
 21.22–21.26
 Coulomb explosion, 21.24–21.25
 nuclear motion and alignment in, 21.22–21.23,
 21.23t, 21.24f

- Molecules, strong field interactions with (*Cont.*):
 triatomic and larger, 21.26
 tunnel ionization and ionization distance,
 21.25–21.26, 21.25f, 21.27f
- Molybdenum:
 absorptance of, 4.41f, 4.48f, 4.48t, 4.50t, 4.51t
 optical properties of, 4.15t–4.16t, 4.24f
 physical properties of, 4.52t, 4.54t
 reflectance of, 4.33t–4.35t, 4.41f
 thermal properties of
 coefficient of linear thermal expansion,
 4.56t, 4.58f
 elastic stiffness, 4.69t
 moduli and Poisson's ratio, 4.69t
 at room temperature, 4.55t
 specific heat, 4.65t, 4.67f
 strength and fracture properties, 4.70t
 thermal conductivity, 4.58t, 4.62f–4.63f
- Momentum, family, 20.38
- Monochromators, 5.59–5.61
- Monoclinic crystals, 2.7t, 2.18, 2.47t, 8.9t, 8.19t
- Morel model of absorption, 1.24, 1.28
- Morse interatomic potential, 2.16
- Mott-Wannier excitons (*see* Wannier excitons)
- Multicomponent polarizers, 7.69
- Multilayer (ML) coatings, 7.96–7.98, 19.4
- Multilayer reflectors, 7.39–7.53
 of absorbing materials, 7.37–7.38, 7.38f
 all-dielectric broadband reflectors, 7.39,
 7.40f, 7.45–7.47, 7.45f–7.47f
 coatings for ultrafast optics, 7.47–7.48, 7.48f
 for far-infrared region, 7.52, 7.52f
 graded reflectivity mirrors, 7.52
 imperfections in, 7.40–7.43, 7.41f–7.43f
 for interferometers and lasers, 7.39–7.40,
 7.39f–7.40f
 narrowband reflection coatings, 7.43, 7.44f
 rejection filters, 7.48–7.50, 7.49f–7.51f
 for soft x-ray and XUV regions, 7.53
 in two-material periodic multilayers theory,
 7.37–7.38, 7.38f
- Multilayers:
 matrix theory for, 7.6–7.10, 7.9f
 periodic
 $[(0.5A)B(0.5A)]^N$, 7.35, 7.36f
 nonabsorbing $[AB]^N$ and $[AB]NA$,
 7.32–7.34, 7.33f–7.35f
 $[xH.(1-x)L]^N.xH$, 7.37
- Multiphonon absorption, 5.16–5.17, 5.17t,
 5.18f, 19.6
- Multiphoton absorption:
 of crystals and glasses, 2.15–2.17, 2.16f, 2.17f
 of solids, 8.30–8.31
- Multiphoton ionization (MPI), 21.10–21.12,
 21.11f
- Multiphoton refraction, of crystals and glasses,
 2.15–2.17, 2.16f, 2.17f
- Multiple anti-Stokes scattering, 15.2t
- Multiple bound excitons, 5.26t
- Multiple cavities, in polymeric optics, 3.10
- Multiple quantum wells (MQWs), 12.22–12.23,
 12.22f
- Multiple Raman scattering, 15.2–15.3, 15.2t
- Multiple Raman Stokes generation, 15.38–15.40,
 15.40f
- Multiple Stokes scattering, 15.2–15.3, 15.2t,
 15.3t
- Multiple-layer surfaces, 6.26, 6.26f–6.27f
- Multiple-reflection filters, 7.111–7.113,
 7.111f–7.113f
- Multiple-resonant oscillators, 17.16–17.21
 doubly resonant, 17.16–17.17
 pump-enhanced singly resonant, 17.17–17.20,
 17.18f–17.20f
 triply resonant, 17.20–17.21, 17.21f
- Multishot laser-induced damage, 19.4
- Mutually incoherent beam couplers, 12.7, 12.8f
- Mutually pumped phase conjugators
 (MPPCs), 12.7
- N (optical constant of water), 1.17, 1.17f
- N on 1 annealing, 19.3, 19.4
- $N00N$ state, 23.9–23.12, 23.10f, 23.11f
- Nanoplasma, in clusters, 21.34–21.35, 21.35f
- Nanostructure semiconductors, 5.52
- Nanostructuring, 6.55, 6.59f
- Nanotubes, 6.55, 6.59f
- Narrowband reflection coatings (narrowband
 rejection filters), 7.43, 7.44f, 7.49
- Narrow-bandpass filters, 7.78–7.83, 7.79f, 7.80f,
 7.82f–7.88f, 7.88–7.89, 7.89f
- NAS polymer, 3.4t
- National Institute of Standards and Technology
 (NIST), 20.16
- Natural waters, 1.3, 1.13–1.15
- N-BaF10 glass (670472), 2.49t, 2.54t, 2.59t, 2.67t
- N-BaK4 glass (569560), 2.49t, 2.54t, 2.59t, 2.66t
- N-BaLF4 glass (580537), 2.49t, 2.54t, 2.59t, 2.66t
- N-BaSF64 glass (704394), 2.49t, 2.54t, 2.59t, 2.67t
- NbF1 glass (743492), 2.49t, 2.54t, 2.59t, 2.67t

- N-BK7 glass (517642), 2.49t, 2.54t, 2.59t, 2.66t
 NBS Black, 6.49–6.50, 6.50f
 Near-IR crystals, visible, 10.21t, 10.22t
 Near-UV spectrum, semiconductor interactions
 with, 5.4f, 5.5
 Negative lenses, 3.13
 Neoprene, 6.32f, 6.34f
 Net irradiance, of water, 1.5t, 1.7f, 1.9
 Networks, 3D photonic crystals and, 9.4–9.5
 Neumann's principle, 2.5
 Neural networks, 12.33–12.35
 Neutral arsenic antisite (As_{Ga}), 18.3
 Neutral atoms, trapping of, 20.21–20.39
 in atomic beam brightening, 20.27–20.28,
 20.27f
 in atomic clocks, 20.28
 in Bose-Einstein condensation, 20.35–20.37,
 20.36f
 in dark states, 20.37–20.39, 20.38f
 with magnetic traps, 20.21–20.23, 20.22f
 with magneto-optical traps, 20.24–20.25,
 20.24f, 20.26f
 in optical lattices, 20.31–20.34, 20.32f–20.34f
 with optical traps, 20.23–20.24, 20.23f
 in ultracold collisions, 20.28–20.31, 20.30f,
 20.31f
 Neutral attenuators, 7.105, 7.105f
 Neutral filters, 7.67, 7.67f–7.68f
 Neutron production, fusion, 21.53
 Nextel 2010, 6.35, 6.37
 Nextel Suede Coating Series 3101-C10, 6.37,
 6.38f, 6.53f
 N-F2 glass (620364), 2.49t, 2.54t, 2.59t, 2.67t
 N-FKS glass (487704), 2.49t, 2.54t, 2.59t, 2.66t
 Nickel:
 absorptance of, 4.41f, 4.48t, 4.50t, 4.51t
 optical properties of, 4.16t–4.17t, 4.24f
 penetration depth, 4.47f
 physical properties of, 4.52t, 4.54t
 reflectance of, 4.35t–4.36t, 4.41f, 4.47f
 thermal properties of
 coefficient of linear thermal expansion,
 4.56t, 4.57f
 elastic stiffness, 4.69t
 moduli and Poisson's ratio, 4.69t
 at room temperature, 4.55t
 specific heat, 4.65t, 4.67f
 strength and fracture properties, 4.70t
 thermal conductivity, 4.58f, 4.62f–4.63f
 Nickel alloys, 4.47f
 Niobium flint glass, 2.42t
 N-K5 glass (522595), 2.49t, 2.54t, 2.59t, 2.66t
 N-KF9 glass (523515), 2.49t, 2.54t, 2.59t, 2.66t
 N-KzFS4 glass (613443), 2.49t, 2.54t, 2.59t, 2.66t
 N-LaF2 glass (744447), 2.49t, 2.54t, 2.59t, 2.67t
 N-LaF33 glass (754523), 2.49t, 2.54t, 2.59t, 2.67t
 N-LaK10 glass (720504), 2.49t, 2.54t, 2.59t, 2.67t
 N-LaSF31A glass (883409), 2.49t, 2.54t, 2.59t, 2.67t
 Noise:
 shot, 23.4–23.6, 23.5f
 and stimulated Raman scattering, 15.35–15.38,
 15.39f
 Nonabsorbing $[AB]^N$ and $[AB]/NA$ multilayers,
 7.32–7.34, 7.33f–7.35f
 Noncritical phase matching (NCPM), 17.1
 Nondegenerate four-wave mixing
 (NDFWM), 14.28
 Nonlinear absorption (NLA), 16.29
 limiters of, 13.6–13.7, 13.6f
 mechanisms of, 16.12–16.13, 16.13f
 nondegenerate, 16.27
 optical limiting by, 13.4–13.7, 13.4f, 13.5f
 and third-order optical nonlinearities,
 16.7–16.9
 Nonlinear atom-field interactions, 11.10
 Nonlinear optical coefficients, of crystals
 and glasses, 2.26–2.27, 2.27t
 Nonlinear optical crystals, 10.19–10.20,
 10.20t–10.22t
 Nonlinear optical (NLO) effects, 19.5,
 19.9–19.11, 19.9f, 19.10f
 Nonlinear optical frequency conversion,
 14.24–14.28, 14.24f, 14.27f
 Nonlinear optical properties (of semiconductors),
 5.52–5.56
 Maxwell's equations and polarization power
 series expansion, 5.52–5.53, 5.54t
 second-order, 5.53–5.55
 third-harmonic generation, 5.56
 third-order, 5.55
 two-photon absorption, 5.56
 Nonlinear optics, 10.3–10.23
 about, 10.4–10.5
 conversion efficiencies, 10.14–10.16
 crystals for, 10.19–10.20, 10.20t–10.22t
 equations for, 10.4–10.5
 microscopic origin of, 10.5–10.10, 10.6f, 10.8f
 and MKS systems, 10.21–10.23
 optical parametric process in, 10.16–10.19,
 10.17f–10.19f

- Nonlinear optics (*Cont.*):
- phase-matching condition in second-order processes, **10.12–10.14**, **10.12f**, **10.13f**
 - second-order susceptibility tensor in, **10.10–10.11**
 - third-order optical nonlinearities, **16.1–16.31**
 - cascaded $x^{(1)}:x^{(1)}$ processes, **16.20–16.22**, **16.21f**
 - cascaded $x^{(2)}:x^{(2)}$ processes, **16.22–16.24**, **16.23f**, **16.24f**
 - four-wave mixing, **16.27–16.28**, **16.28f**
 - interferometry, **16.28–16.29**
 - Kerr effect, **16.11–16.14**, **16.13f**, **16.14f**
 - Kramers-Kronig dispersion relations, **16.9–16.11**
 - nonlinear absorption and nonlinear refraction, **16.7–16.9**
 - propagation effects, **16.24–16.26**
 - and quantum mechanics, **16.4–16.7**, **16.5f**
 - stimulated scattering, **16.14–16.19**, **16.15f**, **16.17f**
 - terms for, **16.1–16.3**, **16.3t**
 - third-harmonic generation, **16.14**
 - time-resolved excite-probe techniques, **16.26–16.27**, **16.26f**
 - two-photon absorption, **16.19–16.20**
 - Z-scan, **16.29–16.30**, **16.30f**
 - ultrashort pulse generation, **18.1–18.23**
 - Kerr effect, **18.11–18.15**, **18.12f**
 - saturable absorbers, **18.5–18.11**, **18.6f–18.8f**
 - semiconductor ultrafast nonlinearities, **18.15–18.23**, **18.16f**, **18.17f**, **18.22f**
 - and ultrafast lasers, **18.3–18.5**, **18.4f**, **18.5f**
- Nonlinear optics in gases, strong field, **21.27–21.31**, **21.28f**
- Nonlinear reflectivity, **18.6–18.7**, **18.7f**
- Nonlinear refraction (NLR), **13.4f**, **13.7–13.8**, **16.7–16.9**, **16.29–16.30**, **16.30f**
- Nonlinear scattering, **13.4f**, **13.8**
- Nonlinear Schrödinger equation (NLSE), **16.25**
- Nonlinear susceptibility, of crystals and glasses, **2.26**
- Nonlinear Thomson scattering, **21.8–21.9**, **21.9f**
- Nonnormal angle of incidence, antireflection coatings at, **7.28–7.31**, **7.28f–7.31f**
- Nonpolarizing beam splitters, **7.63**, **7.64f–7.65f**, **7.65**
- Nonpolarizing edge and bandpass filters, **7.66**, **7.67f**
- Nonreactive evaporation, **7.11**
- Nonresonant degenerate four-wave mixing, **16.27–16.28**, **16.28f**
- No-phonon (NP) photoluminescence, **5.72**, **5.73f**
- Normalized detuning, **14.11f**
- Normalized vector potential, **21.6**
- Normalized water-leaving irradiance, **1.46–1.48**, **1.48f**
- Notch filters, **7.43**
- Novelty filters, **12.32**, **12.33f–12.35f**
- N-PK52A glass (497816), **2.49t**, **2.54t**, **2.59t**, **2.66t**
- N-PSK53A glass (618634), **2.49t**, **2.54t**, **2.59t**, **2.66t**
- N-SF6 (805254), **2.49t**, **2.54t**, **2.59t**, **2.67t**
- N-SK10 glass (623570), **2.49t**, **2.54t**, **2.59t**, **2.67t**
- N-SSK5 glass (658509), **2.49t**, **2.54t**, **2.59t**, **2.67t**
- Nuclear motion, in strong fields, **21.22–21.23**, **21.23t**, **21.24f**
- Null correctors, optical, **3.16**
- Null optics, polymers and, **3.16–3.17**
- N-ZK7 glass (508612), **2.49t**, **2.54t**, **2.59t**, **2.66t**
- Oak Ridge, **6.51**
- Ocean color, **1.46**
- Omnidirectional reflectors, **9.2**
- One-dimensional optical molasses, **20.15–20.16**, **20.15f**
- One-electron transitions, **5.7**
- Open aperture Z-scan, **16.30**
- Open systems, **11.17–11.18**, **11.17f**, **11.18f**
- Optical Bloch equations (OBEs), **11.3–11.6**, **20.3**, **20.6**
- “Optical Characterization in Microelectronics Manufacturing” (S. Perkowitz, D. G. Seiler, W. M. Duncan), **5.61–5.62**
- Optical constants:
- for coatings, **7.13**
 - and dielectric function, **5.8–5.9**
 - of metals, **4.11**, **4.12t–4.19t**, **4.20f–4.26f**
 - of solids, **8.6–8.7**, **8.15**
 - of water, **1.17**, **1.17f**
- Optical frequency synthesis, **17.28**, **17.29**
- Optical indicatrix (index ellipsoid), of crystals and glasses, **2.18–2.19**, **2.19f**
- Optical interconnects, **12.31**, **12.31f**
- Optical Kerr effect (OKE), **16.11–16.14**, **16.13f**, **16.14f**
- Optical lattices, **20.31–20.34**, **20.32f–20.34f**
- Optical limiters, **12.32**

- Optical modes, of crystals and glasses, 2.68t–2.76t
 with cesium chloride structure, 2.68t
 with chalcopyrite structure, 2.74t
 with corundum structure, 2.70t
 with cubic perovskite structure, 2.73t
 with diamond structure, 2.68t
 with fluorite structure, 2.69t
 other structures, 2.74t–2.76t
 with α -quartz structure, 2.71t
 with rutile structure, 2.71t
 with scheelite structure, 2.72t
 with sodium chloride structure, 2.69t
 with spinel structure, 2.73t
 with tetragonal perovskite structure, 2.73t
 with trigonal selenium structure, 2.70t
 with wurtzite structure, 2.70t
 with zinblende structure, 2.69t
- Optical molasses (OM), 20.13–20.17
 defined, 20.3, 20.10
 Doppler cooling, 20.13–20.15, 20.14f
 one-dimensional, 20.15–20.16, 20.15f
 three-dimensional, 20.16–20.17, 20.16f
- Optical monitoring, in thin film manufacturing, 7.11
- Optical parametric amplifiers (OPAs), 23.13–23.14
- Optical parametric chirped pulse amplification (OPCPA), 21.5
- Optical parametric oscillators (OPOs), 10.18–10.19, 10.19f, 14.15 [See also Continuous-wave optical parametric oscillators (cw OPOs)]
- Optical parametric (OP) process, 10.16–10.19, 10.17f–10.19f
- Optical phonons, 5.14
- Optical Ramsey fringes, 11.20–11.22, 11.21f
- Optical responses, classification of, 5.12, 5.13t
- Optical spectroscopy, 11.2
- Optical Stark effect, 16.13
- Optical traps, 20.23–20.24, 20.23f
- Optical tweezers, 20.23
- Optically generated plasmas, 16.20
- Optically polished solid spacers, 7.88–7.89
- Optically induced phase charge, 13.4f, 13.9
- Orange peel, of polymers, 3.11
- Organic black dye, 6.15
- Organic crystals, 12.23–12.25, 12.26t–12.27t
- Organic matter:
 absorption by, 1.22–1.23, 1.23t, 1.25–1.27, 1.25t, 1.26f
 passive limiting in, 13.10
 in water, 1.14
- Organic-inorganic composites, hybrid, 12.27t
- Orlando Black optical coating, 6.54, 6.55f
- Orlando Black surface, 6.8f
- Orthorhombic crystals, 2.7t, 2.18, 2.46t, 8.9t, 8.10, 8.19t
- Oscillator models, of optical nonlinearity, 10.5–10.9, 10.6f, 10.8f
- Oscillators, 12.7–12.9, 12.8f–12.9f (See also specific oscillators, e.g.: Raman oscillators)
- Outer ionization, of cluster, 21.32–21.33
- Outgassing:
 of black surfaces, 6.17
 of polymers, 3.4
- Out-of-plane coupling, 9.11–9.12
- Overdense plasmas, strong field interactions
 with, 21.46–21.52
 high harmonic generation, 21.50–21.52, 21.51f
 $j \times B$ heating and anomalous skin effect, 21.49
 ponderomotive steepening and hole boring, 21.49–21.50, 21.50f
 relativistic effects and induced transparency, 21.52
 resonance absorption, 21.47–21.48, 21.47f
 structure of irradiated plasma, 21.46–21.47, 21.46f
 vacuum heating, 21.47f, 21.48–21.49
- Oxide layer, aluminum reflectance and, 4.44f
- Package-induced laser-induced damage, 19.4–19.5
- Painted surfaces, 6.2t–6.3t
- Paints and surface treatments, 6.35–6.58, 6.37f, 6.43f, 6.53f
 Acktar black coatings, 6.55
 Aeroglaze Z series, 6.36f, 6.37, 6.37f, 6.39, 6.39f–6.42f
 Akzo Nobel paints, 6.39, 6.42f, 6.43f
 anodized processes, 6.44–6.49, 6.47f, 6.48f, 6.51f, 6.53f
 black glass, 6.57
 Black Kapton, 6.57, 6.57f
 carbon nanotubes and nanostructured materials, 6.55, 6.59f
 Cardinal Black, 6.36f, 6.39, 6.44f
 Cat-a-lac Black, 6.39, 6.42f, 6.53f
 conductive/nonconductive, 6.12, 6.12t
 DeSoto Black, 6.37f, 6.39
 DURACON, 6.55–6.56
 electrically conductive black paint, 6.56
 electrodeposited surfaces, 6.53–6.54, 6.54f, 6.55f

- Paints and surface treatments (*Cont.*):
 etching of electroless nickel, 6.49–6.50, 6.50f, 6.51f, 6.53f
 flame-sprayed aluminum, 6.57
 Floquil, 6.44
 gold blacks, 6.57
 high-resistivity coatings, 6.56
 IBM Black (tungsten hexafluoride), 6.56
 ion beam-sputtered surfaces, 6.53
 Parson's Black, 6.44, 6.53f
 plasma-sprayed surfaces, 6.50–6.52, 6.51f–6.53f
 silicon carbide, 6.56
 SolarChem, 6.44, 6.48f, 6.53f
 specular metallic anodized surfaces, 6.57, 6.58f
 sputtered and CVD surfaces, 6.56
 3M paints and derivatives, 6.35–6.37, 6.36f, 6.38f, 6.53f
 ZO-MOD BLACK, 6.56
- Parametric amplification, 10.17–10.18, 16.3t
 Parametric oscillators, 10.18–10.19, 10.18f, 10.19f
 Paratellurite (TeO_2), 2.40t, 2.45t, 2.48t, 2.52t, 2.58t, 2.65t, 2.76t
 Parson's Black, 6.44, 6.53f
- Particles:
 surface coatings and generation of, 6.17–6.18
 in water
 particle size distributions, 1.15–1.16, 1.16f
 refraction index of, 1.20
 scattering by, 1.30–1.35, 1.31t, 1.32f, 1.33f, 1.34t–1.35t
- Particulate matter, in water, 1.14–1.15
 Passband region, transmission in, 7.53, 7.54
 Passive modelocking, 18.8–18.9, 18.8f, 18.12–18.15
 Passive nonlinear optical phenomena, 5.54, 5.54t
 Passive optical limiting, 13.1–13.12
 active vs., 13.1–13.3, 13.2f, 13.3f
 in materials, 13.9–13.12, 13.10f–13.12f
 by nonlinear absorption, 13.4–13.7, 13.4f–13.6f
 by nonlinear refraction, 13.4f, 13.7–13.8
 by nonlinear scattering, 13.4f, 13.8
 optically induced phase charge, 13.4f, 13.9
 by photorefractive, 13.8–13.9
- Pattern recognition, by matched filtering, 12.28–12.29, 12.29f, 12.30f
- Pendular states, 21.23
 Penetration depth, of metals, 4.47f
 Perfect diffuse reflectors (PDRs), 6.7f, 6.27f
 Perfect mirrors, 7.47
- Periclase (MgO), 2.44t, 2.48t, 2.52t, 2.57t, 2.63t, 2.69t
- Periodic multilayers:
 $[(0.5A)B(0.5A)]^N$, 7.35, 7.36f
 nonabsorbing $[AB]^N$ and $[AB]NA$, 7.32–7.34, 7.33f–7.35f
 $[xH.(1-x)L]^N.xH$, 7.37
- Periodically poled lithium niobate (LiNbO_3) (PPLN), 17.1, 17.4–17.13, 17.6f–17.11f
- Periodically poled lithium tantalate (LiTaO_3) (PPLT), 17.14–17.15, 17.15f
- Periodically poled potassium titanyl phosphate (KTiOPO_4) (PPKTP), 17.2, 17.28, 17.29f, 17.30f
- Permittivity, of water, 1.16
- Perovskite, 2.73t
- Petawatt lasers, 21.5
- Petzold volume scattering functions, 1.33, 1.33f
- Petzval sums, 3.8
- Phase charge, optically induced, 13.4f, 13.9
- Phase coatings, 7.101, 7.101f–7.104f, 7.102
- Phase conjugate interferometry, 12.32, 12.33f, 12.34f
- Phase conjugate mirrors, 12.7, 12.8f–12.9f, 12.33–12.35
- Phase conjugation:
 Brillouin, 15.48, 15.52–15.54, 15.52f–15.54f
 photo echo and geometry of, 11.19–11.22, 11.20f, 11.21f
- Phase dispersion filters, 7.89, 7.89f
- Phase distortion, nonlinearly induced, 16.28–16.29
- Phase matching:
 and harmonic yield, 21.30
 and Maxwell-Bloch equations, 11.6
 noncritical, 17.1
 QPM materials, 17.1, 17.13–17.14
 in second-order processes, 10.12–10.14, 10.12f, 10.13f
 and stimulated Raman scattering, 15.7, 15.34, 15.34f
- Phase pulling, of transient Raman scattering, 15.26–15.27, 15.27f
- Phase retarding reflectors, 7.101–7.102, 7.102f, 7.103f
- Phase shifter, cross-Kerr, 23.11
- Phasesonium, 14.3
- Phase-transition temperatures, of crystals and glasses, 2.32, 2.33

- Phonons, 2.11
 acoustic, 5.14–5.16
 coupled plasmon and, 5.35, 5.36, 5.36f, 5.37f
 lattice absorption by, 5.13–5.16, 5.15f
 optical, 5.14
 and Raman scattering, 5.79–5.81, 5.79f–5.81f
 transverse optical (TO), 8.16–8.18
- Phosphate crown glass, 2.41t
- Photoacoustic spectroscopy (PAS), 17.22–17.26, 17.22f–17.24f, 17.26f
- Photo-associative spectroscopy (PAS), 20.30–20.31, 20.31f
- Photoconductors (PCs), 5.61
- Photoelastic coefficients, of crystals and glasses, 2.24
- Photoexcitation, 5.70f
- Photoluminescence (PL), 5.70–5.75, 5.70f, 5.72f–5.75f
- Photoluminescence excitation (PLE) spectroscopy, 5.75, 5.75f
- Photomultipliers (PMTs), 5.61
- Photon absorption (PA), 19.9, 19.10, 19.10f
- Photon echo, 11.11–11.19
 about, 11.11–11.15, 11.12f, 11.13f, 11.15f
 stimulated, 11.15–11.19, 11.16f–11.19f
- Photon echo spectroscopy, 11.26–11.27
- Photon loss, 23.14–23.15
- Photonic bandgaps (PBGs), 9.1–9.17
 fibers with, 2.23
 Maxwell's equations, 9.2–9.3
 in 3D photonic crystals, 9.4–9.12
 criteria for, 9.4–9.5
 examples of, 9.5, 9.5f
 microcavities in, 9.6–9.12, 9.7f, 9.8f
 2D periodicity in microcavities of, 9.8–9.12, 9.9f
 in-plane coupling, 9.10–9.11
 out-of-plane coupling, 9.11–9.12
 and waveguides, 9.12–9.17
 in photonic crystals with 2D periodicity, 9.13–9.14, 9.13f
 waveguide bends, 9.14–9.16, 9.15f, 9.16f
 waveguide intersections, 9.16–9.17, 9.17f
- Photonic crystals, 9.3
- Photonic resonance, 22.9–22.13, 22.10f, 22.12f
- Photon-number eigenstates, 23.7–23.8
- Photons:
 mass attenuation coefficients for, 4.48, 4.48t, 4.49
 Raman interaction with, 15.21–15.22
 and water, 1.11
- Photoreactive scattering, 12.9
- Photoreflectance (PR), 5.66t, 5.67
- Photorefraction:
 optical limiting by, 13.8–13.9
 of third-order optical nonlinearities, 16.22
- Photorefractive effect, 12.1–12.38
 beam spatial profiles, 12.10
 devices using, 12.28–12.38
 associative memories and neural networks, 12.33–12.35
 gain and two-beam coupling, 12.29–12.32, 12.31f
 holographic data storage, 12.37
 holographic storage, 12.36–12.37
 loss and two-beam coupling, 12.31–12.32, 12.33f–12.35f
 phase conjugate interferometry, 12.32, 12.33f, 12.34f
 real-time holography, 12.28–12.29, 12.29f, 12.30f
 solitons, 12.38
 thresholding, 12.35–12.36, 12.36f
 waveguides, 12.37
- grating formation, 12.1–12.3, 12.2f
- oscillators and self-pumped mirrors, 12.7–12.9, 12.8f–12.9f
- standard rate equation model, 12.3–12.4
- stimulated photoreactive scattering, 12.9
- time-dependent effects, 12.10
- wave interactions, 12.4–12.7
 anisotropic scattering, 12.7
 four-wave mixing, 12.6–12.7, 12.6f
 two-beam coupling, 12.4–12.6, 12.4f
- Photorefractive materials, 12.10–12.25
 bulk compound semiconductors, 12.20–12.21, 12.20t, 12.21f
 comparison of, 12.13, 12.13t
 cubic oxides (sillenites), 12.17–12.19, 12.18t, 12.19f
 features of, 12.10–12.11
 ferroelectric, 12.13–12.17, 12.13t
 barium titanate, 12.15–12.16, 12.16f
 lithium niobate and lithium tantalate, 12.14
 potassium niobate, 12.16–12.17
 strontium barium niobate and related compounds, 12.17
 tin hypthiodiphosphate, 12.17, 12.18t
- figures of merit for, 12.11–12.13
 steady-state performance, 12.11–12.12
 transient performance, 12.12, 12.13

- Photorefractive materials (*Cont.*):
 multiple quantum wells, 12.22–12.23, 12.22f
 organic crystals and polymer films,
 12.23–12.25, 12.26t–12.27t
 passive limiting in, 13.11, 13.12f
 Photorefractive self-oscillation, 12.7, 12.9
 Photorelaxation, 5.70f
 Photosynthetically available radiation (PAR),
 1.5t, 1.7f, 1.9
 Photovoltage, 5.66t
 Physical constants, for crystals and glasses, 2.8t
 Physical properties:
 of crystals and glasses, 2.37, 2.38t–2.43t
 classes and symmetry properties, 2.7t
 composition, structure, and density,
 2.38t–2.41t
 physical constants, 2.8t
 specialty, and substrate materials, 2.43t
 symmetry properties, 2.5, 2.6t–2.8t
 of metals, 4.6, 4.49, 4.52t–4.54t
 of polymers, 3.2–3.5, 3.4t
 Phytoplankton:
 absorption by, 1.23–1.25, 1.24f–1.25f, 1.28
 in water, 1.14
 Picosecond, 5.7
 Picosecond solid-state lasers, 18.10
 Piezo-optic coefficients, for crystals
 and glasses, 2.21
 Pikhitin-Yas'kov formula, 2.22
 Plasma beat wave acceleration, 21.41
 Plasma frequency, 2.15, 8.21
 Plasma-ion-assisted deposition, 7.11
 Plasmas, 21.3–21.4
 in fast ignition, 21.54–21.55, 21.54f, 21.55f
 in femtosecond x-ray production, 21.52–21.53,
 21.53f
 in fusion neutron production, 21.53
 in high magnetic field production, 21.53
 in MeV proton acceleration, 21.54
 nano-, 21.34–21.35, 21.35f
 optically generated, 16.20
 overdense, 21.46–21.52
 high harmonic generation, 21.50–21.52,
 21.51f
 irradiated plasma, 21.46–21.47, 21.46f
 $j \times B$ heating and anomalous skin effect, 21.49
 ponderomotive steepening and hole
 boring, 21.49–21.50, 21.50f
 relativistic effects and induced
 transparency, 21.52
 Plasmas, overdense (*Cont.*):
 resonance absorption, 21.47–21.48, 21.47f
 vacuum heating, 21.47f, 21.48–21.49
 wakefield generation and electron
 acceleration, 21.39–21.42, 21.40f, 21.42f
 in Raman amplification, 21.55
 underdense, 21.36–21.46
 direct laser acceleration and betatron
 resonance, 21.42–21.43
 intense laser pulses, 21.38–21.39, 21.38f
 inverse Bremsstrahlung heating, 21.37,
 21.37f
 ionization-induced defocusing,
 21.43–21.44, 21.43f
 ponderomotive channel formation, 21.42
 self-channeling and self-phase modulation,
 21.44–21.46, 21.45f
 Plasma-sprayed surfaces, 6.50–6.52, 6.51f–6.52f
 Plasmons, 5.33–5.36, 5.35f–5.37f, 5.58, 8.23–8.24
 Plastic spacers, for bandpass filters, 7.89
 Plate equations, mirror design and, 4.8–4.9
 Plate polarizers, 7.69–7.70, 7.70f
 Platinum:
 absorptance of, 4.41f, 4.48t, 4.50t, 4.51t
 optical properties of, 4.17t, 4.25f
 physical properties of, 4.54t
 reflectance of, 4.36t–4.37t, 4.41f
 thermal properties of, 4.69t, 4.70t
 Pockels effect, for crystals and glasses, 2.26
 Point defects, 9.12–9.13
 Poisson's equation, for solids, 8.5
 Poisson's ratio:
 for crystals, 2.31
 for metals, 4.7, 4.69t
 Polaritons, 5.29
 Polarization:
 and inelastic scattering, 1.49
 at interface of solid, 8.12–8.13, 8.13f
 and size of electric field, 10.4–10.5
 Polarization dependence, of stimulated Raman
 scattering, 15.41–15.42, 15.42t
 Polarization gradient cooling, 20.17–20.18,
 20.18f, 20.19f
 Polarization power series expansion, 5.52–5.53,
 5.54t
 Polarization splitting, 7.93, 7.93f
 Polarizers:
 embedded, 7.70–7.71, 7.71f, 7.72f
 interference, 7.69–7.73, 7.70f–7.72f
 MacNeille, 7.70–7.72, 7.71f, 7.72f, 7.75f

- Polarizers (*Cont.*):
 multicomponent, 7.69
 plate, 7.69–7.70, 7.70f
- Polarizing beam splitters, 7.70–7.73,
 7.70f–7.72f, 7.74f–7.75f
- Polarons, in magnetic field, 5.50
- Poling process, 12.14
- Polishing, of optical surfaces, 19.3
- Polyallyl diglycol carbonate, 3.4t
- Polyamide (Nylon), 3.4t
- Polyarylate, 3.4t
- Polycarbonate, 3.4t, 3.6, 3.6t, 3.7t, 3.12
- Polychloro-trifluoroethylene, 3.4t
- Polycrystalline materials, 2.3 (*See also* Crystals)
- Polycyclohexyl methacrylate (PCHMA), 3.4t,
 3.6, 3.6t, 3.7t
- Poly-diallylglycol (CR-39) resin, 3.11
- Polyetherimide (PEI), 3.4t, 3.6, 3.6t, 3.7t
- Polyethersulfone, 3.4t
- Polymer composites, 12.26t
- Polymer films, 12.24
- Polymeric optics, 3.1–3.18
 coatings on, 3.17–3.18
 design of, 3.7–3.11
 aberration control, 3.8
 aspheric surfaces, 3.8–3.9
 athermalization, 3.9
 dimensional variations, 3.10
 manufacturing error budget, 3.10
 material selection, 3.8
 multiple cavities, 3.10
 optical figure variations, 3.10
 processing considerations, 3.9
 specification, 3.10–3.11
 strategy, 3.7–3.8
 materials for
 forms of, 3.2
 optical properties, 3.5–3.7, 3.6t, 3.7f
 physical properties, 3.2–3.5, 3.4t
 selection, 3.1–3.2
 processing of, 3.11–3.17
 abrasive forming, 3.11–3.12
 casting, 3.11
 compression molding, 3.12
 geometry considerations, 3.13–3.14
 injection molding, 3.12–3.13
 mechanical assembly, 3.14–3.16, 3.14f, 3.15f
 null optics, 3.16–3.17
 shrinkage, 3.14
 single-point turning, 3.12
- Polymeric optics, processing of (*Cont.*):
 testing and qualification, 3.16
 vendor selection, 3.13
- Polymers, fully functional, 12.27t
- Polymethyl pentene, 3.4t
- Polymethylmethacrylate (PMMA), 3.4t, 3.6,
 3.6t, 3.7t, 3.12
- Polymorphs, 2.27
- Polystyrene, 3.4t, 3.6, 3.6t, 3.7t
- Polystyrene co-butadiene, 3.4t
- Polysulfone, 3.4t
- Polytetrafluoroethylene (Teflon), 6.27
- Polyvinylidene fluoride, 3.4t
- Ponderomotive channel formation, 21.42
- Ponderomotive force, 21.5–21.6
- Ponderomotive steepening, 21.49–21.50, 21.50f
- Population trapping, velocity-selective
 coherent, 20.37
- Positive-powered lenses, 3.13
- Posttreated Martin Black, 6.47
- Potassium bromide (KBr), 2.39t, 2.44t, 2.48t,
 2.51t, 2.56t, 2.62t, 2.69t, 2.77t
- Potassium dihydrogen phosphate (KH_2PO_4)
 (KDP), 2.39t, 2.45t, 2.48t, 2.51t, 2.56t,
 2.62t, 2.75t
- Potassium iodide (KI), 2.39t, 2.44t, 2.48t, 2.51t,
 2.56t, 2.62t, 2.69t, 2.77t
- Potassium niobate (KNbO_3), 2.39t, 2.46t, 2.48t,
 2.51t, 2.56t, 2.62t, 2.75t, 12.16–12.17
- Potassium tantalate (KTaO_3), 2.39t, 2.44t, 2.48t,
 2.51t, 2.56t, 2.62t, 2.73t
- Potassium titanyl arsenate (KTiOAsO_4)
 (KTA), 17.1
- Potassium titanyl phosphate (KTiOPO_4)
 (KTP), 2.39t, 2.46t, 2.51t, 2.56t, 2.63t,
 2.75t, 17.1, 17.2, 17.28, 17.29f, 17.30f
- Pound-Drever-Hall (PDH) technique, 17.18
- Powelite (CaMoO_4), 2.38t, 2.45t, 2.47t, 2.50t,
 2.55t, 2.61t, 2.72t, 2.77t
- Power handling capability, of cw lasers, 7.14
- Poynting vector, 8.7
- Praseodymium, 14.34–14.35, 14.35f, 14.36f
- Pressure force, 20.7
- Principle dispersion (term), 2.23
- Prism-based monochromators, 5.59
- Propagation effects, on third-order optical
 nonlinearities, 16.24–16.26
- Proustite (Ag_3AsS_3), 2.38t, 2.46t, 2.47t, 2.50t,
 2.55t, 2.60t
- Pulse area, 11.8

- Pulse excitation, chirped, 11.25–11.26
Pulse generation, 18.4, 21.31
Pulse propagation, 14.20–14.22
Pulsed lasers, 14.15–14.16, 14.16f
Pulsed transient Raman scattering, 15.22–15.25, 15.24f–15.26f, 15.24t
Pump (incident light), 15.1
Pump depletion, 15.9, 15.10, 15.15, 15.20, 15.20f
Pump-enhanced singly resonant oscillators (PE-SROs), 17.2–17.4, 17.3f, 17.4f, 17.17–17.20, 17.18f–17.20f, 17.27–17.28
Pumping, for cw SROs, 17.2–17.4
Pump-probe (excite-probe) measurements, 16.26–16.27, 16.27f
Pump-probe spectroscopy, 18.18–18.19
Pump-probe techniques, 18.6, 18.6f
Pump-resonant singly resonant oscillators, 17.17–17.20, 17.18f–17.20f
Purity, of optical materials, 2.5
Pyrex glass, 2.43t, 2.49t, 2.54t
- Q-switched modelocking, 18.4, 18.5, 18.5f
Qualification, of polymers, 3.16
Quality factor (of microcavities), 9.7–9.8, 9.8f
Quantum coherence tomography, 23.13
Quantum dots, 12.25
Quantum entanglement, in optical interferometry, 23.1–23.15
 concepts and equations for, 23.1–23.4, 23.2f, 23.4f
 digital approaches to, 23.7–23.9
 Heisenberg limit, 23.6–23.7
 N00N state, 23.9–23.12, 23.10f, 23.11f
 and quantum imaging, 23.13–23.14
 and remote sensing, 23.14–23.15
 shot-noise limit, 23.4–23.6, 23.5f
Quantum fluctuations, 15.38, 15.39f, 23.5, 23.5f
Quantum imaging, 23.13–23.14
Quantum interferences, 14.8
Quantum mechanical model, for solids, 8.4, 8.24–8.25
Quantum mechanics, third-order optical nonlinearities and, 16.4–16.7, 16.5f
Quantum remote sensing, 23.14–23.15
Quantum theory of nonlinear optical susceptibility, 10.9–10.10
Quart-enhanced photoacoustic spectroscopy (QEPAS), 17.25–17.26, 17.26f
- α -Quartz (SiO₂), 2.40t, 2.46t, 2.48t, 2.52t, 2.57t, 2.64t, 2.71t
Quasi-phase-matched (QPM) nonlinear materials, 17.1, 17.13–17.14
Quasi-plane-wave solution (to transverse profile problem), 12.10
- Rabi frequency, 11.4
Radiance, of water, 1.5t, 1.6, 1.7f
Radiation pressure force, 20.7
Radiation resistance, of polymers, 3.5
Radiative escape, 20.29
Radiative transfer theory, 1.4
Radiometric quantities, of water, 1.4–1.9, 1.5t, 1.7f
Raman amplification, 21.55
Raman amplifiers, 15.4, 15.4f, 22.15
Raman cooling, 20.21
Raman cross sections, 15.5
Raman gain coefficients, 15.16t–15.18t
Raman generators, 15.4, 15.4f
Raman induced Kerr effect (RIKE), 16.3t, 16.12, 16.17
Raman linewidths, 15.9, 15.10f, 15.11t–15.20t
Raman modes, of crystals and glasses, 2.11
Raman oscillators, 15.4, 15.4f
Raman scattering, 15.1–15.43
 anti-Stokes, 15.4, 15.4f, 15.32–15.34, 15.33f, 15.35f, 15.42, 15.42t, 15.43f
 backward, 15.41, 21.38–21.39, 21.38f
 Brillouin vs., 15.1
 coherent, 15.3, 15.4, 15.4f, 15.34, 15.42, 15.42t, 15.43f
 for crystals and glasses, 2.27
 forward, 21.38–21.39, 21.38f
 measurement with, 5.76–5.83, 5.76f, 5.78f–5.82f
 and Raman interactions, 15.2–15.3, 15.3t
 regimes of, 15.3
 in solids, 8.16, 8.18, 8.19t–8.20t
 spontaneous, 15.3, 15.5
 stimulated (*see* Stimulated Raman scattering)
 transient, 15.22–15.32
 broadband effects, 15.28–15.32, 15.29f
 phase pulling, 15.26–15.27, 15.27f
 pulsed, 15.22–15.25, 15.24f–15.26f, 15.24t
 solitons, 15.27–15.28, 15.29f
 spectral properties, 15.32
 by water, 1.48, 1.49, 1.49f

- Raman shifts, 15.2
- Raman sidebands, generation of, 14.31, 14.32f
- Raman spectroscopy, 5.57–5.58, 7.48, 7.66, 7.83, 7.96, 7.97f, 7.98
- Raman susceptibility, 15.5–15.6
- Raman threshold, 15.38
- Raman transition frequencies, 15.11t–15.15t
- Ramsey fringes, 11.20–11.22, 11.21f
- Rapid adiabatic passage (RAP), 14.1
- Rayleigh scattering:
 - in crystals and glasses, 2.10, 2.27
 - and scattering by sea water, 1.30
 - in third-order optical nonlinearities, 16.14–16.15
- Rayleigh-wing scattering, 16.13
- Reactive evaporation, 7.11
- Reactive force, 20.8
- Real transitions, 16.4
- Real-time holography, 12.28–12.29, 12.29f, 12.30f
- Recoil limit, 20.5
- Redistribution force, 20.8
- Redspot Paint and Varnish, 6.35
- Reflectance:
 - at interface of solid, 8.12
 - of metals, 4.5, 4.11, 4.27t–4.39t, 4.39, 4.40f–4.47f
 - nonabsorbing $[AB]^N$ and $[AB]NA$ multilayers, 7.32–7.34, 7.33f–7.35f
 - of optical coatings, 7.12–7.13
 - in solids, 8.22f, 8.23
 - specular, of black coatings, 6.26f
 - and surface coatings, 6.20t
- Reflection(s):
 - from all-dielectric broadband reflectors, 7.47
 - of coatings on substrate, 7.3
 - enhancement of, for filters and coatings, 7.108–7.109, 7.109f–7.110f
 - filters with low, 7.104–7.106, 7.104f–7.105f
 - magnetoreflexion, 5.43, 5.44, 5.44f
 - magnitude of, 9.2
 - measurement of, 5.62–5.64, 5.63f
 - and optical performance, 7.15–7.16, 7.16f
 - and phase change of nonabsorbing periodic multilayers, 7.34
 - photoreflexion, 5.66t, 5.67
 - total internal, 8.13
- Reflection coatings, 7.106–7.113, 7.106f–7.113f
 - at angles close to grazing incidence, 7.111
 - enhancement of reflection, 7.108–7.109, 7.109f–7.110f
 - metallic reflectors, 7.106–7.108, 7.106f–7.109f
- Reflection coefficient, for optical constants, 5.9–5.10
- Reflection filters, 7.5, 7.5f, 7.111–7.113, 7.111f–7.113f
- Reflectivity:
 - at interface of solid, 8.12
 - nonlinear, 18.6–18.7, 18.7f
 - of solids, 8.23
- Reflectivity amplitude, for solids, 8.15
- Reflectometers, 5.62–5.64, 5.63f
- Reflectors:
 - heat, 7.58, 7.58f
 - metal-dielectric, 7.81–7.82, 7.108, 7.109, 7.110f
 - multilayer, 7.39–7.53
 - all-dielectric broadband reflectors, 7.39, 7.40f, 7.45–7.47, 7.45f–7.47f
 - coatings for ultrafast optics, 7.47–7.48, 7.48f
 - for far-infrared region, 7.52, 7.52f
 - graded reflectivity mirrors, 7.52
 - imperfections in, 7.40–7.43, 7.41f–7.43f
 - for interferometers and lasers, 7.39–7.40, 7.39f–7.40f
 - narrowband reflection coatings, 7.43, 7.44f
 - rejection filters, 7.48–7.50, 7.49f–7.51f
 - for soft x-ray and XUV regions, 7.53
 - and two-material periodic multilayers theory, 7.37–7.38, 7.38f
 - omnidirectional, 9.2
 - resonant, 7.43–7.45, 7.44f
 - very low loss, 7.41–7.42
- Refraction:
 - enhanced, 14.20
 - nonlinear, 16.7–16.9
- Refractive index:
 - of anisotropic crystals, 8.8
 - complex, 1.16–1.17
 - of crystals and glasses, 2.6, 2.8, 8.10
 - dispersion formulas for, 2.21–2.22
 - in dressed atoms, 14.19–14.20
 - of metals, 4.3, 4.11, 4.12t–4.19t, 4.21f–4.26f
 - of particles in water, 1.20
 - in polymeric optics, 3.6–3.7, 3.6t, 3.7f
 - in solids, 8.22f, 8.23
 - and temperature, 2.24–2.26
 - uniformity of, 2.5
 - of water, 1.16–1.20, 1.18f, 1.19t–1.20t
- Refractive index spectrum, 22.6, 22.6f, 22.7f
- Rejection filters, 7.48–7.50, 7.49f–7.51f
- Rejection ratios, of bandpass filters, 7.77

- Rejection region, for cutoff filters, 7.56
- Relativistic effect(s), in strong field interactions:
 with atoms, 21.19–21.20
 with free electrons, 21.6–21.8, 21.7*f*
 with overdense plasmas, 21.52
 self-focusing and self-channeling as,
 21.44–21.45, 21.45*f*
 self-phase modulation as, 21.45–21.46,
 21.45*f*
- Relativistic electron ATI, 21.20, 21.21*f*
- Relativistic electron beams, strong field
 interactions with, 21.9–21.10
- Relativistic suppression of rescattering, 21.20
- Relaxation:
 longitudinal and transverse, 11.5
 photorelaxation, 5.70*f*
- Remote sensing:
 quantum, 23.14–23.15
 in water, 1.46–1.47
- Reorientational Kerr effect, in solids,
 16.13–16.14, 16.14*f*
- Rescattering effects:
 relativistic suppression of, 21.20
 in strong field interactions with atoms,
 21.18–21.19, 21.18*f*, 21.19*f*
- Resins, 3.2, 3.11
- Resistivity:
 coatings with high, 6.56
 of metals, 4.54*t*
 of polymers, 3.5
- Resonance:
 cyclotron, 5.12, 5.12*f*, 5.40, 5.47–5.50,
 5.48*f*–5.50*f*
 photonic, 22.9–22.13, 22.10*f*, 22.12*f*
 slow light propagation and atomic,
 22.2–22.9, 22.3*f*, 22.6*f*–22.8*f*
- Resonance absorption, in overdense plasmas,
 21.47–21.48, 21.47*f*
- Resonant degenerate four-wave mixing, 16.28
- Resonant modes (RMs), 5.17
- Resonant Raman scattering, 16.15
- Resonant reflectors, 7.43–7.45, 7.44*f*
- Retinal damage, 13.3
- Reverse-proton-exchanged (RPE) PPLN, 17.16
- Reverse-saturable absorption (RSA), 13.5,
 13.6*f*, 13.7
- Rigidity, of polymers, 3.3
- Ring resonators, 12.7, 12.8*f*, 22.11, 22.12*f*
- Rock-salt lattices, 5.16
- Rotating light pipe reflectometers, 5.62, 5.63*f*
- Rotating wave approximation (RWA), 20.7
- Roughened aluminum, 6.21, 6.22*f*, 6.30*f*
- Roughness, surface, 6.15
- Rubidium, 14.17, 14.17*f*
- Rubidium titanyl arsenate (RbTiOAsO₄)
 (RTA), 17.1
- Rubidium titanyl phosphate (RbTiOPO₄)
 (RTP), 2.40*t*, 2.46*t*, 2.48*t*, 2.52*t*,
 2.57*t*, 2.64*t*
- Rugate filters, 7.49, 7.50
- Rutile (TiO₂), 2.40*t*, 2.45*t*, 2.48*t*, 2.53*t*, 2.58*t*,
 2.65*t*, 2.71*t*
- S on 1 annealing, 19.3, 19.4
- Salt (NaCl), 2.40*t*, 2.44*t*, 2.48*t*, 2.52*t*, 2.57*t*,
 2.64*t*, 2.69*t*
- Sandblasted aluminum, 6.45*f*
- Sapphire (Al₂O₃):
 dispersion formulas for, 2.60*t*
 elastic constants of, 2.46*t*
 infrared spectrum of, 2.13, 2.13*f*
 lattice vibration model parameters for, 2.76*t*
 mechanical properties of, 2.47*t*
 optical modes of, 2.70*t*
 optical properties of, 2.55*t*
 properties of, 2.38*t*
 thermal properties of, 2.50*t*
 Ti:sapphire amplifiers, 21.5
 Ti:sapphire lasers, 18.3
- Saturable absorbers, 18.5–18.11
 fast, 18.9–18.10
 self-amplitude modulation, 18.5–18.7,
 18.6*f*, 18.7*f*
 semiconductor saturable absorber mirrors,
 18.3, 18.10–18.11
 slow, 18.7–18.9, 18.8*f*
- Saturable absorption, 13.5
- Saturable Bragg reflectors (SBRs), 18.3, 18.11
- Saturation fluence, 18.5, 18.6
- Scatterance, of water, 1.5*t*, 1.10
- Scattering:
 angular distribution of, 1.12
 anisotropic, 12.7
 in crystals and glasses, 2.27
 Einstein-Smoluchowski theory of, 1.30
 nonlinear, 13.4*f*, 13.8
 rescattering, 21.18–21.20, 21.18*f*, 21.19*f*
 stimulated, 16.14–16.19
 from surface, 6.15
 Thomson, 21.8–21.9, 21.9*f*

- Scattering (*Cont.*):
- by water
 - inelastic and polarization, 1.47–1.49, 1.48f, 1.49f
 - measurement, 1.29–1.30
 - particles, 1.30–1.35, 1.31t, 1.32f, 1.33f, 1.34t–1.35t
 - pure water and pure sea water, 1.30
 - wavelength dependence, 1.35–1.40, 1.35t, 1.36f, 1.37t, 1.38t, 1.39f, 1.40t
 - (*See also specific scattering, e.g.: Anti-Stokes scattering*)
 - Scattering force, 20.7
 - Scattering losses, multilayer reflectors and, 7.41
 - Scattering rate, 20.4
 - Scatterometers, 7.13
 - Scheelite (CaWO_4), 2.38t, 2.45t, 2.47t, 2.51t, 2.56t, 2.61t, 2.72t, 2.77t
 - Schott Glass Technologies, 6.57
 - Schott glasses, 2.22, 2.23, 2.26
 - Schrödinger equation, 2.16, 8.24, 8.25
 - Schrödinger's cat, 23.1, 23.9
 - Sea water, 1.3, 1.18–1.21, 1.18f, 1.19t–1.20t, 1.22t, 1.30
 - Second-harmonic processes, conversion efficiencies for, 10.14–10.16
 - Second-order nonlinear optics:
 - anharmonic oscillator model of susceptibility in, 10.7–10.9, 10.8f
 - phase matching in, 10.12–10.14, 10.12f, 10.13f
 - properties of semiconductors in, 5.53–5.55
 - Second-order susceptibility tensor, 10.10–10.11
 - Selenium, 2.70t
 - Self-action effects, 16.25
 - Self-amplitude modulation (SAM), 18.3, 18.5–18.8, 18.6f, 18.7f
 - Self-channeling, 21.44–21.45, 21.45f
 - Self-defocusing, 13.7, 13.8, 19.9–19.11, 19.10f
 - Self-focusing:
 - and laser-induced damage, 19.5, 19.7
 - as relativistic effect in strong field interactions, 21.44–21.45, 21.45f
 - thermal, 13.7–13.8
 - of third-order optical nonlinearities, 16.25
 - Self-heating, 17.12
 - Self-lensing, 19.5
 - Self-modulated wakefield generation, 21.41
 - Self-oscillation, 12.7–12.9
 - Self-phase modulation (SPM), 21.45–21.46, 21.45f (*See also Longitudinal Kerr effect*)
 - Self-protecting limiters, 13.9
 - Self-pumped phase conjugate mirrors (SPPCMs), 12.7, 12.8f–12.9f
 - Self-trapped excitons, 5.26t
 - Sellaite (MgF_2), 2.40t, 2.45t, 2.48t, 2.52t, 2.57t, 2.63t, 2.71t
 - Sellenite ($\text{Bi}_{12}\text{SiO}_{20}$) BSO, 2.38t, 2.44t, 2.47t, 2.50t, 2.55t, 2.61t, 2.75t
 - Sellmeier dispersion model, for crystals and glasses, 2.14, 2.15, 2.21–2.23, 2.25
 - Sellmeier formula, for refractive index, 8.14
 - Sellmeier model, for crystals and glasses, 2.10, 2.11
 - Semiconductor(s), 5.1–5.83
 - dielectrics and metals vs., 8.4
 - electromagnetic spectrum interactions with, 5.3–5.6, 5.4f
 - ion-implanted, 18.21
 - linear optical properties of, 5.11–5.39, 5.12f, 5.13t, 5.35f–5.37f
 - free carriers, 5.33–5.36
 - impurity and defect absorption, 5.37–5.39, 5.38f, 5.39f
 - interband absorption, 5.21–5.33, 5.21f–5.25f, 5.26t, 5.27f–5.28f, 5.30f–5.34f
 - lattice absorption, 5.13–5.20, 5.15f, 5.17t, 5.18f–5.19f, 5.20t, 5.21f
 - low-temperature, 18.21
 - magneto-optical properties of, 5.39–5.52, 5.41t
 - effect of magnetic field on energy bands, 5.40, 5.42, 5.42f
 - interband effects, 5.42–5.46, 5.43f–5.45f, 5.47f
 - intraband or free-carrier effects, 5.47–5.52, 5.48f–5.51f
 - semiconductor nanostructures, 5.52
 - materials and applications of, 5.84t–5.86t
 - measurement techniques for (*see* Measurement techniques, for semiconductors)
 - nanostructure, 5.52
 - nonlinear, 5.52–5.56, 5.54t
 - optical/dielectric response in, 5.8–5.11
 - passive limiting in, 13.9, 13.10f
 - structure of, 5.6–5.7
 - Semiconductor saturable absorber mirrors (SESAMs), 18.3, 18.10–18.11
 - Semiconductor saturable absorbers, 18.4, 18.5, 18.5f, 18.10

- Semiconductor ultrafast nonlinearities, 18.15–18.23
 applications of, 18.19–18.23
 and carrier trapping, 18.21–18.23, 18.22*f*
 in coherent regime, 18.19–18.20
 and continuum excitations, 18.20
 and excitonic excitations, 18.19–18.20
 properties of, 18.16–18.17, 18.16*f*
 and pump-probe spectroscopy, 18.18–18.19
 in thermalization regime, 18.20–18.21
 and transient four-wave mixing, 18.17–18.18, 18.17*f*
- Semiconductor-doped dielectric films, 18.11
- Semiconductor-doped glasses, 18.11
- Semilinear mirrors, 12.7, 12.8*f*
- Semiquantum four-parameter model, 2.12
- Service temperature, of polymers, 3.3
- Shake off model (of strong field behavior), 21.18
- Shape factors, of bandpass filters, 7.77
- Shear modulus, for metals, 4.69*t*
- Short flint glass, 2.42*t*
- Short pulse, 11.7
- Shot-noise limit (SNL), 23.4–23.6, 23.5*f*
- Shrinkage, polymer, 3.14
- Sikkens Aerospace Finishes, 6.39
- Silica glass, fused, 2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- Silicon:
 absorptance of, 4.48*t*
 crystals of, 2.40*t*, 2.44*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.64*t*, 2.68*t*
 IR absorption due to interstitial oxygen in, 5.19*f*
 multiphonon absorption of vacuum-grown, 5.18*f*
 physical properties, 4.52*t*
 reflectance, 4.46*f*
 thermal properties of
 coefficient of linear thermal expansion, 4.56*t*, 4.58*f*
 elastic stiffness, 4.69*t*
 moduli and Poisson's ratio, 4.69*t*
 at room temperature, 4.55*t*
 specific heat, 4.65*t*, 4.68*f*
 strength and fracture properties, 4.70*t*
 thermal conductivity, 4.58*t*, 4.63*f*–4.64*f*
- Silicon carbide (SiC):
 absorptance of, 4.49*f*, 4.50*f*
 optical properties of, 4.19*t*, 4.25*f*
 particles of, 6.15
 physical properties of, 4.52*t*
 reflectance of, 4.41*f*, 4.42*f*, 4.46*f*
- Silicon carbide (SiC) (*Cont.*):
 as surface material, 6.56
 thermal properties of
 coefficient of linear thermal expansion, 4.56*t*, 4.58*f*
 elastic stiffness, 4.69*t*
 moduli and Poisson's ratio, 4.69*t*
 at room temperature, 4.55*t*
 specific heat, 4.65*t*, 4.68*f*
 strength and fracture properties, 4.70*t*
 thermal conductivity, 4.58*t*, 4.63*f*–4.64*f*
- Silicon-on-insulator planar waveguide, 22.15
- Sillenites (cubic oxides), 12.17–12.19, 12.18*t*, 12.19*f*
- Silver:
 absorptance of, 4.42*f*, 4.48*t*, 4.50*t*, 4.51*t*
 optical properties of, 4.17*t*–4.18*t*, 4.26*f*
 physical properties of, 4.52*t*, 4.54*t*
 reflectance of, 4.37*t*–4.38*t*, 4.42*f*
 resistivity of, 4.54*t*
 thermal properties of
 coefficient of linear thermal expansion, 4.56*t*, 4.57*f*
 elastic stiffness, 4.69*t*
 moduli and Poisson's ratio, 4.69*t*
 at room temperature, 4.55*t*
 specific heat, 4.65*t*, 4.66*f*
 strength and fracture properties, 4.70*t*
 thermal conductivity, 4.58*t*, 4.60*f*–4.61*f*
- Silver gallium sulfide (AgGaS₂) (AGS), 2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.74*t*, 2.76*t*
- Silver mirrors, 7.106*f*–7.108*f*, 7.107–7.109
- Silver selenogallate (AgGaSe₂), 2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.74*t*, 2.76*t*
- Single active electron approximation, 21.10
- Single-particle excitation, 5.81, 5.82*f*
- Single-pass absorption spectroscopy, 17.24–17.27, 17.25*f*, 17.26*f*
- Single-point turning, of polymers, 3.12
- Single-scattering albedo, 1.5*t*, 1.7*f*
- Singly resonant optical parametric (SOP) oscillators, 10.18, 10.18*f*
- Singly resonant oscillators (SROs), 17.2–17.16, 17.3*f*, 17.4*f*
 guided-wave nonlinear structures, 17.15–17.16
 MgO:sPPLT in, 17.14–17.15, 17.15*f*
 PPLN crystals in, 17.4–17.13, 17.6*f*–17.11*f*
 QPM nonlinear materials, 17.13–17.14
- Sink, in molding, 3.14
- Sisyphus laser cooling, 20.19

- Skin depth (term), 4.5
 Skin effect, anomalous, 21.49
 Slope of cutoff, 7.56
 Slow light propagation, 22.1–22.15
 and atomic resonance, 22.2–22.9, 22.3f, 22.6f–22.8f
 in optical fibers, 22.13–22.15, 22.14f
 and photonic resonance, 22.9–22.13, 22.10f, 22.12f
 Slow saturable absorbers, 18.7–18.9, 18.8f
 Slowing, of atomic beams, 20.11–20.13, 20.12f, 20.12t, 20.13f
 Snell's law, 7.7, 8.11, 8.23
 Sodium chloride (NaCl), 2.69t
 Soft x-ray region:
 bandpass filters for, 7.94–7.96, 7.95f–7.96f
 interference polarizers for, 7.73, 7.76f–7.77f
 multilayer reflectors for, 7.42, 7.53
 Solar absorptance, 6.19, 6.21
 Solar-cell cover filters, 7.58, 7.59f
 SolarChem, 6.44, 6.45f, 6.48f, 6.53f
 Sol-gels, 12.27t
 Solid spacers, for bandpass filters, 7.83, 7.88
 Solids:
 band structures and interband transitions
 of, 8.24–8.32
 direct interband absorption, 8.27–8.28
 energy band structures, 8.25–8.27, 8.26f
 excitons, 8.31–8.32
 indirect transitions, 8.29–8.30, 8.29f
 joint density of states, 8.28, 8.29f
 multiphoton absorption, 8.30–8.31
 quantum mechanical model, 8.24–8.25
 selection rules and forbidden transitions, 8.28, 8.29
 bound electronic optical Kerr effect in, 16.12–16.13, 16.13f
 dephasing in, 14.13–14.14
 dispersion relations in, 8.14–8.16
 EIT in, 14.33–14.36, 14.35f, 14.36f
 extrinsic optical properties of, 8.3
 free-electron properties of, 8.21–8.24, 8.22f
 intrinsic optical properties of, 8.3
 lattice interactions in, 8.16–8.18, 8.17f, 8.19t–8.20t
 optical properties of, 8.1–8.32
 propagation of light in, 8.4–8.13
 anisotropic crystals, 8.8–8.11, 8.9t, 8.10f
 energy flow, 8.7–8.8
 interfaces, 8.11–8.13, 8.12f, 8.13f
 Solids, propagation of light in (*Cont.*):
 Maxwell's equations, 8.4–8.6
 wave equations and optical constants, 8.6–8.7
 reorientational Kerr effect in, 16.13–16.14, 16.14f
 Solid-state detectors, 5.61
 Solid-state lasers, 18.3
 Soliton lasers, 18.14
 Soliton modelocking, 18.8f, 18.12–18.14
 Solitons:
 photorefractive, 12.38
 and third-order optical nonlinearities, 16.25–16.26
 and transient Raman scattering, 15.27–15.28, 15.29f
 Space groups, of crystals, 2.27–2.28
 Space-based surfaces, 6.17–6.18, 6.20t, 6.46
 Spacecraft, 6.12, 6.16–6.17
 Spacers, for bandpass filters, 7.83, 7.88–7.89
 Spatial solitons, 16.26
 Specific heat (*see* Heat capacity)
 Specification, in polymeric optics, 3.10–3.11
 Spectral absorptance, of water, 1.9
 Spectral absorption, by detritus in water, 1.26–1.27, 1.26f
 Spectral absorption coefficient:
 for case 1 waters, 1.27–1.28
 for phytoplankton, 1.24, 1.24f, 1.25f, 1.25t
 upper bound, in sea water, 1.21, 1.22t
 for waters, 1.10, 1.17, 1.20–1.21, 1.27f
 Spectral beam attenuation coefficient, for water, 1.10
 Spectral diffuse attenuation coefficient, for water, 1.12
 Spectral diffusion, 11.15
 Spectral downwelling average cosine, of water, 1.12
 Spectral downwelling irradiance, of water, 1.8
 Spectral downward scalar irradiance, of water, 1.8
 Spectral emittance, of metals, 4.6
 Spectral gain narrowing, in steady-state Stokes scattering, 15.21
 Spectral hole, 18.20
 Spectral hole burning, 22.6, 22.7f
 Spectral irradiance reflectance, of water, 1.12, 1.46, 1.47f
 Spectral net irradiance, of water, 1.9
 Spectral radiance, of water, 1.6
 Spectral scatterance, of water, 1.10

- Spectral scattering coefficient, for water, 1.10
- Spectral transmission, in polymeric optics, 3.6
- Spectral transmittance, 7.3–7.4
- Spectral upward plane irradiance, of water, 1.8
- Spectral upward scalar irradiance, of water, 1.8
- Spectral upwelling irradiance, of water, 1.8
- Spectral volume scattering function (VSFs), 1.37–1.38, 1.38*t*
- Spectrometers, 5.59–5.61, 5.60*f*, 5.72, 7.12
- Spectroscopic ellipsometry, 5.66*t*
- Spectroscopic measurement, 5.56–5.58
- Spectroscopy, 17.21–17.27
 - continuous-wave, 11.2
 - defined, 11.2
 - differential transmission, 18.18–18.19
 - high-resolution Doppler-free, 17.27
 - knowledge derived from, 5.6
 - photoacoustic, 17.22–17.26, 17.22*f*–17.24*f*, 17.26*f*
 - photo-associative, 20.30–20.31, 20.31*f*
 - photon echo, 11.26–11.27
 - pump-probe, 18.18–18.19
 - Raman, 7.48, 7.66, 7.83, 7.96, 7.97*f*, 7.98
 - single-pass absorption, 17.24–17.27, 17.25*f*, 17.26*f*
 - time-dependent, 11.2
- Specular baffles, 6.14
- Specular black surfaces, 6.14
- Specular reflectance, of black coatings, 6.26*f*
- Spherelike Brillouin zone, 9.4
- Spinel (MgAl_2O_4), 2.39*t*, 2.44*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.63*t*, 2.73*t*
- Spire Corporation, 6.51
- Split-step beam propagation, 12.10
- Spontaneous emission, 9.8, 20.5
- Spontaneous parametric down-conversion (SPDC), 23.10, 23.13
- Spontaneous parametric process, 10.16–10.17, 10.17*f*
- Spontaneous Raman scattering, 15.3, 15.5
- Sputtered beryllium surface, 6.6*f*
- Sputtered surfaces, 6.56
- Sputtering method, of manufacturing thin films, 7.11, 7.14
- Square-top multivacuity bandpass filters, 7.82–7.83, 7.82*f*–7.88*f*
- Squeezing, in quantum optical interferometry, 23.7
- Stability, of bandpass filters, 7.94
- Stainless steel:
 - physical properties of, 4.52*t*
 - thermal properties of, 4.55*t*–4.58*f*, 4.62*f*–4.63*f*, 4.65*t*, 4.67*f*, 4.69*t*, 4.70*t*
- Standard rate equation model, 12.3–12.4
- Standard Test Method (of ASTM), 6.17
- Star Instruments, 6.57
- Stark chirped rapid adiabatic passage (SCRAP), 14.1
- Stationary spectroscopy, 11.2
- Steady-state Stokes scattering, 15.7–15.22
 - gain coefficients in, 15.16*t*–15.18*t*
 - gain narrowing in, 15.21
 - photon description in, 15.21–15.22
 - pump depletion in, 15.9, 15.10, 15.15, 15.20, 15.20*f*
 - Raman linewidths in, 15.9, 15.10*f*, 15.11*t*–15.20*t*
 - spectral gain narrowing in, 15.21
 - steady-state gain in, 15.8–15.9
- Steel, 4.50*t*, 4.51*t*
- Stiffness, in crystals and glass, 2.30, 2.31*t*
- Stimulated Brillouin scattering (SBS), 15.43–15.54
 - equations for, 15.44–15.48, 15.45*f*
 - phase conjugation, 15.48, 15.52–15.54, 15.52*f*–15.54*f*
 - Raman vs., 15.1
 - scattering parameters of materials, 15.48, 15.49*t*–15.51*t*
 - and slow light, 22.14–22.15
 - and stimulated photorefractive scattering, 12.9
 - third-order, 16.18–16.19
- Stimulated parametric process, 10.17–10.18
- Stimulated photon echo, 11.15–11.19, 11.16*f*–11.19*f*
- Stimulated photoreactive scattering (SPS), 12.9
- Stimulated Raman adiabatic passage (STIRAP), 14.1, 14.4, 14.5, 14.7
- Stimulated Raman anti-Stokes scattering, 15.2*t*
- Stimulated Raman scattering (SRS), 15.3–15.43
 - anti-Stokes, 15.2*t*, 15.32–15.34, 15.33*f*, 15.35*f*
 - backward, 15.41
 - coherent spectroscopy, 15.42, 15.42*t*, 15.43*f*
 - equations for, 15.6–15.7
 - focused beams, 15.41
 - formulas for, 15.20*t*
 - geometries for, 15.4, 15.4*f*

- Stimulated Raman scattering (SRS) (*Cont.*):
 multiple Stokes generation, 15.38–15.40, 15.40f
 and noise, 15.35–15.38, 15.39f
 in plasma, 21.38–21.39, 21.38f
 polarization dependence, 15.41–15.42, 15.42t
 Raman susceptibility, 15.5–15.6
 and slow light, 22.15
 steady-state Stokes, 15.7–15.22
 gain coefficients, 15.16t–15.18t
 gain narrowing, 15.21
 photon description, 15.21–15.22
 pump depletion, 15.9, 15.10, 15.15, 15.20, 15.20f
 Raman linewidths, 15.9, 15.10f, 15.11t–15.20t
 spectral gain narrowing, 15.21
 steady-state gain, 15.8–15.9
- Stokes, 15.2t
 terminology, 16.3t
 third-order, 16.15–16.18, 16.15f, 16.17f
 time-gated imaging, 15.42, 15.43, 15.44f
 transient effects of, 15.22–15.32
 broadband effects, 15.28–15.32, 15.29f
 phase pulling, 15.26–15.27, 15.27f
 pulsed, 15.22–15.25, 15.24f–15.26f, 15.24t
 solitons, 15.27–15.28, 15.29f
 spectral properties, 15.32
- Stimulated Rayleigh-Wing scattering, 16.3t
- Stimulated scattering, 16.14–16.19 [*See also* Stimulated Brillouin scattering (SBS); Stimulated Raman scattering (SRS)]
- Stokes scattering, 15.2–15.3, 15.2t, 15.3t (*See also* Anti-Stokes scattering; Steady-state Stokes scattering)
- Stokes shift, 15.1, 15.43
- Stokes shifted Raman scattering, 16.15, 16.15f
- Stokes wave, 15.1, 15.43
- Strain, in crystals and glasses, 2.6t, 2.30, 2.31t, 2.36
- Stray light, analysis and issues, 6.10, 6.19
- Strength:
 of crystals and glasses, 2.31–2.32, 2.32f, 2.32t
 of metals, 4.8, 4.70, 4.70t
 of scattering in water, 1.12
- Stress:
 in crystals and glasses, 2.6t, 2.30, 2.31t
 uniaxial, 5.66t
- Strong field approximation (SFA), 21.12
- Strong field double ionization, 21.18–21.19
- Strong field interactions, 21.1–21.55
 about, 21.2–21.3
 with atoms, 21.10–21.21
 above threshold ionization, 21.14–21.17, 21.15f, 21.16f
 ionization stabilization, 21.20–21.21, 21.22f
 Keldysh parameter, 21.10
 multiphoton and quasi-classical regimes, 21.10
 multiphoton ionization, 21.10–21.12, 21.11f
 relativistic effects, 21.19–21.20, 21.21f
 rescattering effects, 21.18–21.20, 21.18f, 21.19f
 tunnel ionization, 21.12–21.14, 21.12f, 21.14f
 with clusters, 21.31–21.36
 Coulomb explosion, 21.33–21.34
 intense laser pulse interactions, 21.35–21.36, 21.36f
 ionization, 21.31–21.33, 21.32f
 nanoplasma description, 21.34–21.35, 21.35f
 history of, 21.3–21.4
 laser technology for, 21.4–21.5, 21.4f
 with molecules, 21.22–21.26
 Coulomb explosion, 21.24–21.25
 nuclear motion and alignment, 21.22–21.23, 21.23t, 21.24f
 triatomic and larger molecules, 21.26
 tunnel ionization and ionization distance, 21.25–21.26, 21.25f, 21.27f
 nonlinear optics in gases, 21.27–21.31
 attosecond pulse generation, 21.31
 high order harmonic generation, 21.27–21.30, 21.28f
 with single electrons, 21.5–21.10
 interactions with relativistic electron beams, 21.9–21.10
 nonlinear Thomson scattering, 21.8–21.9, 21.9f
 ponderomotive force, 21.5–21.6
 relativistic effects, 21.6–21.8, 21.7f
 with underdense plasmas, 21.36–21.52
 applications of, 21.52–21.55, 21.53f–21.55f
 direct laser acceleration and betatron resonance, 21.42–21.43
 high harmonic generation, 21.50–21.52, 21.51f
 intense laser pulses, 21.38–21.39, 21.38f
 inverse Bremsstrahlung heating, 21.37, 21.37f

- Strong field interactions, with underdense plasmas (*Cont.*):
 ionization-induced defocusing,
 21.43–21.44, 21.43*f*
 $j \times B$ heating and anomalous skin effect, 21.49
 ponderomotive channel formation, 21.42
 ponderomotive steepening and hole boring,
 21.49–21.50, 21.50*f*
 relativistic effects and induced transparency,
 21.52
 resonance absorption, 21.47–21.48, 21.47*f*
 self-channeling and self-phase modulation,
 21.44–21.46, 21.45*f*
 structure of irradiated plasma, 21.46–21.47,
 21.46*f*
 vacuum heating, 21.47*f*, 21.48–21.49
 wakefield generation and electron acceleration,
 21.39–21.42, 21.40*f*, 21.42*f*
- Strontium, 14.15, 14.16*f*
- Strontium barium niobate (SBN), 12.17
- Strontium fluoride (SrF_2), 2.40*t*, 2.44*t*, 2.48*t*,
 2.52*t*, 2.57*t*, 2.64*t*, 2.69*t*
- Strontium molybdate (SrMoO_4), 2.40*t*, 2.45*t*,
 2.48*t*, 2.52*t*, 2.57*t*, 2.65*t*, 2.72*t*
- Strontium titanate (SrTiO_3), 2.40*t*, 2.44*t*, 2.48*t*,
 2.52*t*, 2.57*t*, 2.65*t*, 2.73*t*
- Structured antireflection coatings, 7.25–7.26,
 7.25*f*, 7.26*f*
- Styrene acrylonitrile (SAN), 3.4*t*, 3.6, 3.6*t*, 3.7*t*
- Sum frequency generation (SFG), 16.2, 16.3*t*
- Sum rules:
 for crystals and glasses, 2.9
 for dispersion in solids, 8.14–8.15
 for semiconductors, 5.11
- Super resolution, 23.10, 23.12
- Super-Beer's law, 23.15
- Superlinear absorption of light, 5.57
- Supermirrors, 7.111
- Surface coatings, 6.4*t*
- Surface damage, laser-induced, 19.2–19.4
- Surface-tension effects, from polymer molding process, 3.14
- Susceptibility:
 nonlinear
 anharmonic oscillator model of
 second-order, 10.7–10.9, 10.8*f*
 of crystals and glasses, 2.26
 quantum theory of, 10.9–10.10
 Raman, 15.5–15.6
- Swept-carrier time-domain optical memory,
 11.25
- Switching, third-order optical nonlinearities and,
 16.30–16.31
- Symmetry properties, of crystals and glasses,
 2.5, 2.6*t*–2.8*t*
- Tandem limiters, 13.6, 13.6*f*
- Tantalum, 4.50*t*, 4.69*t*, 4.70*t*
- Tantalum crown glass, 2.42*t*
- Target normal sheath acceleration, 21.54
- Teflon overcoat, 6.27
- Teflon Wet Lubricant, 6.27
- Telescopes:
 and black surfaces, 6.21
 far-infrared, 6.48
 ground-based, 6.12
- Temperature(s):
 of crystals and glasses, 2.32, 2.33
 in laser cooling, 20.5
 of metals, 4.7
 and refractive index of glasses, 2.24–2.26
- Temperature dependence:
 of bandpass filters, 7.94
 of line broadening parameters, 15.19*t*
 of line shift parameters, 15.19*t*
- Temporal instability, of metals, 4.9
- Tensor properties:
 of crystals and glasses, 2.5, 2.6*t*
 of third-order optical nonlinearities,
 16.2–16.3
- Tensors:
 compliance, 2.30
 dielectric, 2.17–2.18
 d -tensor, 10.11
 inverse dielectric, 2.6*t*, 2.19
 second-order susceptibility, 10.10–10.11
- Tests and testing:
 Knoop, 2.31, 2.32*f*
 of polymers, 3.16
- Tetragonal crystals, 8.9*t*, 8.19*t*
 room-temperature elastic constants,
 2.44*t*–2.45*t*
 symmetries of, 2.7*t*
- Tetragonal perovskite, 2.73*t*
- Tetrahedral lattice site, 5.6
- Textured graphite surface, 6.8*f*
- Thallium bromide (TlBr), 2.40*t*, 2.44*t*, 2.48*t*,
 2.53*t*, 2.58*t*, 2.65*t*, 2.68*t*
- Thallium chloride (TlCl), 2.65*t*, 2.68*t*

- The Theory of Coherent Atomic Excitation*
(B. W. Shore), 14.3
- Thermal blooming, 16.22
- Thermal conductivity:
of crystals and glasses, 2.6t, 2.35–2.36, 2.35f
of metals, 4.7, 4.10t, 4.53, 4.55, 4.55t, 4.58t,
4.60f–4.64f
- Thermal cycling, of metals, 4.10
- Thermal defocusing, 13.8
- Thermal diffusivity, for metals, 4.10t
- Thermal effects, on third-order optical nonlinearities, 16.22
- Thermal expansion:
of crystals and glasses, 2.6t, 2.34–2.35, 2.34f
for metals, 4.10t
of metals, 4.6
- Thermal instability, of metals, 4.10
- Thermal properties:
of crystals and glasses, 2.50t–2.53t, 2.55t
of metals, 4.6–4.7, 4.53, 4.55
coefficient of linear thermal expansion,
4.56t, 4.57f, 4.58f
elastic properties, 4.69, 4.69t
at room temperature, 4.55t
specific heat, 4.65t, 4.66f–4.69f
strength and fracture properties, 4.70, 4.70t
thermal conductivity, 4.58t, 4.59f–4.63f
- Thermal self-focusing, 13.7–13.8
- Thermalization, of free electron and hole distributions, 18.20–18.21
- Thermomodulation, 5.66t
- Thermo-optic coefficients, of crystals and glasses, 2.21, 2.24–2.26, 2.24f
- Thermo-optic effect, 16.22
- Thermoplastic resins, 3.2
- Thermoset resins, 3.2
- Thickness errors, for multilayer reflectors, 7.40
- Thin-film coatings:
and antireflection coatings, 7.27–7.28, 7.28f
laser-induced damage to, 19.4
manufacturing of, 7.10–7.12
of metal, 7.104, 7.104f
for multiple reflection filters, 7.111, 7.112,
7.112f, 7.113f
theory and design of, 7.5–7.10, 7.6f, 7.9f
- Third-order harmonic generation (THG),
16.2, 16.3t
in crystals, 16.14
energy level diagrams for, 16.5f
and semiconductors, 5.56
- Third-order optical nonlinearities, 16.1–16.31
cascaded $x^{(1)}:x^{(1)}$ processes, 16.20–16.22,
16.21f
cascaded $x^{(2)}:x^{(2)}$ processes, 16.22–16.24,
16.23f, 16.24f
and four-wave mixing, 16.27–16.28, 16.28f
and interferometry, 16.28–16.29
Kerr effect, 16.11–16.14, 16.13f, 16.14f
Kramers-Kronig dispersion relations,
16.9–16.11
nonlinear absorption and nonlinear
refraction, 16.7–16.9
propagation effects, 16.24–16.26
and quantum mechanics, 16.4–16.7, 16.5f
and semiconductors, 5.55
stimulated scattering, 16.14–16.19, 16.15f,
16.17f
terms for, 16.1–16.3, 16.3t
third-harmonic generation, 16.14
and time-resolved excite-probe techniques,
16.26–16.27, 16.27f
two-photon absorption, 16.19–16.20
and Z-scan, 16.29–16.30, 16.30f
- Thomson scattering, 21.8–21.9, 21.9f
- 3D bandgap materials, 9.2
- 3D optical molasses, 20.16–20.17, 20.16f
- 3D photonic crystals, 9.4–9.8
criteria for, 9.4–9.5
examples of, 9.5, 9.5f
microcavities in, 9.7–9.8, 9.7f, 9.8f
- 3M Black Velvet, 6.14
- 3M Black Velvet 101-C10, 6.12, 6.35
- 3M Nextel Black Velvet, 6.35, 6.36f
- 3M paints, 6.35–6.37, 6.38f, 6.53f
3M Nextel Black Velvet, 6.35, 6.36f
MH 2200, 6.37
Nextel 2010, 6.35, 6.37
Nextel Suede Coating Series 3101-C10, 6.37,
6.38f, 6.53f
- Three-level atomic systems, 14.4–14.6, 14.4f, 14.6f
- Three-photon absorption (3PA), 19.9, 19.10,
19.10f
- Three-wave mixing, 14.26
- Threshold ionization:
absorbance above, 21.14–21.17, 21.15f, 21.16f
defined, 21.3
- Thresholding devices, 12.35–12.36, 12.36f
- Time-dependent (transient) spectroscopy, 11.2
- Time-domain atom interferometers,
11.22–11.24, 11.24f

- Time-gated imaging, 15.42, 15.43, 15.44f
- Time-integrated intensity, 11.18
- Time-of-flight (TOF) measurement of velocity distribution, 20.13, 20.13f
- Time-resolved excite-probe techniques, 16.26–16.27, 16.27f
- Tin hypthiodiphosphate ($\text{Sn}_2\text{P}_2\text{S}_6$), 12.17, 12.18t
- Tiodize V-E17, 6.49
- Titanium, 4.48t, 4.50t, 4.52t, 4.55t
- Titanium dioxide, 6.15
- Titanium sapphire (Ti:sapphire) amplifiers, 21.5
- Titanium sapphire (Ti:sapphire) lasers, 18.3
- Tomography, quantum coherence, 23.13
- Total emittance, of metals, 4.6
- Total integrated excitation, for crystals and glasses, 2.19
- Total internal reflection (TIR), 8.13
- Total mass loss (TML), 6.17
- Total power law, 2.19–2.20, 2.20f
- Total strain, of crystals and glasses, 2.36
- Transfer matrix solution (to Maxwell's equations), 9.3
- Transient four-wave mixing (TFW), 18.17–18.18, 18.17f
- Transient Raman scattering, 15.22–15.32
broadband effects, 15.28–15.32, 15.29f
phase pulling, 15.26–15.27, 15.27f
pulsed, 15.22–15.25, 15.24f–15.26f, 15.24t
solitons, 15.27–15.28, 15.29f
spectral properties, 15.32
- Transitions, 16.4 (*See also specific transitions, e.g.: One-electron transitions*)
- Transmission:
of coatings on substrate, 7.3
for cutoff filters, 7.54–7.55, 7.55f
measurement of, 5.64
in passband region, 7.53, 7.54
- Transmission coefficient, for optical constants, 5.9–5.10
- Transmission filters, 7.3–7.5, 7.83, 7.88f
- Transmission grating, 12.7, 12.8f
- Transmittance:
at interface of solid, 8.12
of metals, 4.6
of optical coatings, 7.12–7.13
spectral, 7.3–7.4
of water, 1.5t
- Transparency:
and absorption, 2.17
EIT [*see* Electromagnetically induced transparency (EIT)]
induced, 21.52
- Transparent prisms, 5.59
- Transverse acoustic (TA) phonons, 5.24, 5.25f
- Transverse Kerr effect, 18.11, 18.12f, 18.14–18.15
- Transverse optical (TO) frequencies, for crystals and glasses, 2.11, 2.12
- Transverse optical (TO) phonons, 5.24, 5.25f, 5.80, 5.80f, 8.16–8.18
- Transverse relaxation time, 11.5
- Trap loss collisions, 20.29
- Trapping atoms, 20.21–20.39
applications of, 20.26–20.39
and atomic beam brightening, 20.27–20.28, 20.27f
and atomic clocks, 20.28
and Bose-Einstein condensation, 20.35–20.37, 20.36f
and dark states, 20.37–20.39, 20.38f
magnetic traps, 20.21–20.23, 20.22f
magneto-optical traps, 20.24–20.25, 20.24f, 20.26f
and optical lattices, 20.31–20.34, 20.32f–20.34f
optical traps, 20.23–20.24, 20.23f
and ultracold collisions, 20.28–20.31, 20.30f, 20.31f
- Triatomic molecules, in strong fields, 21.26
- Triclinic crystals, 2.7t, 2.18, 8.9t, 8.10
- Trigonal crystals, 8.9t, 8.19t
- Trigonal selenium, 2.70t
- Triply resonant oscillators (TROs), 17.2–17.4, 17.3f, 17.4f, 17.20–17.21, 17.21f
- TRU-Color Diffuse Black, 6.49
- Tunable double resonance (electromagnetically induced transparency), 22.6–22.9, 22.7f, 22.8f
- Tunable phase-dispersion filters, 7.89, 7.89f
- Tungsten:
absorptance of, 4.42f, 4.48t, 4.50t, 4.51t
elastic properties of, 4.69t
extinction coefficient for, 4.18t–4.19t, 4.26f
reflectance of, 4.38t–4.39t, 4.42f
refraction index for, 4.18t–4.19t, 4.26f
resistivity of, 4.54t
strength and fracture properties of, 4.70t
- Tungsten hexafluoride, 6.56

- Tunnel ionization:
 atomic, 21.12–21.14, 21.12*f*, 21.14*f*
 molecular, 21.25–21.26, 21.25*f*, 21.27*f*
 relativistic, 21.20
- Tunneling, collective, 21.18
- Turning, single-point, 3.12
- Tweezers, optical, 20.23
- Twin beams of light, correlated, 17.28, 17.29*f*, 17.30*f*
- Two-beam coupling:
 optical limiting by, 13.8–13.9
 photorefractive gain in, 12.29–12.32, 12.31*f*
 photorefractive loss in, 12.31–12.32, 12.33*f*–12.35*f*
 and wave interactions, 12.4–12.6, 12.4*f*
- 2D photonic crystals, microcavities of, 9.8–9.12, 9.9*f*
 in-plane coupling, 9.10–9.11
 out-of-plane coupling, 9.11–9.12
 waveguides in, 9.13–9.14, 9.13*f*
- Two-level atoms:
 coherence in, 14.4–14.5, 14.4*f*
 force on, 20.6–20.7
 at rest, 20.7–20.8
- Two-level coupling, 14.30, 14.30*f*
- Two-material periodic multilayers theory, 7.32–7.38
 $[(0.5A)B(0.5A)]^N$ multilayers, 7.35, 7.36*f*
 angular sensitivity, 7.37
 multilayer reflectors of absorbing materials, 7.37–7.38, 7.38*f*
 nonabsorbing $[AB]^N$ and $[AB]NA$ multilayers, 7.32–7.34, 7.33*f*–7.35*f*
 width of high-reflectance zone, 7.36–7.37, 7.37*f*
 $[xH.(1-x)L]^N.xH$ multilayers, 7.37
- Two-photon absorption (2PA):
 energy level diagrams for, 16.5*f*
 and laser-induced damage, 19.9, 19.10, 19.10*f*
 and optical limiting, 13.4, 13.5, 13.6*f*
 of semiconductors, 5.56
 symbols, 16.8
 in third-order optical nonlinearities, 16.19–16.20
- Two-photon transitions, 11.22–11.23, 11.24*f*
- ULE glass, 2.43*t*, 2.49*t*, 2.54*t*
- Ultimate strength, of metals, 4.8
- Ultracold collisions, 20.26, 20.28–20.31, 20.30*f*, 20.31*f*
- Ultrafast depletion, of semiconductor band states, 18.21
- Ultrafast lasers, 11.26, 18.3–18.5, 18.4*f*, 18.5*f*
- Ultrafast optics, coatings for, 7.47–7.48, 7.48*f*
- Ultrashort pulse generation, 18.1–18.23
 Kerr effect, 18.11–18.15
 longitudinal, 18.11–18.15, 18.12*f*
 transverse, 18.14–18.15
- saturable absorbers, 18.5–18.11
 fast, 18.9–18.10
 self-amplitude modulation, 18.5–18.7, 18.6*f*, 18.7*f*
 semiconductor saturable absorber mirrors, 18.3, 18.10–18.11
 slow, 18.7–18.9, 18.8*f*
- semiconductor ultrafast nonlinearities, 18.15–18.23
 carrier trapping, 18.21–18.23, 18.22*f*
 in coherent regime, 18.19–18.20
 continuum excitations, 18.20
 excitonic excitations, 18.19–18.20
 experimental techniques, 18.17–18.19, 18.17*f*
 properties, 18.16–18.17, 18.16*f*
 in thermalization regime, 18.20–18.21
 and ultrafast lasers, 18.3–18.5, 18.4*f*, 18.5*f*
- Ultralow light pulses, 14.22–14.23, 14.23*f*
- Ultraviolet (UV) crystals, 10.22*t*
- Ultraviolet (UV) light:
 and black surfaces, 6.21, 6.22*f*–6.25*f*
 metal-dielectric reflectors for, 7.108–7.109, 7.109*f*
 semiconductor interactions with, 5.4*f*, 5.5
 [See also Extreme ultraviolet (XUV) light]
- Uncertainty principle, 23.4, 23.6
- Uncertainty state, 23.6
- Underdense plasmas, strong field interactions with, 21.36–21.46
 direct laser acceleration and betatron resonance, 21.42–21.43
 intense laser pulses, 21.38–21.39, 21.38*f*
 inverse Bremsstrahlung heating, 21.37, 21.37*f*
 ionization-induced defocusing, 21.43–21.44, 21.43*f*
 ponderomotive channel formation, 21.42
 self-channeling and self-phase modulation, 21.44–21.46, 21.45*f*
 wakefield generation and electron acceleration, 21.39–21.42, 21.40*f*, 21.42*f*

- Uniaxial crystals, 8.8, 8.9t, 8.10f
 Uniaxial stress, 5.66t
 Unit cell, crystal, 2.30
 Universal antireflection coatings, 7.26, 7.27f
 Upward plane irradiance, 1.5t, 1.7f, 1.8
 Upward scalar irradiance, 1.5t, 1.7f, 1.8
 Upwelling average cosine, 1.6t, 1.7f
 Upwelling irradiance, 1.7f, 1.8
 Urbach tail model, 2.14–2.15
 Urbach's rule, 5.23
 Urea [(NH₄)₂CO], 2.40t, 2.45t, 2.48t, 2.52t, 2.57t, 2.64t

 Vacuum heating, 21.47f, 21.48–21.49
 Vacuum ultraviolet (VUV) spectrum, 5.4f, 5.5
 Vacuum-metal interfaces, 4.43f
 Vacuum-ultraviolet (VUV) radiation, 14.3
 Valence band (VB), 18.3
 van Hove singularities, 8.28
 Vee (V) coupling, 14.1, 14.6f
 Velocity distribution measurement, 20.13, 20.13f
 Velocity-changing collisions, 11.15
 Velocity-selective coherent population trapping (VSCPT), 14.5, 20.37
 Very dense crown glass, 2.42t
 Very light flint glass, 2.41t
 Very low loss reflectors, 7.41–7.42
 Vibration(s):
 lattice, 2.11–2.12, 2.76t–2.77t
 local, 5.17, 5.18, 5.19f, 5.20, 5.20f, 5.82f, 5.83
 phonon, 5.14–5.16
 Vibrational optic effects, 5.17–5.20, 5.18f–5.19f, 5.20t, 5.21f
 Virtual transitions, 16.4
 Viruses, in water, 1.14
 Visible near-IR nonlinear optical crystals, 10.21t, 10.22t
 Volkov state, 21.12
 Volume scattering function (VSFs):
 for sea water and ocean water, 1.34t–1.35t
 spectral, 1.37–1.38, 1.38t
 of water, 1.5t, 1.7f, 1.31, 1.31t, 1.32f, 1.33
 wavelength dependence of, 1.36f

 W point, 9.4
 Wakefield generation, 21.39–21.42, 21.40f
 “Walking” backward, by solitons, 15.28
 Wannier excitons, 5.26–5.27, 5.26t, 8.31

 Water, 1.3–1.50
 absorption, 1.20–1.29
 bio-optical models for, 1.27–1.29, 1.27f, 1.28t
 by dissolved organic matter, 1.22–1.23, 1.23t
 by organic detritus, 1.25–1.27, 1.25t, 1.26f
 by phytoplankton, 1.23–1.25, 1.24f–1.25f
 by polymers, 3.4
 by sea water, 1.21, 1.22t
 apparent optical properties, 1.12–1.13
 attenuation
 beam, 1.40–1.41, 1.41f, 1.42f
 diffuse and Jerlov water types, 1.42–1.46, 1.43t–1.45t, 1.44f, 1.45f
 constituents of natural waters, 1.13–1.15
 dissolved substances, 1.13
 particulate substances, 1.14–1.15
 electromagnetic properties of, 1.16–1.17, 1.18t
 inherent optical properties, 1.9–1.12, 1.10f
 irradiance reflectance and remote sensing, 1.46–1.47, 1.47f
 particle size distributions, 1.15–1.16, 1.16f
 pure, 1.3
 radiometric quantities, 1.4–1.9, 1.5t–1.6t, 1.7f
 refraction index, 1.18–1.20
 particles, 1.20
 sea water, 1.18–1.20, 1.18f, 1.19t–1.20t
 scattering, 1.30
 inelastic and polarization, 1.47–1.49, 1.48f, 1.49f
 measurement of, 1.29–1.30
 by particles, 1.30–1.35, 1.31t, 1.32f, 1.33f, 1.34t–1.35t
 by pure water and pure sea water, 1.30
 wavelength dependence of, 1.35–1.40, 1.35t, 1.36f, 1.37t, 1.38t, 1.39f, 1.40t
 terminology and notation, 1.3–1.4
 Water vapor regained (WVR), 6.17
 Wave equations, for light propagation in solids, 8.6–8.7
 Wave interactions, photorefractive effect and, 12.4–12.7, 12.4f, 12.6f
 Waveguide bends, 9.14–9.16, 9.15f, 9.16f
 Waveguide intersections, 9.16–9.17, 9.17f

- Waveguides:
 and photonic bandgaps, 9.12–9.17
 in photonic crystals with 2D periodicity,
 9.13–9.14, 9.13f
 waveguide bends, 9.14–9.16, 9.15f, 9.16f
 waveguide intersections, 9.16–9.17, 9.17f
 photorefractive, 12.37
 silicon-on-insulator planar, 22.15
- Wavelength dependence, of scattering, 1.35–1.40,
 1.35t, 1.36f, 1.37t, 1.38t, 1.39f, 1.40t
- Wavelength modulation, 5.66t
- Wavelengths, de-Broglie, 8.4
- Wedge filters, 7.90, 7.91, 7.91f
- Weld-line, in molding, 3.14
- Wide bandpass filters, 7.90, 7.90f
- Wide-angle bandpass filters, 7.93–7.94, 7.93f
- Wulfenite (PbMoO_4), 2.40t, 2.45t, 2.48t, 2.52t,
 2.57t, 2.64t, 2.72t
- Wurtzite (α -ZnS):
 in crystals and glasses, 2.70t
 lattices of, 5.6
 properties of, 2.41t, 2.44t, 2.46t, 2.48t, 2.53t,
 2.58t, 2.66t, 2.69t, 2.70t
- X-ray region:
 beam splitters for, 7.63
 soft
 bandpass filters for, 7.94–7.96, 7.95f–7.96f
 interference polarizers for, 7.73, 7.76f–7.77f
 multilayer reflectors for, 7.42, 7.53
- Yellow matter, 1.13, 1.21–1.23, 1.28
- Yield strength, of metals, 4.8, 4.70, 4.70t
- Young's modulus:
 for crystals, 2.30, 2.31
 for metals, 4.7, 4.10t, 4.69t
 for polymers, 3.3
- Yttria (Y_2O_3), 2.40t, 2.44t, 2.48t, 2.53t, 2.58t,
 2.65t, 2.76t
- Yttrium aluminum garnet ($\text{Y}_3\text{Al}_5\text{O}_{12}$) (YAG),
 2.44t, 2.53t, 2.58t, 2.65t, 2.76t
- Yttrium lithium fluoride (LiYF_4) (YLF), 2.45t,
 2.48t, 2.52t, 2.57t, 2.63t, 2.72t
- Yttrium vanadate (YVO_4), 2.40t, 2.45t, 2.48t,
 2.53t, 2.58t, 2.65t, 2.76t
- Yurke state, 23.8
- ZBLAN glass, 2.43t, 2.49t, 2.54t,
 2.59t, 2.68t
- Zernike dispersion formula, 2.22
- Zero dispersion point, for glasses, 2.23
- Zerodur glass, 2.43t, 2.49t, 2.54t
- Zinc, absorptance of, 4.48t
- Zinc crown glass, 2.41t
- Zinc selenide (ZnSe), 2.41t, 2.48t, 2.53t, 2.58t,
 2.66t, 2.69t
- Zinc telluride (ZnTe), 2.41t, 2.48t, 2.53t, 2.58t,
 2.66t, 2.69t
- Zinblende (β -ZnS):
 in crystals and glasses, 2.69t
 lattices of, 5.6, 5.16, 5.17t
 properties of, 2.41t, 2.44t, 2.46t, 2.48t, 2.53t,
 2.58t, 2.65t, 2.69t, 2.70t
- Zinc-germanium diphosphide (ZnGeP_2), 2.41t,
 2.45t, 2.48t, 2.53t, 2.58t, 2.65t
- Zirconia, cubic ($\text{ZrO}_2 \cdot 0.12\text{Y}_2\text{O}_3$), 2.41t, 2.44t,
 2.48t, 2.69t
- ZO-MOD BLACK, 6.56
- Zooplankton, in water, 1.14
- Z-scan:
 for nonlinear optical parameters, 19.9, 19.9f,
 19.10f
 and third-order optical nonlinearities,
 16.29–16.30, 16.30f

Third Edition

Sponsored by the Optical Society of America

HANDBOOK OF OPTICS

Volume V

*Atmospheric Optics, Modulators, Fiber Optics,
X-Ray and Neutron Optics*



Editor-in-Chief:
Michael Bass

Associate Editors:
Casimer M. DeCusatis
Jay M. Enoch
Vasudevan Lakshminarayanan
Guifang Li
Carolyn MacDonald
Virendra N. Mahajan
Eric Van Stryland

OSA[®]

HANDBOOK OF OPTICS

DO NOT DUPLICATE

ABOUT THE EDITORS

Editor-in-Chief: Dr. Michael Bass is professor emeritus at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Associate Editors:

Dr. Casimer M. DeCusatis is a distinguished engineer and technical executive with IBM Corporation, Poughkeepsie, New York.

Dr. Jay M. Enoch is dean emeritus and professor at the School of Optometry at the University of California, Berkeley.

Dr. Vasudevan Lakshminarayanan is professor of Optometry, Physics, and Electrical Engineering at the University of Waterloo, Ontario, Canada.

Dr. Guifang Li is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

Dr. Carolyn MacDonald is professor and chair of physics at the University at Albany, SUNY, and the director of the Albany Center for X-Ray Optics.

Dr. Virendra N. Mahajan is a distinguished scientist at The Aerospace Corporation.

Dr. Eric Van Stryland is a professor at CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida.

HANDBOOK OF OPTICS

Volume V

Atmospheric Optics, Modulators, Fiber Optics, X-Ray and Neutron Optic

THIRD EDITION

Sponsored by the
OPTICAL SOCIETY OF AMERICA

Michael Bass Editor-in-Chief

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

Carolyn MacDonald Associate Editor

*Department of Physics
University at Albany
Albany, New York*

Guifang Li Associate Editor

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

Casimer M. DeCusatis Associate Editor

*IBM Corporation
Poughkeepsie, New York*

Virendra N. Mahajan Associate Editor

*Aerospace Corporation
El Segundo, California*



New York Chicago San Francisco Lisbon London Madrid
Mexico City Milan New Delhi San Juan Seoul
Singapore Sydney Toronto

This page intentionally left blank.

DO NOT DUPLICATE

COVER ILLUSTRATIONS

Boadband supercontinuum. Generated in a photonic crystal fiber using a mode-locked Ti:Sapphire laser as pump source. The spectrum is much broader than seen in the photograph, extending from 400 nm to beyond 2 μm . (*Photo courtesy of the Optoelectronics Group, University of Bath.*)

Supernova remnant. A Chandra X-Ray Space Telescope image of the supernova remnant G292.0+1.8. The colors in the image encode the X-ray energies emitted by the supernova remnant; the center of G292.0+1.8 contains a region of high energy X-ray emission from the magnetized bubble of high-energy particles that surrounds the pulsar, a rapidly rotating neutron star that remained behind after the original, massive star exploded. (*This image is from NASA/CXC/Penn State/S.Park et al. and more detailed information can be found on the Chandra website: <http://chandra.harvard.edu/photo/2007/g292/>.*)

Crab Nebula. A Chandra X-Ray Space Telescope image of the Crab Nebula—the remains of a nearby supernova explosion first seen on Earth in 1054 AD. At the center of the bright nebula is a rapidly spinning neutron star, or pulsar, that emits pulses of radiation 30 times a second. (*This image is from NASA/CXC/ASU/J.Hester et al. and more detailed information can be found on the Chandra website: <http://chandra.harvard.edu/photo/2002/0052/>.*)

This page intentionally left blank.

DO NOT DUPLICATE

CONTENTS

Contributors	xix
Brief Contents of All Volumes	xxiii
Editors' Preface	xxix
Preface to Volume V	xxx
Glossary and Fundamental Constants	xxxiii

Part 1. Measurements

Chapter 1. Scatterometers	<i>John C. Stover</i>	1.3
----------------------------------	-----------------------	------------

- 1.1 Glossary / 1.3
- 1.2 Introduction / 1.3
- 1.3 Definitions and Specifications / 1.5
- 1.4 Instrument Configurations and Component Descriptions / 1.7
- 1.5 Instrumentation Issues / 1.11
- 1.6 Measurement Issues / 1.13
- 1.7 Incident Power Measurement, System Calibration, and Error Analysis / 1.14
- 1.8 Summary / 1.16
- 1.9 References / 1.16

Chapter 2. Spectroscopic Measurements	<i>Brian Henderson</i>	2.1
--	------------------------	------------

- 2.1 Glossary / 2.1
- 2.2 Introductory Comments / 2.2
- 2.3 Optical Absorption Measurements of Energy Levels / 2.2
- 2.4 The Homogeneous Lineshape of Spectra / 2.13
- 2.5 Absorption, Photoluminescence, and Radiative Decay Measurements / 2.19
- 2.6 References / 2.24

Part 2. Atmospheric Optics

Chapter 3. Atmospheric Optics	<i>Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman</i>	3.3
--------------------------------------	---	------------

- 3.1 Glossary / 3.3
- 3.2 Introduction / 3.4
- 3.3 Physical and Chemical Composition of the Standard Atmosphere / 3.6
- 3.4 Fundamental Theory of Interaction of Light with the Atmosphere / 3.11
- 3.5 Prediction of Atmospheric Optical Transmission: Computer Programs and Databases / 3.22
- 3.6 Atmospheric Optical Turbulence / 3.26
- 3.7 Examples of Atmospheric Optical Remote Sensing / 3.36
- 3.8 Meteorological Optics / 3.40
- 3.9 Atmospheric Optics and Global Climate Change / 3.43
- 3.10 Acknowledgments / 3.45
- 3.11 References / 3.45

Chapter 4. Imaging through Atmospheric Turbulence 4.1
Virendra N. Mahajan and Guang-ming Dai

- Abstract / 4.1
- 4.1 Glossary / 4.1
- 4.2 Introduction / 4.2
- 4.3 Long-Exposure Image / 4.3
- 4.4 Kolmogorov Turbulence and Atmospheric Coherence Length / 4.7
- 4.5 Application to Systems with Annular Pupils / 4.10
- 4.6 Modal Expansion of Aberration Function / 4.17
- 4.7 Covariance and Variance of Expansion Coefficients / 4.20
- 4.8 Angle of Arrival Fluctuations / 4.23
- 4.9 Aberration Variance and Approximate Strehl Ratio / 4.27
- 4.10 Modal Correction of Atmospheric Turbulence / 4.28
- 4.11 Short-Exposure Image / 4.31
- 4.12 Adaptive Optics / 4.35
- 4.13 Summary / 4.36
- 4.14 Acknowledgments / 4.37
- 4.15 References / 4.37

Chapter 5. Adaptive Optics *Robert Q. Fugate* 5.1

- 5.1 Glossary / 5.1
- 5.2 Introduction / 5.2
- 5.3 The Adaptive Optics Concept / 5.2
- 5.4 The Nature of Turbulence and Adaptive Optics Requirements / 5.5
- 5.5 AO Hardware and Software Implementation / 5.21
- 5.6 How to Design an Adaptive Optical System / 5.38
- 5.7 Acknowledgments / 5.46
- 5.8 References / 5.47

PART 3. Modulators

Chapter 6. Acousto-Optic Devices *I-Cheng Chang* 6.1

- 6.1 Glossary / 6.3
- 6.2 Introduction / 6.4
- 6.3 Theory of Acousto-Optic Interaction / 6.5
- 6.4 Acousto-Optic Materials / 6.16
- 6.5 Acousto-Optic Deflector / 6.22
- 6.6 Acousto-Optic Modulator / 6.31
- 6.7 Acousto-Optic Tunable Filter / 6.35
- 6.8 References / 6.45

Chapter 7. Electro-Optic Modulators *Georgianne M. Purvinis and Theresa A. Maldonado* 7.1

- 7.1 Glossary / 7.1
- 7.2 Introduction / 7.3
- 7.3 Crystal Optics and the Index Ellipsoid / 7.3
- 7.4 The Electro-Optic Effect / 7.6
- 7.5 Modulator Devices / 7.16
- 7.6 Applications / 7.36
- 7.7 Appendix: Euler Angles / 7.39
- 7.8 References / 7.40

Chapter 8. Liquid Crystals *Sebastian Gauza and Shin-Tson Wu* 8.1

- Abstract / 8.1
- 8.1 Glossary / 8.1

8.2	Introduction to Liquid Crystals /	8.2
8.3	Types of Liquid Crystals /	8.4
8.4	Liquid Crystals Phases /	8.8
8.5	Physical Properties /	8.13
8.6	Liquid Crystal Cells /	8.25
8.7	Liquid Crystals Displays /	8.29
8.8	Polymer/Liquid Crystal Composites /	8.36
8.9	Summary /	8.37
8.10	References /	8.38
8.11	Bibliography /	8.39

Part 4. Fiber Optics

Chapter 9. Optical Fiber Communication Technology and System Overview	9.3
--	------------

9.1	Introduction /	9.3
9.2	Basic Technology /	9.4
9.3	Receiver Sensitivity /	9.8
9.4	Bit Rate and Distance Limits /	9.12
9.5	Optical Amplifiers /	9.13
9.6	Fiber-Optic Networks /	9.14
9.7	Analog Transmission on Fiber /	9.15
9.8	Technology and Applications Directions /	9.17
9.9	References /	9.17

Chapter 10. Nonlinear Effects in Optical Fibers	10.1
--	-------------

10.1	Key Issues in Nonlinear Optics in Fibers /	10.1
10.2	Self- and Cross-Phase Modulation /	10.3
10.3	Stimulated Raman Scattering /	10.4
10.4	Stimulated Brillouin Scattering /	10.7
10.5	Four-Wave Mixing /	10.9
10.6	Conclusion /	10.11
10.7	References /	10.12

Chapter 11. Photonic Crystal Fibers	11.1
--	-------------

11.1	Glossary /	11.1
11.2	Introduction /	11.2
11.3	Brief History /	11.2
11.4	Fabrication Techniques /	11.4
11.5	Modeling and Analysis /	11.6
11.6	Characteristics of Photonic Crystal Cladding /	11.7
11.7	Linear Characteristics of Guidance /	11.11
11.8	Nonlinear Characteristics of Guidance /	11.22
11.9	Intrafiber Devices, Cutting, and Joining /	11.26
11.10	Conclusions /	11.28
11.11	Appendix /	11.28
11.12	References /	11.28

Chapter 12. Infrared Fibers	12.1
------------------------------------	-------------

12.1	Introduction /	12.1
12.2	Nonoxide and Heavy-Metal Oxide Glass IR Fibers /	12.3
12.3	Crystalline Fibers /	12.7
12.4	Hollow Waveguides /	12.10
12.5	Summary and Conclusions /	12.13
12.6	References /	12.13

Chapter 13. Sources, Modulators, and Detectors for Fiber Optic Communication Systems	<i>Elsa Garmire</i>	13.1
<hr/>		
13.1	Introduction / 13.1	
13.2	Double Heterostructure Laser Diodes / 13.3	
13.3	Operating Characteristics of Laser Diodes / 13.8	
13.4	Transient Response of Laser Diodes / 13.13	
13.5	Noise Characteristics of Laser Diodes / 13.18	
13.6	Quantum Well and Strained Lasers / 13.24	
13.7	Distributed Feedback and Distributed Bragg Reflector Lasers / 13.28	
13.8	Tunable Lasers / 13.32	
13.9	Light-Emitting Diodes / 13.36	
13.10	Vertical Cavity Surface-Emitting Lasers / 13.42	
13.11	Lithium Niobate Modulators / 13.48	
13.12	Electroabsorption Modulators / 13.55	
13.13	Electro-Optic and Electrorefractive Modulators / 13.61	
13.14	<i>PIN</i> Diodes / 13.63	
13.15	Avalanche Photodiodes, MSM Detectors, and Schottky Diodes / 13.71	
13.16	References / 13.74	
<hr/>		
Chapter 14. Optical Fiber Amplifiers	<i>John A. Buck</i>	14.1
<hr/>		
14.1	Introduction / 14.1	
14.2	Rare-Earth-Doped Amplifier Configuration and Operation / 14.2	
14.3	EDFA Physical Structure and Light Interactions / 14.4	
14.4	Other Rare-Earth Systems / 14.7	
14.5	Raman Fiber Amplifiers / 14.8	
14.6	Parametric Amplifiers / 14.10	
14.7	References / 14.11	
<hr/>		
Chapter 15. Fiber Optic Communication Links (Telecom, Datacom, and Analog)	<i>Casimer DeCusatis and Guifang Li</i>	15.1
<hr/>		
15.1	Figures of Merit / 15.2	
15.2	Link Budget Analysis: Installation Loss / 15.6	
15.3	Link Budget Analysis: Optical Power Penalties / 15.8	
15.4	References / 15.18	
<hr/>		
Chapter 16. Fiber-Based Couplers	<i>Daniel Nolan</i>	16.1
<hr/>		
16.1	Introduction / 16.1	
16.2	Achromaticity / 16.3	
16.3	Wavelength Division Multiplexing / 16.4	
16.4	$1 \times N$ Power Splitters / 16.4	
16.5	Switches and Attenuators / 16.4	
16.6	Mach-Zehnder Devices / 16.4	
16.7	Polarization Devices / 16.5	
16.8	Summary / 16.6	
16.9	References / 16.6	
<hr/>		
Chapter 17. Fiber Bragg Gratings	<i>Kenneth O. Hill</i>	17.1
<hr/>		
17.1	Glossary / 17.1	
17.2	Introduction / 17.1	
17.3	Photosensitivity / 17.2	
17.4	Properties of Bragg Gratings / 17.3	
17.5	Fabrication of Fiber Gratings / 17.4	
17.6	The Application of Fiber Gratings / 17.8	
17.7	References / 17.9	

Chapter 18. Micro-Optics-Based Components for Networking	18.1
<i>Joseph C. Palais</i>	
18.1 Introduction / 18.1	
18.2 Generalized Components / 18.1	
18.3 Network Functions / 18.2	
18.4 Subcomponents / 18.5	
18.5 Components / 18.9	
18.6 References / 18.12	
Chapter 19. Semiconductor Optical Amplifiers	19.1
<i>Jay M. Wiesenfeld and Leo H. Spiekman</i>	
19.1 Introduction / 19.1	
19.2 Device Basics / 19.2	
19.3 Fabrication / 19.15	
19.4 Device Characterization / 19.17	
19.5 Applications / 19.22	
19.6 Amplification of Signals / 19.22	
19.7 Switching and Modulation / 19.28	
19.8 Nonlinear Applications / 19.29	
19.9 Final Remarks / 19.36	
19.10 References / 19.36	
Chapter 20. Optical Time-Division Multiplexed Communication Networks	20.1
<i>Peter J. Delfyett</i>	
20.1 Glossary / 20.1	
20.2 Introduction / 20.3	
20.3 Multiplexing and Demultiplexing / 20.3	
20.4 Introduction to Device Technology / 20.12	
20.5 Summary and Future Outlook / 20.24	
20.6 Bibliography / 20.25	
Chapter 21. WDM Fiber-Optic Communication Networks	21.1
<i>Alan E. Willner, Changyuan Yu, Zhongqi Pan, and Yong Xie</i>	
21.1 Introduction / 21.1	
21.2 Basic Architecture of WDM Networks / 21.4	
21.3 Fiber System Impairments / 21.13	
21.4 Optical Modulation Formats for WDM Systems / 21.27	
21.5 Optical Amplifiers in WDM Networks / 21.37	
21.6 Summary / 21.44	
21.7 Acknowledgments / 21.44	
21.8 References / 21.44	
Chapter 22. Solitons in Optical Fiber Communication Systems	22.1
<i>Pavel V. Mamyshev</i>	
22.1 Introduction / 22.1	
22.2 Nature of the Classical Soliton / 22.2	
22.3 Properties of Solitons / 22.4	
22.4 Classical Soliton Transmission Systems / 22.5	
22.5 Frequency-Guiding Filters / 22.7	
22.6 Sliding Frequency-Guiding Filters / 22.8	
22.7 Wavelength Division Multiplexing / 22.9	
22.8 Dispersion-Managed Solitons / 22.12	
22.9 Wavelength-Division Multiplexed Dispersionmanaged Soliton Transmission / 22.15	
22.10 Conclusion / 22.17	
22.11 References / 22.17	

Chapter 23. Fiber-Optic Communication Standards 23.1
Casimer DeCusatis

- 23.1 Introduction / 23.1
- 23.2 ESCON / 23.1
- 23.3 FDDI / 23.2
- 23.4 Fibre Channel Standard / 23.4
- 23.5 ATM/SONET / 23.6
- 23.6 Ethernet / 23.7
- 23.7 Infiniband / 23.8
- 23.8 References / 23.8

Chapter 24. Optical Fiber Sensors 24.1
Richard O. Claus, Ignacio Matias, and Francisco Arregui

- 24.1 Introduction / 24.1
- 24.2 Extrinsic Fabry-Perot Interferometric Sensors / 24.2
- 24.3 Intrinsic Fabry-Perot Interferometric Sensors / 24.4
- 24.4 Fiber Bragg Grating Sensors / 24.5
- 24.5 Long-Period Grating Sensors / 24.8
- 24.6 Comparison of Sensing Schemes / 24.13
- 24.7 Conclusion / 24.13
- 24.8 References / 24.13
- 24.9 Further Reading / 24.14

Chapter 25. High-Power Fiber Lasers and Amplifiers 25.1
Timothy S. McComb, Martin C. Richardson, and Michael Bass

- 25.1 Glossary / 25.1
- 25.2 Introduction / 25.3
- 25.3 Fiber Laser Limitations / 25.6
- 25.4 Fiber Laser Fundamentals / 25.7
- 25.5 Fiber Laser Architectures / 25.9
- 25.6 LMA Fiber Designs / 25.18
- 25.7 Active Fiber Dopants / 25.22
- 25.8 Fiber Fabrication and Materials / 25.26
- 25.9 Spectral and Temporal Modalities / 25.29
- 25.10 Conclusions / 25.33
- 25.11 References / 25.33

PART 5. X-Ray and Neutron Optics

SUBPART 5.1. INTRODUCTION AND APPLICATIONS

Chapter 26. An Introduction to X-Ray and Neutron Optics 26.5
Carolyn MacDonald

- 26.1 History / 26.5
- 26.2 X-Ray Interaction with Matter / 26.6
- 26.3 Optics Choices / 26.7
- 26.4 Focusing and Collimation / 26.9
- 26.5 References / 26.11

Chapter 27. Coherent X-Ray Optics and Microscopy 27.1
Qun Shen

- 27.1 Glossary / 27.1
- 27.2 Introduction / 27.2
- 27.3 Fresnel Wave Propagation / 27.2
- 27.4 Unified Approach for Near- and Far-Field Diffraction / 27.2

- 27.5 Coherent Diffraction Microscopy / 27.4
 27.6 Coherence Preservation in X-Ray Optics / 27.5
 27.7 References / 27.5

Chapter 28. Requirements for X-Ray Diffraction *Scott T. Misture* 28.1

- 28.1 Introduction / 28.1
 28.2 Slits / 28.1
 28.3 Crystal Optics / 28.3
 28.4 Multilayer Optics / 28.5
 28.5 Capillary and Polycapillary Optics / 28.5
 28.6 Diffraction and Fluorescence Systems / 28.5
 28.7 X-Ray Sources and Microsources / 28.7
 28.8 References / 28.7

Chapter 29. Requirements for X-Ray Fluorescence *Walter Gibson and George Havrilla* 29.1

- 29.1 Introduction / 29.1
 29.2 Wavelength-Dispersive X-Ray Fluorescence (WDXRF) / 29.2
 29.3 Energy-Dispersive X-Ray Fluorescence (EDXRF) / 29.3
 29.4 References / 29.12

Chapter 30. Requirements for X-Ray Spectroscopy *Dirk Lützenkirchen-Hecht and Ronald Frahm* 30.1

- 30.1 References / 30.5

Chapter 31. Requirements for Medical Imaging and X-Ray Inspection *Douglas Pfeiffer* 31.1

- 31.1 Introduction to Radiography and Tomography / 31.1
 31.2 X-Ray Attenuation and Image Formation / 31.1
 31.3 X-Ray Detectors and Image Receptors / 31.4
 31.4 Tomography / 31.5
 31.5 Computed Tomography / 31.5
 31.6 Digital Tomosynthesis / 31.7
 31.7 Digital Displays / 31.8
 31.8 Conclusion / 31.9
 31.9 References / 31.10

Chapter 32. Requirements for Nuclear Medicine *Lars R. Furenlid* 32.1

- 32.1 Introduction / 32.1
 32.2 Projection Image Acquisition / 32.2
 32.3 Information Content in SPECT / 32.3
 32.4 Requirements for Optics For SPECT / 32.4
 32.5 References / 32.4

Chapter 33. Requirements for X-Ray Astronomy *Scott O. Rohrbach* 33.1

- 33.1 Introduction / 33.1
 33.2 Trade-Offs / 33.2
 33.3 Summary / 33.4

Chapter 34. Extreme Ultraviolet Lithography *Franco Cerrina and Fan Jiang* 34.1

- 34.1 Introduction / 34.1
 34.2 Technology / 34.2

- 34.3 Outlook / 34.5
- 34.4 Acknowledgments / 34.6
- 34.5 References / 34.7

Chapter 35. Ray Tracing of X-Ray Optical Systems *Franco Cerrina and Manuel Sanchez del Rio* 35.1

- 35.1 Introduction / 35.1
- 35.2 The Conceptual Basis of SHADOW / 35.2
- 35.3 Interfaces and Extensions of SHADOW / 35.3
- 35.4 Examples / 35.4
- 35.5 Conclusions and Future / 35.5
- 35.6 References / 35.6

Chapter 36. X-Ray Properties of Materials *Eric M. Gullikson* 36.1

- 36.1 X-Ray and Neutron Optics / 36.2
- 36.2 Electron Binding Energies, Principal K- and L-Shell Emission Lines,
and Auger Electron Energies / 36.3
- 36.3 References / 36.10

SUBPART 5.2. REFRACTIVE AND INTERFERENCE OPTICS

Chapter 37. Refractive X-Ray Lenses *Bruno Lengeler and Christian G. Schroer* 37.3

- 37.1 Introduction / 37.3
- 37.2 Refractive X-Ray Lenses with Rotationally Parabolic Profile / 37.4
- 37.3 Imaging with Parabolic Refractive X-Ray Lenses / 37.6
- 37.4 Microfocusing with Parabolic Refractive X-Ray Lenses / 37.7
- 37.5 Prefocusing and Collimation with Parabolic Refractive X-Ray Lenses / 37.8
- 37.6 Nanofocusing Refractive X-Ray Lenses / 37.8
- 37.7 Conclusion / 37.11
- 37.8 References / 37.11

Chapter 38. Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region *Malcolm R. Howells* 38.1

- 38.1 Introduction / 38.1
- 38.2 Diffraction Properties / 38.1
- 38.3 Focusing Properties / 38.3
- 38.4 Dispersion Properties / 38.6
- 38.5 Resolution Properties / 38.7
- 38.6 Efficiency / 38.8
- 38.7 References / 38.8

Chapter 39. Crystal Monochromators and Bent Crystals *Peter Siddons* 39.1

- 39.1 Crystal Monochromators / 39.1
- 39.2 Bent Crystals / 39.5
- 39.3 References / 39.6

Chapter 40. Zone Plates *Alan Michette* 40.1

- 40.1 Introduction / 40.1
- 40.2 Geometry of a Zone Plate / 40.1
- 40.3 Zone Plates as Thin Lenses / 40.3
- 40.4 Diffraction Efficiencies of Zone Plates / 40.4
- 40.5 Manufacture of Zone Plates / 40.8

- 40.6 Bragg-Fresnel Lenses / 40.9
 40.7 References / 40.10

Chapter 41. Multilayers *Eberhard Spiller* 41.1

- 41.1 Glossary / 41.1
 41.2 Introduction / 41.1
 41.3 Calculation of Multilayer Properties / 41.3
 41.4 Fabrication Methods and Performance / 41.4
 41.5 Multilayers for Diffractive Imaging / 41.9
 41.6 References / 41.10

Chapter 42. Nanofocusing of Hard X-Rays with Multilayer Laue Lenses *Albert T. Macrander, Hanfei Yan, Hyon Chol Kang, Jörg Maser, Chian Liu, Ray Conley, and G. Brian Stephenson* 42.1

- Abstract / 42.1
 42.1 Introduction / 42.2
 42.2 MLL Concept and Volume Diffraction Calculations / 42.4
 42.3 Magnetron-Sputtered MLLs / 42.5
 42.4 Instrumental Beamline Arrangement and Measurements / 42.9
 42.5 Takagi-Taupin Calculations / 42.12
 42.6 Wedged MLLs / 42.12
 42.7 MMLs with Curved Interfaces / 42.14
 42.8 MLL Prospects / 42.15
 42.9 Summary / 42.17
 42.10 Acknowledgments / 42.17
 42.11 References / 42.18

Chapter 43. Polarizing Crystal Optics *Qun Shen* 43.1

- 43.1 Introduction / 43.1
 43.2 Linear Polarizers / 43.2
 43.3 Linear Polarization Analyzers / 43.4
 43.4 Phase Plates for Circular Polarization / 43.5
 43.5 Circular Polarization Analyzers / 43.6
 43.6 Acknowledgments / 43.8
 43.7 References / 43.8

SUBPART 5.3. REFLECTIVE OPTICS

Chapter 44. Image Formation with Grazing Incidence Optics *James E. Harvey* 44.3

- 44.1 Glossary / 44.3
 44.2 Introduction to X-Ray Mirrors / 44.3
 44.3 Optical Design and Residual Aberrations of Grazing Incidence Telescopes / 44.6
 44.4 Image Analysis for Grazing Incidence X-Ray Optics / 44.12
 44.5 Validation of Image Analysis for Grazing Incidence X-Ray Optics / 44.16
 44.6 References / 44.18

Chapter 45. Aberrations for Grazing Incidence Optics *Timo T. Saha* 45.1

- 45.1 Grazing Incidence Telescopes / 45.1
 45.2 Surface Equations / 45.1
 45.3 Transverse Ray Aberration Expansions / 45.3
 45.4 Curvature of the Best Focal Surface / 45.5
 45.5 Aberration Balancing / 45.5
 45.6 On-Axis Aberrations / 45.6
 45.7 References / 45.8

Chapter 46. X-Ray Mirror Metrology	<i>Peter Z. Takacs</i>	46.1
<hr/>		
46.1	Glossary / 46.1	
46.2	Introduction / 46.1	
46.3	Surface Finish Metrology / 46.2	
46.4	Surface Figure Metrology / 46.3	
46.5	Practical Profile Analysis Considerations / 46.6	
46.6	References / 46.12	
Chapter 47. Astronomical X-Ray Optics	<i>Marshall K. Joy and Brian D. Ramsey</i>	47.1
<hr/>		
47.1	Introduction / 47.1	
47.2	Wolter X-Ray Optics / 47.2	
47.3	Kirkpatrick-Baez Optics / 47.7	
47.4	Hard X-Ray Optics / 47.9	
47.5	Toward Higher Angular Resolution / 47.10	
47.6	References / 47.11	
Chapter 48. Multifoil X-Ray Optics	<i>Ladislav Pina</i>	48.1
<hr/>		
48.1	Introduction / 48.1	
48.2	Grazing Incidence Optics / 48.1	
48.3	Multifoil Lobster-Eye Optics / 48.2	
48.4	Multifoil Kirkpatrick-Baez Optics / 48.3	
48.5	Summary / 48.4	
48.6	References / 48.4	
Chapter 49. Pore Optics	<i>Marco W. Beijersbergen</i>	49.1
<hr/>		
49.1	Introduction / 49.1	
49.2	Glass Micropore Optics / 49.1	
49.3	Silicon Pore Optics / 49.6	
49.4	Micromachined Silicon / 49.7	
49.5	References / 49.7	
Chapter 50. Adaptive X-Ray Optics	<i>Ali Khounsary</i>	50.1
<hr/>		
50.1	Introduction / 50.1	
50.2	Adaptive Optics in X-Ray Astronomy / 50.2	
50.3	Active and Adaptive Optics for Synchrotron- and Lab-Based X-Ray Sources / 50.2	
50.4	Conclusions / 50.8	
50.5	References / 50.8	
Chapter 51. The Schwarzschild Objective	<i>Franco Cerrina</i>	51.1
<hr/>		
51.1	Introduction / 51.1	
51.2	Applications to X-Ray Domain / 51.3	
51.3	References / 51.5	
Chapter 52. Single Capillaries	<i>Donald H. Bilderback and Sterling W. Cornaby</i>	52.1
<hr/>		
52.1	Background / 52.1	
52.2	Design Parameters / 52.1	
52.3	Fabrication / 52.4	
52.4	Applications of Single-Bounce Capillary Optics / 52.5	
52.5	Applications of Condensing Capillary Optics / 52.6	
52.6	Conclusions / 52.6	

- 52.7 Acknowledgments / 52.6
 52.8 References / 52.6

Chapter 53. Polycapillary X-Ray Optics *Carolyn MacDonald and
 Walter Gibson* **53.1**

- 53.1 Introduction / 53.1
 53.2 Simulations and Defect Analysis / 53.3
 53.3 Radiation Resistance / 53.5
 53.4 Alignment and Measurement / 53.5
 53.5 Collimation / 53.8
 53.6 Focusing / 53.9
 53.7 Applications / 53.10
 53.8 Summary / 53.19
 53.9 Acknowledgments / 53.19
 53.10 References / 53.19

SUBPART 5.4. X-RAY SOURCES

Chapter 54. X-Ray Tube Sources *Susanne M. Lee and
 Carolyn MacDonald* **54.3**

- 54.1 Introduction / 54.3
 54.2 Spectra / 54.4
 54.3 Cathode Design and Geometry / 54.10
 54.4 Effect of Anode Material, Geometry, and Source Size on Intensity and Brightness / 54.11
 54.5 General Optimization / 54.15
 54.6 References / 54.17

Chapter 55. Synchrotron Sources *Steven L. Hulbert and
 Gwyn P. Williams* **55.1**

- 55.1 Introduction / 55.1
 55.2 Theory of Synchrotron Radiation Emission / 55.2
 55.3 Insertion Devices (Undulators and Wigglers) / 55.9
 55.4 Coherence of Synchrotron Radiation Emission in the Long Wavelength Limit / 55.17
 55.5 Conclusion / 55.20
 55.6 References / 55.20

Chapter 56. Laser-Generated Plasmas *Alan Michette* **56.1**

- 56.1 Introduction / 56.1
 56.2 Characteristic Radiation / 56.2
 56.3 Bremsstrahlung / 56.8
 56.4 Recombination Radiation / 56.10
 56.5 References / 56.10

Chapter 57. Pinch Plasma Sources *Victor Kantsyrev* **57.1**

- 57.1 Introduction / 57.1
 57.2 Types of Z-Pinch Radiation Sources / 57.2
 57.3 Choice of Optics for Z-Pinch Sources / 57.4
 57.4 References / 57.5

Chapter 58. X-Ray Lasers *Greg Tallents* **58.1**

- 58.1 Free-Electron Lasers / 58.1
 58.2 High Harmonic Production / 58.2
 58.3 Plasma-Based EUV Lasers / 58.2
 58.4 References / 58.4

Chapter 59. Inverse Compton X-Ray Sources *Frank Carroll* 59.1

- 59.1 Introduction / 59.1
 59.2 Inverse Compton Calculations / 59.2
 59.3 Practical Devices / 59.2
 59.4 Applications / 59.3
 59.5 Industrial/Military/Crystallographic Uses / 59.4
 59.6 References / 59.4

SUBPART 5.5. X-RAY DETECTORS
Chapter 60. Introduction to X-Ray Detectors *Walter Gibson and Peter Siddons* 60.3

- 60.1 Introduction / 60.3
 60.2 Detector Type / 60.3
 60.3 Summary / 60.9
 60.4 References / 60.10

Chapter 61. Advances in Imaging Detectors *Aaron Couture* 61.1

- 61.1 Introduction / 61.1
 61.2 Flat-Panel Detectors / 61.3
 61.3 CCD Detectors / 61.7
 61.4 Conclusion / 61.8
 61.5 References / 61.8

Chapter 62. X-Ray Spectral Detection and Imaging *Eric Lifshin* 62.1

- 62.1 References / 62.6

SUBPART 5.6. NEUTRON OPTICS AND APPLICATIONS
Chapter 63. Neutron Optics *David Mildner* 63.3

- 63.1 Neutron Physics / 63.3
 63.2 Scattering Lengths and Cross Sections / 63.5
 63.3 Neutron Sources / 63.12
 63.4 Neutron Optical Devices / 63.15
 63.5 Refraction and Reflection / 63.19
 63.6 Diffraction and Interference / 63.23
 63.7 Polarization Techniques / 63.27
 63.8 Neutron Detection / 63.31
 63.9 References / 63.35

Chapter 64. Grazing-Incidence Neutron Optics *Mikhail Gubarev and Brian Ramsey* 64.1

- 64.1 Introduction / 64.1
 64.2 Total External Reflection / 64.1
 64.3 Diffractive Scattering and Mirror Surface Roughness Requirements / 64.2
 64.4 Imaging Focusing Optics / 64.3
 64.5 References / 64.7

CONTRIBUTORS

- Francisco Arregui** *Public University Navarra, Pamplona, Spain* (CHAP. 24)
- Michael Bass** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 25)
- Marco W. Beijersbergen** *Cosine Research B.V./Cosine Science & Computing B.V., Leiden University, Leiden, Netherlands* (CHAP. 49)
- Donald H. Bilderback** *Cornell High Energy Synchrotron Source, School of Applied and Engineering Physics, Cornell University, Ithaca, New York* (CHAP. 52)
- John A. Buck** *Georgia Institute of Technology, School of Electrical and Computer Engineering, Atlanta, Georgia* (CHAPS. 10, 14)
- Frank Carroll** *MXISystems, Nashville, Tennessee* (CHAP. 59)
- Franco Cerrina** *Department of Electrical and Computer Engineering, University of Wisconsin, Madison, Wisconsin* (CHAPS. 34, 35, 51)
- I-Cheng Chang** *Accord Optics, Sunnyvale, California* (CHAP. 6)
- James H. Churnside** *National Oceanic and Atmospheric Administration, Earth System Research Laboratory, Boulder, Colorado* (CHAP. 3)
- Richard O. Claus** *Virginia Tech, Blacksburg, Virginia* (CHAP. 24)
- Ray Conley** *X-Ray Science Division, Argonne National Laboratory, Argonne, Illinois, and National Synchrotron Light Source II, Brookhaven National Laboratory, Upton, New York* (CHAP. 42)
- Sterling W. Cornaby** *Cornell High Energy Synchrotron Source, School of Applied and Engineering Physics, Cornell University Ithaca, New York* (CHAP. 52)
- Aaron Couture** *GE Global Research Center, Niskayuna, New York* (CHAP. 61)
- Guang-ming Dai** *Laser Vision Correction Group, Advanced Medical Optics, Milpitas, California* (CHAP. 4)
- Casimer DeCusatis** *IBM Corporation, Poughkeepsie, New York* (CHAPS. 15, 23)
- Peter J. Delfyett** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 20)
- Ronald Frahm** *Bergische Universität Wuppertal, Wuppertal, Germany* (CHAP. 30)
- Robert O. Fugate** *Starfire Optical Range, Directed Energy Directorate, Air Force Research Laboratory, Kirtland Air Force Base, New Mexico* (CHAP. 5)
- Lars R. Furenlid** *University of Arizona, Tucson, Arizona* (CHAP. 32)
- Elsa Garmire** *Dartmouth College, Hanover, New Hampshire* (CHAP. 13)
- Sebastian Gauza** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 8)
- Walter Gibson** *X-Ray Optical Systems, Inc., East Greenbush, New York* (CHAPS. 29, 53, 60)
- Mikhail Gubarev** *NASA/Marshall Space Flight Center, Huntsville, Alabama* (CHAP. 64)
- Eric M. Gullikson** *Center for X-Ray Optics, Lawrence Berkeley National Laboratory, Berkeley, California* (CHAP. 36)
- James A. Harrington** *Rutgers University, Piscataway, New Jersey* (CHAP. 12)
- James E. Harvey** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 44)
- George Havrilla** *Los Alamos National Laboratory, Los Alamos, New Mexico* (CHAP. 29)

- Brian Henderson** *Department of Physics and Applied Physics, University of Strathclyde, Glasgow, United Kingdom* (CHAP. 2)
- Kenneth O. Hill** *Communications Research Centre, Ottawa, Ontario, Canada, and Nu-Wave Photonics, Ottawa, Ontario, Canada* (CHAP. 17)
- Malcolm R. Howells** *Advanced Light Source, Lawrence Berkeley National Laboratory, Berkeley, California* (CHAP. 38)
- Steven L. Hulbert** *National Synchrotron Light Source, Brookhaven National Laboratory, Upton, New York* (CHAP. 55)
- Ira Jacobs** *The Bradley Department of Electrical and Computer Engineering, Virginia Polytechnic Institute and State University, Blacksburg, Virginia* (CHAP. 9)
- Fan Jiang** *Electrical and Computer Engineering & Center for Nano Technology, University of Wisconsin, Madison* (CHAP. 34)
- Marshall K. Joy** *National Aeronautics and Space Administration, Marshall Space Flight Center, Huntsville, Alabama* (CHAP. 47)
- Hyon Chol Kang** *Materials Science Division, Argonne National Laboratory, Argonne, Illinois, and Advanced Materials Engineering Department, Chosun University, Gwangju, Republic of Korea* (CHAP. 42)
- Victor Kantsyrev** *Physics Department, University of Nevada, Reno, Nevada* (CHAP. 57)
- Ali Khounsary** *Argonne National Laboratory, Argonne, Illinois* (CHAP. 50)
- Dennis K. Killinger** *Center for Laser Atmospheric Sensing, Department of Physics, University of South Florida, Tampa, Florida* (CHAP. 3)
- Susanne M. Lee** *GE Global Research, Nikayuna, New York* (CHAP. 54)
- Bruno Lengeler** *Physikalisches Institut, RWTH Aachen University, Aachen, Germany* (CHAP. 37)
- Guifang Li** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 15)
- Eric Lifshin** *College of Nanoscale Science and Engineering, University at Albany, Albany, New York* (CHAP. 62)
- Chian Liu** *X-Ray Science Division, Argonne National Laboratory, Argonne, Illinois* (CHAP. 42)
- Dirk Lützenkirchen-Hecht** *Bergische Universität Wuppertal, Wuppertal, Germany* (CHAP. 30)
- Carolyn MacDonald** *University at Albany, Albany, New York* (CHAPS. 26, 53, 54)
- Albert T. Macrander** *X-Ray Science Division, Argonne National Laboratory, Argonne, Illinois* (CHAP. 42)
- Virendra N. Mahajan** *The Aerospace Corporation, El Segundo, California* (CHAP. 4)
- Theresa A. Maldonado** *Department of Electrical and Computer Engineering, Texas A&M University, College Station, Texas* (CHAP. 7)
- Pavel V. Mamyshev** *Bell Laboratories—Lucent Technologies, Holmdel, New Jersey* (CHAP. 22)
- Jörg Maser** *X-Ray Science Division, Argonne National Laboratory, Argonne, Illinois, and Center for Nanoscale Materials, Argonne National Laboratory, Argonne, Illinois* (CHAP. 42)
- Ignacio Matias** *Public University Navarra, Pamplona, Spain* (CHAP. 24)
- Timothy S. McComb** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 25)
- Alan Michette** *King's College, London, United Kingdom* (CHAPS. 40, 56)
- David Mildner** *NIST Center for Neutron Research, National Institute of Standards and Technology, Gaithersburg, Maryland* (CHAP. 63)
- Scott T. Misture** *Kazuo Inamori School of Engineering, Alfred University, Alfred, New York* (CHAP. 28)
- Daniel Nolan** *Corning Inc., Corning, New York* (CHAP. 16)
- Joseph C. Palais** *Ira A. Fulton School of Engineering, Arizona State University, Tempe, Arizona* (CHAP. 18)

- Zhongqi Pan** *University of Louisiana at Lafayette, Lafayette, Louisiana* (CHAP. 21)
- Greg J. Pearce** *Max-Planck Institute for the Science of Light, Erlangen, Germany* (CHAP. 11)
- Douglas Pfeiffer** *Boulder Community Hospital, Boulder, Colorado* (CHAP. 31)
- Ladislav Pina** *Faculty of Nuclear Sciences and Physical Engineering, Czech Technical University, Prague, Holesovickach* (CHAP. 48)
- Georgianne M. Purvinis** *The Battelle Memorial Institute, Columbus, Ohio* (CHAP. 7)
- Brian D. Ramsey** *National Aeronautics and Space Administration, Marshall Space Flight Center, Huntsville, Alabama* (CHAPS. 47, 64)
- Martin C. Richardson** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 25)
- Scott O. Rohrbach** *Optics Branch, Goddard Space Flight Center, NASA, Greenbelt, Maryland* (CHAP. 33)
- Laurence S. Rothman** *Harvard-Smithsonian Center for Astrophysics, Atomic and Molecular Physics Division, Cambridge, Massachusetts* (CHAP. 3)
- Philip St. J. Russell** *Max-Planck Institute for the Science of Light, Erlangen, Germany* (CHAP. 11)
- Timo T. Saha** *NASA/Goddard Space Flight Center, Greenbelt, Maryland* (CHAP. 45)
- Manuel Sanchez del Rio** *European Synchrotron Radiation Facility, Grenoble, France* (CHAP. 35)
- Christian G. Schroer** *Institute of Structural Physics, TU Dresden, Dresden, Germany* (CHAP. 37)
- Qun Shen** *National Synchrotron Light Source II, Brookhaven National Laboratory, Upton, New York* (CHAPS. 27, 43)
- Peter Siddons** *National Synchrotron Light Source, Brookhaven National Laboratory, Upton, New York* (CHAPS. 39, 60)
- Leo H. Spiekman** *Alphion Corp., Princeton Junction, New Jersey* (CHAP. 19)
- Eberhard Spiller** *Spiller X-Ray Optics, Livermore, California* (CHAP. 41)
- G. Brian Stephenson** *Center for Nanoscale Materials, Argonne National Laboratory, Argonne, Illinois, Materials Science Division, Argonne National Laboratory, Argonne, Illinois* (CHAP. 42)
- John C. Stover** *The Scatter Works, Inc., Tucson, Arizona* (CHAP. 1)
- Peter Z. Takacs** *Brookhaven National Laboratory, Upton, New York* (CHAP. 46)
- Greg Tallents** *University of York, York, United Kingdom* (CHAP. 58)
- Jay M. Wiesenfeld** *Bell Laboratories, Alcatel-Lucent, Murray Hill, New Jersey* (CHAP. 19)
- Gwyn P. Williams** *Free Electron Laser, Thomas Jefferson National Accelerator Facility, Newport News, Virginia* (CHAP. 55)
- Alan E. Willner** *University of Southern California, Los Angeles, California* (CHAP. 21)
- Shin-Tson Wu** *CREOL, The College of Optics and Photonics, University of Central Florida, Orlando, Florida* (CHAP. 8)
- Yong Xie** *Texas Instruments Inc., Dallas, Texas* (CHAP. 21)
- Hanfei Yan** *Center for Nanoscale Materials, Argonne National Laboratory, Argonne, Illinois, and National Synchrotron Light Source II, Brookhaven National Laboratory, Upton, New York* (CHAP. 42)
- Changyuan Yu** *National University of Singapore, and A *STAR Institute for Infocomm Research, Singapore* (CHAP. 21)

This page intentionally left blank.

DO NOT DUPLICATE

BRIEF CONTENTS OF ALL VOLUMES

VOLUME I. GEOMETRICAL AND PHYSICAL OPTICS, POLARIZED LIGHT, COMPONENTS AND INSTRUMENTS

PART 1. GEOMETRICAL OPTICS

Chapter 1. General Principles of Geometrical Optics *Douglas S. Goodman*

PART 2. PHYSICAL OPTICS

Chapter 2. Interference *John E. Greivenkamp*

Chapter 3. Diffraction *Arvind S. Marathay and John F. McCalmont*

Chapter 4. Transfer Function Techniques *Glenn D. Boreman*

Chapter 5. Coherence Theory *William H. Carter*

Chapter 6. Coherence Theory: Tools and Applications *Gisele Bennett, William T. Rhodes, and J. Christopher James*

Chapter 7. Scattering by Particles *Craig F. Bohren*

Chapter 8. Surface Scattering *Eugene L. Church and Peter Z. Takacs*

Chapter 9. Volume Scattering in Random Media *Aristide Dogariu and Jeremy Ellis*

Chapter 10. Optical Spectroscopy and Spectroscopic Lineshapes *Brian Henderson*

Chapter 11. Analog Optical Signal and Image Processing *Joseph W. Goodman*

PART 3. POLARIZED LIGHT

Chapter 12. Polarization *Jean M. Bennett*

Chapter 13. Polarizers *Jean M. Bennett*

Chapter 14. Mueller Matrices *Russell A. Chipman*

Chapter 15. Polarimetry *Russell A. Chipman*

Chapter 16. Ellipsometry *Rasheed M. A. Azzam*

PART 4. COMPONENTS

Chapter 17. Lenses *R. Barry Johnson*

Chapter 18. Afocal Systems *William B. Wetherell*

Chapter 19. Nondispersive Prisms *William L. Wolfe*

Chapter 20. Dispersive Prisms and Gratings *George J. Zissis*

Chapter 21. Integrated Optics *Thomas L. Koch, Frederick J. Leonberger, and Paul G. Suchoski*

Chapter 22. Miniature and Micro-Optics *Tom D. Milster and Tomasz S. Tkaczyk*

Chapter 23. Binary Optics *Michael W. Farn and Wilfrid B. Veldkamp*

Chapter 24. Gradient Index Optics *Duncan T. Moore*

PART 5. INSTRUMENTS

Chapter 25. Cameras *Norman Goldberg*

Chapter 26. Solid-State Cameras *Gerald C. Holst*

Chapter 27. Camera Lenses *Ellis Betensky, Melvin H. Kreitzer, and Jacob Moskovich*

Chapter 28. Microscopes *Rudolf Oldenbourg and Michael Shribak*

- Chapter 29. Reflective and Catadioptric Objectives *Lloyd Jones*
Chapter 30. Scanners *Leo Beiser and R. Barry Johnson*
Chapter 31. Optical Spectrometers *Brian Henderson*
Chapter 32. Interferometers *Parameswaran Hariharan*
Chapter 33. Holography and Holographic Instruments *Lloyd Huff*
Chapter 34. Xerographic Systems *Howard Stark*
Chapter 35. Principles of Optical Disk Data Storage *Masud Mansuripur*

VOLUME II. DESIGN, FABRICATION, AND TESTING; SOURCES AND DETECTORS; RADIOMETRY AND PHOTOMETRY

PART 1. DESIGN

- Chapter 1. Techniques of First-Order Layout *Warren J. Smith*
Chapter 2. Aberration Curves in Lens Design *Donald C. O'Shea and Michael E. Harrigan*
Chapter 3. Optical Design Software *Douglas C. Sinclair*
Chapter 4. Optical Specifications *Robert R. Shannon*
Chapter 5. Tolerancing Techniques *Robert R. Shannon*
Chapter 6. Mounting Optical Components *Paul R. Yoder, Jr.*
Chapter 7. Control of Stray Light *Robert P. Breault*
Chapter 8. Thermal Compensation Techniques *Philip J. Rogers and Michael Roberts*

PART 2. FABRICATION

- Chapter 9. Optical Fabrication *Michael P. Mandina*
Chapter 10. Fabrication of Optics by Diamond Turning *Richard L. Rhorer and Chris J. Evans*

PART 3. TESTING

- Chapter 11. Orthonormal Polynomials in Wavefront Analysis *Virendra N. Mahajan*
Chapter 12. Optical Metrology *Zacarias Malacara and Daniel Malacara-Hernández*
Chapter 13. Optical Testing *Daniel Malacara-Hernández*
Chapter 14. Use of Computer-Generated Holograms in Optical Testing *Katherine Creath and James C. Wyant*

PART 4. SOURCES

- Chapter 15. Artificial Sources *Anthony LaRocca*
Chapter 16. Lasers *William T. Silfvast*
Chapter 17. Light-Emitting Diodes *Roland H. Haitz, M. George Craford, and Robert H. Weissman*
Chapter 18. High-Brightness Visible LEDs *Winston V. Schoenfeld*
Chapter 19. Semiconductor Lasers *Pamela L. Derry, Luis Figueroa, and Chi-shain Hong*
Chapter 20. Ultrashort Optical Sources and Applications *Jean-Claude Diels and Ladan Arisian*
Chapter 21. Attosecond Optics *Zenghu Chang*
Chapter 22. Laser Stabilization *John L. Hall, Matthew S. Taubman, and Jun Ye*
Chapter 23. Quantum Theory of the Laser *János A. Bergou, Berthold-Georg Englert, Melvin Lax, Marian O. Scully, Herbert Walther, and M. Suhail Zubairy*

PART 5. DETECTORS

- Chapter 24. Photodetectors *Paul R. Norton*
Chapter 25. Photodetection *Abhay M. Joshi and Gregory H. Olsen*
Chapter 26. High-Speed Photodetectors *John E. Bowers and Yih G. Wey*
Chapter 27. Signal Detection and Analysis *John R. Willison*
Chapter 28. Thermal Detectors *William L. Wolfe and Paul W. Kruse*

PART 6. IMAGING DETECTORS

- Chapter 29. Photographic Films *Joseph H. Altman*
Chapter 30. Photographic Materials *John D. Baloga*

- Chapter 31. Image Tube Intensified Electronic Imaging *C. Bruce Johnson and Larry D. Owen*
 Chapter 32. Visible Array Detectors *Timothy J. Tredwell*
 Chapter 33. Infrared Detector Arrays *Lester J. Kozlowski and Walter F. Kosonocky*

PART 7. RADIOMETRY AND PHOTOMETRY

- Chapter 34. Radiometry and Photometry *Edward F. Zalewski*
 Chapter 35. Measurement of Transmission, Absorption, Emission, and Reflection *James M. Palmer*
 Chapter 36. Radiometry and Photometry: Units and Conversions *James M. Palmer*
 Chapter 37. Radiometry and Photometry for Vision Optics *Yoshi Ohno*
 Chapter 38. Spectroradiometry *Carolyn J. Sher DeCusatis*
 Chapter 39. Nonimaging Optics: Concentration and Illumination *William Cassarly*
 Chapter 40. Lighting and Applications *Anurag Gupta and R. John Koshel*

VOLUME III. VISION AND VISION OPTICS

- Chapter 1. Optics of the Eye *Neil Charman*
 Chapter 2. Visual Performance *Wilson S. Geisler and Martin S. Banks*
 Chapter 3. Psychophysical Methods *Denis G. Pelli and Bart Farell*
 Chapter 4. Visual Acuity and Hyperacuity *Gerald Westheimer*
 Chapter 5. Optical Generation of the Visual Stimulus *Stephen A. Burns and Robert H. Webb*
 Chapter 6. The Maxwellian View with an Addendum on Apodization *Gerald Westheimer*
 Chapter 7. Ocular Radiation Hazards *David H. Sliney*
 Chapter 8. Biological Waveguides *Vasudevan Lakshminarayanan and Jay M. Enoch*
 Chapter 9. The Problem of Correction for the Stiles-Crawford Effect of the First Kind in Radiometry and Photometry, a Solution *Jay M. Enoch and Vasudevan Lakshminarayanan*
 Chapter 10. Colorimetry *David H. Brainard and Andrew Stockman*
 Chapter 11. Color Vision Mechanisms *Andrew Stockman and David H. Brainard*
 Chapter 12. Assessment of Refraction and Refractive Errors and Their Influence on Optical Design *B. Ralph Chou*
 Chapter 13. Binocular Vision Factors That Influence Optical Design *Clifton Schor*
 Chapter 14. Optics and Vision of the Aging Eye *John S. Werner, Brooke E. Scheffrin, and Arthur Bradley*
 Chapter 15. Adaptive Optics in Retinal Microscopy and Vision *Donald T. Miller and Austin Roorda*
 Chapter 16. Refractive Surgery, Correction of Vision, PRK and LASIK *L. Diaz-Santana and Harilaos Ginis*
 Chapter 17. Three-Dimensional Confocal Microscopy of the Living Human Cornea *Barry R. Masters*
 Chapter 18. Diagnostic Use of Optical Coherence Tomography in the Eye *Johannes F. de Boer*
 Chapter 19. Gradient Index Optics in the Eye *Barbara K. Pierscionek*
 Chapter 20. Optics of Contact Lenses *Edward S. Bennett*
 Chapter 21. Intraocular Lenses *Jim Schwiegerling*
 Chapter 22. Displays for Vision Research *William Cowan*
 Chapter 23. Vision Problems at Computers *Jeffrey Anshel and James E. Sheedy*
 Chapter 24. Human Vision and Electronic Imaging *Bernice E. Rogowitz, Thrasyvoulos N. Pappas, and Jan P. Allebach*
 Chapter 25. Visual Factors Associated with Head-Mounted Displays *Brian H. Tsou and Martin Shenker*

VOLUME IV. OPTICAL PROPERTIES OF MATERIALS, NONLINEAR OPTICS, QUANTUM OPTICS

PART 1. PROPERTIES

- Chapter 1. Optical Properties of Water *Curtis D. Mobley*
 Chapter 2. Properties of Crystals and Glasses *William J. Tropf, Michael E. Thomas, and Eric W. Rogala*
 Chapter 3. Polymeric Optics *John D. Lytle*
 Chapter 4. Properties of Metals *Roger A. Paquin*

- Chapter 5. Optical Properties of Semiconductors *David G. Seiler, Stefan Zollner, Alain C. Diebold, and Paul M. Amirtharaj*
- Chapter 6. Characterization and Use of Black Surfaces for Optical Systems *Stephen M. Pompea and Robert P. Breault*
- Chapter 7. Optical Properties of Films and Coatings *Jerzy A. Dobrowolski*
- Chapter 8. Fundamental Optical Properties of Solids *Alan Miller*
- Chapter 9. Photonic Bandgap Materials *Pierre R. Villeneuve*

PART 2. NONLINEAR OPTICS

- Chapter 10. Nonlinear Optics *Chung L. Tang*
- Chapter 11. Coherent Optical Transients *Paul R. Berman and Duncan G. Steel*
- Chapter 12. Photorefractive Materials and Devices *Mark Cronin-Golomb and Marvin Klein*
- Chapter 13. Optical Limiting *David J. Hagan*
- Chapter 14. Electromagnetically Induced Transparency *Jonathan P. Marangos and Thomas Halfmann*
- Chapter 15. Stimulated Raman and Brillouin Scattering *John Reintjes and Mark Bashkansky*
- Chapter 16. Third-Order Optical Nonlinearities *Mansoor Sheik-Bahae and Michael P. Hasselbeck*
- Chapter 17. Continuous-Wave Optical Parametric Oscillators *Majid Ebrahim-Zadeh*
- Chapter 18. Nonlinear Optical Processes for Ultrashort Pulse Generation *Uwe Siegner and Ursula Keller*
- Chapter 19. Laser-Induced Damage to Optical Materials *Marion J. Soileau*

PART 3. QUANTUM AND MOLECULAR OPTICS

- Chapter 20. Laser Cooling and Trapping of Atoms *Harold J. Metcalf and Peter van der Straten*
- Chapter 21. Strong Field Physics *Todd Ditmire*
- Chapter 22. Slow Light Propagation in Atomic and Photonic Media *Jacob B. Khurgin*
- Chapter 23. Quantum Entanglement in Optical Interferometry *Hwang Lee, Christoph F. Wildfeuer, Sean D. Huver, and Jonathan P. Dowling*

VOLUME V. ATMOSPHERIC OPTICS, MODULATORS, FIBER OPTICS, X-RAY AND NEUTRON OPTICS

PART 1. MEASUREMENTS

- Chapter 1. Scatterometers *John C. Stover*
- Chapter 2. Spectroscopic Measurements *Brian Henderson*

PART 2. ATMOSPHERIC OPTICS

- Chapter 3. Atmospheric Optics *Dennis K. Killinger, James H. Churnside, and Laurence S. Rothman*
- Chapter 4. Imaging through Atmospheric Turbulence *Virendra N. Mahajan and Guang-ming Dai*
- Chapter 5. Adaptive Optics *Robert Q. Fugate*

PART 3. MODULATORS

- Chapter 6. Acousto-Optic Devices *I-Cheng Chang*
- Chapter 7. Electro-Optic Modulators *Georgeanne M. Purvinis and Theresa A. Maldonado*
- Chapter 8. Liquid Crystals *Sebastian Gauza and Shin-Tson Wu*

PART 4. FIBER OPTICS

- Chapter 9. Optical Fiber Communication Technology and System Overview *Ira Jacobs*
- Chapter 10. Nonlinear Effects in Optical Fibers *John A. Buck*
- Chapter 11. Photonic Crystal Fibers *Philip St. J. Russell and Greg J. Pearce*
- Chapter 12. Infrared Fibers *James A. Harrington*
- Chapter 13. Sources, Modulators, and Detectors for Fiber Optic Communication Systems *Elsa Garmire*
- Chapter 14. Optical Fiber Amplifiers *John A. Buck*

- Chapter 15. Fiber Optic Communication Links (Telecom, Datacom, and Analog) *Casimer DeCusatis and Guifang Li*
- Chapter 16. Fiber-Based Couplers *Daniel Nolan*
- Chapter 17. Fiber Bragg Gratings *Kenneth O. Hill*
- Chapter 18. Micro-Optics-Based Components for Networking *Joseph C. Palais*
- Chapter 19. Semiconductor Optical Amplifiers *Jay M. Wiesenfeld and Leo H. Spiekman*
- Chapter 20. Optical Time-Division Multiplexed Communication Networks *Peter J. Delfyett*
- Chapter 21. WDM Fiber-Optic Communication Networks *Alan E. Willner, Changyuan Yu, Zhongqi Pan, and Yong Xie*
- Chapter 22. Solitons in Optical Fiber Communication Systems *Pavel V. Mamyshev*
- Chapter 23. Fiber-Optic Communication Standards *Casimer DeCusatis*
- Chapter 24. Optical Fiber Sensors *Richard O. Claus, Ignacio Matias, and Francisco Arregui*
- Chapter 25. High-Power Fiber Lasers and Amplifiers *Timothy S. McComb, Martin C. Richardson, and Michael Bass*

PART 5. X-RAY AND NEUTRON OPTICS

Subpart 5.1. Introduction and Applications

- Chapter 26. An Introduction to X-Ray and Neutron Optics *Carolyn MacDonald*
- Chapter 27. Coherent X-Ray Optics and Microscopy *Qun Shen*
- Chapter 28. Requirements for X-Ray Diffraction *Scott T. Misture*
- Chapter 29. Requirements for X-Ray Fluorescence *George J. Havrilla*
- Chapter 30. Requirements for X-Ray Spectroscopy *Dirk Lützenkirchen-Hecht and Ronald Frahm*
- Chapter 31. Requirements for Medical Imaging and X-Ray Inspection *Douglas Pfeiffer*
- Chapter 32. Requirements for Nuclear Medicine *Lars R. Furenlid*
- Chapter 33. Requirements for X-Ray Astronomy *Scott O. Rohrbach*
- Chapter 34. Extreme Ultraviolet Lithography *Franco Cerrina and Fan Jiang*
- Chapter 35. Ray Tracing of X-Ray Optical Systems *Franco Cerrina and M. Sanchez del Rio*
- Chapter 36. X-Ray Properties of Materials *Eric M. Gullikson*

Subpart 5.2. Refractive and Interference Optics

- Chapter 37. Refractive X-Ray Lenses *Bruno Lengeler and Christian G. Schroer*
- Chapter 38. Gratings and Monochromators in the VUV and Soft X-Ray Spectral Region
Malcolm R. Howells
- Chapter 39. Crystal Monochromators and Bent Crystals *Peter Siddons*
- Chapter 40. Zone Plates *Alan Michette*
- Chapter 41. Multilayers *Eberhard Spiller*
- Chapter 42. Nanofocusing of Hard X-Rays with Multilayer Laue Lenses *Albert T. Macrander, Hanfei Yan, Hyon Chol Kang, Jörg Maser, Chian Liu, Ray Conley, and G. Brian Stephenson*
- Chapter 43. Polarizing Crystal Optics *Qun Shen*

Subpart 5.3. Reflective Optics

- Chapter 44. Image Formation with Grazing Incidence Optics *James Harvey*
- Chapter 45. Aberrations for Grazing Incidence Optics *Timo T. Saha*
- Chapter 46. X-Ray Mirror Metrology *Peter Z. Takacs*
- Chapter 47. Astronomical X-Ray Optics *Marshall K. Joy and Brian D. Ramsey*
- Chapter 48. Multifoil X-Ray Optics *Ladislav Pina*
- Chapter 49. Pore Optics *Marco Beijersbergen*
- Chapter 50. Adaptive X-Ray Optics *Ali Khounsary*
- Chapter 51. The Schwarzschild Objective *Franco Cerrina*
- Chapter 52. Single Capillaries *Donald H. Bilderback and Sterling W. Cornaby*
- Chapter 53. Polycapillary X-Ray Optics *Carolyn MacDonald and Walter Gibson*

Subpart 5.4. X-Ray Sources

- Chapter 54. X-Ray Tube Sources *Susanne M. Lee and Carolyn MacDonald*
Chapter 55. Synchrotron Sources *Steven L. Hulbert and Gwyn P. Williams*
Chapter 56. Laser-Generated Plasmas *Alan Michette*
Chapter 57. Pinch Plasma Sources *Victor Kantsyrev*
Chapter 58. X-Ray Lasers *Greg Tallents*
Chapter 59. Inverse Compton X-Ray Sources *Frank Carroll*

Subpart 5.5. X-Ray Detectors

- Chapter 60. Introduction to X-Ray Detectors *Walter M. Gibson and Peter Siddons*
Chapter 61. Advances in Imaging Detectors *Aaron Couture*
Chapter 62. X-Ray Spectral Detection and Imaging *Eric Lifshin*

Subpart 5.6. Neutron Optics and Applications

- Chapter 63. Neutron Optics *David Mildner*
Chapter 64. Grazing-Incidence Neutron Optics *Mikhail Gubarev and Brian Ramsey*

DO NOT DUPLICATE

EDITORS' PREFACE

The third edition of the *Handbook of Optics* is designed to pull together the dramatic developments in both the basic and applied aspects of the field while retaining the archival, reference book value of a handbook. This means that it is much more extensive than either the first edition, published in 1978, or the second edition, with Volumes I and II appearing in 1995 and Volumes III and IV in 2001. To cover the greatly expanded field of optics, the *Handbook* now appears in five volumes. Over 100 authors or author teams have contributed to this work.

Volume I is devoted to the fundamentals, components, and instruments that make optics possible. Volume II contains chapters on design, fabrication, testing, sources of light, detection, and a new section devoted to radiometry and photometry. Volume III concerns vision optics only and is printed entirely in color. In Volume IV there are chapters on the optical properties of materials, nonlinear, quantum and molecular optics. Volume V has extensive sections on fiber optics and x ray and neutron optics, along with shorter sections on measurements, modulators, and atmospheric optical properties and turbulence. Several pages of color inserts are provided where appropriate to aid the reader. A purchaser of the print version of any volume of the *Handbook* will be able to download a digital version containing all of the material in that volume in PDF format to one computer (see download instructions on bound-in card). The combined index for all five volumes can be downloaded from www.HandbookofOpticsOnline.com.

It is possible by careful selection of what and how to present that the third edition of the *Handbook* could serve as a text for a comprehensive course in optics. In addition, students who take such a course would have the *Handbook* as a career-long reference.

Topics were selected by the editors so that the *Handbook* could be a desktop (bookshelf) general reference for the parts of optics that had matured enough to warrant archival presentation. New chapters were included on topics that had reached this stage since the second edition, and existing chapters from the second edition were updated where necessary to provide this compendium. In selecting subjects to include, we also had to select which subjects to leave out. The criteria we applied were: (1) was it a specific application of optics rather than a core science or technology and (2) was it a subject in which the role of optics was peripheral to the central issue addressed. Thus, such topics as medical optics, laser surgery, and laser materials processing were not included. While applications of optics are mentioned in the chapters there is no space in the *Handbook* to include separate chapters devoted to all of the myriad uses of optics in today's world. If we had, the third edition would be much longer than it is and much of it would soon be outdated. We designed the third edition of the *Handbook of Optics* so that it concentrates on the principles of optics that make applications possible.

Authors were asked to try to achieve the dual purpose of preparing a chapter that was a worthwhile reference for someone working in the field and that could be used as a starting point to become acquainted with that aspect of optics. They did that and we thank them for the outstanding results seen throughout the *Handbook*. We also thank Mr. Taisuke Soda of McGraw-Hill for his help in putting this complex project together and Mr. Alan Tourtlotte and Ms. Susannah Lehman of the Optical Society of America for logistical help that made this effort possible.

We dedicate the third edition of the *Handbook of Optics* to all of the OSA volunteers who, since OSA's founding in 1916, give their time and energy to promoting the generation, application, archiving, and worldwide dissemination of knowledge in optics and photonics.

Michael Bass, Editor-in-Chief

Associate Editors:

Casimer M. DeCusatis

Jay M. Enoch

Vasudevan Lakshminarayanan

Guifang Li

Carolyn MacDonald

Virendra N. Mahajan

Eric Van Stryland

This page intentionally left blank.

DO NOT DUPLICATE

PREFACE TO VOLUME V

Volume V begins with Measurements, Atmospheric Optics, and Optical Modulators. There are chapters on scatterometers, spectroscopic measurements, transmission through the atmosphere, imaging through turbulence, and adaptive optics to overcome distortions as well as chapters on electro- and acousto-optic modulators and liquid crystal spatial light modulators. These are followed by the two main parts of this volume—Fiber Optics and X-Ray and Neutron Optics.

Optical fiber technology is truly an interdisciplinary field, incorporating aspects of solid-state physics, material science, and electrical engineering, among others. In the section on fiber optics, we introduce the fundamentals of optical fibers and cable assemblies, optical connectors, light sources, detectors, and related components. Assembly of the building blocks into optical networks required discussion of the unique requirements of digital versus analog links and telecommunication versus data communication networks. Issues such as optical link budget calculations, dispersion- or attenuation-limited links, and compliance with relevant industry standards are all addressed. Since one of the principle advantages of fiber optics is the ability to create high-bandwidth, long-distance interconnections, we also discuss the design and use of optical fiber amplifiers for different wavelength transmission windows. This leads to an understanding of the different network components which can be fabricated from optical fiber itself, such as splitters, combiners, fiber Bragg gratings, and other passive optical networking elements. We then provide a treatment of other important devices, including fiber sensors, fibers optimized for use in the infrared, micro-optic components for fiber networks and fiber lasers. Note that micro-optics for other applications are covered in Volume I of this *Handbook*. The physics of semiconductor lasers and photodetectors are presented in Volume II. Applications such as time or wavelength-division multiplexing networks provide their own challenges and are discussed in detail. High optical power applications lead us to a consideration of non-linear optical fiber properties. Advanced topics for high speed, future networks are described in this section, including polarization mode dispersion; readers interested in the physical optics underlying dispersion should consult Volume I of this *Handbook*. This section includes chapters on photonic crystal fibers (for a broader treatment of photonic bandgap materials, see Volume IV) and on the growing applications of optical fiber networks.

Part 5 of this volume discusses a variety of X-Ray and Neutron Optics and their use in a wide range of applications. Part 5.1 is an introduction to the use and properties of x rays. It begins with a short chapter summarizing x-ray interactions and optics, followed by a discussion of coherence effects, and then illustrations of application constraints to the use of optics in seven applications, ranging from materials analysis to medicine, astronomy, and chip manufacturing. Because modeling is an important tool for both optics development and system design, Part 5.1 continues with a discussion of optics simulations, followed by tables of materials properties in the x-ray regime. Parts 5.2 and 5.3 are devoted to the discussion of the three classes of x-ray optics. Part 5.2 covers refractive, interference, and diffractive optics, including gratings, crystals (flat, bent, and polarizing), zone plates, and Laue lenses. It also includes a discussion of multilayer coatings, which are based on interference, but often added to reflective x-ray optics. Reflective optics is the topic of Part 5.3. Since reflective optics in the x-ray regime are used primarily in grazing incidence, the first three chapters of Part 5.3 cover the theory of image formation, aberrations, and metrology of grazing incidence mirrors. This is followed with descriptions of mirrors for astronomy and microscopy, adaptive optics for high heat load synchrotron beam lines, glass capillary reflective optics, also generally used for beam lines, and array optics such as multifoils, pore optics, and polycapillaries. The best choice of optic for a particular function depends on the application requirements, but is also influenced by the properties of the available sources and detectors. Part 5.4 describes six different types of x-ray sources. This is followed by Part 5.5, which includes an introduction to detectors and in-depth

discussions of imaging and spectral detectors. Finally, Part 5.6 describes the similarities and differences in the use of comparable optics technologies with neutrons.

In 1998, Walter Gibson designed the expansion of the x-ray and neutron section of the second edition of the *Handbook* from its original single chapter form. The third edition of this section is dedicated to his memory.

Guifang Li, Casimer M. DeCusatis, Virendra N. Mahajan, and Carolyn MacDonald
Associate Editors

In Memoriam Walter Maxwell Gibson (November 11, 1930–May 15, 2009)

After a childhood in Southern Utah working as a sheepherder and stunt rider, Walt received his Ph.D. in nuclear chemistry under Nobel Laureate Glen Seaborg in 1956. He then spent 20 years at Bell Labs, where he did groundbreaking research in semiconductor detectors, particle-solid interactions, and in the development of ion beam techniques for material analysis. His interest in materials analysis and radiation detection naturally led him to an early and ongoing interest in developing x-ray analysis techniques, including early synchrotron beam line development. In 1970, he was named a fellow of the American Physical Society. In 1976, Walter was invited to chair the physics department of the University at Albany, SUNY (where he was fond of noting that they must have been confused as he had been neither an academic nor a physicist). He remained with the department for more than 25 years and was honored with the university's first named professorship, the James W. Corbett Distinguished Service Professor of Physics, in 1998. He later retired from the university to become the full-time chief technical officer of X-Ray Optical Systems, Inc., which he had cofounded coincident with UAlbany's Center for X-Ray Optics in 1991. He was the author of more than 300 technical articles and mentor to more than 48 doctoral graduates.

Walter Gibson's boundless energy, enthusiasm, wisdom, caring, courage, and vision inspired multiple generations of scientists.

GLOSSARY AND FUNDAMENTAL CONSTANTS

Introduction

This glossary of the terms used in the *Handbook* represents to a large extent the language of optics. The symbols are representations of numbers, variables, and concepts. Although the basic list was compiled by the author of this section, all the editors have contributed and agreed to this set of symbols and definitions. Every attempt has been made to use the same symbols for the same concepts throughout the entire *Handbook*, although there are exceptions. Some symbols seem to be used for many concepts. The symbol α is a prime example, as it is used for absorptivity, absorption coefficient, coefficient of linear thermal expansion, and more. Although we have tried to limit this kind of redundancy, we have also bowed deeply to custom.

Units

The abbreviations for the most common units are given first. They are consistent with most of the established lists of symbols, such as given by the International Standards Organization ISO¹ and the International Union of Pure and Applied Physics, IUPAP.²

Prefixes

Similarly, a list of the numerical prefixes¹ that are most frequently used is given, along with both the common names (where they exist) and the multiples of ten that they represent.

Fundamental Constants

The values of the fundamental constants³ are listed following the sections on SI units.

Symbols

The most commonly used symbols are then given. Most chapters of the *Handbook* also have a glossary of the terms and symbols specific to them for the convenience of the reader. In the following list, the symbol is given, its meaning is next, and the most customary unit of measure for the quantity is presented in brackets. A bracket with a dash in it indicates that the quantity is unitless. Note that there is a difference between units and dimensions. An angle has units of degrees or radians and a solid angle square degrees or steradians, but both are pure ratios and are dimensionless. The unit symbols as recommended in the SI system are used, but decimal multiples of some of the dimensions are sometimes given. The symbols chosen, with some cited exceptions, are also those of the first two references.

RATIONALE FOR SOME DISPUTED SYMBOLS

The choice of symbols is a personal decision, but commonality improves communication. This section explains why the editors have chosen the preferred symbols for the *Handbook*. We hope that this will encourage more agreement.

Fundamental Constants

It is encouraging that there is almost universal agreement for the symbols for the fundamental constants. We have taken one small exception by adding a subscript B to the k for Boltzmann's constant.

Mathematics

We have chosen i as the imaginary unit arbitrarily. IUPAP lists both i and j , while ISO does not report on these.

Spectral Variables

These include expressions for the wavelength λ , frequency ν , wave number σ , ω for circular or radian frequency, k for circular or radian wave number and dimensionless frequency x . Although some use f for frequency, it can be easily confused with electronic or spatial frequency. Some use $\tilde{\nu}$ for wave number, but, because of typography problems and agreement with ISO and IUPAP, we have chosen σ ; it should not be confused with the Stefan-Boltzmann constant. For spatial frequencies we have chosen ξ and η , although f_x and f_y are sometimes used. ISO and IUPAP do not report on these.

Radiometry

Radiometric terms are contentious. The most recent set of recommendations by ISO and IUPAP are L for radiance [$\text{Wcm}^{-2}\text{sr}^{-1}$], M for radiant emittance or exitance [Wcm^{-2}], E for irradiance or incidence [Wcm^{-2}], and I for intensity [Wsr^{-2}]. The previous terms, W , H , N , and J , respectively, are still in many texts, notably Smith⁴ and Lloyd⁵ but we have used the revised set, although there are still shortcomings. We have tried to deal with the vexatious term *intensity* by using *specific intensity* when the units are $\text{Wcm}^{-2}\text{sr}^{-1}$, *field intensity* when they are Wcm^{-2} , and *radiometric intensity* when they are Wsr^{-1} .

There are two sets of terms for these radiometric quantities, which arise in part from the terms for different types of reflection, transmission, absorption, and emission. It has been proposed that the *ion* ending indicate a process, that the *ance* ending indicate a value associated with a particular sample, and that the *ivity* ending indicate a generic value for a "pure" substance. Then one also has reflectance, transmittance, absorptance, and emittance as well as reflectivity, transmissivity, absorptivity, and emissivity. There are now two different uses of the word emissivity. Thus the words *exitance*, *incidence*, and *sterance* were coined to be used in place of emittance, irradiance, and radiance. It is interesting that ISO uses radiance, exitance, and irradiance whereas IUPAP uses radiance, exitance [*sic*], and irradiance. We have chosen to use them both, i.e., emittance, irradiance, and radiance will be followed in square brackets by exitance, incidence, and sterance (or vice versa). Individual authors will use the different endings for transmission, reflection, absorption, and emission as they see fit.

We are still troubled by the use of the symbol E for irradiance, as it is so close in meaning to electric field, but we have maintained that accepted use. The spectral concentrations of these quantities, indicated by a wavelength, wave number, or frequency subscript (e.g., L_λ) represent partial differentiations; a subscript q represents a photon quantity; and a subscript ν indicates a quantity normalized to the response of the eye. Thereby, L_ν is luminance, E_ν illuminance, and M_ν and I_ν luminous emittance and luminous intensity. The symbols we have chosen are consistent with ISO and IUPAP.

The refractive index may be considered a radiometric quantity. It is generally complex and is indicated by $\tilde{n} = n - ik$. The real part is the relative refractive index and k is the extinction coefficient. These are consistent with ISO and IUPAP, but they do not address the complex index or extinction coefficient.

Optical Design

For the most part ISO and IUPAP do not address the symbols that are important in this area.

There were at least 20 different ways to indicate focal ratio; we have chosen FN as symmetrical with NA; we chose f and efl to indicate the effective focal length. Object and image distance, although given many different symbols, were finally called s_o and s_i since s is an almost universal symbol for distance. Field angles are θ and ϕ ; angles that measure the slope of a ray to the optical axis are u ; u can also be $\sin u$. Wave aberrations are indicated by W_{ijk} , while third-order ray aberrations are indicated by σ_i and more mnemonic symbols.

Electromagnetic Fields

There is no argument about \mathbf{E} and \mathbf{H} for the electric and magnetic field strengths, Q for quantity of charge, ρ for volume charge density, σ for surface charge density, etc. There is no guidance from Refs. 1 and 2 on polarization indication. We chose \perp and \parallel rather than p and s , partly because s is sometimes also used to indicate scattered light.

There are several sets of symbols used for reflection transmission, and (sometimes) absorption, each with good logic. The versions of these quantities dealing with field amplitudes are usually specified with lower case symbols: r , t , and a . The versions dealing with power are alternately given by the uppercase symbols or the corresponding Greek symbols: R and T versus ρ and τ . We have chosen to use the Greek, mainly because these quantities are also closely associated with Kirchhoff's law that is usually stated symbolically as $\alpha = \epsilon$. The law of conservation of energy for light on a surface is also usually written as $\alpha + \rho + \tau = 1$.

Base SI Quantities

length	m	meter
time	s	second
mass	kg	kilogram
electric current	A	ampere
temperature	K	kelvin
amount of substance	mol	mole
luminous intensity	cd	candela

Derived SI Quantities

energy	J	joule
electric charge	C	coulomb
electric potential	V	volt
electric capacitance	F	farad
electric resistance	Ω	ohm
electric conductance	S	siemens
magnetic flux	Wb	weber
inductance	H	henry
pressure	Pa	pascal
magnetic flux density	T	tesla
frequency	Hz	hertz
power	W	watt
force	N	newton
angle	rad	radian
angle	sr	steradian

Prefixes

Symbol	Name	Common name	Exponent of ten
F	exa		18
P	peta		15
T	tera	trillion	12
G	giga	billion	9
M	mega	million	6
k	kilo	thousand	3
h	hecto	hundred	2
da	deca	ten	1
d	deci	tenth	-1
c	centi	hundredth	-2
m	milli	thousandth	-3
μ	micro	millionth	-6
n	nano	billionth	-9
p	pico	trillionth	-12
f	femto		-15
a	atto		-18

Constants

c	speed of light vacuo [299792458 ms ⁻¹]
c_1	first radiation constant = $2\pi^2 h = 3.7417749 \times 10^{-16}$ [Wm ²]
c_2	second radiation constant = $hc/k = 0.014838769$ [mK]
e	elementary charge [$1.60217733 \times 10^{-19}$ C]
g_n	free fall constant [9.80665 ms ⁻²]
h	Planck's constant [$6.6260755 \times 10^{-34}$ Ws]
k_B	Boltzmann constant [1.380658×10^{-23} JK ⁻¹]
m_e	mass of the electron [$9.1093897 \times 10^{-31}$ kg]
N_A	Avogadro constant [6.0221367×10^{23} mol ⁻¹]
R_∞	Rydberg constant [10973731.534 m ⁻¹]
ϵ_0	vacuum permittivity [$\mu_0^{-1}c^{-2}$]
σ	Stefan-Boltzmann constant [5.67051×10^{-8} Wm ⁻¹ K ⁻⁴]
μ_0	vacuum permeability [$4\pi \times 10^{-7}$ NA ⁻²]
μ_B	Bohr magneton [$9.2740154 \times 10^{-24}$ JT ⁻¹]

General

B	magnetic induction [Wbm ⁻² , kgs ⁻¹ C ⁻¹]
C	capacitance [f, C ² s ² m ⁻² kg ⁻¹]
C	curvature [m ⁻¹]
c	speed of light in vacuo [ms ⁻¹]
c_1	first radiation constant [Wm ²]
c_2	second radiation constant [mK]
D	electric displacement [Cm ⁻²]
E	incidence [irradiance] [Wm ⁻²]
e	electronic charge [coulomb]
E_v	illuminance [lux, lmm ⁻²]
E	electrical field strength [Vm ⁻¹]
E	transition energy [J]
E_g	band-gap energy [eV]
f^g	focal length [m]
f_f	Fermi occupation function, conduction band
f_v	Fermi occupation function, valence band

FN	focal ratio (<i>f</i> /number) [—]
<i>g</i>	gain per unit length [m^{-1}]
g_{th}	gain threshold per unit length [m^{-1}]
H	magnetic field strength [Am^{-1} , $\text{Cs}^{-1} \text{m}^{-1}$]
<i>h</i>	height [m]
<i>I</i>	irradiance (see also <i>E</i>) [Wm^{-2}]
<i>I</i>	radiant intensity [Wsr^{-1}]
<i>I</i>	nuclear spin quantum number [—]
<i>I</i>	current [A]
<i>i</i>	$\sqrt{-1}$
Im()	imaginary part of
<i>J</i>	current density [Am^{-2}]
j	total angular momentum [$\text{kg m}^2 \text{s}^{-1}$]
$J_1()$	Bessel function of the first kind [—]
<i>k</i>	radian wave number $=2\pi/\lambda$ [rad cm^{-1}]
k	wave vector [rad cm^{-1}]
<i>k</i>	extinction coefficient [—]
<i>L</i>	sterance [radiance] [$\text{Wm}^{-2} \text{sr}^{-1}$]
L_v	luminance [cdm^{-2}]
<i>L</i>	inductance [h, $\text{m}^2 \text{kg C}^2$]
<i>L</i>	laser cavity length
<i>L, M, N</i>	direction cosines [—]
<i>M</i>	angular magnification [—]
<i>M</i>	radiant exitance [radiant emittance] [Wm^{-2}]
<i>m</i>	linear magnification [—]
<i>m</i>	effective mass [kg]
MTF	modulation transfer function [—]
<i>N</i>	photon flux [s^{-1}]
<i>N</i>	carrier (number) density [m^{-3}]
<i>n</i>	real part of the relative refractive index [—]
\tilde{n}	complex index of refraction [—]
NA	numerical aperture [—]
OPD	optical path difference [m]
<i>P</i>	macroscopic polarization [C m^{-2}]
Re()	real part of [—]
<i>R</i>	resistance [Ω]
r	position vector [m]
<i>S</i>	Seebeck coefficient [VK^{-1}]
<i>s</i>	spin quantum number [—]
<i>s</i>	path length [m]
S_o	object distance [m]
S_i	image distance [m]
T	temperature [K, C]
<i>t</i>	time [s]
<i>t</i>	thickness [m]
<i>u</i>	slope of ray with the optical axis [rad]
<i>V</i>	Abbe reciprocal dispersion [—]
<i>V</i>	voltage [V , $\text{m}^2 \text{kg s}^{-2} \text{C}^{-1}$]
<i>x, y, z</i>	rectangular coordinates [m]
<i>Z</i>	atomic number [—]

Greek Symbols

α	absorption coefficient [cm^{-1}]
α	(power) absorptance (absorptivity)

ϵ	dielectric coefficient (constant) [—]
ϵ	emittance (emissivity) [—]
ϵ	eccentricity [—]
ϵ_1	Re (ϵ)
ϵ_2	Im (ϵ)
τ	(power) transmittance (transmissivity) [—]
ν	radiation frequency [Hz]
ω	circular frequency = $2\pi\nu$ [rads ⁻¹]
ω	plasma frequency [Hz]
λ	wavelength [μm , nm]
σ	wave number = $1/\lambda$ [cm ⁻¹]
σ	Stefan Boltzmann constant [Wm ⁻² K ⁻¹]
ρ	reflectance (reflectivity) [—]
θ, ϕ	angular coordinates [rad, °]
ξ, η	rectangular spatial frequencies [m ⁻¹ , r ⁻¹]
ϕ	phase [rad, °]
ϕ	lens power [m ⁻²]
Φ	flux [W]
χ	electric susceptibility tensor [—]
Ω	solid angle [sr]

Other

\Re	responsivity
$\exp(x)$	e^x
$\log_a(x)$	log to the base a of x
$\ln(x)$	natural log of x
$\log(x)$	standard log of x : $\log_{10}(x)$
Σ	summation
Π	product
Δ	finite difference
δx	variation in x
dx	total differential
∂x	partial derivative of x
$\delta(x)$	Dirac delta function of x
δ_{ij}	Kronecker delta

REFERENCES

1. Anonymous, *ISO Standards Handbook 2: Units of Measurement*, 2nd ed., International Organization for Standardization, 1982.
2. Anonymous, *Symbols, Units and Nomenclature in Physics*, Document U.I.P. 20, International Union of Pure and Applied Physics, 1978.
3. E. Cohen and B. Taylor, "The Fundamental Physical Constants," *Physics Today*, 9 August 1990.
4. W. J. Smith, *Modern Optical Engineering*, 2nd ed., McGraw-Hill, 1990.
5. J. M. Lloyd, *Thermal Imaging Systems*, Plenum Press, 1972.

William L. Wolfe
 College of Optical Sciences
 University of Arizona
 Tucson, Arizona

PART

1

MEASUREMENTS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

SCATTEROMETERS

John C. Stover

*The Scatter Works, Inc.
Tucson, Arizona*

1.1 GLOSSARY

BRDF	bidirectional reflectance distribution function
BTDF	bidirectional transmittance distribution function
BSDF	bidirectional scatter distribution function
f	focal length
L	distance
P	power
R	length
r	radius
TIS	total integrated scatter
θ	angle
θ_N	vignetting angle
θ_{spec}	specular angle
λ	wavelength
σ	rms roughness
Ω	solid angle

1.2 INTRODUCTION

In addition to being a serious source of noise, scatter reduces throughput, limits resolution, and has been the unexpected source of practical difficulties in many optical systems. On the other hand, its measurement has proved to be an extremely sensitive method of providing metrology information for components used in many diverse applications. Measured scatter is a good indicator of surface quality and can be used to characterize surface roughness as well as locate and size

discrete defects. It is also used to measure the quality of optical coatings and bulk optical materials. This chapter reviews basic issues associated with scatter metrology and touches on various industrial applications.

The pioneering scattering instrumentation^{1–32} work started in the 1960s and extended into the 1990s. This early work (reviewed in 1995)¹¹ resulted in commercially available lab scatterometers and written standards in SEMI and ASTM detailing measurement, calibration and reporting.^{33–36} Understanding the measurements and the ability to repeat results and communicate them led to an expansion of industrial applications, scatterometry has become an increasingly valuable source of noncontact metrology in industries where surface inspection is important. For example, each month millions of silicon wafers (destined to be processed into computer chips) are inspected for point defects (pits and particles) with “particle scanners,” which are essentially just scatterometers. These rather amazing instruments (now costing more than \$1 million each) map wafer defects smaller than 50 nm and can distinguish between pits and particles. In recent years their manufacture has matured to the point where system specifications and calibration are now also standardized in SEMI.^{37–40} Scatter metrology is also found in industries as diverse as medicine, sheet metal production and even the measurement of appearance—where it has been noted that while beauty is in the eye of the beholder, what we see is scattered light. The polarization state of scatter signals has also been exploited^{25–28, 41–44} and is providing additional product information. Many more transitions from lab scatterometer to industry application are expected. They depend on understanding the basic measurement concepts outlined in this chapter.

Although it sounds simple, the instrumentation required for these scatter measurements is fairly sophisticated. Scatter signals are generally small compared to the specular beam and can vary by several orders of magnitude in just a few degrees. Complete characterization may require measurement over a large fraction of the sphere surrounding the scatter source. For many applications, a huge array of measurement decisions (incident angle, wavelength, source and receiver polarization, scan angles, etc.) faces the experimenter. The instrument may faithfully record a signal, but is it from the sample alone? Or, does it also include light from the instrument, the wall behind the instrument, and even the experimenter’s shirt? These are not easy questions to answer at nanowatt levels in the visible and get even harder in the infrared and ultraviolet. It is easy to generate scatter data—lots of it. Obtaining accurate values of appropriate measurements and communicating them requires knowledge of the instrumentation as well as insight into the problem being addressed.

In 1961, Bennett and Porteus¹ reported measurement of signals obtained by integrating scatter over the reflective hemisphere. They defined a parameter called the *total integrated scatter* (TIS) as the integrated reflected scatter normalized by the total reflected light. Using a scalar diffraction theory result drawn from the radar literature,² they related the TIS to the reflector root mean square (rms) roughness. By the mid-1970s, several scatterometers had been built at various university, government, and industry labs that were capable of measuring scatter as a function of angle; however, instrument operation and data manipulation were not always well automated.^{3–6} Scattered power per unit solid angle (sometimes normalized by the incident power) was usually measured. Analysis of scatter data to characterize sample surface roughness was the subject of many publications.^{7–11} Measurement comparison between laboratories was hampered by instrument differences, sample contamination, and confusion over what parameters should be compared. A derivation of what is commonly called BRDF (for bidirectional reflectance distribution function) was published by Nicodemus and coworkers in 1970, but did not gain common acceptance as a way to quantify scatter measurements until after publication of their 1977 NBS monograph.¹² With the advent of small powerful computers in the 1980s, instrumentation became more automated. Increased awareness of scatter problems and the sensitivity of many end-item instruments increased government funding for better instrumentation.^{13–14} As a result, instrumentation became available that could measure and analyze as many as 50 to 100 samples a day instead of just a handful. Scatterometers became commercially available and the number (and sophistication) of measurement facilities increased.^{15–17} Further instrumentation improvements will include more out-of-plane capability, extended wavelength control, and polarization control at both source and receiver. As of 2008 there are written standards for BRDF and TIS in ASTM and SEMI.^{33–36}

This review gives basic definitions, instrument configurations, components, scatter specifications, measurement techniques, and briefly discusses calibration and error analysis.

1.3 DEFINITIONS AND SPECIFICATIONS

One of the difficulties encountered in comparing measurements made on early instruments was getting participants to calculate the same quantities. There were problems of this nature as late as 1988 in a measurement round-robin run at 633 nm.²⁰ But, there are other reasons for reviewing these basic definitions before discussing instrumentation. The ability to write useful scatter specifications (i.e., the ability to make use of quantified scatter information) depends just as much on understanding the defined quantity as it does on understanding the instrumentation and the specific scatter problem. In addition, definitions are often given in terms of mathematical abstractions that can only be approximated in the lab. This is the case for BRDF.

$$\text{BRDF} = \frac{\text{differential radiance}}{\text{differential irradiance}} \approx \frac{dP_s/d\Omega}{P_i \cos \theta_s} \approx \frac{P_s/\Omega}{P_i \cos \theta_s} \quad (1)$$

BRDF has been strictly defined as the ratio of the sample differential radiance to the differential irradiance under the assumptions of a collimated beam with uniform cross section incident on an isotropic surface reflector (no bulk scatter allowed). Under these conditions, the third quantity in Eq. (1) is found, where power P in watts instead of intensity I in W/m^2 has been used. The geometry is shown in Fig. 1. The value θ_s is the polar angle in the scatter direction measured from reflector normal and Ω is the differential solid angle (in steradians) through which dP_s (watts) scatters when P_i (watts) is incident on the reflector. The cosine comes from the definition of radiance and may be viewed as a correction from the actual size of the scatter source to the apparent size (or projected area) as the viewer rotates away from surface normal.

The details of the derivation do not impact scatter instrumentation, but the initial assumptions and the form of the result do. When scattered light power is measured, it is through a finite diameter

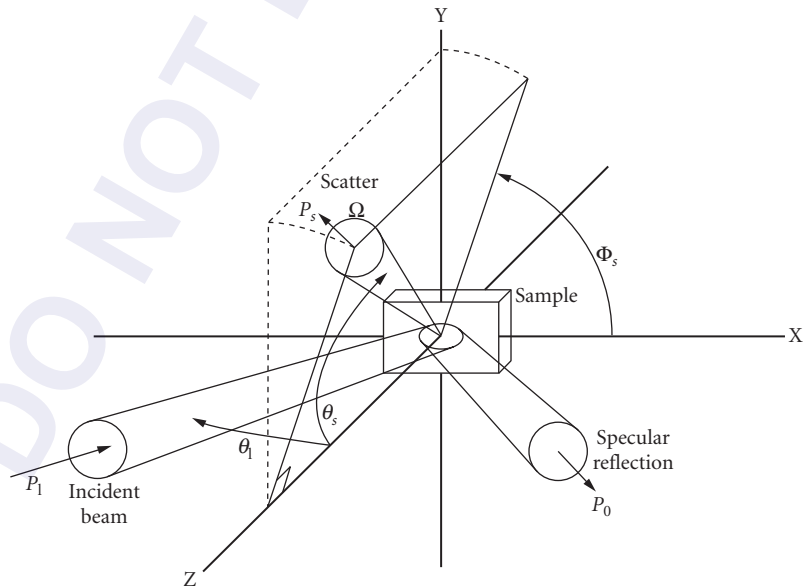


FIGURE 1 Geometry for the definition of BRDF.

aperture; as a result the calculation is for an average BRDF over the aperture. This is expressed in the final term of Eq. (1), where P_s is the measured power through the finite solid angle Ω defined by the receiver aperture and the distance to the scatter source. Thus, when the receiver aperture is swept through the scatter field to obtain angle dependence, the measured quantity is actually the convolution of the aperture over the differential BRDF. This does not cause serious distortion unless the scatter field has abrupt intensity changes, as it does near specular or near diffraction peaks associated with periodic surface structure. But there are even more serious problems between the strict definition of BRDF (as derived by Nicodemus) and practical measurements. There are no such things as uniform cross-section beams and isotropic samples that scatter only from surface structure. So, the third term of Eq. (1) is not exactly the differential radiance/irradiance ratio for the situations we create in the lab with our instruments. However, it makes perfect sense to measure normalized scattered power density as a function of direction [as defined in the fourth term of Eq. (1)] even though it cannot be exactly expressed in convenient radiometric terms.

A slightly less cumbersome definition (in terms of writing scatter specifications) is realized if the cosine term is dropped. This is referred to as “the cosine-corrected BRDF,” or sometimes, “the scatter function.” Its use has caused some of the confusion surrounding measurement differences found in scatter round robins. In accordance with the original definition, accepted practice, and BRDF Standards,^{33,36} the BRDF contains the cosine, as given in Eq. (1), and the cosine-corrected BRDF does not. It also makes sense to extend the definition to volume scatter sources and even make measurements on the transmissive side of the sample. The term BTDF (for bidirectional transmission distribution function) is used for transmissive scatter, and BSDF (bidirectional scatter distribution function) is all-inclusive.

The BSDF has units of inverse steradians and, unlike reflectance and transmission (which vary from 0.0 to 1.0), can take on very large values as well as very small values.^{1,21} For near-normal incidence, a measurement made at the specular beam of a high reflectance mirror results in a BSDF value of approximately $1/\Omega$, which is generally a large number. Measured values at the specular direction on the order of 10^6 sr^{-1} are common for a HeNe laser source. For low-scatter measurements, large apertures are generally used and values fall to the noise equivalent BSDF (or NEBSDF). This level depends on incident power and polar angle (position) as well as aperture size and detector noise, and typically varies from 10^{-4} sr^{-1} to 10^{-10} sr^{-1} . Thus, the measured BSDF can easily vary by over a dozen orders of magnitude in a given angle scan. This large variation results in challenges in instrumentation design as well as data storage, analysis, and presentation, and is another reason for problems with comparison measurements.

Instrument signature is the measured background scatter signal caused by the instrument and not the sample. It is caused by a combination of scatter created within the instrument and by the NEBSDF. Any instrument scatter that reaches the receiver field of view (FOV) will contribute to it. Common causes are scatter from source optics and the system beam dump. It is typically measured without a sample in place; however, careful attention has to be paid to the receiver FOV to ascertain that this is representative of the sample measurement situation. It is calculated as though the signal came from the sample (i.e., the receiver/sample solid angle is used) so that it can be compared to the measured sample BSDF. Near specular, the signal can be dominated by scatter (or diffraction) contributions from the source, called instrument signature. At high scatter angles it can generally be limited to NEBSDF levels. Sample measurements are always a combination of desired signal and instrument signature. Reduction of instrument signature, especially near specular, is a prime consideration in instrument design and use.

BSDF specifications always require inclusion of incident angle, source wavelength, and polarization as well as observation angles, scatter levels, and sample orientation. Depending on the sample and the measurement, they may also require aperture information to account for convolution effects. Specifications for scatter instrumentation should include instrument signature limits and the required NEBSDF. Specifications for the NEBSDF must include the polar angle, the solid angle, and the incident power to be meaningful.

TIS measurements are made by integrating the BSDF over a majority of either the reflective or transmissive hemispheres surrounding the scatter source. This is usually done with instrumentation that gathers (integrates) the scattered light signal. The TIS can sometimes be calculated from BSDF

data. If an isotropic sample is illuminated at near-normal incidence with circularly polarized light, data from a single measurement scan is enough to calculate a reasonably accurate TIS value for an entire hemisphere of scatter. The term “total integrated scatter” is a misnomer in that the integration is never actually “total,” as some scatter is never measured. Integration is commonly performed from a few degrees from specular to polar angles approaching 90° (approaching 1° to more than 45° in the SEMI Standard).³⁴ Measurements can be made of either transmissive or reflective scatter. TIS is calculated by ratioing the integrated scatter to the reflected (or transmitted) power as shown below in Eq. (2). For optically smooth components, the scatter signal is small compared to the specular beam and is often ignored. For reflective scatter the conversion to rms roughness (σ) under the assumption of an optically smooth, clean, reflective surface, via Davies’ scalar theory,² is also given. This latter calculation does not require gaussian surface statistics (as originally assumed by Davies) or even surface isotropy, but will work for other distributions, including gratings and machined optics.^{8,11} There are other issues (polarization and the assumption of mostly near specular scatter) that cause some error in this conversion. Comparison of TIS-generated roughness to profile-generated values is made difficult by a number of issues (bandwidth limits, one-dimensional profiling of a two-dimensional surface, etc.) that are beyond the scope of this section (see Ref. 11 for a complete discussion). One additional caution is that the literature and more than one stray radiation analysis program define TIS as scattered power normalized by incident power (which is essentially diffuse reflectance). This is a seriously incorrect distortion of TIS. Such a definition obviously cannot be related to surface roughness, as a change in reflectance (but not roughness) will change the ratio.

$$\text{TIS} = \frac{\text{integrated scattered power}}{\text{total reflected power}} \cong \left(\frac{4\pi\sigma}{\lambda} \right)^2 \quad (2)$$

TIS is one of three ratios that may be formed from the incident power, the specular reflected (or transmitted) power, and the integrated scatter. The other two ratios are the diffuse reflectance (or transmittance) and the specular reflectance (or transmittance). Typically, all three ratios may be obtained from measurements taken in TIS or BSDF instruments. Calculation, or specification, of any of these quantities that involve integration of scatter, also requires that the integration limits be given, as well as the wavelength, angle of incidence, source polarization, and sample orientation.

1.4 INSTRUMENT CONFIGURATIONS AND COMPONENT DESCRIPTIONS

The scatterometer shown in Fig. 2 is representative of the most common instrument configuration in use. The source is fixed in position. The sample is rotated to the desired incident angle, and the receiver is rotated about the sample in the plane of incidence. Although dozens of instruments have been built following this general design, other configurations are in use. For example, the source and receiver may be fixed and the sample rotated so that the scatter pattern moves past the receiver. This is easier mechanically than moving the receiver at the end of an arm, but complicates analysis because the incident angle and the observation angle change simultaneously. Another combination is to fix the source and sample together, at constant incident angle, and rotate this unit (about the point of illumination on the sample) so that the scatter pattern moves past a fixed receiver. This has the advantage that a long receiver/sample distance can be used without motorizing a long (heavy) receiver arm. It has the disadvantage that heavy (or multiple) sources are difficult to deal with. Other configurations, with everything fixed, have been designed that employ several receivers to merely sample the BSDF and display a curve fit of the resulting data. This is an economical solution if the BSDF is relatively uniform without isolated diffraction peaks. Variations on this last combination are common in industry where the samples are moved through the beam (sometimes during the manufacturing process) and the scatter measured in one or more directions.

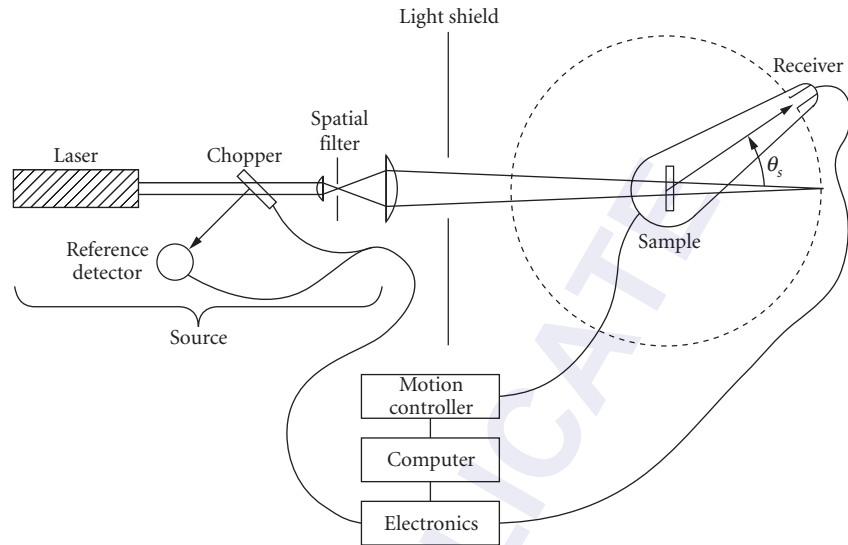


FIGURE 2 Components of a typical BSDF scatterometer.

Computer control of the measurement is essential to maximize versatility and minimize measurement time. The software required to control the measurement plus display and analyze the data can be expected to be a significant portion of total instrument development cost. The following paragraphs review typical design features (and issues) associated with the source, sample mount, and receiver components.

The source in Fig. 2 is formed by a laser beam that is chopped, spatially filtered, expanded, and finally brought to a focus on the receiver path. The beam is chopped to reduce both optical and electronic noise. This is accomplished through the use of lock-in detection in the electronics package which suppresses all signals except those at the chopping frequency. Low-noise, programmable gain electronics are essential to reducing NEBSDF. The reference detector is used to allow the computer to ratio out laser power fluctuations and, in some cases, to provide the necessary timing signal to the lock-in electronics. Polarizers, wave plates, and neutral density filters are also commonly placed prior to the spatial filter when required in the source optics. The spatial filter removes scatter from the laser beam and presents a point source which is imaged by the final focusing element to the detector zero position. Although a lens is shown in Fig. 2, the use of a mirror, which works over a larger range of wavelengths and generally scatters less light, is more common. For most systems the relatively large focal length of the final focusing element allows use of a spherical mirror and causes only minor aberration. Low-scatter spherical mirrors are easier to obtain than other conic sections. The incident beam is typically focused at the receiver to facilitate near specular measurement. Another option (a collimated beam at the receiver) is sometimes used and will be considered in the discussion on receivers. In either case, curved samples can be accommodated by adjusting the position of the spatial filter with respect to the final focusing optic. The spot size on the sample is obviously determined by elements of the system geometry and can be adjusted by changing the focal length of the first lens (often a microscope objective). The source region is completed by a shield that isolates stray laser light from the detector.

Lasers are convenient sources, but are not necessary. Broadband sources are often required to meet a particular application or to simulate the environment where a sample will be used. Monochromators and filters can be used to provide scatterometer sources of arbitrary wavelength.²⁰ Noise floor with these tunable incoherent sources increases dramatically as the spectral bandpass is narrowed, but they have the advantage that the scatter pattern does not contain laser speckle.

The sample mount can be very simple or very complex. In principle, 6° of mechanical freedom are required to fully adjust the sample. Three translational degrees of freedom allow the sample area (or volume) of interest to be positioned at the detector rotation axis and illuminated by the source. Three rotational degrees of freedom allow the sample to be adjusted for angle of incidence, out-of-plane tilt, and rotation about sample normal. The order in which these stages are mounted affects the ease of use (and cost) of the sample holder. In practice, it often proves convenient to either eliminate, or occasionally duplicate, some of these degrees of freedom. Exact requirements for these stages differ depending on whether the sample is reflective or transmissive, as well as with size and shape. In addition, some of these axes may be motorized to allow the sample area to be raster-scanned to automate sample alignment or to measure reference samples. The order in which these stages are mounted affects the ease of sample alignment. As a general rule, the scatter pattern is insensitive to small changes in incident angle but very sensitive to small angular deviations from specular. Instrumentation should be configured to allow location of the specular reflection (or transmission) very accurately.

The receiver rotation stage should be motorized and under computer control so that the input aperture may be placed at any position on the observation circle (dotted line in Fig. 2). Data scans may be initiated at any location. Systems vary as to whether data points are taken with the receiver stopped or “on the fly.” The measurement software is less complicated if the receiver is stopped. Unlike many TIS systems, the detector is always approximately normal to the incoming scatter signal. In addition to the indicated axis of rotation, some mechanical freedom is required to ensure that the receiver is at the correct height and pointed (tilted) at the illuminated sample. Sensitivity, low noise, linearity, and dynamic range are the important issues in choosing a detector element and designing the receiver housing. In general, these requirements are better met with photovoltaic detectors than photoconductive detectors. Small area detectors reduce the NEBSDF.

Receiver designs vary, but changeable apertures, bandpass filters, polarizers, lenses, and field stops are often positioned in front of the detector element. Figure 3 shows two receiver configurations, one designed for use with a converging source and one with a collimated source. In Fig. 3a, the illuminated sample spot is imaged on a field stop in front of the detector. This configuration is commonly

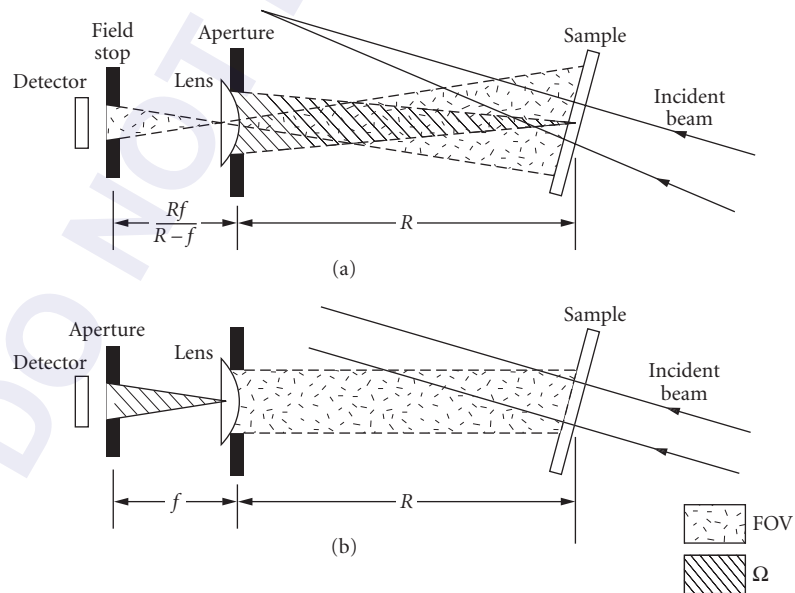


FIGURE 3 Receiver configurations: (a) converging source and (b) collimated source.

used with the source light converging on the receiver path. The field stop determines the receiver FOV. The aperture at the front of the receiver determines the solid angle over which scatter is gathered. Any light entering this aperture that originates from within the FOV will reach the detector and become part of the signal. This includes instrument signature contributions scattered through small angles by the source optics. It will also include light scattered by the receiver lens so that it appears to come from the sample. The configuration in Fig. 3a can be used to obtain near specular measurements by bringing a small receiver aperture close to the focused specular beam. With this configuration, reducing the front aperture does not limit the FOV. The receiver in Fig. 3b is in better accordance with the strict definition of BRDF in that a collimated source can be used. An aperture is located one focal length behind a collecting lens (or mirror) in front of the detector. The intent is to measure bundles of nearly parallel rays scattered from the sample. The angular spread of rays allowed to pass to the detector defines the receiver solid angle, which is equal to the aperture size divided by the focal length of the lens. This ratio (not the front aperture/sample distance) determines the solid angle of this receiver configuration. The FOV is determined by the clear aperture of the lens, which must be kept larger than the illuminated spot on the sample. The Fig. 3b design is unsuitable for near specular measurement because the relatively broad collimated specular beam will scatter from the receiver lens for several degrees from specular. It is also limited in measuring large incident angle situations where the elongated spot may exceed the FOV. If the detector (and its stop) can be moved in relation to the lens, receivers can be adjusted from one configuration to the other. Away from the specular beam, in low instrument signature regions, there is no difference in the measured BSDF values between the two systems. Commercially available research scatterometers are available that measure both in and out of the incident plane and from the mid-IR to the near UV.

The two common methods of approaching TIS measurements are shown in Fig. 4. The first one, employed by Bennett and Porteus in their early instrument,¹ uses a hemispherical mirror (or Coblentz sphere) to gather scattered light from the sample and image it onto a detector. The specular beam enters and leaves the hemisphere through a small circular hole. The diameter of that hole defines the near specular limit of the instrument. The reflected beam (not the incident beam) should be centered in the hole because the BSDF will be symmetrical about it. Alignment of the hemispherical mirror is critical in this approach. The second approach involves the use of an integrating sphere. A section of the sphere is viewed by a recessed detector. If the detector FOV is limited to a section of the sphere

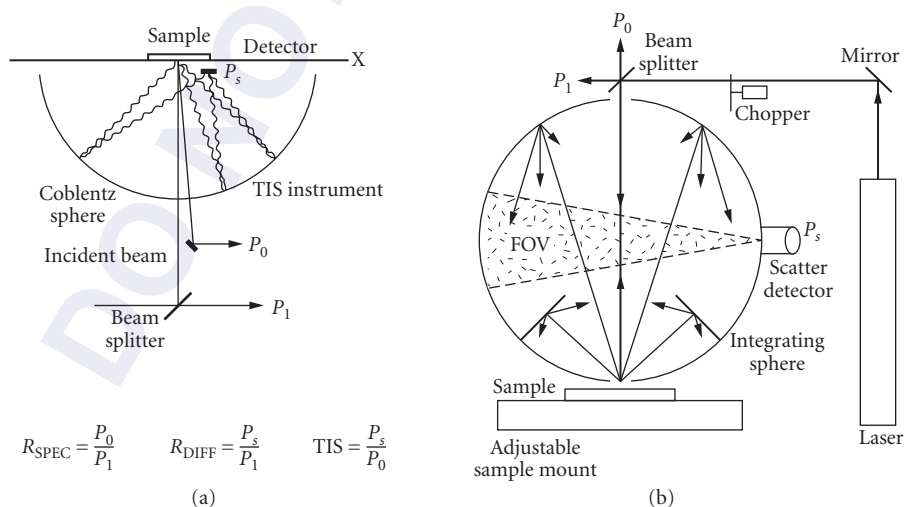


FIGURE 4 TIS measurement with a (a) Coblentz sphere and (b) diffuse integrating sphere.

that is not directly illuminated by scatter from the sample, then the signal will be proportional to total scatter from the sample. Again, the reflected beam should be centered on the exit hole. The Coblenz Sphere method presents more signal to the detector; however, some of this signal is incident on the detector at very high angles. Thus, this approach tends to discriminate against high-angle scatter (which is often much smaller for many samples). The integrating sphere is easier to align, but has a lower signal-to-noise ratio (less signal on the detector) and is more difficult to build in the IR where uniform diffuse surfaces are harder to obtain. Even so, sophisticated integrating sphere systems have become commercially available that can measure down to 0.5 angstroms rms roughness. A common mistake with TIS measurements is to assume that for near-normal incidence, the orientation between source polarization and sample orientation is not an issue. TIS measurements made with a linearly polarized source on a grating at different orientations will quickly demonstrate this dependence.

TIS measurements can be made over very near specular ranges by utilizing a diffusely reflecting plate with a small hole in it. A converging beam is reflected off the sample and through the hole. Scatter is diffusely reflected from the plate to a receiver designed to uniformly view the plate. The reflected power is measured by moving the plate so the specular beam misses the hole. Measurements starting closer than 0.1° from specular can be made in this manner, and it is an excellent way to check incoming optics or freshly coated optics for low scatter.

1.5 INSTRUMENTATION ISSUES

Measurement of near specular scatter is often one of the hardest requirements to meet when designing an instrument and has been addressed in several publications.²¹⁻²³ The measured BSDF may be divided into two regions relative to the specular beam, as shown in Fig. 5. Outside the angle θ_N from specular is a low-signature region where the source optics are not in the receiver FOV. Inside θ_N , at least some of the source optics scatter directly into the receiver and the signature increases rapidly until the receiver aperture reaches the edge of the specular beam. As the aperture moves

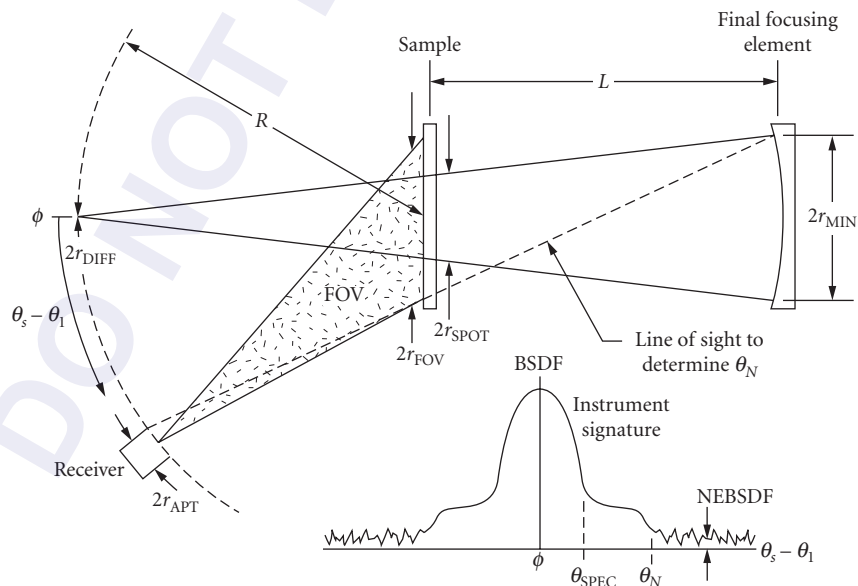


FIGURE 5 Near specular geometry and instrument signature.

closer to specular center, the measurement is dominated by the aperture convolution of the specular beam, and there is no opportunity to measure scatter. The value θ_N is easily calculated (via a small-angle approximation) using the instrument geometry and parameters identified in Fig. 5, where the receiver is shown at the θ_N position. The parameter F is the focal length of the sample.

$$\theta_N = (r_{\text{MIR}} + r_{\text{FOV}})/L + (r_{\text{FOV}} + r_{\text{apt}})/R - r_{\text{spot}}/F \quad (3)$$

It is easy to achieve values of θ_N below 10° and values as small as 1° can be realized with careful design. The offset angle from specular, θ_{spec} , at which the measurement is dominated by the specular beam, can be reduced to less than a tenth of a degree at visible wavelengths and is given by

$$\theta_{\text{spec}} = \frac{r_{\text{diff}} + r_{\text{apt}}}{R} \approx \frac{3\lambda}{D} + \frac{r_{\text{apt}}}{R} \quad (4)$$

Here, r_{diff} and r_{apt} are the radius of the focused spot and the receiver aperture, respectively (see Fig. 5 again). The value of r_{diff} can be estimated in terms of the diameter D of the focusing optic and its distance to the focused spot, $R + L$ (estimated as $2.5R$). The diffraction limit has been doubled in this estimate to allow for aberrations.

To take near specular measurements, both angles and the instrument signature need to be reduced. The natural reaction is to “increase R to increase angular resolution.” Although a lot of money has been spent doing this, it is an unnecessarily expensive approach. Angular resolution is achieved by reducing r_{apt} and by taking small steps. The radius r_{apt} can be made almost arbitrarily small so the economical way to reduce the r_{apt}/R terms is by minimizing r_{apt} —not by increasing R . A little thought about r_{FOV} and r_{diff} reveals that they are both proportional to R , so nothing is gained in the near specular game by purchasing large-radius rotary stages.

The reason for building a large-radius scatterometer is to accommodate a large FOV. This is often driven by the need to take measurements at large incident angles or by the use of broadband sources, both of which create larger spots on the sample. When viewing normal to the sample, the FOV requirements can be stringent. Because the maximum FOV is proportional to detector diameter (and limited at some point by minimum receiver lens FN), increasing R is the only open-ended design parameter available. It should be sized to accommodate the smallest detector likely to be used in the system. This will probably be in the mid-IR where uniform high-detectivity photovoltaic detectors are more difficult to obtain. On the other hand, a larger detector diameter means increased electronic noise and a larger NEBSDF.

Scatter sources of instrument signature can be reduced by these techniques.

1. Use the lowest-scatter focusing element in the source that you can afford and learn how to keep it clean. This will probably be a spherical mirror.
2. Keep the source area as “black” as possible. This especially includes the sample side of the spatial filter pinhole which is conjugate with the receiver aperture. Use a black pinhole.
3. Employ a specular beam dump that rides with your receiver and additional beam dumps to capture sample reflected and transmitted beams when the receiver has left the near specular area. Use your instrument to measure the effectiveness of your beam dumps.²⁶
4. Near specular scatter caused by dust in the air can be significantly reduced through the use of a filtered air supply over the specular beam path. A filtered air supply is essential for measuring optically smooth surfaces.

Away from specular, reduction of NEBSDF is the major factor in measuring low-scatter samples and increasing instrument quality. Measurement of visible scatter from a clean semiconductor wafer will take many instruments right down to instrument signature levels. Measurement of cross-polarized scatter requires a low NEBSDF for even high-scatter optics. For a given receiver solid angle, incident power, and scatter direction, the NEBSDF is limited by the noise equivalent power of the receiver (and associated electronics), once optical noise contributions are eliminated. The electronic contributions to NEBSDF are easily measured by simply covering the receiver aperture during a measurement.

TABLE 1 Comparison of Characteristics for Detectors Used at Different Wavelengths

Detector (2 mm dia.)	NEP (W/Hz)	Wavelength (nm)	P_i (W)	NEBSDF (sr ⁻¹)
PMT	10 ⁻¹⁵	633	0.005	10 ⁻¹⁰
Si	10 ⁻¹³	633	0.005	10 ⁻⁸
Ge	3 × 10 ⁻¹³	1,320	0.001	10 ⁻⁷
InSb	10 ^{-12*}	3,390	0.002	5 × 10 ⁻⁷
HgMgTe	10 ^{-11*}	10,600	2.0	10 ⁻⁸
Pyro	10 ⁻⁸	10,600	2.0	10 ⁻⁵

*Detector at 77 K.

Because the resulting signal varies in a random manner, NEBSDF should be expressed as an rms value (roughly equal to one-third of the peak level). An absolute minimum measurable scatter signal (in watts) can be found from the product of three terms: the required signal-to-noise ratio, the system noise equivalent power (or NEP given in watts per square root hertz), and the square root of the noise bandwidth (BW_n). The system NEP is often larger than the detector NEP and cannot be reduced below it. The detector NEP is a function of wavelength and increases with detector diameter. Typical detector NEP values (2-mm diameter) and wavelength ranges are shown as follows for several common detectors in Table 1. Notice that NEP tends to increase with wavelength. The noise bandwidth varies as the reciprocal of the sum of the system electronics time constant and the measurement integration time. Values of 0.1 to 10 Hz are commonly achieved. In addition to system NEP, the NEBSDF may be increased by contributions from stray source light, room lights, and noise in the reference signal. Table 1 also shows achievable rms NEBSDF values that can be realized at unity cosine, a receiver solid angle of 0.003 sr, 1-second integration, and the indicated incident powers. This column can be used as a rule of thumb in system design or to evaluate existing equipment. Simply adjust by the appropriate incident power, solid angle, and so on, to make the comparison. Adjusted values substantially higher than these indicate there is room for system improvement (don't worry about differences as small as a factor of 2). Further reduction of the instrument signature under these geometry and power conditions will require dramatically increased integration time (because of the square root dependence on noise bandwidth) and special attention to electronic dc offsets. Because the NEP tends to increase with wavelength, higher powers are needed in the mid-IR to reach the same NEBSDFs that can be realized in the visible. Because scatter from many sources tends to decrease at longer wavelengths, a knowledge of the instrument NEBSDF is especially critical in the mid-IR.

As a final configuration comment, the software package (both measurement and analysis) is crucial for an instrument that is going to be used for any length of time. Poor software will quickly cost work-years of effort due to errors, increased measurement and analysis time, and lost business. Expect to expend 1 to 2 man-years with experienced programmers writing a good package—it is worth it.

1.6 MEASUREMENT ISSUES

Sample measurement should be preceded (and sometimes followed) by a measurement of the instrument signature. This is generally accomplished by removing the sample and measuring the apparent BSDF from the sample as a transmissive scan. This is not an exact measure of instrument noise during sample measurement, but if the resulting BSDF is multiplied by sample reflectance (or transmission) before comparison to sample data, it can define some hard limits over which the sample data cannot be trusted. The signature should also be compared to the NEBSDF value obtained with the receiver aperture blocked. Obtaining the instrument signature also presents an opportunity to measure the incident power, which is required for calculation of the BSDF. The ability to see the data displayed as it is taken is an extremely helpful feature when it comes to reducing instrument signature and eliminating measurement setup errors.

Angle scans, which have dominated the preceding discussion, are an obvious way to take measurements. BSDF is also a function of position on the sample, source wavelength, and source polarization, and scans can also be taken at fixed angle (receiver position) as a function of these variables. Obviously, a huge amount of data is required to completely characterize scatter from a sample.

Raster scans are taken to measure sample uniformity or locate (map) discrete defects. A common method is to fix the receiver position and move the sample in its own x - y plane, recording the BSDF at each location. Faster approaches involve using multiple detectors (e.g., array cameras) with large area illumination, and scanning the source over the sample. Results can be presented using color maps or 3D isometric plots. Results can be further analyzed via histograms and various image-processing techniques.

There are three obvious choices for making wavelength scans. Filters (variable or discrete) can be employed at the source or receiver. A monochromator can be used as a source.²⁰ Finally, there is some advantage to using a Fourier transforming infrared spectrometer (FTIR) as a source in the mid-IR.²⁰ Details of these techniques are beyond the scope of this discussion; however, a couple of generalities will be mentioned. Even though these measurements often involve relatively large bandwidths at a given wavelength (compared to a laser), the NEBSDF is often larger by a few orders because of the smaller incident power. Further, because the bandwidths change differently between the various source types given above, meaningful measurement comparisons between instruments are often difficult to make.

Polarization scans are often limited to SS, SP, PS, and PP (source/receiver) combinations. However, complete polarization dependence of the sample requires the measurement of the sample Mueller matrix. This is found by creating a set of Stokes vectors at the source and measuring the resulting Stokes vector in the desired scatter direction.^{10,11,25-28} This is an area of instrumentation development that is the subject of increasing attention.⁴¹⁻⁴⁴

Speckle effects in the BSDF from a laser source can be eliminated in several ways. If a large receiver solid angle is used (generally several hundred speckles in size) there is not a problem. The sample can be rotated about its normal so that speckle is time averaged out of the measurement. This is still a problem when measuring very near the specular beam because sample rotation unavoidably moves the beam slightly during the measurement. In this case, the sample can be measured several times at slightly different orientations and the results averaged to form one speckle-free BSDF.

Scatter measurement in the retrodirection (back into the incident beam) has been of increasing interest in recent years and represents an interesting measurement challenge. Measurement requires the insertion of a beam splitter in the source. This also scatters light and, because it is closer to the receiver than the sample, dramatically raises the NEBSDF. Diffuse samples can be measured this way, but not much else. A clever (high tech) Doppler-shift technique, employing a moving sample, has been reported²⁹ that allows separation of beam-splitter scatter from sample scatter and allows measurement of mirror scatter. A more economical (low tech) approach simply involves moving the source chopper to a location between the receiver and the sample.³⁰ Beam-splitter scatter is now dc and goes unnoticed by the ac-sensitive receiver. Noise floor is now limited by scatter from the chopper which must be made from a low-scatter, specular, absorbing material. Noise floors as low as $3 \times 10^{-8} \text{ sr}^{-1}$ have been achieved.

1.7 INCIDENT POWER MEASUREMENT, SYSTEM CALIBRATION, AND ERROR ANALYSIS

Regardless of the type of BSDF measurement, the degree of confidence in the results is determined by instrument calibration, as well as by attention to the measurement limitations previously discussed. Scatter measurements have often been received with considerable skepticism. In part, this has been due to misunderstanding of the definition of BSDF and confusion about various measurement subtleties, such as instrument signature or aperture convolution. However, quite often the measurements have been wrong and the skepticism is justified.

Instrument calibration is often confused with the measurement of P_p , which is why these topics are covered in the same section. To understand the source of this confusion, it is necessary to first consider the various quantities that need to be measured to calculate the BSDF. From Eq. (1), they

are P_s , θ_s , Ω , and P_i . The first two require measurement over a wide range of values. In particular, P_s , which may vary over many orders of magnitude, is a problem. In fact, linearity of the receiver to obtain a correct value of P_s , is a key calibration issue. Notice that an absolute measurement of P_s is not required, as long as the P_s/P_i ratio is correctly evaluated. P_i and Ω generally take on only one (or just a few) discrete values during a data scan. The value of Ω is determined by system geometry. The value of P_i is generally measured in one of two convenient ways.¹¹

The first technique, sometimes referred to as *the absolute method* makes use of the scatter detector (and sometimes a neutral density filter) to directly measure the power incident upon the sample. This method relies on receiver linearity (as does the overall calibration of BSDF) and on filter accuracy when one is used. The second technique, sometimes referred to as *the reference method* makes use of a known BSDF reference sample (usually a diffuse reflector and unfortunately often referred to as the “calibration sample”) to obtain the value of P_i . Scatter from the reference sample is measured and the result used to infer the value of P_i via Eq. (1). The $P_i\Omega$ product may be evaluated this way. This method depends on knowing the absolute BSDF of the reference. Both techniques become more difficult in the mid-IR, where “known” neutral density filters and “known” reference samples are difficult to obtain. Reference sample uniformity in the mid-IR is often the critical issue and care must be exercised. Variations at 10.6 μm as large as 7:1 have been observed across the face of a diffuse gold reference “of known BRDF.” The choice of measurement methods is usually determined by whether it is more convenient to measure the BSDF of a reference or the total power P_i . Both are equally valid methods of obtaining P_i . However, neither method constitutes a system calibration, because calibration issues such as an error analysis and a linearity check over a wide range of scatter values are not addressed over the full range of BSDF angles and powers when P_i is measured (or calculated). The use of a reference sample is an excellent system check regardless of how P_i is obtained.

System linearity is a key part of system calibration. In order to measure linearity, the receiver transfer characteristic, signal out as a function of light in, must be found. This may be done through the use of a known set of neutral density filters or through the use of a comparison technique³¹ that makes use of two data scans—with and without a single filter. However, there are other calibration problems than just linearity. The following paragraph outlines an error analysis for BSDF systems.

Because the calculation of BSDF is very straightforward, the sources of error can be examined through a simple analysis^{11,32} under the assumption that the four defining parameters are independent.

$$\frac{\Delta\text{BSDF}}{\text{BSDF}} = \left[\left(\frac{\Delta P_s}{P_s} \right)^2 + \left(\frac{\Delta P_i}{P_i} \right)^2 + \left(\frac{\Delta \Omega}{\Omega} \right)^2 + \left(\frac{\Delta \theta_s \sin \theta_s}{\cos^2 \theta_s} \right)^2 \right]^{1/2} \quad (5)$$

In similar fashion, each of these terms may be broken into the components that cause errors in it. When this is done, the total error may be found as a function of angle. Two high-error regions are identified. The first is the near specular region (inside 1°), where errors are dominated by the accuracy to which the receiver aperture can be located in the cross-section direction. Or, in other words, did the receiver scan exactly through the specular beam, or did it just miss it? The second relatively high error region is near the sample plane where $\cos \theta_s$ approaches zero. In this region, a small error in angular position results in a large error in calculated BSDF. These errors are often seen in BSDF data as an abrupt increase in calculated BSDF in the grazing scatter direction, the result of division by a very small cosine into the signal gathered by a finite receiver aperture (and/or a dc offset voltage in the detector electronics). This is another example where use of the cosine-corrected BSDF makes more sense.

Accuracy is system dependent; however, at signal levels well above the NEBSDF, uncertainties less than ± 10 percent can be obtained away from the near specular and grazing directions. With expensive electronics and careful error analysis, these inaccuracies can be reduced to the ± 1 percent level.

Full calibration is not required on a daily basis. Sudden changes in instrument signature are an indication of possible calibration problems. Measurement of a reference sample that varies over several orders of magnitude is a good system check. It is prudent to take such a reference scan with data sets in case the validity of the data is questioned at a later time. A diffuse sample, with nearly constant BRDF, is a good reference choice for the measurement of P_i but a poor one for checking system calibration.

1.8 SUMMARY

The art of scatter measurement has evolved to an established form of metrology within the optics industry. Because scatter measurements tend to be a little more complicated than many other optical metrology procedures, a number of key issues must be addressed to obtain useful information. System specifications and measurements need to be given in terms of accepted, well-defined (and understood) quantities (BSDF, TIS, etc.). All parameters associated with a measurement specification need to be given (such as angle limits, receiver solid angles, noise floors, wavelength, etc.). Measurement of near specular scatter and/or low BSDF values are particularly difficult and require careful attention to instrument signature values; however, if the standardized procedures are followed, the result will be repeatable, accurate data.

TIS and BSDF are widely accepted throughout the industry and their measurement is defined by SEMI and ASTM standards. Scatter measurements are used routinely as a quality check on optical components. BSDF specifications are now often used (as they should be) in place of scratch/dig or rms roughness, when scatter is the issue. Conversion of surface scatter data to other useful formats, such as surface roughness statistics, is commonplace. The sophistication of the instrumentation (and analysis) applied to these problems is still increasing. Out-of-plane measurements and polarization-sensitive measurements are two areas that are experiencing rapid advances. Measurement of scatter outside the optics community is also increasing. Although the motivation for scatter measurement differs in industrial situations, the basic measurement and instrumentation issues encountered are essentially the ones described here.

1.9 REFERENCES

1. H. E. Bennett and J. O. Porteus, "Relation between Surface Roughness and Specular Reflectance at Normal Incidence," *J. Opt. Soc. Am.* **51**:123 (1961).
2. H. Davies, "The Reflection of Electromagnetic Waves from a Rough Surface," *Proc. Inst. Elec. Engrs.* **101**:209 (1954).
3. J. C. Stover, "Roughness Measurement by Light Scattering," in A. J. Glass and A. H. Guenther (eds.), *Laser Induced Damage in Optical Materials*, U. S. Govt. Printing Office, Washington, D. C., 1974, p. 163.
4. J. E. Harvey, "Light Scattering Characteristics of Optical Surfaces," Ph.D. dissertation, University of Arizona, 1976.
5. E. L. Church, H. A. Jenkinson, and J. M. Zavada, "Measurement of the Finish of Diamond-Turned Metal Surfaces By Differential Light Scattering," *Opt. Eng.* **16**:360 (1977).
6. J. C. Stover and C. H. Gillespie, "Design Review of Three Reflectance Scatterometers," *Proc. SPIE* **362** (Scattering in Optical Materials):172 (1982).
7. J. C. Stover, "Roughness Characterization of Smooth Machined Surfaces by Light Scattering," *Appl. Opt.* **14** (N8):1796 (1975).
8. E. L. Church and J. M. Zavada, "Residual Surface Roughness of Diamond-Turned Optics," *Appl. Opt.* **14**:1788 (1975).
9. E. L. Church, H. A. Jenkinson, and J. M. Zavada, "Relationship between Surface Scattering and Microtopographic Features," *Opt. Eng.* **18**(2):125 (1979).
10. E. L. Church, "Surface Scattering," in M. Bass (ed.), *Handbook of Optics*, vol. I, 2d ed., McGraw-Hill, New York, 1994.
11. J. C. Stover, *Optical Scattering Measurement and Analysis*, SPIE Press, (1995—new edition to be published in 2009).
12. F. E. Nicodemus, J. C. Richmond, J. J. Hsia, I. W. Ginsberg, and T. Limperis, *Geometric Considerations and Nomenclature for Reflectance*, NBS Monograph 160, U. S. Dept. of Commerce, 1977.
13. W. L. Wolfe and F. O. Bartell, "Description and Limitations of an Automated Scatterometer," *Proc. SPIE* **362**:30 (1982).
14. D. R. Cheever, F. M. Cady, K. A. Klicker, and J. C. Stover, "Design Review of a Unique Complete Angle-Scatter Instrument (CASI)," *Proc. SPIE* **818** (Current Developments in Optical Engineering II):13 (1987).
15. P. R. Spyak and W. L. Wolfe, "Cryogenic Scattering Measurements," *Proc. SPIE* **967**:15 (1989).

16. W. L. Wolfe, K. Magee, and D. W. Wolfe, "A Portable Scatterometer for Optical Shop Use," *Proc. SPIE* **525**:160 (1985).
17. J. Rifkin, "Design Review of a Complete Angle Scatter Instrument," *Proc. SPIE* **1036**:15(1988).
18. T. A. Leonard and M. A. Pantoliano, "BRDF Round Robin," *Proc. SPIE* **967**:22 (1988).
19. T. F. Schiff, J. C. Stover, D. R. Cheever, and D. R. Bjork, "Maximum and Minimum Limitations Imposed on BSDF Measurements," *Proc. SPIE* **967** (1988).
20. F. M. Cady, M. W. Knighton, D. R. Cheever, B. D. Swimley, M. E. Southwood, T. L. Hundtoft, and D. R. Bjork, "Design Review of a Broadband 3-D Scatterometer," *Proc. SPIE* **1753**:21 (1992).
21. K. A. Klicker, J. C. Stover, D. R. Cheever, and F. M. Cady, "Practical Reduction of Instrument Signature in Near Specular Light Scatter Measurements," *Proc. SPIE* **818**:26 (1987).
22. S. J. Wein and W. L. Wolfe, "Gaussian Apodized Apertures and Small Angle Scatter Measurements," *Opt. Eng.* **28**(3):273–280 (1989).
23. J. C. Stover and M. L. Bernt, "Very Near Specular Measurement via Incident Angle Scaling," *Proc. SPIE* **1753**:16 (1992).
24. F. M. Cady, D. R. Cheever, K. A. Klicker, and J. C. Stover, "Comparison of Scatter Data from Various Beam Dumps," *Proc. SPIE* **818**:21 (1987).
25. W. S. Bickle and G. W. Videen, "Stokes Vectors, Mueller Matrices and Polarized Light: Experimental Applications to Optical Surfaces and All Other Scatterers," *Proc. SPIE* **1530**:02 (1991).
26. T. F. Schiff, D. J. Wilson, B. D. Swimley, M. E. Southwood, D. R. Bjork, and J. C. Stover, "Design Review of a Unique Out-of-Plane Polarimetric Scatterometer," *Proc. SPIE* **1753**:33 (1992).
27. T. F. Schiff, D. J. Wilson, B. D. Swimley, M. E. Southwood, D. R. Bjork, and J. C. Stover, "Mueller Matrix Measurements with an Out-Of-Plane Polarimetric Scatterometer," *Proc. SPIE* **1746**:33 (1992).
28. T. F. Schiff, B. D. Swimley, and J. C. Stover, "Mueller Matrix Measurements of Scattered Light," *Proc. SPIE* **1753**:34 (1992).
29. Z. H. Gu, R. S. Dummer, A. A. Maradudin, and A. R. McGurn, "Experimental Study of the Opposition Effect in the Scattering of Light from a Randomly Rough Metal Surface," *Appl. Opt.* **28**(N3):537 (1989).
30. T. F. Schiff, D. J. Wilson, B. D. Swimley, M. E. Southwood, D. R. Bjork, and J. C. Stover, "Retroreflections on a Low Tech Approach to the Measurement of Opposition Effects," *Proc. SPIE* **1753**:35 (1992).
31. F. M. Cady, D. R. Bjork, J. Rifkin, and J. C. Stover, "Linearity in BSDF Measurement," *Proc. SPIE* **1165**:44 (1989).
32. F. M. Cady, D. R. Bjork, J. Rifkin, and J. C. Stover, "BRDF Error Analysis," *Proc. SPIE* **1165**:13 (1989).
33. SEMI ME1392-0305—*Guide for Angle Resolved Optical Scatter Measurements on Specular or Diffuse Surfaces.*
34. SEMI MF1048-1105—*Test Method for Measuring the Reflective Total Integrated Scatter.*
35. SEMI MF1811-0704—*Guide for Estimating the Power Spectral Density Function and Related Finish Parameters from Surface Profile Data.*
36. ASTM E2387-05 *Standard Practice for Goniometric Optical Scatter Measurements.*
37. SEMI M50-0307—*Test Method for Determining Capture Rate and False Count Rate for Surface Scanning Inspection Systems by the Overlay Method.*
38. SEMI M52-0307—*Guide for Specifying Scanning Surface Inspection Systems for Silicon Wafers for the 130 nm, 90 nm, 65 nm, and 45 nm Technology Generations.*
39. SEMI M53-0706—*Practice for Calibrating Scanning Surface Inspection Systems Using Certified Depositions of Monodisperse Polystyrene Latex Spheres on Unpatterned Semiconductor Wafer Surfaces.*
40. SEMI M58-0704—*Test Method for Evaluating DMA Based Particle Deposition Systems and Processes.*
41. T. A. Germer and C. C. Asmail, "Goniometric Optical Scatter Instrument for Out-of-Plane Ellipsometry Measurements," *Rev. Sci. Instrum.* **70**:3688–3695 (1999).
42. T. A. Germer, "Measuring Interfacial Roughness by Polarized Optical Scattering," in A. A. Maradudin (ed.), *Light Scattering and Nanoscale Surface Roughness*, Springer, New York, 2007, Chap. 10, pp. 259–284.
43. B. DeBoo, J. Sasian, and R. Chipman, "Depolarization of Diffusely Reflecting Manmade Objects," *Appl. Opt.* **44**(26):5434–5445 (2005).
44. B. DeBoo, J. Sasian, and R. Chipman, "Degree of Polarization Surfaces and Maps for Analysis of Depolarization," *Optics Express* **12**(20):4941–4958 (2004).

This page intentionally left blank.

DO NOT DUPLICATE

SPECTROSCOPIC MEASUREMENTS

Brian Henderson

*Department of Physics and Applied Physics
University of Strathclyde
Glasgow, United Kingdom*

2.1 GLOSSARY

A_{ba}	Einstein coefficient for spontaneous emission
a_0	Bohr radius
B_{if}	Einstein coefficient between initial state $ i\rangle$ and final state $ f\rangle$
e	charge on the electron
ED	electric dipole term
E_{DC}	Dirac Coulomb term
E_{hf}	hyperfine energy
E_n	eigenvalues of quantum state n
EQ	electric quadrupole term
$\mathbf{E}(t)$	electric field at time t
$\mathbf{E}(\omega)$	electric field at frequency ω
g_a	degeneracy of ground level
g_b	degeneracy of excited level
g_N	gyromagnetic ratio of nucleus
h	Planck's constant
H_{SO}	spin-orbit interaction Hamiltonian
I	nuclear spin
$I(t)$	emission intensity at time t
\mathbf{j}	total angular momentum vector given by $\mathbf{j} = \mathbf{l} \pm \frac{1}{2}$
l_i	orbital state
m	mass of the electron
MD	magnetic dipole term
M_N	mass of nucleus N
$n_\omega(T)$	equilibrium number of photons in a blackbody cavity radiator at angular frequency ω and temperature T

QED	quantum electrodynamics
$R_n(r)$	radial wavefunction
R_∞	Rydberg constant for an infinitely heavy nucleus
s	spin quantum number with value $\frac{1}{2}$
s_i	electronic spin
T	absolute temperature
W_{ab}	transition rate in absorption transition between states a and b
W_{ba}	transition rate in emission transition from state b to state a
Z	charge on the nucleus
$\alpha = e^2/4\pi\epsilon_0\hbar c$	fine structure constant
$\Delta\omega$	natural linewidth of the transition
$\Delta\omega_D$	Doppler width of transition
ϵ_0	permittivity of free space
$\zeta(r)$	spin-orbit parameter
μ_B	Bohr magneton
$\rho(\omega)$	energy density at frequency ω
τ_R	radiative lifetime
ω	angular frequency
ω_k	mode k with angular frequency ω
$\langle f v^1 i\rangle$	matrix element of perturbation V

2.2 INTRODUCTORY COMMENTS

The conceptual basis of optical spectroscopy and its relationship to the electronic structure of matter as presented in the chapter entitled “Optical Spectroscopy and Spectroscopic Lineshapes” in Vol. I, Chap. 10 of this *Handbook*. The chapter entitled “Optical Spectrometers” in Vol. I, Chap. 31 of this *Handbook* discusses the operating principles of optical spectrometers. This chapter illustrates the underlying themes of the earlier ones using the optical spectra of atoms, molecules, and solids as examples.

2.3 OPTICAL ABSORPTION MEASUREMENTS OF ENERGY LEVELS

Atomic Energy Levels

The interest in spectroscopic measurements of the energy levels of atoms is associated with tests of quantum theory. Generally, the optical absorption and luminescence spectra of atoms reveal large numbers of sharp lines corresponding to transitions between the stationary states. The hydrogen atom has played a central role in atomic physics because of the accuracy with which relativistic and quantum electrodynamic shifts to the one-electron energies can be calculated and measured. Tests of quantum electrodynamics usually involve transitions between low-lying energy states (i.e., states with small principal quantum number). For the atomic states $|a\rangle$ and $|b\rangle$, the absorption and luminescence lines occur at exactly the same wavelength and both spectra have the same gaussian lineshape. The $1s \rightarrow 2p$ transitions on atomic hydrogen have played a particularly prominent role, especially since the development of sub-Doppler laser spectroscopy.¹ Such techniques resulted in values of $R_\infty = 10973731.43 \text{ m}^{-1}$, 36.52 m^{-1} for the spin-orbit splitting in the $n = 2$ state and a Lamb shift of 3.53 m^{-1} in the $n = 1$ state. Accurate isotope shifts have been determined from hyperfine structure measurements on hydrogen, deuterium, and tritium.²

Helium is the simplest of the multielectron atoms, having the ground configuration ($1s^2$). The energy levels of helium are grouped into singlet and triplet systems. The observed spectra arise *within* these systems (i.e., singlet-to-singlet and triplet-to-triplet); normally transitions between singlet and triplet levels are not observed. The lowest-lying levels are 1^1S , 2^3S , 2^1S , 2^3P , and 2^1P in order of increasing energy. The $1^1S \rightarrow 2^1S$ splitting is of order 20.60 eV and transitions between these levels are not excited by photons. Transitions involving the 2^1S and 2^3S levels, respectively, and higher-lying spin singlet and spin triplet states occur at optical wavelengths. Experimental work on atomic helium has emphasized the lower-lying triplet levels, which have long excited-state lifetimes and large quantum electrodynamic (QED) shifts.

As with hydrogen, the spectra of He atoms are inhomogeneously broadened by the Doppler effect. Precision measurements have been made using two-photon laser spectroscopy (e.g., $2^3S \rightarrow n^3S$ ($n = 4 - 6$) and n^3D ($n = 3 - 6$), or laser saturation absorption spectroscopy ($2^3S \rightarrow 2^3P$ and $3^3P \rightarrow 3^3D$).³⁻⁶ The $2^1S \rightarrow 3^1P$ and two photon $2^1S \rightarrow n^1D$ ($n = 3 - 7$) spectra have been measured using dye lasers.^{7,8} The wide tune ranging of the Ti-sapphire laser and the capability for generating frequencies not easily accessible with dye lasers using frequency-generation techniques makes it an ideal laser to probe transitions starting on the 2S levels of He.⁹ Two examples are the two-photon transition $2^3S \rightarrow 3^3S$ at 855 nm and the $2^3S \rightarrow 3^3P$ transition of 389 nm. The power of Doppler-free spectroscopy is shown to advantage in measurements of the $2^3S \rightarrow 2^3P$ transition.¹⁰ Since both 3^3S and 3^3P are excited levels, the homogeneous width is determined by the sum of the reciprocal lifetimes of the two levels. Since both 2^3S and 2^3P levels are long-lived, the resulting homogeneous width is comparatively narrow. Figure 1a shows the Doppler-broadened profile of the $2^3S \rightarrow 2^3P$ transition of ^4He for which the FWHM is about 5.5 GHz. The inhomogeneously broadened line profile shown in Fig. 1a also shows three very weak “holes” corresponding to saturated absorption of the Ti-sapphire laser radiation used to carry out the experiment. These components correspond to the $2^3S_1 \rightarrow 3^3P_2$ and $2^3S_1 \rightarrow 3^3P_1$ transitions and their crossover resonance. The amplitude of the saturated signal is some 1–2 percent of the total absorption. The relativistic splittings including spin-orbit coupling of $3^3P_0 - 3^3P_1$ and $3^3P_1 - 3^3P_2$ are 8.1 GHz and 658.8 MHz, respectively. Frequency modulation of the laser (see Vol. I, Chap. 31 of this *Handbook*) causes the “hole” to be detected as a first derivative of the absorption line, Fig. 1b. The observed FWHM of the Doppler-free signals was only 20 MHz. The uncertainty in the measured $2^3S \rightarrow 3^3P$ interval was two parts in 10^9 (i.e., 1.5 MHz), an improvement by a factor of 60 on earlier measurements. A comparison of the experimental results with recent calculations of the non-QED terms¹¹ gives a value for the one-electron Lamb shift of -346.5 (2.8) MHz, where the uncertainty in the quoted magnitude is in parentheses. The theoretical value is -346.3 (13.9) MHz. Finally, the frequencies of the $2^3S_1 \rightarrow 3^3P_1$, 3^3P_2 transitions were determined to be 25708.60959 (5) cm^{-1} and 25708.58763 (5) cm^{-1} , respectively.

The H^- ion is another two-electron system, of some importance in astrophysics. It is the simplest quantum mechanical three-body species. Approximate quantum mechanical techniques give a wave function and energy eigenvalue which are exact, for all practical purposes. Experimentally, the optical absorption spectrum of H^- is continuous, a property of importance in understanding the opacity of the sun. H^- does not emit radiation in characteristic emission lines. Instead the system sheds its excess (absorbed) energy by ejecting one of the electrons. The radiant energy associated with the ejected electron consists of photons having a continuous energy distribution. Recent measurements with high-intensity pulsed lasers, counterpropagating through beams of 800-MeV H^- ions, have produced spectacular “doubly excited states” of the He^- ion. Such ions are traveling at 84 percent of the velocity of light and, in this situation, the visible laboratory photons are shifted into the vacuum ultraviolet region. At certain energies the H^- ion is briefly excited into a system in which both electrons are excited prior to one of the electrons being ejected. Families of new resonances up to the energy level $N = 8$ have been observed.¹² These resonances are observed as windows in the continuous absorption spectrum, at energies given by a remarkably simple relation reminiscent of the Bohr equation for the Balmer series in hydrogen.¹³ The resonance line with the lowest energy in each family corresponds to both electrons at comparable distances from the proton.

These experiments on H^- are but one facet of the increasingly sophisticated measurements designed to probe the interaction between radiation and matter driven by experiments in laser technology. “Quantum jump” experiments involving single ions in an electromagnetic trap have become almost commonplace. Chaos has also become a rapidly growing subfield of atomic spectroscopy.¹⁴

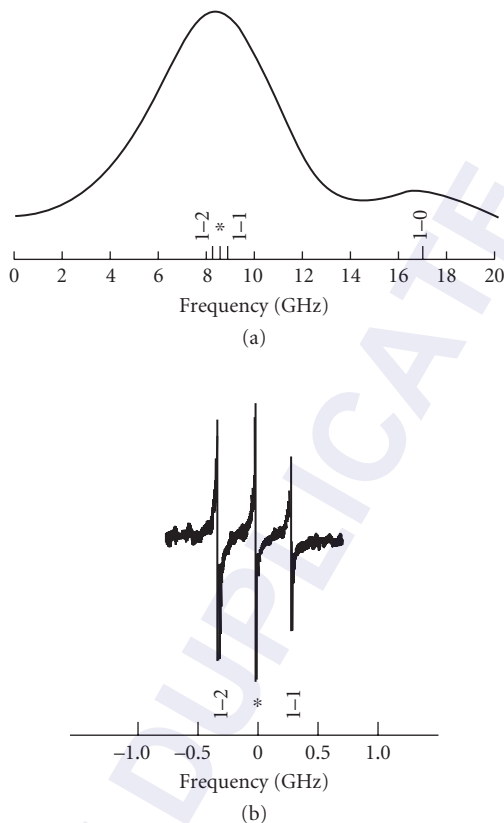


FIGURE 1 (a) Inhomogeneously broadened line profile of the $2^3S \rightarrow 2^3P$ absorption in ^4He , including the weak “holes” due to saturated absorption and the position of the $2^3S \rightarrow 2^3P_2$, $3P_1$, and $3P_0$ components. (b) Doppler-free spectra showing the $2^3S \rightarrow 2^3P_2$, $3P_1$ transitions, and the associated crossover resonance. (After Adams, Riis, and Ferguson.¹⁰)

The particular conditions under which chaos may be observed in atomic physics include hydrogenic atoms in strong homogeneous magnetic fields such that the cyclotron radius of the electron approaches the dimensions of the atomic orbitals. A more easily realizable situation using magnetic field strengths of only a few tesla uses highly excited orbitals close to the ionization threshold. Iu et al.¹⁵ have reported the absorption spectrum of transitions from the $3s$ state of Li to bound and continuum states near the ionization limit in a magnetic field of approximately six tesla. There is a remarkable coincidence between calculations involving thousands of energy levels and experiments involving high-resolution laser spectroscopy.

Atomic processes play an important role in the energy balance of plasmas, whether they be created in the laboratory or in star systems. The analysis of atomic emission lines gives much information on the physical conditions operating in a plasma. In laser-produced plasmas, the densities of charged ions may be in the range $10^{20} - 10^{25}$ ions cm^{-3} , depending on the pulse duration of the laser. The spectra of many-electron ions are complex, and may have the appearance of an unresolved transition

array between states belonging to specific initial and final configurations. Theoretical techniques have been developed to determine the average ionization state of the plasma from the observed optical spectrum. In many cases, the spectra are derived from ionic charge states in the nickel-like configuration containing 28 bound electrons. In normal nickel, the outershell configuration is $(3d^8)(4s^2)$, the 4s levels having filled before 3d because the electron-electron potentials are stronger than electron-nuclear potentials. However, in highly ionized systems, the additional electron-nuclear potential is sufficient to shift the configuration from $(3d^8)(4s^2)$ to the closed shell configuration $(3d^{10})$. The resulting spectrum is then much simpler than for atomic nickel. The Ni-like configuration has particular relevance in experiments to make x-ray lasers. For example, an analog series of collisionally pumped lasers using Ni-like ions has been developed, including a Ta⁴⁵⁺ laser operating at 4.48 nm and a W⁴⁶⁺ laser operating at 4.32 nm.¹⁶

Molecular Spectroscopy

The basic principles of gas-phase molecular spectroscopy were also discussed in “Optical Spectroscopy and Spectroscopic Lineshapes,” Vol. I, Chap. 10 of this *Handbook*. The spectra of even the simplest molecules are complicated by the effects of vibrations and of rotations about an axis. This complexity is illustrated elsewhere in this *Handbook* in Fig. 8 of this chapter, which depicts a photographically recorded spectrum of the $2\Pi \rightarrow 3\Sigma$ bands of the diatomic molecule NO, which was interpreted in terms of progressions and line sequences associated with the *P*-, *Q*-, and *R*-branches. The advent of Fourier transform spectroscopy led to great improvements in resolution and greater efficiency in revealing all the fine details that characterize molecular spectra. Figure 2a is a Fourier-transform infrared spectrum of nitrous oxide, N₂O, which shows the band center at 2462 cm⁻¹ flanked by the *R*-branch and a portion of the *P*-branch; the density of lines in the *P*- and *R*-branch is evident. On an expanded scale, in Fig. 2b, there is a considerable simplification of the rotational-vibrational structure at the high-energy portion of the *P*-branch. The weaker lines are the so-called “hot bands.”

More precise determinations of the transition frequencies in molecular physics are measured using Lamb dip spectroscopy. The spectrum shown in Fig. 3a is a portion of the laser Stark spectrum of methyl fluoride measured using electric fields in the range 20 to 25 KV cm⁻¹ with the 9- μ m *P*(18) line of the CO₂ laser, which is close to the ν_3 band origin of CH₃F.¹⁷ The spectrum in Fig. 3a consists of a set of $\Delta M_J = \pm 1$ transitions, brought into resonance at different values of the static electric field. Results from the alternative high-resolution technique using a supersonic molecular beam and bolometric detector are shown in Fig. 3b: this spectrum was obtained using CH₃F in He mixture expanded through a 35- μ m nozzle. The different M_J components of the *Q*(1, 0), *Q*(2, 1), and *Q*(3, 3) components of the *Q*-branch are shown to have very different intensities relative to those in Fig. 3a on account of the lower measurement temperature.

There has been considerable interest in the interaction of intense laser beams with molecules. For example, when a diatomic molecule such as N₂ or CO is excited by an intense (10^{15} W cm⁻²), ultrashort (0.6 ps) laser pulse, it multiply ionizes and then fragments as a consequence of Coulomb repulsion. The charge and kinetic energy of the resultant ions can be determined by time-of-flight (TOF) mass spectrometry. In this technique, the daughter ions of the “Coulomb explosion” drift to a cathode tube at different times, depending on their weight. In traditional methods, the TOF spectrum is usually averaged over many laser pulses to improve the signal-to-noise ratio. Such simple averaging procedures remove the possibility of correlations between particular charged fragments. This problem was overcome by the *covariance mapping* technique developed by Frasiniski and Codling.¹⁸ Experimentally, a linearly polarized laser pulse with **E**-vector pointing toward the detector is used to excite the molecules, which line up with their internuclear axis parallel to the **E**-field. Under Coulomb explosion, one fragment heads toward the detector and the other away from the detector. The application of a dc electric field directs the “backward” fragment ion to the detector, arriving at some short time after the forward fragment. This temporal separation of two fragments arriving at the detector permits the correlation between molecular fragments to be retained. In essence, the TOF spectrum, which plots molecular weight versus counts, is arranged both horizontally (forward ions) and vertically

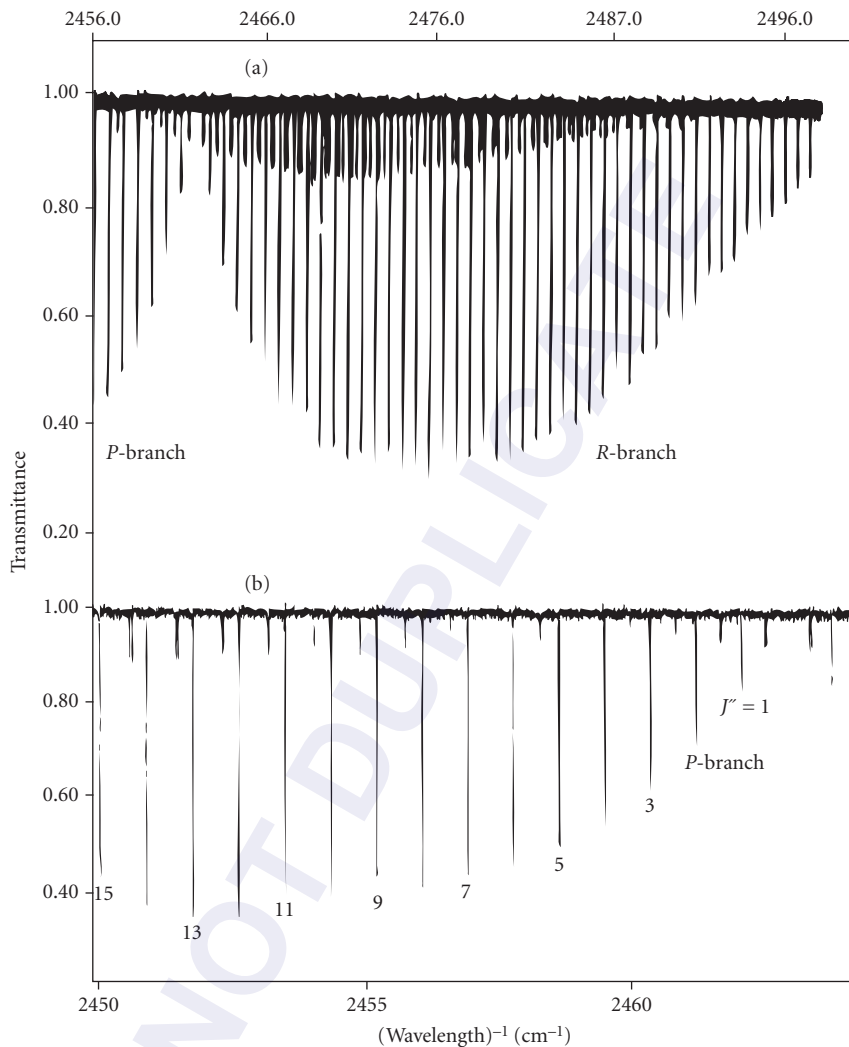


FIGURE 2 (a) *P*- and *R*-branches for the nitrous oxide molecule measured using Fourier-transform infrared spectroscopy. The gas pressure was 0.5 torr and system resolution 0.006 cm⁻¹, (b) on an expanded scale, the high-energy portion of the *P*-branch up to $J'' = 15$.

(backward ions) on a two-dimensional graph. A coordinate point on a preliminary graph consists of two ions along with their counts during a single pulse. Coordinates from 10⁴ pulses or so are then assembled in a final map. Each feature on the final map relates to a specific fragmentation channel, i.e., the pair of fragments and their parent molecule. The strength of the method is that it gives the probability for the creation and fragmentation of the particular parent ion. Covariance mapping experiments on N₂ show that 610- and 305-nm pulses result in fragmentation processes that are predominantly charge-symmetric. In other words, the Coulomb explosion proceeds via the production of ions with the same charge.

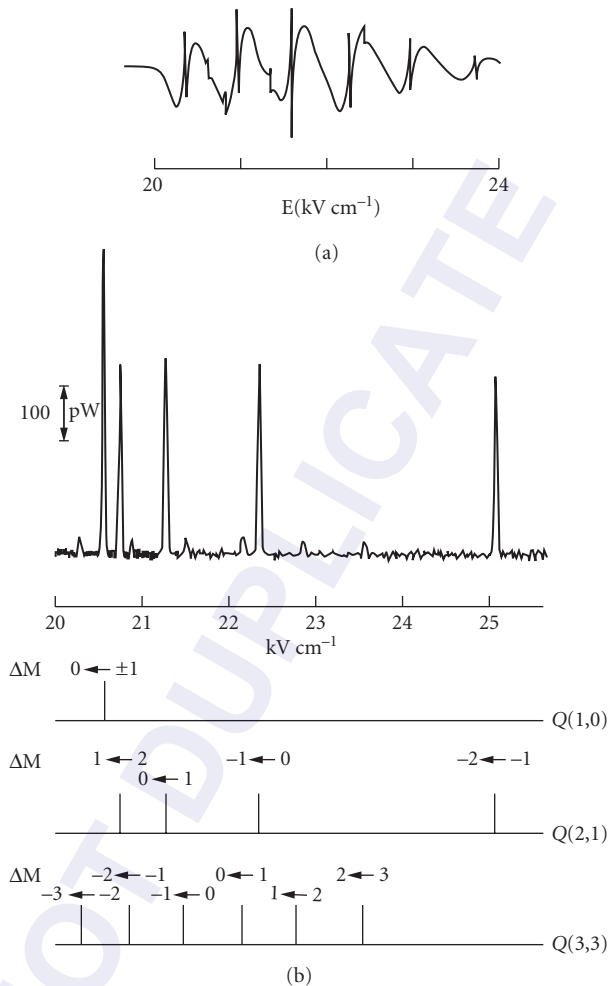


FIGURE 3 (a) Laser Stark absorption spectrum of methyl fluoride measured at 300 K using Lamb dip spectroscopy with a gas pressure of 5 mTorr and (b) the improved resolution obtained using molecular-beam techniques with low-temperature bolometric detection and CH_3F in He expanded through a 35- μm nozzle. (After Douketic and Gough.¹⁷)

Optical Spectroscopy of Solids

One of the more fascinating aspects of the spectroscopy of electronic centers in condensed matter is the variety of lineshapes displayed by the many different systems. Those discussed in Vol. I, Chap. 10 include Nd^{3+} in YAG (Fig. 6), O_2^- in KBr (Fig. 11), Cr^{3+} in YAG (Fig. 12), and F centers in KBr (Fig. 13). The very sharp Nd^{3+} lines (Fig. 6) are zero-phonon lines, inhomogeneously broadened by strain. The abundance of sharp lines is characteristic of the spectra of trivalent rare-earth ions in ionic crystals. Typical low-temperature linewidths for Nd^{3+} : YAG are 0.1–0.2 cm^{-1} . There is particular interest in the spectroscopy of Nd^{3+} because of the efficient laser transitions from the $^4\text{F}_{3/2}$ level into the $^4\text{I}_J$ -manifold. The low-temperature luminescence transitions between $^4\text{F}_{3,2} \rightarrow ^4\text{I}_{15/2}$, $^4\text{I}_{13/2}$, $^4\text{I}_{11/2}$, and $^4\text{I}_{9/2}$ levels are

shown in Fig. 6: all are split by the effects of the crystalline electric field. Given the relative sharpness of these lines, it is evident that the Slater integrals $F^{(k)}$, spin-orbit coupling parameters ζ , and crystal field parameters, B_t^k , may be measured with considerable accuracy. The measured values of the $F^{(k)}$ and ζ vary little from one crystal to another.¹⁹ However, the crystal field parameters, B_t^k , depend strongly on the rare-earth ion-ligand-ion separation. Most of the $4f^n$ ions have transitions which are the basis of solid-state lasers. Others such as Eu^{3+} and Tb^{3+} are important red-emitting and green-emitting phosphor ions, respectively.²⁰

Transition-metal ion spectra are quite different from those of the rare-earth ions. In both cases, the energy-level structure may be determined by solving the Hamiltonian

$$H = H_o + H' + H_{so} + H_c \quad (1)$$

in which H_o is a sum of one-electron Hamiltonians including the central field of each ion, H' is the interaction between electrons in the partially filled $3d^n$ or $4f^n$ orbitals, H_{so} is the spin-orbit interaction, and H_c is the interaction of the outer shell electrons with the crystal field. For rare-earth ions $H', H_{so} \gg H_c$, and the observed spectra very much reflect the free-ion electronic structure with small crystal field perturbations. The spectroscopy of the transition-metal ions is determined by the relative magnitudes of $H' \approx H_c \gg H_{so}$.^{19,21} The simplest of the transition-metal ions is Ti^{3+} : in this $3d^1$ configuration a single $3d$ electron resides outside the closed shells. In this situation, $H' = 0$ and only the effect of H_c needs be considered (Fig. 4a). The Ti^{3+} ion tends to form octahedral complexes, where the $3d$ configuration is split into 2E and 2T_2 states with energy separation $10Dq$. In cation sites with weak, trigonally symmetric distortions, as in Al_2O_3 and $\text{Y}_3\text{Al}_5\text{O}_{12}$, the lowest-lying state, 2T_2 , splits into 2A_1 and 2E states (using the C_{3v} group symmetry labels). In oxides, the octahedral splitting is of order $10Dq \approx 20,000 \text{ cm}^{-1}$ and the trigonal field splitting $\nu \approx 700 - 1000 \text{ cm}^{-1}$. Further splittings of the levels

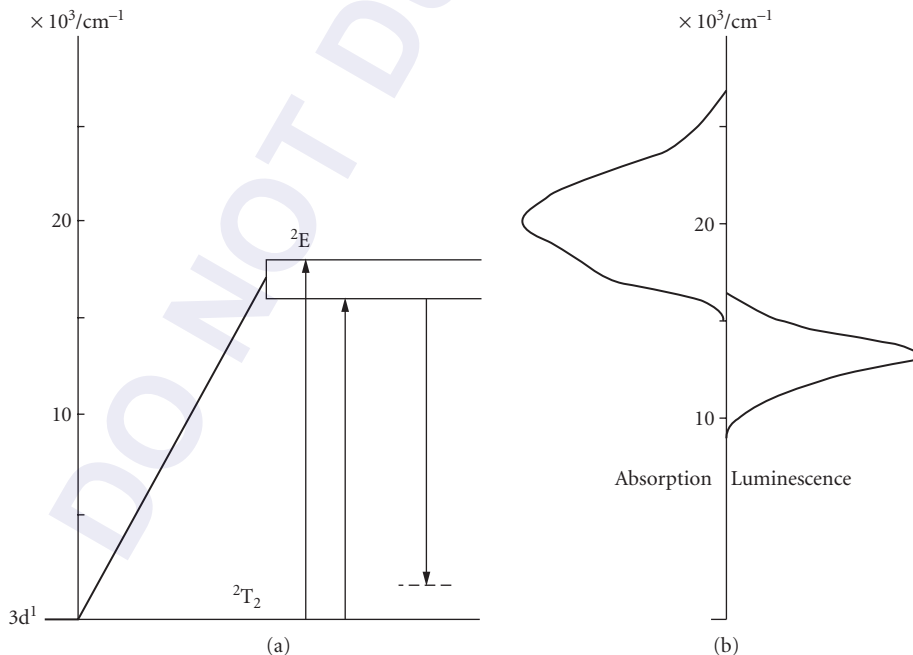


FIGURE 4 Absorption and emission spectra of Ti^{3+} ions [$(3d^1)$ configuration] in Al_2O_3 measured at 300 K.

occur because of spin-orbit coupling and the Jahn-Teller effect. The excited 2E state splits into $2\bar{A}$ and \bar{E} (from 2A_1), and \bar{E} and $2A$ (from 2E). The excited state splitting by a *static* Jahn-Teller effect is large, ~ 2000 to 2500 cm^{-1} , and may be measured from the optical absorption spectrum. In contrast, ground-state splittings are quite small: a *dynamic* Jahn-Teller effect has been shown to strongly quench the spin-orbit coupling ζ and trigonal field splitting ν parameters.²² In $\text{Ti}^{3+}:\text{Al}_2\text{O}_3$ the optical absorption transition, ${}^2T_2 \rightarrow {}^2E$, measured at 300 K, Fig. 4b, consists of two broad overlapping bands separated by the Jahn-Teller splitting, the composite band having a peak at approximately $20,000\text{ cm}^{-1}$. Luminescence occurs only from the lower-lying excited state $2A$, the emission band peak occurring at approximately $14,000\text{ cm}^{-1}$. As Fig. 4 shows, both absorption and emission bands are broad because of strong electron-phonon coupling. At low temperatures the spectra are characterized by weak zero-phonon lines, one in absorption due to transitions from the $2A$ ground state and three in emission corresponding to transitions in the \bar{E} , and $2\bar{A}$ levels of the electronic ground state.^{22,23} These transitions are strongly polarized.

For ions with $3d^n$ configuration it is usual to neglect H_c and H_{so} in Eq. (1), taking into account only the central ion terms and the Coulomb interaction between the $3d$ electrons. The resulting energies of the free-ion LS terms are expressed in terms of the Racah parameters A , B , and C . Because energy differences between states are measured in spectroscopy, only B and C are needed to categorize the free-ion levels. For pure d-functions, $C/B = 4.0$. The crystal field term H_c and H_{so} also are treated as perturbations. In many crystals, the transition-metal ions occupy octahedral or near-octahedral cation sites. The splittings of each free-ion level by an octahedral crystal field depend in a complex manner on B , C , and the crystal field strength Dq given by

$$Dq = \left(\frac{Ze^2}{24\pi\epsilon_0} \right) \frac{\langle r^4 \rangle 3d}{a^5} \quad (2)$$

The parameters D and q always occur as a product. The energy levels of the $3d^n$ transition-metal ions are usually represented on Tanabe-Sugano diagrams, which plot the energies $E(\Gamma)$ of the electronic states as a function of the octahedral crystal field.^{19,21} The crystal field levels are classified by irreducible representations Γ of the octahedral group, O_h . The Tanabe-Sugano diagram for the $3d^3$ configuration, shown in Fig. 5a, was constructed using a C/B ratio = 4.8: the vertical broken line drawn at $Dq/B = 2.8$ is appropriate for Cr^{3+} ions in ruby. If a particular value of C/B is assumed, only two variables, B and Dq , need to be considered: in the diagram $E(\Gamma)/B$ is plotted as a function of Dq/B . The case of ruby, where the 2E level is below 4T_2 , is referred to as the *strong field* case. Other materials where this situation exists include YAlO_3 , $\text{Y}_3\text{Al}_5\text{O}_{12}$ (YAG), and MgO . In many fluorides, Cr^{3+} ions occupy *weak field* sites, where $E({}^4T_2) < E({}^2E)$ and Dq/B is less than 2.2. When the value of Dq/B is close to 2.3, the *intermediate* crystal field, the 4T_2 and 2E states almost degenerate. The value of Dq/B at the level crossing between 4T_2 and 2E depends slightly on the value of C .

The Tanabe-Sugano diagram represents the static lattice. In practice, electron-phonon coupling must be taken into account: the relative strengths of coupling to the states involved in transitions and the consequences may be inferred from Fig. 5a. Essentially ionic vibrations modulate the crystal field experienced by the central ion at the vibrational frequency. Large differences in slope of the E versus Dq graphs indicate large differences in coupling strengths and hence large homogeneous bandwidths. Hence, absorption and luminescence transitions from the 4A_2 ground state to the 4T_2 and 4T_1 states will be broadband due to the large differences in coupling of the electronic energy to the vibrational energy. For the ${}^4A_2 \rightarrow {}^2E, {}^2T_1$ transition, the homogeneous linewidth is hardly affected by lattice vibrations, and sharp line spectra are observed.

The Cr^{3+} ion occupies a central position in the folklore of transition-metal ion spectroscopy, having been studied by spectroscopists for over 150 years. An extensive survey of Cr^{3+} luminescence in many compounds was published as early as 1932.²⁴ The Cr^{3+} ions have the outer shell configuration, $3d^3$, and their absorption and luminescence spectra may be interpreted using Fig. 5. First, the effect of the octahedral crystal field is to remove the degeneracies of the free-ion states 4F and 2G . The ground term of the free ion, ${}^4F_{3/2}$, is split by the crystal field into a ground-state orbital singlet, ${}^4A_{2g}$, and two orbital triplets ${}^4T_{2g}$ and ${}^4T_{1g}$, in order of increasing energy. Using the energy of the 4A_2 ground state as the zero, for all values of Dq/B , the energies of the ${}^4T_{2g}$ and ${}^4T_{1g}$ states are seen to vary strongly as a function of the octahedral crystal field. In a similar vein, the 2G free-ion state splits into 2E , 2T_1 ,

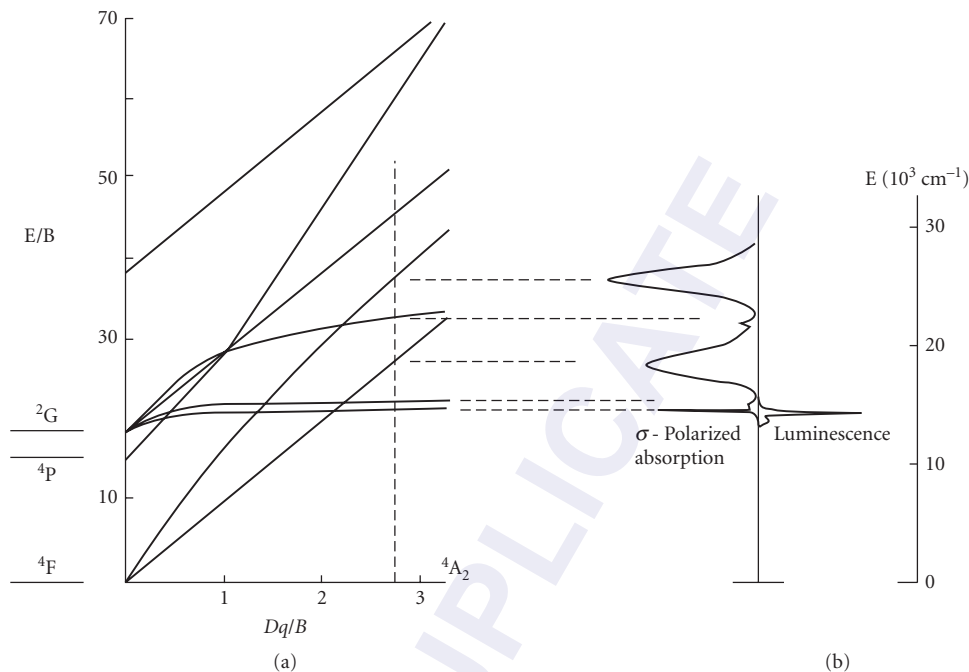


FIGURE 5 Tanabe-Sugano diagram for Cr^{3+} ions with $C/B = 4.8$, appropriate for ruby for which $Dq/B = 2.8$. On the right of the figure are shown the optical absorption and photoluminescence spectrum of ruby measured at 300 K.

2T_2 , and 2A_1 states, the two lowest of which, 2E and 2T_1 , vary very little with Dq . The energies $E(^2T_2)$ and $E(^2A_1)$ are also only weakly dependent on Dq/B . The free-ion term, 4P , which transforms as the irreducible representation 4T_1 of the octahedral group is derived from (e^2t_2) configuration: this term is not split by the octahedral field although its energy is a rapidly increasing function of Dq/B . Low-symmetry distortions lead to strongly polarized absorption and emission spectra.²⁵

The σ -polarized optical absorption and luminescence spectra of ruby are shown in Fig. 5b. The expected energy levels predicted from the Tanabe-Sugano diagram are seen to coincide with appropriate features in the absorption spectrum. The most intense features are the vibronically broadened $^4A_2 \rightarrow ^4T_1, ^4T_2$ transitions. These absorptions occur in the blue and yellow-green regions, thereby accounting for the deep red color of ruby. Many other Cr^{3+} -doped hosts have these bands in the blue and orange-red regions, by virtue of smaller values of Dq/B ; the colors of such materials (e.g., MgO , $\text{Gd}_3\text{Sc}_2\text{Ga}_3\text{O}_{12}$, LiSrAlF_6 , etc.) are different shades of green. The absorption transitions, $^4A_2 \rightarrow ^2E, ^2T_1, ^2T_2$ levels are spin-forbidden and weakly coupled to the phonon spectrum ($S < 0.5$). The spectra from these transitions are dominated by sharp zero-phonon lines. However, the low-temperature photoluminescence spectrum of ruby is in marked contrast to the optical absorption spectrum since only the sharp zero-phonon line (R-line) due to the $^2E \rightarrow ^4A_2$ transition being observed. Given the small energy separations between adjacent states of Cr^{3+} , the higher excited levels decay nonradiatively to the lowest level, 2E , from which photoluminescence occurs across a bandgap of approximately $15,000 \text{ cm}^{-1}$. Accurate values of the parameters Dq , B , and C may be determined from these absorption data. First, the peak energy of the $^4A_2 \rightarrow ^4T_2$ absorption band is equal to $10Dq$. The energy shift between the $^4A_2 \rightarrow ^4T_2, ^4T_1$ bands is dependent on both Dq and B , and the energy separation between the two broad absorption bands is used to determine B . Finally, the position of the R-line varies with Dq , B , and C : in consequence, once Dq and B are known, the magnitude of C may be determined from the position of the $^4A_2 \rightarrow ^2E$ zero-phonon line.

This discussion of the spectroscopy of the Cr^{3+} ion is easily extended to other multielectron configurations. The starting points are the Tanabe-Sugano diagrams collected in various texts.^{19,21} Analogous series of elements occur in the fifth and sixth periods of the periodic table, respectively, where the $4d^n$ (palladium) and $5d^n$ (platinum) groups are being filled. Compared with electrons in the 3d shell, the 4d and 5d shell electrons are less tightly bound to the parent ion. In consequence, charge transfer transitions, in which an electron is transferred from the cation to the ligand ion (or vice versa), occur quite readily. The charge transfer transitions arise from the movement of electronic charge over a typical interatomic distance, thereby producing a large dipole moment and a concomitant large oscillator strength for the absorption process. For the Fe-group ions ($3d^n$ configuration), such charge transfer effects result in the absorption of ultraviolet photons.

For example, the Fe^{3+} ion in MgO absorbs in a broad structureless band with peak at 220 nm and half-width of order 120 nm (i.e., 0.3 eV). The Cr^{2+} ion also absorbs by charge transfer process in this region. In contrast, the palladium and platinum groups have lower-lying charge transfer states. The resulting intense absorption bands in the visible spectrum may overlap spectra due to low-lying crystal field transitions. Rare-earth ions also give rise to intense charge transfer bands in the ultraviolet region.

Various metal cations have been used as broadband visible region phosphors. For example, transitions between the $4f^n$ and $4f^{(n-1)}5d$ levels of *divalent* rare-earth ions give rise to intense broad transitions which overlap many of the sharp $4f^n$ transitions of the *trivalent* rare-earth ions. Of particular interest are Sm^{2+} , Dy^{2+} , Eu^{2+} , and Tm^{2+} . In Sm^{2+} ($4f^6$) broadband absorption transitions from the ground state 7F_1 to the $4f^55d$ level may result in either broadband emission (from the vibrationally relaxed $4f^55d$ level) or sharp line emission from 5D_0 ($4f^6$) depending upon the host crystal. The $4f^55d$ level, being strongly coupled to the lattice, accounts for this variability. There is a similar material-by-material variation in the absorption and emission properties of the Eu^{2+} ($4f^7$ configuration), which has the $^8S_{7/2}$ ground level. The next highest levels are derived from the $4f^65d$ state, which is also strongly coupled to the lattice. This state is responsible for the varying emission colors of Eu^{2+} in different crystals, e.g., violet in $\text{Sr}_2\text{P}_2\text{O}_7$, blue in $\text{BaAl}_{12}\text{O}_{19}$, green in SrAl_2O_4 , and yellow in Ba_2SiO_5 .

The heavy metal ions Tl^+ , In^+ , Ga^+ , Sn^{2+} , and Pb^{2+} may be used as visible-region phosphors. These ions all have two electrons in the ground configuration ns^2 and excited configurations $(ns)(np)$. The lowest-lying excited states, in the limit of Russell Saunders coupling, are then $^1S_0(ns^2)$, $^3P_{0,1,2}$, and 1P_1 from $(ns)(np)$. The spectroscopy of Tl^+ has been much studied especially in the alkali halides. Obviously $^1S_0 \rightarrow ^1P_1$ is the strongest absorption transition, occurring in the ultraviolet region. This is labeled as the C-band in Fig. 6. Next in order of observable intensity is the A-band, which is a spin-forbidden absorption transition $^1S_0 \rightarrow ^3P_1$, in which the relatively large oscillator strength is borrowed from the 1P_1 state by virtue of the strong spin-orbit interaction in these heavy metal ions. The B and D bands, respectively, are due to absorption transitions from 1S_0 to the 3P_2 and 3P_0 states induced by vibronic mixing.²⁶ A phenomenological theory^{26,27} quantitatively accounts for both absorption spectra and the triplet state emission spectra.^{28,29}

The examples discussed so far have all concerned the spectra of ions localized in levels associated with the central fields of the nucleus and closed-shells of electrons. There are other situations which warrant serious attention. These include electron-excess centers in which the positive potential of an anion vacancy in an ionic crystal will trap one or more electrons. The simplest theory treats such a *color center* as a particle in a finite potential well.^{19,27} The simplest such center is the F-center in the alkali halides, which consist of one electron trapped in an anion vacancy. As we have already seen (e.g., Fig. 13 in "Optical Spectroscopy and Spectroscopic Lineshapes", Vol. I, Chap. 10 of this *Handbook*), such centers give rise to broadbands in both absorption and emission, covering much of the visible and near-infrared regions for alkali halides. F-aggregate centers, consisting of multiple vacancies arranged in specific crystallographic relationships with respect to one another, have also been much studied. They may be positive, neutral, or negative in charge relative to the lattice depending upon the number of electrons trapped by the vacancy aggregate.³⁰

Multiquantum wells (MQWs) and strained-layer superlattices (SLs) in semiconductors are yet another type of *finite-well* potential. In such structures, alternate layers of two different semiconductors are grown on top of each other so that the bandgap varies in one dimension with the periodicity

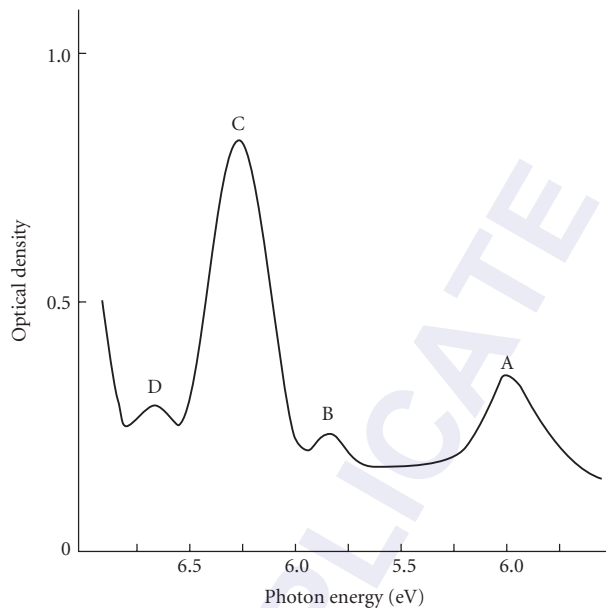


FIGURE 6 Ultraviolet absorption spectrum of Tl^+ ions in KCl measured at 77 K. (After Delbecq *et al.*²⁶)

of the epitaxial layers. A modified Kronig-Penney model is often used to determine the energy eigenvalues of electrons and holes in the conduction and valence bands, respectively, of the narrower gap material. Allowed optical transitions between valence band and conduction band are then subject to the selection rule $\Delta n = 0$, where $n = 0, 1, 2$, etc. The example given in Fig. 7 is for SLSs in the II-VI family of semiconductors ZnS/ZnSe.³¹ The samples were grown by metallo-organic vapor phase epitaxy³²

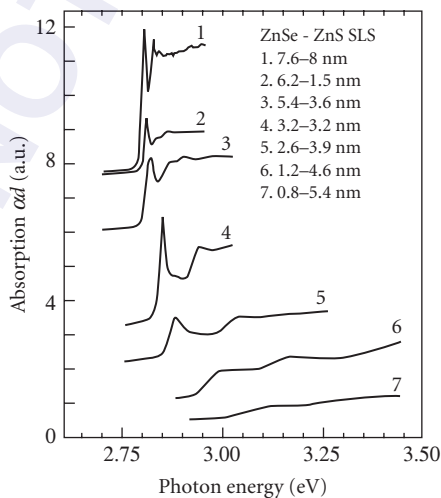


FIGURE 7 Optical absorption spectra of SLSs of ZnS/ZnSe measured at 14 K. (After Fang *et al.*³¹)

with a superlattice periodicity of 6 to 8 nm while varying the thickness of the narrow gap material (ZnSe) between 0.8 and 7.6 nm. The splitting between the two sharp features occurs because the valence band states are split into “light holes” (lh) and “heavy holes” (hh) by spin-orbit interaction. The absorption transitions then correspond to transitions from the $n = 1$ lh- and hh-levels in the valence band to the $n = 1$ electron states in the conduction band. Higher-energy absorption transitions are also observed. After absorption, electrons rapidly relax down to the $n = 1$ level from which emission takes place down to the $n = 1$, lh-level in the valence band, giving rise to a single emission line at low temperature.

2.4 THE HOMOGENEOUS LINESHAPE OF SPECTRA

Atomic Spectra

The homogeneous widths of atomic spectra are determined by the uncertainty principle, and hence by the radiative decaytime, τ_R (as discussed in Vol. I, Chap. 10, “Optical Spectroscopy and Spectroscopic Lineshapes”). Indeed, the so-called natural or homogeneous width of $\Delta\omega$, is given by the Einstein coefficient for spontaneous emission, $A_{ho} = (\tau_R)^{-1}$. The homogeneously broadened line has a Lorentzian lineshape with FWHM given by $(\tau_R)^{-1}$. In gas-phase spectroscopy, atomic spectra are also broadened by the Doppler effect: random motion of atoms broadens the lines in-homogeneously leading to a gaussian-shaped line with FWHM proportional to $(T/M)^{-1/2}$, where T is the absolute temperature and M the atomic mass. Saturated laser absorption or optical hole-burning techniques are among the methods which recover the true homogeneous width of an optical transition. Experimental aspects of these types of measurement were discussed in this *Handbook* in Vol. I, Chap. 31, “Optical Spectroscopy,” and examples of Doppler-free spectra (Figs. 1 and 3, Vol. I, Chap. 10) were discussed in terms of the fundamental tests of the quantum and relativistic structure of the energy levels of atomic hydrogen. Similar measurements were also discussed for the case of He (Fig. 1) and in molecular spectroscopy (Fig. 3). In such examples, the observed lineshape is very close to a true Lorentzian, typical of a lifetime-broadened optical transition.

Zero-Phonon Lines in Solids

Optical hole burning (OHB) reduces the effects of inhomogeneous broadening in solid-state spectra. For rare-earth ions, the homogeneous width amounts to some 0.1–1.0 MHz, the inhomogeneous widths being determined mainly by strain in the crystal. Similarly, improved resolution is afforded by fluorescence line narrowing (FLN) in the R -line of Cr^{3+} (Vol. I, Chap. 10, Fig. 12). However, although the half-width measured using OHB is the true homogeneous width, the observed FLN half-width, at least in resonant FLN, is a convolution of the laser width and twice the homogeneous width of the transition.³³ In solid-state spectroscopy, the underlying philosophy of OHB and FLN experiments may be somewhat different from that in atomic and molecular physics. In the latter cases, there is an intention to relate theory to experiment at a rather sophisticated level. In solids, such high-resolution techniques are used to probe a range of other dynamic processes than the natural decay rate. For example, hole-burning may be induced by photochemical processes as well as by intrinsic lifetime processes.^{34,35} Such photochemical hole-burning processes have potential in optical information storage systems. OHB and FLN may also be used to study distortions in the neighborhood of defects. Figure 8 is an example of Stark spectroscopy and OHB on a zero-phonon line at 607 nm in irradiated NaF.³⁴ This line had been attributed to an aggregate of four F -centers in nearest-neighbor anion sites in the rocksalt-structured lattice on the basis of the polarized absorption/emission measurements. The homogeneous width in zero electric field was only 21 MHz in comparison with the inhomogeneous width of 3 GHz. Interpretation of these results is inconsistent with the four-defect model.

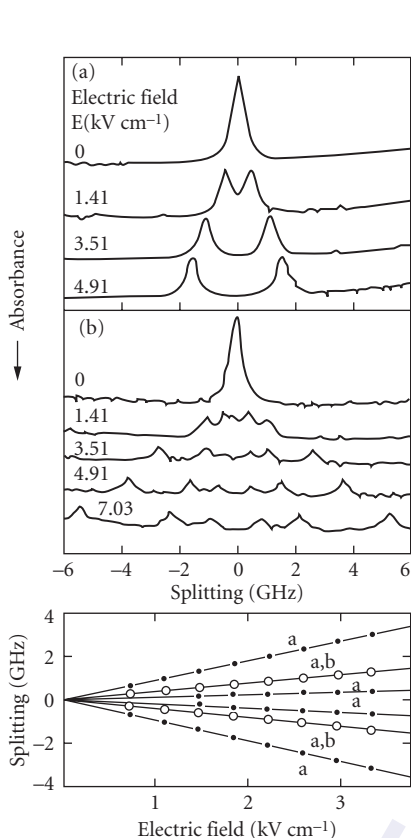


FIGURE 8 Effects of an applied electric field in a hole burned in the 607-nm zero-phonon line observed in irradiated NaF. (After Macfarlane *et al.*²²)

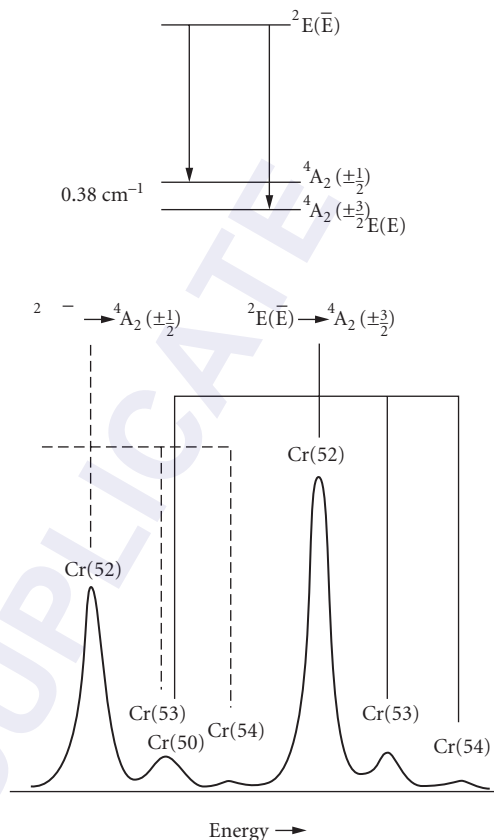


FIGURE 9 Fine structure splitting in the 4A_2 ground state and the isotope shifts of Cr^{3+} in the R_1 -line of ruby measured using FLN spectroscopy. (After Jessop and Szabo.³⁶)

The FLN technique may also be used to measure the effects of phonon-induced relaxation processes and isotope shifts. Isotope and thermal shifts have been reported for $Cr^{3+}:Al_2O_3$ ³⁶ and $Nd^{3+}:LaCl_3$.³⁷ The example given in Fig. 9 shows both the splitting in the ground 4A_2 state of Cr^{3+} in ruby and the shift between lines due to the Cr(50), Cr(52), Cr(53), and Cr(54) isotopes. The measured differential isotope shift of 0.12 cm^{-1} is very close to the theoretical estimate.¹⁹ Superhyperfine effects by the 100 percent abundant Al isotope with $I = 5/2$ also contribute to the homogeneous width of the FLN spectrum of Cr^{3+} in Al_2O_3 (Fig. 12 in Vol. I, Chap. 10).³⁶ Furthermore, in antiferromagnetic oxides such as $GdAlO_3$, $Gd_3Ga_5O_{12}$, and $Gd_3Sc_2Ga_3O_{12}$, spin-spin coupling between the Cr^{3+} ions ($S = 3/2$) and nearest-neighbor Gd^{3+} ions ($S = 3/2$) contributes as much to the zero-phonon R -linewidth as inhomogeneous broadening by strain.³⁸

Configurational Relaxation in Solids

In the case of the broadband ${}^4T_2 \rightarrow {}^4A_2$ transition of Cr^{3+} in YAG (Fig. 12 in Vol. I, Chap. 10) and MgO (Fig. 6), the application of OHB and FLN techniques produce no such narrowing because the

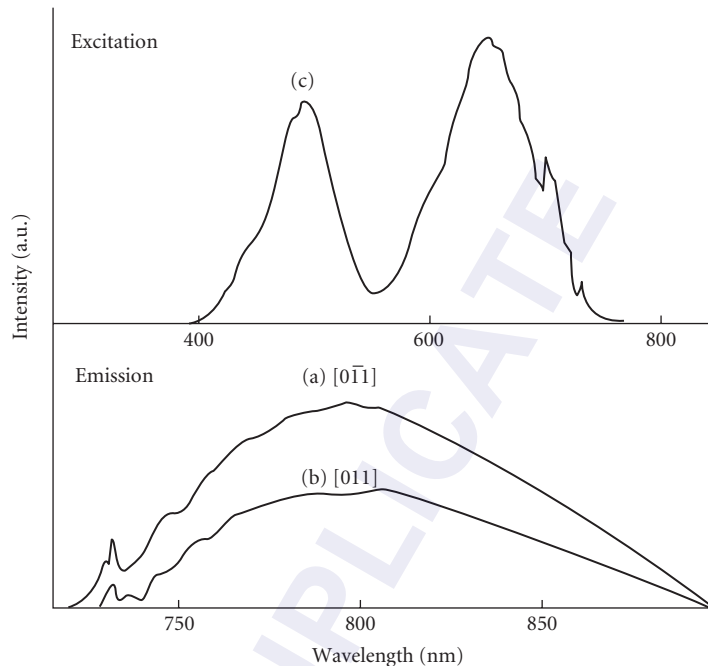


FIGURE 10 Polarized emission of the ${}^4T_2 \rightarrow {}^4A_2$ band from Cr^{3+} ions in orthorhombic sites in MgO . Shown also, (c) is the excitation spectrum appropriate to (a). (After Yamaga *et al.*⁴⁸)

vibronic sideband is the homogeneously broadened shape determined by the phonon lifetime rather than the radiative lifetime. It is noteworthy that the vibronic sideband emission of Cr^{3+} ions in orthorhombic sites in MgO , Fig. 10, shows very little structure. In this case, the Huang-Rhys factor $S \approx 6$, i.e., the strong coupling case, where the multiphonon sidebands tend to lose their separate identities to give a smooth bandshape on the lower-energy side of the peak. By way of contrast, the emission sideband of the R -line transition of Cr^{3+} ions in octahedral sites in MgO is very similar in shape to the known density of one-phonon vibrational modes of MgO ³⁹ (Fig. 11), although there is a difference in the precise positions of the peaks, because the Cr^{3+} ion modifies the lattice vibrations in its neighborhood relative to those of the perfect crystal. Furthermore, there is little evidence in Fig. 11 of higher-order sidebands which justifies treating the MgO R -line process in the weak coupling limit. The absence of such sidebands suggests that $S < 1$, as the discussion in “Optical Spectroscopy and Spectroscopic Lineshapes” (Vol. I, Chap. 10 of this *Handbook*) showed. That the relative intensities of the zero-phonon line and broadband, which should be about e^{-S} , is in the ratio 1:4 shows that the sideband is induced by odd parity phonons. In this case it is partially electric-dipole in character, whereas the zero-phonon line is magnetic-dipole in character.¹⁹

There has been much research on bandshapes of Cr^{3+} -doped spectra in many solids. This is also the situation for F -centers and related defects in the alkali halides. Here, conventional optical spectroscopy has sometimes been supplemented by laser Raman and sub-picosecond relaxation spectroscopies to give deep insights into the dynamics of the optical pumping cycle. The F -center in the alkali halides is a halide vacancy that traps an electron. The states of such a center are reminiscent of a particle in a finite potential well¹⁹ and strong electron-phonon coupling. Huang-Rhys factors in the range $S = 15 - 40$ lead to broad, structureless absorption/luminescence bands with large Stokes

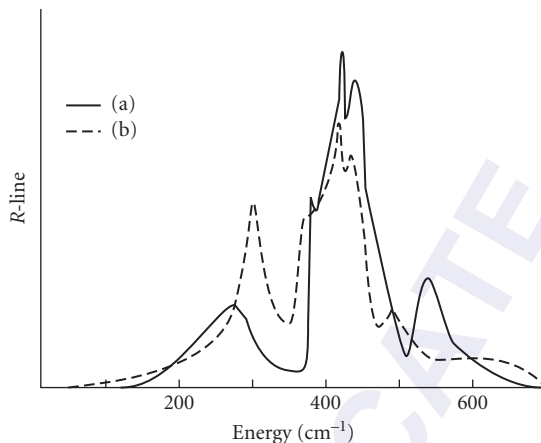


FIGURE 11 A comparison of (a) the vibrational sideband accompanying the ${}^2E \rightarrow {}^4A_2$ R-line of Cr^{3+} ; MgO with (b) the density of phonon modes in MgO as measured by Peckham et al.,³⁹ using neutron scattering. (After Henderson and Imbusch.¹⁹)

shifts (see Fig. 13 in Vol. I, Chap. 10). Raman-scattering measurements on F -centers in NaCl and KCl (Fig. 23 in Vol. I, Chap. 31), showed that the first-order scattering is predominantly due to defect-induced local modes.⁴⁰

The F_A -center is a simple variant on the F -center in which one of the six nearest cation neighbors of the F -center is replaced by an alkali impurity.⁴¹ In KCl the K^+ may be replaced by Na^+ or Li^+ . For the case of the Na^+ substituent, the $F_A(\text{Na})$ center has tetragonal symmetry about a $\langle 100 \rangle$ crystal axis, whereas in the case of Li^+ an off-axis relaxation in the excited state leads to interesting polarized absorption/emission characteristics.^{19,41} The most dramatic effect is the enormous Stokes shift between absorption and emission bands, of order $13,000 \text{ cm}^{-1}$, which has been used to advantage in color center lasers.^{42,43} For $F_A(\text{Li})$ centers, configurational relaxation has been probed using picosecond relaxation and Raman-scattering measurements. Mollenauer et al.⁴⁴ used the experimental system shown in Fig. 10 in Vol. I, Chap. 31 to carry out measurements of the configurational relaxation time of $F_A(\text{Li})$ centers in KCl. During deexcitation many phonons are excited in the localized modes coupled to the electronic states which must be dissipated into the continuum of lattice modes. Measurement of the relaxation time constitutes a probe of possible phonon damping. A mode-locked dye laser producing pulses of 0.7-ps duration at 612 nm was used both to pump the center in the F_{A2} -absorption band and to provide the timing beam. Such pumping leads to optical gain in the luminescence band and prepares the centers in their relaxed state. The probe beam, collinear with the pump beam, is generated by a CW $F_A(\text{Li})$ -center laser operating at 2.62 μm . The probe beam and gated pulses from the dye laser are mixed in a non-linear optical crystal (lithium iodate). A filter allows only the sum frequency at 496 nm to be detected. The photomultiplier tube then measures the rise in intensity of the probe beam which signals the appearance of gain where the $F_A(\text{Li})$ -centers have reached the relaxed excited state. The pump beam is chopped at low frequency to permit phase-sensitive detection. The temporal evolution of $F_A(\text{Li})$ -center gains (Fig. 12a and b) was measured by varying the time delay between pump and gating pulses. In this figure, the solid line is the instantaneous response of the system, whereas in b the dashed line is the instantaneous response convolved with a 1.0-ps rise time.

Measurements of the temperature dependence of the relaxation times of $F_A(\text{Li})$ -centers in potassium chloride (Fig. 12c) show that the process is very fast, typically of order 10 ps at 4 K. Furthermore, configurational relaxation is a multiphonon process which involves mainly the creation of some 20 low-energy phonons of energy $E_p/hc = 47 \text{ cm}^{-1}$. That only about $(20 \times 47/8066) \text{ eV} = 0.1 \text{ eV}$ deposited into the 47 cm^{-1} mode, whereas 1.6 eV of optical energy is lost to the overall relaxation process, indicates

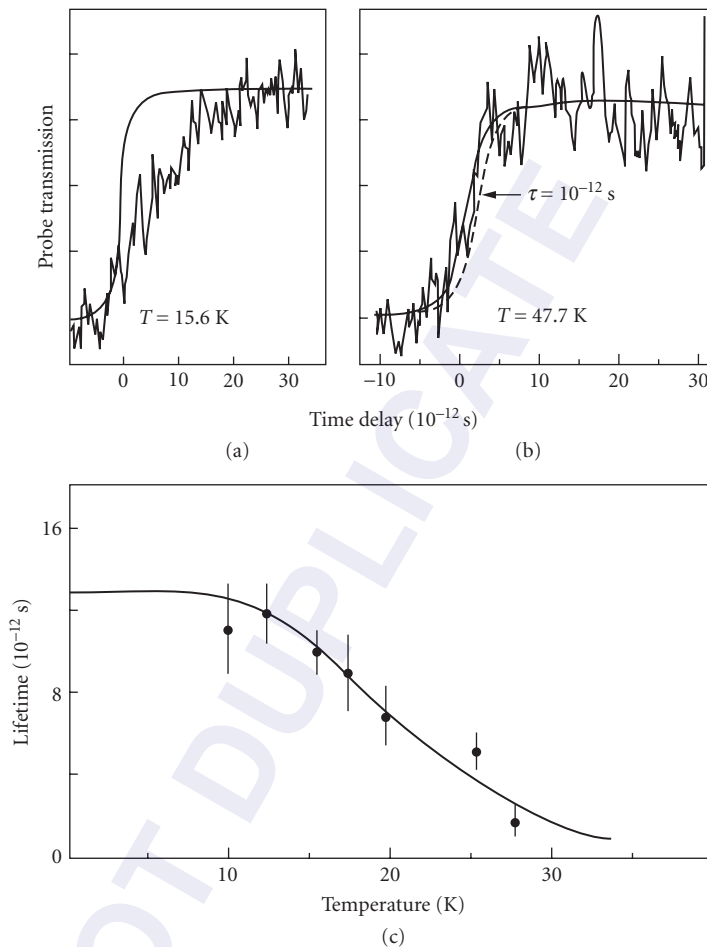


FIGURE 12 (a) Temporal evolution of gain in the $F_A(\text{Li})$ center emission in picosecond pulse probe measurements; (b) the temperature-dependence of the gain process; and (c) temperature dependence of the relaxation time. (After Mollenauer *et al.*⁴⁴)

that other higher-energy modes of vibrations must be involved.⁴⁴ This problem is resolved by Raman-scattering experiments. For $F_A(\text{Li})$ -centers in potassium chloride, three sharp Raman-active local modes were observed with energies of 47 cm^{-1} , 216 cm^{-1} , and 266 cm^{-1} , for the ^7Li isotope.⁴⁵ These results and later polarized absorption/luminescence studies indicated that the Li^+ ion lies in an off-center position in a $\langle 110 \rangle$ crystal direction relative to the z axis of the F_A center. Detailed polarized Raman spectroscopy resonant and nonresonant with the F_A -center absorption bands are shown in Fig. 13.⁴⁶ These spectra show that under resonant excitation in the F_{A1} absorption band, each of the three lines due to the sharp localized modes is present in the spectrum. The polarization dependence confirms that the 266 cm^{-1} mode is due to Li^+ ion motion in the mirror plane and parallel to the defect axis. The 216 cm^{-1} mode is stronger under nonresonant excitation, reflecting the off-axis vibrations of the Li^+ ion vibrating in the mirror plane perpendicular to the z axis. On the other hand, the low-frequency mode is an amplified band mode of the center which hardly involves the motion of the Li^+ ion.

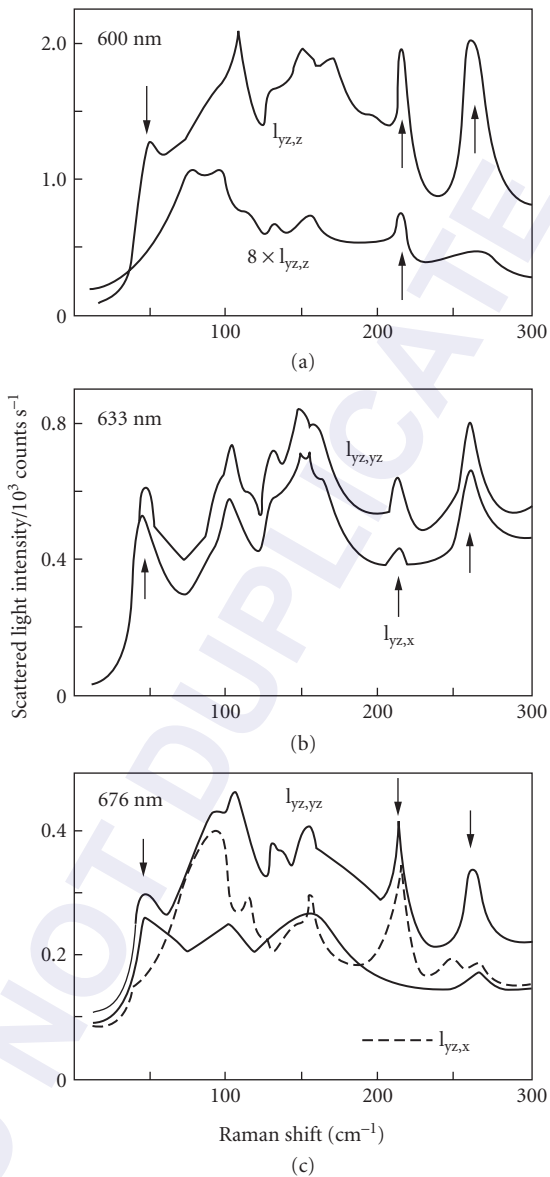


FIGURE 13 Raman spectra of F_A (Li) centers in potassium chloride measured at 10 K for different senses of polarization. In (a) the excitation wavelength $\lambda = 600$ nm is midway between the peaks of the F_{A1} bands; (b) $\lambda = 632.8$ nm is resonant with the F_{A1} band; and (c) $\lambda = 676.4$ nm is nonresonant. (After Joosen *et al.*⁴⁴)

2.5 ABSORPTION, PHOTOLUMINESCENCE, AND RADIATIVE DECAY MEASUREMENTS

The philosophy of solid-state spectroscopy is subtly different from that of atomic and molecular spectroscopies. It is often required not only to determine the nature of the absorbing/emitting species but also the symmetry and structure of the local environment. Also involved is the interaction of the electronic center with other neighboring ions, which leads to lineshape effects as well as time-dependent phenomena. The consequence is that a combination of optical spectroscopic techniques may be used in concert. This general approach to optical spectroscopy of condensed matter phenomena is illustrated by reference to the case of Al_2O_3 and MgO doped with Cr^{3+} .

Absorption and Photoluminescence of Cr^{3+} in Al_2O_3 and MgO

The absorption and luminescence spectra may be interpreted using the Tanabe-Sugano diagram shown in Fig. 5, as discussed previously. Generally, the optical absorption spectrum of Cr^{3+} : Al_2O_3 (Fig. 14) is dominated by broadband transitions from the ${}^4\text{A}_2 \rightarrow {}^4\text{T}_2$ and ${}^4\text{A}_2 \rightarrow {}^4\text{T}_1$. The crystal used in this measurement contained some $10^{18} \text{ Cr}^{3+} \text{ ions cm}^{-3}$. Since the absorption coefficient at the peak of the ${}^4\text{A}_2 \rightarrow {}^4\text{T}_2$ band is only 2 cm^{-1} , it is evident from Eq. (6) in Chap. 31, "Optical Spectrometers," in Vol. I, that the cross section at the band peak is $\sigma_o \cong 5 \times 10^{-19} \text{ cm}^2$. The spin-forbidden absorption transitions ${}^4\text{A}_2 \rightarrow {}^2\text{E}, {}^2\text{T}_1$ are just distinguished as weak absorptions ($\sigma_o \sim 10^{21} \text{ cm}^2$) on the long-wavelength side of the ${}^4\text{A}_2 \rightarrow {}^4\text{T}_2$ band. This analysis strictly applies to the case of octahedral symmetry. Since the cation site in ruby is distorted from perfect octahedral symmetry, there are additional electrostatic energy terms associated with this reduced symmetry. One result of this distortion, as illustrated in Fig. 14, is that the absorption and emission spectra are no longer optically isotropic. By measuring the peak shifts of the ${}^4\text{A}_2 \rightarrow {}^4\text{T}_2$ and ${}^4\text{A}_2 \rightarrow {}^4\text{T}_1$ absorption transitions between π and σ senses of polarization, the trigonal field splittings of the ${}^4\text{T}_2$ and ${}^4\text{T}_1$ levels may be determined.²⁵

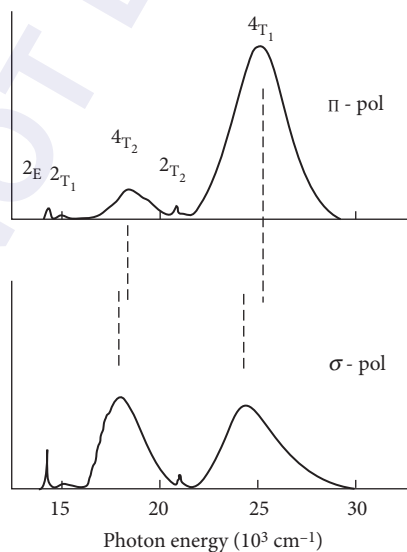


FIGURE 14 Polarized optical absorption spectrum of a ruby $\text{Cr}^{3+}:\text{Al}_2\text{O}_3$ crystal containing $2 \times 10^{18} \text{ Cr}^{3+} \text{ ions cm}^{-3}$ measured at 77 K.

The Cr^{3+} ion enters the MgO substitutionally for the Mg^{2+} ion. The charge imbalance requires that for every two impurity ions there must be one cation vacancy. At low-impurity concentrations, charge-compensating vacancies are mostly remote from the Cr^{3+} ions. However, some 10 to 20 percent of the vacancies occupy sites close to individual Cr^{3+} ions, thereby reducing the local symmetry from octahedral to tetragonal or orthorhombic.¹⁹ The optical absorption spectrum of $\text{Cr}^{3+}:\text{MgO}$ is also dominated by broadband ${}^4\text{A}_2 \rightarrow {}^4\text{T}_2, {}^4\text{T}_1$ transitions; in this case, there are overlapping contributions from Cr^{3+} ions in three different sites. There are substantial differences between the luminescence spectra of Cr^{3+} in the three different sites in MgO (Fig. 15), these overlapping spectra being determined by the ordering of the ${}^4\text{T}_2$ and ${}^2\text{E}$ excited states. For strong crystal fields, $Dq/B > 2.5$, ${}^2\text{E}$ lies lowest and nonradiative decay from ${}^4\text{T}_1$ and ${}^4\text{T}_2$ levels to ${}^2\text{E}$ results in very strong emission in the sharp *R*-lines, with rather weaker vibronic sidebands. This is the situation for Cr^{3+} ions in octahedral and tetragonal sites in MgO.¹⁹ The ${}^2\text{E} \rightarrow {}^4\text{A}_2$ luminescence transition is both spin- and parity-forbidden (see Vol. I, Chap. 10) and this is signaled by relatively long radiative lifetimes—11.3 ms for octahedral sites and 8.5 ms for tetragonal sites at 77 K. This behavior is in contrast to that of Cr^{3+} in orthorhombic sites, for which the ${}^4\text{T}_2$ level lies below the ${}^2\text{E}$ level. The stronger electron-phonon coupling for the ${}^4\text{T}_2 \rightarrow {}^4\text{A}_2$ transition at

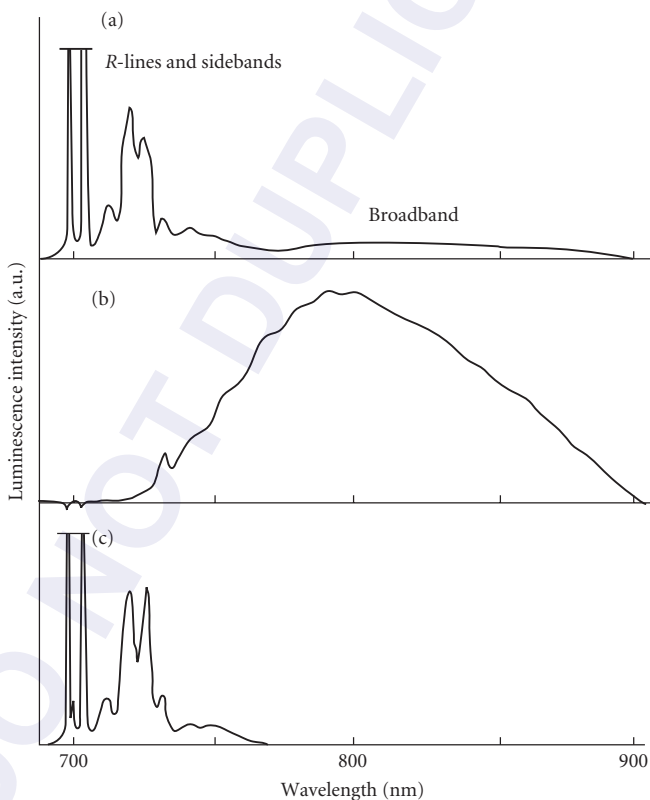


FIGURE 15 Photoluminescence spectra of $\text{Cr}^{3+}:\text{MgO}$ using techniques of phase-sensitive detection. In (a) the most intense features are sharp *R*-lines near 698 to 705 nm due to Cr^{3+} ions at sites with octahedral and tetragonal symmetry; a weak broadband with peak at 740 nm is due to Cr^{3+} ions in sites of orthorhombic symmetry. By adjusting the phase-shift control on the lock-in amplifier (Fig. 4 in Vol. I, Chap. 31), the relative intensities of the three components may be adjusted as in parts (b) and (c).

orthorhombic sites leads to a broadband luminescence with peak at 790 nm. Since this is a spin-allowed transition, the radiative lifetime is much shorter—only 35 μs .⁴⁷

As noted previously the decay time of the luminescence signals from Cr^{3+} ions in octahedral and tetragonal symmetry are quite similar, and good separation of the associated R -lines using the phase-nulling technique are then difficult. However, as Fig. 15 shows, good separation of these signals from the ${}^4\text{T}_2 \rightarrow {}^4\text{A}_2$ broadband is very good. This follows from the applications of Eqs. (8) through (12) in Vol. I, Chap. 31. For Cr^{3+} ions in cubic sites, the long lifetime corresponds to a signal phase angle of 2° : the R -line intensity can be suppressed by adjusting the detector phase angle to $(90^\circ + 2^\circ)$. In contrast, the Cr^{3+} ions in orthorhombic sites give rise to a phase angle of 85° : this signal is reduced to zero when $\phi_D = (90^\circ + 85^\circ)$.

Excitation Spectroscopy

The precise positions of the ${}^4\text{A}_2 \rightarrow {}^4\text{T}_2, {}^4\text{T}_1$ absorption peaks corresponding to the sharp lines and broadbands in Fig. 15 may be determined by excitation spectroscopy (see Vol. I, Chap. 31). An example of the application of this technique is given in Fig. 10, which shows the emission band of the ${}^4\text{T}_2 \rightarrow {}^4\text{A}_2$ transition at centers with orthorhombic symmetry, Figs. 10a and b, and its excitation spectrum, Fig. 10c.⁴⁷ The latter was measured by setting the wavelength of the detection spectrometer at $\lambda = 790$ nm, i.e., the emission band peak, and scattering the excitation monochromator over the wavelength range 350 to 750 nm of the Xe lamp. Figure 10 gives an indication of the power of excitation spectroscopy in uncovering absorption bands not normally detectable under the much stronger absorptions from cubic and tetragonal centers. Another example is given in Fig. 16—in this case, of recombining excitons in the smaller gap material (GaAs) in GaAs/AlGaAs quantum wells. In this case, the exciton luminescence peak energy, $h\nu_x$, is given by

$$h\nu_x = E_G + E_{1e} + E_{1h} + E_b \quad (3)$$

where E_G is the bandgap of GaAs, E_{1e} and E_{1h} are the $n = 1$ state energies of electrons (e) and holes (h) in conduction and valence bands, respectively, and E_b is the electron-hole binding energy. Optical transitions involving electrons and holes in these structures are subject to the $\Delta n=0$ selection rule. In consequence, there is a range of different absorption transitions at energies above the bandgap. Due to the rapid relaxation of energy in levels with $n > 1$, the recombination luminescence occurs between the $n = 1$ electron and hole levels only, in this case at 782 nm. The excitation spectrum in which this luminescence is detected and excitation wavelength varied at wavelengths shorter than 782 nm reveals the presence of absorption transitions above the bandgap. The first absorption transition shown is the $1lh \rightarrow 1e$ transition, which occurs at slightly longer wavelength than the $1hh \rightarrow 1e$ transition. The light hole (lh)-heavy hole (hh) splitting is caused by spin-orbit splitting and strain in these epilayer structures. Other, weaker transitions are also discernible at higher photon energies.

Polarization Spectroscopy

The discussions on optical selection rules, in Vol. 1, Chaps. 10 and 31, showed that when a well-defined axis is presented, the strength of optical transitions may depend strongly on polarization. In atomic physics the physical axis is provided by an applied magnetic field (Zeeman effect) or an applied electric field (Stark effect). Polarization effects in solid-state spectroscopy may be used to give information about the site symmetry of optically active centers. The optical properties of octahedral crystals are normally isotropic. In this situation, the local symmetry of the center must be lower than octahedral so that advantage may be taken of the polarization-sensitivity of the selection rules. Several possibilities exist in noncubic crystals. If the local symmetry of all centers in the crystal point in the same direction, then the crystal as a whole displays an axis of symmetry. Sapphire (Al_2O_3) is an example, in which the Al^{3+} ions occupy trigonally distorted octahedral sites. In consequence, the optical absorption and luminescence spectra of ions in this crystal are naturally polarized. The observed π - and σ -polarized absorption spectra of ruby shown in Fig. 14 are in general agreement with the calculated selection rules,²⁵

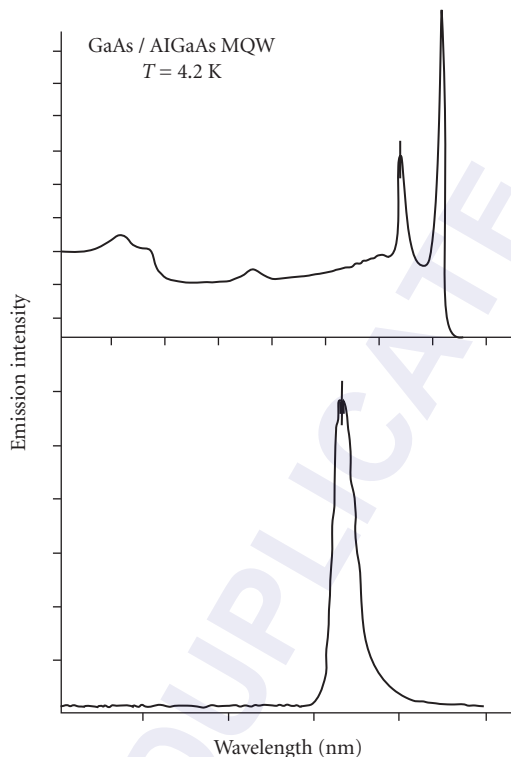


FIGURE 16 Luminescence spectrum and excitation spectrum of multiple quantum wells in GaAs/AlGaAs samples measured at 6 K. (*P. Dawson, 1986 private communication to the author.*)

although there are undoubtedly vibronic processes contributing to these broadband intensities.⁴⁷ The other important ingredient in the spectroscopy of the Cr^{3+} ions in orthorhombic symmetry sites in MgO is that the absorption and luminescence spectra are strongly polarized. It is then quite instructive to indicate how the techniques of polarized absorption/luminescence help to determine the symmetry axes of the dipole transitions. The polarization of the ${}^4T_2 \rightarrow {}^4A_2$ emission transition in Fig. 10 is clear. In measurements employing the “straight-through” geometry, Henry et al.⁴⁷ reported the orientation intensity patterns shown in Fig. 17 for the broadband spectrum. A formal calculation of the selection rules and the orientation dependence of the intensities shows that the intensity at angle θ is given by

$$I(\theta) = (A_\pi - A_\sigma)(E_\pi - E_\sigma) \sin^2\left(\theta + \frac{\pi}{4}\right) + \text{constant} \quad (4)$$

where A and E refer to the absorbed and emitted intensities for π - and σ -polarizations.⁴⁸ The results in Fig. 17 are consistent with the dipoles being aligned along $\langle 110 \rangle$ directions of the octahedral MgO lattice. This is in accord with the model of the structure of the Cr^{3+} ions in orthorhombic symmetry, which locates the vacancy in the nearest neighbor cation site relative to the Cr^{3+} ion along a $\langle 110 \rangle$ direction.

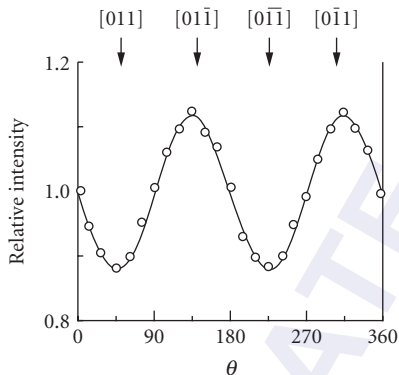


FIGURE 17 The polarization characteristics of the luminescence spectrum of Cr^{3+} ions in orthorhombic sites in $\text{Cr}^{3+}:\text{MgO}$. (After Henry *et al.*⁴⁷)

Zeeman Spectroscopy

The Zeeman effect is the splitting of optical lines by a static magnetic field due to the removal of the spin degeneracy of levels involved in the optical transitions. In many situations the splittings are not much larger than the optical linewidth of zero-phonon lines and much less than the width of vibronically broadened bands. The technique of optically detected magnetic resonance (ODMR) is then used to measure the Zeeman splittings. As we have already shown, ODMR also has the combined ability to link inextricably, an excited-state ESR spectrum with an absorption band and a luminescence band. The spectrum shown in Fig. 16 in Vol. I, Chap. 31 is an example of this unique power, which has been used in such diverse situations as color centers, transition-metal ions, rare-earth ions, phosphor- and laser-active ions (e.g., Ga^+ , Tl°), as well as donor-acceptor and exciton recombination in semiconductors.¹⁹ We now illustrate the relationship of the selection rules and polarization properties of the triplet-singlet transitions.

The F -center in calcium oxide consists of two electrons trapped in the Coulomb field of a negative-ion vacancy. The ground state is a spin singlet, $^1A_{1g}$, from which electric dipole absorption transitions are allowed into a $^1T_{1u}$ state derived from the $(1s2p)$ configuration. Such $^1A_{1g} \rightarrow ^1T_{1u}$ transitions are signified by a strong optical absorption band centered at a wavelength $\lambda \approx 400$ nm (Vol. I, Chap. 31, Fig. 16). Excitation of this $^1T_{1u}$ state does not proceed via $^1T_{1u} \rightarrow ^1A_{1g}$ luminescence. Instead, there is an efficient nonradiative decay from $^1T_{1u}$ into the triplet $^3T_{1u}$ state also derived from the $(1s2p)$ configuration.⁴⁹ The spin-forbidden $^3T_{1u} \rightarrow ^1A_{1g}$ transition gives rise to a striking orange fluorescence, which occurs with a radiative lifetime $\tau_R = 3.4$ ms at 4.2 K. The ODMR spectrum of the F -center and its absorption and emission spectral dependences are depicted in Fig. 16 in Vol. I, Chap. 31, other details are shown in Fig. 18. With the magnetic field at some general orientation in the (100) plane there are six lines. From the variation of the resonant fields with the orientation of the magnetic field in the crystal, Edel *et al.* (1972)⁵⁰ identified the spectrum with the $S = 1$ state of tetragonally distorted F -center. The measured orientation dependence gives $g_{\parallel} \approx g_{\perp} = 1.999$ $D = 60.5$ mT.

Figure 18 shows the selection rules for emission of circularly polarized light by $S = 1$ states in axial crystal fields. We denote the populations of the $M_s = 0, \pm 1$ levels as N_0 and $N_{\pm 1}$. The low-field ESR line, corresponding to the $M_s = 0 \rightarrow M_s = +1$ transition, should be observed as an increase in σ_+ -light because $N_0 > N_{+1}$ and ESR transitions enhance the $M_s = \pm 1$ level. However, the high-field line is observed as a change in intensity of σ_- -light. If spin-lattice relaxation is efficient (i.e., $T_1 < \tau_R$), then the spin states are in thermal equilibrium, $N_0 < N_{-1}$, ESR transitions depopulate the $|M_s = -1\rangle$ level. Thus, the high-field ODMR line is seen as a decrease in the F -center in these crystals (*viz.*, that

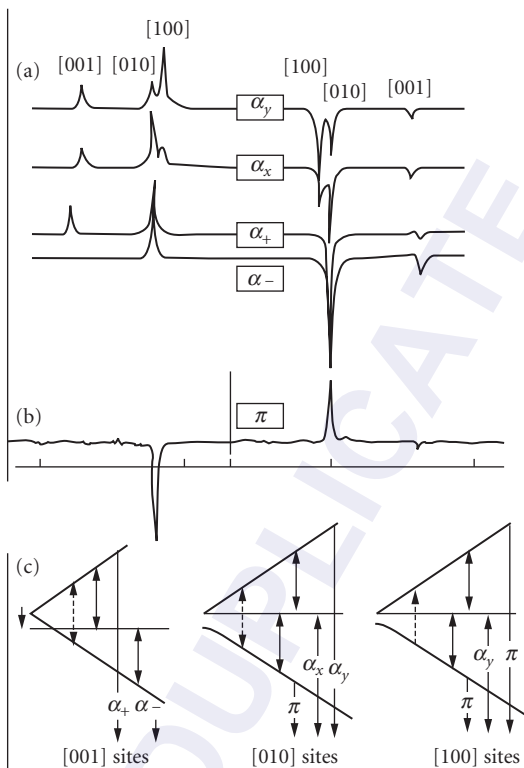


FIGURE 18 Polarization selection rules and the appropriately detected ODMR spectra of the ${}^3T_{1u} \rightarrow {}^1A_{1g}$ transition of F -centers in CaO. (After Edel et al.⁵⁰)

for the lowest 3T_1 state, D is positive and the spin states are in thermal equilibrium). It is worth noting that since the $|M_s = 0\rangle \rightarrow M_s = \pm 1$ ESR transitions occur at different values of the magnetic field, ODMR may be detected simply as a change in the emission intensity at resonance; it is not necessary to measure specifically the sense of polarization of the emitted light. The experimental data clearly establish the tetragonal symmetry of the F -center in calcium oxide: the tetragonal distortion occurs in the excited ${}^3T_{1u}$ state due to vibronic coupling to modes of E_g symmetry resulting in a static Jahn-Teller effect.⁵⁰

2.6 REFERENCES

1. T. W. Hänsch, I. S. Shakin, and A. L. Shawlow, *Nature* (London), **225**:63 (1972).
2. D. N. Stacey, private communication to A. I. Ferguson.
3. E. Giacobino and F. Birabem, *J. Phys.* **B15**:L385 (1982).
4. L. Housek, S. A. Lee, and W. M. Fairbank Jr., *Phys. Rev. Lett.* **50**:328 (1983).
5. P. Zhao, J. R. Lawall, A. W. Kam, M. D. Lindsay, F. M. Pipkin, and W. Lichten, *Phys. Rev. Lett.* **63**:1593 (1989).

6. T. J. Sears, S. C. Foster, and A. R. W. McKellar, *J. Opt. Soc. Am.* **B3**:1037 (1986).
7. C. J. Sansonetti, J. D. Gillaspay, and C. L. Cromer, *Phys. Rev. Lett.* **65**:2539 (1990).
8. W. Lichten, D. Shinen, and Zhi-Xiang Zhou, *Phys. Rev.* **A43**:1663 (1991).
9. C. Adams and A. I. Ferguson, *Opt. Commun.* **75**:419 (1990) and **79**:219 (1990).
10. C. Adams, E. Riis, and A. I. Ferguson, *Phys. Rev. A* (1992) and **A45**:2667 (1992).
11. G. W. F. Drake and A. J. Makowski, *J. Opt. Soc. Amer.* **B5**:2207 (1988).
12. P. G. Harris, H. C. Bryant, A. H. Mohagheghi, R. A. Reeder, H. Sharifian, H. Tootoonchi, C. Y. Tang, J. B. Donahue, C. R. Quick, D. C. Rislove, and W. W. Smith, *Phys. Rev. Lett.* **65**:309 (1990).
13. H. R. Sadeghpour and C. H. Greene, *Phys. Rev. Lett.* **65**:313 (1990).
14. See, for example, H. Friedrich, *Physics World* **5**:32 (1992).
15. Iu et al., *Phys. Rev. Lett.* **66**:145 (1991).
16. B. MacGowan et al., *Phys. Rev. Lett.* **65**:420 (1991).
17. C. Douketic and T. E. Gough, *J. Mol. Spectrosc.* **101**:325 (1983).
18. See, for example, K. Codling et al., *J. Phys.* **B24**:L593 (1991).
19. B. Henderson and G. F. Imbusch, *Optical Spectroscopy of Inorganic Solids*, Clarendon Press, Oxford, 1989.
20. G. Blasse, in B. Di Bartolo (ed.), *Energy Transfer Processes in Condensed Matter*, Plenum Press, New York, 1984.
21. Y. Tanabe and S. Sugano, *J. Phys. Soc. Jap.* **9**:753 (1954).
22. R. M. MacFarlane, J. Y. Wong, and M. D. Sturge, *Phys. Rev.* **166**:250 (1968).
23. B. F. Gachter and J. A. Königstein, *J. Chem. Phys.* **66**:2003 (1974).
24. See O. Deutschbein, *Ann. Phys.* **20**:828 (1932).
25. D. S. McClure, *J. Chem. Phys.* **36**:2757 (1962).
26. After C. J. Delbecq, W. Hayes, M. C. M. O'Brien, and P. H. Yuster, *Proc. Roy. Soc.* **A271**:243 (1963).
27. W. B. Fowler, in Fowler (ed.), *Physics of Color Centers*, Academic Press, New York, 1968. See also G. Boulon, in B. Di Bartolo (ed.), *Spectroscopy of Solid State Laser-Type Materials*, Plenum Press, New York, 1988.
28. Le Si Dang, Y. Merle d'Aubigné, R. Romestain, and A. Fukuda, *Phys. Rev. Lett.* **38**:1539 (1977).
29. A. Ranfagni, D. Mugna, M. Bacci, G. Villiani, and M. P. Fontana, *Adv. in Phys.* **32**:823 (1983).
30. E. Sonder and W. A. Sibley, in J. H. Crawford and L. F. Slifkin (eds.), *Point Defects in Solids*, Plenum Press, New York, vol. 1, 1972.
31. Y. Fang, P. J. Parbrook, B. Henderson, and K. P. O'Donnell, *Appl. Phys. Letts.* **59**:2142 (1991).
32. P. J. Parbrook, B. Cockayne, P. J. Wright, B. Henderson, and K. P. O'Donnell, *Semicond. Sci. Technol.* **6**:812 (1991).
33. T. Kushida and E. Takushi, *Phys. Rev.* **B12**:824 (1975).
34. R. M. Macfarlane, R. T. Harley, and R. M. Shelby, *Radn. Effects* **72**:1 (183).
35. W. Yen and P. M. Selzer, in W. Yen and P. M. Selzer (eds.), *Laser Spectroscopy of Solids*, Springer-Verlag, Berlin, 1981.
36. P. E. Jessop and A. Szabo, *Optics Comm.* **33**:301 (1980).
37. N. Pelletier-Allard and R. Pelletier, *J. Phys.* **C17**:2129 (1984).
38. See Y. Gao, M. Yamaga, B. Henderson, and K. P. O'Donnell, *J. Phys. (Cond. Matter)* (1992) in press (and references therein).
39. G. E. Peckham, *Proc. Phys. Soc. (Lond.)* **90**:657 (1967).
40. J. M. Worlock and S. P. S. Porto, *Phys. Rev. Lett.* **15**:697 (1965).
41. F. Luty, in W. B. Fowler (eds.), *The Physics of Color Centers*, Academic Press, New York, 1968.
42. L. F. Mollenauer and D. H. Olson, *J. App. Phys.* **24**:386 (1974).
43. F. Luty and W. Gellerman, in C. B. Collins (ed.), *Lasers '81*, STS Press, McClean, 1982.
44. L. F. Mollenauer, J. M. Wiesenfeld, and E. P. Ippen, *Radiation Effects* **72**:73 (1983); see also J. M. Wiesenfeld, L. F. Mollenauer, and E. P. Ippen, *Phys. Rev. Lett.* **47**:1668 (1981).

45. B. Fritz, J. Gerlach, and U. Gross, in R. F. Wallis (ed.), *Localised Excitations in Solids*, Plenum Press, New York, 1968, p. 496.
46. W. Joosen, M. Leblans, M. Vahimbeek, M. de Raedt, E. Goovaertz, and D. Schoemaker, *J. Cryst. Def. Amorph. Solids* **16**:341 (1988).
47. M. O. Henry, J. P. Larkin, and G. F. Imbusch, *Phys. Rev.* **B13**:1893 (1976).
48. M. Yamaga, B. Henderson, and K. P. O'Donnell, *J. Luminescence* **43**:139 (1989); see also *ibid.* **46**:397 (1990).
49. B. Henderson, S. E. Stokowski, and T. C. Ensign, *Phys. Rev.* **183**:826 (1969).
50. P. Edel, C. Hennies, Y. Merle d'Aubigné, R. Romestain, and Y. Twarowski, *Phys. Rev. Lett.* **28**:1268 (1972).

DO NOT DUPLICATE

PART

2

ATMOSPHERIC
OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

ATMOSPHERIC OPTICS

Dennis K. Killinger

*Department of Physics
Center for Laser Atmospheric Sensing
University of South Florida
Tampa, Florida*

James H. Churnside

*National Oceanic and Atmospheric Administration
Earth System Research Laboratory
Boulder, Colorado*

Laurence S. Rothman

*Harvard-Smithsonian Center for Astrophysics
Atomic and Molecular Physics Division
Cambridge, Massachusetts*

3.1 GLOSSARY

c	speed of light
C_n^2	atmospheric turbulence strength parameter
D	beam diameter
F	hypergeometric function
$g(\nu)$	optical absorption lineshape function
H	height above sea level
h	Planck's constant
I	Irradiance (intensity) of optical beam (W/m^2)
k	optical wave number
K	turbulent wave number
L	propagation path length
L_0	outer scale size of atmospheric turbulence
l_0	inner scale size of atmospheric turbulence
N	density or concentration of molecules
$p(I)$	probability density function of irradiance fluctuations
P_ν	Planck radiation function
R	gas constant

S	molecular absorption line intensity
T	temperature
v	wind speed
β	backscatter coefficient of the atmosphere
γ_p	pressure-broadened half-width of absorption line
κ	optical attenuation
λ	wavelength
ν	optical frequency (wave numbers)
ρ_0	phase coherence length
σ_I^2	variance of irradiance fluctuations
σ_R	Rayleigh scattering cross section

3.2 INTRODUCTION

Atmospheric optics involves the transmission, absorption, emission, refraction, and reflection of light by the atmosphere and is probably one of the most widely observed of all optical phenomena.¹⁻⁵ The atmosphere interacts with light due to the composition of the atmosphere, which under normal conditions, consists of a variety of different molecular species and small particles like aerosols, water droplets, and ice particles. This interaction of the atmosphere with light is observed to produce a wide variety of optical phenomena including the blue color of the sky, the red sunset, the optical absorption of specific wavelengths due to atmospheric molecules, the twinkling of stars at night, the greenish tint sometimes observed during a severe storm due to the high density of particles in the atmosphere, and is critical in determining the balance between incoming sunlight and outgoing infrared (IR) radiation and thus influencing the earth's climate.

One of the most basic optical phenomena of the atmosphere is the absorption of light. This absorption process can be depicted as in Fig. 1 which shows the transmission spectrum of the atmosphere as

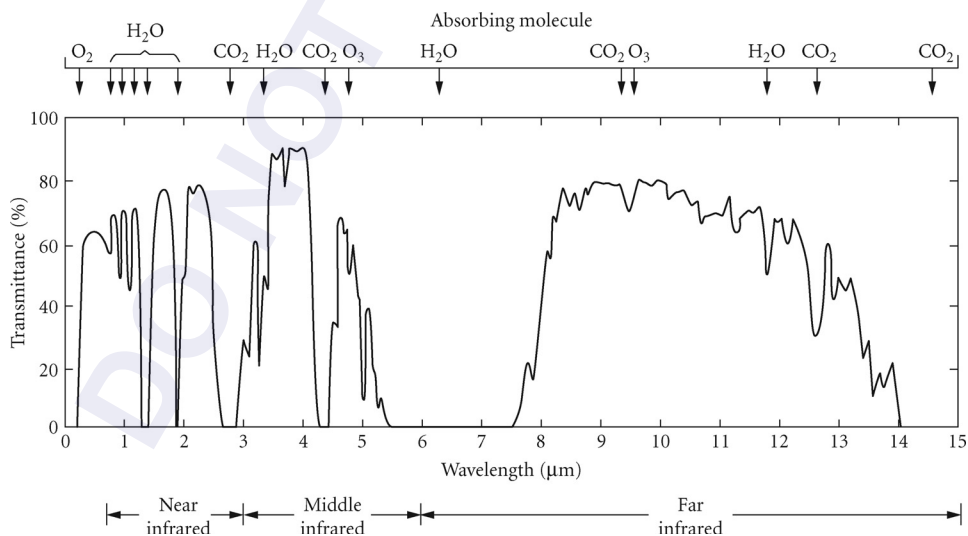


FIGURE 1 Transmittance through the earth's atmosphere as a function of wavelength taken with low spectral resolution (path length 1800 m). (From Measures, Ref. 5.)

a function of wavelength.⁵ The transmission of the atmosphere is highly dependent upon the wavelength of the spectral radiation, and, as will be covered later in this chapter, upon the composition and specific optical properties of the constituents in the atmosphere. The prominent spectral features in the transmission spectrum in Fig. 1 are primarily due to absorption bands and individual absorption lines of the molecular gases in the atmosphere, while a portion of the slowly varying background transmission is due to aerosol extinction and continuum absorption.

This chapter presents a tutorial overview of some of the basic optical properties of the atmosphere, with an emphasis on those properties associated with optical propagation and transmission of light through the earth's atmosphere. The physical phenomena of optical absorption, scattering, emission, and refractive properties of the atmosphere will be covered for optical wavelengths from the ultraviolet (UV) to the far-infrared. The primary focus of this chapter is on *linear* optical properties associated with the transmission of light through the atmosphere. Historically, the study of atmospheric optics has centered on the radiance transfer function of the atmosphere, and the linear transmission spectrum and blackbody emission spectrum of the atmosphere. This emphasis was due to the large body of research associated with passive, electro-optical sensors which primarily use the transmission of ambient optical light or light from selected emission sources. During the past few decades, however, the use of lasers has added a new dimension to the study of atmospheric optics. In this case, not only is one interested in the transmission of light through the atmosphere, but also information regarding the optical properties of the backscattered optical radiation.

In this chapter, the standard linear optical interactions of an optical or laser beam with the atmosphere will be covered, with an emphasis placed on linear absorption and scattering interactions. It should be mentioned that the first edition of the *OSA Handbook of Optics* chapter on "Atmospheric Optics" had considerable nomographs and computational charts to aid the user in numerically calculating the transmission of the atmosphere.² Because of the present availability of a wide range of spectral databases and computer programs (such as the HITRAN Spectroscopy Database, LOWTRAN, MODTRAN, and FASCODE atmospheric transmission computer programs) that model and calculate the transmission of light through the atmosphere, these nomographs, while still useful, are not as vital. As a result, the emphasis on this edition of the "Atmospheric Optics" chapter is on the basic theory of the optical interactions, how this theory is used to model the optics of the atmosphere, the use of available computer programs and databases to calculate the optical properties of the atmosphere, and examples of instruments and meteorological phenomena related to optical or visual remote sensing of the atmosphere.

The overall organization of this chapter begins with a description of the natural, homogeneous atmosphere and the representation of its physical and chemical composition as a function of altitude. A brief survey is then made of the major linear optical interactions that can occur between a propagating optical beam and the naturally occurring constituents in the atmosphere. The next section covers several major computational programs (HITRAN, LOWTRAN, MODTRAN, and FASCODE) and U.S. Standard Atmospheric Models which are used to compute the optical transmission, scattering, and absorption properties of the atmosphere. The next major technical section presents an overview of the influence of atmospheric refractive turbulence on the statistical propagation of an optical beam or wavefront through the atmosphere. Finally, the last few sections of the chapter include a brief introduction to some optical and laser remote sensing experiments of the atmosphere, a brief introduction to the visually important field of meteorological optics, and references to the critical influence of atmospheric optics on global climate change.

It should be noted that the material contained within this chapter has been compiled from several recent overview/summary publications on the optical transmission and atmospheric composition of the atmosphere, as well as from a large number of technical reports and journal publications. These major overview references are (1) *Atmospheric Radiation*, (2) the previous edition of the *OSA Handbook of Optics* (chapter on "Optical Properties of the Atmosphere"), (3) *Handbook of Geophysics and the Space Environment* (chapter on "Optical and Infrared Properties of the Atmosphere"), (4) *The Infrared Handbook*, and (5) *Laser Remote Sensing*.¹⁻⁵ The interested reader is directed toward these comprehensive treatments as well as to the listed references therein for detailed information concerning the topics covered in this brief overview of atmospheric optics.

3.3 PHYSICAL AND CHEMICAL COMPOSITION OF THE STANDARD ATMOSPHERE

The atmosphere is a fluid composed of gases and particles whose physical and chemical properties vary as a function of time, altitude, and geographical location. Although these properties can be highly dependent upon local and regional conditions, many of the optical properties of the atmosphere can be described to an adequate level by looking at the composition of what one normally calls a *standard atmosphere*. This section will describe the background, homogeneous standard composition of the atmosphere. This will serve as a basis for the determination of the quantitative interaction of the molecular gases and particles in the atmosphere with a propagating optical wavefront.

Molecular Gas Concentration, Pressure, and Temperature

The majority of the atmosphere is composed of lightweight molecular gases. Table 1 lists the major gases and trace species of the terrestrial atmosphere, and their approximate concentration (volume fraction) at standard room temperature (296 K), altitude at sea level, and total pressure of 1 atm.⁶ The major optically active molecular constituents of the atmosphere are N_2 , O_2 , H_2O , and CO_2 , with a secondary grouping of CH_4 , N_2O , CO , and O_3 . The other species in the table are present in the atmosphere at trace-level concentrations (ppb, down to less than ppt by volume); however, the concentration may be increased by many orders of magnitude due to local emission sources of these gases.

The temperature of the atmosphere varies both with seasonal changes and altitude. Figure 2 shows the average temperature profile of the atmosphere as a function of altitude presented for the U.S. Standard Atmosphere.⁷⁻⁹ The temperature decreases significantly with altitude until the level of the stratosphere is reached where the temperature profile has an inflection point. The U.S. Standard Atmosphere is one of six basic atmospheric models developed by the U.S. government; these different models furnish a good representation of the different atmospheric conditions which are often encountered. Figure 3 shows the temperature profile for the six atmospheric models.⁷⁻⁹

The pressure of the atmosphere decreases with altitude due to the gravitational pull of the earth and the hydrostatic equilibrium pressure of the atmospheric fluid. This is indicated in Fig. 4 which shows the total pressure of the atmosphere in millibars (1013 mb = 1 atm = 760 torr) as a function of altitude for the different atmospheric models.⁷⁻⁹ The fractional or partial pressure of most of the major gases (N_2 , O_2 , CO_2 , N_2O , CO , and CH_4) follows this profile and these gases are considered uniformly mixed. However, the concentration of water vapor is very temperature-dependent due to freezing and is not uniformly mixed in the atmosphere. Figure 5a shows the density of water vapor as a function of altitude; the units of density are in molecules/cm³ and are related to 1 atm by the appropriate value of Loschmidts number (the number of molecules in 1 cm³ of air) at a temperature of 296 K, which is 2.479×10^{19} molecules/cm³.⁷⁻⁹

The partial pressure of ozone (O_3) also varies significantly with altitude because it is generated in the upper altitudes and near ground level by solar radiation, and is in chemical equilibrium with other gases in the atmosphere which themselves vary with altitude and time of day. Figure 5b shows the typical concentration of ozone as a function of altitude.⁷⁻⁹ The ozone concentration peaks at an altitude of approximately 20 km and is one of the principle molecular optical absorbers in the atmosphere at that altitude. Further details of these atmospheric models under different atmospheric conditions are contained within the listed references and the reader is encouraged to consult these references for more detailed information.^{3,7-9}

Aerosols, Water Droplets, and Ice Particles

The atmospheric propagation of optical radiation is influenced by particulate matter suspended in the air such as aerosols (e.g., dust, haze) and water (e.g., ice or liquid cloud droplets, precipitation). Figure 6 shows the basic characteristics of particles in the atmosphere as a function of altitude,³ and Fig. 7 indicates the approximate size of common atmospheric particles.⁵

TABLE 1 List of Molecular Gases and Their Typical Concentration (Volume Fraction) for the Ambient U.S. Standard Atmosphere

Molecule	Concentration (Volume Fraction)
N ₂	0.781
O ₂	0.209
H ₂ O	0.0775 (variable)
CO ₂	3.3 × 10 ⁻⁴ (higher now: 3.9 × 10 ⁻⁴ , or 390 ppm)
A (argon)	0.0093
CH ₄	1.7 × 10 ⁻⁶
N ₂ O	3.2 × 10 ⁻⁷
CO	1.5 × 10 ⁻⁷
O ₃	2.66 × 10 ⁻⁸ (variable)
H ₂ CO	2.4 × 10 ⁻⁹
C ₂ H ₆	2 × 10 ⁻⁹
HCl	1 × 10 ⁻⁹
CH ₃ Cl	7 × 10 ⁻¹⁰
OCS	6 × 10 ⁻¹⁰
C ₂ H ₂	3 × 10 ⁻¹⁰
SO ₂	3 × 10 ⁻¹⁰
NO	3 × 10 ⁻¹⁰
H ₂ O ₂	2 × 10 ⁻¹⁰
HCN	1.7 × 10 ⁻¹⁰
HNO ₃	5 × 10 ⁻¹¹
NH ₃	5 × 10 ⁻¹¹
NO ₂	2.3 × 10 ⁻¹¹
HOCl	7.7 × 10 ⁻¹²
HI	3 × 10 ⁻¹²
HBr	1.7 × 10 ⁻¹²
OH	4.4 × 10 ⁻¹⁴
HF	1 × 10 ⁻¹⁴
ClO	1 × 10 ⁻¹⁴
HCOOH	1 × 10 ⁻¹⁴
COF ₂	1 × 10 ⁻¹⁴
SF ₆	1 × 10 ⁻¹⁴
H ₂ S	1 × 10 ⁻¹⁴
PH ₃	1 × 10 ⁻²⁰
HO ₂	Trace
O (atom)	Trace
ClONO ₂	Trace
NO ⁺	Trace
HOBr	Trace
C ₂ H ₄	Trace
CH ₃ OH	Trace
CH ₃ Br	Trace
CH ₃ CN	Trace
CF ₄	Trace

Note: The trace species have concentrations less than 1 × 10⁻⁹, with a value that is variable and often dependent upon local emission sources.

Aerosols in the boundary layer (surface to 1 to 2 km altitude) are locally emitted, wind-driven particulates, and have the greatest variability in composition and concentration. Over land, the aerosols are mostly soil particles, dust, and organic particles from vegetation. Over the oceans, they are mostly sea salt particles. At times, however, long-range global winds are capable of transporting land particulates vast distances across the oceans or continents, especially those particulates associated with dust storms or large biomass fires, so that substantial mixing of the different particulate types may occur.

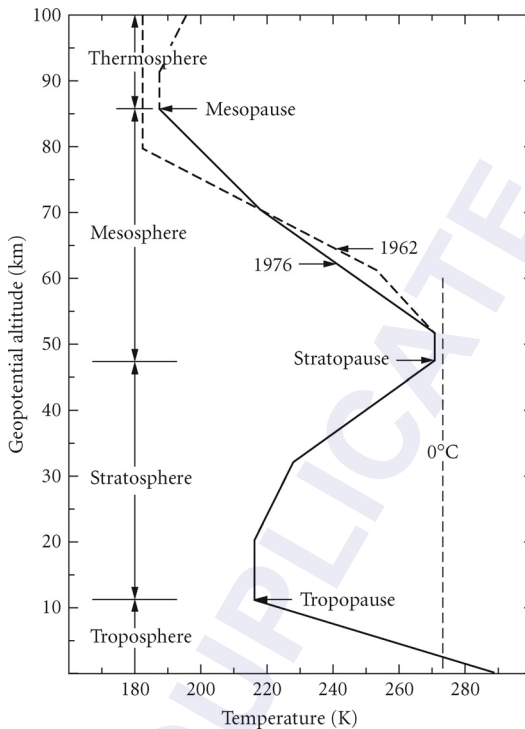


FIGURE 2 Temperature-height profile for U.S. Standard Atmosphere (0 to 86 km).

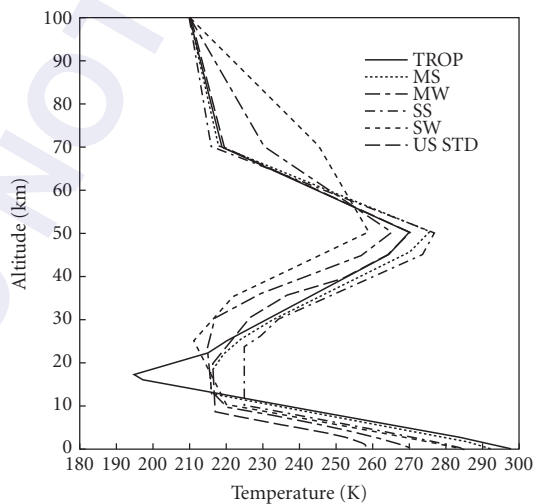


FIGURE 3 Temperature vs. altitude for the six model atmospheres: tropical (TROP), midlatitude summer (MS), midlatitude winter (MW), subarctic summer (SS), subarctic winter (SW), and U.S. standard (US STD).

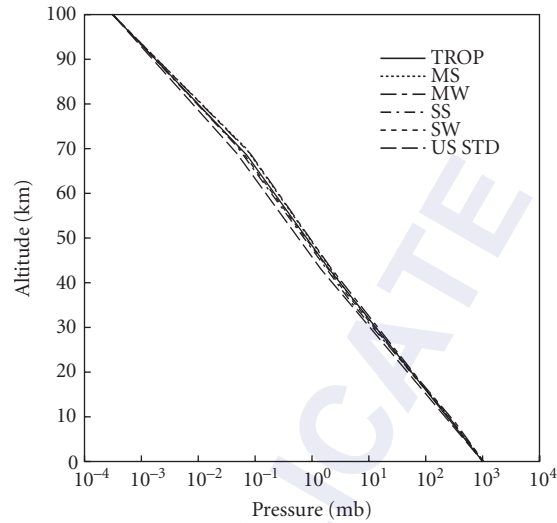


FIGURE 4 Pressure vs. altitude for the six model atmospheres.

In the troposphere above the boundary layer, the composition is less dependent upon local surface conditions and a more uniform, global distribution is observed. The aerosols observed in the troposphere are mostly due to the coagulation of gaseous compounds and fine dust. Above the troposphere, in the region of the stratosphere from 10 to 30 km, the background aerosols are mostly sulfate particles and are uniformly mixed globally. However, the concentration can be perturbed by several orders of magnitude due to the injection of dust and SO_2 by volcanic activity, such as the

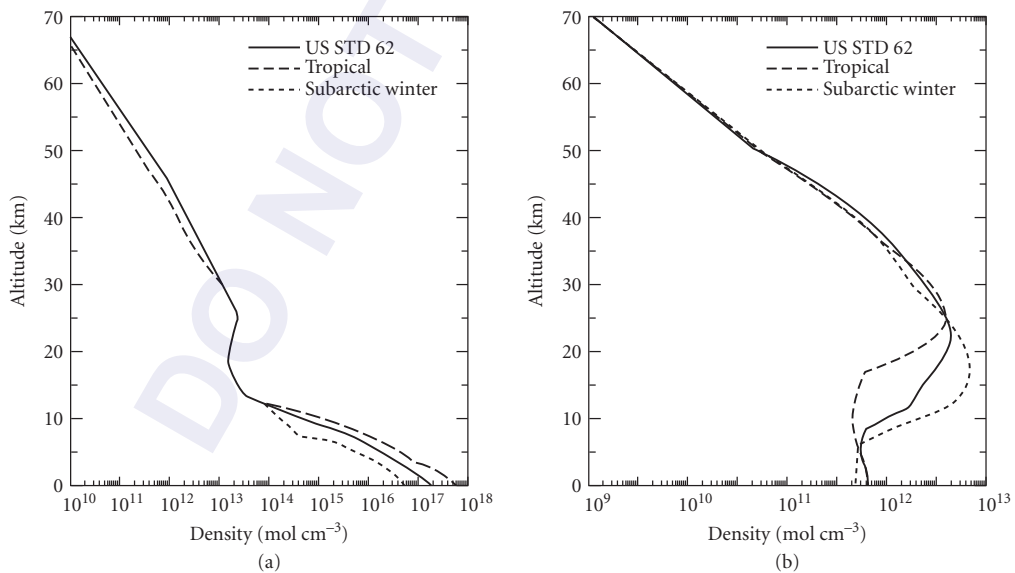


FIGURE 5 (a) Water vapor profile of several models and (b) ozone profile for several models; the U.S. standard model shown is the 1962 model.

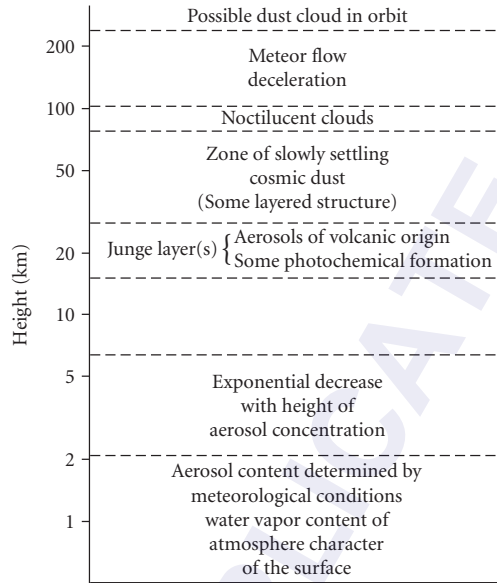


FIGURE 6 Physical characteristics of atmospheric aerosols.

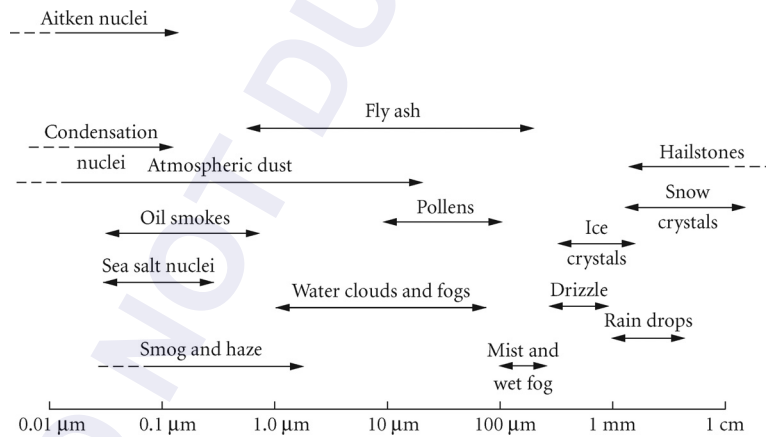


FIGURE 7 Representative diameters of common atmospheric particles. (From Measures, Ref. 5.)

recent eruption of Mt. Pinatubo.¹⁰ Such increases in the aerosol concentration may persist for several years and significantly impact the global temperature of the earth.

Several models have been developed for the number density and size distribution of aerosols in the atmosphere.⁷⁻⁹ Figures 8 and 9 show two aerosol distribution models appropriate for the rural environment and maritime environment, as a function of relative humidity;⁷⁻⁹ the humidity influences the size distribution of the aerosol particles and their growth characteristics. The greatest number density (particles/cm³) occurs near a size of 0.01 μm but a significant number of aerosols are still present even at the larger sizes near 1 to 2 μm. Finally, the optical characteristics of the aerosols can also be dependent upon water vapor concentration, with changes in surface, size, and growth

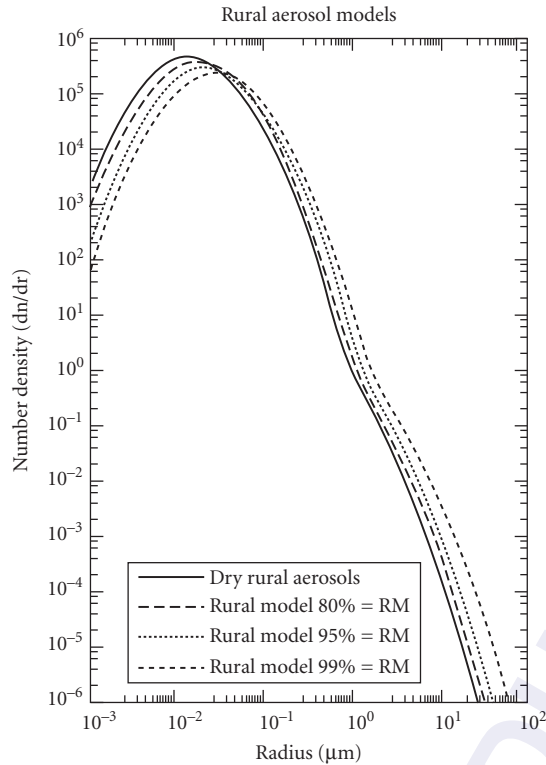


FIGURE 8 Aerosol number density distribution ($\text{cm}^{-3} \mu\text{m}^{-1}$) for the rural model at different relative humidities with total particle concentrations fixed at $15,000 \text{ cm}^{-3}$.

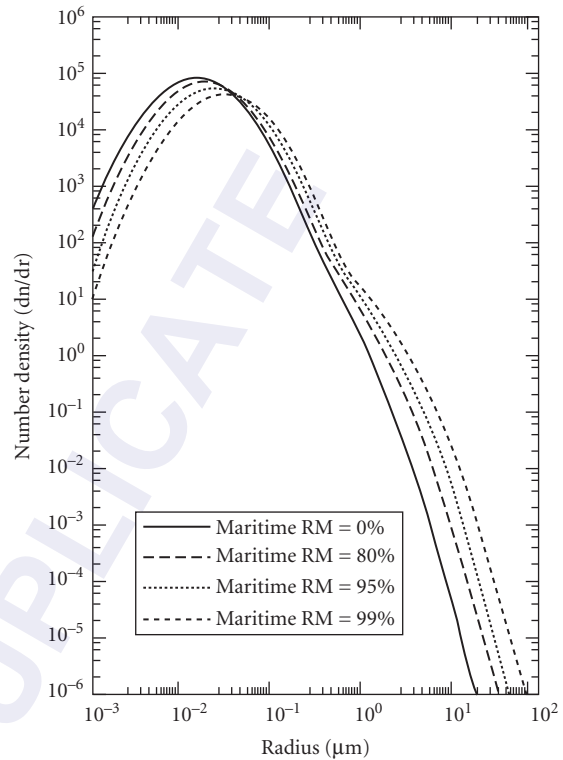


FIGURE 9 Aerosol number density distribution ($\text{cm}^{-3} \mu\text{m}^{-1}$) for the maritime model at different relative humidities with total particle concentrations fixed at 4000 cm^{-3} .

characteristics of the aerosols sometimes observed to be dependent upon the relative humidity. Such humidity changes can also influence the concentration of some pollutant gases (if these gases have been absorbed onto the surface of the aerosol particles).⁷⁻⁹

3.4 FUNDAMENTAL THEORY OF INTERACTION OF LIGHT WITH THE ATMOSPHERE

The propagation of light through the atmosphere depends upon several optical interaction phenomena and the physical composition of the atmosphere. In this section, we consider some of the basic interactions involved in the transmission, absorption, emission, and scattering of light as it passes through the atmosphere. Although all of these interactions can be described as part of an overall radiative transfer process, it is common to separate the interactions into distinct optical phenomena of molecular absorption, Rayleigh scattering, Mie or aerosol scattering, and molecular emission. Each of these basic phenomena is discussed in this section following a brief outline of the fundamental equations for the transmission of light in the atmosphere centered on the Beer-Lambert law.^{1,2}

The linear transmission (or absorption) of monochromatic light by species in the atmosphere may be expressed approximately by the Beer-Lambert law as

$$I(\lambda, t', x) = I(\lambda, t, 0) e^{-\int_0^x \kappa(\lambda) N(x', t) dx'} \quad (1)$$

where $I(\lambda, t', x)$ is the intensity of the optical beam after passing through a path length of x , $\kappa(\lambda)$ is the optical attenuation or extinction coefficient of the species per unit of species density and length, and $N(x, t)$ is the spatial and temporal distribution of the species density that is producing the absorption; λ is the wavelength of the monochromatic light, and the parameter time t' is inserted to remind one of the potential propagation delay. Equation (1) contains the term $N(x, t)$ which explicitly indicates the spatial and temporal variability of the concentration of the attenuating species since in many experimental cases such variability may be a dominant feature.

It is common to write the attenuation coefficient in terms of coefficients that can describe the different phenomena that can cause the extinction of the optical beam. The most dominant interactions in the natural atmosphere are those due to Rayleigh (elastic) scattering, linear absorption, and Mie (aerosol/particulate) scattering; elastic means that the scattered light does not change in wavelength from that which was transmitted while inelastic infers a shift in the wavelength. In this case, one can write $\kappa(\lambda)$ as

$$\kappa(\lambda) = \kappa_a(\lambda) + \kappa_R(\lambda) + \kappa_M(\lambda) \quad (2)$$

where these terms represent the individual contributions due to absorption, Rayleigh scattering, and Mie scattering, respectively. The values for each of these extinction coefficients are described in the following sections along with the appropriate species density term $N(x, t)$. In some of these cases, the reemission of the optical radiation, possibly at a different wavelength, is also of importance. Rayleigh extinction will lead to Rayleigh backscatter, Raman extinction leads to spontaneous Raman scattering, absorption can lead to fluorescence emission or thermal heating of the molecule, and Mie extinction is defined primarily in terms of the scattering coefficient. Under idealized conditions, the scattering processes can be related directly to the value of the attenuation processes. However, if several complex optical processes occur simultaneously, such as in atmospheric propagation, the attenuation and scattering processes are not directly linked via a simple analytical equation. In this case, independent measurements of the scattering coefficient and the extinction coefficient have to be made, or approximation formulas are used to relate the two coefficients.^{4,5}

Molecular Absorption

The absorption of optical radiation by molecules in the atmosphere is primarily associated with individual optical absorption transitions between the allowed quantized energy levels of the molecule. The energy levels of a molecule can usually be separated into those associated with rotational, vibrational, or electronic energy states. Absorption transitions between the rotational levels occur in the far-IR and microwave spectral region, transitions between vibrational levels occur in the near-IR (2 to 20 μm wavelength), and electronic transitions generally occur in the UV-visible region (0.3 to 0.7 μm). Transitions can occur which combine several of these categories, such as rotational-vibrational transitions or electronic-vibrational-rotational transitions.

Some of the most distinctive and identifiable absorption lines of many atmospheric molecules are the rotational-vibrational optical absorption lines in the infrared spectral region. These lines are often clustered together into vibrational bands according to the allowed quantum transitions of the molecule. In many cases, the individual lines are distinct and can be resolved if the spectral resolution of the measuring instrument is fine enough (i.e., $<0.1 \text{ cm}^{-1}$). An example of such a region is the absorption feature near 2.04 μm in Fig. 1 which is actually composed of individual absorption lines if viewed under higher spectral resolution. Figure 10 is a computed high-resolution expansion of the atmospheric spectrum of Fig. 1 near 2.04 μm over a path length of 1800 m which shows these individual lines. In this case, the individual lines are well-separated and appear like a “picket fence” spectrum showing gaps between the absorption lines. Many of the atmospheric gaseous molecules listed earlier in Table 1 have similar spectral structure. These gases are relatively lightweight and have few (less than 5 or 6) atoms per molecule so that their moments of inertia are relatively small. The resulting energy spacing between the allowed rotational-vibrational absorption transitions is large and well-separated in wavelength.

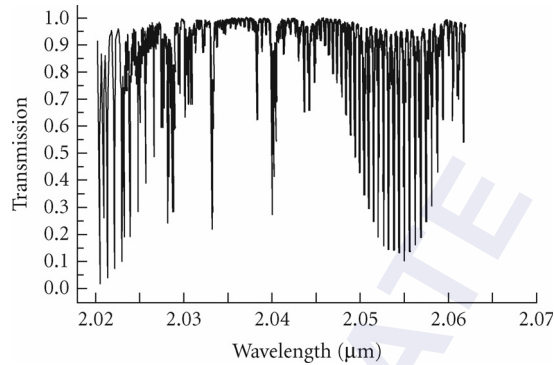


FIGURE 10 High-resolution transmission spectrum of the atmosphere for a horizontal path of 1800 m (similar to that in Fig. 1) for the spectral region near 2.04 μm . The individual rotational absorption lines due to CO_2 and H_2O are easily observed.

In other spectral regions, however, the individual lines overlap or are so strong or saturated that the transmission spectrum displays only broad spectral features; an example of such a spectral region is the strong absorption seen near 2.7 μm , 5 to 8 μm , and beyond 13 μm in Fig. 1. Finally, the molecular absorption observed in the UV is often due to optical transitions to an electronic energy level, molecular energy continuum or a predissociation energy level. In some cases, such as that for O_3 or SO_2 , this results in broad absorption bands that extend throughout the UV region (250 to 350 nm).

More complex, heavier molecules, such as benzene or chlorofluorocarbons, have absorption spectra which are blended or merged together into band spectra due to the complexity and overlap of the rotational-vibrational transitions of these molecules. Figure 11 shows a transmission spectrum of Freon-12 (CCl_2F_2) which has a complex spectrum near 11 μm . As seen, the individual rotational lines are merged into a band spectrum. The band spectrum is unique for each gas and can be used to identify the chemical composition of the gas. Heavy molecules in the atmosphere are not normally part of the natural atmosphere and are usually the result of pollution or gaseous plumes injected into the atmosphere.

The overall transmission or absorption of the atmosphere due to an individual molecular absorption line can be given quantitatively as^{1,2}

$$\kappa_a(\lambda)N(x, t) = Sg(\nu - \nu_0)NP_a \quad (3)$$

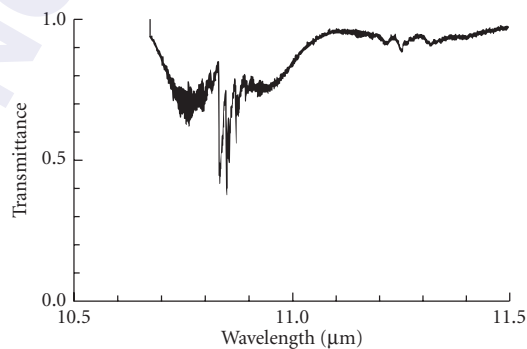


FIGURE 11 High-resolution absorption spectrum of Freon-12 gas showing complex band structure typical of heavy, complex molecules in the atmosphere (path of 100 m with 1 ppm concentration).

where S is the molecular transition line intensity (units of $\text{cm}/\text{molecule}$), $g(\nu - \nu_0)$ is the normalized lineshape function (units of cm or $1/\text{cm}^{-1}$), N is the number of molecules of absorbing species per cm^3 per atm, and P_a is the partial pressure of the absorbing gas in atm. The value of N is equal to the value of Loschmidt's number, which is 2.479×10^{19} molecules cm^{-3} atm $^{-1}$ at a temperature of 296 K; the value of N is inversely proportional to temperature due to the change in gas concentration as a function of temperature for 1 atm of pressure. As will be seen later, the definition of S as given in Eq. (3) is that used in the HITRAN database.¹¹⁻¹³ In this case, S contains the Boltzmann population factor and isotope fraction (natural abundance) as well as the stimulated emission term due to finite population in the upper energy states of the molecule.^{11,12} In Eq. (3), $Sg(\nu - \nu_0)$ is the absorption cross section per molecule ($\text{cm}^2/\text{molecule}$) and NP_a is the number of absorbing molecules in units of molecules/ cm^3 .

The lineshape function can be described by several different models. The two most prevalent are the Lorentzian lineshape associated with pressure broadening and the Gaussian lineshape associated with Doppler broadening which becomes important at elevated temperatures or low pressures.

The Lorentzian/pressure-broadened profile is given by

$$g_L(\nu - \nu_0) = (\gamma_p/\pi) / \left[(\nu - \nu_0)^2 + \gamma_p^2 \right] \quad (4)$$

where γ_p is the pressure-broadened half-width at half-maximum (HWHM) in wave numbers (cm^{-1}). The pressure-broadened half-width is obtained from the air-broadened half-width parameter g as $\gamma_p = gP_t$, where P_t is the total background atmospheric pressure. Under ambient atmospheric conditions, g is approximately 0.05 cm^{-1} (i.e., 1.5 GHz) for many molecules in the atmosphere.

It should be noted that under very low pressure conditions, where the time between collisions with other molecules is relatively long, the intrinsic radiative lifetime of the molecule will determine the lineshape profile. Under these conditions, the linewidth is called the natural linewidth. The natural linewidth of many molecules is on the order of a few MHz (i.e., approximately 0.0001 cm^{-1}) or less.

The Gaussian or Doppler line profile is expressed as

$$g_D(\nu - \nu_0) = (1/\gamma_D) (\ln 2/\pi)^{1/2} \left[-\ln 2(\nu - \nu_0)^2/\gamma_D^2 \right] \quad (5)$$

where γ_D is the Doppler linewidth (HWHM in cm^{-1}) given by

$$\gamma_D = (\nu_0/c) [2RT \ln 2/M]^{1/2} \quad (6)$$

where R is the gas constant, T is the temperature in Kelvin, and M is the molecular weight of the molecule.

The value for the lineshape at the peak (line center) is equal to $1/(\pi\gamma_p) = 0.318/\gamma_p$ for the pressure-broadened case. For the Doppler peak, the maximum value is $(\ln 2/\pi)^{1/2}/\gamma_D = 0.469/\gamma_D$. Under ambient atmospheric conditions, the Doppler linewidth is usually much smaller than the pressure-broadened linewidth.

For those cases where both Lorentzian and Doppler broadening are present in approximately equal amounts, a convolution of the Doppler and Lorentzian profile must be used. This convolution of a Doppler and Lorentzian is called a Voigt profile and involves a double integral for an exact calculation. Fortunately, several numerical approximations are available for the computation of the Voigt profile and lineshape parameters.¹⁴⁻¹⁶ The Voigt profile is important in the spectroscopy of molecules in the upper atmosphere where the ambient pressure is low and the Doppler and pressure-broadened linewidths are of the same order of magnitude. Recent advances in remote-sensing experiments have suggested that further refinements to the description of the line shape in the atmosphere are required. These phenomena include line coupling, speed-dependent corrections, collision-induced narrowing, and other effects on the lineshape.¹⁷

Finally, the large number of transition lines of water vapor and other gases in the atmosphere can produce a significant level of background "quasi-continuum" absorption in the atmosphere. This phenomenon is primarily due to the additive contribution from the wings of the absorption lines even at wavelengths far removed ($>25 \text{ cm}^{-1}$) from the line centers. Such an effect has been studied by Burch and by Clough et al. for water vapor due to strong self-broadening interactions.¹⁸ Figure 12

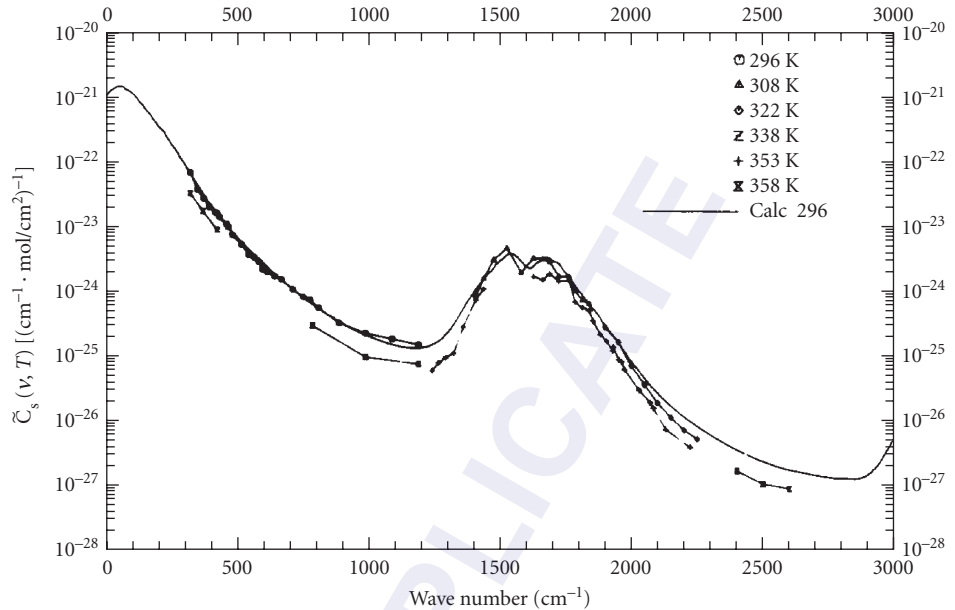


FIGURE 12 Self-density absorption continuum values C_s for water vapor as a function of wave number. The experimental values were measured by Burch. (From Ref. 3.)

shows a plot of the relative continuum coefficient for water vapor as a function of wave number. Good agreement with the experimental data and model calculations is shown. Models for water vapor and nitrogen continuum absorption are contained within many of the major atmospheric transmission programs (such as FASCODE). The typical value for the continuum absorption is negligible in the visible to the near-IR, but can be significant at wavelengths in the range of 5 to 20 μm .

Molecular Rayleigh Scattering

Rayleigh scattering is elastic scattering of the optical radiation due to the displacement of the weakly bound electronic cloud surrounding the gaseous molecule which is perturbed by the incoming electromagnetic (optical) field. This phenomenon is associated with optical scattering where the wavelength of light is much larger than the physical size of the scatterers (i.e., atmospheric molecules). Rayleigh scattering, which makes the color of the sky blue and the setting or rising sun red, was first described by Lord Rayleigh in 1871. The Rayleigh differential scattering cross section for polarized, monochromatic light is given by⁵

$$d\sigma_R/d\Omega = [\pi^2(n^2 - 1^2)/N^2\lambda^4][\cos^2\phi\cos^2\theta + \sin^2\phi] \quad (7)$$

where n is the index of refraction of the atmosphere, N is the density of molecules, λ is the wavelength of the optical radiation, and ϕ and θ are the spherical coordinate angles of the scattered polarized light referenced to the direction of the incident light. As seen from Eq. (7), shorter-wavelength light (i.e., blue) is more strongly scattered out from a propagating beam than the longer wavelengths (i.e., red), which is consistent with the preceding comments regarding the color of the sky or the sunset. A typical value for $d\sigma_R/d\Omega$, at a wavelength of 700 nm in the atmosphere (STP) is approximately $2 \times 10^{-28} \text{ cm}^2 \text{ sr}^{-1}$.³ This value depends upon the molecule and has been tabulated for many of the major gases in the atmosphere.¹⁴⁻¹⁹

The total Rayleigh scattering cross section can be determined from Eq. (7) by integrating over 4π steradians to yield

$$\sigma_R(\text{total}) = [8/3][\pi^2(n^2 - 1)^2/N^2\lambda^4] \quad (8)$$

At sea level (and room temperature, $T = 296$ K) where $N = 2.5 \times 10^{19}$ molecules/cm³, Eq. (8) can be multiplied by N to yield the total Rayleigh scattering extinction coefficient as

$$\kappa_R(\lambda)N(x, t) = N\sigma_R(\text{total}) = 1.18 \times 10^{-8} [550 \text{ nm}/\lambda(\text{nm})]^4 \text{ cm}^{-1} \quad (9)$$

The neglect of the effect of dispersion of the atmosphere (variation of the index of refraction n with wavelength) results in an error of less than 3 percent in Eq. (9) in the visible wavelength range.⁵

The molecular Rayleigh backscatter ($\theta = \pi$) cross section for the atmosphere has been given by Collins and Russell for polarized incident light (and received scattered light of the same polarization) as¹⁹

$$\sigma_R = 5.45 \times 10^{-28} [550 \text{ nm}/\lambda(\text{nm})]^4 \text{ cm}^2 \text{ sr}^{-1} \quad (10)$$

At sea level where $N = 2.47 \times 10^{19}$ molecules/cm³, the atmospheric volume backscatter coefficient, β_R , is thus given by

$$\beta_R = N\sigma_R = 1.39 \times 10^{-8} [550 \text{ nm}/\lambda(\text{nm})]^4 \text{ cm}^{-1} \text{ sr}^{-1} \quad (11)$$

The backscatter coefficient for the reflectivity of a laser beam due to Rayleigh backscatter is determined by multiplying β_R by the range resolution or length of the optical interaction being considered.

For unpolarized incident light, the Rayleigh scattered light has a depolarization factor δ which is the ratio of the two orthogonal polarized backscatter intensities. δ is usually defined as the ratio of the perpendicular and parallel polarization components measured relative to the direction of the incident polarization. Values of δ depend upon the anisotropy of the molecules or scatters, and typical values range from 0.02 to 0.11.²⁰ Depolarization also occurs for multiple scattering and is of considerable interest in laser or optical transmission through dense aerosols or clouds.²¹ The depolarization factor can sometimes be used to determine the physical and chemical composition of the cloud constituents, such as the relative ratio of water vapor or ice crystals in a cloud.

Mie Scattering: Aerosols, Water Droplets, and Ice Particles

Mie scattering is similar to Rayleigh scattering, except the size of the scattering sites is on the same order of magnitude as the wavelength of the incident light, and is, thus, due to aerosols and fine particulates in the atmosphere. The scattered radiation is the same wavelength as the incident light but experiences a more complex functional dependence upon the interplay of the optical wavelength and particle size distribution than that seen for Rayleigh scattering.

In 1908, Mie investigated the scattering of light by dielectric spheres of size comparable to the wavelength of the incident light.²² His analysis indicated the clear asymmetry between the forward and backward directions, where for large particle sizes the forward-directed scattering dominates. Complete treatments of Mie scattering can be found in several excellent works by Deirmendjian and others, which take into account the complex index of refraction and size distribution of the particles.^{23,24} These calculations are also influenced by the asymmetry of the aerosols or particulates which may not be spherical in shape.

The effect of Mie scattering in the atmosphere can be described as in the following figures. Figure 13 shows the aerosol Mie extinction coefficient as a function of wavelength for several atmospheric models, along with a typical Rayleigh scattering curve for comparison.²⁵ Figure 14 shows similar values for the volume Mie backscatter coefficient as a function of wavelength.²⁵ Extinction and backscatter coefficient values are highly dependent upon the wavelength and particulate composition.

Figures 15 and 16 show the calculated extinction coefficient for the rural and maritime aerosol models described in Sec. 3.3 as a function of relative humidity and wavelength.³ Significant changes in the backscatter can be produced by relatively small changes in the humidity.

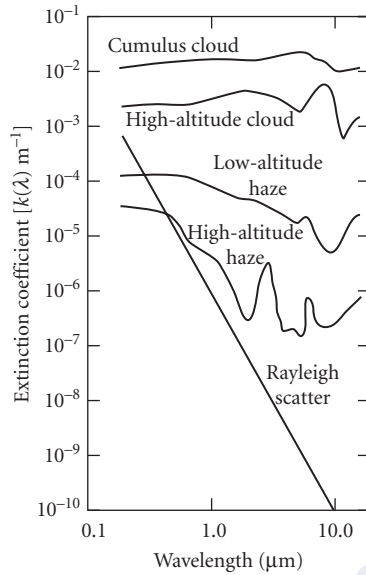


FIGURE 13 Aerosol extinction coefficient as a function of wavelength. (From Measures, Ref. 5.)

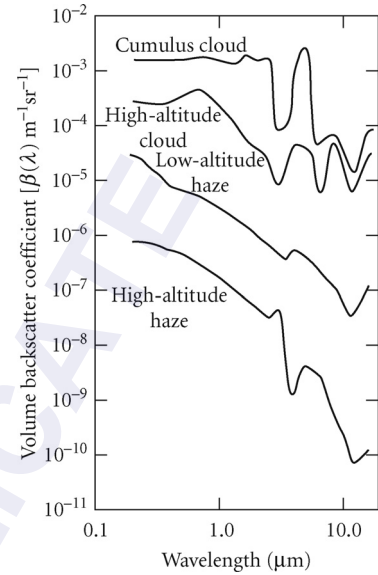


FIGURE 14 Aerosol volume backscattering coefficient as a function of wavelength. (From Measures, Ref. 5.)

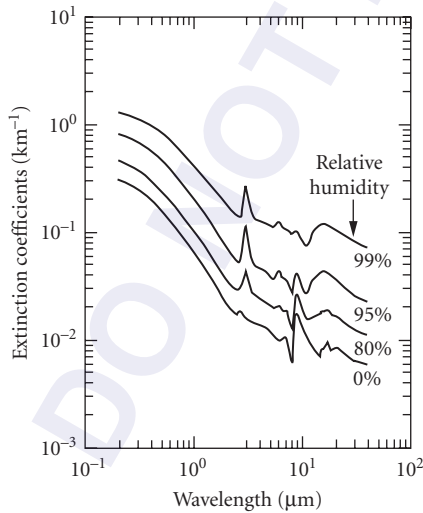


FIGURE 15 Extinction coefficients vs. wavelength for the rural aerosol model for different relative humidities and constant number density of particles.

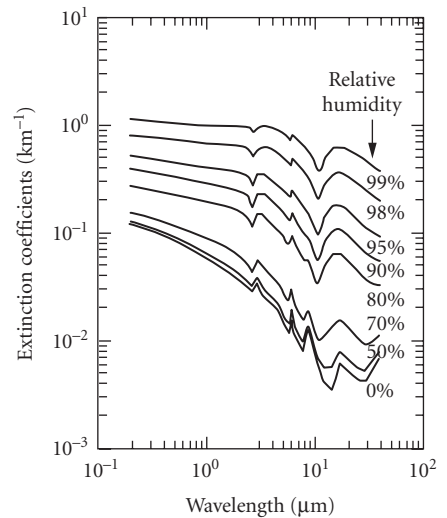


FIGURE 16 Extinction coefficients vs. wavelength for the maritime aerosol model for different relative humidities and constant number density of particles.

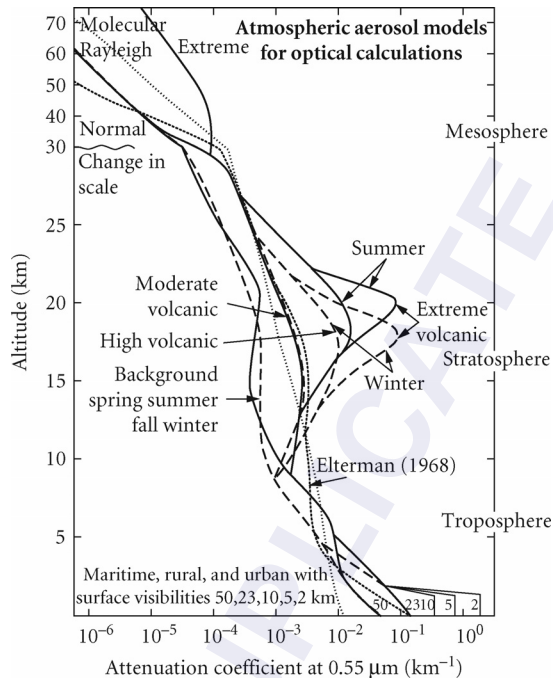


FIGURE 17 The vertical distribution of the aerosol extinction coefficient (at 0.55- μm wavelength) for the different atmospheric models. Also shown for comparison are the Rayleigh profile (dotted line). Between 2 and 30 km, where the distinction on a seasonal basis is made, the spring-summer conditions are indicated with a solid line and fall-winter conditions are indicated by a dashed line. (From Ref. 3.)

The extinction coefficient is also a function of altitude, following the dependence of the composition of the aerosols. Figure 17 shows an atmospheric aerosol extinction model as a function of altitude for a wavelength of 0.55 μm .^{3,10} The influence of the visibility (in km) at ground level dominates the extinction value at the lower altitudes and the composition and density of volcanic particulate dominates the upper altitude regions. The dependence of the extinction on the volcanic composition at the upper altitudes is shown in Fig. 18 which shows these values as a function of wavelength and of composition.^{3,10}

The variation of the backscatter coefficient as a function of altitude is shown in Fig. 19 which displays atmospheric backscatter data obtained by McCormick using a 1.06- μm Nd:YAG Lidar.²⁶ The boundary layer aerosols dominate at the lower levels and the decrease in the atmospheric particulate density determines the overall slope with altitude. Of interest is the increased value near 20 km due to the presence of volcanic aerosols in the atmosphere due to the eruption of Mt. Pinatubo in 1991.

Molecular Emission and Thermal Spectral Radiance

The same optical molecular transitions that cause absorption also emit light when they are thermally excited. Since the molecules have a finite temperature T , they will act as blackbody radiators with optical emission given by the Planck radiation law. The allowed transitions of the molecules

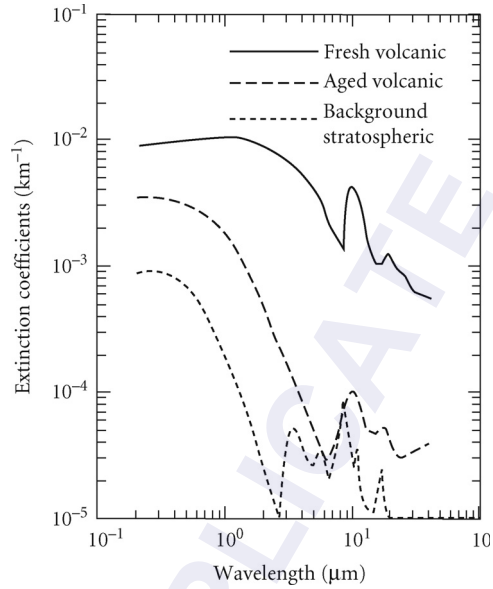


FIGURE 18 Extinction coefficients for the different stratospheric aerosol models (background, volcanic, and fresh volcanic). The extinction coefficients have been normalized to values around peak levels for these models.

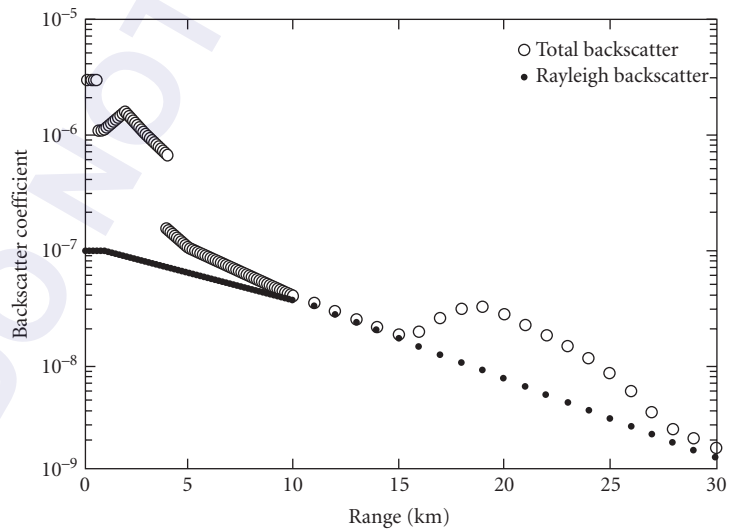


FIGURE 19 1.06- μm lidar backscatter coefficient measurements as a function of vertical altitude. (From McCormick and Winker, Ref. 26.)

will modify the radiance distribution of the radiation due to emission of the radiation according to the thermal distribution of the population within the energy levels of the molecule; it should be noted that the Boltzmann thermal population distribution is essentially the same as that which is described by the Planck radiation law for local thermodynamic equilibrium conditions. As such, the molecular emission spectrum of the radiation is similar to that for absorption. The thermal radiance from the clear atmosphere involves the calculation of the blackbody radiation emitted by each elemental volume of air multiplied by the absorption spectral distribution of the molecular absorption lines, $\kappa_a(s)$ and then this emission spectrum is attenuated by the rest of the atmosphere as the emission propagates toward the viewer. This may be expressed as

$$I_\nu = \int_0^s \kappa_a(s) P_\nu(s) \exp\left[-\int_0^s \kappa_a(s') ds'\right] ds \quad (12)$$

where the exponential term is Beer's law, and $P_\nu(s)$ is the Planck function given by

$$P_\nu(s) = 2h\nu^3/[c^2 \exp([h\nu/kT(s)] - 1)] \quad (13)$$

In these equations, s is the distance from the receiver along the optical propagation path, ν is the optical frequency, h is Planck's constant, c is the speed of light, k is Boltzmann's constant, and $T(s)$ is the temperature at position s along the path. As seen in Eq. (12), each volume element emits thermal radiation of $\kappa_a(s)P_\nu(s)$, which is then attenuated by Beer's law. The total emission spectral density is obtained by summing or integrating over all the emission volume elements and calculating the appropriate absorption along the optical path for each element.

As an example, Fig. 20 shows a plot of the spectral radiance measured on a clear day with 1 cm^{-1} spectral resolution. Note that the regions of strong absorption produce more radiance as the foregoing equation suggests, and that regions of little absorption correspond to little radiance. In the 800- to 1200-wave number spectral region (i.e., 8.3- to 12.5- μm wavelength region), the radiance is relatively low. This is consistent with the fact that the spectral region from 8 to 12 μm is a transmission window of the atmosphere with relatively little absorption of radiation.

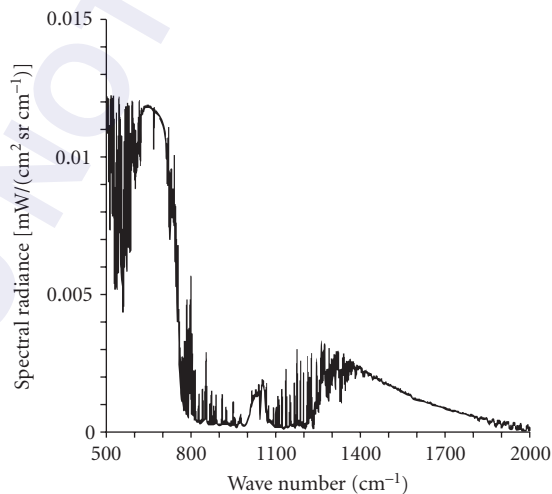


FIGURE 20 Spectral radiance (molecular thermal emission) measured on a clear day showing the relatively low value of radiance near 1000 cm^{-1} (i.e., 10- μm wavelength). (Provided by Churnside.)

Surface Reflectivity and Multiple Scattering

The spectral intensity of naturally occurring light at the earth's surface is primarily due to the incident intensity from the sun in the visible to mid-IR wavelength range, and due to thermal emission from the atmosphere and background radiance in the mid-IR. In both cases, the optical radiation is affected by the reflectance characteristics of the clouds and surface layers. For instance, the fraction of light that falls on the earth's surface and is reflected back into the atmosphere is dependent upon the reflectivity of the surface, the incident solar radiation (polarization and spectral density), and the absorption of the atmosphere.

The reflectivity of a surface, such as the earth's surface, is often characterized using the bidirectional reflectance function (BDRF). This function accounts for the nonspecular reflection of light from common rough surfaces and describes the changes in the reflectivity of a surface as a function of the angle which the incident beam makes with the surface. In addition, the reflectivity of a surface is usually a function of wavelength. This latter effect can be seen in Fig. 21 which shows the reflectance of several common substances for normal incident radiation.² As seen in Fig. 21, the reflectivity of these surfaces is a strong function of wavelength.

The effect of multiple scattering sometimes must be considered when the scattered light undergoes more than one scatter event, and is rescattered on other particles or molecules. These multiple scattering events increase with increasing optical thickness and produce deviations from the Beer-Lambert law. Extensive analyses of the scattering processes for multiple scattering have been conducted and have shown some success in predicting the overall penetration of light through a thick dense cloud. Different computational techniques have been used including the Gauss-Seidel Iterative Method, Layer Adding Method, and Monte-Carlo Techniques.^{3,5}

Additional Optical Interactions

In some optical experiments on the atmosphere, a laser beam is used to excite the molecules in the atmosphere to emit inelastic radiation. Two important inelastic optical processes for atmospheric remote sensing are fluorescence and Raman scattering.^{5,27}

For the case of laser-induced fluorescence, the molecules are excited to an upper energy state and the reemitted photons are detected. In these experiments, the inelastic fluorescence emission is red-shifted in wavelength and can be distinguished in wavelength from the elastic scattered Rayleigh or Mie backscatter. Laser-induced fluorescence is mostly used in the UV to visible spectral region; collisional quenching is quite high in the infrared so that the fluorescence efficiency is higher in the UV-visible than in the IR. Laser-induced fluorescence is sometimes reduced by saturation effects due to stimulated emission from the upper energy levels. However, in those cases where laser-induced fluorescence can be successfully used, it is one of the most sensitive optical techniques for the detection of atomic or molecular species in the atmosphere.

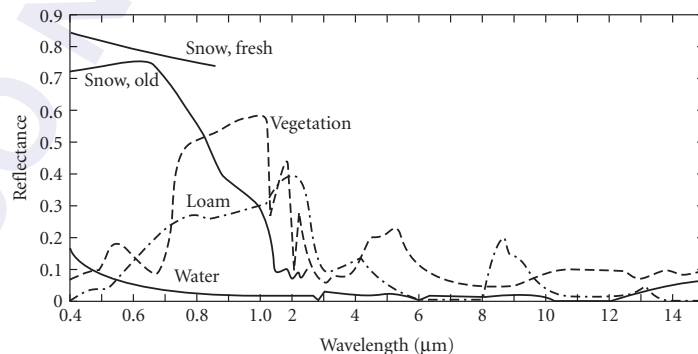


FIGURE 21 Typical reflectance of water surface, snow, dry soil, and vegetation. (From Ref. 2.)

Laser-induced Raman scattering of the atmosphere is a useful probe of the composition and temperature of concentrated species in the atmosphere. The Raman-shifted emitted light is often weak due to the relatively small cross section for Raman scattering. However, for those cases where the distance is short from the laser to the measurement cloud, or where the concentration of the species is high, it offers significant information concerning the composition of the gaseous atmosphere.

The use of an intense laser beam can also bring about nonlinear optical interactions as the laser beam propagates through the atmosphere. The most important of these are stimulated Raman scattering, thermal blooming, dielectric breakdown, and harmonic conversion. Each of these processes requires a tightly focused laser beam to initiate the nonlinear optical process.^{28,29}

3.5 PREDICTION OF ATMOSPHERIC OPTICAL TRANSMISSION: COMPUTER PROGRAMS AND DATABASES

During the past three decades, several computer programs and databases have been developed which are very useful for the determination of the optical properties of the atmosphere. Many of these are based upon programs originally developed at the U.S. Air Force Cambridge Research Laboratories. The latest versions of these programs and databases are the HITRAN database,¹³ FASCODE computer program,^{30–32} and the LOWTRAN or MODTRAN computer code.^{33,34} In addition, several PC (personal computer) versions of these database/computer programs have recently become available so that the user can easily use these computational aids.

Molecular Absorption Line Database: HITRAN

The HITRAN database contains optical spectral data on most of the major molecules contributing to absorption or radiance in the atmosphere; details of HITRAN are covered in several recent journal articles.^{11–13} The 40 molecules contained in HITRAN are given in Table 1, and cover over a million individual absorption lines in the spectral range from 0.000001 cm^{-1} to $25,233\text{ cm}^{-1}$ (i.e., 0.3963 to $10^{10}\text{ }\mu\text{m}$). A free copy of this database can be obtained by filling out a request form in the HITRAN Web site.³⁵ Each line in the database contains 19 molecular data items that consist of the molecule formula code, isotopologue type, transition frequency (cm^{-1}), line intensity S in $\text{cm}/\text{molecule}$, Einstein A -coefficient (s^{-1}), air-broadened half-width ($\text{cm}^{-1}/\text{atm}$), self-broadened half-width ($\text{cm}^{-1}/\text{atm}$), lower state energy (cm^{-1}), temperature coefficient for air-broadened linewidth, air-pressure induced line shift ($\text{cm}^{-1}\text{ atm}^{-1}$), upper-state global quanta index, lower-state global quanta index, upper- and lower-state quanta, uncertainty codes, reference numbers, a flag for line coupling if necessary, and upper- and lower-level statistical weights. The density of the lines of the 2004 HITRAN database is shown in Fig. 22. The recently released 2008 HITRAN database has been expanded to 42 molecules.

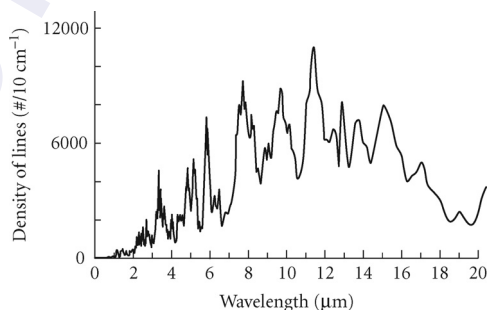


FIGURE 22 Density of absorption lines in HITRAN 2004 spectral database.

<i>M</i>	<i>I</i>	wn (cm ⁻¹)	<i>S</i> (T)	<i>R</i>	<i>g</i> (T)	<i>E</i> ''	<i>Q</i> '	<i>Q</i> ''
NH ₃	1	5000.04530	5.161E-23	2.33E-5	.0599	244.09880		
CO ₂	1	5000.05653	3.823E-27	5.99E-6	.0680	2783.34814		P 33
CO ₂	3	5000.06953	1.766E-26	1.46E-6	.0657	975.59039		P 51
NH ₃	1	5000.20320	3.096E-23	1.54E-6	.1059	16.20000	2 0-1	s 1 1 0 s
CO ₂	1	5000.21444	1.358E-25	8.01E-6	.0627	2240.23755		R 63
H ₂ O	1	5000.22500	2.160E-24	1.86E-4	.0088	2358.30396	14 014	15 015
NH ₃	1	5000.32760	2.854E-23	1.42E-6	.1059	17.00000	2 0-1	a 1 1 0 a
CO ₂	2	5000.33471	1.524E-26	1.00E-5	.0800	1398.12122		R 8
CO ₂	1	5000.34883	3.939E-25	8.02E-6	.0640	2004.66040		R 58
NH ₃	1	5000.39930	2.100E-23	1.18E-5	.0599	289.09781		
N ₂ O	1	5000.42800	2.990E-23	6.82E-7	.0719	316.67001		P 27
CO ₂	2	5000.47621	2.938E-25	4.38E-6	.0739	754.88672		P 16
CO ₂	1	5000.48083	4.762E-22	9.21E-6	.0677	464.17169		R 34
CO ₂	2	5000.48726	5.469E-27	4.31E-6	.0681	1710.65881		P 32
H ₂ O	1	5000.54070	6.959E-27	1.10E-8	.0556	1411.61206	9 5 4	8 6 3
NH ₃	1	5000.65640	2.221E-23	1.54E-5	.0599	333.09299		
H ₂ O	4	5000.75000	4.629E-27	3.10E-6	.0856	809.39502	8 3 5	8 4 4
CO ₂	1	5000.81774	8.595E-26	8.01E-6	.0623	2340.73218		R 65
H ₂ O	4	5000.83600	4.469E-26	3.43E-6	.0922	308.61700	5 0 5	6 1 6
CO ₂	2	5000.85619	7.150E-24	4.82E-6	.0762	60.87380		R 12
NH ₃	1	5000.88490	4.193E-23	1.4E0				

FIGURE 23 Example of data contained within the HITRAN database showing individual absorption lines, frequency, line intensity, and other spectroscopic parameters.

Figure 23 shows an output from a computer program that was used to search the HITRAN database and display some of the pertinent information.³⁶ The data in HITRAN are in sequential order by transition frequency in wave numbers, and list the molecular name, isotope, absorption line strength *S*, transition probability *R*, air-pressure-broadened linewidth γ_g , lower energy state *E*'', and upper/lower quanta for the different molecules and isotopic species in the atmosphere.

Line-by-Line Transmission Program: FASCODE

FASCODE is a large, sophisticated computer program that uses molecular absorption equations (similar to those under "Molecular Absorption") and the HITRAN database to calculate the high-resolution spectra of the atmosphere. It uses efficient algorithms to speed the computations of the spectral transmission, emission, and radiance of the atmosphere at a spectral resolution that can be set to better than the natural linewidth,³² and includes the effects of Rayleigh and aerosol scattering and the continuum molecular extinction. FASCODE also calculates the radiance and transmittance of atmospheric slant paths, and can calculate the integrated transmittance through the atmosphere from the ground up to higher altitudes. Voigt lineshape profiles are also used to handle the transition from pressure-broadened lineshapes near ground level to the Doppler-dominated lineshapes at very high altitudes. Several representative models of the atmosphere are contained within FASCODE, so that the user can specify different seasonal and geographical models. Figure 24 shows a sample output generated from data produced from the FASCODE program and a comparison with experimental data obtained by J. Dowling at NRL.^{30,31} As can be seen, the agreement is very good. There are also several other line-by-line codes available for specialized applications; examples include GENLN2 developed by D.P. Edwards at NCAR (National Center for Atmospheric Research/Boulder), and LBLRTM (Atmospheric and Environmental Research, Inc.).

Broadband Transmission: LOWTRAN and MODTRAN

The LOWTRAN computer program does not use the HITRAN database directly, but uses absorption band models based on degrading spectral calculations based on HITRAN to calculate the moderate resolution (20 cm⁻¹) transmission spectrum of the atmosphere. LOWTRAN uses extensive

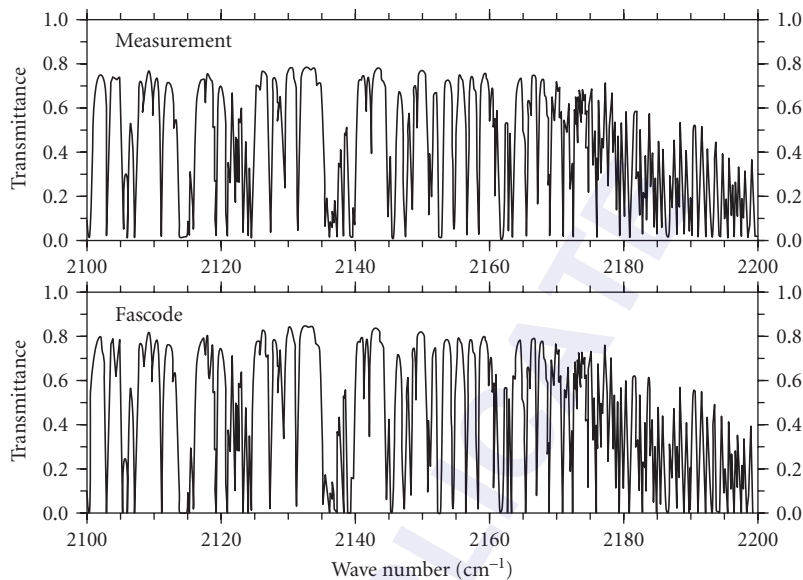


FIGURE 24 Comparison of an FASCOD2 transmittance calculation with an experimental atmospheric measurement (from NRL) over a 6.4-km path at the ground. (Courtesy of Clough, Ref. 30.)

band-model calculations to speed up the computations, and provides an accurate and rapid means of estimating the transmittance and background radiance of the earth's atmosphere over the spectral interval of 350 cm^{-1} to $40,000\text{ cm}^{-1}$ (i.e., 250-nm–28- μm wavelength). The spectral range of the LOWTRAN program extends into the UV. In the LOWTRAN program, the total transmittance at a given wavelength is given as the product of the transmittances due to molecular band absorption, molecular scattering, aerosol extinction, and molecular continuum absorption. The molecular band absorption is composed of four components of water vapor, ozone, nitric acid, and the uniformly mixed gases (CO_2 , N_2O , CH_4 , CO , O_2 , and N_2).

The latest version of LOWTRAN(7) contains models treating solar and lunar scattered radiation, spherical refractive geometry, slant-path geometry, wind-dependent maritime aerosols, vertical structure aerosols, standard seasonal and geographic atmospheric models (e.g., mid-latitude summer), cirrus cloud model, and a rain model.³⁴ As an example, Fig. 25 shows a ground-level solar radiance model used by LOWTRAN, and Fig. 26 shows an example of a rain-rate model and its effect upon the transmission of the atmosphere as a function of rain rate in mm of water per hour.³⁴

Extensive experimental measurements have been made to verify LOWTRAN calculations. Figure 27 shows a composite plot of the LOWTRAN-predicted transmittance and experimental data for a path length of 1.3 km at sea level.³ As can be seen, the agreement is quite good. It is estimated that the LOWTRAN calculations are good to about 10 percent.³ It should be added that the molecular absorption portion of the preceding LOWTRAN (moderate-resolution) spectra can also be generated using the high-resolution FASCOD/HITRAN program and then spectrally smoothing (i.e., degrading) the spectra to match that of the LOWTRAN spectra.

The most recent extension of the LOWTRAN program is the MODTRAN program. MODTRAN is similar to LOWTRAN but has increased spectral resolution. At present, the resolution for the latest version of MODTRAN, called MODTRAN(5), can be specified by the user between 0.1 and 20 cm^{-1} .

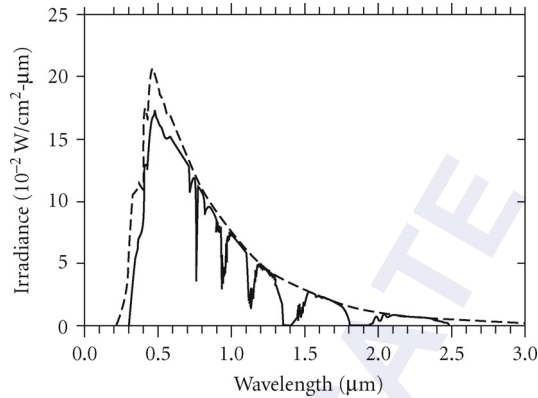


FIGURE 25 Solar radiance model (dashed line) and directly transmitted solar irradiance (solid line) for a vertical path, from the ground (U.S. standard 1962 model, no aerosol extinction) as used by the LOWTRAN program.

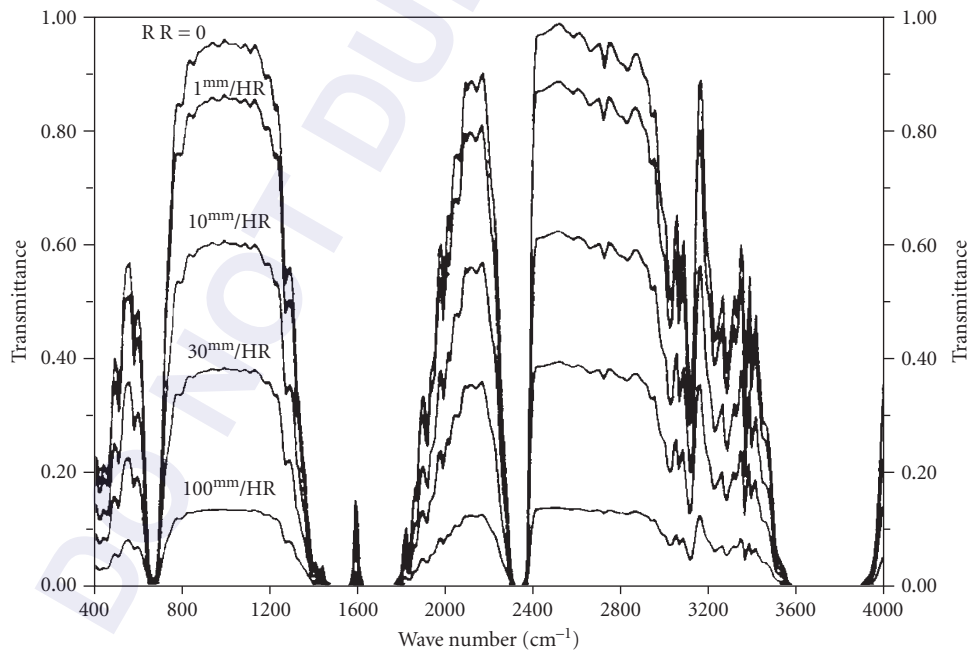


FIGURE 26 Atmospheric transmittance for different rain rates and for spectral frequencies from 400 to 4000 cm^{-1} . The measurement path is 300 m at the surface with $T = T_{\text{dew}} = 10^\circ\text{C}$, with a meteorological range of 23 km in the absence of rain.

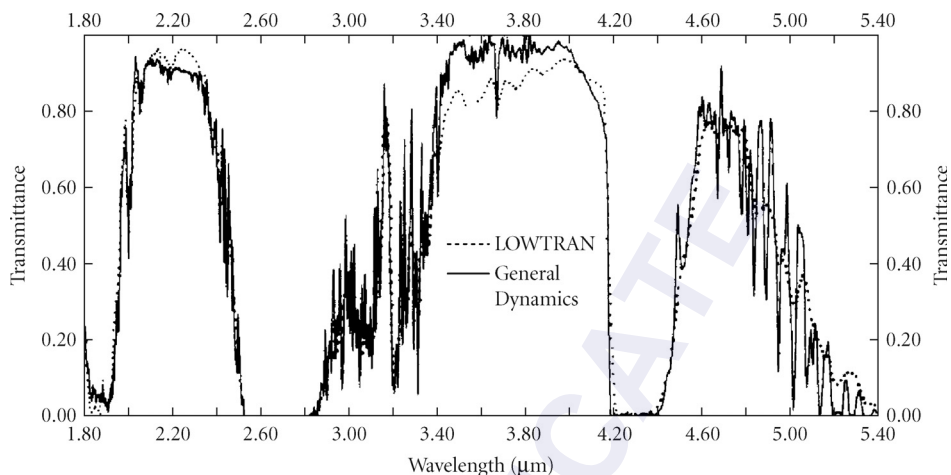


FIGURE 27 Comparison between LOWTRAN predicted spectrum and General Dynamics atmospheric measurements; range = 1.3 km at sea levels. (From Ref. 3.)

Programs and Databases for Use on Personal Computers

The preceding databases and computer programs have been converted or modified to run on different kinds of personal computers.^{35,36} The HITRAN database has been available on CD-ROMs for the past decade, but is now available via the internet. Several related programs are available, ranging from a complete copy of the FASCODE and LOWTRAN programs³⁵ to a simpler molecular transmission program of the atmosphere.³⁶ These programs calculate the transmission spectrum of the atmosphere and some show the overlay spectra of known laser lines. As an example, Fig. 28 shows the transmission spectrum produced by the HITRAN-PC program³⁶ for a horizontal path of 300 m (U.S. Standard Atmosphere) over the wavelength range of 250 nm ($40,000\text{ cm}^{-1}$) to $20\text{ }\mu\text{m}$ (500 cm^{-1}); the transmission spectrum includes water, nitrogen, and CO_2 continuum and urban aerosol attenuation, and was smoothed to a spectral resolution of 1 cm^{-1} to better display the overall transmission features of the atmosphere.

While these PC versions of the HITRAN database and transmission programs have become available only recently, they have already made a significant impact in the fields of atmospheric optics and optical remote sensing. They allow quick and easy access to atmospheric spectral data which was previously only available on a mainframe computer. It should be added that other computer programs are available which allow one to add or subtract different spectra generated by these HITRAN-based programs, from spectroscopic instrumentation such as FT-IR spectrometers or from other IR gas spectra databases. In the latter case, for example, the U.S. National Institute of Standards and Technology (NIST) has compiled a computer database of the qualitative IR absorption spectra of over 5200 different gases (toxic and other hydrocarbon compounds) with a spectral resolution of 4 cm^{-1} .³⁷ The Pacific Northwest National Laboratory also offers absorption cross-section files of numerous gases.³⁸ In addition, higher-resolution quantitative spectra for a limited group of gases can be obtained from several commercial companies.³⁹

3.6 ATMOSPHERIC OPTICAL TURBULENCE

The most familiar effects of refractive turbulence in the atmosphere are the twinkling of stars and the shimmering of the horizon on a hot day. The first of these is a random fluctuation of the amplitude of light also known as scintillation. The second is a random fluctuation of the phase front that

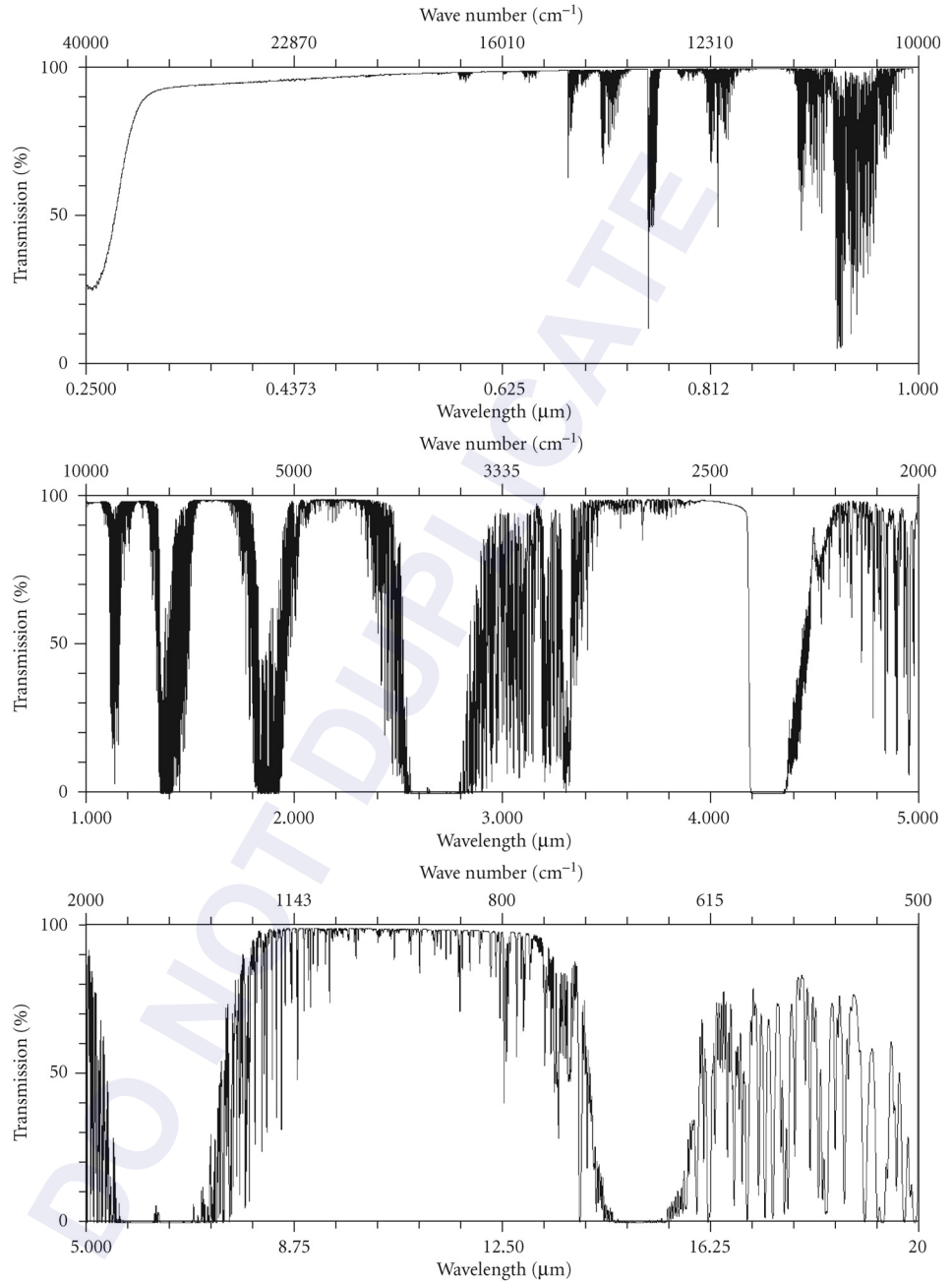


FIGURE 28 Example of generated atmospheric transmission spectrum of the atmosphere for a horizontal path of 300 m for the wavelength range from UV (250 nm or 40,000 cm^{-1}) to the IR (20 μm or 500 cm^{-1}); the spectrum includes water, nitrogen, and CO_2 continuum, urban ozone and NO_2 , and urban aerosol attenuation, and has been smoothed to a resolution of 0.5 cm^{-1} to better show the absorption and transmission windows of the atmosphere. (From Ref. 36 and D. Plütau.)

leads to a reduction in the resolution of an image. Other effects include the wander and break-up of an optical beam. A detailed discussion of all of these effects and the implications for various applications can be found in Ref. 40.

In the visible and near-IR region of the spectrum, the fluctuations of the refractive index in the atmosphere are determined by fluctuations of the temperature. These temperature fluctuations are caused by turbulent mixing of air of different temperatures. In the far-IR region, humidity fluctuations also contribute.

Turbulence Characteristics

Refractive turbulence in the atmosphere can be characterized by three parameters. The outer scale L_0 is the length of the largest scales of turbulent eddies. The inner scale l_0 is the length of the smallest scales. For eddies in the inertial subrange (sizes between the inner and outer scale), the refractive index fluctuations are best described by the structure function. This function is defined by

$$D_n(r_1, r_2) = \langle [n(r_1) - n(r_2)]^2 \rangle \quad (14)$$

where $n(r_1)$ is the index of refraction at point r_1 and the angle brackets denote an ensemble average. For homogeneous and isotropic turbulence it depends only on the distance between the two points r and is given by

$$D_n(r) = C_n^2 r^{2/3} \quad (15)$$

where C_n^2 is a measure of the strength of turbulence and is defined by this equation.

The power spectrum of turbulence is the Fourier transform of the correlation function, which is contained in the cross term of the structure function. For scales within the inertial subrange, it is given by the Kolmogorov spectrum:

$$\Phi_n(K) = 0.033 C_n^2 K^{-11/3} \quad (16)$$

For scales larger than the outer scale, the spectrum actually approaches a constant value, and the result can be approximated by the von Kármán spectrum:

$$\Phi_n(K) = 0.033 C_n^2 (K^2 + K_0^2)^{-11/6} \quad (17)$$

where K_0 is the wave number corresponding to the outer scale. For scales near the inner scale, there is a small increase over the Kolmogorov spectrum, with a large decrease at smaller scales.⁴¹ The resulting spectrum can be approximated by a rather simple function.^{42,43}

In the boundary layer (the lowest few hundred meters of the atmosphere), turbulence is generated by radiative heating and cooling of the ground. During the day, solar heating of the ground drives convective plumes. Refractive turbulence is generated by the mixing of these warm plumes with the cooler air surrounding them. At night, the ground is cooled by radiation and the cooler air near the ground is mixed with warmer air higher up by winds. A period of extremely low turbulence exists at dawn and at dusk when there is no temperature gradient in the lower atmosphere. Turbulence levels are also very low when the sky is overcast and solar heating and radiative cooling rates are low. Measured values of turbulence strength near the ground vary from less than 10^{-17} to greater than $10^{-12} \text{ m}^{-2/3}$ at heights of 2 to 2.5 m.^{44,45}

Figure 29 illustrates typical summertime values near Boulder, Colorado. This is a 24-hour plot of 15-minute averages of C_n^2 measured at a height of about 1.5 m on August 22, 1991. At night, the sky was clear, and C_n^2 was a few parts times 10^{-13} . The dawn minimum is seen as a very short period of low turbulence just after 6:00. After sunrise, C_n^2 increases rapidly to over 10^{-12} . Just before noon, cumulus clouds developed, and C_n^2 became lower with large fluctuations. At about 18:00, the clouds dissipated,

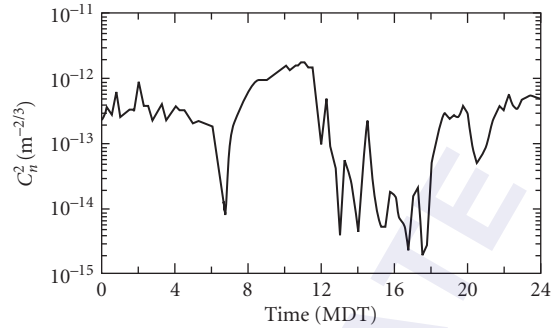


FIGURE 29 Plot of refractive-turbulence structure parameter C_n^2 for a typical summer day near Boulder, Colorado. (Courtesy G. R. Ochs, NOAA WPL.)

and turbulence levels increased. The dusk minimum is evident just after 20:00, and then turbulence strength returns to typical nighttime levels.

From a theory introduced by Monin and Obukhov,⁴⁶ a theoretical dependence of turbulence strength on height in the boundary layer above flat ground can be derived.^{47,48} During periods of convection C_n^2 decreases as the $-4/3$ power of height. At other times (night or overcast days), the power is more nearly $-2/3$. These height dependencies have been verified by a number of experiments over relatively flat terrain.⁴⁹⁻⁵³ However, values measured in mountainous regions are closer to the $-1/3$ power of height day or night.⁵⁴ Under certain conditions, the turbulence strength can be predicted from meteorological parameters and characteristics of the underlying surface.⁵⁵⁻⁵⁷

Farther from the ground, no theory for the turbulence profile exists. Measurements have been made from aircraft^{44,49} and with balloons.⁵⁸⁻⁶⁰ Profiles of C_n^2 have also been measured remotely from the ground using acoustic sounders,⁶¹⁻⁶³ radar,^{49,50,64-69} and optical techniques.^{59,70-73} The measurements show large variations in refractive turbulence strength. They all exhibit a sharply layered structure in which the turbulence appears in layers of the order of 100 m thick with relatively calm air in between. In some cases these layers can be associated with orographic features; that is, the turbulence can be attributed to mountain lee waves. Generally, as height increases, the turbulence decreases to a minimum value that occurs at a height of about 3 to 5 km. The turbulence level then increases to a maximum at about the tropopause (10 km). Turbulence levels decrease rapidly above the tropopause.

Model turbulence profiles have evolved from this type of measurement. Perhaps the best available model for altitudes of 3 to 20 km is the Hufnagel model:^{74,75}

$$C_n^2 = \left\{ \left[(2.2 \times 10^{-53}) H^{10} \left(\frac{W}{27} \right)^2 \right] \exp\left(-\frac{H}{1000}\right) + 10^{-16} \exp\left(-\frac{H}{1500}\right) \right\} \exp[u(H, t)] \quad (18)$$

where H is the height above sea level in meters, W is the vertical average of the square of the wind speed, and u is a random variable that allows the random nature of the profiles to be modeled. W is defined by

$$W^2 = \frac{1}{1500} \int_{5000}^{20,000} v^2(H) dH \quad (19)$$

where $v(H)$ is the wind speed at height H . In data taken over Maryland, W was normally distributed with a mean value of 27 m/s and a standard deviation of 9 m/s. The random variable u is assumed to be a zero-mean, Gaussian variable with a covariance function given by

$$\langle u(H, t)u(H + \delta H, t + \delta t) \rangle = A(\delta H / 100) \exp(-\delta t / 5) + A(\delta H / 2000) \exp(-\delta t / 80) \quad (20)$$

where

$$A(\delta H/L) = 1 - |\delta H/L| \quad \text{for } |H| < L \quad (21)$$

and equals 0 otherwise.

The time interval δt is measured in minutes. The average C_n^2 profile can be found by recognizing that $\langle \exp(u) \rangle = \exp(1)$. To extend the model to local ground level, one should add the surface layer dependence (e.g., $H^{-4/3}$ for daytime).

Another attempt to extend the model to ground level is the Hufnagel-Valley model.⁷⁶ This is given by

$$C_n^2 = 0.00594 \left(\frac{W}{27} \right)^2 (H \times 10^{-5})^{10} \exp\left(-\frac{H}{1000}\right) + 2.7 \times 10^{-16} \exp\left(-\frac{H}{1500}\right) + A \exp\left(-\frac{H}{100}\right) \quad (22)$$

where W is commonly set to 21 and A to 1.7×10^{-14} . This specific model is referred to as the $HV_{5/7}$ model because it produces a coherence diameter r_0 of about 5 cm and an isoplanatic angle of about $7 \mu\text{rad}$ for a wavelength of $0.5 \mu\text{m}$. Although this is not as accurate for modeling turbulence near the ground, it has the advantage that the moments of the turbulence profile important to propagation can be evaluated analytically.⁷⁶

The $HV_{5/7}$ model is plotted as a function of height in the dashed line in Fig. 30. The solid line in the figure is a balloon measurement taken in College Station, Pennsylvania. The data were reported with 20-m vertical resolution and smoothed with a Gaussian filter with a 100-m $\exp(-1)$ full-width. This particular data set was chosen because it has a coherence diameter of about 5 cm and an isoplanatic angle of about $7 \mu\text{rad}$. The layered structure of the real atmosphere is clear in the data. Note also the difference between the model atmosphere and the real atmosphere even when the coherence diameter and the isoplanatic angle are similar.

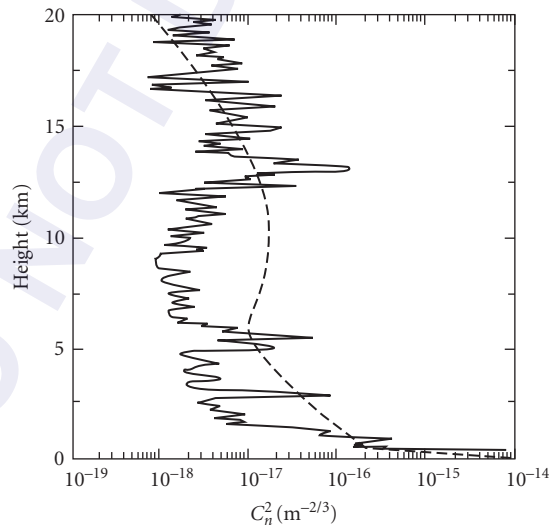


FIGURE 30 Turbulence strength C_n^2 as a function of height. The solid line is a balloon measurement made in College Station, Pennsylvania, and the dashed line is the $HV_{5/7}$ model. (Courtesy R. R. Beland, Geophysics Directorate, Phillips Laboratory, U.S. Air Force.)

Less is known about the vertical profiles of inner and outer scales. Near the ground (1 to 2 m) we typically observe inner scales of 5 to 10 mm over flat grassland in Colorado. Calculations of inner scale from measured values of Kolmogorov microscale range from 0.5 to 9 mm at similar heights.⁷⁷ Aircraft measurements of dissipation rate were used along with a viscosity profile calculated from typical profiles of temperature and pressure to estimate a profile of microscale.⁷⁸ Values increase monotonically to about 4 cm at a height of 10 km and to about 8 cm at 20 km.

Near the ground, the outer scale can be estimated using Monin-Obukhov similarity theory.⁴⁶ The outer scale can be defined as that separation at which the structure function of temperature fluctuations is equal to twice the variance. Using typical surface layer scaling relationships⁷⁹ we see that

$$L_0 = \begin{cases} 7.04H(1-7S_{MO})(1-16S_{MO})^{-32} & \text{for } -2 < S_{MO} \\ 7.04H(1+S_{MO})^{-3}(1+2.75S_{MO}^{2/3})^{-32} & \text{for } 0 < S_{MO} < \end{cases} \quad (23)$$

where S_{MO} is the Monin-Obukhov stability parameter. For typical daytime conditions ($S_{MO} < 0$), L_0 is generally between $H/2$ and H .

Above the boundary layer, the situation is more complex. Barat and Bertin⁸⁰ measured outer scale values of 10 to 100 m in a turbulent layer using a balloonborne instrument. Some recent optical data^{81,82} suggest that the outer scale is generally between 10 and 40 m.

Beam Wander

The first effect of refractive turbulence to consider is the wander of an optical beam in the atmosphere. The deviations of the centroid of a beam in each of the two orthogonal transverse axes will be independent of Gaussian random variables. This wander is generally characterized statistically by the variance of the angular displacement. In isotropic turbulence, the variances in the two axes are equal and the magnitude of the displacement is a Rayleigh random variable with a variance that is twice that of the displacement in a single axis. Both the single-axis and the magnitude variances are reported in the literature.

Three main approaches to the calculation of beam-wander variance have been used: (1) If diffraction effects are negligible and if the path-integrated turbulence is small enough, the geometric optics approximation⁸³⁻⁸⁶ can be used. Diffraction effects are negligible when the aperture diameter is greater than the Fresnel zone size.⁸⁶ The other condition requires that the product of the transverse coherence of the field and the aperture diameter is greater than the square of the Fresnel zone size.⁸⁶ (2) If diffraction must be considered, the Huygens-Fresnel approximation⁸⁷⁻⁹⁰ can be used. (3) The most complete theory uses the Markov random process approximation to the moment equation.⁹¹⁻⁹³ In the two types of calculation that include diffraction, beams with a Gaussian irradiance profile are generally assumed.

In the geometric optics approximation, the variance of the angular displacement in a single axis is given by⁸⁶

$$\sigma_d^2 = 2.92D^{-1/3} \int_0^L dz C_n^2(z) \frac{\left(1 - \frac{z}{L}\right)^2}{\left|1 - \frac{z}{F}\right|^{1/3}} \quad (24)$$

where D is the diameter of the initial beam, z is position along the path, L is the propagation path length, and F is the geometric focal range of the initial beam (negative for a diverging beam). For homogeneous turbulence and a beam that does not come to a focus between the transmitter and observation plane, Eq. (24) reduces to

$$\sigma_d^2 = 0.97C_n^2 D^{-1/3} L_2 F_1\left(\frac{1}{3}, 1; 4; \frac{L}{F}\right) \quad (25)$$

where ${}_2F_1$ is the hypergeometric function. The hypergeometric function is 1 for a collimated beam and 1.125 for a focused beam.

In the Markov approximation, the single-axis variance is given by⁹¹

$$\sigma_a^2 = 4\pi^2 \int_0^L dz \left(1 - \frac{z}{L}\right)^2 \int_0^\infty dK K^3 \Phi_n(K, z) \times \exp\left\{-\frac{K^2 D^2}{8} \left[\left(1 - \frac{z}{F}\right)^2 + \frac{16z^2}{K^2 D^2}\right] - \pi D_\Psi \left(\frac{Kz}{k}\right)\right\} \quad (26)$$

where $\Phi_n(K, z)$ is the path-dependent refractive index spectrum, K is the wave number of turbulence, D is the exp (-1) irradiance diameter of the initial beam, k is the optical wave number, and $D_\Psi(r)$ is the wave structure function for separation r of a spherical wave. The structure function is given by

$$D_\Psi(r) = 8\pi^2 k^2 \int_0^z dz' \int_0^\infty dK' K' \Phi_n(K', z') [1 - J_0(K' r z' / z)] \quad (27)$$

where J_0 is the zero-order Bessel function of the first kind. The last term in the exponential of Eq. (26) is a correction term for strong turbulence. The middle term includes the effects of diffraction.

Beam Spreading

The next effect of refractive turbulence to consider is the spread of an optical beam as it propagates through the atmosphere. There are two types of beam spread denoted as long term and short term. The long-term beam spread is defined as the turbulence-induced beam spread observed over a long time average. It includes the effects of the slow wander of the entire beam. The short-term beam spread is defined as the beam spread observed at an instant of time. It does not include the effects of beam wander and is often approximated by the long-term beam spread with the effects of wander removed, although the two are not identical.

We can consider beam wander to be caused by turbulent eddies that are larger than the beam. Short-term beam spread is caused by turbulent eddies that are smaller than the beam. There are more small eddies than large in the beam at any time which implies that the beam spread at any instant is averaged over more eddies. As a result, the fluctuations of short-term beam spread are much smaller than those of beam wander and are typically neglected. The primary effect of short-term beam spread is to spread the average energy over a larger area. Thus, the average value of the on-axis irradiance is reduced, and the average value of the irradiance at large angles is increased.

The extended Huygens-Fresnel principle can be used to calculate the long-term spread of a Gaussian beam in refractive turbulence.⁹⁴⁻⁹⁸ The exp (-1) irradiance radius of a Gaussian beam is

$$P_1 = \left[\frac{4}{k^2 D^2} + \frac{D^2}{4} \left(1 - \frac{L}{F}\right)^2 + \frac{4}{k^2 \rho_0^2} \right]^{1/2} \quad (28)$$

where D is the exp (-1) irradiance diameter of the source and ρ_0 is the phase coherence length that would be observed for a point source propagating from the receiver to the transmitter. The first term in this equation is the diffraction beam spread, the second is the geometrical optics projection of the transmitter aperture, and the final term is the total turbulence-induced spread.

The phase coherence length is defined as the transverse separation at which the coherence of the field is reduced to exp (-1) . If the coherence length ρ_0 is much larger than the inner scale, the structure function is given by $D_\Psi(r) = 2(r/\rho_0)^{5/3}$ and

$$\rho_0 = \left[1.46 k^2 \int_0^L dz \left(\frac{z}{L}\right)^{5/3} C_n^2(z) \right]^{-3/5} \quad (29)$$

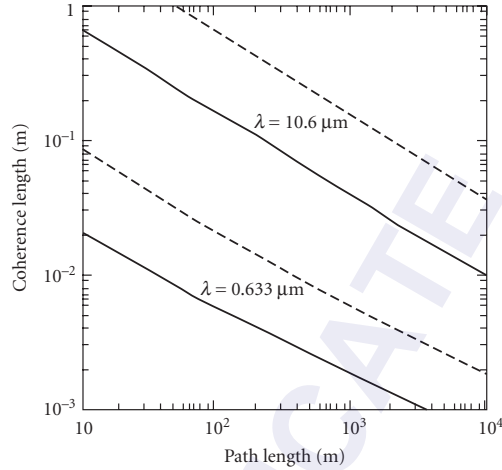


FIGURE 31 Coherence length ρ_0 as a function of path length for two common laser wavelengths. An inner scale of 1 cm, and values of 10^{-12} (solid lines) and 10^{-13} (dashed lines) were used in the calculations.

If the coherence length is much smaller than the inner scale, then $D_\Psi(r) = 2(r/\rho_0)^2$ and

$$\rho_0 = \left[1.86k^2 \int_0^L dz \left(\frac{z}{L} \right)^2 C_n^2(z) l_0^{-1/3}(z) \right]^{-1/2} \quad (30)$$

Typical values of ρ_0 for spherical-wave propagation through homogeneous turbulence are presented in Fig. 31. In these plots there is a slope change from $L^{-3/5}$ to $L^{-1/2}$, where ρ_0 is about equal to the inner scale of 1 cm. This small change is nearly imperceptible in the figure.

The extended Huygens-Fresnel principle has been used to calculate the short-term beam spread by explicitly subtracting the beam wander.⁹⁹⁻¹⁰² If ρ_0 and l_0 are much smaller than D , then the short-term beam spread is approximately given by

$$p_s = \left\{ \frac{4}{k^2 D^2} + \frac{D^2}{4L^2} \left(1 - \frac{L}{F} \right)^2 + \frac{4}{k^2 \rho_0^2} \left[1 - 0.62 \left(\frac{\rho_0}{D} \right)^{1/3} \right]^{6/5} \right\}^{1/2} \quad (31)$$

If ρ_0 is much greater than D , the turbulence-induced component of beam spreading can be neglected.

If inner scale and outer scale effects are included, more complicated integral expressions result.¹⁰³ Numerical calculations were performed for truncated Gaussian beams with central obscurations.¹⁰⁴ The following approximation was obtained by a curve fit to the results:

$$p_s = \left(1 + 0.182 \frac{D_{\text{eff}}^2}{r_0^2} \right)^{1/2} p_d \quad \text{for} \quad \frac{D_{\text{eff}}}{r_0} < 3$$

$$p_s = \left[1 + \left(\frac{D_{\text{eff}}}{r_0} \right)^2 - 1.18 \left(\frac{D_{\text{eff}}}{r_0} \right)^{-5/3} \right]^{1/2} p_d \quad \text{for} \quad 3 < \frac{D_{\text{eff}}}{r_0} < \quad (32)$$

where D_{eff} is the effective diameter of the truncated aperture, $r_0 = 2.099 p_0$, and p_d is the diffraction-limited value. These expressions agree fairly well with available data.^{87,105,106}

Imaging and Heterodyne Detection

The problem of imaging through the turbulent atmosphere is similar to the problem of beam propagation through the atmosphere. The dancing of an image in the focal plane of an imaging system is mathematically equivalent to the wander of a beam focused at the object by the same optical system. The resolution of the short-exposure image is equivalent to the short-term beam spread of the same focused beam. The resolution for a long exposure is related to the long-term beam spread.

Thus, the position of the image of a point object will drift in each axis in the focal plane. The variance of that drift is given by¹⁰⁷

$$\sigma_i^2 = 1.10C_n^2 D^{-1/3} L F^2 \quad (33)$$

where D is the aperture diameter of the imaging system, F is its focal length, and L is the distance to the object.

Fried¹⁰⁸ used the idea of tilt correction to calculate the average image resolution for a short-exposure image in the turbulent atmosphere. This problem is mathematically equivalent to the propagation of a beam in the opposite direction if the imaging aperture replaces the beam width. These results were refined by Lutomirski et al.¹⁰⁹ and applied to a space-to-ground path by Valley.¹⁰³

Image resolution is also related to the signal-to-noise ratio of an optical heterodyne receiver. The long-exposure resolution is equivalent to a staring receiver. The short-exposure resolution is equivalent to a receiver that employs tilt-correction of the signal or of the local oscillator.^{108,110-112}

Scintillation

The refractive index inhomogeneities that distort the optical phase front also produce amplitude scintillation at some distance. The first cases to be considered were plane-wave and spherical-wave propagation through weak path-integrated turbulence where the weak-turbulence condition requires that fluctuations of irradiance be much less than the mean value. Tatarskii¹¹³ used a perturbation approach to the wave equation. Lee and Harp¹¹⁴ used a physical approach to arrive at the same results. These results are summarized in a number of good reviews.^{101,102,115,116}

We will first consider the weak-turbulence results. For propagation from space to the ground, the plane wave formula is generally valid. The variance of irradiance fluctuations (normalized by the mean irradiance value) is given by¹¹⁵

$$\sigma_r^2 = k^{7/6} \sec^{11/6} \theta \int_0^\infty dH H^{5/6} C_n^2(H) \quad (34)$$

where k is the optical wave number, θ is the zenith angle, and H is the height of the receiver above the ground. This expression is valid as long as the path-integrated turbulence is weak enough that the variance is much less than unity. This condition is usually met for near-zenith propagation.

For propagation of diverging waves near the ground, the spherical-wave approximation is often valid. Assuming constant turbulence along the path, the weak-turbulence variance in this case is given by¹¹⁵

$$\begin{aligned} \sigma_r^2 &= \exp [0.5k^{7/6} L^{11/6} C_n^2] - 1 & \text{for } l_0 < \sqrt{L} \\ \sigma_r^2 &= \exp [1.28L^3 l_0^{-7/3} C_n^2] - 1 & \text{for } l_0 > \sqrt{L} \end{aligned} \quad (35)$$

where L is the path length and l_0 is the inner scale of turbulence.

For narrow-beam propagation, the effects of the finite beam must be considered. Kon and Tatarskii¹¹⁷ calculated the amplitude fluctuations of a collimated beam using the perturbation technique. Schmeltzer¹¹⁸ extended these results to include focused beams. These results were used to obtain numerical values for a variety of propagation conditions.¹¹⁹⁻¹²¹ Ishimaru¹²²⁻¹²⁴ used a spectral representation to obtain similar results. Under certain conditions, one sees a reduction in the variance on the optical axis^{119,120,125,126} and an increase off of the optical axis.¹²⁷

The spatial scale of weak scintillations is about equal to the larger of either the Fresnel zone size $(L/k)^{1/2}$ or the inner scale.¹¹⁵ Scintillation will be reduced if an aperture larger than the scale size is used to collect the light. If the aperture diameter D is much larger than the scale size, the reduction factor can be expressed as^{128,129}

$$A = C \left(\frac{D}{D_0} \right)^{-7/3} \quad (36)$$

where D_0 is $(\lambda L)^{1/2}$ and C is 4.71 for a plane wave and 23.5 for a spherical wave when the inner scale is much smaller than the Fresnel zone. If the inner scale is much larger than the Fresnel zone, D_0 is the inner scale and C is 0.45 for a plane wave and 9.17 for a spherical wave.

It is generally accepted that the probability density function for weak scintillation is log normal.^{101,113,115,130} This density function has the form

$$p(I) = \frac{1}{\sqrt{2\pi}\sigma_I I} \exp \left[-\frac{1}{2\sigma_I^2} (\ln I + 0.5\sigma_I^2)^2 \right] \quad (37)$$

The weak-turbulence theory is, in essence, a single-scattering theory. As the path-integrated turbulence becomes larger, multiple-scattering effects become important, and this theory breaks down. Actual observed values of irradiance variance are smaller than predicted in this region as shown in Fig. 32. The circles in this figure are 1-minute averages of irradiance variance for 488-nm laser light propagated across 1200 m of flat grassland. The solid line is the weak-turbulence approximation.

In these and other experiments, variance values reach a peak value of between 3 and 5 for a spherical wave and they begin to decrease with increasing turbulence strength.^{131–134} In the limit of infinite path-integrated turbulence, the normalized variance of irradiance is predicted to approach unity.^{135,136} In the intermediate region near the peak irradiance, numerical evaluation of currently available theories is impractical, and variance values are most easily obtained through numerical simulation.^{137–140}

In very strong turbulence two distinct spatial scales are evident in the scintillation pattern.^{136,141–144} The smaller scale is about the size of the coherence length ρ_0 discussed under “Beam Spreading.” The larger scale is the size of the scattering disk which is the ratio of the square of the Fresnel zone to the coherence length $(L/k\rho_0)$. As turbulence increases, the small scale becomes smaller and the large scale becomes larger. The strength of the small-scale fluctuations is constant in this regime and contributes a value of unity to the variance. The large-scale fluctuations contribute the rest of the variance and become weaker with increasing turbulence strength.

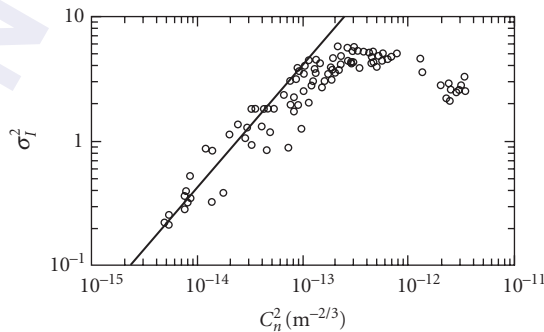


FIGURE 32 Plot of normalized variance of irradiance σ_1^2 as a function of turbulence strength. C_n^2 circles are data taken with a 488-nm wavelength laser over a 1200-m path, and the line is the corresponding weak-turbulence theory.

The two-scale nature of the irradiance fluctuations suggests a two-scale probability density function. The log-normally modulated Rician density function has been shown to agree with experimental data¹³⁰ and simulations.¹⁴⁵ It has the form

$$p(I) = \int_0^{\infty} dz \frac{(1+r)}{z} \exp\left(-r - \frac{1+r}{z} I\right) I_0 \left\{ 2 \left[\frac{(1+r)r}{z} \right]^{1/2} I \right\} \frac{1}{(2\pi)^{1/2} \sigma_z} \exp\left[-\frac{(\ln z + \frac{1}{2}\sigma_z^2)^2}{2\sigma_z^2}\right] \quad (38)$$

where r is a coherence parameter that decreases with increasing turbulence, I_0 is the 0-order modified Bessel function of the first kind, and σ_z^2 is the variance of the logarithm of the modulation factor z . For $r \gg 1$ this function reduces to the lognormal; for $r \ll 1$, it reduces to a log-normally modulated exponential:^{134,146}

$$p(I) = \frac{1}{\sqrt{2\pi}\sigma_z} \int_0^{\infty} \frac{dz}{z^2} \exp\left[-\frac{I}{z} - \frac{(\ln z + \frac{1}{2}\sigma_z^2)^2}{2\sigma_z^2}\right] \quad (39)$$

where the parameter σ_z^2 is related to the irradiance variance by the relationship

$$\sigma_I^2 = 2 \exp(\sigma_z^2) - 1 \quad (40)$$

As this model suggests, the density function of the fluctuations approaches a lognormal, even in strong turbulence, if an aperture $\gg \rho_0$ is used.¹⁴⁷

If the fluctuations at both large and small scales are approximated by gamma distributions, the resulting integral can be evaluated analytically to get the gamma-gamma density function:¹⁴⁸

$$p(I) = \frac{2(\alpha\beta)^{0.5(\alpha+\beta)}}{\Gamma(\alpha)\Gamma(\beta)} I^{0.5(\alpha+\beta)-1} K_{\alpha-\beta}(2\sqrt{\alpha\beta}I) \quad (41)$$

where α and β are related to the variances at the two scales, Γ is the gamma function, and K is the modified Bessel function of the second kind.

3.7 EXAMPLES OF ATMOSPHERIC OPTICAL REMOTE SENSING

One of the more important applications of atmospheric optics is optical remote sensing. Atmospheric optical remote sensing concerns the use of an optical or laser beam to remotely sense information about the atmosphere or a distant target. Optical remote sensing measurements are diverse in nature and include the use of a spectral radiometer aboard a satellite for the detection of trace species in the upper atmosphere, the use of spectral emission and absorption from the earth for the detection of the concentration of water vapor in the atmosphere, the use of lasers to measure the range-resolved distribution of several molecules including ozone in the atmosphere, and Doppler wind measurements. In this section, some typical optical remote sensing experiments will be presented in order to give a flavor of the wide variety of atmospheric optical measurements that are currently being conducted. More in-depth references can be found in several current journal papers, books, and conference proceedings.¹⁴⁹⁻¹⁵⁶

The Upper Atmospheric Research Satellite (UARS) was placed into orbit in September 1991 as part of the Earth Observing System. One of the optical remote sensing instruments aboard UARS is the High Resolution Doppler Imager (HRDI) developed by P. Hays' and V. Abreu's group while at the University of Michigan.^{157,158} The HRDI is a triple etalon Fabry-Perot Interferometer designed to measure Doppler shifts of molecular absorption and emission lines in the earth's atmosphere in order to determine the wind velocity of the atmosphere. A wind velocity of 10 m/s causes a Doppler shift of 2×10^{-5} nm for the oxygen lines detected near a wavelength of 600 to 800 nm. A schematic of the instrument is given in Fig. 33a which shows the telescope, triple Fabry-Perots, and unique imaging Photo-Multiplier tubes to detect the Fabry-Perot patterns of the spectral absorption lines. The HRDI instrument is a passive remote sensing system and uses the reflected or scattered sunlight as its illumination source. Figure 33b shows the wind field measured by UARS (HRDI) for an altitude of 90 km.

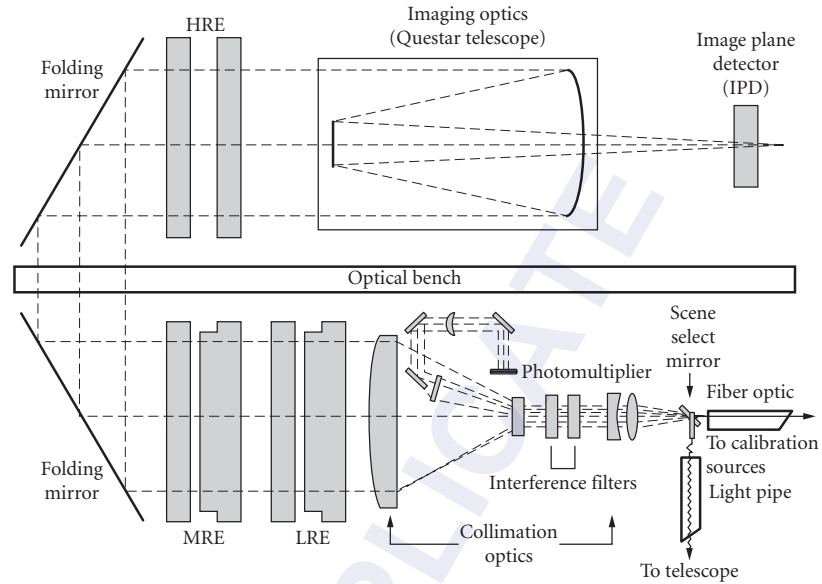


FIGURE 33a Optical layout of the Upper Atmospheric Resolution Satellite (UARS) High Resolution Doppler Imager (HRDI) instrument. FO = fiber optic, LRE = low-resolution etalon, MRE = medium-resolution etalon, HRE = high-resolution etalon. (From Hays, Ref. 157.)

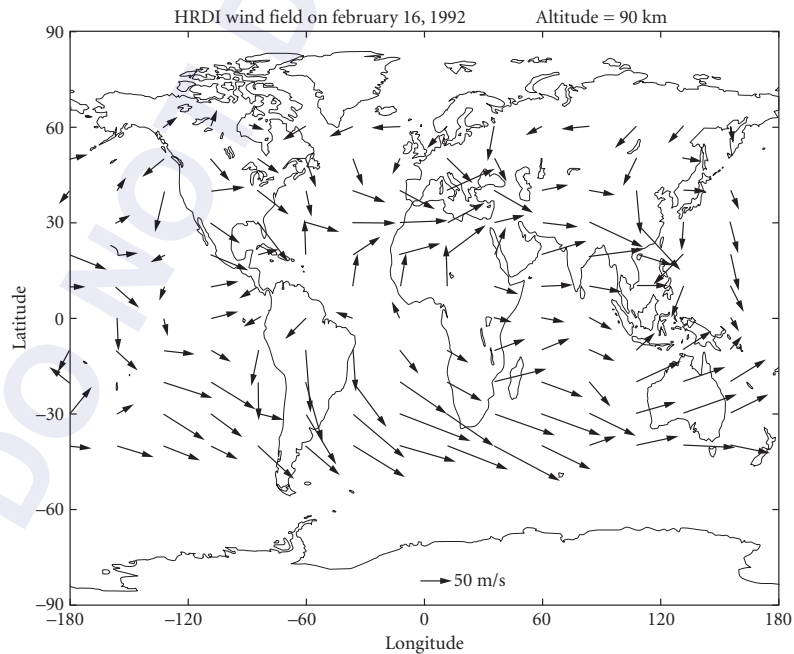


FIGURE 33b Upper atmospheric wind field measured by UARS/HRDI satellite instrument. (From Hays, Ref. 157.)

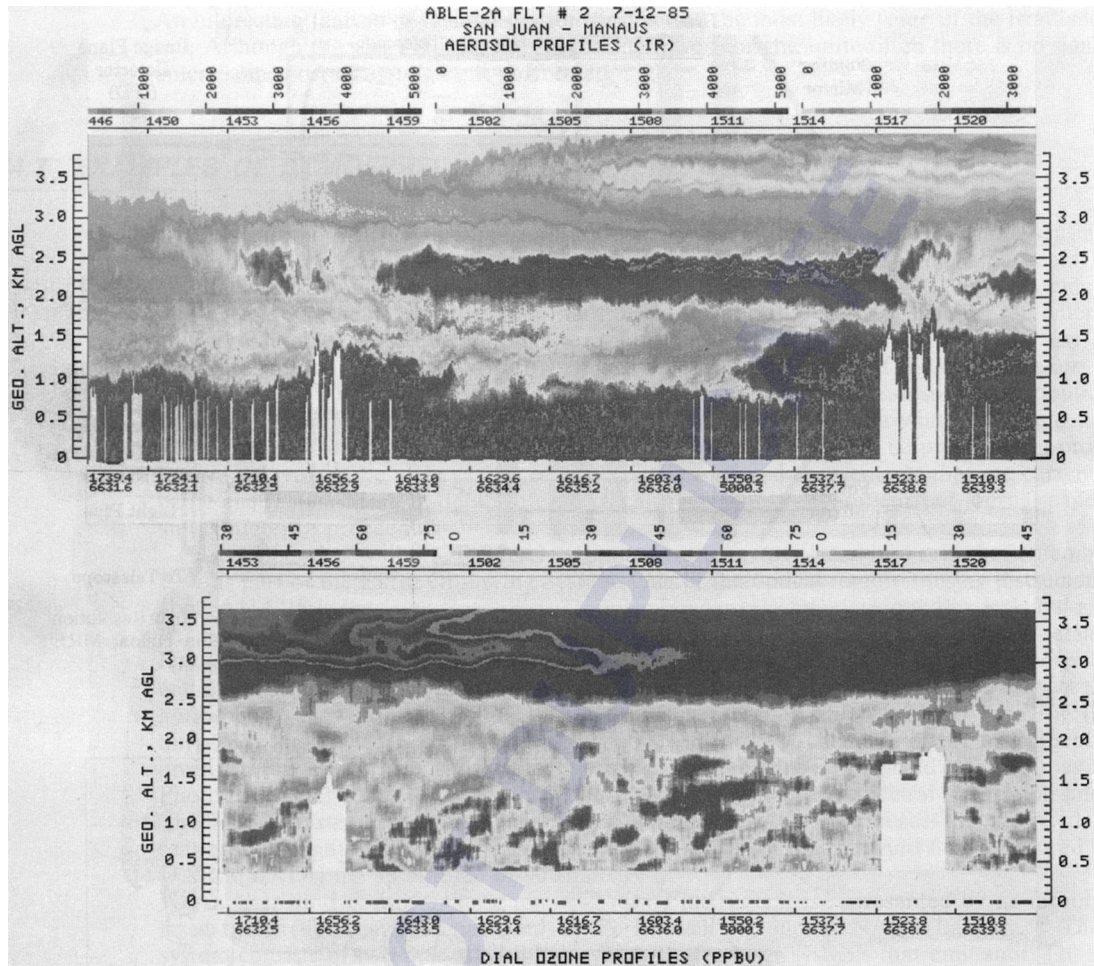


FIGURE 34 Range-resolved lidar measurements of atmospheric aerosols and ozone density. (From Browell, Ref. 159.)

Another kind of atmospheric remote sensing instrument is represented by an airborne laser radar (lidar) system operated by E. Browell's group at NASA/Langley.¹⁵⁹ Their system consists of two pulsed, visible-wavelength dye laser systems that emit short (10 ns) pulses of tunable optical radiation that can be directed toward aerosol clouds in the atmosphere. By the proper tuning of the wavelength of these lasers, the difference in the absorption due to ozone, water vapor, or oxygen in the atmosphere can be measured. Because the laser pulse is short, the timing out to the aerosol scatterers can be determined and range-resolved lidar measurements can be made. Figure 34 shows range-resolved lidar backscatter profiles obtained as a function of the lidar aircraft ground position. The variation in the atmospheric density and ozone distribution as a function of altitude and distance is readily observed.

A Coherent Doppler lidar is one which is able to measure the Doppler shift of the backscattered lidar returns from the atmosphere. Several Doppler lidar systems have been developed which can determine wind speed with an accuracy of 0.1 m/s at ranges of up to 15 km. One such system is operated by M. Hardesty's group at NOAA/WPL for the mapping of winds near airports and for meteorological studies.^{160,161} Figure 35 shows a two-dimensional plot of the measured wind velocity obtained during the approach of a wind gust front associated with colliding thunderstorms; the upper

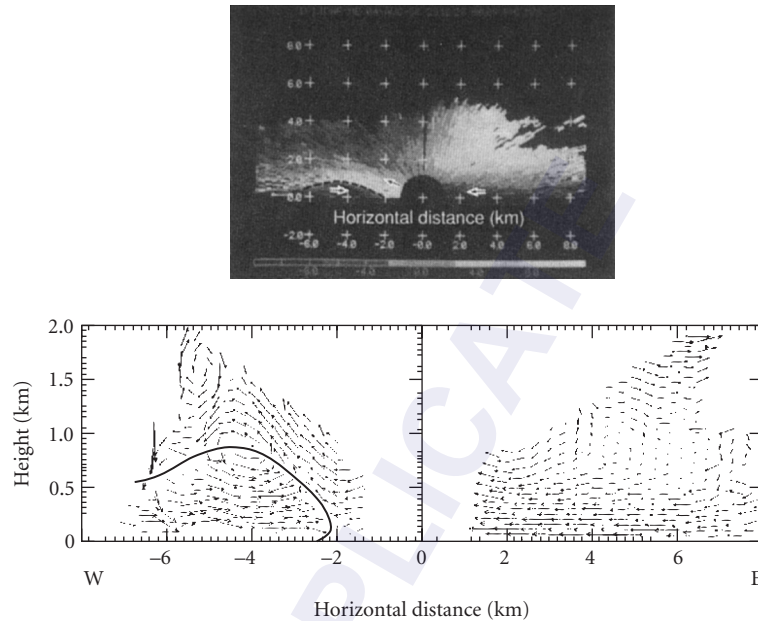


FIGURE 35 Coherent Doppler lidar measurements of atmospheric winds showing velocity profile of gust front. Upper plot is real-time display of Doppler signal and lower plot is range-resolved wind field. (From Hardesty, Ref. 160.)

figure shows the real-time Doppler lidar display of the measured radial wind velocity, and the lower plot shows the computed wind velocity. As seen, a Doppler lidar system is able to remotely measure the wind speed with spatial resolution on the order of 100 m. A similar Doppler lidar system is being considered for the early detection of windshear in front of commercial aircraft.

A further example of atmospheric optical remote sensing is that of the remote measurement of the global concentration and distribution of atmospheric aerosols. P. McCormick's group at NASA/Langley and Hampton University has developed the SAGE II satellite system which is part of a package of instruments to detect global aerosol and selected species concentrations in the atmosphere.¹⁶² This system measures the difference in the optical radiation emitted from the earth's surface and the differential absorption due to known absorption lines or spectral bands of several species in the atmosphere, including ozone. The instrument also provides for the spatial mapping of the concentration of aerosols in the atmosphere, and an example of such a measurement is shown in Fig. 36. This figure shows the measured concentration of aerosols after the eruption of Mt. Pinatubo and demonstrates the global circulation and transport of the injected material into the earth's atmosphere.

More recently, this capability has been refined by David Winker's group at NASA that developed the laser based CALIPSO lidar satellite which has produced continuous high-resolution 3D maps of global cloud and aerosol distributions since its launch in 2006.

There are several ongoing optical remote sensing programs to map and measure the global concentration of CO_2 and other green house gases in the atmosphere. For example, the spaceborne Atmospheric Infrared Sounder (AIRS) from JPL has measured the CO_2 concentration at the mid-troposphere (8 km altitude) beginning in 2003, and the NASA Orbiting Carbon Observatory (OCO) to be launched in 2008 will be the first dedicated spaceborne instrument to measure the sources and sinks of CO_2 globally. Both of these instruments use optical spectroscopy of atmospheric CO_2 lines to measure the concentration of CO_2 in the atmosphere.

Finally, there are related nonlinear optical processes that can also be used for remote sensing. For example, laser-induced-breakdown spectroscopy (LIBS) has been used recently in a lidar system for the remote detection of chemical species by focusing a pulsed laser beam at a remote target, producing a plasma spark at the

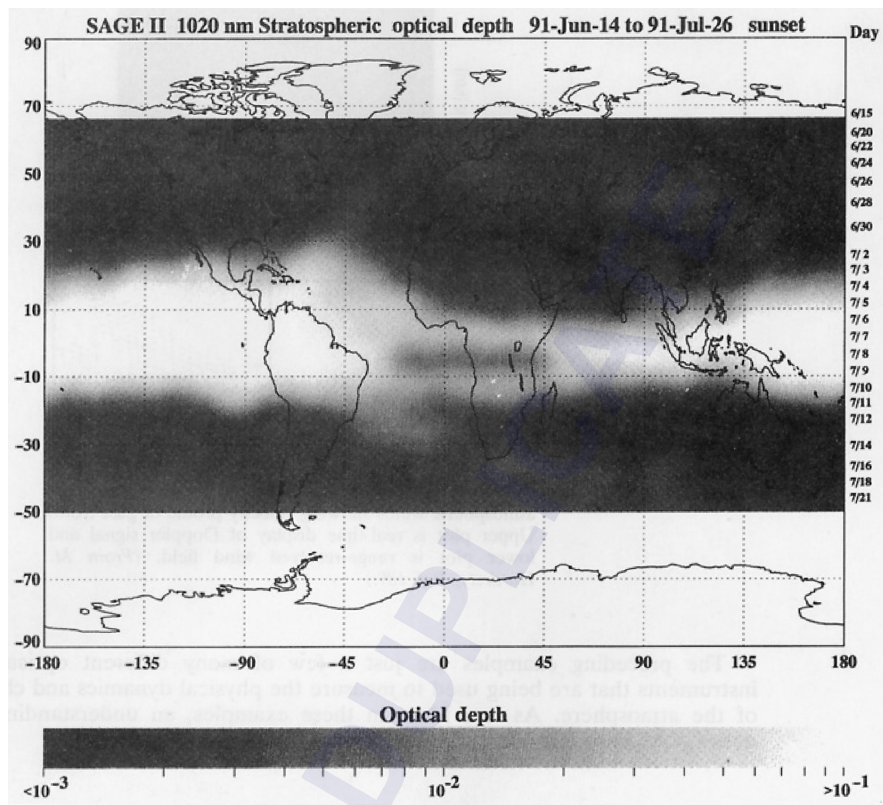


FIGURE 36 Measurement of global aerosol concentration using SAGE II satellite following eruption of Mt. Pinatubo. (From McCormick, Ref. 162.)

target, and analyzing the emitted spectral light after being transmitted back through the atmosphere.^{163,164} Another technique is the use of a high-power femtosecond pulse-length laser to produce a dielectric breakdown (spark) in air that self focuses into a long filament of several 100s of meters in length. The channeling of the laser filament has been used to remotely detect distant targets and atmospheric gases.¹⁶⁵

The preceding examples are just a few of many different optical remote sensing instruments that are being used to measure the physical dynamics and chemical properties of the atmosphere. As is evident in these examples, an understanding of atmospheric optics plays an important and integral part in these measurements.

3.8 METEOROLOGICAL OPTICS

One of the most colorful aspects of atmospheric optics is that associated with meteorological optics. Meteorological optics involves the interplay of light with the atmosphere and the physical origin of the observed optical phenomena. Several excellent books have been written about this subject, and the reader should consult these and the contained references.^{166,167} While it is beyond the scope of this chapter to present an overview of meteorological optics, some specific optical phenomena will be described to give the reader a sampling of some of the interesting effects involved in naturally occurring atmospheric and meteorological optics.

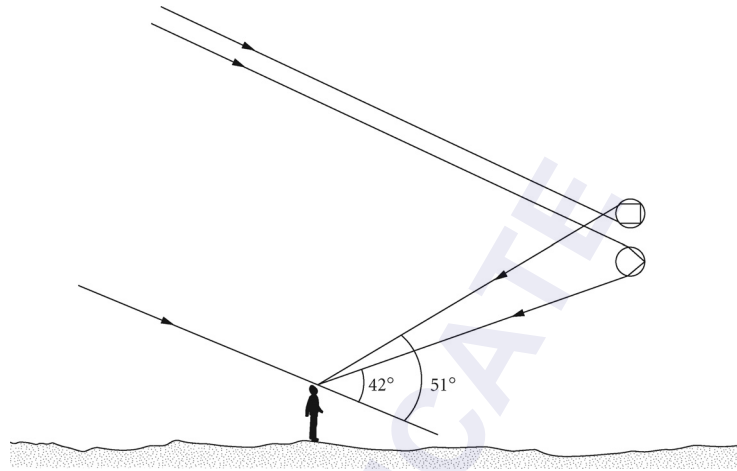


FIGURE 37 Different raindrops contribute to the primary and to the larger, secondary rainbow. (From Greenler, Ref. 167.)

Some of the more common and interesting meteorological optical phenomena involve rainbows, ice-crystal halos, and mirages. The rainbow in the atmosphere is caused by internal reflection and refraction of sunlight by water droplets in the atmosphere. Figure 37 shows the geometry involved in the formation of a rainbow, including both the primary and larger secondary rainbow. Because of the dispersion of light within the water droplet, the colors or wavelengths are separated in the backscattered image. Although rainbows are commonly observed in the visible spectrum, such refraction also occurs in the infrared spectrum. As an example, Fig. 38 shows a natural rainbow in the atmosphere photographed with IR-sensitive film by R. Greenler.¹⁶⁷



FIGURE 38 A natural infrared rainbow. (From Greenler, Ref. 167.)

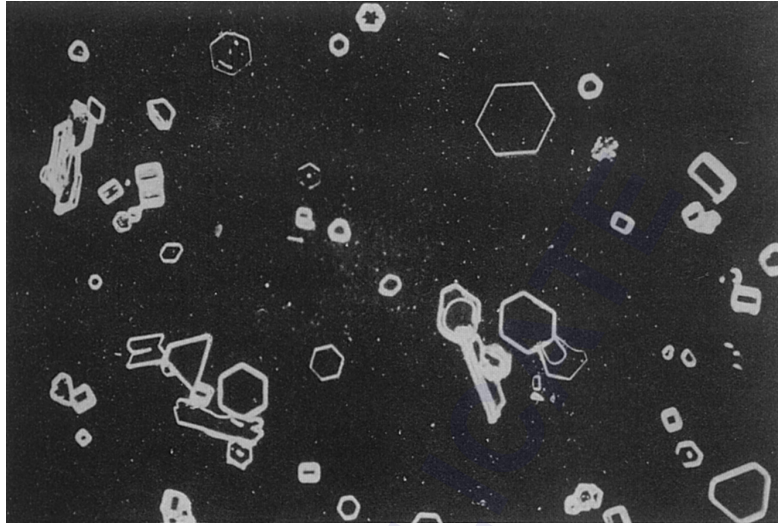


FIGURE 39 Photograph of magnified small ice crystals collected as they fell from the sky. (From Greenler, Ref. 167.)

The phenomena of halos, arcs, and spots are due to the refraction of light by ice crystals suspended in the atmosphere. Figure 39 shows a photograph of collected ice crystals as they fell from the sky. The geometrical shapes, especially the hexagonal (six-sided) crystals, play an important role in the formation of halos and arcs in the atmosphere.

The common optical phenomenon of the mirage is caused by variation in the temperature and thus, the density of the air as a function of altitude or transverse geometrical distance. As an example, Fig. 40 shows the geometry of light-ray paths for a case where the air temperature decreases with height

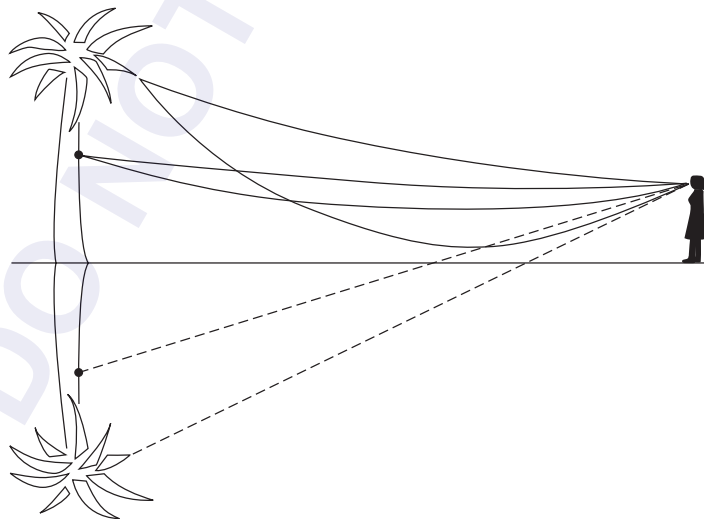


FIGURE 40 The origin of the inverted image in the desert mirage. (From Greenler, Ref. 167.)



FIGURE 41 The desert (or hot-road) mirage. In the inverted part of the image you can see the apparent reflection of motorcycles, cars, painted stripes on the road, and the grassy road edge. (From Greenler, Ref. 167.)

to a sufficient extent over the viewing angle that the difference in the index of refraction can cause a refraction of the image similar to total internal reflection. The heated air (less dense) near the ground can thus act like a mirror, and reflect the light upward toward the viewer. As an example, Fig. 41 shows a photograph taken by Greenler of motorcycles on a hot road surface. The reflected image of the motorcycles “within” the road surface is evident. There are many manifestations of mirages dependent upon the local temperature gradient and geometry of the situation. In many cases, partial and distorted images are observed leading to the almost surreal connotation often associated with mirages.

Finally, another atmospheric meteorological optical phenomenon is that of the green flash. A green flash is observed under certain conditions just as the sun is setting below the horizon. This phenomenon is easily understood as being due to the different relative displacement of each different wavelength or color in the sun’s image due to spatially distributed refraction of the atmosphere.¹⁶⁷ As the sun sets, the last image to be observed is the shortest wavelength color, blue. However, most of the blue light has been Rayleigh scattered from the image seen by the observer so that the last image observed is closer to a green color. Under extremely clear atmospheric conditions when the Rayleigh scattering is not as preferential in scattering the blue light, the flash has been reported as blue in color. Lastly, one of the authors (DK) has observed several occurrences of the green flash and noticed that the green flash seems to be seen more often from sunsets over water than over land, suggesting that a form of water vapor layer induced ducting of the optical beam along the water’s surface may be involved in enhancing the absorption and scattering process.

3.9 ATMOSPHERIC OPTICS AND GLOBAL CLIMATE CHANGE

Of importance, atmospheric optics largely determines the earth’s climate because incoming sunlight (optical energy) is scattered and absorbed and outgoing thermal radiation is absorbed and reemitted by the atmosphere. This net energy flux, either incoming or outgoing, determines if the earth’s

climate will warm or cool, and is directly related to the radiative transfer calculations mentioned in Sec. 3.4. A convenient way to express changes in this energy balance is in terms of the radiative forcing, which is defined as the change in the incoming and outgoing radiative flux at the top of the troposphere. To a good approximation, individual forcing terms add linearly to produce a linear surface-temperature response at regional to global scales. Reference 168 by the Intergovernmental Panel on Climate Change (2007) provides a detailed description of radiative forcing and the scientific basis for these findings.

As an example, Fig. 42 shows that the major radiative forcing mechanisms for changes between the years of 1750 and 2005 are the result of anthropogenic changes to the atmosphere.¹⁶⁸ As can be seen, the largest effect is that of infrared absorption by the greenhouse gases, principally CO_2 . This is the reason that the measurement of the global concentration of CO_2 is important for accurate predictions for future warming or cooling trends of the earth. Other forcing effects are also shown in Fig. 42. Water vapor remains the largest absorber of infrared radiation, but changes in water vapor caused

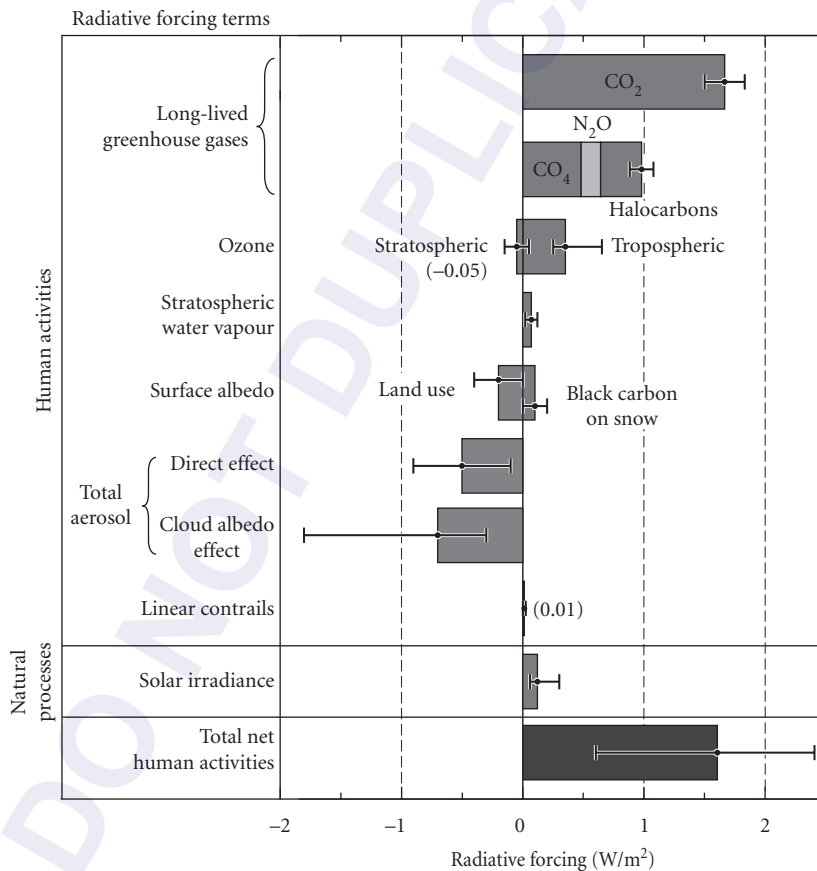


FIGURE 42 Effect of atmospheric optics on global climate change represented by the radiative forcing terms between the years of 1750 and 2005. The change in the energy (W/m^2) balance for the earth over this time period is shown resulting in a net positive energy flow onto the earth, with a potential warming effect. Contributions due to individual terms, such as changes in the CO_2 gas concentration or changes in land use, are shown. (Reprinted with permission from Ref. 168.)

by climate changes are considered a response rather than a forcing term. The exception is a small increase in stratospheric water vapor produced by methane emissions. Similarly, the effects of clouds are part of a climate response, except for the increase in cloudiness that is a direct result of increases in atmospheric aerosols. The linearity of climate response to radiative forcing implies that the most efficient radiative transfer calculations for each term can be used. For example, a high resolution spectral absorption calculation is not needed to calculate the radiative transfer through clouds, and a multiple-scattering calculation is not needed to calculate the effects of absorption by gases.

In summary, Fig. 42 shows the important contributions of radiative forcing effects that contribute to changes in the heat balance of the earth's atmosphere. All of the noted terms are influenced by the optical properties of the atmosphere whether due to absorption of sunlight by the line or band spectrum of molecules in the air, the reflection of light by the earth's surface or oceans, or reabsorption of thermal radiation by greenhouse gases. The interested reader is encouraged to study Ref. 168 for more information, and references therein.

3.10 ACKNOWLEDGMENTS

We would like to acknowledge the contributions and help received in the preparation of this chapter and in the delineation of the authors' work. The authors divided the writing of the chapter sections as follows: D. K. Killinger served as lead author and wrote Secs. 3.2 through 3.4 and Secs. 3.7 and 3.8 on atmospheric interactions with light, remote sensing, and meteorological optics. L. S. Rothman wrote and provided the extensive background information on HITRAN, FASCODE, and LOWTRAN in Sec. 3.5. The comprehensive Secs. 3.6 and 3.9 were written by J. H. Churnside. The data of Fig. 29 were provided by G. R. Ochs of NOAA/WPL and the data in Fig. 30 were provided by R. R. Beland of the Geophysics Directorate, Phillips Laboratory. We wish to thank Prof. Robert Greenler for providing original photographs of the meteorological optics phenomena; Paul Hays, Vincent Abreu, and Wilbert Skinner for information on the HRDI instrument; P. McCormick and D. Winker for SAGE II data; Mike Hardesty for Doppler lidar wind profiles; and Ed Browell for lidar ozone mapping data. We want to thank A. Jursa for providing a copy of the *Handbook of Geophysics*. R. Measures for permission to use diagrams from his book *Laser Remote Sensing*, and M. Thomas and D. Duncan for providing a copy of their chapter on atmospheric optics from *The Infrared Handbook*.

Finally, we wish to thank many of our colleagues who have suggested topics and technical items added to this work. We hope that the reader will gain an overall feeling of atmospheric optics from reading this chapter, and we encourage the reader to use the references cited for further in-depth study.

3.11 REFERENCES

1. R. M. Goody and Y. L. Young, *Atmospheric Radiation*, Oxford University Press, London, 1989.
2. W. G. Driscoll (ed.), Optical Society of America, *Handbook of Optics*, McGraw-Hill, New York, 1978.
3. A. S. Jursa (ed.), *Handbook of Geophysics and the Space Environment*, Air Force Geophysics Lab., NTIS Doc#ADA16700, 1985.
4. W. Wolfe and G. Zissis, *The Infrared Handbook*, Office of Naval Research, Washington D.C., 1978.
5. R. Measures, *Laser Remote Sensing*, Wiley-Interscience, John Wiley & Sons, New York, 1984.
6. "Major Concentration of Gases in the Atmosphere," NOAA S/T 76-1562, 1976; "AFGL Atmospheric Constituent Profiles (0-120 km)," AFGL-TR-86-0110, 1986; U.S. Standard Atmosphere, 1962 and 1976; Supplement 1966, U.S. Printing Office, Washington D.C., 1976.
7. E. P. Shettle and R. W. Fenn, "Models of the Aerosols of the Lower Atmosphere and the Effects of Humidity Variations on Their Optical Properties," AFGL TR-79-0214; ADA 085951, 1979.
8. A. Force, D. K. Killinger, W. DeFeo, and N. Menyuk, "Laser Remote Sensing of Atmospheric Ammonia Using a CO₂ Lidar System," *Appl. Opt.* **24**:2837 (1985).

9. B. Nilsson, "Meteorological Influence on Aerosol Extinction in the 0.2–40 μm Range," *Appl. Opt.* **18**:3457 (1979).
10. J. F. Luhr, "Volcanic Shade Causes Cooling," *Nature* **354**:104 (1991).
11. L. S. Rothman, R. R. Gamache, A. Goldman, L. R. Brown, R. A. Toth, H. Pickett, R. Poynter, et al., "The HITRAN Database: 1986 Edition," *Appl. Optics* **26**:4058 (1986).
12. L. S. Rothman, R. R. Gamache, R. H. Tipping, C. P. Rinsland, M. A. H. Smith, D. C. Benner, V. Malathy Devi, et al., "The HITRAN Molecular Database: Editions of 1991 and 1992," *J. Quant. Spectrosc. Radiat. Transfer* **48**:469 (1992).
13. L. S. Rothman, I. E. Gordan, A. Barbe, D. Chris Benner, P. F. Bernath, M. Birk, L. R. Brown, et al., "The HITRAN 2008 Molecular Spectroscopic Database," *J. Quant. Spectrosc. Radiat. Transfer* **110**:533 (2009).
14. E. E. Whiting, "An Empirical Approximation to the Voigt Profile," *J. Quant. Spectrosc. Radiat. Transfer* **8**:1379 (1968).
15. J. Olivero and R. Longbothum, "Empirical Fits to the Voigt Linewidth: A Brief Review," *J. Quant. Spectrosc. Radiat. Transfer* **17**:233 (1977).
16. F. Schreir, "The Voigt and Complex Error Function—A Comparison of Computational Methods," *J. Quant. Spectrosc. Radiat. Transfer* **48**:734 (1992).
17. J. -M. Hartmann, C. Boulet, and D. Robert, *Collision Effects on Molecular Spectra. Laboratory Experiments and Models, Consequences for Applications*, Elsevier, Paris, 2008.
18. D. E. Burch, "Continuum Absorption by H_2O ," AFGL-TR-81-0300; ADA 112264, 1981; S. A. Clough, F. X. Kneizys, and R. W. Davies, "Lineshape and the Water Vapor Continuum," *Atmospheric Research* **23**:229 (1989).
19. Shardanand and A. D. Prasad Rao, "Absolute Rayleigh Scattering Cross Section of Gases and Freons of Stratospheric Interest in the Visible and Ultraviolet Region," NASA TN 0-8442 (1977); R. T. H. Collins and P. B. Russell, "Lidar Measurement of Particles and Gases by Elastic and Differential Absorption," in D. Hinkley (ed.), *Laser Monitoring of the Atmosphere*, Springer-Verlag, New York, 1976.
20. E. U. Condon and H. Odishaw (eds.), *Handbook of Physics*, McGraw-Hill, New York, 1967.
21. S. R. Pal and A. I. Carswell, "Polarization Properties of Lidar Backscattering from Clouds," *Appl. Opt.* **12**:1530 (1973).
22. G. Mie, "Bertrage Z. Phys. TruberMedien, Speziell kolloidaler Metallosungen," *Ann. Physik* **25**:377 (1908).
23. D. Deirndjian, "Scattering and Polarization Properties of Water Clouds and Hazes in the Visible and Infrared," *Appl. Opt.* **2**:187 (1964).
24. E. J. McCartney, *Optics of the Atmosphere*, Wiley, New York, 1976.
25. M. Wright, E. Proctor, L. Gasiorek, and E. Liston, "A Preliminary Study of Air Pollution Measurement by Active Remote Sensing Techniques," NASA CR-132724, 1975.
26. P. McCormick and D. Winker, "NASA/LaRC: 1 μm Lidar Measurement of Aerosol Distribution," Private communication, 1991.
27. D. K. Killinger and N. Menyuk, "Laser Remote Sensing of the Atmosphere," *Science* **235**:37 (1987).
28. R. W. Boyd, *Nonlinear Optics*, Academic Press, Orlando, Fla., 1992.
29. M. D. Levenson and S. Kano, *Introduction to Nonlinear Laser Spectroscopy*, Academic Press, Boston, 1988.
30. S. A. Clough, F. X. Kneizys, E. P. Shettle, and G. P. Anderson, "Atmospheric Radiance and Transmittance: FASCOD2," *Proc. of Sixth Conf. on Atmospheric Radiation*, Williamsburg, Va., published by Am. Meteorol. Soc., Boston, 1986.
31. J. A. Dowling, W. O. Gallery, and S. G. O'Brian, "Analysis of Atmospheric Interferometer Data," AFGL-TR-84-0177, 1984.
32. R. Isaacs, S. Clough, R. Worsham, J. Moncet, B. Lindner, and L. Kaplan, "Path Characterization Algorithms for FASCOD2," Tech. Report GL-TR-90-0080, AFGL, 1990; ADA#231914.
33. F. X. Kneizys, E. Shettle, W. O. Gallery, J. Chetwynd, L. Abreu, J. Selby, S. Clough, and R. Fenn, "Atmospheric Transmittance/Radiance: Computer Code LOWTRAN6," AFGL TR-83-0187, 1983; ADA#137786.
34. F. X. Kneizys, E. Shettle, L. Abreu, J. Chetwynd, G. Anderson, W. O. Gallery, J. E. A. Selby, and S. Clough, "Users Guide to LOWTRAN7," AFGL TR-88-0177, 1988; ADA#206773.
35. The HITRAN database compilation is available on the internet and can be accessed by filling out a request form at the HITRAN web site <http://cfa.harvard.edu/hitrان>. The MODTRAN code (with LOWTRAN7 embedded within), and a PC version of FASCOD2, can be obtained from the ONTAR Corp, 9 Village Way, North Andover, MA 01845-2000.

36. D. K. Killinger and W. Wilcox, Jr., HITRAN-PC Program; can be obtained from ONTAR Corp. at www.ontar.com.
37. NIST/EPA Gas Phase Infrared Database, U.S. Dept. of Commerce, NIST, Standard Ref. Data, Gaithersburg, MD 20899.
38. Vapor phase infrared spectral library, Pacific Northwest National Laboratory, Richland, WA 99352.
39. LAB_CALC, Galactic Industries, 395 Main St., Salem, NH, 03079 USA; Infrared Analytics, 1424 N. Central Park Ave, Anaheim, CA 92802; Aldrich Library of Spectra, Aldrich Co., Milwaukee, WI 53201; Sadtler Spectra Data, Philadelphia, PA 19104-2596; Coblenz Society, P.O. Box 9952, Kirkwood, MO 63122.
40. L. C. Andrews and R. L. Phillips, *Laser Beam Propagation through Random Media*, 2nd ed., SPIE Press, Washington, 2005.
41. R. J. Hill, "Models of the Scalar Spectrum for Turbulent Advection," *J. Fluid Mech.* **88**:541–662 (1978).
42. J. H. Churnside, "A Spectrum of Refractive-Index Turbulence in the Turbulent Atmosphere," *J. Mod. Opt.* **37**:13–16 (1990).
43. L. C. Andrews, "An Analytical Model for the Refractive Index Power Spectrum and Its Application to Optical Scintillations in the Atmosphere," *J. Mod. Opt.* **39**:1849–1853 (1992).
44. R. S. Lawrence, G. R. Ochs, and S. F. Clifford, "Measurements of Atmospheric Turbulence Relevant to Optical Propagation," *J. Opt. Soc. Am.* **60**:826–830 (1970).
45. M. A. Kallistratova and D. F. Timanovskiy, "The Distribution of the Structure Constant of Refractive Index Fluctuations in the Atmospheric Surface Layer," *Iz. Atmos. Ocean. Phys.* **7**:46–48 (1971).
46. A. S. Monin and A. M. Obukhov, "Basic Laws of Turbulent Mixing in the Ground Layer of the Atmosphere," *Trans. Geophys. Inst. Akad. Nauk. USSR* **151**:163–187 (1954).
47. A. S. Monin and A. M. Yaglom, *Statistical Fluid Mechanics: Mechanics of Turbulence*, MIT Press, Cambridge, 1971.
48. J. C. Wyngaard, Y. Izumi, and S. A. Collins, Jr., "Behavior of the Refractive-Index-Structure Parameter Near the Ground," *J. Opt. Soc. Am.* **61**:1646–1650 (1971).
49. L. R. Tsvang, "Microstructure of Temperature Fields in the Free Atmosphere," *Radio Sci.* **4**:1175–1177 (1969).
50. A. S. Frisch and G. R. Ochs, "A Note on the Behavior of the Temperature Structure Parameter in a Convective Layer Capped by a Marine Inversion," *J. Appl. Meteorol.* **14**:415–419 (1975).
51. K. L. Davidson, T. M. Houlihan, C. W. Fairall, and G. E. Schader, "Observation of the Temperature Structure Function Parameter, C_T^2 , over the Ocean," *Boundary-Layer Meteorol.* **15**:507–523 (1978).
52. K. E. Kunkel, D. L. Walters, and G. A. Ely, "Behavior of the Temperature Structure Parameter in a Desert Basin," *J. Appl. Meteorol.* **15**:130–136 (1981).
53. W. Kohsiek, "Measuring C_T^2 , C_Q^2 , and C_{TQ} in the Unstable Surface Layer, and Relations to the Vertical Fluxes of Heat and Moisture," *Boundary-Layer Meteorol.* **24**:89–107 (1982).
54. M. S. Belen'kiy, V. V. Boronyev, N. Ts. Gomboyev, and V. L. Mironov, *Sounding of Atmospheric Turbulence*, Nauka, Novosibirsk, p. 114, 1986.
55. A. A. M. Holtslag and A. P. Van Ulden, "A Simple Scheme for Daytime Estimates of the Surface Fluxes from Routine Weather Data," *J. Clim. Appl. Meteorol.* **22**:517–529 (1983).
56. T. Thiermann and A. Kohnle, "A Simple Model for the Structure Constant of Temperature Fluctuations in the Lower Atmosphere," *J. Phys. D: Appl. Phys.* **21**:S37–S40 (1988).
57. E. A. Andreas, "Estimating C_n^2 over Snow and Sea Ice from Meteorological Data," *J. Opt. Soc. Am. A* **5**:481–494 (1988).
58. J. L. Bufton, P. O. Minott, and M. W. Fitzmaurice, "Measurements of Turbulence Profiles in the Troposphere," *J. Opt. Soc. Am.* **62**:1068–1070 (1972).
59. F. W. Eaton, W. A. Peterson, J. R. Hines, K. R. Peterman, R. E. Good, R. R. Beland, and J. W. Brown, "Comparisons of VHF Radar, Optical, and Temperature Fluctuation Measurements of C_n^2 , r_θ , and Θ_θ ," *Theor. Appl. Climatol.* **39**:17–29 (1988).
60. F. Dalaudier, M. Crochet, and C. Sidi, "Direct Comparison between in situ and Radar Measurements of Temperature Fluctuation Spectra: A Puzzling Result," *Radio Sci.* **24**:311–324 (1989).
61. D. W. Beran, W. H. Hooke, and S. F. Clifford, "Acoustic Echo-Sounding Techniques and Their Application to Gravity-Wave, Turbulence, and Stability Studies," *Boundary-Layer Meteorol.* **4**:133–153 (1973).
62. M. Fukushima, K. Akita, and H. Tanaka, "Night-Time Profiles of Temperature Fluctuations Deduced from Two-Year Solar Observation," *J. Meteorol. Soc. Jpn.* **53**:487–491 (1975).

63. D. N. Asimakopoulis, R. S. Cole, S. J. Caughey, and B. A. Crease, "A Quantitative Comparison between Acoustic Sounder Returns and the Direct Measurement of Atmospheric Temperature Fluctuations," *Boundary-Layer Meteorol.* **10**:137–147 (1976).
64. T. E. VanZandt, J. L. Green, K. S. Gage, and W. L. Clark, "Vertical Profiles of Refractivity Turbulence Structure Constant: Comparison of Observations by the Sunset Radar with a New Theoretical Model," *Radio Sci.* **13**:819–829 (1978).
65. K. S. Gage and B. B. Balsley, "Doppler Radar Probing of the Clear Atmosphere," *Bull. Am. Meteorol. Soc.* **59**:1074–1093 (1978).
66. R. B. Chadwick and K. P. Moran, "Long-Term measurements of C_n^2 in the Boundary Layer," *Radio Sci.* **15**:355–361 (1980).
67. B. B. Balsley and V. L. Peterson, "Doppler-Radar Measurements of Clear Air Turbulence at 1290 MHz," *J. Appl. Meteorol.* **20**:266–274 (1981).
68. E. E. Gossard, R. B. Chadwick, T. R. Detman, and J. Gaynor, "Capability of Surface-Based Clear-Air Doppler Radar for Monitoring Meteorological Structure of Elevated Layers," *J. Clim. Appl. Meteorol.* **23**:474 (1984).
69. G. D. Nastrom, W. L. Ecklund, K. S. Gage, and R. G. Strauch, "The Diurnal Variation of Backscattered Power from VHF Doppler Radar Measurements in Colorado and Alaska," *Radio Sci.* **20**:1509–1517 (1985).
70. D. L. Fried, "Remote Probing of the Optical Strength of Atmospheric Turbulence and of Wind Velocity," *Proc. IEEE* **57**:415–420 (1969).
71. J. W. Strohbehn, "Remote Sensing of Clear-air Turbulence," *J. Opt. Soc. Am.* **60**:948 (1970).
72. J. Vernin and F. Roddier, "Experimental Determination of Two-Dimensional Power Spectra of Stellar Light Scintillation. Evidence for a Multilayer Structure of the Air Turbulence in the Upper Troposphere," *J. Opt. Soc. Am.* **63**:270–273 (1973).
73. G. R. Ochs, T. Wang, R. S. Lawrence, and S. F. Clifford, "Refractive Turbulence Profiles Measured by One-Dimensional Spatial Filtering of Scintillations," *Appl. Opt.* **15**:2504–2510 (1976).
74. R. E. Hufnagel and N. R. Stanley, "Modulation Transfer Function Associated with Image Transmission through Turbulent Media," *J. Opt. Soc. Am.* **54**:52–61 (1964).
75. R. E. Hufnagel, "Variations of Atmospheric Turbulence," in *Technical Digest of Topical Meeting on Optical Propagation through Turbulence*, Optical Society of America, Washington, D.C., 1974.
76. R. J. Sasiela, *A Unified Approach to Electromagnetic Wave Propagation in Turbulence and the Evaluation of Multiparameter Integrals*, Technical Report 807, MIT Lincoln Laboratory, Lexington, 1988.
77. V. A. Banakh and V. L. Mironov, *Lidar in a Turbulent Atmosphere*, Artech House, Boston, 1987.
78. C. W. Fairall and R. Markson, "Aircraft Measurements of Temperature and Velocity Microturbulence in the Stably Stratified Free Troposphere," *Proceedings of the Seventh Symposium on Turbulence and Diffusion*, November 12–15, Boulder, Co. 1985.
79. J. C. Kaimal, *The Atmospheric Boundary Layer—Its Structure and Measurement*, Indian Institute of Tropical Meteorology, Pune, 1988.
80. J. Barat and F. Bertin, "On the Contamination of Stratospheric Turbulence Measurements by Wind Shear," *J. Atmos. Sci.* **41**:819–827 (1984).
81. A. Ziad, R. Conan, A. Tokovinin, F. Martin, and J. Borgnino, "From the Grating Scale Monitor to the Generalized Seeing Monitor," *Appl. Opt.* **39**:5415–5425 (2000).
82. A. Ziad, M. Schöck, G. A. Chanan, M. Troy, R. Dekany, B. F. Lane, J. Borgnino, and F. Martin, "Comparison of Measurements of the Outer Scale of Turbulence by Three Different Techniques," *Appl. Opt.* **43**:2316–2324 (2004).
83. L. A. Chernov, *Wave Propagation in a Random Medium*, Dover, New York, p. 26, 1967.
84. P. Beckmann, "Signal Degeneration in Laser Beams Propagated through a Turbulent Atmosphere," *Radio Sci.* **69D**:629–640 (1965).
85. T. Chiba, "Spot Dancing of the Laser Beam Propagated through the Atmosphere," *Appl. Opt.* **10**:2456–2461 (1971).
86. J. H. Churnside and R. J. Latatis, "Wander of an Optical Beam in the Turbulent Atmosphere," *Appl. Opt.* **29**:926–930 (1990).
87. G. A. Andreev and E. I. Gelfer, "Angular Random Walks of the Center of Gravity of the Cross Section of a Diverging Light Beam," *Radiophys. Quantum Electron.* **14**:1145–1147 (1971).

88. M. A. Kallistratova and V. V. Pokasov, "Defocusing and Fluctuations of the Displacement of a Focused Laser Beam in the Atmosphere," *Radiophys. Quantum Electron.* **14**:940–945 (1971).
89. J. A. Dowling and P. M. Livingston, "Behavior of Focused Beams in Atmospheric Turbulence: Measurements and Comments on the Theory," *J. Opt. Soc. Am.* **63**:846–858 (1973).
90. J. R. Dunphy and J. R. Kerr, "Turbulence Effects on Target Illumination by Laser Sources: Phenomenological Analysis and Experimental Results," *Appl. Opt.* **16**:1345–1358 (1977).
91. V. I. Klyatskin and A. I. Kon, "On the Displacement of Spatially Bounded Light Beams in a Turbulent Medium in the Markovian-Random-Process Approximation," *Radiophys. Quantum Electron.* **15**:1056–1061 (1972).
92. A. I. Kon, V. L. Mironov, and V. V. Nosov, "Dispersion of Light Beam Displacements in the Atmosphere with Strong Intensity Fluctuations," *Radiophys. Quantum Electron.* **19**:722–725 (1976).
93. V. L. Mironov and V. V. Nosov, "On the Theory of Spatially Limited Light Beam Displacements in a Randomly Inhomogeneous Medium," *J. Opt. Soc. Am.* **67**:1073–1080 (1977).
94. R. F. Lutomirski and H. T. Yura, "Propagation of a Finite Optical Beam in an Inhomogeneous Medium," *Appl. Opt.* **10**:1652–1658 (1971).
95. R. F. Lutomirski and H. T. Yura, "Wave Structure Function and Mutual Coherence Function of an Optical Wave in a Turbulent Atmosphere," *J. Opt. Soc. Am.* **61**:482–487 (1971).
96. H. T. Yura, "Atmospheric Turbulence Induced Laser Beam Spread," *Appl. Opt.* **10**:2771–2773 (1971).
97. H. T. Yura, "Mutual Coherence Function of a Finite Cross Section Optical Beam Propagating in a Turbulent Medium," *Appl. Opt.* **11**:1399–1406 (1972).
98. H. T. Yura, "Optical Beam Spread in a Turbulent Medium: Effect of the Outer Scale of Turbulence," *J. Opt. Soc. Am.* **63**:107–109 (1973).
99. H. T. Yura, "Short-Term Average Optical-Beam Spread in a Turbulent Medium," *J. Opt. Soc. Am.* **63**:567–572 (1973).
100. M. T. Tavis and H. T. Yura, "Short-Term Average Irradiance Profile of an Optical Beam in a Turbulent Medium," *Appl. Opt.* **15**:2922–2931 (1976).
101. R. L. Fante, "Electromagnetic Beam Propagation in Turbulent Media," *Proc. IEEE* **63**:1669–1692 (1975).
102. R. L. Fante, "Electromagnetic Beam Propagation in Turbulent Media: An Update," *Proc. IEEE* **68**:1424–1443 (1980).
103. G. C. Valley, "Isoplanatic Degradation of Tilt Correction and Short-Term Imaging Systems," *Appl. Opt.* **19**:574–577 (1980).
104. H. J. Breaux, *Correlation of Extended Huygens-Fresnel Turbulence Calculations for a General Class of Tilt Corrected and Uncorrected Laser Apertures*, Interim Memorandum Report No. 600, U.S. Army Ballistic Research Laboratory, 1978.
105. D. M. Cordray, S. K. Searles, S. T. Hanley, J. A. Dowling, and C. O. Gott, "Experimental Measurements of Turbulence Induced Beam Spread and Wander at 1.06, 3.8, and 10.6 μm ," *Proc. SPIE* **305**:273–280 (1981).
106. S. K. Searles, G. A. Hart, J. A. Dowling, and S. T. Hanley, "Laser Beam Propagation in Turbulent Conditions," *Appl. Opt.* **30**:401–406 (1991).
107. J. H. Churnside and R. J. Lataitis, "Angle-of-Arrival Fluctuations of a Reflected Beam in Atmospheric Turbulence," *J. Opt. Soc. Am. A* **4**:1264–1272 (1987).
108. D. L. Fried, "Optical Resolution through a Randomly Inhomogeneous Medium for Very Long and Very Short Exposures," *J. Opt. Soc. Am.* **56**:1372–1379 (1966).
109. R. F. Lutomirski, W. L. Woodie, and R. G. Buser, "Turbulence-Degraded Beam Quality: Improvement Obtained with a Tilt-Correcting Aperture," *Appl. Opt.* **16**:665–673 (1977).
110. D. L. Fried, "Statistics of a Geometrical Representation of Wavefront Distortion," *J. Opt. Soc. Am.* **55**:1427–1435 (1965); **56**:410 (1966).
111. D. M. Chase, "Power Loss in Propagation through a Turbulent Medium for an Optical-Heterodyne System with Angle Tracking," *J. Opt. Soc. Am.* **56**:33–44 (1966).
112. J. H. Churnside and C. M. McIntyre, "Partial Tracking Optical Heterodyne Receiver Arrays," *J. Opt. Soc. Am.* **68**:1672–1675 (1978).
113. V. I. Tatarskii, *The Effects of the Turbulent Atmosphere on Wave Propagation*, Israel Program for Scientific Translations, Jerusalem, 1971.
114. R. W. Lee and J. C. Harp, "Weak Scattering in Random Media, with Applications to Remote Probing," *Proc. IEEE* **57**:375–406 (1969).

115. R. S. Lawrence and J. W. Strohbehn, "A Survey of Clear-Air Propagation Effects Relevant to Optical Communications," *Proc. IEEE* **58**:1523–1545 (1970).
116. S. F. Clifford, "The Classical Theory of Wave Propagation in a Turbulent Medium," in J. W. Strohbehn (ed.), *Laser Beam Propagation in the Atmosphere*, Springer-Verlag, New York, pp. 9–43, 1978.
117. A. I. Kon and V. I. Tatarskii, "Parameter Fluctuations of a Space-Limited Light Beam in a Turbulent Atmosphere," *Izv. VUZ Radiofiz.* **8**:870–875 (1965).
118. R. A. Schmeltzer, "Means, Variances, and Covariances for Laser Beam Propagation through a Random Medium," *Quart. J. Appl. Math.* **24**:339–354 (1967).
119. D. L. Fried and J. B. Seidman, "Laser-Beam Scintillation in the Atmosphere," *J. Opt. Soc. Am.* **57**:181–185 (1967).
120. D. L. Fried, "Scintillation of a Ground-to-Space Laser Illuminator," *J. Opt. Soc. Am.* **57**:980–983 (1967).
121. Y. Kinoshita, T. Asakura, and M. Suzuki, "Fluctuation Distribution of a Gaussian Beam Propagating through a Random Medium," *J. Opt. Soc. Am.* **58**:798–807 (1968).
122. A. Ishimaru, "Fluctuations of a Beam Wave Propagating through a Locally Homogeneous Medium," *Radio Sci.* **4**:295–305 (1969).
123. A. Ishimaru, "Fluctuations of a Focused Beam Wave for Atmospheric Turbulence Probing," *Proc. IEEE* **57**:407–414 (1969).
124. A. Ishimaru, "The Beam Wave Case and Remote Sensing," in J. W. Strohbehn (ed.), *Laser Beam Propagation in the Atmosphere*, Springer-Verlag, New York, pp. 129–170, 1978.
125. P. J. Titterton, "Scintillation and Transmitter-Aperture Averaging over Vertical Paths," *J. Opt. Soc. Am.* **63**:439–444 (1973).
126. H. T. Yura and W. G. McKinley, "Optical Scintillation Statistics for IR Ground-to-Space Laser Communication Systems," *Appl. Opt.* **22**:3353–3358 (1983).
127. P. A. Lightsey, J. Anspach, and P. Sydney, "Observations of Uplink and Retroreflected Scintillation in the Relay Mirror Experiment," *Proc. SPIE* **1482**:209–222 (1991).
128. D. L. Fried, "Aperture Averaging of Scintillation," *J. Opt. Soc. Am.* **57**:169–175 (1967).
129. J. H. Churnside, "Aperture Averaging of Optical Scintillations in the Turbulent Atmosphere," *Appl. Opt.* **30**:1982–1994 (1991).
130. J. H. Churnside and R. G. Frehlich, "Experimental Evaluation of Log-Normally Modulated Rician and IK Models of Optical Scintillation in the Atmosphere," *J. Opt. Soc. Am. A* **6**:1760–1766 (1989).
131. G. Parry and P. N. Pusey, "K Distributions in Atmospheric Propagation of Laser Light," *J. Opt. Soc. Am.* **69**:796–798 (1979).
132. G. Parry, "Measurement of Atmospheric Turbulence Induced Intensity Fluctuations in a Laser Beam," *Opt. Acta* **28**:715–728 (1981).
133. W. A. Coles and R. G. Frehlich, "Simultaneous Measurements of Angular Scattering and Intensity Scintillation in the Atmosphere," *J. Opt. Soc. Am.* **72**:1042–1048 (1982).
134. J. H. Churnside and S. F. Clifford, "Lognormal Rician Probability-Density Function of Optical Scintillations in the Turbulent Atmosphere," *J. Opt. Soc. Am. A* **4**:1923–1930 (1987).
135. R. Dashen, "Path Integrals for Waves in Random Media," *J. Math. Phys.* **20**:894–920 (1979).
136. K. S. Gochelashvily and V. I. Shishov, "Multiple Scattering of Light in a Turbulent Medium," *Opt. Acta* **18**:767–777 (1971).
137. S. M. Flatté and G. Y. Wang, "Irradiance Variance of Optical Waves Through Atmospheric Turbulence by Numerical Simulation and Comparison with Experiment," *J. Opt. Soc. Am. A* **10**:2363–2370 (1993).
138. R. Frehlich, "Simulation of Laser Propagation in a Turbulent Atmosphere," *Appl. Opt.* **39**:393–397 (2000).
139. S. M. Flatté and J. S. Gerber, "Irradiance-Variance Behavior by Numerical Simulation for Plane-Wave and Spherical-Wave Optical Propagation through Strong Turbulence," *J. Opt. Soc. Am. A* **17**:1092–1097 (2000).
140. R. Rao, "Statistics of the Fractal Structure and Phase Singularity of a Plane Light Wave Propagation in Atmospheric Turbulence," *Appl. Opt.* **47**:269–276 (2008).
141. K. S. Gochelashvily, V. G. Pevgov, and V. I. Shishov, "Saturation of Fluctuations of the Intensity of Laser Radiation at Large Distances in a Turbulent Atmosphere (Fraunhofer Zone of Transmitter)," *Sov. J. Quantum Electron.* **4**:632–637 (1974).
142. A. M. Prokhorov, F. V. Bunkin, K. S. Gochelashvily, and V. I. Shishov, "Laser Irradiance Propagation in Turbulent Media," *Proc. IEEE* **63**:790–810 (1975).

143. R. L. Fante, "Inner-Scale Size Effect on the Scintillations of Light in the Turbulent Atmosphere," *J. Opt. Soc. Am.* **73**:277–281 (1983).
144. R. G. Frehlich, "Intensity Covariance of a Point Source in a Random Medium with a Kolmogorov Spectrum and an Inner Scale of Turbulence," *J. Opt. Soc. Am. A.* **4**:360–366 (1987).
145. R. J. Hill and R. G. Frehlich, "Probability Distribution of Irradiance for the Onset of Strong Scintillation," *J. Opt. Soc. Am. A.* **14**:1530–1540 (1997).
146. J. H. Churnside and R. J. Hill, "Probability Density of Irradiance Scintillations for Strong Path-Integrated Refractive Turbulence," *J. Opt. Soc. Am. A.* **4**:727–733 (1987).
147. F. S. Vetelino, C. Young, L. Andrews, and J. Reolons, "Aperture Averaging Effects on the Probability Density of Irradiance Fluctuations in Moderate-to-Strong Turbulence," *Appl. Opt.* **46**: 2099–2108 (2007).
148. M. A. Al-Habash, L. C. Andrews, and R. L. Phillips, "Mathematical Model for the Irradiance PDF of a Laser Beam Propagating through Turbulent Media," *Opt. Eng.* **40**:1554–1562 (2001).
149. D. K. Killinger and A. Mooradian (eds.), *Optical and Laser Remote Sensing*, Springer-Verlag, New York, Optical Sciences, vol. 39, 1983.
150. L. J. Radziemski, R. W. Solarz, and J. A. Paisner (eds.), *Laser Spectroscopy and Its Applications*, Marcel Dekker, New York, Optical Eng., vol. 11, 1987.
151. T. Kobayashi, "Techniques for Laser Remote Sensing of the Environment," *Remote Sensing Reviews*, **3**:1–56 (1987).
152. E. D. Hinkley (ed.), *Laser Monitoring of the Atmosphere*, Springer-Verlag, Berlin, 1976.
153. W. B. Grant and R. T. Menzies, "A Survey of Laser and Selected Optical Systems for Remote Measurement of Pollutant Gas Concentrations," *APCA Journal* **33**:187 (1983).
154. "Optical Remote Sensing of the Atmosphere," *Conf. Proceedings, OSA Topical Meeting*, Williamsburg, 1991.
155. Dennis K. Killinger, "Lidar and Laser Remote Sensing", *Handbook of Vibrational Spectroscopy*, John Wiley & Sons, Chichester, 2002.
156. Claus Weitkamp, (ed.), *Lidar: Range Resolved Optical Remote Sensing of the Atmosphere*, Springer-Verlag, New York, 2005.
157. P. B. Hays, V. J. Abreu, D. A. Gell, H. J. Grassl, W. R. Skinner, and M. E. Dobbs, "The High Resolution Doppler Imager on the Upper Atmospheric Research Satellite," *J. Geophys. Res. (Atmosphere)* **98**:10713 (1993).
158. P. B. Hays, V. J. Abreu, M. D. Burrage, D. A. Gell, A. R. Marshall, Y. T. Morton, D. A. Ortland, W. R. Skinner, D. L. Wu, and J. H. Yee, "Remote Sensing of Mesospheric Winds with the High Resolution Imager," *Planet. Space Sci.* **40**:1599 (1992).
159. E. Browell, "Differential Absorption Lidar Sensing of Ozone," *Proc. IEEE* **77**:419 (1989).
160. J. M. Intrieri, A. J. Dedard, and R. M. Hardesty, "Details of Colliding Thunderstorm Outflow as Observed by Doppler Lidar," *J. Atmospheric Sciences* **47**:1081 (1990).
161. R. M. Hardesty, K. Elmore, M. E. Jackson, in *21st Conf. on Radar Meteorology*, American Meteorology Society, Boston, 1983.
162. M. P. McCormick and R. E. Veiga, "Initial Assessment of the Stratospheric and Climatic Impact of the 1991 Mount Pinatubo Eruption—Prolog," *Geophysical Research Lett.* **19**:155 (1992).
163. D. K. Killinger, S.D. Allen, R.D. Waterbury, C. Stefano, and E. L. Dottery, "Enhancement of Nd:YAG LIBS Emission of a Remote Target Using a Simultaneous CO₂ Laser Pulse," *Optics Express* **15**:12905 (2007).
164. A. Miziolek, V. Palleschi, and I. Schechter, (eds.), *Laser Induced Spectroscopy*, Cambridge University Press, Cambridge, 2006.
165. J. Kasparian, R. Sauerbrey, and S.L. Chin, "The Critical Laser Intensity of Self-Guided Light Filaments in Air," *Appl. Phys. B* **71**:877 (2000).
166. R. A. R. Tricker, *Introduction to Meteorological Optics*, American Elsevier, New York, 1970.
167. R. Greenler, *Rainbows, Halos, and Glories*, Cambridge University Press, Cambridge, 1980.
168. P. Forster, V. Ramaswamy, P. Artaxo, T. Berntsen, R. Betts, D.W. Fahey, J. Haywood, et al. "Changes in Atmospheric Constituents and in Radiative Forcing," In *Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the IV Assessment Report of the Intergovernmental Panel on Climate Change*, S. Solomon, D. Qin, M. Manning, Z. Chen, M. Marquis, K. B. Averyt, M. Tignor, and H. L. Miller (eds.), Cambridge University Press, 2007. Online at <http://www.ipcc.ch/pdf/assessment-report/ar4/wg1/ar4-wg1-chapter2.pdf>.

This page intentionally left blank.

DO NOT DUPLICATE

IMAGING THROUGH ATMOSPHERIC TURBULENCE

Virendra N. Mahajan*

*The Aerospace Corporation
El Segundo, California*

Guang-ming Dai

*Laser Vision Correction Group
Advanced Medical Optics
Milpitas, California*

ABSTRACT

In this chapter, how the random aberrations introduced by atmospheric turbulence degrade the image formed by a ground-based telescope with an annular pupil is considered. The results for imaging with a circular pupil are obtained as a special case of the annular pupil. Both the long- and short-exposure images are discussed in terms of their Strehl ratio, point-spread function (PSF), and transfer function. The discussion given is equally applicable to laser beams propagating through turbulence. An atmospheric coherence length is defined and it is shown that, for fixed power of a beam and regardless of the size of its diameter, the central irradiance in the focal plane is smaller than the corresponding aberration-free value for a beam of diameter equal to that of the coherence length. The aberration function is decomposed into Zernike annular polynomials and the autocorrelation and crosscorrelations of the expansion coefficients are given for Kolmogorov turbulence. It is shown that the aberration variance increases with the obscuration ratio of the annular pupil. The angle of arrival is also discussed, both in terms of the wavefront tilt as well as the centroid of the aberrated PSF. It is shown that the difference between the two is small, and the obscuration has only a second-order effect.

4.1 GLOSSARY

a	outer radius of the pupil
a_j	expansion coefficients
$A_L(\vec{r}_p)$	amplitude function of the lens (or imaging system) at a pupil point with position vector \vec{r}_p
C_n^2	refractive index structure parameter
\mathcal{D}	structure function (\mathcal{D}_w —wave, \mathcal{D}_ϕ —phase, \mathcal{D}_l —log amplitude, \mathcal{D}_n —refractive index)
D	diameter of exit pupil or aperture
F	focal ratio of the image-forming light cone ($F = R/D$)

*The author is also an adjunct professor at the College of Optical Sciences at the University of Arizona, Tucson and the Department of Optics and Photonics, National Central University, Chung Li, Taiwan.

$I_i(\vec{r}_i)$	irradiance at a point \vec{r}_i in the image plane
$I(r)$	normalized irradiance in the image plane such that its aberration-free central value is $I(0) = 1$
$\langle I_0 \rangle$	time-averaged irradiance at the exit pupil
j	Zernike aberration mode number
J	number of Zernike aberration modes
$l(\vec{r}_p)$	log-amplitude function introduced by atmospheric turbulence
L	path length through atmosphere or from source to receiver
$n(\vec{r})$	fluctuating part of refractive index $N(\vec{r})$
$P(r_c)$	encircled power in a circle of radius r_c in the image plane
P_{ex}	power in the exit pupil
$P_L(\vec{r}_p)$	lens pupil function
$P_R(\vec{r}_p)$	complex amplitude variation introduced by atmospheric turbulence
r_0	Fried's atmospheric coherence length
R	radius of curvature of the reference sphere with respect to which the aberration is defined
$R_n^m(\rho)$	Zernike circle radial polynomial of degree n and azimuthal frequency m
$\langle S \rangle$	time-averaged Strehl ratio
S_a	coherent area $\pi r_0^2/4$ of atmospheric turbulence
$\langle S_t \rangle$	tilt-corrected time-averaged Strehl ratio
S_{ex}	area of exit pupil
$Z_n^m(\rho, \theta)$	Zernike circle polynomial of degree n and azimuthal frequency m
ϵ	obscuration ratio of an annular pupil
Δ_j	phase aberration variance after correcting $J = j$ aberration modes
η	central irradiance in the image plane normalized by the aberration-free value for a pupil with area S_a but containing the same total power
λ	wavelength of object radiation
\vec{v}_i	spatial frequency vector in the image plane
\vec{v}	normalized spatial frequency vector
ϑ_0	isoplanatic angle of turbulence
$\sigma_\alpha, \sigma_\beta$	tip and tilt angle standard deviations
σ_Φ^2	phase aberration variance
$\tau(\vec{v})$	optical transfer function
ρ	radial variable normalized by the pupil radius a
$\tau_a(\nu)$	long-exposure (LE) atmospheric MTF reduction factor
$\Phi(\vec{r}_p)$	phase aberration
$\Phi_R(\vec{r}_p)$	random phase aberration introduced by atmospheric turbulence

4.2 INTRODUCTION

The resolution of a telescope forming an aberration-free image is determined by its diameter D ; larger the diameter, better the resolution. However, in ground-based astronomy, the resolution is degraded considerably because of the aberrations introduced by atmospheric turbulence. A plane wave of uniform amplitude and phase representing the light from a star propagating through the atmosphere undergoes both amplitude and phase variations due to the random inhomogeneities

in its refractive index. The amplitude variations, called scintillations, result in the twinkling of stars. The purpose of a large ground-based telescope has therefore generally not been better resolution but to collect more light so that dim objects may be observed. Of course, with the advent of adaptive optics,¹⁻³ the resolution can be improved by correcting the phase aberrations with a deformable mirror. The amplitude variations are negligible in near-field imaging, that is, when the far-field distance $D^2/\lambda \gg L$ or $D \gg \sqrt{\lambda L}$, where λ is the wavelength of the starlight and L is the propagation path length through the turbulence.⁴ In principle, a diffraction-limited image can be obtained if the aberrations are corrected completely in real time by the deformable mirror. However, in far-field imaging, that is, when $D^2/\lambda \ll L$ or $D \ll \sqrt{\lambda L}$, significant amplitude variations are introduced in addition to the phase variations. Since the amplitude variations cannot be corrected by adaptive optics, even a complete correction of the phase errors does not yield a diffraction-limited image.

In this chapter, we consider the effect of random aberrations introduced by atmospheric turbulence on the image quality. The discussion given is applicable equally to laser beams propagating through turbulence. First, we derive expressions for the degraded time-averaged point-spread function (PSF) and optical transfer function (OTF). The pupil function of the overall imaging system is written as the product of the pupil function of the optical system and a complex amplitude factor introduced by turbulence. As a result, the time-averaged OTF of the overall system is also equal to the product of the OTF of the optical system and an OTF reduction factor representing the mutual coherence function of the wave propagating through the turbulent atmosphere. The time-averaged images thus obtained are referred to as the *long-exposure images*.⁴ The exposure time may be 1 to 10 s.

Next, the structure functions for refractive index and phase fluctuations are given for the *Kolmogorov turbulence*. We introduce the notion of *atmospheric coherence length* r_0 , and show that it limits the resolution of the image regardless of how large the telescope diameter is.⁴ Since large telescopes have annular pupils, we consider systems with such pupils and describe the Strehl ratio, PSF, and encircled power as a function of the ratio of the pupil diameter and the coherence length. The phase aberration introduced by Kolmogorov turbulence is expanded in terms of the Zernike annular polynomials,⁷⁻⁸ and autocorrelation and crosscorrelation of the expansion coefficients are given.⁹ It is shown that atmospheric turbulence dominates the degradation of the system modulation transfer function (MTF), Strehl ratio, or the angle of arrival; the effect of the pupil obscuration is small even for weak turbulence. The piston-removed aberration variance increases monotonically as the obscuration of the pupil increases. A large portion (87 percent for a circular pupil and somewhat larger for an annular pupil) of the aberration variance is a random wavefront tilt, resulting in a random displacement of the image over a long exposure. Hence, a better quality image is obtained if the exposure time is short enough so that the image does not wander. Such an image is referred to as a *short-exposure image*.⁴ The exposure time may be 0.1 s or less. The location of such an image is random which is not relevant in astronomy. The image wanders, sometimes referred to as the *angle of arrival*, is considered in terms of the wavefront tilt (called the Z-tilt) or the *centroid* of the image (called the G-tilt), and expressions are given for their standard deviations.

The fluctuating image can be stabilized if the tilt is corrected in real time with a steering mirror. A tilt-corrected image is equivalent to a time-averaged short-exposure image. The characteristics of such an image are considered and compared with those of a long-exposure image. The angular resolution of the system for a long-exposure image is about λ/r_0 compared to the diffraction-limited resolution λ/D . For a short-exposure image the resolution is substantially better than that for a long-exposure image. When the coherence length is much smaller than the pupil diameter, the image breaks up into small spots, called *speckles*, whose size depends on the latter, while the overall size of the image is determined by the former. An example of a speckled short-exposure image of a point object is illustrated and compared with the corresponding aberration-free image.

4.3 LONG-EXPOSURE IMAGE

When the medium between an object and the optical system imaging it is homogeneous, a spherical wavefront of uniform amplitude centered on an object point is incident on the entrance pupil of the system. If the system is aberration free and has uniform transmittance, a spherical wavefront of

uniform amplitude centered at the Gaussian image point emerges from its exit pupil. However, if the system is aberrated and has nonuniform transmittance, then the emerging wavefront is distorted with a nonuniform amplitude. Let $\Phi_L(\vec{r}_p)$ and $A_L(\vec{r}_p)$ be the phase and amplitude transmittance corresponding to a point \vec{r}_p on the exit pupil. A dimensionless pupil function of the system may be written

$$P_L(\vec{r}_p) = A_L(\vec{r}_p) \exp[i\Phi_L(\vec{r}_p)] \quad (1)$$

When the medium between the object and the imaging system is inhomogeneous as in ground-based astronomy, it introduces random phase and amplitude variations across the wavefront propagating through it. Let $\Phi_R(\vec{r}_{en})$ and $\ell(\vec{r}_{en})$ be the random phase and (dimensionless) log amplitude introduced by atmospheric turbulence at a point \vec{r}_{en} on the entrance pupil. As the wave propagates through the imaging system, it undergoes additional phase and amplitude variations. The total phase aberration at a point \vec{r}_p on the exit pupil may be written

$$\Phi(\vec{r}_p) = \Phi_L(\vec{r}_p) + \Phi_R(\vec{r}_p) \quad (2)$$

where the position vector \vec{r}_p is related to the position vector \vec{r}_{en} by the pupil magnification m_p (i.e., $\vec{r}_p = m_p \vec{r}_{en}$). The pupil function of the overall system (i.e., including the effects of atmospheric turbulence) representing the wavefront at the exit pupil may be written

$$P(\vec{r}_p) = P_L(\vec{r}_p) P_R(\vec{r}_p) \quad (3)$$

where

$$P_R(\vec{r}_p) = \exp[\ell(\vec{r}_p) + i\Phi_R(\vec{r}_p)] \quad (4)$$

is truncated by the pupil function $P_L(\vec{r}_p)$ and represents the complex amplitude variation introduced by turbulence.

The instantaneous irradiance distribution of the star image formed by the overall system is given by^{10,11}

$$I_i(\vec{r}_i) = \frac{P_{ex}}{S_{ex} \lambda^2 R^2} \left| \int P(\vec{r}_p) \exp\left(-\frac{2\pi i}{\lambda R} \vec{r}_p \cdot \vec{r}_i\right) d\vec{r}_p \right|^2 \quad (5)$$

where \vec{r}_i is the position vector of a point in the image plane, P_{ex} is the total power in the image, S_{ex} is the area of the exit pupil, λ is the wavelength of object radiation, and R is the radius of curvature of the reference sphere with respect to which the aberration is defined.

Substituting Eq. (3) into Eq. (5), the time-averaged distribution, representing a long-exposure image, may be written

$$\begin{aligned} \langle I_i(\vec{r}_i) \rangle &= \frac{\langle P_{ex} \rangle}{S_{ex} \lambda^2 R^2} \iint P_L(\vec{r}_p) P_L^*(\vec{r}'_p) \langle P_R(\vec{r}_p) P_R^*(\vec{r}'_p) \rangle \\ &\quad \times \exp\left[-\frac{2\pi i}{\lambda R} (\vec{r}_p - \vec{r}'_p) \cdot \vec{r}_i\right] d\vec{r}_p d\vec{r}'_p \end{aligned} \quad (6)$$

where * denotes a complex conjugate. Thus, considering Eq. (4), we may write the *mutual coherence function* (MCF) of the wave incident on the system, using abbreviated notation

$$\begin{aligned} \langle P_R(\vec{r}_p) P_R^*(\vec{r}'_p) \rangle &\equiv \langle P_1 P_2^* \rangle \\ &= \langle \exp[(\ell_1 + \ell_2) + i(\Phi_1 - \Phi_2)] \rangle \end{aligned} \quad (7)$$

If the refractive index fluctuations are statistically stationary, that is, they are statistically homogeneous and isotropic, then so are the fluctuations in ℓ and Φ which they generate. Therefore,

$$\langle (\ell_1 + \ell_2) (\Phi_1 - \Phi_2) \rangle = (\langle \ell_1 \Phi_1 \rangle - \langle \ell_2 \Phi_2 \rangle) + (\langle \ell_2 \Phi_1 \rangle - \langle \ell_1 \Phi_2 \rangle) \quad (8a)$$

$$= 0 \quad (8b)$$

where the two averaged quantities in the first term on the right-hand side of Eq. (8a) cancel each other due to homogeneity of turbulence, and those in the second term cancel due to its isotropy. Thus, $\ell_1 + \ell_2$ and $\Phi_1 - \Phi_2$ are uncorrelated random variables. Hence, Eq. (7) may be written

$$\langle P_R(\vec{r}_p) P_R^*(\vec{r}'_p) \rangle = \langle \exp(\ell_1 + \ell_2) \rangle \langle \exp[i(\Phi_1 - \Phi_2)] \rangle \quad (9)$$

For a Gaussian random variable x with a mean value of $\langle x \rangle$ and a standard deviation of σ , it is easy to show that

$$\begin{aligned} \langle \exp(bx) \rangle &= \frac{1}{2\pi\sigma} \int_{-\infty}^{\infty} \exp(bx) \exp[-(x - \langle x \rangle)^2 / 2\sigma^2] dx \\ &= \exp \left[\frac{1}{2} b^2 \langle (x - \langle x \rangle)^2 \rangle + b \langle x \rangle \right] \end{aligned} \quad (10)$$

where b is an arbitrary constant. As the wave propagates through the random atmosphere, it accumulates randomly both phase and amplitude variations. Each of these two random variables can be expressed as an integral along the propagation path. Since the atmospheric path length is much longer than the correlation length of these variable, the accumulation consists of a sum of a large number of terms that are statistically independent. By the central limit theorem, we can infer that they obey Gaussian statistics. Letting both ℓ and Φ_R be Gaussian random variables where Φ_R has a mean value of zero, it can be shown from conservation of power, that^{4,11}

$$\langle \exp(\ell_1 + \ell_2) \rangle = \exp \left[-\frac{1}{2} \mathcal{D}_\ell(|\vec{r}_p - \vec{r}'_p|) \right] \quad (11)$$

where

$$\mathcal{D}_\ell(|\vec{r}_p - \vec{r}'_p|) = \langle [(\ell(\vec{r}_p) - \ell(\vec{r}'_p))]^2 \rangle \quad (12)$$

is the *log-amplitude structure function*.

Since Φ_1 and Φ_2 are Gaussian random variables with zero mean values, $\Phi_1 - \Phi_2$ is also a Gaussian random variable. Hence, following Eq. (10), we obtain

$$\begin{aligned} \langle \exp[i(\Phi_1 - \Phi_2)] \rangle &= \exp \left[-\frac{1}{2} \langle (\Phi_1 - \Phi_2)^2 \rangle \right] \\ &= \exp \left[-\frac{1}{2} \mathcal{D}_\Phi(|\vec{r}_p - \vec{r}'_p|) \right] \end{aligned} \quad (13)$$

where

$$\mathcal{D}_\Phi(|\vec{r}_p - \vec{r}'_p|) = \langle [(\Phi(\vec{r}_p) - \Phi(\vec{r}'_p))]^2 \rangle \quad (14)$$

is the *phase structure function* of turbulence. Substituting Eqs. (11) and (13) into Eq. (9), we may write the MCF

$$\langle P_R(\vec{r}_p) P_R^*(\vec{r}'_p) \rangle = \exp \left[-\frac{1}{2} \mathcal{D}_w(|\vec{r}_p - \vec{r}'_p|) \right] \quad (15)$$

where

$$\mathcal{D}_w(|\vec{r}_p - \vec{r}'_p|) = \mathcal{D}_\ell(|\vec{r}_p - \vec{r}'_p|) + \mathcal{D}_\Phi(|\vec{r}_p - \vec{r}'_p|) \quad (16)$$

is called the *wave structure function* of turbulence.

Substituting Eq. (15) into Eq. (6), we obtain

$$\begin{aligned} \langle I_i(\vec{r}_i) \rangle &= \frac{\langle P_{\text{ex}} \rangle}{S_{\text{ex}} \lambda^2 R^2} \iint P_L(\vec{r}_p) P_L^*(\vec{r}'_p) \exp \left[-\frac{1}{2} \mathcal{D}_w(|\vec{r}_p - \vec{r}'_p|) \right] \\ &\times \exp \left[-\frac{2\pi i}{\lambda R} (\vec{r}_p - \vec{r}'_p) \cdot \vec{r}_i \right] d\vec{r}_p d\vec{r}'_p \end{aligned} \quad (17)$$

Because of a Fourier transform relationship between the PSF and the OTF, we identify $(\vec{r}_p - \vec{r}'_p)/\lambda R$ with a spatial frequency \vec{v}_i . Thus, we let

$$\vec{r}_p - \vec{r}'_p = \lambda R \vec{v}_i \quad (18)$$

Substituting Eq. (18) into Eq. (17) and carrying out the integration over \vec{r}'_p , we may write

$$\langle I_i(\vec{r}_i) \rangle = \langle P_{\text{ex}} \rangle \int \tau_L(\vec{v}_i) \exp \left[-\frac{1}{2} \mathcal{D}_w(\lambda R \vec{v}_i) \right] \exp(-2\pi i \vec{v}_i \cdot \vec{r}_i) d\vec{v}_i \quad (19)$$

where

$$\tau_L(\vec{v}_i) = S_{\text{ex}}^{-1} \int P_L(\vec{r}_p) P_L^*(\vec{r}_p - \lambda R \vec{v}_i) d\vec{r}_p \quad (20)$$

is the OTF of the (turbulence-free) optical system and $v_i = |\vec{v}_i|$. Now, we introduce normalized quantities

$$I(r) = I_i(\vec{r}_i)/I(0) \quad (21a)$$

$$\vec{r} = \vec{r}_i/\lambda F \quad (21b)$$

and

$$\vec{v} = \vec{v}_i/(1/\lambda F) \quad (21c)$$

where

$$I(0) = \frac{\langle P_{\text{ex}} \rangle S_{\text{ex}}}{\lambda^2 R^2} \quad (22)$$

is the aberration-free central irradiance for a uniform-amplitude wavefront with a total power $\langle P_{\text{ex}} \rangle$ and $1/\lambda F$ is the cutoff spatial frequency of the optical system. Here, $F = R/D$ is the focal ratio of the image-forming light cone. In terms of the normalized quantities, Eq. (19) for the time-averaged PSF may be written

$$\langle I(\vec{r}) \rangle = (4/\pi) \int \langle \tau(\vec{v}) \rangle \exp(-2\pi i \vec{v} \cdot \vec{r}) d\vec{v} \quad (23)$$

where

$$\langle \tau(\vec{v}) \rangle = \tau_L(\vec{v}) \exp \left[-\frac{1}{2} \mathcal{D}_w(\nu D) \right] \quad (24)$$

is the time-averaged OTF. This OTF being equal to the product of the OTF of the optical system and the MCF associated with atmospheric turbulence is a consequence of the fact that the pupil function of the overall system is equal to the product of the pupil function of the optical system and the complex amplitude variation introduced by turbulence. It should be evident that, since $\nu D = \nu \lambda R$, the modification of the system OTF due to turbulence is independent of the diameter D of the exit pupil of the system. The MCF acts as a *reduction factor* for the OTF of the overall system.¹¹ It may also be referred to as the Hopkins ratio of the image if $\tau_L(\vec{v})$ is the aberration-free OTF.¹¹

4.4 KOLMOGOROV TURBULENCE AND ATMOSPHERIC COHERENCE LENGTH

The refractive index $N(\vec{r})$ of the turbulent atmosphere at a point \vec{r} in space fluctuates due to fluctuations of its temperature. It can be written as the sum of its mean value $\langle N(\vec{r}) \rangle$ at a point \vec{r} and a fluctuating part $n(\vec{r})$ at that point, in the form

$$N(\vec{r}) = \langle N(\vec{r}) \rangle + n(\vec{r}) \quad (25)$$

Whereas $\langle N(\vec{r}) \rangle \approx 1$, $n(\vec{r})$ is only on the order of 10^{-6} . It should be evident that $n(\vec{r})$ has a mean value $\langle n(\vec{r}) \rangle = 0$. The structure function $\mathcal{D}_n(\vec{r}_1, \vec{r}_2)$ of the refractive index fluctuations represents the mean square value of the difference of refractive index at two points \vec{r}_1 and \vec{r}_2 , that is,

$$\mathcal{D}_n(\vec{r}_1, \vec{r}_2) = \langle [n(\vec{r}_1) - n(\vec{r}_2)]^2 \rangle \quad (26)$$

The turbulent atmosphere consists of packets of air called *eddies* each with a characteristic value of its refractive index. For Kolmogorov turbulence, the structural function is given by

$$\mathcal{D}_n(r) = C_n^2 r^{2/3}, \quad \ell_0 \ll r \ll L_0 \quad (27)$$

where $r = |\vec{r}_2 - \vec{r}_1|$ and C_n^2 (in units of $\text{m}^{-2/3}$) is called the *refractive index structure parameter*. The quantities ℓ_0 and L_0 are called the *inner* and *outer scales of turbulence* representing the smallest and the largest eddies, respectively. Typical values of C_n^2 vary from $10^{-13} \text{ m}^{-2/3}$ for strong turbulence to $10^{-17} \text{ m}^{-2/3}$ for weak turbulence. Values of ℓ_0 are on the order of a few millimeters, and those of L_0 vary from 1 to 100 m.

The wave structure function for a spherical wave propagating through Kolmogorov turbulence is given by¹²

$$\mathcal{D}_w(r) = 2.914 k^2 r^{5/3} \int_0^L C_n^2(z) (z/L)^{5/3} dz \quad (28)$$

where $k = 2\pi/\lambda$ and z varies along the atmospheric path of total length L from a value of zero at the source (or the object plane) to a value of L at the receiver (or the image plane). The wave structure function can also be written in the form^{4,12}

$$\mathcal{D}_w(r) = 6.88 (r/r_0)^{5/3} \quad (29)$$

where, for a spherical wave,

$$r_0 = 0.185 \lambda^{6/5} \left\{ \int_0^L C_n^2(z) (z/L)^{5/3} dz \right\}^{-3/5} \quad (30)$$

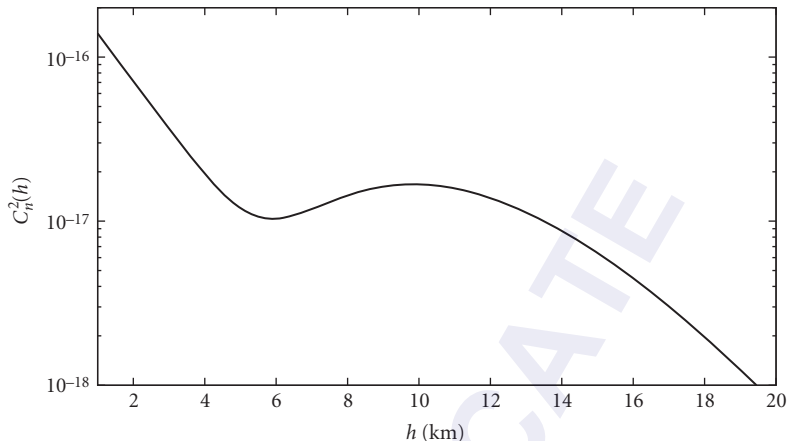


FIGURE 1 Variation of the refractive index structure parameter $C_n^2(h)$ with the altitude h for the Hufnagle-Valley model $H-V_{5/7}$ of Eq. (31).

is a characteristic length of Kolmogorov turbulence called its *coherence length* or *diameter*, often referred to as *Fried's coherence length*. We note that r_0 varies with the optical wavelength as $\lambda^{6/5}$. If the line of sight makes an angle θ with the zenith, then the path length L and the integral in Eq. (30) are increased by $\sec \theta$, or r_0 decreases by a factor of $(\sec \theta)^{3/5}$.

A commonly used model for the $C_n^2(z)$ variation is the Hufnagle-Valley model, often referred to as the $H-V_{5/7}$ model, given by¹³

$$C_n^2(h) = 8.2 \times 10^{-26} W^2 h^{10} e^{-h} + 2.7 \times 10^{-16} e^{-h/1.5} + A e^{-10h} \quad (31)$$

where h is the height from ground in kilometers, $A = 1.7 \times 10^{-14}$, and $W = 21$ m/s is the wind speed. Its variation with h is shown in Fig. 1. It decreases rapidly for the first few kilometers, dips at about 5.9 km, rises slightly due to turbulence in the jet stream peaking at about 9.8 km, and decreases monotonically to negligible values beyond about 20 km.

For a point source on the ground observed by an observer on an aircraft, space shuttle, or a satellite, light propagates upward in a manner similar to the $C_n^2(h)$ profile of increasing height h . The observer looks at the point source down through the atmosphere.¹⁴ Substituting Eq. (31) into Eq. (30) with h replaced by z (keeping in mind that C_n^2 is in units of $\text{m}^{-2/3}$), we obtain the value of r_0 in the plane of the observer as a function of his/her height L , as shown in Fig. 2 by the solid curve. The observation is assumed to be along the zenith at a wavelength of $\lambda = 0.5 \mu\text{m}$. We note that r_0 increases slowly, goes through a peak at about 7.2 km and a valley at about 13.7 km, and then increases linearly with L beyond a height of about 20 km above which there is very little atmosphere.

For a ground observer looking up at a point source on an aircraft, space shuttle, or a satellite, z increases downward with a value of zero at the source and L on the ground. To determine r_0 from Eq. (30) using the $C_n^2(h)$ profile such as given by Eq. (31), we must either replace h by $L-z$, or replace z/L by $1 - z/L$ as may be seen by a change of variable. Thus, Eq. (30) in this case is replaced by

$$\begin{aligned} r_0 &= 0.185 \lambda^{6/5} \left[\int_0^L C_n^2(L-z) (z/L)^{5/3} dz \right]^{-3/5} \\ &= 0.185 \lambda^{6/5} \left[\int_0^L C_n^2(z) \left(1 - \frac{z}{L}\right)^{5/3} dz \right]^{-3/5} \end{aligned} \quad (32)$$

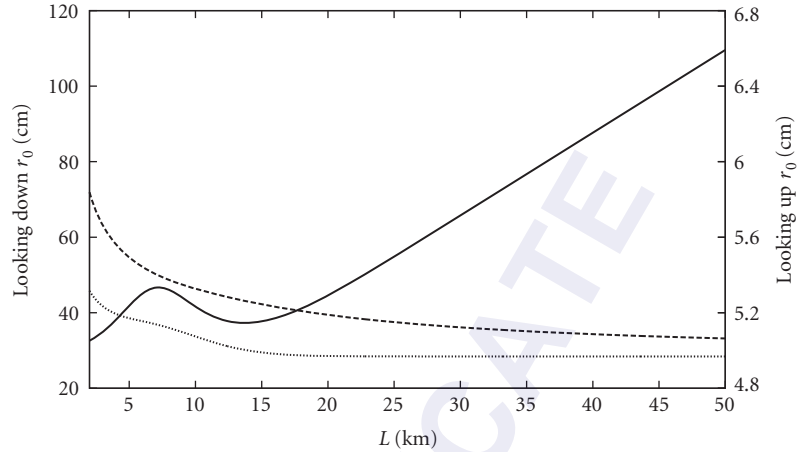


FIGURE 2 Variation of r_0 in the plane of the observer with the path length L from the source to the observer. The propagation of light is along the zenith at a wavelength $\lambda = 0.5 \mu\text{m}$. The solid curve is for a point source on ground such that the spherical wave propagates upward to a space observer looking down directly at the source. The dashed curve is for a point source in space so that a spherical wave propagates down and an observer on ground looks at it directly above. The dotted curve is for plane wave propagation with an observer on ground (or in space) looking up (or down) at a space (or ground) object.

As shown in Fig. 2 by the dashed curve, the value of r_0 in the plane of the observer on ground does not change very much with the height L of the point source.

For a plane wave propagating up or down through the atmosphere, the factor z/L in Eq. (28) and in turn, in Eq. (30), is replaced by unity. It is easy to see this for starlight observed from ground. A plane wave can be thought of as a spherical wave originating at an infinite distance and traveling through a uniform medium for which $C_n^2 = 0$, except for the propagation path through the atmosphere. The value of z/L in the region for which $C_n^2 \neq 0$ is infinitesimally different from unity. Hence, the starlight propagation in ground-based astronomy can be considered as a plane wave propagation with $(z/L)^{5/3}$ in Eq. (30) replaced by unity, or a spherical wave propagating an infinite distance to reach the earth's atmosphere with nonzero C_n^2 value only near the end of its path for which z/L is negligibly different from unity. The variation of r_0 with the separation L of the plane-wave source and the observer is also shown in Fig. 2 by the dotted curve. We note that the variation is small and occurs for small values of L . For large values of L , approaching infinity as in ground-based astronomy, the value of r_0 is close to 5 cm. Moreover, the value of r_0 for a ground observer looking up through the atmosphere at an incoming plane wave is smaller than its value for a corresponding spherical wave, as expected from the smaller factor of $1 - z/L$ compared to z/L in Eq. (32).

We also note that r_0 is smaller when a satellite is observed from ground compared to when a ground object is observed from a satellite. As an observer moves to a higher altitude above the atmosphere, the value of r_0 increases linearly with the altitude L . Consequently, the image degradation is much smaller when a ground object is observed from space than when a space object is observed from ground. In the first case, the object is near the region of turbulence and it is observed from far away. In the second case, the object is away from the region of turbulence, but it is observed from nearby. The fact that the image quality is superior in the first case is similar to when an object behind a diffuse shower glass is observed. One can see some detail in the object when it is in contact with the shower glass. However, as soon as the object is moved slightly away from the shower glass, it appears only as a halo, illustrating complete loss of image resolution. This does not, however, mean that *reciprocity* of wave propagation does not hold. For example, if the wavefront errors of a wave from a point source

in space propagating downward are measured on ground and a conjugate correction is introduced in a beam transmitted upward with a deformable mirror, the beam focus in space will be diffraction limited (neglecting any measurement or correction error), illustrating that the atmosphere introduces the same wavefront errors whether a beam is propagating up or down through it.

Neglecting the variation of C_n^2 for horizontal propagation, we obtain

$$\mathcal{D}_w(r) = \begin{cases} 2.91C_n^2Lk^2r^{5/3} & \text{Plane wave} \\ (3/8)2.91C_n^2Lk^2r^{5/3} & \text{Spherical wave} \end{cases} \quad (33a)$$

$$(33b)$$

and

$$r_0 = \begin{cases} 1.68(C_n^2Lk^2)^{-3/5} & \text{Plane wave} \\ 3.02(C_n^2Lk^2)^{-3/5} & \text{Spherical wave} \end{cases} \quad (34a)$$

$$(34b)$$

In the *near field* ($L \ll D^2/\lambda$), the amplitude variations are negligible and, therefore, $\mathcal{D}_\Phi(r) = \mathcal{D}_w(r)$. In the *far field* ($L \gg D^2/\lambda$), $\mathcal{D}_\Phi(r) = (1/2)\mathcal{D}_w(r)$. For in-between ranges, the multiplying factor between the two structure functions varies smoothly between 1 and 1/2.

Substituting Eq. (29) into Eq. (24), we may write the time-averaged OTF

$$\langle \tau(\bar{v}; D/r_0) \rangle = \tau_L(\bar{v})M(v; D/r_0) \quad (35)$$

where

$$M(v; D/r_0) = \exp[-3.44(vD/r_0)^{5/3}] \quad (36)$$

is the *long-exposure MCF*. Since $vD = v_i\lambda R$ and r_0 varies as $\lambda^{6/5}$, the factor in the exponent in Eq. (36) varies with wavelength as $\lambda^{-1/3}$.

Since $\exp(-3.44) \approx 0.03$, atmospheric turbulence reduces the overall system MTF corresponding to a spatial frequency $v = r_0/D$ by a factor of 0.03. From Eqs. (15) and (29), we find that r_0 represents a correlation length such that the correlation of complex amplitudes at two points on a wave separated by a distance r_0 is 0.03. Moreover, by definition, the MCF represents the *degree of spatial coherence* of the wave at the receiver, and thus the visibility of the fringes formed in a two-pinhole experiment and observed in the vicinity of a point that is equidistant from the two pinholes. Note that because of the random nature of atmospheric turbulence, the time-averaged irradiances at the two pinholes are equal to each other. Hence, r_0 represents a *partial coherence length* of the wave so that its degree of coherence corresponding to two points on it separated by r_0 is only 0.03, or that the visibility of the fringes formed by the secondary waves from these points is 0.03. The value of r_0 on a mountain site may vary from 5 to 10 cm in the visible region of the spectrum, and increases with wavelength as $\lambda^{6/5}$.

4.5 APPLICATION TO SYSTEMS WITH ANNULAR PUPILS

Now we apply Eq. (23) to a system with an annular pupil of outer radius $a = D/2$ and inner radius $a\epsilon$, and thus an obscuration ratio ϵ and an area $S_{\text{ex}}(\epsilon) = \pi(1 - \epsilon^2)a^2$ (see Fig. 3). For simplicity, we assume that the system is aberration free. Accordingly, its OTF is given by^{11,16}

$$\tau(v; \epsilon) = \frac{1}{1 - \epsilon^2} [\tau(v) + \epsilon^2\tau(v/\epsilon) - \tau_{12}(v; \epsilon)], \quad 0 \leq v \leq 1 \quad (37)$$

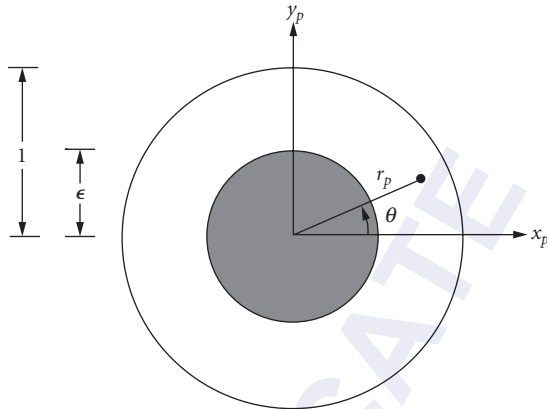


FIGURE 3 Annular pupil of obscuration ratio ϵ , representing the ratio of the inner and outer radii of the pupil.

where $\tau(v)$ is given by

$$\begin{aligned} \tau(v) &= (2/\pi)[\cos^{-1}v - v(1 - v^2)^{1/2}], \quad 0 \leq v \leq 1 \\ &= 0, \text{ otherwise} \end{aligned} \quad (38)$$

and represents the OTF of a corresponding system without any obscuration (i.e., for a circular pupil) and

$$\tau_{12}(v; \epsilon) = 2\epsilon^2, \quad 0 \leq v \leq (1 - \epsilon)/2 \quad (39a)$$

$$= (2/\pi)(\theta_1 + \epsilon^2\theta_2 - 2v\sin\theta_1), \quad (1 - \epsilon)/2 \leq v \leq (1 + \epsilon)/2 \quad (39b)$$

$$= 0, \text{ otherwise} \quad (39c)$$

In Eq. (39b), θ_1 and θ_2 are given by

$$\cos\theta_1 = \frac{4v^2 + 1 - \epsilon^2}{4v} \quad (39d)$$

and

$$\cos\theta_2 = \frac{4v^2 - 1 + \epsilon^2}{4\epsilon v} \quad (39e)$$

respectively. Since the OTF is equal to the autocorrelation of the pupil function,^{10,11,17} the cutoff frequency, which corresponds to the separation of two pupils by their diameter, is independent of ϵ . From Eq. (37)

$$\frac{\tau(v; \epsilon)}{\tau(v)} = \frac{1}{1 - \epsilon^2} \quad (40)$$

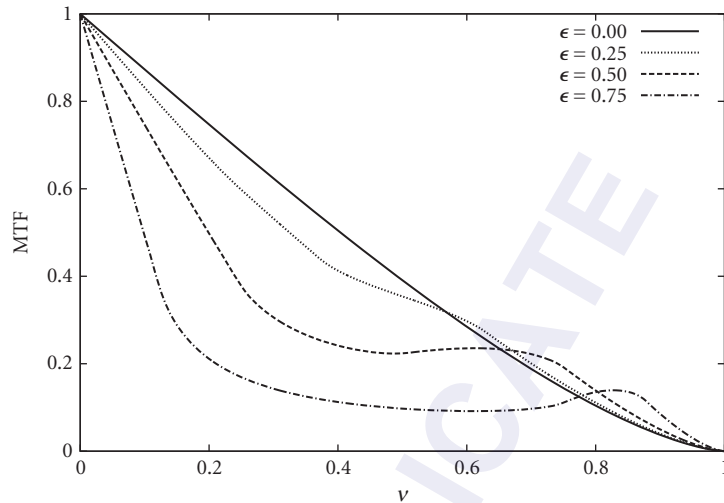


FIGURE 4 Aberration-free OTF $\tau(v; \epsilon)$ of an annular pupil illustrating the effect of its obscuration.

for spatial frequencies in the range $(1 + \epsilon)/2 < v < 1$. The overlap area of two annular pupils displaced from each other by an amount corresponding to a frequency in this frequency range is independent of ϵ , but the fractional area is larger owing to the smaller area of the obscured exit pupil.¹¹

How $\tau(v; \epsilon)$ varies with v is shown in Fig. 4 for various values of ϵ , including zero. It is evident that the OTF of an annular pupil is significantly lower at low frequencies but somewhat higher at high frequencies, compared to that for a corresponding circular ($\epsilon = 0$) pupil. The slope of the OTF for an annular pupil at the origin is given by

$$\tau'(0; \epsilon) = -4/\pi(1 - \epsilon) \quad (41)$$

This slope does not change when aberrations are introduced.¹¹ Moreover,

$$\int_0^1 \tau(v; \epsilon) v dv = (1 - \epsilon^2)/8 \quad (42)$$

The time-averaged irradiance distribution and encircled power (in a circle of radius r_c in units of λF) of the image of a point object are given by

$$\langle I(r; \epsilon; D/r_0) \rangle = \frac{8}{1 - \epsilon^2} \int_0^1 \langle \tau(v; \epsilon; D/r_0) \rangle J_0(2\pi r v) v dv \quad (43)$$

and

$$\langle P(r_c; \epsilon; D/r_0) \rangle = 2\pi r_c \int_0^1 \langle \tau(v; \epsilon; D/r_0) \rangle J_1(2\pi r_c v) dv \quad (44)$$

where

$$\langle \tau(v; \epsilon; D/r_0) \rangle = \tau(v; \epsilon) \exp \left[-3.44(vD/r_0)^{5/3} \right] \quad (45)$$

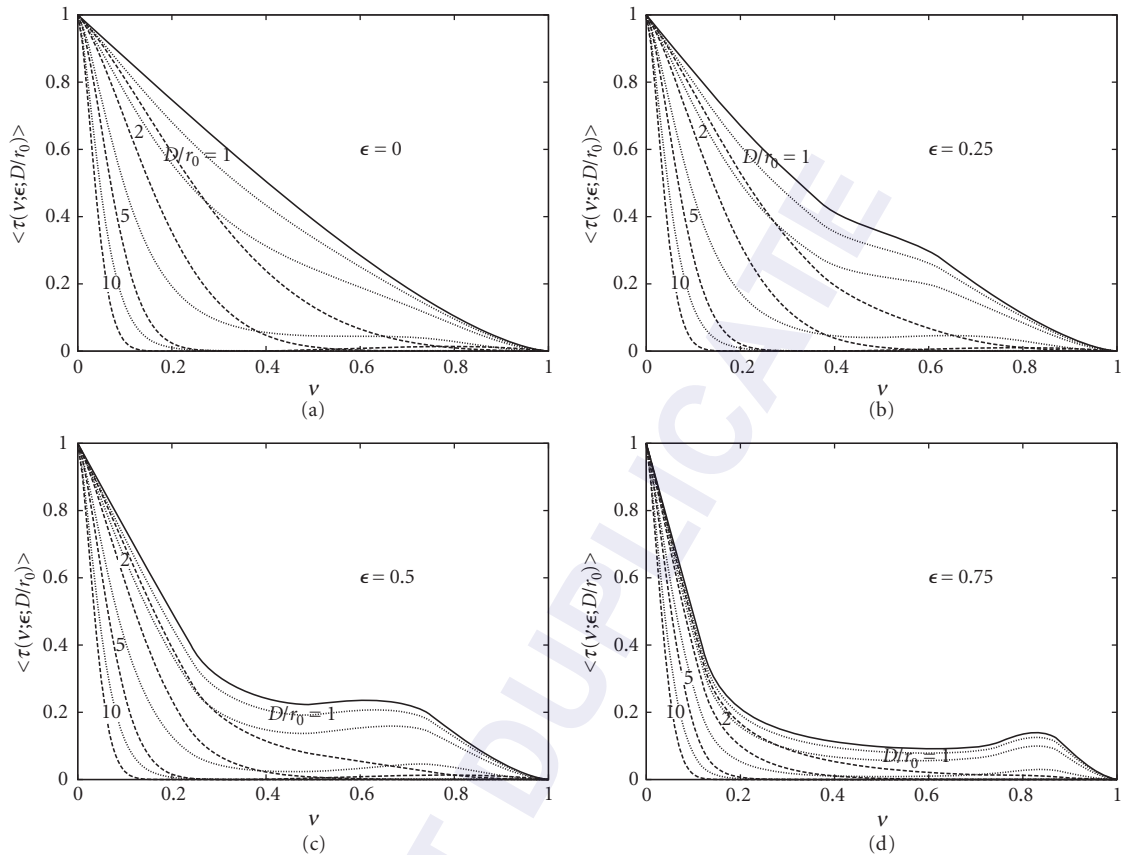


FIGURE 5 Time-averaged OTF for various values of ϵ and D/r_0 . The solid curves represent the aberration-free OTF, and the dashed and dotted curves represent the corresponding long- and short-exposure OTFs.

is the OTF of the system degraded by atmospheric turbulence. The irradiance is normalized by the aberration-free central irradiance $\langle P_{\text{ex}} \rangle S_{\text{ex}}(\epsilon)/\lambda^2 R^2$ and the encircled power is normalized by the total power $\langle P_{\text{ex}} \rangle$. Figure 5 shows the time-averaged OTF for several values of D/r_0 . The OTF gain at high frequencies disappears even for weak turbulence, as is evident from Fig. 5b for $D/r_0 = 1$. The turbulence dominates the OTF for large values of D/r_0 , and the effect of obscuration becomes small, as illustrated in Fig. 5c. Not only is the MTF at any frequency reduced, but the effective cutoff frequency is also reduced, for example, from a value of 1 to about 0.1 when $D/r_0 = 10$.

Letting $r = 0$ in Eq. (43) yields the Strehl ratio of the image

$$\begin{aligned} \langle S(\epsilon; D/r_0) \rangle &\equiv \langle I(0; \epsilon; D/r_0) \rangle \\ &= \frac{8}{1 - \epsilon^2} \int_0^1 \langle \tau(v; \epsilon; D/r_0) \rangle v dv \end{aligned} \quad (46)$$

Figure 6 shows how the time-averaged Strehl ratio varies with D/r_0 for several values of ϵ . The Strehl ratio decreases monotonically as D/r_0 increases. It is slightly lower for an annular pupil for values of $D/r_0 \leq 3$ and somewhat higher for $D/r_0 \geq 3$. As shown in Fig. 6b, it decreases monotonically with increasing

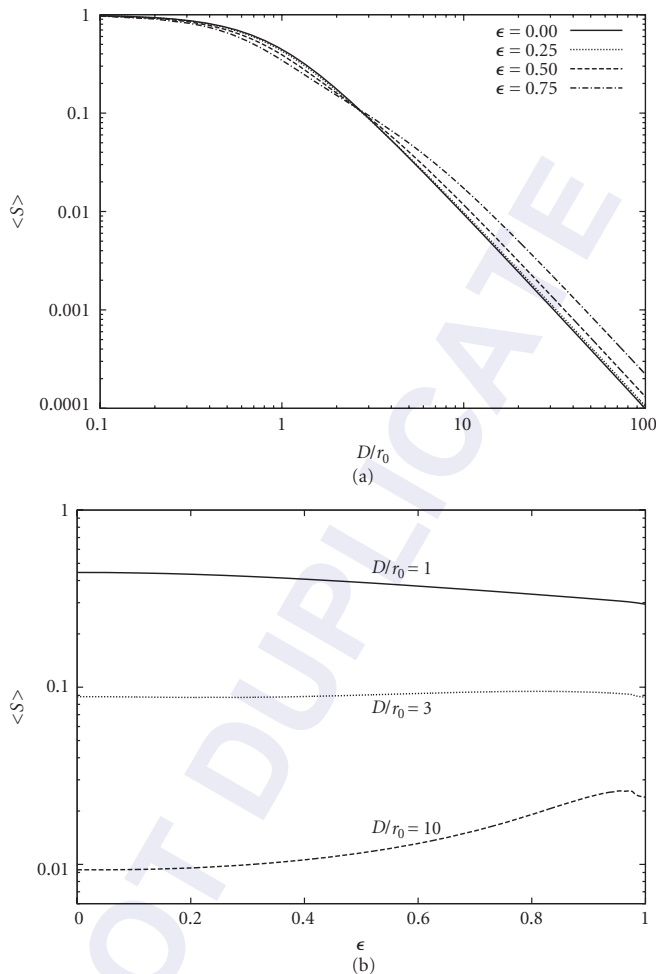


FIGURE 6 Variation of time-averaged long-exposure Strehl ratio $\langle S(\epsilon; D/r_0) \rangle$, representing the central irradiance normalized by its corresponding aberration-free value. (a) As a function of D/r_0 and (b) as a function of ϵ .

value of ϵ for small values of D/r_0 , remains nearly constant for $D/r_0 = 3$, and increases monotonically for $D/r_0 = 10$. Some typical values of the Strehl ratio are listed in Table 1 for several values of ϵ and D/r_0 .

Since $S_{\text{ex}}(\epsilon) = \pi(1 - \epsilon^2)D^2/4$, the aberration-free central irradiance $\langle P_{\text{ex}} \rangle S_{\text{ex}}(\epsilon)/\lambda^2 R^2$ for a fixed total power $\langle P_{\text{ex}} \rangle$ increases with D as D^2 . To see how the aberrated central irradiance varies with D , we consider a quantity¹⁸

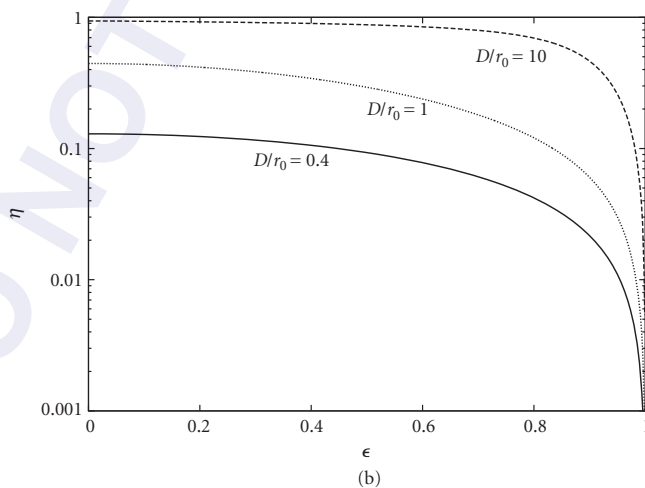
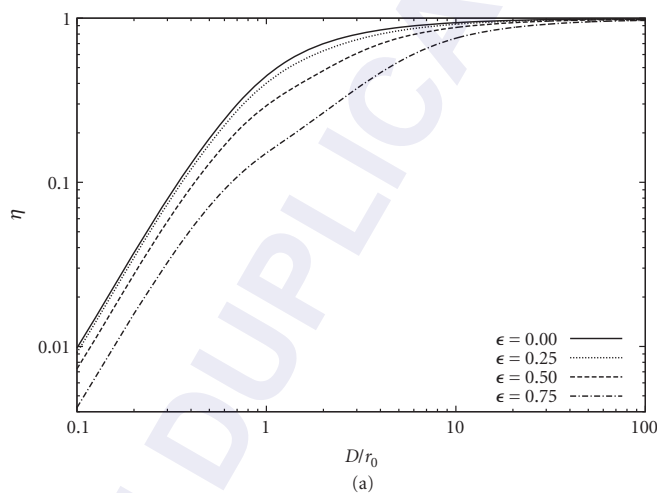
$$\eta(\epsilon; D/r_0) = (1 - \epsilon^2)(D/r_0)^2 \langle S(\epsilon; D/r_0) \rangle \quad (47a)$$

$$= 8(D/r_0)^2 \int_0^1 \langle \tau(v; \epsilon; D/r_0) \rangle \exp[-3.44(vD/r_0)^{5/3}] v dv \quad (47b)$$

Figure 7 shows how η varies with D/r_0 . Its aberration-free (or diffraction-limited) value varies as $(1 - \epsilon^2)(D/r_0)^2$. For very small values of D/r_0 , the atmospheric MTF reduction factor is approximately

TABLE 1 Long (LE) and Short-Exposure (SE) Time-Averaged Strehl Ratio for Various Values of ϵ and D/r_0

ϵ	$D/r_0 = 1$		$D/r_0 = 2$		$D/r_0 = 5$		$D/r_0 = 10$	
	LE	SE	LE	SE	LE	SE	LE	SE
0	0.446	0.889	0.175	0.691	0.035	0.204	0.009	0.023
0.25	0.429	0.890	0.169	0.694	0.036	0.207	0.010	0.024
0.50	0.391	0.892	0.160	0.699	0.040	0.216	0.012	0.027
0.75	0.345	0.889	0.152	0.691	0.050	0.211	0.017	0.033

**FIGURE 7** Variation of long-exposure $\eta(\epsilon; D/r_0)$, representing the time-averaged central irradiance normalized by its aberration-free value for a circular pupil of diameter r_0 but with the same total power as the annular pupil. (a) As a function of D/r_0 and (b) as a function of ϵ .

equal to unity. Accordingly, $\langle S(\epsilon; D/r_0) \rangle$ is also approximately equal to unity, and the aberrated value of η increases with D/r_0 as in the aberration-free case. However, it increases slowly as D/r_0 increases, with a negligible increase beyond a certain value of D/r_0 , depending on the value of ϵ . The contribution to the integral in Eq. (47b) comes from values of v small enough that the factor in the exponent is not vanishingly small. Since, $\tau(v; \epsilon) = 1$ near the origin irrespective of the value of ϵ , Eq. (47b) for large values of D/r_0 may be written

$$\begin{aligned} \eta(\epsilon; D/r_0) &= 8(D/r_0)^2 \int_0^1 \exp[-3.44(vD/r_0)^{5/3}] v dv \\ &= 8(3.44)^{-6/5} (3/5) \int_0^\infty x^{(6/5)-1} \exp(-x) dx \\ &= 1 \end{aligned} \quad (48)$$

where $x = 3.44(vD/r_0)^{5/3}$ and the integral over x is the gamma function $\Gamma(6/5)$. Thus,

$$\eta(\epsilon; D/r_0) \rightarrow 1 \text{ as } D/r_0 \rightarrow \infty \quad (49)$$

independent of the value of ϵ . As is evident from Fig. 7, the saturation effects of atmospheric turbulence occur at larger and larger values of D/r_0 as ϵ increases. The two asymptotes of $\eta(\epsilon; D/r_0)$ for a given value of ϵ intersect at the point given by $(1 - \epsilon^2)(D/r_0)^2 = 1$ or $D/r_0 = (1 - \epsilon^2)^{-1/2}$.¹⁸

The unnormalized aberrated central irradiance is given by

$$\langle I_i(0; \epsilon; D/r_0) \rangle = \frac{\langle P_{\text{ex}} \rangle S_{\text{ex}}(\epsilon)}{\lambda^2 R^2} \langle S(\epsilon; D/r_0) \rangle \quad (50a)$$

$$= \frac{\langle P_{\text{ex}} \rangle S_a}{\lambda^2 R^2} \eta(\epsilon; D/r_0) \quad (50b)$$

where $S_a = \pi r_0^2/4$ is the *coherent area* of the atmosphere. Hence, regardless of how large D is, the central irradiance is less than or equal to the aberration-free central irradiance for a system with a circular exit pupil of diameter r_0 , equality approaching as $D/r_0 \rightarrow \infty$. The limiting value of the central irradiance is independent of the value of ϵ . Since $S_a \sim r_0^2 \sim \lambda^{12/5}$, the limiting value varies with wavelength as $\lambda^{2/5}$. It is evident from Eq. (50b) that η represents the aberrated central irradiance normalized by the aberration-free value $\langle P_{\text{ex}} \rangle S_a / \lambda^2 R^2$ for a circular pupil of diameter r_0 .

In astronomical observations, the time-averaged image power given by $\langle P_{\text{ex}} \rangle = \pi(1 - \epsilon^2)D^2 \langle I_0 \rangle$, where $\langle I_0 \rangle$ is the time-averaged irradiance across the exit pupil, increases as D increases. However, if the observation is made against a uniform background, then the background irradiance in the image also increases as D^2 . Hence, regardless of the value of ϵ , the detectability of a point object is limited by turbulence to a value corresponding to an exit pupil of diameter r_0 , no matter how large the diameter D of the exit pupil is. In the case of a laser transmitter with a fixed value of laser power P_{ex} , the central irradiance on a target will again be limited to its aberration-free value for an exit pupil of diameter r_0 , regardless of how large its beam diameter is. Similarly, the ratio of the signal and noise powers in an optical heterodyne detection of a turbulence-degraded signal is limited to the aberration-free value corresponding to an exit pupil of diameter r_0 .

Figure 8 shows how the time-averaged irradiance distribution or the PSF and encircled power change as D/r_0 increases for several values of ϵ . The PSFs are normalized to unity at the center, the actual central value being the LE Strehl ratio given in Table 1. As ϵ increases, power flows from the central bright spot into the diffraction rings. As D/r_0 increases, the diffraction rings disappear and the PSFs become smooth, and a given fraction of total power is contained in a circle of larger and larger radius.

4.6 MODAL EXPANSION OF ABERRATION FUNCTION

So far we have considered the degradation of system performance by propagation through atmospheric turbulence without any knowledge of the magnitude or the type of the phase aberrations introduced by it. In this section, we expand the Kolmogorov turbulence-degraded aberration function into a complete set of orthonormal polynomials and determine its time-averaged variance as well as the variance and covariance of the expansion coefficients.

We expand the aberration function $\Phi(\rho, \theta; \epsilon)$ in terms of a complete set of Zernike annular polynomials $Z_j(\rho, \theta; \epsilon)$ that are orthonormal over the unit annulus in the form⁵⁻⁸

$$\Phi(\rho, \theta; \epsilon) = \sum_j a_j(\epsilon) Z_j(\rho, \theta; \epsilon) \quad (51)$$

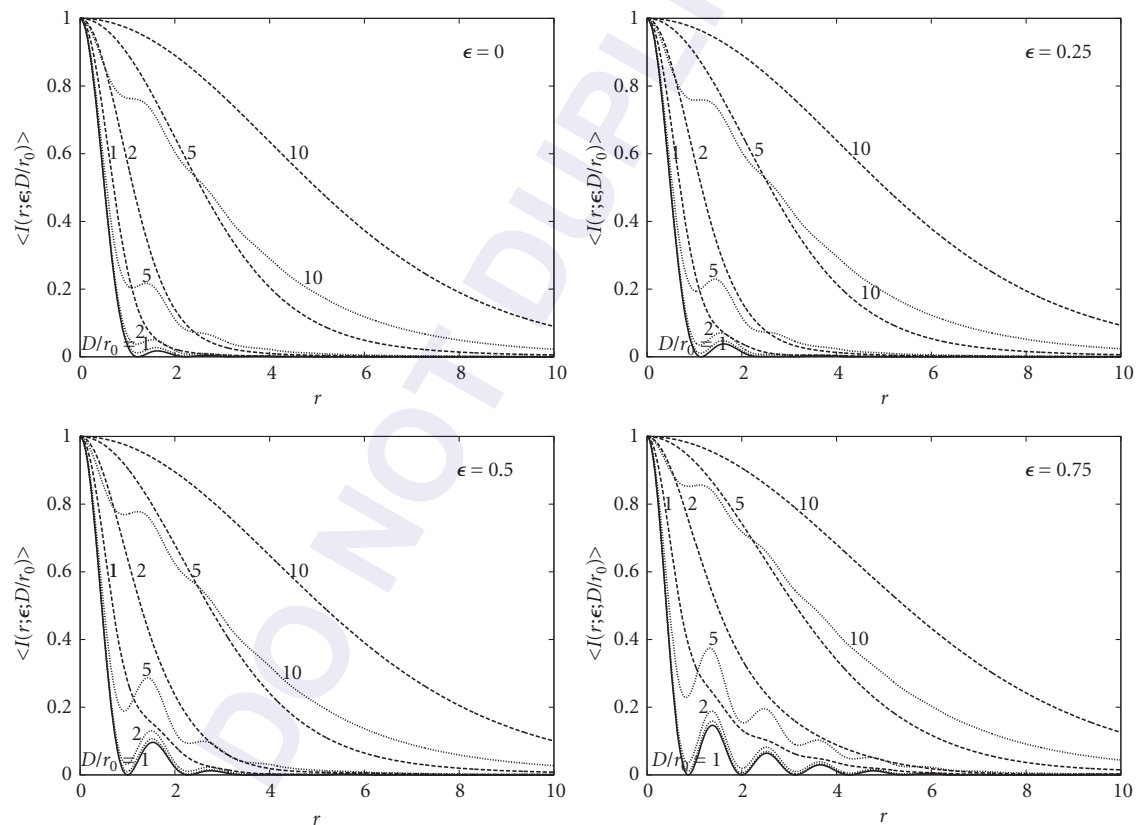


FIGURE 8 Time-averaged irradiance $\langle I(r; \epsilon; D/r_0) \rangle$ and encircled power $\langle P(r; \epsilon; D/r_0) \rangle$ distributions for circular and annular pupils. The solid curves represent the aberration-free distributions, and the dashed and dotted curves represent the corresponding long- and short-exposure distributions.

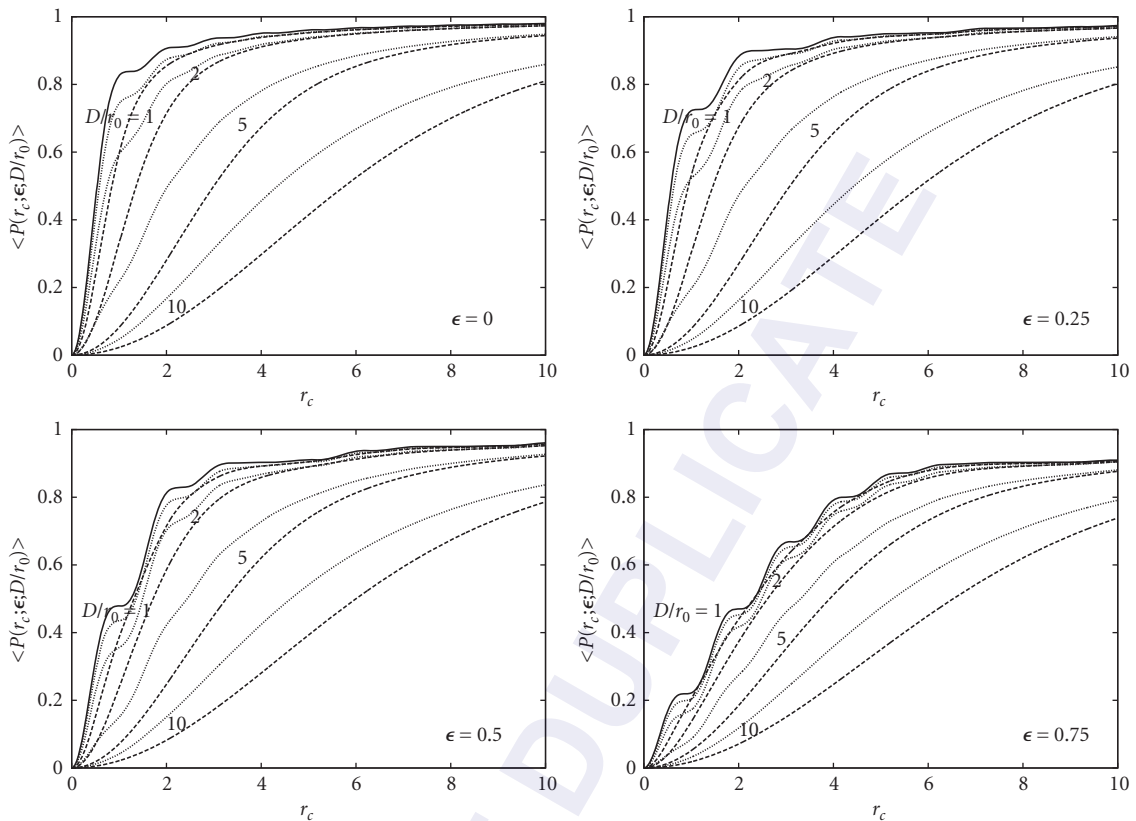


FIGURE 8 (Continued)

where $a_j(\epsilon)$ are the expansion coefficients of the polynomials, $\epsilon \leq \rho \leq 1$ and $0 \leq \theta < 2\pi$. Since the phase aberration is a Gaussian random variable with a zero time-averaged value, so are the expansion coefficients. The annular polynomials are given by

$$Z_{\text{even}j}(\rho, \theta; \epsilon) = \sqrt{2(n+1)}R_n^m(\rho; \epsilon)\cos m\theta, \quad m \neq 0 \quad (52a)$$

$$Z_{\text{odd}j}(\rho, \theta; \epsilon) = \sqrt{2(n+1)}R_n^m(\rho; \epsilon)\sin m\theta, \quad m \neq 0 \quad (52b)$$

$$Z_j(\rho, \theta; \epsilon) = \sqrt{n+1}R_n^0(\rho; \epsilon), \quad m = 0 \quad (52c)$$

where n and m are positive integers (including zero) and $n - m \geq 0$ and even.

The index n represents the radial degree or the order of a polynomial, since it represents the highest power of ρ in the polynomial, and m is called the azimuthal frequency. The index j is a polynomial-ordering number and is a function of both n and m . The polynomials are ordered such that an even j corresponds to a symmetric polynomial varying as $\cos m\theta$, while an odd j corresponds to an antisymmetric polynomial varying as $\sin m\theta$. A polynomial with a lower value of n is ordered first, while for a given value of n , a polynomial with a lower value of m is ordered first.

TABLE 2 Relationship Among the Polynomial Indices n , m , and j

n	m	j	n	m	j
0	0	1	6	0	22
1	1	2, 3	6	2	23, 24
2	0	4	6	4	25, 26
2	2	5, 6	6	6	27, 28
3	1	7, 8	7	1	29, 30
3	3	9, 10	7	3	31, 32
4	0	11	7	5	33, 34
4	2	12, 13	7	7	35, 36
4	4	14, 15	8	0	37
5	1	16, 17	8	2	38, 39
5	3	18, 19	8	4	40, 41
5	5	20, 21	8	6	42, 43
			8	8	44, 45

The relationships among the indices j , n , and m are given in Table 2. For example, when $n = 4$ and $m = 2$, then $j = 12$ for the $\cos 2\theta$ polynomial and $j = 13$ for the $\sin 2\theta$ polynomial. The number of polynomials for a given value of n is $n + 1$, and the number of polynomials up to and including a certain order n is given by

$$N_n = (n + 1)(n + 2)/2 \quad (53)$$

The radial annular polynomials $R_n^m(\rho; \epsilon)$ obey the orthogonality relation

$$\int_{\epsilon}^1 R_n^m(\rho; \epsilon) R_{n'}^m(\rho; \epsilon) \rho \, d\rho = \frac{1 - \epsilon^2}{2(n + 1)} \delta_{nn'} \quad (54)$$

where $\delta_{nn'}$ is a Kronecker delta. Accordingly, the annular polynomials obey the orthonormality condition

$$\int_{\epsilon}^1 \int_0^{2\pi} Z_j(\rho, \theta; \epsilon) Z_{j'}(\rho, \theta; \epsilon) \rho \, d\rho \, d\theta \bigg/ \int_{\epsilon}^1 \int_0^{2\pi} \rho \, d\rho \, d\theta = \delta_{jj'} \quad (55)$$

The Zernike expansion coefficients are given by

$$a_j(\epsilon) = \frac{1}{\pi(1 - \epsilon^2)} \int_{\epsilon}^1 \int_0^{2\pi} \Phi(\rho, \theta; \epsilon) Z_j(\rho, \theta; \epsilon) \rho \, d\rho \, d\theta \quad (56)$$

as may be seen by substituting Eq. (51) into Eq. (56) and using the orthonormality Eq. (55).

The first eleven annular polynomials are listed in Table 3, where

$$R_1^1(\rho; \epsilon) = \frac{\rho}{(1 + \epsilon^2)^{1/2}} \quad (57a)$$

$$R_2^0(\rho; \epsilon) = \frac{(2\rho^2 - 1 - \epsilon^2)}{(1 - \epsilon^2)} \quad (57b)$$

$$R_2^2(\rho; \epsilon) = \frac{\rho^2}{(1 + \epsilon^2 + \epsilon^4)^{1/2}} \quad (57c)$$

TABLE 3 Orthonormal Zernike Annular Polynomials $Z_j(\rho, \theta; \epsilon)$

j	n	m	$Z_j(\rho, \theta; \epsilon)$	Aberration Name*
1	0	0	$R_0^0(\rho; \epsilon) = 1$	Piston
2	1	1	$2R_1^1(\rho; \epsilon) \cos \theta$	x tilt
3	1	1	$2R_1^1(\rho; \epsilon) \sin \theta$	y tilt
4	2	0	$\sqrt{3}R_2^0(\rho; \epsilon)$	Defocus
5	2	2	$\sqrt{6}R_2^2(\rho; \epsilon) \sin 2\theta$	Primary astigmatism at 45°
6	2	2	$\sqrt{6}R_2^2(\rho; \epsilon) \cos 2\theta$	Primary astigmatism at 0°
7	3	1	$\sqrt{8}R_3^1(\rho; \epsilon) \sin \theta$	Primary y coma
8	3	1	$\sqrt{8}R_3^1(\rho; \epsilon) \cos \theta$	Primary x coma
9	3	3	$\sqrt{8}R_3^3(\rho; \epsilon) \sin 3\theta$	
10	3	3	$\sqrt{8}R_3^3(\rho; \epsilon) \cos 3\theta$	
11	4	0	$\sqrt{5}R_4^0(\rho; \epsilon)$	Primary spherical

*The name must precede with "orthonormal annular," e.g., orthonormal annular primary x coma.

$$R_3^1(\rho; \epsilon) = \frac{3(1 + \epsilon^2)\rho^3 - 2(1 + \epsilon^2 + \epsilon^4)\rho}{(1 - \epsilon^2)[(1 + \epsilon^2)(1 + 4\epsilon^2 + \epsilon^4)]^{1/2}} \quad (57d)$$

$$R_3^3(\rho; \epsilon) = \frac{\rho^3}{(1 + \epsilon^2 + \epsilon^4 + \epsilon^6)^{1/2}} \quad (57e)$$

and

$$R_4^0(\rho; \epsilon) = \frac{6\rho^4 - 6(1 + \epsilon^2)\rho^2 + 1 + 4\epsilon^2 + \epsilon^4}{(1 - \epsilon^2)^2} \quad (57f)$$

Of course, they reduce to the corresponding circle polynomials as $\epsilon \rightarrow 0$.

4.7 COVARIANCE AND VARIANCE OF EXPANSION COEFFICIENTS

From Eq. (56), the covariance or cross-correlation of the expansion coefficients may be written

$$\langle a_j(\epsilon) a_{j'}(\epsilon) \rangle = \frac{1}{\pi^2(1 - \epsilon^2)^2} \iint Z_j(\rho, \epsilon) Z_{j'}(\rho'; \epsilon) \langle \Phi(\rho; \epsilon) \Phi(\rho'; \epsilon) \rangle d\rho d\rho' \quad (58)$$

where ρ is the normalized position vector of a pupil point. It can be obtained by using the Fourier transforms of the annular polynomials.⁹ The autocorrelation, that is, the mean square value or the variance of an expansion coefficient, may be obtained from Eq. (58) by letting $j = j'$, that is, by letting $n = n'$ and $m = m'$.

A closed-form analytical solution is obtained when $m = n = n'$

$$\begin{aligned} \langle [a_n^n(\epsilon)]^2 \rangle &= \frac{0.0229(n+1)\Gamma(n-5/6)\pi^{8/3}}{2^{5/3}\Gamma(17/6)(1-\epsilon^2)[1-\epsilon^{2(n+1)}]} \left(\frac{D}{r_0}\right)^{5/3} \\ &\times \left[\frac{(1+\epsilon^{2n+17/3})\Gamma(14/3)}{\Gamma(17/6)\Gamma(n+23/6)} - \frac{2\epsilon^{2(n+1)}}{(n+1)!} F_1\left(n - \frac{5}{6}, -\frac{11}{6}; n+2; \epsilon^2\right) \right] \end{aligned} \quad (59)$$

where $F_1(a, b; c; z)$ is a hypergeometric function. Since $r_0 \propto \lambda^{6/5}$, the variance in terms of the optical path-length errors [obtained by multiplying $\langle [a_n^n(\epsilon)]^2 \rangle$ by $(\lambda/2\pi)^2$] is independent of λ , as expected in the absence of atmospheric dispersion.

For a circular pupil,¹⁹ the following results can be obtained from the corresponding results for an annular pupil by letting $\epsilon \rightarrow 0$:

$$\begin{aligned} \langle a_j a_{j'} \rangle &= 0.1534(-1)^{(n+n'-2m)/2} [(n+1)(n'+1)]^{1/2} (D/r_0)^{5/3} \delta_{mm'} \\ &\times \frac{\Gamma(14/3)\Gamma[(n+n'-5/3)/2]}{\Gamma[(n-n'+17/3)/2]\Gamma[(n'-n+17/3)/2]\Gamma[(n+n'+23/3)/2]}, \quad n, n' \neq 0 \end{aligned} \quad (60)$$

and

$$\langle a_j^2 \rangle = 0.7587 \frac{\Gamma(n-5/6)}{\Gamma(n+23/6)} (n+1) \left(\frac{D}{r_0} \right)^{5/3}, \quad n \neq 0 \quad (61)$$

Numerical values of variance and covariance are given in Tables 4 and 5 for various values of ϵ including zero. We note that when $\epsilon = 0$, the variance of an aberration coefficient depends on

TABLE 4 Variance $\langle a_j^2(\epsilon) \rangle$ of Expansion Coefficients in Units of $(D/r_0)^{5/3}$

$\langle a_j^2 \rangle$	$\epsilon=0$	0.05	0.1	0.15	0.2	0.25	0.35	0.5	0.75
$\langle a_2^2 \rangle, \langle a_3^2 \rangle$	0.44888	0.45000	0.45336	0.45893	0.46670	0.47660	0.50260	0.55610	0.67917
$\langle a_4^2 \rangle$	0.02322	0.02302	0.02245	0.02153	0.02033	0.01891	0.01560	0.01028	0.00292
$\langle a_5^2 \rangle, \langle a_6^2 \rangle$	0.02322	0.02328	0.02345	0.02375	0.02418	0.02475	0.02635	0.02997	0.03888
$\langle a_7^2 \rangle, \langle a_8^2 \rangle$	0.00619	0.00621	0.00624	0.00628	0.00630	0.00626	0.00593	0.00468	0.00162
$\langle a_9^2 \rangle, \langle a_{10}^2 \rangle$	0.00619	0.00621	0.00625	0.00633	0.00645	0.00660	0.00705	0.00815	0.01132
$\langle a_{11}^2 \rangle$	0.00245	0.00241	0.00227	0.00207	0.00183	0.00158	0.00109	0.00052	0.00007
$\langle a_{12}^2 \rangle, \langle a_{13}^2 \rangle$	0.00245	0.00246	0.00248	0.00251	0.00255	0.00260	0.00266	0.00247	0.00108
$\langle a_{14}^2 \rangle, \langle a_{15}^2 \rangle$	0.00245	0.00246	0.00248	0.00251	0.00256	0.00262	0.00280	0.00326	0.00480

TABLE 5 Covariance $\langle a_j(\epsilon) a_{j'}(\epsilon) \rangle$ of Expansion Coefficients in Units of $(D/r_0)^{5/3}$

$\langle a_j a_{j'} \rangle$	$\epsilon = 0$	0.05	0.1	0.15	0.2	0.25	0.35	0.5	0.75
$\langle a_2 a_8 \rangle, \langle a_3 a_7 \rangle$	-0.01416	-0.01420	-0.01430	-0.01444	-0.01460	-0.01474	-0.01481	-0.01392	-0.00885
$\langle a_4 a_{11} \rangle$	-0.00388	-0.00381	-0.00363	-0.00334	-0.00300	-0.00262	-0.00186	-0.00091	-0.00012
$\langle a_4 a_{22} \rangle$	0.00032	0.00030	0.00024	0.00017	0.00009	0.00002	-0.00008	-0.00011	-0.00003
$\langle a_5 a_{13} \rangle, \langle a_6 a_{12} \rangle$	-0.00388	-0.00389	-0.00392	-0.00397	-0.00404	-0.00412	-0.00433	-0.00451	-0.00341
$\langle a_7 a_{17} \rangle, \langle a_8 a_{16} \rangle$	-0.00156	-0.00158	-0.00164	-0.00172	-0.00178	-0.00173	-0.00095	0.00365	0.01180
$\langle a_{11} a_{22} \rangle$	-0.00076	-0.00073	-0.00067	-0.00058	-0.00048	-0.00038	-0.00022	-0.00007	0.00000
$\langle a_9 a_{19} \rangle, \langle a_{10} a_{18} \rangle$	-0.00156	-0.00156	-0.00157	-0.00159	-0.00162	-0.00166	-0.00176	-0.00196	-0.00177
$\langle a_{14} a_{26} \rangle, \langle a_{15} a_{25} \rangle$	-0.00076	-0.00076	-0.00077	-0.00078	-0.00079	-0.00081	-0.00086	-0.00099	-0.00104

the value of n , but is independent of the value of m . Thus, for example, $\langle a_4^2 \rangle = \langle a_5^2 \rangle = \langle a_6^2 \rangle$, or $\langle a_7^2 \rangle = \langle a_8^2 \rangle = \langle a_9^2 \rangle = \langle a_{10}^2 \rangle$. However, when $\epsilon \neq 0$, the variance depends on the values of both n and m , but it is independent of whether it is a cosine or a sine mode. Thus, for example, $\langle a_4^2(\epsilon) \rangle \neq \langle a_5^2(\epsilon) \rangle = \langle a_6^2(\epsilon) \rangle$, or $\langle a_7^2(\epsilon) \rangle = \langle a_8^2(\epsilon) \rangle \neq \langle a_9^2(\epsilon) \rangle = \langle a_{10}^2(\epsilon) \rangle$. Similarly, the covariance for given values of n and n' are independent of their m values when $\epsilon = 0$. For example, $\langle a_4 a_{11} \rangle = \langle a_5 a_{13} \rangle = \langle a_6 a_{12} \rangle$, or $\langle a_7 a_{17} \rangle = \langle a_8 a_{16} \rangle = \langle a_9 a_{19} \rangle = \langle a_{10} a_{18} \rangle$. However, when $\epsilon \neq 0$, equality is obtained only when the m values are also equal. For example, $\langle a_4(\epsilon) a_{11}(\epsilon) \rangle \neq \langle a_5(\epsilon) a_{13}(\epsilon) \rangle = \langle a_6(\epsilon) a_{12}(\epsilon) \rangle$, or $\langle a_7(\epsilon) a_{17}(\epsilon) \rangle = \langle a_8(\epsilon) a_{16}(\epsilon) \rangle \neq \langle a_9(\epsilon) a_{19}(\epsilon) \rangle = \langle a_{10}(\epsilon) a_{18}(\epsilon) \rangle$. The first nonzero cross-correlations are the tilt-coma cross-correlations $\langle a_2(\epsilon) a_8(\epsilon) \rangle$ and $\langle a_3(\epsilon) a_7(\epsilon) \rangle$. Moreover, for a given value of n , the correlation values decrease rapidly as the order difference $n' - n$ increases.

Figures 9 and 10 show the variance and covariance of Zernike coefficients as a function of ϵ for some low-order terms. The variance of terms with $m = n$ increase monotonically with an increasing value of ϵ , e.g., tip and tilt variances $\langle a_2^2 \rangle$ and $\langle a_3^2 \rangle$, or astigmatism variances $\langle a_5^2 \rangle$ and $\langle a_6^2 \rangle$. For other terms, e.g., coma $\langle a_7^2 \rangle$ and $\langle a_8^2 \rangle$, or defocus $\langle a_4^2 \rangle$ and spherical aberration $\langle a_{11}^2 \rangle$, they generally decrease monotonically and approach zero as $\epsilon \rightarrow 1$. Similarly, the covariances of tilt with coma, for example, $\langle a_2(\epsilon) a_8(\epsilon) \rangle$, and defocus with primary or secondary spherical aberration, that is, $\langle a_4(\epsilon) a_{11}(\epsilon) \rangle$ or $\langle a_4(\epsilon) a_{22}(\epsilon) \rangle$, approach zero as $\epsilon \rightarrow 1$.

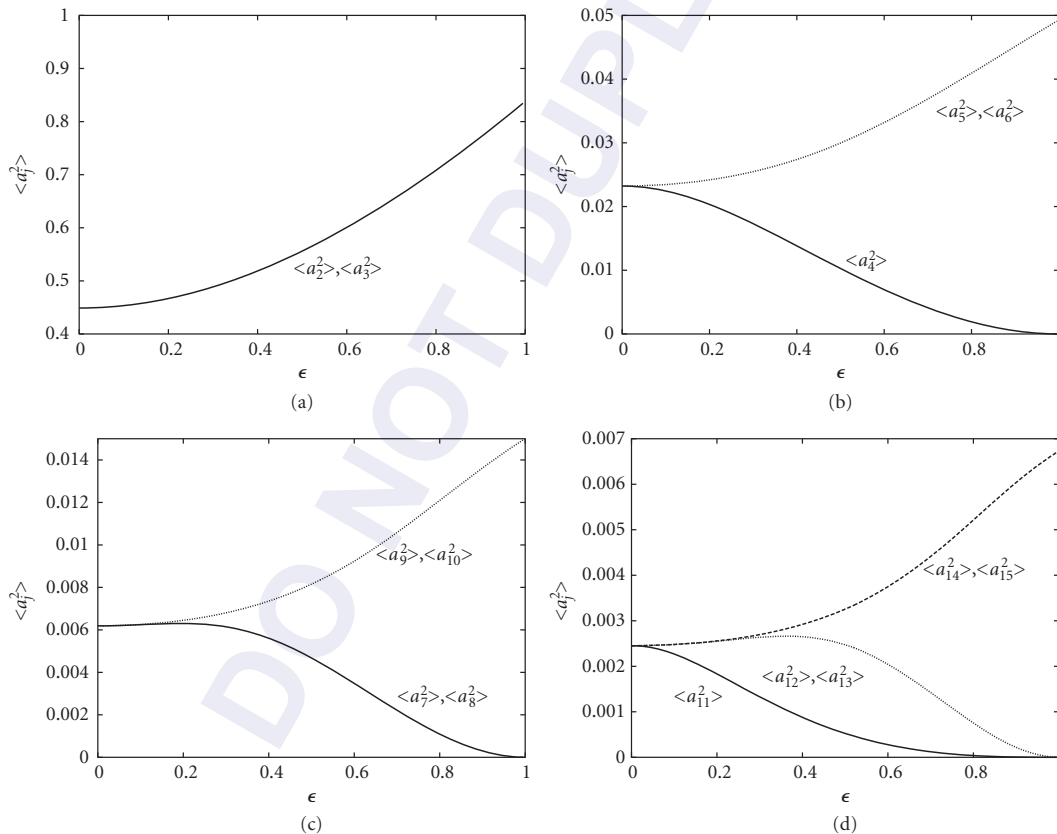


FIGURE 9 Variance $\langle a_j^2(\epsilon) \rangle$ of Zernike coefficients. (a) Tip and tilt; (b) defocus and astigmatism; (c) coma and tilt, etc.; and (d) spherical and higher orders.

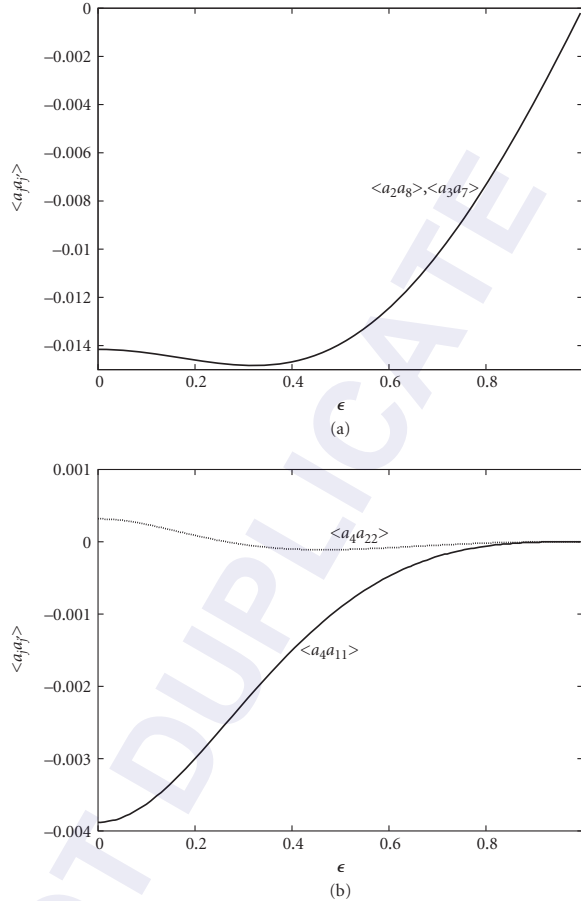


FIGURE 10 Covariance $\langle a_j(\epsilon)a_l(\epsilon) \rangle$ of Zernike coefficients. (a) Tilt and coma and (b) defocus and spherical.

4.8 ANGLE OF ARRIVAL FLUCTUATIONS

Depending on the measurement process, the angle of arrival can be defined by the Zernike tilts of the wavefront or the centroid of the aberrated PSF. The first is referred to as the Z-tilt and represents the least square fit of the linear terms to the aberration function or equivalently the gradient of the wavefront at the center of the pupil. The PSF centroid, however, represents the mean value of the wavefront gradient across the pupil. The tilt associated with the centroid is referred to as the G- or the gradient-tilt.^{1,20}

From Eq. (51), the wavefront tilt aberration introduced by turbulence is given by

$$\begin{aligned}
 \Phi_t(\rho, \theta; \epsilon) &= a_2(\epsilon)Z_2(\rho, \theta; \epsilon) + a_3(\epsilon)Z_3(\rho, \theta; \epsilon) \\
 &= 2[a_2(\epsilon)R_1^1(\rho; \epsilon)\cos\theta + a_3(\epsilon)R_1^1(\rho; \epsilon)\sin\theta] \\
 &= \frac{2}{(1 + \epsilon^2)^{1/2}}[a_2(\epsilon)\rho\cos\theta + a_3(\epsilon)\rho\sin\theta]
 \end{aligned} \tag{62}$$

Or, in Cartesian coordinates,

$$\Phi_t(x, y; \epsilon) = \frac{2}{(1 + \epsilon^2)^{1/2}} [a_2(\epsilon)x + a_3(\epsilon)y] \quad (63)$$

The wavefront tilt displaces the image to a point $[x_t(\epsilon), y_t(\epsilon)]$ given by²¹

$$\begin{aligned} [x_t(\epsilon), y_t(\epsilon)] &= 2F \frac{\lambda}{2\pi} \left(\frac{\partial \Phi_t}{\partial x}, \frac{\partial \Phi_t}{\partial y} \right) \\ &= \frac{2\lambda F}{\pi(1 + \epsilon^2)^{1/2}} [a_2(\epsilon), a_3(\epsilon)] \end{aligned} \quad (64)$$

Accordingly, the angle of arrival of the wave is given by

$$\begin{aligned} [\alpha(\epsilon), \beta(\epsilon)] &= \frac{1}{R} [x_t(\epsilon), y_t(\epsilon)] \\ &= \frac{2\lambda}{\pi D(1 + \epsilon^2)^{1/2}} [a_2(\epsilon), a_3(\epsilon)] \end{aligned} \quad (65)$$

Or, the mean square fluctuation of the angle of arrival is given by

$$[\sigma_\alpha^2(\epsilon), \sigma_\beta^2(\epsilon)] = \left(\frac{2\lambda}{\pi D} \right)^2 \frac{1}{(1 + \epsilon^2)} \langle a_2^2(\epsilon), a_3^2(\epsilon) \rangle \quad (66)$$

where from Eq. (59)

$$\begin{aligned} \langle a_2^2(\epsilon) \rangle &= \langle a_3^2(\epsilon) \rangle \\ &= \frac{0.9858}{(1 + \epsilon^2)(1 - \epsilon^2)^2} \left[0.4554(1 + \epsilon^{2/3}) - \epsilon^4 F_1 \left(\frac{1}{6}, -\frac{11}{6}; 3; \epsilon^2 \right) \right] \left(\frac{D}{r_0} \right)^{5/3} \end{aligned} \quad (67)$$

Accordingly, the standard deviation of the angular fluctuations based on the Zernike tilts is given by

$$[\sigma_\alpha(\epsilon)]_Z = [\sigma_\beta(\epsilon)]_Z = \frac{2}{\pi(1 + \epsilon^2)^{1/2}} \frac{\lambda}{D} \left[\langle a_2^2(\epsilon) \rangle \right]^{1/2} \quad (68)$$

Since $\langle a_2^2(\epsilon) \rangle \sim (D/r_0)^{5/3}$, Eq. (68) shows that the mean square fluctuation of the image displacement decreases with the pupil diameter D as $D^{-1/3}$. Moreover, since $r_0 \sim \lambda^{6/5}$, we find that the fluctuation is independent of the wavelength, as expected in the absence of any atmospheric dispersion. For a circular pupil, substituting $\langle a_2^2 \rangle = \langle a_3^2 \rangle = 0.449(D/r_0)^{5/3}$, we obtain

$$\begin{aligned} (\sigma_\alpha)_Z = (\sigma_\beta)_Z &= \frac{2}{\pi} \frac{\lambda}{D} \left[\langle a_2^2 \rangle \right]^{1/2} \\ &= 0.4265 \frac{\lambda}{D} \left(\frac{D}{r_0} \right)^{5/6} \end{aligned} \quad (69)$$

Since the wavefront tilt produced by turbulence is independent of the shape of the pupil, except for its correlation with the coma terms, $\langle a_2^2(\epsilon) \rangle$ increases with ϵ approximately as $1 + \epsilon^2$. This is illustrated in Fig. 11, where $\langle a_2^2(\epsilon) \rangle$ and $(1 + \epsilon^2)\langle a_2^2 \rangle$ are compared with each other as a function of ϵ . Accordingly, the angular fluctuation depends very weakly on the obscuration ratio. For example, when $\epsilon = 0.5$, the constant 0.4265 on the right-hand side of Eq. (69) is replaced by 0.4246.

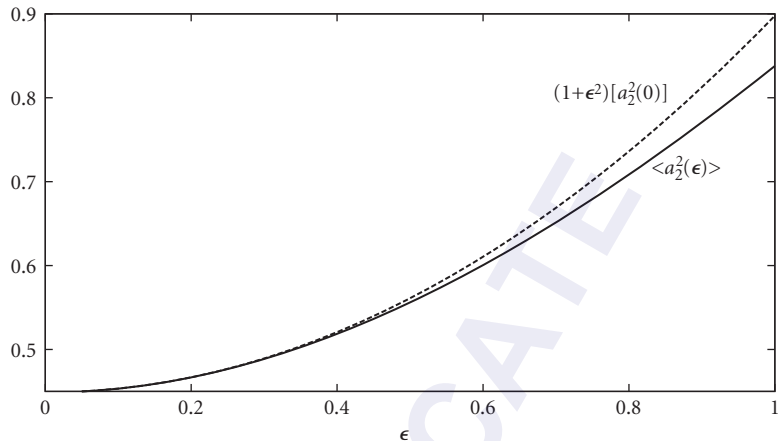


FIGURE 11 Comparison of tilt variance $\langle a_2^2(\epsilon) \rangle$ with $(1 + \epsilon^2)\langle a_2^2(0) \rangle$.

Tatarski gives an expression for the fluctuations of the angle of arrival based on the centroid of the aberrated PSF for a circular pupil.²² His analysis can be extended to systems with annular pupils. However, we consider the centroid in terms of the Zernike annular-polynomial expansion coefficients of the aberration function. This approach gives insight into why the G-tilt is different from the Z-tilt, why the two tilts are approximately equal to each other, and why the G-tilt is smaller.

Neglecting the amplitude variations introduced by turbulence, the centroid (x_c, y_c) of an aberrated PSF is given by the mean value of the gradient of the aberration function^{11,23}

$$(x_c, y_c) = \frac{R}{S_p} \iint \left[\frac{\partial W(x, y)}{\partial x}, \frac{\partial W(x, y)}{\partial y} \right] dx dy \quad (70)$$

where $W(x, y)$ is the wave aberration at a point (x, y) on a pupil of area S_p . If we write the aberration function in terms of the Zernike annular polynomials in the form

$$W(\rho, \theta; \epsilon) = \sum_{n=0}^{\infty} \sum_{m=0}^n [2(n+1)/(1 + \delta_{m0})]^{1/2} R_n^m(\rho; \epsilon) [c_{nm}(\epsilon) \cos m\theta + s_{nm}(\epsilon) \sin m\theta] \quad (71)$$

where $c_{nm}(\epsilon)$ and $s_{nm}(\epsilon)$ are the Zernike annular expansion coefficients, it can be shown that the centroid of the corresponding PSF is given by^{11,23}

$$[x_c(\epsilon), y_c(\epsilon)] = \frac{2F}{1 - \epsilon^2} \sum_{n=1}^{\infty} \sqrt{2(n+1)} [R_n^1(1; \epsilon) - \epsilon R_n^1(\epsilon; \epsilon)] [c_{n1}(\epsilon), s_{n1}(\epsilon)] \quad (72)$$

where a prime on the summation sign indicates a summation over only the odd integral values of n . Thus, the only aberrations that contribute to the centroid are those with $m = 1$. Aberrations of the type $R_n^1(\rho; \epsilon) \cos \theta$ contribute to x_c and those of the type $R_n^1(\rho; \epsilon) \sin \theta$ contribute to y_c , which also follows from the symmetry considerations of the aberrations. Keeping the tilt and only the primary coma terms (and thus neglecting higher orders of coma owing to their small magnitudes) in Eq. (72), we obtain

$$[\sigma_{\alpha}(\epsilon)]_G = [\sigma_{\beta}(\epsilon)]_G = \frac{1}{\pi(1 - \epsilon^2)} \lambda \left[4A^2(\epsilon) \langle a_2^2(\epsilon) \rangle + 8B^2(\epsilon) \langle a_8^2(\epsilon) \rangle + 8\sqrt{2} \frac{1 - \epsilon^2}{(1 + \epsilon^2)^{1/2}} B(\epsilon) \langle a_2(\epsilon) a_8(\epsilon) \rangle \right]^{1/2} \quad (73)$$

where, from Eqs. (57a) and (57d),

$$A(\epsilon) = R_1^1(1; \epsilon) - \epsilon R_1^1(\epsilon; \epsilon) = \frac{1 - \epsilon^2}{(1 + \epsilon^2)^{1/2}} \quad (74a)$$

and

$$B(\epsilon) = R_3^1(1; \epsilon) - \epsilon R_3^1(\epsilon; \epsilon) = \frac{1 + 3\epsilon^2 - 3\epsilon^4 - \epsilon^6}{(1 - \epsilon^2)[(1 + \epsilon^2)(1 + 4\epsilon^2 + \epsilon^4)]^{1/2}} \quad (74b)$$

If we retain only the Zernike tilt terms, we obtain the Z-tilt given by Eq. (68).

For a circular pupil, $A(\epsilon)$ and $B(\epsilon)$ both reduce to unity, and Eq. (73) reduces to

$$(\sigma_\alpha)_G = (\sigma_\beta)_G = \frac{\lambda}{\pi D} \left[4\langle a_2^2 \rangle + 8\langle a_8^2 \rangle + 8\sqrt{2}\langle a_2 a_8 \rangle \right]^{1/2} \quad (75)$$

Substituting for $\langle a_2^2 \rangle$, $\langle a_8^2 \rangle$, and $\langle a_2 a_8 \rangle$ from Tables 4 and 5, we obtain

$$(\sigma_\alpha)_G = 0.4132 \frac{\lambda}{D} \left(\frac{D}{r_0} \right)^{5/6} \quad (76)$$

The dependence of G- and Z-tilts on ϵ is illustrated in Fig. 12. The G-tilt is smaller than the Z-tilt by approximately 3 percent. It is not surprising that the Z- and G-tilts are approximately equal to each other. Since the Zernike tilt aberrations account for 87 percent or more of the aberration variance (see Sec. 4.11 and Fig. 16), the variance of the coma aberration(s) is relatively small. The G-tilt is smaller than the Z-tilt because of the negative correlations of the tilt coefficients with the coma coefficients.

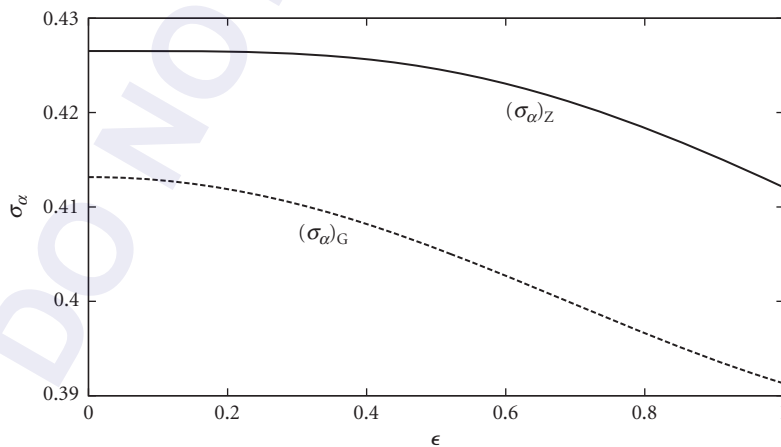


FIGURE 12 Dependence of the angle of arrival fluctuations on the obscuration ratio ϵ , illustrating the small difference between the Zernike wavefront tilt $[\sigma_\alpha(\epsilon)]_Z$ and the centroid-based G-tilt $[\sigma_\alpha(\epsilon)]_G$.

4.9 ABERRATION VARIANCE AND APPROXIMATE STREHL RATIO

The piston autocorrelation $\langle a_1^2(\epsilon) \rangle$ for Kolmogorov turbulence, obtained by letting $j = j' = 1$ in Eq. (58), so that both n and n' are equal to zero, is infinity.⁹ The mean square value $\langle \Phi^2(\epsilon) \rangle$ of the aberration is also infinite, but the difference of the two is finite. Since the mean value $\langle \Phi(\epsilon) \rangle$ of the aberration is zero, $\langle \Phi^2(\epsilon) \rangle$ is also the variance of the aberration. However, the piston aberration $Z_1(\rho, \theta; \epsilon)$, being a constant, does not affect the aberration variance or the image quality. Accordingly, we refer to the piston-removed variance as simply the aberration variance. It is given by

$$\begin{aligned} \sigma_\Phi^2(\epsilon) &\equiv \langle \Phi^2(\epsilon) \rangle - \langle a_1^2(\epsilon) \rangle \\ &= \frac{0.046\pi^{11/3}}{2^{5/3}\Gamma(17/6)\Gamma(11/6)(1-\epsilon^2)^2} \left[\frac{\Gamma(14/3)(1+\epsilon)^{17/3}}{\Gamma(17/6)\Gamma(23/6)} - 2\epsilon^2 F_1\left(-\frac{5}{6}, -\frac{11}{6}, 2; \epsilon^2\right) \right] \left(\frac{D}{r_0}\right)^{5/3} \end{aligned} \quad (77)$$

Letting $\epsilon = 0$ for a circular pupil, Eq. (77) reduces to

$$\begin{aligned} \sigma_\Phi^2 &\equiv \langle \Phi^2 \rangle - \langle a_1^2 \rangle \\ &= 1.0324(D/r_0)^{5/3} \end{aligned} \quad (78)$$

Figure 13 shows how $\Delta_1(\epsilon) \equiv \sigma_\Phi^2(\epsilon)$ in units of $(D/r_0)^{5/3}$ varies with ϵ . Its value increases monotonically from 1.0324 for a circular pupil ($\epsilon = 0$) to 1.843 for an infinitesimally thin ring pupil ($\epsilon \rightarrow 1$).

The approximate expression $S = \exp(-\sigma_\Phi^2)$ for the Strehl ratio^{11,24,25} is not suitable for calculating the time-averaged Strehl ratio for a random aberration, especially when its value is small. For example, the Strehl ratio given by

$$\langle S_1(D/r_0) \rangle \approx \exp[-1.03(D/r_0)^{5/3}] \quad (79)$$

is illustrated in Fig. 14 by the dashed curve, which is quite steep owing to the $(D/r_0)^{5/3}$ dependence of σ_Φ^2 . Even for a small value of $D/r_0 = 1$, it gives a Strehl ratio of 0.357, compared to a true value of

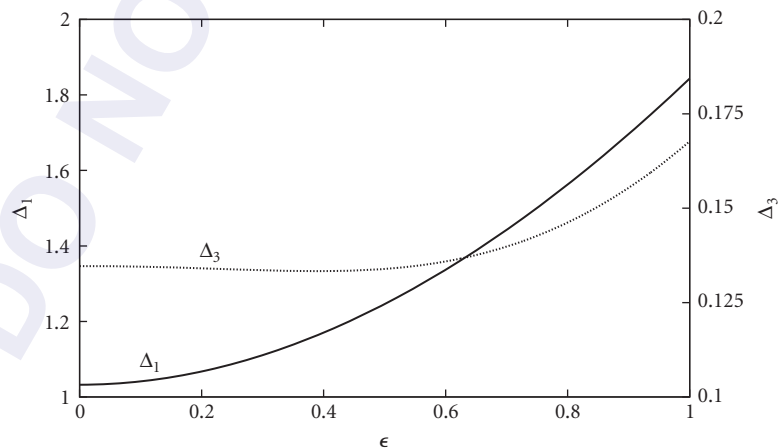


FIGURE 13 Piston-removed aberration variance $\Delta_1(\epsilon) \equiv \sigma_\Phi^2(\epsilon)$ and tilt-corrected variance $\Delta_3(\epsilon)$ in units of $(D/r_0)^{5/3}$ a function of ϵ .

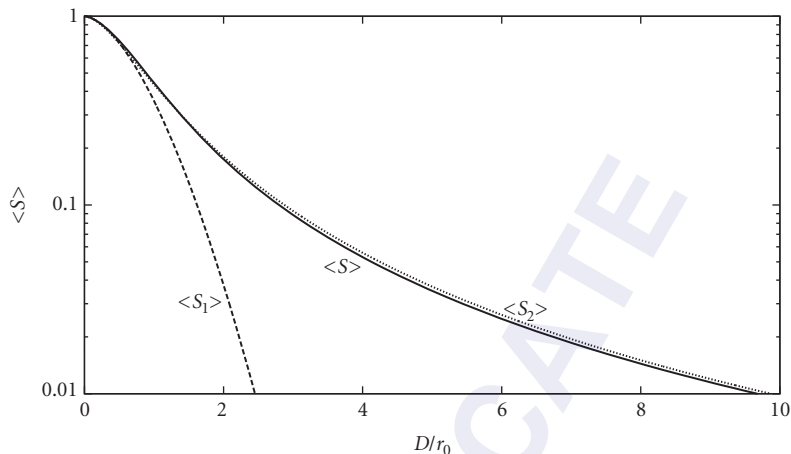


FIGURE 14 Comparison of exact and approximate long-exposure Strehl ratios for a circular pupil.

0.446. For larger values of D/r_0 , it underestimates the Strehl ratio by larger factors. It is not surprising that Eq. (79) yields an underestimate of the true value of the Strehl ratio. If we write Eq. (14) in the form

$$\mathcal{D}_\Phi(r) = \langle [\Phi(0) - \Phi(r)]^2 \rangle = 2[\sigma_\Phi^2 - R_\Phi(r)] \quad (80)$$

where

$$R_\Phi(r) = \langle \Phi(0)\Phi(r) \rangle \quad (81)$$

is the autocorrelation of the phase aberration and substitute in Eq. (24), we find that the Strehl ratio is given by Eq. (79) only when the phase correlation is zero. Otherwise, it will yield a larger value. Unfortunately, the phase correlation function for Kolmogorov turbulence is not defined; only the phase structure function is. A much better approximation is given by

$$\langle S_2(D/r_0) \rangle = [1 + (D/r_0)^{5/3}]^{-6/5} \quad (82)$$

as may be seen from Fig. 14, where it is illustrated by the dotted curve. It overestimates the true value only slightly.

4.10 MODAL CORRECTION OF ATMOSPHERIC TURBULENCE

If the first J modes of the phase aberration are corrected, the corrected phase error may be written

$$\Phi_c(\rho, \theta; \epsilon) = \sum_{j=1}^J a_j(\epsilon) Z_j(\rho, \theta; \epsilon) \quad (83)$$

Accordingly, the residual phase aberration is given by

$$\Phi_j(\rho, \theta; \epsilon) = \Phi(\rho, \theta; \epsilon) - \Phi_c(\rho, \theta; \epsilon) = \sum_{j=J+1}^{\infty} a_j(\epsilon) Z_j(\rho, \theta; \epsilon) \quad (84)$$

Hence, the mean square value of the residual aberration is given by

$$\begin{aligned} \Delta_j(\epsilon) &= \langle [\Phi_j(\rho, \theta; \epsilon)]^2 \rangle \\ &= \pi^{-1} \int_0^1 \int_0^{2\pi} \left\langle \left[\Phi(\rho, \theta; \epsilon) - \sum_{j=1}^J a_j(\epsilon) Z_j(\rho, \theta; \epsilon) \right]^2 \right\rangle \rho d\rho d\theta \\ &= \Delta_1(\epsilon) - \sum_{j=2}^J \langle a_j^2(\epsilon) \rangle \end{aligned} \quad (85)$$

where $\Delta_1(\epsilon) \equiv \sigma_{\Phi}^2(\epsilon)$. Numerical values of $\Delta_j(\epsilon)$ are given in Table 6 for several values of ϵ . Figure 13 shows how $\Delta_3(\epsilon)$ (the piston-removed and tilt-corrected variance) varies with ϵ . We note that it is nearly constant until ϵ is larger than approximately 0.5. However, as shown in Fig. 15, $\Delta_6(\epsilon)$ piston removed and aberrations corrected up to and including astigmatism) and $\Delta_{11}(\epsilon)$ (the piston removed and aberrations corrected up to and including primary spherical aberration) first increase slightly with ϵ and then decrease considerably, and finally increase rapidly as $\epsilon \rightarrow 1$.

TABLE 6 Variance $\Delta_j(\epsilon)$ of Residual Phase Aberration in Units of $(D/r_0)^{5/3}$

Δ_j	$\epsilon = 0$	0.25	0.5	0.75
Δ_1	1.0324	1.0870	1.2461	1.5010
Δ_2	0.5835	0.6104	0.6900	0.8218
Δ_3	0.1347	0.1338	0.1339	0.1426
Δ_4	0.1115	0.1149	0.1236	0.1397
Δ_5	0.0882	0.0902	0.0937	0.1008
Δ_6	0.0650	0.0654	0.0637	0.0620
Δ_7	0.0588	0.0591	0.0590	0.0603
Δ_8	0.0526	0.0529	0.0543	0.0587
Δ_9	0.0464	0.0463	0.0462	0.0474
Δ_{10}	0.0402	0.0397	0.0380	0.0361
Δ_{11}	0.0378	0.0381	0.0375	0.0360
Δ_{12}	0.0353	0.0355	0.0350	0.0349
Δ_{13}	0.0329	0.0329	0.0325	0.0338
Δ_{14}	0.0304	0.0303	0.0293	0.0290
Δ_{15}	0.0280	0.0277	0.0260	0.0242
Δ_{16}	0.0268	0.0264	0.0244	0.0218
Δ_{17}	0.0256	0.0251	0.0229	0.0193
Δ_{18}	0.0244	0.0239	0.0215	0.0186
Δ_{19}	0.0232	0.0226	0.0201	0.0178
Δ_{20}	0.0220	0.0213	0.0185	0.0153
Δ_{21}	0.0208	0.0200	0.0169	0.0129
Δ_{22}	0.0202	0.0197	0.0168	0.0129
Δ_{23}	0.0195	0.0190	0.0159	0.0115
Δ_{24}	0.0189	0.0183	0.0151	0.0101
Δ_{25}	0.0182	0.0176	0.0142	0.0095
Δ_{26}	0.0176	0.0169	0.0134	0.0089
Δ_{27}	0.0169	0.0162	0.0125	0.0075
Δ_{28}	0.0162	0.0155	0.0117	0.0061

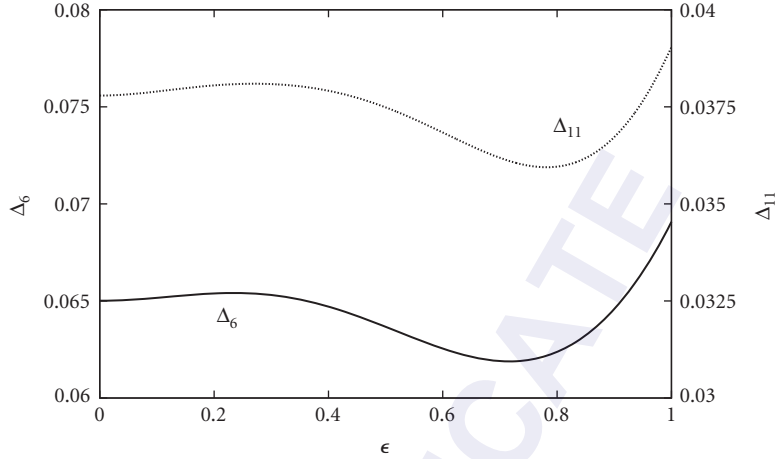


FIGURE 15 $\Delta_j(\epsilon)$ of the residual aberration in units of $(D/r_0)^{5/3}$ after correction of $J=6$ and $J=11$ modes.

The residual phase structure function after correction of the first J modes is given by

$$\mathcal{D}_j(\vec{\rho}, \vec{\rho}'; \epsilon) = \langle [\Phi_j(\vec{\rho}; \epsilon) - \Phi_j(\vec{\rho}'; \epsilon)]^2 \rangle \quad (86)$$

$$\begin{aligned} &= \langle \{ [\Phi(\vec{\rho}; \epsilon) - \Phi(\vec{\rho}'; \epsilon)] - [\Phi_c(\vec{\rho}; \epsilon) - \Phi_c(\vec{\rho}'; \epsilon)] \}^2 \rangle \\ &= \mathcal{D}_\Phi(\vec{\rho}, \vec{\rho}'; \epsilon) - \langle [\Phi_c(\vec{\rho}; \epsilon) - \Phi_c(\vec{\rho}'; \epsilon)]^2 \rangle \\ &\quad - 2\langle [\Phi_c(\vec{\rho}; \epsilon) - \Phi_c(\vec{\rho}'; \epsilon)] [\Phi_j(\vec{\rho}; \epsilon) - \Phi_j(\vec{\rho}'; \epsilon)] \rangle \end{aligned} \quad (87)$$

where, as in Eq. (14),

$$\mathcal{D}_\Phi(\vec{\rho}, \vec{\rho}'; \epsilon) = \langle [\Phi(\vec{\rho}, \epsilon) - \Phi(\vec{\rho}'; \epsilon)]^2 \rangle \quad (88)$$

is the long-exposure phase structure function. From Eq. (14), $\mathcal{D}_\Phi(\vec{\rho}, \vec{\rho}'; \epsilon)$ depends on $|\vec{\rho} - \vec{\rho}'|$, or from Eqs. (18) and (21c) on $v = |\vec{\rho} - \vec{\rho}'|/2$. For Kolmogorov turbulence in the near field, since $\mathcal{D}_\Phi = \mathcal{D}_w$, we may write from Eq. (29)

$$\mathcal{D}_\Phi(v; \epsilon) = 6.88(vD/r_0)^{5/3} \quad (89)$$

It is independent of ϵ , as expected, since it is a characteristic of turbulence.

The MCF for the residual phase aberration after correction of the first J modes is given by

$$M_j(v; \epsilon) = \exp\left\{i\langle [\Phi_j(\vec{\rho}; \epsilon) - \Phi_j(\vec{\rho}'; \epsilon)]^2 \rangle\right\} = \exp[-(1/2)\mathcal{D}_j(v; \epsilon)] \quad (90)$$

The corresponding time-averaged MTF is accordingly given by

$$\langle \tau(v; \epsilon; D/r_0) \rangle_j = \tau(v; \epsilon) \exp[-(1/2)\mathcal{D}_j(v; \epsilon)] \quad (91)$$

The MTF at a certain frequency or the Strehl ratio for a given value of D/r_0 improves as more and more modes are corrected. Of course, $J = \infty$ represents complete correction, that is the aberration-free case.

4.11 SHORT-EXPOSURE IMAGE

It is evident from Tables 4 and 5 that a large portion of the phase aberration introduced by turbulence is a random wavefront tilt caused by large eddies and represented by the coefficients $a_2(\epsilon)$ and $a_3(\epsilon)$. Figure 16 shows the relative tilt-corrected variance $[(\Delta_1(\epsilon) - \Delta_3(\epsilon))/\Delta_1(\epsilon)]$ as a function of ϵ . It increases from a value of approximately 87 percent for a circular pupil to nearly 91 percent as $\epsilon \rightarrow 1$. Thus, if the tilt is corrected in (near) real time with a steering mirror, the variance of the phase aberration for a circular pupil is reduced by a factor of $1.03/0.134 \approx 7.7$. Hence, unless D/r_0 is very large, a significant improvement in the quality of an image is obtained by correcting the wavefront tilt. A short-exposure image is equivalent to a tilt-free image, and a tilt-corrected image yields a time-averaged short-exposure image. For a circular pupil, the phase aberration with a standard deviation of 1 radian is obtained without any correction (except for the piston mode) when $D = r_0$. However, if the x and y tilts are corrected, we find from the expression for Δ_3 , D can be as large as $3.34 r_0$ for one radian of phase aberration.

Now the tilt coefficients $[a_2(\epsilon), a_3(\epsilon)]$ correlate with the primary coma coefficients $[a_8(\epsilon), a_7(\epsilon)]$, secondary coma coefficients $[a_{16}(\epsilon), a_{17}(\epsilon)]$, tertiary coma coefficients $[a_{30}(\epsilon), a_{29}(\epsilon)]$, and so on. Neglecting their correlation with coefficients other than the primary coma coefficients, the tilt-corrected phase structure function may be written⁹

$$\mathcal{D}_3(v; \epsilon) = \mathcal{D}_\Phi(v; \epsilon) - \frac{16v^2}{1 + \epsilon^2} \langle a_2^2(\epsilon) \rangle - \frac{32\sqrt{2}v^2[(1 + \epsilon^2)(6v^2 - 6v + 1) - 2\epsilon^4]}{(1 - \epsilon^2)(1 + 6\epsilon^2 + 10\epsilon^4 + 6\epsilon^6 + \epsilon^8)^{1/2}} \langle a_2(\epsilon)a_8(\epsilon) \rangle \quad (92)$$

It is evident that unlike the uncorrected phase structure function, the tilt-corrected structure function does depend on the diameter D . For a circular pupil, Eq. (92) reduces to

$$\begin{aligned} \mathcal{D}_3(v) &= \mathcal{D}_\Phi(v) - 16v^2 \langle a_2^2 \rangle - 32\sqrt{2}v^2(6v^2 - 6v + 1) \langle a_2 a_8 \rangle \\ &= 6.8839v^{5/3}(1 - 0.9503v^{1/3} - 0.5585v^{4/3} + 0.5585v^{7/3})(D/r_0)^{5/3} \end{aligned} \quad (93)$$

$\mathcal{D}_3(v)$ has a maximum value of $0.344(D/r_0)^{5/3}$ at $v = 1$, which is $1/20$ of the value $6.88(D/r_0)^{5/3}$ without any correction (i.e., for $J = 1$), illustrating a significant benefit of the tilt correction. When $\epsilon \neq 0$, $\mathcal{D}_3(v; \epsilon)$ is larger for larger values of ϵ , indicating more severe fluctuations of the wavefront.

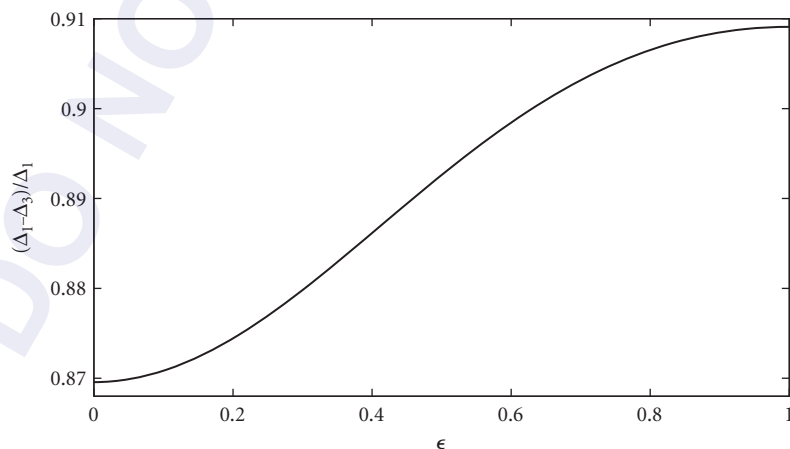


FIGURE 16 Relative tilt-corrected phase variance $[(\Delta_1(\epsilon) - \Delta_3(\epsilon))/\Delta_1(\epsilon)]$ as a function of the obscuration ratio ϵ .

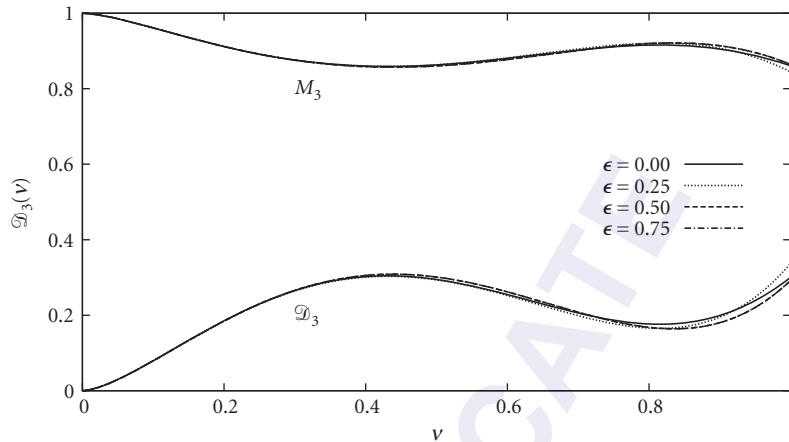


FIGURE 17 Tilt-corrected structure function $\mathcal{D}_3(v; \epsilon)$ and the corresponding mutual coherence function $M_3(v; \epsilon)$.

Using Eq. (92), the short-exposure MTF can be calculated as

$$\langle \tau(v; \epsilon; D/r_0) \rangle_3 = \tau(v; \epsilon) M_3(v; \epsilon; D/r_0) \quad (94)$$

where

$$M_3(v; \epsilon; D/r_0) = \exp[-(1/2)\mathcal{D}_3(v; \epsilon; D/r_0)] \quad (95)$$

is the tilt-corrected or short-exposure MCF or the MTF reduction factor associated with atmospheric turbulence. This factor does depend on the pupil diameter D , in contrast to the long-exposure MTF reduction factor given by Eq. (36), which is independent of D . As expected, it is smaller than the long-exposure reduction factor. Figure 17 shows how the tilt-corrected structure function and the corresponding MCF vary with v . The corresponding long-exposure functions are included for comparison.

The time-averaged Strehl ratio of the tilt-corrected image is given by

$$\langle S(\epsilon; D/r_0) \rangle_3 = \frac{8}{(1 - \epsilon^2)} \int_0^1 \langle \tau(v; \epsilon; D/r_0) \rangle_3 v dv \quad (96)$$

Replacing the long-exposure atmospheric MTF reduction factor in Eq. (47a) by the short-exposure MTF reduction factor, the short-exposure central irradiance η_{SE} may be written

$$\eta_3(\epsilon; D/r_0) = (1 - \epsilon^2) \left(\frac{D}{r_0} \right)^2 \langle S(\epsilon; D/r_0) \rangle_3 \quad (97)$$

The Strehl ratio $\langle S \rangle_3$ and the central irradiance η_3 are shown in Figs. 18 and 19, respectively. Some typical values of the tilt-corrected Strehl ratios are given in Table 1 for various values of D/r_0 and ϵ . For small values of D/r_0 , η_3 increases approximately as for a diffraction-limited system, since the aberration is small; reaches a maximum (e.g., 3.74 for $D/r_0 = 3.5$ for a circular pupil), and then decreases slowly but monotonically to unity. Since the image displacement due to wavefront tilt decreases as $D^{-1/3}$, the effect of tilt correction becomes negligible for large values of D/r_0 due to the large residual phase errors. As in the case of a circular pupil, if the covariance of the tilt coefficients with others is ignored, for example $\langle a_2(\epsilon) a_8(\epsilon) \rangle$ in Eq. (92) is neglected, an unrealistic overcorrection will result.^{11,26,27}

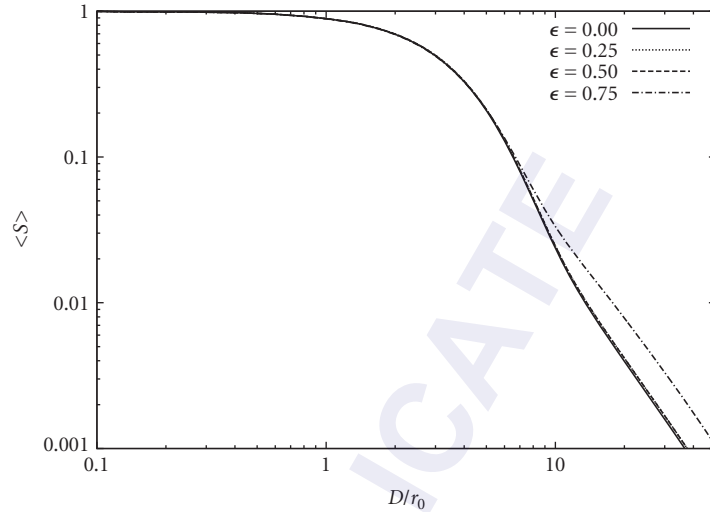


FIGURE 18 Variation of time-averaged short-exposure Strehl ratio $\langle S(\epsilon; D/r_0) \rangle_3$ with D/r_0 .

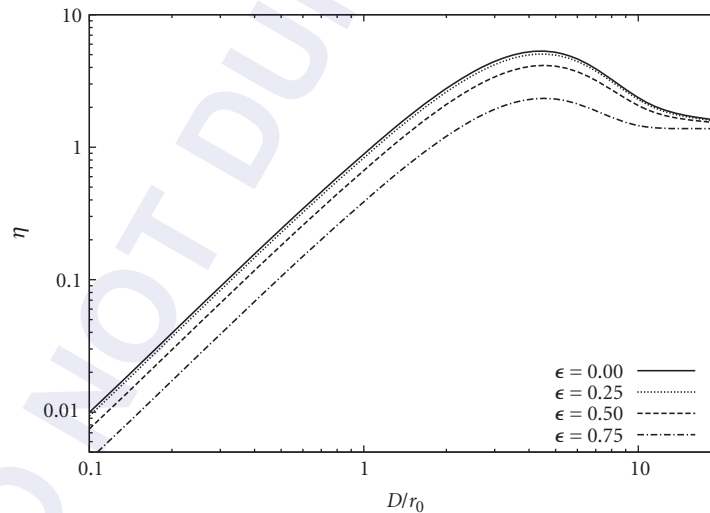


FIGURE 19 Variation of time-averaged short-exposure central irradiance $\langle \eta(\epsilon; D/r_0) \rangle_3$ with D/r_0 .

For a circular pupil, Fig. 20 shows how the tilt-corrected time-averaged Strehl ratio varies with D/r_0 . The uncorrected Strehl ratio is also shown to illustrate the improvement made by the tilt correction. An approximate value of the Strehl ratio given by $\exp(-\Delta_3)$, where

$$\Delta_3 = 0.134(D/r_0)^{5/3} \quad (98)$$

is the tilt-corrected time-averaged phase aberration variance, and shown by the dashed curve underestimates the true value $\langle S \rangle_3$.

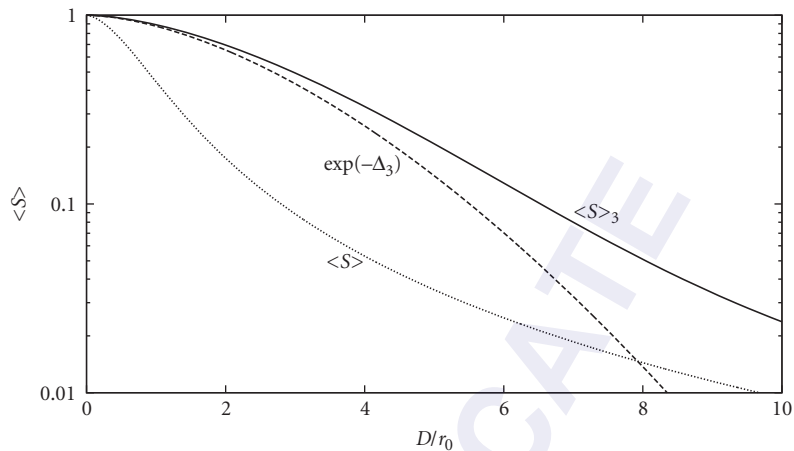


FIGURE 20 Short-exposure Strehl ratio $\langle S \rangle_3$ compared with its approximate value $\exp(-\Delta_3)$. The long-exposure Strehl ratio $\langle S \rangle$ is also shown to illustrate the improvement made by the tilt correction.

An example of an instantaneous short-exposure PSF illustrating the structure of the image of a star as seen by a ground-based telescope with a circular pupil and $D/r_0 = 10$ is shown in Fig. 21, where it is compared with the corresponding aberration-free image. On the average, the standard deviation $\sqrt{\Delta_3}$ of the instantaneous aberration is 2.5 radians or 0.4λ . We note that the image is broken up into small spots called *speckles*, which is a characteristic of large spatially random aberrations. The size of a speckle is determined by D , its angular radius being approximately equal to λ/D . The image size lies between its diffraction-limited value λ/D and its long-exposure value λ/r_0 (varying as $\lambda^{-0.2}$). The image becomes progressively worse as r_0 decreases, showing the effects of what astronomers call *seeing*. However, the size of a speckle decreases as D increases without affecting the long-exposure image size λ/r_0 . Thus, an increase in D does not significantly improve the resolution of the system (as determined by the overall size of the image). The angular spot radius in units of λ/D representing the point at which the irradiance drops to half of its central value is given in Table 7 for both the short- and long-exposure images for various values of D/r_0 and ϵ , and compared with the corresponding

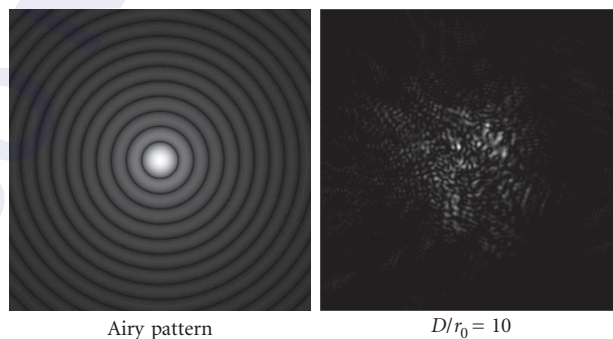


FIGURE 21 Instantaneous short-exposure PSF compared with the corresponding aberration-free PSF.

TABLE 7 Angular Spot Radius of Short-Exposure (SE), Long-Exposure (LE), and Diffraction-Limited (DL) Images in Units of λ/D for Various Values of ϵ and D/r_0

ϵ	$D/r_0=1$			$D/r_0=2$		$D/r_0=5$		$D/r_0=10$	
	DL	LE	SE	LE	SE	LE	SE	LE	SE
0	0.515	0.708	0.516	1.122	0.518	2.552	0.560	4.983	2.650
0.25	0.497	0.698	0.498	1.151	0.501	2.595	0.539	5.021	2.686
0.5	0.456	0.688	0.457	1.280	0.459	2.711	0.498	5.106	2.840
0.75	0.406	0.700	0.407	1.480	0.411	3.130	0.465	5.432	3.469

diffraction-limited value. We note that, whereas the LE radius is roughly D/r_0 times the DL radius, the SE radius for small values of D/r_0 , for example $D/r_0=5$, is roughly equal to the DL radius, but for large values of D/r_0 , for example $D/r_0=10$, it is about half the LE radius. Of course, as ϵ increases and more and more light goes into the diffraction rings, these factors change somewhat. For small values of D/r_0 , the diffraction ring structure of the DL image, smoothed out in the LE image, returns in the SE image. This is more and more evident for larger and larger values of ϵ .

If the image wander is corrected in real time and the image is averaged over time, the speckles disappear and the image becomes smooth. Such PSFs and the corresponding encircled power can be obtained from Eqs. (43) and (44), respectively, by substituting the short-exposure OTF given by Eq. (94). The MTF, PSF, and encircled power thus obtained are shown by the dotted curves in Figs. 5 and 8 for several values of ϵ . The short-exposure image when randomly displaced over time yields a smooth long-exposure image, as illustrated by the dashed curves in Fig. 8. Modal correction of a Gaussian beam (including one with a central obscuration) propagating through turbulence has also been discussed by Wang using the Zernike circle polynomials²⁸ though such polynomials are not orthogonal over the Gaussian weighted (or an annular) pupil.^{5,29}

4.12 ADAPTIVE OPTICS

The correction of wavefront errors in (near) real time is accomplished by using adaptive optics.¹⁻³ In practice, a steering mirror with only three actuators is used to correct the large x and y wavefront tilts (also called tip and tilt). The residual aberration is corrected by a deformable mirror, which is deformed by an array of actuators attached to it. The signals for the actuators are determined either by sensing the wavefront errors with a wavefront sensor in a closed loop to minimize the variance of the residual errors, or the actuators are actuated to produce Zernike modes (e.g., focus, two modes of astigmatism, two modes of coma) iteratively until sharpness of the image is maximized.³⁰⁻³² The signals are independent of the optical wavelength, provided atmospheric dispersion is negligible. The two approaches are referred to as *zonal* and *modal approaches*, respectively. The zonal approach has the advantage that the rate of correction is limited only by the rate at which the wavefront errors can be sensed and the actuators can be actuated. However, the amount of light that is used by the wavefront sensor is lost from the image. In practice, the image beam is split into two parts. The centroid of the image of one part is measured with a quad cell, and the tilt indicated by it is corrected with a steering mirror. The resulting tilt-corrected image of the other part with the residual aberration is corrected with a deformable mirror in a closed-loop manner.³³ In the modal approach, there is no loss of light, but the rate or the bandwidth of correction³⁴⁻³⁶ can be slow due to its iterative nature, especially when turbulence is severe and a large number of modes must be corrected. Moreover, for imaging an extended object, wavefront sensing requires a point source in its vicinity, but the modal approach is applicable to the extended object itself.

Of course, adaptive optics can improve the image quality only if the object lies within an isoplanatic angle of turbulence. In the case of a ground-to-space laser illuminating a satellite, the angular travel

(point-ahead angle) of the satellite during the round-trip time of the beam to the satellite must be less than the isoplanatic angle of turbulence. An estimate of the isoplanatic angle can be obtained from³⁷

$$\vartheta_0 = 0.058 \lambda^{6/5} (\cos \vartheta)^{8/5} \left\{ \int_0^L C_n^2(z) z^{5/3} dz \right\}^{-3/5} \quad (99)$$

Substituting Eq. (31) into Eq. (99), we obtain $\vartheta_0 = 7 \mu\text{rad}$ for $\lambda = 0.5 \mu\text{m}$. Thus, the $H - V_{5/7}$ model for C_n^2 yields $r_0 = 5 \text{ cm}$ and $\vartheta_0 = 7 \mu\text{rad}$ for visible light; hence the 5/7 subscript.

4.13 SUMMARY

We have considered imaging through atmospheric turbulence by a system with an annular pupil, such as the Keck telescope. The results for a system with a circular pupil are obtained as a special case of those for the annular pupil. An atmospheric coherence length r_0 is defined such that the coherence of two points separated by this distance is 0.03. It is calculated using the well-known $H - V_{5/7}$ model for the refractive index structure parameter C_n^2 , assuming Kolmogorov turbulence, for looking up and down at a point source through the atmosphere as well as for plane wave propagation (see Fig. 2). It is shown that turbulence limits the resolution of a system to one with pupil diameter r_0 regardless of the actual pupil diameter. The two-point angular resolution on ground when observed from any point in space varies approximately as λ/r_0 , compared to the diffraction-limited resolution of λ/D . There is reciprocity in wave propagation through the turbulent atmosphere; otherwise adaptive optics couldn't have worked in ground-based astronomy. Accordingly, the PSF observed on an aircraft looking at a ground point object is the same as the irradiance distribution of a beam focused on ground from an aircraft. Similarly, the irradiance distribution of a beam focused in space from ground is the same as the PSF observed on ground looking at a space object. The effect of turbulence is shown in terms of the MTF, Strehl ratio, PSE, and encircled power for both long- and short-exposure images. How the image improves as more and more Zernike aberration modes are corrected is also discussed. The aberration function may also be expanded in terms of the Karhunen-Loève functions whose coefficients are statistically independent of each other.^{27,38} It is found though that the effect of correlation of Zernike coefficients is negligible for $D/r_0 \leq 4$.²⁶

The dependence of the various image-related quantities on the obscuration ratio ϵ may be summarized as follows:

1. Atmospheric turbulence dominates the degradation of MTF; the MTF gain at high frequencies due to the pupil obscuration is lost even for weak turbulence (compare Figs. 4 and 5). The Strehl ratio is similarly dominated by turbulence (see Fig. 6).
2. The piston-removed aberration variance in units of $(D/r_0)^{5/3}$ increases monotonically from a value of 1.0324 for a circular pupil ($\epsilon = 0$) to a value of 1.843 for an infinitesimally thin ring pupil ($\epsilon \rightarrow 1$) (see Fig. 13).
3. The fluctuations in the angle of arrival depend weakly on the pupil obscuration (see Figs. 12 and 13).
4. The time-averaged variance of some aberration coefficients increases with ϵ and decreases for others (Fig. 9). For example the variance of tilt and astigmatism increases, but that of defocus and spherical aberration decreases and approaches zero as $\epsilon \rightarrow 1$. The covariance follows a similar trend (Fig. 10).
5. When a certain number of Zernike aberration modes is corrected, for example with a deformable mirror, the image quality improves. The residual aberration variance when the first J modes are corrected is given in Table 6. If only tip and tilt modes are corrected, say with a steering mirror, the aberration variance is reduced by 87 percent for a circular pupil and slightly larger (up to 90 percent) for an annular pupil. Figure 19 illustrates the improvement in the Strehl ratio. The aberration-free Airy pattern representing the image of a star breaks up into speckles, as illustrated in Fig. 21. The time-averaged short-exposure OTFs are compared in Fig. 5 with the corresponding long-exposure results. The PSFs and encircled-power curves are similarly compared in Fig. 8.

4.14 ACKNOWLEDGMENTS

V. N. Mahajan gratefully acknowledges helpful discussions with H. T. Yura, who also read the manuscript and made useful comments.

4.15 REFERENCES

1. R. Fugate, "Adaptive Optics," *Handbook of Optics*, 3rd ed, M. Bass (ed.), Chap. 5, Vol. V, McGraw-Hill, New York, 2009.
2. J. W. Hardy, *Adaptive Optics for Astronomical Telescopes*, Oxford, New York, 1998.
3. R. K. Tyson, *Introduction to Adaptive Optics*, SPIE Press, Bellingham, Washington, 1999.
4. D. Fried, "Optical Resolution through a Randomly Inhomogeneous Medium for Very Long and Very Short Exposures," *J. Opt. Soc. Am.* **56**:1372–1379 (1966).
5. V. N. Mahajan, "Zernike Annular Polynomials for Imaging Systems with Annular Pupils," *J. Opt. Soc. Am.* **71**: 75–85 (1981).
6. V. N. Mahajan, "Zernike Annular Polynomials for Imaging Systems with Annular Pupils," *J. Opt. Soc. Am.* **71**: 1408 (1981).
7. V. N. Mahajan, "Zernike Annular Polynomials and Optical Aberrations of Systems with Annular Pupils," *Appl. Opt.* **33**:8125–8127 (1994).
8. V. N. Mahajan, "Orthogonal Polynomials in Wavefront Analysis," *Handbook of Optics*, 3rd ed, M. Bass (ed.), Chap. 11. Vol. II, McGraw-Hill, New York, 2009.
9. G.-M Dai and V. N. Mahajan, "Zernike Annular Polynomials and Atmospheric Turbulence," *J. Opt. Soc. Am.* **A24**:139–155 (2007).
10. M. Born and E. Wolf, *Principles of Optics*, 7th ed., Oxford, New York, 1999.
11. V. N. Mahajan, *Optical Imaging and Aberrations, Part II: Wave Diffraction Optics*, SPIE Press, Bellingham, Washington (Second Printing 2004).
12. D. Fried, "Evaluation of r_0 for Propagation Down through the Atmosphere," *Appl. Opt.* **13**:2620–2622 (1974); errata 1, *Appl. Opt.* **14**:2567 (1975); errata 2, *Appl. Opt.* **16**:549 (1977).
13. R. R. Beland, "Propagation through Atmospheric Turbulence," *Atmospheric Propagation of Radiation, The Infrared & Electro-Optical Systems Handbook*, F. G. Smith (ed.), Springer, 1993. There is a typographical error in the power of 10 in the value of A given in this reference.
14. D. Fried, "Limiting Resolution Looking Down through the Atmosphere," *J. Opt. Soc. Am.* **A 56**:1380–1384 (1966).
15. D. L. Walters and L. W. Bradford, "Measurement of r_0 and θ_0 : Two Decades and 18 Sites," *Appl. Opt.* **36**: 7876–7886 (1997).
16. E. L. O'Neill, "Transfer Function for an Annular Aperture," *J. Opt. Soc. Am.* **46**:285–288 (1956). Note that a term of $-2\eta^2$ is missing in the second of O'Neill's Eq. (26).
17. J. W. Goodman, *Introduction to Fourier Optics*, 2nd ed., McGraw-Hill, New York, 1996.
18. V. N. Mahajan and B. K. C. Lum, "Imaging through Atmospheric Turbulence with Annular Pupils," *Appl. Opt.* **20**:3233–3237 (1981).
19. R. J. Noll, "Zernike Polynomials and Atmospheric Turbulence," *J. Opt. Soc. Am.* **66**:207–211 (1976).
20. R. J. Sasiela, *Electromagnetic Wave Propagation in Turbulence*, Springer-Verlag, New York, 1994.
21. V. N. Mahajan, *Optical Imaging and Aberrations, Part I: Ray Geometrical Optics*, SPIE Press, Bellingham, Washington (Second Printing 2001).
22. V. I. Tatarski, *The Effects of the Turbulent Atmosphere on Wave Propagation*, U. S. Department of Commerce, 1971.
23. V. N. Mahajan, "Line of Sight of an Aberrated Optical System," *J. Opt. Soc. Am.* **A2**:833–846 (1985).
24. W. Wetherell, "The Calculation of Image Quality," *Applied Optics and Applied Engineering*, R. R. Shannon and J. C. Wyant (eds.), Academic Press, 1980.
25. V. N. Mahajan, "Strehl Ratio for Primary Aberrations in Terms of Their Aberration Variance," *J. Opt. Soc. Am.* **73**:240–241 (1983).

26. J. Y. Wang, "Optical Resolution through a Turbulent Medium with Adaptive Phase Compensation," *J. Opt. Soc. Am.* **67**:383–390 (1977).
27. G.-M. Dai, "Modal Compensation of Atmospheric Turbulence with the Use of Zernike Polynomials and Karhunen Loève Functions," *J. Opt. Soc. Am.* **A12**:2182–2193 (1995).
28. J. Y. Wang, "Phase-Compensated Optical Beam Propagation through Atmospheric Turbulence," *Appl. Opt.* **17**:2580–2590 (1978).
29. V. N. Mahajan, "Zernike-Gauss Polynomials and Optical Aberrations of Systems with Gaussian Pupils," *Appl. Opt.* **34**:8057–8059 (1995).
30. R. A. Miller and A. Buffington, "Real-Time Wavefront Correction of Atmospherically Degraded Telescopic Images through Image Sharpening," *J. Opt. Soc. Am.* **61**:1200–1210 (1974).
31. A. Buffington, F. S. Crawford, R. A. Miller, A. J. Schwemin, and R. G. Smits, "Correction of Atmospheric Distortion with an Image-Sharpening Telescope," *J. Opt. Soc. Am.* **67**:298–305 (1977).
32. V. N. Mahajan, J. Govignon, and R. J. Morgan, "Adaptive Optics without Wavefront Sensors," *SPIE Proc.* **228**:63–69 (1980).
33. L. C. Roberts, Jr. and C. R. Neyman, "Characterization of the AEOS Adaptive Optics System," *Pub. Astro. Soc. Pacific* **114**:1260–1266 (2002).
34. J. Y. Wang, "Effect of Finite Bandwidth on Far-Field Performance of Modal Wavefront-Compensative Systems," *J. Opt. Soc. Am.* **69**:819–828 (1977).
35. C. B. Hogge and R. R. Butts, "Frequency Spectra for the Geometric Representation of Wavefront Distortions Due to Atmospheric Turbulence," *IEEE Trans. Antennas Propagation* **AP-24**:144–154 (1976).
36. D. P. Greenwood, "Bandwidth Specification of Adaptive Optics Systems," *J. Opt. Soc. Am.* **67**:390–393 (1977).
37. D. L. Fried, "Anisoplanatism in Adaptive Optics," *J. Opt. Soc. Am.* **72**:52–61 (1982).
38. J. Y. Wang and J. K. Markey, "Modal Compensation of Atmospheric Turbulence Phase Distortion," *J. Opt. Soc. Am.* **68**:78–87 (1978).

ADAPTIVE OPTICS

Robert Q. Fugate

*Starfire Optical Range
Directed Energy Directorate
Air Force Research Laboratory
Kirtland Air Force Base, New Mexico*

5.1 GLOSSARY

C_n^2	refractive index structure parameter
d_0	section size for laser beacon adaptive optics
\mathcal{D}_n	refractive index structure function
\mathcal{D}_ϕ	phase structure function
E_f^2	mean square phase error due to fitting error
E_{FA}^2	mean square phase error due to focus anisoplanatism
E_n^2	mean square phase error due to sensor read noise
E_s^2	mean square phase error due to servo lag
f_G	Greenwood frequency
f_{T_G}	Tyler G -tilt tracking frequency
$F_\phi(f)$	phase power spectrum
$F_{\phi_G}(f)$	G -tilt power spectrum
$H(f, f_c)$	servo response function
N_{pde}	number of photo-detected electrons per wavefront subaperture
r_0	Fried's coherence length
SR_{HO}	Strehl ratio resulting from higher-order phase errors
SR_{tilt}	Strehl ratio due to full-aperture tilt
$Z_j(\rho, \theta)$	Zernike polynomials
σ_θ	rms full-aperture tracking error
σ_{θ_G}	rms full-aperture G -tilt induced by the atmosphere
θ_0	isoplanatic angle

5.2 INTRODUCTION

An Enabling Technology

Adaptive optics (AO) is *the* enabling technology for many applications requiring the real-time control of light propagating through an aberrating medium. Today, significant advances in AO system performance for scientific, commercial, and military applications are being made because of improvements in component technology, control algorithms, and signal processing. For example, the highest-resolution images of living human retinas ever made have been obtained using AO.¹⁻³ Increased resolution of retinal features could provide ophthalmologists with a powerful tool for the early detection and treatment of eye diseases. Adaptive optics has brought new capabilities to the laser materials processing industry, including precision machining of microscopic holes and parts.^{4,5} Beaming power to space to raise satellites from low earth orbit to geosynchronous orbit, to maneuver satellites that are already in orbit, and to provide electricity by illuminating solar panels remains a topic of commercial interest.^{6,7} Recent work has shown that high-energy pulsed lasers, corrected for atmospheric distortions with AO, could be used to alter the orbits of space debris objects,⁸ causing them to reenter the atmosphere much sooner than natural atmospheric drag—creating an opportunity for the environmental restoration of space. High-speed optical communication between the ground, aircraft, and spacecraft (including very deep space probes) is a topic of continuing interest that would benefit from the use of AO.⁹ Adaptive optics offers significant benefit to a large variety of military applications, including the U.S. Air Force Airborne Laser program, in which a high-powered chemical oxygen-iodine laser will be mounted in an aircraft to engage boosting missiles at long ranges.¹⁰ Adaptive optics will be used to predistort the laser beam as it leaves the aircraft to compensate for the atmospheric distortions that the beam encounters as it propagates to the missile target.

Perhaps the most widely known and discussed application for AO is imaging through the atmosphere. Military sites are now using AO on medium-sized ground-based telescopes [of the order of 4 meters (m)] to inspect low-earth-orbiting satellites. However, the most dramatic use of AO is in astronomy. Astronomers have been plagued by the atmosphere for centuries—since the invention of the telescope. Indeed, distortion caused by turbulence is the principal motivation for launching large telescopes into space. However, if ground-based telescopes were able to achieve diffraction-limited resolution and high Strehl ratios, a significant obstacle to new discoveries and more productive research would be overcome. Adaptive optics, originally proposed by the astronomer Horace Babcock in 1953,¹¹ may enable that to happen, creating a revolution in ground-based optical astronomy. Adaptive optics is especially important in light of the large new telescope mirrors. Incredible advances in mirror technology have made 8-m-diameter monolithic mirrors almost commonplace (four have seen first light and five more are under construction). These are in addition to the two 10-m Keck telescopes and the 11-m Hobby-Eberly telescope, all of which use segmented primaries. The atmosphere limits the resolution of ground-based telescopes to an equivalent diameter that is equal to Fried's coherence length, r_0 (a few tens of centimeters at visible wavelengths at the best sites).¹² Adaptive optics offers the potential for telescopes to achieve diffraction-limited imaging and high throughput to spectroscopic instruments. Furthermore, interferometric arrays of large telescopes will enable unprecedented imaging resolution of faint objects, given that the unit telescopes can each produce nearly distortion-free wave fronts. Because imaging through the atmosphere is a topic of great interest, and because it embodies all aspects of AO technology, this discussion is oriented toward that application.

The chapter is organized as follows: Sec. 5.3 describes the basic concept of AO as applied to ground-based telescopes. Section 5.4 is a short summary of the classical description of atmospheric turbulence and parameters that are important for designing and evaluation of AO systems. Section 5.5 is a description of the hardware and the software that are needed for practical implementation of AO. This includes tracking, wavefront sensing, processors, and wavefront correctors. Section 5.6 is a discussion of design issues for top-level system performance trades using scaling laws and formulas.

5.3 THE ADAPTIVE OPTICS CONCEPT

This section provides a brief, qualitative overview of how AO imaging systems work. Figure 1 is a highly simplified diagram that illustrates the principles. In this figure, a conventional telescope is shown on the left and one that is equipped with AO is shown on the right. The science object

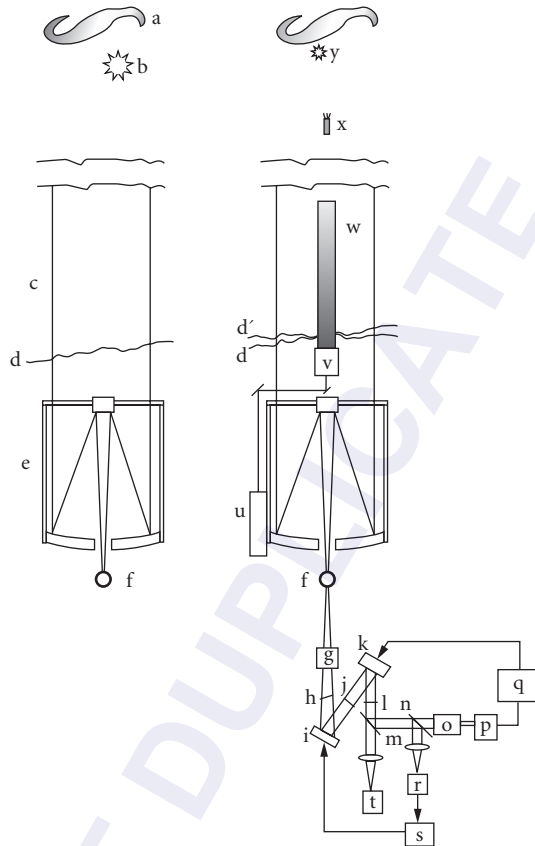


FIGURE 1 Conventional telescope with no AO (left) and additional components needed for AO (right). (a) Object of interest, (b) natural guide star, (c) atmospheric turbulence, (d) aberrated wavefront (including tilt and higher-order distortions) after passing through the turbulent atmosphere, (d') aberrated wavefront from the laser beacon with full-aperture tilt removed, (e) telescope, (f) aberrated image, (g) relay optics, (h) demagnified aberrated wavefront including full-aperture tilt, (i) fast-steering mirror, (j) tilt-removed wavefront with only higher-order aberrations remaining, (k) deformable mirror, (l) corrected wavefront, (m) aperture-sharing element, (n) tilt sensor pickoff beam splitter, (o) relay optics, (p) higher-order wavefront sensor, (q) electronic processor to compute deformable mirror commands, (r) full-aperture tilt sensor, (s) tilt mirror processor and controller, (t) science camera, (u) laser for generating laser beacons, (v) launch telescope for the beacon laser, (w) Rayleigh laser beacon, (x) mesospheric sodium laser beacon, and (y) faint natural guide star for full-aperture tilt sensing.

(a) and a natural guide star (b) are at the top of the figure. The turbulent atmosphere (c) creates higher-order phase distortion and an overall tilt on wavefronts reaching the telescope (d). The telescope (e) forms an aberrated image (f) at the location of the camera or spectrograph. The natural guide star is used, either manually or automatically, to correct pointing errors in the telescope mount and the overall wavefront tilt that are induced by atmospheric turbulence.

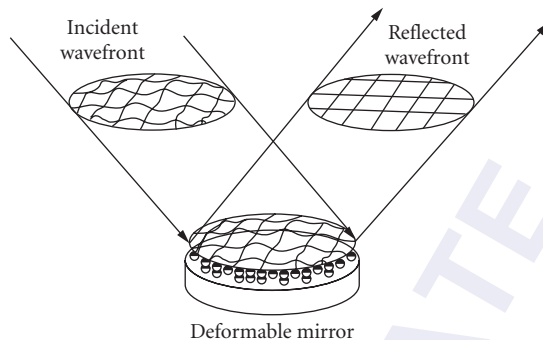


FIGURE 2 Deformable mirror concept in which phase aberrations on an incident wavefront are removed by setting the deformable mirror to the conjugate of the aberration.

The objectives of the AO system are to continuously remove the higher-order distortion and to stabilize the position of the image by removing the overall tilt. The components are shown on the right side of Fig. 1. One can think of the AO system as a dynamic optical system that simultaneously relays the original image to a new focal plane [the camera (*t*)] while removing both the higher-order aberrations (those distortions having spatial frequencies higher than tilt) with the deformable mirror (*k*) and the tilt aberrations with the fast-steering mirror (*i*), leaving only a minor amount of residual error in the wave front (*l*). The appropriate optical relay telescopes (*g* and *o*) are used as required by the particular implementation.

Figure 2 illustrates how a deformable mirror removes phase aberrations that are induced by atmospheric turbulence. The conjugate of the measured wavefront distortion is imposed on the surface of the deformable mirror so that on reflection, the distortions are removed. The AO system is never able to perfectly match the distortions in the wavefront because of a number of error sources that will be discussed in later sections.

To set the figure on the deformable mirror, the AO system must get information about the turbulent atmosphere by measuring its effect on the wavefronts from a beacon—a source of light such as a bright star (see Fig. 1, item *b*) located at or near the science object. The object itself may serve as the beacon if it is bright enough (either by self-emission, as in the case of a star or astronomical object, or by reflection of natural or man-made light, as in the case of artificial satellites). When a natural star does not exist or when the object itself is not bright enough, it may be possible to use an artificial beacon generated by a low-power laser (see Fig. 1, item *u*). A laser beacon can be created either by Rayleigh scattering up to a range of about 20 km, or by resonant scattering of atomic species such as sodium in the mesosphere, at an altitude of 90 km above Earth's surface. In the implementation that is shown in Fig. 1, a laser beam that is capable of exciting the D_2 line in atomic sodium in the mesosphere is projected from behind the secondary mirror of the telescope.

Control signals for the wavefront corrector mirrors are generated by a full-aperture tracking sensor (see Fig. 1, item *r*) and a wavefront sensor (*p*) by observing the residual error in the beacon wavefront (*l*). An aperture-sharing element (*m*) directs light from the beacon into the wavefront sensor, which is located at an image of the entrance pupil of the telescope. The wavefront sensor samples the slope of the wavefront over subaperture regions of the order of r_0 in size—for which the wavefront is essentially an undistorted, but tilted, plane wave—measuring the tilt and reporting that tilt as a wavefront phase gradient. The computer (*q*) combines the subaperture gradient measurements and “reconstructs” a best estimate of the residual phase error at specific points in the aperture, generally the locations of the actuators of the deformable mirror. Error signals derived from the reconstructed wavefront are sent to the deformable mirror to further reduce the residual error.

The tracking sensor measures the overall tilt of the beacon wavefront by computing the centroid of a focused image of the beacon that is formed by the full aperture of the telescope (or by other, more sophisticated means). The tracker processor (*s*) generates an error signal that controls the fast-steering mirror to keep the image of the beacon centered on the tracking sensor. It is also possible to derive full-aperture tilt information from the higher-order wavefront sensor, but optimum tracking performance usually requires a specialized sensor and processor. When a laser beacon is used to obtain higher-order wavefront information, it is still necessary to use a natural guide star to derive full-aperture tilt information: The laser beacon's position in the sky is not known with respect to an inertial reference (like a star), because its path on the upward propagation is random. To first order the beacon does not appear to move at all in the tracker due to near perfect reciprocity on the upward and downward paths, generating a return wavefront having no full aperture tilt as shown at (*d'*) in Fig. 1. The requirement to use a natural guide star for tracking is not as serious as it first sounds, however, because it is much more likely that a natural guide star can be found near the science object that satisfies tracking requirements but is still too faint for higher-order wavefront measurements. This is possible because the full aperture of the telescope and a small number of specialized detectors can be used for tracking so that stars fainter by a factor of $\sim (D/r_0)^2$ that are available to the wavefront sensor become available to the tracking sensor.

Figure 1 is highly simplified. The details of the optics have been left out to emphasize the main principles. For instance, it is customary to form an image of the entrance pupil of the telescope onto the fast-steering mirror, the deformable mirror, and again on the wavefront sensor. At least two powered optical elements are required to generate each of these images. Even more sophisticated approaches have been proposed in which multiple deformable mirrors are reimaged to layers in the turbulence. The nature of the aperture-sharing element will vary significantly depending on the adopted philosophy for the system. For example, if obtaining enough light for the wavefront sensor is a problem (as is usually the case in astronomy), it is prudent to send the entire spectrum over which the detectors are responsive to the wavefront sensor and divert only the absolute minimum needed part of the spectrum to the imaging sensor. This approach also means that the optics must be designed to work over a wide spectral range (e.g., from 450 to 1000 nm for a visible wavefront sensor) and also requires careful design of atmospheric dispersion correction. Furthermore, there are innovative approaches for implementing the wavefront corrector mirrors. The Steward Observatory group is developing a tilting, deformable secondary mirror for the new 6.5-m primary replacement at the Multiple Mirror Telescope (MMT) and 8.4-m Large Binocular Telescope (LBT) observatories.¹³ This arrangement requires only two reflecting surfaces between the sky and the camera.

5.4 THE NATURE OF TURBULENCE AND ADAPTIVE OPTICS REQUIREMENTS

Turbulence Generation and the Kolmogorov Model

This section summarizes the classical description of atmospheric turbulence and its effect on the propagation of light. The literature on this topic is enormous. Our purpose here is to introduce the principles and define a few important parameters that are relevant to the operation of adaptive optical systems. The references cited throughout provide more detail.

Atmospheric turbulence is generated by solar heating of the Earth's surface. Air at the surface of the Earth is warmed in the day and cooled at night. Temperature gradients develop, creating a convective flow of large masses of air. Turbulence develops as large-scale masses of air break up into smaller spatial scales and dissipate their energy to the surrounding air. The effects of turbulence on electromagnetic wave propagation is governed by the nature of the spatial and temporal fluctuations of the index of refraction of air. The refractive index of air at optical wavelengths obeys the formula¹⁴

$$n = 1 + 77.6 \left[1 + \frac{7.52 \cdot 10^{-3}}{\lambda^2} \right] \frac{P}{T} \cdot 10^{-6} \quad (1)$$

where λ is the wavelength of light in micrometers, P is the atmospheric pressure in millibars, and T is the atmospheric temperature in kelvins (K). The effect of atmospheric pressure changes on n are small and can, for problems of interest in this chapter, be neglected. The dominant influence on variations of n is the air temperature. At visible wavelengths, $dn/dT \sim 10^{-6}$ at standard pressure and temperature. This means that two 1-m-long columns of air having a temperature difference of 1 K create roughly one wave of optical path difference for visible light, a very significant effect. We can model the index of refraction as the sum of two parts,

$$n(\mathbf{r}, t, \lambda) = n_0(\mathbf{r}, \lambda) + n_1(\mathbf{r}, t) \quad (2)$$

where n_0 represents the deterministic, slowly changing contribution (such as variation with height above the ground), and n_1 represents the random fluctuations arising from turbulence. The wavelength dependence of n_1 is ignored. Furthermore, typical values of n_1 are several orders of magnitude smaller than unity.

During the 1940s, Andrey Nikolayevich Kolmogorov (1903–1987) developed theories¹⁵ describing how energy is dissipated in the atmosphere and modeled the spatial power spectrum of turbulent velocity fluctuations. V. I. Tatarskii postulated the applicability of the Kolmogorov spatial power spectrum to refractive index fluctuations and then solved the wave equation for the power spectrum of Kolmogorov's model and determined the effect on propagation through weak turbulence.¹⁶ David L. Fried subsequently extended Tatarskii's results to describe phase distortions of turbulence in terms of Zernike polynomials¹⁷ and derived his famous atmospheric coherence diameter parameter r_0 , as a measure of optical resolution.¹² Nearly all subsequent work in statistical atmospheric optics is based on the contributions of Kolmogorov, Tatarskii, and Fried. The effects of turbulence on wave propagation are critical because the spatial and temporal frequency distribution and optical depth of the phase aberrations drive the design requirements for beacon brightness, wavefront sensing, tracking, electronic processing, and wavefront corrector mirror subsystems in a real-time phase compensation system.

For refractive index fluctuations, the Kolmogorov spatial power spectral density, $\Phi_n(\kappa)$, is given by

$$\Phi_n(\kappa) = 0.033 C_n^2 \kappa^{-11/3} \quad (3)$$

where κ is the spatial wavenumber and C_n^2 is the refractive index structure constant, a measure of the optical strength of turbulence. The wavenumber κ is inversely related to the size of the turbulent eddies in the atmosphere. Equation (3) is valid between two limits of κ called the inertial subrange. The scale size of eddies for the smallest values of κ in the inertial subrange is called the outer scale, denoted $L_0 = 2\pi/\kappa_0$, and is of the order of meters to kilometers. The inner scale, $l_0 = 2\pi/\kappa_m$, or smallest eddies in the inertial subrange, is of the order of millimeters. For values of κ outside the inertial subrange, the von Kármán spectrum is often used in place of Eq. (3) to avoid mathematical complications as κ approaches zero. The von Kármán spectrum has the form

$$\Phi_n(\kappa) \cong \frac{0.033 C_n^2}{(\kappa^2 + \kappa_0^2)^{11/6}} \exp\left(-\frac{\kappa^2}{\kappa_m^2}\right) \quad (4)$$

where κ_0 and κ_m are the spatial wavenumbers for the outer and inner scales.

The Variation of n and the C_n^2 Parameter

Tatarskii introduced structure functions to describe the effect of turbulence on wave propagation. Structure functions describe how a physical parameter is *different* between two points in space or time. The structure function of the index of refraction, $\mathcal{D}_n(\mathbf{r})$ is defined as

$$\mathcal{D}_n(\mathbf{r}) = \{[n_1(\mathbf{r}') - n_1(\mathbf{r}' - \mathbf{r})]^2\} \quad (5)$$

where $n_1(\mathbf{r})$ represents the fluctuating part of the index of refraction and $\{ \dots \}$ represents an ensemble average over the turbulent conditions. The autocorrelation function and power spectral density form a three-dimensional Fourier transform pair (a useful relation used many times in statistical optics), and using that relationship, it can be shown that¹⁸

$$\mathcal{D}_n(r) = C_n^2 r^{2/3} \quad (6)$$

where we have also assumed that the fluctuations are isotropic and dropped the vector dependence on r . Equation (6) is the defining equation for the refractive index structure constant, C_n^2 , in that the numerical coefficient on the right side of Eq. (6) is unity. Defined in this way, C_n^2 is a measure of the optical strength of turbulence.

A few meters above the ground, C_n^2 has an average value of the order of $10^{-14} \text{ m}^{-2/3}$, rapidly decreasing by three or four orders of magnitude at heights above 10 km. Several means have been developed to measure the C_n^2 profile, including in situ instruments on balloons as well as remote sensing optical and radar techniques.^{19,20} Several mathematical models of C_n^2 profiles have been developed based on experimental data. One of the most widely used was suggested by Hufnagel²¹ and modified to include boundary layer effects by Valley.²² The expression for the Hufnagel-Valley model for C_n^2 is

$$C_n^2(h) = 5.94 \cdot 10^{-23} h^{10} e^{-h} \left(\frac{W}{27} \right)^2 + 2.7 \cdot 10^{-16} e^{-2h/3} + A e^{-10h} \quad (7)$$

where h is the height above the site in kilometers, W is an adjustable wind correlating parameter, and A is a scaling constant, almost always taken to be $A = 1.7 \cdot 10^{-14}$. Winker²³ suggested $W = 21$ producing the HV_{5/7} model named from the fact that the resulting profile yields a value of Fried's coherence diameter of $r_0 = 5 \text{ cm}$ and an isoplanatic angle of $\theta_0 = 7 \mu\text{rad}$ for zenith propagation at $0.5 \mu\text{m}$. (The parameters r_0 and θ_0 are defined and their relationships to C_n^2 are given in Sec. 5.4.) The HV_{5/7} model has been widely used in the evaluation of AO system performance. Figure 3 shows how C_n^2 , computed with this model, varies with altitude.

The Hufnagel-Valley model is useful for continental sites or sites with well-developed boundary layers. The C_n^2 profile, which is associated with mountaintop sites that are surrounded by water

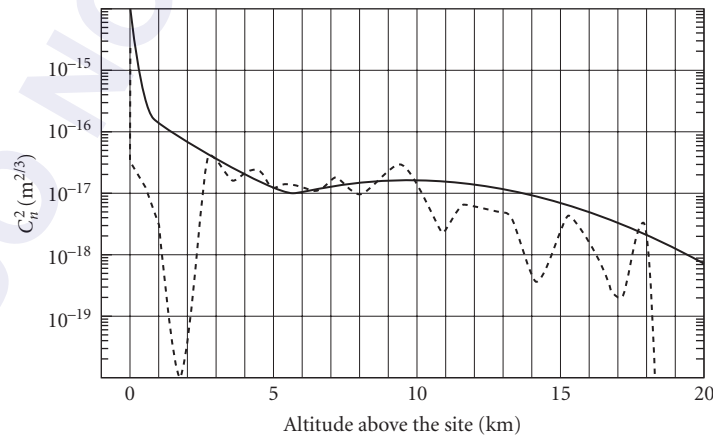


FIGURE 3 C_n^2 profiles for the HV_{5/7} model (solid curve) and for average seeing conditions at Mauna Kea (dotted curve).

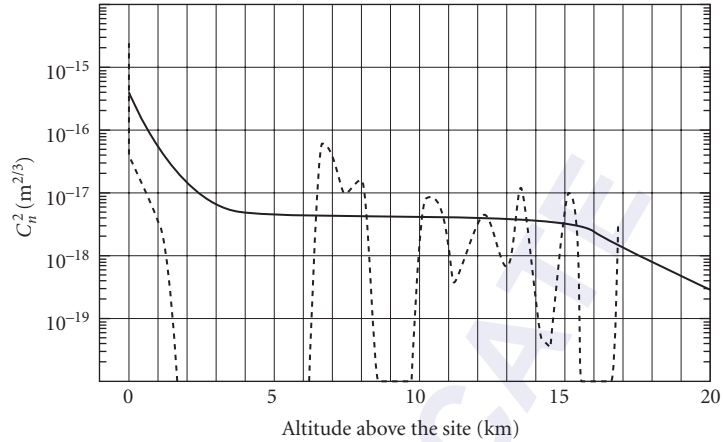


FIGURE 4 Average seeing C_n^2 profile analytical model for Paranal, Chile [site of European Southern Observatory's (ESO's) Very Large Telescope (VLT) (solid curve)], and a C_n^2 profile that is representative of best seeing conditions (occurring less than 10 percent of the time) at Mauna Kea (dotted curve) and represents 0.25 arcsec seeing. Note how the turbulence appears as layers and the lack of a boundary layer in the Mauna Kea profile. Similar conditions have also been observed at Paranal and La Palma.

(e.g., Mauna Kea or La Palma) or mountaintop sites that are close to a coastline [so as to be above the marine boundary layer (e.g., Paranal, Chile)], often exhibits distinct layers of turbulence, but little or no boundary layer in the first few kilometers above the ground. A C_n^2 profile that is representative of average seeing for Mauna Kea is also shown in Fig. 3 and yields an r_0 of 18 cm and θ_0 of 15 μ rad. Figure 4 shows a C_n^2 profile for Mauna Kea representative of the 90th-percentile best seeing (note the distinct layering), giving an r_0 of 45 cm, and an analytical model that is used for average seeing at Paranal,²⁴ giving an r_0 of 18 cm.

Wave Propagation in Turbulence

An electromagnetic wave propagating in the atmosphere must satisfy Maxwell's equations—in particular, the electric and magnetic fields must satisfy the time-independent Helmholtz equation:

$$\nabla^2 U(\mathbf{r}) + n^2 k^2 U(\mathbf{r}) = 0 \quad (8)$$

where $k = 2\pi/\lambda$ and n is the refractive index of air. $U(\mathbf{r})$ is the complex phasor representation of a spatial wave:

$$U(\mathbf{r}) = A(\mathbf{r})e^{i\phi(\mathbf{r})} \quad (9)$$

where the amplitude $A(\mathbf{r})$ and phase $\phi(\mathbf{r})$ are random variables governed by the statistics of the fluctuations of the refractive index, n . There is no closed-form solution for $A(\mathbf{r})$ and $\phi(\mathbf{r})$.

The surfaces in space that are defined by $\phi_0(\mathbf{r}) = \text{constant}$ are wavefronts. For starlight arriving at the top of the atmosphere, the wavefronts are planes of infinite extent. When a plane wave passes

through the atmosphere, whose refractive index varies randomly along each propagation path, the plane wave becomes distorted. The total phase change from that in a vacuum along a path z is

$$\Delta\phi = \phi_0 - k \int_0^z n_1(z) dz \quad (10)$$

where $k = 2\pi/\lambda$, and $n_1(z)$ is the fluctuating part of the index of refraction. If this equation could be evaluated along every direction, one could construct the three-dimensional shape of the wavefront at any position along its propagation path due to phase changes induced by atmospheric turbulence. This is, of course, an intractable problem in terms of a closed-form solution since the fluctuations of the index of refraction are a random process.

A phase deformation that propagates far enough becomes an amplitude fluctuation as different parts of the propagating wavefront combine to create regions of constructive and destructive interference. This effect is called scintillation—the twinkling of starlight is a famous example. Amplitude fluctuations are generally not severe for astronomical observing, and 70 to 80 percent or more of the distortions can be removed by phase-only compensation (by using a steering mirror and deformable mirror). However, both the amplitude and phase of the wave are random variables, and the most general representation of the wave is as given in Eq. (9).

The historically accepted approach to dealing with wave propagation in turbulence was introduced by Tatarskii by use of the Rytov transformation. The Rytov transformation defines a complex quantity ψ as the natural logarithm of U : $\psi = \ln U$. This substitution allows a solution to the wave equation in the form of a multiplicative perturbed version of the free space field:

$$U(\mathbf{r}) = A_0(\mathbf{r}) e^{i\phi_0(\mathbf{r})} e^{\ln[A(\mathbf{r})/A_0(\mathbf{r})] - i[\phi_0(\mathbf{r}) - \phi(\mathbf{r})]} \quad (11)$$

Since the field at any point is the superposition of many independent contributions of propagating waves in the turbulent medium, we expect the fluctuating parts of the field to obey Gaussian statistics (invoking the central limit theorem). This means if $\ln(A/A_0)$ obeys Gaussian statistics, we expect the amplitude fluctuations to be log-normally distributed. Most experimental evidence supports this for weak fluctuations. Tatarskii's results make it possible to compute important characteristic functions and to use the properties of Gaussian random processes to develop practical descriptions of the effects of turbulence.

Fried's Coherence Diameter, r_0 , and the Spatial Scale of Turbulence

Hufnagel and Stanley²⁵ extended Tatarskii's work and developed an expression for the modulation transfer function in terms of the mutual coherence function. Fried²⁶ developed an expression for *phase* structure function, \mathcal{D}_ϕ of a propagating electromagnetic plane wave showing that it is proportional to the 5/3 power of spatial separation, r , and is given by*

$$\mathcal{D}_\phi(r) = \{[\phi(r') - \phi(r' - r)]^2\} = \left[2.91 \left(\frac{2\pi}{\lambda} \right)^2 \int C_n^2(z) dz \right] r^{5/3} \quad (12)$$

where λ is the wavelength of propagation, $C_n^2(z)$ is the position-dependent refractive index structure constant and integration is over the optical path to the source.

*In the more general case where amplitude effects are significant, the *wave* structure function, $\mathcal{D}(r) = \mathcal{D}_x(r) + \mathcal{D}_\phi(r)$, should be used, but since it has exactly the same value as given by Eq. (12), and since conventional AO does not correct intensity fluctuations, the amplitude term will be dropped for the material that is presented here.

Fried²⁶ developed a very useful relationship between the phase structure function and a particular measure of optical resolution—the volume under the two-dimensional optical transfer function. Defined in this way, the seeing-limited resolving power asymptotically approaches a limiting value as the aperture size increases. The limiting value is set by the strength of the turbulence. Fried defined a quantity called the coherence diameter, r_0 , such that the limiting resolution that is obtained in the presence of atmospheric turbulence is the same as that that is obtained by a diffraction-limited lens of diameter r_0 in a vacuum. Its value can be computed from the definition¹⁷

$$r_0 = 0.185\lambda^{6/5}(\cos\psi)^{3/5} \left[\int_0^\infty C_n^2(h) dh \right]^{-3/5} \quad (13)$$

where λ is the wavelength, $C_n^2(h)$ is the vertical profile of the index of refraction structure constant, h is the height above the ground, ψ is the angle between the propagation direction and the zenith, and the integral is evaluated from the ground to an altitude at which $C_n^2(h)$ no longer contributes significantly (typically 20 km). This equation explicitly shows r_0 scales with wavelength as $\lambda^{6/5}$ and with zenith angle as $(\cos\psi)^{3/5}$. By convention, the numerical values of r_0 are usually expressed at 0.5 μm for zenith viewing. At the best astronomical sites, the median value of r_0 is from 15 to 25 cm, corresponding to seeing conditions of $\lambda/r_0 = 0.7$ to 0.4 arcsec. The best intracontinental mountain-top sites exhibit median values of r_0 between 10 and 20 cm, and average intracontinental sites have r_0 values between 5 and 10 cm.²⁰

Given the definition of r_0 in Eq. (13), the phase structure function becomes

$$\mathcal{D}_\phi(r) = 6.88 \left(\frac{r}{r_0} \right)^{5/3} \quad (14)$$

Fried defined r_0 so that the knee in the curve showing resolving power as a function of diameter occurs at $D/r_0 = 1$, where D is the aperture diameter of the imaging telescope. Most scaling laws of interest to AO are expressed in terms of the ratio D/r_0 .

Turbulence Noll²⁷ developed a description of the average spatial content of turbulence-induced wavefront distortion in terms of Zernike polynomials. The Zernike polynomials are often used to describe the classical aberrations of an optical system—tilt, defocus, astigmatism, coma, and spherical aberration. Noll expresses the phase of the distorted wavefront over a circular aperture of radius, R , as

$$\phi(\rho, \theta) = \sum_j a_j Z_j(\rho, \theta) \quad (15)$$

where Z_j is a modified set of the orthogonal Zernike polynomials defined over points in the aperture at reduced radius $\rho = r/R$ and azimuthal angle θ . The Z_j is given by

$$\left. \begin{aligned} Z_{\text{even}j} &= \sqrt{n+1} R_n^m(\rho) \sqrt{2} \cos m\theta \\ Z_{\text{odd}j} &= \sqrt{n+1} R_n^m(\rho) \sqrt{2} \sin m\theta \end{aligned} \right\} m \neq 0 \quad (16)$$

$$Z_j = \sqrt{n+1} R_n^0(\rho) \quad m = 0 \quad (17)$$

where

$$R_n^m(\rho) = \sum_{s=0}^{(n-m)/2} \frac{(-1)^s (n-s)!}{s! [(n+m)/2 - s]! [(n-m)/2 - s]!} \rho^{n-2s} \quad (18)$$

The indices n and m are always positive integers and satisfy $m \leq n$ and $n - m = \text{even}$. The index, j , is a mode-ordering number and depends on n and m .

The a_j are coefficients having mean square values that accurately weight the particular Zernike mode to represent the Kolmogorov distribution of turbulence. Using this approach, Noll computed the values of a_j and showed that the piston-removed wavefront distortion for Kolmogorov turbulence averaged over an aperture of diameter, D , expressed as a mean square value in units of square radians (rad^2) of optical phase is

$$\langle \phi^2 \rangle = 1.0299 \left(\frac{D}{r_0} \right)^{5/3} \quad (19)$$

Table 1 lists the strengths of the first 21 Zernike modes, and the residual mean square distortion as each component is removed. Data in this table show that the two components of tilt (X and Y axes) make up 87 percent of the total wavefront distortion.

Noll's results are useful for estimating the performance of an AO system if one can estimate how many Zernike modes the system will correct (a rule of thumb is approximately one mode is corrected per deformable mirror actuator for systems having large numbers of actuators). Figure 5 shows the mean square residual phase error for a range of values of D/r_0 . Generally, the mean square residual error should be much less than 1 rad^2 . Even if all 21 modes listed in Table 1 are corrected, the mean square residual error for conditions in which D/r_0 is 20, for example, $0.0208(20)^{5/3} = 3.06 \text{ rad}^2$, is a very significant, generally unacceptable error.

When more than 21 modes are corrected, the residual mean square wavefront error is given by²⁷

$$\sigma_M^2 = 0.2944 M^{-(\sqrt{3}/2)} (D/r_0)^{5/3} \quad (20)$$

where M is the number of modes corrected. Figure 6 shows the required number of Zernike modes to be removed as a function of D/r_0 for three image quality conditions. As an example, if r_0 is 12 cm at

TABLE 1 Mean Square Wavefront Distortion Contributions (in rad^2) from the First 21 Zernike Modes of Atmospheric Aberrations and Residual Error As Each Term Is Corrected

Aberration	Contribution	Distortion
All terms	$1.0299(D/r_0)^{5/3}$	$1.0299(D/r_0)^{5/3}$
Z_2 X-tilt	$0.4479(D/r_0)^{5/3}$	$0.582 (D/r_0)^{5/3}$
Z_3 Y-tilt	$0.4479(D/r_0)^{5/3}$	$0.134 (D/r_0)^{5/3}$
Z_4 Defocus	$0.0232(D/r_0)^{5/3}$	$0.111 (D/r_0)^{5/3}$
Z_5 Astigmatism at 45°	$0.0232(D/r_0)^{5/3}$	$0.0880 (D/r_0)^{5/3}$
Z_6 Astigmatism at 0°	$0.0232(D/r_0)^{5/3}$	$0.0649 (D/r_0)^{5/3}$
Z_7 X-coma	$0.0061(D/r_0)^{5/3}$	$0.0587 (D/r_0)^{5/3}$
Z_8 Y-coma	$0.0061(D/r_0)^{5/3}$	$0.0525 (D/r_0)^{5/3}$
Z_9	$0.0062(D/r_0)^{5/3}$	$0.0463(D/r_0)^{5/3}$
Z_{10}	$0.0062(D/r_0)^{5/3}$	$0.0401 (D/r_0)^{5/3}$
Z_{11} Spherical aberration	$0.0024(D/r_0)^{5/3}$	$0.0377 (D/r_0)^{5/3}$
Z_{12}	$0.0025(D/r_0)^{5/3}$	$0.0352 (D/r_0)^{5/3}$
Z_{13}	$0.0024(D/r_0)^{5/3}$	$0.0328 (D/r_0)^{5/3}$
Z_{14}	$0.0024(D/r_0)^{5/3}$	$0.0304 (D/r_0)^{5/3}$
Z_{15}	$0.0025(D/r_0)^{5/3}$	$0.0279 (D/r_0)^{5/3}$
Z_{16}	$0.0012(D/r_0)^{5/3}$	$0.0267 (D/r_0)^{5/3}$
Z_{17}	$0.0012(D/r_0)^{5/3}$	$0.0255 (D/r_0)^{5/3}$
Z_{18}	$0.0012(D/r_0)^{5/3}$	$0.0243 (D/r_0)^{5/3}$
Z_{19}	$0.0011(D/r_0)^{5/3}$	$0.0232 (D/r_0)^{5/3}$
Z_{20}	$0.0012(D/r_0)^{5/3}$	$0.0220 (D/r_0)^{5/3}$
Z_{21}	$0.0012(D/r_0)^{5/3}$	$0.0208 (D/r_0)^{5/3}$

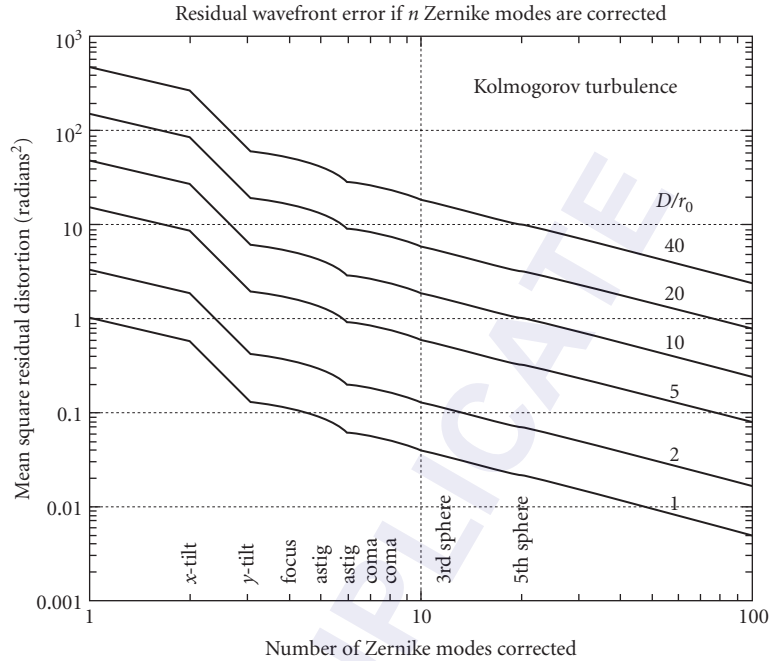


FIGURE 5 Residual mean square phase distortion when n Zernike modes of atmospheric turbulence are corrected for various values of D/r_0 .

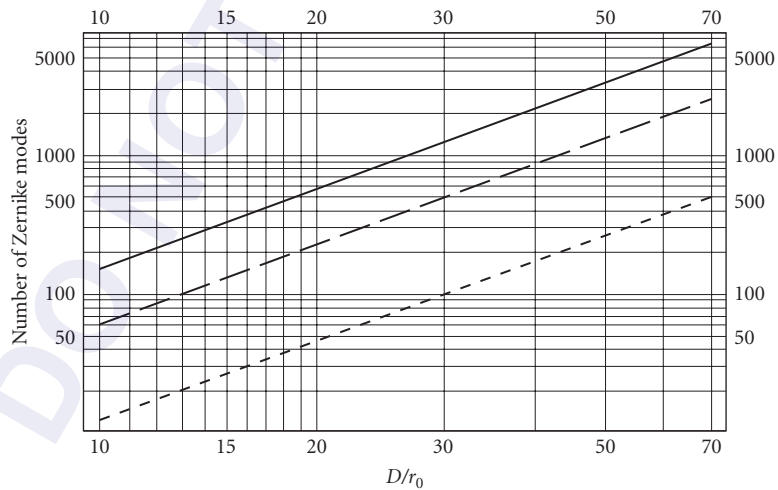


FIGURE 6 Required number of modes to be corrected (roughly the number of actuators or degrees of freedom in the deformable mirror) as a function of D/r_0 . Solid curve: $\lambda/15$ image quality; dashed curve: $\lambda/10$ image quality; dotted curve: $\lambda/5$ image quality. In these curves, r_0 is evaluated at the imaging wavelength.

0.8 μm and the telescope aperture is 3.5 m, $D/r_0 = 29$, and more than 1000 modes must be removed to achieve $\lambda/15$ image quality, but only 100 modes need be removed if $\lambda/5$ image quality is sufficient. If r_0 is 60 cm at 1.6 μm and $D = 10$ m, we must remove nearly 400 modes to achieve $\lambda/15$ image quality. For imaging applications, several techniques for postprocessing are well established²⁸ and allow significant enhancements of images that are obtained in real time with AO. Postprocessing may be justification to allow the relaxation of wavefront quality requirements in an AO system. However, for spectroscopy (or other applications requiring high Strehl ratios in real time), there may be little choice but to design the AO system with the required number of degrees of freedom needed to achieve desired Strehl ratios.

The results in Table 1 can be used to estimate the maximum stroke requirement for the actuators in the deformable mirror. Assuming that tilt will be corrected with a dedicated two-axis beam-steering mirror, the root-mean-square (rms) higher-order distortion is (from Table 1) $0.366(D/r_0)^{5/6}$ rad of phase. A practical rule of thumb is that the deformable mirror surface should be able to correct five times the rms wavefront distortion to get 99 percent of the peak values. Since a reflected wavefront has twice the distortion of the mirror's surface (reducing the stroke requirement by a factor of 2), the maximum actuator stroke requirement for operation becomes

$$\sigma_m (\mu\text{m}) = 0.073 \left(\frac{D}{r_0} \right)^{5/6} \quad (21)$$

where σ_m is the maximum actuator stroke in μm , and r_0 is the value of Fried's coherence diameter at a wavelength of 0.5 μm along the line of sight corresponding to the maximum zenith angle of interest. Figure 7 shows the maximum stroke that is required for several aperture diameters as a function of seeing conditions. This figure assumes the outerscale is much larger than the aperture diameter and Kolmogorov turbulence. If the outerscale is of the order of the aperture size or even smaller, the stroke requirements will be reduced.

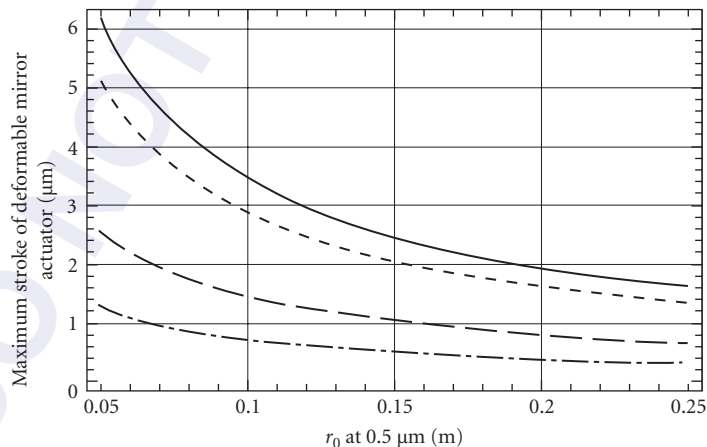


FIGURE 7 Maximum stroke requirement of actuators in the deformable mirror. The curves are for aperture diameters (from top to bottom) of 10, 8, 3.5, and 1.5 m. The value of r_0 is for an observing wavelength of 0.5 μm along the line of sight corresponding to the maximum zenith angle of interest, and the stroke is in micrometers. These curves are independent of wavelength since they do not include any dispersion effects.

Atmospheric Tilt and Its Effect on the Strehl Ratio

Tilt comprises 87 percent of the power in atmospheric turbulence-induced wavefront distortion (see the preceding section). The Strehl ratio of an image that is due to jitter alone is approximated to better than 10 percent by the equation^{29,30}

$$SR_{\text{tilt}} = \frac{1}{1 + \frac{\pi^2}{2} \left(\frac{\sigma_{\theta}}{\lambda/D} \right)^2} \quad (22)$$

where σ_{θ} is the one-axis rms angular jitter. Figure 8 shows how the one-axis rms jitter affects system Strehl ratio. Greenwood and Fried³¹ derived bandwidth requirements for AO systems and developed expressions for the variance of full-aperture tilt that is induced by turbulence. The one-axis angular jitter variance (units of angular rad²) is given by

$$\sigma_{\theta G}^2 = 0.170 \left(\frac{\lambda}{D} \right)^2 \left(\frac{D}{r_0} \right)^{5/3} \quad (23)$$

where $\sigma_{\theta G}$ is the rms G -tilt. G -tilt is the average gradient over the wavefront and is well approximated by most centroid trackers. For Z -tilt (the direction of the normal of the best-fit plane to the wavefront distortion), the coefficient in Eq. (23) is 0.184. Consequently, a tracking sensor that measures the centroid of the focused image improperly estimates the Z -tilt; this has been called centroid anisoplanatism.³⁰

Figure 9 shows the dependence of $\sigma_{\theta G}$ on aperture diameter for two seeing conditions, one representative of a continental site and one of a mountaintop island like Mauna Kea, Hawaii.

The effect of atmospheric tilt on the Strehl ratio can be seen by substituting Eq. (23) into Eq. (22). Figure 10 shows that the loss in the Strehl ratio due to jitter alone is significant. For $D/r_0 = 1$, the atmospherically induced jitter limits the Strehl ratio to 0.54. It is most important to effectively control image motion that is created by full-aperture atmospheric tilt.

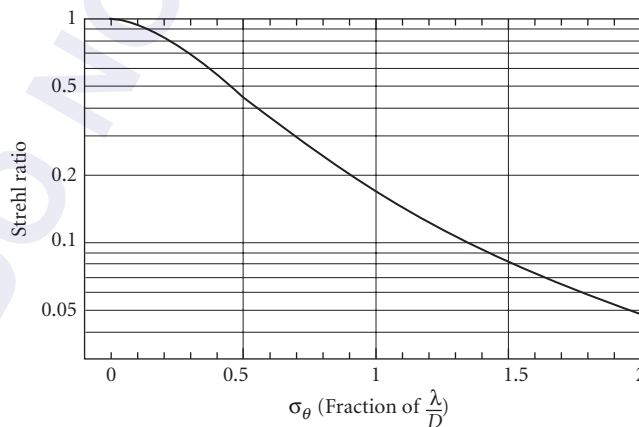


FIGURE 8 Strehl ratio versus one-axis rms jitter expressed as a fraction of λ/D .

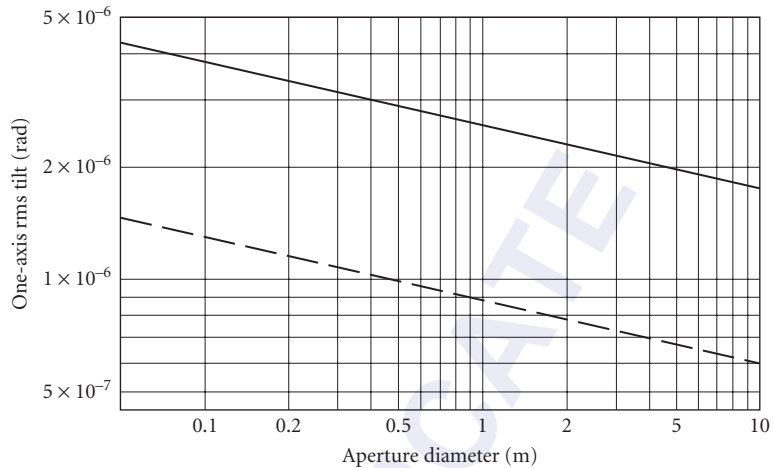


FIGURE 9 One-axis rms tilt for atmospheric turbulence. The upper line corresponds to an r_0 of 5 cm and the lower curve corresponds to an r_0 of 18 cm, both referenced to $0.5 \mu\text{m}$ and zenith. The value of full-aperture tilt is independent of wavelength for a given seeing condition.

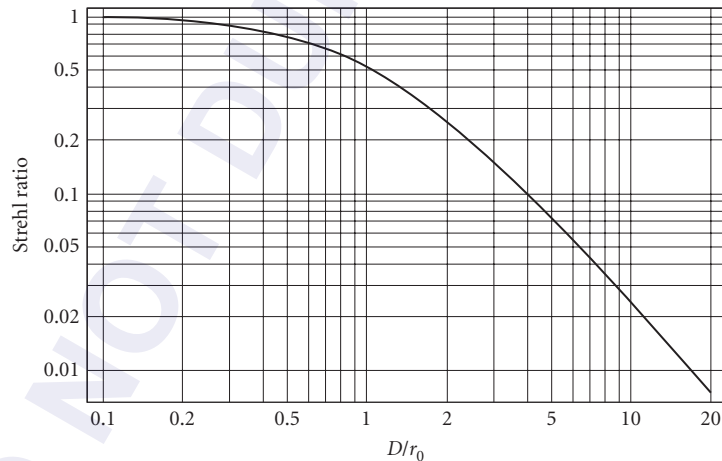


FIGURE 10 The effect of turbulence-induced jitter on system Strehl ratio.

Tracking: Bandwidth and Steering Mirror Stroke Requirements

Tyler³² has considered the problem of required tracking bandwidth by computing the tilt power spectrum and determining the error rejection achieved by a steering mirror under control of an RC-type servo. He considers both G -tilt and Z -tilt. Temporal fluctuations are derived by translating “frozen-turbulence” realizations past the telescope’s aperture. The hypothesis is that the internal structure of the turbulence changes more slowly than its mass transport by the wind (but this assumption may not suffice for today’s giant-aperture telescopes). The results for power spectra are

in integral form, requiring numerical evaluation. Sasiela³³ uses Mellin transforms to derive a solution in terms of a power series of ratios of gamma functions. The asymptotic limits for high and low frequencies are, however, tractable in a form that is easily evaluated. We consider here the results for only G -tilt since that is the most likely implementation in a real system (approximated by a centroid tracker). The low-frequency limit for G -tilt is

$$\lim_{f \rightarrow 0} F_{\phi_G}(f) = 0.804 D^{-1/3} \sec \psi f^{-2/3} \int_0^{\infty} C_n^2(h) [\nu(h)/D]^{-1/3} dh \quad (24)$$

and the high-frequency limit is

$$\lim_{f \rightarrow \infty} F_{\phi_G}(f) = 0.0110 D^{-1/3} \sec \psi f^{-11/3} \int_0^{\infty} C_n^2(h) [\nu(h)/D]^{8/3} dh \quad (25)$$

where $F_{\phi_G}(f)$ is the one-axis G -tilt power spectrum in rad^2/Hz , D is the aperture diameter, ψ is the zenith angle, and $\nu(h)$ is the transverse wind component along a vertical profile. Note that the tilt spectrum goes as $f^{-2/3}$ at low frequencies and as $f^{-11/3}$ at high frequencies.

The disturbances described in Eqs. (24) and (25) can be partially corrected with a fast-steering mirror having an RC filter response, $H(f, f_c)$, of the form

$$H(f, f_c) = \frac{1}{1 + i \frac{f}{f_c}} \quad (26)$$

where f_c represents a characteristic frequency [usually the 3 decibel (dB) response]. The uncorrected power in the residual errors can be computed from control theory using the relation

$$\sigma_{\theta_G}^2 = \int_0^{\infty} |1 - H(f, f_c)|^2 F_{\phi_G}(f) df \quad (27)$$

where $F_{\phi_G}(f)$ is the power spectrum (rad^2/Hz) of G -tilt induced by optical turbulence.

For a steering mirror response function given by Eq. (26), the G -tilt variance is given by

$$\sigma_{\theta_G}^2 = \int_0^{\infty} \frac{(f/f_{3\text{dB}})^2}{1 + (f/f_{3\text{dB}})^2} F_{\phi_G}(f) df \quad (28)$$

In the limit of a small (zero) servo bandwidth, this equation represents the tilt that is induced by atmospheric turbulence. Expressing the variance of the jitter in terms of the diffraction angle (λ/D) results in the same expression as Eq. (23).

What does the servo bandwidth need to be? Tyler defines a G -tilt tracking frequency characteristic of the atmospheric turbulence profile and wind velocity as

$$f_{T_G} = 0.331 D^{-1/6} \lambda^{-1} \sec^{1/2} \psi \left[\int_0^{\infty} C_n^2(h) \nu^2(h) dh \right]^{1/2} \quad (29)$$

such that the variance of the tilt can be expressed as

$$\sigma_{\theta}^2 = \left(\frac{f_{T_G}}{f_{3\text{dB}}} \right)^2 \left(\frac{\lambda}{D} \right)^2 \quad (30)$$

Three wind velocity profiles are shown in Fig. 11. The analytical model of Bufton³⁴ is given by

$$v_B(h) = 5 + 37 e^{-(h-12)^2/25} \quad (31)$$

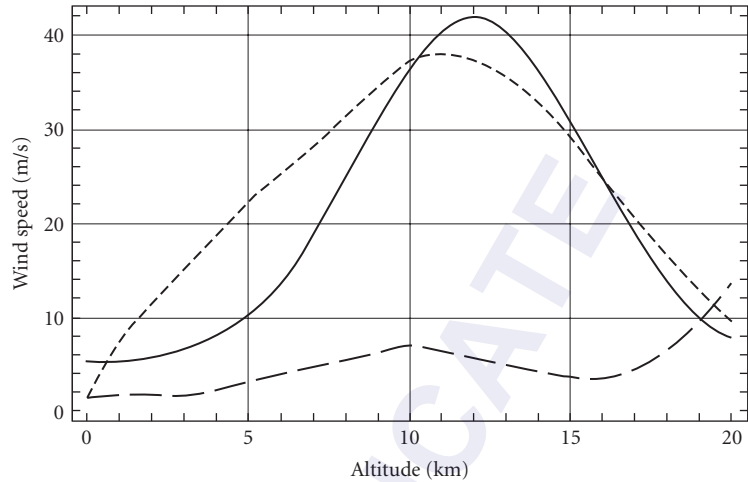


FIGURE 11 Wind profiles. The solid curve is an analytical model by Bufton. The dotted and dashed curves are based on two decades of rawinsonde measurements at Albuquerque, New Mexico, and are known as SOR winter and SOR summer models, respectively. These two models show clearly the absence of a jet stream in the summer months.

where h is the height above ground in km and v_B the wind speed in m/s. The other two curves in Fig. 11 are based on measurements made with National Weather Service sounding balloons launched in Albuquerque, New Mexico. They show a marked contrast between winter and summer months and the effect of the jet stream.

When the 3-dB closed-loop tracking servo frequency equals the G -tilt tracking frequency (the Tyler frequency), the one-axis rms jitter is equal to λ/D . The track bandwidth needs to be four times the Tyler frequency in order for $\sigma_{\theta G}$ to be $1/4(\lambda/D)$, a condition that provides a Strehl ratio due to tilt of greater than 80 percent. Figure 12 shows the required 3-dB tracking bandwidths to

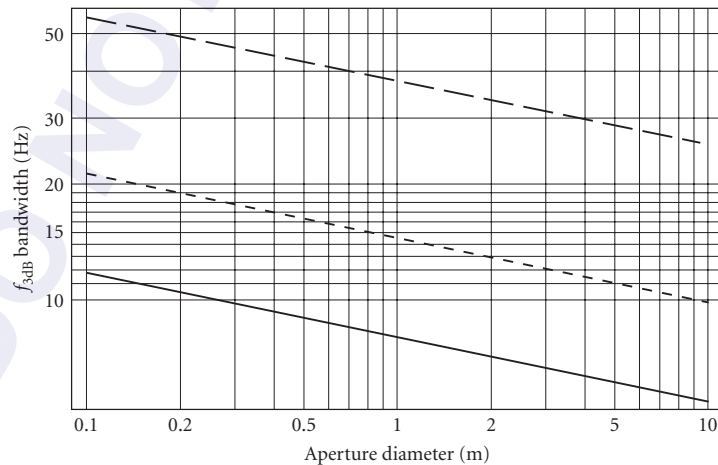


FIGURE 12 Required 3-dB tracking bandwidth to achieve Strehl ratio of 0.82. The dashed curve is for the SOR winter wind and $HV_{5/7}$ profiles; the dotted curve is for the Mauna Kea average C_n^2 profile and Bufton wind profile; the solid curve is for the SOR summer wind and $HV_{5/7}$ C_n^2 profiles. All three curves are for $0.5 \mu\text{m}$ wavelength.

meet this criterion for three combinations of wind and C_n^2 profiles as a function of aperture diameter. Other combinations can be computed easily by using Eq. (29) and the wind profiles that are shown in Fig. 11.

The one-axis rms tilt that is induced by the atmosphere is given by Eq. (23). When specifying the maximum excursion of a tilt mirror, the mirror should be able to move five times σ_{θ_G} in order to accommodate over 99 percent of the tilt spectrum. We also need to account for the optical magnification between the tilt mirror and the primary mirror of the telescope. If, for instance, the tilt mirror has a diameter of 10 cm, the telescope aperture is 10 m, and r_0 is 18 cm (dashed curve in Fig. 9), the required maximum stroke of the tilt mirror is $0.6 \mu\text{rad}$ (from Fig. 9) times 50 ($5\sigma_{\theta_G} \times$ optical magnification of 10) or $30 \mu\text{rad}$.

Higher-Order Phase Fluctuations and Bandwidth Requirements

Time delays that are caused by detector readout and data processing create errors in the phase correction that is applied to the deformable mirror—the correction being applied is relevant to an error that is measured at an earlier time. The control system concepts discussed in the preceding two sections are also relevant to higher-order errors. The residual errors can be computed knowing the closed-loop servo response, $H(f, f_c)$, where f_c represents a characteristic frequency (usually the 3-dB response) and the disturbance. The uncorrected power is similar to Eq. (27) and is given by

$$\sigma_r^2 = \int_0^\infty |1 - H(f, f_c)|^2 F_\phi(f) df \quad (32)$$

where $F_\phi(f)$ is the power spectrum (rad^2/Hz) of the phase distortions that are induced by optical turbulence. At high frequencies, Greenwood³⁵ showed that $F_\phi(f)$ is proportional to $f^{-8/3}$ and is given by

$$\lim_{f \rightarrow \infty} F_\phi(f) = 0.0326 k^2 f^{-8/3} \int_0^\infty C_n^2(z) v^{5/3}(z) dz \quad (33)$$

where $k = 2\pi/\lambda$, f is the frequency in Hertz, and $v(z)$ is the transverse component of the wind profile. The closed-loop response of the AO servo can be modeled as a conventional resistorcapacitor (RC) filter making

$$H(f, f_c) = \frac{1}{1 + i \frac{f}{f_c}} \quad (34)$$

the characteristic frequency becomes

$$f_c = \left[0.102 \left(\frac{k}{\sigma_r} \right)^2 \int_0^\infty C_n^2(z) v^{5/3}(z) dz \right]^{3/5} \quad (35)$$

Figure 11 shows three wind profiles, one analytical and two based on measured data. For the $HV_{5/7} C_n^2$ profile and the SOR summer and winter wind models and $\sigma_r = 0.2\pi$, then $f_c = 23$ and 95 Hz, respectively.

The value of f_c for which $\sigma_r = 1$ is known as the Greenwood frequency and is given explicitly as

$$f_G = \frac{2.31}{\lambda^{6/5}} \left[\int_0^L C_n^2(z) \nu(z)^{5/3} dz \right]^{3/5} \quad (36)$$

Tyler³⁶ shows that the mean square residual phase in an AO control system having a -3 -dB closed-loop bandwidth of $f_{3\text{dB}}$ Hz is

$$\sigma_{\phi_{\text{servo}}}^2 = (f_G / f_{3\text{dB}})^{5/3} \quad (37)$$

In a real system, this means that the closed-loop servo bandwidth should be several times the Greenwood frequency in order to “stay up with the turbulence” and keep the residual wavefront error to acceptable levels. Note that the Greenwood frequency scales as $\lambda^{-6/5}$ (the inverse of r_0 and θ_0 scaling) and that zenith angle scaling depends on the specific azimuthal wind direction.

Anisoplanatism

Consider two point sources in the sky (a binary star) and an instrument to measure the wavefront distortion from each. As the angular separation between the sources increases, the wavefront distortions from each source become decorrelated. The *isoplanatic angle* is that angle for which the difference in mean square wavefront distortion is 1 rad². Fried³⁷ defined the isoplanatic angle, θ_0 , as

$$\theta_0 = 0.058 \lambda^{6/5} (\sec \psi)^{-8/5} \left[\int_0^\infty C_n^2(h) h^{5/3} dh \right]^{-3/5} \quad (38)$$

where the integral is along a vertical path through the atmospheric turbulence. Note that θ_0 scales as $\lambda^{6/5}$ (the same as r_0), but as $(\cos \psi)^{8/5}$ with zenith angle. Fried has shown that if a very large diameter ($D/r_0 \gg 1$) imaging system viewing along one path is corrected by an AO system, the optical transfer function of the system viewing along a different path separated by an angle θ is reduced by a factor $\exp[-(\theta/\theta_0)^{5/3}]$. (For smaller diameters, there is not so much reduction.) The numerical value of the isoplanatic angle for zenith viewing at a wavelength of $0.5 \mu\text{m}$ is of the order of a few arc seconds for most sites. The isoplanatic angle is strongly influenced by high-altitude turbulence (note the $h^{5/3}$ weighting in the preceding definition). The small value of the isoplanatic angle can have very significant consequences for AO by limiting the number of natural stars suitable for use as reference beacons and by limiting the corrected field of view to only a few arc seconds.

Turbulence on Imaging and Spectroscopy

The optical transfer function (OTF) is one of the most useful performance measures for the design and analysis of AO imaging systems. The OTF is the Fourier transform of the optical system's point spread function. For an aberration-free circular aperture, the OTF for a spatial frequency \mathbf{f} is well known³⁸ to be

$$H_0(\mathbf{f}) = \frac{2}{\pi} \left[\arccos\left(\frac{\mathbf{f}\bar{\lambda}F}{D}\right) - \frac{\mathbf{f}\bar{\lambda}F}{D} \sqrt{1 - \left(\frac{\mathbf{f}\bar{\lambda}F}{D}\right)^2} \right] \quad (39)$$

where D is the aperture diameter, F is the focal length of the system, and $\bar{\lambda}$ is the average imaging wavelength. Notice that the cutoff frequency (where the OTF of an aberration-free system reaches 0) is equal to $D/\bar{\lambda}F$.

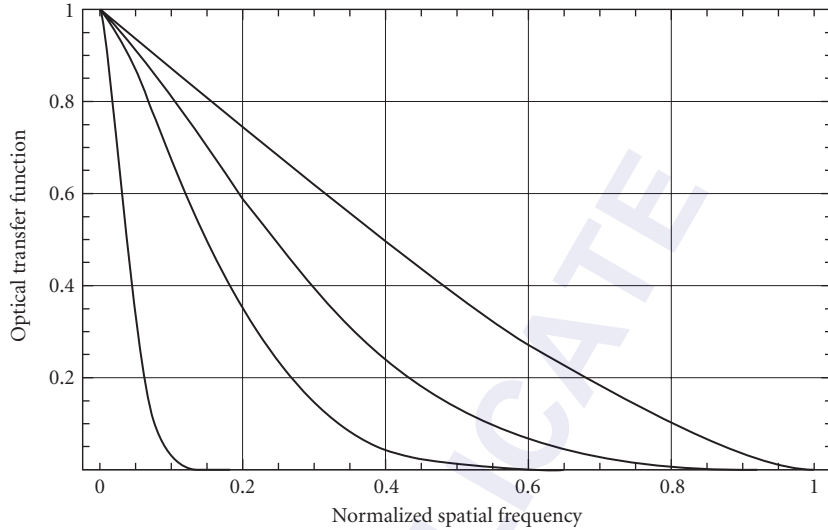


FIGURE 13 The OTF due to the atmosphere. Curves are shown for (top to bottom) $D/r_0 = 0.1, 1, 2,$ and 10 . The cutoff frequency is defined when the normalized spatial frequency is 1.0 and has the value $D/\lambda F$, where F is the focal length of the telescope and D is its diameter.

Fried²⁶ showed that for a long-exposure image, the OTF is equal to $H_0(\mathbf{f})$ times the long-exposure OTF due to turbulence, given by

$$H_{LE}(\mathbf{f}) = \exp\left\{-\frac{1}{2}D(\bar{\lambda}F\mathbf{f})\right\} = \exp\left\{-\frac{1}{2}6.88\left(\frac{\bar{\lambda}F\mathbf{f}}{r_0}\right)^{5/3}\right\} \quad (40)$$

where D is the wave structure function and where we have substituted the phase structure function that is given by Eq. (14). Figure 13 shows the OTF along a radial direction for a circular aperture degraded by atmospheric turbulence for values of D/r_0 ranging from 0.1 to 10 . Notice the precipitous drop in the OTF for values of $D/r_0 > 2$. The objective of an AO system is, of course, to restore the high spatial frequencies that are lost due to turbulence.

Spectroscopy is a very important aspect of observational astronomy and is a major contributor to scientific results. The goals of AO for spectroscopy are somewhat different than for imaging. For imaging, it is important to stabilize the corrected point spread function in time and space so that postprocessing can be performed. For spectroscopy, high Strehl ratios are desired in real time. The goal is flux concentration and getting the largest percentage of the power collected by the telescope through the slit of the spectrometer. A 4-m telescope is typically limited to a resolving power of $R \sim 50,000$. Various schemes have been tried to improve resolution, but the instruments become large and complex.

However, by using AO, the corrected image size decreases linearly with aperture size, and very high resolution spectrographs are, in principle, possible without unreasonable-sized gratings. A resolution of $700,000$ was demonstrated on a 1.5-m telescope corrected with AO.³⁹ Tyler and Ellerbroek⁴⁰ have estimated the sky coverage at the galactic pole for the Gemini North 8-m telescope at Mauna Kea as a function of the slit power coupling percentage for a 0.1-arcsec slit width at J, H, and K bands in the near IR. Their results are shown in Fig. 14 for laser guide star (top curves) and natural guide star (lower curves) operation.

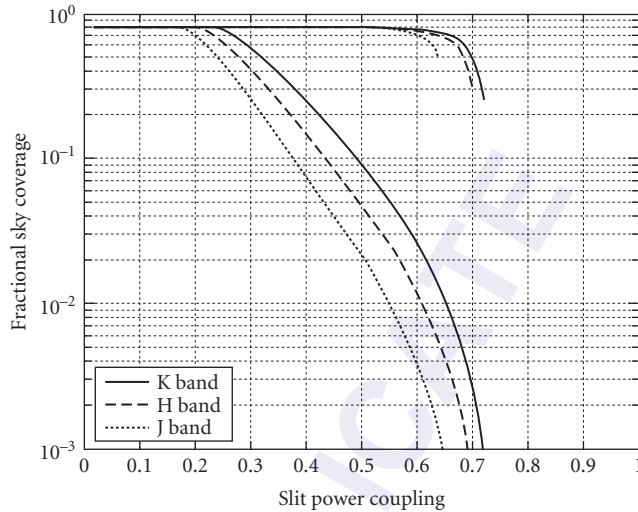


FIGURE 14 Spectrometer slit power sky coverage at Gemini North using NGS (lower curves) and LGS adaptive optics. See text for details.

The higher-order AO system that was analyzed⁴⁰ for the results in Fig. 14 employs a 12-by-12 subaperture Shack-Hartmann sensor for both natural guide star (NGS) and laser guide star (LGS) sensing. The spectral region for NGS is from 0.4 to 0.8 μm and the 0.589- μm LGS is produced by a 15-watt laser to excite mesospheric sodium at an altitude of 90 km. The wavefront sensor charge-coupled device (CCD) array has 3 electrons of read noise per pixel per sample and the deformable mirror is optically conjugate to a 6.5-km range from the telescope. The tracking is done with a 2-by-2-pixel sensor operating at J + H bands with 8 electrons of read noise. Ge et al.⁴¹ have reported similar results with between 40 and 60 percent coupling for LGSs and nearly 70 percent coupling for NGSs brighter than 13th magnitude.

5.5 AO HARDWARE AND SOFTWARE IMPLEMENTATION

Tracking

The wavefront tilt averaged over the full aperture of the telescope accounts for 87 percent of the power in turbulence-induced wavefront aberrations. Full-aperture tilt has the effect of blurring images and reducing the Strehl ratio of point sources. The Strehl ratio due to tilt alone is given by^{29,30}

$$\text{SR}_{\text{tilt}} = \frac{1}{1 + \frac{\pi^2}{2} \left(\frac{\sigma_{\theta}}{\bar{\lambda}/D} \right)^2} \quad (41)$$

where σ_{θ} is the one-axis rms full-aperture tilt error, $\bar{\lambda}$ is the imaging wavelength, and D is the aperture diameter. Figure 8 is a plot showing how the Strehl ratio drops as jitter increases. This figure

shows that in order to maintain a Strehl ratio of 0.8 due to tilt alone, the image must be stabilized to better than $0.25\lambda/D$. For an 8-m telescope imaging at $1.2\ \mu\text{m}$, $0.25\lambda/D$ is 7.75 milliarcsec.

Fortunately, there are several factors that make tilt sensing more feasible than higher-order sensing for faint guide stars that are available in any field. First, we can use the entire aperture, a light gain over higher-order sensing of roughly $(D/r_0)^2$. Second, the image of the guide star will be compensated well enough that its central core will have a width of approximately λ/D rather than λ/r_0 (assuming that tracking and imaging are near the same wavelength). Third, we can track with only four discrete detectors, making it possible to use photon-counting avalanche photodiodes (or other photon-counting sensors), which have essentially no noise at the short integration times required (~ 10 ms). Fourth, Tyler⁴² has shown that the fundamental frequency that determines the tracking bandwidth is considerably less (by as much as a factor of 9) than the Greenwood frequency, which is appropriate for setting the servo bandwidth of the deformable mirror control system. One must, however, include the vibrational disturbances that are induced into the line of sight by high-frequency jitter in the telescope mount, and it is dangerous to construct a simple rule of thumb comparing tracking bandwidth with higher-order bandwidth requirements.

The rms track error is given approximately by the expression

$$\sigma_\theta = 0.58 \frac{\text{angular image size}}{\text{SNR}_V} \quad (42)$$

where SNR_V is the voltage signal-to-noise ratio in the sensor. An error of $\sigma_\theta = \lambda/4D$ will provide a Strehl ratio of approximately 0.76. If the angular image size is λ/D , the SNR_V needs to be only 2. Since we can count on essentially shot-noise-limited performance, in theory we need only four detected photoelectrons per measurement under ideal conditions (in the real world, we should plan on needing twice this number). Further averaging will occur in the control system servo.

With these assumptions, it is straightforward to compute the required guide star brightness for tracking. Results are shown in Fig. 15 for 1.5-, 3.5-, 8-, and 10-m telescope apertures, assuming a 500-Hz sample rate and twice-diffraction-limited AO compensation of higher-order errors of the guide star at the track wavelength. These results are not inconsistent with high-performance tracking

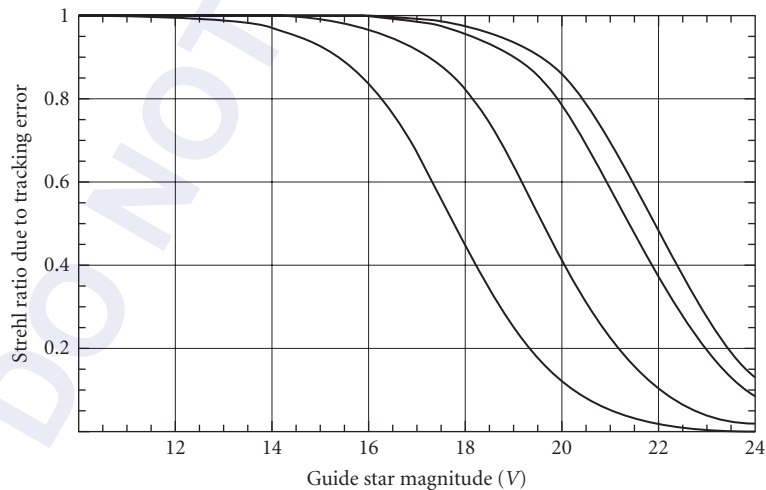


FIGURE 15 Strehl ratio due to tracking jitter. The curves are (top to bottom) for 10-, 8-, 3.5-, and 1.5-m telescopes. The assumptions are photon shot-noise-limited sensors, 25 percent throughput to the track sensor, 200-nm optical bandwidth, and 500-Hz sample rate. The track wavelength is $0.9\ \mu\text{m}$.

systems already in the field, such as the HRCam system that has been so successful on the Canada-France-Hawaii Telescope at Mauna Kea.⁴³

As mentioned previously, the track sensor can consist of just a few detectors implemented in a quadrant-cell algorithm or with a two-dimensional CCD or IR focal plane array. A popular approach for the quad-cell sensor is to use an optical pyramid that splits light in four directions to be detected by individual detectors. Avalanche photodiode modules equipped with photon counter electronics have been used very successfully for this application. These devices operate with such low dark current that they are essentially noise-free and performance is limited by shot noise in the signal and quantum efficiency.

For track sensors using focal plane arrays, different algorithms can be used depending on the tracked object. For unresolved objects, centroid algorithms are generally used. For resolved objects (e.g., planets, moons, asteroids, etc.), correlation algorithms may be more effective. In one instance at the Starfire Optical Range, the highest-quality images of Saturn were made by using an LGS to correct higher-order aberrations and a correlation tracker operating on the rings of the planet provided tilt correction to a small fraction of λ/D .⁴⁴

Higher-Order Wavefront Sensing and Reconstruction: Shack-Hartmann Technique

Higher-order wavefront sensors determine the *gradient*, or slope, of the wavefront measured over subapertures of the entrance pupil, and a dedicated controller maps the slope measurements into deformable mirror actuator voltages. The traditional approach to AO has been to perform these functions in physically different pieces of hardware—the wavefront sensor, the wavefront reconstructor and deformable mirror controller, and the deformable mirror. Over the years, several optical techniques (Shack-Hartmann, various forms of interferometry, curvature sensing, phase diversity, and many others) have been invented for wavefront sensing. A large number of wavefront reconstruction techniques, geometries, and predetermined or even adaptive algorithms have also been developed. A description of all these techniques is beyond the scope of this document, but the interested reader should review work by Wallner,⁴⁵ Fried,⁴⁶ Wild,⁴⁷ and Ellerbroek and Rhoadarmer.⁴⁸

A wavefront sensor configuration in wide use is the Shack-Hartmann sensor.⁴⁹ We will use it here to discuss wavefront sensing and reconstruction principles. Figure 16 illustrates the concept. An array of lenslets is positioned at a relayed image of the exit pupil of the telescope. Each lenslet represents a subaperture—in the ideal case, sized to be less than or equal to r_0 at the sensing wavelength.

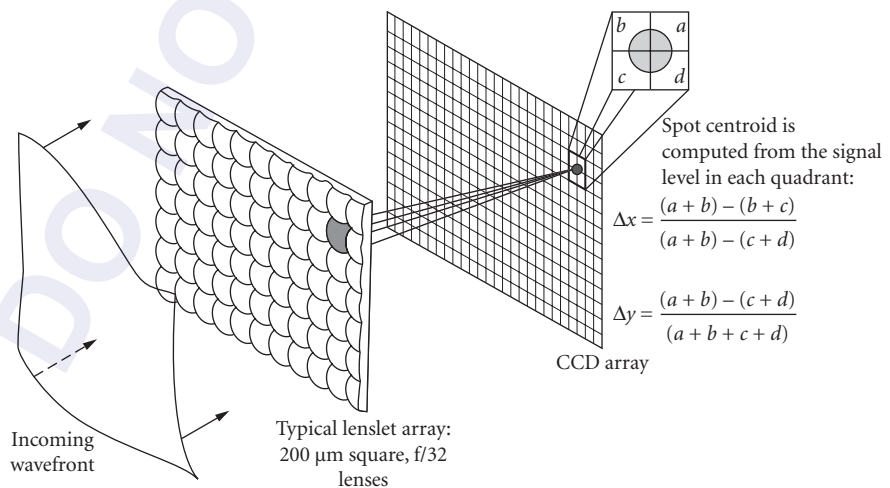


FIGURE 16 Geometry of a Shack-Hartmann sensor.

For subapertures that are roughly r_0 in size, the wavefront that is sampled by the subaperture is essentially flat but tilted, and the objective of the sensor is to measure the value of the subaperture tilt. Light collected by each lenslet is focused to a spot on a two-dimensional detector array. By tracking the position of the focused spot, we can determine the X- and Y-tilt of the wavefront averaged over the subaperture defined by the lenslet. By arranging the centers of the lenslets on a two-dimensional grid, we generate gradient values at points in the centers of the lenslets. Many other geometries are possible resulting in a large variety of patterns of gradient measurements.

Figure 17 is a small-scale example from which we can illustrate the basic equations. In this example, we want to estimate the phase at 16 points on the corners of the subapertures (the Fried geometry of the Shack-Hartmann sensor), and to designate these phases as $\phi_1, \phi_2, \phi_3, \dots, \phi_{16}$, using wavefront gradient measurements that are averaged in the centers of the subapertures, $S_{1x}, S_{1y}, S_{2x}, S_{2y}, \dots, S_{9x}, S_{9y}$. It can be seen by inspection that

$$\begin{aligned} S_{1x} &= \frac{(\phi_2 + \phi_6) - (\phi_1 + \phi_5)}{2d} \\ S_{1y} &= \frac{(\phi_5 + \phi_6) - (\phi_1 + \phi_2)}{2d} \\ &\dots \\ S_{9x} &= \frac{(\phi_{16} + \phi_{12}) - (\phi_{11} + \phi_{15})}{2d} \\ S_{9y} &= \frac{(\phi_{15} + \phi_{16}) - (\phi_{11} + \phi_{12})}{2d} \end{aligned}$$

This system of equations can be written in the form of a matrix as

$$\begin{bmatrix} S_{1x} \\ S_{2x} \\ S_{3x} \\ S_{5x} \\ S_{6x} \\ S_{7x} \\ S_{8x} \\ S_{9x} \\ S_{1y} \\ S_{2y} \\ S_{3y} \\ S_{4y} \\ S_{5y} \\ S_{6y} \\ S_{7y} \\ S_{8y} \\ S_{9y} \end{bmatrix} = \frac{1}{2d} \begin{bmatrix} -1 & 1 & 0 & 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 1 & 0 & 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -1 & 1 & 0 & 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & -1 & 1 & 0 & 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -1 & 1 & 0 & 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & 1 & 0 & 0 & -1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & 1 & 0 & 0 & -1 & 1 & 0 \\ -1 & -1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & -1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & -1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} \phi_1 \\ \phi_2 \\ \phi_3 \\ \phi_4 \\ \phi_5 \\ \phi_6 \\ \phi_7 \\ \phi_8 \\ \phi_9 \\ \phi_{10} \\ \phi_{11} \\ \phi_{12} \\ \phi_{13} \\ \phi_{14} \\ \phi_{15} \\ \phi_{16} \end{bmatrix} \quad (43)$$

These equations express gradients in terms of phases,

$$\mathbf{S} = \mathbf{H}\Phi \quad (44)$$

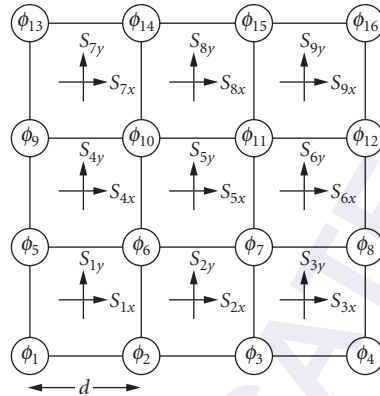


FIGURE 17 A simple Shack-Hartmann sensor in the Fried geometry.

where Φ is a vector of the desired phases, H is the measurement matrix, and S is a vector of the measured slopes.

However, in order to control the actuators in the deformable mirror, we need a control matrix, M , that maps subaperture slope measurements into deformable mirror actuator control commands. In essence, we need to invert Eq. (44) to make it of the form

$$\Phi = MS \quad (45)$$

where M is the desired control matrix. The most straightforward method to derive the control matrix, M , is to minimize the difference in the measured wavefront slopes and the actual slopes on the deformable mirror. We can do this by a maximum a posteriori method accounting for actuator influence functions in the deformable mirror, errors in the wavefront slope measurements due to noise, and statistics of the atmospheric phase distortions. If we do not account for any effects except the geometry of the actuators and the wavefront subapertures, the solution is a least-squares estimate (the most widely implemented to date). It has the form

$$\Phi = [H^T H]^{-1} H^T S \quad (46)$$

and is the *pseudoinverse* of H . (For our simple geometry shown in Fig. 17, the pseudoinverse solution is shown in Fig. 18.) Even this simple form is often problematical since the matrix $[H^T H]$ is often singular or acts as a singular matrix from computational roundoff error and cannot be inverted.

However, in these instances, singular-value-decomposition (SVD) algorithms can be used to directly compute a solution for the inverse of H . Singular value decomposition decomposes an $m \times n$ matrix into the product of an $m \times n$ matrix (U), an $n \times n$ diagonal matrix (D), and an $n \times n$ square matrix (V). So that

$$H = UDV^T \quad (47)$$

and H^{-1} is then

$$H^{-1} = VD^{-1}U^T \quad (48)$$

$-\frac{7}{16}$	$-\frac{9}{80}$	$-\frac{1}{20}$	$\frac{9}{80}$	$-\frac{3}{20}$	$-\frac{1}{80}$	$-\frac{1}{20}$	$\frac{1}{80}$	$-\frac{1}{16}$	$-\frac{7}{16}$	$\frac{9}{80}$	$-\frac{1}{20}$	$-\frac{9}{80}$	$-\frac{3}{20}$	$\frac{1}{80}$	$-\frac{1}{20}$	$-\frac{1}{80}$	$-\frac{1}{16}$
$\frac{11}{60}$	$-\frac{61}{240}$	$-\frac{1}{16}$	$-\frac{29}{240}$	$\frac{1}{20}$	$-\frac{19}{240}$	$\frac{1}{16}$	$-\frac{11}{240}$	$\frac{1}{60}$	$-\frac{11}{60}$	$-\frac{61}{240}$	$\frac{1}{16}$	$-\frac{29}{240}$	$-\frac{1}{20}$	$-\frac{19}{240}$	$-\frac{1}{16}$	$-\frac{11}{240}$	$-\frac{1}{60}$
$\frac{1}{16}$	$\frac{61}{240}$	$-\frac{11}{60}$	$\frac{19}{240}$	$-\frac{1}{20}$	$\frac{29}{240}$	$-\frac{1}{60}$	$\frac{11}{240}$	$-\frac{1}{16}$	$\frac{1}{16}$	$-\frac{61}{240}$	$-\frac{11}{60}$	$-\frac{19}{240}$	$-\frac{1}{20}$	$-\frac{29}{240}$	$-\frac{1}{60}$	$-\frac{11}{240}$	$-\frac{1}{16}$
$\frac{1}{20}$	$\frac{9}{80}$	$\frac{7}{16}$	$\frac{1}{80}$	$\frac{3}{20}$	$-\frac{9}{80}$	$\frac{1}{16}$	$-\frac{1}{80}$	$\frac{1}{20}$	$-\frac{1}{20}$	$\frac{9}{80}$	$-\frac{7}{16}$	$\frac{1}{80}$	$-\frac{3}{20}$	$-\frac{9}{80}$	$-\frac{1}{16}$	$-\frac{1}{80}$	$-\frac{1}{20}$
$-\frac{11}{60}$	$-\frac{29}{240}$	$-\frac{1}{16}$	$-\frac{61}{240}$	$-\frac{1}{20}$	$-\frac{11}{240}$	$\frac{1}{16}$	$-\frac{19}{240}$	$-\frac{1}{60}$	$\frac{11}{60}$	$-\frac{29}{240}$	$\frac{1}{16}$	$-\frac{61}{240}$	$\frac{1}{20}$	$-\frac{11}{240}$	$-\frac{1}{16}$	$-\frac{19}{240}$	$\frac{1}{60}$
$\frac{1}{16}$	$-\frac{9}{80}$	$-\frac{1}{20}$	$\frac{9}{80}$	$-\frac{3}{20}$	$-\frac{1}{80}$	$-\frac{1}{20}$	$\frac{1}{80}$	$-\frac{1}{16}$	$\frac{1}{16}$	$\frac{9}{80}$	$-\frac{1}{20}$	$-\frac{9}{80}$	$-\frac{3}{20}$	$\frac{1}{80}$	$-\frac{1}{20}$	$-\frac{1}{80}$	$-\frac{1}{16}$
$\frac{1}{20}$	$\frac{9}{80}$	$-\frac{1}{16}$	$\frac{1}{80}$	$\frac{3}{20}$	$-\frac{9}{80}$	$\frac{1}{16}$	$-\frac{1}{80}$	$\frac{1}{20}$	$-\frac{1}{20}$	$\frac{9}{80}$	$\frac{1}{16}$	$\frac{1}{80}$	$-\frac{3}{20}$	$-\frac{9}{80}$	$-\frac{1}{16}$	$-\frac{1}{80}$	$-\frac{1}{20}$
$\frac{1}{16}$	$\frac{29}{240}$	$\frac{11}{60}$	$\frac{11}{240}$	$\frac{1}{20}$	$\frac{61}{240}$	$\frac{1}{60}$	$\frac{19}{240}$	$-\frac{1}{16}$	$\frac{1}{16}$	$-\frac{29}{240}$	$\frac{11}{60}$	$-\frac{11}{240}$	$\frac{1}{20}$	$-\frac{61}{240}$	$\frac{1}{60}$	$-\frac{19}{240}$	$-\frac{1}{16}$
$\frac{1}{16}$	$-\frac{19}{240}$	$-\frac{1}{60}$	$-\frac{61}{240}$	$-\frac{1}{20}$	$-\frac{11}{240}$	$-\frac{11}{60}$	$-\frac{29}{240}$	$-\frac{1}{16}$	$\frac{1}{16}$	$\frac{19}{240}$	$-\frac{1}{60}$	$\frac{61}{240}$	$-\frac{1}{20}$	$\frac{11}{240}$	$-\frac{11}{60}$	$\frac{29}{240}$	$-\frac{1}{16}$
$-\frac{1}{20}$	$\frac{1}{80}$	$-\frac{1}{16}$	$\frac{9}{80}$	$-\frac{3}{20}$	$-\frac{1}{80}$	$\frac{1}{16}$	$-\frac{9}{80}$	$-\frac{1}{20}$	$\frac{1}{20}$	$\frac{1}{80}$	$\frac{1}{16}$	$\frac{9}{80}$	$\frac{3}{20}$	$-\frac{1}{80}$	$-\frac{1}{16}$	$-\frac{9}{80}$	$\frac{1}{20}$
$\frac{1}{16}$	$-\frac{1}{80}$	$\frac{1}{20}$	$\frac{1}{80}$	$\frac{3}{20}$	$-\frac{9}{80}$	$\frac{1}{20}$	$\frac{9}{80}$	$-\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{80}$	$\frac{1}{20}$	$-\frac{1}{80}$	$\frac{3}{20}$	$\frac{9}{80}$	$\frac{1}{20}$	$-\frac{9}{80}$	$-\frac{1}{16}$
$\frac{1}{60}$	$\frac{19}{240}$	$-\frac{1}{16}$	$\frac{11}{240}$	$\frac{1}{20}$	$\frac{61}{240}$	$\frac{1}{16}$	$\frac{29}{240}$	$\frac{11}{60}$	$-\frac{1}{60}$	$\frac{19}{240}$	$\frac{1}{16}$	$\frac{11}{240}$	$-\frac{1}{20}$	$\frac{61}{240}$	$-\frac{1}{16}$	$\frac{29}{240}$	$-\frac{11}{60}$
$-\frac{1}{20}$	$\frac{1}{80}$	$-\frac{1}{16}$	$\frac{9}{80}$	$-\frac{3}{20}$	$-\frac{1}{80}$	$-\frac{7}{16}$	$-\frac{9}{80}$	$-\frac{1}{20}$	$\frac{1}{20}$	$\frac{1}{80}$	$\frac{1}{16}$	$\frac{9}{80}$	$\frac{3}{20}$	$-\frac{1}{80}$	$\frac{7}{16}$	$-\frac{9}{80}$	$\frac{1}{20}$
$\frac{1}{16}$	$-\frac{11}{240}$	$\frac{1}{60}$	$-\frac{29}{240}$	$\frac{1}{20}$	$-\frac{19}{240}$	$\frac{11}{60}$	$-\frac{61}{240}$	$-\frac{1}{16}$	$\frac{1}{16}$	$\frac{11}{240}$	$\frac{1}{60}$	$\frac{29}{240}$	$\frac{1}{20}$	$\frac{19}{240}$	$\frac{11}{60}$	$\frac{61}{240}$	$-\frac{1}{16}$
$\frac{1}{60}$	$\frac{11}{240}$	$\frac{1}{16}$	$\frac{19}{240}$	$\frac{1}{20}$	$\frac{29}{240}$	$\frac{1}{16}$	$\frac{61}{240}$	$-\frac{11}{60}$	$\frac{1}{60}$	$\frac{11}{240}$	$\frac{1}{16}$	$\frac{19}{240}$	$\frac{1}{20}$	$\frac{29}{240}$	$-\frac{1}{16}$	$\frac{61}{240}$	$\frac{11}{60}$
1	$-\frac{1}{80}$	$\frac{1}{20}$	$\frac{1}{80}$	$\frac{3}{20}$	$-\frac{9}{80}$	$\frac{1}{20}$	$\frac{9}{80}$	$\frac{7}{16}$	$\frac{1}{16}$	$\frac{1}{80}$	$\frac{1}{20}$	$-\frac{1}{80}$	$\frac{3}{20}$	$\frac{9}{80}$	$\frac{1}{20}$	$-\frac{9}{80}$	$\frac{7}{16}$

FIGURE 18 Least-squares reconstructor matrix for the geometry shown in Fig. 17.

If H is singular, some of the diagonal elements of D will be zero and D^{-1} cannot be defined. However, this method allows us to obtain the closest possible solution in a least-squares sense by zeroing the elements in the diagonal of D^{-1} that come from zero elements in the matrix D . We can arrive at a solution that discards only those equations that generated the problem in the first place. In addition to straightforward SVD, more general techniques have been proposed involving iterative solutions for the phases.^{46, 50, 51}

In addition, several other “tricks” have been developed to alleviate the singularity problem. For instance, piston error is normally ignored, and this contributes to the singularity since there are then an infinite number of solutions that could give the same slope measurements. Adding a row of 1s to the measurement matrix H and setting a corresponding value of 0 in the slope vector for piston allows inversion.⁵²

In the implementation of Shack-Hartmann sensors, a reference wavefront must be provided to calibrate imperfections in the lenslet array and distortions that are introduced by any relay optics to match the pitch of the lenslet array to the pitch of the pixels in the detector array. It is general practice to inject a “perfect” plane wave into the optical train just in front of the lenslet array and to record the positions of all of the Shack-Hartmann spots. During normal operation, the wavefront sensor gradient processor then computes the difference between the spot position of the residual error wavefront and the reference wavefront.

Laser Beacons (Laser Guide Stars)

Adaptive optical systems require a beacon or a source of light to sense turbulence-induced wave distortions. From an anisoplanatic point of view, the ideal beacon is the object being imaged. However, most objects of interest to astronomers are not bright enough to serve as beacons.

It is possible to create artificial beacons (also referred to as synthetic beacons) that are suitable for wavefront sensing with lasers as first demonstrated by Fugate et al.⁵³ and Primmerman et al.⁵⁴ Laser beacons can be created by Rayleigh scattering of focused beams at ranges between 15 and 20 km or by resonant scattering from a layer of sodium atoms in the mesosphere at an altitude of between 90 and 100 km. Examples of Rayleigh beacon AO systems are described by Fugate et al.⁵⁵ for the SOR 1.5-m telescope and by Thompson et al.⁵⁶ for the Mt. Wilson 100-in telescope; examples of sodium beacon AO systems are described by Olivier et al.⁵⁷ for the Lick 3-m Shane telescope and by Butler et al.⁵⁸ for the 3.5-m Calar Alto telescope. Researchers at the W. M. Keck Observatory at Mauna Kea are in the process of installing a sodium dye laser to augment their NGS AO system on Keck II.

The laser beacon concept was first conceived and demonstrated within the U.S. Department of Defense during the early 1980s (see Fugate⁵⁹ for a short summary of the history). The information developed under this program was not declassified until May 1991, but much of the early work was published subsequently.⁶⁰ The laser beacon concept was first published openly by Foy and Labeyrie⁶¹ in 1985 and has been of interest in the astronomy community since.

Even though laser beacons solve a significant problem, they also introduce new problems and have two significant limitations compared with bright NGSs. The new problems include potential light contamination in science cameras and tracking sensors, cost of ownership and operation, and observing complexities that are associated with propagating lasers through the navigable airspace and near-earth orbital space, the home of thousands of space payloads. The technical limitations are that laser beacons provide no information on full-aperture tilt, and that a “cone effect” results from the finite altitude of the beacon, which contributes an additional error called focus (or focal) anisoplanatism. Focus anisoplanatism can be partially alleviated by using a higher-altitude beacon, such as a sodium guide star, as discussed in the following.

Focus Anisoplanatism The mean square wavefront error due to the finite altitude of the laser beacon is given by

$$\sigma_{\text{FA}}^2 = (D/d_0)^{5/3} \quad (49)$$

where d_0 is an effective aperture size corrected by the laser beacon that depends on the height of the laser beacon, the C_n^2 profile, the zenith angle, and the imaging wavelength. This parameter was defined by Fried and Belsher.⁶² Tyler⁶³ developed a method to rapidly evaluate d_0 for arbitrary C_n^2 profiles given by the expression

$$d_0 = \lambda^{6/5} \cos^{3/5}(\psi) \left[\int C_n^2(z) F(z/H) dz \right]^{-3/5} \quad (50)$$

where H is the vertical height of the laser beacon and the function $F(z/H)$ is given by

$$\begin{aligned}
 F(z/H) = & 16.71371210(1.032421640 - 0.8977579487u) \\
 & \times [1 + (1 - z/H)^{5/3}] - 2.168285442 \\
 & \times \left\{ \frac{6}{11} {}_2F_1 \left[-\frac{11}{6}, -\frac{5}{6}; 2; \left(1 - \frac{z}{H}\right)^2 \right] \right. \\
 & \left. - \frac{6}{11} (z/H)^{5/3} - u \frac{10}{11} \left(1 - \frac{z}{H}\right) \right. \\
 & \left. \times {}_2F_1 \left[-\frac{11}{6}, \frac{1}{6}; 3; \left(1 - \frac{z}{H}\right)^2 \right] \right\}
 \end{aligned} \tag{51}$$

for $z < H$ and

$$F(z/H) = 16.71371210(1.032421640 - 0.8977579487u) \tag{52}$$

for $z > H$. In these equations, z is the vertical height above the ground, H is the height of the laser beacon, u is a parameter that is equal to zero when only piston is removed and that is equal to unity when piston and tilt are removed, and ${}_2F_1[a, b; c; z]$ is the hypergeometric function.

Equations (50) and (51) are easily evaluated on a programmable calculator or a personal computer. They are very useful for quickly establishing the expected performance of laser beacons for a particular C_n^2 profile, imaging wavelength, and zenith angle view. Figures 19 and 20 are plots of d_0 representing the best and average seeing at Mauna Kea and a site described by the HV_{57} turbulence profile. Since d_0 scales as $\lambda^{6/5}$, values at $2.2 \mu\text{m}$ are 5.9 times larger, as shown in the plots.

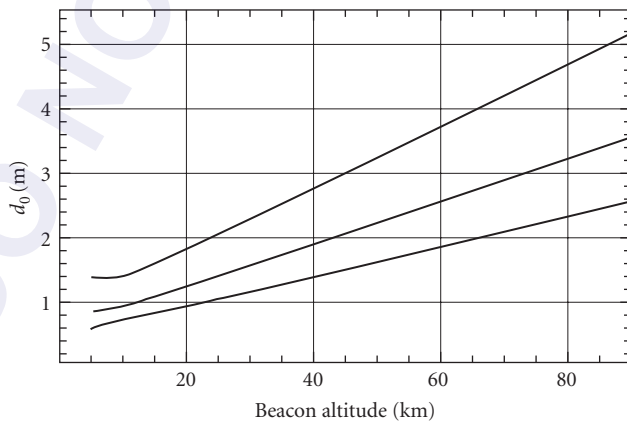


FIGURE 19 Values of d_0 versus laser beacon altitude for zenith imaging at $0.5 \mu\text{m}$ and (top to bottom) best Mauna Kea seeing, average Mauna Kea seeing, and the HV_{57} turbulence profile.

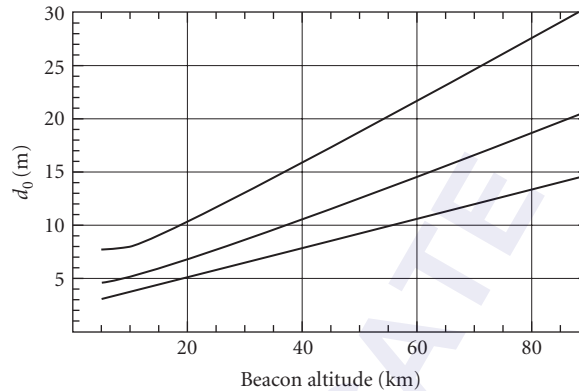


FIGURE 20 Values of d_0 versus laser beacon altitude for zenith imaging at $2.2\ \mu\text{m}$ and (top to bottom) best Mauna Kea seeing, average Mauna Kea seeing, and the $HV_{5/7}$ turbulence profile.

It is straightforward to compute the Strehl ratio due only to focus anisoplanatism using the approximation $SR = e^{-\sigma_{FA}^2}$, where $\sigma_{FA}^2 = (D/d_0)^{5/3}$. There are many possible combinations of aperture sizes, imaging wavelength, and zenith angle, but to illustrate the effect for 3.5- and 10-m apertures, Figs. 21 and 22 show the focus anisoplanatism Strehl ratio as a function of wavelength for 15- and 90-km beacon altitudes and for three seeing conditions. As these plots show, the effectiveness of laser beacons is very sensitive to the aperture diameter, seeing conditions, and beacon altitude. A single Rayleigh beacon at an altitude of 15 km is essentially useless on a 10-m aperture in $HV_{5/7}$ seeing, but it is probably useful at a Mauna Kea—like site. One needs to keep in mind that the curves in Figs. 21 and 22 are the upper limits of performance since other effects will further reduce the Strehl ratio.

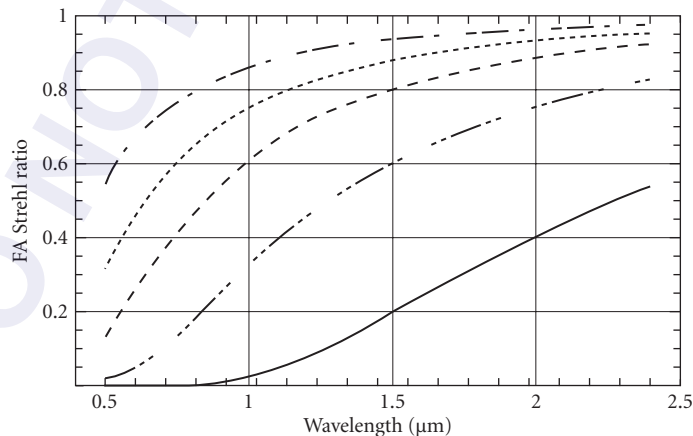


FIGURE 21 The telescope Strehl ratio due to focus anisoplanatism only. Conditions are for a 3.5-m telescope viewing at 30° zenith angle. Curves are (top to bottom): best seeing at Mauna Kea, 90-km beacon; average seeing at Mauna Kea, 90-km beacon; $HV_{5/7}$, 90-km beacon; best seeing at Mauna Kea, 15-km beacon; and $HV_{5/7}$, 15-km beacon.

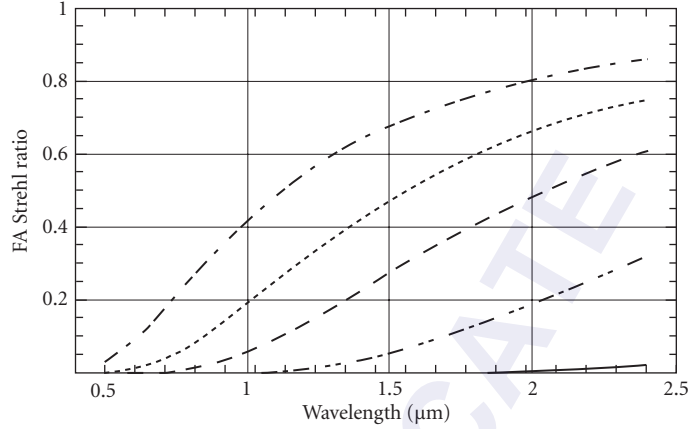


FIGURE 22 The telescope Strehl ratio due to focus anisoplanatism only. Conditions are for a 10.0-m telescope viewing at 30° zenith angle. Curves are (top to bottom): best seeing at Mauna Kea, 90-km beacon; average seeing at Mauna Kea, 90-km beacon; HV_{5/7}, 90-km beacon; best seeing at Mauna Kea, 15-km beacon; and HV_{5/7}, 15-km beacon.

Generation of Rayleigh Laser Beacons For Rayleigh scattering, the number of photo-detected electrons (pdes) per subaperture in the wavefront sensor, N_{pde} , can be computed by using a lidar equation of the form

$$N_{\text{pde}} = \eta_{\text{QE}} T_t T_r T_{\text{atm}}^2 \frac{A_{\text{sub}} \beta_{\text{BS}} \Delta l}{R^2} \frac{E_p \lambda}{hc} \quad (53)$$

where η_{QE} = quantum efficiency of the wavefront sensor

T_t = laser transmitter optical transmission

T_r = optical transmission of the wavefront sensor

T_{atm} = one-way transmission of the atmosphere

A_{sub} = area of a wavefront sensor subaperture

β_{BS} = fraction of incident laser photons backscattered per meter of scattering volume [steradian (sr)⁻¹m⁻¹]

R = range to the midpoint of the scattering volume

Δl = length of the scattering volume—the range gate

E_p = energy per pulse

λ = laser wavelength

h = Planck's constant

c = speed of light

The laser beam is focused at range R , and the wavefront sensor is gated on and off to exclude back-scattered photons from the beam outside a range gate of length Δl centered on range R . The volume-scattering coefficient is proportional to the atmospheric pressure and is inversely proportional to the temperature and the fourth power of the wavelength. Penndorf⁶⁴ developed the details of this relationship, which can be reduced to

$$\beta_{\text{BS}} = 4.26 \times 10^{-7} \frac{P(h)}{T(h)} \text{sr} \cdot \text{m}^{-1} \quad (54)$$

where values of number density, pressure, and temperature at sea level (needed in Penndorf's equations) have been obtained from the U.S. Standard Atmosphere.⁶⁵ At an altitude of 10 km, $\beta_{BS} = 5.1 \times 10^{-7} \text{ sr}^{-1} \text{ m}^{-1}$ and a 1-km-long volume of the laser beam scatters only 0.05 percent of the incident photons per steradian. Increasing the length of the range gate would increase the total signal that is received by the wavefront sensor; however, we should limit the range gate so that subapertures at the edge of the telescope are not able to resolve the projected length of the scattered light. Range gates that are longer than this criterion increase the size of the beacon's image in a subaperture and increase the rms value of the measurement error in each subaperture. A simple geometric analysis leads to an expression for the maximum range gate length:

$$\Delta L = 2 \frac{\lambda R_b^2}{D r_0} \quad (55)$$

where R_b is the range to the center of the beacon and D is the aperture diameter of the telescope. Figure 23 shows the computed signal in detected electrons per subaperture as a function of altitude for a 100-watt (W) average power laser operating at 1000 pulses per second at either 351 nm (upper curve) or 532 nm (lower curve). Other parameters used to compute these curves are listed in the figure caption. When all first-order effects are accounted for, notice that (for high-altitude beacons) there is little benefit to using an ultraviolet wavelength laser over a green laser—even though the scattering goes as λ^{-4} . With modern low-noise detector arrays, signal levels as low as 50 pdes are very usable; above 200 pdes, the sensor generally no longer dominates the performance.

Equations (53) and (55) accurately predict what has been realized in practice. A pulsed copper-vapor laser, having an effective pulse energy of 180 millijoules (mJ) per wavefront sample, produced 190 pdes per subaperture on the Starfire Optical Range 1.5-m telescope at a backscatter range of 10 km, range gate of 2.4 km, and subapertures of 9.2 cm. The total round-trip optical and quantum efficiency was only 0.4 percent.⁵⁵

There are many practical matters on how to design the optical system to project the laser beam. If the beam shares any part of the optical train of the telescope, then it is important to inject the beam

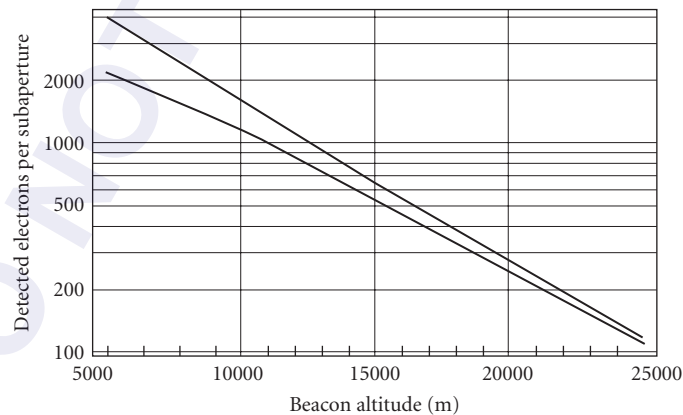


FIGURE 23 Rayleigh laser beacon signal versus altitude of the range gate for two laser wavelengths: 351 nm (top curve) and 532 nm (bottom curve). Assumptions for each curve: range gate length = $2\lambda R_b^2 / (D r_0)$, $D = 1.5$ m, $r_0 = 5$ cm at 532 nm and 3.3 cm at 351 nm, $T_t = 0.30$, $T_r = 0.25$, $\eta_{QE} = 90$ percent at 532 nm and 70 percent at 351 nm, $E_p = 0.10$ J, subaperture size = 10 cm^2 , one-way atmospheric transmission computed as a function of altitude and wavelength.

as close to the output of the telescope as feasible. The best arrangement is to temporally share the aperture with a component that is designed to be out of the beam train when sensors are using the telescope (e.g., a rotating reflecting wheel as described by Thompson et al.⁵⁶). If the laser is injected optically close to wavefront sensors or trackers, there is the additional potential of phosphorescence in mirror substrates and coatings that can emit photons over a continuum of wavelengths longer than the laser fundamental. The decay time of these processes can be longer than the interpulse interval of the laser, presenting a pseudo-continuous background signal that can interfere with faint objects of interest. This problem was fundamental in limiting the magnitude of natural guide star tracking with photon counting avalanche photo diodes at the SOR 1.5-m telescope during operations with the copper-vapor laser.⁶⁶

Generation of Mesospheric Sodium Laser Beacons As the curves in Figs. 21 and 22 show, performance is enhanced considerably when the beacon is at high altitude (90 versus 15 km). The signal from Rayleigh scattering is 50,000 times weaker for a beacon at 90 km compared with one at 15 km. Happer et al.⁶⁷ suggested laser excitation of mesospheric sodium in 1982; however, the generation of a beacon that is suitable for use with AO has turned out to be a very challenging problem. The physics of laser excitation is complex and the development of an optimum laser stresses modern materials science and engineering. This section only addresses how the temporal and spectral format of the laser affects the signal return. The engineering issues of building a laser are beyond the present scope.

The sodium layer in the mesosphere is believed to arise from meteor ablation. The average height of the layer is approximately 95 km above sea level and it is 10 km thick. The column density is only $2\text{--}5 \cdot 10^9$ atoms/cm², or roughly 10^3 atoms/cm³. The temperature of the layer is approximately 200 K, resulting in a Doppler broadened absorption profile having a full width half maximum (FWHM) of about 3 gigahertz (GHz), which is split into two broad resonance peaks that are separated by 1772 MHz, caused by splitting of the $3S_{1/2}$ ground state. The natural lifetime of an excited state is only 16 ns. At high laser intensities (roughly 6 mW/cm^2 for lasers with a natural line width of 10 MHz or 5 W/cm^2 for lasers having spectral content that covers the entire Doppler broadened spectrum), saturation occurs and the return signal does not increase linearly with increasing laser power.

The three types of lasers that have been used to date for generating beacons are the continuous-wave (CW) dye laser, the pulsed-dye laser, and a solid-state sum-frequency laser. Continuous-wave dye lasers are available commercially and generally provide from 2 to 4 W of power. A specialized pulsed-dye lasers installed at Lick Observatory's 3-m Shane Telescope, built at Lawrence Livermore National Laboratory by Friedman et al.,⁶⁸ produces an average power of 20 W consisting of 100- to 150-ns-long pulses at an 11-kHz pulse rate. The sum-frequency laser concept relies on the sum-frequency mixing of the 1.064- and 1.319- μm lines of the Nd:YAG laser in a nonlinear crystal to produce the required 0.589- μm wavelength for the spectroscopic D_2 line. The first experimental devices were built by Jeys at MIT Lincoln Laboratory.^{69, 70} The pulse format is one of an envelope of macropulses, lasting of the order of 100 μs , containing mode-locked pulses at roughly a 100-MHz repetition rate, and a duration of from 400 to 700 ps. The sum-frequency lasers built to date have had average powers of from 8 to 20 W and have been used in field experiments at SOR and Apache Point Observatory.^{71, 72}

A comprehensive study of the physics governing the signal that is generated by laser excitation of mesospheric sodium has been presented by Milonni et al.^{73, 74} for short, intermediate, and long pulses, and CW, corresponding to the lasers described previously. Results are obtained by numerical computations and involve a full-density matrix treatment of the sodium D_2 line. In some specific cases, it has been possible to approximate the results with analytical models that can be used to make rough estimates of signal strengths.

Figure 24 shows results of computations by Milonni et al. for short- and long-pulse formats and an analytical extension of numerical results for a high-power CW laser. Results are presented as the number of photons per square centimeter per millisecond received at the primary mirror of the telescope versus average power of the laser. The curves correspond to (top to bottom) the sum-frequency laser; a CW laser whose total power is spread over six narrow lines that are distributed across the Doppler profile; a CW laser having only one narrow line, and the pulsed-dye laser. The specifics are given in the caption for Fig. 24. Figure 25 extends these results to 200-W average power and shows

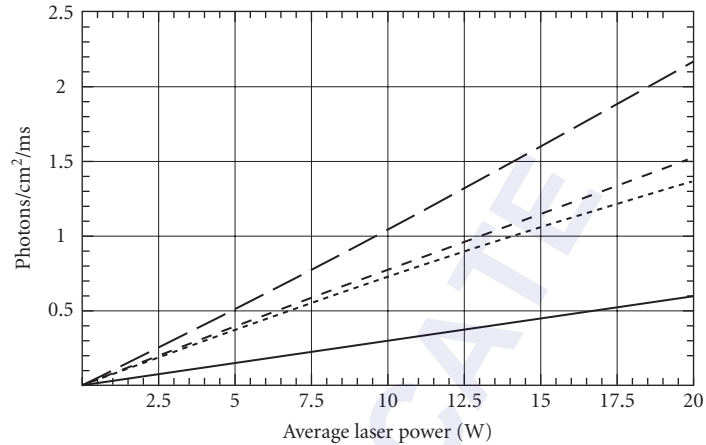


FIGURE 24 Sodium laser beacon signal versus average power of the pump laser. The lasers are (top to bottom) sum-frequency laser with a micropulse-macropulse format (150- μ s macropulses filled with 700-ps micropulses at 100 MHz), a CW laser having six 10-MHz-linewidth lines, a CW laser having a single 10-MHz linewidth, and a pulsed-dye laser having 150-ns pulses at a 30-kHz repetition rate. The sodium layer is assumed to be 90 km from the telescope, the atmospheric transmission is 0.7, and the spot size in the mesosphere is 1.2 arcsec.

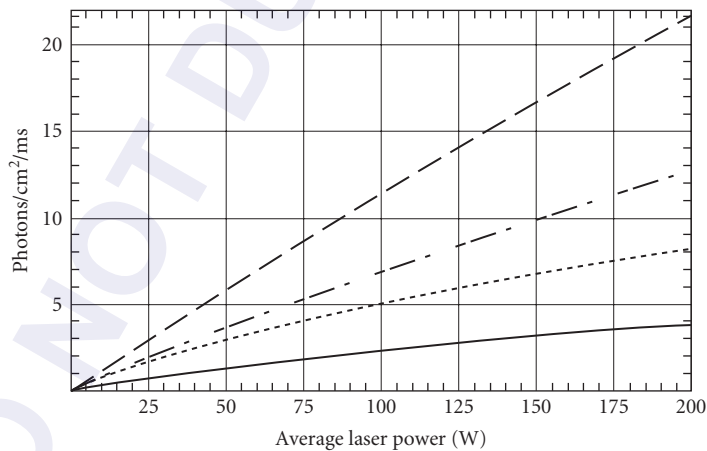


FIGURE 25 Sodium laser beacon signal versus average power of the pump laser. The lasers are (top to bottom) sum-frequency laser with a micropulse-macropulse format (150- μ s macropulses filled with 700-ps micropulses at 100 MHz), a CW laser having six 10-MHz-linewidth lines, a CW laser having a single 10-MHz linewidth, and a pulsed-dye laser having 150-ns pulses at a 30-kHz repetition rate. The sodium layer is assumed to be 90 km from the telescope, the atmospheric transmission is 0.7, and the spot size in the mesosphere is 1.2 arcsec. Saturation of the return signal is very evident for the temporal format of the pulsed-dye laser and for the single-line CW laser. Note, however, that the saturation is mostly eliminated (the line is nearly straight) in the single-line laser by spreading the power over six lines. There appears to be very little saturation in the micropulse-macropulse format.

significant saturation for the pulsed-dye laser format and the single-frequency CW laser as well as some nonlinear behavior for the sum-frequency laser. If the solid-state laser technology community can support high-power, narrow-line CW lasers, this chart says that saturation effects at high power can be ameliorated by spreading the power over different velocity classes of the Doppler profile. These curves can be used to do first-order estimates of signals that are available from laser beacons that are generated in mesospheric sodium with lasers having these temporal formats.

Real-Time Processors

The mathematical process described by Eq. (45) is usually implemented in a dedicated processor that is optimized to perform the matrix multiplication operation. Other wavefront reconstruction approaches that are not implemented by matrix multiplication routines are also possible, but they are not discussed here.

Matrix multiplication lends itself to parallel operations and the ultimate design to maximize speed is to dedicate a processor to each actuator in the deformable mirror. That is, each central processing unit (CPU) is responsible for multiplying and accumulating the sum of each element of a row in the M matrix with each element of the slope vector S . The values of the slope vector should be broadcast to all processors simultaneously to reap the greatest benefit. Data flow and throughput is generally a more difficult problem than raw computing power. It is possible to buy very powerful commercial off-the-shelf processing engines, but getting the data into and out of the engines usually reduces their ultimate performance. A custom design is required to take full advantage of component technology that is available today.

An example of a custom-designed system is the SOR 3.5-m telescope AO system⁷⁵ containing 1024 digital signal processors running at 20 MHz, making a 20-billion-operations-per-second system. This system can perform a $(2048 \times 1024) \times (2048)$ matrix multiply to 16-bit precision (40-bit accumulation), low-pass-filter the data, and provide diagnostic data collection in less than 24 μs . The system throughput exceeds 400 megabytes per second (MB/s). Most astronomy applications have less demanding latency and throughput requirements that can be met with commercial off-the-shelf hardware and software. The importance of latency is illustrated in Fig. 26. This figure shows results

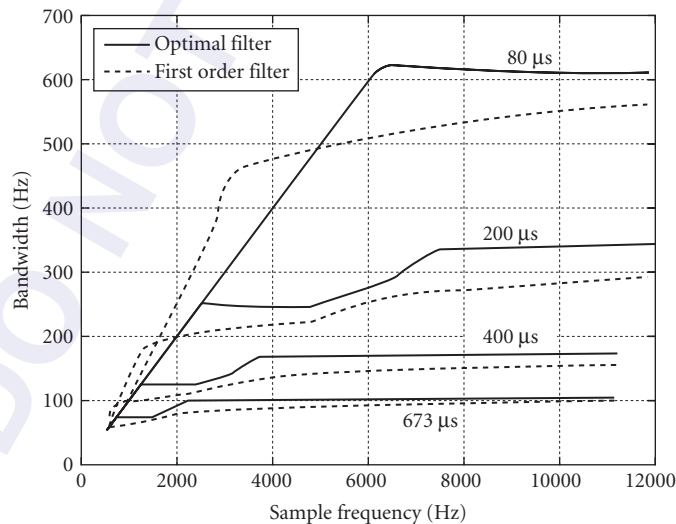


FIGURE 26 Effect of data latency (varying from 80 to 673 μs in this plot) on the control loop bandwidth of the AO system. In this plot, the latency includes only the sensor readout and wavefront processing time.

of analytical predictions showing the relationship between control loop bandwidth of the AO system and the wavefront sensor frame rate for different data latencies.⁷⁶ These curves show that having a high-frame-rate wavefront sensor camera is not a sufficient condition to achieve high control loop bandwidth. The age of the data is of utmost importance. As Fig. 26 shows, the optimum benefit in control bandwidth occurs when the latency is significantly less than a frame time ($\sim 1/2$).

Ensemble-averaged atmospheric turbulence conditions are very dynamic and change on the scale of minutes. Optimum performance of an AO system cannot be achieved if the system is operating on phase estimation algorithms that are based on inaccurate atmospheric information. Optimum performance requires changing the modes of operation in near-real time. One of the first implementations of adaptive control was the system ADONIS, which was implemented on the ESO 3.6-m telescope at La Silla, Chile.⁷⁷ ADONIS employs an artificial intelligence control system that controls which spatial modes are applied to the deformable mirror, depending on the brightness of the AO beacon and the seeing conditions.

A more complex technique has been proposed⁴⁸ that would allow continuous updating of the wavefront estimation algorithm. The concept is to use a recursive least-squares adaptive algorithm to track the temporal and spatial correlations of the distorted wavefronts. The algorithm uses current and recent past information that is available to the servo system to predict the wavefront for a short time in the future and to make the appropriate adjustments to the deformable mirror. A sample scenario has been examined in a detailed simulation, and the system Strehl ratio achieved with the recursive least-squares adaptive algorithm is essentially the same as an optimal reconstructor with *a priori* knowledge of the wind and turbulence profiles. Sample results of this simulation are shown in Fig. 27. The requirements for implementation of this algorithm in a real-time hardware processor have not been worked out in detail. However, it is clear that they are considerable, perhaps greater by an order of magnitude than the requirements for an ordinary AO system.

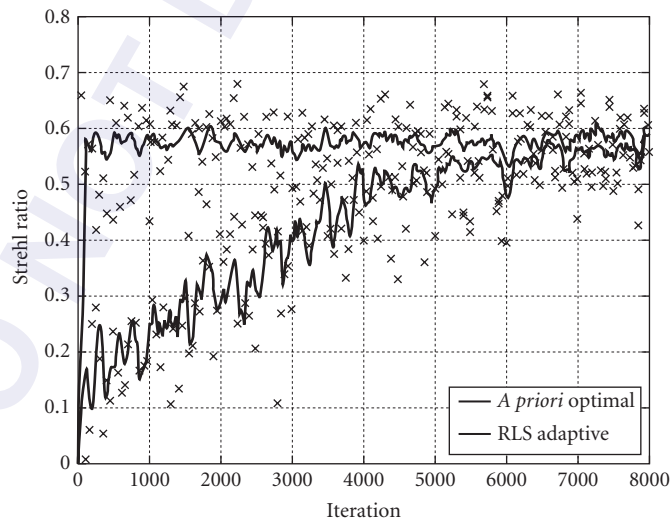


FIGURE 27 Results of a simulation of the Strehl ratio versus time for a recursive least-squares adaptive estimator (lower curve) compared with an optimal estimator having *a priori* knowledge of turbulence conditions (upper curve). The \times 's represent the instantaneous Strehl ratio computed by the simulation, and the lines represent the average values.

Other Higher-Order Wavefront Sensing Techniques

Various forms of shearing interferometers have been successfully used to implement wavefront sensors for atmospheric turbulence compensation.^{78,79} The basic principle is to split the wavefront into two copies, translate one laterally with respect to the other, and then interfere with them. The bright and dark regions of the resulting fringe pattern are proportional to the slope of the wavefront. Furthermore, a lateral shearing interferometer is self-referencing—it does not need a plane wave reference like the Shack-Hartmann sensor. Shearing interferometers are not in widespread use today, but they have been implemented in real systems in the past.⁷⁸

Roddier^{80,81} introduced a new concept for wavefront sensing based on measuring local wavefront curvature (the second derivative of the phase). The concept can be implemented by differencing the irradiance distributions from two locations on either side of the focal plane of a telescope. If the two locations are displaced a distance, l , from the focus of the telescope and the spatial irradiance distribution is given by $I_1(\mathbf{r})$ and $I_2(\mathbf{r})$ at the two locations, then the relationship between the irradiance and the phase is given by

$$\frac{I_1(\mathbf{r}) - I_2(-\mathbf{r})}{I_1(\mathbf{r}) + I_2(-\mathbf{r})} = \frac{\lambda F^2 (F - l)}{2\pi l^2} \left[\frac{\partial \phi}{\partial n}(\mathbf{r}) \delta_c - \nabla^2 \phi(\mathbf{r}) \right] \quad (56)$$

where F is the focal length of the telescope, and the Dirac delta δ_c represents the outward-pointing normal derivative on the edge of the phase pattern. Equation (56) is valid in the geometrical optics approximation. The distance, l , must be chosen such that the validity of the geometrical optics approximation is ensured, requiring that the blur at the position of the defocused pupil image is small compared with the size of the wavefront aberrations desired to be measured. These considerations lead to a condition on l that

$$l \geq \theta_b \frac{F^2}{d} \quad (57)$$

where θ_b is the blur angle of the object that is produced at the positions of the defocused pupil, and d is the size of the subaperture determined by the size of the detector. For point sources and when $d > r_0$, $\theta_b = \lambda/r_0$ and $l \geq \lambda F^2/r_0 d$. For extended sources of angular size $\theta > \lambda/r_0$, l must be chosen such that $l \geq \theta F^2/d$. Since increasing l decreases the sensitivity of the curvature sensor, in normal operations l is set to satisfy the condition $l \geq \lambda F^2/r_0 d$, but once the loop is closed and low-order aberrations are reduced, the sensitivity of the sensor can be increased by making l smaller. To perform wavefront reconstruction, an iterative procedure can be used to solve Poisson's equation. The appeal of this approach is that certain deformable mirrors such as piezoelectric bimorphs deform locally as nearly spherical shapes, and can be driven directly with no intermediate mathematical wavefront reconstruction step. This has not been completely realized in practice, but excellent results have been obtained with two systems deployed at Mauna Kea.⁸² All implementations of these wavefront sensors to date have employed single-element avalanche photodiodes operating in the photon-counting mode for each subaperture. Since these devices remain quite expensive, it may be cost prohibitive to scale the curvature-sensing technique to very high density actuator systems. An additional consideration is that the noise gain goes up linearly with the number of actuators, not logarithmically as with the Shack-Hartmann or shearing interferometer.

The problem of deriving phase from intensity measurements has been studied extensively.^{83–86} A particular implementation of multiple-intensity measurements to derive phase data has become known as phase diversity.⁸⁷ In this approach, one camera is placed at the focus of the telescope and another in a defocused plane with a known amount of defocus. Intensity gathered simultaneously from both cameras can be processed to recover the phase in the pupil of the telescope. The algorithms needed to perform this operation are complex and require iteration.⁸⁸ Such a technique does not presently lend itself to real-time estimation of phase errors in an adaptive optical system, but it could in the future.

Artificial neural networks have been used to estimate phase from intensity as well. The concept is similar to other methods using multiple-intensity measurements. Two focal planes are set up, one at the focus of the telescope and one near focus. The pixels from each focal plane are fed into the nodes of an artificial neural network. The output of the network can be set up to provide almost any desired

information from Zernike decomposition elements to direct-drive signals to actuators of a zonal, deformable mirror. The network must be trained using known distortions. This can be accomplished by using a deterministic wavefront sensor to measure the aberrations and using that information to adjust the weights and coefficients in the neural network processor. This concept has been demonstrated on real telescopes in atmospheric turbulence.⁸⁹ The concept works for low-order distortions, but it appears to have limited usefulness for large, high-density actuator adaptive optical systems.

Wavefront Correctors

There are three major classes of wavefront correctors in use today: segmented mirrors, bimorph mirrors, and stacked-actuator continuous-facesheet mirrors. Figure 28 shows the concept for each of these mirrors.

The segmented mirror can have piston-only or piston-and-tilt actuators. Since the individual segments are completely independent, it is possible for significant phase errors to develop between the segments. It is therefore important that these errors be controlled by real-time interferometry or by strain gauges or other position-measuring devices on the individual actuators. Some large segmented mirrors have been built⁹⁰ for DOD applications and several are in use today for astronomy.^{91,92}

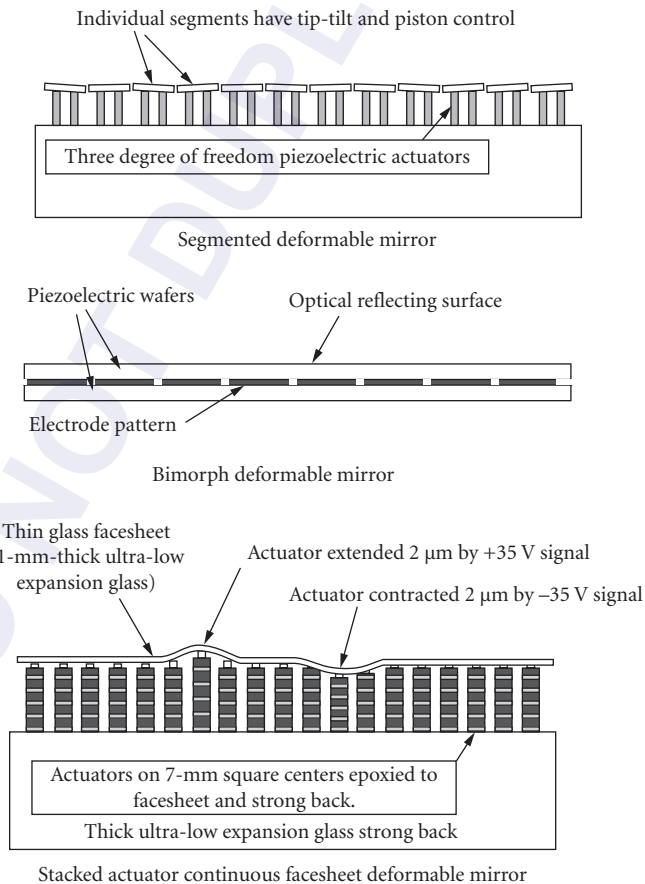


FIGURE 28 Cross sections of three deformable mirror designs.

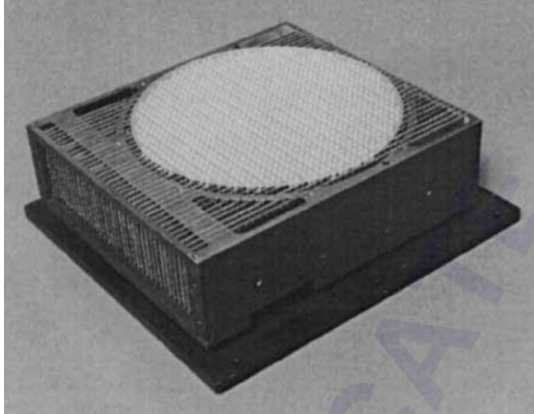


FIGURE 29 The 941-actuator deformable mirror built by Xinetics, Inc., for the SOR 3.5-m telescope.

The stacked actuator deformable mirror is probably the most widely used wavefront corrector. The modern versions are made with lead-magnesium-niobate—a ceramic-like electrostrictive material producing $4\ \mu\text{m}$ of stroke for 70 V of drive.⁹³ The facesheet of these mirrors is typically 1 mm thick. Mirrors with as many as 2200 actuators have been built. These devices are very stiff structures with first resonant frequencies as high as 25 kHz. They are typically built with actuators spaced as closely as 7 mm. One disadvantage with this design is that it becomes essentially impossible to repair individual actuators once the mirror structure is epoxied together. Actuator failures are becoming less likely with today's refined technology, but a very large mirror may have 1000 actuators, which increases its chances for failures over the small mirrors of times past. Figure 29 is a photograph of the 941-actuator deformable mirror in use at the 3.5-m telescope at the SOR.

New Wavefront Corrector Technologies

Our progress toward good performance at visible wavelengths will depend critically on the technology that is available for high-density actuator wavefront correctors. There is promise in the areas of liquid crystals, MEM devices, and nonlinear-optics processes. However, at least for the next few years, it seems that we will have to rely on conventional mirrors with piezoelectric-type actuators or bimorph mirrors made from sandwiched piezoelectric layers. Nevertheless, there is a significant development in the conventional mirror area that has the potential for making mirrors with very large numbers of actuators that are smaller, more reliable, and much less expensive.

5.6 HOW TO DESIGN AN ADAPTIVE OPTICAL SYSTEM

Adaptive optical systems are complex and their performance is governed by many parameters, some controlled by the user and some controlled by nature. The system designer is faced with selecting parameter values to meet performance requirements. Where does he or she begin? One approach is presented here and consists of the following six steps:

1. Determine the average seeing conditions (the mean value of r_0 and f_G) for the site.
2. Determine the most important range of wavelengths of operation.
3. Decide on the minimum Strehl ratio that is acceptable to the users for these seeing conditions and operating wavelengths.
4. Determine the brightness of available beacons.
5. Given the above requirements, determine the optimum values of the subaperture size and the servo bandwidth to minimize the residual wavefront error at the most important wavelength and minimum beacon brightness. (The most difficult parameters to change after the system is built are the wavefront sensor subaperture and deformable mirror actuator geometries. These parameters need to be chosen carefully to address the highest-priority requirements in terms of seeing conditions, beacon brightness, operating wavelength, and required Strehl ratio.) Determine if the associated Strehl ratio is acceptable.
6. Evaluate the Strehl ratio for other values of the imaging wavelength, beacon brightness, and seeing conditions. If these are unsatisfactory, vary the wavefront sensor and track parameters until an acceptable compromise is reached.

Our objective in this section is to develop some practical formulas that can be implemented and evaluated quickly on desktop or laptop computers using programs like Mathematica™ or MatLab™ that will allow one to iterate the six aforementioned steps to investigate the top-level trade space and optimize system performance for the task at hand.

The most important parameters that the designer has some control over are the subaperture size, the wavefront sensor and track sensor integration times, the latency of data in the higher-order and tracker control loops, and the wavelength of operation. One can find optimum values for these parameters since changing them can either increase or decrease the system Strehl ratio. The parameters that the designer has little or no control over are those that are associated with atmospheric turbulence. On the other hand, there are a few parameters that always make performance better the larger (or smaller) we make the parameter: the quantum efficiency of the wavefront and track sensors, the brightness of the beacon(s), the optical system throughput, and read noise of sensors. We can never make a mistake by making the quantum efficiency as large as physics will allow and the read noise as low as physics will allow. We can never have a beacon too bright, nor can we have too much optical transmission (optical attenuators are easy to install). Unfortunately, it is not practical with current technology to optimize the AO system's performance for changing turbulence conditions, for different spectral regions of operation, for different elevation angle, and for variable target brightness by changing the mechanical and optical design from minute to minute. We must be prepared for some compromises based on our initial design choices.

We will use two system examples to illustrate the process that is outlined in the six aforementioned steps: (1) a 3.5-m telescope at an intracontinental site that is used for visible imaging of low-earth-orbiting artificial satellites, and (2) a 10-m telescope that operates at a site of excellent seeing for near-infrared astronomy.

Establish the Requirements

Table 2 lists a set of requirements that we shall try to meet by trading system design parameters. The values in Table 2 were chosen to highlight how requirements can result in significantly different system designs.

In the 3.5-m telescope example, the seeing is bad, the control bandwidths will be high, the imaging wavelength is short, and the required Strehl ratio is significant. Fortunately, the objects (artificial earth satellites) are bright. The 10-m telescope application is much more forgiving with respect to the seeing and imaging wavelengths, but the beacon brightness is four magnitudes fainter and a high Strehl ratio is still required at the imaging wavelength. We now investigate how these requirements determine an optimum choice for the subaperture size.

TABLE 2 Requirements for Two System Examples

Parameter	3.5-m Requirement	10-m Requirement
C_n^2 profile	Modified HV ₅₇	Average Mauna Kea
Wind profile	Slew dominated	Bufton
Average r_0 (0.5 μm)	10 cm	18 cm
Average f_G (0.5 μm)	150 Hz	48 Hz
Imaging wavelength	0.85 μm	1.2–2.2 μm
Elevation angle	45°	45°
Minimum Strehl ratio	0.5	0.5
Beacon brightness	$m_v = 6$	$m_v = 10$

Selecting a Subaperture Size

As was mentioned previously, choosing the subaperture size should be done with care, because once a system is built it is not easily changed. The approach is to develop a mathematical expression for the Strehl ratio as a function of system design parameters and seeing conditions and then to maximize the Strehl ratio by varying the subaperture size for the required operating conditions that are listed in Table 2.

The total-system Strehl ratio is the product of the higher-order and full-aperture tilt Strehl ratios:

$$\text{SR}_{\text{Sys}} = \text{SR}_{\text{HO}} \cdot \text{SR}_{\text{tilt}} \quad (58)$$

The higher-order Strehl ratio can be estimated as

$$\text{SR}_{\text{HO}} = e^{-[E_n^2 + E_f^2 + E_s^2 + E_{\text{FA}}^2]} \quad (59)$$

where E_n^2 , E_f^2 , E_s^2 , and E_{FA}^2 are the mean square phase errors due to wavefront measurement and reconstruction noise, fitting error, servo lag error, and focus anisoplanatism (if a laser beacon is used), respectively. Equation (59) does not properly account for interactions between these effects and could be too conservative in estimating performance. Ultimately, a system design should be evaluated with a detailed computer simulation, which should include wave-optics atmospheric propagations, diffraction in wavefront sensor optics, details of hardware, processing algorithms and time delays, and many other aspects of the system's engineering design. A high-fidelity simulation will properly account for the interaction of all processes and provide a realistic estimate of system performance. For our purposes here, however, we will treat the errors as independent in order to illustrate the processes that are involved in system design and to show when a particular parameter is no longer the dominant contributor to the total error. We will consider tracking effects later [see below Eq. (66)], but for now we need expressions for components of the higher-order errors so that we may determine an optimum subaperture size.

Wavefront Measurement Error, E_n^2 We will consider a Shack-Hartmann sensor for the discussion that follows. The wavefront measurement error contribution, E_n^2 , can be computed from the equation⁷³

$$E_n^2 = \alpha \left[0.09 \ln(n_{\text{sa}}) \sigma_{\theta_{\text{sa}}}^2 d_s^2 \left(\frac{2\pi}{\lambda_{\text{img}}} \right)^2 \right] \quad (60)$$

where α accounts for control loop averaging and is described below, n_{sa} is the number of subapertures in the wavefront sensor, $\sigma_{\theta_{sa}}^2$ is the angular measurement error over a subaperture, d_s is the length of a side of a square subaperture, and λ_{img} is the imaging wavelength. The angular measurement error in a subaperture is proportional to the angular size of the beacon [normally limited by the seeing for unresolved objects but, in some cases, by the beacon itself (e.g., laser beacons or large astronomical targets)] and is inversely proportional to the signal-to-noise ratio in the wavefront sensor. The value of $\sigma_{\theta_{sa}}$ is given by⁹⁴

$$\sigma_{\theta_{sa}} = \pi \left(\frac{\lambda_b}{d_s} \right) \left[\left(\frac{3}{16} \right)^2 + \left(\frac{\theta_b/2}{4\lambda_b/d_s} \right)^2 \right]^{1/2} \left(\frac{1}{n_s} + \frac{4n_e^2}{n_s^2} \right)^{1/2} \quad (61)$$

where λ_b is the center wavelength of the beacon signal, θ_b is the angular size of the beacon, r_0 is the Fried seeing parameter at 0.5 μm at zenith, ψ is the zenith angle of the observing direction, n_s is the number of photodetected electrons per subaperture per sample, and n_e is the read noise in electrons per pixel. We have used the wavelength and zenith scaling laws for r_0 to account for the angular size of the beacon at the observing conditions.

The number of photodetected electrons per subaperture per millisecond, per square meter, per nanometer of spectral width is given by

$$n_s = \frac{(101500)^{\sec\psi}}{(2.51)^{m_v}} d_s^2 T_{ro} n_{QE} t_s \Delta\lambda \quad (62)$$

where m_v is the equivalent visual magnitude of the beacon at zenith, d_s is the subaperture size in meters, T_{ro} is the transmissivity of the telescope and wavefront sensor optics, n_{QE} is the quantum efficiency of the wavefront sensor, t_s is the integration time per frame of the wavefront sensor, and $\Delta\lambda$ is the wavefront sensor spectral bandwidth in nanometers.

The parameter α is a factor that comes from the filtering process in the control loop and can be considered the mean square gain of the loop. Ellerbroek^{73, 95} has derived an expression (discussed by Milonni et al.⁷³ and by Ellerbroek⁹⁵) for α using the Z -transform method and a simplified control system model. His result is $\alpha = g/(2-g)$, where $g = 2\pi f_{3dB} t_s$ is the gain, f_{3dB} is the control bandwidth, and t_s is the sensor integration time. For a typical control system, we sample at a rate that is 10 times the control bandwidth, making $t_s = 1/(10f_{3dB})$, $g = 0.628$, and $\alpha = 0.458$.

Fitting Error, E_f^2 The wavefront sensor has finite-sized subapertures, and the deformable mirror has a finite number of actuators. There is a limit, therefore, to the spatial resolution to which the system can “fit” a distorted wavefront. The mean square phase error due to fitting is proportional to the $-5/6$ th power of the number of actuators and $(D/r_0)^{5/3}$ and is given by

$$E_f^2 = C_f \left(\frac{D}{r_0 \left(\frac{\lambda_{img}}{0.5 \mu\text{m}} \right)^{6/5} (\cos\psi)^{3/5}} \right)^{5/3} n_a^{-5/6} \quad (63)$$

where C_f is the fitting error coefficient that is determined by the influence function of the deformable mirror and has a value of approximately 0.28 for thin, continuous-facesheet mirrors.

Servo Lag, E_s^2 The mean square wavefront error due to a finite control bandwidth has been described earlier by Eq. (37), which is recast here as

$$E_s^2 = \left(\frac{f_G \left(\frac{\lambda_{img}}{0.5 \mu\text{m}} \right)^{-6/5} (\cos\psi)^{-3/5}}{f_{3dB}} \right)^{5/3} \quad (64)$$

where f_G is the Greenwood frequency scaled for the imaging wavelength and a worst-case wind direction for the zenith angle, and $f_{3\text{dB}}$ is the control bandwidth -3-dB error rejection frequency. Barchers⁷⁶ has developed an expression from which one can determine $f_{3\text{dB}}$ for a conventional proportional integral controller, given the sensor sample time (t_s), the additional latency from readout and data-processing time (δ_0), and the design gain margin [G_M (a number between 1 and ∞ that represents the minimum gain that will drive the loop unstable)].⁹⁶ The phase crossover frequency of a proportional integral controller with a fixed latency is $\omega_{\text{cp}} = \pi/(2\delta)$ rad/s, where $\delta = t_s + \delta_0$ is the total latency, which is made up of the sample period and the sensor readout and processing time, δ_0 . The loop gain that achieves the design gain margin is $K = \omega_{\text{cp}}/G_M$ where G_M is the design gain margin. Barchers shows that $f_{3\text{dB}}$ can be found by determining the frequency for which the modulus of the error rejection function is equal to 0.707 or when

$$|S(i\omega_{3\text{dB}})| = \left| \frac{i\omega_{3\text{dB}}}{i\omega_{3\text{dB}} + \frac{\pi/2(t_s + \delta_0)}{G_M} e^{-i(t_s + \delta_0)\omega_{3\text{dB}}}} \right| = 0.707 \quad (65)$$

where $f_{3\text{dB}} = \omega_{3\text{dB}}/2\pi$. This equation can be solved graphically or with dedicated iterative routines that find roots, such as Mathematica. Figure 30 shows the error rejection curves for four combinations of t_s , δ_0 , and G_M , which are detailed in the caption. Note in particular that decreasing the latency δ_0 from 2 ms to 0.5 ms increased $f_{3\text{dB}}$ from 14 to 30 Hz (compare two curves on the left), illustrating the sensitivity to readout and processing time.

Focus Anisoplanatism, E_{FA}^2 If laser beacons are being used, the effects of focus anisoplanatism, which are discussed in Sec. 5.5, must be considered since this error often dominates the wavefront sensor error budget. Focus anisoplanatism error is given by

$$E_{\text{FA}}^2 = \left(\frac{D}{d_0 \left(\frac{\lambda_{\text{img}}}{0.5 \mu\text{m}} \right)^{6/5} (\cos\psi)^{3/5}} \right)^{5/3} \quad (66)$$

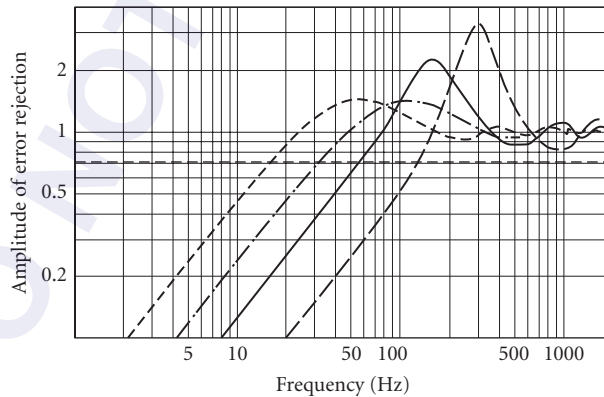


FIGURE 30 Control loop error rejection curves for a proportional integral controller. The curves (left to right) represent the following parameters: dotted ($t_s = 1 \text{ ms}$, $\delta_0 = 2 \text{ ms}$, $G_M = 4$); dash-dot ($t_s = 1 \text{ ms}$, $\delta_0 = 0.5 \text{ ms}$, $G_M = 4$); solid ($t_s = 667 \mu\text{s}$, $\delta_0 = 640 \mu\text{s}$, $G_M = 2$); dashed ($t_s = 400 \mu\text{s}$, $\delta_0 = 360 \mu\text{s}$, $G_M = 1.5$). $f_{3\text{dB}}$ is determined by the intersection of the horizontal dotted line with each of the three curves and has values of 14, 30, 55, and 130 Hz, respectively.

where d_0 is the section size given by Eqs. (50) and (51), scaled for the imaging wavelength and the zenith angle of observation.

Tracking The tracking system is also very important to the performance of an AO system and should not be overlooked as an insignificant problem. In most instances, a system will contain tilt disturbances that are not easily modeled or analyzed arising most commonly from the telescope mount and its movement and base motions coupled into the telescope from the building and its machinery or other seismic disturbances. As described earlier in Eq. (22), the Strehl ratio due to full-aperture tilt variance, σ_θ^2 , is

$$SR_{\text{tilt}} = \frac{1}{1 + \frac{\pi^2}{2} \left(\frac{\sigma_\theta}{\lambda/D} \right)^2} \quad (67)$$

As mentioned previously, the total system Strehl ratio is then the product of these and is given by

$$SR_{\text{sys}} = SR_{\text{HO}} \cdot SR_{\text{tilt}} \quad (68)$$

Results for 3.5-m Telescope AO System

When the preceding equations are evaluated, we can determine the dependence of the system Strehl ratio on d_s for targets of different brightness. Figure 31 shows the Strehl ratio versus subaperture size for four values of target brightness corresponding to (top to bottom) $m_v = 5, 6, 7,$ and 8 and for other system parameters as shown in the figure caption.

Figure 31 shows that we can achieve the required Strehl ratio of 0.5 (see Table 2) for a $m_v = 6$ target by using a subaperture size of about 11 to 12 cm. We could get slightly more performance ($SR = 0.52$) by reducing the subaperture size to 8 cm, but at significant cost since we would have nearly twice as many subapertures and deformable mirror actuators. Notice also that for fainter

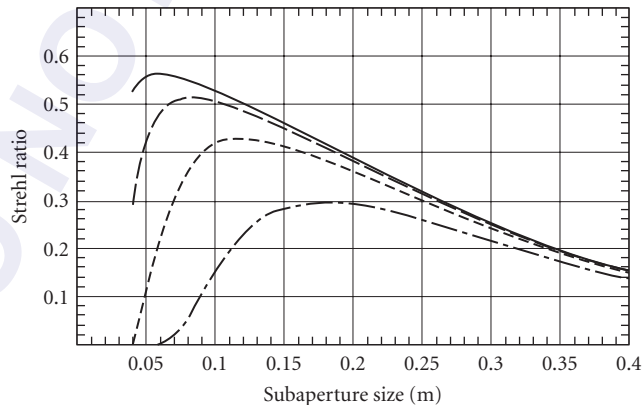


FIGURE 31 System Strehl ratio as a function of subaperture size for the 3.5-m telescope example. Values of m_v are (top to bottom) 5, 6, 7, and 8. Other parameters are: $r_0 = 10$ cm, $f_G = 150$ Hz, $t_s = 400$ μ s, $\delta_0 = 360$ μ s, $G_M = 1.5$, $f_{3\text{dB}} = 127$ Hz, $\lambda_{\text{img}} = 0.85$ μ m, $\psi = 45^\circ$, $T_{r_0} = 0.25$, $\eta_{\text{QE}} = 0.90$, $\Delta\lambda = 400$ nm.

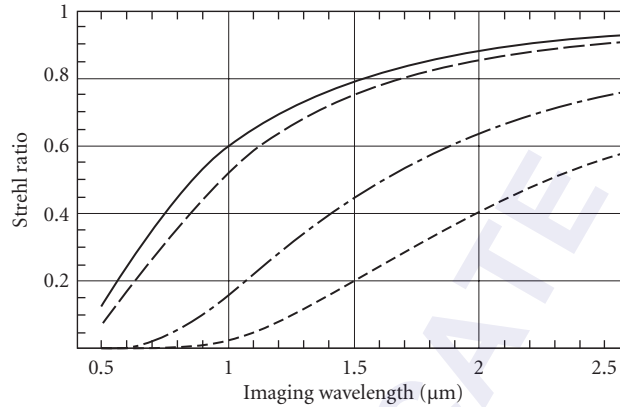


FIGURE 32 System Strehl ratio as a function of imaging wavelength for the 3.5-m telescope example. Curves are for (top to bottom) $m_V = 6$, $d_s = 12$ cm; $m_V = 6$, $d_s = 18$ cm; $m_V = 8$, $d_s = 18$ cm; $m_V = 8$, $d_s = 12$ cm. Other parameters are: $r_0 = 10$ cm, $f_G = 150$ Hz, $t_s = 400$ μ s, $\delta_0 = 360$ μ s, $G_M = 1.5$, $f_{3dB} = 127$ Hz, $\psi = 45^\circ$, $T_{ro} = 0.25$, $\eta_{QE} = 0.90$, $\Delta\lambda = 400$ nm.

objects, a larger subaperture (18 cm) would be optimum for an $m_V = 8$ target, but the SR would be down to 0.4 for the $m_V = 6$ target and would not meet the requirement.

Figure 32 shows the performance of the system for sixth- and eighth-magnitude targets as a function of wavelength. The top two curves in this figure are for 12-cm subapertures (solid curve) and 18-cm subapertures (dashed curve) for the $m_V = 6$ target. Notice that the 12-cm subapertures have better performance at all wavelengths with the biggest difference in the visible. The bottom two curves are for 18-cm subapertures (dash-dot curve) and 12-cm subapertures (dotted curve) for the $m_V = 8$ target. These curves show quantitatively the trade-off between two subaperture sizes and target brightness.

Figure 33 shows how the system will perform for different seeing conditions (values of r_0) and as a function of target brightness. These curves are for a subaperture choice of 12 cm. The curves in

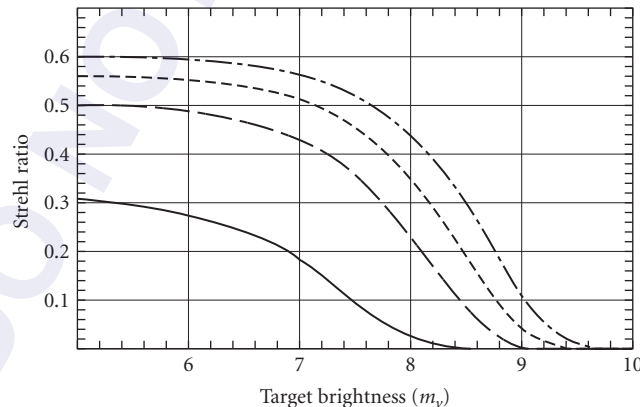


FIGURE 33 System Strehl ratio as a function of target brightness for the 3.5-m telescope example. Curves are for values of r_0 of (top to bottom) 25, 15, 10, and 5 cm. Other parameters are: $d_s = 12$ cm, $f_G = 150$ Hz, $t_s = 400$ μ s, $\delta_0 = 360$ μ s, $G_M = 1.5$, $f_{3dB} = 127$ Hz, $\lambda_{img} = 0.85$ μ m, $\psi = 45^\circ$, $T_{ro} = 0.25$, $\eta_{QE} = 0.90$, $\Delta\lambda = 400$ nm.

Figs. 31 through 33 give designers a good feel of what to expect and how to do top-level system trades for those conditions and parameters for which only they and the users can set the priorities. These curves are easily and quickly generated with modern mathematics packages.

Results for the 10-m AO Telescope System

Similar design considerations for the 10-m telescope example lead to Figs. 34 through 36. Figure 34 shows the Strehl ratio for an imaging wavelength of $1.2\ \mu\text{m}$ versus the subaperture size for four

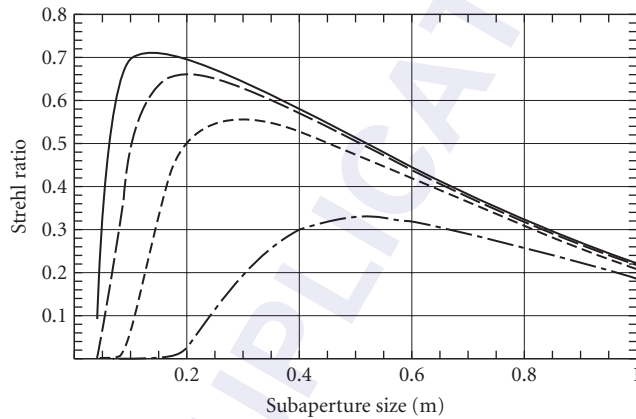


FIGURE 34 System Strehl ratio as a function of subaperture size for the 10-m telescope example. Values of m_v are (top of bottom) 8, 9, 10, and 11. Other parameters are: $r_0 = 18\ \text{cm}$, $f_G = 48\ \text{Hz}$, $t_s = 1000\ \mu\text{s}$, $\delta_0 = 1000\ \mu\text{s}$, $G_M = 2$, $f_{\text{3dB}} = 39\ \text{Hz}$, $\lambda_{\text{img}} = 1.2\ \mu\text{m}$, $\psi = 45^\circ$, $T_{\text{ro}} = 0.25$, $\eta_{\text{QE}} = 0.90$, $\Delta\lambda = 400\ \text{nm}$.

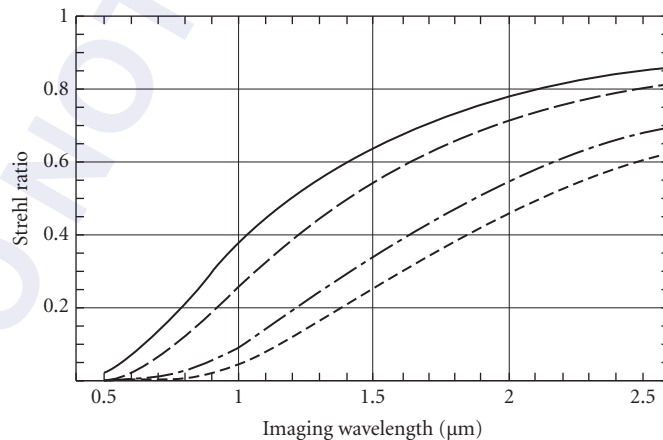


FIGURE 35 System Strehl ratio as a function of imaging wavelength for the 10-cm telescope example. Curves are for (top to bottom) $m_v = 10$, $d_s = 45\ \text{cm}$; $m_v = 10$, $d_s = 65\ \text{cm}$; $m_v = 12$, $d_s = 65\ \text{cm}$; $m_v = 12$, $d_s = 45\ \text{cm}$. Other parameters are: $r_0 = 18\ \text{cm}$, $f_G = 48\ \text{Hz}$, $t_s = 1000\ \mu\text{s}$, $\delta_0 = 1000\ \mu\text{s}$, $G_M = 2$, $f_{\text{3dB}} = 39\ \text{Hz}$, $\psi = 45^\circ$, $T_{\text{ro}} = 0.25$, $\eta_{\text{QE}} = 0.90$, $\Delta\lambda = 400\ \text{nm}$.

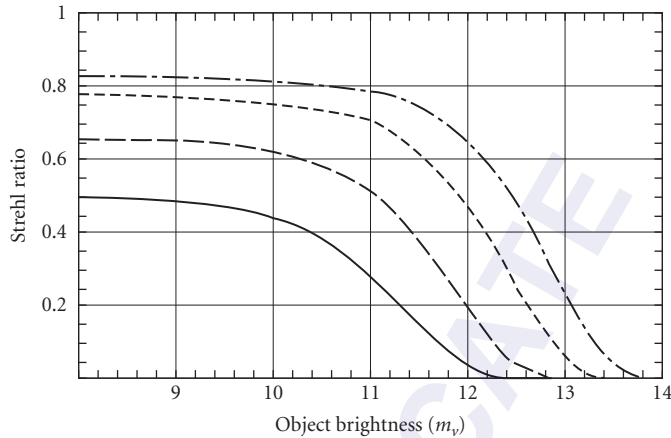


FIGURE 36 System Strehl ratio as a function of target brightness for the 10-m telescope example. Curves are for values of r_0 of (top to bottom) 40, 25, 15, and 10 cm. Other parameters are: $d_s = 45$ cm, $f_G = 48$ Hz, $t_s = 1000$ μ s, $\delta_0 = 1000$ μ s, $G_M = 2$, $f_{3dB} = 39$ Hz, $\lambda_{img} = 1.6$ μ m, $\psi = 45^\circ$, $T_{ro} = 0.25$, $\eta_{QE} = 0.90$, $\Delta\lambda = 400$ nm.

values of beacon brightness corresponding to (top to bottom) $m_v = 8, 9, 10$, and 11 . Other system parameters are listed in the figure caption. These curves show that we can achieve the required Strehl ratio of 0.5 for the $m_v = 10$ beacon with subapertures in the range of 20 to 45 cm. Optimum performance at 1.2 μ m produces a Strehl ratio of 0.56 , with a subaperture size of 30 cm. Nearly 400 actuators are needed in the deformable mirror for 45 -cm subapertures, and nearly 900 actuators are needed for 30 -cm subapertures. Furthermore, Fig. 34 shows that 45 cm is a better choice than 30 cm in the sense that it provides near-optimal performance for the $m_v = 11$ beacon, whereas the Strehl ratio is down to 0.2 for the 30 -cm subapertures.

Figure 35 shows system performance as a function of imaging wavelength. The top two curves are for the $m_v = 10$ beacon, with $d_s = 45$ and 65 cm, respectively. Note that the 45 -cm subaperture gives excellent performance with a Strehl ratio of 0.8 at 2.2 μ m (the upper end of the required spectral range) and a very useful Strehl ratio of 0.3 at 0.8 μ m. A choice of $d_s = 65$ cm provides poorer performance for $m_v = 10$ (in fact, it does not satisfy the requirement) than does $d_s = 45$ cm due to fitting error, whereas a choice of $d_s = 65$ cm provides better performance for $m_v = 12$ due to improved wavefront sensor signal-to-noise ratio.

Figure 36 shows system performance in different seeing conditions as a function of beacon brightness for an intermediate imaging wavelength of 1.6 μ m. These curves are for a subaperture size of 45 cm. This chart predicts that in exceptional seeing, a Strehl ratio of 0.5 can be achieved using an $m_v = 12.5$ beacon.

As in the 3.5 -m telescope case, these curves and others like them with different parameters can be useful in performing top-level design trades and in selecting a short list of design candidates for further detailed analysis and simulation.

5.7 ACKNOWLEDGMENTS

I would like to thank David L. Fried, Earl Spillar, Jeff Barchers, John Anderson, Bill Lowrey, and Greg Peisert for reading the manuscript and making constructive suggestions that improved the content and style of this chapter. I would especially like to acknowledge the efforts of my editor, Bill Wolfe, who relentlessly kept me on course from the first pitiful draft.

5.8 REFERENCES

1. D. R. Williams, J. Liang, and D. T. Miller, "Adaptive Optics for the Human Eye," *OSA Technical Digest 1996* **13**:145–147 (1996).
2. J. Liang, D. Williams, and D. Miller, "Supernormal Vision and High-Resolution Retinal Imaging through Adaptive Optics," *JOSA A* **14**:2884–2892 (1997).
3. A. Roorda and D. Williams, "The Arrangement of the Three Cone Classes in the Living Human Eye," *Nature* **397**:520–522 (1999).
4. K. Bar, B. Freisleben, C. Kozlik, and R. Schmiedl, "Adaptive Optics for Industrial CO₂-Laser Systems," *Lasers in Engineering*, vol. 4, no. 3, unknown publisher, 1961.
5. M. Huonker, G. Waibel, A. Giesen, and H. Hugel, "Fast and Compact Adaptive Mirror for Laser Materials Processing," *Proc. SPIE* **3097**:310–319 (1997).
6. R. Q. Fugate, "Laser Beacon Adaptive Optics for Power Beaming Applications," *Proc SPIE* **2121**:68–76 (1994).
7. H. E. Bennett, J. D. G. Rather, and E. E. Montgomery, "Free-Electron Laser Power Beaming to Satellites at China Lake, California," *Proc SPIE* **2121**:182–202 (1994).
8. C. R. Phipps, G. Albrecht, H. Friedman, D. Gavel, E. V. George, J. Murray, C. Ho, W. Priedhorsky, M. M. Michaels, and J. P. Reilly, "ORION: Clearing Near-Earth Space Debris Using a 20-kW, 530-nm, Earth-Based Repetitively Pulsed Laser," *Laser and Particle Beams* **14**:1–44 (1996).
9. K. Wilson, J. Lesh, K. Araki, and Y. Arimoto, "Overview of the Ground to Orbit Lasercom Demonstration," *Space Communications* **15**:89–95 (1998).
10. R. Szeto and R. Butts, "Atmospheric Characterization in the Presence of Strong Additive Measurement Noise," *JOSA A* **15**:1698–1707 (1998).
11. H. Babcock, "The Possibility of Compensating Astronomical Seeing," *Publications of the Astronomical Society of the Pacific* **65**:229–236 (October 1953).
12. D. Fried, "Optical Resolution through a Randomly Inhomogeneous Medium for Very Long and Very Short Exposures," *J. Opt. Soc. Am.* **56**:1372–1379 (October 1966).
13. R. P. Angel, "Development of a Deformable Secondary Mirror," *Proceedings of the SPIE* **1400**:341–351 (1997).
14. S. F. Clifford, "The Classical Theory of Wave Propagation in a Turbulent Medium," *Laser Beam Propagation in the Atmosphere*, Springer-Verlag, New York, 1978.
15. A. N. Kolmogorov, "The Local Structure of Turbulence in Incompressible Viscous Fluids for Very Large Reynolds' Numbers," *Turbulence, Classic Papers on Statistical Theory*, Wiley-Interscience, New York, 1961.
16. V. I. Tatarski, *Wave Propagation in a Turbulent Medium*, McGraw-Hill, New York, 1961.
17. D. Fried, "Statistics of a Geometrical Representation of Wavefront Distortion," *J. Opt. Soc. Am.* **55**:1427–1435 (November 1965).
18. J. W. Goodman, *Statistical Optics*, Wiley-Interscience, New York, 1985.
19. F. D. Eaton, W. A. Peterson, J. R. Hines, K. R. Peterman, R. E. Good, R. R. Beland, and J. H. Brown, "Comparisons of vhf Radar, Optical, and Temperature Fluctuation Measurements of C_n^2 , r_0 , and θ_0 ," *Theoretical and Applied Climatology* **39**:17–29 (1988).
20. D. L. Walters and L. Bradford, "Measurements of r_0 and θ_0 : 2 Decades and 18 Sites," *Applied Optics* **36**:7876–7886 (1997).
21. R. E. Hufnagel, *Proc. Topical Mtg. on Optical Propagation through Turbulence, Boulder, CO* **1** (1974).
22. G. Valley, "Isoplanatic Degradation of Tilt Correction and Short-Term Imaging Systems," *Applied Optics* **19**:574–577 (February 1980).
23. D. Winker, "Unpublished Air Force Weapons Lab Memo, 1986," U.S. Air Force, 1986.
24. E. Marchetti and D. Bonaccini, "Does the Outer Scale Help Adaptive Optics or Is Kolmogorov Gentler?" *Proc. SPIE* **3353**:1100–1108 (1998).
25. R. E. Hufnagel and N. R. Stanley, "Modulation Transfer Function Associated with Image Transmission through Turbulent Media," *J. Opt. Soc. Am.* **54**:52–61 (January 1964).
26. D. Fried, "Limiting Resolution Looking down through the Atmosphere," *J. Opt. Soc. Am.* **56**:1380–1384 (October 1966).

27. R. Noll, "Zernike Polynomials and Atmospheric Turbulence," *J. Opt. Soc. Am.* **66**:207–211 (March 1976).
28. J. Christou, "Deconvolution of Adaptive Optics Images," *Proceedings, ESO/OSA Topical Meeting on Astronomy with Adaptive Optics: Present Results and Future Programs* **56**:99–108 (1998).
29. G. A. Tyler, *Reduction in Antenna Gain Due to Random Jitter*, The Optical Sciences Company, Anaheim, CA, 1983.
30. H. T. Yura and M. T. Tavis, "Centroid Anisoplanatism," *JOSA A* **2**:765–773 (1985).
31. D. P. Greenwood and D. L. Fried, "Power Spectra Requirements for Wave-Front-Compensative Systems," *J. Opt. Soc. Am.* **66**:193–206 (March 1976).
32. G. A. Tyler, "Bandwidth Considerations for Tracking through Turbulence," *JOSA A* **11**:358–367 (1994).
33. R. J. Sasiela, *Electromagnetic Wave Propagation in Turbulence*, Springer-Verlag, New York, 1994.
34. J. Bufton, "Comparison of Vertical Profile Turbulence Structure with Stellar Observations," *Appl. Opt.* **12**:1785 (1973).
35. D. Greenwood, "Tracking Turbulence-Induced Tilt Errors with Shared and Adjacent Apertures," *J. Opt. Soc. Am.* **67**:282–290 (March 1977).
36. G. A. Tyler, "Turbulence-Induced Adaptive-Optics Performance Degradation: Evaluation in the Time Domain," *JOSA A* **1**:358 (1984).
37. D. Fried, "Anisoplanatism in Adaptive Optics," *J. Opt. Soc. Am.* **72**:52–61 (January 1982).
38. J. W. Goodman, *Introduction to Fourier Optics*, McGraw-Hill, San Francisco, 1968.
39. J. Ge, "Adaptive Optics," *OSA Technical Digest Series* **13**:122 (1996).
40. D. W. Tyler and B. L. Ellerbroek, "Sky Coverage Calculations for Spectrometer Slit Power Coupling with Adaptive Optics Compensation," *Proc. SPIE* **3353**:201–209 (1998).
41. J. Ge, R. Angel, D. Sandler, C. Shelton, D. McCarthy, and J. Burge, "Adaptive Optics Spectroscopy: Preliminary Theoretical Results," *Proc. SPIE* **3126**:343–354 (1997).
42. G. Tyler, "Rapid Evaluation of d_0 ," Tech. Rep. TR-1159, The Optical Sciences Company, Placentia, CA, 1991.
43. R. Racine and R. McClure, "An Image Stabilization Experiment at the Canada-France-Hawaii Telescope," *Publications of the Astronomical Society of the Pacific* **101**:731–736 (August 1989).
44. R. Q. Fugate, J. F. Riker, J. T. Roark, S. Stogsdill, and B. D. O'Neil, "Laser Beacon Compensated Images of Saturn Using a High-Speed Near-Infrared Correlation Tracker," *Proc. Top. Mtg. on Adaptive Optics, ESO Conf. and Workshop Proc.* **56**:287 (1996).
45. E. Wallner, "Optimal Wave-Front Correction Using Slope Measurements," *J. Opt. Soc. Am.* **73**:1771–1776 (December 1983).
46. D. L. Fried, "Least-Squares Fitting a Wave-Front Distortion Estimate to an Array of Phase-Difference Measurements," *J. Opt. Soc. Am.* **67**:370–375 (1977).
47. W. J. Wild, "Innovative Wavefront Estimators for Zonal Adaptive Optics Systems, ii," *Proc. SPIE* **3353**:1164–1173 (1998).
48. B. L. Ellerbroek and T. A. Rhoadarmer, "Real-Time Adaptive Optimization of Wave-Front Reconstruction Algorithms for Closed Loop Adaptive-Optical Systems," *Proc. SPIE* **3353**:1174–1185 (1998).
49. R. B. Shack and B. C. Platt, "Production and Use of a Lenticular Hartman Screen," *J. Opt. Soc. Am.* **61**:656 (1971).
50. B. R. Hunt, "Matrix Formulation of the Reconstruction of Phase Values from Phase Differences," *J. Opt. Soc. Am.* **69**:393 (1979).
51. R. H. Hudgin, "Optimal Wave-Front Estimation," *J. Opt. Soc. Am.* **67**:378–382 (1977).
52. R. J. Sasiela and J. G. Mooney, "An Optical Phase Reconstructor Based on Using a Multiplier-Accumulator Approach," *Proc. SPIE* **551**:170 (1985).
53. R. Q. Fugate, D. L. Fried, G. A. Ameer, B. R. Boeke, S. L. Browne, P. H. Roberts, R. E. Ruane, G. A. Tyler, and L. M. Wopat, "Measurement of Atmospheric Wavefront Distortion Using Scattered Light from a Laser Guide-Star," *Nature* **353**:144–146 (September 1991).
54. C. A. Primmerman, D. V. Murphy, D. A. Page, B. G. Zollars, and H. T. Barclay, "Compensation of Atmospheric Optical Distortion Using a Synthetic Beacon," *Nature* **353**:140–141 (1991).
55. R. Q. Fugate, B. L. Ellerbroek, C. H. Higgins, M. P. Jelonek, W. J. Lange, A. C. Slavin, W. J. Wild, D. M. Winker, J. M. Wynia, J. M. Spinhirne, B. R. Boeke, R. E. Ruane, J. F. Moroney, M. D. Olikier,

- D. W. Swindle, and R. A. Cleis, "Two Generations of Laser-Guide-Star Adaptive Optics Experiments at the Starfire Optical Range," *JOSA A* **11**:310–324 (1994).
56. L. A. Thompson, R. M. Castle, S. W. Teare, P. R. McCullough, and S. Crawford, "Unisis: A Laser Guided Adaptive Optics System for the Mt. Wilson 2.5-m Telescope," *Proc. SPIE* **3353**:282–289 (1998).
 57. S. S. Olivier, D. T. Gavel, H. W. Friedman, C. E. Max, J. R. An, K. Avicola, B. J. Bauman, J. M. Brase, E. W. Campbell, C. Carrano, J. B. Cooke, G. J. Freeze, E. L. Gates, V. K. Kanz, T. C. Kuklo, B. A. Macintosh, M. J. Newman, E. L. Pierce, K. E. Waltjen, and J. A. Watson, "Improved Performance of the Laser Guide Star Adaptive Optics System at Lick Observatory," *Proc. SPIE* **3762**:2–7 (1999).
 58. D. J. Butler, R. I. Davies, H. Fewes, W. Hackenburg, S. Rabien, T. Ott, A. Eckart, and M. Kasper, "Calar Alto Affa and the Sodium Laser Guide Star in Astronomy," *Proc. SPIE* **3762**:184–193 (1999).
 59. R. Q. Fugate, "Laser Guide Star Adaptive Optics for Compensated Imaging," *The Infrared and Electro-Optical Systems Handbook*, S. R. Robinson, (ed.), vol. 8, 1993.
 60. See the special edition of *JOSA A* on Atmospheric-Compensation Technology (January-February 1994).
 61. R. Foy and A. Labeyrie, "Feasibility of Adaptive Telescope with Laser Probe," *Astronomy and Astrophysics* **152**: L29–L31 (1985).
 62. D. L. Fried and J. F. Belsher, "Analysis of Fundamental Limits to Artificial-Guide-Star Adaptive-Optics-System Performance for Astronomical Imaging," *JOSA A* **11**:277–287 (1994).
 63. G. A. Tyler, "Rapid Evaluation of d_0 : The Effective Diameter of a Laser-Guide-Star Adaptive-Optics System," *JOSA A* **11**:325–338 (1994).
 64. R. Penndorf, "Tables of the Refractive Index for Standard Air and the Rayleigh Scattering Coefficient for the Spectral Region Between 0.2 and 20.0 μm and Their Application to Atmospheric Optics," *J. Opt. Soc. Am.* **47**:176–182 (1957).
 65. *U.S. Standard Atmosphere*, National Oceanic and Atmospheric Administration, Washington, D.C., 1976.
 66. R. Q. Fugate, "Observations of Faint Objects with Laser Beacon Adaptive Optics," *Proceedings of the SPIE* **2201**:10–21 (1994).
 67. W. Happer, G. J. MacDonald, C. E. Max, and F. J. Dyson, "Atmospheric Turbulence Compensation by Resonant Optical Backscattering from the Sodium Layer in the Upper Atmosphere," *JOSA A* **11**:263–276 (1994).
 68. H. Friedman, G. Erbert, T. Kuklo, T. Salmon, D. Smauley, G. Thompson, J. Malik, N. Wong, K. Kanz, and K. Neeb, "Sodium Beacon Laser System for the Lick Observatory," *Proceedings of the SPIE* **2534**:150–160 (1995).
 69. T. H. Jeys, "Development of a Mesospheric Sodium Laser Beacon for Atmospheric Adaptive Optics," *The Lincoln Laboratory Journal* **4**:133–150 (1991).
 70. T. H. Jeys, A. A. Brailove, and A. Mooradian, "Sum Frequency Generation of Sodium Resonance Radiation," *Applied Optics* **28**:2588–2591 (1991).
 71. M. P. Jelonek, R. Q. Fugate, W. J. Lange, A. C. Slavin, R. E. Ruane, and R. A. Cleis, "Characterization of artificial guide stars generated in the mesospheric sodium layer with a sum-frequency laser," *JOSA A* **11**:806–812 (1994).
 72. E. J. Kibblewhite, R. Vuilleumier, B. Carter, W. J. Wild, and T. H. Jeys, "Implementation of CW and Pulsed Laser Beacons for Astronomical Adaptive Optics," *Proceedings of the SPIE* **2201**:272–283 (1994).
 73. P. W. Milonni, R. Q. Fugate, and J. M. Telle, "Analysis of Measured Photon Returns from Sodium Beacons," *JOSA A* **15**:217–233 (1998).
 74. P. W. Milonni, H. Fern, J. M. Telle, and R. Q. Fugate, "Theory of Continuous-Wave Excitation of the Sodium Beacon," *JOSA A* **16**:2555–2566 (1999).
 75. R. J. Eager, "Application of a Massively Parallel DSP System Architecture to Perform Wavefront Reconstruction for a 941 Channel Adaptive Optics System," *Proceedings of the ICSPAT* **2**:1499–1503 (1977).
 76. J. Barchers, Air Force Research Laboratory/DES, Starfire Optical Range, Kirtland AFB, NM, private communication, 1999.
 77. E. Gendron and P. Lena, "Astronomical Adaptive Optics in Modal Control Optimization," *Astron. Astrophys.* **291**:337–347 (1994).
 78. J. W. Hardy, J. E. Lefebvre, and C. L. Koliopoulos, "Real-Time Atmospheric Compensation," *J. Opt. Soc. Am.* **67**:360–367 (1977); and J. W. Hardy, *Adaptive Optics for Astronomical Telescopes*, Oxford University Press, Oxford, 1998.

79. J. Wyant, "Use of an AC Heterodyne Lateral Shear Interferometer with Real-Time Wavefront Correction Systems," *Applied Optics* **14**:2622–2626 (November 1975).
80. F. Roddier, "Curvature Sensing and Compensation: A New Concept in Adaptive Optics," *Applied Optics* **27**:1223–1225 (April 1988).
81. F. Roddier, *Adaptive Optics in Astronomy*, Cambridge University Press, Cambridge, England, 1999.
82. F. Roddier and F. Rigault, "The VH-CFHT Systems," *Adaptive Optics in Astronomy*, ch. 9, F. Roddier, (ed.), Cambridge University Press, Cambridge, England, 1999.
83. J. R. Fienup, "Phase Retrieval Algorithms: A Comparison," *Appl. Opt.* **21**:2758 (1982).
84. R. A. Gonsalves, "Fundamentals of wavefront sensing by phase retrieval," *Proc. SPIE* **351**, p. 56, 1982.
85. J. T. Foley and M. A. A. Jalil, "Role of Diffraction in Phase Retrieval from Intensity Measurements," *Proc. SPIE* **351**:80 (1982).
86. S. R. Robinson, "On the Problem of Phase from Intensity Measurements," *J. Opt. Soc. Am.* **68**:87 (1978).
87. R. G. Paxman, T. J. Schultz, and J. R. Fienup, "Joint Estimation of Object and Aberrations by Using Phase Diversity," *JOSA A* **9**:1072–1085 (1992).
88. R. L. Kendrick, D. S. Acton, and A. L. Duncan, "Phase Diversity Wave-Front Sensor for Imaging Systems," *Appl. Opt.* **33**:6533–6546 (1994).
89. D. G. Sandler, T. Barrett, D. Palmer, R. Fugate, and W. Wild, "Use of a Neural Network to Control an Adaptive Optics System for an Astronomical Telescope," *Nature* **351**:300–302 (May 1991).
90. B. Hulburd and D. Sandler, "Segmented Mirrors for Atmospheric Compensation," *Optical Engineering* **29**:1186–1190 (1990).
91. D. S. Acton, "Status of the Lockheed 19-Segment Solar Adaptive Optics System," *Real Time and Post Facto Solar Image Correction*, Proc. Thirteenth National Solar Observatory, Sacramento Peak, Summer Shop Series 13, 1992.
92. D. F. Busher, A. P. Doel, N. Andrews, C. Dunlop, P. W. Morris, and R. M. Myers, "Novel Adaptive Optics with the Durham University Electra System," *Adaptive Optics*, Proc. OSA/ESO Conference Tech Digest, Series 23, 1995.
93. M. A. Ealey and P. A. Davis, "Standard Select Electrostrictive PMN Actuators for Active and Adaptive Components," *Optical Engineering* **29**:1373–1382 (1990).
94. G. A. Tyler and D. L. Fried, "Image-Position Error Associated with a Quadrant Detector," *J. Opt. Soc. Am.* **72**:804–808 (1982).
95. B. L. Ellerbroek, Gemini Telescopes Project, Hilo, Hawaii, private communication, 1998.
96. C. L. Phillips and H. T. Nagle, *Digital Control System: Analysis and Design*, Prentice Hall, Upper Saddle River, NJ, 1990.

PART

3

MODULATORS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

I-Cheng Chang

*Accord Optics
Sunnyvale, California*

6.1 GLOSSARY

$\delta\theta_o, \delta\theta_a$	divergence: optical, acoustic
ΔB_m	impermeability tensor
$\Delta f, \Delta F$	bandwidth, normalized bandwidth
Δn	birefringence
$\Delta\theta$	deflection angle
λ_o, λ	optical wavelength (in vacuum/medium)
Λ	acoustic wavelength
ρ	density
τ	acoustic transit time
ψ	phase mismatch function
A	optical to acoustic divergence ratio
a	optical to acoustic wavelength ratio
D	optical aperture
E_i, E_d	electric field, incident, diffracted light
f, F	acoustic frequency, normalized acoustic frequency
H	acoustic beam height
k_i, k_d, k_a	wavevector: incident, diffracted light, acoustic wave
L, l	interaction length, normalized interaction length
L_o	characteristic length
M	figure of merit
n_o, n_e	refractive index: ordinary, extraordinary
P_a, P_d	acoustic power, acoustic power density
p, p_{mn}, p_{ijkl}	elasto-optic coefficient

S, S_I	strain, strain tensor components
t_r, T	rise time scan time
V	acoustic velocity
W	bandpass function

6.2 INTRODUCTION

When an acoustic wave propagates in an optically transparent medium, it produces a periodic modulation of the index of refraction via the elasto-optical effect. This provides a moving phase grating which may diffract portions of an incident light into one or more directions. This phenomenon, known as the acousto-optic (AO) diffraction, has led to a variety of optical devices that can be broadly grouped into AO deflectors, modulators, and tunable filters to perform spatial, temporal, and spectral modulations of light. These devices have been used in optical systems for light-beam control, optical signal processing, and optical spectrometry applications.

Historically, the diffraction of light by acoustic waves was first predicted by Brillouin¹ in 1921. Nearly a decade later, Debye and Sears² and Lucas and Biquard³ experimentally observed the effect. In contrast to Brillouin's prediction of a single diffraction order, a large number of diffraction orders were observed. This discrepancy was later explained by the theoretical work of Raman and Nath.⁴ They derived a set of coupled wave equations that fully described the AO diffraction in unbounded isotropic media. The theory predicts two diffraction regimes; the Raman-Nath regime, characterized by the multiple of diffraction orders, and the Bragg regime, characterized by a single diffraction order. Discussion of the early work on AO diffraction can be found in Ref. 5.

The earlier theoretically work tend to treat AO diffraction from a mathematical point of view, and for decades, solving the multiple-order Raman-Nath diffraction has been the primary interest on acousto-optics research. As such, the early development did not lead to any AO devices for practical applications prior to the invention of the laser. It was the need of optical devices for laser beam modulation and deflection that stimulated extensive research on the theory and practice of AO devices. Significant progress has been made in the decade from 1966 to 1976, due to the development of superior AO materials and efficient broadband ultrasonic transducers. During this period several important research results of AO devices and techniques were reported. These include the works of Gordon⁶ on the theory of AO diffraction in finite interaction geometry, by Korpel et al. on the use of acoustic beam steering,⁷ the study of AO interaction in anisotropic media by Dixon,⁸ and the invention of AO tunable filter by Harris and Wallace⁹ and Chang.¹⁰ As a result of these basic theoretical works, various AO devices were developed and demonstrated its use for laser beam control and optical spectrometer applications. Several review papers during this period are listed in Refs. 11 to 14.

Intensive research programs in the 1980s and early 1990s further advanced the AO technology in order to explore the unique potential as real-time spatial light modulators (SLMs) for optical signal processing and remote sensing applications. By 1995, the technology had matured, and a wide range of high performance AO devices operating from UV to IR spectral regions had become commercially available. These AO devices have been integrated with other photonic components and deployed into optical systems with electronic technology in diverse applications.

It is the purpose of this chapter to review the theory and practice of bulk-wave AO devices and their applications. In addition to bulk AO, there have also been studies based on the interaction of optical guided waves and surface acoustic waves (SAW). Since the basic AO interaction structure and fabrication process is significantly different from that of the bulk acousto-optics, this subject is treated separately in Chap. 7.

This chapter is organized as follows: Section 6.3 discusses the theory of acousto-optic interaction. It provides the necessary background for the design of acousto-optic devices. The subject of acousto-optic materials is discussed in Sec. 6.4. The next three sections deal with the three basic types of acousto-optic devices. Detailed discussion of AO deflectors, modulators, and tunable filters are presented in Section 6.5, 6.6, and 6.7, respectively.

6.3 THEORY OF ACOUSTO-OPTIC INTERACTION

Elasto-Optic Effect

The elasto-optic effect is the basic mechanism responsible for the AO interaction. It describes the change of refractive index of an optical medium due to the presence of an acoustic wave. To describe the effect in crystals, we need to introduce the elasto-optic tensor based on Pockels' phenomenological theory.¹⁵

An elastic wave propagating in a crystalline medium is generally described by the strain tensor S , which is defined as the symmetric part of the deformation gradient

$$S_{ij} = \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right) / 2 \quad i, j = 1 \text{ to } 3 \quad (1)$$

where u_i is the displacement. Since the strain tensor is symmetric, there are only six independent components. It is customary to express the strain tensor in the contracted notation

$$S_1 = S_{11} \quad S_2 = S_{22} \quad S_3 = S_{33} \quad S_4 = S_{23} \quad S_5 = S_{13} \quad S_6 = S_{12} \quad (2)$$

The conventional elasto-optic effect introduced by Pockels states that the change of the impermeability tensor ΔB_{ij} is linearly proportional to the symmetric strain tensor.

$$\Delta B_{ij} = p_{ijkl} S_{kl} \quad (3)$$

where p_{ijkl} is the elasto-optic tensor. In the contracted notation

$$\Delta B_m = p_{mn} S_n \quad m, n = 1 \text{ to } 6 \quad (4)$$

Most generally, there are 36 components. For the more common crystals of higher symmetry, only a few of the elasto-optic tensor components are nonzero.

In the above classical Pockels' theory, the elasto-optic effect is defined in terms of the change of the impermeability tensor ΔB_{ij} . In the more recent theoretical work on AO interactions, analysis of the elasto-optic effect has been more convenient in terms of the nonlinear polarization resulting from the change of dielectric tensor $\Delta \epsilon_{ij}$. We need to derive the proper relationship that connects the two formulations.

Given the inverse relationship of ϵ_{ij} and B_{ij} in a principal axis system $\Delta \epsilon_{ij}$ is

$$\Delta \epsilon_{ij} = -\epsilon_{ii} \Delta B_{ij} \epsilon_{jj} = -n_i^2 n_j^2 \Delta B_{ij} \quad (5)$$

where n_i is the refractive index. Substituting Eq. (3) into Eq. (5), we can write

$$\Delta \epsilon_{ij} = \chi_{ijkl} S_{kl} \quad (6)$$

where we have introduced the elasto-optic susceptibility tensor

$$\chi_{ijkl} = -n_i^2 n_j^2 p_{ijkl} \quad (7)$$

For completeness, two additional modifications of the basic elasto-optic effect are discussed as follows.

Roto-Optic Effect Nelson and Lax¹⁶ discovered that the classical formulation of elasto-optic effect was inadequate for birefringent crystals. They pointed out that there exists an additional roto-optic susceptibility due to the antisymmetric rotation part of the deformation gradient.

$$\Delta B'_{ij} = p'_{ijkl} R_{kl} \quad (8)$$

where $R_{ij} = (S_{ij} - S_{ji})/2$.

It turns out that the roto-optic tensor components can be predicted analytically. The coefficient of p'_{ijkl} is antisymmetric in kl and vanishes except for shear waves in birefringent crystals. In a uniaxial crystal the only nonvanishing components are $p_{2323} = p_{2313} = (n_o^{-2} - n_e^{-2})/2$, where n_o and n_e are the principal refractive indices for the ordinary and extraordinary wave, respectively. Thus, the roto-optic effect can be ignored except when the birefringence is large.

Indirect Elasto-Optic Effect In the piezoelectric crystal, an indirect elasto-optic effect occurs as the result of the piezoelectric effect and electro-optic effect in succession. The effective elasto-optic tensor for the indirect elasto-optic effect is given by¹⁷

$$p_{ij}^* = p_{ij} - \frac{r_{im} S_m e_{jn} S_n}{\epsilon_{mn} S_m S_n} \quad (9)$$

where p_{ij} is the direct elasto-optic tensor, r_{im} is the electro-optic tensor, e_{jn} is the piezoelectric tensor, ϵ_{mn} is the dielectric tensor, and S_m is the unit acoustic wavevector. The effective elasto-optic tensor thus depends on the direction of the acoustic mode. In most crystals the indirect effect is negligible. A notable exception is LiNbO₃. For instance, along the z axis, $r_{33} = 31 \times 10^{-12}$ m/v, $e_{33} = 1.3$ c/m², $E_{33}^s = 29$, thus $p^* = 0.088$, which differs notably from the contribution $p_{33} = 0.248$.

Plane Wave Analysis of Acousto-Optic Interaction

We now consider the diffraction of light by acoustic waves in an optically transparent medium. As pointed out before, in the early development, the AO diffraction in isotropic media was described by a set of coupled wave equations known as the Raman-Nath equations.⁴ In this model, the incident light is assumed to be a plane wave of infinite extent. It is diffracted by a rectangular sound column into a number of plane waves propagating along different directions. Solution of the Raman-Nath equations gives the amplitudes of these various orders of diffracted optical waves.

In general, the Raman-Nath equations can be solved only numerically and judicious approximations are required to obtain analytic solutions. Using a numerical procedure computation Klein and Cook¹⁸ calculated the diffracted light intensities for various diffraction orders in this regime. Depending on the interaction length L relative to a characteristic length $L_o = n\Lambda^2/\lambda_o$, where n is the refractive index and Λ and λ_o are wavelengths of the acoustic and optical waves, respectively, solutions of the Raman-Nath equations can be classified into three different regimes.

In the Raman-Nath regime, where $L \ll L_o$, the AO diffraction appears as a large number of different orders. The diffraction is similar to that of a thin phase grating. The direction of the various diffraction orders are given by the familiar grating equation, $\sin \theta_m = m\lambda_o/n\Lambda$, where m is the diffraction order. Solution of the Raman-Nath equations shows that the amplitude of the m th-order diffracted light is proportional to the m th-order Bessel functions. The maximum intensity of the first-order diffracted light (relative to the incident light) is about 34 percent. Due to this relatively low efficiency, AO diffraction in the Raman-Nath regime is of little interest to practical device applications.

In the opposite limit, $L \gg L_o$, the AO diffraction appears as a predominant first order and is said to be in the Bragg regime. The effect is called Bragg diffraction since it is similar to that of the x-ray diffraction in crystals. In the Bragg regime the acoustic column is essentially a plane wave of infinite extent. An important feature of the Bragg diffraction is that the maximum first-order diffraction efficiency obtainable is 100 percent. Therefore, practically all of today's AO devices are designed to operate in the Bragg regime.

In the immediate case, $L \leq L_o$, the AO diffraction appears as a few dominant orders. This region is referred as the near Bragg region since the solutions can be explained based on the near field effect of the finite length and height of the acoustic transducer.

Many modern AO devices are based on the light diffraction in anisotropic media. The Raman-Nath equations are no longer adequate and a new formulation is required. We have previously presented a plane wave analysis of AO interaction in anisotropic media.¹³ The analysis was patterned after that of Klienman¹⁹ used in the theory of nonlinear optics. Unlike the Raman-Nath equations wherein the optical plane waves are diffracted by an acoustic column, the analysis assumes that the acoustic wave is also a plane wave of infinite extent. Results of the plane wave AO interaction in anisotropic media are summarized as follows.

The AO interaction can be viewed as a parametric process where the incident optical plane wave mixes with the acoustic wave to generate a number of polarization waves, which in turn generate new optical plane waves at various diffraction orders. Let the angular frequency and optical wavevector of the incident optical wave be denoted by ω_m and \vec{k}_o , respectively, and those of the acoustic waves by ω_a and \vec{k}_a . The polarization waves and the diffracted optical waves consist of waves with angular frequencies $\omega_m = \omega_o + m\omega_a$ and wavevectors $\vec{K}_m = \vec{k}_o + m\vec{k}_a$ ($m \pm 1, \pm 2, \dots$). The diffracted optical waves are new normal modes of the interaction medium with the angular frequencies $\omega_m = \omega_o + m\omega_a$ and wavevectors \vec{k}_m making angles θ_m with the z axis.

The total electric field of the incident and diffracted light be expanded in plane waves as

$$\vec{E}(r, t) = \frac{1}{2} \sum \hat{e}_m E_m(z) \exp j(\omega_m t - \vec{k}_m \cdot \vec{r}) + \text{c.c.} \quad (10)$$

where \hat{e}_m is a unit vector of the electric field of the m th wave, E_m is the slowly varying amplitude of the electric field and c.c. stands for the complex conjugate. The electric field of the optical wave satisfies the wave equation,

$$\nabla \times \nabla \times \vec{E} + \frac{1}{c^2} \left(\vec{\epsilon} \cdot \frac{\partial^2 \vec{E}}{\partial t^2} \right) = -\mu_o \frac{\partial^2 \vec{P}}{\partial t^2} \quad (11)$$

where $\vec{\epsilon}$ is the relative dielectric tensor and \vec{P} is the acoustically induced polarization. Based on Pockels' theory of the elasto-optic effect,

$$\vec{P}(r, t) = \epsilon_o \vec{\chi} \vec{S}(r, t) \vec{E}(r, t) \quad (12)$$

where $\vec{\chi}$ is the elasto-optical susceptibility tensor defined in Eq. (7). $\vec{S}(r, t)$ is the strain of the acoustic wave

$$\vec{S}(r, t) = 1/2 (\hat{s} S e^{j(\omega_a t - \vec{k}_a \cdot \vec{r})} + \text{c.c.}) \quad (13)$$

where \hat{s} is a unit strain tensor of the acoustic wave and S is the acoustic wave amplitude. Substituting Eqs. (10), (12), and (13) into Eq. (11) and neglecting the second-order derivatives of electric-field amplitudes, we obtain the coupled wave equations for AO Bragg diffraction.

$$\frac{dE_m}{dz} = \frac{j(\omega_o/c)^2}{4k_m \cos \gamma_m} (\chi_m S E_{m-1} e^{-j\Delta \vec{k}_m \cdot \vec{r}} + \chi_{m+1} S^* E_{m+1} e^{j\Delta \vec{k}_{m+1} \cdot \vec{r}}) \quad (14)$$

where $\chi_m = n_m^2 n_{m-1}^2 p_m$, $p_m = \hat{e}_m \cdot \vec{p} \cdot \vec{s} \cdot \hat{e}_{m-1}$, γ_m is the angle between \vec{k}_m and the z axis, and $\Delta \vec{k}_m = \vec{K}_m - \vec{k}_m = \vec{k}_o + m\vec{k}_a - \vec{k}_m$ is the momentum mismatch between the optical polarization waves and m th-order normal modes of the medium. Equation (14) is the coupled wave equation describing the AO interaction in an anisotropic medium. Solution of the equation gives the field of the optical waves in various diffraction orders.

Two-Wave AO Interaction In the Bragg limit, the coupled wave equation reduces to the two-wave interaction between the incident and the first-order diffracted light ($m = 0, 1$):

$$\frac{dE_d}{dz} = \frac{j\pi n_i^2 n_d P_e}{2\lambda_o \cos\gamma_o} S E_i e^{j\Delta\vec{k}\cdot\hat{z}} \quad (15)$$

$$\frac{dE_i}{dz} = \frac{j\pi n_d^2 n_i P_e}{2\lambda_o \cos\gamma_o} S^* E_d e^{-j\Delta\vec{k}\cdot\hat{z}} \quad (16)$$

where n_i and n_d are the refractive indices for the incident and diffracted light, $p_e = \hat{e}_d \cdot \hat{p} \cdot \hat{s} \cdot \hat{e}_i$ is the effective elasto-optic constant for the particular mode of AO interaction, γ_o is the angle between the z axis and the median of incident and diffracted light and, $\Delta\vec{k} \cdot \hat{z} = \Delta k_z$ is the component of the momentum mismatch $\Delta\vec{k}$ along the z axis, and $\Delta\vec{k}$ is the momentum mismatch between the polarization wave \vec{K}_d and the free wave \vec{k}_d of the diffracted light.

$$\Delta\vec{k} = \vec{K}_d - \vec{k}_d = \vec{k}_i + \vec{k}_a - \vec{k}_d \quad (17)$$

Equations (15) and (16) admit simple analytic solutions. At $z = L$, the intensity of the first-order diffracted light (normalized to the incident light) is

$$I_1 = \frac{I_d(L)}{I_i(0)} = \eta \operatorname{sinc}^2 \frac{1}{\pi} \left(\eta + \left(\frac{\Delta k_z L}{2} \right)^2 \right)^{1/2} \quad (18)$$

where $\operatorname{sinc}(x) = (\sin\pi x)/\pi x$, and

$$\eta = \frac{\pi^2}{2\lambda_o^2} \left(\frac{n^6 p^2}{2} \right) S^2 L^2 = \frac{\pi^2}{2\lambda_o^2} M_2 P_a \left(\frac{L}{H} \right) \quad (19)$$

In the above equation, we have used the relation $P_a = \rho V^3 S^2 LH/2$, where P_a is the acoustic power, H is the acoustic beam height, ρ is the mass density, V is the acoustic wave velocity, and $M_2 = n^6 p^2 / \rho V^3$ is a material figure of merit.

Far Bragg Regime Equation (18) shows that in the far-field limit ($L \rightarrow \infty$) the diffracted light will build up finite amplitude only when the exact phase matching is met. When $\Delta k = 0$, the diffracted light intensity becomes

$$I_1 = \sin^2 \sqrt{\eta} \quad (20)$$

where $\eta \ll 1$, $I_1 \approx \eta$ and the diffraction efficiency is linearly proportional to acoustic power. This is referred to as the weak interaction approximation (or low efficiency regime). As acoustic power increases, the diffraction efficiency approaches 100 percent. However, the acoustic power required is about 2.5 times that predicted by the low efficiency regime.

Near Bragg Regime In the near field the growth of the diffraction is determined by the accumulated phase mismatch over the interaction length. The fractional diffracted light I_1 can be approximated by

$$I_1 = \eta \sin^2 \psi \quad (21)$$

where $\psi = \Delta k_z L / 2\pi$ is the phase mismatch (normalized to 2π). In the following, we shall discuss first AO diffraction in the far Bragg limit, that is, when exact phase matching is satisfied.

Phase Matching

Particle Picture of AO Diffraction In the far Bragg limit ($L \rightarrow \infty$), the diffracted light will build up finite amplitude only when the exact phase matching is met.

$$\vec{k}_d = \vec{k}_i \pm \vec{k}_a \quad (22)$$

In this limit case, the AO diffraction can be viewed as the interaction of optical and acoustic plane waves of infinite extent. To analyze AO interaction, it is more conveniently to use the particle picture of plane waves. The optical and acoustic plane waves can be thought of as made of photons and phonons of well-defined momentum and energy. Equation (22) simply states the principle of conservation of total momentum in the collision process. The principle of conservation of energy may be written as

$$\omega_d = \omega_i \pm \omega_a \quad (23)$$

The optical frequency of the diffracted light is shifted by the frequency of the acoustic wave. In analog to the Doppler effect, the optical frequency is up or down shifted if the direction of the acoustic wave is the same as or opposite to that of light wave.

In the following, we shall apply the conservation principles to derive the basic characteristics of light diffraction by acoustic wave in the plane wave formulation. For the general case of AO interaction in an anisotropic medium, the magnitudes of the incident optical wavevector can be written as

$$k_i = \frac{2\pi n_i}{n_o \lambda} \quad k_d = \frac{2\pi n_d}{n_o \lambda} \quad k_a = \frac{2\pi}{\Lambda} \quad (24)$$

where $\lambda = \lambda_o / n_o$ is the wavelength of the o-wave in the medium, n_i and n_d are refractive indices for the incident and diffracted light, $\Lambda = V/f$ is the acoustic wavelength and V and f are the acoustic velocity and frequency, respectively. Since the acoustic frequency is typically below 10^{-4} of the optical frequency, the small optical frequency difference (due to acoustic frequency shift) of the incident and diffracted light beams are neglected in Eq. (24).

Isotropic Diffraction Consider first the case of isotropic diffraction. Figure 1a shows the wavevector interaction geometry in the interaction plane. The loci of the incident and diffracted optical wavevectors fall on a circle of radius n_o . The principle of momentum conservation requires that the acoustic and optical wavevectors form a closed triangle. For isotropic AO diffraction the triangle is isosceles, and the incident and diffracted optical wavevectors make the same angle with the acoustic wavefront at the Bragg angle θ_b .

$$\sin \theta_b = \frac{\lambda}{2\Lambda} = \frac{\lambda_o f}{2n_o V} \quad (25)$$

Consider, for example, the diffraction of a laser beam at 633 nm by a longitudinal mode acoustic wave in TeO₂, $n = 2.216$, $V = 4.2$ mm/ μ sec, Eq. (25) yields $\theta_b = 1.9^\circ$ at an acoustic frequency $f = 500$ MHz. For all practical purpose the Bragg angle θ_b is small and the incident optical beam is always nearly

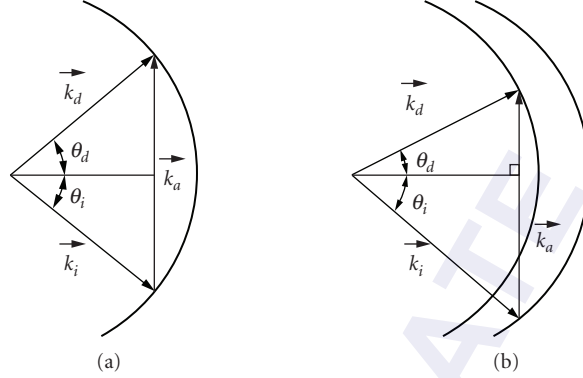


FIGURE 1 Vector diagram of acousto-optic interaction: (a) isotropic diffraction and (b) birefringent diffraction.

perpendicular to the acoustic wavefront. Thus, for an isotropic AO diffraction the primary effect of the acoustically driving polarization is to provide a transverse spatial light modulation (SLM). Similar to a grating, the transverse SLM acts as an acoustically driven optical beam deflector that scans the output light beam as the acoustic frequency is changed. The deflection angle θ_D , defined as the angle of separation between the incident and diffracted light, is given by (in the small-angle approximation)

$$\theta_D = \frac{\lambda}{\Lambda} = \frac{\lambda f}{V} \quad (26)$$

Anisotropic Diffraction Next, consider the AO diffraction in an optically anisotropic medium such as a birefringent crystal. In this case there are two distinct loci of the optical wavevectors (normal surfaces), an ordinary optical wave (polarized perpendicular to the c axis) and an extraordinary optical wave (polarized parallel to the c axis). The refraction indices are in general dependent on the propagation direction and polarization of the optical wave. The acoustic wave could couple the incident and diffracted optical wave of the same polarization ($o \leftrightarrow o$, $e \leftrightarrow e$) or orthogonal polarization ($o \leftrightarrow e$). In the latter case, the AO diffraction occurs between two loci of unequal refractive indices; this process is referred as the birefringent diffraction.

Phase Matching Equations Consider the AO interaction in a positive uniaxial crystal. Figure 1b shows the wavevector diagram for $o \leftrightarrow e$ type birefringent diffraction in the interaction plane. Applying the law of cosine to the triangle yields the following pair of equations for the incident and diffraction angle.

$$\sin\theta_i = \frac{\lambda_o}{2n_i\Lambda} \left[1 + \frac{\Lambda^2}{\lambda_o^2} (n_i^2 - n_d^2) \right] \quad (27)$$

$$\sin\theta_d = \frac{\lambda_o}{2n_d\Lambda} \left[1 - \frac{\Lambda^2}{\lambda_o^2} (n_i^2 - n_d^2) \right] \quad (28)$$

where n_i and n_d are the refractive indices for the incident and diffracted light, respectively. Equations (27) and (28) were first derived by Dixon⁸ for AO diffraction in a uniaxial crystal wherein the interaction plane is perpendicular to the c axis. The wavevector loci for the ordinary polarized light (o -wave) and the extraordinary polarized light (e -wave) are concentric circles; thus the two

equations can be separately solved. In general the refractive index for an extraordinary polarized light (e -wave) is a function of the polar angle of the incident and diffracted light. Thus, Dixon's equations are not applicable in practice since Eqs. (27) and (28) are coupled.

In order to decouple the input and diffracted angles from this pair of equations, we have chosen to express the wavevectors in the exact momentum matching condition using an elliptical parametric representation. Since the polarization is switched in the birefringent diffraction, the output characteristics of the diffracted light are same as the input characteristics of the incident light with the orthogonal polarization. Therefore, we need to consider only the input frequency-angle characteristics of the birefringent diffraction. The phase matching equations may be written in a form similar to the Dixon equations

$$e \rightarrow o: \quad \hat{s}_e \cdot \hat{s}_a = \frac{\lambda}{2\mu_e \Lambda} \left[1 + \frac{\Lambda^2}{\lambda^2} (\mu_e^2 - 1) \right] \quad (29)$$

$$o \rightarrow e: \quad \bar{\sigma}_o \cdot \bar{\sigma}_a = \frac{\lambda}{2\Lambda} \left[\bar{\sigma}_a^2 + \frac{\Lambda^2}{\lambda^2} (\bar{\sigma}_o^2 - 1) \right] \quad (30)$$

where

$$\begin{aligned} \hat{s}_a &= (\cos\theta_a, \sin\theta_a) & \hat{s}_e &= (\cos\theta_e, \sin\theta_e) & \mu_e &= (\cos^2\theta_e + e^{-2}\sin^2\theta_e)^{-1/2} \\ \bar{\sigma}_a &= (\cos\theta_a, e^{-1}\sin\theta_a) & \bar{\sigma}_o &= (\cos\theta_o, e^{-1}\sin\theta_o) \end{aligned} \quad (31)$$

where $e = n_e/n_o$ is the ratio of the principal refractive indices, θ_a , θ_o , and θ_e are the polar angles of the acoustic wave, ordinary, and extraordinary optical wave, respectively. The phase matching Eqs. (29) and (30) are the basic equations for AO diffraction in the Bragg regime.

To proceed with the analysis we introduce a dimensionless parameter a , defined as the ratio of optical wavelength (inside the medium) to the acoustic wavelength, $a = \lambda/\Lambda$. In terms of the wavelength ratio Eqs. (29) and (30) can be written in the form of quadratic equations.

$$q_{eo}(a) = a^2 - 2a\mu_e(\hat{s}_e \cdot \hat{s}_a) + \mu_e^2 - 1 = 0 \quad (32)$$

$$q_{oe}(a) = \bar{\sigma}_a^2 a^2 - 2a(\bar{\sigma}_a \cdot \bar{\sigma}_o) + \bar{\sigma}_o^2 - 1 = 0 \quad (33)$$

Equations (32) and (33) may be considered the dispersion relations for the acousto-optic "grating" for an extraordinary and ordinary light input, respectively. Consider, for example, the diffraction of an extraordinary wave into an ordinary wave in a uniaxial crystal. Solving the quadratic equation [Eq. (32)] yields the wavelength ratio a for the $e \rightarrow o$ type diffraction

$$a = \mu_e \left\{ \sin\theta_i \pm \sqrt{\sin^2\theta_i - (1 - \mu_e^{-2})} \right\} \quad (34)$$

where $\theta_i = \theta_e - \theta_a - \pi/2$ is the incidence angle. Equation (34) admits two real roots if $|\sin\theta_i| \geq \sqrt{1 - \mu_e^{-2}}$. Depending on the separation of the two roots, the AO diffraction can be classified into two regimes; referred as transverse and longitudinal spatial modulation.²⁰

Transverse Spatial Modulation By selecting the incidence angle near its critical value, the two roots are close to each other. This happens when direction of the acoustic wave is nearly perpendicular to that of the optical waves. The acoustically driven polarization provides a transverse spatial modulation (TSM) to the incident light beam. Similar to isotropic AO diffraction, the TSM acts as an optical beam deflector.

The characteristics birefringent diffraction can be utilized to provide significant performance advantages for of the AO deflector are. For instance, when the light beam is chosen to incident at the critical angle, Eq. (32) yields two equal roots and the acoustic roots. The AO diffraction is said to operate at the tangential phase matching (TPM) since acoustic wavevector is tangential to the locus of the diffracted wavevector. At TPM the AO deflector acquires a wideband frequency characteristics since a range of acoustic frequencies nearly satisfy the exact phase matching condition. The optic and acoustic wavelengths at TPM are related by

$$\lambda_o = \Lambda_t \sqrt{\mu_e^2 - 1} \approx \Lambda_t \sqrt{2n_o \Delta n} \sin \theta_e \quad (35)$$

where $\Delta n = n_e - n_o$ is the birefringence. At $\theta_e = 90^\circ$, the acoustic frequency reaches a maximum value.

$$\lambda_o \approx \Lambda_t \sqrt{2n_o \Delta n} \quad (36)$$

Longitudinal Spatial Modulation By selecting the optical incidence much larger than the critical angle, the two terms in Eq. (34) are about equal. Only the solution that corresponds to the difference of the two terms yields a low acoustic frequency near its minimum value. In this case the acoustic wavevector has a larger component along the direction of the optical wave, the AO diffraction thus provides a longitudinal spatial modulation (LSM) to the incident light. Unlike the previous case the LSM acts as an acoustically tuned optical bandpass filter. When the acoustic frequency is changed, the passband wavelength is tuned accordingly while the diffracted optical beam remains at a fixed deflection. This LSM-based AO device forms the basis of the acousto-optic tunable filter (AOTF).

Since the wavelength ratio is small, in this case the second order in Eq. (32) or (33) can be neglected. For an incident light with either polarization, an approximate solution of the center wavelength the AOTF passband is

$$\lambda_o \approx \frac{\Lambda \Delta n \sin^2 \theta_o}{\cos(\theta_a - \theta_o)} \quad (37)$$

where $\Delta n = n_e - n_o$ is the birefringence. At $\theta_a = 90^\circ$, the incident and diffracted optical waves are collinear with the acoustic wave, Eq. (37) reduced to

$$\lambda_o = \Lambda \Delta n \quad (38)$$

Equation (38) is the momentum-matching condition for the collinear AOTF.

Frequency Characteristics of AO Interaction

In the near-Bragg regime, the AO diffraction is determined by the z -component of phase mismatch function $\psi = \Delta k_z L / 2\pi = \Delta \sigma_z L$. In our analysis, this z -component mismatch is caused by the deviation of the light incidence from exact momentum matching. Based on this model it can be shown that the momentum mismatch (normalized to 2π) is $\Delta \sigma_z = (n_o/2)q$, where q is defined in Eqs. (32) and (33) for $e \rightarrow o$ and $o \rightarrow e$ type diffraction, respectively. Substituting $\psi = Ln_o q/2$ from Eq. (32) or (33) into Eq. (21), the bandpass function of AO diffraction can be expressed as a function of the variation of the acoustic angle θ_a and acoustic frequency f . It is convenient to normalize the acoustic frequency to a center frequency f_o (wavelength Λ_o). In terms of the normalized acoustic frequency ($F = ff_o$), the phase mismatch ψ can be written as

$$\psi = \frac{l}{2}(F^2 - FF_m + F_t^2) \quad (39)$$

where

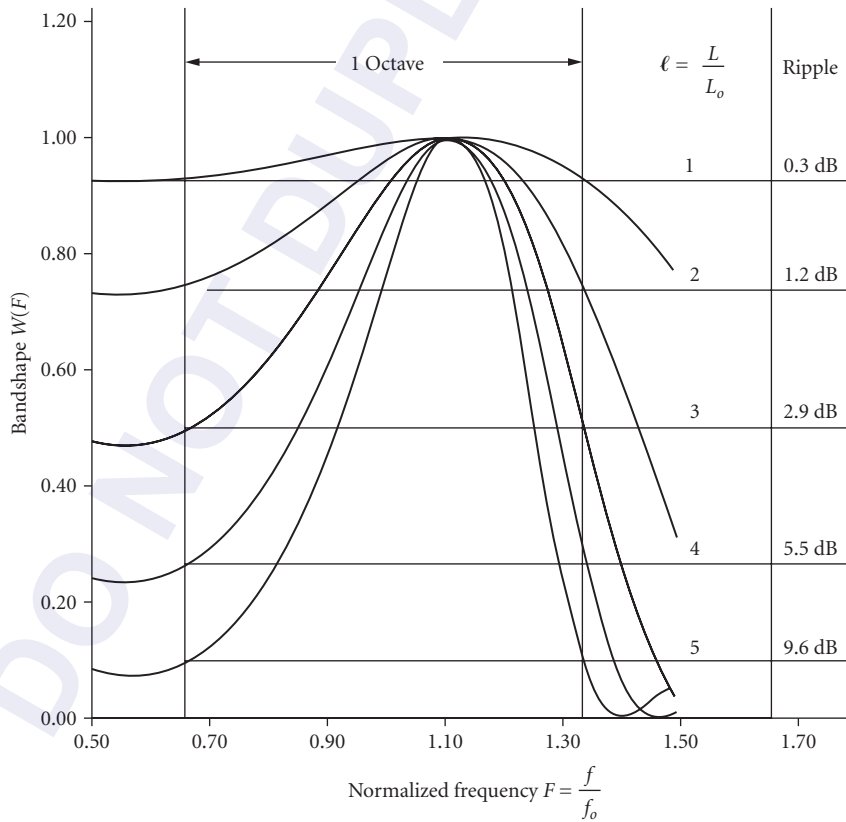
$$l = L/L_o \quad L_o = n_o \Lambda^2 / \lambda_o \quad F_m = 2(n_o \Lambda_o / \lambda_o) \mu_e \sin \theta_i \quad F_i = f_i / f_o \quad (40)$$

Isotropic Diffraction Bandshape By choosing $F_m = 1 + (\Delta F/2)^2$, the bandshape function $W(\psi)$ has equal loss at the two ends of the passband $F_1 = 1 + (\Delta F/2)$, where ΔF is the fractional bandwidth of the AO interaction. The diffraction efficiency reduces to 0.5, where $\psi = 0.45$. This corresponds to a fractional bandwidth

$$\Delta F = 1.8/l \quad (41)$$

To realize octave bandwidth, ($\Delta F = 2/3$) for instance, the normalized interaction length l is equal to 2.7. Figure 2a shows the bandshape of isotropic AO diffraction.

Birefringent Diffraction Bandshape For $e \rightarrow o$ diffraction, it is possible to obtain a wide bandpass response by operating near tangential phase matching (TPM). At TPM the two



(a)

FIGURE 2 Acousto-optic bandshapes: (a) isotropic diffraction and (b) birefringent diffraction.

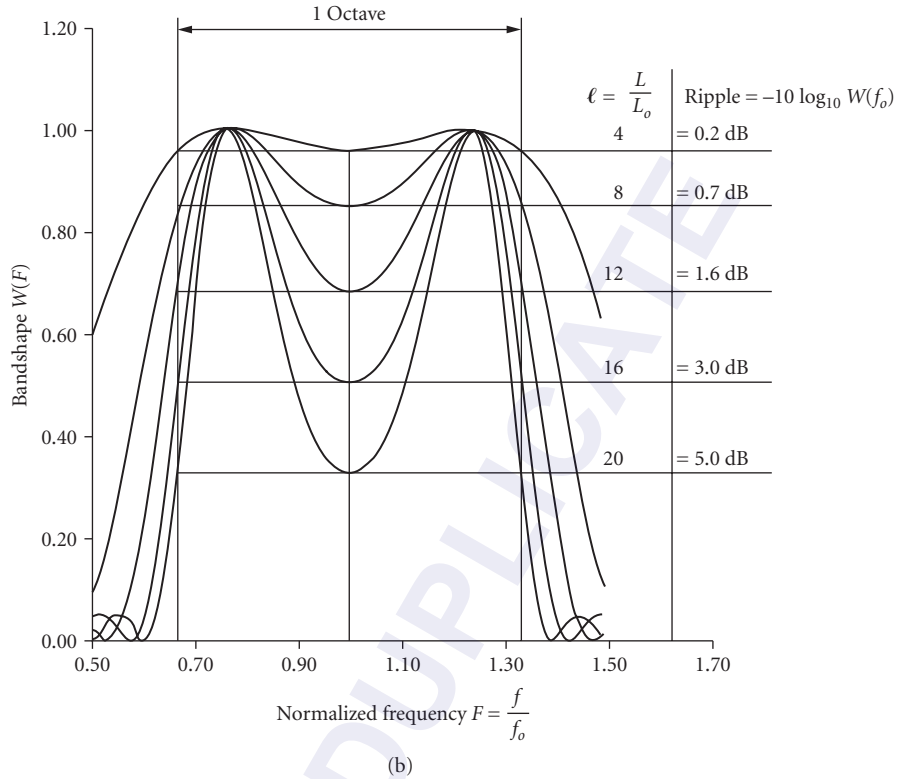


FIGURE 2 (Continued)

phase-matching frequencies coincide and the AO bandpass exhibits a flat-top shape with a bandwidth

$$\Delta F = \frac{\Delta f}{f} \approx \left(\frac{3.6}{l} \right)^{1/2} \quad (42)$$

For octave bandwidth ($\Delta F = 2/3$) the normalized interaction length l is equal to 8.1. This represents an efficiency advantage factor of 3 compared to isotropic AO diffraction. Figure 2b shows the bandshape of the birefringent AO diffraction.

Acousto-Optic Interaction of Finite Geometry

The plane wave analysis of AO interaction in the Bragg regime appears to be inadequate for real devices since the interaction geometry usually involves optical and acoustic beams of finite sizes with nonuniform amplitude distributions, for example, a Gaussian optical beam and a divergent acoustic beam. Gordon⁶ considered the AO diffraction in finite interaction geometry using the Green's function solution for the wave equation. To make the approach analytically tractable, the first-order Born's approximation is used, which is same as the weak interaction assumption.

Instead of the Green's function solution, an equivalent approach is to decompose the optical and acoustic waves into angular spectrum of plane waves (ASW) and add all the spatial frequency components of the

diffracted light that satisfy the exact phase-matching condition.^{21,22} Since the AO interaction is modeled as a filtering process in the spatial frequency domain, this approach is referred to as the frequency domain analysis.

As an example we use the frequency domain approach to determine the frequency and angular characteristics of an AO deflector. To simplify the analysis we will make a two-dimensional model for AO interaction in the interaction (x - z) plane. The reduction of diffraction efficiency in the height (y - z) plane will be calculated in the real domain as an overlapping integral of the acoustic and optical beams in the transverse plane along the height direction.

In the case of weak interaction, there is negligible depletion of the incident light. Assume that the incident light amplitude is a constant, Eq. (15) can be integrated to yield

$$E_d(\vec{\sigma}) = j\kappa S(\vec{\sigma}_a) E_i(\vec{\sigma}_i) \delta(\vec{\sigma} - \vec{\sigma}_i - \vec{k}_a) = j\kappa S(\vec{\sigma}_a) E_i(\sigma - \vec{\sigma}_a) \quad (43)$$

where $\vec{\sigma}$, $\vec{\sigma}_i$, and $\vec{\sigma}_a$ are the wavevector components in the interaction plane for the diffracted light, incident light and the acoustic waves, respectively.

We assume that the acoustic beam consists of acoustic plane waves propagating along the x axis in the interaction (x - z) plane. It is generated by an acoustic transducer with a top electrode of effective length L along the z axis and height H along the y axis. We also assume that the incident light beam consists of optical plane waves propagating in the x - z plane at an angle of incidence γ_i with the z axis. The incident light beam has a Gaussian with a beam waist $2\omega_1$ and $2\omega_2$ at e^{-2} of the intensity profile in the interaction and height plane, respectively. By expanding the optical and acoustic field in terms of their spatial frequency components in the interaction plane and applying Eq. (43), the total power spectrum of the diffracted light intensity can be written as

$$I_d(E, \xi) = \eta W(F) U(\xi) \cdot J \quad (44)$$

where η is the small signal diffraction efficiency given by Eq. (19). The next three multiplying factors are briefly discussed as follows.

Overlapping Integral The overlapping integral J of optical and acoustic profile in the height plane is given by

$$J = \frac{2}{\pi\omega_1\omega_2} \int_0^{H/2} \int_{-H/2}^{H/2} |V(x, y)|^2 e^{-y^2/\omega_2^2} e^{-(x-D/2)^2/\omega_1^2} dy dx \quad (45)$$

where $V(x, y)$ is the acoustic field distribution and can be determined by the Fresnel integral

$$V(x, y) = \int_{-H/2}^{H/2} \frac{1}{\sqrt{jB\Lambda x}} \left[\exp\left(\frac{j\pi(y_o - y)^2}{B\Lambda x}\right) \right] dy_o \quad (46)$$

where B is the curvature of the acoustic slowness surface in the transverse plane. The overlapping integral represents a reduction of diffraction efficiency of the AO diffraction due to acoustic diffraction in the transverse plane along the height direction. For a transducer of height H , the acoustic radiation in a region near the transducer is approximately collimated. This collimated region, known as the near field, extends a distance D from transducer and is given by

$$H_o = \sqrt{B\Lambda D} = V\sqrt{\tau B/f} \quad (47)$$

Numerical calculation of the overlapping integral J in Eq. (43) shows its value is on the order of unity if the transducer height H is chosen to be equal to H_o given by Eq. (47).

Bandpass Response The bandpass function $W(F)$ as a function of the normalized frequency F is equal to the magnitude square of the Fourier transform, that is, the power spectrum of the acoustic

field along the z axis. For a standard rectangular window of the transducer field, the bandpass function is same as that given by

$$W(\psi) = \text{sinc}^2(\psi) \quad (48)$$

where $\psi = \Delta\sigma_z L = (n_o/2)q$ is the phase mismatch along optical direction (z axis) and is given by Eq. (39) as a function of the normalized frequency F .

Spatial Frequency Spectrum The spatial frequency spectrum $U(\xi)$ expressed as a function of spatial frequency component ξ signal direction (x axis) is determined by the power spectrum of the incident optical beam profile along the x axis, the direction perpendicular to the optical direction. For input Gaussian light with a beam waist $2\omega_1$ along the z axis in the interaction plane, the spatial frequency spectrum $U(\xi)$ is given by

$$U(\xi) = \exp(-\Delta\sigma_x^2/\pi^2\omega_1^2) \quad (49)$$

where $\Delta\sigma_x = \xi - \sigma_i \sin \theta_i - \sigma_o$, ξ is the spatial frequency of the diffracted light along x axis.

6.4 ACOUSTO-OPTIC MATERIALS

The significant progress of AO devices has been due largely to the development of superior materials such as TeO_2 and GaP. In this section we shall review the material issues related to AO device applications. A comprehensive list of tables summarizing the properties of AO materials is presented at the end of this section. The AO material issues and the selection guidelines have been discussed in a previous publication.¹²

The selection of AO materials depends on the specific device application.²³ An AO material suited for one type of device may not even be applicable for another. For example, GaP is perhaps the best choice for making wideband AO deflectors or modulators. However, since GaP is optically isotropic, it cannot be used for tunable filter purposes.

Some of the requirements for materials' properties apply to the more general cases of optical device applications, for example, high optical transparency over wavelength range of interest, availability in large single crystals, and the like. We shall restrict the discussion to material properties that are particularly required for AO device applications

Acousto-Optic Figures of Merit

A large AO figure of merit is desired for device applications. There are several AO figures of merit that have been used for judging the usefulness of an AO material. The relevant one to be used depends on the specific applications. Several AO figures of merit are defined in the literature. These include

$$M_1 = \frac{n^7 p^2}{\rho V} \quad M_2 = \frac{n^6 p^2}{\rho V^3} \quad M_3 = \frac{n^7 p^2}{\rho V^2} \quad M_4 = \frac{n^8 p^2 V}{\rho} \quad M_5 = \frac{n^8 p^2}{\rho V^3} \quad (50)$$

where n is the index of refraction, p is the relevant elasto-optic coefficient, ρ is the density, and V is the acoustic wave velocity. These figures of merit are generally listed as the normalized quantities M (normalized to values for fused silica).

The figure of merit M_1 relates the diffraction efficiency η , to the acoustic power P_a for a given device aspect ratio L/H in accordance with Eq. (18). M_2 is the AO figure of merit most often referred to in the literature and is widely used for the comparison of AO materials. This is a misconception, since from the viewpoint of device applications, M_2 is usually not appropriate. Comparison of AO materials (or modes) based on M_2 can lead to erroneous conclusions. M_2 is used only when efficiency is the only parameter of

concern. In most practical cases, other parameters such as bandwidth and resolution must also be considered. To optimize the efficiency bandwidth product, the relevant figure of merit is $M_1 = nV^2M_2$.

In the design of AO deflectors, or Bragg cells, besides efficiency and bandwidth, a third parameter of interest is the aperture time τ . A minimum acoustic beam height H must be chosen to ensure that the aperture is within the near field of the acoustic radiation. With this choice, the relevant figure of merit for optimized efficiency for a specified bandwidth and aperture time is $M_3 = nVM_2$.

For wideband AO modulators, the acoustic power density P_d is often the limiting factor. The appropriate AO figure of merit is then M_4 , that is, $M_4 = nV^2M_1$.

In the design of AO tunable filters, the parameters to be optimized are the product of efficiency η , the resolving power $\lambda_o/\Delta\lambda$, and the solid angular aperture $\Delta\Omega$. In this case the appropriate AO figure of merit is $M_5 = n^2M_2$.

Acoustic Propagation and Attenuation

The performance of AO devices also depends on the acoustic properties of the interaction medium. As seen from the figures of merit listed above, a low acoustic velocity is one of the most important desired material parameters for AO devices. It not only provides the advantage of lower drive power, but it also allows the realization of large resolution for AO deflectors and tunable filters for a specified maximum crystal length.

The anisotropic acoustic propagation characteristic also plays an important role in the performance of AO devices. For instance, the small curvature of slowness surface of the transverse plane in GaP allows a smaller transducer height and greatly lowers the drive power. The use of large acoustic beam walk-off in TeO₂ allows the realization of a preferred interaction geometry to extend interaction length or optical aperture.

A low acoustic attenuation is also desired for increased resolution of deflectors or aperture of tunable filters. Generally, at room temperature $\omega\tau_{th} \ll 1$, where ω is the angular frequency of the acoustic wave and τ_{th} is the thermal phonon relaxation, the dominant contribution to acoustic attenuation is due to Akhieser loss caused by relaxation of the thermal phonon distribution toward equilibrium. A widely used result of this theory is the relation derived by Woodruff and Erhenrich. It states that acoustic attenuation measured in nepers per unit time is given by¹² $\alpha = \gamma^2 f^2 \kappa T / \rho V^4$, where γ is the Grüneisen constant, T is the temperature, and κ is the thermal conductivity. From this relation, it is seen that the AO figure of merit and the acoustic attenuation approximately follow an inverse relation. Notice also that the acoustic attenuation has a quadratic frequency-dependence for most crystals. In practice, in some crystals (such as GaP), it has been found that the frequency dependence of attenuation $a \sim f^n$, where n varies from 1 to 2 as the frequency increases ranges. The deviation from a quadratic dependence may be attributed to the additional extrinsic attenuation caused by scattering from lattice imperfections.

Optical Birefringence

Optical birefringence is a requirement for materials used in AO tunable filters. The requirement is met by optically birefringent crystals so that the AO filter interaction can operate as a longitudinal spatial modulator (LSM). For AO deflectors and modulators, optical birefringence is not necessary. AO devices with high efficiency, wide bandwidth, and large resolution are realizable with superior isotropic materials such as GaP. However, in an optically birefringent crystal it is possible to operate at tangential phase matching (TPM), which provides an enhancement of interaction length l for a given fractional bandwidth. The optical birefringence in LiNbO₃ and TeO₂ has been largely responsible for the superior device performance.

Tabulation of Acousto-Optic Material Properties

To aid in the selection of AO materials, the relevant properties of some promising materials are listed below. Table 1 lists the values of elasto-optical tensor components. Table 2, taken from Ref. 24

TABLE 1 Elasto-Optic Coefficients of Materials

(a) Isotropic			
Material	λ (μm)	P_{11}	P_{12}
Fused silica (SiO ₂)	0.63	+0.121	+0.270
As ₂ S ₃ glass	1.15	+0.308	+0.299
Water	0.63	0.31	0.31
Ge ₃₃ Se ₅₅ As ₁₂ (glass)	1.06	0.21	0.21
Lucite	0.63	0.30	0.28
Polystyrene	0.63	0.30	0.31
SF-59	0.63	0.27	0.24
SF-8	0.63	0.198	0.262
Tellurite glass	0.63	0.257	0.241

(b) Cubic classes $\bar{4}3m$, 432 , and $m\bar{3}m$			
Material	λ (μm)	P_{11}	P_{44}
CdTe	10.60	-0.152	-0.017
GaAs	1.15	-0.165	-0.140
GaP	0.633	-0.151	-0.082
Ge	2.0-2.2	-0.063	-0.0535
	10.60	0.151	0.124
InP	1.5	0.18	0.06
NaCl	0.55-0.65	0.115	0.159
NaF	0.633	0.08	0.20
	0.589		-0.021
Si	1.15	-0.101	0.0094
	3.39	-0.094	0.017
Y ₃ Fe ₅ O ₁₂ (YIG)	1.15	0.025	0.073
Y ₃ Al ₅ O ₁₂ (YAG)	0.633	-0.029	+0.0091
KRS5	0.633	0.18	0.27
KRS6	0.633	0.28	0.25
β -ZnS	0.633	0.091	-0.01
Y ₃ Ga ₅ O ₁₂	0.63	0.091	0.019
Diamond	0.59	-0.31	-0.03
	0.59	-0.43	+0.19
LiF	0.59	+0.02	+0.128
MgO	0.59	-0.32	-0.08
KBr	0.59	+0.22	+0.71
KCl	0.59	+0.17	+0.124
KI	0.59	+0.210	0.169

(c) Cubic classes 23 and m3

Material	λ (μm)	P_{11}	P_{12}	P_{13}	P_{14}	P_{15}	P_{16}
Ba(NO ₃) ₂	0.63	0.15	0.35	0.29	0.02	0.02	0.02
NaBrO ₃	0.59	0.185	0.218	0.213	-0.0139	-0.0139	-0.0139
NaClO ₃	0.59	0.162	0.24	0.2	-0.198	-0.198	-0.198
BA(NO ₃) ₂	0.63	0.15	0.35	0.29	0.02	0.02	0.02
Bi ₂ GeO ₂₀	0.63	0.12			0.04	0.04	0.04
Bi ₁₂ SiO ₂₀	0.63	0.13			0.04	0.04	0.04

(d) Hexagonal system: classes 6m2, 6mm, 622, and 6/m

Material	λ (μm)	P_{11}	P_{12}	P_{13}	P_{31}	P_{33}	P_{44}
CdS	0.63	-0.142	-0.066	-0.057	-0.041	-0.20	± 0.054
	10.60	0.104	0.011	0.011			
SnO	0.63	0.222	0.099	-0.111	0.088	-0.235	-0.0585
α -ZnS	0.63	-0.115	0.017	0.025	0.0271	-0.13	-0.0627
ZnO	0.63	0.222	0.199	-0.111	0.088	-0.235	-0.061

(e) Hexagonal system: classes 6, $\bar{6}$, and 6/m

Substance	λ (μm)	P_{11}	P_{12}	P_{13}	P_{31}	P_{33}	P_4	P_{45}	P_{16}
LiI O ₃	0.63	0.32			0.31				0.03

(f) Trigonal system: classes 3m, $\bar{3}2$, and $\bar{3}m$

Substance	λ (μm)	P_{11}	P_{12}	P_{13}	P_{14}	P_{31}	P_{33}	P_{41}	P_{44}
Al ₂ O ₃	0.644	-0.23	-0.03	0.02	0.00	-0.04	-0.20	0.01	-0.10
LiNbO ₃	0.633	-0.026	0.090	0.133	-0.075	0.179	0.071	-0.151	0.146
LiTaO ₃	0.633	-0.081	0.081	0.093	-0.026	0.089	-0.044	-0.085	0.028
SiO ₂ (quartz)	0.589	0.16	0.27	0.27	-0.030	0.29	0.10	-0.047	-0.079
Ag ₃ As ₃	0.633	± 0.10	± 0.19	± 0.22	± 0.22	± 0.24	± 0.20	± 0.100	± 0.01
(proustite)	1.15	± 0.056	± 0.082	± 0.068	± 0.103	± 0.103	± 0.100	± 0.100	± 0.01
Te	10.6	0.164	0.138	0.146	0.086	0.086	0.038	0.038	0.038
CaCO ₃		± 0.095	± 0.189	± 0.215	-0.006	0.309	+0.178	+0.01	-0.090
HgS	0.63			0.445			0.115		
Tl ₃ AsSe ₃	3.39	0.4	0.22	0.24	0.04	0.2	0.22	0.018	0.15
Tl ₃ AsS ₃	3.39	0.36	0.13	0.2	0.15	0.15	0.36	0.02	0.02

(Continued)

TABLE 1 Elasto-Optic Coefficients of Materials (*Continued*)

(g) Tetragonal system: classes 4mm, $\bar{4}2m$, 422, and 4/m										
Substance	λ (μm)	P_{11}	P_{12}	P_{13}	P_{31}	P_{33}	P_{44}	P_{66}	P_{44}	P_{66}
$(\text{NH}_4)_2\text{H}_2\text{P}_2\text{O}_7$ (ADP)	0.589	0.319	0.277	0.169	0.197	0.167	-0.058	-0.091	-0.058	-0.091
	0.63	0.296	0.243	0.208	0.188	0.228				
$\text{KH}_2\text{P}_2\text{O}_7$ (KDP)	0.589	0.287	0.282	0.174	0.241	0.122	-0.019	-0.064	-0.019	-0.064
	0.63	0.254	0.230	0.233	0.221	0.212				
$\text{Sr}_{0.75}\text{Ba}_{0.25}\text{Nb}_2\text{O}_6$	0.63	0.16	0.10	0.08	0.11	0.47				
$\text{Sr}_{0.5}\text{Ba}_{0.5}\text{Nb}_2\text{O}_6$	0.63	0.06	0.08	0.17	0.09	0.23				
TeO_2	0.63	0.0074	0.187	0.340	0.0905	0.240	0.04	-0.0463	0.04	-0.0463
TiO_2 (rutile)	0.514	-0.001	0.113	-0.167	-0.106	-0.064	0.0095	-0.066	0.0095	-0.066
	0.63	-0.011	0.172	-0.168	-0.0965	-0.058		± 0.072		± 0.072
ZrSiO_4	0.63	0.06		0.13	0.07	0.09		0.10		0.10
CdGeP_2	0.63	0.21	-0.09	0.09		0.4	0.1	0.12		0.12
Hg_2Cl_2	0.63	0.551	0.44	0.256	0.137	0.01		0.047		0.047
Hg_2Br_2	0.63	0.262	0.175	0.148	0.177	0.116				

(h) Tetragonal system: classes 4, $\bar{4}$, 422, and 4/m										
Material	λ (μm)	P_{11}	P_{12}	P_{13}	P_{31}	P_{33}	P_{44}	P_{45}	P_{61}	P_{66}
PbMoO_4	0.63	0.24	0.24	0.255	0.175	0.3	0.067	-0.01	0.013	0.05
CdMoO_4	0.63	0.12	0.10	0.13	0.18					
NaBiMoO_4	0.63	0.243	0.265	0.25	0.21	0.29				
LiBiMoO_4	0.63	0.265	0.201	0.244	0.227	0.309				
CaMoO_4	0.63	0.17	-0.15	-0.08	0.10	0.08	0.06	0.06	0.1	0.026

(i) Orthorhombic all classes													
Material	λ (μm)	P_{11}	P_{12}	P_{13}	P_{21}	P_{22}	P_{23}	P_{31}	P_{32}	P_{33}	P_{44}	P_{55}	P_{66}
α - HfO_3	0.63	0.406	0.277	0.304	0.279	0.343	0.305	0.503	0.310	0.334			0.092
PbCO_3	0.63	0.15	0.12	0.16	0.05	0.06	0.21	0.14	0.16	0.12			
BaSO_4	0.59	0.21	0.25	0.16	0.34	0.24	0.19	0.27	0.22	0.31	0.22	-0.012	0.037
$\text{Gd}_2(\text{MoO}_4)_3$	0.63	0.19	0.31	0.175	0.215	0.235	0.175	0.185	0.23	0.115	-0.033	-0.028	0.035

TABLE 2 Selected AO Materials

Material	Optical Transmission (μm)	Density (g/cm^3)	Acoustic Mode	Acoustic Velocity ($\text{mm}/\mu\text{s}$)	Acoustic Attenuation ($\text{dB}/\mu\text{s}\text{-GHz}^2$)	Optical Polarization*	Refraction Index	Figure of Merit		
								M_1	M_2	M_3
Fused silica	0.2–4.5	2.2	L	5.96	7.2	—	1.46	1.0	1.0	1.0
LiNbO ₃	0.4–4.5	4.64	L[100]	6.57	1.0	35° Y Rot.	2.2	8.5	4.6	7.7
	0.4–4.5	0.4–4.5	S(100)35°	~3.6	~1.0	[100]	2.2	2.3	4.2	3.8
TiO ₂	0.45–6.0	4.23	L[110]	7.93	~1.0	—	2.58	18.6	6.0	14
PbMoO ₄	0.42–5.5	6.95	L[001]	3.63	5.5	—	2.39	14.6	23.9	24
BGO	0.45–7.5	9.22	L[110]	3.42	1.6	Arb.	2.55	3.8	6.7	6.7
TeO ₂	0.35–5.0	6.0	L[001]	4.2	6.3	—	2.26	17.6	22.9	25
	0.4–4.5	6.0	S[110]	0.62	17.9	Cir.	2.26	13.1	795	127
GaP	0.6–10.0	4.13	L[110]	6.32	8.0	—	3.31	75.3	29.5	71
	0.6–10.0	4.13	S[110]	4.13	2.0	—	3.31	—	16.6	26
Tl ₃ AsS ₄	0.6–12.0	6.2	L[001]	2.15	5.0	—	2.83	152	523	416
Tl ₃ AsSe ₃	1.26–13.0	7.83	L[100]	2.05	14.0	—	3.34	607	2259	1772
Hg ₂ Cl ₂	0.38–28.0	7.18	L[100]	1.62	—	—	2.62	34	337	125
			S[110]	0.347	8.0	—	2.27	4.3	703	73
Ge ₃₃ As ₁₂ Se ₃₃	1.0–14.0	4.4	L	2.52	1.7	—	2.7	54.4	164	129
GaAs	1.0–11.0	5.34	L[110]	5.15	15.5	—	3.37	118	69	137
As ₂ Se ₃	0.9–11.0	4.64	L	2.25	27.5	—	2.89	204	722	539
Ge	2.0–20.0	5.33	L[111]	5.5	16.5	—	4.0	278	120	302

L: longitudinal; S: shear; arb.: arbitrary; cir.: circular birefringence.

*Optical polarization is normal or in the interaction plane.

lists the relevant properties of selected AO materials. The listed figures of merit M_1 , M_2 , and M_3 are normalized relative to that of fused silica, which has the following absolute values:

$$\begin{aligned} M_1 &= 7.83 \times 10^{-7} & [\text{cm}^2\text{sg}^{-1}] \\ M_2 &= 1.51 \times 10^{-18} & [\text{s}^3\text{g}^{-1}] \\ M_3 &= 1.3 \times 10^{-12} & [\text{cm}^2\text{s}^2\text{g}^{-1}] \end{aligned}$$

6.5 ACOUSTO-OPTIC DEFLECTOR

Acousto-optic interaction provides a simple means to deflect an optical beam in a sequential or random-access manner. As the driving frequency of the acoustic wave is changed, the direction of the diffracted beam can also be varied. The angle between the first-order diffracted beam and the incident beam for a frequency range Δf is approximately given by

$$\Delta\theta_d = \frac{\lambda\Delta f}{V} \quad (51)$$

In a deflector, the most important performance parameters are resolution and speed. The resolution, or the maximum number of resolvable spots, is defined as the ratio of the range of deflection angle divided by the angular spread of the diffracted beam, that is,

$$N = \frac{\Delta\theta_d}{\delta\theta_o} \quad (52)$$

The angular divergence of the incident optical beam is given by

$$\delta\theta_o = \xi \frac{\lambda}{D} \quad (53)$$

where D is the width of the incident beam and ξ is the a factor (near unity) that depends on the incident beam's amplitude distribution. For a nontruncated Gaussian beam $\xi = 4/\pi$. From Eqs. (51), (52), and (53) it follows

$$N \approx \tau \Delta f \quad (54)$$

where $\tau = D/V \cos \theta_o$ is the acoustic-transit time across the optical aperture. Notice that the acoustic-transit time also represents the (random) access time and is a measure of the speed of the deflector. Equation (54) shows that the resolution is equal to time (aperture) bandwidth product. This is the basic tradeoff relation between resolution and speed (or bandwidth) of AO deflectors.

It is instructive to specify the interaction geometry for AO deflectors using a dimensionless parameter Λ , the ratio of divergence angle between the optical and acoustic beams. For an acoustic wave generated from a flat transducer of width L , the acoustic beam divergence is given by

$$\delta\theta_a = \frac{\Lambda}{L} \quad (55)$$

TABLE 3 Interaction Geometry and Performance of AO Deflectors

Type I-Laser beam scanner
Regime of AO diffraction: Transverse SLM
Divergence ratio: $A \ll 1, N \gg 1$
Primary performance: High resolution, high peak efficiency
Type II-Signal processing Bragg cell
Regime of AO diffraction: Transverse SLM
Divergence ratio/Resolution: $A \ll 1, N \gg 1$
Primary performance: Wide bandwidth, large dynamic range

From the above equations, the divergence ratio A and resolution N are related by

$$A = \frac{\delta\theta_o}{\delta\theta_a} = \frac{l\Delta F}{N} \quad (56)$$

where $l = L/L_o$ is the normalized interaction length and $\Delta F = \Delta f/f_o$ is the fractional bandwidth. For isotropic AO diffraction F , $l\Delta F = 1.8$. Thus, an AO deflector is characterized by the requirement $A \ll 1$. As an example, to realize a high resolution AO deflector with 1000 resolvable spots, the divergence ratio must be less than 0.002. The use of expanded incident beam with diffracted limited optics becomes the most critical requirement in the design high-resolution AO deflectors.

In summary, the AO deflector can be viewed as a transverse spatial light modulator (SLM) with a large number of resolvable angular channels. Based on the primary performance goals and the interaction geometry the AO deflector can be further divided into type I-laser beam scanners and type II-signal processing Bragg cells. Table 3 lists the range of the divergence ratio and the key performance requirements for the two types of AO deflectors.

Deflector Interaction Geometry

Principle of Operation

Isotropic AO diffraction In the following we consider the design of AO deflectors.²⁵ Figure 3a and b shows the device geometry of an AO deflector in the interaction plane and transverse plane, respectively. A piezoelectric transducer is bonded to the appropriate face of a selected AO medium, oriented for efficient interaction. A top electrode deposited on the transducer defines the active area with equivalent acoustic beam length L and height H . An optical beam generally truncated to an aperture width, D is incident in the interaction plane at a proper Bragg angle with respect to the acoustic wavefront. The incident beam has a Gaussian profile with a beam waist $2\omega_1$ at $1/e^2$ of the intensity in the interaction plane and a beam waist $2\omega_2$ at $1/e^2$ intensity in the transverse plane.

In accordance to Eqs. (19) and (44) the peak diffraction efficiency of an AO diffraction under the weak interaction approximation is given by

$$\eta = \frac{\pi^2}{2\lambda_o^2} M_2 P_a \left(\frac{L}{H} \right) \cdot J \quad (57)$$

where L and H are the equivalent acoustic beam length and height and J is the overlapping integral given by Eq. (45). As an example, assume an optical Gaussian beam with a truncation ratio of 1.5 and a rectangular shaped electrode, numerical calculation shows when the electrode height is chosen to be equal to H_o so that the interaction region coincide with the acoustic near field. With this choice by Eq. (47), so that the AO interaction region extends from the transducer to the acoustic near field, the peak diffraction efficiency becomes

$$\eta_o = \frac{\pi^2 P_a}{2\lambda_o^3 f^{3/2} \tau^{1/2}} \left(\frac{M_3}{\sqrt{B}} \right) \cdot l \quad (58)$$

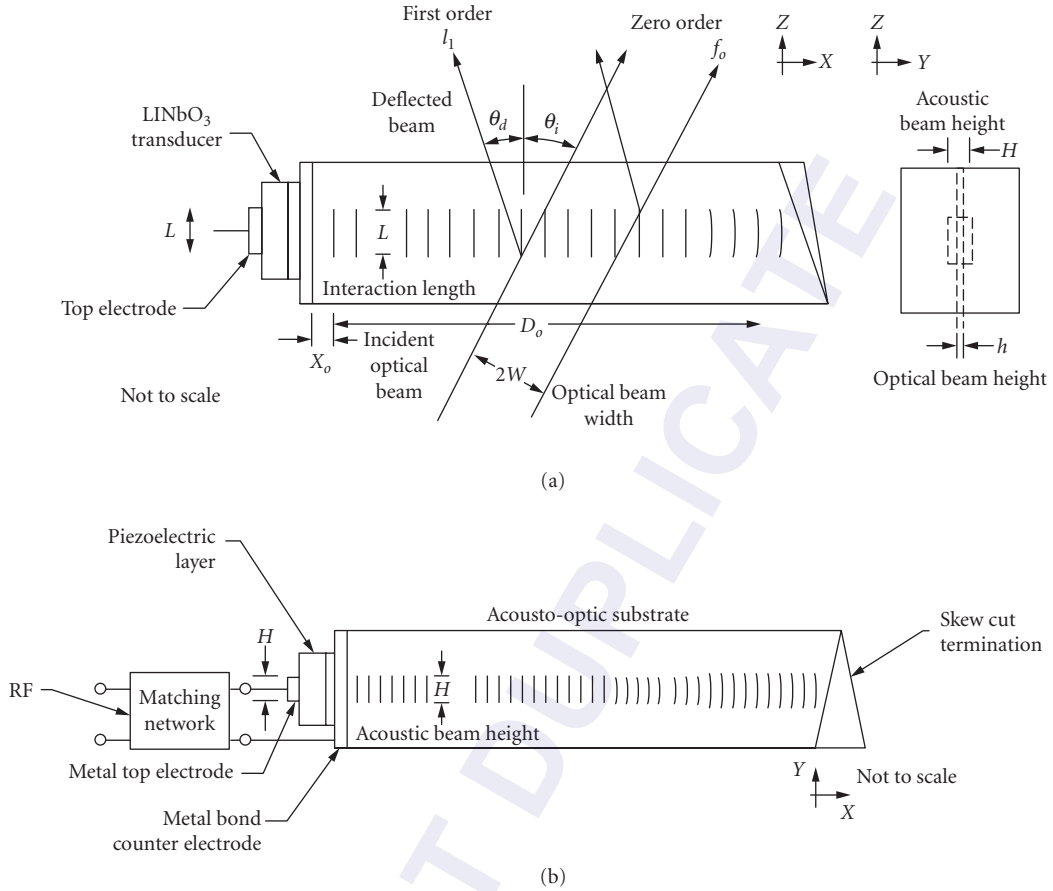


FIGURE 3 Acousto-optic deflector configuration: (a) interaction plane and (b) transverse plane.

Equation (58) is the basic equation for the design of the deflector. The multiplying factor J/h in Eq. (58) depends on the the optical beam geometry and the transducer electrode pattern in the transverse plane. For example, for an optical Gaussian beam with a truncation ratio $D/2\omega_1 = 1.5$, and rectangular electrode, the value of J/h exhibits a broad maximum about unity when $h \approx 0.7$. The choice of normalized interaction length $l = L/L_o$ is determined by the specified bandwidth. For isotropic diffraction it is related to the fractional bandwidth by, from Eq. (41):

$$l = \frac{1.8}{\Delta F} \tag{59}$$

With the selection of l , all the key performance parameters can be determined. The frequency response and spatial frequency distribution are obtained from Eqs. (48) and (49).

Anisotropic acoustic beam collimation Equation (58) shows that increased diffraction efficiency can be obtained by selecting AO materials and modes with small slowness curvature B . This is referred to as anisotropic acoustic beam collimation. An example is the shear mode in GaP propagating

along [110] direction.²⁶ In this case $B = 0.026$. Compared to an acoustically isotropic direction, the transducer height can be reduced by a factor of 6.2.

Resolution The resolution of the AO deflector obtainable is limited by the maximum available crystal size, requirement of large optics, and acoustic attenuation across the aperture. For most crystalline solids, the acoustic attenuation is proportional to f^2 . If we allow a change of average attenuation of ξ (dB) across the band, the maximum deflector resolution is given by

$$N_{\max} = \frac{\xi}{\alpha_o f_o} \quad (60)$$

where α_o (dB/ μ sec-GHz²) is the acoustic attenuation coefficient.

It is instructive to estimate the maximum resolution achievable of AO deflectors. Consider, for example, a slow shear wave TeO₂ AO deflector with the following parameters: the octave bandwidth $f_o = 1.5 \Delta f$, acoustic velocity $V = 0.62$ mm/ μ sec, acoustic attenuation coefficient $\alpha_o = 17.9$ (dB/ μ sec-GHz²), maximum acoustic loss $\xi = 4$ dB, aperture size $D = 5$ cm. From Eq. (60), the maximum resolution N_{\max} is estimated to be about 3000 resolvable spots. This is close to the limit of a practical single stage AO deflector.

Birefringent tangential phase matching The maximum interaction length is limited by the bandwidth requirement. One approach of increasing while maintaining the fractional bandwidth ΔF is the use of birefringent diffraction. It was shown in Sec. 6.3 that by choosing the birefringent AO interaction near the tangential phase matching (TPM), a large band of acoustic frequencies will simultaneously satisfy the momentum matching condition. Equivalently a larger l can be used without narrowing the bandwidth. Under the birefringent TPM condition l is related to the fractional bandwidth by, from Eq. (42):

$$l = \frac{3.6}{\Delta F^2} \quad (61)$$

Compared to isotropic diffraction, the birefringent phase matching achieves an efficiency advantage factor of $2/\Delta F$, which becomes particularly significant for smaller fractional bandwidths. For instance, at 50 percent fractional bandwidth, the birefringent phase matching has an advantage factor of 4.

Since the first report on the use of birefringent TPM,²⁷ this performance enhancement technique has been used in the design of high resolution AO deflectors. All of these TPM deflectors are based on the birefringent diffraction in TeO₂ by a slow shear wave propagating on or near the [110] axis. Because of the exceptional low velocity and large figure of merit, these TeO₂ deflectors can achieve a large number of resolvable spots with a relatively low drive power.

The birefringent TPM deflectors can be classified as on-axis, optically rotated (OR) and acoustically rotated (AR) types. The design of these high resolution (HR) AO deflectors shown are briefly discussed as follows.

ON-axis TPM The first TPM type of AO deflector uses a slow shear wave propagating along the [110] axis of the TeO₂ crystal. This is referred as 90° or on-axis design, since the acoustic angle θ_a is chosen to be equal to 90°. For lowering TPM frequency, the AO diffraction birefringent diffraction uses circular birefringence in the TeO₂ crystal along the c axis.²⁸ Due to the slow velocity $V = 0.62$ mm/ μ sec, large resolution is obtainable with a relatively small size crystal. For instance, a 70- μ sec cell realize more than 2000 resolvable spots with an optical aperture size of only 4.3 cm.

The on-axis design has several drawbacks. First, operation of this mode requires the use of circular polarizers. Second, the tangential frequency f_t is fixed and thus other performance parameters such as figure of merit cannot be independently chosen. Third, at the tangential phase-matching frequency, the second-order diffraction is also matched. The first-order diffracted light is rediffracted into the second order and results in a dip at the middle of the deflector passband.

To resolve this, two types of birefringent phase matching (TPM) are proposed. These include optically rotated (OR)²⁵ and acoustically rotated (AR)²⁹ TPM. For both types of phase-matching schemes the tangential matching frequency, $f_t = F_t f_o$, is given by, from Eq. (35):

$$f_t \approx \frac{V(\theta_a)}{\lambda_o} \sqrt{2n_o \Delta n \sin \theta_e} \quad (62)$$

Figure 4 shows the wavevector diagram for the two types of tangential phase matching.

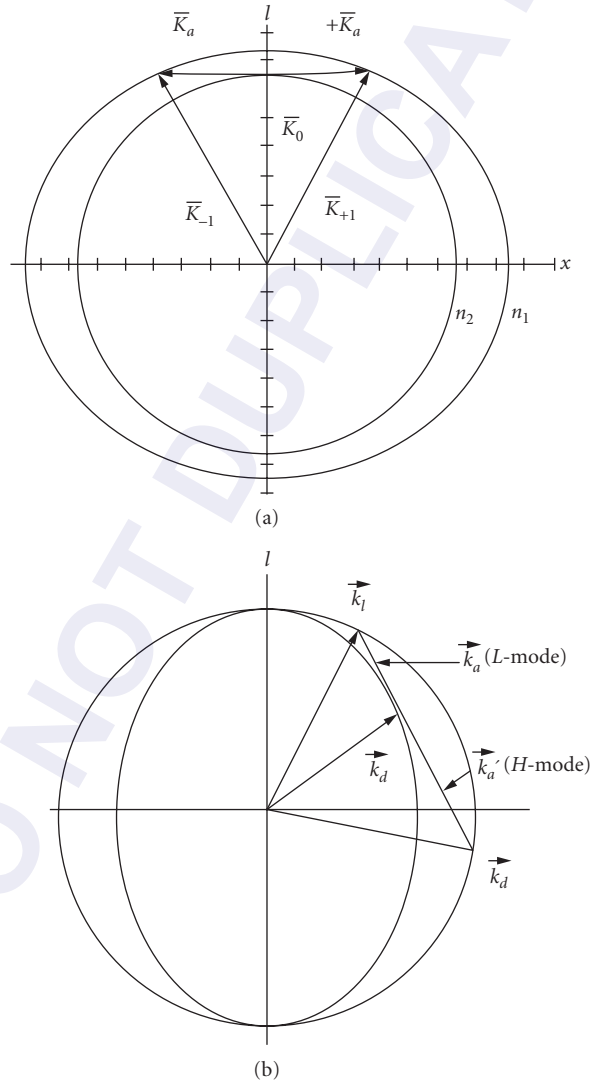


FIGURE 4 Wavevector diagram for tangential phase matching: (a) optically rotated geometry and (b) acoustically rotated geometry.

Optically rotated TPM The wavevector diagram shown in Fig. 4a describes the OR tangential phase matching. In this case the acoustic wavevector is chosen to be perpendicular to the optic axis. The incident light wavevector is allowed to rotate in different polar angles to achieve TPM at the selected frequency. Operation of this cell is more convenient than the on-axis device since the optical beam is linearly polarized. The OR type can realize long time aperture and higher efficiency since the acoustic wave is propagating along the isotropic direction. However, rediffraction into second order occurs since optical locus is symmetrical to the acoustic wave direction, the peak maximum efficiency remains limited by the dip in the center of the frequency response.

Acoustically rotated TPM Figure 4b shows the wavevector diagram for AR type of tangential phase matching. The constant azimuth plane is chosen as the interaction plane and the acoustic wavevector is rotated in the plane to be tangential to the locus of the diffracted light wavevector. Because of the rotation of the acoustic wavevector, the asymmetry of optical wavevector locus removes the degeneracy of the second-order diffraction. The dip in the bandpass is eliminated and thus can operate at the high-efficiency regime. This is the major advantage of the AR type cell. Since the acoustic wave is propagating along an anisotropic path, the large walk-off will reduce the optical aperture and thus limits the maximum resolution. A higher TPM frequency is chosen in order to achieve 2000 resolvable spots.

Phase array beam steering The birefringent type of tangential phase-matching technique is limited to AO diffraction in birefringent crystals. For isotropic materials such as wideband GaP Bragg cells an alternate technique for realizing TPM is to use acoustic beam steering.⁷ In this approach, a phase array of transducers is used so that the composite acoustic wavefront will effectively track the Bragg angle. The simplest phase array employs a fixed inter-relevant phase difference that corresponds to an acoustic delay of $PA/2$, where P is an integer. This is referred to as first-order beam steering and can be realized in either a stepped array^{7,30} or a planar configuration.³¹ The two types of phased array configuration are shown in Fig. 5. The stepped array configuration is more efficient, but fabrication difficulty makes it less practical. The following analysis addresses only the planar configuration of first-order beam steering.

Consider the simplest geometry of a planar first-order beam-steered transducer array where each element is driven with an inter-element phase difference of 180° . The acousto-optic bandpass response of this transmitter configuration is equal to the single-element bandshape multiplied by the interference (array) function,²⁵

$$W(F) = \left(\frac{\sin \pi X}{\pi X} \right)^2 \left(\frac{\sin N \pi Y}{N \sin \pi Y} \right)^2 \quad (63)$$

$$X = \frac{l_e F}{2} (F_b - F) \quad (64)$$

$$Y = \frac{d}{2} \left\{ F(F_b - F) + \frac{1}{d} \right\} \quad (65)$$

where $l_e = L_e/L_o$, $d = D/L_o$, L_e is the length of one element, N is the number of elements, and D is the center-to-center distance between adjacent elements. Since radiation pattern for a single element is broad, the bandpass function is primarily determined by the array functions. The bandpass characteristics of the phased array are similar to birefringent diffraction with an equivalent interaction length $\ell = Nd$ and tangential matching frequency $F_t = 1/\sqrt{d}$.

Referring to Eq. (59), notice that at the peak of the grating lobe for the phase-array radiation, the value of the single element radiation $\text{sinc}^2 x$ is approximately equal to 0.5. There is thus an additional 3 dB loss due to the planar phase array. This is the major disadvantage of the phase array approach.

Since the power density is proportional to ℓ^2 for both types of TPM techniques, the increase of ℓ will significantly reduce the power density. This is of practical importance, since high power density

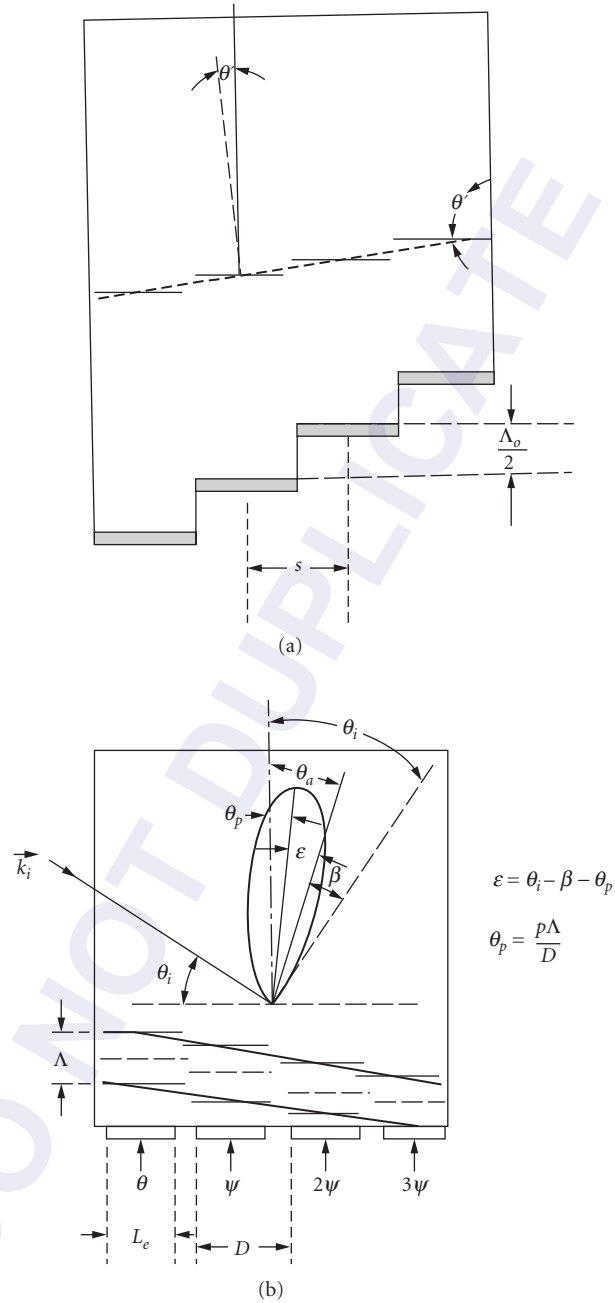


FIGURE 5 Acoustic beam steering using phased array: (a) stepped phased array and (b) planar phased array.

has been the dominant factor for performance degradation and even catastrophic device failure due to thermal gradient, nonlinear acoustics and overheating of transducers.

Birefringent phased array An interesting design combining the two performance enhancement techniques was reported in birefringent AO device using phased-array transducers.³² Theoretical analysis of the new interaction geometry was described and experimentally demonstrated in a wide-band Bragg cell. The new approach provides design flexibility in the choice of acoustic and optic modes for increased efficiency, bandwidth, suppression of nonlinear effects, and optimizing other performance characteristics. The approach was further explored in the design of AO deflectors for specific performance improvement. By applying the birefringent phased array design to on-axis,³³ OR and AR configuration,³⁴ some of the deficiencies in each configuration can be removed. The design flexibility of this approach was also utilized to improve other key performance features, including, for example, the development of wide angle AO Bragg cells.³⁵

AO Laser Beam Scanner The early development of AO deflectors was aimed at laser beam scanning applications. The primary goal is to provide a simple solid-state laser scanner that eliminates the inherent drawbacks of mechanical scanners due to moving parts, such as facet errors and the requirement of realignment because of bearing wear. For certain applications, such as beam-addressed optical memory, the rapid random access capability of AO deflectors offers a distinct advantage.

High Resolution AO Deflector The AO device for laser beam scanning is referred as the high resolution (HR) or type I AO deflector. As shown in Table 3, the primary design objective is to realize high resolution with the choice of a long-time aperture. Another key performance specification is high throughput or peak diffraction efficiency. Table 4 lists the performance of a few representative HR AO deflectors operated in the visible range. To satisfy the high peak efficiency specification, the degradation due to in-band dip must be minimized if OR type TPM rediffraction at high efficiency.

Resolution and Scan Rate In the design of the AO deflector for linear scanning, the most important system specification is the resolution, or maximum number of resolvable spots, N . Another key system parameter to specify is the scan time required to scan a single line and the flyback time. The acoustic transit time across the optical aperture sets the minimum flyback time for the scanner. The finite scan rate also degrades the resolution of the AO deflector. When the acoustic transit time τ (which is equal to flyback time) becomes an appreciable portion of the scan time T , there is a reduction of the effective aperture by the factor $1 - \tau/T$. This results in a resolvable number of spots for the scanning mode

$$N = \tau \Delta f [1 - \tau/(T + \tau)] \quad (66)$$

To minimize the loss of resolution, the transit time τ should be much less than the scan time T .

TABLE 4 Performance of High Resolution AO Deflectors

Material (Mode)	Center Frequency (MHz)	Bandwidth (MHz)	Aperture (μ s)	TB Product	Efficiency (%/W)
TeO ₂ (S)	50	30	70	2100	200
TeO ₂ (S)	90	50	40	2000	110
TeO ₂ (S)	160	100	20	2000	95
GaP(S)	300	200	2.5	500	50

Resolution Enhancement One early design for increased resolution involves the use of cascaded deflectors.³⁶ Another technique is to use higher-order AO diffraction. One interesting design was to utilize the second-order diffraction in an OR-type birefringent cell.³⁷ Since both the first- and second-order diffractions are degenerately phase-matched, efficient rediffraction of the first order into the second order was obtained. However, the use of the second-order diffraction allows the deflector resolution to be doubled for a given bandwidth and time aperture.

AO Signal Processing Bragg Cell The deflection of an incident optical wave according to a single RF frequency input is just one of the characteristics of the acoustically driven transverse spatial modulator (TSM). More generally, the TSM encodes the light beam with the complete spectrum information contained in the RF signal. An AO deflector simultaneously driven by multifrequencies provides a simple but powerful technique for processing wideband electronic signals. This multifrequency AO device used for signal processing is referred to as the AO Bragg cell.

Dynamic Range The design of Bragg cells is similar to that of the AO deflector. However, instead of optimization of the peak diffraction efficiency, the primary objective of the Bragg cell design is to achieve a large dynamic range rather than the peak diffraction. The dynamic range is defined as the ratio of the intensity of the diffracted light to the intensity of the spurious light in the Bragg cell. The spurious optical beams are caused by the nonlinear mixing of various diffraction orders when multiple frequencies are simultaneously present. Hecht³⁸ analyzed isotropic AO diffraction with multiple acoustic waves with different carrier frequencies based on the coupled wave theory. His analysis shows that the dominant in-band nonlinear signals for two simultaneous signals at f_1 and f_2 are the third-order intermodulation (IM) products occurring at $2f_1 - f_2$ and $2f_2 - f_1$. The magnitude of this spurious signal due to multiple AO diffraction is

$$I_{\text{AO}}(2, -1) = \frac{I_1^2 I_2}{36} \quad (67)$$

where I_1 and I_2 are the diffracted light efficiency at f_1 and f_2 , respectively. To obtain a spurious-free dynamic range of 40 dB, for instance, the Bragg cell must be operated with diffraction efficiency less than 5 percent.

In our experimental work on wideband Bragg cells, we have found another type of IM product in the nonlinear response of Bragg cells.³⁹ This second source of IM products is caused by the nonlinear acoustic (NA) interaction that occurs at high acoustic power densities. A key parameter for the NA type IM products process is the critical interaction length $L_c = (\pi \gamma_2 f \tau)$, where γ_2 is the ratio of second-order nonlinear elastic coefficient to the first-order elastic coefficient. The dominant NA modes are second harmonics and two-tone third-order IM product, which grow according to $(L_c/L)^2 P_1^2/4$ and $(L_c/L)^4 P_1^2 P_2/9$, respectively. Based on the new theory, the AO type of IM products dominates when the ratio L_c/L is small, as the acoustic power increases. The dynamic range of the Bragg cell is reduced by the NA type of IM product, initially by the successive AO diffraction from fundamental and second harmonics and finally by the single AO diffraction from the acoustic IM product.

Wideband Bragg Cell Two candidate materials for wideband Bragg cells are GaP and LiNbO₃. Both have very low acoustic attenuation and large M_3^* , making them ideal for this purpose. A number of high performance wideband cells using these two materials have been developed with bandwidth in excess of 1000 MHz. These include the AR-type birefringent LiNbO₃ cell⁴⁰ and the phased-array GaP cell.⁴¹ The best-known design is a LiNbO₃ device, which uses a γ - z propagating off-axis x -polarized shear wave. The device demonstrated an overall bandwidth of 2 GHz, a peak diffraction efficiency of 6 percent per watt and about 600 resolvable spots. A summary of the representative wideband AO Bragg cells is shown in Table 5.

Multichannel Cell As an extension to the usual single-channel configuration there has been considerable activity in the development of multichannel Bragg cells (MCBC). The MCBC uses a pattern

TABLE 5 Performance of Acousto-Optic Bragg Cells ($\lambda_o = 830$ nm)

Material (Mode)	Center Frequency (MHz)	Bandwidth (MHz)	Aperture (μ s)	TB Product	Efficiency (%/W)
GaP(S)	1000	500	2.0	1000	30
LiNbO ₃ (S)	2500	1000	1.0	1000	10
GaP(L)	2500	1000	0.25	250	44
GaP(S)	3000	2000	0.25	500	8
LiNbO ₃ (S)	3000	2000	0.30	600	6

of multiple individually addressed transducer electrodes in the transverse plane. In addition to the design rules used for single-channel cells, other performance parameters such as crosstalk, amplitude, and phase tracking must be considered.

Early work on MCBC uses the y - z cut shear-wave wideband LiNbO₃ cell design.⁴² Good amplitude and phase tracking were obtained over the operating bandwidth of 1 GHz. The channel-to-channel isolation, typically about 25 dB, was limited by RF cross talk. In a later development of GaP MCBC,⁴³ reduction of electrical crosstalk (to -40 dB) was obtained by using stripline interconnection structures. The use of the anisotropic collimating modes in GaP also brings the advantage of lowering acoustic crosstalk.

6.6 ACOUSTO-OPTIC MODULATOR

Acousto-optic interaction has also been used to modulate light. The AO modulator uses the transverse spatial modulation (TSM) of AO interaction to convert an amplitude modulated RF signal into an intensity modulation light beam. Unlike AO deflectors, however, at the output end the diffracted light spectral components are focused into one channel to reproduce the video information via a collinear mixing process. Functionally the most important requirement for the AO modulators is its temporal response. Thus a basic design consideration for AO modulator is the limitation of temporal response caused by the finite acoustic transit time across the light beam. To reduce the transit time effect on the fidelity of the input video information, the AO modulator is operated in a focused beam geometry.

Following similar design procedure as the deflectors, we shall describe the characteristics of AO modulator by the basic parameters. As a modulator the key performance parameters are rise time or video modulation bandwidth, contrast ratio or rejection against the undeflected light beam, and the optical throughput, the diffracted light intensity integrated over the acceptance angular aperture. A thorough discussion on the tradeoff of key performance parameters of AO modulators is reviewed by Johnson.⁴⁴

Based on the optics configuration AO modulators can be divided into type-I focused beam modulator and type II-wide beam modulator. Table 6 lists the interaction geometry and key performance requirement for the two types of AO modulators.

TABLE 6 Interaction geometry and key performance parameters of AO modulators

Type I: Focused beam modulator
Regime of AO diffraction: Transverse SLM
Divergence ratio/Resolution: $A = 1, N = 2$
Primary performance: Wide modulation bandwidth, large optical throughput
Type II: Wide beam modulator
Regime of AO diffraction: Transverse SLM
Divergence ratio: $A \ll 1, N \gg 1$
Prime performance: High peak efficiency, high rejection ratio

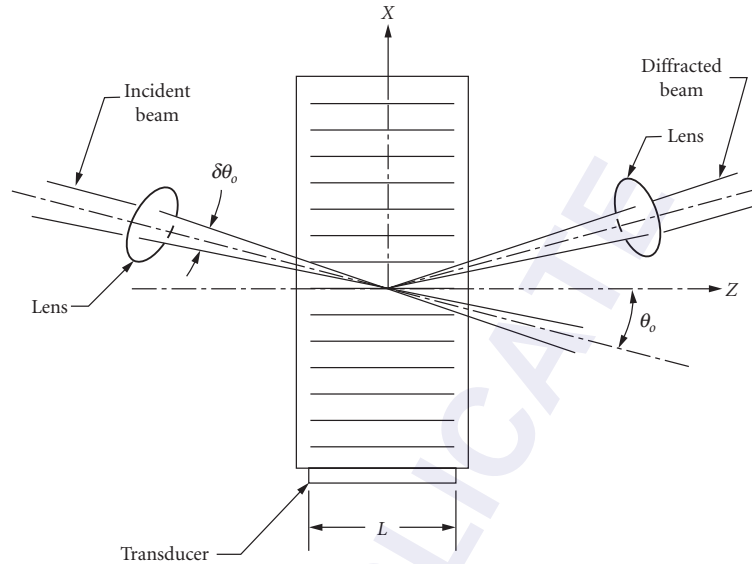


FIGURE 6 Diffraction geometry of acousto-optic modulator.

Principle of Operation Figure 6 shows the diffraction geometry of a focused-beam AO modulator. Unlike the AO deflector the optical beam is focused in both directions into the interaction region near the transducer. For an incident Gaussian beam with a beam waist d , ideally the rise time (10 to 90 percent) of the AO modulator response to a step function input pulse is given by

$$t_r = 0.64\tau \quad (68)$$

where $\tau = d/V$ is the acoustic transit time across the optical beam. To reduce rise time, the optical beam is focused to a spot size as small as possible. However, focusing of the optical beam will increase the optical beam divergence $\delta\theta_o$, which may exceed the acoustic beam divergence $\delta\theta_a$, $A > 1$. This will result in a decrease of the diffracted light since the Bragg condition will no longer be satisfied. To make a compromise between the frequency response (spatial frequency bandwidth) and temporal response (rise time or modulation bandwidth) the optical and acoustic divergence should be approximately equal, $A = \delta\theta_o/\delta\theta_a \approx 1$. The actual value of the divergence ratio depends on the tradeoff between key performance parameters as dictated by the specific application.

Analog Modulation In the following, we shall consider the design of a focused beam AO modulator for analog modulation, the diffraction of an incident Gaussian beam by an amplitude-modulated (AM) acoustic wave. The carrier, upper, and lower sidebands of the AM acoustic wave will generate three correspondingly diffracted light waves traveling in separate directions. The modulated light intensity is determined by the overlapping collinear heterodyning of the diffracted optical carrier beam and the two sidebands. Using the frequency domain analysis, the diffracted light amplitudes of an AO modulator were calculated. The numerical results are summarized in the plot of modulation bandwidth and peak diffracted light intensity as a function of the optical to acoustic divergence ratio A .¹³ Design of the AO modulator based on the choice of the divergence ratio is discussed below.

Unity Divergence Ratio ($A \leq 1$) The characteristics of the AO modulator can be best described by the modulation transfer function (MTF), defined as the frequency domain response to a sinusoidal video signal. In the limit of $A \leq 1$, the MTF takes a simple form:

$$\text{MTF}(f) = \exp - (\pi f \tau)^2/8 \quad (69)$$

The modulation bandwidth f_m , the frequency at -3 dB is given by

$$f_m = \frac{0.75}{\tau} \quad (70)$$

From Eqs. (62) and (64), the modulator rise time and the modulation bandwidth are related by

$$f_m t_r = 0.48 \quad (71)$$

Minimum Profile Distortion ($0.5 < A < 0.67$) Equation (70) shows that the modulation bandwidth can be increased by further reducing the acoustic transit time τ . When the optical divergence $\delta\theta_o$ exceeds the acoustic divergence $\delta\theta_a$, that is, $A > 1$, the Bragg condition is no longer satisfied for all the optical plane waves. The light at the edges of the incident light beam will not be diffracted. This will result in an elliptically shaped diffracted beam. In many laser modulation systems this distortion of the optical beam is not acceptable. The effect of the parameter A on the eccentricity of the diffracted beam was analyzed based on numerical calculation.⁴⁵ The result shows that to limit the eccentricity to less than 10 percent, the divergence ratio value for A is about 0.67.

This distortion of the diffracted beam profile is caused by the finite acceptance angle of the isotropic AO diffraction. The limited angular aperture also results in the lowering the optical throughput. Based on curve fitting the numerical results the peak optical throughput can be expressed as a function of the divergence ratio⁴⁴

$$I_1(A) = 1 - 0.211A^2 + 0.026A^4 \quad (72)$$

As an example, at $A = 0.67$, the peak throughput I_1 is about 94.9 percent. However, with the choice of unity divergence ratio, $A = 1$, it reduces to 81.5 percent.

Digital Modulation ($1.5 < A < 2$) Another important case is the digital, or pulse modulation. Maydan⁴⁶ calculated the rise time and efficiency of pulsed AO modulators. His results show that an optimized choice of A is equal to 1.5, and that the corresponding rise time (10 to 90 percent) is

$$t_r = 0.85\tau \quad (73)$$

Contrast Ratio Another key performance parameter for AO modulator is the extinction ratio or the rejection against the undeflected light. To obtain an adequate extinction ratio, 30 dB, for instance, the angle of separation θ_D is chosen to be equal to twice that of the optical beam divergence. It follows that the acoustic frequency must be greater than

$$f_o = \frac{8}{\pi\tau} \quad (74)$$

Comparing to Eq. (70), the center frequency of the AO modulator is about 4 times that of the modulation frequency.

Birefringent Modulation In the above analysis of focused beam modulator it is assumed that the AO modulator is operated in the isotropic diffraction mode. Similar to the deflector case a large frequency bandwidth can be obtained by using birefringent diffraction. However, the input angular aperture according to Eq. (55) is proportional to $1/L$, an increase of L will in effect reduce the modulation bandwidth as well as the optical throughput. Based on this argument, it was well recognized that the birefringent diffraction offers no advantage for AO modulators. A careful reexamination of this basic limitation shows that a wide angle AO modulation is obtainable if the birefringent diffraction satisfies the noncritical phase matching (NPM) condition.⁴⁸ The concept of NPM will be discussed in Sec. 6.7, "Acousto-Optic Tunable Filter."

TABLE 7 Performance of Acousto-Optic Modulators

Material	Wavelength (μm)	Center Frequency (MHz)	Modulation Bandwidth (MHz)	Rise Time (nsec)	Efficiency (%)
TeO ₂	0.44–0.64	200	40	10	80
GaP	0.63–0.83	200	40	10	80
GaP	1.06	160	32	15	70
GaAs	1.3–1.55	120	24	20	70
Ge	10.6	100	20	25	70

Focused Beam AO Modulator

Intensity modulator Simple intensity AO modulators using the focused beam interaction geometry have been standard approach for external modulation of lasers. Since the development of laser diode with wideband internal modulation capability, the use of AO modulators has been limited to gas and solid-state lasers for moderate modulation bandwidth f_m . Because of the high power requirement, in practice f_m of AO modulators is limited to about 50 MHz. Table 7 lists a few selected AO modulators and typical performances.

Intracavity modulator Since AO diffraction occurs in all crystal or amorphous solids, high optical quality, low optical absorption AO materials are readily available for intracavity applications. These intracavity modulator applications include Q-switching, mode locker, and cavity dumping. All of these intracavity AO modulator have been discussed in the early review paper.¹³ In the following we shall only briefly discuss the AO Q-switches.

Q-switching of YAG and other similar solid state lasers has been an importance requirement for industrial applications such as cutting, scribing, marking, and other material processing process. In a Q-switching operation, the optical loss introduced by the intracavity AO modulator keeps the laser below threshold. When the loss is suddenly removed by switching off the acoustic pulse, the laser bursts into a short pulse with extremely high intensity. During Q-switching, the AO modulator should add minimum loss to the laser cavity.

Besides low loss, other key performance requirement of the AO Q-switch include (a) high damage threshold, (b) fast pulse rise time, (c) good thermal and mechanical property, and (d) polarization-insensitive. The last requirement is due to a special characteristics of YAG lasers that it reaches maximum gain when operated in a unpolarized mode. To meet these requirement the standard design for high power laser is a shear wave AO modulator using UV grade fused silica as the interaction medium. Longitudinal mode Q-switch will have the advantage of lower drive power and faster rise time; however, since the diffracted efficiency for perpendicular polarization is five times that of the parallel polarization. the TAG laser will tend to operate in the mode with the lower diffraction efficiency. Thus the use of longitudinal mode Q-switch is limited to lower power polarized laser. More efficient materials such as PbMoO₄ or TeO₂ have been used but because of lower damage threshold these crystal Q-switch is also limited to low power lasers.

Widebeam AO Modulator

Image (Scophony) AO Modulator The focused-beam-type AO modulator has certain disadvantages. The diffraction spread associated with the narrow optical beam tends to lower the diffraction efficiency. More importantly, the focusing of the incident beam results in high peak intensity that can cause optical damage for even relatively low laser power levels. For these reasons, it is desirable to open up the optical aperture. Due to the basic issue of acoustic transit time, the temporal bandwidth of the wide-beam AO modulator will be severely degraded.

In certain applications, such as the laser display system, it is possible to use a much broader optical beam in the modulator than that which would be allowed by the transit time limitation. The operation of the wide-beam modulator is based on the ingenious technique of Scophony light modulation. A brief description of the Scophony light modulator is given below; refer to Johnson⁴⁷ for a thorough treatment.

The Scophony technique is applicable to any system that scans a line at a uniform scan velocity. The basic idea is to illuminate a number of picture elements, or pixels, in the modulator (window) onto the output line, such that the moving video signal in the modulator produces a corresponding image of the pixels that travels across the beam at sound velocity. The image can be made stationary by directing the image through a deflector that scans with equal and opposite velocity. Now, if the window contains N spots, then N picture elements can be simultaneously exposed in the image at any instant, and each picture element will be built up over the access time of the window that is equal to N times the spot time. Since the spots are immobilized, there is no loss of resolution in the image, provided that the modulator bandwidth is sufficient to produce the required resolution. The design of the wide-beam AO modulator is thus the same as that of the Bragg cell.

Acousto-Optic Frequency Shifters Because of the frequency shift associated with the acousto-optic diffraction, the AO frequency shifter (AOFS) provides a capability of shifting the optical frequency of an incident light beam, a unique ability shared by no other technique. Usually the frequency shift specified is either fixed or tunable over a small frequency range. Since the AOFS is often intended for use in interferometry, an important requirement is to minimize the distortion of diffracted light beam profile. For isotropic AO diffraction the use of an input optics design with low divergence ratio $A \ll 1$ is the preferred choice.

Another important requirement is the rejection of incident light, other diffraction order with opposite frequency shift, as well as various spurious signals caused by multiple acoustic reflections within the AO medium. To achieve high suppression of the residual AM components and obtain the ideal spectral purity, the acoustic wave reflection and diffraction in the medium must be considered in the device design and fabrications of the AOFS. Similar to focused beam AO modulator the wide angle design using birefringent phase array can provide large degree of freedom in the optimization of performance parameters such as the optical throughput. The use of crossed polarizers also provides the advantage of increased rejection ratio.

In the design of AOFS, the desired frequency shift f_s is generally chosen to be near the center frequency f_o . This constraint is obviously impractical for small frequency shift $f_s \ll f_o$. The AOFS in this case is typically constructed by putting two copropagating AO cells in cascade so that the two frequencies subtract to give the small frequency shift. The cascaded AO cells configuration can also be used to obtain variable frequency shifts without change of direction. For example, two counterpropagating AO cells can be configured so that the two frequencies are added while the directions cancel each other.

6.7 ACOUSTO-OPTIC TUNABLE FILTER

The acousto-optic tunable filter (AOTF) is an all-solid-state optical filter that operates on the principle of acousto-optic diffraction in an anisotropic medium. The center wavelength of the filter passband can be rapidly tuned across a wide spectral range by changing the frequency of the applied radio frequency (RF) signal. In addition to the electronic tunability, other outstanding features of the AOTF include large angular aperture while maintaining high spectral resolution, inherent intensity, and wavelength modulation capability.

The first AOTF, proposed by Harris and Wallace⁹ used a configuration in which the interacting optical and acoustic waves were collinear. Later, Chang¹⁰ extended the AOTF concept to a noncollinear configuration. The theory and practice of the AOTF have been discussed in a review paper.⁴⁹

Based on the viewpoint of device physics, the AOTF is a new AO device of a fundamentally different nature compared to AO deflectors and modulators. It operates on longitudinal spatial light modulation (SLM) that occurs only in an anisotropic medium. In accordance with the interaction geometry

TABLE 8 Interaction Geometry and Key Performance Goals of AO Tunable Filters

Regime of AO diffraction: Longitudinal SLM
Divergence ratio: $a \gg 1, R \gg 1$
Primary performance: Large angular aperture, high optical throughput
Type II Critical phase matching (CPM) AOTF
Regime of AO diffraction: Longitudinal SPM
Primary performance requirement: High resolution, low drive power

and primary performance requirement the AOTF can be divided into two types. Type I noncritical phase matching (NPM) AOTF exhibits large angular aperture and high optical throughput, is best suited for spectral imaging applications. Type II critical phase matching (CPF) AOTF emphasizes high spectral resolution and low drive power, has shown to be best suited for use as a dynamic-wavelength division multiplexing (WDM) component. Table 8 shows the range of divergence ratio and key performance requirements for the two types of AOTFs.

Principle of Operation

Collinear acousto-optic tunable filter Consider the collinear AO interaction in a birefringent crystal, a linearly polarized light beam will be diffracted into the orthogonal polarization if the momentum matching condition is satisfied. For a given acoustic frequency the diffracted light beam contains only for a small band of optical frequencies centered at the passband wavelength,

$$\lambda_o = \frac{V \Delta n}{f} \quad (75)$$

where Δn is the birefringence. Equation (75) shows that the passband wavelength of the filter can be tuned simply by changing the frequency of the RF signal. To separate the filtered light from the broadband incident light, a pair of orthogonal polarizers is used at the input and output of the filter. The spectral resolution of the collinear AOTF is given by

$$R = \frac{\Delta n L}{\lambda_o} \quad (76)$$

A significant feature of this type of electronically tunable filter is that the spectral resolution is maintained over a relatively large angular distribution of incident light. The total angular aperture (outside the medium) is

$$\Delta \psi = 2n \sqrt{\frac{\lambda_o}{\Delta n L}} \quad (77)$$

where n is the refractive index for the incident light wave.

This unique capability of collinear AOTF for obtaining high spectral resolution within a large angular aperture was experimentally demonstrated using a transmissive-type configuration shown in Fig. 7.⁵⁰ A longitudinal wave is mode converted at the prism interface into a shear wave that propagates down along the x axis of the CaMoO_4 crystal. An incident light with a wavelength satisfying Eq. (75) is diffracted into the orthogonal polarization and is coupled out by the output polarizer. The center wavelength of the filter passband was changed from 400 to 700 nm when the RF frequency was tuned from 114 to 38 MHz. The full width at half-maximum (FWHM) of the filter passband was measured to be 8 Å with an input light cone angle of $\pm 4.8^\circ$ ($F/6$). This angular aperture is more than one order of magnitude larger than that of a grating for the same spectral resolution.

Considering the potential of making electronically driven rapid-scan spectrometers, a dedicated effort was initiated to develop collinear AOTFs for the UV using quartz as the filter medium.⁵¹ Because of the collinear structure quartz AOTF demonstrated very high spectral resolution. With a 5-cm-long crystal, a filter

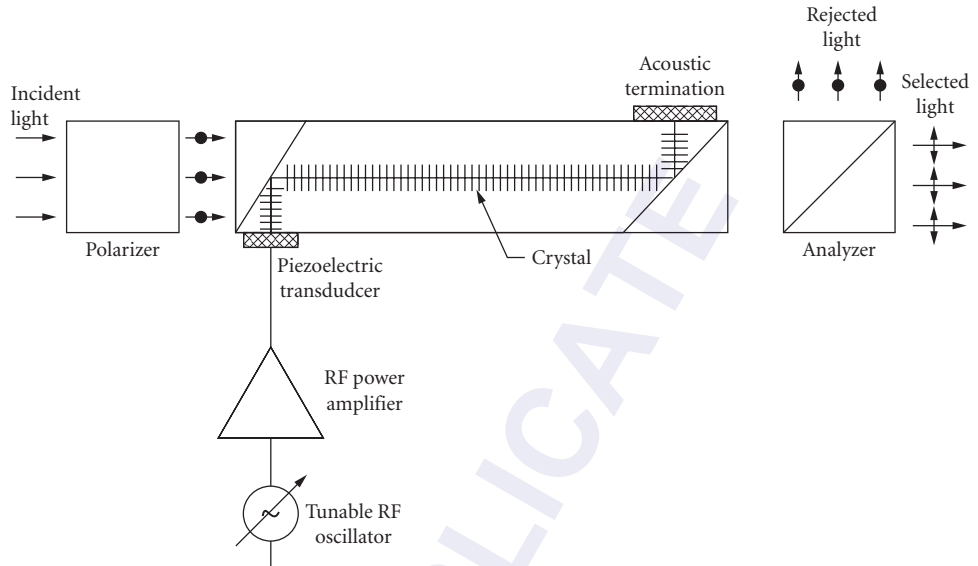


FIGURE 7 Collinear acousto-optic tunable filter with transmissive configuration.

bandwidth of 0.39 nm was obtained at 250 nm. One limitation of this mode conversion type configuration is the complicated fabrication procedures in the filter construction. To resolve this issue, a simpler configuration using acoustic beam walk-off was proposed and demonstrated.⁵² The walk-off AOTF allows the use of multiple transducers and thus could realize a wide tuning range. Experimentally, the passband wavelength was tunable from about 250 to 650 nm by changing the acoustic frequency from 174 to 54 MHz. The simple structure of the walk-off filter is particularly attractive for manufacturing.

Noncollinear AOTF The collinearity requirement limits the AOTF materials to rather restricted classes of crystals. Some of the most efficient AO materials (e.g., TeO_2) are excluded since the pertinent elasto-optic coefficient for collinear AO interaction is zero. Early work has demonstrated a noncollinear TeO_2 AOTF operation using a S[110] on-axis design. However, since the phase-matching condition of the noncollinear AO interaction is critically dependent on the direction of the incident angle of the light beam, this type of filter has a very small angular aperture (on the order of milliradians), and its use must be restricted to well-collimated light sources. To overcome this deficiency, a new method was proposed to obtain a large-angle filter operation in a noncollinear configuration.

The basic concept of the noncollinear AOTF is shown in the wavevector diagram in Fig. 8. The acoustic wavevector is chosen so that the tangents to the incident and diffracted light wavevector loci are parallel. When the parallel tangents condition is met, the phase mismatch due to the change of angle incidence is compensated for by the angular change of birefringence. The AO diffraction thus becomes relatively insensitive to the angle of light incidence, a process referred to as the noncritical phase-matching (NPM) condition. The figure also shows the wavevector diagram for the collinear AOTF as a special case of noncritical phase-matching.

Figure 9 shows the schematic of a noncollinear acousto-optic tunable filter. The first experimental demonstration of the noncollinear AOTF was reported for the visible spectral region using TeO_2 as the filter medium. The filter had a FWHM of 4 nm at an $F/6$ aperture. The center wavelength is tunable from 700 to 450 nm as the RF frequency is changed from 100 to 180 MHz. Nearly 100 percent of the incident light is diffracted with a drive power of 120 mW. The filtered beam is separated from the incident beam with an angle of about 6° . The experimental result is in agreement

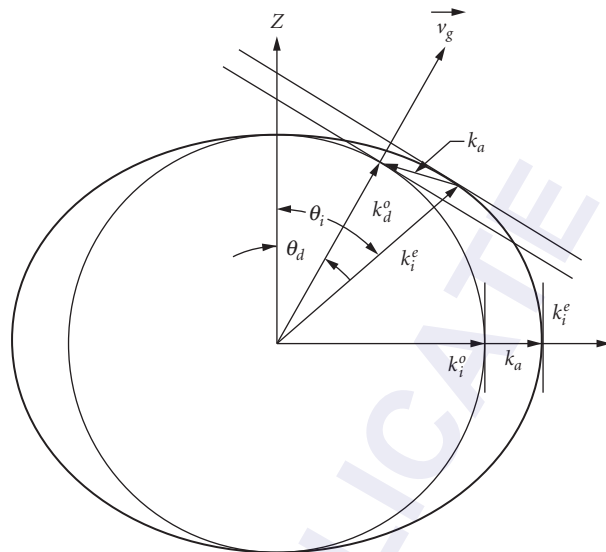


FIGURE 8 Wavevector diagram for noncollinear AOTF showing noncritical phase matching.

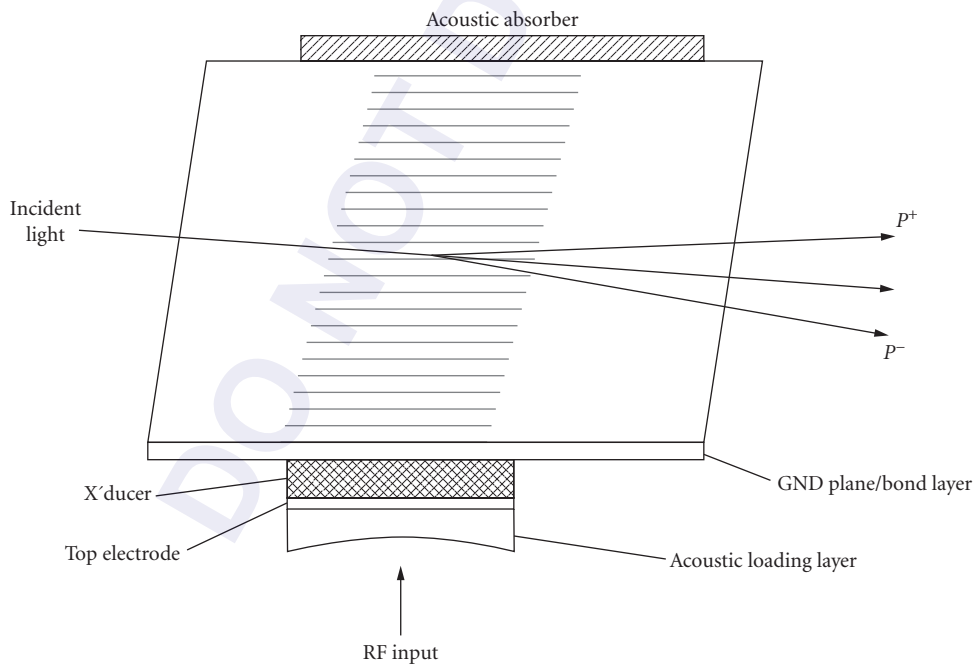


FIGURE 9 Schematic of noncollinear AOTF.

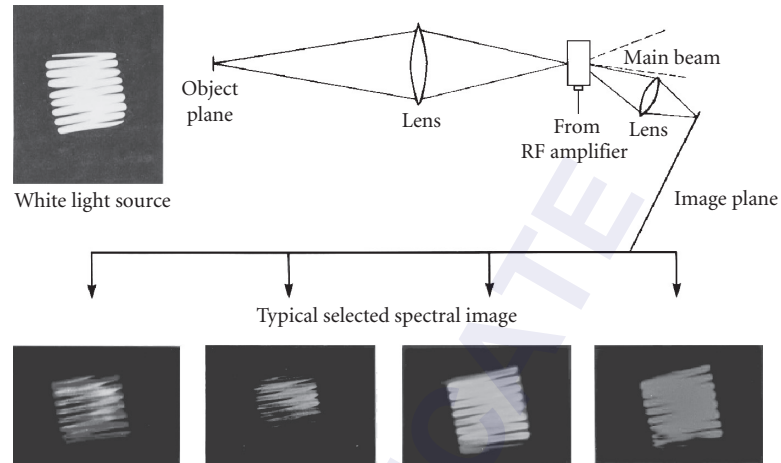


FIGURE 10 Color images of lamp filament through the first noncollinear AOTF.

with a theoretical analysis.⁵³ The early work on the theory and practice of the AOTF is given in a review paper.⁵⁴

After the first noncollinear TeO_2 AOTF, it was recognized that because of its larger aperture and simpler optical geometry, this new type of AOTF would be better suited to spectral imaging applications. A multispectral spectral experiment using a TeO_2 AOTF was performed in the visible region.⁵⁵ The white light beam was spatially separated from the filtered light and blocked by an aperture stop in the immediate frequency plane. A resolution target was imaged through the AOTF and relayed onto the camera. At a few wavelengths in the visible region selected by the driving RF frequencies, the spectral imaging of the resolution was measured. The finest bar target had a horizontal resolution of 144 lines/mm and a vertical resolution of 72 lines/mm. Figure 10 shows the measured spectral images of the lamp filament through the AOTF at a few selected wavelengths in the visible.

Other application concepts of the new AOTF have also been demonstrated. These include the detection of weak laser lines in the strong incoherent background radiation using the extreme difference in the temporal coherence⁵⁶ and the operation of multiwavelength AOTF driven by multifrequencies simultaneously.⁵⁷

It is instructive to summarize the advantages of the noncollinear AOTF: (a) it uses efficient AO materials; (b) it affords the design freedom of choosing the direction of optical and acoustic wave for optimizing efficiency, resolution, angular aperture, and so on; (c) it can be operated without the use of polarizers; and (d) it allows simple filter construction ease for manufacturing. As a result, practically all AOTFs are noncollinear types satisfying the NPM condition.

Noncritical Phase-Matching AOTF

AOTF characteristics The key performance characteristics of the AOTF operated in the noncritical phase-matching (NPM) mode will be reviewed. These characteristics include tuning relation, angle of deflection and imaging resolution, passband response, spectral resolution, angular aperture, transmission and drive power, out-of-band rejection, and sidelobe suppression.

Tuning relation For an AOTF operated at NPM, the tangents to the incident and diffracted optical wavevector surfaces must be parallel. The parallel tangent condition implies that the diffracted extraordinary ray is collinear with the incident ordinary ray,¹⁹

$$\tan \theta_e = e^2 \tan \theta_o \quad (78)$$

For a given incident light angle θ_i , the diffracted light angle θ_d is readily determined from the above equation. Thus, without loss of generality we can assume the incident optical beam to be either ordinary or extraordinary polarized. To minimize wavelength dispersion an ordinary polarized light (*o*-wave) is chosen as the incident light in the following analysis.

Substituting Eq. (78) into the phase matching Eq. (22), the acoustic wave direction and center wavelength of the passband can be expressed as a function of the incident light angle.

$$\tan \theta_a = (\cos \theta_o + \rho_o) / \sin \theta_o \quad (79)$$

$$\lambda_o = [n_o V(\theta_a) / f_a] [(\rho_o - 1)^2 \cos^2 \theta_o + (\delta \rho_o - 1)^2 \sin^2 \theta_o]^{1/2} \quad (80)$$

where $\rho_o = (1 + \delta \sin^2 \theta_o)^{1/2}$ and $\delta = (e^2 - 1)/2$. The above expressions are exact. For small birefringence $\Delta n = n_o |\delta| \ll n_o$, Eqs. (81) and (82) yield the approximate solution for the acoustic wave direction and acoustic frequency:

$$\tan \theta_a = 2 / \tan \theta_o - (1 + \delta) \tan \theta_o \quad (81)$$

$$\lambda_o = (V \Delta n / f_a) (\sin^4 \theta_o + \sin^2 2\theta_o)^{1/2} \quad (82)$$

Angle of deflection In a noncollinear AOTF the diffracted light is spatially separated from the incident light. The deflection angle θ_D is defined as the deviation of the diffracted light from that of the incident light. For small birefringence, θ_D is approximately given by

$$\theta_D \approx \Delta n \sin 2\theta_o \quad (83)$$

Equation (83) shows that θ_D depends only on the input optical angle. It reaches maximum value Δn when $\theta_o = 45^\circ$. For instance, in a TeO₂ AOTF operated at 633 nm, $\Delta n \approx 0.15$, the maximum deflection angle for TeO₂ AOTF is about 8.6°. As a result of the finite angle of deflection angle the noncollinear AOTF can operate without the use of polarizers.

Bandpass response For an acoustic column of uniform amplitude, the bandpass response is given by $T(\lambda_o) = T_o \sin^2 \Delta\sigma L$, where T_o is the peak transmission at exact momentum matching, L is the interaction length, and $\Delta\sigma$ is the momentum mismatch. It can be shown¹⁹ that

$$\Delta\sigma = -b \sin^2 \theta_i (\Delta\lambda / \lambda_o^2) + (\Delta n / 2\lambda_o) [F_1 (\Delta\theta)^2 + F_2 (\sin \theta_i \Delta\phi)^2 + F_3 (\Delta\theta)^3] \quad (84)$$

where

$$F_1 = 2 - 3 \sin^2 \theta_o \quad F_2 = 2 - \sin^2 \theta_o \quad F_3 = -\sin 2\theta_o \quad (85)$$

$\Delta\lambda$, $\Delta\theta$, and $\Delta\phi$ are deviations in wavelength, polar, and azimuth angles of the incident light beam, b is the dispersion constant defined by $b = \Delta n - \lambda_o \delta(\Delta n) / \delta \lambda_o$.

Resolution Equation (84) shows that half peak transmission occurs when $\Delta\sigma L \approx 0.44$. The full width at half-maximum (FWHM) of the AOTF is given by

$$\Delta\lambda = 0.9 \lambda_o^2 / b L \sin^2 \theta_i \quad (86)$$

The spectral resolution (i.e., longitudinal resolution) of the AOTF is given by

$$R_l = \lambda_o / \Delta\lambda = \Delta n L \sin^2 \theta_i / \lambda_o \quad (87)$$

Angular aperture The acceptance (half) angles in the polar plane ($\theta_o \neq 54.7$) and azimuth planes are given by

$$\Delta\theta = \pm n(\lambda_o/\Delta n L F_1)^{1/2} \quad (88)$$

$$\Delta\phi = \pm n(\lambda_o/\Delta n L F_2)^{1/2} \quad (89)$$

At $\theta_o \approx 54.7$, $F_1 = 0$, and the angular aperture in the polar plane reaches a maximum.

$$\Delta\theta = \pm n(\lambda_o/\Delta n L)^{1/3} \quad (90)$$

As a general rule, the noncollinear AOTF is operated without the use of polarizers. Thus the deflection angle θ_D sets the upper limit for the angular aperture for the AOTF. Assuming the incident is a rectangular beam of width D . This yields a transverse or imaging resolution

$$R_i = \frac{\theta_D}{\delta\theta_o} = \frac{D\Delta n}{\lambda_o} \sin 2\theta_o \quad (91)$$

Optical throughput The filter transmission for normalized input light distribution $I(\theta_r, \phi_i)$ is obtained by integrating the plane-wave transmission over the solid angle aperture

$$I_d(\Delta\lambda) = \int_{\alpha_i, \beta_i} T(\Delta\lambda, \theta_i, \phi_i) I(\theta_i, \phi_i) \sin\theta_i d\theta_i d\phi_i \quad (92)$$

Transmission and drive power The peak transmission of an AOTF is given by

$$T_o = \sin^2 \left(\frac{\pi^2}{2\lambda_o^2} M_2 P_d L^2 \right)^{1/2} \quad (93)$$

where P_d is the acoustic power density. Maximum transmission occurs when the drive power reaches the value

$$P_d = \frac{\lambda_o^2 A}{2M_2 L^2} \quad (94)$$

where A is the optical aperture. Because of the λ_o^2 dependence, the drive power of an AOTF increases rapidly as the wavelength increases. For instance, at $\lambda_o = 4 \mu\text{m}$, the required acoustic power for a 1 cm infrared TeO_2 AOTF exceeds about 10 W. The required high drive power is perhaps the most severe limitation of the infrared AOTF.

Sidelobe suppression As an optical filter, one of the most undesired characteristics of the AOTF is the nonuniform in-band frequency response and the high sidelobes/fall-off shape near the edge of the passband. The AOTF bandpass response is proportional to the amplitude square of the Fourier transform of the acoustic field. For uniform acoustic excitation, the AOTF exhibits a sinc²-type bandpass characteristic with many sidelobes. The nearest sidelobe is only about 13 dB below the main lobe. A more serious problem is that the envelope of the sidelobe decays slowly at a fall-off rate of -6 dB per octave. In principle, suppression of the high sidelobes can be obtained by the technique of amplitude apodization, a process that slowly reduces the acoustic amplitude profile at the transducer according to a prescribed weighting function. The real problem is the practicality of implementing the apodization. A number of transducer apodization techniques have been experimentally investigated.⁵⁸ These include (1) the use of segmented transducers weighted with designed resistive coupling or varying electrode area; and (2) tapered windows with an air gap or dielectric layer.

The effectiveness of each method varies depending on the size and thickness of the acoustic frequency, among other factors. Typically, the highest sidelobes of apodized devices are about 20 to 25 db below main peak. However, the result is very sensitive to the fabrication tolerance and design parameters. The price to pay is the reliability and high cost of the fabrication.

Spectral Tuning Range Prior to 1980, both the collinear and noncollinear types of AOTFs have been successfully demonstrated covering the spectral range from ultraviolet (UV) to long-wave infrared (LWIR). However, due to the various performance limitations such as polarization dependence, limited optical aperture and resolution and in particular high drive power, the TeO₂ AOTF is the only type that has been deployed in practical optical systems. Considerable effort were devoted to overcome the basic limitations of AOTF technology and extend the usable spectral range.⁵⁸ The progress of AOTFs for the various spectral regions are summarized as follows.

Ultraviolet (UV) AOTF Until recently, collinear quartz AOTF has been the only type for the UV region. The performance limitations of the quartz filter include (1) limited UV operating range: due to the optical absorption of calcite polarizers, the short wavelength cutoff of the UV AOTF is about 240 nm, (2) high drive power: because of the small figure of merit, the drive power is typically above 10 W; and (3) low optical throughput: the fixed on-axis of the collinear interaction geometry lacks the design flexibility with typically small solid angle \times aperture characteristics. To overcome these drawbacks, two noncollinear UV AOTFs were developed.⁵⁹

These include (1) A KDP UV AOTF operated without polarizers was demonstrated at 253.7 nm using a Hg lamp. The device was tunable from 190 to 340 nm. The measured result include a full-width at half-maximum (FWHM) of 0.79 nm and a peak diffraction efficiency of 50 percent at 1.3 W input power. The deflection angle of the filtered light separated from the broadband incident light is 1.9°. (2) A Far UV AOTF made of MgF₂ was tested at 196.03 nm from a hollow cathode lamp. A Brewster-type polarizer made of a pile of thin MgF₂ plates was used in the test. The FWHM of the AOTF was about 3 nm. The diffracted light intensity was about 45 percent at a drive power of 6 W. Because of the UV transmission, the AOTF is extendable to 150 nm.

Mid-infrared AOTF The most severe technological issues that limit the extension of spectral range to thermal infrared region are (1) high drive power requirement due to the λ^2 dependence and (2) catastrophic device failure when operated at low temperatures. To resolve these critical issues a number of the TeO₂ AOTFs feasibility models for the mid-wave infrared (MWIR) have been built and tested.⁶⁰ Performance characteristics of the 1 cm² size acoustically resonant AOTFs measured at 3.39 μ m include tuning range: 2.5 to 5 μ m, bandwidth 15 nm, and drive power (at 80 percent efficiency) less than 0.5 W, with a resonance gain about 20. Repeated operation of the low-power MWIR AOTF was obtainable for temperature above 150 K. However, because of the power thermal conductivity of TeO₂ below the critical temperature the crystal or the transducer may crack. The reliability of the resonant AOTF operated at cryogenic temperature remains a major engineering issue.

Long wave infrared AOTF Operation of the AOTF were demonstrated in the LWIR region of 8 to 12 μ m were demonstrated prior to 1980. Stimulated by the important potential application for hyperspectral imaging, there has been considerable effort to develop practical LWIR AOTF for optical remote sensing systems.

For AOTF operated in the long-wave infrared (LWIR) the high drive power is even a much serious problem. Because of its high figure of merit Tl₃AsSe₃ (TAS) is considered to be the best LWIR AOTF material. Even with its exceedingly large AO figure of merit (>2000), the TAS AOTF has to operate on pulsed basis with low duty cycle. However, because of extremely low thermal conductivity and the brittle nature of the TAS crystals, the potential of practical TAS AOTF operated at low temperature does not appear promising.

Performance of NPM AOTF The AOTF possesses many salient features that make it attractive for a variety of optical system applications. To realize the benefits of these merits the limitations of the AOTF must be considered when compared with the competing technologies. Based on the

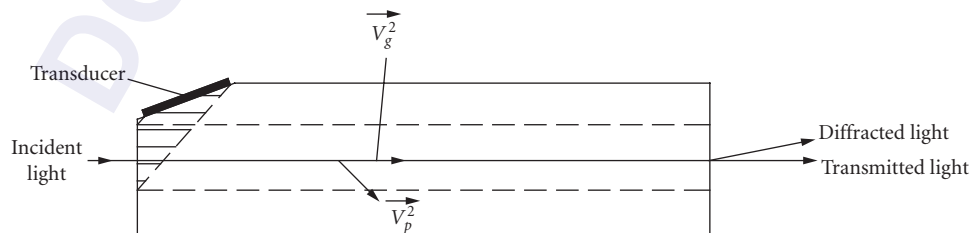
TABLE 9 Broadband and High Resolution Type NPM AOTF

Type	Broadband (8 × 8 mm)		High-Resolution AOTF	
	10 × 10 mm		5 × 5 mm	
Aperture	10 × 10 mm		5 × 5 mm	
Wavelength	400–1100	700–2500	400–650	650–1100
Bandwidth (cm ⁻¹)	25	20	5	7
$\Delta\lambda$ (nm)	1	5	0.12	0.5
At λ_o (μm)	633	1550	442	830
Efficiency (%)	80	70	95	95
RF power (W)	1	2	1	2

previous discussion of the usable spectral range, it appears that with the exception of possible new development in the UV, only the TeO₂ AOTFs operated in the visible to SWIR can be considered as a matured technology ready for system deployment. Considering the primary niche, it is pertinent to improve the basic performance of the AOTFs for meeting the system requirement. Table 9 shows two selected high performance NPM AOTFs. These include (1) broadband imaging AOTF with two octave tuning range in a single unit and (2) high resolution high efficiency AOTF suited as rapid random access laser tuner.

The AOTF has the unique capability of being able to simultaneously and independently add or drop multiwavelength signals. As such, it can serve as a WDM cross-connect for routing multiwavelength optical signals along a prescribed connection path determined by the signal's wavelength. Because of this unique attractive feature, the AOTF appears to be suited for use as dynamical reconfigurable components for the WDM network. However, due to the relatively high drive power and low resolution requirement, the AOTF has not been able to meet the requirement for dense wavelength division multiplexing (DWDM) applications. This basic drawback of the AOTF is the result of finite interaction length limited by the large acoustic beam walk-off in the TeO₂ crystal. To overcome this intrinsic limitation, a new type of noncollinear AOTF showing significant improvement of resolution and diffraction efficiency was proposed and demonstrated.^{60,61} The filter is referred to as the collinear beam (CB) AOTF, since the group velocity of the acoustic wave and light beams are chosen to be collinear. An extended interaction length was realized in a TeO₂ AOTF with narrow bandwidth and significantly lower drive power. Figure 11 shows the schematic of an in-line TeO₂ CBAOTF using internal mode conversion.

An initial test of the CBAOTF was performed at 1532 nm using a HeNe laser as the light source.⁶¹ Figure 12 shows the bandpass response of the CBAOTF obtained by monitoring the diffracted light intensity as the acoustic frequency was swept through the laser line. As shown in the figure, the slowly decaying bandpass response appears to be a Lorentzian shape with observable sidelobes. However, the falloff rate of the bandpass at -6 dB per octave wavelength change is the same as the envelope of the sinc² response of the conventional AOTF.

**FIGURE 11** Collinear beam AOTF using internal mode conversion.

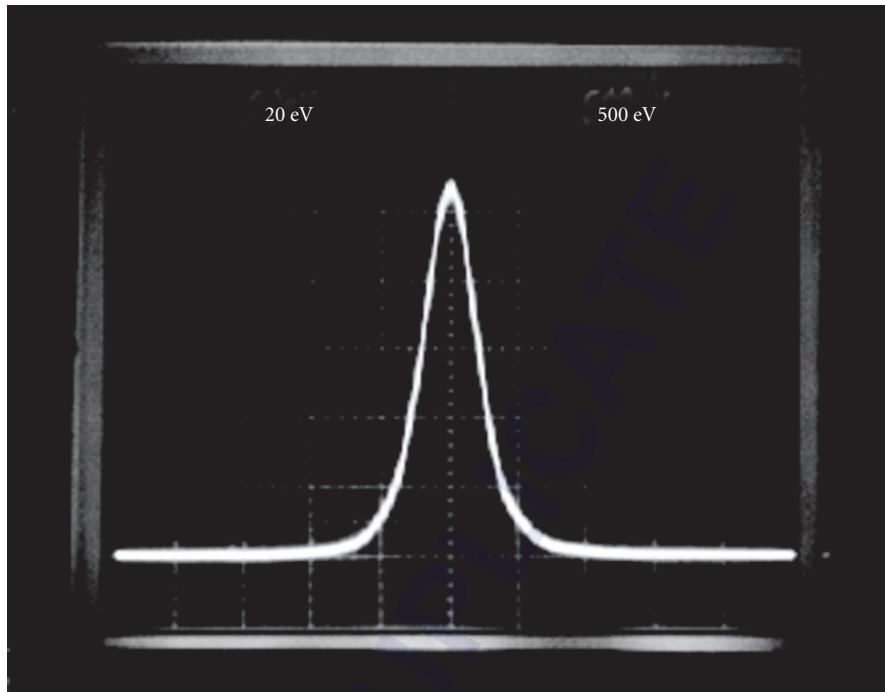


FIGURE 12 Bandpass of CBAOTF (RF swept through a 1.550- μm laser line).

The half-power bandwidth was measured to be 25 kHz, which corresponds to an optical FWHM of 1 nm or 4.3 cm^{-1} . The diffracted light reaches a peak value of 95 percent when the drive power is increased to about 55 mW. Compared to the state-of-the-art high resolution NPM, the measured result showed that the drive power was about 50 times smaller. The low drive power advantage of the CBAOTF is most important for WDM application, which requires simultaneously a large number of channels.

Although the CBAOTF has resolved the most basic limitation of high drive power requirement, to be used as a dynamic DWDM component, there still remains several critical technical bottleneck. For a 100-GHz (0.8 nm at 1550 nm) wavelength spacing system, the AOTF has to satisfy a set of performance goals. These include: polarization independent operation, full width at half-maximum (FWHM) of 0.4 nm, drive power of 150 mW per signal, and the sidelobe must be suppressed to be lower than at least 30 dB at 100 GHz away from the center wavelength. Significant progress has been made in the effort to overcome these critical issues. These are discussed below.

Polarization independence The operation of AOTF is critically dependent on the polarization state of the incoming light beam. To make it polarization-independent, polarization diversity configurations (PDL) are used.⁶² The scheme achieves polarization independence by dividing into two beams of orthogonal polarization, o- and e-rays with a polarization beam splitter (PBS) passing through two single polarization AOTFs in two paths, then combing the two diffracted o- and e-rays of selected wavelengths with a second PBS. Half-wave plates are used to convert the polarization of the light beam so that o-rays are incident onto the AOTF.

Sidelobe suppression Due to low crosstalk requirement, the sidelobe at the channel must be sufficiently low ($\sim 35 \text{ dB}$). This is the most basic limitation and as such it will be discussed in some detail. There are three kinds of crosstalk. These include the interchannel crosstalk caused by the

TABLE 10 Measured Performance of CBAOTF

Tuning range	1300–1600 nm
Measured wavelength	1550 nm
FWHM (3 dB)	0.4 nm @ 1.55 nm
Sidelobe @ ch. spacing	–20 dB @ –0.8 nm
Sidelobe @ ch. spacing	–25 dB @ +0.8 nm
Peak efficiency	95(%)
Drive power	80 mW

sidelobe level of the AOTF bandpass; and the extinction ratio, an intrachannel crosstalk due to the finite extinction ratio. The most severe crosstalk is the coherent type that originates from the mixing of sidelobe of light beam λ_1 shifted by frequencies f_1 and f_2 . The optical interference of the two light beams will result in an amplitude-modulated type crosstalk at the difference frequency of $f_1 - f_2$. This type of interchannel crosstalk is much more severe since the modulation is proportional to the amplitude or square root of the sidelobe power level.

Several apodization techniques have been developed in order to suppress the high sidelobe. A simple approach of tilted configuration to simulate a various weighting function appears to be most practical.⁵⁸ A major advantage of this approach is the design flexibility to obtain a desired tradeoff between sidelobe suppression and bandwidth.

Another technique for reducing the sidelobe is to use two or more AOTFs in an incoherent optical cascade. The bandpass response of the incoherently cascaded AOTFs is equal to the product of the single stage and thus can realize a significantly reduced sidelobe. This doubling of sidelevel by cascaded cells has been experimentally demonstrated.⁵⁸

A number of prototype devices of polarization independent (PI) CBAOTF using the tilt configuration have been built. Typical measured results at 1550 nm include 1.0 nm FWHM, peak efficiency 95 percent at 80 mW drive power insertion loss; –3 dB polarization-independent loss; 0.1 dB polarization mode dispersion (PDL); 1 psec, and sidelobe below –27 dB at 3 nm from the center wavelength.⁵⁸

To further reduce the half power bandwidth, a higher angle cut design was chosen in the follow-on experimental work. A 65°, apodized CBAOTF was designed and fabricated. The primary design goal was to meet the specified narrow bandwidth and accept the sidelobe level based on the tradeoff relation. Test results of the 65° devices measured at 1550 are summarized in Table 10. Except the sidelobe suppression goal, the 65° device essentially satisfies all other specifications. The specified sidelobe level of –35 dB for the 100-GHz channel spacing can be met by using two CBAOTFs in cascade.

In conclusion, it is instructive to emphasize the unique advantage of the AOTF. Because of its random access wavelength tunability over large spectral range, the CBAOTF provides a low-cost implementation of a dynamic multiwavelength component for DWDM and other noncommunication type of application with nonuniform distribution of wavelengths.

6.8 REFERENCES

1. L. Brillouin, "Diffusion de la lumière et des ray x par un corps transparent homogène," *Ann. Phys.* **17**:80–122 (1992).
2. P. Debye and F. W. Sears, "On the Scattering of Light by Supersonic Waves," *Proc. Nat. Acad. Sci. (U.S.)* **18**:409–414 (1932).
3. R. Lucas and P. Biquard, "Propriétés optiques des milieux solides et liquides soumis aux vibration élastiques ultra sonores," *J. Phys. Rad.* **3**:464–477 (1932).
4. C. V. Raman and N. S. Nagendra Nath, "The Diffraction of Light by High Frequency Sound Waves," *Proc. Ind. Acad. Sci.* **2**:406–420 (1935); **3**:75–84 (1936); **3**:459–465 (1936).

5. M. Born and E. Wolf, *Principles of Optics*, 3rd ed., Pergamon Press, New York, 1965, Chap. 12.
6. E. I. Gordon, "A Review of Acoustooptical Deflection and Modulation Devices," *Proc. IEEE* **54**:1391–1401 (1966).
7. A. Korpel, R. Adler, Desmares, and W. Watson, "A Television Display Using Acoustic Deflection and Modulation of Coherent Light," *Proc. IEEE* **54**:1429–1437 (1966).
8. R. W. Dixon, "Acoustic Diffraction of Light in Anisotropic Media," *IEEE J. Quantum Electron.* **QE-3**:85–93 (Feb. 1967).
9. S. E. Harris and R. W. Wallace, "Acousto-Optic Tunable Filter," *J. Opt. Soc. Am.* **59**:744–747 (June 1969).
10. I. C. Chang, "Noncollinear Acousto-Optic Filter with Large Angular Aperture," *Appl. Phys. Lett.* **25**:370–372 (Oct. 1974).
11. E. K. Sittig, "Elasto-Optic Light Modulation and Deflection," Chap. VI, in *Progress in Optics*, vol. X, E. Wolf (ed.), North-Holland, Amsterdam, 1972.
12. N. Uchida and N. Niizeki, "Acoustooptical Deflection Materials and Techniques," *Proc. IEEE* **61**:1073–1092 (1973).
13. I. C. Chang, "Acoustooptical Devices and Applications," *IEEE Trans. Sonics Ultrason.* **SU-23**:2–22 (1976).
14. A. Korpel, "Acousto-Optics," Chap. IV, in *Applied Optics and Optical Engineering*, R. Kingslake and B. J. Thompson (eds.), Academic Press, New York, vol. VI, 1980.
15. J. F. Nye, *Physical Properties of Crystals*, Clarendon Press, Oxford, England, 1967.
16. F. Nelson and M. Lax, "New Symmetry for Acousto-Optic Scattering," *Phys. Rev., Lett.*, **24**:378–380 (Feb. 1970); "Theory of Photoelastic Interaction," *Phys. Rev.* **B3**:2778–2794 (Apr. 1971).
17. J. Chapelle and L. Tael, "Theorie de la diffusion de la lumiere par les cristaux fortement piezoelectriques," *C. R. Acad. Sci.* **240**:743 (1955).
18. W. R. Klein and B. D. Cook, "Unified Approach to Ultrasonic Light Diffraction," *IEEE Trans. Sonics Ultrason.* **SU-14**:123–134 (1967).
19. D. A. Klienman, "Optical Second Harmonic Generation," *Phys. Rev.* **128**:1761 (1962).
20. I. C. Chang, "Acousto-Optic Tunable Filters," *Opt. Eng.* **20**:824–828 (1981).
21. A. Korpel "Acoustic Imaging by Diffracted Light. I. Two-Dimensional Interaction," *IEEE Trans. Sonics Ultrason.* **SU-15**(3):153–157 (1968).
22. I. C. Chang and D. L. Hecht, "Device Characteristics of Acousto-Optic Signal Processors," *Opt. Eng.* **21**:76–81 (1982).
23. I. C. Chang, "Selection of Materials for Acoustooptical Devices," *Opt. Eng.* **24**:132–137 (1985).
24. I. C. Chang, "Acousto-Optic Devices and Applications," in *Optic Society of America Handbook of Optics*, 2nd ed, M. Bass (ed.), Vol. II, Chap. 12, pp. 12.1–54, 1995.
25. I. C. Chang, "Design of Wideband Acoustooptical Bragg Cells," *Proc. SPIE* **352**:34–41 (1983).
26. D. L. Hecht and G. W. Petrie, "Acousto-Optic Diffraction from Acoustic Anisotropic Shear Modes in GaP," *IEEE Ultrason. Symp. Proc.* p. 474, Nov. 1980.
27. E. G. H. Lean, C. F. Quate, and H. J. Shaw, "Continuous Deflection of Laser Beams," *Appl. Phys. Lett.* **10**:48–50 (1967).
28. W. Warner, D. L. White, and W. A. Bonner, "Acousto-Optic Light Deflectors Using Optical Activity in Praterullurite," *J. Appl. Phys.* **43**:4489–4495 (1972).
29. T. Yano, M. Kawabuchi, A. Fukumoto, and A. Watanabe, "TeO₂ Anisotropic Bragg Light Deflector Without Midband Degeneracy," *Appl. Phys. Lett.* **26**:689–691 (1975).
30. G. A. Couquin, J. P. Griffin, and L. K. Anderson, "Wide-Band Acousto-Optic Deflectors Using Acoustic Beam Steering," *IEEE Trans. Sonics Ultrason.* **SU-18**:34–40 (Jan. 1970).
31. D. A. Pinnow, "Acousto-Optic Light Deflection: Design Considerations for First Order Beamsteering Transducers," *IEEE Trans. Sonics Ultrason.* **SU-18**:209–214 (1971).
32. I. C. Chang, "Birefringent Phased Array Bragg Cells," *IEEE Ultrason. Symp. Proc.*, pp. 381–384, 1985).
33. E. H. Young, H. C. Ho, S. K. Yao, and J. Xu, "Generalized Phased Array Bragg Interaction in Birefringent Materials," *Proc. SPIE* **1476**: (1991).
34. A. J. Hoffman and E. Van Rooyen, "Generalized Formulation of Phased Array Bragg Cells in Uniaxial Crystals," *IEEE Ultrason. Symp. Proc.* p. 499, 1989.

35. R. T. Waverka and K. Wagner "Wide Angle Aperture Acousto-Optic Bragg Cell," *Proc. SPIE* **1562**:66–72 (1991).
36. W. H. Watson and R. Adler, "Cascading Wideband Acousto-Optic Deflectors," *IEEE Conf. Laser Eng. Appl.* Washington., D.C., June 1969.
37. I. C. Chang and D. L. Hecht, "Doubling Acousto-Optic Deflector Resolution Utilizing Second Order Birefringent Diffraction," *Appl. Phys. Lett.* **27**:517–518 (1975).
38. D. L. Hecht, "Multifrequency Acousto-Optic Diffraction," *IEEE Trans. Sonics and Ultrason.* **SU-24**:7 (1977).
39. I. C. Chang and R. T. Waverka, "Multifrequency Acousto-Optic Diffraction," *IEEE Ultrason. Symp. Proc.* p. 445, Oct. 1983.
40. I. C. Chang and S. Lee, "Efficient Wideband Acousto-Optic Bragg Cells," *IEEE Ultrason. Symp. Proc.* p. 427, Oct. 1983.
41. I. C. Chang et al., "Progress of Acousto-Optic Bragg Cells," *IEEE Ultrason. Symp. Proc.* p. 328, 1984.
42. I. C. Chang and R. Cadieux, "Multichannel Acousto-Optic Bragg Cells," *IEEE Ultrason. Symp. Proc.* p. 413, 1982.
43. W. R. Beaudot, M. Popek, and D. R. Pape, "Advances in Multichannel Bragg Cell Technology," *Proc. SPIE* **639**:28–33 (1986).
44. R. V. Johnson, "Acousto-Optic Modulator," in *Design and Fabrication of Acousto-Optic Devices*, A. Goutzoulis and D. Pape, (eds.), Marcel Dekker, New York, 1994.
45. E. H. Young and S. K. Yao, "Design Considerations for Acousto-Optic Devices," *Proc. IEEE* **69**:54–64 (1981).
46. D. Maydan, "Acousto-Optic Pulse Modulators," *J. Quantum Electron.* **QE-6**:15–24 (1967).
47. R. V. Johnson, "Scophony Light Valve," *Appl. Opt.* **18**:4030–4038 (1979).
48. I. C. Chang, "Large Angular Aperture Acousto-Optic Modulator," *IEEE Ultrason. Symp. Proc.* pp. 867–870, 1994.
49. I. C. Chang, "Acoustic-Optic Tunable Filters," in *Acousto-Optic Signal Processing*, N. Berg and J. M. Pellegrino (eds.), Marcel Dekker, New York, 1996.
50. S. E. Harris, S. T. K. Nieh, and R. S. Feigelson, "CaMoO₄ Electronically Tunable Acousto-Optical Filter," *Appl. Phys. Lett.* **17**:223–225 (Sep. 1970).
51. J. A. Kusters, D. A. Wilson, and D. L. Hammond, "Optimum Crystal Orientation for Acoustically Tuned Optic Filters," *J. Opt. Soc. Am.* **64**:434–440 (Apr. 1974).
52. I. C. Chang, "Tunable Acousto-Optic Filter Utilizing Acoustic Beam Walk-Off in Crystal Quartz," *Appl. Phys. Lett.* **25**:323–324 (Sep. 1974).
53. I. C. Chang, "Analysis of the Noncollinear Acousto-Optic Filter," *Electron. Lett.* **11**:617–618 (Dec. 1975).
54. D. L. Hecht, I. C. Chang, and A. Boyd, "Multispectral Imaging and Photomicroscopy Using Tunable Acousto-Optic Filters," *OSA Annual Meeting*, Boston, Mass., Oct. 1975.
55. I. C. Chang, "Laser Detection Using Tunable Acousto-Optic Filter," *J. Quantum Electron.* **14**:108 (1978).
56. I. C. Chang, et al., "Programmable Acousto-Optic Filter," *IEEE Ultrason. Symp. Proc.* p. 40, 1979.
57. R. T. Waverka, P. Katzka, and I. C. Chang, "Bandpass Apodization Techniques for Acousto-Optic Tunable Filters," *IEEE Ultrason. Symp.*, San Francisco, CA, Oct. 1985.
58. I. C. Chang, "Progress of Acoustooptic Tunable Filter," *IEEE Ultrason. Symp. Proc.* p. 819, 1996.
59. I. C. Chang and J. Xu, "High Performance AOTFs for the Ultraviolet," *IEEE Ultrason. Proc.* p. 1289, 1988.
60. I. C. Chang, "Collinear Beam Acousto-Optic Tunable Filter," *Electron Lett.* **28**:1255 (1992).
61. I. C. Chang et al, "Bandpass Response of Collinear Beam Acousto-Optic Filter," *IEEE Ultrason. Symp. Proc.* vol. 1, pp. 745–748.
62. I. C. Chang, "Acousto-Optic Tunable Filters in Wavelength Division Multiplexing (WDM) Networks," *1997 Conf. Laser Electro-Optics (CLEO)*, Baltimore, MD, May 1997.

This page intentionally left blank.

DO NOT DUPLICATE

ELECTRO-OPTIC MODULATORS

Georganne M. Purvinis

*The Battelle Memorial Institute
Columbus, Ohio*

Theresa A. Maldonado

*Department of Electrical and Computer Engineering
Texas A&M University
College Station, Texas*

7.1 GLOSSARY

$A, [A]$	general symmetric matrix
$[a]$	orthogonal transformation matrix
b	electrode separation of the electro-optic modulator
D	displacement vector
d	width of the electro-optic crystal
E	electric field
H	magnetic field
IL	insertion loss
k	wavevector, direction of phase propagation
L	length of the electro-optic crystal
L/b	aspect ratio
N	number of resolvable spots
n_m	refractive index of modulation field
n_x, n_y, n_z	principal indices of refraction
r_{ijk}	third-rank linear electro-optic coefficient tensor
S	Poynting (ray) vector, direction of energy flow
S_{ijkl}	fourth-rank quadratic electro-optic coefficient tensor
T	transmission or transmissivity
V	applied voltage
V_π	half-wave voltage
v_π	phase velocity
v_m	modulation phase velocity

v_s	ray velocity
w	beamwidth
ω_o	resonant frequency of an electro-optic modulator circuit
\mathbf{X}	position vector in cartesian coordinates
\mathbf{X}'	electrically perturbed position vector in cartesian coordinates
(x, y, z)	unperturbed principal dielectric coordinate system
(x', y', z')	new electro-optically perturbed principal dielectric coordinate system
(x'', y'', z'')	wavevector coordinate system
(x''', y''', z''')	eigenpolarization coordinate system
β_1	polarization angle between x''' and x''
β_2	polarization angle between y''' and x''
Γ	phase retardation
Γ_m	amplitude modulation index
$\Delta\eta$	electro-optically induced change in the index of refraction or birefringence
$\Delta(1/\eta^2)$	electro-optically induced change in an impermeability tensor element
$\Delta\phi$	angular displacement of beam
$\Delta\nu$	bandwidth of a lumped electro-optic modulator
δ	phase modulation index
$[\boldsymbol{\epsilon}]$	permittivity tensor
ϵ_0	permittivity of free space
$[\boldsymbol{\epsilon}]^{-1}$	inverse permittivity tensor
$\epsilon_x, \epsilon_y, \epsilon_z$	principal permittivities
$[\boldsymbol{\epsilon}]$	dielectric constant tensor
$[\boldsymbol{\epsilon}]^{-1}$	inverse dielectric constant tensor
$\epsilon_x, \epsilon_y, \epsilon_z$	principal dielectric constants
η_m	extinction ratio
θ	half-angle divergence
θ_k, ϕ_k	orientation angles of the wavevector in the (x, y, z) coordinate system
ϑ	optic axis angle in biaxial crystals
λ	wavelength of the light
$[\lambda]$	diagonal matrix
ν_{tw}	bandwidth of a traveling wave modulator
ξ	modulation efficiency
ρ	modulation index reduction factor
τ	transit time of modulation signal
ϕ	phase of the optical field
Ω	plane rotation angle
ϖ	beam parameter for bulk scanners
ω_d	frequency deviation
ω_e	stored electric energy density
ω_m	modulation radian frequency
$[1/n^2]'$	electro-optically perturbed impermeability tensor
$[1/n^2]$	inverse dielectric constant (impermeability) tensor
Γ_{wg}	overlap correction factor
β	waveguide effective propagation constant

7.2 INTRODUCTION

Electro-optic modulators are used to control the amplitude, phase, polarization state, or position of an optical beam, or light wave carrier, by application of an electric field. The electro-optic effect is one of several means to impose information on, or modulate, the light wave. Other means include acousto optic, magneto optic, thermo optic, electroabsorption, mechanical shutters, and moving mirror modulation and are not addressed in this chapter, although the fundamentals presented in this chapter may be applied to other crystal optics driven modulation techniques.

There are basically two types of modulators: bulk and integrated optic. Bulk modulators are made of single pieces of optical crystals, whereas the integrated optic modulators are constructed using waveguides fabricated within or adjacent to the electro-optic material. Electro-optic devices have been developed for application in communications,¹⁻⁴ analog and digital signal processing,⁵ information processing,⁶ optical computing,^{6,7} and sensing.^{5,7} Example devices include phase and amplitude modulators, multiplexers, switch arrays, couplers, polarization controllers, and deflectors.^{1,2} Given are rotation devices,⁸ correlators,⁹ A/D converters,¹⁰ multichannel processors,¹¹ matrix-matrix and matrix-vector multipliers,¹¹ sensors for detecting temperature, humidity, radio-frequency electrical signals,^{5,7} and electro-optic sampling in ultrashort laser applications.^{12,13} The electro-optic effect allows for much higher modulation frequencies than other methods, such as mechanical shutters, moving mirrors, or acousto-optic devices, due to a faster electronic response time of the material.

The basic idea behind electro-optic devices is to alter the optical properties of a material with an applied voltage in a controlled way. The direction dependent, electrically induced physical changes in the optical properties are mathematically described by changes in the second rank permittivity tensor. The tensor can be geometrically interpreted using the index ellipsoid, which is specifically used to determine the refractive indices and polarization states for a given direction of phase propagation. The changes in the tensor properties translate into a modification of some parameter of a light wave carrier, such as phase, amplitude, frequency, polarization, or position of the light as it propagates through the device. Therefore, understanding how light propagates in these materials is necessary for the design and analysis of electro-optic devices. The following section gives an overview of light propagation in anisotropic materials that are homogeneous, nonmagnetic, lossless, optically inactive, and nonconducting. Section 7.4 gives a geometrical and mathematical description of the linear and quadratic electro-optic effects. This section illustrates how the optical properties described by the index ellipsoid change with applied voltage. A mathematical approach is offered to determine the electrically perturbed principal dielectric axes and indices of refraction of any electro-optic material for any direction of the applied electric field as well as the phase velocity indices and eigenpolarization orientations for a given wavevector direction. Sections 7.5 and 7.6 describe basic bulk electro-optic modulators and integrated optic modulators, respectively. Finally, example applications, common materials, design considerations, and performance criteria are discussed.

The discussion presented in this chapter applies to any electro-optic material, any direction of the applied voltage, and any direction of the wavevector. Therefore, no specific materials are described explicitly, although materials such as lithium niobate (LiNbO_3), potassium dihydrogen phosphate (KDP), and gallium arsenide (GaAs) are just a few materials commonly used. Emphasis is placed on the general fundamentals of the electro-optic effect, bulk modulator devices, and practical applications.

7.3 CRYSTAL OPTICS AND THE INDEX ELLIPSOID

With an applied electric field, a material's anisotropic optical properties will be modified, or an isotropic material may become optically anisotropic. Therefore, it is necessary to understand light propagation in these materials. For any anisotropic (optically inactive) crystal class there are two allowed orthogonal linearly polarized waves propagating with differing phase velocities for a given wavevector \mathbf{k} . Biaxial crystals represent the general case of anisotropy. Generally, the allowed waves exhibit *extraordinary-like* behavior; the wavevector and ray (Poynting) vector directions differ.

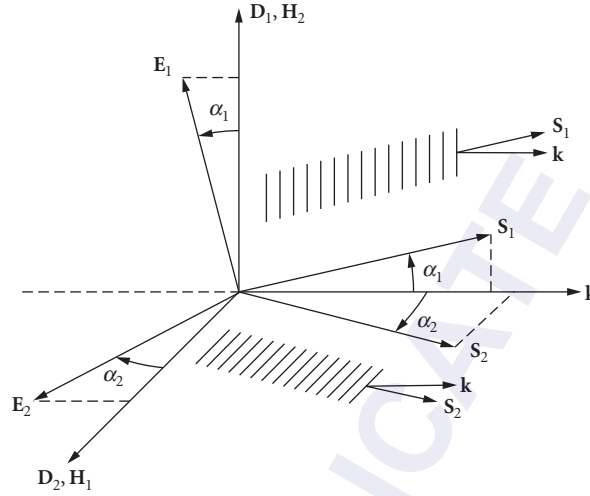


FIGURE 1 The geometric relationships of the electric quantities \mathbf{D} and \mathbf{E} and the magnetic quantities \mathbf{B} and \mathbf{H} to the wavevector \mathbf{k} and the ray vector \mathbf{S} are shown for the two allowed extraordinary-like waves propagating in an anisotropic medium.¹⁴

In addition, the phase velocity, polarization orientation, and ray vector of each wave change distinctly with wavevector direction. For each allowed wave, the electric field \mathbf{E} is not parallel to the displacement vector \mathbf{D} (which defines polarization orientation) and, therefore, the ray vector \mathbf{S} is not parallel to the wavevector \mathbf{k} as shown in Fig. 1. The angle α_1 between \mathbf{D} and \mathbf{E} is the same as the angle between \mathbf{k} and \mathbf{S} , but for a given \mathbf{k} , $\alpha_1 \neq \alpha_2$. Furthermore, for each wave $\mathbf{D} \perp \mathbf{k} \perp \mathbf{H}$ and $\mathbf{E} \perp \mathbf{S} \perp \mathbf{H}$, forming orthogonal sets of vectors. The vectors \mathbf{D} , \mathbf{E} , \mathbf{k} , and \mathbf{S} are coplanar for each wave.¹⁴

The propagation characteristics of the two allowed orthogonal waves are directly related to the fact that the optical properties of an anisotropic material depend on direction. These properties are represented by the constitutive relation $\mathbf{D} = [\boldsymbol{\varepsilon}] \mathbf{E}$, where $[\boldsymbol{\varepsilon}]$ is the permittivity tensor of the medium and \mathbf{E} is the corresponding optical electric field vector. For a homogeneous, nonmagnetic, lossless, optically inactive, and nonconducting medium, the permittivity tensor has only real components. Moreover, the permittivity tensor and its inverse, $[\boldsymbol{\varepsilon}]^{-1} = 1/\varepsilon_0 [1/n^2]$, where n is the refractive index, are symmetric for all crystal classes and for any orientation of the dielectric axes.^{15–17} Therefore the matrix representation of the permittivity tensor can be diagonalized, and in principal coordinates the constitutive equation has the form

$$\begin{pmatrix} D_x \\ D_y \\ D_z \end{pmatrix} = \begin{pmatrix} \varepsilon_x & 0 & 0 \\ 0 & \varepsilon_y & 0 \\ 0 & 0 & \varepsilon_z \end{pmatrix} \begin{pmatrix} E_x \\ E_y \\ E_z \end{pmatrix} \quad (1)$$

where reduced subscript notation is used. The principal permittivities lie on the diagonal of $[\boldsymbol{\varepsilon}]$.

The index ellipsoid is a construct with geometric characteristics representing the phase velocities and the vibration directions of \mathbf{D} of the two allowed plane waves corresponding to a given optical wave-normal direction \mathbf{k} in a crystal. The index ellipsoid is a quadric surface of the stored electric energy density ω_e of a dielectric,^{15,18}

$$\omega_e = \frac{1}{2} \mathbf{E} \times \mathbf{D} = \frac{1}{2} \sum_i \sum_j \varepsilon_i \varepsilon_{ij} E_j = \frac{1}{2} \varepsilon_0 \mathbf{E}^T [\boldsymbol{\varepsilon}] \mathbf{E} \quad i, j = x, y, z \quad (2a)$$

where T indicates the transpose.

In principal coordinates, that is, when the dielectric principal axes are parallel to the reference coordinate system, the ellipsoid is simply

$$\omega_e = \frac{1}{2} \epsilon_0 (E_x^2 \epsilon_x + E_y^2 \epsilon_y + E_z^2 \epsilon_z) \quad (2b)$$

where the convention for the dielectric constant subscript is $\epsilon_{xx} = \epsilon_x$, and so on. The stored energy density is positive for any value of electric field; therefore, the quadric surface is always given by an ellipsoid.^{15,18–20}

Substituting the constitutive equation, Eq. (2b) assumes the form $(D_x^2/\epsilon_x) + (D_y^2/\epsilon_y) + (D_z^2/\epsilon_z) = 2\omega_e \epsilon$. By substituting $x = D_x/(2\omega_e \epsilon_0)^{1/2}$ and $n_x^2 = \epsilon_x$ and similarly for y and z , the ellipsoid is expressed in cartesian principal coordinates as

$$\frac{x^2}{n_x^2} + \frac{y^2}{n_y^2} + \frac{z^2}{n_z^2} = 1 \quad (3a)$$

In a general orthogonal coordinate system, that is, when the reference coordinate system is not aligned with the principal dielectric coordinate system, the index ellipsoid of Eq. (3a) can be written in summation or matrix notation as

$$\sum_i \sum_j X_i (1/n_{ij}^2) X_j = \mathbf{X}^T [1/n^2] \mathbf{X} = (x \ y \ z) \begin{pmatrix} 1/n_{xx}^2 & 1/n_{xy}^2 & 1/n_{xz}^2 \\ 1/n_{xy}^2 & 1/n_{yy}^2 & 1/n_{yz}^2 \\ 1/n_{xz}^2 & 1/n_{yz}^2 & 1/n_{zz}^2 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} \quad (3b)$$

where $\mathbf{X} = [x, y, z]^T$, and all nine elements of the inverse dielectric constant tensor, or impermeability tensor, may be present. For sections that follow, the index ellipsoid in matrix notation will be particularly useful.

Equation (3b) is the general index ellipsoid for an optically biaxial crystal. If $n_{xx} = n_{yy}$, the surface becomes an ellipsoid of revolution, representing a uniaxial crystal. In this crystal, one of the two allowed eigenpolarizations will always be an *ordinary* wave with its Poynting vector parallel to the wavevector and \mathbf{E} parallel to \mathbf{D} for any direction of propagation. An isotropic crystal ($n_{xx} = n_{yy} = n_{zz}$) is represented by a sphere with the principal axes having equal length. Any wave propagating in this crystal will exhibit ordinary characteristics. The index ellipsoid for each of these three optical symmetries is shown in Fig. 2.

For a general direction of phase propagation \mathbf{k} , a cross section of the ellipsoid through the origin perpendicular to \mathbf{k} is an ellipse, as shown in Fig. 2. The major and minor axes of the ellipse represent the

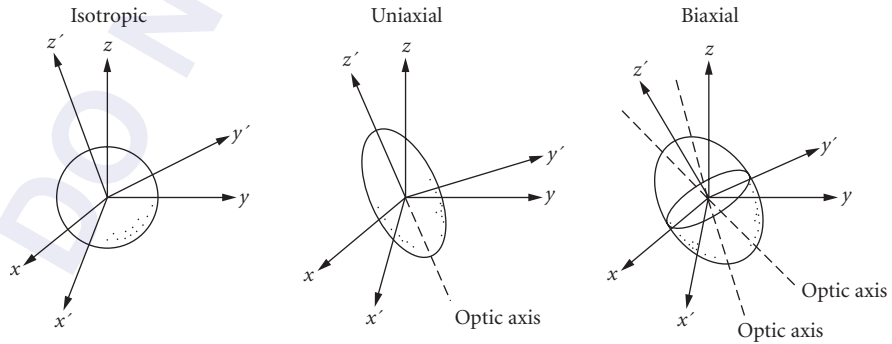


FIGURE 2 The index ellipsoids for the three crystal symmetries are shown in nonprincipal coordinates (x', y', z') relative to the principal coordinates (x, y, z). For isotropic crystals, the surface is a sphere. For uniaxial crystals, it is an ellipsoid of revolution. For biaxial crystals, it is a general ellipsoid.²¹

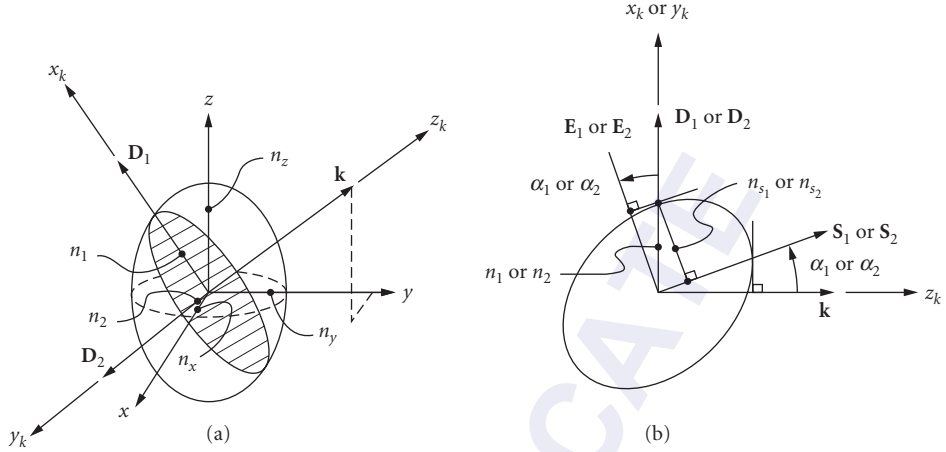


FIGURE 3 (a) The index ellipsoid cross section (crosshatched) that is normal to the wavevector \mathbf{k} has the shape of an ellipse. The major and minor axes of this ellipse represent the directions of the allowed polarizations \mathbf{D}_1 and \mathbf{D}_2 and (b) for each eigenpolarization (1 or 2) the vectors \mathbf{D} , \mathbf{E} , \mathbf{S} , and \mathbf{k} are coplanar.²¹

orthogonal vibration directions of \mathbf{D} for that particular direction of propagation. The lengths of these axes correspond to the phase velocity refractive indices. They are, therefore, referred to as the “fast” and “slow” axes. Figure 3b illustrates the field relationships with respect to the index ellipsoid. The line in the $(\mathbf{k}, \mathbf{D}_i)$ plane ($i=1$ or 2) that is tangent to the ellipsoid at \mathbf{D}_i is parallel to the ray vector \mathbf{S}_i ; the electric field \mathbf{E}_i also lies in the $(\mathbf{k}, \mathbf{D}_i)$ plane and is normal to \mathbf{S}_i . The line length denoted by n_{s_i} gives the ray velocity as $v_{s_i} = c/n_{s_i}$ for \mathbf{S}_i . The same relationships hold for either vibration, \mathbf{D}_1 or \mathbf{D}_2 .

In the general ellipsoid for a biaxial crystal there are two cross sections passing through the center that are circular. The normals to these cross sections are called the *optic axes* (denoted in Fig. 2 in a nonprincipal coordinate system), and they are coplanar and symmetric about the z principal axis in the x - z plane. The angle ϑ of an optic axis with respect to the z axis in the x - z plane is

$$\tan \vartheta = \frac{n_z}{n_x} \sqrt{\frac{n_y^2 - n_x^2}{n_z^2 - n_y^2}} \quad (4)$$

The phase velocities for \mathbf{D}_1 and \mathbf{D}_2 are equal for these two directions: $v_1 = v_2 = c/n_y$. In an ellipsoid of revolution for a uniaxial crystal, there is one circular cross section perpendicular to the z principal axis. Therefore, the z axis is the optic axis, and $\vartheta = 0^\circ$ in this case.

7.4 THE ELECTRO-OPTIC EFFECT

At an atomic level, an electric field applied to certain crystals causes a redistribution of bond charges and possibly a slight deformation of the crystal lattice.¹⁶ In general, these alterations are not isotropic; that is, the changes vary with direction in the crystal. Therefore, the dielectric tensor and its inverse, the impermeability tensor, change accordingly. The linear electro-optic effect, or Pockels, is a change in the impermeability tensor elements that is proportional to the magnitude of the externally applied electric field. Only crystals lacking a center of symmetry or macroscopically ordered dipolar molecules exhibit the Pockels effect. On the other hand, all materials, including amorphous materials and liquids, exhibit a quadratic (Kerr) electro-optic effect. The changes in the impermeability tensor elements are proportional to the square of the applied field. When the linear effect is present, it generally dominates over the quadratic effect.

Application of the electric field induces changes in the index ellipsoid and the impermeability tensor of Eq. (3b) according to

$$\mathbf{X}^T \left[\frac{1}{n^2} + \Delta \frac{1}{n^2} \right] \mathbf{X} = 1 \quad (5)$$

where the perturbation is

$$\Delta \left[\frac{1}{n^2} \right] = \left[\frac{1}{n^2} \right]_{E \neq 0} - \left[\frac{1}{n^2} \right]_{E=0} = \sum_{k=1}^3 r_{ijk} E_k + \sum_{k=1}^3 \sum_{l=1}^3 s_{ijkl} E_k E_l \quad (6)$$

Since n is dimensionless, and the applied electric field components are in units of V/m, the units of the linear r_{ijk} coefficients are in m/V and the quadratic coefficients s_{ijkl} are in m^2/V^2 .

The linear electro-optic effect is represented by a third rank tensor r_{ijk} with $3^3 = 27$ independent elements, that if written out in full form, will form the shape of a cube. The permutation symmetry of this tensor is $r_{ijk} = r_{ikj}$, $i, j, k = 1, 2, 3$ and this symmetry reduces the number of independent elements to 18.²² Therefore, the tensor can be represented in contracted notation by a 6×3 matrix; that is, $r_{ijk} \Rightarrow r_{ij}$, $i = 1, \dots, 6$ and $j = 1, 2, 3$. The first suffix is the same in both the tensor and the contracted matrix notation, but the second two tensor suffixes are replaced by a single suffix according to the following relation.

Tensor notation	11	22	33	23,32	31,13	12,21
Matrix notation	1	2	3	4	5	6

Generally, the r_{ij} coefficients have very little dispersion in the optical transparent region of a crystal.²³ The electro-optic coefficient matrices for all crystal classes are given in Table 1. References 16, 23, 24, and 25, among others, contain extensive tables of numerical values for indices and electro-optic coefficients for different materials.

The quadratic electro-optic effect is represented by a fourth rank tensor s_{ijkl} . The permutation symmetry of this tensor is $s_{ijkl} = s_{jikl} = s_{ijlk}$, $i, j, k, l = 1, 2, 3$. The tensor can be represented by a 6×6 matrix; that is, $s_{ijkl} \Rightarrow s_{kl}$, $k, l = 1, \dots, 6$. The quadratic electro-optic coefficient matrices for all crystal classes are given in Table 2. Reference 16 contains a table of quadratic electro-optic coefficients for several materials.

The Linear Electro-Optic Effect

An electric field applied in a general direction to a noncentrosymmetric crystal produces a linear change in the constants $(1/n^2)_i$, due to the linear electro-optic effect according to

$$\Delta(1/n^2)_i = \sum_j r_{ij} E_j \quad \begin{matrix} i = 1, \dots, 6 \\ j = x, y, z = 1, 2, 3 \end{matrix} \quad (7)$$

where r_{ij} is the ij th element of the linear electro-optic tensor in contracted notation. In matrix form Eq. (7) is

$$\begin{pmatrix} \Delta(1/n^2)_1 \\ \Delta(1/n^2)_2 \\ \Delta(1/n^2)_3 \\ \Delta(1/n^2)_4 \\ \Delta(1/n^2)_5 \\ \Delta(1/n^2)_6 \end{pmatrix} = \begin{pmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \\ r_{41} & r_{42} & r_{43} \\ r_{51} & r_{52} & r_{53} \\ r_{61} & r_{62} & r_{63} \end{pmatrix} \begin{pmatrix} E_x \\ E_y \\ E_z \end{pmatrix} \quad (8)$$

TABLE 1 The Linear Electro-Optic Coefficient Matrices in Contracted Form for All Crystal Symmetry Classes¹⁶

Centrosymmetric ($\bar{1}$, 2/m, mmm, 4/m, 4/mmm, $\bar{3}$, $\bar{3}$ m6/m, 6/mmm, m3, m3m):

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

Triclinic:

$$I^* \begin{pmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \\ r_{41} & r_{42} & r_{43} \\ r_{51} & r_{52} & r_{53} \\ r_{61} & r_{62} & r_{63} \end{pmatrix}$$

Cubic:

$$\begin{matrix} \bar{4}3m,23 & 432 \\ \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ r_{41} & 0 & 0 \\ 0 & r_{41} & 0 \\ 0 & 0 & r_{41} \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

Monoclinic:

$$\begin{matrix} 2(2||x_2) & 2(2||x_3) \\ \begin{pmatrix} 0 & r_{12} & 0 \\ 0 & r_{22} & 0 \\ 0 & r_{32} & 0 \\ r_{41} & 0 & r_{43} \\ 0 & r_{52} & 0 \\ 6_{61} & 0 & r_{63} \end{pmatrix} & \begin{pmatrix} 0 & 0 & r_{13} \\ 0 & 0 & r_{23} \\ 0 & 0 & r_{33} \\ r_{41} & r_{42} & 0 \\ r_{51} & r_{52} & 0 \\ 0 & 0 & r_{63} \end{pmatrix} \\ m(m \perp x_2) & m(m \perp x_3) \\ \begin{pmatrix} r_{11} & 0 & r_{13} \\ r_{21} & 0 & r_{23} \\ r_{31} & 0 & r_{33} \\ 0 & r_{42} & 0 \\ r_{51} & 0 & r_{53} \\ 0 & r_{62} & 0 \end{pmatrix} & \begin{pmatrix} r_{11} & r_{12} & 0 \\ r_{21} & r_{22} & 0 \\ r_{31} & r_{32} & 0 \\ 0 & 0 & r_{43} \\ 0 & 0 & r_{53} \\ r_{61} & r_{62} & 0 \end{pmatrix} \end{matrix}$$

Tetragonal:

$$\begin{matrix} 4 & \bar{4} & 422 \\ \begin{pmatrix} 0 & 0 & r_{13} \\ 0 & 0 & r_{13} \\ 0 & 0 & r_{33} \\ r_{41} & r_{51} & 0 \\ r_{51} & -r_{41} & 0 \\ 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & r_{13} \\ 0 & 0 & -r_{13} \\ 0 & 0 & 0 \\ r_{41} & -r_{51} & 0 \\ r_{51} & r_{41} & 0 \\ 0 & 0 & r_{63} \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ r_{41} & 0 & 0 \\ 0 & -r_{41} & 0 \\ 0 & 0 & 0 \end{pmatrix} \\ 4mm & \bar{4}2m(2||x_1) \\ \begin{pmatrix} 0 & 0 & r_{13} \\ 0 & 0 & r_{13} \\ 0 & 0 & r_{33} \\ 0 & r_{51} & 0 \\ r_{51} & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ r_{41} & 0 & 0 \\ 0 & r_{41} & 0 \\ 0 & 0 & r_{63} \end{pmatrix} \end{matrix}$$

Orthorhombic:

$$\begin{matrix} 222 & 2mm \\ \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ r_{41} & 0 & 0 \\ 0 & r_{52} & 0 \\ 0 & 0 & r_{63} \end{pmatrix} & \begin{pmatrix} 0 & 0 & r_{13} \\ 0 & 0 & r_{23} \\ 0 & 0 & r_{33} \\ 0 & r_{42} & 0 \\ r_{51} & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

Trigonal:

$$\begin{matrix} 3 & 32 \\ \begin{pmatrix} r_{11} & -r_{22} & r_{13} \\ -r_{11} & -r_{22} & r_{13} \\ 0 & 0 & r_{33} \\ r_{41} & r_{51} & 0 \\ r_{51} & -r_{41} & 0 \\ -r_{22} & -r_{11} & 0 \end{pmatrix} & \begin{pmatrix} r_{11} & 0 & 0 \\ -r_{11} & 0 & 0 \\ 0 & 0 & 0 \\ r_{41} & 0 & 0 \\ 0 & -r_{41} & 0 \\ 0 & -r_{11} & 0 \end{pmatrix} \end{matrix}$$

Hexagonal:

$$\begin{matrix} 6 & 6mm & 622 \\ \begin{pmatrix} 0 & 0 & r_{13} \\ 0 & 0 & r_{13} \\ 0 & 0 & r_{33} \\ r_{41} & r_{51} & 0 \\ r_{51} & -r_{41} & 0 \\ 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & r_{13} \\ 0 & 0 & r_{13} \\ 0 & 0 & r_{33} \\ 0 & r_{51} & 0 \\ r_{51} & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ r_{41} & 0 & 0 \\ 0 & -r_{41} & 0 \\ 0 & 0 & 0 \end{pmatrix} \\ \bar{6} & \bar{6}m2(m \perp x_1) & \bar{6}m2(m \perp x_2) \\ \begin{pmatrix} r_{11} & -r_{22} & 0 \\ -r_{11} & r_{22} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ -r_{22} & r_{11} & 0 \end{pmatrix} & \begin{pmatrix} 0 & -r_{22} & 0 \\ 0 & r_{22} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ -r_{22} & 0 & 0 \end{pmatrix} & \begin{pmatrix} r_{11} & 0 & 0 \\ -r_{11} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & -r_{11} & 0 \end{pmatrix} \\ 3m(m \perp x_1) & 3m(m \perp x_2) \\ \begin{pmatrix} 0 & -r_{22} & r_{13} \\ 0 & r_{22} & r_{13} \\ 0 & 0 & r_{33} \\ 0 & r_{51} & 0 \\ r_{51} & 0 & 0 \\ -r_{22} & 0 & 0 \end{pmatrix} & \begin{pmatrix} r_{11} & 0 & r_{13} \\ -r_{11} & 0 & r_{13} \\ 0 & 0 & r_{33} \\ 0 & r_{51} & 0 \\ r_{51} & 0 & 0 \\ 0 & -r_{11} & 0 \end{pmatrix} \end{matrix}$$

¹⁶The symbol over each matrix is the conventional symmetry-group designation.

TABLE 2 The Quadratic Electro-Optic Coefficient Matrices in Contracted Form for All Crystal Symmetry Classes¹⁶

Triclinic:

$$1, \bar{1}$$

$$\begin{pmatrix} s_{11} & s_{12} & s_{13} & s_{14} & s_{15} & s_{16} \\ s_{21} & s_{22} & s_{23} & s_{24} & s_{25} & s_{26} \\ s_{31} & s_{32} & s_{33} & s_{34} & s_{35} & s_{36} \\ s_{41} & s_{42} & s_{43} & s_{44} & s_{45} & s_{46} \\ s_{51} & s_{52} & s_{53} & s_{54} & s_{55} & s_{56} \\ s_{61} & s_{62} & s_{63} & s_{64} & s_{65} & s_{66} \end{pmatrix}$$

Monoclinic:

$$2, m, 2/m$$

$$\begin{pmatrix} s_{11} & s_{12} & s_{13} & 0 & s_{15} & 0 \\ s_{21} & s_{22} & s_{23} & 0 & s_{25} & 0 \\ s_{31} & s_{32} & s_{33} & 0 & s_{35} & 0 \\ 0 & 0 & 0 & s_{44} & 0 & s_{46} \\ s_{51} & s_{52} & s_{53} & 0 & s_{55} & 0 \\ 0 & 0 & 0 & s_{64} & 0 & s_{66} \end{pmatrix}$$

Orthorhombic:

$$2mm, 222, mmm$$

$$\begin{pmatrix} s_{11} & s_{12} & s_{13} & 0 & 0 & 0 \\ s_{21} & s_{22} & s_{23} & 0 & 0 & 0 \\ s_{31} & s_{32} & s_{33} & 0 & 0 & 0 \\ 0 & 0 & 0 & s_{44} & 0 & 0 \\ 0 & 0 & 0 & 0 & s_{55} & 0 \\ 0 & 0 & 0 & 0 & 0 & s_{66} \end{pmatrix}$$

Tetragonal:

$$4, \bar{4}, 4/m$$

$$\begin{pmatrix} s_{11} & s_{12} & s_{13} & 0 & 0 & s_{16} \\ s_{12} & s_{11} & s_{13} & 0 & 0 & -s_{16} \\ s_{31} & s_{31} & s_{33} & 0 & 0 & 0 \\ 0 & 0 & 0 & s_{44} & s_{45} & 0 \\ 0 & 0 & 0 & -s_{45} & s_{44} & 0 \\ s_{61} & -s_{61} & 0 & 0 & 0 & s_{66} \end{pmatrix}$$

$$422, 4mm, \bar{4}2m, 4/m\bar{2}$$

$$\begin{pmatrix} s_{11} & s_{12} & s_{13} & 0 & 0 & 0 \\ s_{12} & s_{11} & s_{13} & 0 & 0 & 0 \\ s_{31} & s_{31} & s_{33} & 0 & 0 & 0 \\ 0 & 0 & 0 & s_{44} & 0 & 0 \\ 0 & 0 & 0 & 0 & s_{44} & 0 \\ 0 & 0 & 0 & 0 & 0 & s_{66} \end{pmatrix}$$

Hexagonal:

$$6, \bar{6}, 6/m$$

$$\begin{pmatrix} s_{11} & s_{12} & s_{13} & 0 & 0 & -s_{61} \\ s_{12} & s_{11} & s_{13} & 0 & 0 & s_{61} \\ s_{31} & s_{31} & s_{33} & 0 & 0 & 0 \\ 0 & 0 & 0 & s_{44} & s_{45} & 0 \\ 0 & 0 & 0 & -s_{45} & s_{44} & 0 \\ s_{61} & -s_{61} & 0 & 0 & 0 & \frac{1}{2}(s_{11} - s_{12}) \end{pmatrix}$$

$$622, 6mm, \bar{6}m2, 6/mmm$$

$$\begin{pmatrix} s_{11} & s_{12} & s_{13} & 0 & 0 & 0 \\ s_{12} & s_{11} & s_{13} & 0 & 0 & 0 \\ s_{31} & s_{31} & s_{33} & 0 & 0 & 0 \\ 0 & 0 & 0 & s_{44} & 0 & 0 \\ 0 & 0 & 0 & 0 & s_{44} & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{2}(s_{11} - s_{12}) \end{pmatrix}$$

Cubic:

$$23, m\bar{3}$$

$$\begin{pmatrix} s_{11} & s_{12} & s_{13} & 0 & 0 & 0 \\ s_{13} & s_{11} & s_{12} & 0 & 0 & 0 \\ s_{12} & s_{13} & s_{11} & 0 & 0 & 0 \\ 0 & 0 & 0 & s_{44} & 0 & 0 \\ 0 & 0 & 0 & 0 & s_{44} & 0 \\ 0 & 0 & 0 & 0 & 0 & s_{44} \end{pmatrix}$$

(Continued)

TABLE 2 The Quadratic Electro-optic Coefficient Matrices in Contracted Form for All Crystal Symmetry Classes¹⁶ (Continued)

Trigonal:	Cubic:
$3, \bar{3}$ $\begin{pmatrix} s_{11} & s_{12} & s_{13} & s_{14} & s_{15} & -s_{61} \\ s_{12} & s_{11} & s_{13} & -s_{14} & -s_{15} & s_{61} \\ s_{31} & s_{31} & s_{33} & 0 & 0 & 0 \\ s_{41} & -s_{41} & 0 & s_{44} & s_{45} & -s_{51} \\ s_{51} & -s_{51} & 0 & -s_{45} & s_{44} & s_{41} \\ s_{61} & -s_{61} & 0 & -s_{15} & s_{14} & \frac{1}{2}(s_{11} - s_{12}) \end{pmatrix}$	$432, m\bar{2}m, \bar{4}3m$ $\begin{pmatrix} s_{11} & s_{12} & s_{12} & 0 & 0 & 0 \\ s_{12} & s_{11} & s_{12} & 0 & 0 & 0 \\ s_{12} & s_{12} & s_{11} & 0 & 0 & 0 \\ 0 & 0 & 0 & s_{44} & 0 & 0 \\ 0 & 0 & 0 & 0 & s_{44} & 0 \\ 0 & 0 & 0 & 0 & 0 & s_{44} \end{pmatrix}$
$32, 3m, \bar{3}m$ $\begin{pmatrix} s_{11} & s_{12} & s_{13} & s_{14} & 0 & 0 \\ s_{12} & s_{11} & s_{13} & -s_{14} & 0 & 0 \\ s_{13} & s_{13} & s_{33} & 0 & 0 & 0 \\ s_{41} & -s_{41} & 0 & s_{44} & 0 & 0 \\ 0 & 0 & 0 & 0 & s_{44} & s_{41} \\ 0 & 0 & 0 & 0 & s_{14} & \frac{1}{2}(s_{11} - s_{12}) \end{pmatrix}$	Isotropic: $\begin{pmatrix} s_{11} & s_{12} & s_{12} & 0 & 0 & 0 \\ s_{12} & s_{11} & s_{12} & 0 & 0 & 0 \\ s_{12} & s_{12} & s_{11} & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{2}(s_{11} - s_{12}) & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{2}(s_{11} - s_{12}) & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{2}(s_{11} - s_{12}) \end{pmatrix}$

E_x , E_y , and E_z are the components of the applied electric field in principal coordinates. The magnitude of $\Delta(1/n^2)$ is typically on the order of less than 10^{-5} . Therefore, these changes are mathematically referred to as perturbations. The new impermeability tensor $[1/n^2]'$ in the presence of an applied electric field is no longer diagonal in the reference principal dielectric axes system. It is given by

$$[1/n^2]' = \begin{pmatrix} 1/n_x^2 + \Delta(1/n^2)_1 & \Delta(1/n^2)_6 & \Delta(1/n^2)_5 \\ \Delta(1/n^2)_6 & 1/n_y^2 + \Delta(1/n^2)_2 & \Delta(1/n^2)_4 \\ \Delta(1/n^2)_5 & \Delta(1/n^2)_4 & 1/n_z^2 + \Delta(1/n^2)_3 \end{pmatrix} \quad (9)$$

and is determined by the unperturbed principal refractive indices, the electro-optic coefficients, and the direction of the applied field relative to the principal coordinate system. However, the field-induced perturbations are symmetric, so the symmetry of the tensor is not disturbed. The new index ellipsoid is now represented by

$$(1/n^2)'_1 x^2 + (1/n^2)'_2 y^2 + (1/n^2)'_3 z^2 + 2(1/n^2)'_4 yz + 2(1/n^2)'_5 xz + 2(1/n^2)'_6 xy = 1 \quad (10)$$

or equivalently, $\mathbf{X}^T [1/n^2]' \mathbf{X} = 1$, where $\mathbf{X} = [x \ y \ z]^T$.^{19,26} The presence of cross terms indicates that the ellipsoid is rotated and the lengths of the principal dielectric axes are changed. Determining the new orientation and shape of the ellipsoid requires that $[1/n^2]'$ be diagonalized, thus determining its eigenvalues and eigenvectors. After diagonalization, in a suitably rotated new coordinate system $\mathbf{X}' = [x' \ y' \ z']^T$ the perturbed ellipsoid will then be represented by a square sum:

$$\frac{x'^2}{n_x'^2} + \frac{y'^2}{n_y'^2} + \frac{z'^2}{n_z'^2} = 1 \quad (11)$$

The eigenvalues of $[1/n^2]'$ are $1/n_x'^2$, $1/n_y'^2$, $1/n_z'^2$. The corresponding eigenvectors are $\mathbf{x}' = [x_x' \ y_x' \ z_x']^T$, and $\mathbf{y}' = [x_y' \ y_y' \ z_y']^T$, respectively.

The Quadratic or Kerr Electro-Optic Effect

An electric field applied in a general direction to any crystal, centrosymmetric or noncentrosymmetric, produces a quadratic change in the constants $(1/n^2)_i$ due to the quadratic electro-optic effect according to

$$\begin{pmatrix} \Delta(1/n^2)_1 \\ \Delta(1/n^2)_2 \\ \Delta(1/n^2)_3 \\ \Delta(1/n^2)_4 \\ \Delta(1/n^2)_5 \\ \Delta(1/n^2)_6 \end{pmatrix} = \begin{pmatrix} s_{11} & s_{12} & s_{13} & s_{14} & s_{15} & s_{16} \\ s_{21} & s_{22} & s_{23} & s_{24} & s_{25} & s_{26} \\ s_{31} & s_{32} & s_{33} & s_{34} & s_{35} & s_{36} \\ s_{41} & s_{42} & s_{43} & s_{44} & s_{45} & s_{46} \\ s_{51} & s_{52} & s_{53} & s_{54} & s_{55} & s_{56} \\ s_{61} & s_{62} & s_{63} & s_{64} & s_{65} & s_{66} \end{pmatrix} \begin{pmatrix} E_x^2 \\ E_y^2 \\ E_z^2 \\ E_y E_z \\ E_x E_z \\ E_x E_y \end{pmatrix} \quad (12)$$

E_x , E_y , and E_z are the components of the applied electric field in principal coordinates. The perturbed impermeability tensor and the new index ellipsoid have the same form as Eqs. (9) and (10).

Normally, there are two distinctions made when considering the Kerr effect: the ac Kerr effect and the dc Kerr effect. The induced changes in the optical properties of the material can occur as the result of a slowly varying applied electric field, or it can result from the electric field of the light itself. The former is the dc Kerr effect and the latter is the ac or optical Kerr effect. The dc Kerr effect is given by

$$\Delta n = \lambda K E^2 \quad (13)$$

where K is the Kerr constant in units of m/V^2 and λ is the freespace wavelength. Some polar liquids such as nitrobenzene ($\text{C}_6\text{H}_5\text{NO}_2$), which is poisonous, exhibit very large Kerr constants which are much greater than those of transparent crystals.²⁷ In contrast to the linear Pockels electro-optic effect, larger voltages are required for any significant Kerr modulation.

The ac or optical Kerr effect occurs when a very intense beam of light modulates the optical material. The optical Kerr effect is given by

$$n = n_o + n_2 I \quad (14)$$

which describes the intensity dependent refractive index n , where n_o is the unmodulated refractive index, n_2 is the second order nonlinear refractive index (m^2/W), and I is the intensity of the wave (W). Equation 14 is derived from the expression for the electric field induced polarization in a material as a function of the linear and nonlinear susceptibilities. This intensity dependent refractive index is used as the basis for Kerr-lens mode-locking of ultrashort lasers. It is also responsible for nonlinear effects of self-focusing and self-phase modulation.²⁸

A Mathematical Approach: The Jacobi Method

The analytical design and study of electro-optic modulators require robust mathematical techniques due to the small, anisotropic perturbations to the refractive index profile of a material. Especially with the newer organic crystals, polymers, and tailored nanostructured materials, the properties are often biaxial before and after applied voltages. The optimum modulator configuration may not be along principal axes. In addition, sensitivities in modulation characteristics of biaxial materials (natural and/or induced) can negatively impact something as simple as focusing a beam onto the material. Studying the electro-optic effect is basically an eigenvalue problem.

Although the eigenvalue problem is a familiar one, obtaining accurate solutions has been the subject of extensive study.^{29–32} A number of formalisms are suggested in the literature to address the specific problem of finding the new set of principal dielectric axes relative to the zero-field principal

dielectric axes. Most approaches, however, do not provide a consistent means of labeling the new axes. Also, some methods are highly susceptible to numerical instabilities when dealing with very small off-diagonal elements as in the case of the electro-optic effect. In contrast to other methods,^{15,22,26,30,33,34} a similarity transformation is an attractive approach for diagonalizing a symmetric matrix for the purpose of determining its eigenvalues and eigenvectors.^{21,29,30,32,35}

The Jacobi method utilizes the concepts of rigid-body rotation and the properties of ellipsoids to determine the principal axes and indices of a crystal by constructing a series of similarity transformations that consist of elementary plane rotations. The method produces accurate eigenvalues and orthogonal eigenvectors for matrices with very small off-diagonal elements, and it is a systematic procedure for ordering the solutions to provide consistent labeling of the principal axes.^{21,31} The sequence of transformations are applied to the perturbed index ellipsoid and convert from one set of orthogonal axes $\mathbf{X} = [x, y, z]$ to another set $[x', y', z']$, until a set of axes coincides with the new principal dielectric directions and the impermeability matrix is diagonalized. Since similarity is a transitive property, several transformation matrices can be multiplied to generate the desired cumulative matrix.^{21,29} Thus, the problem of determining the new principal axes and indices of refraction of the index ellipsoid in the presence of an external electric field is analogous to the problem of finding the cumulative transformation matrix $[a] = [a_m] \cdots [a_2] [a_1]$ that will diagonalize the perturbed impermeability tensor. The transformation required matrix, $[a]$, is simply the product of the elementary plane rotation matrices multiplied in the order in which they are applied.

When plane rotations are applied to the matrix representation of tensors, the magnitude of a physical property can be evaluated in any arbitrary direction. When the matrix is transformed to diagonal form, the eigenvalues lie on the diagonal and the eigenvectors are found in the rows or columns of the corresponding transformation matrices. Specifically, a symmetric matrix $[\mathbf{A}]$ can be reduced to diagonal form by the transformation $[a][\mathbf{A}][a]^T = [\lambda]$, where $[\lambda]$ is a 3×3 diagonal matrix and $[a]$ is the orthogonal transformation matrix. Since the eigenvalues of $[\mathbf{A}]$ are preserved under similarity transformation, they lie on the diagonal of $[\lambda]$, as in Eq. (1).

In terms of the index ellipsoid, first recall that the perturbed index ellipsoid in the original (zero field) coordinate system is $\mathbf{X}^T [1/n^2]' \mathbf{X} = 1$, where $[1/n^2]'$ is given by Eq. (9). A suitable matrix, $[a]$, will relate the "new" principal axes \mathbf{X}' of the perturbed ellipsoid to the "old" coordinate system; that is, $\mathbf{X}' = [a]\mathbf{X}$, or $\mathbf{X} = [a]^T \mathbf{X}'$. Substituting these relationships into the index ellipsoid results in

$$\begin{aligned} ([a]^T \mathbf{X}')^T [1/n^2]' [a]^T \mathbf{X}' &= 1 \\ \mathbf{X}'^T [a] [1/n^2]' [a]^T \mathbf{X}' &= 1 \\ \mathbf{X}'^T [1/n^2]'' \mathbf{X}' &= 1 \end{aligned} \quad (15)$$

where $[a] [1/n^2]' [a]^T = [1/n^2]''$ is the diagonalized impermeability matrix in the new coordinate system and $[a]$ is the cumulative transformation matrix.

Using the Jacobi method, each simple elementary plane rotation that is applied at each step will zero an off-diagonal element of the impermeability tensor. The goal is to produce a diagonal matrix by minimizing the norm of the off-diagonal elements to within a desired level of accuracy. If m transformations are required, each step is represented by

$$[1/n^2]''_m = [a_m] [1/n^2]''_{m-1} [a_m]^T \quad (16)$$

To determine the form of each elementary plane rotation, the Jacobi method begins by first selecting the largest off-diagonal element $(1/n^2)''_{ij}$ and executing a rotation in the (i, j) plane, $i < j$, so as to zero that element. The required rotation angle Ω is given by

$$\tan(2\Omega) = \left(\frac{2(1/n^2)''_{ij}}{(1/n^2)''_{ii} - (1/n^2)''_{jj}} \right) \quad i, j = 1, 2, 3 \quad (17)$$

For example, if the largest off-diagonal element is $(1/n_{12}^2) = (1/n_{21}^2)$, then the plane rotation is represented by

$$[a] = \begin{pmatrix} \cos \Omega & \sin \Omega & 0 \\ -\sin \Omega & \cos \Omega & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (18)$$

which is a counter clockwise rotation about the three-axis. If $(1/n^2)_{ii} = (1/n^2)_{jj}$, which can occur in isotropic and uniaxial crystals, then $|\Omega|$ is taken to be 45° , and its sign is taken to be the same as the sign of $(1/n^2)_{ij}$. The impermeability matrix elements are updated with the following equations, which are calculated from the transformation of Eq. (15):

$$\begin{pmatrix} (1/n_{11}^2)' & 0 & (1/n_{13}^2)' \\ 0 & (1/n_{22}^2)' & (1/n_{23}^2)' \\ (1/n_{13}^2)' & (1/n_{23}^2)' & (1/n_{33}^2)' \end{pmatrix} = \begin{pmatrix} \cos \Omega & \sin \Omega & 0 \\ -\sin \Omega & \cos \Omega & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1/n_{11}^2 & 1/n_{12}^2 & 1/n_{13}^2 \\ 1/n_{12}^2 & 1/n_{22}^2 & 1/n_{23}^2 \\ 1/n_{13}^2 & 1/n_{23}^2 & 1/n_{33}^2 \end{pmatrix} \begin{pmatrix} \cos \Omega & -\sin \Omega & 0 \\ \sin \Omega & \cos \Omega & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (19)$$

Once the new elements are determined, the next iteration step is performed, selecting the new largest off-diagonal element and repeating the procedure with another suitable rotation matrix. The process is terminated when all of the off-diagonal elements are reduced below the desired level (typically 10^{-10}). The next step is to determine the cumulative transformation matrix $[a]$. One way is to multiply the plane rotation matrices in order, either as $[a] = [a_m] \cdots [a_2] [a_1]$ or equivalently for the transpose of $[a]$ as

$$[a]^T = [a_1]^T [a_2]^T \cdots [a_m]^T \quad (20)$$

The set of Euler angles, which also defines the orientation of a rigid body, can be obtained from the cumulative transformation matrix $[a]$.^{36,37} These angles are given in the Appendix. Several examples for using the Jacobi method are given in Ref. 21.

Determining the Eigenpolarizations and Phase Velocity Indices of Refraction

After the perturbed impermeability matrix is diagonalized, the polarization directions of the two allowed linear orthogonal waves \mathbf{D}_1 and \mathbf{D}_2 that propagate independently for a given wavevector \mathbf{k} can be determined along with their respective phase velocity refractive indices $n_{x''}$ and $n_{y''}$. These waves are the only two that can propagate with unchanging orientation for the given wavevector direction. Figure 4a depicts these axes for a crystal in the absence of an applied field. Figure 4b depicts the x''' and y''' axes, which define the fast and slow axes, when an electric field is applied in a direction so as to reorient the index ellipsoid. The applied field, in general, rotates the allowed polarization directions in the plane perpendicular to the direction of phase propagation as shown in Fig. 4b. Determining these "eigenpolarizations," that is, \mathbf{D}_1 , \mathbf{D}_2 , n_1 , and n_2 , is also an eigenvalue problem.

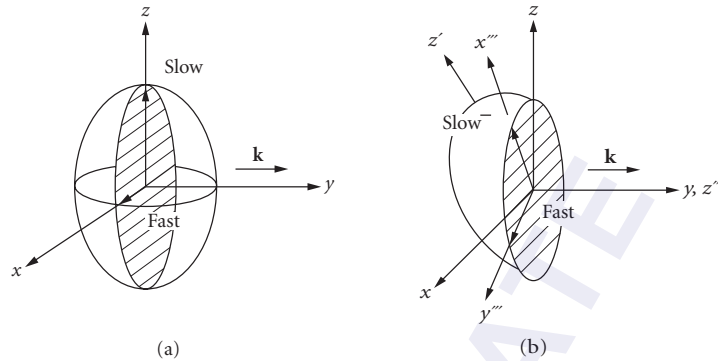


FIGURE 4 (a) The cross-section ellipse for a wave propagating along the y principal axis is shown with no field applied to the crystal; (b) with an applied electric field the index ellipsoid is reoriented, and the eigenpolarizations in the plane transverse to \mathbf{k} are rotated, indicated by x''' and y''' .

The perturbed index ellipsoid resulting from an external field was given by Eq. (10) in the original principal-axis coordinate system. For simplicity, the coefficients may be relabeled as

$$Ax^2 + By^2 + Cz^2 + 2Fyz + 2Gxz + 2Hxy = 1$$

or

$$\mathbf{X}^T \begin{pmatrix} A & H & G \\ H & B & F \\ G & F & C \end{pmatrix} \mathbf{X} = 1 \quad (21)$$

where x , y , and z represent the original dielectric axes with no applied field and $\mathbf{X}^T = [x \ y \ z]$. However, before the eigenpolarizations can be determined, the direction of light propagation, \mathbf{k} , through the material whose index ellipsoid has been perturbed, must be defined. The problem is then to determine the allowed eigenpolarizations and phase velocity refractive indices associated with this direction of propagation. The optical wavevector direction \mathbf{k} is conveniently specified by the spherical coordinates angles θ_k and ϕ_k in the (x, y, z) coordinate system as shown in Fig. 5. Given \mathbf{k} , the cross

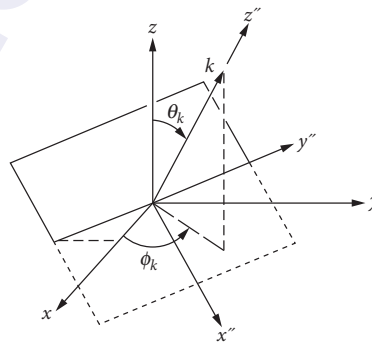


FIGURE 5 The coordinate system (x'', y'', z'') of the wavevector \mathbf{k} is defined with its angular relationship (ϕ_k, θ_k) with respect to the unperturbed principal dielectric axes coordinate system (x, y, z) .²¹

section ellipse through the center of the perturbed ellipsoid of Eq. (21) may be drawn. The directions of the semiaxes of this ellipse represent the fast and slow polarization directions of the two waves \mathbf{D}_1 and \mathbf{D}_2 that propagate independently. The lengths of the semiaxes are the phase velocity indices of refraction. The problem is to determine the new polarization directions x''' of \mathbf{D}_1 and y''' of \mathbf{D}_2 relative to the (x, y, z) axes and the corresponding new indices of refraction $n_{x'''}'$ and $n_{y'''}'$.

The first step is to do a transformation from the (x, y, z) (lab or principal axis) coordinate system to a coordinate system (x'', y'', z'') aligned with the direction of phase propagation. In this example, (x'', y'', z'') is chosen such that $z'' \parallel \mathbf{k}$, and x'' is lying in the (z, z'') plane. The (x'', y'', z'') system is, of course, different from the (x', y', z') perturbed principal axes system. Using the spherical coordinate angles of \mathbf{k} , the (x'', y'', z'') system may be produced first by a counterclockwise rotation ϕ_k about the z axis followed by a counterclockwise rotation θ_k about y'' as shown in Fig. 5. This transformation is described by $\mathbf{X}'' = [a]\mathbf{X}$, or $[a]^T \mathbf{X}'' = \mathbf{X}$ and is explicitly,

$$\mathbf{X} = \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} \cos \phi_k & -\sin \phi_k & 0 \\ \sin \phi_k & \cos \phi_k & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \cos \theta_k & 0 & -\sin \theta_k \\ 0 & 1 & 0 \\ \sin \theta_k & 0 & \cos \theta_k \end{pmatrix} \begin{pmatrix} x'' \\ y'' \\ z'' \end{pmatrix} \quad (22)$$

The equation for the cross section ellipse normal to \mathbf{k} is determined by substituting Eq. (22) into Eq. (21) and setting $z'' = 0$ or by matrix substitution as follows:

$$\begin{aligned} ([a]^T \mathbf{X}'')^T [1/n^2]' ([a]^T \mathbf{X}'') &= 1 \\ \mathbf{X}''^T \underbrace{[a] [1/n^2]' [a]^T}_{[1/n^2]''} \mathbf{X}'' &= 1 \end{aligned} \quad (23)$$

which results in

$$\mathbf{X}''^T [1/n^2]'' \mathbf{X}'' = 1$$

or

$$\begin{pmatrix} x'' & y'' & z'' \end{pmatrix} \begin{pmatrix} A'' & H'' & G'' \\ H'' & B'' & F'' \\ G'' & F'' & C'' \end{pmatrix} \begin{pmatrix} x'' \\ y'' \\ z'' \end{pmatrix} = 1 \quad (24)$$

The coefficients of the cross section ellipse equation described by Eq. (24), with z'' set to zero, are used to determine the eigenpolarization directions and the associated phase velocity refractive indices for the chosen direction of propagation. The cross section ellipse normal to the wavevector direction $\mathbf{k} \parallel z''$ is represented by the 2×2 submatrix of $[1/n^2]''$:

$$\begin{pmatrix} x'' & y'' \end{pmatrix} \begin{pmatrix} A'' & H'' \\ H'' & B'' \end{pmatrix} \begin{pmatrix} x'' \\ y'' \end{pmatrix} = A''x''^2 + B''y''^2 + 2H''x''y'' = 1 \quad (25)$$

The polarization angle β_1 of x''' (\mathbf{D}_1) with respect to x'' , as shown in Fig. 6, is given by

$$\beta_1 = \frac{1}{2} \tan^{-1} \left[\frac{2H''}{(A'' - B'')} \right] \quad (26)$$

The polarization angle β_2 of y''' (\mathbf{D}_2) with respect to x'' is $\beta_1 + \pi/2$. The axes are related by a plane rotation $\mathbf{X}''' = [a_{\beta_1}]\mathbf{X}''$ or

$$\begin{pmatrix} x''' \\ y''' \end{pmatrix} = \begin{pmatrix} \cos \beta_1 & -\sin \beta_1 \\ \sin \beta_1 & \cos \beta_1 \end{pmatrix} \begin{pmatrix} x'' \\ y'' \end{pmatrix} \quad (27)$$

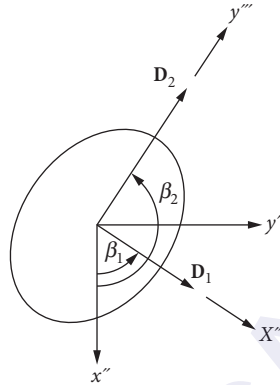


FIGURE 6 The polarization axes (x''' , y''') are the fast and slow axes and are shown relative to the (x'' , y'') axes of the wavevector coordinate system. The wavevector \mathbf{k} and the axes z'' and z''' are normal to the plane of the figure.²¹

The refractive indices, $n_{x'''}$ and $n_{y'''}$ may be found by performing one more rotation in the plane of the ellipse normal to \mathbf{k} , using the angle β_1 or β_2 and the rotation of Eq. (27). The result is a new matrix that is diagonalized,

$$[1/n_2]''' = \begin{pmatrix} 1/n_{x'''} & 0 \\ 0 & 1/n_{y'''} \end{pmatrix} = [a_{\beta_1}]^T [1/n_2]'' [a_{\beta_1}] \quad (28)$$

The larger index corresponds to the slow axis and the smaller index to the fast axis.

7.5 MODULATOR DEVICES

An electro-optic modulator is a device with operation based on an electrically induced change in index of refraction or change in natural birefringence. Depending on the device configuration, the following properties of the light wave can be varied in a controlled way: phase, polarization, amplitude, frequency, or direction of propagation. The device is typically designed for optimum performance at a single wavelength, with some degradation in performance with wideband or multimode lasers.^{16,38,39}

Electro-optic devices can be used in analog or digital modulation formats. The choice is dictated by the system requirements and the characteristics of available components (optical fibers, sources/detectors, etc.). Analog modulation requires large signal-to-noise ratios (SNR), thereby limiting its use to narrow-bandwidth, short-distance applications. Digital modulation, on the other hand, is more applicable to large-bandwidth, medium to long distance systems.^{38,39}

Device Geometries

A bulk electro-optic modulator can be classified as one of two types, *longitudinal* or *transverse*, depending on how the voltage is applied relative to the direction of light propagation in the device. Basically a bulk modulator consists of an electro-optic crystal sandwiched between a pair of electrodes and, therefore, can be modeled as a capacitor. In general, the input and output faces are parallel for the beam to undergo a uniform phase shift over the beam cross section.¹⁶ Waveguide modulators are discussed later in the section "Waveguide or Integrated-Optic Modulators" and have a variety of electrode configurations that are analogous to longitudinal and transverse orientations, although the distinction is not as well defined.

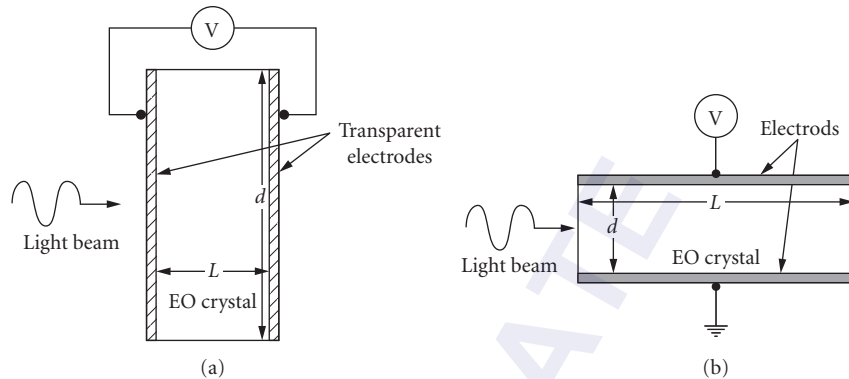


FIGURE 7 (a) A longitudinal electro-optic modulator has the voltage applied parallel to the direction of light propagation and (b) a transverse modulator has the voltage applied perpendicular to the direction of light propagation.¹⁶

In the bulk longitudinal configuration, the voltage is applied parallel to the wavevector direction as shown in Fig. 7a.^{16,25,40-43} The electrodes must be transparent to the light either by the choice of material used for them (metal-oxide coatings of SnO, InO, or CdO) or by leaving a small aperture at their center at each end of the electro-optic crystal.^{25,41-43} The ratio of the crystal length L to the electrode separation b is defined as the *aspect ratio*. For this configuration $b = L$, and, therefore, the aspect ratio is always unity. The magnitude of the applied electric field inside the crystal is $E = V/L$. The induced phase shift is proportional to V and the wavelength λ of the light but not the physical dimensions of the device. Therefore, for longitudinal modulators, the required magnitude of the applied electric field for a desired degree of modulation cannot be reduced by changing the aspect ratio, and it increases with wavelength. However, these modulators can have a large acceptance area and are useful if the light beam has a large cross-sectional area.

In the transverse configuration, the voltage is applied perpendicular to the direction of light propagation as shown in Fig. 7b.^{16,40-43} The electrodes do not obstruct the light as it passes through the crystal. For this case, the aspect ratio can be very large. The magnitude of the applied electric field is $E = V/d$, ($b = d$), and d can be reduced to increase E for a given applied voltage, thereby increasing the aspect ratio L/b . The induced phase shift is inversely proportional to the aspect ratio; therefore, the voltage necessary to achieve a desired degree of modulation can be greatly reduced. Furthermore,

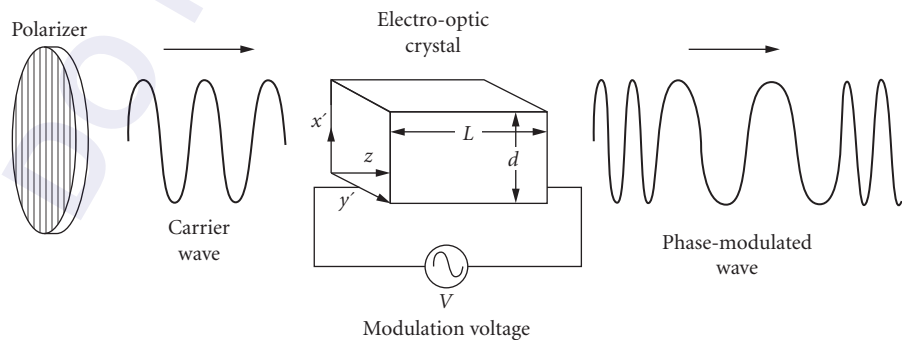


FIGURE 8 A longitudinal phase modulator is shown with the light polarized along the new x' principal axis when the modulation voltage V is applied.¹⁶

the interaction length can be long for a given field strength. However, the transverse dimension d is limited by the increase in capacitance, which affects the modulation bandwidth or speed of the device, and by diffraction for a given length L , since a beam with finite cross section diverges as it propagates.^{16,41,44}

Bulk Modulators

The modulation of phase, polarization, amplitude, frequency, and position of light can be implemented using an electro-optic bulk modulator with polarizers and passive birefringent elements. Three assumptions are made in this section. First, the modulating field is *uniform* throughout the length of the crystal; the change in index or birefringence is uniform unless otherwise stated. Second, the modulation voltage is dc or very low radian frequency ω_m ($\omega_m \ll 2\pi/\tau$); the light experiences the same induced Δn during its transit time τ through the crystal of length L , and the capacitance is negligible. Finally, light propagation is taken to be along a principal axis, before and after the voltage is applied; therefore, the equations for the eigenpolarizations are presented in terms of the *optical electric field* \mathbf{E} , rather than the *displacement vector* \mathbf{D} , which is a common practice in various optical references. For other general configurations, the equations should be expressed in terms of the eigenpolarizations \mathbf{D}_1 and \mathbf{D}_2 . However, the electric field will determine the direction of energy flow, which is generally not in the same direction as the wavevector. References 16 and 41, among others, provide examples of modulator devices using potassium dihydrogen phosphate (KDP), lithium niobate (LiNbO_3), lithium tantalate (LiTaO_3), gallium arsenide (GaAs), and barium titanate (BaTiO_3).

Phase Modulator A light wave can be phase modulated, without change in polarization or intensity, using an electro-optic crystal and an input polarizer in the proper configuration. This is the simplest electro-optic modulator. As an example, consider a longitudinal device that is made of a LiNbO_3 crystal as shown in Fig. 8, with the voltage applied in the z direction. In general, an applied voltage V will rotate the principal axes in the crystal cross section. However, for phase modulation, the input polarizer must be aligned parallel to one of the principal axes that will have a preserved orientation when the voltage is on or off. The LiNbO_3 crystal is uniaxial with symmetry $3m$, and it is a common material for electro-optic modulators. Substituting into Eq. (7) results in

$$\begin{pmatrix} 0 & -r_{22} & r_{13} \\ 0 & r_{22} & r_{13} \\ 0 & 0 & r_{33} \\ 0 & r_{51} & 0 \\ r_{51} & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ E_z \end{pmatrix} = \begin{pmatrix} r_{13}E_z \\ r_{13}E_z \\ r_{33}E_z \\ 0 \\ 0 \\ 0 \end{pmatrix} \quad (29)$$

Figure 8 indicates the polarizer along x' provides an input optical electric field $E_{\text{in}}(t) = E_x \cos \omega t$. The voltage is applied in the z direction and thus, has only one component, E_z . In this case, the perturbed index ellipsoid will be

$$x^2(1/n_o^2 + r_{13}E_z) + y^2(1/n_o^2 + r_{13}E_z) + z^2(1/n_e^2 + r_{33}E_z) = 1 \quad (30)$$

where $n_x = n_y = n_o$ is the ordinary refractive index and $n_z = n_e$ is the extraordinary refractive index of this uniaxial crystal. The orientation of the principal axes remains unchanged, but the lengths of the axes, and hence the phase velocity refractive indices, have been modified by the applied electric field. The optical wave at the output of the crystal at $z = L$ is

$$E_{\text{out}}(t) = E_x \cos(\omega t - \phi) \quad (31)$$

where

$$\phi = \frac{2\pi}{\lambda}(n_{x'} + \Delta n_{x'})L = \phi_o + \Delta\phi_{x'} \quad (32)$$

is the total phase shift consisting of a natural phase term $\phi_o = (2\pi/\lambda)Ln_{x'}$, with $n_{x'}$ being the unperturbed index in the x' ($= x$, in this example) direction, and an electrically induced phase term $\Delta\phi_{x'} = (2\pi/\lambda)L\Delta n_{x'}$ for a polarization along x' . The new x' axis has length $2n_{x'}$, where

$$(1/n_{x'})^2 = (1/n_o^2) + r_{13}E_z \quad (33)$$

Using the derivative $\Delta(1/n^2)' = -2n^{-3}\Delta n$ results in

$$\Delta n_{x'}' \approx (-1/2)n_o^3 r_{13} E_z \quad (34)$$

an approximation often seen in literature.

For a longitudinal modulator the applied electric field is $E_z = V/L$, and the induced phase shift is $\Delta\phi_{x'} = \pi/\lambda n_{x'}^3 r V$, which is independent of L and is linearly related to V . For a transverse modulator $E = V/d$ and the induced phase shift is $\Delta\phi_{x'} = \pi/\lambda n_{x'}^3 r V(L/d)$, which is a function of the aspect ratio L/d and V . The voltage that would produce an induced phase shift of $\Delta\phi_{x'} = \pi$ is the *half-wave voltage*. The half-wave voltage is $V_\pi = \lambda/n_{x'}^3 r$ for a longitudinal modulator and $\tilde{V}_\pi = (\lambda/n_{x'}^3 r)(d/L)$ for a transverse modulator.

Whenever possible, it is desired to take advantage of the largest electro-optic coefficient. In LiNbO_3 , $r_{13} = 10$ and $r_{33} = 32$. To use r_{33} , the optical input signal should be polarized in the z direction.

If a dc applied voltage is used, one of two possibilities is required for a crystal and its orientation. The first possibility is a crystal having principal axes which will not rotate with applied voltage V ; an example is LiNbO_3 with V applied in the z direction and an input polarization along the $x' = x$ axis propagating along $z' = z$. The second possibility is a crystal having a characteristic plane, that is, a plane exhibiting a constant refractive index perpendicular to the direction of propagation. If a field is applied such that the axes rotate in this plane, the input wave must be polarized along one of the new principal axes. Therefore, it will always be polarized along a principal axis, whether the voltage is on or off. An example is KDP with V along the z axis and an input wave polarized along the new principal axis x' and propagating along $z' = z$. Phase modulation is then achieved by turning the voltage on and off.

If the applied modulation voltage is sinusoidal in time ($V = V_m \sin \omega_m t$), the corresponding electric field can be represented by $E = E_m \sin \omega_m t$.

The magnitude of the field varies only with time, not space; it is a stationary wave applied in the same direction for all the time. In other words, this time-varying voltage signal is to be distinguished from a traveling wave voltage which will be discussed in the next section. In this case,

$$\begin{aligned} \phi &= \left(\frac{2\pi}{\lambda}\right) \left(n_{x'} - \frac{1}{2} n_{x'}^3 r E_m \sin \omega_m t \right) L \\ &= \left(\frac{2\pi}{\lambda}\right) n_{x'} L - \delta \sin \omega_m t \end{aligned} \quad (35)$$

The parameter $\delta = (\pi/\lambda)n_{x'}^3 r E_m L = \pi V_m / V_\pi$, where V_π is the half-wave voltage for a given configuration, is the *phase modulation index* or *depth-of-phase modulation*. By neglecting the constant phase term ϕ_o , applying the identity $\cos(\delta \sin \omega_m t) + j \sin(\delta \sin \omega_m t) = \exp[j\delta \sin \omega_m t] = \sum_{l=-\infty}^{\infty} J_l(\delta) \exp[jl\omega_m t]$, and equating the real and imaginary parts, the output light wave becomes

$$\begin{aligned} E_o(t) &= E_i [j_0(\delta) \cos \omega t + J_1(\delta) \cos(\omega + \omega_m)t - J_1(\delta) \cos(\omega - \omega_m)t \\ &\quad + J_2(\delta) \cos(\omega + 2\omega_m)t + J_2(\delta) \cos(\omega - 2\omega_m)t + \dots] \end{aligned} \quad (36)$$

The output consists of components at frequencies ω and $(\omega + n\omega_m)$, $n = \pm 1, \pm 2, \dots$. For no modulation, $\delta = 0$ and $J_0(0) = 1$, $J_n(0) = 0$ for $n \neq 0$ and $E_o(t) = E_i \cos \omega t = E_{i_x}(t)$.¹⁶ For $\delta \approx 2.4048$, $J_0(\delta) = 0$ all the power is transferred to harmonic frequencies.⁴¹ For the case of small modulation index, $\delta \ll 1$, most of the power resides in the carrier frequency at ω and a small amount resides in the first two sidebands at frequencies $\omega \pm \omega_m$. This condition makes phase modulators useful in laser mode-locking. Increasing δ results in the presence of more sidebands.

Polarization Modulator (Dynamic Retardation) Polarization modulation involves the coherent addition of two orthogonal waves, resulting in a change of the input polarization state at the output. As with a phase modulator, the basic components for an electro-optic polarization modulator (or dynamic retardation plate or polarization state converter) is an electro-optic crystal and an input polarizer. The crystal and applied voltage V (dc assumed) are configured to produce dynamically the fast and slow axes in the crystal cross section. In this case, however, the polarizer is positioned such that the input light wave is decomposed equally into the two orthogonal linear eigenpolarizations along these axes as shown in Fig. 9. If the light is polarized along the x axis and propagates along the z principal axis, for example, the propagating fields are

$$\begin{aligned} E_{x'} &= E_o \cos[\omega t - (2\pi/\lambda)n_x z] \\ E_{y'} &= E_o \cos[\omega t - (2\pi/\lambda)n_y z] \end{aligned} \quad (37)$$

where the fast and slow axes are x' and y' . The corresponding refractive indices are

$$\begin{aligned} n_{x'} &\approx n_x - \frac{1}{2}r_x n_x^3 E = n_x - \Delta n_x \\ n_{y'} &\approx n_y - \frac{1}{2}r_y n_y^3 E = n_y - \Delta n_y \end{aligned} \quad (38)$$

where n_x and n_y are the indices in the absence of an applied field and r_x, r_y are the appropriate electro-optic coefficients for the material being used and the orientation of the applied voltage. As the two polarizations propagate at different speeds through the crystal, a phase difference (relative phase) or *retardation* Γ evolves between them as a function of length:

$$\begin{aligned} \Gamma &= \frac{2\pi}{\lambda}(n_{x'} - n_{y'})L \\ &= \frac{2\pi}{\lambda}(n_x - n_y)L - \frac{\pi}{\lambda}(r_x n_x^3 - r_y n_y^3)EL = \Gamma_o + \Gamma_i \end{aligned} \quad (39)$$

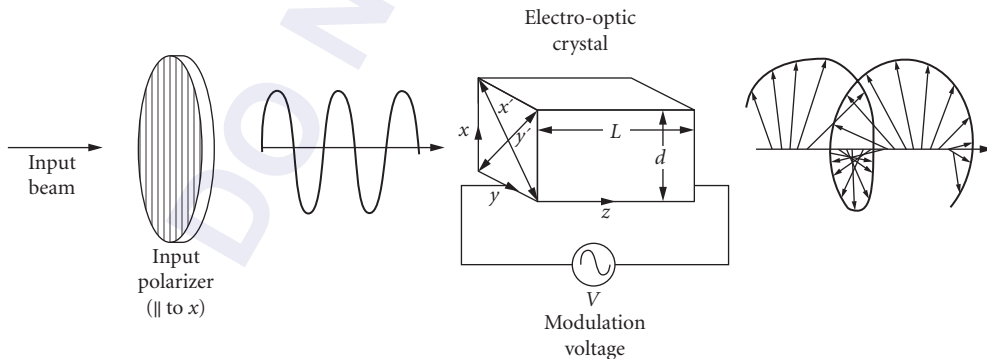


FIGURE 9 A longitudinal polarization modulator is shown with the input polarizer oriented along the x principal axis at 45° with respect to the perturbed x' and y' axes.

where Γ_ϕ is the natural phase retardation in the absence of an applied voltage and Γ_i is the induced retardation linearly related to V .

For a longitudinal modulator the applied electric field is $E = V/L$, and the induced retardation is $\Gamma_i = (\pi/\lambda)(r_x n_y^3 - r_x n_x^3)V$, which is independent of L and linearly related to V .

For a transverse modulator $E = V/d$, and the induced retardation is $\Gamma_i = (\pi/\lambda)(r_y n_y^3 - r_x n_x^3)V(L/d)$, which is dependent on the aspect ratio L/d and V .

The optical fields at the output can be expressed in terms of Γ :

$$\begin{aligned} E_{x'} &= \cos \omega t \\ E_{y'} &= \cos(\omega t - \Gamma) \end{aligned} \quad (40)$$

Therefore, the desired output polarization is obtained by applying the appropriate voltage magnitude. Figure 10 illustrates the evolution of the polarization state as a function of propagation distance z . In terms of an active device, Fig. 10 also can be interpreted as a change in polarization state as a function of applied voltage for fixed length. The eigenpolarizations $E_{x'}$ and $E_{y'}$ are in phase at $z = 0$. They have the same frequency but different wavelengths. Light from one polarization gradually couples into the other. In the absence of natural birefringence, $n_x - n_y = 0$, the voltage that would produce a retardation of $\Gamma = \Gamma_i = \pi$, such that a vertical polarization input becomes a horizontal polarization

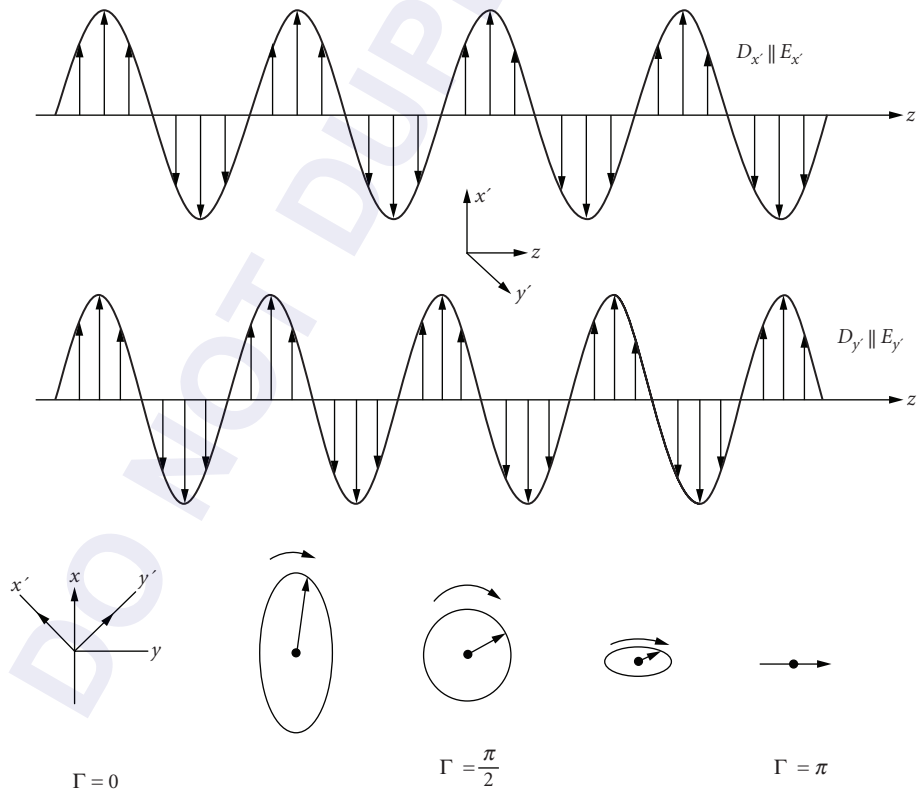


FIGURE 10 The polarization state of an input vertical linear polarization is shown as a function of crystal length L or applied voltage V . The retardation $\Gamma = \pi$ for a given length L_π in a passive $\lambda/2$ wave plate or applied voltage V_π in an electro-optic polarization modulator.¹⁶

output, is the *half-wave voltage* V_π . For a longitudinal modulator $V_\pi = \lambda/(r_x n_x^3 - r_y n_y^3)$, which is independent of L . For a transverse modulator $V_\pi = \lambda/(r_x n_x^3 - r_y n_y^3)(d/L)$, which is dependent on the aspect ratio L/d . The total retardation in terms of V_π (calculated assuming no birefringence) is

$$\Gamma = \Gamma_o + \pi \left(\frac{v}{v_\pi} \right) \quad (41)$$

To cancel the effect of natural birefringence, the phase retardation Γ_o can be made a multiple of 2π by slightly polishing the crystal to adjust the length or by applying a bias voltage. If birefringence is present, an effective V_π can be calculated that would give a total retardation of $\Gamma = \pi$.

To achieve polarization modulation, a birefringence must exist in the crystal cross section. If the cross section is a characteristic plane, then the input polarization propagates through the crystal unchanged when $V = 0$. If an applied voltage causes the axes to rotate 45° in this cross section with respect to the input polarization, as in Fig. 9, then the input will decompose into two equal components and change polarization state at the output. If the cross section has natural birefringence, then the input polarization state will change with $V = 0$ as well as with an applied voltage.

Amplitude Modulator The intensity (optical energy) of a light wave can be modulated in several ways. Some possibilities include using (1) a dynamic retarder configuration with a crossed polarizer at the output, (2) a dynamic retarder configuration with a parallel polarizer at the output, (3) a phase modulator configuration in a branch of a Mach-Zehnder interferometer, or (4) a dynamic retarder with push-pull electrodes. The intensity modulator parameter of interest is the *transmission* $T = I_o/I_p$, the ratio of output to input intensity.

An intensity modulator constructed using a dynamic retarder with crossed polarizers is shown in Fig. 11. The transmission for this modulator is

$$T(V) = \sin^2 \left(\frac{\Gamma}{2} \right) = \sin^2 \left(\frac{\Gamma_o}{2} + \frac{\pi V}{2V_\pi} \right) \quad (42)$$

For linear modulation, where the output is a replica of the modulating voltage signal, a fixed bias of $\Gamma_o = \pi/2$ must be introduced either by placing an additional phase retarder, a $\lambda/4$ wave plate (Fig. 11), at the output of the electro-optic crystal or by applying an additional dc voltage of $V_\pi/2$. This bias produces a transmission of $T = 0.5$ in the absence of a modulating voltage. If the crystal cross section has natural birefringence, then a variable compensator (Babinet-Soleil) or a voltage less than $V_\pi/2$ must be used to tune the birefringence to give a fixed retardation of $\pi/2$.

For a sinusoidal modulation voltage $V = V_m \sin \omega_m t$, the retardation at the output of the crystal, including the bias, is

$$\Gamma = \Gamma_o + \Gamma_i = \frac{\pi}{2} + \Gamma_m \sin \omega_m t \quad (43)$$

where $\Gamma_m = \pi V_m / V_\pi$ is the *amplitude modulation index* or *depth-of-amplitude modulation* and V_π is the half-wave voltage. The transmission becomes

$$\begin{aligned} T(V) &= \sin^2 \left(\frac{\pi}{4} + \frac{\Gamma_m}{2} \sin \omega_m t \right) \\ &= \frac{1}{2} \left[1 - \cos \left(\frac{\pi}{2} + \Gamma_m \sin \omega_m t \right) \right] \end{aligned} \quad (44)$$

If the modulation voltage is small ($V_m \ll 1$), then the modulation depth is small ($\Gamma_m \ll 1$) and

$$T(V) = \frac{1}{2} [1 + \Gamma_m \sin \omega_m t] \quad (45)$$

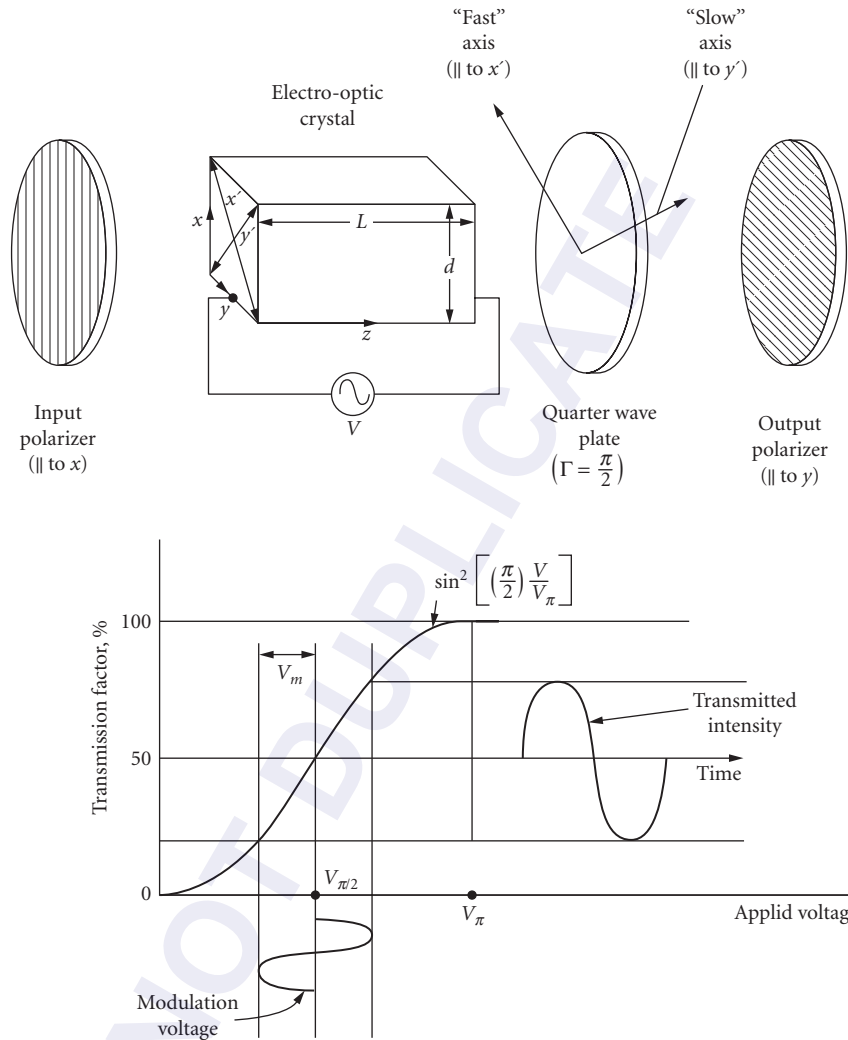


FIGURE 11 A longitudinal intensity modulator is shown using crossed polarizers with the input polarization along the x principal axis. A $\lambda/4$ wave plate is used as a bias to produce linear modulation.¹⁶

Therefore, the transmission or output intensity is linearly related to the modulating voltage. If the signal is large, then the output intensity becomes distorted and higher-order odd harmonics appear.¹⁶ The dynamic retarder with parallel polarizers has a transmission of¹⁶

$$\begin{aligned}
 t &= \cos^2\left(\frac{\Gamma}{2}\right) = \cos^2\left(\frac{\pi}{4} + \frac{\Gamma_m}{2} \sin\omega_m t\right) \\
 &= \frac{1}{2} \left[1 + \cos\left(\frac{\pi}{2} + \Gamma_m \sin\omega_m t\right) \right]
 \end{aligned}
 \tag{46}$$

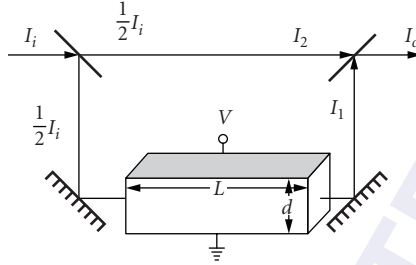


FIGURE 12 An intensity modulator is shown implementing a Mach-Zehnder interferometer configuration with a phase modulator in one branch.⁴²

For small modulation, $T(V) = 1/2[-\Gamma_m \sin \omega_m t]$ and again, the output is a replica of the modulating voltage.

Similarly, the output of a Mach-Zehnder interferometer is given by

$$I_o = I_1 + I_2 + \frac{1}{2}[I_i \cos \Gamma_o + I_i] = I_i \cos^2 \left(\frac{\Gamma_o}{2} \right) \quad (47)$$

where Γ_o is the relative phase shift between the two branches. An intensity modulator is produced by placing a phase modulator in one branch as shown in Fig. 12. The total retardation is $\Gamma = \Gamma_o + \Gamma_p$, as before. The transmission is

$$T = \frac{I_o}{I_i} = \cos^2 \left(\frac{\Gamma}{2} \right) \quad (48)$$

The push-pull modulator is based on the Mach-Zehnder interferometer. In this case, a phase modulator is placed in each branch with opposite polarity voltages applied to the arms; the phase modulators are driven 180° out of phase. This configuration requires lower drive voltages and provides a shorter transit time for the light for a defined degree of modulation.⁴⁵

Frequency Modulator In frequency modulation a shift or deviation in the frequency by ω_d from the optical carrier instantaneous frequency ω is desired. One approach to achieve a shift in frequency is to use an intensity modulator configuration of an electro-optic crystal between left- and right-hand circular polarizers. The electrodes on the modulator must be designed to produce an applied circular electric field.⁴⁶

A crystal and its orientation are selected such that there is no birefringence in the crystal cross section at zero voltage. When a circular electric field with frequency ω_m , is applied; however, a birefringence is induced in the crystal cross section, and the induced principal axes rotate with angular velocity $\omega_m/2$ in the opposite sense with respect to the modulating field. The relative rotation between the axes and the modulating field creates a frequency shift in the optical electric field as it propagates through the crystal.

An example of such a device⁴⁶ is shown in Fig. 13. There are two sets of electrodes in the transverse configuration. The applied voltages are 90° out of phase to produce a left circular modulating field:

$$\begin{aligned} E_x &= E_m \cos \omega_m t \\ E_y &= E_m \sin \omega_m t \end{aligned} \quad (49)$$

The principal axes in the crystal cross section rotate through an angle

$$\beta_1(t) = -\frac{1}{2}(\omega_m t + \Phi) \quad (50)$$

where Φ is a fixed angle that depends on the electro-optic coefficients r_{ij} of the crystal and not on the electric field magnitude.

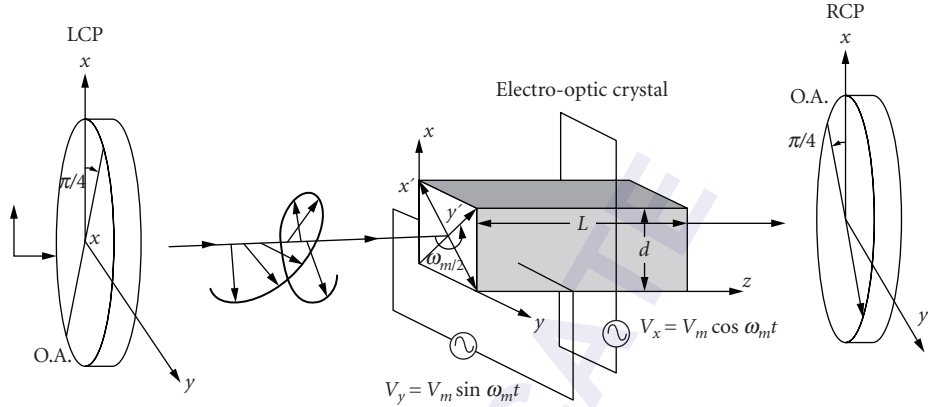


FIGURE 13 A frequency modulator using a phase modulator with two pairs of transverse electrodes set 90° out of phase to produce a circular applied electric field. The phase modulator is placed between left and right circular polarizers to create a frequency deviation in the output optical field.

The input optical wave (after the left circular polarizer) has field components

$$\begin{aligned} E_{i_x} &= E_i \cos \omega t \\ E_{i_y} &= E_i \sin \omega t \end{aligned} \quad (51)$$

The induced retardation $\Gamma_i = \Gamma = (2\pi/\lambda)\Delta nL$ is independent of time, since Δn is constant although the principal axes are rotating at constant angular velocity. The optical field components along the induced principal axes at the output of the crystal are

$$\begin{aligned} E_{o_{x'}} &= E_i \cos\left(\omega t - \beta_1 + \frac{\Gamma}{2}\right) \\ E_{o_{y'}} &= E_i \sin\left(\omega t - \beta_1 + \frac{\Gamma}{2}\right) \end{aligned} \quad (52)$$

In terms of the original stationary x, y axes, the optical field components at the output are

$$\begin{aligned} E_{o_x} &= E_i \cos(\Gamma/2) \cos \omega t - E_i \sin(\Gamma/2) \sin[(\omega + \omega_m)t + \Phi] \\ E_{o_y} &= E_i \cos(\Gamma/2) \sin \omega t - E_i \sin(\Gamma/2) \cos[(\omega + \omega_m)t + \Phi] \end{aligned} \quad (53)$$

The first terms in E_{o_x} and E_{o_y} represent left circular polarization at the original optical frequency ω with constant amplitude $E_i \cos(\Gamma/2)$ and phase independent of the rotating principal axes (that is, of β_1). The second terms represent right circular polarization at frequency $(\omega + \omega_m)$ with constant amplitude $E_i \sin(\Gamma/2)$ and phase proportional to $2\beta_1 = \omega_m t$. Therefore, the frequency shift ω_d is the modulation frequency ω_m . If the retardation $\Gamma = \pi$, the frequency of the light has complete deviation to $(\omega + \omega_m)$. If Γ is very small, the component optical fields at frequency $(\omega + \omega_m)$ are linearly related to Γ and therefore, to the applied voltage.

A shift in frequency to $(\omega - \omega_m)$ is obtained if the optical and the applied modulating electric fields rotate in the opposite sense.

Scanners The position of an optical beam can be changed dynamically by using an electro-optic deflecting device or scanner. Analog scanners are based on refraction phenomena: (1) refraction at a dielectric interface (prism) and (2) refraction by an index gradient that exists perpendicular to the direction of light propagation. Digital scanners or switches are based on birefringence.

One of the most important parameters characterizing the performance of a scanner is its resolution, the number of independent resolvable spots it can scan, which is limited by the diffraction occurring at the aperture of the device. The Rayleigh criterion states that two spots are just resolved when the angular displacement of the beam $\Delta\varphi$ is equal to the half-angle divergence θ due to diffraction.¹⁵ Therefore, the total number of resolvable spots is given by the ratio of the total deflection angle φ to the half-angle divergence θ

$$N = \frac{\varphi}{\theta} \quad (54)$$

The half-angle divergence is $\theta = \overline{\sigma}/w$, where w is the beamwidth, $\overline{\sigma} = 1$ for a rectangular beam of uniform intensity, $\overline{\sigma} = 1.22$ for a circular beam of uniform intensity, and $\overline{\sigma} = 1.27$ for a beam of Gaussian intensity distribution.⁴⁷

An analog scanner can be constructed by a prism of electro-optic material with electrodes on the crystal faces as shown in Fig. 14. The resolution is maximum in this isosceles-shaped prism when the beam is transmitted through at the minimum deviation angle and is^{47,48}

$$N = \Delta n \left(\frac{l}{\overline{\sigma}\lambda} \right) \left(\frac{w}{W} \right) \quad (55)$$

where l is the base length of the prism and w/W is the ratio of the beamwidth to input aperture of the prism. This type of device typically requires a voltage much higher than the half-wave voltage V_{π} to resolve N spots.⁴²

An analog scanner based on a gradient index of refraction⁴⁹ is shown in Fig. 15. The voltage is applied to a crystal such that the change in index is linear with distance perpendicular to the direction of light propagation; that is, $n(x) = n_0 + (2\Delta n/W)x$, where W is the width of the crystal. For linear gradient and small refraction angles, the wavefront remains planar. The (small) deflection angle of the ray after propagating a distance L in the crystal is approximated to be⁴⁷

$$\varphi = L \frac{dn}{dx} \quad (56)$$

and the resolution is

$$N = \frac{\varphi}{\theta} = 2\Delta n \left(\frac{L}{\overline{\sigma}\lambda} \right) \left(\frac{w}{W} \right) \quad (57)$$

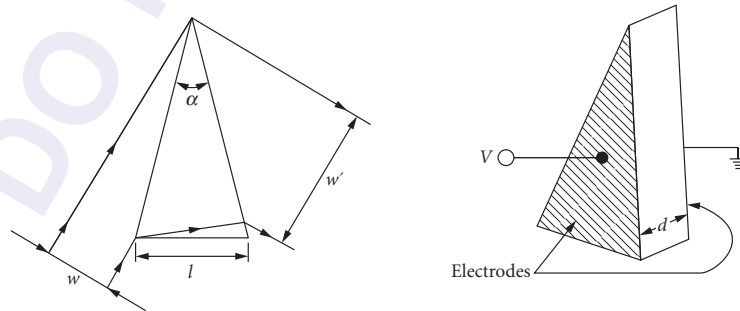


FIGURE 14 An analog scanner can be constructed using an isosceles-shaped prism with the beam transmitted at the minimum deviation angle.⁴⁸

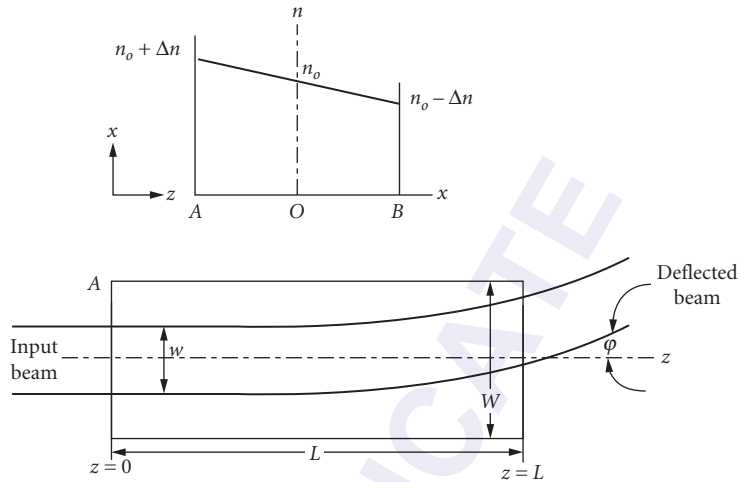


FIGURE 15 An analog scanner based on an index gradient is shown.⁴⁹

A large crystal length L is needed to obtain appreciable deflection, since Δn is very small ($\sim 10^{-4}$). Laser beams, however, are very narrow to make such a device practical.

Digital light deflection can be implemented with a number of *binary units*,^{42,50} each of which consists of a polarization modulator followed by a birefringent crystal (or discriminator) as shown in Fig. 16. The polarizer is oriented to give a horizontal polarization input to the birefringent crystal, such as calcite, and passes through undeflected. With an applied voltage, the principal axes rotate 45° and the input is then decomposed into orthogonal components of equal amplitude. The voltage

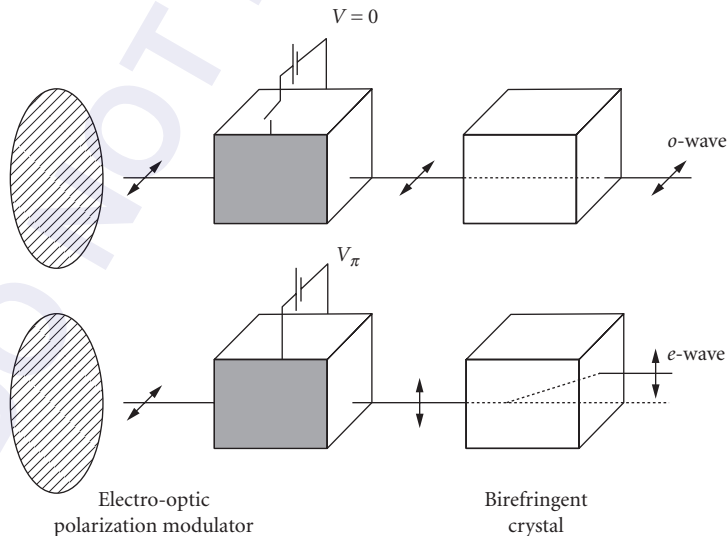


FIGURE 16 The component of a digital scanner is a binary unit. It consists of a polarization modulator followed by a birefringent crystal which serves as a discriminator.⁴²

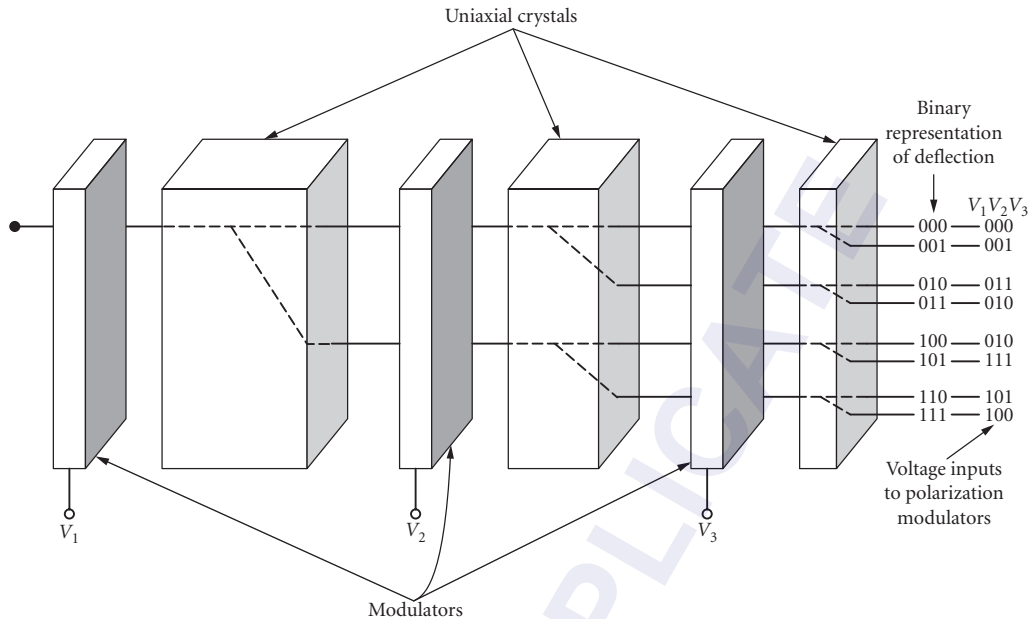


FIGURE 17 A three-stage digital scanner produces 2^3 possible deflected positions at the output.⁵⁰

magnitude is set to produce a retardation of π , thereby rotating the polarization by 90° . The vertical polarization then enters the birefringent crystal which deflects it.^{47,51}

The number of binary units n would produce a deflector of 2^n possible deflected positions at the output. An example of a three-stage deflector is shown in Fig. 17.⁵⁰⁻⁵² The on-off states of the three voltages determine at what position the output beam would be deflected. For example, if all three voltages are off, the input polarization remains horizontal through the system and is undeflected at the output. However, if all three voltages are on, the horizontal input is rotated 90° , becoming vertical after the first polarization modulator, and is deflected. The polarization is rotated 90° after the second modulator, becoming horizontal, and therefore, propagates straight through the birefringent crystal. Finally, the polarization is rotated 90° after the third modulator, becoming vertical, and is deflected. The corresponding output represents a binary five, that is, 101.

Traveling Wave Modulators The bandwidth of bulk modulators is limited by the transit time of the light through the crystal compared to the modulation frequency. The transit time of the light is the time for it to pass through the crystal

$$\tau = \frac{nL}{c} \quad (58)$$

where n is the index seen by the light. This parameter has no relevance for modulation frequencies $\omega_m \ll 2\pi/\tau$; the modulation field appears uniform in the crystal at very low frequencies. A rule of thumb for defining the limiting frequency such that τ can be neglected is that the length of the crystal be less than one-tenth the wavelength of the modulating field⁵³ or $L \ll 2\pi c/\omega_m \sqrt{\epsilon}$.⁴¹ The electro-optic crystal is modeled as a lumped capacitor at low frequencies.

As the modulating frequency becomes larger, the transit time must be taken into account in evaluating modulator performance. The modulation electric field has the form $E = E_m \sin \omega_m t$, and the optical phase can no longer follow the time-varying index of refraction adiabatically. The result is a

reduction in the modulation index parameters, δ for phase modulation and Γ_m for amplitude modulation, by a factor^{16,41,44}

$$\rho = \frac{\sin\left(\frac{1}{2}\omega_m \tau\right)}{\frac{1}{2}\omega_m \tau} \quad (59)$$

Therefore, the phase modulation index at high frequencies becomes

$$\delta_{\text{RF}} = \delta \cdot \rho = \delta \cdot \left[\frac{\sin\left(\frac{1}{2}\omega_m \tau\right)}{\frac{1}{2}\omega_m \tau} \right] \quad (60)$$

and the amplitude modulation index at high frequencies becomes

$$\Gamma_{m\text{RF}} = \Gamma_m \cdot \rho = \Gamma_m \cdot \left[\frac{\sin\left(\frac{1}{2}\omega_m \tau\right)}{\frac{1}{2}\omega_m \tau} \right] \quad (61)$$

If $\tau = 2\pi/\omega_m$ such that the transit time of the light is equal to the time period of the modulation signal, then there is no retardation; the retardation produced in the first half of the crystal is exactly canceled by the retardation produced in the second half.⁴¹ The maximum modulation frequency for a given crystal length L is determined by the allowable ρ parameter set by the designer.

The limitation of the transit time on the bandwidth of the modulator can be overcome by applying the voltage as a traveling wave, propagating collinearly with the optical wave. Figure 18 illustrates a transverse traveling wave configuration. The electrode is designed to be an extension of the driving transmission line to eliminate electrode charging time effects on the bandwidth. Therefore, the transit time problem is addressed by adjusting the phase velocity of the modulation signal to be equal to the phase velocity of the optical signal.^{16,41,44}

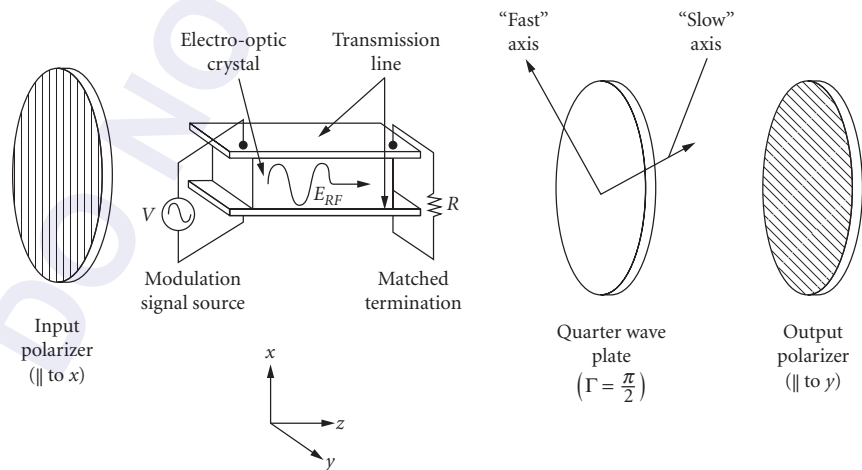


FIGURE 18 A transverse traveling wave modulator has a modulating voltage polarized in the same orientation as the input light and propagating collinearly.¹⁶

The applied modulation electric field has the form

$$E_{\text{RF}}(t, z) = E_m \sin(\omega_m t - k_m z) \quad (62)$$

for propagation along the z axis of a crystal. The direction of field vibration is the direction of the applied field, not the direction it is traveling, and is along x in Fig. 18. The parameter k_m is the wave-vector magnitude of the modulation field and is

$$k_m = \frac{\omega_m}{v_m} = \frac{\omega_m n_m}{c} \quad (63)$$

where v_m and $n_m = \sqrt{\epsilon}$ are the phase velocity and index of refraction of the modulating signal. A mismatch in the phase velocities of the modulating signal and optical wave will produce a reduction in the modulation index δ or Γ_m , by a factor

$$\rho_{\text{tw}} = \frac{\sin\left[\frac{\omega_m(n - n_m)L}{2c}\right]}{\frac{\omega_m(n - n_m)L}{2c}} = \frac{\sin(\Delta L)}{\Delta L} \quad (64)$$

In the case of amplitude modulation, Eq. (64) holds only if there is no natural birefringence in the cross section of the crystal. The eigenpolarization magnitudes are functions of time and space. Therefore, they satisfy the coupled wave equations, which are more complicated when birefringence is present.¹⁶

With no phase velocity mismatch, the phase modulation index is $\delta_{\text{RF}} = \delta$ and is linearly proportional to the crystal length L . Likewise, $\Gamma_{\text{RF}} = \Gamma_m$ for amplitude modulation. With a mismatch, the maximum possible phase modulation index is $\delta_{\text{RF,max}} = \delta/(\Delta L)$ likewise for the amplitude modulation index $\Gamma_{\text{RF,max}}$. The modulation index becomes a sinusoidal function of L . This maximum index can be achieved for crystal lengths $L = \pi/2\Delta, 3\pi/2\Delta$, and so on. The ratio n/n_m is approximately 1/2 for LiNbO₃, producing a walk-off between the optical and modulation waves.⁴⁴ Therefore, for a given length L the modulation frequency is greatly affected by the velocity mismatch.

Waveguide or Integrated-Optic Modulators Electro-optic modulators are often fabricated as integrated optical devices. Because of the small dimensions of the waveguides, these modulators require lower drive voltages and operate at higher frequencies than bulk modulators. However, the disadvantages of integrated-optic modulators are that they have lower optical power handling capabilities; they present coupling difficulties when interfaced with free-space beams; and they are designed to operate at a specific wavelength or narrow spectral window rather than working with broad optical spectra. On the other hand, because of their size, they are compatible with single mode optical fibers, and as a result, they have globally transformed communications with optical data rates up to 40 Gbps. Integrated-optic modulators can also find usage as sensors; such as gyroscopes, voltage sensors, or other phase-sensitive applications.

Most integrated optic waveguide modulators are formed on crystalline wafer surfaces, such as LiNbO₃, using standard photolithography techniques. The waveguides are formed by changing the refractive index in the guiding region using methods such as annealed proton exchange (APE) or titanium in-diffusion, and then depositing electrodes on the wafer surface. Integrated optic modulators have also been made by applying organic electro-optic material over glass waveguides, such that the application of the RF signal voltage modulates the waveguide cladding, thereby changing the effective propagation constant of the guided mode.⁵⁴

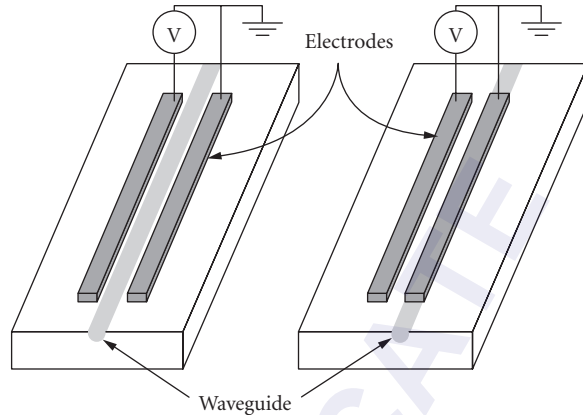


FIGURE 19 Schematic cross sections of two typical configurations of imbedded optical waveguides.

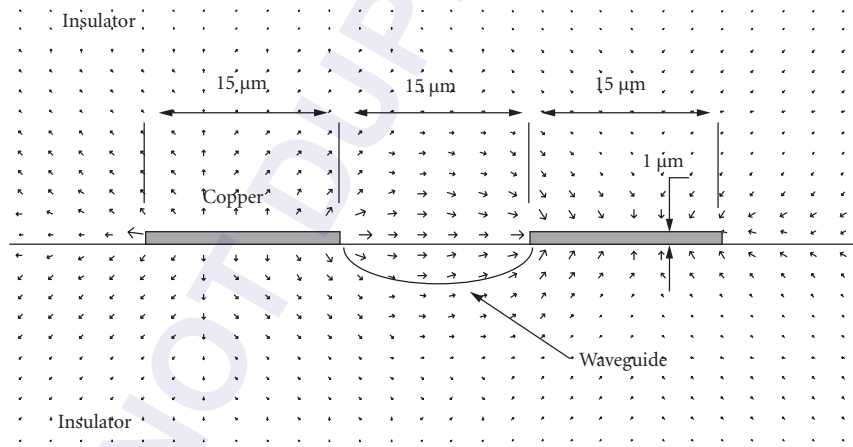


FIGURE 20 Electric field lines from electrodes across the embedded electro-optic waveguide. The electric field distribution can not be assumed to be uniform, as in bulk modulators.

The basic operation principles of waveguide modulators are the same as for bulk modulators; except that the spatial electrical and optical fields may not overlap completely. Figure 19 illustrates schematic cross sections of two typical configurations of embedded optical waveguides, and Fig. 20 shows that the field lines can not be assumed uniform as in the case of bulk modulators. Thus, the equation describing the electro-optically induced change in refractive index for a given configuration must include an overlap correction factor, Γ_{wg} .⁵⁵ For example, with a waveguide modulator, Eq. (34) would become

$$\Delta n'_x = \Delta\beta/k = -n_o^3 r_{13} \Gamma_{\text{wg}} E_z / 2 \quad (65)$$

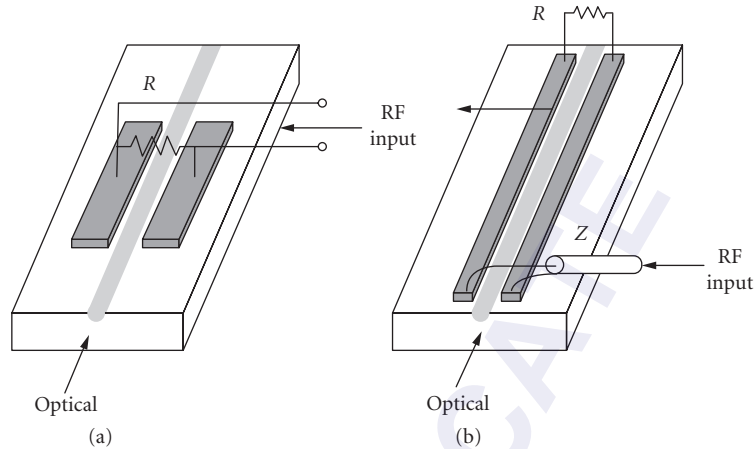


FIGURE 21 Electrode configurations on waveguides create either (a) lumped circuit element modulators or (b) traveling wave modulators.

where β is the propagation constant of the guided optical wave for a particular waveguide geometry. The overlap correction factor must be derived numerically. The total phase shift resulting from Δn over the propagation length L is again the sum of the natural phase shift and an electrically induced phase shift, but now includes a correction factor according to

$$\phi = \beta L + \Delta\beta L = \frac{2\pi}{\lambda}(n'_x + \Delta n'_x)L = \frac{2\pi}{\lambda}(n'_x - n_o^3 r_{13} \Gamma_{wg} E_z / 2)L \quad (66)$$

Waveguide modulators can be either lumped circuit element modulators or traveling wave modulators, based on their electrode configurations, which are illustrated by the phase modulators shown in Fig. 21. Similar to low frequency bulk modulators, lumped circuit element modulators have a modulation bandwidth that is limited by the parallel load resistance and electrode capacitance, which is proportional to the electrode length and gap. Decreasing the electrode length or increasing the electrode gap will decrease capacitance, and hence the RC time constant, but at the expense of requiring a higher drive voltage for a given voltage-L product. Traveling wave modulators, on the other hand, are bandwidth limited by the velocity mismatch between the propagating optical wave and the traveling RF or microwave signal.

Figure 21 illustrates phase modulators; however, amplitude (intensity) modulators are prevalent in communications applications, where they have found widespread use in fiber optical telecommunications and optical transmission of cable TV over fiber. Figure 22 shows a variety of electrode and waveguide configurations that implement intensity modulators. The Mach-Zehnder interferometer in Fig. 22a consists of traveling wave electrodes that modulate one or both arms and incur a phase difference between them. The phase-modulated optical signal recombines at the output Y-junction either constructively or destructively, so that intensity modulation results. There is an inherent 3 dB loss at the recombining Y-junction, due to an antisymmetric radiation mode that can not remain guided, which is a characteristic that does not exist for the bulk interferometric modulator. The directional coupler in Fig. 22b is designed to bring the two arms close to each other such that coupling occurs between the respective guided modes. Application of modulation controls the amount of cross coupling between the two guides and intensity modulation again results. In Fig. 22c, a device is shown that consists of a silicon substrate, a silica undercladding, and a doped-silica waveguide covered with an electro-optic polymer that has a refractive index slightly less than that of the core. When a voltage is applied to the electrodes, an electric field is produced through the electro-optic polymer and the passive optical waveguide, thus affecting the phase of the propagating light.

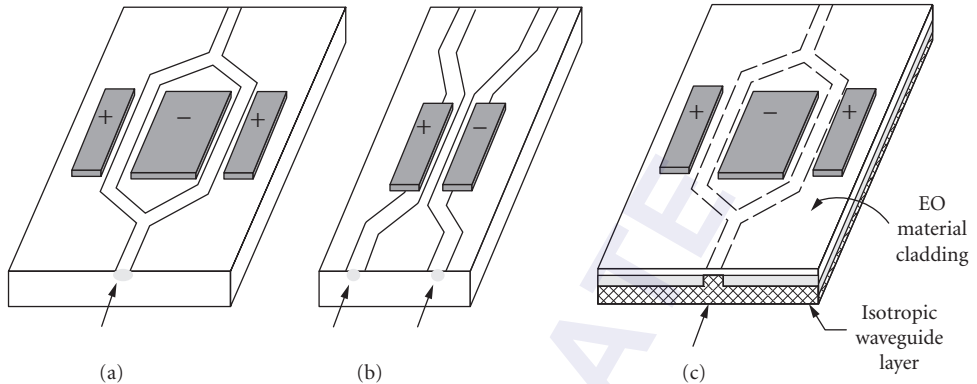


FIGURE 22 (a) Mach-Zehnder interferometer, (b) directional coupler, and (c) electro-optic cladding modulator.

Materials and Design Considerations

Two fundamental considerations to design an electro-optic modulator are signal voltage (and V_{π}) and frequency response. The goal is to achieve a desired depth of modulation with the smallest amount of drive power and at the same time obtain the largest possible bandwidth. In general, when designing bulk modulators, the field of the modulation signal and optical wave are assumed to have 100 percent overlap, whereas for waveguide modulators the overlap is not complete or uniform and a correction factor is needed. Bulk modulators are also capable of handling higher optical powers than waveguides, but at the expense of higher drive voltages and generally less modulation bandwidth. These are the trades that must be considered in choosing an electro-optic modulator.

This section discusses the various materials used to make electro-optic modulators and how the properties of electro-optic materials can affect the device performance. The various criteria and parameters that describe device performance are summarized.

Choice of Material An electro-optic material must be selected that would perform ideally at a given wavelength or a broad spectral range, and in a given environment. The material should have a large electro-optic coefficient, good homogeneity, low absorption, low dispersion, high resistivity; and thermal, mechanical, and photochemical stability. The fabricated modulator should have good ohmic contacts, and the fabrication techniques should be economically and technically feasible.

From a practical point of view, there are several limits on the performance of electro-optic modulators that must be considered. For example, some materials with natural birefringence, such as KDP and ferroelectric LiTaO_3 and LiNbO_3 , are sensitive to temperature variations and would induce a temperature-dependent phase shift in the light, called thermal drift. Some of these materials also have acoustic resonances due to the piezo-electric effect that cause undesirable peaks in the frequency response.^{24,45,52,56} The optical power handling capability, which is wavelength dependent, should be known so as not to induce photorefractive damage to the material.

Electro-optic materials include inorganic crystals, organic crystals, semiconductors, macroscopically oriented amorphous organics (dye-doped polymers), and liquids. Each of these materials are suited to particular applications based on their material properties and are briefly summarized as follows:

1. *Inorganic crystals* have high optical power handling ability and superior optical quality. Lithium niobate, LiNbO_3 , continues to be the “workhorse” of electro-optic modulators, both as a bulk modulator and especially as a waveguide. For laser applications, such as Pockels cells, crystals with high optical power handling capabilities are used. KDP and BBO are typical Pockels cell crystals for visible to near infrared applications, while CdTe is used for longer wavelengths.
2. *Organic crystals* have been explored for their large electro-optic coefficients that exceed those of inorganic crystals and their fast response time. The organic crystal 4-dimethylamino-4-stilbazolium tosylate

(DAST) presently has the highest known electro-optic coefficient and is commercially grown. While the electro-optic coefficients are large, organic crystals cannot be grown into large single crystals like inorganics, and they generally have lower power handling capabilities.

3. *Semiconductors* have been investigated as electro-optic materials because of the potential for these materials to be integrated with electronic integrated circuits and because of their transparency in the infrared; however, the effect is typically weak. Materials that have been studied include GaAs, ZnSe, ZnTe, CdS, CdTe:InP, CuCl, InSb, InGaAsP, and CdHgTe.⁵⁷
4. *Ordered dye-doped polymers* are a class of amorphous organic materials whose individual molecules display extremely large electro-optic properties. The electronic properties are based on field-induced delocalized π -electrons and additional donor and acceptor groups on opposite sides of the molecules. A mixture of these molecules does not exhibit a macroscopic electro-optic effect until the molecules have generally been oriented in the same direction. This is accomplished by several means: an electric field is applied to the assembly of molecules while above the glass transition temperature, then the material and orientation is solidified; or the orientation may be obtained by self-assembly techniques such as Langmuir-Blodgett. While the electro-optic effect can be large in these materials, over time the ordering of assembly of molecules can degrade, causing the effect to weaken. Nevertheless, research in this area remains active.
5. *Liquids* can exhibit the Pockels effect under applied electric field if the molecules are dipolar with donor acceptor groups, similar to dye-doped polymers. Some liquids also have very large Kerr coefficients and are the basis of Kerr cells.

Performance Criteria

The following parameters are indicators of modulator performance. Basically, the tradeoffs occur between the aspect ratio L/b , the drive voltage, and the electrode configuration (lumped or traveling wave) to achieve a desired depth of modulation with the smallest amount of power and at the same time obtain the largest possible bandwidth.

Modulation Bandwidth The bandwidth is the difference between the two closest frequencies at which the modulation index, δ or Γ , falls to 50 percent of the maximum value.^{55,58} Therefore, the 3 dB optical bandwidth is determined by setting the modulation reduction factor ρ or ρ_{tw} to 0.5.^{44,55,58} Modulation speed depends on the electrode type, lumped or traveling wave, and the modulator capacitance per length, which depends on the RF dielectric constant ϵ and the geometry of the electrodes. For a lumped modulator the bandwidth is limited by the optical (c/Ln) or electrical ($c/L\sqrt{\epsilon}$) transit time, whichever is smaller, or the time constant of the lumped-circuit parameters ($1/RC$), where R is the resistance in the circuit and C is the capacitance of the modulator. For a traveling wave modulator the bandwidth ν_{tw} is limited by the velocity mismatch between the optical and modulation waves, $\nu_{tw} = c/Ln(1 - \sqrt{3/\eta})$.¹⁶

Power per Unit Bandwidth (Specific Energy) and Resonant Circuits To apply the modulating signal efficiently to a *lumped* electro-optic modulator, a resonant *RLC* circuit can be created by adding an inductance and resistance in parallel to the modulator, which is modeled by a capacitor C .^{16,41} The terminating resistance R ensures that more of the voltage drop occurs over the modulator rather than the internal resistance of the source R_s .

The impedance of the resonant circuit is high over a bandwidth (determined by the *RC* time constant) of $\Delta\nu = \Delta\omega_m/2\pi = 1/2\pi RC$ centered at the resonant frequency of $\omega_o = 1/\sqrt{LC}$, where the impedance of the parallel *RLC* circuit is exactly equal to R . A peak voltage V_m must be applied to achieve a desired peak retardation $\Gamma_m = \pi V_m/V_\pi$. Therefore, $V_m = \Gamma_m(V_\pi/\pi)$. The power required to achieve Γ_m is $P = V_m^2/2R = \pi C\Delta\nu$, giving

$$P/\Delta\nu = \frac{1}{\pi} \Gamma_m^2 V_\pi^2 C \quad (67)$$

The power per unit bandwidth (mW/MHz) depends on the modulator capacitance C and the peak modulation voltage V_m to give the desired depth of modulation Γ_m . The required drive power increases with modulation frequency.^{41,55} Since the capacitance is a function of the modulator active area dimensions, the required power is also a function of the modulator dimensions.⁵⁵ Furthermore, since V_π is directly proportional to the wavelength λ , a higher input power (voltage) is required at longer wavelengths for a given peak retardation.

Extinction Ratio The extinction ratio is the maximum depth of intensity modulation for a time-varying voltage at the output when the optical bias is adjusted properly.^{56,58} If for no voltage the output intensity is I_o and for maximum applied voltage the output intensity is I_m , then the extinction ratio is defined as⁵⁸

$$\eta_m = \frac{|I_m - I_o|}{I_o} \quad I_m \leq I_o$$

$$\eta_m = \frac{|I_m - I_o|}{I_m} \quad I_m \geq I_o \quad (68)$$

Material effects, such as crystal imperfections, temperature sensitivities, and birefringence, can degrade the extinction ratio,^{16,38,45} thereby affecting the signal-to-noise ratio at the detector.⁵⁸ Another definition is in terms of the transmission T , $\eta_m = T_{\max}/T_{\min}$.^{42,59} In general, $T_{\max} < 1$ due to absorption, reflection, and scattering losses, and $T_{\min} > 0$ due to beam-divergence angle, residual crystal birefringence, crystal inhomogeneity, electric field uniformity, background scattered light, and polarizer-analyzer alignment.⁴⁰ Extinction ratio also can be applied to phase modulation, since phase changes can be related to equivalent changes in intensity.⁵⁸

Maximum Frequency Deviation A similar figure of merit as exists for intensity modulators likewise exists for frequency modulators. The maximum deviation of a frequency modulator is defined as⁵⁵

$$D_{\max} = \frac{|\omega_d - \omega|}{\omega} \quad (69)$$

where ω_d is the frequency shift when the maximum voltage is applied.

Percent Modulation Percent modulation is an intensity modulation parameter. Basically, it is the transmission $T = I_o/I_i$ times 100 percent at a specific wavelength. For a device with total retardation of $\Gamma = \pi/2$ radians at no voltage, the transmission $T = 0.5$. Then 70 percent modulation is achieved with an analog signal if a voltage is applied such that $\Gamma = 2$ radians; that is, $\sin^2(1) = 0.70$.⁵⁸ Hundred percent intensity modulation is the output intensity varying between the input intensity and zero intensity.⁴¹ A peak voltage of $V = V_\pi/2$ is required to achieve this level of performance for a linear modulator.

Another definition of percent modulation is the ratio of the peak output intensity at the modulation frequency to the maximum output intensity at dc voltage.⁵⁶ Reference 45 defines percent modulation as the peak modulation voltage output per dc voltage output, assuming no temperature variations.

Degree of Modulation For a sinusoidal reference of an optical electric field, $E(t) = E_i \sin \omega t$, an amplitude-modulated signal is typically expressed as

$$E(t) = E_i(1 + m \sin \omega_m t) \sin \omega t \quad (70)$$

where m is the degree of modulation at a specific wavelength.¹⁷ For $m \ll 1$ the intensity is

$$I = \frac{1}{2} E_i^2 (1 + 2m \sin \omega_m t) \quad (71)$$

The amplitude modulation index is then $\Gamma_m = 2m$, a function of the degree of modulation. The degree of modulation is often referred to as percent modulation for intensity modulation (70 to 100 percent nominal) and as modulation index for phase modulation (1 radian nominal).³⁸

Modulation Efficiency Modulation efficiency is defined as the percentage of total power which conveys information.^{60,61} The total power of an amplitude-modulated optical carrier is a function of the modulation frequency and is proportional to $1 + m^2 \langle \sin^2 \omega_m t \rangle = 1 + (1/2)m^2$. Therefore, the modulation efficiency is⁶⁰

$$\xi = \frac{\frac{1}{2}m^2}{1 + \frac{1}{2}m^2} \quad (72)$$

The maximum efficiency is achieved when m is maximum; that is, $m = (1/2)\Gamma_m = \pi V/2V_\pi$. In terms of the input power P the modulation efficiency is⁶¹

$$\xi = \frac{\Gamma^2}{P} \quad (73)$$

where $\Gamma = \pi V/V_\pi$ for maximum efficiency.

Optical Insertion Loss For an external modulator, in particular, the optical insertion loss must be minimized when coupling light into and out of the device. For an input light intensity I_{in} the insertion loss IL is⁵⁸

$$\begin{aligned} IL &= 1 - \frac{I_m}{I_{in}} & I_m &\geq I_o \\ IL &= 1 - \frac{I_o}{I_{in}} & I_m &\leq I_o \end{aligned} \quad (74)$$

where I_o and I_m are the output intensities at no voltage and maximum voltage, respectively. If a beam has a large cross-sectional area, the modulator must have a large acceptance area to couple all or most of the light into the device. A longitudinal modulator typically has a large acceptance area. Minimizing insertion loss in a transverse modulator is more of a challenge due to its narrow width, particularly for wide beams. However, a transverse modulator requires less voltage for a desired degree of modulation.

7.6 APPLICATIONS

Section 7.5 has discussed various device geometries and how the geometries translate into phase, amplitude, frequency, or displacement modulators. In this section, several applications of these modulators are discussed and examples are given. Waveguide modulators are discussed in another chapter of this *Handbook* and will not be further discussed here. Instead, high-speed electro-optic sampling, sensors applications, and laser mode-locking are discussed as classical application of both Pockels and Kerr electro-optic modulation.

Electro-Optic Sampling

Electro-optic sampling is a technique that can be used to sample ultrashort temporal phenomena that cannot be otherwise measured by conventional means. The sensing of terahertz radiation to

perform spectroscopy or to measure other short electrical phenomena in electrical circuits are two examples. The technique has been enabled by the advent of pulsed laser systems, creating picosecond to femtosecond optical pulses of duration shorter than the electrical pulse to be measured. It consists of a “probe” pulsed optical light beam that is modulated by the temporary short electric field, both of which are coincident within an electro-optic crystal or material. Because the probe duration is less than the modulating electrical pulse (which may be a terahertz pulse or other electrical transient), the probe polarization is modified through the electro-optic effect, the degree to which is based on the strength of the electric field at that instant of time. A polarizer is placed after the electro-optic crystal, so that the change in intensity of one particular polarization can be measured. The temporal position of the probe relative to the electrical signal is slowly varied. However, at a particular instant of time, many optical pulses are perturbed, thus permitting the use of slower speed detectors to measure the variation of the optical power.

Figure 23a illustrates a typical terahertz time-domain spectroscopy (TDS) system that utilizes a pump-probe method of electro-sampling to measure terahertz pulses. A ZnTe crystal or an organic DAST crystal is used as the electric field sampling sensor. However, since these both are noncentrosymmetric crystals, they can also be used to generate the terahertz pulse by the second-order process of optical rectification. Figure 23b shows the reconstructed terahertz pulse.

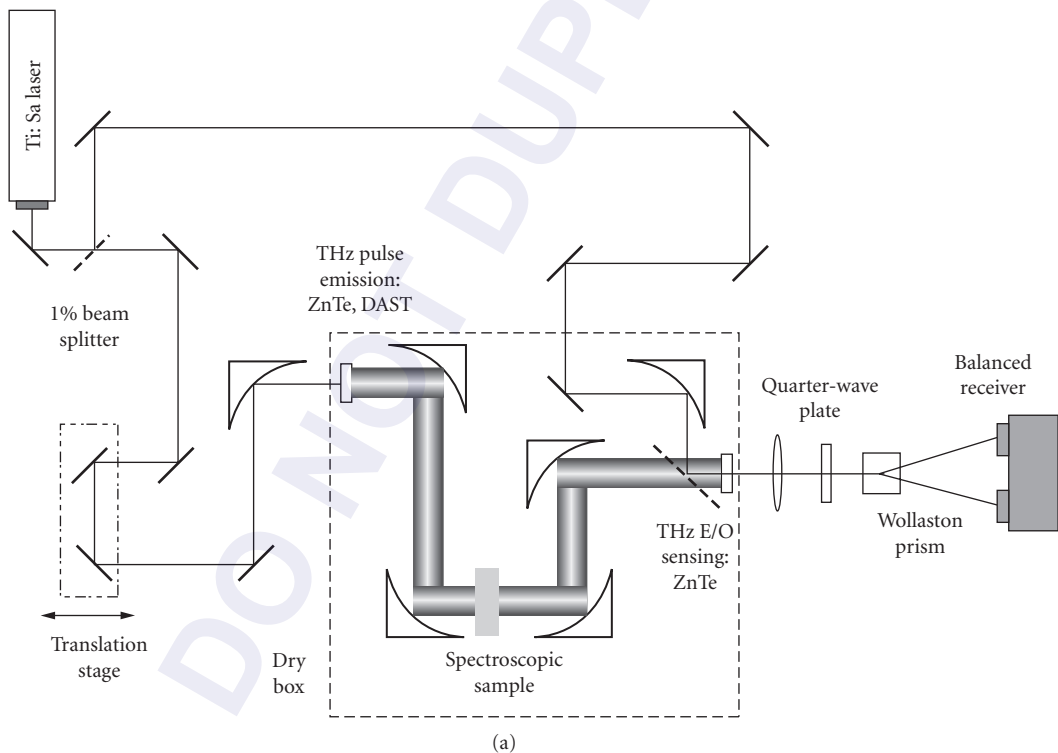


FIGURE 23 (a) A terahertz TDS system uses a Ti:Sapphire laser to emit ultrashort optical pulses, which are split into a pump beam to generate terahertz pulses, and a probe beam used to sample the terahertz pulses. An electro-optic crystal is used as a sensor. (b) The terahertz pulse reconstructed by sampling the electric field strength at different points in time.

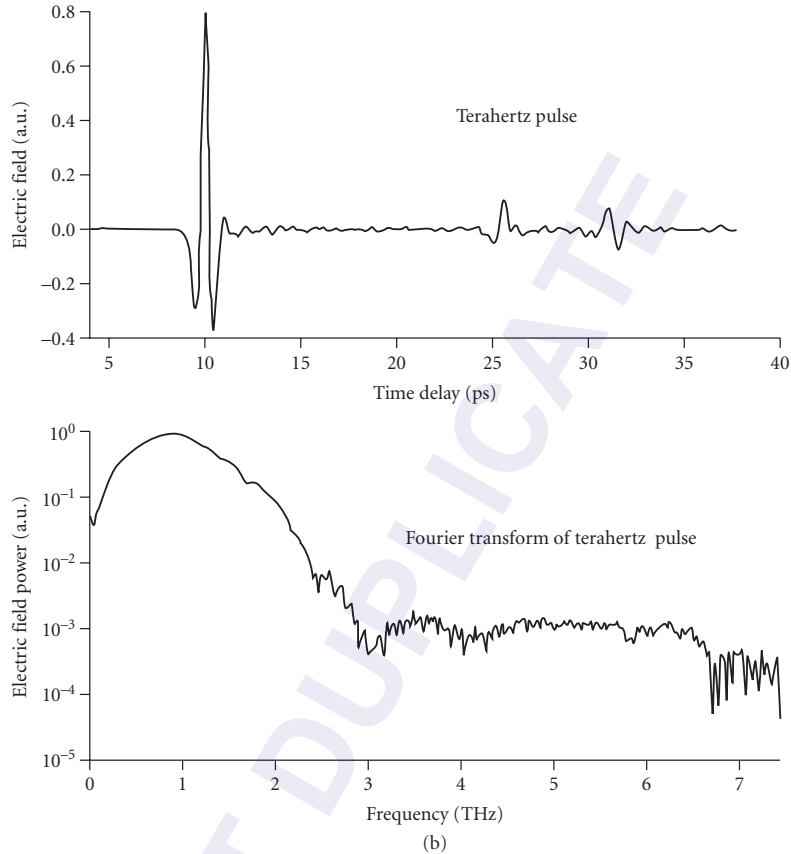


FIGURE 23 (Continued)

Sensors

A practical commercial sensor based on the electro-optic effect is a high-voltage sensor. This sensor uses the Pockels effect to sense high voltages without contact with a conductor or a voltage source. It finds application in improving reliability of electric utility systems. In contrast to other modulators discussed in this chapter, wherein a low-drive voltage was desired, high-voltage sensors accurately measure voltages as high as 765 kV. The electro-optic material for the sensor is being used as a bulk modulator, and the excitation voltage is of low frequency.

In contrast to the low-frequency high-voltage sensor just discussed, an integrated optic Mach Zehnder interferometer (MZI) can be used as a high-power microwave sensor or RF sensor. This device has been shown to operate over a wide frequency range of 200 MHz to 12 GHz.⁶³

Laser Mode-Locking

Mode-locking is a resonant phenomenon that results in ultrashort laser pulses. Mode-locking occurs by a relatively weak modulation of the radiation passing through the resonator, which results in initiating a pulse that becomes shorter with each pass through the resonator. The pulse compression technique uses the Kerr nonlinearity of an optical medium. The pulses propagating

through the medium experience nonlinear phase shifts that lead to spectral broadening, resulting in a spread of frequencies, or chirp. The different frequency components are superimposed by propagation in a dispersive medium, or by reflection from a dispersive element. The effects of a fast saturable absorber to produce mode-locking can be simulated by Kerr focusing. The Kerr effect within the resonator causes the high-intensity part of the optical beam to be focused, whereas the low-intensity components remain unfocused. If an aperture is placed in the beam path, the low-intensity parts are attenuated, and the pulse is shortened. This effect is now called Kerr-lens mode-locking (KLM).²⁸ The KLM technique is commonly used for mode-locking of Ti:Sapphire lasers.

7.7 APPENDIX: EULER ANGLES

Euler angles represent a set of three independent parameters which specify the orientation of a rigid body, in this case the index ellipsoid. Orthogonal transformations are used to convert from one set of axes (x, y, z) to another set (x', y', z') . The two sets of axes are related to each other by nine direction cosines. An orthogonal transformation matrix consisting of direction cosines and having a determinant of +1 corresponds to defining the orientation of a rigid body.

The transformation matrix is developed by a specific sequence of three plane rotations (not in principal planes) in a defined order. There are 12 possible conventions for defining a set of Euler angles in a right-handed coordinate system.³⁶ One convention that has been proposed as a standard is the y convention.⁶⁴ The transformation evolves by an initial counterclockwise rotation ζ about the z axis, followed by a counterclockwise rotation η about the intermediate y' axis, and finally a counterclockwise rotation ω about z'' . The resulting transformation matrix is

$$\bar{a} = \begin{pmatrix} \cos \omega & \sin \omega & 0 \\ -\sin \omega & \cos \omega & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \cos \eta & 0 & -\sin \eta \\ 0 & 1 & 0 \\ \sin \eta & 0 & \cos \eta \end{pmatrix} \begin{pmatrix} \cos \zeta & \sin \zeta & 0 \\ -\sin \zeta & \cos \zeta & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (\text{A.1})$$

To find the Euler angles (ζ, η, ω) , Eq. (A.1) is rearranged as³⁷

$$\begin{pmatrix} \cos \omega & -\sin \omega & 0 \\ \sin \omega & \cos \omega & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} \cos \eta & 0 & -\sin \eta \\ 0 & 1 & 0 \\ \sin \eta & 0 & \cos \eta \end{pmatrix} \begin{pmatrix} \cos \zeta & \sin \zeta & 0 \\ -\sin \zeta & \cos \zeta & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (\text{A.2})$$

where \bar{a} is the cumulative transformation matrix of the normalized eigenvectors. Multiplying the matrices, the Euler angles are related to the elements of \bar{a} :

$$a_{11} \cos \omega - a_{21} \sin \omega = \cos \eta \cos \zeta \quad (\text{A.3a})$$

$$a_{12} \cos \omega - a_{22} \sin \omega = \cos \eta \sin \zeta \quad (\text{A.3b})$$

$$a_{13} \cos \omega - a_{23} \sin \omega = -\sin \eta \quad (\text{A.3c})$$

$$a_{11} \sin \omega + a_{21} \cos \omega = -\sin \zeta \quad (\text{A.3d})$$

$$a_{12} \sin \omega + a_{22} \cos \omega = \cos \zeta \quad (\text{A.3e})$$

$$a_{13} \sin \omega + a_{23} \cos \omega = 0 \quad (\text{A.3f})$$

$$a_{31} = \sin \eta \cos \zeta \quad (\text{A.3g})$$

$$a_{32} = \sin \eta \sin \zeta \quad (\text{A.3h})$$

$$a_{33} = \cos \eta \quad (\text{A.3i})$$

From Eq. (A.3f), the angle ω is $\omega = \tan^{-1}(-a_{23}/a_{13})$. From Eq. (A.3c), the angle η is $\eta = \sin^{-1}(a_{23} \sin \omega - a_{13} \cos \omega)$. From Eq. (A.3d), the angle ζ is $\zeta = \sin^{-1}(-a_{11} \sin \omega - a_{21} \cos \omega)$.³⁷

7.8 REFERENCES

1. W. J. Tomlinson and C. A. Brackett, "Telecommunications Applications of Integrated Optics and Optoelectronics," *Proc. IEEE* **75**(11):1512–1523 (1987).
2. E. Vogues and A. Neyer, "Integrated-Optic Devices on LiNbO₃ for Optical Communication," *IEEE/OSA J. Lightwave Technol.* **LT-5**:1229–1238 (1987).
3. L. Thylen, "Integrated Optics in LiNbO₃: Recent Developments in Devices in Telecommunications," *IEEE/OSA J. Lightwave Technol.* **6**(6):847–861 (1988).
4. R. C. Alferness, "Guided-Wave Devices for Optical Communication," *IEEE J. Quantum Electron.* **QE-17**(6):946–959 (1981).
5. H. F. Taylor, "Application of Guided-Wave Optics in Signal Processing and Sensing," *Proc. IEEE* **75**(11):1524–1535 (1987).
6. See, for example, "Special Issue on Optical Computing," *Proc. IEEE* **72** (1984).
7. See, for example, "Special Feature on Integrated Optics: Evolution and Prospects," *Opt. News* **14** (1988).
8. T. K. Gaylord and E. I. Verriest, "Matrix Triangularization Using Arrays of Integrated Optical Givens Rotation Devices," *Computer* **20**:59–66 (1987).
9. C. M. Verber, R. P. Kenan, and J. R. Busch, "Design and Performance of an Integrated Optical Digital Correlator," *IEEE/OSA J. Lightwave Technol.* **LT-1**:256–261 (1983).
10. C. L. Chang and C. S. Tsai, "Electro-Optic Analog-to-Digital Conversion Using Channel Waveguide Fabry-Perot Modulator Array," *Appl. Phys. Lett.* **43**:22 (1983).
11. C. M. Verber "Integrated-Optical Approaches to Numerical Optical Processing," *Proc. IEEE* **72**:942–953 (1984).
12. X. Zheng, S. Wu, and R. Sobolewski, "Electro-Optic Sampling System with a Single-Crystal 4-N, N-Dimethylamino-4'-N'-Methyl-4-Stibazolium Tosylate Sensor," *Appl. Phys. Lett.* **82**(15):2383–2385 (2003).
13. J. Ruan, H. Edwards, C.-Y. Tan, R. Thurman-Keup, and V. Scarpine, "Design of an Electro-Optic Sampling Experiment at the AWA Facility," *Proc. IEEE PAC'07*, Albuquerque, NM, 3901–3903 (2007).
14. T. A. Maldonado and T. K. Gaylord, "Light Propagation Characteristics for Arbitrary Wavevector Directions in Biaxial Crystals by a Coordinate-Free Approach," *Appl. Opt.* **30**:2465–2480 (1991).
15. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon Press, Oxford, UK, 1980.
16. A. Yariv and P. Yeh, *Optical Waves in Crystals*, Wiley, New York, 1984.
17. L. D. Landau and E. M. Lifshitz, *Electrodynamics of Continuous Media*, Pergamon, London, 1960.
18. A. I. Borisenko and I. E. Tarapov, in *Vector and Tensor Analysis with Applications*, R. A. Silverman (ed.), Prentice-Hall, Englewood Cliffs, NJ, 1968.
19. I. P. Kaminow, *An Introduction to Electrooptic Devices*, Academic Press, New York, 1974.
20. T. C. Phemister, "Fletcher's Indicatrix and the Electromagnetic Theory of Light," *Am. Mineralogist* **39**:173–192 (1954).
21. T. A. Maldonado and T. K. Gaylord, "Electro-Optic Effect Calculations: Simplified Procedure for Arbitrary Cases," *Appl. Opt.* **27**:5051–5066 (1988).
22. J. F. Nye, *Physical Properties of Crystals, Their Representation by Tensors and Matrices*, Oxford Univ. Press, New York, 1985.
23. I. P. Kaminow, in *Handbook of Laser Science and Technology*, vol. IV, part 2, M. J. Weber (ed.), CRC Press, Boca Raton, FL, 1986, pp. 253–278.
24. I. P. Kaminow and E. H. Turner, "Electro-Optic Light Modulators," *Proc. IEEE* **54**(10):1374–1390 (1966).
25. J. M. Bennett and H. E. Bennett, "Polarization," in *Handbook of Optics*, W. G. Driscoll and W. Vaughan (eds.), McGraw-Hill, New York, 1978, chap. 10.
26. A. Yariv, *Optical Electronics*, Holt, Rinehart, and Winston, New York, 1976.
27. D. C. Cronmeyer and L. R. Spanberger, "Infrared Transmittance and Kerr Effect of Nitrobenzene," *J. Opt. Soc. Am.* **51**:1061–1066 (1961).
28. H. A. Haus, "Mode-Locking of Lasers," *IEEE J. Select. Topics Quantum Electronics* **66**(6):1173–1185 (2000).
29. J. H. Wilkinson, *The Algebraic Eigenvalue Problem*, Oxford Univ. Press, London, 1965.

30. W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterlirig, *Numerical Recipes*, Cambridge Univ. Press, New York, 1986.
31. J. H. Wilkinson and C. Reinsch, *Handbook for Automatic Computation*, Springer-Verlag, New York, 1971.
32. G. H. Golub and C. F. Van Loan, *Matrix Computations*, John Hopkins Univ. Press, Baltimore, 1983.
33. D. E. Sands, *Vectors and Tensors in Crystallography*, Addison-Wesley, Reading, MA, 1982.
34. D. R. Hartree, *Numerical Analysis*, Clarendon Press, Oxford, 1952.
35. B. N. Parlett, *The Symmetric Eigen Value Problem*, Prentice-Hall, Englewood Cliffs, NJ, 1980.
36. H. Goldstein, *Classical Mechanics*, Addison-Wesley, Reading, MA, 1981.
37. R. P. Paul, *Robot Manipulators*, The MIT Press, Cambridge, MA, 1981.
38. I. P. Kaminow and T. Li, "Modulation Techniques," in *Optical Fiber Telecommunications*, S. E. Miller and A. G. Chynoweth (eds.), Academic Press, New York, 1979, chap. 17.
39. D. F. Nelson, "The Modulation of Laser Light," *Scientific Am.* **218**(6):17–23 (1968).
40. E. Hartfield and B. J. Thompson, "Optical Modulators," in *Handbook of Optics*, W. G. Driscoll and W. Vaughan (eds.), McGraw-Hill, New York, 1978, chap. 17.
41. A. Ghatak and K. Thyagarajan, *Optical Electronics*, Cambridge Univ. Press, New York, 1989.
42. B. E. A. Saleh and M. C. Teich, *Fundamentals of Photonics*, Wiley, New York, 1991, chap. 18.
43. E. Hecht, *Optics*, 2nd ed., Addison-Wesley, Reading, MA, 1990.
44. R. C. Alferness, in *Guided-Wave Optoelectronics*, T. Tamir (ed.), Springer-Verlag, New York, 1990, chap. 4.
45. W. H. Steier, "A Push-Pull Optical Amplitude Modulator," *IEEE J. Quantum Electron.* **QE-3**(12):664–667 (1967).
46. C. F. Buhner, D. Baird, and E. M. Conwell, "Optical Frequency Shifting by Electro-Optic Effect," *Appl. Phys. Lett.* **1**(2):46–49 (1962).
47. V. J. Fowler and J. Schlafer, "A Survey of Laser Beam Deflection Techniques," *Appl. Opt.* **5**:1675–1682 (1966).
48. F. S. Chen, et al., "Light Modulation and Beam Deflection with Potassium Tantalate Niobate Crystals," *J. Appl. Phys.* **37**:388–398 (1966).
49. V. J. Fowler, C. F. Buhner, and L. R. Bloom, "Electro-Optic Light Beam Deflector," *Proc. IEEE* (correspondence) **52**:193–194 (1964).
50. T. J. Nelson, "Digital Light Deflection," *B.S.T.I.* 821–845 (1964).
51. W. Kulcke, et al., "A Fast, Digital-Indexed Light Deflector," *IBM J.* **8**:64–67 (1964).
52. M. Gottlieb, C. L. M. Ireland, and J. M. Ley, *Electro-Optic and Acousto-Optic Scanning and Deflection*, Marcel Dekker, New York, 1983.
53. D. M. Pozar, *Microwave Engineering*, Addison-Wesley, Reading, MA, 1990.
54. R. W. Ridgway, V. McGinniss, D. W. Nippa, and S. M. Risser, "Electrooptic Control of Functionally-Clad Silica Waveguides," in *Integrated Photonics Research*, A. Sawchuk (ed.), vol. 91 of OSA Trends in Optics and Photonics (Optical Society of America, 2003), paper ITuH5.
55. R. G. Hunsperger, *Integrated Optics: Theory and Technology*, 2nd ed., Springer-Verlag, Berlin, 1983.
56. R. T. Denton, et al., "Lithium Tantalate Light Modulators," *Appl. Phys.* **38**(4):1611–1617 (1967).
57. T. G. Brown, K. Creath, H. Kogelnik, M. A. Kriss, J. Schmit, M. J. Weber (eds.), *The Optics Encyclopedia: Basic Foundations and Practical Applications*, vol. 1, Electro-Optic Devices, Wiley, 2004.
58. J. M. Hammer, in *Integrated Optics*, T. Tamir (ed.), Springer-Verlag, Berlin, 1979, chap. 4, pp. 140–200.
59. R. Simon, *Optical Control of Microwave Devices*, Artech House, Boston, MA, 1990.
60. R. E. Ziemer and W. H. Tranter, *Principles of Communications*, Houghton Mifflin Co., Boston, MA, 1976.
61. I. P. Kaminow and J. Liu, "Propagation Characteristics of Partially Loaded Two-Conductor Transmission Line for Broadband Light Modulators," *Proc. IEEE* **51**(1):132–136 (1963).
62. R. Forber, W. C. Wang, and De-Yu Zang, "Dielectric EM Field Probes for HPM Test and Evaluation," 2006 Annual ITEA Technology Review, August 7–10, Cambridge, MA 1 Non-Intrusive Instrumentation Technologies, 2006.
63. W. L. Bond "The Mathematics of the Physical Properties of Crystals," *Bell Sys. Tech. J.* **23**:1–72 (1943).

This page intentionally left blank.

DO NOT DUPLICATE

LIQUID CRYSTALS

Sebastian Gauza and Shin-Tson Wu

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

ABSTRACT

This chapter introduces the basic properties of liquid crystal materials, electro-optic properties of liquid crystal cells, and operation principles of liquid crystal display (LCD) devices, including the large screen thin-film-transistor LCD TVs and transmissive LCDs for mobile displays.

8.1 GLOSSARY

C_{ii}	reduced elastic constants
d	liquid crystal thickness
d_R	reflective region cell gap
d_T	transmissive region cell gap
E	activation energy
F	Onsager reaction field
f_i	oscillator strength
f_c	crossover frequency
G	isotropic intermolecular attraction energy
H_f	heat fusion enthalpy
h	cavity field factor
I	moment of inertia
J	mean-field coupling constant
K_{ii}	elastic constants
k	Boltzmann constant
L	length of the molecule
l	aspect ratio

n	liquid crystal director
$n_{e,o}$	refractive indices
N	molecular packing density
P_4	order parameter of the fourth rank
P_s	spontaneous polarization
R	gas constant
p	molecular length-to-width ratio
S	order parameter of the second rank ($= P_2$)
T_c	clearing temperature
V_n	mole volume
V_b	bias voltage
V_{th}	threshold voltage
Z	number of electrons
α_{lt}	molecular polarizability
α_i	Leslie viscosity coefficients
β	angle between dipole and molecular axis
γ_1	rotational viscosity
δ	phase retardation
ϵ	dielectric constant
ϵ_0	vacuum permittivity
$\Delta\epsilon$	dielectric anisotropy
Δn	birefringence
μ	dipole moment
η_i	Miesowicz viscosity coefficient
θ	twist angle
Λ	friction coefficient
Φ	volume fraction of the molecules
λ	wavelength
ϕ	tilt angle
ϕ_s	pretilt angle
τ	response time
ω	frequency

8.2 INTRODUCTION TO LIQUID CRYSTALS

Liquid crystal (LC) was discovered in 1888 when an Austrian botanist named Friedrich Reinitzer observed that a material known as cholesteryl benzoate exhibiting two distinct melting points.¹ In his experiments, Reinitzer increased the temperature of a solid sample and observed that the crystal changes into a hazy liquid. As he continued to increase the temperature further, the hazy liquid turned into a clear state. Because of this early work, Reinitzer is often credited as the discoverer of a new phase of matter—the LC phase which is also known as liquid crystalline or mesophase. Most people are familiar with the fact that matter can exist in three different states: solid, liquid, and gas (vapor). However, this is a simplification, and under extreme conditions other forms of matter can exist, for example, plasma at very high temperatures or superfluid helium at very low temperatures. The difference between these states of matter is the degree of order in the material, which is directly related to the surrounding temperature and pressure. If the temperature is raised,

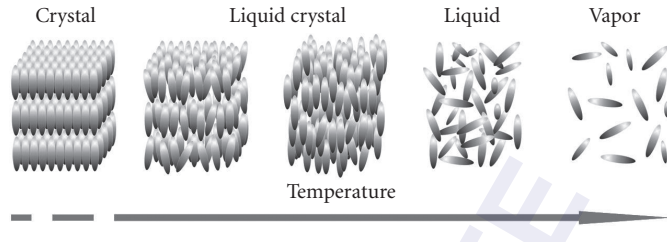


FIGURE 1 Schematic arrangement of the molecules in different phases.

more energy is transferred into the system, leading to increasingly stronger vibrations. Finally, at the transition temperatures between the solid and liquid states, the long range positional order is broken and the constituents may move in a random fashion (Fig. 1), constantly bumping into one another and abruptly changing the direction of motion. However, the thermal energy is not high enough to completely overcome the attractive forces between the constituents, so there is still some positional order at short range. Because of the remaining cohesion, the density of the liquid is constant even though the liquid takes the shape of its container, as opposed to a solid. The liquid and solid phases are called condensed phases. If we keep on raising the temperature until the next phase change, as shown in Fig. 1, the substance enters its gas (or vapor) state and its constituents are no longer bounded to each other.

The molecular order of an LC lies between those of the isotropic liquid and the crystal. The classification of LC materials is based on the degree of orientational and positional order. When the LC molecules are sandwiched between two glass substrates with their inner surfaces coated by an alignment agent, such as rubbed polyimide, the LC molecules tend to form a statistically preferred direction called the director, denoted by the vector \mathbf{n} . To specify the amount of orientational order in an LC phase, an order parameter is defined as

$$S = \langle P_2(\cos \theta) \rangle = \frac{2}{3} \langle \cos^2 \theta - 1 \rangle \quad (1)$$

where $\langle \rangle$ denotes a thermal averaging and θ is the angle between each molecule and the director \mathbf{n} . If the molecules are perfectly oriented (crystal state), that is, if $\theta = 0$ with the director, then $S = 1$. On the contrary, if the molecules are randomly oriented about \mathbf{n} , that is, an isotropic state, then $S = 0$ because there is no orientational order. So the higher the order parameter, the more ordered the liquid crystal phase is. In a typical LC system, order parameter decreases as the temperature increases. Temperature dependence of the order parameter S is shown in Fig. 2. Most common values

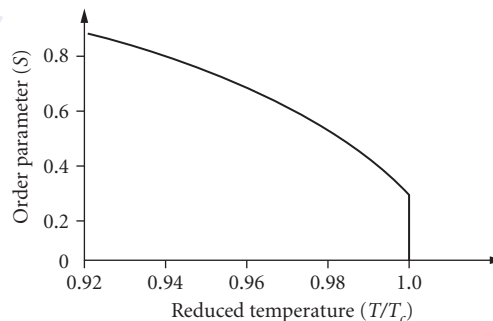


FIGURE 2 Temperature dependent order parameter of the liquid crystal phase.

of S are between 0.3 and 0.8. The order parameter is typically determined by the measured macroscopic properties like optical birefringence or diamagnetism. At certain temperatures, the LC material may gain a certain amount of positional order. When it happens, the center of mass of the LC molecules, although still forming a fluid, prefers to lie, on average, in layers as Fig. 1 depicts. This positional ordering may be described in terms of the density of the center of the mass of the molecules as

$$\rho(z) = \rho_0 \left[1 + \Psi \cos\left(\frac{2\pi z}{d}\right) \right] \quad (2)$$

where z is the coordinate parallel to the layer normal, the average density of the fluid is ρ_0 , d is the distance between layers, and ψ is the order parameter. The modulus of ψ , $|\psi|$, represents the amplitude of the oscillation of the density. When $|\psi| = 0$, there is no layering, but if $|\psi| > 0$ then some amount of sinusoidal layering exists. Different LC phases are formed as a consequence, which has been described later in this chapter.

From a basic physics point of view, LC materials are of great interest and have contributed to the modern understanding of phase transitions and order phenomena in one, two, and three dimensions. To most of the people today, LC is almost synonymous to flat panel displays (LCDs) for computers, mobile phones, and other electronic equipments. However, there is a rapid development of the other types of applications. For example, there have been developments in telecommunications, pattern recognition, real time holography, light shutters, nonmechanical beam steering, and so on.² Liquid crystals constitute a unique form of soft matter and are becoming increasingly more important in pure material science in the development of polymer materials and biomaterials.

8.3 TYPES OF LIQUID CRYSTALS

Considering the geometrical structure of the mesogenic molecules, the liquid crystals can be classified into several types. The widespread LCDs are using rod-shaped nematic molecules, as shown in Fig. 3, where one molecular axis is much longer than the other two. These molecules are called calamitic liquid crystals. Molecules have to be rigid in one part (rigid core) and flexible in another (terminal flexible hydrocarbon chain). Different physical properties may be affected by exchanging one of the terminal chains with a group having different polarities. Liquid crystals formed by disc-shaped molecules with one molecular axis much shorter than the other two are called discotic liquid crystals, as shown in Fig. 3. In this case, the rigid part of the molecule is disc-shaped, typically having multiple aromatic rings. There are some possible intermediates between rod-shaped and disc-shaped molecules known as lath-like LC molecules. Transition to the mesophase (liquid crystalline phase) may be caused by either temperature or influence of solvents. The LCs obtained from the former method are called thermotropic. If the influence of solvents induces the transition, the LCs are called lyotropic. The LC materials capable of forming thermotropic as well as lyotropic mesophases are called amphotropic liquid crystals, as Fig. 4 shows. The lyotropic liquid crystal phases are formed by dissolving amphiphilic molecules in a suitable solvent. Amphiphilic molecules consist of

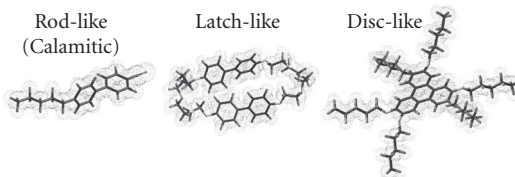


FIGURE 3 Molecular geometry of different types of LC.

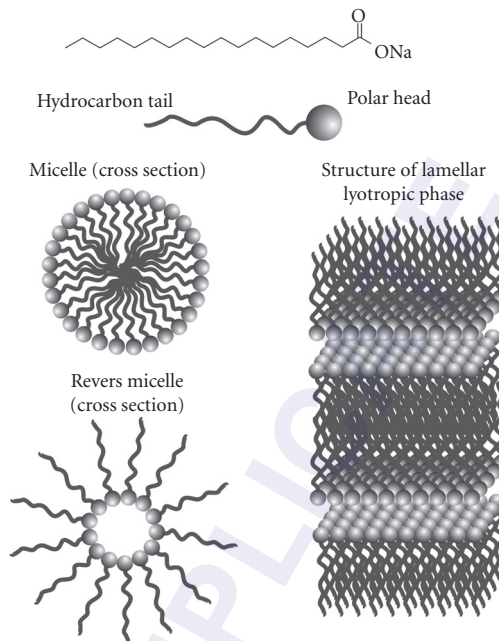
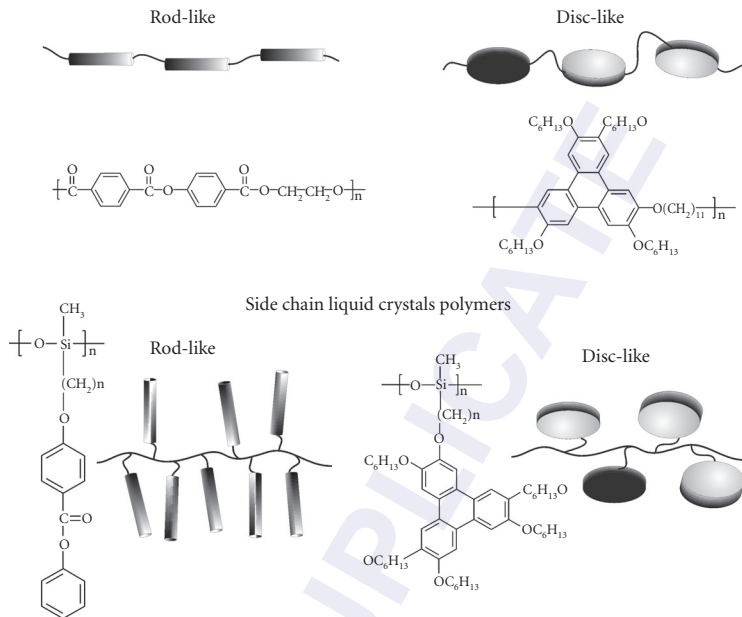
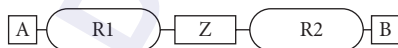


FIGURE 4 Examples of the lyotropic LCs.

hydrophobic group at one end and hydrophilic group at another. Such molecules form self-organized structures in solvents with different polarities. A very good example is soap. When dissolved in water (polar solvent), hydrophobic tails face each other and hydrophilic parts face solvent, forming a micelle. Another type, which is biologically important, is phospholipids which form bilayers. Such aggregations are the building units of biological membranes. Another important structure is liquid crystal polymers. As thermotropic mesogens, these structures consist of mesogenic subunits (rod-like or disc-like) which are linked together with flexible linkage units, as shown in Fig. 5. If rigid mesogenic units are linked directly by flexible links, the main-chain polymer LC is formed. Mesogenic subunits can also be attached to the polymer chain as a side group which are known as side-chain polymer liquid crystals. Merged structure, called combined polymer liquid crystal has build-in main-chain as well as side-chain mesogenic units. Since the discovery of the first LC substances, there have been many research activities to determine what kind of structure forms desired liquid crystalline phase. Theory suggests that mesophases can be achieved when the molecules have elongated shape and some flexibility. Typically mesogenic molecules consist of several building blocks. If we concentrate on the most common calamitic LC block systems, it will contain elements shown in Fig. 6. As shown in Fig. 6, calamitic structures consist of a rigid rod formed by two (or more) ring systems (R1 and R2). These are connected together by a single bond or a linking group known as central bridge group. Molecular constituents at the ends of the rigid core (*para* position to the central group) are called terminal groups or chains. Usually at least one of the terminal groups must be a flexible carbon chain. Tables 1 to 3 show the most popular choices for the central bridge, ring systems, and terminal substituents.

The introduction of a linking group into a mesogenic molecule often determines linearity, increases overall molecular length, and changes polarizability anisotropy, thus influences mesomorphic and physical properties in general. The presence of a central linking bridge often widens the temperature range of the mesophase by reducing the melting point. However, sometimes, this may also affect the thermal and photochemical stability by increasing π -electron conjugation of the molecules. In some cases, linking groups may induce coloration.

Main chain liquid crystals polymers

**FIGURE 5** Types of LC polymeric materials.**FIGURE 6** Building blocks of calamitic LCs. A-R1-Z-R2-B.**TABLE 1** Typical Central Bridge Groups

Central Bridge Group	Central Bridge Group
$[-\text{CH}_2-]_n$ alkane	$-\text{CH}_2-\text{O}-$ ether
$-\text{C}=\text{C}-$ alkene (olefin)	$-\text{N}=\text{N}-$ azo
$-\text{C}\equiv\text{C}-$ alkyne (tolane, acetylene)	$-\text{N}(\text{O})=\text{N}-$ azoxy
$-\text{C}(\text{O})\text{O}-$ ester	$-\text{C}\equiv\text{C}-\text{C}\equiv\text{C}-$ diacetylene

TABLE 2 Typical Ring Systems

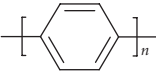

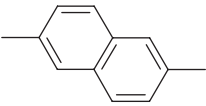
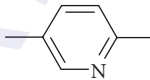
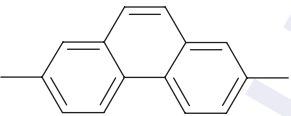
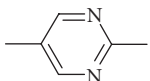

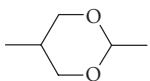
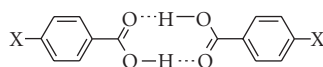
Ring System	Ring System
 $n = 1$ benzene $n = 2$ biphenyl $n = 3$ terphenyl	 bicyclooctane
 naphthalene	 pyridine
 phenanthrene	 pyrimidine
 cyclohexane	 dioxane

TABLE 3 Terminal Units

Terminal Unit A	Terminal Unit B
$C_n H_{2n+1}$ Alkyl	-F, -Cl, -Br, -I Single elements
$OC_n H_{2n+1}$ Alkoxy	-CN, -NO ₂ , -NCS Highly polar moieties

Most calamitic liquid crystals possess aromatic and alicyclic rings in the molecular structures. These include single benzene ring (commonly 1,4-phenyl), alicyclic ring (commonly cyclohexane), heterocyclic ring (e.g., 1,3-pyrimidine), and a wide variety of the combinations. Individual benzene rings do not yield good mesogens. At least two rings are needed to create a rod-like LC. The exception is the 4-alkyl-benzoic acid, which forms a dimer by hydrogen bond, as depicted in Fig. 7. Beside the heterocyclic rings, ring systems of the rigid core could be laterally substituted by elements other than hydrogen (most popular are F and Cl) or groups of elements such as CH₃, NH₂, CN, and so on. This

**FIGURE 7** Dimeric structure formed by 4-alkyl-benzoic acid.

will reduce the melting point of a mesogen, however, increase its viscosity. Additionally, optical and dielectric anisotropies will be affected. In general, the core of the mesogen determines the mesogenic properties by establishing the primary shape of the molecule and its rigidity. Many terminal units (A, B) have been employed in the creation of mesogenic molecules. The most successful route is to use either small polar substituents or a fairly long, straight, hydrocarbon chain. The role of these groups is to act as either a flexible extension to the rigid core or as a dipolar moiety to introduce anisotropy in physical properties. Molecules can obtain chirality from the flexible chain if it is branched and chiral. Further, the terminal moieties are believed to be responsible for stabilizing the molecular order essential for the mesophase generation. All physical properties strongly depend on the terminal units employed in the molecular system of mesogenic molecules. Some of the typical calamitic LC single compounds are listed in Table 4.³

8.4 LIQUID CRYSTALS PHASES

In this section, we focus on the calamitic liquid crystals which have been widely used in display industries. Other LC materials, such as discotic LCs, polymeric LCs, and lyotropic LCs which have great scientific values will not be discussed further. By doing so, we do not intend to minimize the importance of these materials. For example, Kevlar reveals the importance of lyotropic LC polymer material, while others are important in the biological sciences. The readers can find more precise information about these groups of LCs in Ref. 4.

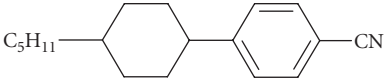
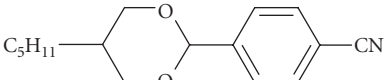
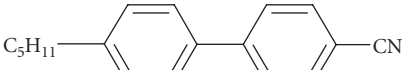
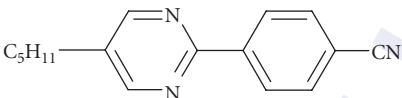
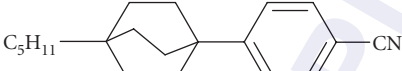
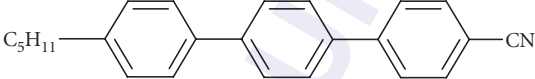
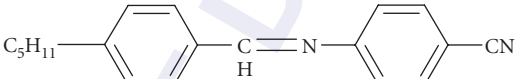
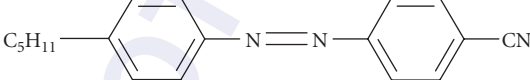
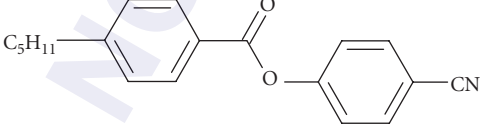
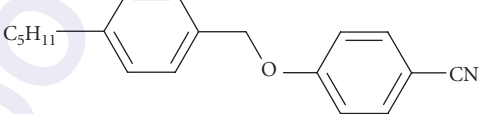
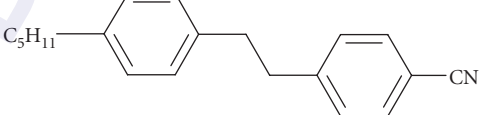
Calamitic Liquid Crystals

When classifying LC phases after G. Friedel,⁵ we first distinguish between two main types: one with nematic order and the other with smectic order (Fig. 8). In the nematic phase (of which tilted nematic is known as the cholesteric phase and is a special case) the molecules are free to move in all directions (e.g., there is no positional order of the center of mass). But on average, they keep their long axes locally parallel. In a smectic state, a number of structural variations exist and there is a positional order along one dimension (some smectic phases have positional order in more than one dimensions). A smectic LC is a layered structure with the molecules oriented parallel or tilted relative to the layer normal. Two smectic phases, called smectic A and smectic C, have acquired a special importance and are now relatively well understood. They are characterized by an absence of positional order within the layers. The molecules have some freedom to move within the layers, as in all smectic phases, but are much less free to move between the layers. These smectics can therefore be described as stacks of two-dimensional fluids, but behave as crystalline across the layers. The absence of in-layer order contributes to their high potential for future electro-optic applications. There are several smectic phases different from one another in areas such as the tilt angle of the director with regard to the layer normal, and the arrangement of molecules within each layer. The simplest is the smectic A phase (Fig. 8) characterized by a director parallel to the layer normal and a random positional order within the plane. Substances featuring the A phase often exhibit the smectic C phase at a lower temperature (Fig. 8). In this phase, the molecules have the same random order within the layer but tilt relative to the layer normal. The tilt angle normally increases with decreasing temperature. The other smectic phases are even more crystalline as they also feature some positional order within the layers. They may, for instance, exhibit hexagonal packing of the molecules.

Chiral Liquid Crystals




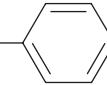

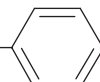
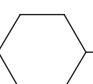


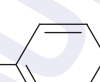
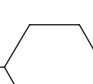
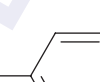

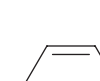

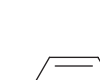

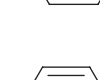
Molecules which are not identical to their mirror image are said to be chiral. A depiction of the simple concept of a chiral molecule is shown in Fig. 9. Another example is the human hand. Chiral molecules are able to form liquid crystals with structures related to those of nonchiral materials but with

TABLE 4 Physical Properties of Popular LC Compounds³

Compound	Phase Transition Temperature(°C)
Rigid core	
	cr 31 N 55 is
	cr 56 (N 52) is
	cr 31 N 55 is
	cr 71 (N 52) is
	cr 62 N 100 is
	cr 130 N 239 is
Central bridge link	
	cr 46 N 75 is
	cr 89 (N 87) is
	cr 65 (N 56) is
	cr 49 (N-20) is
	cr 62 (N -24) is

(Continued)

TABLE 4 Physical Properties of Popular LC Compounds³ (Continued)

Compound	Phase Transition Temperature (°C)
C_5H_{11} —  —C≡C—  —CN	cr 80 (N71) is
C_5H_{11} —  —CH=CH—  —CN	cr 55 N 101 is
Terminal moiety	
C_5H_{11} —  —C(=O)O— 	cr 30 is
C_5H_{11} —  —C(=O)O—  —CH ₃	cr 48 (N 45) is
C_5H_{11} —  —C(=O)O—  —OMe	cr 41 N 71 is
C_5H_{11} —  —C(=O)O—  —Br	cr 78 (N 48) is
C_5H_{11} —  —C(=O)O—  —CN	cr 47 N 79 is
C_5H_{11} —  —C(=O)O—  —NO ₂	cr 53 (N 38) is
C_5H_{11} —  —C(=O)O—  —NCS	cr 87 (N 86) is

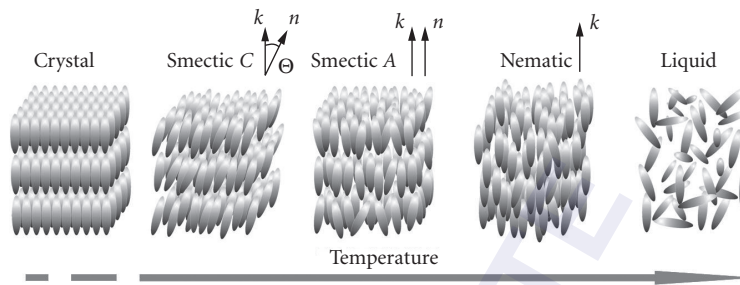


FIGURE 8 Schematic arrangement of the molecules in basic nematic and smectic phases.

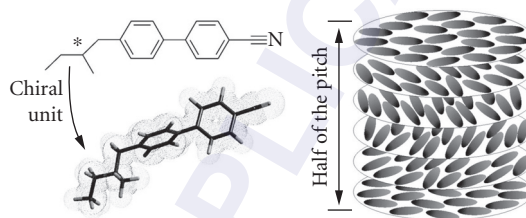


FIGURE 9 Chirality of the calamitic molecule and chiral nematic structure.

different properties. The cholesteric (or chiral nematic) liquid crystal phase is typically composed of nematic mesogenic molecules containing a chiral center which produces intermolecular forces that favor alignment between molecules at a slight angle to one another. When the molecules that make up a nematic liquid crystal are chiral, the chiral nematic phase will exist instead of the normal nematic. In this phase, the molecules prefer to lie next to each other in a slightly skewed orientation. This leads to the formation of a structure which can be visualized as a stack of very thin 2-D nematic-like layers with the director in each layer twisted with respect to those above and below, as Fig. 9 depicts. In this structure, the directors actually form a continuous helical pattern around the layer normal. The molecules shown here are merely representations of the many chiral nematic mesogens lying in the slabs of infinitesimal thickness with a distribution of orientation around the director. This is not to be confused with the planar arrangement found in smectic mesophases.

An important characteristic of the cholesteric mesophase is the pitch. The pitch, p , is defined as the distance it takes for the director to rotate one full turn in the helix as illustrated in Fig. 9. A by-product of the helical structure of the chiral nematic phase is its ability to selectively reflect light with central wavelength located $\lambda_0 = p \cdot \langle n \rangle$ and bandwidth $\Delta\lambda = p \cdot \Delta n$, where $\langle n \rangle$ represents the average refractive index and Δn , the birefringence of the LC. The cholesteric display reflects color without using color filters. Thus, its brightness is about $3\times$ higher than those using color filters. This Bragg reflection is established by the helical pitch of the cholesteric layers. The angle the LC director changes can be made larger by increasing the temperature of the sample, and thus tighten the pitch. Hence, more thermal energy will result from the increased temperature. Similarly, decreasing the temperature of the chiral sample increases the pitch length of the chiral nematic liquid crystal. Similar pitch changes are possible by applying electromagnetic field to aligned cholesteric samples. An interesting phenomenon of chiral nematic phase is that chirality can be introduced to the nonchiral nematic material by adding a small amount of chiral nematic mesogens. Not necessarily all the molecules have to be chiral. Sometimes, slightly below phase transition to the isotropic state (clearing point), some anomalous phases appear. They are known as blue phases (BPs).^{6,7} In many chiral compounds, with sufficiently

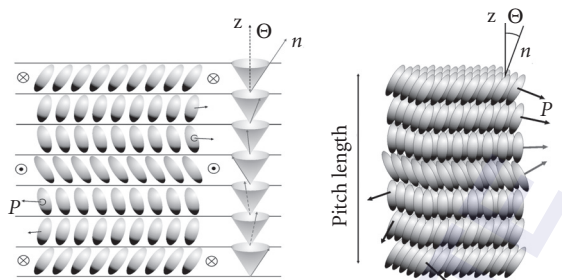


FIGURE 10 Structure of the ferroelectric smectic C^* phase.

high twist, up to three distinct blue phases appear. The two low-temperature phases, blue phase I (BP-I) and blue phase II (BP-II), have cubic symmetry. The highest temperature phase (closest to the clearing point), blue phase III (BP-III) appears to be amorphous.

Similar to chiral nematics, there are chiral forms of smectic phases. The only untilted chiral smectic phase is S_{C^*} . The most important feature of the chiral smectic phases with a tilted structure is ferroelectricity.^{8,9} Due to their low symmetry, chiral smectic phases are able to exhibit spontaneous polarization (P_S) that is oriented perpendicular to the director n and parallel to the smectic layer plane. Ferroelectric smectic C^* phase (S_{C^*}) is the most known tilted chiral smectic phase. The structure of the S_{C^*} has layers of molecules which are tilted in each layer at a temperature-dependent angle (θ) to the layer normal. Additionally, there is a slight and continuous change in the direction of the molecular tilt between adjacent layers as described in Fig. 10. In a macroscopic sample, without surfaces or external electric field, tilt will follow the helix and result in zero total spontaneous polarization. However, if a strong electric field is applied, the helix is unwound and nonzero spontaneous polarization can be observed. Recently, more phases closely related to the ferroelectric S_{C^*} have been discovered. In these phases, the layer spacing and the polarization direction are related in a different manner than in S_{C^*} phase. In an antiferroelectric smectic C^* phase ($S_{C^*}^{anti}$),¹⁰ the spontaneous polarization and the molecular tilt in adjacent layers are pointing in alternating directions, as Fig. 11 shows. Thus, for the macroscopic sample, both average spontaneous polarization and ideally, average molecular tilts are zero. Sufficiently strong electric fields will switch the antiferroelectric order to ferroelectric order. An important difference between ferroelectric and antiferroelectric smectic C^* phase is that antiferroelectric repeats its helical structure every 180° of rotation about the layer normal, compared to 360° for the ferroelectric phase. Many existing ferroelectric phases are different in the proportion of number of the layers with opposite directions. Similar to blue phases in the chiral nematic, the twist grain boundary (TGB) phases appear in chiral smectics.^{6,11} Just like in the blue

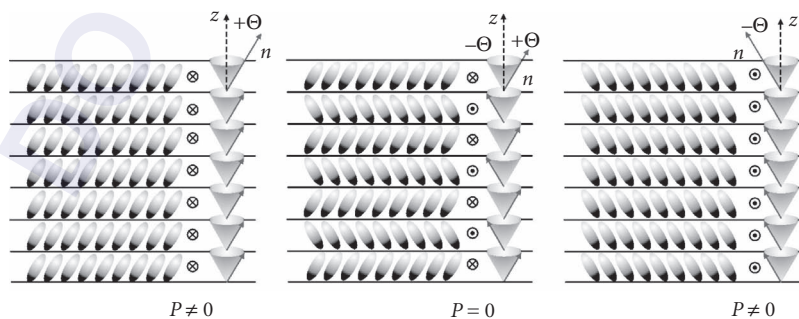


FIGURE 11 Structure of the antiferroelectric smectic C^* phase.

phases, high chirality subtly changes the energy of the system leading to a different type of structure. In this case, the free energy is minimized by introducing grain boundaries at periodic intervals. We may observe different TGB phases depending on the base of the phase from which they were developed. For example, Smectic C^* will generate a TGB $_C$ phase.^{6,12}

8.5 PHYSICAL PROPERTIES

The molecular order existing in the liquid crystalline phases induces anisotropy in the system. What it means is that all directions in the system are not equivalent to each other due to the shape of the molecules and the molecular distribution along the director \mathbf{n} . The anisotropy of the physical property is the most useful feature of liquid crystals and enables electro-optical application of LC materials. The physical properties can be divided into scalar and nonscalar quantities. Typical scalar properties are the thermodynamic transition parameters such as transition temperatures, transition enthalpy, and entropy changes. The dielectric, diamagnetic, optical, elastic, and viscous properties are the most important nonscalar properties. We will concentrate on the physical properties of the nematic phase because we commonly employ this phase and its properties in electro-optical applications.

Phase Transitions

The difference in the transition temperatures between melting and clearing points give the range of stability of the mesophases. For polymorphous (more than one phase) substances, the higher ordered phase exhibits the lower transition temperatures. When a material melts, a change of state occurs from solid to liquid (mesophase) and this process requires energy (endothermic) from the surroundings. If several mesophases exist, then several transitions will occur. The melting transition has typically $\sim 10\times$ larger enthalpy change (30 to 40 kJ/mol) than a transition between different mesophases (3 to 5 kJ/mol). The large enthalpy change is due to the drastic structural changes during the melting process; however, there is only a small difference between the different mesophases. One of the richest polymorphism of the single LC compound is shown by Terephthalylidene-bis-*p*-*n*-pentylaniline (TBPA) with six different mesophases, as shown in Fig. 12. A number of techniques (optical microscopy with hot stage, polarizing optical microscopy, differential thermal analysis, differential scanning calorimetry, etc.) may be used for determining the phase transition temperatures. Polarizing microscopy with hot stage and differential scanning calorimetry is particularly useful for the measurement of the phase transition temperatures. By using both methods, one can determine the number and type of mesophases and also the exact phase transitions, temperature, and enthalpy change associated with each transition. These are the crucial parameters to determine the components of eutectic mixture formulation for devices. A majority of the single liquid crystal components do not possess adequate range of required mesophase. Therefore, the

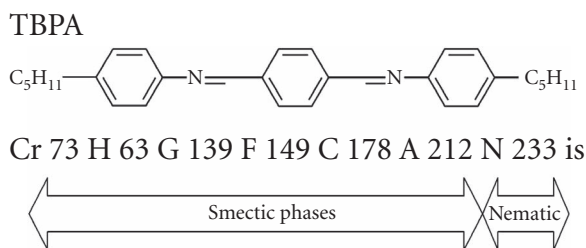


FIGURE 12 Molecular structure and mesomorphic properties of the TBPA.

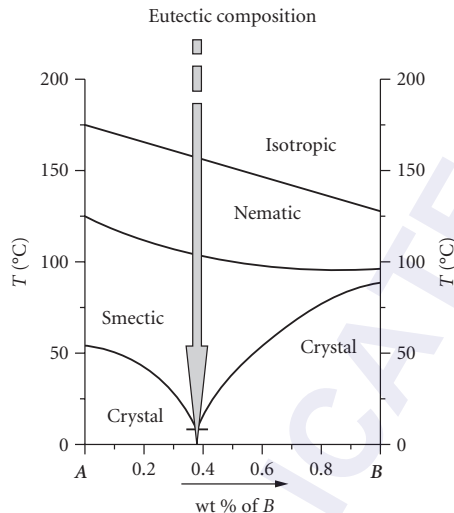


FIGURE 13 Simple phase diagram of the binary LC mixture.

eutectic compositions are required to lower the melting point of the material. It is known that the melting point of a binary (or higher number components) mixture is less than either of its constituent compounds. Figure 13 shows the phase diagram of a binary mixture. The mesogenic range for components 1 and 2 are shown in the two boundary vertical axes. The horizontal axis represents the molar concentration (X_2) of component 2. At a certain molar concentration, the melting point of the mixture will reach its minimum. Meanwhile, the clearing point of the mixture is linearly proportional to the molar concentration. The eutectic mixture calculation is based on the Schroder-Van Laar equation:^{13,14}

$$T_i = \frac{\Delta H_{fi}}{\frac{\Delta H_{fi}}{T_{fi}} - R \ln(X_i)} \quad (3)$$

where T_i is the temperature at which the pure component melts in the mixture, ΔH_{fi} is the heat fusion enthalpy of the pure component i , T_{fi} is the melting point of the pure component i , R is the gas constant (1.98 cal/mol-K), and X_i is the mole fraction of the pure component i . The clearing point (T_c) of the eutectic mixture can be estimated from the clearing points (T_{ci}) of the individual components (X_i) as

$$T_c = \sum_i X_i T_{ci} \quad (4)$$

Dielectric Properties

It was already mentioned, that liquid crystals exhibit anisotropy in many of their physical properties. It is well known that liquid crystals as dielectric and diamagnetic materials are sensitive to the external electric and magnetic fields. Dielectric studies are in general concerned with the response

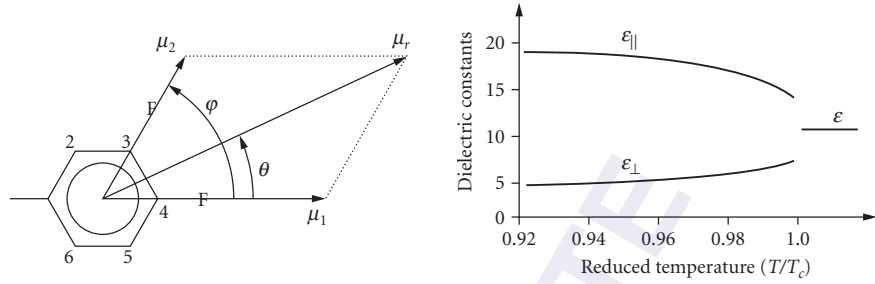


FIGURE 14 Dipole moment of the simple mesogenic structure.

of materials to the application of an electric field. Dielectric constants and their anisotropy affect the sharpness of the voltage-dependent optical transmission curve of a LC device and its threshold voltage.

The nematic liquid crystals are uniaxially symmetric to the axes of the director \mathbf{n} , and the dielectric constants differ in value along the preferred axis (ϵ_{\parallel}) and perpendicular to this axis (ϵ_{\perp}). They are mainly determined by the dipole moment, μ , its angle, θ , with respect to the principal molecular axis, and order parameter, S . When the molecule has two (or more) polar groups with dipole moment μ_1 and μ_2 , its effective dipole can be calculated by the vector addition method, see Fig. 14. In Fig. 14, the first dipole is along the principal molecular axis and the second dipole μ_2 is at an angle ϕ with respect to the principal molecular axis (and μ_1). The resultant dipole moment μ_r can be calculated from following equation:

$$\mu_r = (\mu_1^2 + \mu_2^2 + 2\mu_1\mu_2 \cos\phi)^{1/2} \quad (5)$$

Maier and Meier¹⁵ have developed a mean field theory to correlate the microscopic molecular parameters with the macroscopic dielectric constants of anisotropic LCs:

$$\epsilon_{\parallel} = NhF\{\langle\alpha_{\parallel}\rangle + (F\mu^2/3kT)[1 - (1 - 3\cos^2\theta)S]\} \quad (6a)$$

$$\epsilon_{\perp} = NhF\{\langle\alpha_{\perp}\rangle + (F\mu^2/3kT)[1 + (1 - 3\cos^2\theta)S/2]\} \quad (6b)$$

$$\Delta\epsilon = NhF\{(\langle\alpha_{\parallel}\rangle - \langle\alpha_{\perp}\rangle) - (F\mu^2/2kT)(1 - 3\cos^2\theta)S\} \quad (6c)$$

where N is the molecular packing density, $h = 3\epsilon/(2\epsilon + 1)$ is the cavity field factor, $\epsilon = (\epsilon_{\parallel} + 2\epsilon_{\perp})/3$ is the averaged dielectric constant, F is the Onsager reaction field, $\langle\alpha_{\parallel}\rangle$ and $\langle\alpha_{\perp}\rangle$ are the principal elements of the molecular polarizability tensor, β is the angle between the dipole moment μ and the principal molecular axis, and S is the order parameter of the second rank. From Eq. (6), the dielectric constants of anisotropic liquid crystals are influenced by the molecular structure, temperature, and frequency. These individual effects are discussed separately.

Structural Effect In general, $\Delta\epsilon$ depends on the molecular constituents, temperature, and frequency. For a nonpolar liquid crystal compound, its dipole moment $\mu \sim 0$. Thus, its $\Delta\epsilon$ is expected to be small and its magnitude is proportional to the differential molecular polarizability, similar to birefringence. On the other hand, for a LC molecule containing a polar group, such as cyano, isocyanate, fluoro, or chloro group, its $\Delta\epsilon$ can be positive or negative depending on the position(s) of the polar group(s). If $\beta \sim 0$, that is, the dipole moment of the polar group is along the principal molecular axis, $\Delta\epsilon$ is large and positive. Cyano-biphenyls are such examples. The $\Delta\epsilon$ of 5CB is about 10 at $T = 20^\circ\text{C}$ and $f = 1$ kHz. These positive $\Delta\epsilon$ materials are useful for parallel or twist alignment. On the contrary, if $\beta > 55^\circ$, $1 - 3\cos^2\beta > 0$ and $\Delta\epsilon$ may become negative depending on the dipole moment as indicated in Eq. (6c).

From Eq. (6c), for a nonpolar compound, $\mu \sim 0$ and its dielectric anisotropy is very small ($\Delta\epsilon < 0.5$). In this case, $\Delta\epsilon$ is determined mainly by the differential molecular polarizability, that is, the first term in Eq. (6c). For a polar compound, the dielectric anisotropy depends on the dipole moment, angle θ , temperature (T) and applied frequency. If a LC has more than one dipole, then the resultant dipole moment is their vector summation. In a phenyl ring, the position of the dipole is defined as

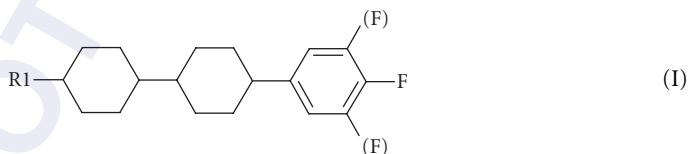


From Eq. (6c), if a polar compound has an effective dipole at $\theta < 55^\circ$, then its $\Delta\epsilon$ is positive. On the other hand, $\Delta\epsilon$ becomes negative if $\theta > 55^\circ$.

Fluoro (F)¹⁶, cyano (CN)¹⁷, and isothiocyanato (NCS)¹⁸ are the three commonly employed polar groups. Among them, fluoro group possess a modest dipole moment ($\mu \sim 1.5$ Debye), high resistivity, and low viscosity. However, its strong negativity compresses the electron clouds and, subsequently, lowers the compound's birefringence. For direct-view LCDs, the required birefringence is around 0.1, depending on the LC alignment and cell gap (d) employed. On the other hand, the cyano and isothiocyanato groups not only exhibit a large dipole moment ($\mu \sim 3.9$ Debye for $C\equiv N$ and ~ 3.7 Debye for $N=C=S$) but also contribute to lengthen the π -electron conjugation. As a result, their birefringence is much higher than their fluorinated counterpart. High birefringence is favorable for long wavelength, such as infrared applications in order to use a thin cell gap to achieve the same phase change while keeping a fast response time. Under strong anchoring condition, the LC response time is proportional to d^2 . However, the CN compounds are more viscous than the corresponding NCS and fluoro compounds. Therefore, their major applications are in low-end displays such as wrist watches and calculators where response time is not crucial.

Example 1: Positive $\Delta\epsilon$ LCs Positive $\Delta\epsilon$ LCs have been used in twisted nematic (TN)¹⁹ and in-plane switching (IPS)^{20,21} displays, although IPS can also use negative $\Delta\epsilon$ LCs. For TFT LCD, the employed LC material must also possess a high resistivity ($>10^{13}$ Ω -cm) in order to steadily hold the charges and avoid image flickering.²² The resistivity of a LC mixture depends heavily on the impurity contents, for example, ions. Purification process plays an important role in removing the ions for achieving high resistivity. Fluorinated compounds exhibit a high resistivity and are the natural choices for TFT LCDs.^{23,24}

A typical fluorinated LC structure is shown below:



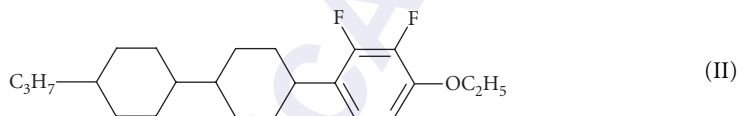
Most liquid crystal compounds discovered so far possess at least two rings, either cyclohexane-cyclohexane, cyclohexane-phenyl or phenyl-phenyl, and a flexible alkyl or alkoxy chain. The compound shown in structure (I) has two cyclohexane and one phenyl rings. The R_1 group represents a terminal alkyl chain, and a single or multiple fluoro substitutions take place in the phenyl ring. For multiple dipoles, the net dipole moment can be calculated from their vector sum. From Eq. (6c), to obtain the largest $\Delta\epsilon$ for a given dipole, the best position for the fluoro substitution is along the principal molecular axis, that is, in the fourth position. The single fluoro compound should have $\Delta\epsilon \sim 5$. To further increase $\Delta\epsilon$, more fluoro groups can be added. For example, compound (I) has two more fluoro groups in the third and fifth positions.²⁴ Its $\Delta\epsilon$ is about 10, but its birefringence would slightly decrease (because of the lower molecular packing density) and viscosity increases substantially (because of the higher moment of inertia). The birefringence of compound (I) is around 0.07. If a higher birefringence is needed, the middle cyclohexane ring can be replaced by a phenyl ring. The elongated electron cloud will enhance the birefringence to approximately 0.12 without increasing the viscosity noticeably.

The phase transition temperatures of a LC compound are difficult to predict before the compound is synthesized. In general, the lateral fluoro substitution lowers the melting temperature of the parent

compound because the increased intermolecular separation leads to a weaker molecular association. Thus, a smaller thermal energy is able to separate the molecules which implies to a lower melting point. A drawback of the lateral substitution is the increased viscosity.

Example 2: Negative $\Delta\epsilon$ LCs From Eq. (6c), in order to obtain a negative dielectric anisotropy, the dipoles should be in the lateral (2,3) positions. For the interest of obtaining high resistivity, lateral difluoro group is a favorable choice. The negative $\Delta\epsilon$ LCs are useful for vertical alignment.²⁵ The VA cell exhibits an unprecedented contrast ratio when viewed at normal direction between two crossed linear polarizers.^{26,27} However, a single domain VA cell has a relatively narrow viewing angle and is only useful for projection displays. For wide-view LCDs, a multidomain (4 domains) vertical alignment (MVA) cell is required.²⁸

The following structure is an example of the negative $\Delta\epsilon$ LC:²⁹



Compound (II) has two lateral fluoro groups in the (2,3) positions so that their dipoles in the horizontal components are perfectly cancelled whereas the vertical components add up. Thus, the net $\Delta\epsilon$ is negative. A typical $\Delta\epsilon$ of lateral difluoro compounds is -4 . The neighboring alkoxy group also has a dipole in the vertical direction. Therefore, it contributes to enlarge the dielectric anisotropy ($\Delta\epsilon \sim -6$). However, the alkoxy group has a higher viscosity than its alkyl counterpart and it also increases the melting point by $\sim 20^\circ$.

Temperature Effect In general, as temperature rises ϵ_{\parallel} decreases but ϵ_{\perp} increases gradually resulting in a decreasing $\Delta\epsilon$. From Eq. (6c) the temperature dependence of $\Delta\epsilon$ is proportional to S for the nonpolar LCs and S/T for the polar LCs. At $T > T_c$, the isotropic phase is reached and dielectric anisotropy vanishes, as Fig. 14 shows.

Frequency Effect From Eq. (6), two types of polarizations contribute to the dielectric constant: (1) induced polarization (the first term), and (2) orientation polarization (the dipole moment term). The field-induced polarization has a very fast response time, and it follows the alternating external field. But the permanent dipole moment associated orientation polarization exhibits a longer decay time, τ . If the external electric field frequency is comparable to $1/\tau$, the time lag between the average orientation of the dipole moments and the alternating field becomes noticeable. At a frequency $\omega (=2\pi f)$ which is much higher than $1/\tau$, the orientation polarization cannot follow the variations of the external field any longer. Thus, the dielectric constant drops to ϵ_{\parallel} which is contributed solely by the induced polarization:³⁰

$$\epsilon_{\parallel}(\omega) = \epsilon_{\infty} + \frac{\epsilon_{\parallel} - \epsilon_{\infty}}{1 + \omega^2 \tau^2} \quad (7)$$

where $\epsilon_{\parallel} = \epsilon_{\parallel}(\omega = 0)$ and $\epsilon_{\infty} = \epsilon_{\parallel}(\omega = \infty)$ are the parallel component of the dielectric constant at static and high frequencies, respectively.

In an aligned LC, the molecular rotation around their short axis is strongly hindered. Thus, the frequency dispersion occurs mainly at ϵ_{\parallel} , while ϵ_{\perp} remains almost constant up to mega-Hertz region. Figure 15 shows the frequency dependent dielectric constants of the M1 LC mixture (from Roche) at various temperatures.³¹ As the frequency increases ϵ_{\parallel} decreases and beyond the crossover frequency f_c , $\Delta\epsilon$ changes sign. The dielectric anisotropies of M1 are fairly symmetric at low and high frequencies. The crossover frequency is sensitive to the temperature. As temperature rises, the ϵ_{\parallel} and ϵ_{\perp} of M1 both decrease slightly. However, the frequency-dependent ϵ_{\parallel} is strongly dependent on the temperature, but ϵ_{\perp} is inert. Thus, the crossover frequency increases exponentially with the temperature as: $f_c \sim \exp(-E/kT)$, where E is the activation energy. For M1 mixture, $E = 0.96$ eV.³¹

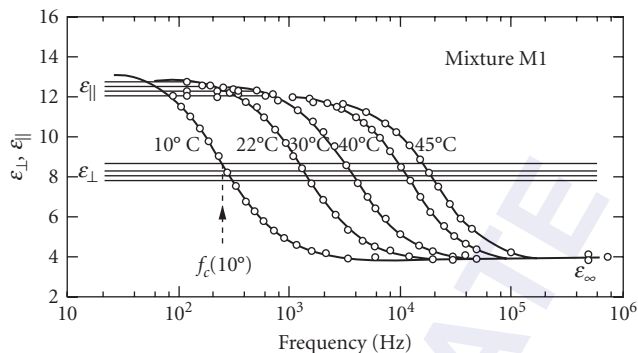


FIGURE 15 The frequency dependent dielectric constants of the M1 LC mixture (from Roche). (Redrawn from Ref. 31.)

Dual frequency effect is a useful technique for improving the response times of a LC device.³² In the dual frequency effect, a low frequency ($f < f_c$, where $\Delta\epsilon > 0$) an electric field is used to drive the device to its ON state, and during the decay period a high frequency ($f > f_c$, where $\Delta\epsilon < 0$) electric field is applied to speed up the relaxation process. From material standpoint, a LC mixture with low f_c and large $|\Delta\epsilon|$ at both low and high frequencies is beneficial. But for a single LC substance (such as cyanobiphenyls), its f_c is usually too high ($>10^6$ Hz) to be practically employed.

In such a high frequency regime, the imaginary part of dielectric constant (which is responsible for absorption) becomes so significant that the dielectric heating effect is amplified and heats up the LC. The electro-optic properties of the cell are then altered. The imaginary part of dielectric constant contributes to heat. Thus, if a LC device is operated at MHz frequency region, significant heating effect due to the applied voltage will take place. This heating effect may be large enough to change all the physical properties of the LC. Dielectric heating is more severe if the crossover frequency is high.^{33–35}

Thus, liquid crystals are useful electro-optic media in the spectral range covering from UV, visible, IR to microwave. Of course, in each spectral region, an appropriate LC material has to be selected. For instance, the totally saturated LC compounds should be chosen for UV application because of photostability. For flat panel displays, LCs with a modest conjugation are appropriate. On the other hand, highly conjugated LCs are favorable for IR and microwave applications for the interest of keeping fast response time.

Optical Properties

Refractive indices and absorption are fundamentally and practically important parameters of a LC compound or mixture.³⁶ Almost all the light modulation mechanisms are involved with refractive index change. The absorption has a crucial impact on the photostability or lifetime of the liquid crystal devices. Both refractive indices and absorption are determined by the electronic structures of the liquid crystal studied.

The major absorption of a LC compound occurs in ultraviolet (UV) and infrared (IR) regions. The $\sigma \rightarrow \sigma^*$ electronic transitions take place in the vacuum UV (100 to 180 nm) region whereas the $\pi \rightarrow \pi^*$ electronic transitions occur in the UV (180 to 400 nm) region. Figure 16 shows the measured polarized UV absorption spectra of 5CB.³⁷ The λ_1 band which is centered at ~ 200 nm consists of two closely overlapped bands. The λ_2 band shifts to ~ 282 nm. The λ_2 band should occur in the vacuum UV region ($\lambda_0 \sim 120$ nm) which is not shown in the figure.

Refractive Indices Refractive index has great impact on LC devices. Almost every electro-optic effect of LC modulators, no matter amplitude or phase modulation, involves refractive index

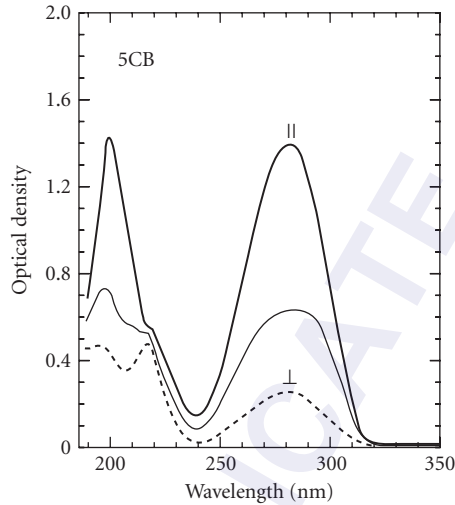


FIGURE 16 The measured polarized absorption spectra of 5CB. The middle trace is for unpolarized light. $\lambda_1 \sim 200$ nm and $\lambda_2 \sim 282$ nm.

change. An aligned LC exhibits anisotropic properties, including dielectric, elastic, and optical anisotropies. Let us take a homogeneous alignment as an example.³⁸ Assume a linearly polarized light is incident to the LC cell at normal direction. If the polarization axis is parallel to the LC alignment axis (i.e., LC director which represents an average molecular distribution axis), then the light experiences the extraordinary refractive index, n_e . If the polarization is perpendicular to the LC directors, then the light sees the ordinary refractive index n_o . The difference between n_e and n_o is called birefringence, defined as $\Delta n = n_e - n_o$. Refractive indices are dependent on the wavelength and temperature. For a full-color LCD, RGB color filters are employed. Thus, the refractive indices at these wavelengths need to be known in order to optimize the device performance. Moreover, about 50 percent of the backlight is absorbed by the polarizer. The absorbed light turns into heat and causes the LCD panel's temperature to increase. As the temperature increases, refractive indices decrease gradually. The following sections will describe how the wavelength and temperature affect the LC refractive indices.

Wavelength Effect Based on the electronic absorption, a three-band model which takes one $\sigma \rightarrow \sigma^*$ transition (the λ_0 -band) and two $\pi \rightarrow \pi^*$ transitions (the λ_1 - and λ_2 -bands) into consideration has been developed. In the three band model, the refractive indices (n_e and n_o) are expressed as follows:^{39,40}

$$n_{e,o} \cong 1 + g_{0e,o} \frac{\lambda^2 \lambda_0^2}{\lambda^2 - \lambda_0^2} + g_{1e,o} \frac{\lambda^2 \lambda_1^2}{\lambda^2 - \lambda_1^2} + g_{2e,o} \frac{\lambda^2 \lambda_2^2}{\lambda^2 - \lambda_2^2} \quad (8)$$

The three-band model clearly describes the origins of refractive indices of LC compounds. However, a commercial mixture usually consists of several compounds with different molecular structures in order to obtain a wide nematic range. The individual λ_i 's are therefore different. Under such a circumstance, Eq. (8) would have too many unknowns to quantitatively describe the refractive indices of a LC mixture.

In the off-resonance region, the right three terms in Eq. (8) can be expanded by a power series to the λ^{-4} terms to form the extended Cauchy equations for describing the wavelength-dependent refractive indices of *anisotropic* LCs:^{40,41}

$$n_{e,o} \equiv A_{e,o} + \frac{B_{e,o}}{\lambda^2} + \frac{C_{e,o}}{\lambda^4} \quad (9)$$

In Eq. (9), $A_{e,o}$, $B_{e,o}$, and $C_{e,o}$ are known as Cauchy coefficients. Although Eq. (9) is derived based on a LC compound, it can be extended easily to include eutectic mixtures by taking the superposition of each compound. From Eq. (9) if we measure the refractive indices at three wavelengths, the three Cauchy coefficients ($A_{e,o}$, $B_{e,o}$, and $C_{e,o}$) can be obtained by fitting the experimental results. Once these coefficients are determined, the refractive indices at any wavelength can be calculated. From Eq. (9) both refractive indices and birefringence decrease as the wavelength increases. In the long wavelength (IR and millimeter wave) region, n_e and n_o are reduced to A_e and A_o , respectively. The coefficients A_e and A_o are constants; they are independent of wavelength, but dependent on the temperature. That means, in the IR region the refractive indices are insensitive to wavelength, except for the resonance enhancement effect near the local molecular vibration bands. This prediction is consistent with many experimental evidences.⁴²

Figure 17 depicts the wavelength-dependent refractive indices of E7 at $T = 25^\circ\text{C}$. Open squares and circles represent the n_e and n_o of E7 in the visible region while the downward- and upward-triangles stand for the measured data at $\lambda = 1.55$ and $10.6 \mu\text{m}$, respectively. Solid curves are fittings to the experimental n_e and n_o data in the visible spectrum by using the extended Cauchy equations [Eq. (9)]. The fitting parameters are listed as follows: ($A_e = 1.6933$, $B_e = 0.0078 \mu\text{m}^2$, $C_e = 0.0028 \mu\text{m}^4$) and ($A_o = 1.4994$, $B_o = 0.0070 \mu\text{m}^2$, $C_o = 0.004 \mu\text{m}^4$). In Fig. 17, the extended Cauchy model is extrapolated to the near- and far-infrared regions. The extrapolated lines almost strike through the center of the experimental data measured at $\lambda = 1.55$ and $10.6 \mu\text{m}$. The largest difference between the extrapolated and experimental data is only 0.4 percent.

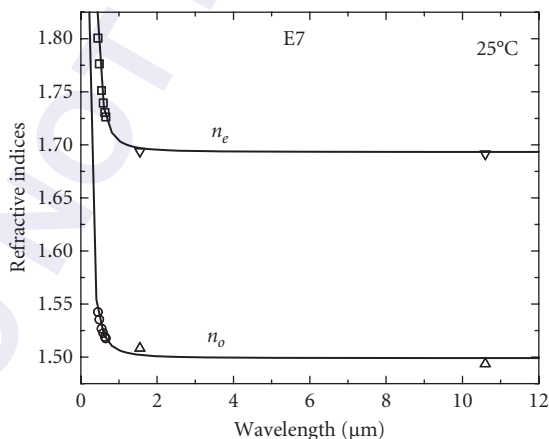


FIGURE 17 Wavelength-dependent refractive indices of E7 at $T = 25^\circ\text{C}$. Open squares and circles are the n_e and n_o measured in the visible spectrum. Solid lines are fittings to the experimental data measured in the visible spectrum by using the extended Cauchy equation [Eq. (4.9)]. The downward- and upward-triangles are the n_e and n_o measured at $T = 25^\circ\text{C}$ and $\lambda = 1.55$ and $10.6 \mu\text{m}$, respectively.

Equation (9) applies equally well to both high and low birefringence LC materials in the off-resonance region. For low birefringence ($\Delta n < 0.12$) LC mixtures, the λ^{-4} terms are insignificant and can be omitted and the extended Cauchy equations are simplified as:⁴³

$$n_{e,o} \cong A_{e,o} + \frac{B_{e,o}}{\lambda^2} \quad (10)$$

Thus, n_e and n_o each has only two fitting parameters. By measuring the refractive indices at two wavelengths, we can determine $A_{e,o}$ and $B_{e,o}$. Once these two parameters are determined, n_e and n_o can be calculated at any wavelength of interest. Because most of TFT LC mixtures have $\Delta n \sim 0.1$, the two-coefficient Cauchy model is adequate to describe the refractive index dispersions. Although the extended Cauchy equation fits experimental data well,⁴⁴ its physical origin is not clear. A better physical meaning can be obtained by the three-band model which takes three major electronic transition bands into consideration.

Temperature Effect The temperature effect is particularly important for projection displays.⁴⁵ Due to the thermal effect of the lamp, the temperature of the display panel could reach 50°C. It is important to know the LC properties at the anticipated operating temperature beforehand.

Birefringence Δn is defined as the difference between the extraordinary and ordinary refractive indices, $\Delta n = n_e - n_o$ and the average refractive indices $\langle n \rangle$ is defined as $\langle n \rangle = (n_e + 2n_o)/3$. Based on these two definitions, n_e and n_o can be rewritten as

$$n_e = \langle n \rangle + \frac{2}{3} \Delta n \quad (11)$$

$$n_o = \langle n \rangle - \frac{1}{3} \Delta n \quad (12)$$

To describe the temperature dependent birefringence, the Haller approximation can be employed when the temperature is not too close to the clearing point:

$$\Delta n(T) = (\Delta n)_o (1 - T/T_c)^\beta \quad (13)$$

In Eq. (13), $(\Delta n)_o$ is the LC birefringence in the crystalline state (or $T = 0$ K), the exponent α is a material constant, and T_c is the clearing temperature of the LC material under investigation. On the other hand, the average refractive index decreases linearly with increasing temperature as⁴⁶

$$\langle n \rangle = A - BT \quad (14)$$

because the LC density decreases with increasing temperature. By substituting Eqs. (14) and (13) back to Eqs. (11) and (12), the 4-parameter model for describing the temperature dependence of the LC refractive indices is given as⁴⁷

$$n_e(T) \approx A - BT + \frac{2(\Delta n)_o}{3} \left(1 - \frac{T}{T_c}\right)^\beta \quad (15)$$

$$n_o(T) \approx A - BT - \frac{(\Delta n)_o}{3} \left(1 - \frac{T}{T_c}\right)^\beta \quad (16)$$

The parameters $[A, B]$ and $[(\Delta n)_o, \beta]$ can be obtained separately by two-stage fittings. To obtain $[A, B]$, one can fit the average refractive index $\langle n \rangle = (n_e + 2n_o)/3$ as a function of temperature using Eq. (14). To find $[(\Delta n)_o, \beta]$, one can fit the birefringence data as a function of temperature using Eq. (13). Therefore, these two sets of parameters can be obtained separately from the same set of refractive indices but at different forms.

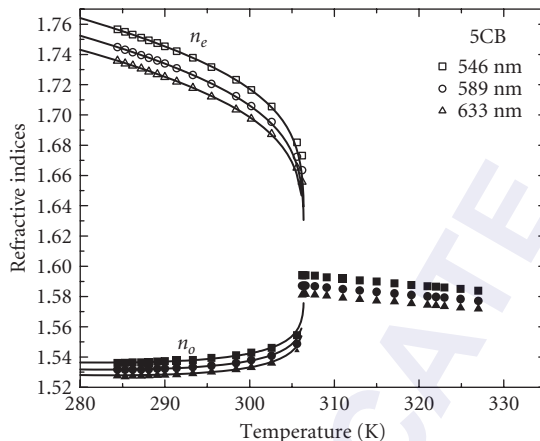


FIGURE 18 Temperature-dependent refractive indices of 5CB at $\lambda = 546, 589,$ and 633 nm. Squares, circles, and triangles are experimental data for refractive indices measured at $\lambda = 546, 589,$ and 633 nm, respectively.

Figure 18 is a plot of the temperature dependent refractive indices of 5CB at $\lambda = 546, 589,$ and 633 nm. As the temperature increases, n_e decreases, but n_o gradually increases. In the isotropic state, $n_e = n_o$ and the refractive index decreases linearly as the temperature increases. This correlates with the density effect.

Elastic Properties

The molecular order existing in liquid crystals has interesting consequences on the mechanical properties of these materials. They exhibit elastic behavior. Any attempt to deform the uniform alignments of the directors and the layered structures (in case of smectics) results in an elastic restoring force. The constants of proportionality between deformation and restoring stresses are known as elastic constants.

Elastic Constants Both threshold voltage and response time are related to the elastic constant of the LC used. There are three basic elastic constants involved in the electro-optics of liquid crystals depending on the molecular alignment of the LC cell: the splay (K_{11}), twist (K_{22}), and bend (K_{33}), as Fig. 19 shows. Elastic constants affect the liquid crystal's electro-optical cell in two aspects: the threshold voltage and the response time. The threshold voltage in the most common case of homogeneous electro-optical cell is expressed as follows:

$$V_{th} = \pi \sqrt{\frac{K_{11}}{\epsilon_0 \Delta \epsilon}} \quad (17)$$

Several molecular theories have been developed for correlating the Frank elastic constants with molecular constituents. Here we only introduce two theories: (1) the mean-field theory,^{48,49} and (2) the generalized van der Waals theory.⁵⁰

Mean-Field Theory In the mean-field theory, the three elastic constants are expressed as

$$K_{ii} = C_{ii} V_n^{-7/3} S^2 \quad (18a)$$

$$C_{ii} = (3A/2)(Lm^{-1}\chi_{ii}^{-2})^{1/3} \quad (18b)$$

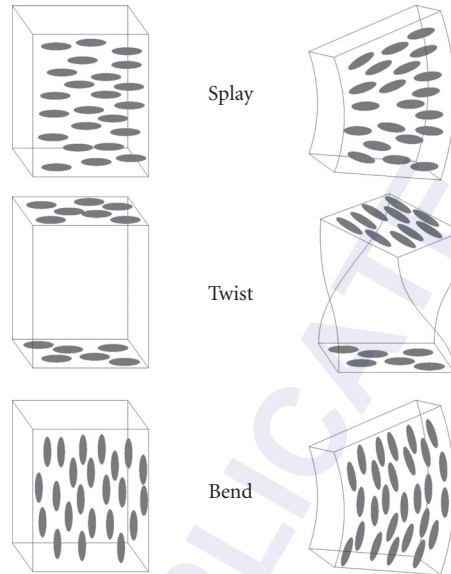


FIGURE 19 Elastic constants of the liquid crystals.

where C_{ii} is called the reduced elastic constant, V_n is the mole volume, L is the length of the molecule, m is the number of molecules in a steric unit in order to reduce the steric hindrance, $\chi_{11} = \chi_{22} = z/x$ and $\chi_{33} = (x/z)^2$, where x ($y = x$) and z are the molecular width and length, respectively, and $A = 1.3 \times 10^{-8}$ erg \cdot cm⁶.

From Eq. (18), the ratio of $K_{11} : K_{22} : K_{33}$ is equal to 1: 1: $(z/x)^2$ and the temperature dependence of elastic constants is basically proportional to S^2 . This S^2 dependence has been experimentally observed for many LCs. However, the prediction for the relative magnitude of K_{ii} is correct only to the first order. Experimental results indicate that K_{22} often has the lowest value, and the ratio of K_{33}/K_{11} can be either greater or less than unity.

Generalized van der Waals Theory Gelbart and Ben-Shaul⁵⁰ extended the generalized van der Waals theory for explaining the detailed relationship between elastic constants and molecular dimensions, polarizability, and temperature. They derived the following formula for nematic liquid crystals:

$$K_{ii} = a_i \langle P_2 \rangle \langle P_2 \rangle + b_i \langle P_2 \rangle \langle P_4 \rangle \quad (19)$$

where a_i and b_i represent sums of contributions of the energy and the entropy terms; they depend linearly on temperature, $\langle P_2 \rangle$ ($=S$) and $\langle P_4 \rangle$ are the order parameter of the second and the fourth rank, respectively. In general, the second term may not be negligible in comparison with the S^2 term depending on the value of $\langle P_4 \rangle$. As temperature increases, both S and $\langle P_4 \rangle$ decrease. If the $\langle P_4 \rangle$ of a LC is much smaller than S in its nematic range, Eq. (19) is reduced to the mean-field theory, or $K_{ii} \sim S^2$. The second term in Eq. (19) is responsible for the difference between K_{11} and K_{33} .

Viscosities

The resistance of fluid system to flow when subjected to a shear stress is known as viscosity. In liquid crystals several anisotropic viscosity coefficients may result, depending on the relative orientation

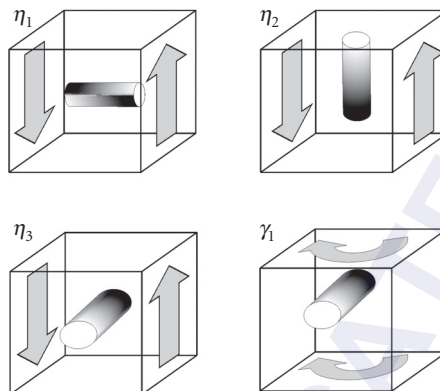


FIGURE 20 Anisotropic viscosity coefficients required to characterize a nematic.

of the director with respect to the flow of the LC material. When an oriented nematic liquid crystal is placed between two plates which are then sheared, four cases shown in Fig. 20 are to be studied. Three, are known as Miesowicz viscosity coefficients. They are η_1 —when director of LC is perpendicular to the flow pattern and parallel to the velocity gradient, η_2 —when director is parallel to the flow pattern and perpendicular to the velocity gradient and η_3 —when director is perpendicular to the flow pattern and to the velocity gradient. Viscosity, especially rotational viscosity γ_1 , plays a crucial role in the liquid crystals displays (LCD) response time. The response time of a nematic liquid crystals device is linearly proportional to γ_1 .

The rotational viscosity of an aligned LC is a complicated function of molecular shape, moment of inertia, activation energy, and temperature. Several theories, including both rigorous and semi-empirical, have been developed in an attempt to account for the origin of the LC viscosity. However, owing to the complicated anisotropic attractive and steric repulsive interactions among LC molecules, these theoretical results are not yet completely satisfactory. Some models fit certain LCs, but fail to fit others.⁵¹ In the molecular theory developed by Osipov and Terentjev,⁵² all six Leslie viscosity coefficients are expressed in terms of microscopic parameters:

$$\alpha_1 = -N\Lambda \frac{p^2 - 1}{p^2 + 1} P_4 \quad (20a)$$

$$\alpha_2 = -\frac{N\Lambda}{2} S - \frac{1}{2} \gamma_1 \cong -\frac{N\Lambda}{2} \left[S + \frac{1}{6} \sqrt{J/kT} e^{j/kT} \right] \quad (20b)$$

$$\alpha_3 = -\frac{N\Lambda}{2} S + \frac{1}{2} \gamma_1 \cong -\frac{N\Lambda}{2} \left[S - \frac{1}{6} \sqrt{J/kT} e^{j/kT} \right] \quad (20c)$$

$$\alpha_4 = \frac{N\Lambda}{35} \frac{p^2 - 1}{p^2 + 1} [7 - 5S - 2P_4] \quad (20d)$$

$$\alpha_5 = \frac{N\Lambda}{2} \left[\frac{p^2 - 1}{p^2 + 1} \frac{3S + 4P_4}{7} + S \right] \quad (20e)$$

$$\alpha_6 = \frac{N\Lambda}{2} \left[\frac{p^2 - 1}{p^2 + 1} \frac{3S + 4P_4}{7} - S \right] \quad (20f)$$

where N represents molecular packing density, p is the molecular length-to-width (w) ratio, $J \equiv J_o S$ is the Maier-Saupe mean-field coupling constant; k is the Boltzmann constant, T is the Kelvin temperature, S and P_4 are the nematic order parameters of the second and fourth order, respectively, and Λ is the friction coefficient of LC:

$$\Lambda \approx 100 (1 - \Phi) N^2 w^6 p^{-2} [(kT)^5 / G^3] (I_{\perp} / kT)^{1/2} \cdot \exp [3(G + I_o) / kT] \quad (21)$$

where Φ is the volume fraction of molecules ($\Phi = 0.5$ to 0.6 for dense molecular liquids), $G \gg I_o$ is the isotropic intermolecular attraction energy, and I_{\perp} and I_o are the inertia tensors.

From the above analysis, the parameters affecting the rotational viscosity of a LC are⁵³

1. Activation energy ($E_a \sim 3G$): a LC molecule with low activation energy leads to a low viscosity.
2. Moment of inertia: a LC molecule with linear shape and low molecular weight would possess a small moment of inertia and exhibit a low viscosity.
3. Inter-molecular association: a LC with weak inter-molecular association, for example, not form dimer, would reduce the viscosity significantly.
4. Temperature: elevated temperature operation of a LC device may be the easiest way to lower viscosity. However, birefringence, elastic, and dielectric constants are all reduced as well.

8.6 LIQUID CRYSTAL CELLS

Three kinds of LC cells have been widely used for display applications. They are (1) twisted-nematic (TN) cell, (2) in-plane switching (IPS) cell, and multidomain vertical alignment (MVA) cell. For phase-only modulation, homogeneous cell is preferred. The TN cell dominates notebook market because of its high transmittance and low cost. However, its viewing angle is limited. For wide-view applications, for example, LCD TVs, optical film-compensated IPS and MVA are the two major camps. In this section, we will discuss the basic electro-optics of TN, IPS, and MVA cells.

Twisted-Nematic (TN) Cell

The first liquid crystal displays that became successful on the market were the small displays in digital watches. These were simple twisted nematic (TN) devices with characteristics satisfactory for such simple applications.¹⁹ The basic operation principals are shown in Fig. 21. In the TN display, each LC cell consists of a LC material sandwiched between two glass plates separated by a gap of 5 to 8 μm .

The inner surfaces of the plates are deposited with transparent electrodes made of conducting coatings of indium tin oxide (ITO). These transparent electrodes are overcoated with a thin layer of polyimide with a thickness of about 80 angstroms. The polyimide films are unidirectional rubbed with the rubbing direction of the lower substrate perpendicular to the rubbing direction of the upper surface. Thus, in the inactivated state (voltage OFF), the local director undergoes a continuous twist of 90° in the region between the plates. Sheet polarizers are laminated on the outer surfaces of the plates. The transmission axes of the polarizers are aligned parallel to the rubbing directions of the adjacent polyimide films. When light enters the cell, the first polarizer lets through only the component oscillating parallel to the LC director next to the entrance substrate. During the passage through the cell, the polarization plane is turned along with the director helix, so that when the light wave arrives at the exit polarizer, it passes unobstructed. The cell is thus transparent in the OFF state; this mode is called normally white (NW).

Figure 22 depicts the normalized voltage-dependent light transmittance (T_{\perp}) of the 90° TN cell at three primary wavelengths: $R = 650$ nm, $G = 550$ nm, and $B = 450$ nm. Since human eye has the greatest

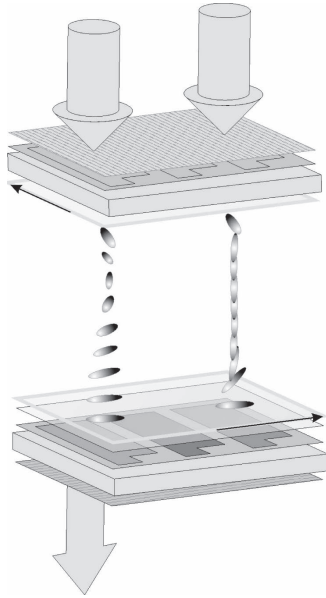


FIGURE 21 The basic operation principles of the TN cell.

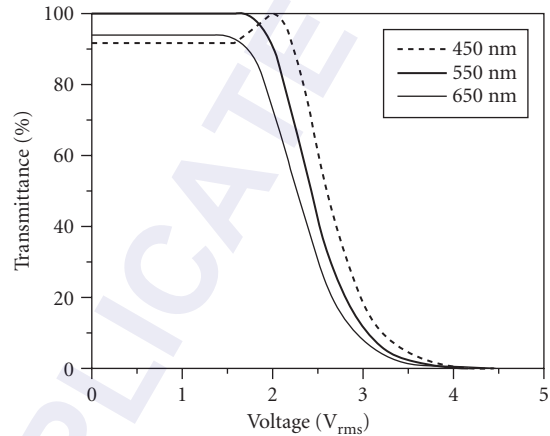


FIGURE 22 Voltage-dependent transmittance of a normally white 90° TN cell. $d\Delta n = 480$ nm.

sensitivity at green, we normally optimize the cell design at $\lambda \sim 550$ nm. To meet the Gooch-Tarry first minimum condition, for example, $d\Delta n = (\sqrt{3}/2)\lambda$, the employed cell gap is $5\text{-}\mu\text{m}$ and the LC birefringence is $\Delta n \sim 0.096$. From Fig. 22, the wavelength effect on the transmittance at $V = 0$ is within 8 percent. Therefore, the TN cell can be treated as an “achromatic” half-wave plate.

The response time of a TN LCD depends on the cell gap and the γ_1/K_{22} of the LC mixture employed. For a $5\text{-}\mu\text{m}$ cell gap, the optical response time is ~ 30 to 40 ms. At $V = 5 V_{rms}$, the contrast ratio (CR) reaches $\sim 500:1$. These performances, although not perfect, are acceptable for notebook computers. A major drawback of the TN cell is its narrow viewing angle and grayscale inversion originated from the LC directors tilting out of the plane. Because of this molecular tilt, the viewing angle in the vertical direction is narrow and asymmetric, and has grayscale inversion.⁵⁴ Despite a relatively slow switching time and limited viewing angle, the TN cell is still widely used for many applications because of its simplicity and low cost. Recently, a fast-response (~ 2 ms gray-to-gray response time) TN notebook computer has been demonstrated by using a thin cell gap ($2\text{ }\mu\text{m}$), low viscosity LC mixture, and overdrive and undershoot voltage method.

In-Plane Switching (IPS) Cell

In an IPS cell, the transmission axis of the polarizer is parallel to the LC director at the input plane.⁵⁵ The optical wave traversing through the LC cell is an extraordinary wave whose polarization state remains unchanged. As a result, a good dark state is achieved since this linearly polarized light is completely absorbed by the crossed analyzer. When an electric field is applied to the LC cell, the LC directors are reoriented toward the electric field (along the y axis). This leads to a new director distribution with a twist $f(z)$ in the xy plane as Fig. 23 shows.

Figure 24 depicts the voltage-dependent light transmittance of an IPS cell. The LC employed is MLC-6686, whose $\Delta\epsilon = 10$ and $\Delta n = 0.095$, the electrode width is $4\text{ }\mu\text{m}$, electrode gap $8\text{ }\mu\text{m}$, and cell

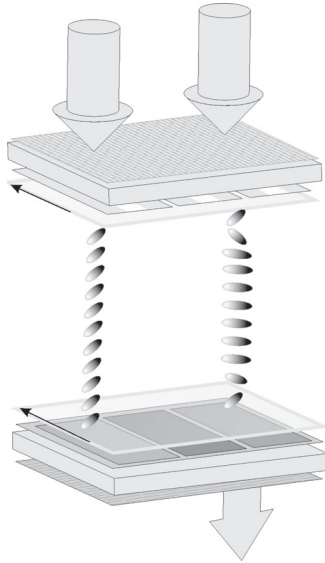


FIGURE 23 IPS mode LC display.

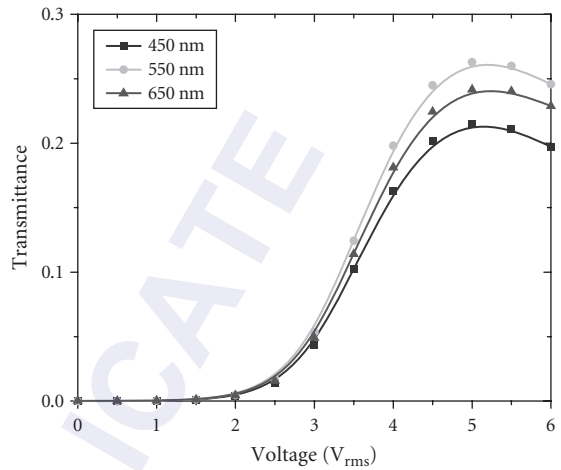


FIGURE 24 Voltage-dependent light transmittance of the IPS cell. LC: MLC-6686, $\Delta\epsilon = 10$, electrode width = $4 \mu\text{m}$, gap = $8 \mu\text{m}$, and cell gap $d = 3.6 \mu\text{m}$.

gap $d = 3.6 \mu\text{m}$. The threshold voltage occurs at $V_{th} \sim 1.5 V_{rms}$ and maximum transmittance at $\sim 5 V_{rms}$ for both wavelengths. Due to absorption, the maximum transmittance of the two polarizers (without LC cell) is 35.4, 33.7, and 31.4 percent for RGB wavelengths, respectively.

Vertical Alignment (VA) Cell

The vertically aligned LC cell exhibits an unprecedentedly high contrast ratio among all the LC modes developed.⁵⁶ Moreover, its contrast ratio is insensitive to the incident light wavelength, cell thickness, and operating temperature. In principle, in voltage-OFF state the LC directors are perpendicular to the cell substrates, as Fig. 25 shows. Thus, its contrast ratio at normal angle is limited by the crossed polarizers. Application of a voltage to the ITO electrodes causes the directors to tilt away from the normal to the glass surfaces. This introduces birefringence and subsequently light transmittance because the refractive indices for the light polarized parallel and perpendicular to the directors are different.

Figure 26 shows the voltage-dependent transmittance of a VA cell with $d\Delta n = 350 \text{ nm}$ between two crossed polarizers. Here, a single domain VA cell employing Merck high resistivity MLC-6608 LC mixture is simulated. Some physical properties of MLC-6608 are listed as follows: $n_e = 1.562$, $n_o = 1.479$ (at $\lambda = 546 \text{ nm}$ and $T = 20^\circ\text{C}$); clearing temperature $T_c = 90^\circ\text{C}$; $\Delta\epsilon = -4.2$, and rotational viscosity $\gamma_1 = 186 \text{ mPas}$ at 20°C .

From Fig. 26, an excellent dark state is obtained at normal incidence. As the applied voltage exceeds the Freederickz threshold voltage ($V_{th} \sim 2.1 V_{rms}$), LC directors are reoriented by the applied electric field resulting in light transmission from the crossed analyzer. From the figure, RGB wavelengths reach their peak at different voltages, blue at $\sim 4 V_{rms}$, and green at $\sim 6 V_{rms}$. The on-state dispersion is more forgiven than the dark state. A small light leakage in the dark state would degrade the contrast ratio significantly, but less noticeable from the bright state.

It should be mentioned that in Figs. 25 and 26 only a single domain is considered, thus, its viewing angle is quite narrow. To achieve wide view, four domains with film compensation are required. Several approaches, such as Fujitsu's MVA (multidomain VA) and Samsung's PVA (Patterned VA),

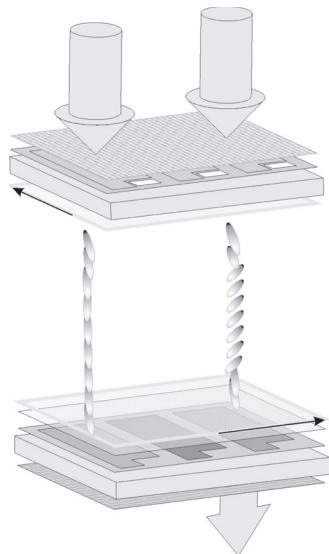


FIGURE 25 VA mode LC display.

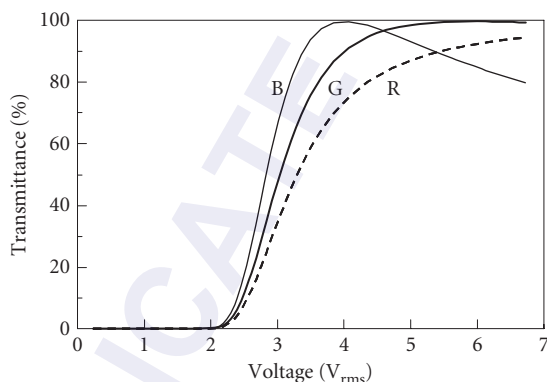


FIGURE 26 Voltage-dependent normalized transmittance of a VA cell. LC: MLC-6608. $d\Delta n = 350$ nm. $R = 650$ nm, $G = 550$ nm, and $B = 450$ nm.

have been developed for obtaining four complementary domains. Figure 27 depicts the PVA structure developed by Samsung.

As shown in Fig. 27, PVA has no pretilt angle. The four domains are induced by the fringe electric fields. On the other hand, in Fujitsu's MVA, physical protrusions are used to provide initial pretilt direction for forming four domains. The protrusions not only reduce the aperture ratio but also cause light leakage in the dark state because the LCs on the edges of protrusions are tilted so that they exhibit birefringence. It would be desirable to eliminate protrusions for MVA and create a pretilt angle in each domain for both MVA and PVA to guide the LC reorientation direction. Based on this concept, surface polymer sustained alignment (PSA) technique has been developed.⁵⁷ A very small percentage (~ 0.2 wt %) of reactive mesogen monomer and photoinitiator are mixed in a negative $\Delta\epsilon$ LC host and injected into a LCD panel. While a voltage is applied to generate four domains, a UV light is used to cure the monomers. As a result, the monomers are adsorbed onto the surfaces. These cured polymers, although in low density, will provide a pretilt angle within each domain to guide the LC reorientation. Thus, the rise time is reduced by nearly $2\times$ while the decay time remains more or less unchanged.⁵⁸

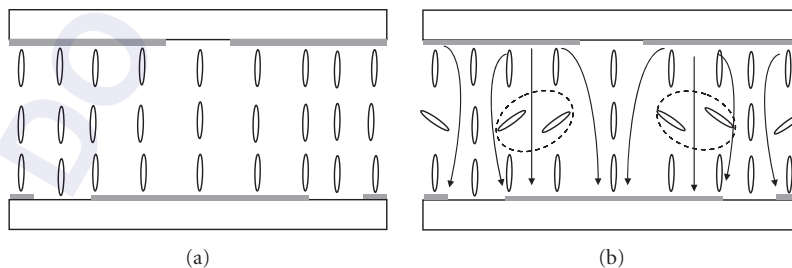


FIGURE 27 (a) LC directors of PVA at $V = 0$ and (b) LC directors of PVA at a voltage-on state. The fringe fields generated by the top and bottom slits create two opposite domains in this cross section. When the zigzag electrodes are used, four domains are generated.

8.7 LIQUID CRYSTALS DISPLAYS

The most common and well recognized applications of liquid crystals nowadays are displays. It is the most natural way to utilize extraordinary electro-optical properties of liquid crystals together with its liquid-like behavior. All the other applications are called as nondisplay applications of liquid crystals. Nondisplay applications are based on the liquid crystals molecular order sensitivity to the external incentive. This can be an external electric and magnetic field, temperature, chemical agents, mechanical stress, pressure, irradiation by different electromagnetic wave, or radioactive agents. Liquid crystals sensitivity for such wide spectrum of factors results in tremendous diversity of nondisplay applications. It starts from spatial light modulators for laser beam steering, adaptive optics through telecommunication area (light shutters, attenuators, and switches), cholesteric LC filters, LC thermometers, stress meters, dose meters to ends up with liquid crystals paints, and cosmetics. Another field of interest employing lyotropic liquid crystals is biomedicine where it plays an important role as a basic unit of the living organisms by means of plasma membranes of the living cells. More about existing nondisplays applications can be found in preferred reading materials.

Three types of liquid crystal displays (LCDs) have been developed: (1) transmissive, (2) reflective, and (3) transfective. Each one has its own unique properties. In the following sections, we will introduce the basic operation principles of these display devices.

Transmissive TFT LCDs

A transmissive LCD uses a backlight to illuminate the LCD panel to achieve high brightness (300 to 500 nits) and high contrast ratio ($>2000:1$). Some transmissive LCDs, such as twisted-nematic (TN), do not use phase compensation films or multidomain structures so that their viewing angle is limited and they are more suitable for single viewer applications, for example, mobile displays and notebook computers. With phase compensation films and multidomain structures, the direct-view transmissive LCDs exhibit a wide viewing angle and high contrast ratio, and have been widely used for desktop computers and televisions. However, the cost of direct-view large screen LCDs is still relatively expensive. To obtain a screen diagonal larger than 2.5 m, projection displays, such as data projector, using transmissive microdisplays are still a favorable choice. There, a high power arc lamp or light emitting diode (LED) arrays are used as light source. Using a projection lens, the displayed image is magnified by more than 50 \times . To reduce the size of optics and cost, the LCD panel is usually made small (<25 mm in diagonal) and each pixel size is ~ 20 to 40 μm . Thus, poly-silicon-based TFT LCD is the common choice.

Figure 28 shows the device structure of a transmissive thin-film-transistor (TFT) LCD using amorphous silicon (a-Si) transistors. LCD is a nonemissive display, that is, it does not emit light, instead, it functions as a two-dimensional spatial light modulator. Thus, a backlight is needed. A diffuser is used to homogenize the backlight in order to avoid hot spots. Some optical films are stacked to steer the Lambertian backlight to the central $\pm 40^\circ$ for improving display brightness. Since most LCDs require a linearly polarized light in order to obtain a high contrast ratio, two sheets of stretched dichroic polarizers are used for large size direct-view displays. The first glass substrate contains TFT arrays, which serve as independent light switches. Each display pixel is controlled by a TFT. Since TFT is sensitive and should be shielded from backlight illumination, the actual aperture ratio (the transparent ITO electrode area) is reduced to ~ 80 percent, depending on the pixel density. As the pixel density increases, the aperture ratio decreases. The LC layer is sandwiched between two ITO substrates whose inner surface is coated with a thin (80 to 100 nm) polyimide layer. Some LCDs (twisted-nematic, in-plane switching and fringe field switching) require rubbing but some (multidomain vertical alignment and patterned vertical alignment) do not. The cell gap is usually controlled at ~ 3.5 to 4.0 μm for a transmissive LCD. The performance of the display such as light throughput, response time, and viewing angle are all influenced by the LC configuration employed.

For direct-view LCDs, compact size, lightweight, and low power consumption are equally important as viewing angle, color, and contrast ratio. For direct-view LCDs, color filters are embedded in the inner side of the top (second) substrate. Three subpixels (red, green, and blue) form a color pixel. The size of

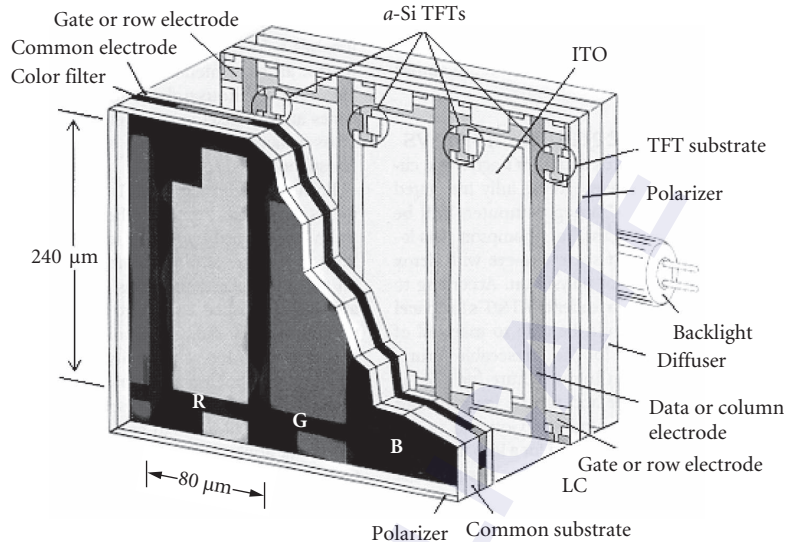


FIGURE 28 Device structure of a color pixel of a transmissive TFT LCD.

each subpixel is $\sim 80 \mu\text{m} \times 240 \mu\text{m}$. Each subpixel transmits only one color; the other colors are absorbed. Figure 29 depicts the emission spectra of a backlight (cold cathode fluorescent lamp, CCFL) and RGB light emitting diodes (LEDs), and the transmission spectra of RGB color filters.

From Fig. 29, the transmission spectra of RGB color filters are relatively broad. The advantage is to transmit more light however its color purity is degraded. The peak transmission of the RGB color filters is ~ 70 , 80 , and 90 percent, respectively. Roughly speaking, each color filter only transmits ~ 25 percent of the incident white light. The rest ~ 75 percent is absorbed by the color pigments. Moreover, CCFL emits two unwanted lines: blue-green ($\sim 480 \text{ nm}$) and orange ($\sim 580 \text{ nm}$). The blue-green light will transmit through the blue and green color filters simultaneously. Similarly, the orange light will transmit through the green and red color filters simultaneously. These leaked lights will downgrade the color purity (or color saturation) of the display. Therefore, the color gamut of a typical transmissive TFT LCD is ~ 75 percent of NTSC (National Television System Committee) standard. With improved CCFL spectra, the color gamut can reach ~ 92 percent. LEDs have narrower emission spectra that also match

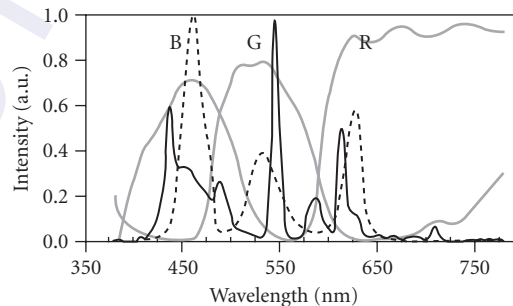


FIGURE 29 Transmission spectra of RGB color filters (thick gray lines), and emission spectra of CCFL backlight (thin black lines) and LEDs (dashed lines).

better with the transmission spectra of the color filters, thus the color gamut reaches ~ 120 percent.⁵⁹ For a display device, a wider color gamut is not necessarily better; natural colors are also important. After all, display is an art where perception plays an important role.

After taking into account the optical losses from polarizers, color filters, and TFT aperture ratio, the overall system optical efficiency is about 6 to 7 percent for a direct-view LCD. If wide view technology is included,⁶⁰ the total light efficiency is decreased to ~ 5 percent. Low optical efficiency implies high power consumption and more heat generation inside the display chassis. For a thin LCD, thermal dissipation is a critical issue. For portable displays, low power consumption is desirable because it lengthens the battery operating hours. Several approaches have been developed to reduce power consumption, for example, polarization conversion of backlight⁶¹ and two-dimensional LED backlight with local dimming capability.⁶²⁻⁶⁴ The use of LED backlight offers several additional advantages such as wide color gamut, high dynamic contrast ratio ($>50,000:1$), $\sim 2\times$ reduction in power consumption, and fast turn-on and off times (~ 10 ns) for reducing the motion picture image blurs.⁶⁵ Some technological concerns are color and power drifting as the junction temperature changes, and cost.

Reflective LCDs

Figure 30 shows a device structure of a TFT-based reflective LCD. The top linear polarizer and a broadband quarter-wave film forms an equivalent crossed polarizer for the incident and exit beams. This is because the LC modes work better under crossed-polarizer condition. The bumpy reflector not only reflects but also diffuses the ambient light to the observer in order to avoid specular reflection and widen the viewing angle. This is a critical part for reflective LCDs. The TFT is hidden beneath the bumpy reflector, thus the R-LCD can have a large aperture ratio ($\sim 90\%$). The light blocking layer (LBL) is used to absorb the scattered light from neighboring pixels. Two popular LCD modes have been widely used for R-LCDs: (1) VA cell, and (2) mixed-mode twisted nematic (MTN) cell. The VA cell utilizes the phase retardation effect while the MTN cell uses the combination of polarization rotation and birefringence effects.

In a reflective LCD, there is no built-in backlight unit; instead, it utilizes ambient light for reading out the displayed images. In comparison to transmissive LCDs, reflective LCDs have advantages in lower power consumption, lighter weight, and better sunlight readability. However, a reflective LCD is inapplicable under low or dark ambient conditions. Therefore, the TFT-based reflective LCD is gradually losing its ground.

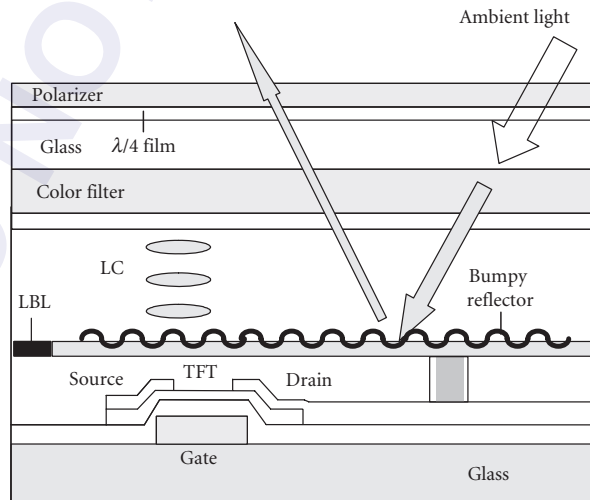


FIGURE 30 Device structure of a direct-view reflective LCD.

Flexible reflective LCDs using cholesteric liquid crystal display (Ch-LCD) and bistable nematic are gaining momentum because they can be used as electronic papers. These reflective LCDs use ambient light to readout the displayed images. Ch-LCD has helical structure which reflects color so that the display does not require color filters, neither polarizer. Thus, the reflectance for a given color band which depends on the pitch length and refractive index of the employed LC is relatively high (~30%). Moreover, it does not require a backlight so that its weight is light and the total device thickness can be thinner than 200 μm . Therefore, it is a strong contender for color flexible displays. Ch-LCD is a bistable device so that its power consumption is low, provided that the device is not refreshed too frequently. A major drawback of a reflective direct-view LCD is its poor readability under low ambient light.

Another reflective LCD developed for projection TVs is liquid-crystal-on-silicon (LCoS) microdisplay. Unlike a transmissive microdisplay, LCoS is a reflective device. Here the reflector employed is an aluminum metallic mirror. Crystalline silicon has high mobility so that the pixel size can be made small (<10 μm) and aperture ratio >90 percent. Therefore, the image not only has high resolution but also is seamless. By contrast, a transmissive microdisplay's aperture ratio is about 65 percent. The light blocked by the black matrices show up in the screen as dark patterns (also known as screen door effect). Viewing angle of a LCD is less critical in projection than direct-view displays because in a projection display the polarizing beam splitter has a narrower acceptance angle than the employed LCD.

Transflective LCDs

In a transflective liquid crystal display (TR-LCD), a pixel is divided into two parts: transmissive (T) and reflective (R). The T/R area ratio can vary from 80/20 to 20/80, depending on the applications. In dark ambient the backlight is on and the display works as a transmissive one, while at bright ambient the backlight is off and only the reflective mode is operational.

Dual-Cell-Gap Transflective LCDs In a TR-LCD, the backlight traverses the LC layer once, but the ambient light passes through twice. As a result, the optical path length is unequal. To balance the optical path difference between the T and R regions for a TR-LCD, dual-cell-gap device concept is introduced. The basic requirement for a TR-LCD is to find equal phase retardation between the T and R modes, which is

$$d_T(\Delta n)_T = 2d_R(\Delta n)_R \quad (22)$$

If T and R modes have the same effective birefringence, then the cell gap should be different. This is the so-called dual cell gap approach. On the other hand, if the cell gap is uniform (single cell gap approach), then we should find ways to make $(\Delta n)_T = 2(\Delta n)_R$. Let us discuss the dual cell gap approaches first.

Figure 31a shows the schematic device configuration of a dual-cell-gap TR-LCD. Each pixel is divided into a reflective region with cell gap d_R and a transmissive region with cell gap d_T . The LC employed could be homogeneous alignment (also known as ECB, electrically controlled birefringence) or vertical alignment, as long as it is a phase retardation type. To balance the phase retardation between the single and double pass of the T and R parts, we could set $d_T = 2d_R$. Moreover, to balance the color saturation due to single and double-pass discrepancy, we could use thinner or holed color filters in the R part. The top quarter wave plate is needed mainly for the reflective mode to obtain a high contrast ratio. Therefore, in the T region, the optic axis of the bottom quarter-wave plate should be aligned perpendicular to that of the top one so that their phase retardations are canceled.

A thin homogeneous cell is difficult to find a good common dark state for RGB wavelengths without a compensation film.⁵⁵ The compensation film can be designed into the top quarter-wave film shown in Fig. 31a to form a single film. Here, let us take a dual cell gap TR-LCD using VA (or MVA for wide-view) and MLC-6608 ($\Delta\epsilon = -4.2$, $\Delta n = 0.083$) as an example. We set $d_R = 2.25 \mu\text{m}$ in the R region and $d_T = 4.5 \mu\text{m}$ in the T region. Figure 31b depicts the voltage-dependent transmittance (VT) and reflectance (VR) curves at normal incidence. As expected, both VT and VR curves perfectly overlap with each other. Here $d_R\Delta n = 186.8 \text{ nm}$ and $d_T\Delta n = 373.5 \text{ nm}$ are intentionally designed to be larger than $\lambda/4$ (137.5 nm) and $\lambda/2$ (275 nm), respectively, in order to reduce the on-state voltage to $\sim 5 V_{\text{rms}}$.

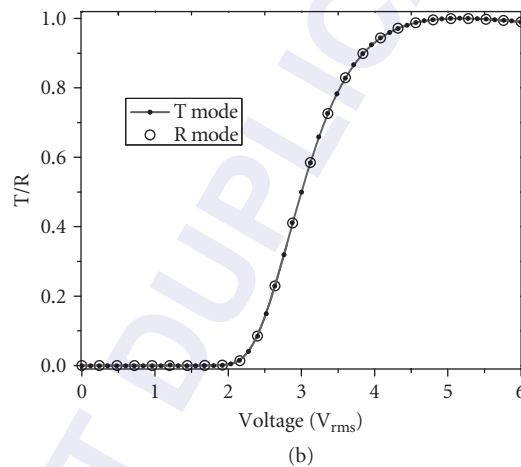
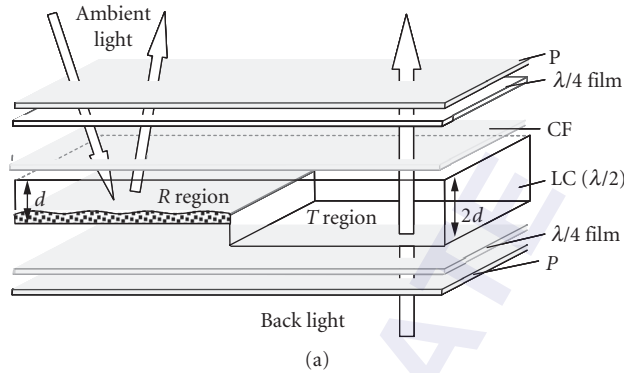


FIGURE 31 (a) Schematic device configuration of the dual-cell-gap TR-LCD. (b) Simulated VT and VR curves using VA (or MVA) cells. LC: MLC-6608, $d_T = 4.5 \mu\text{m}$, $d_R = 2.25 \mu\text{m}$, and $\lambda = 550 \text{ nm}$.

Three problems of dual-cell-gap TR-LCDs are encountered: (1) Due to the cell gap difference the LC alignment is distorted near the T and R boundaries. The distorted LCs will cause light scattering and degrade the device contrast ratio. Therefore, these regions should be covered by black matrices in order to retain a good contrast ratio. (2) The thicker cell gap in the T region results in a slower response time than the R region. Fortunately, the dynamic response requirement in mobile displays is not as strict as those for video applications. This response time difference, although not perfect, is still tolerable. (3) The view angle of the single-domain homogeneous cell mode is relatively narrow because the LC directors are tilted out of the plane by the longitudinal electric field. To improve view angle, a biaxial film⁶⁶ or a hybrid aligned nematic polymeric film⁶⁷ is needed. Because the manufacturing process is compatible with the LCD fabrication lines, the dual-cell-gap TR-LCD is widely used in commercial products, such as iPhones.

Single-Cell-Gap Transflective LCDs As its name implies, the single-cell-gap TR-LCD has a uniform cell gap in the T and R regions. Therefore, we need to find device concepts to achieve $(\Delta n)_T = 2(\Delta n)_R$. Several approaches have been proposed to solve this problem. In this section, we will discuss two

examples: (1) dual-TFT method in which one TFT is used to drive the T mode and another TFT to drive the R mode at a lower voltage,⁶⁸ and (2) divided-voltage method:⁶⁹ to have multiple R parts and the superimposed VR curve matches the VT curve.

Example 3: Dual-TFT Method Figure 32a shows the device structure of a TR-LCD using two TFTs to separately control the gamma curves of the T and R parts. Here, TFT-1 is connected to the bumpy reflector and TFT-2 is connected to the ITO of the transmissive part. Because of the double passes, the VR curve has a sharper slope than the VT curve and it reaches the peak reflectance at a lower voltage, as shown in Fig. 32b. Let us use a 4.5- μm vertically aligned LC layer with 88° pretilt angle as an example. The LC mixture employed is Merck MLC-6608 and wavelength is $\lambda = 550$ nm. From Fig. 32b, the peak reflectance occurs at $3 V_{\text{rms}}$ and transmittance at $5.5 V_{\text{rms}}$. Thus, the maximum voltage of TFT-1 should be set at 5.5 V and TFT-2 at 3 V. This driving scheme is also called double-gamma method.⁷⁰ The major advantage of this dual-TFT approach is its simplicity. However, each TFT takes up some real estate so that the aperture ratio for the T mode is reduced.

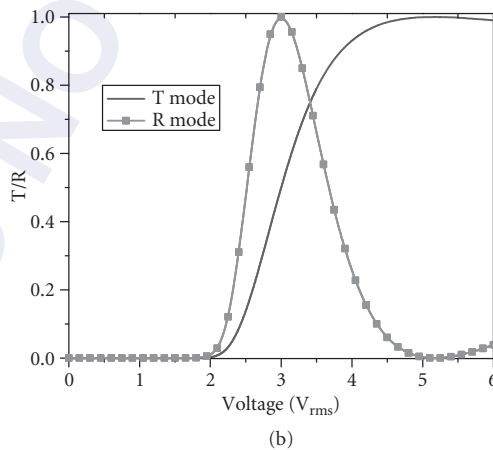
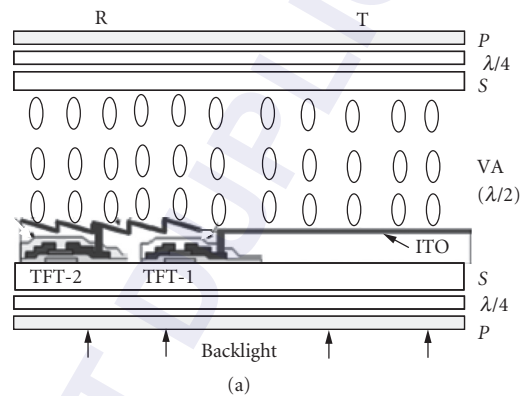


FIGURE 32 (a) Device structure of a dual-TFT TR-LCD and (b) simulated VT and VR curves using a 4.5- μm , MLC-6608 LC layer.

For a TR-LCD, the T mode should have priority over R mode. The major function of R mode is to preserve sunlight readability. In general, the viewing angle, color saturation, and contrast ratio of R mode are all inferior to T mode. In most lighting conditions except under direct sunlight, T mode is still the primary display.

Example 4: Divided-Voltage Method Figure 33a shows the device structure of a TR-LCD using divided voltage method.⁷¹ The R region consists of two sub-regions: R-I and R-II. Between R-II and bottom ITO, there is a passivation layer to weaken the electric field in the R-II region. As plotted in Fig. 33b, the VR-II curve in the R-II region has a higher threshold voltage than the VT curve due to this voltage shielding effect. To better match the VT curve, a small area in the T region is also used for R-I. The bumpy reflector in the R-I region is connected to the bottom ITO through a channeled electrode. Because of the double passes of ambient light, the VR-I curve is sharper than the VT curve. By properly choosing the R-I and R-II areas, we can match the VT and VR curves well, as shown in Fig. 33b.

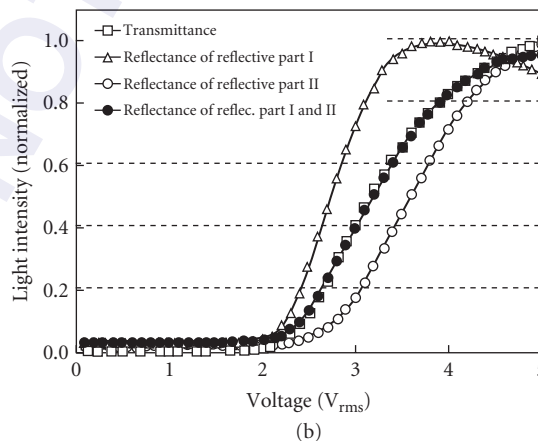
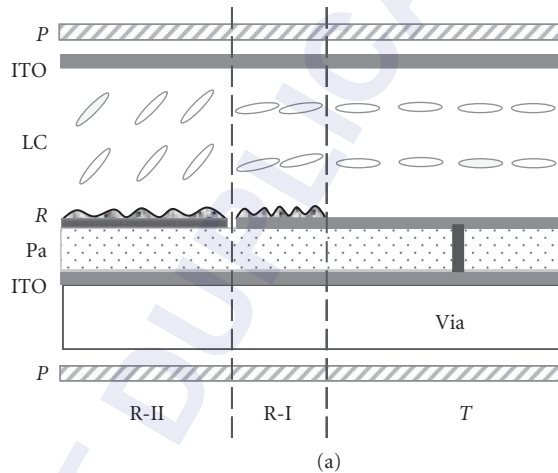


FIGURE 33 A transreflective LCD using divided-voltage approach. (a) Device structure and (b) VT and VR curves at different regions. Here, *P*: polarizer, *R*: bumpy reflector, and *Pa*: passivation layer. (Redrawn from Ref. 71.)

8.8 POLYMER/LIQUID CRYSTAL COMPOSITES

Some types of liquid crystal display combine LC material and polymer material in a single device. Such LC/polymer composites are relatively new class of materials used for displays but also for light shutters or switchable windows. Typically LC/polymer composites consist of calamitic low mass LCs and polymers and can be either polymer-dispersed liquid crystal (PDLC) or polymer-stabilized liquid crystals (PSLC).⁷²⁻⁷⁴ The basic difference between these two types comes from concentration ratio between LC and polymer. In case of PDLC there is typically around 1:1 percentage ratio of LC and polymer. In PSLC, the LC occupies 90 percent or more of the total composition. Such difference results in different phase separation process during composites polymerization. For equal concentration of LC and polymer, the LC droplets will form. But in case the LC is a majority, polymer will build up only walls or strings which divide LC into randomly aligned domains. Both types of composites operate between transparent state and scattering state. There are two requirements on the polymer for PDLC or PSLC device to work. First, refraction index of the polymer, n_p , must be equal to the refraction index for light polarized perpendicular to the director of the liquid crystal, (ordinary refractive index of the LC). Second, the polymer must induce the director of the LC in the droplets (PDLC) or domains (PSLC) to orient parallel to the surface of the polymer (Fig. 34). In the voltage OFF state the LC molecules in the droplets are partially aligned. In addition, the average director orientation \mathbf{n} of the droplets exhibits a random distribution of orientation within the cell.

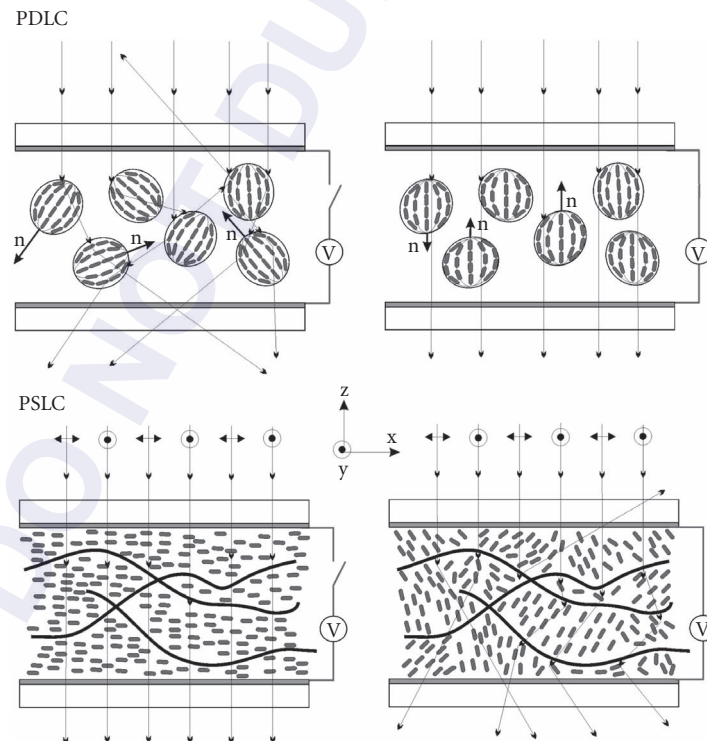


FIGURE 34 Schematic view and working principles of polymer/LC composites.

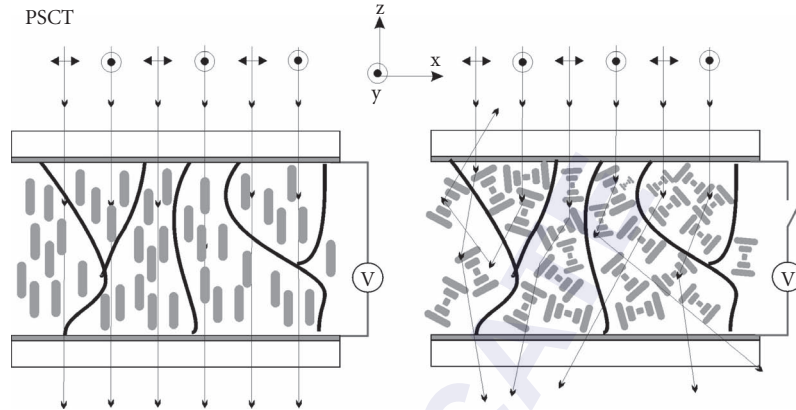


FIGURE 35 Schematic view of working principles of polymer-stabilized cholesteric LC light valve.

The incident unpolarized light is scattered if it goes through such a sample. When a sufficiently strong electric field (typically above $1 V_{rms}/\mu m$) is applied to the cell, all the LC molecules align parallel to the electric field. If the light is also propagating in the direction parallel to the field, then the beam of light is affected by ordinary refractive index of LC which is matched with refractive index of polymer, thus cell appears transparent. When the electric field is OFF, again the LC molecules go back to the previous random positions. A polymer mixed with a chiral liquid crystal is a special case of PSLC called as polymer stabilized cholesteric texture (PSCT). The ratio between polymer and liquid crystal remains similar to the one necessary for PSLC. When the voltage is not applied to the PSCT cell liquid crystals tends to have helical structure while the polymer network tends to keep LC director parallel to it (normal-mode PSCT). Therefore, the material has a poly-domain structure, as Fig. 35 shows. In this state the incident beam is scattered. When a sufficiently high electric field is applied across the cell, the liquid crystal is switched to the homeotropic alignment and, as a result, the cell becomes transparent.

8.9 SUMMARY

Liquid crystal was discovered more than 100 years ago and is finding widespread applications. This class of organic material exhibits some unique properties, such as good chemical and thermal stabilities, low operation voltage and low power consumption, and excellent compatibility with semiconductor fabrication processing. Therefore, it has dominated direct-view and projection display markets. The forecasted annual TFT LCD market is going to exceed \$100 billion by year 2011. However, there are still some technical challenges need to be overcome, for example, (1) faster response time for reducing motion picture blurs, (2) higher optical efficiency for reducing power consumption and lengthening battery life, (3) smaller color shift when viewed at oblique angles, (4) wider viewing angle with higher contrast ratio (ideally its viewing characteristics should be as good as an emissive display), and (5) lower manufacturing cost, especially for large screen TVs. In addition to displays, LC materials also share an important part of emerging photonics applications, such as spatial light modulators for laser beam steering and adaptive optics, adaptive-focus lens, variable optical attenuator for fiber-optic telecommunications, and LC-infiltrated photonic crystal fibers, just to name a few. Often neglected, lyotropic liquid crystals are important materials in biochemistry of the cell membranes.

8.10 REFERENCES

1. F. Reinitzer, *Monatsh. Chem.*, **9**:421 (1888); for English translation see, *Liq. Cryst.* **5**:7 (1989).
2. P. F. McManamon, T. A. Dorschner, D. L. Corkum, L. Friedman, D. S. Hobbs, M. Holz, S. Liberman, et al., *Proc. IEEE* **84**:268 (1996).
3. V. Vill, *Database of Liquid Crystalline Compounds for Personal Computers*, Ver. 4.6 (LCI Publishers, Hamburg 2005).
4. C. S. O'Hern and T. C. Lubensky, *Phys. Rev. Lett.* **80**:4345 (1998).
5. G. Friedel, *Ann. Physique* **18**:173 (1922).
6. P. G. De Gennes and J. Prost, *The Physics of Liquid Crystals*, 2nd ed. (Clarendon, Oxford, 1993).
7. P. R. Gerber, *Mol. Cryst. Liq. Cryst.* **116**:197 (1985).
8. N. A. Clark and S. T. Lagerwall, *Appl. Phys. Lett.* **36**:889 (1980).
9. A. D. L. Chandani, T. Hagiwara, Y. Suzuki, Y. Ouchi, H. Takezoe, and A. Fukuda, *Jpn. J. Appl. Phys.* **27**:L1265 (1988).
10. A. D. L. Chandani, E. Górecka, Y. Ouchi, H. Takezoe, and A. Fukuda, *Jpn. J. Appl. Phys.* **27**:L729 (1989).
11. J. W. Goodby, M. A. Waugh, S. M. Stein, E. Chin, R. Pindak, and J. S. Patel, *Nature* **337**:449 (1989).
12. T. C. Lubensky and S. R. Renn, *Mol. Cryst. Liq. Cryst.* **209**:349 (1991).
13. L. Schröder, *Z. Phys. Chem.* **11**:449 (1893).
14. J. J. Van Laar, *Z. Phys. Chem.* **63**:216 (1908).
15. W. Maier, G. Meier, and Z. Naturforsch. *Teil. A* **16**:262 (1961).
16. M. Schadt, *Displays* **13**:11 (1992).
17. G. Gray, K. J. Harrison, and J. A. Nash, *Electron. Lett.* **9**:130 (1973).
18. R. Dabrowski, *Mol. Cryst. Liq. Cryst.* **191**:17 (1990).
19. M. Schadt and W. Helfrich, *Appl. Phys. Lett.* **18**:127 (1971).
20. R. A. Soref, *Appl. Phys. Lett.* **22**:165 (1973).
21. M. Oh-e and K. Kondo, *Appl. Phys. Lett.* **67**:3895 (1995).
22. Y. Nakazono, H. Ichinose, A. Sawada, S. Naemura, and K. Tarumi, *Int'l Display Research Conference*, p. 65 (1997).
23. R. Tarao, H. Saito, S. Sawada, and Y. Goto, *SID Tech. Digest* **25**:233 (1994).
24. T. Geelhaar, K. Tarumi, and H. Hirschmann, *SID Tech. Digest* **27**:167 (1996).
25. Y. Goto, T. Ogawa, S. Sawada and S. Sugimori, *Mol. Cryst. Liq. Cryst.* **209**:1 (1991).
26. M. F. Schiekel and K. Fahrenschon, *Appl. Phys. Lett.* **19**:391 (1971).
27. Q. Hong, T. X. Wu, X. Zhu, R. Lu, and S.T. Wu, *Appl. Phys. Lett.* **86**:121107 (2005).
28. C. H. Wen, S. Gauza, and S.T. Wu, *Appl. Phys. Lett.* **87**:191909 (2005).
29. R. Lu, Q. Hong, and S.T. Wu, *J. Display Technology* **2**:217 (2006).
30. R. Eidenschink and L. Pohl, *US Patent* 4, 415, 470 (1983).
31. W. H. de Jeu, "The Dielectric Permittivity of Liquid Crystals" *Solid State Phys. Suppl.* **14**: "Liquid Crystals" Edited by L. Liebert. (Academic Press, New York, 1978); also, *Mol. Cryst. Liq. Cryst.* **63**:83 (1981).
32. M. Schadt, *Mol. Cryst. Liq. Cryst.* **89**:77 (1982).
33. H. K. Bucher, R. T. Klingbiel, and J. P. VanMeter, *Appl. Phys. Lett.* **25**:186 (1974).
34. H. Xianyu, Y. Zhao, S. Gauza, X. Liang, and S. T. Wu, *Liq. Cryst.* **35**:1129 (2008).
35. T. K. Bose, B. Campbell, and S. Yagihara, *Phys. Rev. A* **36**:5767 (1987).
36. C. H. Wen and S. T. Wu, *Appl. Phys. Lett.* **86**:231104 (2005).
37. I. C. Khoo and S. T. Wu, *Optics and Nonlinear Optics of Liquid Crystals* (World Scientific, Singapore, 1993).
38. S. T. Wu, E. Ramos, and U. Finkenzeller, *J. Appl. Phys.* **68**:78 (1990).
39. S. T. Wu, U. Efron, and L. D. Hess, *Appl. Opt.* **23**:3911 (1984).
40. S. T. Wu, J. *Appl. Phys.* **69**:2080 (1991).
41. S. T. Wu, C. S. Wu, M. Warengem, and M. Ismaili, *Opt. Eng.* **32**:1775 (1993).

41. J. Li and S. T. Wu, *J. Appl. Phys.* **95**:896 (2004).
42. S. T. Wu, U. Efron and L. D. Hess, *Appl. Phys. Lett.* **44**:1033 (1984).
43. J. Li and S. T. Wu, *J. Appl. Phys.* **96**:170 (2004).
44. H. Mada and S. Kobayashi, *Mol. Cryst. Liq. Cryst.* **33**:47 (1976).
45. E. H. Stupp and M. S. Brennessoltz, *Projection Displays* (Wiley, New York, 1998).
46. J. Li, S. Gauza, and S. T. Wu, *Opt. Express* **12**:2002 (2004).
47. J. Li and S. T. Wu, *J. Appl. Phys.* **96**:19 (2004).
48. W. Maier and A. Saupe, *Z. Naturforsch. Teil A* **15**:287 (1960).
49. H. Gruler, and Z. Naturforsch. *Teil A* **30**:230 (1975).
50. W. M. Gelbart and A. Ben-Shaul, *J. Chem. Phys.* **77**:916 (1982).
51. S. T. Wu and C. S. Wu, *Liq. Cryst.* **8**:171 (1990). Seven commonly used models (see the references therein) have been compared in this paper.
52. M. A. Osipov and E. M. Terentjev, and Z. Naturforsch. *Teil A* **44**:785 (1989).
53. S. T. Wu and C. S. Wu, *Phys. Rev. A* **42**:2219 (1990).
54. S. T. Wu and C. S. Wu, *J. Appl. Phys.* **83**:4096 (1998).
55. M. Oh-e and K. Kondo, *Appl. Phys. Lett.* **67**:3895 (1995).
56. M.F. Schiekel and K. Fahrenschoen, *Appl. Phys. Lett.* **19**:391 (1971).
57. K. Hanaoka, Y. Nakanishi, Y. Inoue, S. Tanuma, Y. Koike, and K. Okamoto, *SID Tech. Digest* **35**:1200 (2004).
58. S. G. Kim, S. M. Kim, Y. S. Kim, H. K. Lee, S. H. Lee, G. D. Lee, J. J. Lyu, and K. H. Kim, *Appl. Phys. Lett.* **90**:261910 (2007).
59. R. Lu, Q. Hong, Z. Ge, and S. T. Wu, *Opt. Express* **14**:6243 (2006).
60. D. K. Yang and S. T. Wu, *Fundamentals of Liquid Crystal Devices* (Wiley, New York, 2006).
61. J. M. Jonza, M. F. Weber, A. J. Ouderkirk, and C. A. Stover, *U.S. Patent* 5, 962,114 (1999).
62. P. de Greef and H. G. Hulze, *SID Symp. Digest* **38**:1332 (2007).
63. H. Chen, J. Sung, T. Ha, and Y. Park, *SID Symp. Digest* **38**:1339 (2007).
64. F. C. Lin, C. Y. Liao, L. Y. Liao, Y. P. Huang, and H. P. Shieh, *SID Symp. Digest* **38**:1343 (2007).
65. M. Anandan, *J. SID* **16**:287 (2008).
66. M. Shibazaki, Y. Ukawa, S. Takahashi, Y. Iefuji, and T. Nakagawa, *SID Tech. Digest* **34**:90 (2003).
67. T. Uesaka, S. Ikeda, S. Nishimura, and H. Mazaki, *SID Tech. Digest* **28**:1555 (2007).
68. K. H. Liu, C. Y. Cheng, Y. R. Shen, C. M. Lai, C. R. Sheu, and Y. Y. Fan, *Proc. Int. Display Manuf. Conf.* p. 215 (2003).
69. C. Y. Tsai, M. J. Su, C. H. Lin, S. C. Hsu, C. Y. Chen, Y. R. Chen, Y. L. Tsai, C. M. Chen, C. M. Chang, and A. Lien, *Proc. Asia Display* **24** (2007).
70. C. R. Sheul, K. H. Liu, L. P. Hsin, Y. Y. Fan, I. J. Lin, C. C. Chen, B. C. Chang, C. Y. Chen, and Y. R. Shen, *SID Tech. Digest* **34**:653 (2003).
71. Y. C. Yang, J. Y. Choi, J. Kim, M. Han, J. Chang, J. Bae, D. J. Park, et al., *SID Tech. Digest* **37**:829 (2006).
72. J. L. Fergason, *SID Symp. Digest* **16**:68 (1985).
73. J. W. Doane, N. A. Vaz, B. G. Wu, and S. Zumer, *Appl. Phys. Lett.* **48**:269 (1986).
74. R. L. Sutherland, V. P. Tondiglia, and L. V. Natarajan, *Appl. Phys. Lett.* **64**:1074 (1994).

8.11 BIBLIOGRAPHY

1. Chandrasekhar, S., *Liquid Crystals*, 2nd edn. (Cambridge University Press, Cambridge, England, 1992).
2. Collings, P. J., *Nature's Delicate Phase of Matter*, 2nd ed. (Princeton University Press, Princeton, N.J., 2001).
3. Collings, P. J., and M. Hird, *Introduction to Liquid Crystals Chemistry and Physics*, (Taylor & Francis, London 1997).
4. de Gennes, P. G., and J. Prost, *The Physics of Liquid Crystals*, 2nd ed. (Oxford University Press, Oxford, 1995).

5. de Jeu, W. H., *Physical Properties of Liquid Crystalline Materials* (Gorden and Breach, New York, 1980).
6. Demus, D., J. Goodby, G. W. Gray, H.-W. Spiess, and V. Vill, *Handbook of Liquid Crystals Vol. 1–4* (Wiley-VCH, Weinheim, New York, 1998).
7. Khoo, I. C., and S. T. Wu, *Optics and Nonlinear Optics of Liquid Crystals* (World Scientific, Singapore, 1993).
8. Kumar, S., *Liquid Crystals* (Cambridge University Press, Cambridge, England, 2001).
9. Oswald, P. and P. Pieranski, *Nematic and Cholesteric Liquid Crystals: Concepts and Physical Properties Illustrated by Experiments* (Taylor & Francis CRC Press, Boca Raton, FL, 2005).
10. Wu, S. T., and D.K. Yang, *Reflective Liquid Crystal Displays* (Wiley, New York, 2001).
11. Wu, S. T., and D. K. Yang, *Fundamentals of Liquid Crystal Devices* (Wiley, Chichester, England, 2006).
12. Yeh, P., and C. Gu, *Optics of Liquid Crystals* (Wiley, New York, 1999).

PART

4

FIBER OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

OPTICAL FIBER COMMUNICATION TECHNOLOGY AND SYSTEM OVERVIEW

Ira Jacobs

*The Bradley Department of Electrical and Computer Engineering
Virginia Polytechnic Institute and State University
Blacksburg, Virginia*

9.1 INTRODUCTION

Basic elements of an optical fiber communication system include the transmitter [laser or light-emitting diode (LED)], fiber (multimode, single-mode, or dispersion-shifted), and the receiver [positive-intrinsic-negative (PIN) diode, avalanche photodiode (APD) detectors, coherent detectors, optical preamplifiers, receiver electronics]. Receiver sensitivities of digital systems are compared on the basis of the number of photons per bit required to achieve a given bit error probability, and eye degradation and error floor phenomena are described. Laser relative intensity noise and nonlinearities are shown to limit the performance of analog systems. Networking applications of optical amplifiers and wavelength-division multiplexing are considered, and future directions are discussed.

Although the light-guiding property of optical fibers has been known and used for many years, it is only relatively recently that optical fiber communications has become both a possibility and a reality.¹ Following the first prediction in 1966² that fibers might have sufficiently low attenuation for telecommunications, the first low-loss fiber (20 dB/km) was achieved in 1970.³ The first semiconductor laser diode to radiate continuously at room temperature was also achieved in 1970.⁴ The 1970s were a period of intense technology and system development, with the first systems coming into service at the end of the decade. The 1980s saw both the growth of applications (service on the first transatlantic cable in 1988) and continued advances in technology. This evolution continued in the 1990s with the advent of optical amplifiers and with the applications emphasis turning from point-to-point links to optical networks. The beginning of the 21st century has seen extensive fiber-to-the-home deployment as well as continued technology advances.

This chapter provides an overview of the basic technology, systems, and applications of optical fiber communication. It is an update and compression of material presented at a 1994 North Atlantic Treaty Organization (NATO) Summer School.⁵ Although there have been significant advances in technology and applications in subsequent years, the basics have remained essentially the same.

9.2 BASIC TECHNOLOGY

This section considers the basic technology components of an optical fiber communications link, namely the fiber, the transmitter, and the receiver, and discusses the principal parameters that determine communications performance.

Fiber

An optical fiber is a thin filament of glass with a central core having a slightly higher index of refraction than the surrounding cladding. From a physical optics standpoint, light is guided by total internal reflection at the core-cladding boundary. More precisely, the fiber is a dielectric waveguide in which there are a discrete number of propagating modes.⁶ If the core diameter and the index difference are sufficiently small, only a single mode will propagate. The condition for single-mode propagation is that the normalized frequency V be less than 2.405, where

$$V = \frac{2\pi a}{\lambda} \sqrt{n_1^2 - n_2^2} \quad (1)$$

and a is the core radius, λ is the free space wavelength, and n_1 and n_2 are the indexes of refraction of the core and cladding, respectively. Multimode fibers typically have a fractional index difference (Δ) between core and cladding of between 1 and 1.5 percent and a core diameter of between 50 and 100 μm . Single-mode fibers typically have $\Delta \approx 0.3$ percent and a core diameter of between 8 and 10 μm .

The fiber numerical aperture (NA), which is the sine of the half-angle of the cone of acceptance, is given by

$$\text{NA} = \sqrt{n_1^2 - n_2^2} = n_1 \sqrt{2\Delta} \quad (2)$$

Single-mode fibers typically have an NA of about 0.1, whereas the NA of multimode fibers is in the range of 0.2 to 0.3.

From a transmission system standpoint, the two most important fiber parameters are attenuation and bandwidth.

Attenuation There are three principal attenuation mechanisms in fiber: absorption, scattering, and radiative loss. Silicon dioxide has resonance absorption peaks in the ultraviolet (electronic transitions) and in the infrared beyond 1.6 μm (atomic vibrational transitions), but is highly transparent in the visible and near-infrared.

Radiative losses are generally kept small by using a sufficiently thick cladding (communication fibers have an outer diameter of 125 μm), a compressible coating to buffer the fiber from external forces, and a cable structure that prevents sharp bends.

In the absence of impurities and radiation losses, the fundamental attenuation mechanism is Rayleigh scattering from the irregular glass structure, which results in index of refraction fluctuations over distances that are small compared to the wavelength. This leads to a scattering loss

$$\alpha = \frac{B}{\lambda^4} \quad \text{with } B \approx 0.9 \frac{\text{dB}}{\text{km}} \mu\text{m}^4 \quad (3)$$

for “best” fibers. Attenuation as a function of wavelength is shown in Fig. 1. The attenuation peak at $\lambda = 1.4 \mu\text{m}$ is a resonance absorption due to small amounts of water in the fiber, although fibers are available in which this peak is absent. Initial systems operated at a wavelength around 0.85 μm owing to the availability of sources and detectors at this wavelength. Present systems (other than some short-distance data links) generally operate at wavelengths of 1.3 or 1.55 μm . The former, in addition to being low in attenuation (about 0.32 dB/km for best fibers), is the wavelength of minimum intramodal dispersion (see next section) for standard single-mode fiber. Operation at 1.55 μm allows even lower attenuation (minimum is about 0.16 dB/km) and the use of erbium-doped-fiber amplifiers (see Sec. 9.5), which operate at this wavelength.

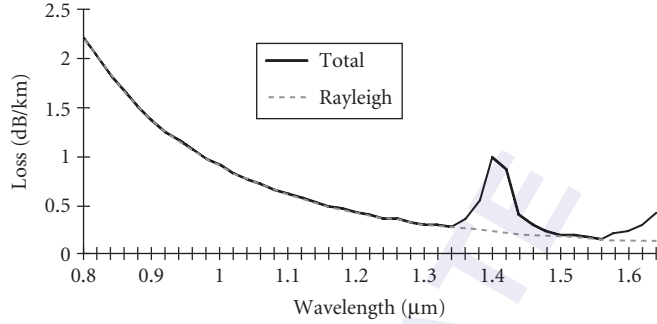


FIGURE 1 Fiber attenuation as a function of wavelength. Dashed curve shows Rayleigh scattering. Solid curve indicates total attenuation including resonance absorption at 1.38 μm from water and tail of infrared atomic resonances above 1.6 μm .

Dispersion Pulse spreading (dispersion) limits the maximum modulation bandwidth (or maximum pulse rate) that may be used with fibers. There are two principal forms of dispersion: intermodal dispersion and intramodal dispersion. In multimode fiber, the different modes experience different propagation delays resulting in pulse spreading. For graded-index fiber, the lowest dispersion per unit length is given approximately by⁷

$$\frac{\delta\tau}{L} = \frac{n_1\Delta^2}{10c} \quad (\text{intermodal}) \quad (4)$$

[Grading of the index of refraction of the core in a nearly parabolic function results in an approximate equalization of the propagation delays. For a step-index fiber, the dispersion per unit length is $\delta\tau/L = n_1\Delta/c$, which for $\Delta = 0.01$ is 1000 times larger than that given by Eq. (4).]

Bandwidth is inversely proportional to dispersion, with the proportionality constant dependent on pulse shape and how bandwidth is defined. If the dispersed pulse is approximated by a Gaussian pulse with $\delta\tau$ being the full width at the half-power point, then the -3-dB bandwidth B is given by

$$B = 0.44/\delta\tau \quad (5)$$

Multimode fibers are generally specified by their bandwidth in a 1-km length. Typical specifications are in the range from 200 MHz to 1 GHz. Fiber bandwidth is a sensitive function of the index profile and is wavelength dependent, and the scaling with length depends on whether there is mode mixing.⁸ Also, for short-distance links, the bandwidth is dependent on the launch conditions. Multimode fibers are generally used only when the bit rates and distances are sufficiently small that accurate characterization of dispersion is not of concern, although this may be changing with the advent of graded-index plastic optical fiber for high-bit-rate short-distance data links.

Although there is no intermodal dispersion in single-mode fibers,* there is still dispersion within the single mode (intramodal dispersion) resulting from the finite spectral width of the source and the dependence of group velocity on wavelength. The intramodal dispersion per unit length is given by

$$\begin{aligned} \delta\tau/L &= D \delta\lambda & \text{for } D \neq 0 \\ &= 0.2S_o(\delta\lambda)^2 & \text{for } D = 0 \end{aligned} \quad (6)$$

*A single-mode fiber actually has two degenerate modes corresponding to the two principal polarizations. Any asymmetry in the transmission path removes this degeneracy and results in polarization dispersion. This is typically very small (in the range of 0.1 to 1 ps/km^{1/2}), but is of concern in long-distance systems using linear repeaters.⁹

where D is the dispersion coefficient of the fiber, $\delta\lambda$ is the spectral width of the source, and S_o is the dispersion slope

$$S_o = \frac{dD}{d\lambda} \text{ at } \lambda = \lambda_0 \quad \text{where} \quad D(\lambda_0) = 0 \quad (7)$$

If both intermodal and intramodal dispersion are present, the square of the total dispersion is the sum of the squares of the intermodal and intramodal dispersions. For typical digital systems, the total dispersion should be less than half the interpulse period T . From Eq. (5) this corresponds to an effective fiber bandwidth that is at least $0.88/T$.

There are two sources of intramodal dispersion: material dispersion, which is a consequence of the index of refraction being a function of wavelength, and waveguide dispersion, which is a consequence of the propagation constant of the fiber waveguide being a function of wavelength.

For a material with index of refraction $n(\lambda)$, the material dispersion coefficient is given by

$$D_{\text{mat}} = -\frac{\lambda}{c} \frac{d^2n}{d\lambda^2} \quad (8)$$

For silica-based glasses, D_{mat} has the general characteristics shown in Fig. 2. It is about -100 ps/km · nm at a wavelength of 820 nm, goes through zero at a wavelength near 1300 nm, and is about 20 ps/km · nm at 1550 nm.

For step-index single-mode fibers, waveguide dispersion is given approximately by¹⁰

$$D_{\text{wg}} \approx -\frac{0.025\lambda}{a^2cn_2} \quad (9)$$

For conventional single-mode fiber, waveguide dispersion is small (about -5 ps/km · nm at 1300 nm). The resultant $D(\lambda)$ is then slightly shifted (relative to the material dispersion curve) to longer wavelengths, but the zero-dispersion wavelength (λ_0) remains in the vicinity of 1300 nm. However, if the waveguide dispersion is made larger negative by decreasing a or equivalently by tapering the index of refraction in the core the zero-dispersion wavelength may be shifted to the vicinity of 1550 nm

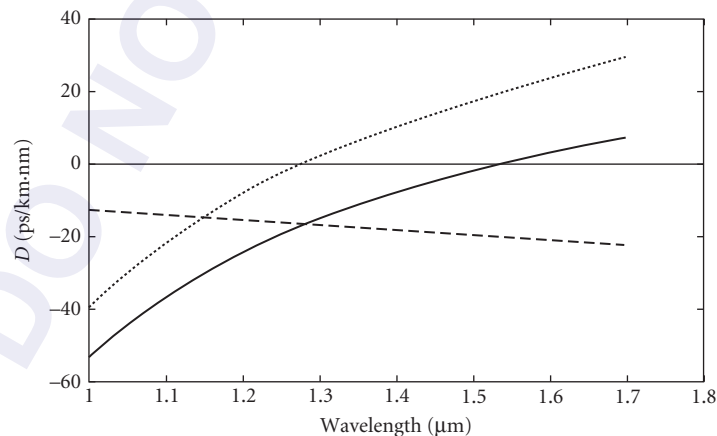


FIGURE 2 Intramodal dispersion coefficient as a function of wavelength. Dotted curve shows D_{mat} ; dashed curve shows D_{wg} to achieve D (solid curve) with zero dispersion at 1.55 μm .

(see Fig. 2). Such fibers are called *dispersion-shifted fibers* and are advantageous because of the lower fiber attenuation at this wavelength and the advent of erbium-doped-fiber amplifiers (see Sec. 9.5). Note that dispersion-shifted fibers have a smaller slope at the dispersion minimum ($S_0 \approx 0.06$ ps/km · nm² compared to $S_0 \approx 0.09$ ps/km · nm² for conventional single-mode fiber).

With more complicated index of refraction profiles, it is possible, at least theoretically, to control the shape of the waveguide dispersion such that the total dispersion is small in both the 1300- and 1550-nm bands, leading to dispersion-flattened fibers.¹¹

Transmitting Sources

Semiconductor light-emitting diodes (LEDs) or lasers are the primary light sources used in fiber-optic transmission systems. The principal parameters of concern are the power coupled into the fiber, the modulation bandwidth, and (because of intramodal dispersion) the spectral width.

Light-Emitting Diodes (LEDs) LEDs are forward-biased positive-negative (PN) junctions in which carrier recombination results in spontaneous emission at a wavelength corresponding to the energy gap. Although several milliwatts may be radiated from high-radiance LEDs, the radiation is over a wide angular range, and consequently there is a large coupling loss from an LED to a fiber. Coupling efficiency (η = ratio of power coupled to power radiated) from an LED to a fiber is given approximately by¹²

$$\begin{aligned}\eta &\approx (\text{NA})^2 && \text{for } r_s < a \\ \eta &\approx (a/r_s)^2 (\text{NA})^2 && \text{for } r_s > a\end{aligned}\quad (10)$$

where r_s is the radius of the LED. Use of large-diameter, high-NA multimode fiber improves the coupling from LEDs to fiber. Typical coupling losses are 10 to 20 dB for multimode fibers and more than 30 dB for single-mode fibers.

In addition to radiating over a large angle, LED radiation has a large spectral width (about 50 nm at $\lambda = 850$ nm and 100 nm at $\lambda = 1300$ nm) determined by thermal effects. Systems employing LEDs at 850 nm tend to be intramodal-dispersion-limited, whereas those at 1300 nm are intermodal-dispersion-limited.

Owing to the relatively long time constant for spontaneous emission (typically several nanoseconds), the modulation bandwidths of LEDs are generally limited to several hundred MHz. Thus, LEDs are generally limited to relatively short-distance, low-bit-rate applications.

Lasers In a laser, population inversion between the ground and excited states results in stimulated emission. In edge-emitting semiconductor lasers, this radiation is guided within the active region of the laser and is reflected at the end faces.* The combination of feedback and gain results in oscillation when the gain exceeds a threshold value. The spectral range over which the gain exceeds threshold (typically a few nanometers) is much narrower than the spectral width of an LED. Discrete wavelengths within this range, for which the optical length of the laser is an integer number of half-wavelengths, are radiated. Such a laser is termed a *multilongitudinal mode Fabry-Perot laser*. Radiation is confined to a much narrower angular range than for an LED, and consequently may be efficiently coupled into a small-NA fiber. Coupled power is typically about 1 mW.

The modulation bandwidth of lasers is determined by a resonance frequency caused by the interaction of the photon and electron concentrations.¹⁴ Although this resonance frequency was less than 1 GHz in early semiconductor lasers, improvements in materials have led to semiconductor lasers with resonance frequencies (and consequently modulation bandwidths) in excess of 10 GHz. This not only is important for very high-speed digital systems, but now also allows semiconductor lasers to be directly modulated with microwave signals. Such applications are considered in Sec. 9.7.

*In vertical cavity surface-emitting lasers (VCSELs), reflection is from internal "mirrors" grown within the semiconductor structure.¹³

Although multilongitudinal-mode Fabry-Perot lasers have a narrower spectral spread than LEDs, this spread still limits the high-speed and long-distance capability of such lasers. For such applications, single-longitudinal-mode (SLM) lasers are used. SLM lasers may be achieved by having a sufficiently short laser (less than 50 μm), by using coupled cavities (either external mirrors or cleaved coupled cavities¹⁵), or by incorporating a diffraction grating within the laser structure to select a specific wavelength. The latter has proven to be most practical for commercial application, and includes the distributed feedback (DFB) laser, in which the grating is within the laser active region, and the distributed Bragg reflector (DBR) laser, where the grating is external to the active region.¹⁶

There is still a finite line width for SLM lasers. For lasers without special stabilization, the line width is on the order of 0.1 nm. Expressed in terms of frequency, this corresponds to a frequency width of 12.5 GHz at a wavelength of 1550 nm. (Wavelength and frequency spread are related by $\delta f/f = -\delta\lambda/\lambda$, from which it follows that $\delta f = -c\delta\lambda/\lambda^2$.) Thus, unlike electrical communication systems, optical systems generally use sources with spectral widths that are large compared to the modulation bandwidth.

The finite line width (phase noise) of a laser is due to fluctuations of the phase of the optical field resulting from spontaneous emission. In addition to the phase noise contributed directly by the spontaneous emission, the interaction between the photon and electron concentrations in semiconductor lasers leads to a conversion of amplitude fluctuations to phase fluctuations, which increases the line width.¹⁷ If the intensity of a laser is changed, this same phenomenon gives rise to a change in the frequency of the laser (*chirp*). Uncontrolled, this causes a substantial increase in line width when the laser is modulated, which may cause difficulties in some system applications, possibly necessitating external modulation. However, the phenomenon can also be used to advantage. For appropriate lasers under small signal modulation, a change in frequency proportional to the input signal can be used to frequency-modulate and/or to tune the laser. Tunable lasers are of particular importance in networking applications employing wavelength-division multiplexing (WDM).¹⁸

Photodetectors

Fiber-optic systems generally use PIN or APD photodetectors. In a reverse-biased PIN diode, absorption of light in the intrinsic region generates carriers that are swept out by the reverse-bias field. This results in a photocurrent (I_p) that is proportional to the incident optical power (P_R), where the proportionality constant is the responsivity (\mathfrak{R}) of the photodetector; that is, $\mathfrak{R} = I_p/P_R$. Since the number of photons per second incident on the detector is power divided by the photon energy, and the number of electrons per second flowing in the external circuit is the photocurrent divided by the charge of the electron, it follows that the quantum efficiency ($\eta = \text{electrons/photons}$) is related to the responsivity by

$$\eta = \frac{hc}{q\lambda} \frac{I_p}{P_R} = \frac{1.24(\mu\text{m} \cdot \text{V})}{\lambda} \mathfrak{R} \quad (11)$$

For wavelengths shorter than 900 nm, silicon is an excellent photodetector, with quantum efficiencies of about 90 percent. For longer wavelengths, InGaAs is generally used, with quantum efficiencies typically around 70 percent. Very high bandwidths may be achieved with PIN photodetectors. Consequently, the photodetector does not generally limit the overall system bandwidth.

In an avalanche photodetector (APD), a larger reverse voltage accelerates carriers, causing additional carriers by impact ionization resulting in a current $I_{\text{APD}} = MI_p$, where M is the current gain of the APD. As we will note in Sec. 9.3, this can result in an improvement in receiver sensitivity.

9.3 RECEIVER SENSITIVITY

The receiver in a direct-detection fiber-optic communication system consists of a photodetector followed by electrical amplification and signal-processing circuits intended to recover the communications signal. Receiver sensitivity is defined as the average received optical power needed to

achieve a given communication rate and performance. For analog communications, the communication rate is measured by the bandwidth of the electrical signal to be transmitted (B), and performance is given by the signal-to-noise ratio (SNR) of the recovered signal. For digital systems, the communication rate is measured by the bit rate (R_b) and performance is measured by the bit error probability (P_e).

For a constant optical power transmitted, there are fluctuations of the received photocurrent about the average given by Eq. (11). The principal sources of these fluctuations are signal shot noise (quantum noise resulting from random arrival times of photons at the detector), receiver thermal noise, APD excess noise, and relative intensity noise (RIN) associated with fluctuations in intensity of the source and/or multiple reflections in the fiber medium.

Digital On-Off-Keying Receiver

It is instructive to define a normalized sensitivity as the average number of photons per bit (\bar{N}_p) to achieve a given error probability, which we take here to be $P_e = 10^{-9}$. Given \bar{N}_p , the received power when a 1 is transmitted is obtained from

$$P_R = 2\bar{N}_p R_b \frac{hc}{\lambda} \quad (12)$$

where the factor of 2 in Eq. (12) is because P_R is the peak power, and \bar{N}_p is the average number of photons per bit.

Ideal Receiver In an ideal receiver individual photons may be counted, and the only source of noise is the fluctuation of the number of photons counted when a 1 is transmitted. This is a Poisson random variable with mean $2\bar{N}_p$. No photons are received when a 0 is transmitted. Consequently, an error is made only when a 1 is transmitted and no photons are received. This leads to the following expression for the error probability

$$P_e = \frac{1}{2} \exp(-2\bar{N}_p) \quad (13)$$

from which it follows that $\bar{N}_p = 10$ for $P_e = 10^{-9}$. This is termed the *quantum limit*.

PIN Receiver In a PIN receiver, the photodetector output is amplified, filtered, and sampled, and the sample is compared with a threshold to decide whether a 1 or 0 was transmitted. Let I be the sampled current at the input to the decision circuit scaled back to the corresponding value at the output of the photodetector. (It is convenient to refer all signal and noise levels to their equivalent values at the output of the photodetector.) I is then a random variable with means and variances given by

$$\mu_1 = I_p \quad \mu_0 = 0 \quad (14a)$$

$$\sigma_1^2 = 2qI_p B + \frac{4kTB}{R_e} \quad \sigma_0^2 = \frac{4kTB}{R_e} \quad (14b)$$

where the subscripts 1 and 0 refer to the bit transmitted, kT is the thermal noise energy, and R_e is the effective input noise resistance of the amplifier. Note that the noise values in the 1 and 0 states are different owing to the shot noise in the 1 state.

Calculation of error probability requires knowledge of the distribution of I under the two hypotheses. Under the assumption that these distributions may be approximated by gaussian distributions

with means and variances given by Eq. (14), the error probability may be shown to be given by (Chap. 4 in Ref. 19)

$$P_e = K \left(\frac{\mu_1 - \mu_0}{\sigma_1 + \sigma_0} \right) \quad (15)$$

where

$$K(Q) = \frac{1}{\sqrt{2\pi}} \int_Q^\infty dx \exp(-x^2/2) = \frac{1}{2} \operatorname{erfc}(Q/\sqrt{2}) \quad (16)$$

It can be shown from Eqs. (11), (12), (14), and (15) that

$$\bar{N}_p = \frac{B}{\eta R_b} Q^2 \left[1 + \frac{1}{Q} \sqrt{\frac{8\pi kTC_e}{q^2}} \right] \quad (17)$$

where

$$C_e = \frac{1}{2\pi R_e B} \quad (18)$$

is the effective noise capacitance of the receiver, and from Eq. (16), $Q = 6$ for $P_e = 10^{-9}$. The minimum bandwidth of the receiver is half the bit rate, but in practice B/R_b is generally about 0.7.

The gaussian approximation is expected to be good when the thermal noise is large compared to the shot noise. It is interesting, however, to note that Eq. (17) gives $\bar{N}_p = 18$, when $C_e = 0$, $B/R_b = 0.5$, $\eta = 1$, and $Q = 6$. Thus, even in the shot noise limit, the gaussian approximation gives a surprisingly close result to the value calculated from the correct Poisson distribution. It must be pointed out, however, that the location of the threshold calculated by the gaussian approximation is far from correct in this case. In general, the gaussian approximation is much better in estimating receiver sensitivity than in establishing where to set receiver thresholds.

Low-input-impedance amplifiers are generally required to achieve the high bandwidths required for high-bit-rate systems. However, a low input impedance results in high thermal noise and poor sensitivity. High-input-impedance amplifiers may be used, but this narrows the bandwidth, which must be compensated for by equalization following the first-stage amplifier. Although this may result in a highly sensitive receiver, the receiver will have a poor dynamic range owing to the high gains required in the equalizer.²⁰ Receivers for digital systems are generally implemented with transimpedance amplifiers having a large feedback resistance. This reduces the effective input noise capacitance to below the capacitance of the photodiode, and practical receivers can be built with $C_e \approx 0.1 pF$. Using this value of capacitance and $B/R_b = 0.7$, $\eta = 0.7$, and $Q = 6$, Eq. (17) gives $\bar{N}_p \approx 2600$. Note that this is about 34 dB greater than the value given by the quantum limit.

APD Receiver In an APD receiver, there is additional shot noise owing to the excess noise factor F of the avalanche gain process. However, thermal noise is reduced because of the current multiplication gain M before thermal noise is introduced. This results in a receiver sensitivity given approximately by*

$$\bar{N}_p = \frac{B}{\eta R_b} Q^2 \left[F + \frac{1}{Q} \sqrt{\frac{8\pi kTC_e}{q^2 M^2}} \right] \quad (19)$$

*The gaussian approximation is not as good for an APD as for a PIN receiver owing to the nongaussian nature of the excess APD noise.

The excess noise factor is an increasing function of M , which results in an optimum M to minimize \bar{N}_p .²⁰ Good APD receivers at 1300 and 1550 nm typically have sensitivities of the order of 1000 photons per bit. Owing to the lower excess noise of silicon APDs, sensitivity of about 500 photons per bit can be achieved at 850 nm.

Impairments There are several sources of impairment that may degrade the sensitivity of receivers from the values given by Eqs. (17) and (19). These may be grouped into two general classes: eye degradations and signal-dependent noise.

An eye diagram is the superposition of all possible received sequences. At the sampling point, there is a spread of the values of a received 1 and a received 0. The difference between the minimum value of a received 1 and the maximum value of the received 0 is known as the *eye opening*. This is given by $(1 - \epsilon)I_p$ where ϵ is the eye degradation. The two major sources of eye degradation are intersymbol interference and finite laser extinction ratio. Intersymbol interference results from dispersion, deviations from ideal shaping of the receiver filter, and low-frequency cutoff effects that result in direct current (DC) offsets.

Signal-dependent noises are phenomena that give a variance of the received photocurrent that is proportional to I_p^2 and consequently lead to a maximum signal-to-noise ratio at the output of the receiver. Principal sources of signal-dependent noise are laser relative intensity noise (RIN), reflection-induced noise, mode partition noise, and modal noise. RIN is a consequence of inherent fluctuations in laser intensity resulting from spontaneous emission. This is generally sufficiently small that it is not of concern in digital systems, but is an important limitation in analog systems requiring high signal-to-noise ratios (see Sec. 9.7). Reflection-induced noise is the conversion of laser phase noise to intensity noise by multiple reflections from discontinuities (such as at imperfect connectors.) This may result in a substantial RIN enhancement that can seriously affect digital as well as analog systems.²¹ Mode partition noise occurs when Fabry-Perot lasers are used with dispersive fiber. Fiber dispersion results in changing phase relation between the various laser modes, which results in intensity fluctuations. The effect of mode partition noise is more serious than that of dispersion alone.²² Modal noise is a similar phenomenon that occurs in multimode fiber when relatively few modes are excited and these interfere.

Eye degradations are accounted for by replacing Eq. (14a) by

$$\mu_1 - \mu_0 = (1 - \epsilon)I_p \quad (20a)$$

and signal-dependent noise by replacing Eq. (14b) by

$$\sigma_1^2 = 2qI_pB + \frac{4kTB}{R_e} + \alpha^2 I_p^2 B \quad \sigma_0^2 = \frac{4kTB}{R_e} + \alpha^2 I_p^2 B \quad (20b)$$

where α^2 is the relative spectral density of the signal-dependent noise. (It is assumed that the signal-dependent noise has a large bandwidth compared to the signal bandwidth B .) With these modifications, the sensitivity of an APD receiver becomes

$$\bar{N}_p = \frac{\frac{B}{\eta R_p} \left(\frac{Q}{1 - \epsilon} \right)^2 \left[F + \left(\frac{1 - \epsilon}{Q} \right) \sqrt{\frac{8\pi kTC_e}{q^2 M^2}} \right]}{1 - \alpha^2 B \left(\frac{Q}{1 - \epsilon} \right)^2} \quad (21)$$

The sensitivity of a PIN receiver is obtained by setting $F = 1$ and $M = 1$ in Eq. (21). It follows from Eq. (21) that there is a minimum error probability (*error floor*) given by

$$P_{e,\min} = K(Q_{\max}) \text{ where } Q_{\max} = \frac{1 - \epsilon}{\alpha\sqrt{B}} \quad (22)$$

The existence of eye degradations and signal-dependent noise causes an increase in the receiver power (called *power penalty*) required to achieve a given error probability.

9.4 BIT RATE AND DISTANCE LIMITS

Bit rate and distance limitations of digital links are determined by loss and dispersion limitations. The following example is used to illustrate the calculation of the maximum distance for a given bit rate. Consider a 2.5 Gbit/s system at a wavelength of 1550 nm. Assume an average transmitter power of 0 dBm coupled into the fiber. Receiver sensitivity is taken to be 3000 photons per bit, which from Eq. (12) corresponds to an average receiver power of -30.2 dBm. Allowing a total of 8 dB for margin and for connector and cabling losses at the two ends gives a loss allowance of 22.2 dB. If the cabled fiber loss, including splices, is 0.25 dB/km, this leads to a loss-limited transmission distance of 89 km.

Assuming that the fiber dispersion is $D = 15$ ps/km · nm and source spectral width is 0.1 nm, this gives a dispersion per unit length of 1.5 ps/km. Taking the maximum allowed dispersion to be half the interpulse period, this gives a maximum dispersion of 200 ps, which then yields a maximum dispersion-limited distance of 133 km. Thus, the loss-limited distance is controlling.

Consider what happens if the bit rate is increased to 10 Gbit/s. For the same number of photons per bit at the receiver, the receiver power must be 6 dB greater than that in the preceding example. This reduces the loss allowance by 6 dB, corresponding to a reduction of 24 km in the loss-limited distance. The loss-limited distance is now 65 km (assuming all other parameters are unchanged). However, dispersion-limited distance scales inversely with bit rate, and is now 22 km. The system is now dispersion-limited. Dispersion-shifted fiber would be required to be able to operate at the loss limit.

Increasing Bit Rate

There are two general approaches for increasing the bit rate transmitted on a fiber: time-division multiplexing (TDM), in which the serial transmission rate is increased, and wavelength-division multiplexing (WDM), in which separate wavelengths are used to transmit independent serial bit streams in parallel. TDM has the advantage of minimizing the quantity of active devices but requires higher-speed electronics as the bit rate is increased. Also, as indicated by the preceding example, dispersion limitations will be more severe.

WDM allows use of existing lower-speed electronics, but requires multiple lasers and detectors as well as optical filters for combining and separating the wavelengths. Technology advances, including tunable lasers, transmitter and detector arrays, high-resolution optical filters, and optical amplifiers (Sec. 9.5) have made WDM more attractive, particularly for networking applications (Sec. 9.6).

Longer Repeater Spacing

In principal, there are three approaches for achieving longer repeater spacing than that calculated in the preceding text: lower fiber loss, higher transmitter powers, and improved receiver sensitivity (smaller \bar{N}_p). Silica-based fiber is already essentially at the theoretical Rayleigh scattering loss limit. There has been research on new fiber materials that would allow operation at wavelengths longer than 1.6 μm , with consequent lower theoretical loss values.²³ There are many reasons, however, why achieving such losses will be difficult, and progress in this area has been slow.

Higher transmitter powers are possible, but there are both nonlinearity and reliability issues that limit transmitter power. Since present receivers are more than 30 dB above the quantum limit, improved receiver sensitivity would appear to offer the greatest possibility. To improve the receiver sensitivity, it is necessary to increase the photocurrent at the output of the detector without introducing significant excess loss. There are two main approaches for doing so: optical amplification and optical mixing. Optical preamplifiers result in a theoretical sensitivity of 38 photons per bit²⁴ (6 dB above the quantum limit), and experimental systems have been constructed with sensitivities of about 100 photons per bit.²⁵ This will be discussed further in Sec. 9.5. Optical mixing (coherent receivers) will be discussed briefly in the following text.

Coherent Systems A photodetector provides an output current proportional to the magnitude square of the electric field that is incident on the detector. If a strong optical signal (*local oscillator*) coherent in phase with the incoming optical signal is added prior to the photodetector, then the photocurrent will contain a component at the difference frequency between the incoming and local oscillator signals. The magnitude of this photocurrent, relative to the direct detection case, is increased by the ratio of the local oscillator to the incoming field strengths. Such a coherent receiver offers considerable improvement in receiver sensitivity. With on-off keying, a heterodyne receiver (signal and local oscillator frequencies different) has a theoretical sensitivity of 36 photons per bit, and a homodyne receiver (signal and local oscillator frequencies the same) has a sensitivity of 18 photons per bit. Phase-shift keying (possible with coherent systems) provides a further 3-dB improvement. Coherent systems, however, require very stable signal and local oscillator sources (spectral linewidths need to be small compared to the modulation bandwidth) and matching of the polarization of the signal and local oscillator fields.²⁶ Differentially coherent systems (e.g., DPSK) in which the prior bit is used as a phase reference are simpler to implement and are beginning to find application.²⁷

An advantage of coherent systems, more so than improved receiver sensitivity, is that because the output of the photodetector is linear in the signal field, filtering for WDM demultiplexing may be done at the difference frequency (typically in the microwave range).^{*} This allows considerably greater selectivity than is obtainable with optical filtering techniques. The advent of optical amplifiers has slowed the interest in coherent systems, but there has been renewed interest in recent years.²⁸

9.5 OPTICAL AMPLIFIERS

There are two types of optical amplifiers: laser amplifiers based on stimulated emission and parametric amplifiers based on nonlinear effects (Chap. 10 in Ref. 32). The former are currently of most interest in fiber-optic communications. A laser without reflecting end faces is an amplifier, but it is more difficult to obtain sufficient gain for amplification than it is (with feedback) to obtain oscillation. Thus, laser oscillators were available much earlier than laser amplifiers.

Laser amplifiers are now available with gains in excess of 30 dB over a spectral range of more than 30 nm. Output saturation powers in excess of 10 dBm are achievable. The amplified spontaneous emission (ASE) noise power at the output of the amplifier, in each of two orthogonal polarizations, is given by

$$P_{\text{ASE}} = n_{\text{sp}} \frac{hc}{\lambda} B_o (G - 1) \quad (23)$$

where G is the amplifier gain, B_o is the bandwidth, and the spontaneous emission factor n_{sp} is equal to 1 for ideal amplifiers with complete population inversion.

Comparison of Semiconductor and Fiber Amplifiers

There are two principal types of laser amplifiers: semiconductor laser amplifiers (SLAs) and doped-fiber amplifiers. The erbium-doped-fiber amplifier (EDFA), which operates at a wavelength of 1.55 μm , is of most current interest.

The advantages of the SLA, similar to laser oscillators, are that it is pumped by a DC current, it may be designed for any wavelength of interest, and it can be integrated with electrooptic semiconductor components.

^{*}The difference frequency must be large compared to the modulation bandwidth. As modulation bandwidths have increased beyond 10 GHz this may necessitate difference frequencies greater than 100 GHz which may be difficult to implement.

The advantages of the EDFA are that there is no coupling loss to the transmission fiber, it is polarization-insensitive, it has lower noise than SLAs, it can be operated at saturation with no intermodulation owing to the long time constant of the gain dynamics, and it can be integrated with fiber devices. However, it does require optical pumping, with the principal pump wavelengths being either 980 or 1480 nm.

Communications Application of Optical Amplifiers

There are four principal applications of optical amplifiers in communication systems.^{29,30}

1. Transmitter power amplifiers
2. Compensation for splitting loss in distribution networks
3. Receiver preamplifiers
4. Linear repeaters in long-distance systems

The last application is of particular importance for long-distance networks (particularly undersea systems), where a bit-rate-independent linear repeater allows subsequent upgrading of system capacity (either TDM or WDM) with changes only at the system terminals. Although amplifier noise accumulates in such long-distance linear systems, transoceanic lengths are achievable with amplifier spacings of about 60 km corresponding to about 15-dB fiber attenuation between amplifiers.

However, in addition to the accumulation of ASE, there are other factors limiting the distance of linearly amplified systems, namely dispersion and the interaction of dispersion and nonlinearity.³¹ There are two alternatives for achieving very long-distance, very high-bit-rate systems with linear repeaters: *solitons*, which are pulses that maintain their shape in a dispersive medium,³² and dispersion compensation.³³

9.6 FIBER-OPTIC NETWORKS

Networks are communication systems used to interconnect a number of terminals within a defined geographic area—for example, local area networks (LANs), metropolitan area networks (MANs), and wide area networks (WANs). In addition to the transmission function discussed throughout the earlier portions of this chapter, networks also deal with the routing and switching aspects of communications.

Passive optical networks utilize couplers to distribute signals to users. In an $N \times N$ ideal star coupler, the signal on each input port is uniformly distributed among all output ports. If an average power P_T is transmitted at a transmitting port, the power received at a receiving port (neglecting transmission losses) is

$$P_R = \frac{P_T}{N}(1 - \delta_N) \quad (24)$$

where δ_N is the excess loss of the coupler. If N is a power of 2, an $N \times N$ star may be implemented by $\log_2 N$ stages of 2×2 couplers. Thus, it may be conservatively assumed that

$$1 - \delta_N = (1 - \delta_2)^{\log_2 N} = N^{\log_2(1 - \delta_2)} \quad (25)$$

The maximum bit rate per user is given by the average received power divided by the product of the photon energy and the required number of photons per bit (N_p). The throughput Y is the product of the number of users and the bit rate per user, and from Eqs. (24) and (25) is therefore given by

$$Y = \frac{P_T}{N_p} \frac{\lambda}{hc} N^{\log_2(1 - \delta_2)} \quad (26)$$

Thus, the throughput (based on power considerations) is independent of N for ideal couplers ($\delta_2 = 0$) and decreases slowly with $N(N^{-0.17})$ for $10 \log(1 - \delta_2) = 0.5$ dB. It follows from Eq. (26) that for a power of 1 mW at $\lambda = 1.55 \mu\text{m}$ and with $N_p = 3000$, the maximum throughput is 2.6 Tbit/s.

This may be contrasted with a tapped bus, where it may be shown that optimum tap weight to maximize throughput is given by $1/N$, leading to a throughput given by³⁴

$$Y = \frac{P_T}{N_p} \frac{\lambda}{hc} \frac{1}{Ne^2} \exp(-2N\delta) \quad (27)$$

Thus, even for ideal ($\delta = 0$) couplers, the throughput decreases inversely with the number of users. If there is excess coupler loss, the throughput decreases exponentially with the number of users and is considerably less than that given by Eq. (26). Consequently, for a power-limited transmission medium, the star architecture is much more suitable than the tapped bus. The same conclusion does not apply to metallic media, where bandwidth rather than power limits the maximum throughput.

Although the preceding text indicates the large throughput that may be achieved in principle with a passive star network, it doesn't indicate how this can be realized. Most interest is in WDM networks.³⁵ The simplest protocols are those for which fixed-wavelength receivers and tunable transmitters are used. However, the technology is simpler when fixed-wavelength transmitters and tunable receivers are used, since a tunable receiver may be implemented with a tunable optical filter preceding a wideband photodetector. Fixed-wavelength transmitters and receivers involving multiple passes through the network are also possible, but this requires utilization of terminals as relay points. Protocol, technology, and application considerations for gigabit networks (networks having access at gigabit rates and throughputs at terabit rates) is an extensive area of research.³⁶

9.7 ANALOG TRANSMISSION ON FIBER

Most interest in fiber-optic communications is centered around digital transmission, since fiber is generally a power-limited rather than a bandwidth-limited medium. There are applications, however, where it is desirable to transmit analog signals directly on fiber without converting them to digital signals. Examples are cable television (CATV) distribution and microwave links such as entrance links to antennas and interconnection of base stations in mobile radio systems.

Carrier-to-Noise Ratio (CNR)

Optical intensity modulation is generally the only practical modulation technique for incoherent-detection fiber-optic systems. Let $f(t)$ be the carrier signal that intensity modulates the optical source. For convenience, assume that the average value of $f(t)$ is equal to 0, and that the magnitude of $f(t)$ is normalized to be less than or equal to 1. The received optical power may then be expressed as

$$P(t) = P_o[1 + mf(t)] \quad (28)$$

where m is the optical modulation index

$$m = \frac{P_{\max} - P_{\min}}{P_{\max} + P_{\min}} \quad (29)$$

The carrier-to-noise ratio is then given by

$$\text{CNR} = \frac{\frac{1}{2} m^2 \mathfrak{R}^2 P_o^2}{\text{RIN} \mathfrak{R}^2 P_o^2 B + 2q \mathfrak{R} P_o B + \langle i_{\text{th}}^2 \rangle B} \quad (30)$$

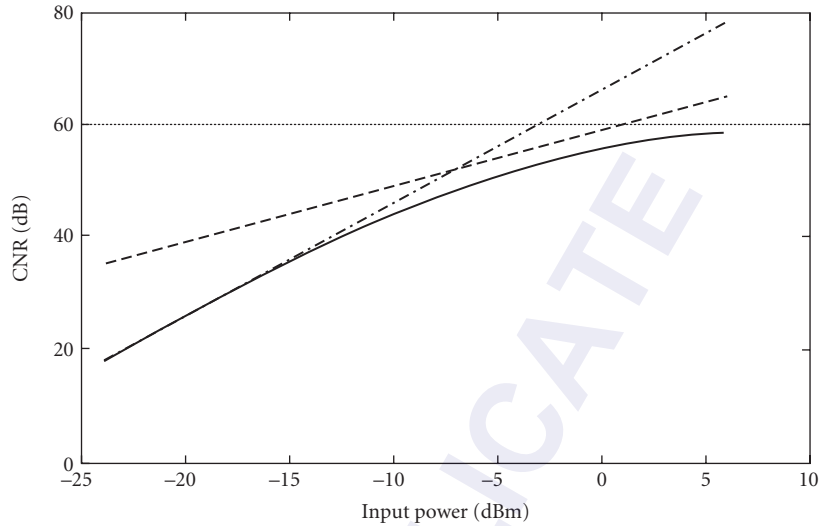


FIGURE 3 CNR as a function of input power. Straight lines indicate thermal noise (-----), shot noise (—), and RIN (.....) limits.

where \mathfrak{R} is the photodetector responsivity, RIN is the relative intensity noise spectral density (denoted by α^2 in Sec. 9.3), and $\langle i_{\text{th}}^2 \rangle$ is the thermal noise spectral density (expressed as $4kT/R_e$ in Sec. 9.3). CNR is plotted in Fig. 3 as a function of received optical power for a bandwidth of $B = 4$ MHz (single video channel), optical modulation index $m = 0.05$, $\mathfrak{R} = 0.8$ A/W, $\text{RIN} = -155$ dB/Hz, and $\sqrt{\langle i_{\text{th}}^2 \rangle} = 7$ pA/ $\sqrt{\text{Hz}}$. At low received powers (typical of digital systems) the CNR is limited by thermal noise. However, to obtain the higher CNR generally needed by analog systems, shot noise and then ultimately laser RIN become limiting.

Analog Video Transmission on Fiber³⁷

It is helpful to distinguish between single-channel and multiple-channel applications. For the single-channel case, the video signal may directly modulate the laser intensity [amplitude-modulated (AM) system], or the video signal may be used to frequency-modulate an electrical subcarrier, with this subcarrier then intensity-modulating the optical source [frequency-modulated (FM) system]. Equation (30) gives the CNR of the recovered subcarrier. Subsequent demodulation of the FM signal gives an additional increase in signal-to-noise ratio. In addition to this FM improvement factor, larger optical modulation indexes may be used than in AM systems. Thus FM systems allow higher signal-to-noise ratios and longer transmission spans than AM systems.

Two approaches have been used to transmit multichannel video signals on fiber. In the first (AM systems), the video signals undergo electrical frequency-division multiplexing (FDM), and this combined FDM signal intensity modulates the optical source. This is conceptually the simplest system, since existing CATV multiplexing formats may be used.

In FM systems, the individual video channels frequency-modulate separate microwave carriers (as in satellite systems). These carriers are linearly combined and the combined signal intensity modulates a laser. Although FM systems are more tolerant than AM systems to intermodulation distortion and noise, the added electronics costs have made such systems less attractive than AM systems for CATV application.

Multichannel AM systems are of interest not only for CATV application but also for mobile radio applications to connect signals from a microcellular base station to a central processing station.

Relative to CATV applications, the mobile radio application has the additional complication of being required to accommodate signals over a wide dynamic power range.

Nonlinear Distortion

In addition to CNR requirements, multichannel analog communication systems are subject to intermodulation distortion. If the input to the system consists of a number of tones at frequencies ω_j , then nonlinearities result in intermodulation products at frequencies given by all sums and differences of the input frequencies. Second-order intermodulation gives intermodulation products at frequencies $\omega_j \pm \omega_k$, whereas third-order intermodulation gives frequencies $\omega_j \pm \omega_k \pm \omega_l$. If the signal frequency band is such that the maximum frequency is less than twice the minimum frequency, then all second-order intermodulation products fall outside the signal band, and third-order intermodulation is the dominant nonlinearity. This condition is satisfied for the transport of microwave signals (e.g., mobile radio signals) on fiber, but is not satisfied for wideband CATV systems, where there are requirements on composite second-order (CSO) and composite triple-beat (CTB) distortion.

The principal causes of intermodulation in multichannel fiber-optic systems are laser threshold nonlinearity,³⁸ inherent laser gain nonlinearity, and the interaction of chirp and dispersion.

9.8 TECHNOLOGY AND APPLICATIONS DIRECTIONS

Fiber-optic communication application in the United States began with metropolitan and short-distance intercity trunking at a bit rate of 45 Mbit/s, corresponding to the DS-3 rate of the North American digital hierarchy. Technological advances, primarily higher-capacity transmission and longer repeater spacings, extended the application to long-distance intercity transmission, both terrestrial and undersea. Also, transmission formats are now based on the synchronous digital hierarchy (SDH), termed synchronous optical network (SONET) in the U.S. OC-192 system* operating at 10 Gbit/s are widely deployed, with OC-768 40 Gbit/s systems also available. All of the signal processing in these systems (multiplexing, switching, performance monitoring) is done electrically, with optics serving solely to provide point-to-point links.

For long-distance applications, 10 Gbit/s dense wavelength-division multiplexing (DWDM), with channel spacings of 50 GHz and with upward of 100 wavelength channels, has extended the bit rate capability of fiber to greater than 1 Tbit/s in commercial systems and more than 3 Tbit/s in laboratory trials.³⁹ For local access, there is extensive interest in fiber directly to the premises⁴⁰ as well as hybrid combinations of optical and electronic technologies and transmission media.^{41,42}

The huge bandwidth capability of fiber optics (measured in tens of terahertz) is not likely to be utilized by time-division techniques alone, and DWDM technology and systems are receiving considerable emphasis, although work is also under way on optical time-division multiplexing (OTDM) and optical code-division multiplexing (OCDM).

Nonlinear phenomena, when uncontrolled, generally lead to system impairments. However, controlled nonlinearities are the basis of devices such as parametric amplifiers and switching and logic elements. Nonlinear optics will consequently continue to receive increased emphasis.

9.9 REFERENCES

1. J. Hecht, *City of Light: The Story of Fiber Optics*, Oxford University Press, New York, 1999.
2. C. K. Kao and G. A. Hockham, "Dielectric-Fiber Surface Waveguides for Optical Frequencies," *Proc. IEEE* **113**:1151–1158 (July 1966).
3. F. P. Kapron, et al., "Radiation Losses in Glass Optical Waveguides," *Appl. Phys. Lett.* **17**:423 (November 15, 1970).

*OC- n systems indicate optical channel at a bit rate of $(51.84)n$ Mbit/s.

4. I. Hayashi, M. B. Panish, and P. W. Foy, "Junction Lasers which Operate Continuously at Room Temperature," *Appl. Phys. Lett.* **17**:109 (1970).
5. I. Jacobs, "Optical Fiber Communication Technology and System Overview," in O. D. D. Soares (ed.), *Trends in Optical Fibre Metrology and Standards*, NATO ASI Series, vol. 285, pp. 567–591, Kluwer Academic Publishers, Dordrecht, The Netherlands, 1995.
6. D. Gloge, "Weakly Guiding Fibers," *Appl. Opt.* **10**:2252–2258 (October 1971).
7. R. Olshansky and D. Keck, "Pulse Broadening in Graded Index Fibers," *Appl. Opt.* **15**:483–491 (February 1976).
8. D. Gloge, E. A. J. Marcatili, D. Marcuse, and S. D. Personick, "Dispersion Properties of Fibers," in S. E. Miller and A. G. Chynoweth (eds.), *Optical Fiber Telecommunications*, chap. 4, Academic Press, New York, 1979.
9. Y. Namihiro and H. Wakabayashi, "Fiber Length Dependence of Polarization Mode Dispersion Measurements in Long-Length Optical Fibers and Installed Optical Submarine Cables," *J. Opt. Commun.* **2**:2 (1991).
10. W. B. Jones Jr., *Introduction to Optical Fiber Communication Systems*, pp. 90–92, Holt, Rinehart and Winston, New York, 1988.
11. L. G. Cohen, W. L. Mammel, and S. J. Jang, "Low-Loss Quadruple-Clad Single-Mode Lightguides with Dispersion Below 2 ps/km · nm over the 1.28 μm –1.65 μm Wavelength Range," *Electron. Lett.* **18**:1023–1024 (1982).
12. G. Keiser, *Optical Fiber Communications*, 3d ed., chap. 5, McGraw-Hill, New York, 2000.
13. N. M. Margalit, S. Z. Zhang, and J. E. Bowers, "Vertical Cavity Lasers for Telecom Applications," *IEEE Commun. Mag.* **35**:164–170 (May 1997).
14. J. E. Bowers and M. A. Pollack, "Semiconductor Lasers for Telecommunications," in S. E. Miller and I. P. Kaminow (eds.), *Optical Fiber Telecommunications II*, chap. 13, Academic Press, San Diego, CA, 1988.
15. W. T. Tsang, "The Cleaved-Coupled-Cavity (C^3) Laser," in *Semiconductors and Semimetals*, vol. 22, part B, chap. 5, pp. 257–373, 1985.
16. K. Kobayashi and I. Mito, "Single Frequency and Tunable Laser Diodes," *J. Lightwave Technol.* **6**:1623–1633 (November 1988).
17. T. Mukai and Y. Yamamoto, "AM Quantum Noise in 1.3 μm InGaAsP Lasers," *Electron. Lett.* **20**:29–30 (January 5, 1984).
18. J. Buus and E. J. Murphy, "Tunable Lasers in Optical Networks," *J. Lightwave Technol.* **24**:5–11 (January 2006).
19. G. P. Agrawal, *Fiber-Optic Communication Systems*, 3d ed., Wiley Interscience, New York, 2002.
20. S. D. Personick, "Receiver Design for Digital Fiber Optic Communication Systems I," *Bell Syst. Tech. J.* **52**:843–874 (July–August 1973).
21. J. L. Gimlett and N. K. Cheung, "Effects of Phase-to-Intensity Noise Conversion by Multiple Reflections on Gigabit-per-Second DFB Laser Transmission Systems," *J. Lightwave Technol.* **LT-7**:888–895 (June 1989).
22. K. Ogawa, "Analysis of Mode Partition Noise in Laser Transmission Systems," *IEEE J. Quantum Electron.* **QE-18**:849–855 (May 1982).
23. D. C. Tran, G. H. Sigel, and B. Bendow, "Heavy Metal Fluoride Fibers: A Review," *J. Lightwave Technol.* **LT-2**:566–586 (October 1984).
24. P. S. Henry, "Error-Rate Performance of Optical Amplifiers," *Optical Fiber Communications Conference (OFC'89 Technical Digest)*, THK3, Houston, Texas, February 9, 1989.
25. O. Gautheron, G. Grandpierre, L. Pierre, J.-P. Thiery, and P. Kretzmeyer, "252 km Repeaterless 10 Gbits/s Transmission Demonstration," *Optical Fiber Communications Conference (OFC'93) Post-deadline Papers*, PD11, San Jose, California, February 21–26, 1993.
26. I. W. Stanley, "A Tutorial Review of Techniques for Coherent Optical Fiber Transmission Systems," *IEEE Commun. Mag.* **23**:37–53 (August 1985).
27. A. H. Gnauck, S. Chandrasekhar, J. Leutholdt, and L. Stulz, "Demonstration of 42.7 Gb/s DPSK Receiver with 45 Photons/Bit Sensitivity," *Photonics Technol. Letts.* **15**:99–101 (January 2003).
28. E. Ip, A. P. T. Lau, D. J. F. Barros, and J. M. Kahn, "Coherent Detection of Optical Fiber Systems," *Optics Express* **16**:753–791 (January 9, 2008).
29. Bellcore, "Generic Requirements for Optical Fiber Amplifier Performance," Technical Advisory TA-NWT-001312, Issue 1, December 1992.
30. T. Li, "The Impact of Optical Amplifiers on Long-Distance Lightwave Telecommunications," *Proc. IEEE* **81**:1568–1579 (November 1993).

31. A. Naka and S. Saito, "In-Line Amplifier Transmission Distance Determined by Self-Phase Modulation and Group-Velocity Dispersion," *J. Lightwave Technol.* **12**:280–287 (February 1994).
32. G. P. Agrawal, *Nonlinear Fiber Optics*, 3d ed., chap. 5, Academic Press, San Diego, CA, 2001.
33. Bob Jopson and Alan Gnauck, "Dispersion Compensation for Optical Fiber Systems," *IEEE Commun. Mag.* **33**:96–102 (June 1995).
34. P. E. Green, Jr., *Fiber Optic Networks*, chap. 11, Prentice Hall, Englewood Cliffs, NJ, 1993.
35. M. Fujiwara, M. S. Goodman, M. J. O'Mahony, O. K. Tonguz, and A. E. Willner (eds.), Special Issue on Multiwavelength Optical Technology and Networks, *J. Lightwave Technology* **14**(6):932–1454 (June 1996).
36. P. J. Smith, D. W. Faulkner, and G. R. Hill, "Evolution Scenarios for Optical Telecommunication Networks Using Multiwavelength Transmission," *Proc. IEEE* **81**:1580–1587 (November 1993).
37. T. E. Darcie, K. Nawata, and J. B. Glabb, Special Issue on Broad-Band Lightwave Video Transmission, *J. Lightwave Technol.* **11**(1) (January 1993).
38. A. A. M. Saleh, "Fundamental Limit on Number of Channels in SCM Lightwave CATV System," *Electron. Lett.* **25**(12):776–777 (1989).
39. A. H. Gnauck, G. Charlet, P. Tran, et al., "25.6 Tb/s WDM Transmission of Polarization Multiplexed RZ-DQPSK Signals," *J. Lightwave Technol.* **26**:79–84 (January 1, 2008).
40. T. Koonen, "Fiber to the Home/Fiber to the Premises: What, Where, and When?" *Proc. IEEE* **94**:911–934 (May 2006).
41. C. Baack and G. Walf, "Photonics in Future Telecommunications," *Proc. IEEE* **81**:1624–1632 (November 1993).
42. G. C. Wilson, T. H. Wood, J. A. Stiles, et al., "FiberVista: An FTTH or FTTC System Delivering Broadband Data and CATV Services," *Bell Labs Tech. J.* **4**:300–322 (January–March 1999).

This page intentionally left blank.

DO NOT DUPLICATE

NONLINEAR EFFECTS IN OPTICAL FIBERS

John A. Buck

*Georgia Institute of Technology
School of Electrical and Computer Engineering
Atlanta, Georgia*

Fiber nonlinearities are important in optical communications, both as useful attributes and as characteristics to be avoided. They must be considered when designing long-range high-data-rate systems that involve high optical power levels and in which signals at multiple wavelengths are transmitted. The consequences of nonlinear transmission can include (1) the generation of additional signal bandwidth within a given channel, (2) modifications of the phase and shape of pulses, (3) the generation of light at other wavelengths at the expense of power in the original signal, and (4) crosstalk between signals at different wavelengths and polarizations. The first two, arising from self-phase modulation, can be used to advantage in the generation of *solitons*—pulses whose nonlinear phase modulation compensates for linear group dispersion in the fiber link¹ or in fiber gratings,² leading to pulses that propagate without changing shape or width (see Chap. 22). The third and fourth effects arise from stimulated Raman or Brillouin scattering or four-wave mixing. These can be used to advantage when it is desired to generate or amplify additional wavelengths, but they must usually be avoided in systems.

10.1 KEY ISSUES IN NONLINEAR OPTICS IN FIBERS

Optical fiber waveguides, being of glass compositions, do not possess large nonlinear coefficients. Nonlinear processes can nevertheless occur with high efficiencies since intensities are high and propagation distances are long. Even though power levels are usually modest (a few tens of milliwatts), intensities within the fiber are high due to the small cross-sectional areas involved. This is particularly true in single-mode fiber, where the LP_{01} mode typically presents an effective cross-sectional area of between 10^{-7} and 10^{-8} cm², thus leading to intensities on the order of MW/cm². Despite this, long interaction distances are usually necessary to achieve nonlinear mixing of any significance, so processes must be phase matched, or nearly so. Strategies to avoid unwanted nonlinear effects usually involve placing upper limits on optical power levels, and if possible, choosing other parameters such that phase mismatching occurs. Such choices may include wavelengths or wavelength spacing in wavelength-division multiplexed systems, or may be involved in special fiber waveguide designs.³

The generation of light through nonlinear mixing arises through polarization of the medium, which occurs through its interaction with intense light. The polarization consists of an array of phased dipoles in which the dipole moment is a nonlinear function of the applied field strength. In the classical picture, the dipoles, once formed, reradiate light to form the nonlinear output. The medium polarization is conveniently expressed through a power series expansion involving products of real electric fields:

$$\mathcal{P} = \epsilon_0 [\chi^{(1)} \cdot \mathcal{E} + \chi^{(2)} \cdot \mathcal{E}\mathcal{E} + \chi^{(3)} \cdot \mathcal{E}\mathcal{E}\mathcal{E} + \dots] = \mathcal{P}_L + \mathcal{P}_{NL} \quad (1)$$

in which the χ terms are the linear, second-, and third-order susceptibilities. Nonlinear processes are described through the product of two or more optical fields to form the nonlinear polarization, \mathcal{P}_{NL} , consisting of all terms of second order and higher in Eq. (1).

The second-order term in Eq. (1) [involving $\chi^{(2)}$] describes three-wave mixing phenomena, such as second-harmonic generation. The third-order term describes four-wave mixing (FWM) processes and stimulated scattering phenomena. In the case of optical fibers, second-order processes are generally not possible, since these effects require noncentrosymmetric media.⁴ In amorphous fiber waveguides, third-order effects [involving $\chi^{(3)}$] are usually seen exclusively, although second-harmonic generation can be observed in special instances.⁵

The interactions between fields and polarizations are described by the nonlinear wave equation:

$$\nabla^2 \mathcal{E} + n_0^2 \mu_0 \epsilon_0 \frac{\partial^2 \mathcal{E}}{\partial t^2} = \mu_0 \frac{\partial^2 \mathcal{P}_{NL}}{\partial t^2} \quad (2)$$

where \mathcal{E} and \mathcal{P} are the sums of all electric fields and nonlinear polarizations that are present, and n_0 is the refractive index of the medium. The second-order differential equation is usually reduced to first order through the slowly varying envelope approximation (SVEA):

$$\left| \frac{\partial^2 E}{\partial z^2} \right| \ll \left| \frac{2\pi}{\lambda} \frac{\partial E}{\partial z} \right| \quad (3)$$

where E is the complex field amplitude. The interpretation of the SVEA is that the changes in field amplitude that occur over distances on the order of a wavelength are very large compared to variations in the rate of change over the same distance. The wave equation will separate according to frequencies or propagation directions, yielding sets of coupled differential equations that, under the SVEA, are first order. These describe the growth or decay of fields involved in the mixing process.

The requirement for phase matching is that the nonlinear polarization wave and the electric field associated with the generated wave propagate with the same phase constant; that is, their phase velocities are equal. Phase-matched processes in fiber include those that involve (1) interacting waves at the same wavelength and polarization, such as self- and cross-phase modulation, as well as other degenerate Kerr-type interactions, and (2) stimulated scattering processes (Raman and Brillouin), in addition to cross-phase modulation involving two wavelengths. Four-wave mixing processes involving light at different wavelengths can occur that are not precisely phase matched but that can nevertheless yield high efficiencies. Matters are further complicated by the fact that different nonlinear processes can occur simultaneously, with each affecting the performance of the other. Nonlinear effects are usually favored to occur under pulsed operation, since high peak powers can be achieved with comparatively modest average powers. Consequently, group velocity matching is desirable (although not always required) to achieve efficient mixing between pulses.

10.2 SELF- AND CROSS-PHASE MODULATION

Self-phase modulation (SPM) can occur whenever a signal having a time-varying amplitude is propagated in a nonlinear material. The origin of the effect is the refractive index of the medium, which will change with the instantaneous signal intensity. The complex nonlinear polarization for the process is:

$$P_{\text{NL}} = \frac{3}{4} \epsilon_0 \chi^{(3)} |E_0(z, t)|^2 E_0(z, t) \exp[i(\omega t - \beta z)] \quad (4)$$

where $E_0(t)$ is the time-varying electric field amplitude that describes the pulse or signal envelope, and where the frequency ω is the same as that of the input light. Incorporating this polarization and the field into the wave equation leads to a modified refractive index over the original zero-field value n_0 . The net index becomes:⁶

$$n = n_0 + n'_2 |E_0(z, t)|^2 \quad (5)$$

where the nonlinear refractive index is given by $n'_2 = \text{Re}\{3\chi^{(3)}/8n_0\}$. In fused silica it has the value $n'_2 = 6.1 \times 10^{-23} \text{ m}^2/\text{V}^2$.⁷ Equation (5) can also be expressed in terms of light intensity through $n(I) = n_0 + n_2 I(z, t)$, where $n_2 = 3.2 \times 10^{-20} \text{ m}^2/\text{W}$. In optical fibers the index is modified from the effective mode index of the single-mode fiber n_{eff} (which assumes the role of n_0).

The complex field as it propagates through the medium can be expressed as:

$$E = E_0(z, t) \exp(i\{\omega_0 t - [n_0 + n_2 I(z, t)]k_0 z\}) \quad (6)$$

which exhibits phase modulation that follows the shape of the intensity envelope. The instantaneous frequency is found through the time derivative of the phase:

$$\omega' = \omega_0 - n_2 k_0 z \frac{\partial I}{\partial t} \quad (7)$$

The effects of self-phase modulation on pulse propagation can be qualitatively observed from Eqs. (6) and (7). First, additional frequency components are placed on the pulse, thus increasing its spectral width. Second, a frequency sweep (chirp) is imposed on the pulse, the direction of which depends on the sign of $\partial I/\partial t$. The latter feature is particularly important in optical fibers, since the imposed frequency sweep from SPM will either add to or subtract from the chirp imposed by linear group dispersion. If the chirp directions for self-phase modulation and group dispersion are opposite, an effective cancellation may occur, leading to the formation of an optical soliton. In more conventional systems in which solitons are not employed, SPM must be considered as a possible benefit or detriment to performance, as some pulse shaping (which could include broadening or compression) can occur;^{8,9} however, such systems can in theory yield excellent performance.¹⁰ Furthermore, in systems employing fiber amplifiers, the change in refractive index associated with the signal-induced upper state population in erbium has been shown to be an important performance factor.¹¹ An additional effect can occur when pulse spectra lie within the anomalous group dispersion regime of the fiber; pulse breakup can occur as a result of *modulation instability*, in which the interplay between dispersive and nonlinear contributions to pulse shaping becomes unstable.¹²

Cross-phase modulation (XPM) is similar to SPM, except that two overlapping but distinguishable pulses (having, for example, different frequencies or polarizations) are involved. One pulse will modulate the index of the medium, which then leads to phase modulation of an overlapping pulse. XPM thus becomes a cross-talk mechanism between two channels if phase encoding is employed or if intensity modulation is used in dispersive systems.^{13,14} No transfer of energy occurs between channels, however, which distinguishes the process from other crosstalk mechanisms in which growth

of signal power in one channel occurs at the expense of power in another. The strength of the effect is enhanced by a factor of 2 over that which can be obtained by a single field acting on itself (the nonlinear refractive index n_2 is effectively doubled in XPM). The XPM process, while twice as strong as SPM, is effectively weakened by the fact that pulses of differing frequencies or polarizations are generally not group velocity matched, and so cannot maintain overlap indefinitely. The efficiency is further reduced if the interaction occurs between cross-polarized waves; in this case the nonlinear tensor element (and thus the effective nonlinear index) is a factor of $1/3$ less than the tensor element that describes copolarized waves.⁶

Self- and cross-phase modulation are analyzed by way of coupled equations of the *nonlinear Schrödinger* form,¹⁶ which describes the evolution over time and position of the electric field envelopes of two pulses, E_{0a} and E_{0b} , where SVEA is used and where pulse widths are on the order of 1 ps or greater:

$$\frac{\partial E_{0a}}{\partial z} + \beta_{1a} \frac{\partial E_{0a}}{\partial t} = -\frac{i}{2} \beta_{2a} \frac{\partial^2 E_{0a}}{\partial t^2} + i \gamma_a |E_{0a}|^2 E_{0a} + i \delta \gamma_a |E_{0b}|^2 E_{0a} - \frac{\alpha_a}{2} E_{0a} \quad (8)$$

$$\frac{\partial E_{0b}}{\partial z} + \beta_{1b} \frac{\partial E_{0b}}{\partial t} = -\frac{i}{2} \beta_{2b} \frac{\partial^2 E_{0b}}{\partial t^2} + i \gamma_b |E_{0b}|^2 E_{0b} + i \delta \gamma_b |E_{0a}|^2 E_{0b} - \frac{\alpha_b}{2} E_{0b} \quad (9)$$

In these equations, β_{1j} ($j = a, b$) are the group delays of the pulses at the two frequencies or polarizations over a unit distance; β_{2j} are the group dispersion parameters associated with the two pulses; and $\gamma_j = n_2' \omega_j / (c A_{\text{eff}})$, where A_{eff} is the effective cross-sectional area of the fiber mode. The coefficient δ is equal to 2 for copolarized pulses of different frequencies and is $2/3$ if the pulses are cross-polarized. Propagation loss characterized by coefficients α_j is assumed. The equation form that describes the propagation with SPM of a single pulse— E_{0a} , for example—is found from Eq. (8) by setting $E_{0b} = 0$. The terms on the right sides of Eqs. (8) and (9) describe in order the effects of group dispersion, SPM, XPM, and loss. The equations can be solved using numerical techniques that are described in Refs. 15 and 16.

For subpicosecond pulses, the accuracy of Eqs. (8) and (9) begins to degrade as pulse band-widths increase with decreasing temporal width. Additional terms are usually incorporated in the equations as pulse widths are reduced to the vicinity of 100 fs. These embody (1) cubic dispersion, which becomes important as bandwidth increases, and (2) changes in group velocity with intensity. This latter effect can result in *self-steepening*, in which the pulse trailing edge shortens to the point of forming an optical shock front under appropriate conditions. An additional consequence of broad pulse spectra is that power conversion from high-frequency components within a pulse to those at lower frequencies can occur via stimulated Raman scattering, provided the interacting components are sufficiently separated in wavelength. The effect is an overall red shift of the spectrum. At sufficiently high intensities, cross-coupling between pulses having different center wavelengths can also occur through Raman scattering, regardless of pulse width.

10.3 STIMULATED RAMAN SCATTERING

In *stimulated Raman scattering* (SRS), coupling occurs between copropagating light waves whose frequency difference is in the vicinity of resonances of certain molecular oscillation modes. In silica-based fibers, stretch vibrational resonances occur between Si and O atoms in several possible modes within the glass matrix (see Ref. 17 for illustrations of the important modes in pure silica). In the Stokes process, light at frequency ω_2 (pump wave) is downshifted to light at ω_1 (Stokes wave), with the excess energy being absorbed by the lattice vibrational modes (manifested in the generation of optical phonons). The process is either spontaneous, in which the Stokes wave builds up from noise, or is stimulated, in which both waves are present in sufficient strength to generate a beat frequency that excites the oscillators and promotes coupling. A fiber Raman amplifier works on this principle, in which an input signal at ω_1 experiences gain in the presence of pump light at ω_2 . Figure 1 shows

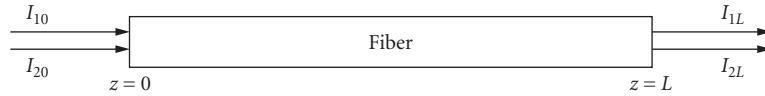


FIGURE 1 Beam geometry for stimulated Raman scattering in an optical fiber.

the beam geometry in which an input wave at ω_1 and intensity I_{10} can emerge at the far end with amplified value I_{1L} . This occurs in the presence of the pump wave at ω_2 that has initial intensity I_{20} and that emerges with depleted intensity I_{2L} .

Back-conversion from ω_1 to ω_2 (the inverse Raman effect) will also occur once the Stokes wave reaches sufficient intensity, but gain will only occur for the Stokes wave. Both processes are phase matched, and so occur with high efficiency in long fibers. The back-conversion process is to be distinguished from anti-Stokes scattering, in which pump light at ω_2 is upshifted to frequency ω_3 , with the additional energy being supplied by optical phonons associated with the previously excited medium. The anti-Stokes process is rarely seen in fiber transmission because (1) it is phase mismatched and (2) it requires a substantial population of excited oscillators, which is not the case at thermal equilibrium.

Figure 2 shows the measured Raman gain for the Stokes wave in fused silica. The gain is plotted as a function of difference frequency between the interacting waves measured in cm^{-1} (to convert this to wavelength shift, use the formula $\Delta\lambda = \lambda_p^2 \Delta f (\text{cm}^{-1})$, where λ_p is the pump wavelength). Other fiber constituents such as GeO_2 , P_2O_5 , and B_2O_3 exhibit their own Raman resonances, which occur at successively greater wavelength shifts;¹⁹ the effects of these will be weak, since their concentration in the fiber is generally small. Thus the dominant Raman shifts in optical fiber are associated with SiO_2 , and occur within the range of 440 to 490 cm^{-1} , as is evident in Fig. 2.

Nonlinear polarizations at frequencies ω_1 and ω_2 can be constructed that are proportional to products of the Stokes and pump fields, $E_1^{\omega_1}$ and $E_2^{\omega_2}$. These are of the form $P_{\text{NL}}^{\omega_1} \propto |E_2^{\omega_2}|^2 E_1^{\omega_1}$ (Stokes

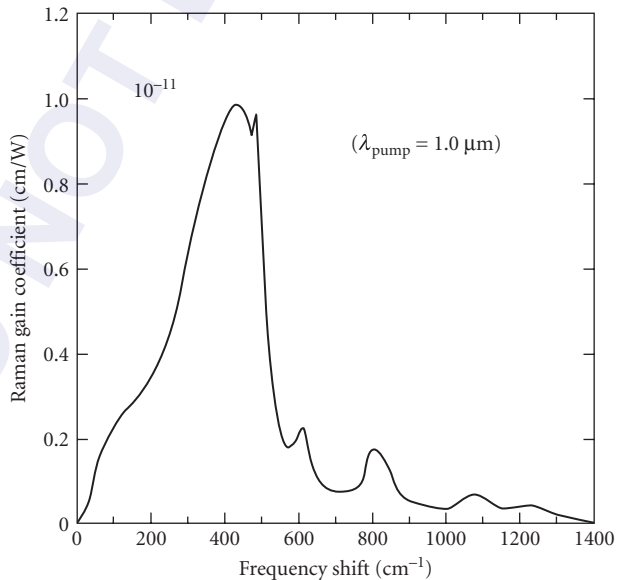


FIGURE 2 Raman gain spectrum in fused silica. (Adapted from Ref. 22. © 1980 IEEE.)

generation) and $P_{\text{NL}}^{\omega_2} \propto |E_1^{\omega_1}|^2 E_2^{\omega_2}$ (the inverse Raman effect). Substituting these polarizations and the two fields into the wave equation, using the SVEA, and assuming copolarized fields leads to the following coupled equations involving the Stokes and pump wave intensities I_1 and I_2 :¹⁸

$$\frac{dI_1}{dz} = g_r I_1 I_2 - \alpha I_1 \quad (10)$$

$$\frac{dI_2}{dz} = -\frac{\omega_2}{\omega_1} g_r I_1 I_2 - \alpha I_2 \quad (11)$$

where the loss terms involving α (the fiber loss per unit distance) are added phenomenologically. The Raman gain function g_r is expressed in a general way as

$$g_r = \frac{A}{\lambda_2} f(\lambda_1 - \lambda_2) \quad (12)$$

where A is a function of the material parameters and $f(\lambda_1 - \lambda_2)$ is a normalized line shape function, which is either derived from theory or experimentally measured (determined from Fig. 2, for example). With λ_2 expressed in μm , $A = 1.0 \times 10^{-11} \text{ cm} - \mu\text{m}/\text{W}$.²⁰ The solutions of Eqs. (10) and (11) are:

$$I_1(z) = \frac{\omega_1}{\omega_2} I_0 \exp(-\alpha z) \frac{\psi_r}{1 + \psi_r} \quad (13)$$

$$I_2(z) = I_0 \exp(-\alpha z) \frac{1}{1 + \psi_r} \quad (14)$$

In these equations, $I_0 = I_{20} + (\omega_1/\omega_2)I_{10}$, where I_{10} and I_{20} are the Stokes and pump intensities at the fiber input. The coupling parameter ψ_r assumes different forms, depending upon whether the input Stokes intensity I_{10} is present or not. If I_{10} is present, and if its magnitude is much greater than light from spontaneous Raman scattering, we have:

$$\psi_r = \frac{\omega_1}{\omega_2} \frac{I_{10}}{I_{20}} \exp(G_0) \quad (15)$$

When no Stokes input is present, the signal builds up from spontaneous Raman scattering, and the coupling parameter in this case becomes:

$$\psi_r = \frac{\hbar \omega_2 \Delta \omega_r}{4\sqrt{\pi}} \frac{1}{I_{20} A_{\text{eff}}} G_2^{-1/2} \exp(G_2) \quad (16)$$

with the gain parameters defined through $G_0 = g_r I_0 L_{\text{eff}}$ and $G_2 = g_r I_{20} L_{\text{eff}}$. The effective length of the fiber accounts for the reduction of Stokes and pump intensities as a result of loss, and is defined as

$$L_{\text{eff}} = \int_0^L \exp(-\alpha z) dz = \frac{1 - \exp(-\alpha L)}{\alpha} \quad (17)$$

The effective area of a single-mode fiber A_{eff} calculated through πr_0^2 , where r_0 is the mode field radius. For a multimode fiber, A_{eff} is usually taken as the core area, assuming that the power is uniformly distributed over the core. The power in the fiber is then $P_{1,2} = I_{1,2} A_{\text{eff}}$.

Two basic issues concerning SRS are of interest in fiber communication systems. First, pump-to-Stokes coupling provides a mechanism for crosstalk from short- to long-wavelength channels. This will occur most efficiently if the channel frequency spacing is in the vicinity of that associated with the maximum Raman gain. The Raman gain peak at approximately 500 cm^{-1} corresponds to a frequency spacing of 15 THz, meaning that operation at 1.55 μm produces a Stokes wave of about 1.67 μm

wavelength. Two-channel operation at these wavelengths would lead to a maximum allowable signal level of about 50 mW.²¹ In WDM systems, within the 1.53- to 1.56- μm erbium-doped fiber amplifier window, channel spacings of 50 or 100 GHz are used. Raman gain is thus considerably reduced, but is still sufficient to cause appreciable crosstalk, which can lead to system penalties of between 1 and 3 dB depending on the number of channels.²² Second, and of more importance to single-wavelength systems, is the conversion to Stokes power from the original signal—a mechanism by which signal power can be depleted. A related problem is walkoff²³ occurring between the signal and Stokes pulses, since these will have different group delays. Walkoff is a means for aliasing to occur in digital transmission, unless the signal is filtered at the output. If pulses are of subpicosecond widths, additional complications arise due to the increased importance of SPM and XPM.²⁴ In any event, an upper limit must be placed on the signal power if significant conversion to Stokes power is to be avoided. In single-wavelength systems, where crosstalk is not an issue, pulse peak powers must be kept below about 500 mW to avoid significant SRS conversion.²⁵

A useful criterion is the so-called critical condition (or Raman threshold), defined as the condition under which the output Stokes and signal powers are equal. This occurs when $\psi_r = 1$, which, from Eq. (16), leads to $G_s \approx 16$. SRS can also be weakened by taking advantage of the gain reduction that occurs as signal (pump) wavelengths increase, as shown in Eq. (12). For example, operation at 1.55 μm yields less SRS for a given signal power than operation at 1.3 μm .

Apart from the need to reduce SRS, the effect can be used to advantage in wavelength conversion and in amplification. Fiber Raman lasers have proven to be good sources of tunable radiation and operate at multiple Stokes wavelengths.²⁶ Specifically, a Stokes wave can serve as a pump to generate an additional (higher-order) Stokes wave at a longer wavelength.²⁷ Fiber Raman amplifiers are used routinely in long-haul systems.²⁸

10.4 STIMULATED BRILLOUIN SCATTERING

The stimulated Brillouin scattering (SBS) process involves the input of a single intense optical wave at frequency ω_2 , which initiates a copropagating acoustic wave at frequency ω_p . The acoustic wave is manifested as a traveling index grating in the fiber, which back-diffracts a portion of the original input. The backward (Stokes) wave is Doppler-shifted to a lower frequency ω_1 and is proportional to the phase conjugate of the input.²⁹ The backward wave is amplified as it propagates, with the gain increasing with increasing input (pump) power.

The beam interaction geometry is shown in Fig. 3. Usually, the Stokes wave builds up spontaneously, but can be inputted at the far end. The effect can be understood by considering a case in which counter-propagating Stokes and pump waves exist that together form a moving interference pattern whose velocity is proportional to the difference frequency $\omega_2 - \omega_1$. Coupling between the waves will occur via SBS when the interference pattern velocity is in the vicinity of the acoustic wave velocity v_p . It is the interference pattern that forms and reinforces the acoustic wave through electrostriction. With a single input, spontaneous scattering from numerous shock waves occurs, with preferential feedback from the acoustic wave that matches the condition just described. With the Stokes wave generated (although it is initially weak), the acoustic wave is reinforced, and so backscattering increases.

In terms of the wave vector magnitudes, the condition for phase matching is given by $k_p = k_1 + k_2$. Since the sound frequency is much less than those of the two optical waves, we can write $k_p \approx 2k_2$. Then, since $k_p = \omega_p/v_p$, it follows that $\omega_p \approx 2n\omega_2 v_p/c$, where n is the refractive index (assumed to

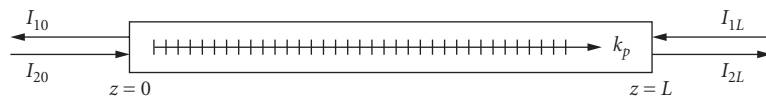


FIGURE 3 Beam geometry for stimulated Brillouin scattering in an optical fiber.

be the same value at both optical frequencies). The Brillouin frequency shift under phase-matched conditions thus becomes

$$\omega_2 - \omega_1 \approx 2n\omega_2 \frac{v_p}{c} \quad (18)$$

This yields a value of about 11 GHz, with $v_p \approx 6$ km/s in fused silica and $\lambda_2 = 1.55 \mu\text{m}$.

The process can be described by the nonlinear polarization produced by the product of complex fields, E_1 , E_2^* , and E_2 ; this yields a polarization at ω_1 that propagates with wavevector k_1 in the direction of the Stokes wave. Another polarization, describing back-coupling from Stokes to pump, involves the product $E_1 E_1^* E_2$. Substituting fields and polarizations into the wave equation yields the following coupled equations that describe the evolution of the optical intensities with distance (pp. 287–290 of Ref. 18):

$$\frac{dI_1}{dz} = -g_b I_1 I_2 + \alpha I_1 \quad (19)$$

$$\frac{dI_2}{dz} = -g_b I_1 I_2 - \alpha I_2 \quad (20)$$

where α is the linear loss coefficient. The Brillouin gain is given by

$$g_b = g_{b0} \left(1 + \frac{4(\omega_1 - \omega_{10})^2}{v_p^2 \alpha_p^2} \right)^{-1} \quad (21)$$

where ω_{10} is the Stokes frequency at precise phase matching, α_p is the loss coefficient for the acoustic wave, and the peak gain g_{b0} is a function of the material parameters. The Brillouin line width, defined as the full width at half-maximum of g_b , is $\Delta\omega_b = v_p \alpha_p$. In optical fibers, $\Delta f_b = \Delta\omega_b / 2\pi$ is typically between 10 and 30 MHz and $g_{b0} = 4.5 \times 10^{-9} \text{ cm/W}$.²⁰ Signal bandwidths in high-data-rate communication systems greatly exceed the Brillouin line width, and so SBS is typically too weak to be considered a source of noise or signal depletion. This is to be compared to stimulated Raman scattering, which supports considerable gain over approximately 5 THz. Consequently, SRS is a much more serious problem in high-data-rate systems.

Using analysis methods similar to those employed in SRS, a critical condition (or threshold) can be defined for SBS, at which the backscattered power is equal to the input power:³⁰

$$\frac{\omega_1 k_B T \Delta\omega_b}{4\sqrt{\pi} \omega_p I_{20} A_{\text{eff}}} G_b^{-3/2} \exp(G_b) = 1 \quad (22)$$

where k_B is the Boltzmann's constant and T is the temperature in Kelvin. The gain parameter is:

$$G_b = g_b I_{20} L_{\text{eff}} \quad (23)$$

with L_{eff} as defined in Eq. (17). Equation (22) is approximately satisfied when $G_b \approx 21$.³⁰ In practice, the backscattered power will always be less than the input power, since pump depletion will occur. Nevertheless, this condition is used as a benchmark to determine the point at which SBS becomes excessive in a given system.³¹ In one study, it was found that $G_b \approx 21$ yields the pump power required to produce an SBS output that is at the level of Rayleigh back-scattering.³² Pump powers required to achieve threshold can be on the order of a few milliwatts for CW or narrowband signals, but these increase substantially for broadband signals.³³ Reduction of SBS is accomplished in practice by

lowering the input signal power (I_{20}) or by taking advantage of the reduction in g_b that occurs when signal bandwidths ($\Delta\omega$) exceed the Brillouin line width. Specifically, if $\Delta\omega \gg \Delta\omega_b$,

$$g_b(\Delta\omega) \approx g_b \frac{\Delta\omega_b}{\Delta\omega} \quad (24)$$

10.5 FOUR-WAVE MIXING

The term *four-wave mixing* in fibers is generally applied to wave coupling through the electronic nonlinearity in which at least two frequencies are involved and in which frequency conversion is occurring. The fact that the electronic nonlinearity is involved distinguishes four-wave mixing interactions from stimulated scattering processes because in the latter the medium was found to play an active role through the generation or absorption of optical phonons (in SRS) or acoustic phonons (in SBS). If the nonlinearity is electronic, bound electron distributions are modified according to the instantaneous optical field configurations. For example, with light at two frequencies present, electron positions can be modulated at the difference frequency, thus modulating the refractive index. Additional light will encounter the modulated index and can be up- or downshifted in frequency. In such cases, the medium plays a passive role in the interaction, as it does not absorb applied energy or release energy previously stored. The self- and cross-phase modulation processes also involve the electronic nonlinearity, but in those cases, power conversion between waves is not occurring—only phase modulation.

As an illustration of the process, consider the interaction of two strong waves at frequencies ω_1 and ω_2 , which mix to produce a downshifted (Stokes) wave at ω_3 and an upshifted (anti-Stokes) wave at ω_4 . The frequencies have equal spacing, that is, $\omega_1 - \omega_3 = \omega_2 - \omega_1 = \omega_4 - \omega_2$ (Fig. 4). All fields assume the real form:

$$\mathcal{E}_j = \frac{1}{2} E_{oj} \exp[i(\omega_j t - \beta_j z) + c.c.], j = 1 - 4 \quad (25)$$

The nonlinear polarization will be proportional to \mathcal{E}^3 , where $\mathcal{E} = \mathcal{E}_1 + \mathcal{E}_2 + \mathcal{E}_3 + \mathcal{E}_4$. With all fields copolarized, complex nonlinear polarizations at ω_3 and ω_4 appear that have the form:

$$P_{NL}^{\omega_3} = \frac{3}{4} \epsilon_0 \chi^{(3)} E_{o1}^2 E_{o2}^* \exp[i(2\omega_1 - \omega_2)t] \exp[-i(2\beta^{\omega_1} - \beta^{\omega_2})z] \quad (26)$$

$$P_{NL}^{\omega_4} = \frac{3}{4} \epsilon_0 \chi^{(3)} E_{o2}^2 E_{o1}^* \exp[i(2\omega_2 - \omega_1)t] \exp[-i(2\beta^{\omega_2} - \beta^{\omega_1})z] \quad (27)$$

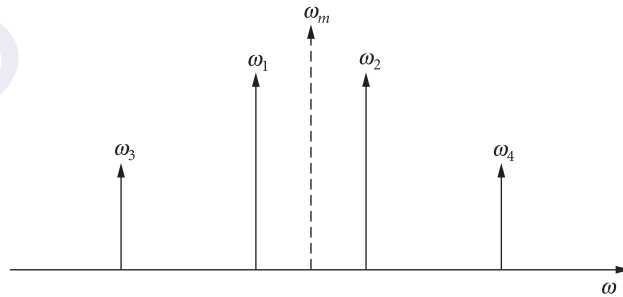


FIGURE 4 Frequency diagram for four-wave mixing, showing pump frequencies (ω_1 and ω_2) and sideband frequencies (ω_3 and ω_4).

where $\omega_3 = 2\omega_1 - \omega_2$, $\omega_4 = 2\omega_1 - \omega_1$, and $\chi^{(3)}$ is proportional to the nonlinear refractive index n_2' . The significance of these polarizations lies not only in the fact that waves at the sideband frequencies ω_3 and ω_4 can be generated, but that preexisting waves at those frequencies can experience gain in the presence of the two pump fields at ω_1 and ω_2 , thus forming a parametric amplifier. The sideband waves will contain the amplitude and phase information on the pumps, thus making this process an important crosstalk mechanism in multiwavelength communication systems. Under phase-matched conditions, the gain associated with FWM is more than twice the peak gain in SRS.³⁴

The wave equation, when solved in steady state, yields the output intensity at either one of the sideband frequencies.³⁵ For a medium of length L , having loss coefficient α , the sideband intensities are related to the pump intensities through

$$I^{\omega_3} \propto \left(\frac{n_2 L_{\text{eff}}}{\lambda_m} \right)^2 I^{\omega_2} (I^{\omega_1})^2 \eta \exp(-\alpha L) \quad (28)$$

$$I^{\omega_4} \propto \left(\frac{n_2 L_{\text{eff}}}{\lambda_m} \right)^2 I^{\omega_1} (I^{\omega_2})^2 \eta \exp(-\alpha L) \quad (29)$$

where L_{eff} is defined in Eq. (17), and where

$$\eta = \frac{\alpha^2}{\alpha^2 + \Delta\beta^2} \left(1 + \frac{4 \exp(-\alpha L) \sin^2(\Delta\beta L/2)}{(1 - \exp(-\alpha L))^2} \right) \quad (30)$$

Other FWM interactions can occur, involving products of intensities at three different frequencies rather than two as demonstrated here. In such cases, the output wave intensities are increased by a factor of 4 over those indicated in Eqs. (28) and (29).

One method of suppressing four-wave mixing in WDM systems includes the use of unequal channel spacing.³⁶ This ensures, for example, that $\omega_3 \neq 2\omega_1 + \omega_2$, where ω_1 , ω_2 , and ω_3 are assigned channel frequencies. More common methods involve phase-mismatching the process in some way. This is accomplished by increasing $\Delta\beta$, which has the effect of decreasing η in Eqs. (28) and (29). Note that in the low-loss limit, where $\alpha \rightarrow 0$, Eq. (30) reduces to $\eta = (\sin^2(\Delta\beta L/2) / (\Delta\beta L/2)^2)$. The $\Delta\beta$ expressions associated with wave generation at ω_3 and ω_4 are given by

$$\Delta\beta(\omega_3) = 2\beta^{\omega_1} - \beta^{\omega_2} - \beta^{\omega_3} \quad (31)$$

and

$$\Delta\beta(\omega_4) = 2\beta^{\omega_2} - \beta^{\omega_1} - \beta^{\omega_4} \quad (32)$$

It is possible to express Eqs. (31) and (32) in terms of known fiber parameters by using a Taylor series for the propagation constant, where the expansion is about frequency ω_m as indicated in Fig. 4, where $\omega_m = (\omega_2 + \omega_1)/2$

$$\beta \approx \beta_0 + (\omega - \omega_m)\beta_1 + \frac{1}{2}(\omega - \omega_m)^2\beta_2 + \frac{1}{6}(\omega - \omega_m)^3\beta_3 \quad (33)$$

In Eq. (33), β_1 , β_2 , and β_3 are, respectively, the first, second, and third derivatives of β with respect to ω , evaluated at ω_m . These in turn relate to the fiber dispersion parameter D (ps/nm · km) and its first derivative with respect to wavelength through $\beta_2 = -(\lambda_m^2/2\pi c)D(\lambda_m)$ and $\beta_3 = (\lambda_m^3/2\pi^2 c^2)[D(\lambda_m) + (\lambda_m/2)(dD/d\lambda)|_{\lambda_m}]$ where $\lambda_m = 2\pi c/\omega_m$. Using these relations along with Eq. (33) in Eqs. (31) and (32) results in:

$$\Delta\beta(\omega_3, \omega_4) \approx 2\pi c \frac{\Delta\lambda^2}{\lambda_m^2} \left[D(\lambda_m) \pm \frac{\Delta\lambda}{2} \frac{dD}{d\lambda} \Big|_{\lambda_m} \right] \quad (34)$$

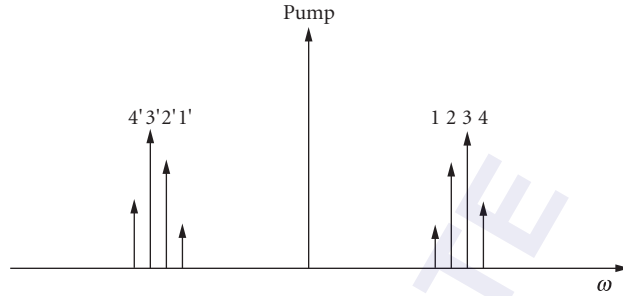


FIGURE 5 Frequency diagram for spectral inversion using four-wave mixing with a single pump frequency.

where the plus sign is used for $\Delta\beta(\omega_3)$, the minus sign is used for $\Delta\beta(\omega_4)$, and $\Delta\lambda = \lambda_1 - \lambda_2$. Phase matching is not completely described by Eq. (34), since cross-phase modulation plays a subtle role, as discussed in Ref. 16. Nevertheless, Eq. (34) does show that the retention of moderate values of dispersion D is a way to reduce FWM interactions that would occur, for example, in WDM systems. As such, modern commercial fiber intended for use in WDM applications will have values of D that are typically in the vicinity of 4 ps/nm · km.³⁷ With WDM operation in conventional dispersion-shifted fiber (with the dispersion zero near 1.55 μm), having a single channel at the zero dispersion wavelength can result in significant four-wave mixing.³⁸ Methods that were found to reduce four-wave mixing in such cases include the use of cross-polarized signals in dispersion-managed links³⁹ and operation within a longer-wavelength band near 1.6 μm ⁴⁰ at which dispersion is appreciable and where gain-shifted fiber amplifiers are used.⁴¹

Examples of other cases involving four-wave mixing include single-wavelength systems, in which the effect has been successfully used in a demultiplexing technique for TDM signals.⁴² In another case, coupling through FWM can occur between a signal and broadband amplified spontaneous emission (ASE) in links containing erbium-doped fiber amplifiers.⁴³ As a result, the signal becomes spectrally broadened and exhibits phase noise from the ASE. The phase noise becomes manifested as amplitude noise under the action of dispersion, producing a form of modulation instability.

An interesting application of four-wave mixing is spectral inversion. Consider a case that involves the input of a strong single-frequency pump wave along with a relatively weak wave having a spectrum of finite width positioned on one side of the pump frequency. Four-wave mixing leads to the generation of a wave whose spectrum is the “mirror image” of that of the weak wave, in which the mirroring occurs about the pump frequency. Figure 5 depicts a representation of this, where four frequency components comprising a spectrum are shown along with their imaged counterparts. An important application of this is pulses that have experienced broadening with chirping after propagating through a length of fiber exhibiting linear group dispersion.⁴⁴ Inverting the spectrum of such a pulse using four-wave mixing has the effect of reversing the direction of the chirp (although the pulse center wavelength is displaced to a different value). When the spectrally inverted pulse is propagated through an additional length of fiber having the same dispersive characteristics, the pulse will compress to nearly its original input width. Compensation for nonlinear distortion has also been demonstrated using this method.⁴⁵

10.6 CONCLUSION

An overview of fiber nonlinear effects has been presented here in which emphasis is placed on the basic concepts, principles, and perspectives on communication systems. Space is not available to cover the more subtle details of each effect or the interrelations between effects that often occur. The

text by Agrawal¹⁶ is recommended for further in-depth study, which should be supplemented by the current literature. Nonlinear optics in fibers and in fiber communication systems comprises an area whose principles and implications are still not fully understood. It thus remains an important area of current research.

10.7 REFERENCES

1. L. F. Mollenauer and P. V. Mamyshev, "Massive Wavelength-Division Multiplexing with Solitons," *IEEE J. Quantum Electron.* **34**:2089–2102 (1998).
2. C. M. de Sterke, B. J. Eggleton, and P. A. Krug, "High-Intensity Pulse Propagation in Uniform Gratings and Grating Superstructures," *IEEE J. Lightwave Technol.* **15**:1494–1502 (1997).
3. L. Clark, A. A. Klein, and D. W. Peckham, "Impact of Fiber Selection and Nonlinear Behavior on Network Upgrade Strategies for Optically Amplified Long Interoffice Routes," *Proceedings of the 10th Annual National Fiber Optic Engineers Conference*, vol. 4, 1994.
4. Y. R. Shen, *The Principles of Nonlinear Optics*, Wiley-Interscience, New York, 1984, p. 28.
5. R. H. Stolen and H. W. K. Tom, "Self-Organized Phase-Matched Harmonic Generation in Optical Fibers," *Opt. Lett.* **12**:585–587 (1987).
6. R. W. Boyd, *Nonlinear Optics*, Academic Press, San Diego, 2008.
7. R. H. Stolen and C. Lin, "Self-Phase Modulation in Silica Optical Fibers," *Phys. Rev. A* **17**:1448–1453 (1978).
8. G. Bellotti, A. Bertaina, and S. Bigo, "Dependence of Self-Phase Modulation Impairments on Residual Dispersion in 10-Gb/s-Based Terrestrial Transmission Using Standard Fiber," *IEEE Photon. Technol. Lett.* **11**:824–826 (1999).
9. M. Stern, J. P. Heritage, R. N. Thurston, and S. Tu, "Self-Phase Modulation and Dispersion in High Data Rate Fiber Optic Transmission Systems," *IEEE J. Lightwave Technol.* **8**:1009–1015 (1990).
10. D. Marcuse and C. R. Menyuk, "Simulation of Single-Channel Optical Systems at 100 Gb/s," *IEEE J. Lightwave Technol.* **17**:564–569 (1999).
11. S. Reichel and R. Zengerle, "Effects of Nonlinear Dispersion in EDFA's on Optical Communication Systems," *IEEE J. Lightwave Technol.* **17**:1152–1157 (1999).
12. M. Karlsson, "Modulational Instability in Lossy Optical Fibers," *J. Opt. Soc. Am. B* **12**:2071–2077 (1995).
13. R. Hui, K. R. Demarest, and C. T. Allen, "Cross-Phase Modulation in Multispan WDM Optical Fiber Systems," *IEEE J. Lightwave Technol.* **17**:1018–1026 (1999).
14. S. Bigo, G. Billotti, and M. W. Chbat, "Investigation of Cross-Phase Modulation Limitation over Various Types of Fiber Infrastructures," *IEEE Photon. Technol. Lett.* **11**:605–607 (1999).
15. L. F. Mollenauer and J. P. Gordon, *Solitons in Optical Fibers*, Academic Press, Boston, 2006.
16. G. P. Agrawal, *Nonlinear Fiber Optics*, 4th ed., Academic Press, San Diego, 2006.
17. G. Herzberg, *Infra-Red and Raman Spectroscopy of Polyatomic Molecules*, Van Nostrand, New York, 1945, pp. 99–101.
18. J. A. Buck, *Fundamentals of Optical Fibers*, 2nd ed., Wiley-Interscience, Hoboken, 2004.
19. F. L. Galeener, J. C. Mikkelsen Jr., R. H. Geils, and W. J. Mosby, "The Relative Raman Cross Sections of Vitreous SiO₂, GeO₂, B₂O₃, and P₂O₅," *Appl. Phys. Lett.* **32**:34–36 (1978).
20. R. H. Stolen, "Nonlinear Properties of Optical Fibers," in S. E. Miller and A. G. Chynoweth (eds.), *Optical Fiber Telecommunications*, Academic Press, New York, 1979.
21. A. R. Chraplyvy, "Optical Power Limits in Multi-Channel Wavelength Division Multiplexed Systems due to Stimulated Raman Scattering," *Electron. Lett.* **20**:58–59 (1984).
22. F. Forghieri, R. W. Tkach, and A. R. Chraplyvy, "Effect of Modulation Statistics on Raman Crosstalk in WDM Systems," *IEEE Photon. Technol. Lett.* **7**:101–103 (1995).
23. R. H. Stolen and A. M. Johnson, "The Effect of Pulse Walkoff on Stimulated Raman Scattering in Fibers," *IEEE J. Quantum Electron.* **22**:2154–2160 (1986).
24. C. H. Headley III and G. P. Agrawal, "Unified Description of Ultrafast Stimulated Raman Scattering in Optical Fibers," *J. Opt. Soc. Am. B* **13**:2170–2177 (1996).

25. R. H. Stolen, J. P. Gordon, W. J. Tomlinson, and H. A. Haus, "Raman Response Function of Silica Core Fibers," *J. Opt. Soc. Am. B* **6**:1159–1166 (1988).
26. L. G. Cohen and C. Lin, "A Universal Fiber-Optic (UFO) Measurement System Based on a Near-IR Fiber Raman Laser," *IEEE J. Quantum Electron.* **14**:855–859 (1978).
27. K. X. Liu and E. Garmire, "Understanding the Formation of the SRS Stokes Spectrum in Fused Silica Fibers," *IEEE J. Quantum Electron.* **27**:1022–1030 (1991).
28. J. Bromage, "Raman Amplification for Fiber Communication Systems," *IEEE J. Lightwave Technol.* **22**:79–93 (2004).
29. A. Yariv, *Quantum Electronics*, 3d ed., Wiley, New York, 1989, pp. 513–516.
30. R. G. Smith, "Optical Power Handling Capacity of Low Loss Optical Fibers as Determined by Stimulated Raman and Brillouin Scattering," *Appl. Opt.* **11**:2489–2494 (1972).
31. A. R. Chraplyvy, "Limitations on Lightwave Communications Imposed by Optical Fiber Nonlinearities," *IEEE J. Lightwave Technol.* **8**:1548–1557 (1990).
32. X. P. Mao, R. W. Tkach, A. R. Chraplyvy, R. M. Jopson, and R. M. Derosier, "Stimulated Brillouin Threshold Dependence on Fiber Type and Uniformity," *IEEE Photon. Technol. Lett.* **4**:66–68 (1992).
33. C. Edge, M. J. Goodwin, and I. Bennion, "Investigation of Nonlinear Power Transmission Limits in Optical Fiber Devices," *Proc. IEEE* **134**:180–182 (1987).
34. R. H. Stolen, "Phase-Matched Stimulated Four-Photon Mixing in Silica-Fiber Waveguides," *IEEE J. Quantum Electron.* **11**:100–103 (1975).
35. R. W. Tkach, A. R. Chraplyvy, F. Forghieri, A. H. Gnauck, and R. M. Derosier, "Four-Photon Mixing and High-Speed WDM Systems," *IEEE J. Lightwave Technol.* **13**:841–849 (1995).
36. F. Forghieri, R. W. Tkach, and A. R. Chraplyvy, and D. Marcuse, "Reduction of Four-Wave Mixing Crosstalk in WDM Systems Using Unequally-Spaced Channels," *IEEE Photon. Technol. Lett.* **6**:754–756 (1994).
37. AT&T Network Systems data sheet 4694FS-Issue 2 LLC, "TrueWave Single Mode Optical Fiber Improved Transmission Capacity," December 1995.
38. D. Marcuse, A. R. Chraplyvy, and R. W. Tkach, "Effect of Fiber Nonlinearity on Long-Distance Transmission," *IEEE J. Lightwave Technol.* **9**:121–128 (1991).
39. E. A. Golovchenko, N. S. Bergano, and C. R. Davidson, "Four-Wave Mixing in Multispan Dispersion Managed Transmission Links," *IEEE Photon. Technol. Lett.* **10**:1481–1483 (1998).
40. M. Jinno et al, "1580 nm Band, Equally-Spaced 8×10 Gb/s WDM Channel Transmission Over 360 km (3×120 km) of Dispersion-Shifted Fiber Avoiding FWM Impairment," *IEEE Photon. Technol. Lett.* **10**:454–456 (1998).
41. H. Ono, M. Yamada, and Y. Ohishi, "Gain-Flattened Er^{3+} -Fiber Amplifier for A WDM Signal in the 1.57–1.60 μm Wavelength Region," *IEEE Photon. Technol. Lett.* **9**:596–598 (1997).
42. P. O. Hedekvist, M. Karlsson, and P. A. Andrekson, "Fiber Four-Wave Mixing Demultiplexing with Inherent Parametric Amplification," *IEEE J. Lightwave Technol.* **15**:2051–2058 (1997).
43. R. Hui, M. O'Sullivan, A. Robinson, and M. Taylor, "Modulation Instability and Its Impact on Multispan Optical Amplified IMDD Systems: Theory and Experiments," *IEEE J. Lightwave Technol.* **15**:1071–1082 (1997).
44. A. H. Gnauck, R. M. Jopson, and R. M. Derosier, "10 Gb/s 360 km Transmission over Dispersive Fiber Using Midsystem Spectral Inversion," *IEEE Photon. Technol. Lett.* **5**:663–666 (1993).
45. A. H. Gnauck, R. M. Jopson, and R. M. Derosier, "Compensating the Compensator: A Demonstration of Nonlinearity Cancellation in a WDM System," *IEEE Photon. Technol. Lett.* **7**:582–584 (1995).

This page intentionally left blank.

DO NOT DUPLICATE

Philip St. J. Russell and Greg J. Pearce

*Max-Planck Institute for the Science of Light
Erlangen, Germany*

11.1 GLOSSARY

A_i	nonlinear effective area of subregion i
c	velocity of light in vacuum
d	hole diameter
D	$\frac{\partial}{\partial \lambda} \frac{1}{v_g} = \frac{\partial^2 \beta_m}{\partial \lambda \partial \omega}$ the group velocity dispersion of mode m in engineering units (ps/nm · km)
k	vacuum wavevector $2\pi/\lambda = \omega/c$
n_2^i	nonlinear refractive index of subregion i
n_{\max}	maximum index supported by the PCF cladding (fundamental space-filling mode)
n_z	z-component of refractive index
n_q^∞	$\sqrt{\sum_i n_i^2 A_i / \sum_i A_i}$ the area-averaged refractive index of an arbitrary region q of a microstructured fiber, where n_i and A_i are respectively the refractive indices and total area of subregion i
V	$k\rho\sqrt{n_{\text{co}}^2 - n_{\text{cl}}^2}$ the normalized frequency for a step-index fiber
V_{gen}	$k\Lambda\sqrt{n_1^2 - n_2^2}$ a generalized form of V for a structure made from two materials of index n_1 and n_2
β	component of wavevector along the fiber axis
β_m	axial wavevector of guided mode
β_{\max}	maximum possible axial component of wavevector in the PCF cladding
Δ	$(n_{\text{co}} - n_{\text{cl}})/n_{\text{cl}}$, where n_{co} and n_{cl} are, respectively, the core and cladding refractive indices of a conventional step-index fiber
Δ_∞	$(n_{\text{co}}^\infty - n_{\text{cl}}^\infty)/n_{\text{cl}}^\infty$, where n_{co}^∞ and n_{cl}^∞ are the refractive indices of core and cladding in the long wavelength limit
γ	$\text{W}^{-1}\text{km}^{-1}$ nonlinear coefficient of an optical fiber
$\epsilon(\mathbf{r}_T)$	relative dielectric constant of a PCF as a function of transverse position $\mathbf{r}_T = (x, y)$
λ	vacuum wavelength
λ_{eff}^i	$2\pi/\sqrt{k^2 n_i^2 - \beta^2}$ the effective transverse wavelength in subregion i

- Λ interhole spacing, period, or pitch
- ρ core radius
- ω angular frequency of light

11.2 INTRODUCTION

Photonic crystal fibers (PCFs)—fibers with a periodic transverse microstructure—have been in practical existence as low loss waveguides since early 1996.^{1–4} The initial demonstration took 4 years of technological development, and since then the fabrication techniques have become more and more sophisticated. It is now possible to manufacture the microstructure in air-glass PCF to accuracies of 10 nm on the scale of 1 μm , which allows remarkable control of key optical properties such as dispersion, birefringence, nonlinearity, and the position and width of the photonic bandgaps (PBGs) in the periodic “photonic crystal” cladding. PCF has in this way extended the range of possibilities in optical fibers, both by improving well-established properties and introducing new features such as low loss guidance in a hollow core.

Standard single-mode telecommunications fiber (SMF), with a normalized core-cladding refractive index difference Δ approximately 0.4 percent, a core radius of ρ approximately 4.5 μm , and of course a very high optical clarity (better than 5 km/dB at 1550 nm), is actually quite limiting for many applications. Two major factors contribute to this. The first is the smallness of Δ , which causes bend loss (0.5 dB at 1550 nm in Corning SMF-28 for one turn around a mandrel 32 mm in diameter¹⁰) and limits the degree to which group velocity dispersion and birefringence can be manipulated. Although much higher values of Δ can be attained (modified chemical vapor deposition yields an index difference of 0.00146 per mol % GeO_2 , up to a maximum $\Delta \sim 10\%$ for 100 mol %¹¹), the single-mode core radius becomes very small and the attenuation rises through increased absorption and Rayleigh scattering. The second factor is the reliance on total internal reflection (TIR), so that guidance in a hollow core is impossible, however useful it would be in fields requiring elimination of glass-related nonlinearities or enhancement of laser interactions with dilute or gaseous media. PCF has made it possible to overcome these limitations, and as a result many new applications of optical fibers are emerging.

Outline of Chapter

In the next section a brief history of PCF is given, and in Sec. 11.4 fabrication techniques are reviewed. Numerical modeling and analysis are covered in Sec. 11.5 and the optical properties of the periodic photonic crystal cladding in Sec. 11.6. The characteristics of guidance are discussed in Sec. 11.7. In Sec. 11.8 the nonlinear characteristics of guidance are reviewed and intrafiber devices (including cleaving and splicing) are discussed in Sec. 11.9. Brief conclusions, including a list of applications, are drawn in Sec. 11.10.

11.3 BRIEF HISTORY

The original motivation for developing PCF was the creation of a new kind of dielectric waveguide—one that guides light by means of a two-dimensional PBG. In 1991, the idea that the well-known “stop-bands” in multiply periodic structures (for a review see Ref. 12) could be extended to eliminate all photonic states¹³ was leading to attempts worldwide to fabricate three-dimensional PBG materials. At that time the received wisdom was that the refractive index difference needed to create a PBG

in two dimensions was large—of order 2.2:1. It was not widely recognized that the refractive index difference requirements for PBG formation in two dimensions are greatly relaxed if, as in a fiber, propagation is predominantly along the third axis—the direction of invariance.

Photonic Crystal Fibers

The original 1991 idea, then, was to trap light in a hollow core by means of a two-dimensional “photonic crystal” of microscopic air capillaries running along the entire length of a glass fiber.¹⁴ Appropriately designed, this array would support a PBG for incidence from air, preventing the escape of light from a hollow core into the photonic crystal cladding and avoiding the need for TIR.

The first 4 years of work on understanding and fabricating PCF were a journey of exploration. The task of solving Maxwell’s equations numerically made good progress, culminating in a 1995 paper that showed that photonic bandgaps did indeed exist in two-dimensional silica-air structures for “conical” incidence from vacuum—this being an essential prerequisite for hollow-core guidance.¹⁵ Developing a suitable fabrication technique took rather longer. After 4 years of trying different approaches, the first successful silica-air PCF structure was made in late 1995 by stacking 217 silica capillaries (8 layers outside the central capillary), specially machined with hexagonal outer and a circular inner cross sections. The diameter-to-pitch ratio d/Λ of the holes in the final stack was approximately 0.2, which theory showed was too small for PBG guidance in a hollow core, so it was decided to make a PCF with a solid central core surrounded by 216 air channels (Fig. 1a).^{2,5,6}

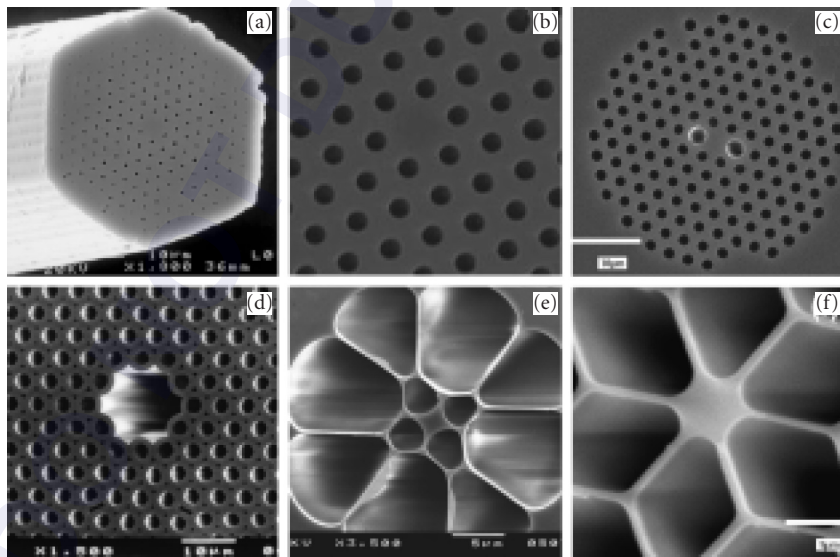


FIGURE 1 Selection of scanning electron micrographs of PCF structures. (a) The first working PCF—the solid glass core is surrounded by a triangular array of 300 nm diameter air channels, spaced 2.3 μm apart;^{5,6} (b) detail of a low loss solid-core PCF (interhole spacing $\sim 2 \mu\text{m}$); (c) birefringent PCF (interhole spacing 2.5 μm); (d) the first hollow-core PCF (core diameter $\sim 10 \mu\text{m}$);⁷ (e) a PCF extruded from Schott SF6 glass with a core approximately 2 μm in diameter;⁸ and (f) PCF with very small core (diameter 800 nm) and zero GVD wavelength 560 nm.⁹

This led to the discovery of endlessly single-mode PCF, which, if it guides at all, only supports the fundamental guided mode.¹⁶ The success of these initial experiments led rapidly to a whole series of new types of PCF—large mode area,¹⁷ dispersion-controlled,^{9,18} hollow core,⁷ birefringent,¹⁹ and multicore.²⁰

These initial breakthroughs led quickly to applications, perhaps the most celebrated being the report in 2000 of supercontinuum generation from unamplified Ti:sapphire fs laser pulses in a PCF with a core small enough to give zero dispersion at 800 nm wavelength (subsection “Supercontinuum Generation” in Sec. 11.8).²¹

Bragg Fibers

In the late 1960s and early 1970s, theoretical proposals were made for another kind of fiber with a periodically structured cross section.^{22,23} This was a cylindrical “Bragg” fiber that confines light within an annular array of rings of high and low refractive index arranged concentrically around a central core. A group in France has made a solid-core version of this structure using modified chemical vapor deposition.²⁴ Employing a combination of polymer and chalcogenide glass, researchers in the United States have realized a hollow-core version of a similar structure,²⁵ reporting 1 dB/m loss at 10 μm wavelength (the losses at telecom wavelengths are as yet unspecified). This structure guides light in the TE_{01} mode, used in microwave telecommunications because of its ultralow loss; the field moves away from the attenuating waveguide walls as the frequency increases, resulting in very low losses, although the guide must be kept very straight to avoid the fields entering the cladding and experiencing high absorption.

11.4 FABRICATION TECHNIQUES

Photonic crystal fiber (PCF) structures are currently produced in many laboratories worldwide using a variety of different techniques (see Fig. 2 for schematic drawings of some example structures). The first stage is to produce a “preform”—a macroscopic version of the planned microstructure in the drawn PCF. There are many ways to do this, including stacking of capillaries and rods,⁵ extrusion,^{8,26–28} sol-gel casting,²⁹ injection molding and drilling. The materials used range from silica to compound glasses, chalcogenide glasses, and polymers.³⁰

The most widely used technique is stacking of circular capillaries (Fig. 3). Typically, meter-length capillaries with an outer diameter of approximately 1 mm are drawn from a starting tube of high-purity synthetic silica with a diameter of approximately 20 mm. The inner/outer diameter of the starting tube, which typically lies in the range from 0.3 up to beyond 0.9, largely determines the d/Λ value in the drawn fiber. The uniformity in diameter and circularity of the capillaries must be controlled to at least 1 percent of the diameter. They are stacked horizontally, in a suitably shaped jig, to form the desired crystalline arrangement. The stack is bound with wire before being inserted into a jacketing tube, and the whole assembly is then mounted in the preform feed unit for drawing down to fiber. Judicious use of pressure and vacuum during the draw allows some limited control over the final structural parameters, for example the d/Λ value.

Extrusion offers an alternative route to making PCF, or the starting tubes, from bulk glass; it permits formation of structures that are not readily made by stacking. While not suitable for silica (no die material has been identified that can withstand the $\sim 2000^\circ\text{C}$ processing temperatures without contaminating the glass), extrusion is useful for making PCF from compound silica glasses, tellurites, chalcogenides, and polymers—materials that melt at lower temperatures. Figure 1e shows the cross section of a fiber extruded, through a metal die, from a commercially available glass (Schott SF6).⁸ PCF has also been extruded from tellurite glass, which has excellent IR transparency out to beyond 4 μm , although the reported fiber losses (a few dB/m) are as yet rather high.^{27,31–33} Polymer PCFs, first developed in Sydney, have been successfully made using many different approaches, for example extrusion, casting, molding, and drilling.³⁰

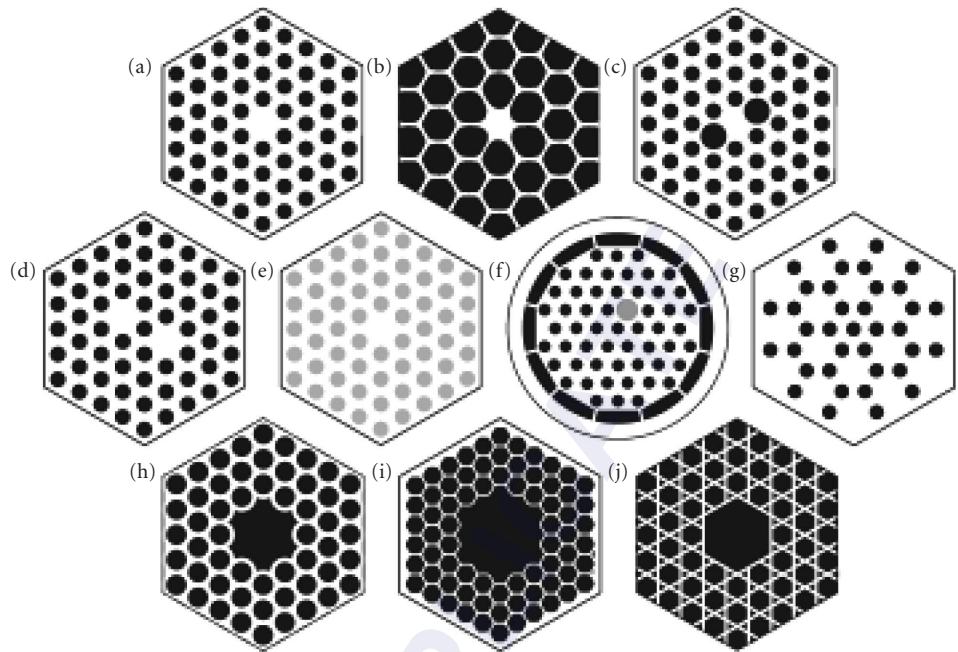


FIGURE 2 Representative sketches of different types of PCF. The black regions are hollow, the white regions are pure glass, and the gray regions doped glass. (a) Endlessly single-mode solid core; (b) highly nonlinear (high air-filling fraction, small core, characteristically distorted holes next to the core); (c) birefringent; (d) dual-core; (e) all-solid glass with raised-index doped glass strands (colored gray) in the cladding; (f) double-clad PCF with off-set doped lasing core and high numerical aperture inner cladding for pumping (the photonic crystal cladding is held in place by thin webs of glass); (g) “carbon-ring” array of holes for PBG guidance in core with extra hole; (h) seven-cell hollow core; (i) 19-cell hollow core with high air-filling fraction (in a real fiber, surface tension smoothes out the bumps on the core surround); and (j) hollow-core with Kagomé lattice in the cladding.

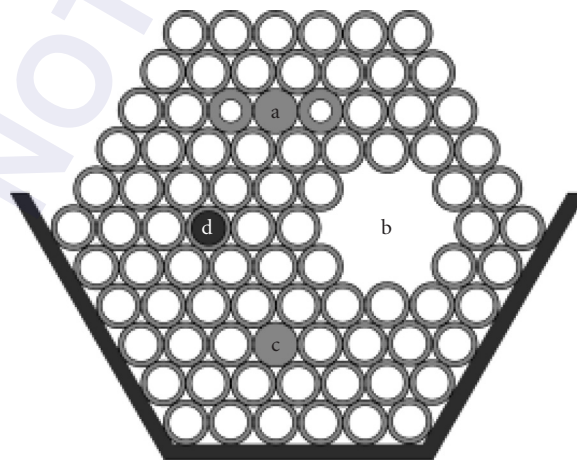


FIGURE 3 Preform stack containing (a) birefringent solid core; (b) seven-cell hollow core; (c) solid isotropic core; and (d) doped core. The capillary diameters are approximately 1 mm—large enough to ensure that they remain stiff for stacking.

Design Approach

The successful design of a PCF for a particular application is not simply a matter of using numerical modeling (see next section) to calculate the parameters of a structure that yields the required performance. This is because the fiber drawing process is not lithographic, but introduces its own highly reproducible types of distortion through the effects of viscous flow, surface tension, and pressure. As a result, even if the initial preform stack precisely mimics the theoretically required structure, several modeling and fabrication iterations are usually needed before a successful design can be reached.

11.5 MODELING AND ANALYSIS

The complex structure of PCF—in particular the large refractive index difference between glass and air—makes its electromagnetic analysis challenging. Maxwell's equations must usually be solved numerically, using one of a number of specially developed techniques.^{15,34–38} Although standard optical fiber analyses and number of approximate models are occasionally helpful, these are only useful as rough guidelines to the exact behavior unless checked against accurate numerical solutions.

Maxwell's Equations

In most practical cases, a set of equal frequency modes is more useful than a set of modes of different frequency sharing the same value of axial wavevector component β . It is therefore convenient to arrange Maxwell's equations with β^2 as eigenvalue

$$\left(\nabla^2 + k^2\epsilon(\mathbf{r}_T) + [\nabla \ln \epsilon(\mathbf{r}_T)] \wedge \nabla \wedge\right) \mathbf{H}_T = \beta^2 \mathbf{H}_T \quad (1)$$

where all the field vectors are taken in the form $\mathbf{Q} = \mathbf{Q}_T(\mathbf{r}_T)e^{-j\beta z}$, $\epsilon_T(\mathbf{r}_T)$ is the dielectric constant, $\mathbf{r}_T = (x, y)$ is position in the transverse plane, and $k = \omega/c$ is the vacuum wavevector. This form allows material dispersion to be easily included, something which is not possible if the equations are set up with k^2 as eigenvalue. Written out explicitly in cartesian coordinates Eq. (1) yields two equations relating h_x and h_y

$$\begin{aligned} \frac{\partial^2 h_x}{\partial y^2} + \frac{\partial^2 h_x}{\partial x^2} - \frac{\partial \ln \epsilon}{\partial y} \left(\frac{\partial h_x}{\partial y} - \frac{\partial h_y}{\partial x} \right) + (\epsilon k^2 - \beta^2) h_x &= 0 \\ \frac{\partial^2 h_y}{\partial x^2} + \frac{\partial^2 h_y}{\partial y^2} + \frac{\partial \ln \epsilon}{\partial x} \left(\frac{\partial h_x}{\partial y} - \frac{\partial h_y}{\partial x} \right) + (\epsilon k^2 - \beta^2) h_y &= 0 \end{aligned} \quad (2)$$

and a third differential equation relating h_x , h_y , and h_z , which is however not required to solve Eq. (2).

Scalar Approximation

In the paraxial scalar approximation the second term inside the operator in Eq. (1), which gives rise to the middle terms in Eq. (2) that couple between the vector components of the field, can be neglected, yielding a scalar wave equation

$$\nabla^2 \mathbf{H}_T + [k^2 \epsilon(\mathbf{r}_T) - \beta^2] \mathbf{H}_T = 0 \quad (3)$$

This leads to a scaling law, similar to the one used in standard analyses of step-index fiber,³⁹ that can be used to parameterize the fields.⁴⁰ Defining Λ as the interhole spacing and n_1 and n_2 the refractive indices of the two materials used to construct a particular geometrical shape of photonic crystal, the mathematical forms of the fields and the dispersion relations are identical provided the generalized V -parameter

$$V_{\text{gen}} = k\Lambda\sqrt{n_1^2 - n_2^2} \quad (4)$$

is held constant. This has the interesting (though in the limit not exactly practical) consequence that bandgaps can exist for vanishingly small index differences, provided the structure is made sufficiently large (see subsection “All-Solid Structures” in Sec. 11.7).

Numerical Techniques

A common technique for solving Eq. (1) employs a Fourier expansion to create a basis set of plane waves for the fields, which reduces the problem to the inversion of a matrix equation, suitable for numerical computation.³⁷ Such an implicitly periodic approach is especially useful for the study of the intrinsic properties of PCF claddings. However, in contrast to versions of Maxwell’s equations with k^2 as eigenvalue,⁴¹ Eq. (1) is non-Hermitian, which means that standard matrix inversion methods for Hermitian problems cannot straightforwardly be applied. An efficient iterative scheme can, however, be used to calculate the inverse of the operator by means of fast Fourier transform steps. This method is useful for accurately finding the modes guided in a solid-core PCF, which are located at the upper edge of the eigenvalue spectrum of the inverted operator. In hollow-core PCF, however (or other fibers relying on a cladding bandgap to confine the light), the modes of interest lie in the interior of the eigenvalue spectrum. A simple transformation can, however, be used to move the desired interior eigenvalues to the edge of the spectrum, greatly speeding up the calculations and allowing as many as a million basis waves to be incorporated.^{42,43} To treat PCFs with a central guiding core in an otherwise periodic lattice, a supercell is constructed, its dimensions being large enough so that, once tiled, the guided modes in adjacent cores do not significantly interact.

The choice of a suitable numerical method often depends on fiber geometry, as some methods can exploit symmetries or regularity of structure to increase efficiency. Other considerations are whether material dispersion is significant (more easily included in fixed-frequency methods), and whether leakage losses or a treatment of leaky modes (requiring suitable boundaries on the computational domain) are desired. If the PCF structure consists purely of circular holes, for example, the multipole or Rayleigh method is a particularly fast and efficient method.^{34,35} It uses Mie theory to evaluate the scattering of the field incident on each hole. Other numerical techniques include expanding the field in terms of Hermite–Gaussian functions,^{38,44} the use of finite-difference time-domain (FDTD) analysis (a simple and versatile tool for exploring waveguide geometries⁴⁵) or finite-difference method in the frequency domain,⁴⁶ and the finite-element approach.⁴⁷ Yet another approach is a source-model technique which uses two sets of fictitious elementary sources to approximate the fields inside and outside circular cylinders.⁴⁸

11.6 CHARACTERISTICS OF PHOTONIC CRYSTAL CLADDING

The simplest photonic crystal cladding is a biaxially periodic, defect-free, composite material with its own well-defined dispersion and band structure. These properties determine the behavior of the guided modes that form at cores (or “structural defects” in the jargon of photonic crystals).

A convenient graphical tool is the propagation diagram—a map of the ranges of frequency and axial wavevector component β where light is evanescent in all transverse directions regardless of its polarization state (Fig. 4).¹⁵ The vertical axis is the normalized frequency $k\Lambda$, and the horizontal axis is the normalized axial wavevector component $\beta\Lambda$. Light is unconditionally cutoff from propagating (due to either TIR or a PBG) in the black regions.

In any subregion of isotropic material (glass or air) at fixed optical frequency, the maximum possible value of $\beta\Lambda$ is given by $k\Lambda n$, where n is the refractive index (at that frequency) of the region under consideration. For $\beta < kn$ light is free to propagate, for $\beta > kn$ it is evanescent and at $\beta = kn$ the critical angle is reached—denoting the onset of TIR for light incident from a medium of index larger than n .

The slanted guidelines (Fig. 4) denote the transitions from propagation to evanescence for air, the photonic crystal, and glass. At fixed optical frequency for $\beta < k$, light propagates freely in every subregion of the structure. For $k < \beta < kn_g$ (n_g is the index of the glass), light propagates in the glass substrands and is evanescent in the hollow regions. Under these conditions the “tight binding” approximation holds, and the structure may be viewed as an array of coupled glass waveguides.

The photonic bandgap “fingers” in Fig. 4 are most conveniently investigated by plotting the photonic density of states (DOS),⁴³ which shows graphically the density of allowed modes in the PCF cladding relative to vacuum (Fig. 5). Regions of zero DOS are photonic bandgaps, and by plotting a quantity such as $(\beta - k)\Lambda$ it is possible to see clearly how far a photonic bandgap extends below the light line.

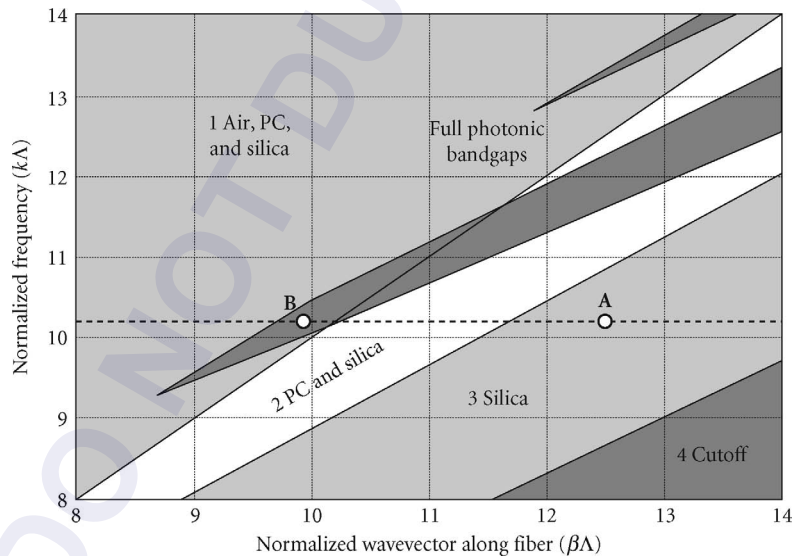


FIGURE 4 Propagation diagram for a triangular array of circular air holes (radius $\rho = 0.47\Lambda$) in silica glass, giving an air-filling fraction of 80 percent. Note the different regions where light is (4) cutoff completely (dark), (3) able to propagate only in silica glass (light gray), (2) able to propagate also in the photonic crystal cladding (white), and (1) able to propagate in all regions (light gray). Guidance by total internal reflection in a silica core is possible at point A. The “fingers” indicate the positions of full two-dimensional photonic bandgaps, which can be used to guide light in air at positions such as B where a photonic bandgap crosses the light line $k = \beta$.

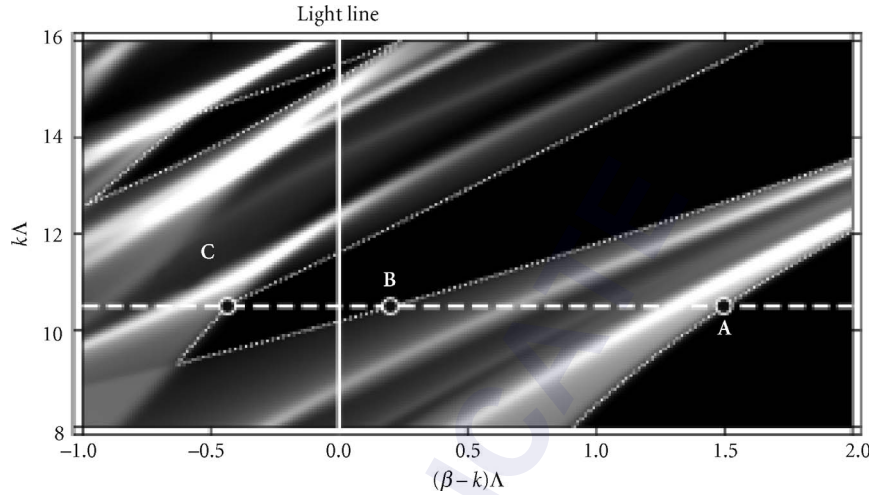


FIGURE 5 Photonic density of states (DOS) for the fiber structure described in Fig. 4, where black regions show zero DOS and lighter regions show higher DOS. The edges of full two-dimensional photonic bandgaps and the band edge of the fundamental space-filling mode are highlighted with thin dotted white lines. The vertical white line is the light line, and the labeled points mark band edges at the frequency of the thick dashed white line, discussed in the section “Maximum Refractive Index and Band Edges.”

Maximum Refractive Index and Band Edges

The maximum axial refractive index $n_{\max} = \beta_{\max}/k$ in the photonic crystal cladding lies in the range $k < \beta < kn_g$ as expected of a composite glass-air material. This value coincides with the z -pointing “peaks” of the dispersion surfaces in reciprocal space, where multiple values of transverse wavevector are allowed, one in each tiled Brillouin zone. For a constant value of β slightly smaller than β_{\max} , these wavevectors lie on small approximately circular loci, with a transverse component of group velocity that points normal to the circles in the direction of increasing frequency. Thus, light can travel in *all* directions in the transverse plane, even though its wavevectors are restricted to values lying on the circular loci. The real-space distribution of this field is shown in Fig. 6, together with the fields at two other band edges.

The maximum axial refractive index n_{\max} depends strongly on frequency, even though neither the air nor the glass are assumed dispersive in the analysis; microstructuring itself creates dispersion, through a balance between transverse energy storage and energy flow that is highly dependent upon frequency.

By averaging the square of the refractive index in the photonic crystal cladding it is simple to show that

$$n_{\max} \rightarrow \sqrt{(1-F)n_g^2 - Fn_a^2} \quad k\Lambda \rightarrow 0 \quad (5)$$

in the long-wavelength limit for a scalar approximation, where F is the air-filling fraction and n_a is the index in the holes (which we take to be 1 for the rest of this subsection).

As the wavelength of the light falls, the optical fields are better able to distinguish between the glass regions and the air. The light piles up more and more in the glass, causing the effective n_{\max} “seen” by

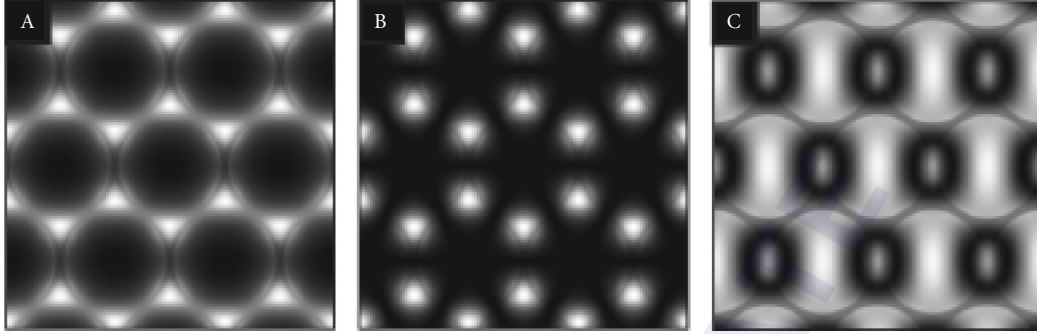


FIGURE 6 Plots showing the magnitude of the axial Poynting vector at band edges A, B, and C as shown in Fig. 5. White regions have large Poynting vector magnitude. A is the fundamental space-filling mode, for which the field amplitudes are in phase between adjacent unit cells. In B (the “dielectric” edge) the field amplitudes change sign between anti-nodes, and in C (the “air” edge) the central lobe has the opposite sign from the six surrounding lobes.

it to change. In the limit of small wavelength $k\Lambda \rightarrow \infty$ light is strongly excluded from the air holes by TIR, and the field profile “freezes” into a shape that is independent of wavelength. The variation of n_{\max} with frequency may be estimated by expanding fields centered on the air holes in terms of Bessel functions and applying symmetry.¹⁶ Defining the normalized parameters

$$u = \Lambda\sqrt{k^2 n_g^2 - \beta^2} \quad v = k\Lambda\sqrt{n_g^2 - 1} \quad (6)$$

the analysis yields the polynomial fit (see Sec. 11.11):

$$u(v) = (0.00151 + 2.62v^{-1} + 0.0155v - 0.000402v^2 + 3.63 \times 10^{-6}v^3)^{-1} \quad (7)$$

for $d/\Lambda = 0.4$ and $n_g = 1.444$. This polynomial is accurate to better than 1 percent in the range $0 < v < 50$. The resulting expression for n_{\max} is plotted in Fig. 7 against the parameter v .

Transverse Effective Wavelength

The transverse effective wavelength in the i th material is defined as follows:

$$\lambda_{\text{eff}}^i = \frac{2\pi}{\sqrt{k^2 n_i^2 - \beta^2}} \quad (8)$$

where n_i is its refractive index. This wavelength can be many times the vacuum value, tending to infinity at the critical angle $\beta \rightarrow kn_i$, and being imaginary when $\beta > kn_i$. It is a measure of whether or not the light is likely to be resonant within a particular feature of the structure, for example, a hole or a strand of glass, and defines PCF as a wavelength-scale structure.

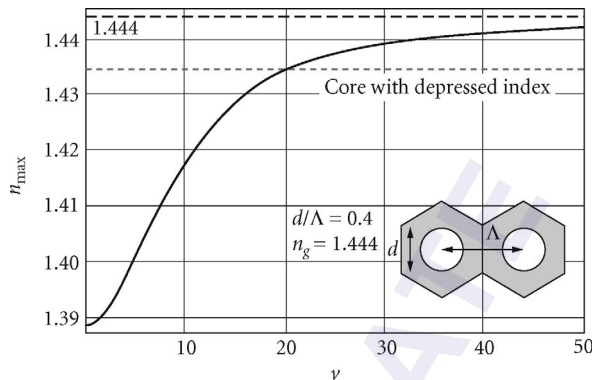


FIGURE 7 Maximum axial refractive index in the photonic crystal cladding as a function of the normalized frequency parameter ν for $d/\lambda = 0.4$ and $n_g = 1.444$. For this filling fraction of air (14.5%), the value at long wavelength ($\nu \rightarrow 0$) is $n_{\max} = 1.388$, in agreement with Eq. (5). The horizontal dashed gray line represents the case when the core is replaced with a glass of refractive index $n_{\text{co}} = 1.435$ (below that of silica), when guidance ceases for $\nu > 20$.

Photonic Bandgaps

Full two-dimensional PBGs exist in the black finger-shaped regions on Fig. 4. Some of these extend into the region $\beta < k$ where light is free to propagate in vacuum, confirming the feasibility of trapping light within a hollow core.

The bandgap edges coincide with points where resonances in the cladding unit cells switch on and off, that is, the eigenvalues of the unitary inter-unit-cell field transfer matrices change from $\exp(\pm j\phi)$ (propagation) to $\exp(\pm\gamma)$ (evanescence). At these transitions, depending on the band edge, the light is to a greater or lesser degree preferentially redistributed into the low or high index subregions. For example, at fixed optical frequency and small values of β , leaky modes peaking in the low index channels form a pass-band that terminates when the standing wave pattern has 100 percent visibility (Fig. 6c). For the high index strands (Fig. 6a and b) on the other hand, the band of real states is bounded by a lower value of β where the field amplitude changes sign between selected pairs of adjacent strands (depending on the lattice geometry), and an upper bound where the field amplitude does not change sign between the strands (this field distribution yields n_{\max}).

11.7 LINEAR CHARACTERISTICS OF GUIDANCE

In SMF, guided modes form within the range of axial refractive indices $n_{\text{cl}} < n_z < n_{\text{co}}$, when light is evanescent in the cladding ($n_z = \beta/k$; core and cladding indices are n_{co} and n_{cl}). In PCF, three distinct guidance mechanisms exist: a modified form of TIR,^{16,49} photonic bandgap guidance^{7,50} and a low leakage mechanism based on a Kagomé cladding structure.^{51,52} In the following subsections we explore the role of resonance and antiresonance, and discuss chromatic dispersion, attenuation mechanisms, and guidance in cores with refractive indices raised and lowered relative to the “mean” cladding value.

Resonance and Antiresonance

It is helpful to view the guided modes as being confined (or not confined) by resonance and antiresonance in the unit cells of the cladding crystal. If the core mode finds no states in the cladding with which it is phase-matched, light cannot leak out. This is a familiar picture in many areas of photonics. What is perhaps not so familiar is the use of the concept in two dimensions, where a repeating unit is tiled to form a photonic crystal cladding. This allows the construction of an intuitive picture of “cages,” “bars,” and “windows” for light and actually leads to a blurring of the distinction between guidance by modified TIR and photonic bandgap effects.

Positive Core-Cladding Index Difference

This type of PCF may be defined as one where the mean cladding refractive index in the long wavelength limit, $k \rightarrow 0$, [Eq. (5)] is lower than the core index (in the same limit). Under the correct conditions (high air-filling fraction), PBG guidance may also occur in this case, although experimentally the TIR-guided modes will dominate.

Controlling Number of Modes A striking feature of this type of PCF is that it is “endlessly single-mode” (ESM), that is, the core does not become multimode in the experiments, no matter how short the wavelength of the light.¹⁶ Although the guidance in some respects resembles conventional TIR, it turns out to have some interesting and unique features that distinguish it markedly from step-index fiber. These are due to the piecewise discontinuous nature of the core boundary—sections where (for $n_z > 1$) air holes strongly block the escape of light are interspersed with regions of barrier-free glass. In fact, the cladding operates in a regime where the transverse effective wavelength, Eq. (8), in silica is comparable with the glass substructures in the cladding. The zone of operation in Fig. 4 is $n_{\max} < n_z < n_g$ (point A).

In a solid-core PCF, taking the core radius $\rho = \Lambda$ and using the analysis in Ref. 16, the effective V -parameter can be calculated. This yields the plot in Fig. 8, where the full behavior from very low to very high frequency is predicted (the glass index was kept constant at 1.444). As expected, the number of guided modes approximately $V_{\text{PCF}}^2/2$ is almost independent of wavelength at high frequencies;

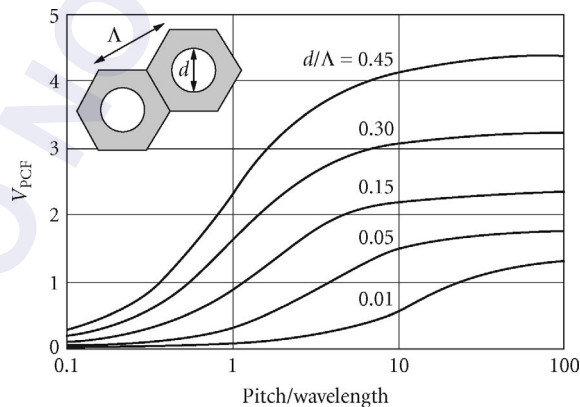


FIGURE 8 V -parameter for solid-core PCF (triangular lattice) plotted against the ratio of hole spacing to vacuum wavelength for different values of d/Λ . Numerical modeling shows that ESM behavior is maintained for $V_{\text{PCF}} \leq 4$ or $d/\Lambda \leq 0.43$.

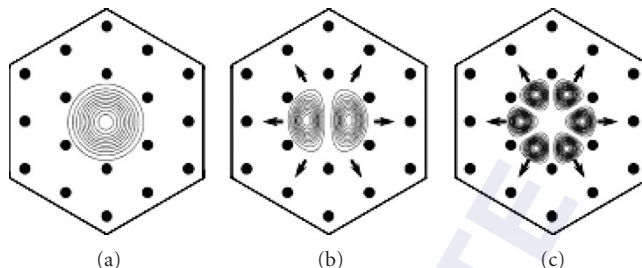


FIGURE 9 Schematic of modal filtering in a solid-core PCF: (a) The fundamental mode is trapped whereas (b) and (c) higher-order modes leak away through the gaps between the air holes.

the single-mode behavior is determined solely by the geometry. Numerical modeling shows that if $d/\Lambda < 0.43$ the fiber never supports any higher-order guided modes, that is, it is ESM.

This behavior can be understood by viewing the array of holes as a modal filter or “sieve” (Fig. 9). The fundamental mode in the glass core has a transverse effective wavelength $\lambda_{\text{eff}}^s \approx 4\Lambda$. It is thus unable to “squeeze through” the glass channels between the holes, which are $\Lambda - d$ wide and thus below the Rayleigh resolution limit $\approx \lambda_{\text{eff}}^s / 2 = 2\Lambda$. Provided the relative hole size d/Λ is small enough, higher-order modes are able to escape their transverse effective wavelength is shorter so they have higher resolving power. As the holes are made larger, successive higher-order modes become trapped.

ESM behavior may also be viewed as being caused by strong wavelength dispersion in the photonic crystal cladding, which forces the core-cladding index step to fall as the wavelength gets shorter (Fig. 7).^{16,49} This counteracts the usual trend toward increasingly multimode behavior at short wavelengths. In the limit of very short wavelength the light strikes the glass-air interfaces at glancing incidence, and is strongly rejected from the air holes. In this regime the transverse single-mode profile does not change with wavelength. As a consequence the angular divergence (roughly twice the numerical aperture) of the emerging light is proportional to wavelength; in SMFs it is approximately constant owing to the appearance of more and more higher-order guided modes as the frequency increases.

Thus, the refractive index of the photonic crystal cladding increases with optical frequency, tending toward the index of silica glass in the short wavelength limit. If the core is made from a glass of refractive index lower than that of silica (e.g., fluorine-doped silica), guidance is lost at wavelengths shorter than a certain threshold value (see Fig. 7).⁵³ Such fibers have the unique ability to prevent transmission of short wavelength light—in contrast to conventional fibers which guide more and more modes as the wavelength falls.

Ultra-Large Area Single-Mode The modal filtering in ESM-PCF is controlled only by the geometry (d/Λ for a triangular lattice). A corollary is that the behavior is quite independent of the absolute size of the structure, permitting single-mode fiber cores with arbitrarily large areas. A single-mode PCF with a core diameter of 22 μm at 458 nm was reported in 1998.¹⁷ In conventional step-index fibers, where $V < 2.405$ for single-mode operation, this would require uniformity of core refractive index to approximately part in 10^5 —very difficult to achieve if MCVD is used to form the doped core. Larger mode areas allow higher power to be carried before the onset of intensity-related nonlinearities or damage, and have obvious benefits for delivery of high laser power, fiber amplifiers, and fiber lasers. The bend-loss performance of such large-core PCFs is discussed in subsection “Bend Loss” in Sec. 11.7.

Fibers with Multiple Cores The stacking procedure makes it straightforward to produce multicore fiber. A preform stack is built up with a desired number of solid (or hollow) cores, and drawn down to fiber in the usual manner.²⁰ The coupling strength between the cores depends on

the sites chosen, because the evanescent decay rate of the fields changes with azimuthal direction. Applications include curvature sensing.⁵⁴ More elaborate structures can be built up, such as fibers with a central single-mode core surrounded by a highly multimode cladding waveguide are useful in applications such as high power cladding-pumped fiber lasers^{55,56} and two-photon fluorescence sensors⁵⁷ (see Fig. 2f).

Negative Core-Cladding Index Difference

Since TIR cannot operate under these circumstances, low loss waveguiding is only possible if a PBG exists in the range $\beta < kn_{\text{co}}$.

Hollow-Core Silica/Air In silica-air PCF, larger air-filling fractions and small interhole spacings are necessary to achieve photonic bandgaps in the region $\beta < k$. The relevant operating region on Fig. 4 is to the left of the vacuum line, inside one of the bandgap fingers (point B). These conditions ensure that light is free to propagate, and form guided modes, within the hollow core while being unable to escape into the cladding. The number N of such modes is controlled by the depth and width of the refractive index “potential well” and is approximately given by

$$N \approx k^2 \rho^2 (n_{\text{high}}^2 - n_{\text{low}}^2) / 2 \quad (9)$$

where n_{high} and n_{low} are the refractive indices at the edges of the PBG at fixed frequency and ρ is the core radius. Since the bandgaps are quite narrow ($n_{\text{high}}^2 - n_{\text{low}}^2$ is typically a few percent) the hollow core must be sufficiently large if a guided mode is to exist at all. In the first hollow-core PCF, reported in 1999,⁷ the core was formed by omitting seven capillaries from the preform stack (Fig. 1d). An electron micrograph of a more recent structure, with a hollow core made by removing 19 missing capillaries from the stack, is shown in Fig. 10.⁵⁸

In hollow-core PCF, guidance can only occur when a photonic bandgap coincides with a core resonance. This means that only restricted bands of wavelength are guided. This feature can be very useful for suppressing parasitic transitions by filtering away the unwanted wavelengths, for example, in fiber lasers and in stimulated Raman scattering in gases.⁵⁹

Higher Refractive Index Glass Achieving a bandgap in higher refractive index glasses for $\beta < k$ presents at first glance a dilemma. Whereas a larger refractive index contrast generally yields wider bandgaps, the higher “mean” refractive index seems likely to make it more difficult to achieve bandgaps

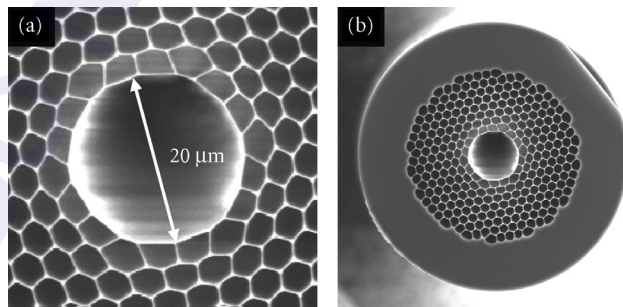


FIGURE 10 Scanning electron micrographs of a low loss hollow PCF (manufactured by BlazePhotonics Ltd.) with attenuation approximately 1 dB/km at 1550 nm wavelength: (a) detail of the core (diameter 20.4 μm) and (b) the complete fiber cross section.

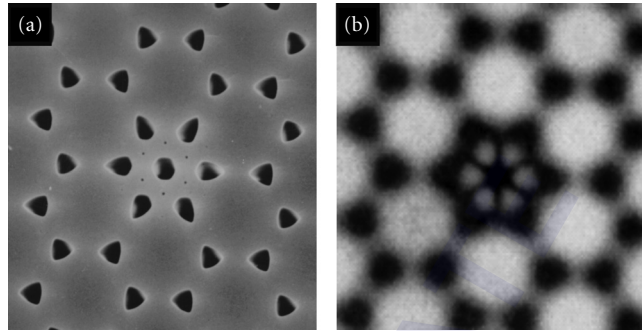


FIGURE 11 (a) A PCF with a “carbon-ring” lattice of air holes and an extra central hole to form a low index core. (b) When white light is launched, only certain bands of wavelength are transmitted in the core—here a six-lobed mode (in the center of the image, blue in the original near-field image) emerges from the end-face.⁵⁰

for incidence from vacuum. Although this argument holds in the scalar approximation, the result of calculations show that vector effects become important at higher levels of refractive index contrast (e.g., 2:1 or higher) and a new species of bandgap appears for smaller filling fractions of air than in silica-based structures. The appearance of this new type of gap means that it is actually easier to obtain wide bandgaps with higher index glasses such as tellurites or chalcogenides.⁴³

Surface States on Core-Cladding Boundary The first PCF that guided by photonic bandgap effects consisted of a lattice of air holes arranged in the same way as the carbon rings in graphite. The core was formed by introducing an extra hole at the center of one of the rings, its low index precluding the possibility of TIR guidance.⁵⁰ When white light was launched into the core region, a colored mode was transmitted—the colors being dependent on the absolute size to which the fiber was drawn. The modal patterns had six equally strong lobes, disposed in a flower-like pattern around the central hole. Closer examination revealed that the light was guided not in the air holes but in the six narrow regions of glass surrounding the core (Fig. 11). The light remained in these regions, despite the close proximity of large “rods” of silica, full of modes. This is because, for particular wavelengths, the phase velocity of the light in the core is not coincident with any of the phase velocities available in the transmission bands created by nearest-neighbor coupling between rod modes. Light is thus unable to couple over to them and so remains trapped in the core.

Similar guided modes are commonly seen in hollow-core PCF, where they form surface states (analogous with electronic surface states in semiconductor crystals) on the rim of the core, confined on the cladding side by photonic bandgap effects. These surface states become phase-matched to the air-guided mode at certain wavelengths, and if the two modes share the same symmetry they couple to form an anticrossing on the frequency-wavevector diagram (Figs. 12 and 13). Within the anticrossing region, the modes share the characteristics of both an air-guided mode and a surface mode, and this consequently perturbs the group velocity dispersion and contributes additional attenuation (see subsection “Absorption and Scattering” in Sec. 11.7).^{60–62}

All-Solid Structures In all-solid bandgap guiding fibers the core is made from low index glass and is surrounded by an array of high index glass strands.^{63–65} Since the mean core-cladding index contrast is negative, TIR cannot operate, and photonic bandgap effects are the only possible guidance mechanism. These structures have some similarities with one-dimensional “ARROW” structures, where antiresonance plays an important role.⁶⁶

When the cladding strands are antiresonant, light is confined to the central low index core by a mechanism not dissimilar to the modal filtering picture in subsection “Controlling Number of

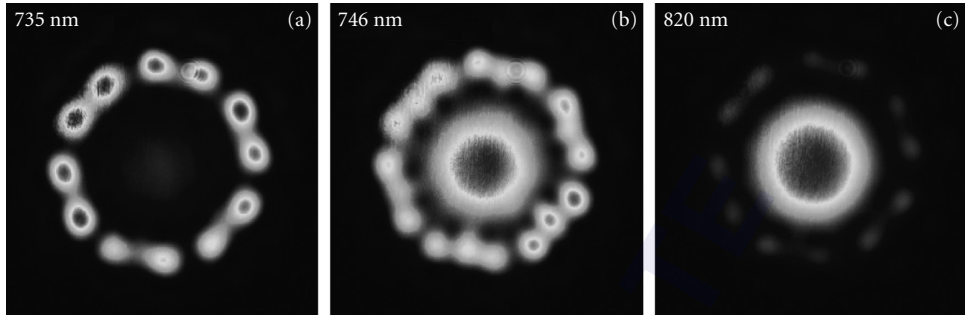


FIGURE 12 Near-field end-face images of the light transmitted in hollow-core PCF designed for 800 nm transmission. For light launched in the core mode, at 735 nm an almost pure surface mode is transmitted, at 746 nm a coupled surface-core mode, and at 820 nm an almost pure core mode.⁶⁰ (The ring-shaped features are an artifact caused by converting from false color to gray scale; the intensity increases toward the dark centers of the rings.) (Images courtesy G. Humbert, University of Bath).⁶⁰

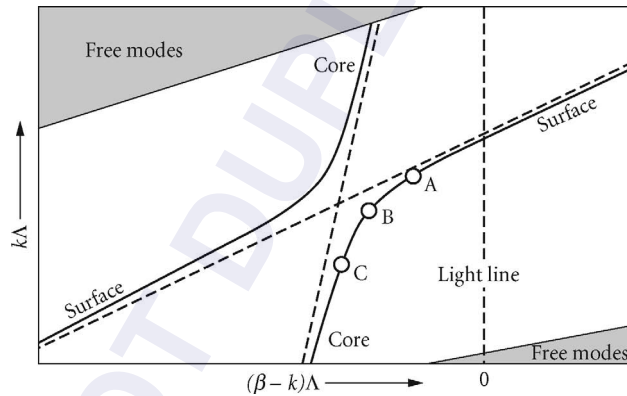


FIGURE 13 Example mode trajectories showing the anticrossing of a core-guided mode with a surface mode in a hollow-core PCF. The dotted lines show the approximate trajectories of the two modes in the absence of coupling (for instance if the modes are of different symmetries), and the vertical dashed line is the air line. The gray regions, within which the mode trajectories are not shown, are the band edges; the white region is the photonic bandgap. Points A, B, and C are the approximate positions of the modes shown in Fig. 12.

Modes” in Sec. 11.7;⁵² the high index cores act as the “bars of a cage,” so that no features in the cladding are resonant with the core mode, resulting in a low loss guided mode. Guidance is achieved over wavelength ranges that are punctuated with high loss windows where the cladding “bars” become resonant (Fig. 14). Remarkably, it is possible to achieve photonic bandgap guidance by this mechanism even at index contrasts of 1 percent,^{63,67} with losses as low as 20 dB/km at 1550 nm.⁶⁸

Low Leakage Guidance The transmission bands are greatly widened in hollow-core PCFs with a kagomé lattice in the cladding⁵² (Fig. 2j). The typical attenuation spectrum of such a fiber has a loss of order 1 dB/m over a bandwidth of 1000 nm or more. Numerical simulations show that, while the cladding structure supports no bandgaps, the density of states is greatly reduced near the vacuum

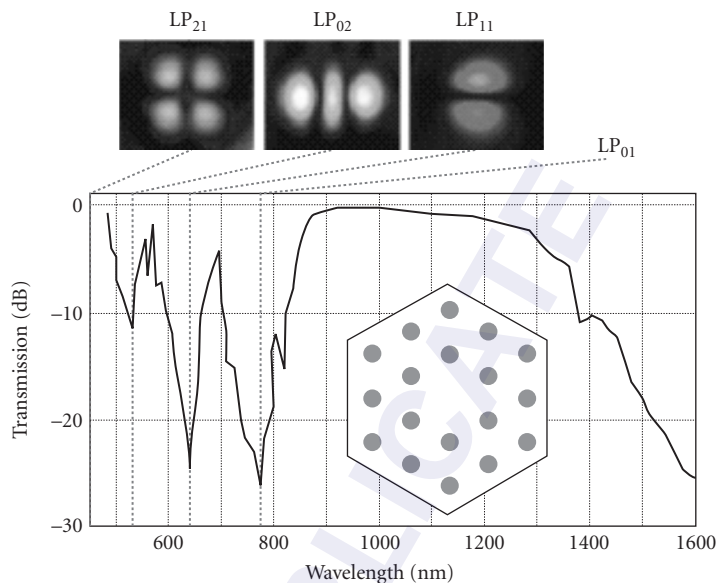


FIGURE 14 Lower: Measured transmission spectrum (using a white-light supercontinuum source) for a PCF with a pure silica core and a cladding formed by an array of Ge-doped strands ($d/\Lambda = 0.34$ hole spacing $\sim 7 \mu\text{m}$, index contrast 1.05:1). The transmission is strongly attenuated when the core mode becomes phase-matched to different LP_{nm} “resonances” in the cladding strands. Upper: Experimental images [left to right, taken with blue (500 nm), green (550 nm), and red (650 nm) filters] of the near-field profiles in the cladding strands at three such wavelengths. The fundamental LP_{01} resonance occurs at approximately 820 nm and the four-lobed blue resonance lies off the edge of the graph.

line. The consequential poor overlap between the core states, together with the greatly reduced number of cladding states, appears to slow down the leakage of light—though the precise mechanism is still a matter of debate.⁵²

Birefringence

The modes of a perfect sixfold symmetric core and cladding structure are not birefringent.⁶⁹ In practice, however, the large glass-air index difference means that even slight accidental distortions in the structure yield a degree of birefringence. Therefore, if the core is deliberately distorted so as to become twofold symmetric, extremely high values of birefringence can be achieved. For example, by introducing capillaries with different wall thicknesses above and below a solid glass core (Figs. 1c and 2g), values of birefringence some 10 times larger than in conventional fibers can be obtained.⁷⁰ It is even possible to design and fabricate strictly single-polarization PCFs in which only one polarization state is guided.⁷¹ By filling selected holes with a polymer, the birefringence can be thermally tuned.⁷² Hollow-core PCF with moderate levels of birefringence ($\sim 10^{-4}$) can be realized either by forming an elliptical core or by adjusting the structural design of the core surround.^{73,74}

Experiments show that the birefringence in PCF is some 100 times less sensitive to temperature variations than in conventional fibers, which is important in many applications.^{75–77} This is because traditional “polarization maintaining” fibers (bow-tie, elliptical core, or Panda) contain at least two different glasses, each with a different thermal expansion coefficient. In such structures, the resulting temperature-dependent stresses make the birefringence a strong function of temperature.

Group Velocity Dispersion

Group velocity dispersion (GVD)—which causes different frequencies of light to travel at different group velocities—is a factor crucial in the design of telecommunications systems and in all kinds of nonlinear optical experiments. PCF offers greatly enhanced control of the magnitude and sign of the GVD as a function of wavelength. In many ways this represents an even greater opportunity than a mere enhancement of the effective nonlinear coefficient.

Solid Core As the optical frequency increases, the GVD in SMF changes sign from anomalous ($D > 0$) to normal ($D < 0$) at approximately $1.3 \mu\text{m}$. In solid-core PCF as the holes get larger, the core becomes more and more isolated, until it resembles an isolated strand of silica glass (Fig. 15). If the whole structure is made very small (core diameters $< 1 \mu\text{m}$ have been made) the zero dispersion point of the fundamental guided mode can be shifted to wavelengths in the visible.^{9,18} For example, the PCF in Fig. 1f has a dispersion zero at 560 nm .

By careful design, the wavelength dependence of the GVD can also be reduced in PCFs at much lower air-filling fractions. Figure 16 shows the flattened GVD profiles of three PCFs with cores several μm in diameter.^{78,79} These fibers operate in the regime where SMF is multimoded. Although the fundamental modes in both SMF and PCF have similar dispersion profiles, the presence of higher-order modes (not guided in the PCF, which is endlessly single mode) makes the use of SMF impractical.

A further degree of freedom in GVD design may be gained by working with multicomponent glasses, such as Schott SF6, where the intrinsic zero dispersion point occurs at approximately $1.8 \mu\text{m}$.⁸ In highly nonlinear small-core PCF, this shifts the whole dispersion landscape to longer wavelengths than in a silica-based PCF with the same core size and geometry.

Hollow Core Hollow-core fiber behaves in many respects rather like a circular-cylindrical hollow metal waveguide, which has anomalous dispersion (the group velocity increases as the frequency rises). The main difference, however, is that the dispersion changes sign at the high frequency edge, owing to the approach of the photonic band edge and the weakening of the confinement (Fig. 17).⁸⁰

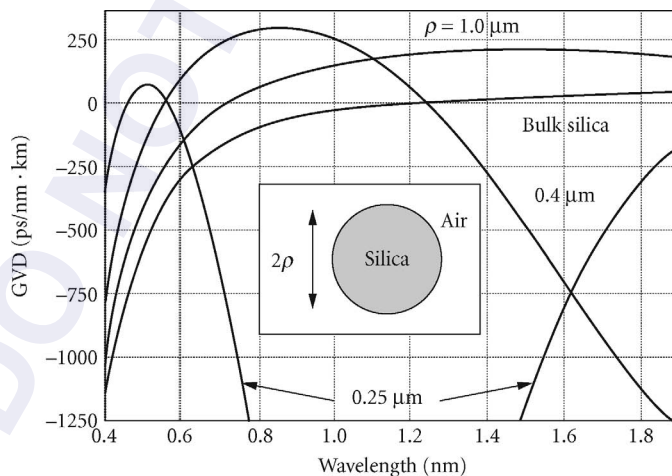


FIGURE 15 The calculated group velocity dispersion of three circular strands of silica glass, radii 0.25 , 0.4 , and $1.0 \mu\text{m}$, compared with the dispersion of bulk glass. The narrower strands have two dispersion zeros within the transparency window of silica.

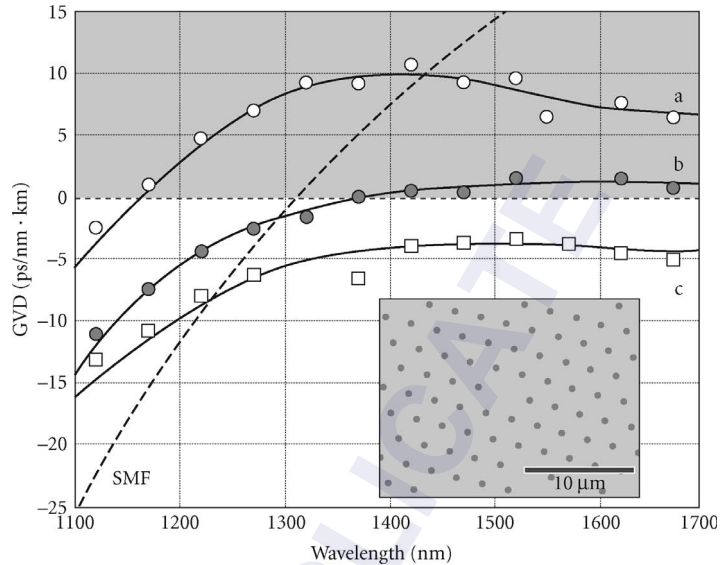


FIGURE 16 Group velocity dispersion profiles, against wavelength, for three different PCFs designed to have lowlevel ultraflattened GVD.^{78,79} The curve for Corning SMF-28 is included for comparison.

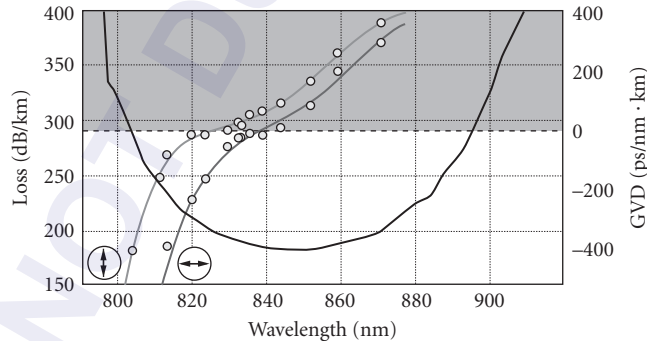


FIGURE 17 Measured attenuation and GVD spectra for a hollow-core PCF designed for 850 nm transmission.⁸⁰ The core is slightly elliptical, so the dispersion in each eigenstate of polarization is different.

Attenuation Mechanisms

An advantage common to all fibers is the very large extension ratio from preform to fiber, which has the effect of smoothing out imperfections, resulting in a transverse structure that is extremely invariant with distance along the fiber. This is the chief reason for the ultralow attenuation displayed by fibers, compared to other waveguide structures. In PCF the losses are governed by two main parameters: the fraction of light in glass and the roughness at the glass-air interfaces. The light-in-glass fraction can be controlled by judicious design, and ranges from close to 100 percent in solid core fibers to less than 1 percent in the best hollow-core fibers.

Absorption and Scattering The best reported loss in solid-core PCF, from a group in Japan, stands at 0.28 dB/km at 1550 nm, with a Rayleigh scattering coefficient of $0.85 \text{ dB} \cdot \text{km}^{-1} \cdot \mu\text{m}^{-4}$. A 100 km length of this fiber was used in the first PCF-based penalty-free dispersion-managed soliton transmission system at 10 Gb/s.^{81,82} The slightly higher attenuation compared to SMF is due to roughness at the glass-air interfaces.⁶¹

It is hollow-core PCF, however, that has the greatest potential for extremely low loss, since the light is traveling predominantly in empty (or gas-filled) space. Although the best reported attenuation in hollow-core PCF stands at 1.2 dB/km,⁵⁸ values below 0.2 dB/km or even lower seem feasible with further development of the technology. The prospect of improving on conventional fiber, at the same time greatly reducing the nonlinearities associated with a solid glass core, is intriguing. By using infrared glasses, transmission can be extended into the infrared⁸³ and recent work shows that silica hollow-core PCF can even be used with acceptable levels of loss in the mid-IR,⁸⁴ owing to the very low overlap between the optical field and the glass.

In the latest hollow-core silica PCF, with loss levels approaching 1 dB/km at 1550 nm, very small effects can contribute significantly to the attenuation floor. The ultimate loss limit in such fibers is determined by surface roughness caused by thermally driven capillary waves, which are present at all length scales. These interface ripples freeze in when the fiber cools, introducing high scattering losses for modes that are concentrated at the interfaces, such as surface modes guided on the edge of the core. The pure core mode does not itself “feel” the ripples very strongly, except at anticrossing wavelengths where it becomes phase-matched to surface modes, causing light to move to the surface and experience enhanced scattering.

The result is a transmission spectrum consisting of windows of high transparency punctuated with bands of high attenuation (Fig. 18). This picture has been confirmed by measurements of the surface roughness in hollow-core PCFs, the angular distribution of the power scattered out of the core, and the wavelength dependence of the minimum loss of fibers drawn to different scales.⁵⁸ The thin glass shell surrounding the hollow core can be designed to be antiresonant with the core mode, permitting further exclusion of light from the glass.⁸⁵

Ignoring material dispersion, the whole transmission landscape shifts linearly in wavelength in proportion to the overall size of the structure (a consequence of Maxwell’s equations). This means

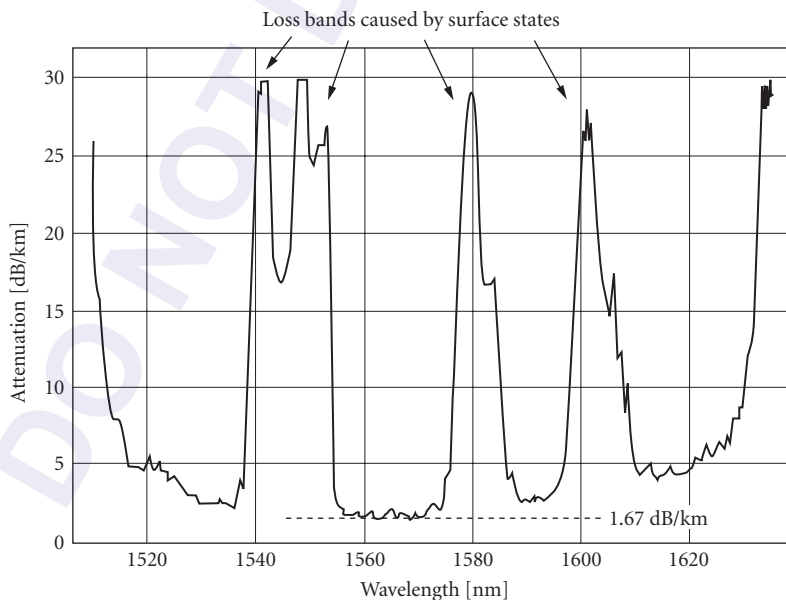


FIGURE 18 Attenuation spectrum of a typical ultralow loss hollow-core PCF designed for operation in the 1550 nm telecommunications band (see micrographs in Fig. 10).

that the smallest loss at a given wavelength will be obtained by drawing a fiber to a particular diameter. The optical overlap with the surface roughness scales inversely with the size with the fiber, and the scattering itself may be regarded as being governed by the density of states into which scattering can occur, which in three-dimensions scales as λ^{-2} . Thus the wavelength of minimum loss scales as λ^{-3} , in contrast to the λ^{-4} dependence of Rayleigh scattering in bulk glass.

Bend Loss Conventional fibers suffer additional loss if bent beyond a certain critical radius R_c , which depends on wavelength, core-cladding refractive index step and—most notably—the third power of core radius ρ .³⁹ For wavelengths longer than a certain value (the “long wavelength bend edge”) all guidance is effectively lost.

A starting point for understanding bend loss in solid core ESM-PCF (perhaps the most interesting case) is the long wavelength limit. ESM behavior occurs when $d/\Lambda < 0.43$, which sets the highest air-filling fraction at 16.8 percent and yields an area-averaged cladding refractive index of 1.388 (silica index of 1.444)—valid in the limit $k \rightarrow 0$. This index step is some 10 times higher than in Corning SMF-28, making ESM-PCF relatively much less susceptible to bend loss at long wavelengths. For a step-index fiber with a Ge-doped core, 40 mol % of GeO_2 would be needed to reach the same index step (assuming 0.0014 index change per mol % GeO_2).¹¹ The result is that the long-wavelength bend edge in ESM-PCF is in the infrared beyond the transparency window of silica glass, even when the core radius is large.⁸⁶

ESM-PCF also exhibits a *short* wavelength bend edge, caused by bend-induced coupling from the fundamental to higher-order modes, which of course leak out of the core.¹⁶ The critical bend radius for this loss varies as

$$R_c \sim \Lambda^3/\lambda^2 \quad (10)$$

compared to R_c approximately λ for SMF. The reciprocal dependence on λ^2 makes it inevitable that a short-wavelength bend edge will appear in ESM-PCF. Following⁸⁶ in taking the pre-factor in Eq. (10) as unity, R_c can be plotted against wavelength for different values of core radius (\approx interhole spacing). This is illustrated in Fig. 19. A step-index fiber, with the same core-cladding index step as

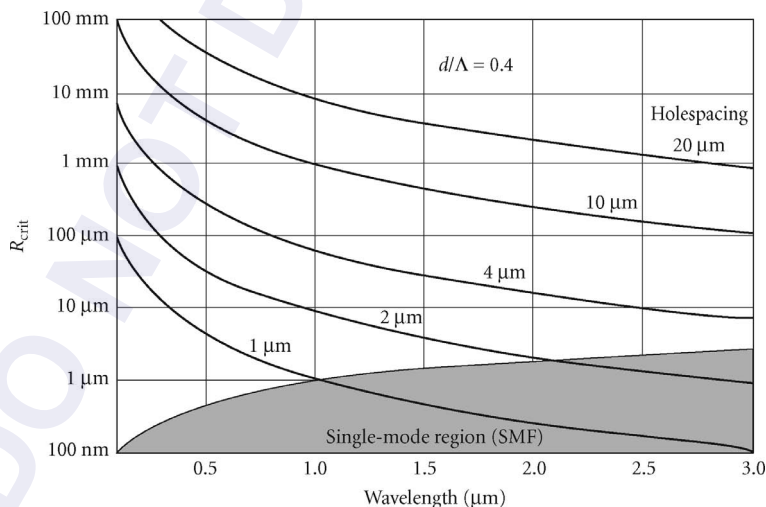


FIGURE 19 Short-wavelength critical bend radii for ESM-PCF with $d/\Lambda = 0.4$ plotted against vacuum wavelength for different values of hole-spacing (approximately equal to the core radius). As the wavelength increases at constant core size, a step-index fiber with the same core:cladding index as the PCF in the long-wavelength limit [1.444:1.388 using Eq. (5)] becomes single-mode when the curves enter the shaded region. The step-index fiber is multi-mode over wide parameter ranges where ESM-PCF has negligible short-wavelength bend loss.

ESM-PCF in the long wavelength limit, is multimode over wide parameter ranges where ESM-PCF has negligible bend loss.

In contrast, hollow-core PCF is experimentally very insensitive to bend loss—in many cases no appreciable drop in transmission is observed until the fiber breaks. This is because the effective depth of “potential well” for the guided light (see section “Negative Core-Cladding Index Difference”), given by the distance $\Delta\beta$ between the edges of the photonic bandgap, is substantially larger than in SMF.

Confinement Loss The photonic crystal cladding in a realistic PCF is of course finite in extent. For a guided mode, the Bloch waves in the cladding are evanescent, just like the evanescent plane waves in the cladding of a conventional fiber. If the cladding is not thick enough, the evanescent field amplitudes at the cladding/coating boundary can be substantial, causing attenuation. In the solid core case for small values of d/Λ the resulting loss can be large unless a sufficiently large number of periods is used.⁷⁸

Very similar losses are observed in hollow-core fibers, where the “strength” of the photonic bandgap (closely related to its width in β) determines how many periods are needed to reduce confinement loss to acceptable levels. Numerical modeling is useful for giving an indication of how many periods are needed to reach a required loss level. The cladding field intensity in the ultralow loss PCF reported in⁵⁸ falls by approximately 9 dB per period, reaching -63 dB at the edge of the photonic crystal region.

11.8 NONLINEAR CHARACTERISTICS OF GUIDANCE

The ability to enhance or reduce the effective nonlinear coefficients, and at the same time control the magnitude and wavelength dependence of the GVD, makes PCF a versatile vehicle for studies of nonlinear effects.

Kerr Nonlinearities

The optical Kerr effect drives effects such as four-wave mixing, self-phase modulation, modulation instability, and soliton formation. To take account of the differing proportions of light in glass and air,⁸⁷ it is necessary to derive an expression for the effective nonlinear γ coefficient, which in the scalar approximation takes the form

$$\gamma = \frac{k \iint n(x, y) n_2(x, y) \psi^4(x, y) dx dy}{n_m \left(\iint \psi^2(x, y) dx dy \right)^2} \text{ W}^{-1} \text{ m}^{-1} \quad (11)$$

where ψ is the transverse field amplitude profile, $n(x, y)$ and $n_2(x, y)$ are the linear and nonlinear refractive index profiles, and n_m is the modal phase index. The integrals must be evaluated over the entire transverse plane, which in practice means over the area where the field amplitudes are nonnegligible. Equation (11) can be reexpressed for a composite PCF made from two or more homogeneous materials

$$\gamma = k \sum_i \frac{n_2^i}{A_i} \quad (12)$$

where n_2^i is the nonlinear refractive index of material i and A_i its nonlinear effective area

$$A_i = \frac{n_m \left(\iint \psi^2(x, y) dx dy \right)^2}{\iint u_i(x, y) \psi^4(x, y) dx dy} \quad (13)$$

where $u_i(x, y)$ equals n_i in regions made from material i and zero elsewhere.

The nonlinear refractive indices of silica glass and air are, respectively, $2.5 \times 10^{-20} \text{ m}^2\text{W}^{-1}$ and $2.9 \times 10^{-23} \text{ m}^2\text{W}^{-1}$. Multicomponent glasses typically have values an order of magnitude or more higher, and the hollow channels can of course be filled with another gas with a different nonlinear coefficient.

The highest nonlinear coefficient available in conventional step-index fibers is $\gamma: 20 \text{ W}^{-1}\text{km}^{-1}$ at 1550 nm.⁸⁸ By comparison, a solid-core PCF similar to the one in Fig. 1f but with a core diameter 1 μm has a nonlinearity of $\gamma: 240 \text{ W}^{-1}\text{km}^{-1}$, at 850 nm, and values as high as $\gamma = 550 \text{ W}^{-1}\text{km}^{-1}$ at 1550 nm have been measured for PCFs made from multicomponent glasses.⁸⁹ In complete contrast, hollow-core PCF has extremely low levels of nonlinearity, owing to the small overlap between the glass and the light. In a recent example, a fiber was reported with a nonlinear coefficient $\gamma = 0.023 \text{ W}^{-1}\text{km}^{-1}$ (some $10^4 \times$ smaller than in a typical highly nonlinear solid-core PCF).^{90,91}

Although the level of nonlinearity is clearly important, the actual nonlinear effects that appear in a particular case are also strongly dependent on the magnitude, sign and wavelength dependence of the GVD as well as on the characteristics of the pump laser.⁹² They are determined by the relative values of the nonlinear length $L_{\text{nl}} = \gamma^{-1}P_0^{-1}$, where P_0 is the peak power, the dispersion length $L_{\text{D}} = \tau^2/|\beta_2|$, where τ is the pulse duration and the effective fiber length $L_{\text{eff}} = [1 - \exp(-\alpha L)]/\alpha$ where $\alpha \text{ m}^{-1}$ is the power attenuation coefficient.⁹³ For a solid-core PCF with $\gamma: 240 \text{ W}^{-1}\text{km}^{-1}$, a peak power of 10 kW yields $L_{\text{nl}} < 0.5 \text{ mm}$. For typical values of loss (usually between 1 and 100 dB/km) $L_{\text{eff}} \gg L_{\text{nl}}$ and the nonlinearity dominates. For dispersion values in the range $-300 < \beta_2 < 300 \text{ ps}^2/\text{km}$ and pulse durations $\tau = 200 \text{ fs}$, $L_{\text{D}} > 0.1 \text{ m}$. Since both these lengths are much longer than the nonlinear length, it is easy to observe strong nonlinear effects.

Supercontinuum Generation One of the most successful applications of nonlinear PCF is to supercontinuum (SC) generation from ps and fs laser pulses. When high power pulses travel through a material, their frequency spectrum can be broadened by a range of interconnected nonlinear effects.⁹⁴ In bulk materials, the preferred pump laser is a regeneratively amplified Ti:sapphire system producing high (mJ) energy fs pulses at 800 nm wavelength and kHz repetition rate. Supercontinua have also previously been generated in SMF by pumping at 1064 or 1330 nm,⁹⁵ the spectrum broadening out to longer wavelengths mainly due to stimulated Raman scattering (SRS). Then in 2000, it was observed that highly nonlinear PCF, designed with zero GVD close to 800 nm, massively broadens the spectrum of low (few nJ) energy unamplified Ti:sapphire pulses launched into just a few cm of fiber.^{21,96,97} Removal of the need for a power amplifier, the hugely increased ($\sim 100 \text{ MHz}$) repetition rate, and the spatial and temporal coherence of the light emerging from the core, makes this source unique. The broadening extends both to higher and to lower frequencies because four-wave mixing operates more efficiently than SRS when the dispersion profile is appropriately designed. This SC source has applications in optical coherence tomography,^{98,99} frequency metrology,^{100,101} and all kinds of spectroscopy. It is particularly useful as a bright lowcoherence source in measurements of group delay dispersion based on a Mach-Zehnder interferometer.

A comparison of the bandwidth and spectrum available from different broad-band light sources is shown in Fig. 20; the advantages of PCF-based SC sources are evident. Supercontinua have been generated in different PCFs at 532 nm,¹⁰² 647 nm,¹⁰³ 1064 nm,¹⁰⁴ and 1550 nm.⁸ Using inexpensive microchip lasers at 1064 or 532 nm with an appropriately designed PCF, compact SC sources are now available with important applications across many areas of science. A commercial ESM-PCF based source uses a 10-W fiber laser delivering 5 ps pulses at 50 MHz repetition rate, and produces an average spectral power density of approximately 4.5 mW/nm in the range 450 to 800 nm.¹⁰⁵ The use of multicomponent glasses such as Schott SF6 or tellurite glass allows the balance of nonlinearity and dispersion to be adjusted, as well as offering extended transparency into the infrared.¹⁰⁶

Parametric Amplifiers and Oscillators In step-index fibers the performance of optical parametric oscillators and amplifiers is severely constrained owing to the limited scope for GVD engineering. In PCF these constraints are lifted, permitting flattening of the dispersion profile and control of higher-order dispersion terms. The wide range of experimentally available group-velocity dispersion profiles has, for example, allowed studies of ultrashort pulse propagation in the 1550 nm wavelength band with flattened dispersion.^{78,79} The effects of higher-order dispersion in such PCFs are subtle.^{107,108} Parametric devices have been designed for pumping at 647, 1064, and 1550 nm, the small

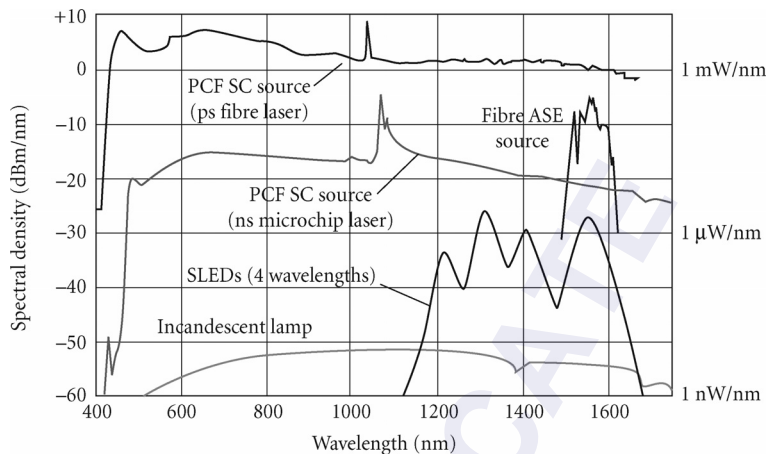


FIGURE 20 Comparison of the brightness of various broad-band light sources (SLED—superluminescent light-emitting diode; ASE—amplified spontaneous emission; SC—supercontinuum). The microchip laser SC spectrum was obtained by pumping at 1064 nm with 600 ps pulses. (Updated version of a plot by Hendrik Sabert.)

effective mode areas offering high gain for a given pump intensity, and PCF-based oscillators synchronously pumped by fs and ps pump pulses have been demonstrated at relatively low power levels.^{109–112} Dispersion-engineered PCF is being successfully used in the production of bright sources of correlated photon pairs, by allowing the signal and idler side-bands to lie well outside the noisy Raman band of the glass. In a recent example, a PCF with zero dispersion at 715 nm was pumped by a Ti:sapphire laser at 708 nm (normal dispersion).¹¹³ Under these conditions phase-matching is satisfied by signal and idler waves at 587 and 897 nm, and 10 million photon pairs per second were generated and delivered via single-mode fiber to Si avalanche detectors, producing approximately 3.2×10^5 coincidences per second for a pump power of 0.5 mW. These results point the way to practical and efficient sources entangled photon pairs that can be used as building blocks in future multiphoton interference experiments.

Soliton Self-Frequency Shift Cancellation The ability to create PCFs with negative dispersion slope at the zero dispersion wavelength (in SMF the slope is positive, i.e., the dispersion becomes more anomalous as the wavelength increases) has made it possible to observe Čerenkov-like effects in which solitons (which form on the anomalous side of the dispersion zero) shed power into dispersive radiation at longer wavelengths on the normal side of the dispersion zero. This occurs because higher-order dispersion causes the edges of the soliton spectrum to phase-match to linear waves. The result is stabilization of the soliton self-frequency shift, at the cost of gradual loss of soliton energy.¹¹⁴ The behavior of solitons in the presence of wavelength-dependent dispersion is the subject of many recent studies.^{115.}

Raman Scattering

The basic characteristics of glass-related Raman scattering in PCF, both stimulated and spontaneous, do not noticeably differ compared to SMF. One must of course take account of the differing proportions of light in glass and air (see section “Kerr Nonlinearities”) to calculate the effective strength of the Raman response. A very small solid glass core allows one to enhance stimulated Raman scattering, whereas in a hollow core it is strongly suppressed.

Brillouin Scattering

The periodic micro/nanostructuring in ultrasmall core glass-air PCF strongly alters the acoustic properties compared to conventional SMF.^{116–119} Sound can be guided in the core both as leaky and as tightly confined acoustic modes. In addition, the complex geometry and “hard” boundaries cause coupling between all three displacement components (radial, azimuthal, and axial), with the result that each acoustic mode has elements of both shear (S) or longitudinal (L) strain. This complex acoustic behavior strongly alters the characteristics of forward and backward Brillouin scattering.

Backward Scattering When a solid-core silica-air PCF has a core diameter of around 70 percent of the vacuum wavelength of the launched laser light, and the air-filling fraction in the cladding is very high, the spontaneous Brillouin signal displays multiple bands with Stokes frequency shifts in the 10 GHz range. These peaks are caused by discrete guided acoustic modes, each with different proportions of longitudinal and shear strain, strongly localized to the core.¹²⁰ At the same time the threshold power for stimulated Brillouin scattering increases fivefold—a rather unexpected result, since conventionally one would assume that higher intensities yield lower nonlinear threshold powers. This occurs because the effective overlap between the tightly confined acoustic modes and the optical mode is actually smaller than in a conventional fiber core; the sound field contains a large proportion of shear strain, which does not contribute significantly to changes in refractive index. This is of direct practical relevance to parametric amplifiers, which can be pumped 5 times harder before stimulated Brillouin scattering appears.

Forward Scattering The very high air-filling fraction in small-core PCF also permits sound at frequencies of a few GHz to be trapped purely in the transverse plane by phononic bandgap effects (Fig. 21). The ability to confine acoustic energy at zero axial wavevector $\beta_{ac} = 0$ means that the ratio

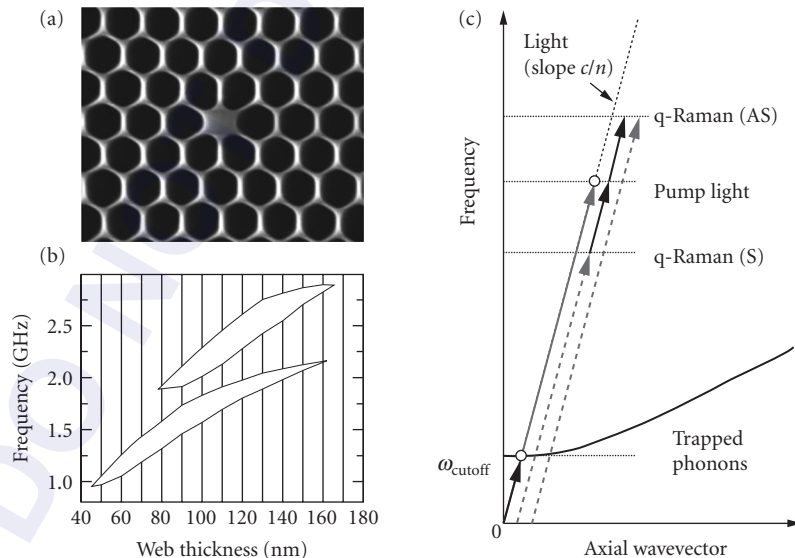


FIGURE 21 (a) Example of PCF used in studies of Brillouin scattering (core diameter 1.1 μm); (b) the frequencies of full phononic bandgaps (in-plane propagation, pure in-plane motion) in the cladding of the PCF in (b); and (c) illustrating how a trapped acoustic phonon can phase-match to light at the acoustic cutoff frequency. The result is a quasi-Raman scattering process that is automatically phase-matched. (After Ref. 121.)

of frequency ω_{ac} to wavevector β_{ac} becomes arbitrarily large as $\beta_{ac} \rightarrow 0$, and thus can easily match the value for the light guided in the fiber, c/n . This permits phase-matched interactions between the acoustic mode and two spatially identical optical modes of different frequency.¹²¹ Under these circumstances the acoustic mode has a well-defined cutoff frequency ω_{cutoff} above which its dispersion curve—plotted on an (ω, β) diagram—is flat, similar to the dispersion curve for optical phonons in diatomic lattices. The result is a scattering process that is Raman-like (i.e., the participating phonons are optical-phonon-like), even though it makes use of acoustic phonons; Brillouin scattering is turned into Raman scattering, power being transferred into an optical mode of the same order, frequency shifted from the pump frequency by the cutoff frequency. Used in stimulated mode, this effect may permit generation of combs of frequencies spaced by approximately 2 GHz at 1550 nm wavelength.

11.9 INTRAFIBER DEVICES, CUTTING, AND JOINING

As PCF becomes more widely used, there is an increasing need for effective cleaves, low loss splices, multipoint couplers, intrafiber devices, and mode-area transformers. The air holes provide an opportunity not available in standard fibers: the creation of dramatic morphological changes by altering the hole size by collapse (under surface tension) or inflation (under internal overpressure) when heating to the softening temperature of the glass. Thus, not only can the fiber be stretched locally to reduce its cross-sectional area, but the microstructure can itself be radically altered.

Cleaving and Splicing

PCF cleaves cleanly using standard tools, showing slight end-face distortion only when the core crystal is extremely small (interhole spacing $\sim 1 \mu\text{m}$) and the air-filling fraction very high ($>50\%$). Solid glass end-caps can be formed by collapsing the holes (or filling them with sol-gel glass) at the fiber end to form a core-less structure through which light can be launched into the fiber. Solid-core PCF can be fusion-spliced successfully both to itself and to step-index fiber using resistive heating elements (electric-arcs do not allow sufficient control). The two fiber ends are placed in intimate contact and heated to softening point. With careful control, they fuse together without distortion. Provided the mode areas are well matched, splice losses of <0.2 dB can normally be achieved except when the core is extremely small ($<\sim 1.5 \mu\text{m}$). Fusion splicing hollow-core fiber is feasible when there is a thick solid glass outer sheath (e.g., as depicted in Fig. 10b), although very low splice losses can be obtained simply by placing identical fibers end-to-end and clamping them (the index-matching “fluid” for hollow-core PCF is vacuum).

The ability to hermetically splice gas-filled hollow-core PCF to SMF has made it possible to produce in-line gas cells for stimulated Raman scattering in hydrogen and frequency measurement and stabilization (using acetylene). These developments may lead for the first time to practical miniature gas-laser devices that could even be coiled up inside a credit card.¹²²

Mode Transformers

In many applications it is important to be able to change the mode area without losing light. This is done traditionally using miniature bulk optics—tiny lenses precisely designed to match to a desired numerical aperture and spot-size. In PCF an equivalent effect can be obtained by scanning a heat source (flame or carbon dioxide laser) along the fiber. This causes the holes to collapse, the degree of collapse depending on the dwell-time of the heat. Drawing the two fiber ends apart at the same time provides additional control. Graded transitions can fairly easily be made—mode diameter reductions as high as 5:1 have been realized with low loss.

Ferrule methods have been developed for making low loss interfaces between conventional single-mode fibers and photonic crystal fibers.¹²³ Adapted from the fabrication of PCF preforms from stacked tubes and rods, these techniques avoid splicing and are versatile enough to interface to virtually any type of index-guiding silica PCF. They are effective for coupling light into and out of all of the individual cores of a multicore fiber without input or output crosstalk. The technique also creates another opportunity—the use of taper transitions to couple light between a multimode fiber and several single-mode fibers. When the number of single-mode fibers matches the number of spatial modes in the multimode fiber, the transition can have low loss in both directions. This means that the high performance of single-mode fiber devices can be reached in multimode systems, for example a multimode fiber filter with the transmission spectrum of a single-mode fiber Bragg grating,¹²⁴ a device that has applications in earth-based astronomy, where the high throughput of a multimode fiber can be retained while unwanted atmospheric emission lines are filtered out.

A further degree of freedom may be gained by pressurizing the holes during the taper process.¹²⁵ The resulting hole inflation permits radical changes in the guidance characteristics. It is possible for example to transform a PCF, with a relatively large core and small air-filling fraction, into a PCF with a very small core and a large air-filling fraction, the transitions having very low loss.¹²⁶

Heating and stretching PCF can result in quite remarkable changes in the scale of the micro/nano structures, without significant distortion. Recently a solid-core PCF was reduced 5 times in linear scale, resulting in a core diameter of 500 nm (Fig. 22). This permitted formation of a PCF with a zero dispersion wavelength that matched the 532 nm emission wavelength of a frequency doubled Nd:YAG laser¹⁰² (this is important for supercontinuum generation—as discussed in Sec. 11.8). A further compelling advantage of the tapering approach is that it neatly side-steps the difficulty of launching light into submicron sized cores; light is launched into the entry port (in this case with core diameter 2.5 μm) and adiabatically evolves, with negligible loss, into the mode of the 500 nm core.

In-Fiber Devices

Precise use of heat and pressure induces large changes in the optical characteristics of PCF, giving rise to a whole family of new intrafiber components. Microcouplers can be made in a PCF with two optically isolated cores by collapsing the holes so as to allow the mode fields to expand and

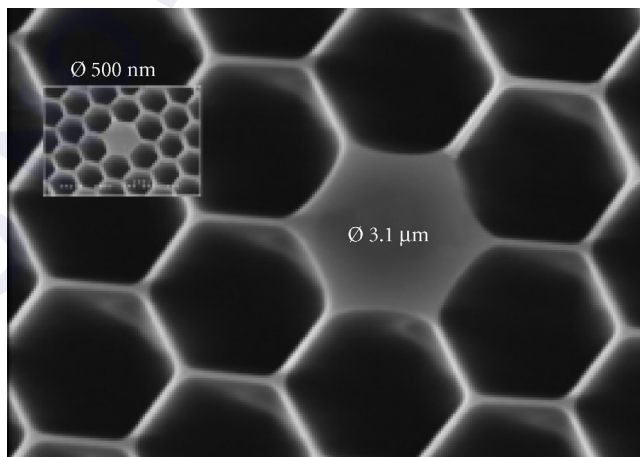


FIGURE 22 Scanning electron micrographs (depicted to the same scale) of the fiber cross sections produced by tapering a solid-core PCF. The structure is very well preserved, even down to core diameters of 500 nm.¹⁰²

interact with each other, creating local coupling.¹²⁷ Long period gratings, which scatter the core light into cladding modes within certain wavelength bands, can be made by periodic modulation of hole size.¹²⁸ By rocking a birefringent PCF to and fro, while scanning a carbon dioxide laser along it, so-called “rocking filters” can be made, which transfer power from one polarization state to the other within a narrow band of wavelengths.¹²⁹ All these components have one great advantage over equivalent devices made in conventional fiber: being formed by permanent changes in morphology, they are highly stable with temperature and over time.

11.10 CONCLUSIONS

In moving away from the constraints of conventional fibers, photonic crystal fibers have created new opportunities spanning many areas of science and technology. Some of the more important from an every widening range of applications (as of late 2006) are discussed in Ref. 1. Compact high brightness supercontinuum sources and frequency comb systems for optical clocks are already available as commercial products. Emerging applications include constrained photochemistry or biochemistry and microfluidics, biophotonic and biomedical devices, medical imaging, astronomy, particle delivery, high power fiber lasers, gas-based fiber devices, and fiber delivery of high power laser light in manufacturing. The next decade should see many of these applications mature into fields of scientific research or even commercial products.

11.11 APPENDIX

The analysis in Ref. 16 to estimate n_{\max} in a triangular lattice of air holes leads to the equation:

$$w I_1(a_n w)[J_1(bu)Y_0(a_n u) - J_0(a_n u)Y_1(bu)] + u I_0(a_n w)[J_1(bu)Y_1(a_n u) - J_1(a_n u)Y_1(bu)] = 0 \quad (14)$$

where $a_n = d/2\Lambda$, $b = (\sqrt{3}/2\pi)^{1/2}$ (needed to ensure the correct value of n_{\max} in the long wavelength limit) and $u^2 + w^2 = v^2$. The leading root of Eq. (14), evaluated for $a_n = 0.2$ and $n_g = 1.444$ (i.e., neglecting dispersion), yields the polynomial fit (7).

11.12 REFERENCES

1. P. St.J. Russell, “Photonic-Crystal Fibers,” *Journal of Lightwave Technology*, **24**:4729–4749, 2006.
2. J. C. Knight, T. A. Birks, P. St.J. Russell, and D. M. Atkin, “Pure Silica Single-Mode Fiber with Hexagonal Photonic Crystal Cladding,” in *Conference on Optical Fiber Communications* San Jose, California: Optical Society of America, 1996.
3. P. St.J. Russell, “Photonic Crystal Fibers,” *Science*, **299**:358–362, Jan. 2003.
4. J. C. Knight, “Photonic Crystal Fibres,” *Nature*, **424**:847–851, Aug. 2003.
5. J. C. Knight, T. A. Birks, P. St.J. Russell, and D. M. Atkin, “All-Silica Single-Mode Optical Fiber with Photonic Crystal Cladding,” *Optics Letters*, **21**:1547–1549, Oct. 1996.
6. J. C. Knight, T. A. Birks, P. St.J. Russell, and D. M. Atkin, “All-Silica Single-Mode Optical Fiber with Photonic Crystal Cladding: Errata,” *Optics Letters*, **22**:484–485, Apr. 1997.
7. R. F. Cregan, B. J. Mangan, J. C. Knight, T. A. Birks, P. St.J. Russell, P. J. Roberts, and D. C. Allan, “Single-Mode Photonic Band Gap Guidance of Light in Air,” *Science*, **285**:1537–1539, Sep. 1999.
8. V. V. R. K. Kumar, A. K. George, W. H. Reeves, J. C. Knight, P. St.J. Russell, F. G. Omenetto, and A. J. Taylor, “Extruded Soft Glass Photonic Crystal Fiber for Ultrabroad Supercontinuum Generation,” *Optics Express*, **10**:1520–1525, Dec. 2002.

9. J. C. Knight, J. Arriaga, T. A. Birks, A. Ortigosa-Blanch, W. J. Wadsworth, and P. St.J. Russell, "Anomalous Dispersion in Photonic Crystal Fiber," *IEEE Photonics Technology Letters*, **12**:807–809, Jul. 2000.
10. www.corning.com/opticalfiber/.
11. E. M. Dianov and V. M. Mashinsky, "Germania-Based Core Optical Fibers," *Journal of Lightwave Technology*, **23**:3500–3508, 2005.
12. P. St.J. Russell, "Designing Photonic Crystals," in *Electron and Photon Confinement in Semiconductor Nanostructures*, Amsterdam: IOS Press, pp. 79–103, 2003.
13. C. M. Bowden, J. P. Dowling, and H. O. Everitt, "Development and Applications of Materials Exhibiting Photonic Band Gaps," *Journal of the Optical Society of America B—Optical Physics*, **10**:280–413, 1993.
14. P. St.J. Russell, in *NATO Advanced Study Institute on Confined Electrons and Holes*, Erice, Sicily, 1993.
15. T. A. Birks, P. J. Roberts, P. St.J. Russell, D. M. Atkin, and T. J. Shepherd, "Full 2-D Photonic Band Gaps in Silica/Air Structures," *Electronics Letters*, **31**:1941–1942, 1995.
16. T. A. Birks, J. C. Knight, and P. St.J. Russell, "Endlessly Single-Mode Photonic Crystal Fiber," *Optics Letters*, **22**:961–963, Jul. 1997.
17. J. C. Knight, T. A. Birks, R. F. Cregan, P. St.J. Russell, and J. P. de Sandro, "Large Mode Area Photonic Crystal Fibre," *Electronics Letters*, **34**:1347–1348, Jun. 1998.
18. D. Mogilevtsev, T. A. Birks, and P. St.J. Russell, "Group-Velocity Dispersion in Photonic Crystal Fibers," *Optics Letters*, **23**:1662–1664, Nov. 1998.
19. A. Ortigosa-Blanch, A. Diez, M. Delgado-Pinar, J. L. Cruz, and M. V. Andres, "Ultrahigh Birefringent Non-linear Microstructured Fiber," *IEEE Photonics Technology Letters*, **16**:1667–1669, Jul. 2004.
20. B. J. Mangan, J. C. Knight, T. A. Birks, P. St.J. Russell, and A. H. Greenaway, "Experimental Study of Dual-core Photonic Crystal Fibre," *Electronics Letters*, **36**:1358–1359, Aug. 2000.
21. J. K. Ranka, R. S. Windeler, and A. J. Stentz, "Visible Continuum Generation in Air–silica Microstructure Optical Fibers with Anomalous Dispersion at 800 nm," *Optics Letters*, **25**:25–27, 2000.
22. V. N. Melekin and A. B. Manenkov, "Dielectric Tube as a Low-Loss Waveguide," *Soviet Physics—Technical Physics*, **13**:1698–1699, 1968.
23. P. Yeh and A. Yariv, "Bragg Reflection Waveguides," *Optics Communications*, **19**:427–430, 1976.
24. F. Brechet, P. Roy, J. Marcou, and D. Pagnoux, "Single-Mode Propagation in Depressed-Core-Index Photonic-Bandgap Fiber Designed for Zero-Dispersion Propagation at Short Wavelengths," *Electronics Letters*, **36**:514–515, 2000.
25. S. G. Johnson, M. Ibanescu, M. Skorobogatiy, O. Weisberg, T. D. Engeness, M. Soljacic, S. A. Jacobs, J. D. Joannopoulos, and Y. Fink, "Low-Loss Asymptotically Single-Mode Propagation in Large-Core Omniguide Fibers," *Optics Express*, **9**:748–779, 2001.
26. D. C. Allan, J. A. West, J. C. Fajardo, M. T. Gallagher, K. W. Koch, and N. F. Borrelli, "Photonic Crystal Fibers: Effective Index and Bandgap Guidance," in *Photonic Crystals and Light Localisation in the 21st Century*, C. M. Soukoulis (ed.), Kluwer Academic Publishers, Netherlands, pp. 305–320, 2001.
27. V. V. R. K. Kumar, A. K. George, J. C. Knight, and P. St.J. Russell, "Tellurite Photonic Crystal Fiber," *Optics Express*, **11**:2641–2645, Oct. 2003.
28. K. M. Kiang, K. Frampton, M. Monro, R. Moore, J. Tucknott, D. W. Hewak, D. J. Richardson, and H. N. Rutt, "Extruded Single-Mode Non-Silica Glass Holey Optical Fibers," *Electronics Letters*, **38**:546–547, 2002.
29. R. Bise and D. J. Trevor, "Sol-Gel Derived Microstructured Fiber: Fabrication and Characterization," in *Conference on Optical Fiber Communication*, Anaheim, Paper OWL6 2005.
30. M. C. J. Large, A. Argyros, F. Cox, M. A. van Eijkelenborg, S. Ponrathnam, N. S. Pujari, I. M. Bassett, R. Lwin, and G. W. Barton, "Microstructured Polymer Optical Fibres: New Opportunities and Challenges," *Molecular Crystals and Liquid Crystals*, **446**:219–231, 2006.
31. E. F. Chillce, C. M. B. Cordeiro, L. C. Barbosa, and C. H. B. Cruz, "Tellurite Photonic Crystal Fiber Made by a Stack-and-Draw Technique," *Journal of Non-Crystalline Solids*, **352**:3423–3428, Sep. 2006.
32. X. Feng, T. M. Monro, V. Finazzi, R. C. Moore, K. Frampton, P. Petropoulos, and D. J. Richardson, "Extruded Single-Mode, High-Nonlinearity, Tellurite Glass Holey Fibre," *Electronics Letters*, **41**:835–837, Jul. 2005.
33. D. A. Gaponov and A. S. Biryukov, "Optical Properties of Microstructure Tellurite Glass Fibres," *Quantum Electronics*, **36**:343–348, Apr. 2006.

34. T. P. White, B. Kuhlmeiy, R. C. McPhedran, D. Maystre, G. Renversez, C. M. De Sterke, and L. C. Botten, "Multipole Method for Microstructured Optical Fibers: 1. Formulation," *Journal of the Optical Society of America B—Optical Physics*, **19**:2322–2330, 2002.
35. R. C. McPhedran, L. C. Botten, A. A. Asatryan, N. A. P. Nicorovici, P. A. Robinson, and C. M. De Sterke, "Calculation of Electromagnetic Properties of Regular and Random Arrays of Metallic and Dielectric Cylinders," *Physical Review E*, **60**:7614–7617, 1999.
36. P. J. Roberts and T. J. Shepherd, "The Guidance Properties of Multi-Core Photonic Crystal Fibers," *Journal of Optics a—Pure and Applied Optics*, **3**:S1–S8, 2001.
37. A. Ferrando, E. Silvestre, J. J. Miret, P. Andres, and M. V. Andres, "Full-Vector Analysis of a Realistic Photonic Crystal Fiber," *Optics Letters*, **24**:276–278, Mar. 1999.
38. D. Mogilevtsev, T. A. Birks, and P. St.J. Russell, "Localized Function Method for Modeling Defect Modes in 2-D Photonic Crystals," *Journal of Lightwave Technology*, **17**:2078–2081, Nov. 1999.
39. A. W. Snyder and J. D. Love, *Optical Waveguide Theory*, London: Chapman & Hall, 1983.
40. T. A. Birks, D. M. Bird, T. D. Hedley, J. M. Pottage, and P. St.J. Russell, "Scaling Laws and Vector Effects in Bandgap-Guiding Fibres," *Optics Express*, **12**:69–74, Jan. 2004.
41. R. D. Meade, A. M. Rappe, K. D. Brommer, J. D. Joannopoulos, and O. L. Alerhand, "Accurate Theoretical Analysis of Photonic Band-Gap Materials," *Physical Review B*, **48**:8434–8437, 1993.
42. G. J. Pearce, T. D. Hedley, and D. M. Bird, "Adaptive Curvilinear Coordinates in a Plane-Wave Solution of Maxwell's Equations in Photonic Crystals," *Physical Review B*, **71**:195108, 2005.
43. J. M. Pottage, D. M. Bird, T. D. Hedley, T. A. Birks, J. C. Knight, P. St.J. Russell, and P. J. Roberts, "Robust Photonic Band Gaps for Hollow Core Guidance in PCF Made from High Index Glass," *Optics Express*, **11**:2854–2861, Nov. 2003.
44. T. M. Monro, D. J. Richardson, N. G. R. Broderick, and P. J. Bennett, "Holey Optical Fibers: An Efficient Modal Model," *Journal of Lightwave Technology*, **17**:1093–1102, 1999.
45. C. T. Chan, Q. L. Yu, and K. M. Ho, "Order N Spectral Method for Electromagnetic Waves," *Physical Review B*, **51**:16635–16642, 1995.
46. V. Dangui, M. J. F. Digonnet, and G. S. Kino, "A Fast and Accurate Numerical Tool to Model the Modal Properties of Photonic-Bandgap Fibers," *Optics Express*, **14**:2979–2993, Apr. 2006.
47. C. Mias, J. P. Webb, and R. L. Ferrari, "Finite Element Modelling of Electromagnetic Waves in Doubly and Triply Periodic Structures," *IEE Proceedings—Optoelectronics*, **146**:111–118, 1999.
48. A. Hochman and Y. Leviatan, "Analysis of Strictly Bound Modes in Photonic Crystal Fibers by Use of a Source-Model Technique," *Journal of the Optical Society of America A—Optics Image Science and Vision*, **21**:1073–1081, 2004.
49. J. C. Knight, T. A. Birks, P. St.J. Russell, and J. P. de Sandro, "Properties of Photonic Crystal Fiber and the Effective Index Model," *Journal of the Optical Society of America A—Optics Image Science and Vision*, **15**:748–752, Mar. 1998.
50. J. C. Knight, J. Broeng, T. A. Birks, and P. St.J. Russell, "Photonic Band Gap Guidance in Optical Fibers," *Science*, **282**:1476–1478, Nov. 1998.
51. F. Benabid, J. C. Knight, G. Antonopoulos, and P. St.J. Russell, "Stimulated Raman Scattering in Hydrogen-Filled Hollow-Core Photonic Crystal Fiber," *Science*, **298**:399–402, Oct. 2002.
52. F. Couny, F. Benabid, and P. S. Light, "Large-Pitch Kagome-Structured Hollow-Core Photonic Crystal Fiber," *Optics Letters*, **31**:3574–3576, Dec. 2006.
53. B. J. Mangan, J. Arriaga, T. A. Birks, J. C. Knight, and P. St.J. Russell, "Fundamental-Mode Cutoff in a Photonic Crystal Fiber with a Depressed-Index Core," *Optics Letters*, **26**:1469–1471, Oct. 2001.
54. W. N. MacPherson, M. J. Gander, R. McBride, J. D. C. Jones, P. M. Blanchard, J. G. Burnett, A. H. Greenaway, et al., "Remotely Addressed Optical Fibre Curvature Sensor Using Multicore Photonic Crystal Fibre," *Optics Communications*, **193**:97–104, Jun. 2001.
55. J. Limpert, N. D. Robin, I. Manek-Honninger, F. Salin, F. Roser, A. Liem, T. Schreiber, et al., "High-Power Rod-Type Photonic Crystal Fiber Laser," *Optics Express*, **13**:1055–1058, Feb. 2005.
56. W. J. Wadsworth, R. M. Percival, G. Bouwmans, J. C. Knight, and P. St.J. Russell, "High Power Air-Clad Photonic Crystal Fiber Laser," *Optics Express*, **11**:48–53, Jan. 2003.
57. M. T. Myaing, J. Y. Ye, T. B. Norris, T. Thomas, J. R. Baker, W. J. Wadsworth, G. Bouwmans, J. C. Knight, and P. St.J. Russell, "Enhanced Two-Photon Biosensing with Double-Clad Photonic Crystal Fibers," *Optics Letters*, **28**:1224–1226, Jul. 2003.

58. P. J. Roberts, F. Couny, H. Sabert, B. J. Mangan, D. P. Williams, L. Farr, M. W. Mason, et al., "Ultimate Low Loss of Hollow-Core Photonic Crystal Fibres," *Optics Express*, **13**:236–244, Jan. 2005.
59. F. Benabid, G. Bouwmans, J. C. Knight, P. St.J. Russell, and F. Couny, "Ultrahigh Efficiency Laser Wavelength Conversion in a Gas-Filled Hollow Core Photonic Crystal Fiber by Pure Stimulated Rotational Raman Scattering in Molecular Hydrogen," *Physical Review Letters*, **93**:123903 Sep. 2004.
60. G. Humbert, J. C. Knight, G. Bouwmans, P. St.J. Russell, D. P. Williams, P. J. Roberts, and B. J. Mangan, "Hollow Core Photonic Crystal Fibers for Beam Delivery," *Optics Express*, **12**:1477–1484, Apr. 2004.
61. P. J. Roberts, F. Couny, H. Sabert, B. J. Mangan, T. A. Birks, J. C. Knight, and P. St.J. Russell, "Loss in Solid-Core Photonic Crystal Fibers Due to Interface Roughness Scattering," *Optics Express*, **13**:7779–7793, Oct. 2005.
62. J. A. West, C. Smith, N. F. Borrelli, D. C. Allan, and K. W. Koch, "Surface Modes in Air-Core Photonic Band-Gap Fibers," *Optics Express*, **12**:1485–1496, 2004.
63. A. Argyros, T. A. Birks, S. G. Leon-Saval, C. M. B. Cordeiro, F. Luan, and P. St.J. Russell, "Photonic Bandgap with an Index Step of One Percent," *Optics Express*, **13**:309–314, Jan. 2005.
64. J. C. Knight, F. Luan, G. J. Pearce, A. Wang, T. A. Birks, and D. M. Bird, "Solid Photonic Bandgap Fibres and Applications," *Japanese Journal of Applied Physics Part 1—Regular Papers Short Notes & Review Papers*, **45**:6059–6063, Aug. 2006.
65. F. Luan, A. K. George, T. D. Hedley, G. J. Pearce, D. M. Bird, J. C. Knight, and P. St.J. Russell, "All-Solid Photonic Bandgap Fiber," *Optics Letters*, **29**:2369–2371, Oct. 2004.
66. N. M. Litchinitser, S. C. Dunn, B. Usner, B. J. Eggleton, T. P. White, R. C. McPhedran, and C. M. De Sterke, "Resonances in Microstructured Optical Waveguides," *Optics Express*, **11**:1243–1251, 2003.
67. A. Argyros, T. A. Birks, S. G. Leon-Saval, C. M. B. Cordeiro, and P. St.J. Russell, "Guidance Properties of Low-Contrast Photonic Bandgap Fibres," *Optics Express*, **13**:2503–2511, Apr. 2005.
68. G. Bouwmans, L. Bigot, Y. Quiquempois, F. Lopez, L. Provino, and M. Douay, "Fabrication and Characterization of an All-Solid 2D Photonic Bandgap Fiber with a Low-Loss Region (<20 dB/km) around 1550 nm," *Optics Express*, **13**:8452–8459, 2005.
69. M. J. Steel, T. P. White, C. M. De Sterke, R. C. McPhedran, and L. C. Botten, "Symmetry and Degeneracy in Microstructured Optical Fibers," *Optics Letters*, **26**:488–490, 2001.
70. A. Ortigosa-Blanch, J. C. Knight, W. J. Wadsworth, J. Arriaga, B. J. Mangan, T. A. Birks, and P. St.J. Russell, "Highly Birefringent Photonic Crystal Fibers," *Optics Letters*, **25**:1325–1327, Sep. 2000.
71. H. Kubota, S. Kawanishi, S. Koyanagi, M. Tanaka, and S. Yamaguchi, "Absolutely Single Polarization Photonic Crystal Fiber," *IEEE Photonics Technology Letters*, **16**:182–184, Jan. 2004.
72. C. Kerbage, P. Steinvurzel, P. Reyes, P. S. Westbrook, R. S. Windeler, A. Hale, and B. J. Eggleton, "Highly Tunable Birefringent Microstructured Optical Fiber," *Optics Letters*, **27**:842–844, May. 2002.
73. P. J. Roberts, D. P. Williams, H. Sabert, B. J. Mangan, D. M. Bird, T. A. Birks, J. C. Knight, and P. St.J. Russell, "Design of Low-Loss and Highly Birefringent Hollow-Core Photonic Crystal Fiber," *Optics Express*, **14**:7329–7341, Aug. 2006.
74. D. R. Chen and L. F. Shen, "Ultrahigh Birefringent Photonic Crystal Fiber with Ultralow Confinement Loss," *IEEE Photonics Technology Letters*, **19**:185–187, Jan.–Feb. 2007.
75. A. Michie, J. Canning, K. Lytikainen, M. Aslund, and J. Digweed, "Temperature Independent Highly Birefringent Photonic Crystal Fibre," *Optics Express*, **12**:5160–5165, Oct. 2004.
76. T. Martynkien, M. Szpulak, and W. Urbanczyk, "Modeling and Measurement of Temperature Sensitivity in Birefringent Photonic Crystal Holey Fibers," *Applied Optics*, **44**:7780–7788, Dec. 2005.
77. D. H. Kim and J. U. Kang, "Sagnac Loop Interferometer Based on Polarization Maintaining Photonic Crystal Fiber with Reduced Temperature Sensitivity," *Optics Express*, **12**:4490–4495, Sep. 2004.
78. W. H. Reeves, J. C. Knight, P. St.J. Russell, and P. J. Roberts, "Demonstration of Ultra-Flattened Dispersion in Photonic Crystal Fibers," *Optics Express*, **10**:609–613, Jul. 2002.
79. W. H. Reeves, D. V. Skryabin, F. Biancalana, J. C. Knight, P. St.J. Russell, F. G. Omenetto, A. Efimov, and A. J. Taylor, "Transformation and Control of Ultra-Short Pulses in Dispersion-Engineered Photonic Crystal Fibres," *Nature*, **424**:511–515, Jul. 2003.
80. G. Bouwmans, F. Luan, J. C. Knight, P. St.J. Russell, L. Farr, B. J. Mangan, and H. Sabert, "Properties of a Hollow-Core Photonic Bandgap Fiber at 850 nm Wavelength," *Optics Express*, **11**:1613–1620, Jul. 2003.
81. K. Kurokawa, K. Tajima, and K. Nakajima, "10-GHz 0.5-ps Pulse Generation in 1000-nm Band in PCF for High-Speed Optical Communication," *Journal of Lightwave Technology*, **25**:75–78, Jan. 2007.

82. K. Kurokawa, K. Tajima, K. Tsujikawa, and K. Nakagawa, "Penalty-Free Dispersion-Managed Soliton Transmission Over a 100-km Low-Loss PCF," *Journal of Lightwave Technology*, **24**:32–37, 2006.
83. G. J. Pearce, J. M. Pottage, D. M. Bird, P. J. Roberts, J. C. Knight, and P. St.J. Russell, "Hollow-Core PCF for Guidance in the Mid to Far Infrared," *Optics Express*, **13**:6937–6946, Sep. 2005.
84. J. D. Shephard, W. N. MacPherson, R. R. J. Maier, J. D. C. Jones, D. P. Hand, M. Mohebbi, A. K. George, P. J. Roberts, and J. C. Knight, "Single-Mode Mid-IR Guidance in a Hollow-Core Photonic Crystal Fiber," *Optics Express*, **13**:7139–7144, Sep. 2005.
85. P. J. Roberts, D. P. Williams, B. J. Mangan, H. Sabert, F. Couny, W. J. Wadsworth, T. A. Birks, J. C. Knight, and P. St.J. Russell, "Realizing Low Loss Air Core Photonic Crystal Fibers by Exploiting an Antiresonant Core Surround," *Optics Express*, **13**:8277–8285, Oct. 2005.
86. M. D. Nielsen, N. A. Mortensen, M. Albertsen, J. R. Folkenberg, A. Bjarklev, and D. Bonaccini, "Predicting Macrobanding Loss for Large-Mode Area Photonic Crystal Fibers," *Optics Express*, **12**:1775–1779, 2004.
87. J. Laegsgaard, N. A. Mortensen, J. Riishede, and A. Bjarklev, "Material Effects in Air-Guiding Photonic Bandgap Fibers," *Journal of the Optical Society of America B—Optical Physics*, **20**:2046–2051, 2003.
88. M. Onishi, T. Okuno, T. Kashiwada, S. Ishikawa, N. Akasaka, and M. Nishimura, "Highly Nonlinear Dispersion-Shifted Fibers and Their Application to Broadband Wavelength Converter," *Optical Fiber Technology*, **4**:204–214, 1998.
89. P. Petropoulos, H. Ebendorff-Heidepriem, V. Finazzi, R. Moore, K. Frampton, D. J. Richardson, and M. Monro, "Highly Nonlinear and Anomalous Dispersive Lead Silicate Glass Holey Fibers," *Optics Express*, **11**:3568–3573, 2003.
90. F. Luan, J. C. Knight, P. St.J. Russell, S. Campbell, D. Xiao, D. T. Reid, B. J. Mangan, D. P. Williams, and P. J. Roberts, "Femtosecond Soliton Pulse Delivery at 800 nm Wavelength in Hollow-Core Photonic Bandgap Fibers," *Optics Express*, **12**:835–840, Mar. 2004.
91. C. J. Hensley, D. G. Ouzounov, A. L. Gaeta, N. Venkataraman, M. T. Gallagher, and K. W. Koch, "Silica-Glass Contribution to the Effective Nonlinearity of Hollow-Core Photonic Band-Gap Fibers," *Optics Express*, **15**:3507–3512, Mar. 2007.
92. A. Efimov, A. J. Taylor, F. G. Omenetto, A. V. Yulin, N. Y. Joly, F. Biancalana, D. V. Skryabin, J. C. Knight, and P. St.J. Russell, "Time-Spectrally-Resolved Ultrafast Nonlinear Dynamics in Small-Core Photonic Crystal Fibers: Experiment and Modelling," *Optics Express*, **12**:6498–6507, Dec. 2004.
93. G. P. Agrawal, *Nonlinear Fiber Optics*, 4th ed., Academic Press, San Diego, CA, 2007.
94. R. R. Alfaro, *The Supercontinuum Laser Source*, Springer-Verlag, New York, 1989.
95. S. V. Chernikov, Y. Zhu, J. R. Taylor, and V. P. Gapontsev, "Supercontinuum Self-Q-Switched Ytterbium Fiber Laser," *Optics Letters*, **22**:298–300, 1997.
96. J. M. Dudley, G. Genty, and S. Coen, "Supercontinuum Generation in Photonic Crystal Fiber," *Reviews of Modern Physics*, **78**:1135–1184, Oct.–Dec. 2006.
97. W. J. Wadsworth, A. Ortigosa-Blanch, J. C. Knight, T. A. Birks, T. P. M. Man, and P. St.J. Russell, "Supercontinuum Generation in Photonic Crystal Fibers and Optical Fiber Tapers: A Novel Light Source," *Journal of the Optical Society of America B—Optical Physics*, **19**:2148–2155, Sep. 2002.
98. G. Humbert, W. J. Wadsworth, S. G. Leon-Saval, J. C. Knight, T. A. Birks, P. St.J. Russell, M. J. Lederer, et al., "Supercontinuum Generation System for Optical Coherence Tomography Based on Tapered Photonic Crystal Fibre," *Optics Express*, **14**:1596–1603, Feb. 2006.
99. I. Hartl, X. D. Li, C. Chudoba, R. K. Ghanta, T. H. Ko, J. G. Fujimoto, J. K. Ranka, and R. S. Windeler, "Ultra-high-Resolution Optical Coherence Tomography Using Continuum Generation in an Air-Silica Microstructure Optical Fiber," *Optics Letters*, **26**:608–610, May. 2001.
100. H. Hundertmark, D. Kracht, D. Wandt, C. Fallnich, V. V. R. K. Kumar, A. K. George, J. C. Knight, and P. St.J. Russell, "Supercontinuum Generation with 200 pJ Laser Pulses in an Extruded SF6 Fiber at 1560 nm," *Optics Express*, **11**:3196–3201, Dec. 2003.
101. R. Holzwarth, T. Udem, T. W. Haensch, J. C. Knight, W. J. Wadsworth, and P. St.J. Russell, "Optical Frequency Synthesizer for Precision Spectroscopy," *Physical Review Letters*, **85**:2264–2267, Sep. 2000.
102. S. G. Leon-Saval, T. A. Birks, W. J. Wadsworth, P. St.J. Russell, and M. W. Mason, "Supercontinuum Generation in Submicron Fibre Waveguides," *Optics Express*, **12**:2864–2869, Jun. 2004.
103. S. Coen, A. H. L. Chau, R. Leonhardt, J. D. Harvey, J. C. Knight, W. J. Wadsworth, and P. St.J. Russell, "Supercontinuum Generation by Stimulated Raman Scattering and Parametric Four-Wave Mixing in Photonic Crystal Fibers," *Journal of the Optical Society of America B—Optical Physics*, **19**:753–764, Apr. 2002.

104. W. J. Wadsworth, N. Joly, J. C. Knight, T. A. Birks, F. Biancalana, and P. St.J. Russell, "Supercontinuum and Four-Wave Mixing with Q-Switched Pulses in Endlessly Single-Mode Photonic Crystal Fibres," *Optics Express*, **12**:299–309, Jan. 2004.
105. www.fianium.com.
106. F. G. Omenetto, N. A. Wolchover, M. R. Wehner, M. Ross, A. Efimov, A. J. Taylor, V. Kumar, et al., "Spectrally Smooth Supercontinuum from 350 nm to 3 μ m in Sub-Centimeter Lengths of Soft-Glass Photonic Crystal Fibers," *Optics Express*, **14**:4928–4934, May 2006.
107. M. Yu, C. J. McKinstrie, and G. P. Agrawal, "Modulational Instabilities in Dispersion-Flattened Fibers," *Physical Review E*, **52**:1072–1080, 1995.
108. A. Y. H. Chen, G. K. L. Wong, S. G. Murdoch, R. Leonhardt, J. D. Harvey, J. C. Knight, W. J. Wadsworth, and P. St.J. Russell, "Widely Tunable Optical Parametric Generation in a Photonic Crystal Fiber," *Optics Letters*, **30**:762–764, Apr. 2005.
109. Y. J. Deng, Q. Lin, F. Lu, G. P. Agrawal, and W. H. Knox, "Broadly Tunable Femtosecond Parametric Oscillator Using a Photonic Crystal Fiber," *Optics Letters*, **30**:1234–1236, May 2005.
110. J. Lasri, P. Devgan, R. Y. Tang, J. E. Sharping, and P. Kumar, "A Microstructure-Fiber-Based 10-GHz Synchronized Tunable Optical Parametric Oscillator in the 1550-nm Regime," *IEEE Photonics Technology Letters*, **15**:1058–1060, Aug. 2003.
111. J. E. Sharping, M. Fiorentino, P. Kumar, and R. S. Windeler, "Optical Parametric Oscillator Based on Four-Wave Mixing in Microstructure Fiber," *Optics Letters*, **27**:1675–1677, Oct. 2002.
112. J. D. Harvey, R. Leonhardt, S. Coen, G. K. L. Wong, J. C. Knight, W. J. Wadsworth, and P. St.J. Russell, "Scalar Modulation Instability in the Normal Dispersion Regime by Use of a Photonic Crystal Fiber," *Optics Letters*, **28**:2225–2227, Nov. 2003.
113. J. Fulconis, O. Alibart, W. J. Wadsworth, P. St.J. Russell, and J. G. Rarity, "High Brightness Single Mode Source of Correlated Photon Pairs Using a Photonic Crystal Fiber," *Optics Express*, **13**:7572–7582, Sep. 2005.
114. D. V. Skryabin, F. Luan, J. C. Knight, and P. St.J. Russell, "Soliton Self-Frequency Shift Cancellation in Photonic Crystal Fibers," *Science*, **301**:1705–1708, Sep. 2003.
115. N. Y. Joly, F. G. Omenetto, A. Efimov, A. J. Taylor, J. C. Knight, and P. St.J. Russell, "Competition Between Spectral Splitting and Raman Frequency Shift in Negative-Dispersion Slope Photonic Crystal Fiber," *Optics Communications*, **248**:281–285, Apr. 2005.
116. V. Laude, A. Khelif, S. Benchbane, M. Wilm, T. Sylvestre, B. Kibler, A. Mussot, J. M. Dudley, and H. Maillotte, "Phononic Band-Gap Guidance of Acoustic Modes in Photonic Crystal Fibers," *Physical Review B*, **71**:045107, 2005.
117. S. Guenneau and A. B. Movchan, "Analysis of Elastic Band Structures for Oblique Incidence," *Archive for Rational Mechanics and Analysis*, **171**:129–150, 2004.
118. P. St.J. Russell, E. Marin, A. Diez, S. Guenneau, and A. B. Movchan, "Sonic Band Gaps in PCF Preforms: Enhancing the Interaction of Sound and Light," *Optics Express*, **11**:2555–2560, Oct. 2003.
119. P. St.J. Russell, "Light in a Tight Space: Enhancing Matter-Light Interactions Using Photonic Crystals," in *Conference on Nonlinear Optics*, Hawaii, pp. 377–379, 2002.
120. P. Dainese, P. St.J. Russell, N. Joly, J. C. Knight, G. S. Wiederhecker, H. L. Fragnito, V. Laude, and A. Khelif, "Stimulated Brillouin Scattering from Multi-GHz-Guided Acoustic Phonons in Nanostructured Photonic Crystal Fibres," *Nature Physics*, **2**:388–392, Jun. 2006.
121. P. Dainese, P. St.J. Russell, G. S. Wiederhecker, N. Joly, H. L. Fragnito, V. Laude, and A. Khelif, "Raman-Like Light Scattering from Acoustic Phonons in Photonic Crystal Fiber," *Optics Express*, **14**:4141–4150, May 2006.
122. F. Benabid, F. Couny, J. C. Knight, T. A. Birks, and P. St.J. Russell, "Compact, Stable and Efficient All-Fibre Gas Cells Using Hollow-Core Photonic Crystal Fibres," *Nature*, **434**:488–491, Mar. 2005.
123. S. G. Leon-Saval, T. A. Birks, N. Y. Joly, A. K. George, W. J. Wadsworth, G. Kakarantzias, and P. St.J. Russell, "Splice-Free Interfacing of Photonic Crystal Fibers," *Optics Letters*, **30**:1629–1631, Jul. 2005.
124. S. G. Leon-Saval, T. A. Birks, J. Bland-Hawthorn, and M. Englund, "Multimode Fiber Devices with Single-Mode Performance," *Optics Letters*, **30**:2545–2547, 2005.
125. T. A. Birks, G. Kakarantzias, P. St.J. Russell, and D. F. Murphy, "Photonic Crystal Fiber Devices," *Proceedings of the Society of Photo-Instrumentation Engineers*, **4943**:142–151, 2002.

126. W. J. Wadsworth, A. Witkowska, S. Leon-Saval, and T. A. Birks, "Hole Inflation and Tapering of Stock Photonic Crystal Fibers," *Optics Express*, **13**:6541–6549, 2005.
127. G. Kakarantzas, T. E. Dimmick, T. A. Birks, R. Le Roux, and P. St.J. Russell, "Miniature All-Fiber Devices Based on CO₂ Laser Microstructuring of Tapered Fibers," *Optics Letters*, **26**:1137–1139, Aug. 2001.
128. G. Kakarantzas, T. A. Birks, and P. St.J. Russell, "Structural Long-Period Gratings in Photonic Crystal Fibers," *Optics Letters*, **27**:1013–1015, Jun. 2002.
129. G. Kakarantzas, A. Ortigosa-Blanch, T. A. Birks, P. St.J. Russell, L. Farr, F. Couny, and B. J. Mangan, "Structural Rocking Filters in Highly Birefringent Photonic Crystal Fiber," *Optics Letters*, **28**:158–160, Feb. 2003.

DO NOT DUPLICATE

James A. Harrington

*Rutgers University
Piscataway, New Jersey*

12.1 INTRODUCTION

Infrared (IR) optical fibers may be defined as fiber optics transmitting radiation with wavelengths greater than approximately $2\ \mu\text{m}$. The first IR fibers were fabricated in the mid-1960s from chalcogenide glasses such as arsenic trisulfide and had losses in excess of $10\ \text{dB/m}$.¹ During the mid-1970s, the interest in developing an efficient and reliable IR fiber for short-haul applications increased, partly in response to the need for a fiber to link broadband, long-wavelength radiation to remote photodetectors in military sensor applications. In addition, there was an ever-increasing need for a flexible fiber delivery system for transmitting CO_2 laser radiation in surgical applications. Around 1975, a variety of IR materials and fibers were developed to meet these needs. These included the heavy metal fluoride glass (HMFG) and polycrystalline fibers as well as hollow rectangular waveguides. While none of these fibers had physical properties even approaching those of conventional silica fibers, they were nevertheless useful in lengths less than 2 to 3 m for a variety of IR sensor and power delivery applications.²

IR fiber optics may logically be divided into three broad categories: glass, crystalline, and hollow waveguides. These categories may be further subdivided based on fiber material, structure, or both, as shown in Table 1. Over the past 30 years many novel IR fibers have been made in an effort to fabricate a fiber optic with properties as close as possible to those of silica, but only a relatively small number have survived. A good source of general information on these various IR fiber types may be found in the literature.²⁻⁶ In this review only the best, most viable, and, in most cases, commercially available IR fibers are discussed. In general, both the optical and mechanical properties of IR fibers remain inferior to those of silica fibers, and therefore the use of IR fibers is still limited primarily to nontelecommunication, short-haul applications requiring only tens of meters of fiber rather than the kilometer lengths common to telecommunication applications. The short-haul nature of IR fibers results from the fact that most IR fibers have losses in the range of a few decibels per meter. An exception is fluoride glass fibers, which can have losses as low as a few decibels per kilometer. In addition, IR fibers are much weaker than silica fiber and, therefore, more fragile. These deleterious features have slowed the acceptance of IR fibers and restricted their use today to applications in chemical sensing, thermometry, and laser power delivery.

A key feature of current IR fibers is their ability to transmit longer wavelengths than most oxide glass fibers can. In some cases the transmittance of the fiber can extend well beyond $20\ \mu\text{m}$, but most applications do not require the delivery of radiation longer than about $12\ \mu\text{m}$. In Fig. 1 we give the

TABLE 1 Categories of IR Fibers with a Common Example to Illustrate Each Subcategory

Main	Subcategory	Examples
Glass	Heavy metal fluoride (HMFG)	ZrF ₄ -BaF ₂ -LaF ₃ -AlF ₃ -NAF(ZBLAN)
	Germanate	GeO ₂ -PbO
	Chalcogenide	As ₂ S ₃ and AsGeTeSe
Crystal	Polycrystalline (PC)	AgBrCl
	Single crystal (SC)	Sapphire
Hollow waveguide	Metal/dielectric film	Hollow glass waveguide
	Refractive index <1	Hollow sapphire at 10.6 μm

attenuation values for some of the most common IR fibers as listed in Table 1. From the data it is clear that there is a wide variation in range of transmission for the different IR fibers and that there is significant extrinsic absorption that degrades the overall optical response. Most of these extrinsic bands can be attributed to various impurities, but, in the case of the hollow waveguides, they are due to interference effects resulting from the thin film coatings used to make the guides.

Some of the other physical properties of IR fibers are listed in Table 2. For comparison, the properties of silica fibers are also listed. The data in Table 2 and in Fig. 1 reveal that, compared to silica, IR fibers usually have higher losses, larger refractive indices and dn/dT values, lower melting or softening points, and greater thermal expansion. For example, chalcogenide and polycrystalline Ag halide fibers have refractive indices greater than 2. This means that the Fresnel loss exceeds 20 percent for two fiber ends. The higher dn/dT and low melting or softening point lead to thermal lensing and low laser-induced damage thresholds for some of the fibers. Finally, a number of these fibers do not have cladding analogous to clad oxide glass fibers. Nevertheless, core-only IR fibers such as sapphire and chalcogenide fibers can still be useful because their refractive indices are sufficiently high. For these high-index fibers, the energy is largely confined to the core of the fiber as long as the unprotected fiber core does not come in contact with an absorbing medium.¹²

The motivation to develop a viable IR fiber stems from many proposed applications. A summary of the most important current and future applications and the associated candidate IR fiber that will best meet each need is given in Table 3. We may note several trends from this table. The first is that

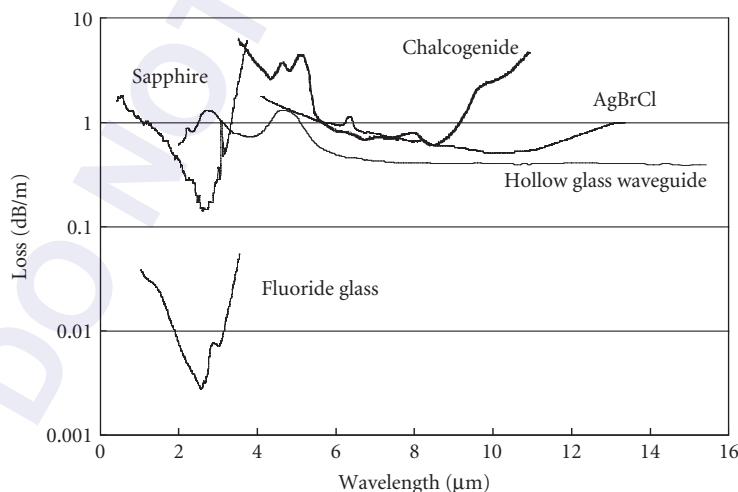


FIGURE 1 Composite loss spectra for some common IR fiber optics: ZBLAN fluoride glass,⁷ SC sapphire,⁸ chalcogenide glass,⁹ PC AgBrCl,¹⁰ and hollow glass waveguide.¹¹

TABLE 2 Selected Physical Properties of Key IR Fibers Compared to Conventional Silica Fiber

Property	Glass			Crystal		Hollow
	Silica	HMFG ZBLAN	Chalcogenide AsGeSeTe	PC AgBrCl	SC Sapphire	Hollow Silica Waveguide
Glass transition or melting point, °C	1175	265	245	412	2030	150 (usable <i>T</i>)
Thermal conductivity, W/m °C	1.38	0.628	0.2	1.1	36	1.38
Thermal expansion coefficient, 10 ⁻⁶ °C ⁻¹	0.55	17.2	15	30	5	0.55
Young's modulus, GPa	70.0	58.3	21.5	0.14	430	70.0
Density, g/cm ³	2.20	4.33	4.88	6.39	3.97	2.20
Refractive index (λ , μm)	1.455 (0.70)	1.499 (0.589)	2.9 (10.6)	2.2 (10.6)	1.71 (3.0)	NA
dn/dT , 10 ⁻⁵ °C ⁻¹ (λ , μm)	+1.2 (1.06)	-1.5 (1.06)	+10 (10.6)	-1.5 (10.6)	+1.4 (1.06)	NA
Fiber transmission range, μm	0.24–2.0	0.25–4.0	4–11	3–16	0.5–3.1	0.9–25
Loss* at 2.94 μm , dB/m	~800	0.08	5	3	0.4	0.5
Loss* at 10.6 μm , dB/m	NA	NA	2	0.5	NA	0.4

*Typical measured loss.
NA = not applicable.

TABLE 3 Examples of IR Fiber Candidates for Various Sensor and Power Delivery Applications

Application	Comments	Suitable IR fibers
Fiber-optic chemical sensors	Evanescent wave principle—liquids Hollow core waveguides—gases	AgBrCl, sapphire, chalcogenide, HMFG Hollow glass waveguides
Radiometry	Blackbody radiation, temperature measurements	Hollow glass waveguides, AgBrCl, chalcogenide, sapphire
Er:YAG laser power delivery	3- μm transmitting fibers with high damage threshold	Hollow glass waveguides, sapphire, germanate glass
CO ₂ laser power delivery	10- μm transmitting fibers with high damage threshold	Hollow glass waveguides
Thermal imaging	Coherent bundles	HMFG, chalcogenide
Fiber amplifiers and lasers	Doped IR glass fibers	HMFG, chalcogenide

hollow waveguides are an ideal candidate for laser power delivery at all IR laser wavelengths. The air core of these special fibers or waveguides gives an inherent advantage over solid-core fibers, whose damage threshold is frequently very low for these IR-transmissive materials. The high refractive index of chalcogenide fibers is ideal for chemical sensing via evanescent wave coupling of a small portion of the light from the core into an IR-absorbing medium. For the measurement of temperature through the simple transmission of blackbody radiation, IR fibers that transmit beyond about 8 μm , such as the Ag halide, chalcogenide, and hollow waveguides, are excellent candidates for use in measuring temperatures below 50°C. This is because the peak for room-temperature blackbody radiation is about 10 μm .

12.2 NONOXIDE AND HEAVY-METAL OXIDE GLASS IR FIBERS

There are two IR-transmitting glass fiber systems that are relatively similar to conventional silica-containing glass fibers. One is the HMFG and the other is heavy-metal germanate glass fibers based on GeO₂. The germanate glass fibers generally do not contain fluoride compounds; instead, they

contain heavy metal oxides to shift the IR absorption edge to longer wavelengths. The advantage of germanate fibers over HMFG fibers is that germanate glass has a higher glass transition temperature and, therefore, higher laser damage thresholds. But the level of loss for the HMFG fibers is lower. Finally, chalcogenide glass fibers made from chalcogen elements such as As, Ge, S, and Te contain no oxides or halides, making them a good choice for nonlaser power delivery applications.

HMFG Fibers

Poulain et al.¹³ discovered HMFGs or fluoride glasses accidentally in 1975 at the University of Rennes. In general, the typical fluoride glass has a glass transition temperature T_g four times less than that of silica, is considerably less stable than silica, and has failure strains of only a few percent compared to greater than 5 percent for silica. While an enormous number of multicomponent fluoride glass compositions have been fabricated, comparably few have been drawn into fiber. This is because the temperature range for fiber drawing is normally too small in most HMFGs to permit fiberization of the glass. The most popular HMFGs for fabrication into fibers are the fluorozirconate and fluoroaluminate glasses, of which the most common are ZrF_4 - BaF_2 - LaF_3 - AlF_3 - NaF (ZBLAN) and AlF_3 - ZrF_4 - BaF_2 - CaF_2 - YF_3 , respectively. The key physical properties of these glasses are summarized in Table 4. An important feature of the fluoroaluminate glass is its higher T_g , which largely accounts for the higher laser damage threshold for the fluoroaluminate glasses compared to ZBLAN at the Er:YAG laser wavelength of 2.94 μm .

The fabrication of HMFG fiber is similar to any glass fiber drawing technology except that the preforms are made using some type of melt-forming method rather than by a vapor deposition process as is common with silica fibers. Specifically, a casting method based on first forming a clad glass tube and then adding the molten core glass is used to form either multimode or single-mode fluorozirconate fiber preforms. The cladding tube is made either by a rotational casting technique in which the tube is spun in a metal mold or by merely inverting and pouring out most of the molten cladding glass contained in a metal mold to form a tube.¹⁴ The cladding tube is then filled with a higher-index core glass. Other preform fabrication techniques include rod-in-tube and crucible techniques. The fluoroaluminate fiber preforms have been made using an unusual extrusion technique in which core and cladding glass plates are extruded into a core-clad preform.¹⁵ All methods, however, involve fabrication from the melted glass rather than from the more pristine technique of vapor deposition used to form SiO_2 -based fibers. This process creates inherent problems such as the formation of bubbles, core-cladding interface irregularities, and small preform sizes. Most HMFG fiber drawing is done using preforms rather than the crucible method. A ZBLAN preform is drawn at about 310°C in a controlled atmosphere (to minimize contamination by moisture or oxygen impurities, which can significantly weaken the fiber) using a narrow heat zone compared to silica. Either ultraviolet (UV) acrylate or Teflon coatings are applied to the fiber. In the case of Teflon, heat-shrink Teflon (FEP) is generally applied to the glass preform prior to the draw.

The attenuation in HMFG fibers is predicted to be about 10 times less than that for silica fibers.¹⁶ Based on extrapolations of the intrinsic losses resulting from Rayleigh scattering and multiphonon

TABLE 4 Fluorozirconate vs. Fluoroaluminate Glasses

Property	Fluorozirconate (ZBLAN)	Fluoroaluminate (AlF_3 - ZrF_4 - BaF_2 - CaF_2 - YF_3)
Glass transition temperature, °C	265	400
Durability	Medium	Excellent
Loss at 2.94 μm , dB/m	0.01	0.1
Er:YAG laser peak output energy, mJ	300 (300- μm core)	850 (500- μm core)

Comparison between fluorozirconate and fluoroaluminate glasses of some key properties that relate to laser power transmission and durability of the two HMFG fibers. Other physical properties are relatively similar.

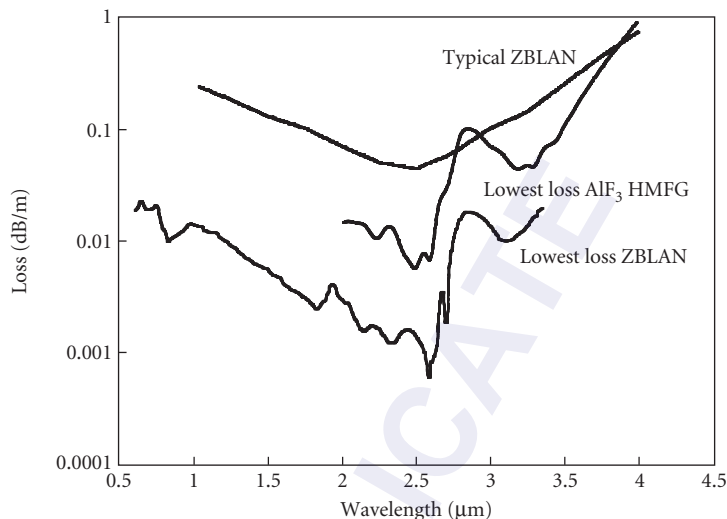


FIGURE 2 Losses in the best BTRL⁷ and typical (Infrared Fiber Systems, Silver Spring, Maryland) ZBLAN fluoride glass fibers compared to those for fluoroaluminate glass fibers.¹⁵

absorption, the minimum in the loss curves or V-curves is projected to be about 0.01 dB/km at 2.55 μm . Recent refinements of the scattering loss have modified this value slightly to be 0.024 dB/km, or about 8 times less than that for silica fiber.⁷ In practice, however, extrinsic loss mechanisms still dominate fiber loss. In Fig. 2, losses for two ZBLAN fibers are shown. The data from British Telecom (BTRL) represents state-of-the-art fiber 110 m in length.⁷ The other curve is more typical of commercially available (Infrared Fiber Systems, Silver Spring, Maryland) ZBLAN fiber. More recently commercially available fiber (IRphotonics, Montreal, CA) has become available with a loss of 0.02 dB/m at 3 μm .¹⁷ The lowest measured loss for a BTRL 60-m-long fiber is 0.45 dB/km at 2.3 μm . Some of the extrinsic absorption bands that contribute to the total loss shown in Fig. 2 for the BTRL fiber are Ho^{3+} (0.64 and 1.95 μm), Nd^{3+} (0.74 and 0.81 μm), Cu^{2+} (0.97 μm), and OH^- (2.87 μm). Scattering centers such as crystals, oxides, and bubbles have also been found in the HMFG fibers. In their analysis of the data in Fig. 2, the BTRL group separated the total minimum attenuation coefficient (0.65 dB/km at 2.59 μm) into an absorptive loss component equal to 0.3 dB/km and a scattering loss component equal to 0.35 dB/km. The losses for the fluoroaluminate glass fibers are also shown for comparison in Fig. 2.¹⁵ Clearly, the losses are not as low as for the BTRL-ZBLAN fiber, but the AlF_3 -based fluoride fibers do have the advantage of higher glass transition temperatures and therefore are better candidates for laser power delivery.

The reliability of HMFG fibers depends on protecting the fiber from attack by moisture and on pretreatment of the preform to reduce surface crystallization. In general, the HMFGs are much less durable than oxide glasses. The leach rates for ZBLAN glass range between 10^{-3} and 10^{-2} g/cm²/day. This is about five orders of magnitude higher than the leach rate for Pyrex glass. The fluoroaluminate glasses are more durable, with leach rates that are more than three times lower than those for the fluorozirconate glasses. The strength of HMFG fibers is less than that of silica fibers. From Table 2 we see that Young's modulus E for fluoride glass is 51 GPa compared to 73 GPa for silica glass. Taking the theoretical strength to be about one-fifth that of Young's modulus gives a theoretical value of strength of 11 GPa for fluoride glass. The largest bending strength measured has been about 1.4 GPa, well below the theoretical value. To estimate the bending radius R , we may use the approximate expression $R = 1.198r(E/\sigma_{\text{max}})$, where σ_{max} is the maximum fracture stress and r is the fiber radius.¹⁸

Germanate Fibers

Heavy metal oxide glass fibers based on GeO_2 have recently shown great promise as an alternative to HMFG fibers for 3- μm laser power delivery.¹⁹ Today, GeO_2 -based glass fibers are composed of GeO_2 (30–76 percent)–RO (15–43 percent)–XO (3–20 percent), where R represents an alkaline earth metal and X represents an element of Group IIIA.²⁰ In addition, small amounts of heavy metal fluorides may be added to the oxide mixture.²¹ The oxide-only germanate glasses have glass transition temperatures as high as 680°C, excellent durability, and a relatively high refractive index of 1.84. In Fig. 3, loss data is given for a typical germanate glass fiber. While the losses are not as low as they are for the fluoride glasses shown in Fig. 2, these fibers have an exceptionally high damage threshold at 3 μm . Specifically, over 20 W (2 J at 10 Hz) of Er:YAG laser power has been launched into these fibers.

Chalcogenide Fibers

Chalcogenide glass fibers were drawn into essentially the first IR fiber in the mid-1960s.¹ Chalcogenide fibers fall into three categories: sulfide, selenide, and telluride.²² One or more chalcogen elements are mixed with one or more elements such as As, Ge, P, Sb, Ga, Al, Si, and so on to form a glass having two or more components. From the data in Table 2 we see that the glasses have low softening temperatures more comparable to those of fluoride glass than to those of oxide glasses. Chalcogenide glasses are very stable, durable, and insensitive to moisture. A distinctive difference between these glasses and the other IR fiber glasses is that they do not transmit well in the visible region and their refractive indices are quite high. Additionally, most of the chalcogenide glasses, except for As_2S_3 , have a rather large value of dn/dT .⁹ This fact limits the laser power handling capability of the fibers. In general, chalcogenide glass fibers have proven to be an excellent candidate for evanescent wave fiber sensors and for IR fiber image bundles.²³

Chalcogenide glass is made by combining highly purified (>6 nines purity) raw elements in an ampoule that is heated in a rocking furnace for about 10 hours. After melting and mixing, the glass is quenched and a glass preform is fabricated using rod-in-tube or rotational casting methods. Preform

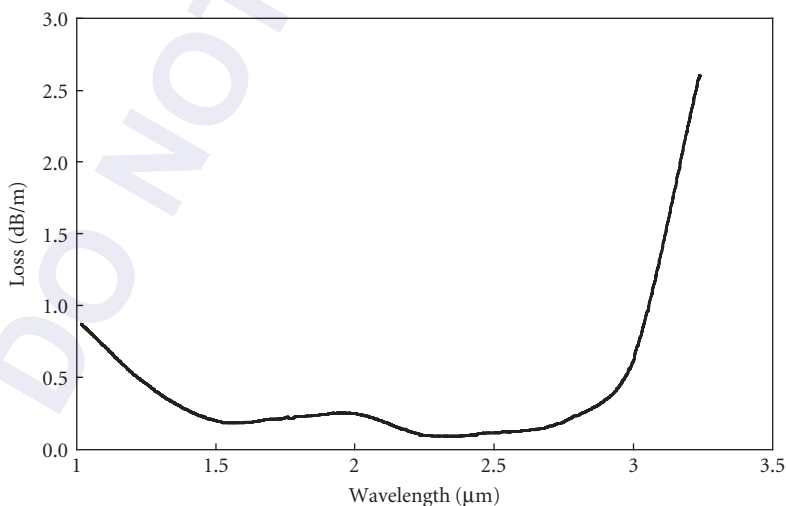


FIGURE 3 Germanate glass fiber manufactured by Infrared Fiber Systems, Silver Spring, Maryland.

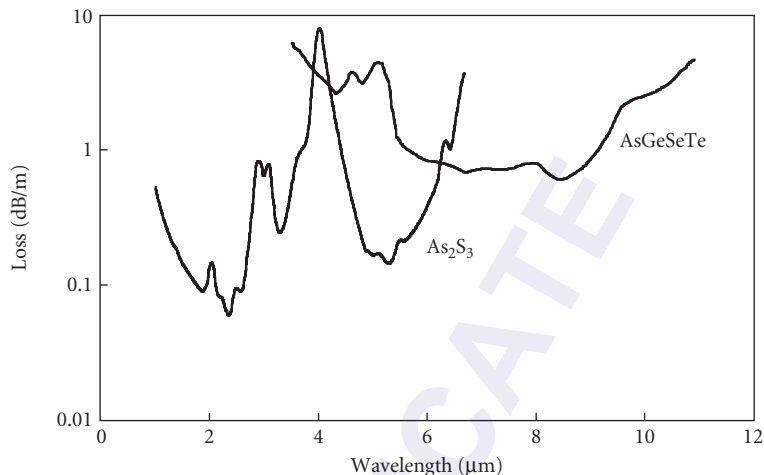


FIGURE 4 Two common chalcogenide glass fibers: As_2S_3 and an AsGeSeTe fiber.⁹ Note the many impurity bands pervasive in these fiber systems.

fiber draws involve drawing a core-clad preform or a core-only preform. For the core-only preform draw, either a soft chalcogenide cladding can be extruded over the fiber as it is drawn or the preform can be Teflon clad. Crucible drawing is also possible.

The losses for the most important chalcogenide fibers are given in Fig. 4. Arsenic trisulfide (As_2S_3) fiber, one of the simplest and oldest chalcogenide fibers, has a transmission range from 0.7 to about 6 μm .²² This fiber is red in color and therefore transmits furthest into the visible region but cuts off in the long-wavelength end well before the heavier chalcogenide fibers.⁹ Longer wavelengths are transmitted through the addition of heavier elements like Te, Ge, and Se, as shown in Fig. 4. A key feature of essentially all chalcogenide glasses is the strong extrinsic absorption resulting from the bonding of contaminants such as hydrogen, H_2O , and OH^- to the elemental cations. In particular, absorption peaks between 4.0 and 4.6 μm are due to S-H or Se-H bonds, and those at 2.78 and 6.3 μm are due to OH^- (2.78 μm) and/or molecular water. The hydride impurities are often especially strong and can be deleterious when these fibers are used in chemical sensing applications where the desired chemical signature falls in the region of extrinsic absorption. Another important feature of most of the chalcogenide fibers is that their losses are in general much higher than those of the fluoride glasses. In fact, at the important CO_2 laser wavelength of 10.6 μm , the lowest loss is still above 1 dB/m for the Se-based fibers.²² More recently single-mode chalcogenide fibers have been drawn from preforms made both using rotational casting and rod-in-tube methods.²⁴ It is also possible to form a chalcohalide glass which is composed of both chalcogen elements and a halide usually I⁻. One example is the TeSeAsI glass which has been drawn into fiber by Lucas' group in Rennes, France.²⁵

12.3 CRYSTALLINE FIBERS

Crystalline IR fibers are an attractive alternative to glass IR fibers because most nonoxide crystalline materials can transmit longer-wavelength radiation than IR glasses and, in the case of sapphire, exhibit some superior physical properties as well.² The disadvantage is that crystalline fibers are difficult to fabricate. There are two types of crystalline fiber: single-crystal (SC)⁸ and polycrystalline (PC).^{26,27} Historically, the first crystalline fiber made was hot-extruded KRS-5 fiber fabricated at Hughes Research Labs in 1975,²⁸ KRS-5 or TlBrI was chosen because it is very ductile and because

it can transmit beyond the 20- μm range required for the intended military surveillance satellite application. In fact, crystalline fibers such as KRS-5 and other halide crystals were initially thought to hold great potential as next-generation ultra-low-loss fibers because their intrinsic loss was predicted to be as low as 10^{-3} dB/m.²⁸ Unfortunately, this loss was not only never achieved but not even approached experimentally.

PC Fibers

There are many halide crystals that have excellent IR transmission, but only a few have been fabricated into fiber optics. The technique used to make PC fibers is hot extrusion. As a result, only the silver and thallium halides have the requisite physical properties (such as ductility, low melting point, and independent slip systems) to be successfully extruded into fiber. In the hot extrusion process, a single-crystal billet or preform is placed in a heated chamber and the fiber is extruded to net shape through a diamond or tungsten carbide die at a temperature equal to about half the melting point. The final PC fibers are usually from 500 to 900 μm in diameter with no buffer jacket. The polycrystalline structure of the fiber consists of grains on the order of 10 μm or larger in size. The billet may be clad using the rod-in-tube method. In this method, a mixed silver halide such as AgBrCl is used as the core and then a lower-index tube is formed using a Cl⁻-rich AgBrCl crystal. The extrusion of a core-clad fiber is not as easy to achieve as it is in glass drawing, but Artjushenko et al.¹⁰ at the General Physics Institute (GPI) in Moscow have achieved clad Ag halide fibers with losses nearly as low as those for the core-only Ag halide fiber. Single-mode PC fibers have been extruded from a rod-in-tube preform by Katzir and his group at Tel Aviv University. These fibers have a core diameter of 50 to 60 μm and they are single-mode at 10.6 μm .²⁹ Today, the PC Ag halide fibers represent the best PC fibers. KRS-5 is no longer a viable candidate due largely to the toxicity of Tl and the greater flexibility of the Ag halide fibers.

The losses for the Ag halide fibers are shown in Fig. 5. Both the core-only and core-clad fibers are shown, and, as with the other IR fibers, we again see that there are several extrinsic absorption bands. Water is often present at 3 and 6.3 μm and there is sometimes an SO₄⁻ absorption near 9.6 μm . Furthermore, we note the decreasing attenuation as the wavelength increases. This is a result of λ^{-2}

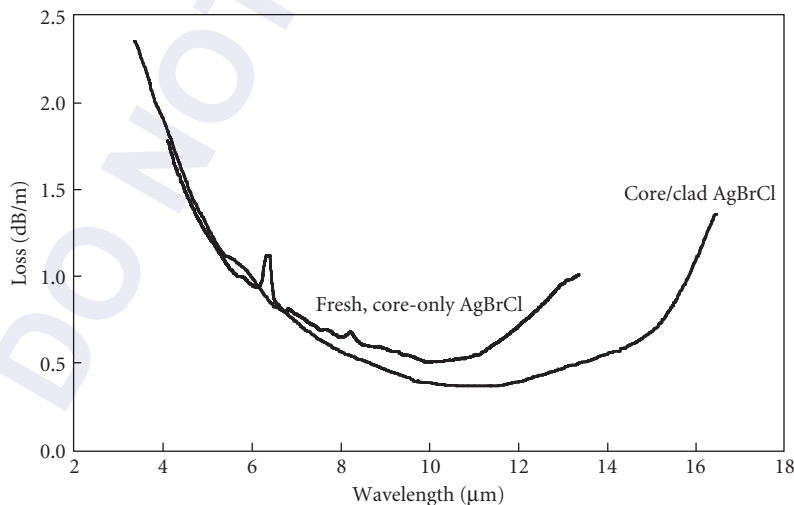


FIGURE 5 Losses in a typical PC silver halide fiber²⁶ compared to those in recently developed core-clad silver halide fiber.¹⁰

scattering from strain-induced defects in the extruded fiber. An important feature of the data is that the loss at 10.6 μm can be as low as 0.2 dB/m for the core-only fiber and these fibers will transmit to almost 20 μm . These fibers have been used to transmit about 100 W of CO_2 laser power, but the safe limit seems to be 20 to 25 W.³⁰ This is due to the low melting point of the fibers.

There are several difficulties in handling and working with PC fibers. One is an unfortunate aging effect in which the fiber transmission is observed to decrease over time.³¹ Normally the aging loss, which increases uniformly over the entire IR region, is a result of strain relaxation and possible grain growth as the fiber is stored. Another problem is that Ag halides are photosensitive; exposure to visible or UV radiation creates colloidal Ag, which in turn leads to increased losses in the IR. Finally, AgBrCl is corrosive to many metals. Therefore, the fibers should be packaged in dark jackets and connectorized with materials such as Ti, Au, or ceramics.

The mechanical properties of these ductile fibers are quite different from those of glass fibers. The fibers are weak, with ultimate tensile strengths of about 80 MPa for a 50–50 mixture of AgBrCl. However, the main difference between the PC and glass fibers is that the PC fibers plastically deform well before fracture. This plastic deformation leads to increased loss as a result of increased scattering from separated grain boundaries. Therefore, in use, the fibers should not be bent beyond their yield point; too much bending can lead to permanent damage and a region of high loss in the fiber.

SC Fibers

Meter-long lengths of SC fibers have been made from only a small number of the over 80 IR transmissive crystalline materials. Initially some SC fibers were grown by zone-refining methods from the same metal halides used to extrude PC fibers. The idea was that removal of the grain boundaries in the PC fibers would improve the optical properties of the fiber. This did not occur, so most of the crystalline materials chosen today for SC fiber fabrication have been oxides. Compared to halides, oxide materials like Al_2O_3 (sapphire) have the advantage of high melting points, chemical inertness, and the ability to be conveniently melted and grown in air. Currently, sapphire is the most popular SC fiber.^{8,32,33}

Sapphire is an insoluble, uniaxial crystal (trigonal structure) with a melting point of over 2000°C. It is an extremely hard and robust material with a usable fiber transmission from about 0.5 to 3.2 μm . Other important physical properties shown in Table 2 include a refractive index equal to 1.75 at 3 μm , a thermal expansion about 10 times higher than that of silica, and a Young's modulus approximately six times greater than that of silica. These properties make sapphire an almost ideal IR fiber candidate for applications less than about 3.2 μm . In particular, this fiber has been used to deliver over 10 W of average power from an Er:YAG laser operating at 2.94 μm .³⁴

Sapphire fibers are fabricated using either the edge-defined, film-fed growth (EFG) or the laser-heated pedestal growth (LHPG) techniques.³⁵ In either method, some or all of the starting sapphire material is melted and an SC fiber is pulled from the melt. In the EFG method, a capillary tube is used to conduct the molten sapphire to a seed fiber, which is drawn slowly into a long fiber. Multiple capillary tubes, which also serve to define the shape and diameter of the fiber, may be placed in one crucible of molten sapphire so that many fibers can be drawn at one time. The LHPG process is a crucibleless technique in which a small molten zone at the tip of an SC sapphire source rod (< 2 mm diameter) is created using a CO_2 laser. A seed fiber slowly pulls the SC fiber as the source rod continuously moves into the molten zone to replenish the molten material. Both SC fiber growth methods are very slow (several millimeters per minute) compared to glass fiber drawing. The EFG method, however, has an advantage over LHPG methods because more than one fiber can be continuously pulled at a time. LHPG methods, however, have produced the cleanest and lowest-loss fibers owing to the fact that no crucible is used that can contaminate the fiber. The sapphire fibers grown by these techniques are unclad, pure Al_2O_3 with the *C* axis usually aligned along the fiber axis. Fiber diameters range from 100 to 300 μm and lengths are generally less than 2 m. Postcladding techniques mostly involve a Teflon coating using heat-shrink tubing.

The optical properties of the as-grown sapphire fibers are normally inferior to those of the bulk starting material. This is particularly evident in the visible region and is a result of color-center-type

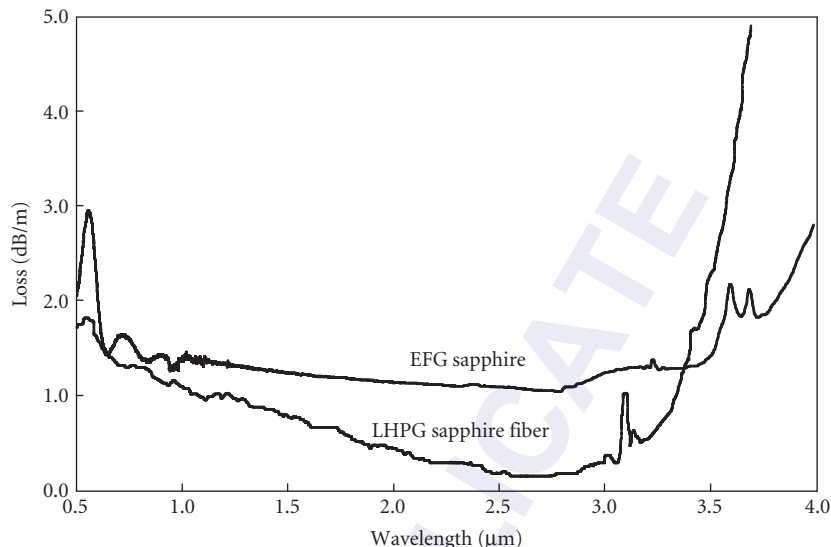


FIGURE 6 SC sapphire fibers grown by the EFG³⁵ (Saphikon, Inc., Milford, New Hampshire) and LHPG³⁴ methods.

defect formation during the fiber drawing. These defects and the resulting absorption can be greatly reduced if the fibers are postannealed in air or oxygen at about 1000°C. In Fig. 6, the losses for LHPG fiber grown at Rutgers University⁸ and EFG fiber grown by Photran, Inc. (Milford, New Hampshire) are shown. Both fibers have been annealed at 1000°C to reduce short-wavelength losses. We see that the LHPG fiber has the lowest overall loss. In particular, LHPG fiber loss at the important Er:YAG laser wavelength of 2.94 μm is less than 0.3 dB/m, compared to the intrinsic value of 0.15 dB/m. There are also several impurity absorptions beyond 3 μm that are believed to be due to transition metals like Ti or Fe. Sapphire fibers have been used at temperatures of up to 1400°C without any change in their transmission.

12.4 HOLLOW WAVEGUIDES

The first optical-frequency hollow waveguides were similar in design to microwave guides. Garmire et al.³⁶ made a simple rectangular waveguide using aluminum strips spaced 0.5 mm apart by bronze shim stock. Even when the aluminum was not well polished, these guides worked surprisingly well. Losses at 10.6 μm were well below 1 dB/m, and Garmire early demonstrated the high power handling capability of an air-core guide by delivering over 1 kW of CO₂ laser power through this simple structure. These rectangular waveguides, however, never gained much popularity, primarily because their overall dimensions (about 0.5 × 10 mm) were quite large in comparison to circular-cross-section guides and also because the rectangular guides cannot be bent uniformly in any direction. As a result, hollow circular waveguides with diameters of 1 mm or less fabricated using metal, glass, or plastic tubing are the most common guides today. In general, hollow waveguides are an attractive alternative to conventional solid-core IR fibers for laser power delivery because of the inherent advantage of their air core. Hollow waveguides not only enjoy the advantage of high laser power thresholds but also low insertion loss, no end reflection, ruggedness, and small beam divergence. A disadvantage, however, is a loss on bending, which varies as 1/R where R is the bending radius.

In addition, the losses for these guides vary as $1/a^3$ where a is the radius of the bore; therefore the loss can be arbitrarily small for a sufficiently large core. The bore size and bending radius dependence of all hollow waveguides are characteristics of these guides not shared by solid-core fibers. Initially these waveguides were developed for medical and industrial applications involving the delivery of CO_2 laser radiation, but more recently they have been used to transmit incoherent light for broadband spectroscopic and radiometric applications.^{37–39} Today they are one of the best alternatives for power delivery in IR laser surgery and industrial laser delivery systems, with losses as low as 0.1 dB/m and transmitted CW laser powers as high as 2.7 kW.⁴⁰

Hollow core waveguides may be grouped into two categories: (1) those whose inner core materials have refractive indices greater than 1 (leaky guides) and (2) those whose inner wall materials have refractive indexes less than 1 [attenuated total reflectance (ATR) guides]. Leaky or $n > 1$ guides have metallic and dielectric films deposited on the inside of metallic,⁴¹ plastic,⁴² or glass¹¹ tubing. ATR guides are made from dielectric materials with refractive indices of less than 1 in the wavelength region of interest.⁴³ Therefore, $n < 1$ guides are fiber-like in that the core index ($n \approx 1$) is greater than the clad index. Hollow sapphire fibers operating at 10.6 μm ($n = 0.67$) are an example of this class of hollow guide.⁴⁴

Hollow Metal and Plastic Waveguides

The earliest circular-cross-section hollow guides were formed using metallic and plastic tubing as the structural members. Miyagi and colleagues in Japan used sputtering methods to deposit Ge,⁴⁵ ZnSe, and ZnS⁴¹ coatings on aluminum mandrels. Then a final layer of Ni was electroplated over these coatings before the aluminum mandrel was removed by chemical leaching. The final structure was then a flexible Ni tube with optically thick dielectric layers on the inner wall to enhance the reflectivity in the infrared. Croitoru and colleagues⁴⁶ at Tel Aviv University applied Ag followed by AgI coatings on the inside of polyethylene and Teflon tubing to make a very flexible waveguide. Similar Ag and Ag halide coatings were deposited inside Ag tubes by Morrow and colleagues.⁴⁷

Hollow Glass Waveguides

The most popular structure today is the hollow glass waveguide (HGW) developed initially at Rutgers University.⁴⁸ The advantage of glass tubing is that it is much smoother than either metal and plastic tubing and, therefore, the scattering losses are less. HGWs are fabricated using wet chemistry methods to first deposit a layer of Ag on the inside of silica glass tubing and then to form a dielectric layer of AgI over the metallic film by converting some of the Ag to AgI. The silica tubing used has a polymer coating of UV acrylate or polyimide on the outside surface to preserve the mechanical strength. The thickness of the AgI is optimized to give high reflectivity at a particular laser wavelength or range of wavelengths. Using these techniques, HGWs have been fabricated with lengths as long as 13 m and bore sizes ranging from 250 to 1300 μm .

The spectral loss for an HGW with a 530- μm bore is given in Fig. 2. This HGW was designed for an optimal response at 10 μm . The peaks at about 3 and 5 μm are not absorption peaks but rather interference bands due to thin-film optical effects. For broadband applications and shorter-wavelength applications, a thinner AgI coating would be used to shift the interference peaks to shorter wavelengths. In fact, for the thinnest AgI films these HGWs not only transmit IR radiation but also shorter wavelengths down to nearly 0.5 μm .⁴⁹ For such HGWs the optical response will be nearly flat without interference bands in the far IR fiber region of the spectrum. The data in Fig. 7 shows the straight loss measured using a CO_2 and Er:YAG laser for different bore sizes. An important feature of this data is the $1/a^3$ dependence of loss on bore size predicted by the theory of Marcatili and Schmeltzer.⁵⁰ In general, the losses are less than 0.5 dB/m at 10 μm for bore sizes larger than ~ 400 μm . Furthermore, the data at 10.6 μm agrees well with the calculated values, but at 3 μm the measured losses are somewhat above those predicted by Marcatili and Schmeltzer. This is a result of increased scattering at the shorter wavelengths from the metallic and dielectric films. The bending loss depends

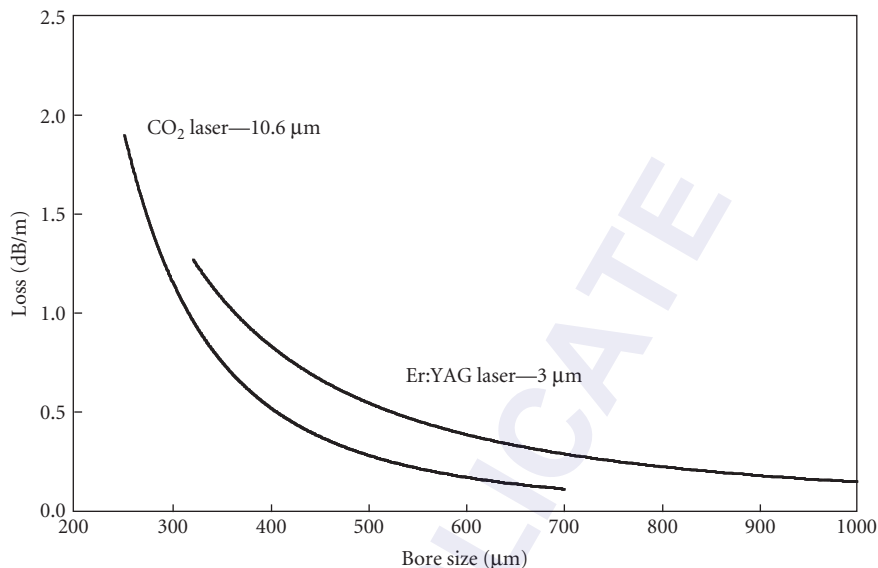


FIGURE 7 Straight losses measured in hollow glass waveguides with Ag/AgI films. The guide labeled “CO₂ laser” was designed for optimal transmission at 10.6 μm while that labeled “Er:YAG laser” was designed for optimal transmission at 3 μm. Note that the loss varies approximately as $1/a^3$.

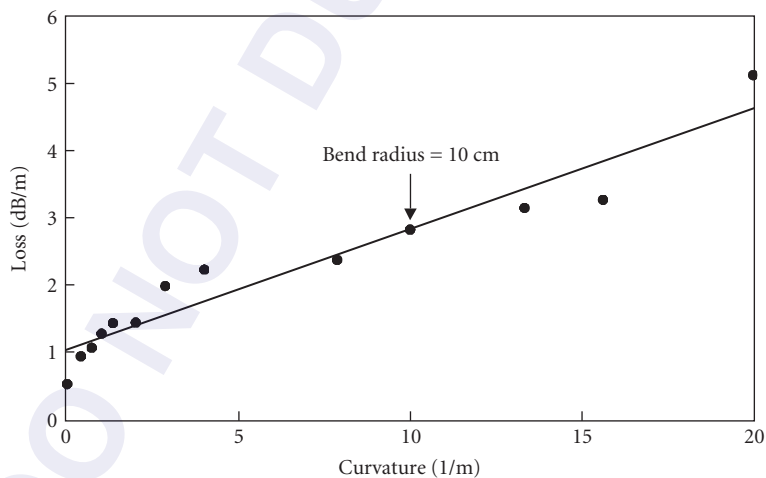


FIGURE 8 Additional loss on bending an HGW with a 530-μm bore, measured at 10.6 μm. The loss is seen to increase as the curvature increases.

on many factors such as the quality of the films, the bore size, and the uniformity of the silica tubing. A typical bending loss curve for an HGW with a 530-μm bore measured with a CO₂ laser is given in Fig. 8. The losses are seen to increase linearly with increasing curvature as predicted. It is important to note that while there is an additional loss on bending for any hollow guide, it does not necessarily mean that this restricts the use of hollow guides in power delivery or sensor applications. Normally

most fiber delivery systems have rather large bend radii and therefore a minimal amount of the guide is under tight bending conditions and the bending loss is low. From the data in Fig. 8 one can calculate the bending loss contribution for an HGW link by assuming some modest bends over a small section of guide length. An additional important feature of hollow waveguides is that they are nearly single mode. This is a result of the strong dependence of loss on the fiber mode parameter. That is, the loss of high-order modes increases as the square of the mode parameter, so even though the guides are very multimode, in practice only the lowest-order modes propagate. This is particularly true for the small-bore ($<300\ \mu\text{m}$) guides, in which virtually only the lowest-order HE_{11} mode is propagated.

HGWs have been used quite successfully in IR laser power delivery and, more recently, in some sensor applications. Modest CO_2 and Er:YAG laser powers below about 80 W can be delivered without difficulty. At higher powers, water-cooling jackets have been placed around the guides to prevent laser damage. The highest CO_2 laser power delivered through a water-cooled hollow metallic waveguide with a bore of $1800\ \mu\text{m}$ was 2700 W, and the highest power through a water-cooled HGW with a $700\text{-}\mu\text{m}$ bore was 1040 W.⁵¹ Sensor applications include gas and temperature measurements. A coiled HGW filled with gas can be used in place of a more complex and costly White cell to provide an effective means for gas analysis. Unlike evanescent wave spectroscopy, in which light is coupled out of a solid-core-only fiber into media in contact with the core, all of the light is passing through the gas in the hollow guide cell, making this a sensitive, quick-response fiber sensor. Temperature measurements may be aided by using an HGW to transmit blackbody radiation from a remote site to an IR detector. Such an arrangement has been used to measure jet engine temperatures.

12.5 SUMMARY AND CONCLUSIONS

During the past 30 years of the development of IR fibers, there has been a great deal of fundamental research designed to produce a fiber with optical and mechanical properties close to those of silica. We can see that today we are still far from that Holy Grail, but some viable IR fibers have emerged that, as a class, can be used to address some of the needs for a fiber that can transmit greater than $2\ \mu\text{m}$. Yet we are still limited with the current IR fiber technology by high loss and low strength. Nevertheless, more applications are being found for IR fibers as users become aware of their limitations and, more importantly, how to design around their properties.

There are two near-term applications of IR fibers: laser power delivery and sensors. An important future application for these fibers, however, may be more in active fiber systems like the Er- and Pr-doped fluoride fibers and emerging doped chalcogenide fibers. In regards to power delivery fibers, the best choice seems to be hollow waveguides for CO_2 lasers and SC sapphire, germanate glass, or HGWs for Er:YAG laser delivery. Chemical, temperature, and imaging bundles make use mostly of solid-core fibers. Evanescent wave spectroscopy (EWS) using chalcogenide and fluoride fibers is quite successful. A distinct advantage of an IR-fiber EWS sensor is that the signature of the analyte is often very strong in the infrared or fingerprint region of the spectrum. Temperature sensing generally involves the transmission of blackbody radiation. IR fibers can be very advantageous at low temperatures, especially near room temperature, where the peak in the blackbody radiation is near $10\ \mu\text{m}$. Finally, there is an emerging interest in IR imaging using coherent bundles of IR fibers. Several thousand chalcogenide fibers have been bundled by Amorphous Materials (Garland, Texas) to make an image bundle for the 3- to $10\text{-}\mu\text{m}$ region.

12.6 REFERENCES

1. N. S. Kapany and R. J. Simms, "Recent Developments of Infrared Fiber Optics," *Infrared Phys.* 5:69 (1965).
2. J. A. Harrington, *Infrared Fiber Optics and Their Applications*, SPIE Press, Bellingham, WA, 2004.
3. T. Katsuyama and H. Matsumura, *Infrared Optical Fibers*, Adam Hilger, Bristol, UK, 1989.
4. I. Aggarwal and G. Lu, *Fluoride Glass Optical Fiber*, Academic Press, San Diego, 1991.

5. P. France, M. G. Drexhage, J. M. Parker, M. W. Moore, S. F. Carter, and J. V. Wright, *Fluoride Glass Optical Fibers*, CRC Press, Boca Raton, Florida, 1990.
6. J. Sanghera and I. Aggarwal, *Infrared Fiber Optics*, CRC Press, Boca Raton, Florida, 1998.
7. S. F. Carter, M. W. Moore, D. Szebesta, D. Ransom, and P. W. France, "Low Loss Fluoride Fibre by Reduced Pressure Casting," *Electron. Lett.* **26**:2115–2117 (1990).
8. R. Nubling and J. A. Harrington, "Optical Properties of Single-Crystal Sapphire Fibers," *Appl. Opt.* **36**:5934–5940 (1997).
9. J. Nishii, S. Morimoto, I. Inagawa, R. Iizuka, T. Yamashita, and T. Yamagishi, "Recent Advances and Trends in Chalcogenide Glass Fiber Technology: A Review," *J. Non-Cryst. Sol.* **140**:199–208 (1992).
10. V. Artjushenko, V. Ionov, K. I. Kalaidjian, A. P. Kryukov, E. F. Kuzin, A. A. Lerman, A. S. Prokhorov, E. V. Stepanov, K. Bakhshpour, K. B. Moran, and W. Neuberger, "Infrared Fibers: Power Delivery and Medical Applications," *Proc. SPIE* **2396**:25–36 (1995).
11. Y. Matsuura, T. Abel, and J. A. Harrington, "Optical Properties of Small-Bore Hollow Glass Waveguides," *Appl. Opt.* **34**:6842–6847 (1995).
12. P. Kaiser, A. C. Hart Jr., and L. L. Blyler, "Low Loss FEP-Clad Silica Fibers," *Appl. Opt.* **14**:156 (1975).
13. M. Poulain, M. Chanthanasinh, and J. Lucas, "New Fluoride Glasses," *Mat. Res. Bull.* **12**:151–156 (1977).
14. D. Tran, G. H. Sigel, and B. Bendow, "Heavy Metal Fluoride Glasses and Fibers: A Review," *J. Light-Wave Technol.* **LT-2**:566–586 (1984).
15. K. Itoh, K. Miura, M. Masuda, M. Iwakura, and T. Yamagishi, "Low-Loss Fluorozirco-Aluminate Glass Fiber," in *Proceedings of 7th International Symposium on Halide Glass*, Center for Advanced Materials Technology, Monash University, Lorne, Australia, 1991, pp. 2.7–2.12.
16. P. W. France, S. F. Carter, M. W. Moore, and C. R. Day, "Progress in Fluoride Fibres for Optical Communications," *Brit. Telecom Tech. J.* **5**:28–44 (1987).
17. F. Sequin, M. Saad, P. Orsini, and D. Baierl, "Fluoride Glass Fiber for Reliable Er:YAG and Er, Cr:YSGG Laser Power Delivery," I. Gannot, ed., *Proc. SPIE* **6852**, (2008).
18. M. J. Matthewson, C. R. Kurkjian, and S. T. Gulati, "Strength Measurement of Optical Fibers by Bending," *J. Am. Cer. Soc.* **69**:815–821 (1986).
19. S. Kobayashi, N. Shibata, S. Shibata, and T. Izawa, "Characteristics of Optical Fibers in Infrared Wavelength Region," *Rev. Electrical Comm. Lab.* **26**:453–467, 1978.
20. D. Tran, Heavy Metal-Oxide Glass Optical Fibers for Use in Laser Medical Surgery, U.S. Patent 5,274,728, December 28, 1993.
21. G. Cao, F. Lin, H. Hu, and F. Gan, "A New Fluorogermanate Glass," *J. Non-Cryst. Solids* **326 & 327**:170–176, (2003).
22. Y. Kanamori, Y. Terunuma, and T. Miyashita, "Preparation of Chalcogenide Optical Fiber," *Rev. Electrical Comm. Lab.* **32**:469–477 (1984).
23. J. Nishii, T. Yamashita, T. Tamagishi, C. Tanaka, and H. Sone, "As₂S₃ Fibre for Infrared Image Bundle," *Int. J. Optoelectron.* **7**:209–216 (1992).
24. C. Boussard-Pledel, V. S. Shiryaev, P. Houizot, T. Jouan, J. L. Adam, and J. Lucas, "Single-Mode Infrared Fibers Based on TeAsSe Glass Systems," *Mater. Sci. Eng. B (Solid-State Materials for Advanced Technology)* **127**:138–143 (2006).
25. C. Blanchetiere, K. LeFoulgoc, H. L. Ma, X. H. Zhang, and J. Lucas, "Tellurium Halide Glass Fibers: Preparation and Application," *J. Non-Cryst. Solids* **184**:200–203 (1995).
26. V. G. Artjushenko, L. N. Butvina, V. V. Vojtsekhovskiy, E. M. Dianov, and J. G. Kolesnikov, "Mechanisms of Optical Losses in Polycrystalline KRS-5 Fibers," *J. Lightwave Technol.* **LT-4**:461–465 (1986).
27. A. Sa'ar, F. Moser, S. Akselrod, and A. Katzir, "Infrared Optical Properties of Polycrystalline Silver Halide Fibers," *Appl. Phys. Lett.* **49**:305–307 (1986).
28. D. A. Pinnow, A. L. Gentile, A. G. Standlee, A. J. Timper, and L. M. Hobrock, "Polycrystalline Fiber Optical Waveguides for Infrared Transmission," *Appl. Phys. Lett.* **33**:28–29 (1978).
29. S. Shalem, A. Tsun, E. Rave, A. Millo, I. Nagli, and A. Katzir, "Silver Halide Single-Mode Fibers for the Middle Infrared," *Appl. Phys. Lett.* **87**:91103-1–91103-3 (2005).
30. K. Takahashi, N. Yoshida, and M. Yokota, "Optical Fibers for Transmitting High-Power CO₂ Laser Beam," *Sumitomo Electric Tech. Rev.* **23**:203–210 (1984).

31. J. A. Wysocki, R. G. Wilson, A. G. Standlee, A. C. Pastor, R. N. Schwartz, A. R. Williams, G.-D. Lei, and L. Kevan, "Aging Effects in Bulk and Fiber TlBr-TlI," *J. Appl. Phys.* **63**:4365–4371 (1988).
32. D. H. Jundt, M. M. Fejer, and R. L. Byer, "Characterization of Single-Crystal Sapphire Fibers for Optical Power Delivery Systems," *Appl. Phys. Lett.* **55**:2170–2172 (1989).
33. R. S. F. Chang, V. Phomsakha, and N. Djeu, "Recent Advances in Sapphire Fibers," *Proc. SPIE* **2396**:48–53 (1995).
34. R. Nubling and J. A. Harrington, "Single-Crystal LHPG Sapphire Fibers for Er:YAG Laser Power Delivery," *Appl. Opt.* **37**:4777–4781 (1998).
35. H. E. LaBelle, "EFG, the Invention and Application to Sapphire Growth," *J. Cryst. Growth* **50**:8–17 (1980).
36. E. Garmire, T. McMahon, and M. Bass, "Flexible Infrared Waveguides for High-Power Transmission," *J. Quant. Elect.* **QE-16**:23–32 (1980).
37. S. J. Saggese, J. A. Harrington, and G. H. Sigel Jr., "Attenuation of Incoherent Infrared Radiation in Hollow Sapphire and Silica Waveguides," *Opt. Lett.* **16**:27–29 (1991).
38. M. Saito, Y. Matsuura, M. Kawamura, and M. Miyagi, "Bending Losses of Incoherent Light in Circular Hollow Waveguides," *J. Opt. Soc. Am. A* **7**:2063–2068 (1990).
39. M. Saito and K. Kikuchi, "Infrared Optical Fiber Sensors," *Opt. Rev.* **4**:527–538 (1997).
40. A. Hongo, K. Morosawa, K. Matsumoto, T. Shiota, and T. Hashimoto, "Transmission of Kilowatt-Class CO₂ Laser Light Through Dielectric-Coated Metallic Hollow Waveguides for Material Processing," *Appl. Opt.* **31**:5114–5120 (1992).
41. Y. Matsuura, M. Miyagi, and A. Hongo, "Fabrication of Low-Loss Zinc-Selenide Coated Silver Hollow Waveguides for CO₂ Laser Light," *J. Appl. Phys.* **68**:5463–5466 (1990).
42. M. Alaluf, J. Dror, R. Dahan, and N. Croitoru, "Plastic Hollow Fibers as a Selective Infrared Radiation Transmitting Medium," *J. Appl. Phys.* **72**:3878–3883 (1992).
43. C. C. Gregory and J. A. Harrington, "Attenuation, Modal, Polarization Properties of $n < 1$, Hollow Dielectric Waveguides," *Appl. Opt.* **32**:5302–5309 (1993).
44. J. A. Harrington and C. C. Gregory, "Hollow Sapphire Fibers for the Delivery of CO₂ Laser Energy," *Opt. Lett.* **15**:541–543 (1990).
45. M. Miyagi, Y. Shimada, A. Hongo, K. Sakamoto, and S. Nishida, "Fabrication and Transmission Properties of Electrically Deposited Germanium-Coated Waveguides for Infrared Radiation," *J. Appl. Phys.* **60**:454–456 (1986).
46. O. Morhaim, D. Mendlovic, I. Gannot, J. Dror, and N. Croitoru, "Ray Model for Transmission of Infrared Radiation through Multibent Cylindrical Waveguides," *Opt. Eng.* **30**:1886–1891 (1991).
47. P. Bhardwaj, O. J. Gregory, C. Morrow, G. Gu, and K. Burbank, "Performance of a Dielectric-Coated Monolithic Hollow Metallic Waveguide," *Mat. Lett.* **16**:150–156 (1993).
48. T. Abel, J. Hirsch, and J. A. Harrington, "Hollow Glass Waveguides for Broadband Infrared Transmission," *Opt. Lett.* **19**:1034–1036 (1994).
49. K. E. Sui, W. Shi, X. L. Tang, X. S. Shu, K. Iwai, and M. Miyagi, "Optical Properties of AgI/Ag Infrared Hollow Fiber in the Visible Wavelength Region," *Opt. Lett.* **33**:318–320 (2008).
50. E. A. J. Marcanti and R. A. Schmeltzer, "Hollow Metallic and Dielectric Waveguides for Long Distance Optical Transmission and Lasers," *Bell Syst. Tech. J.* **43**:1783–1809 (1964).
51. R. K. Nubling and J. A. Harrington, "Hollow-Waveguide Delivery Systems for High-Power, Industrial CO₂ Lasers," *Appl. Opt.* **34**:372–380 (1996).

This page intentionally left blank.

DO NOT DUPLICATE

SOURCES, MODULATORS, AND DETECTORS FOR FIBER OPTIC COMMUNICATION SYSTEMS

Elsa Garmire

*Dartmouth College
Hanover, New Hampshire*

13.1 INTRODUCTION

Optical communication systems utilize fiber optics to transmit the light that carries the signals. Such systems require optoelectronic devices as sources and detectors of such light, and they need modulators to impress the telecommunication signals onto the light. This chapter outlines the basics of these devices. Characteristics of devices designed for both high-performance, high-speed telecommunication systems (*telecom*) and low-cost, more modest performance local area networks (LAN) and data communication systems (*datacom*) are presented. Sources for telecom are edge-emitting laser diodes (LDs), including double heterostructure (DH), quantum well (QW), strained layer (SL), distributed feedback (DFB), and distributed Bragg reflector (DBR) laser diodes. Operating characteristics of these edge-emitting LDs include threshold, light-out versus current-in, spatial, and spectral characteristics. The transient response includes relaxation oscillations, turn-on delay, and modulation response. The noise characteristics are described by relative intensity noise (RIN), signal-to-noise ratio (SNR), mode partition noise (in multimode LDs), and phase noise (which determines linewidth). Frequency chirping broadens the linewidth, described in the small- and large-signal regime; external optical feedback may profoundly disturb the stability of the LDs and may lead to coherence collapse.

Semiconductor lasers usually have a laser cavity in the plane of the semiconductor device, with light emitted out through a cleaved edge in an elliptical output pattern. This output is not ideally suited to coupling into fibers, which have circular apertures. Low-cost systems, such as datacom, put a premium on simplicity in optical design. These systems typically use multimode fibers and surface-emitting, light-emitting diodes (LEDs). The LEDs are less temperature dependent than LDs and are more robust, but they typically are slower and less efficient. Those LEDs applicable to fiber optics are described here, along with their operating and transient response characteristics. Edge-emitting LEDs have some niche fiber-optic applications and are briefly described along with the latest advance in LEDs, which is to incorporate a resonant cavity to narrow both the linewidth and increase the output efficiency that can couple into a fiber.

Vertical cavity surface-emitting lasers (VCSELs) have vertical laser cavities that emit light vertically out of the plane of the semiconductor device. Fibers couple more easily to these surface-emitting sources, but VCSEL performance is usually degraded compared to that of the edge-emitting sources. This chapter outlines typical VCSEL designs (material, optical, and electrical); their spatial, spectral, and

polarization characteristics; and their light-out versus current-in characteristics. The most common VCSELs are based on gallium arsenide, operating from 750 to 980 nm, but longer wavelength VCSELs are becoming more practical with GaInNAs quantum wells on GaAs or five-component InAlGaAsP on indium phosphide (InP).

The most common modulators used in fiber-optic systems today are external lithium niobate modulators. These are usually used in Y-branch interferometers, creating intensity modulation as a result of phase modulation by the electro-optic effect. These modulators are introduced in this chapter along with a discussion of high-speed modulation, losses, and polarization dependence, as well as a brief description of optical damage and other modulator geometries. These devices provide chirpfree modulation that can be made very linear for applications such as cable TV.

A recent development in waveguide modulators is the multimode interference (MMI) modulator, based on interference between multiple in-plane spatial modes. This leads to ultra compact modulators, because of much reduced optical path length requirements.

Until recently, lithium niobate was the electro-optic material of choice. Recently polymers have become a viable option for some applications, particularly low-cost. These modulators can be ultra-high-speed and have less insertion loss because of their lower refractive index. Initial challenges of long-term reliability are being met with additional research.

An alternative modulator uses semiconductors, particularly QWs, which allow for more compact devices and monolithic integration. Typically, these are intensity modulators using electroabsorption. By careful design, the chirp in these modulators can be controlled and even used to counteract pulse spreading from chromatic dispersion in fibers. The quantum-confined Stark effect (QCSE) is described, along with the *pin* waveguides used as modulators and techniques for their integration with lasers. Their operating characteristics as intensity modulators, their chirp, and improvements available by using strained QWs are presented.

Some semiconductor modulators use phase change rather than absorption change. The electro-optic effect in III-V semiconductors is discussed, along with the enhanced refractive index change that comes from the QCSE, termed *electrorefraction*. Particularly large refractive index changes occur if available quantum well states are filled by electrons. Phase-change modulators based on this principle can be used in interferometers to yield intensity modulators.

Detectors used in fiber systems are primarily *pin* diodes, although short descriptions of avalanche photodetectors (APDs) and metal-semiconductor-metal (MSM) detectors are provided. The geometry, sensitivity, speed, dark current, and noise characteristics of the most important detectors used in fiber systems are described.

Most of the devices discussed in this chapter are based on semiconductors, and their production relies on the ability to tailor the material to design specifications through *epitaxial growth*. This technology starts with a bulk crystal substrate (usually the binary compounds GaAs or InP) and employs the multilayered growth upon this substrate of a few micrometers of material with a different composition, called a *heterostructure*. *Ternary* layers substitute a certain fraction x for one of the two binary components.

Thus, $\text{Al}_x\text{Ga}_{1-x}\text{As}$ is a common ternary alloy used in laser diodes. Another common ternary is $\text{In}_x\text{Ga}_{1-x}\text{As}$. Layers are *lattice-matched* when the ternary layers have the same size lattice as the binary; otherwise, the epitaxial layer will have *strain*. Lattice-matched epitaxial layers require the substituting atom to have approximately the same size as the atom it replaces. This is true of Al and Ga, so that $\text{Al}_x\text{Ga}_{1-x}\text{As}$ ternary layers are lattice-matched to GaAs. The lowest-cost lasers are those based on GaAs substrates with $\text{Al}_x\text{Ga}_{1-x}\text{As}$ ternary layers surrounding the active layer. These lasers operate at wavelengths near the bandgap of GaAs, about 850 nm, and are typically used in low-cost data communications (as well as in CD players).

The wavelengths required for laser sources in telecommunications applications are those at which the fiber has the lowest loss and/or dispersion, traditionally 1.55 and 1.3 μm . There is no binary semiconductor with a bandgap at these wavelengths, nor is there a lattice-matched ternary. The $\text{In}_x\text{Ga}_{1-x}\text{As}$ ternary will be strained under compression when it is grown on either GaAs or InP, because indium is a much bigger atom than gallium, and arsenic is much bigger than phosphorous. The way to eliminate this strain is to use a fourth small atom to reduce the size of the lattice back to that of the binary. This forms a *quaternary*. The heterostructure most useful for fiber optics applications is based on InP substrates.

The quaternary $\text{In}_x\text{Ga}_{1-x}\text{As}_y\text{P}_{1-y}$ is commonly used, with the compositions x and y chosen to simultaneously provide the desired wavelength and lattice match. These quaternary heterostructures are the basis for much of the long-wavelength technology: sources, modulators, and detectors.

Earlier volumes of this *Handbook* discuss the basics of lasers (Chap. 16, "Lasers," Vol. II), LEDs (Chap. 17, "Light-Emitting Diodes," Vol. II), modulators (Pt. 3, "Modulators," Vol. V), and detectors (Pt. 5, "Detectors," Vol. II). The reader is referred there for general information. This chapter is specific to characteristics that are important for fiber communication systems.

13.2 DOUBLE HETEROSTRUCTURE LASER DIODES

Telecommunication sources are usually edge-emitting lasers, grown with an active laser layer that has a bandgap near either 1.55 or 1.3 μm . The active layer has a quaternary composition consisting of $\text{In}_x\text{Ga}_{1-x}\text{As}_y\text{P}_{1-y}$, grown lattice-matched to InP. The materials growth and fabrication technology had to be developed specifically for telecommunication applications and is now mature. These LDs are more temperature sensitive than GaAs lasers, and this fact has to be incorporated into their use. For telecom applications they are often attached to a thermoelectric cooler and typically provided with a monitoring photodiode in the laser package in order to provide a signal for temperature and/or current control.

Today's telecom systems use single-mode fibers, which require lasers with a single spatial mode. In order to avoid dispersion over long distances, a single frequency mode is necessary. These requirements constrain the geometry of laser diodes used for telecom applications, as discussed in the next section. Following sections discuss the operating characteristics of these LDs and their transient response and noise characteristics, both as isolated diodes and when subject to small reflections from fiber facets. The modulation characteristics of these diodes are discussed, along with frequency chirping. Advanced laser concepts, such as quantum well lasers, strained layer lasers, and lasers with distributed reflection (DFB and DBR lasers) are also introduced.

A typical geometry of an edge-emitting InGaAsP/InP laser is shown in Fig. 1. The active quaternary laser region is shown cross-hatched. It is from this region that light will be emitted. Light travels back and forth between cleaved mirror facets, confined to the active InGaAsP region by the buried heterostructure, and is emitted out of the cross-hatched region, where it diffracts to the far field. The current is confined to the stripe region by the current-blocking *npn* structure on either side. Traditionally, the active regions have uniform composition and are lattice-matched to the substrate. More advanced laser diodes, often used for telecom applications, have active regions containing one or more quantum wells and may be grown to incorporate internal strain in the active region. Both these characteristics are described in Sec. 13.6.

The design of a double heterostructure laser diode requires optimization of the issues discussed in the following subsections.

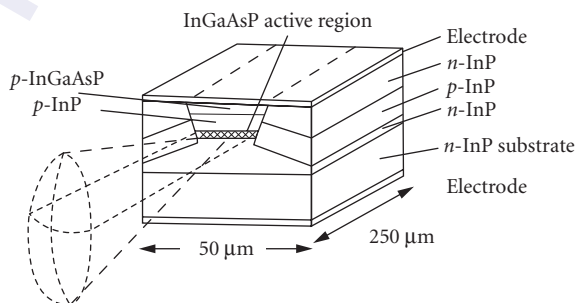


FIGURE 1 Typical geometry for an edge-emitting long-wavelength laser diode, as used in telecommunication systems.

Injection of a Population Inversion into the Active Region

Stimulated emission requires a forward-biased diode, created by sandwiching the active region between p and n layers. Electrons are injected into the active region from the n side and holes are injected from the p side; they become *free carriers*. Efficient electrical injection requires high-quality ohmic contacts attached to the n and p layers; electrical current through the junction then drives the laser.

Confinement of Carriers within the Plane of the Active Layer

Confinement is achieved by growing the n and p layers on either side of the active region with a larger bandgap, as shown in Fig. 2a. In quaternary lasers, wider bandgap material is provided by decreasing x and y relative to their values in the active region. Stimulated emission during electron-hole recombination in the narrow bandgap active layer provides the laser light. The thinner the active layer, the higher its gain. When the active layer thickness is as small as a few 10s of nanometers, the free electron and hole energy levels become quantized in the growth direction, and the active layer becomes a *quantum well*. Quantum wells have higher gain than bulk semiconductor active layers, and thus one or more QWs are often used as the active layers (see Sec. 13.6).

Confinement of Light Near the Active Layer

Stimulated emission gain is proportional to the product of the carrier and photon densities, so that edge-emitting lasers require the highest possible light intensity. This is done by containing the light in an optical waveguide, with a typical near-field light profile, as shown in Fig. 2b. To achieve optical confinement, the layers surrounding the waveguide must have lower refractive indices. It is fortunate that higher-bandgap materials that confine carriers also have smaller refractive index, and so the active layer automatically becomes a waveguide.

Proper optical confinement requires a single waveguide mode. This means that the waveguide layer must be thinner than the cutoff value for higher-order modes. The waveguide thickness d_g must be small enough that

$$d_g k_o \sqrt{n_g^2 - n_c^2} \equiv V < \pi \quad (1)$$

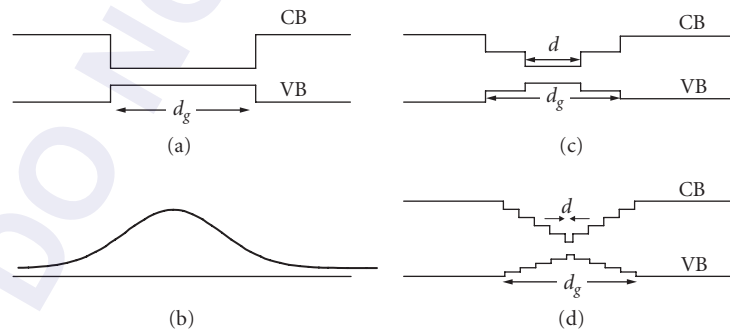


FIGURE 2 Typical diode laser energy band structures and guided mode profile: (a) conduction band (CB) and valence band (VB) of double heterostructure; (b) corresponding near-field spatial profile for light guided in layer of width d_g ; (c) separate confinement heterostructure (SCH); and (d) graded index separate confinement heterostructure (GRINSCH).

where n_g is the refractive index of the waveguide layer (usually the active layer), n_c is the refractive index of the surrounding cladding (usually the p and n layers), and $k_o = 2\pi/\lambda$, where λ is the free-space wavelength of the laser light. Typically, $n_g - n_c \sim 0.2$ and $d_g < 0.56 \mu\text{m}$ for $\lambda = 1.3 \mu\text{m}$. The parameter V is usually introduced to characterize a waveguide.

If the waveguide is too thin, however, the waveguided optical mode spreads out beyond the waveguide layer. The fraction of optical power Γ_g (called the *waveguide confinement factor*) that remains in the waveguide layer of thickness d_g is given approximately by¹

$$\Gamma_g = \frac{V^2}{V^2 + 2} \quad (2)$$

As d_g becomes small, the confinement factor becomes small. When the carriers are confined in very thin layers, such as in quantum wells, the electrical carrier confinement layer cannot serve as an effective optical waveguide because the confinement factor is too small. Then a thicker waveguide region is used to confine the photons, and the carriers are separately confined in a thinner active region, a geometry called a *separate confinement heterostructure* (SCH), as shown in Fig. 2c. In this case the optical confinement factor, defined by the fraction of photons in the active layer of thickness d , is $\Gamma = \Gamma_g(d/d_g)$.

The light can be more effectively focused into a thin active layer by grading the refractive index in the separate confinement region, called a *graded index SCH* (GRINSCH) laser, shown in Fig. 2d. This graded refractive index is produced by growing material with varying bandgaps within the waveguide layer. Grading can be achieved by several discrete layers, as shown, or by grading many ultrathin layers with slight compositional differences. In either case, the focusing property of a GRINSCH structure can be approximated by fitting the graded refractive index to a parabolic refractive index profile $n(x)$ such that $n(x)^2 = n_g^2(1 - x^2/x_0^2)$, where x_0 is related to the curvature of the refractive index near $x = 0$: $x_0 = (n_g/n'')^{1/2}$, where $n'' \equiv \partial^2 n/\partial x^2$ near $x = 0$. The mode guided by this profile has a gaussian beam-intensity profile $I(x) = I_o \exp(-x^2/w^2)$, where $w^2 = x_o^2/k_g$ and $k_g = 2\pi n_g/\lambda$.

Limiting Carrier Injection to Stripe Geometry

Lasers are most efficient when the drive current is limited to the width of the optically active laser area. This requires defining a narrow stripe geometry electrode by means of a window etched in an isolating oxide layer or by ion implantation to render either side of the stripe resistive. More complex laser structures, such as those used in telecommunication applications, often define the conductive stripe electrode by using current-blocking *npn* layers grown on either side of the electrode, as shown in Fig. 1. The *npn* layers, consisting of back-to-back diodes, do not conduct current.

Injected carriers do not usually need lateral confinement, except to achieve the highest possible efficiency. Lateral free-carrier confinement may occur as a by-product of lateral optical confinement, which is discussed next.

Lateral Confinement of Light

The simplest laser diode structures do not specifically confine light laterally, except as the result of the stripe geometry carrier injection. These are called *gain-guided* because high gain in the stripe region, due to the presence of free carriers, introduces a complex refractive index that guides the light laterally. Gain-guided LDs tend to be multimode (both lateral spatial modes and longitudinal frequency modes) unless the stripe is very narrow ($< 10 \mu\text{m}$). In this case, the spatial far-field pattern has “rabbit ears,” a double-lobed far-field pattern that is typically not very useful for coupling into single-mode fibers. Thus, gain-guided LDs are not usually used for telecommunications.

High-quality single-mode LDs for telecom applications typically require a real refractive index difference laterally across the laser. The lowest threshold lasers use *buried heterostructure* (BH) lasers, the geometry shown in Fig. 1. After most of the layers are grown, the sample is taken out of the

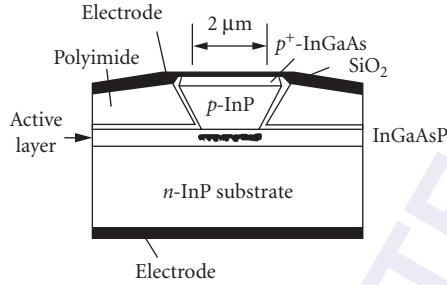


FIGURE 3 Geometry for a ridge waveguide (RWG) laser.

growth chamber and a stripe geometry mesa is etched. Then the sample is returned to the growth chamber, and one or more cladding layers with lower refractive index (higher bandgap) are grown, typically InP, as shown in Fig. 1. When the regrowth is planar, these are called *planar buried heterostructure* (PBH) lasers. The result is a real refractive index guide in the lateral dimension. The width of these *index-guided* laser stripes may be anything from $1\ \mu\text{m}$ to more than $10\ \mu\text{m}$, depending on the refractive index difference between the active stripe and the lateral cladding material. Equation (1), which specifies the condition for single-mode, applies here, with d_g^g as the width of the lateral index guide and n_c defined by the regrown material. A typical lateral width for low-threshold BH lasers is $3\ \mu\text{m}$.

A laser geometry that is much simpler to fabricate and has a higher reliability in production than that of BH lasers is the *ridge waveguide* (RWG) laser, shown in Fig. 3. Light is confined to the region under the p -InP etched mesa by strip loading, which increases the effective refractive index in the waveguide region under the etched mesa. The fabrication starts with the growth of a separate confinement heterostructure (sometimes with the addition of a thin etch-stop layer just after the top waveguide layer), followed by a stripe mesa etch down to the waveguide layer, finishing with planarization and contacting to the stripe. The etch leaves a ridge of p -cladding material above the waveguide layer, which causes *strip loading*, raising the effective refractive index locally in the stripe region, thereby creating lateral confinement of the light. Although the RWG laser is attractive because it requires only a single epitaxial growth, its threshold current is relatively high.

Retroreflection of Guided Light Along the Stripe

Light is usually reflected back and forth inside the laser cavity by Fresnel reflection from cleaved end facets. Since the waveguide refractive index is $n_g \sim 3.5$, the natural Fresnel reflectivity at an air interface, $R = [(n_g - 1)/(n_g + 1)]^2$, is ~ 0.3 . This rather low reflectivity means that LDs are high gain, requiring enough amplification that 70 percent of the light is regenerated on each pass through the active medium.

Relying on Fresnel reflection means that both facets emit light. The light emitted out the back facet may be recovered by depositing a high-reflectivity multilayer coating on the back facet, as is typically done in most telecom lasers. Sometimes a coating is also provided on the front facet to alter its reflectivity, typically to lower it, which increases the output power (as long as the gain is high enough to overcome the large loss upon reflection). The reflectivities must be such that the laser can obey the laser operating condition, which states that in a single round-trip through a laser of length L , the increase in optical power from gain must balance the reduction from finite reflectivity, so that their product is unity. That is,

$$R_1 R_2 \exp(2g_L L) = 1 \quad (3)$$

where R_1 and R_2 are the reflectivities of the two facets and g_L is the *modal gain per unit length* (as experienced by the waveguided laser mode), with a subscript L to represent that the gain is measured with respect to length. If $R_1 = R_2 = 0.3$, then $g_L L = 1.2$. Typical laser diodes have lengths of 400 μm , so $g_L \sim 30 \text{ cm}^{-1}$.

In-plane retroreflection can also be achieved by using distributed feedback created from a grating fabricated near top of the active layer. This method enables the construction of *distributed feedback* (DFB) lasers and *distributed Bragg reflector* (DBR) lasers, which are discussed in Sec. 13.5.

Mounting so that Light Is Edge-Emitted

Because the light is emitted out of the facet laterally, there must be a clear optical path for the light as it exits the laser. In many cases, the light is mounted with the active layer down, very close to its copper (or diamond) heat sink, in order to maximize cooling.² In this case, the laser chip must be placed at the very edge of the heat-sink block, as shown in Fig. 4a.

In some cases, the laser is mounted with its active region up with its substrate next to the heat sink. The edge alignment is not so critical in this case, but of course the laser light will still be emitted in a direction parallel to the plane of the heat sink. Because the thermal conductivity of the heat sink is much higher than that of the substrate, only the lowest threshold lasers, operating at moderate power levels, are operated with the active region up.

Suitable Packaging in Hermetic Enclosure

Water vapor can degrade bare facets of a semiconductor laser when it is operating; therefore, LDs are usually passivated (i.e., their facets are coated with protective layers) and/or they are placed in sealed packages. The LD may be placed in a standard three-pin semiconductor device package, such as a TO-46 can with an optical window replacing the top of the can, as shown in Fig. 4a. The LD will be situated near the package window because the light diverges rapidly after it is emitted from the laser facet. The package window should be antireflection coated because any light reflected back into the laser can have serious consequences on the stability of the output (see Sec. 13.5).

Many high-end applications require an on-chip power monitor and/or a controllable thermoelectric cooler. In this case a more complex package will be used, typically a 14-pin “butterfly” package, often aligned to a fiber pigtail, such as shown in Fig. 4b. In the less expensive datacom applications, nonhermetic packages may be acceptable with proper capping and passivation of the laser surfaces.

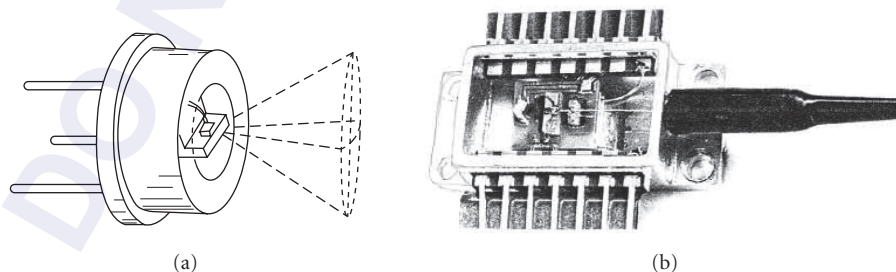


FIGURE 4 Packaging laser diodes: (a) typical hermetically sealed package showing heat sink and emission pattern for a laser diode with its active region placed down on a copper (or diamond) heat sink and (b) typical butterfly package, showing laser in the middle, monitoring photodiode (behind), and fiber alignment chuck in front, all mounted on a thermoelectric cooler. (Photo provided by Spectra-Diode Laboratories.)

Fiber Pigtail Connection

Because light diverges at a rather large angle as it comes out of an edge-emitting laser (as discussed later), it is often desirable to use a laser provided with a *fiber pigtail*, which is a prealigned length of fiber that can be spliced or connected to the telecom fiber in the field. There will be an inevitable reduction in output power (compared to that of a laser with no pigtail) because of finite coupling efficiency into the pigtail, but the output will be immediately useful in a telecom system. The alternative to a fiber pigtail is the use of a microlens, often a graded index (GRIN) lens, discussed in Chap. 18 in this volume.

Long Life

Early lasers showed degradation with running time, but those problems have been solved, and the LDs used in telecom systems should last hundreds of thousands of hours. However, long life requires that care be taken in their use because large reverse-bias static voltages can break down the *pn* diode. Thus, protection from electrostatic shock while handling and from reflected reverse-bias electrical currents during operation should be maintained. In addition, if LDs are driven with too much forward-bias current, the optical output can be so large that the light may damage the facet out of which it is emitted. Since the threshold is strongly temperature dependent, a laser driven at constant current that becomes too cold can emit too much light, with resulting optical damage. Thus, many telecom lasers have monitoring photodiodes to control the laser output and ensure that it stays within acceptable bounds.

13.3 OPERATING CHARACTERISTICS OF LASER DIODES

The principles of semiconductor laser operation are shown in Chap. 19 in Vol. II of this *Handbook*. A forward-biased *pn* junction injects carriers into the active region. As the drive current increases, the carrier density in the active region increases. This reduces the absorption from an initially high value (at thermal equilibrium the absorption coefficient $\alpha \approx 500 \text{ cm}^{-1}$) to zero, at which point the active layer becomes transparent at the prospective laser wavelengths. An active layer is characterized by its carrier density at transparency, N_{tr} . Typically, $N_{tr} \approx 10^{18} \text{ cm}^{-3}$. Above this carrier density, stimulated emission occurs, with a gain proportional to the diode carrier density above transparency. The gain depends on the detailed device design, taking into account the issues enumerated in the preceding section and the materials involved. The gain is sizeable only in direct-band semiconductors (semiconductors based on the III-V or II-VI columns of the periodic table).

Laser Threshold

Threshold is given by the requirement that the round-trip optical gain due to stimulated emission must equal the round-trip optical loss due to the sum of the transmission out the end facets and any residual distributed loss. Gain occurs only for light that is actually in the active region, and not for the fraction of waveguided light that extends outside the active region. Typically, the *local* gain per unit length G_L is defined as that experienced locally by light inside the active region. (The *modal* gain per unit length is $g_L = \Gamma G_L$.) Near transparency, the local gain depends linearly on carrier density N : $G_L = a_L (N - N_{tr})/N_{tr}$, where a_L is the proportionality constant in units of length ($a_L \equiv N \partial G_L / \partial N$ near N_{tr}). Assuming G_L is linear in N , an unpumped region ($N = 0$) has $G_L = -a_L$, which is its loss per unit length. Typically, $a_L \sim 250 \text{ cm}^{-1}$.

The current density (J) is related to the carrier density through $J = eNd/\tau$, where τ is the lifetime of the electron-hole pairs. The transparency current density is 1200 A/cm^2 when $d = 0.15 \text{ }\mu\text{m}$,

$N_{tr} = 10^{18} \text{ cm}^{-3}$, and $\tau = 2 \text{ ns}$. The threshold condition on gain can be found by taking the natural logarithm of Eq. (3) to obtain $g_{L,th} = G_{L,th} \Gamma = \alpha_i + \alpha_m$, where α_m is the mirror reflectivity amortized over length, such that $2\alpha_m L = \ln(1/R_1 R_2)$; and α_i represents any internal losses for the laser mode, also amortized over length. Combining these relations, along with the fact that the current I in a LD with stripe width w and length L is $I = JwL$, gives

$$I_{th} = I_{tr} + \frac{ewN_{tr}}{\tau a_L} \left[\frac{1}{2} \ln \left(\frac{1}{R_1 R_2} \right) + \alpha_i L \right] d \left(1 + \frac{2}{V^2} \right) \quad (4)$$

where the waveguide V parameter is from Eq. (1) with $d_g = d$. Note that the threshold current is independent of device length L when the internal losses are small. The longer the spontaneous lifetime, the lower the threshold current density (although this may lengthen the turn-on time, as discussed in Sec. 13.4). Finally, as expected by the relation between current and current density, a thinner stripe width w will lower the threshold current (consistent with appropriate spatial output, as discussed later). The current density at transparency N_{tr} is a basic property of the gain curve of the active region. It is smaller for QW lasers (Sec. 13.6) than for thicker active regions.

Because V is linearly proportional to d , there is an optimal active layer thickness, a trade-off between increasing the carrier density as much as possible, but not so much as to lose optical confinement. The optimum thickness for $\lambda = 1.3\text{-}\mu\text{m}$ LDs is $\sim 0.15 \mu\text{m}$, and comparable for $\lambda = 1.55\text{-}\mu\text{m}$ LDs (0.15 to 0.18 μm). Threshold currents for broad-area DH lasers can be under 500 A/cm^2 at $\lambda = 1.3 \mu\text{m}$ and $\sim 1000 \text{ A/cm}^2$ at $\lambda = 1.55 \mu\text{m}$. Confining carriers and light separately can beat this requirement, a trick used in designing QW lasers.

Light Out versus Current In (the L-I Curve)

Below laser threshold only spontaneous emission is observed, which is the regime of the LED, as discussed in Sec. 13.8. In the spontaneous regime, the output varies linearly with input current and is emitted in all directions within the active region. As a result, a negligible amount of light is captured by the single-mode fiber of telecom below threshold.

Above threshold, the electrical power is converted to optical power. In general, the light will come out of both facets, so the amount of light reflected out the front facet depends on the rear facet reflectivity. When a 100 percent reflective mirror is deposited on the back facet, the optical power at photon energy $h\nu$ (wavelength $\lambda = c/\nu$) emitted out the front facet is

$$P_{out} = \frac{h\nu}{e} \frac{\alpha_m}{\alpha_m + \alpha_i} (I - I_{th} - I_L) \eta_i \quad (5)$$

where η_i is the *internal quantum efficiency*, which is the fraction of injected carriers that recombine by radiative recombination (usually close to unity in a well-designed semiconductor laser), and I_L is any leakage current. This equation indicates a linear dependence between light out and current above threshold (for constant quantum efficiency). The power out will drop by a factor of 2 if the back facet has a reflectivity equal to that of the front facet, since half the light will leave out the back.

From Eq. (5) can be calculated the *external slope efficiency* of the LD, given by $\partial P_{out} / \partial I$. This allows the *differential quantum efficiency* η_D to be calculated in terms of the power out (both facets):

$$\eta_D \equiv \frac{e}{h\nu} \frac{\partial P_{out}}{\partial I} = \eta_i \frac{\alpha_m}{\alpha_m + \alpha_i} \quad (6)$$

The quantity η_D is also called the *external quantum efficiency*. This efficiency depends on intrinsic loss, possibly due to scattering from roughness in the edges of the waveguide. In long-wavelength lasers, this loss is primarily due to intervalence band absorption.

The *internal* quantum efficiency η_i depends on the way carriers recombine. The rate of carrier loss is the sum of spontaneous processes, expressed in terms of carrier density divided by a lifetime τ_e , and stimulated emission, expressed in terms of local gain per unit time G_T and local photon density P :

$$R(N) = \frac{N}{\tau_e} + G_T(N)P \quad (7)$$

The spontaneous carrier lifetime is given by

$$\frac{1}{\tau_e} = A_{\text{nr}} + BN + CN^2 \quad (8)$$

The term BN is due to spontaneous radiative recombination; its dependence on N results from needing the simultaneous presence of both an electron and hole (which have the same charge densities because of charge neutrality in undoped active regions). The nonradiative recombination terms that decrease the quantum efficiency below unity are a constant term A_{nr} (accounting for all background nonradiative recombination) and an Auger recombination term (with coefficient C) that depends on the square of the carrier density (Auger processes involve several carriers simultaneously). This term is particularly important in long-wavelength lasers where the Auger coefficient C is large.

Stimulated emission is accounted for by gain in the time domain G_T , which depends on N (approximately linearly near threshold). The group velocity v_g converts gain per unit length G_L into a rate \dot{G}_T (gain per unit time), $G_T \equiv v_g G_L$. We can define a_T in the time domain by $a_T = v_g a_L$ so that $\dot{G}_T = a_T (N - N_{\text{tr}})/N_{\text{tr}}$.

The *internal* quantum efficiency in a laser is the fraction of the recombination processes that emit light:

$$\eta_i = \frac{BN^2 + G_T(N)P}{A_{\text{nr}}N + BN^2 + CN^3 + G_T(N)P} \quad (9)$$

At high powers in well-designed lasers, the internal quantum efficiency approaches 1.

Figure 5 shows a typical experimental result³ for the light out of a LD as a function of applied current (the so-called *L-I* curve). It can be seen that the linear relation between light out and current

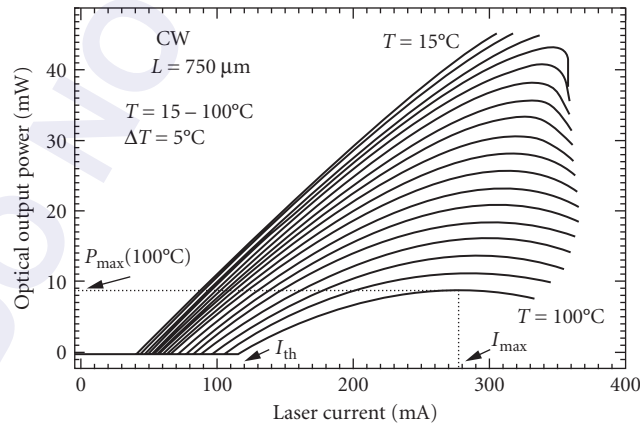


FIGURE 5 Typical experimental result for light out versus current in (the *L-I* curve). These results are for diodes operating at 1.3 μm , consisting of strained layer multiple quantum well InGaAsP lasers measured at a series of elevated temperatures.³

saturates as the current becomes large enough, particularly at high temperatures. The decrease in external slope efficiency with increasing current has contributions from the increase in leakage current with injection current, from junction heating that reduces recombination lifetime and increases threshold current, and because the internal absorption increases with injection current.

When there is more than one mode (longitudinal or transverse) in the LD, the L - I curve has *kinks* at certain current levels. These are slight abrupt reductions in light out as the current increases. After a kink the external slope efficiency may be different, along with different spatial and spectral features of the laser. These multimode lasers may be acceptable for low-cost communication systems, but high-quality systems require singlemode lasers that do not exhibit such kinks in their L - I curves.

Temperature Dependence of Laser Properties

Long-wavelength lasers are typically more sensitive to temperature than are GaAs lasers. This sensitivity is usually expressed as an experimentally measured exponential dependence of threshold on temperature T through $I_{th}(T) = I_o \exp(T/T_o)$, where T_o is a characteristic temperature (in kelvin) that expresses the measured thermal sensitivity. This formula is valid only over a limited temperature range, because it has no real physical derivation, but it has proved convenient and is often quoted. The data in Fig. 5 correspond to $T_o \approx 80$ K.

The mechanisms for LD sensitivity to temperature depend on the material system. In InP-based long-wavelength DH lasers, T_o is dominated by Auger recombination. However, in short-wavelength GaAs lasers and in strained layer QWs, Auger recombination is suppressed, T_o is higher and is attributed to intervalence band absorption and/or carrier leakage over the heterostructure barrier, depending on the geometry. Typical long-wavelength DH lasers have T_o in the range of 50 to 70 K. Typical strained layer QW lasers have T_o in the range of 70 to 90 K, although higher T_o can be achieved by incorporating aluminum in the barriers, with as high as 143 K reported.⁴ This temperature dependence limits the maximum optical power that can be obtained because of the phenomenon of *thermal runaway*, as shown at the highest temperatures in Fig. 5. While increasing the current usually increases the power, the junction temperature also increases (due to ohmic losses), so the threshold may increase and the output power may tend to decrease.

Various means for increasing T_o have been explored; the most effective solution is an active layer of tensile strained QWs (discussed in Sec. 13.6). This has increased T_o from approximately 50 K to as high as 140 K, comparable to that measured in GaAs. In double heterostructures, losses by carrier leakage can be reduced with a dual active region for double carrier confinement, which was demonstrated to achieve T_o values as high as 180 K in 1.3- μ m InP lasers.⁵

The temperature dependence of long-wavelength lasers may limit their performance at high temperatures, which in turn limits where they can be used in the field. In practice, many long-wavelength lasers require thermoelectric coolers, particularly at higher power levels.

Spatial Characteristics of Emitted Light

Light is emitted out the facet of the laser diode after it has been waveguided in both directions. It will diverge by diffraction, more strongly in the out-of-plane dimension, where it has been more strongly guided. The diffracting output is sketched in Fig. 1. The spatial characteristics of the output can be estimated by fitting the guided light to a gaussian profile and then calculating the far-field pattern. The *out-of-plane* near-field profile for the lowest-order mode in an optical confinement layer of width d_g has been fit to a gaussian distribution $\exp(-x^2/w^2)$ with⁶

$$w = d_g \left(0.321 + \frac{2.1}{V^{3/2}} + \frac{4}{V^6} \right) \quad \text{for } 1.8 < V < 6 \quad (10)$$

where V is from Eq. (1). The far-field diffraction angle θ_{ff} has a slightly different gaussian fit, with a half-angle given by $\tan \theta_{ff} = \lambda/\pi w_o$, with⁶

$$w_o = d_g \left(0.31 + \frac{3.15}{V^{3/2}} + \frac{2}{V^6} \right) \quad \text{for } 1.5 < V < 6 \quad (11)$$

Experimental data can be compared to the gaussian beam formulation by remembering that the full-width half-maximum power FWHM = $w(2 \ln 2)^{1/2}$. For a typical strongly index-guided buried heterostructure laser, the far-field FWHM angle out-of-plane is approximately 1 rad and in-plane is approximately 1/2 rad. These angles are independent of current for index-guided lasers. Separate confinement heterostructure (SCH) lasers can have smaller out-of-plane beam divergences, more typically approximately 30°.

Single-mode lasers that are index-guided in the lateral direction (buried heterostructure and ridge waveguide) will obey the preceding equations, with lateral divergence angles varying from 30° to 10°, depending on design. This beam width will also be independent of current. When LDs are gain-guided laterally, the spatial variation of the gain leads to a complex refractive index and a curved wavefront. The result is that the equivalent gaussian lateral beam appears to have been emitted from somewhere inside the laser facet. The out-of-plane beam, however, is still index-guided and will appear to be emitted from the end facet. This means that the output of a gain-guided laser has *astigmatism*, which must be compensated for by a suitably designed external lens if the laser is to be focused effectively into a fiber (as discussed elsewhere in this *Handbook*).

If the laser emits a diverging gaussian beam with waist w , a lens can be used to focus it into a fiber. An effective thin lens of focal length f placed a distance d_1 after the laser facet will focus to a new waist w' given by

$$w'^2 = w^2 \frac{f^2}{b^2 + X_1^2} \quad (12)$$

where $b \equiv \pi w^2/\lambda$ and $X_1 \equiv d_1 - f$.

The new waist will be located a distance d_2 from the lens, where

$$X_2 = X_1 \frac{f^2}{X_1^2 + b^2} \quad (13)$$

and $X_2 \equiv d_2 - f$.

This new waist must be matched to the spatial mode of the fiber. Because of the large numerical aperture of LD emission, simple lenses may exhibit severe spherical aberration and typical coupling efficiencies may be only a few percent. Fiber systems usually utilize pigtailed fiber, butt-coupled as close as possible to the laser, without any intervening lens. Alternatively, a ball lens may be melted directly onto a fiber tip and placed near the laser facet. Sometimes *graded index* (GRIN) lenses are used to improve coupling into fibers.

Gain-guided LDs with electrode stripe widths of more than 5 μm usually emit multiple in-plane spatial modes that interfere laterally, producing a spatial output with multiple maxima and nulls. Such spatial profiles are suitable for multimode fiber applications, but cannot be efficiently coupled into single-mode fibers. They will diffract at an angle given by setting w equal to the minimum near-field feature size. If the stripe is narrow enough, gain-guided LDs are always single-mode, but these devices have a double-lobed far-field spatial profile (from the complex refractive index in the gain medium) that cannot be conveniently coupled into single-mode fibers.

Spectral Characteristics of Laser Light

In principle, a Fabry-Perot laser has many frequency modes with frequencies ν_m , given by requiring standing waves within the laser cavity. Since the m th mode obeys $m\lambda/2n = L$, where n is the refractive index experienced by the guided laser mode, then $\nu_m = mc/2nL$. Taking the differential, the frequency

difference between modes is $\Delta\nu = c/2n_{\text{eff}}L$, where the *effective group refractive index* $n_{\text{eff}} = n + \nu(\partial n/\partial \nu)$. For typical semiconductor lasers, $n = 3.5$ and $n_{\text{eff}} = 4$, so that when $L = 250 \mu\text{m}$, the frequency difference between modes is $\Delta\nu = 150 \text{ GHz}$, and since $\Delta\lambda = (\lambda^2/c)\Delta\nu$, when $\lambda = 1.5 \mu\text{m}$, the wavelength spacing is $\Delta\lambda \approx 1 \text{ nm}$.

At any given instant in time, a single spatial mode emits only one spectral mode. However, in multimode lasers, considerable *mode hopping* occurs, in which the LD jumps from one spectral mode to another very rapidly. Most spectral measurements are time averages and do not resolve this mode hopping, which can occur in nanoseconds or less. Explanations for the mode hopping typically involve *spatial hole burning* or *spectral hole burning*. Hole burning occurs when the available carrier density is momentarily depleted, either spatially or spectrally. At that time an adjacent mode with a different (longitudinal or lateral) spatial profile providing a different resonance wavelength may be more advantageous for laser action. Thus, the laser jumps to this new mode. The competition for available gain between different modes often induces semiconductor lasers to emit light at several frequencies.

One way to provide a single spectral mode is to ensure a single (lateral) spatial mode. It has been found that single spatial mode lasers usually have single spectral modes, at least at moderate power levels. However, this may change upon modulation. The only way to *ensure* a single-frequency LD is to ensure a single longitudinal mode by using distributed feedback, as discussed in Sec. 13.7.

Polarization

The light emitted from a typical LD is usually linearly polarized in the plane of the heterostructure. While gain in the semiconductor medium has no favored polarization dependence, the *transverse electric (TE) waveguide mode* (polarized in-plane) is favored for two reasons. First, the TE mode is slightly more confined than the *transverse magnetic (TM) mode* (polarized out-of-plane). Second, the Fresnel reflectivity off the cleaved end facets is strongly polarization sensitive. As waveguided light travels along the active stripe region, it can be considered to follow a zig-zag path, being totally internally reflected by the cladding layers. The total internal reflection angle for these waves is about 10° off the normal to the cleaved facets of the laser. This is enough to cause the TM waveguide mode to experience less reflectivity, while the TE-polarized mode experiences higher reflectivity and has a lower threshold. Thus, laser light from LDs is traditionally polarized in the plane of the junction.

The introduction of strain (Sec. 13.6) in the active layer changes the polarization properties, and the particular polarization will depend on the details of device geometry. In addition, DFB and DBR lasers (Sec. 13.7) do not have strong polarization preferences, and they must be carefully designed and fabricated if well-defined single polarization is required.

13.4 TRANSIENT RESPONSE OF LASER DIODES

When laser diodes are operated by direct current, the output is constant and follows the L - I curve discussed previously. When the LD is rapidly switched, however, there are transient phenomena that must be taken into account. Such considerations are important for any high-speed communication system, especially digital systems. The study of these phenomena comes from solving the semiconductor rate equations.⁷

Turn-On Delay

When a semiconductor laser is turned on abruptly by applying forward-biased current to the diode, it takes time for the carrier density to reach its threshold value and for the photon density to build up, as shown in the experimental data of Fig. 6.⁸ This means that a laser has an unavoidable turn-on time. The delay time depends on applied current and on carrier lifetime, which depends on carrier density

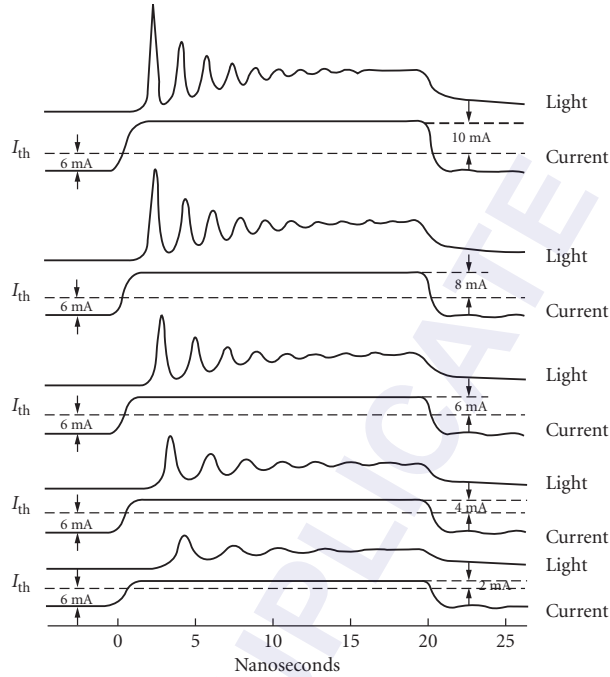


FIGURE 6 Experimental example of turn-on delay and relaxation oscillations in a laser diode when the operating current is suddenly switched from 6 mA below the threshold current of 177 mA to varying levels above threshold (from 2 to 10 mA). The GaAs laser diode was 50 μm long, with a SiO_2 -defined stripe 20 μm wide. Light output and current pulse are shown for each case.⁸

N , as shown in Eq. (8). Using a differential analysis, the turn-on time for a laser that is switched from an initial current I_i just below threshold to I just above threshold is

$$\tau_d = \tau'(N_{\text{th}}) \frac{I_{\text{th}} - I_i}{I - I_{\text{th}}} \quad (14)$$

where $\tau'(N)$ is a differential carrier recombination given by $1/\tau' = A_{\text{nr}} + 2BN + 3CN^2$. When the LD is initially off ($I_i = 0$) and $I \gg I_{\text{th}}$, the turn-on delay has an inverse current dependence: $\tau_d = \tau_e(N_{\text{th}}) I_{\text{th}}/I$.

When radiative recombination dominates, then $1/\tau_e \approx BN$ and $1/\tau' \approx 2BN \approx 2/\tau_e$, as seen by comparing the middle terms of Eqs. (11) and (22). For a 1.3- μm laser, $A_{\text{nr}} = 10^8/\text{s}$, $B = 10^{-10} \text{ cm}^3/\text{s}$, $C = 3 \times 10^{-29} \text{ cm}^6/\text{s}$, and $N_{\text{th}} \approx N_{\text{tr}} = 10^{18} \text{ cm}^{-3}$. Thus, $\tau_e = 5 \text{ ns}$ and a typical turn-on time at 1.5 times threshold current is 3 ns. The increase in delay time as the current approaches threshold is clearly seen in the data of Fig. 6. To switch a laser rapidly, it is necessary to switch it from just below threshold to far above threshold. However, Fig. 6 shows that under these conditions there are large transient oscillations, discussed in the following section.

Relaxation Oscillations

An important characteristic of the output of any rapidly switched laser (not just LDs) is the *relaxation oscillations* that can be observed in Fig. 6. These overshoots and undershoots occur as the photon and

carrier dynamics are coming into equilibrium. Such oscillations are characteristic of the nonlinear coupled laser rate equations and can be found by simple perturbation theory. These relaxation oscillations have a radian frequency Ω_R given, to first order, by⁹

$$\Omega_R^2 = \frac{1 + \chi}{\tau_e \tau_p} \frac{I - I_{th}}{I_{th}} \quad (15)$$

where I is the current, I_{th} is the current at threshold, τ_p is the photon lifetime in the cavity, given by $v_g \tau_p = (\alpha_i + \alpha_m)^{-1}$, and $\chi = \Gamma a_L v_g \tau_p = \Gamma a_T \tau_p = \Gamma [a_L / (\alpha_i + \alpha_m)]$, where a_L was previously defined and is approximately the unpumped absorption loss in the active medium. The factor χ is then the ratio of the unpumped absorption loss to the cavity loss, typically 1 to 3 for semiconductor lasers. It can also be shown that $\chi = I_{tr} / (I_{th} - I_{tr})$, where I_{tr} is the current at transparency.

At 50 percent above threshold, and when $\chi \approx 1$, the time between successive relaxation oscillation maxima is approximately the geometric mean of the carrier and photon lifetimes: $\Omega_R^{-2} \approx 1 / \tau_e \tau_p$. Typical LD numbers are $\tau_e = 10$ ns and $\tau_p = 3$ ps, so at 1.5 times threshold current, the relaxation oscillation frequency is $f_R = \Omega_R / 2\pi = 1$ GHz, and the time between relaxation oscillation peaks is 1 ns.

The decay rate of these relaxation oscillations γ_R is given by

$$2\tau_e \gamma_R = 1 + \chi \frac{I - I_{th}}{I_{tr}} = 1 + (1 + \chi) \frac{I - I_{th}}{I_{th}} = 1 + \Omega_R^2 \tau_e \tau_p \quad (16)$$

and is roughly 2 ns at twice threshold for typical heterostructure lasers. At 1.5 times threshold and when $\chi \approx 1$, Eq. (16) gives $\gamma_R \approx 1 / \tau_e$. The relaxation oscillations last approximately as long as the spontaneous emission lifetime of the carriers.

This first-order analysis has several assumptions that do not seriously affect the relaxation oscillation frequency, but will overestimate the time that relaxation oscillations are present. The analysis ignores *gain saturation*, which reduces gain with increased photon density P and is important at high optical powers. It also ignores the rate of spontaneous emission in the cavity, R_{sp} , which is important at small optical powers. Finally, it ignores the impact of changing carrier density on spontaneous emission lifetime. A more typical experimental decay rate for lasers at 1.3 μm wavelength is $\gamma \approx 3 / \tau_e$. A more exact theoretical formulation can be found in Ref. 10.

The number of relaxation oscillations (before they die out) in an LD at 1.5 times threshold is roughly $\Omega_R / \gamma_R \propto (\tau_e / \tau_p)^{1/2}$. Shorter carrier lifetimes mean fewer relaxation oscillations, because carriers reach steady state more rapidly. Shorter carrier lifetimes also mean faster turn-on times. Carrier lifetimes are shortened by high carrier densities, an important regime for high-speed semiconductor lasers; high carrier densities can be achieved by using as small an active region as possible (such as QWs) and by reducing the reflectivity of the laser facets to raise the threshold carrier density.

The relaxation oscillations disappear if the current is just at threshold. However, we've also seen that under this situation the turn-on time becomes very long. It is more advantageous to turn the laser on fast, suffering the relaxation oscillations and using LDs designed to achieve a high decay rate, which means using LDs with the highest possible relaxation oscillation frequency.

The relaxation oscillation can also be expressed in terms of the photon density P inside the laser cavity, or in terms of the power out, P_{out} :

$$\Omega_R^2 = g_T g'_T P = \frac{g'_T P}{\tau_m} = g'_T \left(\frac{P_{out}}{h\nu V_a} \frac{\alpha_m + \alpha_i}{\alpha_m} \right) \quad (17)$$

where g_T is the modal gain per unit time and $g'_T \equiv \partial g_T / \partial N = \Gamma a_T / N_{tr}$. This formulation makes use of the relationship between output power and internal photon density in a cavity of volume V_a :

$$P_{out} = h\nu \left(\frac{P}{\tau_m} \right) V_a \quad (18)$$

where $\tau_m = (v_g \alpha_m)^{-1}$ is the time it takes light to bleed out the mirror. Note that the relaxation oscillation frequency increases as the photon density increases, confirming that smaller laser dimensions are better for high-speed modulation.

Relaxation oscillations can be avoided by biasing the laser just below threshold (letting the laser “simmer”). Fiber optic systems are designed to reduce the impact of these relaxation oscillators, whether they are digital telecommunication systems, or high-speed analog CATV systems. Many communication applications avoid these issues by operating the LD at steady state with DC current applied and use external modulators.

Modulation Response and Gain Saturation

The *modulation response* describes the amplitude of the modulated optical output as a function of frequency under small-signal current modulation. Laser diodes have a resonance in the modulation response at the relaxation oscillation frequency, as indicated by the experimental data in Fig. 7.¹¹ It is more difficult to modulate the laser above the relaxation oscillation frequency. Carrying out a small-signal expansion of the rate equations around photon density P , the modulation response (in terms of current density J) is¹²

$$\frac{\partial P}{\partial J} = \frac{(1/ed)(g'_T P + \beta_{sp}/\tau_e)}{(g'_T P/\tau_p + \beta_{sp}/\tau_e \tau_p - \omega^2) + j\omega(g'_T P + 1/\tau_e)} \quad (19)$$

where β_{sp} is the fraction of spontaneous emission that radiates into the mode. This modulation response has the form of a second-order low-pass filter. Resonance occurs when $\omega^2 \approx g'_T P/\tau_p = \Omega_R^2$ (from Eq. (17), with negligible internal loss); that is, at the relaxation oscillation frequency.

The modulation response at a frequency well below the relaxation oscillation frequency can be expressed as $\partial P_{out}/\partial I$ using Eq. (19) in the limit when $\omega \rightarrow 0$. From $\partial P/\partial J = \tau_p/ed$ and using Eq. (18), the low frequency modulation response is

$$\frac{\partial P_{out}}{\partial I} = \frac{hv}{\tau_m} V^a \frac{\partial P}{\partial J} \cdot \frac{1}{wL} = \frac{hv}{e} \cdot \frac{\tau_p}{\tau_m} = \frac{hv}{e} \cdot \frac{\alpha_m}{\alpha_m + \alpha_i} \quad (20)$$

which is expected from Eq. (5) for $\eta_i \rightarrow 1$.

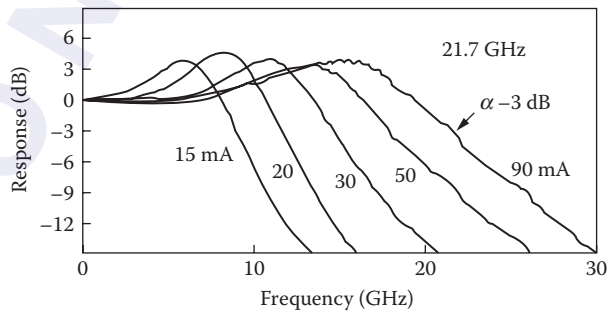


FIGURE 7 Measured small-signal modulation response of a high-speed DFB laser at several bias levels. Zero-dB modulation response is defined in terms of the low-frequency modulation response, Eq. (32).¹¹

The 3-dB modulation radian frequency bandwidth ω_B can be expressed in terms of the relaxation oscillation parameters by¹³

$$\omega_B^2 = \Omega_R^2 - \gamma_R^2 + 2\sqrt{\Omega_R^2(\Omega_R^2 + \gamma_R^2) + \gamma_R^4} \quad (21)$$

The parameters are strongly power dependent and ω_B increases with optical power. When $\gamma_R \ll \Omega_R$, which means the time carriers remain in the active region is longer than the time photons remain in the cavity, the 3-dB bandwidth $\omega_B \approx \sqrt{3\Omega_R} \propto \sqrt{P}$. At high optical powers the presence of gain saturation (reduced gain at high optical power densities) must be included; the modulation bandwidth saturates, and the limiting value depends on the way that the gain saturates with photon density.

Using a heuristic expression for modal gain saturation of the form $g_T(N, P) = g_T'(N - N_o)/(1 + P/P_s)^{1/2}$, the limiting value of the modulation bandwidth at high optical powers is $\omega_B^2 = 3g_T'P_s/2\tau_p$, where P_s is the saturation photon density. A typical modulation bandwidth for an InGaAsP laser operating at 1.55- μm wavelength is $\omega_B = 20$ to 40 GHz.

Frequency Chirping

When the carrier density in the active region is rapidly changed, the refractive index n also changes rapidly, causing a frequency shift proportional to $\partial n/\partial t$. This broadens the laser linewidth from its original width of approximately 100 MHz into a double-peaked profile with a gigahertz linewidth, as shown in the experimental results of Fig. 8.¹⁴ The frequency spread is directly proportional to the dependence of the refractive index n on carrier density N . This is a complex function that depends on wavelength and degree of excitation, but for simplicity a Taylor expansion around the steady-state carrier density N_o can be assumed: $n = n_o + n_1(N - N_o)$, where $n_1 \equiv \partial n/\partial N$. The (normalized) ratio of this slope to that of the modal gain per unit length g_L is called the cavity *linewidth enhancement factor* β_c :

$$\beta_c \equiv -2k_o \frac{\partial n/\partial N}{\partial g_L/\partial N} = -2k_o \frac{n_1}{g_L'} \quad (22)$$

Sometimes the linewidth enhancement factor is called the α -factor. Typical values lie between 2 and 6.

The magnitude of the frequency spread between the double lobes of a chirped pulse, $2\delta\omega_{\text{CHF}}$, can be estimated in the small-signal and large-signal regimes from analyzing the time dependence of a modulated pulse in terms of the sum of all frequency components, shown next.¹⁵

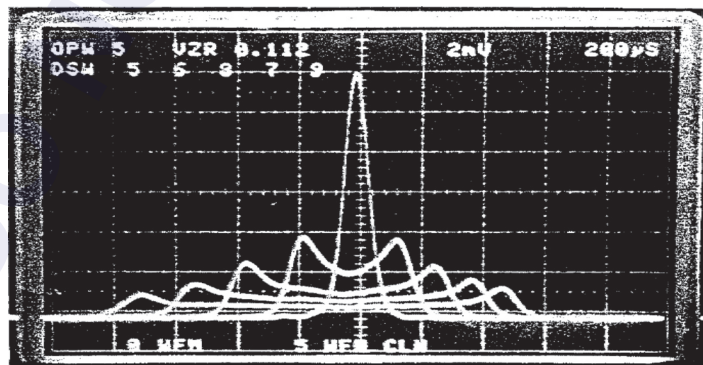


FIGURE 8 Time-averaged power spectra of 1.3- μm InGaAsP laser under sinusoidal modulation at 100 MHz. Horizontal scale is 0.05 nm per division. Spectrum broadens with increase in modulation current due to frequency chirping.¹⁴

Small Signal Modulation For a modulation frequency ω_m that is less than the relaxation oscillation frequency, and assuming that the carrier lifetime is longer than the photon lifetime, $\gamma_R \ll \Omega_R$, a small modulation current I_m will cause a frequency chirp of magnitude

$$\delta\omega_{\text{CH}} = \frac{\beta_c I_m h_\nu}{2eP_{\text{out}}} \left(\frac{\alpha_m}{\alpha_m + \alpha_i} \right) \sqrt{\omega_m^2 + \gamma_p^2} \quad (23)$$

where $\gamma_p = N/P(\beta_{\text{sp}}/\tau_c) - (\partial g_T/\partial P)P$ (remembering that $\partial g_T/\partial P$ is negative). The origin of chirp is the linewidth enhancement factor β_c . It will be largest for gain-guided devices where β_c is the maximum. The chirp will be smaller in lasers with $\alpha_m \ll \alpha_p$, such as will occur for long lasers, where mirror loss is amortized over a longer length, but such lasers will have a smaller differential quantum efficiency and smaller relaxation oscillation frequency. Typical chirp bandwidths (using $\delta\lambda_{\text{CH}} = \lambda^2 \delta\omega_{\text{CH}}/2\pi c$) at 25-mA modulation current can vary from $\delta\lambda_{\text{CH}} = 0.2$ nm for gain-guided lasers to $\delta\lambda_{\text{CH}} = 0.03$ nm for ridge waveguide lasers. At a wavelength of 1 micrometer in frequency units this would be 60 GHz for gain-guided lasers and 9 GHz for ridge waveguide lasers.

Large Signal Modulation There is a transient frequency shift during large-signal modulation given by

$$\delta\omega_{\text{CH}} = \frac{\beta_c}{2} \left(\frac{1}{P} \frac{\partial P}{\partial t} \right) \quad (24)$$

When a gaussian shape pulse is assumed as $\exp(-t^2/T^2)$, the chirp becomes $\delta\omega_{\text{CH}} \approx \beta_c/T$.

The importance of the linewidth enhancement factor β_c is evident from this section; its existence will inevitably broaden modulated laser linewidths. Because it is a basic characteristic of the laser medium itself, there are no design freedoms to reduce chirp under large signal modulation.

13.5 NOISE CHARACTERISTICS OF LASER DIODES

Noise in LDs results from fluctuations in spontaneous emission and from the carrier generation-recombination process (shot noise). To analyze the response of LDs to noise, one starts with rate equations, introduces Langevin noise sources as small perturbations, and linearizes (performs a small-signal analysis). Finally, one solves in the frequency domain using Fourier analysis.^{16,17} Only the results are given here.

Relative Intensity Noise

Noise at a given frequency is described in terms of *relative intensity noise* (RIN), defined by

$$\text{RIN} = \frac{\Delta P^2}{P_T^2} \quad (25)$$

where ΔP^2 is the photon noise spectral density (noise per unit frequency interval), and P_T is the total photon number, $P_T = PV_a$. The solution to the analysis for RIN in an LD is

$$\text{RIN} = \frac{2\beta_{\text{sp}} I_{\text{th}}}{ePV_a} \cdot \frac{1/\tau'^2 + \omega^2 + (\partial g_T/\partial N)^2 P/(\beta_{\text{sp}} V_a)}{[(\Omega_R - \omega)^2 + \gamma_R^2][(\Omega_R + \omega)^2 + \gamma_R^2]} \quad (26)$$

where β_{sp} is the fraction of spontaneous emission emitted into the laser cavity. As before, the photon density P can be related to the optical power out both facets by $P_{\text{out}} = (h\nu) PV_a/\tau_m$. Note the significant enhancement of noise near the relaxation oscillation frequency $\omega = \Omega_R$, where the noise has

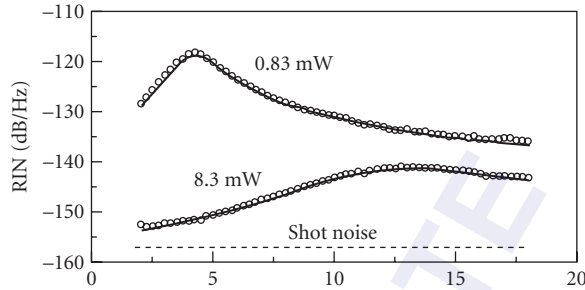


FIGURE 9 Measured relative intensity noise as a function of frequency in a multiple quantum well 1.5- μm laser diode, for optical power near threshold and high above threshold. The shot noise level for the higher power measurement is also shown.¹⁸

its maximum value. An example of RIN as a function of frequency ω is shown in Fig. 9,¹⁸ for both low power and high power, showing that the RIN goes up as the total optical power decreases.

At low frequencies, and for $\gamma_R \ll \Omega_R$, the noise is proportional to the inverse fourth power of the relaxation oscillation frequency. Clearly, it is advantageous to use as high a relaxation oscillation frequency as possible to reduce RIN. Since the relaxation oscillation frequency is proportional to the square root of the power P , the RIN increases as $1/P^3$ as the power decreases. Inserting the expression for Ω_R into Eq. (26) gives

$$\text{RIN}_{\text{lr}} = \frac{2\beta_{\text{sp}} I_{\text{th}} \tau_m^2 V_a^2}{eP_T^3 (\partial g_T / \partial N)^2} \left[\frac{1}{\tau'^2} + \left(\frac{\partial g_T}{\partial N} \right)^2 \frac{P}{\beta_{\text{sp}} V_a} \right] \quad (27)$$

where τ' is the differential carrier recombination time and P_T is the total photon number. Usually, the first term dominates. It can be seen that the volume of the active laser region V_a should be as small as possible, consistent with maintaining a significant power out.

Signal-to-Noise Ratio

The *signal-to-noise ratio* (SNR) can be found in terms of the relaxation oscillation parameters using the simplified expression for RIN (which assumes $\tau_e \Omega_R \gg \gamma_R \Omega_R \gg 1$) and the total photon number:

$$\text{SNR}^2 = \frac{2\gamma_R e}{\beta_{\text{sp}} I_{\text{th}}} P_T = \frac{2\gamma_R e \tau_m}{\beta_{\text{sp}} I_{\text{tr}} h\nu} P_{\text{out}} \quad (28)$$

As expected, the SNR ratio increases linearly with output power and increases as the amount of spontaneous emission decreases. This expression can be simplified further by recalling that $\gamma_R \approx \tau_e/3$.

Gain saturation at high optical powers eventually limits the SNR to about 30 dB; while at powers of a few milliwatts it is 20 dB, with intensity fluctuations typically close to 1 percent.

Mode Partition Noise in Multimode Laser Diodes

The preceding discussion of noise holds qualitatively for multimode lasers as long as all the laser modes are included. However, measurements made on any one mode show much more noise, particularly at low frequencies. This is due to the mode-hopping discussed previously, and is referred

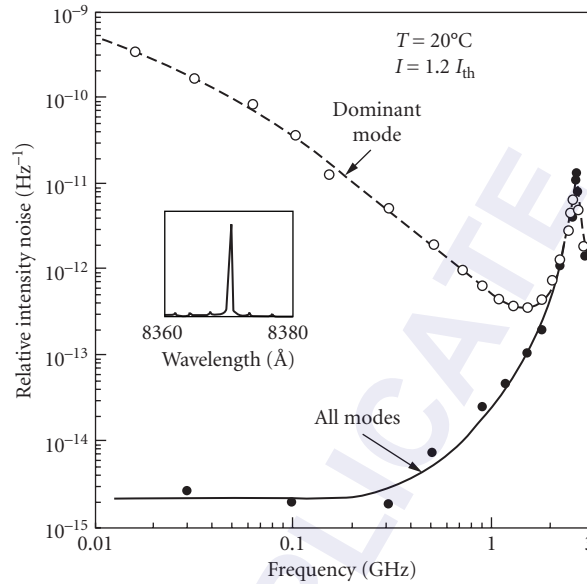


FIGURE 10 Effect of mode partition noise on relative intensity noise in multimode lasers. Experimentally observed intensity-noise spectra in all modes (solid curve) or in dominant mode (dashed curve). Inset shows spectrum of average mode power.¹⁹

to as *mode partition noise*. That is, the power partitions itself between different laser modes in a way that keeps the overall intensity relatively constant, as shown by the solid line in Fig. 10. The power in each mode is not a steady function of time, because the power distribution among the modes changes with time. Whenever the distribution changes, the power output undergoes fluctuation, leading to a noise term on a nominally stable output. This leads to the enhanced RIN on the dominant mode in Fig. 10.¹⁹ Even an output whose spectrum looks nominally single-mode, as shown in the inset of Fig. 10, can have a large RIN on the dominant mode. This is because the spectrum is time averaged. A sidemode does not contain 5 percent of the power all the time; for example, it contains 100 percent of the power for 5 percent of the time. This causes the very large RIN observed. The solution to avoiding this noise is to insist on a single longitudinal mode by using distributed feedback. Since lasers for telecommunication applications are typically single mode, we will not consider mode partition noise further. It becomes important for data communications based on multimode lasers, however.

Phase Noise—Linewidth

The fundamental linewidth of a laser is given by the stochastic process of spontaneous emission, as first derived by Schawlow and Townes in the very early days of lasers. In a laser diode, additional noise enters from the stochastic process of carrier injection. Because the refractive index is a function of the carrier density, changes in carrier density cause changes in refractive index, which in turn create phase noise.

The formula for the radian frequency linewidth of a LD includes the *linewidth enhancement factor* β_c :

$$\delta\omega = (1 + \beta_c^2)\delta\omega_0 \quad (29)$$

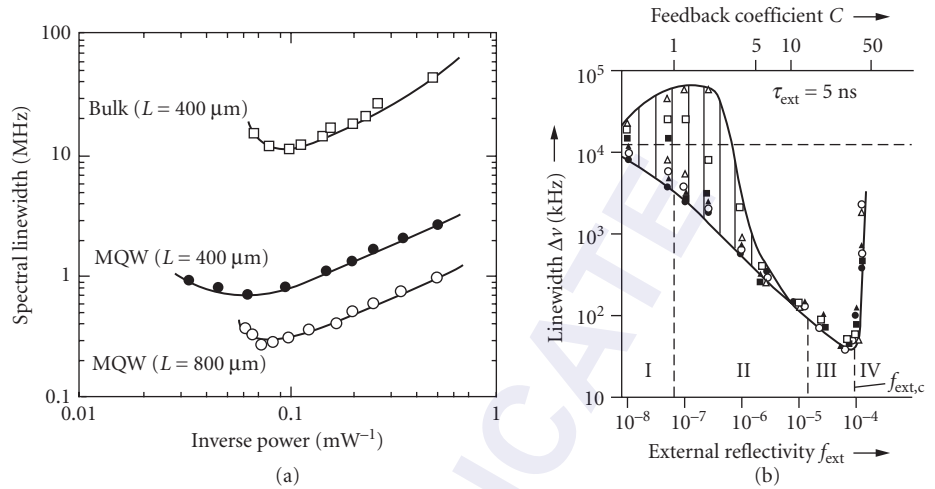


FIGURE 11 Linewidth of a laser diode: (a) DFB lasers as a function of inverse power, comparing bulk active regions and multiple quantum well active regions²⁰ and (b) Fabry-Perot laser as a function of feedback parameter, with corresponding feedback regimes.²¹

where $\delta\omega_o$ the original Schawlow-Townes linewidth is given by

$$\delta\omega_o = \frac{\beta_{sp} I_{th} h\nu \tau_m}{2eP_{out}} \quad (30)$$

Typical value of the linewidth enhancement factor is $\beta_c = 5$. The linewidth decreases inversely as the laser power increases. However, as shown in the experimental data in Fig. 11a,²⁰ at high enough power (above 10 mW) the linewidth narrowing saturates at approximately 1 to 10 MHz and then begins to broaden again at even higher power levels. It is also possible to reduce the linewidth by using QWs and increasing the cavity length (to decrease N_{th} and increase P).

External Optical Feedback and Coherence Collapse

Laser diodes are extremely sensitive to weak time-delayed feedback, such as from reflections off the front ends of fiber pigtails. These fed-back signals can result in mode hopping, strong excess noise, and chaotic behavior in the *coherence collapse* regime. Some of the features of feedback-induced noise are outlined here.

Various regimes of feedback can be characterized by a *feedback parameter* C :

$$C = \sqrt{f_{ext}} C_e \frac{\tau_{ext}}{\tau_L} \sqrt{1 + \beta_c^2} \quad (31)$$

where f_{ext} is the fraction of the emitted power that is externally reflected back into the laser. The external coupling factor $C_e = (1 - R)/\sqrt{R}$, where R is the reflectivity of the laser facet and β_c is the linewidth enhancement factor. The external round-trip time delay is τ_{ext} and the laser round-trip time is τ_L . Figure 11b²¹ gives an example of the linewidth of a semiconductor laser versus the external coupling factor C in which various feedback regimes are indicated.

Regimes of Feedback The following provides a useful classification scheme:²²

Regime I ($C < 1$): At the lowest levels of feedback, narrowing or broadening of the emission line is observed, depending on the phase of the feedback.

Regime II ($C > 1$): At higher levels of feedback, mode hopping between different external cavity modes may appear.

Regime III ($C \gg 1$): Further increasing the levels of feedback, the laser is observed to operate in a mode with the smallest linewidth.

Regime IV (coherence collapse): At yet higher feedback levels, satellite modes appear, separated from the main mode by the relaxation oscillation frequency. These grow as the feedback increases and the laser line eventually broadens, due to the collapse of the coherence of the laser. This regime does not depend on distance from the laser to the reflector.

Regime V (external cavity laser): This regime of stable operation can be reached only with an antireflection-coated laser output facet to ensure a two-mirror cavity with the largest possible coupling back into the laser, and is not of concern here.

A quantitative discussion of these regimes follows.²¹ Assume that the coupling efficiency from the laser into the fiber is η . Because feedback requires a double pass, the fraction of emitted light fed back into the laser is $f_{\text{ext}} = \eta^2 R_c$, where R_c is the reflectivity from the end of the fiber. The external reflection changes the overall cavity reflectivity and therefore the threshold gain, depending on its phase relative to the phase inside the cavity. Possible modes are defined by the threshold gain and the requirement that an effective round-trip phase = $m\pi$. But a change in the threshold gain also changes the refractive index and the phase through the linewidth enhancement factor β_c .

Regime I For very weak feedback, there is only one solution when the laser phase is set equal to $m\pi$, so that the laser frequency ω is at most slightly changed and its linewidth $\Delta\omega_o$ will be narrowed or broadened, as the external reflection adds to or subtracts from the output of the laser. The linewidth lies between extremes:

$$\frac{\Delta\omega_o}{(1+C)^2} \ll \Delta\omega \ll \frac{\Delta\omega_o}{(1-C)^2} \quad (32)$$

The system performance moves toward regime II as $C \rightarrow 1$. Note that at $C = 1$ the maximum value predicts an infinite linewidth. This indicates that even very small feedback can cause wide spectral response, as long as $C \sim 1$.

Regime II For higher feedback with $C > 1$, several solutions with the round-trip laser phase = $m\pi$ may exist. Linewidth broadening occurs because the single external cavity mode now has split into a dual mode, accompanied by considerable phase noise. Mode hopping gives linewidth broadening and intensity noise. This is a low-frequency noise with a cutoff frequency of about 10 MHz.

Regime III With increasing feedback, the mode splitting increases up to a frequency $\Delta\omega = 1/\tau_{\text{ext}}$ and the cutoff frequency for mode hopping noise decreases until only one of the split modes survives. To understand which mode survives, it is important to realize that in regime III, stable and unstable modes alternate with increasing phase. Because $\beta_c \neq 0$, the mode with the best phase stability (corresponding to the minimum linewidth mode) does not coincide with the mode with minimum threshold gain. The minimum linewidth mode predominates, due to its phase stability. This mode remains relatively stable and has the same emission frequency as the laser without feedback. This mode predominates as long as the inverse of the linewidth of the solitary laser is larger than the external cavity round-trip time. In this regime, the laser is stably phase locked by the feedback.

Regime IV The stable linewidth of regime III collapses as the fraction of power fed back f_{ext} increases to a critical value. There is considerable discussion of the physical mechanism that leads to

this coherence collapse. The existence of this regime has been demonstrated by simulation, through numerical solution of the rate equations. Fitting theoretical analyses to experimental results indicates that the onset of coherence collapse occurs when the feedback factor is larger than a critical value given by²³

$$C \geq C_{\text{crit}} = 2\gamma_R \tau_{\text{ext}} \frac{1 + \beta_c^2}{\beta_c^2} \quad (33)$$

where γ_R is the damping rate of the relaxation oscillations. As the feedback level approaches the critical value C , undamped relaxation oscillations appear and oscillations of carrier density through the linewidth enhancement factor β_c induce the phase of the field to oscillate. This analytical result assumes that $\tau_{\text{ext}} \gg 1/\gamma_R$. An approximate solution for γ_R gives $C_{\text{crit}} = \tau_{\text{ext}}/\tau_e \times (1 + \beta_c^2)/\beta_c^2$.

Cavity Length Dependence and RIN In some regimes the stable regions depend on the length of the external cavity, that is, the distance from the extra reflection to the laser diode. These regions have been mapped out for two different laser diodes, as shown in Fig. 12.²⁴ The qualitative dependence on the distance of the laser to the reflection should be similar for all LDs.

The RIN is low for weak to moderate levels of feedback but increases tremendously in regime IV. The RIN and the linewidth are strongly related (see Fig. 11); the RIN is suppressed in regimes III and V.

Low-Frequency Fluctuations When a laser operating near threshold is subject to a moderate amount of feedback, chaotic behavior evolves into *low-frequency fluctuations* (LFF). During LFF the average laser intensity shows sudden dropouts, from which it gradually recovers, only to drop out again after some variable time, typically on the order of tens of external cavity round-trips. This occurs in regimes of parameter space where at least one stable external cavity mode exists, typically at the transition between regimes IV and V. Explanations differ as to the cause of LFF, but they appear to originate in strong-intensity pulses that occur during the build-up of average intensity, as a form of mode locking, being frustrated by the drive toward maximum gain. Typical frequencies for LFF are 20 to 100 MHz, although feedback from reflectors very close to the laser has caused LFF at frequencies as high as 1.6 GHz.

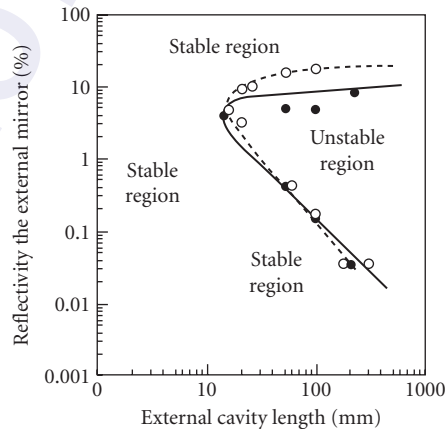


FIGURE 12 Regimes of stable and unstable operation for two laser diodes (○ and ●) when subject to external feedback at varying distances and of varying amounts.²⁴

Conclusions A laser diode subject to optical feedback exhibits a rich and complex dynamic behavior that can enhance or degrade the laser's performance significantly. Feedback can occur through unwanted back reflections—for instance, from a fiber facet—and can lead to a severe degradation of the spectral and temporal characteristics, such as in the coherence collapse or in the LFF regime. In both regimes, the laser intensity fluctuates erratically and the optical spectrum broadens, showing large sidebands. Because these unstable regimes can occur for even minute levels of feedback, optical isolators or some other means of reflection prevention are often used in systems applications.

13.6 QUANTUM WELL AND STRAINED LASERS

Introducing quantum wells and strain into the active region of diode lasers has been shown to provide higher gain, greater efficiency, and lower threshold. Essentially all high-quality lasers for optical communications use one or both of these means to improve performance over bulk heterostructure lasers.

Quantum Well Lasers

We have seen that the optimum design for low-threshold LDs uses the thinnest possible active region to confine free carriers, as long as the laser light is waveguided. When the active layer has a thickness less than a few tens of nanometers (hundreds of angstroms), it becomes a *quantum well*. That is, the layer is so thin that the confined carriers have energies that are quantized in the growth direction z , as described in Chap. 19 in Vol. II of this *Handbook*. This changes the density of states and the gain (and absorption) spectrum. While bulk semiconductors have an absorption spectrum near the band edge that increases with photon energy above the bandgap energy E_g as $(h\nu - E_g)^{1/2}$, quantum wells have an absorption spectrum that is steplike in photon energy at each of the allowed quantum states. Riding on this steplike absorption is a series of exciton resonances at the absorption steps that occur because of the Coulomb interaction between free electrons and holes, which can be seen in the spectra of Fig. 13.²⁵ These abrupt absorption features result in much higher gain for QW lasers than for bulk semiconductor lasers. The multiple spectra in Fig. 13 record the reduction in absorption as the QW states are filled with carriers. When the absorption goes to zero, transparency is reached. Figure 13 also shows that narrower wells push the bandgap to higher energies, a result of quantum confinement. The QW thickness is another design parameter in optimizing lasers for telecommunications.

Because a single quantum well (SQW) is so thin, its optical confinement factor is small. It is necessary to use either multiple QWs (separated by heterostructure barriers that contain the electron wave functions within individual wells) or to use a guided wave structure that focuses the light into a SQW. The latter is usually a GRIN structure, as shown in Fig. 2*d*. Band diagrams as a function of distance in the growth direction for typical QW separate confinement heterostructures are shown in Fig. 14. The challenge is to properly confine carriers and light using materials that can be reliably grown and processed by common crystal growth methods.

Quantum wells have provided significant improvement over bulk active regions, as originally observed in GaAs lasers. In InP lasers, Auger recombination and other losses come into play at the high carrier densities that occur in quantum-confined structures. However, it has been found that strain in the active region can improve the performance of quaternary QW lasers to a level comparable with GaAs lasers. Strained QW lasers are described in the following section.

The LD characteristics described in Secs. 13.2 to 13.5 hold for QW lasers as well as for bulk lasers. While the local gain is larger, the optical confinement factor will be much smaller. Equation (1) shows that the parameter V becomes very small when d_g is small, and Eq. (2) shows Γ_g is likewise small. With multiple quantum wells (MQWs), d_g can be the thickness of the entire region containing the MQWs and their barriers, but Γ must now be multiplied by the filling factor Γ_f of the QWs within the MQW region. If there are N_w wells, each of thickness d_w , then $\Gamma_f = N_w d_w / d_g$. With a GRINSCH structure, the optical confinement factor depends on the curvature of its refractive gradient near the center of the guide.

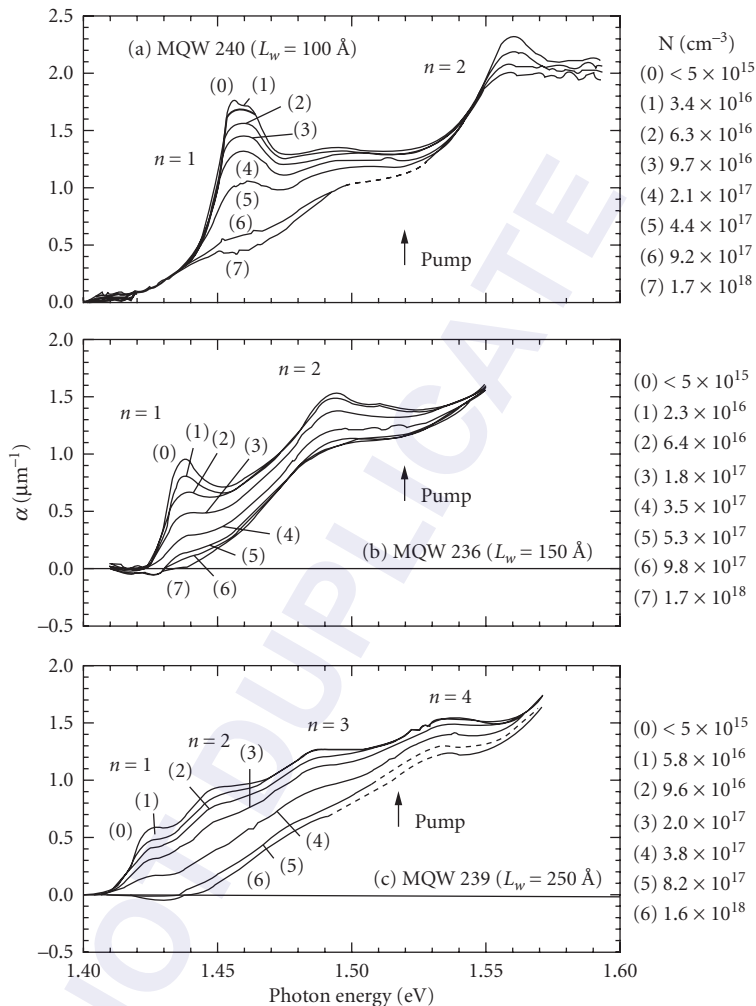


FIGURE 13 Absorption spectrum for multiple quantum wells of three different well sizes, for varying levels of optically induced carrier density, showing the decrease in absorption toward transparency. Note the stronger excitonic resonances and increased bandgap with smaller well size.²⁵

Different geometries have subtle differences in performance, depending on how many QWs are used and the extent to which a GRINSCH structure is dominant. The lowest threshold current densities have been reported for the highest Q cavities (longest lengths or highest reflectivities) using SQWs. However, for lower Q cavities the lowest threshold current densities are achieved with MQWs, even though they require higher carrier densities to achieve threshold. This is presumably because Auger recombination depends on the cube of the carrier density, so that SQW lasers will have excess losses, due to their higher carrier densities. In general, MQWs are a better choice in long-wavelength lasers, while SQWs have the advantage in GaAs lasers. However, with MQW lasers it is important to realize that the transport of carriers moving from one well to the next during high-speed modulation must be taken into account. In addition, improved characteristics of strained layer QWs make SQW devices more attractive.

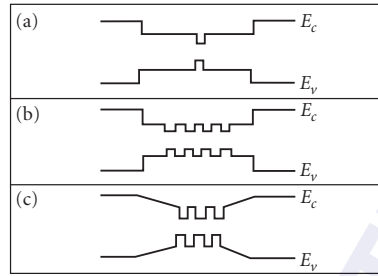


FIGURE 14 Typical band diagrams (energy of conduction band E_c and valence band E_v versus growth direction) for quantum wells in separate confinement laser heterostructures: (a) single quantum well; (b) multiple quantum wells; and (c) graded index separate confinement heterostructure (GRIN SCH) and multiple quantum wells.

Strained Layer Quantum Well Lasers

Active layers containing *strained quantum wells* have proven to be an extremely valuable advance in high-performance long-wavelength InP lasers. They have lower thresholds, enhanced differential quantum efficiency η_D , larger characteristic temperature T_0 , reduced linewidth enhancement factor β_c (less chirp), and enhanced high-speed characteristics (larger relaxation oscillation frequency Ω_R), compared to unstrained QW and bulk devices. This results from the effect of strain on the energy-versus-momentum band diagram. Bulk semiconductors have two valence bands that are degenerate at the potential well minimum (at momentum $k_x = 0$), as shown in Fig. 15a. They are called

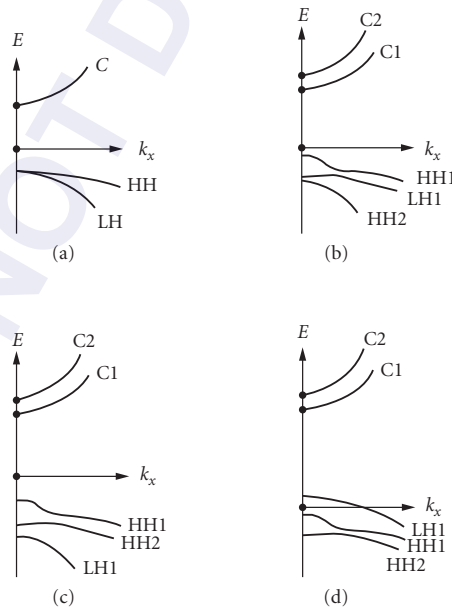


FIGURE 15 The effect of strain on the band diagram (energy E versus in-plane momentum k_x) of III-V semiconductors: (a) no strain; (b) quantum wells; (c) compressive strain; and (d) tensile strain.

heavy-hole (HH) and light-hole (LH) bands, since the smaller curvature means a heavier effective mass. Quantum wells lift this degeneracy, and interaction between the two bands near momentum $k = 0$ causes a local distortion in the formerly parabolic bands, also shown in Fig. 15b. There are now separately quantized conduction bands (C_1 and C_2) and a removal of the valence band degeneracy, with the lowest energy heavy holes HH_1 no longer having the same energy as the lowest energy light holes LH_1 at $k = 0$. The heavy hole effective mass becomes smaller, more nearly approaching that of the conduction band. This allows population inversion to become more efficient, increasing the differential gain; this is one factor in the reduced threshold of QW lasers.²⁶

Strain additionally alters this structure in a way that can improve performance even more. Compressive strain in the QW moves the heavy-hole and light-hole valence bands further apart and further reduces the hole effective mass (Fig. 15c). Strain also decreases the heavy-hole effective mass by a factor of 2 or more, further increasing the differential gain and reducing the threshold carrier density. Higher differential gain also results in a smaller linewidth enhancement factor. Tensile strain moves the heavy-hole and light-hole valence bands closer together (Fig. 15d). In fact, at one particular tensile strain value these bands become degenerate at $k = 0$. Further tensile strain results in the light hole having the lowest energy at $k = 0$. These lasers will be polarized TM, because of the angular momentum properties of the light-hole band. This polarization has a larger optical matrix element, which can enhance the gain within some wavelength regions.

In addition to the heavy- and light-hole bands, there is an additional, higher-energy valence band (called the *split-off band*, not shown in Fig. 15) which participates in Auger recombination and intervalence band absorption, both of which reduce quantum efficiency. In unstrained material there is a near-resonance between the bandgap energy and the difference in energy between the heavy-hole and split-off valence bands, which enhances these mechanisms for nonradiative recombination. Strain removes this near-resonance and reduces those losses that are caused by Auger recombination and intervalence band absorption. This means that incorporating strain is essential in long-wavelength laser diodes intended to be operated at high carrier densities. The reliability of strained layer QW lasers is excellent, when properly designed. However, strain does increase the intraband relaxation time, making the gain compression factor worse, so strained lasers tend to be more difficult to modulate at high-speed.

Specific performance parameters depend strongly on the specific material, amount of strain, size and number of QWs, and device geometry, as well as the quality of crystal growth. Calculations show that compressive strain provides the lowest transparency current density, but tensile strain provides the largest gain (at sufficiently high carrier densities), as shown in Fig. 16.²⁷ The lowest threshold lasers, then, will typically be compressively strained. Nonetheless, calculations show that, far enough above the band edge, the differential gain is 4 times higher in tensile strain compared to compressive strain. This results in a smaller linewidth enhancement factor, even if the refractive index changes per carrier density are larger. It has also been found that tensile strain in the active region reduces

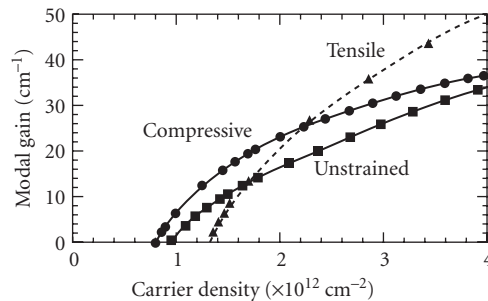


FIGURE 16 Modal gain at 1.55 μm in InGaAs QW lasers calculated as a function of the carrier density per unit area contained in the quantum well. Well widths were determined by specifying wavelength.²⁷

the Auger recombination, decreasing the losses introduced at higher temperatures. This means that T_o can increase with strain, particularly tensile strain. Strained QWs enable performance at 1.55 μm comparable with that of GaAs lasers. Deciding between compressively and tensilely strained QWs will be a matter of desired performance for specific applications.

Threshold current densities under 200 A/cm² have been reported at 1.55 μm ; T_o values on the order of 140 K have been reported—3 times better than bulk lasers. Strained QW lasers have improved modulation properties compared with bulk DH lasers. Because the gain coefficient can be almost double, the relaxation oscillation frequency is expected to be almost 50 percent higher, enhancing the modulation bandwidth and decreasing the relative intensity noise for the same output power. Even the frequency chirp under modulation will be less, because the linewidth enhancement factor is less. The typical laser geometry, operating characteristics, transient response, noise, frequency chirping, and the effects of external optical feedback are all similar in the strained QW lasers to what has been described previously for bulk lasers. Only the experimentally derived numerical parameters will be somewhat different; strained long-wavelength InP-based semiconductor lasers have performance parameters comparable to those of GaAs lasers. One difference is that the polarization of the light emitted from strained lasers may differ from that emitted from bulk lasers. As explained in Sec. 13.3, the gain in bulk semiconductors is independent of polarization, but lasers tend to be polarized in-plane because of higher facet reflectivity for that polarization. The use of QWs causes the gain for the TE polarization to be slightly (~10 percent) higher than for the TM polarization, so lattice-matched QW lasers operate with in-plane polarization. Compressive strain causes the TE polarization to have significantly more gain than the TM polarization (typically 50 to 100 percent more), so these lasers are also polarized in-plane. However, tensile strain severely depresses the TE gain, and these lasers have the potential to operate in TM polarization.

Typical 1.3- and 1.5- μm InP lasers today use from 5 to 15 wells that are grown with internal strain. By providing strain-compensating compressive barriers, there is no net buildup of strain. Typical threshold current densities today are ~1000 A/cm², threshold currents ~10 mA, T_o ~ 50 to 70 K, maximum powers ~40 mW, differential efficiencies ~0.3 W/A, and maximum operating temperatures ~70°C before the maximum power drops by 50 percent. There are trade-offs on all these parameters; some can be made better at the expense of some of the others.

13.7 DISTRIBUTED FEEDBACK AND DISTRIBUTED BRAGG REFLECTOR LASERS

Rather than cleaved facets, some lasers use distributed reflection from corrugated waveguide surfaces. Each groove provides some slight reflectivity, which adds up coherently along the waveguide at the wavelength given by the corrugation. This has two advantages. First, it defines the wavelength (by choice of grating spacing) and can be used to fabricate single-mode lasers. Second, it is an in-plane technology (no cleaves) and is therefore compatible with monolithic integration with modulators and/or other devices.

Distributed Bragg Reflector Lasers

The *distributed Bragg reflector* (DBR) laser replaces one or both laser facet reflectors with a waveguide diffraction grating located outside the active region, as shown in Fig. 17.²⁸ The reflectivity of a Bragg mirror is the square of the reflection coefficient (given here for the assumption of lossless mirrors):²⁹

$$r = \frac{\kappa}{\delta - iS \coth(SL)} \quad (34)$$

where κ is the *coupling coefficient* due to the corrugation (which is real for corrugations that modify the effective refractive index in the waveguide, but would be imaginary for periodic modulations in the gain

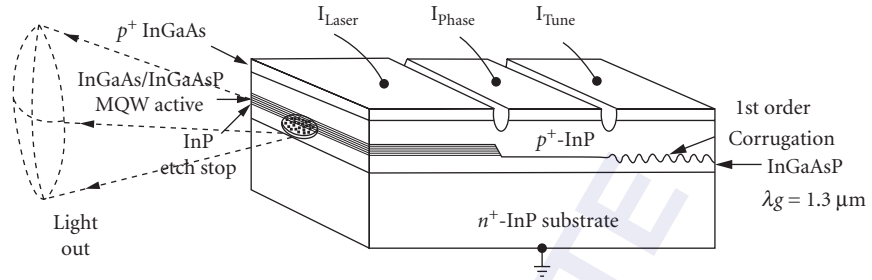


FIGURE 17 Schematic for DBR laser configuration in a geometry that includes a phase portion for phase tuning and a tunable DBR grating. Fixed-wavelength DBR lasers do not require this tuning region. Designed for 1.55- μm output, light is waveguided in the transparent layer below the MQW that has a bandgap at a wavelength of 1.3 μm . The guided wave reflects from the rear grating, sees gain in the MQW active region, and is partially emitted and partially reflected from the cleaved front facet. Fully planar integration is possible if the front cleave is replaced by another DBR grating.²⁸

and could, indeed, be complex). Also, δ is a *detuning parameter* that measures the offset of the optical wavelength λ from the grating periodicity Λ . When the grating is used in the m th order, the detuning δ relates to the optical wavelength λ by

$$\delta = \frac{2\pi n_g}{\lambda} - \frac{m\pi}{\Lambda} \quad (35)$$

where n_g is the effective group refractive index of the waveguide mode, and m is any integer. Also, S is given by $S^2 = \kappa^2 - \delta^2$. When detuning $\delta > \kappa$, Eq. (34) is still valid, and is analytically evaluated with S as imaginary.

The Bragg mirror has its maximum reflectivity on resonance when $\delta \rightarrow 0$ and the wavelength on resonance λ_m is determined by the m th order of the grating spacing through $\Lambda = m\lambda_m/2n_g$. The reflection coefficient on resonance is $r_{\max} = -i \tanh(KL)$, and the Bragg reflectivity on resonance is

$$R_{\max} = \tanh^2(KL) \quad (36)$$

where K is the coupling per unit length, $K = |\kappa|$, and is larger for deeper corrugations or when the refractive index difference between the waveguide and the cladding is larger. The reflectivity falls off as the wavelength moves away from resonance and the detuning increases. Typical resonant reflectivities are $R_{\max} = 0.93$ for $KL = 2$ and $R_{\max} = 0.9987$ for $KL = 4$.

A convenient formula for the shape of the reflectivity as a function of detuning near resonance is given by the reflection loss:

$$1 - R = \frac{1 - (\delta L)^2 / (\kappa L)^2}{\cosh^2(SL) - (\delta L)^2 / (\kappa L)^2} \quad (37)$$

The reflection loss doubles when off resonance by an amount $\delta L = 1.6$ for $\kappa L = 2$ and when $\delta L = 2$ for $\kappa L = 4$. The wavelength half-bandwidth is related to the detuning δL by $\Delta\lambda = \lambda_o^2 (\delta L / 2\pi) / n_g L$. The calculated FWHM of the resonance is 0.6 nm (when $L = 500 \mu\text{m}$, $\lambda = 1.3 \mu\text{m}$) for a 99.9 percent reflective mirror. The wavelength of this narrow resonance is fixable by choosing the grating spacing and can be modulated by varying the refractive index (with, for example, carrier injection), properties that make the DBR laser very favorable for use in optical communication systems.

The characteristics of Fabry-Perot lasers described previously still hold for DBR lasers, except that the narrow resonance can ensure that these lasers are single-mode, even at high excitation levels.

Distributed Feedback Lasers

When the corrugation is put directly on the active region or its cladding, this is called *distributed feedback* (DFB). One typical BH example is shown in Fig. 18, with a buried grating waveguide that was grown on top of a grating-etched substrate, which forms the separate confinement heterostructure laser. The cross-hatched region contains the MQW active layer. A stripe mesa was etched and regrown to bury the heterostructure. Reflection from the cleaved facets must be suppressed by means of an antireflection coating. As before, the grating spacing is chosen such that, for a desired wavelength near λ_o , $\Lambda = m\lambda_o/2n_{ge}$, where n_{ge} is the effective group refractive index of the laser mode inside its waveguiding active region, and m is any integer. A laser operating under the action of this grating has feedback that is distributed throughout the laser gain medium. In this case, Eq. (34) is generalized to allow for the gain: $\delta = \delta_o + ig_L$, where g_L is the laser gain and $\delta_o = 2\pi n_{ge}/\lambda - 2\pi n_{ge}/\lambda_o$. Equations (34) to (37) remain valid, understanding that now δ is complex.

The laser oscillation condition requires that after a round-trip inside the laser cavity, a wave must have the same phase that it started out with, so that successive reflections add in phase. Thus, the phase of the product of the complex reflection coefficients (which now include gain) must be an integral number of 2π . This forces r^2 to be a positive real number. So, laser oscillation requires $r^2 > 0$.

On resonance $\delta_o = 0$ and $S_o^2 = \kappa^2 + g_L^2$, so that S_o is pure real for simple corrugations (κ real). Since the denominator in Eq. (34) is now pure imaginary, r^2 is negative and the round-trip laser oscillation condition cannot be met. Thus, there is no on-resonance solution to a simple DFB laser with a corrugated waveguide and/or a periodic refractive index. The DFB laser oscillates slightly off-resonance.

DFB Threshold We look for an off-resonance solution to the DFB. A laser requires sufficient gain that the reflection coefficient becomes infinite. That is, $\delta_{th} = iS_{th} \coth(S_{th}L)$, where $S_{th}^2 = \kappa^2 - \delta_{th}^2$. By simple algebraic manipulation, the expression for δ_{th} can be inverted. For large gain, $\delta_{th}^2 \gg K^2$, so that $S_{th} = i\delta_{th} = i\delta_o - g_L$, and the inverted equation becomes³⁰

$$\exp(2S_{th}) \frac{4(g_L - i\delta_o)^2}{K^2} = -1 \quad (38)$$

This is a complex eigenvalue equation that has both a real and an imaginary part, which give both the detuning δ_o and the required gain g_L .

The required laser gain is found from the magnitude of Eq. (38) through

$$2 \tan^{-1} \left(\frac{\delta_o}{g_L} \right) - 2\delta_o L + \delta_o L \frac{K^2}{g_L^2 + \delta_o^2} = (2m + 1)\pi \quad (39)$$

There is a series of solutions, depending on the value of m .

For the largest possible gains, $\delta_o L = -(m + 1/2)\pi$. There are two solutions, $m = -1$ and $m = 0$, giving $\delta_o L = -\pi/2$ and $\delta_o L = +\pi/2$. These are two modes equally spaced around the Bragg resonance.

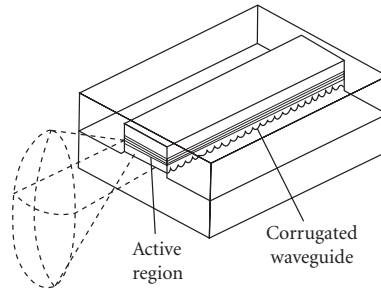


FIGURE 18 Geometry for a buried grating heterostructure DFB laser.

Converting to wavelength units, the mode detuning becomes $\delta_o L = -2\pi n_g L (\delta\lambda/\lambda^2)$, where $\delta\lambda$ is the deviation from the Bragg wavelength. Considering $\delta_o L = \pi/2$, for $L = 500 \mu\text{m}$, $n_g = 3.5$, and $\lambda = 1.55 \mu\text{m}$, this corresponds to $\delta\lambda = 0.34 \text{ nm}$. The mode spacing is twice this, or 0.7 nm .

The required laser gain is found from the magnitude of Eq. (38) through

$$\frac{K^2}{4} = (g_L^2 L^2 + \delta_o^2 L^2) \exp(-2g_L L) \quad (40)$$

For the required detuning $\delta_o L = -\pi/2$, the gain can be found by plotting Eq. (40) as a function of gain g_L , which gives $K(g_L)$, which can be inverted to give $g_L(K)$.

These results show that there is a symmetry around $\delta_o = 0$, so that there will tend to be *two* modes, equally spaced around λ_o . Such a multimode laser is not useful for communication systems so something must be done about this. The first reality is that there are usually cleaved facets, at least at the output end of the DFB laser. This changes the analysis from that given here, requiring additional Fresnel reflection to be added to the analysis. The additional reflection will usually favor one mode over the other, and the DFB will end up as a single-mode. However, there is very little control over the exact positioning of these additional cleaved facets with respect to the grating, and this has not proven to be a reliable way to achieve single-mode operation. The most common solution to this multimode problem is to use a *quarter-wavelength-shifted grating*, as shown in Fig. 19. Midway along the grating, the phase is made to change by $\pi/2$ and the two-mode degeneracy is lifted. This is the way DFB lasers are made today.

Quarter-Wavelength-Shifted Grating Introducing an additional phase shift of π to the round-trip optical wave enables an on-resonance DFB laser. This is done by interjecting an additional phase region of length $\Lambda/2$, or $\lambda/4n_g$, as shown in Fig. 19. This provides an additional $\pi/2$ phase each way, so that the high-gain oscillation condition becomes $\delta_o L = -m\pi$. Now there is a unique solution at $m = 0$, given by Eq. (40) with $\delta_o = 0$:

$$KL = g_L L \exp(-g_L L) \quad (41)$$

Given a value for the DFB coupling parameter KL , the gain can be calculated. Alternatively, the gain can be varied, and the coupling coefficient that must be used with that gain can be calculated. It can be seen that if there are internal losses α_p , the laser must have sufficient gain to overcome them as well: $g_L \rightarrow g_L + \alpha_p$.

Quarter-wavelength-shifted DFB lasers are commonly used in telecommunication applications. DFB corrugations can be placed in a variety of ways with respect to the active layer. Most common is to place the corrugations laterally on either side of the active region, where the evanescent wave of the guided mode experiences sufficient distributed feedback for threshold to be achieved. Alternative methods place the corrugations on a thin cladding above the active layer. Because the process of corrugation may introduce defects, it is traditional to avoid corrugating the active layer directly. Once a

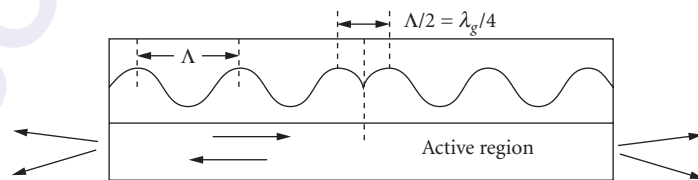


FIGURE 19 Side view of a quarter-wavelength-shifted grating, etched into a separate confinement waveguide above the active laser region. Light with wavelength in the medium λ_g sees a $\pi/4$ phase shift, resulting in a single-mode DFB laser operating on line-center.

DFB laser has been properly designed, it will be single-mode at essentially all power levels and under all modulation conditions. Then the single-mode laser characteristics described in the early part of this chapter will be well satisfied. However, it is crucial to avoid reflections from fibers back into the laser, because instabilities may arise, and the output may cease to be single-mode.

A different technique that is sometimes used is to spatially modulate the gain. This renders κ complex and enables an on-resonance solution for the DFB laser, since S will then be complex on-resonance. Corrugation directly on the active region makes this possible, but care must be taken to avoid introducing centers for nonradiative recombination.

More than 35 years of research and development have gone into semiconductor lasers for telecommunications. Today it appears that the optimal sources for these applications are strained QW distributed feedback lasers operating at 1.3 or 1.55 μm wavelength.

13.8 TUNABLE LASERS

The motivation to use tunable lasers in optical communication systems comes from wavelength division multiplexing (WDM), in which a number of independent signals are transmitted simultaneously, each at a different wavelength. The first WDM systems used wavelengths far apart (so-called *coarse* WDM) and settled on a standard of 20-nm wavelength spacing (~ 2500 GHz). But interest grew rapidly toward *dense* wavelength division multiplication (DWDM), with much closer wavelength spacing. The International Telecommunications Union (ITU) defined a standard for a grid of optical frequencies, each referring to a reference frequency which has been fixed at 193.10 THz (1552.5 nm). The grid separation can be as narrow as 12.5 GHz or as wide as 100 GHz. Tuning range can extend across the conventional erbium amplifier window C band (1530 to 1565 nm) and ideally extends to either side. Tuning to longer wavelengths will extend through the L band out to 1625 nm or even farther through the ultra-long U band to 1675 nm. On the short wavelength side, the S band goes to 1460 nm, after which the extended E band transmits only in fibers without water absorption. The original O band lies between 1260 and 1360 nm. An ideal tuning range would extend throughout all these optical fiber transmission bands.

Two kinds of tunable lasers have application to fiber optical communications. The first is a laser with a set of fixed wavelengths at the ITU frequencies that can be tuned to any wavelength on the grid and operated permanently at that frequency. This approach may be cost-effective because network operators do not have to stock-pile lasers at each of the ITU frequencies; they can purchase a few identical tunable lasers and set them to the required frequency when replacements are needed. The other kind of tunable laser is agile in frequency; it can be tuned in real time to whatever frequency is open to use within the system. These agile tunable lasers offer the greatest systems potential, perhaps someday enabling wavelength switching even at high-speed packet rates. These agile lasers also tend to be more expensive, and at the present time somewhat less reliable.

Most tunable laser diodes in fiber optics communications can be divided into three categories: an array of different frequency lasers with a moveable external mirror, a tunable external cavity laser (ECL), and a monolithic tunable laser. Each of these will be discussed in the following sections.

Array with External Mirror

Many applications do not require rapid change in frequency, such as for replacement lasers. In this case it is sufficient to have an array of DFB lasers, each operating at a different frequency on the ITU grid, and to move an external mirror to align the desired laser to the output fiber. Fujitsu, NTT, Furukawa, and Santur have all presented this approach at various conferences. A typical system might contain a collimating lens, a tilting mirror [often a MEMS (micro-electromechanical structure)] and a lens that focuses into the fiber. The challenge is to develop a miniature device that is cost-effective. A MEMS mirror may be fast enough to enable agile wavelength switching at the circuit level; speeds are typically milliseconds but may advance to a few tenths of a millisecond.

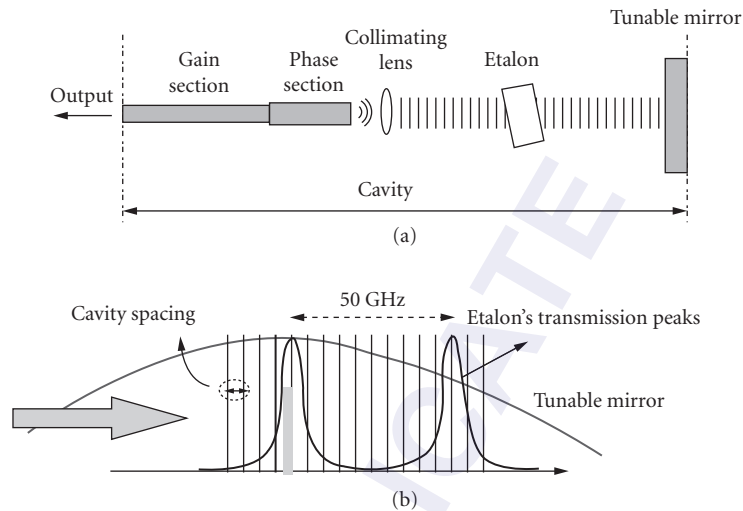


FIGURE 20 Tunable external cavity laser: (a) geometry, showing etalon set at the ITU spacing of 50 GHz and a tunable mirror and (b) spectrum of mirror reflectivity. Spectral maximum can be tuned to pick the desired etalon transmission maximum. The cavity spacing must match the two other maxima, which is done with the phase section. (Adapted from Ref. 31.)

External Cavity Laser

Tunable lasers may consist of a single laser diode in a tiny, wavelength-tunable external cavity, as shown in Fig 20a. A tunable single-mode is achieved by inserting two etalons inside the cavity. One is set at the 50-GHz spacing of the ITU standard; the other provides tuning across the etalon grid, as shown in Fig. 20b. A tunable phase section is also required to ensure that the mode selected by the overall laser cavity length adds constructively to the mode selection from the etalon and mirror. The mechanism for tuning has varied from a liquid crystal mirror,³¹ to thermal tuning of two Fabry-Perot filters within the cavity, reported by researchers at Intel. Pirelli is another company that uses an ECL; they have not reported their method of tuning, but previous work from their laboratory suggests that a polished fiber coupler with variable core separation could be used to tune a laser wavelength.

An alternative tunable filter is acousto-optic, fabricated in lithium niobate, which has been shown to have a tuning range of 132 nm, covering the entire L, C, and S bands.³² Stable oscillation was achieved for 167 channels, each separated by 100 GHz, although there is no evidence that this has become a commercial device.

The speed of tuning external cavity lasers to date is on the order of milliseconds, perhaps fast enough to enable circuit switching to different wavelengths; additional research is underway to achieve faster switching times.

Monolithic Tunable Lasers

Integration of all elements on one substrate offers the greatest potential for compact, inexpensive devices that can switch rapidly from one wavelength to another. The aim is to tune across the entire ITU frequency grid, as far as the laser gain spectrum will allow. All monolithic tunable lasers reported to date involve some sort of grating vernier so that a small amount of tuning can result in a large spectral shift. When the periods of two reflection spectra are slightly mismatched, lasing will occur at that pair of reflectivity maxima that are aligned. Inducing a small index change in one mirror

relative to the other causes adjacent reflectivity maxima to come into alignment, shifting the lasing wavelength a large amount for a small index change. However, to achieve continuous tuning between ITU grid frequencies, the phase within the two-mirror cavity must be adjusted so that its mode also matches the chosen ITU frequency.

The refractive index in the gratings is usually changed by current injection, since changing the free-carrier density in a semiconductor alters its refractive index. Free-carrier injection also introduces loss, which is made up for by the semiconductor optical amplifier (SOA). In principle this modulation speed can be as fast as carriers can be injected and removed. Thermal tuning of the refractive index is an important alternative, because of the low thermal conductivity of InP-based materials. Local resistive heating is enough to create the 0.2 percent change in refractive index needed for effective tuning. The electro-optic effect under reverse bias is not used at present, because the effect does not create large-enough index change at moderate voltages.

In order for the grating to be retroreflective, its periodicity must be half the wavelength of light in the medium (or an odd integral of that); this spacing was discussed in the section on DBR gratings (Sec. 13.7). The other requirement is that there be a periodicity Λ at a scale that provides a comb of possible frequencies at the ITU-T grid (like the etalon in the ECL). This is done in monolithic grating devices by installing an overall periodicity to the grating at the grid spacing. One way to do this is with a *sampled grating* (SG), as shown in Fig. 21a. Only samples of the grating are provided, periodically at frequency Λ ; this is usually done by removing periodic regions of a continuous grating. A laser with DBR mirrors containing sampled gratings is called a *sampled grating DBR* (SG-DBR) laser. If the sampling is abrupt, the reflectance spectrum of the overall comb of frequencies will have the conventional sinc function. The comb reflectance spectrum can be made flat by adding a semiconductor optical amplifier in-line with the DBR laser, or inserting an electroabsorption modulator (EAM) (which will be described in Sec. 13.12), or both, as shown in Fig. 21b.³³ The EAM can be used to modulate the laser output, rather than using direct modulation of the laser. JDS Uniphase tunable diode lasers apparently have this geometry.

An alternative approach to achieving a flat spectrum over the tuning range has been to divide the grating into identical elements, each Λ long, and each containing its own structure.³⁴ This concept has been titled the *superstructure grating* (SSG). The periodicity Λ provides the ITU grid frequencies. When the structure of each element is a phase grating that is chirped quadratically (Fig. 22a), the overall reflectance spectrum is roughly constant and the number of frequencies can be very large. Without the quadratic phase grating structure, the amplitude of the overall reflectance spectrum would have the typical sinc function. Figure 22b shows how the desired quadratic phase shift can be achieved with uneven spacing of the grating teeth, and Fig. 22c shows the resulting measured flat reflectance spectrum.³⁴

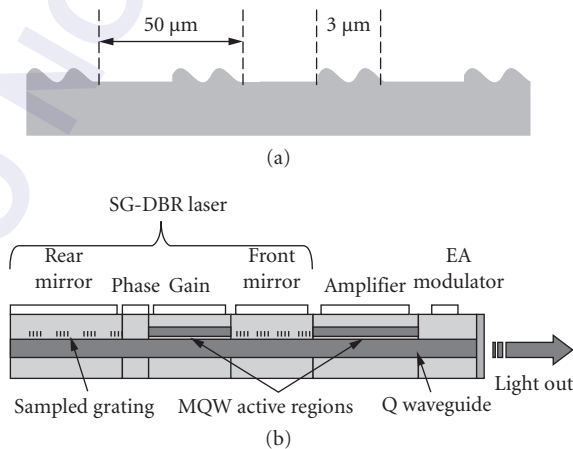


FIGURE 21 Sampled grating: (a) the geometry and (b) sampled grating DBR laser integrated with SOA and EAM.³³

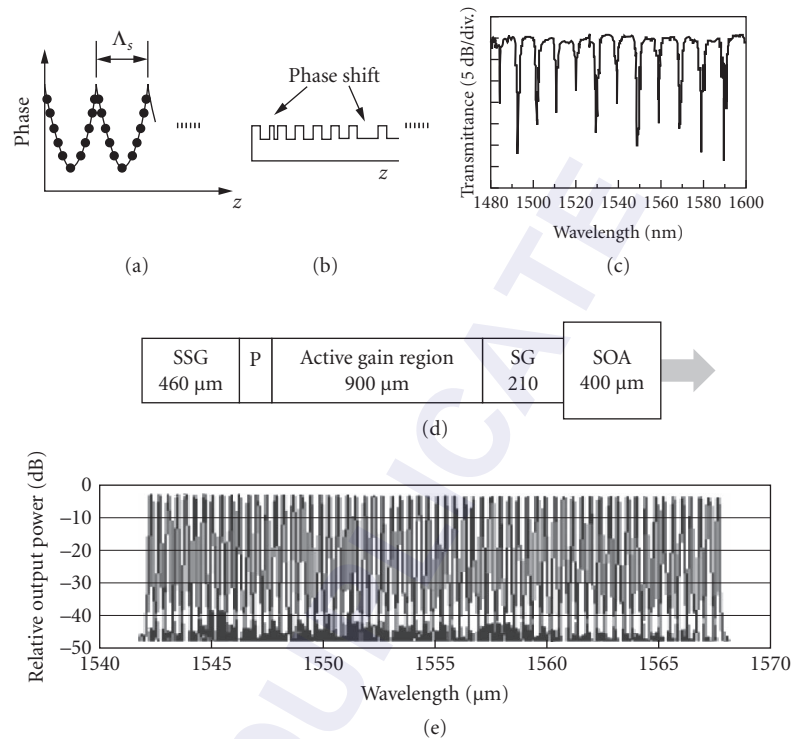


FIGURE 22 Sampled superstructure grating: One period of (a) quadratic phase shift and (b) quadratic phase superstructure; (c) SSG reflectance spectrum;³⁴ (d) geometry for wide-bandwidth integrated tunable laser; and (e) measured tuning output spectrum.³⁵

A relatively new tunable laser diode that can be tuned for DWDM throughout most of the C band contains a front SG-DBR and a rear SSG-DBR.³⁵ Figure 22d shows the design of this device, with lengths given in micrometers (P represents a phase shift region 150 μm long). At the bottom is the spectral output of such a monolithic laser, tuned successively across each wavelength of the ITU grid. A short, low-reflectivity front mirror enables high output power, while keeping the minimum reflectivity that enables wavelength selection based on the Vernier mechanism. A long SSG-DBR was adopted as the rear mirror along with phase control inside the laser cavity, to provide a uniform reflectivity spectrum envelope with a high peak reflectivity (greater than 90%). This monolithic tunable laser includes an integrated SOA for high output power.

A laser tunable for coarse WDM uses a rear reflector comprising a number of equal lengths of uniform phase grating separated by π -phase shifts and a single continuously chirped grating at the front.³⁶ By the correct choice of the number and positions of the phase shifts, the response of the grating can be tailored to produce flat comb of reflection peaks throughout the gain bandwidth. The phase grating was designed to provide seven reflection peaks, each with a 6.8-nm spacing. Over the 300- μm -long front grating is placed a series of short contacts for injecting current into different parts of the grating in a controlled way (Fig. 23). This enables the enhancement of the reflection at a desired wavelength simply by injecting current into a localized part of the chirped grating. The chirp rate was chosen to yield a total reflector bandwidth of around 70 nm. As with other devices, the monolithic chip includes a phase control region and a SOA. As seen in Fig. 23, the waveguide path through the SOA is curved to avoid retroreflection back into the laser cavity—a standard technique.

A different approach is to use a multimode interferometer (MMI) as a Y branch (which will be explained in Sec. 13.11 on modulators).³⁷ This separates backward-going light into two branches, each

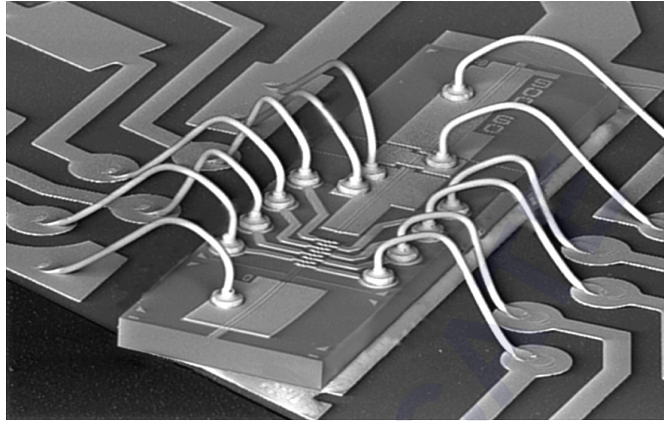


FIGURE 23 Scanning electron microscope image of a monolithically integrated tunable laser for coarse WDM.³⁶

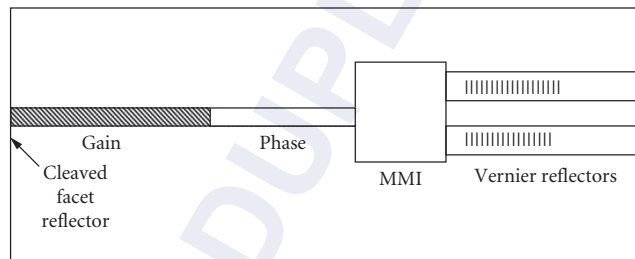


FIGURE 24 Conceptual design of tunable modulated grating DBR laser, including multimode interferometric beam splitter as a Vernier.

reflecting from a grating of a different periodicity, as shown in Fig. 24. The Vernier effect extends the tuning range through parallel coupling of these two *modulated grating* (MG) reflectors with slightly different periods; both reflections are combined at the MMI. The aggregate reflection seen from the input port of the MMI coupler gives a large reflection only when the reflectivity peaks of both gratings align. A large tuning range (40 nm) is obtained for relatively small tuning of a single reflector (by an amount equal to the difference in peak separation). A phase section aligns a longitudinal cavity mode with the overlapping reflectivity peaks. Tunable high-speed direct modulation at 10 Gb/s has been demonstrated with low injection current, with low power consumption and little heat dissipation.

Commercial tunable laser diodes for optical communications are still in their infancy; it is still unclear which of these technologies will be optimum for practical systems. The ultimate question is whether for any given application it is worth the added cost and complexity of tunability.

13.9 LIGHT-EMITTING DIODES

Sources for low-cost fiber communication systems, such as used for communicating data, have traditionally used light-emitting diodes (LEDs). These may be *edge-emitting LEDs* (E-LEDs), which resemble laser diodes, or *surface-emitting LEDs* (S-LEDs), which emit light from the surface of the diode and can be butt-coupled to multimode fibers. The S-LEDs resemble today's VCSELs, discussed

in the following section. The LED can be considered a laser diode operated below threshold, but it must be specially designed to maximize its output.

When a pn junction is forward biased, electrons are injected from the n region and holes are injected from the p region into the active region. When free electrons and free holes coexist with comparable momentum, they will combine and may emit photons of energy near that of the bandgap, resulting in an LED. The process is called *injection-* or *electroluminescence*, since injected carriers recombine and emit light by spontaneous emission. A semiconductor laser diode below threshold acts as an LED. Indeed, a laser diode without mirrors *is* an LED. Because LEDs have no threshold, they usually are not as critical to operate and are often much less expensive because they do not require the fabrication step to provide optical feedback (in the form of cleaved facets or DFB). Because the LED operates by spontaneous emission, it is an incoherent light source, typically emitted from a larger aperture (out the top surface) with a wider far-field angle and a much wider wavelength range (30 to 50 nm). In addition, LEDs are slower to modulate than laser diodes because stimulated emission does not remove carriers. Nonetheless, they can be excellent sources for inexpensive multimode fiber communication systems, as they use simpler drive circuitry. They are longer lived, exhibit more linear input-output characteristics, are less temperature sensitive, and are essentially noise-free electrical-to-optical converters. The disadvantages are lower power output, smaller modulation bandwidths, and pulse distortion in fiber systems because of the wide wavelength band emitted. Some general characteristics of LEDs are discussed in Chap. 17 in Vol. II of this *Handbook*.

In fiber communication systems, LEDs are used for low-cost, high-reliability sources typically operating with graded index multimode fibers (core diameters approximately 62 μm) at data rates up to 622 Mb/s. For short fiber lengths they may be used with step-index plastic fibers. The emission wavelength will be at the bandgap of the active region in the LED; different alloys and materials have different bandgaps. For medium-range distances up to ~ 10 km (limited by modal dispersion), LEDs of $\text{In}_x\text{Ga}_y\text{As}_{1-x-y}\text{P}_{1-y}$ grown on InP and operating at $\lambda = 1.3 \mu\text{m}$ offer low-cost, high-reliability transmitters. For short-distance systems, up to 2 km, GaAs LEDs operating near $\lambda = 850$ nm are used, because they have the lowest cost, both to fabricate and to operate, and the least temperature dependence. The link length is limited to ~ 2 km because of chromatic dispersion in the fiber and the finite linewidth of the LED. For lower data rates (a few megabits per second) and short distances (a few tens of meters), very inexpensive systems consisting of red-emitting LEDs with $\text{Al}_x\text{Ga}_{1-x}\text{As}$ or $\text{GaIn}_y\text{P}_{1-y}$ active regions emitting at 650 nm can be used with plastic fibers and standard silicon detectors. The 650-nm wavelength is a window in the absorption in acrylic plastic fiber, where the loss is ~ 0.3 dB/m; a number of companies now offer 650-nm LEDs.

A typical GaAs LED heterostructure is shown in Fig. 25, with (a) showing the device geometry and (b) showing the heterostructure bandgap under forward bias. The forward-biased pn junction injects electrons and holes into the narrowband GaAs active region. The $\text{Al}_x\text{Ga}_{1-x}\text{As}$ cladding layers confine the carriers in the active region. High-speed operation requires high levels of injection (and/or doping) so that the spontaneous recombination rate of electrons and holes is very high. This means that the active region should be very thin. However, nonradiative recombination increases at high carrier concentrations, so there is a trade-off between internal quantum efficiency and speed. Under some conditions, LED performance is improved by using QWs or strained layers. The improvement is not as marked as with lasers, however.

Spontaneous emission causes light to be emitted in all directions inside the active layer, with an internal quantum efficiency that may approach 100 percent in these direct band semiconductors. However, only the light that gets out of the LED and into the fiber is useful in a communication system, as illustrated in Fig. 25a. The challenge, then, is to collect as much light as possible into the fiber end. The simplest approach is to butt-couple a multimode fiber to the S-LED surface as shown. Light emitted at too large an angle will not be guided in the fiber core, or will miss the core altogether. Light from the edge-emitting E-LED (in the geometry of Fig. 1, with antireflection-coated cleaved facets) is more directional and can be focused into a single-mode fiber. Its inexpensive fabrication and integration process makes the S-LED common for inexpensive data communication. The E-LED has a niche in its ability to couple with reasonable efficiency into single-mode fibers. Both LED types can be modulated at bit rates up to 622 Mbps, an asynchronous transfer mode (ATM) standard, but many commercial LEDs have considerably smaller modulation bandwidths.

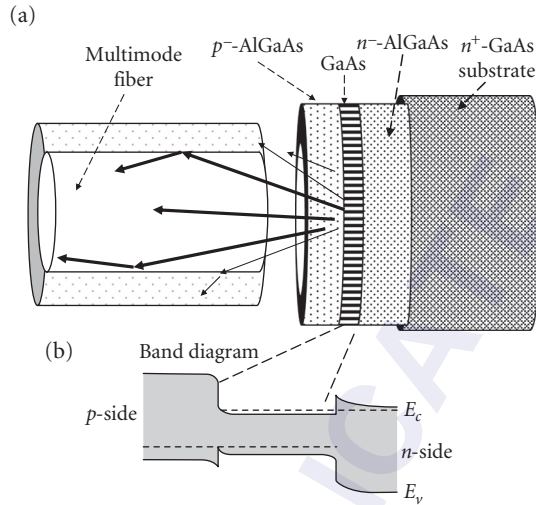


FIGURE 25 GaAs light-emitting diode (LED) structure: (a) cross-section of surface-emitting LED aligned to a multimode fiber, showing rays that are guided by the fiber core and rays that cannot be captured by the fiber and (b) conduction band E_c and valence band E_v as a function of distance through the LED.

Surface-Emitting LEDs

The coupling efficiency of an S-LED butt-coupled to a multimode fiber (shown in Fig. 26a) is typically small, unless methods are employed to optimize it. Because light is spontaneously emitted in all internal directions, only half is emitted toward the top surface. In addition, light emitted at too great an angle to the surface normal is totally internally reflected back down and is lost (although it may be reabsorbed, creating more electron-hole pairs). The critical angle for total internal reflection between the semiconductor of refractive index n_s and the output medium (air or plastic encapsulant) of refractive index n_o is given by $\sin \theta_c = n_o/n_s$. The refractive index of GaAs is $n_s \sim 3.3$, when the output medium is air, the critical angle $\theta_c \sim 18^\circ$. Because this angle is so small, less than 2 percent of the total internal spontaneous emission can come out the top surface, at any angle. A butt-coupled fiber can accept only spontaneous emission at those external angles that are smaller than its numerical aperture. For a typical fiber $NA \approx 0.25$, this corresponds to an external angle (in air) of 14° , which corresponds to only 4.4° inside the GaAs. This means

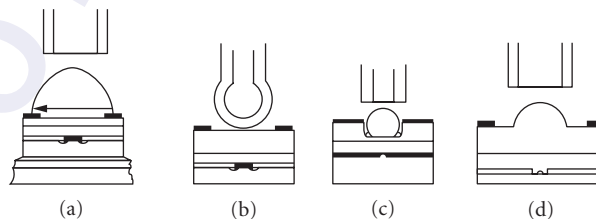


FIGURE 26 Typical geometries for coupling from LEDs into fibers: (a) hemispherical lens attached with encapsulating plastic; (b) lensed fiber tip; (c) microlens aligned through use of an etched well; and (d) spherical semiconductor surface formed on the substrate side of the LED.

that the cone of spontaneous emission that can be accepted by the fiber in this simple geometry is only ~0.2 percent of the entire spontaneous emission! Fresnel reflection losses make this number even smaller.

For InP-based LEDs, operating in the 1.3- or 1.55- μm wavelength region, the substrate is transparent and LED light can be emitted out the substrate. In this geometry the top contact to the p -type material no longer need be a ring; it can be solid and reflective, so light emitted backward can be reflected toward the substrate, increasing the efficiency by a factor of 2.

The coupling efficiency can be increased in a variety of other ways, as shown in Fig. 26. The LED source is incoherent, a lambertian emitter, and follows the law of imaging optics: A lens can be used to reduce the angle of divergence of LED light, but this will enlarge the apparent source. The image of the LED source must be smaller than the fiber into which it is to be coupled. Unlike a laser, the LED has no modal interference and the output of a well-designed LED has a smooth intensity distribution that lends itself to imaging.

The LED can be encapsulated in materials such as plastic or epoxy, with direct attachment to a focusing lens (Fig. 26a). The output cone angle will depend on the design of this encapsulating lens. Even with a parabolic surface, the finite size of the emitting aperture and resulting aberrations will be the limiting consideration. In general, the user must know both the area of the emitting aperture and the angular divergence in order to optimize coupling efficiency into a fiber. Typical commercially available LEDs at $\lambda = 850\text{ nm}$ for fiber optic applications have external half-angles of $\sim 25^\circ$ without a lens and $\sim 10^\circ$ with a lens, suitable for butt-coupling to multimode fiber.

Improvement can also be achieved by lensing the pigtailed fiber to increase its acceptance angle (Fig. 26b); this example shows the light emitted through and out the substrate. Another alternative is to place a microlens between the LED and the fiber (Fig. 26c), possibly within an etched well. Substrate-side emission enables a very effective geometry for capturing light by means of a domed surface fabricated directly on the substrate, as shown in in Fig. 26d. Because the refractive index of encapsulating plastic is < 1.5 , compared to 3.3 of the semiconductor, only a dome etched within the semiconductor can entirely eliminate total internal reflection. Integrated semiconductor domes require advanced semiconductor fabrication technology, but have proven effective. In GaAs diodes the substrate is absorptive, but etching a well in the substrate and inserting a fiber can serve to collect backside emission. For any of these geometries, improvement in efficiency by as much as a factor of 2 can be obtained if a mirror is provided to reflect backward-emitted light. This mirror can be either metal or a dielectric stack deposited at the air-semiconductor interface, or even a DBR mirror grown within the semiconductor structure.

Current must be confined to the surface area of emission, which is typically 25 to 75 μm in diameter. This is done by constricting the flow of injection current by mesa etching or by using an oxide-defined electrode. Regrowth using $n\text{-p-n}$ blocking layers or semi-insulating material in the surrounding areas (as in lasers) has the advantage of reducing thermal heating. Surface-emitting LEDs require that light be emitted out of the surface in a gaussian-like pattern; it must not be obscured by the contacting electrode. Typically, a highly conductive cap layer brings the current in from a ring electrode. In InP-based devices, when light is collected out of the substrate side, a solid top electrode can be reflective and electrical contact may be made to the substrate surrounding an optical output aperture.

A typical S-LED at 1310-nm wavelength might couple 35 μW of light into a pigtailed multimode fiber of 62.5 μm diameter, if driven with 62-mA input current at $\sim 1\text{-V}$ forward bias, for an efficiency of less than 0.06 percent. The spectral width is 160 nm, and the rise or fall time of 3 ns means 200-MHz modulation capability. By contrast, a similar S-LED at 850-nm wavelength under similar drive conditions has comparable power at half the modulation speed and a fifth the bandwidth.³⁸

Improved performance has been obtained by sandwiching the active layer between DBR mirrors to form a resonant cavity (RC-LED). This reduces the spontaneous emission linewidth, thereby increasing the modulation bandwidth that can be transmitted through dispersive fibers. The RC-LEDs look very much like VCSELs operated below threshold (see the following section). The resonant cavity promotes emission into resonances supported by the cavity while suppressing off-resonance emission. The cavity narrows the emission angles and increases the external quantum efficiency to greater than 20 percent. By controlling reflectivities and drive current, the output can be tailored anywhere from a wide spectrum LED to a narrow spectrum laser. Below threshold the RC-LED will have the operating characteristics of a spontaneous LED, as described in this section. Above threshold, the resonant structure becomes a VCSEL. The more exacting fabrication of these devices increases their cost, however.

A typical GaAs LED has a spectral width of 50 nm, while an RC-LED has a spectral width of 6.7 nm and 5 times the output power when coupled into a multimode fiber. The narrower spectrum also means much less pulse-broadening in fiber communication systems. The RC-LED, with a spot size of 20 μm , will have an output power typically a quarter that of a VCSEL laser. In comparison, a typical VCSEL has a spot size of 8 μm and a spectral width of 3.2 nm.

RC-LEDs have enhanced modulation bandwidth compared to conventional LEDs, due to a higher carrier density for a given current, which leads to a reduction in the spontaneous lifetime. However, the RC-LEDs bandwidth is current dependent. The fastest modulation speed is obtained by lasers, of course, but they have a highly nonlinear light-current relationship, unlike LEDs.

Edge-Emitting LEDs

Edge-emitting LEDs (E-LEDs or EELEDs) have a geometry that is similar to that of a conventional LD (Fig. 1), but without a feedback cavity. That is, light travels back and forth in the plane of the active region of an E-LED and is emitted out of one antireflection-coated cleaved end. As in a laser, the active layer is 0.1 to 0.2 μm thick. Because the light in an E-LED is waveguided in the out-of-plane dimension and is lambertian in-plane, the output radiation pattern will be elliptical, with the largest divergence in-plane with full width at half-maximum (FWHM) angle of 120° . The out-of-plane guided direction typically radiates with a 30° half-angle. An elliptical collimating lens will be needed to optimally couple light into a fiber. The efficiency can be doubled by depositing a reflector on the back facet of the E-LED, just as in the case of a laser.

Edge-emitting LEDs can be coupled into fibers with greater efficiency because their source area is smaller than that of a S-LED. However, the alignment and packaging is more cumbersome than with S-LEDs. Typically, E-LEDs can be obtained already pigtailed to fibers. Edge-emitting diodes can be coupled into *single-mode* fiber with modest efficiency. A single-mode fiber pigtailed to an E-LED can typically transmit 30 μW at 150-mA drive at 1 V, for an overall efficiency of 0.04 percent.³⁹ This efficiency is comparable to the emission of surface-emitting lasers into *multimode* fiber (with an area 50 times larger than single-mode fibers). Because of their wide emission wavelength bandwidth, E-LEDs are typically used as low-coherence sources for fiber sensor applications, rather than in communications applications.

Operating Characteristics of LEDs

In an LED, the output optical power P_{opt} is linearly proportional to the drive current I ; the relation defines the output efficiency η :

$$P_{\text{out}} = \frac{\eta h \nu I}{e} \quad (42)$$

This efficiency is affected by the geometry of the LED. The power coupled into a fiber is further reduced by the coupling efficiency between the LED emitter and the fiber, which depends on the location, size, and numerical aperture of the fiber as well as on the spatial distribution of the LED output light and the optics of any intervening lens. The *internal* quantum efficiency (ratio of emitted photons to incident electrons) is usually close to 100 percent.

Figure 27 shows a typical result for power coupled into a graded index multimode fiber as a function of current for various temperatures. The nonlinearity in the light out versus current, which is much less than in a laser diode, nevertheless causes some nonlinearity in the modulation of LEDs. This LED nonlinearity arises from both material properties and device configuration; it may be made worse by ohmic heating at high drive currents. The residual nonlinearity is an important characteristic of any LED used in communication systems. Edge emitters are typically less linear because they operate nearer the amplified spontaneous limit.

The InP-based S-LED shows approximately 10 percent reduction in output power for a 25°C increase in temperature (compared to ~ 50 percent reduction for a typical laser). Unlike a laser, there is no temperature-dependent threshold. Also, the geometric factors that determine the fraction of light

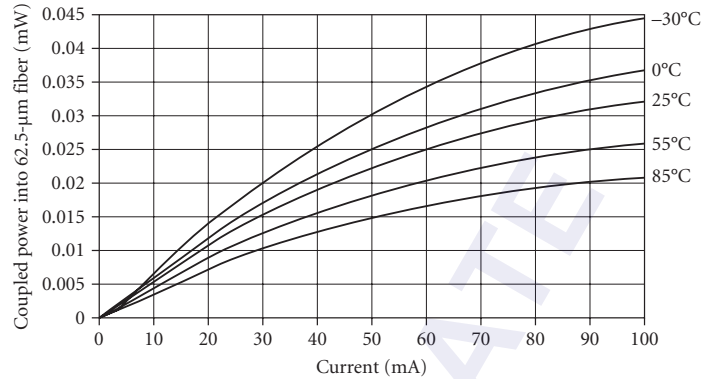


FIGURE 27 Optical power coupled from an *InGaAsP* S-LED into graded index fiber at 1.3- μm wavelength as a function of drive current, for several temperatures.⁴⁰

emitted from the LED are not temperature dependent. Nonetheless, InP-based LEDs have a stronger temperature dependence than GaAs-based LEDs, because of more nonradiative recombination, particularly at the high injection levels required by high-speed LEDs.

The spectrum of the incoherent light emitted from an LED has a roughly gaussian shape with an FWHM around 40 nm in the case of a typical GaAs/AlGaAs LED operating near $\lambda = 0.8 \mu\text{m}$. This bandwidth, along with chromatic dispersion in graded index fibers, limits the distance over which these LEDs can be used in fiber systems. InGaAsP/InP LEDs have wider linewidths (due to alloy scattering, heavy doping, and temperature fluctuations), which depend on the details of their design. As temperature increases, the peak of the spectrum shifts to longer wavelength and the spectrum widens. The variation of the central wavelength with temperature at $\lambda = 1.3 \mu\text{m}$, is approximately 5 meV/°C. However, graded index fibers have negligible chromatic dispersion at this wavelength, so this is not usually a problem; if it is, heat sinking and/or cooling can be provided. Resonant cavity LEDs can provide narrower linewidths.

LEDs do not suffer from the catastrophic optical damage that lasers do, because of their lower optical power densities. However, they do degrade with time; 10^6 to 10^9 hours can be expected. Because degradation processes have an exponential dependence on temperature, LED life can be shortened by operating at excessive temperatures. Experiments with thermally accelerated life testing suggest that the power out P varies with time t as

$$P(t) = P(0) \exp(-qt) \quad (43)$$

where $q = q_0 \exp(-W_a/k_B T)$, with W_a as the activation energy, k_B as Boltzman's constant, and T as temperature. In GaAs LEDs, W_a is 0.6 to 1 eV. Of course this assumes that the LEDs are placed in a proper electrical circuit.

LED light is typically unpolarized, since there is no preferred polarization for spontaneous emission.

Transient Response Most LEDs respond in times faster than 1 μs ; with optimization, they can reach the nanosecond response times needed for optical communication systems. To achieve the 125 Mb/s rate of the fiber distributed data interface (FDDI) standard requires maximum rise and fall times of 3.5 ns; to achieve the 622 Mb/s rate of the asynchronous transfer mode standard, the necessary times drop to 0.7 ns.

The speed of an LED is limited by the recombination time of injected carriers; it does not have the turn-on delay of lasers, nor any relaxation oscillations, but it also does not have the fast decay of stimulated emission. The LED intrinsic frequency response (defined as the ratio of the AC components of the emitted light to the DC value) is⁴¹

$$r(\omega) = (1 + \omega^2 \tau^2)^{-1/2} \quad (44)$$

where τ is the minority carrier lifetime in the injected region. This shows that high-speed LEDs require small minority carrier lifetimes. The square-root dependence comes out of the rate equation solution.

When the active region is doped higher than the injected carrier density, the lifetime τ_L decreases as the background doping density N_o increases:

$$\frac{1}{\tau_L} = BN_o \quad (45)$$

The challenge is to provide high levels of doping without increasing the nonradiative recombination. Typical high-speed response is about 1 ns, although doping with beryllium (or carbon) at levels as high as $7 \times 10^{19} \text{ cm}^{-3}$ has allowed speeds to increase to as much as 0.1 ns, resulting in a cutoff frequency of 1.7 GHz (at the sacrifice of some efficiency).⁴²

When operating in the high-injection regime, the injected carrier density N can be much larger than the doping density, and $1/\tau_H = BN$, where N is created by a current density J such that $N = J\tau/ed$. Combining these two equations

$$\frac{1}{\tau_H} = \left(\frac{BJ}{ed} \right)^{1/2} \quad (46)$$

The recombination time may be reduced by thinning the active region and by increasing the drive current. However, too much injection may lead to thermal problems, which in turn may cause modulation nonlinearity. LEDs with thin active layers operated in the high-injection regime will have the fastest response. Bandwidths in excess of 1 GHz have been achieved in practical LEDs.

Because LEDs have such wide wavelength spectra, frequency chirping is negligible. Because LEDs do not have optical cavities, as do lasers, they will not have modal interference and noise. Also, there will not be strong feedback effects coming from external fiber facets, such as the coherence collapse. Because of their inherent light-current linearity at moderate drive levels, the modulation response of LEDs should be a direct measure of their frequency response. They add no noise to the circuit, and they add distortion only at the highest drive levels.

Drive Circuitry and Packaging The LED is operated under sufficient forward bias to flatten the bands of the *pn* junction. This voltage depends on the bandgap and doping and is typically between 1 and 2 V. The current will be converted directly to light; typically a hundred milliamperes is required to produce a few milliwatts of output into a fiber, with a series resistor used to limit the current.

The LED is modulated by varying the drive current. A typical circuit might apply the signal to the base circuit of a transistor connected in series with the LED and a current-limiting resistor. The variation in current flowing through the LED (and therefore in the light out) is proportional to the input voltage in the base circuit. LEDs are typically mounted on standard headers such as TO-18 or TO-46 cans; SMA and ST connectors are also used. The header is covered by a metal cap with a clear glass top through which light can pass.

13.10 VERTICAL CAVITY SURFACE-EMITTING LASERS

The *vertical cavity surface-emitting laser* (VCSEL) has advantages for low-cost data transmission. The use of a laser means that multi-gigahertz modulation is possible, and the stimulated emission is directional, rather than the isotropic spontaneous emission of LEDs, leading to much higher efficiencies. Because the light is emitted directly from the surface, a single or multimode fiber can be directly butt-coupled with an inexpensive mounting technology, and the coupling efficiency can be very high. VCSELs can be fabricated in linear arrays that can be coupled inexpensively to linear arrays of fibers for parallel fiber interconnects with aggregate bit rates of several gigabits per second,

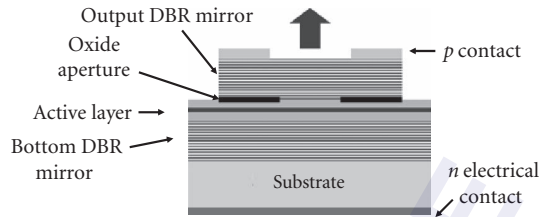


FIGURE 28 Cross-sectional view of an oxidized GaAs VCSEL. The oxidized AlAs layer aperture is shown in black, as is the active layer.⁴³

amortizing the alignment cost over the number of elements in the array. VCSELs lend themselves to two-dimensional arrays as well, which makes them attractive to use with smart modulating pixels. The planar fabrication of VCSELs allows for wafer-scale testing, another cost saving.

Advantages of VCSELs include single longitudinal-mode operation and a circular emission pattern that allows for better coupling into optical fiber without the need for beam-shaping optics. The planar structure also allows much easier fabrication and testing, resulting in higher device yield. VCSELs can also be fabricated in two-dimensional laser arrays with each laser operating at a slightly different frequency, for WDM.

The VCSEL requires mirrors on the top and bottom of the active layer, forming a vertical cavity, as shown in Fig. 28. These lasers use the fact that a DBR (or a multilayer quarter-wavelength dielectric stack) can achieve very high reflectance. Thus, the very short path length through a few QWs (at normal incidence to the plane) is sufficient to reach threshold.

The first VCSELs were based on GaAs: either GaAs active regions that emitted at $\lambda = 850$ nm, or strained InGaAs active regions that emitted at $\lambda = 980$ nm. The former are of greater interest in low-cost communication systems because they are compatible with inexpensive silicon detectors. This section explains VCSEL concepts in terms of GaAs-based devices operating in the 850-nm region. With an active region of $\text{Al}_x\text{In}_y\text{Ga}_{1-x-y}\text{P}$, GaAs-based devices can be fabricated in the 650-nm wavelength regime that is optimum for use with plastic fiber. Highly strained InGaAs QWs enable GaAs-based devices to operate near the 1.3- μm wavelength zero-dispersion regime. The addition of a small amount of nitrogen in the active region extends GaAs-based devices to even longer wavelengths. The development of GaAs-based VCSELs has been relatively straightforward because of the existence of epitaxial GaAs/AlGaAs high-reflectivity DBR mirrors.

VCSELs based on InP technology have been more challenging because the reflectivity of InP-based DBR mirrors is quite low. Also, InP-based QWs have lower gain and significant free carrier absorption, especially in the p layers. Also, the quaternaries have poor thermal conductivity, so the active layer has a relatively high temperature. Several novel technologies have been developed that look like they might overcome these challenges, which will be described in subsequent sections.

Number of Quantum Wells

A single quantum well of GaAs requires input current densities of approximately 100 A/cm² to achieve transparency; N wells require N times this current. To keep the threshold current less than 1 kA/cm², then, means active regions with less than 10 QWs. The VCSEL requires an optical standing wave that has a period of a half-wavelength, which is approximately 120 nm in GaAs. The gain region should be confined to the quarter-wavelength region at the peak of the optical standing wave, a region of about 60 nm. Thus, a typical active region might consist of 3 QWs of 8 μm thickness, each separated by approximately 10 nm. Strained QWs have higher gain than unstrained QWs, and can be grown by adding indium or phosphorous into the composition.

InP-based QWs for 1.3 μm and 1.55 μm operation have somewhat lower gain, but typical geometries still use 3 to 5 QWs; higher reflectivity mirrors compensate for lower gain.

Mirror Reflectivity

When the mirror reflectivity R in a laser is close to 1, a simple expression for the threshold gain-length product, $G_L L$ is

$$G_L L = (1 - R_1)(1 - R_2) \quad (47)$$

Typical GaAs lasers have gains $G_L \sim 1000 \text{ cm}^{-1}$. For a quantum well thickness of 10 nm, the gain per quantum well is 10^{-3} , so that with 3 QW, reflectivities of ~ 98 percent for each mirror should be sufficient to achieve threshold. Very often, however, in order to lower the threshold much higher reflectivities are used, particularly on the back mirror.

The on-resonance Bragg mirror reflectivity is the square of the reflection coefficient r , given by

$$r = \frac{1 - (n_f/n_i)(n_l/n_h)^{2N}}{1 + (n_f/n_i)(n_l/n_h)^{2N}} \quad (48)$$

where there are N pairs of quarter-wavelength layers that alternate high-index and low-index (n_h and n_l respectively), and n_f and n_i are the refractive index of the final and initial media, respectively.⁴⁴

For high-reflectance Bragg mirrors, the second term in the numerator and denominator is small, and the reflectivity can be simplified to

$$\epsilon = 1 - R = 1 - r^2 = 4 \left(\frac{n_f}{n_i} \right) \left(\frac{n_l}{n_h} \right)^{2N} \quad (49)$$

Higher reflectivity (smaller ϵ) is provided by either more layer pairs or a larger refractive index difference between the two compositions in the layer pairs. Also, Eq. (49) shows that internal mirrors ($n_f = n_i$) will have a smaller reflectivity than external mirrors ($n_f = 1$), for the same number of layer pairs. If the layer pair consists of GaAs ($n \sim 3.6$) and AlAs ($n \sim 3.0$), a mirror of 15 layer pairs will have an internal reflectivity $R = 98$ percent and an external reflectivity $R = 99.5$ percent. Thirty layer pairs are required to increase the internal reflectivity to 99.96 percent. Bragg mirrors with a smaller fraction of AlAs in the low-index layers will require even more layer pairs to achieve the same reflectivity.

In long-wavelength InP-based VCSELs, the $\text{In}_x\text{Ga}_{1-x}\text{As}_y\text{P}_{1-y}/\text{InP}$ refractive index difference is small and many layers are required to achieve sufficiently high reflectivity. Often the top DBR layer is replaced by a deposited mirror, such as alternating quarter-wave layers of ZnSe and MgF, or amorphous silicon and aluminum oxide. Hybrid mirrors use fewer dielectric layers and can terminate with a gold layer. Because the substrate is transparent, these devices are mounted upside-down, with the gold layer attaching to a thermal heat sink.

Another option has been to replace epitaxial InP-based DBR mirrors with fusion-bonded traditional GaAs-AlAs DBR mirrors that have been grown on GaAs substrates. Some commercial VCSELs have used this approach. Wafer fusion occurs when pressing the two wafers together (after removing oxide off their surfaces) at 15 atm and heating to 630°C under hydrogen for 20 min. Mirrors can be wafer-fused on both sides of the VCSEL by etching away the InP substrate and one of the GaAs substrates.

Finally, new multicomponent compositions show promise of alternative mirrors, especially by adding aluminum. Forty lattice-matched pairs of $\text{Al}_x\text{Ga}_y\text{In}_{1-x-y}\text{As}$ and InP quarter-wave layers have been shown to yield a reflectivity of 99.9 percent.⁴⁵ Some commercial VCSEL manufacturers have favored this all-epitaxial fabrication process.

An entirely different approach, which is not yet commercial, is to push GaAs-based devices out to longer wavelength. Adding dilute amounts of nitrogen yields an active layer of $\text{Ga}_x\text{In}_{1-x}\text{N}_y\text{As}_{1-y}$, which has been shown to produce VCSELs near the 1.3- μm wavelength range. It turns out that adding a small amount of antimony produces QWs of higher quality. This has suggested that the five-component material $\text{Ga}_x\text{In}_{1-x}\text{N}_y\text{As}_z\text{Sb}_{1-y-z}$ has the potential to achieve even lower bandgaps; VCSELs of this composition have been demonstrated, but not yet as far out in wavelength as 1.55 μm .

Electrical Injection and Current Confinement

As shown in Fig. 28, light is emitted through a window hole in the top electrode. It is difficult to inject carriers through the Bragg reflector because the wide bandwidth layers provide potential wells that trap carriers; but this can be overcome with increased operating voltage. Sometimes graded layers are used in the DBR structure to reduce carrier trapping. Current confinement to a finite area was originally done by proton implantation or by etching mesas and then planarizing with polyimide. VCSELs of fairly large diameter ($>10\ \mu\text{m}$) are straight-forward to make and are useful when low threshold and single-mode are not required. Recently a selective oxidation technique has been developed that enables a small oxide-defined current aperture. A high-aluminum fraction $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layer ($\sim 98\%$) is grown above the active layer and a mesa is etched to below that layer. Then a long, soaking, wet-oxidization process selectively creates a ring of native oxide that can stop carrier transport. The chemical reaction moves in from the side of the etched pillar and is stopped when the desired diameter is achieved. Such a resistive aperture confines current only where needed, and can double the maximum conversion efficiency to almost 60 percent.

Threshold voltages less than 3 V are common in diameters $\sim 12\ \mu\text{m}$. The oxide-defined current channel increases the efficiency, but tends to cause multiple transverse modes due to relatively strong oxide-induced index guiding. This may introduce modal noise into fiber communication systems. Single-mode requirements force the diameter to be very small (below 4 to 5 μm) or for the design to incorporate additional features, as discussed in the following section.

Transverse injection eliminates the need for the current to travel through the DBR region, but typically requires even higher voltage. This approach has been proven useful when highly conductive layers are grown just above and below the active region. Because carriers have to travel farther with transverse injection, it is important that these layers have as high mobility as possible. This has been achieved by injecting carriers from both sides through n -type layers. Such a structure can still inject holes into the active layer if a buried tunnel junction is provided, as shown in Fig. 29. The tunnel junction consists of a single layer-pair of very thin highly doped n^{++} and p^{++} layers and must be located near the active layer so that its holes can be utilized. After the first As-based growth step, the tunnel junction is laterally structured by means of standard photolithography and chemical dry etching. It is then regrown with phosphorous-containing n -layers. The lateral areas surrounding the tunnel junction contain an $n\text{pn}$ electronic structure and do not conduct electricity. Only within the area of the tunnel junction are electrons from the n -type InP spacer converted into holes.

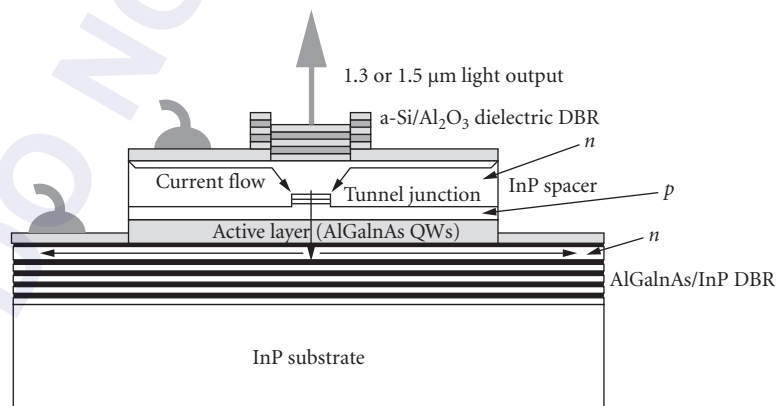


FIGURE 29 BHT geometry for single-mode VCSEL, showing flow of current. (Adapted from Ref. 45.)

Spatial Characteristics of Emitted Light

Single transverse mode operation remains a challenge for VCSELs, particularly at the larger diameters and higher current levels. When modulated, lateral spatial instabilities tend to set in and spatial hole burning causes transverse modes to jump. This can introduce considerable modal noise when coupling VCSEL light into fibers.

The two most common methods to control transverse modes are the same as used to control current spreading: ion implantation and an internal oxidized aperture. Ion implantation keeps the threshold relatively low, but thermal lensing coupled with weak index guiding is insufficient to prevent multilateral-mode operation due to spatial hole burning; also the implanted geometry does not provide inherent polarization discrimination. The current confining aluminum oxide aperture formed by selective oxidization acts as a spatial filter and encourages the laser to operate in low-order modes. Devices with small oxide apertures (2×2 to $4 \times 4 \mu\text{m}^2$) can operate in a single-mode. Devices with $3.5\text{-}\mu\text{m}$ diameters have achieved single-mode output powers on the order of 5 mW, but devices with larger apertures will rapidly operate in multiple transverse modes as the current is raised.⁴⁶

VCSELs with small diameter are limited in the amount of power they can emit while remaining single-mode, and their efficiency falls off as the diameter becomes smaller. A variety of designs have been reported for larger-aperture VCSELs that emit single-mode. This section lists several approaches; some have become commercially available and others are presently at the research stage: (1) Ion implantation and oxide-defined spatial filters have been *combined* with some success at achieving single-mode. (2) Etched pillar mesas favor single-mode operation because they have sidewall scattering losses that are higher for higher-order modes. The requirement is that the mode selective losses must be large enough to overcome the effects of spatial hole burning.⁴⁷ (3) As with traditional lasers, an etched pillar mesa can be overgrown to create a buried heterostructure (BH), providing a real index guide that can be structured to be single-mode.⁴⁸ (4) A BH design can be combined with ion implantation and/or selectively oxidized apertures, for the greatest design flexibility. (5) Surface relief has been integrated on top of the cladding layer (before depositing a top dielectric mirror), physically structuring it so as to eliminate higher-order modes; the surface relief incorporates a quarter-wave ring structure that decreases the reflectivity for higher-order modes.⁴⁹ (6) The VCSEL can be surrounded with a second growth of *higher* refractive semiconductor material, which causes an *antiguide* that preferentially confines the lowest-order mode.⁵⁰ (7) A photonic crystal has been incorporated under the top dielectric mirror, which provides an effective graded index structure that favors maintaining a single-mode.⁵¹

All of these approaches can be used with current injection through the DBR mirrors, or with lateral injection, usually through an oxide aperture. A number of single-mode geometries use buried tunnel junctions (BTJ), which may be made with small enough area to create single-mode lasers without incorporating any other features.⁴⁵ Often higher powers are achieved by combining wider area BTJ along with some of the other approaches outlined above.

Light Out versus Current In

The VCSEL will, in general, have similar L - I performance to edge-emitting laser diodes, with some small differences. Because the acceptance angle for the mode is higher than in edge-emitting diodes, there will be more spontaneous emission, which will show up as a more graceful turn-on of light out versus voltage in. As previously mentioned, the operating voltage is 2 to 3 times that of edge-emitting lasers. Thus, Eq. (5) must be modified to take into account the operating voltage drop across the resistance R of the device. The operating power efficiency is

$$\eta_{\text{eff}} = \eta_D \frac{I_{\text{op}} - I_{\text{th}}}{I_{\text{th}}} \frac{V_g}{V_g + I_{\text{op}} R} \quad (50)$$

Small diameter single-mode VCSELs would typically have a $5\text{-}\mu\text{m}$ radius, a carrier injection efficiency of 80 to 90 percent, an internal optical absorption loss $\alpha_i L$ of 0.003, an optical scattering loss

of 0.001, and a net transmission through the front mirror of 0.005 to 0.0095. Carrier losses reducing the quantum efficiency are typically due to spontaneous emission in the wells, spontaneous emission in the barriers, Auger recombination, and carrier leakage.

Typical commercial VCSELs designed for compatibility with single-mode fiber incorporate an 8- μm proton implantation window and 10- μm -diameter window in the top contact. Such diodes may have threshold voltages of ~ 3 V and threshold currents of a few milliamperes. These lasers may emit up to ~ 2 mW maximum output power. Devices will operate in zero-order transverse spatial mode with gaussian near-field profile when operated with DC drive current less than about twice threshold. When there is more than one spatial mode, or both polarizations, there will usually be kinks in the L - I curve, as with multimode edge-emitting lasers.

Spectral Characteristics

Since the laser cavity is short, the longitudinal modes are much farther apart in wavelength than in a cleaved cavity laser, typically separated by 50 nm, so only one longitudinal mode will appear, and there is longitudinal mode purity. The problem is with *lateral* spatial modes, since at higher power levels the laser does not operate in a single spatial mode. Each spatial mode will have a slightly different wavelength, perhaps separated by 0.01 to 0.02 nm. Lasers that start out as single-mode and single frequency at threshold will often demonstrate frequency broadening at higher currents due to multiple modes, as shown in Fig. 30a.⁵²

Even when the laser operates in a single spatial mode, it may have two orthogonal directions of polarization (discussed next), that will exhibit different frequencies, as shown in Fig. 30b.⁵³ Thus both single-mode and polarization stability are required to obtain a true single-mode.

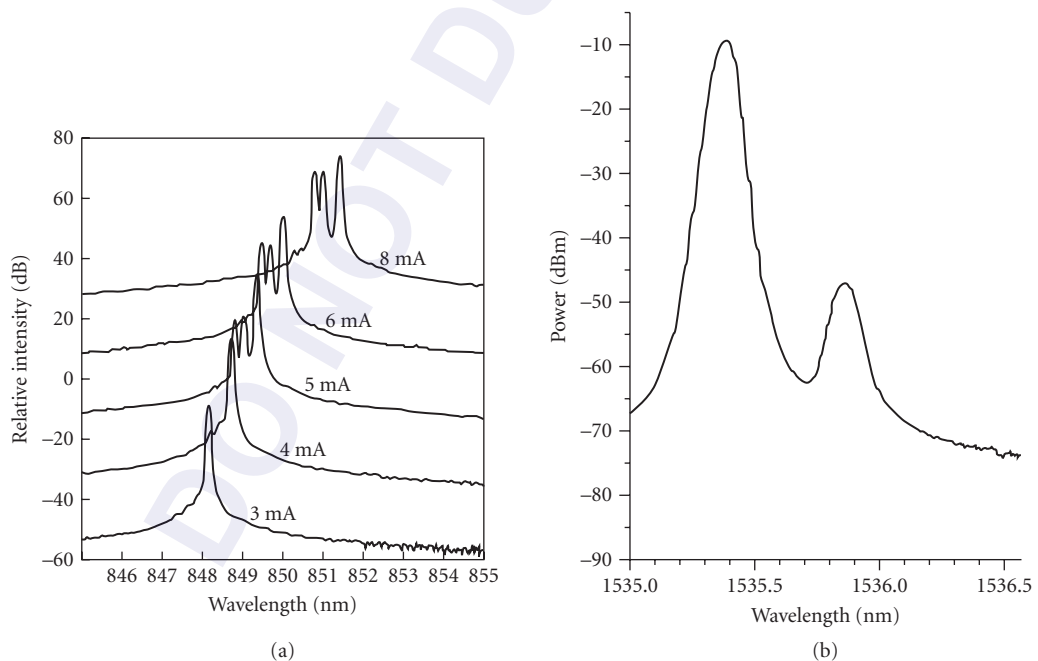


FIGURE 30 VCSEL spectra: (a) emission spectra recorded at different injection currents for BH-VCSELs of 10 μm diameter⁵² and (b) different emission spectra due to different polarizations of a BJT single-mode VCSEL.⁵³

When a VCSEL is modulated, lateral spatial instabilities may set in and spatial hole burning may cause transverse modes to jump. This can broaden the spectrum. In addition, external reflections can cause instabilities and increased relative intensity noise, just as in edge-emitting lasers.⁵⁴ For very short cavities, such as between the VCSEL and a butt-coupled fiber (with ~ 4 percent reflectivity), instabilities do not set in, but the output power can be affected by the additional mirror, which forms a Fabry-Perot cavity with the output mirror and can reduce or increase its effective reflectivity, depending on the round-trip phase difference. When the external reflection comes from ~ 1 cm away, bifurcations and chaos can be introduced with a feedback parameter $F > 10^{-4}$, where $F = C_e \sqrt{f_{\text{ext}}}$, with C_e and f_{ext} as defined in the discussion surrounding Eq. (31). For $R_o = 0.995$, $R_{\text{ext}} = 0.04$, the feedback parameter $F \sim 10^{-3}$, instabilities can be observed if reflections get back into the VCSEL.

Polarization

A VCSEL with a circular aperture has no preferred polarization state. The output tends to oscillate in linear but random polarization states, which may wander with time (and temperature) and may have slightly different emission wavelengths (Fig. 30*b*). Polarization-preserving VCSELs require breaking the symmetry by introducing anisotropy in the optical gain or loss. Some polarization selection may arise from an elliptical current aperture. The strongest polarization selectivity has come from growth on (311) GaAs substrates, which causes anisotropic gain.

Commercial VCSELs

The most readily available VCSELs are GaAs-based, emitting at 850-nm wavelength. Commercial specifications for these devices list typical multimode output powers from 1 to 2.4 mW and single-mode output powers from 0.5 to 2 mW, depending on design. Drive voltages vary from 1.8 to 3 V, with series resistance typically about 100 Ω . Spectral width for multimode lasers is about 0.1 nm, while for single-mode lasers it can be as narrow as 100 MHz. Beam divergence FWHM is typically 18° to 25° for multimode lasers and 8° to 12° for single-mode VCSELs, with between 20 to 30 dB sidemode suppression.⁵⁵

Red VCSELs, emitting at 665 nm, are available from fewer suppliers, and have output powers of 1 mW, threshold currents between 0.6 and 2.5 mA, and operating voltages of 2.8 to 3.5 V. Their divergence angle is 14° to 20° and slope efficiency is 0.9 mW/mA, with sidemode suppression between 14 and 50 dB. Reported bandwidths are 3 to 3.5 GHz.⁵⁶

Long-wavelength VCSELs have output powers between 0.7 and 1 mW, with threshold currents between 1.1 and 2.5 mA. Series resistance is 100 Ω , with operating voltage between 2 and 3 V. Single-mode spectral width is 30 MHz and modulation bandwidth is 3 GHz. Sidemode suppression is between 30 and 40 dB, with slope efficiency of 0.2 mW/mA and an angular divergence between 9° and 20°.⁵⁷

13.11 LITHIUM NIOBATE MODULATORS

The most direct way to create a modulated optical signal for communication application is to directly modulate the current driving the laser diode. However, as discussed in Secs. 13.4 and 13.5, this may cause turn-on delay, relaxation oscillation, mode-hopping, and/or chirping of the optical wavelength. Therefore, an alternative often used is to operate the laser in a continuous manner and to place a modulator after the laser. This modulator turns the laser light off and on without impacting the laser itself. The modulator can be butt-coupled directly to the laser, located in the laser chip package and optically coupled by a microlens, or remotely attached by means of a fiber pigtail between the laser and modulator.

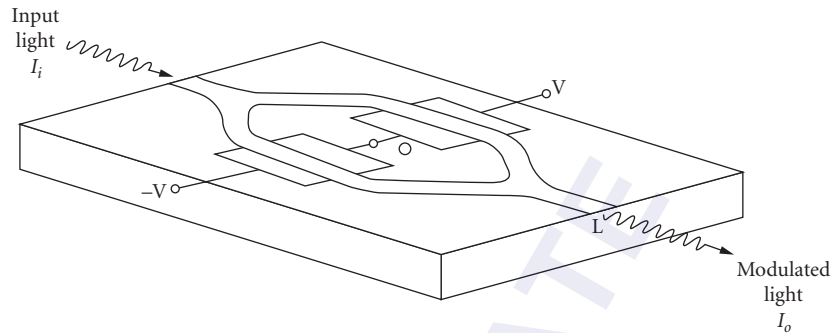


FIGURE 31 Y-branch interferometric modulator in the “push-pull” configuration. Center electrodes are grounded. Opposite polarity electrodes are placed on the outsides of the waveguides. Light is modulated by applying positive or negative voltage to the outer electrodes.

Lithium niobate modulators have become one of the main technologies used for high-speed modulation of continuous-wave (CW) diode lasers, particularly in applications (such as cable television) where extremely linear modulation is required, or where chirp is to be avoided at all costs. These modulators operate by the electro-optic effect, in which an applied electric field changes the refractive index. Integrated optic waveguide modulators are fabricated by diffusion into lithium niobate substrates. The end faces are polished and butt-coupled (or lens-coupled) to a single-mode fiber pigtail (or to the laser driver itself). This section describes the electro-optic effect in lithium niobate, its use as a phase modulator and an intensity modulator, considerations for high-speed operation, and the difficulties in achieving polarization independence.⁵⁸

Most common is the Y-branch interferometric modulator shown in Fig. 31, discussed in a following subsection. The waveguides that are used for these modulators are fabricated in lithium niobate either by diffusing titanium into the substrate from a metallic titanium strip or by means of ion exchange. The waveguide pattern is obtained by photolithography. The standard thermal indiffusion process takes place in air at 1050°C over 10 hours. An 8- μm -wide, 50-nm thick strip of titanium creates a fiber-compatible single-mode at $\lambda = 1.3 \mu\text{m}$. The process introduces ~ 1.5 percent titanium at the surface, with a diffusion profile depth of $\sim 4 \mu\text{m}$. The result is a waveguide with increased extraordinary refractive index of 0.009 at the surface and an ordinary refractive index change of ~ 0.006 . A typical modulator will incorporate aluminum electrodes 2 cm long, deposited on either side of the waveguides, with a gap of 10 μm .

In the case of ion exchange, the lithium niobate sample is immersed in a melt containing a large proton concentration (typically benzoic acid or pyrophosphoric acid at $>170^\circ\text{C}$), with nonwaveguide areas protected from diffusion by masking; the lithium near the surface of the substrate is replaced by protons, which increases the refractive index. Ion-exchange alters only the extraordinary polarization; that is, only light polarized parallel to the z axis is waveguided. Thus, it is possible in lithium niobate to construct a polarization-independent modulator with titanium indiffusion, but not with proton-exchange. Nonetheless, ion exchange creates a much larger refractive index change (~ 0.12), which provides more flexibility in modulator design. Annealing after diffusion can reduce insertion loss and restore the degraded electro-optic effect. Interferometric modulators with moderate index changes ($\Delta n < 0.02$) are insensitive to aging at temperatures of 95°C or below. Using higher index change devices, or higher temperatures, may lead to some degradation with time. Tapered waveguides can be fabricated easily by ion exchange for high coupling efficiency.⁵⁹

Electro-Optic Effect

The *electro-optic effect* is the change in refractive index that occurs in a noncentrosymmetric crystal in the presence of an applied electric field. The linear electro-optic effect is represented by a

third-rank tensor for the refractive index. However, using symmetry rules it is sufficient to define a reduced tensor r_{ij} , where $i = 1, \dots, 6$ and $j = x, y, z$, denoted as 1, 2, 3. Then, the linear electro-optic effect is traditionally expressed as a linear change in the inverse refractive index tensor squared (see Chap. 7 in this volume):

$$\Delta\left(\frac{1}{n^2}\right)_i = \sum_j r_{ij} E_j \quad j = x, y, z \quad (51)$$

where E_j is the component of the applied electric field in the j th direction. In isotropic materials, r_{ij} is a diagonal tensor. An applied electric field can introduce off-diagonal terms in r_{ij} , as well as change the lengths of the principle dielectric axes. The general case is treated in Chap. 13, Vol. II. In lithium niobate (LiNbO_3), the material of choice for electro-optic modulators, the equations are simplified because the only nonzero components and their magnitudes are⁶⁰

$$\begin{aligned} r_{33} &= 31 \times 10^{-12} \text{ m/V} & r_{13} &= r_{23} = 8.6 \times 10^{-12} \text{ m/V} \\ r_{51} &= r_{42} = 28 \times 10^{-12} \text{ m/V} & r_{22} &= -r_{12} = -r_{61} = 3.4 \times 10^{-12} \text{ m/V} \end{aligned}$$

The crystal orientation is usually chosen so as to obtain the largest electro-optic effect. This means that if the applied electric field is along z , then light polarized along z sees the largest field-induced change in refractive index. Since $\Delta(1/n^2)_3 = \Delta(1/n^2)_z = r_{33}E_z$, performing the difference gives

$$\Delta n_z = -\frac{n^3}{2} r_{33} E_z \Gamma \quad (52)$$

A *filling factor* Γ (also called an *optical-electrical field overlap parameter*) has been included due to the fact that the applied field may not be uniform as it overlaps the waveguide, resulting in an effective field that is somewhat less than 100 percent of the maximum field.

In the general case for the applied electric field along z , the tensor remains diagonal and $\Delta(1/n^2)_1 = r_{13}E_z = \Delta(1/n^2)_2 = r_{23}E_z$, and $\Delta(1/n^2)_3 = r_{33}E_z$. This means that the index ellipsoid has not rotated, its axes have merely changed in length. Light polarized along any of these axes will see a pure phase modulation. Because r_{33} is largest, polarizing the light along z and providing the applied field along z will provide the largest phase modulation for a given field. Light polarized along either x or y will have the same index change, which might be a better direction if polarization-independent modulation is desired. However, this would require light to enter along z , which is the direction in which the field is applied, so it is not practical.

As another example, consider the applied electric field along y . In this case the nonzero terms are

$$\Delta\left(\frac{1}{n^2}\right)_1 = r_{12}E_y \quad \Delta\left(\frac{1}{n^2}\right)_2 = r_{22}E_y = -r_{12}E_y \quad \Delta\left(\frac{1}{n^2}\right)_4 = r_{42}E_y \quad (53)$$

There is now a yz cross-term, coming from r_{42} . Diagonalization of the perturbed tensor finds new principal axes, only slightly rotated about the z axis. Therefore, the principal refractive index changes are essentially along the x and y axes, with the same values as $\Delta(1/n^2)_1$ and $\Delta(1/n^2)_2$ in Eq. (53). If light enters along the z axis without a field applied, both polarizations (x and y) see an ordinary refractive index. With a field applied, both polarizations experience the same phase change (but opposite sign). In a later section titled "Polarization Independence," we describe an interferometric modulator that does not depend on the sign of the phase change. This modulator is polarization independent, using this crystal and applied-field orientation, at the expense of operating at somewhat higher voltages, because $r_{22} < r_{33}$.

Since lithium niobate is an insulator, the direction of the applied field in the material depends on how the electrodes are applied. Figure 32 shows a simple phase modulator. Electrodes that straddle the modulator provide an in-plane field as the field lines intersect the waveguide, as shown in Fig. 32b.

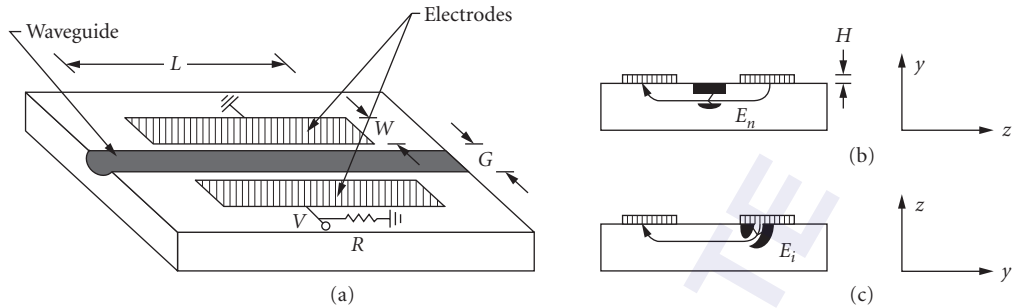


FIGURE 32 (a) Geometry for phase modulation in lithium niobate with electrodes straddling the channel waveguide. (b) End view of (a), showing how the field in the channel is parallel to the surface. (c) End view of a geometry placing one electrode over the channel, showing how the field in the channel is essentially normal to the surface.

This requires the modulator to be *y-cut* LiNbO₃ (the *y* axis is normal to the wafer plane), with the field lines along the *z* direction; *x-cut* LiNbO₃ will perform similarly. Figure 32c shows a modulator in *z-cut* LiNbO₃. In this case, the electrode is placed over the waveguide, with the electric field extending downward through the waveguide (along the *z* direction). The field lines will come up at a second, more distant electrode. In either case, the field may be fringing and nonuniform, which is why the filling factor Γ has been introduced.

Phase Modulation

Applying a field to one of the geometries shown in Fig. 32 results in pure *phase modulation*. The field is roughly V/G , where G is the gap between the two electrodes. For an electrode length L , the phase shift is

$$\Delta\phi = \Delta n_z kL = -\frac{n_o^3}{2} r_{33} \left(\frac{V}{G}\right) \Gamma kL \quad (54)$$

The refractive index for bulk LiNbO₃ is given by⁶¹

$$n_o = 2.195 + \frac{0.037}{[\lambda (\mu\text{m})]^2} \quad \text{and} \quad n_e = 2.122 + \frac{0.031}{[\lambda (\mu\text{m})]^2}$$

Inserting numbers for $\lambda = 1.55 \mu\text{m}$ gives $n_o = 2.21$. When $G = 10 \mu\text{m}$ and $V = 5 \text{ V}$, a π phase shift is expected in a length $L \sim 1 \text{ cm}$.

It can be seen from Eq. (54) that the electro-optic phase shift depends on the product of the length and voltage. Longer modulators can use smaller voltages to achieve a π phase shift. Shorter modulators require higher voltages. Thus, the figure of merit for phase modulators is typically the product of the voltage required to reach π times the length. The modulator just discussed has a 5-V·cm figure of merit.

The electro-optic phase shift has a few direct uses, such as providing a *frequency shifter* (since $\partial\phi/\partial t \propto \Delta\nu$). However, in communication systems this phase shift is generally used in an interferometric configuration to provide intensity modulation, discussed in the following section.

Y-Branch Interferometric (Mach-Zehnder) Modulator

The *interferometric modulator* is shown schematically in Fig. 31. This geometry allows waveguided light from the two branches to interfere, forming the basis of an intensity modulator. The amount of

interference is tunable by providing a relative phase shift on one arm with respect to the other. Light entering a single-mode waveguide is equally divided into the two branches at the Y junction, initially with zero relative phase difference. The guided light then enters the two arms of the waveguide interferometer, which are sufficiently separated that there is no coupling between them. If no voltage is applied to the electrodes, and the arms are exactly the same length, the two guided beams arrive at the second Y junction in phase and enter the output single-mode waveguide in phase. Except for small radiation losses, the output is equal in intensity to the input. However, if a π phase difference is introduced between the two beams via the electro-optic effect, the combined beam has a lateral amplitude profile of odd spatial symmetry. This is a second-order mode and is not supported in a single-mode waveguide. The light is thus forced to radiate into the substrate and is lost. In this way, the device operates as an electrically driven optical intensity on-off modulator. Assuming perfectly equal splitting and combining, the fraction of light transmitted is

$$\eta = \left[\cos\left(\frac{\Delta\phi}{2}\right) \right]^2 \quad (55)$$

where $\Delta\phi$ is the difference in phase experienced by the light in the different arms of the interferometer: $\Delta\phi = \Delta nkL$, where $k = 2\pi/\lambda$, Δn is the difference in refractive index between the two arms, and L is the path length of the field-induced refractive index difference. The voltage at which the transmission goes to zero ($\Delta\phi = \pi$) is usually called V_π . By operating in a push-pull manner, with the index change increasing in one arm and decreasing in the other, the index difference Δn is twice the index change in either arm. This halves the required voltage.

Note that the transmitted light is periodic in phase difference (and therefore voltage). The response depends only on the integrated phase shift and not on the details of its spatial evolution. Therefore, nonuniformities in the electro-optically induced index change that may occur along the interferometer arms do not affect the extinction ratio. This property has made the Mach Zehnder (MZ) modulator the device of choice in communication applications.

For analog applications, where linear modulation is required, the modulator is prebiased to the quarter-wave point (at voltage $V_b = \pi/2$), and the transmission efficiency becomes linear in $V - V_b$ (for moderate excursions):

$$\eta = \frac{1}{2} \left[1 + \sin \frac{\pi(V - V_b)}{2V_\pi} \right] \approx \frac{1}{2} + \frac{\pi}{2} \frac{(V - V_b)}{V_\pi} \quad (56)$$

The electro-optic effect depends on the polarization. For the electrode configuration shown here, the applied field is in the plane of the lithium niobate wafer, and the polarization of the light to be modulated must also be in that plane. This will be the case if a TE-polarized laser diode is butt-coupled (or lens-coupled) with the plane of its active region parallel to the plane of the lithium niobate wafer, and if the wafer is Y -cut. Polarization-independent modulation requires a different orientation, to be described later. First, however, we discuss the electrode requirements for high-speed modulation.

High-Speed Operation

The optimal electrode design depends on how the modulator is to be driven. Because the electrode is on the order of 1 cm long, the fastest devices require traveling-wave electrodes rather than lumped electrodes. Lower-speed modulators can use lumped electrodes, in which the modulator is driven as a capacitor terminated in a parallel resistor matched to the impedance of the source line. The modulation speed depends primarily on the RC time constant determined by the electrode capacitance and the terminating resistance. To a smaller extent, the speed also depends on the resistivity of the electrode itself. The capacitance per unit length is a critical design parameter. This depends on the material dielectric constant, the electrode gap G and the electrode width W . With

increasing G , the capacitance per unit length decreases and the bandwidth-length product increases essentially logarithmically. In LiNbO_3 , when the electrode widths and gap are equal, the capacitance per unit length is 2.3 pF/cm and the bandwidth-length product is $\Delta f_{RC}L = 2.5 \text{ GHz} \cdot \text{cm}$. The trade-off is between large G/W to reduce capacitance and small G/W to reduce drive voltage and electrode resistance. The ultimate speed of lumped electrode devices is limited by the electric signal transit time, with a bandwidth-length product of $2.2 \text{ GHz} \cdot \text{cm}$. The way to achieve higher speed modulation is to use traveling-wave electrodes.

The traveling-wave electrode is a miniature transmission line. Ideally, the impedance of this coplanar line is matched to the electrical drive line and is terminated in its characteristic impedance. In this case, the modulator bandwidth is determined by the difference in velocity between the optical and electrical signals (velocity mismatch or walk-off), and any electrical propagation loss. Because of competing requirements between a small gap to reduce drive voltage and a wide electrode width to reduce RF losses, as well as reflections at any impedance transition, subtle trade-offs must be considered in designing traveling-wave devices.

Lithium niobate MZ modulators operating out to 35 GHz at $\lambda = 1.55 \mu\text{m}$ are commercially available, with $V_\pi = 10 \text{ V}$, with $<5 \text{ dB}$ insertion loss and $>20 \text{ dB}$ extinction ratio.⁶² To operate near quadrature, which is the linear modulation point, a bias voltage of $\sim 4 \text{ V}$ is required. Direct coupling from a laser or polarization-maintaining fiber is required, since these modulators operate on only one polarization.

Insertion Loss

Modulator insertion loss can be due to Fresnel reflection at the lithium niobate-air interfaces, which can be reduced by antireflection coatings or index matching (which only helps, but does not eliminate this loss, because of the very high refractive index of lithium niobate). The other cause of insertion loss is mode mismatch. To match the spatial profile of the fiber mode, a deep and buried waveguide must be diffused. Typically, the waveguide will be $9 \mu\text{m}$ wide and $5 \mu\text{m}$ deep. While the in-plane mode can be gaussian and can match well to the fiber mode, the out-of-plane mode is asymmetric, and its depth must be carefully optimized. In an optimized modulator, the coupling loss per face is about 0.35 dB and the propagation loss is about 0.3 dB/cm. This result includes a residual index-matched Fresnel loss of 0.12 dB.

Misalignment can also cause insertion loss. An offset of $2 \mu\text{m}$ typically increases the coupling loss by 0.25 dB. The angular misalignment must be maintained below 0.5° in order to keep the excess loss below 0.25 dB.⁶³

Propagation loss comes about from absorption, metallic overlay, scattering from the volume or surface, bend loss, and excess loss in the Y-branches. Absorption loss at 1.3 and $1.55 \mu\text{m}$ wavelengths appears to be $<0.1 \text{ dB/cm}$. Bend loss can be large, unless any curvature of guides is small. The attenuation coefficient in a bend has the form $\alpha = C_1 \exp(-C_2 R)$, where $C_1 = 15 \text{ mm}^{-1}$ and $C_2 = 0.4 \text{ mm}^{-1}$ in titanium indiffused lithium niobate, at wavelengths around 1.3 to $1.5 \mu\text{m}$. This means that a 5-mm-long section of constant radius 20 mm will introduce only 0.1 dB of excess loss.⁶³

A final source of loss in Y-branches is excess radiation introduced by sharp transitions at the Y junction. These junctions must be fabricated carefully to avoid such losses, since tolerances on waveguide roughness are critically small.

Polarization Independence

As previously shown, if the light is incident along the z axis and the field is along the y axis, then light polarized along x and y experience the same phase shift, but with opposite signs. This requires an x -cut crystal, with an in-plane field along y , which provides polarization-independent interferometric modulation at the sacrifice of somewhat higher half-wave voltage (e.g., 17 V).⁶⁴ Because of the difficulty of achieving exactly reproducible lengths in the two arms of the Y-branch interferometer, it was found useful to do a postfabrication phase correction using laser ablation.

Photorefractivity and Optical Damage

Lithium niobate exhibits *photorefractivity*, also called *optical damage*, when it is a nuisance. This phenomenon is a change in refractive index as a result of photoconduction originating in weak absorption by deep traps and a subsequent redistribution of charges within the lithium niobate. Because the photoconductive crystal is electro-optic, the change in electric field resulting from charge motion shows up as a change in refractive index, altering the phase shift as well as the waveguiding properties. While photorefractivity seriously limits the performance of lithium niobate modulators at shorter wavelengths (even at 850 nm),⁶⁵ it is not a serious concern at 1.3 μm and 1.55 μm .

However, partial screening by photocarriers may cause a drift in the required bias voltage of modulators, and systems designers may need to be sensitive to this.

Multimode Interferometric Mach-Zehnder Modulator

An alternative geometry has been developed for the waveguide Mach-Zehnder interferometer that replaces the Y branch by a multimode interferometric (MMI) power splitter, which is shown in Fig. 33. The MMI modulator works on the principle that there are distances at which light traveling down a step-index waveguide self-images. These self-imaging planes are analogous to Talbot planes for plane waves. The process of self-imaging is indicated in Fig. 33a.⁶⁶ It can be seen that if light is incident in a spatial distribution localized to only one-half of the guide, at a distance of $1/2 (3L_\pi)$ the power will be split equally into two spatial distributions. This enables a simple power splitter, as shown in Fig 33b.⁶⁶

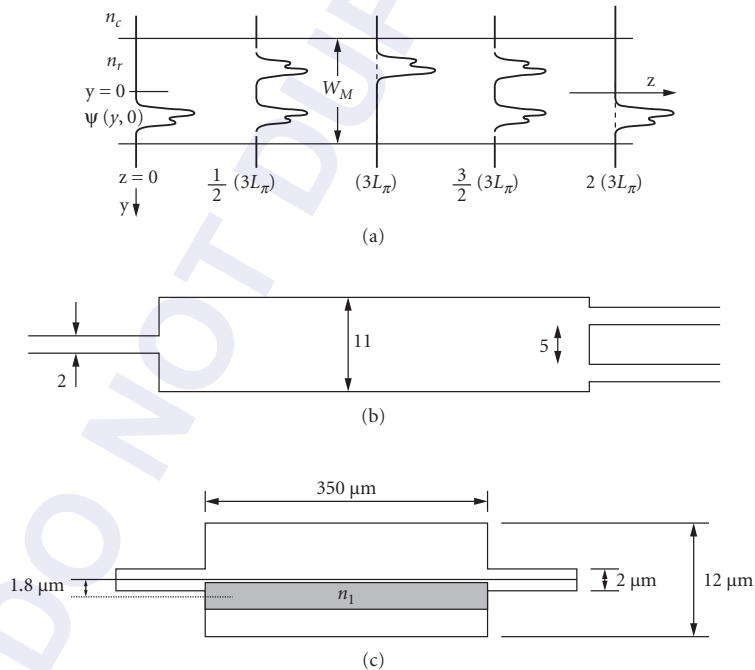


FIGURE 33 Multimode interferometric (MMI) devices as observed from above: (a) multimode waveguide showing the input field $\Psi(y,0)$, a mirrored single image at $(3L_\pi)$, a direct single image at $2(3L_\pi)$, and two-fold images at $1/2 (3L_\pi)$ and $3/2 (3L_\pi)$;⁶⁶ (b) power splitter, achieved at $1/2 (3L_\pi)$ with typical dimensions as shown in micrometers;⁶⁶ and (c) interference modulator, achieved at $(3L_\pi)$, with voltage applied by the darkened electrode; ground plane is under the substrate.⁶⁷

In a step-index multimode guide, because the modes are equally spaced, the beat length L_π is given by

$$L_\pi = \pi/(\beta_0 - \beta_1) \approx 4n_{co}W_e^2/3\lambda_0 \quad (57)$$

where W_e is the effective width of the guide, including Goos-Hanschen shifts. For $W_e = 11 \mu\text{m}$, this gives $L_\pi \sim 0.5 \text{ mm}$ in semiconductors, so a power splitter could be less than a millimeter long.

Ultracompact modulators can be constructed by the method shown in Fig. 33c,⁶⁷ with an electrode applied to only part of the structure, which will nonuniformly change the phases and destroy the coherent mixing that self-images the light to the output waveguide. Different electrode structures and electro-optic media have been demonstrated, including prism electrodes in electro-optic polymers.⁶⁸

Electro-Optic Polymer Modulators

For a number of years researchers have been working to make polymer modulators a reality. Electro-optic polymers contain nonlinear chromophores capable of providing large electro-optic effects at transmission wavelengths of 1.3 and 1.55 μm , with electro-optic coefficients of $\sim 70 \text{ pm/V}$, with the potential of $\sim 100 \text{ pm/V}$, which is 3 to 4 times greater than of lithium niobate. The technology typically involves spin coating of the polymer onto a silicon substrate, poling it, and applying electrodes. Electro-optic polymers are compatible with semiconductor fabrication methods and have low dielectric constants, so they can be easily velocity matched when applying traveling-wave electrodes for high-speed switching.

However, electro-optic polymers have had serious reliability issues. Gradual relaxation of the electro-optic coefficient is due to the slow misalignment of the chromophores, which form antiparallel pairs and effectively cancel out the electro-optic effect. This problem has been largely resolved using chemical anchoring schemes and/or poling in inert environments. There continues to be a problem with gradual changes in refractive index; however, some device designs are relatively robust to refractive index changes. The MMI is proposed to be such a structure, with simple design tolerance toward fluctuation of refractive indices.

For these modulators, highly active chromophores are dispersed at high concentrations in host polymers. A traditional commercially available chromophore is Disperse Red 1 (DR1), while much higher activity chromophores have been researched and are rapidly becoming commercially available. The polymer host is amorphous polycarbonate or polysulfone. The chromophore-doped polymer is dissolved in cyclohexanone (CHN) and tetrahydrofuran and then spin coated on the substrate. Cladding is typically spun on using commercially available acrylate. Relief structures, such as ridge waveguides, are plasma-etched into the films. After electrodes are deposited, the chromophore-doped polymer is made electro-optic by poling. In this process the sample is heated (e.g., at 130°C for 15 min) with a constant voltage applied (e.g., 600 V) between the bottom and the top electrodes. The temperature is dropped back to room temperature with the electric field still applied to ensure that the chromophores are oriented to achieve noncentrosymmetric alignment.⁶⁹

13.12 ELECTROABSORPTION MODULATORS

Semiconductors exhibit field-dependent absorption and refractive index. Such modulators can be integrated directly on the same chip as the laser, or placed external to the laser chip. External modulators may be butt-coupled to the laser, coupled by means of a microlens or by means of a fiber pigtail.

Electroabsorption

The electric field dependence of the absorption near the band edge of a semiconductor is called *electroabsorption*, and is particularly strong in quantum wells, where it is often called the *quantum-confined*

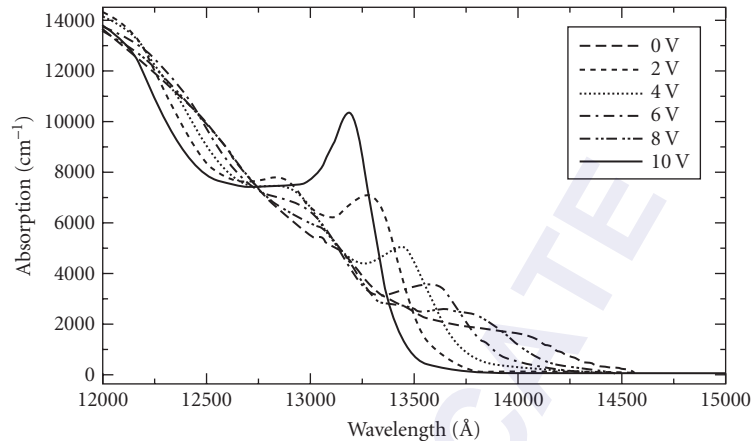


FIGURE 34 Spectrum of quantum-confined Stark effect (QCSE) in InAsP/InP-strained MQWs. The absorption changes with applied field.⁷⁰

Stark Effect (QCSE) An example of the wavelength dependence of the QCSE is shown in Fig. 34. The absorption spectrum of QWs exhibits a peak at the *exciton resonance*: when a field is applied, the exciton resonance moves to longer wavelengths, becomes weaker, and broadens. This means that on the long-wavelength side, the absorption increases with field, as the exciton resonance moves to longer wavelengths. At wavelengths closer to the exciton resonance, the absorption will first increase with field, then plateau, and finally decrease, as the field continues to grow. At wavelengths shorter than the zero-field exciton resonance, the absorption will decrease with increasing field, as the resonance moves to longer wavelengths.

While electroabsorption in QWs is much larger than in bulk, due to the sharpness of the excitonic-enhanced absorption edge, the useful absorption change must be multiplied by the filling factor of the QW in the waveguide, which reduces its effective magnitude. Under some conditions, electroabsorption near the band edge in bulk semiconductors (typically called the *Franz-Keldysh effect*) may also be useful in electroabsorption modulators.

Waveguide Modulators

When light traverses a length of QW material, the transmission will be a function of applied voltage. An electroabsorption modulator consists of a length of waveguide containing QWs. The waveguide is necessary to confine the light to the QW region so that it does not diffract away. Thus, low-refractive-index layers must surround the layer containing the QWs. Discrete electroabsorption modulators are typically made by using geometries similar to those of edge-emitting lasers (Fig. 1), but without mirrors. They are cleaved, antireflection coated, and then butt-coupled to the laser chip. They are operated by a reverse bias, rather than the forward bias of a laser. Or the modulator can be integrated on the laser chip, with the electroabsorption modulation region following a DFB or DBR laser in the optical train, as shown in Fig. 35. This figure shows the simplest electroabsorption modulator, with the same MQW composition as the DFB laser. This ridge waveguide device has been demonstrated with a 3-dB bandwidth of 30 GHz. The on-off contrast ratio is 12.5 dB for a 3-V drive voltage in a 90- μm -long modulator.⁷¹ The use of the same QWs is possible by setting the grating that determines the laser wavelength to well below the exciton resonance. Because of the inherently wide-gain spectrum exhibited by strained layer MQWs, this detuning is possible for the laser and still allows it to operate in the optimal wavelength region for the electroabsorption modulator.

Other integrated electroabsorption modulators use a QW composition in the electroabsorption region that is different from that of the laser medium. Techniques for integration are discussed later.

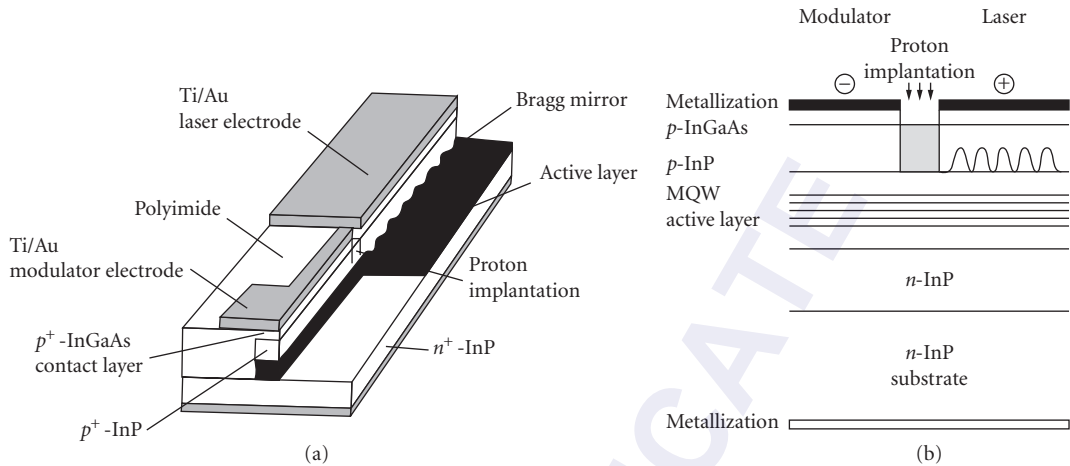


FIGURE 35 Electroabsorption modulator: (a) geometry for a channel electroabsorption modulator (foreground) integrated on the same chip with a DFB laser (background, under the Bragg mirror) and (b) side view, showing how the same MQW active layer can be used under forward bias with a grating to provide a DFB laser, and in a separate region under reverse bias for modulation, with the two regions electrically separated by proton implantation.⁵⁰

Intensity Modulation by Electroabsorption

In an electroabsorption waveguide modulator of length L , where the absorption is a function of applied field E , the transmission is a function of field: $T(E) = \exp[-\alpha(E)L]$, where α is the absorption per unit length, averaging the QW absorption over the entire waveguide. (That is, α is the QW absorption multiplied by the filling factor of the QW in the waveguide.) Performance is usually characterized by two quantities: *insertion loss* (throughput at high transmission) and *contrast ratio* (ratio of high transmission to low transmission). Assume that the loss in the QW, initially at low value α_0 , increases by $\delta\alpha$. The *contrast ratio* is given by: $CR \equiv T_{\text{high}}/T_{\text{low}} = \exp(\delta\alpha L)$. The *insertion loss* is given by $A \equiv 1 - T_{\text{high}} = 1 - \exp(-\alpha_0 L) \approx \alpha_0 L$. A long path length L means a high contrast ratio but also a large insertion loss and large capacitance, which results in a slower speed. Choosing the most practical length for any given application requires trading off the contrast ratio against insertion loss and speed.

To keep a moderate insertion loss, waveguide lengths should be chosen so that $L \approx 1/\alpha_0$, which sets the contrast ratio as $CR = \exp(\delta\alpha/\alpha_0)$. The contrast ratio depends on the ratio of the change in absorption to the absorption in the low-loss state; this fact is used to design the QW composition and dimensions relative to the wavelength of operation. In general, the contrast ratio improves farther from the band edge, but the maximum absorption is smaller there, so the modulator must be longer, which increases its capacitance, decreases its speed, and increases its loss. Contrast ratios may reach 10/1 or more with <2 V applied for optimized electroabsorption modulators. In a waveguide, the contrast ratio does not depend on the filling factor of the QW in the waveguide, but the required length L does. Since high-speed modulators require small capacitance and small length, the filling factor should be as high as possible.

Waveguide modulators are used at wavelengths where the absorption is not too large, well below the band edge. In this wavelength region, electroabsorption at a fixed wavelength can be modeled by a pure quadratic dependence on field. Thus $\alpha(E) \approx \alpha_0 + \alpha_2 E^2$, where α_2 will typically depend on the wavelength, the QW and barrier dimensions and composition, and the waveguide filling factor. Intimately connected with this change in absorption is a change in refractive index with a similar field dependence: $\delta n(E) \approx n_2 E^2$, where n_2 is also strongly dependent on wavelength. Both electroabsorption and electrorefraction are about an order of magnitude larger in QWs than in bulk material. Specific

numerical values depend on the detailed design, but typical values are on the order of $\alpha_p \sim 100 \text{ cm}^{-1}$, $\delta\alpha \sim 1000 \text{ cm}^{-1}$, $L \sim 200 \text{ }\mu\text{m}$ for 2 V applied across an i region 2.5- μm thick, for a field of $\sim 10 \text{ kV/cm}$. This means $\alpha_2 \sim 2 \times 10^{-5} \text{ cm/V}^2$. Also, $n_2 \sim 2 \times 10^{-11} \text{ V}^{-2}$.

Applying a Field in a Semiconductor

The electric field is usually applied by reverse biasing a *pin* junction. The electric field is supported by the semiconductor depletion region that exists within a *pin* junction, or at a metal-semiconductor junction (Schottky barrier). Charge carrier depletion in the n and p regions may play a role in determining the electric field across thin intrinsic regions. Taking this into account while assuming an undoped i region, the electric field across the i region of an ideally abrupt *pin* junction, when the undoped layer d_i is sufficiently large, can be expressed as⁷²

$$E = \frac{V_{\text{tot}}}{d_i} \left[1 - \frac{\epsilon V_{\text{tot}} (1 + N_d/N_a)}{2eN_d d_i^2} \right] \quad (58)$$

where N_d is the (donor) doping density in the n region, N_a is the (acceptor) doping density in the p region, e is the elementary charge, ϵ is the dielectric constant, d_i is the thickness of the intrinsic region, and V_{tot} is the sum of the applied and built-in field (defined positive for reverse bias). This assumes the n and p regions are highly doped and the i region is undoped, so that most of the voltage is dropped across the i region. To lowest order, this is just the field across a capacitor of thickness d_i . For a typical applied voltage of 5 V and $d_i = 0.25 \text{ }\mu\text{m}$, with $N_a = 10^{18} \text{ cm}^{-3}$ and $N_d = 10^{18} \text{ cm}^{-3}$, $E = 2 \times 10^5 \text{ V/cm}$. How much this will change the absorption and refractive index depends on wavelength and, of course, material design.

Operating Characteristics

In addition to contrast ratio, insertion loss, and required voltage, performance of electroabsorption modulators depends on speed, chirp, polarization dependence, optical power-handling capabilities, and linearity. These factors all depend on the wavelength of operation, the materials, the presence of strain, the QW and waveguide geometry, and the device design. Extensive trade-offs must be made to achieve the best possible operation for a given application. Modulators will differ, depending on the laser and the proposed applications.

Chirp Because a change in refractive index usually occurs during any absorption change $\delta\alpha$, electroabsorption modulators, in general, exhibit *chirp* (frequency broadening due to the time-varying refractive index, also observed in modulated lasers), which can seriously limit their usefulness. As with semiconductor lasers, the figure of merit is $\beta = k_o \delta n / \delta\alpha$. Unlike lasers, however, there are particular wavelengths of sizable absorption change at which $\delta n = 0$. Studies have shown that these nulls in index change can be positioned where $\delta\alpha$ is large by using coupled quantum wells (CQWs).⁷⁴ These structures provide two, three, or more wells so closely spaced that the electron wave functions overlap between them. If desired, several sets of these CQWs may be used in a single waveguide, if they are separated by large enough barriers that they do not interact. Chirp-free design is an important aspect of electroabsorption modulators.

On the other hand, since the chirp can be controlled in electroabsorption modulators, there are conditions under which it is advantageous to provide a negative chirp to cancel out the positive chirp introduced by fibers. This allows 1.55- μm laser pulses to travel down normally dispersive fiber (with zero material dispersion at 1.3 μm) without the pulses unduly spreading.

Polarization Dependence In general, the QCSE is strongly polarization dependent, although there may be specific wavelengths at which TE and TM polarized light experience the same values of

electroabsorption (and/or electrorefraction). It turns out that polarization-independent modulation is more readily achieved by using strained QWs. In addition, the contrast ratio of electroabsorption change at long wavelengths can be improved by using strained QWs.

Optical Power Dependence During the process of electroabsorption, the modulators can absorb some of the incident light, creating electron-hole pairs. If these pairs remain in the QWs, at high optical powers they will introduce a free carrier plasma field that can screen the exciton resonance. This broadens the absorption spectrum and reduces the contrast ratio. In some cases, electroabsorption modulators operating at the band edge of bulk semiconductors (the *Franz-Keldysh effect*) may be able to operate with higher laser power. A common approach is to use shallow QWs, so that the electrons and holes may escape easily.

Even when the electron-hole pairs created by absorption escape the QWs, they will move across the junction to screen the applied field. This will tend to reduce the applied field, so the performance will depend on the magnitude of absorbed light. Photogenerated carriers must also be removed, or they will slow down the modulator's response time. Carriers may be removed by leakage currents in the electrodes or by recombination.

Built-in Bias Because *pn* junctions have built-in fields, even at zero applied voltage, electroabsorption modulators have a prebias. Some applications use a small forward bias to achieve even larger modulation depths. However, the large forward current resulting from the forward bias limits the usefulness of this approach. There are, at present, some research approaches to remove the internal fields using an internal strain-induced piezoelectric effect to offset the *pn* junction intrinsic field.

Commercial Discrete Electroabsorption Modulators Discrete electroabsorption modulators can be fast: 30 GHz or more. They can be designed either for low chirp (α can be < 0.5), or low polarization-dependent loss (< 0.5), but not both. Typically they may operate with voltages from -4 to $+1$ V with 20-dB extinction ratio.⁷⁵ They are limited in optical power to 20 mW and require temperature-stabilization with a thermoelectric cooler. Their greatest disadvantage, however, is their large insertion loss (typically 7.5 to 10 dB), because mode profiles from single-mode fibers do not match the electroabsorptive semiconductor waveguides.

One approach to reduce the insertion loss is to use the electroabsorption modulator in reflection (R-EAM). Devices have been reported with 4.5 dB typical insertion loss and 0.5 dB polarization dependent loss. These devices have 11 dB modulation depth with a voltage range from 0 to -3 V, and 14 GHz bandwidth. Other devices operate at 35 GHz with a 1 dB polarization dependent loss and 10 dB contrast ratio.⁷⁶

Integrating the Modulator

Stripe-geometry modulators can be cleaved from a wafer, antireflection coated, and butt-coupled to either a laser or a fiber pigtail. Typical insertion losses may be ~ 10 dB. Or, the modulator may be monolithically integrated with the laser. A portion of the same epitaxial layer grown for the laser's active region can be used as an electroabsorption modulator by providing a separate contact and applying a reverse bias. When such a modulator is placed inside the laser cavity, a multielement laser results that can have interesting switching properties, including wavelength tunability. When the electroabsorption modulator is placed outside the laser cavity, it is necessary to operate an electroabsorption modulator with a higher energy bandgap than the laser medium. Otherwise, the incident light will be absorbed, creating free electron-hole pairs that will move to screen the applied field and ruin the modulator.

Etch and Regrowth Typically, a first set of epitaxial layers is grown everywhere, which includes the laser structure up through the QW layer. Then the QW layer is etched away from the regions where it is not needed. The structure is then overgrown everywhere with the same upper cladding layers.

This typically results in a bulk electroabsorption modulator, consisting of laser-cladding material. A more complex fabrication process might mask the laser region during the regrowth process and grow a different QW composition that would provide an integrated butt-coupled modulator for the DFB (or DBR) QW laser.

Vertical Coupling between Layers This approach makes it possible to use a QW modulator as well as a QW laser, with a different QW composition in each. Two sets of QWs can be grown one on top of the other and the structures can be designed so that light couples vertically from one layer to the other, using, for example, grating-assisted coupling. This may involve photolithographically defining a grating followed by a regrowth of cladding layers, depending on the design.

Selective Area Epitaxy Growth on a patterned substrate allows the width of the QWs to be varied across the wafer during a single growth. The substrate is usually coated with a SiO_2 mask in which slots are opened. Under a precise set of growth conditions no growth takes place on top of the dielectric, but surface migration of the group III species (indium) can take place for some distance across the mask to the nearest opening. The growth rate in the opened area depends on the width of the opening and the patterning on the mask. Another approach is epitaxial growth on faceted mesas, making use of the different surface diffusion lengths of deposited atomic species on different crystal facets.

Well and Barrier Intermixing The bandgap of a QW structure can be modified after growth by intermixing the well and barrier materials to form an alloy. This causes a rounding of the initially square QW bandgap profile and, in general, results in an increase of the bandgap energy. This provides a way to fabricate lasers and bandgap-shifted QCSE modulators using only one epitaxial step. Intermixing is greatly enhanced by the presence of impurities or defects in the vicinity of the QW interfaces. Then the bandgap is modified using impurity-induced disordering, laser beam-induced disordering, impurity-free vacancy diffusion, or ion implantation-enhanced interdiffusion. The challenge is to ensure that the electrical quality of the *pin* junction remains after interdiffusion; sometimes regrowth of a top p layer helps.⁷³

Electroabsorption Modulators Integrated with Lasers

An electroabsorption modulator (EAM) integrated with a DFB or DBR laser is a good choice as a transmitter device in systems because of its compactness and ease of operation. Selective epitaxy enables separate optimization of the laser and the modulator. The laser may have either a DFB or DBR geometry and the EAM is a reverse biased *pin* structure. The two devices may be fabricated in a butt-coupled geometry, or separated and a waveguide grown to connect them.

Electroabsorption-modulated lasers are now commercially available from such companies as Oki, Cyoptics, Mitsubishi, JDS Uniphase, SVEDICE in Sweden, among others. They typically provide modulation to 40 Gb/s, with parameters such as 20-dB extinction ratio, 6-mW optical power, and 10-MHz spectral width with 35-dB sidemode suppression. These devices require thermoelectric coolers that typically dissipate 2 W. The specific designs are usually company-proprietary, but enough operational data is provided that they can be operated in a relatively straightforward manner. The great advantage of the electroabsorption modulator is the severe reduction in chirp, compared to modulating the laser current.

An alternative to selective epitaxy is the use of a localized quantum well intermixing process that has been shown to increase the bandgap. Intermixing can occur by impurity-induced disordering, ion-implantation enhanced intermixing, and impurity-free vacancy diffusion. The latter method involves deposition of a dielectric capping material and subsequent thermal annealing. In GaAs–AlGaAs QW materials, SiO_x induces outdiffusion of Ga during annealing, causing vacancies that enhance the intermixing of Ga and Al in the QWs. The bandgap becomes larger in the QW because of partial disordering of the two materials. In the InGaAs(P)–InP QW materials, SiN is used as a cap.

13.13 ELECTRO-OPTIC AND ELECTROREFRACTIVE MODULATORS

Some semiconductor modulators are based on phase modulation that is converted to amplitude modulation by using a Mach-Zehnder interferometer, in the same manner as discussed in Sec. 13.11. Such modulators can be integrated on the same substrate as the laser, but do not have the chirp issues that electroabsorption modulators exhibit.

Electro-Optic Effect in Semiconductors

The III-V semiconductors are electro-optic. Although not initially anisotropic, they become so when an electric field is applied. Referring to the discussion of the electro-optic effect in Sec. 13.11 for definitions, the GaAs electro-optic coefficients have only one nonzero term: $r_{41} = r_{52} = r_{63} = 1.4 \times 10^{-12}$ m/V. Crystals are typically grown on the (001) face, with the z axis normal to the surface, and the field is usually applied along z . The only electro-optically induced index change will be $\Delta(1/n^2)_4 = r_{41}E_z$, which introduces nondiagonal terms in the (x, y, z) coordinate system. Diagonalizing the matrix results in new axes (x', y', z) that are at 45° to the (x, y, z) crystal axes and new values of the inverse squared refractive index along these axes: $1/n_{x'}^2 = 1/n_o^2 + r_{41}E_z$ and $1/n_{y'}^2 = 1/n_o^2 - r_{41}E_z$.

Differentiation for small refractive change provides refractive index changes for light polarized at 45° to crystal axes:

$$n_{x'} = n_o + \frac{n_o^3}{2} r_{41} E_z \quad \text{and} \quad n_{y'} = n_o - \frac{n_o^3}{2} r_{41} E_z \quad (59)$$

The direction of these new optic axes (45° to the crystal axes) turns out to be in the direction that the zincblende material cleaves. Thus, TE-polarized light traveling down a waveguide normal to a cleave experiences the index change shown in Eq. (59). Depending on whether light goes along x' or y' , the index will increase or decrease.⁷⁷ Light polarized along z will not see any index change.

With an electro-optic coefficient of $r_{41} = 1.4 \times 10^{-10}$ cm/V, in a field of 10 kV/cm (2 V across 2 μm), and since $n_o = 3.3$, the index change for the TE polarization in GaAs will be 2.5×10^{-5} . The index change in InP-based materials is comparable. The phase shift in a sample of length L is $\Delta n k L$. At 1- μm wavelength, this will require a sample of length 1 cm to achieve a π phase shift, so that the voltage-length product for electro-optic GaAs (or other semiconductor) will be ~ 20 V-mm. Practical devices require larger refractive index changes, which can be achieved by using the electrorefractive effect in QWs and choosing the exciton resonance at a shorter wavelength than that of the light to be modulated.

Electrorefraction in Semiconductors

Near the band edge in semiconductors, the change in refractive index with applied field can be particularly large, especially in QWs, and is termed *electrorefraction*, or the *electrorefractive effect*. Electrorefraction is calculated from the spectrum of electroabsorption using the Kramers-Kronig relations. The existence of electroabsorption means there will inevitably be electrorefraction at wavelengths below the band edge. Electrorefraction may be larger than the electro-optic effect and may significantly reduce the length and drive voltages required for phase modulation in semiconductor waveguides. The voltage-length product depends on how close to the absorption resonance the modulator is operating. It also depends on device design. As with electroabsorption modulators, the field is usually applied across a *pin* junction. Some reported π voltage-length products are 2.3 V-mm in GaAs/AlGaAs QWs (at 25-V bias), 1.8 V-mm in InGaAs/InAlAs QWs, and the same in GaAs/AlGaAs double heterostructures.⁷⁸ These voltage-length products depend on wavelength detuning from the exciton resonance and therefore on insertion and electroabsorption losses. The larger the voltage-length product, the greater the loss.

Typical Performance Electrorefraction is polarization dependent, because the QCSE is polarization dependent. In addition, the TE polarization experiences the electro-optic effect, which may add to or subtract from the electrorefractive effect, depending on the crystal orientation. Typically, Δn for TE polarization (in an orientation that sums these effects) will be 8×10^{-4} at 82 kV/cm (7 V across a waveguide with an *i* layer 0.85- μm thick). Of this, the contribution from the electro-optic effect is 2×10^{-4} . Thus, electrorefraction is about 4 times larger than the electro-optic effect and the voltage-length product will be reduced by a factor of 5. Of course, this ratio depends on the field, since the electro-optic effect is linear in field and electrorefraction is quadratic in the field. This ratio also depends on wavelength; electrorefraction can be larger at wavelengths closer to the exciton resonance, but the residual losses go up. The TE polarization of electrorefractive modulators integrates well with a TE-polarized laser, and they can be grown on the same substrate.

The TM polarization, which experiences electrorefraction alone, will be 5×10^{-4} at the same field, slightly smaller than the TE electrorefraction, because QCSE is smaller for TM than TE polarization.⁷⁹

Polarization-independent electrorefraction modulators have been demonstrated using suitably strained quantum wells.

Advanced QW Concepts Compressive strain increases electrorefraction, as it does QCSE. Measurements at the same 82 kV/cm show an increase from 2.5×10^{-4} to 7.5×10^{-4} by increasing compressive strain.⁸⁰ Strained QWs also make it possible to achieve polarization-independent electrorefractive modulators (although when integrated with a semiconductor laser, which typically has a well-defined polarization, this should not be a necessity).

Advanced QW designs have the potential to increase the refractive index change below the exciton resonance. One example analyzes asymmetric coupled QWs and finds more than 10 times enhancement in Δn below the band edge, at least at small biases. However, when fabricated and incorporated into Mach-Zehnder modulators, the complex three-well structure lowered V_π by only a factor of 3, attributed to the growth challenges of these structures.⁸¹

Nipi Modulators One way to obtain a particularly low voltage-length product is to place MQWs in a hetero-*nipi* waveguide. These structures incorporate multiple *pin* junctions (alternating *n-i-p-i-n-i-p*) and include QWs in each *i* layer. Selective contacts to each electrode are required, which limits how fast the modulator can be switched. A voltage-length product of 0.8 V·mm was observed at a wavelength 115 meV below the exciton resonance. The lowest voltage InGaAs modulator had $V_\pi = 0.5$ V, at speeds up to 110 MHz. Faster speeds require shorter devices and higher voltages.⁷⁸

Band-Filling Modulators When one operates sufficiently far from the band edge so that the absorption is not large, then the electrorefractive effect is only 2 to 3 times larger than the bulk electro-optic effect. This is because oscillations in the change in absorption with wavelength tend to cancel out their contributions to the change in refractive index at long wavelengths. By contrast, during band filling, the long-wavelength refractive index change is much larger, because band filling decreases the absorption at all wavelengths. However, because band filling relies on real carriers, it lasts as long as the carriers do, and it is important to find ways to remove these carriers to achieve high-speed operation.

Voltage-controlled transfer of electrons into and out of QWs (BRAQWET modulator) can yield large electrorefraction by band filling. The refractive index change at 1.55 μm can be as large as $\Delta n = 0.02$ for 6 V. One structure consists of 12 repeating elements,⁸² with the single QW replaced by three closely spaced strongly coupled QWs, demonstrating $V_\pi L = 3.2$ V·mm with negligible loss. Researchers have been looking at these QW devices, and others driven by heterostructure buried transistors (HBTs) that drive current in and out of the QWs. None seems to have reached practicality by this time, however.

An important direction for bandfilling nonlinearities is the holy grail of all-silicon photonics. In silicon the refractive index changes because of electron-hole plasma injection. One approach is to integrate a SiGe HBT together with silicon waveguide. Including a distributed Bragg reflector can convert the index change into a transmission change. The use of SiGe quantum wells enables the quantum-confined Stark effect to be used, leading to larger electrorefraction than from band-filling alone.

Semiconductor Interferometric Modulators

The issues for Mach-Zehnder modulators fabricated in semiconductors are similar to those for modulators in lithium niobate, but the design and fabrication processes in semiconductors are by no means as well developed. Fabrication tolerances, polarization dependence, interaction with lasers, and operation at high optical input powers are just some of the issues that need to be addressed. While the power splitters in the interferometer can be Y branches, fabricated by etching to form ridge waveguides, they are more usually formed by *multimode interference* (MMI) power splitters, discussed previously.

One example reports a Mach-Zehnder interferometer at 1.55 μm in InGaAlAs QWs with InAlAs barriers.⁸⁵ A polarization-independent extinction ratio of 30 dB was reported, over a 20-nm wavelength range without degradation at input powers of 18 dBm (63 mW). The interferometer phase-shifting region was 1000- μm long, and each MMI was 200- μm long. The insertion loss of 13 dB was due to the mismatch between the mode of the single-mode optical fiber and of the semiconductor waveguide, which was 2- μm wide and 3.5- μm high. Various semiconductor structures to convert spot size should bring this coupling loss down.

13.14 PIN DIODES

The detectors used for fiber optic communication systems are usually *pin photodiodes*. In high-sensitivity applications, such as long-distance systems operating at 1.55- μm wavelength, *avalanche photodiodes* are sometimes used because they have internal gain. Occasionally, *metal-semiconductor-metal* (MSM) *photoconductive detectors* with interdigitated electrode geometry are used because of ease of fabrication and integration. For the highest-speed applications, *Schottky photodiodes* may be chosen. This section reviews properties of *pin* photodiodes. The next section outlines the other photodetectors.

The material of choice for these photodiodes depends on the wavelength at which they will be operated. The short-wavelength *pin* silicon photodiode is perfectly suited for GaAs wavelengths (850 nm); these inexpensive detectors are paired with GaAs/AlGaAs LEDs for low-cost data communication applications. They are also used in the shorter-wavelength plastic fiber applications at 650 nm. The longer-wavelength telecommunication systems at 1.3 and 1.55 μm require longer-wavelength detectors. These are typically *pin* diodes composed of lattice-matched ternary $\text{In}_{0.47}\text{Ga}_{0.53}\text{As}$ grown on InP. Silicon is an indirect bandgap semiconductor while InGaAs is a direct-band material; this affects each material's absorption properties and therefore its photodiode design. In particular, silicon photodiodes tend to be slower than those made of GaAs or InGaAs, because silicon intrinsic regions must be thicker. Carriers absorbed in the *n* region that diffuse into the *i* layer also slow silicon *pin* detectors down because diffusion lengths in silicon are much longer than in GaAs. Speeds are also determined by carrier mobilities, which are lower in silicon than in the III-V materials.

Previous volumes in this *Handbook* have outlined the concepts behind the photodetectors discussed here. Chapter 24 in Vol. II places *pin* photodetectors in context with other detectors, and gives specific characteristics of some commercially available detectors, allowing direct comparison of silicon, InGaAsP, and germanium detectors. Chapter 25 in Vol. II describes the principles by which the *pin* and the avalanche photodiodes operate. The properties of greatest interest to fiber communications will be repeated here. Chapter 26 in Vol. II concentrates on high-speed photodetectors and provides particularly useful information on their design for high-speed applications.

The *pin* junction consists of a thin, heavily doped *n* region (called n^+), a near-intrinsic *n* region (the *i* region), and a heavily doped *p* region (called p^+). An incident photon with energy greater than or equal to the bandgap of the semiconductor generates electron-hole pairs. In a well-designed photodiode, this generation takes place in the space-charge region of the *pn* junction. As a result of the electric field in this region, the electrons and holes separate and drift in opposite directions, causing current to flow in the external circuit. This current is monitored as a change in voltage across a load resistor. The *pin* photodiode is the workhorse of fiber communication systems.

Typical Geometry

Typically, the electric field is applied across the pn junction and photocarriers are collected across the diode. A typical geometry for a silicon photodiode is shown in Fig. 36a. A pn junction is formed by a thin p^+ diffusion into a lightly doped n^- or i layer (also called the i layer since it is almost *intrinsic*) through a window in a protective SiO_2 film. Then super minus or i layer i layer between the p^+ and n^+ regions supports a *space-charge region*, which, in the dark, is depleted of free carriers (this is sometimes called the *depletion region*) and supports the voltage drop that results from the pn junction. When light is absorbed in this space-charge region, the absorption process creates electron-hole pairs that separate in the electric field (field lines are shown in Fig. 36a), the electrons falling down the potential hill to the n region and the holes moving to the p region. This separation of charge produces a current in the external circuit, which is read out as a measure of the light level. Free carriers generated within a diffusion length of the junction may diffuse into the space charge region, adding to the measured current.

Long-wavelength detectors utilize n^- or i layers that are grown with a composition that will absorb efficiently in the wavelength region of interest. The ternary $\text{In}_{0.47}\text{Ga}_{0.53}\text{As}$ can be grown lattice-matched to InP and has a spectral response that is suitable for both the 1.3- and 1.55- μm wavelength regions. Thus, this ternary is usually the material of choice, rather than the more difficult to grow quaternary InGaAsP, although the latter provides more opportunity to tune the wavelength response. Figure 36b shows a typical geometry. Epitaxial growth provides lightly doped material on a heavily doped substrate. The InP buffer layer is grown to keep the dopants from diffusing into the lightly doped absorbing InGaAs layer. The required thin p region is formed by diffusion through a silicon nitride insulating window. Because InP is transparent to 1.3- and 1.55- μm light, the photodiode can be back-illuminated, which makes electrical contacting convenient. In some embodiments, a well is etched in the substrate and the fiber is glued in place just below the photosensitive region.

Carriers generated outside the space-charge region may enter into the junction by diffusion, and can introduce considerable delay time. In silicon, the diffusion length is as long as 1 cm, so photocarriers generated anywhere within the silicon photodiode can contribute to the photocurrent. Because the diffusion velocity is much slower than the transit time across the space-charge region, diffusion currents slow down silicon photodiodes. This is particularly true in pn diodes. Thus, high-speed applications typically use pin diodes with a narrower bandgap in the i layer, limiting absorption to this region.

To minimize diffusion from the p^+ entrance region, the junction should be formed very close to the surface. The space-charge region should be sufficiently thick so that most of the light will be absorbed

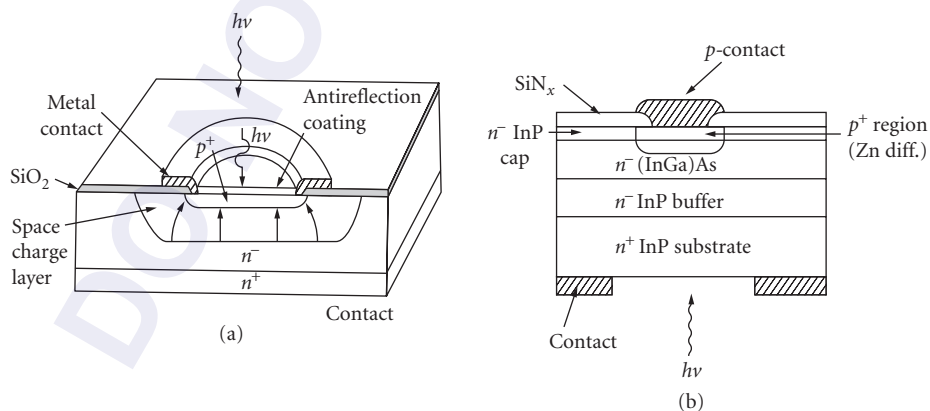


FIGURE 36 Geometry for pin photodiodes: (a) cut-away of silicon, illuminated from the top, showing the ring electrode and static electric field lines in the space-charge region and (b) cross-section of InGaAs/InP, illuminated from the bottom. The p^+ region is formed by diffusion. The low-doped n -layer is the i or nearly intrinsic layer.

there (thickness $\approx 1/\alpha$). With sufficient reverse bias, carriers will drift at their scattering-limited velocity. The space-charge layer must not be too thick, however, or transit-time will limit the frequency response. Neither should it be too thin, or excessive capacitance will result in a large RC time constant. The optimum compromise occurs when the modulation period is on the order of twice the transit time. For example, for a modulation frequency of 10 GHz, the optimum i layer thickness in silicon is about 5 μm . However, this is not enough thickness to absorb more than ~ 50 percent of the light at 850-nm wavelength. Thus, there is a trade-off between sensitivity and speed. If the wavelength drops to 980 nm, only 10 percent of the light is absorbed in a 10- μm thick i layer of silicon.

The doping must be sufficiently small so that the i region can support the voltage drop of the built-in voltage V_{bi} plus the applied voltage. When the doping density of the p^+ region is much higher than the doping density of the i layer (actually lightly doped with N_D donor density), the thickness of the space-charge layer is

$$W_s = \sqrt{\frac{2\epsilon_s(V_{\text{bi}} - V)}{e N_D}} \quad (60)$$

To achieve space-charge layer thickness $W_s = 10 \mu\text{m}$ in a silicon photodiode with 10 V applied requires an extremely pure i layer with doping density $N_D \approx 10^{14} \text{cm}^{-3}$. If the doping is not this low, the voltage drops more rapidly, and the field will not extend fully across the low-doped region i region.

GaAs photodiodes are faster than silicon, because their diffusion length in the highly doped n region is only $\sim 100 \mu\text{m}$. Also, their transit time across the i layer is faster than the transit time in silicon. The fastest diodes have transparent n^+ and p^+ regions, such as InGaAs/GaAs or InGaAs/InP photodiodes (also GaAs/AlGaAs photodiodes). These diodes have highly doped n and p regions with wider bandgap material than the i layer, so they contribute no photocurrent. The thickness of the i layer is chosen thin enough to achieve the desired speed (trading off transit time and capacitance), with a possible sacrifice of sensitivity.

Typically, light makes a single pass through the active layer. In silicon photodiodes, the light usually enters through the p contact at the surface of the diode (Fig. 36a). The top metal contact must have a window for light to enter (or be a transparent contact, such as indium tin oxide). The InGaAs photodiodes may receive light from the p side or the n side, because neither is absorbing. In addition, some back-illuminated devices use a double pass, reflecting off a mirrored top surface, to double the absorbing length. Some more advanced detectors, *resonant photodiodes*, use integrally grown Fabry-Perot cavities (using DBR mirrors, as in VCSELs) that resonantly reflect the light back and forth across the i region, enhancing the quantum efficiency. These are typically used only at the highest bandwidths (>20 GHz) or for wavelength division multiplexing (WDM) applications, where wavelength-selective photodetection is required. In addition, photodiodes designed for integration with other components are illuminated through a waveguide in the plane of the pn junction. The reader is directed to Chap. 26, Vol. II to obtain more information on these advanced geometries.

Sensitivity (Responsivity)

To operate a pin photodiode, it is sufficient to place a load resistor between ground and the n side and apply reverse voltage to the p side ($V < 0$). The photocurrent is monitored as a voltage drop across this load resistor. The photodiode current in the presence of an optical signal of power P_s is negative, with a magnitude given by

$$I = \eta_D \left(\frac{e}{h\nu} \right) P_s + I_D \quad (61)$$

where I_D is the magnitude of the (negative) current measured in the dark. The detector *quantum efficiency* η_D (electron-hole pairs detected per photon) is determined by the optical transmission of

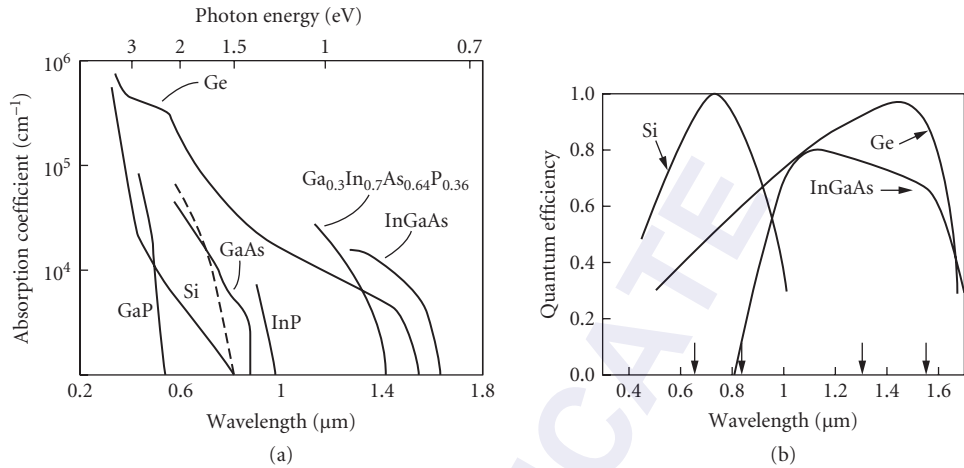


FIGURE 37 (a) Absorption coefficient as a function of wavelength for several semiconductors used in *pin* diode detectors and (b) spectral response of typical photodetectors.

the top electrode T , by the reflection R of light from the top surface of the photodiode (which can be reduced by adding an antireflective coating), and by how much light is absorbed in the region of thickness W , that is, at or within a diffusion length of the i layer (which depends on the absorption coefficient α). The detector quantum efficiency can be expressed as $\eta_D = (1 - R)T [1 - \exp(-\alpha W)]$.

The *sensitivity* (or *responsivity* \mathcal{R}) of a detector is the ratio of milliamperes of current out per milliwatt of light in. Thus, the responsivity is

$$\mathcal{R} = \frac{I_{PD}}{P_s} = \eta_D \frac{e}{h\nu} \quad (62)$$

For detection of a given wavelength, the photodiode material must be chosen with a bandgap sufficient to provide suitable sensitivity. The absorption spectra and quantum efficiency of candidate detector materials are shown in Fig. 37. Silicon photodiodes provide low-cost detectors for most data communications applications, with acceptable sensitivity at 850 nm (absorption coefficient $\sim 500 \text{ cm}^{-1}$). These detectors work well with the GaAs lasers and LEDs that are used in the inexpensive datacom systems and for short-distance or low-bandwidth local area network (LAN) applications. GaAs detectors are faster, both because their absorption can be larger and because their electron mobility is higher, but they are more expensive. Systems that require longer-wavelength InGaAsP/InP lasers typically use InGaAs photodiodes. Germanium has a larger dark current, so it is not usually employed for optical communications applications. Essentially all commercial photodetectors use bulk material, not quantum wells, as these are simpler, are less wavelength sensitive, and have comparable performance. Table 1 gives the sensitivity of typical detectors of interest in fiber communications, measured in units of amperes per watt, along with speed and relative dark current.

TABLE 1 Characteristics of Typical Photodiodes

	Wavelength (nm)	Sensitivity \mathcal{R} (A/W)	Speed t (ns)	Dark Current (Normalized Units)
Silicon	850	0.55	3	1
	650	0.4	3	
GaInAs	1300–1600	0.95	0.2	3
Ge (<i>pn</i>)	1550	0.9	3	66

Speed

Contributions to the speed of a *pin* diode come from the transit time across the space-charge region and from the RC time constant of the diode circuit in the presence of a load resistor R_L . Finally, in silicon there may be a contribution from the diffusion of carriers generated in undepleted regions.

The transit time across the space-charge region depends on its thickness. When the photodiode is properly operated, the space-charge region should extend fully across the i layer, which has a thickness W_i . Equation (59) gives the thickness of the space-charge region W_s , as long as it is less than the thickness of the i layer W_i . Define V_i as that voltage at which $W_s = W_i$; then $-V_i = W_i^2 e N_D / 2 \epsilon_s - V_{bi}$. For any voltage with magnitude larger than this, the space-charge width is essentially W_i (since the space charge extends a negligible distance into highly doped regions).

If the electric field across the space-charge region is high enough for the carriers to reach their saturation velocity v_s and high enough to fully deplete the i region, then the carrier transit time will be $\tau_i = W_i / v_s$. For $v_s = 10^7$ cm/s and $W_i = 4$ μ m, the transit time $\tau_i = 40$ ps. A finite transit time τ_i reduces the response at modulation frequency ω , such that:⁸⁴

$$\mathfrak{R}(\omega) = \mathfrak{R}_0 \frac{\sin(\omega\tau_i/2)}{\omega\tau_i/2} \quad (63)$$

Defining the 3-dB bandwidth as that modulation frequency at which the electrical power decreases by 50 percent, it can be shown that the transit-limited 3-dB bandwidth is $\delta\omega_i = 2.8/\tau_i = 2.8 v_s/W_i$. (Because electrical power is proportional to I^2 and \mathfrak{R}^2 , the half-power point is achieved when the current is reduced by $1/\sqrt{2}$.) There is a trade-off between sensitivity and transit time, since, for thin layers the quantum efficiency is $\eta_D \approx (1 - R)T\alpha W_i$. Thus, the quantum efficiency–bandwidth product is: $\eta_D \delta\omega_i \approx 2.8\alpha v_s (1 - R)T$. When the top electrode has no loss or reflection, this product depends only on the semiconductor's absorption loss and saturation velocity, not the geometry.

The speed of a *pin* photodiode is also limited by its capacitance, through the RC of the load resistor. Sandwiching an intrinsic layer, which is depleted of carriers, between conductive n and p regions causes a diode capacitance proportional to the detector area A : $C_D = \epsilon_s A/W_i$.

For a given load resistance, the smaller the area, the smaller the RC time constant, and the higher the speed. We will see also that the dark current I_s decreases as the detector area decreases. The detector should be as small as possible, as long as all the light from the fiber can be collected onto the detector. Multimode fibers easily butt-couple to detectors whose area matches the fiber core size. High-speed detectors compatible with single-mode fibers can be extremely small, which increases the alignment difficulty, so typically high-speed photodetectors are already pigtailed to single-mode fiber. A low load resistance may be needed to keep the RC time constant small, which may result in a small signal that needs amplification. Speeds in excess of 1 GHz are straightforward to achieve, and speeds of 50 GHz are not uncommon.

Thicker space-charge regions provide smaller capacitance, but too thick a space-charge region causes the speed to be limited by the carrier transit time. The bandwidth with a load resistor R_L is given by

$$\omega_{3\text{ dB}} = \frac{2.8}{\tau_i} + \frac{1}{R_L C} = \frac{2.8v_s}{W_i} + \frac{W_i}{\omega_s AR_L} \quad (64)$$

There is an optimum thickness W_i for high-speed operation. Any additional series resistance R_s or parasitic capacitance C_p must be added by using $R \rightarrow R_L + R_s$ and $C \rightarrow C + C_p$. The external connections to the photodetector can also limit speed. The gold-bonding wire may provide additional series inductance. It is important to realize that the photodiode is a high impedance load, with very high electrical reflection, so that an appropriate load resistor must be used. As pointed out in Chap. 26 in Vol. II, it is possible to integrate a matching load resistor inside the photodiode device, with a reduction in sensitivity of a factor of 2 (since half the photocurrent goes through the load resistor), but double the speed (since the RC time constant is halved). A second challenge is to build external bias circuits without high-frequency electrical resonances. Innovative design of the photodetector may integrate the necessary bias capacitor and load resistor, ensuring smooth electrical response.

Silicon photodetectors are inherently not as fast because of diffusion into the i layer of carriers absorbed in the highly doped p and n regions. Their photoresponse has a component with a slower response time governed by the carrier diffusion time: $T_D = W_D^2/2D$, where W_D is the width of the absorbing doped region, and D is the diffusion constant for whichever carrier is dominant (usually holes in the n region). For silicon, $D = 12 \text{ cm}^2/\text{s}$, so that when $W_D = 100 \text{ }\mu\text{m}$, $\tau_D = 400 \text{ ns}$.

Unitraveling-Carrier Photodiode

In conventional high-speed pin photodiodes, high optical power can saturate the current output. The cause is the space charge due to photogenerated holes, which introduces a variation of the electric field in the depletion region. This space charge effect can be relieved by using a unitraveling-carrier (UTC) design, shown in Fig. 38.

This $p^+n^-n^+$ structure utilizes only electrons as the active carriers, offering both high-speed and high-output simultaneously. The UTC structure separates the p^+ absorption layer from the depleted n^- collector layer, which is transparent because it has a wider bandgap. In the narrow bandgap p^+ photoabsorption layer, electrons and holes are generated by photoexcitation. Minority carrier electrons diffuse toward the collector, where the depletion field causes them to drift across the collector at high-speed. Because the electrical field vanishes in the p^+ absorbing layer, intense illumination has, in principle, no effect on the responsivity and dynamics of the photodetector. The speed of the UTC structure is determined only by the electron transit times in the photoabsorption and collector layers; the maximum available output voltage is much higher than that of the conventional structure. The price to pay is a low responsivity for top-illuminated devices, since the p^+ absorbing layer must be thin for fast collection of electrons.

These high-speed photodetectors have accelerated the progress of optical communication systems and measurement systems, reaching data rates as much as 40 Gb/s in long-haul transmission systems. These devices simplify and improve the performance of high-speed receivers, offering reliability and stability. The 3-dB bandwidth has been reported as high as 310 GHz; photoreceivers of up to 160 Gb/s have been reported.⁸⁵

When the photodetection mechanism takes place along the length of a waveguide, the field screening due to intense optical fields can be much less. An optical waveguide can also increase the absorbance of UTC detectors. Waveguide devices use either *edge-coupling* or *evanescent-wave* coupling. The challenge is for the optical waveguide to match the light's spatial profile from the fiber and

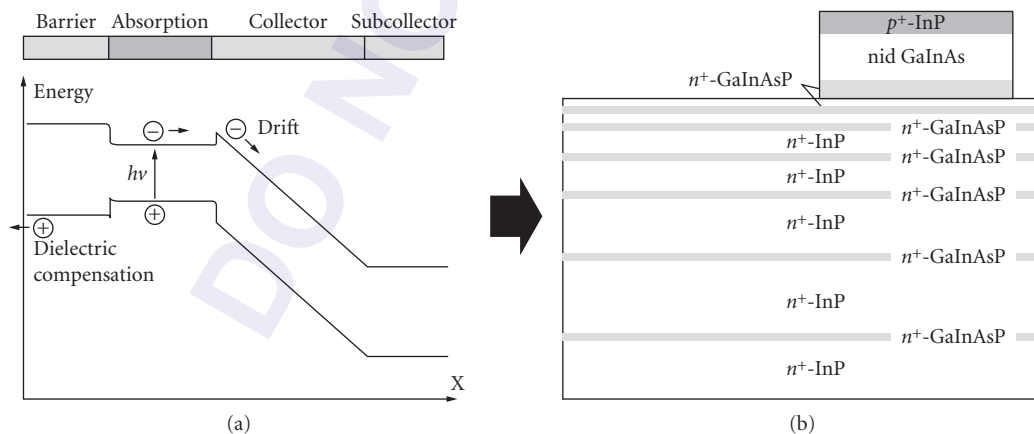


FIGURE 38 (a) Band structure of an UTC photodiode. Only electrons pass through the collector. (b) Multimode evanescent coupled PIN waveguide.

for this light to be coupled effectively into the absorption region of the detector. A variety of methods have been explored, but it is unclear that any design has found universal acceptance. One interesting approach developed at Alcatel is shown in Fig. 38b. The graded layer structure causes incident light to undertake an oscillatory path and be “focused” to the UTC $p^+n^-n^+$ structure, where it is evanescently coupled through the transparent n^+ GaInAsP and GaInAs collecting layers into the absorbing InP layer.⁸⁶

Dark Current

Semiconductor diodes can pass current even in the dark, giving rise to *dark current* that provides a background present in any measurement. This current comes primarily from the thermally generated diffusion of minority carriers out of the n and p regions into the depleted junction region, where they recombine. The current-voltage equation for a pn diode (in the dark) is

$$I = I_s \left[\exp\left(\frac{eV}{\beta kT}\right) - 1 \right] \quad (65)$$

where I_s is the *saturation current* that flows at large back bias (V large and negative). This equation represents the current that passes through any biased pn junction. Photodiodes use pn junctions reverse biased ($V < 0$) to avoid large leakage current.

Here β is the *ideality factor*, which varies from 1 to 2, depending on the diode structure. In a metal-semiconductor junction (Schottky barrier) or an ideal pn junction in which the only current in the dark is due to minority carriers that diffuse from the p and n regions, then $\beta = 1$. However, if there is thermal generation and recombination of carriers in the space-charge region, then β tends toward the value 2. This is more likely to occur in long-wavelength detectors.

The saturation current I_s is proportional to the area A of the diode in an ideal junction: $I_s = e(D_p p_{n0}/L_p + D_n n_{p0}/L_n)A$, where D_n, D_p are diffusion constants, L_n, L_p are diffusion lengths, and n_{p0}, p_{n0} are equilibrium minority carrier densities, all of electrons and holes, respectively. The saturation current I_s can be related to the diode resistance at $V = 0$, measured in the dark through $1/R_0 = -(dI/dV)|_{V=0}$ to give $R_0 = \beta kT/eI_s$. This implies that the dark resistance is inversely proportional to the diode area.

The diffusion current through the diode in Eq. (64) has two components that are of opposite sign in a forward-biased diode: a forward current $I_s \exp(eV/\beta kT)$ and a backward current $-I_s$. Each of these components is statistically independent, coming from diffusive contributions to the forward current and backward current, respectively. This fact is important in understanding the noise properties of photodiodes.

In photodiodes under reverse bias, $V \leq 0$, so $\exp(eV/\beta kT) < 1$; both currents are negative and add. Negative dark current has a direction that is opposite to the current flow in a forward-biased diode. Holes move toward the p region and electrons move toward the n region, the same as photogenerated carriers. These dark currents are thermally generated. Assuming $e|V| \gg kT$, the dark current becomes $I_D = I_s \approx \beta kT/eR_0$. The dark current increases linearly with temperature and is independent of (large enough) reverse bias. *Trap-assisted* thermal generation current increases β ; in this process, carriers trapped in impurity levels can be thermally elevated to the conduction band. The temperature of photodiodes should be kept moderate in order to avoid excess dark current.

When light is present in a reverse-biased photodiode with $V \equiv -V'$, the photocurrent is negative, moving in the direction of the applied voltage, and adding to the negative dark current. The net effect of carrier motion will be to tend to screen the internal field. Defining the magnitude of the photocurrent as $I_{PC} = \eta_D(e/h\nu)P_S$, then the total current is negative:

$$I = -[I_D + I_{PC}] = -I_s \left[1 - \exp\left(\frac{-eV'}{\beta kT}\right) \right] - I_{PC} \quad (66)$$

Noise in Photodiodes

Successful fiber optic communication systems depend on a large signal-to-noise ratio. This requires photodiodes with high sensitivity and low noise. Background noise comes from shot noise due to the discrete process of photon detection, from thermal processes in the load resistor (Johnson noise), and from generation-recombination noise due to carriers within the semiconductor. When used with a field-effect transistor (FET) amplifier, there will also be shot noise from the amplifier and $1/f$ noise in the drain current.

Shot Noise Shot noise is fundamental to all photodiodes and is due to the discrete nature of the conversion of photons to free carriers. The shot noise current is a statistical process. If N photons are detected in a time interval Δt , Poisson noise statistics cause the uncertainty in N to be \sqrt{N} . Using the fact that N electron-hole pairs create a current I through $I = eN/\Delta t$, then the signal-to-noise ratio is $N/\sqrt{N} = \sqrt{N} = \sqrt{(I\Delta t/e)}$. Writing the frequency bandwidth Δf in terms of the time interval through $\Delta f = 1/(2\Delta t)$, the signal-to-noise ratio is: $\text{SNR} = (I/2e\Delta f)^{1/2}$. The root mean square (rms) photon noise, given by \sqrt{N} , creates an rms shot noise current of $i_{\text{SH}} = e(N/\Delta t)^{1/2} = (eI/\Delta t)^{1/2} = (2eI\Delta f)^{1/2}$.

Shot noise depends on the average current I ; therefore, for a given photodiode, it depends on the details of the current-voltage characteristic. Expressed in terms of the photocurrent I_{PC} or the optical signal power P_s (when the dark current is small enough to be neglected) and the responsivity (or sensitivity) \mathfrak{R} , the rms shot noise current is

$$i_{\text{SH}} = \sqrt{2eI_{\text{PC}}\Delta f} = \sqrt{2e\mathfrak{R}P_s\Delta f} \quad (67)$$

The shot noise can be expressed directly in terms of the properties of the diode when all sources of noise are included. Since they are statistically independent, the contributions to the noise currents will be additive. Noise currents can exist in both the forward and backward directions, and these contributions must add, along with the photocurrent contribution. The entire noise current squared becomes

$$i_{\text{N}}^2 = 2e \left\{ I_{\text{PC}} + \left(\frac{\beta kT}{eR_0} \right) \left[1 + \exp\left(\frac{-eV'}{\beta kT} \right) \right] \right\} \Delta f \quad (68)$$

Clearly, noise is reduced by increasing the reverse bias. When the voltage is large, the shot noise current squared becomes $i_{\text{N}}^2 = 2e[I_{\text{PC}} + I_D]\Delta f$. The dark current adds linearly to the photocurrent in calculating the shot noise.

Thermal (Johnson) Noise In addition to shot noise due to the random variations in the detection process, the random thermal motion of charge carriers contributes to a *thermal noise* current, often called *Johnson* or *Nyquist noise*. It can be calculated by assuming thermal equilibrium with $V = 0$, $\beta = 1$, so that Eq. (67) becomes

$$i_{\text{th}}^2 = 4 \left(\frac{kT}{R_0} \right) \Delta f \quad (69)$$

This is just *thermal* or *Johnson noise* in the resistance of the diode. The noise appears as a fluctuating voltage, independent of bias level.

Johnson Noise from External Circuit An additional noise component will be from the load resistor R_L and resistance from the input to the preamplifier, R_i :

$$i_{\text{Nj}}^2 = 4kT \left(\frac{1}{R_L} + \frac{1}{R_i} \right) \Delta f \quad (70)$$

Note that the resistances add in parallel as they contribute to noise current.

Noise Equivalent Power The ability to detect a signal requires having a photocurrent equal to or higher than the noise current. The amount of noise that detectors produce is often characterized by the *noise equivalent power* (NEP), which is the amount of optical power required to produce a photocurrent just equal to the noise current. Define the noise equivalent photocurrent I_{NEP} , which is set equal to the noise current i_{SH} . When the dark current is negligible, the *noise equivalent photocurrent* is $i_{SH} = \sqrt{2eI_{NE}\Delta f} = I_{NE}$.

Thus, the noise equivalent current is $I_{NE} = 2e\Delta f$, and depends only on the bandwidth Δf . The noise equivalent power can now be expressed in terms of the noise equivalent photo current:

$$NEP = \frac{I_{NE}}{\eta} \frac{hv}{e} = 2 \frac{hv}{\eta} \Delta f \quad (71)$$

The second equality assumes the absence of dark current. In this case, the NEP can be decreased only by increasing the quantum efficiency (for a fixed bandwidth). In terms of sensitivity (amperes per watt): $NEP = 2(e/\mathcal{R})\Delta f = I_{NE} \Delta f$. This expression is usually valid for photodetectors used in optical communication systems, which have small dark currents.

If dark current dominates, $i_N = \sqrt{2eI_D \Delta f}$, and

$$NEP = \sqrt{\frac{2I_D \Delta f}{e}} \frac{hv}{\eta} \quad (72)$$

This is often the case in infrared detectors such as germanium. Note that the dark-current-limited noise equivalent power is proportional to the square root of the area of the detector, because the dark current is proportional to the detector area. The NEP is also proportional to the square root of the bandwidth Δf . Thus, in photodetectors whose noise is dominated by dark current, NEP divided by the square root of area times bandwidth should be a constant. The inverse of this quantity has been called the *detectivity* D^* and is often used to describe infrared detectors. In photodiodes used for communications, dark current usually does not dominate and it is better to use Eq. (70), an expression which is independent of area, but depends linearly on bandwidth.

13.15 AVALANCHE PHOTODIODES, MSM DETECTORS, AND SCHOTTKY DIODES

The majority of optical communication systems use photodiodes, sometimes integrated with a preamplifier. Avalanche photodiodes offer an alternative way to create gain. Other detectors sometimes used are low-cost MSM detectors or ultrahigh-speed Schottky diodes. Systems decisions, such as signal-to-noise, cost, and reliability will dictate the choice.

Avalanche Detectors

When large voltages are applied to photodiodes, the avalanche process produces gain, but at the cost of excess noise and slower speed. In fiber telecommunication applications, where speed and signal-to-noise are of the essence, avalanche photodiodes (APDs) are frequently at a disadvantage. Nonetheless, in long-haul systems at 2488 Mb/s, APDs may provide up to 10 dB greater sensitivity in receivers limited by amplifier noise. While APDs are inherently complex and costly to manufacture, they are less expensive than optical amplifiers and may be used when signals are weak.

Gain (Multiplication) When a diode is subject to a high reverse-bias field, the process of impact ionization makes it possible for a single electron to gain sufficient kinetic energy to knock another electron from the valence to the conduction band, creating another electron-hole pair. This enables the quantum efficiency to be >1 . This internal multiplication of photocurrent could be compared to the gain in photomultiplier tubes. The *gain* (or *multiplication*) M of an APD is the ratio of the

photocurrent divided by that which would give unity quantum efficiency. Multiplication comes with a penalty of an excess noise factor, which multiplies shot noise. This excess noise is function of both the gain and the ratio of impact ionization rates between electrons and holes.

Phenomenologically, the low-frequency multiplication factor is

$$M_{\text{DC}} = \frac{1}{1 - (V/V_B)^n} \quad (73)$$

where the parameter n varies between 3 and 6, depending on the semiconductor, and V_B is the breakdown voltage. Gains of $M > 100$ can be achieved in silicon APDs, while they are more typically 10 to 20 for longer-wavelength detectors, before multiplied noise begins to exceed multiplied signal. A typical voltage will be 75 V in InGaAs APDs, while in silicon it can be 400 V.

The avalanche process involves using an electric field high enough to cause carriers to gain enough energy to accelerate them into ionizing collisions with the lattice, producing electron-hole pairs. Then, both the original carriers and the newly generated carriers can be accelerated to produce further ionizing collisions. The result is an avalanche process.

In an intrinsic i layer (where the electric field is uniform) of width W_i , the gain relates to the fundamental avalanche process through $M = 1/(1 - aW_i)$, where a is the *impact ionization coefficient*, which is the number of ionizing collisions per unit length. When $aW_i \rightarrow 1$, the gain becomes infinite and the diode breaks down. This means that avalanche multiplication appears in the regime before the probability of an ionizing collision is 100 percent. The gain is a strong function of voltage, and these diodes must be used very carefully. The total current will be the sum of avalanching electron current and avalanching hole current.

In most *pin* diodes the i region is really low n -doped. This means that the field is not exactly constant, and an integration of the avalanche process across the layer must be performed to determine a . The result depends on the relative ionization coefficients; in III-V materials they are approximately equal. In this case, aW_i is just the integral of the ionizing coefficient that varies rapidly with electric field.

Separate Absorber and Multiplication APDs In this design the long-wavelength infrared light is absorbed in an intrinsic narrow-bandgap InGaAs layer, and photocarriers move to a separate, more highly n -doped InP layer that supports a much higher field. This layer is designed to provide avalanche gain in a separate region without excessive dark currents from tunneling processes. This layer typically contains the *pn* junction, which traditionally has been diffused. Fabrication procedures such as etching a mesa, burying it, and introducing a guard ring electrode are all required to reduce noise and dark current. All-epitaxial structures provide low-cost batch-processed devices with high performance characteristics.⁸⁷

Speed When the gain is low, the speed is limited by the RC time constant. As the gain increases, the avalanche buildup time limits the speed, and for modulated signals the multiplication factor decreases. The multiplication factor as a function of modulation frequency is

$$M(\omega) = \frac{M_{\text{DC}}}{\sqrt{1 + M_{\text{DC}}^2 \omega^2 \tau_1^2}} \quad (74)$$

where $\tau_1 = p\tau$, with τ as the multiplication-region transit time and p as a number that changes from 2 to 1/3 as the gain changes from 1 to 1000. The gain decreases from its low-frequency value when $M_{\text{DC}}\omega = 1/\tau_1$. The gain-bandwidth product describes the characteristics of an avalanche photodiode in a communication system.

Noise The shot noise in an APD is that of a *pin* diode multiplied by M^2 times an excess noise factor F_e :

$$i_s^2 = 2eI_{\text{PC}} \Delta f M^2 F_e \quad (75)$$

where

$$F_e(M) = \beta M + (1 - \beta) \left(2 - \frac{1}{M} \right)$$

In this expression, β is the ratio of the ionization coefficient of the opposite type divided by the ionization coefficient of the carrier type that initiates multiplication. In the limit of equal ionization coefficients of electrons and holes (usually the case in III-V semiconductors), $F_e = M$ and $F_h = 1$. Typical numerical values for enhanced APD sensitivity are given in Chap. 26 in Vol. II, Fig. 15.

Dark Current and Shot Noise In an APD, dark current is the sum of the unmultiplied current I_{du} , mainly due to surface leakage, and the bulk dark current experiencing multiplication I_{dm} , multiplied by the gain: $I_d = I_{du} + MI_{dm}$. The shot noise from dark (leakage) current i_d is $i_d^2 = 2e[I_{du} + I_{dm}M^2F_e(M)]\Delta f$.

The proper use of APDs requires choosing the proper design, carefully controlling the voltage, and using the APD in a suitably designed system, since the noise is so large.

MSM Detectors

Volume II, Chap. 26, Fig. 1 of this *Handbook* shows that interdigitated electrodes on top of a semiconductor can provide a planar configuration for electrical contacts. Either a *pn* junction or bulk semiconductor material can reside under the interdigitated fingers. The MSM geometry has the advantage of lower capacitance for a given cross-sectional area, but the transit times may be longer, limited by the lithographic ability to produce very fine lines. Typically, MSM detectors are photoconductive. Volume II, Chap. 26, Fig. 17 shows the geometry of high-speed interdigitated photoconductors. These are simple to fabricate and can be integrated in a straightforward way onto MESFET preamplifiers.

Consider parallel electrodes deposited on the surface of a photoconductive semiconductor with a distance L between them. Under illumination, the photocarriers will travel laterally to the electrodes. The photocurrent in the presence of P_s input optical flux at photon energy $h\nu$ is: $I_{ph} = q\eta GP h\nu$. The photoconductive gain G is the ratio of the carrier lifetime τ to the carrier transit time τ_t : $G = \tau/\tau_t$. Decreasing the carrier lifetime increases the speed but decreases the sensitivity.

The output signal is due to the time-varying resistance that results from the time-varying photo-induced carrier density $N(t)$:

$$R_s(t) = \frac{L}{eN(t)\mu w d_e} \quad (76)$$

where μ is the sum of the electron and hole mobilities, w is the length along the electrodes excited by light, and d_e is the effective absorption depth into the semiconductor.

Usually, MSM detectors are not the design of choice for high-quality communication systems. Nonetheless, their ease of fabrication and integration with other components makes them desirable for some low-cost applications—for example, when there are a number of parallel channels and dense integration is required.

Schottky Photodiodes

A Schottky photodiode uses a metal-semiconductor junction rather than a *pin* junction. An abrupt contact between metal and semiconductor can produce a space-charge region. Absorption of light in this region causes photocurrent that can be detected in an external circuit. Because metal-semiconductor diodes are majority carrier devices they may be faster than *pin* diodes (they rely on drift currents only; there is no minority carrier diffusion). Modulation speeds up to 100 GHz have been reported in a $5 \times 5\text{-}\mu\text{m}$ area detector with a $0.3\text{-}\mu\text{m}$ thin drift region using a semitransparent platinum film 10 nm thick to provide the abrupt Schottky contact. Resonant reflective enhancement of the light has been used to improve sensitivity.

13.16 REFERENCES

1. D. Botez, *IEEE J. Quant. Electr.* **17**:178 (1981).
2. See, for example, E. Garmire and M. Tavis, *IEEE J. Quant. Electr.* **20**:1277 (1984).
3. B. B. Elenkrig, S. Smetona, J. G. Simmons, T. Making, and J. D. Evans, *J. Appl. Phys.* **85**:2367 (1999).
4. M. Yamada, T. Anan, K. Tokutome, and S. Sugou, *IEEE Photon. Technol. Lett.* **11**:164 (1999).
5. T. C. Hasenberg and E. Garmire, *IEEE J. Quant. Electr.* **23**:948 (1987).
6. D. Botez and M. Ettenberg, *IEEE J. Quant. Electr.* **14**:827 (1978).
7. G. P. Agrawal and N. K. Dutta, *Semiconductor Lasers*, 2d ed., Van Nostrand Reinhold, New York, 1993, Sec. 6.4.
8. G. H. B. Thompson, *Physics of Semiconductor Laser Devices*, John Wiley & Sons, New York, 1980, Fig. 7.8.
9. K. Tatah and E. Garmire, *IEEE J. Quant. Electr.* **25**:1800 (1989).
10. G. P. Agrawal and N. K. Dutta, Sec. 6.4.3.
11. W. H. Cheng, A. Mar, J. E. Bowers, R. T. Huang, and C. B. Su, *IEEE J. Quant. Electr.* **29**:1650 (1993).
12. J. T. Verdeyen, *Laser Electronics*, 3d ed., Prentice Hall, Englewood Cliffs, N.J., 1995, p. 490.
13. G. P. Agrawal and N. K. Dutta, Eq. 6.6.32.
14. N. K. Dutta, N. A. Olsson, L. A. Koszi, P. Besomi, and R. B. Wilson, *J. Appl. Phys.* **56**:2167 (1984).
15. G. P. Agrawal and N. K. Dutta, Sec. 6.6.2.
16. G. P. Agrawal and N. K. Dutta, Sec. 6.5.2.
17. L. A. Coldren and S. W. Corizine, *Diode Lasers and Photonic Integrated Circuits*, John Wiley & Sons, New York, 1995, Sec. 5.5.
18. M. C. Tatham, I. F. Lealman, C. P. Seltzer, L. D. Westbrook, and D. M. Cooper, *IEEE J. Quant. Electr.* **28**:408 (1992).
19. H. Jackel and G. Guekos, *Opt. Quant. Electr.* **9**:223 (1977).
20. M. K. Aoki, K. Uomi, T. Tsuchiya, S. Sasaki, M. Okai, and N. Chinone, *IEEE J. Quant. Electr.* **27**:1782 (1991).
21. K. Petermann, *IEEE J. Sel. Top. Quant. Electr.* **1**:480 (1995).
22. R. W. Tkach and A. R. Chaplyvy, *J. Lightwave Technol.* **LT-4**:1655 (1986).
23. T. Hirono, T. Kurosaki, and M. Fukuda, *IEEE J. Quant. Electr.* **32**:829 (1996).
24. Y. Kitaoka, *IEEE J. Quant. Electr.* **32**:822 (1996) Fig. 2.
25. M. Kawase, E. Garmire, H. C. Lee, and P. D. Dapkus, *IEEE J. Quant. Electr.* **30**:981 (1994).
26. L. A. Coldren and S. W. Corizine, Sec. 4.3.
27. S. L. Chuang, *Physics of Optoelectronic Devices*, John Wiley & Sons, New York, 1995, Fig. 10.33.
28. T. L. Koch and U. Koren, *IEEE J. Quant. Electr.* **27**:641 (1991).
29. B. G. Kim and E. Garmire, *J. Opt. Soc. Am.* **A9**:132 (1992).
30. A. Yariv, *Optical Electronics*, 4th ed., Saunders, Philadelphia, Pa., 1991, Eq. 13.6–19.
31. J. De Merlier, K. Mizutani, S. Sudo, K. Naniwae, Y. Furushima, S. Sato, K. Sato, and K. Kudo, *IEEE Photon. Technol. Lett.* **17**:681 (2005).
32. K. Takabayashi, K. Takada, N. Hashimoto, M. Doi, S. Tomabechei, G. Nakagawa, H. Miyata, T. Nakazawa, and K. Morito, *Electron. Lett.* **40**:1187 (2004).
33. Y. A. Akulova, G. A. Fish, P.-C. Koh, C. L. Schow, P. Kozodoy, A. P. Dahl, S. Nakagawa, et al., *IEEE J. Sel. Top. Quant. Electr.* **8**:1349 (2002).
34. H. Ishii, Y. Tohmori, Y. Yoshikuni, T. Tamamura, and Y. Kondo, *IEEE Photon. Technol. Lett.* **5**:613 (1993).
35. M. Gotoda, T. Nishimura, and Y. Tokuda, *J. Lightwave Technol.* **23**:2331 (2005).
36. Information from <http://www.bookham.com/pr/20070104.cfm>, accessed July, 2008.
37. L. B. Soldano and E. C. M. Pennings, *IEEE J. Lightwave Technol.* **13**:615 (1995).
38. Information from <http://www.pd-ld.com/pdf/PSLEDSeries.pdf>, accessed July, 2008.
39. Information from <http://www.qphotonics.com/catalog/SINGLE-MODE-FIBER-COUPLED-LASER-DIODE-30mW--1550-nm-p-204.html>, accessed July 2008.

40. C. L. Jiang and B. H. Reysen, *Proc. SPIE* **3002**:168 (1997) Fig. 7.
41. See, for example, P. Bhattacharya, *Semiconductor Optoelectronic Devices*, Prentice-Hall, New York, 1998, Chap. 6.
42. C. H. Chen, M. Hargis, J. M. Woodall, M. R. Melloch, J. S. Reynolds, E. Yablonovitch, and W. Wang, *Appl. Phys. Lett.* **74**:3140 (1999).
43. Y. A. Wu, G. S. Li, W. Yuen, C. Caneau, and C. J. Chang-Hasnain, *IEEE J. Sel. Topics Quant. Electr.* **3**:429 (1997).
44. See, for example, F. L. Pedrotti and L. S. Pedrotti, *Introduction to Optics*, Prentice-Hall, Englewood Cliffs, N.J., 1987.
45. N. Nishiyama, C. Caneau, B. Hall, G. Guryanov, M. H. Hu, X. S. Liu, M.-J. Li, R. Bhat, and C. E. Zah, *IEEE J. Sel. Top. Quant. Electr.* **11**:990 (2005).
46. M. Orenstein, A. Von Lehmen, C. J. Chang-Hasnain, N. G. Stoffel, J. P. Harbison, L. T. Florez, E. Clausen, and J. E. Jewell, *Appl. Phys. Lett.* **56**:2384 (1990).
47. Y. A. Wu, G. S. Li, R. F. Nabiev, K. D. Choquette, C. Caneau, and C. J. Chang-Hasnain, *IEEE J. Sel. Top. Quant. Electr.* **1**:629 (1995).
48. K. Mori, T. Asaka, H. Iwano, M. Ogura, S. Fujii, T. Okada, and S. Mukai, *Appl. Phys. Lett.* **60**:21 (1992).
49. E. Söderberg, J. S. Gustavsson, P. Modh, A. Larsson, Z. Zhang, J. Berggren, and M. Hammar, *IEEE Photon. Technol. Lett.* **19**:327 (2007).
50. C. J. Chang-Hasnain, M. Orenstein, A. Von Lehmen, L. T. Florez, J. P. Harbison, and N. G. Stoffel, *Appl. Phys. Lett.* **57**:218 (1990).
51. F. Romstad, S. Bischoff, M. Juhl, S. Jacobsen and D. Birkedal, *Proc. SPIE*, **C-1-14**: 69080 (2008).
52. C. Carlsson, C. Angulo Barrios, E. R. Messmer, A. Löfqvist, J. Halonen, J. Vukusic, M. Ghisoni, S. Lourduos, and A. Larsson, *IEEE J. Quant. Electr.* **37**:945 (2001).
53. A. Valle, M. Gómez-Molina, and L. Pesquera, *IEEE J. Sel. Top. Quant. Electr.* **14**:895 (2008).
54. J. W. Law and G. P. Agrawal, *IEEE J. Sel. Top. Quant. Electr.* **3**:353 (1997).
55. Data from the Web sites of Lasermate, Finisar, ULM, Raycan, JDS Uniphase and LuxNet, accessed on July 24, 2008.
56. Data from the Web sites of Firecomms and Vixar, accessed on July, 2008.
57. Data from the Web sites of Vertilas and Raycan, accessed on July, 2008.
58. S. K. Korotky and R. C. Alferness, "Ti:LiNbO₃ Integrated Optic Technology," in L. D. Hutcheson (ed.), *Integrated Optical Circuits and Components*, Dekker, New York, 1987.
59. G. Y. Wang and E. Garmire, *Opt. Lett.* **21**:42 (1996).
60. A. Yariv, Table 9.2.
61. G. D. Boyd, R. C. Miller, K. Nassau, W. L. Bond, and A. Savage, *Appl. Phys. Lett.* **5**:234 (1964).
62. Information from www.covega.com, accessed on June, 2008.
63. F. P. Leonberger and J. P. Donnelly, "Semiconductor Integrated Optic Devices," in T. Tamir (ed.), *Guided Wave Optoelectronics*, Springer-Verlag, 1990, p. 340.
64. C. C. Chen, H. Porte, A. Carenco, J. P. Goedgebuer, and V. Armbruster, *IEEE Photon. Technol. Lett.* **9**:1361 (1997).
65. C. T. Mueller and E. Garmire, *Appl. Opt.* **23**:4348 (1984).
66. L. B. Soldano and E. C. M. Pennings, *IEEE J. Lightwave Technol.* **13**:615 (1995).
67. D. A. May-Arrijoja, P. LiKamWa, R. J. Selvas-Aguilar, and J. J. Sánchez-Mondragón, *Opt. Quant. Electr.* **36**:1275 (2004).
68. R. Thapliya, T. Kikuchi, and S. Nakamura, *Appl. Opt.* **46**:4155 (2007).
69. R. Thapliya, S. Nakamura, and T. Kikuchi, *Appl. Opt.* **45**:5404 (2006).
70. H. Q. Hou, A. N. Cheng, H. H. Wieder, W. S. C. Chang, and C. W. Tu, *Appl. Phys. Lett.* **63**:1833 (1993).
71. A. Ramdane, F. Devauz, N. Souli, D. Dalprat, and A. Ougazzaden, *IEEE J. Sel. Top. Quant. Electr.* **2**:326 (1996).
72. S. D. Koehler and E. M. Garmire, in T. Tamir, H. Bertoni, and G. Griffel (eds.), *Guided-Wave Optoelectronics: Device Characterization, Analysis and Design*, Plenum Press, New York, 1995.
73. See, for example, S. Carbonneau, E. S. Koteles, P. J. Poole, J. J. He, G. C. Aers, J. Haysom, M. Buchanan, et al., *IEEE J. Sel. Top. Quantum Electron.* **4**:772 (1998).

74. J. A. Trezza, J. S. Powell, and J. S. Harris, *IEEE Photon. Technol. Lett.* **9**:330 (1997).
75. Information from www.okioptical.com, accessed on June 30, 2008.
76. Information from www.ciphotonics.com, accessed on June 30, 2008.
77. M. Jupina, E. Garmire, M. Zembutsu, and N. Shibata, *IEEE J. Quant. Electr.* **28**:663 (1992).
78. S. D. Koehler, E. M. Garmire, A. R. Kost, D. Yap, D. P. Doctor, and T. C. Hasenberg, *IEEE Photon. Technol. Lett.* **7**:878 (1995).
79. A. Sneh, J. E. Zucker, B. I. Miller, and L. W. Stultz, *IEEE Photon. Technol. Lett.* **9**:1589 (1997).
80. J. Pamulapati, J. P. Loehr, J. Singh, and P. K. Bhattacharya, *J. Appl. Phys.* **69**:4071 (1991).
81. H. Feng, J. P. Pang, M. Sugiyama, K. Tada, and Y. Nakano, *IEEE J. Quant. Electr.* **34**:1197 (1998).
82. J. Wang, J. E. Zucker, J. P. Leburton, T. Y. Chang, and N. J. Sauer, *Appl. Phys. Lett.* **65**:2196 (1994).
83. N. Yoshimoto, Y. Shibata, S. Oku, S. Kondo, and Y. Noguchi, *IEEE Photon. Technol. Lett.* **10**:531 (1998).
84. A. Yariv, Sec. 11.7.
85. Y. Muramoto, K. Kato, M. Mitsuohara, O. Nakajima, Y. Matsuoka, N. Shimizu and T. Ishibashi, *Electro. Lett.* **34**(1):122 (1998).
86. M. Achouche, V. Magnin, J. Harari, F. Lelarge, E. Derouin, C. Jany, D. Carpentier, F. Blache, and D. Decoster. *IEEE Photon. Technol. Lett.* **16**:584, 2004.
87. E. Hasnain et al., *IEEE J. Quant. Electr.* **34**:2321 (1998).

John A. Buck

Georgia Institute of Technology
School of Electrical and Computer Engineering
Atlanta, Georgia

14.1 INTRODUCTION

The development of optical fiber amplifiers has led to dramatic increases in the transport capacities of fiber communication systems. At the present time, fiber amplifiers are used in practically every long-haul optical fiber link and advanced large-scale network. Additional applications of fiber amplifiers include their use as gain media in fiber lasers, as wavelength converters, and as stand-alone high-intensity light sources.

The original intent in fiber amplifier development was to provide a simpler alternative to the electronic repeater, chiefly by allowing the signal to remain in optical form throughout a link or network. Fiber amplifiers as repeaters offer additional advantages, which include the ability to change system data rates as needed, or to simultaneously transmit multiple rates—all without the need to modify the transmission link. A further advantage is that signal power at multiple wavelengths can be simultaneously boosted by a single amplifier—a task that would otherwise require a separate electronic repeater for each wavelength. This latter feature contributed to the realization of dense wavelength division multiplexed (DWDM) systems that provide terabit per second data rates.¹ As an illustration, the useful gain in a normally configured erbium-doped fiber amplifier (EDFA) occupies a wavelength range spanning 1.53 to 1.56 μm , which defines the C band. In DWDM systems this allows, for example, the use of some 40 channels having 100 GHz spacing.

A fundamental disadvantage of the fiber amplifier as a repeater is that dispersion is not reset. This requires additional efforts in dispersion management, which may include optical or electronic equalization methods.^{2,3} More recently, special fiber amplifiers have been developed that provide compensation for dispersion.^{4,5} The development⁶ and deployment⁷ of long-range systems that employ optical solitons have also occurred. The use of solitons (pulses that maintain their shape by balancing linear group velocity dispersion with nonlinear self-phase modulation) requires fiber links in which optical power levels can be adequately sustained over long distances. The use of fiber amplifiers allows this possibility. The deployment of fiber amplifiers in commercial networks demonstrates the move toward *transparent* fiber systems, in which signals are maintained in optical form, and in which multiple wavelengths, data rates, and modulation formats are supported.

Successful amplifiers can be grouped into three main categories. These include (1) rare-earth-doped fibers (including EDFAs), in which dopant ions in the fiber core provide gain through stimulated emission, (2) Raman amplifiers, in which gain for almost any optical wavelength can be formed

TABLE 1 Fiber Transmission Bands Showing Fiber Amplifier Coverage

Designation	Meaning	Wavelength Range (μm)	Amplifier (Pump Wavelength) [Ref]
O	Original	1.26–1.36	PDFA (1.02) ⁸
E	Extended	1.36–1.46	Raman (1.28–1.37)
S	Short	1.46–1.53	TDFA (0.8, 1.06, or 1.56 with 1.41) ¹¹
C	Conventional	1.53–1.56	EDFA (0.98, 1.48), EYDFA (1.06) ¹⁰
L	Long	1.56–1.63	Reconfigured EDFA (0.98, 1.48) ¹³
U (XL)	Ultralong	1.63–1.68	Raman (1.52–1.56)

in conventional or specialty fiber through stimulated Raman scattering, and (3) parametric amplification, in which signals are amplified through nonlinear four-wave mixing in fiber.

Within the first category, the most widely used are the erbium-doped fiber amplifiers, in which gain occurs at wavelengths in the vicinity of 1.53 μm . The amplifiers are optically pumped using light at either 1.48 μm or (more commonly) 0.98 μm wavelengths. Other devices include praseodymium-doped fiber amplifiers (PDFAs), which provide gain at 1.3 μm and which are pumped at 1.02 μm .⁸ Ytterbium-doped fibers (YDFAs)⁹ amplify from 0.98 to 1.15 μm , using pump wavelengths between 0.91 and 1.06 μm ; erbium-ytterbium codoped fibers (EYDFAs) enable use of pump light at 1.06 μm while providing gain at 1.55 μm .¹⁰ Additionally, thulium and thulium/terbium-doped fluoride fibers have been constructed for amplification at 0.8, 1.4, and 1.65 μm .¹¹

In the second category, gain from stimulated Raman scattering (SRS) develops as power couples from an optical pump wave to a longer-wavelength signal (Stokes) wave, which is to be amplified. This process occurs as both waves, which either copropagate or counterpropagate, interact with vibrational resonances in the glass material. Raman fiber amplifiers have the advantage of enabling useful gain to occur at any wavelength (or multiple wavelengths) at which the fiber exhibits low loss, and for which the required pump wavelength is available. Long spans (typically several kilometers) are usually necessary for Raman amplifiers, whereas only a few meters are needed for a rare-earth-doped amplifier. Incorporating Raman amplifiers in systems is therefore often done by introducing the required pump power into a section of the existing transmission fiber.

Finally, in nonlinear parametric amplification (arising from four-wave mixing) signals are again amplified in the presence of a strong pump wave at a different frequency. The process differs in that the electronic (catalytic) nonlinearity in fiber is responsible for wave coupling, as opposed to vibrational resonances (optical phonons) that mediate wave coupling in SRS. Parametric amplifiers can exhibit exponential gain coefficients that are twice the value of those possible for SRS in fiber.¹² Achieving high parametric gain is a more complicated problem in practice, however, as will be discussed. A useful by-product of the parametric process is a wavelength-shifted replica of the signal wave (the idler) which is proportional to the phase conjugate of the signal.

Table 1 summarizes doped fiber amplifier usage in the optical fiber transmission bands, in which it is understood that Raman or parametric amplification can in principle be performed in any of the indicated bands. In cases where doped fibers are unavailable, Raman amplification is indicated, along with the corresponding pump wavelengths.

14.2 RARE-EARTH-DOPED AMPLIFIER CONFIGURATION AND OPERATION

Pump Configuration and Optimum Fiber Length

A typical rare-earth-doped fiber amplifier configuration consists of the doped fiber positioned between polarization-independent optical isolators. Pump light is input by way of a wavelength-selective coupler (WSC) which can be configured for forward, backward, or bidirectional pumping (Fig. 1). Pump absorption throughout the amplifier length results in a population inversion that

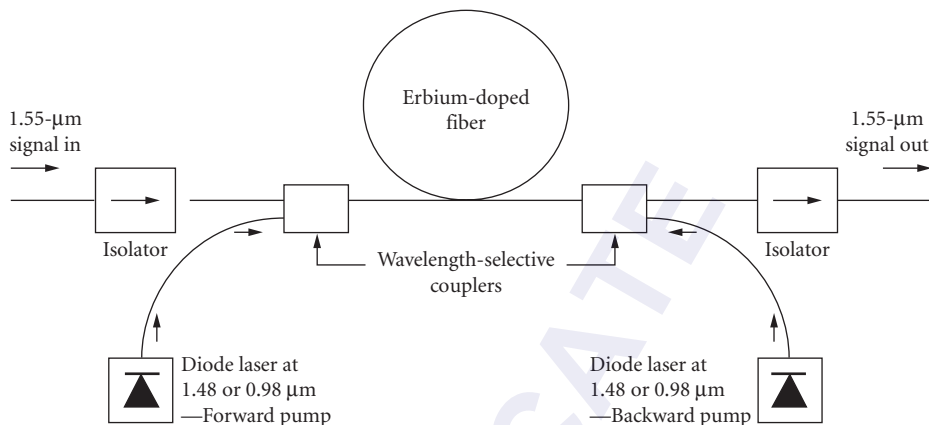


FIGURE 1 General erbium-doped fiber amplifier configuration showing bidirectional pumping.

varies with position along the fiber; this reaches a minimum at the fiber end opposite the pump laser for unidirectional pumping, or minimizes at midlength for bidirectional pumping, using equal pump powers. To achieve the highest overall gain for unidirectional pumping, the fiber length is chosen so that at the output (the point of minimum pump power), the exponential gain coefficient is zero—and no less. If the amplifier is too long, some reabsorption of the signal will occur beyond the transparency point, as the gain goes negative. With lengths shorter than the optimum, full use is not made of the available pump energy, and the overall gain factor is reduced. Other factors may modify the optimum length, particularly if substantial gain saturation occurs, or if amplified spontaneous emission (ASE), which can result in additional gain saturation and noise, is present.¹⁴

Isolators maintain unidirectional light propagation so that, for example, no Rayleigh backscattered or reflected light from further down the link can re-enter the amplifier and cause gain quenching, noise enhancement, or possibly lasing. Double-pass and segmented configurations are also used; in the latter, isolators are positioned between two or more lengths of amplifying fiber which are separately pumped. The result is that gain quenching and noise arising from backscattered light or from amplified spontaneous emission (ASE) are reduced over that of a single fiber amplifier of the combined lengths.

Regimes of Operation

There are roughly three operating regimes, the choice between which is determined by the use intended for the amplifier.^{15,16} These are (1) small-signal or linear regime, (2) saturation, and (3) deep saturation.

In the linear regime, low input signal levels ($< 1 \mu\text{W}$) are amplified with negligible gain saturation, using the optimum amplifier length as discussed in the last section. EDFA gains that range between 25 and 35 dB are possible in this regime.¹⁶ Amplifier gain in decibels is defined in terms of the input and output signal powers as $G \text{ (dB)} = 10 \log_{10} (P_s^{\text{out}}/P_s^{\text{in}})$.

In the saturation regime, the input signal level is high enough to cause a measurable reduction (compression) in the net gain. A useful figure of merit is the *input saturation power*, $P_{\text{sat}}^{\text{in}}$, defined as the input signal power required to compress the net amplifier gain by 3 dB. Specifically, the gain in this case is $G = G_{\text{max}} - 3 \text{ dB}$, where G_{max} is the small-signal gain. A related parameter is the *saturation output power*, $P_{\text{sat}}^{\text{out}}$, defined as the amplifier output that is achieved when the overall gain is compressed by 3 dB. The two quantities are thus related through $G_{\text{max}} - 3 \text{ dB} = 10 \log_{10} (P_{\text{sat}}^{\text{out}}/P_{\text{sat}}^{\text{in}})$. Again, it is assumed that the amplifier length is optimized when defining these parameters.

The *dynamic range* of the amplifier is defined through $P_s^{\text{in}} \leq P_{\text{sat}}^{\text{in}}$, or equivalently $P_s^{\text{out}} \leq P_{\text{sat}}^{\text{out}}$. For an N -channel wavelength-division multiplexed signal, the dynamic range is reduced accordingly by a factor of $1/N$, assuming a flat gain spectrum.¹⁶

With the amplifier operating in deep saturation, gain compressions on the order of 20 to 40 dB occur.¹⁵ This is typical of *power amplifier* applications, in which input signal levels are high, and where the maximum output signal power is desired. In this application, the concept of *power conversion efficiency* (PCE) between pump and signal becomes important. It is defined as $PCE = (P_s^{\text{out}} - P_s^{\text{in}})/P_p^{\text{in}}$, where P_p^{in} is the input pump power.

Another important quantity that is pertinent to the deep saturation regime is the *saturated output power*, $P_s^{\text{out}}(\text{max})$ (not to be confused with the saturation output power described above). $P_s^{\text{out}}(\text{max})$ is the maximum output signal power that can be achieved for a given input signal level and available pump power. This quantity would maximize when the amplifier, having previously been fully inverted, is then completely saturated by the signal. Maximum saturation, however, requires the input signal power to be extremely high, such that ultimately, $P_s^{\text{out}}(\text{max}) \approx P_s^{\text{in}}$, representing a net gain of nearly 0 dB. Clearly the more important situations are those in which moderate signal powers are to be amplified; in these cases the choice of pump power and pumping configuration can substantially influence $P_s^{\text{out}}(\text{max})$.

14.3 EDFA PHYSICAL STRUCTURE AND LIGHT INTERACTIONS

Energy Levels in the EDFA

Gain in the erbium-doped fiber system occurs when an inverted population exists between parts of the ${}^4I_{13/2}$ and ${}^4I_{15/2}$ states, as shown in Fig. 2a.¹⁷ This notation uses the standard form, $(2S+1)L_J$, where L , S , and J are the orbital, spin, and total angular momenta, respectively. EDFAs are manufactured by incorporating erbium ions into the glass matrix that forms the fiber core. Interactions between the ions and the host matrix induces Stark splitting of the ion energy levels, as shown in Fig. 2a. This produces an average spacing between adjacent Stark levels of 50 cm^{-1} , and an overall spread of 300 to 400 cm^{-1} within each state. A broader emission spectrum results, since more de-excitation pathways are produced, which occur at different transition wavelengths.

Other mechanisms further broaden the emission spectrum. First, the extent to which ions interact with the glass varies from site to site, as a result of the nonuniform structure of the amorphous glass matrix. This produces some degree of inhomogeneous broadening in the emission spectrum, the extent of which varies with the type of glass host used.¹⁸ Second, thermal fluctuations in the material lead to homogeneous broadening of the individual Stark transitions. The magnitudes of the two broadening mechanisms are 27 to 60 cm^{-1} for inhomogeneous, and 8 to 49 cm^{-1} for homogeneous.¹⁸

The choice of host material strongly affects the shape of the emission spectrum, owing to the character of the ion-host interactions. For example, in pure silica (SiO_2), the spectrum of the Er-doped system is narrowest and has the least smoothness. Use of an aluminosilicate host ($\text{SiO}_2\text{-Al}_2\text{O}_3$), produces slight broadening and smoothing.¹⁹ The broadest spectra, however, occur when using fluoride-based glass, such as ZBLAN ($\text{ZrF}_4\text{-BaF}_2\text{-LaF}_3\text{-AlF}_3\text{-NaF}$).²⁰

Gain Formation

Figure 2b shows how the net emission spectrum is constructed from the superposition of the individual Stark spectra; the latter are associated with the transitions shown in Fig. 2a. Similar diagrams can be constructed for the upward (absorptive) transitions, from which the absorption spectrum can be developed.²⁰ The shapes of both spectra are further influenced by the populations within the Stark-split levels, which assume a Maxwell-Boltzman distribution. The sequence of events in the population dynamics is (1) pump light boosts population from the ground state, ${}^4I_{15/2}$, to the upper Stark levels in the first excited state, ${}^4I_{13/2}$; (2) the upper state Stark level populations thermalize; and (3) de-excitation from ${}^4I_{13/2}$ to ${}^4I_{15/2}$ occurs through either spontaneous or stimulated emission.

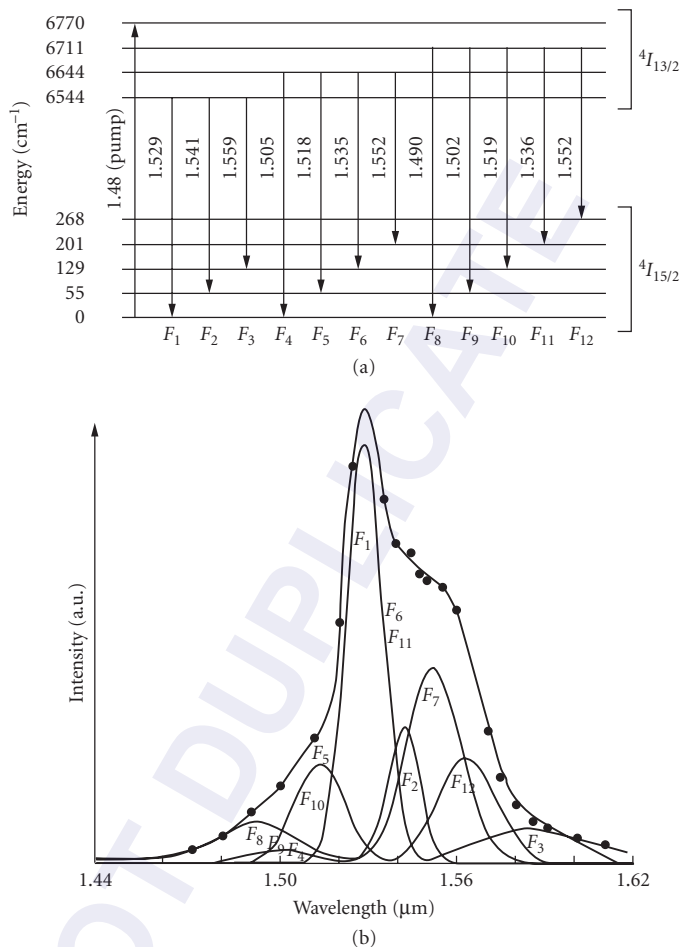


FIGURE 2 (a) Emissive transitions between Stark-split levels of erbium in an aluminosilicate glass host. Values on transition arrows indicate wavelengths in micrometers. (Adapted from Ref. 17.) (b) EDFA fluorescence spectrum arising from the transitions in Fig. 2a. (Reprinted with permission from Ref. 18.)

The system can be treated using a simple two-level (1.48 μm pump) or three-level model (0.98 μm pump), from which rate equations can be constructed that incorporate the actual wavelength- and temperature-dependent absorption and emission crosssections. These models have been formulated with and without inhomogeneous broadening. In most cases, results that are in excellent agreement with experiment have been achieved by assuming only homogeneous broadening.^{21–23}

Pump Wavelength Options in EDFAs

The 1.48 μm pump wavelength corresponds to the energy difference between the two most widely spaced Stark levels, as shown in Fig. 2a. A better alternative is to pump with light at 0.98 μm , which boosts the ground state population to the second excited state, $4I_{11/2}$, which lies above $4I_{13/2}$.

This is followed by rapid nonradiative decay into ${}^4I_{13/2}$ and gain is formed as before. The pumping efficiency suffers slightly at 0.98 μm , owing to some excited state absorption (ESA) between ${}^4I_{11/2}$ and the higher-lying ${}^4F_{7/2}$ state at this wavelength.²⁴ Use of 0.98 μm pump light as opposed to 1.48 μm , will nevertheless yield a more efficient system, since the 0.98 μm pump will not contribute to the de-excitation process, as occurs when 1.48 μm is used.

The *gain efficiency* of a rare-earth-doped fiber is defined as the ratio of the maximum small signal gain to the input pump power, using the optimized fiber length. EDFA efficiencies are typically on the order of 10 dB/mW for pumping at 0.98 μm . For pumping at 1.48 μm , efficiencies are about half the values obtainable at 0.98 μm , and require about twice the fiber length. Other pump wavelengths can be used,²⁴ but with some penalty to be paid in the form of excited state absorption from the ${}^4I_{13/2}$ state into various upper levels, thus depleting the gain that would otherwise be available. This problem is minimized when using either 0.98 or 1.48 μm , and so these two wavelengths are almost exclusively used in erbium-doped fibers.

Noise

Performance is degraded by the presence of noise from two fundamental sources. These are (1) amplified spontaneous emission (ASE), and (2) Rayleigh scattering. Both processes lead to additional light that propagates in the forward and backward directions, and which encounters considerable gain over long amplifier lengths. The more serious of the two noise sources is ASE. In severe cases, involving high-gain amplifiers of long lengths, ASE can be of high enough intensity to partially saturate the gain, thus reducing the available gain for signal amplification. This *self-saturation* effect has been reduced by using the backward pumping geometry.²⁵

In general, ASE can be reduced by (1) assuring that the population inversion is as high as possible (ideally, completely inverted), (2) operating the amplifier in the deep saturation regime, or (3) using two or more amplifier stages rather than one continuous length of fiber, and positioning bandpass filters and isolators between stages. Rayleigh scattering noise can be minimized by using multistage configurations, in addition to placing adequate controls on dopant concentration and confinement during the manufacturing stage.²⁶

The *noise figure* of a rare-earth-doped fiber amplifier is stated in a manner consistent with the IEEE standard definition for a general amplifier (Friis definition). This is the signal-to-noise ratio of the fiber amplifier input divided by the signal-to-noise ratio of the output, expressed in decibels, where the input signal is shot noise limited. Although this definition is widely used, it has become the subject of a debate, arising from the physical nature of ASE noise, and the resulting awkwardness in applying the definition to cascaded amplifier systems.²⁷ An in-depth review of this subject is in Ref. 28.

The best noise figures for EDFAs are achieved by using pump configurations that yield the highest population inversions. Again, the use of 0.98 μm is preferred, yielding noise figures that approach the theoretical limit of 3 dB. Pumping at 1.48 μm gives best results of about 4 dB.¹⁵

Gain Flattening

Use of multiple wavelength channels in WDM systems produces a strong motivation to construct a fiber amplifier in which the gain is uniform for all wavelengths. Thus some means needs to be employed which will effectively flatten the emission spectrum as depicted in Fig. 2b. Flattening techniques can be classified into roughly three categories. First, intrinsic methods can be used; these involve choices of fiber host materials such as fluoride glass²⁹ that yield smoother and broader gain spectra. In addition, by carefully choosing pump power levels, a degree of population inversion can be obtained which will allow some cancellation to occur between the slopes of the absorption and emission spectra,³⁰ thus producing a flatter gain spectrum. Second, spectral filtering at the output of a single amplifier or between cascaded amplifiers can be employed; this effectively produces higher loss for wavelengths that have achieved higher gain. Examples of successful filtering devices include long-period fiber gratings³¹ and Mach-Zehnder filters.³² Third, hybrid amplifiers that use cascaded

configurations of different gain media can be used to produce an overall gain spectrum that is reasonably flat. Flattened gain spectra have been obtained having approximate widths that range from 12 to 85 nm. Reference 33 is recommended for an excellent discussion and comparison of the methods.

14.4 OTHER RARE-EARTH SYSTEMS

Praseodymium-Doped Fiber Amplifiers (PDFAs)

In the praseodymium-doped fluoride system, the strongest gain occurs in the vicinity of 1.3 μm , with the pump wavelength at 1.02 μm . Gain formation is described by a basic three-level model, in which pump light excites the system from the ground state, 3H_4 , to the metastable excited state, 1G_4 . Gain for 1.3 μm light is associated with the downward $^1G_4 \rightarrow ^3H_5$ transition, which peaks in the vicinity of 1.32 to 1.34 μm . Gain diminishes at longer wavelengths, principally as a result of ground state absorption from 3H_4 to 3F_3 .¹⁸

The main problem with the PDFA system has been the reduction of the available gain through the competing $^1G_4 \rightarrow ^3F_4$ transition (2900 cm^{-1} spacing), occurring through multiphonon relaxation. The result is that the *radiative quantum efficiency* (defined as the ratio of the desired transition rate to itself plus all competing transition rates) can be low enough in conventional glass host materials to make the system impractical. The multiphonon relaxation rate is reduced when using hosts having low phonon energies, such as fluoride or chalcogenide glasses. Use of the latter material has essentially solved the problem, by yielding radiative quantum efficiencies on the order of 90%.³⁴ For comparison, erbium systems exhibit quantum efficiencies of nearly 100% for the 1.5 μm transition. Other considerations such as broadening mechanisms and excited state absorption are analogous to the erbium system. References 1 and 35 are recommended for further reading.

Ytterbium-Doped Fiber Amplifiers (YDFAs)

Ytterbium-doping provides the most efficient fiber amplifier system, as essentially no competing absorption and emission mechanisms exist.⁹ This is because in ytterbium, there are only two energy states that are resonant at the wavelengths of interest. These are the ground state $^2F_{7/2}$ and the excited state $^2F_{5/2}$. When doped into the host material, Stark splitting within these levels occurs as described previously, which leads to strong absorption at wavelengths in the vicinity of 0.92 μm , and emission between 1.0 and 1.1 μm , maximizing at around 1.03 μm . Pump absorption is very high, which makes side-pumping geometries practical. Because of the extremely high gain that is possible, Yb-doped fibers are attractive as power amplifiers for 1.06- μm light, and have been employed in fiber laser configurations. YDFAs have also proven attractive as superfluorescent sources,³⁶ in which the output is simply amplified spontaneous emission, and there is no signal input.

Accompanying the high power levels in YDFAs are unwanted nonlinear effects, which are best reduced (at a given power level) by lowering the fiber mode intensity. This has been accomplished to an extent in special amplifier designs that involve large mode effective areas (A_{eff}). Such designs have been based on either conventional fiber³⁷ or photonic crystal fiber.^{38,39} The best results in both cases involve dual-core configurations, in which the pump light propagates in a large core (or inner cladding) which surrounds a smaller concentric core that contains the dopant ions, and that propagates the signal to be amplified. The large inner cladding region facilitates the input coupling of high-power diode pump lasers that have large output beam cross sections. Such designs have proven successful with other amplifiers (including erbium-doped) as well.

Erbium/Ytterbium-Doped Fiber Amplifiers (EYDFAs)

Erbium/ytterbium codoping takes advantage of the strong absorption of ytterbium at the conventional 0.98 μm erbium pump wavelength. When codoped with erbium, ytterbium ions in their excited state transfer their energy to the erbium ions, and gain between 1.53 and 1.56 μm is formed

as before.³ Advantages of such a system include the following: With high pump absorption, side-pumping is possible, thus allowing the use of large-area diode lasers as pumps. In addition, high gain can be established over a shorter propagation distance in the fiber than is possible in a conventional EDFA. As a result, shorter length amplifiers having lower ASE noise can be constructed. An added benefit is that the absorption band allows pumping by high-power lasers such as Nd: YAG (at 1.06 μm) or Nd: YLF (at 1.05 μm), and there is no excited state absorption. Yb-sensitized fibers are attractive for use as C or L band power amplifiers, and in the construction of fiber lasers, in which a short-length, high-gain medium is needed.⁴⁰

14.5 RAMAN FIBER AMPLIFIERS

Amplification by stimulated Raman scattering has proven to be a successful alternative to rare-earth-doped fiber, and has found wide use in long-haul fiber communication systems.^{41,42} Key advantages of Raman amplifiers include (1) improvement in signal-to-noise ratio over rare-earth-doped fiber, and (2) wavelengths to be amplified are not restricted to lie within a specific emission spectrum, but only require a pump wavelength that is separated from the signal wavelength by the Raman resonance. In this way, the entire low-loss spectral range of optical fiber can in principle be covered by Raman amplification, as in, for example, O band applications.⁴³ The main requirement is that a pump laser is available having power output on the order of 0.5 W, and whose frequency is up-shifted from that of the signal by the primary Raman resonance frequency of 440 cm^{-1} or about 13.2 THz. The required pump wavelength, λ_2 , is thus expressed in terms of the signal wavelength, λ_1 , through

$$\lambda_2 = \frac{\lambda_1}{(1 + 0.044\lambda_1)} \quad (1)$$

where the wavelengths are expressed in micrometers.

Another major difference from rare-earth-doped fiber is that Raman amplifiers may typically require lengths on the order of tens of kilometers to achieve the same gain that can be obtained, for example, in 10 m of erbium-doped fiber. The long Raman amplifier span is not necessarily a disadvantage because (1) Raman amplification can be carried out within portions of the existing fiber link, and (2) the long span may contribute to an improvement in the signal-to-noise ratio for the entire link. This happens if Raman amplification is used *after* amplifying using erbium-doped fiber; the Raman amplifier provides gain for the signal, while the long span attenuates the spontaneous emission noise from the EDFA. A Raman amplifier can be implemented at the receiver end of a link by introducing a backward pump at the output end.

The basic configuration of a Raman fiber amplifier is shown in Fig. 3. Pumping can be done in either the forward direction (with input pump power P_{20}) or in the backward direction (with input P_{2L}). The governing equations are Eqs. (10) and (11) in Chap. 10 of this volume, rewritten here in terms of wave power:

$$\frac{dP_1}{dz} = \frac{g_r}{A_{\text{eff}}} P_1 P_2 - \alpha P_1 \quad (2)$$

$$\frac{dP_2}{dz} = \pm \frac{\omega_2}{\omega_1} \frac{g_r}{A_{\text{eff}}} P_1 P_2 \pm \alpha P_2 \quad (3)$$

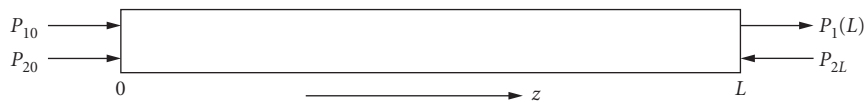


FIGURE 3 Beam configuration for a Raman fiber amplifier using forward or backward pumping.

where the plus and minus signs apply to backward and forward pumping, respectively, and where A_{eff} is the fiber mode cross-sectional area, as before. The peak Raman gain in silica occurs when pump and signal frequencies are spaced by 440 cm^{-1} , corresponding to a wavelength spacing of about $0.1 \text{ }\mu\text{m}$ (see Fig. 2 in Chap. 10 of this volume). The gain in turn is inversely proportional to pump wavelength, λ_2 , and to a good approximation is given by

$$g_r(\lambda_2) = \frac{10^{-11}}{\lambda_2} \text{ cm/W} \quad (4)$$

with λ_2 expressed in micrometers.

The simplest case is the small-signal regime, in which the Stokes power throughout the amplifier is sufficiently low such that negligible pump depletion will occur. The pump power dependence on distance is thus determined by loss in the fiber, and is found by solving Eq. (3) under the assumption that the first term on the right-hand side is negligible. It is further assumed that there is no spontaneous scattering, and that the Stokes and pump fields maintain parallel polarizations. For forward pumping, the solutions of Eqs. (2) and (3) thus simplified are

$$P_1(z) = P_{10} \exp(-\alpha z) \exp \left[\frac{g_r P_{20}}{A_{\text{eff}}} \left(\frac{1 - \exp(-\alpha z)}{\alpha} \right) \right] \quad (5a)$$

$$P_2(z) = P_{20} \exp(-\alpha z) \quad (5b)$$

For backward pumping with no pump depletion, and with fiber length, L :

$$P_1(z) = P_{10} \exp(-\alpha z) \exp \left[\frac{g_r P_{2L}}{A_{\text{eff}}} \exp(-\alpha L) \left(\frac{\exp(\alpha z) - 1}{\alpha} \right) \right] \quad (6a)$$

$$P_2(z) = P_{2L} \exp[\alpha(z - L)] \quad (6b)$$

Implicit throughout is the assumption that there is no spontaneous scattering, and that the Stokes and pump fields maintain parallel polarizations.

The z -dependent signal gain in decibels is $\text{Gain (dB)} = 10 \log_{10}(P_1(z)/P_{10})$. Figure 4⁴⁴ shows this evaluated for several choices of fiber attenuation. It is evident that for forward pumping, an optimum length may exist at which the gain maximizes. For backward pumping, gain will always increase with amplifier length. It is also evident that for a specified fiber length, input pump power, and loss, the same gain is achieved over the *total* length for forward *and* backward pumping, as must be true with no pump depletion.

It is apparent in Eqs. (5a) and (6a) that increased gain is obtained for lower-loss fiber, and by increasing the pump *intensity*, given by P_2/A_{eff} . For a given available pump power, fibers having smaller effective areas, A_{eff} , will yield higher gain, but the increased intensity may result in additional unwanted nonlinear effects.

In actual systems, additional problems arise. Among these is pump depletion, reducing the overall gain as Stokes power levels increase. This is effectively a gain saturation mechanism, and occurs as some of the Stokes power is back-converted to the pump wavelength through the inverse Raman effect. Again, maintaining pump power levels that are significantly higher than the Stokes levels maintains the small signal approximation, and minimizes saturation. In addition, noise may arise from several sources. These include spontaneous Raman scattering, Rayleigh scattering,⁴⁵ pump intensity noise,⁴⁶ and Raman-amplified spontaneous emission from rare-earth-doped fiber amplifiers elsewhere in the link.⁴⁷ Finally, *polarization-dependent gain* (PDG) arises as pump and Stokes field polarizations randomly move in and out of parallelism, owing to the usual changes in fiber birefringence.⁴⁸ This last effect can be reduced by using depolarized pump inputs.

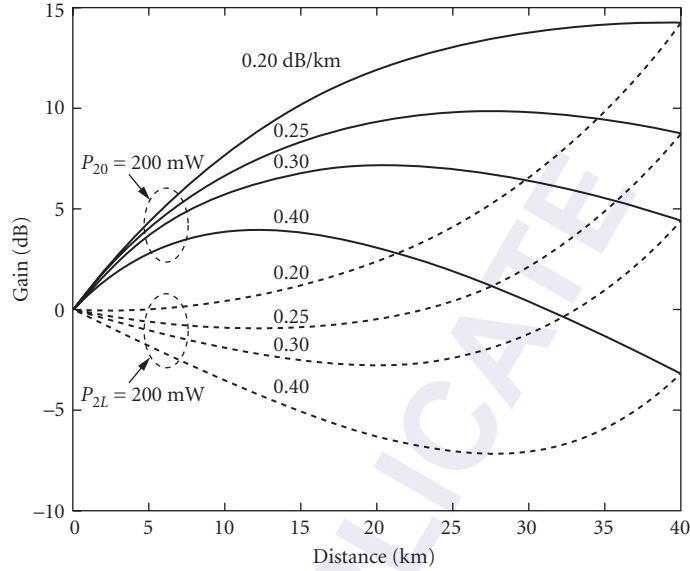


FIGURE 4 Small signal Raman gain as a function of distance z in a single-mode fiber, as calculated using Eqs. (5a) and (6a) for selected values of distributed fiber loss. Plots are shown for cases of forward pumping (solid curves) and backward pumping (dotted curves). Parameter values are $P_{20} = P_{2L} = 200$ mW, $g_r = 7 \times 10^{-12}$ cm/W, and $A_{\text{eff}} = 5 \times 10^{-7}$ cm². (After Ref. 44.)

14.6 PARAMETRIC AMPLIFIERS

Parametric amplification uses the nonlinear four-wave mixing interaction in fiber, as described in Sec. 10.5 (Chap. 10). Two possible configurations are used that involve either a single pump wave, or two pumps at different wavelengths (Fig. 5). The signal to be amplified copropagates with the pumps and is of a different wavelength than either pump. In addition to the amplified signal, the process also generates a fourth wave known as the idler, which is a wavelength-shifted and phase-conjugated replica of the signal. In view of this, parametric amplification is also attractive for use in wavelength conversion applications and—owing to the phase conjugate nature of the idler—in dispersion compensation.

The setup is essentially the same as described in Sec. 10.5, in which we allow the possibility of two distinct pump waves, carrying powers P_1 and P_2 at frequencies ω_1 and ω_2 , or a single pump at frequency ω_0 having power P_0 . The pumps interact with a relatively weak signal, having input power $P_3(0)$, and frequency ω_3 . The signal is amplified as the pump power couples to it, while the idler (power P_4 and frequency ω_4) is generated and amplified. The frequency relations are $\omega_3 + \omega_4 = \omega_1 + \omega_2$ for dual pumps, or $\omega_3 + \omega_4 = 2\omega_0$ for a single pump. In the simple case of a single pump that is nondepleted,

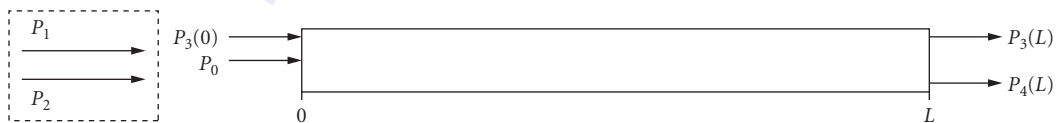


FIGURE 5 Beam configuration for a parametric fiber amplifier using single or dual-wavelength pumping.

and assuming continuous wave operation with a signal input of power $P_3(0)$, the power levels at the amplifier output (length L) are given by:⁴⁹

$$P_3(L) = P_3(0) \left[1 + \left(1 + \frac{\kappa^2}{4g^2} \right) \sinh^2(gL) \right] \quad (7)$$

$$P_4(L) = P_3(0) \left(1 + \frac{\kappa^2}{4g^2} \right) \sinh^2(gL) \quad (8)$$

The parametric gain g is given by

$$g = \left[P_0^2 \left(\frac{2\pi n_2}{\lambda A_{\text{eff}}} \right)^2 - \left(\frac{\kappa}{2} \right)^2 \right]^{1/2} \quad (9)$$

where n_2 is the nonlinear index in m^2/W , λ is the average of the three wavelengths, and A_{eff} is the fiber mode cross-sectional area. The phase mismatch parameter κ includes linear and power-dependent terms:

$$\kappa = \Delta\beta + \frac{2\pi n_2}{\lambda A_{\text{eff}}} P_0 \quad (10)$$

where the linear part, $\Delta\beta = \beta_3 + \beta_4 - 2\beta_0$ has the usual interpretation as the difference in phase constants of a nonlinear polarization wave (at either ω_3 or ω_4) and the field (at the same frequency) radiated by the polarization. The phase constants are expressed in terms of the unperturbed fiber mode indices \tilde{n}_i through $\Delta\beta_i = \tilde{n}_i \omega_i / c$. The second (nonlinear) term on the right hand side of Eq. (10) represents the mode index change arising from the intense pump field through the optical Kerr effect. This uniformly changes the mode indices of all waves, and is thus important to include in the phase mismatch evaluation.

If two pumps are used, having powers P_1 and P_2 , with frequencies ω_1 and ω_2 , Eqs. (9) and (10) are modified by setting $P_0 = P_1 + P_2$. Also, in evaluating P_0^2 in Eq. (9), only the cross term ($2P_1P_2$) is retained in the expression. The linear term in Eq. (10) becomes $\Delta\beta = \beta_3 + \beta_4 - \beta_1 - \beta_2$. With these modifications, Eqs. (7) and (8) may not strictly apply because the two-pump interaction is complicated by the generation of multiple idler waves, in addition to contributions from optically induced Bragg gratings, as discussed in Ref. 50. The advantage of using two pumps of different frequencies is that the phase mismatch is possible to reduce significantly over a much broader wavelength spectrum than is possible using a single pump. In this manner, relatively flat gain spectra over several tens of nanometers have been demonstrated with the pump wavelengths positioned on opposite sides of the zero dispersion wavelength.⁵¹ In single-pump operation, gain spectra of widths on the order of 20 nm have been achieved with the pump wavelength positioned at or near the fiber zero dispersion wavelength.⁵² In either pumping scheme, amplification factors on the same order or greater than those available in Raman amplifiers are in principle obtainable, and best results have exceeded 40 dB.⁵³ In practice, random fluctuations in fiber dimensions and in birefringence represent significant challenges in avoiding phase mismatch, and in maintaining alignment of the interacting field polarizations.⁵⁴

14.7 REFERENCES

1. See for example: B. Hoang and O. Perez, "Terabit Networks," www.ieee.org/portal/site/emergingtech/index, 2007.
2. For background on optical methods of dispersion compensation, see the "Special Mini-Issue on Dispersion Compensation," *IEEE Journal of Lightwave Technology* **12**:1706–1765 (1994).
3. B. J. C. Schmidt, A. J. Lowery, and J. Armstrong, "Experimental Demonstration of Electronic Dispersion Compensation for Long-Haul Transmission Using Direct-Detection Optical OFDM," *IEEE Journal of Lightwave Technology* **26**:196–204 (2008).

4. J. Maury, J. L. Auguste, S. Fevrier, and J. M. Blondy, "Conception and Characterization of a Dual-Eccentric-Core Erbium-Doped Dispersion-Compensating Fiber," *Optics Letters* **29**:700–702 (2004).
5. A. C. O. Chan and M. Premaratne, "Dispersion-Compensating Fiber Raman Amplifiers with Step, Parabolic, and Triangular Refractive Index Profiles," *IEEE Journal of Lightwave Technology* **25**:1190–1197 (2007).
6. L. F. Mollenauer and P. V. Mamyshev, "Massive Wavelength-Division Multiplexing with Solitons," *IEEE Journal of Quantum Electronics* **34**:2089–2102 (1998).
7. Alcatel-Lucent 1625 *LambdaXtreme* Transport, www.alcatel-lucent.com, 2006.
8. Y. Ohishi, T. Kanamori, T. Kitagawa, S. Takahashi, E. Snitzer, and G. H. Sigel, Jr., "Pr³⁺-Doped Fluoride Fiber Amplifier Operation at 1.31 μm ," *Optics Letters* **16**:1747–1749 (1991).
9. R. Paschotta, J. Nilsson, A. C. Tropper, D. C. Hanna, "Ytterbium-Doped Fiber Amplifiers," *IEEE Journal of Quantum Electronics* **33**:1049–1056 (1997).
10. J. E. Townsend, K. P. Jedrezewski, W. L. Barnes, and S. G. Grubb, "Yb³⁺ Sensitized Er³⁺ Doped Silica Optical Fiber with Ultra High Efficiency and Gain," *Electronics Letters* **27**:1958–1959 (1991).
11. S. Sudo, "Progress in Optical Fiber Amplifiers," in *Current Trends in Optical Amplifiers and Their Applications*, T. P. Lee, ed. (World Scientific, New Jersey, 1996), see pp. 19–21 and references therein.
12. R. H. Stolen, "Phase-Matched Stimulated Four-Photon Mixing in Silica-Fiber Waveguides," *IEEE Journal of Quantum Electronics* **11**:100–103 (1975).
13. Y. Sun, J. L. Zyskind, and A. K. Srivastava, "Average Inversion Level, Modeling, and Physics of Erbium-Doped Fiber Amplifiers," *IEEE Journal of Selected Topics in Quantum Electronics* **3**:991–1007 (1997).
14. P. C. Becker, N. A. Olsson, and J. R. Simpson, *Erbium-Doped Fiber Amplifiers, Fundamentals and Technology* (Academic Press, San Diego, 1999), pp. 139–140.
15. J.-M. P. Delavaux and J. A. Nagel, "Multi-Stage Erbium-Doped Fiber Amplifier Design," *IEEE Journal of Lightwave Technology* **13**:703–720 (1995).
16. E. Desurvire, *Erbium-Doped Fiber Amplifiers, Principles and Applications* (Wiley-Interscience, New York, 1994), pp. 337–340.
17. E. Desurvire, *Erbium-Doped Fiber Amplifiers, Principles and Applications* (Wiley-Interscience, New York, 1994), p. 238.
18. S. Sudo, "Outline of Optical Fiber Amplifiers," in *Optical Fiber Amplifiers: Materials, Devices, and Applications* (Artech House, Norwood, 1997), see pp. 81–83 and references therein.
19. W. J. Miniscalco, "Erbium-Doped Glasses for Fiber Amplifiers at 1500 nm," *IEEE Journal of Lightwave Technology* **9**:234–250 (1991).
20. S. T. Davey and P. W. France, "Rare-Earth-Doped Fluorozirconate Glass for Fibre Devices," *British Telecom Technical Journal* **7**:58 (1989).
21. C. R. Giles and E. Desurvire, "Modeling Erbium-Doped Fiber Amplifiers," *IEEE Journal of Lightwave Technology* **9**:271–283 (1991).
22. E. Desurvire, "Study of the Complex Atomic Susceptibility of Erbium-Doped Fiber Amplifiers," *IEEE Journal of Lightwave Technology* **8**:1517–1527 (1990).
23. Y. Sun, J. L. Zyskind, and A. K. Srivastava, "Average Inversion Level, Modeling, and Physics of Erbium-Doped Fiber Amplifiers," *IEEE Journal of Selected Topics in Quantum Electronics* **3**:991–1007 (1997).
24. M. Horiguchi, K. Yoshino, M. Shimizu, M. Yamada, and H. Hanafusa, "Erbium-Doped Fiber Amplifiers Pumped in the 660- and 820-nm Bands," *IEEE Journal of Lightwave Technology* **12**:810–820 (1994).
25. E. Desurvire, "Analysis of Gain Difference between Forward- and Backward-Pumped Erbium-Doped Fibers in the Saturation Regime," *IEEE Photonics Technology Letters* **4**:711–713 (1992).
26. M. N. Zervas and R. I. Laming, "Rayleigh Scattering Effect on the Gain Efficiency and Noise of Erbium-Doped Fiber Amplifiers," *IEEE Journal of Quantum Electronics* **31**:469–471 (1995).
27. H. A. Haus, "The Noise Figure of Optical Amplifiers," *IEEE Photonics Technology Letters* **10**:1602–1604 (1998).
28. E. Desurvire, D. Bayart, B. Desthieux, and S. Bigo, *Erbium-Doped Fiber Amplifiers, Devices and System Developments* (Wiley-Interscience, Hoboken, 2002), Chap. 2.
29. D. Bayart, B. Clesca, L. Hamon, and J. L. Beylat, "Experimental Investigation of the Gain Flatness Characteristics for 1.55 μm Erbium-Doped Fluoride Fiber Amplifiers," *IEEE Photonics Technology Letters* **6**:613–615 (1994).

30. E. L. Goldstein, L. Eskildsen, C. Lin, and R. E. Tench, "Multiwavelength Propagation in Light-Wave Systems with Strongly-Inverted Fiber Amplifiers," *IEEE Photonics Technology Letters* **6**:266–269 (1994).
31. C. R. Giles, "Lightwave Applications of Fiber Bragg Gratings," *IEEE Journal of Lightwave Technology* **15**:1391–1404 (1997).
32. J.-Y. Pan, M. A. Ali, A. F. Elrefaie, and R. E. Wagner, "Multiwavelength Fiber Amplifier Cascades with Equalization Employing Mach-Zehnder Optical Filters," *IEEE Photonics Technology Letters* **7**:1501–1503 (1995).
33. P. C. Becker, N. A. Olsson, and J. R. Simpson, op. cit., pp. 285–295.
34. D. M. Machewirth, K. Wei, V. Krasteva, R. Datta, E. Snitzer, and G. H. Sigel, Jr., "Optical Characterization of Pr³⁺ and Dy³⁺ Doped Chalcogenide Glasses," *Journal of Noncrystalline Solids* **213–214**:295–303 (1997).
35. T. J. Whitley, "A Review of Recent System Demonstrations Incorporating Praseodymium-Doped Fluoride Fiber Amplifiers," *IEEE Journal of Lightwave Technology* **13**:744–760 (1995).
36. P. Wang and W. A. Clarkson, "High-Power Single-Mode, Linearly Polarized Ytterbium-Doped Fiber Superfluorescent Source," *Optics Letters* **32**:2605–2607 (2007).
37. Y. Jeong, J. Sahu, D. Payne, and J. Nilsson, "Ytterbium-Doped Large-Core Fiber Laser with 1.36 kW Continuous-Wave End-Pumped Optical Power," *Optics Express* **12**:6088–6092 (2004).
38. P. Russel, "Photonic Crystal Fibers," *Science* **299**:358–362 (2003).
39. O. Schmidt J. Rothhardt, T. Eidam, F. Röser, J. Limpert, A. Tünnermann, K. P. Hansen, C. Jakobsen, and J. Broeng, "Single-Polarization Ultra-Large Mode Area Yb-Doped Photonic Crystal Fiber," *Optics Express* **16**:3918–3923 (2008).
40. G. G. Vienne J. E. Caplen, L. Dong, J. D. Minelly, J. Nilsson, and D. N. Payne, "Fabrication and Characterization of Yb³⁺:Er³⁺ Phosphosilicate Fibers for Lasers," *IEEE Journal of Lightwave Technology* **16**:1990–2001 (1998).
41. M. N. Islam, "Raman Amplifiers for Telecommunications," *IEEE Journal of Selected Topics in Quantum Electronics* **8**:548–559 (2002).
42. J. Bromage, "Raman Amplification for Fiber Communication Systems," *IEEE Journal of Lightwave Technology* **22**:79–93 (2004).
43. T. N. Nielsen, P. B. Hansen, A. J. Stentz, V. M. Aguaro, J. R. Pedrazzani, A. A. Abramov, and R. P. Espindola, "8 × 10 Gb/s 1.3- μ m Unrepeated Transmission over a Distance of 141 km with Raman Post- and Pre-Amplifiers," *IEEE Photonics Technology Letters* **10**:1492–1494 (1998).
44. J. A. Buck, *Fundamentals of Optical Fibers*, 2nd ed., (Wiley-Interscience, Hoboken, 2004), Chap. 8.
45. P. B. Hansen, L. Eskildsen, A. J. Stentz, T. A. Strasser, J. Judkins, J. J. DeMarco, R. Pedrazzani, and D. J. DiGiovanni, "Rayleigh Scattering Limitations in Distributed Raman Pre-Amplifiers," *IEEE Photonics Technology Letters* **10**:159–161 (1998).
46. C. R. S. Fludger, V. Handerek, and R. J. Mears, "Pump to Signal RIN Transfer in Raman Fiber Amplifiers," *IEEE Journal of Lightwave Technology* **19**:1140–1148 (2001).
47. N. Takachio and H. Suzuki, "Application of Raman-Distributed Amplification to WDM Transmission Systems Using 1.55- μ m Dispersion-Shifted Fiber," *IEEE Journal of Lightwave Technology* **19**:60–69 (2001).
48. H. H. Kee, C. R. S. Fludger, and V. Handerek, "Statistical Properties of Polarization Dependent Gain in Fibre Raman Amplifiers," *Optical Fiber Communications Conference*, 2002, paper WB2 (TOPS, vol. 70, Optical Society of America, Washington, D.C.)
49. R. H. Stolen and J. E. Bjorkholm, "Parametric Amplification and Frequency Conversion in Optical Fibers," *IEEE Journal of Quantum Electronics* **18**:1062–1072 (1982).
50. C. J. McKinstrie, S. Radic, and A. R. Chraplyvy, "Parametric Amplifiers Driven by Two Pump Waves," *IEEE Journal of Selected Topics in Quantum Electronics* **8**:538–547 (2002). Erratum, **8**:956.
51. S. Radic, C. J. McKinstrie, R. M. Jopson, J. C. Centanni, Q. Lin, and G. P. Agrawal, "Record Performance of Parametric Amplifier Constructed with Highly-Nonlinear Fibre" *Electronics Letters* **39**:838–839 (2003).
52. J. Hansryd, P. A. Andrekson, M. Westlund, J. Li, and P. O. Hedekvist, "Fiber-Based Optical Parametric Amplifiers and Their Applications," *IEEE Journal of Selected Topics in Quantum Electronics* **8**:506–520 (2002).
53. J. Hansryd and P. A. Andrekson, "Broad-Band Continuous-Wave-Pumped Fiber Optical Parametric Amplifier with 49 dB Gain and Wavelength-Conversion Efficiency," *IEEE Photonics Technology Letters* **13**:194–196 (2001).
54. F. Yaman, Q. Lin, and G. P. Agrawal, "A Novel Design for Polarization-Independent Single-Pump Parametric Amplifiers," *IEEE Photonics Technology Letters* **18**:2335–2337 (2006).

This page intentionally left blank.

DO NOT DUPLICATE

FIBER OPTIC COMMUNICATION LINKS (TELECOM, DATACOM, AND ANALOG)

Casimer DeCusatis

*IBM Corporation
Poughkeepsie, New York*

Guifang Li

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

There are many different applications for fiber optic communication systems, each with its own unique performance requirements. For example, analog communication systems may be subject to different types of noise and interference than digital systems, and consequently require different figures of merit to characterize their behavior. At first glance, telecommunication and data communication systems appear to have much in common, as both use digital encoding of data streams; in fact, both types can share a common network infrastructure. Upon closer examination, however, we find important differences between them. First, datacom systems must maintain a much lower bit error rate (BER), defined as the number of transmission errors per second in the communication link (we will discuss BER in more detail in the following sections). For telecom (voice) communications, the ultimate receiver is the human ear and voice signals have a bandwidth of only about 4 kHz; transmission errors often manifest as excessive static noise such as encountered on a mobile phone, and most users can tolerate this level of fidelity. In contrast, the consequences of even a single bit error to a datacom system can be very serious; critical data such as medical or financial records could be corrupted, or large computer systems could be shut down. Typical telecom systems operate at a BER of about 10^{-9} , compared with about 10^{-12} to 10^{-15} for datacom systems. Another unique requirement of datacom systems is eye safety versus distance trade-offs. Most telecommunication equipment is maintained in a restricted environment and accessible only to personnel trained in the proper handling of high power optical sources. Datacom equipment is maintained in a computer center and must comply with international regulations for inherent eye safety; this limits the amount of optical power which can safely be launched into the fiber, and consequently limits the maximum distances which can be achieved without using repeaters or regenerators. For

the same reason, datacom equipment must be rugged enough to withstand casual use while telecom equipment is more often handled by specially trained service personnel. Telecom systems also tend to make more extensive use of multiplexing techniques, which are only now being introduced into the data center, and more extensive use of optical repeaters.

In the following sections, we will examine the technical requirements for designing fiber optic communication systems suitable for these different environments. We begin by defining some figures of merit to characterize the system performance. Then, concentrating on digital optical communication systems, we will describe how to design an optical link loss budget and how to account for various types of noise sources in the link.

15.1 FIGURES OF MERIT

There are several possible figures of merit which may be used to characterize the performance of an optical communication system. Furthermore, different figures of merit may be more suitable for different applications, such as analog or digital transmission. In this section, we will describe some of the measurements used to characterize the performance of optical communication systems. Even if we ignore the practical considerations of laser eye safety standards, an optical transmitter is capable of launching a limited amount of optical power into a fiber; similarly, there is a limit as to how weak a signal can be detected by the receiver in the presence of noise and interference. Thus, a fundamental consideration in optical communication systems design is the optical link power budget, or the difference between the transmitted and received optical power levels. Some power will be lost due to connections, splices, and bulk attenuation in the fiber. There may also be optical power penalties due to dispersion, modal noise, or other effects in the fiber and electronics. The optical power levels define the signal-to-noise ratio (SNR) at the receiver, which is often used to characterize the performance of analog communication systems. For digital transmission, the most common figure of merit is the bit error rate (BER), defined as the ratio of received bit errors to the total number of transmitted bits. Signal-to-noise ratio is related to the bit error rate by the Gaussian integral

$$\text{BER} = \frac{1}{\sqrt{2\pi}} \int_Q^{\infty} e^{-Q^2/2} dQ \cong \frac{1}{Q\sqrt{2\pi}} e^{-Q^2/2} \quad (1)$$

where Q represents the SNR for simplicity of notation.¹⁻⁴ From Eq. (1), we see that a plot of BER versus received optical power yields a straight line on semilog scale, as illustrated in Fig. 1. Nominally, the slope is about 1.8 dB/decade; deviations from a straight line may indicate the presence of nonlinear or non-Gaussian noise sources. Some effects, such as fiber attenuation, are linear noise sources; they can be overcome by increasing the received optical power, as seen from Fig. 1, subject to constraints on maximum optical power (laser safety) and the limits of receiver sensitivity. There are other types of noise sources, such as mode partition noise or relative intensity noise (RIN), which are independent of signal strength. When such noise is present, no amount of increase in transmitted signal strength will affect the BER; a noise floor is produced, as shown by curve B in Fig. 1. This type of noise can be a serious limitation on link performance. If we plot BER versus receiver sensitivity for increasing optical power, we obtain a curve similar to Fig. 2 which shows that for very high power levels, the receiver will go into saturation. The characteristic “bathtub”-shaped curve illustrates a window of operation with both upper and lower limits on the received power. There may also be an upper limit on optical power due to eye safety considerations.

We can see from Fig. 1 that receiver sensitivity is specified at a given BER, which is often too low to measure directly in a reasonable amount of time (e.g., a 200 Mbit/s link operating at a BER of 10^{-15} will only take one error every 57 days on average, and several hundred errors are recommended for a reasonable BER measurement). For practical reasons, the BER is typically measured at much higher error rates, where the data can be collected more quickly (such as 10^{-4} to 10^{-8}) and then extrapolated

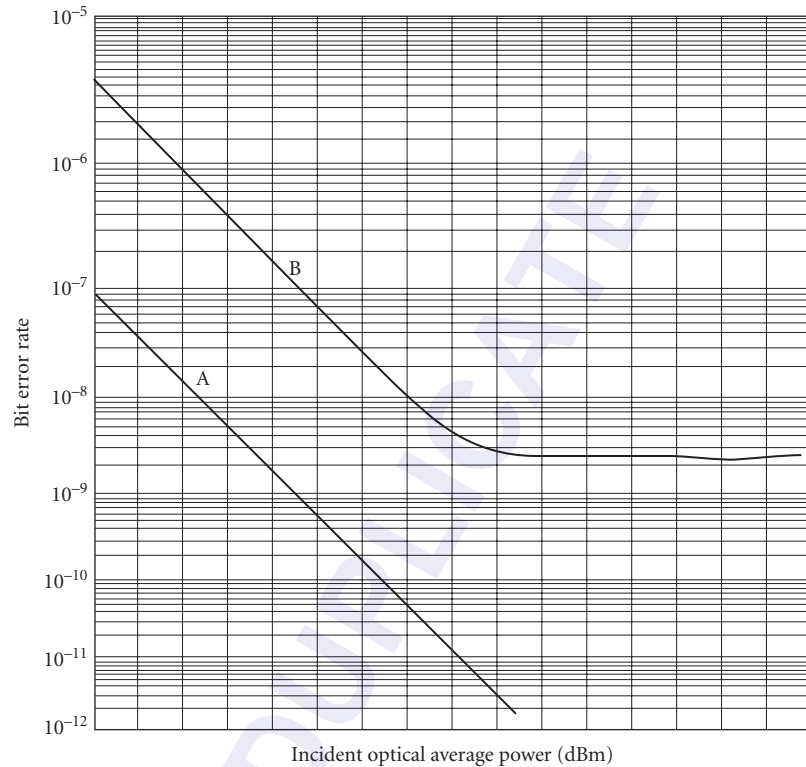


FIGURE 1 Bit error rate as a function of received optical power. Curve A shows typical performance, whereas curve B shows a BER floor.⁵

to find the sensitivity at low BER. This assumes the absence of nonlinear noise floors, as cautioned previously. The relationship between optical input power, in watts, and the BER, is the complimentary Gaussian error function

$$\text{BER} = 1/2 \operatorname{erfc}(P_{\text{out}} - P_{\text{signal}}/\text{RMS noise}) \quad (2)$$

where the error function is an open integral that cannot be solved directly. Several approximations have been developed for this integral, which can be developed into transformation functions that yield a linear least squares fit to the data.¹ The same curve fitting equations can also be used to characterize the eye window performance of optical receivers. Clock position/phase versus BER data are collected for each edge of the eye window; these data sets are then curve fitted with the above expressions to determine the clock position at the desired BER. The difference in the two resulting clock position on either side of the window gives the clear eye opening.¹⁻⁴

In describing Figs. 1 and 2, we have also made some assumptions about the receiver circuit. Most data links are asynchronous, and do not transmit a clock pulse along with the data; instead, a clock is extracted from the incoming data and used to retime the received data stream. We have made the assumption that the BER is measured with the clock at the center of the received data bit; ideally, this is when we compare the signal with a preset threshold to determine if a logical “1” or “0” was sent. When the clock is recovered from a receiver circuit such as a phase lock loop, there is always some uncertainty about the clock position; even if it is centered on the data bit, the relative clock position

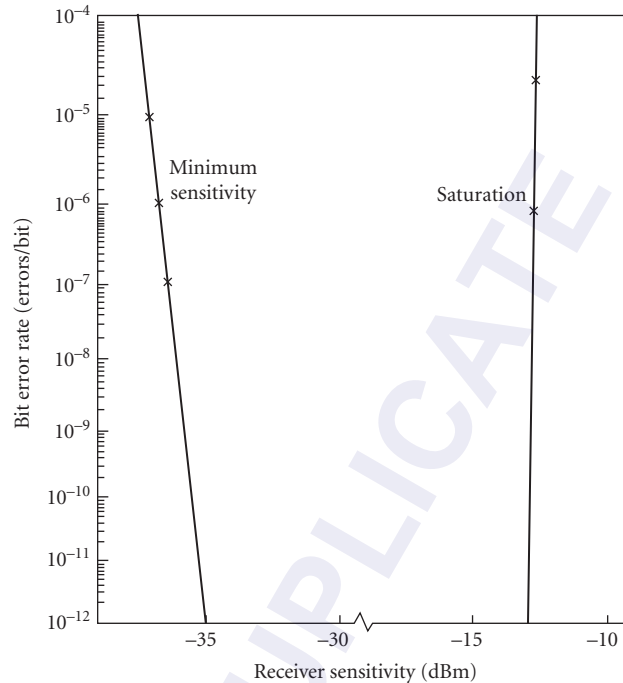


FIGURE 2 Bit error rate as a function of received optical power illustrating range of operation from minimum sensitivity to saturation.

may drift over time. The region of the bit interval in the time domain where the BER is acceptable is called the eyewidth; if the clock timing is swept over the data bit using a delay generator, the BER will degrade near the edges of the eye window. Eyewidth measurements are an important parameter in link design, which will be discussed further in the section on jitter and link budget modeling.

In the design of some analog optical communication systems, as well as some digital television systems (e.g., those based on 64-bit Quadrature Amplitude Modulation), another possible figure of merit is the modulation error ratio (MER). To understand this metric, we will consider the standard definition of the Digital Video Broadcasters (DVB) Measurements Group.⁵ First, the video receiver captures a time record of N received signal coordinate pairs, representing the position of information on a two-dimensional screen. The ideal position coordinates are given by the vector (X_j, Y_j) . For each received symbol, a decision is made as to which symbol was transmitted, and an error vector $(\Delta X_j, \Delta Y_j)$ is defined as the distance from the ideal position to the actual position of the received symbol. The MER is then defined as the sum of the squares of the magnitudes of the ideal symbol vector divided by the sum of the squares of the magnitudes of the symbol error vectors:

$$\text{MER} = 10 \log \frac{\sum_{j=1}^N (X_j^2 + Y_j^2)}{\sum_{j=1}^N (\Delta X_j^2 + \Delta Y_j^2)} \text{dB} \quad (3)$$

when the signal vectors are corrupted by noise, they can be treated as random variables. The denominator in Eq. (3) becomes an estimate of the average power of the error vector (in other words, its second moment) and contains all signal degradation due to noise, reflections, transmitter quadrature errors, etc. If the only significant source of signal degradation is additive white Gaussian noise, then MER and SNR are equivalent. For communication systems which contain other noise sources, MER offers some advantages; in particular, for some digital transmission systems there may be a

very sharp change in BER as a function of SNR (a so-called “cliff effect”) which means that BER alone cannot be used as an early predictor of system failures. MER, on the other hand, can be used to measure signal-to-interference ratios accurately for such systems. Because MER is a statistical measurement, its accuracy is directly related to the number of vectors N used in the computation; an accuracy of 0.14 dB can be obtained with $N = 10,000$, which would require about 2 ms to accumulate at the industry standard digital video rate of 5.057 Msymbols/s.

In order to design a proper optical data link, the contribution of different types of noise sources should be assessed when developing a link budget. There are two basic approaches to link budget modeling. One method is to design the link to operate at the desired BER when all the individual link components assume their worst case performance. This conservative approach is desirable when very high performance is required, or when it is difficult or inconvenient to replace failing components near the end of their useful lifetimes. The resulting design has a high safety margin; in some cases, it may be overdesigned for the required level of performance. Since it is very unlikely that all the elements of the link will assume their worst case performance at the same time, an alternative is to model the link budget statistically. For this method, distributions of transmitter power output, receiver sensitivity, and other parameters are either measured or estimated. They are then combined statistically using an approach such as the Monte Carlo method, in which many possible link combinations are simulated to generate an overall distribution of the available link optical power. A typical approach is the 3-sigma design, in which the combined variations of all link components are not allowed to extend more than 3 standard deviations from the average performance target in either direction. The statistical approach results in greater design flexibility, and generally increased distance compared with a worst-case model at the same BER.

Harmonic Distortions, Intermodulation Distortions, and Dynamic Range

Fiber-optic analog links are in general nonlinear. That is, if the input electrical information is a harmonic signal of frequency f_0 , the output electrical signal will contain the fundamental frequency f_0 as well as high-order harmonics of frequencies nf_0 ($n > 2$). These high-order harmonics comprise the harmonic distortions of analog fiber-optic links.⁶ The nonlinear behavior is caused by nonlinearities in the transmitter, the fiber, and the receiver. The same sources of nonlinearities in the fiber-optic links lead to intermodulation distortions (IMD), which can be best illustrated in a two-tone transmission scenario. If the input electrical information is a superposition of two harmonic signals of frequencies f_1 and f_2 , the output electrical signal will contain second-order intermodulation at frequencies $f_1 + f_2$ and $f_1 - f_2$ as well as third-order intermodulation at frequencies $2f_1 - f_2$ and $2f_2 - f_1$.

Most analog fiber-optic links require bandwidth of less than one octave ($f_{\max} < 2f_{\min}$). As a result harmonic distortions as well as second-order IMD products are not important as they can be filtered out electronically. However, third-order IMD products are in the same frequency range (between f_{\min} and f_{\max}) as the signal itself and therefore appear in the output signal as the spurious response. Thus the linearity of analog fiber-optic links is determined by the level of third-order IMD products. In the case of analog links where third-order IMD is eliminated through linearization circuitry, the lowest odd-order IMD determines the linearity of the link.

To quantify IMD distortions, a two-tone experiment (or simulation) is usually conducted where the input RF powers of the two tones are equal. The linear and nonlinear power transfer functions—the output RF power of each of two input tones and the second or third-order IMD product as a function of the input RF power of each input harmonic signal—are schematically presented in Fig. 3. When plotted on a log-log scale, the fundamental power transfer function should be a line with a slope of unity. The second- (third-) order power transfer function should be a line with a slope of two (three). The intersections of the power transfer functions are called second- and third-order intercept points, respectively. Because of the fixed slopes of the power transfer functions, the intercept points can be calculated from measurements obtained at a single input power level. Suppose at a certain input level, the output power of each of the two fundamental tones, the second-order IMD product and third-order

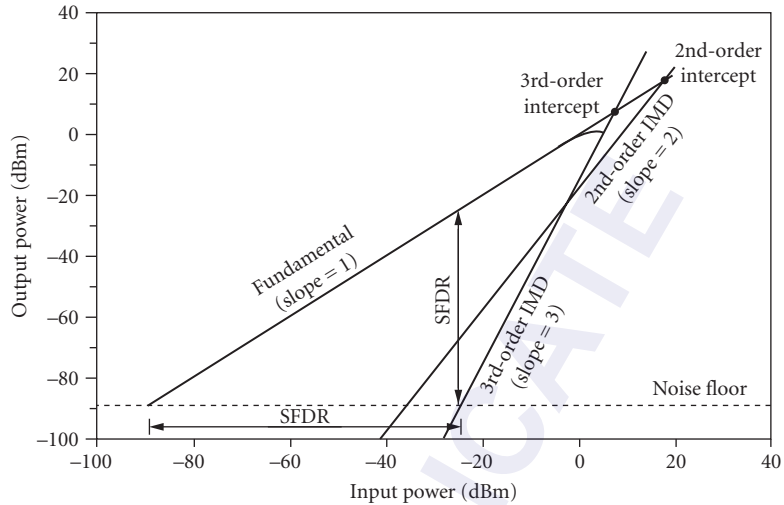


FIGURE 3 Intermodulation and dynamic range of analog fiberoptic links.

IMD products are P_1 , P_2 , and P_3 , respectively. When the power levels are in units of dB or dBm, the second-order and third-order intercept points are

$$IP_2 = 2P_1 - P_2 \quad (4)$$

and

$$IP_3 = (3P_1 - P_3)/2 \quad (5)$$

The dynamic range is a measure of the ability of an analog fiber-optic link to faithfully transmit signals at various power levels. At the low input power end, the analog link can fail due to insufficient power level so that the output power is below the noise level. At the high input power end, the analog link can fail due to the fact that the IMD products become the dominant source of signal degradation. In terms of the output power, the dynamic range (of the output power) is defined as the ratio of the fundamental output to the noise power. However, it should be noted that the third-order IMD products increase three times faster than the fundamental signal. After the third-order IMD products exceed the noise floor, the ratio of the fundamental output to the noise power is meaningless as the dominant degradation of the output signal comes from IMD products. So a more meaningful definition of the dynamic range is the so-called spurious-free dynamic range (SFDR),^{6,7} which is the ratio of the fundamental output to the noise power at the point where the IMD products are at the noise level. The spurious-free dynamic range is then practically the maximum dynamic range. Since the noise floor depends on the bandwidth of interest, the unit for SFDR should be dB-Hz^{2/3}. The dynamic range decreases as the bandwidth of the system is increased. The spurious-free dynamic range is also often defined with reference to the input power, which corresponds to SFDR with reference to the output power if there is no gain compression.

15.2 LINK BUDGET ANALYSIS: INSTALLATION LOSS

It is convenient to break down the link budget into two areas: installation loss and available power. Installation or DC loss refers to optical losses associated with the fiber cable plant, such as connector loss, splice loss, and bandwidth considerations. Available optical power is the difference between the

transmitter output and receiver input powers, minus additional losses due to optical noise sources on the link (also known as AC losses). With this approach, the installation loss budget may be treated statistically and the available power budget as worst case. First, we consider the installation loss budget, which can be broken down into three areas, namely, transmission loss, fiber attenuation as a function of wavelength, and connector or splice losses.

Transmission Loss

Transmission loss is perhaps the most important property of an optical fiber; it affects the link budget and maximum unrepeated distance. Since the maximum optical power launched into an optical fiber is determined by international laser eye safety standards,⁸ the number and separation between optical repeaters and regenerators is largely determined by this loss. The mechanisms responsible for this loss include material absorption as well as both linear and nonlinear scattering of light from impurities in the fiber.¹⁻⁵ Typical loss for single-mode optical fiber is about 2 to 3 dB/km near 800-nm wavelength, 0.5 dB/km near 1300 nm, and 0.25 dB/km near 1550 nm. Multimode fiber loss is slightly higher, and bending loss will only increase the link attenuation further.

Attenuation versus Wavelength

Since fiber loss varies with wavelength, changes in the source wavelength or use of sources with a spectrum of wavelengths will produce additional loss. Transmission loss is minimized near the 1550-nm wavelength band, which unfortunately does not correspond with the dispersion minimum at around 1310 nm. An accurate model for fiber loss as a function of wavelength has been developed by Walker;⁹ this model accounts for the effects of linear scattering, macrobending, and material absorption due to ultraviolet and infrared band edges, hydroxide (OH) absorption, and absorption from common impurities such as phosphorous. Using this model, it is possible to calculate the fiber loss as a function of wavelength for different impurity levels; the fiber properties can be specified along with the acceptable wavelength limits of the source to limit the fiber loss over the entire operating wavelength range. Design tradeoffs are possible between center wavelength and fiber composition to achieve the desired result. Typical loss due to wavelength-dependent attenuation for laser sources on single-mode fiber can be held below 0.1 dB/km.

Connector and Splice Losses

There are also installation losses associated with fiber optic connectors and splices; both of these are inherently statistical in nature and can be characterized by a Gaussian distribution. There are many different kinds of standardized optical connectors, some of which have been discussed previously; some industry standards also specify the type of optical fiber and connectors suitable for a given application.¹⁰ There are also different models which have been published for estimating connection loss due to fiber misalignment;^{11,12} most of these treat loss due to misalignment of fiber cores, offset of fibers on either side of the connector, and angular misalignment of fibers. The loss due to these effects is then combined into an overall estimate of the connector performance. There is no general model available to treat all types of connectors, but typical connector loss values average approximately 0.5 dB worst case for multimode, slightly higher for singlemode (see Table 1).

Optical splices are required for longer links, since fiber is usually available in spools of 1 to 5 km, or to repair broken fibers. There are two basic types, mechanical splices (which involve placing the two fiber ends in a receptacle that holds them close together, usually with epoxy) and the more commonly used fusion splices (in which the fiber are aligned, then heated sufficiently to fuse the two ends together). Typical splice loss values are given in Table 1.

TABLE 1 Typical Cable Plant Optical Losses⁵

Component	Description	Size (μm)	Mean Loss	Variance (dB^2)
Connector*	Physical contact	62.5–62.5	0.40 dB	0.02
		50.0–50.0	0.40 dB	0.02
		9.0–9.0 [†]	0.35 dB	0.06
		62.5–50.0	2.10 dB	0.12
Connector*	Nonphysical contact (multimode only)	50.0–62.5	0.00 dB	0.01
		62.5–62.5	0.70 dB	0.04
		50.0–50.0	0.70 dB	0.04
		62.5–50.0	2.40 dB	0.12
Splice	Mechanical	50.0–62.5	0.30 dB	0.01
		62.5–62.5	0.15 dB	0.01
		50.0–50.0	0.15 dB	0.01
Splice	Fusion	9.0–9.0 [†]	0.15 dB	0.01
		62.5–62.5	0.40 dB	0.01
		50.0–50.0	0.40 dB	0.01
Cable	IBM multimode jumper	62.5	1.75 dB/km	NA
	IBM multimode jumper	50.0	3.00 dB/km at 850 nm	NA
	IBM single-mode jumper	9.0	0.8 dB/km	NA
	Trunk	62.5	1.00 dB/km	NA
	Trunk	50.0	0.90 dB/km	NA
	Trunk	9.0	0.50 dB/km	NA

*The connector loss value is typical when attaching identical connectors. The loss can vary significantly if attaching different connector types.

[†]Single-mode connectors and splices must meet a minimum return loss specification of 28 dB.

15.3 LINK BUDGET ANALYSIS: OPTICAL POWER PENALTIES

Next, we will consider the assembly loss budget, which is the difference between the transmitter output and receiver input powers, allowing for optical power penalties due to noise sources in the link. We will follow the standard convention in the literature of assuming a digital optical communication link which is best characterized by its BER. Contributing factors to link performance include the following:

- Dispersion (modal and chromatic) or intersymbol interference
- Mode partition noise
- Mode hopping
- Extinction ratio
- Multipath interference
- Relative intensity noise (RIN)
- Timing jitter
- Radiation induced darkening
- Modal noise

Higher order, nonlinear effects including Stimulated Raman and Brillouin scattering and frequency chirping will be discussed elsewhere.

Dispersion

The most important fiber characteristic after transmission loss is dispersion, or intersymbol interference. This refers to the broadening of optical pulses as they propagate along the fiber. As pulses broaden, they tend to interfere with adjacent pulses; this limits the maximum achievable data rate. In multimode fibers, there are two dominant kinds of dispersion, modal and chromatic. Modal dispersion refers to the fact that different modes will travel at different velocities and cause pulse broadening. The fiber's modal bandwidth in units of MHz-km, is specified according to the expression

$$BW_{\text{modal}} = BW_1 / L^\gamma \quad (6)$$

where BW_{modal} is the modal bandwidth for a length L of fiber, BW_1 is the manufacturer-specified modal bandwidth of a 1-km section of fiber, and γ is a constant known as the modal bandwidth concatenation length scaling factor. The term γ usually assumes a value between 0.5 and 1, depending on details of the fiber manufacturing and design as well as the operating wavelength; it is conservative to take $\gamma = 1.0$. Modal bandwidth can be increased by mode mixing, which promotes the interchange of energy between modes to average out the effects of modal dispersion. Fiber splices tend to increase the modal bandwidth, although it is conservative to discard this effect when designing a link.

The other major contribution is chromatic dispersion BW_{chrom} which occurs because different wavelengths of light propagate at different velocities in the fiber. For multimode fiber, this is given by an empirical model of the form

$$BW_{\text{chrom}} = \frac{L^{\gamma_c}}{\sqrt{\lambda_w (a_o + a_1 |\lambda_c - \lambda_{\text{eff}}|)}} \quad (7)$$

where L is the fiber length in km; λ_c is the center wavelength of the source in nm; λ_w is the source FWHM spectral width in nm; γ_c is the chromatic bandwidth length scaling coefficient, a constant; λ_{eff} is the effective wavelength, which combines the effects of the fiber zero dispersion wavelength and spectral loss signature; and the constants a_1 and a_o are determined by a regression fit of measured data. From Ref. 13, the chromatic bandwidth for 62.5/125- μm fiber is empirically given by

$$BW_{\text{chrom}} = \frac{10^4 L^{-0.69}}{\sqrt{\lambda_w (1.1 + 0.0189 |\lambda_c - 1370|)}} \quad (8)$$

For this expression, the center wavelength was 1335 nm and λ_{eff} was chosen midway between λ_c and the water absorption peak at 1390 nm; although λ_{eff} was estimated in this case, the expression still provides a good fit to the data. For 50/125- μm fiber, the expression becomes

$$BW_{\text{chrom}} = \frac{10^4 L^{-0.65}}{\sqrt{\lambda_w (1.01 + 0.0177 |\lambda_c - 1330|)}} \quad (9)$$

For this case, λ_c was 1313 nm and the chromatic bandwidth peaked at $\lambda_{\text{eff}} = 1330$ nm. Recall that this is only one possible model for fiber bandwidth.¹ The total bandwidth capacity of multimode fiber BW_t is obtained by combining the modal and chromatic dispersion contributions, according to

$$\frac{1}{BW_t^2} = \frac{1}{BW_{\text{chrom}}^2} + \frac{1}{BW_{\text{modal}}^2} \quad (10)$$

Once the total bandwidth is known, the dispersion penalty can be calculated for a given data rate. One expression for the dispersion penalty in decibel is

$$P_d = 1.22 \left[\frac{\text{bit rate (Mb/s)}}{BW_t \text{ (MHz)}} \right]^2 \quad (11)$$

For typical telecommunication grade fiber, the dispersion penalty for a 20-km link is about 0.5 dB.

Dispersion is usually minimized at wavelengths near 1310 nm; special types of fiber have been developed which manipulate the index profile across the core to achieve minimal dispersion near 1550 nm, which is also the wavelength region of minimal transmission loss. Unfortunately, this dispersion-shifted fiber suffers from some practical drawbacks, including susceptibility to certain kinds of nonlinear noise and increased interference between adjacent channels in a wavelength multiplexing environment. There is a new type of fiber which minimizes dispersion while reducing the unwanted crosstalk effects, called dispersion optimized fiber. By using a very sophisticated fiber profile, it is possible to minimize dispersion over the entire wavelength range from 1300 nm to 1550 nm, at the expense of very high loss (around 2 dB/km); this is known as dispersion flattened fiber. Yet another approach is called dispersion compensating fiber; this fiber is designed with negative dispersion characteristics, so that when used in series with conventional fiber it will offset the normal fiber dispersion. Dispersion compensating fiber has a much narrower core than standard singlemode fiber, which makes it susceptible to nonlinear effects; it is also birefringent and suffers from polarization mode dispersion, in which different states of polarized light propagate with very different group velocities. Note that standard singlemode fiber does not preserve the polarization state of the incident light; there is yet another type of specialty fiber, with asymmetric core profiles, capable of preserving the polarization of incident light over long distances.

By definition, single-mode fiber does not suffer modal dispersion. Chromatic dispersion is an important effect, though, even given the relatively narrow spectral width of most laser diodes. The dispersion of single-mode fiber corresponds to the first derivative of group velocity τ_g with respect to wavelength, and is given by

$$D = \frac{d\tau_g}{d\lambda} = \frac{S_o}{4} \left(\lambda_c - \frac{\lambda_o^4}{\lambda_c^3} \right) \quad (12)$$

where D is the dispersion in ps/(km-nm) and λ_c is the laser center wavelength. The fiber is characterized by its zero dispersion wavelength, λ_o , and zero dispersion slope, S_o . Usually, both center wavelength and zero dispersion wavelength are specified over a range of values; it is necessary to consider both upper and lower bounds in order to determine the worst case dispersion penalty. This can be seen from Fig. 4 which plots D versus wavelength for some typical values of λ_o and λ_c ; the largest absolute value of D occurs at the extremes of this region. Once the dispersion is determined,

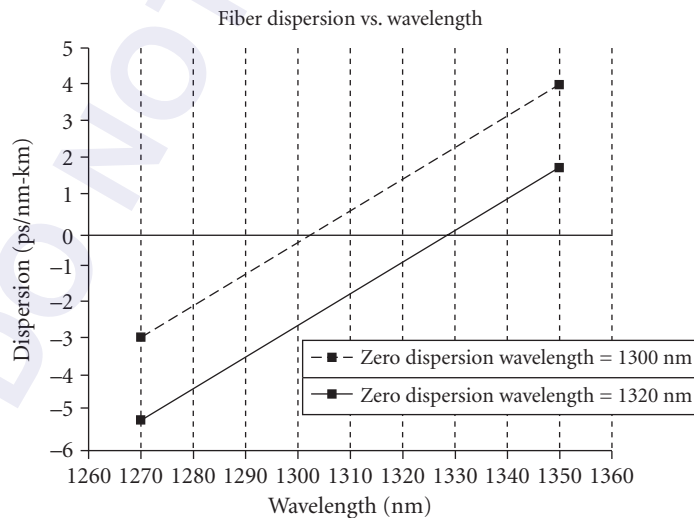


FIGURE 4 Single-mode fiber dispersion as a function of wavelength.⁵

the intersymbol interference penalty as a function of link length L can be determined to a good approximation from a model proposed by Agrawal:¹⁴

$$P_d = 5 \log(1 + 2\pi(BD\Delta\lambda)^2 L^2) \tag{13}$$

where B is the bit rate and $\Delta\lambda$ is the root-mean-square (RMS) spectral width of the source. By maintaining a close match between the operating and zero dispersion wavelengths, this penalty can be kept to a tolerable 0.5 to 1.0 dB in most cases.

Mode Partition Noise

Group velocity dispersion contributes to another optical penalty, which remains the subject of continuing research, mode partition noise and mode hopping. This penalty is related to the properties of a Fabry-Perot type laser diode cavity; although the total optical power output from the laser may remain constant, the optical power distribution among the laser's longitudinal modes will fluctuate. This is illustrated by the model depicted in Fig. 5; when a laser diode is directly modulated with

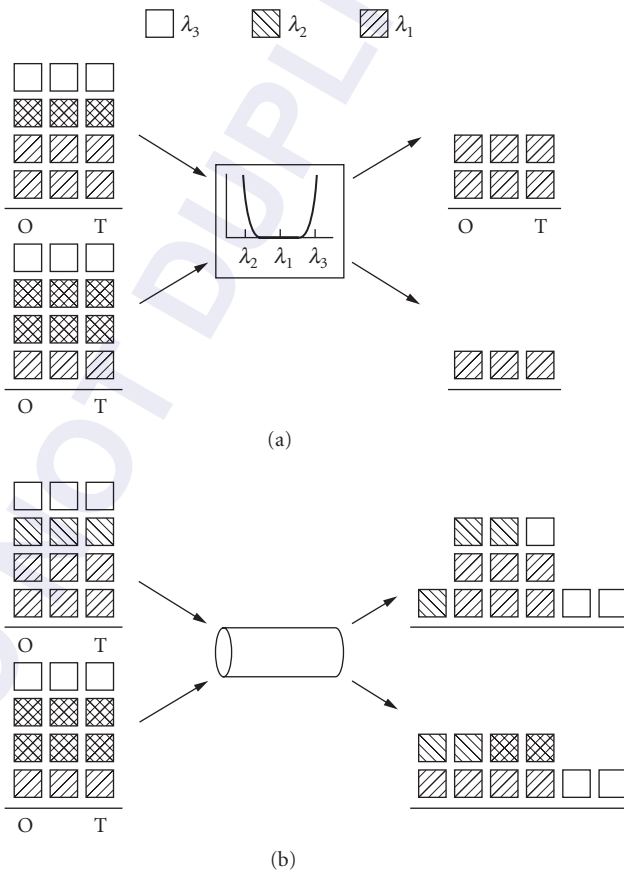


FIGURE 5 Model for mode partition noise; an optical source emits a combination of wavelengths, illustrated by different color blocks: (a) wavelength-dependent loss and (b) chromatic dispersion.

injection current, the total output power stays constant from pulse to pulse; however, the power distribution among several longitudinal modes will vary between pulses. We must be careful to distinguish this behavior of the instantaneous laser spectrum, which varies with time, from the time-averaged spectrum which is normally observed experimentally. The light propagates through a fiber with wavelength-dependent dispersion or attenuation, which deforms the pulse shape. Each mode is delayed by a different amount due to group velocity dispersion in the fiber; this leads to additional signal degradation at the receiver, in addition to the intersymbol interference caused by chromatic dispersion alone, discussed earlier. This is known as mode partition noise; it is capable of generating bit error rate floors, such that additional optical power into the receiver will not improve the link BER. This is because mode partition noise is a function of the laser spectral fluctuations and wavelength-dependent dispersion of the fiber, so the signal-to-noise ratio due to this effect is independent of the signal power. The power penalty due to mode partition noise was first calculated by Ogawa¹⁵ as

$$P_{\text{mp}} = 5 \log(1 - Q^2 \sigma_{\text{mp}}^2) \quad (14)$$

where

$$\sigma_{\text{mp}}^2 = \frac{1}{2} k^2 (\pi B)^4 [A_1^4 \Delta \lambda^4 + 42 A_1^2 A_2^2 \Delta \lambda^6 + 48 A_2^4 \Delta \lambda^8] \quad (15)$$

$$A_1 = DL \quad (16)$$

and

$$A_2 = \frac{A_1}{2(\lambda_c - \lambda_o)} \quad (17)$$

The mode partition coefficient k is a number between 0 and 1 which describes how much of the optical power is randomly shared between modes; it summarizes the statistical nature of mode partition noise. According to Ogawa, k depends on the number of interacting modes and rms spectral width of the source, the exact dependence being complex. However, subsequent work has shown¹⁶ that Ogawa's model tends to underestimate the power penalty due to mode partition noise because it does not consider the variation of longitudinal mode power between successive baud periods, and because it assumes a linear model of chromatic dispersion rather than the nonlinear model given in the above equation. A more detailed model has been proposed by Campbell,¹⁷ which is general enough to include effects of the laser diode spectrum, pulse shaping, transmitter extinction ratio, and statistics of the data stream. While Ogawa's model assumed an equiprobable distribution of zeros and ones in the data stream, Campbell showed that mode partition noise is data dependent as well. Recent work based on this model¹⁸ has re-derived the signal variance:

$$\sigma_{\text{mp}}^2 = E_{\text{av}} (\sigma_o^2 + \sigma_{+1}^2 + \sigma_{-1}^2) \quad (18)$$

where the mode partition noise contributed by adjacent baud periods is defined by

$$\sigma_{+1}^2 + \sigma_{-1}^2 = \frac{1}{2} k^2 (\pi B)^4 [1.25 A_1^4 \Delta \lambda^4 + 40.95 A_1^2 A_2^2 \Delta \lambda^6 + 50.25 A_2^4 \Delta \lambda^8] \quad (19)$$

and the time-average extinction ratio $E_{\text{av}} = 10 \log(P_1/P_0)$, where P_1, P_0 represent the optical power by a 1 and 0, respectively. If the operating wavelength is far away from the zero dispersion wavelength, the noise variance simplifies to

$$\sigma_{\text{mp}}^2 = 2.25 \frac{k^2}{2} E_{\text{av}} (1 - e^{-\beta l^2})^2 \quad (20)$$

which is valid provided that,

$$\beta = (\pi BD\Delta\lambda)^2 \ll 1 \quad (21)$$

Many diode lasers exhibit mode hopping or mode splitting in which the spectrum appears to split optical power between 2 or 3 modes for brief periods of time. The exact mechanism is not fully understood, but stable Gaussian spectra are generally only observed for CW operation and temperature stabilized lasers. During these mode hops the above theory does not apply since the spectrum is non-Gaussian, and the model will overpredict the power penalty; hence, it is not possible to model mode hops as mode partitioning with $k = 1$. There is no currently published model describing a treatment of mode hopping noise, although recent papers¹⁹ suggest approximate calculations based on the statistical properties of the laser cavity. In a practical link, some amount of mode hopping is probably unavoidable as a contributor to burst noise; empirical testing of link hardware remains the only reliable way to reduce this effect. A practical rule of thumb is to keep the mode partition noise penalty less than 1.0 dB maximum, provided that this penalty is far away from any noise floors.

Extinction Ratio

The receiver extinction ratio also contributes directly to the link penalties. The receiver BER is a function of the modulated AC signal power; if the laser transmitter has a small extinction ratio, the DC component of total optical power is significant. Gain or loss can be introduced in the link budget if the extinction ratio at which the receiver sensitivity is measured differs from the worst case transmitter extinction ratio. If the extinction ratio E_t at the transmitter is defined as the ratio of optical power when a one is transmitted versus when a zero is transmitted,

$$E_t = \frac{\text{Power}(1)}{\text{Power}(0)} \quad (22)$$

then we can define a modulation index at the transmitter M_t according to

$$M_t = \frac{E_t - 1}{E_t + 1} \quad (23)$$

Similarly, we can measure the linear extinction ratio at the optical receiver input and define a modulation M_r . The extinction ratio penalty is given by

$$P_{er} = -10 \log \left(\frac{M_t}{M_r} \right) \quad (24)$$

where the subscripts T and R refer to specifications for the transmitter and receiver, respectively. Usually, the extinction ratio is specified to be the same at the transmitter and receiver, and is large enough so that there is no power penalty due to extinction ratio effects.

Multipath Interference

Another important property of the optical link is the amount of reflected light from the fiber end-faces which returns up the link back into the transmitter. Whenever there is a connection or splice in the link, some fraction of the light is reflected back; each connection is thus a potential noise generator, since the reflected fields can interfere with one another to create noise in the detected optical signal. The phenomenon is analogous to the noise caused by multiple atmospheric reflections of radio waves, and is known as *multipath interference noise*. To limit this noise, connectors and splices are specified with a minimum return loss. If there are a total of N reflection points in a link and the geometric mean of the connector reflections is alpha, then based on the model of Duff et al.²⁰ the

power penalty due to multipath interference (adjusted for bit error rate and bandwidth) is closely approximated by

$$P_{\text{mpi}} = 10 \log(1 - 0.7Na) \quad (25)$$

Multipath noise can usually be reduced well below 0.5 dB with available connectors, whose return loss is often better than 25 dB.

Relative Intensity Noise

Stray light reflected back into a Fabry-Perot type laser diode gives rise to intensity fluctuations in the laser output. This is a complicated phenomena, strongly dependent on the type of laser; it is called either reflection-induced intensity noise or relative intensity noise (RIN). This effect is important since it can also generate BER floors. The power penalty due to RIN is the subject of ongoing research; since the reflected light is measured at a specified signal level, RIN is data dependent although it is independent of link length. Since many laser diodes are packaged in windowed containers, it is difficult to correlate the RIN measurements on an unpackaged laser with those of a commercial product. There have been several detailed attempts to characterize RIN;^{21,22} typically, the RIN noise is assumed Gaussian in amplitude and uniform in frequency over the receiver bandwidth of interest. The RIN value is specified for a given laser by measuring changes in the optical power when a controlled amount of light is fed back into the laser; it is signal dependent, and is also influenced by temperature, bias voltage, laser structure, and other factors which typically influence laser output power.²² If we assume that the effect of RIN is to produce an equivalent noise current at the receiver, then the additional receiver noise σ_r may be modeled as

$$\sigma_r = \gamma^2 S^2 g B \quad (26)$$

where S is the signal level during a bit period, B is the bit rate, and g is a noise exponent which defines the amount of signal-dependent noise. If $g = 0$, noise power is independent of the signal, while for $g = 1$ noise power is proportional to the square of the signal strength. The coefficient γ is given by

$$\gamma^2 = S_i^{2(1-g)} 10^{(\text{RIN}_i/10)} \quad (27)$$

where RIN_i is the measured RIN value at the average signal level S_i , including worst case backreflection conditions and operating temperatures. The Gaussian BER probability due to the additional RIN noise current is given by

$$P_{\text{error}} = \frac{1}{2} \left[P_e^1 \left(\frac{S_1 - S_o}{2\sigma_1} \right) + P_e^0 \left(\frac{S_1 - S_o}{2\sigma_o} \right) \right] \quad (28)$$

where σ_1 , σ_o represent the total noise current during transmission of a digital 1 and 0, respectively, and P_e^1 , P_e^0 are the probabilities of error during transmission of a 1 or 0, respectively. The power penalty due to RIN may then be calculated by determining the additional signal power required to achieve the same BER with RIN noise present as without the RIN contribution. One approximation for the RIN power penalty is given by

$$P_{\text{rin}} = -5 \log \left[1 - Q^2(\text{BW})(1 + M_r)^{2g} (10^{\text{RIN}/10}) \left(\frac{1}{M_r} \right)^2 \right] \quad (29)$$

where the RIN value is specified in dB/Hz, BW is the receiver bandwidth, M_r is the receiver modulation index, and the exponent g is a constant varying between 0 and 1 which relates the magnitude of RIN noise to the optical power level. The maximum RIN noise penalty in a link can usually be kept below 0.5 dB.

Jitter

Although it is not strictly an optical phenomena, another important area in link design deals with the effects of timing jitter on the optical signal. In a typical optical link, a clock is extracted from the incoming data signal which is used to retime and reshape the received digital pulse; the received pulse is then compared with a threshold to determine if a digital 1 or 0 was transmitted. So far, we have discussed BER testing with the implicit assumption that the measurement was made in the center of the received data bit; to achieve this, a clock transition at the center of the bit is required. When the clock is generated from a receiver timing recovery circuit, it will have some variation in time and the exact location of the clock edge will be uncertain. Even if the clock is positioned at the center of the bit, its position may drift over time. There will be a region of the bit interval, or eye, in the time domain where the BER is acceptable; this region is defined as the eyewidth.¹⁻³ Eyewidth measurements are an important parameter for evaluation of fiber optic links; they are intimately related to the BER, as well as the acceptable clock drift, pulse width distortion, and optical power. At low optical power levels, the receiver signal-to-noise ratio is reduced; increased noise causes amplitude variations in the received signal. These amplitude variations are translated into time domain variations in the receiver decision circuitry, which narrows the eyewidth. At the other extreme, an optical receiver may become saturated at high optical power, reducing the eyewidth and making the system more sensitive to timing jitter. This behavior results in the typical “bathtub” curve shown in Fig. 2; for this measurement, the clock is delayed from one end of the bit cell to the other, with the BER calculated at each position. Near the ends of the cell, a large number of errors occur; toward the center of the cell, the BER decreases to its true value. The eye opening may be defined as the portion of the eye for which the BER remains constant; pulse width distortion occurs near the edges of the eye, which denotes the limits of the valid clock timing. Uncertainty in the data pulse arrival times causes errors to occur by closing the eye window and causing the eye pattern to be sampled away from the center. This is one of the fundamental problems of optical and digital signal processing, and a large body of work has been done in this area.^{23,24} In general, multiple jitter sources will be present in a link; these will tend to be uncorrelated. However, jitter on digital signals, especially resulting from a cascade of repeaters, may be coherent.

International standards on jitter were first published by the CCITT (Central Commission for International Telephony and Telegraphy, now known as the International Telecommunications Union, or ITU). This standards body has adopted a definition of jitter²⁴ as short-term variations of the significant instants (rising or falling edges) of a digital signal from their ideal position in time. Longer-term variations are described as wander; in terms of frequency, the distinction between jitter and wander is somewhat unclear. The predominant sources of jitter include the following:

- Phase noise in receiver clock recovery circuits, particularly crystal-controlled oscillator circuits; this may be aggravated by filters or other components which do not have a linear phase response. Noise in digital logic resulting from restricted rise and fall times may also contribute to jitter.
- Imperfect timing recovery in digital regenerative repeaters, which is usually dependent on the data pattern.
- Different data patterns may contribute to jitter when the clock recovery circuit of a repeater attempts to recover the receive clock from inbound data. Data pattern sensitivity can produce as much as 0.5-dB penalty in receiver sensitivity. Higher data rates are more susceptible (>1 Gbit/s); data patterns with long run lengths of 1s or 0s, or with abrupt phase transitions between consecutive blocks of 1s and 0s, tend to produce worst case jitter.
- At low optical power levels, the receiver signal-to-noise ratio Q is reduced; increased noise causes amplitude variations in the signal, which may be translated into time domain variations by the receiver circuitry.
- Low frequency jitter, also called wander, resulting from instabilities in clock sources and modulation of transmitters.
- Very low frequency jitter caused by variations in the propagation delay of fibers, connectors, etc., typically resulting from small temperature variations. (This can make it especially difficult to perform long-term jitter measurements.)

In general, jitter from each of these sources will be uncorrelated; jitter related to modulation components of the digital signal may be coherent, and cumulative jitter from a series of repeaters or regenerators may also contain some well correlated components.

There are several parameters of interest in characterizing jitter performance. Jitter may be classified as either random or deterministic, depending on whether it is associated with pattern-dependent effects; these are distinct from the duty cycle distortion which often accompanies imperfect signal timing. Each component of the optical link (data source, serializer, transmitter, encoder, fiber, receiver, retiming/clock recovery/deserialization, decision circuit) will contribute some fraction of the total system jitter. If we consider the link to be a “black box” (but not necessarily a linear system) then we can measure the level of output jitter in the absence of input jitter; this is known as the “intrinsic jitter” of the link. The relative importance of jitter from different sources may be evaluated by measuring the spectral density of the jitter. Another approach is the maximum tolerable input jitter (MTIJ) for the link. Finally, since jitter is essentially a stochastic process, we may attempt to characterize the jitter transfer function (JTF) of the link, or estimate the probability density function of the jitter. When multiple traces occur at the edges of the eye, this can indicate the presence of data dependent jitter or duty cycle distortion; a histogram of the edge location will show several distinct peaks. This type of jitter can indicate a design flaw in the transmitter or receiver. By contrast, random jitter typically has a more Gaussian profile and is present to some degree in all data links.

The problem of jitter accumulation in a chain of repeaters becomes increasingly complex; however, we can state some general rules of thumb. It has been shown²⁵ that jitter can be generally divided into two components, one due to repetitive patterns and one due to random data. In receivers with phase-lock loop timing recovery circuits, repetitive data patterns will tend to cause jitter accumulation, especially for long run lengths. This effect is commonly modeled as a second-order receiver transfer function. Jitter will also accumulate when the link is transferring random data; jitter due to random data is of two types, systematic and random. The classic model for systematic jitter accumulation in cascaded repeaters was published by Byrne.²⁶ The Byrne model assumes cascaded identical timing recovery circuits, and then the systematic and random jitter can be combined as rms quantities so that total jitter due to random jitter may be obtained. This model has been generalized to networks consisting of different components,²⁷ and to nonidentical repeaters.²⁸ Despite these considerations, for well designed practical networks the basic results of the Byrne model remain valid for N nominally identical repeaters transmitting random data; systematic jitter accumulates in proportion to $N^{1/2}$ and random jitter accumulates in proportion to $N^{1/4}$. For most applications the maximum timing jitter should be kept below about 30 percent of the maximum receiver eye opening.

Modal Noise

An additional effect of lossy connectors and splices is modal noise. Because high capacity optical links tend to use highly coherent laser transmitters, random coupling between fiber modes causes fluctuations in the optical power coupled through splices and connectors; this phenomena is known as *modal noise*.²⁹ As one might expect, modal noise is worst when using laser sources in conjunction with multimode fiber; recent industry standards have allowed the use of short-wave lasers (750 to 850 nm) on 50 μm fiber which may experience this problem. Modal noise is usually considered to be nonexistent in single-mode systems. However, modal noise in single-mode fibers can arise when higher-order modes are generated at imperfect connections or splices. If the lossy mode is not completely attenuated before it reaches the next connection, interference with the dominant mode may occur. The effects of modal noise have been modeled previously,²⁹ assuming that the only significant interaction occurs between the LP_{01} and LP_{11} modes for a sufficiently coherent laser. For N sections of fiber, each of length L in a single-mode link, the worst case sigma for modal noise can be given by

$$\sigma_m = \sqrt{2N\eta(1 - \eta)}e^{-aL} \quad (30)$$

where a is the attenuation coefficient of the LP₁₁ mode, and η is the splice transmission efficiency, given by

$$\eta = 10^{-(\eta_0/10)} \quad (31)$$

where η_0 is the mean splice loss (typically, splice transmission efficiency will exceed 90%). The corresponding optical power penalty due to modal noise is given by

$$P = -5 \log(1 - Q^2 \sigma_m^2) \quad (32)$$

where Q corresponds to the desired BER. This power penalty should be kept to less than 0.5 dB.

Radiation Induced Loss

Another important environmental factor as mentioned earlier is exposure of the fiber to ionizing radiation damage. There is a large body of literature concerning the effects of ionizing radiation on fiber links.^{30,31} There are many factors which can affect the radiation susceptibility of optical fiber, including the type of fiber, type of radiation (gamma radiation is usually assumed to be representative), total dose, dose rate (important only for higher exposure levels), prior irradiation history of the fiber, temperature, wavelength, and data rate. Optical fiber with a pure silica core is least susceptible to radiation damage; however, almost all commercial fiber is intentionally doped to control the refractive index of the core and cladding, as well as dispersion properties. Trace impurities are also introduced which become important only under irradiation; among the most important are Ge dopants in the core of graded index (GRIN) fibers, in addition to F, Cl, P, B, OH content, and the alkali metals. In general, radiation sensitivity is worst at lower temperatures, and is also made worse by hydrogen diffusion from materials in the fiber cladding. Because of the many factors involved, there does not exist a comprehensive theory to model radiation damage in optical fibers. The basic physics of the interaction has been described;^{30,31} there are two dominant mechanisms, radiation induced darkening and scintillation. First, high energy radiation can interact with dopants, impurities, or defects in the glass structure to produce color centers which absorb strongly at the operating wavelength. Carriers can also be freed by radiolytic or photochemical processes; some of these become trapped at defect sites, which modifies the band structure of the fiber and causes strong absorption at infrared wavelengths. This radiation-induced darkening increases the fiber attenuation; in some cases, it is partially reversible when the radiation is removed, although high-levels or prolonged exposure will permanently damage the fiber. A second effect is caused if the radiation interacts with impurities to produce stray light, or scintillation. This light is generally broad-band, but will tend to degrade the BER at the receiver; scintillation is a weaker effect than radiation-induced darkening. These effects will degrade the BER of a link; they can be prevented by shielding the fiber, or partially overcome by a third mechanism, photobleaching. The presence of intense light at the proper wavelength can partially reverse the effects of darkening in a fiber. It is also possible to treat silica core fibers by briefly exposing them to controlled levels of radiation at controlled temperatures; this increases the fiber loss, but makes the fiber less susceptible to future irradiation. These so-called radiation hardened fibers are often used in environments where radiation is anticipated to play an important role. Recently, several models have been advanced³¹ for the performance of fiber under moderate radiation levels; the effect on BER is a power law model of the form

$$\text{BER} = \text{BER}_0 + A(\text{dose})^b \quad (33)$$

where BER_0 is the link BER prior to irradiation, the dose is given in rads, and the constants A and b are empirically fitted. The loss due to normal background radiation exposure over a typical link lifetime can be held approximately below 0.5 dB.

15.4 REFERENCES

1. S. E. Miller and A. G. Chynoweth, editors, *Optical Fiber Telecommunications*, Academic Press, Inc., New York, N.Y. (1979).
2. J. Gowar, *Optical Communication Systems*, Prentice Hall, Englewood Cliffs, N.J. (1984).
3. C. DeCusatis, editor, *Handbook of Fiber Optic Data Communication*, Elsevier/Academic Press, New York, N.Y. (first edition 1998, second edition 2002); see also *Optical Engineering* special issue on optical data communication (December 1998).
4. R. Lasky, U. Osterberg, and D. Stigliani, editors, *Optoelectronics for Data Communication*, Academic Press, New York, N.Y. (1995).
5. Digital Video Broadcasting (DVB) Measurement Guidelines for DVB Systems, "European Telecommunications Standards Institute ETSI Technical Report ETR 290, May 1997;" Digital Multi-Programme Systems for Television Sound and Data Services for Cable Distribution, "International Telecommunications Union ITU-T Recommendation J.83, 1995;" Digital Broadcasting System for Television, Sound and Data Services; Framing Structure, Channel Coding and Modulation for Cable Systems, "European Telecommunications Standards Institute ETSI 300 429," 1994.
6. W. E. Stephens and T. R. Hoseph, "System Characteristics of Direct Modulated and Externally Modulated RF Fiber-Optic Links," *IEEE J. Lightwave Technol.* **LT-5(3)**:380–387 (1987).
7. C. H. Cox, III and, E. I. Ackerman, "Some Limits on the Performance of an Analog Optical Link," *Proc. SPIE—Int. Soc. Opt. Eng.* **3463**:2–7 (1999).
8. United States laser safety standards are regulated by the Dept. of Health and Human Services (DHHS), Occupational Safety and Health Administration (OSHA), Food and Drug Administration (FDA), Code of Radiological Health (CDRH), 21 Code of Federal Regulations (CFR) subchapter J; the relevant standards are ANSI Z136.1, "Standard for the Safe Use of Lasers" (1993 revision) and ANSI Z136.2, "Standard for the Safe Use of Optical Fiber Communication Systems Utilizing Laser Diodes and LED Sources" (1996–97 revision); elsewhere in the world, the relevant standard is International Electrotechnical Commission (IEC/CEI) 825 (1993 revision).
9. S. S. Walker, "Rapid Modeling and Estimation of Total Spectral Loss in Optical Fibers," *IEEE J. Lightwave Technol.* **4**:1125–1132 (1996).
10. Electronics Industry Association/Telecommunications Industry Association (EIA/TIA) Commercial Building Telecommunications Cabling Standard (EIA/TIA-568-A), Electronics Industry Association/Telecommunications Industry Association (EIA/TIA) Detail Specification for 62.5 micron Core Diameter/125 micron Cladding Diameter Class 1a Multimode Graded Index Optical Waveguide Fibers (EIA/TIA-492AAAA), Electronics Industry Association/Telecommunications Industry Association (EIA/TIA) Detail Specification for Class IV-a Dispersion Unshifted Single-Mode Optical Waveguide Fibers Used in Communications Systems (EIA/TIA-492BAAA), Electronics Industry Association, New York, N.Y.
11. D. Gloge, "Propagation Effects in Optical Fibers," *IEEE Trans. Microwave Theory Technol.* **MTT-23**:106–120 (1975).
12. P. M. Shanker, "Effect of Modal Noise on Single-Mode Fiber Optic Network," *Opt. Comm.* **64**:347–350 (1988).
13. J. J. Refi, "LED Bandwidth of Multimode Fiber as a Function of Source Bandwidth and LED Spectral Characteristics," *IEEE J. Lightwave Technol.* **LT-14**:265–272 (1986).
14. G. P. Agrawal et al., "Dispersion Penalty for 1.3 Micron Lightwave Systems with Multimode Semiconductor Lasers," *IEEE J. Lightwave Technol.* **6**:620–625 (1988).
15. K. Ogawa, "Analysis of Mode Partition Noise in Laser Transmission Systems," *IEEE J. Quantum Elec.* **QE-18**:849–9855 (1982).
16. K. Ogawa, *Semiconductor Laser Noise; Mode Partition Noise*, in *Semiconductors and Semimetals* (R. K. Willardson and A. C. Beer, editors), vol. 22C, Academic Press, New York, N. Y. (1985).
17. J. C. Campbell, "Calculation of the Dispersion Penalty of the Route Design of Single-Mode Systems," *IEEE J. Lightwave Technol.* **6**:564–573 (1988).
18. M. Ohtsu et al., "Mode Stability Analysis of Nearly Single-Mode Semiconductor Laser," *IEEE J. Quantum Elec.* **24**:716–723 (1988).
19. M. Ohtsu and Y. Teramachi, "Analysis of Mode Partition and Mode Hopping in Semiconductor Lasers," *IEEE Quantum Elec.* **25**:31–38 (1989).

20. D. Duff et al., "Measurements and Simulations of Multipath Interference for 1.7 Gbit/s Lightwave Systems Utilizing Single and Multifrequency Lasers," *Proc. OFC* p. 128 (1989).
21. J. Radcliffe, "Fiber Optic Link Performance in the Presence of Internal Noise Sources," *IBM Technical Report*, Glendale Labs, Endicott, New York, N.Y. (1989).
22. L. L. Xiao, C. B. Su, and R. B. Lauer, "Increase in Laser RIN due to Asymmetric Nonlinear Gain, Fiber Dispersion, and Modulation," *IEEE Photon. Tech. Lett.* **4**:774–777 (1992).
23. P. Trischitta and P. Sannuti, "The Accumulation of Pattern Dependent Jitter for a Chain of Fiber Optic Regenerators," *IEEE Trans. Comm.* **36**:761–765 (1988).
24. *CCITT Recommendations G.824, G.823, O.171, and G.703 on timing jitter in digital systems* (1984).
25. R. J. S. Bates, "A Model for Jitter Accumulation in Digital Networks," *IEEE Globecom Proc.* pp. 145–149 (1983).
26. C. J. Byrne, B. J. Karafin, and D. B. Robinson, Jr., "Systematic Jitter in a Chain of Digital Regenerators," *Bell Sys. Tech. J.* **43**:2679–2714 (1963).
27. R. J. S. Bates and L. A. Sauer, "Jitter Accumulation in Token Passing Ring LANs," *IBM J. Res. Dev.* **29**:580–587 (1985).
28. C. Chamzas, "Accumulation of Jitter: A Stochastic Model," *AT&T Tech. J.* p. 64 (1985).
29. D. Marcuse and H. M. Presby, "Mode Coupling in an Optical Fiber with Core Distortion," *Bell Sys. Tech. J.* **1**:3 (1975).
30. E. J. Frieble et al., "Effect of Low Dose Rate Irradiation on Doped Silica Core Optical Fibers," *App. Opt.* **23**:4202–4208 (1984).
31. J. B. Haber et al., "Assessment of Radiation Induced Loss for AT&T Fiber Optic Transmission Systems in the Terrestrial Environment," *IEEE J. Lightwave Technol.* **6**:150–154 (1988).

This page intentionally left blank.

DO NOT DUPLICATE

Daniel Nolan

*Corning Inc.
Corning, New York*

16.1 INTRODUCTION

Fiber-optic couplers, including splitters and wavelength division multiplexing components, have been used extensively over the last two decades. This use continues to grow both in quantity and in the ways in which the devices are used. The uses today include, among other applications, simple splitting for signal distribution and the wavelength multiplexing and demultiplexing multiple wavelength signals.

Fiber-based splitters and wavelength division multiplexing (WDM) components are among the simplest devices. Other technologies that can be used to fabricate components that exhibit similar functions include the planar waveguide and micro-optic technologies. Planar waveguides are most suitable for highly integrated functions. Micro-optic devices are often used when complex multiple wavelength functionality is required. In this chapter, we will show the large number of optical functions that can be achieved with simple tapered fiber components. We will also describe the physics of the propagation of light through tapers in order to better understand the breadth of components that can be fabricated with this technology. The phenomenon of coupling includes an exchange of power that can depend both on wavelength and on polarization. Beyond the simple 1×2 power splitter, other devices that can be fabricated from tapered fibers include $1 \times N$ devices, wavelength multiplexing, polarization multiplexing, switches, attenuators, and filters.

Fiber-optic couplers have been fabricated since the early 1970s. The fabrication technologies have included fusion tapering,¹⁻³ etching,⁴ and polishing.⁵⁻⁷ The tapered single-mode fiber-optic power splitter is perhaps the most universal of the single-mode tapered devices.⁸ It has been shown that the power transferred during the tapering process involves an initial adiabatic transfer of the power in the input core to the cladding air interface.⁹ The light is then transferred to the adjacent core-cladding mode. During the uptapering process, the input light will transfer back onto the fiber cores. In this case, it is referred to as a cladding mode coupling device. Light that is transferred to a higher-order mode of the core-cladding structure leads to an excess loss. This is because these higher-order modes are not bounded by the core and are readily stripped by the higher index of the fiber coating.

In the tapered fiber coupler process, two fibers are brought into close proximity after the protective plastic jacket is removed. Then, in the presence of a torch, the fibers are fused and stretched (Fig. 1.) The propagation of light through this tapered region is described using Maxwell's vector

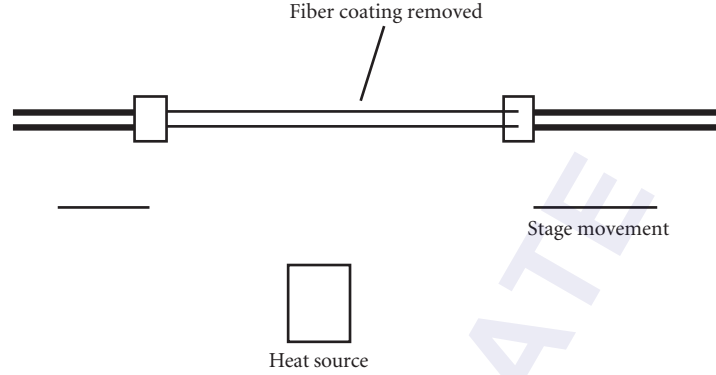


FIGURE 1 Fusing and tapering process.

equations, but to a good approximation, the scalar wave approximation is valid. The scalar wave equation written in cylindrical coordinates is expressed as

$$[1/r\partial/\partial r r\partial/\partial r - v/r^2 + k^2 n_1^2 - (V/a)^2 f(r/a)]\psi = \epsilon\mu\partial^2\psi/\partial t^2 \quad (1)$$

In Eq. (1), n_1 is the index value at $r = 0$, β is the propagation constant, which is to be determined, a is the core radius, $f(r/a)$ is a function describing the index distribution with radius, and V is the modal volume

$$V = \frac{2\pi n_1}{\lambda\sqrt{2\delta}} \quad (2)$$

with

$$\delta = \frac{[n_1^2 - n_2^2]}{2n_1^2} \quad (3)$$

As light propagates in the single-mode fiber, it is not confined to the core region, but extends out into the surrounding region. As the light propagates through the tapered region, it is bounded by the shrinking, air-cladding boundary.

In the simplest case, the coupling from one cladding to the adjacent one can be described by perturbation theory.⁹ In this case, the cladding air boundary is considered as the waveguide outer boundary, and the exchange of power along z is described as

$$P = \sin^2[Cz] \quad (4)$$

where¹⁰

$$C = \frac{[\pi\delta/Wd\rho]^2 U^2 \exp(-Wd/\rho)}{[V^3 K_1^2(W)]} \quad (5)$$

with

$$\alpha = 2\pi n_1/\lambda \quad U = \rho(k^2 n_1^2 - \beta^2) \quad W = \rho(\beta^2 - k^2 n_2^2) \quad (6)$$

In Eq. (6), the waveguide parameters are defined in the tapered region. Here the core of each fiber is small and the cladding becomes the effective core, while air becomes the cladding. Also, it

is important to point out that Eqs. (4) and (5) are only a first approximation. These equations are derived using first-order perturbation theory. Also, the scalar wave equation is not strictly valid under the presence of large index differences, such as at a glass-air boundary. However, these equations describe a number of important effects. The sinusoidal dependence of the power coupled with wavelength, as well as the dependence of power transfer with cladding diameter and other dependencies, is well described with the model.

Equation (4) can be described by considering the light input to one core as a superposition of symmetric and antisymmetric.⁹ These modes are eigen solutions to the composite two-core structure. The proper superposition of these two modes enables one to impose input boundary conditions for the case of a two-core structure. The symmetric and antisymmetric modes are written as

$$\psi_s = \frac{\psi_1 + \psi_2}{\sqrt{2}} \quad (7)$$

$$\psi_a = \frac{\psi_1 - \psi_2}{\sqrt{2}} \quad (8)$$

Light input onto one core is described with ψ_1 at $z = 0$,

$$\psi_1 = \frac{\psi_s + \psi_a}{\sqrt{2}} \quad (9)$$

Propagation through the coupler is characterized with the superposition of ψ_s and ψ_a . This superposition describes the power transfer between the two guides along the direction of propagation.¹⁰ The propagation constants of ψ_s and ψ_a are slightly different, and this value can be used to estimate excess loss under certain perturbations.

16.2 ACHROMATICITY

The simple sinusoidal dependence of the coupling with wavelength as described above is not always desired, and often a more achromatic dependence of the coupling is required. This can be achieved when dissimilar fibers¹⁰ are used to fabricate the coupler. Fibers are characterized as dissimilar when the propagation constants of the guides are of different values. When dissimilar fibers are used, Eqs. (4) and (5) can be replaced with

$$P_1(x) = P_1(0) + F^2(P_2(0) - P_1(0) + [\delta\beta/C][P_1(0)P_2(0)]^2)\sin^2(Cz/F) \quad (10)$$

where

$$F = 1/[1 + \delta\beta/(4C^2)] \quad (11)$$

In most cases, the fibers are made dissimilar by changing the cladding diameter of one of the fibers. Etching or pretapering one of the fibers can do this. Another approach is to slightly change the cladding index of one of the fibers.¹¹ When dissimilar fibers are used, the total amount of power coupled is limited. As an example, an achromatic 3-dB coupler is made achromatic by operating at the sinusoidal maximum with wavelength rather than at the power of maximum power change with wavelength. Another approach to achieve achromaticity is to taper the device such that the modes expand well beyond the cladding boundaries.¹² This condition greatly weakens the wavelength dependence of the coupling. This has been achieved by encapsulating the fibers in a third-matrix glass with an index very close to that of the fiber's cladding index. The difference in index between the cladding and the matrix glass is in the order of 0.001. The approach of encapsulating the fibers in a third-index material^{13,14} is also useful for reasons other than achromaticity. One reason is that

the packaging process is simplified. Also, a majority of couplers made for undersea applications use this method because it is a proven approach to ultrahigh reliability.

The wavelength dependence of the couplers described above is most often explained using mode coupling and perturbation theory. Often, numerical analysis is required to explain the effects that the varying taper angles have on the overall coupling. An important numerical approach is the beam propagation method.¹⁵ In this approach, the propagation of light through a device is solved by an expansion of the evolution operator using a Taylor series and with the use of fast Fourier transforms to evaluate the appropriate derivatives. In this way, the propagation of the light can be studied as it couples to the adjacent guides or to higher-order modes.

16.3 WAVELENGTH DIVISION MULTIPLEXING

Besides power splitting, tapered couplers can be used to separate wavelengths. To accomplish this separation, we utilize the wavelength dependence of Eqs. (4) and (5). By proper choice of the device length and taper ratio, two predetermined wavelengths can be put out onto two different ports. Wavelengths from 60 to 600 nanometers can be split using this approach. Applications include the splitting and/or combining of 1480 nm and 1550 nm light, as well as multiplexing 980 and 1550 nm onto an erbium fiber for signal amplification. Also important is the splitting of the 1310- to 1550-nm wavelength bands, which can be achieved using this approach.

16.4 $1 \times N$ POWER SPLITTERS

Often it is desirable to split a signal onto a number of output ports. This can be achieved by concatenating 1×2 power splitters. Alternatively, one can split the input simultaneously onto multiple output ports.^{16,17} Typically, the output ports are of the form 2^N (i.e., 2, 4, 8, 16, . . .). The configuration of the fibers in the tapered region affects the distribution of the output power per port. A good approach to achieve uniform 1×8 splitting is described in Ref. 18.

16.5 SWITCHES AND ATTENUATORS

In a tapered device, the power coupled over to the adjacent core can be significantly affected by bending the device at the midpoint. By encapsulating two fibers before tapering in a third index medium, the device is rigid and can be reliably bent in order to frustrate the coupling.¹⁹ The bending establishes a difference in the propagation constants of the two guiding media, preventing coupling or power transfer.

This approach can be used to fabricate both switches and attenuators. Switches with up to 30 dB crosstalk and attenuators with variable crosstalk up to 30 dB as well over the erbium wavelength band have been fabricated. Displacing one end of a 1 cm taper by 1 mm is enough to alter the crosstalk by the 30-dB value. Applications for attenuators have been increasing significantly over the last few years. An important reason is to maintain the gain in erbium-doped fiber amplifiers. This is achieved by limiting the amount of pump power into the erbium fiber. Over time, as the pump degrades, the power output of the attenuator is increased to compensate for the pump degradation.

16.6 MACH-ZEHNDER DEVICES

Devices to split narrowly spaced wavelengths are very important. As mentioned above, tapers can be designed such that wavelengths from 60 to 600 nm can be split in a tapered device. Dense WDM networks require splitting of wavelengths with separations on the order of nanometers. Fiber-based

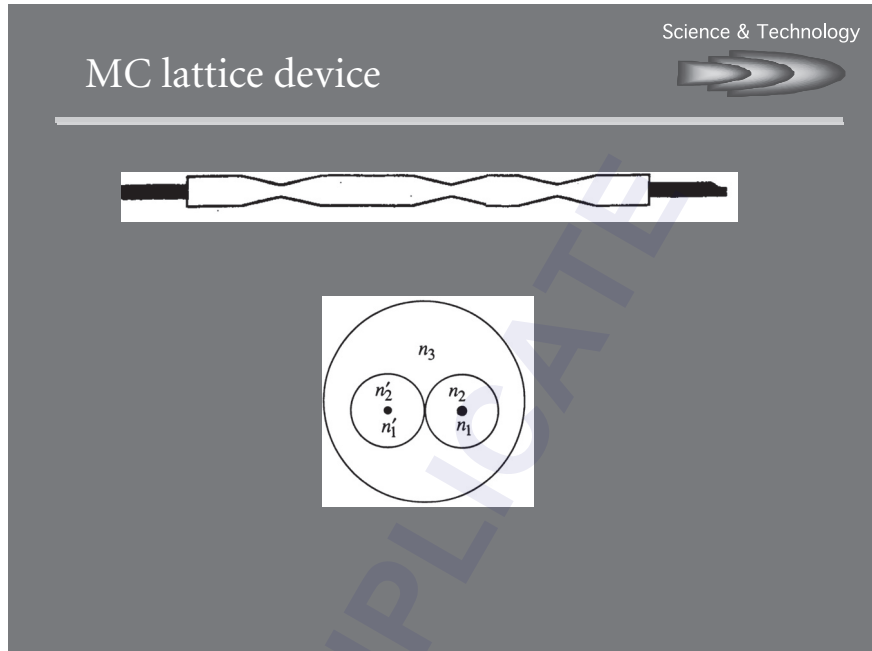


FIGURE 2 Fiber-based Mach-Zehnder devices.²⁴

Mach-Zehnder devices enable such splitting. Monolithic fiber-based Mach-Zehnders can be fabricated using fibers with different cores,^{20,21} i.e., different propagation constants. Two or more tapers can be used to cause light from two different optical paths to interfere (Fig. 2). The dissimilar cores enable light to propagate at different speeds between the tapers, causing the required constructive and destructive interference. These devices are environmentally stable due to the monolithic structure. Mach-Zehnders can also be fabricated using fibers with different lengths between the tapers.²² In this approach, it is the packaging that enables an environmentally stable device.

Mach-Zehnders and lattice filters can also be fabricated by tapering single-fiber devices.^{23,24} In the tapered regions, the light couples to a cladding mode. The cladding mode propagates between tapers since a lower-index overcladding replaces the higher-index coating material. An interesting application for these devices is gain-flattening filters for amplifiers.

16.7 POLARIZATION DEVICES

It is well known that two polarization modes propagate in single-mode fiber. Most optical fiber modules allow both polarizations to propagate, but specify that the performance of the components be insensitive to the polarization states of the propagating light. However, this is often not the situation for fiber-optic sensor applications. Often, the state of polarization is important to the operation of the sensor itself. In these situations, polarization-maintaining fiber is used. Polarization components such as polarization-maintaining couplers and also single polarization devices are used. In polarization-maintaining fiber, a difference in propagation constants of the polarization modes prevents mode coupling or exchange of energy. This is achieved by introducing stress or shape birefringence within the fiber core. A significant difference between the two polarization modes is maintained as the fiber twists in a cable or package.

In many fiber sensor systems, tapered fiber couplers are used to couple light from one core to another. Often the couplers are composed of birefringent fibers.^{24,25} This is done to maintain the

alignment of the polarizations to the incoming and outgoing fibers and also to maintain the polarization states within the device. The axes of the birefringent fibers are aligned before tapering, and care is taken not to excessively twist the fibers during the tapering process.

The birefringent fibers contain stress rods, elliptical core fibers, or inner claddings to maintain the birefringence. The stress rods in some birefringent fibers have an index higher than the silica cladding. In the tapering process, this can cause light to be trapped in these rods, resulting in an excess loss in the device. Stress rods with an index lower than that of silica can be used in these fibers, resulting in very low loss devices.

16.8 SUMMARY

Tapered fiber couplers are extremely useful devices. Such devices include 1×2 and $1 \times N$ power splitters, wavelength division multiplexers and filters, and polarization-maintaining and -splitting components. Removing the fiber's plastic coating and then fusing and tapering two or more fibers in the presence of heat forms these devices. The simplicity and flexibility of this fabrication process is in part responsible for the widespread use of these components. The mechanism involved in the fabrication process is reasonably understood and simple, which is in part responsible for the widespread deployment of these devices. These couplers are found in optical modules for the telecommunication industry and in assemblies for the sensing industry. They are also being deployed as standalone components for fiber-to-home applications.

16.9 REFERENCES

1. T. Ozeki and B. S. Kawaski, "New Star Coupler Compatible with Single Multimode Fiber Links," *Electron. Lett.* **12**:151–152, 1976.
2. B. S. Kawaski and K. O. Hill, "Low Loss Access Coupler for Multimode Optical Fiber Distribution Networks," *Appl. Opt.* **16**:1794–1795, 1977.
3. G. E. Rawson and M. D. Bailey, "Bitaper Star Couplers with up to 100 Fiber Channels," *Electron. Lett.* **15**:432–433, 1975.
4. S. K. Sheem and T. G. Giallorenzi, "Single-Mode Fiber Optical Power Divided; Encapsulated Etching Technique," *Opt. Lett.* **4**:31, 1979.
5. Y. Tsujimoto, H. Serizawa, K. Hatori, and M. Fukai, "Fabrication of Low Loss 3 dB Couplers with Multimode Optical Fibers," *Electron. Lett.* **14**:157–158, 1978.
6. R. A. Bergh, G. Kotler, and H. J. Shaw, "Single-Mode Fiber Optic Directional Coupler," *Electron. Lett.* **16**:260–261, 1980.
7. O. Parriaux, S. Gidon, and A. Kuznetsov, "Distributed Coupler on Polished Single-Mode Fiber," *Appl. Opt.* **20**:2420–2423, 1981.
8. B. S. Kawaski, K. O. Hill, and R. G. Lamont, "Biconical—Taper Single-Mode Fiber Coupler," *Opt. Lett.* **6**:327, 1981.
9. R. G. Lamont, D. C. Johnson, and K. O. Hill, "Power Transfer in Fused Biconical Single Mode Fiber Couplers: Dependence on External Refractive Index," *Appl. Opt.* **24**:327–332, 1984.
10. A. Snyder and J. D. Love, *Optical Waveguide Theory*, London: Chapman and Hall, 1983.
11. W. J. Miller, C. M. Truesdale, D. L. Weidman, and D. R. Young, "Achromatic Fiber Optic Coupler," U.S. Patent 5,011,251, Apr. 1991.
12. D. L. Weidman, "Achromat Overclad Coupler," U.S. Patent, 5,268,979, Dec. 1993.
13. C. M. Truesdale and D. A. Nolan, "Core-Clad Mode Coupling in a New Three-Index Structure," *European Conference on Optical Communications*, Barcelona Spain, 1986.
14. D. B. Keck, A. J. Morrow, D. A. Nolan, and D. A. Thompson, "Passive Optical Components in the Subscriber Loop," *J. Lightwave Technol.* **7**:1623–1633, 1989.

15. M. D. Feit and J. A. Fleck, "Simple Spectral Method for Solving Propagation Problems in Cylindrical Geometry with Fast Fourier Transforms," *Opt. Lett.* **14**:662–664, 1989.
16. D. B. Mortimore and J. W. Arkwright, "Performance of Wavelength-Flattened 1×7 Fused Couplers," Optical Fiber Conference, TUG6, 1990.
17. D. L. Weidman, "A New Approach to Achromaticity in Fused $1 \times N$ Couplers," Optical Fiber Conference, Post Deadline papers, 1994.
18. W. J. Miller, D. A. Nolan, and G. E. Williams, "Method of Making a $1 \times N$ Coupler," US Patent, 5,017,206, 1991.
19. M. A. Newhouse and F. A. Annunziata, "Single-Mode Optical Switch," Technical Digest of the National Fiber Optic Conference, 1990.
20. D. A. Nolan and W. J. Miller, "Wavelength Tunable Mach-Zehnder Device," Optical Conference, 1994.
21. B. Malo, F. Bilodeau, K. O. Hill, and J. Albert, "Unbalanced Dissimilar—Fiber Mach—Zehnder Interferometer: Application as Filter," *Electron. Lett.* **25**:1416, 1989.
22. C. Huang, H. Luo, S. Xu, and P. Chen, "Ultra Low Loss, Temperature Insensitive 16 Channel 100 GHz Dense WDMs Based on Cascaded All Fiber Unbalanced Mach-Zehnder Structure," Optical Fiber Conference, TUH2, 1999.
23. D. A. Nolan, W. J. Miller, and R. Irion, "Fiber Based Band Splitter," Optical Fiber Conference, 1998.
24. D. A. Nolan, W. J. Miller, G. Berkey, and L. Bhagavatula, "Tapered Lattice Filters," Optical Fiber Conference, TUH4, 1999.
25. I. Yokohama, M. Kawachi, K. Okamoto, and J. Noda, *Electron. Lett.* **22**:929, 1986.

This page intentionally left blank.

DO NOT DUPLICATE

FIBER BRAGG GRATINGS

Kenneth O. Hill

*Communications Research Centre
Ottawa, Ontario, Canada, and
Nu-Wave Photonics
Ottawa, Ontario, Canada*

17.1 GLOSSARY

FBG	fiber Bragg grating
FWHM	full width measured at half-maximum intensity
N_{eff}	effective refractive index for light propagating in a single mode
pps	pulses per second
β	propagation constant of optical fiber mode
Δn	magnitude of photoinduced refractive index change
κ	grating coupling coefficient
Λ	spatial period (or pitch) of spatial feature measured along optical fiber
λ	vacuum wavelength of propagating light
λ_B	Bragg wavelength
L	length of grating

17.2 INTRODUCTION

A fiber Bragg grating (FBG) is a periodic variation of the refractive index of the fiber core along the length of the fiber. The principal property of FBGs is that they reflect light in a narrow bandwidth that is centered about the Bragg wavelength λ_B which is given by $\lambda_B = 2N_{\text{eff}}\Lambda$, where Λ is the spatial period (or pitch) of the periodic variation and N_{eff} is the effective refractive index for light propagating in a single mode, usually the fundamental mode of a monomode optical fiber. The refractive index variations are formed by exposure of the fiber core to an intense optical interference pattern of ultraviolet light. The capability of light to induce permanent refractive index changes in the core of an optical fiber has been named photosensitivity. Photosensitivity was discovered by Hill et al. in 1978 at the Communications Research Centre in Canada (CRC).^{1,2} The discovery has led to techniques for fabricating Bragg gratings in the core of an optical fiber and a means for manufacturing a

wide range of FBG-based devices that have applications in optical fiber communications and optical sensor systems.

This chapter reviews the characteristics of photosensitivity, the properties of Bragg gratings, the techniques for fabricating Bragg gratings in optical fibers, and some FBG devices. More information on FBGs can be found in the following references, which are reviews on Bragg grating technology,^{3,4} the physical mechanisms underlying photosensitivity,⁵ applications for fiber gratings,⁶ and the use of FBGs as sensors.⁷

17.3 PHOTSENSITIVITY

When ultraviolet light radiates an optical fiber, the refractive index of the fiber is changed permanently; the effect is termed *photosensitivity*. The change in refractive index is permanent in the sense that it will last for several years (lifetimes of 25 years are predicted) if the optical waveguide after exposure is annealed appropriately; that is, by heating for a few hours at a temperature of 50°C above its maximum anticipated operating temperature.⁸ Initially, photosensitivity was thought to be a phenomenon that was associated only with germanium-doped-core optical fibers. Subsequently, photosensitivity has been observed in a wide variety of different fibers, many of which do not contain germanium as dopant. Nevertheless, optical fiber with a germanium-doped core remains the most important material for the fabrication of Bragg grating-based devices.

The magnitude of the photoinduced refractive index change (Δn) obtained depends on several different factors: the irradiation conditions (wavelength, intensity, and total dosage of irradiating light), the composition of glassy material forming the fiber core, and any processing of the fiber prior and subsequent to irradiation. A wide variety of different continuous-wave and pulsed-laser light sources, with wavelengths ranging from the visible to the vacuum ultraviolet, have been used to photoinduce refractive index changes in optical fibers. In practice, the most commonly used light sources are KrF and ArF excimer lasers that generate, respectively, 248- and 193-nm light pulses (pulse width ~10 ns) at pulse repetition rates of 50 to 100 pps. Typically, the fiber core is exposed to laser light for a few minutes at pulse levels ranging from 100 to 1000 mJ cm⁻² pulse⁻¹. Under these conditions, Δn is positive in germanium-doped monomode fiber with a magnitude ranging between 10⁻⁵ and 10⁻³.

The refractive index change can be enhanced (photosensitization) by processing the fiber prior to irradiation using such techniques as *hydrogen loading*⁹ or *flame brushing*.¹⁰ In the case of hydrogen loading, a piece of fiber is put in a high-pressure vessel containing hydrogen gas at room temperature; pressures of 100 to 1000 atmospheres (atm; 101 kPa/atm) are applied. After a few days, hydrogen in molecular form has diffused into the silica fiber; at equilibrium the fiber becomes saturated (i.e., loaded) with hydrogen gas. The fiber is then taken out of the high-pressure vessel and irradiated before the hydrogen has had sufficient time to diffuse out. Photoinduced refractive index changes up to 100 times greater are obtained by hydrogen loading a Ge-doped-core optical fiber. In flame brushing, the section of fiber that is to be irradiated is mounted on a jig and a hydrogen-fueled flame is passed back and forth (i.e., brushed) along the length of the fiber. The brushing takes about 10 minutes, and upon irradiation, an increase in the photoinduced refractive index change by about a factor of 10 can be obtained.

Irradiation at intensity levels higher than 1000 mJ/cm² marks the onset of a different non-linear photosensitive process that enables a single irradiating excimer light pulse to photo-induce a large index change in a small localized region near the core/cladding boundary of the fiber. In this case, the refractive index changes are sufficiently large to be observable with a phase contrast microscope and have the appearance of physically damaging the fiber. This phenomenon has been used for the writing of gratings using a single-excimer light pulse.

Another property of the photoinduced refractive index change is *anisotropy*. This characteristic is most easily observed by irradiating the fiber from the side with ultraviolet light that is polarized perpendicular to the fiber axis. The anisotropy in the photoinduced refractive index change results in the fiber becoming birefringent for light propagating through the fiber. The effect is useful for fabricating polarization mode-converting devices or rocking filters.¹¹

The physical processes underlying photosensitivity have not been fully resolved. In the case of germanium-doped glasses, photosensitivity is associated with GeO color center defects that have strong absorption in the ultraviolet (~242 nm) wavelength region. Irradiation with ultraviolet light bleaches the color center absorption band and increases absorption at shorter wavelengths, thereby changing the ultraviolet absorption spectrum of the glass. Consequently, as a result of the Kramers-Kronig causality relationship,¹² the refractive index of the glass also changes; the resultant refractive index change can be sensed at wavelengths that are far removed from the ultraviolet region extending to wavelengths in the visible and infrared. The physical processes underlying photosensitivity are, however, probably much more complex than this simple model. There is evidence that ultraviolet light irradiation of Ge-doped optical fiber results in structural rearrangement of the glass matrix leading to densification, thereby providing another mechanism for contributing to the increase in the fiber core refractive index. Furthermore, a physical model for photosensitivity must also account for the small anisotropy in the photoinduced refractive index change and the role that hydrogen loading plays in enhancing the magnitude of the photoinduced refractive change. Although the physical processes underlying photosensitivity are not completely known, the phenomenon of glass-fiber photosensitivity has the practical result of providing a means, using ultraviolet light, for photoinducing permanent changes in the refractive index at wavelengths that are far removed from the wavelength of the irradiating ultraviolet light.

17.4 PROPERTIES OF BRAGG GRATINGS

Bragg gratings have a periodic index structure in the core of the optical fiber. Light propagating in the Bragg grating is backscattered slightly by Fresnel reflection from each successive index perturbation. Normally, the amount of backscattered light is very small except when the light has a wavelength in the region of the Bragg wavelength λ_B , given by

$$\lambda_B = 2N_{\text{eff}}\Lambda$$

where N_{eff} is the modal index and Λ is the grating period. At the Bragg wavelength, each back reflection from successive index perturbations is in phase with the next one. The back reflections add up coherently and a large reflected light signal is obtained. The reflectivity of a strong grating can approach 100 percent at the Bragg wavelength, whereas light at wavelengths longer or shorter than the Bragg wavelength pass through the Bragg grating with negligible loss. It is this wavelength-dependent behavior of Bragg gratings that makes them so useful in optical communications applications. Furthermore, the optical pitch ($N_{\text{eff}}\Lambda$) of a Bragg grating contained in a strand of fiber is changed by applying longitudinal stress to the fiber strand. This effect provides a simple means for sensing strain optically by monitoring the concomitant change in the Bragg resonant wavelength.

Bragg gratings can be described theoretically by using coupled-mode equations.^{4,6,13} Here, we summarize the relevant formulas for tightly bound monomode light propagating through a uniform grating. The grating is assumed to have a sinusoidal perturbation of constant amplitude Δn . The reflectivity of the grating is determined by three parameters: (1) the coupling coefficient κ (2) the mode propagation constant $\beta = 2\pi N_{\text{eff}}/\lambda$, and (3) the grating length L . The coupling coefficient κ which depends only on the operating wavelength of the light and the amplitude of the index perturbation Δn is given by $\kappa = (\pi/\lambda)\Delta n$. The most interesting case is when the wavelength of the light corresponds to the Bragg wavelength. The grating reflectivity R of the grating is then given by the simple expression, $R = \tan^2(\kappa L)$, where κ is the coupling coefficient at the Bragg wavelength and L is the length of the grating. Thus, the product κL can be used as a measure of grating strength. For $\kappa L = 1, 2, 3$, the grating reflectivity is, respectively, 58, 93, and 99 percent. A grating with a κL greater than one is termed a strong grating, whereas a weak grating has κL less than one. Figure 1 shows the typical reflection spectra for weak and strong gratings.

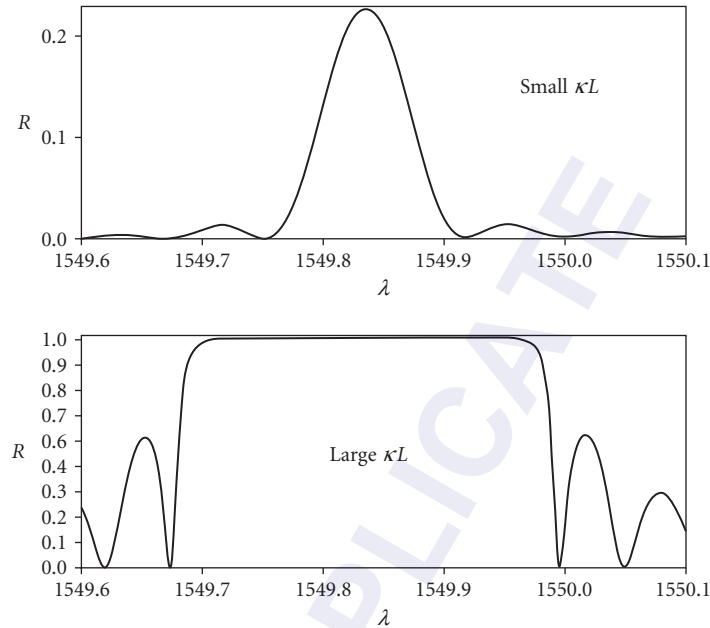


FIGURE 1 Typical reflection spectra for weak (small κL) and strong (large κL) fiber gratings.

The other important property of the grating is its bandwidth, which is a measure of the wavelength range over which the grating reflects light. The bandwidth of a fiber grating that is most easily measured is its full width at half-maximum, $\Delta\lambda_{\text{FWHM}}$, of the central reflection peak, which is defined as the wavelength interval between the 3-dB points. That is the separation in the wavelength between the points on either side of the Bragg wavelength where the reflectivity has decreased to 50 percent of its maximum value. However, a much easier quantity to calculate is the bandwidth, $\Delta\lambda_0 = \lambda_0 - \lambda_B$, where λ_0 is the wavelength where the first zero in the reflection spectra occurs. This bandwidth can be found by calculating the difference in the propagation constants, $\Delta\beta_0 = \beta_0 - \beta_B$, where $\beta_0 = 2\pi N_{\text{eff}}/\lambda_0$ is the propagation constant at wavelength λ_0 for which the reflectivity is first zero, and $\beta_B = 2\pi N_{\text{eff}}/\lambda_B$ is the propagation constant at the Bragg wavelength for which the reflectivity is maximum.

In the case of weak gratings ($\kappa L < 1$), $\Delta\beta_0 = \beta_0 - \beta_B = \pi/L$, from which it can be determined that $\Delta\lambda_{\text{FWHM}} \sim \Delta\lambda_0 = \lambda_B^2/2N_{\text{eff}}L$; the bandwidth of a weak grating is inversely proportional to the grating length L . Thus, long, weak gratings can have very narrow bandwidths. The first Bragg grating written in fibers^{1,2} was more than 1 m long and had a bandwidth less than 100 MHz, which is an astonishingly narrow bandwidth for a reflector of visible light. On the other hand, in the case of a strong grating ($\kappa L > 1$), $\Delta\beta_0 = \beta_0 - \beta_B = 4\kappa$ and $\Delta\lambda_{\text{FWHM}} \sim 2\Delta\lambda_0 = 4\lambda_B^2\kappa/\pi N_{\text{eff}}$. For strong gratings, the bandwidth is directly proportional to the coupling coefficient κ and is independent of the grating length.

17.5 FABRICATION OF FIBER GRATINGS

Writing a fiber grating optically in the core of an optical fiber requires irradiating the core with a periodic interference pattern. Historically, this was first achieved by interfering light that propagated in a forward direction along an optical fiber with light that was reflected from the fiber end and propagated in a backward direction.¹ This method for forming fiber gratings is known as the *internal*

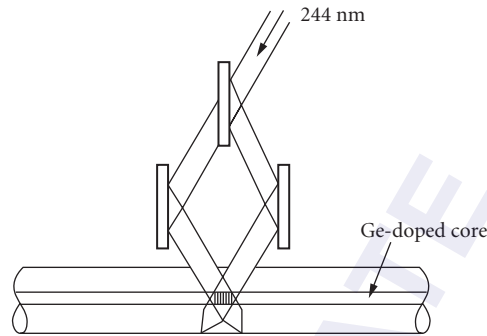


FIGURE 2 Schematic diagram illustrating the writing of an FBG using the transverse holographic technique.

writing technique, and the gratings were referred to as *Hill gratings*. The Bragg gratings, formed by internal writing, suffer from the limitation that the wavelength of the reflected light is close to the wavelength at which they were written (i.e., at a wavelength in the blue-green spectral region).

A second method for fabricating fiber gratings is the *transverse holographic technique*,¹⁴ which is shown schematically in Fig. 2. The light from an ultraviolet source is split into two beams that are brought together so that they intersect at an angle θ . As Fig. 2 shows, the intersecting light beams form an interference pattern that is focused using cylindrical lenses (not shown) on the core of the optical fiber. Unlike the internal writing technique, the fiber core is irradiated from the side, thus giving rise to its name *transverse holographic technique*. The technique works because the fiber cladding is transparent to the ultraviolet light, whereas the core absorbs the light strongly. Since the period Λ of the grating depends on the angle θ between the two interfering coherent beams through the relationship $\Lambda = \lambda_{UV}/2 \sin(\theta/2)$, Bragg gratings can be made that reflect light at much longer wavelengths than the ultraviolet light that is used in the fabrication of the grating. Most important, FBGs can be made that function in the spectral regions that are of interest for fiber-optic communication and optical sensing.

A third technique for FBG fabrication is the *phase mask technique*,¹⁵ which is illustrated in Fig. 3. The phase mask is made from a flat slab of silica glass, which is transparent to ultraviolet light. On one of the flat surfaces, a one-dimensional periodic surface relief structure is etched using photolithographic techniques. The shape of the periodic pattern approximates a square wave in profile. The optical fiber is placed almost in contact with and at right angles to the corrugations of the phase mask, as shown in Fig. 3. Ultraviolet light, which is incident normal to the phase mask, passes through and is diffracted by the periodic corrugations of the phase mask. Normally, most of the diffracted light is contained in the 0, +1, and -1 diffracted orders. However, the phase mask is designed to suppress the diffraction into the zero order by controlling the depth of the corrugations in the phase mask. In practice, the amount of light in the zero order can be reduced to less than 5 percent with approximately 80 percent of the total light intensity divided equally in the ± 1 orders. The two ± 1 diffracted-order beams interfere to produce a periodic pattern that photoimprints a corresponding grating in the optical fiber. If the period of the phase mask grating is Λ_{mask} , the period of the photoimprinted index grating is $\Lambda_{\text{mask}}/2$. Note that this period is independent of the wavelength of ultraviolet light that irradiates the phase mask.

The phase mask technique has the advantage of greatly simplifying the manufacturing process for Bragg gratings, while yielding high-performance gratings. In comparison with the holographic technique, the phase mask technique offers easier alignment of the fiber for photoimprinting, reduced stability requirements on the photoimprinting apparatus, and lower coherence requirements on the ultraviolet laser beam, thereby permitting the use of a cheaper ultraviolet excimer laser source. Furthermore, there is the possibility of manufacturing several gratings at once in a single exposure by irradiating parallel fibers through the phase mask. The capability to manufacture high-performance gratings at a

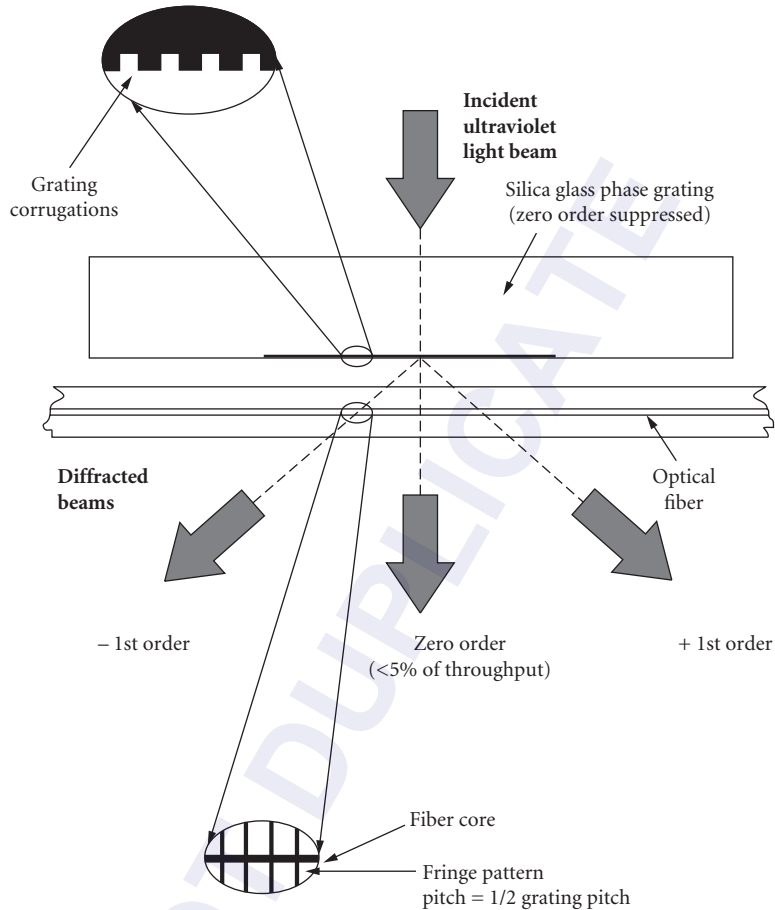


FIGURE 3 Schematic diagram of the phase mask technique for the manufacture of fiber Bragg gratings.

low per-unit grating cost is critical for the economic viability of using gratings in some applications. A draw-back of the phase mask technique is that a separate phase mask is required for each different Bragg wavelength. However, some wavelength tuning is possible by applying tension to the fiber during the photoimprinting process; the Bragg wavelength of the relaxed fiber will shift by ~ 2 nm.

The phase mask technique not only yields high-performance devices, but is also very flexible in that it can be used to fabricate gratings with controlled spectral response characteristics. For instance, the typical spectral response of a finite-length grating with a uniform index modulation along the fiber length has secondary maxima on both sides of the main reflection peak. In applications like wavelength-division multiplexing, this type of response is not desirable. However, if the profile of the index modulation Δn along the fiber length is given a bell-like functional shape, these secondary maxima can be suppressed.¹⁶ The procedure is called *apodization*. Apodized fiber gratings have been fabricated using the phase mask technique, and suppressions of the sidelobes of 30 to 40 dB have been achieved.^{17,18}

Figure 4 shows the spectral response of two Bragg gratings with the same full width at half-maximum (FWHM). One grating exhibits large sidebands, whereas the other has much-reduced sidebands. The one with the reduced sidebands is a little longer and has a coupling coefficient κ apodized as a second-degree cosine (\cos^2) along its length. Apodization has one disadvantage:

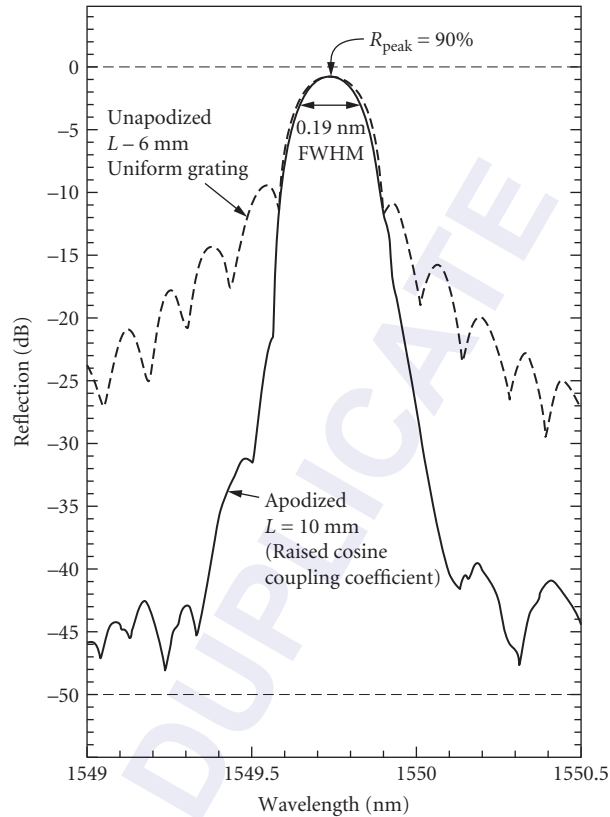


FIGURE 4 Comparison of an unapodized fiber grating's spectral response with that of an apodized fiber grating having the same bandwidth (FWHM).

It decreases the effective length of the Bragg grating. Therefore, to obtain fiber gratings having the same FWHM, the apodized fiber grating has a longer length than the equivalent-bandwidth unapodized fiber grating.

The phase mask technique has been extended to the fabrication of chirped or aperiodic fiber gratings. *Chirping* means varying the grating period along the length of the grating in order to broaden its spectral response. Aperiodic or chirped gratings are desirable for making dispersion compensators¹⁹ or filters having broad spectral responses. The first chirped fiber gratings were made using a double-exposure technique.²⁰ In the first exposure, an opaque mask is positioned between the fiber and the ultraviolet beam blocking the light from irradiating the fiber. The mask is then moved slowly out of the beam at a constant velocity to increase continuously the length of the fiber that is exposed to the ultraviolet light. A continuous change in the photoinduced refractive index is produced that varies linearly along the fiber length with the largest index change occurring in the section of fiber that is exposed to ultraviolet light for the longest duration. In a second exposure, a fiber grating is photoimprinted in the fiber by using the standard phase mask technique. Because the optical pitch of a fiber grating depends on both the refractive index and the mechanical pitch (i.e., optical pitch = $N_{\text{eff}}\Lambda$), the pitch of the photoimprinted grating is effectively chirped, even though its mechanical period is constant. Following this demonstration, a variety of other methods have been developed to manufacture gratings that are chirped permanently^{21,22} or that have an adjustable chirp.^{23,24}

The phase mask technique can also be used to fabricate tilted or blazed gratings. Usually, the corrugations of the phase mask are oriented normal to the fiber axis, as shown in Fig. 3. However, if the corrugations of the phase mask are oriented at an angle to the axis of the fiber, the photoimprinted grating is tilted or blazed. Such fiber gratings couple light out from the bound modes of the fiber to either the cladding modes or the radiation modes. Tilted gratings have applications in fabricating fiber taps.²⁵ If the grating is simultaneously blazed and chirped, it can be used to fabricate an optical spectrum analyzer.²⁶

Another approach to grating fabrication is the *point-by-point technique*,²⁷ also developed at CRC. In this method, each index perturbation of the grating is written point by point. For gratings with many index perturbations, the method is not very efficient. However, it has been used to fabricate micro-Bragg gratings in optical fibers,²⁸ but it is most useful for making coarse gratings with pitches of the order of 100 μm that are required for LP₀₁ to LP₁₁ mode converters²⁷ and polarization mode converters.¹¹ The interest in coarse period gratings has increased lately because of their use in long-period fiber-grating band-rejection filters²⁹ and fiber-amplifier gain equalizers.³⁰

17.6 THE APPLICATION OF FIBER GRATINGS

Hill and Meltz⁶ provide an extensive review of the many potential applications of fiber gratings in lightwave communication systems and in optical sensor systems. Our purpose here is to note that a common problem in using FBGs is that a transmission device is usually desired, whereas FBGs function as reflection devices. Thus, means are required to convert the reflection spectral response into a transmission response. This can be achieved using a Sagnac loop,³¹ a Michelson (or Mach-Zehnder) interferometer,³² or an optical circulator. Figure 5 shows an example of how this is achieved for the case of a multichannel dispersion compensator using chirped or aperiodic fiber gratings.

In Fig. 5a, the dispersion compensator is implemented using a Michelson interferometer. Each wavelength channel ($\lambda_1, \lambda_2, \lambda_3$) requires a pair of identically matched FBGs, one in each arm of the interferometer. Since it is difficult to fabricate identical Bragg gratings (i.e., having the same

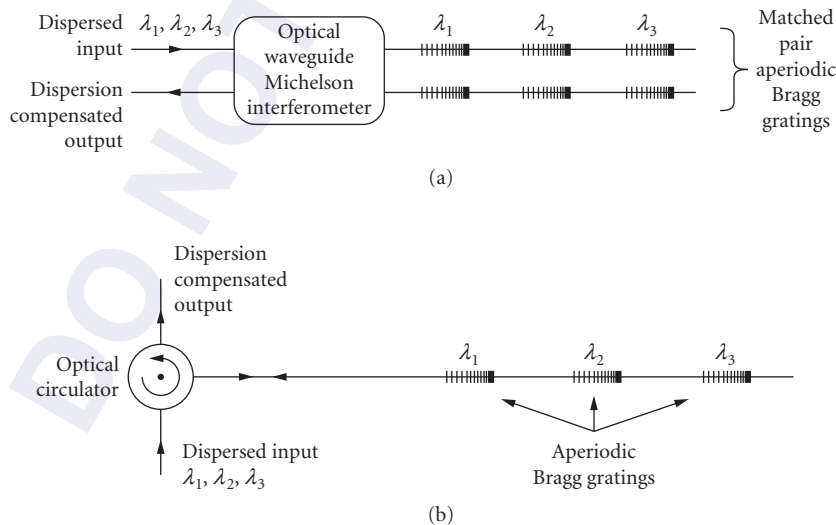


FIGURE 5 Schematic diagram of a multichannel dispersion compensator that is formed by using (a) a Michelson interferometer and (b) an optical circulator.

resonant wavelength and chirp), this configuration for the dispersion compensator has not yet been demonstrated. However, a wavelength-selective device that requires matched grating pairs has been demonstrated.^{33,34} An additional disadvantage of the Michelson interferometer configuration being an interferometric device is that it would require temperature compensation. The advantage of using a Michelson interferometer is that it can be implemented in all-fiber or planar-integrated optics versions.

Figure 5b shows the dispersion compensator implemented using an optical circulator. In operation, light that enters through the input port is routed by the circulator to the port with the Bragg gratings. All of the light that is reflected by the FBGs is routed to the output channel. This configuration requires only one chirped FBG per wavelength channel and is the preferred method for implementing dispersion compensators using FBGs. The only disadvantage of this configuration is that the optical circulator is a bulk optic device (or microoptic device) that is relatively expensive compared with the all-fiber Michelson interferometer.

17.7 REFERENCES

1. K. O. Hill, Y. Fujii, D. C. Johnson, et al., "Photosensitivity in Optical Fiber Waveguides: Application to Reflection Filter Fabrication," *Applied Physics Letters* **32**(10):647–649 (1978).
2. B. S. Kawasaki, K. O. Hill, D. C. Johnson, et al., "Narrow-Band Bragg Reflectors in Optical Fibers," *Optics Letters* **3**(8):66–68 (1978).
3. K. O. Hill, B. Malo, F. Bilodeau, et al., "Photosensitivity in Optical Fibers," *Annual Review of Material Science* **23**:125–157 (1993).
4. I. Bennion, J. A. R. Williams, L. Zhang, et al., "Tutorial Review, UV-Written In-Fibre Bragg Gratings," *Optical and Quantum Electronics* **28**:93–135 (1996).
5. B. Pommellec, P. Niay, M. Douay, et al., "The UV-Induced Refractive Index Grating in Ge:SiO₂ Preforms: Additional CW Experiments and the Macroscopic Origin of the Change in Index," *Journal of Physics D, Applied Physics* **29**:1842–1856 (1996).
6. Kenneth O. Hill and Gerald Meltz, "Fiber Bragg Grating Technology Fundamentals and Overview," *Journal of Lightwave Technology* **15**(8):1263–1276 (1997).
7. A. D. Kersey, M. A. Davis, H. J. Patrick, et al., "Fiber Grating Sensors," *Journal of Lightwave Technology* **15**(8):1442–1463 (1997).
8. T. Erdogan, V. Mizrahi, P. J. Lemaire, et al., "Decay of Ultraviolet-Induced Fiber Bragg Gratings," *Journal of Applied Physics* **76**(1):73–80 (1994).
9. P. J. Lemaire, R. M. Atkins, V. Mizrahi, et al., "High Pressure H₂ Loading as a Technique for Achieving Ultrahigh UV Photosensitivity and Thermal Sensitivity in GeO₂ Doped Optical Fibres," *Electronics Letters* **29**(13):1191–1193 (1993).
10. F. Bilodeau, B. Malo, J. Albert, et al., "Photosensitization of Optical Fiber and Silica-on-Silica Waveguides," *Optics Letters* **18**(12):953–955 (1993).
11. K. O. Hill, F. Bilodeau, B. Malo, et al., "Birefringent Photosensitivity in Monomode Optical Fibre: Application to the External Writing of Rocking Filters," *Electronic Letters* **27**(17):1548–1550 (1991).
12. Alan Miller, "Fundamental Optical Properties of Solids," in *Handbook of Optics*, edited by M. Bass, McGraw-Hill, New York, 1995, vol. 1, pp. 9–15.
13. D. K. W. Lam and B. K. Garside, "Characterization of Single-Mode Optical Fiber Filters," *Applied Optics* **20**(3):440–445 (1981).
14. G. Meltz, W. W. Morey, and W. H. Glenn, "Formation of Bragg Gratings in Optical Fibers by a Transverse Holographic Method," *Optics Letters* **14**(15):823–825 (1989).
15. K. O. Hill, B. Malo, F. Bilodeau, et al., "Bragg Gratings Fabricated in Monomode Photosensitive Optical Fiber by UV Exposure Through a Phase Mask," *Applied Physics Letters* **62**(10):1035–1037 (1993).
16. M. Matsuhara and K. O. Hill, "Optical-Waveguide Band-Rejection Filters: Design," *Applied Optics* **13**(12):2886–2888 (1974).

17. B. Malo, S. Thériault, D. C. Johnson, et al., "Apodised In-Fibre Bragg Grating Reflectors Photoimprinted Using a Phase Mask," *Electronics Letters* **31**(3):223–224 (1995).
18. J. Albert, K. O. Hill, B. Malo, et al., "Apodisation of the Spectral Response of Fibre Bragg Gratings Using a Phase Mask with Variable Diffraction Efficiency," *Electronics Letters* **31**(3):222–223 (1995).
19. K. O. Hill, "Aperiodic Distributed-Parameter Waveguides for Integrated Optics," *Applied Optics* **13**(8): 1853–1856 (1974).
20. K. O. Hill, F. Bilodeau, B. Malo, et al., "Chirped In-Fibre Bragg Grating for Compensation of Optical-Fiber Dispersion," *Optics Letters* **19**(17):1314–1316 (1994).
21. K. Sugden, I. Bennion, A. Molony, et al., "Chirped Gratings Produced in Photosensitive Optical Fibres by Fibre Deformation during Exposure," *Electronics Letters* **30**(5):440–442 (1994).
22. K. C. Byron and H. N. Rourke, "Fabrication of Chirped Fibre Gratings by Novel Stretch and Write Techniques," *Electronics Letters* **31**(1):60–61 (1995).
23. D. Garthe, R. E. Epworth, W. S. Lee, et al., "Adjustable Dispersion Equaliser for 10 and 20 Gbit/s over Distances up to 160 km," *Electronics Letters* **30**(25):2159–2160 (1994).
24. M. M. Ohn, A. T. Alavie, R. Maaskant, et al., "Dispersion Variable Fibre Bragg Grating Using a Piezoelectric Stack," *Electronics Letters* **32**(21):2000–2001 (1996).
25. G. Meltz, W. W. Morey, and W. H. Glenn, "In-Fiber Bragg Grating Tap," presented at the *Conference on Optical Fiber Communications*, OFC'90, San Francisco, CA, 1990 (unpublished).
26. J. L. Wagnier, T. A. Strasser, J. R. Pedrazzani, et al., "Fiber Grating Optical Spectrum Analyzer Tap," presented at the *IOOC-ECOC'97*, Edinburgh, UK, 1997 (unpublished).
27. K. O. Hill, B. Malo, K. A. Vineberg, et al., "Efficient Mode Conversion in Telecommunication Fibre Using Externally Written Gratings," *Electronics Letters* **26**(16):1270–1272 (1990).
28. B. Malo, K. O. Hill, F. Bilodeau, et al., "Point-by-Point Fabrication of Micro-Bragg Gratings in Photosensitive Fibre Using Single Excimer Pulse Refractive Index Modification Techniques," *Electronic Letters* **29**(18):1668–1669 (1993).
29. A. M. Vengsarkar, P. J. Lemaire, J. B. Judkins, et al., "Long-Period Fiber Gratings as Band-Rejection Filters," presented at the *Optical Fiber Communication conference*, OFC'95, San Diego, CA, 1995 (unpublished).
30. A. M. Vengsarkar, J. R. Pedrazzani, J. B. Judkins, et al., "Long-Period Fiber-Grating-Based Gain Equalizers," *Optics Letters* **21**(5):336–338 (1996).
31. K. O. Hill, D. C. Johnson, F. Bilodeau, et al., "Narrow-Bandwidth Optical Waveguide Transmission Filters: A New Design Concept and Applications to Optical Fibre Communications," *Electronics Letters* **23**(9):465–466 (1987).
32. D. C. Johnson, K. O. Hill, F. Bilodeau, et al., "New Design Concept for a Narrowband Wavelength-Selective Optical Tap and Combiner," *Electronics Letters* **23**(13):668–669 (1987).
33. F. Bilodeau, K. O. Hill, B. Malo, et al., "High-Return-Loss Narrowband All-Fiber Bandpass Bragg Transmission Filter," *IEEE Photonics Technology Letters* **6**(1):80–82 (1994).
34. F. Bilodeau, D. C. Johnson, S. Thériault, et al., "An All-Fiber Dense-Wavelength-Division Multiplexer/Demultiplexer Using Photoimprinted Bragg Gratings," *IEEE Photonics Technology Letters* **7**(4):388–390 (1995).

MICRO-OPTICS-BASED COMPONENTS FOR NETWORKING

Joseph C. Palais

*Ira A. Fulton School of Engineering
Arizona State University
Tempe, Arizona*

18.1 INTRODUCTION

The optical portion of many fiber networks requires a number of functional devices, some of which can be fabricated using small optical components (so-called *micro-optic* components). Micro-optic components are made up of parts which have linear dimensions on the order of a few millimeters. The completed functional device may occupy a space a few centimeters on a side. Components to be described in this chapter have the common feature that the fiber transmission link is opened and small (micro-optic) devices are inserted into the gap between the fiber ends to produce a function component. Network components constructed entirely of fibers or constructed in integrated-optic form are described elsewhere in this handbook.

The following sections describe, in order: a generalized component, specific useful network functions, micro-optic subcomponents required to make up the final component, and complete components.

18.2 GENERALIZED COMPONENTS

A generalized fiber-optic component is drawn in Fig. 1. As indicated, input fibers are on the left and output fibers are on the right. Although some components have only a single input port and a single output port, many applications require more than one input and/or output ports. In fact, the number of ports in some devices can be more than 100. The *coupling loss* between any two ports is given, in decibels, by

$$L = -10 \log(P_{\text{out}}/P_{\text{in}}) \quad (1)$$

With respect to Fig. 1, P_{in} refers to the input power at any of the ports on the left and P_{out} refers to the output power at any of the ports on the right. Because we are only considering passive components in this section, P_{out} will be less than P_{in} and the loss will be a positive number.

Insertion loss refers to the coupling loss between any two ports where coupling is desired and *isolation* (or *directionality*) refers to the coupling loss between any two ports where coupling is unwanted. *Excess loss* is the fraction of input power that does not emerge from any of the desired output ports, as expressed in decibels. It is the sum of all the useful power out divided by the input power.



FIGURE 1 The generalized component.

18.3 NETWORK FUNCTIONS

Functions useful for many fiber-optic communications applications are described in the following sections.

Attenuators

Attenuators reduce the amount of power flowing through the fiber system. Both fixed and variable attenuators are available. The applications include testing of receiver sensitivities (varying the attenuation changes the amount of power incident on the receiver) and protecting a receiver from saturating due to excess incident power. Attenuation from a few tenths of a decibe to more than 50 dB are sometimes required.

Power Splitters and Directional Couplers

These devices distribute input power from a single fiber to two or more output fibers. The component design controls the fraction of power delivered to each of the output ports. Applications include power distribution in local area networks (LANs) and in subscriber networks. The most common splitters and couplers have a single input and equal distribution of power among each of two outputs, as shown schematically in Fig. 2*a*. For an ideal three-port splitter (one with no excess loss), half the input power emerges from each of the two output ports. The insertion loss, as calculated from Eq. (1) with a ratio of powers of 0.5, yields a 3-dB loss to each of the two output ports. Any excess loss is added to the 3 dB.

A splitter with more than two output ports can be constructed by connecting several three-port couplers in a tree pattern as indicated schematically in Fig. 2*b*. Adding more splitters in the same manner allows coupling from one input port to 8, 16, 32 (and more) output ports.

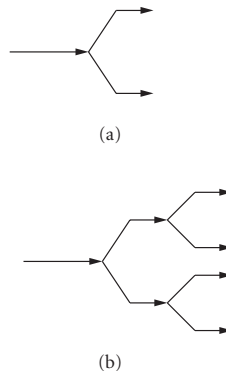


FIGURE 2 Power splitters. (a) 1:2 split and (b) 1:4 split.

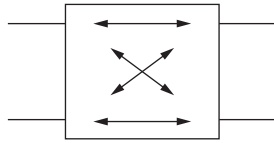


FIGURE 3 Four-port directional coupler.

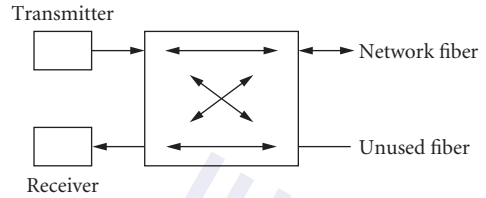


FIGURE 4 LAN terminal illustrating application of the directional coupler.

Adding a fourth port, as in Fig. 3, creates a *directional coupler*. The arrows in the figure show the allowed directions of wave travel through the coupler. An input beam is split between two output ports and is isolated from the fourth. By proper component design, any desired power-splitting ratio can be obtained. One application of the directional coupler is to the distribution network of a local area network, where simultaneous transmission and reception are required. Figure 4 illustrates this usage at one LAN terminal.

Isolators

An isolator is a one-way transmission line. It permits the flow of optical power in just one direction (the forward direction). Applications include protection of a transmitting laser diode from back reflections. Such reflections increase the noise in the system by disrupting the diode's operation. Isolators also improve the stability of fiber amplifiers by minimizing the possibility of feedback, which causes unwanted oscillations in such devices.

Circulators

In a circulator, power into the first port emerges from the second, while power into the second port emerges from a third. This behavior repeats at each successive input port until power into the last port emerges from the first. Practical circulators are typically three- or four-port devices.

Using a circulator, efficient two-way transmission (*full-duplex*) along a single fiber at a single wavelength is possible. The circulator separates the transmitting and receiving beams of light at each terminal, as illustrated in Fig. 5.

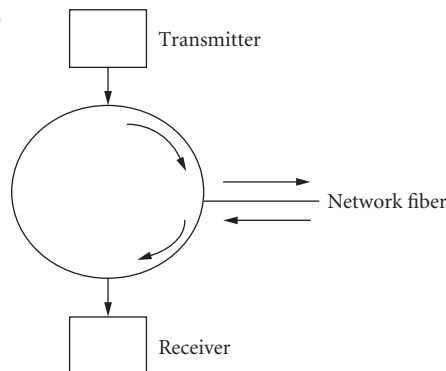


FIGURE 5 An optical circulator separates transmitted and received messages at a terminal.

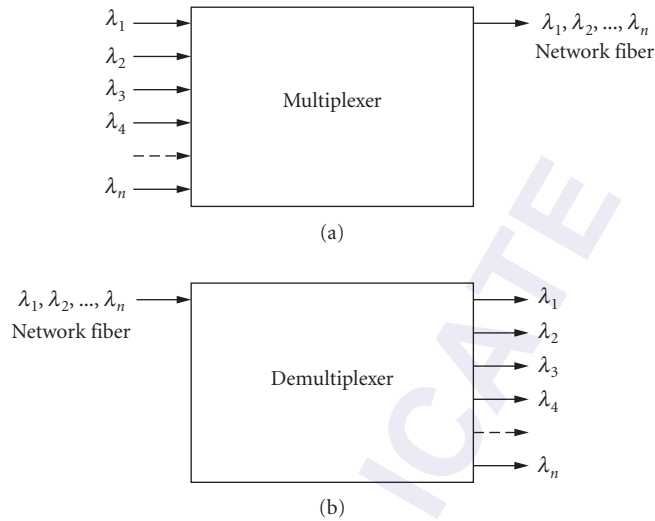


FIGURE 6 (a) A multiplexer combines different wavelength channels onto a single fiber for transmission. (b) A demultiplexer separates several incoming channels at different wavelengths and directs them to separate receivers.

Multiplexers/Demultiplexers/Duplexers

The multiplexer and demultiplexer are heavily utilized in fiber-optic wavelength-division multiplexed (WDM) systems. The *multiplexer* combines beams of light from the different transmitters (each at a slightly shifted wavelength) onto the single transmission fiber. The *demultiplexer* separates the individual wavelengths transmitted and guides the separate channels to the appropriate optical receivers. These functions are illustrated in Fig. 6. Requirements for multiplexers/demultiplexers include combining and separating independent channels less than a nanometer apart, accommodating numerous (in some cases over 100) channels. A frequency spacing between adjacent channels of 100 GHz corresponds to a wavelength spacing of 0.8 nm for wavelengths near 1.55 μm . Insertion losses can be as low as a few tenths of a decibe and isolations of 40 dB or more.

The *duplexer* allows for simultaneous two-way transmission along a single fiber. The wavelengths are different for the transmitting and receiving light beam. The duplexer separates the beams as indicated in Fig. 7, where λ_1 is the transmitting wavelength and λ_2 is the receiving wavelength.

Mechanical Switches

Operationally, an optical switch acts just like an electrical switch. Mechanical movement of some part (as implied schematically in Fig. 8) causes power entering one port to be directed to one of two

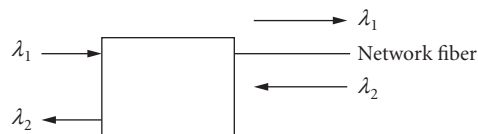


FIGURE 7 A duplexer allows two-way transmission along a single network fiber.

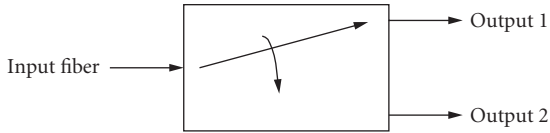


FIGURE 8 Mechanical optical switch.

or more output ports. Such devices are useful in testing of fiber components and systems and in other applications, such as bypassing inoperative nodes in a local area network. Insertion losses less than 0.10 dB and isolations greater than 50 dB are reasonable requirements.

18.4 SUBCOMPONENTS

Micro-optic subcomponents which form part of the design of many complete micro-optic component are described in this section.

Prisms

Because of the dispersion in glass prisms, they can operate as multiplexers, demultiplexers, and duplexers. The dispersive property is illustrated in Fig. 9.

Right-angle glass prisms also act as excellent reflectors, as shown in Fig. 10, owing to perfect reflection (total internal reflection) at the glass-to-air interface. The critical angle for the glass-to-air interface is about 41° and the incident ray is beyond that at 45° .

The beam-splitting cube, drawn in Fig. 11, consists of two right-angle prisms cemented together with a thin reflective layer between them. This beam splitter has the advantage over a flat reflective plate in that no angular displacement occurs between the input and output beam directions. This simplifies the alignment of the splitter with the input and output fibers.

Gratings

Ruled reflection gratings are also used in multiplexers and demultiplexers. As illustrated in Fig. 12, the dispersion characteristics of the grating perform the wavelength separation function required of a demultiplexer. The grating has much greater dispersive power than a prism, permitting increased

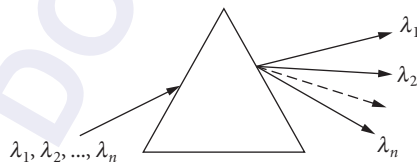


FIGURE 9 A dispersive prism spatially separates different wavelengths. This represents demultiplexing. Reversing the directions of the arrows illustrates combining of different wavelengths. This is multiplexing.

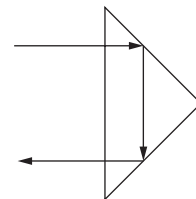


FIGURE 10 Totally reflecting prism.

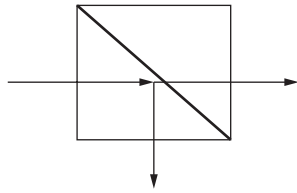


FIGURE 11 Beam-splitting cube.

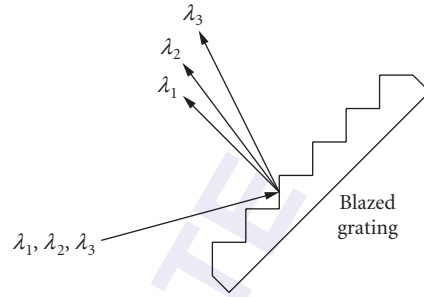


FIGURE 12 Blazed reflection grating operated as a demultiplexer.

wavelength spatial separation. The relationship between the incident and reflected beams, for an incident collimated light beam, is given by the diffraction equation

$$\sin \theta_i + \sin \theta_r = m\lambda/d \quad (2)$$

where θ_i and θ_r are the incident and reflected beam angles, d is the separation between adjacent reflecting surfaces, and m is the *order* of the diffraction. Typically, gratings are blazed so as to maximize the power into the first-order beams. As deduced from Eq. (2) for $m = 1$, the diffracted peak occurs at a different angle for different wavelengths. This feature produces the demultiplexing function needed in WDM systems. Reversing the arrows in Fig. 12 illustrates the multiplexing capability of the grating.

Filters

Dielectric-layered filters, consisting of very thin layers of various dielectrics deposited onto a glass substrate, are used to construct multiplexers, demultiplexers, and duplexers. Filters have unique reflectance/transmittance characteristics. They can be designed to reflect at certain wavelengths and transmit at others, thus spatially separating (or combining) different wavelengths as required for WDM applications.

Beam Splitters

A beam-splitting plate, shown in Fig. 13, is a partially silvered glass plate. The thickness of the silvered layer determines the fraction of light transmitted and reflected. In this way, the input beam can be divided in two parts of any desired ratio.

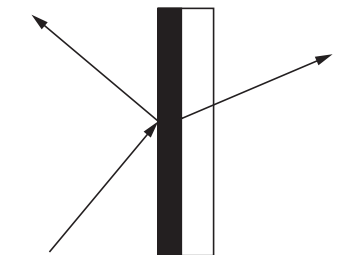


FIGURE 13 Beam-splitting plate.

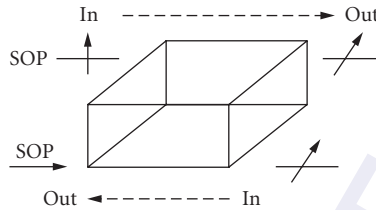


FIGURE 14 Faraday rotator. The dashed arrows indicate the direction of beam travel. The solid arrows represent the wave polarization in the plane perpendicular to the direction of wave travel.

Faraday Rotators

The Faraday rotator produces a nonreciprocal rotation of the plane of polarization. The amount of rotation is given by

$$\theta = VHL \quad (3)$$

where θ is the rotation angle, V is the *Verdet constant* (a measure of the strength of the Faraday effect), H is the applied magnetic field, and L is the length of the rotator. A commonly used rotator material is YIG (yttrium-iron garnet), which has a high value of V .

Figure 14 illustrates the nonreciprocal rotation of the state of polarization (SOP) of the wave. The rotation of a beam traveling from left to right is 45° , while the rotation for a beam traveling from right to left is an additional 45° .

The Faraday rotator is used in the isolator and the circulator.

Polarizers Polarizers based upon dichroic absorbers and polarization prisms using birefringent materials are common. The polarizing beam splitter, illustrated in Fig. 15, is useful in micro-optics applications, such as the optical circulator. The polarizing splitter separates two orthogonally polarized beams.

GRIN-Rod Lens

The subcomponents discussed in the last few sections perform the operations indicated in their descriptions. The problem is that they cannot be directly inserted into a fiber transmission line. To insert one of the subcomponents into the fiber link requires that the fiber be opened to produce a gap. The subcomponent would then fit into the gap. Because the light emerging from a fiber diverges, with

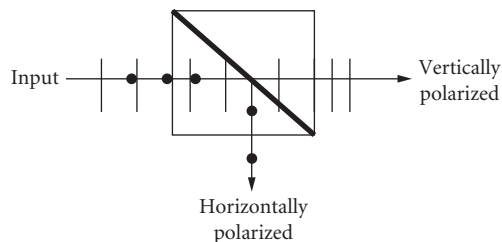


FIGURE 15 Polarizing beam splitter.

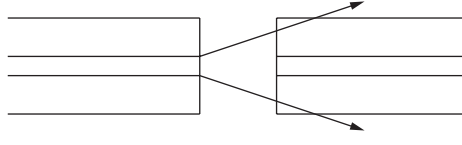


FIGURE 16 Diverging wave emitted from an open fiber couples poorly to the receiving fiber.

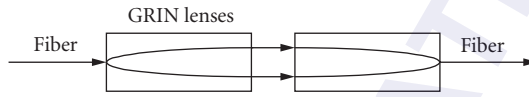


FIGURE 17 Collimating light between two fibers using GRIN-rod lenses.

a gap present the receiving fiber does not capture much of the transmitted light. This situation is illustrated in Fig. 16. The emitted diverging light must be collimated, the required subcomponent (e.g., beamsplitter, grating, etc.) inserted, and the light refocused. A commonly used device for performing this function is the *graded-index rod lens* (GRIN-rod lens). Its use is illustrated in Fig. 17. The diverging light emitted by the transmitting fiber is collimated by the first GRIN-rod lens. The collimated beam is refocused onto the receiving fiber by the second GRIN-rod lens. The collimation is sufficient such that a gap of 20 mm introduces less than 0.5 dB excess loss.¹ This allows for the insertion of beam-modifying devices of the types described in the preceding sections (e.g., prisms, gratings, and beamsplitters) in the gap with minimum added loss.

MEMS Mirrors

Tiny mirrors are the foundation of the *micro-electromechanical systems* (MEMS) optical switch. These mirrors can be fabricated in several ways. Two examples are thin-film mirrors and bulk mirrors, sketched respectively in Figs. 18 and 19. In the thin-film mirror, epitaxial layers are deposited on a silicon substrate. The moveable mirror is formed by removing material from underneath these layers but leaving silicon hinges as indicated in the figure. The bulk mirror is formed by etching it from the silicon substrate.

Mirror movement can be controlled by electrostatic, electromagnetic, piezoelectric, and thermal effects. In electrostatic control, two plates are oppositely charged by placing a voltage across. The resulting attraction causes them to attract and move toward each other. Electromagnetic control uses the forces of attraction between two magnetic circuits. Piezoelectric control is obtained when a voltage placed across a material causes the dimensions of that body to change. Thermal control uses the deformation of a resistive body that occurs when it is heated by passing a current through it.

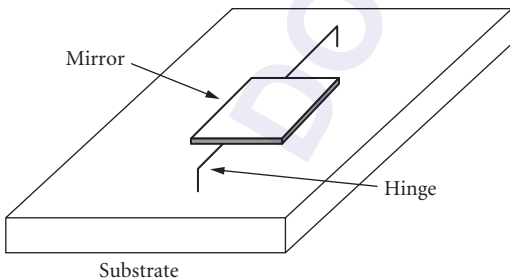


FIGURE 18 Thin-film MEMS mirror.

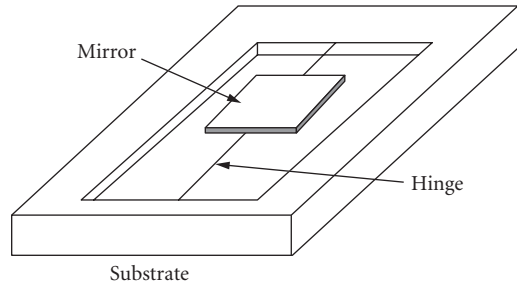


FIGURE 19 Bulk-silicon MEMS mirror.

18.5 COMPONENTS

The subcomponents introduced in the last section are combined into useful fiber devices in the manner described in this section.

Attenuators

The simplest attenuator is produced by a gap introduced between two fibers, as in Fig. 20. As the gap length increases, so does the loss. Loss is also introduced by a lateral displacement. A variable attenuator is produced by allowing the gap (or the lateral offset) to be changeable. A disc whose absorption differs over different parts may also be placed between the fibers. The attenuation is varied by rotating the disk.

In another attenuator design, a small thin flat reflector is inserted at variable amounts into the gap to produce the desired amount of loss.²

Power Splitters and Directional Couplers

A power splitter³ can be constructed as illustrated in Fig. 21. A beam-splitting cube (or a beam-splitting plate) is placed in the gap between two GRIN-rod lenses to connect ports 1 and 2. A third combination of lens and fiber collects the reflected light at port 3. The division of power between the two output fibers is determined by the reflective properties of the splitter itself. Any desired ratio of outputs can be obtained.

If a fourth port is added (Port 4 in Fig. 21) the device is a four-port directional coupler.

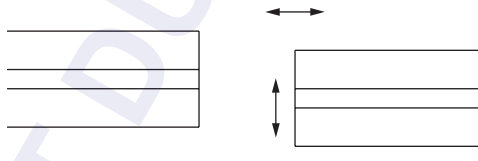


FIGURE 20 Gap attenuator showing relative displacement of the fibers to vary the insertion loss.

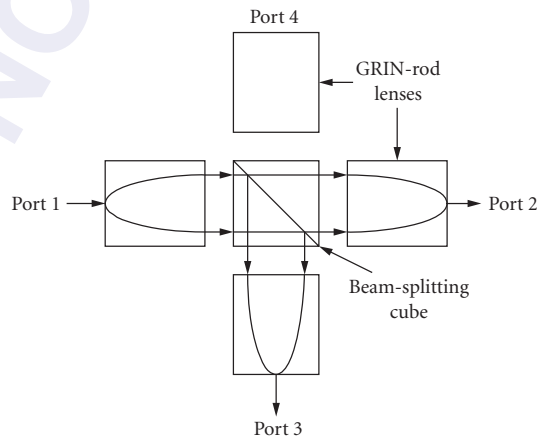


FIGURE 21 Power splitter and (with Port 4 added) four-port directional coupler.

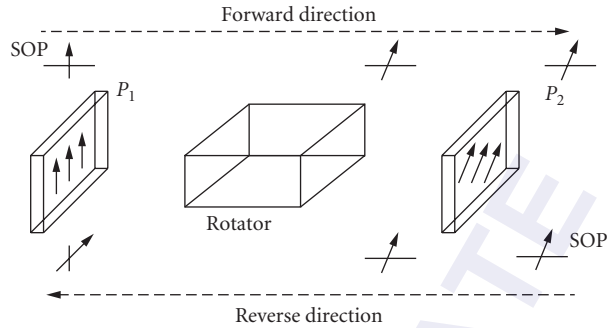


FIGURE 22 Optical isolator. P_1 and P_2 are polarizers.

Isolators and Circulators

The isolator combines the Faraday rotator and two polarizers⁴ as indicated in Fig. 22. The input and output fibers can be coupled to the isolator using GRIN lenses. The vertically polarized beam at the input is rotated by 45° and passed through the output polarizer. Any reflected light is rotated an additional 45° , emerging cross-polarized with respect to the polarizer on the left. In this state, the reflected light will not pass back into the transmitting fiber. Similarly, a light beam traveling from right to left will be cross-polarized at the input polarizer and will not travel further in that direction. The polarizers can be polarizing beam splitters, dichroic polarizers, or polarizing fibers.

A circulator also requires a Faraday rotator and polarizers (polarizing beam splitters or polarizing fiber). Additional components include reflecting prisms, reciprocal 45° rotators, and fiber-coupling devices such as GRIN-rod lenses.⁵

Multiplexers/Demultiplexers/Duplexers

The multiplexer, demultiplexer, and duplexer are fundamentally the same device. The application determines which of the three descriptions is most appropriate. One embodiment is illustrated in Fig. 23 for a two-channel device. As a demultiplexer, the GRIN lens collimates the diverging beam from the network fiber and guides it onto the diffraction grating. The grating redirects the beam according to its wavelength. The GRIN lens then focuses the various wavelengths onto the output fibers for reception. As a multiplexer, the operation is just reversed with the “receiver fibers” replaced by transmitter fibers and all arrows reversed. As a duplexer, one of the two “receiver fibers” becomes a transmitter fiber.

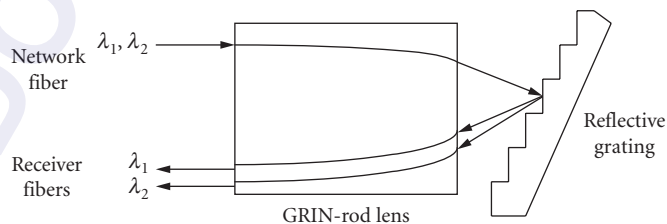


FIGURE 23 Two-channel demultiplexer. Only the beam's central rays are drawn. To operate as a multiplexer the arrows are reversed. To operate as a duplexer, the arrows for just one of the two wavelengths is reversed.

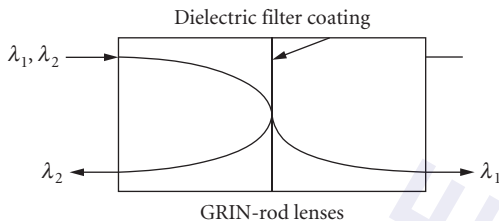


FIGURE 24 Filter-based multiplexer/demultiplexer.

Other configurations also use the diffraction grating, including one incorporating a concave reflector for properly collimating and focusing the beams between input and output fibers.⁶ Micro-optic grating-based devices can accommodate more than 100 WDM channels, with wavelength spacing on the order of 0.4 nm.

A filter-based multiplexer/demultiplexer appears in Fig. 24. The reflective coating transmits wavelength λ_1 and reflects wavelength λ_2 . The device is illustrated as a demultiplexer. Again, by reversing the directions of the arrows, the device becomes a multiplexer. Filter-based multiplexers/demultiplexers can be extended to several channels in the micro-optical form, essentially by cascading several devices of the type just described.

Mechanical Switches⁷

Switching the light beam from one fiber to another one is basically easy. Simply move the transmitting fiber to mechanically align it with the desired receiving fiber. The problem is that even very small misalignments between the two fiber cores introduce unacceptable transmission losses. Several construction strategies have been utilized. Some incorporate a moving fiber and other incorporate a moveable reflector.⁸ In a moveable fiber switch, the fiber can be positioned either manually or by using an electromagnetic force. The switching action in Fig. 25 occurs when the totally reflecting prism moves to align the beam with one or the other of the two output fibers.

A two-dimensional MEMS switch can be constructed as drawn in Fig. 26.⁹ As indicated, light enters the switch from any of the fibers on the left. The mirrors are arranged in an array. Each mirror can be raised into the light path to deflect the beam or can be lowered to allow beam passage. In this manner, light entering any of the input fibers can be directed to any of the output fibers. Collimators (either GRIN or conventional lenses) collimate the entering beams and refocus those exiting.

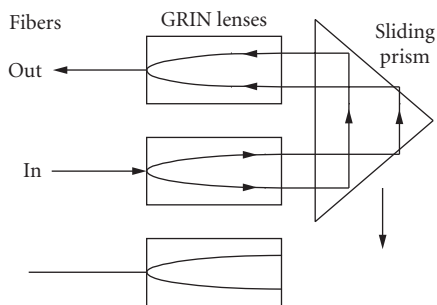


FIGURE 25 Moveable reflecting prism switch.

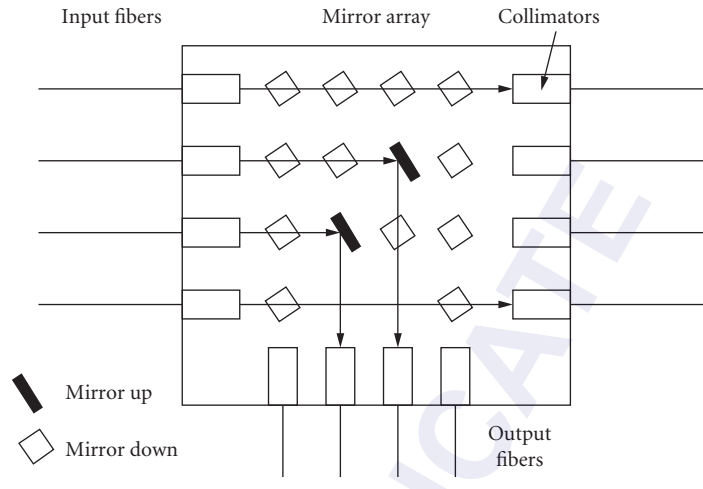


FIGURE 26 Two-dimensional MEMS switch.

18.6 REFERENCES

1. R. W. Gilsdorf, and J. C. Palais, "Single-Mode Fiber Coupling Efficiency with Graded-Index Rod Lenses," *Appl. Opt.* **33**:3440–3445 (1994).
2. C. Marxer, P. Griss, and N. F. de Rooij, "A Variable Optical Attenuator Based on Silicon Micromechanics," *IEEE Photon. Technol. Lett.* **11**:233–235 (1999).
3. C.-L. Chen, *Elements of Optoelectronics and Fiber Optics* (Irwin, Chicago, 1996).
4. R. Ramaswami and K. N. Sivarajan, *Optical Networks: A Practical Perspective* (Morgan Kaufmann, San Francisco, 1998).
5. N. Kashima, *Passive Optical Components for Optical Fiber Transmission* (Artech House, Boston, 1995).
6. J. P. Laude and J. M. Lerner, "Wavelength Division Multiplexing/Demultiplexing (WDM) Using Diffraction Gratings," *SPIE-Application, Theory and Fabrication of Periodic Structures* **503**:22–28 (1984).
7. J. C. Palais, *Fiber Optic Communications*, 5th ed. (Prentice-Hall, Upper Saddle River, N. J., 2005).
8. W. J. Tomlinson, "Applications of GRIN-Rod Lenses in Optical Fiber Communications Systems," *Appl. Opt.* **19**:1123–1138 (1980).
9. Special Issue on Optical MEMS and its Future Trends. *J. Lightwave Tech.* **21**(3) (March 2003).

SEMICONDUCTOR OPTICAL AMPLIFIERS

Jay M. Wiesenfeld

*Bell Laboratories, Alcatel-Lucent
Murray Hill, New Jersey*

Leo H. Spiekman

*Alphion Corp.
Princeton Junction, New Jersey*

19.1 INTRODUCTION

Amplification is a useful, even necessary, function for optical systems, and laser amplifiers have been considered since the early days of quantum electronics. In photonic systems, semiconductor lasers are a fundamental light source. They have the advantage of compactness and efficiency and the ability to emit at many wavelengths, depending on the material system from which they are fabricated. It is therefore natural to make optical amplifiers based on the technology created for semiconductor lasers. Indeed, the first optical amplifier based on a semiconductor laser was reported as early as 1966.¹

A semiconductor optical amplifier (SOA) is a device that provides optical gain based on inversion in a semiconductor medium. Most SOAs are electrically pumped and rely on a p-n junction structure to create a spatial region of optical inversion. As is generally the case for a semiconductor laser, the optical field in a SOA is confined within an optical waveguide, which is an integral part of the structure. In general, the SOA functions as a gain block within an optical system. Again sharing common attributes with the semiconductor laser, SOAs are compact and can be integrated with other elements into photonic integrated circuits (PICs).

SOAs were studied in the early 1980s in the AlGaAs/GaAs material system.² As optical communication systems shifted to the “long wavelength” 1300- and 1550-nm bands, SOAs were intensely studied in the InGaAsP/InP material system.³⁻⁵ This chapter concentrates on SOAs for telecommunication systems, and so on devices based on the InGaAsP/InP material system. However, the general principles and features discussed herein apply to other material systems as well.

In this chapter, Secs. 19.2 to 19.4 describe the basic principles and properties of individual SOAs. Section 19.2 covers the basic device physics. Section 19.3 describes the fabrication of devices and Sec. 19.4 provides a description of the characterization of the properties of SOAs. Sections 19.5 to 19.8 describe some applications for SOAs, primarily related to optical communication systems and networks. Applications related to simple gain for optical communication signals are covered in Sec. 19.6. Section 19.7 covers applications where SOAs are used as gates, with moderate (for switching) and fast (for modulation) gating speeds. Section 19.8 covers some applications that are enabled by the nonlinear gain and index properties that occur within the SOAs.

The illustrative examples in this chapter correspond to devices that operate within the 1300- or 1550-nm regions. In this chapter, when there is a difference between upper- and lower-case symbols, lower case corresponds to linear units and upper case corresponds to decibel (dB) units.

19.2 DEVICE BASICS

The amplification mechanism in the semiconductor optical amplifier is based on the stimulated emission of photons. Just like in any laser, the arrival of a photon in the excited medium prompts the generation of a second photon—a perfect copy in wavelength, phase, and polarization. The basic mechanism of all stimulated emission is the same, but the physical interaction with matter depends on the specific system. For example, in neodymium-doped YAG lasers a transition between different energy levels of an excited Nd ion is used as the energy source for the new photon, and the emission wavelength is determined by the energy difference of those levels. In the SOA on the other hand, the light quantum is generated by an electron in the conduction band of a semiconductor that recombines with a hole in the valence band. For this reason, the wavelength at which photons can be amplified in a SOA is determined fully by the band structure of the semiconductor material. The energy of emitted photons is in a range beginning at the *band gap energy*, as illustrated in Fig. 1. A larger band gap results in shorter wavelength emission, and vice versa.

In order for amplification to occur, a *population inversion* must be present, that is, more carriers must be in the excited state (electrons in the conduction band, leaving holes in the valence band) than in the ground state. Otherwise, emitted photons would be readily reabsorbed. This inversion is usually accomplished by electrical pumping, applying a forward current to a semiconductor diode: electrons are injected from the n-doped side and holes from the p-doped side, and where they meet, they recombine.

In a good optical amplifier this recombination needs to occur in a region of the device through which the optical input signal is propagating, so that it can be efficiently amplified. To this end, the SOA consists of an optical waveguide in which the input signal is confined to a core with a higher index of refraction than the surrounding material. This higher index core has at the same time a lower band gap than the cladding, so that injected carriers are trapped in a potential well and overlap spatially with the signal photons (see Fig. 2). The configuration just described is equivalent to that of a semiconductor laser, and in fact a SOA can be viewed as a minor variation on a laser.

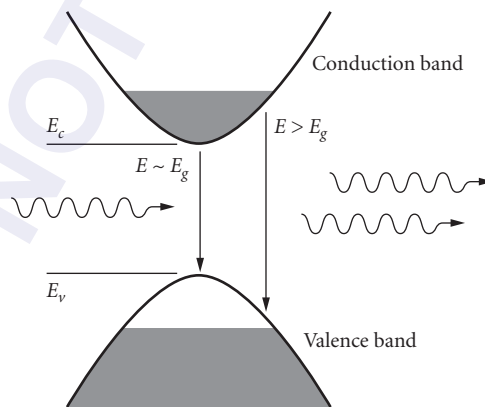


FIGURE 1 Stimulated emission of photons in a semiconductor as a result of population inversion. Recombination of electrons and holes close to the band edges results in emission of photons with an energy close to that of the band gap, while recombination of carriers from higher occupied states within the bands produces photons with a shorter wavelength.

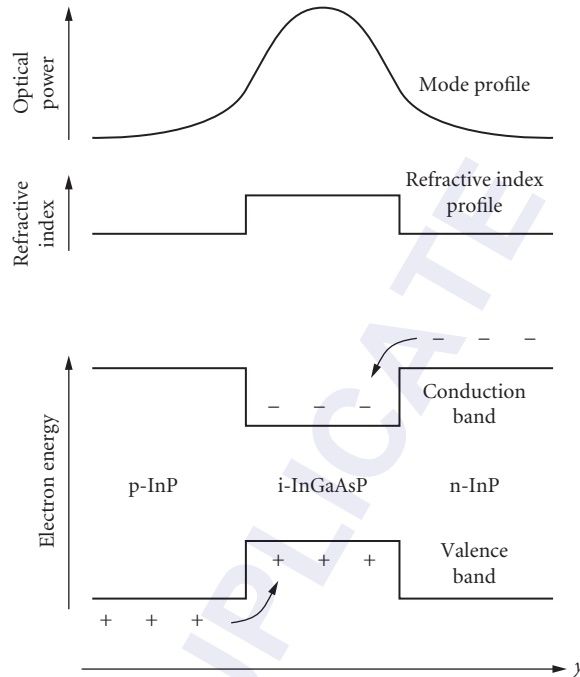


FIGURE 2 Confinement of carriers and photons in a *double heterostructure* waveguide consisting of a lower band gap, higher refractive index material, the *active layer*, sandwiched between layers of higher band gap material with a lower index. By appropriately applying p- and n-doping, a diode structure is formed in which excited states are easily created by injecting a forward current. Note that the lower index cladding material also has a larger band gap and is generally transparent to the emission from the active region.

ASE Noise

In the absence of an input signal, photons are generated in an excited medium by *spontaneous emission*. In a pumped semiconductor this occurs due to the spontaneous recombination of electron-hole pairs. Without these random events, in a laser the lasing action would never start; in a SOA they are a source of noise.

Spontaneous emission occurs over a range of wavelengths corresponding to the occupied excited states of the semiconductor bands, and in all spatial directions. The fraction that couples to the waveguide will subsequently give rise to stimulated emission, and for this reason we speak of *amplified spontaneous emission* (ASE).

In a laser, a feedback mechanism is present that causes the initial ASE to make round trips through the device. When enough current is injected to make this process self-sustaining, lasing action starts. In an optical amplifier, on the other hand, we go to great lengths to avoid optical feedback, so that amplification occurs in a single pass through the device, in a strictly *traveling-wave* fashion. In this case, an ASE spectrum emanates from the device without signs of lasing. An example is shown in Fig. 3. Any residual feedback, for example in the form of reflections from both ends of the amplifier, appears in the spectrum as a ripple, caused by constructive and destructive resonance.

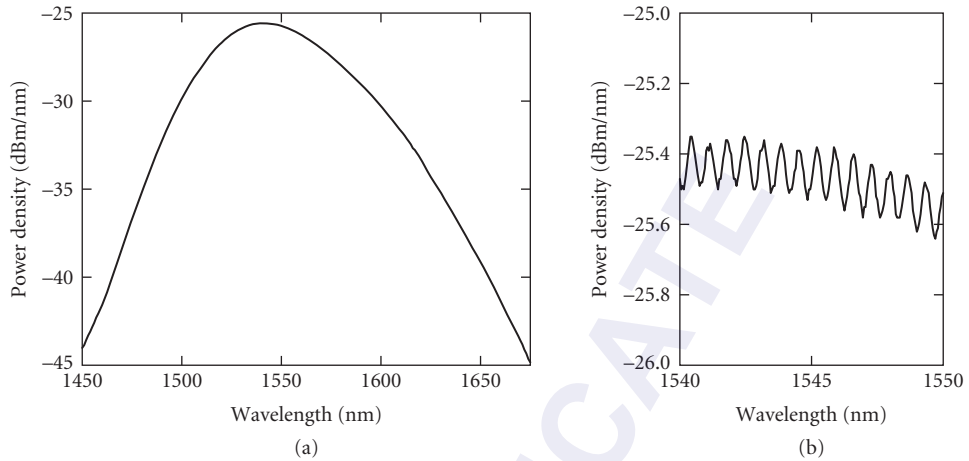


FIGURE 3 (a) Typical amplified spontaneous emission spectrum of a traveling wave SOA. (b) ASE ripple caused by residual reflections from the SOA chip facets.

Gain

A complete SOA typically consists of a semiconductor chip with a waveguide in which the amplification occurs, and two fibers that couple the signal into and out of the chip using lenses, as shown in Fig. 4. A signal coupled into the chip experiences gain as it propagates along the waveguide, according to the process of stimulated emission described above. The chip gain can be written as

$$g_{\text{chip}} = \frac{p_{\text{out}}}{p_{\text{in}}} = e^{(g_{\text{wg}} - \alpha)L} \quad (1)$$

with p_{in} and p_{out} the chip-coupled input and output powers, g_{wg} the gain of the active waveguide per unit length, α a loss term that includes propagation loss and absorption through mechanisms other than producing electron-hole pairs in the active layer, and L the length of the waveguide.

We already saw in Fig. 1 that gain at different wavelengths is generated by electron-hole pairs from different occupied states in the bands. The gain spectrum of the SOA is determined by the semiconductor band structure, and the extent to which it is filled with free carriers.

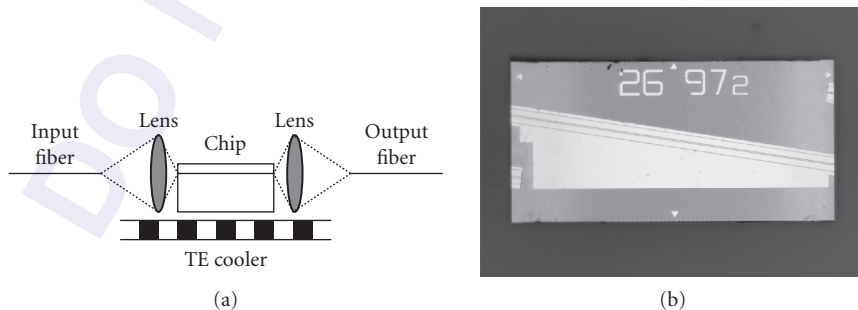


FIGURE 4 (a) Typical semiconductor optical amplifier configuration: lenses or lensed fibers couple the signal into/out of the SOA chip. A thermoelectric cooler (TEC) controls the operating temperature. (b) Photograph of a SOA chip, with the active waveguide visible, as well as the p-side metallization.

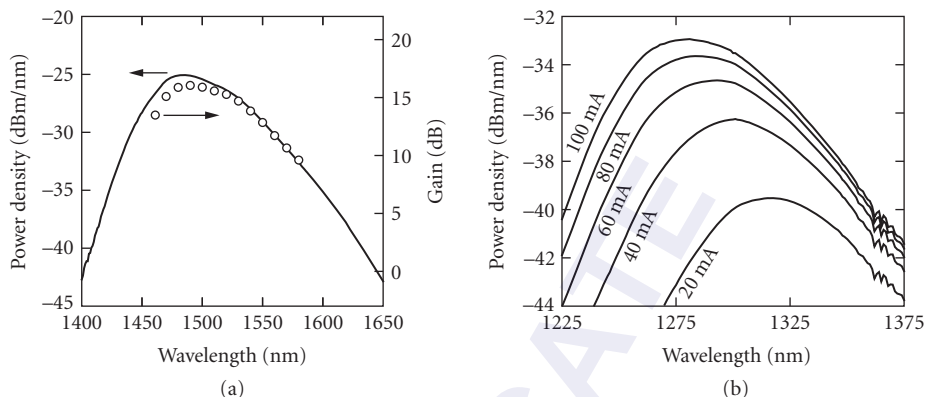


FIGURE 5 (a) Gain and ASE spectra of a SOA plotted to the same scale; only a vertical translation has been applied to match the curves. The mismatch toward the left conveys the larger noise figure of the device at shorter wavelengths. (b) ASE spectrum (of a different SOA) as a function of injected current. The gain near the band edge wavelength hardly changes, but higher current induces gain at shorter wavelength.

Since the ASE spectrum represents spontaneous emission that has been amplified by the same gain that amplifies incoming signal, ASE spectrum and gain spectrum are strongly related. Figure 5 shows the two plotted together. Gain spectrum and ASE spectrum depend on injected current, with higher current filling more states higher in the bands, which extends the gain to shorter wavelengths.

The gain of a SOA depends strongly on temperature. This is the reason why a thermoelectric cooler (TEC) is applied to keep the chip at a nominal operating temperature, often 20 or 25°C (see Fig. 4). At high temperature, free carriers higher in the bands can be ejected out of the potential well formed by the double heterostructure without recombining radiatively (Fig. 2), which has the same effect as lowering the injection current. Figure 6 shows the effect on gain and peak wavelength of varying the chip temperature.

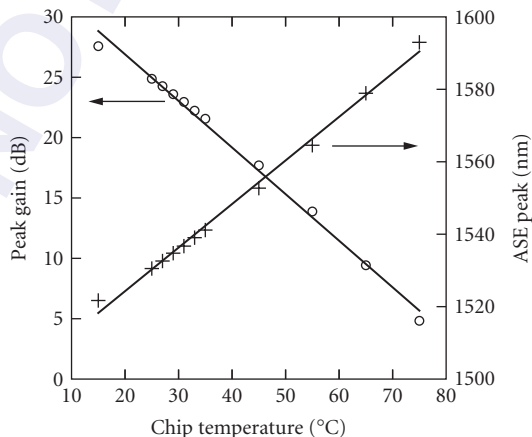


FIGURE 6 Variation of the gain and the ASE peak wavelength with chip temperature. Typical coefficients of -0.4 dB/°C and 1 nm/°C are found in a SOA fabricated in the InGaAsP/InP material system with a bulk active layer.

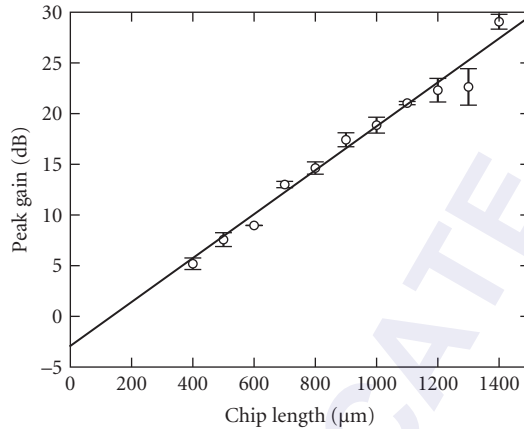


FIGURE 7 Gain versus chip length, for equal chip design and injection current density. The slope of the fit shows that, for this particular chip design and injection current, the on-chip gain is 2.2 dB per 100 μm , and at the 0 μm mark we find the loss caused by fiber-chip coupling, which in this case is about 1.5 dB per side.

In Figs. 5 and 6, gain has been plotted as a fiber-to-fiber number. The on-chip gain is a more fundamental quantity, but it cannot be readily measured without knowledge of the optical loss incurred by coupling a signal from the fiber into the chip and vice versa. Figure 7 shows fiber-to-fiber gain for many chips of different length but otherwise equal design, which allows us to deduce on-chip gain per unit length, as well as obtain an estimate for the fiber-chip coupling loss.

Confinement Factor

The gain per unit length as mentioned in the preceding paragraph is related to the material gain of the active region through the *confinement factor* Γ . This represents the fraction of the total power propagating in the waveguide that is confined to the active region, that is, it can be written as

$$\Gamma = \frac{\text{power in active region}}{\text{total power}} \quad (2)$$

The confinement factor thus relates the waveguide gain to the material gain as $g_{\text{wg}} = \Gamma g_{\text{mat}}$. In the example shown in Fig. 7, the confinement factor is $\Gamma = 0.1$. Therefore the waveguide gain of 22 dB/mm implies a material gain of 220 dB/mm.

The relation between material gain and the free carrier density n can be written in first order approximation as

$$g_{\text{mat}} = g_0(n - n_0) \quad (3)$$

where $g_0 = dg/dn$ is the differential gain, and n_0 is the free carrier density needed for material transparency, that is, the number of free carriers for which the rate of absorption just equals the rate of stimulated emission. g_0 and n_0 can vary with wavelength and temperature. The waveguide loss α varies only very weakly with wavelength.

Polarization Dependence

A SOA does not always amplify light in different polarization states by the same amount. The reason for this is a dependence on the polarization state of the confinement factor, which in turn is caused by the different waveguide boundary conditions for the two polarization directions giving rise to different mode profiles.

The principal polarization states of a planar waveguide are those in which the light is linearly polarized in the horizontal and vertical direction. The waveguide mode in which the electric field vector is predominantly in the plane of the substrate of the device is called the transverse electric (TE) polarization, while the mode in which the electric field is normal to the substrate is called transverse magnetic (TM).

In isotropic active material, such as the zincblende structure for common III-V crystals, the material gain is independent of the polarization direction, and is the same for the TE and TM modes. However, the rectangular shape of the waveguide, which is usually much more wide than it is high (typically around $2\ \mu\text{m}$ wide while only around $100\ \text{nm}$ thick), causes the confinement factor to be smaller for the TM mode than for the TE mode, sometimes by 50 percent or more. Without any mitigating measures, this would result in a significant polarization-dependent gain (PDG) (see Fig. 8).

Several methods exist with which polarization-independent gain can be achieved. The most straightforward one is to use a square active waveguide. This ensures symmetry between the TE and TM modes, and thus equal gain for both.⁶ This approach is not very practical, though, because the waveguide dimensions have to be kept very small (around $0.5 \times 0.5\ \mu\text{m}$) to keep it from becoming multimode, and even though thin layers can be produced in crystal growth with high accuracy, the same accuracy is not available in the lithographic processes that define the waveguide width.

Another method is to introduce a material anisotropy that causes the material gain to become more favorable for the TM polarization direction in the exact amount needed to compensate for the waveguide anisotropy. This can be done by introducing tensile crystal strain in the active layer during material growth, a fact that was first discovered when lasers based on tensile-strained quantum wells were found to emit in the TM polarization direction.⁷

Introducing tensile strain in bulk active material modifies the shape of the light-hole and the heavy-hole bands comprising the valence band of the semiconductor in such a way that TE gain is somewhat reduced, while TM gain stays more or less constant. As a result, with increasing strain the PDG is reduced, and it can even overshoot, yielding devices that exhibit a TM gain higher than their TE gain. Appropriately optimized, this method can yield devices with a PDG close to 0 dB.⁸

In a multiquantum well (MQW) active layer, the same method can be used by introducing tensile strain into the quantum wells.⁹ Alternatively, a stack of QWs alternating between tensile and compressive strain can be used. The compressive wells amplify only TE, while the tensile wells predominantly amplify TM. This way, the gain of TE and TM can be separately optimized, and low-PDG structures can be obtained.^{10,11}

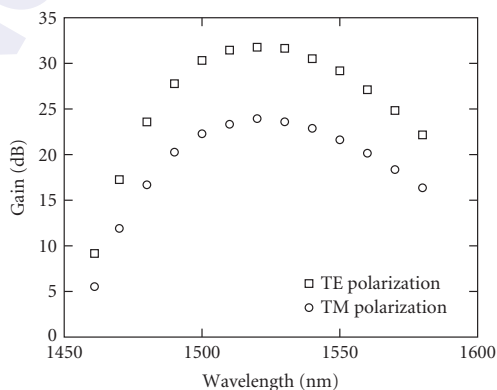


FIGURE 8 Gain versus wavelength in a SOA with significantly undercompensated polarization dependence.

Gain Ripple and Feedback Reduction

Reflections from the chip facets can cause resonant or antiresonant amplification, depending on whether a whole number of wavelengths fit in the cavity. This behavior shows up in the ASE spectrum, as shown in Fig. 3b. The depth of this gain ripple is given by⁴

$$\text{Ripple} = \frac{(1 + gr)^2}{(1 - gr)^2} \quad (4)$$

in which g is the on-chip gain experienced by the guided mode, and r is the facet reflectivity. Obviously, gain ripple becomes more of a problem for high-gain SOAs. A device with an on-chip gain of 30 dB will need the facet reflectivities to be suppressed to as low as 5×10^{-6} in order to show less than 0.1 dB ripple.

Several methods are used to suppress facet reflections (see Fig. 9); the most well-known one is to apply antireflection (AR) coatings onto the facets. An AR coating is a dielectric layer or stack of layers that is designed such that destructive interference occurs among the reflections of all its interfaces.

For a planar wave, a quarter-wave layer with a refractive index that is the geometric mean of the indices of the two regions it separates is a perfect AR coating. Such a design is also a reasonable first approximation for guided waves, but for ultralow reflectivity, careful optimization to the actual mode field needs to be done, taking account of the fact that the optimum for the TE and TM modes may be different.

Since any light that is reflected at an interface is not transmitted through it, AR coatings also help lowering the coupling loss to fiber. The approximately 30 percent reflection of an InP-air interface in absence of a coating would represent a loss of 1.5 dB.

A second method is to angle the SOA waveguide on the chip. Such an *angled stripe* design makes the reflected field propagate backward at double the angle with respect to the waveguide, rather than being reflected directly into the waveguide, and therefore only a small fraction will couple back into the guided mode. Note that a consequence of using an angled stripe is that the output light will emanate from the SOA at a larger angle according to Snell's law, and an appropriate angle of the fiber assemblies will have to be provisioned, which can be as large as 35° for an on-chip angle of 10° .

Another way to reduce reflections back into the waveguide is to end the waveguide a few micrometers before the facet and continue with only cladding material; this is often done in combination with a taper that enhances mode matching to the fiber. This so-called *window structure* allows the modal field to diverge before it hits the facet, so that the reflected field couples poorly back into the waveguide.

For low-gain SOAs, applying only an AR coating suffices. It has also been shown that only an angled stripe (without AR coating) can be sufficient.¹² But when ultralow reflectivity is needed, it is common to find an AR coating being combined with an angled stripe, or a window region, or both. Reflectivities as low as 2×10^{-6} have been obtained in this way.¹³

It has to be noted that mode-matching techniques to enhance fiber-chip coupling efficiency, such as lateral tapers, have an effect on reflections. A larger (usually better-matched) mode diffracts at a smaller angle, improving the effectiveness of an angled stripe, but reducing the effectiveness of a window region.

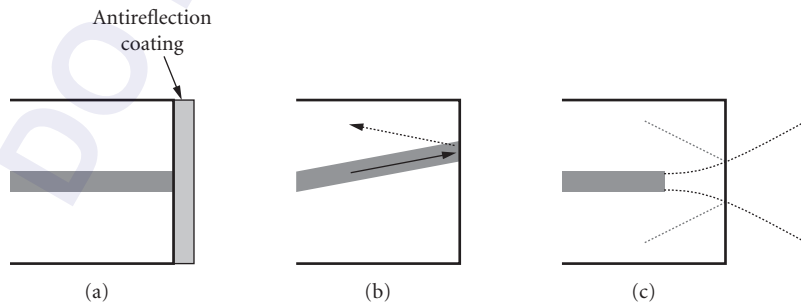


FIGURE 9 Suppressing facet reflections: (a) antireflection coating; (b) angled stripe; and (c) window structure.

Noise Figure

An optical amplifier's noise figure depends on its inversion factor as $nf = 2n_{sp}/\eta_i$, with n_{sp} the inversion factor (equal to one for full carrier inversion) and η_i the optical transmission ($1 - \text{loss}$) at the input of the amplifier. An ideal fully inverted amplifier with zero input loss ($\eta_i = 1$) would have a $nf = 2$ (or, expressed in decibels, $NF = 3$ dB). In a SOA, η_i consists of the fiber-chip coupling coefficient (see Fig. 4a), which can be significantly different from unity. This is the reason why NF is usually somewhat higher for SOAs compared to fiber amplifiers.

The conventional interpretation of the noise figure is that it is the signal-to-noise ratio at an amplifier's output divided by that at its input, for a shot-noise-limited input signal. For optical amplifiers this definition is not very practical, since signals in optical networks are seldom shot-noise-limited.

A more practical definition is based on the approximation that in an optically amplified, optically filtered transmission line, the noise in the receiver is dominated by signal-spontaneous emission beat noise.¹⁴ This results in a definition of noise figure as

$$nf = \frac{2\rho_{\text{ASE}}}{gh\nu} \quad (5)$$

in which g is the gain, ρ_{ASE} is the power spectral density of the amplifier's ASE noise, and $h\nu$ is the photon energy. Only the noise power copolarized with the signal is taken into account, since noise with a polarization orthogonal to that of the signal does not give rise to beat noise in the detector. All quantities in this expression can be easily measured, which allows for straightforward characterization of a device's noise figure (see "Noise Figure" in Sec. 19.4).

Saturation

In an amplifier, the gain depends on the amplified signal power, which at high values causes the output power to saturate. A strong input signal causes the stimulated emission to reduce the carrier density, which decreases the gain and at the same time shifts the gain peak to longer wavelengths, closer to the band gap emission wavelength of the active stripe. This gain compression can be written as a function of the output power p_o as follows:¹⁵

$$g = g_{ss} e^{-p_o/p_{\text{sat}}} \quad (6)$$

with g_{ss} the small-signal gain (assumed to be ≥ 1), and

$$p_{\text{sat}} = \frac{h\nu A \eta_o}{\tau \Gamma dg/dn} \quad (7)$$

the characteristic saturation output power, which depends on the carrier lifetime τ , the confinement factor Γ , the differential gain dg/dn , the cross-section area A of the active stripe, and the output coupling efficiency η_o .

A convenient description of the saturation power of an optical amplifier is given by the output power at which the gain is reduced by a factor of two, or 3 dB. This so-called 3-dB saturation output power can now simply be written as

$$p_{3\text{dB}} = \ln 2 \cdot p_{\text{sat}} \quad (8)$$

Figure 10 shows an example of a measured gain versus output power curve, from which the small-signal gain of the amplifier and the 3-dB saturation power can be directly determined.

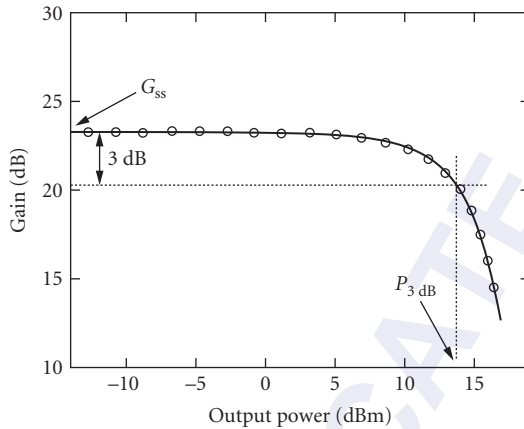


FIGURE 10 Gain compression curve of a SOA. At small powers, the gain approaches the value of the small-signal gain G_{ss} . The output power corresponding to a gain compression of a factor of two is the 3-dB saturation output power $P_{3\text{dB}}$.

Increasing the injection current into the SOA increases the saturation power through reduction of the carrier lifetime τ and reduction of the differential gain dg/dn . An increase can also be accomplished by using a SOA with gain peak significantly shorter than the signal wavelength. Due to band filling, the differential gain is smaller on the red side of the gain peak, which causes the saturation power to be larger at longer wavelengths (see Fig. 11).

Structurally, the saturation power of a SOA can be increased by reducing the thickness of the active layer.¹⁶ The optical field expands widely in the vertical direction, which decreases the optical confinement factor Γ much faster than it decreases the active cross-section A .

In the horizontal direction, the field is usually much better confined. Therefore another effective way to increase P_{sat} is to use a flared gain stripe, that tapers to a much larger width at the output. This increases the active cross-section much faster than it does the confinement factor. Since the waveguide at the output will typically be multimode, care has to be taken in the design of the taper.

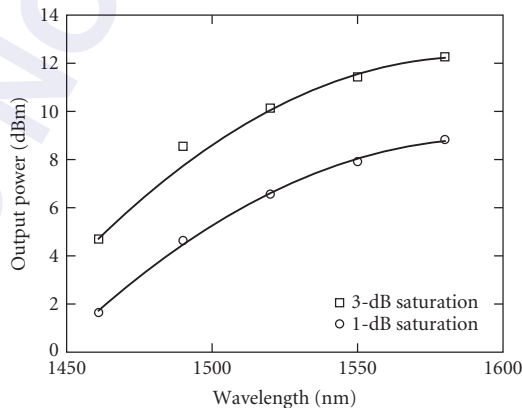


FIGURE 11 Saturation power of a SOA versus signal wavelength. The smaller differential gain at longer wavelengths causes an increase in P_{sat} .

Material Systems

Semiconductor optical amplifiers are most commonly fabricated in the InGaAsP/InP material system. The active and other waveguiding layers, as well as electrical contacting layers are epitaxially grown on an InP substrate. InGaAsP is chosen because it allows the emission wavelength to be chosen in the range 1250 to 1650 nm, which contains a number of bands that are important for telecommunications.

Since the gallium atom is slightly smaller than the indium atom, whereas the arsenic atom is slightly larger than the phosphorus atom, by choosing the element ratios In:Ga and As:P properly, a crystal lattice can be formed with the same lattice constant as InP. The remaining degree of freedom among these *lattice-matched* compositions is used to tune the band gap, which is direct over the full range from binary InP to ternary InGaAs, hence the ability to form 1250 to 1650 nm emitters. For emission at wavelengths such as 850 or 980 nm, the GaAs/AlGaAs material is commonly used.

A *quantum well* may be created by sandwiching a thin layer of material between two layers with wider band gap. This forms a potential well in which the free carriers may be confined, leaving them only the plane of the quantum-well layer to move freely (see Fig. 12). Even though quantum wells are used almost exclusively for the fabrication of semiconductor lasers, both quantum well and bulk

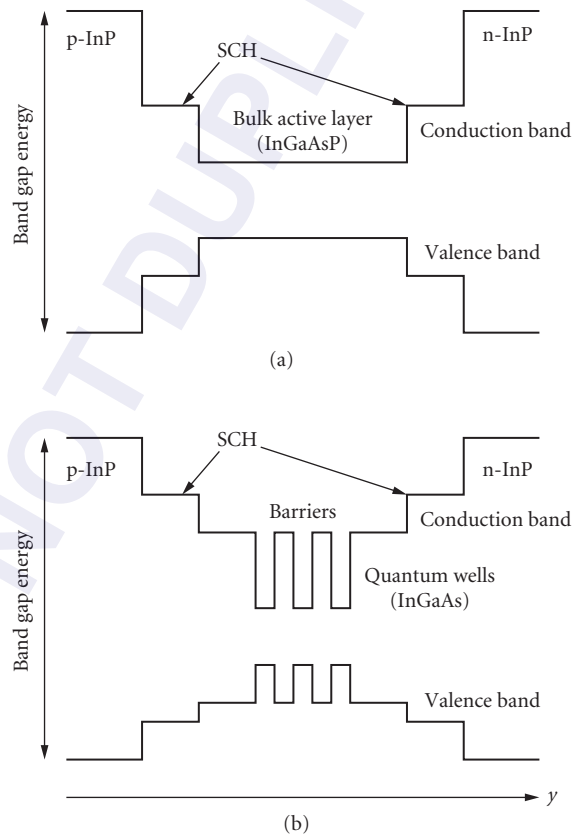


FIGURE 12 (a) Layer structure of a SOA waveguide with a bulk active layer and (b) structure of a MQW SOA. Note that the band offsets in the conduction and valence bands are not drawn to scale.

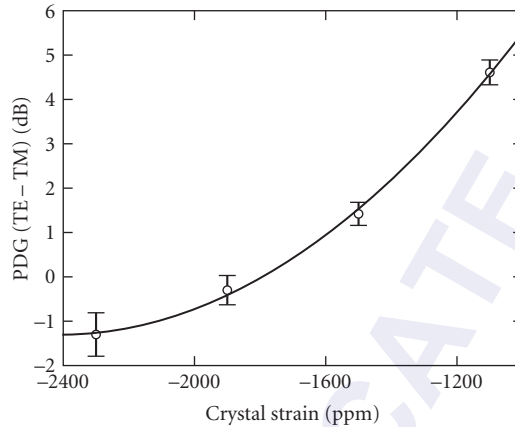


FIGURE 13 Polarization dependence in SOAs with bulk active layers with varying amounts of tensile strain.

active structures are used for SOAs, as the advantages of quantum wells are less pronounced in amplifiers than they are in lasers. Recently, active structures have been demonstrated based on *quantum dots*. These dots confine the electrons not in one, but in all three spatial directions, giving rise to delta function-like density of states. This property is expected to lead to devices with reduced temperature dependence with respect to bulk or quantum well devices.

As mentioned earlier, strain can be used in the active layer to tune the polarization-dependent gain of a SOA. Strain is introduced by deviating from the layer compositions that would yield the active layer lattice-matched, either by introducing a larger fraction of the larger elements, to produce compressive strain, or by emphasizing the smaller atoms, to produce tensile strain. Figure 13 shows the effect of strain on PDG in a bulk active SOA.

Gain Dynamics

The dynamic behavior of a SOA is governed by the time constants associated with the various processes its free carriers can undergo. The carrier lifetime τ has already been mentioned. This is the characteristic time associated with *interband* processes such as spontaneous emission and the electrical pumping of the active layer, that is, with the movement of electrons between the valence and the conduction band. The carrier lifetime is of the order of 25 to 250 ps.

Intraband processes such as spectral hole burning and carrier heating, on the other hand, govern the (re)distribution of carriers inside the semiconductor bands. These processes are much faster than the carrier lifetime.¹⁷

Dynamic effects can be a nuisance when one only wants to amplify modulated signals, because they introduce nonlinear behavior that leads to intersymbol interference. But they can be used advantageously in various forms of all-optical processing. Using a strong signal to influence the gain of the amplifier, one can affect the amplitude of other signals being amplified at the same time. An example of this *cross-gain modulation* (XGM) is shown in Fig. 14. In the gain-recovery measurement, the gain of the SOA is reduced almost instantaneously as the pump pulse sweeps the free carriers out of the active region. After the pulse has passed, the gain slowly recovers back to its original value.

The gain-recovery time depends on the design of the SOA and the injection current, as shown in Fig. 15. Cross-gain modulation can support all-optical processing for signals with data rates higher than 100 Gb/s.

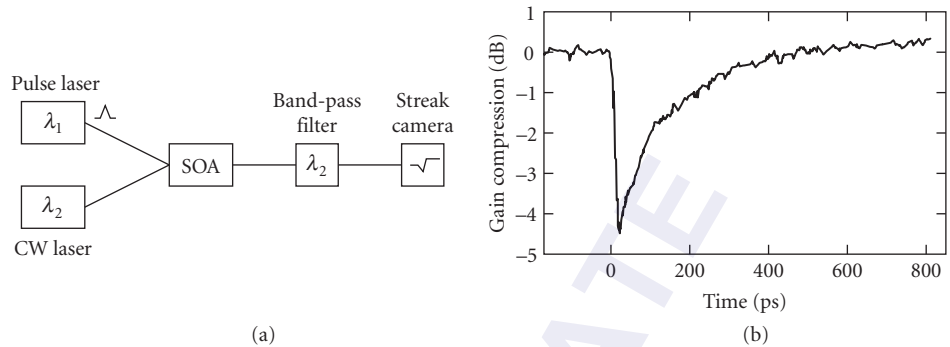


FIGURE 14 (a) Gain recovery experiment in which an intense pulse (the pump) compresses the gain of a SOA, which is measured by a weak probe beam and (b) gain compression and recovery at λ_2 .

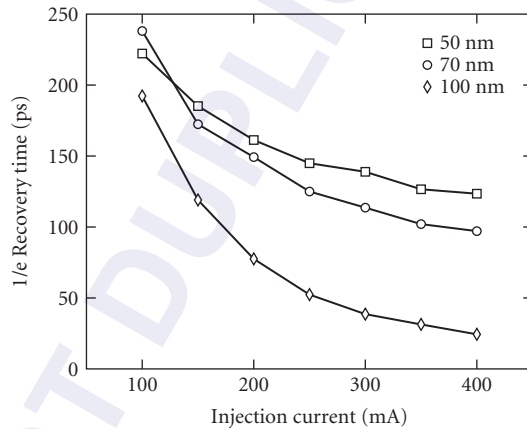


FIGURE 15 Cross-gain modulation recovery times versus active-layer thickness and current injection. The SOAs have a bulk active layer and a gain peak at 1550 nm. Chip length is 1 mm.

Along with the gain change caused by a strong input signal, there is a phase change associated with the refractive index difference caused by the removal of free carriers, which results in heavy chirping of signals optically modulated by the XGM. However, this *cross-phase modulation* (XPM) can also be used to advantage. Only a small gain change is needed to obtain a π phase shift, so all-optical phase modulation can be obtained without adding much amplitude modulation. Using a waveguide interferometer, this phase modulation can be converted back to on-off keying.

Intraband processes give rise to effects like four-wave mixing (FWM). This is an interaction between wavelengths injected into a SOA that creates photons at different wavelengths (see Fig. 16). A straightforward way to understand FWM is as follows. Two injected pump beams create a moving beat pattern of intensity hills and valleys, which interacts with the SOA nonlinearities to set up a moving grating of minima and maxima in refractive index. Photons in either beam can be scattered by that moving grating, creating beams at lower or higher frequency, spaced by the frequency difference between the two pump beams.

More detail on applications of the nonlinearities will be described in Sec. 19.8.

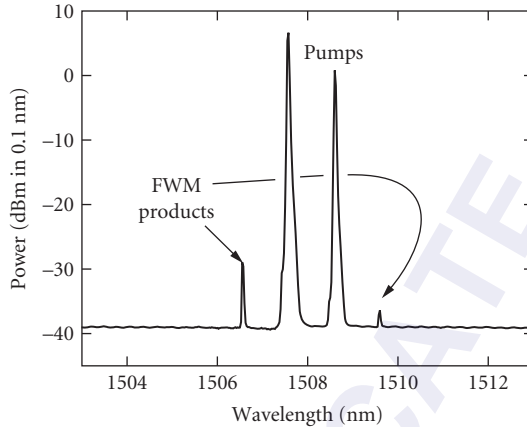


FIGURE 16 Four-wave mixing in a SOA. The two center pump beams give rise to mixing products on both sides.

Gain Clamping

One possible solution to limit intersymbol interference caused by SOA nonlinearities is to resort to gain clamping. When controlled lasing is introduced in a SOA, the gain is clamped by virtue of the lasing condition, and no gain variations are caused by modulated input signals.

Lasing can be introduced in a SOA by etching short gratings at both ends of the active waveguide.^{18,19} The feedback wavelength is set to fall outside the wavelength band of interest for amplification, and the grating strength defines the round-trip loss, and therefore the level at which the gain is clamped (see Fig. 17). Now when an input signal is introduced into the gain-clamped SOA (GC-SOA), the gain will not change as long as the device is lasing. As the amplified input signal takes up more power, less carriers are available to support the lasing action. Only when so many carriers are used by the signal to make the laser go below threshold, will gain start to drop.

The steady-state picture sketched needs to be augmented to account for the dynamic behavior of the clamping laser. Relaxation oscillations limit the effectiveness of gain clamping. In order to support amplification of 10-Gb/s NRZ on-off keying modulated signals, the GC-SOA is designed with its relaxation oscillation peak at 10 GHz, where the modulation spectrum has a null.

A different type of GC-SOA has its clamping laser operating vertically. This device, called linear optical amplifier (LOA),²⁰ has a vertical cavity surface emitting laser (VCSEL) integrated along the full length of the active stripe. This design has the advantage that the clamping laser line is not present in the amplifier output, as it emits orthogonally to the propagation direction of the amplified signals.

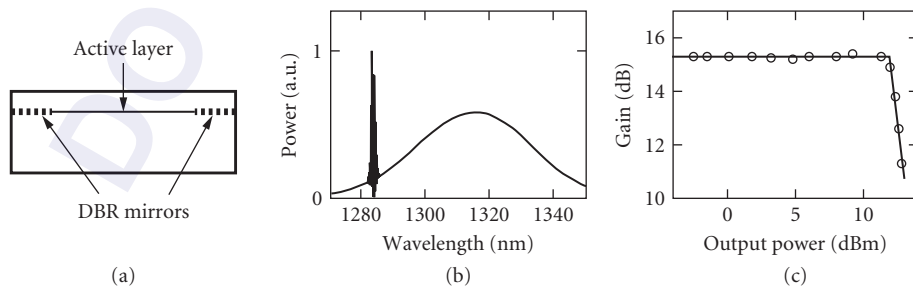


FIGURE 17 Gain-clamped SOA from Ref. 19: (a) schematic; (b) output spectrum; and (c) gain versus output power curve.

The gain clamping laser in this case has a relaxation oscillation that varies over the length of the device. For this reason, no hard relaxation oscillation peak is observed. At the same time, the large spontaneous emission factor of the vertical laser makes the clamping level less well-defined: rather than staying absolutely constant up to the point of going below threshold, it causes a softer knee in the gain versus output power curve, somewhat in between the horizontal GC-SOA case (Fig. 17c) and the case of an unclamped SOA (Fig. 10).

In practice, the “sweet spot” for amplification of digitally modulated signals is at an output power corresponding to around 1 to 3 dB of gain compression,²¹ depending on data rate and modulation format. At this point in the gain versus output power curve, the LOA has no output power advantage over standard, unclamped SOAs. For analog transmission, the gain has to stay absolutely constant; this has only been attempted with horizontally clamped GC-SOAs.²²

19.3 FABRICATION

The fabrication of SOAs is a wafer-scale process that is very similar to the manufacturing of semiconductor laser diodes. First, epitaxial layers are grown on a semiconductor substrate. Then, waveguides are formed by etching, followed by one or more optional regrowth steps. Finally, p- and n-side metalization is applied.

Waveguide Processing

InGaAsP/InP SOAs, like their laser counterparts, mostly use buried waveguide structures that are fabricated in a standard buried heterostructure (BH) process: A mesa is etched in the epitaxial layer stack containing the gain stripe using a dielectric mask. Using this mask, selective regrowth is performed in the regions next to the waveguide, in order to form current blocking layers that force all injection current to flow through the active stripe. This is accomplished either using semi-insulating material, for example, Fe:InP, or with a p-n structure that forms a reverse-biased diode. After removing the etching and regrowth mask, a p-doped InP top layer is grown to provide for the p-contact (see Fig. 18).

GaAs/AlGaAs structures are not easily overgrown. In this material system, ridge waveguides are usually used. Usually the epi layers are grown in a single growth step with a thick p-doped top layer already in place. Waveguides are formed by etching, after which the whole structure is covered with a passivation layer, in which contact openings are etched on top of the ridge. (See Fig. 19 for an example of a ridge structure grown on InP.)

The waveguide pattern on a SOA wafer usually includes angled stripes to reduce feedback into the waveguide from facet reflections. For the same reason, a buried waveguide may contain window structures, in which the waveguide ends a few micrometers before the facet, the remaining distance

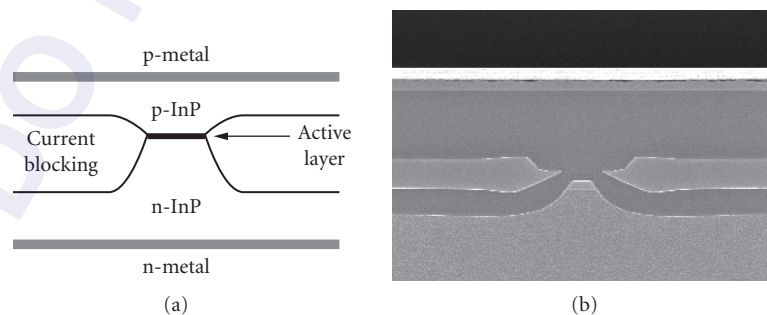


FIGURE 18 (a) Schematic of planar buried heterostructure and (b) SEM photograph of BH waveguide.

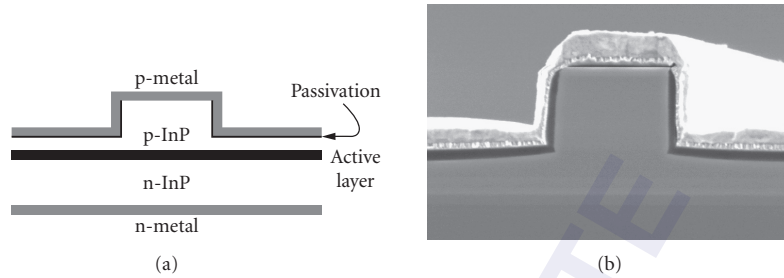


FIGURE 19 (a) Schematic of ridge waveguide structure and (b) SEM photograph of ridge waveguide.

being bridged by nonguided propagation in the regrown InP. The waveguides may contain tapers near the facets to shape the mode for improved fiber-chip coupling efficiency, and may be flared as described earlier to obtain higher saturation power.

Metallization

After waveguide processing, p-side metal is applied on top of the wafer, which is patterned to form contact pads. Before n-side metal is applied to the back side of the wafer, it is usually thinned in a lapping and polishing process to improve cleaving yield, as well as thermal conductivity between the active layer and the heatsink on which the devices will be mounted. Figure 20 shows a photograph of a finished SOA chip on submount.

Postprocessing and Package Assembly

A fully processed 2-in wafer can contain thousands of individual SOA devices. Before these can be used, a facet coating needs to be applied for passivation and to reduce reflections. To this end, the wafer is cleaved into bars containing many SOAs, which are subsequently antireflection coated. An AR coating is nominally a quarter-wave layer of dielectric material, which reduces reflections from the active waveguide back into the chip. In combination with an angled stripe and possibly a window region, both of

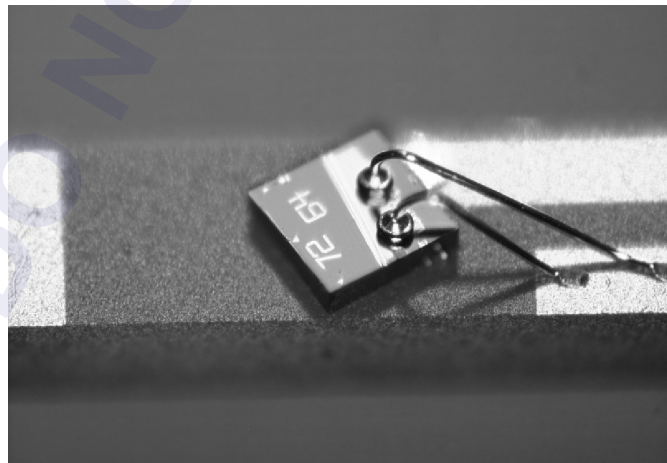


FIGURE 20 SOA chip—mounted on a heatsink.

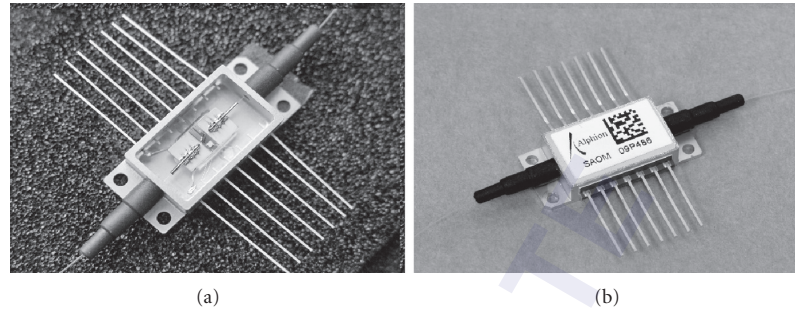


FIGURE 21 (a) Fiber-coupled chip and (b) SOA in industry-standard butterfly package.

which reduce the fraction of backreflected light that is coupled back into the waveguide, the effective facet reflectivity may be lower than 10^{-5} .¹³ Finally, the coated bars are diced into individual chips.

Most modern SOAs come packaged in a 14-pin industry standard butterfly package (see Fig 21). To increase fiber-chip coupling efficiency, either lensed fibers are used, or separate lenses are aligned between chip and fibers. In packaging schemes with multiple lenses per side, a collimated free-space section can be provided that allows for inclusion of polarization-independent isolators inside the package.

19.4 DEVICE CHARACTERIZATION

Chip Screening

The initial characterization of a SOA chip comprises the diode characteristic and the ASE light output. This is accomplished using a L-I-V measurement, which measures the light output versus current (L-I) and the I-V curve at the same time.

The measurement procedure is identical to the L-I-V measurement that is used to determine the threshold of a laser diode, the main differences being that no threshold is measured for properly antireflection coated SOA chips, and facet emission is monitored on both sides. Figure 22 shows an example measurement: A current sweep is applied to the chip, and the voltage across the chip is measured, while facet emission is monitored using broad-area photodiodes.

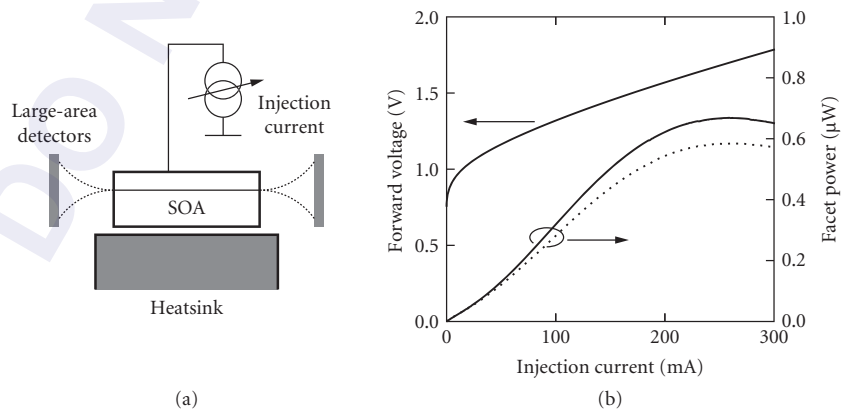


FIGURE 22 (a) Schematic of L-I-V measurement setup and (b) example measurement.

Gain Measurement

Measuring device gain requires coupling an input signal into the chip, and measuring output power divided by input power. This is easily done in a packaged device. Measuring an unpackaged chip requires aligning lenses or lensed fibers in front of the facets using nanopositioning stages.

In small signal gain measurements, the ASE noise the amplifier produces will typically overwhelm the output signal power. Therefore, a filter has to be used in the output. Depending on filter width, signal power, and accuracy required, it might even be necessary to measure noise level and subtract out the noise power transmitted through the filter, to avoid overestimating the gain.

A gain versus output power measurement like Fig. 10 is easily accomplished by stepping the input power. From Eq. (6) it appears that small signal gain and 3-dB saturation power can be easily extracted by a linear fit to

$$G = G_{ss} - 3p_o/p_{3\text{dB}} \quad (9)$$

where the measured gain values G are expressed in decibels and the measured output powers p_o are expressed in milliwatts. The fit parameter G_{ss} is the small-signal gain and $P_{3\text{dB}}$ is the 3-dB saturation output power as defined before. However, this expression assumes that the differential gain and the carrier lifetime remain constant upon saturation, which is not the case in practice. A second-order fit is necessary to accurately characterize the saturation behavior of the amplifier, in which case the 1-dB and 3-dB saturation powers can be used as figures of merit.¹³

Polarization-Dependent Gain Measurement

Signals of controlled polarization need to be applied to the SOA in order to measure its PDG. If a device has been pigtailed using polarization-maintaining fiber, all that is needed is to connect a horizontally polarized signal and then a vertically polarized signal. However, standard single-mode fiber does not preserve the polarization state, so more complicated measurement techniques are required in most cases.

Methods and systems available to measure polarization-dependent loss (PDL) in passive devices are usually applicable to PDG measurements, with one caveat: the device under test generates ASE noise, which has to be filtered out prior to the detector stage of the system. Because one is usually interested in the absolute value of the gain in addition to the PDG, the insertion loss of the filter needs to be well calibrated.

An often used method is polarization scrambling. A fast polarization scrambler scans through many different polarization states, evenly covering the Poincaré sphere. By measuring the output power of the SOA using a fast power meter with a min/max hold function, the highest and lowest gain over all polarizations can be found, and thus the PDG. Note that if an optical spectrum analyzer (OSA) is used as the power meter (in which case the detector *is* the filter), often a min/max function is not available, and having to coordinate the scans of the polarization scrambler and the OSA makes this method of PDG measurement inconvenient.

The *Mueller matrix method*²³ is an alternative that does not suffer from this problem. In this case, a polarization controller is needed that can synthesize any desired polarization state in the input signal, although a limited variant of the method works with a polarization controller that can access only the principal states on the Poincaré sphere: horizontal, vertical, 45°, -45°, right-hand circular, left-hand circular. Note that these states are defined as launched in the fiber directly following the polarization controller. Due to the nonpolarization-maintaining nature of the fiber, the states launched into the SOA chip are unknown, although their relative properties are maintained, that is, orthogonal states remain orthogonal, etc.

The essence of the method is to first measure the gain corresponding to an unpolarized signal by averaging the gain of two orthogonally polarized states, for example horizontal and vertical.

Next, gain is measured for three states that are orthogonal to each other *on the Poincare sphere*, for example, horizontal, 45° , and right-hand circular. The differences between these three gain values on the one hand, and the average gain measured in the first step on the other hand, contain the information needed to determine the polarization transformation the light has undergone in the fiber between the polarization controller and the chip. The gains in the maximum and minimum polarization states can now be calculated, or alternatively those two polarization states can be synthesized in the polarization controller and the corresponding gains measured directly.

In practice, therefore, the Mueller method allows for rapid measurement of both gain and PDG of a SOA by measuring output signal powers for four predetermined polarization states. Note that both this method and the polarization scrambling method only yield the highest gain and lowest gain over polarization; not information on which of these gains belongs to TE and TM polarization in the waveguide, respectively. In other words, the methods yield the absolute value of the PDG but not its sign. The Mueller method, but not the scrambling method, can yield the sign of the PDG by comparing the calculated polarization transformations for a device under test (DUT) relative to a device with a large PDG of which the sign is known, for example, a device with unstrained or compressively strained active layer. The fiber leading up to the facet needs to remain relatively undisturbed when replacing the DUT with the reference chip, in order to preserve the polarization transformation.

Gain Ripple

Gain ripple can be measured either directly, by sampling gain at closely spaced wavelength points, or by proxy by measuring ASE ripple. If the gain ripple period is known (e.g., by deriving it from the length of the chip), the number of gain measurements can be greatly reduced. Since the ripple is caused by a Fabry-Pérot cavity formed by the gain stripe and the two facets, the curve of reciprocal gain ($1/g$) versus wavelength is a sinusoid. Measuring three points on this sinusoid separated by one-third of the gain ripple period, the average and standard deviation are invariant for translation over wavelength. The inverted average now corresponds to the single pass gain, while the standard deviation can be worked into the ripple amplitude¹³ (see Fig. 23a).

Alternatively, the amplitude of the ripple on the ASE spectrum can be found by measuring a small wavelength span with sufficiently narrow resolution bandwidth (see Fig. 23b.)

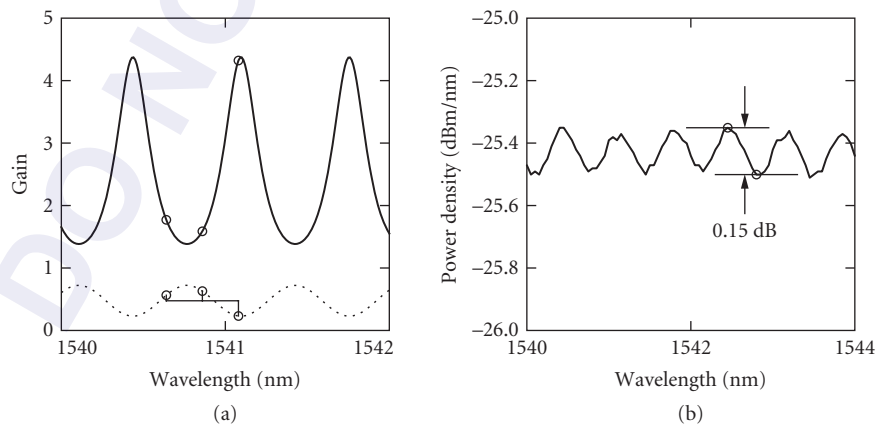


FIGURE 23 (a) Measurement of gain ripple and (b) measurement of ASE ripple.

Noise Figure

If the noise level has been measured during a gain measurement, all ingredients are present to calculate noise figure. If $G = P_o - P_i$ is converted to linear units and ρ_{ASE} is given in watts per hertz, noise figure in linear units is given by

$$\text{nf} = \frac{2(\rho_{\text{ASE}}/2)}{gh\nu} + \frac{1}{g} \quad (10)$$

The rightmost term was neglected in Eq. (5); it is only of importance for low-gain devices. The factor 1/2 expresses the assumption that half the total ASE power is copolarized with the signal. This is a good approximation for a polarization-insensitive SOA. For devices with considerable polarization dependence, the value of the PDG can be used to partition the total ASE power into two polarization components. Combined with the corresponding gain, the above expression will then yield two equal values for the NF. An actual measurement of the NF associated with each polarization state requires measurement of the output signal and ASE power through a polarization analyzer.

Complete Characterization

Figure 24 schematically shows a measurement system that allows full steady-state characterization of a SOA that includes all of the above measurements. The input source is a tunable laser to allow for measurement of gain versus wavelength. A variable attenuator and a polarization controller enable measurement of a saturation curve and of the PDG. Input power is referenced to a power meter, while output power and noise levels are measured on an optical spectrum analyzer.

An additional pair of power meters is provided on both sides for measuring the total ASE power. Comparing the ASE powers thus measured with direct power measurement of the DUT fibers with a reference power meter allows calibration of the loss of the connectors indicated in Fig. 24. This ensures a true fiber-to-fiber measurement. For chip measurements, the two extra power meters are helpful during optimization of the fiber position in front of the facets.

An example of a complete dataset measured for a SOA is shown in Fig. 25.

High-Performance SOA Properties

The gain of a SOA is only limited by the quality of the antireflection coatings, and by the rising ASE noise power which may cause the amplifier to autosaturate. Devices with fiber-to-fiber gain as high as 36 dB (at $\lambda = 1310$ nm, using a strained-layer MQW active layer) have been reported in the literature.¹³ High gain is not hard to obtain, by extending the device length and proportionally increasing the injection current. The difficult parameters to optimize are usually PDG, saturation output power, and noise figure. In this section, we will cite some state of the art results; unless stated otherwise, numbers quoted are fiber-coupled values.

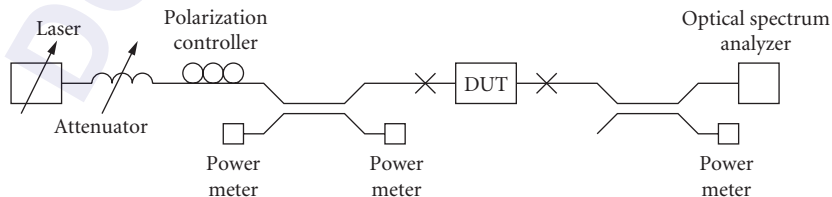


FIGURE 24 Measurement system to characterize gain versus wavelength, polarization, and input power.

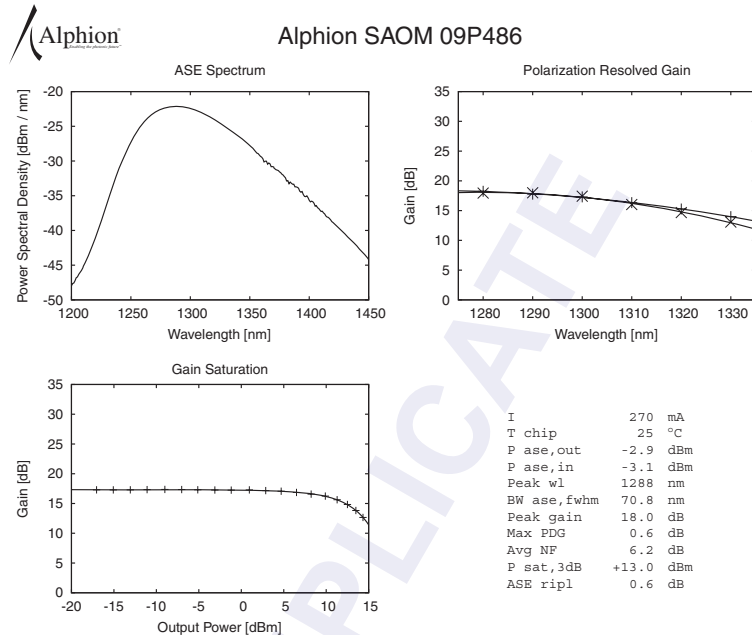


FIGURE 25 Example of a datasheet as delivered with a commercial SOA, showing the complete characterization results of the device.

A tensile strained bulk active device with a saturation output power of 17 dBm at an injection current of $I = 500$ mA was reported, having a gain of 19 dB at $\lambda = 1550$ nm.²⁴ Polarization dependence was 0.2 dB thanks to the tensile strain, and NF was 7 dB. The high saturation power was reached mainly thanks to the thin active layer (50 nm). At 600 mA injection current, a similar device, implemented using a MQW structure with tensile strained barriers and unstrained wells, reached a saturation power of 20 dBm.²⁵ Gain was 11 dB and PDG was 0.6 dB at 1550 nm, and NF was 6 dB.

A compressively strained MQW structure reached a chip output power of 24 dBm at 1590 nm.²⁶ This is a single-polarization device, having a peak chip gain of 18 dB. The 1.8 mm long structure was injected with a current of 1 A. The minimum chip NF was 3.6 dB, only 0.6 dB above the theoretical minimum, thanks to low cladding absorption. This chip was built into a polarization diversity module, in which a gain of 15 dB with a PDG of 0.5 dB was obtained. $P_{3\text{dB}}$ and NF were 22 dBm and 5.7 dB, respectively.²⁷

A very large saturation output power of 29 dBm has been obtained in a slab-coupled optical waveguide amplifier (SCOWA), thanks to an ultralow confinement factor ($\Gamma < 0.005$).²⁸ At an injection current of 5 A, the 10 mm long device exhibits a gain of 13.8 dB. The single-polarization device has a NF of 5.5 dB.²⁹

The above results have been obtained with MQW or bulk active layers. Quantum dot (QD) active layers show promise because of their ultrafast gain dynamics (a few picoseconds) and ultrawide gain bandwidth. A high-performance device has been demonstrated with a gain of 25 dB, a NF of 5 dB, and a $P_{3\text{dB}}$ of 19 dBm, all over a bandwidth of 1410 to 1500 nm.³⁰ The 6-mm long device was pumped with a current of 2.5 A. These are chip-referenced numbers, and the device amplifies one polarization only. Recently, a polarization-independent QD device was demonstrated, which uses dots arranged in columns, surrounded by tensile strained barriers, to obtain a PDG of 0.5 dB.³¹ At a length of 6 mm and an injection current of 1.2 A, a fiber-to-fiber gain of 4 dB and a saturation power of 16.5 dBm were obtained. The chip NF of 9.5 dB was relatively high due to the low gain.

19.5 APPLICATIONS

Applications for SOAs make use of the gain provided by the device, as well as the high-speed internal dynamics, which can be used for various optical signal-processing applications. The gain and index dynamics within the SOA must be managed to enhance the desired application. In a first set of applications, the SOA is used for linear amplification. Here, the operating regime of the device is managed to reduce nonlinear effects to acceptable levels, which benefits from devices with large saturation powers. For nonlinear applications, operating conditions are chosen to enhance the nonlinear gain and/or index nonlinearities. Both the SOA design and its configuration within a system are chosen to produce the desired nonlinear functionality. For nonlinear applications, smaller values of saturation power can be a benefit. Also, incorporation of the SOA into a subsystem with other optical elements, such as filters, interferometers, and the like, is used to create the nonlinear functional devices.

In the following sections, we characterize the applications of SOAs under three main headings: amplification of signals, switching and modulation, and nonlinear applications.

19.6 AMPLIFICATION OF SIGNALS

The function of amplification in a transmission system is performed in several places, as shown in Fig. 26, in which the amplifier can be a power amplifier, an in-line amplifier, or a preamplifier. As shown, the transmission system is divided into three parts: the transmitter, the transmission line, and the receiver. The power amplifier is used in the transmitter, the in-line amplifier is used in the transmission line, and the preamplifier is used in the receiver.

Single-Channel Systems

The power amplifier is located in the transmitter to boost the optical signal level before it enters the transmission system. (In Fig. 26, two possible locations are shown before or after the modulator.) In most cases, the power amplifier is placed before the signal modulator, so that the optical signal is either continuous wave (CW) or a pulse train. In this case, the properties of concern for the SOA are its output power and broadband ASE. The signal injected into the SOA is relatively large, so if the broadband ASE is filtered to allow only a narrow band around the signal channel, noise degradation by the SOA is minimal in the power amplifier. Because the polarization of the system is controlled at this location, polarization sensitivity of the amplifier is not of concern. Because the signal is either CW or a train of equally spaced pulses, the gain recovery time will not influence the temporal properties of the amplified signal. The SOA can be operated in saturation to maximize the output power.

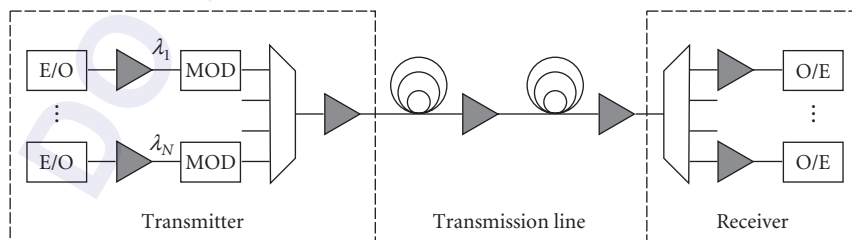


FIGURE 26 Placement of amplifiers in a WDM optical transmission system. The system is divided into transmitter, receiver, and transmission line blocks. SOAs here and in other figures are indicated by shaded triangles. E/O and O/E are electronic-to-optical converters (e.g., lasers) and optical-to-electronic converters (e.g., optical receivers), respectively.

However, when the SOA is used after the signal has been modulated and operated in saturation, it can distort the signal because the finite gain recovery time will introduce a time-dependent gain for each bit, which depends on the pattern of bits preceding it. For single-channel, intensity-modulated systems, this produces intersymbol interference (ISI). Because the gain recovery time is on the order of 100 ps, this will be a problem for multi-gigabit/s systems, where the bit period is comparable to the gain recovery time. An illustration of this is shown in Fig. 27, which shows an RZ bit sequence at 12.5 Gb/s, under conditions of minimal and 6-dB gain compression. The pattern-dependent gain produces significant eye closure, as shown. When the SOA is used as a power amplifier for analog-modulated signals, gain saturation can introduce second and higher-order distortions.¹⁹ To control ISI and analog distortions, such as composite second order (CSO) and composite triple beat (CTB), requires operating the SOA in a region of small gain compression, commensurate with the overall system application. Gain-clamped SOAs have been used for this application.¹⁹ Clearly, SOAs with large saturation powers offer advantages because they operate at higher output power for a given level of gain compression.

The pre-amp SOA is used in front of the receiver, as shown in Fig 26. At this point, the signal level is low so saturation of the amplifier is not an issue. The role of the SOA is to increase the optical power in the signal to make the received electrical power far exceed the thermal noise power in the receiver. Hence, the gain and noise figure are the important issues for the SOA when used as a pre-amp. To achieve maximum receiver sensitivity, the noise in the receiver should be dominated by signal-spontaneous beat noise, which requires a narrow-band optical filter after the SOA that is centered on the signal wavelength. At the receiver, the signal polarization has become randomized, so the polarization-dependence of the gain of the SOA pre-amp should be minimal (practically, less than about 1 dB).

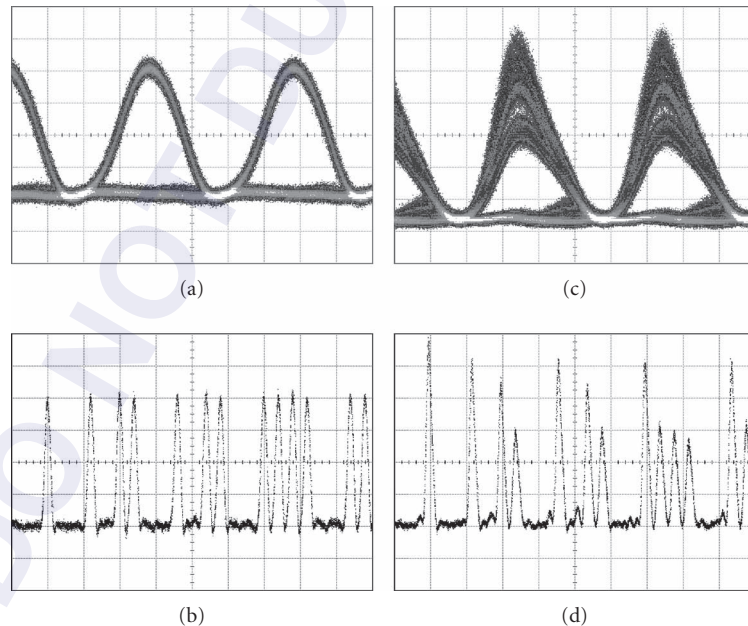


FIGURE 27 Effects of gain compression on bit sequences. Bottom traces show single-channel bit sequences at 12.5 Gb/s and top traces are corresponding eye diagrams. (a) and (b) correspond to the SOA operating under < 0.5 dB gain compression, while (c) and (d) correspond to the case of 6 dB gain compression. The effect of finite gain recovery on the pattern sequence and the corresponding eye diagrams is significant.

The function of the in-line SOA is to compensate for the loss in the fiber and other components in the transmission link. Issues of gain, ASE noise, gain recovery and its effect on ISI, and polarization-dependent gain arise for the in-line amplifiers. The output power level is determined by the operating characteristics for the transmission line, which require power levels compatible with management of fiber nonlinearities. Nevertheless, high-gain permits amplification to overcome the loss in long fiber links. The saturation power should be large, so that the amplifier can be operated at most lightly into the saturation regime. The ASE noise added to the signal by the SOA will degrade the optical signal-to-noise ratio and will ultimately limit the system. The noise figure of the SOA can be a limit on the number of SOAs that could be cascaded.

The ASE noise power in a single polarization from an optical amplifier is given by Eq. (11), which follows from Eq. (5), where

$$p_{\text{ASE}} = \frac{1}{2} n f (g - 1) h \nu B_0 \quad (11)$$

where p_{ASE} is the ASE power ($p_{\text{ASE}} = \rho_{\text{ASE}} B_0$), $n f$ is the noise figure, g is the gain, $h \nu$ is the photon energy, and B_0 is the optical bandwidth within which the ASE power is measured. Figure 28 shows a transmission line, which is composed of n fiber spans and n amplifiers. The loss of each fiber span, l , is exactly compensated by the gain g of each amplifier, so that the signal power is p_s at the output of each amplifier. Each amplifier adds ASE noise given by Eq. (11), so that the optical signal-to-noise ratio (OSNR), at the end of the line is given by

$$\text{OSNR} = \frac{P_s}{n p_{\text{ASE}}} = \frac{P_s}{n n f (g - 1) h \nu B_0} \quad (12)$$

Converting Eq. (12) to logarithmic units and referencing power levels to 1 mW (power is given in dBm), produces the very useful equation:

$$\text{OSNR}(\text{dB}) = 58 + P_s - G - \text{NF} - 10 \log n \quad (13)$$

In Eq. (13), P_s is the signal output power in dBm, G is the gain (and span loss L) in decibels, and NF is the noise figure in decibels. In arriving at Eq. (13), the gain is assumed much larger than unity and the optical bandwidth B_0 is taken to be 0.1 nm (at 1550-nm wavelength). Equations (12) and (13) also assume that the dominant optical noise is signal-spontaneous beat noise.¹⁴

Single-channel systems using SOAs to compensate fiber and other loss illustrate the general comments above. An early single-channel demonstration of a link amplified using SOAs was reported in 1996,³² for a 10-Gb/s RZ signal transmitted over 420 km at 1310 nm. In this experiment, there were 10 in-line SOAs (compensating ≈ 17 -dB loss from 40-km spans and other components), a power amplifier, and a preamplifier. The signal level was maintained below 10 dBm, well below the amplifier saturation power of 18 dBm. The signal degradation for this experiment is fully described by noise accumulation from the cascaded SOAs.

Successful demonstration of analog signals has been reported when the signal level is kept well below the saturation power. When used as a preamplifier for a CATV demonstration, which used a single wavelength carrying 23 QPSK subcarriers, a polarization-insensitive non-gain-clamped SOA enabled an 11-dB improvement for an 85-km link, (compared to a receiver without preamplifier).³³ Again, the main degradation arises from the SOA ASE noise.

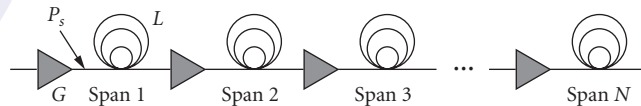


FIGURE 28 An amplified transmission line. Each fiber span has loss l (L in decibels), each SOA has gain g (G in decibels), and the power at the output of each SOA is p_s (P_s in decibels).

DWDM Systems

In DWDM systems, channels are typically spaced by 50 to 200 GHz. In the range around 1.55 μm , an eight-channel system spaced at 200 GHz can cover a spectral wavelength range of 11.2 nm. Over this wavelength range, while the gain and saturation power of the SOA can vary by a few decibels, as can be seen from Figs. 5a and 11, insight can be obtained by considering the gain and saturation powers to have averaged values.

In DWDM systems, power amplifiers can be used before combining channels, so their use is equivalent to the single-channel systems above. If the power amplifier is used after combining intensity-modulated signals, cross-gain modulation from one channel to the others will impress crosstalk onto the signals. This effect is mitigated by operating the amplifier under the conditions of small (at most ≈ 1 dB) gain compression. For equally spaced (in frequency) channels, four-wave mixing can potentially also cause crosstalk between channels. Except for systems that use short pulses at high bit rates (and large peak powers), four-wave mixing in the power amplifier is small compared to cross-gain modulation.³⁴

At the end of the system where the SOA can be used as a preamplifier, it is unlikely that the signals will saturate the amplifier, whether a single SOA is used to amplify many channels before they are separated or used for amplifying a single channel. The issues are the same as for a single-channel preamplifier: noise addition and polarization-dependent gain.

In-line amplifiers compensate for the span losses, which include losses due to dispersion compensation and other elements in nodes. As for the single-channel systems discussed above, the output power should be as large as possible to ensure overall large OSNR for the system. As the output power approaches the saturation power, the amplifier gain dynamics lead to time-dependent gain and crosstalk between channels caused by cross-gain compression. A variety of mitigation techniques have been proposed, including operation in a region of low saturation (≈ 1 dB gain compression),^{21,35,36} use of a reservoir channel that biases the SOA into a region of gain compression for which cross-gain modulation is reduced,^{36,37} use of a gain-clamped amplifier,³⁸ and the use of a pair of complementary signals to maintain a constant signal power into the amplifier.^{39,40}

An illustration of the effects of gain compression on systems using SOAs is shown in Fig. 29, which shows results for a system of eight channels at 40Gb/s, operating over two 80-km spans. As the

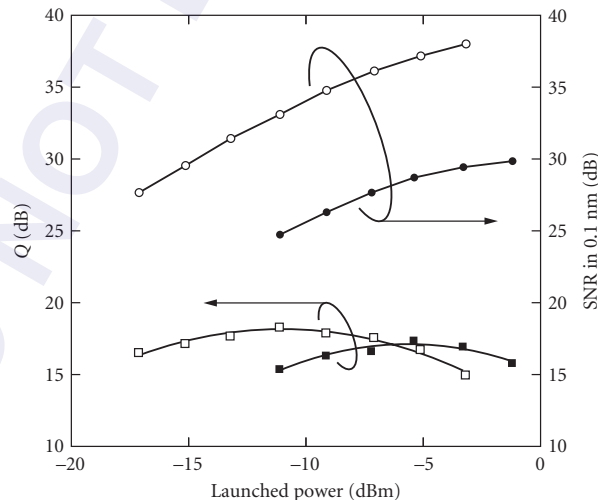


FIGURE 29 System results for 40-Gb/s transmission. The transmission line consists of two amplified 80-km fiber spans. Open symbols correspond to single-channel transmission, closed symbols correspond to transmission of eight channels, spaced at 200 GHz. Optimum Q is achieved in both cases for launch power in a single channel of ≈ -13 dBm.

optical power is increased, the OSNR increases, although the system performance, as measured by the Q-factor⁴¹ reaches a peak and declines at higher power levels, due to cross-gain modulation effects. The peak in input powers for one-channel and eight-channel systems is separated by ≈ 9 dB, and in both cases corresponds to the same total input power to the SOA of -5 dBm, which corresponds to about 1-dB gain compression for the SOA that is used in the experiment.

The system capacity is limited by the OSNR at about 1 dB of gain compression. In a rough way, for similar systems before the onset of nonlinearities, Q is almost linearly related to the OSNR. Also, for different bit rates, the optical bandwidth for the optimized receiver depends linearly on the bit rate. Assuming these linear relationships, and noting that the relevant signal output power in Eq. (13) is the sum of powers for all channels (and assuming a single value for gain and saturation power), Eq. (13) can be manipulated to give

$$10 \log c = 58 - k - P_{1\text{ dB}} - L - NF - 10 \log n \quad (14)$$

In Eq. (14) c is the overall system capacity in number of channels \times bit-rate, $P_{1\text{ dB}}$ is the 1-dB gain compression power for the SOA, L is the span loss in dB (spans \times loss/span), and k is a constant relating optimal Q to OSNR. The value of $P_{1\text{ dB}}$ for the output power has been used because this is the operating point for which nonlinearities caused by gain compression are minimal. Figure 30 shows a set of experiments for a variety of bit rates and links using similar SOAs. For these experiments, each systems capacity is determined for $Q \approx 16$ dB, corresponding to a bit-error rate $\approx 1 \cdot 10^{-9}$. The horizontal axis is the total link loss, which is the number of spans times the loss per span. In these experiments, the span losses were in the range of 15 to 17 dB and channel separations were either 100 or 200 GHz. The results

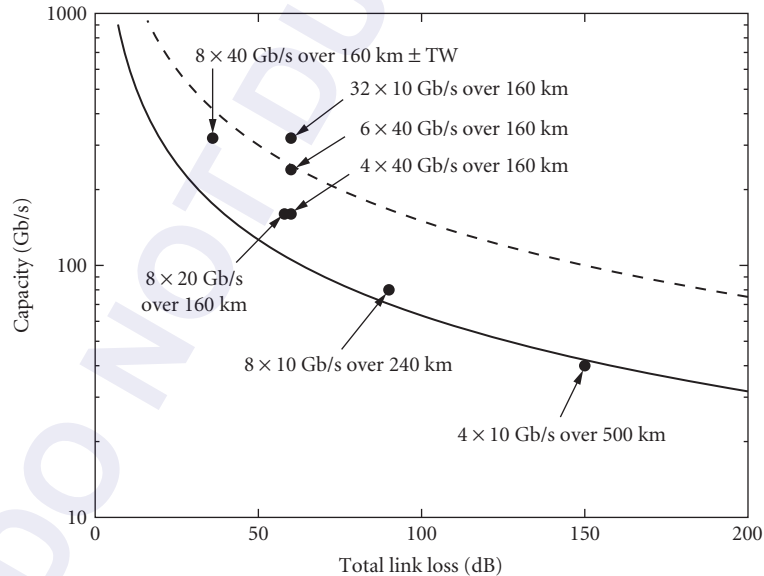


FIGURE 30 Results of a series of DWDM system experiments over multiple fiber spans, using a set of nominally similar SOAs. Overall system capacity, which is the number of channels times the bit-rate per channel, is plotted against the total link loss, which includes loss of dispersion compensating fiber for each span. The systems use NRZ intensity modulation and capacity is determined where all signals achieve a bit-error rate of 10^{-9} ($Q \approx 15.8$ dB). The two curves show expectations for all spans with 20-dB loss (solid line) and 15-dB loss (dashed line). TW: TrueWave fiber.

show that the analysis in terms of OSNR captures the main systems limitations and describes the overall system capacity.

When SOA systems are used with forward error correction (FEC) coding, Q values as low as 8 dB can be corrected to essentially error-free operation. This permits operation at lower OSNR (about 8 dB) and from Eq. (14), this will enable the overall system loss (and therefore length or capacity) to be increased by about 8 dB from the results shown in Fig. 30.

For phase-modulated signals, the lack of temporal intensity-dependence of the input signals will reduce the effects of time-dependent gain and cross-gain compression. In this case, the dominant nonlinearity, which will limit transmission, will be four-wave mixing. Thus, large system capacity and distance have been achieved using differential phase shift keying (DPSK) modulation, but larger channel separation is required to reduce four-wave mixing effects. Using DPSK modulation and SOAs operating under 3-dB gain compression, it was possible to operate an eight-channel, 10.7-Gb/s system over 1050 km, with a FEC margin of 6 dB.⁴² Because of four-wave mixing, the channels were spaced by 400 GHz. Four-wave mixing can also lead to system penalties comparable to those from cross-gain modulation for high bit-rate systems using short RZ pulses.³⁸

CWDM Systems

Coarse WDM (CWDM) systems are envisioned to be inexpensive and environmentally tolerant for access and metro applications, with a reach of up to about 100 km. To this end, the channel spacing is set to 20 nm in a wavelength range that spans the window from about 1250 to 1650 nm—a window that could include 20 channels. The coarse spacing allows for very reduced temperature and wavelength stability requirements on the components, and requires filters with large, almost 20-nm bandwidths. Amplification can extend the reach of such systems.

The broad gain bandwidth of SOAs is advantageous for application to CWDM systems, although a single SOA is incapable of covering the entire 1200- to 1600-nm window. Nevertheless, a broadband LOA has been used to amplify eight CWDM channels covering the range from 1470 to 1610 nm.⁴³ To cover the full range, a set of SOAs designed to amplify different gain bands can be used.

Because the channel wavelength is specified only to within a band of 20 nm, filters that separate the different channels are much wider than the filters used in DWDM systems. Therefore, in addition to signal-spontaneous beat noise, the component of optical noise due to beating of the ASE from the SOA with itself, within the signal bandwidth, (spontaneous-spontaneous beat noise) is also significant. Equations (12) and (13) are no longer valid. For filter bandwidths of 13 nm, the addition of the spontaneous-spontaneous beat noise will reduce the sensitivity of a preamplified receiver by about 3 dB, assuming about 25-dB amplifier gain and a bit rate ≈ 5 Gb/s.¹⁴ Other issues for SOAs used in CWDM systems are the wavelength variation of the gain and saturation powers. Because the gain will vary by more than 3 dB over the SOA bandwidth used, a cascade of amplifiers would quickly lead to very unbalanced channels. With only one amplifier in the system, and only ≈ 100 km of fiber, this does not lead to much difficulty. For a cascade of amplifiers and fiber spans, however, weaker channels would rapidly degrade in OSNR and fiber nonlinearities would build up rapidly for the stronger channels.

Cross-gain modulation becomes an impairment when the amplifier is operated in saturation. All input wavelengths can contribute to saturation and reduction of the carrier density, but for equal input power in each channel, the channels nearest the gain peak will cause the largest effect. Because the differential gain is largest on the short-wavelength side of the gain peak, the carrier density changes will cause the largest gain compression for these channels. Therefore, the short wavelength channels will be most affected by cross-gain modulation. The nonlinearity caused by four-wave mixing will not be significant in CWDM systems, because the wavelengths are so widely separated and the four-wave mixing interaction falls off rapidly with wavelength (see Fig. 37).

For broader bandwidth coverage, a hybrid amplifier composed of an SOA and a fiber Raman amplifier, has been applied to CWDM systems.⁴⁴ Here, the Raman pump is chosen so that the Raman gain peak lies to longer wavelength than the SOA peak. The composite amplifier can have a 3-dB gain bandwidth of over 120 nm. Wireless over fiber transmission has also been demonstrated in a CWDM system, employing a bidirectional SOA.⁴⁵

19.7 SWITCHING AND MODULATION

An unbiased SOA is strongly attenuating. The difference in output of a signal emerging from an SOA with gain compared to one attenuated by an unpumped SOA can be as high as 60 dB. This forms the basis of an optical switch, with very high extinction ratio.⁴⁶ Figure 31 shows a 2×2 switch fabric, where the switching elements are SOAs. For this simple switch matrix, the bar state is enabled by supplying bias, and therefore gain, to SOAs 1 and 4, while SOAs 2 and 3 are unbiased. Thus, the cross state is highly attenuated. For the cross state, the opposite set of biases are applied—SOAs 2 and 3 are biased and SOAs 1 and 4 are not biased.

The SOA-based switch is capable of moderately fast reconfiguration. The temporal response of an SOA between full on and full off states is several times the gain recovery time (to achieve high extinction), but can generally be ≈ 1 ns. To switch the state of the SOA by adjusting its drive current, the response must also include electrical parasitics, which can have bandwidths of up to a few gigahertz (or less), depending on device design. Thus, gigahertz rearrangements for switching fabrics based on SOAs are possible, and this is one of their significant benefits. For a switching application, the signal transmitting the fabric should not be distorted and should have minimal noise degradation. Thus, signals passing through SOAs in a switch fabric are in these regards similar to signals passing through cascades of SOAs in the transmission systems described above.

The fabric above is intensive in the number of SOA gates required per port. Other switch geometries, involving combinations with AWGs, can reduce this requirement.^{47–50} Because of the high-speed response of the SOA switch fabrics, they are often used in space-switching stages of larger all-optical cross-connects, which may have additional features of wavelength interchange (using wavelength converters, such as discussed below). Impressive all-optical nodes, including all-optical packet switching nodes, have been constructed using these switches.^{51–53}

Because the gain of the SOA varies with the bias current, the SOA could be used as an optical modulator. The response speed of an SOA as a modulator is significantly slower than the response of a directly modulated laser, because there is no optical resonator and hence the response has no relaxation oscillation peak. The modulation bandwidth can be ≈ 2 GHz, and direct modulation has been achieved for a 2 Gb/s bit rate.⁵⁴ Of course, the optical signal would be highly chirped. The amplifier does provide gain and its simplicity could be advantageous for some applications. The concept of a reflective SOA (RSOA) as a modulator was introduced for simple access systems, where the RSOA at the user node could amplify and modulate a wavelength that originates from the central office.⁵⁵ Various improvements and applications use RSOAs in access networks with bit rates beyond 1 Gb/s.⁵⁶

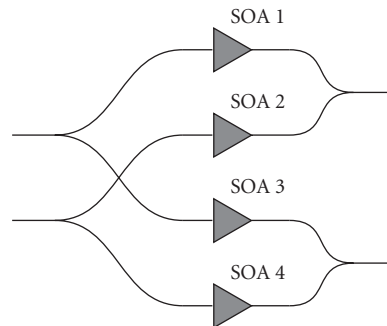


FIGURE 31 2×2 switch fabric based on SOA gates.

19.8 NONLINEAR APPLICATIONS

The gain and index of refraction within the SOA depend on several different physical processes with different time scales. The magnitudes and particular time responses will depend on material, structures (multiple quantum well vs. quantum dot vs. bulk) and operating conditions, such as wavelength, power, bit-rate, modulation format, and the like. There is a rich variety of effects and their applications. The general applications of these nonlinearities lead to optically controlled optical gates. The three most significant nonlinearities in SOAs are cross-gain modulation, cross-phase modulation and four-wave mixing. XGM and XPM depend on the optical intensity, while FWM depends on the optical field. The following sections describe the application of these nonlinearities.

Cross-Gain Modulation

The principle of XGM is shown in Fig. 32. In this process, a strong, saturating input (pump) signal at λ_1 and a weaker input at λ_2 (probe) are simultaneously injected into the SOA. The injection of the strong optical signal at λ_1 reduces the carrier density, which compresses the gain at all wavelengths. When the pump signal is turned off, the gain recovers. Thus, an inverted copy of the intensity modulation on the pump signal is copied onto the probe signal, as in Fig. 32. The device has functioned as a wavelength converter.^{57,58} Note that if probes at several wavelengths are input to the SOA, each probe experiences gain compression and the pump signal can be copied simultaneously onto each.⁵⁹ The schematic shows that the pump and probe signals copropagate and are separated by an external filter. It is possible to avoid the filter by counter-propagating the two signals. This reduces the highest speed capabilities of XGM, however, due to traveling-wave effects.

As the free carrier density is reduced, the gain compression depends on wavelength, as can be seen by comparing ASE curves for different pumping currents (related to carrier density) in Fig. 5b.

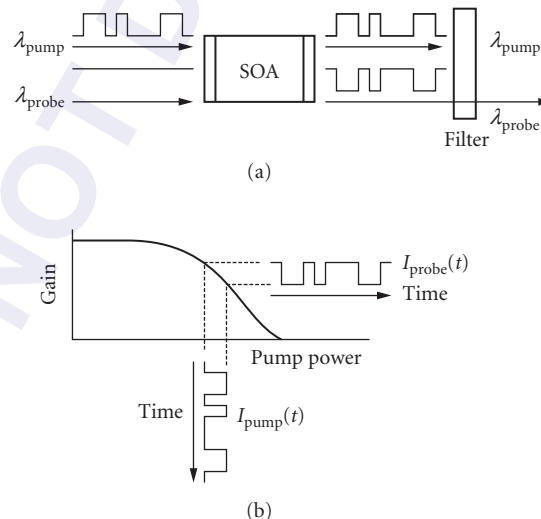


FIGURE 32 Principle of cross-gain modulation in a SOA. (a) In basic operation, a strong, modulated pump beam and a CW probe beam are input to the SOA. (b) The pump power is intense enough to compress the gain of the SOA, which creates an inverse of its intensity modulation on the probe. In this collinear arrangement, the outputs are separated by a filter (a).

Thus, the gain is compressed more for shorter wavelengths than for longer wavelengths, and so the extinction ratio (which is the gain compression) for the converted signal will be larger for conversion from longer to shorter wavelength than vice versa. At different wavelengths, the power required for gain compression will scale with the saturation power, which increases with longer wavelength (Fig. 11). To achieve an extinction ratio of 10 dB or larger, the pump output power must be several times the saturation power, as can be seen from Fig. 10.

The polarization sensitivity for wavelength conversion by XGM depends on the polarization dependence of the SOA gain. For SOAs with minimal polarization dependence, there is minimal polarization sensitivity for XGM. Dynamical processes that are faster than carrier recombination can lead to anisotropies even in nominally polarization-insensitive SOAs,⁶⁰ but these processes are fast and are not significant in cross-gain modulation wavelength conversion for multi-gigabit/s signals.

The temporal response of a wavelength converter based on XGM depends on the gain recovery time, which is dominated by the carrier recovery time. For SOAs, which operate with carrier densities in the range of 5×10^{18} to 10×10^{18} cm⁻³, the gain recovery is dominated by Auger recombination and is typically in the range of 100 ps (which is consistent with the data of Fig. 14). This would seem to support an operating bit rate of up to 10 Gb/s. To enable operation at higher bit-rates requires reducing the gain recovery time, which has been accomplished in several ways. Adding another significant recombination process speeds up gain recovery and can be accomplished by using stimulated emission. If the input signal powers are high enough, stimulated emission can become strong enough to exceed the Auger recombination rate and reduce the gain recovery time to enable operation at 40 Gb/s.⁶¹ Use of a long amplifier also helps, as the overall gain increases the stimulated emission rate and relatively enhances the gain for high-frequency signal components.^{62,63} In practice, stimulated emission can be increased by using a relatively high-power probe beam⁶⁴ or using a third beam, the control beam, to set the rate of stimulated emission.⁶⁵ Different materials can also affect the gain recovery time by carrier recombination. Recently, quantum dot material has shown promise for high-speed operation, because of its short-gain recovery time.⁶⁶ SOAs based on p-type modulation-doped MQW active regions have shown a fast intrinsic carrier recombination time, which can be exploited for high-speed nonlinear devices.⁶⁷ Another method for higher-speed operation of XGM involves filtering, which will be described next.

The carrier density affects not only the gain, but also the refractive index of the inverted semiconductor material. Thus, the temporal variation of the gain will impose a temporal variation on the refractive index within the SOA, which will lead to chirp on the wavelength-converted signal.⁶⁸ For large gain compression, the refractive index change can be several times π radians. The phase change corresponding to a gain change can be given approximately by¹⁷

$$\frac{\Delta\phi}{\Delta G} = -\alpha/8.68 \quad (15)$$

where ΔG is the gain change in decibels and α is the linewidth enhancement factor for semiconductor lasers and amplifiers.⁶⁹ For a value of $\alpha \approx 8$, a phase-change of π requires about 3.4 dB of gain compression.

Because the wavelength-converted output at λ_2 is chirped, filtering can reduce both the excess spectral content of the output and its temporal response. Thus, narrow-band filtering at the output of the SOA can be used to effectively speed up the response of the wavelength converter.⁷⁰ This approach has been used to achieve 100-Gb/s operation⁷¹ and up to 320-Gb/s operation of a wavelength converter based on XGM.⁷²

Cross-Phase Modulation

The refractive index changes with carrier density (and gain) led to chirp for signals created by XGM. The phase change caused by the refractive index change is the basis for another set of optical-optical gates and wavelength converters based on cross-phase modulation. To use XPM with intensity-modulated signals requires converting phase to amplitude and thus requires an interferometric structure.

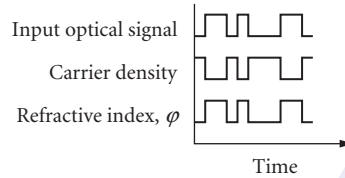


FIGURE 33 Principle underlying cross-phase modulation. An intense input optical signal (top) passes through the SOA. This signal reduces the carrier density by stimulated emission (middle). The refractive index of the active region depends on the free carrier density, increasing when the carrier density decreases (bottom). The phase change experienced by the probe propagating through the SOA is used for cross-phase modulation switching. Note that the temporal response is assumed instantaneous for this figure.

The operating principle is shown in Fig. 33. An input intensity-modulated signal will reduce the carrier density within the SOA by stimulated emission and reduces the gain as well, resulting in XGM. The refractive index within the active region depends on the free carrier density and increases as the free carrier density decreases. This produces a phase change $\Delta\phi = \Delta nL$, for a device of length L (assuming uniform Δn), as shown in the bottom of Fig. 33. This phase change can be read out by a second optical signal, if the SOA is placed in some kind of interferometric structure through which the second signal propagates.

The Mach-Zehnder structure, with an SOA in each arm, is used for many applications based on XPM.⁷³ Fig. 34a shows a Mach-Zehnder interferometer with an SOA in each arm and three input waveguides. A CW probe at λ_2 is input from the left. The relative phases in each arm can be adjusted by controlling the carrier density in each SOA by its bias current. If the phases of the two arms are

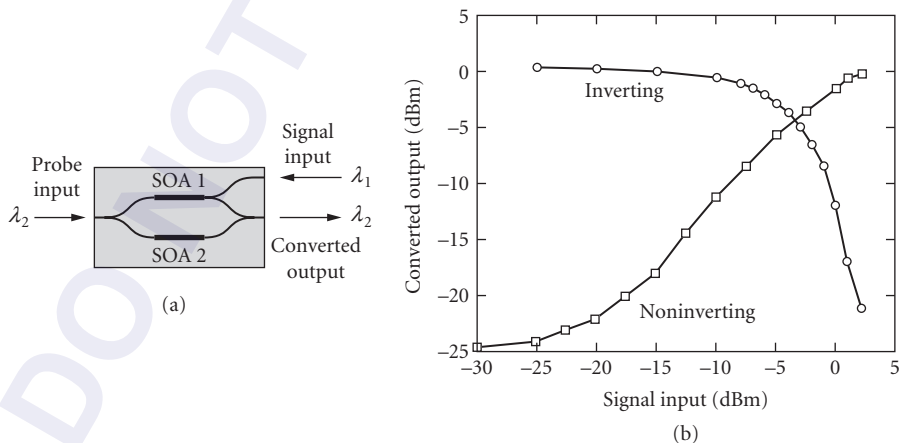


FIGURE 34 Mach-Zehnder interferometer used for cross-phase modulation wavelength conversion. (a) The arrangement shows the input (pump) signal and the probe signal counter-propagating. The input signal modulates the phase in SOA 1, switching the output state of the interferometer for the probe signal. (b) The output signal for inverting and noninverting operation, under CW conditions. Because of the nonlinear shape of the switch curves, it is possible to increase the extinction ratio from input to converted output.

equal, the two arms constructively interfere at the output port, producing a maximum for the probe at the output port. Note that each SOA also has gain, so that the output probe has greater intensity than the input. If the two arms have a phase difference of π , the probe experiences destructive interference at the output port and a minimum intensity. The values of the maxima and minima for points of constructive and destructive interference also depend on the gain in each arm. The greatest static extinction ratio occurs when the gains of the arms are equal. When an intensity-modulated signal at λ_1 is input into one arm of the interferometer (Fig. 34a), it changes the phase in that arm. If the phase is changed by π , the output of the probe signal can be switched from on to off, for the case of initial constructive interference, and from off to on, for the case of initial destructive interference. These cases correspond to inverting and noninverting operation, respectively. Figure 34b shows switching curves for inverting and noninverting operation. The extinction ratio for such static switching can be in excess of 25 dB (sometimes as large as 40 dB).

The free carrier density dependence of the refractive index has a very weak dependence on wavelength. Therefore, for a given gain compression, the extinction ratio and wavelength conversion penalties are almost independent of the direction and magnitude of wavelength shift (as long as it remains within the gain bandwidth of the SOA). In operation, the dynamic extinction ratio is not as large as the static extinction ratio, because of carrier dynamics which, at multi-gigabit/s bit rates, do not permit the device to reach a steady state. Nevertheless, dynamic extinction ratios of 15 dB are generally achieved.

The carrier density is varying during the pump intensity transitions, so chirp is imposed on the converted signal transitions. For the case of noninverting operation, the sign of the chirp leads to pulse compression in standard single-mode fibers in the 1550-nm regime.⁶² For inverting operation, the sign of the chirp leads to immediate pulse broadening in standard SMF.

The polarization dependence of XPM in the Mach-Zehnder structure depends on the polarization-dependence of its components. For polarization-independent operation, the SOAs and the waveguides must be made polarization-insensitive. It is possible to optimize the Mach-Zehnder structures in many other ways, to permit counter-propagation of pump and probe signals, to separate pump and probe signals without need of an external filter, and to maximize the overlaps of optical fields inside the SOA active regions.^{74,75} Other interferometer structures, such as Michelson and Sagnac, have also been used for nonlinear applications of XPM.

The gain recovery time in the SOA also affects the phase recovery time, which limits the response speed for XPM. To increase the operating speed of the optical-optical gate based on XPM in a Mach-Zehnder interferometer, it is possible to operate the device in a differential mode,⁷⁶ as illustrated in Fig. 35. By supplying pump inputs into the SOA in each arm, but with a delay τ between them, it is possible to carve a gating window of duration τ . Figure 35b shows that the phase dynamics of the SOA in the upper arm can be balanced by the phase dynamics in the lower arm. The output phase is determined by the difference in phase between the two arms, which creates the gating window. Of course, the relative amplitudes of the pumps in each arm must also be properly balanced so that full cancellation of the time-varying phase in each arm can be achieved, to create a gate with high extinction. Figure 35c shows the gating windows that can be carved by varying the relative delays between pump inputs to the upper and lower SOAs. Using this device as an optical demultiplexer, a 10.5-Gb/s signal has been demultiplexed from a 168-Gb/s data stream.⁷⁶

XPM in interferometric structures that require only a single SOA have also been demonstrated. The ultrafast nonlinear interferometer (UNI) device⁷⁷ operates by creating an interferometer based on orthogonal polarization states propagating through the SOA. By proper timing of the pump pulse with respect to temporally delayed probe pulses in the orthogonal polarizations, the pump can affect the index of only one polarization state. Thus, when the two orthogonal polarizations are recombined, switching is effected. In a different arrangement, a single SOA in which a pump and probe copropagate is combined with a 1-bit delay interferometer to convert XPM to intensity modulation, effectively copying the signal from the pump to the probe.⁷⁸ This can be fast and the SOA and delay interferometer can be integrated monolithically, but it does operate for a fixed bit rate, set by the delay interferometer.

In all of the above discussions, intensity modulated signals were required for both XGM and XPM. Wavelength conversion of phase-modulated signals has also been demonstrated, in a sophisticated two-stage device based on XPM, using a delay interferometer and a Mach-Zehnder wavelength converter.⁷⁹

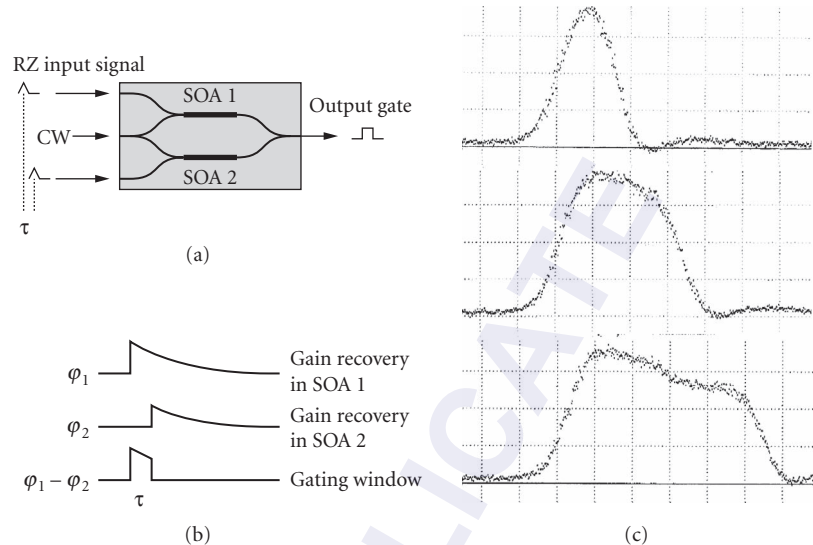


FIGURE 35 Differential operation of the Mach-Zehnder wavelength converter. (a) The input pulses, separated by a time delay τ , are sent into SOA 1 and SOA 2. The first pulse, to SOA 1, creates a phase change in the top arm, which switches the interferometer output for the CW probe input. The second pulse, to SOA 2, creates a phase change in the bottom arm, which cancels the remaining phase change in the top arm, closing the interferometer output. (b) Schematic of the gain recovery and the creation of a gating window. (c) Examples of tuning the switching window by varying the separation of pulses. The horizontal scale is 20 ps/division.

Four-Wave Mixing

Four-wave mixing is the fastest nonlinearity in SOAs, but it is very sensitive to polarization and to the wavelength separation between inputs. This is a $\chi^{(3)}$ nonlinearity, which depends on the optical field rather than optical intensity.⁸⁰ Figure 36 shows a schematic for four-wave mixing in a SOA with two input wavelengths. As explained previously, in four-wave mixing the two input fields interfere in the SOA to produce phase and index gratings. Each of the input fields then scatters off these gratings, producing longer and shorter wavelength sidebands, as shown in Fig. 16. If the two

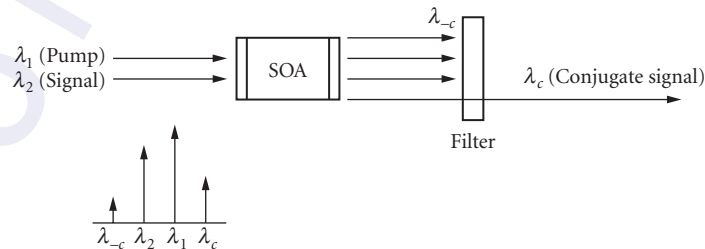


FIGURE 36 Four-wave mixing in a SOA. A strong pump and a signal are input to the SOA. Nonlinear mixing creates sidebands separated from the input pair by the difference in frequency between inputs (bottom). The desired signal, the conjugate signal, is separated from the inputs and the other sideband by a filter (top).

input wavelengths are λ_p and λ_s , for pump and signal, respectively, conservation of energy produces a converted signal λ_c at

$$\frac{1}{\lambda_c} = \frac{2}{\lambda_p} - \frac{1}{\lambda_s} \quad (16)$$

If the pump signal is CW, the converted signal will carry both the amplitude and the phase-conjugate information of the original signal. Thus, FWM is an optical-optical gate that, when used for wavelength conversion, straightforwardly works for both amplitude and phase-modulated signals. The converted signal described by Eq. (16) is created using one photon from the input signal and two from the input pump. The other sideband at λ_{-c} , which is created from two signal photons and one pump photon ($1/\lambda_{-c} = 2/\lambda_s - 1/\lambda_p$), has a more complicated amplitude and phase relation to the original signal.

An example of the dependence of four-wave mixing efficiency as a function of wavelength separation between the signal and pump beams is shown in Fig. 37. The conversion efficiency is weak, depends strongly on wavelength separation, and is asymmetric to longer or shorter wavelength. The shape and asymmetry arises because the nonlinear gratings formed within the SOA arise from several dynamical processes.⁸¹ Carrier density modulation, which has a characteristic lifetime of ≈ 100 ps, is responsible for the highest efficiency mixing process, but only to frequency shifts of ≈ 10 GHz (≈ 0.1 nm at 1550 nm). Carrier heating and cooling processes, which have a time constant ≈ 1 ps, are responsible for frequency shifts to ≈ 1 THz (≈ 8 nm), while the interband Kerr effect, with a characteristic time constant of ≈ 1 fs, is responsible for frequency shifts beyond 1 THz. Each dynamical process has its own amplitude-phase coupling constant and the overall four-wave mixing process is the coherent addition of nonlinear interactions from each process. Because of the different amplitude-phase coupling constants, the conversion efficiency is asymmetric with wavelength. Four-wave mixing requires phase-matching as well as energy conservation [Eq. (16)]. For typical SOA lengths below ≈ 1 mm

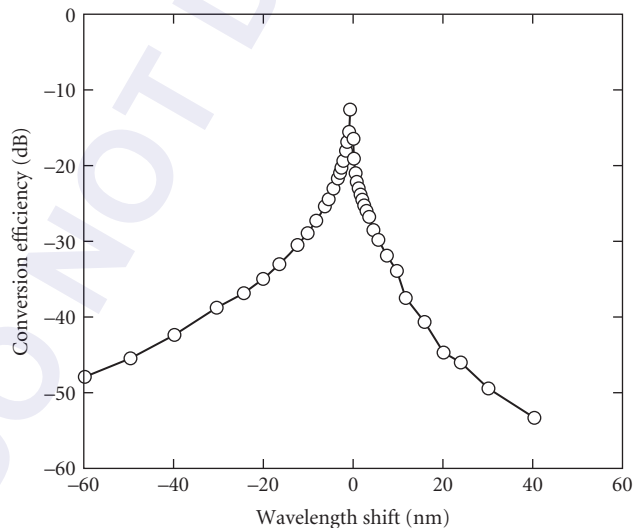


FIGURE 37 Four-wave mixing efficiency as a function of wavelength separation between pump and signal inputs, for a SOA operating near 1550 nm. For this experiment, the two inputs have equal intensity. Note the asymmetry and that conversion efficiency is larger to the shorter wavelength (higher frequency) side.

and for frequency shifts of less than several terahertz, phase-matching requirements are satisfied. For wavelength shifts greater than ≈ 100 GHz, the dynamical processes responsible are on the 1-ps or shorter time scale, so the response speed of FWM is very fast.

The FWM process is capable of shifting multiple signal wavelengths simultaneously, with a single pump wavelength. Each signal produces a converted output as described by Eq. (16). Thus, a band of wavelengths may be shifted simultaneously (but the wavelengths are not shifted independently).

The conversion process is weak: typically for input signals of a few dBm, conversion efficiency varies from -10 to -40 dB, as the wavelength shift increases. Because the SOA also produces ASE noise, the signal-to-noise ratio of the converted output can be affected. For conditions of small saturation, in addition to the wavelength shift, the converted power depends on $g^3 p_p^2 p_s$, where g , p_p , and p_s are the gain, pump power, and signal power, respectively, while the ASE noise depends on g .⁸² For saturated conditions, the converted power can increase more rapidly with pump power.⁸³ The efficiency of FWM can be enhanced by using long SOAs, up to 2 mm in length.^{84,85} This effect can be understood by noting that in a long SOA, the later portion is saturated and provides little gain. However, the nonlinearity is active in this region and the FWM fields grow coherently and quadratic in the SOA length, while the ASE adds incoherently and grows linearly in SOA length.⁸⁶ Therefore, the efficiency and the signal-to-noise ratio for FWM are enhanced in a long SOA.

FWM is very polarization sensitive—the results in Fig. 37 are obtained with all polarizations parallel. More complicated arrangements, involving two nondegenerate pump wavelengths, have been used to make the FWM process polarization-insensitive and to flatten the wavelength dependence of the conversion efficiency.^{87, 88}

FWM in SOAs has been used for a variety of optical-optical gate functions, including wavelength conversion,⁸⁹ optical sampling,⁹⁰ optical logic, and the like. Because the converted signal preserves the phase information of the input signal, but is its phase-conjugate, FWM in SOAs has also been applied to dispersion compensation in fiber transmission systems.^{91,92}

Further Comments for Nonlinear Applications

The optimum length of the SOA depends on the application.⁸⁶ For linear applications, high-gain, high saturation power, and low noise figure are desirable properties. For short amplifiers with high gain, the carrier density (and therefore pumping current density) is high. This will lead to a large inversion factor and therefore low noise figure. Also, the carrier lifetime will be short because of the large carrier density, and the differential gain coefficient may be small because of band-filling, resulting in a good saturation power. However, if the SOA is too short, the required pumping current density to produce the desired gain may be unsustainable without damaging the SOA. Alternatively, large saturation power for a given gain can be achieved in a long amplifier with a lower gain coefficient. This can be achieved using a thin active region, which increases the saturation power by decreasing the mode confinement factor, Γ , increasing the carrier density, which reduces the carrier lifetime, and again decreasing the differential gain coefficient because of band-filling. In long SOAs, however, the ends of the SOA can be saturated by the amplifier's own ASE, which results in decreased inversion and higher noise figure. Additionally, the effects of internal loss in the SOA are more significant when the gain coefficient is low (because of Γ), which also effects the carrier inversion and noise figure. SOAs of various lengths have therefore been used for linear applications.

For nonlinear applications, it is desirable to have amplifiers with lower saturation input powers (the input power that would produce 3-dB gain compression of the amplifier gain), to enable XGM and XPM with smaller pump power requirements. This is achieved in high-gain amplifiers, with relatively longer lengths. Because the probe signal is input with high power, the ASE noise is of less significance for XGM and XPM than for linear systems. Also, for longer SOAs the high rate of stimulated emission can shorten the carrier lifetime, leading to shorter gain-recovery times. To fully understand the temporal response of an SOA for XGM or XPM, however, requires modeling the SOA as a distributed amplifier, which shows that the modulation frequency response varies along the length of the amplifier and depends on the magnitude of the internal loss.^{62,63,93} The final, output section of the SOA, where the gain is compressed significantly, has a gain response with larger bandwidth than

the preceding sections, which leads to high-speed operation. In general, the temporal response of the SOA for XGM or XPM is enhanced for a long amplifier and enhanced by large Γ , large differential gain coefficient, high current (for large carrier density), and large input power. For FWM applications, a long amplifier is also beneficial. The front portion of the SOA can amplify the input signals, while the latter, saturated section can build up the nonlinear interactions, as discussed in “Four-Wave Mixing” in Sec. 19.8.

The nonlinear applications for SOAs were illustrated above mostly for the function of all-optical wavelength conversion. But in general, the application is that of an optical-optical gate, which can be used for various all-optical logic functions, such as optical demultiplexers (as illustrated for a Mach-Zehnder based device, above,⁷⁶ XOR and other logic gates, optical header processors, optical sampling, and all-optical regeneration. The optical transfer function for an interferometer, such as the Mach-Zehnder wavelength converter in Fig. 34, is decidedly nonlinear, so it can be used as a 2R (reamplifying and reshaping) regenerator. When it is combined with optical retiming, it becomes a key element in a 3R optical regenerator (reamplify, reshape, retime). There is a significant body of work on all-optical regeneration, which often uses nonlinear SOAs as a fundamental gating and shaping mechanism.⁹⁴

19.9 FINAL REMARKS

This chapter has described the properties, technology, and applications for SOAs. It is not meant to be a complete literature review on the subject, but rather to describe the principles underlying the devices and their basic applications. The technology of semiconductor optical amplifiers has reached a level of maturity where devices with well-defined properties are supplied and characterized. They are capable of supplying linear gain for a number of applications and systems of low to moderate complexity. As stand-alone gain blocks, the fact that, unlike fiber amplifiers, their gain spectrum can be tuned by their stoichiometry, permits them to be used in systems with many wavelength ranges or broad wavelength ranges, such as CWDM.

Increasingly, SOAs are key elements in photonic integrated circuits, both for applications of their nonlinear functionality as well as for gain blocks. PICs of moderate size are being developed and can incorporate tens of active elements, including SOAs. For example, SOAs have been used as power amplifiers for ten-channel 40 Gb/s per channel transmitter-PICs.⁹⁵ SOAs are used as gain blocks and also the basic nonlinear elements in medium-size PICs for functionality such as optical packet forwarding,⁹⁶ wavelength selection,⁵⁰ and wavelength conversion.⁹⁷

19.10 REFERENCES

1. J. Crowe and W. Ahearn, “9B8—Semiconductor Laser Amplifier,” *IEEE J. Quantum Electron.* 2(8):283–289, 1966.
2. T. Mukai and Y. Yamamoto, “Gain, Frequency Bandwidth, and Saturation Output Power of AlGaAs DH Laser Amplifiers,” *IEEE J. Quantum Electron.* 17(6):1028–1034, 1981.
3. G. Eisenstein, R. M. Jopson, R. A. Linke, C. A. Burrus, U. Koren, M. S. Whalen, and K. L. Hall, “Gain Measurements of InGaAsP 1.5 μm Optical Amplifiers,” *Electron. Lett.* 21(23):1076–1077, 1985.
4. M. J. O’Mahony, “Semiconductor Laser Optical Amplifiers for Use in Future Fiber Systems,” *J. Lightwave Technol.* 6(4):531–544, 1988.
5. T. Saitoh and T. Mukai, “Recent Progress in Semiconductor Laser Amplifiers,” *J. Lightwave Technol.* 6(11):1656–1664, 1988.
6. P. Doussi re, P. Garabedian, C. Graver, D. Bonnevie, T. Fillion, E. Derouin, M. Monnot, J. G. Provost, D. Leclerc, and M. Klenk, “1.55 μm Polarisation Independent Semiconductor Optical Amplifier with 25 dB Fiber to Fiber Gain,” *IEEE Photon. Technol. Lett.* 6(2):170–172, 1994.
7. P. J. A. Thijs, L. F. Tiemeijer, P. I. Kuindersma, J. J. M. Binsma, and T. van Dongen, “High-Performance 1.5 μm Wavelength InGaAs-InGaAsP Strained Quantum Well Lasers and Amplifiers,” *IEEE J. Quantum Electron.* 27(6):1426–1439, 1991.

8. J. Y. Emery, T. Ducellier, M. Bachmann, P. Doussière, F. Pommereau, R. Ngo, F. Gaborit, L. Goldstein, G. Laube, and J. Barrau, "High Performance 1.55 μm Polarisation-Insensitive Semiconductor Optical Amplifier Based on Low-Tensile Strained Bulk GaInAsP," *Electron. Lett.* **33**(12):1083–1084, 1997.
9. K. Magari, M. Okamoto, and Y. Noguchi, "1.55 μm Polarization-Insensitive High Gain Tensile-Strained-Barrier MQW Optical Amplifier," *IEEE Photon. Technol. Lett.* **3**(11):998–1000, 1991.
10. L. F. Tiemeijer, P. J. A. Thijs, T. van Dongen, R. W. M. Slootweg, J. M. M. van der Heijden, J. J. M. Binsma, and M. P. C. M. Krijn, "Polarization Insensitive Multiple Quantum Well Laser Amplifiers for the 1300 nm Window," *Appl. Phys. Lett.* **62**(8):826–828, 1993.
11. M. A. Newkirk, B. I. Miller, U. Koren, M. G. Young, M. Chien, R. M. Jopson, and C. A. Burrus, "1.5 μm Multi-quantum-Well Semiconductor Optical Amplifier with Tensile and Compressively Strained Wells for Polarization-Independent Gain," *IEEE Photon. Technol. Lett.* **5**(4):406–408, 1993.
12. A. E. Kelly, I. F. Lealman, L. J. Rivers, S. D. Perrin, and M. Silver, "Polarisation Insensitive, 25 dB Gain Semiconductor Laser Amplifier without Antireflection Coatings," *Electron. Lett.* **32**(19):1835–1836, 1996.
13. L. F. Tiemeijer, P. J. A. Thijs, T. van Dongen, J. J. M. Binsma, and E. J. Jansen, "Polarization Resolved, Complete Characterization of 1310 nm Fiber Pigtailed Multiple-Quantum-Well Optical Amplifiers," *J. Lightwave Technol.* **14**(6):1524–1533, 1996.
14. N. A. Olsson, "Lightwave Systems with Optical Amplifiers," *J. Lightwave Technol.* **7**(7):1071–1082, 1989.
15. A. E. Siegman, *Lasers*, Section 7.7, pages 297–303. University Science Books, Mill Valley, 1986 (from Eq. 82).
16. K. Morito, M. Ekawa, T. Watanabe, and Y. Kotaki, "High-Output-Power Polarization-Insensitive Semiconductor Optical Amplifier," *J. Lightwave Technol.* **21**(1):176–181, 2003.
17. J. M. Wiesenfeld, "Gain Dynamics and Associated Nonlinearities in Semiconductor Optical Amplifiers," *Int. J. High-Speed Electron. Sys.* **7**(1):179–222, 1996.
18. M. Bachmann, P. Doussière, J. Y. Emery, R. N'Go, F. Pommereau, L. Goldstein, G. Soulage, and A. Jourdan, "Polarisation-Insensitive Clamped-Gain SOA with Integrated Spot-Size Converter and DBR Gratings for WDM Applications at 1.55 μm Wavelength," *Electron. Lett.* **32**(22):2076–2078, 1996.
19. L. F. Tiemeijer, G. N. van den Hoven, P. J. A. Thijs, T. van Dongen, J. J. M. Binsma, and E. J. Jansen, "1310-nm DBR-type MQW Gain-Clamped Semiconductor Optical Amplifiers with AM-CATV-Grade Linearity," *IEEE Photon. Technol. Lett.* **8**(11):1453–1455, 1996.
20. D. A. Francis, S. P. Djaili, and J. D. Walker, "A Single-Chip Linear Optical Amplifier," In *Proc. Optical Fiber Communication Conference*, Anaheim, California, 2001. Paper PD13.
21. L. H. Spiekman, J. M. Wiesenfeld, A. H. Gnauck, L. D. Garrett, G. N. van den Hoven, T. van Dongen, M. J. H. Sander-Jochem, and J. J. M. Binsma, "Transmission of 8 DWDM Channels at 20 Gb/s over 160 km of Standard Fiber Using a Cascade of Semiconductor Optical Amplifiers," *IEEE Photon. Technol. Lett.* **12**(6):717–719, 2000.
22. V. G. Mutalik, G. van den Hoven, and L. Tiemeijer, "Analog Performance of 1310-nm Gain-Clamped Semiconductor Optical Amplifiers," In *Proc. Optical Fiber Communication Conference*, pages 266–267, Dallas, Texas, 1997.
23. B. Nyman, D. Favin, and G. Wolter, "Automated System for Measuring Polarization Dependent Loss," In *Proc. Optical Fiber Communication Conference*, pages 230–231, San Jose, California, 1994.
24. K. Morito, M. Ekawa, T. Watanabe, and Y. Kotaki, "High-Output-Power Polarization-Insensitive Semiconductor Optical Amplifier," *J. Lightwave Technol.* **21**(1):176–181, 2003.
25. S. Tanaka, S. Tomabechi, A. Uetake, M. Ekawa, and K. Morito, "Record High Saturation Output Power (+20 dBm) and Low NF (6.0 dB) Polarisation-Insensitive MQW-SOA Module," *Electron. Lett.* **42**(18):1059–1060, 2006.
26. K. Morito, S. Tanaka, S. Tomabechi, and A. Kuramata, "A Broad-Band MQW Semiconductor Optical Amplifier with High Saturation Output Power and Low Noise Figure," *IEEE Photon. Technol. Lett.* **17**(5):974–976, 2005.
27. K. Morito and S. Tanaka, "Record High Saturation Power (+22 dBm) and Low Noise Figure (5.7 dB) Polarization-Insensitive SOA Module," *IEEE Photon. Technol. Lett.* **17**(6):1298–1300, 2005.
28. P. W. Juodawlkis, J. J. Plant, L. J. Missaggia, K. E. Jensen, and F. J. O'Donnell, "Advances in 1.5- μm InGaAsP/InP Slab-Coupled Optical Waveguide Amplifiers (SCOWAs)," In *LEOS 2007. The 20th Annual Meeting of the IEEE*, pages 309–310, Oct. 2007.
29. W. Loh, J. J. Plant, F. J. O'Donnell, and P. W. Juodawlkis, "Noise Figure of a Packaged, High-Power Slab-Coupled Optical Waveguide Amplifier (SCOWA)," In *LEOS 2008. The 21st Annual Meeting of the IEEE*, pages 852–853, Nov. 2008.

30. T. Akiyama, M. Ekawa, M. Sugawara, K. Kawaguchi, H. Sudo, A. Kuramata, H. Ebe, and Y. Arakawa, "An Ultrawide-Band Semiconductor Optical Amplifier Having an Extremely High Penalty-Free Output Power of 23 dBm Achieved with Quantum Dots," *IEEE Photon. Technol. Lett.* **17**(8):1614–1616, 2005.
31. N. Yasuoka, K. Kawaguchi, H. Ebe, T. Akiyama, M. Ekawa, K. Morito, M. Sugawara, and Y. Arakawa, "Quantum-Dot Semiconductor Optical Amplifiers with Polarization-Independent Gains in 1.5- μ m Wavelength Bands," *IEEE Photon. Technol. Lett.* **20**(23):1908–1910, 2008.
32. P. L. Kuindersma, G. P. J. M. Cuijpers, J. G. L. Jennen, J. J. E. Reid, L. F. Teimeijer, H. de Waardt, and A. J. Boot, "10 Gbit/s RZ Transmission at 1309 nm Over 420 km Using a Chain of Multiple Quantum Well Semiconductor Optical Amplifier Modules at 38 km Intervals," In *Proc. European Conference on Optical Communications* 1996. Paper TuD.2.1.
33. K. D. LaViolette, "CTB Performance of Cascaded Externally Modulated and Directly Modulated CATV Transmitters," *IEEE Photon. Technol. Lett.* **8**(2):281–283, 1996.
34. Y. Awaji, J. Inoue, H. Sotobayashi, F. Kubota, and T. Ozeki, "Nonlinear Interchannel Cross Talk of Linear Optical Amplifier (LOA) in DWDM Applications," In *Optical Fiber Communications Conference*, 2003. OFC 2003, pages 441–443 vol. 2, Mar. 2003.
35. L. H. Spiekman, J. M. Wiesenfeld, A. H. Gnauck, L. D. Garrett, G. N. van den Hoven, T. van Dongen, M. J. H. Sander-Jochem, and J. J. M. Binsma, " 8×10 Gb/s DWDM Transmission over 240 km of Standard Fiber Using a Cascade of Semiconductor Optical Amplifiers," *IEEE Photon. Technol. Lett.* **12**(8):1082–1084, 2000.
36. L. H. Spiekman, A. H. Gnauck, J. M. Wiesenfeld, and L. D. Garrett, "DWDM Transmission of Thirty Two 10 Gbit/s Channels through 160 km Link Using Semiconductor Optical Amplifiers," *Electron. Lett.* **36**(12):1046–1047, 2000.
37. Y. Sun, A. K. Srivastava, S. Banerjee, J. W. Sulhoff, R. Pan, K. Kantor, R. M. Jopson, and A. R. Chraplyvy, "Error-Free Transmission of 32×2.5 Gbit/s DWDM Channels over 125 km Using Cascaded In-Line Semiconductor Optical Amplifiers," *Electron. Lett.* **35**(21):1863–1865, 1999.
38. Y. Awaji, H. Sotobayashi, and F. Kubota, "Transmission of 80 Gb/s \times 6 WDM over 100 km Using Linear Optical Amplifiers," *IEEE Photon. Technol. Lett.* **17**(3):699–701, 2005.
39. H. K. Kim and S. Chandrasekhar, "Reduction of Cross-Gain Modulation in the Semiconductor Optical Amplifier by Using Wavelength Modulated Signal," *IEEE Photon. Technol. Lett.* **12**(10):1412–1414, 2000.
40. A. K. Srivastava, S. Banerjee, B. R. Eichenbaum, C. Wolf, Y. Sun, J. W. Sulhoff, and A. R. Chraplyvy, "A Polarization Multiplexing Technique to Mitigate WDM Crosstalk in SOAs," *IEEE Photon. Technol. Lett.* **12**(10):1415–1416, 2000.
41. N. S. Bergano, F. W. Kerfoot, and C. R. Davidson, "Margin Measurements in Optical Amplifier Systems," *IEEE Photon. Technol. Lett.* **5**(3):304–306, 1993.
42. Z. Li, Y. Dong, J. Mo, Y. Wang, and C. Lu, "1050-km WDM Transmission of 8×10.709 Gb/s DPSK Signal Using Cascaded In-Line Semiconductor Optical Amplifier," *IEEE Photon. Technol. Lett.* **16**(7):1760–1762, 2004.
43. P. P. Iannone, K. C. Reichmann, and L. Spiekman, "In-Service Upgrade of an Amplified 130-km Metro CWDM Transmission System Using a Single LOA with 140-nm Bandwidth," In *Opt. Fiber Commun. Conf.*, vol. 2, pages 548–549. OSA, 2003. Paper ThQ3.
44. P. P. Iannone, H. H. Lee, K. C. Reichmann, X. Zhou, M. Du, B. Pálsdóttir, K. Feder, P. Westbrook, K. Brar, J. Mann, and L. Spiekman, "Four Extended-Reach TDM PONs Sharing a Bidirectional Hybrid CWDM Amplifier," *J. Lightwave Technol.* **26**(1):138–143, 2008.
45. T. Ismail, C. P. Liu, J. E. Mitchell, A. J. Seeds, X. Qian, A. Wonfor, R. V. Penty, and I. H. White, "Transmission of 37.6-GHz QPSK Wireless Data over 12.8-km Fiber with Remote Millimeter-Wave Local Oscillator Delivery Using a Bi-Directional SOA in a Full-Duplex System with 2.2-km CWDM Fiber Ring Architecture," *IEEE Photon. Technol. Lett.* **17**(9):1989–1991, 2005.
46. E. Almstrom, C. P. Larsen, L. Gillner, W. H. van Berlo, M. Gustavsson, and E. Berglind, "Experimental and Analytical Evaluation of Packaged 4×4 InGaAsP/InP Semiconductor Optical Amplifier Gate Switch Matrices for Optical Networks," *J. Lightwave Technol.* **14**(6):996–1004, 1996.
47. G. A. Fish, B. Mason, L. A. Coldren, and S. P. DenBaars, "Compact, 4×4 InGaAsP-InP Optical Crossconnect with a Scalable Architecture," *IEEE Photon. Technol. Lett.* **10**(9):1256–1258, 1998.
48. Y. Maeno, Y. Suemura, A. Tajima, and N. Henmi, "A 2.56-Tb/s Multiwavelength and Scalable Switch-Fabric for Fast Packet-Switching Networks," *IEEE Photon. Technol. Lett.* **10**(8):1180–1182, 1998.
49. A. Tajima, N. Kitamura, S. Takahashi, S. Kitamura, Y. Maeno, Y. Suemura, and N. Henmi, "10-Gb/s/Port Gated Divider Passive Combiner Optical Switch with Single-Mode to Multimode Combiner," *IEEE Photon. Technol. Lett.* **10**(1):162–164, 1998.

50. N. Kikuchi, Y. Shibata, H. Okamoto, Y. Kawaguchi, S. Oku, H. Ishii, Y. Yoshikuni, and Y. Tohmori, "Error-Free Signal Selection and High-Speed Channel Switching by Monolithically Integrated 64-channel WDM Channel Selector," *Electron. Lett.* **38**(15):823–824, 2002.
51. T. Sakamoto, A. Okada, M. Hirayama, Y. Sakai, O. Moriwaki, I. Ogawa, R. Sato, K. Noguchi, and M. Matsuoka, "Optical Packet Synchronizer Using Wavelength and Space Switching," *IEEE Photon. Technol. Lett.* **14**(9):1360–1362, 2002.
52. D. Chiaroni, "Packet Switching Matrix: A Key Element for the Backbone and the Metro," *IEEE J. Sel. Areas Commun.* **21**(7):1018–1025, 2003.
53. L. Dittmann, C. Develder, D. Chiaroni, F. Neri, F. Callegati, W. Koerber, A. Stavdas, et al., "The European IST Project DAVID: A Viable Approach toward Optical Packet Switching," *IEEE J. Sel. Areas Commun.*, **21**(7):1026–1040, 2003.
54. U. Koren, B. I. Miller, M. G. Young, T. L. Koch, R. M. Jopson, A. H. Gnauck, J. D. Evankow, and M. Chien, "High-Frequency Modulation of Strained Layer Multiple Quantum Well Optical Amplifiers," *Electron. Lett.* **27**(1):62–64, 1991.
55. M. D. Feuer, J. M. Wiesenfeld, J. S. Perino, C. A. Burrus, G. Raybon, S. C. Shunk, and N. K. Dutta, "Single-Port Laser-Amplifier Modulators for Local Access," *IEEE Photon. Technol. Lett.* **8**(9):1175–1177, 1996.
56. K. Y. Cho, Y. Takuchima, and Y. C. Chung, "10-Gb/s Operation of RSOA for WDM PON," *IEEE Photon. Technol. Lett.* **20**(18):1533–1535, 2008.
57. M. Koga, N. Tokura, and K. Nawata, "Gain-Controlled All-Optical Inverter Switch in a Semiconductor Laser Amplifier," *Appl. Opt.* **27**(19):3964–3965, 1988.
58. B. Glance, J. M. Wiesenfeld, U. Koren, A. H. Gnauck, H. M. Presby, and A. Jourdan, "High-Performance Optical Wavelength Shifter," *Electron. Lett.* **28**(18):1714–1715, 1992.
59. J. M. Wiesenfeld and B. Glance, "Cascadability and Fanout of Semiconductor Optical Amplifier Wavelength Shifter," *IEEE Photon. Technol. Lett.* **4**(10):1168–1171, 1992.
60. G. Lenz, E. P. Ippen, J. M. Wiesenfeld, M. A. Newkirk, and U. Koren, "Femtosecond Dynamics of the Nonlinear Anisotropy in Polarization Insensitive Semiconductor Optical Amplifiers," *Appl. Phys. Lett.* **68**:2933–2935, 1996.
61. S. L. Danielsen, C. Joergensen, M. Vaa, B. Mikkelsen, K. E. Stubkjaer, P. Doussiere, F. Pommerau, L. Goldstein, R. Ngo, and M. Goix, "Bit Error Rate Assessment of 40 Gbit/s All Optical Polarization Independent Wavelength Converter," *Electron. Lett.* **32**(18):1688–1690, 1996.
62. T. Durhuus, B. Mikkelsen, C. Joergensen, S. L. Danielsen, and K. E. Stubkjaer, "All-Optical Wavelength Conversion by Semiconductor Optical Amplifiers," *J. Lightwave Technol.* **14**(6):942–954, 1996.
63. D. Marcenac and A. Mecozzi, "Switches and Frequency Converters Based on Crossgain Modulation in Semiconductor Optical Amplifiers," *IEEE Photon. Technol. Lett.* **9**(6):749–751, 1997.
64. J. M. Wiesenfeld, B. Glance, J. S. Perino, and A. H. Gnauck, "Wavelength Conversion at 10 Gb/s Using a Semiconductor Optical Amplifier," *IEEE Photon. Technol. Lett.* **5**(11):1300–1303, 1993.
65. R. J. Manning and D. A. O. Davies, "Three-Wavelength Device for All-Optical Signal Processing," *Opt. Lett.* **19**(12):889–891, 1994.
66. T. Akiyama, N. Hatori, Y. Nakata, H. Ebe, and M. Sugawara, "Pattern-Effect-Free Semiconductor Optical Amplifier Achieved Using Quantum Dots," *Electron. Lett.* **38**(19):1139–1140, 2002.
67. L. Zhang, I. Kang, A. Bhardwaj, N. Sauer, S. Cabot, J. Jaques, and D. T. Neilson, "Reduced Recovery Time Semiconductor Optical Amplifier Using p-Type-Doped Multiple Quantum Wells," *IEEE Photon. Technol. Lett.* **18**(22):2323–2325, 2006.
68. J. S. Perino, J. M. Wiesenfeld, and B. Glance, "Fiber Transmission of 10 Gbit/s Signals Following Wavelength Conversion Using a Travelling-Wave Semiconductor Optical Amplifier," *Electron. Lett.* **30**(3):256–258, 1994.
69. G. P. Agrawal and N. K. Dutta, *Long-Wavelength Semiconductor Lasers*. Van Nostrand Reinhold, New York, 1986.
70. J. Leuthold, D. M. Marom, S. Cabot, J. J. Jaques, R. Ryf, and C. R. Giles, "All-Optical Wavelength Conversion Using a Pulse Reformating Optical Filter," *J. Lightwave Technol.* **22**(1):186–192, 2004.
71. A. D. Ellis, A. E. Kelly, D. Nasset, D. Pitcher, D. G. Moodie, and R. Kashyup, "Error Free 100 Gbit/s Wavelength Conversion Using Grating Assisted Cross Gain Modulation in 2 mm Long Semiconductor Amplifier," *Electron. Lett.* **34**(20):1958–1959, 1998.

72. Y. Liu, E. Tangdiongga, Z. Li, H. de Waardt, A. M. J. Koonen, G. D. Khoe, X. Shu, I. Bennion, and H. J. S. Dorren, "Error-Free 320-Gb/s All-Optical Wavelength Conversion Using a Single Semiconductor Optical Amplifier," *J. Lightwave Technol.* **25**(1):103–108, 2007.
73. T. Durhuus, C. Joergensen, B. Mikkelsen, R. J. S. Pedersen, and K. E. Stubkjaer, "All Optical Wavelength Conversion by SOAs in a Mach-Zehnder Configuration," *IEEE Photon. Technol. Lett.* **6**(1):53–55, 1994.
74. J. Leuthold, P. A. Besse, E. Gamper, M. Dulk, St. Fischer, and H. Melchior, "Cascadable Dual-Order Mode All-Optical Switch with Integrated Data- and Control Signal Separators," *Electron. Lett.* **34**(16):1598–1600, 1998.
75. B. Dagens, C. Janz, D. Leclerc, V. Verdrager, F. Poingt, I. Guillemot, F. Gaborit, and D. Ottenwalder, "Design Optimization of All-Active Mach-Zehnder Wavelength Converters," *IEEE Photon. Technol. Lett.* **11**(4):424–426, 1999.
76. S. Nakamura, Y. Ueno, K. Tajima, J. Sasaki, T. Sugimoto, T. Kato, T. Shimoda, M. Itoh, H. Hatakeyama, T. Tamanuki, and T. Sasaki, "Demultiplexing of 168-Gb/s Data Pulses with a Hybrid-Integrated Symmetric Mach-Zehnder All-Optical Switch," *IEEE Photon. Technol. Lett.* **12**(4):425–427, 2000.
77. N. S. Patel, K. L. Hall, and K. A. Rauschenbach, "40-Gbit/s Cascadable All-Optical Logic with an Ultrafast Nonlinear Interferometer," *Opt. Lett.* **21**(18):1466–1468, 1996.
78. J. Leuthold, C. H. Joyner, B. Mikkelsen, G. Raybon, J. L. Pleumeekers, B. I. Miller, K. Dreyer, and C. A. Burrus, "100 Gbit/s All-Optical Wavelength Conversion with Integrated SOA Delayed-Interference Configuration," *Electron. Lett.* **36**(13):1129–1130, 2000.
79. I. Kang, C. Dorrer, L. Zhang, M. Rasras, L. Buhl, A. Bhardway, S. Cabot, et al., "Regenerative All Optical Wavelength Conversion of 40-Gb/s DPSK Signals Using a Semiconductor Optical Amplifier Mach-Zehnder Interferometer," In *Proc. European Conference on Optical Communications* 6:29–30, 2005. Paper Th4.3.3.
80. G. P. Agrawal, "Population Pulsations and Nondegenerate Four-Wave Mixing in Semiconductor Lasers and Amplifiers," *J. Opt. Soc. Am. B* **5**(1):147–159, 1988.
81. J. Zhou, N. Park, J. W. Dawson, K. J. Vahala, M. A. Newkirk, and B. I. Miller, "Terahertz Four-Wave Mixing Spectroscopy for Study of Ultrafast Dynamics in a Semiconductor Optical Amplifier," *Appl. Phys. Lett.* **63**(9):1179–1181, 1993.
82. J. Zhou, N. Park, J. W. Dawson, K. J. Vahala, M. A. Newkirk, and B. I. Miller, "Efficiency of Broadband Four-Wave Mixing Wavelength Conversion Using Semiconductor Traveling-Wave Amplifiers," *IEEE Photon. Technol. Lett.* **6**(1):50–52, 1994.
83. A. Mecozzi, "Analytical Theory of Four-Wave Mixing in Semiconductor Amplifiers," *Opt. Lett.* **19**(12):892–894, 1994.
84. F. Girardin, J. Eckner, G. Guekos, R. Dall'Ara, A. Mecozzi, A. D'Ottavi, F. Martelli, S. Scotti, and P. Spano, "Low-Noise and Very High-Efficiency Fourwave Mixing in 1.5-mm-Long Semiconductor Optical Amplifiers," *IEEE Photon. Technol. Lett.* **9**(6):746–748, 1997.
85. A. E. Kelly, D. D. Marcenac, and D. Nasset, "40 Gbit/s Wavelength Conversion over 24.6 nm Using FWM in a Semiconductor Optical Amplifier with an Optimised MQW Active Region," *Electron. Lett.* **33**(25):2123–2124, 1997.
86. A. Mecozzi and J. M. Wiesenfeld, "The Roles of Semiconductor Optical Amplifiers in Optical Networks," *Opt. Photonics News* **12**(3):36–42, 2001.
87. R. M. Jopson and R. E. Tench, "Polarisation-Independent Phase Conjugation of Lightwave Signals," *Electron. Lett.* **29**(25):2216–2217, 1993.
88. G. Contestabile, A. D'Ottavi, F. Martelli, A. Mecozzi, P. Spano, and A. Tersigni, "Polarization-Independent Four-Wave Mixing in a Bidirectional Traveling wave Semiconductor Optical Amplifier," *Appl. Phys. Lett.* **75**(25):3914–3916, 1999.
89. R. Ludwig and G. Raybon, "BER Measurements of Frequency Converted Signals Using Four-Wave Mixing in a Semiconductor Laser Amplifier at 1, 2.5, 5 and 10 Gbit/s," *Electron. Lett.* **30**(4):338–339, 1994.
90. S. Diez, R. Ludwig, C. Schmidt, U. Feiste, and H. G. Weber, "160-Gb/s Optical Sampling by Gain-Transparent Four-Wave Mixing in a Semiconductor Optical Amplifier," *IEEE Photon. Technol. Lett.* **11**(11):1402–1404, 1999.
91. W. Pieper, C. Kurtzke, R. Schnabel, D. Breuer, R. Ludwig, K. Petermann, and H. G. Weber, "Nonlinearity-Insensitive Standard-Fibre Transmission Based on Optical-Phase Conjugation in a Semiconductor-Laser Amplifier," *Electron. Lett.* **30**(9):724–726, 1994.
92. D. D. Marcenac, D. Nasset, A. E. Kelly, M. Brierley, A. D. Ellis, D. G. Moodie, and C. W. Ford, "40 Gbit/s Transmission over 406 km of NDSF Using Mid-Span Spectral Inversion by Four-Wave-Mixing in a 2 mm Long Semiconductor Optical Amplifier," *Electron. Lett.* **33**(10):879–880, 1997.

93. M. L. Nielsen, D. J. Blumenthal, and J. Mork, "A Transfer Function Approach to the Small-Signal Response of Saturated Semiconductor Optical Amplifiers," *J. Lightwave Technol.* **18**(12):2151–2157, 2000.
94. O. Leclerc, B. Lavigne, E. Balmefrezol, P. Brindel, L. Pierre, D. Rouvillain, and F. Segueineau, "Optical Regeneration at 40 Gb/s and Beyond," *J. Lightwave Technol.* **21**(11):2779–2790, 2003.
95. R. Nagarajan, M. Kato, V. G. Dominic, C. H. Joyner, Jr. R. P. Schneider, A. G. Dentai, T. Desikan, et al., "400 Gbit/s (10 Channel \times 40 Gbit/s) DWDM Photonic Integrated Circuits," *Electron. Lett.* **41**(6):347–349, 2005.
96. V. Lal, M. Masanovic, D. Wolfson, G. Fish, C. Coldren, and D. J. Blumenthal, "Monolithic Widely Tunable Optical Packet Forwarding Chip in InP for All-Optical Label Switching with 40 Gbps Payloads and 10 Gbps Labels," In *Proc. European Conference on Optical Communications*, 6:25–26, 2005. Paper Th4.3.1.
97. P. Bernasconi, L. Zhang, W. Yang, N. Sauer, L. L. Buhl, J. H. Sinsky, I. Kang, S. Chandrasekhar, and D. T. Neilson, "Monolithically Integrated 40-Gb/s Switchable Wavelength Converter," *J. Lightwave Technol.* **24**(1):71–76, 2006.

This page intentionally left blank.

DO NOT DUPLICATE

OPTICAL TIME-DIVISION MULTIPLEXED COMMUNICATION NETWORKS

Peter J. Delfyett

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

20.1 GLOSSARY

Definitions

- Bandwidth.** A measure of the frequency spread of a signal, system, or information-carrying capacity.
- Chirping.** The time dependence of the instantaneous frequency of a signal.
- Commutator/decommutator.** Devices that assist in the sampling, multiplexing, and demultiplexing of time domain signals.
- Homogeneous broadening.** Physical mechanism that broadens the linewidth of a laser transition. The amount of broadening is exactly the same for all excited states.
- Kerr effect.** Dependence of a material's index of refraction on the square of an applied electric field.
- Mode partition noise.** Noise associated with mode competition in a multimode laser.
- Multiplexing/demultiplexing.** Process of combining and separating several independent signals that share a common communication channel.
- Passband.** Range of frequencies allowed to pass in a linear system.
- Picosecond.** One trillionth of a second.
- Pockel's effect.** Dependence of a material's index of refraction on an applied electric field.
- Photon lifetime.** Time associated with the decay in light intensity within an optical resonator.
- p-n* junction.** Region that joins two materials of opposite doping. This occurs when n-type and p-type materials are joined to form a continuous crystal.
- Quantum-confined Stark effect (QCSE).** Optical absorption induced by an applied electric field across a semiconductor quantum well.
- Quantum well.** Thin semiconductor layer sandwiched between material with a larger bandgap. The relevant dimension of the layer is on the order of 10 nm.
- Sampling.** The process of acquiring discrete values of a continuous signal.
- Spontaneous emission.** Energy decay mechanism to reduce the energy of excited states by the emission of light.

Spatial hole burning. The resultant nonuniform spatial distribution of optical gain in a material, owing to standing waves in an optical resonator.

Stimulated emission. Energy decay mechanism that is induced by the presence of light in matter to reduce the energy of excited states by the emission of light.

Terabit. Information equivalent to one trillion bits of information.

Abbreviations

ADC	analog to digital converter
APD	avalanche photodetector
CMI	code mark inversion
DFB	distributed feedback
DBR	distributed Bragg reflector
DS	digital signal
EDFA	erbium-doped fiber amplifier
FP	fabry-Perot
FDM	frequency division multiplexing
LED	light-emitting diode
NRZ	non-return-to-zero
OOK	on-off keying
OC-N	optical carrier (<i>N</i> th level)
PCM	pulse code modulation
PLM	pulse length modulation
PAM	pulse amplitude modulation
PPM	pulse position modulation
RZ	return-to-zero
SONET	synchronous optical network
SDH	synchronous digital hierarchy
STS	synchronous transmission signal
SPE	synchronous payload envelope
TOAD	terahertz optical asymmetric demultiplexer
SLALOM	semiconductor laser amplifier loop optical mirror
UNI	unbalanced nonlinear interferometer
TDM	time-division multiplexing
TDMA	time domain multiple access
WDM	wavelength division multiplexing
VCO	voltage-controlled oscillator

Symbols

$W(\text{Hz})$	bandwidth of a signal in units of hertz
$x_s(t)$	sampled version of a continuous function of time
$x(t)$	continuous analog signal
T	period
f_s	sampling frequency
p_T	periodic sampling pulse train

δ	delta function
n	index of refraction
n	integer
$X(\omega)$	frequency spectrum of the signal $x(t)$
ω	angular frequency (rad/s)
N	number of levels in an analog-to-digital converter
B	number of bits representing N levels in an analog-to-digital converter
Λ	grating period
λ	wavelength
φ	phase shift
τ_d or τ_p	photon decay time or photon lifetime
τ_{RT}	round trip propagation time of an optical cavity
$R_{1,2}$	mirror reflectivities

20.2 INTRODUCTION

Information and data services, such as voice, data, video, and Internet, are integral parts of our everyday personal and business lives. In 2001, the worldwide telecommunication traffic began being dominated by the fast-growing Internet traffic, as compared to the slow-growing voice traffic that dominated networks in years prior. In 2008, the amount of data traffic being transported over fiber optic networks is about 100 Tb/s, up from 1 Tb/s in 2001. The growth of and demand for bandwidth and communication services is growing steadily at about 80 percent per year globally, and high-speed photonic technologies must evolve to meet this demand. Currently deployed high-speed optical data transmission links are based on electrical time domain multiplexed (ETDM) transmission, where only electronic-signal-processing is performed at the transmitter and receiver to multiplex and demultiplex multiple users to and from the communication channel, respectively. In contrast, optical time domain multiplexing (OTDM) uses optical components and subsystems to multiplex and demultiplex user information at the transmitter and receiver.

The general organization of this chapter is to initially provide the reader with a brief review of digital signals and sampling. Following this introduction, time-division multiplexing (TDM) is introduced. These two sections provide the reader with a firm understanding of the overall system perspective as to how these networks are designed. To provide an understanding of the current state of the art, a review of selected high-speed optical and optoelectronic device technologies is given. Before a final summary and outlook toward future directions, a specific ultrahigh-speed optical time-division optical link is discussed, to coalesce the concepts with the discussed device technology.

20.3 MULTIPLEXING AND DEMULTIPLEXING

Fundamental Concepts

Multiplexing is a technique used to combine the information of multiple communication sites or users over a common communication medium and sending the information over a communication channel where the bandwidth, or information carrying capacity, is shared between each user. In the case where the shared medium is time, a communication link is created by combining the information from several independent sources and transmitting the information from each source simultaneously. This is done by temporally interleaving the bits of each source of information so that each user sends data for a very short period of time over the communication channel. The user waits

until all other users transmit their data before the first user can transmit another bit of information. At the receiver end, the data is demultiplexed for an intended user, while the rest of the information continues to its destination.

Sampling

An important concept in time-division multiplexing is being able to have a simple and effective method for converting real-world information into a form that is suitable for transmission by light. This process of transforming real signals into a form that is suitable for reliable transmission requires one to sample the analog signal to be sent and digitize and convert the analog signal to a stream of 1s and 0s. This process is usually performed by a sample and hold circuit, followed by an analog to digital converter (ADC). In the following section the concepts of signal sampling and digitization is reviewed, with the motivation to convey the idea of the robustness of digital communications.

Sampling Theorem

Time-division multiplexing relies on the fact that an analog bandwidth-limited signal may be exactly specified by taking samples of the signal, if the samples are taken sufficiently frequently. Time multiplexing is achieved by interleaving the samples of the individual signals. To see that any signal can be exactly represented by a sequence of samples, an understanding of the sampling theorem is needed. The theorem states that a real-valued bandwidth-limited signal that has no spectral components above a frequency of W hertz is determined uniquely by its value at uniform intervals spaced no greater than $1/(2W)$ s apart. This means that an analog signal can be completely reconstructed from a set of uniformly spaced discrete samples in time. The signal samples $x_s(t)$ are usually obtained by multiplying the signal $x(t)$ by a train of narrow pulses $p_T(t)$, with a time period $T = 1/f_s \leq 1/2W$. The process of sampling can be mathematically represented as

$$\begin{aligned} x_s(t) &= x(t) \cdot p_T(t) \\ &= x(t) \cdot \sum_{n=-\infty}^{+\infty} \delta(t - nT) \\ &= \sum_{n=-\infty}^{+\infty} x(nT) \delta(t - nT) \end{aligned} \quad (1)$$

where it is assumed that the sampling pulses are ideal impulses and n is an integer. Defining the Fourier transform and its inverse as

$$X(\omega) = \int_{-\infty}^{+\infty} x(t) \exp(-j\omega t) dt \quad (2)$$

and

$$x(t) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} X(\omega) \exp(+j\omega t) d\omega \quad (3)$$

one can show that the spectrum $X_s(\omega)$ of the signal $x_s(t)$ is given by

$$\begin{aligned} X_s(\omega) &= \frac{1}{T} \sum P\left(\frac{2\pi n}{T}\right) \cdot X\left(\omega - \frac{2\pi n}{T}\right) \\ &= \frac{1}{T} P(\omega) \cdot \sum X\left(\omega - \frac{2\pi n}{T}\right) \end{aligned} \quad (4)$$

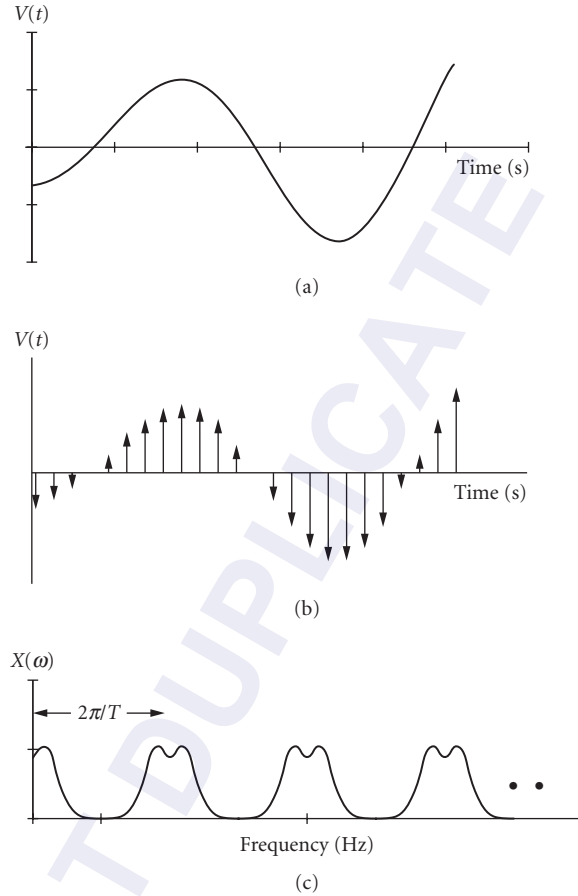


FIGURE 1 An analog bandwidth-limited signal (a), along with its sampled counterpart, sampled at a rate of ~ 8 times the Nyquist rate (b), and (c) frequency spectrum of a band-limited signal that has been sampled at a rate $T = 1/2W$, where W is the bandwidth of the signal.

In the case of the sampling pulses, p , being perfect delta functions, and given that the Fourier transform of $\delta(t)$ is 1, the signal spectrum is given by

$$X_s = \sum X\left(\omega - \frac{2\pi n}{T}\right) \quad (5)$$

This is represented pictorially in Fig. 1a to c. In Fig. 1a and b an analog signal and its sampled version, where the sample intervals is ~ 8 times the nominal sample rate of $1/(2W)$ are represented. From Fig. 1c it is clear that the spectrum of the signal is repeated in frequency every $2\pi/T$ Hz, if the sample rate T is $1/(2W)$. By employing (passing the signal through) an ideal rectangular low-pass filter, that is, a uniform (constant) passband with a sharp cutoff, centered at dc with a bandwidth of $2\pi/T$, the signal can be completely recovered. This filter characteristic implies an impulse response of

$$h(t) = 2W \frac{\sin(2\pi Wt)}{2\pi Wt} \quad (6)$$

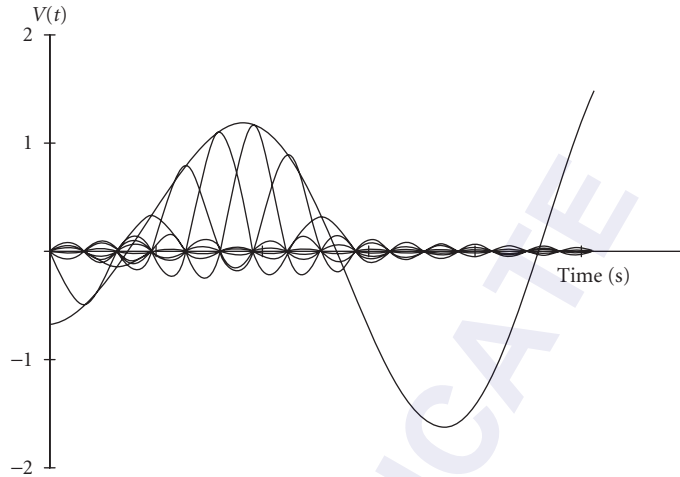


FIGURE 2 Temporal reconstruction of the sampled signal after passing the samples through a rectangular filter.

The reconstructed signal can now be given as

$$\begin{aligned}
 x(t) &= 2W \sum_{n=-\infty}^{+\infty} x(nT) \cdot \frac{\sin[2\pi W(t - nT)]}{2\pi W(t - nT)} \\
 &= x(t)/T \quad T = \frac{1}{2W}
 \end{aligned} \tag{7}$$

This reconstruction is shown in Fig. 2. It should be noted that the oscillating nature of the impulse response $h(t)$ interferes destructively with other sample responses, for times away from the centroid of each reconstructed sample.

Interleaving

The sampling principle can be exploited in time-division multiplexing by considering the ideal case of a single point-to-point link connecting N users to N other users over a single communication channel (see Fig. 3). As the rotary arm of the switch swings around, it samples each signal sequentially. The rotary switch at the receiving end is in synchronism with the switch at the sending end. The two switches make contact simultaneously at a similar number of contacts. With each revolution of the switch one sample is taken of each input signal and presented to the correspondingly numbered contact of the receiving end switch. The train of samples at Receiver 1, pass through a low-pass filter and at the filter output the original signal $m(t)$ appears reconstructed.

When the signals to be multiplexed vary rapidly in time, electronic switching systems are employed, as opposed to simple mechanical switches. The transmitter commutator samples and combines samples, while the receiver decommutator separates or demultiplexes samples belonging to individual signals so that these signals may be reconstructed.

The interleaving of the samples that allow multiplexing is shown in Fig. 4. For illustrative purposes, only two analog signals are considered. Both signals are repetitively sampled at a sample rate of T ; however, the instant at which the samples of each signal are taken are different. The input signal to Receiver 1 in Fig. 3, is the train of samples from Transmitter 1 and the input signal to Receiver 2 is

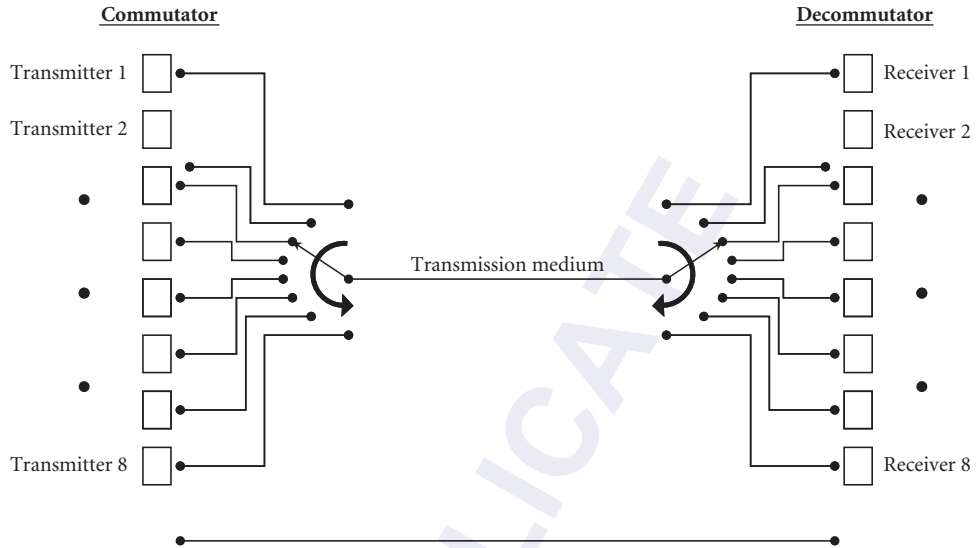


FIGURE 3 Illustration of a time multiplexer/demultiplexer based on simple mechanical switches, called commutators and decommutators.

the train of samples from Transmitter 2. The relative timing of the sampled signals of Transmitter 1 has been drawn to be exactly between the samples of Transmitter 2 for clarity; however, in practice, these samples would be separated by a smaller timing interval to accommodate additional temporally multiplexed signals.

Demultiplexing—Synchronization of Transmitter and Receiver

In any type of time-division multiplexing system, it is required that the sampling at the transmitter end and the demultiplexing at the receiver end be synchronized to each other. As an example, consider the diagram of the commutator of Fig. 3. When the transmitting multiplexer is set in a position that samples and transmits information from Transmitter 1, the receiving demultiplexer must be in a position to demultiplex and receive information that is intended for Receiver 1. To accomplish this timing synchronization, the receiver has a local clock signal that controls the timing of the commutator to switch from one time slot to the next. The clock signal may be a narrow band sinusoidal signal from which an appropriate clocking signal, with sufficiently fast rising edges of the appropriate signal strength, can be derived. The repetition rate of the clock in a simple configuration would then be equal to the sample rate of an individual channel times the number of channels being multiplexed, thereby assigning one time slot per clock cycle.

At the receiver end, the clock signal is required to keep the decommutator synchronized to the commutator, that is, running at the same rate. In addition, there must be additional timing information to provide agreement as to the relative positions, or phase of the commutator-decommutator pair, which ensures that information from Transmitter 1 is guaranteed to be received at the desired destination of Receiver 1. The time interval from the beginning of the time slot allocated to a particular channel, until the next recurrence of that particular time slot is commonly referred to as a *frame*. As a result, timing information is required at both the bit (time slot) and frame levels. A common arrangement in time-division-multiplexed systems is to allow for one or more time slots per frame to provide timing information, depending on the temporal duration of the transmitted frame. It should

be noted that there are a variety of methods for providing timing information, such as directly using a portion of the allocated bandwidth, as mentioned above, or alternatively, recovering a clock signal by deriving timing information directly from the transmitted data.

Digital Signals—Analog to Digital Conversion

The sampled signal, as shown in Fig. 4, represent the actual values of the analog signal at the sampling instants. In a practical communication system or in a realistic measurement setup, the received or measured values can never be absolutely correct, because of the noise introduced by the transmission channel or small inaccuracies impressed on the received data owing to the detection or measurement process. It turns out that it is sufficient to transmit and receive only the quantized values of the signal samples. The quantized values of sampled signals represented to the nearest digit, may be represented in a binary form or in any coded form using only 1s and 0s. For example, sampled values of a signal between 2.5 and 3.4 would be represented by the quantized value of 3, and could be represented as 11, using 2 bits (in base 2 arithmetic). This method of representing a sampled analog signal is known as *pulse code modulation*. An error is introduced on the signal by this quantization process. The magnitude of this error is given by

$$\epsilon = \frac{0.4}{N} \quad (8)$$

where N is the number of levels determined by $N = 2^B$, and B is the B -bit binary code, for example, $B = 8$ for 8-bit words representing 256 levels. Thus one can minimize the error by increasing the number of levels, which is achieved by reducing the step size in the quantization process. It is interesting to note that using only 4 bit (16 levels), a maximum error of 2.5 percent is achieved, while increasing the number of bits to 8 (256 levels) gives a maximum error of 0.15 percent.

Optical Representation of Binary Digits and Line Coding

The binary digits can be represented and transmitted on an optical beam and passed through an optical fiber or transmitted in free space. The optical beam is modulated to form pulses to represent the sampled and digitized information. A family of four such representations is shown in Fig. 5. There are two particular forms of data transmission that are quite common in optical communications owing to the fact that their modulation formats occur naturally in both direct and externally modulated optical sources. These two formats are referred to as “non-return-to-zero” (NRZ), or “return-to-zero” (RZ). In addition to NRZ and RZ data formats, pulse-code-modulated data signals are transmitted in other codes that are designed to optimize the link performance, owing to channel constraints.

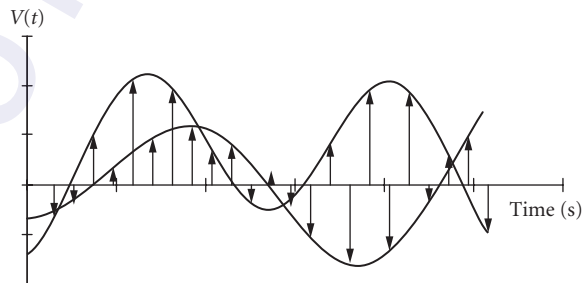


FIGURE 4 Two band-limited analog signals and their respective samples occurring at a rate of approximately six times the highest frequency, or three times the Nyquist rate.

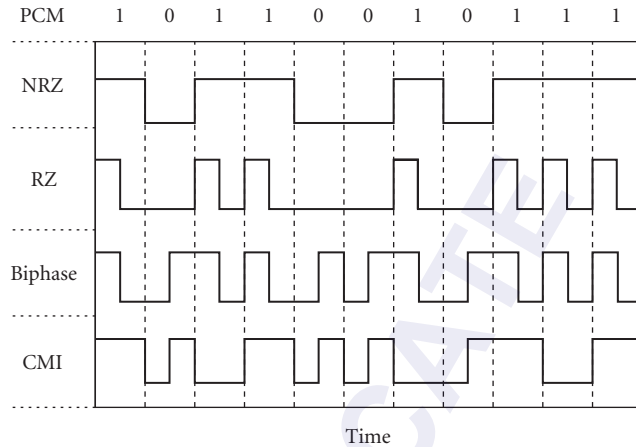


FIGURE 5 Line-coded representations of the pulse-code-modulated logic signal “10110010111.” NRZ: non-return to zero; RZ: return to zero; biphase, also commonly referred to as Manchester coding; CMI: code mark inversion.

Some important data transmission formats for optical time-division-multiplexed networks are code mark inversion (CMI), and Manchester coding or biphase coding. In CMI, the coded data has no transitions for logical 1 levels. Instead, the logic level alternates between a high and low level. For logical 0, on the other hand, there is always a transition from low to high at the middle of the bit interval. This transition for every logical 0 bit ensures proper timing recovery. For Manchester coding, logic 1 is represented by a return-to-zero pulse with a 50 percent duty cycle over the bit period (a half-cycle square wave), and logic 0 is represented by a similar return-to-zero waveform of “opposite phase,” hence the name biphase. The salient feature of both biphase and CMI coding is that their power spectra have significant energy at the bit rate, owing to the guarantee of a significant number of transitions from logic 1 to 0. This should be compared to the power spectra of RZ and NRZ data, which are shown in Fig. 6. The NRZ has no energy at the bit rate, while the RZ power spectrum does have energy at the bit rate, but the spectrum is also broad, having a width twice as large as NRZ.

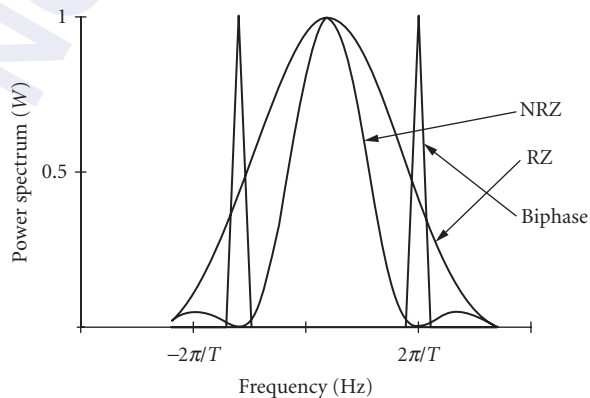


FIGURE 6 Power spectra of NRZ, RZ, and biphase line-coded data. Note the relative power at the bit rate.

The received data power spectra is important for TDM transmission links, where at the receiver end, a clock or synchronization signal is required to demultiplex the data. It is useful to be able to recover a clock or synchronization signal derived from the transmitted data, instead of using a portion of the channel bandwidth to send a clock signal. Therefore by choosing a transmission format with a large power spectral component at the transmitted bit rate, provides an easy method to recover a clock signal.

Timing Recovery

Time-division multiplexing and time-division multiple access networks inherently require timing signals to assist in demultiplexing individual signals from their multiplexed counterparts. One possible method is to utilize a portion of the communication bandwidth to transmit a timing signal. Technically, this is feasible, however, (1) this approach requires hardware dedicated to timing functions distributed at each network node that performs multiplexing and demultiplexing functions, and (2) network planners want to optimize the channel bandwidth without resorting to dedicating a portion of the channel bandwidth to timing functions. The desired approach is to derive a timing signal directly from the transmitted data. This allows the production of the required timing signals for multiplexing and demultiplexing without the need of using valuable channel bandwidth.

As suggested by Fig. 7, a simple method for recovering a timing signal from transmitted return-to-zero data is to use a bandpass filter to pass a portion of the power spectrum of the transmitted data. The filtered output from the tank circuit is a pure sinusoid that provides the timing information.

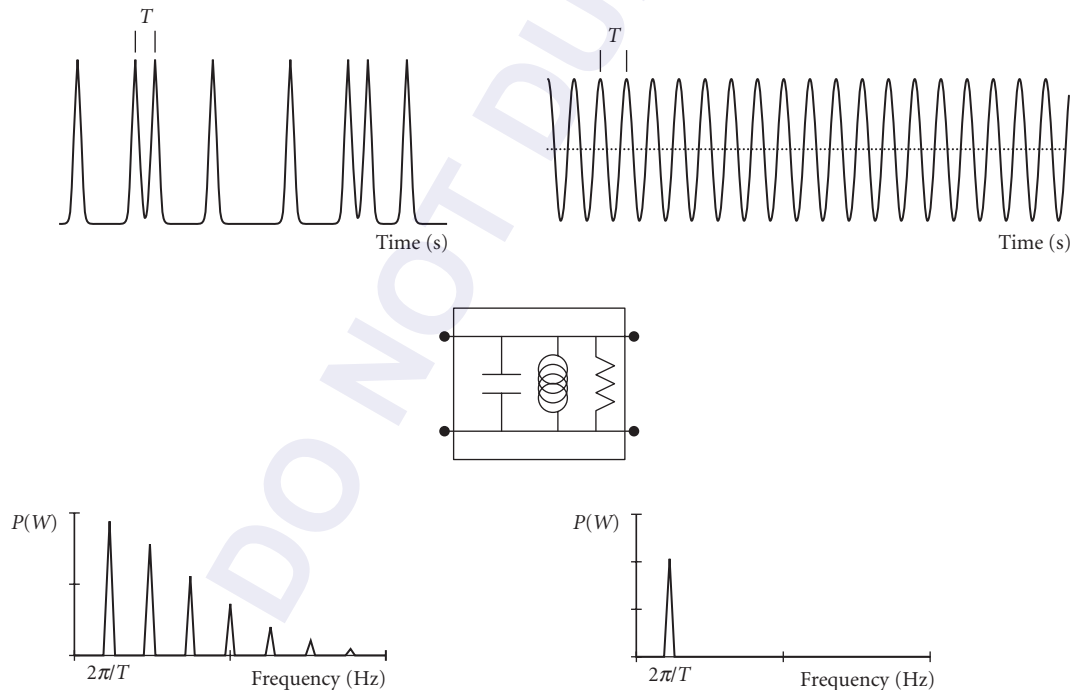


FIGURE 7 Principle of clock recovery using line filtering. Upper left: Input RZ data stream. Lower left: Power spectrum of a periodic RZ sequence. Center: Schematic of an electrical tank circuit for realizing a bandpass filter. Lower right: Power spectrum of the filtered signal. Upper right: Filtered time domain clock signal.

An important parameter to consider in line filtering is the quality factor, designated as the filter Q . Generally, the Q factor is defined as

$$Q = \frac{\omega_0}{\Delta\omega} \quad (9)$$

where ω_0 is the resonant frequency and $\Delta\omega$ is the bandwidth of the filter. It should also be noted that the Q is a measure of the amount of energy stored in the bandpass filter, such that the output from the filter decays exponentially at a rate directly proportional to Q . In addition, for bandpass filters based on passive electrical circuits, the output peak signal is directly proportional to Q . These two important physical features of passive line filtering imply that the filter output will provide a large and stable timing signal if the Q factor is large. However, since Q is inversely proportional to the filter bandwidth, a large Q typically implies a small filter bandwidth. As a result, if the transmitter bit rate and the resonant frequency of the tank circuit do not coincide, the clock output could be zero. In addition the clock output is very sensitive to the frequency offset between the transmitter and resonant frequency. Therefore, line filtering can provide a large and stable clock signal for large filter Q , but the same filter will not perform well when the bit rate of the received signal has a large frequency variation. In TDM bit timing recovery, the ability to recover the clock of an input signal over a wide frequency range is called frequency acquisition or locking range, and the ability to tolerate timing jitter and a long interval of zero transitions is called frequency tracking or hold over time. Therefore, the trade-off exists between the locking range (low Q) and hold over time (large Q) in line filtering.

A second general scheme to realize timing recovery and overcome the drawbacks of line filtering using passive linear components is the use of a phase-locked loop in conjunction with a voltage-controlled oscillator (VCO) (see Fig. 8a). In this case, two signals are fed into the mixer. One signal is derived from the data, for example, from a line-filtered signal possessing energy at the bit rate, while the second signal is a sinusoid generated from the VCO. The mixer is used as a phase detector and produces a DC voltage that is applied to the VCO to adjust its frequency of operation. The overall design of the PLL is to adjust the voltage of the PLL to track the frequency and phase of the input data signal. Owing to the active components in the PLL, this approach for timing recovery can realize a broad locking range, low insertion loss, and good phase-tracking capabilities. It should be noted that

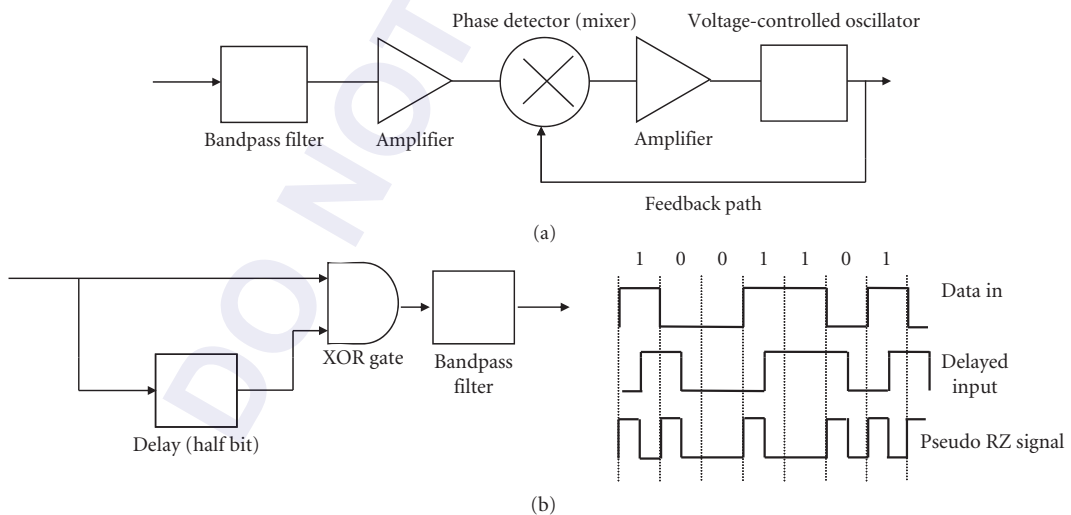


FIGURE 8 (a) Schematic diagram of a phase-locked loop using a mixer as a phase detector and a voltage-controlled oscillator to provide the clock signal that can track phase wander in the data stream. (b) Data format conversion between input NRZ data to RZ output data using an electronic logic gate. The subsequent RZ output is then suitable for use in a clock recovery device.

while the concepts for timing recovery described in this section were illustrated using techniques that are not directly applicable to ultrahigh-speed optical networking, the underlying principles will still hold for high-speed all-optical techniques. These approaches are discussed in more detail later in the section on device technology.

While both these techniques require the input data to be in the return-to-zero format, many data transmission links use nonreturn-to-zero line coding owing to its bandwidth efficiency. Unfortunately, in NRZ format there is no component in the power spectrum at the bit rate. As a result, some preprocessing of the input data signal is required before clock recovery can be performed. A simple method for achieving this is illustrated in Fig. 8b. The general concept is to present the data signal with a delayed version of the data at the input ports of a logic gate that performs the exclusive OR operation. The temporal delay, in this case, should be equal to 1/2 bit. The output of the XOR gate is a pseudo RZ data stream that can then be line filtered for clock recovery.

20.4 INTRODUCTION TO DEVICE TECHNOLOGY

Thus far, a general description of the concepts of digital communications and the salient features of TDM and TDMA have been presented. Next we address specific device technology that is employed in OTDM networks, for example, sources, modulators, receivers, clock recovery oscillators, demultiplexers, to provide an understanding of how and why specific device technology may be employed in a system to either optimize the network performance, to minimize cost, or to provide the maximum flexibility in supporting a wide variety of user applications.

Optical Time division Multiplexing—Serial versus Parallel

Optical time division multiplexing can generally be achieved by two main methods. The first method is referred as *parallel multiplexing*, while the second method classified as *serial multiplexing*. These two approaches are schematically illustrated in Fig. 9. The advantage of the parallel type of multiplexer is that it employs simple, linear passive optical components, not including the intensity modulator, and that the limitation in the transmission speed is not limited by the modulator or any other high-speed switching element. The drawback, however, is that the relative temporal delays between each channel must be accurately controlled and stabilized, which increases the complexity of this approach. Alternatively, the serial approach to multiplexing is simple to configure. In this approach a high-speed optical clock pulse train and modulation signal pulses are combined and introduced into an all-optical switch to create a modulated channel on the high bit-rate clock signal. Cascading this process allows all the channels to be independently modulated, with the requirement that the relative delay between each channel must be appropriately adjusted.

Device Technology—Transmitters

For advanced lightwave systems and networks, it is the semiconductor laser that dominates as the primary optical source that is used to generate the light that is modulated and transmitted as information. The reason for their dominance is that these devices are very small, typically a few hundred micrometers on a side, have excellent efficiency in converting electrons to photons, and are low cost. In addition, semiconductor diode lasers can generate optical signals at wavelengths of 1.3 and 1.55 μm . These wavelengths are important because they correspond to the spectral regions where optical signals experience minimal dispersion (spreading of the optical data bits) and minimal loss.

These devices initially evolved from simple light-emitting diodes, comprising a simple *p-n* junction, to Fabry-Perot (FP) semiconductor lasers, to distributed feedback (DFB) lasers and distributed Bragg reflector (DBR) lasers, and finally to mod-locked semiconductor diode lasers and optical fiber lasers. Below, a simple description of each of these devices is given, along with advantages and disadvantages that influence how these optical transmitter are deployed in current optical systems and networks.

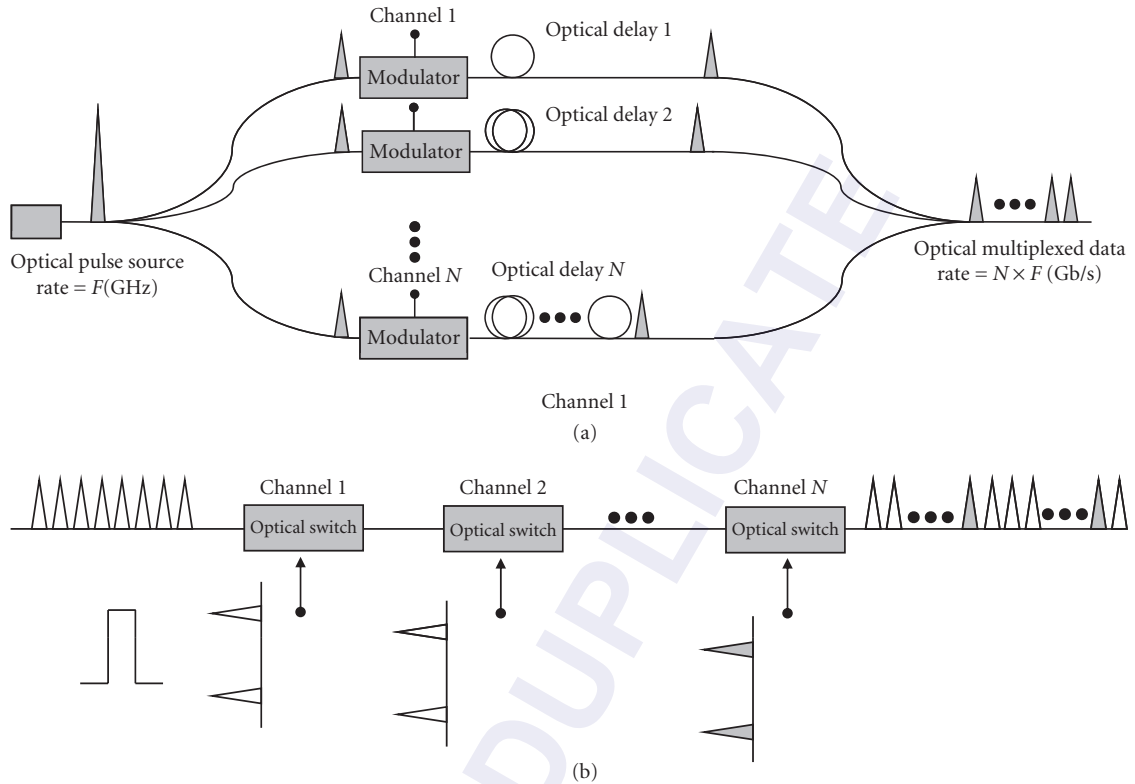


FIGURE 9 Schematic of optical time-division multiplexing for interleaving high-speed RZ optical pulses: (a) Parallel implementation and (b) serial implementation.

Fabry-Perot Semiconductor Lasers The Fabry-Perot semiconductor laser diode comprises a semiconductor p - n junction that is heavily doped and fabricated from a direct-gap semiconductor material. The injected current is sufficiently large to provide optical gain. The optical feedback is provided by mirrors, which are usually obtained by cleaving the semiconductor material along its crystal planes. The large refractive index difference between the crystal and the surrounding air causes the cleaved surfaces to act as reflectors. As a result, the semiconductor crystal acts as both the gain medium and as an optical resonator, or cavity (see Fig. 10). Provided that the gain coefficient is sufficiently large, the feedback transforms the device into an optical oscillator or laser diode.

Considering that the physical dimensions of the semiconductor diode laser are quite small, the short length of the diode forces the longitudinal mode spacing $c/2nL$ to be quite large. Here, c is the speed of light, L is the length of the diode chip, and n is the refractive index. Nevertheless, many of these modes can generally fit within the broad gain bandwidth allowed in a semiconductor diode laser. As an example, consider a FP laser diode operating at $1.3\ \mu\text{m}$, fabricated from the InGaAsP material system. If $n = 3.5$ and $L = 400\ \mu\text{m}$, the modes are spaced by 107 GHz, and corresponds to a wavelength spacing of 0.6 nm. In this device, the gain bandwidth can be 1.2 THz, corresponding to a wavelength spread of 7 nm, and as many as 11 modes can oscillate. Given that the mode spacing can be modified by cleaving the device so that only one axial mode exists within the gain bandwidth, the resulting device length would be approximately $36\ \mu\text{m}$, which is difficult to achieve. It should be noted that if the bias current is increased well above threshold, the device can tend to oscillate on a single longitudinal mode. However for telecommunications, it is very desirable to directly modulate

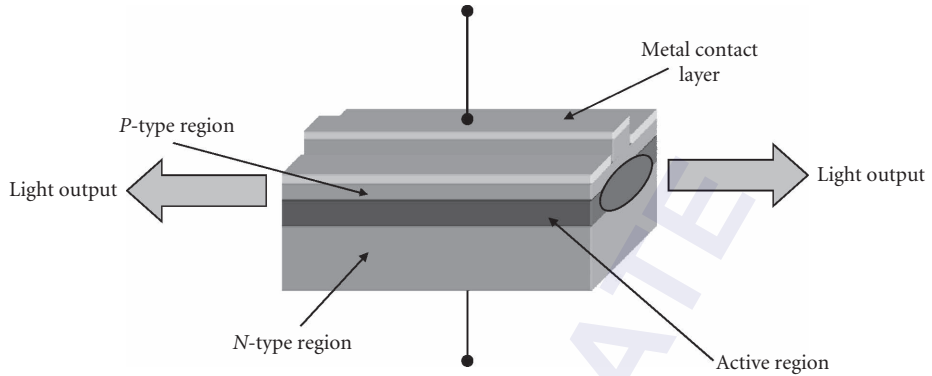


FIGURE 10 Schematic illustration of a simple Fabry-Perot semiconductor diode laser.

the laser, thus avoiding the cost of an external modulator. However, in the case of direct modulation, the output emission spectrum will be multimode, and as a result, effects of dispersion will broaden the optical data bits, and force the data rate to be reduced to avoid intersymbol interference. Given this effect, Fabry-Perot lasers tend to have a more limited use in longer optical links.

Distributed Feedback Lasers The effects of dispersion and the broad spectral emission from semiconductor LEDs and semiconductor Fabry-Perot laser diodes tend to reduce the overall optical data transmission rate. Thus, methods have been developed to design novel semiconductor laser structures that will only operate on a single longitudinal mode. This then will allow these devices to be directly modulated and allow for longer transmission paths since the overall spectral width is narrowed, and the effect of dispersion is minimized.

The preferred method of achieving single frequency operation from semiconductor diode lasers is to incorporate frequency-selective reflectors at both end of the diode chip or, alternately, fabricate the grating directly adjacent to the active layer. These two approaches result in devices referred to as distributed Bragg reflector lasers and distributed feedback lasers, respectively. In practice, it is easier to fabricate a single grating structure above the active layer, as opposed to two separate gratings at each end and as a result, the DFB laser has become the laser of choice for telecommunication applications. These devices operate with spectral widths on the order of a few megahertz, and have modulation bandwidths over 10 GHz. Clearly, the high modulation bandwidth and low spectral width make these devices well suited for direct modulation, or on-off keyed (OOK) optical networks. It should be noted that the narrow linewidth of a few megahertz is for the device operating in a continuous wave mode, while modulating the device will necessarily broaden the spectral width.

In DFB lasers, Bragg reflection gratings are employed along the longitudinal direction of the laser cavity and are used to suppress the lasing of additional longitudinal modes. As shown in Fig. 11a, a periodic structure, similar to a corrugated washboard, is fabricated over the active layer, where the periodic spacing is denoted as Λ . Owing to this periodic structure, both forward and backward traveling waves must interfere constructively with each other. In order to achieve this constructive interference between the forward and backward waves, the round trip phase change over one period should be $2\pi m$, where m is an integer and is called the order of the Bragg diffraction. With $m = 1$, the first-order Bragg wavelength, λ_B , is

$$2\pi = 2\Lambda(2\pi n/\lambda_B) \quad (10)$$

or

$$\lambda_B = 2\Lambda n \quad (11)$$

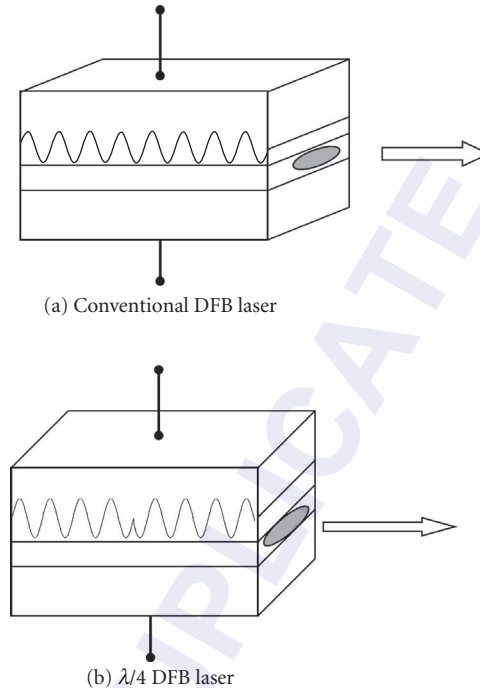


FIGURE 11 Schematic illustrations of distributed feedback lasers (DFB): (a) Conventional DFB and (b) Quarter-wave DFB, showing the discontinuity of the Bragg grating structure to achieve single wavelength operation.

where n is the refractive index of the semiconductor. Therefore, the period of the periodic structure determines the wavelength for the single mode output. In reality, a periodic DFB structure generates two main modes, symmetrically placed on either side of the Bragg wavelength λ_B . In order to suppress this dual frequency emission, and generate only one mode at the Bragg wavelength, a phase shift of $\lambda/4$ can be used to remove the symmetry. As shown in Fig. 11b, the periodic structure has a phase discontinuity of $\pi/2$ at the middle, which gives an equivalent $\lambda/4$ phase shift. Owing to the ability of the $\lambda/4$ DFB structure to generate a single frequency, narrow spectral linewidth, these devices are the preferred device for present telecommunications.

Mode-Locked Lasers Mode locking is a technique for obtaining very short bursts of light from lasers, and can be easily achieved employing both semiconductor and fiber gain media. As a result of mode locking, the light that is produced is automatically in a pulsed form that produces RZ data if passed through an external modulator being electrically driven with NRZ data. More importantly, *the temporal duration of the optical bits produced by mode locking is much shorter than the period of the driving signal*. To contrast this, consider a DFB laser whose light is externally modulated. In this case, the temporal duration of the optical bits will be equal to the temporal duration of the electrical pulses driving the external modulator. As a result, the maximum possible data transmission rate achievable from the DFB will be limited to the speed of the electronic driving signal. With mode locking; however, a low frequency electrical drive signal can be used to generate ultrashort optical bits. By following the light production with external modulation and optical bit interleaving, one can realize the ultimate in OTDM transmission rates. To show the difference between a mode-locked

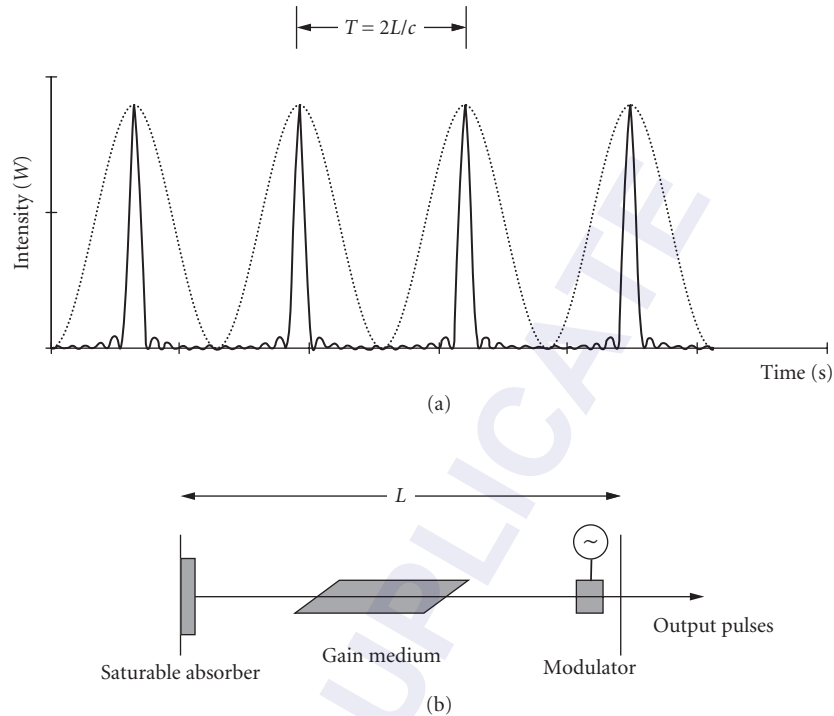


FIGURE 12 Optical intensity distribution of five coherent, phase-locked modes of a laser (a), and a schematic diagram of an external cavity mode-locked laser (b). Superimposed on the optical pulse train is a typical sinusoid that could be used to mode-lock the laser, showing that much shorter optical pulses can be obtained from a low-frequency signal.

pulse train and its drive, Fig. 12 plots a sinusoid and a mode-locked pulse train consisting of five locked optical modes.

To understand the process of mode locking, it should be recalled that a laser can oscillate on many longitudinal modes that are equally spaced by the longitudinal mode spacing, $c/(2nL)$. Normally, these modes oscillate independently; however, techniques can be employed to couple and lock their relative phases together. The modes can then be regarded as the components of a Fourier-series expansion of a periodic function of time of period $T = (2nL)/c$, which represents a periodic train of optical pulses. Consider, for example, a laser with multiple longitudinal modes separated by $c/2nL$. The output intensity of a perfectly mode-locked laser, as a function of time t , and axial position z , with M locked longitudinal modes, each with equal intensity, is given by

$$I(t, z) = M^2 |A|^2 \frac{\sin c^2[M(t - z/c)/T]}{\sin c^2[(t - z/c)T]} \quad (12)$$

where T is the periodicity of the optical pulses, and $\sin c(x)$ is $\sin(x)/x$. In practice, there are several methods to generate optical pulse trains by mode locking and generally fall into two categories: (1) active mode locking and (2) passive mode locking. In both cases, to lock the longitudinal modes in phase, the gain of the laser is allowed to increase above its threshold for a short duration, by opening and closing a shutter that is placed within the optical cavity. This allows a pulse of light to form. By allowing the light to propagate around the cavity and continually reopening and closing the shutter

at a rate inversely proportional to the round-trip time forms a stable, well-defined optical pulse is formed. If the shutter is realized by using an external modulator, the technique is referred to as *active mode locking*, where as if the shutter is realized by a device or material that is activated by the light intensity itself, the process is called *passive mode locking*. Both techniques can be used simultaneously and is referred to as *hybrid mode locking* (see Fig. 12b).

Direct and Indirect Modulation

To transmit information in OTDM networks, the light output of the laser source must be modulated in intensity. Depending on whether the output light is modulated by directly modulating the current source to the laser or whether the light is modulated external to, or after the light has been generated, the process of modulation can be classified as either (1) direct modulation or (2) indirect or external modulation (see Fig. 13a and b). With direct modulation, the light is directly modulated inside the light source, while external modulation, uses a separate external modulator placed after the laser source.

Direct modulation is used in many optical communication systems owing to its simple and cost-effective implementation. However, owing to the physics of laser action and the finite response of populating the lasing levels owing to current injection, the light output under direct modulation cannot respond to the input electrical signal instantaneously. Instead, there are turn-on delays and oscillations that occur when the modulating signal, which is used as the pumping current, has large and fast changes. As a result, direct modulation has several undesirable effects, such as frequency chirping and linewidth broadening. In frequency chirping, the spectrum of the output-generated light is time varying, that is, the wavelength and spectrum changes in time. This is owing to the fact

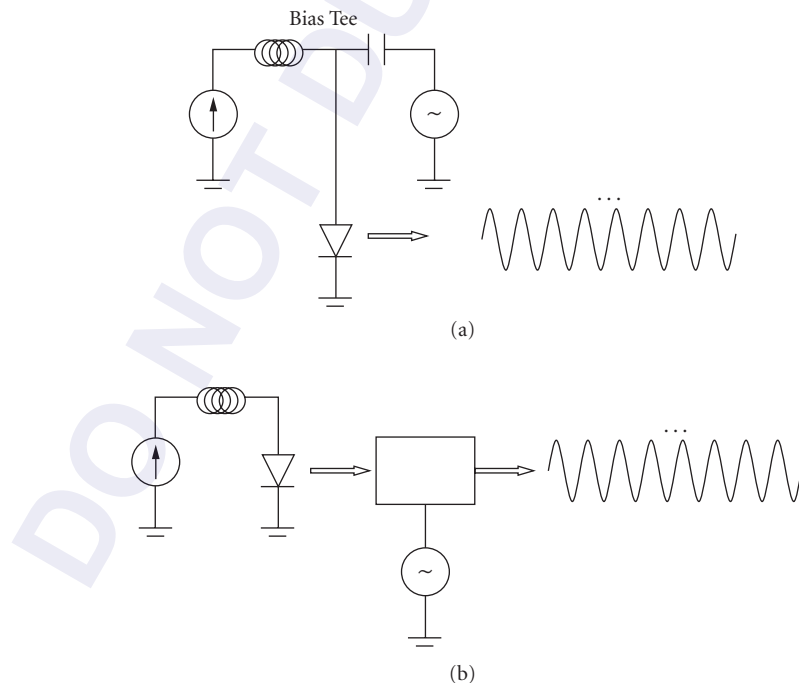


FIGURE 13 Illustrative example of direct modulation (a) and external modulation (b) of a laser diode.

that as the laser is turned on and off, the gain is changed from a very low value to a high value. Since the index of refraction of the laser diode is closely related to the optical gain of the device, as the gain changes, so does its index. It is this time varying refractive index that leads to frequency chirping, and is sometimes referred to as phase modulation.

External Modulation

To avoid the undesirable frequency chirping effects in DFB lasers and mode partition noise in FP lasers associated with direct modulation, external modulation provides an alternative approach to achieve light modulation with the added benefit of avoiding these detrimental effects. A typical external modulator consists of an optical waveguide in which the incident light propagates through and the refractive index or absorption of the medium is modulated by a signal that represents the data to be transmitted. Depending on the specific device, one can realize three basic types of external modulators: (1) electro-optic, (2) acousto-optic, and (3) electroabsorption. Generally, acousto-optic modulators respond slowly, on the order of several nanoseconds, and as a result are not used for external modulators in telecommunication applications. Electroabsorption (EA) modulators rely on the fact that the band edge of a semiconductor can be frequency shifted to realize an intensity modulation for a well-defined wavelength that is close to the band edge of the modulator. Linear frequency responses up to 50 GHz are possible; however, the fact that both the wavelength of the laser and the modulator must be accurately matched makes this approach more difficult to implement with individual devices. It should be noted, however, that EA modulators and semiconductor lasers can be integrated in the same device, helping to remove restrictions on matching the transmitter's and modulator's wavelength.

The typical desirable properties of an external modulator, from a communications perspective, are that these devices should possess a large modulation bandwidth, a large depth of modulation, a small insertion loss or loss of the signal light passing through the device, and a low electrical drive power. In addition, for some types of communication TDM links, a high degree of linearity between the drive signal and modulated light signal is required (typical for analog links), and an independence of input polarization (polarization diversity) is desired. Finally, low cost and small size of these devices are extremely useful for cost-effective and wide-area deployment.

Electro-Optic Modulators An electro-optic modulator can be a simple optical channel or waveguide that the light to be modulated propagates. The material that is chosen to realize the electro-optic modulator must possess an optical birefringence that can be controlled or adjusted by an external electrical field that is applied along or transverse to the direction of propagation of the light to be modulated. This birefringence means that the index of refraction is different for light that propagates in different directions in the crystal. If the input light has a well-defined polarization state, this light can be made to see, or experience, a different refractive index for different input polarization states. By adjusting the applied voltage to the electro-optic modulator, the polarization can be made to rotate, or the speed of the light can be slightly varied. This modification of the input light property can be used to realize a change of the output light intensity, by using a crossed polarizer or by interfering the modulated light with an exact copy of the unmodulated light. This can easily be achieved by using a waveguide interferometer, such as a Mach-Zehnder interferometer. If the refractive index is directly proportional to the applied electric field, the effect is referred to as the *Pockel's effect*.

Generally, for high-speed telecommunication applications, device designers employ the use of the electro-optic effect as a phase modulator in conjunction with an integrated Mach-Zehnder interferometer or an integrated directional coupler. Phase modulation (or delay/retardation modulation) does not effect the intensity of the input light beam. However, by incorporating a phase modulator in one branch of an interferometer, the resultant output light from the interferometer will be intensity modulated. Consider an integrated Mach-Zehnder interferometer in Fig. 14. If the waveguide divides the input optical power equally, the transmitted intensity is related to the output intensity by the well-known interferometer equation $I_o = I_i \cos^2(\phi/2)$, where ϕ is the phase difference between the two light beams, and the transmittance function is defined as $I_o/I_i = \cos^2(\phi/2)$.

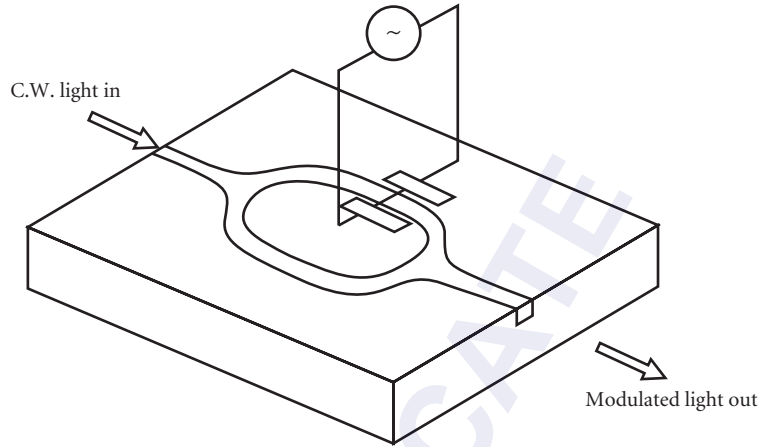


FIGURE 14 Illustration of an integrated lithium niobate Mach-Zehnder modulator.

Owing to the presence of the phase modulator in one of the interferometer arms, and with the phase being controlled by the applied voltage in accordance with a linear relation for the Pockel's effect, for example, $\varphi = \varphi_o - \pi V/V\pi$. In this equation, φ_o is determined by the optical path difference between the two beams and $V\pi$ is the voltage required to achieve a π phase shift between the two beams. The transmittance of the device, therefore becomes a function of the applied voltage V ,

$$T(V) = \cos^2(\varphi/2 - \pi V/2V\pi) \tag{13}$$

This function is plotted in Fig. 15 for an arbitrary value of φ_o . Commercially available integrated devices operate at speeds up to 40 GHz and are quite suitable for OTDM applications such as modulation and demultiplexing.

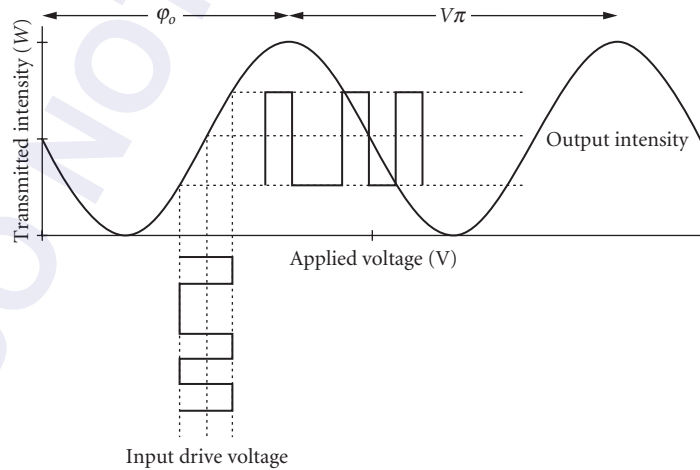


FIGURE 15 Input–output relations of an external modulator based on the Pockel's effect. Superimposed on the transfer function is a modulated drive signal and the resultant output intensity from the modulator.

Electroabsorption Modulators Electroabsorption modulators are intensity modulators that rely on the quantum confined Stark effect. In this device, thin layers of semiconductor material are grown on a semiconductor substrate to generate a multiplicity of semiconductor quantum wells or multiple quantum wells (MQW). For telecommunication applications, the semiconductor material family that is generally used is InGaAsP/InP. The number of quantum wells can vary, but is typically on the order of 10, with an overall device length of a few hundred micrometers. Owing to the dimensions of the thin layers, typically 100 Å or less, the electron and holes bind to form excitons. These excitons have sharp, and well-defined, optical absorption peaks that occur near the bandgap of the semiconductor material. By applying an electric field, or bias voltage, in a direction perpendicular to the quantum well layers, the relative position of the exciton absorption peak can be made to shift to longer wavelengths. As a result, an optical field that passes through these wells can be preferentially absorbed, if the polarization of the light field is parallel to the quantum well layers. Therefore, by modulating the bias voltage across the MQWs, the input light can be modulated. These devices can theoretically possess modulation speeds as high as 50 GHz, with contrasts approaching 50 dB. A typical device schematic and absorption curve is show in Fig. 16a and b.

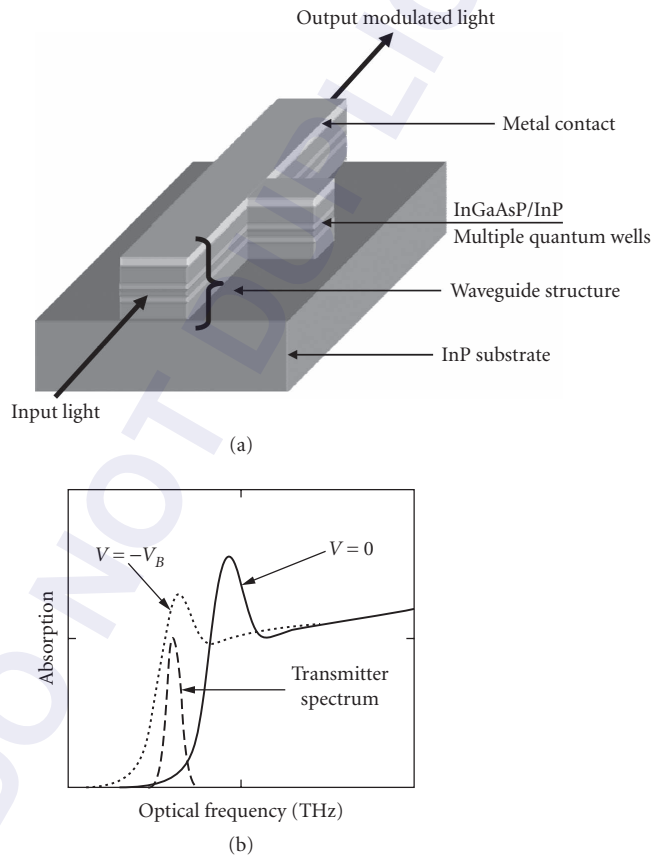


FIGURE 16 (a) Schematic diagram of an electro-absorption modulator. Light propagation occurs along the fabricated waveguide structure, in the plane of the semiconductor multiple quantum wells. (b) Typical absorption spectrum of a multiple quantum well stack under reverse bias and zero bias. Superimposed is a spectrum of a laser transmitter, showing how the shift in the absorption edge can either allow passage or attenuate the transmitted light.

Optical Clock Recovery

In time-division-multiplexed and multiple-access networks, it is necessary to regenerate a timing signal to be used for demultiplexing. Above, a general discussion of clock extraction was given, and in this section, an extension to those concepts are outlined for clock recovery in the optical domain. As in the conventional approaches to clock recovery, optical clock extraction has three general approaches: (1) the optical tank circuit, (2) high-speed phase-locked loops, and (3) injection locking of pulsed optical oscillators. The optical tank circuit can be easily realized by using a simple Fabry-Perot cavity. For clock extraction, the length L of the cavity must be related to the optical transmission bit rate. For example, if the input optical bit rate is 10 Gb/s, the effective length of the optical tank cavity is 15 mm. While the concept of the optical tank circuit is intuitively pleasing, since it has many of the same features as electrical tank circuits, that is, a cavity Q and its associated decay time. In the case of a simple Fabry-Perot cavity as the optical tank circuit, the optical decay time, or photon lifetime, is given by

$$\tau_D = \frac{\tau_{RT}}{1 - R_1 R_2} \quad (14)$$

where τ_{RT} is the round-trip time given as $2L/c$, and R_1 and R_2 are the reflection coefficients of the cavity mirrors. One major difference between the optical tank circuit and its electrical counterpart is that the output of the optical tank circuit never exceeds the input optical intensity (see Fig. 17a).

A second technique that builds on the concept of the optical tank is optical injection seeding or injection locking. In this technique, the optical data bits are injected into a nonlinear device such as a passively mode-locked semiconductor laser diode (see Fig. 17b). The key difference between this approach and the optical tank circuit approach is that the injection-locking technique has internal gain to compensate for the finite photon lifetime, or decay, of the empty cavity. In addition to the gain, the cavity also contains a nonlinear element, for example, a saturable absorber to initiate and sustain pulsed operation. Another important characteristic of the injection-locking technique, using passively mode-locked laser diodes is that clock extraction can be prescaled, that is, a clock signal can

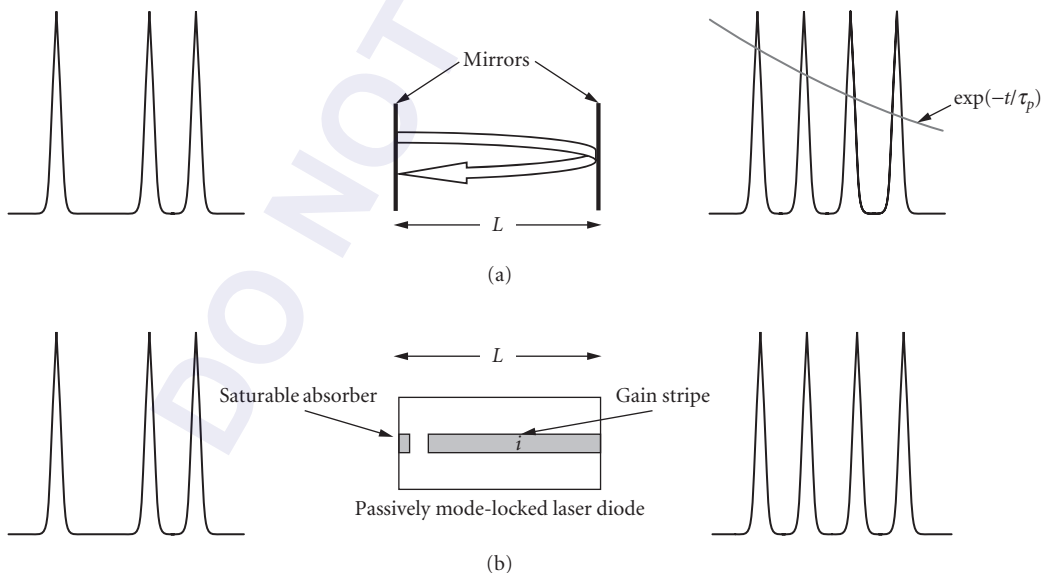


FIGURE 17 All optical clock recovery based on optical injection of (a) an optical tank circuit (Fabry-Perot cavity), and (b) a mode-locked semiconductor diode laser.

be obtained at bit rates exactly equal to the input data bit rate or at harmonics or subharmonics of the input bit rate. In this case of generating a prescaled clock signal at a subharmonic of the input data stream, the resultant signal can directly be used for demultiplexing, without any addition signal processing.

The operation of the injection-seeded optical clock is as follows. The passively mode-locked laser produces optical pulses at its natural rate, which is proportional to the longitudinal mode spacing of the device cavity, $c/(2L)$. Optical data bits from the transmitter are injected into the mode-locked laser, where the data transmission rate is generally a harmonic of the clock rate. This criteria immediately provides the prescaling required for demultiplexing. The injected optical bits serve as a seeding mechanism to allow the clock to build up pulses from the injected optical bits. As the injected optical bits and the internal clock pulse compete for gain, the continuous injection of optical bits force the internal clock pulse to evolve and shift in time to produce pulses that are synchronized with the input data. It should be noted that it is not necessary for the input optical bit rate to be equal to or greater than the nominal pulse rate of the clock, for example, the input data rate can be lower than the nominal bit rate of the clock. This is analogous to the transmitter sending data with primarily 0s, with logic 1 pulses occurring infrequently. The physical operating mechanism can also be understood by examining the operation in the frequency domain. From a frequency domain perspective, since the injected optical data bits are injected at a well-defined bit rate, the optical spectrum has a series of discrete line spectra centered around the laser emission wavelength, and separated in frequency by the bit rate. Since the optical clock emits a periodic train of optical pulses, its optical spectra is also a series of discrete line spectra, separated by the clock-repetition frequency. If the line spectra of the injected data bits fall within optical gain bandwidth of the optical clock, the injected line spectra will serve as seeding signals to force the optical clock to emit line spectra similar to the injected signals. Since the injected data bits are repetitively pulsed, the relative phase relation between the discrete line spectra have the proper phase relation to force the clock to emit synchronously with the injected data.

All-Optical Switching for Demultiplexing

In an all-optical switch, light controls light with the aid of a nonlinear optical material. It should be noted here that all materials will exhibit a nonlinear optical response, but the strength of the response will vary widely depending on the specific material. One important effect to realize an all-optical switch is the optical Kerr effect, where the refractive index of a medium is proportional to the square of the incident electric field. Since light is inducing the nonlinearity, or in other words, providing the incident electric field, the refractive index becomes proportional to the light intensity. Since the intensity of a light beam can change the refractive index, the speed of a second weaker beam can be modified owing to the presence of the intense beam. This effect is used extensively in combination with an optical interferometer to realize all-optical switching (see section on electro-optic modulation using a Mach-Zehnder interferometer). Consider, for example, a Mach-Zehnder interferometer that includes a nonlinear optical material that possess the optical Kerr effect (see Fig. 18). If data to be demultiplexed is injected into the interferometer, the relative phase delay in each arm can be adjusted so that the entire injected data signal is present only at one output port. If an intense optical control beam is injected into the nonlinear optical medium, synchronized with a single data bit passing through the nonlinear medium, that bit can be slowed down, such that destructive interference occurs at the original output port and constructive interference occurs at the secondary output port. In this case, the single bit has been switched out of the interferometer, while all other bits are transmitted.

Optical switches have been realized using optical fiber in the form of a Sagnac interferometer, and the fiber itself is used as the nonlinear medium. These devices are usually referred to as nonlinear optical loop mirrors (NOLM). Other versions of all-optical switches may use semiconductor optical amplifiers as the nonlinear optical element. In this case, it is the change in gain induced by the control pulse that changes the refractive index owing to the Kramers-Kronig relations. Devices such as these are referred to as a terahertz asymmetric optical demultiplexers (TOAD), semiconductor laser amplifier loop optical mirrors (SLALOM), and unbalanced nonlinear interferometers (UNI).

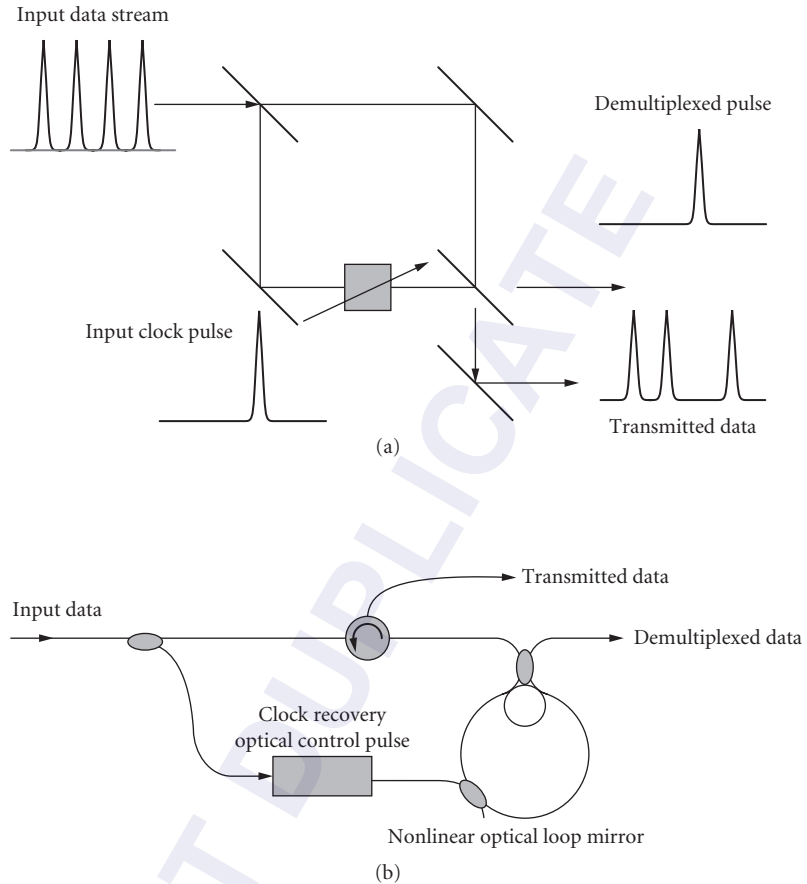


FIGURE 18 Schematic diagram of an all optical switch. (a) Simple configuration based on a Mach-Zehnder interferometer and a separate nonlinear material activated by an independent control pulse. (b) An optical fiber implementation of an all optical switch. This implementation relies on the inherent nonlinearity of the fiber that is induced by an independent control pulse.

Ultrahigh-Speed Optical Time-Division-Multiplexed Optical Link—A Tutorial Example

To demonstrate how the ultrafast device technology can realize state-of-the-art performance in optical time domain multiplexed systems, the following is an example that incorporates the device technology discussed in this chapter. Figure 19 shows a schematic illustration of an ultrahigh-speed OTDM transmission experiment that was used by a collaborative research team of scientists at the Heinrich Hertz Institute in Germany and Fujitsu Laboratories in Japan, to demonstrate 2.56-Tb/s transmission over an optical fiber span of 160 km. The transmitter comprised a 10-GHz mode-locked laser operating at 1550-nm wavelength, generating pulses of 0.42 ps. The pulse train was modulated and temporally multiplexed in two multiplexing stages to realize a maximum data transmission rate of 2.56 Tb/s. This signal is transmitted through 160 km of dispersion-managed optical fiber. At the receiver, the data generates a timing reference signal that drives a mode-locked fiber

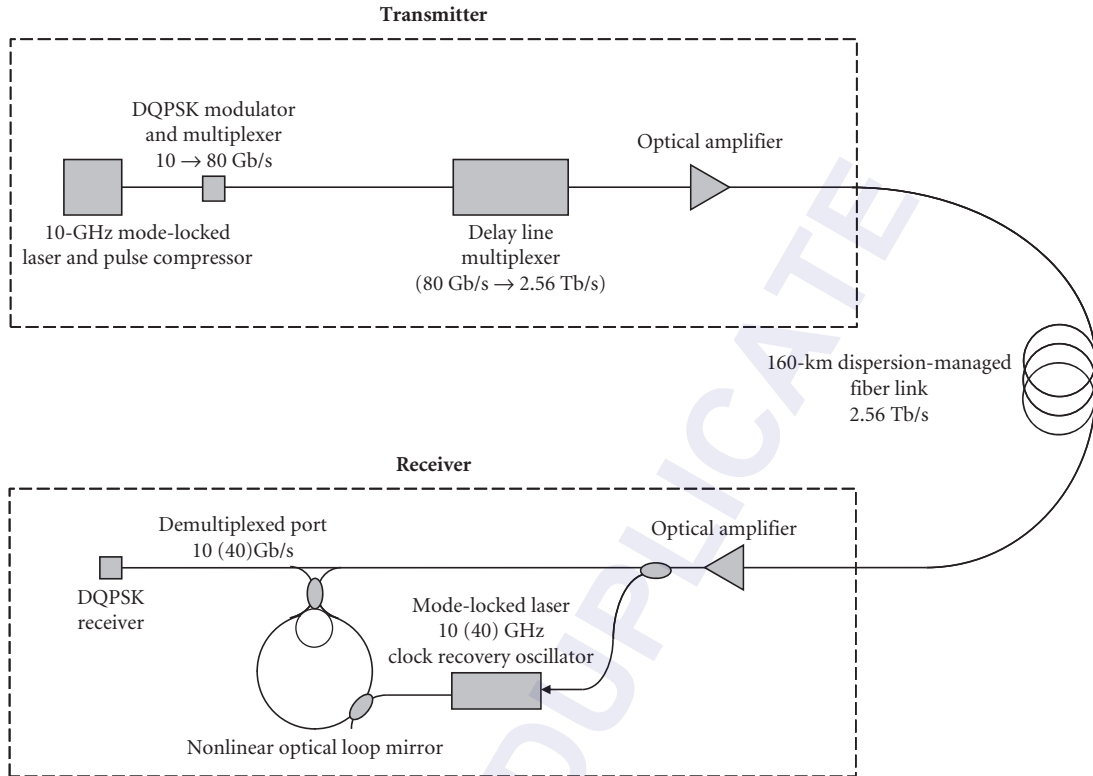


FIGURE 19 Schematic diagram of a 640-Gb/s optical fiber link, using all the critical components described in this chapter, for example, mode-locked laser, high-speed modulator, temporal interleaver, optical clock recovery, all optical switching for demultiplexing and an optical photoreceiver.

laser that is used as an optical gate in a nonlinear optical loop mirror. The resulting demultiplexed data is then detected using a DQPSK photodetection scheme. The robustness of this experiment is quantified by the ratio of the number of errors received to the data pulse period. Error-free operation is defined as less than 10^{-9} . The system performed error free with a received optical power of less than 0.1 mW.

20.5 SUMMARY AND FUTURE OUTLOOK

This chapter reviewed the fundamental basics of optical time-division-multiplexed communication networks, starting from an elementary perspective of digital sampling. Given this as an underlying background, specific device technology was introduced to show how the system functionality can be realized using ultrahigh-speed optics and photonic technologies. Finally, as an example of how these system and device technologies are incorporated into a functioning ultrahigh-speed optical time-division-multiplexed system, a 2.56-Tb/s link was discussed.

In the introduction, the difference between ETDM and OTDM was described. To have an idea of what the future may provide, it should be noted that 100-Gb/s ETDM are now being currently tested in laboratories, while the same data rates were investigated using OTDM techniques nearly

a decade earlier. As electronic technology continues to improve in speed, OTDM techniques are replaced by their electronic counterparts in commercially deployed systems. As a result, we have seen OTDM transmission technology pushing the boundaries of ultrahigh-speed data transmission in optical fibers, thus providing a roadmap for the development of commercial ETDM transmission over optical fiber. As we have shown, OTDM techniques are demonstrating transmission rates in excess of 2.5 Tb/s today, and perhaps suggest similar data rates for ETDM fiber optic transmission in the next decade.

20.6 BIBLIOGRAPHY

- Das, J., "Fundamentals of Digital Communication," in Bishnu P. Pal, (ed.), Section 18 *Fundamentals of Fiber Optics in Telecommunication and Sensor Systems*, Wiley Eastern, New Delhi, 1992, pp. 7-415–7-451.
- Kawanishi, S., "Ultrahigh-Speed Optical Time-Division-Multiplexed Transmission Technology Based on Optical Signal Processing," *J. Quant. Electron. IEEE* **34**(11):2064–2078 (1998).
- Saruwatari, M., "All Optical Time Division Multiplexing Technology," in N. Grote and H. Venghaus, (eds.), *Fiber Optic Communication Devices*, Springer-Verlag, Berlin 2001.
- Schuh, K. and E. Lach, E., "High-Bit-Rate ETDM Transmission Systems," in I. P. Kaminow, T. Li, and A. E. Willner (eds.), *Optical Fiber Telecommunications VB*, Elsevier-Academic Press, New York, 2008, pp. 179–200.
- Weber, H. G., Ferber, S., Kroh, M., Schmidt-Langhorst, C., Ludwig, R., Marembert, V., Boerner, C., Futami, F., Watanabe, S., and C. Schubert, C., "Single Channel 1.28 Tbit/s and 2.56 Tbit/s DQPSK Transmission," *Electron. Lett.* **42**:178–179 (2006).
- Weber, H. G. and Ludwig, R., "Ultra-High-Speed OTDM Transmission Technology," in I. P. Kaminow, T. Li, and A. E. Willner (eds.), *Optical Fiber Telecommunications VB*, Elsevier-Academic Press, New York, 2008, pp. 201–232.
- Weber, H. G. and Nakazawa, N., (eds.), *Optical and Fiber Communication Reports 3*, Springer Science + Business Media, LLC, 2007.

This page intentionally left blank.

DO NOT DUPLICATE

WDM FIBER-OPTIC COMMUNICATION NETWORKS

Alan E. Willner

*University of Southern California
Los Angeles, California*

Changyuan Yu

*National University of Singapore, and
A*STAR Institute for Infocomm Research
Singapore*

Zhongqi Pan

*University of Louisiana at Lafayette
Lafayette, Louisiana*

Yong Xie

*Texas Instruments Inc.
Dallas, Texas*

21.1 INTRODUCTION

The progress in optical communications over the past 30 years has been astounding. It has experienced many revolutionary changes since the days of short-distance multimode transmission at $0.8 \mu\text{m}$.¹ In 1980, AT&T could transmit 672 two-way conversations along a pair of optical fibers.² In 1994, an AT&T network connecting Florida with the Virgin Islands was able to carry 320,000 two-way conversations along two pairs of optical fibers. The major explosion came after the maturity of fiber amplifiers and wavelength-division multiplexing (WDM) technologies. By 2003, the transoceanic system of Tyco Telecommunications is able to transmit 128 wavelengths per fiber pair at 10 Gb/s/wavelength with the total capacity 10 Tb (eight fiber pairs)—a capability of transmitting more than 100 million simultaneous voice circuits on a eight-fiber pair cable. In experiments, a recent notable report demonstrated a record of 25.6-Tb/s WDM transmission using 160 channels within 8000 GHz fiber bandwidth.³

High-speed single-channel systems may offer compact optical systems with a minimized footprint and a maximized cost-per-bit efficiency.⁴ Transmission of single-channel 1.28 Tb/s signal over 70 km was achieved by traditional return-to-zero (RZ) modulation format using optical time-division multiplexing (OTDM) and polarization multiplexing.⁵ Although single-channel results are quite impressive, they have

two disadvantages: (1) they make use of only a very small fraction of the enormous bandwidth available in an optical fiber, and (2) they connect two distinct end points, not allowing for a multiuser environment. Since the required rates of data transmission among many users have been increasing at an impressive pace for the past several years, it is a highly desirable objective to connect many users with a high-bandwidth optical communication system.

WDM-related technologies have been growing rapidly and clearly dominated the research field and the telecommunications market. For instance, the new fiber bands (S and L) are being opened up, new modulation schemes are being deployed, unprecedented bit-rate/wavelength (≥ 40 Gb/s) is being carried over all-optical distance. Over 100 WDM channels are being simultaneously transmitted over ultralong distances.⁶⁻⁹ It is difficult to overstate the impact of fiber amplifiers and WDM technologies in both generating and supporting the telecommunications revolution during last 10 years.

Due to its high capacity and performance, optical fiber communications have already replaced many conventional communication systems in point-to-point transmission and networks. Driven by the rising capacity demand and the need to reduce cost, in addition to the fixed WDM transmission links with point-to-point configuration, the next step of the network migration is to insert the reconfigurable optical add/drop multiplexers (OADMs) in the links and network nodes. An automatically switching network can support the flexibility of the transmission. Hence, the necessary bandwidth can be provided to the customers on demand.^{10,11}

This chapter intends to review the state-of-the-art WDM technologies and reflect the tremendous progress over the past few years.

Fiber Bandwidth

The driving factor for the use of multichannel optical systems is the abundant bandwidth available in the optical fiber. The attenuation curve as a function of optical carrier wavelength is shown in Fig. 1.¹² There are two low-loss windows: one near 1.3 μm and an even lower-loss one near 1.55 μm . Consider the window at 1.55 μm , which is approximately 25,000 GHz wide. (Note that due to the extremely desirable characteristics of the erbium-doped fiber amplifier (EDFA), which amplifies only near 1.55 μm , most systems would use EDFAs and therefore not use the dispersion-zero 1.3- μm band of the existing embedded conventional fiber base.) The high-bandwidth characteristic of the optical fiber implies that a single optical carrier at 1.55 μm can be baseband-modulated at approximately 25,000 Gb/s, occupying 25,000 GHz surrounding 1.55 μm , before transmission losses of the optical fiber would limit transmission. Obviously, this bit rate is impossible for present-day electrical and optical devices to achieve, given that even heroic lasers, external modulators, switches,

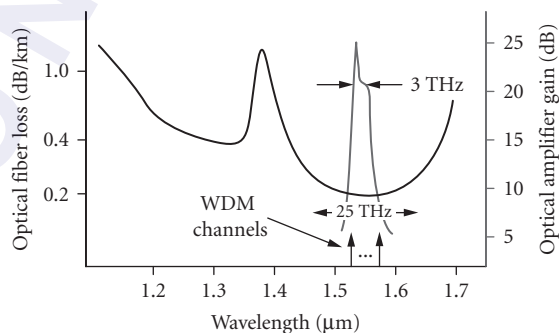


FIGURE 1 Fiber loss as a function of wavelength in conventional single-mode silica fiber. The gain spectrum of the EDFA is also shown.¹²

and detectors all have bandwidths less than or equal to 100 GHz. Practical data links today are significantly slower, perhaps no more than tens of gigabits per second per channel. Since the single high-speed channel makes use of an extremely small portion of the available fiber bandwidth, an efficient multiplexing method is needed to take full advantage of the huge bandwidth offered by optical fibers. As we will see in this chapter, WDM has been proven to be the most appropriate approach.

Introduction to WDM Technology

In real systems, even a single channel will probably be a combination of many lower-speed signals since few individual applications today utilize this high bandwidth. These lower-speed channels are multiplexed together in time to form a higher-speed channel. This time-division multiplexing (TDM) can be accomplished in either the electrical or optical domain. In TDM, each lower-speed channel transmits a bit (or a collection of bits known as a packet) in a given time slot and then waits its turn to transmit another bit (or packet) after all the other channels have had their opportunity to transmit. Until the late 1980s, fiber communication was mainly confined to transmitting a single optical channel using TDM technology. Due to fiber attenuation, this channel required periodic regeneration which included detection, electronic processing, and optical retransmission. Such regeneration causes a high-speed optoelectronic bottleneck, is bit-rate specific, and can only handle a single wavelength. The need for these single-channel regenerators (i.e., repeaters) was replaced when the EDFA was developed, enabling high-speed repeaterless single-channel transmission. We can think of this single approximate gigabits per second channel as a single high-speed “lane” in a highway in which the cars represent packets of optical data and the highway represents the optical fiber. It seems natural to dramatically increase the system capacity by transmitting several different independent wavelengths simultaneously down a fiber in order to more fully utilize this enormous fiber bandwidth.^{13,14} Therefore, the intent was to develop a multiple-lane highway, with each lane representing data traveling on a different wavelength.

In the most basic WDM arrangement as shown in Fig. 2, the desired number of lasers, each emitting a different wavelength, are multiplexed together by a wavelength multiplexer (or a combiner) into the same high-bandwidth fiber.^{13–16} Each of N different-wavelength lasers is operating at the slower gigabits per second speeds, but the aggregate system is transmitting at N times the individual

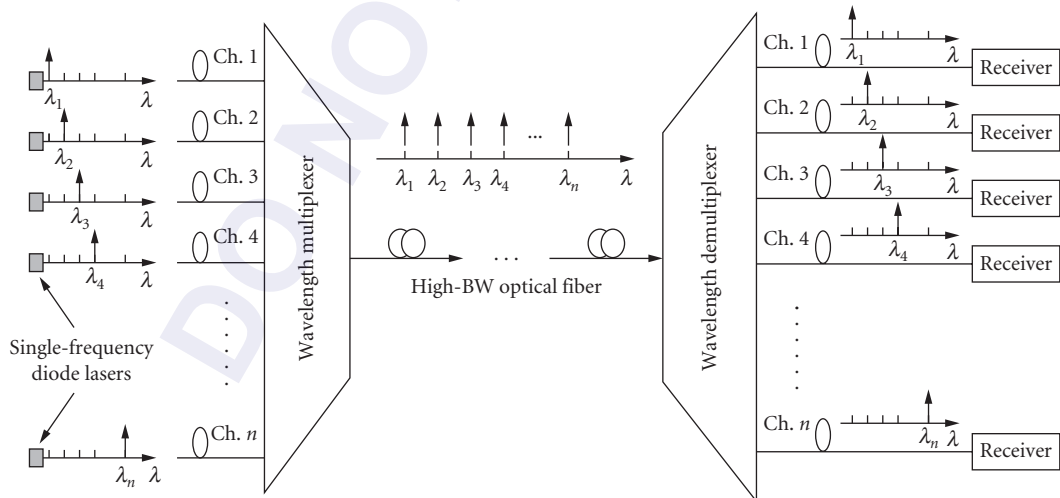


FIGURE 2 Diagram of a simple WDM system.

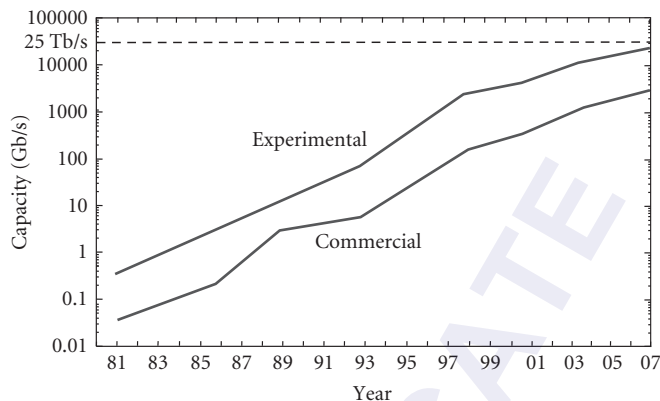


FIGURE 3 Continuous capacity growth in optical fiber transmission systems.

laser speed providing a significant capacity enhancement. After being transmitted through a high-bandwidth optical fiber, the combined optical signals must be demultiplexed by a wavelength demultiplexer at the receiving end by distributing the total optical power to each output port and then requiring that each receiver selectively recovers only one wavelength. Therefore, only one signal is allowed to pass and establish a connection between source and destination. WDM allows us to make use much of the available fiber bandwidth, although various device, system, and network issues will still limit utilization of the full fiber bandwidth.

Figure 3 shows the continuous capacity growth in optical fiber systems over the past 30 years. The highest capacity has been achieved using WDM. One interesting point about this trend predicted two decades ago by T. Li of AT&T Bell Labs is that the transmission capacity doubles every 2 years. WDM technology has provided the platform for this trend to continue, and there is no reason to assume that WDM won't continue to produce dramatic progress.

21.2 BASIC ARCHITECTURE OF WDM NETWORKS

We have explained how WDM enables the utilization of a significant portion of the available fiber bandwidth by allowing many independent signals to be transmitted simultaneously in one fiber. In fact, WDM technology also enables wavelength routing and switching of data paths in an optical network. By utilizing wavelength-selective components, each data channel's wavelength can be routed and detected independently through the network. The wavelength determines the communication path by acting as the signature address of the origin, destination, or routing. Data can then be presumed not traveling on optical fiber but on wavelength-specific "light-paths" from source to destination that can be arranged by a network controller to optimize throughput. Therefore, the basic system architecture that can take the full advantage of WDM technology is an important issue, and will be discussed in this section.

Point-to-Point Links

As shown in Fig. 2, in a simple point-to-point WDM system, several channels are multiplexed at one node, the combined signals are then transmitted across some distance of fiber, and the channels are demultiplexed at a destination node. This point-to-point WDM link facilitates the high-bandwidth fiber transmission without routing or switching in the optical data path.

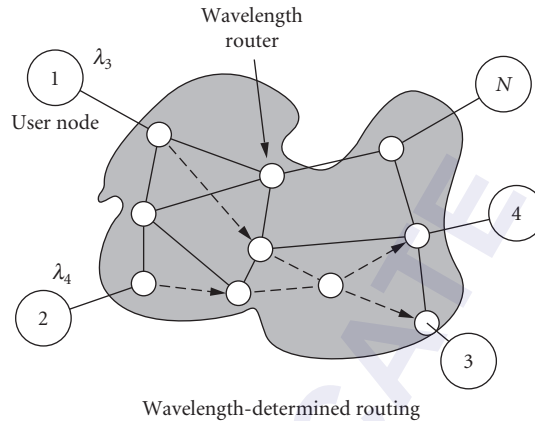


FIGURE 4 A generic multiuser network in which the communications links and routing paths are determined by the wavelengths used within the optical switching fabric.

Wavelength-Routed Networks

Figure 4 shows a more complex multiuser WDM network structure, where the wavelength is used as the signature address for either the transmitters or the receivers, and determines the routing path through an optical network. In order for each node to be able to communicate with any other node and facilitate proper link setup, the transmitters or the receivers must be wavelength tunable; we have arbitrarily chosen the transmitters to be tunable in this network example. Note that the wavelengths are routed passively in wavelength-routed networks.

WDM Stars, Rings, and Meshes

Three common WDM network topologies are star, ring, and mesh networks.¹⁷⁻¹⁹ In the star topology, each node has a transmitter and receiver, with the transmitter connected to one of the passive central star's inputs and the receiver connected to one of the star's outputs, as is shown in Fig. 5a. Rings, as shown in Fig. 5b, are also popular because: (1) many electrical networks use this topology, and (2) rings are easy to implement for any geographical network configuration. In this example,

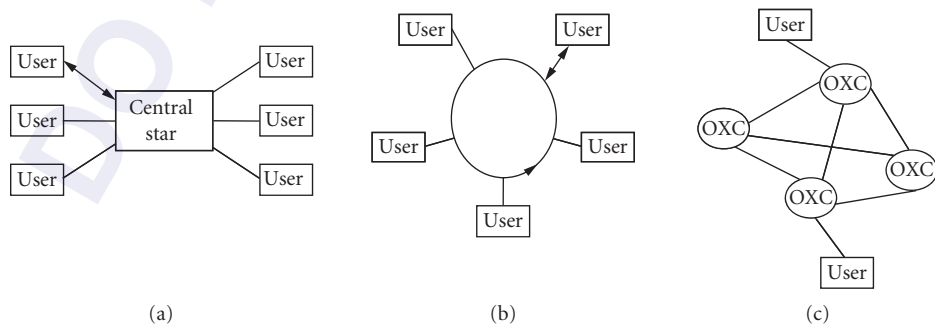


FIGURE 5 WDM: (a) stars; (b) rings; and (c) meshes.

each node in the unidirectional ring can transmit on a specific signature wavelength, and each node can recover any other node's wavelength signal by means of a wavelength-tunable receiver. Although not depicted in the figure, each node must recover a specific channel. This can be performed: (1) where a small portion of the combined traffic is tapped off by a passive optical coupler, thereby allowing a tunable filter to recover a specific channel, or (2) in which a channel-dropping filter completely removes only the desired signal and allows all other channels to continue propagating around the ring. Furthermore, a synchronous optical network (SONET) dual-ring architecture, with one ring providing service and the other protection, can provide automatic fault detection and protection switching.²⁰

In both the star and ring topologies, each node has a designated wavelength, and any two nodes can communicate with each other by transmitting and recovering that wavelength. This implies that N wavelengths are required to connect N nodes. The obvious advantage of this configuration, known as a single-hop network, is that data transfer occurs with an uninterrupted optical path between the origin and destination; the optical data starts at the originating node and reaches the destination node without stopping at any other intermediate node. A disadvantage of this single-hop WDM network is that the network and all its components must accommodate N wavelengths, which may be difficult (or impossible) to achieve in a large network, that is, the present fabrication technology cannot provide and the transmission capability cannot accommodate 1000 distinct wavelengths for a 1000-user network.

It is important to address that the reliability is a problem in fiber ring. If a station is disabled or if a fiber breaks, the whole network goes down. To address this problem, a double-ring optical network, also called a "self-healing" ring, is used to bypass the defective stations and loops back around a fiber break, as shown in Fig. 6. Each station has two inputs and two outputs connected to two rings that operate in opposite directions.

An alternative to require N wavelengths to accommodate N nodes is to have a multihop network (mesh network) in which two nodes can communicate with each other by sending data through a third node, with many such intermediate hops possible, shown in Fig. 5c. In the mesh network, the nodes are connected by reconfigurable optical crossconnects (OXC).²¹ The wavelength can be dynamically switched and routed by controlling the OXCs. Therefore, the required number of wavelengths and the tunable range of the components can be reduced in this topology. Moreover, the mesh topology can also provide multiple paths between two nodes to make the network protection and restoration easier to realize. If a failure occurs in one of the paths, the system can automatically find another path and restore communications between any two nodes. However, OXCs with large numbers of ports are extremely difficult to obtain, which limits the scalability of the mesh network.

In addition, there exist several other network topologies, such as tree network, which is widely used in broadcasting or distributing systems. At the "base" of the tree is the source transmitter from which emanates the signal to be broadcast throughout the network. From this base, the tree splits

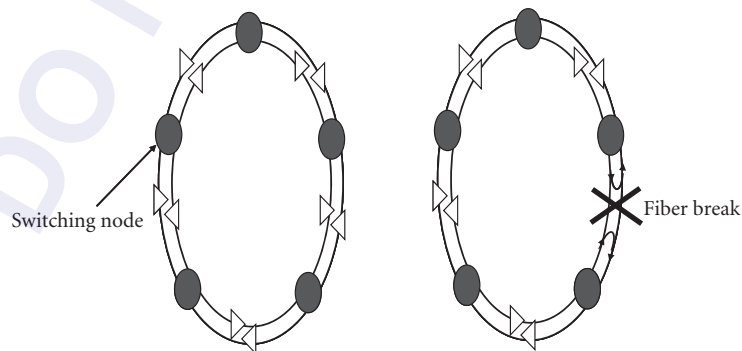


FIGURE 6 A self-healing ring network.

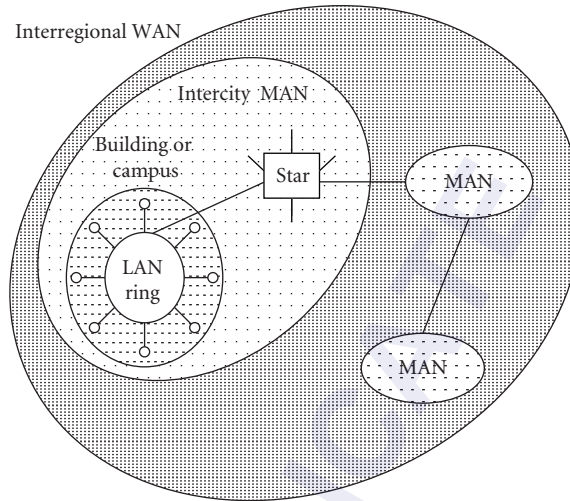


FIGURE 7 Hybrid network topologies and architectures woven together to form a large network.

many times into different “branches,” with each branch either having nodes connected to it or further dividing into subbranches. This continues until all the nodes in the network can access the base transmitter. Whereas the other topologies are intended to support bidirectional communication among the nodes, this topology is useful for distributing information unidirectionally from a central point to a multitude of users. This is a very straightforward topology and is in use in many systems, most notably cable television (CATV).

By introducing Fig. 7 in which a larger network is composed of smaller ones, we have also introduced the subject of the architecture of the network which depends on the network’s geographical extent. The three main architectural types are the local-, metropolitan-, and wide-area networks, denoted by LAN, MAN, and WAN, respectively.²² Although no rule exists, the generally accepted understanding is that a LAN interconnects a small number of users covering a few kilometers (i.e., intra- and interbuilding), a MAN interconnects users inside a city and its outlying regions, and a WAN interconnects significant portions of a country (100 of kilometers). Based on Fig. 7, the smaller networks represent LANs, the larger ones MANs, and the entire figure would represent a WAN. In other words, a WAN is composed of smaller MANs, and a MAN is composed of smaller LANs. Hybrid systems exist, and typically a wide-area network will consist of smaller local-area networks, with mixing and matching between the most practical topologies for a given system. For example, stars and rings may be desirable for LANs whereas buses may be the only practical solution for WANs. It is, at present, unclear which network topology and architecture will ultimately and most effectively take advantage of high-capacity optical systems.

Circuit and Packet Switching

The two fundamental types of underlying telecommunication network infrastructures, based on how traffic is multiplexed and switched, are *circuit-switched* and *packet switched*. A circuit-switched network provides circuit-switched connections to its customers. Once the connection is established, a guaranteed amount of bandwidth is allocated to each connection and become available to the connection anytime. The network is also transparent and the nodes seem to be directly connected. In addition, circuit switching requires a lower switching speed (more than milliseconds). Many types of communication links and distribution systems may satisfactorily be interconnected by

circuit switching which is relatively simple to operate. The problems with circuit switching include (1) it is not efficient at handling bursty data traffic (low utilization for traffic with changing intensity or short lived connections); (2) there is a delay before the connection can be used; (3) the resources are permanently allocated to a connection and cannot be used for any other users, and (4) circuit-switched networks are more sensitive to faults (e.g., if a part of the connection fails, the whole transfer fails).

Since optical circuit switching is the relatively mature technology today, the current deployed WDM wavelength-routing network is generally circuit-switched. The basic mechanism of communication in a wavelength-routed network is a lightpath (corresponding to a circuit), which is an all optical connection (communication channel) linking multiple optical segments from a source node to a destination node over a wavelength on each intermediate link. At each intermediate node of the network, the lightpath is routed and switched from one link to another link. A lightpath can use either the same wavelength throughout the whole link or a concatenation of different wavelengths after undergoing wavelength conversion at intermediate optical nodes. In the absence of any wavelength conversion device, a lightpath is required to be on the same wavelength channel throughout its path in the network; this requirement is referred to as the wavelength continuity property of the lightpath. Once the setup of a lightpath is completed, the whole lightpath is available during the connection. Note that different lightpaths can use the same wavelength as long as they do not share any common links (i.e., same wavelength can be reused spatially in different parts of the network).

As shown in Fig. 8, the key network elements in the wavelength-routing network are optical line terminal, OADM (see Fig. 9), and OXC (see Fig. 10). There are many different OADM structures such as parallel or serial, fixed or reconfigurable. In general, an ideal OADM would add/drop any channel and any number of channels, and would be remotely controlled and reconfigured without disturbance to unaffected channels. There will be more discussion on reconfigurable in “Network Reconfigurability” section. The other requirements for an OADM include low and fixed loss, independent of set of wavelengths dropped, and low cost. An OXC can switch the channels from input to output ports and input to output wavelengths. The functions in an OXC node include providing lightpath, rerouting (switching protection), restoring failed lightpath, monitoring performances, accessing to test signals, wavelength conversion, and multiplexing and grooming. An OXC can be either

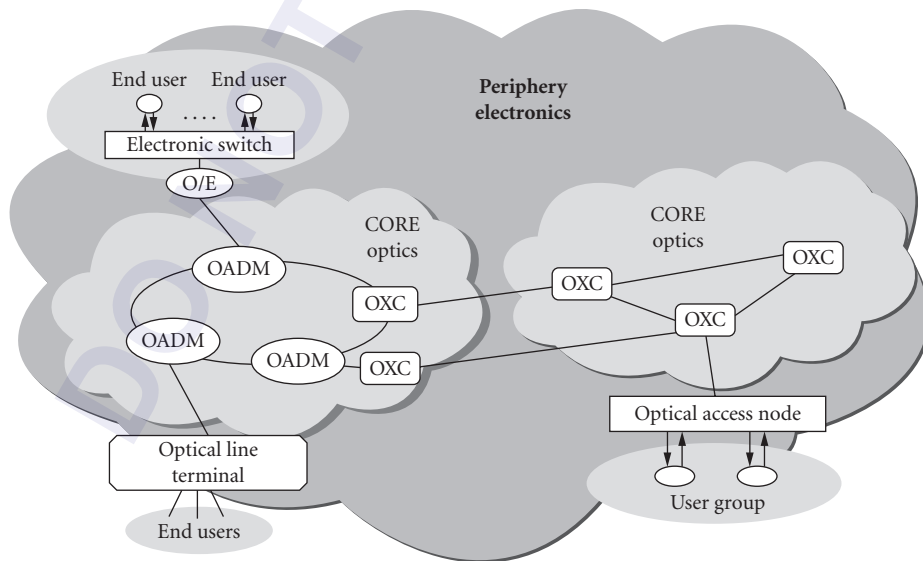


FIGURE 8 An optical network showing optical line terminal, OADM, and OXC.

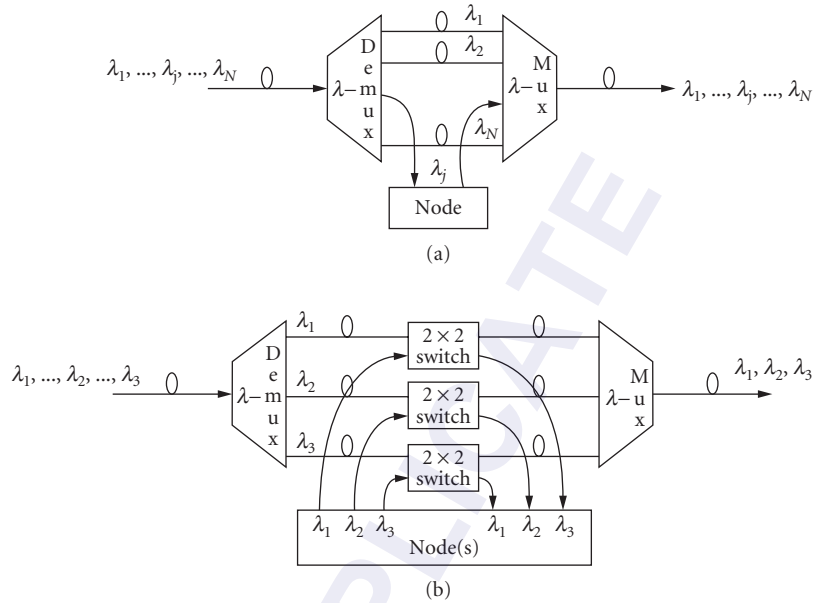


FIGURE 9 Optical add-drop multiplexing (OADM) systems: (a) fixed versus (b) reconfigurable.

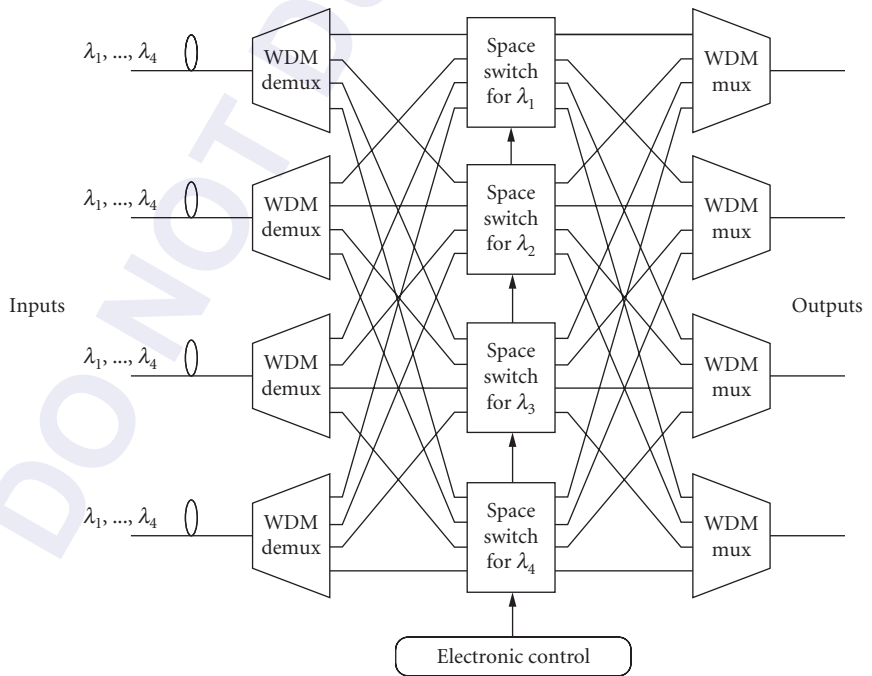


FIGURE 10 An optical crossconnect system in reconfigurable optical networks.

electrical (performing O-E-O conversion for each WDM channels) or optical, transparent or opaque. To accomplish all these functions, the OXC needs three building blocks: (1) fiber switching, to route all of the wavelengths on an incoming fiber to a different outgoing fiber (optical space switches); (2) wavelength switching, to switch specific wavelengths from an incoming fiber to multiple outgoing fibers (multiplexing/demultiplexing); and (3) wavelength conversion, to take incoming wavelengths and convert them to another optical frequency on the outgoing port. This is essential to achieve strictly nonblocking architectures when using wavelength switching.

In packet switched networks, the data stream is broken into small packets. These data packets are multiplexed together with packets from other data stream inside the network. The packets are switched inside the network based on their destination. To facilitate this switching, a packet header is added to the payload in each packet. The header carries addressing information. The switching nodes read the header and then determine where to switch the packet. At the receiver end, packets belonging to a particular stream are put back together. Therefore, packet switching requires switching speeds of microseconds or less. For high-speed optical transmission, packet switching holds the promise for more efficient data transfer, for which no long-distance handshaking is required. The high-bandwidth links are used more efficiently.

As shown in Fig. 11, an optical packet switching node is generally composed of three parts: control unit, switching unit, and input/output interfaces. The control unit retains information about network topology, the forwarding table, scheduling, and buffering. It decides the switching time and is in charge of resolving contentions at a node. The switching unit allows the data to remain in the optical domain during the routing process. It is especially important for optical packet-switched networks that the switching speed be fast enough to minimize overhead. The input/output interface is where optical technologies are utilized to deal with contention problems. Optical buffers and wavelength converters are the building blocks of time and wavelength domain contention resolution modules, respectively, and are housed in the interface units. In addition, other physical layer functionalities required for an optical switching node such as synchronization are realized at the interface.

Note that network packet switching can be accomplished in a conceptually straightforward manner by requiring a node to optoelectronically detect and retransmit each and every incoming optical data packet. The control and routing information is contained in the newly detected electronic packet, and all the switching functions can occur in the electrical domain prior to optical retransmission of the signal. Unfortunately, this approach implies that an optoelectronic speed bottleneck will eventually occur in this system. On the other hand, it is extremely difficult to accomplish the signal processing in the optical domain currently. Alternatively, much research is focused toward maintaining an all-optical data path and performing the switching functions all-optically with only some electronic control of the optical components. The reason is that the control unit detects and processes only the

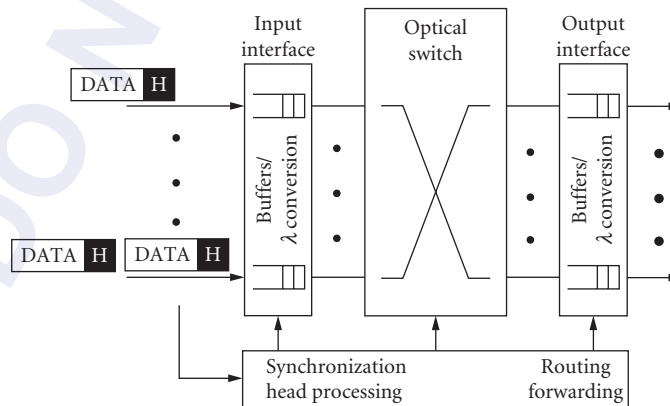


FIGURE 11 Schematic diagram of an optical packet-switching node.

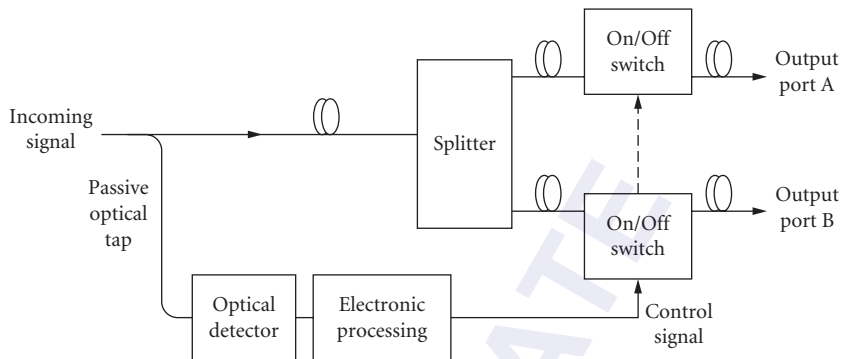


FIGURE 12 Passive optical tapping of an optical packet in order to determine routing information and allow a node to electronically control an optical switch.

header of a packet (not the payload). Therefore it is possible to transmit the header at a lower bit rate to facilitate the processing. As a result, the presence of electronics at the control unit does not necessarily pose a limitation on the data transmission rate. Figure 12 shows a generic solution which passively taps an incoming optical signal. Information about the signal is made known electrically to the node but the signal itself remains in the optical domain. The routing information may be contained in the packet header or in some other form (i.e., wavelength, etc.). Header information can be transmitted in several different ways. For example, the baseband header can be transmitted as a data field inside a packet either at the same bit rate as the data, or at a lower rate to relax the speed requirement on the electronics for header detection. It can also be located out-of-band either on a subcarrier on the same wavelength or on a different wavelength, altogether.

The various functions that may be performed in an optical switching node include (1) address/label recognition to determine the intended output port, (2) header updating or label swapping to prepare the packet header for the next node, (3) bit and/or packet synchronization to the local node to time the switching process, (4) routing-table caching as a reference for routing decisions, (5) output-port contention resolution via buffering and/or wavelength conversion, (6) signal monitoring to assess the signal quality, (7) signal regeneration to combat the accumulated distortion of the signal, and (8) optical switching to direct packets to the appropriate output ports.

It is very important to mention that in the past several years, data traffic has been growing at a much faster rate than voice traffic. Data traffic is “bursty” in nature, and reserving bandwidth for bursty traffic could be very inefficient. Circuit-switching networks are not optimized for this type of traffic. On the other hand, packet-switching networks require a complex processing unit and fast optical switches. Therefore, optical burst switching (OBS) was introduced to reduce the processing required for switching at each node and to avoid optical buffering. OBS is somewhere between packet switching and circuit switching with the switching period on the order of many packets. OBS takes advantage of time domain statistical multiplexing to utilize the large bandwidth of a single channel to transmit several lower-bandwidth bursty channels. At the network edge, packets are aggregated to generate bursts that are sent over the network core. In almost all OBS schemes, the header is sent separately from the payload with an offset time that is dependent on the scheme.^{23,24} The header is sent to the switches, and the path is then reserved for the payload that follows. The loose coupling between the control packet and the burst alleviates the need for optical buffering and relaxes the requirement for switching speed.

Many challenging optical switching issues require solutions,²⁵ such as routing control,^{26–28} contention resolution,^{29–40} and optical header processing.^{41–48} Most of these issues will relate to packet switching, although circuit switching will also require attention. Due to the limited space, we give only some references to facilitate further study by the reader.

Network Reconfigurability

In addition to high capacity, the reconfigurable WDM network could offer flexibility and availability. A reconfigurable network is highly desirable to meet the requirements of high bandwidth and bursty traffic in future networks. Through the reconfigurable network, service providers and network operators could respond quickly and cost-effectively to new revenue opportunities. As shown in Fig. 9, a fixed add/drop multiplexing node can only process the signal(s) at a given wavelength or a group of wavelengths. While in the dynamic reconfigurable node, operators could add or drop any number of wavelengths. This added flexibility would save operating and maintenance costs and improve network efficiency. In general, a network is reconfigurable if it can provide the following functionality for multichannel operations: (1) channel add/drop and (2) path reconfiguration for bandwidth allocation or restoration. It appears that a reconfigurable network is highly desirable to meet the requirements of high bandwidth and bursty traffic in future networks.

Since a reconfigurable network allows dynamic network optimization to accommodate changing traffic patterns, it provides more efficient use of network resources. Figure 13 shows blocking probability as a function of call arrival rate in a WDM ring network with 20 nodes.⁴⁹ A configurable topology can support 6 times the traffic of a fixed WDM topology for the same blocking probability.

Among many different solutions, the reconfigurable optical add/drop multiplexers (ROADMs) have emerged as key building blocks for the next-generation WDM systems with the goal of any wavelength anywhere. ROADMs add the ability to remotely switch traffic at the wavelength layer in WDM network, thus allowing individual data channels to be added or dropped without optical-electrical-optical (O-E-O) conversion of all WDM channels.^{50–55} Figure 14 shows a few ROADM architectures based on different switching technologies: discrete switches or switch matrix plus filters [mux/switch/demux, variable optical attenuators (VOA)]; wavelength blockers (WB); integrated planar lightwave circuits (iPLC); and wavelength-selective switches (WSS). With the exception of mux-switch-demux design, the devices are typically implemented in broadcast-and-select optical architectures with passive splitters in the pass-through path. A relevant attribute of ROADM technology is the integration of multiplexing/demultiplexing and switching into a single component. This integration can significantly lower pass-through losses when compared with multiple discrete components.⁵³

ROADM networks can deliver considerable operation benefits such as simplified planning and engineering, and improved network utilization. Its applications include optical multicast/broadcast, scalable colorless add/drop capacity, capacity/service upgrade without traffic disruption, cost-effective optical protection and restoration, and low cost of ownership.^{50,53}

The key component technologies enabling network reconfigurability include wavelength-tunable lasers and laser arrays, wavelength routers, optical switches, OXCs, OADMs, optical amplifiers, and tunable optical filters, and the like. Although huge benefits are possible with a reconfigurable topology,

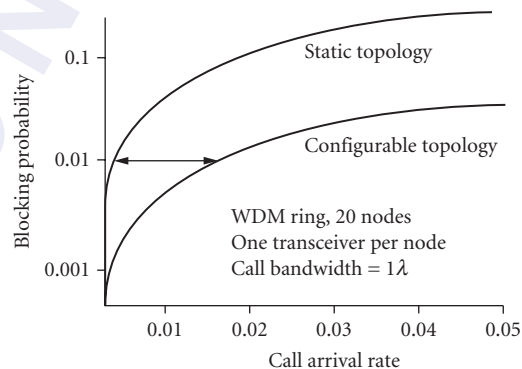


FIGURE 13 Blocking probability as a function of call arrival rate in a WDM ring.⁴⁹

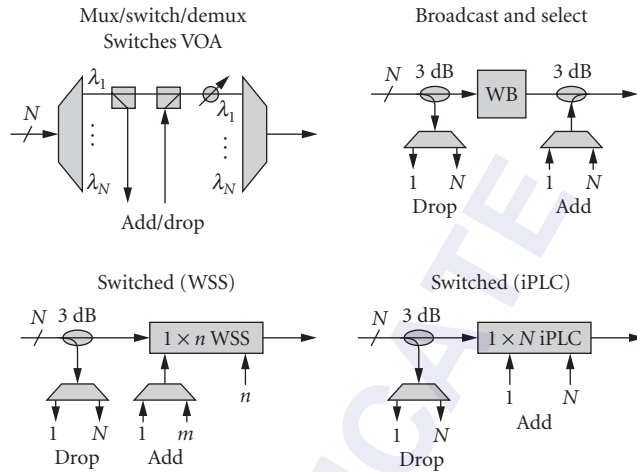


FIGURE 14 ROADMs architectures overview.⁵³

the path to reconfigurability is paved with various degrading effects. As shown in Fig. 4, the signal may pass through different lengths of fiber links due to the dynamic routing, causing some degrading effects in reconfigurable networks to be more critical than in static networks, such as nonstatic dispersion and nonlinearity accumulation due to reconfigurable paths, EDFA gain transients, channel power nonuniformity, cross-talk in optical switching and cross-connects, and wavelength drift of components. We will discuss some of these important effects in Sec. 21.3, followed by selected advanced technologies to dealing with these effects in WDM systems.

21.3 FIBER SYSTEM IMPAIRMENTS

One key benefit of reconfigurable WDM networks might be the transparency to bit rate, protocol and modulation format of all the various wavelength channels propagating in the system. However, key challenges exist when determining an optimum path through the network, since an optical wavelength might accumulate different physical impairments as it is switched through the network. These nonidealities will be imposed by both the transmission links and the optical switching nodes. Since system performance depends on many different optical impairments, a network-layer routing and wavelength assignment algorithm might rapidly provision a lightpath that cannot meet the signal-quality requirement.^{56–59} In this section, we will give a brief review of different physical-layer impairments, including fiber attenuation and power loss, fiber chromatic dispersion, fiber polarization mode dispersion, and fiber nonlinear effects. The management of both fiber dispersion and nonlinearities will also be discussed at the end of this section. Note that EDFA-related impairments, such as noise, fast power transients, and gain peaking in EDFA cascades will be described in Sec. 21.5.

Fiber Attenuation and Optical Power Loss

The most basic characteristic of a link is the power loss, which is caused by fiber attenuation and connections.⁶⁰ *Attenuation*, defined as the ratio of the input power to the output power, is the loss of optical power as light travels along the fiber. Attenuation in an optical fiber is caused by absorption, scattering, and bending losses. The fundamental physical limits imposed on the fiber attenuation are

due to scattering off the silica atoms at shorter wavelengths and the material absorption at longer wavelengths. There are two minima in the loss curve: one near $1.3\ \mu\text{m}$ and an even lower one near $1.55\ \mu\text{m}$ (see Fig. 1). Fiber bending can also induce power loss because radiation escapes through its bends. The bending loss is inversely proportional to the bend radius and is wavelength dependent.

Power loss is also present at fiber connections, such as connectors, splices, and couplers. Coupling of light into and out of a small-core fiber is much more difficult to achieve than coupling electrical signals in copper wires since: (1) photons are weakly confined to the waveguide whereas electrons are tightly bound to the wire, and (2) the core of a fiber is typically much smaller than the core of an electrical wire. First, light must be coupled into the fiber from a diverging laser beam, and two fibers must be connected to each other. Second, connecting two different fibers in a system must be performed with great care due to the small size of the cores. One wishes to achieve connections exhibiting: (1) low loss, (2) low back reflection, (3) repeatability, and (4) reliability. Two popular methods of connecting fibers are the permanent splice and the mechanical connector. The permanent “fusion” splice can be accomplished by placing two fiber ends near each other, generating a high-voltage electric arc which melts the fiber ends, and “fusing” the fibers together. Losses and back reflections tend to be extremely low, being less than $0.1\ \text{dB}$ and less than $-60\ \text{dB}$, respectively. Disadvantages of these fusion splices include (1) the splice is delicate and must be protected, and (2) the splice is permanent. Alternatively, there are several types of mechanical connectors, such as ST and FC/PC. Losses and back reflections are still fairly good, and are typically less than $0.3\ \text{dB}$ and less than $-45\ \text{dB}$, respectively.

Low loss is extremely important since a light pulse must contain a certain minimum amount of power in order to be detected, such that “0” or “1” data bit can be unambiguously detected. If not for dispersion, we would clearly prefer to operate with $1.55\ \mu\text{m}$ light due to its lower loss for long-distance systems.

Chromatic Dispersion

In fact, in any medium (other than vacuum) and in any waveguide structure (other than ideal infinite free space), different electromagnetic frequencies travel at different speeds. This is the essence of chromatic dispersion. As the real fiber-optic world is rather distant from the ideal concepts of both vacuum and infinite free space, dispersion will always be a concern when one is dealing with the propagation of electromagnetic radiation through fiber. The velocity in fiber of a single monochromatic wavelength is constant. However, data modulation causes a broadening of the spectrum of even the most monochromatic laser pulse. Thus, all modulated data has a nonzero spectral width which spans several wavelengths, and the different spectral components of modulated data travel at different speeds. In particular, for digital data intensity modulated on an optical carrier, chromatic dispersion leads to pulse broadening—which in turn leads to chromatic dispersion limiting the maximum data rate that can be transmitted through optical fiber (see Fig. 15).

Considering that the chromatic dispersion in optical fibers is due to the frequency-dependent nature of the propagation characteristics for both the material and the waveguide structure, the speed of light of a particular wavelength λ can be expressed as follows using a Taylor series expansion of the value of the refractive index as a function of the wavelength.

$$v(\lambda) = \frac{c_o}{n(\lambda)} = \frac{c_o}{n_o(\lambda_o) + \frac{\partial n}{\partial \lambda} \delta\lambda + \frac{\partial^2 n}{\partial \lambda^2} (\delta\lambda)^2} \quad (1)$$

where c_o is the speed of light in vacuum, λ_o is a reference wavelength, and the terms in $\partial n/\partial \lambda$ and $\partial^2 n/\partial \lambda^2$ are associated with the chromatic dispersion and the dispersion slope (i.e., the variation of the chromatic dispersion with wavelength), respectively. Transmission fiber has positive dispersion, that is, longer wavelengths see longer propagation delays. The units of chromatic dispersion are picoseconds per nanometer per kilometer, meaning that shorter time pulses, wider frequency spread due to data modulation, and longer fiber lengths will each contribute linearly to temporal dispersion. Figure 16 shows the dispersion coefficient, D (ps/nm·km), of a conventional single-mode fiber with the material and waveguide contributions plotted separately.⁶⁰ For a given system, a pulse will disperse

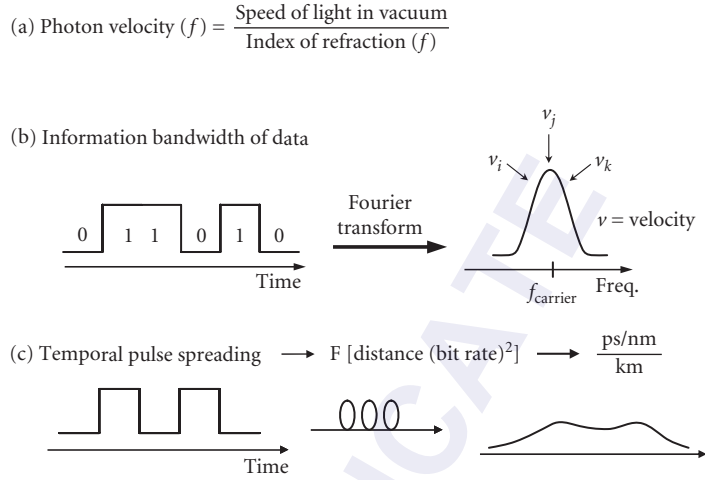


FIGURE 15 The origin of chromatic dispersion in data transmission. (a) Chromatic dispersion is caused by the frequency-dependent refractive index in fiber. (b) The nonzero spectral width due to data modulation. (c) Dispersion leads to pulse broadening, proportional to the transmission distance and the data rate.

more in time for a wider frequency distribution of the light and for a longer length of fiber. Higher data rates inherently have both shorter pulses and wider frequency spreads. Therefore, as network speed increases, the impact of chromatic dispersion rises precipitously as the square of the increase in data rate. The quadratic increase with the data rate is a result of two effects, each with a linear contribution. On one hand, a doubling of the data rate makes the spectrum twice as wide, doubling the effect of dispersion. On the other hand, the same doubling of the data rate makes the data pulses only half as long (hence twice as sensitive to dispersion). The combination of a wider signal spectrum and a shorter pulse width is what leads to the overall quadratic impact—when the bit rate increases by a factor of 4, the effects of chromatic dispersion increase by a whopping factor of 16!⁶¹

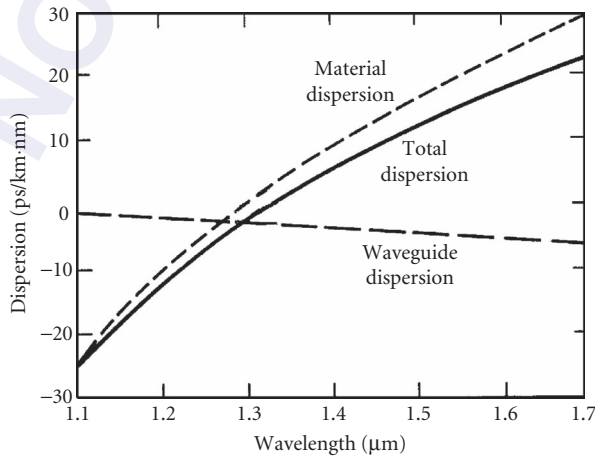


FIGURE 16 Dispersion coefficient, D , as a function of wavelength in the conventional silica single-mode fiber.⁶⁰

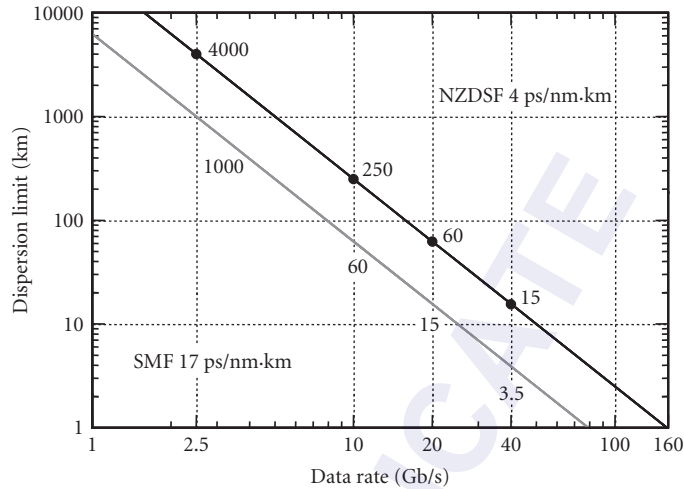


FIGURE 17 Transmission distance limitations due to uncompensated dispersion in SMF as a function of data rate for intensity modulated optical signals.⁶²

The data rate and the data-modulation format can significantly affect the sensitivity of a system to chromatic dispersion. For example, the common non-return-to-zero (NRZ) data format, in which the optical power stays high throughout the entire time slot of a “1” bit, is more robust to chromatic dispersion than is the return-to-zero (RZ) format, in which the optical power stays high in only part of the time slot of a “1” bit. This difference is due to the fact that RZ data has a much wider channel frequency spectrum compared to NRZ data, thus incurring more chromatic dispersion. However, in a real WDM system, the RZ format increases the maximum allowable transmission distance by virtue of its reduced duty cycle (compared to the NRZ format) making it less susceptible to fiber nonlinearities. We will discuss some robust modulation formats in Sec. 21.4 too.

A rule for the maximum distance over which data can be transmitted is to consider a broadening of the pulse equal to the bit period. For a bit period B , a dispersion value D and a spectral width $\Delta\lambda$, the dispersion-limited distance is given by

$$L_D = \frac{1}{D \cdot B \cdot \Delta\lambda} = \frac{1}{D \cdot B \cdot (cB)} \propto \frac{1}{B^2} \quad (2)$$

(see Fig. 17). For example, for single mode fiber, $D = 17$ ps/nm·km, so for 10 Gb/s data the distance is L_D equal to 52 km. In fact, a more exact calculation shows that for 60 km, the dispersion induced power penalty is less than 1 dB.⁶² The power penalty for uncompensated dispersion rises exponentially with transmission distance, and thus to maintain good signal quality, dispersion compensation is required.

Polarization-Mode Dispersion

Single-mode fibers actually support two perpendicular polarizations of the original transmitted signal (fundamental mode). In an ideal fiber (perfect) these two modes are indistinguishable, and have the same propagation constants owing to the cylindrical symmetry of the waveguide. However, the core of an optical fiber may not be perfectly circular, and the resultant ellipse has two orthogonal axes. The index-of-refraction of a waveguide, which determines the speed of light, depends on the shape of the waveguide as well as the glass material itself. Therefore, light polarized along one fiber axis travels at a

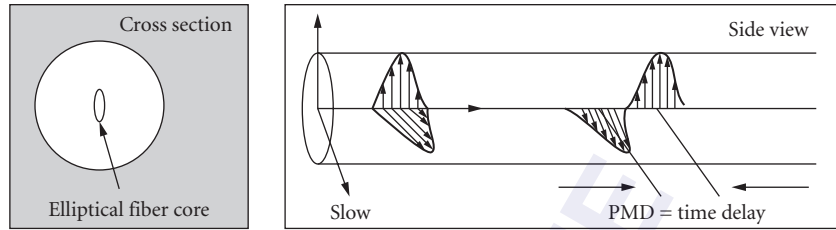


FIGURE 18 Illustration of polarization mode dispersion caused by imperfect round fiber core. An input optical pulse has its power transmitted on two orthogonal polarization modes, each arriving at different times.

different speed as does light polarized along the orthogonal fiber axis (see Fig. 18). This phenomenon is called polarization mode dispersion (PMD). Fiber asymmetry may be inherent in the fiber from the manufacturing process, or it may be a result of mechanical stress on the deployed fiber. The inherent asymmetries of the fiber are fairly constant over time, while the mechanical stress due to movement of the fiber can vary, resulting in a dynamic aspect to PMD. Since the light in the two orthogonal axes travel with different group velocities, to the first order, this differential light speed will cause a temporal spreading of signals, which is termed the differential group delay (DGD).

Because of random variations in the perturbations along a fiber span, PMD in long fiber spans accumulates in a random-walk-like process that leads to a square root of transmission-length dependence.⁶³ Moreover, PMD does not have a single value for a given span of fiber. Rather, it is described in terms of average DGD, and a fiber has a distribution of DGD values over time. The probability of the DGD of a fiber section being a certain value at any particular time follows a maxwellian distribution (see Fig. 19). The probability of $DGD = \Delta \tau$ given by

$$\text{prob}(\Delta \tau) = \sqrt{\frac{2}{\pi}} \frac{\Delta \tau^2}{\alpha^3} \exp\left(-\frac{\Delta \tau^2}{2\alpha^2}\right) \tag{3}$$

with mean value $\langle \Delta \tau \rangle = \sqrt{8/\pi} \alpha$. PMD is usually expressed in $\text{ps}/\text{km}^{1/2}$ in long fiber spans, and the typical PMD parameter (D_p) is 0.1 to $10 \text{ ps}/\text{km}^{1/2}$.^{64,65}

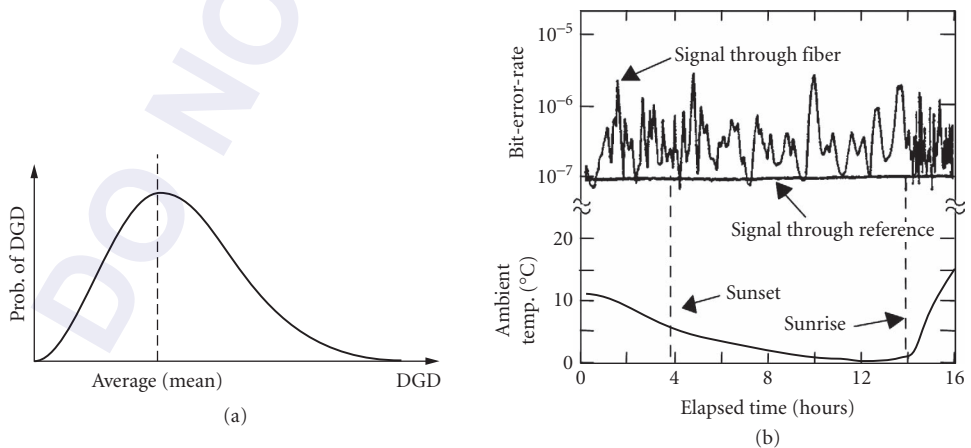


FIGURE 19 (a) Probability distribution of DGD in a typical fiber. (b) System performance (bit-error rate) fluctuations due to changes in temperature caused by PMD.⁶⁶

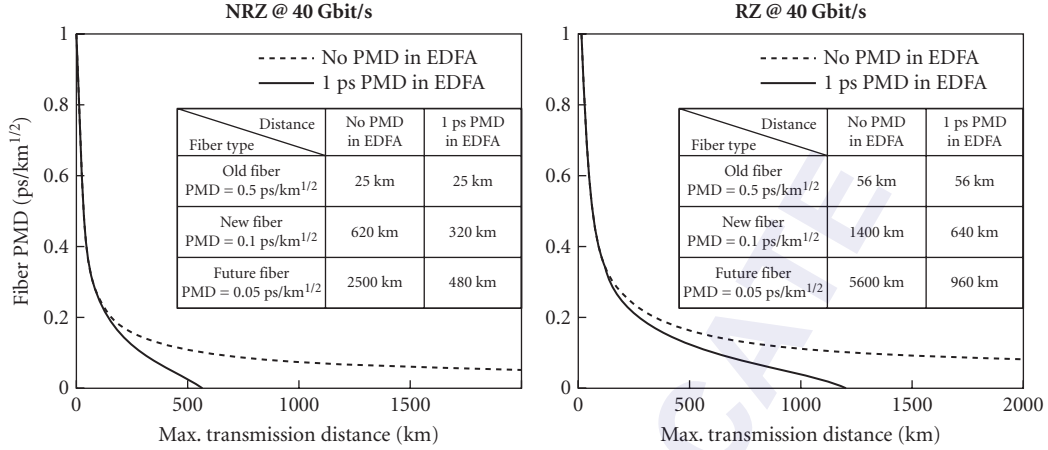


FIGURE 20 Limitations of transmission distances caused by fiber PMD.

Today's fiber has a very low PMD value and is well characterized. But there is still a small residual asymmetry in the fiber core. Moreover, slight polarization dependencies exist in discrete inline components such as isolators, couplers, filters, erbium-doped fiber, modulators, and multiplexers. Therefore, even under the best of circumstances, PMD still significantly limit the deployment of more than or equal to 40 Gb/s systems (see Fig. 20).⁶⁷

Other polarization-related impairments, such as polarization-dependent loss (PDL), and polarization-dependent gain (PDG), may also cause deleterious effects in a fiber transmission link.⁶⁸ Moreover, the interaction between PMD and PDL/PDG may lead to significant overall performance degradation, which dramatically surpasses the result of adding the degradations induced by the two impairments independently.^{69–73} The readers can find more in-depth discussion in the related literatures.

Fiber Nonlinearities

Most nonlinear effects originate from the nonlinear refractive index of fiber. The refractive index not only depends on the frequency of light but also on the intensity (optical power), and it is related to the optical power as⁷⁴

$$\tilde{n}(\omega, P) = n_o(\omega) + n_2 I = n_o(\omega) + n_2 \frac{P}{A_{\text{eff}}} \quad (4)$$

where $n_o(\omega)$ is the linear refractive index of silica, n_2 is the intensity-dependent refractive index coefficient, P is the optical power inside the fiber, and A_{eff} is the effective mode area of the fiber. The typical value of n_2 is $2.6 \times 10^{-20} \text{ m}^2/\text{W}$. This number takes into account the averaging of the polarization states of the light as it travels in the fiber. The intensity dependence of the refractive index gives rise to three major nonlinear effects.

Self-Phase Modulation A million photons “see” a different glass than does a single photon, and a photon traveling along with many other photons will slow down. Self-phase modulation (SPM) occurs because of the varying intensity profile of an optical pulse on a single WDM channel (see Fig. 21a). This intensity profile causes a refractive index profile and, thus, a photon speed differential. The resulting phase change for light propagating in an optical fiber is expressed as

$$\Phi_{\text{NL}} = \gamma P L_{\text{eff}} \quad (5)$$

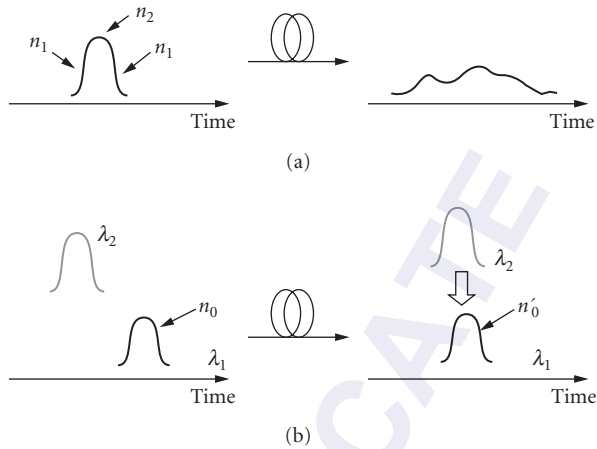


FIGURE 21 (a) Self-phase modulation: the photons in the pulse “see” different refractive index and (b) cross-phase modulation: the glass that a photon in the λ_2 pulse “sees” changes as other channels (with potentially varying power) move to coincide with the λ_2 pulse.

where the quantities γ and L_{eff} are defined as

$$\gamma = \frac{2\pi n_2}{\lambda A_{\text{eff}}} \quad \text{and} \quad L_{\text{eff}} = \frac{1 - e^{-\alpha L}}{\alpha} \quad (6)$$

where α is the fiber attenuation loss, L_{eff} is the effective nonlinear length of the fiber that accounts for fiber loss, and γ is the nonlinear coefficient measured in radians per kilometer per watt. A typical range of values for γ is about 10^{-30} rad/km · W. Although the nonlinear coefficient is small, the long transmission lengths and high optical powers that have been made possible by the use of optical amplifiers can cause a large enough nonlinear phase change to play a significant role in state-of-the-art lightwave systems.

Cross-Phase Modulation When considering many WDM channels copropagating in a fiber, photons from channels 2 through N can distort the index profile that is experienced by channel 1. The photons from the other channels “chirp” the signal frequencies on channel 1, which will interact with fiber chromatic dispersion and cause temporal distortion (see Fig. 21b). This effect is called cross-phase modulation (XPM). In a two-channel system, the frequency chirp in channel 1 due to power fluctuation within both channels is given by

$$\Delta B = \frac{d\Phi_{\text{NL}}}{dt} = \gamma L_{\text{eff}} \frac{dP_1}{dt} + 2\gamma L_{\text{eff}} \frac{dP_2}{dt} \quad (7)$$

where, dP_1/dt and dP_2/dt are the time derivatives of the pulse powers of channels 1 and 2, respectively. The first term on right hand side of the above equation is due to SPM, and the second term is due to XPM. Note that the XPM-induced chirp term is double that of the SPM-induced chirp term. As such, XPM can impose a much greater limitation on WDM systems than can SPM, especially in systems with many WDM channels.

Four-Wave-Mixing The optical intensity propagating through the fiber is related to the electric field intensity squared. In a WDM system, the total electric field is the sum of the electric fields of each individual channel. When squaring the sum of different fields, products emerge that are beat terms

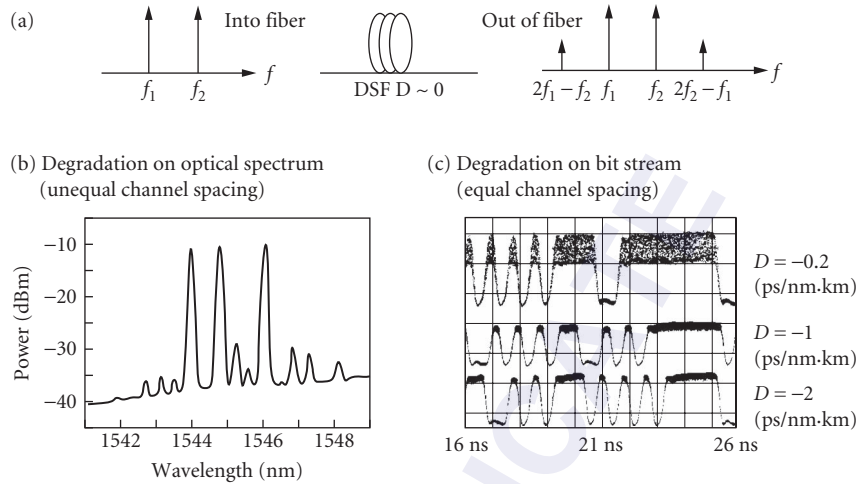


FIGURE 22 (a) and (b) FWM induces new spectral components via nonlinear mixing of two wavelength signals. (c) The signal degradation due to FWM products falling on a third data channel can be reduced by even small amounts of dispersion.⁷⁵

at various sum and difference frequencies to the original signals. Figure 22 depicts that if a WDM channel exists at one of the four-wave-mixing (FWM) beat-term frequencies, the beat term will interfere coherently with this WDM channel and potentially destroy the data.

Other Nonlinear Effects The nonlinear effects described above are governed by the power dependence of refractive index, and are elastic in the sense that no energy is exchanged between the electromagnetic field and the dielectric medium. A second class of nonlinear effects results from stimulated inelastic scattering in which the optical field transfers part of its energy to the nonlinear medium. Two important nonlinear effects fall in this category:⁷⁴ (1) stimulated Raman scattering (SRS) and (2) stimulated Brillouin scattering (SBS). The main difference between the two is that optical phonons participate in SRS, while acoustic phonons participate in SBS. In a simple quantum-mechanical picture applicable to both SRS and SBS, a photon of the incident field is annihilated to create a photon at a downshifted frequency. The downshifted frequency range where new photons can be generated is approximately 30 THz in SRS and only approximately 30 MHz in SBS.

The fiber nonlinearities, including SPM, XPM, FWM as well as stimulated scattering, will start to degrade the optical signals when the optical power in fiber becomes high. An important parameter when setting up spans in optical systems is the launch power to the fiber. The power must be large enough to provide an acceptable optical signal-to-noise ratio (OSNR) at the output of the span but below the limit where excited fiber nonlinearities distort the signal. The specific limit depends on several different factors such as the type of fiber used, the bit rate, amplifier spacing, and the applied dispersion map. In dense WDM systems, the trade-off relationship between OSNR degradation by accumulation of amplified spontaneous emission (ASE) noise from optical amplifiers and nonlinear waveform-distortion in transmission fibers determines the optimum transmission power and together they limit the regenerative repeater spacing.⁷⁶

Dispersion and Nonlinearities Management

In this section, we will address the concepts of chromatic dispersion and fiber nonlinearities management followed by some examples highlighting the need for tunability to enable robust optical WDM systems in dynamic environments.

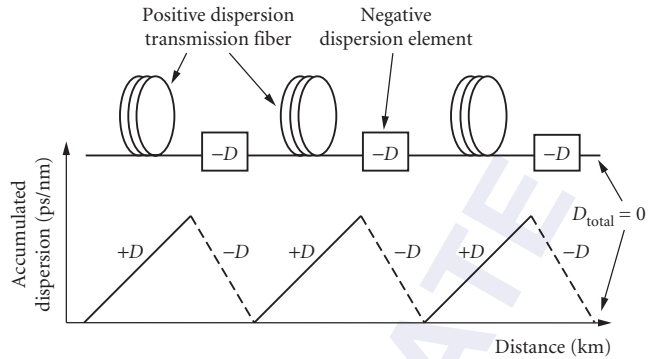


FIGURE 23 Dispersion map of a basic dispersion-managed system. Positive dispersion transmission fiber alternates with negative dispersion compensation elements such that the total dispersion is zero end-to-end.

In the preceding section, we saw that although chromatic dispersion is generally considered a negative characteristic, it is not always bad for fiber transmission. It is, in fact, a necessary evil for the deployment of WDM systems. When the fiber dispersion is near zero in a WDM system, different channels travel at almost the same speed. Any nonlinear effects that require phase matching between the different wavelength channels will accumulate at a higher rate than if wavelengths travel at widely different speeds (the case of higher dispersion fiber). Therefore, it may not be a good idea to reduce the fiber dispersion to zero by using dispersion-shifted fiber, which has both the dispersion zero and the loss minimum located at 1.55 μm . As an alternative, we keep the local dispersion along the transmission link high enough to suppress nonlinear effects, while managing the total dispersion of the link to be close to zero, as shown in Fig. 23. This is a very powerful concept: at each point along the fiber the dispersion has some nonzero value, eliminating FWM and XPM, but the total dispersion at the end of the fiber link is zero, so that no pulse broadening is induced. The most advanced systems require periodic dispersion compensation, as well as pre- and postcompensation (before and after the transmission fiber).

The addition of negative dispersion to a standard fiber link has been traditionally known as “dispersion compensation,” however, the term “dispersion management” is more appropriate. Standard single-mode fiber (SMF) has positive dispersion, but some new varieties of nonzero dispersion-shifted fiber (NZDSF) come in both positive and negative dispersion varieties, as shown in Fig. 24. Reverse dispersion

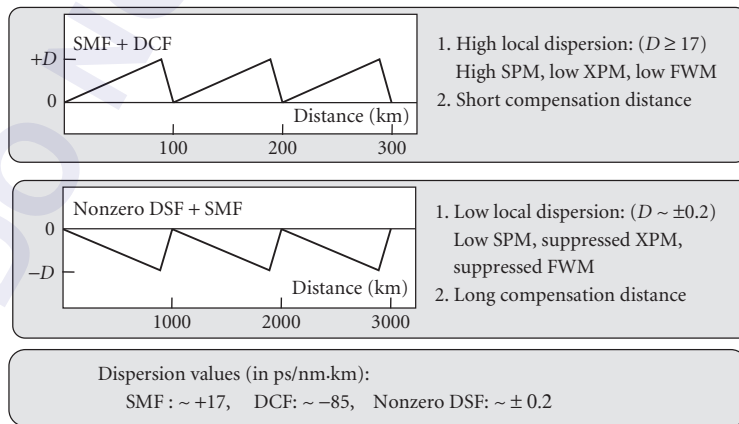


FIGURE 24 Various dispersion maps for SMF-DCF and NZDSF-SMF.

fiber is also now available, with a large dispersion comparable to that of SMF, but with the opposite sign. When such flexibility is available in choosing both the magnitude and sign of the dispersion of the fiber in a link, dispersion-managed systems can fully be optimized to the desired dispersion map using a combination of fiber and dispersion compensation devices (see Fig. 24). Dispersion is a linear process, so the first-order-dispersion maps can be understood as linear systems. However, the effects of nonlinearities cannot be ignored, especially in WDM systems with many tens of channels where the launch power may be very high. In particular, in systems deploying dispersion compensating fiber (DCF), the large nonlinear coefficient of the DCF can dramatically affect the dispersion map. We will review and highlight a few different dispersion management solutions in the following sections.

Fixed Dispersion Compensation From a systems point of view, there are several requirements for a dispersion compensating module: low loss, low optical nonlinearity, broadband (or multichannel) operation, small footprint, low weight, low power consumption, and clearly low cost. It is unfortunate that the first dispersion compensation modules, based on DCF only met two of these requirements: broadband operation and low power consumption. On the other hand, several solutions have emerged that can complement or even replace these first-generation compensators.

Dispersion compensating fiber One of the first dispersion compensation techniques was to deploy specially designed sections of fiber with negative chromatic dispersion. The technology for DCF emerged in the 1980s and has developed dramatically since the advent of optical amplifiers in 1990. DCF is the most widely deployed dispersion compensator, providing broadband operation and stable dispersion characteristics, and the lack of a dynamic, tunable DCF solution has not reduced its popularity.⁷⁷ In general, the core of the average DCF is much smaller than that of standard SMF, and beams with longer wavelengths experience relatively large changes in mode size (due to the waveguide structure) leading to greater propagation through the cladding of the fiber, where the speed of light is greater than that of the core. This leads to a large negative dispersion value. Additional cladding layers can lead to improved DCF designs that can include negative dispersion slope to counteract the positive dispersion slope of standard SMF.

In spite of its many advantages, DCF has a number of drawbacks. First of all, it is limited to a fixed compensation value. In addition, DCF has a weakly guiding structure and has a much smaller core cross-section, approximately $19 \mu\text{m}^2$, compared to the $85 \mu\text{m}^2$ (approximately) of SMF. This leads to higher nonlinearity, higher splice losses, as well as higher bending losses. Secondly, the length of DCF required to compensate for SMF dispersion is rather long, about one-fifth of the length of the transmission fiber for which it is compensating. Thus DCF modules induce loss, and are relatively bulky and heavy. The bulk is partly due to the mass of fiber, but also due to the resin used to hold the fiber securely in place. Another contribution to the size of the module is the higher bending loss associated with the refractive index profile of DCF; this limits the radius of the DCF loop to 6 to 8 inches, compared to the minimum bend radius of 2 inches for SMF.

Traditionally, DCF-based dispersion compensation modules are usually located at amplifier sites. This serves several purposes. First, amplifier sites offer relatively easy access to the fiber, without requiring any digging or unbraiding of the cable. Second, DCF has high loss (usually at least double that of standard SMF), so a gain stage is required before the DCF module to avoid excessively low signal levels. DCF has a cross section 4 times smaller than SMF, hence a higher nonlinearity, which limits the maximum launch power into a DCF module. The compromise is to place the DCF in the midsection of a two-section EDFA. This way, the first stage provides pre-DCF gain, but not to a power level that would generate excessive nonlinear effects in the DCF. The second stage amplifies the dispersion compensated signal to a power level suitable for transmission through the fiber link. This launch power level is typically much higher than the one that could be transmitted through DCF without generating large nonlinear effects. Many newer dispersion compensation devices have better performance than DCF, in particular lower loss and lower nonlinearities. For this reason, they may not have to be deployed at the midsection of an amplifier.

Chirped fiber bragg gratings Fiber Bragg gratings (FBGs) have emerged as major components for dispersion compensation because of their low loss, small footprint, and low optical nonlinearities.⁷⁸

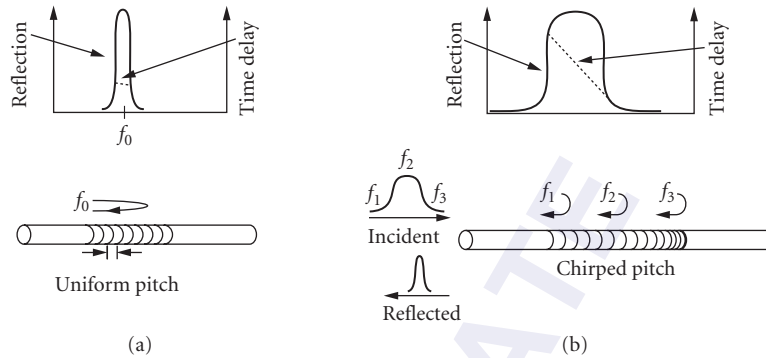


FIGURE 25 Uniform and chirped FBGs: (a) a grating with uniform pitch has a narrow reflection spectrum and a flat time delay as a function of wavelength and (b) a chirped FBG has a wider bandwidth, a varying time delay, and a longer grating length. Chirped gratings reflect different frequency components at different locations within the grating.

When the periodicity of the grating is varied along its length, the result is a chirped grating which can be used to compensate for chromatic dispersion. The chirp is understood as the rate of change of the spatial frequency as a function of position along the grating. In chirped gratings the Bragg matching condition for different wavelengths occurs at different positions along the grating length. Thus, the roundtrip delay of each wavelength can be tailored by designing the chirp profile appropriately. Figure 25 compares the chirped FBG with uniform FBG. In a data pulse that has been distorted by dispersion, different frequency components arrive with different amounts of relative delay. By tailoring the chirp profile such that the frequency components see a relative delay which is the inverse of the delay of the transmission fiber, the pulse can be compressed back. The dispersion of the grating is the slope of the time delay as a function of wavelength, which is related to the chirp. An optical circulator is traditionally used to separate the reflected output beam from the input beam.

The main drawback of Bragg gratings is that the amplitude profile and the phase profile as a function of wavelength have some amount of ripple. Ideally, the amplitude profile of the grating should have a flat (or rounded) top in the passband, and the phase profile should be linear (for linearly chirped gratings) or polynomial (for nonlinearly chirped gratings). The grating ripple is the deviation from the ideal profile shape. Considerable effort has been expended on reducing the ripple. While early gratings were plagued by more than 100 ps of ripple, published results have shown vast improvement to values close to ± 3 ps.

Ultimately, dispersion compensators should accommodate multichannel operation. Several WDM channels can be accommodated by a single chirped FBG in one of two ways: fabricating a much longer (i.e., meters-length) grating, or using a sampling function when writing the grating, thereby creating many replicas of transfer function of the FBG in the wavelength domain.⁷⁹

Tunable Dispersion Compensation

The need for tunability In a perfect world, all fiber links would have a known, discrete, and unchanging value of chromatic dispersion. Network operators would then deploy fixed dispersion compensators periodically along every fiber link to exactly match the fiber dispersion. Unfortunately, several vexing issues may necessitate that dispersion compensators have tunability, that is, they have the ability to adjust the amount of dispersion to match system requirements.

First, there is the most basic business issue of inventory management. Network operators typically do not know the exact length of a deployed fiber link nor its chromatic dispersion value. Moreover, fiber plants periodically undergo upgrades and maintenance, leaving new and nonexact lengths of fiber behind. Therefore, operators would need to keep in stock a large number

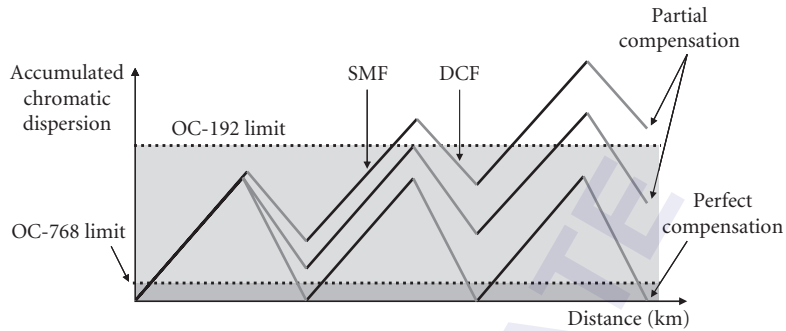


FIGURE 26 The need for tunability. The tolerance of OC-768 systems to chromatic dispersion is 16 times lower than that of OC-192 systems. Approximate compensation by fixed in-line dispersion compensators for a single channel may lead to rapid accumulation of unacceptable levels of residual chromatic dispersion.

of different compensator models, and even then the compensation would only be approximate. Second, we must consider the sheer difficulty of 40 Gb/s signals. The tolerable threshold for accumulated dispersion for a 40 Gb/s data channel is 16 times smaller than that at 10 Gb/s. If the compensation value does not exactly match the fiber to within a few percent of the required dispersion value, then the communication link will not work. Tunability is considered a key enabler for this bit rate (see Fig. 26). Third, the accumulated dispersion changes slightly with temperature, which begins to be an issue for 40 Gb/s systems and 10 Gb/s ultralong haul systems. In fiber, the zero-dispersion wavelength changes with temperature at a typical rate of 0.03 nm/°C. It can be shown that a not-uncommon 50°C variation along a 1000-km 40-Gb/s link can produce significant degradation (see Fig. 27). Fourth, we are experiencing the dawn of reconfigurable optical networking.

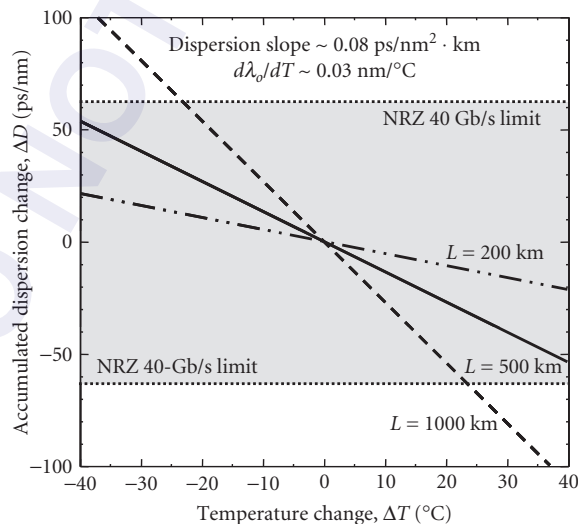


FIGURE 27 Accumulated dispersion changes as a function of the link length and temperature fluctuation along the fiber link.

In such systems, the network path, and therefore the accumulated fiber dispersion, can change. It is important to note that even if the fiber spans are compensated for span-by-span, the pervasive use of compensation at the transmitter and receiver suggests that optimization and tunability based on path will still be needed.

Other issues that increase the need for tunability include (1) laser and (de)mux wavelength drifts for which a data channel no longer resides on the flat-top portion of a filter, thereby producing a chirp on the signal that interacts with the fiber's chromatic dispersion, (2) changes in signal power that change both the link's nonlinearity and the optimal system dispersion map, and (3) small differences that exist in transmitter-induced signal chirp.

Approaches to tunable dispersion compensation A host of techniques for tunable dispersion compensation have been proposed in recent years. Some of these ideas are just interesting research ideas, but several have strong potential to become viable technologies. We will discuss FBG-based technology as an example.

If a FBG has a refractive-index periodicity that varies nonlinearly along the length of the fiber, it will produce a time delay that also varies nonlinearly with wavelength (see Fig. 28). Herein lays the key to tunability. When a linearly chirped grating is stretched uniformly by a single mechanical element, the time delay curve is shifted toward longer wavelengths, but the slope of the ps versus nm curve remains constant at all wavelengths within the passband. When a nonlinearly chirped grating is stretched, the time delay curve is shifted toward longer wavelengths, but the slope of the ps versus nm curve at a specific channel wavelength changes continuously.⁸⁰

Another solution was also reported, which is based on differential heating of the substrate. The thermal gradient induced a chirp gradient, which could be altered electrically⁸¹ and has a major advantage: no moving parts. However, this is countered by the disadvantage of slow tuning, limited to seconds or minutes. Additionally, the technology requires accurate deposition of a thin film of tapered thickness. The process of deposition of the tapered film seems to have some yield issues,

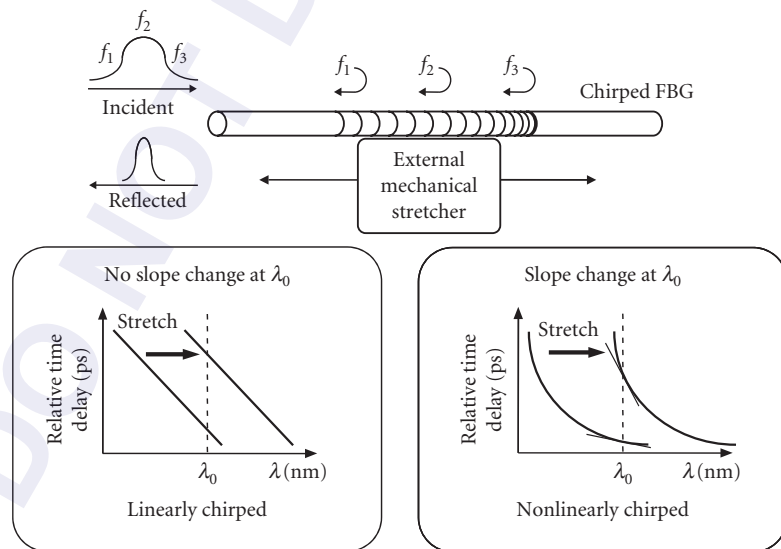


FIGURE 28 Tuning results for both linearly and nonlinearly chirped FBGs using uniform stretching elements. The slope of the dispersion curve at a given wavelength λ_0 is constant when the linearly chirped grating is stretched, but changes as the nonlinearly chirped grating is stretched.

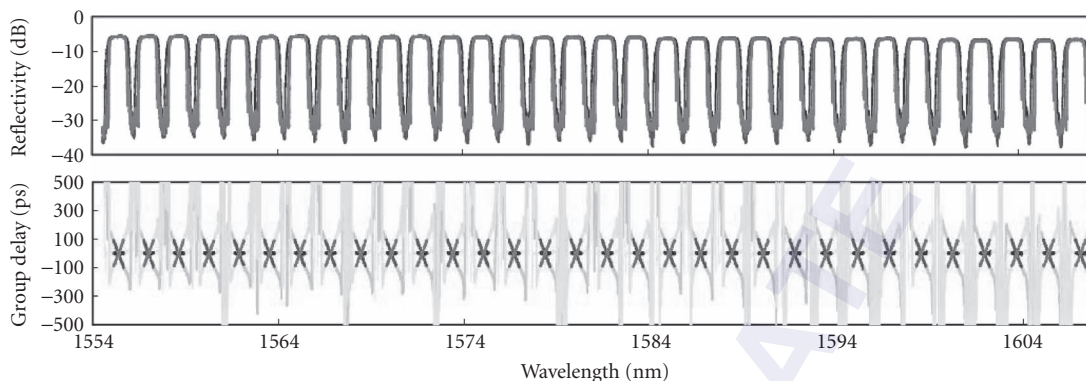


FIGURE 29 Thirty-two-channel, 100-GHz channel spacing, FBG-based tunable dispersion compensator made by Teraxion Inc.⁸²

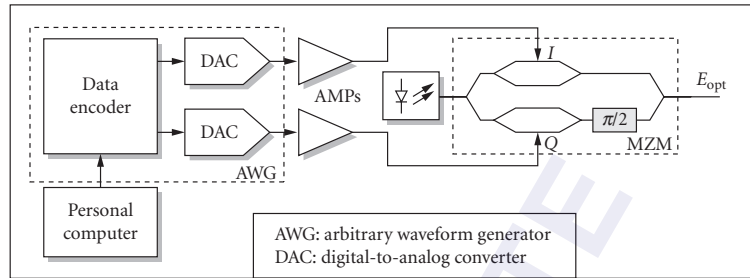
making it rather difficult to manufacture. A 32-channel, 100-GHz channel spacing, FBG-based tunable dispersion compensator is demonstrated recently, with ± 400 ps/nm range.⁸² The parameters of the compensator are shown in Fig. 29. We can see that the tunable dispersion compensator exhibits uniform channel profiles, with flat top, steep edges, and low crosstalk.

Although currently no technology is a clear winner, the trend of dispersion compensation is toward tunable devices, or even actively self-tunable compensators, and such devices will allow system designers to cope with the shrinking system margins and with the emerging rapidly reconfigurable optical networks.

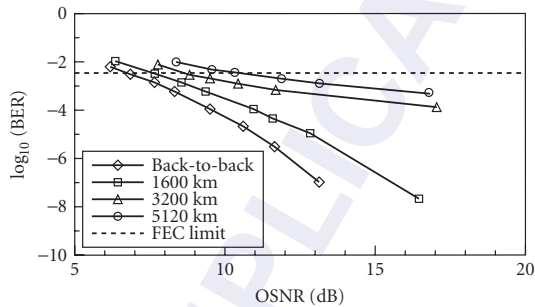
Electronic Solutions It is worth to mention that some of the most promising solutions of dispersion are electronic signal processing such as electronic equalizers and forward error correction (FEC) coding.^{83,84} The electronic equalizers rely on post-detection signal processing including filtering and adaptive signal processing to sharpen up distorted data pulses. Because the detection itself is nonlinear, the job of compensating for the linear distortions of chromatic dispersion is quite a bit more complicated. Many of the high-performance 40 Gb/s systems also incorporate FEC coding. Such coding adds some redundancy into the bits of a data stream to more easily find and correct errors. FEC is implemented using electronic chips, and it adds a system power margin that can ease the deleterious problems associated with fiber nonlinearities, chromatic dispersion, signal-to-noise ratio, and PMD. Note that electronic processing is potentially very cheap, and much easily scalable to large volume production, at least at data less than 10 Gb/s.

Analog equalizers that can combat the distortion produced by chromatic dispersion and PMD have been demonstrated at bit rates up to 43 Gb/s.⁸⁵ Hardware-implemented maximum likelihood sequence estimation (MLSE) has been demonstrated to reduce penalties from intersymbol interference and to extend the chromatic-dispersion-limited transmission length.^{86,87} Using electrical pre-equalization, 10-Gb/s transmission is demonstrated over 5120-km SMF without optical dispersion compensation.⁸⁸ The electrically precompensating transmitter is shown in Fig. 30a and Fig. 30b shows the BER measurement after transmission.

Another greatest trend in electronic signal processing for optical transmission is digital sampling and signal processing techniques for optical receivers.⁸⁹ Digital sampling techniques move several of the largest problems in optical transmission into the electrical domain. High-speed sampling followed by offline processing has been used in most experimental demonstrations of coherent detection with digital signal processing.^{90–93} It is still challenging for the benefits of coherent detection and digital signal processing to outweigh the cost of implementing real-time processing at high data rates.



(a)



(b)

FIGURE 30 Ten Gb/s transmission over 5120-km SMF without optical dispersion compensation using electrical preequalization: (a) electrically precompensating transmitter and (b) BER versus OSNR after transmission.⁸⁸

21.4 OPTICAL MODULATION FORMATS FOR WDM SYSTEMS

Most of current fiber systems use binary modulation with error-control coding schemes. The spectral efficiency cannot exceed $1 \text{ b/s} \cdot \text{Hz}$ per polarization regardless of detection technique. With the increase of bit rate and decrease of channel spacing in WDM systems, higher spectral efficiency is required. To achieve spectral efficiencies above $1 \text{ b/s} \cdot \text{Hz}$ and increase overall capacity of a WDM transmission system, more advanced modulation formats will be needed. The types of data modulation formats also have substantial impacts on fiber impairments. Due to the fact that optical signals propagating in fibers offer several degrees of freedom, including amplitude, frequency, phase, polarization, and time, intense research efforts have been made toward the combination coding over these degrees of freedom as a means to increase fiber transmission capacity, especially as a way to combat or benefit from fiber impairments. In this section, we will discuss the optical modulation formats of digital signals in WDM systems. We will highlight a few examples and present their advantages and disadvantages based on fiber system performance characterization.

Basic Concepts

The digital signal, which may be modulated at approximately gigabits per second rates, is being transmitted on an optical carrier wave whose frequency is in the multiterahertz regime. This optical carrier wave, $A(t)$, has an intensity amplitude, A_0 , an angular frequency, ω_c , and a phase ϕ :⁹⁴

$$A(t) = A_0 \cos(\omega_c t + \phi) \quad (8)$$

A binary digital signal implies transmitting two different quantities of anything which can subsequently be detected as representing a “1” and “0,” that is, we can transmit blue and red, and this can represent “1” and “0” in the receiver electronics if blue and red can be distinguished. We can therefore modulate either the amplitude, frequency, or phase of the optical carrier between two different values to represent either a “1” or “0,” known respectively as amplitude-, frequency-, and phase-shift keying (ASK, FSK, and PSK), with the other two variables remaining constant:

$$\begin{aligned}
 A_{\text{ASK}}(t) &= [A_0 + m(\Delta A)]\cos(\omega_c t + \phi) & m &= \begin{cases} +1, "1" \\ -1, "0" \end{cases} \\
 A_{\text{FSK}}(t) &= A_0 \cos\{[\omega_c + m(\Delta\omega)]t + \phi\} \\
 A_{\text{PSK}} &= A_0 \cos[\omega_c t + (\phi + m\pi)]
 \end{aligned} \tag{9}$$

where ΔA is the amplitude modulation and is less than A_0 , and $\Delta\omega$ is the FSK frequency deviation. ASK has two different light amplitude levels; FSK has two different optical carrier wavelengths; and PSK has two different phases which can be detected as an amplitude change in the center of the bit time for which a “1” or “0” bit can be determined. It is important to emphasize that the differential-phase-shift-keying (DPSK) format, in which the phase of the preceding bit is used as a relative phase reference, has been reemerged in the last few years due to its less OSNR requirement and robustness to fiber nonlinearities.^{95–97} The DPSK modulation signal is not the binary code itself, but a code that records changes in the binary stream. The PSK signal can be converted to a DPSK signal by the following rules: a “1” in the PSK signal is denoted by no change in the DPSK, a “0” in the PSK signal is denoted by a change in the DPSK signal. For a DPSK signal, optical power appears in each bit slot, and can occupy the entire bit slot (NRZ-DPSK) or can appear as an optical pulse (RZ-DPSK).

Figure 31 shows the impression of a simple digital signal on the optical carrier. These three formats can be implemented by appropriately changing the optical source, whether by modulating the

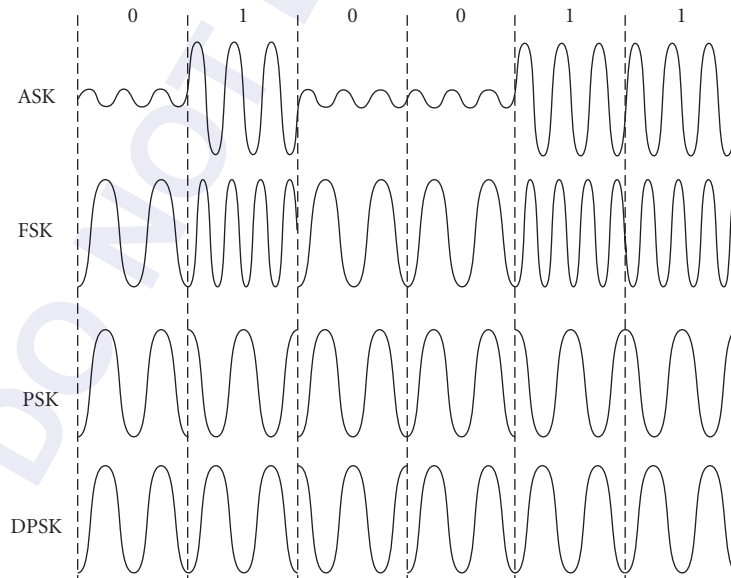


FIGURE 31 ASK, FSK, and PSK (DPSK) time modulation while employing an optical carrier wave.

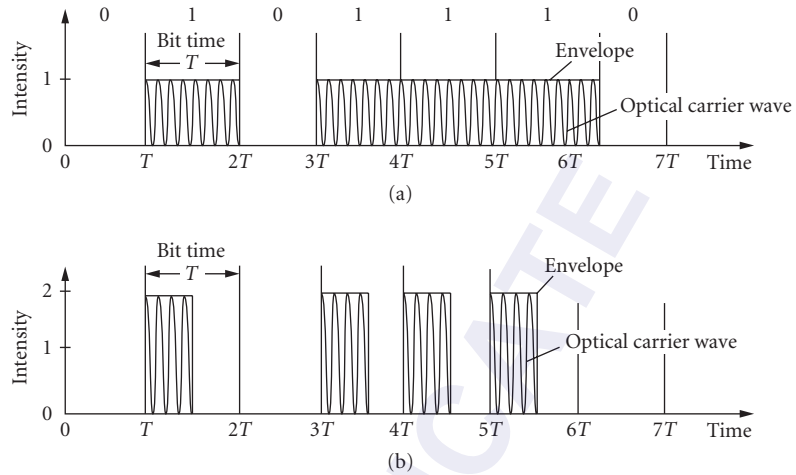


FIGURE 32 Modulation formats: (a) NRZ (non-return-to-zero) modulation format and (b) RZ (return-to-zero) modulation format.

light amplitude, laser output wavelength, or using an external phase shifter. ASK is important since it is the simplest to implement, FSK is important because a smaller chirp is incurred when direct modulation of a laser is used, and PSK is important because it, in theory, requires the least amount of optical power to enjoy error-free data recovery.⁹⁸

It should be mentioned that ASK, which is by far the most common form, is called on-off-keying (OOK) if the “0” level is really at zero amplitude. Figure 32 shows NRZ and RZ OOK. NRZ is the simplest format in which the amplitude level is high during the bit-time if a “1” is transmitted and is low if a “0” is transmitted. RZ format requires that a “1” always return to the low state during the bit-time even when two “1”s are transmitted in sequence, whereas the “0” bit remains at a low level. This eliminates the possibility of a long string of “1”s producing a constant high level but does not eliminate a long string of “0”s from producing a constant low level. The main attraction of the RZ format is its demonstrated improved immunity to fiber nonlinearities relative to NRZ. Note that the frequency of the optical carrier is so high that the optical detector, whose electronic bandwidth is usually $\ll 100$ GHz, will not detect it and only the envelope of the gigabits per second data will be electrically recovered.

Once a signal has been modulated and transmitted, it must be accurately detected. Two detection schemes include direct detection and coherent detection.^{99,100} Direct detection can be used for ASK and FSK modulation. PSK generally requires complicated coherent detection, which is based on mixing of two light waves to detect the phase information.

Direct detection, which is by far the simpler of the two schemes, involves detecting the amount of optical power incident on the optical detector. Direct detection of an ASK (or OOK) signal, that is, light “on” or light “off,” is extremely simple to accomplish using a detector and a high-bandwidth power meter (see Fig. 33a). This ASK signal is recovered by a detector of a certain electrical low-pass-filtering bandwidth. The electrical spectrum of the recovered ASK signal, which is sent in random NRZ format and can be measured on an RF spectrum analyzer, is shown in Fig. 33b. Only the first lobe of RF spectrum is necessary for recovering the data since only the transition edges would be affected by cutting off the higher lobes. The first lobe is considered the baseband signal representing the data stream.¹⁰¹ It is important to mention that the sensitivity of a direct detection receiver can be improved significantly by using a low noise optical preamplifier that is just before the photodetector. The output of the optical amplifier needs to be high enough so that at the photodetector the noise is dominated by the signal-spontaneous beat noise of the optical preamplifier.

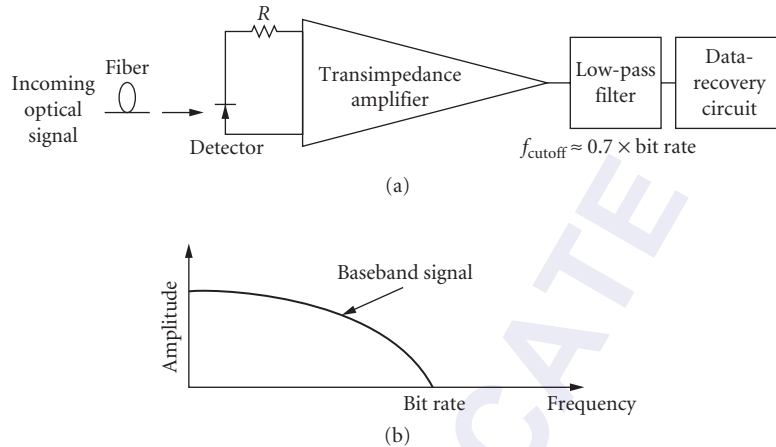


FIGURE 33 (a) Direct detection optical system and (b) baseband signal of a directly-detected NRZ signal.

Carrier-Suppressed Return-to-Zero and Duobinary

It is important to note that the type of data modulation formats have substantial impacts on the fiber dispersive and nonlinear effects. Some of these formats carry information through OOK, but also modulate the optical phase in a noninformation-bearing way in order to enhance the signals' robustness to chromatic dispersion, optical filtering, and/or nonlinearities. This group includes formats such as optical duobinary, chirped return-to-zero (CRZ), and alternating-phase OOK formats such as carrier-suppressed return-to-zero (CS-RZ). We will highlight a few examples here and present their advantages and disadvantages based on fiber system performance characterization.

CS-RZ optical signals have the feature of presenting bits that are π phase-shifted relative to neighboring bits such that on average all the phases cancel each other out for a net phase of zero. A CSRZ optical data stream may consist of a plurality of pulses where half of the pulses have an alternating phase relationship with the other half of the pulses, which leads to carrier suppression, as shown in Fig. 34a. Because of this phase inversion between adjacent bit periods that reduces interbit interference, CSRZ signals show increased tolerance dispersion and nonlinear penalties (see Fig. 34b).¹⁰²

The CS-RZ format can be generated by using a sinusoidal signal to drive the Mach-Zehnder modulator; the drive frequency is one half of the bit rate ($f_r/2$) and the amplitude is $2V_\pi$. As shown in Fig. 35, the generated pulse has the repetition rate of f_r , and the phase of the pulses alternates between 0 and π .

Optical duobinary has been proven more resilient to dispersion for more than 10 Gb/s data and is reasonably simple to implement. As shown in Fig. 36, optical duobinary signal is commonly generated by applying a baseband duobinary (three-level) electrical signal to the RF input of a lithium niobate (LiNbO_3) Mach-Zehnder modulator biased at maximum distinction (at V_π). The result is a binary, intensity-modulated optical signal, with a duobinary-modulated optical electric field caused by π -radian shift. With a zero ($c_k = 1$) input, no light is transmitted, but the +1 ($c_k = 2$) and -1 ($c_k = 0$) inputs are transmitted as +E and -E electrical fields. While this is a three-level signal in terms of the electric field, it is a two-level signal in terms of optical power. The same receiver that is used for a NRZ modulation scheme can be used for duobinary modulation. The power detector squares the electric field to detect power and hence the +E and -E outputs of the fiber get mapped to the same power level and are detected as logical 1s. Note that the received data is the invert of the original binary input. This choice significantly reduces the complexity of the receiver (the first optical duobinary system used a mapping that requires three levels of optical power). One of the key components is a driver that can produce a voltage swing of $2V_\pi$ V at high data rates such as more than 10 Gb/s in order to drive the Mach-Zehnder modulator.

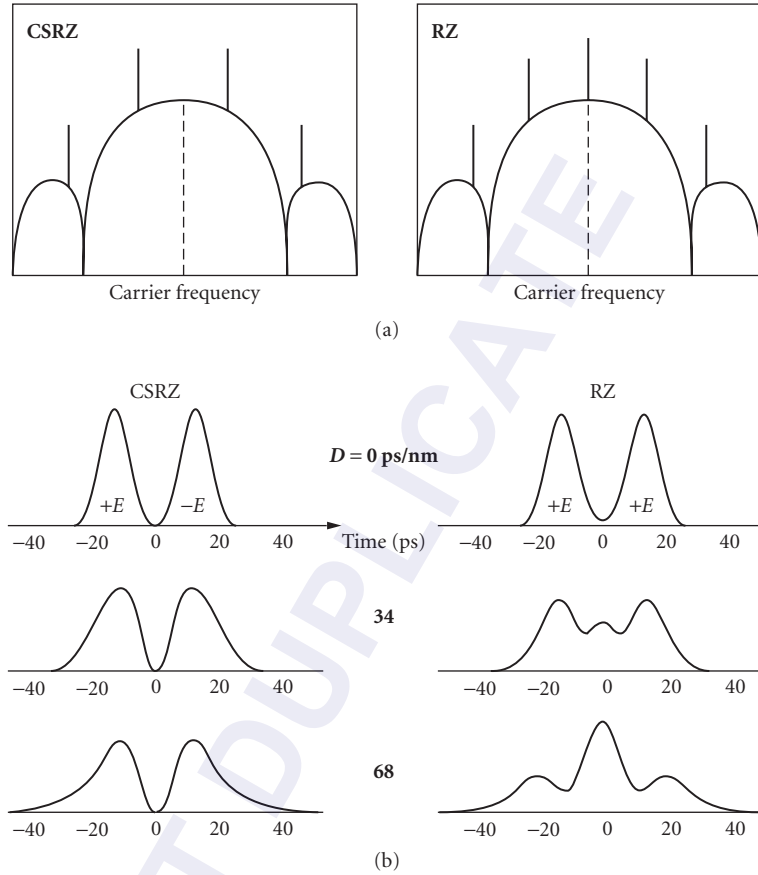


FIGURE 34 (a) Optical spectra of CS-RZ and RZ formats and (b) 40 Gb/s CS-RZ and RZ pulses under different dispersion values.¹⁰²

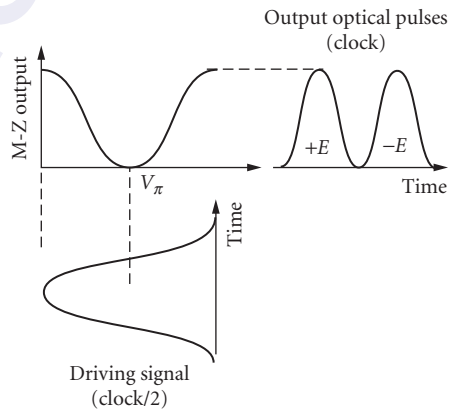


FIGURE 35 Generation of CS-RZ format using Mach-Zehnder modulator.

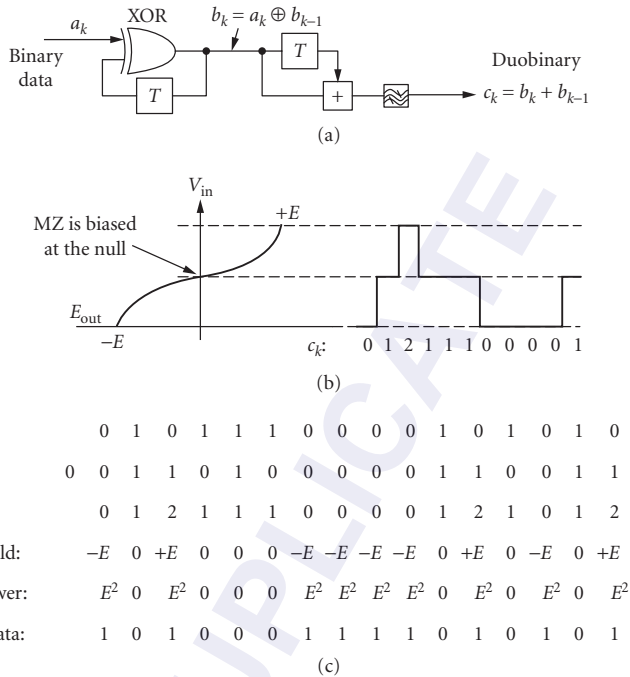


FIGURE 36 Optical duobinary modulation format: (a) duobinary encoder; (b) biasing of Mach-Zehnder modulator; and (c) an example of data transformation.

The combination of the duobinary encoder and the above mapping of electric fields help reduce the effects of dispersion in the fiber. The pulses spread out as they travel down the fiber. In an NRZ scheme, a data sequence of 1 0 1 is mapped onto the optical domain as $+E$ 0 $+E$. In the encoded duobinary sequence, a 1 0 1 sequence cannot occur, but a 1 0 -1 does occur, which is mapped to $+E$ 0 $-E$ in the optical domain. The effect of dispersion in the two cases is shown in Fig. 37, which depicts why the

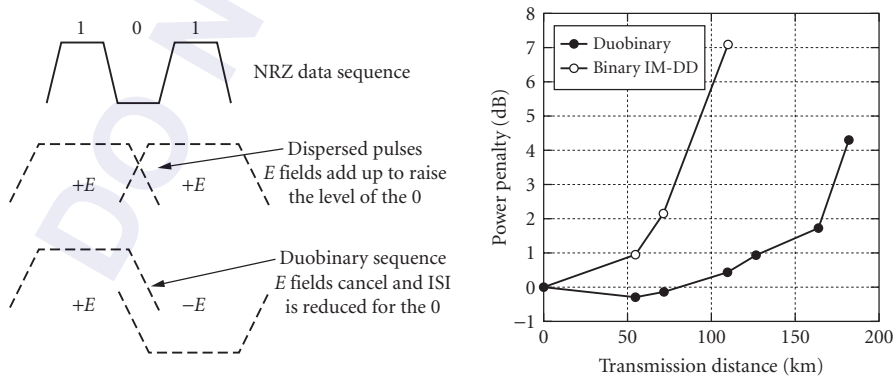


FIGURE 37 Dispersion effects on NRZ format and duobinary format. The power penalty due to fiber chromatic dispersion was measured at the BER of 10^{-9} .¹⁰³

resulting dispersion is less in the case of duobinary modulation. Figure 37 also shows an experimental result that optical duobinary technique expands the transmission distance to more than 150 km SMF for 10 Gb/s data.¹⁰³

The optical duobinary technique has been proven in experiment to expand the usable bandwidth and transmission distance in a four-channel multiplexed 40-Gb/s-based WDM system without individual channel dispersion slope compensation. That 160-Gb/s WDM transmission using four 40-Gb/s optical duobinary channels over a 100-km DSF is successful by virtue of the high dispersion tolerance of the optical duobinary signal.¹⁰⁴ Furthermore, the narrow optical spectrum of optical duobinary signals provides high spectral efficiency and reduces the coherent crosstalks in the ultradense WDM system.¹⁰⁵

DPSK and DQPSK

With the transmission capacity increases in WDM systems, PSK systems regain much interest in recent years. As we mentioned in the early section, PSK formats carry the information in the optical phase itself and DPSK formats carry the information in optical phase changes between bits. Since photodetector is inherently insensitive to the optical phase; a detector only converts the optical signal power into an electrical signal, directly detecting PSK signal is impractical due to the lack of an absolute phase reference. Therefore, PSK systems generally require complicated coherent detection which needs a local oscillator (laser) to mix with the received signal light.¹⁰⁶ However, we can detect the DPSK signal using a 1-bit delayed Mach-Zehnder interferometer followed by a balanced direct-detection system. A DPSK signal records the phase changes in the binary stream. Thus the demodulator only needs to determine these changes in the coming signal phase. A typical balanced DPSK receiver is shown in Fig. 38. The optical DPSK signal is first sent to a Mach-Zehnder interferometer (MZI) with the 1-bit period differential delay between the two arms. The MZI lets two adjacent bits interfere with each other at its output ports. This interference leads to the presence (absence) of power at an MZI output port if two adjacent bits interfere constructively (or destructively) with each other. Thus, the preceding bit in a DPSK-encoded bit stream acts as the phase reference for the

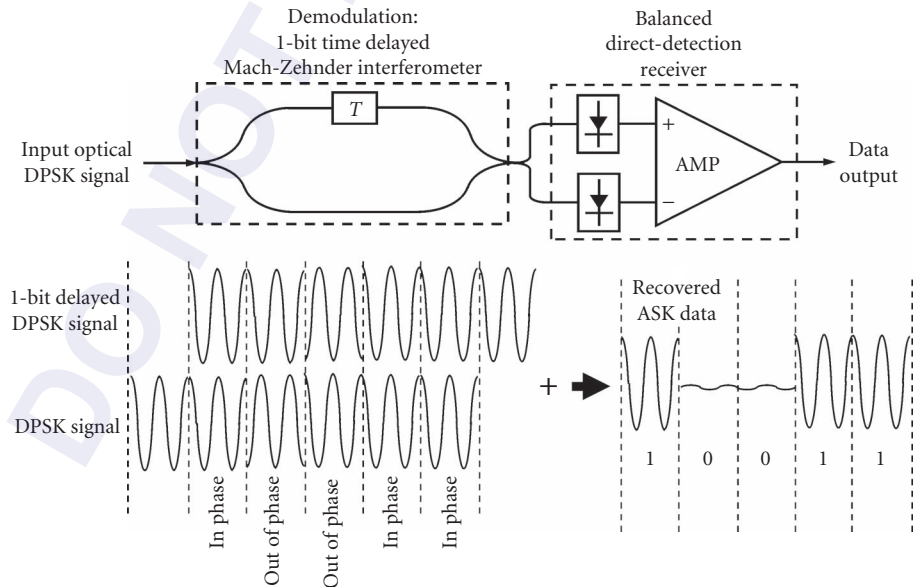


FIGURE 38 A typical DPSK receiver.

current bit. No phase change between these two bits will generate a “1” (or “0”) at the constructive port (or destructive port). Recall the generation of DPSK from PSK signals, this process is exactly the reverse process. Ideally, one of the MZI output ports is adjusted for destructive interference in the absence of phase modulation (*destructive port*), while the other output port then automatically exhibits constructive interference due to energy conservation (*constructive port*). For the same reason, the two MZI output ports will carry identical, but logically inverted data streams under DPSK modulation.⁹⁵

Using the balanced direct-detection scheme, the DPSK system has the advantage of requiring a lower OSNR than OOK to reach a given BER. Intuitively, this can be understood by comparing the signal constellations for DPSK and OOK.⁹⁹ To achieve the same symbol distance, the average optical power in DPSK is only half as compared to OOK. At 40 Gb/s, a sensitivity of about 38 photons/bit has been reported using RZ-DPSK.¹⁰⁷ This is approximately 3 dB better than the best OOK results of 78 photons/bit.¹⁰⁸ The lower OSNR requirement of DPSK can be used to extend transmission distance, reduce optical power requirements, or relax component specifications. Both numerical simulations and experiments have also shown DPSK to be more robust to some nonlinear effects than OOK⁹⁷ due to the following facts: (1) the optical power is evenly distributed (power is present in every bit slot for DPSK, which reduces bit-pattern-dependent nonlinear effects), and (2) the optical peak power is 3 dB lower for DPSK than for OOK for the same average optical power. Figure 39 shows the simulation results of comparison of the transmission of 43 Gb/s signals, NRZ, CSRZ and RZ-DPSK.¹⁰⁹

To further increase the bit rates and spectral efficiency in WDM systems, an extension to differential quadrature phase-shift keying (DQPSK) is introduced.⁹⁵ It transmits the four phase shifts $\{0, \pi/2, -\pi/2, \pi\}$ at a symbol rate of half the aggregate bit. DQPSK requires relatively complicated transmitter and receiver. As shown in Fig. 40, the transmitter consists of two parallel DPSK modulators that are integrated together in order to achieve phase stability. The receiver essentially consists of two DPSK receivers, although the phase difference in the arms of the delay interferometers is now set to $\pi/4$ and $-\pi/4$.

As compared to DPSK, the required OSNR to reach a given BER is increased by about 1 to 2 dB, depending on the BER.¹¹⁰ Also, the frequency offset tolerance between the laser and the delay interferometer is about 6 times less than for DPSK,¹¹¹ making the delay interferometer design and stabilization somewhat challenging. The benefit of DQPSK is that, for the same data rate, the symbol rate is reduced by a factor of two. Consequently, the spectral occupancy is reduced, the transmitter and receiver bandwidth requirements are reduced, and the chromatic dispersion and PMD limitations are

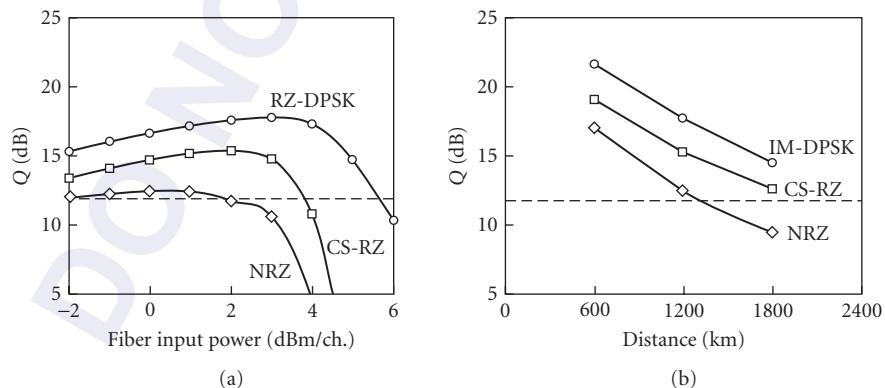


FIGURE 39 Comparison of 43 Gb/s signals, NRZ, CSRZ, and RZ-DPSK, with channel spacing of 75 GHz, (a) after 1200 km NZ-DSF (dispersion = 8 ps/nm · km) transmission and (b) at the optimum fiber input power. The broken line indicates the 11.8-dB Q-factor that corresponds a BER of 10^{-15} after Reed Solomon (255, 239) forward error correction.¹⁰⁹

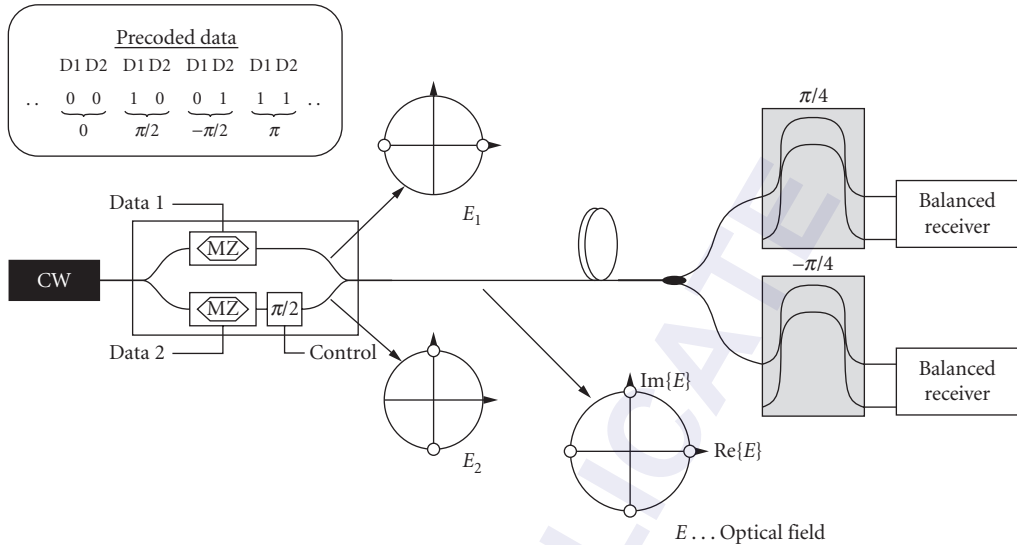


FIGURE 40 A typical DQPSK transmitter and receiver.⁹⁵

extended. The tolerance of DQPSK signal to PMD enables polarization-multiplexed data transmission to increase the spectral efficiency further. As shown in Fig. 41, a recent experiment demonstrated of 1.6-b/s · Hz spectrally efficient transmission over 1700-km single-mode fiber, using 40-channel 85.6-Gb/s (totally 3.2-Tb/s) polarization-multiplexed RZ-DQPSK.¹¹²

Even higher spectral efficiency can be achieved using various combinations of phase- and amplitude-shift keying.^{113–115} Such multilevel modulation can also improve system tolerance to chromatic dispersion and PMD. However, these schemes quickly become quite complicated to implement, require higher OSNR, and are sensitive to nonlinear phase noise.

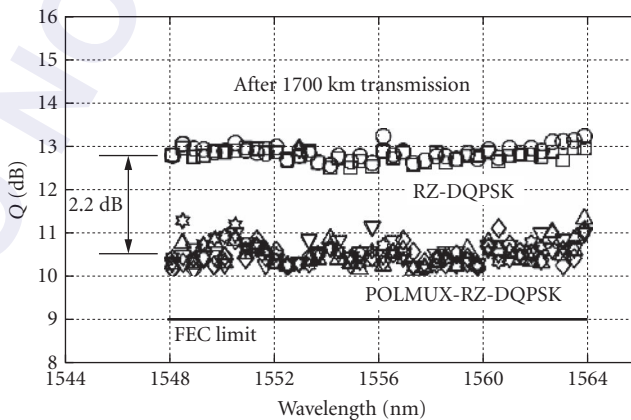


FIGURE 41 Experimental demonstration of 1.6-b/s · Hz spectrally efficient transmission over 1700-km single-mode fiber, using 40-channel 85.6-Gb/s (totally 3.2-Tb/s) polarization-multiplexed RZ-DQPSK.¹¹²

TABLE 1 Overview of Modulation Formats and Some Performance Values at 42.7 Gb/s (Required OSNR at BER = 10^{-3})¹¹⁶

Modulation Format	TX Complexity	RX Complexity	Required OSNR (dB)	CD (ps/nm) (2-dB Penalty)	DGD (ps) (1-dB Penalty)
NRZ-OOK	1 MZM	1 PD	15.9	54	8
50% RZ-OOK	1–2 MZMs	1 PD	14.4	48	10
67% CSRZ-OOK	2 MZMs	1 PD	14.9	42	11
Duobinary	1 MZM	1 PD	16.6	211	6
NRZ-DPSK	1 MZM	1 DI + 2 PDs	11.7	74	10
50% RZ-DPSK	1–2 MZMs	1 DI + 2 PDs	11.1	50	10
NRZ-DQPSK	2 nested MZMs	2 DI + 4 PDs	13.2	168	20
50% RZ-DQPSK	2 nested MZMs	2 DI + 4 PDs	12.2	161	21

PD: photodiode; OF: optical filter; DI: delay interferometer; MZM: Mach–Zehnder modulator; PC: pulse carver; CD: chromatic dispersion; DGD: differential group delay.

Table 1 gives an overview of some key characteristics of the optical modulation formats.¹¹⁶ The second and third columns summarize the transmitter and receiver hardware complexities, respectively, in terms of the optoelectronic component requirements. The fourth column specifies required OSNR for BER equal to 10^{-3} based on BER simulations that properly take into account the non-gaussian noise statistics of beat-noise-limited detection. The assumed 42.7 Gb/s are representative of a 40-Gb/s per-channel bit rate, including a 7 percent overhead for FEC, as standardized for terrestrial fiber transmission systems. The fifth column quantifies the accumulated chromatic dispersion that yields a 2-dB penalty in OSNR, and the sixth column quantifies the tolerance (1-dB OSNR penalty) of different modulation formats to the first-order PMD.

While actually measured values for required OSNR may differ somewhat from the numbers given in Table 1 due to various optical and electronic hardware implementation aspects, some general facts are

1. RZ formats in general require 1 to 3 dB less OSNR for identical BER than their NRZ equivalents, which is mostly due to the reduced impact of intersymbol interference (ISI) on RZ formats.
2. Using DPSK instead of intensity modulation, OSNR requirements are significantly reduced. The gain of balanced-detection DPSK over OOK is generally independent of the target BER and typically amounts to around 3 dB. Depending on the modulation waveforms, extinction ratios, and optical as well as electrical filters, the gain of DPSK can also even exceed 3 dB.
3. DQPSK requires only 1 to 1.5 dB higher OSNR than DPSK at poor BER (e.g., 10^{-3}), though the OSNR gap between DPSK and DQPSK increases at good BER (e.g., 10^{-12}).¹¹⁷ The good OSNR performance at FEC error ratios makes DQPSK an attractive candidate for optically routed networks that require narrow optical signal spectra.¹¹⁸
4. Duobinary and DQPSK have significantly better dispersion tolerance than other formats due to their narrower spectra.
5. For most modulation formats, a 1-dB penalty occurs at a DGD between 30 and 40 percent of the symbol duration, with RZ formats being in general more resilient to PMD than NRZ formats.¹¹⁹ Since the tolerance to the first-order PMD scales linearly with symbol duration, DQPSK has about twice the PMD tolerance of binary modulation formats at the same bit rate.

Table 2 shows a summary comparison of the various modulation formats relative to NRZ-OOK.¹²⁰ Note that the performance of a modulation format may depend significantly on its implementation details. The choice of the most suitable format depends on the application (metropolitan, regional, or long-haul), bit rate, wavelength spacing, optical power level, fiber type and dispersion map, the number of pass-through nodes and associated multiplexing filtering technique, the amplification scheme, as well as other system requirements.¹²¹ The “optimum” implementation characteristics of a modulation format often depend on a subtle interplay of several parameters and require extensive studies to maximize system performance.¹²²

TABLE 2 Comparison of the Various Modulation Formats Relative to NRZ-OOK at 40 Gb/s¹²⁰

Mod. Format	CS-RZ	Duobinary	DPSK	DQPSK	DPSK-RZ	DQPSK-RZ
Performance vs. NRZ						
OSNR sensitivity	Slightly better	Slightly worse	Much better	Slightly better	Much better	Better
CD tolerance and spectral efficiency	Slightly worse	Much better	Slightly better	Much better	Slightly worse	Much better
PMD tolerance	Better	Equivalent	Slightly better	Much better	Better	Much better
Nonlinearity tolerance	Better	Equivalent	Better	Equivalent	Much better	Equivalent
Cost and complexity	Slightly worse	Equivalent	Slightly worse	Much worse	Worse	Much worse

21.5 OPTICAL AMPLIFIERS IN WDM NETWORKS

The optical amplifier is ideally a transparent box which provides gain and is also insensitive to the bit-rate, modulation-format, power, and wavelengths of the signal(s) passing through it. The signals remain in optical form during amplification. Optical amplifiers have played a key role in the optical telecommunications world due to rapid device progress and revolutionary systems results.¹²³ In fact, much of advances in optical communications can be traced to the incorporation of optical amplifiers.

EDFA is the most widely used optical amplifier. In this section, we will consider some important issues about EDFAs with regards to their implementation in WDM systems. EDFAs have been used in multichannel WDM systems to compensate for: (1) fiber attenuation losses in transmission, (2) component excess losses, and (3) optical network splitting losses. These optical splitting losses can occur in a passive star, in which the optical power is divided by the number of users (N), or in a ring/bus in which there may possibly be optical tapping losses at each node.

Figure 42a shows WDM transmission in a conventional electrically regenerated system. Regenerators can correct for fiber attenuation and chromatic dispersion by detecting an optical signal and then

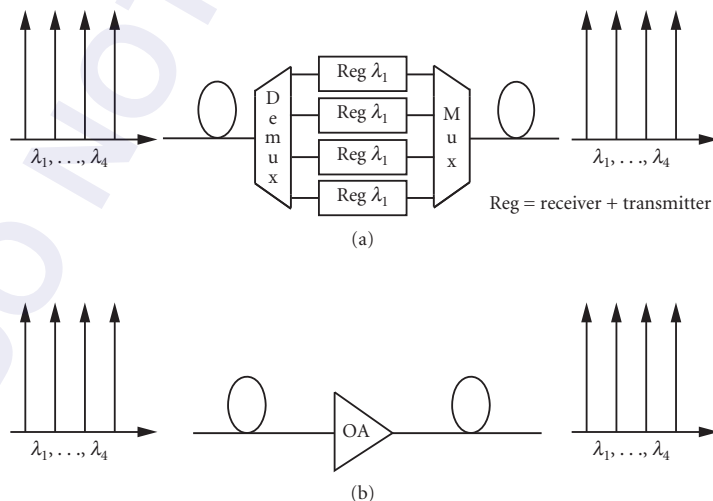


FIGURE 42 Wideband amplifiers enable WDM: (a) regeneration and (b) optical amplification.

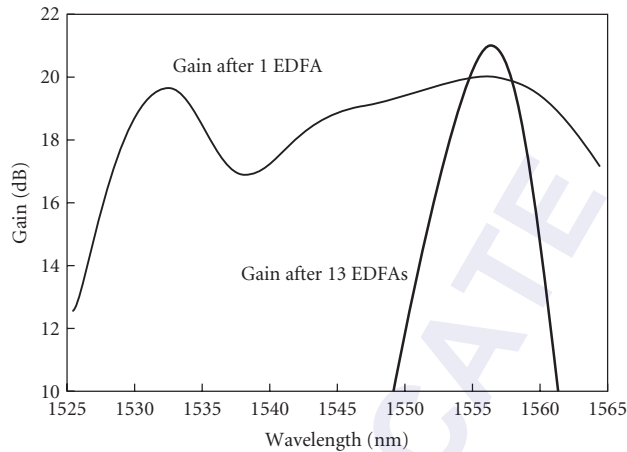


FIGURE 43 EDFA gain nonuniformity accumulation.

retransmitting it as a new signal using an internal laser. However, regenerators (being a hybrid of optics and electronics) are expensive, bit-rate and modulation-format specific, and waste much power and time in converting from photons to electrons and back again to photons. In contrast, as shown in Fig. 42b, the EDFA is ideally a transparent box which is insensitive to the bit-rate, modulation-format, power, and wavelengths of the signal(s) passing through it, and most importantly, provides gain for all the WDM channels simultaneously. Since all the channels remain in optical form during amplification, optically amplified WDM systems are potentially cheaper and more reliable than electrically regenerated systems.

Gain Peaking in EDFA Cascades

The EDFA is an almost ideal optical amplifier for WDM systems except for one major flaw: the gain is not uniform with wavelength, whereas the inter-amplifier losses are nearly wavelength independent.^{124–127} For a single amplifier, as shown in Fig. 43, the gain exhibits a peak at 1530 nm and a relatively flat region near 1555 nm. Moreover, the gain shape of an EDFA is dependent on the inversion of Er^{3+} in the Erbium-doped fiber.¹²⁸ When the inversion is low, which can be achieved by operating the amplifier in deep saturation, the gain peak at 1530 nm can be suppressed and the gain flatness around 1555 nm would become quite flat.

If several channels are located on the relatively flat shoulder region of the gain spectrum, then the gain differential after a single amplifier will be within a few decibels. However, when a cascade of EDFAs is used to periodically compensate for losses, the differential in gain and resultant OSNR can become quite severe. A large differential in SNR among many channels can be deleterious for proper system performance. Figure 43 shows the gain spectrum after a single amplifier and after 13 cascaded amplifiers. The gain does not accumulate linearly from stage to stage, and the resultant wavelength-dependent gain shape dramatically changes in a cascade. Along the cascade, gain is gradually “pulled” away from the shorter wavelengths and made available at the longer wavelengths, resulting in a usable bandwidth of only several nanometers.

EDFA Gain Flattening

We have shown the bandwidth reduction due to nonuniform gain in a cascade of EDFAs. It is clear that gain flattening is an important issue in optically amplified networks. Several methods have been reported for equalizing nonuniform EDFA gain. These methods include

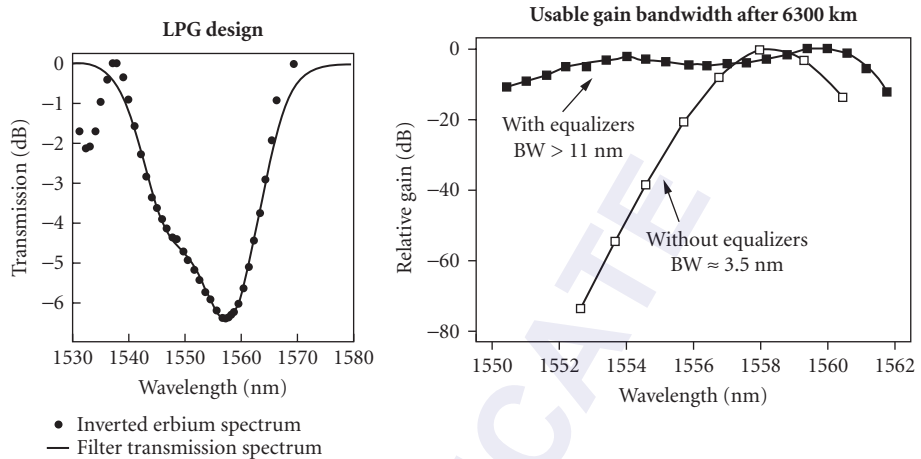


FIGURE 44 LPG design and gain equalization results.¹²⁹

1. Long period grating filters: A long period grating (LPG) with an index-varying period of approximately $100\ \mu\text{m}$ provides coupling between the core modes and the cladding modes, creating a wavelength-dependent loss to equalize the EDFA gain shape,^{129–131} as shown in Fig. 44.
2. Mach-Zehnder filters: The wavelength dependent transmission characteristics of cascaded Mach-Zehnder filters can be tailored to compensate for the gain nonuniformity of EDFAs.¹³²
3. Special designed EDFAs: A new coaxial dual-core gain-flattened EDF refractive index profile (RIP) is demonstrated recently, which is based on resonant coupling analogous to that in an asymmetric directional coupler. It has median gains more than 28 dB and gain excursion within ± 2 dB across the C-band.¹³³

Fast Power Transients

The lifetime of a stimulated erbium ion is generally approximately 10 ms, which seems to be long enough to be transparent to signals modulated by data at the rates of several gigabits per second or higher. However, the EDFAs could be critically affected by the adding or dropping of WDM channels, network reconfiguration, or link failures, as illustrated in Fig. 45. To achieve optimal channel SNRs, the EDFAs are typically operated in the gain-saturation regime where all channels must share

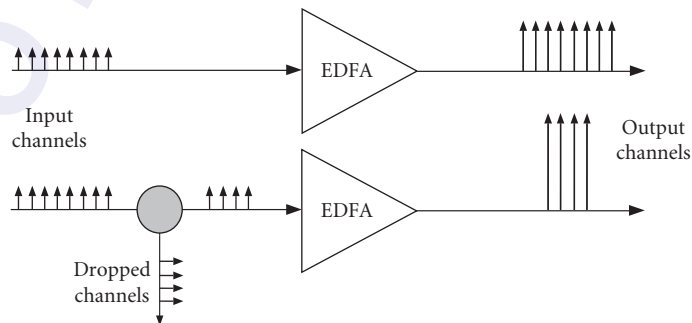


FIGURE 45 EDFA gain transients.

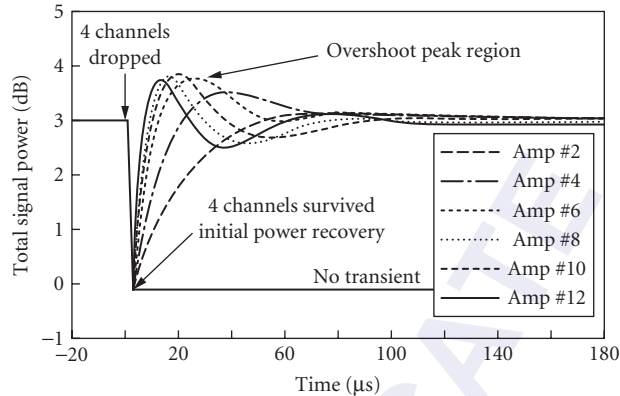


FIGURE 46 Fast power transients in EDFA cascades.¹³⁶

the available gain.^{134,135} Therefore, when channels are added or dropped, the power of the remaining channels will increase resulting transient effects.

The transients can be very fast in EDFA cascades.¹³⁶ As shown in Fig. 46, with an increase in the number of cascaded EDFAs, the transients can occur in approximately 2 μs. These fast power transients in chain-amplifier systems should be controlled dynamically, and the response time required scales as the size of the network. For large-scale networks, response times shorter than 100 ns may be necessary.

From a system point of view, fiber nonlinearity may become a problem when too much channel power exists, and a small SNR at the receiver may arise when too little power remains.¹³⁷ The corresponding fiber transmission penalty of the surviving channel is shown in Fig. 47 in terms of the Q factor, for varying numbers of cascaded EDFAs. When 15 channels are dropped or added, the penalties are quite severe. Note that this degradation increases with the number of channels N simply because of enhanced SPM due to a large power excursion as a result of dropping $N - 1$ channels.

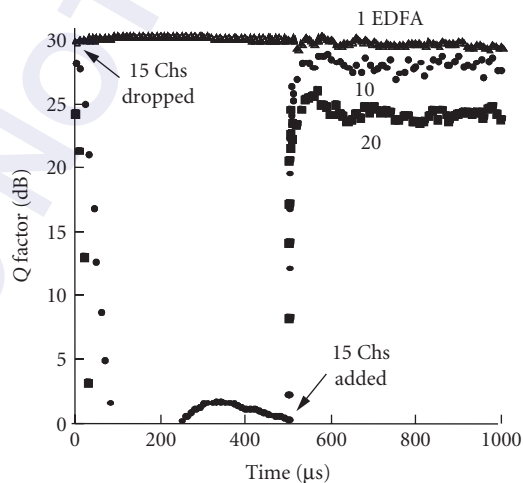


FIGURE 47 Q factor versus time for adding and dropping 15 channels of a 16-channel system at a bit rate of 10 Gb/s.¹³⁷

In order to maintain the quality of service, the surviving channels must be protected when channel add or drop or network reconfiguration occurs. The techniques include (1) optical attenuation, by adjusting optical attenuators between the gain stages in the amplifier to control the amplifier gain,¹³⁸ (2) pump power control, by adjusting the drive current of the pump lasers to control the amplifier gain,¹³⁹ (3) link control, using a power-variable control channel propagating with the signal channels to balance the amplifier gain,¹⁴⁰ and (4) EDFA gain clamping, by an automatic optical feedback control scheme to achieve all-optical gain clamping.¹⁴¹

Static Gain Dynamic and Channel Power Equalization

We just discussed EDFA gain flattening, which is a passive channel power equalization scheme effective only for a static link. However, in the nonstatic optical networks, the power in each channel suffers from dynamic network changes, including wavelength drift of components, changes in span loss, and channel add or drop. As an example, Fig. 48 shows how the gain shape of a cascaded EDFA chain varies significantly with link loss changes due to environmental problems. This is because the EDFA gain spectra are dependent on the saturation level of the amplifiers. The results in Fig. 48 are for a cascade of 10 gain-flattened EDFAs, each with 20-dB gain, saturated by 16 input channels with -18 dBm per channel.

System performance can be degraded due to unequalized WDM channel power. These degrading effects include SNR differential (reduced system dynamic range), widely varying channel crosstalk, nonlinear effects, and low signal power at the receiver. Therefore, channel power needs to be equalized dynamically in WDM networks to ensure stable system performance. To obtain feedback for control purposes, a channel power monitoring scheme is very important. A simple way to accomplish this is to demultiplex all the channels and detect the power in each channel using different photodetectors or detector arrays. To avoid the high cost of many discrete components in WDM systems with large numbers of channels, other monitoring techniques that take advantage of wavelength-to-time mapping have also been proposed including the use of concatenated FBGs or swept acousto-optic tunable filters.

Various techniques have been proposed for dynamic channel power equalization, including parallel loss elements,¹⁴² individual bulk devices (e.g., AOTFs),¹⁴³ serial filters,¹⁴⁴ micro-optomechanics (MEMS),¹⁴⁵ and integrated devices.^{146,147} As an example, Fig. 49 shows the parallel loss element scheme, where the channels are demultiplexed and attenuated by separate loss elements. An additional advantage of this scheme is that ASE noise is reduced by the WDM multiplexer and demultiplexer. Possible candidates for the loss elements in this scheme include optomechanical attenuators, acousto-optic modulators, and FBGs.

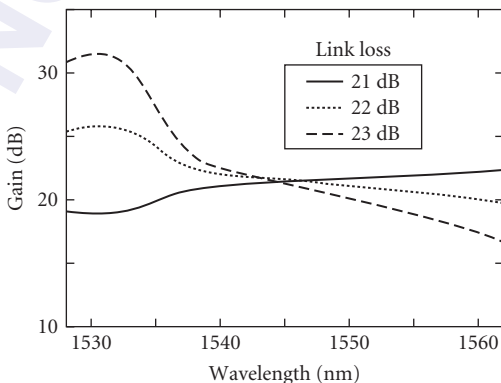


FIGURE 48 Gain spectra variation due to link loss changes for a cascade of 10 EDFAs.

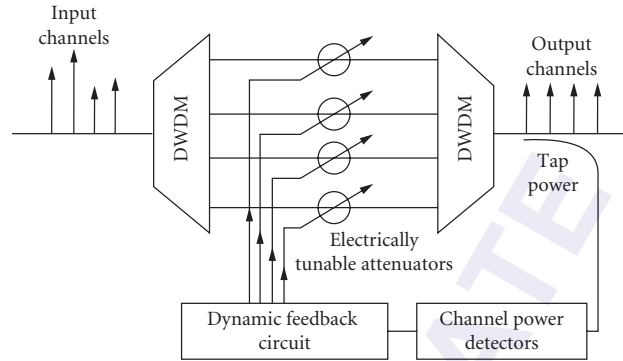


FIGURE 49 Parallel loss element scheme for dynamic channel power equalization.

Raman Amplifier

Raman amplifier is another important type of optical amplifier for WDM systems. The fundamental principles are based on Raman scattering as follows. The pump light photon is absorbed and sets the fiber molecules into mechanical vibrations. A photon is again radiated at the Stokes frequency, but since mechanical vibrations are not uniform in a fiber, the Stokes frequency is not a set number. Furthermore, the pump and signal may co- or counter-propagate in the fiber. It is worth to mention that practical, efficient, and high-power pump sources have diminished the disadvantage of the relatively poor efficiency of the Raman process over the last few years. Interest in Raman amplification has steadily increased.^{148,149}

The most important feature of Raman amplifiers is their capability to provide gain at any signal wavelength, as opposed to EDFAs based on the doped ions in the fibers. The position of the gain bandwidth within the wavelength domain can be adjusted simply by tuning the pump wavelength. Thus, Raman amplification potentially can be achieved in every region of the transmission window of the optical transmission fiber. It only depends on the availability of powerful pump sources at the required wavelengths. Figure 50 illustrates the Raman gain coefficient in a few different fibers.

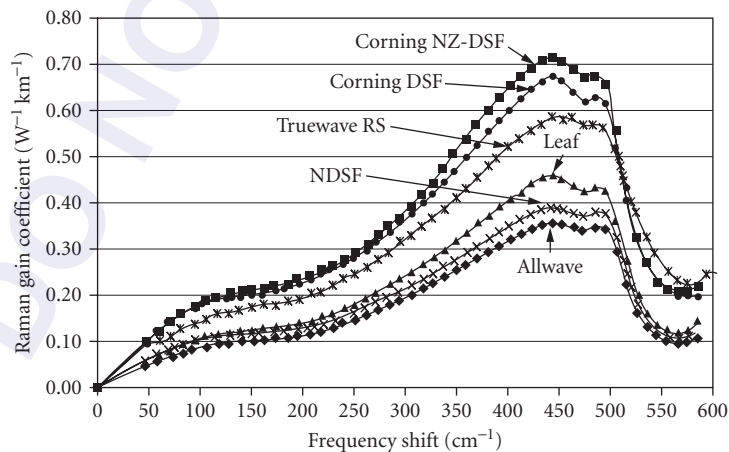


FIGURE 50 Raman gain spectra for different commercial fibers (gain peak is shifted 13 THz from the pump wavelength toward longer wavelength).¹⁴⁸

The disadvantage of Raman amplification is the need for high pump powers to provide a reasonable gain. However, the Raman effect can be used for signal amplification in transmission windows that cannot be covered properly by EDFAs. The upgrade of already existing systems by opening another transmission window where Raman amplification is applied could be an attractive application. Another application of the Raman effect is given with hybrid EDFA/Raman amplifiers characterized by a flat gain over especially large bandwidths. Repeaters can be built that compensate for the nonflatness of the EDFA gain with a more flexible Raman gain. Multiwavelength pumping could be used to shape the Raman gain such that it equalizes for the EDFA gain shaping.

Figure 51 shows a typical Raman amplifier that is backward pumped and the gain is distributed over the long transmission fibers.^{148,149} The spectral flexibility of Raman amplification allows the gain spectrum to be shaped by combining multiple pump wavelengths to make a polychromatic pump spectrum. There have been many studies searching for optimization approaches that give the flattest gain with the fewest number of pumps. Using this broadband pumping approach, amplifiers with gain bandwidths greater than 100 nm have been demonstrated.¹⁵⁰ When designing such broadband Raman amplifiers, one must consider the strong Raman interaction between the pumps. The short wavelength pumps amplify the longer wavelengths, and so more power is typically needed at the

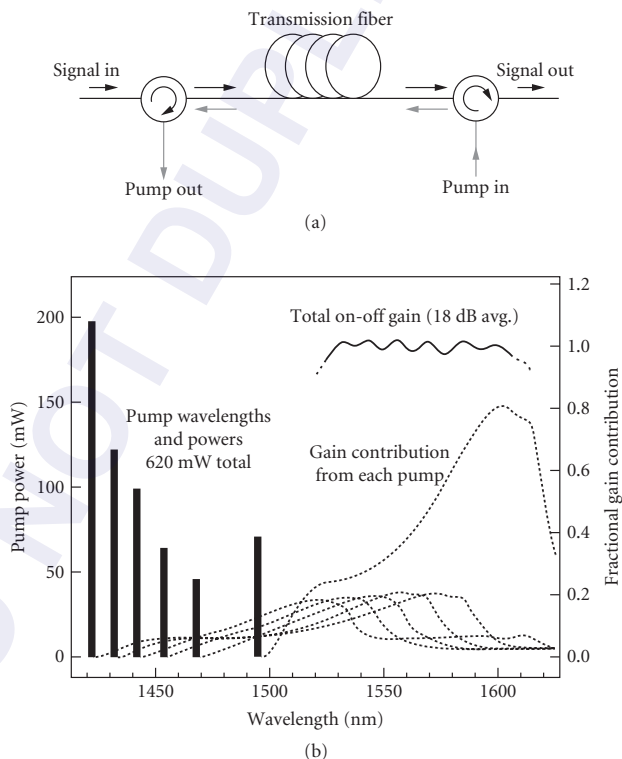


FIGURE 51 (a) Basic setup of backward pumped Raman amplifiers. (b) A numerical example of broadband Raman gain obtained using a broadband spectrum to pump a NZDSF. Bars show the counter-pump wavelengths and its power. Solid line shows the total small-signal on-off gain. Dashed lines show the fractional gain contribution from each pump wavelength.¹⁴⁹

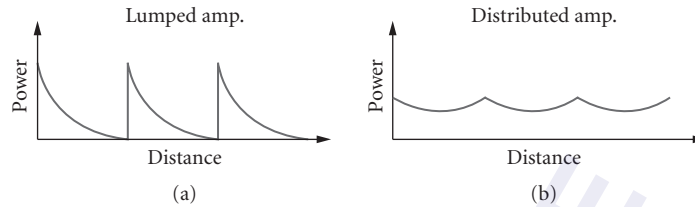


FIGURE 52 Power evolution along distance with (a) lump amplification (EDFA) only and (b) distributed amplification.

shortest wavelengths (see Fig. 51b). This interaction between the pumps also affects the noise properties of broadband amplifiers.

Another advantage of Raman amplifier is the feature of distributed amplification, since the transmission fiber itself can be used as a gain medium. As shown in Fig. 52a, in the conventional EDFA repeater systems, the signal monotonically attenuates in the fiber span, which is amplified at a point of the EDFA (lumped amplifier) location to recover the original level before entering the next fiber span. On the other hand, Raman amplifiers are mostly used in a distributed configuration, as shown in Fig. 52b. The transmission impairments are caused mostly by signal quality degradation due to optical nonlinearity in the transmission fiber and ASE noise entailed by optical amplifiers. In the presence of distributed Raman amplifier, the magnitude of the signal level excursion is smaller than the case with EDFA only, which can reduce both nonlinearity and degradation of OSNR due to ASE noise.^{151,152} A transmission of 6.4 Tb/s (160×42.7 Gb/s) over 3200 km of fiber has been demonstrated in a distributed Raman amplified system.¹⁵³

21.6 SUMMARY

In this chapter, we have covered many different aspects of high-speed WDM fiber-optic communication networks. We have endeavored to treat the most important topics—those that will likely impact these networks for years to come. The enormous growth of these systems is due to the revolutionary introduction of the EDFA. With the increasing knowledge, more development, higher data rates, and increasing channel count, WDM network limitations are being continually redefined. Network reconfigurability can offer great benefits for future WDM networks. However, a number of new degrading effects must be solved before reconfigurable networks become a reality. Yet the push for more bandwidth in WDM systems continues due to the enormous inherent potential of the optical fiber.

21.7 ACKNOWLEDGMENTS

We would like to extend our gratitude to Phillip Regan, Jing Yang, and Leroy Chee for their generous help to this chapter.

21.8 REFERENCES

1. F. P. Kapron, "Fiber-Optic System Tradeoffs," *IEEE Spectrum Magazine* **22**:68–75, 1985.
2. J. MacMillan, "Advanced Fiber Optics," *U.S. News and World Report*, p. 58, May 1994.
3. A. H. Gnauck, G. Charlet, P. Tran, P. J. Winzer, C. R. Doerr, J. C. Centanni, E. C. Burrows, T. Kawanishi, T. Sakamoto, and K. Higuma, "25.6-Tb/s WDM Transmission of Polarization-Multiplexed RZ-DQPSK Signals," *IEEE Journal of Lightwave Technology* **26**:79–84, 2008.

4. J. Seoane, A. T. Clausen, L. K. Oxenlowe, M. Galili, T. Tökle, and P. Jeppesen, "Enabling Technologies for OTDM Networks at 160 Gbit/s and beyond," pp. 22–28, Paper MG1, Orlando, Fla., October 2005.
5. M. Nakazawa, T. Yamamoto, and K. R. Tamura, "1.28 Tbit/s-70 km OTDM Transmission Using Third- and Fourth-Order Simultaneous Dispersion Compensation with a Phase Modulator," *IEEE Electronics Letters* **36**:2027–2029, 2000.
6. P. J. Winzer and R. -J. Essiambre, "Advanced Modulation Formats for High-Capacity Optical Transport Networks," *IEEE Journal of Lightwave Technology* **24**:4711–4728, 2006.
7. G. Varella, F. Pitel, and J. F. Marcero, "3 Tbit/s (300 × 11.6 Gbit/s) Transmission over 7380 km Using C+L Band with 25 GHz Spacing and NRZ Format," *Optical Fiber Communication Conference and Exhibit*, Paper PD22, Anaheim, Calif, March 2001.
8. J. -X. Cai, M. Nissov, C. R. Davidson, Y. Cai, A. N. Pilipetskii, H. Li, M. A. Mills, et al., "Transmission of Thirty-Eight 40 Gb/s Channels (> 1.5 Tb/s) over Transoceanic Distance," *Conference on Optical Fiber Communication (OFC) '02*, Paper PD FC-4, Anaheim, Calif, March 2002.
9. Y. Frignac, G. Charlet, W. Idler, R. Dischler, P. Tran, S. Lanne, S. Borne, et al., "Transmission of 256 Wavelength-Division and Polarization-Division-Multiplexed Channels at 42.7 Gb/s (10.2 Tbit/s capacity) over 3 × 100 km TeraLight™ Fiber," *Conference on Optical Fiber Communication (OFC) '02*, Paper: PD FC-5, Anaheim, Calif, March 2002.
10. J. Berthold, A. A. M. Saleh, L. Blair, and J. M. Simmons, "Optical Networking: Past, Present, and Future," *IEEE Journal of Lightwave Technology* **26**:1104–1118, 2008.
11. L. G. Kazovsky, W. -T. Shaw, D. Gutierrez, N. Cheng, and S. -W. Wong, "Next-Generation Optical Access Networks," *IEEE Journal of Lightwave Technology* **25**:3428–3442, 2007.
12. P. Kaiser and D. B. Keck, "Fiber Types and Their Status," *Optical Fiber Telecommunications II*, Chap. 2, p. 40, S. E. Miller and I. P. Kaminow, eds., Academic Press, New York, 1988.
13. C. A. Brackett, "Dense Wavelength Division Multiplexing: Principles and Applications," *IEEE Journal on Selected Areas in Communications* **8**(6):948–964, 1990.
14. I. P. Kaminow, "FSK with Direct Detection in Optical Multiple-Access FDM Networks," *IEEE Journal on Selected Areas in Communications* **8**:1005–1014, 1990.
15. P. E. Green, Jr., *Fiber Optic Networks*, Prentice Hall, Englewood Cliffs, N.J., 1993.
16. N. K. Cheung, K. Nosu, and G. Winzer, Special Issue on Wavelength Division Multiplexing, *IEEE Journal on Selected Areas in Communications* **8**, 1990.
17. A. E. Willner, I. P. Kaminow, M. Kuznetsov, J. Stone, and L. W. Stulz, "1.2 Gb/s Closely-Spaced FDMA-FSK Direct-Detection Star Network," *IEEE Photonics Technology Letters* **2**:223–226, 1990.
18. N. R. Dono, P. E. Green, K. Liu, R. Ramaswami, and F. F. Tong, "A Wavelength Division Multiple Access Network for Computer Communication," *IEEE Journal on Selected Areas in Communications* **8**:983–994, 1990.
19. W. I. Way, D. A. Smith, J. J. Johnson, and H. Izadpanah, "A Self-Routing WDM High-Capacity SONET Ring Network," *IEEE Photonics Technology Letters* **4**:402–405, 1992.
20. T. -H. Wu, *Fiber Network Service Survivability*, Artech House, Boston, Mass., 1992.
21. A. S. Acampora, M. J. Karol, and M. G. Hluchyj, "Terabit Lightwave Networks: The Multihop Approach," *AT&T Technical Journal* **66**:21–34, November/December 1987.
22. M. Schwartz, *Telecommunication Networks, Protocols, Modeling, and Analysis*, Addison Wesley, New York, 1987.
23. M. Jeong, H. C. Cankaya, and C. Qiao, "On a New Multicasting Approach in Optical Burst Switched Networks," *IEEE Communications Magazine* **40**:96–103, 2002.
24. I. Baldine, H. G. Perros, G. N. Rouskas, and D. Stevenson, "JumpStart: A Just-in-Time Signaling Architecture for WDM Burst-Switched Networks," *IEEE Communications Magazine* **40**:82–89, 2002.
25. J. E. Berthold, "Networking Fundamentals," *Conference on Optical Fiber Communications (OFC) '94*, Tutorial TuK, San Jose, Calif, February 1994.
26. J. Y. Hui, *Switching and Traffic Theory for Integrated Broadband Networks*, Kluwer Academic Publishers, Boston, 1990.
27. J. B. Yoo and G. K. Chang, "High-Throughput, Low-Latency Next Generation Internet Using Optical-Tag Switching," *U.S. Patent 6,111,673*, 1997.
28. M. W. Maeda, A. E. Willner, J. R. Wullert II, J. Patel, and M. Allersma, "Wavelength-Division Multiple-Access Network Based on Centralized Common-Wavelength Control," *IEEE Photonics Technology Letters* **5**:83–86, 1993.

29. K. K. Goel, "Nonrecirculating and Recirculating Delay Line Loop Topologies of Fiber-Optic Delay Line Filters," *IEEE Photonics Technology Letters* **5**:1086–1088, 1993.
30. I. Chlamtac, A. Fumagalli, and S. Chang-Jin, "Multibuffer Delay Line Architectures for Efficient Contention Resolution in Optical Switching Nodes," *IEEE Transactions on Communications* **48**:2089–2098, 2000.
31. A. Agrawal, L. Wang, Y. Su, and P. Kumar, "All-Optical Erasable Storage Buffer Based on Parametric Nonlinearity in Fiber," *Conference on Optical Fiber Communication (OFC) '01*, Paper ThH5, Anaheim, Calif., March 2001.
32. A. Rader and B. L. Anderson, "Demonstration of a Linear Optical True-Time Delay Device by Use of a Microelectromechanical Mirror Array," *Applied Optics* **42**:1409–1416, 2003.
33. D. R. Pape and A. P. Goutzoulis, "New Wavelength Division Multiplexing True-Time-Delay Network for Wideband Phased Array Antennas," *Journal of Optics A: Pure Applied Optics* **1**:320–323, 1999.
34. C. J. Chang-Hasnain, P. Ku, and J. Kim, S. Chuang, "Variable Optical Buffer Using Slow Light in Semiconductor Nanostructures," *Proceedings of the IEEE*. **91**:1884–1897, 2003.
35. S. Rangarajan, H. Zhaoyang, L. Rau, and D. J. Blumenthal, "All-Optical Contention Resolution with Wavelength Conversion for Asynchronous Variable-Length 40 Gb/s Optical Packets," *IEEE Photonics Technology Letters* **16**:689–691, 2004.
36. J. Elmirghani and H. Mouftah, "All-Optical Wavelength Conversion: Techniques and Applications in DWDM Networks," *IEEE Communications Magazine* **38**:86–92, 2000.
37. D. Nessel, T. Kelly, and D. Marcenac, "All-Optical Wavelength Conversion Using SOA Nonlinearities," *IEEE Communications Magazine* **36**:56–61, 1998.
38. I. Brener, M. H. Chou, and M. M. Fejer, "Efficient Wideband Wavelength Conversion Using Cascaded Second-Order Nonlinearities in LiNbO_3 Waveguides," *Conference on Optical Fiber Communication (OFC) '99*, Paper FB6, San Diego, Calif., February 1999.
39. A. Hsu and S. L. Chuang, "Wavelength Conversion by Cross-Absorption Modulation Using an Integrated Electroabsorption Modulator/Laser," *Summaries of Papers Presented at the Conference on Lasers and Electro-Optics (CLEO) '99*, Paper CThV3, Baltimore, Md., May 1999.
40. M. Baresi, S. Bregni, A. Pattavina, and G. Vegetti, "Deflection Routing Effectiveness in Full-Optical IP Packet Switching Networks," *IEEE International Conference on Communications* **2**:1360–1364, Anchorage, Alaska, May 2003.
41. F. -S. Choa, X. Zhao, Y. Xiuqin, J. Lin, J. P. Zhang, Y. Gu, G. Ru, et al., "An Optical Packet Switch Based on WDM Technologies," *IEEE Journal of Lightwave Technology* **23**:994–1014, 2005.
42. A. E. Willner, D. Gurkan, A. B. Sahin, J. E. McGeehan, and M. C. Hauer, "All-Optical Address Recognition for Optically-Assisted Routing in Next-Generation Optical Networks," *IEEE Communications Magazine* **41**: S38–S44, May 2003.
43. N. Calabretta, H. de Waardt, G. D. Khoe, and H. J. S Dorren, "Ultrafast Asynchronous Multioutput All-Optical Header Processor," *IEEE Photonics Technology Letters* **16**:1182–1184, 2004.
44. D. Gurkan, M. C. Hauer, A. B. Sahin, Z. Pan, S. Lee, A. E. Willner, K. R. Parameswaran, and M. M. Fejer, "Demonstration of Multi-Wavelength All-Optical Header Recognition Using a PPLN and Optical Correlates," *European Conference on Optical Communication (ECOC) '01*, Paper We.B.2.5, Amsterdam, NL, September/October 2001.
45. J. Bannister, J. Touch, P. Kamath, and A. Patel, "An Optical Booster for Internet Routers," *8th International Conference. High Performance Computing*, pp. 399–413, Hyderabad, India, December 2001.
46. D. J. Blumenthal, J. E. Bowers, L. Rau, Hsu-Feng Chou, S. Rangarajan, Wei Wang, and K. N. Poulsen, "Optical Signal Processing for Optical Packet Switching Networks," *IEEE Communications Magazine* **41**:S23–S29, 2003.
47. L. Rau, S. Rangarajan, D. J. Blumenthal, H. -F. Chou, Y. -J. Chiu, and J. E. Bowers, "Two-Hop All-Optical Label Swapping with Variable Length 80 Gb/s Packets and 10 Gb/s Labels Using Nonlinear Fiber Wavelength Converters, Unicast/Multicast Output and a Single EAM for 80- to 10 Gb/s Packet Demultiplexing," *Optical Fiber Communication Conference (OFC) '02*, Pages FD2-1, Anaheim, Calif., March 2002.
48. S. Yao, B Mukherjee, and S. Dixit, "Advances in Photonic Packet Switching: An Overview," *IEEE Communications Magazine* **38**:84–94, 2000.
49. V. W. S. Chan, K. L. Hall, E. Modiano, and K. A. Rauschenbach, "Architectures and Technologies for High-Speed Optical Data Networks," *IEEE Journal of Lightwave Technology* **16**:2146–2168, 1998.
50. Winston I. Way, "Rules of the ROADM, A Deployment Guide," *Electronic Engineering Times*, pp. 60–64, November 2005. <http://www.eetimes.com/news/latest/showArticle.jhtml?articleID=173601985>.

51. R. S. Bernhey and M. Kanaan, "ROADM Deployment, Challenges, and Applications," *Optical Fiber Communication and the National Fiber Optic Engineers Conference (OFC/NFOEC)'07*, Paper NWD1, Anaheim, CA, USA, March 2007.
52. B. P. Keyworth, "ROADM Subsystems and Technologies," *Optical Fiber Communication and the National Fiber Optic Engineers Conference (OFC/NFOEC)'05*, Paper OWB5, Anaheim, CA, USA, March 2005.
53. K. Grobe, "Applications of ROADMS and Metro Planes in Control and Regional Networks," *Optical Fiber Communication and the National Fiber Optic Engineers Conference (OFC/NFOEC)'07*, Paper NTuC1, Anaheim, CA, USA, March 2007.
54. P. Roorda and B. Collings, "Evolution to Colorless and Directionless ROADM Architectures," *Optical Fiber Communication and the National Fiber Optic Engineers Conference (OFC/NFOEC)'08*, Paper NWE2, San Diego, CA, USA, February 2008.
55. <http://www.nortel.com/solutions/optical/collateral/nn115940.pdf>, April 2006.
56. K. Guild, "Impairment-Aware Routing for OBS and OPS Networks," *International Conference on Transparent Optical Networks, 2006* 3:61, Nottingham, UK, June 2006.
57. Y. Huang, J. P. Heritage, and B. Mukherjee, "Connection Provisioning with Transmission Impairment Consideration in Optical WDM Networks with High-Speed Channels," *IEEE Journal of Lightwave Technology* 23:982–993, 2005.
58. T. Carpenter, D. Shallcross, J. Gannett, J. Jackel, and A. Von Lehmen, "Maximizing the Transparency Advantage in Optical Networks," *Optical Fiber Communications Conference (OFC)'03* 2:616–617, Atlanta, Ga., March 2003.
59. N. Andriolli, P. Castoldi, J. Cornelias, F. Cugini, G. Junyent, R. Martinez, C. Pinart, L. Vakarenghi, and L. Wosinska, "Challenges and Requirements for Introducing Impairment-Awareness into the Management and Control Planes of ASON/GMPLS WDM Networks," *IEEE Communications Magazine* 44:76–85, 2006.
60. G. P. Agrawal, *Fiber-Optics Communication Systems*, 2nd ed., Wiley, New York, 2002.
61. A. E. Willner and B. Hoanca, "Fixed and Tunable Management of Fiber Chromatic Dispersion," *Optical Fiber Telecommunications IV*, Ivan P. Kaminow and Tingye Li, eds., New York: Academic Press, 2002.
62. L. D. Garrett, "All about Chromatic Dispersion in Dense WDM Optical Fiber Transmission," Invited Short Course, *Optical Fiber Communication Conference (OFC)'01*, Anaheim, Calif., March 2001.
63. C. D. Poole, "Statistical Treatment of Polarization Dispersion in Single-Mode Fiber," *Optics Letters* 13:687–689, 1988.
64. C. D. Poole, "Measurement of Polarization-Mode Dispersion in Single-Mode Fibers with Random Mode Coupling," *Optics Letters* 14:523–525, 1989.
65. Y. Namihira and H. Wakabayashi, "Fiber Length Dependence of Polarization Mode Dispersion Measurement in Long-Length Optical Fibers and Installed Optical Submarine Cables," *Journal of Optical Communications* 12:2–9, 1991.
66. C. D. Poole, R. W. Tkach, A. R. Chraplyvy, and D. A. Fishman, "Fading in Lightwave Systems due to Polarization-Mode Dispersion," *IEEE Photonics Technology Letters* 3:68–70, 1991.
67. A. E. Willner, "Polarization Mode Dispersion: Playing Russian Roulette with Your Network," *Lightwave Magazine* 19:79–82, 2002.
68. V. J. Mazurczyk and J. L. Zyskind, "Polarization Hole Burning in Erbium Doped Fiber Amplifiers," *Lasers and Electro-Optics/Quantum Electronics and Laser Science (CLEO/QELS)'93*, Paper CPD26, Baltimore, Md., May 1993.
69. H. F. Haunstein and H. M. Kallert, "Influence of PMD on the Performance of Optical Transmission Systems in the Presence of PDL," *Conference on Optical Fiber Communication (OFC)'01*, Paper WT4, Anaheim, Calif., March 2001.
70. N. Gisin and B. Huttner, "Combined Effects of Polarization Mode Dispersion and Polarization Dependent Loss," *Optics Communications* 142:119–125, 1997.
71. B. Huttner, C. Geiser, and N. Gisin, "Polarization-Induced Distortions in Optical Fiber Networks with Polarization-Mode Dispersion and Polarization-Dependent Losses," *IEEE Journal of Selected Topics in Quantum Electronics* 6:317–329, 2000.
72. L. S. Yan, Q. Yu, Y. Xie, and A. E. Willner, "Experimental Demonstration of the System Performance Degradation due to the Combined Effect of Polarization Dependent Loss with Polarization Mode Dispersion," *IEEE Photonics Technology Letters* 14:224–226, 2002.

73. L. -S. Yan, Q. Yu, T. Luo, J. E. McGeehan, and A. E. Willner, "Deleterious System Effects due to Low-Frequency Polarization Scrambling in the Presence of Nonnegligible Polarization-Dependent Loss," *IEEE Photonic Technology Letters* **15**:464–466, 2003.
74. G. P. Agrawal, *Nonlinear Fiber Optics*, Academic Press, New York, 1990.
75. R. W. Tkach, A. R. Chraplyvy, F. Forghieri, A. H. Gnauck, and R. M. Derosier, "Four-Photo Mixing and High-Speed WDM System," *IEEE Journal of Lightwave Technology* **13**:841–849, 1995.
76. P. P. Mitra and K. B. Stark, "Nonlinear Limits to the Information Capacity of Optical Fiber Communications," *Nature* **411**:1027–1030, 2001.
77. A. J. Antos and D. K. Smith, "Design and Characterization of Dispersion Compensating Fiber Based on the LP₀₁ Mode," *IEEE Journal of Lightwave Technology* **12**:1739–1745, 1994.
78. A. H. Gnauck, L. D. Garrett, F. Forghieri, V. Gusmeroli, and D. Scarano, "8 × 20 Gb/s 315-km, 8 × 10 Gb/s 480-km WDM Transmission Over Conventional Fiber Using Multiple Broad-Band Fiber Gratings," *IEEE Photonics Technology Letters* **10**:1495–1497, 1998.
79. M. Ibsen, M. K. Durkin, M. J. Cole, and R. I. Laming, "Sinc-Sampled Fiber Bragg Gratings for Identical Multiple Wavelength Operation," *IEEE Photonics Technology Letters* **10**:842–845, 1998.
80. K. -M. Feng, J. -X. Cai, V. Grubsky, D. S. Starodubov, M. I. Hayee, S. Lee, X. Jiang, A. E. Willner, and J. Feinberg, "Dynamic Dispersion Compensation in a 10-Gb/s Optical System Using a Novel Voltage Tuned Nonlinearly Chirped Fiber Bragg Grating," *IEEE Photonics Technology Letters* **11**:373–375, 1999.
81. Benjamin J. Eggleton, John A. Rogers, Paul S. Westbrook, and Thomas A. Strasser, "Electrically Tunable Power Efficient Dispersion Compensating Fiber Bragg Grating," *IEEE Photonics Technology Letters* **11**:854–856, 1999.
82. Y. Painchaud, "Dispersion Compensation Module," <http://www.teraxion.com>, 2008.
83. J. H. Winters and R. D. Gitlin, "Electrical Signal Processing Techniques in Long-Haul Fiber-Optic Systems," *IEEE Transactions on Communications* **38**:1439–1453, 1990.
84. H. Bülow, "Electronic Equalization of Transmission Impairments," *Conference on Optical Fiber Communication (OFC) '02*, Invited Paper, TuE4, Anaheim, Calif, March 2002.
85. B. Franz, D. Rosener, F. Buchali, H. Bulow, and G. Veith, "Adaptive Electronic Feed-Forward Equaliser and Decision Feedback Equaliser For Mitigation of Chromatic Dispersion and PMD in 43 Gbit/s Optical Transmission Systems," *European Conference on Optical Communication (ECOC) '06*, Paper We1.5.1, Cannes, France, September 2006.
86. N. Alic, G. C. Papen, R. E. Saperstein, R. Jiang, C. Marki, Y. Fainman, and S. Radic, "Experimental Demonstration of 10 Gb/s NRZ Extended Dispersion-Limited Reach Over 600 km-SMF Link Without Optical Dispersion Compensation," *Conference on Optical Fiber Communication (OFC) '06*, Paper OWB7, Anaheim, Calif. March 2006.
87. P. Poggiolini, G. Bosco, S. Savory, Y. Benlachar, R. I. Killey, and J. Prat, "1040 km Uncompensated IMDD Transmission Over G.652 Fiber at 10 Gbit/s Using a Reduced-State SQRT-Metric MLSE Receiver," *European Conference on Optical Communication (ECOC) '06*, Postdeadline Paper Th4.4.6., Cannes, France, September 2006.
88. D. McGhan, C. Laperle, A. Savchenko, C. Li, G. Mak, and M. O'Sullivan, "5120-km RZ-DPSK Transmission Over G.652 Fiber at 10 Gb/s without Optical Dispersion Compensation," *IEEE Photonics Technology Letters* **18**:400–402, 2006.
89. A. H. Gnauck, R. W. Tkach, A. R. Chraplyvy, and T. Li, "High-Capacity Optical Transmission Systems," *IEEE Journal of Lightwave Technology* **26**:1032–1045, 2008.
90. C. Laperle, B. Villeneuve, Z. Zhang, D. McGhan, H. Sun, and M. O'Sullivan, "Wavelength Division Multiplexed (WDM) and Polarization Mode Dispersion (PMD) Performance of a Coherent 40 Gbit/s Dual Polarization Quadrature Phase Shift Keying (DP-QPSK) Transceiver," *Conference on Optical Fiber Communication (OFC) '07*, Postdeadline Paper PDP16, Anaheim, Calif., March 2007.
91. G. Charlet, J. Renaudier, H. Mardoyan, O. B. Pardo, F. Cerou, P. Tran, and S. Bigo, "12.8 Tbit/s Transmission of 160 PDM-QPSK (160 × 2 × 40 Gbit/s) Channels with Coherent Detection over 2550 km," *European Conference on Optical Communication (ECOC) '07*, Postdeadline Paper PD1.6., Berlin, Germany, 2007.
92. C. R. S. Fludger, T. Duthel, D. v. d. Borne, C. Schulien, E.-D. Schmidt, T. Wuth, E. D. Man, G. D. Khoe, and H. D. Waardt, "10 × 111 Gbit/s, 50 GHz Spaced, POLMUX-RZ-DQPSK Transmission over 2375 km Employing Coherent Equalisation," *Conference on Optical Fiber Communication (OFC) '07*, Postdeadline Paper PDP22, Anaheim, Calif., USA, March 2007.

93. X. Liu, S. Chandrasekhar, A. H. Gnauck, C. R. Doerr, I. Kang, D. Kilper, L. L. Buhl, and J. Centanni, "DSP-Enabled Compensation of Demodulator Phase Error and Sensitivity Improvement in Direct-Detection 40-Gb/s DQPSK," *European Conference on Optical Communication (ECOC) '06*, Postdeadline Paper Th4.4.5, Cannes, France, 2006.
94. H. B. Killen, *Fiber Optic Communications*, Prentice Hall, Englewood Cliffs, N.J., 1991.
95. A. H. Gnauck and P. J. Winzer, "Optical Phase-Shift-Keyed Transmission," *IEEE Journal of Lightwave Technology* **23**:115–130, 2005.
96. S. R. Chinn, D. M. Boroson, and J. C. Livas, "Sensitivity of Optically Preamplified DPSK Receivers with Fabry-Perot Filters," *IEEE Journal of Lightwave Technology* **14**:370–376, 1996.
97. C. Xu, X. Liu, and L. Mollenauer, "Comparison of Return-to-Zero Phase Shift Keying and On-Off Keying in Long Haul Dispersion Managed Transmissions," *Conference on Optical Fiber Communication (OFC) '03*, Paper ThE3, Atlanta, Ga., March 2003.
98. P. S. Henry, R. A. Linke, and A. H. Gnauck, "Introduction to Lightwave Systems," *Optical Fiber Telecommunications II*, Chap. 21, S. E. Miller and I. P. Kaminow, eds., Academic Press, New York, 1988.
99. J. M. Senior, *Optical Fiber Communications*, Prentice Hall, New York, 1985.
100. S. Betti, G. de Marchis, and E. Iannone, *Coherent Optical Communication Systems*, chaps. 5 and 6, Wiley, New York, 1995.
101. M. S. Roden, *Analog and Digital Communication Systems*, 4th ed., Prentice Hall, Upper Saddle River, N. J., 1995.
102. K. Sato, S. Kuwahara, Y. Miyamoto, K. Murata, and H. Miyazawa, "Carrier-Suppressed Return-to-Zero Pulse Generation Using Mode-Locked Lasers for 40-Gbit/s Transmission," *IEICE Transactions on Communications* **E85-B**:410–415, 2002.
103. K. Yonenaga and S. Kuwano, "Dispersion-Tolerant Optical Transmission System Using Duobinary Transmitter and Binary Receiver," *IEEE Journal of Lightwave Technology* **15**:1530–1537, 1997.
104. K. Yonenaga, M. Yoneyama, Y. Miyamoto, K. Hagimoto, and K. Noguchi, "160 Gbit/s WDM Transmission Experiment Using Four 40 Gbit/s Optical Duobinary Channels," *IEE Electronics Letters* **34**:1506–1507, 1998.
105. T. Ono, Y. Yano, K. Fukuchi, T. Ito, H. Yamazaki, M. Yamaguchi, and K. Emura, "Characteristics of Optical Duobinary Signals in Terabit/s Capacity, High-Spectral Efficiency WDM Systems," *IEEE Journal of Lightwave Technology* **16**:788–797, 1998.
106. A. E. Willner, "Simplified Model of a FSK-to-ASK Direct-Detection System Using a Fabry-Perot Demodulator," *IEEE Photonics Technology Letters* **2**:363–366, 1990.
107. J. H. Sinsky, A. Adamiecki, A. Gnauck, C. Burrus, J. Leuthold, O. Wohlgenuth, and A. Umbach, "A 42.7-Gb/s Integrated Balanced Optical Front End with Record Sensitivity," *Conference on Optical Fiber Communication (OFC) '03*, Postdeadline Paper PD39, Atlanta, Ga., March 2003.
108. P. J. Winzer, A. H. Gnauck, G. Raybon, S. Chandrasekhar, Y. Su, and J. Leuthold, "40-Gb/s Return-to-Zero Alternate-Mark-Inversion (RZ-AMI) Transmission over 2000 km," *IEEE Photonics Technology Letters* **15**:766–768, 2003.
109. T. Hoshida, O. Vassilieva, K. Yamada, S. Choudhary, R. Pecqueur, and H. Kuwahara, "Optimal 40 Gb/s Modulation Formats for Spectrally Efficient Long-Haul DWDM Systems," *IEEE Journal of Lightwave Technology* **20**:1989–1996, 2002.
110. G. Kramer, A. Ashikhmin, A. J. van Wijngaarden, and X. Wei, "Spectral Efficiency of Coded Phase-Shift Keying for Fiber-Optic Communication," *IEEE Journal of Lightwave Technology* **21**:2438–2445, 2003.
111. H. Kim and P. J. Winzer, "Robustness to Laser Frequency Offset in Direct Detection DPSK and DQPSK Systems," *IEEE Journal of Lightwave Technology* **21**:1887–1891, 2003.
112. D. van den Borne, S. L. Jansen, E. Gottwald, P. M. Krummrich, G. D. Khoe, and H. de Waardt, "Spectral Efficiency of Coded Phase-Shift Keying for Fiber-Optic Communication," *IEEE Journal of Lightwave Technology* **25**:222–232, 2007.
113. M. Ohm and J. Speidel, "Quaternary Optical ASK-DPSK and Receivers with Direct Detection," *IEEE Photonics Technology Letters* **15**:159–161, 2003.
114. S. Hayase, N. Kikuchi, K. Sekine, and S. Sasaki, "Proposal of 8-State per Symbol (Binary ASK and QPSK) 30-Gbit/s Optical Modulation/Demodulation Scheme," *European Conference on Optical Communication (ECOC) '03*, Paper Th2.6.4, Rimini, Italy, September 2003.

115. X. Liu, X. Wei, Y. -H. Kao, J. Leuthold, C. R. Doerr, and L. F. Mollenauer, "Quaternary Differential-Phase Amplitude-Shift-Keying for DWDM Transmission," *European Conference on Optical Communication (ECOC) '03*, Paper Th2.6.5, Rimini, Italy, September 2003.
116. P. J. Winzer, and R. -J. Essiambre, "Advanced Modulation Formats for High-Capacity Optical Transport Networks," *IEEE Journal of Lightwave Technology* **24**:4711–4728, 2006.
117. J. G. Proakis, *Digital Communications*, 4th ed., McGraw-Hill, New York, 2001.
118. A. H. Gnauck, P. J. Winzer, S. Chandrasekhar, and C. Dorrer, "Spectrally Efficient (0.8 b/s/Hz) 1-Tb/s (25 × 42.7 Gb/s) RZ-DQPSK Transmission over 28 100-km SSMF Spans with 7 Optical Add/Drops," *Conference on Optical Fiber Communication (OFC) '04*, Paper Th4.4.1, Los Angeles, Calif., February 2004.
119. P. J. Winzer, H. Kogelnik, C. H. Kim, H. Kim, R. M. Jopson, L. E. Nelson, and K. Ramanan, "Receiver Impact on First-Order PMD Outage," *IEEE Photonics Technology Letters* **15**:1482–1484, 2003.
120. E. B. Basch, R. Egorov, S. Gringeri, and S. Elby, "Architectural Tradeoffs for Reconfigurable Dense Wavelength-Division Multiplexing Systems," *IEEE Journal of Selected Topics in Quantum Electronics* **12**:615–626, 2006.
121. A. H. Gnauck, P. J. Winzer, and S. Chandrasekhar, "Hybrid 10/40-G Transmission on a 50-GHz Grid through 2800 km of SSMF and Seven Optical Add-Drops," *IEEE Photonics Technology Letters* **17**:2203–2205, 2005.
122. S. Appathurai, V. Mikhailov, R. I. Killey, and P. Bayvel, "Investigation of the Optimum Alternate-Phase RZ Modulation Format and Its Effectiveness in the Suppression of Intrachannel Nonlinear Distortion in 40-Gbit/s Transmission Over Standard Single-Mode Fiber," *IEEE Journal of Selected Topics in Quantum Electronics* **10**:239–249, 2004.
123. T. Li, "The Impact of Optical Amplifiers on Long-Distance Lightwave Telecommunications," *Proceedings of the IEEE* **81**:1568–1579, 1993.
124. E. L. Goldstein, A. F. Elrefaie, N. Jackman, and S. Zaidi, "Multiwavelength Fiber-Amplifier Cascades in Unidirectional Interoffice Ring Networks," *Conference on Optical Fiber Communication (OFC) '93*, Paper TuJ3, San Jose, Calif., February 1993.
125. J. P. Blondel, A. Pitel, and J. F. Marcero, "Gain-Filtering Stability in Ultralong-Distance Links," *Conference on Optical Fiber Communication (OFC) '93*, paper TuJ3, San Jose, Calif., February 1993.
126. H. Taga, N. Edagawa, Y. Yoshida, S. Yamamoto, and H. Wakabayashi, "IM-DD Four-Channel Transmission Experiment over 1500 km Employing 22 Cascaded Optical Amplifiers," *IEEE Electronics Letters* **29**:485–486, 1993.
127. A. E. Willner and S. -M. Hwang, "Transmission of Many WDM Channels through a Cascade of EDFAs in Long-Distance Link and Ring Networks," *IEEE Journal of Lightwave Technology* **13**:802–816, 1995.
128. I. Kaminow and T. Koch, *Optical fiber Telecommunications IIIB*, Academic Press, San Diego, 1997.
129. A. M. Vengsarkar, P. J. Lemaire, J. B. Judkins, V. Bhatia, T. Erdogan, and J. E. Sipe, "Long-Period Fiber Gratings as Band-Rejection Filters," *IEEE Journal of Lightwave Technology* **14**:58–65, 1996.
130. Y. Sun, J. B. Judkins, A. K. Srivastava, L. Garrett, J. L. Zyskind, J. W. Suloff, C. Wolf, et. al., "Transmission of 32-WDM 10-Gb/s Channels over 640 km Using Broad-Band, Gain-Flattened Erbium-Doped Silica Fiber Amplifiers," *IEEE Photonics Technology Letters* **9**:1652–1654, 1997.
131. P. F. Wysocki, J. B. Judkins, R. P. Espindola, M. Andrejco, and A. M. Vengsarkar, "Broad-Band Erbium-Doped Fiber Amplifier Flattened beyond 40 nm Using Long-Period Grating Filter," *IEEE Photonics Technology Letters* **9**:1343–1345, 1997.
132. J. Y. Pan; M. A. Ali, A.F. Elrefaie, and R. E. Wagner, "Multiwavelength Fiber-Amplifier Cascades with Equalization Employing Mach-Zehnder Optical Filter," *IEEE Photonics Technology Letters* **7**:1501–1503, 1995.
133. B. Nagaraju, M. C. Paul, M. Pal, A. Pal, R. K. Varshney, B. P. Pal, S. K. Bhadra, G. Monnom, and B. Dussardier, "Design and Realization of an Inherently Gain Flattened Erbium Doped Fiber Amplifier," *Conference on Lasers and Electro-Optics (CLEO) '08*, Paper JTuA86, San Jose, Calif., May 2008.
134. E. Desurvire, R. Giles, and J. Simpson, "Gain Saturation Effects in High-Speed, Multi Channel Erbium-Doped Fiber Amplifiers at 1.53 μm ," *IEEE Journal of Lightwave Technology* **7**:2095–2104, 1989.
135. D. C. Kilper, S. Chandrasekhar, and C. A. White, "Transient Gain Dynamics of Cascaded Erbium Doped Fiber Amplifiers with Re-Configured Channel Loading," *Conference on Optical Fiber Communication (OFC) '06*, Paper OTuK6, Anaheim, Calif., March 2006.

136. J. Zyskind, Y. Sun, A. Srivastava, J. Sulhoff, A. Lucero, C. Wolf, and R. Tkach, "Fast Power Transients in Optically Amplified Multiwavelength Optical Networks," *Conference on Optical Fiber Communication (OFC) '96*, Paper PD-31, San Jose, Calif., February 1996.
137. M. Hayee and A. Willner, "Fiber Transmission Penalties due to EDFA Power Transients Resulting from Fiber Nonlinearity and ASE Noise in Add/Drop Multiplexed WDM Networks," *IEEE Photonics Technology Letters* **11**:889–891, 1999.
138. J. -X. Cai, K. -M. Feng, and A. E. Willner, "Simultaneous Compensation of Fast Add/Drop Power-Transients and Equalization of Inter-Channel Power Differentials for Robust WDM Systems with EDFAs," *Conference on Optical Amplifiers and Their Applications*, Paper MC6, Victoria, Canada, July 1997.
139. K. Motoshima, L. Leba, D. Chen, M. Downs, T. Li, and E. Desurvire, "Dynamic Compensation of Transient Gain Saturation in Erbium-Doped Fiber Amplifiers by Pump Feedback Control," *IEEE Photonics Technology Letters* **5**:1423–1426, 1993.
140. A. Srivastava, J. Zyskind, Y. Sun, J. Ellson, G. Newsome, R. Tkach, A. Chraplyvy, J. Sulhoff, T. Strasser, C. Wolf, and J. Pedrazzani, "Fast-Link Control Protection of Surviving Channels in Multiwavelength Optical Networks," *IEEE Photonics Technology Letters* **9**:1667–1669, 1997.
141. G. Luo, J. Zyskind, Y. Sun, A. Srivastava, J. Sulhoff, and M. Ali, "Performance Degradation of All-Optical Gain-Clamped EDFA's due to Relaxation-Oscillations and Spectral-Hole Burning in Amplified WDM Networks," *IEEE Photonics Technology Letters* **9**:1346–1348, 1997.
142. J. Cai, K. Feng, X. Chen, and A. Willner, "Experimental Demonstration of Dynamic High-Speed Equalization of Three WDM Channels Using Acousto-Optic Modulators and a Wavelength Demultiplexer," *IEEE Photonics Technology Letters* **9**:678–680, 1997.
143. S. Huang, X. Zou, A. Willner, Z. Bao, and D. Smith, "Experimental Demonstration of Active Equalization and ASE Suppression of Three 2.5 Gb/s WDM-Network Channels over 2500 km Using AOTF as Transmission Filters," *IEEE Photonics Technology Letters* **9**:389–391, 1997.
144. D. Starodubov, V. Grubsky, J. Feinberg, J. Cai, K. Feng, and A. Willner, "Novel Fiber Amplitude Modulators for Dynamic Channel Power Equalization in WDM Systems," *Conference on Optical Fiber Communication (OFC) '98*, Paper PD-8, San Jose, Calif., February 1998.
145. J. Ford and J. Walker, "Dynamic Spectral Power Equalization Using Micro-Opto-Mechanics," *IEEE Photonics Technology Letters* **10**:1440–1442, 1998.
146. C. Doerr, C. Joyner, and L. Stulz, "Integrated WDM Dynamic Power Equalizer with Low Insertion Loss," *IEEE Photonics Technology Letters* **10**:1443–1445, 1998.
147. C. Doerr, C. Joyner, and L. Stulz, "16-Band Integrated Dynamic Gain Equalization Filter with Less than 2.8-dB Insertion Loss," *IEEE Photonics Technology Letters* **14**:334–336, 2002.
148. C. Fludger, A. Maroney, and N. Jolley, "An Analysis of the Improvements in OSNR from Distributed Raman Amplifiers Using Modern Transmission Fibres," *Conference on Optical Fiber Communication (OFC) '00*, Paper FF2, Baltimore, Md., March 2000.
149. J. Bromage, "Raman Amplification for Fiber Communications Systems," *IEEE Journal of Lightwave Technology* **22**:79–93, 2004.
150. Y. Emori, K. Tanaka, and S. Namiki, "100 nm Bandwidth Flat-Gain Raman Amplifiers Pumped and Gain-Equalised by 12 Wavelength-Channel WDM Laser Diode Unit," *IEEE Electronics Letters* **35**:1355–1356, 1999.
151. S. Namiki and Y. Emori, "Ultrabroad-Band Raman Amplifiers Pumped and Gain-Equalized by Wavelength-Division-Multiplexed Highpower Laser Diodes," *IEEE Journal of Selected Topics in Quantum Electronics* **7**:3–16, 2001.
152. S. Namiki, K. Seo, N. Tsukiji, and S. Shikii, "Challenges of Raman Amplification," *Proceedings of IEEE* **94**:1024–1035, 2006.
153. B. Zhu, L. E. Nelson, S. Stulz, A. H. Gnauck, C. Doerr, J. Leuthold, L. Grüner-Nielsen, M. O. Pedersen, J. Kim, R. Lingle Jr., et al, "6.4-Tb/s (160 × 42.7 Gb/s) Transmission with 0.8 bit/s/Hz Spectral Efficiency over 32 × 100 km of Fiber Using CSRZ-DPSK Format," *Conference on Optical Fiber Communication (OFC) '03*, Paper PD-19, Atlanta, Ga., March 2003.

This page intentionally left blank.

DO NOT DUPLICATE

SOLITONS IN OPTICAL FIBER COMMUNICATION SYSTEMS

Pavel V. Mamyshev

*Bell Laboratories—Lucent Technologies
Holmdel, New Jersey*

22.1 INTRODUCTION

To understand why optical solitons are needed in optical fiber communication systems, we should consider the problems that limit the distance and/or capacity of optical data transmission. A fiber-optic transmission line consists of a transmitter and a receiver connected with each other by a transmission optical fiber. Optical fibers inevitably have chromatic dispersion, losses (attenuation of the signal), and nonlinearity. Dispersion and nonlinearity can lead to the distortion of the signal. Because the optical receiver has a finite sensitivity, the signal should have a high-enough level to achieve error-free performance of the system. On the other hand, by increasing the signal level, one also increases the nonlinear effects in the fiber. To compensate for the fiber losses in a long distance transmission, one has to periodically install optical amplifiers along the transmission line. By doing this, a new source of errors is introduced into the system—an amplifier spontaneous emission noise. (Note that even ideal optical amplifiers inevitably introduce spontaneous emission noise.) The amount of noise increases with the transmission distance (with the number of amplifiers). To keep the signal-to-noise ratio (SNR) high enough for the error-free system performance, one has to increase the signal level and hence the potential problems caused by the nonlinear effects. Note that the nonlinear effects are proportional to the product of the signal power P and the transmission distance L , and both of these multipliers increase with the distance. Summarizing, we can say that all the problems—dispersion, noise, and nonlinearity—grow with the transmission distance. The problems also increase when the transmission bit rate (speed) increases. It is important to emphasize that it is very difficult to deal with the signal distortions when the nonlinearity is involved, because the nonlinearity can couple all the detrimental effects together [nonlinearity, dispersion, noise, polarization mode dispersion (i.e., random birefringence of the fiber), polarization-dependent loss/gain, etc.]. That happens when the nonlinear effects are out of control. The idea of soliton transmission is to guide the nonlinearity to the desired direction and use it for your benefit. When soliton pulses are used as an information carrier, the effects of dispersion and nonlinearity balance (or compensate) each other and thus don't degrade the signal quality with the propagation distance. In such a regime, the pulses propagate through the fiber without changing their spectral and temporal shapes. This mutual compensation of dispersion and nonlinear effects takes place continuously with the distance

in the case of “classical” solitons and periodically with the so-called dispersion map length in the case of dispersion-managed solitons. In addition, because of the unique features of optical solitons, soliton transmission can help to solve other problems of data transmission, like polarization mode dispersion. Also, when used with frequency guiding filters (sliding guiding filters in particular), the soliton systems provide continuous all-optical regeneration of the signal suppressing the detrimental effects of the noise and reducing the penalties associated with wavelength-division multiplexed (WDM) transmission. Because the soliton data looks essentially the same at different distances along the transmission, the soliton type of transmission is especially attractive for all-optical data networking. Moreover, because of the high quality of the pulses and return-to-zero (RZ) nature of the data, the soliton data is suitable for all-optical processing.

22.2 NATURE OF THE CLASSICAL SOLITON

Signal propagation in optical fibers is governed by the nonlinear Schroedinger equation (NSE) for the complex envelope of the electric field of the signal.^{1–3} This equation describes the combined action of the self-phase modulation and dispersion effects, which play the major role in the signal evolution in most practical cases. Additional linear and nonlinear effects can be added to the modified NSE.⁴ Mathematically, one can say that solitons are stable solutions of NSE.^{1,2} In this paper, however, we will give a qualitative physical description of the soliton regimes of pulse propagation, trying to avoid mathematics as much as possible.

Consider first the effect of dispersion. An optical pulse of width τ has a finite spectral bandwidth $BW \approx 1/\tau$. When the pulse is transform limited, or unchirped, all the spectral components have the same phase. In time domain, one can say that all the spectral components overlap in time, or sit on top of each other (see Fig. 1). Because of the dispersion, different spectral components propagate in the fiber with different group velocities, V_{gr} . As a result of the dispersion action alone, the initial

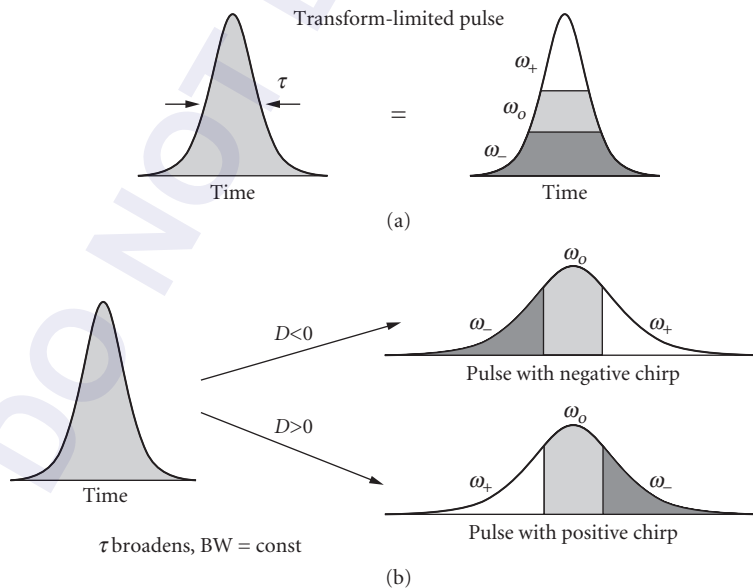


FIGURE 1 (a) Transform-limited pulse: all spectral components of the pulse “sit” on top of each other. (b) Effect of group velocity dispersion on a transform-limited pulse.

unchirped pulse broadens and gets chirped (frequency modulated). The sign of the chirp depends on the sign of the fiber group velocity dispersion (see Fig. 1).

$$D = d \left(\frac{1}{V_{gr}} \right) / d\lambda \quad (1)$$

(λ is the light wavelength). A characteristic fiber length called the *dispersion length*, at which the pulse broadens by a factor sqrt (2), is determined both by the fiber dispersion and the pulse width:

$$z_d = \frac{2\pi c 0.322 \tau^2}{\lambda^2 D} \quad (2)$$

(c is the speed of light). Note that the pulse spectral bandwidth remains unchanged because the dispersion is a linear effect.

Consider now the nonlinear effect of self-phase modulation (SPM).⁵ Due to the Kerr effect, the fiber refractive index depends on the signal intensity, $n(I) = n_0 + n_2 I$, where n_2 is the nonlinear refractive index and intensity is $I = P/A$, P is the signal power and A is the fiber effective cross-section mode area. During a pulse propagation through the fiber, different parts of the pulse acquire different values of the nonlinear phase shift: $\phi(t) = 2\pi/\lambda n_2 I(t)L$. Here $I(t)$ is the intensity pulse shape in time domain and L is the transmission distance. This time-dependent nonlinear phase shift means that different parts of the pulse experience different frequency shifts:

$$\delta\omega(t) = \frac{d\phi}{dt} = -\frac{2\pi}{\lambda} n_2 L \frac{dI(t)}{dt} \quad (3)$$

As one can see, the frequency shift is determined by the time derivative of the pulse shape. Because the nonlinear refractive index in silica-based fibers is positive, the self-phase modulation effect always shifts the front edge of the pulse to the “red” spectral region (downshift in frequency), and the trailing edge of the pulse to the “blue” spectral region (upshift in frequency). This means that an initially unchirped pulse spectrally broadens and gets negatively chirped (Fig. 2). A characteristic fiber length called the *nonlinear length*, at which the pulse spectrally broadens by a factor of two, is

$$z_{NL} = \left(\frac{2\pi}{\lambda} n_2 I_0 \right)^{-1} \quad (4)$$

Note that, when acting alone, SPM does not change the temporal intensity profile of the pulse.

As it was mentioned earlier, when under no control, both SPM and dispersion may be very harmful for the data transmission distorting considerably the spectral and temporal characteristics of the signal. Consider now how to control these effects by achieving the soliton regime of data transmission when the combined action of these effects results in a stable propagation of data pulses without changing their spectral and temporal envelopes.

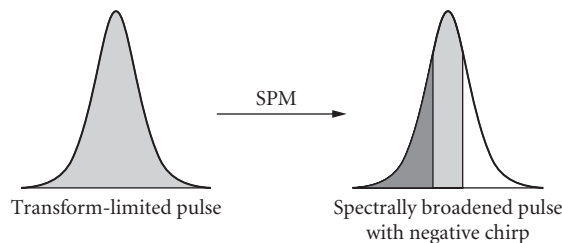


FIGURE 2 Effect of self-phase modulation on a transform-limited pulse.

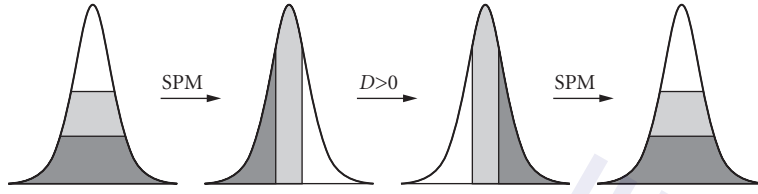


FIGURE 3 Qualitative explanation of classical soliton. Combined action of dispersion and nonlinearity (self-phase modulation) results in a stable pulse propagation with constant spectral and temporal widths. See text.

In our qualitative consideration, consider the combined action of dispersion and nonlinearity (SPM) as an alternative sequence of actions of dispersion and nonlinearity. Assume that we start with a chirp-free pulse (see Fig. 3). The self-phase modulation broadens the pulse spectrum and produces a negative frequency chirp: The front edge of the pulse becomes red-shifted, and the trailing edge becomes blue-shifted. When positive GVD is then applied to this chirped pulse, the red spectral components are delayed in time with respect to the blue ones. If the right amount of dispersion is applied, the sign of the pulse chirp can be reversed to negative: The blue spectral components shift in time to the front pulse edge, while the red spectral components move to the trailing edge. When the nonlinearity is applied again, it shifts the frequency of the front edge to the red spectral region and upshifts the frequency of the trailing edge. That means that the blue front edge becomes green again, the red trailing edge also becomes green, and the pulse spectrum bandwidth narrows to its original width. The described regime of soliton propagation is achieved when the nonlinear and dispersion effect compensate each other exactly. In reality, the effects of dispersion and SPM act simultaneously, so that the pulse spectral and temporal widths stay constant with the distance, and the only net effect is a (constant within the entire pulse) phase shift of 0.5 rad per dispersion length of propagation.⁶ The condition of the soliton regime is equality of the nonlinear and dispersion lengths: $z_d = z_{NL}$. One can rewrite this expression to find a relationship between the soliton peak power, pulse width, and fiber dispersion:

$$P_0 = \frac{\lambda^3 DA}{0.3224 \pi^2 c n_2 \tau^2} \quad (5)$$

Here, P_0 is the soliton peak power and τ is the soliton FWHM. Soliton pulses have a sech^2 form. Note that as it follows from our previous consideration, classical soliton propagation in fibers requires a positive sign of the fiber's dispersion, D (assuming that n_2 is positive). Consider a numerical example. For a pulse of width $\tau = 20$ ps propagating in a fiber with $D = 0.5$ ps nm⁻¹ km⁻¹, fiber cross-section mode area $A = 50 \mu\text{m}^2$, $\lambda = 1.55 \mu\text{m}$, and typical value of $n_2 = 2.6 \text{ cm}^2/\text{W}$, one can find the soliton peak power is 2.4 mW. The dispersion length is $z_d = 200$ km in this case.

22.3 PROPERTIES OF SOLITONS

The most important property of optical solitons is their robustness.⁶⁻²⁰ Consider what robustness means from a practical point of view. When a pulse is injected into the fiber, the pulse does not have to have the exact soliton shape and parameters [Eq. (5)] to propagate as a soliton. As long as the input parameters are not too far from the optimum, during the nonlinear propagation the pulse “readjusts” itself, shaping into a soliton and shedding off nonsoliton components. For example, an unchirped pulse of width τ will be reshaped into a single soliton as long as its input power P is greater than $P_0/4$ and less than $2.25P_0$. Here, P_0 is the soliton power determined by Eq. (5).³

Solitons are also robust with respect to the variations of the pulse energy and of the fiber parameters along the transmission line. As long as these variations are fast enough (period of perturbations is much smaller than the soliton dispersion length z_d), the soliton “feels” only the average values of these parameters. This feature is extremely important for practical systems. In particular, it makes it possible to use solitons in long distance transmission systems where fiber losses are periodically compensated by lumped amplifiers. As long as the amplifier spacing is much less than the soliton dispersion length $L_{\text{amp}} \ll z_d$, classical solitons work very well in these systems. Note that all soliton perturbations result in a loss of some part of the soliton energy, which is radiated into dispersive waves.

Consider now a case of slow variations of parameters along the transmission when a characteristic length at which a fiber parameter (or pulse energy) changes considerably is much longer than the soliton dispersion length. Soliton parameters follow adiabatically these changes. That means that all the parameters in Eq. (5) can be considered as distance dependent, and Eq. (5) remains valid. It can be rewritten in the following form:

$$\tau(z) = \text{const} \frac{D(z)A(z)}{P(z)\tau(z)} = \text{const} \frac{D(z)A(z)}{\text{energy}(z)} \quad (6)$$

One can derive many important consequences from this equation.^{13–22} One example would be the pulse broadening (and spectral narrowing) in a fiber with loss [assuming $D(z)$ and $A(z)$ are constant].^{13–15} Note that the soliton broadening can be used in repeaterless data transmission systems when high-input signal power is required.¹⁵ On the other hand, one can get a pulse compression in a fiber with adiabatic gain. Similar effects can be obtained by changing the fiber dispersion and/or mode area along the length. For example, adiabatic soliton compression can be obtained in a fiber with slowly decreasing dispersion (dispersion-tapered fiber).^{16–22}

It is important to emphasize that the adiabatic soliton propagation does not necessarily require that each of these parameters—pulse energy, fiber dispersion, and mode area—changes adiabatically with the distance, as long as the whole expression, $[D(z)A(z)]/[\text{energy}(z)]$ changes adiabatically with the distance. For example, soliton propagation in a dispersion-tapered fiber with losses is equivalent to transmission in a lossless, constant-dispersion fiber if the dispersion decreases with the same rate with the distance as the pulse energy [i.e., if $D(z)/\text{energy}(z) = \text{const}$]. Note that this is true no matter what the fiber loss and the pulse width are.

So far, we’ve been discussing a single pulse propagation. In communication systems, one has to deal with streams of pulses. When two or more soliton pulses propagate in the fiber at the same wavelength, they can interact with each other: Tails from one soliton pulse may overlap with the other pulse. Due to the cross-phase modulation effect, this overlap leads to the frequency shifts of the interacting solitons. The signs of the frequency shifts are opposite for the two solitons. Through the fiber dispersion, the frequency changes result in the changes of the soliton group velocities. The strength of the interaction decreases very fast with the soliton separation and for most practical applications can be considered to be negligible when the separation is 4 to 5 times greater than the soliton pulse width τ .^{23,24} The character of interaction depends on the mutual optical phases of the solitons: When they are the same, the solitons attract to each other; when they are out of phase, the solitons repel from each other; when the phase difference is $\pi/2$, the solitons do not interact.

22.4 CLASSICAL SOLITON TRANSMISSION SYSTEMS

The soliton properties described earlier determine the engineering rules for designing the soliton-based transmission systems. First, to make sure that every individual pulse is stable in the transmission line with constant fiber dispersion and loss periodically compensated by lump amplifiers, the amplifier spacing L_{amp} should be much smaller than the soliton dispersion length z_d . To avoid considerable pulse-to-pulse interaction, the minimum distance between adjacent pulses should be $T \geq 4\tau$, where $1/T$ is the transmission bit rate and τ is the soliton pulse width. The pulse power determined from Eq. (5) should be considered as a path-average power P_{av} . If the signal energy

decreases with the distance in the fiber spans between the amplifiers as $\exp(-\gamma z)$ (here, γ is the loss rate), the path-average power is related to the pulse power at the output of each amplifier (input to the fiber span) P_{in} , as

$$P_0 = P_{\text{in}} \frac{1 - \exp(-\gamma L_{\text{amp}})}{|\gamma| L_{\text{amp}}} \quad (7)$$

Here, L_{amp} is the amplifier spacing. As it was stated earlier, the dispersion and nonlinear effects “compensate” each other in the soliton regime of transmission, so that the pulses propagate practically without changing their temporal and spectral shapes. As long as the length scale of perturbations of the transmission parameters is much shorter than the soliton dispersion length, the pulses “feel” only the average parameters. Note, however, that perturbations may lead to shedding of dispersive waves by solitons.¹²

There are two main sources of errors in the soliton transmission systems: fluctuations of the pulse energies and fluctuations of the pulse arrival times.²⁵ The origin of the energy fluctuations is the same as in the other types of systems—spontaneous emission noise generated by the amplifiers. Each amplifier contributes a noise with a spectral density (power per unit bandwidth):

$$P_v = (G-1) n_{\text{sp}} h\nu \quad (8)$$

Here, G is the power gain of the amplifier, $h\nu$ is the photon energy, and $n_{\text{sp}} \geq 1$ is the spontaneous emission factor that characterizes the quality of the amplifier. In the best case, when the amplifier is highly inverted, n_{sp} is close to unity. In a broadband transmission system (i.e., without in-line spectral filters), when the lumped amplifiers compensate exactly for the fiber loss, the noise grows linearly with the distance (with the number of amplifiers). At the output of a transmission line of length L , the path-averaged spectral density is

$$P_{\text{vav}} = |\gamma| L n_{\text{sp}} h\nu F(G) \quad (9)$$

Here, function $F(G)$ describes the penalty one has to pay for having high-gain amplifiers (or long amplifier spacing):

$$F(G) = \frac{(G-1)^2}{G \ln^2 G} \quad (10)$$

The penalty function has its minimum [$F(G) = 1$] in the case of distributed amplification (when $G \rightarrow 1$) and grows with G . The SNR at the output of transmission should be high enough to have error-free transmission. Note that the noise spectral density P_v has units of energy. It is also the noise energy received in any time T in a spectral bandwidth $1/T$. That is why P_v is also called the *equipartition energy*. To have the error probability less than 10^{-9} and 10^{-15} , the ratio of the pulse energy to the equipartition energy should be, correspondently, 100 and 160. For example, consider a transmission system with the average loss of 0.21 dB/km, $n_{\text{sp}} = 1.5$, amplifier spacing of 50 km. The minimum pulse energy at the input of each fiber span to have the error probability less than 10^{-9} in such a system of length $L = 5000$ km is 20 fJ, and for $L = 10,000$ km, it is 40 fJ.

Another type of error in the soliton systems is the fluctuation in the pulse arrival times, or timing jitter. The timing jitter can be caused by several factors. The adjacent pulse-to-pulse interactions can cause the pulses to shift in time. As we have stated earlier, interaction problems can be practically eliminated by spacing the solitons in time by more than 4 or 5 of their width. A very important source of the timing jitter is the spontaneous emission noise. Every time the noise is added to the signal, it modulates the carrier frequencies of the solitons at random. The chromatic dispersion of the fiber then converts these frequency variations in a variation of the pulses’ arrival times. This effect

is known as the *Gordon-Haus effect*.^{6,26} The variance of the timing jitter produced by the Gordon-Haus effect is

$$\sigma_{\text{GH}}^2 \approx 0.2n_2hn_{\text{sp}}F(G)\frac{|\gamma|}{A}\frac{D}{\tau}L^3 \quad (11)$$

An error occurs when a pulse arrives outside of the acceptance time window W of the detection system (this window is usually slightly less than the bit slot, T). To have the error probability less than 10^{-9} , the acceptance window should be greater than 12 standard deviations of the timing jitter:

$$W \geq 12\sigma_{\text{GH}} \quad (12)$$

The Gordon-Haus jitter limits the maximum bit rate and transmission distance. As one can see from Eq. (11), the jitter increases very fast with the distance; it also increases when τ decreases. Another factor that limits the maximum transmission distance is that σ_{GH}^2 is proportional to the pulse energy [because the pulse energy is proportional to (D/τ)], and long-distance transmission systems should have high-enough pulse energies to keep the SNR high. Consider a numerical example, $L=9000$ km, $\tau=20$ ps, $n_{\text{sp}}=1.4$, $\gamma=-0.048$ km⁻¹, amplifier spacing = 30 km, $D=0.5$ ps/(nm⁻¹ km⁻¹), $A=50$ μm². Equation (11) then gives the standard deviation of the Gordon-Haus timing jitter $\sigma=11.7$ ps. As one can see, according to Eq. (12), this jitter is too high for 10 Gbit/s transmission ($1/T=100$ ps) to be error-free, because $12\sigma_{\text{GH}} > 1/T$ in this case.

Another source of the timing jitter is the acoustic interaction of pulses.²⁷⁻³⁰ Due to the electrostriction effect in the fiber, each propagating pulse generates an acoustic wave in the fiber. Other pulses experience the refractive index change caused by the acoustic wave. The resultant frequency changes of the pulses lead, through the effect of the fiber chromatic dispersion, to the fluctuation in the arrival times. The acoustic effect causes a “long-range” interaction: Pulses separated by a few nanoseconds can interact through this effect. One can estimate the acoustic timing jitter from the following simplified equation:

$$\sigma_a \approx 4.3\frac{D^2}{\tau}(R-0.99)^{1/2}L^2 \quad (13)$$

Here, standard deviation σ_a is in picoseconds; dispersion D is in picoseconds per nanometer per kilometer; the bit rate $R=1/T$, is in gigabits per second; and the distance L , is in megameters. Equation (13) also assumes the fiber mode area of $A=50$ μm². The acoustic jitter increases with the bit rate, and it has even stronger dependence on the distance than the Gordon-Haus jitter.

As it follows from the previous considerations, the timing jitter can impose severe limitations on the distance and capacity of the systems, and it has to be controlled.

22.5 FREQUENCY-GUIDING FILTERS

The Gordon-Haus and acoustic timing jitters originate from the frequency fluctuations of the pulses. That means that by controlling the frequency of the solitons, one can control the timing jitter as well. The frequency control can be done by periodically inserting narrowband filters (so-called frequency-guiding filters) along the transmission line, usually at the amplifier locations.^{31,32} If, for some reason, the center frequency of a soliton is shifted from the filter peak, the filter-induced differential loss across the pulse spectrum “pushes” the pulse frequency back to the filter peak. As a result, the pulse spectrum returns back to the filter peak in a characteristic damping length Δ . If the damping length is considerably less than the transmission distance L the guiding filters dramatically reduce the timing jitter. To calculate the timing jitter in a filtered system, one should replace L^3 by

$3L\Delta^2$ in Eq. (11), and L^2 in Eq. (13) should be replaced by $2L\Delta$. Then, we get the following expression for the Gordon-Haus jitter:

$$\sigma_{\text{GH},f}^2 \approx 0.6n_2hn_{\text{sp}}F(G)\frac{|\gamma|}{A}\frac{D}{\tau}L\Delta^2 \quad (14)$$

The damping properties of the guiding filters are determined mainly by the curvature of the filter response in the neighborhood of its peak. That means that shallow Fabry-Perot etalon filters can be used as the guiding filters. Fabry-Perot etalon filters have multiple peaks, and different peaks can be used for different WDM channels. The ability of the guiding filters to control the frequency jitter is determined both by the filter characteristics and by the soliton spectral bandwidth. In the case of Fabry-Perot filters with the intensity mirror reflectivity R , and the free spectral range (FSR), the damping length is

$$\Delta = 0.483(\tau \text{FSR})^2 \frac{(1-R)^2}{R} L_f \quad (15)$$

Here, L_f is the spacing between the guiding filters; usually, L_f equals the amplifier spacing L_{amp} .

Note that the Gordon-Haus and acoustic jitters are not specific for soliton transmission only. Any kind of transmission systems, including so-called linear transmission, are subject to these effects. However, the guiding filters can be used in the soliton systems only. Every time a pulse passes through a guiding filter, its spectrum narrows. Solitons can quickly recover their bandwidth through the fiber nonlinearity, whereas for a linear transmission the filter action continuously destroys the signal.

Note that even a more effective reduction of the timing jitter can be achieved if, in addition to the frequency-guiding filters, an amplitude and/or phase modulation at the bit rate is applied to the signal periodically with the distance. "Error-free" transmission over practically unlimited distances can be achieved in this case (1 million kilometers at 10 Gbit/s has been demonstrated).^{33,34} Nevertheless, this technique is not passive, high-speed electronics is involved, and the clock recovery is required each time the modulation is applied. Also, in the case of WDM transmission, all WDM channels have to be demultiplexed before the modulation and then multiplexed back afterward; each channel has to have its own clock recovery and modulator. As one can see, this technique shares many drawbacks of the electronic regeneration schemes.

The frequency-guiding filters can dramatically reduce the timing jitter in the systems. At the same time, though, in some cases they can introduce additional problems. Every time a soliton passes through the filter, it loses some energy. To compensate for this loss, the amplifiers should provide an additional (excess) gain. Under this condition, the spontaneous emission noise and other nonsoliton components with the spectrum in the neighborhood of the filter peak experience exponential growth with the distance, which reduces the SNR and can lead to the soliton instabilities. As a result, one has to use weak-enough filters to reduce the excess gain. In practice, the filter strength is chosen to minimize the total penalty from the timing jitter and the excess gain.

22.6 SLIDING FREQUENCY-GUIDING FILTERS

As one can see, the excess gain prevents one from taking a full advantage of guiding filters. By using the sliding frequency-guiding filters,³⁵ one can essentially eliminate the problems associated with the excess gain. The trick is very simple: The transmission peak of each guiding filter is shifted in frequency with respect to the peak of the previous filter, so that the center frequency slides with the distance with the rate of $f' = df/dz$. Solitons, thanks to the nonlinearity, can follow the filters and slide in frequency with the distance. But all unwanted linear radiation (e.g., spontaneous emission noise, nonsoliton components shedded from the solitons, etc.) cannot slide and eventually is killed by the filters. The sliding allows one to use strong guiding filters and even to reduce the amount of

noise at the output of transmission in comparison with the broadband (no guiding filters) case. The maximum filter strength³⁶ and maximum sliding rate³⁵ are determined by the soliton stability. The error-free transmission of 10 Gbit/s signal over 40,000 km and 20 Gbit/s over 14,000 km was demonstrated with the sliding frequency-guiding filters technique.^{37,38}

It is important to emphasize that by introducing the sliding frequency-guiding filters into the transmission line, one converts this transmission line into an effective, all-optical passive regenerator (compatible with WDM). Solitons with only one energy (and pulse width) can propagate stably in such a transmission line. The parameters of the transmission line (the filter strength, excess gain, fiber dispersion, and mode area) determine the unique parameters of these stable solitons. The system is opaque for a low-intensity radiation (noise, for example). However, if the pulse parameters at the input of the transmission line are not too far from the optimum soliton parameters, the transmission line reshapes the pulse into the soliton of that line. Note, again, that the parameters of the resultant soliton do not depend on the input pulse parameters, but only on the parameters of the transmission line. Note also that all nonsoliton components generated during the pulse reshaping are absorbed by the filters. That means, in particular, that the transmission line removes the energy fluctuations from the input data signal.⁶ Note that the damping length for the energy fluctuations is close to the frequency damping length of Eq. (15). A very impressive demonstration of regenerative properties of a transmission line with the frequency-guiding filters is the conversion of a nonreturn-to-zero (NRZ) data signal (frequency modulated at the bit rate) into a clean soliton data signal.³⁹ Another important consequence of the regenerative properties of a transmission line with the frequency-guiding filters is the ability to self-equalize the energies of different channels in WDM transmission.⁴⁰ Negative feedback provided by frequency-guiding filters locks the energies of individual soliton channels to values that do not change with distance, even in the face of considerable variation in amplifier gain among the different channels. The equilibrium values of the energies are independent of the input values. All these benefits of sliding frequency-guiding filters are extremely valuable for practical systems. Additional benefits of guiding filters for WDM systems will be discussed later.

22.7 WAVELENGTH DIVISION MULTIPLEXING

Due to the fiber chromatic dispersion, pulses from different WDM channels propagate with different group velocities and collide with each other.⁴¹ Consider a collision of two solitons propagating at different wavelengths (different channels). When the pulses are initially separated and the fast soliton (the soliton at shorter wavelength, with higher group velocity) is behind the slow one, the fast soliton eventually overtakes and passes through the slow soliton. An important parameter of the soliton collision is the collision length L_{coll} , the fiber length at which the solitons overlap with each other. If we let the collision begin and end with the overlap of the pulses at half power points, then the collision length is

$$L_{\text{coll}} = \frac{2\tau}{D\Delta\lambda} \quad (16)$$

Here, $\Delta\lambda$ is the solitons wavelengths difference. Due to the effect of cross-phase modulation, the solitons shift each other's carrier frequency during the collision. The frequency shifts for the two solitons are equal in amplitudes (if the pulse widths are equal) and have opposite signs. During the first half of collision, the fast accelerates even faster (carrier frequency increases), while the slow soliton slows down. The maximum frequency excursion δf_{max} of the solitons is achieved in the middle of the collision, when the pulses completely overlap with each other:

$$\delta f_{\text{max}} = \pm \frac{1}{3\pi^2 \cdot 0.322 \Delta f \tau^2} = \pm \frac{1.18n_2 \varepsilon}{A \tau D \lambda \Delta \lambda} \quad (17)$$

Here, $\Delta f = -c \Delta \lambda / \lambda^2$ is the frequency separation between the solitons, and $\varepsilon = 1.13 P_0 \tau$ is the soliton energy. In the middle of collision, the accelerations of the solitons change their signs. As a result, the

frequency shifts in the second half of collision undo the frequency shifts of the first half, so that the soliton frequency shifts go back to zero when the collision is complete. This is a very important and beneficial feature for practical applications. The only residual effect of complete collision in a lossless fiber is the time displacements of the solitons:

$$\delta t_{cc} = \pm \frac{0.1786}{\Delta f^2 \tau} = \pm \frac{2\epsilon n_2 \lambda}{cDA \Delta \lambda^2} \quad (18)$$

The symmetry of the collision can be broken if the collision takes place in a transmission line with loss and lumped amplification. For example, if the collision length L_{coll} is shorter than the amplifier spacing L_{amp} , and the center of collision coincides with the amplifier location, the pulses intensities are low in the first half of collision and high in the second half. As a result, the first half of collision is practically linear. The soliton frequency shifts acquired in the first half of collision are very small and insufficient to compensate for the frequency shifts of opposite signs acquired by the pulses in the second half of collision. This results in nonzero residual frequency shifts. Note that similar effects take place when there is a discontinuity in the value of the fiber dispersion as a function of distance. In this case, if a discontinuity takes place in the middle of collision, one half of the collision is fast (where D is higher) and the other half is slow. The result is nonzero residual frequency shifts. Nonzero residual frequency shifts lead, through the dispersion of the rest of the transmission fiber, to variations in the pulses arrival time at the output of transmission. Nevertheless, if the collision length is much longer than the amplifier spacing and of the characteristic length of the dispersion variations in the fiber, the residual soliton frequency shifts are zero, just like in a lossless uniform fiber. In practice, the residual frequency shifts are essentially zero as long as the following condition is satisfied:⁴¹

$$L_{\text{coll}} \geq 2L_{\text{amp}} \quad (19)$$

Another important case is so-called half-collisions (or partial collisions) at the input of the transmission.⁴² These collisions take place if solitons from different channels overlap at the transmission input. These collisions result in residual frequency shifts of δf_{max} and the following pulse timing shifts δt_{pc} at the output of transmission of length L :

$$\delta t_{\text{pc}} \approx \delta f_{\text{max}} \frac{\lambda^2}{c} D(L - L_{\text{coll}}/4) = \pm \frac{1.18\epsilon n_2 \lambda}{c\tau A \Delta \lambda} (L - L_{\text{coll}}/4) \quad (20)$$

One can avoid half-collisions by staggering the pulse positions of the WDM channels at the transmission input.

Consider now the time shifts caused by all complete collisions. Consider a two-channel transmission, where each channel has a $1/T$ bit rate. The distance between subsequent collisions is

$$l_{\text{coll}} = \frac{T}{D \Delta \lambda} \quad (21)$$

The maximum number of collisions that each pulse can experience is L/l_{coll} . This means that the maximum time shift caused by all complete collisions is

$$\delta t_{\Sigma cc} \approx \delta t_{cc} L/l_{\text{coll}} = \pm \frac{2\epsilon n_2 \lambda}{cTA \Delta \lambda} L \quad (22)$$

It is interesting to note that $\delta t_{\Sigma cc}$ does not depend on the fiber dispersion. Note also that Eq. (22) describes the worst case when the pulse experiences the maximum number of possible collisions. Consider a numerical example. For a two-channel transmission, 10 Gbit/s each ($T = 100$ ps), pulse energy ($\epsilon = 50$ fJ), channel wavelength separation ($\Delta \lambda = 0.6$ nm), fiber mode area ($A = 50 \mu\text{m}^2$ and $L = 10$ mm), we find $\delta t_{\Sigma cc} = 45$ ps. Note that this timing shift can be reduced by increasing the

channel separation. Another way to reduce the channel-to-channel interaction by a factor of two is to have these channels orthogonally polarized to each other. In WDM transmission, with many channels, one has to add timing shifts caused by all other channels. Note, however, that as one can see from Eq. (22), the maximum penalty comes from the nearest neighboring channels.

As one can see, soliton collisions introduce additional jitter to the pulse arrival time, which can lead to considerable transmission penalties. As we saw earlier, the frequency-guiding filters are very effective in suppressing the Gordon-Haus and acoustic jitters. They can also be very effective in suppressing the timing jitter induced by WDM collisions. In the ideal case of parabolical filters and the collision length being much longer than the filter spacing $L_{\text{coll}} \gg L_f$, the filters make the residual time shift of a complete collision δt_{cc} exactly zero. They also considerably reduce the timing jitter associated with asymmetrical collisions and half-collisions. Note that for the guiding filters to work effectively in suppressing the collision penalties, the collision length should be at least a few times greater than the filter spacing. Note also that real filters, such as etalon filters, do not always perform as good as ideal parabolic filters. This is true especially when large-frequency excursions of solitons are involved, because the curvature of a shallow etalon filter response reduces with the deviation of the frequency from the filter peak. In any case, filters do a very good job in suppressing the timing jitter in WDM systems.

Consider now another potential problem in WDM transmission, which is the four-wave mixing. During the soliton collisions, the four-wave mixing spectral sidebands are generated. Nevertheless, in the case of a lossless, constant-dispersion fiber, these sidebands exist only during the collision, and when the collision is complete, the energy from the sidebands regenerates back into the solitons. That is why it was considered for a long time that the four-wave mixing should not be a problem in soliton systems. But this is true only in the case of a transmission in a lossless fiber. In the case of lossy fiber and periodical amplification, these perturbations can lead to the effect of the pseudo-phase-matched (or resonance) four-wave mixing.⁴³ The pseudo-phase-matched four-wave mixing lead to the soliton energy loss to the spectral sidebands and to a timing jitter (we called that effect an *extended Gordon-Haus effect*).⁴³ The effect can be so strong that even sliding frequency-guiding filters are not effective enough to suppress it. The solution to this problem is to use dispersion-tapered fiber spans. As we have discussed earlier, soliton propagation in the condition:

$$\frac{D(z)A(z)}{\text{Energy}(z)} = \text{const} \quad (23)$$

is identical to the case of lossless, constant-dispersion fiber. That means that the fiber dispersion in the spans between the amplifiers should decrease with the same rate as the signal energy. In the case of lumped amplifiers, this is the exponential decay with the distance. Note that the dispersion-tapered spans solve not just the four-wave mixing problem. By making the soliton transmission perturbation-free, they lift the requirements to have the amplifier spacing much shorter than the soliton dispersion length. The collisions remain symmetrical even when the collision length is shorter than the amplifier spacing. (Note, however, that the dispersion-tapered fiber spans do not lift the requirement to have guiding filter spacing as short as possible in comparison with the collision length and with the dispersion length.) The dispersion-tapered fiber spans can be made with the present technology.²² Stepwise approximation of the exact exponential taper made of fiber pieces of constant dispersion can also be used.⁴³ It was shown numerically and experimentally that by using fiber spans with only a few steps one can dramatically improve the quality of transmission.^{44,45} In the experiment, each fiber span was dispersion tapered typically in three or four steps, the path-average dispersion value was $0.5 \pm 0.05 \text{ ps nm}^{-1} \text{ km}^{-1}$ at 1557 nm. The use of dispersion-tapered fiber spans together with sliding frequency-guiding filters allowed transmission of eight 10-Gbit/s channels with the channel spacing, $\Delta\lambda = 0.6 \text{ nm}$, over more than 9000 km. The maximum number of channels in this experiment was limited by the dispersion slope, $dD/d\lambda$, which was about $0.07 \text{ ps nm}^{-2} \text{ km}^{-1}$. Because of the dispersion slope, different WDM channels experience different values of dispersion. As a result, not only the path average dispersion changes with the wavelength, but the dispersion tapering has exponential behavior only in a vicinity of one particular wavelength in the center of the transmission band. Wavelength-division multiplexed channels located far from that wavelength propagate in far from

the optimal conditions. One solution to the problem is to use dispersion-flattened fibers (i.e., fibers with $dD/d\lambda = 0$). Unfortunately, these types of fibers are not commercially available at this time. This and some other problems of classical soliton transmission can be solved by using dispersion-managed soliton transmission.⁴⁶⁻⁶³

22.8 DISPERSION-MANAGED SOLITONS

In the dispersion-managed (DM) soliton transmission, the transmission line consists of the fiber spans with alternating signs of the dispersion. Let the positive and negative dispersion spans of the map have lengths and dispersions, L_+ , D_+ and L_- , D_- , respectively. Then, the path-average dispersion D_{av} is

$$D_{av} = (D_+L_+ + L_-D_-)/L_{map} \quad (24)$$

Here, L_{map} is the length of the dispersion map:

$$L_{map} = L_+ + L_- \quad (25)$$

Like in the case of classical soliton, during the DM soliton propagation, the dispersion and nonlinear effects cancel each other. The difference is that in the classical case, this cancellation takes place continuously, whereas in the DM case, it takes place periodically with the period of the dispersion map length L_{map} . The strength of the DM is characterized by a parameter S , which is determined as^{47,50,52}

$$S = \frac{\lambda^2}{2\pi c} \frac{(D_+ - D_{av})L_+ - (D_- - D_{av})L_-}{\tau^2} \quad (26)$$

The absolute values of the local dispersion are usually much greater than the path average dispersion: $|D_+|, |D_-| \gg |D_{av}|$. As one can see from Eq. (26), the strength of the map is proportional to the number of the local dispersion lengths of the pulse in the map length: $S \approx L_{map}/z_{d,local}$. The shape of the DM solitons are close to Gaussian. A very important feature of DM solitons is the so-called power enhancement. Depending on the strength of the map, the pulse energy of DM solitons, ϵ_{DM} , is greater than that of classical solitons, ϵ_0 , propagating in a fiber with constant dispersion, $D = D_{av}$.^{47,50}

$$\epsilon_{DM} \approx \epsilon_0(1 + 0.7S^2) \quad (27)$$

Note that this equation assumes lossless fiber. The power enhancement effect is very important for practical applications. It provides an extra degree of freedom in the system design by giving the possibility to change the pulse energy while keeping the path-average fiber dispersion constant. In particular, because DM solitons can have adequate pulse energy (to have a high-enough SNR) at or near zero path average dispersion, timing jitter from the Gordon-Haus and acoustic effects is greatly reduced (e.g., the variance of the Gordon-Haus jitter, σ^2 , scales almost as $1/\epsilon_{DM}$).⁴⁹ Single-channel high-bit-rate DM soliton transmission over long distances with weak guiding filters and without guiding filters was experimentally demonstrated.^{46,51}

Dispersion-managed soliton transmission is possible not only in transmission lines with positive dispersion, $D_{av} > 0$, but also in the case of $D_{av} = 0$ and even $D_{av} < 0$.⁵² To understand this, consider qualitatively the DM soliton propagation (Fig. 4). Locally, the dispersive effects are always stronger than the nonlinear effect (i.e., the local dispersion length is much shorter than the nonlinear length). In the zero approximation, the pulse propagation in the map is almost linear. Let's call the middle of

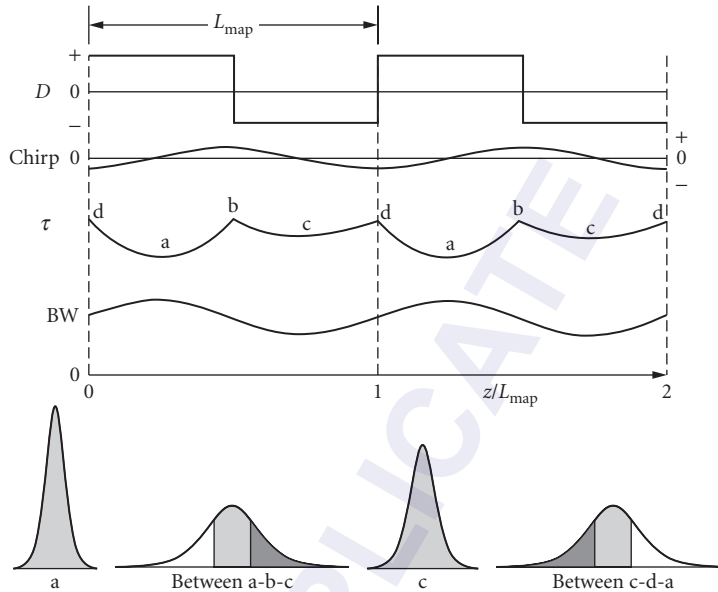


FIGURE 4 Qualitative description of dispersion-managed (DM) soliton transmission. Distance evolution of the fiber dispersion $[D(z)]$, pulse chirp, pulse width $[\tau(z)]$, and pulse bandwidth $[BW(z)]$. Evolution of the pulse shape in different fiber sections is shown in the bottom.

the positive D sections “point a,” the middle of the negative sections “point c,” transitions between positive and negative sections “point b,” and transitions between negative and positive sections “point d.” The chirp-free (minimum pulse width) positions of the pulse are in the middle of the positive- and negative- D sections (points a and c). The pulse chirp is positive between points a, b, and c (see Fig. 4). That means that the high-frequency (blue) spectral components of the pulse are at the front edge of the pulse, and the low-frequency (red) components are at the trailing edge. In the section c-d-a, the pulse chirp is negative. The action of the nonlinear SPM effect always downshifts in frequency the front edge of the pulse and up shifts in frequency the trailing edge of the pulse. That means that the nonlinearity decreases the spectral bandwidth of positively chirped pulses (section a-b-c) and increases the spectral bandwidth of negatively chirped pulses (section c-d-a). This results in the spectral bandwidth behavior also shown in Fig. 4: The maximum spectral bandwidth is achieved in the chirp-free point in the positive section, whereas the minimum spectral bandwidth is achieved in the chirp-free point in the negative section. The condition for the pulses to be DM solitons is that the nonlinear phase shift is compensated by the dispersion-induced phase shift over the dispersion map length. That requires that $\int DBW^2 dz > 0$ (here, BW is the pulse spectral bandwidth). Note that in the case of classical solitons, when spectral bandwidth is constant, this expression means that dispersion D must be positive. In the DM case, however, the pulse bandwidth is wider in the positive- D section than in the negative- D section. As a result, the integral can be positive, even when $D_{av} = \int D dz / L_{map}$ is zero or negative. Note that the spectral bandwidth oscillations explain also the effect of power enhancement of DM solitons.

Consider interaction of adjacent pulses in DM systems.⁵⁴ The parameter that determines the strength of the interaction is the ratio τ/T (here, τ is the pulse width and T is the spacing between adjacent pulses). As in the case of classical soliton transmission, the cross-phase modulation effect (XPM) shifts the frequencies of the interacting pulses, Δf_{XPM} , which, in turn, results in timing jitter at the output of the transmission. As it was discussed earlier, the classical soliton interaction increases

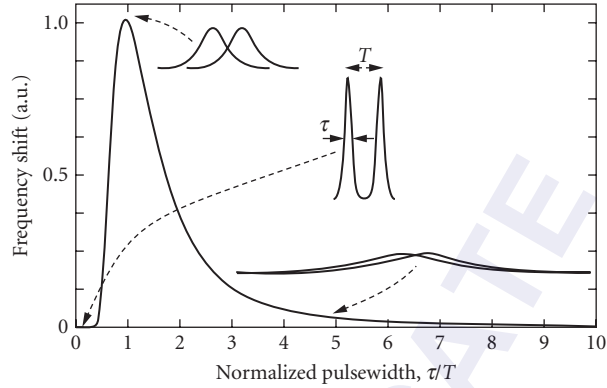


FIGURE 5 Dimensionless function $\Phi(\tau/T)$ describing the XPM-induced frequency shift of two interacting chirped Gaussian pulses as a function of the pulse width normalized to the pulse separation.

very quickly with τ/T . To avoid interaction-induced penalties in classical soliton transmission systems, the pulses should not overlap significantly with each other: τ/T should be less than 0.2 to 0.3. In the DM case, the situation is different. The pulse width in the DM case oscillates with the distance $\tau(z)$; that means that the interaction also changes with distance. Also, because the pulses are highly chirped when they are significantly overlapped with each other, the sign of the interaction is essentially independent of the mutual phases of the pulses. Cross-phase modulation always shifts the leading pulse to the red spectral region, and the trailing pulse shifts to the blue spectral region. The XPM-induced frequency shifts of interacting solitons per unit distance is

$$\frac{d\Delta f_{\text{XPM}}}{dz} \approx \pm 0.15 \frac{2\pi n_2 \mathcal{E}}{\lambda T^2 A} \Phi(\tau/T) \quad (28)$$

The minus sign in Eq. (28) corresponds to the leading pulse, and the plus sign corresponds to the trailing pulse. Numerically calculated dimensionless function, $\Phi(\tau/T)$, is shown in Fig. 5. As it follows from Eq. (28), $\Phi(\tau/T)$ describes the strength of the XPM-induced interaction of the pulses as a function of the degree of the pulse overlap. One can see that the interaction is very small when τ/T is smaller than 0.4 (i.e., when the pulses barely overlap), which is similar to the classical soliton propagation. The strength of the interaction of DM solitons also increases with τ/T , but only in the region $0 < \tau/T < 1$. In fact, the interaction reaches its maximum at $\tau/T \approx 1$ and then decreases and becomes very small again when $\tau/T \gg 1$ (i.e., when the pulses overlap nearly completely). There are two reasons for such an interesting behavior at $\tau/T \gg 1$. The XPM-induced frequency shift is proportional to the time derivative of the interacting pulse's intensity, and the pulse derivative reduces with the pulse broadening. Also, when the pulses nearly completely overlap, the sign of the derivative changes across the region of overlap so that the net effect tends to be canceled out.

Based on Eq. (28) and Fig. 5, one can distinguish three main regimes of data transmission in DM systems. In all these regimes, the minimum pulse width is, of course, less than the bit slot, T . The regimes differ from each other by the maximum pulse breathing with the distance. In the first, “non-pulse-overlapped,” regime, adjacent pulses barely overlap during most of the transmission, so that the pulse interaction is not a problem in this case. This is the most stable regime of transmission. In the “partially-pulse-overlapped” regime, the adjacent pulses spend a considerable portion of the transmission being partially overlapped [$\tau(z)$ being around T]. Cross-phase modulation causes the frequency and timing jitter in this case. In the third, “pulse-overlapped,” regime, the adjacent pulses are almost completely overlapped with each other during most of the transmission [$\tau_{\min}(L_{\text{map}}/z_{d,\text{local}}) \gg T$].

The XPM-induced pulse-to-pulse interaction is greatly reduced in this case in comparison with the previous one. The main limiting factor for this regime of transmission is the intrachannel four-wave mixing taking place during strong overlap of adjacent pulses.⁵⁴ The intrachannel four-wave mixing leads to the amplitude fluctuations of the pulses and “ghost” pulse generation in the “zero” slots of the data stream.

22.9 WAVELENGTH-DIVISION MULTIPLEXED DISPERSION-MANAGED SOLITON TRANSMISSION

One of the advantages of DM transmission over classical soliton transmission is that the local dispersion can be very high ($|D_+|, |D_-| \gg |D_{av}|$), which efficiently suppresses the four-wave mixing from soliton-soliton collisions in WDM. Consider the timing jitter induced by collisions in the non-pulse-overlapped DM transmission. The character of the pulse collisions in DM systems is quite different from the case of a transmission line with uniform dispersion: In the former, the alternating sign of the high local dispersion causes the colliding solitons to move rapidly back and forth with respect to each other, with the net motion determined by D_{av} .^{56–59} Because of this rapid breathing of the distance between the pulses, each net collision actually consists of many fast or “mini” collisions. The net collision length can be estimated as⁵⁹

$$L_{\text{coll}} \approx \frac{2\tau}{D_{av}\Delta\lambda} + \frac{(D_+ - D_{av})L_+}{D_{av}} \approx \frac{2\tau}{D_{av}\Delta\lambda} + \frac{\tau_{\text{eff}}}{D_{av}\Delta\lambda} \quad (29)$$

Here, τ is the minimum (unchirped) pulse width. Here, we also defined the quantity $\tau_{\text{eff}} \equiv L_+ D_+ \Delta\lambda$, which plays the role of an effective pulse width. For strong dispersion management, τ_{eff} is usually much bigger than τ . Thus, L_{coll} becomes almost independent of $\Delta\lambda$ and much longer than it is for classical solitons subject to the same D_{av} . As a result, the residual frequency shift caused by complete pulse collisions tends to become negligibly small for transmission using strong maps.⁵⁸ The maximum frequency excursion during the DM soliton collision is⁵⁹

$$\delta f_{\text{max}} \approx \pm \frac{2n_2\epsilon}{L_+ D_+ A D_{av} \lambda \Delta\lambda^2} = \pm \frac{2n_2\epsilon}{A D_{av} \lambda \Delta\lambda \tau_{\text{eff}}} \quad (30)$$

Now, we can estimate the time shift of the solitons per complete collision:

$$\delta t_{\text{cc}} \approx D_{av} \lambda^2 / c \int \delta f dz \approx \alpha D_{av} L_{\text{coll}} \delta f_{\text{max}} \lambda^2 / c \approx \pm \alpha \frac{2n_2\epsilon\lambda}{c A D_{av} \Delta\lambda^2} \quad (31)$$

Here, $\alpha \leq 1$ is a numerical coefficient that takes into account the particular shape of the frequency shift as a function of distance. Consider now the time shifts caused by all collisions. In a two-channel transmission, the distance between subsequent collisions is $l_{\text{coll}} = T / (D_{av} \Delta\lambda)$. The maximum number of complete collisions at the transmission distance L is $(L - L_{\text{coll}}) / l_{\text{coll}}$ (we assume that $L > L_{\text{coll}}$), and the number of incomplete collisions at the end of transmission is $L_{\text{coll}} / l_{\text{coll}}$. The timing shift caused by all these collisions can be estimated as

$$\delta t_{\Sigma c} \approx \delta t_{\text{cc}} (L - L_{\text{coll}} / 2) / l_{\text{coll}} = \pm \alpha \frac{2n_2\epsilon\lambda}{c A T \Delta\lambda} (L - L_{\text{coll}} / 2) \quad (32)$$

Consider the problem of initial partial collisions. As it was discussed earlier for the case of classical solitons, initial partial collisions can be a serious problem by introducing large timing jitter at the output of transmission. On the other hand, for the classical case, one could avoid the half-collisions

by staggering the pulse positions of the WDM channels at the transmission input. The situation is very different for the DM case. In the DM case, the collision length is usually longer than the distance between subsequent collisions (i.e., $L_{\text{coll}} > l_{\text{coll}}$). Thus, a pulse can collide *simultaneously* with several pulses of another channel. The maximum number of such simultaneous collisions is $N_{\text{sc}} \approx L_{\text{coll}}/l_{\text{coll}} = 2\tau/T + [(D_+ - D_{\text{av}})L_+ \Delta\lambda]/T$. Note that N_{sc} increases when the channel spacing $\Delta\lambda$ increases. The fact that the collision length is greater than the distance between collisions also means that initial partial collisions are inevitable in DM systems. Moreover, depending on the data pattern in the interacting channel, each pulse can experience up to N_{sc} initial partial collisions with that channel (not just one as in the classical case). As a consequence, the residual frequency shifts can be bigger than δf_{max} . The total time shift caused by the initial partial collisions at distance $L > L_{\text{coll}}$ can be estimated as

$$\delta\tau_{\text{pc}} = \beta \delta f_{\text{max}} N_{\text{sc}} (L - L_{\text{coll}}/2) D_{\text{av}} \lambda^2 / c \approx \pm \beta \frac{2n_2 \varepsilon \lambda}{cAT \Delta\lambda} (L - L_{\text{coll}}/2) \quad (33)$$

Here, $\beta \leq 1$ is a numerical coefficient that takes into account the particular shape of the frequency shift as a function of distance for a single collision.

Equations (32) and (33) assume that the transmission distance is greater than the collision length. When $L > L_{\text{coll}}$, these equations should be replaced by

$$\delta t_{\Sigma, \text{pc}} \approx (\alpha, \beta) D_{\text{av}} \delta f_{\text{max}} \frac{\lambda^2}{c} \frac{L^2}{2l_{\text{coll}}} \approx \pm (\alpha, \beta) \frac{n_2 \varepsilon \lambda}{cAT \Delta\lambda} \frac{L^2}{L_{\text{coll}}} \quad (34)$$

Note that the signs of the timing shifts caused by initial partial collisions and by complete collisions are opposite. Thus, the maximum (worst-case) spread of the pulse arriving times caused by pulse collisions in the two-channel WDM transmission is described by:

$$\delta t_{\text{max}} = |\delta t_{\text{pc}}| + |\delta t_{\Sigma}| \quad (35)$$

In a WDM transmission with more than two channels, one has to add contributions to the time shift from all the channels. Note that the biggest contribution makes the nearest neighboring channels, because the time shift is inversely proportional to the channel spacing, $\Delta\lambda$. Now, we can summarize the results of Eqs. (32) through (35) as follows. When $L > L_{\text{coll}}$ [Eqs. (32) to (33)], corresponding to very long distance transmission, δt_{max} increases linearly with the distance and almost independently of the path-average dispersion D_{av} . When $L < L_{\text{coll}}$ [Eq. (34)], which corresponds to short-distance transmission and/or very low path-average dispersion, δt_{max} increases quadratically with the distance and in proportion to D_{av} . Note also that the WDM data transmission at near zero path-averaged dispersion $D_{\text{av}} = 0$, may not be desirable, because $L_{\text{coll}} \rightarrow \infty$ and frequency excursions $\delta f_{\text{max}} \rightarrow \infty$ when $D \rightarrow 0$ [see Eq. (30)]. Thus, even though Eq. (34) predicts the time shift to be zero when D_{av} is exactly zero, the frequency shifts of the solitons can be unacceptably large and Eq. (34) may be no longer valid. There are also practical difficulties in making maps with $D_{\text{av}} < 0.1$ ps nm⁻¹ km⁻¹ over the wide spectral range required for dense WDM transmission.

It is interesting to compare these results with the results for the case of classical solitons [Eqs. (17) to (22)]. The time shifts per complete collisions [Eqs. (18) and (31)] are about the same, the time shifts from all initial partial collisions [Eqs. (20) and (33)] are also close to each other. The total maximum time shifts from all collisions are also close to each other for the case of long distance transmission. That means that, similar to the classical case, one has to control the collision-induced timing jitter when it becomes too large. As it was discussed earlier, the sliding frequency-guiding filters are very effective in suppressing the timing jitter. Because the collision length in DM systems is much longer than in classical systems, and, at the same time, it is almost independent of the channel wavelength separation, the requirement that the collision length is much greater than the filter spacing, $L_{\text{coll}} \gg L_p$, is easy to meet. As a result, the guiding filters suppress the timing jitter in DM systems even more

effective than in classical soliton systems. The fact that the frequency excursions during collisions are much smaller in DM case, also makes the filters to work more effectively.

As we have discussed previously, many important features of DM solitons come from the fact that the soliton spectral bandwidth oscillates with the distance. That is why guiding filters alter the dispersion management itself and give an additional degree of freedom in the system design.⁶⁰ Note also that the position of the filters in the dispersion map can change the soliton stability in some cases.⁶¹ It should also be noted that because of the weak dependence of the DM soliton spectral bandwidth on the soliton pulse energy, the energy fluctuations damping length provided by the guided filters is considerably longer than the frequency damping length.⁶² This is the price one has to pay for many advantages of DM solitons. From the practical point of view, the most important advantage is the flexibility in system design and freedom in choosing the transmission fibers. For example, one can upgrade existing systems by providing an appropriate dispersion compensation with dispersion compensation fibers or with lumped dispersion compensators (fiber Bragg gratings, for example). The biggest advantage of DM systems is the possibility to design dispersion maps with essentially zero dispersion slope of the path-average dispersion, $dD_{av}/d\lambda$, by combining commercially available fibers with different signs of dispersion and dispersion slopes. (Note that it was a nonzero dispersion slope that limited the maximum number of channels in classical soliton long distance WDM transmission.) This was demonstrated in the experiment where almost flat average dispersion $D_{av} = 0.3 \text{ ps nm}^{-1} \text{ km}^{-1}$ was achieved by combining standard, dispersion-compensating, and True-Wave (Lucent nonzero dispersion-shifted) fibers.⁶³ By using sliding frequency-guiding filters and this dispersion map, “error-free” DM soliton transmission of twenty-seven 10-Gbit/s WDM channels was achieved over more than 9000 km without using forward error correction. It was shown that once the error-free transmission with about 10 channels is achieved, adding additional channels practically does not change performance of the system. (This is because, for each channel, only the nearest neighboring channels degrade its performance.) The maximum number of WDM channels in this experiment was limited only by the power and bandwidth of optical amplifiers used in the experiment. One can expect that the number of channels can be increased by a few times if more powerful and broader-bandwidth amplifiers are used.

22.10 CONCLUSION

We considered the basic principles of soliton transmission systems. The main idea of the “soliton philosophy” is to put under control, balance, and even to extract the maximum benefits from otherwise detrimental effects of the fiber dispersion and nonlinearity. The “soliton approach” is to make transmission systems intrinsically stable. Soliton technology is a very rapidly developing area of science and engineering, which promises a big change in the functionality and capacity of optical data transmission and networking.

22.11 REFERENCES

1. V. E. Zaharov and A. B. Shabat, “Exact Theory of Two Dimensional Self Focusing and One-Dimensional Self-Modulation of Waves in Nonlinear Media,” *Zh. Eksp. Teor. Fiz.* **61**:118–134 (1971) [*Sov. Phys. JETP* **34**:62–69 (1972)].
2. A. Hasegawa and F. D. Tappert, “Transmission of Stationary Nonlinear Optical Pulses in Dispersive Dielectric Fibers. I. Anomalous Dispersion,” *Appl. Phys. Lett.* **23**:142–144 (1973).
3. J. Satsuma and N. Yajima, “Initial Value Problem of One-Dimensional Self-Modulation of Nonlinear Waves in Dispersive Media,” *Prog. Theor. Phys. Suppl.* **55**:284–306 (1980).
4. P. V. Mamyshev and S. V. Chernikov, “Ultrashort Pulse Propagation in Optical Fibers,” *Opt. Lett.* **15**: 1076–1078 (1990).

5. R. H. Stolen, in *Optical Fiber Telecommunications*, S. E. Miller and H. E. Chynoweth (eds.), Academic Press, New York, 1979, Chap. 5.
6. L. F. Mollenauer, J. P. Gordon, and P. V. Mamyshev, "Solitons in High Bit Rate, Long Distance Transmission," in *Optical Fiber Telecommunications III*, Academic Press, New York, 1997, Chap. 12.
7. L. F. Mollenauer, J. P. Gordon, and M. N. Islam, "Soliton Propagation in Long Fibers with Periodically Compensated Loss," *IEEE J. Quantum Electron.* **QE-22**:157 (1986).
8. L. F. Mollenauer, M. J. Neubelt, S. G. Evangelides, J. P. Gordon, J. R. Simpson, and L. G. Cohen, "Experimental Study of Soliton Transmission Over More Than 10,000 km in Dispersion Shifted Fiber," *Opt. Lett.* **15**:1203 (1990).
9. L. F. Mollenauer, S. G. Evangelides, and H. A. Haus, "Long Distance Soliton Propagation Using Lumped Amplifiers and Dispersion Shifted Fiber," *J. Lightwave Technol.* **9**:194 (1991).
10. K. J. Blow and N. J. Doran, "Average Soliton Dynamics and the Operation of Soliton Systems with Lumped Amplifiers," *Photonics Tech. Lett.* **3**:369 (1991).
11. A. Hasegawa and Y. Kodama, "Guiding-Center Soliton in Optical Fibers," *Opt. Lett.* **15**:1443 (1990).
12. G. P. Gordon, "Dispersive Perturbations of Solitons of the Nonlinear Schroedinger Equation," *JOSA B* **9**:91–97 (1992).
13. A. Hasegawa and Y. Kodama, "Signal Transmission by Optical Solitons in Monomode Fiber," *Proc. IEEE* **69**:1145 (1981).
14. K. J. Blow and N. J. Doran, "Solitons in Optical Communications," *IEEE J. Quantum Electron.*, **QE-19**:1883 (1982).
15. P. B. Hansen, H. A. Haus, T. C. Damen, J. Shah, P. V. Mamyshev, and R. H. Stolen, "Application of Soliton Spreading in Optical Transmission," Dig. ECOC, Vol. 3, Paper WeC3.4, pp. 3.109–3.112, Oslo, Norway, September 1996.
16. K. Tajima, "Compensation of Soliton Broadening in Nonlinear Optical Fibers with Loss," *Opt. Lett.* **12**:54 (1987).
17. H. H. Kuehl, "Solitons on an Axially Nonuniform Optical Fiber," *J. Opt. Soc. Am. B* **5**:709–713 (1988).
18. E. M. Dianov, L. M. Ivanov, P. V. Mamyshev, and A. M. Prokhorov, "High-Quality Femtosecond Fundamental Soliton Compression in Optical Fibers with Varying Dispersion," *Topical Meeting on Nonlinear Guided-Wave Phenomena: Physics and Applications*, 1989, Technical Digest Series, vol. 2, OSA, Washington, D.C., 1989, pp. 157–160, paper FA-5.
19. P. V. Mamyshev, "Generation and Compression of Femtosecond Solitons in Optical Fibers," *Bull. Acad. Sci. USSR, Phys. Ser.*, **55**(2):374–381 (1991) [*Izv. Acad. Nauk, Ser. Phys.* **55**(2):374–381 (1991)].
20. S. V. Chernikov and P. V. Mamyshev, "Femtosecond Soliton Propagation in Fibers with Slowly Decreasing Dispersion," *J. Opt. Soc. Am. B* **8**(8):1633–1641 (1991).
21. P. V. Mamyshev, S. V. Chernikov, and E. M. Dianov, "Generation of Fundamental Soliton Trains for High-Bit-Rate Optical Fiber Communication Lines," *IEEE J. of Quantum Electron.* **27**(10):2347–2355 (1991).
22. V. A. Bogatyrev, M. M. Bubnov, E. M. Dianov, et al., "Single-Mode Fiber with Chromatic Dispersion Varying along the Length," *IEEE J. of Lightwave Technology* **LT-9**(5):561–566 (1991).
23. V. I. Karpman and V. V. Solov'ev, "A Perturbation Approach to the Two-Soliton System," *Physica D* **3**:487–502 (1981).
24. J. P. Gordon, "Interaction Forces among Solitons in Optical Fibers," *Opt. Lett.* **8**:596–598 (1983).
25. J. P. Gordon and L. F. Mollenauer, "Effects of Fiber Nonlinearities and Amplifier Spacing on Ultra Long Distance Transmission," *J. Lightwave Technol.* **9**:170 (1991).
26. J. P. Gordon and H. A. Haus, "Random Walk of Coherently Amplified Solitons in Optical Fiber," *Opt. Lett.* **11**:665 (1986).
27. K. Smith and L. F. Mollenauer, "Experimental Observation of Soliton Interaction over Long Fiber Paths: Discovery of a Long-Range Interaction," *Opt. Lett.* **14**:1284 (1989).
28. E. M. Dianov, A. V. Luchnikov, A. N. Pilipetskii, and A. N. Starodumov, "Electrostriction Mechanism of Soliton Interaction in Optical Fibers," *Opt. Lett.* **15**:314 (1990).
29. E. M. Dianov, A. V. Luchnikov, A. N. Pilipetskii, and A. M. Prokhorov, "Long-Range Interaction of Solitons in Ultra-Long Communication Systems," *Soviet Lightwave Communications* **1**:235 (1991).
30. E. M. Dianov, A. V. Luchnikov, A. N. Pilipetskii, and A. M. Prokhorov "Long-Range Interaction of Picosecond Solitons Through Excitation of Acoustic Waves in Optical Fibers," *Appl. Phys. B* **54**:175 (1992).

31. A. Mecozzi, J. D. Moores, H. A. Haus, and Y. Lai, "Soliton Transmission Control," *Opt. Lett.* **16**:1841 (1991).
32. Y. Kodama and A. Hasegawa, "Generation of Asymptotically Stable Optical Solitons and Suppression of the Gordon-Haus Effect," *Opt. Lett.* **17**:31 (1992).
33. M. Nakazawa, E. Yamada, H. Kubota, and K. Suzuki, "10 Gbit/s Soliton Transmission over One Million Kilometers," *Electron. Lett.* **27**:1270 (1991).
34. T. Widdowson and A. D. Ellis, "20 Gbit/s Soliton Transmission over 125 Mm," *Electron. Lett.* **30**:1866 (1994).
35. L. F. Mollenauer, J. P. Gordon, and S. G. Evangelides, "The Sliding-Frequency Guiding Filter: An Improved Form of Soliton Jitter Control," *Opt. Lett.* **17**:1575 (1992).
36. P. V. Mamyshev and L. F. Mollenauer, "Stability of Soliton Propagation with Sliding Frequency Guiding Filters," *Opt. Lett.* **19**:2083 (1994).
37. L. F. Mollenauer, P. V. Mamyshev, and M. J. Neubelt, "Measurement of Timing Jitter in Soliton Transmission at 10 Gbits/s and Achievement of 375 Gbits/s-Mm, Error-Free, at 12.5 and 15 Gbits/s," *Opt. Lett.* **19**:704 (1994).
38. D. LeGuen, F. Fave, R. Boittin, J. Debeau, F. Devaux, M. Henry, C. Thebault, and T. Georges, "Demonstration of Sliding-Filter-Controlled Soliton Transmission at 20 Gbit/s over 14 Mm," *Electron. Lett.* **31**:301 (1995).
39. P. V. Mamyshev and L. F. Mollenauer, "NRZ-to-Soliton Data Conversion by a Filtered Transmission Line," in *Optical Fiber Communication Conference OFC-95*, Vol. 8, 1995 OSA Technical Digest Series, OSA, Washington, D.C., 1995, Paper FB2, pp. 302–303.
40. P. V. Mamyshev and L. F. Mollenauer, "WDM Channel Energy Self-Equalization in a Soliton Transmission Line Using Guiding Filters," *Opt. Lett.* **21**(20):1658–1660 (1996).
41. L. F. Mollenauer, S. G. Evangelides, and J. P. Gordon, "Wavelength Division Multiplexing with Solitons in Ultra Long Distance Transmission Using Lumped Amplifiers," *J. Lightwave Technol.* **9**:362 (1991).
42. P. A. Andrekson, N. A. Olsson, J. R. Simpson, T. Tanbun-ek, R. A. Logan, P. C. Becker, and K. W. Wecht, *Electron. Lett.* **26**:1499 (1990).
43. P. V. Mamyshev, and L. F. Mollenauer, "Pseudo-Phase-Matched Four-Wave Mixing in Soliton WDM Transmission," *Opt. Lett.* **21**:396 (1996).
44. L. F. Mollenauer, P. V. Mamyshev, and M. J. Neubelt, "Demonstration of Soliton WDM Transmission at 6 and 7×10 GBit/s, Error-Free over Transoceanic Distances," *Electron. Lett.* **32**:471 (1996).
45. L. F. Mollenauer, P. V. Mamyshev, and M. J. Neubelt, "Demonstration of Soliton WDM Transmission at up to 8×10 GBit/s, Error-Free over Transoceanic Distances," OFC-96, Postdeadline paper PD-22.
46. M. Suzuki, I. Morita, N. Edagawa, S. Yamamoto, H. Taga, and S. Akiba, "Reduction of Gordon-Haus Timing Jitter by Periodic Dispersion Compensation in Soliton Transmission," *Electron. Lett.* **31**:2027–2029 (1995).
47. N. J. Smith, N. J. Doran, F. M. Knox, and W. Forsysiak, "Energy-Scaling Characteristics of Solitons in Strongly Dispersion-Managed Fibers," *Opt. Lett.* **21**:1981–1983 (1996).
48. I. Gabitov and S. K. Turitsyn, "Averaged Pulse Dynamics in a Cascaded Transmission System with Passive Dispersion Compensation," *Opt. Lett.* **21**:327–329 (1996).
49. N. J. Smith, W. Forsysiak, and N. J. Doran, "Reduced Gordon-Haus Jitter Due to Enhanced Power Solitons in Strongly Dispersion Managed Systems," *Electron. Lett.* **32**:2085–2086 (1996).
50. V. S. Grigoryan, T. Yu, E. A. Golovchenko, C. R. Menyuk, and A. N. Pilipetskii, "Dispersion-Managed Soliton Dynamics," *Opt. Lett.* **21**:1609–1611 (1996).
51. G. Carter, J. M. Jacob, C. R. Menyuk, E. A. Golovchenko, and A. N. Pilipetskii, "Timing Jitter Reduction for a Dispersion-Managed Soliton System: Experimental Evidence," *Opt. Lett.* **22**:513–515 (1997).
52. S. K. Turitsyn, V. K. Mezentsev and E. G. Shapiro, "Dispersion-Managed Solitons and Optimization of the Dispersion Management," *Opt. Fiber Tech.* **4**:384–452 (1998).
53. J. P. Gordon and L. F. Mollenauer, "Scheme for Characterization of Dispersion-Managed Solitons," *Opt. Lett.* **24**:223–225 (1999).
54. P. V. Mamyshev and N. A. Mamysheva, "Pulse-Overlapped Dispersion-Managed Data Transmission and Intra-Channel Four-Wave Mixing," *Opt. Lett.* **24**:1454–1456 (1999).
55. D. Le Guen, S. Del Burgo, M. L. Moulinard, D. Grot, M. Henry, F. Favre, and T. Georges, "Narrow Band 1.02 Tbit/s (51×20 gbit/s) Soliton DWDM Transmission over 1000 km of Standard Fiber with 100 km Amplifier Spans," OFC-99, postdeadline paper PD-4.
56. S. Wabnitz, *Opt. Lett.* **21**:638–640 (1996).
57. E. A. Golovchenko, A. N. Pilipetskii, and C. R. Menyuk, *Opt. Lett.* **22**:1156–1158 (1997).

58. A. M. Niculae, W. Forysiak, A. G. Gloag, J. H. B. Nijhof, and N. J. Doran, "Soliton Collisions with Wavelength-Division Multiplexed Systems with Strong Dispersion Management," *Opt. Lett.* **23**:1354–1356 (1998).
59. P. V. Mamyshev and L. F. Mollenauer, "Soliton Collisions in Wavelength-Division-Multiplexed Dispersion-Managed Systems," *Opt. Lett.* **24**:448–450 (1999).
60. L. F. Mollenauer, P. V. Mamyshev, and J. P. Gordon, "Effect of Guiding Filters on the Behavior of Dispersion-Managed Solitons," *Opt. Lett.* **24**:220–222 (1999).
61. M. Matsumoto, *Opt. Lett.* **23**:1901–1903 (1998).
62. M. Matsumoto, *Electron. Lett.* **33**:1718 (1997).
63. L. F. Mollenauer, P. V. Mamyshev, J. Gripp, M. J. Neubelt, N. Mamysheva, L. Gruner-Nielsen, and T. Veng, "Demonstration of Massive WDM over Transoceanic Distances Using Dispersion Managed Solitons," *Opt. Lett.* **25**:704–706 (2000).

FIBER-OPTIC COMMUNICATION STANDARDS

Casimer DeCusatis

*IBM Corporation
Poughkeepsie, New York*

23.1 INTRODUCTION

This chapter presents a brief overview of several major industry standards for optical communications, including the following:

- ESCON/SBCON (Enterprise System Connection/Serial Byte Connection)
- FDDI (Fiber Distributed Data Interface)
- Fibre Channel Standard
- ATM (Asynchronous Transfer Mode)/SONET (Synchronous Optical Network)
- Ethernet (including Gigabit, 10 Gigabit, and other variants)
- InfiniBand

23.2 ESCON

The Enterprise System Connection (ESCON)* architecture was introduced on the IBM System/390 family of mainframe computers in 1990 as an alternative high-speed I/O channel attachment.^{1,2} The ESCON interface specifications were adopted in 1996 by the ANSI X3T1 committee as the serial byte connection (SBCON) standard.³

The ESCON/SBCON channel is a bidirectional, point-to-point 1300-nm fiber-optic data link with a maximum data rate of 17 Mbytes/s (200 Mbit/s). ESCON supports a maximum unrepeated distance of 3 km using 62.5 μm multimode fiber and LED transmitters with an 8-dB link budget, or a maximum unrepeated distance of 20 km using single-mode fiber and laser transmitters with a 14-dB link budget. The laser channels are also known as the ESCON extended distance feature (XDF). Physical connection is provided by an ESCON duplex connector, illustrated in Fig. 1. Recently, the single-mode

*ESCON is a registered trademark of IBM Corporation, 1991.

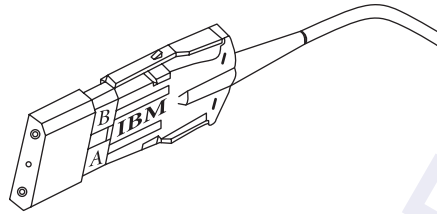


FIGURE 1 ESCON duplex fiber-optic connector.

ESCON links have adopted the SC duplex connector as standardized by Fibre Channel. With the use of repeaters or switches, an ESCON link can be extended up to 3 to 5 times these distances, and using wavelength division multiplexing (WDM) they can be extended even further, up to 100 km or more. However, performance of the attached devices in channel-to-channel applications typically falls off quickly at longer distances due to the longer round-trip latency of the link, making this approach suitable only for applications that can tolerate a lower effective throughput, such as remote backup of data for disaster recovery. There are some applications which can run over an ESCON physical layer without experiencing this performance degradation, such as virtual tape servers (VTS). ESCON devices and CPUs may communicate directly through a channel-to-channel attachment, but more commonly attach to a central nonblocking dynamic crosspoint switch. The resulting network topology is similar to a star-wired ring, which provides both efficient bandwidth utilization and reduced cabling requirements. The switching function is provided by an ESCON director, a nonblocking circuit switch. Although ESCON uses 8B/10B encoded data, it is not a packet-switching network; instead, the data frame header includes a request for connection that is established by the director for the duration of the data transfer. An ESCON data frame includes a header, payload of up to 1028 bytes of data, and a trailer. The header consists of a 2-character start-of-frame delimiter, 2-byte destination address, 2-byte source address, and 1 byte of link control information. The trailer is a 2-byte cyclic redundancy check (CRC) for errors, and a three-character end-of-frame delimiter. ESCON uses a DC-balanced 8B/10B coding scheme developed by IBM.

23.3 FDDI

The fiber distributed data interface (FDDI) was among the first open networking standards to specify optical fiber. It was an outgrowth of the ANSI X3T9.5 committee proposal in 1982 for a high-speed token passing ring as a back-end interface for storage devices. While interest in this application waned, FDDI found new applications as the backbone for local area networks (LANs). The FDDI standard was approved in 1992 as ISO standards IS 9314/1-2 and DIS 9314-3; it follows the architectural concepts of IEEE standard 802 (although it is controlled by ANSI, not IEEE, and therefore has a different numbering sequence) and is among the family of standards (including token ring and ethernet) that are compatible with a common IEEE 802.2 interface. FDDI is a family of four specifications, namely, the physical layer (PHY), physical media dependent (PMD), media access control (MAC), and station management (SMT). These four specifications correspond to sublayers of the data link and physical layer of the OSI reference model; as before, we will concentrate on the physical layer implementation.

The FDDI network is a 100-Mbit/s token passing ring, with dual counterrotating rings for fault tolerance. The dual rings are independent fiber-optic cables; the primary ring is used for data transmission, and the secondary ring is a backup in case a node or link on the primary ring fails. Bypass switches are also supported to reroute traffic around a damaged area of the network and prevent the ring from fragmenting in case of multiple node failures. The actual data rate is 125 Mbit/s, but this is reduced to an effective data rate of 100 Mbit/s by using a 4B/5B coding scheme. This high speed

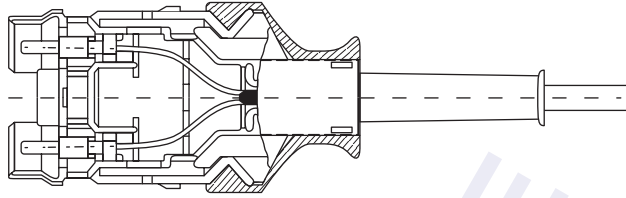


FIGURE 2 FDDI duplex fiber-optic connector.

allows FDDI to be used as a backbone to encapsulate lower speed 4, 10, and 16 Mbit/s LAN protocols; existing ethernet, token ring, or other LANs can be linked to an FDDI network via a bridge or router. Although FDDI data flows in a logical ring, a more typical physical layout is a star configuration with all nodes connected to a central hub or concentrator rather than to the backbone itself. There are two types of FDDI nodes, either *dual attach* (connected to both rings) or *single attach*; a network supports up to 500 dual-attached nodes, 1000 single-attached nodes, or an equivalent mix of the two types. FDDI specifies 1300-nm LED transmitters operating over 62.5 μm multimode fiber as the reference media, although the standard also provides for the attachment of 50, 100, 140, and 185 μm fiber. Using 62.5 μm fiber, a maximum distance of 2 km between nodes is supported with an 11-dB link budget; since each node acts like a repeater with its own phase-lock loop to prevent jitter accumulation, the entire FDDI ring can be as large as 100 km. However, an FDDI link can fail due to either excessive attenuation or dispersion; for example, insertion of a bypass switch increases the link length and may cause dispersion errors even if the loss budget is within specifications. For most other applications, this does not occur because the dispersion penalty is included in the link budget calculations or the receiver sensitivity measurements. The physical interface is provided by a special media interface connector (MIC), illustrated in Fig. 2. The connector has a set of three color-coded keys which are interchangeable depending on the type of network connection;¹ this is intended to prevent installation errors and assist in cable management.

An FDDI data frame is variable in length and contains up to 4500 8-bit bytes, or octets, including a preamble, start of frame, frame control, destination address, data payload, CRC error check, and frame status/end of frame. Each node has a MAC sublayer that reviews all the data frames looking for its own destination address. When it finds a packet destined for its node, that frame is copied into local memory; a copy bit is turned on in the packet; and it is then sent on to the next node on the ring. When the packet returns to the station that originally sent it, the originator assumes that the packet was received if the copy bit is on; the originator will then delete the packet from the ring. As in the IEEE 802.5 token ring protocol, a special type of packet called a *token* circulates in one direction around the ring, and a node can only transmit data when it holds the token. Each node observes a token retention time limit, and also keeps track of the elapsed time since it last received the token; nodes may be given the token in equal turns, or they can be given priority by receiving it more often or holding it longer after they receive it. This allows devices having different data requirements to be served appropriately.

Because of the flexibility built into the FDDI standard, many changes to the base standard have been proposed to allow interoperability with other standards, reduce costs, or extend FDDI into the MAN or WAN. These include a single-mode PMD layer for channel extensions up to 20 to 50 km. An alternative PMD provides for FDDI transmission over copper wire, either shielded or unshielded twisted pairs; this is known as *copper distributed data interface*, or CDDI. A new PMD was also developed to adapt FDDI data packets for transfer over a SONET link by stuffing approximately 30 Mbit/s into each frame to make up for the data rate mismatch (we will discuss SONET as an ATM physical layer in a later section). An enhancement called FDDI-II uses time-division multiplexing to divide the bandwidth between voice and data; it accommodates isochronous, circuit-switched traffic as well as existing packet traffic. An option known as *low cost* (LC) FDDI uses the more common SC duplex connector instead of the more expensive MIC connectors, and a lower-cost transceiver with a 9-pin footprint similar to the single-mode ESCON parts.

23.4 FIBRE CHANNEL STANDARD

Development of the ANSI Fibre Channel (FC) Standard began in 1988 under the X3T9.3 Working Group, as an outgrowth of the Intelligent Physical Protocol Enhanced Physical Project. The motivation for this work was to develop a scaleable standard for the attachment of both networking and I/O devices, using the same drivers, ports, and adapters over a single channel at the highest speeds currently achievable. The standard applies to both copper and fiber-optic media, and uses the English spelling *fib*re to denote both types of physical layers. In an effort to simplify equipment design, FC provides the means for a large number of existing upper-level protocols (ULPs), such as IP, SCSI, and HIPPI, to operate over a variety of physical media. Different ULPs are mapped to FC constructs, encapsulated in FC frames, and transported across a network; this process remains transparent to the attached devices. The standard consists of five hierarchical layers,⁴ namely a physical layer, an encode/decode layer which has adopted the DC-balanced 8B/10B code, a framing protocol layer, a common services layer (at this time, no functions have been formally defined for this layer), and a protocol-mapping layer to encapsulate ULPs into FC. Physical layer specifications for 1, 2, 4, and 8 Gbit/s links have been defined (refer to the ANSI standard for the most recent specifications). If the two link endpoints have different data rate capabilities, the links will auto-negotiate to the highest available rate between either 1, 2, and 4 Gbit/s rates or between 2, 4, and 8 Gbit/s rates. Note that the 10 Gbit/s data rate specifies a 64B/66B encoding scheme, rather than 8B/10B, and consequently is not backward compatible with lower data rates; this rate is typically reserved for inter-switch links (ISLs). The second layer defines the Fibre Channel data frame; frame size depends upon the implementation and is variable up to 2148 bytes long. Each frame consists of a 4-byte start-of-frame delimiter, a 24-byte header, a 2112-byte payload containing from 0 to 64 bytes of optional headers and 0 to 2048 bytes of data, a 4-byte CRC, and a 4-byte end-of-frame delimiter. In October 1994, the Fibre Channel physical and signaling interface standard FC-PH was approved as ANSI standard X3.230-1994.

Logically, Fibre Channel is a bidirectional point-to-point serial data link. Physically, there are many different media options (see Table 1) and three basic network topologies. The simplest, *default topology*, is a point-to-point direct link between two devices, such as a CPU and a device controller. The second, *Fibre Channel Arbitrated Loop* (FC-AL), connects between 2 and 126 devices in a loop configuration. Hubs or switches are not required, and there is no dedicated loop controller; all nodes on the loop share the bandwidth and arbitrate for temporary control of the loop at any given time. Each node has equal opportunity to gain control of the loop and establish a communications path; once the node relinquishes control, a fairness algorithm ensures that the same node cannot win control of the loop again until all other nodes have had a turn. As networks become larger, they may grow into the third topology, an *interconnected switchable network* or *fabric*, in which all network management functions are taken over by a switching point, rather than each node. An analogy for a switched fabric is the telephone network; users specify an address (phone number) for a device with which they want to communicate, and the network provides them with an interconnection path. In theory there is no limit to the number of nodes in a fabric; practically, there are only about 16 million unique addresses. Fibre Channel also defines three classes of connection service, which offer options such as guaranteed delivery of messages in the order they were sent and acknowledgment of received messages.

As shown in Table 1, FC provides for both single-mode and multimode fiber-optic data links using longwave (1300-nm) lasers and LEDs as well as short-wave (780 to 850 nm) lasers. The physical connection is provided by an SC duplex connector defined in the standard (see Fig. 3), which is keyed to prevent misplugging of a multimode cable into a single-mode receptacle. This connector design has since been adopted by other standards, including ATM, low-cost FDDI, and single-mode ESCON. The requirement for international class 1 laser safety is addressed using open fiber control (OFC) on some types of multimode links with shortwave lasers. This technique automatically senses when a full duplex link is interrupted, and turns off the laser transmitters on both ends to preserve laser safety. The lasers then transmit low-duty cycle optical pulses until the link is reestablished; a handshake sequence then automatically reactivates the transmitters.

TABLE 1 Examples of the Fiber Channel Standard Physical Layer

Media Type	Data Rate (Mbytes/s)	Maximum Distance	Signaling Rate (Mbaud)	Transmitter
SMF	800	10 km	8500.0	LW laser
	400	10 or 4 km	4250.0	LW laser
	200	10 km	2125.0	LW laser
	100	10 km	1062.5	LW laser
	50	10 km	1062.5	LW laser
	25	10 km	1062.5	LW laser
50- μ m multimode fiber	800	10 km	8500.0	SW laser
	400	10 or 4 km	4250.0	SW laser
	200	10 km	2125.0	SW laser
	100	500 m	1062.5	SW laser
	50	1 km	531.25	SW laser
	25	2 km	265.625	SW laser
	12.5	10 km	132.8125	LW LED
62.5- μ m multimode fiber	100	300 m	1062.5	SW laser
	50	600 m	531.25	SW laser
	25	1 km	265.625	LW LED
	12.5	2 km	132.8125	LW LED
105- Ω type 1 shielded twisted pair electrical	25	50 m	265.125	ECL
75 Ω mini coax	12.5	100 m	132.8125	ECL
	100	10 m	1062.5	ECL
75 Ω video coax	50	20 m	531.25	ECL
	25	30 m	265.625	ECL
	12.5	40 m	132.8125	ECL
	100	25 m	1062.5	ECL
	50	50 m	531.25	ECL
150 Ω twinax or STP	25	75 m	265.625	ECL
	12.5	100 m	132.8125	ECL
	100	30 m	1062.5	ECL
	50	60 m	531.25	ECL
	25	100 m	265.625	ECL

LW = long wavelength, SW = short wavelength, ECL = emitter-coupled logic.

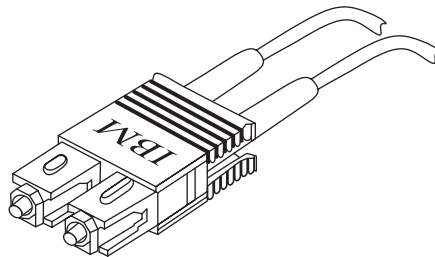


FIGURE 3 Single-mode SC duplex fiber-optic connector, per ANSI FC Standard specifications, with one narrow key and one wide key. Multimode SC duplex connectors use two wide keys.

23.5 ATM/SONET

Developed by the ATM Forum, this protocol promised to provide a common transport media for voice, data, video, and other types of multimedia. ATM is a high-level protocol that can run over many different physical layers including copper; part of ATM's promise to merge voice and data traffic on a single network comes from plans to run ATM over the synchronous optical network (SONET) transmission hierarchy developed for the telecommunications industry. SONET is really a family of standards defined by ANSI T1.105-1988 and T1.106-1988, as well as by several CCITT recommendations.⁵⁻⁸ Several different data rates are defined as multiples of 51.84 Mbit/s, known as OC-1. The numerical part of the OC-level designation indicates a multiple of this fundamental data rate, thus 155 Mbit/s is called OC-3. The standard provides for incremental data rates including OC-3, OC-9, OC-12, OC-18, OC-24, OC-36, and OC-48 (2.48832 Gbit/s). Both single-mode links with laser sources and multimode links with LED sources are defined for OC-1 through OC-12; only single-mode laser links are defined for OC-18 and beyond. SONET also contains provisions to carry sub-OC-1 data rates, called *virtual tributaries*, which support telecom data rates including DS-1 (1.544 Mbit/s), DS-2 (6.312 Mbit/s), and 3.152 Mbit/s (DS1C). The basic SONET data frame is an array of nine rows with 90 bytes per row, known as a synchronous-transport signal level 1 (STS-1) frame. In an OC-1 system, an STS-1 frame is transmitted once every 125 μ s (810 bytes per 125 μ s yields 51.84 Mbit/s). The first three columns provide overhead functions such as identification, framing, error checking, and a pointer which identifies the start of the 87-byte data payload. The payload floats in the STS-1 frame, and may be split across two consecutive frames. Higher speeds can be obtained either by concatenation of N frames into an STS- N c frame (the "c" stands for *concatenated*) or by byte-interleaved multiplexing of N frames into an STS- N frame.

ATM technology incorporates elements of both circuit and packet switching. All data is broken down into a 53-byte cell, which may be viewed as a short fixed-length packet. Five bytes make up the header, providing a 48-byte payload. The header information contains routing information (cell addresses) in the form of virtual path and channel identifiers; a field to identify the payload type; an error check on the header information; and other flow control information. Cells are generated asynchronously; as the data source provides enough information to fill a cell, it is placed in the next available cell slot. There is no fixed relationship between the cells and a master clock, as in conventional time-division multiplexing schemes; the flow of cells is driven by the bandwidth needs of the source. ATM provides bandwidth on demand; for example, in a client-server application the data may come in bursts; several data sources could share a common link by multiplexing during the idle intervals. Thus, the ATM adaptation layer allows for both constant and variable bit rate services. The combination of transmission options is sometimes described as a *pleiosynchronous network*, meaning that it combines some features of multiplexing operations without requiring a fully synchronous implementation. Note that the fixed cell length allows the use of synchronous multiplexing and switching techniques, while the generation of cells on demand allows flexible use of the link bandwidth for different types of data, characteristic of packet switching. Higher-level protocols may be required in an ATM network to ensure that multiplexed cells arrive in the correct order, or to check the data payload for errors (given the typical high reliability and low BER of modern fiber-optic technology, it was considered unnecessary overhead to replicate data error checks at each node of an ATM network). If an intermediate node in an ATM network detects an error in the cell header, cells may be discarded without notification to either end user. Although cell loss priority may be defined in the ATM header, for some applications the adoption of unacknowledged transmission may be a concern.

ATM data rates were intended to match SONET rates of 51, 155, and 622 Mbit/s; an FDDI-compliant data rate of 100 Mbit/s was added, in order to facilitate emulation of different types of LAN traffic over ATM. In order to provide a low-cost copper option and compatibility with 16-Mbit/s token ring LANs to the desktop, a 25-Mbit/s speed has also been approved. For premises wiring applications, ATM specifies the SC duplex connector, color coded beige for multimode links and blue for single-mode links. At 155 Mbit/s, multimode ATM links support a maximum distance of 3 km while single-mode links support up to 20 km.

23.6 ETHERNET

Ethernet was originally a local area network (LAN) communication standard developed for copper interconnections on a common data bus; it is an IEEE standard 802.3.⁹ The basic principle used in Ethernet is carrier sense multiple access with collision detection (CSMA/CD). Ethernet LANs may be configured as a bus, often wired radially through a central hub. A device attached to the LAN that intends to transmit data must first sense whether another device is transmitting. If another device is already sending, then it must wait until the LAN is available; thus, the intention is that only one device will be using the LAN to send data at a given time. When one device is sending, all other attached devices receive the data and check to see if it is addressed to them; if it is not, then the data is discarded. If two devices attempt to send data at the same time (e.g., both devices may begin transmission at the same time after determining that the LAN is available; there is a gap between when one device starts to send and before another potential sender can detect that the LAN is in use), then a collision occurs. Using CSMA/CD as the media access control protocol, when a collision is detected attached devices will detect the collision and must wait for different lengths of time before attempting retransmission. Since it is not always certain that data will reach its destination without errors or that the sending device will know about lost data, each station on the LAN must operate an end-to-end protocol for error recovery and data integrity. Data frames begin with an 8-byte preamble used for determining start-of-frame and synchronization, and a header consisting of a 6-byte destination address, 6-byte source address, and 2-byte length field. User data may vary from 46 to 1500 bytes, with data shorter than the minimum length padded to fit the frame; the user data is followed by a 2-byte CRC error check. Thus, an Ethernet frame may range from 70 to 1524 bytes.

The original Ethernet standard, known also as 10Base-T (10 Mbit/s over unshielded twisted pair copper wires) was primarily a copper standard, although a specification using 850-nm LEDs was also available. Subsequent standardization efforts increased this data rate to 100 Mbit/s over the same copper media (100Base-T), while once again offering an alternative fiber specification (100Base-FX). Recently, the standard has continued to evolve with the development of Gigabit Ethernet (1000Base-FX), which operates over fiber as the primary medium. Standardized as IEEE 802.3z, Gigabit Ethernet includes changes to the MAC layer in addition to a completely new physical layer operating at 1.25 Gbit/s. Switches rather than hubs predominate, since at higher data rates throughput per end user and total network cost are both optimized by using switched rather than shared media. The minimum frame size has increased to 512 bytes; frames shorter than this are padded with idle characters (carrier extension). The maximum frame size remains unchanged, although devices may now transmit multiple frames in bursts rather than single frames for improved efficiency. The physical layer will use standard 8B/10B data encoding. The standard allows several different physical connector types for fiber, including the SC duplex and various small-form-factor connectors about the size of a standard RJ-45 jack, although the LC duplex has become the most commonly used variant. Early transceivers were packaged as gigabit interface converters, or GBICs, which allows different optical or copper transceivers to be plugged onto the same host card. This has been replaced by small form factor pluggable transceivers (either SFP or SFP+). Some variants of the standard allow operation of long-wave (1300 nm) laser sources over both single-mode and multimode fiber. When a transmitter is optimized for a single-mode launch condition, it will underfill the multimode fiber; this causes some modes to be excited and propagate at different speeds than others, and the resulting differential mode delay significantly degrades link performance. One solution involves the use of special optical cables known as *optical mode conditioners* with offset ferrules to simulate an equilibrium mode launch condition into multimode fiber.

Ethernet continues to evolve as one of the predominant protocols for data center networking. Standards are currently being defined for both 40 Gbit/s and 100 Gbit/s Ethernet links. Other emerging standards allow transport of Fibre Channel over Ethernet, in an effort to converge two common types of optical networks. In an effort to make the Ethernet links more robust, other new standards known as Converged Enhanced Ethernet (CEE) are under development. The CEE standards add features such as new types of fine grained flow control, enhanced quality of service, and lossless transmission.

TABLE 2 Examples of the InfiniBand Physical Layer

Media Type	Per Lane Data Rate (Gbits/s)	Number of Lanes	Unidirectional Signaling Rate (Mbytes)	Transmitter
50- μ m multimode fiber	2.5	1X	250	VCSEL
		4X	10,000	VCSEL
		8X	20,000	VCSEL
		12X	30,000	VCSEL
50- μ m multimode fiber	5.0	1X	500	VCSEL
		4X	20,000	VCSEL
		8X	40,000	VCSEL
		12X	60,000	VCSEL
50- μ m multimode fiber	10.0	1X	10,000	VCSEL
SMF	10.0	1X	10,000	LW laser

LW = long wavelength, VCSEL = vertical cavity surface emitting laser.

23.7 INFINIBAND

The InfiniBand standards was developed by the InfiniBand Trade Association (IBTA) in an attempt to converge multiple protocol networks. Currently, it is most widely used for low-latency, high-performance applications in data communication. InfiniBand specifies 8B/10B encoded data, and both serial and parallel optical links, some of which are illustrated in Table 2.¹⁰ These are referred to by the number of lanes in the physical layer interface; for example, a 4X link employs four optical fibers in each direction of a bidirectional link. The industry standard multifiber push-on (MPO) connector is specified for parallel optical links, while the SC duplex connector is commonly used for single fiber links (note that the 10 Gbit/s serial IB link physical layer is very similar to the 10-Gbit Ethernet link specification). InfiniBand is a switched point-to-point protocol, although some data communication applications employ the InfiniBand physical layer only, and are therefore not compatible with InfiniBand switches (e.g., the Parallel Sysplex IB links developed for IBM mainframes). Although InfiniBand links are not designed for operation at distances beyond 10 km, with sufficiently large receive buffers or other flow control management techniques they can be extended to much longer distances. This can be done using protocol independent wavelength multiplexing, or by encapsulating InfiniBand into another protocol such as SONET.

23.8 REFERENCES

1. D. Stigliani, "Enterprise Systems Connection Fiber Optic Link," Chap. 13, in *Handbook of Optoelectronics for Fiber Optic Data Communications*, C. DeCusatis, R. Lasky, D. Clement, and E. Mass (eds.), Academic Press, San Diego, California (1997).
2. "ESCON I/O Interface Physical Layer Document" *IBM Document Number SA23-0394*, 3rd ed. IBM Corporation, Mechanicsburg, Pennsylvania (1995).
3. Draft ANSI Standard X3T11/95-469 (rev. 2.2) "ANSI Single Byte Command Code Sets Connection Architecture (SBCON)," ANSI, Washington, DC (1996).
4. ANSI X3.230-1994 (rev. 4.3), "Fibre Channel—Physical and Signaling Interface (FC-PH)," ANSI X3.272-199x (rev. 4.5), "Fibre Channel—Arbitrated Loop (FC-AL)," (June 1995); ANSI X3.269-199x, (rev. 012), "Fiber Channel Protocol for SCSI (FCP)," ANSI, Washington, DC (May 30, 1995).
5. ANSI T1.105-1988, "Digital Hierarchy Optical Rates and Format Specification," ANSI, Washington, DC (1988).
6. CCITT Recommendation G.707, "Synchronous Digital Hierarchy Bit Rates," CCITT (1991).

7. CCITT Recommendation G.708, "Network Node Interfaces for the Synchronous Digital Hierarchy," CCITT, Geneva, Switzerland (1991).
8. CCITT Recommendation G.709, "Synchronous Multiplexing Structure," CCITT, Geneva, Switzerland (1991).
9. IEEE 802.3z, "Draft Supplement to Carrier Sense Multiple Access with Collision Detection (CSMA/CD) Access Method and Physical Layer Specifications: Media Access Control (MAC) Parameters, Physical Layer, Repeater and Management Parameters for 1000 Mb/s Operation," IEEE, Piscataway, New Jersey (June 1997).
10. A. Ghiasi, "InfiniBand," in *Handbook of Fiber Optic Data Communication*, 3rd ed., C. DeCusatis (ed.), Academic Press, San Diego (2008).

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

Richard O. Claus

*Virginia Tech
Blacksburg, Virginia*

Ignacio Matias and Francisco Arregui

*Public University Navarra
Pamplona, Spain*

24.1 INTRODUCTION

Optical fiber sensors are a broad topic. The objective of this chapter is to briefly summarize the fundamental properties of representative types of optical fiber sensors and how they operate. Four different types of sensors are evaluated systematically on the basis of performance criteria such as resolution, dynamic range, cross-sensitivity to multiple ambient perturbations, fabrication, and demodulation processes. The optical fiber sensing methods that will be investigated include well-established technologies such as fiber Bragg grating (FBG)-based sensors, and rapidly evolving measurement techniques such as those involving long-period gratings (LPGs). Additionally, two popular versions of Fabry-Perot interferometric sensors (intrinsic and extrinsic) are evaluated.

The outline of this chapter is as follows. The principles of operation and fabrication processes of each of the four sensors are discussed separately. The sensitivity of the sensors to displacement and simultaneous perturbations such as temperature is analyzed. The overall complexity and performance of a sensing technique depends heavily on the signal demodulation process. Thus, the detection schemes for all four sensors are discussed and compared on the basis of their complexity. Finally, a theoretical analysis of the cross-sensitivities of the four sensing schemes is presented and their performance is compared.

Measurements of a wide range of physical measurands by optical fiber sensors have been investigated for more than 20 years. Displacement measurements using optical fiber sensors are typical of these, and both embedded and surface-mounted configurations have been reported by researchers in the past.¹ Fiber-optic sensors are small in size, are immune to electromagnetic interference, and can be easily integrated with existing optical fiber communication links. Such sensors can typically be easily multiplexed, resulting in distributed networks that can be used for health monitoring of integrated, high-performance materials and structures.

Optical fiber sensors of displacement are perhaps the most basic of all fiber sensor types because they may be configured to measure many other related environmental factors. They should possess certain important characteristics. First, they should either be insensitive to ambient fluctuations in temperature and pressure or should employ demodulation techniques that compensate for changes in the output signal due to these additional perturbations. In an embedded configuration, the sensors for axial strain measurements should have minimum cross-sensitivity to other strain states. The sensor signal should itself be simple and easy to demodulate. Nonlinearities in the output require

expensive decoding procedures or necessitate precalibration and sensor-to-sensor incompatibility. The sensor should ideally provide an absolute and real-time displacement or strain measurement in a form that can be easily processed. For environments where large strain magnitudes are expected, the sensor should have a large dynamic range while at the same time maintaining the desired sensitivity. We now discuss each of the four sensing schemes individually and present their relative advantages and shortcomings.

24.2 EXTRINSIC FABRY-PEROT INTERFEROMETRIC SENSORS

The extrinsic Fabry-Perot interferometric (EFPI) sensor, proposed by a number of groups and authors, is one of the most popular fiber-optic sensors used for applications in health monitoring of smart materials and structures.^{2,3} As the name suggests, the EFPI is an interferometric sensor in which the detected intensity is modulated by the parameter under measurement. The simplest configuration of an EFPI is shown in Fig. 1.

The EFPI system consists of a single-mode laser diode that illuminates a Fabry-Perot cavity through a fused biconical tapered coupler. The cavity is formed between an input single-mode fiber and a reflecting target element that may be a fiber. Since the cavity is external to the lead-in/lead-out fiber, the EFPI sensor is independent of transverse strain and small ambient temperature fluctuations. The input fiber and the reflecting fiber are typically aligned using a hollow core tube as shown in Fig. 10. For optical fibers with uncoated ends, Fresnel reflection of approximately 4 percent results at the glass-to-air and air-to-glass interfaces that define the cavity. The first reflection at the glass-air interface R_1 , called the *reference reflection*, is independent of the applied perturbation. The second reflection at the air-glass interface R_2 , termed the *sensing reflection*, is dependent on the length of the cavity d , which in turn is modulated by the applied perturbation. These two reflections interfere

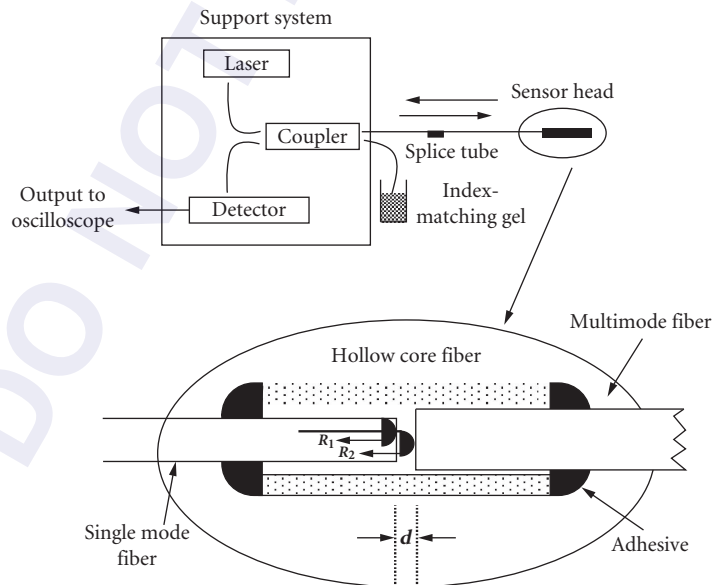


FIGURE 1 Extrinsic Fabry-Perot interferometric sensor and system.

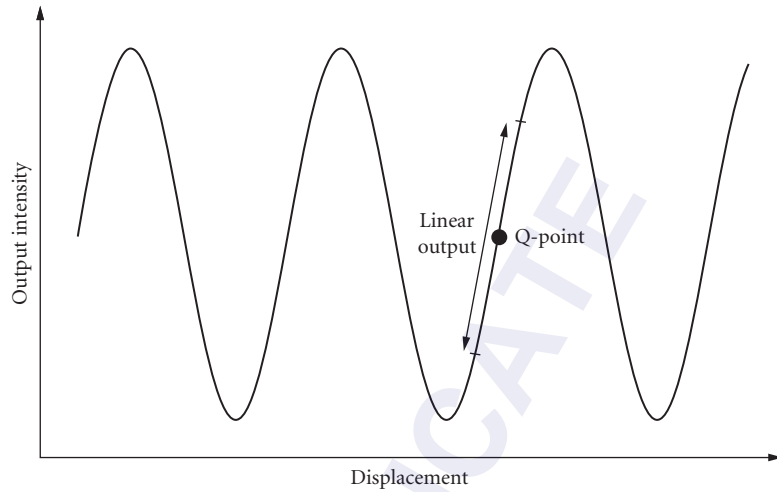


FIGURE 2 EFPI transfer function curve.

(provided $2d < L_c$, the coherence length of the light source) and the intensity I at the detector varies as a function of the cavity length,

$$I = I_0 \cos\left(\frac{4\pi}{\lambda}d\right) \quad (1)$$

where I_0 is the maximum value of the output intensity and λ is the center wavelength of the light source, here assumed to be a laser diode.

The typical intensity-versus-displacement transfer function curve [Eq. (1)] for an EFPI sensor is shown in Fig. 2. Small perturbations that result in operation around the quiescent or Q point of the sensor lead to an approximately linear variation in output intensity versus applied displacement. For larger displacements, the output signal is not a linear function of the input signal, and the output signal may vary over several sinusoidal periods. In this case, a *fringe* in the output signal is defined as the change in intensity from a maximum to a maximum, or from a minimum to a minimum, so each fringe corresponds to a change in the cavity length by half of the operating wavelength λ . The change in the cavity length Δd is then employed to calculate the strain using the expression

$$\varepsilon = \frac{\Delta d}{L} \quad (2)$$

where L is defined as the gauge length of the sensor and is typically the distance between two points where the input and reflecting fibers are bonded to the hollow-core support tube.

The EFPI sensor has been used for the analysis of materials and structures.^{1,3} The relatively low temperature sensitivity of the sensor element, due to the opposite directional expansion of the fiber and tube elements, makes it attractive for the measurement of strain and displacement in environments where the temperature is not anticipated to change over a wide range. The EFPI sensor is capable of measuring subangstrom displacements with strain resolution better than $1 \mu\varepsilon$ and a dynamic range greater than $10,000 \mu\varepsilon$. Moreover, the large bandwidth simplifies the measurement of highly cyclical strain. The sensor also allows single-ended operation and is hence suitable for applications where ingress to and egress from the sensor location are important. The sensor requires simple and inexpensive fabrication equipment and an assembly time of a few minutes. Additionally, since the cavity is external to the fibers, transverse strain components that tend to influence the response of similar intrinsic sensors through Poisson-effect cross-coupling have negligible effect on the EFPI sensor output.

24.3 INTRINSIC FABRY-PEROT INTERFEROMETRIC SENSORS

The intrinsic Fabry-Perot interferometric (IFPI) sensor is similar in operation to its extrinsic counterpart, but significant differences exist in the configurations of the two sensors.⁴ The basic IFPI sensor is shown in Fig. 3. An optically isolated laser diode is used as the optical source to one of the input arms of a bidirectional 2×2 coupler. The Fabry-Perot cavity is formed internally by fusing a small length of single-mode fiber to one of the output legs of the coupler. As shown in Fig. 3, the reference (R) and sensing (S) reflections interfere at the detector to again provide a sinusoidal intensity variation versus cavity path length modulation. The cavity can also be implemented by introducing two Fresnel or other reflectors along the length of a single fiber. The photosensitivity effect in germanosilicate fibers has been used in the past to fabricate broadband grating-based reflector elements to define such an IFPI cavity.⁵ Since the cavity is formed within an optical fiber, changes in the refractive index of the fiber due to the applied perturbation can significantly alter the phase of the sensing signal S . Thus the intrinsic cavity results in the sensor being sensitive to ambient temperature fluctuations and all states of strain.

The IFPI sensor, like all other interferometric signals, has a nonlinear output that complicates the measurement of large-magnitude strain. This can again be overcome by operating the sensor in the linear regime around the Q point of the sinusoidal transfer function curve. The main limitation of the IFPI strain sensor is that the photoelastic-effect-induced change in index of refraction results in a nonlinear relationship between the applied perturbation and the change in cavity length. For most IFPI sensors, the change in the propagation constant of the fundamental mode dominates the change in cavity length. Thus IFPIs are highly susceptible to temperature changes and transverse strain components.⁶ In embedded applications, the sensitivity to all of the strain components can result in complex signal output. The process of fabricating an IFPI strain sensor is more complicated than that for the EFPI sensor since the sensing cavity of the IFPI sensor must be formed within the optical fiber by some special procedure. The strain resolution of IFPI sensors is approximately $1 \mu\epsilon$ with an operating range greater than $10,000 \mu\epsilon$. IFPI sensors also suffer from drift in the output signal due to variations in the polarization state of the input light.

Thus the preliminary analysis shows that the extrinsic version of the Fabry-Perot optical fiber sensor seems to have an overall advantage over its intrinsic counterpart. The extrinsic sensor has negligible cross-sensitivity to temperature and transverse strain. Although the strain sensitivities, dynamic ranges, and bandwidths of the two sensors are comparable, the IFPIs can be expensive and cumbersome to fabricate due to the intrinsic nature of the sensing cavity.

The extrinsic and intrinsic Fabry-Perot interferometric sensors possess nonlinear sinusoidal outputs that complicate signal processing at the detector. Although intensity-based sensors have a simple output variation, they suffer from limited sensitivity to strain or other perturbations of interest. Grating-based sensors have recently become popular as transducers that provide wavelength-encoded output signals that can typically be easily modulated to derive information about the perturbation under investigation. We next discuss the advantages and drawbacks of Bragg grating sensing technology. The basic operating mechanism of Bragg grating-based strain sensors is then reviewed

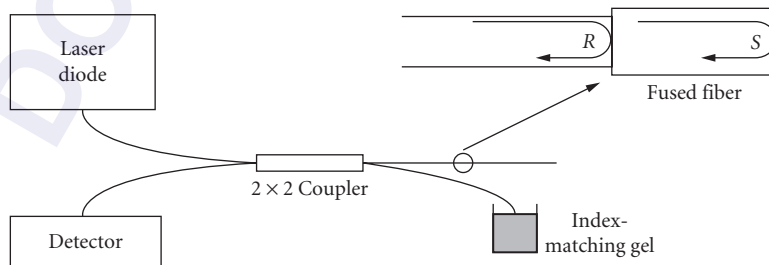


FIGURE 3 The intrinsic Fabry-Perot interferometric (IFPI) sensor.

and the expressions for strain resolution are obtained. These sensors are then compared to the recently developed long-period grating devices in terms of fabrication process, cross-sensitivity to multiple measurands, and simplicity of signal demodulation.

24.4 FIBER BRAGG GRATING SENSORS

The phenomenon of photosensitivity was discovered by Hill and coworkers in 1978.⁷ It was found that permanent refractive index changes could be induced in optical fibers by exposing the germanium-doped core of a fiber to intense light at 488 or 514 nm. Hill found that a sinusoidal modulation of index of refraction in the core created by the spatial variation of such an index-modifying beam gives rise to refractive index grating that can be used to couple the energy in the fundamental guided mode to various guided and lossy modes. Later Meltz et al.⁸ suggested that photosensitivity is more efficient if the fiber is side-exposed to a writing beam at wavelengths close to the absorption wavelength (242 nm) of the germanium defects in the fiber. The side-writing process simplified the fabrication of Bragg gratings, and these devices have recently emerged as highly versatile components for optical fiber communication and sensing systems. Recently, loading of the fibers with hydrogen prior to writing has been used to produce order-of-magnitude larger changes in index in germanosilicate fibers.⁹

Principle of Operation

Bragg gratings in optical fibers are based on a phase-matching condition between propagating optical modes. This phase-matching condition is given by

$$k_g + k_c = k_B \quad (3)$$

where k_g , k_c , and k_B are, respectively, the wave vectors of the coupled guided mode, the resulting coupling mode, and the grating. For a first-order interaction, $k_B = 2\pi/\Lambda$, where Λ is the spatial period of the grating. In terms of propagation constants, this condition reduces to the general form of interaction for mode coupling due to a periodic perturbation

$$\Delta\beta = \frac{2\pi}{\Lambda} \quad (4)$$

where $\Delta\beta$ is the difference in the propagation constants of the two modes involved in mode coupling, where both modes are assumed to travel in the same direction.

Fiber Bragg gratings (FBGs) involve the coupling of the forward-propagating fundamental LP₀₁ in a single-mode fiber to the reverse-propagating LP₀₁ mode.¹⁰ Here, consider a single-mode fiber with β_{01} and $-\beta_{01}$ as the propagation constants of the forward- and reverse-propagating fundamental LP₀₁ modes. To satisfy the phase-matching condition,

$$\Delta\beta = \beta_{01} - (-\beta_{01}) = \frac{2\pi}{\Lambda} \quad (5)$$

where $\beta_{01} = 2\pi n_{\text{eff}}/\lambda$, n_{eff} is the effective index of the fundamental mode, and λ is the free-space wavelength of the source. Equation (5) reduces to¹⁰

$$\lambda_B = 2\Lambda n_{\text{eff}} \quad (6)$$

where λ_B is termed the *Bragg wavelength*—the wavelength at which the forward-propagating LP₀₁ mode couples to the reverse-propagating LP₀₁ mode. Such coupling is wavelength dependent, since the propagation constants of the two modes are a function of the wavelength. Hence, if an FBG element is interrogated using a broadband optical source, the wavelength at which phase matching occurs is back-reflected. This back-reflected wavelength is a function of the grating period Λ and the

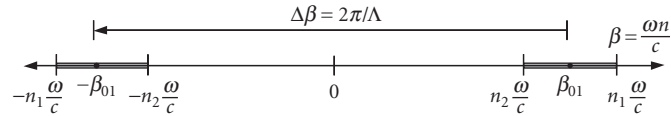


FIGURE 4 Mode-coupling mechanism in fiber Bragg gratings. The large value of $\Delta\beta$ in FBGs requires a small value of the grating periodicity Λ . The hatched regions represent the guided modes in the forward ($\beta > 0$) and reverse ($\beta < 0$) directions.

effective index n_{eff} of the fundamental mode as shown in Eq. (6). Since strain and temperature effects can modulate both of these parameters, the Bragg wavelength is modulated by both of these external perturbations. The resulting spectral shifts are utilized to implement FBGs for sensing applications.

Figure 4 shows the mode-coupling mechanism in fiber Bragg gratings using a β -plot. Since the difference in propagation constants ($\Delta\beta$) between the modes involved in coupling is large, we see from Eq. (4) that only a small period, Λ , is needed to induce this mode coupling. Typically for optical fiber communication system applications the value of λ_B is approximately $1.5 \mu\text{m}$. From Eq. (6), Λ is thus approximately $0.5 \mu\text{m}$ for $n_{\text{eff}} = 1.5$, the approximate index of refraction of the glass in a fiber. Due to the small period, on the order of $1 \mu\text{m}$, FBGs are typically classified as short-period gratings (SPGs).

Bragg Grating Sensor Fabrication

Fiber Bragg gratings have commonly been manufactured using two side-exposure techniques, namely *interferometric* and *phase mask* methods. The interferometric method, shown in Fig. 5, uses an ultraviolet (UV) writing beam at 244 or 248 nm, split into two parts of approximately the same intensity by a beam splitter.⁸ The two beams are focused on a portion of the Ge-doped fiber, whose protective coating has been removed using cylindrical lenses. The period of the resulting interference pattern, and hence the period of the Bragg grating element to be written, is varied by altering the mutual angle θ . A limitation of this method is that any relative vibration of the pairs of mirrors and lenses can lead to the degradation of the quality of the fringe pattern and the fabricated grating; thus the entire system has a stringent stability requirement. To overcome this drawback, Kashyap¹⁰ proposed a novel interferometer technique in which the path difference between the interfering UV beams is produced by propagation through a right-angled prism, as shown in Fig. 6. This geometry is inherently stable because both beams are perturbed similarly by any prism vibration.

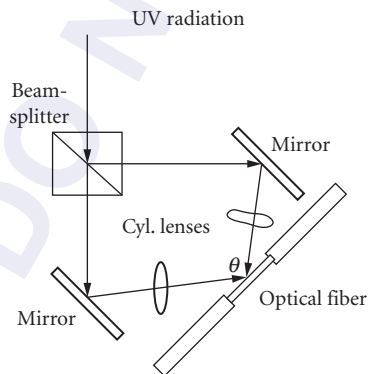


FIGURE 5 Fabrication of Bragg gratings using interferometric scheme.

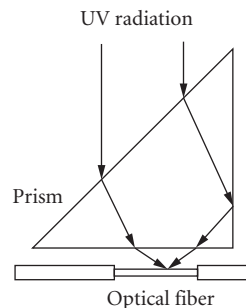


FIGURE 6 Bragg grating fabrication using prism method.

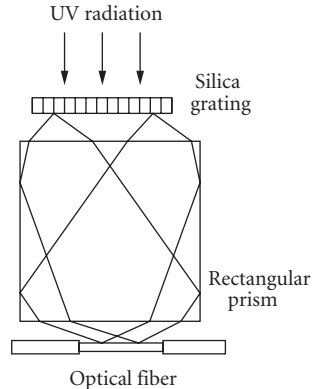


FIGURE 7 Phase mask method of fabricating Bragg gratings.

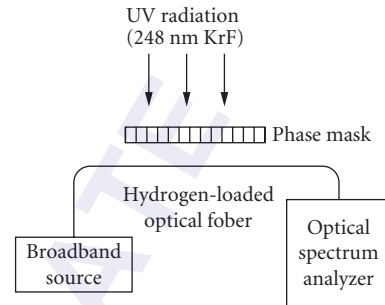


FIGURE 8 Setup to write Bragg gratings in germanosilicate fibers.

The phase mask technique has gained popularity as an efficient holographic side-writing procedure for grating fabrication.¹¹ In this method, shown in Fig. 7, an incident UV beam is diffracted into -1 , 0 , and $+1$ orders by a relief grating typically generated on a silica plate by electron beam exposure and plasma etching. The two first diffraction orders undergo total internal reflection at the glass-air interface of a rectangular prism and interfere at the location of the fiber placed directly behind the mask. This technique is wavelength specific, since the period of the resulting two-beam interference pattern is uniquely determined by the diffraction angle of -1 and $+1$ orders and thus the properties of the phase mask. Obviously, different phase masks are required for the fabrication of gratings at different Bragg wavelengths. A setup for actively monitoring the growth of a grating in the transmission mode during fabrication is shown in Fig. 8.

Bragg Grating Sensors

From Eq. (6) we see that a change in the value of n_{eff} and/or Λ can cause the Bragg wavelength λ to shift. This fractional change in the resonance wavelength $\Delta\lambda/\lambda$ is given by

$$\frac{\Delta\lambda}{\lambda} = \frac{\Delta\Lambda}{\Lambda} + \frac{\Delta n_{\text{eff}}}{n_{\text{eff}}} \quad (7)$$

where $\Delta\Lambda/\Lambda$ and $\Delta n_{\text{eff}}/n_{\text{eff}}$ are the fractional changes in the period and the effective index, respectively. The relative magnitudes of the two changes depend on the type of perturbation to which the grating is subjected. For most applications the effect due to change in effective index is the dominating mechanism.

An axial strain ε in the grating changes the grating period and the effective index and results in a shift in the Bragg wavelength, given by

$$\frac{1}{\lambda} \frac{\Delta\lambda}{\varepsilon} = \frac{1}{\Lambda} \frac{\Delta\Lambda}{\varepsilon} + \frac{1}{n_{\text{eff}}} \frac{\Delta n_{\text{eff}}}{\varepsilon} \quad (8)$$

The first term on the right side of Eq. (8) is unity, while the second term has its origin in the photoelastic effect. An axial strain in the fiber serves to change the refractive index of both the core and the cladding. This results in a variation in the value of the effective index of glass. The photoelastic or strain-optic coefficient is approximately -0.27 . Thus, the variations in n_{eff} and Λ due to strain have contrasting effects on the Bragg peak. The fractional change in the Bragg wavelength due to axial strain is 0.73ε , or 73 percent of the applied strain. At 1550 and 1300 nm, the shifts in the resonance wavelength are 11 nm/% ε and 9 nm/% ε , respectively. An FBG at 1500 nm shifts by 1.6 nm for every 100°C rise in temperature.⁷

Limitations of Bragg Grating Strain Sensors

The primary limitation of Bragg grating sensors is the complex and expensive fabrication technique. Although side-writing is commonly being used to manufacture these gratings, the requirement of expensive phase masks increases the cost of the sensing system. In the interferometric technique, stability of the setup is a critical factor in obtaining high-quality gratings. Since index changes of the order of 10^{-3} are required to fabricate these gratings, laser pulses of high energy levels are necessary.

The second primary limitation of Bragg gratings is their limited bandwidth. The typical value of the full width at half-maximum (FWHM) is between 0.1 and 1 nm. Although higher bandwidths can be obtained by chirping the index or period along the grating length, this adds to the cost of the grating fabrication. The limited bandwidth requires high-resolution spectrum analysis to monitor the grating spectrum. Kersey and Berkoff¹² have proposed an unbalanced Mach-Zehnder interferometer to detect the perturbation-induced wavelength shift. Two unequal arms of the Mach-Zehnder interferometer are excited by the back reflection from a Bragg grating sensor element. Any change in the input optical wavelength modulates the phase difference between the two arms and results in a time-varying sinusoidal intensity at the output. This interference signal can be related to the shift in the Bragg peak and the magnitude of the perturbation can be obtained. Recently, modal interferometers have also been proposed to demodulate the output of a Bragg grating sensor.¹³ The unbalanced interferometers are also susceptible to external perturbations and hence need to be isolated from the parameter under investigation. Moreover, the nonlinear output may require fringe counting, which can be complicated and expensive. Additionally, a change in the perturbation polarity at the maxima or minima of the transfer function curve will not be detected by this demodulation scheme. To overcome this limitation, two unbalanced interferometers may be employed for dynamic measurements.

Cross-sensitivity to temperature leads to erroneous displacement measurements in applications where the ambient temperature has a temporal variation. So a reference grating used to measure temperature change may be utilized to compensate for the output of the strain sensor. Recently, temperature-independent sensing has been demonstrated using chirped gratings written in tapered optical fibers.¹⁴

Finally, the sensitivity of fiber Bragg grating strain sensors may not be adequate for certain applications. This sensitivity of the sensor depends on the minimum detectable wavelength shift at the receiver. Although excellent wavelength resolution can be obtained with unbalanced interferometric detection techniques, standard spectrum analysis systems typically provide a resolution of 0.1 nm. At 1300 nm, this minimum detectable change in wavelength corresponds to a strain resolution of 111 $\mu\epsilon$. Hence, in applications where strains smaller than 100 $\mu\epsilon$ are anticipated, Bragg grating sensors may not be practical. The dynamic range of strain measurement can be as much as 15,000 $\mu\epsilon$.

24.5 LONG-PERIOD GRATING SENSORS

This section discusses the use of novel long-period gratings (LPGs) as displacement-sensing devices. We analyze the principle of operation of these gratings, their fabrication process, typical experimental evaluation, their demodulation process, and their cross-sensitivity to ambient temperature.

Principle of Operation

Long-period gratings that couple the fundamental guided mode to different guided modes have been demonstrated.^{15,16} Gratings with longer periodicities that involve coupling of a guided mode to forward-propagating cladding modes were recently proposed by Vengsarkar et al.^{17,18} As discussed previously, fiber gratings satisfy the Bragg phase-matching condition between the guided and cladding or radiation modes or another guided mode. This wavelength-dependent phase-matching condition is given by

$$\beta_{01} - \beta = \Delta\beta = \frac{2\pi}{\Lambda} \quad (9)$$

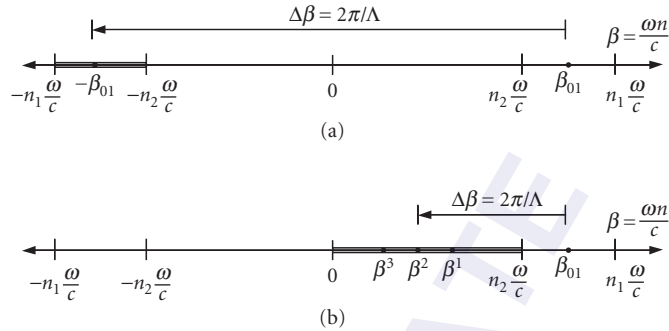


FIGURE 9 Depiction of mode coupling in (a) Bragg gratings and (b) long-period gratings. The differential propagation constant $\Delta\beta$ determines the grating periodicity.

where Λ is the period of the grating and β_{01} and β are the propagation constant of the fundamental guided mode and the mode to which coupling occurs, respectively.

For conventional fiber Bragg gratings, the coupling of the forward-propagating LP_{01} mode occurs to the reverse-propagating LP_{01} mode ($\beta = -\beta_{01}$). Since $\Delta\beta$ is large in this case, as shown in Fig. 9a, the grating periodicity is small, typically on the order of $1\ \mu\text{m}$. Unblazed long-period gratings couple the fundamental mode to the discrete and circularly symmetric, forward-propagating cladding modes ($\beta = \beta^n$), resulting in smaller values of $\Delta\beta$, as shown in Fig. 9b, and hence periodicities ranging in the hundreds of micrometers.¹⁷ The cladding modes attenuate rapidly as they propagate along the length of the fiber, due to the lossy cladding-coating interface and bends in the fiber. Since $\Delta\beta$ is discrete and a function of the wavelength, this coupling to the cladding modes is highly selective, leading to a wavelength-dependent loss. As a result, any modulation of the core and cladding guiding properties modifies the spectral response of long-period gratings, and this phenomenon can be utilized for sensing purposes. Moreover, since the cladding modes interact with the fiber jacket or any other material surrounding the cladding, changes in the index of refraction or other properties of these effective coatings materials can also be detected.

LPG Fabrication Procedure

To fabricate long-period gratings, hydrogen-loaded (3.4 mole %), germanosilicate fibers may be exposed to 248-nm UV radiation from a KrF excimer laser through a chrome-plated amplitude mask possessing a periodic rectangular transmittance function. Figure 10 shows a typical setup used

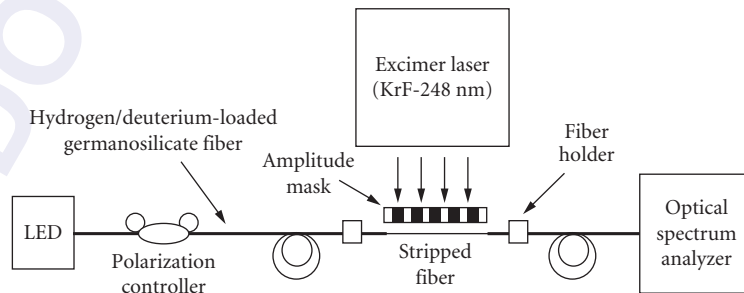


FIGURE 10 Setup to fabricate long-period gratings using an amplitude mask.

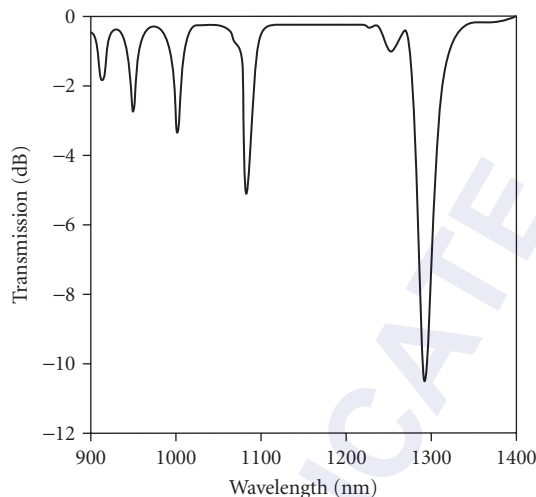


FIGURE 11 Transmission spectrum of a long-period grating written in Corning FLEXCOR fiber with period $\Lambda = 198 \mu\text{m}$. The discrete, spiky loss bands correspond to the coupling of the fundamental guided mode to discrete cladding modes.

to fabricate such gratings. The laser is pulsed at approximately 20 Hz with a pulse duration of several nanoseconds. The typical writing times for an energy of $100 \text{ mJ}/\text{cm}^2/\text{pulse}$ and a 2.5-cm exposed length vary between 6 and 15 min for different fibers. The coupling wavelength λ_p shifts to higher values during exposure due to the photoinduced enhancement of the refractive index of the fiber core and the resulting increase in β_{01} . After writing, the gratings are annealed at 150°C for several hours to remove the unreacted hydrogen. This high-temperature annealing causes λ_p to move to shorter wavelengths due to the decay of UV-induced defects and the diffusion of molecular hydrogen from the fiber. Figure 11 depicts the typical transmittance of a grating. Various attenuation bands correspond to coupling to discrete cladding modes of different orders. A number of gratings can be fabricated at the same time by placing more than one fiber behind the amplitude mask. Due to the relatively long spatial periods, the stability requirements during the writing process are not so severe as those for short-period Bragg gratings.

For coupling to the highest-order cladding mode, the maximum isolation (loss in transmission intensity) is typically in the 5- to 20-dB range on wavelengths depending on fiber parameters, duration of UV exposure, and mask periodicity. The desired fundamental coupling wavelength can easily be varied by using inexpensive amplitude masks of different periodicities. The insertion loss, polarization mode dispersion, backreflection, and polarization-dependent loss of a typical grating are 0.2 dB, 0.01 ps, -80 dB , and 0.02 dB, respectively. The negligible polarization sensitivity and backreflection of these devices eliminates the need for expensive polarizers and isolators.

We now look at representative experiments that have been performed and discussed to examine the displacement sensitivity of long-period gratings written in different fibers.^{19,20} For example, gratings have been fabricated in four different types of fibers—standard dispersion-shifted fiber (DSF), standard 1550-nm fiber, and conventional 980- and 1050-nm single-mode fibers. For the sake of brevity, these will be referred to as fibers A, B, C, and D, respectively. The strain sensitivity of gratings written in different fibers was determined by axially straining the gratings between two longitudinally separated translation stages. The shift in the peak loss wavelength of the grating in fiber D as a function of the applied strain is depicted in Fig. 12 along with that for a Bragg grating (about $9 \text{ nm}/\% \epsilon$ at 1300 nm).⁷ The strain coefficients of wavelength shift β for fibers A, B, C, and D are shown in Table 1.

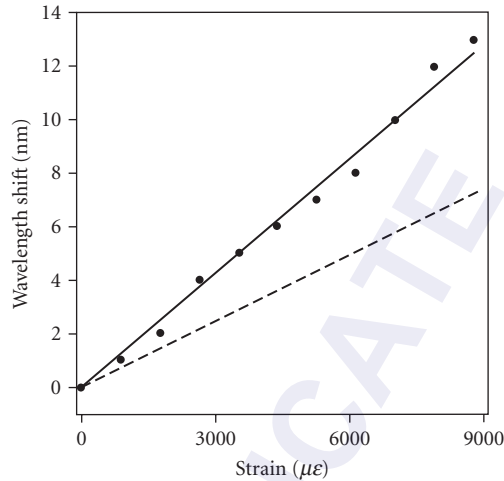


FIGURE 12 Shift in the highest order resonance band with strain for a long-period grating written in fiber D (circles). Also depicted is the shift for a conventional Bragg grating (dashed line).

TABLE 1 Strain Sensitivity of Long-Period Gratings Written in Four Different Types of Fibers

Type of Fiber	Strain Sensitivity (nm/% ϵ)
A—standard dispersion-shifted fiber (DSF)	-7.27
B—standard 1550-nm communication fiber	4.73
C—conventional 980-nm single-mode fiber	4.29
D—conventional 1060-nm single-mode fiber	15.21

Values correspond to the shift in the highest order resonance wavelength.

Fiber D has a coefficient of 15.2 nm/% ϵ , which gives it a strain-induced shift that is 50 percent larger than that for a conventional Bragg grating. The strain resolution of this fiber for a 0.1-nm detectable wavelength shift is 65.75 $\mu\epsilon$.

The demodulation scheme of a sensor determines the overall simplicity and sensitivity of the sensing system. Short-period Bragg grating sensors were shown to possess signal processing techniques that are complex and expensive to implement. We now present a simple demodulation method to extract information from long-period gratings. The wide bandwidth of the resonance bands enables the wavelength shift due to the external perturbation to be converted into an intensity variation that can be easily detected.

Figure 13 shows the shift induced by strain in a grating written in fiber C. The increase in the loss at 1317 nm is about 1.6 dB. A laser diode centered at 1317 nm was used as the optical source, and the change in transmitted intensity was monitored as a function of applied strain. The transmitted intensity is plotted in Fig. 14 for three different trials. The repeatability of the experiment demonstrates the feasibility of using this simple scheme to utilize the high sensitivity of long-period gratings. The transmission of a laser diode centered on the slope of the grating spectrum on either side of the resonance wavelength can be used as a measure of the applied perturbation. A simple detector and amplifier combination at the output can be used to determine the transmission through the detector.

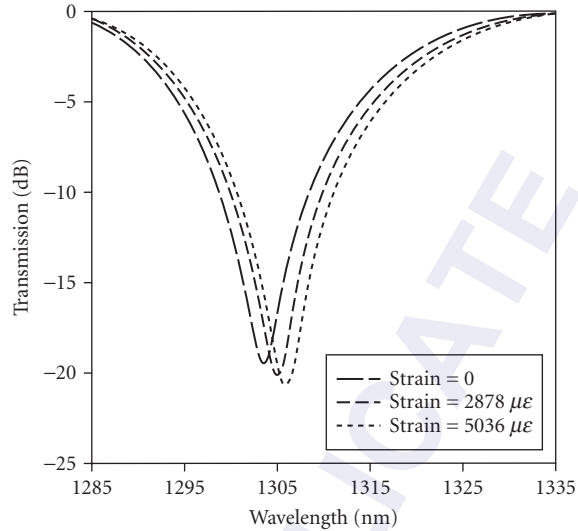


FIGURE 13 Strain-induced shift in a long-period grating fabricated in fiber C. The loss at 1317 nm increases by 1.6 dB due to the applied strain (5036 $\mu\epsilon$).

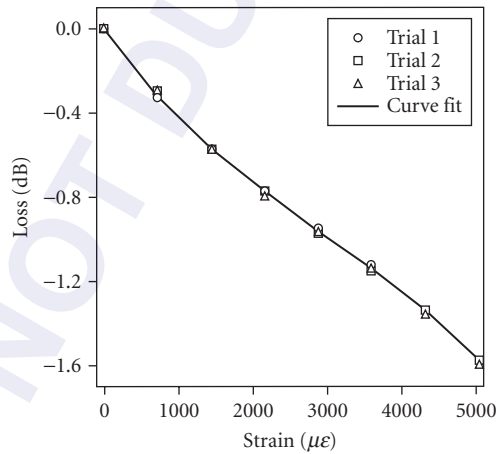


FIGURE 14 The change in the grating transmission at 1317 nm as a function of strain for three different trials. The increase in loss by 1.6 dB at 5036 $\mu\epsilon$ provides evidence of the feasibility of the simple setup used to measure strain.

On the other hand, a broadband source can also be used to interrogate the grating. At the output an optical bandpass filter can be used to transmit only a fixed bandwidth of the signal to the detector. The bandpass filter should again be centered on either side of the peak loss band of the resonance band. These schemes are easy to implement, and unlike the case for conventional Bragg gratings, complex and expensive interferometric demodulation schemes are not necessary.²⁰

TABLE 2 Temperature Sensitivity of Long-Period Gratings Written in Four Different Types of Fibers

Type of Fiber	Temperature Sensitivity (nm/°C)
A—standard dispersion-shifted fiber (DSF)	0.062
B—standard 1550-nm communication fiber	0.058
C—conventional 980-nm single mode fiber	0.154
D—conventional 1060 nm single mode fiber	0.111

Values correspond to the shift in the highest order resonance wavelength.

Temperature Sensitivity of Long-Period Gratings

Gratings written in different fibers were also tested for their cross-sensitivity to temperature.²⁰ The temperature coefficients of wavelength shift for different fibers are shown in Table 2. The temperature sensitivity of a fiber Bragg grating is 0.014 nm/°C. Hence the temperature sensitivity of a long-period grating is typically an order of magnitude higher than that of a Bragg grating. This large cross-sensitivity to ambient temperature can degrade the strain sensing performance of the system unless the output signal is adequately compensated. Multiparameter sensing using long-period gratings has been proposed to obtain precise strain measurements in environments with temperature fluctuations.¹⁹

In summary, long-period grating sensors are highly versatile. These sensors can easily be used in conjunction with simple and inexpensive detection techniques. Experimental results prove that these methods can be used effectively without sacrificing the enhanced resolution of the sensors. Long-period grating sensors are insensitive to input polarization and do not require coherent optical sources. Cross-sensitivity to temperature is a major concern while using these gratings for strain measurements.

24.6 COMPARISON OF SENSING SCHEMES

Based on these results, interferometric sensors have a high sensitivity and bandwidth but are limited by nonlinearity in their output signals. Conversely, intrinsic sensors are susceptible to ambient temperature changes, while grating-based sensors are simpler to multiplex. Each may be used in specific applications.

24.7 CONCLUSION

We have briefly summarized the performance of four different interferometric and grating-based sensors as representative of the very wide range of possible optical fiber sensor instrumentation and approaches. This analysis was based on the sensor head fabrication and cost, signal processing, cross-sensitivity to temperature, resolution, and operating range. Relative merits and demerits of the various sensing schemes were discussed.

24.8 REFERENCES

1. R. O. Claus, M. F. Gunther, A. Wang, and K. A. Murphy, "Extrinsic Fabry-Perot Sensor for Strain and Crack Opening Displacement Measurements from Minus 200 to 900°C," *Smart Mat. Struct.* **1**:237–242 (1992).
2. K. A. Murphy, M. F. Gunther, A. M. Vengsarkar, and R. O. Claus, "Fabry-Perot Fiber Optic Sensors in Full-Scale Fatigue Testing on an F-15 Aircraft," *App. Opt.* **31**:431–433 (1991).
3. V. Bhatia, C. A. Schmid, K. A. Murphy, R. O. Claus, T. A. Tran, J. A. Greene, and M. S. Miller, "Optical Fiber Sensing Technique for Edge-Induced and Internal Delamination Detection in Composites," *J. Smart Mat. Struct.* **4** (1995).

4. C. E. Lee and H. F. Taylor, "Fiber-Optic Fabry-Perot Temperature Sensor Using a Low-Coherence Light Source," *J. Lightwave Technol.* **9**:129–134 (1991).
5. J. A. Greene, T. A. Tran, K. A. Murphy, A. J. Plante, V. Bhatia, M. B. Sen, and R. O. Claus, "Photo Induced Fresnel Reflectors for Point-Wise and Distributed Sensing Applications," in *Proceedings of the Conference on Smart Structures and Materials*, SPIE'95, paper 2444-05, February 1995.
6. J. Sirkis, "Phase-Strain-Temperature Model for Structurally Embedded Interferometric Optical Fiber Strain Sensors with Applications," *Fiber Opt. Smart Struct. Skins IV*, SPIE **1588** (1991).
7. K. O. Hill, Y. Fujii, D. C. Johnson, and B. S. Kawasaki, "Photosensitivity in Optical Fiber Waveguides: Applications to Reflection Filter Fabrication," *Appl. Phys. Lett.* **32**:647 (1978).
8. G. Meltz, W. W. Morey, and W. H. Glenn, "Formation of Bragg Gratings in Optical Fibers by Transverse Holographic Method," *Opt. Lett.* **14**:823 (1989).
9. P. J. Lemaire, A. M. Vengsarkar, W. A. Reed, V. Mizrahi, and K. S. Kranz, "Refractive Index Changes in Optical Fibers Sensitized with Molecular Hydrogen," in *Proceedings of the Conference on Optical Fiber Communications*, OFC'94, Technical Digest, paper TuL1, 1994, p. 47.
10. R. Kashyap, "Photosensitive Optical Fibers: Devices and Applications," *Opt. Fiber Technol.* **1**:17–34 (1994).
11. D. Z. Anderson, V. Mizrahi, T. Ergodan, and A. E. White, "Phase-Mask Method for Volume Manufacturing of Fiber Phase Gratings," in *Proceedings of the Conference on Optical Fiber Communication*, post-deadline paper PD16, 1993, p. 68.
12. A. D. Kersey and T. A. Berkoff, "Fiber-Optic Bragg-Grating Differential-Temperature Sensor," *IEEE Phot. Techno. Lett.* **4**:1183–1185 (1992).
13. V. Bhatia, M. B. Sen, K. A. Murphy, A. Wang, R. O. Claus, M. E. Jones, J. L. Grace, and J. A. Greene, "Demodulation of Wavelength-Encoded Optical Fiber Sensor Signals Using Fiber Modal Interferometers," *SPIE Photon. East*, Philadelphia, Pa, paper 2594-09, October 1995.
14. M. G. Xu, L. Dong, L. Reekie, J. A. Tucknott, and J. L. Cruz, "Chirped Fiber Gratings for Temperature-Independent Strain Sensing," in *Proceedings of the First OSA Topical Meeting on Photosensitivity and Quadratic Nonlinearity in Glass Waveguides: Fundamentals and Applications*, paper PMB2, 1995.
15. K. O. Hill, B. Malo, K. Vineberg, F. Bilodeau, D. Johnson, and I. Skinner, "Efficient Mode Conversion in Telecommunication Fiber Using Externally Written Gratings," *Electron. Lett.* **26**:1270–1272 (1990).
16. F. Bilodeau, K. O. Hill, B. Malo, D. Johnson, and I. Skinner, "Efficient Narrowband LP₀₁ ↔ LP₀₂ Mode Converters Fabricated in Photosensitive Fiber: Spectral Response," *Electron. Lett.* **27**:682–684 (1991).
17. A. M. Vengsarkar, P. J. Lemaire, J. B. Judkins, V. Bhatia, J. E. Sipe, and T. E. Ergodan, "Long-Period Fiber Gratings as Band-Rejection Filters," in *Proceedings of Conference on Optical Fiber Communications*, OFC'95, post-deadline paper, PD4-2, 1995.
18. A. M. Vengsarkar, P. J. Lemaire, J. B. Judkins, V. Bhatia, J. E. Sipe, and T. E. Ergodan, "Long-Period Fiber Gratings as Band-Rejection Filters," *J. Lightwave Technol.* **14**(1):58–65(1996).
19. V. Bhatia, M. B. Burford, K. A. Murphy, and A. M. Vengsarkar, "Long-Period Fiber Grating Sensors," in *Proceedings of the Conference on Optical Fiber Communication*, paper ThP1, February 1996.
20. V. Bhatia and A. M. Vengsarkar, "Optical Fiber Long-Period Grating Sensors," *Opt. Lett.* **21**:692–694(1996).

24.9 FURTHER READING

- Bhatia, V., M. J. de Vries, K. A. Murphy, R. O. Claus, T. A. Tran, and J. A. Greene, "Extrinsic Fabry-Perot Interferometers for Absolute Measurements," *Fiberoptic Prod. News* **9**:12–3 (December 1994).
- Bhatia, V., M. B. Sen, K. A. Murphy, and R. O. Claus, "Wavelength-Tracked White Light Interferometry for Highly Sensitive Strain and Temperature Measurements," *Electron. Lett.*, 1995, submitted.
- Butter, C. D., and G. B. Hocker, "Fiber Optics Strain Gage," *Appl. Opt.* **17**:2867–2869 (1978).
- Sirkis J. S. and H. W. Haslach, "Interferometric Strain Measurement by Arbitrarily Configured, Surface Mounted, Optical Fiber," *J. Lightwave Technol.*, **8**:1497–1503 (1990).

HIGH-POWER FIBER LASERS AND AMPLIFIERS

Timothy S. McComb, Martin C. Richardson, and Michael Bass

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

25.1 GLOSSARY

Symbols

a	fiber core radius
F	filling factor of air holes in a PCF
G	signal gain in a fiber
$g(z)$	gain as a function of length
k_0	wavenumber in vacuum
M^2	beam quality factor; a value of 1 indicates a diffraction-limited beam
n	refractive index
N	population density of an energy level
P	power of laser emission
R	mirror reflectivity
V_{eff}	effective V parameter for a PCF
w_L	mode field radius
α	propagation loss
η_{pump}	pump overlap with doped region
η_{signal}	signal overlap with doped region
λ	wavelength
Λ	pitch or spacing of air holes in a PCF
ν	frequency
σ	emission or absorption cross section
τ	upper-state lifetime

Abbreviations and Definitions

Air-cladding	Region of air holes connected by thin glass bridges to the cladding region of a photonic crystal fiber forming an air-glass boundary used to guide pump radiation with high numerical aperture
AO	Acousto-optical
AR	Antireflective
ASE	Amplified spontaneous emission
CCC	Chirally coupled core, a type of fiber with a small satellite core chirally wrapped around the signal core; the small core couples higher-order modes out from the central core allowing larger mode areas
Cladding	The region of an optical fiber that surrounds the core; in conventional fibers this region is lower refractive index than the core to allow total internal reflection guidance
Core	The central region in a fiber where signal light is guided in most cases by total internal reflection
CPA	Chirped pulse amplification, a technique used to amplify ultrashort pulses whereby the pulse is temporally stretched before being amplified and recompressed after amplification in order to avoid high peak powers in the amplifier
DFB	Distributed feedback laser
Dichroic	An optical element that exhibits desired properties at two separate wavelengths (for instance, a mirror that is HR at one wavelength and HT at another)
EO	Electro-optical
FBG	Fiber Bragg grating
GG IAG	Gain-guided index antiguided fibers are fibers with core index less than that of the cladding but which can have gain in the core to compensate for losses
GMRF	Guided-mode resonance filter, a device consisting of a subwavelength grating on top of a waveguide layer used to form a narrow band reflectivity
HOM	Higher-order mode, any mode of a fiber of higher order than the fundamental mode
HR	High reflectivity
HT	High transmission
LMA	Large mode area; any of a number of techniques or technologies used to increase the mode field diameter of a fiber intended for use in an amplifier or laser system where high power operation requires larger mode field diameters to avoid nonlinear effects
MCVD	Modified chemical vapor deposition; a technique for fiber preform fabrication involving depositing chemical “soots” on the inside of a glass tube and subsequently collapsing the tube
MFA	Mode field adaptor; device used to change the mode field diameter of a fiber to match that of a second fiber
MFD	Mode field diameter; diameter of the distribution of radiation within an optical fiber usually at $1/e^2$ value of power
MOPA	Master oscillator power amplifier; system involving a low-power laser source (the oscillator) amplified by one or many amplifier chains
Multicore fiber	Fibers possessing several cores within a single cladding with each core designed to separately guide light and with these cores arrayed in such a pattern as to produce desired beam profile in the far field
NA	Numerical aperture of a fiber; a function of square root of the difference of the squares of the refractive index of core and cladding
OVD	Outside vapor deposition; a technique for fiber preform manufacture
PCF	Photonic crystal fiber; one of several fiber types with a latticelike structure of different refractive indices to create guidance in a fiber

PM	Polarization-maintaining fiber; fiber designed with one of a variety of stress-inducing structures to introduce birefringence in the fiber core
Pump cladding	The region of a double-clad fiber within which pump radiation is guided so it can cross through the doped fiber core, surrounded by a lower index glass or polymer layer
SBS	Stimulated Brillouin scattering
SESAM	Semiconductor saturable absorber mirror; a mirror with built-in saturable absorption that is often used to mode-lock fiber laser systems
SMET	Single-mode excitation technique; a technique where light from a laser source is launched into a fiber with appropriate care to launch only the fundamental mode of the fiber, even if the fiber itself is multimode
SPM	Self-phase modulation
SRS	Stimulated Raman scattering
TEC	Thermally expanded core; a technique for heating a fiber core causing the core dopants to diffuse, thus expanding it
TFB	Tapered fiber bundle; device combining multiple pump fibers into a bundle to deliver pump radiation to a fiber laser or amplifier
USP	Ultrashort pulse
VAD	Vapor axial deposition; a technique for fiber preform manufacture
VBG	Volume Bragg grating
V-parameter	Parameter that indicates guidance properties of a fiber; a value less than 2.405 indicates single-mode guidance; V is a function of wavelength, NA, and core radius
ZBLAN	A fluoride glass named for its chemical composition containing ZnF_4 , BaF_2 , LaF_3 , AlF_3 , and NaF

25.2 INTRODUCTION

Introductory Remarks

Although fiber lasers were first demonstrated at the dawn of the laser age,¹ it is only recently that fiber lasers have risen in visibility on the landscape of laser development. This resurgence arose as a consequence of transformational changes in pump technology, and fiber design and fabrication techniques. Thus nowadays a fiber laser can be thought of as a device for the conversion of light from low brightness laser diodes to high brightness, highly coherent laser light. Their rise in power and brightness capabilities has been so significant that they are now beginning to invade the applications space once dominated by solid-state lasers. As these devices permeate many fields of laser applications in manufacturing medicine and defense, there are growing demands and constraints placed on the fiber laser's spectrum, beam quality, and pulse duration. These demands can only be met with solid-state lasers at the expense of efficiency, complexity, and cost. However, they are in many cases a natural consequence, or relatively simple modification of modern fiber lasers.

As industrial, medical, and defense applications of the high-power lasers force increasing levels of electrical efficiency, beam quality, and ruggedness with commensurate reductions in cost, complexity, and footprint, fiber lasers will increasingly meet these needs. In this chapter we summarize the basic principles of fiber lasers, the latest developments in design and fabrication, their different modalities in output characteristics, and their potential for future growth in output power and overall utility.

A Brief History of Fiber Lasers

Although most of the developments in fiber lasers did not occur until the onset of the fiber optics telecommunications and the development of high-power optical diodes, it should not be forgotten

that the basic concept of the fiber laser, that of a doped fiber core surrounded by optically transparent cladding was devised by Snitzer in the early 1960s at the very birth of the laser age.¹⁻³ The current rapid growth in high-power fiber laser achievements owes its origins to (1) the development of high quality silica fibers and (2) the development of high-power diode laser technology.

The rapid growth in telecommunications and the large investments in improving optical fiber technology spurred new interest in fiber lasers in the 1970s. Multimode core-pumped fiber lasers were demonstrated, pumped by both diode lasers and bulk lasers.^{4,5} Silica became the standard host material for most fiber lasers, a direct extension of fiber optics telecommunications technology. Despite the improving optical quality of silica fibers, the first single-mode fiber laser was not demonstrated until the mid-1980s. This laser was still of low power, due to the unavailability of high brightness pump sources for the directly core-pumped scheme.⁶ The 1980s also saw the first tunable and Q-switched fiber lasers.^{7,8} At that time fiber laser development was largely driven by the potential of (erbium) Er-doped fiber at approximately 1.5 μm ,^{9,10} to serve as low-loss, low-dispersion, optical amplifiers¹¹ transformed the telecommunications industry.

High laser power single-mode fiber lasers beyond approximately 1 W of output power, at this time were limited by available single-mode diode pump sources required for the core-pumping scheme. It was difficult to efficiently couple light from then-available diode sources into single-mode fiber cores due to the diode's elliptical beam output shape and limitations on single transverse mode diode output powers. The most significant advance to overcoming this problem has become the basis for all high-power fiber lasers manufactured today, the so called "double clad fiber" was first proposed in 1988 by Snitzer.¹² Shown schematically in Fig. 1, the introduction of a second, undoped inner cladding of much larger diameter (typically $> 100 \mu\text{m}$) than the doped core allowed for effective coupling of much higher pump powers. So long as this inner cladding layer had a refractive index less than that of the core, there was effective coupling of the pump light into the core region. Similarly the outer cladding layer must have a refractive index less than that of the inner cladding region, in order to limit the loss of pump light from the fiber.

This innovation led to a period of rapid increase in fiber laser output power, limited only by pump diode power.¹²⁻¹⁵ The double clad fiber enabled scaling to 110 W average powers and approximately 100 μJ peak powers in purely single-mode fiber cores. At these power levels, however, further increases in the power of fiber lasers were impeded by a set of new problems associated with nonlinear optical effects and fiber damage in the very small single-mode cores.^{16,17} New concepts for increasing the size of the fiber laser core, now commonly referred to as large-mode-area (LMA) fibers were demonstrated in Refs. 17 to 19. LMA technologies have given birth to the state-of-the-art high-power fiber lasers that exist today, enabling continuous wave (CW) lasers to reach powers of more than 3 kW and pulsed fiber laser to reach more than 4 mJ output power in nanosecond pulses with diffraction-limited beam quality.^{20,21}

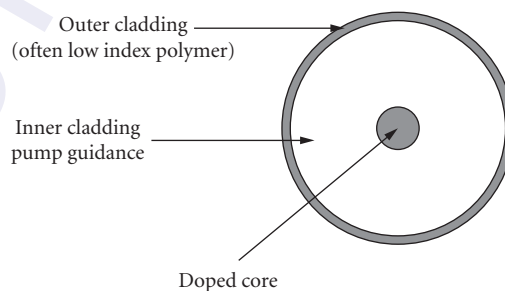


FIGURE 1 Simple schematic of a double clad fiber. The refractive index profile should have the following index relation: $n_{\text{core}} > n_{\text{inner cladding}} > n_{\text{outer cladding}}$

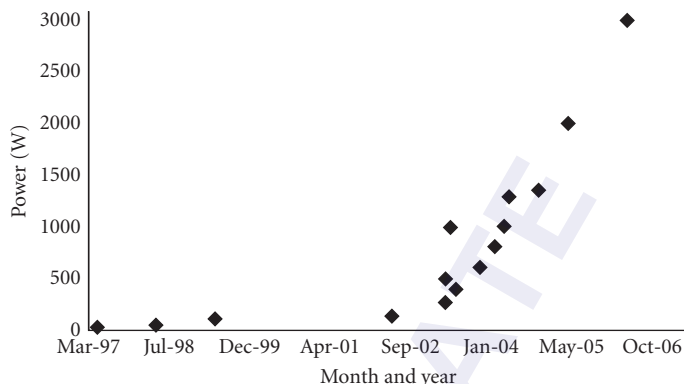


FIGURE 2 Plot of highest achieved CW output powers in fiber lasers operating at 1 μm wavelength dating from 1997 to the present day based on Refs. 16 and 21 to 34.

Rapid Growth of Fiber Lasers

Beginning with the development of the first double clad fiber laser and continuing with LMA fiber laser technologies, the output powers of fiber lasers has undergone near-exponential growth in the last 10 years. A convenient benchmark for such growth can be seen in the output power of CW Yb fiber lasers beginning around 1997. Figure 2 shows the highest reported output powers of Yb CW fiber lasers over a period of years.

The LMA fiber technologies described here are critical to the future growth of high-power fiber lasers. In the future, the development and refinement of LMA technologies will lead to further growth in output powers, with 10 kW being a reasonable goal in the not-too-distant future. Indeed, recently a paper outlining the limitations in output power from single-mode Yb-based fiber lasers has approximated the upper limit on output power at 36 kW from a single fiber based on today's technology and reasonable assumptions for future growth.³⁵

Comparison of Fiber Lasers to Bulk Lasers

High-power fiber lasers possess several advantages when compared to bulk solid-state lasers. These include their compact, simple construction, simplified thermal management, extremely high beam quality, and high optical-to-optical efficiency. These beneficial properties are forces driving the use of fiber lasers for high-power applications.

Fiber laser systems offer a clear advantage in terms of their size and performance in harsh environments. This follows from the fact that fiber lasers are often completely monolithic. That is, they are comprised of an unbroken chain of all-fiber-based components with no need for realignment and no potential for contamination. Though some very high-power kilowatt-class fiber-based systems have not yet been made completely monolithic and still require some free-space components, the fact that a majority of the cavity exists "in fiber" still provides a stability benefit compared to bulk solid-state lasers. In fact the highest power systems demonstrated have been completely monolithic systems.^{21,32}

Fiber lasers also have a distinct advantage over bulk lasers in terms of thermal management resulting from the much larger surface area of fiber lasers compared to bulk lasers.^{36,37} This enables the heat deposition resulting from the quantum defect of the pump light to be distributed over the full length of a fiber. Most fiber lasers can operate with only minimal attention to thermal management, though cooling and attention to thermal issues in fiber lasers cannot be completely ignored for reasons including polymer coating integrity and laser efficiency as suggested by studies done in Refs. 38 to 40.

Often, high-power applications call for the power to be delivered at long distances from the laser output itself or for light to be focused to extremely small spot sizes leading to a requirement for diffraction-limited beams. In a fiber laser, the transverse modes are defined by the fiber itself. Correct fiber design ensures that no other modes besides the lowest-order transverse mode are allowed to exist in the waveguide, the core of the fiber, leading to excellent output beam quality. Bulk lasers are often victims of optical distortions due to thermal lensing and birefringence in the gain media caused by temperature gradients and the temperature sensitive nature of the refractive index. In general, fiber lasers are mostly immune to thermal gradient induced optical distortion because of more efficient heat removal and the wave-guiding nature of the fiber.

25.3 Fiber Laser Limitations

Despite the many advantages of fiber lasers, fiber-based systems also have some limitations that require further research and development to mitigate. These limitations all stem from the small core sizes of fibers.

Optical Damage

The most obvious limitation is the damage threshold of the fiber core material due to the laser high-power density in the relatively small fiber core area. The bulk damage threshold of silica is extremely high ($\sim 600 \text{ GW/cm}^2$ at $\sim 1000 \text{ nm}$ wavelength) though tightly focused pulses can damage the bulk material. However, damage most easily occurs when light exiting the fiber reaches the surface damage threshold, which in silica is approximately 40 GW/cm^2 .⁴¹ The most obvious way to mitigate this damage threshold is to increase the core size; however, due to the need to maintain beam quality this technique has limitations. An alternative method for damage mitigation in fiber amplifiers is end capping, a process which involves splicing a coreless short section of fiber onto the end of a fiber, allowing the expansion of the beam before reaching the glass-air interface.

Nonlinear Effects

A second issue in high-power fiber lasers is the result of detrimental nonlinear effects. Such effects are based on third-order nonlinearities in the glass such as self phase modulation, stimulated Raman scattering (SRS), stimulated Brillouin scattering (SBS), and self focusing. Details on the origins and background of such effects can be found in Chap. 10, "Nonlinear Effects in Optical Fibers," in this volume and Chap. 15, "Stimulated Raman and Brillouin Scattering" in Vol. IV. The impact and severity of these nonlinear effects varies depending on the laser type. Narrow linewidth lasers suffer from unwanted spectral broadening at high powers caused by stimulated Brillouin scattering. Pulsed lasers can experience the effects of Raman scattering. Ultrashort pulse lasers can experience pulse distortions based on self phase modulation. In addition, high average powers cause a hard limitation due to the onset of self focusing in bulk glass leading eventually to catastrophic damage. With the exception of self-focusing, which is core size independent; other nonlinear effects can again be mitigated by simply increasing core diameter. In addition, other techniques such as acoustic design of fiber and thermal control of a fiber can be used to reduce such effects as SBS.⁴²⁻⁴⁴

Energy Storage

A limitation in high power fiber lasers with high energy pulses is the small gain volume of the doped core. Even in a very long fiber the total volume of gain medium is small. Fiber lasers that must have high pulse energies find a limitation in terms of the capability to store sufficient energy in the fiber core.

In addition, the leading edge of a pulse can “steal” or saturate the gain in fiber amplifiers and the pulse itself can become temporally distorted. Such issues can be mitigated by of course increasing the fiber core diameter, leading to more gain medium volume and, consequently, higher energy storage. In addition, pulse deformation can be circumvented by using an input pulse designed to compensate for the deformation with a shape designed to accommodate for the nonuniform temporal gain.

CW Damage Threshold

The damage threshold of extremely high-power CW fiber lasers is reviewed in Ref. 35. There each of the previously discussed damage mechanisms, as well as additional damage considerations are considered in order to make an estimate of maximum achievable power from a single-mode single-fiber device operating at a wavelength near 1000 nm.³⁵ Using this reference as a basis the reader can also make some extrapolations about the damage threshold of fiber lasers at longer, eyesafe wavelengths. In many cases it appears that such longer-wavelength fiber laser (most notably at 1.5 and 2 μm) may have improved power handling capabilities in many situations compared to their 1- μm counterparts.

25.4 Fiber Laser Fundamentals

Fiber Laser Operation

The fundamentals of fiber laser operation and amplification are the same as those in any laser system. Equations describing laser operation of a bulk laser can be adapted to fiber lasers by taking into account the wave-guiding nature of the fiber and using the appropriate parameters. A detailed description of laser operation as a whole can be found in Chap. 16, “Lasers” and Chap. 23, “Quantum Theory of the Laser,” in Vol. II of this *Handbook*. Included here are only a few equations useful to the design and operation of fiber lasers.

It is convenient to describe the operation of a fiber gain medium by a set of rate equations. The following define the operation of a simple fiber laser system that, depending on the sign selected, can have pump and signal light traveling in either direction in the fiber,

$$\begin{aligned} \pm \frac{\partial P_{\text{pump}}^{\pm}}{\partial z} &= \eta_{\text{pump}} (\sigma_{10}(\lambda_{\text{pump}}) N_1 - \sigma_{01}(\lambda_{\text{pump}}) N_0) P_{\text{pump}}^{\pm} \\ \pm \frac{\partial P_{\text{signal}}^{\pm}}{\partial z} &= \eta_{\text{signal}} (\sigma_{10}(\lambda_{\text{signal}}) N_1 - \sigma_{01}(\lambda_{\text{signal}}) N_0) P_{\text{signal}}^{\pm} \end{aligned} \quad (1)$$

Here the \pm indicates direction of propagation, N_i is the energy of a given level i , σ_{ij} is the emission or absorption cross section from level i to level j , P is power of a given signal, η_{pump} is the pump overlap with the doped region, and η_{signal} is the signal overlap with the doped region.^{35,46} This model does not take amplified spontaneous emission (ASE) or temporal effects into account. Equation (1) can be solved by applying slightly more sophisticated modeling discussed in Refs. 45 and 46. Solving the coupled differential Eq. (1) numerically, one can obtain the change in pump and laser signal power over distance along the fiber and in time. In addition, by using the appropriate boundary conditions these equations can be used to model amplifiers in the co-, counter- and bidirectional propagation of pump light and also oscillators.

Considering the signal power in the fiber, the total gain in a fiber is the integral of gain along the fiber length where the gain at any given point along the fiber is given by

$$g(z) = [\sigma_{10}(\lambda_{\text{signal}}) N_1(z) - \sigma_{01}(\lambda_{\text{signal}}) N_0(z)] \quad (2)$$

Thus, the signal gain in a fiber can be expressed as

$$G = \exp \left\{ \frac{[\sigma_{10}(\lambda_{\text{signal}}) + \sigma_{01}(\lambda_{\text{signal}})] \eta_q \tau_{10} P_{\text{absorbed}}}{h\nu_{\text{pump}} A_{\text{core}}} - \sigma_{01}(\sigma_{\text{signal}}) N_0(z) L - \alpha_L L \right\} \quad (3)$$

with α_L the core propagation loss, and all other terms as defined earlier.⁴⁷ The final expression to be defined is the laser saturation power given by⁴⁷

$$P_{\text{satsignal}} = \frac{h\nu_{\text{signal}} A_{\text{core}}}{[\sigma_{10}(\lambda_{\text{signal}}) + \sigma_{01}(\lambda_{\text{signal}})] \tau_{10}} \quad (4)$$

This saturation power can be used to describe the change in the small signal gain as the power is increased. In addition, in order for the amplifier to have efficient energy extraction, the signal input must be on the order of the saturation power.

Another important relationship for fiber lasers describes the relative output power from either end of a fiber. Knowledge of the optimum amount of feedback required into one end of the fiber is useful to ensure a majority of output from the other end. To determine the exact output power performance of a fiber laser, either solutions to coupled differential equations for signal round trip propagation in a fiber or a detailed Rigrod analysis⁴⁸ must be computed. However a simplified expression can be written if certain assumptions are made, and the ratio of output power from each fiber end as a function of their reflectivity can be written as follows:

$$\frac{P_1}{P_2} = \frac{1 - R_1}{1 - R_2} \sqrt{\frac{R_2}{R_1}} \quad (5)$$

where $R_{1 \text{ or } 2}$ is the reflectivity of end 1 or 2. Clearly if the reflectivity of one end is far greater than the other, the output power will be predominantly from the lower reflectivity end.^{46,48} Since most simple fiber lasers use approximately 4 percent Fresnel reflection at one end as an output coupler (and even those that employ fiber Bragg grating output coupler use reflectivities typically in the range of 4 to 15%), only relatively low feedback (much less than 90%) is required on the opposing end to ensure that a majority of power leaves from the low reflectivity end.

Important Fiber Equations

A second set of equations related to fiber lasers govern the guidance of light in the fibers themselves. The fundamentals of light guidance in optical fibers can be found in Part 4, "Fiber Optics," Vol. V and Ref. 49. Since LMA fiber lasers are currently the most promising for substantial output powers the following discussion deals with guided light propagation in this type of fiber laser.

The fiber V parameter is defined as

$$V = k_0 a \sqrt{n_1^2 - n_2^2} \quad (6)$$

where is k_0 the wavenumber and a the core radius and with square root of the difference of the squares of n_1 and n_2 alternately expressed as the fiber NA (numerical aperture). The V parameter can be used to determine the guiding properties of a fiber with a given core and numerical aperture. A fiber core with a V parameter of less than 2.405 can sustain only a single lowest-order mode. Most LMA fibers have a V parameter closer to 3 or 4, with large core radii a and small numerical apertures NA. Weak guidance of higher-order modes in LMA fibers allows them to be stripped out by using techniques that provide preferential loss to higher-order modes discussed later and hence enables single-mode propagation in a multimode fiber. The control of the NA and hence the V parameter is therefore a principal factor in LMA fiber laser design.

It is important to note that the true mode area of a fiber does not correspond to the actual core size of the fiber, due to the multimode nature of LMA fiber cores fundamental mode diameter (also known as mode field diameter) may be smaller than the actual core diameter. In the case of truly single-mode fibers the mode field diameter may actually be larger than core diameter. Hence, when quoting mode area (the important factor in determining damage threshold and nonlinearities in fiber lasers) one should use the actual size of the lowest-order mode, not the core size. This is especially important when dealing with fiber splicing or mode field adaptation between dissimilar fibers where splice or device losses must be minimized. This mode diameter can be approximated by the equation

$$w_L \approx a \left(0.65 + \frac{1.619}{V^{3/2}} + \frac{2.879}{V^6} \right) \quad (7)$$

where w_L is the mode radius, V is the V parameter and a is core radius, it should be noted that this is valid for the mode field radius of the fundamental mode in the cases where V is more than 2.405.⁵⁰

We now consider briefly the highly multimode regime. In this case hundreds to thousands of fiber modes can exist in different mode families. The actual number of modes in any fiber is proportional to V^2 , so for fibers with very large core and/or NA, a very large number of modes can propagate.⁵¹ It is a property of light in fibers that as one goes to higher mode number, the location of energy in the higher-order modes moves outward from the center of the fiber. As a consequence, when light is propagating in a highly multimode fiber (for instance, the pump cladding of a double clad fiber), there is a significant portion of energy that may never cross into the core. This can be both a positive effect in terms of evanescent pump coupling and a negative effect in terms of helical modes not being absorbed in the core of double clad fibers. As consequence, one must be aware of not only the number of modes but their shape within the fiber when making design considerations in a fiber laser system.

25.5 FIBER LASER ARCHITECTURES

Fiber laser architectures have evolved as improvements in pump coupling technologies and laser diode brightness, and fiber-based components have driven higher power, and more efficient, compact, stable, and robust platforms. Originally so-called “free-space coupling” using conventional optical elements to image the light output from diode lasers or diode laser bars (lenses, prisms, etc.) was the only effective fiber laser architecture for high-power systems and it still offers advantages for developmental systems. However, there are inherent advantages to “all-fiber” or “monolithic” systems in simplicity, efficiency, stability, compactness, and cost. Here the relevant diode and fiber technologies are described before free-space pumping architectures are laid out, and the latest all-fiber approaches are summarized.

Pumping Techniques

After the fiber itself, the most critical element of a fiber laser system is the pump delivery technique. Since diode lasers are used to pump fiber lasers, this section describes how low-brightness diode laser light is coupled into the core of a fiber laser or amplifier. Efficient fiber lasers rely on effective launch of high average diode laser pump power into the fiber. Unabsorbed or improperly launched pump light results in waste heat that can cause thermal damage to polymer coatings on the fiber, to the pump combiners, to other devices in the system or to the fiber itself. It is also a waste of expensive diode laser pump power. As a result in any fiber laser system, whether free space coupled or monolithic, the design of the pumping system is critical. Fiber-pumping design can be divided into issues of core-clad configuration, choice of diode laser source, and actual pump power delivery method into the fiber.

Pump Launch Schemes By far the simplest and most widely used pump scheme is end pumping by imaging the diode laser delivery fiber (or diode output facet itself) onto the end of the fiber laser or

amplifier. Two lenses are used, designed for minimal aberrations, resulting in small pump spot size for efficient coupling. The pump light must be delivered into a spot less than the diameter of the fiber and within the fiber pump cladding's numerical aperture. The output of bare diode bars or stacks can be shaped to take their long, narrow distributions and make them approximately circular while conserving brightness. Two methods for this both involve using fast and slow axis micro lenses on bars and then using either stacked glass plates to slice the beam into sections and stack the sections on top of each other,⁵² or by using two mirrors to restack the different parts of the beam.⁵³ These methods allow for the use of bare diode bars or stacks rather than fiber-coupled diodes, thus providing lower cost alternatives compared to fiber-coupled bars, and potentially higher available powers than fiber-coupled methods. Experiments have demonstrated the use of diode bars to directly launch more than 1 kW into a fiber end.²⁷

Though end-pumped schemes are capable of providing kilowatt-level incident pump powers, other methods allowing more uniform distribution of pump light along the length of a fiber have been investigated. One approach uses various techniques to inject pump light along the length of a fiber by either micropisms or V-groves etched into the side of the fiber.^{54,55} This approach allows pump power to be distributed along a very large length of fiber and requires minimal alignment, compared to free-space end-pumping techniques, but its implementation is more difficult and time consuming with the need for multiple etches or other fiber preparation. Moreover, in the event of fiber damage, the entire array must be replaced, rather than simply the fiber as is the case in end-pumped configurations.

Directly side pumping a fiber core has also been attempted, however, obtaining sufficient absorption over a few micron fiber core is challenging. Some fibers have been created with flattened sides allowing them to be coiled tightly in a spiral, thus permitting pump light from diode bars to enter from the side and pass through multiple fiber cores in the spiral, allowing for higher efficiency pumping with relatively simple to use diode bars.²³ While allowing minimal alignment, simple packaging and easy thermal management, difficulties with this technique arise in manufacture, assembly, and repair of such fibers.

Free-space end-pumping schemes work very well but have the disadvantage of being limited to only being able to pump the two ends of the fiber, limiting the amount of pump power to that available in a single diode bar, stack (or two in polarization combined cases), or fiber-coupled device. In addition, launching very high pump power into small fiber ends may cause thermal damage to the fiber. To avoid this, fibers may need to be end capped or have actively cooled ends, especially when using nonsilica glass fibers which have lower melting points.

Monolithic or "all-fiber" fiber lasers best realize fiber's potential advantages over bulk solid-state lasers. There are several main techniques to achieve monolithic pumping, each with particular advantages and disadvantages. Telecommunication lasers utilize waveguide-based "y" couplers or wavelength multiplexers to achieve pump and signal combination; however, these devices are optimized for low powers and for operation of core-pumped lasers with single-mode pump components not capable of the power levels needed for high-power systems. Alternative techniques must be used in high-power double-clad fibers.

The most straightforward technique for coupling pump light is simply splicing a single pump fiber to the end of a gain fiber (of course with a fiber Bragg reflector in between to form a resonator). This method is low loss and very simple to implement, however, the main downfall of this technique is that only one fiber can be spliced to the laser, and hence power delivery is limited to the power available in a reasonably small (100 to 800 μm) delivery fiber. To achieve higher pump powers, devices with multiple input ports are required. One of the most common current methods for launching pump light from multiple sources is by way of tapered fiber bundle (TFB). TFB's are simple in concept as seen in Fig. 3, involving bundling a number of undoped, coreless pump delivery fibers around a central fiber (possibly with a core in amplifier configurations), fusing the whole bundle together with high heat, and tapering it down to an appropriate diameter to allow it to be spliced into a system.⁵⁶⁻⁵⁸

Several methods for TFB design and manufacture and their evolution in time are pointed out in Refs. 56 to 58. TFBs are capable of handling upward of 1 kW of optical power with the addition of proper thermal management techniques as demonstrated in Ref. 59.

A concept similar to TFBs, the nonfused fiber coupler, involves angle polishing a fiber to an appropriate angle and butting it to a flattened portion of a second fiber.⁶⁰ The benefit of this technique is that no fusion splicing is required; the fibers must only be polished and butted together. Again the concept of such fibers is simple and the number of inputs is scalable and power can be distributed along the whole fiber length, however, the implementation is difficult as precise angle polishing is required.

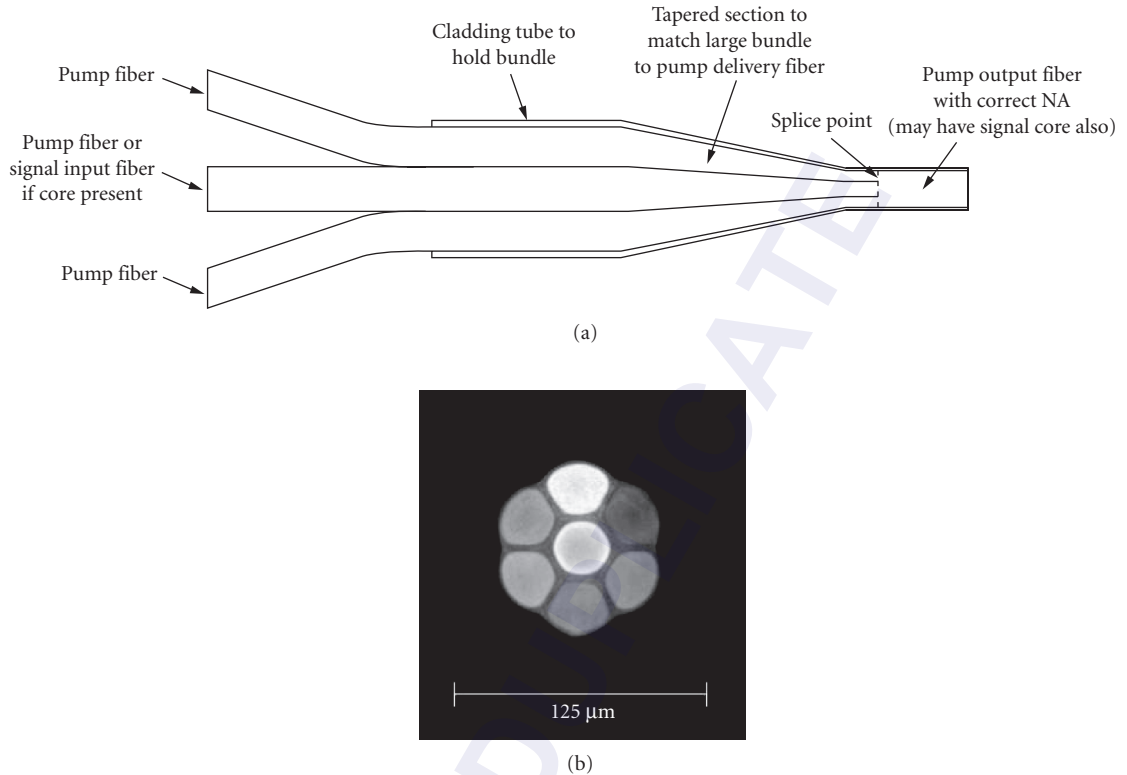


FIGURE 3 (a) Schematic of tapered fiber bundle. (b) End view of actual tapered fiber bundle.

TFBs are excellent options for most free-space fiber systems. However, in many cases lasers need to be very short, or of soft glass materials that cannot be spliced to silica glass TFBs efficiently. It may also be desirable to distribute pump light over the length of the fiber more uniformly to reduce “spot heating” effects. Thus an alternative all-fiber technique is evanescent coupling, also called GTWave technology. This technique consists of placing two fibers in optical contact by either stripping polymer layers and manufacturing by hand, or by pulling two fibers together into optical contact and covering them in one polymer coating.^{61–63} This technique’s ability to work for short, highly doped fibers, or distribute pump light very evenly makes it an attractive technique in some applications.^{62,64,65} The method also does not require any special equipment, only bare fibers and some method to hold them together such as heat-shrink tubing (though custom fiber drawing techniques can also be used), in contrast to the splicers and cleavers needed to make and implement TFBs. Evanescent coupling can also work in situations where TFBs cannot be made or purchased, such as with soft glass fibers or very large mode area fibers. In addition when fabricated as all-fiber systems pulled together on a draw tower this technology is used in several high-power commercial systems.⁶⁵

As development of high-power fiber lasers continues, even higher-power handling “all-fiber” components will become available and this will result in the monolithic fiber laser system increasing in power and beginning to function at power levels comparable to high-power free-space systems currently available. Currently the development of pump-combining devices and techniques continues as new methods are developed or existing methods improved or combined together to provide greater power handling, efficiency, and reliability of pump components. This section is only a brief introduction to the basic categories of pump components and techniques.

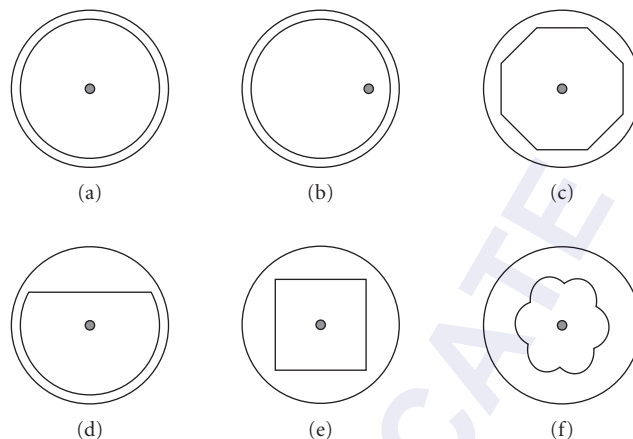


FIGURE 4 Different fiber core geometries for enhanced pump absorption: (a) unmodified fiber; (b) offset core fiber; (c) octagonal (or otherwise polygonal) fiber; (d) “D” shaped fiber; (e) square or rectangular fiber; and (f) “flower” shaped fiber.⁶⁶

Pump Cladding Design Pumping in double clad optical fibers relies on the absorption of pump light by a doped core as it propagates in an outer cladding. However, the circular symmetry of conventional fiber poses a challenge to pump absorption due to excitation of helical pump modes that allow light to propagate down the fiber without crossing through the core. In the double clad pumping structure this is an issue because pump light will never have the opportunity to be absorbed by the doped core. The amount of pump light in these higher-order modes is significant and must be dealt with for efficient laser operation. The simplest technique is to coil the fiber though many fiber types cannot be bent to sufficiently small diameters. A polarization maintaining (PM) fiber has stress rods built into it which have the effect of breaking the symmetry of the fiber to reduce helical modes in cladding pumped fibers. For particular applications PM fibers are not always available. As a result other methods have been developed.⁶⁶ The most common way to reduce higher-order modes is to build a design into the fiber which breaks the circular symmetry by using any number of different cladding shapes such as shown in Fig. 4.

While all these designs are effective in causing higher pump absorption, some have more useful features than others as noted in Ref. 66. Overall, choice of shape is a matter of preference and exactly what purpose the fiber laser will have. It involves making compromises between effectiveness of design, ease of fabrication, and ability of the fiber to be spliced to conventional round fibers. It is, nevertheless, an important choice in the final design of a high-power fiber laser.

Choice of Diode Pump Choice of the appropriate type of diode pump for a particular fiber laser application is critical to optimizing the performance of the laser system. There are two choices to consider when determining what type of diode to use in the construction of a fiber laser. One can use single emitters and combine them to reach the desired powers using pump combiners or other techniques outlined earlier or one can use diode bars or stacks of bars for pumping. A more detailed discussion of this topic can be found in Ref. 59. Diode bars tend to simplify a system, as they are available packaged and deliver high powers, while a larger number of single emitters must first be spliced together via a tapered fiber bundle or other method before high powers can be achieved. Single emitters may be lower in cost per watt and also may be more efficient, meaning lower cooling loads and often need only air cooling. However, the cost savings of single emitter units versus bars may be offset by the time and effort required to splice and wire individual single emitters together.⁵⁹ Fiber-coupled bars tend to be less sensitive to thermal fluctuations causing wavelength shifts, with bars wavelength only fluctuating a few nanometers over temperatures, while single emitters can change by 10 to 20 nm.⁵⁹ Optical damage can occur in single emitters since the fiber pigtailed to them leads directly to the diode facet. Because of

the free-space optics inherent in fiber-coupled bars it is more difficult for stray back-reflections to cause damage and easier to integrate dichroic elements to mitigate the problem.⁵⁹ Single emitters have the advantage of being able to be electrically modulated far more rapidly than diode bars and so are superior in systems where pump diodes must be pulsed, especially at high repetition rates. Single emitters also tend to offer slightly longer lifetimes than diode bars.⁵⁹ Overall, there is no clear choice between diode bars and single emitters. The choice depends on the particular laser application and is usually an engineering decision that weighs all factors involved in using a particular technology.

Free-Space Fiber Laser Designs

Free-space fiber systems can take two basic forms: resonators and amplifiers. Laser resonators are common when moderate- to high-power systems are desired with the maximum simplicity and compactness. High-power resonators offer the ability to use a minimal number of components to achieve desired laser performance. However, they are often limited by difficulties in achieving high powers in particular spectral or temporal modalities. When an oscillator cannot fulfill a particular need, multistage amplifier systems with low-power master oscillators, providing the pulse or CW linewidth characteristics, and one or several power amplifier stages, providing the energy, are turned to. Known as master oscillator power amplifiers (MOPAs), such systems have higher potential for performance but are more complex and expensive to construct due to the need for additional components and complexities. The basic architectures of both types of systems are discussed in the following sections.

Free-Space Oscillator Architectures The most basic of all free-space cavities can be formed by use of a diode pump source and polished or cleaved fiber.^{36,37} The resonator is formed between approximately 4 percent Fresnel reflections (in the case of silica fibers with refractive index of ~ 1.45) of the two cleaved or polished fiber facets. The high gain of the fiber laser allows it to overcome very high losses and oscillate despite the small amount of feedback. (In fact, one of the main difficulties of making high-power fiber amplifiers is keeping the fibers from “parasitically lasing” since with high enough pump powers even angled fiber ends can give enough feedback to promote oscillation.) This simple cavity is not particularly suited for many applications due to its low efficiency, high threshold, and laser emission being split almost evenly between both ends. However, this configuration is useful for simply testing the free-running characteristics of prototype fibers and is also useful for providing initial alignment of a more complicated cavity.

The next step in sophistication of fiber laser resonators involves placing a mirror at one end of the cavity as shown in several forms in Fig. 5.^{36,37}

The simplest form of the single mirror resonator; Fig. 5a uses a dichroic mirror; a mirror that is highly reflective at the laser wavelength and highly transmitting at the pump wavelength to enhance the cavity Q, lower laser threshold, and provide output from only one end of the resonator. This technique works well though at high power there is a potential for mirror damage due to the high intensity incident on the mirror.

Damage to mirrors can be avoided by separating them from the end of the fiber and placing a lens between the mirror and fiber facet as in Fig. 5b and 5c. The scheme in Fig. 5b is clearly the simpler of the two schemes with pump light and laser light both traveling through the same lens. For many fiber lasers systems, especially those based on erbium or thulium where the pump and laser wavelength are significantly separated, chromatic aberration in the lens causes the pump and signal to focus at different points. As a result it is difficult to perfectly align the laser beam to complete the resonator in Fig. 5b and efficiently couple the pump light with a single lens (unless the pump delivery fiber is significantly smaller in diameter than the double clad laser fiber). Cavities based on Fig. 5c are also useful when adding additional elements to the resonator such as Q-switches. It should be noted that fiber lasers usually can operate with only Fresnel feedback from their output ends; however, sometimes a higher reflectivity output coupler must be added to the resonator to provide sufficient feedback to allow laser oscillation. Other more exotic and less common resonators for fiber lasers are free-space ring type resonators which are often used in mode-locked fiber lasers as in Refs. 67 and 68. Rings are often more difficult to keep stably aligned and require additional elements such as isolators and wave plates to ensure unidirectional operation.

Basic Free-Space MOPA System Designs Amplifiers for high-power lasers differ greatly in concept and design from those intended for telecommunications applications. Communications amplifiers

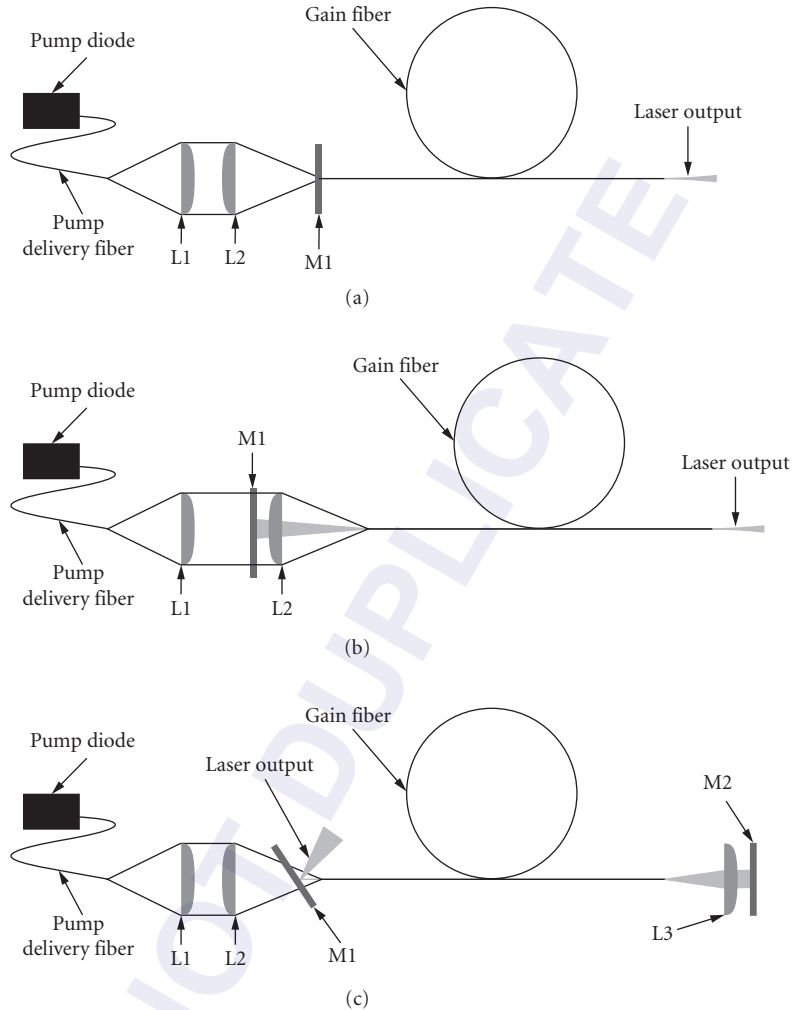


FIGURE 5 Note in all images L1 and L2 are pump coupling optics, L3 is a collimating lens for the signal, and M1 and M2 are dichroic mirrors. (a) Simple butt-joint mirror cavity. (b) External resonator for lasers with pump and signal close in wavelength, mirror outside cavity provided feedback through lens. (c) External feedback cavity on opposite end from pump to compensate for large wavelength difference between pump and signal. Bulk cavity elements like Q-switches elements can be placed between L3 and M2.

operate in the small-signal regime, amplifying small signals from noise with high fidelity and thus demand very high gains. High-power lasers demand efficiency from amplifiers with a desire for maximum extraction of energy. Master oscillator or power amplifier designs (MOPA) are often used to avoid difficulties associated with the maintenance of particular laser parameters in resonators at high powers (i.e., narrow linewidths or short pulse durations).

There are two basic configurations for free-space MOPAs designated as copropagating and counter-propagating MOPAs. These two schemes should produce identical results as they both involve a fixed gain per unit length for an input signal. In practice, however, the two schemes are different in their efficiency at high power and their ASE characteristics. The essential difference between the two

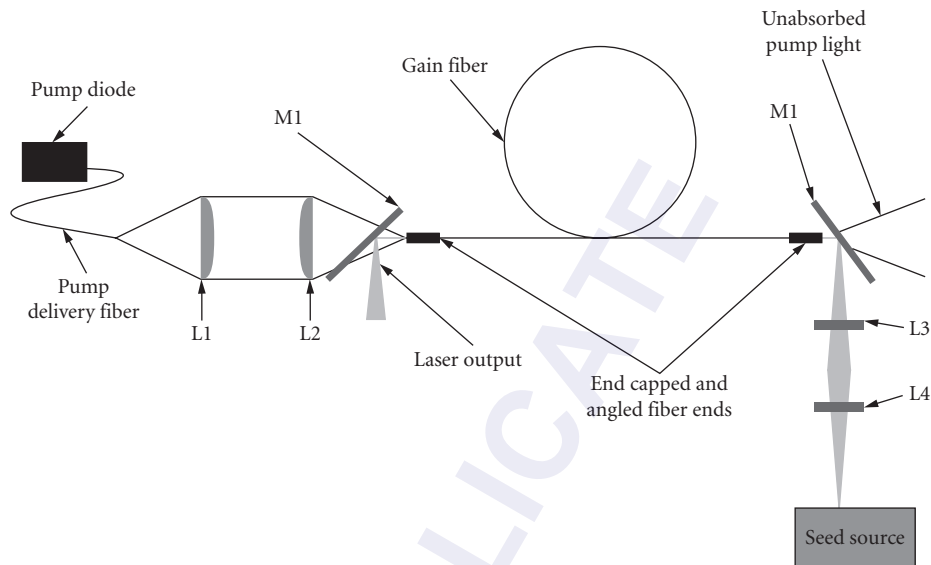


FIGURE 6 Counter-propagating MOPA scheme. Lenses L1 and L2 are matched pump delivery lenses. M1 is a dichroic mirror and L3 and L4 are signal delivery lenses.

schemes is the direction of the pump light relative to the signal light. The counter-propagating MOPA scheme is the most commonly used configuration in high-power amplifiers as it has higher gain when operating in saturation compared to a copropagating scheme.^{10,45,69} A counter-propagating MOPA is shown in Fig. 6.

Aside from the direction of pump propagation relative to the signal the counter-propagating MOPA is identical in construction to a copropagating MOPA. Copropagating MOPAs are usually only used in high-power lasers when one wants to be conservative and protect pump diodes from leaked high-power signals or when one wants improved signal to noise ratio for amplifying very small signals.^{10,45,69} Counter-propagating amplifiers have higher gain and efficiency when generating high powers while copropagating amplifiers, though lower in overall output, tend to have less ASE and better signal to noise ratios on the output direction. Explanations for this phenomenon lie in the fact that gain is not uniform along a fiber laser and hence a signal sees higher gain either earlier or later in its propagation (depending on pump direction) as discussed qualitatively and analytically in Refs. 10, 45, and 69. Amplifier gain and output power can be increased by setting up either a bidirectionally pumped amplifier or a double-passed amplifier. Unfolding the MOPA through the mirror shows that both cases are essentially identical (though bidirectional amplifiers allow more pump power due to double end pumping), and both cases offer the combination of benefits and issues associated with counter- and copropagating schemes.^{10,45,69}

Angle cleaving or polishing is a critical consideration in fiber amplifier systems due to the fact that fibers have such high gain that they can lase, even without mirrors, based solely on Fresnel reflection feedback. As a consequence one or both ends of any fiber intended for MOPA use must be angled to prevent oscillations. Even with this angling, oscillations may still occur at very high pumping powers. Techniques beyond angle cleaving for reducing these oscillations can become quite complex.⁷⁰

Often one amplifier stage for the low-power seed is insufficient to reach the desired output powers since high-power MOPAs require reasonably powerful seeds to fully saturate leading to efficient energy extraction. As a consequence these systems require many stages of preamplification before the final power stage is reached. Often between stages there will also be devices to act as pulse pickers or gates to reduce ASE between pulses. There will also be optical isolators to keep reflected pulse energy and ASE from moving backward in the system, stealing system gain or reaching significantly high powers to cause optical damage to amplifier components.

All-Fiber Monolithic Systems

For fiber lasers to reach their full potential in terms of stability and ease of alignment the laser resonator or amplifier chain must be completely within a fiber. In actuality, no fiber or laser system can be truly monolithic as some components including fiber-coupled diodes, isolators, modulators, and other components require some measure of free-space propagation. However, the nature of these components is that, though not truly monolithic, they have been packaged in such a way that they are extremely stable and are fiber coupled. Here a general overview of the basic components of monolithic fiber lasers and amplifiers is considered.

Monolithic Fiber Laser Resonators The fiber laser resonator in the “all-fiber” configuration consists of “fiberized versions” of the same components as a free-space fiber laser resonator. Such resonators consist of the pump source and a fiber Bragg grating (FBG) spliced to each end of a double clad gain medium. It should be noted that two FBGs are used. One is a high reflector and one is partially reflecting. The latter can be designed to have a very low reflectance to mimic the effects of Fresnel reflection from a fiber end facet. The grating is preferred to directly using the end facet reflection, to avoid degradation of the laser if the end facet becomes contaminated or to provide specific laser performance in terms of linewidth or extraction efficiency by optimizing output coupler parameters.

An alternative form of fiber laser resonator is the ring resonator, consisting of a loop of fiber spliced back onto itself through a so-called tap coupler or splitter. This coupler splits off some portion of the light as a useful beam and allows the rest of the light to remain in the cavity. Pump light is spliced into the ring by way of a TFB or other pump multiplexing device into the gain fiber portion of the ring. Rings may also contain sections of undoped fiber to complete loops, to provide nonlinear effects needed for mode locking, or to provide dispersion compensation. Optical isolators are used to allow unidirectional oscillation. “all-fiber” rings are not commonly used in high-power resonator systems due to the lack of tap couplers and splitters designed for operation with LMA fibers at high powers. On the other hand, the ring configuration is still a useful cavity for seed lasers that can be directly spliced into fiber MOPA systems.

Monolithic MOPA Configurations Monolithic fiber MOPA systems have potential in terms of their ability to extend the stability characteristics of small low-power lasers to higher powers. The fiber MOPA in Fig. 7 is simple, and can be “chained” together as needed to form multiple stages of amplification.

Pump light is injected by TFBs or similar devices in the optimum direction since pump direction influences amplifier performance. Figure 7 shows a three stage system to demonstrate the location where two different pumping directions might be considered (copropagation early on to minimize

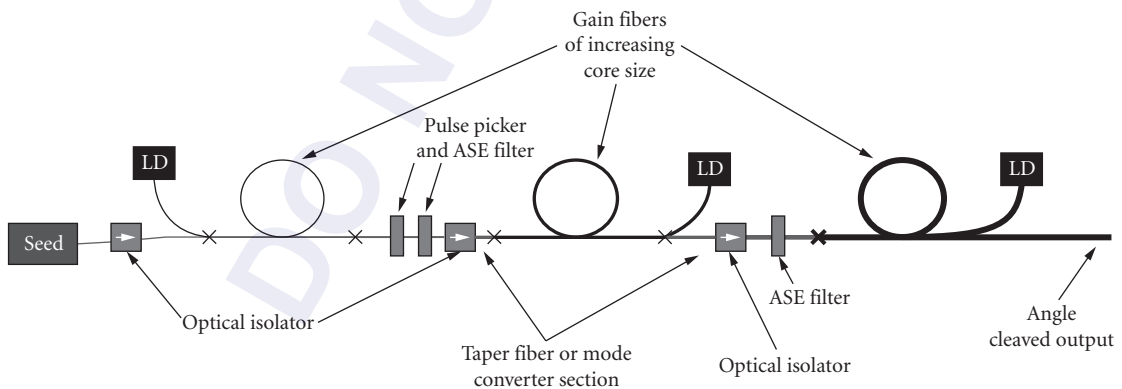


FIGURE 7 All-fiber-spliced MOPA system. LD are varying power pump laser diodes. The seed can be either another fiber laser or simply a fiber coupled laser diode. Sections are fusion spliced together with tapered sections of fiber providing the mode scaling between sections.

forward ASE and counter-propagation in the power stages to be sure of saturation). The only other elements that are included in the MOPA aside from undoped fiber between stages are the in-line isolators, ASE filters, and pulse picker for choosing repetition rate in pulsed systems.

All-fiber MOPAs work well up to reasonably high power levels, however, fibers for handling extreme power levels are not yet practically spliced into MOPAs as often the components such as isolators and pump combiners to deal with such large core sizes are not available.

Mode Field Adaptors The concept of mode matching and expansion can also be used to aid the design of MOPA-based systems. Usually the early stages of a MOPA system comprise small core telecommunications grade fibers while later stages will have larger cores. Sections can be connected together with fiber cores accurately aligned along their centerlines as seen in Fig. 8a. For large jumps in fiber core size a fiber section with a core diameter in between the two being matched may be spliced in as well, the so-called “bridge fiber method.” Often an appropriate bridge fiber is not readily available and thus other techniques have been developed, both similar in their general concept of gradually changing the mode field diameter to match two dissimilar fibers.

One technique involves simply tapering down a section of passive large core fiber matching the LMA section until its core dimensions match that of the smaller core fiber as seen in Fig. 8b.⁵⁶ The second method involves heating a fiber along its length so that the dopants in its core migrate to larger diameters via diffusion forming a mode matched core. This is known as the *thermally expanded core* (TEC) method⁷¹ (Fig. 8c). Choosing the appropriate fiber types by calculating mode field diameter of each using Eq. (7) and altering one or the other to make a reasonable match is the most straightforward way to make MFAs. In some cases, as seen in Eq. (7), the MFD of a fiber is not simply a linear function of the core diameter but rather has a minimum or maximum achievable value. Consequently, a small MFD cannot be increased sufficiently to match a larger MFD or a large MFD cannot be tapered down to a small MFD. A combination of both techniques (or even all three including a bridge fiber) must be

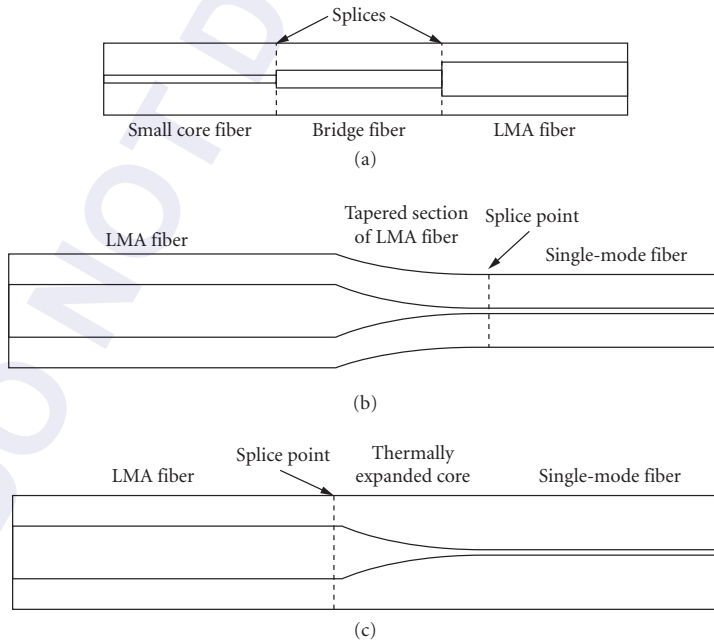


FIGURE 8 Methods for producing mode field adaptors: (a) bridge fiber method; (b) tapered fiber method; and (c) thermally expanded core (TEC) method.

employed to fabricate an appropriate mode field adaptor.⁵⁶ Often such MFAs are included as part of the construction of a TFB used to inject pump light to minimize the number of system components.⁵⁶

Fiber Bragg Gratings Fiber Bragg gratings are Bragg stacked layers of varying index of refraction written directly into the core of an optical fiber. Their operating principle is the same as any dielectric mirror utilizing a stack of quarter-wave layers as discussed in Chap. 17, “Fiber Bragg Gratings” and Ref. 72. Currently FBGs for large-mode area fibers are becoming available and several methods for their manufacture have been investigated including using photosensitive glass and writing holographically⁷² or by writing using direct inscription with femtosecond laser pulses.^{73–75}

Fusion Splicing The key advantage of an all-fiber-based laser system is the fact that the fibers are permanently bonded or “fused” together eliminating the potential for misalignment. Two fiber tips are microaligned in close proximity using a camera or other user-operated or automated feedback systems. The fiber ends are then heated rapidly with an arc discharge or electrical filament to melt the glass and fuse the two ends together. This method allows for minimal losses (< 0.1 dB) while giving high mechanical strength and stability.⁷⁶ Such a splicing technique was used many years ago for telecommunication fibers but only recently have these techniques been devised to handle LMA fibers requiring increased fusing temperatures and facet uniformity. New LMA splicers coming onto the market are one of the main reasons for the advance of LMA fiber technology as a whole. For a complete treatment of fusion splicing and fiber cleaving the reader is referred to Ref. 76.

25.6 LMA Fiber Designs

The desire for higher power or pulse energy fiber lasers requires the increase of core diameter in order to avoid damage and unwanted nonlinear effects. A main advantage of fiber lasers lies in their inherent beam quality stemming from single transverse mode operation. As a result methods for creating large mode area fibers while preserving beam quality are sought.

Figure 9 shows a schematic of different techniques for achieving LMA single-mode operation and will be referred to throughout this section.

Conventional LMA Techniques

Conventional LMA techniques utilize relatively standard fiber designs to enable high beam quality and only require low-core numerical apertures to keep the guidance of higher-order modes weak so that they can be stripped by introducing preferential losses. The weaker guidance causes higher bend loss for the fundamental mode and very low fiber NAs leads to fabrication challenges.^{19,77} Using the weak guidance concept, several techniques for LMA fibers have been developed.

Early LMA Fibers and Dopant Profiling The first designs for LMA fibers were used in Refs. 18, 77, and 78 to achieve high powers from Q-switched lasers. To improve beam quality in addition to using very small NAs the fiber core was given a tailored dopant and index profile to help provide preferential gain to the lowest-order mode.^{18,78} A theoretical study of tailored gain is given in Ref. 79. Using this technique the core diameter was extended to approximately 21 μm and output powers of up to 0.5 mJ at 500 kHz were achieved with M^2 of approximately 1.3.¹⁸

Coiling and Weak NA Guidance The technique of fiber coiling, first used in Ref. 80, takes advantage of the added bend losses in a weakly guided LMA fiber. Losses in higher-order modes increase with bend radius due to their weaker guidance in low NA fibers and by using tight fiber coiling, higher-order modes can be stripped out with minimal expense in power loss to the fundamental mode. Such fibers' index profiles resemble those in Fig. 9a, which is essentially a standard fiber design with very low NA (< 0.1). The theory of how bend loss in LMA fibers has been treated in Refs. 81 and 84, but in general optimal bending is determined experimentally.

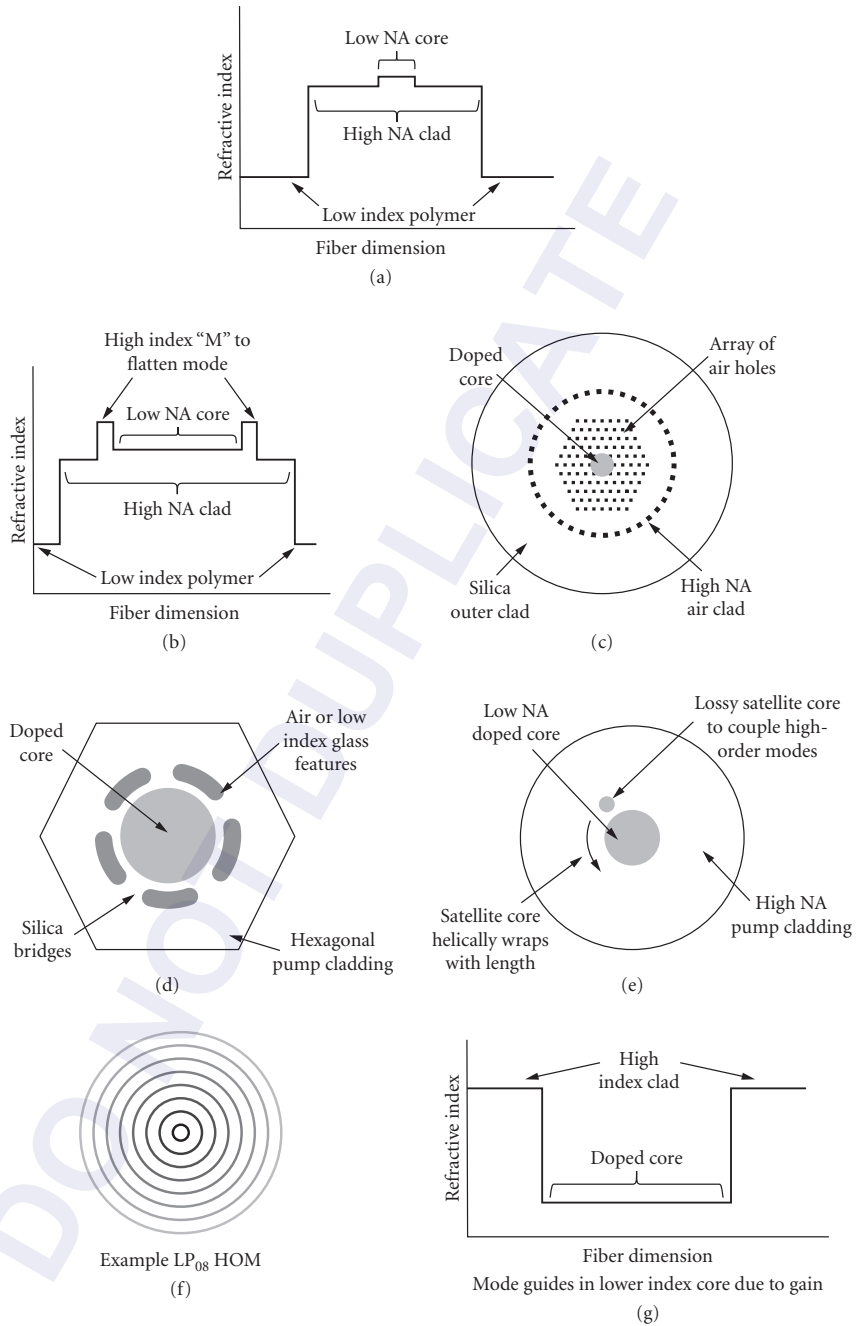


FIGURE 9 Sketches of various types of fiber designs discussed throughout this section: (a) conventional fiber design; (b) large flattened mode or "Batman" fiber; (c) typical photonic crystal fiber design; (d) leakage channel fibers; (e) chirally coupled core fiber; (f) example of the radiation pattern of a higher-order mode; and (g) index profile of a gain-guided index antiguided fiber.

Large Flat-Mode Fibers An extension of the use of specialized index profiles to enhance the performance of LMA type fibers involves the use of “M” shaped (or “Batman”) doping, whereby a slightly higher index step is placed on either side of the fiber core (Fig. 9b).^{85,86} This addition causes “flattening” of the mode field diameter and thus a larger but more flat topped fundamental mode is achieved. The mode area is increased (by up to a factor of 2 or 3) while maintaining near single-mode operation. This design described by Refs. 85 and 86 has been successfully used in several cases, mostly for achieving high peak power ultrashort pulses. By implementing these cores, the threshold for nonlinear interactions in the fiber core was increased by a factor of 2.5.⁸⁵

Single-Mode Excitation A technique that is related to conventional techniques is the single-mode excitation technique (SMET). This method, outlined by Refs. 87, 88, and 89, is most applicable to high-power fiber-based free-space MOPA systems. The mode field diameter of a small-core seed fiber is matched perfectly to that of the fundamental mode of a large-core fiber by using specially selected lenses and careful alignment. Using care, the fundamental mode can then be exclusively launched in a fiber and using this technique fibers with cores as large as 300 μm with mode field diameters (at $1/e^2$ of the peak on axis intensity) as large as 195 μm have been excited⁸⁹ (though not made into lasers) and fiber MOPAs with cores of upwards of 80 μm have been achieved.⁹⁰ SMET can be a very effective technique and has been used in many high peak power fiber lasers including in Refs. 90, 91, and 92.

Inclusion of Single-Mode Sections via Tapering Tapered fiber sections can be used to limit the propagation of higher-order modes in conventional type fibers. A tapered section fiber was demonstrated in Ref. 93, where a small section of multimode fiber was tapered to a diameter that allowed it to become single mode. When inserted into an oscillator (at the output end), the section improved the beam quality, reducing M^2 from 2.6 to 1.4; however, the slope efficiency was reduced from 85 to 67 percent.⁹³ Distributing the mode filtering over the length of the fiber (by such techniques as coiling) is more effective since loss is better distributed.^{93,94} This approach has allowed powers of up to 100 W with near perfect M^2 by using a 27- μm diameter core with an 834- μm clad diameter fiber and tapering it gradually down by factor of 4.8 over approximately 10 m of length.⁹⁵

Photonic Crystal Fibers

Limitations of precision in fabrication of conventional step-index fiber preforms leads to difficulties in obtaining fiber core numerical apertures less than approximately 0.06.^{19,96} In response, a class of fibers called photonic crystal fibers (PCF) has been developed.^{97,98} These fibers utilize an organized structure of refractive index differences (air holes or different high- or low-index glasses) to either tailor the core numerical aperture to sufficiently small values or to create a photonic bandgap where only certain wavelengths of light are allowed to propagate. Versions of these fibers have been developed for their beneficial nonlinear, dispersive, and loss properties in studies since the mid-1990s.^{97,99–101} The first fiber lasers based on PCFs were demonstrated in the early 2000s and had core sizes around 10 to 15 μm .^{102,103} Current state of the art in PCF technology allows core sizes of up to 100 μm ^{20,104} with powers of up to 4.3 mJ pulsed, corresponding to 4.5 MW peak power in approximately nanosecond pulses.

Air-Filled or Holey Fibers Air-filled PCFs are the most common and first demonstrated PCFs for fiber lasers.^{97,103} The air holes function to alter the NA giving it the desired value to maintain single-mode guidance. Figure 9c shows a sketch of a typical PCF design. The V parameter is still the limiting factor in the guidance of single mode beams in a PCF and its value is approximately calculated by

$$V_{\text{eff}} = \frac{2\pi}{\lambda} \Lambda F^{1/2} (n_0^2 - n_a^2)^{1/2} \quad (8)$$

where V_{eff} is the effective V parameter (still < 2.405 for single-mode propagation), Λ is the spacing of the air holes (pitch), F is the filling factor of air to glass, n_0 is the index of the glass, and n_a is the

index of the air (or whatever material fills the holes).¹⁰⁵ This is only an approximate model, giving rule of thumb results, and Ref. 106 covers PCF analysis in further detail. Thus control of the pitch and hole size of the PCF structure allow production of endlessly single mode structures.¹⁰⁷

As most PCFs manufactured for high-power laser purposes are double clad designs, they must also be designed with cladding pumping in mind. Since the largest core PCFs must not be bent sharply they must be able to absorb pump power in relatively short lengths.¹⁰⁸ To achieve this with reasonable doping levels in the core the core-to-cladding ratio must be kept as large as possible which is usually achieved by using an air cladding. The *air cladding* is simply a ring of very thin silica “bridges” or alternatively very large air holes with small connections to the outer glass fiber.¹⁰⁸ Using this technique, air claddings with NAs as large as 0.8 have been realized.¹⁰⁹ The one detrimental effect of the air cladding is that it does not allow for efficient thermal conduction of heat generated in the fiber core out to the cladding. An analysis of heat transport in PCFs is given in Ref. 110. Despite this, the highest average power from a rodlike air clad PCF reported to date is 320 W¹¹¹ and 1.53 kW from a longer, coilable PCF.¹¹² The highest pulse energy is 4.3 mJ.²⁰ A polarization maintaining PCF with 2300 μm^2 mode area was also demonstrated to have 161-W output power in a single polarization.¹¹³

All Solid-Core PCFs The main issue with the use of PCFs, especially in their adoption to all fiber systems lies in their air holes. It is challenging and complicated to form and fabricate a preform and fiber with air holes¹¹⁴ (though extrusion techniques may make this process simpler for instance in Refs. 115 and 116). It is also difficult to splice and cleave a fiber with air holes. There are further difficulties associated with (1) contamination presented by particles working their way into the air holes and (2) finding fiber-based components compatible with PCFs.¹⁰⁴ All-solid PCF fibers simply replace the air holes with a low or high index material to avoid the air-hole issues. In addition the solid defects can be engineered to provide a photonic bandgap which is a more complex way of guidance providing single-mode performance with the potential additional benefits of polarization maintenance and dispersion management. Novel uses for bandgap PCFs such as low-dispersion femtosecond lasers and 900-nm Nd lasers (using a bandgap to suppress 1064 nm operation) have been constructed^{117,118} and other lasers have also been proposed.¹¹⁹

Leakage Channel Fibers So-called leakage channel fibers^{120–122} sketched in Fig. 9d, employ low index of refraction regions surrounding a large core with “bridge regions” of the same index as core glass connecting core and cladding. The result is preferential leakage of the higher-order fiber modes providing more robust lowest-order single-mode operation. The early types of these fibers resembled PCFs as they used air-filled regions. Recently all solid designs have been demonstrated with core diameters as large as 170 μm .¹²⁰ Advantages of such fibers include the ability to be effectively cleaved, spliced, and bent. However, they are more challenging to manufacture than conventional fibers.

Chirally Coupled Core Fibers

A new type of LMA fiber laser is the chirally coupled core (CCC) fiber (Fig. 9e), which utilizes specially designed satellite cores helically wrapped around the central LMA gain core to couple out the higher-order modes into the lossy satellite while maintaining the lowest order mode in the core.¹²³ Designing the helical core or multiple cores with proper pitch and size leads to strong coupling of high-order modes to the intentionally lossy satellite.¹²³ The losses in individual modes can be calculated and optimized for a particular helix pitch and size. Several CCC fibers have been fabricated and tested in two common wavelength regimes, 1060 and 1550 nm exhibiting excellent beam qualities for inner core diameters as large as 50 μm . CCC fibers are attractive also because they do not require small core NAs, can handle strong bending, and can be used for very long lengths to achieve high CW powers with less thermal load. In fabrication CCC fibers require the preform core to be formed by standard modified chemical vapor deposition (MCVD), but a hole must be bored for insertion of the satellite rod. Then during drawing the preform must be spun at a set speed to achieve the chirality of the satellite at the desired pitch. An approximately 40-W fiber laser from a 33- μm core CCC fiber with very large V parameter but high beam quality has been demonstrated at 1064 nm.¹²⁴

Multicore Fibers Fiber lasers with multiple doped cores within one central pump cladding have received some attention in recent years due to their potential to scale mode area by allowing the addition of light from several individual cores. With proper design of the fiber, the fiber laser can have high beam quality in the far field.^{125,126} Multicore fibers are difficult to manufacture due to the need for drilling preforms and adding several cores.

Higher-Order Mode Fibers Another newly emerging class of LMA fiber designs is the higher-order mode fiber (HOM). These rely on using mode-conversion techniques usually based on long period fiber Bragg gratings.^{127,128} HOM beams have Bessel function spatial distribution.¹²⁹ The gratings are specifically designed so that they efficiently couple energy to only one HOM and are then used in tandem to subsequently de-excite the HOM back in to a more useful LP₀₁ mode at the fiber output. (though leaving the beam in a HOM may also be useful). Typically the LP_{0X} mode (where X is an integer >1) is excited. A sketch of such a mode is seen in Fig. 9f.

Modes as high as LP₀₈ have been launched and through this process mode areas of approximately 3200 μm^2 achieved.¹²⁸ A review of the theory is described in Ref. 127 in terms of their use in cladding modes of fibers, where this same theory applies to launching of HOMs into large-core multimode fibers as well. An experimental investigation into launching HOMs in custom large core fibers is reported in Ref. 128. Currently HOMs have only been used once in an actual laser resonator;¹³⁰ however, HOM-based modules have also been exploited for their ability to produce anomalous dispersion in silica fiber allowing generation of pulses as short as 60 fs in Yb fiber lasers.^{131,132} An investigation of HOMs ability to reduce SBS in fibers has been carried out in Ref. 133 where it was determined that SBS in HOM fibers can be reduced by using higher-order modes.

Gain-Guiding Index Antiguiding Gain-guiding index antiguiding (GG IAG) fibers are another new class of potential LMA technologies first proposed by Siegman in 2003.¹³⁴ These fibers consist of low refractive index cores surrounded by higher refractive index claddings, hence their so-called antiguiding nature. The core does not support conventionally index-guided modes. Instead it does support these modes in their “leaky” form. These the modes exist in the core, but constantly leak energy to the cladding if the fiber core is not excited. These modes can be thought of as existing due to Fresnel reflections at the core-cladding interface. Though these reflections can be low loss they can never be lossless as are total internal reflections. When the fiber is pumped gain in the core can compensate for the loss in the leaky modes resulting in lossless modes confined in the core. To maintain single-mode operation, gain must be supplied such that it makes up for the loss in the lowest-order mode but not for higher-order modes. The basic theory for GG IAG fibers is laid out in two papers by Siegman.^{134,135} Based on this principle laser oscillators can be designed with their gain and oscillation threshold conditions optimized to attain single-mode operation by using analysis given in Refs. 136 to 138. The first gain guided fiber lasers were demonstrated in flashlamp pumped cavities with core sizes from 100 to 400 μm ^{136,137} and M^2 less than 1.5. Diode pumping through the fiber end was subsequently demonstrated.¹³⁸ The mode areas reported by these papers are the largest reported in any fiber laser; however, efficiency and output power scaling of GG IAG fibers using new techniques and pumping schemes must still be addressed before the technology can become competitive with other LMA technologies. Gain-guiding effects have also been seen in conventional fibers. Recently a dramatic change in the guidance properties in highly doped optical fibers was reported at very high pump powers where the gain and refractive index changed due to very strong pumping.¹³⁹

25.7 ACTIVE FIBER DOPANTS

Apart from Raman fiber lasers which are not part of this review, all high-power fiber lasers utilize rare earth ion dopants. Though almost every rare earth ion has been doped into fibers for fiber laser applications high-power diode pump sources are not yet available for some. As seen in Table 1, almost every rare earth ion has been doped into fibers for fiber laser applications; however, high-power diode pump sources are not yet available for some dopants.

TABLE 1 Operating Range, Dopant Ion, and Transitions of Various Rare Earth Ions Fabricated in Fiber Lasers¹⁴⁰

Operating Range (nm)	Dopant (km)	Transition	Type of Host		Type of Transition
			Oxide	Fluoride	
∞ 455	Tm ³⁺	¹ D ₁ → ³ F ₄		Yes	UC, ST
∞ 480	Tm ³⁺	¹ G ₄ → ³ H ₅	Yes	Yes	UC, 3L
∞ 490	Pr ³⁺	³ P ₀ → ³ H ₄		Yes	UC, 3L
∞ 520	Pr ³⁺	³ P ₁ → ³ H ₅		Yes	UC, 4L
∞ 550	Ho ³⁺	³ S ₂ , ¹ F ₄ → ⁵ I ₁	No	Yes	UC, 3L
∞ 550	Er ³⁺	⁴ S _{1/2} → ⁴ I _{15/2}	No	Yes	UC, 3L
601–618	Pr ³⁺	³ P ₀ → ³ H ₆		Yes	UC, 4L
631–641	Pr ³⁺	³ P ₀ → ³ F ₂		Yes	UC, 4L
∞ 651	Sm ¹⁺	⁴ G _{1/2} → ⁶ H _{2/2}	Yes		4L
707–725	Pr ³⁺	³ P ₀ → ³ F		Yes	UC, 4L
∞ 753	Ho ³⁺	³ S ₁₂ , ³ F ₄ → ³ I ₁	No	Yes	UC, ST?
803–825	Tm ³⁺	³ H ₄ → ³ H ₆	No	Yes	3L
∞ 850	Er ³⁺	⁴ S _{5/1} → ⁴ I _{15/2}	No	Yes	4L
880–886	Pr ³⁺	³ P ₁ → ¹ G ₄		Yes	4L
902–916	Pr ³⁺	³ P ₁ → ¹ G ₄		Yes	4L
900–950	Nd ³⁺	⁴ F _{3/1} → ⁴ I _{4/2}	Yes		3L
970–1040	Yb ³⁺	⁵ F _{3/2} → ⁵ F _{7/2}	Yes		3L
980–1000	Er ³⁺	⁴ I _{11/12} → ⁴ I _{15/2}	No	Yes	3L
1000–1150	Nd ³⁺	⁴ F _{1/1} → ⁴ I _{11/1}	Yes	Yes	4L
1060–1110	Pr ³⁺	¹ D ₂ → ³ F ₄	Yes		4L
1260–1350	Pr ³⁺	¹ G ₄ → ¹ H ₅	No	Yes	4L
1320–1400	Nd ³⁺	⁴ F _{1/1} → ⁴ I _{1/2}	Yes	Yes	4L
∞ 1380	Ho ³⁺	³ H ₄ → ³ F ₄	?	Yes	4L
1460–1510	Tm ³⁺	³ D ₂ → ³ F ₄	No	Yes	ST
∞ 1510	Tm ³⁺	⁵ D ₂ → ¹ G ₄		Yes	UC, 4L
1500–100	Er ³⁺	⁴ I _{11/4} → ⁴ I _{15/2}	Yes	Yes	3L
∞ 1660	Er ³⁺	² H _{11/2} → ⁴ I _{0/2}	No	Yes	4L
∞ 1720	Er ³⁺	⁴ S _{1/2} → ⁴ I _{0/2}	No	Yes	4L
1700–2015	Tm ³⁺	³ F ₄ → ³ H ₆	Yes	Yes	3L
2040–2080	Ho ³⁺	³ I ₁ → ⁵ I ₈	Yes	Yes	3L
2250–2400	Tm ³⁺	³ H ₄ → ³ H ₅	No	Yes	4L
∞ 2700	Er ³⁺	⁴ I _{14/2} → ⁴ I _{13/2}	No	Yes	ST
∞ 2900	Ho ³⁺	⁵ I ₆ → ⁵ I ₇	No	Yes	ST

3L = three-level; 4L = four-level; UC = up-conversion; ST = apparent self-terminating.

Discussed here are only those dopants, which are now capable of high powers using efficient high power, high brightness AlGaAs, and InGaAs laser pump diodes. These dopant ions include neodymium, erbium, ytterbium, thulium, and holmium. In contrast to these dopants used in a crystalline matrix for bulk solid-state lasers, required for good thermal conductivity, in the glassy hosts used for fibers they have broad absorption and emission spectral linewidths.¹⁴⁰ The broad absorption bands relax the needs for strict control of diode pump wavelength while the broad-emission spectra allow for broad tunability at high powers and the generation of extremely short pulses.

Neodymium-Doped Fibers

Neodymium (Nd) is perhaps the most common laser dopant of bulk solid-state lasers in both crystalline and glass hosts, and it is not surprising that it was the dopant of choice for many early fiber

lasers.^{1,4,12} Its four-level excitation scheme permitted low lasing thresholds necessary in early fiber lasers when glass composition was poor and available pump powers were minimal. The two most common Nd transitions in fiber lasers are the 1.06- μm transition and the 1.3- μm transition, though Nd can also lase directly to the ground state on an approximately 900-nm transition, pump bands for Nd are in the 800-nm range.

The most prevalent transition used in Nd fiber lasers at 1.06 μm has achieved powers in the 300-W range and has been multiplexed into fiber systems to more than 1 kW.^{23,141} Although wavelength tuning over more than 60 nm with watt-level output powers³⁶ has also been shown, Nd is only efficient over 10 to 20 nm around the peak of the 1060-nm emission band.³⁶ As an amplifier Nd has been used at 1.06 and 1.3 μm . The latter wavelength has seen minimal use in high-power lasers due to its larger quantum defect, excited state absorption, and as a consequence of ASE at 900 and 1060 nm stealing gain from the 1.3- μm laser system.^{36,140} In recent years Nd has fallen out of favor for high-power fiber lasers since Yb³⁺ fiber lasers operate in about the same wavelength region and with more efficiency.

Erbium Fiber Lasers

Erbium is well known as the gain medium for most telecommunications amplifiers operating in the 1.5- μm range. It has also been used in high-power fiber laser systems producing output powers as high as 297 W (in Er:Yb codoped lasers).¹⁴² Erbium possesses several potential operating wavelengths, but with available high-power pump sources around 1400 and 980 nm (980 nm is most commonly used when Er is codoped with Yb), only the 1550-nm transitions generate high powers. At these pump wavelengths the absorption cross section is small making a double clad scheme impractical for direct pumping of erbium. To obtain high-power-efficient lasing, erbium is most often sensitized by codoping with Yb, allowing for much higher pump light absorption and, consequently, high power lasing. However, codoping leads to problems, as emission of 1.06- μm light from the Yb ions can take place at high pump powers.¹⁴² The development of high-power diode lasers at 1480 nm may allow Er-doped fibers to be directly pumped with minimal quantum defect to very high power levels and efficiencies akin to those achieved at 1.06 μm with Yb-doped fibers. High-power Er fiber lasers pumped in the 1480-nm region will have many applications because of their relative eye safety. Erbium has a large emission bandwidth with tunability demonstrated from 1533 to 1600 nm.³⁶ Erbium has also been doped in a fluoride host fiber (ZBLAN) giving rise to laser transitions in the 2.7- to 2.9- μm region (silica is not transparent here due to OH⁻ absorption), with powers of approximately 10 W. However, the low melting point of the fluoride glasses limits the potential for very high-power operation.^{143,144} In addition, tunability from 2.7 to 2.83 μm has been achieved with a power more than 2 W.¹⁴⁵ The large quantum defect from its 980-nm pump bands is another challenge though such fiber lasers are still one of the most direct ways to access the 2.7- to 2.9- μm range.

Ytterbium-Doped Fibers

The ytterbium (Yb) ion is by far the most commonly used ion in fiber lasers. Yb worked its way into the fiber laser mainstream for several reasons including its lower quantum defect, the ability to dope high levels of Yb into silica fibers with lower tendency toward concentration quenching, and no excited state absorption or upconversion at longer wavelengths.^{36,47} These benefits outweigh the three-level operation of Yb, which leads to higher laser thresholds than Nd.^{36,47} Yb has only two main energy levels, and is able to lase because these two levels are split in energy due the Stark effect.⁴⁷

The broad (~80 nm) absorption spectrum of Yb is peaked at 915 and 976 nm permitting wide flexibility in the choice of pump source wavelengths. With the low quantum defect there is substantial overlap between emission and absorption. Though the peak of the Yb emission is at 1030 nm, ground state absorption to the upper laser level causes longer fibers to “red-shift” their peak gain wavelength.³⁶ This length-dependent gain limits the tuning potential of Yb lasers. Longer fiber lengths are required to efficiently absorb pump power resulting in larger peak gain wavelength “red-shifts” and narrowing of

the tuning range.³⁶ Tunability has been demonstrated for over 150 nm of bandwidth³⁶ and the minimal quantum defect in Yb has made it the medium of choice for the operation in multi-kilowatts regime.²¹

Thulium-Doped Fiber Lasers

Thulium-doped fiber lasers are of interest because of (eye-safer) emission at approximately 2000 nm. This is desirable at high powers for defense applications and remote sensing. In addition, many applications call specifically for light in the 2- μm regime including difference frequency generation with 1 μm to create other IR wavelengths (3 to 5 μm), difference frequency generation of two closely spaced 2- μm beams to create terahertz radiation, and light detection and ranging (LIDAR).³⁶ Approximately 1.9- to 2- μm output also allows thulium (Tm) fiber lasers to be used to pump holmium fiber lasers to reach further into the IR.

Tm is a three-level laser system which terminates on the ground state leading to high laser threshold pump powers similar to those in both Er and Yb. In addition, termination at the ground state makes Tm lasers extremely temperature sensitive. Tm lasers are usually actively cooled for efficient operation. Tm has the potential advantage of having several available pump bands at wavelengths that can be reached by high-power sources. The 1200-nm absorption line was recently explored with direct diode pumping,¹⁴⁶ however, sufficiently high-power diodes at this wavelength are not yet available. Tm lasers can be pumped at 1060 nm, though with very weak absorption. Similar to Er, Tm can also be codoped with Yb;¹⁴⁷ however, this leads to similar 1- μm emission issues as in the Yb:Er codoped laser and also results in a huge reduction in potential efficiency. The two most commonly used pump wavelengths for Tm are approximately 1550 and 790 nm. The 1550-nm band is pumped by multiplexed Er Yb-codoped fiber lasers. Though this allows for very high power pumping this transition is relatively inefficient. It first requires the construction of a number of Er:Yb fiber lasers. Though pump to signal quantum defect is reasonably low in the Tm laser using Er:Yb laser pumping the overall efficiency suffers due to the Er:Yb lasers. Fiber laser pumping of thulium allows direct core pumping in situations where very short fiber lengths are required. Many commercial systems use cladding or core-pumped schemes with 1550-nm pumping.¹⁴⁸

Another promising scheme for Tm pumping is the use of 790-nm laser diodes. At first glance this process might seem extremely inefficient since there is a large quantum defect between signal and pump (maximum of ~40 percent efficiency). However, there is an additional process that can be exploited. This is called “cross-relaxation” or two-for-one pumping,¹⁴⁹ and allows the transformation of one pump photon into (theoretically) two laser photons. The result is a maximum efficiency of approximately 80 percent.¹⁴⁹ Experimental efficiencies of up to 68 percent (optical to optical power) have been reported¹⁵⁰ with 60 percent easily achieved at high powers.¹⁵¹ The challenge for cross-relaxation pump process is that it is dependent on temperature, dopant-concentration, and fiber composition.^{149,152}

Fiber laser emission from Tm is useful with an extremely large potential emission bandwidth stretching from 1700 nm to beyond 2100 nm.¹⁵³ This wide bandwidth gives Tm potential for application in both ultrashort pulse lasers and highly tunable mid-IR sources. Active tuning of Tm has been demonstrated over 230 nm.³⁶ The fiber length affects tuning range due to the same three-level reabsorption processes discussed for Yb lasers.

Tm fiber lasers have been demonstrated with single-mode CW powers of 268 and 415 W for 790- and 1550-nm pumping, respectively.^{148,151} In addition, an 885-W system with slightly multimode beam quality was also reported.¹⁵⁴ One additional benefit of Tm lasing is that the longer 2- μm wavelength allows larger core sizes with better beam qualities, lower nonlinearities, and higher damage thresholds since all these properties improve with longer wavelength.

Holmium-Doped Fibers

Ho-doped fibers are one of the only current option to achieve wavelengths longer than 2 μm . Doped in the proper fiber material (silica loses transparency beyond ~2.1 μm), holmium (Ho) fiber lasers have reached watt-level output powers.^{155–157} One of the most prominent absorption regions for Ho is around 1.9- to 2- μm region, where high-power pumping can be provided by Tm fiber lasers.

Tm pumping allows for extremely efficient conversion since the laser wavelength of Ho is approximately 2.1 μm and the quantum defect is minimal. This technique provides a way for stretching the Tm bandwidth to longer wavelengths in the near IR. Other Ho pump bands include the 1160-nm band which can be pumped directly by Yb fiber lasers or as demonstrated by Refs. 156 and 157 with direct laser diodes. In addition, Ho can be sensitized with ions such as Yb or Tm, taking advantage of their broad pump bands. Using this approach a Tm:Ho laser has achieved 83-W output power at 2.1 μm .^{158,159}

25.8 FIBER FABRICATION AND MATERIALS

High-power fiber lasers most commonly employ silica-based fibers because of its strength and thermal stability. However, other materials must be considered for fiber lasers because they offer different merits in terms of transparency, ability to be doped, laser parameters, and manufacturability.

Fiber Fabrication

Fiber fabrication can be divided into two main stages: fiber pulling and preform manufacture. The pulling stage is the actual “fiber making” stage and is common to any preform and fiber material. Preform manufacture is dependent on many fabrication methods. Further details on many of the steps of fiber fabrication can be found in Part 2, “Fabrication,” in Vol. II.

Fiber Pulling The process of fiber pulling is rather similar for all fiber materials with the exception of the operating temperature required, whether a polymer coating is deposited, and the final pulling diameter. This stage of manufacture involves taking a fiber preform and lowering it into an oven at the top of a fiber draw tower. The oven, if set to the proper temperature (varying from $\sim 800^\circ\text{C}$ for soft glasses to $\sim 2000^\circ\text{C}$ for silica) causes the preform to heat and eventually a globule of glass will drop down with a solid glass “string” attached.^{11,140,160} This “string” is the beginning of fiber itself and by wrapping this strand around a mandrel at a constant speed, the diameter of the fiber can be controlled precisely and long lengths of fiber can be formed in large spools.^{11,140,160} Coatings can also be applied during this stage and subsequently hardened (this is an important step for double clad fibers, as low-index polymers allow guidance of the pump light in the fiber cladding).

Preform Manufacture The heart of fiber manufacture and the current and future progress in high-power fiber lasers depend upon the development of “fiber pullable,” defect-free, low-loss preforms. It is one of the main challenges of fiber laser development today.^{11,140,160} The four main preform development techniques are summarized below.

Modified chemical vapor deposition Chemical vapor deposition (CVD) is widely used for the rapid fabrication of large preforms with very accurate index steps and compositions. It is implemented in several approaches, including traditional modified chemical vapor deposition method (MCVD), outside vapor deposition (OVD), and vapor-axial deposition (VAD).^{11,140,160} All three methods involve depositing a soot of chemical oxides onto some kind of rotating silica substrate mounted in a lathe by heating various gasses flowing over the substrate. The soot deposition builds up the desired core and clad layers and the preform tube is collapsed to form the actual preform. OVD applies the soot to the outside of a silica rod, VAD applies the soot to the end of a silica rod which acts as a “seed,” while the most common MCVD introduces the soot to the inside of a silica tube which is subsequently collapsed.^{11,140,160} CVD is extremely flexible in terms of the index profiles and dopants; however, there are challenges in making very large uniform cores for high-power fiber lasers.

Core drilling and preform machining MCVD is an excellent way to obtain radially symmetric doping profiles; however, often extra features that are not radially symmetric are desired such as in polarization-maintaining fibers or fibers with odd-shaped claddings to promote pump absorption.

Holes can be drilled into the preform and subsequently filled with glass or left to be air filled as in PCFs. The preform can also be given flat surfaces to provide the nonuniform pump claddings. In principle the core-drilling technique is quite simple, however, the difficulty with this method is twofold; first, obtaining pullable glass for this type of preform can be challenging (as usually this method is used for nonsilica glasses) and second, drilling holes in glass without breaking it is not a simple task; it takes a good deal of time, care, and skill to make such a preform not to mention the expense.

The benefits of core drilling lie in its flexibility to be feasible with any glass (especially, soft glasses for which MCVD is not available). It also allows the core to be inserted with precise doping characteristics since it is prepared as bulk glass separately with precise control over its composition. In addition, core drilling allows the use of smaller samples of glass which is critical for many types of glasses that are expensive or difficult to make in large quantities. The core index profile is also very uniform since it is one solid piece of glass, hence fibers with very large cores and high dopant concentrations can be more effectively manufactured.

Stack and draw This method is commonly used to manufacture photonic crystal fibers, which require complex structures that would be too expensive or risky to manufacture using hole-drilling techniques.¹⁶¹ The basic premise is to use assorted rods (or cane) and tubes of the desired glass material and stack them in the desired pattern inside a larger glass tube.¹⁶¹ A core region can be added by way of introducing a doped core rod into the pattern of rods or tubes.

This method is useful for PCF manufacture, though, as with the core drilling approach, it depends on finding fiber-pullable glass components that are the appropriate size, particularly, in glasses other than silica. Furthermore, fusing of the rods and tubes together into a solid preform without collapsing the tubes is not trivial and requires much practice.

Extrusion In the extrusion technique glass is pressed through a die to obtain a preform with the desired air holes or glass dopants.¹¹⁶ The technique has a great potential in the realm of PCFs, as arbitrary shapes can be readily generated with a suitable die. A doped core region can be included by simply using two types of glass in an extrusion (not dissimilar to the way different colors are obtained in one tube of toothpaste).¹¹⁶ To date this method has only been used with so-called “soft glasses” having low melting temperatures.

Fiber Materials Though silica fiber is the most dominant material for high-power fiber lasers, several other materials find niches for specific applications. Varying the doped host material may provide benefits in one of three important areas, its laser properties such as upper-state lifetime and emission cross section, its doping properties and potential for higher doping, and finally transparency considerations when operating at mid-IR wavelengths. These different materials and their potential benefits will be discussed in subsequent sections. Table 2 which contains glass data for several glass types will be referred to throughout this section.

TABLE 2 Sample Properties of Different Types of Glasses Used for Fibers^{162–172}

Glass Type	T_x (°C)	Bulk Damage (GW/cm ²)	Thermal Conductivity (W/mK)	dn/dt (10 ⁻⁵ /°C)	CTE (10 ⁻⁷)	Trans. (μm)	Young's Modulus (GPa)	Knoop Hardness	n_2 (10 ⁻²⁰ m ² /W)
Silica	1175	600	1.3	11.9	0.55	0.3–2.1	72	600	3.4
Phosphate	366	25	0.84	-4.5	104	0.4–2	71.23	418	1.2
Germanate	741	—	0.55	1.2	63.4	0.5–3.9	85.77	560	—
Tellurite	482	10	1.25	—	—	0.5–5	54.5	—	30
Chalchogenide	180	6	0.37	9.3	21.4	0.6–8	15.9	109	400
ZBLAN	385	0.025	0.628	-14.75	17.2	0.5–5	52.7	225	2

Note that data even among one glass type is scattered among different compositions, hence these values may only be taken as approximate for comparison sake.

Silica fiber Nearly all fibers used in telecommunications are based on silica and as a result the manufacture, splicing, cleaving, and polishing of this fiber has been optimized. As a result, adapting its use for high-power fiber lasers has been a natural transition. Silica is the only fiber commonly manufactured by the MCVD technique which is the fastest, simplest, and cheapest way to manufacture fiber preforms.¹¹ Silica's high damage threshold and melting point of approximately 2000°C make it especially suitable for high-power operation. A final benefit of silica is its physical durability to mechanical and thermal stresses. Silica is able to keep cleaved or polished surface well, is stable under vibration and strong enough to coil tightly, making silica fibers suitable for environmentally taxing packaged fiber laser applications.

As useful as it is in many applications silica does also have some cons including a positive refractive index change with temperature, a slightly higher n_2 than some glasses, a limited transparency window in the mid-IR due to its high phonon energy of 1100 cm^{-1} , as seen in Table 2.¹⁴⁰ Sometimes other host glasses provide superior laser parameters compared to silica.^{173,174} A final, and perhaps most important, limitation on silica is the relatively low dopant concentrations it can handle (a few wt. % before the onset of clustering and other detrimental effects).^{140,152} This maximum doping threshold makes it difficult to form highly doped fibers that may be advantageous for dopant concentration-dependent processes such as up conversion and cross relaxation. Highly doped short fibers, not practical in silica, have applications in ultrashort pulse amplification and single-frequency generation.

Phosphates, germanates, and tellurites Glasses such as phosphates with open glass structures are capable of handling far higher dopant concentrations compared to silica.^{140,152} Fiber lasers have been demonstrated with doping percentages as high as 10 times that allowable in silica^{175,176} enabling very short fiber lengths so that these fibers are suitable for high peak power amplification, narrow line-width generation, and for use in core designs that limit the length of the fiber due to bending losses.

Fiber end melting in end-pumped configurations is a significant challenge to these types of fibers as pump powers in the range of 20 to 100 W can cause catastrophic melting. However, using other more evenly distributed pump schemes heat has been more uniformly distributed and higher powers have been achieved.⁶⁴ Phosphate fiber lasers have achieved as high as 20-W output power with Yb doping¹⁷⁶ and 4-kW peak power in single frequency Q-switched systems,¹⁷⁷ germanate fiber lasers with pulse energies of 0.25 mJ and output powers of 104 W have also been reported.^{150,178}

Some of the soft glasses possess better laser characteristics for some dopant ions in terms of cross section and lifetime (there are many studies for different dopants and fiber types, however, see for example,¹⁷⁴ for tellurite fibers). They are unfortunately more difficult to manufacture in terms of cost of materials and required time for core drilling. Often, though not always, their difficulty in manufacture and limited power-handling outweighs their superior laser, thermal, and doping properties compared to silica.

Other Mid-IR Glasses: Fluorides and Others Though the bulk of current interest in high-power fiber lasers is concentrated in the near IR (1 to 2 μm), some applications require high powers outside of this relatively narrow band.

The most common glass material for producing high-power fiber lasers outside of the traditional wavelength band is the fluoride glass family. The most widespread of these glasses is so called ZBLAN, named for its chemical composition containing ZnF_2 , BaF_2 , LaF_3 , AlF_3 , and NaF .¹⁴⁰ Compared to silica, ZBLAN allows more laser wavelengths in both the visible and farther into the IR.¹⁴⁰ Several reasonably high-power lasers have been reported using ZBLAN doped with Ho with outputs of 0.38 W,^{155,156} Er with outputs of approximately 10 W at 2700 nm,¹⁴⁴ and in Tm with outputs of 20-W CW and 9-W average power pulsed operation with pulse energies of 90 μJ .^{173,179,180} ZBLAN has potential in particular laser applications, where its lifetime and cross-section properties make it advantageous in terms of its having a lower threshold behavior. The characteristics of Tm:ZBLAN and Tm:silica at approximately 50-W pumping levels have been compared and it is found that ZBLAN is superior in terms of efficiency and threshold.¹⁷³ Despite the clear benefits of ZBLAN (and other fluorides) in some situations, there are also limitations stemming from fabrication difficulties, low melting point, and damage threshold making it difficult to produce fiber lasers with this material above the 50-W level.¹⁷³ This precludes it from generating the extreme powers of silica-based fiber lasers and, as a result, keeps the mid-IR wavelengths it is able to produce limited to the sub 100-W level.

There are other potential glasses for mid-IR operation including chalcogenides; however, these glasses share the difficulties of ZBLAN in terms of damage threshold and cost of manufacture. Chalcogenides are difficult and expensive to make in quantities large enough to fabricate a fiber preform. Lasers based on these materials (a Raman fiber laser) have achieved power levels of 0.64 W.¹⁸¹

25.9 SPECTRAL AND TEMPORAL MODALITIES

High output powers from fiber lasers have reached the multi-kilowatts level with diffraction limited beam qualities.²¹ However, many applications demand narrow spectral bandwidths or tunable spectral bandwidths for spectral beam combining and spectroscopy applications. Other applications call for short pulse durations (nanosecond, picoseconds, and femtosecond) with high peak powers at high repetition rate such as LIDAR and materials processing.

High-Power Spectrally Controlled Fiber Lasers

Spectral control is one of the most critical aspects of high-power fiber laser design. Fiber lasers constructed without spectral control tend to display chaotic spectral behavior with lasing occurring in multiple regions of the gain spectrum simultaneously as observed in Refs. 173, 182, and 183. This wide spectral variation and indeterminacy is unsuitable for many applications including pumping of other lasers and for spectral beam combining where significantly higher spectral brightness is desired. Some situations also call for active wavelength selectivity.

Fiber Laser Spectral Tunability Fiber laser spectral tuning is most easily accomplished by the inclusion of a dispersive element in the resonator such as a diffraction grating, prism, or volume or fiber Bragg grating. Free-space fiber laser cavities are easily configured to be wavelength tuned by simply replacing an end mirror with the tunable optical element such as a conventional diffraction grating in the Littrow configuration. Only a small amount of feedback is required to efficiently control the wavelength in fiber lasers because of its high gain, spectral narrowing of laser linewidth, due to angular dispersion of the spectrum in space is caused by the spectrally selective elements, though in many cases at the cost of efficiency.³⁶

Yb-doped fiber lasers with approximately 50-W output and more than 50 nm of tunability were reported, as were Er:Yb-doped fiber lasers with more than 100-W output and more than 40 nm of tunability and Tm-doped fiber lasers with more than 10 W output and more than 200 nm of tunability. These were in tunable oscillator configurations and showed relatively constant output powers over the tuning range.^{36,184,185}

Spectral tuning can also be achieved via volume Bragg gratings (VBGs). These are diffractive holographic grating structures that give very narrow band feedback with higher efficiency than metal or ruled gratings.¹⁸⁶ A Yb fiber laser with 4.3-W output power and a 30-nm tuning range was demonstrated.¹⁸⁷ VBGs do not work in the Littrow configuration, and thus require an integrated feedback mirror. However, they offer the benefit of higher potential efficiency due to lower losses compared to traditional gratings.¹⁸⁷

Fiber lasers can also be tuned with fiber Bragg gratings. Because FBGs are simply layers of differing photoinduced refractive index change arranged in a fiber, changing the distance between these layers by mechanical stretching or thermal tuning can change the reflectivity of the grating. This has been done in many systems. Reference 195 provides an example of 30-nm tunability and output of 43 W.¹⁸⁸ The tuning range was limited by the amount of mechanical or thermal change a grating can tolerate.

Narrow Linewidth Fiber Lasers

Narrow wavelength and wavelength control can be achieved in fiber lasers either by using narrow linewidth spectral control elements in high-power laser cavities or by seeding a high-power amplifier chain with a separate spectrally controlled light source.

Fiber Bragg gratings are a common line narrowing mechanism which can be fabricated in photo-sensitive glass or via direct femtosecond writing of FBGs into fiber laser glass. The later type of grating has the advantage of being written directly in the gain fibers with no special glass needed. This type of FBG enabled 104-W output power and 260-pm linewidth.^{74,189} As seen in Chap. 17 by choosing correct design parameters, a FBG can be tailored to the desired wavelength, bandwidth, and reflectivity with linewidths as narrow as 0.01 nm achieved in single-mode fibers. FBGs have been proven able to handle very high powers, as they were incorporated into several of the highest power lasers reported.²¹

FBGs are not easily compatible with large-mode area fibers because FBGs in large core sizes cannot easily be made to the tight tolerances demanded for laser applications. The use of volume Bragg gratings (VBGs) and guided mode resonance filters (GMRFs) do not suffer this limitation. VBGs can be designed to be nearly 100 percent reflective at normal incidence in extremely narrow wavelength ranges. Their fabrication and design is detailed in Ref. 186. The first use of VBGs in fiber lasers involved low powers in large-core PCFs.¹⁹⁰ This was extended in Yb-doped fiber reaching output power of 4.3 W and tunable linewidth of 5 GHz.¹⁸⁷ A 103-W Er:Yb laser was also demonstrated and was tuned using the VBG over approximately 30 nm at an output power of approximately 30 W.¹⁹¹ VBGs in Tm fiber lasers have also been demonstrated, exhibiting powers of up to 5 W and linewidths as small as 300 pm.¹⁸³ The highest power VBG fiber laser demonstrated was linearly polarized and reached 138-W from Yb-doped fiber. Thermal limitations to the use of VBGs is discussed in Ref. 192.

Guided mode resonant filters (GMRFs) are based on writing a layer of subwavelength gratings on top of a waveguide layer in order to use waveguide coupling effects to cause very narrow band reflectivity. The details of their operation are given in Ref. 193. These elements have not been used at high powers though they have been used to cause significant linewidth narrowing of a watt-level fiber laser.¹⁹⁴

Many applications call for linewidths less than the picometer range. Since fiber lasers usually have long cavity lengths, with closely spaced longitudinal modes, spectral control with a single dispersive element to a single mode is challenging. Short fiber laser resonators using highly doped soft glasses at the watt level have been demonstrated,^{177,195} but are not scalable to higher power.

The most effective way to achieve high powers and extremely narrow linewidths is to use a MOPA seeded by either narrow linewidth diode lasers or distributed feedback fiber lasers.¹⁴⁰ Despite their low power, DFB lasers are ideal seeds as they offer minimal temperature sensitivity with low noise, 1 to 100 kHz linewidth and single polarization. High power, single frequency lasers have been built using all the major gain media. A Yb-doped single frequency laser reached 264 W at less than 60 kHz linewidth.¹⁹⁶ The highest power Tm MOPA reported is 20 W at less than 50 kHz linewidth based on a DFB and only limited by available pump power.¹⁹⁷ High-power Er:Yb MOPAs have also been reported reaching 151 W.¹⁹⁸ The onset of SBS which causes linewidth broadening at high powers is often a limiting factor in these systems. To mitigate this, limitation fiber cores must be made larger, fiber lengths must be made shorter or special techniques must be used to eliminate SBS.^{42,199}

Nanosecond Fiber Systems

High pulse energies in fiber lasers are limited by the optical damage threshold of the fiber. Fiber lasers, so far, are limited to nanosecond pulses with energies of a few millijoules. Even in 100- μm core fibers such nanosecond pulses approach the optical damage threshold of silica glass.¹⁷ However, fiber lasers have the advantage of being able to operate at very high repetition rates (100s of kHz) so they complement bulk lasers in the high pulsed power regime. Nanosecond laser architectures are either Q-switched or use low-power seed pulses.

Q-Switched Oscillators Conventional Q-switched fiber lasers use a light modulator (passive, electro-optical (EO) or acousto-optical (AO)), adjacent to an angle cleaved fiber facet to avoid parasitic lasing between pulses, and a feedback element in the resonator. The main challenge with Q-switched fiber lasers is maintaining hold-off between pulses as ASE can build up to the detriment of laser efficiency.^{140,200}

Higher peak powers can be reached using large core fibers, end capped with coreless caps to prevent surface damage by expanding the beam before it exits the fiber. Several high-power Q-switched

oscillator systems have been reported with millijoule-level output powers based on both conventional LMA and PCF technologies. An LMA-based laser doped with Yb produced 8.4 mJ at 500 Hz and 0.6 mJ at 200 kHz with 120 W of average power at the higher repetition rate. The beam was slightly multi-mode giving M^2 of approximately 4.²⁰¹ A PCF-based Yb-doped fiber laser produced 10 ns pulses with energies up to 2 mJ and an average power of approximately 100 W.^{202,203} Tm-doped fiber lasers have also been Q-switched producing 30-W output power with 270- μ J pulse energies at 125-kHz repetition rates in conventional LMA fiber.²⁰⁴

Several potential monolithic Q-switching solutions have been tested, using both active and passive switching. Passive Q-switching usually involves using some kind of saturable absorber in the cavity either as an end mirror (making the system essentially monolithic), bulk crystal, or by splicing a section of saturable absorber into the fiber.¹⁴⁰ Saturable absorber mirrors and bulk saturable absorbers are limited in output power due to damage concerns in the saturable absorber elements. Nevertheless, saturable absorber Q-switched fiber lasers have achieved watt-level powers and 100- μ J pulse energies.^{205–207} Doped fiber with absorption at the desired operation wavelength has also been used as a saturable absorber. For example, a 10-W average power, approximately 100-kHz Tm laser with microsecond pulses was Q-switched with a Ho-doped fiber section.²⁰⁸ Other alternative Q-switching methods involve using mismatched FBGs, where a resonator is formed between two FBGs and the gratings are altered in length piezoelectrically to change their reflectivity peak and modulate cavity Q.^{140,209} Many other novel methods for Q-switching fiber lasers have been proposed including passive self Q-switching using SBS and using a piezoelectrically modulated high-Q microsphere or electro-optic-based metal-filled FBG.^{210–212} None of these techniques have been tested at high powers.

There is still a general challenge to attaining very short sub-20 nanosecond pulse durations in fiber lasers due to their long length.^{140,213} Typical fiber laser pulse durations are 100s of nanoseconds, while sub-100 ns are achievable with care. In very short PCF or highly doped soft glass lasers sub-10 ns pulses are achievable. A further challenge for Q-switched oscillators is the long intracavity length in a fiber laser causing unwanted nonlinear effects and even undesirable mode locking which affect the pulse shape and energy.^{140,214,215}

Nanosecond MOPA Systems When the most consistent and shortest possible pulse durations with highest achievable peak pulse powers are desired fiber MOPA systems seeded by Q-switched fiber lasers, microchip lasers or modulated laser diodes are a common solution. Microchip lasers possess small cavity sizes and can achieve very short pulse durations with reasonable peak output powers. They operate with passive Q-switching at fixed repetition rates anywhere from 3 to 100 kHz and produce seed energies more than 5 μ J.²⁰⁰ Diode laser seeds have the advantage of being flexible in terms of pulse shape and repetition rate when driven by arbitrarily shaped current waveforms. Modifications to the pulse shape in an amplifier can be calibrated out by tailoring the input pulse shape.^{90–92,216} Direct diode laser seeds require further stages of preamplification to achieve desired seed powers to saturate a power amplifier.

The multistage nature of the MOPA system configuration allows ASE between pulses to be filtered out by narrowband spectral filters or AO gates. MOPA systems currently able to achieve the highest peak powers usually rely on either conventional single-mode excitation in LMA fibers^{90–92,216} or PCF technologies^{20,108,200} to achieve high-beam qualities and large-mode areas.

Some of the downsides to MOPA systems are their complexity, the extra components they require and their longer time for assembly and optimization. MOPAs usually require at least one preamplifier (and in the case of diode systems two or more) to boost the seed power to a point where it can efficiently extract gain from a power amplifier. In addition, MOPAs require high-power optical isolators to protect earlier stages from back-reflected pulses as well as filters to remove ASE and mode field adapters to transfer signals from small-core to large-core stages.

End caps, sections of undoped, coreless fiber (solid silica) spliced to the end of a conventional fiber or sections of PCF with the air holes thermally collapsed allow the fiber mode to expand to much larger sizes before exiting the end face where fiber damage thresholds are far lower than in the bulk.^{20,200}

A Yb-doped system was reported which produced 6.2 mJ in sub-10-ns shaped pulses at approximately 2-kHz repetition rate using conventional LMA fiber with core diameter of 80 μ m. It had M^2 of approximately 1.3.^{92,217} Another similar system based on 200- μ m LMA fiber produced 27 mJ (82 mJ)

in 50 ns pulses, but with M^2 of 6.5 (though this is significantly better beam quality than expected from such a fiber due to the use of coiling techniques). Numerous PCF-based MOPA systems have been constructed with core sizes varying from 40 to 100 μm achieving upwards of 4.4 mJ of pulse energy at 10 kHz in 100- μm core with near-diffraction-limited beam quality.^{20,218,219} At these power levels, the pulse itself began to break up due to nonlinear effects in even such a large fiber core. Additional interest in high-power eyesafe MOPA systems has led to work on Er:Yb systems capable of more than 300- μJ pulses at 6 kHz and (or 100 μJ) at 100 kHz.^{220–222} Tm-based systems in fluoride fibers have reached 5-kW peak power in 30-ns pulses with 33-kHz repetition rates and 1-kW peak power with 125-kHz repetition rates.¹⁷⁹

High-Power Ultrashort Pulse Technologies

Ultrashort pulses (USP) and their applications are an area of increasing interest in the laser community and because of their large gain bandwidths, potential compact size and inherent stability, fiber lasers are an ideal platform for the generation and amplification of high-power USPs for use in frequency conversion, material processing, remote sensing, high harmonic generation, and production of high-power stable frequency combs. As with nanosecond pulses, fiber lasers have fundamental limitations on pulse output energy caused by nonlinearities and damage thresholds within small fiber cores at energies over a few millijoules, but fiber lasers have the ability to provide these pulses at very high average powers and repetition rates. By using pulse stretching and compressing techniques, fiber lasers are capable of producing output energies in ultrashort pulses on the same order of magnitude of nanosecond pulses. Dispersion-management techniques in fibers allow them to readily achieve pulse durations on the order of many classes of bulk ultrashort lasers. System architectures for fiber-laser-based ultrashort systems take the same two basic forms as their nanosecond counterparts: direct generation of pulses by high-power oscillator systems or amplification based on chirped pulse amplification (CPA) MOPA systems.

High-Power USP Oscillators Ultrashort pulses (USPs) have been produced in fiber lasers for many years dating back to the early interest in the generation of pulses for communications applications. Most USP systems in fibers are capable of only modest output powers of less than 100 mW and pulse durations of picoseconds.¹⁴⁰ The earliest lasers were based on temporal solitons where pulse duration was constant throughout the resonator and the high peak pulse powers of even short pulses limited potential output power. Stretched pulse additive pulse mode locking techniques allowed pulse duration to be shortened and pulse energies to increase somewhat; however, they are still relatively low compared to LMA fiber laser standards.^{140,223–225} The advent of stretched pulse, self similar, or all normal dispersion fibers where the pulse is compressed external to the resonator allowed an increase of output energies to nearly 20 nJ. This is a limitation in such lasers due to their small-core conventional fiber rather than LMA construction.^{226–229} High-power LMA fiber lasers based on all-normal dispersion techniques have been produced and have achieved record output powers.^{68,230–233}

High-power mode-locked fiber laser oscillators rely on using external cavities and very large-mode area fibers to achieve their output powers. Cavities usually take the form of ring or sigma resonators containing the gain medium, dispersion compensation, and required polarization control elements. Most of these systems use external dispersion compensation such as chirped mirrors and gratings or use none at all and rely on extra-cavity pulse compression.

Detailed descriptions of many types of mode locking used in fiber lasers are found in Refs. 140, 213, 234, and 235. The most common technique is use of saturable absorbers such as SESAMs (semiconductor saturable absorber mirrors) or carbon nanotubes; SESAM operation is described in Ref. 236, and saturable absorbers based on the use of carbon nanotubes are described in Refs. 237 to 240.

Mode-locked fiber lasers will also require dispersion compensation to compress pulses to their minimum duration. Many fiber lasers rely on traditional bulk optics making them less useful as stable fiber-based systems. Other dispersion compensation alternatives may involve the use of PCF or HOM fibers which allow for dispersion correction with very large-mode areas, and thus all-fiber LMA oscillators.^{131,241}

Despite the high-power achievements in LMA-based mode-locked oscillators, there are difficulties associated with such systems. Oscillator systems run at very high repetition rates in the megahertz regime due to the cavity length leading to low pulse energies even for 100-W-level systems. In addition, pulse durations are more difficult to control and stability can be an issue when operating at high average powers.

Ultrashort MOPA Systems There are very few differences between the construction of MOPAs for ultrashort pulse operation and nanosecond operation. The systems' architectures are quite similar; the differences lie in the way the systems are seeded. USPs have very high peak power. To amplify USPs in fiber lasers the chirped pulse amplification (CPA) technique first demonstrated in fibers in Ref. 242 must be used. Pulses from a low-power mode-locked seed laser are amplified by one or more preamplifier stages. The pulse is then stretched in duration by giving it a linear chirp using the dispersive effects of bulk grating stretchers, chirped mirrors, prisms, or even simply lengths of fiber (including potentially PCF or HOM fiber). This stretched pulse is next injected into an amplifier system in the same way as a nanosecond pulse. The bandwidth of fiber amplifiers is suitably large to handle the wide bandwidth of even sub-100 fs pulses. After amplification, the pulse is recompressed to its shortest possible duration. An added advantage of the MOPA system is that optical modulators can be incorporated after the seed laser to act as pulse pickers to reduce the megahertz repetition rates to kilohertz level rates more suitable to high pulse energy amplification.

With the use of these various types of seed lasers Yb systems have reached 100- μ J pulse energies at 90-W average powers with pulses as short as 500 fs, thus leading to 120-MW peak powers.²⁴³ Higher average power megahertz repetition rate systems have reached 131 W of average power.²⁴⁴ An even higher power system using two large core PCF amplifiers produced 1.45 mJ at 100-kHz repetition rate with more than 100-W output power in approximately 800 fs compressed pulses.²⁴⁵ Though PCFs have produced the highest output powers, LMA conventional fibers are also very capable and have produced high powers. The highest reported outputs from LMA USP lasers being 50- μ J pulses from 65- μ m core fibers which are used for x-ray generation.^{246,247} Many Er:Yb-based systems have also been constructed, the largest of which manage to produce upward of 200 μ J at 5-kHz repetition rates.²⁴⁶ New VBG-based compression and stretching techniques in Er:Yb and Yb fiber lasers have also allowed an increase in efficiency in such CPA systems.^{248,249} In addition CPA systems have been commercialized and are available for use in materials-processing applications.²⁵⁰

USP, CPAMOPAs suffer similar issues as other MOPA systems including increased parts count and complexity, the need for mode field diameter adaptation from stage to stage, and the lack of all-fiber components suitable for high-power levels, therefore necessitating free-space operation. Still, MOPAs are the most effective way to generate high peak energy and peak-power ultrashort pulses from fiber lasers.

25.10 CONCLUSIONS

Based on the discussions, data, and results presented here, fiber lasers are an effective technology for the production of high-power laser light. Despite many excellent results in many regimes, there is still a need for further development of fiber laser technologies to make the leap further in the realm currently dominated by bulk lasers. The many techniques, technologies, concepts, and systems discussed are the groundwork for enabling the advancement of fiber lasers in the near and far future.

25.11 REFERENCES

1. E. Snitzer, "Optical Maser Action of Nd³⁺ in a Barium Crown Glass," *Phys. Rev. Lett.* 7(12):444–446, 1961.
2. E. Snitzer, "Proposed Fiber Cavities for Optical Masers," *J. Appl. Phys.* 32:36, 1961.
3. C. J. Koester and E. Snitzer, "Amplification in a Fiber Laser," *Appl. Opt.* 3(10):1182–1186, 1964.

4. J. Stone and C. A. Burrus, "Neodymium-Doped Silica Lasers in End-Pumped Fiber Geometry," *Appl. Phys. Lett.* **23**:388, 1973.
5. J. Stone and C. Burrus, "Neodymium-Doped Fiber Lasers: Room Temperature CW Operation with an Injection Laser Pump," *IEEE J. Quant. Electron.* **10**(9):794–794, 1974.
6. R. J. Mears, L. Reekie, S. B. Poole, D. N. Payne, "Neodymium-Doped Silica Single-Mode Fibre Lasers," *Electron. Lett.* **21**:738, 1985.
7. R. J. Mears, L. Reekie, S. B. Poole, and D. N. Payne, "Low-Threshold Tunable CW and Q-Switched Fibre Laser Operating at 1.55 μm ," *Electron. Lett.* **22**:159, 1986.
8. L. Reekie, R. Mears, S. Poole, and D. N. Payne, "Tunable Single-Mode Fiber Lasers," *J. Lightwave Technol.* **4**(7):956–960, 1986.
9. E. Desurvire, J. R. Simpson, and P. C. Becker, "High-Gain Erbium-Doped Traveling-Wave Fiber Amplifier," *Opt. Lett.* **12**(11):888–890, 1987.
10. R. J. Mears, L. Reekie, I. M. Jauncey, and D. N. Payne, "Low-Noise Erbium-Doped Fibre Amplifier Operating at 1.54 μm ," *Electron. Lett.* **23**:1026, 1987.
11. S. Poole, D. N. Payne, R. Mears, M. Fermann, and R. Laming, "Fabrication and Characterization of Low-Loss Optical Fibers Containing Rare-Earth Ions," *J. Lightwave Technol.* **4**(7):870–876, 1986.
12. E. Snitzer, H. Po, F. Hakimi, R. Tumminelli, and B. C. McCollum, "Double-Clad, Offset Core Nd fiber Laser," *Opt. Fiber Sensors I*, 1988.
13. V. P. Gapontsev, P. I. Sawvsky, and I. E. Samartsev, "1.5 μm Erbium Glass Lasers," in *Conference on Lasers and ElectroOptics*, San Jose, Calif., 1990.
14. D. Minelly, E. R. Taylor, K. P. Iedrzuewski, J. Wang, and D. N. Payne, "Laser-Diode Pumped Neodymium-Doped Fibre Laser with Output Power >1 W," in *Conference on Lasers and ElectroOptics*, 1992.
15. H. Po, J. D. Cao, B. M. Laliberte, R. A. Minns, R. F. Robinson, B. H. Rockney, R. R. Tricca, and Y. H. Zhang, "High Power Neodymium-Doped Single Transverse Mode Fibre Laser," *Electron. Lett.* **29**(17):1500–1501, 1993.
16. V. Dominic, S. MacCormack, R. Waarts, S. Sanders, S. Bicknese, R. Dohle, E. Wolak, P. Yeh, and E. Zucker, "110 W Fibre Laser," *Electron. Lett.* **35**(14):1158–1160, 1999.
17. C. C. Ranaud, H. L. Offerhaus, J. A. Alvarez-Chavez, J. Nilsson, W. A. Clarkson, P. Turner, D. J. Richardson, and A. B. Grudinin, "Characteristics of Q-Switched Cladding-Pumped Ytterbium-Doped Fiber Lasers with Different High-Energy Fiber Designs," *IEEE J. Quant. Electron.* **37**(2):199–206, 2001.
18. J. A. Alvarez-Chavez, H. L. Offerhaus, J. Nilsson, P. Turner, W. A. Clarkson, and D. J. Richardson, "High-Energy, High-Power Ytterbium-Doped Q-Switched Fiber Laser," *Opt. Lett.* **25**(1):37–39, 2000.
19. N. G. R. Broderick, H. L. Offerhaus, D. J. Richardson, R. A. Sammut, J. Caplen, and L. Dong, "Large Mode Area Fibers for High Power Applications," *Opt. Fiber Technol.* **5**(2):185–196, 1999.
20. F. DiTeodoro and C. D. Brooks, "Multi-mJ Energy, Multi-MW Peak Power Photonic Crystal Fiber Amplifiers with Near Diffraction Limited Output," in *CLEO*, Baltimore, 2007.
21. V. Fomin, A. Mashkin, M. Abramov, A. Ferin, and V. Gapontsev, "3 kW Yb Fiber Lasers with a Single Mode Output," in *International Symposium on High Power Fiber Lasers and their Applications*, St. Petersburg, 2006.
22. H. Zellmer, A. Tünnermann, H. Welling, and V. Reichel, "Double-Clad Fiber Laser with 30 W Output Power," in *Optical Amplifiers and Their Applications*, M. Zervas, A. Willner, and S. Sasaki, eds., Vol. 16 of OSA Trends in Optics and Photonics Series (Optical Society of America, 1997), paper FAW18.
23. K. Ueda, H. Sekiguchi, and H. Kan, "1 kW CW Output from Fiber-Embedded Disk Lasers," in *Proc. CLEO*. Long Beach, CA, 2002.
24. Y. Jeong, J. Sahu, R. B. Williams, D. J. Richardson, K. Furusawa, and J. Nilsson, "Ytterbium-Doped Largecore Fibre Laser with 272 W Output Power," *Electron. Lett.* **39**:977–978, 2003.
25. J. Limpert, A. Liem, H. Zellmer, and A. Tünnermann, "500 W Continuous-Wave Fibre Laser with Excellent Beam Quality," *Electron. Lett.* **39**(8):645–647, 2003.
26. V. P. Gapontsev, N. S. Platonov, O. Shkurihin, and I. Zaitsev, "400 W Low Noise Single-Mode CW Ytterbium Fiber Laser with an Integrated Fiber Delivery," in *Proc. CLEO*, Baltimore, 2003.
27. Y. Jeong, J. Sahu, D. Payne, and J. Nilsson, "Ytterbium-Doped Large-Core Fiber Laser with 1.36 kW Continuous-Wave Output Power," *Opt. Exp.* **12**(25):6088–6092, 2004.
28. Y. Jeong, J. Sahu, S. Baek, C. Alegria, D. B. S. Soh, C. Codemard, and J. Nilsson, "Cladding-Pumped Ytterbium-Doped Large-Core Fiber Laser with 610 W of Output Power," *Opt. Commun.* **234**(1–6):315–319, 2004.

29. C. H. Liu, B. Ehlers, F. Doerfel, S. Heinemann, A. Carter, K. Tankala, J. Farroni, and A. Galvanauskas, "810W Continuous-Wave and Single-Transverse-Mode Fibre Laser Using 20 μm Core Yb-Doped Double-Clad Fibre," *Electron. Lett.* **40**(23):1471–1472, 2004.
30. Y. Jeong, J. K. Sahu, D. N. Payne, and J. Nilsson, "Ytterbium-Doped Large-Core Fibre Laser with 1 kW of Continuous-Wave Output Power," *Electron. Lett.* **40**(8):470–472, 2004.
31. A. Liem, J. Limpert, H. Zellmer, A. Tunnermann, V. Reichel, K. Morl, S. Jetschke, H. R. Unger, Muller, J. Kirchof, T. Sandrock, and T. A. Harschak, "1.3 kW Yb-Doped Fiber Laser with Excellent Beam Quality," in *Proc. CLEO*, 2004. San Francisco, CA, USA.
32. V. P. Gapontsev, D. V. Gapontsev, N. S. Platonov, O. Shkurihin, V. Fomin, A. Mashkin, M. Abramov, and A. Ferin, "2 kW CW Ytterbium Fiber Laser with Record Diffraction Limited Brightness," in *Proc. CLEO Europe*, Munich, Germany, 2005.
33. C. H. Liu, A. Galvanauskas, B. Ehlers, F. Doerfel, S. Heinemann, A. Carter, A. Tanaka, and J. Farroni, "810-W Single Transverse Mode Yb-Doped Fiber Laser," in *ASSP*, Santa Fe, N. Mex., 2004.
34. N. S. Platonov, V. Gapontsev, V. P. Gapontsev, and V. Shumilin, "135 W CW Fiber Laser with Perfect Single Mode Output," in *Proc. CLEO*, Long Beach, Calif., 2002.
35. J. W. Dawson, M. J. Messerly, R. J. Beach, M. Y. Shverdin, E. A. Stappaerts, A. K. Sridharan, P. H. Pax, J. E. Heebner, C. W. Siders, and C. P. J. Barty, "Analysis of the Scalability of Diffraction-Limited Fiber Lasers and Amplifiers to High Average Power," *Opt. Exp.* **16**(17):13240–13266, 2008.
36. J. Nilsson, W. A. Clarkson, R. Selvas, J. K. Sahu, P. W. Turner, S. U. Alam, and A. B. Grudinin, "High-Power Wavelength-Tunable Cladding-Pumped Rare-Earth-Doped Silica Fiber Lasers," *Opt. Fiber Technol.* **10**(1): 5–30, 2004.
37. J. Nilsson, J. K. Sahu, Y. Jeong, W. A. Clarkson, R. Selvas, A. B. Grudinin, and S. U. Alam, "High Power Fiber Lasers: New Developments," *Proc. of SPIE* **4974**:51, 2003.
38. W. Yong, "Heat Dissipation in Kilowatt Fiber Power Amplifiers," *IEEE J. Quant. Electron.* **40**(6):731–740, 2004.
39. W. Yong, X. Chang-Qing, and P. Hong, "Thermal Effects in Kilowatt Fiber Lasers," *Photonics Technol. Lett. IEEE* **16**(1):63–65, 2004.
40. M. K. Davis, M. J. F. Digonnet, and R. H. Pantell, "Thermal Effects in Doped Fibers," *J. Lightwave Technol.* **16**(6):1013–1023, 1998.
41. B. C. Stuart, M. D. Feit, S. Herman, A. M. Rubenchik, B. W. Shore, and M. D. Perry, "Nanosecond-to-Femtosecond Laser-Induced Breakdown in Dielectrics," *Phys. Rev. B* **53**(4):1749, 1996.
42. V. I. Kovalev and R. G. Harrison, "Suppression of Stimulated Brillouin Scattering in High-Power Single-Frequency Fiber Amplifiers," *Opt. Lett.* **31**(2):161–163, 2006.
43. D. N. Payne, Y. Jeong, J. Nilsson, J. K. Sahu, D. B. S. Soh, C. Alegria, P. Dupriez, et al., *Kilowatt-Class Single-Frequency Fiber Sources*, (Invited Paper): SPIE Volume 5709: Fiber Lasers II: Technology, Systems and Applications 1, 2005.
44. P. D. Dragic, L. Chi-Hung, G. C. Papen, and A. Galvanauskas, "Optical Fiber with an Acoustic Guiding Layer for Stimulated Brillouin Scattering Suppression," in *Conference on Lasers and Electro-Optics/Quantum Electronics and Laser Science and Photonic Applications Systems Technologies*: Optical Society of America, 2005.
45. R. Paschotta, J. Nilsson, A. C. Tropper, D. C. Hanna, "Ytterbium-Doped Fiber Amplifiers," *IEEE J. Quant. Electron.* **33**(7):1049–1056, 1997.
46. W. Yong and P. Hong, "Dynamic Characteristics of Double-Clad Fiber Amplifiers for High-Power Pulse Amplification," *J. Lightwave Technol.* **21**(10):2262–2270, 2003.
47. H. M. Pask, R. J. Carman, D. C. Hanna, A. C. Tropper, C. J. Mackenchnie, P. R. Barber, and J. M. Dawes, "Ytterbium-Doped Silica Fiber Lasers: Versatile Sources for the 1–1.2 μm Region," *IEEE J. Sel. Top. Quant. Electron.* **1**(1):2–13, 1995.
48. W. W. Rigrod, "Saturation Effects in High-Gain Lasers," *J. Appl. Phys.* **36**(8):2487–2490, 1965.
49. J. A. Buck, *Fundamentals of Optical Fibers*, New York:Wiley-IEEE Press, 2004.
50. D. Marcuse, "Loss Analysis of Single Mode Fiber Splices," *Bell System Technical Journal*, **56**:703–718, 1977.
51. A. Yariv, *Optical Electronics in Modern Communications*, 5th ed., New York: Oxford University Press, 1997.
52. C. Ullmann and V. Krause, (eds.), *Diode Optics and Diode Lasers*, U.P. Office, United States, 1999.
53. W. A. Clarkson and D. C. Hanna, "Two Mirror Beam Shaping Technique for High Power Diode Bars," *Opt. Lett.* **21**:375–377, 1996.

54. L. Goldberg, J. P. Koplow, and D. A. V. Kliner, "Highly Efficient 4-W Yb-Doped Fiber Amplifier Pumped by a Broad-Stripe Laser Diode," *Opt. Lett.* **24**(10):673–675, 1999.
55. D. J. Ripin, and L. Goldberg, "High Efficiency Side-Coupling of Light into Optical Fibres Using Imbedded V-Grooves," *Electron. Lett.* **31**(25):2204–2205, 1995.
56. F. Gonthier, L. Martineau, N. Azami, M. Faucher, F. Seguin, D. Stryckman, and A. Villeneuve, *High-Power All-Fiber Components: the Missing Link for High-Power Fiber Lasers*, in *Fiber Lasers: Technology, Systems, and Applications*: SPIE, San Jose, Calif., 2004.
57. C. Headley III, M. Fishteyn, A. D. Yablon, M. J. Andrejco, K. Brar, J. Mann, M. D. Mermelstein, and D. J. DiGiovanni, "Tapered Fiber Bundles for Combining Laser Pumps (Invited Paper)," in *Fiber Lasers II: Technology, Systems, and Applications*: SPIE, San Jose, Calif., 2005.
58. A. Kosterin, V. Temyanko, M. Fallahi, and M. Mansuripur, "Tapered Fiber Bundles for Combining High-Power Diode Lasers," *Appl. Opt.* **43**(19):3893–3900, 2004.
59. B. Samson, and G. Frith, *Diode Pump Requirements for High Power Fiber Lasers*, in *ICALEO*, Orlando, Fla, 2007.
60. J. Xu, J. Lu, G. Kumar, J. Lu, and K. Ueda, "A Non-Fused Fiber Coupler for Side-Pumping of Double-Clad Fiber Lasers," *Opt. Commun.* **220**(4–6):389–395, 2003.
61. C. Codemard, K. Yla-Jarkko, J. Singleton, P. W. Turner, I. Godfrey, S. U. Alam, J. Nilsson, J. Sahu, and A. B. Grudinin, "Low-Noise Intelligent Cladding-Pumped L-Band EDFA," *Photon. Technol. Lett. IEEE* **15**(7):909–911, 2003.
62. Alam, S.-U. J. Nilsson, P. W. Turner, M. Ibsen, and A. B. Grudinin, "Low Cost Multi-Port Reconfigurable Erbium Doped Cladding Pumped Fiber Amplifier," in *Proc. ECOC'00*, vol. 2, 2000, pp. 119–120. Munich, Germany, 2000.
63. X. J. Gu, and Y. Liu, "The Efficient Light Coupling in a Twin-Core Fiber Waveguide," *Photonics Technol. Lett. IEEE* **17**(10):2125–2127, 2005.
64. P. Polynkin, V. Temyanko, M. Mansuripur, and N. Peyghambarian, "Efficient and Scalable Side Pumping Scheme for Short High-Power Optical Fiber Lasers and Amplifiers," *Photonics Technol. Lett. IEEE*, **16**(9):2024–2026, 2004.
65. "SPI Lasers," available: www.spilasers.com, accessed on: Dec. 2008.
66. P. Even, and D. Pureur. "High-Power Double-Clad Fiber Lasers: A Review" in *Optical Devices for Fiber Communication III*: SPIE, San Jose, CA, 2002.
67. F. Haxsen, A. Ruehl, M. Engelbrecht, D. Wandt, U. Morgner, and D. Kracht, "Stretched-Pulse Operation of Athulium-Doped Fiber Laser," *Opt. Exp.* **16**(25):20471–20476, 2008.
68. B. Ortaç, C. Lecaplain, A. Hideur, T. Schreiber, J. Limpert, and A. Tünnermann, "Passively Mode-Locked Single-Polarization Microstructure Fiber Laser," *Opt. Exp.* **16**(3):2122–2128, 2008.
69. J. Nilsson, and B. Jaskorzynska, "Modeling and Optimization of Low-Repetition-Rate High-Energy Pulse Amplification in CW-Pumped Erbium-Doped Fiber Amplifiers," *Opt. Lett.* **18**(24):2099, 1993.
70. P. Wang, J. K. Sahu, and W. A. Clarkson, "High-Power Broadband Ytterbium-Doped Helical-Core Fiber Superfluorescent Source," *Photonics Technol. Lett. IEEE*, **19**(5):300–302, 2007.
71. M. Kihara, and S. Tomita, "Loss Characteristics of Thermally Diffused Expanded Core Fibers," *Photonics Technol. Lett.* **4**(12):1390–1391, 1992.
72. K. O. Hill, and G. Meltz, "Fiber Bragg Grating Technology Fundamentals and Overview," *J. Lightwave Technol.* **15**(8):1263–1276, 1997.
73. L. B. Fu, G. D. Marshall, J. A. Bolger, P. Steinvurzel, E. C. Magi, M. J. Withford, and B. J. Eggleton, "Femtosecond Laser Writing Bragg Gratings in Pure Silica Photonic Crystal Fibres," *Electron. Lett.* **41**(11):638–640, 2005.
74. N. Jovanovic, A. Fuerbach, G. D. Marshall, M. J. Withford, and S. D. Jackson, "Stable High-Power Continuous-Wave Yb³⁺-Doped Silica Fiber Laser Utilizing a Point-by-Point Inscribed Fiber Bragg Grating," *Opt. Lett.* **32**(11):1486–1488, 2007.
75. Y. Lai, A. Martinez, I. Khrushchev, and I. Bennion, "Distributed Bragg Reflector Fiber Laser Fabricated by Femtosecond Laser Inscription," *Opt. Lett.* **31**(11):1672–1674, 2006.
76. A. D. Yablon, *Optical Fiber Fusion Splicing*, Springer, New York, 2005.
77. D. J. Richardson, P. Britton, and D. Taverner, "Diode-Pumped, High-Energy, Single Transverse Mode Q-Switch Fibre Laser," *Electron. Lett.* **33**(23):1955–1956, 1997.

78. H. L. Offerhaus, N. G. Broderick, D. J. Richardson, R. Sammut, J. Caplen, and L. Dong, "High-Energy Single-Transverse-Mode Q-Switched Fiber Laser Based on a Multimode Large-Mode-Area Erbium-Doped Fiber," *Opt. Lett.* **23**(21):1683–1685, 1998.
79. T. Bhutta, J. I. Mackenzie, D. P. Shepherd, and R. J. Beach, "Spatial Dopant Profiles for Transverse-Mode Selection in Multimode Waveguides," *J. Opt. Soc. Am. B* **19**(7):1539–1543, 2002.
80. J. P. Koplow, D. A. V. Kliner, and L. Goldberg, "Single-Mode Operation of a Coiled Multimode Fiber Amplifier," *Opt. Lett.* **25**(7):442–444, 2000.
81. D. Marcuse, "Curvature Loss Formula for Optical Fibers," *J. Opt. Soc. Am.* **66**(3):216, 1976.
82. D. Marcuse, "Field Deformation and Loss Caused by Curvature of Optical Fibers," *J. Opt. Soc. Am.* **66**(4):311, 1976.
83. D. Marcuse, "Influence of Curvature on the Losses of Doubly Clad Fibers," *Appl. Opt.* **21**(23):4208, 1982.
84. R. T. Schermer, "Mode Scalability in Bent Optical Fibers," *Opt. Exp.* **15**(24):15674–15701, 2007.
85. J. W. Dawson, R. J. Beach, I. Jovanovic, W. Benoit, Z. Liao, and S. Payne, "Large Flattened Mode Optical Fiber for Reduction of Nonlinear Effects," in *Proc. SPIE Volume 5335 Fiber Lasers: Technology, Systems and Applications*. 2004.
86. A. K. Ghatik, I. C. Goyal, and R. Jindal, "Design of a Waveguide Refractive Index Profile to Obtain a Flat Modal Field," in *International Conference on Fiber Optics and Photonics: Selected Papers from Photonics India '98*, New Delhi, India: SPIE, 1999.
87. P. Facq, F. de Fornel, and F. Jean, "Tunable Single-Mode Excitation in Multimode Fibres," *Electron. Lett.* **20**(15):613–614, 1984.
88. M. E. Fermann, "Single-Mode Excitation of Multimode Fibers with Ultrashort Pulses," *Opt. Lett.* **23**(1):52–54, 1998.
89. C. D. Stacey, R. M. Jenkins, J. Banerji, and A. R. Davies, "Demonstration of Fundamental Mode only Propagation in Highly Multimode Fibre for High Power EDFAs," *Opt. Commun.* **269**(2):310–314, 2007.
90. K. C. Hou, M. Y. Cheng, D. Engin, R. Changkakoti, P. Mamidipudi, and A. Galvanauskas, "Multi-MW Peak Power Scaling of Single Transverse Mode Pulses Using 80 μm Core Yb-Doped LMA Fibers," Optical Society of America *Advanced Solid State Photonics*, Tahoe, Calif., 2006.
91. A. Galvanauskas, C. Ming-Yuan, H. Kai-Chung, and A.K.-H.L. Kai-Hsiu Liao, "High Peak Power Pulse Amplification in Large-Core Yb-Doped Fiber Amplifiers," *IEEE J. Sel. Top. in Quant. Electron.* **13**(3):559–566, 2007.
92. K. -C. Hou, S. George, A. G. Mordovanakis, K. Takenoshita, J. Nees, B. Lafontaine, M. Richardson, and A. Galvanauskas, "High Power Fiber Laser Driver for Efficient EUV Lithography Source with Tin-Doped Water Droplet Targets," *Opt. Exp.* **16**(2):965–974, 2008.
93. J. A. Alvarez-Chavez, A. B. Grudinin, J. Nilsson, P. W. Turner, and W. A. Clarkson, "Mode Selection in High Power Cladding Pumped Fiber Lasers with Tapered Section," p. 247–248 in *CLEO*. Washington D.C., 1999.
94. W. A. Clarkson, "Short Course Notes SC270: High Power Fiber Lasers and Amplifiers," in *Short Courses at CLEO 2007*, Optoelectronics Research Center, University of Southampton, Baltimore, Md., 2007.
95. V. Filippov, Y. Chamorovskii, J. Kerttula, K. Golant, M. Pessa, and O. G. Okhotnikov, "Double Clad Tapered Fiber for High Power Applications," *Opt. Express* **16**(3):1929–1944, 2008.
96. D. Taverner, D. J. Richardson, L. Dong, J. E. Caplen, K. Williams, and R. V. Penty, "158- μJ Pulses from a Single-Transverse-Mode, Large-Mode-Area Erbium-Doped Fiber Amplifier," *Opt. Lett.* **22**(6):378–380, 1997.
97. P. S. J. Russell, "Photonic-Crystal Fibers," *J. Lightwave Technol.* **24**(12):4729–4749, 2006.
98. J. C. Knight, "Photonic Crystal Fibers and Fiber Lasers (Invited)," *J. Opt. Soc. Am. B* **24**(8):1661–1668, 2007.
99. T. A. Birks, P. J. Roberts, P. S. J. Russell, D. Atkin, and T. Shepherd, "Full 2-D Photonic Bandgaps in Silica/air Structures," *Electron. Lett.* **31**(22):1941–1943, 1995.
100. J. Broeng, D. Mogilevstev, S. E. Barkou, and A. Bjarklev, "Photonic Crystal Fibers: A New Class of Optical Waveguides," *Optical Fiber Technology*, **5**(3):305–330, 1999.
101. J. C. Knight, T. A. Birks, P. St. J. Russell, and D. M. Atkin, "All-Silica Single-Mode Optical Fiber with Photonic Crystal Cladding," *Opt. Lett.* **21**, 1547–1549 (1996).
102. W. J. Wadsworth, J. C. Knight, and P. S. and J. Russell, "Large Mode Area Photonic Crystal Fibre Laser," in *Conference in Lasers and Electrooptics*, Washington D.C., 2001.

103. W. J. Wadsworth, J. C. Knight, and P. St. J. Russell, "Yb³⁺-Doped Photonic Crystal Fibre Laser," *Electron. Lett.* **36**:1452–1453, 2000.
104. Crystal-Fibre. *PCF Technology Tutorial*, available: www.crystal-fibre.com, accessed on: Dec. 15, 2008.
105. T. A. Birks, J. C. Knight, and P. S. J. Russell, "Endlessly Single-Mode Photonic Crystal Fiber," *Opt. Lett.* **22**(13):961–963, 1997.
106. T. M. Monro, D. J. Richardson, N. G. R. Broderick, and P. J. Bennett, "Holey Optical Fibers: An Efficient Modal Model," *J. Lightwave Technol.* **17**(6):1093, 1999.
107. K. Furusawa, A. Malinowski, J. Price, T. Monro, J. Sahu, J. Nilsson, and D. Richardson, "Cladding Pumped Ytterbium-Doped Fiber Laser with Holey Inner and Outer Cladding," *Opt. Exp.* **9**(13):714–720, 2001.
108. K. P. Hansen, J. Broeng, P. M. W. Skovgaard, J. P. Folkenberg, M. D. Nielsen, A. Petersson, T. P. Hansen, et al., "High Power Photonic Crystal Fiber Lasers: Design, Handling and Subassemblies," in *Photonics West.: San Jose, Calif.*, 2005.
109. W. Wadsworth, R. Percival, G. Bouwmans, J. Knight, and P. Russell, "High Power Air-Clad Photonic Crystal Fibre Laser," *Opt. Exp.* **11**(1):48–53, 2003.
110. B. Zintzen, T. Langer, J. Geiger, D. Hoffmann, and P. Loosen, "Heat Transport in Solid and Air-Clad Fibers for High-Power Fiber Lasers," *Opt. Exp.* **15**(25):16787–16793, 2007.
111. J. Limpert, O. Schmidt, J. Rothhardt, F. Röser, T. Schreiber, A. Tünnermann, S. Ermeneux, P. Yvernault, and F. Salin, "Extended single-mode Photonic Crystal Fiber Lasers," *Opt. Exp.* **14**(7):2715–2720, 2006.
112. G. Bonati, H. Voelkel, T. Gabler, U. Krause, A. Tuennermann, A. Liem, T. Schreiber, S. Nolte, and H. Zellmer, "1.53 kW from a Single Yb-Doped Photonic Crystal Fiber Laser," in *Late Breaking News, Photonics West*, San Jose, Calif., 2005.
113. O. Schmidt, J. Rothhardt, T. Eidam, F. Röser, J. Limpert, A. Tünnermann, K.P. Hansen, C. Jakobsen, and J. Broeng, "Single-Polarization Ultra-Large-Mode-Area Yb-Doped Photonic Crystal Fiber," *Opt. Exp.* **16**(6):3918–3923, 2008.
114. J. Broeng, G. Vienne, A. Petersson, P. M. W. Skovgaard, J. P. Folkenberg, M. D. Nielsen, C. Jakobsen, H. R. Simonsen, and N. A. Mortensen, "Air-Clad Photonic Crystal Fibers for High-Power Single-Mode Lasers," in *Photonics West.*, San Jose, Calif., 2004.
115. V. V. R. Kumar, A. George, W. Reeves, J. Knight, P. Russell, F. Omenetto, and A. Taylor, "Extruded Soft Glass Photonic Crystal Fiber for Ultrabroad Supercontinuum Generation," *Opt. Exp.* **10**(25):1520–1525, 2002.
116. H. Ebendorff-Heidepriem, and T. M. Monro, "Extrusion of Complex Preforms for Microstructured Optical Fibers," *Opt. Exp.* **15**(23):15086–15092, 2007.
117. A. Isomäki, and O. G. Okhotnikov, "Femtosecond Soliton Mode-Locked Laser Based on Ytterbium-Doped Photonic Bandgap Fiber," *Opt. Exp.* **14**(20):9238–9243, 2006.
118. A. Wang, A. K. George, and J. C. Knight, "Three-Level Neodymium Fiber Laser Incorporating Photonic Bandgap Fiber," *Opt. Lett.* **31**(10):1388–1390, 2006.
119. F. Qiang, W. Zhi, K. Guiyun, Long Jin., Yang Yue., Jiangbing Du, Qing Shi., Zhanyuan Liu., Bo Liu., Yange Liu., Shuzhong Yuan., and Xiaoyi Dong. "Proposal for All-Solid Photonic Bandgap Fiber with Improved Dispersion Characteristics," *Photonics Technol. Lett. IEEE* **19**(16):1239–1241, 2007.
120. L. Dong, J. Li, H. McKay, A. Marcinkevicius, B. Thomas, M. Moore, L. Fu, and M. E. Fermann, "Robust and Practical Optical Fibers for Single Mode Operation with Core Diameters up to 170 μm ," in *CLEO 2008*, San Jose, Calif., 2008.
121. L. Dong, X. Peng, and J. Li, "Leakage Channel Optical Fibers with Large Effective Area," *J. OSA B*, **24**(8):1689–1697, 2007.
122. Wong, William S., X. Peng, Mclaughlin, M. Joseph, and L. Dong, "Breaking the Limit of Maximum Effective Area for Robust Single-Mode Propagation in Optical Fibers," *Opt. Lett.* **30**(21):2855–2857, 2005.
123. C. H. Liu, G. Chang, N. Litchinitser, D. Guertin, N. Jacobsen, K. Tankala, and A. Galvanauskas, "Chirally Coupled Core Fibers at 1550-nm and 1064-nm for Effectively Single-Mode Core Size Scaling," in *CLEO*, Baltimore, 2007.
124. M. C. Swan, C. H. Liu, D. Guertin, N. Jacobsen, A. Tanaka, and A. Galvanauskas, "33 μm Core Effectively Single-Mode Chirally-Coupled-Core Fiber Laser at 1064-nm," in *Optical Fiber Communication Conference & Exposition and the National Fiber Optic Engineers Conference*. San Diego, Calif., 2008.
125. R. J. Beach, M. D. Feit, R. H. Page, L. D. Brasure, R. Wilcox, and S. A. Payne, "Scalable Antiguided Ribbon Laser," *J. Opt. Soc. Am. B*, **19**(7):1521–1534, 2002.

126. M. Wrage, P. Glas, M. Leitner, T. Sandrock, N. N. Elkin, A. P. Napartovich, and A. G. Sukharev, "Experimental and Numerical Determination of Coupling Constant in a Multicore Fiber," *Opt. Commun.* **175**(1–3):97–102, 2000.
127. T. Erdogan, "Cladding-Mode Resonances in Short- and Long-Period Fiber Grating Filters," *J. Opt. Soc. Am. A*, **14**(8):1760–1773, 1997.
128. S. Ramachandran, J. W. Nicholson, S. Ghalmi, M. F. Yan, P. Wisk, E. Monberg, and F. V. Dimarcello, "Light Propagation with Ultralarge Modal Areas in Optical Fibers," *Opt. Lett.* **31**(12):1797–1799, 2006.
129. S. Ramachandran, and S. Ghalmi, "Diffraction Free" Self Healing Bessel Beams from Fibers," in *CLEO 2008*, San Jose, Calif., 2008.
130. S. Suzuki, A. Schülzgen, and N. Peyghambarian, "Single-Mode Fiber Laser Based on Core-Cladding Mode Conversion," *Opt. Lett.* **33**(4):351–353, 2008.
131. S. Ramachandran, S. Ghalmi, J. W. Nicholson, M. F. Yan, P. Wisk, E. Monberg, and F. V. Dimarcello, "Anomalous Dispersion in a Solid, Silica-Based Fiber," *Opt. Lett.* **31**(17):2532–2534, 2006.
132. M. Schultz, O. Prochnow, A. Ruehl, D. Wandt, D. Kracht, S. Ramachandran, and S. Ghalmi, "Sub-60-fs Ytterbium-Doped Fiber Laser with a Fiber-Based Dispersion Compensation," *Opt. Lett.* **32**(16):2372–2374, 2007.
133. M. D. Mermelstein, S. Ramachandran, J. M. Fini, and S. Ghalmi, "SBS Gain Efficiency Measurements and Modeling in a 1714 μm^2 Effective Area LP₀₈ Higher-Order Mode Optical Fiber," *Opt. Exp.* **15**(24):15952–15963, 2007.
134. A. E. Siegman, "Propagating Modes in Gain-Guided Optical Fibers," *J. Opt. Soc. Am. A*, **20**(8):1617–1628, 2003.
135. A. E. Siegman, "Gain-Guided, Index-Antiguide Fiber Lasers," *J. Opt. Soc. Am. B* **24**(8):1677–1682, 2007.
136. Y. Chen, T. McComb, V. Sudesh, M. Richardson, and M. Bass, "Very Large-Core, Single-Mode, Gain-Guided, Index-Antiguide Fiber Lasers," *Opt. Lett.* **32**(17):2505–2507, 2007.
137. Y. Chen, T. McComb, V. Sudesh, M. C. Richardson, M. Bass, "Lasing in a Gain-Guided Index Antiguide Fiber," *J. Opt. Soc. Am. B* **24**(8):1683–1688, 2007.
138. V. Sudesh, T. McComb, Y. Chen, M. Bass, M. Richardson, J. Ballato, and A. E. Siegman, "Diode-Pumped 200 μm Diameter Core, Gain-Guided, Index-Antiguide Single Mode Fiber Laser," *Appl. Phys. B: Lasers and Opt.* **90**(3):369–372, 2008.
139. P. Pavel, T. Valery, M. Jerome, and N. Peyghambarian, "Dramatic Change of Guiding Properties in Heavily Yb-Doped, Soft-Glass Active Fibers Caused by Optical Pumping," *Appl. Phys. Lett.* **90**(24):241106-1–241106-3, 2007.
140. M. J. F. Digonnet, *Rare-Earth-Doped Fiber Lasers and Amplifiers*, CRC Press, New York, 2001.
141. P. Hamamatsu, "The Fiber Disk Laser Explained," *Nat Photon*, **sample**(sample):14–15, 2006.
142. Y. Jeong, S. Yoo, C. A. Codemard, J. Nilsson, J. K. Sahu, D. N. Payne, R. Horley, P. W. Turner, L. M. B. Hickey, A. Harker, M. Lovelady, and A. Piper, "Erbium:Ytterbium Co-Doped Large-Core Fiber Laser with 297 W Continuous-Wave Output Power," *IEEE J. Sel. Top. Quant. Electron.* **13**(3):573–579.
143. X. Zhu, "Mid-IR ZBLAN Fiber Laser Approaches 10 W Output," in *Laser Focus World*, **44**(3), PennWell Corp. Tulsa Okla., 2007.
144. X. Zhu, and R. Jain, "10-W-Level Diode-Pumped Compact 2.78 μm ZBLAN Fiber Laser," *Opt. Lett.* **32**(1):26–29, 2007.
145. X. Zhu, and R. Jain, "Compact 2 W Wavelength-Tunable Er:ZBLAN Mid-Infrared Fiber Laser," *Opt. Lett.* **32**(16):2381–2383, 2007.
146. S. D. Jackson, F. Bugge, and G. Erbert, "High-Power and Highly Efficient Tm³⁺-Doped Silica Fiber Lasers Pumped with Diode Lasers Operating at 1150 nm," *Opt. Lett.* **32**(19), 2007.
147. S. D. Jackson, "Power Scaling Method for 2-mm Diode-Cladding-Pumped Tm³⁺-Doped Silica Fiber Lasers that Uses Yb³⁺ Codoping," *Opt. Lett.* **28**(22):2192, 2003.
148. D. V. Gapontsev, N. S. Platonov, M. Meleshkevich, A. Drozhzhin, and V. Sergeev, "415 W Single Mode CW Thulium Fiber Laser in all Fiber Format," in *CLEO Europe 2007 Munich, Germany*, 2007.
149. S. D. Jackson, "Cross Relaxation and Energy Transfer Upconversion Processes Relevant to the Functioning of 2 μm Tm³⁺-Doped Silica Fibre Lasers," *Opt. Commun.* **230**(1–3):197–203, 2004.
150. J. Wu, Z. Yao, J. Zong, and S. Jiang, "Highly Efficient High-Power Thulium-Doped Germanate Glass Fiber Laser," *Opt. Lett.* **32**(6):638–640, 2007.

151. E. Slobodtchikov, P. F. Moulton, and G. Frith, "Efficient, High Power, Tm-Doped Silica Fiber Laser," in *ASSP 2007 Postdeadline*. Vancouver, Canada, 2007.
152. S. D. Jackson, and S. Mossman, "Efficiency Dependence on the Tm and Al Concentrations for Tm-Doped Silica Double-Clad Fiber Lasers," *Appl. Opt.* **42**(15), 2003.
153. S. D. Jackson, and T. A. King, "Theoretical Modeling of Tm-Doped Silica Fiber Lasers," *J. Lightwave Technol.* **17**(5):948–956, 1999.
154. P. F. Moulton, G. A. Rines, E. V. Slobodtchikov, K. F. Wall, G. Frith, B. Samson, and A. L. G. Carter, "Tm-Doped Fiber Lasers: Fundamentals and Power Scaling," *IEEE J. Sel. Top. Quant. Electron.* **15**(1): 85–92, 2009.
155. S. D. Jackson, "Midinfrared Holmium Fiber Lasers," *J. Quant. Electron. IEEE* **42**(2):187–191, 2006.
156. S. D. Jackson, F. Bugge, and G. Erbert, "Directly Diode-Pumped Holmium Fiber Lasers," *Opt. Lett.* **32**(17):2496–2498, 2007.
157. S. D. Jackson, F. Bugge, and G. Erbert, "High-Power and Highly Efficient Diode-Cladding-Pumped Ho³⁺-Doped Silica Fiber Lasers," *Opt. Lett.* **32**(22):3349–3351, 2007.
158. S. D. Jackson, and S. Mossman, "Diode-Cladding-Pumped Yb³⁺, Ho³⁺-Doped Silica Fiber Laser Operating at 2.1- μ m," *Appl. Opt.* **42**(18):3546–3549, 2003.
159. S.D. Jackson, A. Sabella, A. Hemming, S. Bennetts, and D. G. Lancaster, "High-Power 83 W Holmium-Doped Silica Fiber Laser Operating with High Beam Quality," *Opt. Lett.* **32**(3):241–243, 2007.
160. D. Hewak, (ed.), *Properties, Processing and Applications of Glass and Rare Earth Doped Glasses for Optical Fibers*, INSPEC Publications, London, 1998.
161. D. J. Digiovanni, A. M. Vengsarkar, J. L. Wagener, and R. S. Windeler, U.S. Patent 5,802,236 *Article Comprising a Micro-Structured Optical Fiber, and Method of Making Such Fiber*: United States Patent Office. USA, 1998.
162. H. Bookey, K. Bindra, A. Kar, and B. A. Wherrett, "Telluride Glass Fibres for all Optical Switching: Nonlinear Optical Properties and Fibre Characterisation," in *Workshop on Fibre and Optical Passive Components, Proc. 2002 IEEE/LEOS*. Glasgow, Scotland 2002.
163. G. Boudebs, W. Berlatier, S. Cherukulappurath, F. Smektala, M. Guignard, and J. Troles, "Nonlinear Optical Properties of Chalcogenide Glasses at Telecommunication Wavelength Using Nonlinear Imaging Technique," in *Transparent Optical Networks, 2004, Proc. 2004 6th International Conf. on Transparent Optical Networks*, 2004.
164. C. C. Chen, Y. J. Wu, and L. G. Hwa, "Temperature Dependence of Elastic Properties of ZBLAN Glasses," *Materials Chemistry and Physics* **65**(3):306–309, 2000.
165. J. A. Harrington, *Infrared Fibers and Their Applications*, SPIE Press, Bellingham, Wash. 2004.
166. Y. T. Hayden, S. Payne, J. S. Hayden, J. Campbell, M. K. Aston, and M. Elder, U.S. Patent 5526369 *Phosphate Glass Useful in High Energy Lasers*, United States Patent Office, 1996.
167. "Laser Glass Properties," available: www.kigre.com, accessed on: Oct. 15 2008.
168. A. Kut'in, V. Polyakov, A. Gibin, and M. Churbanov, "Thermal Conductivity of (TeO₂)_{0.7}(WO₃)_{0.2}(La₂O₃)_{0.1} Glass," *Inorg. Mater.* **42**(12):1393–1396, 2006.
169. M. D. O'Donnell, K. Richardson, R. Stolen, A. B. Seddon, D. Furniss, V. K. Tikhomirov, C. Rivero, et al., "Tellurite and Fluorotellurite Glasses for Fiberoptic Raman Amplifiers: Glass Characterization, Optical Properties, Raman Gain, Preliminary Fiberization, and Fiber Characterization," *J. Am. Ceram. Soc.*, **90**(5):1448–1457, 2007.
170. T. Töpfer, J. Hein, J. Philipps, D. Ehrhart, and R. Sauerbrey, "Tailoring the Nonlinear Refractive Index of Fluoride-Phosphate Glasses for Laser Applications," *Appl. Phys. B: Lasers Opt.* **71**(2):203–206, 2000.
171. M. Yamane, and Y. Asahara, *Glasses for Photonics*, Cambridge University Press, Cambridge, UK, 2000.
172. G. Chen, Q. Zhang, G. Yang, and Z. Jiang, "Mid-Infrared Emission Characteristic and Energy Transfer of Ho³⁺-Doped Tellurite Glass Sensitized by Tm³⁺," *Journal of Fluorescence* **17**(3): 301–307, 2007.
173. M. Eichhorn, and S. D. Jackson, "Comparative Study of Continuous Wave Tm³⁺-Doped Silica and Fluoride Fiber Lasers," *Appl. Phys. B: Lasers Opt.* **90**(1):35–41, 2008.
174. B. Richards, Y. Tsang, D. Binks, J. Lousteau, and A. Jha, "Efficient ~2 μ m Tm³⁺ Doped Tellurite Fiber Laser," *Opt. Lett.* **33**(4):402–404, 2008.
175. S. Jiang, M. J. Myers, D. L. Rhonehouse, S. J. Hamlin, J. D. Myers, U. Griebner, R. Koch, and H. Schonagel, "Ytterbium-Doped Phosphate Laser Glasses," in *Solid State Lasers VI: SPIE*, San Jose, Calif., 1997.

176. Y. W. Lee, S. Sinha, M. J. F. Digonnet, R. L. Byer, and S. Jiang, "20 W Single-Mode Yb³⁺-Doped Phosphate Fiber Laser," *Opt. Lett.* **31**(22):3255–3257, 2006.
177. S. Wei, M. Leigh, Z. Jie, A. Zhidong Yao, and A. Shibin Jiang, "Photonic Narrow Linewidth GHz Source Based on Highly Codoped Phosphate Glass Fiber Lasers in a Single MOPA Chain," *Photon. Technol. Lett. IEEE* **20**(2):69–71, 2008.
178. N. P. Barnes, B. M. Walsh, D. J. Reichle, R. J. Deyoung, and S. Jiang, "Tm:Germanate Fiber Laser: Tuning and Q-Switching," *Appl. Phys. B: Lasers Opt.* **89**(2):299–304, 2007.
179. M. Eichhorn, "High-Peak-Power Tm-Doped Double-Clad Fluoride Fiber Amplifier," *Opt. Lett.* **30**(24):3329–3331, 2005.
180. M. Eichhorn, "Development of a High-Pulse-Energy Q-Switched Tm-Doped Double-Clad Fluoride Fiber Laser and Its Application to the Pumping of Mid-IR Lasers," *Opt. Lett.* **32**(9):1056–1058, 2007.
181. S. D. Jackson, and G. Anzueto-Sanchez, "Chalcogenide Glass Raman Fiber Laser," *Appl. Phys. Lett.* **88**(22):221106-3, 2006.
182. A. F. El-Sherif, and T. A. King, "Dynamics and Self-Pulsing Effects in Tm³⁺-Doped Silica Fibre Lasers," *Opt. Commun.* **208**(4–6):381–389, 2002.
183. T. McComb, V. Sudesh, and M. Richardson, "Volume Bragg Grating Stabilized Spectrally Narrow Tm Fiber Laser," *Opt. Lett.* **33**(8):881–883, 2008.
184. W. A. Clarkson, N. P. Barnes, P. W. Turner, J. Nilsson, and D. C. Hanna, "High-Power Cladding-Pumped Tm-Doped Silica Fiber Laser with Wavelength Tuning from 1860 to 2090 nm," *Opt. Lett.* **27**(22):1989–1991, 2002.
185. D. Y. Shen, J. K. Sahu, and W. A. Clarkson, "Highly Efficient Er, Yb-Doped Fiber Laser with 188W Free-Running and & gt; 100 W Tunable Output Power," *Opt. Exp.* **13**(13):4916–4921, 2005.
186. O. M. Efimov, L. B. Glebov, and V. I. Smirnov, "High-Frequency Bragg Gratings in a Photothermorefractive Glass," *Opt. Lett.* **25**(23):1693–1695 2000.
187. P. Jelger, and F. Laurell, "Efficient Skew-Angle Cladding-Pumped Tunable Narrow-Linewidth Yb-Doped Fiber Laser," *Opt. Lett.* **32**(24):3501–3503, 2007.
188. J. Yoonchan, C. Alegria, J. K. Sahu, L. Fu, M. Ibsen, C. Codemard, M. Mokhtar, and J. Nilsson, "A 43-W C-Band Tunable Narrow-Linewidth Erbium-Ytterbium Codoped Large-Core Fiber Laser," *Photon. Technol. Lett. IEEE* **16**(3):756–758, 2004.
189. N. Jovanovic, M. Åslund, A. Fuerbach, S. D. Jackson, G. D. Marshall, and M. J. Withford, "Narrow Linewidth, 100 W cw Yb³⁺-Doped Silica Fiber Laser with a Point-by-Point Bragg Grating Inscribed Directly into the Active Core," *Opt. Lett.* **32**(19):2804–2806, 2007.
190. P. Jelger, and F. Laurell, "Efficient Narrow-Linewidth Volume-Bragg Grating-Locked Nd:Fiber Laser," *Opt. Exp.* **15**(18):11336–11340, 2007.
191. J. W. Kim, P. Jelger, J. K. Sahu, F. Laurell, and W. A. Clarkson, "High-Power and Wavelength-Tunable Operation of an Er, Yb Fiber Laser Using a Volume Bragg Grating," *Opt. Lett.* **33**(11):1204–1206, 2008.
192. P. Jelger, P. Wang, J. K. Sahu, F. Laurell, and W. A. Clarkson, "High-Power Linearly-Polarized Operation of a Cladding-Pumped Yb Fibre Laser Using a Volume Bragg Grating for Wavelength Selection," *Opt. Exp.* **16**(13):9507–9512, 2008.
193. S. Tibuleac, and R. Magnusson, "Reflection and Transmission Guided-Mode Resonance Filters," *J. Opt. Soc. Am. A* **14**(7):1617–1626, 1997.
194. A. Mehta, R. C. Rumpf, Z. A. Roth, and E. J. Johnson, "Guided Mode Resonance Filter as a Spectrally Selective Feedback Element in a Double-Cladding Optical Fiber Laser," *Photon. Technol. Lett. IEEE* **19**(24):2030–2032, 2007.
195. J. Geng, J. Wu, S. Jiang, and J. Yu, "Efficient Operation of Diode-Pumped Single-Frequency Thulium-Doped Fiber Lasers near 2 μm," *Opt. Lett.* **32**(4):355–357, 2007.
196. Y. Jeong, J. Nilsson, J. K. Sahu, D. B. S. Soh, C. Alegria, P. Dupriez, C. A. Codemard, et al., "Single-Frequency, Single-Mode, Plane-Polarized Ytterbium-Doped Fiber Master Oscillator Power Amplifier Source with 264 W of Output Power," *Opt. Lett.* **30**(5):459–461, 2005.
197. D. V. Gapontsev, N. S. Platonov, M. Meleshkevich, O. Mishechkin, O. Shikurikin, S. Agger, P. Varming, and J. H. Poysen, "20 W Single-Frequency Fiber Laser Operating at 1.93 μm," in *CLEO 2007*, Baltimore, Md., 2007.
198. Y. Jeong, J. K. Sahu, D. B. S. Soh, C. A. Codemard, and J. Nilsson, "High-Power Tunable Single-Frequency Single-Mode Erbium:Ytterbium Codoped Large-Core Fiber Master-Oscillator Power Amplifier Source," *Opt. Lett.* **30**(22):2997–2999, 2005.

199. G. P. Agrawal, *Nonlinear Fiber Optics*, Academic, San Diego, Calif., 1995.
200. M. O'Connor, and F. DiTeodoro, "Fiber Lasers in Defense: Fibers, Components and System Design Considerations," in *DEPS SSDLTR Conference Short Course*, Los Angeles, Calif., 2007.
201. Y. Jeong, J. Sahu, M. Laroche, W. A. Clarkson, K. Furusawa, D. J. Richardson, and J. Nilsson, "120 W Q-Switched Cladding Pumped Yb Doped Fiber Laser," in *CLEO Europe*, 2003. Munich, Germany.
202. J. Limpert, N. Deguil-Robin, S. Petit, I. Manek-Hönniger, F. Salin, P. Rigal, C. Hönniger, and E. Mottay, "High Power Q-Switched Yb-Doped Photonic Crystal Fiber Laser Producing Sub-10 ns Pulses," *Appl. Phys. B: Lasers Opt.* **81**(1):19–21, 2005.
203. O. Schmidt, J. Rothhardt, F. Röser, S. Linke, T. Schreiber, K. Rademaker, J. Limpert, S. Ermeneux, P. Yvernault, F. Salin, and A. Tünnermann, "Millijoule Pulse Energy Q-Switched Short-Length Fiber Laser," *Opt. Lett.* **32**(11):1551–1553, 2007.
204. M. Eichhorn, and S. D. Jackson, "High-Pulse-Energy Actively Q-Switched Tm³⁺-Doped Silica 2 μ m Fiber Laser Pumped at 792 nm," *Opt. Lett.* **32**(19):2780–2782, 2007.
205. J. Y. Huang, S. C. Huang, H. L. Chang, K. W. Su, Y. F. Chen, and K. F. Huang, "Passive Q Switching of Er-Yb Fiber Laser with Semiconductor Saturable Absorber," *Opt. Exp.* **16**(5):3002–3007, 2008.
206. R. Paschotta, R. Häring, E. Gini, H. Melchior, U. Keller, H. L. Offerhaus, and D. J. Richardson, "Passively Q-Switched 0.1-mJ Fiber Laser System at 1.53 μ m," *Opt. Lett.* **24**(6):388–390, 1999.
207. F. Z. Qamar, and T. A. King, "Passive Q-Switching of the Tm-Silica Fibre Laser near 2 μ m by a Cr²⁺:ZnSe Saturable Absorber Crystal," *Opt. Commun.* **248**(4–6):501–508, 2005.
208. S. D. Jackson, "Passively Q-Switched Tm³⁺-Doped Silica Fiber Lasers," *Appl. Opt.* **46**(16):3311–3317, 2007.
209. M. V. Andrés, "Actively Q-Switched All-Fiber Lasers," *Laser Phys. Lett.* **5**(2):93–99, 2008.
210. Z. Yu, W. Margulis, O. Tarasenko, H. Knape, and P. Y. Fonjallaz, "Nanosecond Switching of Fiber Bragg Gratings," *Opt. Exp.* **15**(22):14948–14953, 2007.
211. H. Shuling, "Stable NS Pulses Generation from Cladding-Pumped Yb-Doped Fiber Laser," *Microwave and Opt. Technol. Lett.* **48**(12):2442–2444, 2006.
212. K. Kieu, and M. Mansuripur, "Active Q Switching of a Fiber Laser with a Microsphere Resonator," *Opt. Lett.* **31**(24):3568–3570, 2006.
213. W. Koehchner, *Solid-State Laser Engineering*, Springer, New York, 1999.
214. P. Myslinski, J. Chrostowski, J. A. Koningstein, and J. R. Simpson, "Self Mode-Locking in a Q-Switched Erbium-Doped Fiber Laser," *Appl. Opt.* **32**(3):286, 1993.
215. B. N. Upadhyaya, U. Chakravarty, A. Kuruvilla, K. Thyagarajan, M. R. Shenoy, and S. M. Oak, "Mechanisms of Generation of Multi-Peak and Mode-Locked Resembling Pulses in Q-Switched Yb-Doped Fiber Lasers," *Opt. Exp.* **15**(18):11576–11588, 2007.
216. M. -Y. Cheng, Y. -C. Chang, A. Galvanauskas, P. Mamidipudi, R. Changkakoti, and P. Gatchell, "High-Energy and High-Peak-Power Nanosecond Pulsegeneration with Beam Quality Control in 200- μ m Core Highly Multimode Yb-Doped Fiberamplifiers," *Opt. Lett.* **30**(4):358–360, 2005.
217. S. A. George, K. -C. Hou, K. Takenoshita, A. Galvanauskas, and M. C. Richardson, "13.5 nm EUV Generation from Tin-Doped Droplets Using a Fiber Laser," *Opt. Exp.* **15**(25):16348–16356, 2007.
218. A. Liem, J. Limpert, H. Zellmer, and A. Tünnermann, "100-W Single-Frequency Master-Oscillator Fiber Power Amplifier," *Opt. Lett.* **28**(17):1537–1539, 2003.
219. J. Limpert, S. Höfer, A. Liem, H. Zellmer, A. Tünnermann, S. Knoke, and H. Voelckel, "100-W Average-Power, High-Energy Nanosecond Fiber Amplifier," *Appl. Phys. B: Lasers Opt.* **75**(4):477–479, 2002.
220. P. E. Britton, H. L. Offerhaus, D. J. Richardson, P. G. R. Smith, G. W. Ross, and D. C. Hanna, "Parametric Oscillator Directly Pumped by a 1.55- μ m Erbium-Fiber Laser," *Opt. Lett.* **24**(14):975–977, 1999.
221. S. Desmoulins, and F. Di Teodoro, "Watt-Level, High-Repetition-Rate, Mid-Infrared Pulses Generated by Wavelength Conversion of an Eye-Safe Fiber Source," *Opt. Lett.* **32**(1):56–58, 2007.
222. M. Savage-Leuchs, E. Eisenberg, A. Liu, J. Henrie, and M. Bowers, "High-Pulse Energy Extraction with High Peak Power from Short-Pulse Eye Safe All-Fiber Laser System," in *Fiber Lasers III: Technology, Systems, and Applications*: SPIE, San Jose, Calif., 2006.
223. L. E. Nelson, D. J. Jones, K. Tamura, H. A. Haus, and E. P. Ippen, "Ultrashort Pulse Fiber Ring Lasers," *Appl. Phys. B* **65**(2):277–294, 1997.

224. K. Tamura, H. A. Haus, and E. P. Ippen, "Self-Starting Additive Pulse Mode-Locked Erbium Fibre Ring Laser," *Electron. Lett.* **28**(24):2226–2228, 1992.
225. K. Tamura, E. P. Ippen, H. A. Haus, and L. E. Nelson, "77-fs Pulse Generation from a Stretched-Pulse Mode-Locked All-Fiber Ring Laser," *Opt. Lett.* **18**(13):1080, 1993.
226. A. Chong, W. H. Renninger, and F. W. Wise, "All-Normal-Dispersion Femtosecond Fiber Laser with Pulse Energy above 20 nJ," *Opt. Lett.* **32**(16):2408–2410, 2007.
227. L. M. Zhao, D. Y. Tang, and J. Wu, "Gain-Guided Soliton in a Positive Group-Dispersion Fiber Laser," *Opt. Lett.* **31**(12):1788–1790, 2006.
228. F. Ö. Ilday, J. R. Buckley, W. G. Clark, and F. W. Wise, "Self-Similar Evolution of Parabolic Pulses in a Laser," *Phys. Rev. Lett.* **92**(21):213902, 2004.
229. K. Tamura, E. P. Ippen, and H. A. Haus, "Pulse Dynamics in Stretched-Pulse Fiber Lasers," *Appl. Phys. Lett.* **67**(2):158–160, 1995.
230. C. Lecaplain, C. Chédot, A. Hideur, B. Ortaç, and J. Limpert, "High-Power All-Normal-Dispersion Femtosecond Pulse Generation from a Yb-Doped Large-Mode-Area Microstructure Fiber Laser," *Opt. Lett.* **32**(18):2738–2740, 2007.
231. B. Ortaç, J. Limpert, and A. Tünnermann, "High-Energy Femtosecond Yb-Doped Fiber Laser Operating in the Anomalous Dispersion Regime," *Opt. Lett.* **32**(15):2149–2151, 2007.
232. B. Ortaç, O. Schmidt, T. Schreiber, J. Limpert, A. Tünnermann, and A. Hideur, "High-Energy Femtosecond Yb-Doped Dispersion Compensation Free Fiber Laser," *Opt. Exp.* **15**(17):10725–10732, 2007.
233. R. Herda, S. Kivist, O. G. Okhotnikov, A. Kosolapov, A. Levchenko, S. Semjonov, and E. Dianov, "Environmentally Stable Mode-Locked Fiber Laser with Dispersion Compensation by Index-Guided Photonic Crystal Fiber," *Photon. Technol. Lett. IEEE* **20**(3):217–219, 2008.
234. J. C. Diels and W. Rudolph, *Ultrashort Laser Pulse Phenomena*, Academic Press, New York, 1996.
235. A. E. Siegman, *Lasers*, University Science Books, New York, 1986.
236. U. Keller, K. J. Weingarten, F. X. Kartner, D. Kopf, B. Braun, I. Jung, R. Fluck, C. Honninger, N. Matuschek, and J. Aus Der Au, "Semiconductor Saturable Absorber Mirrors (SESAM) for Femtosecond to Nanosecond Pulse Generation in Solid-State Lasers," *IEEE J. Sel. Top. Quant. Electron.* **2**(3):435–453, 1996.
237. H. Kataura, Y. Kumazawa, Y. Maniwa, I. Umezu, S. Suzuki, Y. Ohtsuka, and Achiba, "Optical Properties of Single-Wall Carbon Nanotubes," *Synthetic Metals* **103**:2555–2558, 1999.
238. J. W. Nicholson, R. S. Windeler, and D. J. DiGiovanni, "Optically Driven Deposition of Single-Walled Carbon-Nanotube Saturable Absorbers on Optical Fiber End-Faces," *Opt. Exp.* **15**(15):9176–9183, 2007.
239. S. Y. Set, H. Yaguchi, Y. Tanaka, and M. Jablonski, "Laser Mode Locking Using a Saturable Absorber Incorporating Carbon Nanotubes," *J. Lightwave Technol.* **22**(1):51, 2004.
240. L. Vivien, P. Lancon, F. Hache, D. A. Riehl, and E. A. Anglaret, "Pulse Duration and Wavelength Effects on Optical Limiting Behaviour in Carbon Nanotube Suspensions," in *Lasers and Electro-Optics Europe, 2000. Conference Digest.* 2000.
241. L. P. Shen, W. P. Huang, G. X. Chen, and S. Jian, "Design and Optimization of Photonic Crystal Fibers for Broad-Band Dispersion Compensation," *Photon. Technol. Lett. IEEE* **15**(4):540–542, 2003.
242. D. Strickland, and G. Mourou, "Compression of Amplified Chirped Optical Pulses," *Opt. Commun.* **55**(6):447–449, 1985.
243. F. Röser, D. Schimpf, O. Schmidt, B. Ortaç, K. Rademaker, J. Limpert, and Tünnermann, "90 W Average Power 100 fJ Energy Femtosecond Fiber Chirped-Pulse Amplification System," *Opt. Lett.* **32**(15):2230–2232, 2007.
244. F. Röser, J. Rothhardt, B. Ortac, A. Liem, O. Schmidt, T. Schreiber, J. Limpert, and Tünnermann, "131 W 220 Fs Fiber Laser System," *Opt. Lett.* **30**(20):2754–2756, 2005.
245. F. Röser, T. Eidam, J. Rothhardt, O. Schmidt, D. N. Schimpf, J. Limpert, and A. Tünnermann, "Millijoule Pulse Energy High Repetition Rate Femtosecond Fiber Chirped-Pulse Amplification System," *Opt. Lett.* **32**(24):3495–3497, 2007.
246. A. Galvanauskas, "Mode-Scalable Fiber-Based Chirped Pulse Amplification Systems," *IEEE J. Sel. Top. Quant. Electron.* **7**(4):504–517, 2001.
247. K.-H. Liao, et al., "Generation of Hard X-Rays Using an Ultrafast Fiber Laser System," *Opt. Exp.* **15**(21):13942–13948, 2007.

248. G. Chang, et al., "50-W Chirped Volume Bragg Grating Based Fiber CPA at 1055 nm," in *SSDLTR*. Los Angeles, Calif., 2007.
249. K.-H. Liao, A. G. Mordovanakis, B. Hou, G. Chang, M. Rever, G. A. Mourou, J. Nees, and Galvanauskas, "Large-Aperture Chirped Volume Bragg Grating Based Fiber CPA System," *Opt. Exp.* **15**(8):4876–4882, 2007.
250. "Pioneering Ultrafast Fiber Laser Technology," available: www.imra.com, accessed on: Dec. 2008.

DO NOT DUPLICATE

PART

5

X-RAY AND
NEUTRON OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

SUBPART

5.1

INTRODUCTION AND APPLICATIONS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

AN INTRODUCTION TO X-RAY AND NEUTRON OPTICS

Carolyn MacDonald

*University at Albany
Albany, New York*

26.1 HISTORY

X rays have a century of history of medical and technological applications. Optics have come to play an important role in many capacities. One of Roentgen's earliest observations, shortly after his discovery of x rays in 1895,¹ was that while the rays were easily absorbed by some materials, they did not strongly refract. Standard optical lenses, which require weak absorption and strong refraction, were therefore not useful for manipulating the rays. New techniques were quickly developed. In 1914, van Laue was awarded the Nobel Prize for demonstrating the diffraction of x rays by crystals. By 1929, total reflection at grazing incidence had been used to deflect x rays.² The available optics for x rays still can be classified by those three phenomena: refraction, diffraction, and total reflection. Surprisingly, given that the physics governing these optics has been well known for nearly a century, there has been a recent dramatic increase in the availability, variety, and performance of x-ray optics for a wide range of applications.

An increasing interest in x-ray astronomy was one of the major forces for the development of x-ray optics in the latter half of the last century. Mirror systems similar to those developed for astronomy also proved useful for synchrotron beam lines. Just as x-ray tubes were an accidental offshoot of cathode ray research, synchrotron x-ray sources were originally a parasite of particle physics. The subsequent development of synchrotrons with increasing brightness and numbers of beam lines have created whole new arrays of x-ray tools and a consequent demand for an increasing array of optics. The rapid development of x-ray optics has also been symbiotic with the development of detectors and of compact sources. Detectors developed for particle physics, medicine, and crystallography have found applications across fields. Similarly, the increasing capability of x-ray systems has stimulated the development of new science with evergrowing requirements for intensity, coherence, and spatial and energy resolution. X-ray diffraction and fluorescence were early tools of the rapid development of materials science after World War II, but have been greatly advanced to meet the demands of the shrinking feature sizes and allowed defect levels in semiconductors. X-ray diffraction, especially the development of dedicated synchrotron beam lines, has also been stimulated by the growing demands for rapid protein crystallography for biophysics and pharmaceutical development.

The new abundance of x-ray optics, sources, and detectors requires a fresh look at the problem of optimizing a wide range of x-ray and neutron applications. The development of x-ray technology has also advanced neutron science because a number of the optics and detectors are either applicable to neutrons or have inspired the development of neutron technology.

One question that arises immediately in a discussion of x-ray phenomena is precisely what spectral range is included in the term. Usage varies considerably by discipline, but for the purposes of this volume, the x-ray spectrum is taken to be roughly from 1 to 100 keV in photon energy (1.24 to 0.0124-nm wavelength). The range is extended down into the hard EUV to include some microscopy and astronomical optics, and upward to include nuclear medicine.

26.2 X-RAY INTERACTION WITH MATTER

X rays are applicable to such a wide variety of areas because they are penetrating but interacting, have wavelengths on the order of atomic spacings, and have energies on the order of core electronic energy levels for atoms.

X-Ray Production

X rays are produced primarily by the acceleration of charged particles, the knock out of core electrons, or black body and characteristic emission from very hot sources such as laser-generated plasmas or astronomical objects. The production of x rays by accelerated charges includes incoherent emission such as bremsstrahlung radiation in tube sources and coherent emission by synchrotron undulators or free electron lasers. Highly coherent emission can also be created by pumping transitions between levels in ionic x-ray lasers.

The creation of x rays by the knock out of core electrons is the mechanism for the production of the characteristic lines in the spectra from conventional x-ray tube sources. The incoming electron knocks out a core electron, creating a vacancy, which is quickly filled by an electron dropping down from an outer shell. The energy difference between the outer shell energy level and the core energy level is emitted in the form of an x-ray photon. This is also the origin of the characteristic lines used to identify elemental composition in x-ray fluorescence, described in Chap. 29, and for x-ray spectroscopy, described in Chap. 30. X-ray sources are described in Subpart 5.4 of this section, in Chaps. 54 to 59. Source coherence, and coherence requirements, are discussed in Chap. 27.

Refraction

In the x-ray regime, the real part of the index of refraction of a solid can be simply approximated by³

$$n_r = \text{Re}\{\sqrt{\kappa}\} = \text{Re}\left\{\sqrt{\frac{\epsilon}{\epsilon_0}}\right\} \cong \sqrt{1 - \frac{\omega_p^2}{\omega^2}} \cong 1 - \delta \quad (1)$$

where n is the index of refraction, ϵ is the dielectric constant of the solid, ϵ_0 is the vacuum dielectric constant, κ is their ratio, ω is the photon frequency, and ω_p is the plasma frequency of the material. The plasma frequency, which typically corresponds to tens of electron volts, is given by

$$\omega_p^2 = \frac{Ne^2}{m\epsilon_0} \quad (2)$$

where N is the electron density of the material, and e and m are the charge and mass of the electron. For x rays, the relevant electron density is the total density, including core electrons. Thus changes in x-ray optical properties cannot be accomplished by changes in the electronic levels, in the manner in which optical properties of materials can be manipulated for visible light. The x-ray optical properties

are determined by the density and atomic numbers of the elemental constituents. Tables of the properties of most of the elements are presented in Chap. 36. Because the plasma frequency is very much less than the photon frequency, the index of refraction is slightly less than one.

Absorption and Scattering

The imaginary part of the index of refraction for x rays arises from photoelectric absorption. This absorption is largest for low energy x rays and has peaks at energies resonant with core ionization energies of the atom. X rays can also be deflected out of the beam by incoherent or coherent scattering from electrons. Coherent scattering is responsible for the decrement to the real part of the index of refraction given in Eq. (1). Scattering from nearly free electrons is called *Thompson scattering*, and from tightly bound electrons, *Rayleigh scattering*. The constructive interference of coherent scattering from arrays of atoms constitutes *diffraction*. *Incoherent*, or Compton, scattering occurs when the incident photon imparts energy and momentum to the electron. Compton scattering becomes increasingly important for high-energy photons and thick transparent media.

26.3 OPTICS CHOICES

Slits and Pin Holes

Almost all x-ray systems, whether or not they use more complex optics, contain slits or apertures. Some aperture systems have considerable technological development. Most small sample diffraction systems employ long collimators designed to reduce the background noise from scattered direct beam reaching the detector. For θ - 2θ measurements, Soller slits, arrays of flat metal plates arranged parallel to the beam direction, are often placed after the sample to further reduce the background. Diffraction applications are described in Chap. 28.

Lead hole collimators and pinholes specifically designed for high photon energies are employed for nuclear medicine, with pinhole sizes down to tens of microns. Nuclear medicine is discussed in Chap. 32.

Refractive Optics

One consequence of Eq. (1) is that the index of refraction of all materials is very close to unity in the x-ray regime. Thus Snell's law implies that there is very little refraction at the interface,

$$(1) \sin\left(\frac{\pi}{2} - \theta_1\right) = n \sin\left(\frac{\pi}{2} - \theta_2\right) \quad (3)$$

where the first medium is vacuum and the second medium has index n , as shown in Fig. 1. In x-ray applications the angles θ are measured from the surface, not the normal to the surface. If n is very close to one, θ_1 is very close to θ_2 . Thus, refractive optics are more difficult to achieve in the x-ray regime. The lens maker's equation,

$$f = \frac{R/2}{n-1} \quad (4)$$

gives the relationship between the focal length of a lens f and the radius of curvature of the lens R . The symmetric case is given. Because n is very slightly less than one, the focal length produced even by very small negative radii (convex) lenses is rather long. This can be overcome by using large

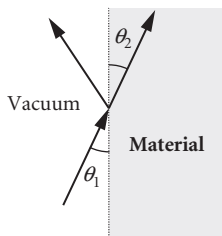


FIGURE 1 Refraction at a vacuum-to-material interface.

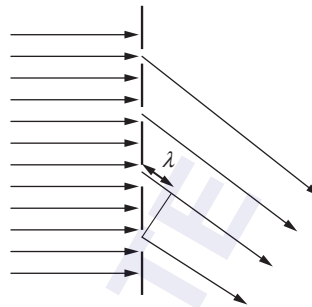


FIGURE 2 Schematic of constructive interference from a transmission grating. λ is the wavelength of the radiation. The extra path length indicated must be an integral multiple of λ . Most real gratings are used in reflection mode.

numbers of surfaces in a compound lens. Because the radii must still be small, the aperture of the lens cannot be large, and refractive optics are generally better suited to narrow synchrotron beams than isotropic point sources. Refractive optics is described in Chap. 37.

Diffractive and Interference Optics

The coherent addition of radiation from multiple surfaces or apertures can only occur for a very narrow wavelength bandwidth. Thus, diffractive and interference optics such as gratings, crystals, zone plates, multilayers, and Laue lenses, described in Chaps. 38 to 43, respectively, are all wavelength selective. Such optics are often used as monochromator to select a particular wavelength range from a white beam source.

To achieve constructive interference, the spacing of a grating must be arranged so that the radiation from successive sources is in phase, as shown in Fig. 2. The circular apertures of zone plates are similar to the openings in a transmission grating. The superposition of radiation from the circular apertures results in focusing of the radiation to points on the axis.

Diffractive optics operate on the same principal, coherent superposition of many rays.⁴ The most common diffractive optic is the crystal. Arranging the beam angle to the plane, θ , and plane spacing d as shown in Fig. 3, so that the rays reflecting from successive planes are in phase (and ignoring refraction) yields the familiar Bragg's law,

$$n\lambda = 2d \sin \theta \quad (5)$$

where n is an integer and λ is the wavelength of the x ray. Bragg's law cannot be satisfied for wavelengths greater than twice the plane spacing, so crystal optics are limited to low wavelength, high energy, x rays. Multilayers are "artificial crystals" of alternating layers of materials. The spacing, which is the thickness of a layer pair, replaces d in Eq. (5), and can be much larger than crystalline plane spacings. Multilayers are therefore effective for lower energy x rays. Diffraction can also be used in transmission, or Laue, mode, as shown in Fig. 4.

Reflective Optics

Using Snell's law, Eq. (3), the angle inside a material is smaller than the incident angle in vacuum. For incident angles less than a critical angle θ_c , no wave can be supported in the medium and the

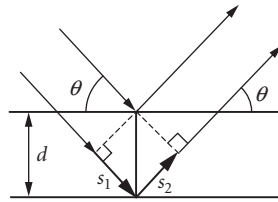


FIGURE 3 Pictorial representation of Bragg diffraction from planes with spacing d . The extra path length, $(s_1 + s_2)$, must be an integral multiple of λ .

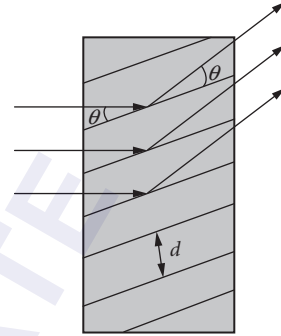


FIGURE 4 Laue transmission diffraction.

incident ray is totally externally reflected. The critical angle, the largest incident angle for total reflection, is given by

$$\sin\left(\frac{\pi}{2} - \theta_c\right) = n \sin\left(\frac{\pi}{2}\right) \quad (6)$$

Thus, using small angle approximations and Eq. (1)

$$\theta_c \approx \frac{\omega_p}{\omega} \quad (7)$$

Because the plasma frequency is very much less than the photon frequency, total external reflection occurs for very small grazing incidence angles. This phenomenon is used extensively for single reflection mirrors for synchrotrons, x-ray microscopes, and x-ray telescopes. Mirrors are described in Chaps. 44 to 47 and 51. To increase the angle of incidence, mirrors are often coated with metals, or with multilayers (although, because the multilayer depends on interference, the optic will no longer have a broadband response).

Arrays of mirrors, such as multifoil or pore optics, are described in Chaps. 48 and 49. Glass capillaries, described in Chap. 52, can be used as single bounce-shaped mirrors, or as multiple-bounce light pipes transport or focus the beam. Multiple reflections are used to transport x rays in polycapillary arrays, described in Chap. 53. Because a mirror is often the first optic in a synchrotron beam line, heat load issues are significant. For a number of synchrotron and other high-flux applications, radiation hardness and thermal stability are important considerations. A large body of experience has been developed for high-heat-load synchrotron mirrors and crystals.^{5,6} Many optics such as zone plates, microscopy objectives, and glass capillary tubes are routinely used in synchrotron beam lines and are stable over acceptable flux ranges. Adaptive optics used to mitigate the effects of thermal changes are described in Chap. 50.

26.4 FOCUSING AND COLLIMATION

Optics Comparisons

A variety of x-ray optics choices exist for collimating or focusing x-ray beams. The best choice of the optic depends to a large extent on the geometry of the sample to be measured, the information desired from the measurement, and the source geometry and power. No one optic can provide

the best resolution, highest intensity, easiest alignment, and shortest data acquisition time for all samples. A true comparison of two optics for a particular application requires careful analysis of the sample and measurement requirements, and adjustment of the source and optic design for the application. No global comparison of all optics for all applications is possible.

Four optics commonly used for collimation or focusing from x-ray tubes are bent crystals, multilayers, nested cones, and polycapillary optics. Bent crystals collect radiation from a point source and diffract it into a nearly monochromatic collimated beam. Doubly bent crystals or two singly bent crystals are required for two-dimensional collimation or focusing. Multilayer optics also work by diffraction, although in this case from the periodicity of the artificial compositional variation imposed in the multilayer. The beam is less monochromatic than for bent crystals. "Supermirrors," or "Goebel mirrors," are multilayers with graded or irregular spacing, and have wider energy bandwidths than multilayers.

Mirrors are broad band optics, but because curved mirrors are single-bounce optics, their maximum angular deflection is limited to the critical angle for their metallic coating, which can be about 10 mrad at 8 keV. The total capture angle is then determined by the length of the mirror. The output divergence is determined by the length of the optic and the source size. One solution to increase the capture angle is to "nest" multiple optics. Nested parabolic mirrors have smaller output divergences than nested cones. Polycapillary optics are also array optics, containing hundreds of thousands glass tubes. Because the focal spot is produced by overlap, it cannot be smaller than the channel size.

Comparison of focusing optics for diffraction also requires a detailed analysis of the effect of the convergence angle on the diffracted signal intensity and so is very sample and measurement dependent. Decreasing the angle of convergence onto the sample improves the resolution and decreases the signal to noise ratio, but also decreases the diffracted signal intensity relative to a large angle.⁷ Conventional practice is to use a convergence angle less than the mosaicity of the sample. It is necessary that the beam cross-section at the sample be larger than the sample to avoid the difficulty of correcting for intensity variations with sample angle.

Focusing for Spatial Resolution

Diffraction effects limit the resolution of visible light optical systems to within an order of magnitude of the wavelength of the light. Thus, in principle, x-ray microscopy systems could have resolution many orders of magnitude better than optical light systems. Electrons also have very small wavelength, and electron microscopes have extremely high resolution. However, electrons are charged particles and necessarily have low penetration lengths into materials. X-ray microscopy is capable of very high resolution imaging of relatively thick objects, including wet samples. In practice, x-ray microscopy covers a range of several orders of magnitude in wavelength and spot size. Very high resolution is commonly obtained with Schwarzschild objectives,⁸ or zone plates.⁹ Because Schwarzschild objectives are used in normal incidence, they are essentially limited to the EUV region. Synchrotron sources are required to provide adequate flux. The resolution of zone plates is determined by the width of the outermost zone. Zone plates are easiest to make for soft x rays, where the thickness required to absorb the beam is small. However, zone plates with high aspect ratio have been demonstrated for hard x rays. Because the diameters of imaging zone plates are small, and the efficiencies are typically 10 percent for amplitude zone plates and 40 percent for phase zone plates, synchrotron sources are required. Very small spot sizes have also been demonstrated with multilayer Laue lenses and with refractive optics.

Single capillary tubes can also have output spot sizes on the order of 100 nm or less, and thus can be used for scanning microscopy or microanalysis. Capillary tubes with outputs this small also require synchrotron sources.

For larger spot sizes, a wider variety of sources and optics are applicable. Microscopes have been developed for laser-plasma sources with both Wolter¹⁰ and bent crystal optics,¹¹ with resolutions of a few microns. These optics, and also capillary and polycapillary optics and nested mirrors, with or without multilayer coatings, can produce spot sizes of a few tens of microns with laboratory tube sources. Optics designed to collect over large solid angles, such as bent crystals, graded multilayers or polycapillary optics, will produce the highest intensities. Bent crystals will yield monochromatic radiation; polycapillary optics can be used to produce higher intensity, but broader band radiation. Polycapillary

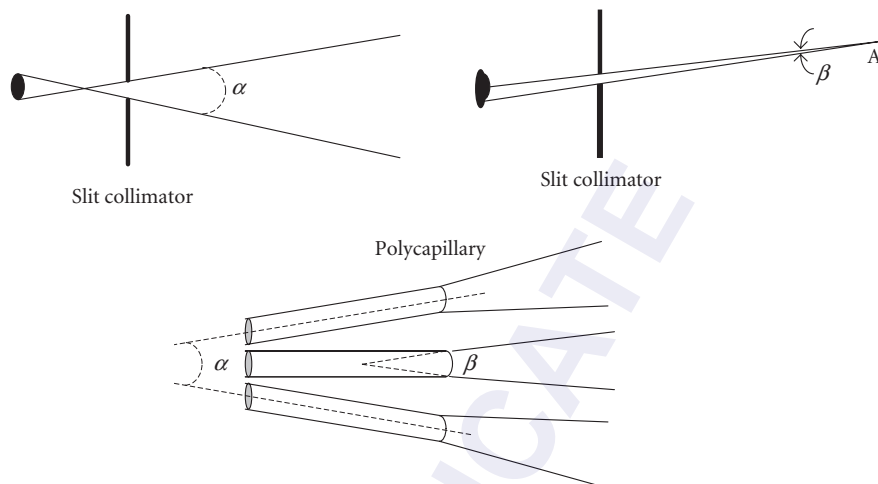


FIGURE 5 For either a slit collimator (top left and right) or an optic such as a polycapillary optic (bottom), the local divergence β seen at a point A is different from the angle subtended by the beam α .

optics have spot sizes no smaller than tens of microns, independent of source size. Mirrors and crystals will have smaller focal spot sizes for smaller sources and larger spot sizes for larger sources. Refractive optics are true imaging optics, with the potential for very small spots. Clearly, the optimal optic depends on the details of the measurement requirements and the sources available.

Collimation

As shown in Fig. 5, the output from a collimating optic has both global divergence α and local divergence β . Even if the global divergence is made very small by employing an optic, the local divergence is usually not zero. For grazing incidence reflection optics, the local divergence is generally given by the critical angle for reflection and can be increased by profile errors in the optic. The degree of collimation achieved by the optic is limited not only by technology, but by the thermodynamic constraint that the beam brightness cannot be increased by any optic.¹² Liouville's theorem states that increasing the density of states in phase space is a violation of the second law of thermodynamics. This implies that the six dimensional real space/momentum space volume occupied by the photons cannot be decreased. More simply, the angle-area product, that is, the cross-sectional area of the beam multiplied by the divergence of the beam, cannot be decreased without losing photons. The beam can never be brighter than the source. An idealized point source of x rays occupies zero area. X rays from such a source could theoretically be perfectly collimated into any chosen cross section. However, a real source has a finite size and so limits the degree of possible collimation. Small spot sources are required to produce bright, well-collimated beams.

26.5 REFERENCES

1. W. C. Röntgen, "On a New Form of Radiation," *Nature* 53:274–276, January 23, 1896, English translation from *Sitzungsberichte der Würzburger Physik-med. Gesellschaft*, 1895.
2. Werner Ehrenberg and Felix Jentsch, "Über die Auslösung von Photoelektronen durch Röntgenstrahlen aus Metallspiegeln an der Grenze der Totalreflexion," *Zeitschrift für Physik* 54(3, 4):227–235, March, 1929.

3. J. D. Jackson, *Classical Electrodynamics*, John Wiley & Sons, New York, pp. 227, 1962.
4. B. D. Cullity and B. D. Cullity, *Elements of X-Ray Diffraction*, Addison Wesley Longman, Reading, Mass., 1978.
5. A. Khounsary, *Advances in Mirror Technology for Synchrotron X-Ray and Laser Applications*, SPIE, Bellingham, Wash., vol. 3447, 1998.
6. A. T. Macrander and A. M. Khounsary (eds.), *High Heat Flux and Synchrotron Radiation Beamlines (Proceedings Volume)* SPIE, Bellingham, Wash., vol. 3151, 11 December, 1997, ISBN: 9780819425737.
7. U. W. Arndt, "Focusing Optics for Laboratory Sources in X-Ray Crystallography," *J. Appl. Cryst.* **23**:161–168, 1990.
8. F. Cerrina, "The Schwarzschild Objective," in *Handbook of Optics*, 3d ed., vol. V, M. Bass (ed.), McGraw-Hill, New York, 2009.
9. A. Michette, "Zone Plates," in *Handbook of Optics*, 3d ed., vol. V, M. Bass (ed.), McGraw-Hill, New York, 2009.
10. P. Trousses, P. Munsch, and J. J. Ferme, "Microfocusing between 1 and 5 keV with Wolter Type Optic," in *X-Ray Optics Design, Performance and Applications*, A. M. Khounsary, A. K. Freund, T. Ishikawa, G. Srajer, J. Lang, (eds.), SPIE, Bellingham, Wash., vol. 3773, pp. 60–69, 1999.
11. T. A. Pikuz, A. Y. Faenov, M. Fraenkel, et al., "Large-Field High Resolution X-Ray Monochromatic Microscope, Based on Spherical Crystal and High Repetition-Rate Laser-Produced Plasmas," in *EUV, X-Ray and Neutron Optics and Sources*, C. A. MacDonald, K. A. Goldberg, J. R. Maldonado, H. H. Chen-Mayer, S. P. Vernon (eds.), SPIE, Bellingham, Wash., vol. 3767, pp. 67–78, 1999.
12. D. L. Goodstein, *States of Matter*, Prentice-Hall, Englewood Cliffs, N.J., 1975.

COHERENT X-RAY OPTICS AND MICROSCOPY

Qun Shen

*National Synchrotron Light Source II
Brookhaven National Laboratory
Upton, New York*

27.1 GLOSSARY

$F(x, y)$	diffracted wave field amplitude on the detector image plane (x, y)
$q(X, Y)$	transmission function through a thin object
r	length of the position vector from point (X, Y) on the object plane to point (x, y) on the detector image plane
λ	x-ray wavelength
k	$2\pi/\lambda$ is the wave number
$\bar{q}(X, Y)$	distorted object with Fresnel zone phase factors embedded in the original object
$N_z = a^2/(4\lambda z)$	number of Fresnel zones on the object when looking back from the image plane

Traditionally the development of x-ray optics is based on ray-tracing wavefield propagation in geometric optics. This is true especially in the hard x-ray regime with energy in the multiple keV regime. With the availability of partially coherent x-ray sources such as those at modern synchrotrons, this situation has changed completely. Very often, some “artifacts” or features beyond ray-tracing can be observed in experiments due to the interference or phase effects from substantial spatial coherence of the synchrotron x-ray source. In this chapter, a brief outline is presented of a simple version of the wave propagation theory that can be used to evaluate coherent propagation of x-ray waves to take into account these effects. Based on optical reciprocity theorem, this type of coherent propagation is also required when evaluating x-ray focusing optics with the intention to achieve diffraction-limited performance.

In addition to coherent x-ray optics, there are substantial interests in the scientific community to use x-rays for imaging microscopic structures based on coherent wave propagation. For example, the success of structural science today is largely based on x-ray diffraction from crystalline materials. However, not all materials of interest are in crystalline forms; examples include the majority of membrane proteins and larger multidomain macromolecular assemblies, as well as many nanostructure specimens at their functioning levels. For these noncrystalline specimens, imaging at high spatial resolution offers an alternative, or the only alternative, to obtain any information on their internal structures. This topic is covered in the second part of this chapter.

27.2 INTRODUCTION

The spatial coherence length is defined as the transverse distance across the beam over which two parts of the beam have a fixed phase relationship. Using classical optics, the coherence length is $L_{\text{coherence}} = (\lambda/\beta)$ where λ is the wavelength of the radiation, and β is the local divergence, the angle subtended by the source. For typical laboratory sources the coherence length is only a few microns at several meters. However, for modern synchrotron sources β is much smaller and the coherence length is large enough to encompass optics such as zone plates which require coherent superposition, or entire samples for phase imaging.

27.3 FRESNEL WAVE PROPAGATION

In principle, imaging and diffraction or scattering are two optical regimes that are intrinsically interrelated based on Fresnel diffraction for wave propagation, which, under the first-order Born approximation,^{1,2} is

$$F(x, y) = \frac{i}{\lambda} \iint q(X, Y) \frac{e^{-ikr}}{r} dXdY \quad (1)$$

where $F(x, y)$ is the diffracted wave field amplitude, $q(X, Y)$ is the transmission function through a thin object, $r = [z^2 + (x - X)^2 + (y - Y)^2]^{1/2}$ is the length of the position vector from point (X, Y) on the object plane to point (x, y) on the detector image plane, λ is the x-ray wavelength, and $k = 2\pi/\lambda$ is the wave number.

Although widely used in optical and electron diffraction and microscopy,¹ the concept of Fresnel diffraction Eq. (1) has only recently been recognized in the broader x-ray diffraction community where traditionally far-field diffraction plus conventional radiography dominated the x-ray research field for the past century. This is because an essential ingredient for Fresnel-diffraction-based wave propagation is a substantial degree of transverse coherence in an x-ray beam, which had not been easily available until recent advances in partially coherent synchrotron and laboratory-based sources.

27.4 UNIFIED APPROACH FOR NEAR- AND FAR-FIELD DIFFRACTION

Coherent wave field propagation based on Fresnel diffraction Eq. (1) is usually categorized into two regimes: the near-field Fresnel or in-line holography regime, and the far-field Fraunhofer regime. A unified method for the evaluation of wave-field propagation in both regimes (Fig. 1) has been developed using the concept of *distorted object* in Fresnel Eq. (1), which can be applied both to Fraunhofer and Fresnel diffractions.

To introduce this method, we expand in Eq. (1), $r = [z^2 + (x - X)^2 + (y - Y)^2]^{1/2} \approx z + [(x - X)^2 + (y - Y)^2]/2z$ so that Eq. (1) becomes

$$F(x, y) = \frac{i e^{-ikz}}{\lambda z} \iint q(X, Y) e^{-ik \frac{(x-X)^2 + (y-Y)^2}{2z}} dXdY$$

Further expanding the terms in the exponential results in

$$F(x, y) = \frac{i e^{-ikR}}{\lambda R} \iint q(X, Y) e^{-\frac{i\pi}{\lambda z}(X^2+Y^2)} e^{-\frac{i2\pi}{\lambda z}(xX+yY)} dXdY$$

where $R = (x^2 + y^2 + z^2)^{1/2}$. We now define a new *distorted object* $\bar{q}(X, Y)$ as follows

$$\bar{q}(X, Y) \equiv q(X, Y) e^{-\frac{i\pi}{\lambda z}(X^2+Y^2)} \quad (2)$$

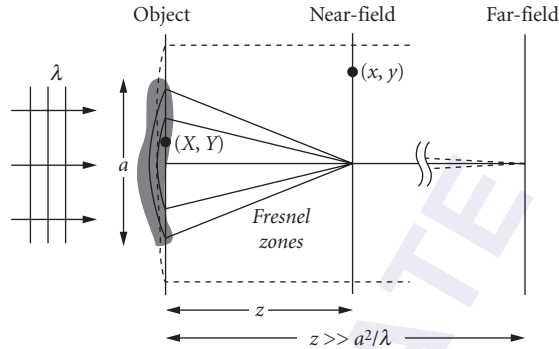


FIGURE 1 Schematic illustration of coherent x-ray wave propagation with a distorted object approach both for near-field Fresnel diffraction, where an object extends into multiple Fresnel zones (solid lines), and for far-field Fraunhofer diffraction, where an object occupies only the center of the first Fresnel zone (dashed lines). (See also color insert.)

and the scattered wave field $F(x, y)$ can then be expressed by a direct Fourier transform of this distorted object

$$F(x, y) = \frac{i e^{-ikR}}{\lambda R} \iint \bar{q}(X, Y) e^{-\frac{ik}{z}(xX+yY)} dXdY \quad (3)$$

Eq. (3) clearly shows that by embedding Fresnel zone construction into the distorted object, Eq. (2), a near-field diffraction pattern can be simply evaluated by a Fourier transform just like in the far-field approximation, with a momentum transfer $(Q_x, Q_y) = (kx/z, ky/z)$. Furthermore, it reduces to the familiar far-field result when $z \gg a^2/(4\lambda)$, where a is the transverse size of the object, since the extra Fresnel phase factor in Eq. (2) can be then approximated to unity. In general, the number of Fresnel phase zones of width π depends on distance z and is given by $N_z = a^2/(4\lambda z)$. Therefore, Eq. (3) can be used both in the near-field and in the far-field regimes, and this traditional but somewhat artificial partition of these two regimes is easily eliminated. Figure 2 shows some examples of calculated diffraction patterns at different detector to specimen distances.

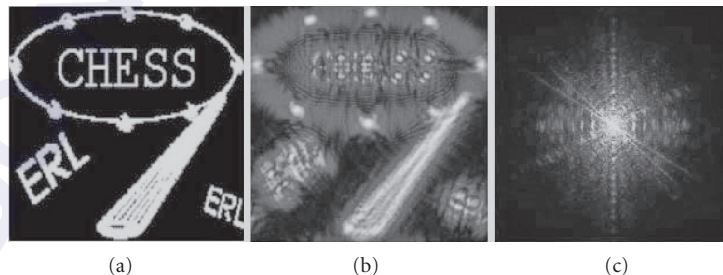


FIGURE 2 Simulated diffraction amplitudes $|F(x, y)|$, of an amplitude object (a) of $10 \mu\text{m} \times 10 \mu\text{m}$, with $\lambda = 1 \text{ \AA}$ x rays, at image-to-object distance (b) $z = 2 \text{ mm}$ and (c) $z = \infty$, using the unified distorted object approach Eq. (3) with $N_z = 500$ zones in (b) and $N_z = 0$ in (c). Notice that the diffraction pattern changes from noncentrosymmetric in the near-field (b) to centrosymmetric in the far-field (c). (See also color insert.)

27.5 COHERENT DIFFRACTION MICROSCOPY

It has been shown in recent years³ that an oversampled continuous coherent diffraction pattern from a nonperiodic object can be phased directly based on real space and reciprocal space constraints using an iterative phasing technique originally developed in optics.^{4,5} The oversampling condition requires a diffraction pattern be measured in reciprocal space at a Fourier interval finer than the Nyquist frequency used in all discrete fast Fourier transforms. Once such an oversampled diffraction pattern is obtained, as shown in Fig. 3, the iterative phasing method starts with a random set of phases for diffraction amplitudes, and Fourier transforms back and forth between diffraction amplitudes in reciprocal space and density in real space. In each iteration, the real space density is confined to within the finite specimen size and the square of diffraction amplitudes in reciprocal space is made equal to the experimentally measured intensities. This iterative procedure has proved to be a powerful phasing method for coherent diffraction imaging of nonperiodic specimens as a form of lensless x-ray microscopy.

One of the applications of the distorted object approach is that it extends the Fourier transform-based iterative phasing technique that works well in the far-field coherent diffraction imaging, into the regime of phasing near-field Fresnel diffraction or holographic images.⁶ Because the distorted object $\bar{q}(X, Y)$ differs from the original object $q(X, Y)$ by only a phase factor, which is known once the origin on the object is chosen, all real-space constraints applicable on $q(X, Y)$ can be transferred onto $\bar{q}(X, Y)$ in a straightforward fashion. In fact, most existing iterative phasing programs may be easily modified to accommodate the distorting phase factor in Eq. (2). A similar technique developed by Nugent et al.⁷ makes use of a curved wave illumination from a focusing x-ray optic in coherent diffraction imaging experiments and has demonstrated that the iterative phasing algorithm may converge much faster with a curved-beam illumination.

A significant recent development in coherent diffraction microscopy is the introduction of a scanning probe so that the coherent diffraction method can be applied to extended specimens.⁸ It has

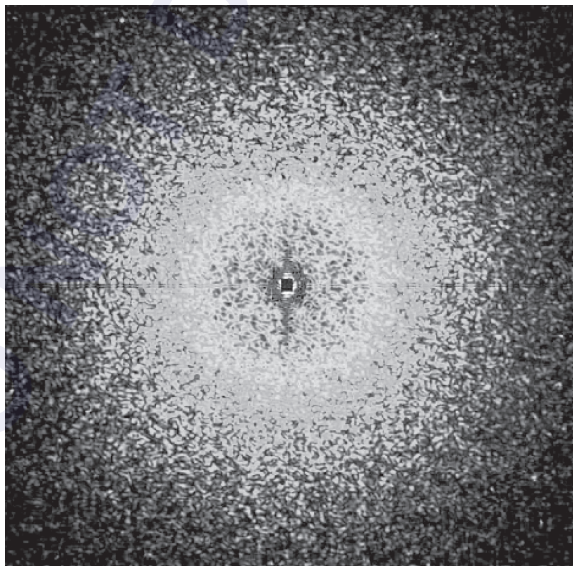


FIGURE 3 Example of a coherent x-ray diffraction pattern from a gold nanofoam specimen of $\sim 2 \mu\text{m}$ in size, using 7.35-keV coherent x rays. The corner of the image corresponds to $\sim 8 \text{ nm}$ spatial frequency. (See also color insert.)

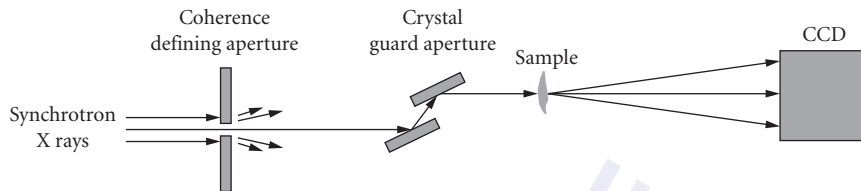


FIGURE 4 Concept of a perfect-crystal guard aperture in coherent diffraction imaging experiments for the purpose of eliminating unwanted parasitic scattering background in order to achieve high signal-to-noise in a diffraction pattern. (See also color insert.)

become apparent that the a combination of scanning x-ray microscopy with coherent diffraction may be the ultimate tool⁹ that scientists will be using in the coming years to image high resolution structures on nonperiodic specimens.

27.6 COHERENCE PRESERVATION IN X-RAY OPTICS

Coherent x-ray wavefield propagation has become an important consideration in many aspects of x-ray optics developments. There are many examples already published in the literature. For example, in order to evaluate the ultimate performance of x-ray focusing optics, it is essential to employ a coherent wave propagation theory from the x-ray optic to the focal spot which ultimately is diffraction limited.¹⁰

Another example where coherent wave propagation is needed is to preserve well-defined wavefronts through an x-ray optical system, which is often referred to as coherence preservation. For instance, one crucial issue in coherent x-ray diffraction imaging is how to increase the signal-to-noise ratio when measuring relatively weak diffraction intensities from a nonperiodic object. Based on coherent wave propagation, a crystal guard aperture concept has been developed¹¹ which makes use of a pair of multiple-bounce crystal optics to eliminate unwanted parasitic scattering background from the upstream coherence defining aperture (see Fig. 4). Recent experimental observation and theoretical analysis confirm the effectiveness of the crystal guard aperture method with coherence-preserved wave propagation through the crystal guard aperture and dramatically reduced scattering background in coherent x-ray diffraction images.¹¹

In summary, the development of coherent x-ray optics and x-ray analysis has become increasingly more important as the state-of-the-art x-ray sources and x-ray microscopic tools are becoming more readily available. It is expected that this field will continue to grow rapidly in order to satisfy the strong scientific interests in the community.

27.7 REFERENCES

1. J. M. Cowley, *Diffraction Physics*, 2nd ed., Elsevier Science Publisher, New York, 1990.
2. F. van der Veen and F. Pfeiffer, *J. Phys.: Condens. Matter* **16**:5003 (2004).
3. Original idea was proposed by D. Sayre, "Prospects for Long-Wavelength X-Ray Microscopy and Diffraction" in *Imaging Processes and Coherence in Physics*, M. Schlenker, M. Fink, J. P. Goedgebuer, C. Malgrange, J. C. Vienot, and R. H. Wade (eds.), *Springer Lect. Notes Phys.* **112**:229–35 (1980), Berlin: Springer. For a recent review, see J. Miao, et al., *Annu. Rev. Phys. Chem.* **59**:387–409 (2008).
4. J. R. Fienup, *Appl. Opt.* **21**:2758 (1982).

5. R. W. Gershberg and W. O. Saxton, *Optik* 25:237 (1972).
6. X. Xiao and Q. Shen, *Phys. Rev. B* 72:033101 (2005).
7. G. J. Williams, H. M. Quiney, B. B. Dhal, et al., *Phys. Rev. Lett.* 97:025506 (2006).
8. J. M. Rodenburg, A. C. Hurst, A. G. Cullis, et al., *Phys. Rev. Lett.* 98:034801 (2007).
9. P. Thibault et al., *Science* 321:379–382 (2008).
10. H. Yan et al., *Phys. Rev. B* 76:115438 (2007).
11. X. Xiao et al., *Opt. Lett.* 31:3194 (2006).

DO NOT DUPLICATE

REQUIREMENTS FOR X-RAY DIFFRACTION

Scott T. Misture

*Kazuo Inamori School of Engineering
Alfred University
Alfred, New York*

28.1 INTRODUCTION

Many analytical tools involve the use of x-rays in both laboratory synchrotron settings. X-ray imaging is a familiar technique, with x-ray diffraction (XRD) and x-ray fluorescence (XRF) nearly ubiquitous in the materials analysis laboratory.^{1–3} A long list of additional tools incorporate x-ray optics, especially at synchrotron sources where a continuous range of x-ray wavelengths is readily accessible.

The optical components used in x-ray analysis range from simple slits and collimators to diffractive elements including crystals and multilayers to reflective elements including capillaries and mirrors (see Chaps. 39, 41, 44, 52, and 53). Regardless of the specific application, a description of x-ray optics can be divided into three components:

- Definition of the beam path
- Definition of the beam divergence
- Definition of beam conditioning, or in other words defining the energy spectrum transmitted to the sample or detector by the various optical components under given conditions

Regardless of the quantity measured—intensity, energy, or angle—the interplay of these three parameters is critical to understanding the instrument response. In order to understand the use of optical components we shall begin by describing the simplest of systems which involves slits only.

28.2 SLITS

Using simple apertures, for example, slits, pinholes, and parallel plate collimators, is often sufficient to obtain high-quality data. The most common example of such a system is the powder diffractometer in Bragg-Brentano geometry as shown in Fig. 1. Figure 1 demonstrates that the divergent beam is achieved using a system of slits to control the divergence in the direction normal to the axis of rotation of the goniometer or the *equatorial* divergence. Note that slits are generally used in a pairs (see Fig. 1), with a primary slit and an antiscatter slit designed to block any photons scattered from the edges of the primary slit.

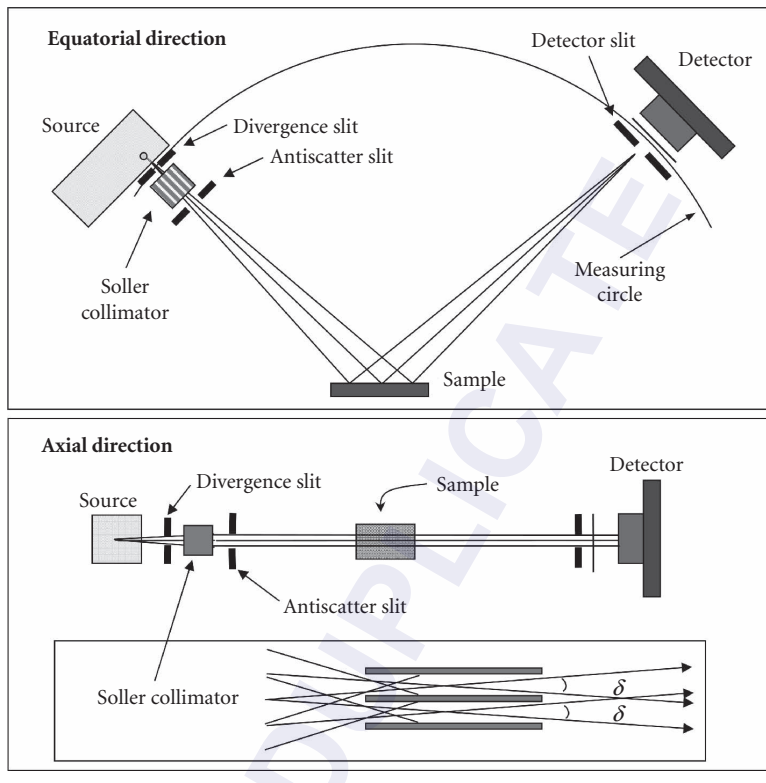


FIGURE 1 Schematic views of the Bragg-Brentano powder diffractometer.

Also shown in Fig. 1 is control of the beam divergence along the axis of the goniometer (*axial* divergence) using parallel plate collimators. The construction of the parallel plate collimator, also called a Soller collimator, includes closely spaced plates that limit the angular range of photons transmitted through the device (Fig. 1). Soller collimators are often used to limit the axial divergence to a few degrees or less, as shown in Fig. 1, but can also be used to achieve “parallel beam” conditions. By “parallel” we mean divergence ranging from the practical limit for a collimator of $\sim 0.05^\circ$ to $\sim 0.2^\circ$.

As an example, Fig. 2 shows the construction of simple parallel beam powder diffractometer incorporating a long Soller collimator on the diffracted beam side. Comparison of data collected in

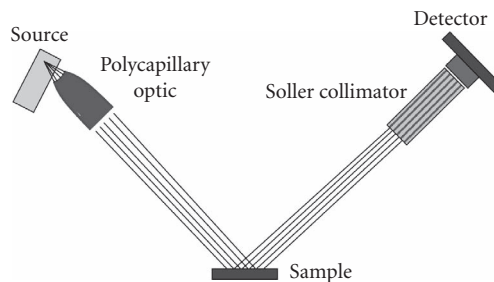


FIGURE 2 A parallel beam diffractometer employing a polycapillary optic and Soller collimator.

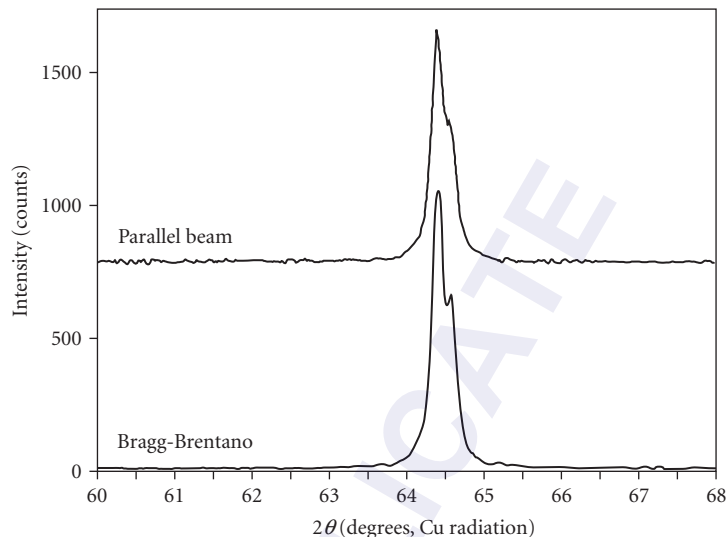


FIGURE 3 Comparison of powder diffraction data for a sample of Ag powder collected using two different instrumental geometries.

the para-focusing Bragg-Brentano configuration to data collected in the parallel beam configuration is shown in Fig. 3 which instantly reveals that instrumental resolution is about the same for both configurations.

The use of the Soller collimator to achieve “parallel” conditions is an important concept, because the beam divergence and instrumental resolution are linked. Consider first the diffractometer in Fig. 1 that works using a divergent beam that diffracts from the sample and then focuses on the receiving slit. As a focusing system, the width of the receiving slit width plays a large role in the instrumental resolution and angular precision, with wide slits worsening both and vice versa. In sharp contrast, the parallel beam system (Fig. 2) relies on the divergence of the Soller collimator to define the measured angle. In other words, only the diffracted x rays that travel through the collimator are detected and these include only the x rays that propagate within the acceptance angle of the collimator (say 0.05° or 180 asec). In the next section, we shall reduce the divergence further by using crystal optics that can reach to $\sim 0.005^\circ$ of divergence, improving the instrumental resolution.

28.3 CRYSTAL OPTICS

Crystal optics are routinely used in two modes: as energy discriminators (monochromators) that define the range of wavelengths used in an experiment and as angular filters that define the beam divergence. The shape and cut of the crystal is critical and allows beam focusing, compression, expansion, and so on by diffracting in one or two dimensions.

Figure 4 shows four crystals, one flat, one bent and cut, one channel-cut, and the fourth doubly curved. In the case of the flat crystal, used either in diffraction or spectroscopy applications, one relies on Bragg’s law to select the wavelength of interest. Focusing crystals (in either 2-D or 3-D, Fig. 4) take on many forms but in general are bent then cut to transform a divergent beam into a focusing beam while selecting some particular wavelength. The channel-cut crystal, is so named because one channel is cut into a single crystal to facilitate diffraction from both inside faces of the channel using a single device. The channels can be cut either symmetrically so that the incident and diffracted beams make

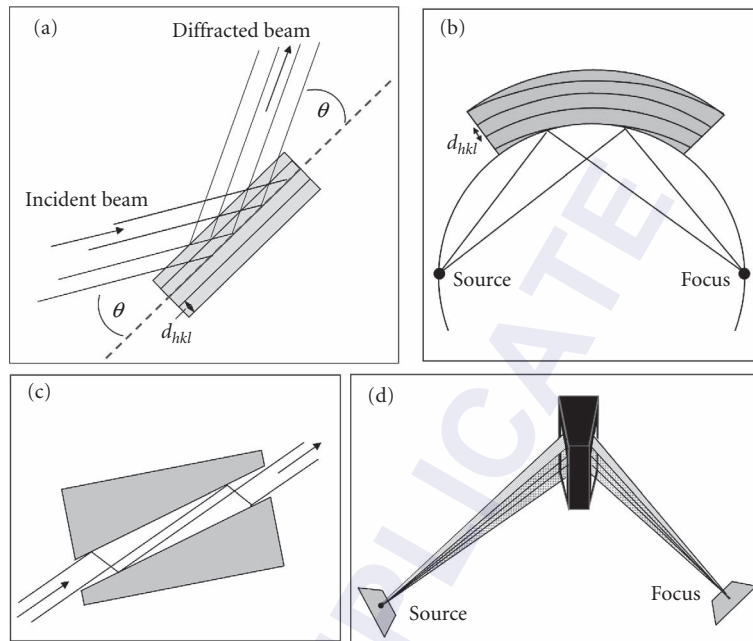


FIGURE 4 The function of several crystal optics including: (a) flat crystal; (b) a bent and cut focusing crystal; (c) an asymmetric channel-cut crystal; and (d) a 2-D focusing crystal.

the same angle with the inside of the channel or asymmetrically where the angle of incidence or diffraction is a small angle with the second angle large. The advantage of asymmetric crystals is higher throughput because of broadening of the rocking curve width.

Selection of crystals involves balancing intensity with resolution (angular or energy), with the latter defined by the rocking curve width. Measuring the intensity diffracted for a particular wavelength as a function of angle provides the rocking curve—a quantitative measure of the perfection of a crystal. The rocking curve width is the most critical aspect of any crystal optic as it defines the range of angles or energies transmitted by the crystal. In the context of energy discrimination, smaller rocking curves result in smaller ranges of energy diffracted by the crystal at some particular angle. Rocking curve widths are a function of wavelength and Miller index, but generally range from ~ 10 asec for high-perfection Si or Ge crystals to ~ 250 asec for LiF or even ~ 1000 for pyrolytic graphite. Graphite and LiF crystals are often used as diffracted beam monochromators in laboratory diffractometers, while Si and/or Ge are reserved for high-resolution epitaxial thin-film analysis or synchrotron beam lines.⁴

In the case of most laboratory diffraction experiments, one or two wavelengths are typically selected, $K\alpha_1$ and/or $K\alpha_2$, using crystal optics. Incorporating a graphite crystal that is highly defected (mosaic) can trim the energy window to include only the $K\alpha_1$ and $K\alpha_2$ components at ~ 60 percent efficiency. Using a crystal of higher perfection facilitates rejection of all but the $K\alpha_1$ radiation, but at a substantially lower efficiency. In order to improve upon the spectral purity and/or beam divergence even further, one can employ multiple crystal monochromators and/or multiple diffraction events using the channel-cut crystal described above. The number of crystals and diffraction events can become quite large for the study epitaxial films in particular, with 4-bounce monochromators on both the incident and diffracted beam sides of the specimen. The reader is referred to recent texts for a more comprehensive review.^{3,4}

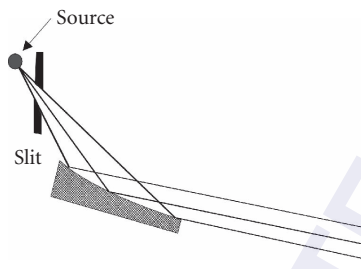


FIGURE 5 A parabolic graded multilayer that transforms a divergent beam into a parallel beam, or vice versa.

28.4 MULTILAYER OPTICS

Multilayer x-ray optics were commercialized in the late 1990s, and offer very specific advantages compared to crystal optics. As shown in Fig. 5 they are man-made crystals that are composed of alternating layers of high and low atomic number materials. They are diffractive optics, generally with large d -spacings and small diffraction angles to provide high efficiency. The rocking curve widths and efficiencies are on the order of 100 to 200 asec with 50 to 70 percent efficiency. As such, they represent a compromise between perfect crystal monochromators (Si, Ge, ~ 10 to 30 asec) and highly mosaic crystals such as graphite (~ 1000 asec).

Multilayer optics are available, like crystals, in a variety of geometrical configurations to provide focused beams and parallel beams, again in one and two dimensions. In addition, cross-coupled multilayers can be used to create point-focused or parallel beams that are today used extensively for single crystal diffraction experiments. The spectral selectivity of multilayers is a function of not only the rocking curve width but also the materials composing the multilayer that can selectively absorb, for example, beta radiation.

28.5 CAPILLARY AND POLYCAPILLARY OPTICS

Drawing hollow glass tubes to a small diameter and with smooth internal surfaces yields single capillary (monocapillary) optics that can be built into arrays to form polycapillary optics. The function of the capillary is total internal reflection of the incident photons that allows the capillary to behave as a “light pipe” to direct x rays in some particular direction. Within limits of the physics of internal x-ray reflection, a variety of beam focusing, collimating, and angular filtering can be achieved, as summarized in Fig. 6. From Fig. 6, it is clear that either mono or polycapillary optics can be used to create small x-ray spot sizes by focusing the beam. Modern capillary optics can provide beam sizes as small as $10\ \mu\text{m}$ routinely, facilitating micro diffraction and micro fluorescence applications.

Another advantage of capillary optics is the ability to improve the intensity in a measurement. Harnessing a large solid angle of x rays emitted from the source or sample in an efficient manner results in 10- or even 100-fold increases in intensity. Similar principles are used in x-ray microsources described below.

28.6 DIFFRACTION AND FLUORESCENCE SYSTEMS

All of the optical components described above can be variously integrated into systems for high intensity, small spot size, large illuminated area, or high energy or angular resolution. One can enhance a diffraction or fluorescence instrument for a specific application by appropriate use of

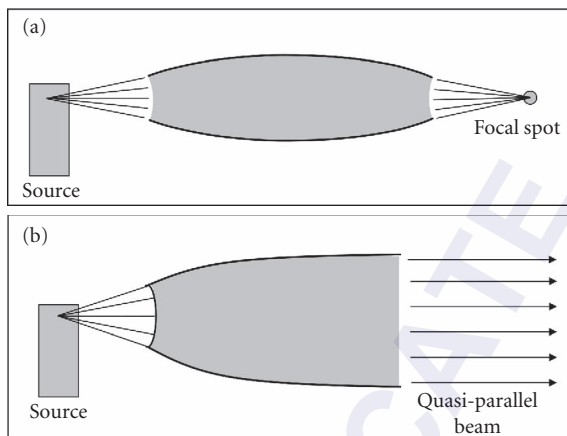


FIGURE 6 Schematics of the function of (a) focusing and (b) collimating polycapillary optics.

optics. Indeed, some modern systems employ prealigned optics that are turn-key interchangeable, affording spectacular flexibility in a single instrument. Naturally, the number of permutations of instrumental arrangements is very large, but in general one attempts to optimize the signal within the limits of the required instrumental resolution.

A clever approach to optimizing both resolution and intensity is the use of “hybrid” optics. Figure 7 shows a schematic of a high-resolution diffractometer applicable for epitaxial film characterization. The defining features of the optics in this case are very high angular resolution provided by the channel cut crystals. However, incorporating a parabolic multilayer before the first crystal monochromator notably improves the intensity. The hybrid design takes advantage of the fact that the multilayer can capture $\sim 0.5^\circ$ of divergent radiation from the x-ray source and convert it to a beam with only ~ 100 asec divergence. Thus, a substantially larger number of photons reach the channel-cut crystal within its rocking curve width of ~ 20 asec from the multilayer than would directly from the x-ray source, improving the overall intensity.

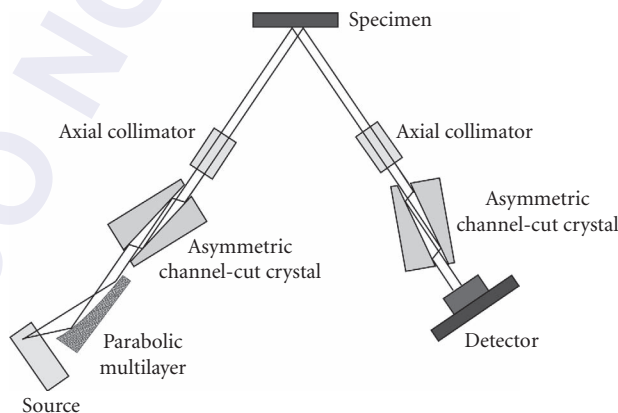


FIGURE 7 Schematic of a high-resolution diffractometer applicable for epitaxial thin film characterization.

28.7 X-RAY SOURCES AND MICROSOURCES

A notable application of x-ray optics is the x-ray “microsource.” Microsource devices in general comprise any low-power and high-flux x-ray source, technology that was enabled by clever application of optical components. The x-ray flux on a particular specimen from a standard x-ray tube is limited by the ability to cool the anode metal, limiting the input power to ~ 2 kW. Rotating anode sources allow for ~ 10 -fold increases in input power, but are again limited by cooling. In either case, traditional systems use a series of slits to guide x rays from the source to the sample in a linear fashion, discarding most of the x rays produced by the source. One can use x-ray optics to harness a larger solid angle of x rays produced at the source and guide those photons to the experiment. Such approaches have been highly successful, leading to commercialization of microsourses that run at power settings as low as 20 W, but provide x-ray flux comparable to traditional x-ray tubes and even rotating anode generators.

28.8 REFERENCES

1. R. Jenkins and R. L. Snyder, *Introduction to X-Ray Powder Diffractometry*, Vol. 138, J. D. Winefordner (ed.), John Wiley & Sons, New York, 1996.
2. H. P. Klug and L. E. Alexander, *X-Ray Diffraction Procedures*, 2nd ed., John Wiley & Sons, New York, 1974, p. 966.
3. B. D. Cullity and S. R. Stock, *Elements of X-Ray Diffraction*, Prentice Hall, NJ, 2001, p. 664.
4. D. K. Bowen and B. K. Tanner, *High Resolution X-Ray Diffractometry and Topography*, Taylor & Francis, London, 1998, p. 252.

This page intentionally left blank.

DO NOT DUPLICATE

REQUIREMENTS FOR X-RAY FLUORESCENCE

Walter Gibson*

*X-Ray Optical Systems
East Greenbush, New York*

George Havrilla

*Los Alamos National Laboratory
Los Alamos, New Mexico*

29.1 INTRODUCTION

The use of secondary x rays that are emitted from solids bombarded by x rays, electrons, or positive ions to measure the composition of the sample is widely used as a nondestructive elemental analysis tool. Such secondary x rays are called *fluorescence x rays*. The “characteristic rays” emitted from a solid irradiated by x rays or electrons¹ were shown in 1913 by Moseley to have characteristic wavelengths (energies) corresponding to the atomic number of specific elements in the target.² Measurement of the wavelength of the characteristic x rays, as well as observation of a continuous background of wavelengths, was made possible by use of the single-crystal diffraction spectrometer first demonstrated by Bragg.³ There was active development by a number of workers and by the late 1920s x-ray techniques were well developed. In 1923, Coster and von Hevesy⁴ used x-ray fluorescence to discover the unknown element hafnium by measurement of its characteristic line in the radiation from a Norwegian mineral, and in 1932 Coster and von Hevesy published the classical text *Chemical Analysis by X-Ray and Its Applications*. Surprisingly, there was then almost no further activity until after World War II. In 1947 Friedman and Birks converted an x-ray diffractometer to an x-ray spectrometer for chemical analysis,⁵ taking advantage of work on diffraction systems and detectors that had gone on in the previous decade. An x-ray fluorescence measurement in which the energy (or wavelength) spectrum is carried out by the use of x-ray diffraction spectrometry is called *wavelength-dispersive x-ray fluorescence* (WDXRF). There was then rapid progress with a number of companies developing commercial x-ray fluorescence (XRF) instruments. The early developments have been discussed in detail by Gilfrich.⁶ During the 1960s, the development of semiconductor particle detectors that could measure the energy spectrum of emitted x rays with much higher energy resolution than possible with gas proportional counters or scintillators resulted in an explosion of applications of energy-dispersive x-ray fluorescence (EDXRF). Until recently, except for the flat or curved diffraction crystals used in WDXRF, x-ray optics have not played an important role in x-ray fluorescence measurements. This situation has changed markedly during the past decade. We will now review the status of both WDXRF and EDXRF with emphasis on the role of x-ray optics without attempting to document the historical development.

*This volume is dedicated in memory of Walter Gibson.

29.2 WAVELENGTH-DISPERSIVE X-RAY FLUORESCENCE (WDXRF)

There are thousands of XRF systems in scientific laboratories, industrial laboratories, and in manufacturing and process facilities worldwide. Although most of these are EDXRF systems, many use WDXRF spectrometry to measure the intensity of selected characteristic x rays. Overwhelmingly, the excitation mechanism of choice is energetic electrons, and many are built onto scanning electron microscopes (SEMs). Electron excitation is simple and can take advantage of electrostatic and magnetic electron optics to provide good spatial resolution and, in the case of the SEM, to give elemental composition maps of the sample with high resolution. In general, the only x-ray optics connected with these systems are the flat or curved analyzing crystals. Sometimes there is a single analyzing crystal that is scanned to give the wavelength spectrum (although multiple crystals, usually two or three, are used to cover different wavelength ranges). However, some systems are multichannel with different (usually curved) crystals placed at different azimuthal angles, each designed to simultaneously measure a specific wavelength corresponding a selected element or background wavelength. Sometimes a scanning crystal is included to give a less sensitive but more inclusive spectral distribution. Such systems have the benefit of high resolution and high sensitivity in cases where the needs are well defined. In general, WDXRF systems have not been designed to take advantage of recent developments in x-ray optics, although there are a number of possibilities and it is expected that such systems will be developed. One important role that x ray optics can be used in such systems is shown in Fig. 1.

In this arrangement, a broad angular range of x-ray emission from the sample is converted into a quasi-parallel beam with a much smaller angular distribution. A variety of collimating optics could be used, for example, polycapillary (as shown), multilayer, nested cone, and so on. The benefit, represented as the gain in diffracted intensity, will depend on the optic used, the system design (e.g., the diffracting crystal or multilayer film), and the x-ray energy. With a polycapillary collimator, 8 keV x rays from the sample with divergence of up to approximately 12° , can be converted to a beam with approximately 0.2° divergence. With a diffraction width of 0.2° and a transmission efficiency for the optic of 50 percent, the gain in the diffracted beam intensity for flat crystal one-dimensional diffraction would be typically more than 30. Further discussion of the gains that can be obtained in x-ray diffraction measurements can be found in Sec. 29.5. Another potential benefit from the arrangement shown in Fig. 1 is confinement of the sampling area to a small spot defined by the collection properties

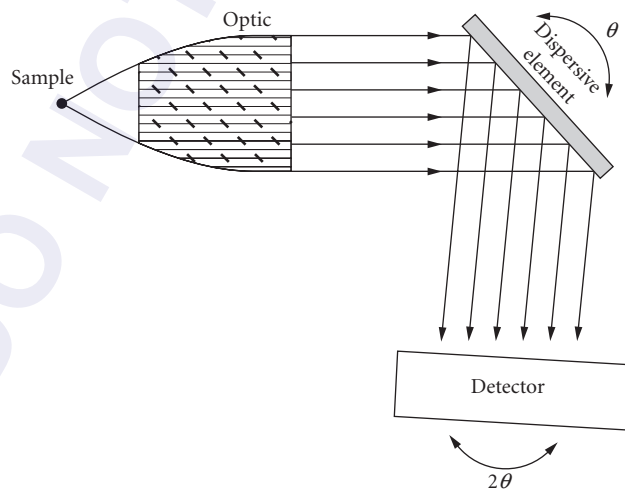


FIGURE 1 Schematic representation of collimating optic in WDXRF system.

of the optic as discussed in Chap. 53. Scanning of the sample will then give the spatial distribution of selected elements. This is also useful in so-called environmental, or high-pressure, SEMs where the position of the exciting electron beam is not so well defined.

Fine Structure in WDXRF Measurements

As noted previously, most of the WDXRF systems in use involve electron excitation either in SEM systems or in dedicated electron micro-probe systems. Photon emission from electron excitation systems contains, in addition to the characteristic lines, a continuous background due to bremsstrahlung radiation resulting from slowing down of the electrons in the solid. Although this background does not seriously interfere with many measurements of elemental composition, it can limit the measurement sensitivity, and can preclude observation of very low intensity features. If the fluorescence x rays are excited by incident x rays, or energetic charged particles, the bremsstrahlung background can be avoided. It should be noted that a continuous background still is present when a broad x-ray spectrum is used as the exciting beam due to scattering of low-energy x rays. This can be largely avoided if monoenergetic x rays are used.⁷

A dramatic illustration of the value of a low background in XRF measurements is contained in recent studies in Japan of fine structure in fluorescence spectra.⁸⁻¹¹ Accompanying each characteristic x-ray fluorescence peak is an Auger excitation peak displaced typically approximately 1 keV in energy and lower in intensity by nearly 1000 times. This peak is not usually observable in the presence of bremsstrahlung background from electron excitation. By using x-ray excitation to get a low background and WDXRF to get high-energy resolution, Kawai and coworkers⁸⁻¹¹ measured the Auger excitation peaks from Silicon in elemental Si and SiO₂ and from Al. They showed that the observed structure corresponds to x-ray absorption fine structure (EXAFS) and x-ray absorption near-edge structure (XANES) that has been observed in high-resolution synchrotron studies. Very long measurement time was necessary to obtain sufficient statistical accuracy. This type of measurement could presumably be considerably enhanced by the use of collimating optics as shown in Fig. 1.

29.3 ENERGY-DISPERSIVE X-RAY FLUORESCENCE (EDXRF)

During the 1960s, semiconductor detectors were developed with dramatic impact on energy and later position measurement of energetic charged particles, electrons, and x rays.^{12,13} Because of their high efficiency, high count rate capability, high resolution compared with gas counters, and their improved energy resolution compared with scintillation counters, these new detectors virtually revolutionized radiation detection and applications including x-ray fluorescence. Initially, the semiconductor junction detectors had a thin active area and, therefore, were not very efficient for x rays. However, by use of lithium compensation in the active area of the detector, it was possible to make very thick depletion layers¹⁴ (junctions) and, therefore, to reach a detection efficiency of 100 percent for x rays. Later, high-purity germanium was used to make thick semiconductor junctions. Such large-volume semiconductor junctions need to be cooled to obtain the highest energy resolution (typically, 130 to 160 eV).

There are now thousands of XRF systems that use cooled semiconductor detectors, most of them mounted on SEMs. Most SEM-based XRF systems do not use any x-ray optics. The detectors can be made large enough (up to 1 to 2 cm²) and can be placed close enough to the sample that they collect x-rays over a relatively large solid angle.

As pointed out previously, electron excitation produces a background of bremsstrahlung radiation that sets a limit on the signal-to-background ratio and, therefore, the minimum detection limit for impurities. This background can be avoided by using x rays or energetic charged particles as the excitation source. Consequently, it is common to see cooled lithium-drifted silicon Si(Li) or high-purity germanium (HpGe) detectors mounted on accelerator beamlines for materials analysis. The ion

beam-based (usually proton or helium ion) technique is called *particle-induced x-ray emission* (PIXE). Again, these do not require optics because the detector can be relatively close to the sample. As with the electron-based systems, optics necessary for controlling or focusing the exciting beam are electrostatic or magnetic and will not be discussed here. When the exciting beam is composed of photons from a synchrotron or free-electron laser (FEL) source, the situation is virtually the same, with no optics required between the sample and the detector. Mirrors and monochromators used to control the exciting beam are discussed in Chaps. 39 and 44.

Monocapillary Micro-XRF (MXRF) Systems

However, if the excitation is accomplished by x rays from a standard laboratory-based x-ray generator, x-ray optics have a very important role to play. In general, the need to obtain a high flux of exciting photons from a laboratory x-ray source requires that the sample be as close as possible to the source. Even then, if the sample is small or if only a small area is irradiated, practical geometrical considerations usually limit the x-ray flux. The solution has been to increase the total number of x rays from the source by increasing the source power, with water-cooled rotating anode x-ray generators becoming the laboratory-based x-ray generator of choice. (For a discussion of x-ray sources, see Chap. 54.) To reduce the geometrical $1/d^2$ reduction of x-ray intensity as the sample is displaced from the source (where d is the sample/source separation), capillaries (hollow tubes) have been used since the 1930s.¹⁵ Although metal capillaries have been used,¹⁶ glass is the overwhelming material of choice,^{17–19} because of its easy formability and smooth surface. In most of the studies reported earlier, a straight capillary was placed between the x-ray source and the sample, and aligned to give the highest intensity on the sample, the capillary length (6 to 20 mm) being chosen to accommodate the source/sample spacing in a commercial instrument. An early embodiment of commercial micro x-ray fluorescence employed metal foil apertures with a variety of dimensions which created spatially resolved x-ray beams. While these crude “optics” provided x-ray beams as small as 50 μm , the x-ray flux was quite limited due to the geometrical constraints.

In 1988, Stern et al.²⁰ described the use of a linearly tapered or conical optic that could be used to produce a smaller, more intense but more divergent beam. This has stimulated a large number of studies of shaped monocapillaries. Many of these are designed for use with synchrotron beams for which they are especially well suited, but they have also been used with laboratory sources. A detailed discussion of monocapillary optics and their applications is given in Chap. 52.

In 1989, Carpenter²¹ built a dedicated system with a very small and controlled source spot size, close-coupling between the capillary and source and variable distance to the sample chamber. The sample was scanned to obtain spatial distribution of observed elemental constituents. A straight 10- μm diameter, 119-mm-long capillary showed a gain of 180 compared to a 10- μm pinhole at the same distance and measurements were carried out with a much lower power x-ray source (12 W) than had been used before.

More recently, a commercial x-ray guide tube or formed monocapillary has been employed to produce x-ray flux gain around 50 times that of straight monocapillary. This modest flux gain enables the more rapid spectrum acquisition and elemental mapping of materials offering new spatially resolved elemental analysis capabilities at the 10 s of micrometers scale.

Polycapillary-Based MXRF

As discussed in Chap. 53, a large number of capillaries can be combined to capture x rays over a large angle from a small, divergent source and focus them onto a small spot. This is particularly useful for microfocus x-ray fluorescence (MXRF) applications.^{22–24}

Using the system developed by Carpenter, a systematic study was carried out by Gao²⁵ in which standard pinhole collimation, straight-capillary, tapered-capillary, and polycapillary focusing optics could be compared. A schematic representation of this system with a polycapillary focusing optic is shown in Fig. 2. The x rays were generated by a focused electron beam, which could be positioned electronically to provide optimum alignment with whatever optical element was being used.²¹

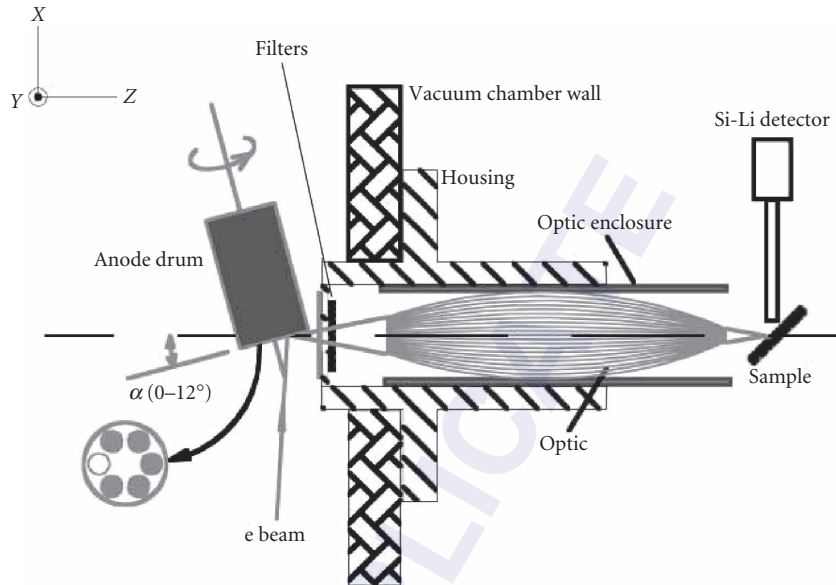


FIGURE 2 Schematic representation of microfocus x-ray fluorescence system. (From Ref. 22.)

The target material could be changed by rotating the anode as shown in Fig. 2. Measurement of the focal spot size produced by the polycapillary focusing optic was carried out by measuring the direct beam intensity while moving a knife edge across the focal spot. The result of such measurements for Cu $K\alpha$ and Mo $K\alpha$ x rays are shown in Fig. 3.

The intensity gain obtained from the polycapillary focusing optic depends on the size of the x-ray emission spot in the x-ray generator, on the x-ray energy, and on the input focal distance (distance between the source spot and the optic). This is because the effective collection angle for each of the transmitting channels is controlled by the critical angle for total external reflection (see Chap. 53). For the system shown, the flux density gain relative to the direct beam of the same size at 100 mm from the source, is shown in Fig. 4 for Cu $K\alpha$ and Mo $K\alpha$ x rays. The maximum gain is about 4400 at 8.0 keV and 2400 at 17.4 keV, respectively.

A secondary x-ray spectrum obtained by irradiating a standard NIST thin-film XRF standard sample, SRM1833, is shown in Fig. 5. Zirconium and aluminum filters were used before the optic to reduce the low-energy bremsstrahlung background from the source. The flux density of the beam at the focus was calculated to be 1.5×10^5 photons \cdot s $\cdot\mu\text{m}^2$ for Mo $K\alpha$ from the 12-W source operated at 40 kV. The minimum detection limits (MDLs) in picograms for 100-s measurement time were as follows: K, 4.1; Ti, 1.5; Fe, 0.57; Zn, 0.28; and Pb, 0.52. The MDL values are comparable with those obtained by Engstrom et al.²⁶ who used a 200- μm diameter straight monocapillary and an x-ray source with two orders of magnitude more power (1.7 kW) than the 12-W source used in the polycapillary measurements.

By scanning the sample across the focal spot of the polycapillary optic, the spatial distribution of elemental constituents was obtained for a rhyolitic glass inclusion in a quartz phenocryst found in a layer of Paleozoic altered volcanic ash. This information is valuable in stratigraphic correlation studies.^{27,28} The results are shown in Fig. 6. Also shown are the images obtained from Compton scattering (Comp) and Rayleigh scattering (Ray).

There are a number of commercial instruments employing the monolithic polycapillary optics to spatially form the excitation beam. Their commercial success lies in being able to generate an increase in x-ray flux 2 to 3 orders of magnitude greater than that can be obtained without the optic at a given

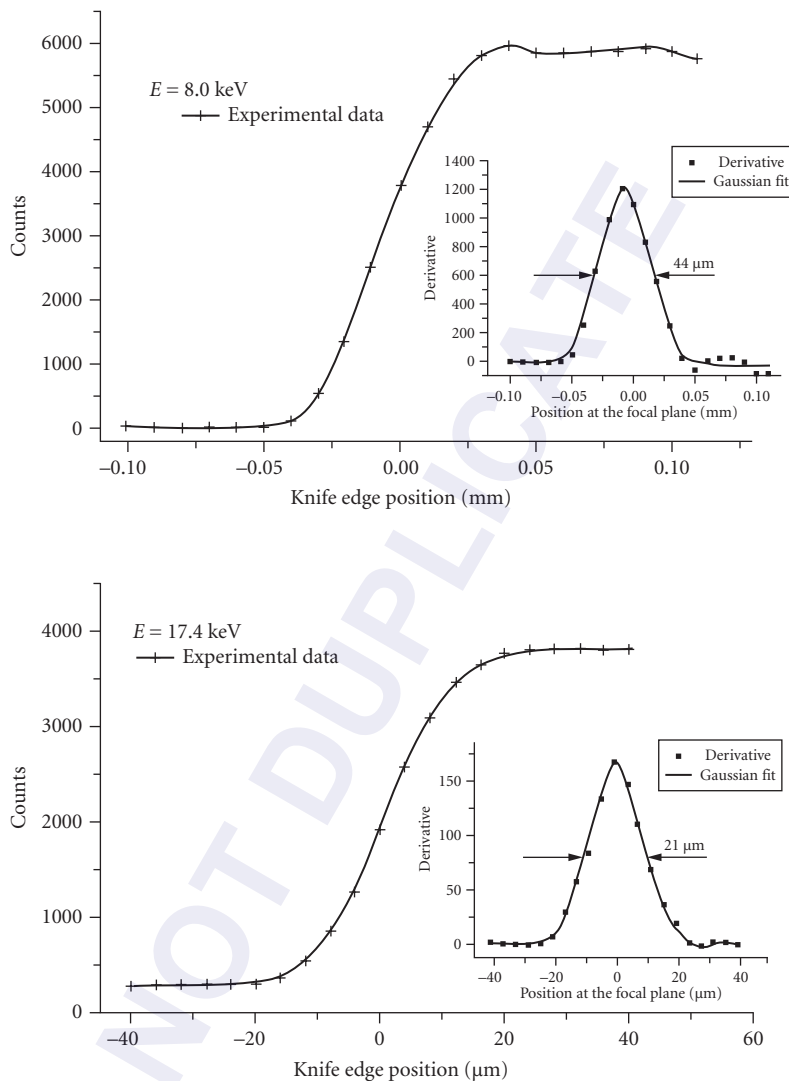


FIGURE 3 Measurement of the focal spot size for Cu and Mo x rays. (From Ref. 25.)

spot size. The future development and potential growth of MXRF rests with the continued innovation of x-ray optics in general and polycapillary optics in particular.

MXRF with Doubly Curved Crystal Diffraction

Although x-ray-induced fluorescence has a significantly lower background than electron-induced fluorescence, there is still background arising from scattering of the continuous bremsstrahlung radiation in the sample. This can be reduced by filtering of high-energy bremsstrahlung in polycapillary focusing optics (see Chap. 53) and by use of filters to reduce the low-energy bremsstrahlung as done for

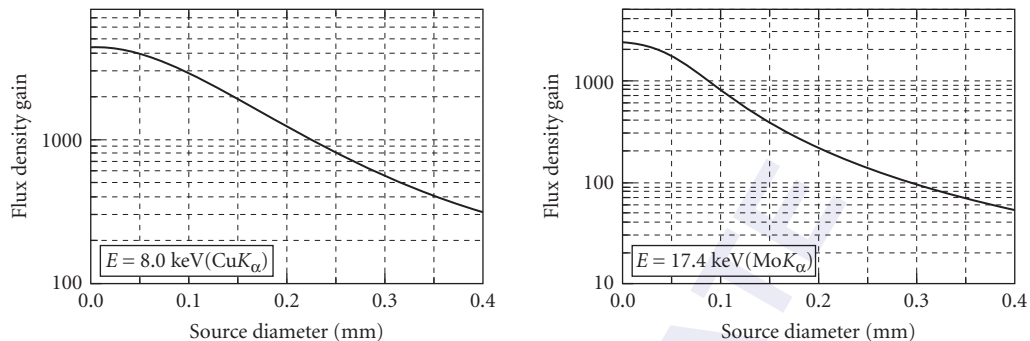


FIGURE 4 Flux density gain as a function of source size. (From Ref. 29.)

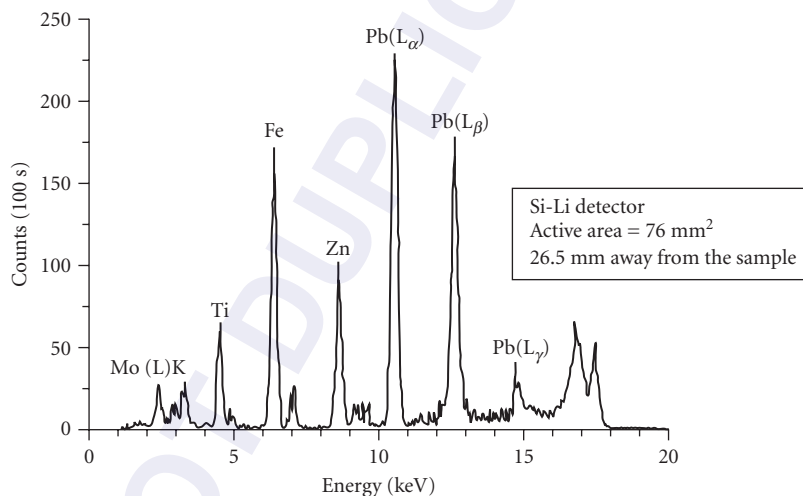


FIGURE 5 Spectrum of SRM 1833 standard XRF thin-film sample. (From Ref. 22.)

the spectrum shown in Fig. 5. Recently, efficient collection and focusing of characteristic x rays by Bragg diffraction with doubly bent single crystals has been demonstrated by Chen and Wytry.³⁰ The arrangement for this is shown in Fig. 7.³¹

Energy spectra taken with a thin hydrocarbon (acrylic) film with a polycapillary focusing optic, and with doubly curved crystal optics with a mica and with a silicon crystal are shown in Fig. 8.³¹ The background reduction for the monoenergetic excitation is evident. Various order reflections are observed with the mica crystal. The angle subtended by the mica crystal is approximately $20^\circ \times 5^\circ$, giving an x-ray intensity only about a factor of three lower than that obtained with the polycapillary lens.

The use of DCCs (doubly curved crystals) in commercial instrumentation has met with commercial success in specific elemental applications. A dual DCC instrument where a DCC is used on the excitation side to create a monochromatic beam for excitation and another DCC on the detection side to limit the region of interest of x-ray fluorescence impinging on the detector provides a highly sensitive and selective detection of sulfur in petroleum streams. Several different embodiments include benchtop, online, and handportable instruments. It is apparent that continued development of DCC-based applications will continue to increase.

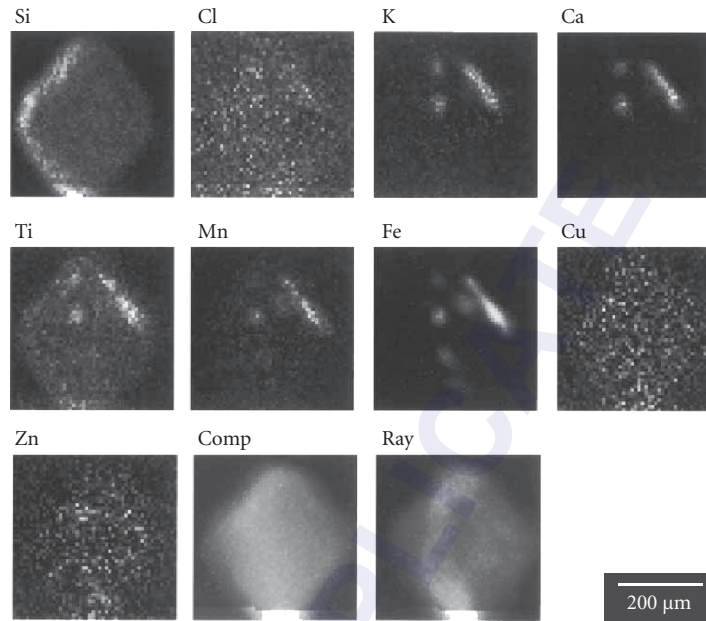


FIGURE 6 MXRF images of various elements in a geological sample that contains small volcanic glass inclusions (tens of micrometers in dimension) within a quartz phenocryst. The last two images are the Compton (energy-shifted) and Rayleigh (elastic) scattering intensity maps. (From Ref. 25.)

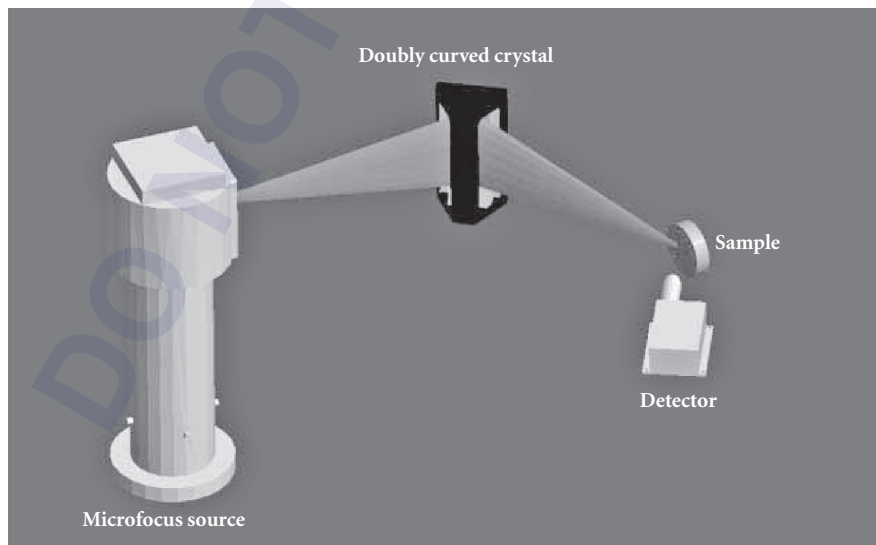


FIGURE 7 Doubly curved crystal MMEDXRF setup. (Courtesy of XOS Inc. From Ref. 31.)

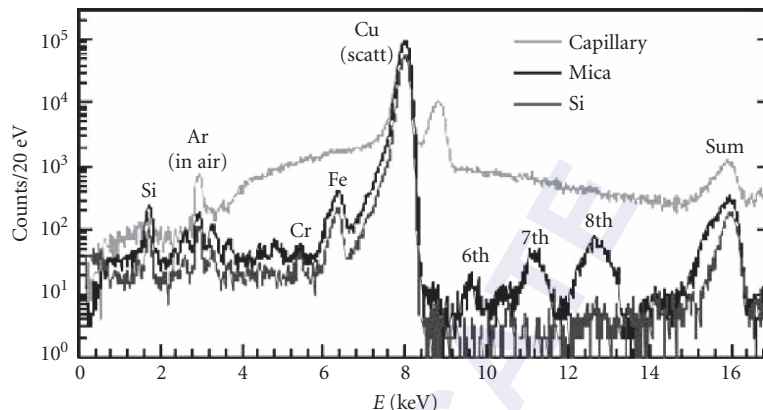


FIGURE 8 Energy spectrum obtained from a thin acrylic sample, measured with a mica doubly bent crystal, a silicon doubly bent crystal, and a polycapillary focusing optic. (From Ref. 31.)

Ultra-high Resolution EDMXRF

As discussed earlier, EDMXRF utilizing cooled semiconductor junction detectors is widely used in science and industry. The energy resolution of semiconductor detectors is typically 140 to 160 eV. During the past few years, very high resolution x-ray detectors based on superconducting transition-edge sensor (TES) microcalorimeters³² semiconductor thermistor microcalorimeters,^{33–35} and superconducting tunnel junctions³⁶ have been developed. Although these detectors are still under active development, there have been demonstrated dramatic benefits for MXRF applications.

TES microcalorimeter detectors have the best reported energy resolution (~ 2 eV at 1.5 keV).³⁷ A schematic representation of a TES microcalorimeter detector is shown in Fig. 9 and energy spectra for a titanium nitride thin film is shown in Fig. 10³⁷ and for a tungsten silicide thin film is shown in

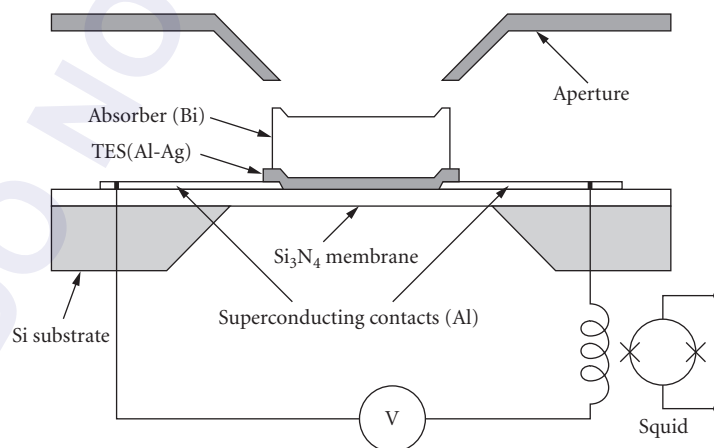


FIGURE 9 A schematic representation of a TES microcalorimeter detector. The operating temperature is ~ 50 mK. (From Ref. 32.)

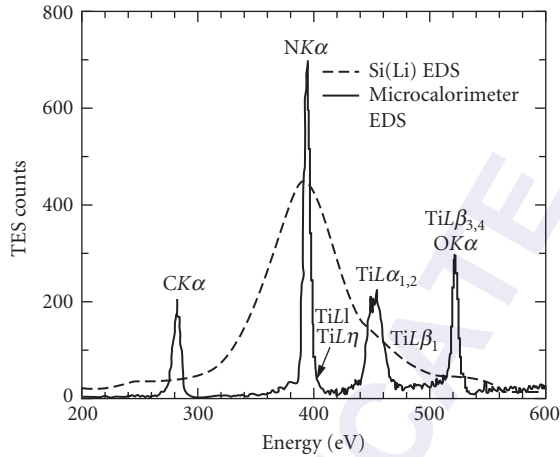


FIGURE 10 Energy spectrum for a TiN thin film on silicon. (From Ref. 37.)

Fig. 11.³² In each case, the spectrum in the same energy region from a silicon junction detector is also shown for comparison.

Because the absorbing element on microcalorimeter detectors must have a low thermal capacitance to achieve high resolution and short recovery time (for higher counting rates), and cannot operate closer than 5 mm to the sample (because of thermal and optical shielding), the sensitivity is low. However, by using a collecting and focusing optic between the sample and the detector, the effective area can be greatly increased.³⁸ A schematic of such an arrangement with a polycapillary focusing optic is shown in Fig. 12. With such a system, the effective area can be increased to approximately 7 mm², comparable with the area of high-resolution semiconductor detectors.

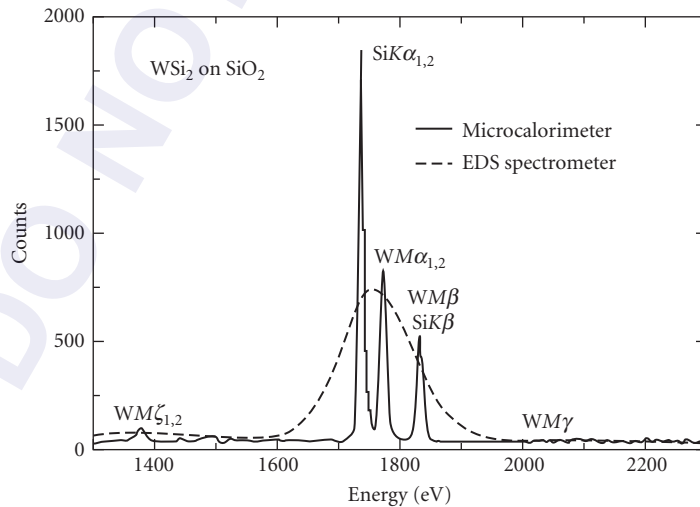


FIGURE 11 Energy spectrum for a WSi₂ film on SiO₂. (From Ref. 32.)

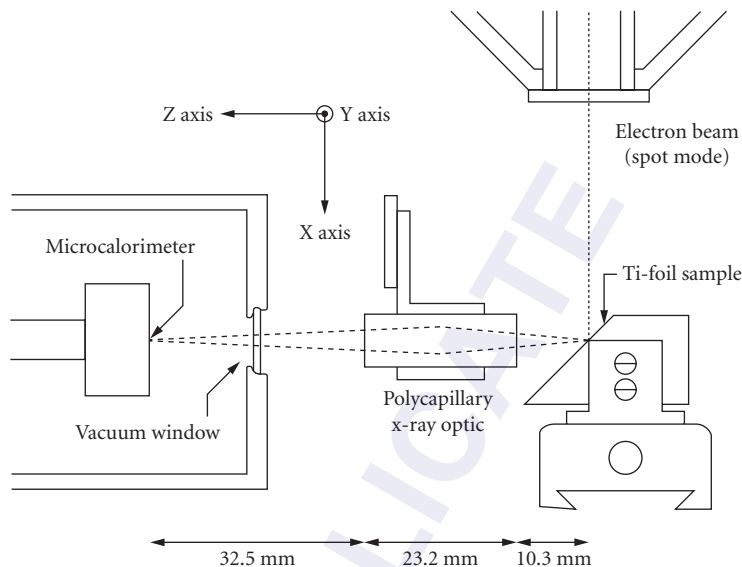


FIGURE 12 Schematic representation of focusing optic for microcalorimeter detector. (From Ref. 38.)

Figure 13 shows a logarithmic plot of the energy spectrum for an aluminum-gallium-arsenide sample.³² This shows the bremsstrahlung background that is present when electrons (in this case, 5 keV) are used as the exciting beam. It is clear that x-ray excitation (especially with monochromatic x rays) will be important in order to use such detectors to observe the low-intensity fine structure to get microstructure and microchemical information with high-resolution energy-dispersive detectors.

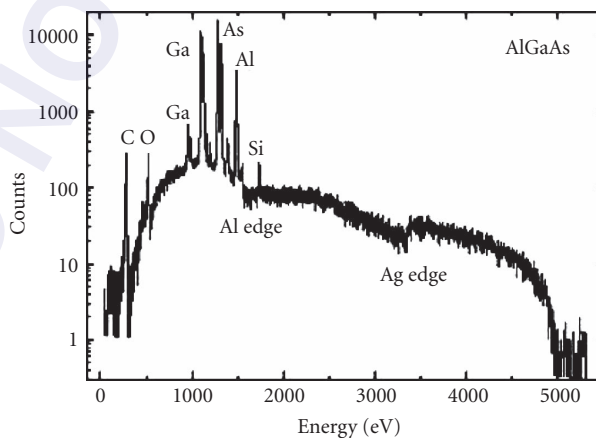


FIGURE 13 Energy spectrum from aluminum-gallium-arsenide sample measured with a TES microcalorimeter spectrometer. (From Ref. 32.)

29.4 REFERENCES

1. R. T. Beatty, "The Direct Production of Characteristic Rontgen Radiations by Cathode Particles," *Proc. Roy. Soc.* **87A**:511 (1912).
2. H. G. J. Moseley, "The High Frequency Spectra of the Elements," *Phil. Mag.* **26**:1024 (1913); **27**:703 (1914).
3. W. H. Bragg and W. L. Bragg, "The Reflection of X-Rays by Crystals," *Proc. Roy. Soc.* **88A**:428 (1913).
4. D. Coster and G. von Hevesey, "On the Missing Element of Atomic Number 72," *Nature* **111**:79, 182 (1923).
5. H. Friedman and L. S. Birks, "A Geiger Counter Spectrometer for X-Ray Fluorescence Analysis," *Rev. Sci. Instr.* **19**:323 (1948).
6. J. V. Gilfrich, "Advances in X-Ray Analysis," *Proc. of 44th Annual Denver X-Ray Analysis Conf.*, Vol. 39, Plenum Press, New York, 1997, pp. 29–39.
7. Z. W. Chen and D. B. Wittry, "Microanalysis by Monochromatic Microprobe X-Ray Fluorescence—Physical Basis, Properties, and Future Prospects," *J. Appl. Phys.* **84**:1064–1073 (1998).
8. J. Kawai, K. Hayashi, and Y. Awakura, "Extended X-Ray Absorption Fine Structure (EXAFS) in X-Ray Fluorescence Spectra," *J. Phys. Soc. Jpn.* **66**:3337–3340 (1997).
9. K. Hayashi, J. Kawai, and Y. Awakura, "Extended Fine Structure in Characteristic X-Ray Fluorescence: A Novel Structural Analysis Method of Condensed Systems," *Spectrochimica Acta* **B52**:2169–2172 (1997).
10. J. Kawai, K. Hayashi, K. Okuda, and A. Nisawa, "Si X-Ray Absorption near Edge Structure (XANES) in X-Ray Fluorescence Spectra," *Chem. Lett. (Japan)* 245–246 (1998).
11. J. Kawai, "Theory of Radiative Auger Effect—An Alternative to X-Ray Absorption Spectroscopy," *J. Electron Spectr. Rel. Phenomona* 101–103, 847–850 (1999).
12. G. L. Miller, W. M. Gibson, and P. F. Donovan. "Semiconductor Particle Detectors," *Ann. Rev. Nucl. Sci.* **33**:380 (1962).
13. W. M. Gibson, G. L. Miller, and P. F. Donovan, "Semiconductor Particle Spectrometers," in *Alpha, Beta, and Gamma Spectroscopy*, 2nd ed., K. Siegbahn (ed.), North Holland Publishing Co. Amsterdam, 1964, p. 345.
14. E. M. Pell, "Ion Drift in an N-P Junction," *J. Appl. Phys.* **31**:291 (1960).
15. F. Jentzsch and E. Nahring, "Reflexion Von Röntgenstrahlen," *Z. The. Phys* **12**:185 (1931); *Z. Tech. Phys.* **15**:151 (1934).
16. L. Marton, "X-Ray Fiber Optics," *Appl. Phys. Lett.* **9**:194 (1966).
17. W. T. Vetterling and R. V. Pound. "Measurements on an X-Ray Light Pipe at 5.9 and 14.4 keV," *J. Opt. Soc. Am.* **66**:1048 (1976).
18. P. S. Chung and R. H. Pantell, "Properties of X-Ray Guides Transmission of X-Rays through Curved Waveguides," *Electr. Lett.* **13**:527 (1977); *IEEE J. Quant. Electr.* **QE-14**:694 (1978).
19. A. Rindby, "Applications of Fiber Technique in the X-Ray Region," *Nucl. Instr. Meth.* **A249**:536 (1986).
20. E. A. Stern, Z. Kalman, A. Lewis, and K. Lieberman, "Simple Method for Focusing X Rays Using Tapered Capillaries," *Appl. Optics* **27**:5135 (1988).
21. D. A. Carpenter, "Improved Laboratory X-Ray Source for Microfluorescence Analysis," *X-Ray Spectrometry* **18**:253–257 (1989).
22. N. Gao, I. Yu. Ponomarev, Q. F. Xiao, W. M. Gibson, and D. A. Carpenter, "Enhancement of Microbeam X-Ray Fluorescence Analysis Using Monolithic Polycapillary Focusing Optics," *Appl. Phys. Lett.* **71**:3441–3443 (1997).
23. Y. Yan and X. Ding, "An Investigation of X-Ray Fluorescence Analysis with an X-Ray Focusing System (X-Ray Lens)," *Nucl. Inst. Meth. Phys. Res.* **B82**:121 (1993).
24. M. A. Kumakhov and F. F. Komarov, "Multiple Reflection from Surface X-Ray Optics," *Phys. Reports* **191**:289–350 (1990).
25. N. Gao, "Capillary Optics and Their Applications in X-Ray Microanalysis," Ph.D. Thesis, University at Albany, SUNY, Albany, NY, 1998.
26. P. Engstrom, S. Larsson, A. Rindby, and B. Stocklassa, "A 200 mm X-Ray Microbeam Spectrometer," *Nucl. Instrum. Meth.* **B36**:222 (1989).
27. J. W. Delano, S. E. Tice, C. E. Mitchell, and D. Goldman, "Rhyolitic Glass in Ordovician K-Bentonites: A New Stratigraphic Tool," *Geology* **22**:115 (1994).

28. B. Hanson, J. W. Delano, and D. J. Lindstrom, "High-Precision Analysis of Hydrous Rhyolitic Glass Inclusions in Quartz Phenocrysts Using the Electron Microprobe and INAA," *American Mineralogist* **81**:1249 (1996).
29. J. X. Ho, E. H. Snell, C. R. Sisk, J. R. Ruble, D. C. Carter, S. M. Owens, and W. M. Gibson, "Stationary Crystal Diffraction with a Monochromatic Convergent X-Ray Beam Source and Application for Macromolecular Crystal Data Collection," *Acta Cryst.* **D54**:200–214 (1998).
30. Z. W. Chen and D. B. Wittry, "Microanalysis by Monochromatic Microprobe X-Ray Fluorescence—Physical Basis, Properties, and Future Prospects," *J. Appl. Phys.* **84**:1064–1073 (1998).
31. Ze Wu Chen, R. Youngman, T. Bievenue, Qi-Fan Xiao, I. C. E. Turcu, R. K. Grygier, and S. Mrowka, "Polycapillary Collimator for Laser-Generated Plasma Source X-Ray Lithography," in C. A. MacDonald, K. A. Goldberg, J. R. Maldonado, H. H. Chen-Mayer, and S. P. Vernon (eds.), *EUV, X-Ray and Neutron Optics and Sources*, SPIE **3767**:52–58 (1999).
32. D. A. Wollman, K. D. Irwin, G. C. Hilton, L. L. Dulcie, D. E. Newbury, and J. M. Martinis, "High-Resolution, Energy-Dispersive Microcalorimeter Spectrometer for X-Ray Microanalysis," *J. Microscopy* **188**:196–223 (1997).
33. D. McCammon, W. Cui, M. Juda, P. Plucinsky, J. Zhang, R. L. Kelley, S. S. Holt, G. M. Madejski, S. H. Moseley, and A. E. Szymkowiak, "Cryogenic Microcalorimeters for High Resolution Spectroscopy: Current Status and Future Prospects," *Nucl. Phys.* **A527**:821 (1991).
34. L. Lesyna, D. Di Marzio, S. Gottesman, and M. Kesselman, "Advanced X-Ray Detectors for the Analysis of Materials," *J. Low. Temp. Phys.* **93**:779 (1993).
35. E. Silver, M. LeGros, N. Madden, J. Beeman, and E. Haller, "High-Resolution, Broad-Band Microcalorimeters for X-Ray Microanalysis," *X-Ray Spectrom.* **25**:115 (1996).
36. M. Frank, L. J. Hiller, J. B. le Grand, C. A. Mears, S. E. Labov, M. A. Lindeman, H. Netel, and D. Chow, "Energy Resolution and High Count Rate Performance of Superconducting Tunnel Junction X-Ray Spectrometers," *Rev. Sci. Instrum.* **69**:25 (1998).
37. G. C. Hilton, D. A. Wollman, K. D. Irwin, L. L. Dulcie, N. F. Bergren, and J. M. Martinis, "Superconducting Transition-Edge Microcalorimeters for X-Ray Microanalysis," *IEEE Trans. Appl. Superconductivity* **9**:3177–3181 (1999).
38. D. A. Wollman, C. Jezewski, G. C. Hilton, Q. F. Xiao, K. D. Irwin, L. L. Dulcie, and J. M. Martinis, *Proc. Microscopy and Microanalysis* **3**:1075 (1997).

This page intentionally left blank.

DO NOT DUPLICATE

REQUIREMENTS FOR X-RAY SPECTROSCOPY

Dirk Lützenkirchen-Hecht and Ronald Frahm

*Bergische Universität Wuppertal
Wuppertal, Germany*

The basic process related to x-ray absorption spectroscopy (XAS, which includes x-ray absorption near edge spectroscopy, XANES, and extended x-ray absorption fine structure, EXAFS) is the absorption by an atom of a photon with sufficient energy to excite a core electron to unoccupied levels (bands) or to the continuum. An XAS experiment comprises the measurement of the absorption coefficient $\mu(E)$ in the vicinity of an absorption edge, where a more or less oscillatory behavior of $\mu(E)$ is observed. Monochromatic radiation is used, and the photon energy is increased to the value at which core electrons can be excited to unoccupied states close to the continuum. As an example, an absorption spectrum of a platinum metal foil is shown in Fig. 1 for the photon energy range in the vicinity of the Pt L_3 edge. The steep increase of the absorption at about 11.564 keV corresponds to the excitation of electrons from the Pt $2p_{3/2}$ level into unoccupied states above the Fermi level of the Pt material. In general, the exact energy of the absorption edge is a sensitive function of the chemical valence of the excited atom, that is, a shift of the absorption edge toward higher photon energies is observed as the chemical valence of the absorber atom is increased. In many cases, a more or less linear shift of the edge with typically 1 to 3 eV per valence unit can be found in the literature for different elements, and thus, an accurate determination of the edge position is essential for a proper valence determination.¹ Furthermore, in some cases there are also sharp features in the absorption spectrum even below the edge. These so called pre-edge peaks can be attributed to transitions from the excited photoelectron into unoccupied discrete electronic levels of the sample which can be probed by an x-ray absorption experiment. Due to the discrete nature of these transitions, again a high resolution of the spectrometer is required in order to be able to investigate such structures in the spectrum. In the case of Pt, a sharp whiteline-like feature is observed directly at the edge. As can be seen in the inset of Fig. 1, where the derivative of the Pt absorption spectrum is shown, the sharpest features in this spectrum have a full width of only a few electron volts and the energy resolution ΔE of the spectrometer has to be superior for such experiments.

In general, synchrotron radiation is used for x-ray absorption spectroscopy, because the energy range of interest can be easily selected from the broad and intense distribution provided by storage rings (see Chap. 55). X-ray monochromators select the desired photon energy by Bragg's law of diffraction

$$n\lambda = 2d \sin\Theta \quad (1)$$

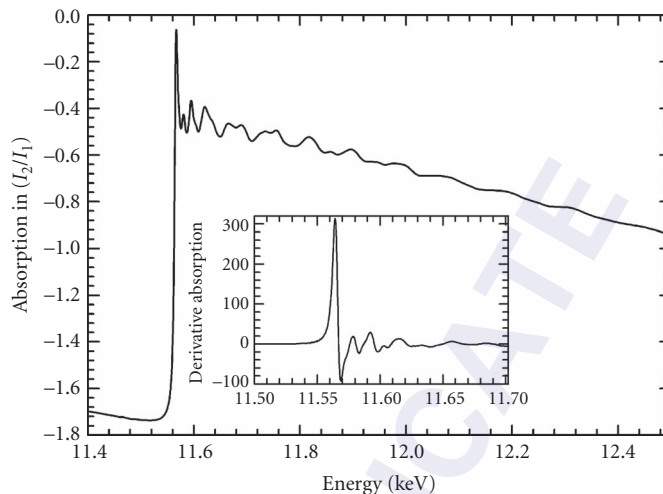


FIGURE 1 High resolution XANES spectrum of a Pt-metal foil at the L_3 -absorption edge measured in transmission mode at the DELTA-XAS beamline using a Si(111)-double crystal monochromator with adaptive optics. In the inset, the derivative of the absorption spectrum is presented in the near edge region.

where λ is the x-ray wavelength, d the lattice spacing of the monochromator crystal, and Θ the angle between the impinging radiation and the lattice planes. However, not only the fundamental wave ($n = 1$) but also higher harmonics ($n > 1$) are transmitted by the monochromator. Selection rules depending on the structure of the crystals can forbid certain harmonics, for example, in the case of the mostly used Si(111) and Si(311) monochromator crystals the second order are not allowed, whereas the third harmonics are present.

However, high-intensity radiation impinging on a monochromator crystal induces a variety of surface slope errors, such as thermal bumps and thermal bending as well as lattice constant variations. This is especially true if insertion devices (wigglers and undulators) are used at third-generation storage rings. Typical heatloads can reach the kilowatt regime, and only a small fraction of typically 10^{-4} is Bragg reflected. Therefore, different concepts have been developed to compensate these unwanted distortions of the Bragg-reflecting surfaces.²⁻⁶ It is far beyond the scope of this paper to describe these efforts in detail, here we will only refer to the literature (see, e.g., Chap. 39), however, stressing the importance of an appropriate compensation. In Fig. 2, we present rocking curves of a Si(111) double-crystal monochromator in order to demonstrate the effect of a compensation on the width (and thereby the energy resolution) of the measured curves. Having in mind that the theoretical value for the energy resolution amounts to ca. 1.3 eV at 9 keV photon energy, those results clearly demonstrate the need for a compensation mechanism. The energy resolution of the noncompensated monochromator with a rocking curve width of more than 100 arc sec, corresponding to more than 14 eV, would be useless for XANES spectroscopy.

Above an absorption edge, a series of wiggles with an oscillatory structure is visible, which modulates the absorption typically by a few percent of the overall absorption cross section as can be seen in Fig. 1. These wiggles are caused by the scattering of the ejected photoelectrons by neighboring atoms and contain quantitative information regarding the structure of the first few coordination shells around the absorbing atom such as bond distances, coordination numbers, and Debye-Waller factors.⁷ Typically, the EXAFS region extends from approximately 40 eV above the x-ray absorption edge up to about 1000 eV or even more.

Higher harmonics in the monochromatic beam may disturb the actual measurement, that is, the absorption coefficients and thus also the EXAFS will be erroneous, so that a suppression of higher

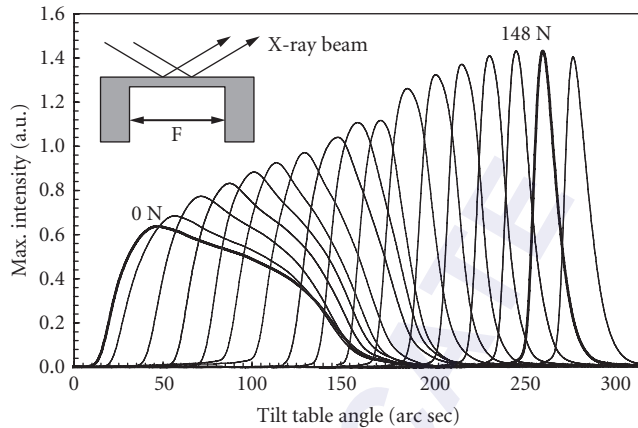


FIGURE 2 Rocking curves measured for a Si(111) crystal pair at 8.9 keV photon energy and a heat load of about 370 W for different bending forces F on the legs of a U-shaped crystal bender as depicted in the inset. The intrinsic width of the rocking curve is 8.9 arc sec (ca. 1.3 eV) compared to about 12 arc sec (ca. 1.8 eV) for the optimum compensated crystal. For comparison, the width of the noncompensated crystal amounts to more than 100 arc sec with less than half of the peak intensity.⁶

harmonics is highly desirable. This is especially important for investigations at lower edge energies, for example, at the Ti K-edge or even lower. Here, any parasitically absorbing elements in the beam path such as the x-ray windows from the beamline to ambient conditions, or even short air pathways, increasingly absorb the photons of interest while the absorption of higher harmonics is negligible. For example, if experiments at the Ti K-edge (4.966 keV) are considered, a beam path of 20 cm in air would result in an absorption of 62 percent for the fundamental wave, in contrast to only 3.5 percent absorption for the third harmonic. In addition, using a Be window of 750 μm thickness as separator between the ultrahigh vacuum system of the beamline and the experimental setup which is usually in ambient air, the transmission increases continuously from only 55 percent at 5 keV to more than 96 percent at 15 keV. In general, thus, the relative intensity of the third harmonic will increase in the course of the beam bath. This is the reason why one has to consider carefully all contributions from the third harmonic, and their influence on the measured spectra. This is illustrated in Fig. 3 for a transmission spectrum of a TiO_2 sample measured in transmission using a Si(111) double crystal monochromator and nitrogen-filled ionization chambers as detectors. Both the different absorption of the low- and the high-energy beams as well as the different ionization of the two beams have been considered. Even in the case of a beam coming from the double crystal monochromator in the vacuum section of a beamline with only 5 percent harmonic content, a significant reduction of the edge jump is visible, as well as a slight damping of the pre-edge peaks. However, this situation, in which the photons of the fundamental wave are strongly absorbed by the environment in contrast to those of the third harmonic, may be regarded as the worst case. This hardening of the beam has a major impact for x-ray studies at lower energies.

From the presented simulation it can be concluded that the x-ray optics have to ensure that the beam is really monochromatic. This may be achieved using different techniques such as detuning of the monochromator crystals⁸ or the use of a mirrored beam. Detuning makes use of the fact that the crystal rocking curve width of higher harmonics is generally much smaller compared to the fundamental wave,⁹ so that a slight tilt of the crystals with respect to each other effectively suppresses the transmitted intensity of higher harmonics. Such a detuning procedure is however not applicable in the case of a channel-cut monochromator, where the two reflecting crystal surfaces originate from a monolithic single crystal with a fixed geometric relation to each other. Here, one makes use of the fact that the x-ray reflectivity of any surface generally decreases with photon energy, that is, a mirror can also be used as

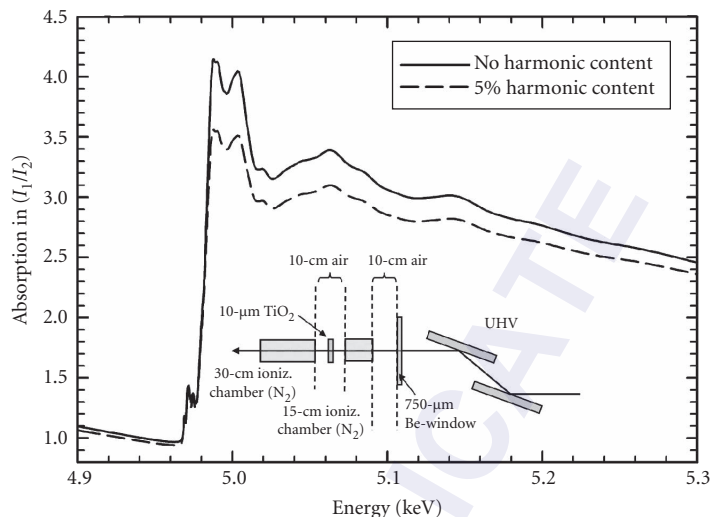


FIGURE 3 Calculation of a XANES spectrum of a nanocrystalline TiO_2 sample (10- μm thickness) for x-ray beams with varying harmonic content. In the calculations, the different absorption of the first and the third harmonic of the Si(111)-double crystal monochromator in the different materials in the optical path were included. We have used a 750- μm Be window which terminates the ultrahigh vacuum (UHV) section of the beamline as well as a 10-cm beam path through laboratory air between the Be window and the first, 15-cm long, N_2 -filled ionization chamber. A second 10-cm air path and the TiO_2 sample were present between the first and the second, (30-cm N_2) ionization chamber.

low pass filter for the harmonic rejection, if the critical energy is smaller than the energy of the corresponding harmonic wave (Refs. 10 and 11 and Chap. 44). It should be mentioned here that the heat load on the monochromator crystals is also reduced if a mirror in front of the monochromator is used, and thus the related problems mentioned above are less important. Furthermore, the mirrors can also be used for the focussing of the x-ray beam:^{11–13} Depending on the surface of the mirror (flat, spherical, cylindrical) and its bending radius, the impinging radiation is either divergent or focused in one or two directions. Using an undulator as source, the point focus of a mirror system can be as small as a few tens of microns, so that it is possible to investigate small specimen (e.g., single crystallites), and spectromicroscopy or microspectroscopic investigations are feasible, even under nonambient conditions.^{14,15}

Up to now, we have not dealt with the problem of lateral stability of the beam on the sample. Even in an idealized setup, we have to consider vertical beam movements on the sample during an EXAFS scan because the beam offset changes significantly for an extended scan as illustrated schematically in Fig. 4. More quantitatively, for a fixed distance D of the x-ray reflecting planes, the beam offset amounts to $h = 2D \cos \Theta$. Thus, h will be larger for a smaller Bragg angle Θ , as can be seen in Fig. 4. Given a typical distance of $D = 50$ mm, a Bragg-angle variation from about 16.5° to 14.5° (which would roughly correspond to an EXAFS scan at the Fe K-edge from ca. 6.9 keV to 7.9 keV) would result in a variation of the beam offset by about 1 mm. Such beam movement is not acceptable for certain experiments, for example, for the investigation of small or inhomogeneous samples or in the case of grazing incidence x-ray experiments, where the height of the sample in the beam is often limited to only a few microns.¹⁶ Furthermore, if the beam downstream the monochromator is subjected to focussing, e.g., for XANES tomography¹⁷ or spectromicroscopic investigations, such beam movements exceed the acceptance width of the focussing optics (see also Chaps. 37, 40, 42, 44, 45, 52, and 53 of this *Handbook*). Thus, a fixed exit geometry by an additional movement of the second crystal (i.e., a variation of the distance D between the Bragg-reflecting surfaces) or, alternatively, a controlled

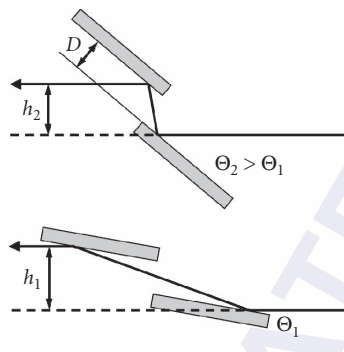


FIGURE 4 Schematic representation of vertical beam offset (h) changes during Bragg-angle variation for two different Bragg angles Θ_1 and Θ_2 , as indicated.

correction of the vertical sample position by means of a lifting table as a function of the Bragg-angle are required. In the case of a channel-cut crystal, a special form of the crystals' reflecting surfaces may also ensure that there is a constant beam height.^{18,19}

We want to conclude here by pointing out that all the topics mentioned above also apply in the case of time-resolved x-ray absorption experiments, where the Bragg angle of the monochromator crystals is moved continuously and the spectrum is collected on the fly (quick-scanning EXAFS^{20, 21}). It should be mentioned here, that a high precision of all movements can be reached so that XANES data can be measured in ca. 5 ms, while EXAFS spectra spanning about 1.5 keV of real samples such as catalysts under working conditions are feasible in about 50 ms with sufficient data quality.^{21,22}

30.1 REFERENCES

1. B. Lengeler, "X-Ray Absorption and Reflection in Materials Science," *Advances in Sol. State Phys.* **29**:53–73 (1989).
2. J. Arthur, W. H. Tompkins, C. Troxel, Jr., R. J. Contolini, E. Schmitt, D. H. Bilderback, C. Henderson, J. White, and T. Settersten, "Microchannel Water Cooling of Silicon X-Ray Monochromator Crystals," *Rev. Sci. Instrum.* **63**:433–436 (1992).
3. J. P. Quintana, M. Hart, D. Bilderback, C. Henderson, D. Richter, T. Setterston, J. White, D. Hausermann, M. Krumrey, and H. Schulte-Schrepping, "Adaptive Silicon Monochromators for High-Power Insertion Devices. Tests at CHESS, ESRF and HASYLAB," *J. Synchrotron Rad.* **2**:1–5 (1995).
4. R. K. Smither, G. A. Forster, D. H. Bilderback, M. Bedzyk, K. Finkelstein, C. Henderson, J. White, L.E. Berman, P. Stefan, and T. Oversluizen, "Liquid Gallium Cooling of Silicon Crystals in High Intensity Photon Beams," *Rev. Sci. Instrum.* **60**:1486–1492 (1989).
5. C. S. Rogers, D. M. Mills, W.-K. Lee, P. B. Fernandez, and T. Graber, "Experimental Results with Cryogenically Cooled Thin Silicon Crystal X-Ray Monochromators on High Heat Flux Beamlines," *Proc. SPIE* **2855**:170–179 (1996).
6. R. Zaepfer, M. Richwin, D. Lützenkirchen-Hecht, and R. Frahm, "A Novel Crystal Bender for X-Ray Synchrotron Radiation Monochromators," *Rev. Sci. Instrum.* **73**:1564–1567 (2002).
7. D. Koningsberger and R. Prins, *X-Ray Absorption: Principles, Applications, Techniques of EXAFS, SEXAFS and XANES*, New York: John Wiley and Sons (1988).
8. A. Krolzig, G. Materlik, and J. Zegenhagen, "A Dynamic Control and Measuring System for X-Ray Rocking Curves," *Nucl. Instrum. Meth.* **208**:613–619 (1983).
9. G. Materlik and V. O. Kostroun, "Monolithic Crystal Monochromators for Synchrotron Radiation with Order Sorting and Polarizing Properties," *Rev. Sci. Instrum.* **51**:86–94 (1980).

10. B. W. Batterman and D. H. Bilderback, "X-Ray Monochromators and Mirrors," In: *Handbook of Synchrotron Radiation* Vol. 3: *X-Ray Scattering Techniques and Condensed Matter Research* G. Brown and D. Moncton (eds.). Amsterdam: North Holland, pp. 105–120 (1991).
11. S. M. Heald and J. B. Hastings, "Grazing Incidence Optics for Synchrotron Radiation X-Ray Beamlines," *Nucl. Instrum. Meth.* **187**:553–561 (1981).
12. P. Kirkpatrick and A. V. Baez, "Formation of Optical Images by X-Rays," *J. Opt. Soc. Am.* **38**:766–774 (1948).
13. J. A. Howell and P. Horowitz, "Ellipsoidal and Bent Cylindrical Condensing Mirrors for Synchrotron Radiation," *Nucl. Instrum. Meth.* **125**:225–230 (1975).
14. M. Newville, S. Sutton, M. Rivers, and P. Eng, "Micro-Beam X-Ray Absorption and Fluorescence Spectroscopies at GSECARS: APS Beamline 13ID," *J. Synchrotron Rad.* **6**:353–355 (1999).
15. U. Kleineberg, G. Haindl, A. Hütten, G. Reiss, E. M. Gullikson, M. S. Jones, S. Mrowka, S. B. Rekawa, and J. H. Underwood, "Microcharacterization of the Surface Oxidation of Py/Cu Multilayers by Scanning X-Ray Absorption Spectromicroscopy," *Appl. Phys. A* **73**:515–519 (2001).
16. D. Hecht, R. Frahm, and H.-H. Strehblow, "Quick-Scanning EXAFS in the Reflection Mode as a Probe for Structural Information of Electrode Surfaces with Time Resolution: An In Situ Study of Anodic Silver Oxide Formation," *J. Phys. Chem.* **100**:10831–10833 (1996).
17. C. G. Schroer, M. Kuhlmann, T. F. Günzler, et al., "Tomographic X-Ray Absorption Spectroscopy," *Proc. SPIE* **5535**:715–723 (2004).
18. P. Spieker, M. Ando, and N. Kamiya, "A Monolithic X-Ray Monochromator with Fixed Exit Beam Position," *Nucl. Instrum. Meth.* **222**:196–201 (1984).
19. S. Oestreich, B. Kaulich, R. Barrett, and J. Susini, "One-Movement Fixed-Exit Channel-Cut Monochromator," *Proc. SPIE* **3448**:176–188 (1998).
20. R. Frahm, "New Method for Time Dependent X-Ray Absorption Studies," *Rev. Sci. Instrum.* **60**:2515–2518 (1989).
21. H. Bornebusch, B. S. Clausen, G. Steffensen, D. Lützenkirchen-Hecht, and R. Frahm, "A New Approach for QEXAFS Data Acquisition," *J. Synchrotron Rad.* **6**:209–211 (1999).
22. R. Frahm, B. Griesebock, M. Richwin, and D. Lützenkirchen-Hecht, "Status and New Applications of Time-Resolved X-Ray Absorption Spectroscopy," *AIP Conf. Proc.* **705**:1411–1414 (2004).

REQUIREMENTS FOR MEDICAL IMAGING AND X-RAY INSPECTION

Douglas Pfeiffer

*Boulder Community Hospital
Boulder, Colorado*

31.1 INTRODUCTION TO RADIOGRAPHY AND TOMOGRAPHY

In the early years following Roentgen's discovery of x rays in 1895, applications and implementations of x rays in medical imaging and therapy were already being developed. The development of detectors, scatter control devices, and associated medical imaging paraphernalia has continued unabated. While more complex optics are being researched, the x-ray optics commonly used in medical imaging are apertures, filters, collimators, and scatter rejection grids. The use of x rays in nondestructive testing was actually mentioned in Roentgen's seminal paper, though this application became active only in the 1920s.

31.2 X-RAY ATTENUATION AND IMAGE FORMATION

It is generally true that medical and industrial radiography are transmission methods. X rays are passed through the object being imaged and detected on the opposite side. The image is formed via the differential attenuation of the x rays as they pass through the object. Variations in thickness or density within the object result in corresponding fluctuations in the intensity of the x-ray flux exiting the object. Most simply, this is defined by the linear attenuation coefficient of the material in the path of a narrow beam of radiation,

$$I = I_0 e^{-\mu x} \quad (1)$$

where I_0 and I are the incident and transmitted x-ray beam intensities, respectively, x is the thickness of the material, and μ is the linear attenuation coefficient, which has units of cm^{-1} and is dependent upon the physical characteristics of the material and the energy of the x-ray beam, as shown in Fig. 1. For heterogeneous objects, Eq. (1) is modified to account for each material, such as differing organs or a cancer, in the path of the beam, as shown in Fig. 2, so that

$$I = I_0 e^{-(\mu_1 x_1 + \mu_2 x_2 \dots)} \quad (2)$$

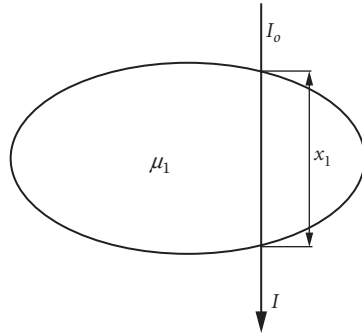


FIGURE 1 Attenuation through a uniform object.

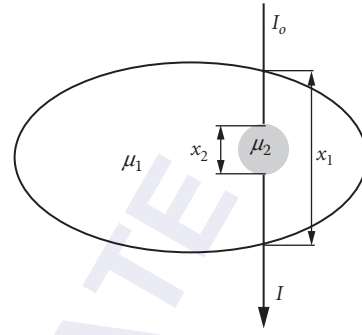


FIGURE 2 Attenuation through an object having an inclusion of differing attenuation.

For a small object, such as a void in a weld, to be visualized, there must be sufficient contrast in the attenuation between the beam passing through the object and beam passing just adjacent to it, as shown in Fig. 3. The contrast is

$$C = \frac{I_2 - I_1}{I_1} \quad (3)$$

where I_2 and I_1 are the intensities of the beams exiting from behind the small object and just adjacent to it, respectively. The amount of contrast required to confidently visualize a given object is determined through the Rose model,¹ a discussion of which is beyond the scope of this chapter. All transmission-based medical and industrial imaging is based upon these fundamental concepts.

As stated earlier, the attenuation coefficient is energy dependent and highly nonlinear. Attenuation in the medical and industrial energy range is dominated mainly by coherent and incoherent (Compton) scattering and photoelectric absorption. Mass attenuation coefficients, are linear attenuation coefficients divided by the density of the material, and thus have units of cm^2/g . Representative mass attenuation curves are demonstrated in Fig. 4, where (a) demonstrates the attenuation curve of water and (b) demonstrates the attenuation curve of iodine.² The discontinuities in the curve for iodine, a commonly used contrast medium in medical imaging, are due to the K and L electron shell absorption edges.

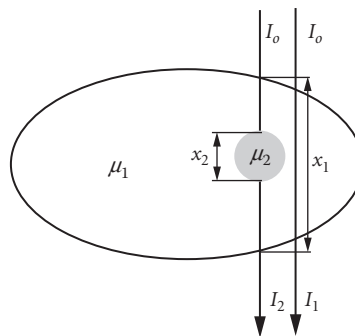


FIGURE 3 The physical basis of radiographic contrast, the difference in attenuation between two adjacent regions.

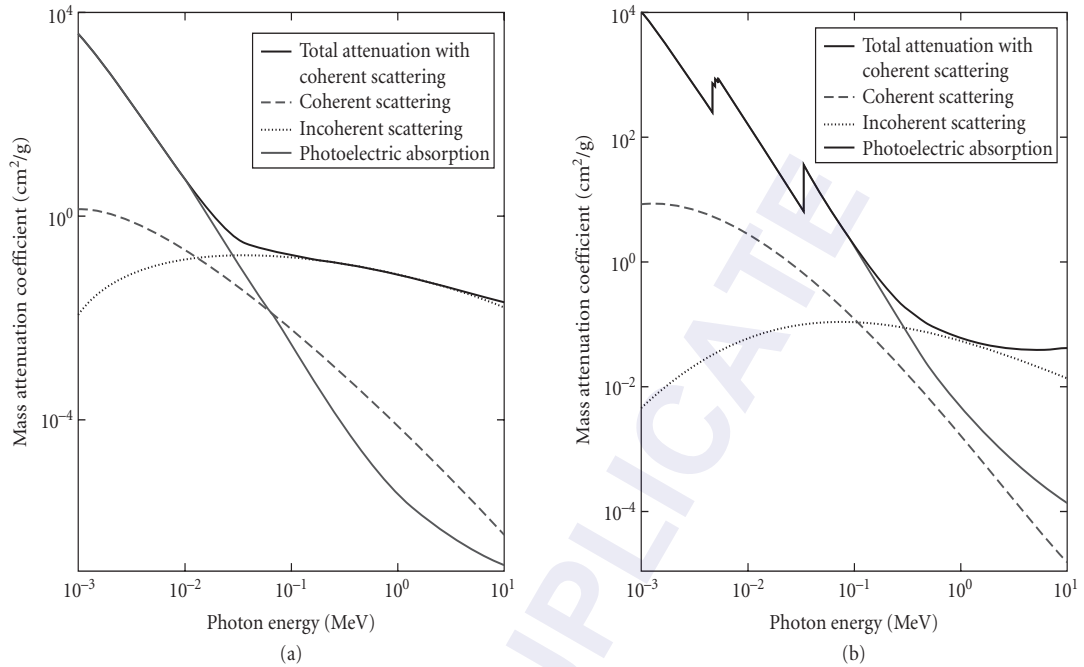


FIGURE 4 X-ray attenuation in (a) water and (b) iodine.

Because of these attenuation characteristics, the energy of the radiation used for imaging must be matched to the task. Roentgen made mention of this in his early work. Diagnostic x rays are typically generated via an x-ray tube having a tungsten anode, operated at an accelerating potential of 60 to 120 kV. A typical unfiltered spectrum is shown in Fig. 5. The spectrum displays the roughly 99 percent component that is bremsstrahlung radiation and the more intense but narrow characteristic peaks of

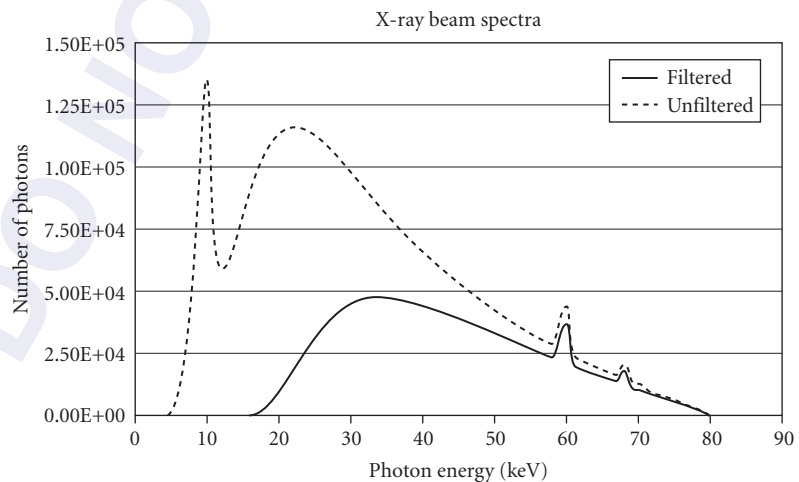


FIGURE 5 Unfiltered and filtered x-ray spectra from a tungsten anode x-ray tube operated at 80 kVp.

tungsten. The accelerating potential for this spectrum is 80 kV. The lower energy radiation is easily attenuated in tissue, therefore, contributing only to patient dose and not to the image. For this reason, medical-use x-ray beams are filtered at the x-ray tube, typically with several millimeters of aluminum. Such a filtered beam is shown in the lower curve of Fig. 5.

31.3 X-RAY DETECTORS AND IMAGE RECEPTORS

Image receptors used for medical and industrial imaging have undergone great development over the last century also. Starting from photographic glass plates, dedicated x-ray film emerged early in the twentieth century. Intensifying screens, which convert x-ray photons to light photons, were introduced for medical imaging soon thereafter, greatly reducing both patient and operator radiation dose. While film-based image receptors are still widely used in industrial and medical imaging, digital detectors came to the scene in the late twentieth century and are gaining broad acceptance. In 2008, digital mammography (breast x-ray imaging) replaced film-screen mammography at a rate of about 6 percent per month.³

The use of film as a radiographic receptor has the advantage of very high spatial resolution. Film alone has resolution greater than 20 lp/mm (20 line pairs per millimeter, equivalent to 25 μm); use of an intensifying screen reduces this to between 10 and 20 lp/mm for most systems. Industrial imagers, for which radiation dose is not a concern, continue to use film directly as the image receptor due to the high resolution.

Particularly for medical imaging, digital image receptors offer a number of advantages. While the dynamic range of film is on the order of 10^2 , for digital receptors the dynamic range is 10^4 or more. Further, the characteristic curve of film is highly nonlinear, leading to variability in the contrast of the image depending on the specific exposure reaching each point of the film. Contrary to this, the response of digital detectors is linear throughout the dynamic range, as shown in Fig. 6. Film serves as x-ray detector, display, and storage medium, meaning that compromises must be made to each role. Since data acquisition and display are separated for digital receptors, each stage can be optimized independently. With physical storage space always at a premium, digital images also have the advantage of requiring little physical space for each image. Multiterabyte storage systems have a small footprint, are relatively inexpensive, and hold large numbers of images.

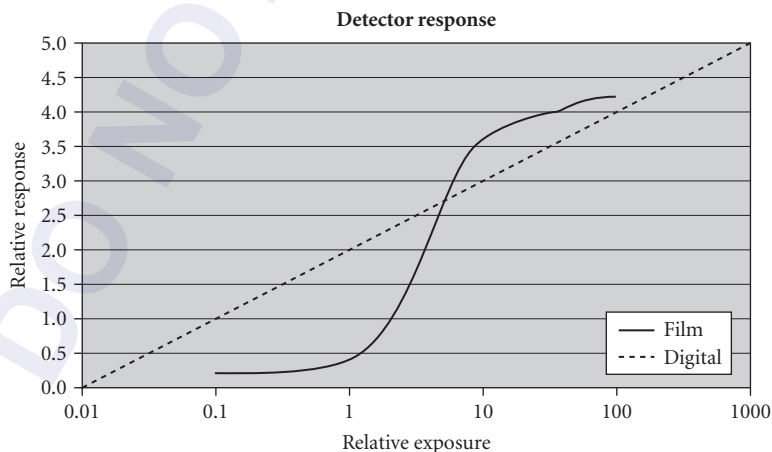


FIGURE 6 Comparison of the relative response of film compared to common digital detectors. Note the linear response and the much wider dynamic range of digital receptors compared to film.

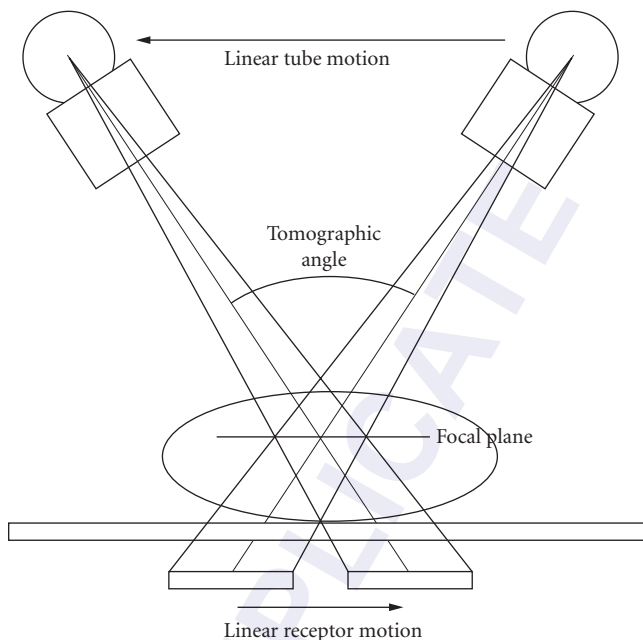


FIGURE 7 Linear tomography. As the tube moves one direction, the image receptor moves the other, with the fulcrum at the plane of interest.

31.4 TOMOGRAPHY

Tomography grew out of conventional x-ray imaging. In its simplest form, linear tomography, the x-ray tube and image receptor move in opposite directions about a fulcrum, which is placed at the height of the object of interest, as shown in Fig. 7. The paired motion of the tube and detector creates an image in which the structures above and below the fulcrum are blurred while leaving a sharp image at the level of the fulcrum. Due to the two-dimensional nature of the x-ray beam, the tomographic plane has a discreet thickness. The wider the tomographic angle, the thinner is the tomographic plane. In the mid to late twentieth century, more complicated paths were used to provide more complete blurring of the off-plane anatomy, but these devices became obsolete by the end of the century due to their bulk and dedicated use.

31.5 COMPUTED TOMOGRAPHY

In 1967, Sir Godfrey Newbold Hounsfield developed the first viable computed tomography scanner,⁴ while Allan McLeod Cormack developed a similar device in parallel. Both were awarded the Nobel Prize in Medicine in 1979 for their work.

Computed tomography (CT) is fundamentally different from conventional tomography. The concept is perhaps best understood by considering the original scanner, which was produced by EMI. This early unit used a pencil beam of radiation and a photomultiplier tube (PMT) detector. The beam and PMT were translated across the object being scanned and a single line of transmission data was collected. The beam and PMT were then rotated 1° and then translated back across the

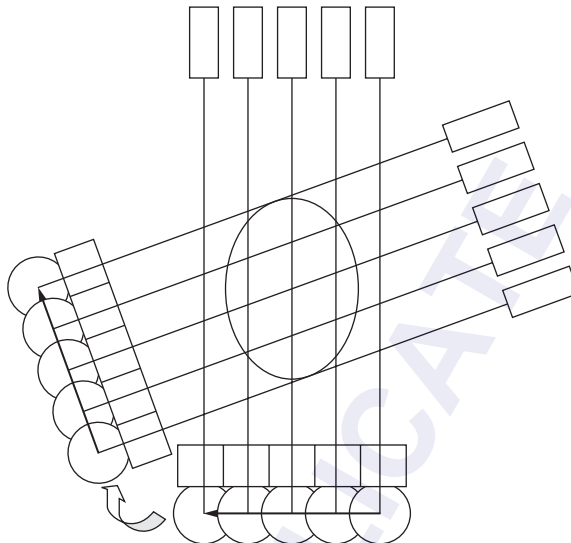


FIGURE 8 Simplified computed tomography data collection, as used in the first CT scanner. A pencil-shaped x-ray beam and detector translated across the object being imaged. This was repeated after the gantry was rotated by a small amount until the full data set had been collected.

patient, generating another line of transmission data as shown in Fig. 8. This process was repeated for 180° until the entire data set had been collected. This is known as “translate-rotate” geometry. Algebraic reconstruction techniques on a computer were then used to reconstruct an image depicting the attenuation of each matrix element. In the first scanner, the brain was divided into slices, each represented by a matrix of just 80×80 elements. Taking several hours to create, this coarse image, however, literally changed how physicians viewed the body.

Developments of this original scanner quickly ensued. More robust reconstruction algorithms, such as filtered back-projection, were created. This algorithm is not as sensitive to large discontinuities in the data, allowing for the technology to be applied not just in the head, but in the body as well. A fan beam of x rays and a linear array of detectors, each of which rotate around the object being scanned, known as “rotate-rotate geometry,” allowed for much more rapid data collection as shown in Fig. 9. Due to the fan beam, a full data set required $180^\circ + \text{fan beam angle} \approx 240^\circ$.

With the development of slip ring technology, the rotation speed of the x-ray tube and detector increased, allowing for improved imaging due to more rapid coverage of the anatomic volume. Increases in computer power and better reconstruction algorithms enabled the development of helical scanning, wherein the table supporting the patient translates through the scanner as the data collection takes place, as seen in Fig. 10.⁵ This has become the standard mode of operation for most medical CT imaging. Further advances in detector arrays, computing power, and reconstruction algorithms have led to multislice scanning. With multislice scanning, multiple data channels, representing multiple axial slices, are collected simultaneously. A continuous increase in the number of slices has been realized since the introduction of this technology, such that modern scanners have up to 320 slices, with 512 slices on the horizon.

One of the most important aspects of current multislice scanners is that the z-axis dimension of the image pixel is now reduced from approximately 12 to 0.625 mm or less. The implication of this is that scanners now provide a cubic voxel, or equivalent resolution in all three directions. With this key development, images may now be constructed in any plane desired. From a single volume of data, the physician may see axial, coronal, sagittal, or arbitrary angle images.

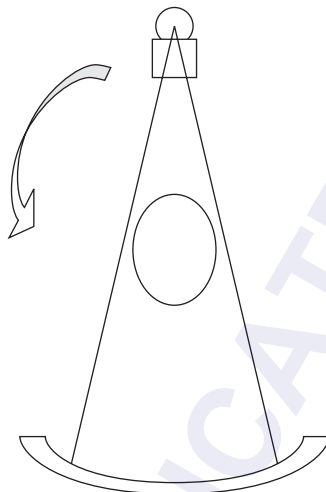


FIGURE 9 Modern computed tomography. A thin, fan beam of x rays is projected through the object and recorded by an array of detectors. The x-ray tube and detectors rotate continuously around the object.

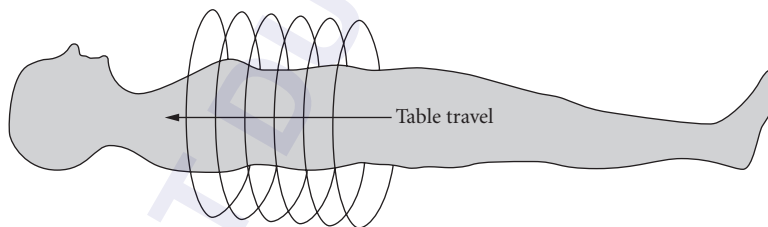


FIGURE 10 Helical computed tomography. The tube and detectors rotate continuously around the patient as the table is translated through the gantry.

Much work is going into the development of “cone-beam CT.” In these devices, the large number of detectors is replaced by a single two-dimensional panel detector with dimensions of tens of centimeters on a side. This is another promising approach to the acquisition of volumetric rather than slice-based data sets.

While most of the focus of computed tomography is on medical imaging, it has also found application in industrial imaging. The imaging principles remain unchanged, the geometry, however, may be adjusted. For example, systems are frequently configured such that the object being imaged is rotated while the imaging system remains fixed. Additionally, energies may range from the approximately 100 kV to over 10 MV, depending on the material to be inspected.

31.6 DIGITAL TOMOSYNTHESIS

The advent of digital projection radiography has led to the development of a technology known as *digital tomosynthesis*. Like conventional tomography, the x-ray tube is translated across or rotated around the object for a specified angle, although the image receptor may remain fixed. Typically,

several tens of images are acquired during the exposure. With this volumetric data set, an image can be reconstructed of any plane within the object. The observer may then scroll through these planes with the overlying anatomy removed. This allows for low contrast objects to be more readily visualized.

31.7 DIGITAL DISPLAYS

Of importance to all digital imaging is the display of the images. Most systems provide 12-bit images having 4096 shades of gray. Most computer monitors, even high-quality medical displays, display at 8 bits, or 256 shades of gray. The human eye is capable of discerning between 6 and 9 bits, or approximately 60 shades of gray. Mapping of the 4096 shades of gray in the image to the 60 shades usable by the observer is achieved through adjusting the window width and level. Figure 11a demonstrates

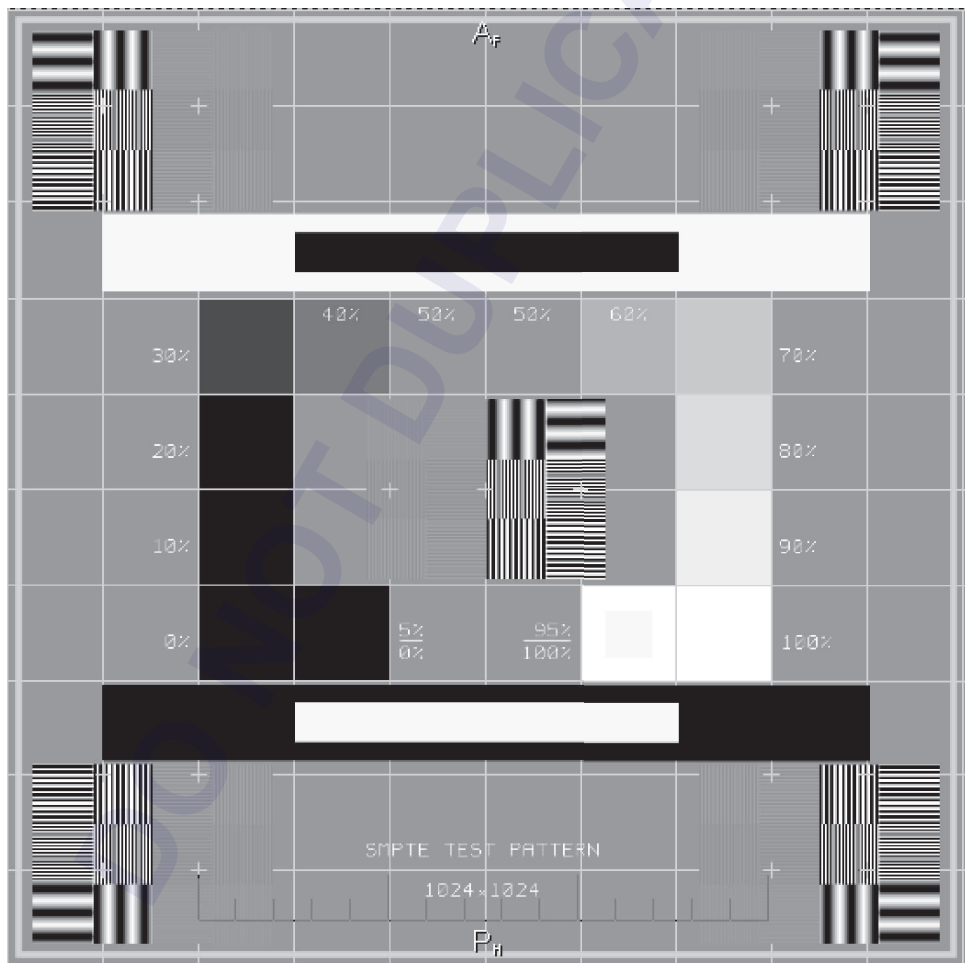


FIGURE 11 The SMPTE test pattern demonstrating (a) good contrast and good display calibration, (b) very high contrast and poor display calibration, and (c) very low contrast and poor display calibration.

good contrast. A narrow window width means that relatively few shades of gray are displayed, yielding a very high contrast image as seen in Fig. 11b. Increasing the width compresses the shades of gray in the original data, yielding a low contrast image, shown in Fig. 11c. The test pattern created by the Society of Motion Picture and Television Engineers (SMPTE), as shown in these images, is commonly used for the calibration of video displays.⁶

31.8 CONCLUSION

Progress in imaging technology has driven increases in the demands of the applications of that technology. Simple projection radiography will continue to be important to both medicine and industry. Computerized techniques, however, are giving additional tools for diagnosis to both. For example, the speed and image quality of multidetector helical computed tomography are allowing for screening of heart conditions that had been solely the realm of interventional procedures with their associated risks

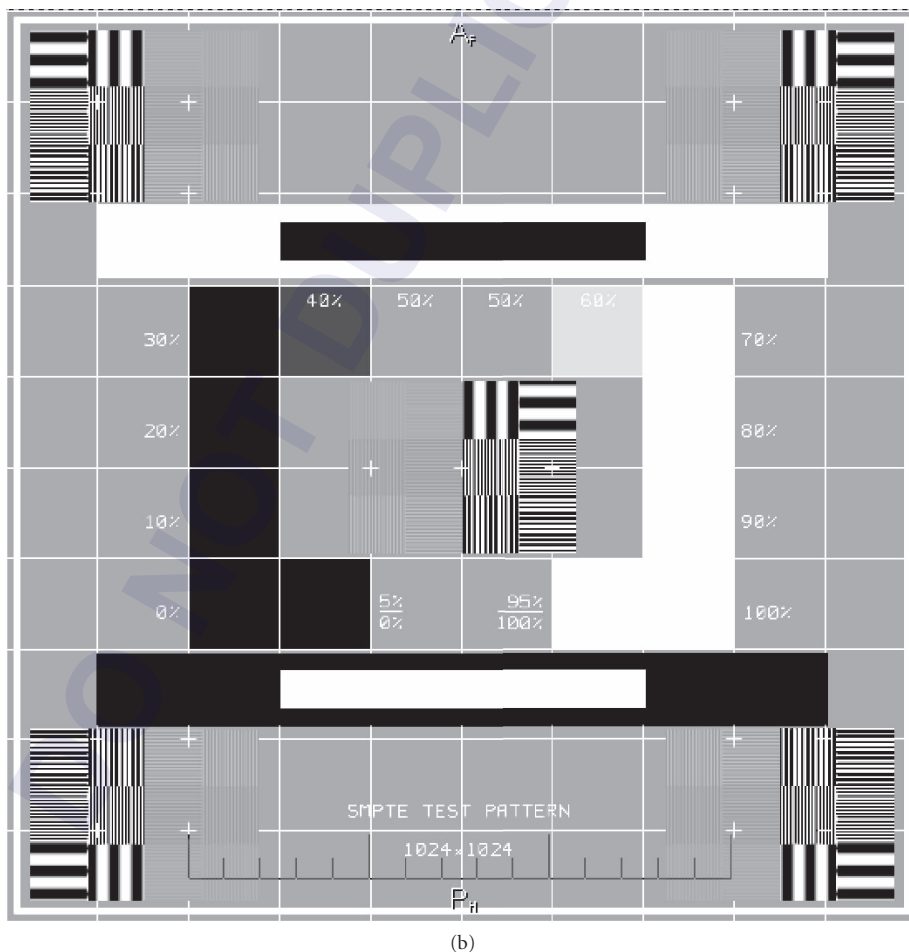


FIGURE 11 (Continued)

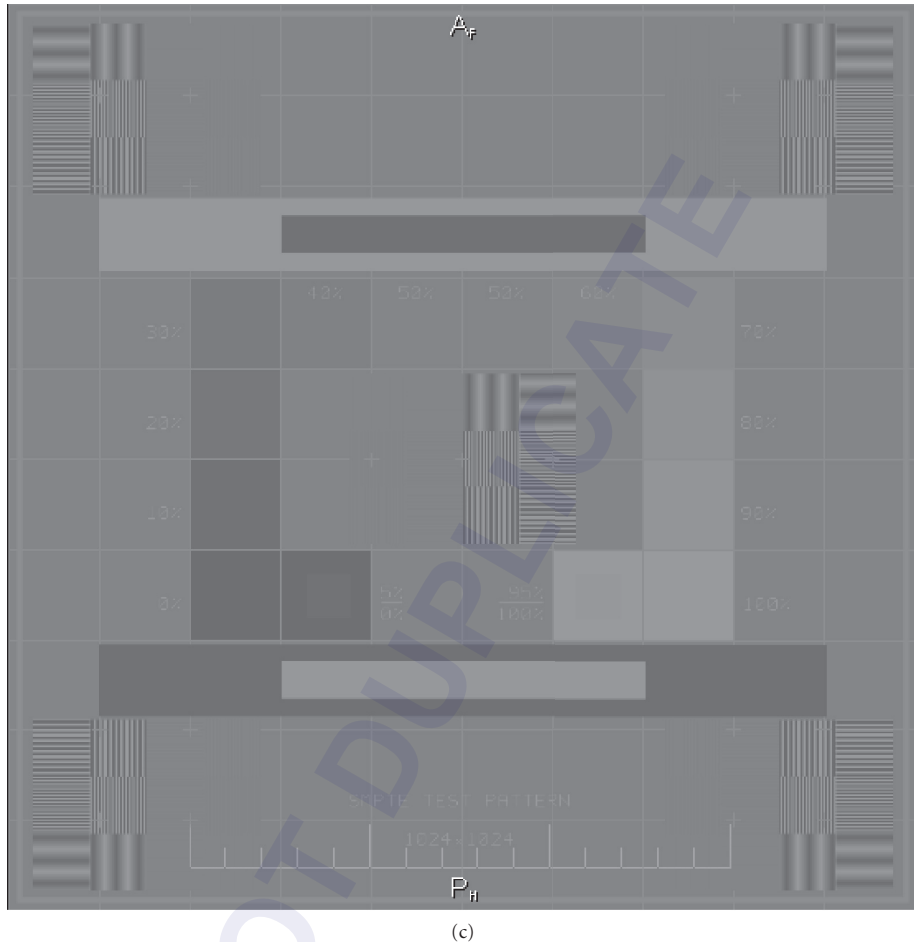


FIGURE 11 (Continued)

of complications. Similarly, traditional colonoscopic screening may be giving way to virtual colonoscopy through CT imaging.

Developments such as these show no sign of diminishment. However, for their continued progress, parallel development of x-ray sources, optics, and detectors must also continue.

31.9 REFERENCES

1. O. Glasser, *Wilhelm Conrad Roentgen and the Early History of the Roentgen Rays*, Springfield, IL, Charles C Thomas, 1934.
2. M. J. Berger, J. H. Hubbell, S. M. Seltzer, J. Chang, J. S. Coursey, R. Sukumar, and D. S. Zucker, "XCOM: Photon Cross Sections Database," NIST Standard Reference Database 8 (XGAM), <http://physics.nist.gov/PhysRefData/Xcom/Text/XCOM.html> (retrieved 5/8/2007).

3. P. Butler, April 2007, private communication.
4. G. N. Hounsfield, "Computerized Transverse Axial Scanning (Tomography): Part I. Description of System," *Br. J. Radiol.* **46**(552):1016–1022 (Dec. 1973).
5. W. A. Kalender, P. Vock, A. Polacin, and M. Soucek, "Spiral-CT: A New Technique for Volumetric Scans. I. Basic Principles and Methodology," *Rontgenpraxis* **43**(9):323–330 (Sep. 1990).
6. SMPTE RP133, "Specifications for Medical Diagnostic Imaging Test Pattern for Television Monitors and Hardcopy Recording Cameras," Society of Motion Picture & Television Engineers (SMPTE), White Plains, New York.

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

REQUIREMENTS FOR NUCLEAR MEDICINE

Lars R. Furenlid

University of Arizona
Tucson, Arizona

32.1 INTRODUCTION

Single Photon Emission Computed Tomography (SPECT) is a cross-sectional imaging modality that makes use of gamma rays emitted by radiopharmaceuticals (or equivalently, radiotracers) introduced into the imaging subject. As a tomographic imaging technique, the resulting images represent the concentration of radiotracer as a function of three-dimensional location in the subject. A related two-dimensional imaging technique is known as *scintigraphy*. The designation *Single Photon* in the SPECT acronym is used to distinguish the technique from a related emission tomography known as *Positron Emission Tomography* (PET) in which a correlated pair of annihilation photons together comprise the detected signal.

SPECT is one of the *molecular imaging* techniques and its strength as a diagnostic and scientific research tool derives from two key attributes: (1) properly designed radiotracers can be very specific and thus bind preferentially or even exclusively at sites where particular *molecular targets* are present and (2) gamma-ray signals in detectors originate from radioisotope tags on individual molecules, leading to potentially very high sensitivity in units of counts per tracer concentration. Virtually all SPECT systems are photon counting, i.e., respond to and record signals from individual photons, in order to extract the maximum possible information about the location and energy of each detected gamma ray.

Gamma rays are photons that are emitted as a result of a nuclear decay and are thus distinguished from x rays, which result from inner-shell electronic transitions, by the physical process they originate from. Gamma rays and x rays used for imaging overlap in energy ranges, with x rays as a rule of thumb covering the range between approximately 1 to 100 keV, and gamma rays covering the range of 10 keV up to several MeV. Most clinical SPECT imaging is performed with gamma rays (or secondary x rays) with energies of 80 keV (^{201}Tl), 140 keV ($^{99\text{m}}\text{Tc}$), 159 keV (^{123}I), or 171 keV and 245 keV (^{111}In). Preclinical imaging of murine species can be carried out with substantially lower energies, such as 30 keV (^{125}I), due to the smaller amount of tissue that needs to be traversed with a low probability of scatter or absorption.

An ensemble of SPECT radioisotopes emit their photons isotropically, i.e., with equal probability in all directions. In order to form a useful projection image of an extended source on a two-dimensional detector, an image forming principle, generally based on a physical optic, must be employed. The purpose

of the optic is to establish a relationship between locations in the object volume and pixels on the detector. In equation form, this can be expressed as

$$\bar{g}_m = \int h_m(\mathbf{r})f(\mathbf{r})d^3r$$

where $f(\mathbf{r})$ is the gamma-ray photon emission rate (which is proportional to tracer concentration) as a function of 3D position \mathbf{r} , $h_m(\mathbf{r})$ is the sensitivity of pixel m to activity in different regions of the object volume, and \bar{g}_m is the mean signal rate ultimately registered in pixel m . A variety of physical factors contribute to the three-dimensional *shape* (and *magnitude*) of the sensitivity functions, but the most important are the parameters of the image-forming optic, and the intrinsic resolution and detection efficiency of the detector.

All modern gamma-ray detectors convert the energy of individual gamma rays into electrical signals that are conditioned and digitized. Most clinical gamma cameras utilize an intermediate step in which the gamma-ray energy excites a burst of secondary lower-energy scintillation photons and are closely related to the scintillation camera designed by Hal Anger in the mid to late 1950s. The design comprises a relatively large slab of inorganic scintillation crystal which is viewed by an array of photomultiplier tubes whose signals are processed to estimate gamma-ray interaction location and energy. Achieved detector resolutions are on the order of 2- to 3-mm FWHM. Research and preclinical SPECT imagers are making increasing use of semiconductor detectors which directly convert gamma-ray energy into electron-hole pairs that migrate under the influence of a bias potential and induce signals in pixel or strip electrodes that can have dimensions down to approximately 50 μm . This is roughly the diameter of the region of space in which the energy of a gamma ray is deposited when it interacts with a solid in a complicated cascade of secondary photons and energetic electrons. SPECT imagers generally incorporate multiple cameras to increase the efficiency of tomographic acquisition, which can involve the measurement of typically 60 to 180 planar projections depending on the system and application.

Tomographic imaging almost always involves a reconstruction operation in which a collection of observations, projection images in the case of SPECT, are processed to recover an estimate of the underlying object. This can be expressed as

$$\hat{f}(\mathbf{r}) = O(\mathbf{g})$$

where O represents the reconstruction operation that acts on data vector \mathbf{g} , and $\hat{f}(\mathbf{r})$ the estimate of the object, generally in the form of activity in voxels. In SPECT, the reconstruction operation is currently most often either a filtered backprojection (FBP)¹ or an iterative statistical algorithm such as maximum-likelihood expectation maximization (MLEM)^{2,3} or ordered subsets expectation maximization (OSEM).⁴ The latter reconstruction methods incorporate a forward model of the imaging process that makes it possible to compensate for imperfections in the imaging system, especially if careful calibration measurements are made, and also make it possible to approximately account for absorption and scatter occurring within the object.

32.2 PROJECTION IMAGE ACQUISITION

Although it is possible and sometimes advantageous to perform SPECT imaging without ever forming intermediate projection images, using so-called listmode reconstruction methods that utilize unprocessed lists of gamma-ray event attributes,⁵ most SPECT systems do acquire projection images at regular angular intervals about the imaging subject. Virtually any of the physical processes that redirect or constrain light, such as absorption, refraction, reflection, and diffraction, can in principle be used to form an image in the conventional sense—mapping light from an object source point or plane to an image plane. The need to map an object volume to an image plane, and the requirement to function at the short wavelengths characteristic of gamma rays, place severe restrictions on the types and geometries of optics that can be considered and used successfully for SPECT.

Conventional SPECT systems use pinholes and parallel-hole collimators, or closely related variations such as converging or diverging collimators, slits, and slats, to form images. By permitting only the fraction of light traveling through the open area of the aperture along a restricted range of angles to reach a detector pixel, pinholes and parallel hole collimators define sensitivity functions that are nonzero in conical regions of the object volume. Parallel-hole collimators are the most commonly employed image-forming elements in clinical applications, in part because the projections they produce are the easiest to understand and process. Key features of parallel-hole collimators are that resolution degrades as a function of distance away from the collimator face while sensitivity is nearly unchanged. The parallel bores ensure that most gamma rays enter into the detector with nearly normal incidence angle, minimizing the parallax errors associated with uncertainty in depth of interaction in the detector. Parallel-hole collimators have no magnifying properties (the converging or diverging versions do) and require a tradeoff between distance dependent loss of resolution and overall sensitivity as governed by the collimator aspect ratio. There is a further design tradeoff between resolution loss from leakage between collimator bores and sensitivity loss from reduced fill factor of open collimator area.⁶

SPECT imager design with pinhole apertures also involves a set of design tradeoffs involving resolution, sensitivity, and field of view. Pinhole cameras have magnifications that depend on the ratio of pinhole to detector and pinhole to source distances. Since a three-dimensional object necessarily has a range of pinhole to source distances, resulting in different magnifications for different parts of the object, and the sensitivity cones intersect the object volume with different angles that depend on individual detector pixel locations, pinhole SPECT projections can appear complicated to a human observer. There are also further design tradeoffs between blur from leakage through the pinhole boundaries versus loss of sensitivity from vignetting at oblique angles. Nonetheless, pinhole and multipinhole apertures are currently providing the highest-resolution preclinical SPECT images, as reconstruction algorithms have no difficulty unraveling the geometric factors in the projections.

New optics, and systems, developed for SPECT should be evaluated in comparison to the conventional absorptive apertures discussed above, ideally with objective measures of system performance such as the Fourier crosstalk matrix.

32.3 INFORMATION CONTENT IN SPECT

Barrett et al.⁷ have suggested the use of the Fourier crosstalk matrix as a means of characterizing SPECT system performance. This analysis derives a measure of resolution, similar to a modulation transfer function (MTF), from the diagonal elements of the crosstalk matrix, and a second measure of system performance from the magnitude of the off-diagonal elements. An ideal system has significant amplitudes on the diagonal elements out to high spatial frequencies, representing high spatial resolution, and very small off-diagonal elements representing minimal aliasing between three-dimensional spatial frequencies (denoted by ρ_k). The crosstalk matrix analysis is best carried out with measured calibration data, but it can also be carried out with accurate system models. The diagonal elements,

$$\beta_{kk} = \sum_{m=1}^M \left| \int_S h_m(\mathbf{r}) \exp(2\pi i \rho_k \cdot \mathbf{r}) d\mathbf{r} \right|^2$$

are easily understood to represent the magnitude of the k th Fourier component of the overall system sensitivity function.

One of the primary motivations for optics and aperture development for SPECT is to increase the optical efficiency. Conventional pinholes and parallel hole collimators image by excluding most of the photon distribution function, the phase-space description of photon trajectories that encodes both locations and directions. They therefore are very inefficient at collecting light. Since SPECT radiopharmaceuticals need to adhere to the “tracer” principle, i.e., be administered in limited amounts,

SPECT images almost always have relatively low total counts per voxel compared to many imaging techniques, and therefore significant Poisson noise. New system designs with large numbers of pinholes are helping to address these concerns.⁸

32.4 REQUIREMENTS FOR OPTICS FOR SPECT

Image-forming elements for SPECT systems need to meet criteria to be useful for imaging that can be formulated as a set of requirements. One of the primary requirements is that the fraction of gamma rays that are able to reach the detector via undesired paths is small relative to the fraction that traverses the desired optical path. For example, for pinholes, the open area of the pinhole must be large relative to the product of the probability of gamma rays penetrating the shield portion of the aperture and the relevant area of that shield. Interestingly, as the number of pinholes increases, the tolerance to leakage also increases, which is an aid in system design. Given the relatively high energies of the gamma rays emitted by the most useful SPECT radioisotopes, this is not a trivial requirement to meet. A comparable condition for parallel hole collimators is that only a small fraction of the photons reach the detector surface after having traversed (or penetrated) at least one bore's septal wall.

The overall system sensitivity function, which represents a concatenation of all of the sensitivity functions in all of the projections in the tomographic acquisition, needs to have significant values for high spatial frequencies in the diagonal elements of the Fourier crosstalk matrix, but have relatively small off-diagonal terms. A variety of conditions can break this requirement. For example, certain geometric features in the sensitivity functions, such as might arise from a symmetric arrangement of multiple pinholes with respect to the imager axis, can introduce correlations that result in enhanced off-diagonal elements. Significant holes in the sensitivity in regions of the object space can also be problematic, as can issues whenever there is an uncertainty as to which of several alternative paths a photon took to arrive at the detector.

Since SPECT imaging generally involves three-dimensional objects, the image-forming element needs to work with photons emitted from a volume of space consistent with the dimensions of the objects or subjects being imaged. In some very specialized applications, for example, in preclinical scanning applications involving known locations of uptake in limbs, the volume required can be small, on the order of 1 cm³ or less. But for most general preclinical imaging, the object volume is measured in units of tens of cubic centimeters. In clinical applications, the object volume is often measured in units of hundreds of cubic centimeters. For practical application with living imaging subjects, SPECT acquisitions should not exceed roughly 30 to 60 minutes in total time, which makes it impractical, or at least undesirable, to employ optics that require rastering to build up single planar projections.

The materials used for gamma-ray optics can sometimes cause problems with secondary fluorescence, depending on the energy spectra of interest. For example, collimators made of lead give off fluorescence k_{α} and k_{β} photons at around 75 and 85 keV, respectively, that can be hard to distinguish from the 80 keV emissions from ²⁰¹Tl labeled tracers. Care in geometric design can often minimize the potential for collimator fluorescence reaching the detector.

32.5 REFERENCES

1. P. E. Kinahan, M. Defrise, and R. Clackdoyle, "Analytic Image Reconstruction Methods," in *Emission Tomography*, M. N. Wernick and J. N. Aarsvold, eds., chap. 20 Elsevier, San Diego, 2004.
2. L. A. Shepp and Y. Vardi, "Maximum Likelihood Estimation for Emission Tomography," *IEEE Trans. Med. Imaging* **1**:113–121 (1982).
3. K. Lange and R. Carson, "EM Reconstruction Algorithms for Emission and Transmission Tomography," *J. Comput. Assist. Tomogr.* **8**:306–316 (1984).

4. H. M. Hudson and R. S. Larkin, "Accelerated Image Reconstruction Using Ordered Subsets of Projection Data," *IEEE Trans. Med. Imaging* **13**:601–609 (1994).
5. H. H. Barrett, T. White, and L. C. Parra, "List-Mode Likelihood," *J. Opt. Soc. Amer. A* **14**:2914–2923 (1997).
6. D. L. Gunter, "Collimator Characteristics and Design," in *Nuclear Medicine*, R. E. Henkin, M. A. Boles, G. L. Dillehay, J. R. Halama, S. M. Karesh, R. H. Wagner, and A. M. Zimmer, eds., Mosby, St. Louis, MO, 96–124 (1996).
7. H. H. Barrett, J. L. Denny, R. F. Wagner, and K. J. Myers, "Objective Assessment of Image Quality. II. Fisher Information, Fourier Cross-Talk, and Figures of Merit for Task Performance," *J. Opt. Soc. Amer.* **12**:834–852 (1995).
8. N. U. Schramm, G. Ebel, U. Engeland, T. Schurrat, M. Behe, and T. M. Behr, "High-Resolution SPECT Using Multipinhole Collimation," *IEEE Trans. Nucl. Sci.* **50**(3): 315–320 (2003).

This page intentionally left blank.

DO NOT DUPLICATE

REQUIREMENTS FOR X-RAY ASTRONOMY

Scott O. Rohrbach

*Optics Branch
Goddard Space Flight Center, NASA
Greenbelt, Maryland*

33.1 INTRODUCTION

X-ray astronomy can be generally defined as the observation of x-rays from extraterrestrial sources. A variety of methods are used to determine, with more or less accuracy, where extraterrestrial photons come from, but the most common are simple collimation, coded aperture masks and focusing mirror systems. Collimator systems are used to achieve imaging resolution down to approximately 1° , and are independent of energy. Coded aperture masks, such as the system on the Swift Burst Alert Telescope (BAT), can localize the direction of strong signals to within a few minutes of arc, and are similarly energy independent. Finally, true imaging systems can achieve sub-arc-second angular resolution, the best demonstration of which is the Chandra X-Ray Observatory, with 0.5 arc-second half-power-diameter (HPD). While the technical details of space-based x-Ray observatories are covered in Chap. 47, the various requirements and trade-offs are outlined here. For x-ray imaging systems, a large number of factors come into play, including

1. Imaging resolution
2. Effective collecting area
3. Cost
4. Focal length
5. Field of view
6. Mass

For missions based on targeted objects, the first three are usually the primary factors and trade offs must be made between them. For survey missions, where the telescope is constantly scanning the sky, field of view also comes into play, but then imaging resolution is not as critical. For either type of observatory, the telescope mass is defined more by the mission mass budget, which is set by the choice of launch vehicle, than any other factor. On top of these mirror requirements are considerations of what detector system is being used, since a higher quality (and thus more expensive) imaging system coupled to a low-resolution detector would be a waste of money.

Current technological limitations in optics and detectors mean that the x-ray regime is split into two observational regions, referred to as the soft and hard x-ray bands. The soft band spans the 0.1 to ~ 10 keV range, while the hard band spans the ~ 10 to ~ 300 keV range. In the soft band, detector efficiencies and filters required to absorb visible light effectively define the low energy cutoff, while the throughput of focusing systems and the average critical angle of total external reflection for simple thin films (see Chap. 26) define the high energy cutoff. In the hard band, the upper limit is generally determined by the quantum efficiency of the detector system being used. Due to the focal length limitations for spacecraft (~ 10 m without an extendable optical bench) and the fact that multilayer coatings that can extend the bandpass of focusing optics have only recently become available, hard x-ray observation on space-based observatories has been limited to simple collimation and coded-aperture mask instruments.

33.2 TRADE-OFFS

As noted above, the most obvious competing factors are imaging resolution, effective area, and cost. Both imaging quality and effective area can be limited by the available financial budget the mission has. Other trade-offs exist. For example, one consideration is mirror substrate thickness. Thicker mirror substrates allow for conventional grinding and polishing, yielding excellent imaging quality. But all of the projected area occupied by the cross section of each substrate means less effective area for the given optical aperture of the telescope. Also, grinding and polishing of each surface can be very costly. On the other hand, the desire for large effective collecting area drives the thickness of each substrate down, in order to minimize the optical aperture obscured by the substrate. This leads to substrates that are effectively impossible to individually polish, both due to time and cost constraints as well as the problems of polishing very thin substrates without introducing print-through errors or making the flexible substrates weaker. And, thinner, more flexible substrates are harder to integrate into a flight housing without introducing significant figure distortions.

For each mission concept that gains significant community and financial support, a team of astronomers studies the various performance trade-offs with respect to the science the mission will pursue. Will the primary targets be bright sources or dim, or a mixture? Will they be compact or span many minutes of arc? Will they change quickly like the rapid oscillations of a pulsar, or be more static like the gas between stars and galaxies? The answers to these questions depend to some degree on what the most compelling scientific questions of the day are. Once the scientific scope of the mission has been defined, the functional requirements begin to take shape, but there are still many other questions to be answered. High-resolution spectroscopy requires more collecting area than pure imaging, but a telescope with good imaging quality also reduces the background inherent in any measurement, improving the spectroscopic results. And since even the most modest observatory will cost more than \$100 million to design, launch, and operate, any new observatory must represent a significant improvement over previous missions in at least one area, be it imaging quality, collecting area, spectral resolution, bandpass, timing accuracy, and so on.

In a short review such as this, it would be impossible to list all of the various targets, science areas, and trade-offs that need to be made for any particular observatory. Instead, a short case study, using three current missions and one upcoming mission study, is presented. Each observatory has different goals and budgets, and as such, approaches the question of optimization from a different point of view. The missions are summarized in Table 1.

At one corner of the imaging/effective area/cost trade space is the Chandra X-Ray Observatory, which was primarily designed to yield the best x-ray imaging ever achieved, and is also the most expensive x-ray observatory ever made. It uses four nested shells with the largest being 1.4 m in diameter, and each of the primary and secondary pairs is figured out of thick monolithic shells of Zerodur, nearly 1 m long. The emphasis on imaging quality translated into a small number of mirror elements that could be figured and polished precisely, yielding a point spread function of 0.5 arc-second HPD, even after all four shells were aligned to one another. During the conceptual development of the mission, the observatory included a microcalorimeter detector that would not only yield high-quality

TABLE 1 Comparison of Mirror Assemblies of Various X-Ray Observatories

Observatory, Approximate Cost	Effective Area @ 1, 6 keV, cm ²	Angular Resolution (arc-second, half-power-diameter)	No. of Shells per Telescope, Maximum Diameter (m), f-number	Mirror Substrate Material	Total No. of Mirror Elements in Observatory
Chandra \$2.5B	800, 300	<0.5	4, 1.4, f/7	16–24 mm Zerodur	4
XMM-Newton \$640M	4,500, 2,700	12–15	58, 0.7, f/11	0.5–1-mm electroless nickel	174
Suzaku \$140M	1,750, 1,200	110	180, 0.4, f/11	0.17-mm segmented aluminum	5,760
Constellation-X ~\$2.5B	15,000, 6,000	5–15	~200, 1.4, f/7	0.4-mm segmented glass	12,000–15,000

imaging, but unparalleled energy resolution in a nondispersive system. But budget restrictions led to the use of a less expensive grating array in order to achieve the spectroscopic energy resolution required. Budget restrictions also meant that only four out of the originally envisioned six mirror shells could be manufactured and flown. The focal length of the system was chosen to be 10 m—to give an f-number of approximately f/7. This is a relatively fast compared to most x-ray observatories, but being an imaging mission, there was also a need for a wide field of view in order to fully capture extended sources in a single exposure.

During the development of the XMM-Newton observatory, the science team stressed effective collecting area over imaging. To quote the XMM-Newton web page,* “The design of the optics was driven by the requirement of obtaining the highest possible effective area over a wide range of energies, with particular emphasis in the region around 7 keV.” In contrast to Chandra, XMM-Newton was designed with 58 nested shells per telescope, and has three identical telescopes on board the spacecraft. That, in conjunction with a much slower f/11 optical system, means a higher effective collecting area—6 times more at 1 keV and nearly 10 times more at 6 keV—due to the shallower grazing angles improving the reflectivity at the higher energy range. But, the drawback is that each mirror shell needed to be between 0.5 and 1 mm thick in order to get them all to fit in the spacecraft volume, compared to the 16- to 24-mm-thick shells on Chandra. Figuring and polishing so many thin substrates would have taken far too long and cost too much, so a replication process was used to manufacture the shells. Fifty-eight thick negative masters, termed “mandrels,” were figured and polished with their outer surface matching the desired inner surface of each shell. The masters were then coated with a gold release layer, and electroless nickel was deposited on top of the gold, forming the mirror shells. The mirrors and mandrels were then separated through a cryogenic process that did not distort the mirror shells. But, since each mirror shell is so thin, and the time and budget available for figuring each mandrel was not unlimited, the imaging quality of final telescopes range from 12 to 15 arc-second HPD.

The Suzaku observatory is the third case to be examined. It is a small mission, with a relatively limited budget. The primary goal of this mission was to perform high-resolution spectroscopy by flying five telescopes with upward of 200 nested shells. Four of the mirror assemblies use CCD detectors as their imaging focal plane, similar to both Chandra and XMM-Newton, but the fifth uses a microcalorimeter to achieve unprecedented 12-eV resolution in a nondispersive system. Due to the relatively small volume of the spacecraft, each of the five mirror assemblies is contained within a 40-cm diameter. Each shell is made up of four 90° segments of revolution, and the reflecting substrates are 0.17-mm-thick aluminum. Obviously, it would be impractical to figure each of thousands of very thin mirrors, and with the budget available it would even be too costly to figure the forming

*http://xmm.esac.esa.int/external/xmm_user_support/documentation/technical/Mirrors/index.shtml

mandrels used to thermally shape the aluminum substrates, so a conical approximation to the Wolter-I geometry is used. This fact and the deformations associated with mounting and integration of the very thin aluminum, result in poorer imaging quality (arc-minute scale HPD), but the system yields more collecting area as Chandra in a significantly smaller package and at $\sim 1/20$ th of the cost.

Finally, the next large x-ray observatory under development at NASA is Constellation-X. This mission will concentrate on high-resolution spectroscopy in the soft x-ray band. Proposed “science enhancement packages” include more sophisticated multiple layer coatings to increase the low-energy effective area, dedicated hard x-ray telescopes, various grating arrays to improve low-energy spectral resolution, and so on. The core mission, however, includes four large Wolter-I telescopes, each 1.4 m in diameter coupled to imaging microcalorimeter arrays. The mission is designed to be a revolutionary step in capability, as opposed to an evolutionary step, as it will be 3 to 4 times larger than any previous mission and use an imaging detector with >30 times better energy resolution than any previously successfully flown. The primary trade-off is more effective collecting area, at the expense of poorer angular resolution than missions with individually polished mirror shells. In light of the science results achieved by Chandra, however, at the time of this writing, the science team is strongly considering changing the imaging requirement from 15 to 5 arc-second HPD. Whether any of the other requirements are adjusted to compensate for the tighter imaging requirement is still under debate.

33.3 SUMMARY

As can be seen in Table 1, there are a variety of competing restrictions in building a large x-ray observatory. Among the science questions that drive the design of any telescope are the breadth, signal strength, energy spectrum, and variability of the observation targets in question. Within the currently envisioned upcoming observatories^{*,†,‡,§} the angular resolution requirement varies from 2 arc-seconds to 1 arc-minute. The energy spectra of interest are as low as 0.25 to 10 keV and as high as 6 to 80 keV. The resulting telescopes have focal distances ranging from 10 to 30 m, and have 600 to 15,000 cm² of effective collecting area. From these contrasts, it is clear that there is no one-size-fits-all solution to the questions of telescope imaging quality, field of view, energy resolution, and so on, especially when the various factors of cost, launch vehicle mass budget, energy budget, volume, and schedule are taken into account. The number of variables than can be traded-off against one another mean that for each mission concept, a significant effort needs to be made to balance all of these factors in a way that results in a telescope tailored to the science goals at hand and within the given budgets. A very short list is given above that demonstrates the wide range of mission concepts, budgets and capabilities. But, while any space flight mission is costly, as a result of careful configuration trade-off studies, all are valuable to the astronomical community.

^{*}Simbol-X—<http://www.asdc.asi.it/simbol-x/>

[†]XEUS—<http://sci.esa.int/science-e/www/area/index.cfm?fareaid=103>

[‡]Constellation-X—<http://constellation.gsfc.nasa.gov/>

[§]<http://www.nustar.caltech.edu/>

EXTREME ULTRAVIOLET LITHOGRAPHY

Franco Cerrina and Fan Jiang

*Electrical and Computer Engineering & Center for NanoTechnology
University of Wisconsin
Madison, Wisconsin*

34.1 INTRODUCTION

The evolution of the integrated circuit is based on the continuous shrinking of the dimensions of the devices, generation after generation of products. As the size of the transistors shrink, the density increases and so does the functionality of the chip. This trend is illustrated in Table 1, derived from the Semiconductor Industry Association (SIA) forecasts of the Semiconductor Technology Roadmap.¹ During the fabrication process, optical imaging (optical lithography [OL]) is used to project the pattern of the circuit on the silicon wafer. These are exceedingly complex images, often having more than 10^{12} pixels. OL provides very high throughput thanks to the parallel nature of the imaging process with exposure times of less than 1 sec. In semiconductor manufacturing the ArF laser line (at 193 nm) is the main source of radiation used in OL. The resolution of the imaging systems has been pushed to very high levels, and it is very remarkable how it is possible to image patterns on large areas with dimensions smaller than 1/5th of the wavelength: Rayleigh's criterion ($\Delta \approx \lambda/2NA$) can be defeated by using super-resolution techniques, commonly referred to as resolution enhancement techniques (RET).² This is possible because in patterning super-resolution techniques can be applied more effectively than in microscopy or astronomy.³ All of this—throughput and resolution—makes OL today's dominant technology. However, even with the best of wavefront engineering it is very unlikely that features as small as 19 nm and less can be imaged with 193 nm—ultimately, the wavelength is the limiting factor in our ability to image and pattern. Hence, a shorter wavelength is needed to continue to use optical patterning techniques, or some variation thereof based on the paradigm of projecting the pattern created on a mask. For various reasons the 157-nm line of the F_2 laser is not suitable, and no other effective sources nor optical solutions exist in wavelength region from 150 to 15 nm—the vacuum ultra violet region, today often called extreme ultra-violet (EUV). Notably, in this part of the spectrum the index of refraction of all materials has a very large imaginary part due to the photoexcitation of valence band states. Thus all materials are very absorptive, and as a consequence transmission optics cannot be implemented. This has far-reaching implications for lithography.

Extreme ultraviolet lithography (EUV-L) is one of the latest additions to the list of patterning techniques available to the semiconductor industry;^{4,5} it aims to maintain the basic paradigm of pattern projection in the context of EUV-based technology. Table 2 lists some of the main goals of these imaging systems.

TABLE 1 Critical Dimensions of Devices in Integrated Circuits

Year	2008	2010	2012	2014
CD (nm)	38	30	24	19
Wavelength (nm)	193	193	13.4	13.4

TABLE 2 EUV Lithography Exposure Tool Parameters

Generation	25 – 12-nm CD
Wavelength	13.4 nm
NA	0.3 – 0.4
Field size	10–25 mm
Power density	10 mW/cm ²
Wafer size	300 mm

Because it is not possible to use any form of transmission optics (i.e., lenses) the imaging systems have to rely on mirrors. Very-high-resolution optical systems based on near-normal incidence designs are possible that use at least 4 to 6 reflecting surfaces,⁶ with aspheres further improving imaging. However, the reflectivity of a surface in the VUV is of the order of 10^{-3} to 10^{-4} , too small by order of magnitudes. As pointed out by Spiller,⁷ it is possible to use multilayer coatings to enhance the reflectivity even if the materials are highly absorbing. Indeed, the choice of the wavelength of 13.4 nm for EUV-L is dictated by the availability of efficient multilayer reflectors. These are essentially $\lambda/4$ stacks of molybdenum and silicon, producing a reflectivity higher than 70% in the region of $\lambda = 10 - 15$ nm. As mentioned above, in the VUV all materials are strongly absorbing, so that the propagation in the stack is limited to 40 to 60 bilayers,⁷ and this in turn limits the maximum reflectivity of the stack. Contrary to $\lambda/4$ dielectric mirrors, the bandpass is fairly large ($\sim 10\%$) and the maximum reflectivity does not go beyond 70 to 75%. Thus, the optical system is very lossy, with a transmission of $T \sim 0.75^6 = 0.18$, and high-power sources are required to achieve the final power density. Additionally, the multilayers must have extremely smooth interfaces to reduce incoherent scattering that worsens the quality of imaging by introducing flare.⁸ We note that while glancing angle mirrors can be used very effectively, the aberrations of these strongly off-axis systems forbid their use in large area imaging.

34.2 TECHNOLOGY

Like all imaging systems, EUV lithography requires a source of radiation, a condenser to uniformly illuminate the mask, a mask with the layout to be exposed, imaging optics and mechanical stage for step and repeat operation.^{4,5,9}

- While synchrotrons are by far the most efficient sources of EUV radiation, for manufacturing a more compact source based on the emission from plasma is more appealing. There are two main types of sources, one based on electrical discharge in a low-pressure gas and the other based on a laser-initiated plasma.¹⁰
- The condenser is a fairly complex optical system designed to match the source phase space to the optics. It includes both near-normal incidence mirrors (Mo-Si) and glancing mirrors. Since the sources are still relatively weak, it is important that the condenser collects as large a fraction as possible of the source phase space.¹¹

- The mask is formed by a patterned absorber (equivalent to 20- to 40-nm Cr) deposited on a multilayer stack (60 to 80 bilayers) formed on a thick quartz or zerodur substrate. The pattern is $4\times$ the final image, thus relaxing somewhat the constraints on the resolution and accuracy of the layout.
- The imaging optics include a set of near-normal incidence mirrors. In the first designs, the number of mirrors was 4, becoming 6 in more recent designs to increase the final NA of the system. Designs with up to 8 mirrors are also being considered. The first generation of tools operated at 0.25 NA, and the second generation works at 0.3 NA, with a goal of reaching NA ~ 0.4 .^{5,6,12}
- The step-and-scan stage is an essential part of the exposure system, since the relatively small field of view of the optics must be repeated over the whole wafer surface. The wafer is typically held in place by an electrostatic chuck, and the motion is controlled by laser interferometers.⁹
- The whole assembly is enclosed in a high-vacuum system, as mandated by the high-absorption coefficient of EUV radiation in any media, including air.
- The photoresist is the material used to record the image projected on the wafer. These materials are very similar to those already in use in OL, and are based on the concept of chemical amplification to increase their sensitivity to the exposing radiation.¹³

Figure 1 shows such a tool, developed at Sandia National Labs¹⁴ and installed at the Advanced Light Source, Berkeley. Figure 2 shows a photograph of one of the ASML Alpha Tools, installed at SUNY-Albany and at IMEC.⁸ Notice the vacuum enclosure and the overall large size of the lithography system. The ASML tool uses a plasma source. These tools are used to study advanced devices, and to develop the processes required for device manufacturing.

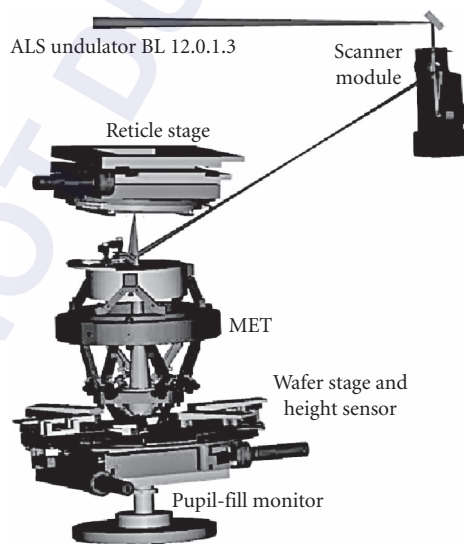


FIGURE 1 EUV exposure tool. The design includes 4 mirrors, and the mask (reticle) and wafer location. The whole system is under vacuum, and the condenser optical system is not shown.¹⁴ The whole system is enclosed in a vacuum chamber, not shown. The overall size is of several meters. This tool is installed at the Advanced Light Source, Berkeley. (See also color insert.)



FIGURE 2 ASML alpha demo tool during installation. The whole system is under vacuum, and a plasma source is used to generate the EUV radiation (not shown).⁸ (See also color insert.)

EUV-Interferometric Lithography (EUV-IL)

In addition to the imaging steppers described above, there are other ways of creating very-high-resolution patterns, in particular for materials-oriented studies where simple structures are sufficient. Interference can be used to produce well-defined fringe patterns such as linear gratings, or even zone plates—in one word, to pattern.^{15–18} Interferometric lithography (IL) is the process of using the interference of two or more beams to form periodic patterns of fringes to be recorded in an imaging material—originally a photographic film, today often a photoresist. IL is useful for the study of the properties of the recording material because it forms well-defined, simple, and high-resolution periodic fringe patterns over relatively large areas. The recording properties of a photoresist material can be studied by measuring its response to a simple sinusoidal intensity distribution; a series of exposures of different period allows us to simply and quickly determine its ability to record increasingly finer images.¹⁹ From this, we can project its ability to later resolve more complex and nonperiodic patterns. IL allows the study and optimization of resist materials without the need of high-resolution optics and complex masks—the diffraction provides the high-resolution fringe patterns needed. These arguments led us to the development of high-resolution EUV-IL as a platform for the development of the materials needed for nanolithography, well beyond the reach of the imaging ability of the current generation of imaging tools (steppers).^{16,17}

In Fig. 3 we present some of the high-resolution exposures from our system.^{16,17} One of the advantages of EUV-IL is that the interference pattern period is half that of the original grating, thus considerably facilitating the fabrication process—a 25-nm period (12.5-nm lines and spaces) is generated from

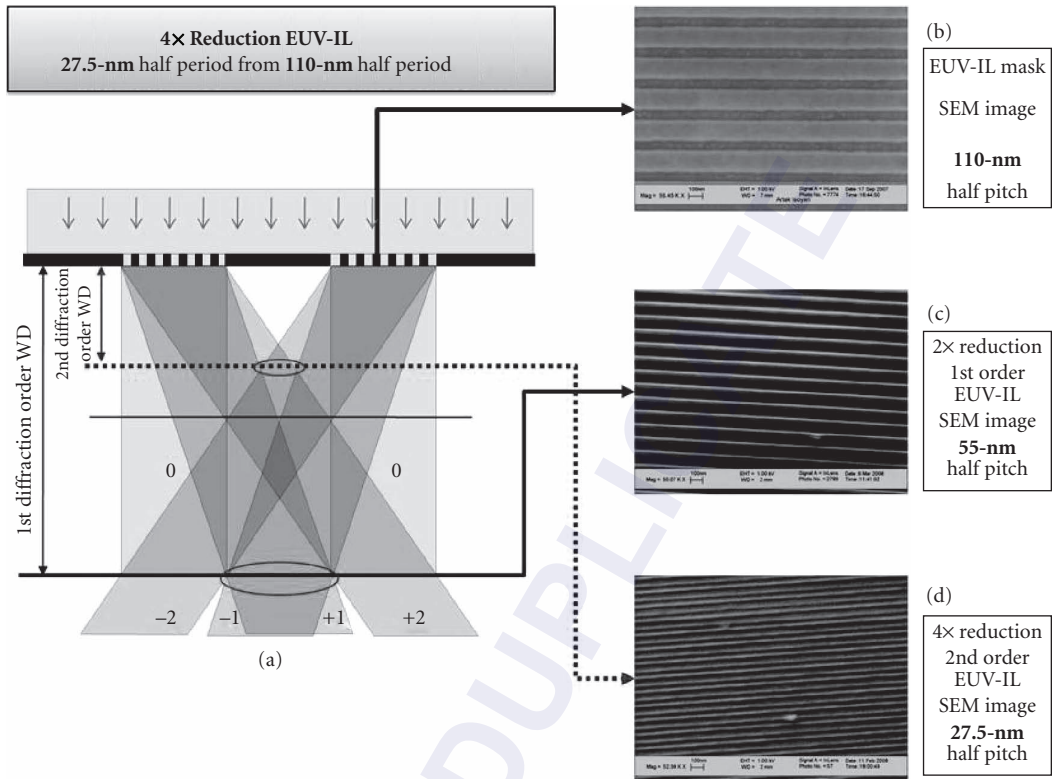


FIGURE 3 EUV interferometric lithography. The diffraction gratings are illuminated by a synchrotron, and the diffracted beams interfere as shown. The beams overlap creating 1st and 2nd order interference patterns of excellent visibility. Right, SEM images of the grating, and of the first- and second-order exposures. Notice the relative period of the images—the 1 \times period is half of that of the diffracting grating, and the 2 \times is 1/4.¹⁷ (See also color insert.)

a 50-nm period grating. The use of second-order interference increase this leverage by another factor of 2, so that a 50-nm period grating could in principle generate 12.5-nm periodic fringes, that is, 6.25-nm lines and spaces without the need of complex optical systems.

34.3 OUTLOOK

The continuing evolution of optical lithography has pushed EUV farther in the future. At the time of writing (2009), industry sources do not expect that EUV will be needed until the 17-nm node.^{4,5} Among the most active developer of systems are Nikon (Japan) and ASML (Belgium). There are several issues that are still unsolved, or only partially solved. Specifically:

- **Sources**—The power delivered by current plasma sources is too low by at least one order of magnitude. The photoresist should have a sensitivity of 5 mJ/cm², thus requiring at least 1 to 5 mW/cm² of power delivered to the resist surface.¹⁰
- **Source debris**—Since the EUV radiation is generated from hot plasmas (with the exception of synchrotrons) the elimination of reduction of debris ejected by the plasma is a major concern.

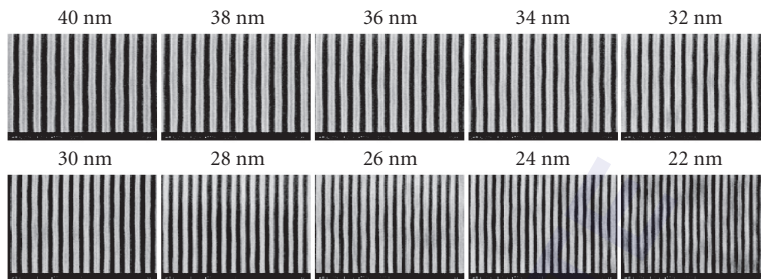


FIGURE 4 Twenty-two-nm dense lines printing of chemically amplified resists achieved at the Advanced Light Source facility. High-quality imaging is demonstrated, but line-edge roughness appears on the higher-resolution patterns.¹⁴

- **Masks**—Much progress has been made in this area, but the issue of defects has not yet been fully solved. Specifically, defects submerged in the reflecting stack underlying the absorber pattern create phase errors that affect severely the image.²⁰ These “buried defects,” only a few nanometer in size, are very difficult to detect short of a full-mask at-wavelength inspection, an expensive and time-consuming task.
- **Photoresist materials**—There is as yet no resist material that can balance the contrasting requirements of sensitivity, resolution, and low-edge roughness.²¹ As shown in Figs. 3 and 4 at the highest resolution the photoresist lines become less smooth, showing a degree of graininess (line edge roughness [LER]) that makes the pattern less clearly defined. The LER is a major problem at these small dimensions, since current materials have LER of ~5 to 7 nm, unacceptably large for patterns of 20 nm. Much work is going on in trying to reduce the LER in resist.
- **Cost and infrastructure**—The costs of the exposure tools and masks have increased dramatically from the original projections of the mid-90s. Today, the cost of a stepper is projected to be well above US\$ 100 million, with the cost of a single EUV mask projected to more than US\$ 200 thousand. These costs are substantially higher than those of OL. Another challenge is the development of the required infrastructure network of suppliers for mask blanks, sources, tools, and so forth. The very specialized nature of EUV technology is making the development of this network slow and difficult.

There is a considerable amount of activity in all these areas, particularly in the development of more powerful sources, better resist materials and more effective mask defect inspection and repair techniques. The optics, multilayer coatings, and mechanical stages are well developed, and do not appear to be show stoppers. Industrial, academic, and national laboratory research remains strong.

In summary, at the time of writing EUV lithography has not yet reached its full potential. The appeal of a parallel imaging system capable of delivering a patterning resolution able to support sub-20-nm imaging is very strong, and explains the continuing interest of the industry in supporting the development of the EUV-L technology. Other techniques that rely mostly on electron beams in various forms of parallelization have not yet been demonstrated as a convincing alternative, while earlier precursors like proximity x-ray lithography have been discontinued. It is reasonable to expect that EUV will mature in the time frame 2008 to 2012, reaching a fully developed system in 2014 to 2015.

34.4 ACKNOWLEDGMENTS

We thank P. Naulleau (U.C. Berkeley) and G. Lorusso (IMEC, Belgium) for sharing the pictures presented in the chapter. This work would not have been possible without the support of the staff and students of the Center for Nano Technology at the University of Wisconsin.

34.5 REFERENCES

The most updated sources of information on EUV-IL can be found in the proceedings of the main lithography conferences: SPIE (www.spie.org), EIPBN (www.eipbn.org), MNE (mne08.org) and MNC (imnc.jp), as well as in the Sematech (www.sematech.org) and SIA (www.itrs.net) Web pages.

1. http://www.itrs.net/Links/2007ITRS/2007_Chapters/2007_Lithography.pdf.
2. C. Mack, *Fundamental Principles of Optical Lithography* Wiley, New York, 2008.
3. G. T. di Francia, "Resolving Power and Information," *J. Opt. Soc. Amer.* **45**:497, 1955.
4. S. Wurm, C. U. Jeon, and M. Lercel, "SEMATECH's EUV Program: A Key Enable for EUVL Introduction," *SPIE Proceeding* **6517**(4), 2007.
5. I. Mori, O. Suga, H. Tanaka, I. Nishiyama, T. Terasawa, H. Shigemura, T. Taguchi, T. Tanaka, and T. Itani, "Selete's EUV Program: Progress and Challenges," *SPIE Proceeding* **6921**(1), 2008.
6. T. E. Jewell, *OSA Proceeding on EUV Lithography* in F. Zernike and D. Attwood (eds.) **23**:98, 1994 and reference therein.
7. E. Spiller, *Soft X-Ray Optics*, SPIE Optical Engr. Press, Bellingham, Wash., 1994.
8. G. F. Lorusso, J. Hermans, A. M. Goethals, et al., "Imaging Performance of the EUV Ipha Demo Tool at IMEC," *SPIE Proceeding* **6921**(24), 2008.
9. N. Harned, M. Goethals, R. Groeneveld, et al., "EUV Lithography with the Alpha Demo Tools: Status and Challenges," *SPIE Proceeding* **6517**(5), 2007.
10. U. Stamm, "Extreme Ultraviolet Light Sources for Use in Semiconductor Lithography—State of the Art and Future Development," *Journal of Physics D: Applied Physics* **37**(23):3244–3253, 2004.
11. K. Murakami, T. Oshino, H. Kondo, H. Chiba, H. Komatsuda, K. Nomura, and H. Iwata, "Development Status of Projection Optics and Illumination Optics for EUV1," *SPIE Proceeding* **6921**(26), 2008.
12. T. Miura, K. Murakami, K. Suzuki, Y. Kohama, K. Morita, K. Hada, Y. Ohkubo, and H. Kawai, "Nikon EUVL Development Progress Update," *SPIE Proceeding* **6921**(22), 2008.
13. T. Wallow, C. Higgins, R. Brainard, K. Petrillo, W. Montgomery, C. Koay, G. Denbeaux, O. Wood, and Y. Wei, "Evaluation of EUV Resist Materials for Use at the 32 nm Half-Pitch Node," *SPIE Proceeding* **6921**(56), 2008.
14. P. Naulleau, UC Berkeley, private communication.
15. H. H. Solak, D. He, W. Li, S. Singh-Gasson, F. Cerrina, B. H. Sohn, X. M. Yang, and P. Nealey, "Exposure of 38 nm Period Grating Patterns with Extreme Ultraviolet Interferometric Lithography," *Applied Physics Letters* **75**(15):2328–2330, 1999.
16. F. Cerrina, A. Isoyan, F. Jiang, Y. C. Cheng, Q. Leonard, J. Wallace, K. Heinrich, A. Ho, M. Efremov, and P. Nealey, "Extreme Ultraviolet Interferometric Lithography: A Path to Nanopatterning," *Synchrotron Radiation News* **21**(4):12–24, 2008.
17. A. Isoyan, A. Wüest, J. Wallace, F. Jiang, and F. Cerrina, "4× Reduction Extreme Ultraviolet Interferometric Lithography," *Optics Express* **16**(12):9106–9111, 2008.
18. H. H. Solak, "Nanolithography with Coherent Extreme Ultraviolet Light," *Journal of Physics D: Applied Physics* **39**(10):R171–R188, 2006.
19. W. M. Moreau, *Semiconductor Lithography: Principles and Materials*, New York: Plenum, 1988.
20. H. Han, K. Goldberg, A. Barty, E. Gullikson, T. Ikuta, Y. Uno, O. Wood II, and S. Wurm, "EUV MET Printing and Actinic Imaging Analysis on the Effects of Phase Defects on Wafer CDs," *SPIE Proceeding* **6517**(10), 2007.
21. G. M. Gallatin, P. Naulleau, D. Niakoula, R. Brainard, E. Hassanein, R. Matyi, J. Thackeray, K. Spear, and K. Dean, "Resolution, LER, and Sensitivity Limitations of Photoresist," *SPIE Proceeding* **6921**(55), 2008.

This page intentionally left blank.

DO NOT DUPLICATE

RAY TRACING OF X-RAY OPTICAL SYSTEMS

Franco Cerrina

*Electrical and Computer Engineering & Center for NanoTechnology
University of Wisconsin
Madison, Wisconsin*

Manuel Sanchez del Rio

*European Synchrotron Radiation Facility
Grenoble, France*

35.1 INTRODUCTION

The first step before the construction of any x-ray system, such as a synchrotron beamline, is an accurate conceptual design of the optics. The beam should be transported to a given image plane (usually the sample position) and its characteristics should be adapted to the experimental requirements, in terms of flux monochromatization, focus, time structure, and so forth. The designer's goal is not only to verify compliance to a minimum set of requirements, but also to optimize matching between the source and the optics étendue (angle-area product) to obtain the highest possible flux.

The small difference between the index of refraction (see Chap. 36) of most materials and vacuum leads to critical angles of a few milliradians, and thus to the need for reflection glancing optics or diffraction systems (crystals and glancing gratings and multilayers). This complicates the job of the optical designer, because the aberrations become asymmetric and the power unequal. Yet, complex optical systems must be designed, with evermore stringent requirements. This is well exemplified in the quest for high-resolution x-ray microprobes and microscopes, or high-resolution phase contrast imaging systems. The traditional approach used in optical design often fails in x rays because of the differences in the approach to designing a visible versus an x-ray optical system. With few exceptions, most optical systems are either dioptric (e.g., lens-based) or catoptric (mirrors only) with a few catadioptric (mixed) systems for special applications. But most importantly, in the visible region the angle of incidence of the principal ray is always near the normal of the lens (or mirror), in what is often a paraxial optical system. By contrast, x-ray optical systems are strong off-axis systems, with angles of incidence close to 90 degrees.

Today, optical design relies more and more on computer simulation and optimization, and indeed very powerful programs such as CodeV (www.opticalres.com) and Zemax (www.zemax.com), to name only two, are widely used. These programs are however unwieldy when applied to the x-ray domain, particularly because of the off-axis geometry that makes cumbersome to define the geometry. A modeling code dedicated to the x-ray domain is thus a powerful tool for the optical designer. In addition, x-ray sources are peculiar, ranging from synchrotron bending magnets, to insertion devices and free electron lasers.

Often, when dealing with synchrotron-based optical systems, the optical elements themselves are just a subsystem of a complex beamline. Thus, the ability of modeling the progress of the radiation through the beamline in a realistic way is essential. The main question that an optical designer must answer is "Will the optical system deliver the performance needed by the experimental station?" The answer often

generates a second question, that is, “What is the effect of optics imperfections, thermal loading, material changes on the performance?” Clearly, the second question is “practical,” and yet essential in determining the performance of the beamline in a realistic world. This is where SHADOW was born: as a Monte-Carlo simulation code capable of ray tracing the progress of a beam of x rays (or any other photon beam) through a complex, sequential optical system. SHADOW models and predicts the properties of a beam of radiation from the source, through multiple surfaces, taking into account the physics in all reflection, refraction, and diffraction processes. While it is of general use, SHADOW was designed from the ground-up for the study of the glancing and diffractive optics typical of x-ray systems.

For synchrotron radiation applications, the code SHADOW has become the *de facto* standard because (1) it is modular and flexible, capable of adapting to any optical configurations, (2) it has demonstrated its reliability during more of 20 years of use, as shown in hundred of publications, (3) it is simple to use, and (4) it is in the public domain. Indeed, almost all of the synchrotron beamlines today in existence have in some way benefited from the help of SHADOW.^{*,†} From the list of selected references at the end of this chapter, the interested reader can form a good idea of the many uses of SHADOW.¹⁻¹¹

35.2 THE CONCEPTUAL BASIS OF SHADOW

The computational model used in SHADOW follows the evolution of a “beam of radiation.” This beam is a collection of independent “rays.” A “ray” is a geometrical entity defined by two vectors: a starting position $\mathbf{x} = (x_0, y_0, z_0)$ and the direction vector $\mathbf{v} = (v_x, v_y, v_z)$. In addition, a scalar, $k = 2\pi/\lambda$, defines the wavenumber and the electric field is described by two vectors \mathbf{A}_σ and \mathbf{A}_π for the two polarizations, with two scalars for the phases ϕ_σ and ϕ_π .

The first step in simulating an optical device with a ray-tracing code is the generation of the source, that is, the beam at the source position. The beam is a collection of rays, usually many thousands, which are created by a Monte Carlo sampling of the spatial (for starting positions) and divergences (for the direction) source distributions. SHADOW includes models for synchrotron (bending magnets, wigglers, and undulators) and geometrical (box, Gaussian, and so forth) sources. A detailed discussion of SHADOW’s source models can be found in Cerrina.¹² Although a single ray is monochromatic, with a well-defined wavelength, white and polychromatic beams are formed by a collection of rays with nonequal wave numbers, giving an overall spectral distribution. Essentially, SHADOW samples the wavefront of an ensemble of point sources. Rays are propagated (in vacuum or air) in straight lines, until they interact with the optical elements (mirrors or gratings). The set of optical elements constitute the optical system. The tracing is sequential: the beam goes to the first optical element, which modifies it, then it goes to a second element, and so on, until arriving at the final detection plane. The optical elements are defined by the equation of the mathematical surface (plane, sphere, ellipsoid, polynomials, and so forth), and the intercept point between each ray and the surface is calculated by solving the equation system of the straight line for the ray and for the surface (in general, a quadric equation, a toroidal or a polynomial surface). At the intercept point, the normal to the surface \mathbf{n}_D is computed from the gradient. The intercept point is the new starting point for the ray after the interaction with the optical element. The direction is changed following either the specular reflection law (for mirrors) or the boundary conditions at the surface (for gratings and crystals). These equations are written in vector form to allow calculations in 3D and improve efficiency (vector operations are faster in a computer than trigonometric calculations). A (compact) vector notation for the specular reflectivity can be written as $\mathbf{v}_{\text{out}} = \mathbf{v}_{\text{in}} - 2(\mathbf{v}_{\text{in}} \cdot \mathbf{n}_D)\mathbf{n}_D$, where all the vectors are unitary. The change in the direction of a monochromatic beam *diffracted* by an optical surface can be calculated (i) using the boundary condition at the surface $\mathbf{k}_{\text{out},\parallel} = \mathbf{k}_{\text{in},\parallel} + \mathbf{G}_\parallel$, where \mathbf{G}_\parallel is the projection of the reciprocal lattice vector \mathbf{G} onto the optical surface, and (ii) selecting $\mathbf{k}_{\text{out},\perp}$ parallel to $\mathbf{k}_{\text{in},\perp}$, which guarantees the elastic scattering in the diffraction process, that is, the conservation of momentum $|\mathbf{k}_{\text{out}}| = |\mathbf{k}_{\text{in}}|$.

^{*}The executables for SHADOW and related files can be downloaded from http://www.nanotech.wisc.edu/CNT_LABS/shadow.html.

[†]The graphical interface to SHADOW developed at ESRF can be downloaded from <http://www.esrf.eu/UsersAndScience/Experiments/TBS/SciSoft/xop2.3/shadowvui/>.

In a diffraction grating the scattering vector is originated by the ruling, with $\mathbf{G}_{\parallel} = (m2\pi/p)\mathbf{u}_g$, p being the local period (*SHADOW* takes into account that p may not be constant along the grating surface), m the spectral order, and \mathbf{u}_g a unitary vector tangent to the optical surface and pointing perpendicular to the grating grooves. It is straightforward to verify that these the scattering equations give the grating equation in its usual form: $m\lambda = p(\sin\alpha + \sin\beta)$ where $\alpha(\beta)$ is the angle of incidence (reflection) defined from the normal to the surface (see Chap. 38).

For crystals, \mathbf{G}_{\parallel} is zero for the Bragg-symmetric crystals, when the crystal surface is parallel to the atomic planes. There is no additional scattering vector here, thus the Bragg-symmetric crystals are nondispersive systems. Any other crystal (asymmetric Bragg or any Laue) gives nonzero scattering: $\mathbf{G}_{\parallel} = (2\pi\sin\gamma/d_{hkl})\mathbf{u}_g$ being d_{hkl} the d-spacing of the selected crystal reflection, γ the angle between Bragg planes and surface, and now \mathbf{u}_g a unitary vector tangent to the optical surface and pointing perpendicular to the lines created by the termination of the crystal planes at the surface. In other words, in crystal with a given asymmetry ($\gamma \neq 0$), the truncation of the Bragg planes with the crystal surface mimics a grating that is used to compute the changes in the direction of the beam. As \mathbf{G}_{\parallel} is a function of λ , Bragg-asymmetric and Laue crystals are dispersive systems.

In addition to the change of direction of the ray at the optical surface, there is a change in amplitude and phase that can be computed using an adequate physical model. The ray electric-field vectors are then changed using Fresnel equations (for mirrors and lenses) or the dynamical theory of diffraction for crystals. In addition, nonidealities are easily prescribed. Roughness is described by the addition of a random scattering vector generated by a stochastic process; given a power spectrum of the roughness, the spectrum is sampled to generate a local “grating” corresponding to the selected spatial frequency. Finally, deterministic surface errors are included by adding a small correction term to the ideal surface, and finding the “real” intercept with an iterative method.

Slits are easily implemented: if the ray goes through the slit, it survives, otherwise it is discarded or rather labeled as “lost”. The beam (i.e., collection of rays) is scored at the detector plane, and several statistical tools (integration, histogramming, scatter and contour plots) may be used for calculating the required parameters (intensity, spatial and angular distributions, resolution, and so forth).

35.3 INTERFACES AND EXTENSIONS OF SHADOW

It is desirable to add a user-friendly graphical user interface (GUI) to *SHADOW* to help the user in analyzing the considerable output from the modeling. The *SHADOW* package includes a complete menu for defining the source and system parameters and basic plotting and conversion utilities, for displaying the results. A GUI written with free software (TCL-TK) is shipped with *SHADOW*,* but much more sophisticated and powerful GUIs can be written using very powerful graphical commercial packages. Some users developed displaying utilities in Matlab, Mathematica, or IDL. A complete user interface (*SHADOWVUI*) written in IDL is freely available under the XOP package,^{13†} and it allows the user to run *SHADOW* using a multiwindow environment. This helps the user to modify the optical system, rerun *SHADOW* with modified inputs, and quickly refresh all the screens showing interesting information for the user (XY plots, histograms, etc.). Another powerful feature in *SHADOWVUI* is the availability of macros. These macros permit the user to run *SHADOW* in a loop, to perform powerful postprocessing, to make parametric calculations, and to compute a posteriori some basic operations (tracing in vacuum, vignetting, etc.). A beamline viewer application (BLViewer) helps in creating three-dimensional schematic views of the optical system. A tutorial with many examples is available.

SHADOW is an open code, and user contributions are not only possible but also encouraged. Whereas the *SHADOW* developers control the “official” code revisions and releases, anyone can contribute routines or interfaces. User contributions that are especially useful and of general applicability may be incorporated (with the obvious author’s agreement and collaboration) to the official version. This is effectively possible due to the file-oriented structure of *SHADOW*.

*The executables for *SHADOW* and related files can be downloaded from http://www.nanotech.wisc.edu/CNT_LABS/shadow.html.

†The graphical interface to *SHADOW* developed at ESRF can be downloaded from <http://www.esrf.eu/UsersAndScience/Experiments/TBS/SciSoft/xop2.3/shadowvui/>.

35.4 EXAMPLES

More than 250 papers have been published describing the use of SHADOW in a broad range of applications. Here, we limit the discussion to a few examples to illustrate the use of program.

Monochromator Optics

Almost all synchrotron beamlines in the world have a monochromator to select and vary the energy of the photons delivered to the sample. Soft x-ray beamlines use grating monochromators, and hard x-ray beamlines use crystal monochromators. The performance of a monochromator is linked not only to the optical element itself, but also to other parameters of the beamline such as the source dimensions and divergences, slits and focusing that, in turn, can be done externally (upstream mirrors) or by using curved optical elements at the monochromator.

SHADOW calculates very accurately the aberrations and resolving power for all possible grating monochromators (PGM, SGM, TGM, SX700, DRAGON, etc.). For hard x-ray beamlines, an Si (alternatively Ge or diamond) double-crystal monochromator is commonly used. Bending magnet beamlines commonly use sagittal focusing with the second crystal. It is well known that the efficiency of the sagittal focusing depends on the magnification factor, and this can be efficiently computed by SHADOW (Fig. 1).

SHADOW also models more complex crystal monochromators, using crystals in transmission geometry for splitting the beam or for increasing efficiency at high energies. Polychromatic focusing (as in XAFS dispersive beamlines, as described in Chap. 30) or monochromatic focusing (including highly asymmetric back-scattering) can be analyzed with SHADOW. A detailed discussion of the crystal optics with SHADOW can be found in M. Sanchez del Rio.¹³

Mirror Optics

Mirrors are used essentially to focus or collimate the x-ray beam. They can also be used as low pass filters to reject the higher harmonics from diffractive elements or as x-ray filters to avoid high heat load on other devices such as the monochromators. Together with the monochromator, they constitute the most commonly used element in a beamline.

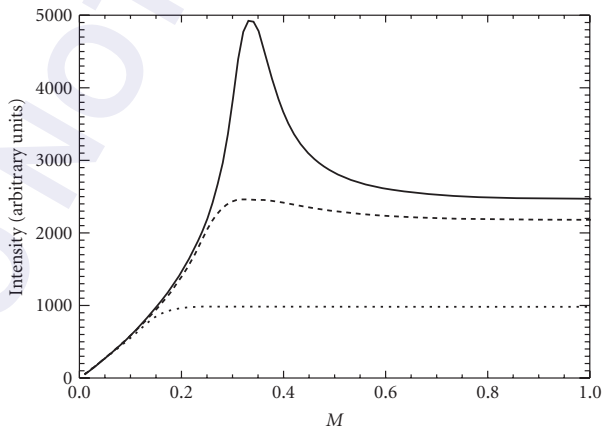


FIGURE 1 Intensity (in arbitrary units) versus magnification factor M for a point and monochromatic ($E = 20$ keV) source placed at 30 m from the sagittally Si bent crystals. Three beam divergences are considered: 1 mrad (lower), 2.5 mrad (middle), and 5 mrad (upper). We clearly observe the maximum of the transmission at $M = 0.33$ when focusing the 5-mrad beam, as predicted by the theory. (C. J. Sparks, B. S. Borie, and J. B. Hastings.¹⁴)

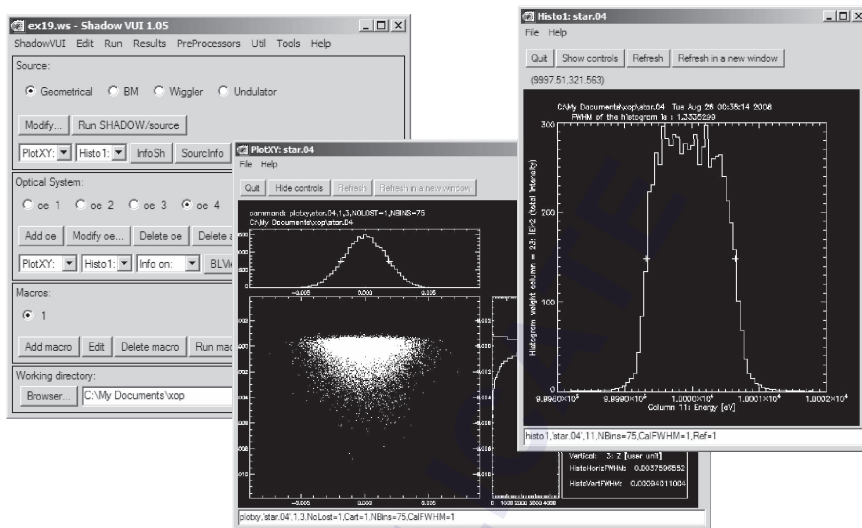


FIGURE 2 SHADOWVUI windows with the outputs of the ray tracing for the hard x-ray beamline (see text).

Reflection of x rays requires glancing incidence, which implies the use of very long mirrors (typically 30 to 100 cm). Glancing angles also magnify the effect of geometrical aberrations and surface irregularities (figure errors, slope errors, and roughness, as described in Chaps. 44 and 45). These effects are difficult to study with fully analytical methods, and it is even more difficult to analyze their combination with the geometrical and spectral characteristics of the synchrotron sources. The mirrors may be bent spherically (cylinders, spheres, and toroids) but elliptical curvature would be ideal in most cases. Spherical mirrors are usually preferred because of their lower cost and higher finishing quality. Dynamically bent mirrors are becoming very popular, and can also be used to compensate figure errors. SHADOW is well suited to perform reliable and accurate calculations of all these mirror systems under synchrotron radiation.

An important part of a ray-tracing calculation is the study of the tolerances to source movements, alignments, sample displacements, etc. SHADOW is well suited for these calculations because it allows the user to freely displace the source and the optical elements whilst conserving the same initial reference frame.

A Hard X-Ray Beamline

As an example, a full hard x-ray beamline has been raytraced with SHADOW. The undulator source can be simplified using Gaussian spatial ($\sigma_x = 0.57 \cdot 10^{-2}$ cm, $\sigma_z = 0.104 \cdot 10^{-2}$ cm) and angular distributions (horizontal $\sigma_x = 88.5 \mu\text{rad}$, vertical $\sigma_z = 7.2 \mu\text{rad}$), at $10,000 \pm 2$ eV (box distribution), with theoretical flux at this energy 5×10^{13} ph/sec/0.1%bw. The beamline has two mirrors, M1, a Rh-coated cylindrical collimating mirror in the vertical plane at 25 m from the source (glancing angle 0.12 mrad), and M2, a refocusing mirror (same coating and angle as M1) at 35 m from the source, focusing at the sample position. A Si (111) double crystal monochromator with second crystal sagittally bent is placed between the mirrors, at 30 m from source. The SHADOWVUI windows are shown in Fig. 2 and the resulting parameters are (i) Beam size at the sample position ($37.6 \times 9.4 \mu\text{m}^2$), (ii) energy resolution (1.33 eV), (iii) transmissivity of the whole beamline ($T = 0.85 \text{ eV}^{-1}$) and number of photons at the sample position (5.65×10^{12} photons/s).

35.5 CONCLUSIONS AND FUTURE

The SHADOW code is now about 20 years old, a statement to forward looking software design and continuing evolution. It has helped generations of postdocs and optical designers in developing

incredibly complex x-ray beamlines. While SHADOW performs very well, it begins to suffer from limitations, inherent in its kernel. For instance, the development of free electron laser sources requires a code capable of dealing with the sophisticated time structure, and derives coherence properties of these novel sources. The manifold increase in power of computers from the mid-80s, both in term of speed and memory space, makes this kind of calculation possible. Thus, SHADOW's developers are planning to rewrite completely the code of the kernel (2008). The new structure will overcome technical and physical limitations of the current version. From the physical point of view, it is important to include in the simulation coherent beams, and the effect of the optical elements on them. Specifically, we need to have efficient tools for computing imaging and propagation of diffracted fields. This is important for many techniques already in use (e.g., phase contrast imaging) and will be essential for the new generation sources, which will be almost completely coherent. In the meantime SHADOW remains available to the x-ray optical community.*†

35.6 REFERENCES

1. L. Alianelli, M. Sanchez del Rio, M. Khan, and F. Cerrina, "A Comment on A New Ray-Tracing Program RIGTRACE for X-Ray Optical Systems," [*J. Synchrotron Rad.* 8:1047–1050 (2001)]. *Journal of Synchrotron Rad.* **10**:191–192, 2003.
2. B. Lai, K. Chapman, and F. Cerrina, "SHADOW: New Developments," *Nuclear Instruments & Methods in Physics Research Section A: Accelerators Spectrometers Detectors and Associated Equipment* **266**:544–549 (1988).
3. M. Sanchez del Rio, "Experience with Ray-Tracing Simulations at the European Synchrotron Radiation Facility," *Review of Scientific Instruments* **67**(9) (1996) [+CD-ROM].
4. M. Sanchez del Rio, "Ray Tracing Simulations for Crystal Optics," *Proceedings of the SPIE—The International Society for Optical Engineering* **3448**:230–245 (1998).
5. M. Sanchez del Rio, S. Bernstorff, A. Savoia, and F. Cerrina, "A Conceptual Model for Ray Tracing Calculations with Mosaic Crystals," *Review of Scientific Instruments* pt.11B **63**(1):932–935 (1992).
6. M. Sanchez del Rio and F. Cerrina, "Asymmetrically Cut Crystals for Synchrotron Radiation Monochromators," *Review of Scientific Instruments* pt. 11B **63**(1):936–940 (1992).
7. M. Sanchez del Rio and F. Cerrina, "Comment on 'Comments on the Use of Asymmetric Monochromators for X-Ray Diffraction on a Synchrotron Source,'" [*Rev. Sci. Instrum.* 66:2174 (1995)]. *Review of Scientific Instruments* **67**:3766–3767 (1996).
8. M. Sanchez del Rio and R. J. Dejus, "XOP 2.1—A New Version of the X-Ray Optics Software Toolkit," *American Institute of Physics Conference Proceedings* 705:784 (2004).
9. M. Sanchez del Rio, C. Ferrero, G. J. Chen, and F. Cerrina, "Modeling Perfect Crystals in Transmission Geometry for Synchrotron Radiation Monochromator Design," *Nuclear Instruments & Methods in Physics Research, Section A: Accelerators, Spectrometers, Detectors, and Associated Equipment* **347**:338–43 (1994).
10. C. Welnak, P. Anderson, M. Khan, S. Singh, and F. Cerrina, "Recent Developments In SHADOW," *Review of Scientific Instruments* **63**:865–868 (1992).
11. C. Welnak, G. J. Chen, and F. Cerrina, "SHADOW—A Synchrotron-Radiation and X-Ray Optics Simulation Tool," *Nuclear Instruments & Methods in Physics Research Section A: Accelerators Spectrometers Detectors and Associated Equipment* **347**:344–347 (1994).
12. F. Cerrina, "Ray Tracing of X-Ray Optical Systems: Source Models," *SPIE Proceedings* **1140**:330–336 (1989).
13. M. Sanchez del Rio and R. J. Dejus, "XOP: Recent Developments," *SPIE Proceedings* **3448**:230–245, 340–345 (1998).
14. C. J. Sparks, B. S. Borie, and J. B. Hastings, "X-Ray Monochromator Geometry For Focusing Synchrotron Radiation above 10 keV," *Nuclear Instruments & Methods* **172**:237–224 (1980).

*The executables for SHADOW and related files can be downloaded from http://www.nanotech.wisc.edu/CNT_LABS/shadow.html.

†The graphical interface to SHADOW developed at ESRF can be downloaded from <http://www.esrf.eu/UsersAndScience/Experiments/TBS/SciSoft/xop2.3/shadowvui/>.

X-RAY PROPERTIES OF MATERIALS

Eric M. Gullikson

*Center for X-Ray Optics
Lawrence Berkeley National Laboratory
Berkeley, California*

The primary interaction of low-energy x rays within matter, namely, photoabsorption and coherent scattering, have been described for photon energies outside the absorption threshold regions by using atomic scattering factors, $f = f_1 + if_2$. The atomic photoabsorption cross section, μ , may be readily obtained from the values of f_2 using the relation

$$\mu = 2r_0\lambda f_2 \quad (1)$$

where r_0 is the classical electron radius, and λ is the wavelength. The transmission of x rays through a slab of thickness, d , is then given by

$$T = \exp(-N\mu d) \quad (2)$$

where N is the number of atoms per unit volume in the slab. The index of refraction, n , for a material is calculated by

$$n = 1 - Nr_0\lambda^2(f_1 + if_2)/(2\pi) \quad (3)$$

The (semiempirical) atomic scattering factors are based upon photoabsorption measurements of elements in their elemental state. The basic assumption is that condensed matter may be modeled as a collection of noninteracting atoms. This assumption is, in general, a good one for energies sufficiently far from absorption thresholds. In the threshold regions, the specific chemical state is important and direct experimental measurements must be made. Note also that the Compton scattering cross section is not included. The Compton cross section may be significant for the light elements ($Z < 10$) at the higher energies considered here (10 to 30 keV).

The atomic scattering factors are plotted in Figs. 1 and 2 for every 10th element. Tables 1 through 3 are based on a compilation of the available experimental measurements and theoretical calculations. For many elements there is little or no published data, and, in such cases, it was necessary to rely on theoretical calculations and interpolations across Z . To improve the accuracy in the future, considerably more experimental measurements are needed.¹ More data and useful calculation engines are available at www.cxro.lbl.gov/optical_constants.

36.1 X-RAY AND NEUTRON OPTICS

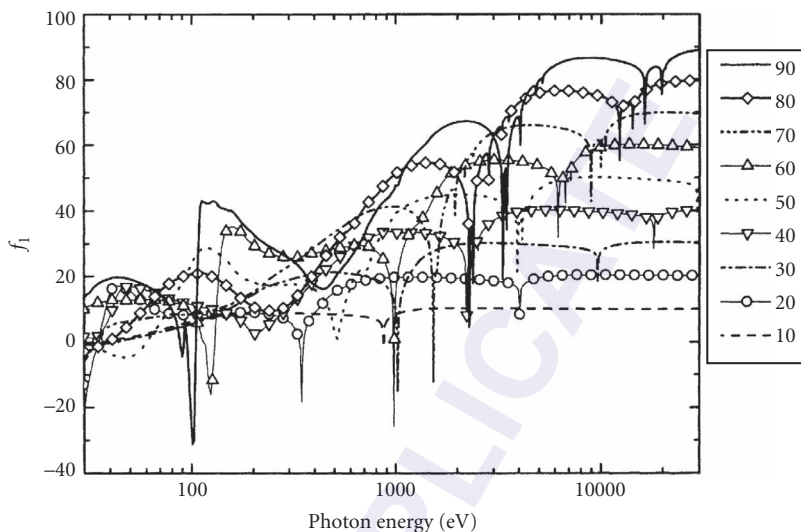


FIGURE 1 The atomic scattering factor f_1 as a function of photon energy from 30 to 30,000 eV for atomic number 10(Ne), 20(Ca), 30(Zn), 40(Zr), 50(Sn), 60(Nd), 70(Yb), 80(Hg), and 90(Th).

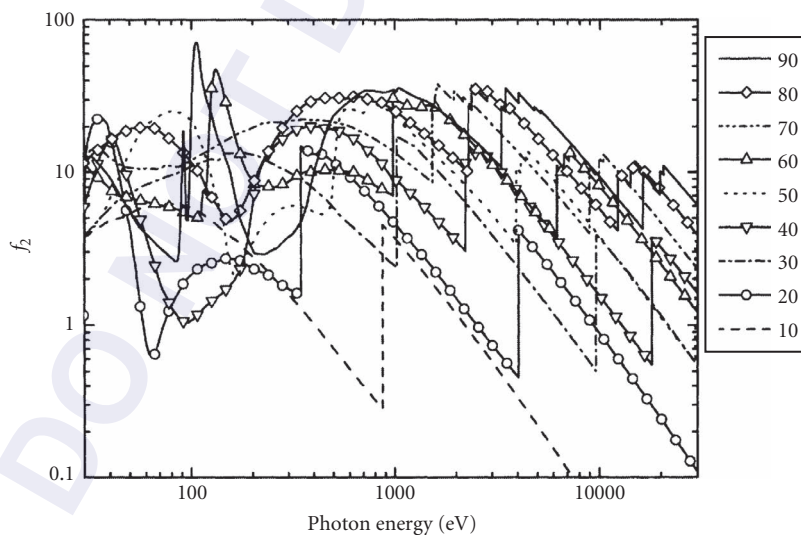


FIGURE 2 The atomic scattering factor f_2 as a function of photon energy from 30 to 30,000 eV for atomic number 10(Ne), 20(Ca), 30(Zn), 40(Zr), 50(Sn), 60(Nd), 70(Yb), 80(Hg), and 90(Th).

36.2 ELECTRON BINDING ENERGIES, PRINCIPAL K- AND L-SHELL EMISSION LINES, AND AUGER ELECTRON ENERGIES

TABLE 1 Electron Binding Energies in Electron Volts (eV) for the Elements in Their Natural Forms

Element	K _{1s}	L ₁ 2s	L ₂ 2p _{1/2}	L ₃ 2p _{3/2}	M ₁ 3s	M ₂ 3p _{1/2}	M ₃ 3p _{3/2}	M ₄ 3d _{3/2}	M ₅ 3d _{5/2}	N ₁ 4s	N ₂ 4p _{1/2}	N ₃ 4p _{3/2}
1 H	13.6											
2 He	24.6*											
3 Li	54.7*											
4 Be	111.5*											
5 B	188*											
6 C	284.2*											
7 N	409.9*	37.3*										
8 O	543.1*	41.6*										
9 F	696.7*											
10 Ne	870.2*	48.5*	21.7*	21.6*								
11 Na	1070.8†	63.5†	30.4†	30.5*								
12 Mg	1303.0†	88.6*	49.2†	49.2†								
13 Al	1559.6	117.8*	72.9*	72.5*								
14 Si	1838.9	149.7*	99.8*	99.2*								
15 P	2145.5	189*	136*	135*								
16 S	2472	230.9*	163.6*	162.5*								
17 Cl	2822.4	270.2*	202*	200*								
18 Ar	3205.9*	326.3*	250.6*	248.4*	29.3*	15.9*	15.7*					
19 K	3608.4*	378.6*	297.3*	294.6*	34.8*	18.3*	18.3*					
20 Ca	4038.5*	438.4†	349.7†	346.2†	44.3†	25.4†	25.4†					
21 Sc	4492.8	498.0*	403.6*	398.7*	51.1*	28.3*	28.3*					
22 Ti	4966.4	560.9†	461.2†	453.8†	58.7†	32.6†	32.6†					
23 V	5465.1	626.7†	519.8†	512.1†	66.3†	37.2†	37.2†					
24 Cr	5989.2	695.7†	583.8†	574.1†	74.1†	42.2†	42.2†					
25 Mn	6539.0	769.1†	649.9†	638.7†	82.3†	47.2†	47.2†					
26 Fe	7112.0	844.6†	719.9†	706.8†	91.3†	52.7†	52.7†					
27 Co	7708.9	925.1†	793.3†	778.1†	101.0†	58.9†	58.9†					
28 Ni	8332.8	1008.6†	870.0†	852.7†	110.8†	68.0†	66.2†					
29 Cu	8978.9	1096.7†	952.3†	932.5†	122.5†	77.3†	75.1†					
30 Zn	9658.6	1196.2*	1044.9*	1021.8*	139.8*	91.4*	88.6*	10.2*	10.1*			
31 Ga	10367.1	1299.0*	1143.2†	1116.4†	159.5†	103.5†	103.5†	18.7†	18.7†			
32 Ge	11103.1	1414.6*	1248.1*	1217.0*	180.1*	124.9*	120.8*	29.0*	29.0*			
33 As	11866.7	1527.0*	1359.1*	1323.6*	204.7*	146.2*	141.2*	41.7*	41.7*			
34 Se	12657.8	1652.0*	1474.3*	1433.9*	229.6*	166.5*	160.7*	55.5*	54.6*			

(Continued)

TABLE 1 Electron Binding Energies in Electron Volts (eV) for the Elements in Their Natural Forms (Continued)

Element	K1s	L ₁ -2s	L ₂ -2p _{1/2}	L ₃ -2p _{3/2}	M ₁ 3s	M ₂ -3p _{1/2}	M ₃ -3p _{3/2}	M ₄ -3d _{3/2}	M ₅ -3d _{5/2}	N ₁ 4s	N ₂ -4p _{1/2}	N ₃ -4p _{3/2}
35 Br	13473.7	1782.0*	1596.0*	1549.9*	257*	189*	182*	70*	69*	27.5*	14.1*	14.1*
36 Kr	14325.6	1921.0	1730.9*	1678.4*	292.8*	222.2*	214.4	95.0*	93.8*	30.5*	16.3*	15.3*
37 Rb	15199.7	2065.1	1863.9	1804.4	326.7*	248.7*	239.1*	113.0*	112*	38.9†	20.3†	20.3†
38 Sr	16104.6	2216.3	2006.8	1939.6	358.7†	280.3†	270.0†	136.0†	134.2†	43.8*	24.4*	23.1*
39 Y	17038.4	2372.5	2155.5	2080.0	392.0*	310.6*	298.8*	157.7†	155.8*	50.6†	28.5†	27.7†
40 Zr	17997.6	2531.6	2306.7	2222.3	430.3†	343.5†	329.8†	181.1†	178.8†	56.4†	32.6†	30.8†
41 Nb	18985.6	2697.7	2464.7	2370.5	466.6†	376.1†	360.6†	205.0†	202.3†	63.2†	37.6†	35.5†
42 Mo	19999.5	2865.5	2625.1	2520.2	506.3†	411.6†	394.0†	231.1†	227.9†	68*	39*	39*
43 Tc	21044.0	3042.5	2793.2	2676.9	544*	445*	425*	257*	253*	75.0†	46.5†	43.2†
44 Ru	22117.2	3224.0	2966.9	2837.9	586.2†	483.5†	461.4†	284.2†	280.0†	81.4*	50.5†	47.3†
45 Rh	23219.9	3411.9	3146.1	3003.8	628.1†	521.3†	496.5†	311.9†	307.2†	87.6*	55.7†	50.9†
46 Pd	24350.3	3604.3	3330.3	3173.3	671.6†	559.9†	532.3†	340.5†	335.2†	97.0†	63.7†	58.3†
47 Ag	25514.0	3805.8	3523.7	3351.1	719.0†	603.8†	573.0†	374.0†	368.0†	109.8†	63.9†	63.9†
48 Cd	26711.2	4018.0	3727.0	3537.5	772.0†	652.6†	618.4†	411.9†	405.2†	122.7†	73.5†	73.5†
49 In	27939.9	4237.5	3938.0	3730.1	827.2†	703.2†	665.3†	451.4†	443.9†	137.1†	83.6†	83.6†
50 Sn	29200.1	4464.7	4156.1	3928.8	884.7†	756.5†	714.6†	493.2†	484.9†	153.2†	95.6†	95.6†
51 Sb	30491.2	4698.3	4380.4	4132.2	946†	812.7†	766.4†	537.5†	528.2†	169.4†	103.3†	103.3†
52 Te	31813.8	4939.2	4612.0	4341.4	1006†	870.8†	820.8†	583.4†	573.0†	186*	123*	123*
53 I	33169.4	5188.1	4852.1	4557.1	1072*	931*	875*	631*	620*	213.2*	146.7	145.5*
54 Xe	34561.4	5452.8	5103.7	4782.2	1148.7*	1002.1*	940.6*	689.0*	676.4*	232.3*	172.4*	161.3*
55 Cs	35984.6	5714.3	5359.4	5011.9	1211*	1071*	1003*	740.5*	726.6*	192	178.6†	178.6†
56 Ba	37440.6	5988.8	5623.6	5247.0	1293*	1137*	1063*	795.7*	780.5*	205.8	206.5*	206.5*
57 La	38924.6	6266.3	5890.6	5482.7	1362*	1209*	1128*	853*	836*	247.7*	243.3	242
58 Ce	40443.0	6548.8	6164.2	5723.4	1436*	1274*	1187*	902.4*	883.8*	304.5	236.3	217.6
59 Pr	41990.6	6834.8	6440.4	5964.3	1511.0	1374	1242.2	948.3*	928.8*	—	242	242
60 Nd	43568.9	7126.0	6721.5	6207.9	1575.3	1402.8	1297.4	1003.3*	980.4*	347.2*	265.6	247.4
61 Pm	45184.0	7427.9	7012.8	6459.3	—	1471.4	1356.9	1051.5	1026.9	360	284	257
62 Sm	46834.2	7736.8	7311.8	6716.2	1722.8	1540.7	1419.8	1110.9*	1083.4*	378.6*	286	270.9
63 Eu	48519.0	8052.0	7617.1	6976.9	1800.0	1613.9	1480.6	1158.6†	1127.5*	396.0*	322.4*	284.1*
64 Gd	50239.1	8375.6	7930.3	7242.8	1880.8	1688.3	1544.0	1221.9*	1189.6*	414.2*	333.5*	293.2*
65 Tb	51995.7	8708.0	8251.6	7514.0	1967.5	1767.7	1611.3	1276.9*	1241.1*	432.4*	343.5	308.2*
66 Dy	53788.5	9045.8	8580.6	7790.1	2046.8	1841.8	1675.6	1332.5	1292.6*	449.8*	366.2*	320.2*
67 Ho	55617.7	9394.2	8917.8	8071.1	2128.3	1922.8	1741.2	1391.5	1351.4	470.9*	385.9*	332.6*
68 Er	57485.5	9751.3	9264.3	8357.9	2206.5	2005.8	1811.8	1453.3	1409.3	480.5*	388.7*	339.7*
69 Tm	59398.6	10115.7	9616.9	8648.0	2306.8	2089.8	1884.5	1514.6	1467.7	538	412.4*	359.2*
70 Yb	61332.3	10486.4	9978.2	8943.6	2398.1	2173.0	1949.8	1576.3	1527.8	563.4†	438.2†	380.7†
71 Lu	63313.8	10870.4	10348.6	9244.1	2491.2	2263.5	2023.6	1639.4	1588.5	634.4†	463.4†	400.9†
72 Hf	65350.8	11270.7	10739.4	9560.7	2600.9	2365.4	2107.6	1716.4	1661.7	735.1	594.1†	423.6†
73 Ta	67416.4	11681.5	11136.1	9881.1	2708.0	2468.7	2194.0	1793.2	1735.1	870.0†	490.4†	423.6†
74 W	69525.0	12099.8	11544.0	10206.8	2819.6	2574.9	2281.6	1871.6	1809.2	—	—	—

Element	N ₄ 4d _{3/2}	N ₄ 4d _{5/2}	N ₄ 4f _{7/2}	O ₅ s	O ₂ 5p _{1/2}	O ₃ 5p _{3/2}	O ₄ 5d _{3/2}	O ₅ d _{5/2}				
75 Re	71676.4	12526.7	11958.7	10535.3	2931.7	2681.6	2367.3	1948.9	1882.9	625.4	518.7 [†]	446.8 [†]
76 Os	73870.8	12968.0	12385.0	10870.9	3048.5	2792.2	2457.2	2030.8	1960.1	658.2 [†]	549.1 [†]	470.7 [†]
77 Ir	76111.0	13418.5	12824.1	11215.2	3173.7	2908.7	2550.7	2116.1	2040.4	691.1 [†]	577.8 [†]	495.8 [†]
78 Pt	78394.8	13879.9	13272.6	11563.7	3296.0	3026.5	2645.4	2201.9	2121.6	725.4 [†]	609.1 [†]	519.4 [†]
79 Au	80724.9	14352.8	13733.6	11918.7	3424.9	3147.8	2743.0	2291.1	2205.7	762.1 [†]	642.7 [†]	546.3 [†]
80 Hg	83102.3	14839.3	14208.7	12283.9	3561.6	3278.5	2847.1	2384.9	2294.9	802.2 [†]	680.2 [†]	576.6 [†]
81 Tl	85530.4	15346.7	14697.9	12657.5	3704.1	3415.7	2956.6	2485.1	2389.3	846.2 [†]	720.5 [†]	609.5 [†]
82 Pb	88004.5	15860.8	15200.0	13035.2	3850.7	3554.2	3066.4	2585.6	2484.0	891.8 [†]	761.9 [†]	643.5 [†]
83 Bi	90525.9	16387.5	15711.1	13418.6	3999.1	3696.3	3176.9	2687.6	2579.6	939 [†]	805.2 [†]	678.8 [†]
84 Po	93105.0	16939.3	16244.3	13813.8	4149.4	3854.1	3301.9	2798.0	2683.0	995 [†]	851 [†]	705 [†]
85 At	95729.9	17493	16784.7	14213.5	4317	4008	3426	2908.7	2786.7	1042 [†]	886 [†]	740 [†]
86 Rn	98404	18049	17337.1	14619.4	4482	4159	3538	3021.5	2892.4	1097 [†]	929 [†]	768 [†]
87 Fr	101137	18639	17906.5	15031.2	4652	4327	3663	3136.2	2999.9	1153 [†]	980 [†]	810 [†]
88 Ra	103921.9	19236.7	18484.3	15444.4	4822.0	4489.5	3791.8	3248.4	3104.9	1208 [†]	1057.6 [†]	879.1 [†]
89 Ac	106755.3	19840	19083.2	15871.0	5002	4656	3909	3370.2	3219.0	1269 [†]	1080 [†]	890 [†]
90 Th	109650.9	20472.1	19693.2	16300.3	5182.3	4830.4	4046.1	3490.8	3332.0	1330 [†]	1168 [†]	966.4 [†]
91 Pa	112601.4	21104.6	20313.7	16733.1	5366.9	5000.9	4173.8	3611.2	3441.8	1387 [†]	1224 [†]	1007 [†]
92 U	115606.1	21757.4	20947.6	17166.3	5548.0	5182.2	4303.4	3727.6	3551.7	1439 [†]	1271 [†]	1043.0 [†]

TABLE 1 Electron Binding Energies in Electron Volts (eV) for the Elements in Their Natural Forms (Continued)

Element	$N_4d_{3/2}$	$N_4d_{5/2}$	$N_6f_{5/2}$	$N_7f_{7/2}$	O_1s	$O_2p_{1/2}$	$O_3p_{3/2}$	$O_4d_{3/2}$	$O_5d_{5/2}$
70 Yb	191.2*	182.4*	—	—	52.0*	30.3*	24.1*	—	—
71 Lu	206.1*	196.3†	8.9*	7.5*	57.3*	33.6*	26.7*	—	—
72 Hf	220.0†	211.5†	15.9†	14.2†	64.2†	38*	29.9*	—	—
73 Ta	237.9†	226.4†	23.5†	21.6†	69.7†	42.2*	32.7*	—	—
74 W	255.9†	243.5†	33.6*	31.4†	75.6†	45.3*	36.8*	—	—
75 Re	273.9†	260.5†	42.9*	40.5†	83†	45.6*	34.6*	—	—
76 Os	293.1†	278.5†	53.4†	50.7†	84†	58*	44.5†	—	—
77 Ir	311.9†	296.3†	63.8†	60.8†	95.2*	63.0*	48.0†	—	—
78 Pt	331.6†	314.6†	74.5†	71.2†	101†	65.3*	51.7†	—	—
79 Au	353.2†	335.1†	87.6†	83.9†	107.2*	74.2†	57.2†	—	—
80 Hg	378.2†	358.8†	104.0†	99.9†	127†	83.1†	64.5†	—	7.8†
81 Tl	405.7†	385.0†	122.2†	117.8†	136*	94.6†	73.5†	—	12.5†
82 Pb	434.3†	412.2†	141.7†	136.9†	147*	106.4†	83.3†	—	18.1†
83 Bi	464.0†	440.1†	162.3†	157.0†	159.3*	119.0†	92.6†	—	23.8†
84 Po	500*	473*	184*	184*	177*	132*	104*	—	31*
85 At	533*	507*	210*	210*	195*	148*	115*	—	40*
86 Rn	567*	541*	238*	238*	214*	164*	127*	—	48*
87 Fr	603*	577*	268*	268*	234*	182*	140*	—	58*
88 Ra	635.9*	602.7*	299*	299*	254*	200*	153*	—	68*
89 Ac	675*	639*	319*	319*	272*	215*	167*	—	80*
90 Th	712.1†	675.2†	342.4†	333.1	290*	229*	182*	—	85.4†
91 Pa	743*	708*	371*	360*	310*	232*	232*	—	94*
92 U	778.3†	736.2†	388.2*	377.4†	321*	257*	192*	—	94.2†

A compilation by G. P. Williams of Brookhaven National Laboratory, "Electron Binding Energies," in *X-Ray Data Booklet*, Lawrence Berkeley National Laboratory Pub-490 Rev. 2 (2001), based largely on values given by J. A. Bearden and A. F. Barr, "Re-evaluation of X-Ray Atomic Energy Levels," *Rev. Mod. Phys.* **39**:125 (1967); corrected in 1998 by E. Guillemin (LBNL, unpublished). The energies are given in electron volts relative to the vacuum level for the rare gases and for H_2 , N_2 , O_2 , F_2 , and Cl_2 ; relative to the Fermi level for the metals; and relative to the top of the valence bands for semiconductors.

*From M. Cardona and L. Lay (eds.), *Photoemission in Solids I: General Principles*, Springer-Verlag, Berlin (1978).

†From J. C. Fuggle and N. Mårtensson, "Core-Level Binding Energies in Metals," *J. Electron. Spectrosc. Relat. Phenom.* **21**:275 (1980).

For further updates, consult the Web site <http://xdb.lbl.gov/>.

TABLE 2 Photon Energies, in Electronvolts (eV), of Principal K- and L-Shell Emission Lines*

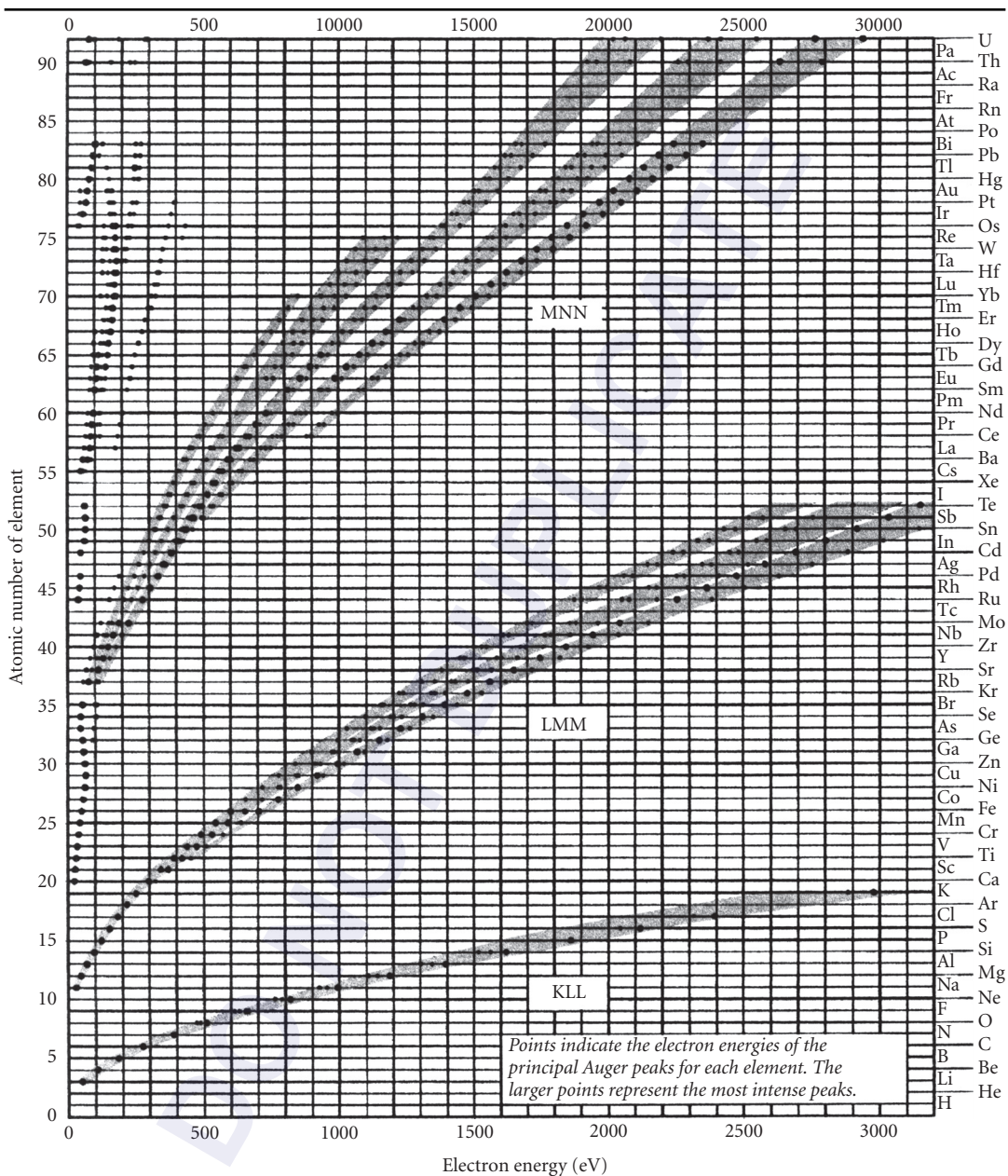
Element	$K\alpha_1$	$K\alpha_2$	$K\beta_1$	$L\alpha_1$	$L\alpha_2$	$L\beta_1$	$L\beta_2$	$L\gamma_1$
3 Li	54.3							
4 Be	108.5							
5 B	183.3							
6 C	277							
7 N	392.4							
8 O	524.9							
9 F	676.8							
10 Ne	848.6	848.6						
11 Na	1,040.98	1,040.98	1,071.1					
12 Mg	1,253.60	1,253.60	1,302.2					
13 Al	1,486.70	1,486.27	1,557.45					
14 Si	1,739.98	1,739.38	1,835.94					
15 P	2,013.7	2,012.7	2,139.1					
16 S	2,307.84	2,306.64	2,464.04					
17 Cl	2,622.39	2,620.78	2,815.6					
18 Ar	2,957.70	2,955.63	3,190.5					
19 K	3,313.8	3,311.1	3,589.6					
20 Ca	3,691.68	3,688.09	4,012.7	341.3	341.3	344.9		
21 Sc	4,090.6	4,086.1	4,460.5	395.4	395.4	399.6		
22 Ti	4,510.84	4,504.86	4,931.81	452.2	452.2	458.4		
23 V	4,952.20	4,944.64	5,427.29	511.3	511.3	519.2		
24 Cr	5,414.72	5,405.509	5,946.71	572.8	572.8	582.8		
25 Mn	5,898.75	5,887.65	6,490.45	637.4	637.4	648.8		
26 Fe	6,403.84	6,390.84	7,057.98	705.0	705.0	718.5		
27 Co	6,930.32	6,915.30	7,649.43	776.2	776.2	791.4		
28 Ni	7,478.15	7,460.89	8,264.66	851.5	851.5	868.8		
29 Cu	8,047.78	8,027.83	8,905.29	929.7	929.7	949.8		
30 Zn	8,638.86	8,615.78	9,572.0	1,011.7	1,011.7	1,034.7		
31 Ga	9,251.74	9,224.82	10,264.2	1,097.92	1,097.92	1,124.8		
32 Ge	9,886.42	9,855.32	10,982.1	1,188.00	1,188.00	1,218.5		
33 As	10,543.72	10,507.99	11,726.2	1,282.0	1,282.0	1,317.0		
34 Se	11,222.4	11,181.4	12,495.9	1,379.10	1,379.10	1,419.23		
35 Br	11,924.2	11,877.6	13,291.4	1,480.43	1,480.43	1,525.90		
36 Kr	12,649	12,598	14,112	1,586.0	1,586.0	1,636.6		
37 Rb	13,395.3	13,335.8	14,961.3	1,694.13	1,692.56	1,752.17		
38 Sr	14,165	14,097.9	15,835.7	1,806.56	1,804.74	1,871.72		
39 Y	14,958.4	14,882.9	16,737.8	1,922.56	1,920.47	1,995.84		
40 Zr	15,775.1	15,690.9	17,667.8	2,042.36	2,039.9	2,124.4	2,219.4	2,302.7
41 Nb	16,615.1	16,521.0	18,622.5	2,165.89	2,163.0	2,257.4	2,367.0	2,461.8
42 Mo	17,479.34	17,374.3	19,608.3	2,293.16	2,289.85	2,394.81	2,518.3	2,623.5
43 Te	18,367.1	18,250.8	20,619	2,424.0	—	2,536.8	—	—
44 Ru	19,279.2	19,150.4	21,656.8	2,558.55	2,554.31	2,683.23	2,836.0	2,964.5
45 Rh	20,216.1	20,073.7	22,723.6	2,696.74	2,692.05	2,834.41	3,001.3	3,143.8
46 Pd	21,177.1	21,020.1	23,818.7	2,838.61	2,833.29	2,990.22	3,171.79	3,328.7
47 Ag	22,162.92	21,990.3	24,942.4	2,984.31	2,978.21	3,150.94	3,347.81	3,519.59
48 Cd	23,173.6	22,984. j	26,095.5	3,133.73	3,126.91	3,316.57	3,528.12	3,716.86
49 In	24,209.7	24,002.0	27,275.9	3,286.94	3,279.29	3,487.21	3,713.81	3,920.81
50 Sn	25,271.3	25,044.0	28,486.0	3,443.98	3,435.42	3,662.80	3,904.86	4,131.12
51 Sb	26,359.1	26,110.8	29,725.6	3,604.72	3,595.32	3,843.57	4,100.78	4,347.79
52 Te	27,472.3	27,201.7	30,995.7	3,769.33	3,758.8	4,029.58	4,301.7	4,570.9

(Continued)

TABLE 2 Photon Energies, in Electronvolts (eV), of Principal K- and L-Shell Emission Lines* (Continued)

Element	$K\alpha_1$	$K\alpha_2$	$K\beta_1$	$L\alpha_1$	$L\alpha_2$	$L\beta_1$	$L\beta_2$	$L\gamma_1$
53 I	28,612.0	28,317.2	32,294.7	3,937.65	3,926.04	4,220.72	4,507.5	4,800.9
54 Xe	29,779	29,458	33,624	4,109.9	—	—	—	—
55 Cs	30,972.8	30,625.1	34,986.9	4,286.5	4,272.2	4,619.8	4,935.9	5,280.4
56 Ba	32,193.6	31,817.1	36,378.2	4,466.26	4,450.90	4,827.53	5,156.5	5,531.1
57 La	33,441.8	33,034.1	37,801.0	4,650.97	4,634.23	5,042.1	5,383.5	5,788.5
58 Ce	34,719.7	34,278.9	39,257.3	4,840.2	4,823.0	5,262.2	5,613.4	6,052
59 Pr	36,026.3	35,550.2	40,748.2	5,033.7	5,013.5	5,488.9	5,850	6,322.1
60 Nd	37,361.0	36,847.4	42,271.3	5,230.4	5,207.7	5,721.6	6,089.4	6,602.1
61 Pm	38,724.7	38,171.2	43,826	5,432.5	5,407.8	5,961	6,339	6,892
62 Sm	40,118.1	39,522.4	45,413	5,636.1	5,609.0	6,205.1	6,586	7,178
63 Eu	41,542.2	40,901.9	47,037.9	5,845.7	5,816.6	6,456.4	6,843.2	7,480.3
64 Gd	42,996.2	42,308.9	48,697	6,057.2	6,025.0	6,713.2	7,102.8	7,785.8
65 Tb	44,481.6	43,744.1	50,382	6,272.8	6,238.0	6,978	7,366.7	8,102
66 Dy	45,998.4	45,207.8	52,119	6,495.2	6,457.7	7,247.7	7,635.7	8,418.8
67 Ho	47,546.7	46,699.7	53,877	6,719.8	6,679.5	7,525.3	7,911	8,747
68 Er	49,127.7	48,221.1	55,681	6,948.7	6,905.0	7,810.9	8,189.0	9,089
69 Tm	50,741.6	49,772.6	57,517	7,179.9	7,133.1	8,101	8,468	9,426
70 Yb	52,388.9	51,354.0	59,377	7,415.6	7,367.3	8,401.8	8,758.8	9,780.1
71 Lu	54,069.8	52,965.0	61,283	7,655.5	7,604.9	8,709.0	9,048.9	10,143.4
72 Hf	55,790.2	54,611.4	63,234	7,899.0	7,844.6	9,022.7	9,347.3	10,515.8
73 Ta	57,532	56,277	65,223	8,146.1	8,087.9	9,343.1	9,651.8	10,895.2
74 W	59,318.24	57,981.7	67,244.3	8,397.6	8,335.2	9,672.35	9,961.5	11,285.9
75 Re	61,140.3	59,717.9	69,310	8,652.5	8,586.2	10,010.0	10,275.2	11,685.4
76 Os	63,000.5	61,486.7	71,413	8,911.7	8,841.0	10,355.3	10,598.5	12,095.3
77 Ir	64,895.6	63,286.7	73,560.8	9,175.1	9,099.5	10,708.3	10,920.3	12,512.6
78 Pt	66,832	65,112	75,748	9,442.3	9,361.8	11,070.7	11,250.5	12,942.0
79 Au	68,803.7	66,989.5	77,984	9,713.3	9,628.0	11,442.3	11,584.7	13,381.7
80 Hg	70,819	68,895	80,253	9,988.8	9,897.6	11,822.6	11,924.1	13,830.1
81 Tl	72,871.5	70,831.9	82,576	10,268.5	10,172.8	12,213.3	12,271.5	14,291.5
82 Pb	74,969.4	72,804.2	84,936	10,551.5	10,449.5	12,613.7	12,622.6	14,764.4
83 Bi	77,107.9	74,814.8	87,343	10,838.8	10,730.91	13,023.5	12,979.9	15,247.7
84 Po	79,290	76,862	89,800	11,130.8	11,015.8	13,447	13,340.4	15,744
85 At	8,152	7,895	9,230	11,426.8	11,304.8	13,876	—	16,251
86 Rn	8,378	8,107	9,487	11,727.0	11,597.9	14,316	—	16,770
87 Fr	8,610	8,323	9,747	12,031.3	11,895.0	14,770	1,445	17,303
88 Ra	8,847	8,543	10,013	12,339.7	12,196.2	15,235.8	14,841.4	17,849
89 Ac	90,884	8,767	10,285	12,652.0	12,500.8	15,713	—	18,408
90 Th	93,350	89,953	105,609	12,968.7	12,809.6	16,202.2	15,623.7	18,982.5
91 Pa	95,868	92,287	108,427	13,290.7	13,122.2	16,702	16,024	19,568
92 U	98,439	94,665	111,300	13,614.7	13,438.8	17,220.0	16,428.3	20,167.1
93 Np	—	—	—	13,944.1	13,759.7	17,750.2	16,840.0	20,784.8
94 Pu	—	—	—	14,278.6	14,084.2	18,293.7	17,255.3	21,417.3
95 Am	—	—	—	14,617.2	14,411.9	18,852.0	17,676.5	22,065.2

*Photon energies in electronvolts (eV) of some characteristic emission lines of the elements of atomic number $3 \leq Z \leq 95$, as compiled by J. Kortright. "Characteristic X-Ray Energies," in *X-Ray Data Booklet* (Lawrence Berkeley National Laboratory Pub-490, Rev. 2, 1999). Values are largely based on those given by J. A. Bearden, "X-Ray Wavelengths," *Rev. Mod. Phys.* **39**:78 (1967), which should be consulted for a more complete listing. Updates may also be noted at the Web site www.cxro.lbl.gov/optical_constants.

TABLE 3 Curves Showing Auger Energies,* in Electronvolts (eV), for Elements of Atomic Number $3 \leq Z \leq 92$ 

*Only dominant energies are given, and only for principal Auger peaks. The literature should be consulted for detailed tabulations, and for shifted values in various common compounds.²⁻⁴ (Courtesy of Physical Electronics, Inc.²)

36.3 REFERENCES

1. B. L. Henke, E. M. Gullikson, and J. C. Davis. "X-Ray Interactions: Photoabsorption, Scattering, Transmission, and Reflection at $E = 50\text{--}3000$ eV, $Z = 1\text{--}92$." *Atomic Data and Nuclear Data Tables* **54**(2):181–342 (July 1993).
2. K. D. Childs, B. A. Carlson, L. A. Vanier, J. F. Moulder, D. F. Paul, W. F. Stickle, and D. G. Watson. *Handbook of Auger Electron Spectroscopy*, C. L. Hedberg (ed.), Physical Electronics, Eden Prairie, MN, 1995.
3. J. F. Moulder, W. F. Stickle, P. E. Sobol, and K. D. Bomben, *Handbook of X-Ray Photoelectron Spectroscopy*, Physical Electronics, Eden Prairie, MR 1995.
4. D. Briggs, *Handbook of X-Ray and Ultraviolet Photoelectron Spectroscopy*, Heyden, London, 1977.

DO NOT DUPLICATE

SUBPART

5.2

REFRACTIVE AND INTERFERENCE OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

Bruno Lengeler

*Physikalisches Institut
RWTH Aachen University
Aachen, Germany*

Christian G. Schroer

*Institute of Structural Physics
TU Dresden
Dresden, Germany*

37.1 INTRODUCTION

The last ten years have seen a remarkable progress in the development of new x-ray optics and in the improvement of existing devices. In this chapter, we describe the properties of one type of these new optics: refractive x-ray lenses.

For a long time these lenses were considered as not feasible, due to the weak refraction and the relatively strong absorption of x rays in matter. However, in 1996 it was shown experimentally that focusing by x-ray lenses is possible if the radius of curvature R of an individual lens is chosen to be small (e.g., below 0.5 mm, cf. Fig. 1a), if many such lenses are stacked behind one another in a row, and if a lens material with low atomic number Z , such as aluminium, is chosen.^{1,2} The first lenses of this type consisted of a row of holes, 1 mm in diameter, drilled in a block of aluminium.

In the meantime, many different types of refractive lenses made of various materials have been developed.^{3–19} One of the most important developments was to make these optics aspherical,⁴ reducing spherical aberration to a minimum and thus making these optics available for high-resolution x-ray microscopy. As each individual lens is thin in the optical sense, the ideal aspherical shape is a paraboloid of rotation. In the following, we focus on two types of high resolution x-ray optics, rotationally parabolic lenses made of beryllium and aluminium^{6,7,13,14,16,20} and cylindrically parabolic lenses with particularly short focal length.^{21,22} An example of the first type of lenses developed and made at Aachen University is described in Secs. 37.2 to 37.5. These lenses allow x-ray imaging nearly free of distortions and can be used as an objective lens in an x-ray microscope for efficient focusing in scanning microscopy, and for a variety of beam conditioning applications at third generation synchrotron radiation sources. The second type of lenses, so-called nanofocusing lenses (NFLs), have a short focal distance and large numerical aperture and are thus particularly suited to focus hard x rays to sub-100 nm dimensions for scanning microscopy applications. While focusing of hard x rays down to 50 nm has been demonstrated experimentally, these optics have the potential of focusing hard x rays to about 10 nm^{21,22} and perhaps below.²³ The development of nanofocusing lenses currently pursued at TU Dresden is described in Sec. 37.6.

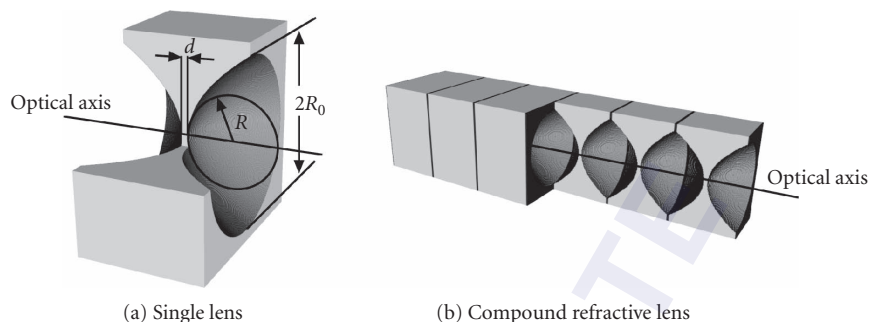


FIGURE 1 (a) Individual refractive x-ray lens with rotationally parabolic profile and (b) stack of individual lenses forming a refractive x-ray lens. (Reused with permission from Ref. 24.)

37.2 REFRACTIVE X-RAY LENSES WITH ROTATIONALLY PARABOLIC PROFILE

Refractive x-ray lenses with rotationally parabolic profiles^{6,7,13,16} allow for focusing in both directions, free of spherical aberration and other distortions. Aluminium and beryllium are the lens materials most commonly used. Beryllium is especially suitable for x-ray energies between 7 and 40 keV due to its low attenuation of x rays (low atomic number $Z = 4$). Between about 40 and 90 keV aluminium ($Z = 13$) is more appropriate as a lens material. Having been able to solve the problems with handling and plastically deforming beryllium, Be lenses can be manufactured with rotationally parabolic profiles.^{13,14} Figure 1a shows a schematic drawing of an individual lens. Note the concave shape of a focusing lens, which is a result of the diffractive part $1 - \delta$ of the index of refraction n being smaller than 1 in the x-ray range. In Fig. 1b a number N of individual lenses is stacked behind each other to form a refractive x-ray lens. Figure 2 shows a stack of Be lenses in their casing with protective atmosphere.

In the thin lens approximation, the focal length of a stack of N lenses is $f_0 = R/2N\delta$. Here, R is the radius of curvature at the apex of the paraboloid (cf. Fig. 1a). For paraboloids, R and the geometric aperture $2R_0$ are independent of one another, in contrast to the case for spherical lenses. Most lenses up to now had the parameters $R = 0.2$ mm and $2R_0 \approx 1$ mm. Up to several hundred lenses can be aligned

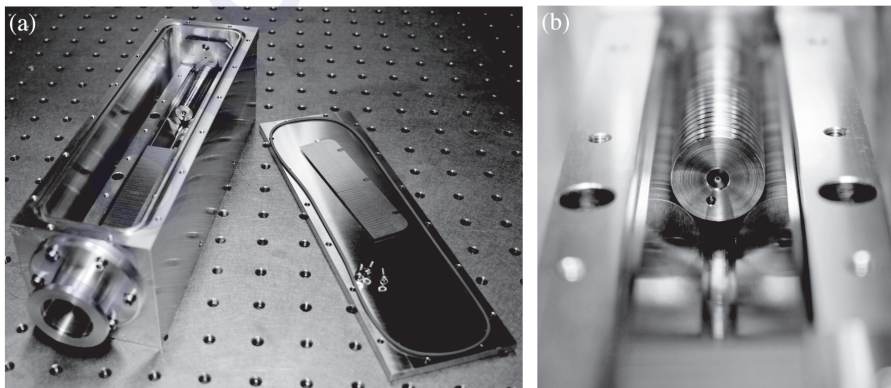


FIGURE 2 (a) Housing with partly assembled Be lens and (b) stack of Be lenses. Each individual lens is centered inside of a hard metal coin. The lenses are aligned along the optical axis by stacking the coins in a high precision v-groove. (See also color insert.)

in a lens stack in such a way that the optical axes of the individual lenses agree on the micrometer scale. The form fidelity of the paraboloids is better than $0.3 \mu\text{m}$ and surface roughness is below $0.1 \mu\text{m}$.

In the meantime, lenses with different radii of curvature R ($R = 50, 100, 200, 300, 500, 1000$, and $1500 \mu\text{m}$) have been developed to optimize the optics for various applications. They are available in three lens materials, i.e., beryllium, aluminium, and nickel. Lenses with small radii R are especially suited for microscopy applications with high lateral resolution, whereas those with a large radius R are designed for beam conditioning purposes, such as prefocusing and collimation.

In general, the total length L of a lens stack is not negligible compared to the focal length f . Then, a correction has to be applied to the thin lens approximation for the focal length. For a thick lens the focal length is

$$f = f_0 \sqrt{\frac{L}{f_0}} \frac{1}{\sin \sqrt{L/f_0}} \quad (1)$$

as measured from the principal planes located at

$$H_{1,2} = \pm \left[f_0 \sqrt{\frac{L}{f_0}} \frac{1 - \cos \sqrt{\frac{L}{f_0}}}{\sin \sqrt{\frac{L}{f_0}}} - \frac{L}{2} \right] \quad (2)$$

behind and before the center of the lens, respectively.

The attenuation of x rays in matter is a key parameter in the design of refractive x-ray lenses. As the thickness of the lens material increases with increasing distance from the optical axis, the lens becomes more and more absorbing toward its periphery. Thus a refractive x-ray lens has no sharp aperture, but a Gaussian transmission that is responsible for the diffraction at the lens. As a result, we can assign an effective aperture D_{eff} to the lens that is smaller than the geometric aperture $2R_0$ ^{6,7} and that determines the diffraction at the lens, its numerical aperture, and the achievable diffraction-limited spot size.

At low energies (below about 10 keV for beryllium), the attenuation is dominated by photoabsorption. The mass photoabsorption coefficient τ/ρ varies approximately like Z^3/E^3 with atomic number Z and with photon energy E . When τ/ρ drops to below about $0.15 \text{ cm}^2/\text{g}$ at higher x-ray energies, the mass attenuation coefficient $\mu = \tau + \mu_C$ is dominated by Compton scattering (μ_C) and stays more or less constant, independent of energy and atomic number Z . For beryllium the cross-over between the photoabsorption and Compton scattering dominated attenuation is at about 17 keV. The performance of beryllium lenses is optimal in this energy range.

Compton scattering ultimately limits the performance (lateral resolution) of refractive x-ray lenses. Compton scattering has a twofold detrimental influence. Photons which are Compton scattered no longer contribute to the image formation. In addition, they generate a background which reduces the signal-to-background-ratio in the image.

Synchrotron radiation sources of the third generation can create a considerable heat load in the first optical element hit by the beam. This is expected to be even more true at x-ray free-electron laser sources which are being developed at present. The compatibility with such a high heat load was tested for refractive lenses made of beryllium at the undulator beamline ID10 at the European Synchrotron Radiation Facility (ESRF) in Grenoble, France. The power density and power of the white beam generated by 3 undulators in a row was about $100 \text{ W}/\text{mm}^2$ and 40 W, respectively. A stack of 12 Be lenses was exposed to the beam. The lenses were housed in an evacuated casing. They were indirectly cooled via a thermal link to a copper plate which in turn was water cooled. The temperature was measured by three thermocouples, one at each end and one at the center of the lens stack. The highest temperature was measured at the center, increasing within a few minutes to 65°C and staying constant afterward, except for small variations due to changes of the electron current in the ring. A temperature of 65°C poses no problem for Be lenses. The melting point of Be is at 1285°C and recrystallization of Be occurs only above 600°C . At present, rotationally parabolic Be lenses have been installed in the front ends of several undulator beamlines at ESRF, being routinely used in the undulator "white" beam. At

the present undulator beamlines no deterioration of Be or Al x-ray lenses has been observed in the monochromatic beam, even after many years of operation. In terms of stability metallic lens materials are far superior to insulators, like plastics or glass. The high density of free electrons in metals prevents radiation damage by bond breaking or local charging.

The heat load resistance of these optics is of utmost importance for focusing applications at future x-ray free-electron lasers. Model calculations suggest that these optics are stable in the hard x-ray beam (8 to 12 keV) generated by x-ray free-electron lasers, such as the LCLS in Stanford and the future European X-Ray Free-Electron Laser Project XFEL in Hamburg.^{24–27}

37.3 IMAGING WITH PARABOLIC REFRACTIVE X-RAY LENSES

Refractive x-ray lenses with parabolic profile are especially suited for hard x-ray full-field microscopy since they are relatively free of distortions compared to crossed lenses with cylindrical symmetry. For this purpose, a refractive x-ray lens is placed a distance L_1 behind the object that is illuminated from behind by monochromatic synchrotron radiation. The image of the object is formed at a distance $L_2 = L_1 f / (L_1 - f)$ behind the lens on a position-sensitive detector. To achieve large magnifications $M = L_2 / L_1 = f / (L_1 - f)$ up to 100, L_1 should be chosen to be slightly larger than the focal distance f . Figure 3 shows the image of a Ni mesh (periodicity 12.7 μm) imaged with a Be lens ($N = 91$, $f = 493$ mm, $L_2 / L_1 = 10$) at 12 keV onto high resolution x-ray film. Details of the contrast formation are described in Refs. 28 and 29. There are no apparent distortions visible in the image. This is a consequence of the parabolic lens profile. Figure 4 compares imaging with parabolic and spherical lenses in a numerical simulation. Spherical aberration dominates the image formed by the spherical lens, clearly demonstrating the need for a lens surface in the form of a paraboloid of rotation. A comparison of Fig. 3 with Fig. 4a shows that the experimental result is very close to a numerical calculation with idealized parabolic lenses.

The first x-ray microscope of this kind was built using aluminium lenses.⁶ However, using beryllium as a lens material rather than aluminium has several advantages. The reduced attenuation inside the lens results in a larger effective aperture that leads to a higher spatial resolution and a larger field of view. In addition, the efficiency of the setup is improved, since the transmission of the lens is higher.

For a refractive lens with a parabolic profile the lateral resolution is limited by the diffraction at its Gaussian aperture, giving rise to a Gaussian shape of the Airy disc.⁷ The full width at half maximum of the Airy disc is given by

$$d_t = 0.75 \frac{\lambda}{2NA} = 0.75 \frac{\lambda L_1}{D_{\text{eff}}} \quad (3)$$

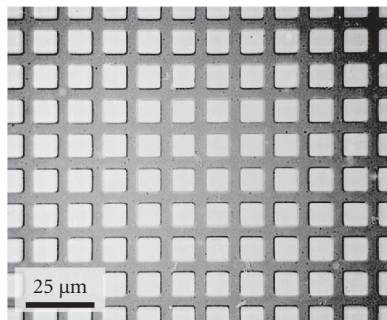


FIGURE 3 Hard x-ray micrograph of a Ni mesh.¹⁴ (Reused with permission from Ref. 20.)

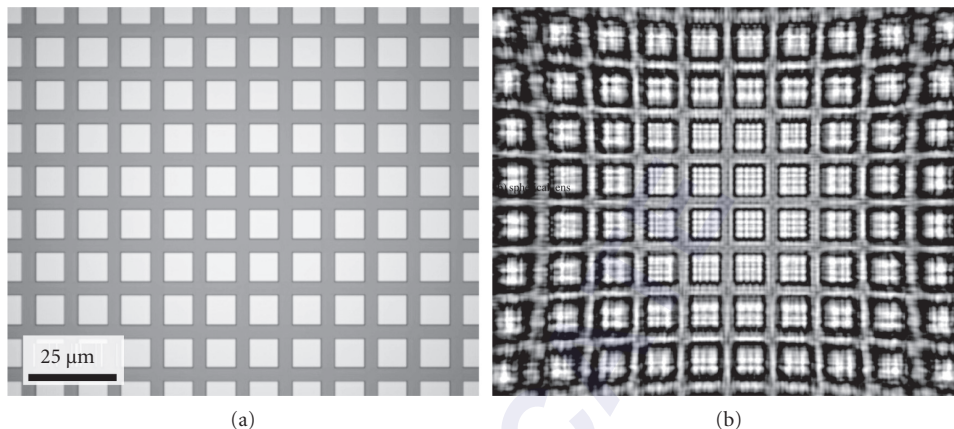


FIGURE 4 Numerical simulation of the imaging process using (a) parabolic and (b) spherical lens.

The numerical aperture NA is defined by $\sin \alpha$, where 2α is the angle spanned by the effective aperture D_{eff} of the lens as seen from an object point.⁷ This result is well known from optics, the factor 0.75 being different from the usual factor 1.22. The difference can be traced back to the fact that in normal optics apertures are sharply delimited whereas for x-ray lenses the attenuation changes smoothly as described earlier. The effective aperture D_{eff} is limited by x-ray attenuation and that is ultimately limited by Compton scattering. An estimate of d_t for lenses with large apertures shows that it scales with

$$d_t = a\lambda \sqrt{\frac{\mu f}{\delta}}$$

where a is a factor of order one. This implies that a low value for d_t needs a small focal length f and a low mass attenuation coefficient μ , in other words a low Z material. Since δ is proportional to λ^2 the main x-ray energy dependence enters via μ . With the present day technology for fabrication of refractive lenses with rotationally parabolic profile, focal lengths between 10 and 20 cm can be achieved for energies between 10 and 20 keV, resulting for Be lenses in a lateral resolution down to about 50 nm. We estimate that it will be difficult to reach values below 30 nm.

The main strength of the x-ray microscope is the large penetration depth of hard x rays in matter that allows one to investigate non-destructively inner structures of an object. In combination with tomographic techniques, it allows one to reconstruct the three-dimensional inner structure of the object with submicrometer resolution.³⁰ In addition, full-field imaging in demagnifying geometry can be used for hard x-ray lithography.³¹ The high quality of refractive lenses is also demonstrated by the preservation of the lateral coherence.^{16,32}

37.4 MICROFOCUSING WITH PARABOLIC REFRACTIVE X-RAY LENSES

Refractive x-ray lenses with parabolic profiles can also be used for generating a (sub-)micrometer focal spot for x-ray microanalysis and tomography. For that purpose, the synchrotron radiation source is imaged by the lens onto the sample in a strongly demagnifying way, i.e., the source-lens distance L_1 is chosen to be much larger than the lens-sample distance L_2 . At a synchrotron radiation source the horizontal source size is typically larger than the vertical one. As the lens images this horizontally elongated source to the sample position, the focal spot is larger in the horizontal direction than in the vertical direction. With Be lenses ($R = 200 \mu\text{m}$), a vertical spot size well below $1 \mu\text{m}$ is

routinely achieved, while the horizontal spot size is typically limited to a few micrometers by the horizontal source size and the demagnification of the setup. The diffraction limit of these optics is usually not reached in typical microfocusing geometries, i.e., 40 to 70 m from a typical undulator source at a synchrotron radiation source of the third generation. The spot size is dominated by the geometric image of the source and diffraction at the lens aperture and aberrations are negligible.

By means of new Be lenses with smaller radii of curvature, e.g., $R = 50 \mu\text{m}$, a focal length of 15 cm can be reached, thus resulting in a demagnification of the source size by about a factor 400 in 60-m distance from the source. At a low- β undulator source, this demagnification allows one to reach the sub-micrometer regime also in the horizontal direction. Close to diffraction-limited focusing, however, is still only possible at very long beamlines. For example, at a distance of $L_1 = 145 \text{ m}$ from a low- β source at the ESRF [effective source size $60 \times 125 \mu\text{m}$ ($V \times H$)], a microbeam of $60 \times 125 \text{ nm}^2$ is expected, approaching the diffraction limit of these optics in the vertical direction. To generate foci well below 100 nm at short distances from the source, focal distances in the centimeter range are needed to generate large demagnifications. This can be achieved with the nanofocusing lenses described in Sec. 37.6.^{21,22}

Hard x-ray microbeams find a large number of applications in scanning microscopy and have been used for a variety of experiments. They include, for example, microdiffraction,³³ microfluorescence mapping³⁴ and tomography,^{35–37} and x-ray absorption spectroscopic³⁸ and small-angle x-ray scattering tomography.³⁹ In materials science there is a great interest in using very hard x rays above about 80 keV, because many samples with high Z metallic components and thicknesses of many millimeters show strong x-ray absorption. At 80 keV parabolic aluminium lenses (preferably with $R = 50 \mu\text{m}$) are well suited.⁴⁰ For higher x-ray energies, lenses made of nickel become advantageous, due to the strong refraction in nickel resulting from its relatively high density ($\rho_{\text{Ni}} = 8.9 \text{ g/cm}^3$). For energies above 80 keV, parabolic cylinder lenses made with a LIGA technique have been successfully tested.⁴¹ Rotationally parabolic nickel lenses are in the process of development. The challenge is to produce lenses with minimal thickness d (cf. Fig. 1a) in order to minimize absorption in the stack of lenses. A value of $d = 10 \mu\text{m}$ is tolerable and feasible.

Be refractive lenses appear to be well suited to focus the beam from an x-ray free-electron-laser.²⁴

37.5 PREFOCUSING AND COLLIMATION WITH PARABOLIC REFRACTIVE X-RAY LENSES

Most experimental setups on synchrotron radiation beamlines are located between 40 and 150 m from the source. Depending on the experiment, the beam size, flux, divergence, or lateral coherence length may not be optimal at the position of the experiment. Using appropriate lenses upstream of the experiment, a given parameter can be optimized.

For example, the divergence of the beam may lead to a significant reduction of the flux at the sample position, in particular at a low- β undulator source. In this case, the beam can be moderately focused with refractive lenses with large radius of curvature and thus large geometric and effective aperture. For instance, lenses with $R = 1500 \mu\text{m}$ have a geometric aperture $2R_0$ of 3 mm. In a one-to-one imaging geometry at 50 m from the source they have an angular acceptance of $85 \mu\text{rad}$ at 17 keV, thus capturing a large fraction of the beam in horizontal direction. This opens excellent possibilities to increase the photon flux, in particular as the lenses can be easily moved in and out of the beam without affecting the optical axis and the alignment of the experiment. In the new future, parabolic cylinder lenses made of Be and Al will also become available. With a height of 3.5 mm and radii of curvature between 200 and 1500 μm , they can be used for one-dimensional focusing and collimation.

37.6 NANOFOCUSING REFRACTIVE X-RAY LENSES

High quality magnified imaging with x rays requires optical components free of distortion, like rotationally parabolic refractive x-ray lenses. However, for technical reasons, their radii of curvature cannot be made smaller than about $50 \mu\text{m}$. This limits the focal distance from below and thus the achievable

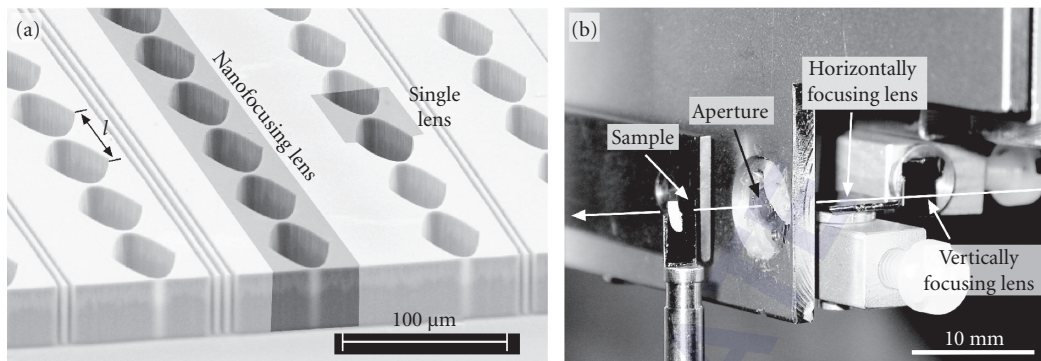


FIGURE 5 (a) Array of nanofocusing lenses made of silicon. A large number of single lenses are aligned behind each other to form a nanofocusing lens. Several nanofocusing lenses with different radius of curvature R are placed in parallel onto the same substrate. (b) Scanning microprobe setup with two crossed nanofocusing lenses. An aperture defining pinhole is placed behind the second lens. (See also color insert.)

demagnification in microfocus experiments. Therefore, another approach has been pursued for the generation of particularly small focal spots. These are nanofocusing cylinder lenses with parabolic profile and a focal length of the order of 1 cm.^{21,22} This is achieved by choosing the radius of curvature of the parabolas as small as 1 to 5 μm . Figure 5a shows an array of nanofocusing lenses made of silicon. When two lenses are used in crossed geometry as shown in Fig. 5b, two-dimensional focusing can be achieved. So far, focal spot sizes down to $47 \times 55 \text{ nm}^2$ ($H \times V$) have been reached at $L_1 = 47 \text{ m}$ from the low- β undulator source at beamline ID13 of the ESRF ($E = 21 \text{ keV}$).²²

While this spot size is close to the ideal performance of silicon lenses in this particular imaging geometry, significant improvements can be made in the future by further optimization of the optics and the imaging geometry. Figure 6 shows the optimal diffraction limit as a function of x-ray energy for different lens materials. Over a wide range of energies (from $E = 8 \text{ keV}$ to over $E = 100 \text{ keV}$, a diffraction

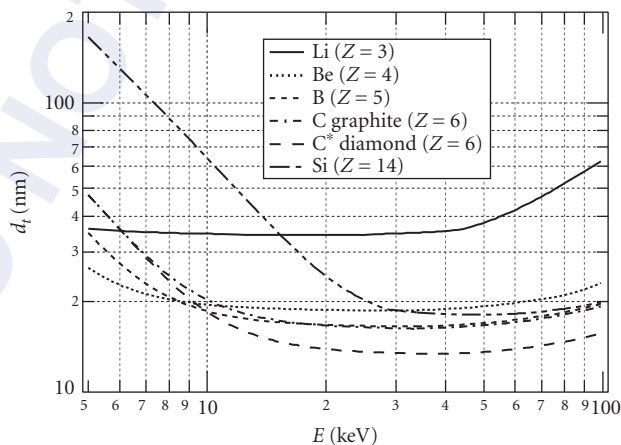


FIGURE 6 Minimal diffraction limits of nanofocusing lenses made of different lens materials and having a working distance of 1 mm. The radius of curvature R and the length of a single lens l were varied within a range accessible by modern microfabrication techniques.²¹ (Copyright 2003 by the American Institute of Physics. Reused with permission from Ref. 21.)

limit below 20 nm is expected. Best performance is obtained for low Z materials with high density. The reason for this is that attenuation is no longer limiting the aperture of nanofocusing lenses for low Z materials, as the overall length of the lens is short. For a given focal length, the geometric aperture is, however, limited by the refractive strength per unit length inside the lens. The higher δ , the larger can be the radius of curvature R and thus the geometric aperture $R_0 = \sqrt{R(l-d)}$, if the thickness l of an individual lens is kept constant (cf. Fig. 5a). In the limit, the numerical aperture approaches $\sqrt{2\delta}$ that coincides with the critical angle of total reflection. At highly brilliant sources, such as the ESRF, these diffraction limits are expected to be reached with fluxes above 10^9 photons per second.

For these optics, prefocusing as described in Sec. 37.5 is of utmost importance to obtain optimal performance. Optimal diffraction-limited focusing is obtained when the lateral coherence length at the optic is slightly larger than the effective aperture. This requirement is usually not fulfilled at the position of the experiment. By appropriate prefocusing, the lateral coherence length can be adapted to the aperture of the NFL, thus optimally focusing the coherent flux from the source onto the sample. This scheme is pursued in modern hard x-ray scanning microscopes, both at ESRF and at the future synchrotron radiation source PETRA III at DESY in Hamburg, Germany.

For refractive lenses made of identical single lenses, the numerical aperture is fundamentally limited by the critical angle of total reflection $\sqrt{2\delta}$. This limitation can be overcome with refractive optics by gradually (adiabatically) adjusting the aperture of the individual lenses to the converging beam inside the optic. For these so-called adiabatically focusing lenses, the numerical aperture can exceed $\sqrt{2\delta}$ leading to diffraction limits well below 10 nm.²³

The main applications of nanofocusing lenses lie in scanning microscopy and microanalysis with hard x rays. They allow one to perform x-ray analytical techniques, such as diffraction,⁴² fluorescence analysis, and absorption spectroscopy, with high spatial resolution. Also, coherent x-ray diffraction imaging greatly benefits from focusing the coherent beam with NFLs.⁴³ The current performance of a hard x-ray scanning microscope based on nanofocusing lenses is illustrated in Fig. 7. In collaboration with W. H. Schröder from the Research Center Jülich the distribution of physiologically relevant ions and heavy metals was mapped inside the tip of a leaf hair (trichome) of the model plant *Arabidopsis thaliana*. Figure 7c shows the two-dimensional map of a variety of elements obtained by scanning the

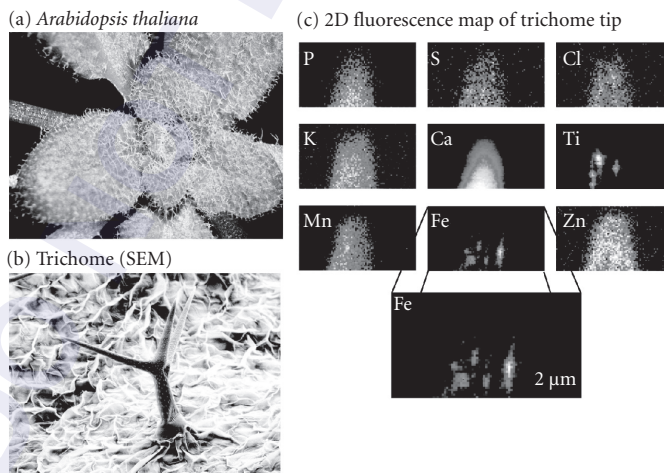


FIGURE 7 (a) Photograph of the plant *Arabidopsis thaliana*, (b) secondary electron micrograph of a leaf hair (trichome), and (c) two-dimensional fluorescence map of the tip of the trichome at 100-nm spatial resolution. While most elements are homogeneously distributed, iron (Fe) and titanium (Ti) are localized on the level of 100 nm. (See also color insert.) (Sample provided by W. H. Schröder, Research Center Jülich.)

tip of a trichome with a hard x-ray nanobeam ($E = 15$ keV). The step size was 100 nm in both dimensions, clearly showing a strong localization of iron and titanium. While the reason for this localization remains unknown, it impressively demonstrates the high spatial resolution obtained with nanofocusing lenses. While these optics are ideal for microbeam applications, they are not well suited for high quality full-field imaging due to distortions in the image due to the crossing to two cylinder lenses with different focal lengths.

37.7 CONCLUSION

Since their first experimental realization about one decade ago, refractive x-ray lenses have developed into a high-quality x-ray optic. Similar to glass lenses for visible light, they have a broad range of applications and can be used in very much the same way. Due to their good imaging properties, refractive optics are particularly suited for hard x-ray microscopy and microanalysis. Due to the weak refraction of hard x rays in matter, they are generally slim, operating in the paraxial regime with typical numerical apertures below a few times 10^{-3} . They can be used in the whole hard x ray range from about five to several hundred keV. As the refractive index depends on energy, refractive lenses are chromatic. Thus, they are mostly used with monochromatic radiation. Today, spatial resolutions down to 50 nm have been reached in hard x-ray microscopy. Potentially, these optics can generate hard x-ray beams down to below 10 nm. Their straight optical path makes them easy to use and align and enhances the stability of x-ray microscopes, as angular instabilities do not affect the focus. In addition, they are extremely robust, both mechanically and thermally. Therefore, they can be used as front end optics at third-generation synchrotron radiation sources and are good candidates to focus the radiation from free-electron lasers. Today, refractive x-ray lenses are routinely used at many beamlines of different synchrotron radiation sources.

37.8 REFERENCES

1. A. Snigirev, V. Kohn, I. Snigireva, and B. Lengeler, "A Compound Refractive Lens for Focusing High Energy X-Rays," *Nature (London)* **384**:49 (1996).
2. B. Lengeler, J. Tümmler, A. Snigirev, I. Snigireva, and C. Raven, "Transmission and Gain of Singly and Doubly Focusing Refractive X-Ray Lenses," *J. Appl. Phys.* **84**:5855–5861 (1998).
3. A. Snigirev, B. Filseth, P. Elleaume, T. Klocke, V. Kohn, B. Lengeler, I. Snigireva, A. Souvorov, and J. Tümmler, "Refractive Lenses for High Energy X-Ray Focusing," In A. M. K. A. T. Macrander, ed., "High Heat Flux and Synchrotron Radiation Beamlines," *Proc. SPIE*, **3151**:164–170 (1997).
4. B. Lengeler, "Linsensysteme für Röntgenstrahlen," *Spektrum der Wissenschaft* 25–30 (1997).
5. A. Snigirev, V. Kohn, I. Snigireva, A. Souvorov, and B. Lengeler, "Focusing High-Energy X-Rays by Compound Refractive Lenses," *Appl. Opt.* **37**:653–662 (1998).
6. B. Lengeler, C. G. Schroer, M. Richwin, J. Tümmler, M. Drakopoulos, A. Snigirev, and I. Snigireva, "A Microscope for Hard X-Rays Based on Parabolic Compound Refractive Lenses," *Appl. Phys. Lett.* **74**:3924–3926 (1999).
7. B. Lengeler, C. Schroer, J. Tümmler, B. Benner, M. Richwin, A. Snigirev, I. Snigireva, and M. Drakopoulos, "Imaging by Parabolic Refractive Lenses in the Hard X-Ray Range," *J. Synchrotron Rad.* **6**:1153–1167 (1999).
8. Y. Kohmura, M. Awaji, Y. Suzuki, T. Ishikawa, Y. I. Dudchik, N. N. Kolchewsky, and F. F. Komarow, "X-Ray Focusing Test and X-Ray Imaging Test by a Microcapillary X-Ray Lens at an Undulator Beamline," *Rev. Sci. Instrum.* **70**:4161–4167 (1999).
9. J. T. Cremer, M. A. Piestrup, H. R. Beguiristain, C. K. Gary, R. H. Pantell, and R. Tatchyn, "Cylindrical Compound Refractive X-Ray Lenses Using Plastic Substrates," *Rev. Sci. Instrum.* **70** (1999).
10. B. Cederström, R. N. Cahn, M. Danielsson, M. Lundqvist, and D. R. Nygren, "Focusing Hard X-Rays with Old LP's," *Nature* **404**:951 (2000).
11. V. Aristov, M. Grigoriev, S. Kuznetsov, et al., "X-Ray Refractive Planar Lens with Minimized Absorption," *Appl. Phys. Lett.* **77**:4058–4060 (2000).

12. E. M. Dufresne, D. A. Arms, R. Clarke, N. R. Pereira, S. B. Dierker, and D. Foster, "Lithium Metal for X-Ray Refractive Optics," *Appl. Phys. Lett.* **79**:4085–4087 (2001).
13. B. Lengeler, C. G. Schroer, B. Benner, A. Gerhardus, T. F. Günzler, M. Kuhlmann, J. Meyer, and C. Zimprich, "Parabolic Refractive X-Ray Lenses," *J. Synchrotron Rad.* **9**:119–124 (2002).
14. C. G. Schroer, M. Kuhlmann, B. Lengeler, T. F. Günzler, O. Kurapova, B. Benner, C. Rau, A. S. Simionovici, A. Snigirev, and I. Snigireva, "Beryllium Parabolic Refractive X-Ray Lenses," In D. C. Mancini, ed., "Design and Microfabrication of Novel X-Ray Optics," *Proc. SPIE*, **4783**:10–18, SPIE, Bellingham (2002).
15. B. Cederström, M. Lundqvist, and C. Ribbing, "Multi-Prism X-Ray Lens," *Appl. Phys. Lett.* **81**:1399–1401 (2002).
16. B. Lengeler, C. G. Schroer, M. Kuhlmann, B. Benner, T. F. Günzler, O. Kurapova, A. Somogyi, A. Snigirev, and I. Snigireva, "Beryllium Parabolic Refractive X-Ray Lenses," In T. Warwick, J. Arthur, H. A. Padmore, and J. Stohr, eds., *Synchrotron Radiation Instrumentation, Proc. AIP Conference*, **705**:748–751 (2004).
17. B. Nöhammer, J. Hoszowska, A. K. Freund, and C. David, "Diamond Planar Refractive Lenses for Third- and Forth-Generation X-Ray Sources," *J. Synchrotron Rad.* **10**:168–171 (2003).
18. V. Nazmov, E. Reznikova, M. Boerner, et al., "Refractive Lenses Fabricated by Deep SR Lithography and LIGA Technology for X-Ray Energies from 1 keV to 1 MeV," In T. Warwick, J. Arthur, H. A. Padmore, and J. Stöhr, eds., *Synchrotron Radiation Instrumentation, Proc. AIP Conference*, **705**:752–755 (2004).
19. V. Nazmov, E. Reznikova, A. Somogyi, J. Mohr, and V. Saile, "Planar Sets of Cross X-Ray Refractive Lenses from SU-8 Polymer," In A. S. Snigirev and D. C. Mancini, eds., "Design and Microfabrication of Novel X-Ray Optics II," *Proc. SPIE*, **5539**:235–243 (2004).
20. B. Lengeler, C. G. Schroer, M. Kuhlmann, B. Benner, T. F. Günzler, O. Kurapova, F. Zontone, A. Snigirev, and I. Snigireva, "Refractive X-Ray Lenses," *J. Phys. D: Appl. Phys.* **38**:A218–A222 (2005).
21. C. G. Schroer, M. Kuhlmann, U. T. Hunger, et al., "Nanofocusing Parabolic Refractive X-Ray Lenses," *Appl. Phys. Lett.* **82**:1485–1487 (2003).
22. C. G. Schroer, O. Kurapova, J. Patommel, et al., "Hard X-Ray Nanoprobe Based on Refractive X-Ray Lenses," *Appl. Phys. Lett.* **87**:124103 (2005).
23. C. G. Schroer and B. Lengeler, "Focusing Hard X Rays to Nanometer Dimensions by Adiabatically Focusing Lenses," *Phys. Rev. Lett.* **94**:054802 (2005).
24. C. G. Schroer, J. Tümmler, B. Lengeler, M. Drakopoulos, A. Snigirev, and I. Snigireva, "Compound Refractive Lenses: High Quality Imaging Optics for the XFEL," In D. M. Mills, H. Schulte-Schrepping, and J. R. Arthur, eds., "X-Ray FEL Optics and Instrumentation," *Proceedings of the SPIE*, **4143**:60–68 (2001).
25. G. Materlik and T. Tschentscher, "TESLA Technical Design Report, Part V, The X-ray Free Electron Laser," *Tech. Rep.* DESY 2001-011, DESY, Hamburg (2001).
26. R. M. Bionta, "Controlling Dose to Low Z Solids at LCLS," *Tech. Rep.* Lawrence Livermore National Laboratory, UCRL-ID-137222, January 3, 2000, LCLS-TN-00-4 (2000).
27. C. G. Schroer, B. Benner, M. Kuhlmann, O. Kurapova, B. Lengeler, F. Zontone, A. Snigirev, I. Snigireva, and H. Schulte-Schrepping, "Focusing Hard X-Ray FEL Beams with Parabolic Refractive Lenses," In S. G. Biedron, W. Eberhardt, T. Ishikawa, and R. O. Tatchyn, eds., "Fourth Generation X-Ray Sources and Optics II," *Proceedings of the SPIE*, **5534**:116–124 (2004).
28. C. G. Schroer, B. Benner, T. F. Günzler, M. Kuhlmann, B. Lengeler, C. Rau, T. Weitkamp, A. Snigirev, and I. Snigireva, "Magnified Hard X-Ray Microtomography: Toward Tomography with Sub-Micron Resolution," In U. Bonse, ed., "Developments in X-Ray Tomography III," *Proceedings of the SPIE*, **4503**:23–33 (2002).
29. V. Kohn, I. Snigireva, and A. Snigirev, "Diffraction Theory of Imaging with X-Ray Compound Refractive Lens," *Opt. Commun.* **216**:247–260 (2003).
30. C. G. Schroer, J. Meyer, M. Kuhlmann, B. Benner, T. F. Günzler, B. Lengeler, C. Rau, T. Weitkamp, A. Snigirev, and I. Snigireva, "Nanotomography Based on Hard X-Ray Microscopy with Refractive Lenses," *Appl. Phys. Lett.* **81**:1527–1529 (2002).
31. C. G. Schroer, B. Benner, T. F. Günzler, et al., "High Resolution Imaging and Lithography With Hard X-Rays Using Parabolic Compound Refractive Lenses," *Rev. Sci. Instrum.* **73**:1640 (2002).
32. B. Lengeler, C. G. Schroer, M. Kuhlmann, B. Benner, T. F. Günzler, O. Kurapova, F. Zontone, A. Snigirev, and I. Snigireva, "Beryllium Parabolic Refractive X-Ray Lenses," In A. S. Snigirev and D. C. Mancini, eds., "Design and Microfabrication of Novel X-Ray Optics II," *Proceedings of the SPIE*, **5539**:1–9 (2004).

33. O. Castelnaud, M. Drakopoulos, C. G. Schroer, I. Snigireva, A. Snigirev, and T. Ungar, "Dislocation Density Analysis in Single Grains of Steel by X-Ray Scanning Microdiffraction," *Nucl. Instrum. Methods A* **467–468**: 1245–1248 (2001).
34. S. Bohic, A. Simionovici, A. Snigirev, R. Ortega, G. Devès, D. Heymann, and C. G. Schroer, "Synchrotron Hard X-Ray Microprobe: Fluorescence Imaging of Single Cells," *Appl. Phys. Lett.* **78**:3544–3546 (2001).
35. A. S. Simionovici, M. Chukalina, C. Schroer, M. Drakopoulos, A. Snigirev, I. Snigireva, B. Lengeler, K. Janssens, and F. Adams, "High-Resolution X-Ray Fluorescence Microtomography of Homogeneous Samples," *IEEE Trans. Nucl. Sci.* **47**:2736–2740 (2000).
36. C. G. Schroer, J. Tümmeler, T. F. Günzler, B. Lengeler, W. H. Schröder, A. J. Kuhn, A. S. Simionovici, A. Snigirev, and I. Snigireva, "Fluorescence Microtomography: External Mapping of Elements Inside Biological Samples," In F. P. Doty, H. B. Barber, H. Roehrig, and E. J. Morton, eds., "Penetrating Radiation Systems and Applications II," *Proceedings of the SPIE*, **4142**:287–296 (2000).
37. C. G. Schroer, "Reconstructing X-Ray Fluorescence Microtomograms," *Appl. Phys. Lett.* **79**:1912–1914 (2001).
38. C. G. Schroer, M. Kuhlmann, T. F. Günzler, et al., "Mapping the Chemical States of an Element Inside a Sample Using Tomographic X-Ray Absorption Spectroscopy," *Appl. Phys. Lett.* **82**:3360–3362 (2003).
39. C. G. Schroer, M. Kuhlmann, S. V. Roth, R. Gehrke, N. Stribeck, A. Almandarez-Camarillo, and B. Lengeler, "Mapping the Local Nanostructure Inside a Specimen by Tomographic Small Angle X-Ray Scattering," *Appl. Phys. Lett.* **88**:164102 (2006).
40. H. Reichert, V. Honkimaki, A. Snigirev, S. Engemann, and H. Dosch, "A New X-ray Transmission-Reflection Scheme for the Study of Deeply Buried Interfaces Using High Energy Microbeams," *Physica B* **336**:46–55 (2003).
41. V. Nazmov, E. Resnikova, A. Last, J. Mohr, V. Saile, R. Simon, and M. DiMichiel, "X-Ray Lenses Fabricated by LIGA Technology," In J. -Y. Choi and S. Rah, eds., "Synchrotron Radiation Instrumentation: Ninth International Conference on Synchrotron Radiation Instrumentation," *AIP Conference Proceedings*, **879**:770–773 (2007).
42. M. Hanke, M. Dubsloff, M. Schmidbauer, T. Boeck, S. Schöder, M. Burghammer, C. Riekkel, J. Patommel, and C. G. Schroer, "Scanning X-Ray Diffraction with 200 nm Spatial Resolution," *Applied Physics Letters* **92**:193109 (2008).
43. C. G. Schroer, P. Boye, J. Feldkamp, J. Patommel, A. Schropp, A. Schwab, S. Stephan, M. Burghammer, S. Schoder, and C. Riekkel, "Coherent X-Ray Diffraction Imaging with Nanofocused Illumination," *Phys. Rev. Lett.* **101**:090801 (2008).

This page intentionally left blank.

DO NOT DUPLICATE

GRATINGS AND MONOCHROMATORS IN THE VUV AND SOFT X-RAY SPECTRAL REGION

Malcolm R. Howells

*Advanced Light Source
Lawrence Berkeley National Laboratory
Berkeley, California*

38.1 INTRODUCTION

Spectroscopy in the photon energy region from the visible to about 1 to 2 keV is generally done using reflection gratings. In the region above 40 eV, reasonable efficiency is only obtained at grazing angles and in this article we concentrate mainly on that case. Flat gratings were the first to be used and even today are still important. However, the advantages of spherical ones were recognized very early.¹ The first type of focusing grating to be analyzed theoretically was that formed by the intersection of a substrate surface with a set of parallel equispaced planes: the so-called “Rowland grating.” The theory of the spherical case was established first,¹⁻³ and was described comprehensively in the 1945 paper of Beutler.⁴ Treatments of toroidal⁵ and ellipsoidal⁶ gratings came later, and the field has been reviewed by Welford,⁷ Samson,⁸ Hunter,⁹ and Namioka.¹⁰

The major developments in the last three decades have been in the use of nonuniformly spaced grooves. The application of holography to spectroscopic gratings was first reported by Rudolph and Schmah^{11,12} and by Labeyrie and Flamand.¹³ Its unique opportunities for optical design were developed initially by Jobin-Yvon¹⁴ and by Namioka and coworkers.^{15,16} A different approach was followed by Harada¹⁷ and others, who developed the capability to produce gratings with variable-line spacing through the use of a computer-controlled ruling engine. The application of this class of gratings to spectroscopy has been developed still more recently, principally by Hettrick.¹⁸

In this chapter we will give a treatment of grating theory up to fourth order in the optical path, which is applicable to any substrate shape and any groove pattern that can be produced by holography or by ruling straight grooves with (possibly) variable spacing. The equivalent information is available up to sixth order at the website of the Center for X-Ray Optics at the Lawrence Berkeley National Laboratory.¹⁹

38.2 DIFFRACTION PROPERTIES

Notation and Sign Convention

We adopt the notation of Fig. 1 in which α and β have opposite signs if they are on opposite sides of the normal.

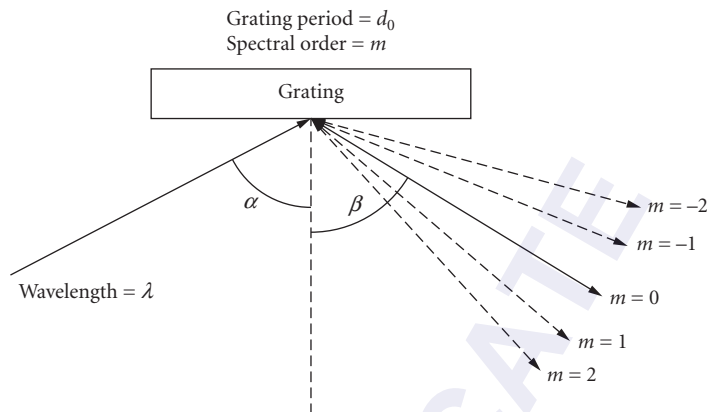


FIGURE 1 Grating equation notation.

Grating Equation

The grating equation may be written

$$m\lambda = d_0(\sin\alpha + \sin\beta) \quad (1)$$

The angles α and β are both arbitrary, so it is possible to impose various conditions relating them. If this is done, then for each λ , there will be a unique α and β . The following conditions are used:

1. *On-blaze condition:*

$$\alpha + \beta = 2\theta_B \quad (2)$$

where θ_B is the blaze angle (the angle of the sawtooth). The grating equation is then

$$m\lambda = 2d_0 \sin\theta_B \cos(\beta + \theta_B) \quad (3)$$

2. *Fixed in and out directions:*

$$\alpha - \beta = 2\theta \quad (4)$$

where 2θ is the (constant) included angle. The grating equation is then

$$m\lambda = 2d_0 \cos\theta \sin(\theta + \beta) \quad (5)$$

In this case, the wavelength scan ends when α or β reaches 90° , which occurs at the horizon wavelength $\lambda_H = 2d_0 \cos^2\theta$.

3. *Constant incidence angle:* Equation (1) gives β directly.
4. *Constant focal distance (of a plane grating):*

$$\frac{\cos\beta}{\cos\alpha} = \text{a constant } c_{ff} \quad (6)$$

leading to a grating equation

$$1 - \left(\frac{m\lambda}{d_0} - \sin\beta \right)^2 = \frac{\cos^2\beta}{c_{ff}^2} \quad (7)$$

Equations (3), (5), and (7) give β (and thence α) for any λ . Examples where the above α - β relationships may be used are as follows:

1. Kunz et al. plane-grating monochromator (PGM),²⁰ Hunter et al. double PGM,²¹ collimated-light SX700.²²
2. Toroidal-grating monochromators (TGMs),^{23,24} spherical-grating monochromators (SGMs, also known as the *Dragon* system),²⁵ Seya-Namioka,^{26,27} most aberration-reduced holographic SGMs,²⁸ and certain PGMs.^{18,29,30} The variable-angle SGM³¹ follows Eq. (4) approximately.
3. Spectrographs, *Grasshopper* monochromator.³²
4. SX700 PGM³³ and variants.^{22,34}

38.3 FOCUSING PROPERTIES³⁵

Calculation of the Path Function F

Following normal practice, we provide an analysis of the imaging properties of gratings by means of the path function F .¹⁶ For this purpose we use the notation of Fig. 2, in which the zeroth groove (of width d_0) passes through the grating pole O , while the n th groove passes through the variable point $P(\xi, w, l)$.

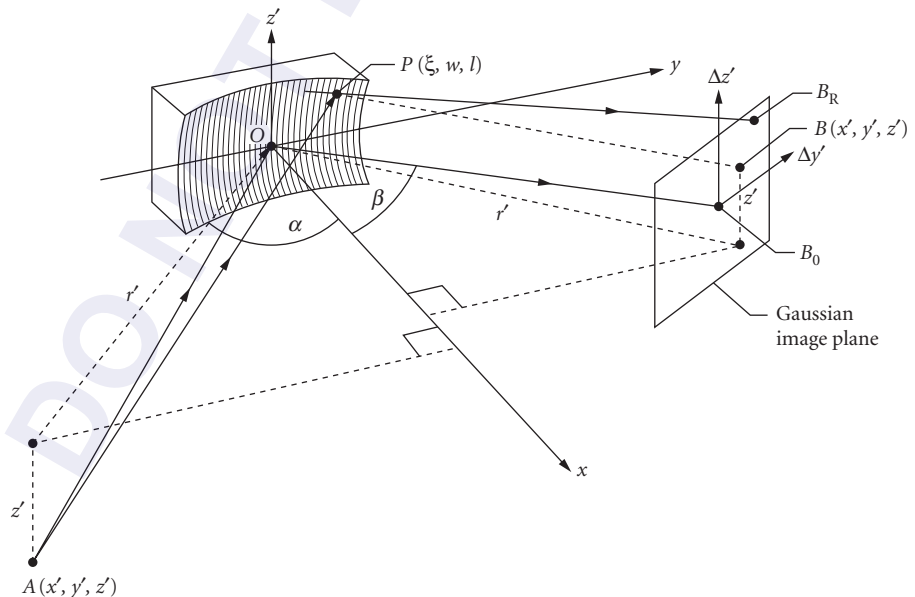


FIGURE 2 Focusing properties notation.

F is expressed as

$$F = \sum_{ijk} F_{ijk} w^i l^j$$

where

$$F_{ijk} = z^k C_{ijk}(\alpha, r) + z'^k C_{ijk}(\beta, r') + \frac{m\lambda}{d_0} f_{ijk} \quad (8)$$

and the f_{ijk} term, originating from the groove pattern, is given by one of the following expressions:

$$f_{ijk} = \begin{cases} 1 & \text{when } ijk = 100, 0 \text{ otherwise} & \text{Rowland} \\ \frac{d_0}{\lambda_0} [z^k C_{ijk}(\gamma, r_C) \pm z^k C_{ijk}(\delta, r_D)] & \text{holographic} \\ n_{ijk} & \text{varied line spacing} \end{cases} \quad (9)$$

The holographic groove pattern in Eq. (9) is assumed to be made using two coherent point sources C and D with cylindrical polar coordinates (r_C, γ, z_C) , (r_D, δ, z_D) relative to O . The lower (upper) sign refers to C and D , both real or both virtual (one real and one virtual), for which case the equiphas surfaces are confocal hyperboloids (ellipses) of revolution about CD . The grating with varied line spacing $d(w)$ is assumed to be ruled according to $d(w) = d_0(1 + v_1 w + v_2 w^2 + \dots)$. We consider all the gratings to be ruled on the general surface $x = \sum_{ij} a_{ij} w^i l^j$ and the a_{ij} coefficients³⁶ are given for the important substrate shapes in Tables 1 and 2.

TABLE 1 Ellipsoidal Mirror a_{ij} 's³⁶

$$\begin{aligned} a_{20} &= \frac{\cos \theta}{4} \left(\frac{1}{r} + \frac{1}{r'} \right) & a_{02} &= \frac{a_{20}}{\cos^2 \theta} & a_{22} &= \frac{a_{20}(2A^2 + C)}{2\cos^2 \theta} \\ a_{30} &= a_{20}A & a_{12} &= \frac{a_{20}A}{\cos^2 \theta} & a_{04} &= \frac{a_{20}C}{8\cos^2 \theta} \\ a_{40} &= \frac{a_{20}(4A^2 + C)}{4} \end{aligned}$$

Other a_{ij} 's with $i + j \leq 4$ are zero.

* r , r' and θ are the object distance, image distance, and incidence angle to the normal, respectively, and

$$A = \frac{\sin \theta}{2} \left(\frac{1}{r} - \frac{1}{r'} \right), \quad C = A^2 + \frac{1}{rr'}$$

The a_{ij} 's for spheres; circular, parabolic, or hyperbolic cylinders; paraboloids; and hyperboloids can also be obtained from Tables 1 and 2 by suitable choices of the input parameters r , r' , and θ .

TABLE 2 Toroidal Mirror a_{ij} 's³⁶

$$\begin{aligned} a_{20} &= \frac{1}{2R} & a_{02} &= \frac{1}{2\rho} & a_{04} &= \frac{1}{8R^3} \\ a_{30} &= \frac{1}{8R^3} & a_{22} &= \frac{1}{4\rho R^2} \end{aligned}$$

Other a_{ij} 's with $i + j \leq 4$ are zero

* R and ρ are the major and minor radii of the bicycle-tire toroid.

TABLE 3 Coefficients C_{ijk} of the Expansion of $F^{*,16}$

$C_{011} = -\frac{1}{r}$	$C_{020} = \frac{S}{2}$
$C_{022} = -\frac{S}{4r^2} - \frac{1}{2r^3}$	$C_{031} = \frac{S}{2r^2}$
$C_{040} = \frac{4a_{02}^2 - S^2}{8r} - a_{04} \cos \alpha$	$C_{100} = -\sin \alpha$
$C_{120} = \frac{S \sin \alpha}{2r} - a_{12} \cos \alpha$	$C_{200} = \frac{T}{2}$
$C_{300} = -a_{30} \cos \alpha + \frac{T \sin \alpha}{2r}$	$C_{102} = \frac{\sin \alpha}{2r^2}$
$C_{220} = -a_{22} \cos \alpha + \frac{1}{4r}(4a_{20}a_{02} - TS - 2a_{12} \sin 2\alpha) + \frac{S \sin^2 \alpha}{2r^2}$	$C_{202} = -\frac{T}{4r^2} + \frac{\sin^2 \alpha}{2r^3}$
$C_{400} = -a_{40} \cos \alpha + \frac{1}{8r}(4a_{20}^2 - T^2 - 4a_{30} \sin 2\alpha) + \frac{T \sin^2 \alpha}{2r^2}$	$C_{211} = \frac{T}{2r^2} - \frac{\sin^2 \alpha}{r^3}$

*The coefficients for which $i \leq 4, j \leq 4, k \leq 2, i + j + k \leq 4$, and $j + k = \text{even}$ are included in these tables.

TABLE 4 Coefficients n_{ijk} of the Expansion of F

$n_{ijk} = 0$ for $j, k \neq 0$	
$n_{100} = 1$	$n_{300} = \frac{v_1^2 - v_2}{3}$
$n_{200} = \frac{-v_1}{2}$	$n_{400} = \frac{-v_1^3 + 2v_1v_2 - v_3}{4}$

The coefficient F_{ijk} is related to the strength of the i, j, k aberration of the wavefront diffracted by the grating. The coefficients C_{ijk} and n_{ijk} are given in Tables 3 and 4 in which the following notation is used:

$$T = T(r, \alpha) = \frac{\cos^2 \alpha}{r} - 2a_{20} \cos \alpha, \quad S = S(r, \alpha) = \frac{1}{r} - 2a_{02} \cos \alpha \quad (10)$$

Determination of the Gaussian Image Point

By definition the principal ray AOB_0 arrives at the Gaussian image point $[B_0(r'_0, \beta'_0, z'_0)]$ in Fig. 2. Its direction is given by Fermat's principle which implies $(\partial F / \partial w)_{w=0, l=0} = 0, (\partial F / \partial l)_{w=0, l=0} = 0$, from which

$$\frac{m\lambda}{d_0} = \sin \alpha + \sin \beta_0 \quad \frac{z}{r} + \frac{z'_0}{r'_0} = 0 \quad (11)$$

which are the grating equation and the law of magnification in the vertical direction. The tangential focal distance r'_0 is obtained by setting the focusing term F_{200} equal to zero and is given by

$$T(r, \alpha) + T(r'_0, \beta_0) = \begin{cases} 0 & \text{Rowland} \\ -\frac{m\lambda}{\lambda_0} [T(r_C, \gamma) \pm T(r_D, \delta)] & \text{holographic} \\ \frac{v_1 m \lambda}{d_0} & \text{varied line spacing} \end{cases} \quad (12)$$

Equations (11) and (12) determine the Gaussian image point B_0 and, in combination with the sagittal focusing condition ($F_{020} = 0$), describe the focusing properties of grating systems under the paraxial approximation. For a Rowland spherical grating the focusing condition [Eq. (12)] is

$$\left(\frac{\cos^2 \alpha}{r} - \frac{\cos \alpha}{R} \right) + \left(\frac{\cos^2 \beta}{r'_0} - \frac{\cos \beta}{R} \right) = 0 \quad (13)$$

which has the following important special cases:

1. A plane grating ($R = \infty$) implying $r'_0 = -r \cos^2 \beta / \cos^2 \alpha = -r/c_{ff}^2$, so that the focal distance and magnification are fixed if c_{ff} is held constant.³⁷
2. Object and image on the Rowland circle; $r = R \cos \alpha$, $r'_0 = R \cos \beta$, and $M = -1$.
3. $\beta = 0$ (Wadsworth condition).

The tangential focal distances of TGMs and SGMs with or without moving slits are also determined by Eq. (13).

In an aberrated system, the outgoing ray will arrive at the Gaussian image plane at a point B_R displaced from the Gaussian image point B_0 by the ray aberrations $\Delta y'$ and $\Delta z'$ (Fig. 2). The latter are given by³⁸⁻⁴⁰

$$\Delta y' = \frac{r'_0}{\cos \beta_0} \frac{\partial F}{\partial w} \quad \Delta z' = r'_0 \frac{\partial F}{\partial l} \quad (14)$$

where F is to be evaluated for $A = (r, \alpha, z)$, $B = (r'_0 = \beta_0, z'_0)$. By means of the series expansion of F , these equations allow the ray aberrations to be calculated separately for each aberration type, as follows:

$$\Delta y'_{ijk} = \frac{r'_0}{\cos \beta_0} F_{ijk} i w^{i-1} l^j \quad \Delta z'_{ijk} = r'_0 F_{ijk} w^i j l^{j-1} \quad (15)$$

Moreover, provided the aberrations are not too large, they are additive, so that they may either reinforce or cancel.

38.4 DISPERSION PROPERTIES

Angular Dispersion

$$\left(\frac{\partial \lambda}{\partial \beta} \right)_\alpha = \frac{d \cos \beta}{m} \quad (16)$$

Reciprocal Linear Dispersion

$$\left(\frac{\partial \lambda}{\partial (\Delta y')} \right)_{\alpha} = \frac{d \cos \beta}{mr'} \equiv \frac{10^{-3} d [\text{\AA}] \cos \beta}{mr' [\text{m}]} \text{\AA}/\text{mm} \quad (17)$$

Magnification (M)

$$M(\lambda) = -\frac{\cos \alpha}{\cos \beta} \frac{r'}{r} \quad (18)$$

Phase-Space Acceptance (ε)

$$\varepsilon = N \Delta \lambda_{S_1} = N \Delta \lambda_{S_2} \quad (\text{assuming } S_2 = M S_1) \quad (19)$$

where N is the number of participating grooves.

38.5 RESOLUTION PROPERTIES

The following are the main contributions to the width of the instrumental line spread function (an estimate of the total width is the vector sum):

1. Entrance slit (width S_1):

$$\Delta \lambda_{S_1} = \frac{S_1 d \cos \alpha}{mr} \quad (20)$$

2. Exit slit (width S_2):

$$\Delta \lambda_{S_2} = \frac{S_2 d \cos \beta}{mr'} \quad (21)$$

3. Aberrations (of a perfectly made grating):

$$\Delta \lambda_A = \frac{\Delta y' d \cos \beta}{mr'} = \frac{d}{m} \left(\frac{\partial F}{\partial w} \right) \quad (22)$$

4. Slope error $\Delta \phi$ (of an imperfectly made grating):

$$\Delta \lambda_{SE} = \frac{d(\cos \alpha + \cos \beta) \Delta \phi}{m} \quad (23)$$

Note that, provided the grating is large enough, diffraction at the entrance slit always guarantees a coherent illumination of enough grooves to achieve the slit-width limited resolution. In such cases, a diffraction contribution to the width need not be added to those listed.

38.6 EFFICIENCY

The most accurate way to calculate grating efficiencies is by the full electromagnetic theory for which code is available from Neviere.^{41,42} However, approximate scalar-theory calculations are often useful and, in particular, provide a way to choose the groove depth (h) of a laminar grating. According to Bennett,⁴³ the best value of the groove-width-to-period ratio (r) is the one for which the area of the usefully illuminated groove bottom is equal to that of the top. The scalar theory efficiency of a laminar grating with $r = 0.5$ is given by Franks et al.⁴⁴ as the following:

$$E_0 = \frac{R}{4} \left[1 + 2(1 - P) \cos \left(\frac{4\pi h \cos \alpha}{\lambda} \right) + (1 - P)^2 \right]$$

$$E_m = \begin{cases} R[1 - 2 \cos Q^+ \cos(Q^- + \delta) + \cos^2 Q^+] / m^2 \pi^2 & m = \text{odd} \\ R \cos^2 Q^+ / m^2 \pi^2 & m = \text{even} \end{cases} \quad (24)$$

where

$$P = \frac{4h \tan \alpha}{d_0} \quad Q^\pm = \frac{m\pi h}{d_0} (\tan \alpha \pm \tan \beta) \quad \delta = \frac{2\pi h}{\lambda} (\cos \alpha + \cos \beta)$$

and R is effective reflectance given by $R = \sqrt{R(\alpha)R(\beta)}$ where $R(\alpha)$ and $R(\beta)$ are the intensity reflectances at α and β , respectively.

38.7 REFERENCES

1. H. A. Rowland, "On Concave Gratings for Optical Purposes," *Phil. Mag.* **16**(5th ser.):197–210 (1883).
2. J. E. Mack, J. R. Stehn, and B. Edlen, "On the Concave Grating Spectrograph, Especially at Large Angles of Incidence," *J. Opt. Soc. Am.* **22**:245–264 (1932).
3. H. A. Rowland, "Preliminary Notice of the Results Accomplished in the Manufacture and Theory of Gratings for Optical Purposes," *Phil. Mag.* **13**(supp.) (5th ser.):469–474 (1882).
4. H. G. Beutler, "The Theory of the Concave Grating," *J. Opt. Soc. Am.* **35**:311–350 (1945).
5. H. Haber, "The Torus Grating," *J. Opt. Soc. Am.* **40**:153–165 (1950).
6. T. Namioka, "Theory of the Ellipsoidal Concave Grating: I," *J. Opt. Soc. Am.* **51**:4–12 (1961).
7. W. Welford, "Aberration Theory of Gratings and Grating Mountings," in E. Wolf (ed.), *Progress in Optics*, vol. 4, North-Holland, Amsterdam, pp. 243–282, 1965.
8. J. A. R. Samson, *Techniques of Vacuum Ultraviolet Spectroscopy*, John Wiley & Sons, New York, 1967.
9. W. R. Hunter, "Diffraction Gratings and Mountings for the Vacuum Ultraviolet Spectral Region," in *Spectrometric Techniques*, G. A. Vanasse (ed.), vol. IV, Academic Press, Orlando, pp. 63–180, 1985.
10. T. Namioka and K. Ito, "Modern Developments in VUV Spectroscopic Instrumentation," *Physica Scripta* **37**:673–681 (1988).
11. D. Rudolph and G. Schmahl, "Verfahren zur Herstellung von Röntgenlinsen und Beugungsgittern," *Umsch. Wiss. Tech.* **67**:225 (1967).
12. D. Rudolph and G. Schmahl, "Holographic Gratings," in *Progress in Optics*, E. Wolf (ed.), vol. 14, North-Holland, Amsterdam, pp. 196–244, 1977.
13. A. Laberie and J. Flamand, "Spectrographic Performance of Holographically Made Diffraction Grating," *Opt. Comm.* **1**:5–8 (1969).

14. G. Pieuchard and J. Flamand, "Concave Holographic Gratings for Spectrographic Applications," Final report on NASA contract number NASW-2146, GSFC 283-56,777, Jobin Yvon, Longjumeau, France, 1972.
15. T. Namioka, H. Noda, and M. Seya, "Possibility of Using the Holographic Concave Grating in Vacuum Monochromators," *Sci. Light* **22**:77-99 (1973).
16. H. Noda, T. Namioka, and M. Seya, "Geometrical Theory of the Grating," *J. Opt. Soc. Am.* **64**:1031-1036 (1974).
17. T. Harada and T. Kita, "Mechanically Ruled Aberration-Corrected Concave Gratings," *Appl. Opt.* **19**:3987-3993 (1980).
18. M. C. Hettrick, "Aberration of Varied Line-Space Grazing Incidence Gratings," *Appl. Opt.* **23**:3221-3235 (1984).
19. <http://www-cxro.lbl.gov/>, 1999.
20. C. Kunz, R. Haensel, and B. Sonntag, "Grazing Incidence Vacuum Ultraviolet Monochromator with Fixed Exit Slit for Use with Distant Sources," *J. Opt. Soc. Am.* **58**:1415 (1968).
21. W. R. Hunter, R. T. Williams, J. C. Rife, J. P. Kirkland, and M. N. Kaber, "A Grating/Crystal Monochromator for the Spectral Range 5 eV to 5 keV," *Nucl. Instr. Meth.* **195**:141-154 (1982).
22. R. Follath and F. Senf, "New Plane-Grating Monochromators for Third Generation Synchrotron Radiation Light Sources," *Nucl. Instrum. Meth.* **A390**:388-394 (1997).
23. D. Lepere, "Monochromators with Single Axis Rotation and Holographic Gratings on Toroidal Blanks for the Vacuum Ultraviolet," *Nouvelle Revue Optique* **6**:173 (1975).
24. R. P. Madden and D. L. Ederer, "Stigmatic Grazing Incidence Monochromator for Synchrotrons (abstract only)," *J. Opt. Soc. Am.* **62**:722 (1972).
25. C. T. Chen, "Concept and Design Procedure for Cylindrical Element Monochromators for Synchrotron Radiation," *Nucl. Instr. Meth.* **A256**:595-604 (1987).
26. T. Namioka, "Construction of a Grating Spectrometer," *Sci. Light* **3**:15-24 (1954).
27. M. Seya, "A New Mounting of Concave Grating Suitable for a Spectrometer," *Sci. Light* **2**:8-17 (1952).
28. T. Namioka, M. Seya, and H. Noda, "Design and Performance of Holographic Concave Gratings," *Jap. J. Appl. Phys.* **15**:1181-1197 (1976).
29. W. Eberhardt, G. Kalkoffen, and C. Kunz, "Grazing Incidence Monochromator FLIPPER," *Nucl. Instr. Meth.* **152**:81-4 (1978).
30. K. Miyake, P. R. Kato, and H. Yamashita, "A New Mounting of Soft X-Ray Monochromator for Synchrotron Orbital Radiation," *Sci. Light* **18**:39-56 (1969).
31. H. A. Padmore, "Optimization of Soft X-Ray Monochromators," *Rev. Sci. Instrum.* **60**:1608-1616 (1989).
32. F. C. Brown, R. Z. Bachrach, and N. Lien, "The SSRL Grazing Incidence Monochromator: Design Considerations and Operating Experience," *Nucl. Instrum. Meth.* **152**:73-80 (1978).
33. H. Petersen and H. Baumgartel, "BESSY SX/700: A Monochromator System Covering the Spectral Range 3 eV-700 eV," *Nucl. Instrum. Meth.* **172**:191-193 (1980).
34. W. Jark, "Soft X-Ray Monochromator Configurations for the ELETTRA Undulators: A Stigmatic SX700," *Rev. Sci. Instrum.* **63**:1241-1246 (1992).
35. H. A. Padmore, M. R. Howells, and W. R. McKinney, "Grazing Incidence Monochromators for Third-Generation Synchrotron Radiation Light Sources," in J. A. R. Samson and D. L. Ederer (eds.), *Vacuum Ultraviolet Spectroscopy*, vol. 31, Academic Press, San Diego, pp. 21-54, 1998.
36. S. Y. Rah, S. C. Irick, and M. R. Howells, "New Schemes in the Adjustment of Bendable Elliptical Mirrors Using a Long-Trace Profiler," in P. Z. Takacs and T. W. Tonnessen (eds.), *Materials Manufacturing and Measurement for Synchrotron-Radiation Mirrors*, Proc. SPIE, vol. 3152, SPIE, Bellingham, WA, 1997.
37. H. Petersen, "The Plane Grating and Elliptical Mirror: A New Optical Configuration for Monochromators," *Opt. Comm.* **40**:402-406 (1982).
38. M. Born and E. Wolf, *Principles of Optics*, Pergamon, Oxford, 1980.
39. T. Namioka and M. Koike, "Analytical Representation of Spot Diagrams and Its Application to the Design of Monochromators," *Nucl. Instrum. Meth.* **A319**:219-227 (1992).
40. W. T. Welford, *Aberrations of the Symmetrical Optical System*, Academic Press, London, 1974.
41. M. Nevriere, P. Vincent, and D. Maystre, "X-Ray Efficiencies of Gratings," *Appl. Opt.* **17**:843-845 (1978). (Nevriere can be reached at michel.nevriere@fresnel.fr.)

42. R. Petit (ed.), *Electromagnetic Theory of Gratings (Topics in Current Physics, Vol. 22)*, Springer Verlag, Berlin, 1980.
43. J. M. Bennett, "Laminar X-Ray Gratings," Ph.D. Thesis, London University, London, 1971.
44. A. Franks, K. Lindsay, J. M. Bennett, R. J. Speer, D. Turner, and D. J. Hunt, "The Theory, Manufacture, Structure and Performance of NPL X-Ray Gratings," *Phil. Trans. Roy. Soc.* **A277**:503–543 (1975).

DO NOT DUPLICATE

CRYSTAL MONOCHROMATORS AND BENT CRYSTALS

Peter Siddons

National Synchrotron Light Source
Brookhaven National Laboratory
Upton, New York

39.1 CRYSTAL MONOCHROMATORS

For x-ray energies higher than 2 keV or so, gratings become extremely inefficient, and it becomes necessary to utilize the periodicity naturally occurring in a crystal to provide the dispersion. Since the periodicity in a crystal is 3-dimensional, the normal single grating equation must be replaced by the three grating equations, one for each dimension, called the Laue equations,¹ as follows:

$$\begin{aligned} \mathbf{a}_1 \cdot (\mathbf{k}_{H_1 H_2 H_3} - \mathbf{k}_0) &= H_1 \\ \mathbf{a}_2 \cdot (\mathbf{k}_{H_1 H_2 H_3} - \mathbf{k}_0) &= H_2 \\ \mathbf{a}_3 \cdot (\mathbf{k}_{H_1 H_2 H_3} - \mathbf{k}_0) &= H_3 \end{aligned} \quad (1)$$

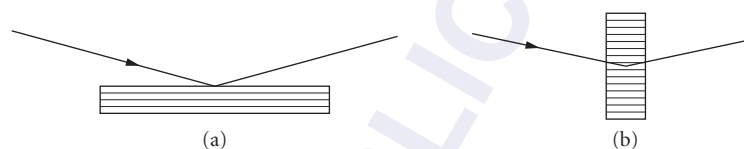
where the \mathbf{a} 's are the repeat vectors in the three dimensions, the \mathbf{k} 's are the wave vectors for the incident and scattered beams, and the H 's are integers denoting the diffraction order in the three dimensions. All of them must be simultaneously satisfied in order to have an interference maximum (commonly called a *Bragg reflection*). One can combine these equations into the well-known Bragg's law² for one component of the crystalline periodicity (usually referred to as a set of *Bragg planes*),

$$n\lambda = 2d \sin \theta \quad (2)$$

where n is an integer indicating the order of diffraction from planes of spacing d , and θ is the angle between the incident beam and the Bragg planes. This equation is the basis for using crystals as x-ray monochromators. By choosing one such set of Bragg planes and setting the crystal so that the incident x rays fall on these planes, the wavelength of the light reflected depends on the angle of incidence of the light. Table 1 shows some commonly used crystals and the spacings of some of their Bragg planes. The most common arrangement is the symmetric Bragg case (Fig. 1a), in which the useful surface of the crystal is machined so that it is parallel to the Bragg planes in use. Under these conditions, the incident and reflected angles are equal. The literature on crystal diffraction differs from that on grating instruments in that the incidence angle for x rays is called the *glancing angle* for

TABLE 1 Some Common Monochromator Crystals, Selected d-Spacings, and Reflection Widths at 1 Å (12.4 keV)

Crystal	Reflection	d-Spacing (nm)	Refl. Width (μrad)	Energy Resolution ($\Delta E/E$)
Silicon	(1 1 1)	0.31355	22.3	1.36×10^{-4}
	(2 2 0)	0.19201	15.8	5.37×10^{-5}
Germanium	(1 1 1)	0.32664	50.1	3.1×10^{-4}
	(2 2 0)	0.20002	37.4	1.37×10^{-4}
Diamond	(1 1 1)	0.20589	15.3	5.8×10^{-5}
	(4 0 0)	0.089153	5.2	7.4×10^{-6}
Graphite	(0 0 0 . 2)	0.3354	Sample-dependent	Sample-dependent

**FIGURE 1** The two most usual x-ray diffraction geometries used for monochromator applications: (a) the symmetric Bragg case and (b) the symmetric Laue case.

gratings. As the x-ray wavelength gets shorter and the angles get smaller, it can be difficult to obtain large enough crystals of good quality. In such cases it is possible to employ the Laue case (Fig. 1*b*), in which the surface is cut perpendicular to the Bragg planes and the x rays are reflected through the bulk of the crystal plate. Of course, the wavelength should be short enough or the crystal thin enough so that the x rays are not absorbed by the monochromator. This is true, for example, in silicon crystals around 1 mm thick above an x-ray energy of around 30 keV ($\lambda = 0.04$ nm).

The detailed calculation of the response of a crystal to x rays depends on the degree of crystalline perfection of the material in use, as well as its chemical composition. Two main theoretical treatments are commonly used: the kinematical theory of diffraction and the dynamical theory. The kinematical theory assumes single-scattering of the x rays by the crystal, and is appropriate for crystals with a high concentration of defects. Such crystals are called *mosaic crystals*, following C. G. Darwin.³ The dynamical theory, in contrast, explicitly treats the multiple scattering that arises in highly perfect crystals and is commonly used for the semiconductor monochromator materials that can be grown to a high degree of perfection. Of course, many crystals fall between these two idealized pictures and there exist approximations to both theories to account for some of their failures. We will not describe these theories in detail here, but will refer the reader to texts on the subject⁴⁻⁶ and content ourselves with providing some of the key formulas that result from them.

Both theories attempt to describe the variation in reflectivity of a given crystal as a function of the incidence angle of the x-ray beam near a Bragg reflection. The neglect or inclusion of multiple scattering changes the result quite dramatically. In the kinematical case, the width of the reflectivity profile is inversely related to the size of the coherently diffracting volume, and for an infinite perfect crystal it is a delta function. The integrated reflectivity increases linearly with the illuminated crystal volume, and is assumed to be small. In the dynamical case, the x-ray beam diffracted by one part of the crystal is exactly oriented to be diffracted by the same Bragg planes, but in the opposite sense. Thus, there coexist two waves in the crystal, one with its wave vector along the incident beam direction and the other with its vector along the diffracted beam direction. These waves are coherent and can interfere. It is these interferences that give rise to all the interesting phenomena that arise from this theory. The integrated reflectivity in this case initially increases with volume, as in the kinematical case, but eventually saturates to a constant value that depends on which of the geometries in Fig. 1 is taken.

For the kinematical theory, the reflectivity curve is approximated by⁷

$$r(\Delta) = \frac{a}{[1 + a + \sqrt{1 + 2a} \cdot \cot h(b\sqrt{1 + 2a})]} \quad (3)$$

where

$$a = \frac{w(\Delta)Q}{\mu}$$

$$b = \frac{\mu t}{\sin \theta_B}$$

$$Q = \left(\frac{e^2}{mc^2V} \right)^2 \cdot \frac{|F_H|^2 \lambda^3}{\sin 2\theta_B}$$

in which e is the electronic charge and m its mass, c is the velocity of light, θ_B is the Bragg angle for the planes whose structure factor is F_H at wavelength λ . $w(\Delta)$ represents the angular distribution of the mosaic blocks, which depends on the details of the sample preparation and/or growth history. The reflection width will primarily reflect the width of this parameter $w(\Delta)$, but often the curve will be further broadened by extinction.

The equivalent equation for the dynamical theory is sensitive to the exact geometry under consideration and so will not be given here. The crystal becomes birefringent near the Bragg angle, and this causes some interesting interference effects that are worthy of study in their own right. Reference 8 is a review of the theory with a physical approach that is very readable. However, the main results for the purposes of this section are summarized in Fig. 2; the key points are the following:

1. In the Bragg case (Fig. 1a) the peak reflectivity reaches nearly unity over a small range of angles near the Bragg angle.

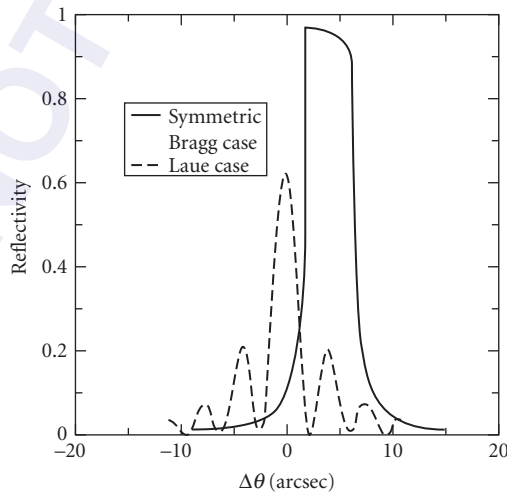


FIGURE 2 The perfect crystal reflectivity curves for the (111) reflection at 1 Å wavelength. The solid curve is for the thick Bragg case, and the dashed curve is for the Laue case for a 20 μm thick crystal

2. The range of angles over which this occurs is given by

$$\Omega = 2R\lambda^2 |C| \frac{\sqrt{|\gamma| F_h F_{-h}}}{\pi V \sin 2\theta} \quad (4)$$

where R is the classical electron radius, 2.81794×10^{-15} m, λ is the x-ray wavelength, C is the polarization factor (1 for σ -polarized light and $\cos 2\theta$ for π -polarized light). γ is the asymmetry parameter (1 in the symmetric case), the F s are the structure factors for reflections (hkl) and $(-h-k-l)$, V is the unit cell volume, and θ is the Bragg angle.

3. The effect of absorption in the Bragg case is to reduce the peak reflectivity and to make the reflectivity curve asymmetric.
4. In the Laue case, the reflectivity is oscillatory, with an average value of 1/2, and the effect of absorption (including increasing thickness) is to damp out the oscillations and reduce the peak reflectivity.

Even small distortions can greatly perturb the behavior of perfect crystals, and there is still no general theory that can handle arbitrary distortions in the dynamical regime. There are approximate treatments that can handle some special cases of interest.^{4,5}

In general, the performance of a given monochromator system is determined by the range of Bragg angles experienced by the incident x-ray beam. This can be determined simply by the incident beam collimation or by the crystal reflectivity profile, or by a combination of both factors. From Bragg's law we have

$$\frac{\delta E}{E} = \frac{\delta \lambda}{\lambda} = \cot \theta \delta \theta + \delta \tau \quad (5)$$

where δE and E are the passband and center energy of the beam, and $\delta \tau$ is the intrinsic energy resolution of the crystal reflection as given in Table 1—essentially the dynamical reflection width as given previously, expressed in terms of energy. This implies that, even for a perfectly collimated incident beam, there is a limit to the resolution achievable, which depends on the monochromator material and the diffraction order chosen.

The classic Bragg reflection monochromator uses a single piece of crystal set to the correct angle to reflect the energy of interest. This is in fact a very inconvenient instrument, since its output beam direction changes with energy. For perfect-crystal devices, the reflectivity is sufficiently high that one can afford to use two of them in tandem, deviating in opposite senses in order to bring the useful beam back into the forward direction, independent of energy (Fig. 3). Such two-crystal monochromators have become the standard design, particularly for use at accelerator-based (synchrotron) radiation sources (as discussed in Chap. 55). Since these sources naturally generate an energy continuum, a monochromator is a common requirement for a beamline at such a facility. There is a wealth of literature arising from such applications.⁹

A particularly convenient form of this geometry was invented by Bonse and Hart,¹⁰ in which the two reflecting surfaces are machined from a single-crystal block of material (in their case germanium, but more often silicon in recent years).

For the highest resolution requirements, it is possible to take two crystals that deviate in the same direction, the so-called ++ geometry. In this case the first crystal acts as a collimator, generating a fan of beams with each angular component having a particular energy. The second one selects which angular (and hence energy) component of this fan to transmit. In this arrangement the deviation is doubled over the single-crystal device, and so the most successful arrangement for this so-called dispersive arrangement is to use two monolithic double-reflection devices like that in Fig. 2, with the ++ deviation taking place between the second and third reflections (Fig. 4).

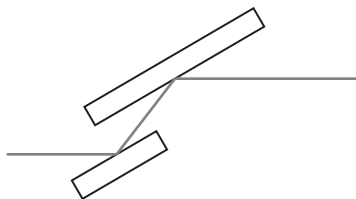


FIGURE 3 The \pm double-crystal x-ray monochromator.

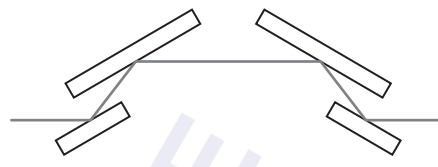


FIGURE 4 Two monolithic double reflectors arranged in a high-resolution configuration.

39.2 BENT CRYSTALS

There are two reasons to consider applying a uniform curvature to a crystal plate for a monochromator system. One is to provide some kind of beam concentration or focussing, and the other is to improve the energy resolution in circumstances where a divergent incident beam is unavoidable. The most common geometry is the Bragg-case one based on the Rowland circle principle, modified to account for the 3-dimensional nature of the crystal periodicity. This principle relies on the well-known property of a circle that an arc segment subtends a constant angle for any point on the circle. For a grating this is straightforwardly applied to a focussing spectrometer, but for crystals one has the added complication that the incidence angle for the local Bragg plane must also be constant. The result was shown by Johansson¹¹ (Fig. 5) to require the radius of curvature to be twice the radius of the crystal surface. This can be achieved by a combination of bending and machining the crystal. For applications in which the required optical aperture is small, the aberrations introduced by omitting the machining operation are small and quite acceptable.¹² Since this is a rather difficult operation, it is attractive to avoid it if possible.

Although Fig. 5 shows the symmetrical setting, it is possible to place the crystal anywhere on the circle and achieve a (de)magnification other than unity. In this case the surface must be cut at an angle to the Bragg planes to maintain the geometry. The Laue case can also be used in a similar arrangement, but the image becomes a virtual image (or source). For very short wavelengths (e.g., in a gamma-ray spectrometer) the Laue case can be preferable since the size of the crystal needed for a given optical aperture is much reduced. In both Laue and Bragg cases, there are changes in the reflection properties on bending. Depending on the source and its collimation geometry, and on the asymmetry angle of the crystal cut, the bending can improve or degrade the monochromator resolution. Each case must

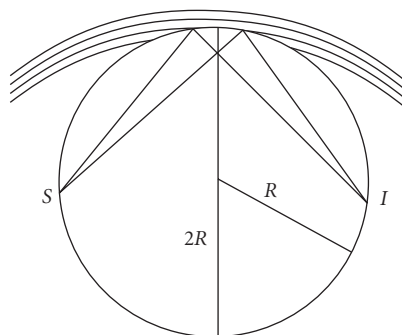


FIGURE 5 The Johansson bent/ground focussing monochromator.

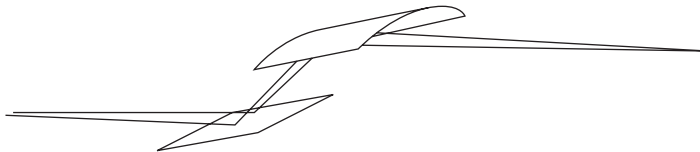


FIGURE 6 The usual arrangement of a sagittally focussing monochromator, with the bent element as the second one of a two-crystal monochromator.

be considered in detail. Again, within the scope of this section we will content ourselves with some generalities:

1. If the angular aperture of the crystal as seen from the source location is large compared to the reflectivity profile of the crystal, then the Rowland circle geometry will improve things for the Bragg case and the symmetric Laue case.
2. When the absorption length of the incident radiation becomes large compared to the extinction distance, then the x rays can travel deep into the bent crystal even in the Bragg case, and the deformation means that the incidence angle changes with depth, leading to a broadening of the bandpass.
3. If the bending radius becomes very small, such that the Bragg angle changes more than the perfect-crystal reflection width within the extinction depth, then the peak reflectivity will fall and the reflection width will increase, and consequently the resolution will deteriorate.
4. In the asymmetric Laue case, the reflectivity profile width for perfect crystals also depends on the curvature,⁶ and so for strongly collimated beams such as synchrotron radiation sources, the bending may well degrade the resolution at the same time as it increases the intensity.

Arrangements of multiple consecutive curved crystals are unusual, but have found application in high-throughput synchrotron radiation (SR) monochromators, where the two-crystal device in Fig. 3 is modified by placing the first crystal on the Rowland circle and adjusting the convex curvature of the second to maximize its transmission of the convergent beam from the first.¹³ There are also examples of combinations of curved Laue and curved Bragg reflectors.¹⁴

Another geometry in which the technique of bending a monochromator crystal can be used is in the so-called sagittal-focussing geometry.¹⁵ This geometry, as its name indicates, has the bending radius perpendicular to that in the Rowland-circle geometry, and provides a focussing effect in the plane perpendicular to the diffraction plane. In the diffraction plane the crystal behaves essentially as though it were flat (Fig. 6). Antielastic bending of the bent crystal can greatly reduce its efficiency, since this curvature is in the meridional plane so that parts of the crystal move out of the Bragg-reflection range. Attempts to counteract this by using stiffening ribs machined into the plate¹⁶ or by making use of the elastic anisotropy of single crystals¹⁷ have been made, with some success.

39.3 REFERENCES

1. M. von Laue, *Münchener Sitzungsberichte* 363 (1912), and *Ann. der Phys.* **41**:989 (1913).
2. W. H. Bragg and W. L. Bragg, *Proc. Roy. Soc. London* **88**:428 (1913), and **89**:246 (1913).
3. C. G. Darwin, *Phil. Mag.* **27**:315, 657 (1914), and **43**:800 (1922).
4. D. Taupin, *Bull. Soc. Franc. Miner. Crist.* **87**:469–511 (1964); S. Takagi, *Acta Cryst.* **12**:1311 (1962).
5. Penning and Polder, *Philips Res. Rep.* **16**:419 (1961).
6. R. Caciuffo, S. Melone, F. Rustichelli, and A. Boeuf, *Phys. Rep.* **152**:1 (1987); P. Suortti and W. C. Thomlinson, *Nuclear Instrum. and Meth.* **A269**:639–648 (1988).

7. A. K. Freund, A. Munkholm, and S. Brennan, *Proc. SPIE* **2856**:68–79 (1996).
8. B. W. Batterman and H. Cole, *Rev. Mod. Phys.* **36**:681–717 (1964). (For a fuller treatment, see W. H. Zachariasen, *Theory of X-Ray Diffraction in Crystals*, Dover, New York, 1967.)
9. See for example, the *Handbook of Synchrotron Radiation* or browse the *Proceedings of the International Conference on Synchrotron Radiation Instrumentation*, published in *Nucl. Instr. and Meth.*, *Rev. Sci. Instrum.*, and *J. Synchr. Rad.* at various times.
10. U. Bonse and M. Hart, *Appl. Phys. Lett.* **7**:238–240 (1965).
11. T. Johansson, *Z. Phys.* **82**:507 (1933).
12. H. H. Johann, *Z. Phys.* **69**:185 (1931).
13. M. J. Van Der Hoek, W. Berne, and P. Van Zuylen, *Nucl. Instrum. and Meth.* **A246**:190–193 (1986).
14. U. Lienert, C. Schulze, V. Honkimäki, Th. Tschentschor, S. Garbe, O. Hignette, A. Horsewell, M. Lingham, H. F. Poulsen, W. B. Thomsen, and E. Ziegler, *J. Synchrot. Rad.* **5**:226–231 (1998).
15. L. Von Hamos, *J. Sci. Instr.* **15**:87 (1938).
16. C. J. Sparks, Jr., B. S. Borie, and J. B. Hastings, *Nucl. Instrum. and Meth.* **172**:237 (1980).
17. V. I. Kushnir, J. P. Quintana, and P. Georgopoulos, *Nucl. Instrum. and Meth.* **A328**:588 (1983).

This page intentionally left blank.

DO NOT DUPLICATE

Alan Michette

King's College London
United Kingdom

40.1 INTRODUCTION

Some radiation incident on a linear transmission grating passes straight through (the *zero order*), some is diffracted to one side of the zero order (the *positive orders*), and some is diffracted to the other side (the *negative orders*). In the first order, the diffraction angle is $\beta = \sin^{-1}(\lambda/d) \approx \lambda/d$ in the small angle approximation; d is the grating period. Thus, for smaller periods, radiation is diffracted through larger angles. A *circular* grating with a constant period would therefore form an axial line focus of a point source (Fig. 1a), and the distance from a radial point r on the grating to a point on the axis is $z = r/\tan\beta \approx rd/\lambda$.

If the period is made to decrease as the radius increases (Fig. 1b) then the distance z can be made constant. The grating then acts as a lens in that radiation from a point source is brought to an axial focus (Fig. 1c). The positive diffraction orders are now defined as being on the opposite side to the source, with the negative orders on the same side.

This is the basis of zone plates, the focusing properties of which depend on

- The relationship between d and r
- The number of zones (For x-ray zone plates the usual convention is that the area between successive boundaries is a zone. Strictly speaking, and in keeping with the terminology used for diffraction gratings, this area should be called a *half-period zone* but zone is usually used.)
- The zone heights and profiles

40.2 GEOMETRY OF A ZONE PLATE

Referring to Fig. 1c, radiation from an object point A is brought to a focus, via the zone plate, to an image point B. To obtain constructive interference at B the optical path difference between successive zone boundaries must be $\pm m\lambda/2$, where m is the diffraction order. Thus, for the first order,

$$a_n + b_n = z_a + z_b + \frac{n\lambda}{2} \quad (1)$$

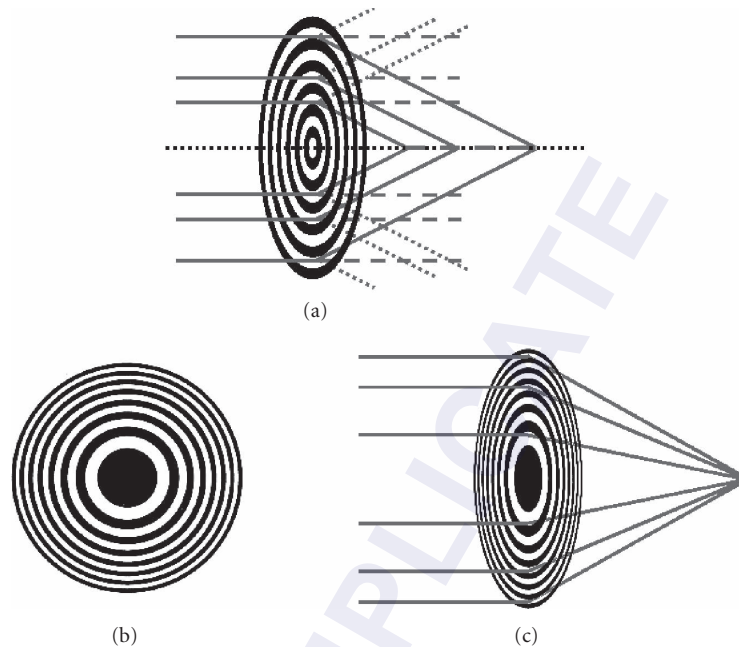


FIGURE 1 Diffraction by (a) a circular grating of constant period and (b, c) a zone plate. (See also color insert.)

where n is the zone number, counting outward from the center, and Δ is the optical path difference introduced by the central zone of radius r_0 . For a distant source ($a_n, z_a \rightarrow \infty$ with $a_n - z_a \rightarrow 0$) and with

$$b_n = \sqrt{z_b^2 + r_n^2} = \sqrt{f_1^2 + r_n^2} \quad (2)$$

where r_n is the radius of the n th zone and f_1 is the first-order focal length, squaring and simplifying leads to

$$r_n^2 = p\lambda f_1 + \left(\frac{p\lambda}{2}\right)^2 \quad (3)$$

where $p = n + 2\Delta/\lambda$.

For a finite source or object distance Eq. (3) still holds with the addition of higher-order terms in λ and if the term in λ^2 is multiplied by $(M^3+1)/(M+1)^3$, where M is the magnification. In most practical cases, terms in λ^2 and above are negligible and so, to a good approximation,

$$r_n^2 = n\lambda f_1 + 2\Delta f_1 = n\lambda f_1 + r_0^2 \quad (4)$$

since, for the central zone, $n = 0$ and $r_0^2 = 2\Delta f_1$. Equation (4) describes the *Fresnel zone plate* and, for $r_0 = 0$, the *Fresnel-Soret zone plate* (often referred to as the Fresnel zone plate). The latter is the most commonly used, with

$$r_n^2 = n\lambda f_1 = nr_1^2 \quad (5)$$

The higher-order terms ignored in deriving Eq. (5) result in aberrations. In particular, the term in λ^2 describes spherical aberration but only becomes comparable to the first term when $n \sim 4f_1/\lambda$, which

is rarely the case for x-ray zone plates since focal lengths are typically several orders of magnitude larger than the wavelength.

Equation (5) shows that the focal length is inversely proportional to the wavelength, so that monochromatic radiation with $\lambda/\Delta\lambda \sim N$, where N is the total number of zones, is needed to avoid chromatic aberration. The area of the n th zone is

$$\pi(r_n^2 - r_{n-1}^2) = \pi[n\lambda f_1 - (n-1)\lambda f_1] = \pi\lambda f_1 \quad (6)$$

which is constant, so that each zone contributes equally to the amplitude at the focus if the zone plate is evenly illuminated. The width d_n of the n th zone is

$$d_n = r_n - r_{n-1} = \sqrt{n\lambda f_1} - \sqrt{(n-1)\lambda f_1} = \sqrt{n\lambda f_1} \left[1 - \left(1 - \frac{1}{n}\right)^{1/2} \right] \approx \frac{r_n}{2n} \quad (7)$$

leading to an expression for the first-order focal length

$$f_1 = \frac{r_n^2}{n\lambda} \approx \frac{D_n d_n}{\lambda} \quad (8)$$

where D_n is the diameter of the n th zone. If D is the overall zone plate diameter and d is the outer zone width then

$$f_1 = \frac{Dd}{\lambda} \quad (9)$$

Since zone plates are diffractive optics they have many foci, corresponding to different diffraction orders. The m th-order focus can be described by m zones acting in tandem, so that the effective period is md and the focal lengths are given by

$$f_m = f_1/m \quad m=0, \pm 1, \pm 2, \pm 3, \dots \quad (10)$$

Positive values of m give real foci, while negative values give virtual foci and $m = 0$ corresponds to undiffracted, that is, unfocused radiation.

40.3 ZONE PLATES AS THIN LENSES

The sizes of the focal spots for a point object—the diffraction pattern at a focus—should be determined by successively adding (for an open zone) and subtracting (for a closed zone) the diffraction patterns of circular apertures of radii r_n .¹ However, when N is large enough (theoretically greater than ~ 100 , but in practice much less) a zone plate acts as a thin lens, so that the object, u , and image, v_m (in the m th order), distances are related by

$$\frac{1}{u} + \frac{1}{v_m} = \frac{1}{f_m} \quad (11)$$

and the diffraction pattern at a focus approximates to an Airy pattern.

For a lens of diameter D and focal length f the first zero of the Airy distribution, at a radius $f \tan(1.22\lambda/D)$, defines the lateral resolution ρ via the Rayleigh criterion. For a zone plate, using the expressions for the focal lengths and the small angle approximation, this gives the resolution in the m th order

$$\rho_m = 1.22 \frac{d}{m} \quad (12)$$

Equation (12) shows that, for high resolution, the outermost zone width must be small and that better resolutions can be obtained from higher diffraction orders. However, the lower diffraction efficiencies (see Sec. 40.4) in the higher orders can negate this advantage.

The depth of focus Δf_m is also determined using the thin lens analogy; for a thin lens $\Delta f = \pm 2(f/D)^2\lambda$, which, for a zone plate, leads to

$$\Delta f_m = \pm \frac{f}{2mN} \quad (13)$$

40.4 DIFFRACTION EFFICIENCIES OF ZONE PLATES

The zone plate properties discussed so far depend solely on the relative placement of the zone boundaries; how much radiation can be focused into the various diffraction orders depends additionally on the zone heights and profiles.

Amplitude Zone Plates

A full analysis of the efficiency requires taking the Fourier transform of the zone distribution.² However, if the zone boundaries are in the correct positions, as discussed above, and for an amplitude zone plate in which alternate zones are totally absorbing or transmitting, a simpler discussion suffices. In this case half of the incident radiation is absorbed and half of the rest goes into the zeroth, undiffracted, order. The other even orders vanish since the amplitudes from adjacent zones cancel. The only orders which contribute are $0, \pm 1, \pm 3, \dots$ and, from symmetry, it is clear that the $+m$ th and $-m$ th diffraction efficiencies are equal.

Thus 25 percent of the incident radiation remains to be distributed between the odd orders. The peak amplitudes in each diffraction order are equal, but Eq. (12) shows that the focal spot areas decrease as m^2 . Hence, if ϵ_m is the diffraction efficiency in the m th order,

$$0.25 = 2 \sum_{\substack{m=1 \\ m \text{ odd}}}^{\infty} \epsilon_m = 2\epsilon_1 \sum_{\substack{m=1 \\ m \text{ odd}}}^{\infty} \frac{1}{m^2} = 2\epsilon_1 \frac{\pi^2}{8} \quad (14)$$

so that

$$\epsilon_0 = 0.25; \quad \epsilon_m = \frac{1}{m^2\pi^2} \quad m = \pm 1, \pm 3, \pm 5, \dots; \quad \epsilon_m = 0 \quad m = \pm 2, \pm 4, \dots \quad (15)$$

The first order therefore gives the highest focused intensity, but even so it is only ≈ 10 percent efficient. If the zone boundaries are displaced from the optimum positions then intensity is distributed into the even orders, at the expense of the odd, to a maximum of $1/m^2\pi^2$ (Fig. 2). If the clear zones are not totally transmitting but have amplitude transmission A_1 —because of, for example, a supporting substrate—and the other zones have amplitude transmission A_2 , then the diffraction efficiencies are reduced by a factor $(A_1^2 - A_2^2)$.

The multiplicity of diffraction orders means that this type of zone plate must normally be used with an axial stop and a pinhole, the *order selecting aperture* (OSA), as shown in Fig. 3, to prevent loss of image contrast. The axial stop, typically with a diameter $\approx 0.4D$, reduces the focused intensity and the width of the central maximum of the diffraction pattern, while putting more intensity into the outer lobes. The pinhole also removes any other wavelengths present, so that zone plates can be used as linear monochromators.³

An alternative type of amplitude zone plate, the Gabor zone plate, has, instead of a square wave amplitude transmittance $T(r)$, an approximately sinusoidal one

$$T(r) = \frac{1}{2} \left[1 + \sin \frac{\pi r^2}{\lambda f_1} \right] \quad (16)$$

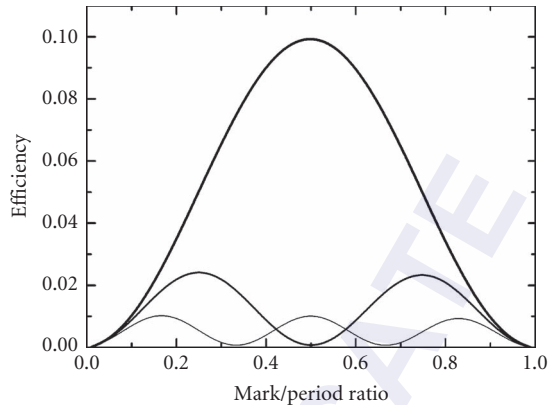


FIGURE 2 Amplitude zone plate diffraction efficiencies: heavy curve, first order; medium curve, second order; light curve, third order.

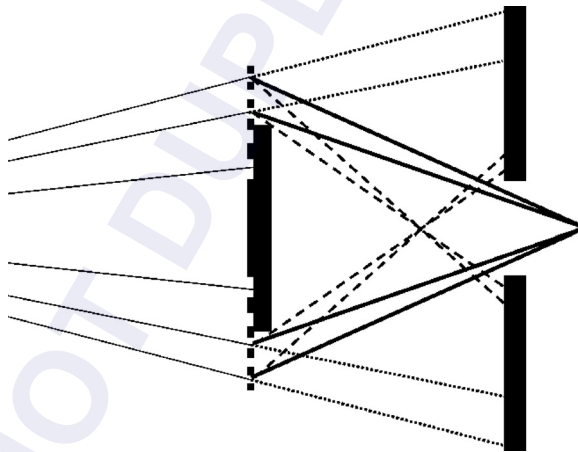


FIGURE 3 Removal of unwanted diffraction orders by use of an axial stop and a pinhole.

The diffraction efficiencies are then 0.25 in the zero order, $1/16$ in the positive and negative first orders and zero in all other orders; the remaining $5/8$ of the incoming intensity is absorbed. The OSA is no longer needed, but the central stop is, and the first-order diffraction efficiency is less than that for an ordinary amplitude zone plate. Gabor zone plates, with the correct profiles, are also more difficult to make.

Phase Zone Plates

If alternate zones can be made to change the phase of the radiation rather than (just) absorbing it, then the amplitude at a focus can be increased. In the absence of absorption, a phase change of π radians would double the focused amplitude so that the diffraction efficiency in the first order would be increased to ≈ 40 percent for rectangular zones. This is not possible for x rays since there

is always some absorption, but a significant improvement in diffraction efficiency can be made if zones of the correct thickness, determined as in the following analysis,⁴ are made.

Pairs of adjacent zones contribute equally to the overall amplitude in a given diffraction order, and so only one pair needs to be considered. The first zone of a pair is assumed to be open and the second has thickness t so that the amplitude is attenuated by a factor $\exp(-2\pi\beta t/\lambda)$ and the phase is retarded by $\Delta\phi = 2\pi\delta t/\lambda$, where δ and β are the optical constants defined by the complex refractive index

$$\tilde{n} = 1 - \delta - i\beta \quad (17)$$

The amplitude at the first-order focus from an open zone is

$$A_o = iC/\pi \quad (18)$$

where $C^2 = I_0$ is the intensity incident on the zone pair. From the phase-shifting zone,

$$A_p = -\frac{iC}{\pi} \exp(-i\Delta\phi) \exp\left(-\frac{2\pi\beta t}{\lambda}\right) \quad (19)$$

so that the contribution from a pair of zones to the intensity at the focus is

$$I_{f_1} = |A_o + A_p|^2 = \left(\frac{C}{\pi}\right)^2 [1 + \exp(-2\eta\Delta\phi) - 2\cos\Delta\phi \exp(-\eta\Delta\phi)] \quad (20)$$

where $\eta = \beta/\delta = 2\pi\beta t/\lambda\Delta\phi$. As for a square-wave amplitude zone plate, the focused intensities in the higher orders are $1/m^2$ in the first order, for odd positive and negative values of m . The maximum intensities are then determined by differentiating Eq. (20) with respect to $\Delta\phi$,

$$\frac{\partial I_{f_m}}{\partial(\Delta\phi)} = 0 = 2\left(\frac{C}{m\pi}\right)^2 [-\eta \exp(-2\eta\Delta\phi) + (\sin\Delta\phi + \eta \cos\Delta\phi) \exp(-\eta\Delta\phi)] \quad (21)$$

Equation (21) shows that the optimum phase shift $\Delta\phi_{\text{opt}}$ is given by the nontrivial solution of

$$\eta \exp(-\eta\Delta\phi_{\text{opt}}) = \sin\Delta\phi_{\text{opt}} + \eta \cos\Delta\phi_{\text{opt}} \quad (22)$$

with two limiting cases $\eta \rightarrow \infty$ for an amplitude zone plate and $\eta \rightarrow 0$ for a phase zone plate with no absorption. Substituting for $\eta \exp(-\eta\Delta\phi_{\text{opt}})$ in Eq. (21) and dividing by C^2 gives the m th-order diffraction efficiency for the optimum phase shift

$$\epsilon_m = \frac{1}{m^2\pi^2} \left(1 + \frac{1}{\eta^2}\right) \sin^2\Delta\phi_{\text{opt}} \quad (23)$$

The undiffracted amplitudes through the open and phase-shifting zones are

$$A_{o_u} = \frac{C}{2} \quad A_{p_u} = \frac{C}{2} \exp(-i\Delta\phi_{\text{opt}}) \exp(-\eta\Delta\phi_{\text{opt}}) \quad (24)$$

so that the zero-order intensity is

$$I_u = |A_{o_u} + A_{p_u}|^2 = \left(\frac{C}{2}\right)^2 [1 + \exp(-2\eta\Delta\phi_{\text{opt}}) + 2\cos\Delta\phi_{\text{opt}} \exp(-\eta\Delta\phi_{\text{opt}})] \quad (25)$$

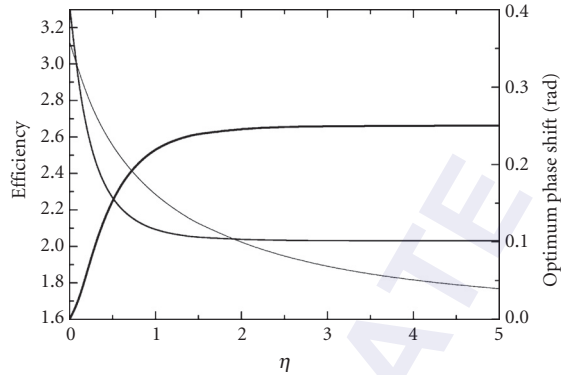


FIGURE 4 Zero (heavy curve) and first-order (medium curve) diffraction efficiencies of a phase zone plate at the optimum phase shift (light curve).

leading to the zero-order efficiency

$$\varepsilon_0 = 0.25 \left[\sin^2 \Delta\phi_{\text{opt}} + \left(2\cos \Delta\phi_{\text{opt}} + \frac{\sin \Delta\phi_{\text{opt}}}{\eta} \right)^2 \right] \quad (26)$$

Since $I_0 = C^2$ is the intensity incident on a zone pair, $I_0/2$ is transmitted by the open zone and $(I_0/2) \exp(-2\eta\Delta\phi_{\text{opt}})$ by the phase-shifting one, so that the total transmitted intensity is

$$I_t = \frac{C^2}{2} [1 + \exp(-2\eta\Delta\phi_{\text{opt}})] \quad (27)$$

leading to the total fractional transmitted intensity at the optimum phase shift

$$\varepsilon_t = 0.5 \left[2 - \left(1 - \frac{1}{\eta^2} \right) \sin^2 \Delta\phi_{\text{opt}} + \frac{1}{\eta} \sin 2\Delta\phi_{\text{opt}} \right] \quad (28)$$

Figure 4 shows the variation of the zero and first-order diffraction efficiencies as functions of η , and Fig. 5 gives an example of the variation of the first-order efficiency with thickness, calculated using Eq. (28) for nickel at a wavelength of 3.37 nm. These figures demonstrate the significant enhancement in efficiency possible over that of an amplitude zone plate.

Applying a similar analysis to a Gabor zone plate gives a corresponding increase in the diffraction efficiency. Higher efficiencies could be obtained by using zone profiles in which the phase shift varies continuously across each zone.⁵ In the absence of absorption it is then possible, in principle, for any given diffraction order to contain 100 percent of the incident intensity. It is not yet possible to make such structures at high resolution, but stepped approximations to the profile have demonstrated efficiencies of ≈ 55 percent at an energy of 7 keV.⁶

Volume Effects

In order to achieve the optimum phase shift discussed in the section “Phase Zone Plates,” the zone thickness required is

$$t_{\text{opt}} = \frac{\Delta\phi_{\text{opt}} \lambda}{2\pi\delta} \quad (29)$$

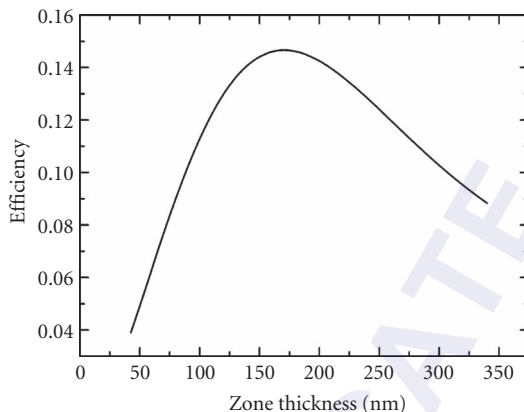


FIGURE 5 The first-order diffraction efficiency of a nickel zone plate at a wavelength of 3.37 nm.

Figure 5 shows that for a nickel zone plate at 3.37 nm, t_{opt} is around 170 nm. For high spatial resolution this means that the aspect ratio, t_{opt}/d , is large, and increases for shorter wavelengths. As well as the resulting technological problems, the previous discussion of spatial resolution in terms of the minimum zone width is no longer valid. The minimum zone width, introduced as the validity criterion of scalar diffraction theory,⁷ is

$$d_{\text{min}} = \sqrt{m\lambda t_{\text{opt}}} \quad (30)$$

this approximation being in good agreement with rigorous electromagnetic theory and with the theory of volume holograms.⁸ For zones with spacing less than d_{min} scalar diffraction theory is not valid due to multiple diffraction of radiation at the zone plate structure. Thus, for the nickel zone plate optimized for a wavelength of 3.37 nm in the first order, the spatial resolution can be no better than about 30 nm.

40.5 MANUFACTURE OF ZONE PLATES

Since the spatial resolution is determined primarily by the outer zone width, taking the discussion of the section “Volume Effects,” into account, zone plates must have small linewidths, along with large areas to provide large apertures and correct zone thicknesses to give optimum efficiencies. In addition, boundaries must be placed within about 1/3 of the outer zone width to maintain efficiencies and focusing properties.⁹ Electron-beam lithography (EBL), which routinely gives zone plates with diameters of around 200 μm and outer zone widths of 25 nm, is now the main method of manufacture but two other techniques—interference (holographic)¹⁰ and sputter and slice¹¹ methods—have historical significance.

In EBL the zone plate pattern is recorded in a polymer resist such as polymethyl methacrylate, followed by etching or electroplating to reproduce the pattern in, for example, nickel with a thickness of ~ 100 to 200 nm for the best efficiency at a few hundred electronvolts or gold or tungsten with thicknesses of ~ 0.5 to 1 μm for a few kiloelectronvolts.¹² Experimental efficiencies are lower than the theoretical optimum values due to manufacturing inaccuracies, primarily misplaced zone boundaries and profile errors.

40.6 BRAGG-FRESNEL LENSES

As discussed in the preceding sections of this chapter, x-ray zone plates work in transmission. However, like gratings they can also be used in reflection, and in this case the resolution-limiting effects of high aspect ratios can be alleviated.¹³ However, since near-normal incidence reflectivities are very small, to allow (near) circular symmetry to be maintained in-phase addition of many reflections is needed, as in crystals and multilayer mirrors. Optics which combine the Bragg reflection of crystals or multilayers with the Fresnel diffraction of gratings or zone plates are known as Bragg-Fresnel lenses.^{14,15} Their properties may be described by considering combinations of zone plates with multilayers or crystals; the generalisation to gratings is obvious.

Properties of Bragg-Fresnel Lenses

The diffraction pattern at a focus is determined as for an ordinary zone plate and the focused intensity is given by the diffraction efficiency combined with the Bragg reflectivity. Spherical waves from point sources S_1 and S_2 (Fig. 6) produce an elliptical interference pattern with S_1 and S_2 at the foci. A slice across the diffraction pattern, perpendicular to the line S_1S_2 , gives the structure of a circular transmission zone plate that focuses radiation emitted at S_1 to S_2 (Fig. 6a). If S_1 is moved to infinity then the interference pattern becomes parabolic and a standard zone plate is formed.

Taking the slice at an angle to the S_1S_2 axis produces an elliptical zone plate which forms a reflected image of S_1 at S_2 (Fig. 6b). The reflectivity is enhanced if the reflecting surface is a crystal or multilayer, with period d equal to the distance between the peaks of the interference pattern (Fig. 6c). Since the Bragg equation must be satisfied, the radiation is monochromatised with a bandpass $\Delta\lambda \sim \lambda/N_L$, where N_L is the number of layer pairs; the monochromaticity requirement of the zone plate, $\lambda/\Delta\lambda$ larger than the number of zones, must also be satisfied.

Defining the origin of the coordinate system to be at the center of lens, with the x and z axes parallel to the multilayer and the y axis perpendicular to the multilayers, the amplitude E of the reflected wave is

$$E(x, y) = r_M \sum_{l=1}^L \int_{Z_l} \exp\left[\frac{2\pi i}{\lambda}(R+r)\right] dr \quad (31)$$

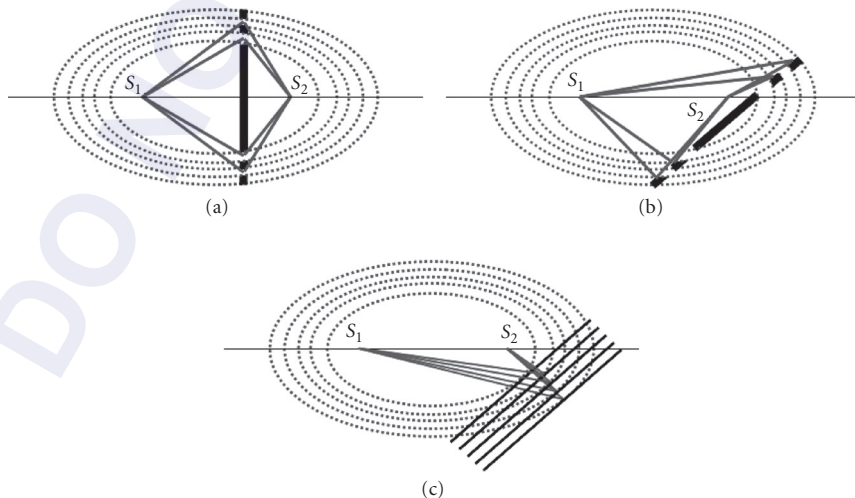


FIGURE 6 Structure of Bragg-Fresnel lenses. (See also color insert.)

where r_M is the peak amplitude reflectivity of the multilayer; the summation is over all layer pairs (l) and the integration is over the zone plate structure for each layer pair. If the source is far from the lens the distances R and r are given by

$$R = R_1 - x \frac{x_1}{R_1} - y \frac{y_1}{R_1} \quad r = r_2 - x \frac{x_2}{r_2} - y \frac{y_2}{r_2} + \frac{x^2}{2r_2} + \frac{y^2}{2r_2} \quad (32)$$

where $R_1 = (x_1^2 + y_1^2)^{1/2}$ is the distance from radiation source at $S_1(x_1, y_1)$ to the center of the lens, and $r_2 = (x_2^2 + y_2^2)^{1/2}$ is the distance from the center of the lens to the focal point $S_2(x_2, y_2)$. Since x varies along the multilayer surface and y varies into the multilayer, with $y = ld$ at the layer interfaces, x and y can be separated and the amplitude at the focal point is

$$E(S_2) = r_M \sum_{l=1}^L \exp \left\{ \frac{2\pi i}{\lambda} \left[-y \left(\frac{y_1}{R_1} + \frac{y_2}{r_2} \right) + \frac{y^2}{2r_2} \right] \right\} \times \int_{z_l} \exp \left\{ \frac{2\pi i}{\lambda} \left[-x \left(\frac{x_1}{R_1} + \frac{x_2}{r_2} \right) + \frac{x^2}{2r_2} \right] \right\} dx \quad (33)$$

The summation describes the wavelength-selecting properties of the multilayer and the integral describes the focusing property of the zone plate. With

$$P_l = \frac{2\pi}{\lambda} \left(\frac{l^2 d^2}{2r_2} - 2ld \sin \theta_0 \right) \quad (34)$$

where θ_0 is the incidence angle giving the maximum reflection at the center of the lens, the summation reduces to

$$G = \sum_{l=1}^L \exp \{iP_l\} \quad (35)$$

and the angular distribution of the reflected radiation is given by

$$\frac{1}{L^2} |G|^2 = \frac{1}{L^2} \left[\left(\sum_{l=1}^L \sin P_l \right)^2 + \left(\sum_{l=1}^L \cos P_l \right)^2 \right] \quad (36)$$

Manufacture of Bragg-Fresnel Lenses

Bragg-Fresnel lenses may be made by masking the surface of a multilayer mirror with an absorbing zone plate or by etching a zone plate pattern into the multilayer.¹⁶ Similar methods can be used for crystal based Bragg-Fresnel lenses.¹⁷ In order to obtain high efficiencies, phase-modulating effects can be used to enhance the efficiency of the zone plate part of the lens. This requires, for example, profiling the multilayer or depositing it on an anisotropically etched substrate.

40.7 REFERENCES

1. A. G. Michette, *Optical Systems for Soft X-Rays*, New York: Plenum, pp. 170–176 (1986).
2. A. G. Michette, *Optical Systems for Soft X-Rays*, New York: Plenum, pp. 178–179 (1986).
3. B. Niemann, D. Rudolph, and G. Schmahl, "Soft X-Ray Imaging Zone Plates with Large Zone Numbers for Microscopic and Spectroscopic Applications," *Opt. Comm.* **12**:160–163 (1974).
4. J. Kirz, "Phase Zone Plates for X-Rays and the Extreme UV," *J. Opt. Soc. Am.* **64**:301–309 (1974).

5. R. O. Tatchyn, "Optimum Zone Plate Theory and Design," in *X-Ray Microscopy, Springer Series in Optical Sciences*, Heidelberg: Springer, **43**:40–50 (1984).
6. E. Di Fabrizio, and M. Gentili, "X-Ray Multilevel Zone Plate Fabrication by Means of Electron-Beam Lithography: Toward High-Efficiency Performances," *J. Vac. Sci. Technol. B* **17**:3439–3443 (1999).
7. A. I. Erko, V. V. Aristov, and B. Vidal, *Diffraction X-Ray Optics*, Bristol: IOP Publishing, pp. 98–101 (1996).
8. R. J. Collier, Ch. B. Burckhardt, and L. H. Lin, *Optical Holography*, New York & London: Academic Press, (1971).
9. M. J. Simpson and A. G. Michette, "The Effects of Manufacturing Inaccuracies on the Imaging Properties of Fresnel Zone Plates," *Optica Acta* **30**:1455–1462 (1983).
10. P. Guttmann, "Construction of a Micro Zone Plate and Evaluation of Imaging Properties," in *X-Ray Microscopy, Springer Series in Optical Sciences*, Heidelberg: Springer, **43**:75–90 (1984).
11. D. Rudolph, B. Niemann, and G. Schmahl, "High Resolution X-Ray Optics," *Proc. SPIE* **316**:103–105 (1982).
12. P. Charalambous, "Fabrication and Characterization of Tungsten Zone Plates for Multi KeV X-Rays," *AIP Conf. Proc.* **507**:625–630 (2000).
13. A. G. Michette, S. J. Pfauntsch, A. Erko, A. Firsov, and A. Svintsov, "Nanometer Focusing of X-Rays with Modified Reflection Zone Plates," *Opt. Commun.* **245**:249–253 (2005).
14. V. V. Aristov, A. I. Erko, and V. V. Martynov, "Principles of Bragg-Fresnel Multilayer Optics," *Rev. Phys. Appl.* **23**:1623–1630 (1988).
15. A. Erko, Y. Agafonov, La. Panchenko, A. Yakshin, P. Chevallier, P. Dhez, and F. Legrand, "Elliptical Multilayer Bragg-Fresnel Lenses with Submicron Spatial Resolution for X-Rays," *Opt. Commun.* **106**:146–150 (1994).
16. A. I. Erko, La. Panchenko, A. A. Firsov, and V. I. Zinenko, "Fabrication and Tests of Multilayer Bragg-Fresnel X-Ray Lenses," *Microelectron Eng.* **13**:335–338 (1991).
17. A. Firsov, A. Svintsov, A. Erko, W. Gudat, A. Asryan, M. Ferstl, S. Shapoval, and V. Aristov, "Crystal-Based Diffraction Focusing Elements for Third-Generation Synchrotron Radiation Sources," *Nucl. Instrum. Methods Phys. Res. A* **467-468**:366–369 (2001).

This page intentionally left blank.

DO NOT DUPLICATE

Eberhard Spiller

*Spiller X-Ray Optics
Livermore, California***41.1 GLOSSARY**

$\tilde{n} = 1 - \delta + i\beta$	refractive index
θ_i	grazing angle of propagation in layer i
φ_i	phase at boundary i
λ	wavelength
q	x-ray wavevector perpendicular to the surface
d	layer thickness
Λ	multilayer period, $\Lambda = d_1 + d_2$
D	total thickness of the multilayer
$f, 1/f$	spatial frequency and spatial period along the surface or boundary
r_{nm}, t_{nm}	amplitude reflection and transmission coefficients
PSD	2-dimensional power spectral density
$a(f)$	roughness replication factor

41.2 INTRODUCTION

The reflectivity of all mirror materials is small beyond the critical grazing angle, and multilayer coatings are used to enhance this small reflectivity by adding the amplitudes reflected from many boundaries coherently, as shown in Fig. 1. Multilayers for the VUV and x-ray region can be seen as an extension of optical coatings toward shorter wavelengths or as artificial one-dimensional Bragg crystals with larger lattice spacings Λ than the d -spacings of natural crystals (see Chap. 39). In contrast to the visible region, no absorption-free materials are available for wavelengths $\lambda < 110$ nm. In addition, the refractive indices of all materials are very close to one, resulting in a small reflectance at each boundary and requiring a large number of boundaries to obtain substantial reflectivity. For absorption-free materials a reflectivity close to 100 percent can always be obtained, independent of the reflectivity r of an individual boundary by making the number of boundaries N sufficiently

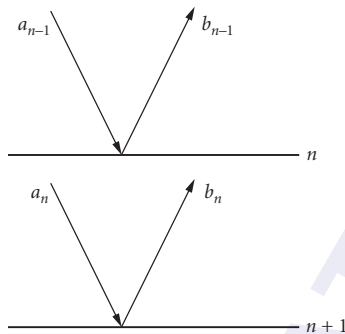


FIGURE 1 Two layers and their boundaries in a multilayer structure with the incoming wave amplitudes a_n and reflected amplitudes b_n in each layer.

large, $Nr \gg 1$. Absorption limits the number of periods that can contribute to the reflectivity in a multilayer to

$$N_{\max} = \frac{\sin^2 \theta}{2\pi\beta} \quad (1)$$

where β is the average absorption index of the coating materials. Multilayers became useful for $N_{\max} \gg 1$, and very high reflectivities can be obtained if N_{\max} is much larger than $N_{\min} = 1/r$, the number required for a substantial reflectivity enhancement.

The absorption index β is in the order of 1 for wavelengths around 100 nm and decreases to very small values at shorter wavelengths in the x-ray range. Materials that satisfy the condition $N_{\max} > 1$ become available for $\lambda < 50$ nm. For $\lambda < 20$ nm, in wavelength regions not too close to absorption edges, β decreases rapidly with decreasing wavelength; $\beta \propto \lambda^3$. The reflected amplitude r from a single boundary also decreases with wavelength, albeit at a slower rate as $r \propto \lambda^2$, and one can compensate for this decrease by increasing the number of boundaries or periods N in a multilayer stack, by using $N \propto 1/\lambda^2$. With this method reflectivities close to 100 percent are theoretically possible at very short wavelengths in the hard x-ray range. A perfect crystal showing Bragg reflection can be seen as a multilayer that realizes this high reflectivity, with the atomic planes located at the nodes of the standing wave within the crystal.

The condition that all periods of a multilayer reflect in phase leads to the Bragg condition for the multilayer period Λ ,

$$m\lambda = 2\Lambda \sin \theta \Rightarrow \Lambda = \frac{m\lambda}{2\sin \theta} \approx \frac{m\lambda}{2\sin \theta_0 \sqrt{1 - 2\delta/\sin^2 \theta_0}} \quad (2)$$

where θ is the effective grazing angle of propagation within the multilayer material, and θ_0 the corresponding angle in vacuum. The refraction correction represented by the square root in Eq. (2) becomes large for small grazing angles, even for small values of δ . The path difference between the amplitude maxima from adjacent periods is $m\lambda$, and multilayers are most of the time used in order $m = 1$. The shortest period Λ for good multilayer performance is limited by the quality of the boundaries (the roughness σ should be smaller than $\Lambda/10$) and this quality is in turn limited by the size of the atoms for noncrystalline materials. Practical values for the shortest period are around $\Lambda = 1.5$ nm. Thus the high reflectivities that are theoretically possible for very hard x rays can only be realized with multilayers of periods $\Lambda > 1.5$ nm, which must be used at small grazing angles. By introducing the momentum transfer q at reflection,

$$q_i = \frac{4\pi}{\lambda} \tilde{n} \sin \theta_i \quad (3)$$

We can express the condition $\Lambda > 1.5 \text{ nm}$ as $|q| < 4 \text{ nm}^{-1}$. It is convenient to introduce the variable q in x-ray optics because, as long as the x-ray wavelength is not too close to an absorption edge, the performance of optical components is determined mainly by q and not by the specific values of λ and θ .

Attempts to observe x-ray interference from thin films started around 1920; by 1931 Kiessig¹ observed and analyzed the x-ray interference structure from Ni films. Multilayer structures produced around 1940² lost their x-ray reflectivity due to diffusion in the structure, and this fact became a tool to measure small diffusion coefficients.³⁻⁵ The usefulness of multilayers for normal incidence optics in the XUV and EUV was recognized in 1972.⁶ The deposition processes to produce these structures were developed by many groups in the next two decades. Today multilayer telescopes on the orbiting observatories SOHO and TRACE provide high-resolution EUV pictures of the solar corona (see <http://umbra.nascom.nasa.gov/eit/> and <http://vestige.lmsal.com/TRACE/>). Cameras with multilayer coated mirrors for the EUV are a main contender for the fabrication of the next generation of computer chips, and multilayer mirrors are found at the beamlines of all synchrotron radiation facilities and in x-ray diffraction equipment as beam deflectors, collimators, filters, monochromators, polarizers, and imaging optics. (See Ref. 7 and Chaps. 28, 43 to 46, and 54 in this volume for more details.)

41.3 CALCULATION OF MULTILAYER PROPERTIES

The theoretical treatment of the propagation of x rays in layered structures does not differ from that for other wavelength regions as discussed by Dobrowolski in Chap. 7 in Vol. IV of this *Handbook* and in many textbooks and review articles. At the boundary of two materials or at an atomic plane an incident wave is split into a transmitted and reflected part and the total field amplitude in the structure is the superposition of all these waves. Figure 1 sketches two layers as part of a multilayer and the amplitudes of the forward a_n and backward b_n running waves in each layer. The amplitudes in each layer are coupled to those of the adjacent layers by linear equations that contain the amplitude transmission $t_{n,m}$ and reflection coefficients r_{nm} ,

$$\begin{aligned} a_n &= a_{n-1} t_{n-1,n} e^{i\varphi_n} + b_n e^{2i\varphi_n} r_{n,n-1} \\ b_n &= a_n r_{n,n+1} + b_{n+1} e^{i\varphi_{n+1}} t_{n+1,n} \end{aligned} \quad (4)$$

The phase delay φ_n due to propagation through layer n of thickness d and angle θ_n from grazing incidence is given by

$$\varphi_n = \frac{2\pi}{\lambda} \tilde{n}_n d_n \sin \theta_n \quad (5)$$

and the transmitted and reflected amplitudes at each boundary are obtained from the Fresnel equations

$$r_{n,n+1} = \frac{q_n - q_{n+1}}{q_n + q_{n+1}} \quad (6)$$

$$t_{n,n+1} = \frac{2q_n}{q_n + q_{n+1}} \quad (7)$$

with the q values as defined in Eq. (3) for s-polarization and $q_i = (4\pi/\lambda\tilde{n}_i) \sin \theta_i$ for p-polarization.

Matrix methods^{8,9} are convenient to calculate multilayer performance using high-level software packages that contain complex matrix manipulation. The transfer of the amplitudes over a boundary and through a film is described by⁹

$$\begin{pmatrix} a_i \\ b_i \end{pmatrix} = \frac{1}{t_{i-1,i}} \begin{pmatrix} e^{i\varphi_i} & r_{i-1,i} e^{-i\varphi_i} \\ r_{i-1,i} e^{i\varphi_i} & e^{-i\varphi_i} \end{pmatrix} \begin{pmatrix} a_{i+1} \\ b_{i+1} \end{pmatrix} \quad (8)$$

and the transfer between the incident medium ($i = 0$) and the substrate ($i = n + 1$) by

$$\begin{pmatrix} a_0 \\ b_0 \end{pmatrix} = \frac{\prod_{i=1}^{i=n+1} M_i}{\prod_{i=1}^{i=n+1} t_{i,i+1}} \begin{pmatrix} a_{n+1} \\ b_{n+1} \end{pmatrix} \quad (9)$$

with the matrices M_i as defined in Eq. (8). For incident radiation from the top ($b_{n+1} = 0$) we can calculate the reflected and transmitted amplitudes from the elements of the product matrix m_{ij} as

$$\begin{aligned} r_{\text{ML}} &= b_0/a_0 = m_{21}/m_{11} \\ t_{\text{ML}} &= a_{n+1}/a_0 = (\prod t_{i,i+1})/m_{11} \end{aligned} \quad (10)$$

The reflected intensity is $R_{\text{ML}} = r_{\text{ML}} r_{\text{ML}}^*$ and the transmitted intensity $T_{\text{ML}} = t_{\text{ML}} t_{\text{ML}}^* (q_{n+1}/q_0)$ with the q values of Eq. (3).

In another convenient matrix formalism due to Abeles each matrix contains only the parameters of a single film;^{8,10} however, the transfer matrix method given earlier⁹ is more convenient when one wants to include the imperfections of the boundaries.

The achievable boundary quality is always of great concern in the fabrication of x-ray multilayers. The first effect of boundary roughness or diffusion is the reduction of the reflected amplitude due to dephasing of the contributions from different depth within the boundary layer. If one describes the reflected amplitude $r(z)$ as a function of depth by a Gaussian of width σ , one obtains a Debye-Waller factor for the reduction of boundary reflectivity,

$$r/r_o = \exp(-0.5 q_1 q_2 \sigma^2) \quad (11)$$

where r_o is the amplitude reflectivity for a perfect boundary and the q values are those of the films on either side of the boundary.^{7,11-14} Reducing the amplitude reflection coefficients in Eq. (8) by the Debye-Waller factor of Eq. (11) gives a good estimate of the influence of boundary roughness on multilayer performance. The roughness σ has by most authors been used as a fitting parameter that characterizes a coating. Usually the Debye-Waller factor is given for the intensity ratio [absolute value of Eq. (11) squared] with the vacuum q -values for both materials (refraction neglected). The phase shift in Eq. (11) due to the imaginary part of q produces a shift of the effective boundary toward the less absorbing material.

For a boundary without scattering (gradual transition of the optical constants or roughness at very high spatial frequencies) the reduction in reflectivity is connected with an increase in transmission according to $t_{12} t_{21} + r_{12}^2 = 1$. Even if the reflectivity is reduced by a substantial factor, the change in transmission can in most cases be neglected, because reflectivities of single boundaries are typically in the 10^{-4} range and transmissions close to one. Roughness that scatters a substantial amount of radiation away from the specular beam can decrease both the reflection and transmission coefficients, and the power spectral density (PSD) of the surface roughness over all relevant spatial frequencies has to be measured to quantify this scattering.

It is straightforward to translate Eqs. (4) to (11) into a computer program, and a personal computer gives typically a reflectivity curve of a 100 layer coating for 100 wavelengths or angles within a few seconds. Multilayer programs and the optical constants of all elements and of many compounds can be accessed on (<http://www-cxro.lbl.gov>) or downloaded¹⁵ (www.rxcollc.com/idl)(see Chap. 36). The first site also has links to other relevant sites, and the Windt programs can also calculate nonspecular scattering.

41.4 FABRICATION METHODS AND PERFORMANCE

Multilayer x-ray mirrors have been produced by practically any deposition method. Thickness errors and boundary roughness are the most important parameters; they have to be kept smaller than $\Lambda/10$ for good performance. Magnetron sputtering, pioneered by Barbee,¹⁶ is most widely used: sputtering systems are very stable, and one can obtain the required thickness control simply by timing. The same is true for Ion beam deposition systems.¹⁷ Thermal evaporation systems usually use an in situ soft x-ray

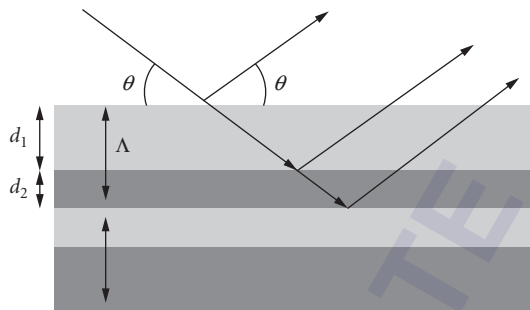


FIGURE 2 A quasiperiodic multilayer. The individual layer thicknesses change while the sum Λ remains constant.

reflectometer to control the thickness and an additional ion gun for smoothing of the boundaries.^{18–21} In the “quarter wave stack” of multilayer mirrors for the visible, all boundaries add their amplitudes in-phase to the total reflectivity, and each layer has the same optical thickness and extends between a node and an antinode of the standing wave generated by the incident and reflected waves. To mitigate the effect of absorption in the VUV and x-ray region one first selects a “spacer” material with the lowest available absorption and combines it with a “reflector” layer with good contrast with the optical constants of the spacer. The optimum design minimizes absorption by reducing the thickness of the reflector layer and by attempting to position it close to the nodes of the standing wave field, while for the spacer layer the thickness is increased and it is centered around the antinodes. The design accepts some dephasing of the contributions from adjacent boundaries and the optimum is a compromise between the effects of this dephasing and the reduction in absorption. The best values for $\gamma = d_r/\Lambda$ are between 0.3 and 0.4.^{7,22,23} Optimum multilayers for longer x-ray wavelengths ($\lambda > 20$ nm) are quasiperiodic, as shown in Fig. 2, with decreasing γ from the bottom to the top of a multilayer stack. One always attempts to locate the stronger absorber close to the nodes of the standing wavefield within the coating to reduce absorption, and the absorption reduction is greater near the top of the stack where the standing wave between incident and reflected radiation has more contrast.⁶

The paper by Rosenbluth²³ gives a compilation of the best multilayer materials for each photon energy in the range from 100 to 2000 eV. Absorption is the main performance limit for multilayers of larger period $\Lambda > 80$ nm used at longer wavelength near normal incidence. In this region it is important to select the material with the smallest absorption as the first component of a structure (usually a light element at the long wavelength side of an absorption edge) and the absorption of this material becomes the main limit for the multilayer performance. A list of materials of low absorption at selected wavelengths with their values for β and N_{\max} is given in Table 1.

Boundary Quality

The quality of the boundary is the main limitation for the performance of a multilayer with short periods. Roughness of a boundary scatters radiation away from the specular beam and reduces both reflectivity and transmission at each boundary. Diffusion of the two materials at a boundary (or roughness at very high spatial frequencies that scatter into evanescent waves) also reduces the reflectivity but can increase the transmission. Deposition at high energy of the incident materials enhances relaxation of atoms of the growing surface and reduces roughness but also increases diffusion at the boundary. For each set of coating materials the deposition energy for the best compromise has to be found. A good solution is to separate the two problems, depositing each film at low energy to minimize diffusion and then polishing the top of each layer after the deposition. Examples are thermal deposition of each layer and ion polishing after deposition,²⁴ or low-energy deposition at the start of each layer and higher energy near the top.²⁵ Polishing of thin films by an ion beam has also been used to remove defects from the substrates of masks in EUV lithography. In a Mo/Si multilayer one can either deposit a thicker Si film than

TABLE 1 Absorption Index β and N_{\max} for the Largest q Values (Normal Incidence for $\lambda > 3.15$ nm and $\sin \theta = \lambda/\pi$ for $\lambda < 3.14$ nm) of Good Spacer Materials near Their Absorption Edges and Absorption of Some Materials at $\lambda = 0.154$ nm

	λ_{edge} (nm)	β	N_{\max}
Mg-L	25.1	7.5×10^{-3}	21
Al-L	17.1	4.2×10^{-3}	38
Si-L	12.3	1.6×10^{-3}	99
Be-K	11.1	1.0×10^{-3}	155
Y-M	8.0	3.5×10^{-3}	45
B-K	6.6	4.1×10^{-4}	390
C-K	4.37	1.9×10^{-4}	850
Ti-L	3.14	4.9×10^{-4}	327
N-K	3.1	4.4×10^{-5}	3,580
Sc-L	3.19	2.9×10^{-4}	557
V-L	2.43	3.4×10^{-4}	280
O-K	2.33	2.2×10^{-5}	3,980
Mg-K	0.99	6.6×10^{-6}	2,395
Al-K	0.795	6.5×10^{-6}	1,568
Si-K	0.674	4.2×10^{-6}	1,744
SiC	0.674	6.2×10^{-6}	1,182
TiN	3.15	4.9×10^{-4}	327
Mg ₂ Si	25.1	7.4×10^{-3}	21
Mg ₂ Si	0.99	6.8×10^{-6}	2,324
Be	0.154	2.0×10^{-9}	189,000
B	0.154	5.7×10^{-9}	67,450
C	0.154	1.2×10^{-8}	32,970
Si	0.154	1.7×10^{-7}	5,239
Ni	0.154	5.1×10^{-7}	750
W	0.154	3.9×10^{-6}	98

needed and etch the excess thickness away with the ion beam or use a more aggressive deposition/etch process on the Si substrate before the multilayer reflecting coating is applied.²⁶

Most good multilayer systems can be described by a simple growth model: particles or atoms arrive randomly and can relax sideways to find locations of lower energy. The random deposition produces a flat power spectrum at low frequencies for the surface roughness with most of the roughness at very high spatial frequencies. The relaxation then reduces roughness at the highest spatial frequencies.

Roughness is characterized by a power spectral density (PSD) such that the intensity of the scattered light is proportional to the PSD. Consistent values for the PSD can be obtained by atomic microscopy and from scatter measurements.²⁷⁻²⁹ The roughness height σ from the spatial frequency range from f_1 to f_2 is related to the PSD by

$$\sigma^2 = 2\pi \int_{f_1}^{f_2} \text{PSD}(f) f df \quad (12)$$

We assume that the surface is isotropic in f and the roughness in Eq. (12) is obtained by integrating over rings with area $2\pi f df$. During the development phase of multilayer x-ray mirrors it was fortuitous that practically perfect smooth substrates were available as Si wafers, float glass, and mirrors fabricated for laser gyros. The 2-dimensional power spectral density (see Chap. 8 by Eugene L. Church and Peter Z. Takacs in Vol. I and Chap. 44 of this volume) of a film on a perfectly smooth substrate is given by^{7,30-32}

$$\text{PSD}(q_s, d) = \Omega \frac{1 - \exp(-2l_r^{n-1} dq_s^n)}{2l_r^{n-1} q_s^n} \quad (13)$$

TABLE 2 Growth Parameters of Multilayer Systems*

System	Λ (nm)	N	Ω (nm ³)	l_r (nm)	n	d (nm)	σ (nm)
Co/C	3.2	150	0.016	1.44	4	480	0.14
Co/C	2.9	144	0.016	1.71	4	423	0.12
Co/C	2.4	144	0.016	1.14	4	342	0.15
Co/C	3.2	85	0.016	1.71	4	274	0.19
Mo/Si	7.2	24	0.035	1.71	4	172	0.14
Ni/C (a)	5.0	30	0.035	1.71	4	150	0.14
W/B ₄ C (b)	1.22	350	0.01	1.22	4	427	0.10
W/B ₄ C (b)	1.78	255	0.018	1.22	3	454	0.12
Mo/Si (c)	6.8	40	0.035	1.36	4	280	0.19
Mo/Si (d)	6.8	40	0.055	1.20	2	280	0.12
Si in Mo/Si (c)	6.8	40	0.02	1.36	4	280	0.14
Mo in Mo/Si (c)	6.8	40	0.5	1.36	4	280	0.23

*Multilayer period Λ , number of periods N , growth parameters Ω , l_r , and n , total thickness d , and roughness σ of the multilayer film calculated for a perfectly smooth substrate within the spatial frequency range $f = 0.001 - 0.25 \text{ nm}^{-1}$.

Coating (a) produced by dual ion beam sputtering courtesy of J. Pedulla (NIST), (b) of Y. Platonov (Osmic). Data for Mo/Si produced by magnetron sputtering (c) from Ref. 36, and ion beam sputtering (d) from P. Kearney and D. Stearns (LLNL). Parameters are the average value for both coating materials, for sputtered Mo/Si we give also the parameters for Mo and Si separately.

where Ω is the particle volume, l_r is a relaxation length, d is the total thickness of the coating, and f a spatial frequency on the surface. The parameter n characterizes the relaxation process; $n = 1$ indicates viscous flow, 2 condensation and re-evaporation, 3 bulk diffusion, and 4 surface diffusion. Equation (13) yields the flat PSD of the random deposition with roughness $\sigma \propto \sqrt{d}$ for small spatial frequencies and a power law with exponent n that is independent of thickness at high spatial frequencies. Roughness is replicated throughout a stack with a replication factor

$$a(f) = \exp(-l_r^{n-1} d q_s^n) \quad (14)$$

so that the PSD of a film on a substrate has a power spectral density PSD^{tot}

$$\text{PSD}^{\text{tot}} = \text{PSD}^{\text{film}} + a^2 \text{PSD}^{\text{sub}} \quad (15)$$

The growth parameters of some multilayer structures are given in Table 2.³³ Note that the roughness values in the last column do not include contributions from spatial frequencies $f > 0.25 \text{ nm}^{-1}$ and from diffuse transition layers. These values have to be added in the form $\sigma_{\text{tot}}^2 = \sigma_1^2 + \sigma_2^2$ in Eq. (11) for the calculation of the reflectivity.

High-Reflectivity Mirrors

A compilation of the reflectivities of multilayers obtained by different groups can be found at www-cxro.lbl.gov/multilayer/survey.html. The highest normal incidence reflectivities, around 70 percent, have been obtained with Mo/Be and Mo/Si wavelengths around $\lambda = 11.3$ and 13 nm, near the absorption edges of Be and Si.³⁴ The peak reflectivity drops at longer wavelengths due to the increased absorption. At shorter wavelengths, roughness becomes important and reduces the reflectivity of normal incidence mirrors. One can, however, still obtain very high reflectivity at short wavelengths by keeping the multilayer period above 2 nm and using grazing angles of incidence to reach short wavelengths.

For hard x rays, where N_{max} is very much larger than the number, $N_{\text{min}} \sim 1/r$ needed to obtain good reflectivity, one can position multilayers with different periods on top of each other (depth-graded multilayers, as shown in Fig. 3) to produce a coating with a large spectral or angular bandwidth. Such “supermirrors” are being proposed to extend the range of grazing incidence telescopes up to the 100 keV.^{35,36} “Supermirrors” are common for cold neutrons where absorption-free materials are

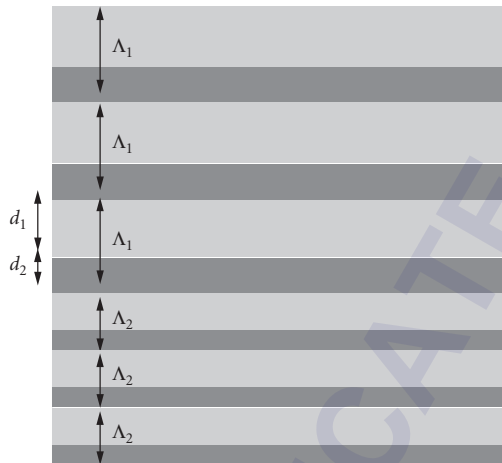


FIGURE 3 A depth-graded multilayer. The repeat distance Λ changes with depth. (The change can be abrupt, as shown, or more gradual.)

available.³⁷ The low absorption for hard x rays makes it also possible to produce coatings with very narrow bandwidth. The spectral width is determined by the effective number of periods contributing to the reflectivity, $\lambda/\Delta\lambda = N_{\text{eff}}$. Reducing the reflectivity from a single period, for example, by using a very small value of $\gamma = d_H/\Lambda$, allows radiation to penetrate deeper into the stack, allowing more layers to contribute, thus reducing the bandwidth.⁷

Multilayer Coated Optics

Multilayers for imaging optics usually require a lateral grading of the multilayer period Λ across the face of the optic to adapt to the varying angle of incidence according to Eq. (2), as shown in Fig. 4. Many methods have been used to produce the proper grading of the multilayer period during deposition, among them shadow masks in front of the rotating optics, substrate tilt, speed control of the platter that moves the optics over the magnetrons, and computer-controlled shutters.^{38–44} The reproducibility that has been achieved in EUV lithography from run to run is better than 0.1 percent. Multilayer mirrors

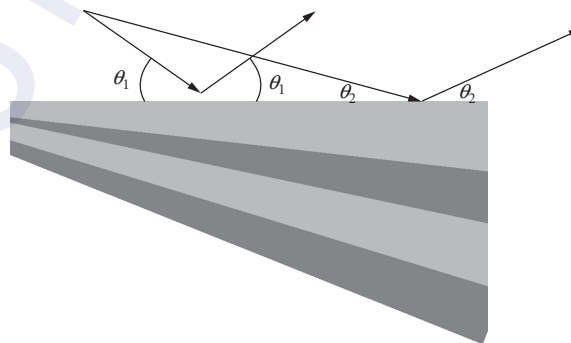


FIGURE 4 Laterally graded multilayer.

with a linear grading of the thickness on parabolically bent Si wafers are commercially available (e.g. www.rigaku.com/optics/index.html) and are being used as collimators in x-ray diffraction experiments at $\lambda = 0.154$ nm. At this wavelength the reflectivity is over 75 percent and the flux to the diffractometer can be increased by about an order of magnitude.⁴⁶

It is still a challenge to produce figured substrates that have both good figure and finish in the 0.1-nm range.⁴⁷ Multilayer mirrors can remove high spatial roughness with spatial periods of less than 50 nm and graded coatings can be used to correct low-order figure errors; however, mirrors that do not meet specifications for high resolution imaging usually have considerable roughness at spatial frequencies between these limits that cannot be modified by thin films.³³

Diffraction-limited performance of multilayer coated mirrors has been achieved in the cameras for EUV lithography. The mirrors of these cameras usually require lateral grading of the layer thickness due to the change of the angle of incidence over the surface of a mirror. These graded coatings modify the shape of the mirror by two components: the total thickness of the multilayer and the shift in the phase of the reflected wave. The two effects have opposite sign in their sensitivity to thickness errors, and for the mirrors used in EUV lithography around $\lambda = 13.5$ nm the total error in the figure is around 75 percent of the error produced by an error in the multilayer thickness alone. After subtracting the changes that can be compensated by alignment the remaining figure error added by the coating is well below 0.1 nm.^{48,49}

Polarizers and Phase Retarders

The Brewster angle in the VUV and x-ray region occurs close to 45° for all materials, so all reflectors used near 45° are effective polarizers. Multilayer coatings do not change the ratio of the reflectivities for s- and p-polarization, but enhance the reflectivity for s-polarization to useful values.^{50,51} The reflectivity for p-polarization at the Brewster angle is zero for absorption-free materials but increases with absorption. Therefore the achievable degree of polarization is higher at shorter wavelengths where absorption is lower. Typical values for the reflectivity ratio R_s/R_p are 10 around $\lambda = 30$ nm and over 1000 around $\lambda = 5$ nm.

It is not possible to design effective 90° phase retarding multilayer reflectors with high reflectivity for both polarizations. The narrower bandwidth of the reflectivity curve for p-polarization allows one to produce a phase delay at incidence angles or wavelengths that are within the high reflectivity band for s- but not for p-polarization; however, because of the greatly reduced p-reflectivity of such a design, one cannot use them to transform linear polarized into circular polarized radiation.⁵² Multilayers used in transmission offer a better solution. Near the Brewster angle p-polarized radiation is transmitted through a multilayer structure without being reflected by the internal boundaries and has a transmission determined mainly by the absorption in the multilayer stack. A high transmission for s-polarization is also obtained from a multilayer stack that is used off-resonance near a reflectivity minimum or transmission maximum. However, the transmitted radiation is delayed due to the internal reflection within the multilayer stack. Calculated designs for wavelength of 13.4 nm^{53–56} produce a phase retardation of 80° at a grazing angle of 50° . One can tune such a phase retarder in wavelength by producing a graded coating or by using different incidence angles.⁵⁷ A 90° phase retarder is possible with high-quality multilayers that can obtain peak reflectivities above 70 percent. Boundary roughness reduces the phase retardation because it reduces the effective number of bounces within the structure. The maximum phase retardation measured at x-ray wavelengths $\lambda < 5$ nm are in the 10° range.⁵⁸

41.5 MULTILAYERS FOR DIFFRACTIVE IMAGING

The development of free electron lasers (FELs) for x rays promises imaging at high resolution of general specimens; the specimens do not have to be crystallized. An image of the specimen is reconstructed from diffraction patterns produced by powerful, short coherent pulses. The specimen is destroyed

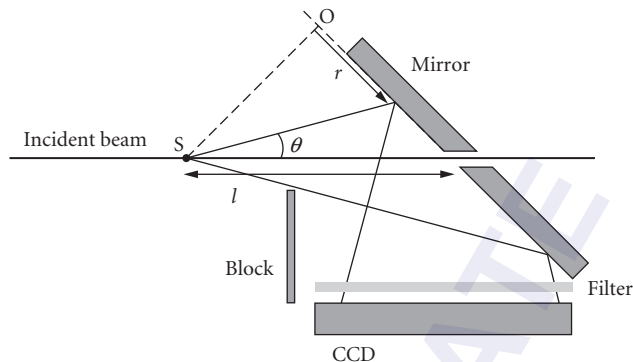


FIGURE 5 Camera used for diffraction imaging. The FEL pulse illuminates the sample S. Radiation diffracted by the sample S is reflected by the multilayer mirror to the CCD detector while the direct beam passes a hole in the mirror. (From Ref. 60.)

by the radiation, but the diffraction pattern recorded during the short pulse represents the specimen before it explodes. Thousands of such diffraction patterns from identical specimens at different rotation angles are needed for a high-resolution 3-D reconstruction of the specimen. Multilayer mirrors have been used to safely direct the diffracted radiation to the detector, transmitting the direct beam through a hole in the center and suppressing the light produced by the exploding specimen (Fig. 5). The thickness of the layers in the multilayer mirror is graded in such a way as to reflect only the wavelength of the incident beam at the correct angle to the detector. First experiments by Chapman et al.^{59,60} have demonstrated the principle using EUV radiation, but considerable challenges remain to transfer the technique to the x-ray region and to interesting 3-D specimen. (See Chap. 27.)

41.6 REFERENCES

1. H. Kiessig, "Interferenz von Röntgenstrahlen an dünnen Schichten," *Ann. der Physik* 5. Folge **10**:769–788 (1931).
2. J. DuMond and J. P. Joutz, "An X-Ray Method of Determining Rates of Diffusion in the Solid State," *J. Appl. Phys.* **11**:357–365 (1940).
3. H. E. Cook and J. E. Hilliard, "Effect of Gradient Energy on Diffusion in Gold–Silver Alloys," *J. Appl. Phys.* **40**:2191–2198 (1969).
4. J. B. Dinklage, "X-Ray Diffraction by Multilayered Thin-Film Structures and Their Diffusion," *J. Appl. Phys.* **38**:3781–3785 (1967).
5. A. L. Greer and F. Spaepen, "Diffusion," in *Modulated Structures*, edited by L. Chang and B. C. Giess (Academic Press, New York, 1985), p. 419.
6. E. Spiller, "Low-Loss Reflection Coatings Using Absorbing Materials," *Appl. Phys. Lett.* **20**:365–367 (1972).
7. E. Spiller, *Soft X-Ray Optics* (SPIE Optical Engineering Press, Bellingham, WA, 1994).
8. F. Abelès, "Recherches sur la propagation des ondes électromagnétique inusoidales dans les milieux stratifiés. Application aux couches minces," *Ann. de Physique* **5**:596–639 (1950).
9. O. S. Heavens, *Optical Properties of Thin Solid Films* (Dover, New York, 1966).
10. M. Born and E. Wolf, *Principles of Optics*, 5th ed. (Pergamon Press, Oxford, 1975).
11. P. Croce, L. Névoit, and B. Pardo, "Contribution à l'étude des couches mince par réflexion spéculaire de rayon X," *Nouv. Rev. d'Optique appliquée* **3**:37–50 (1972).

12. P. Croce and L. Névot, "Étude des couches mince et des surfaces par réflexion rasante, spéculaire ou diffuse, de rayon X," *J. De Physique Appliquée* **11**:113–125 (1976).
13. L. Névot and P. Croce, "Caractérisation des surfaces par réflexion rasante de rayon X, Application à étude du polissage de quelques verres silicates," *Revue Phys. Appl.* **15**:761–779 (1980).
14. F. Stanglmeier, B. Lengeler, and W. Weber, "Determination of the Dispersive Correction $f''(E)$ to the Atomic Form Factor from X-Ray Reflection," *Acta Cryst.* **A48**:626–639 (1992).
15. D. Windt, "IMD—Software for Modeling the Optical Properties of Multilayer Films," *Computers in Physics* **12** (4):360–370 (1998).
16. T. W. Barbee, Jr., "Multilayers for X-Ray Optics," *Opt. Eng.* **25**:893–915 (1986).
17. E. Spiller, S. Baker, P. Mirkarimi, et al., "High Performance Mo/Si Multilayer Coatings for EUV Lithography Using Ion Beam Deposition," *Appl. Opt.* **42**:4049–4058 (2003).
18. E. Spiller, A. Segmüller, J. Rife, et al., "Controlled Fabrication of Multilayer Soft X-Ray Mirrors," *Appl. Phys. Lett.* **37**:1048–1050 (1980).
19. E. Spiller, "Enhancement of the Reflectivity of Multilayer X-Ray Mirrors by Ion Polishing," *Opt. Eng.* **29**:609–613 (1990).
20. E. J. Puik, M. J. van der Wiel, H. Zeijlemaker, et al., "Ion Bombardment of X-Ray Multilayer Coatings: Comparison of Ion Etching and Ion Assisted Deposition," *Appl. Surface Science* **47**:251–260 (1991).
21. E. J. Puik, M. J. van der Wiel, H. Zeijlemaker, et al., "Ion Etching of Thin W Layers: Enhanced Reflectivity of W-C Multilayer Coatings," *Appl. Surface Science* **47**:63–76 (1991).
22. A. V. Vinogradov and B. Ya. Zel'dovich, "X-Ray and Far UV Multilayer Mirrors; Principles and Possibilities," *Appl. Optics* **16**:89–93 (1977).
23. A. E. Rosenbluth, "Computer Search for Layer Materials that Maximize the Reflectivity of X-Ray Multilayers," *Revue Phys. Appl.* **23**:1599–1621 (1988).
24. E. Spiller, "Smoothing of Multilayer X-Ray Mirrors by Ion Polishing," *Appl. Phys. Lett.* **54**:2293–2295 (1989).
25. Fredrik Eriksson, Goeran A. Johansson, Hans M. Hertz, et al., "Enhanced Soft X-Ray Reflectivity of Cr/Sc Multilayers by Ion-Assisted Sputter Deposition," *Opt. Eng.* **41** (11):2903–2909 (2002).
26. P. B. Mirkarimi, E. Spiller, S. L. Baker, et al., "A Silicon-Based, Sequential Coat- and- Etch Process to Fabricate Nearly Perfect Substrate Surfaces," *J. Nanoscience and Nanotechnology* **6**:28–35 (2006).
27. E. M. Gullikson, D. G. Stearns, D. P. Gaines, et al., "Non-Specular Scattering from Multilayer Mirrors at Normal Incidence," presented at the *Grazing Incidence and Multilayer X-Ray Optical Systems*, 1997 (unpublished).
28. E. M. Gullikson, "Scattering from Normal Incidence EUV Optics," *Proc. SPIE* **3331**:72–80 (1998).
29. V. Holý, U. Pietsch, and T. Baumbach, *High Resolution X-Ray Scattering from Thin Films and Multilayers* (Springer, Berlin, 1999).
30. D. G. Stearns, D. P. Gaines, D. W. Sweeney, et al., "Nonspecular X-Ray Scattering in a Multilayer-Coated Imaging System," *J. Appl. Phys.* **84** (2):1003–1028 (1998).
31. E. Spiller, D. G. Stearns, and M. Krumrey, "Multilayer X-Ray Mirrors: Interfacial Roughness, Scattering, and Image Quality," *J. Appl. Phys.* **74**:107–118 (1993).
32. W. M. Tong and R. S. Williams, "Kinetics of Surface Growth: Phenomenology, Scaling, and Mechanisms of Smoothing and Roughening," *Annu. Rev. Phys. Chem.* **45**:401–438 (1994).
33. E. Spiller, S. Baker, E. Parra, et al., "Smoothing of Mirror Substrates by Thin Film Deposition," *Proc. SPIE* **3767**:143–153 (1999).
34. C. Montcalm, R. F. Grabner, R. M. Hudyma, et al., "Multilayer Coated Optics for an Alpha-Class Extreme-Ultraviolet Lithography System," *Proc. SPIE* **3767** (1999).
35. P. Hoghoj, E. Ziegler, J. Susini, et al., "Focusing of Hard X-Rays with a W/Si Supermirror," *Nucl Instrum Meth Phys Res B* **132** (3):528–533 (1997).
36. K. D. Joensen, P. Voutov, A. Szentgyorgyi, et al., "Design of Grazing-Incidence Multilayer Supermirrors for Hard-X-Ray Reflectors," *Appl Opt* **34** (34):7935–7944 (1995).
37. F. Mezei, "Multilayer Neutron Optical Devices," in *Physics, Fabrication, and Applications of Multilayered Structures*, edited by P. Dhez and C. Weisbuch (Plenum Press, New York, 1988), pp. 311–333.
38. D. J. Nagel, J. V. Gilfrich, and T. W. Barbee, Jr., "Bragg Diffractors with Graded-Thickness Multilayers," *Nucl. Instrum. Methods* **195**:63–65 (1982).

39. D. G. Stearns, R. S. Rosen, and S. P. Vernon, "Multilayer Mirror Technology for Soft-X-Ray Projection Lithography," *Appl. Opt.* **32** (34):6952–6960 (1993).
40. S. P. Vernon, M. J. Carey, D. P. Gaines, et al., "Multilayer Coatings for the EUV Lithography Test Bed," presented at the *POS Proc. on Extreme Ultraviolet Lithography*, Monterey, Calif., 1994 (unpublished).
41. E. Spiller and L. Golub, "Fabrication and Testing of Large Area Multilayer Coated X-Ray Optics," *Appl. Opt.* **28**:2969–2974 (1989).
42. E. Spiller, J. Wilczynski, L. Golub, et al., "The Normal Incidence Soft X-Ray, $\lambda = 63.5$ Å Telescope of 1991," *Proc. SPIE* **1546**:168–174 (1991).
43. D. L. Windt and W. K. Waskiewicz, "Multilayer Facilities Required for Extreme-Ultraviolet Lithography," *J. Vac. Sci. Technol. B* **12** (6):3826–3832 (1994).
44. M. P. Bruijn, P. Chakraborty, H. W. van Essen, et al., "Automatic Electron Beam Deposition of Multilayer Soft X-Ray Coatings with Laterally Graded d Spacing," *Opt. Eng.* **25**:916–921 (1986).
45. D. W. Sweeney, R. M. Hudyma, H. N. Chapman, et al., "EUV Optical Design for a 100-nm CD Imaging System," *Proc. SPIE* **3331**:2–10 (1998).
46. M. Schuster and H. Göbel, "Parallel-Beam Coupling into Channel-Cut Monochromators Using Curved Graded Multilayers," *J. Phys. D: Appl. Phys.* **28**:A270–A275 (1995).
47. J. S. Taylor, G. E. Sommargren, D. W. Sweeney, et al., "The Fabrication and Testing of Optics for EUV Projection Lithography," *Proc. SPIE* **3331**:580–590 (1998).
48. R. Soufli, E. Spiller, M. A. Schmidt, et al., "Multilayer Optics for an Extreme Ultraviolet Lithography Tool with 70 nm Resolution," *Proc. SPIE* **4343**:51–59 (2001).
49. R. Soufli, R. M. Hudyma, E. Spiller, et al., "Sub-Diffraction-Limited Multilayer Coatings for the 0.3 Numerical Aperture Micro-Exposure Tool for Extreme Ultraviolet Lithography," *Appl. Opt.* **46** (18):3736–3746 (2007).
50. E. Spiller, "Multilayer Interference Coatings for the Vacuum Ultraviolet," in *Space Optics*, edited by B J Thompson and R R Shannon (National Academy of Sciences, Washington, D.C., 1974), pp. 581–597.
51. A. Khandar and P. Dhez, "Multilayer X Ray Polarizers," *Proc. SPIE* **563**:158–163 (1985).
52. E. Spiller, "The Design of Multilayer Coatings for Soft X Rays and Their Application for Imaging and Spectroscopy," in *New Techniques in X-ray and XUV Optics*, edited by B. Y. Kent and B. E. Patchett (Rutherford Appleton Lab., Chilton, U.K., 1982), pp. 50–69.
53. J. B. Kortright and J. H. Underwood, "Multilayer Optical Elements for Generation and Analysis of Circularly Polarized X Rays," *Nucl. Instrum. Meth.* **A291**:272–277 (1990).
54. J. B. Kortright, H. Kimura, V. Nikitin, et al., "Soft X-Ray (97-eV) Phase Retardation Using Transmission Multilayers," *Appl. Phys. Lett.* **60**:2963–2965 (1992).
55. J. B. Kortright, M. Rice, S. K. Kim, et al., "Optics for Element-Resolved Soft X-Ray Magneto-Optical Studies," *J. Magnetism and Magnetic Materials* **191**:79–89 (1999).
56. S. Di Fonzo, B. R. Muller, W. Jark, et al., "Multilayer Transmission Phase Shifters for the Carbon K Edge and the Water Window," *Rev. Sci. Instrum.* **66** (2):1513–1516 (1995).
57. J. B. Kortright, M. Rice, and K. D. Franck, "Tunable Multilayer EUV Soft-X-Ray Polarimeter," *Rev Sci Instr* **66** (2):1567–1569 (1995).
58. F. Schäfers, H. C. Mertins, A. Gaupp, et al., "Soft-X-Ray Polarimeter with Multilayer Optics: Complete Analysis of the Polarization State of Light," *Appl. Opt.* **38**:4074–4088 (1999).
59. N. N. Chapman, A. Barty, M. J. Bogan, et al., "Femtosecond Diffraction Imaging with a Soft-X-Ray Free-Electron Laser," *Nat. Phys.* **2**:839–843 (2006).
60. S. Bajt, H. N. Chapman, E. Spiller, et al., "A Camera for Coherent Diffractive Imaging and Holography with a Soft-X-Ray Free Electron Laser," *Appl. Opt.* **47**:1673 (2008).

NANOFOCUSING OF HARD X-RAYS WITH MULTILAYER LAUE LENSES

Albert T. Macrander,¹ Hanfei Yan,^{2,3} Hyon Chol Kang,^{4,5}
Jörg Maser,^{1,2} Chian Liu,¹ Ray Conley,^{*1,3} and
G. Brian Stephenson^{2,4}

¹*X-Ray Science Division
Argonne National Laboratory
Argonne, Illinois*

²*Center for Nanoscale Materials
Argonne National Laboratory
Argonne, Illinois*

³*National Synchrotron Light Source II
Brookhaven National Laboratory
Upton, New York*

⁴*Materials Science Division
Argonne National Laboratory
Argonne, Illinois*

⁵*Advanced Materials Engineering Department
Chosun University
Gwangju, Republic of Korea*

ABSTRACT

Multilayer Laue lenses (MLLs) have the potential to provide hard x-ray beams focused to unprecedented dimensions that approach the atomic scale. A focus of 5 nm or below is on the horizon. We review the diffraction theory as well as the experimental results that support this vision, and we present reasons to prefer harder x rays in attempting to achieve this goal.

*Ray Conley is now at National Synchrotron Light Source II, Brookhaven National Laboratory.

42.1 INTRODUCTION

Soon after he discovered x rays and explored their strong penetrating power, William Roentgen also found that they were only very weakly deflected. That is, the new rays traveled almost straight through the materials Roentgen put in front of them.¹ Small angular deflections via refraction arise from a value of the index refraction close to that of the vacuum. In the case of x rays, the index of refraction of any material is only slightly smaller than unity.² Just as for other electromagnetic radiation, focusing of x rays inherently involves deflecting the rays, that is, deflecting the Poynting vector of wavefronts. The consequence of this fundamental property of matter, namely, that the index of refraction for x rays is almost unity, implies that x rays are very difficult to focus to dimensions approaching the x-ray wavelength. (Here we consider a typical x-ray wavelength as 1.24 nm, corresponding to an x-ray energy of 1 keV, which we take as the border between soft and hard x rays.) The numerical aperture (NA) can be increased by placing many refractive lenses in series, and the net effect can approach the sum of the NAs of the individual lenses.^{3,4} However, we concentrate in this review on another means of achieving high NA by employing Bragg diffraction. For Bragg diffraction from crystals, an angle of 45° is not unusual. If a lens can be made that employs diffraction to approach this Bragg angle, it would come close to a NA of unity. A multilayer Laue lens (MLLs), shown schematically in Fig. 1, is an optic that can, in principle, achieve high diffraction angles, and correspondingly large NA, efficiently. The diffracted beam is transmitted through the lens. When crystal diffraction occurs in transmission the diffraction geometry is known as a Laue case,⁵ and, analogously, the name of Max von Laue was used to name this type of lens.

The Rayleigh criterion is well known in classical optics and determines the diffraction-limited focus of a lens. In the focal plane this limit corresponds closely to the distance d of the first minimum of the Fraunhofer diffraction pattern of the lens aperture from the optical axis,

$$d = \alpha [\lambda/\text{NA}] \quad (1)$$

where λ is the wavelength and α is a constant on the order of unity. For two-dimensional focusing by a round lens α equals 0.61, whereas for a linear (or rectangular) lens, α equals 0.5.⁶ Since hard

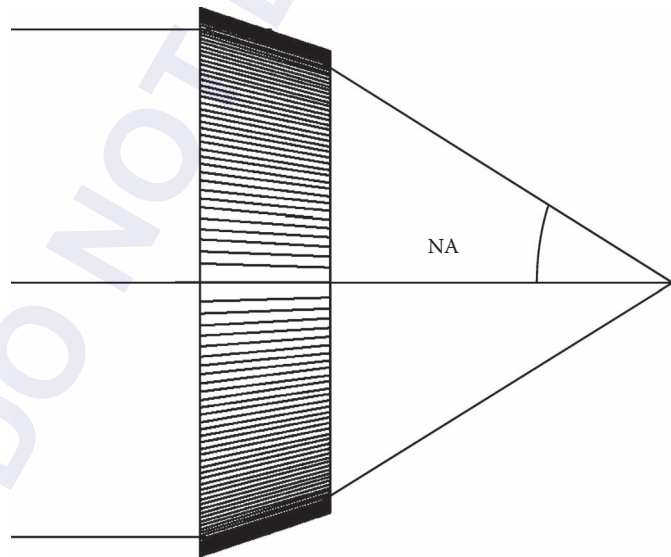


FIGURE 1 Multilayer Laue lens. A wedged version is shown. A local Bragg angle is made between incoming x rays and multilayer interfaces. The numerical aperture (NA) is shown.

x rays have wavelengths at or near atomic dimensions, focus sizes approaching atomic dimensions are thereby, in principle, feasible for values of the NA near unity.

If one can fabricate the layers in an MLL with control over the layers that is of atomic dimensions, one should be able to achieve a focus approaching atomic dimensions.⁷⁻⁹ Such control is available with several thin-film techniques, such as magnetron sputtering, atomic layer deposition,¹⁰ and many types of epitaxy.¹¹ However, there are other important factors to be considered in choosing the deposition technique, as discussed below. To date, only magnetron sputtering has been shown to allow one to usefully deposit a very large number of total zones as required for a useful linear Fresnel zone plate.

We note that phase-reversal zone plates that focus x rays have been designed since 1974,¹² and that these are also diffractive optics (see Chap. 40). In principle, zone plates made by the traditional photolithographic steps are also capable of very large NAs. However, in practice, the photolithographic process is limited to maximum aspect ratios of ~ 20 . That is, for an outermost zone width of 1 nm, the depth of the zone plate cannot be larger than ~ 20 nm. (See Fig. 2 for the definition of the dimension of a zone plate or MLL that we presently refer to as the depth.) This situation limits the efficiency very severely for hard x rays.

The efficiency of phase zone plates depends on the phase shift difference of waves transmitted through adjacent zones. If absorption losses can be ignored, the x-ray waves propagating through adjacent zones should be perfectly out of phase in order to achieve maximum efficiency.¹² The phase-shift difference increases both with an increasing index of refraction contrast between adjacent zones as well as with increasing zone plate depth. The index of refraction contrast is a function of the wavelength of the x rays and, aside from absorption edges, decreases with decreasing x-ray wavelength, that is, with increasing energy. Consequently, zone plates designed for optimum efficiency at high energies must have an increased depth compared to ones designed for optimum efficiency at low energies in order to compensate for the reduced index of refraction contrast at high energies. This is the fundamental reason that efficient focusing of hard x rays requires zone plates with very large aspect ratios. Although zone plates have been stacked on one another with some success,¹³ the MLL technology provides essentially limitless aspect ratios since practically any lens thickness can be chosen. We note that an aspect ratio of 2000 has been demonstrated recently.¹⁴

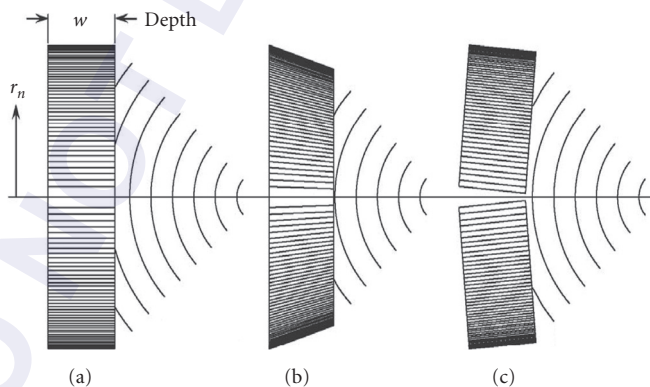


FIGURE 2 Three different types of MLLs. (a) Flat, equivalent to a linear Fresnel zone plate, in which the interfaces are parallel to an the optical axis. Here a Bragg condition is not satisfied. (b) Ideal (also called wedged), as in Fig. 1, in which each interface is angled so as to meet its own local Bragg condition. (c) Tilted, in which a flat case lens is split into two halves, each of which is tilted to an angle. This angle is set to meet a Bragg condition for one of the multilayer pair spacings within the lens. The optical depth dimension w and the “radius” to the n th zone r_n are shown. (Reprinted from Ref. 9. Copyright 2006, with permission from the American Physical Society. <http://link.aps.org/abstract/PRL/v96/e127401>.)

As first explored by Maser and Schmahl¹⁵ and as discussed below, MLLs can also be used in a novel diffractive mode by fulfilling the Bragg condition for Laue-case diffraction from layers. This innovative idea is implemented by tilting the layers to the Bragg angle, and has the consequence that efficiencies are greatly enhanced and wavefront aberrations minimized. Unlike the case for “thin” zone plates,¹² a regime of diffraction known as volume diffraction, akin to x-ray dynamical diffraction in crystals, is needed to model the performance of MLLs.¹⁵ A key result of this theoretical description is that efficiencies in the range of 60 to 70 percent can be achieved. Furthermore, wavefront aberrations can be kept to a level where the Rayleigh resolution can be achieved. For hard x rays, lens thicknesses of tens of micrometers are needed for optimum efficiency. This implies that aspect ratios of 10,000 or larger are needed for a lens capable of focusing x rays to 1 nm. MLLs offer the promise that such a focus can actually be achieved with excellent diffraction efficiency.

42.2 MLL CONCEPT AND VOLUME DIFFRACTION CALCULATIONS

MLLs are made by deposition of bilayers as shown schematically in Fig. 2. As shown in Fig. 2a, a “flat” MLL can be viewed as a linear Fresnel zone plate made by deposition of bilayers. The zone positions r_n must follow the well-known zone plate law given by¹²

$$r_n^2 = n \lambda f + n^2 \lambda^2 / 4 \quad (2)$$

Here λ is the x-ray wavelength and f is the focal length. The second term on the right can be omitted when $n\lambda \ll 4f$, which is still a good approximation for MLLs designed for angstrom wavelengths and millimeter focal lengths, which are the primary subject of this review. The bilayer pair must have significant contrast in charge density resulting in a significant phase and/or intensity difference upon x-ray wave transmission, as in the case of a thin lens,¹² as well as a significant structure factor for local diffraction as a result of satisfying the Bragg condition given by

$$\lambda = 2 (2\Delta r_n) \sin(\theta_n) \quad (3)$$

where $\Delta r_n = (r_n - r_{n-1})$ is the spacing between two successive interfaces. (We note that this equation is for first-order Bragg diffraction. For the sake of simplicity we have not denoted higher orders of diffraction.)

Arranging the diffraction geometry so that $2\Delta r_n$ acts as the Bragg spacing is an essential innovation relative to normal zone-plate usage and theory. Such a diffraction geometry is shown in Fig. 2c, which is called *tilted*. It is achieved by oppositely tilting the two halves to some angle. Normal x-ray zone plates are illuminated parallel to the axis of the zone plate. Their properties are well described by a theory that treats the waves at the exit surface, after accounting for a phase shift and absorption upon transmission through the bilayers, as sources of wavelets in a Huygens construction. These wavelets then propagate to form the focus.¹² This is known as the thin-lens case. The thin-lens treatment fails when multiple scattering of x rays takes place inside the MLL. If that occurs, dynamical diffraction approaches have to be used, and the deviations from the Bragg condition need to be considered. One must then apply a diffraction theory that can also account for multiple internal reflections in a manner akin to dynamical diffraction in crystals.⁵ This more sophisticated theory is known as volume diffraction and was first applied to MLLs by Maser and Schmahl using a coupled-wave approach.¹⁵ Results for this theory are shown in Fig. 3 for an MLL tilted so that the Bragg condition is satisfied for an intermediate zone.⁸ Here the local diffraction efficiency is plotted versus increasing zone number for MLLs having outermost zone widths of 2 and 10 nm. For the central zones of low order, the results are as expected for thin lenses and yield a local diffraction efficiency of 26 percent. However, for the zone where dynamic diffraction occurs (at 2 μm radius and above), the local diffraction efficiency is increased significantly, otherwise the diffraction efficiency becomes very small, similar to x-ray diffraction from a crystal rocked through a diffraction peak. That is, a transition in diffraction properties occurs as one approaches the dynamical diffraction regime, at which point diffraction properties become very sensitive to the Bragg condition.

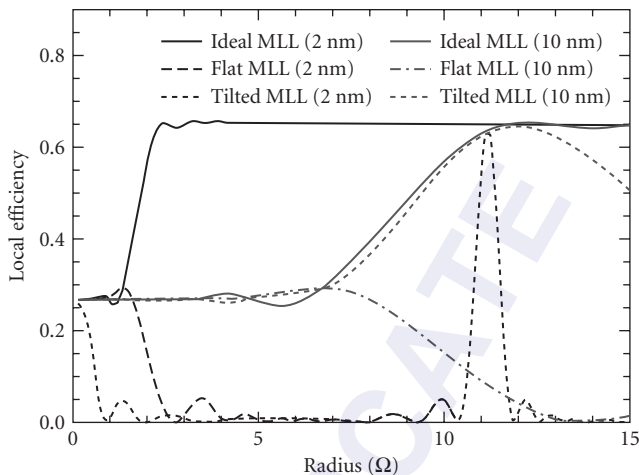


FIGURE 3 Calculated diffraction efficiency at 0.064 nm (19.5 keV) for outermost zone widths of both 10 nm (gray) and 2 nm (black) as a function of radius for ideal (wedged), flat, and tilted MLLs. For flat MLLs efficiencies do not exceed 26 percent and only very low-order zones diffract in the 2-nm case. For an outermost zone width of 10 nm, the ideal and tilted cases have almost the same performance, but for a 2-nm outermost zone, the ideal case is far superior. For the 2-nm tilted case, only a sharp Bragg peak is seen at the radius for which a Bragg condition is satisfied. This figure shows that meeting the Bragg condition everywhere becomes increasingly important for outermost zones less than 10 nm in order to ensure a high efficiency throughout most of the MLL. (See also color insert.) (Reprinted from Ref. 2. Copyright 2006, with permission from the American Physical Society. <http://link.aps.org/abstract/PRL/v96/e127401>.)

As also shown in Fig. 3, the number of zones having a high diffraction efficiency narrows as one considers smaller outermost zones. This reduces the overall efficiency. This would seem to rule out the usefulness of MLLs to structures with outermost zones of 2 nm or less. However, as depicted in Fig. 2*b*, if, instead of depositing an MLL for which all interfaces are parallel, one can build a local Bragg angle into the structure itself, a wide ranging local efficiency of 65 percent can be achieved. As discussed below, this structure, at first called “ideal” and later “wedged” (since even more ideal MLLs with curved interfaces can be considered as discussed below), has recently been demonstrated and now appears feasible.

42.3 MAGNETRON-SPUTTERED MLLS

Preliminary Study on Periodic Multilayers

It has been known for many years that magnetron sputtering can be used to make periodic multilayers useful as diffraction optics for hard x rays.^{16–18} It has only recently been discovered that layer placement can be well controlled over thousands of layers to form multilayers that are many micrometers thick.¹⁴ The ability to deposit such thick multilayers in a reasonable length of time with the requisite perfection is what sets magnetron sputtering apart from other thin-film deposition technologies such as molecular beam epitaxy.

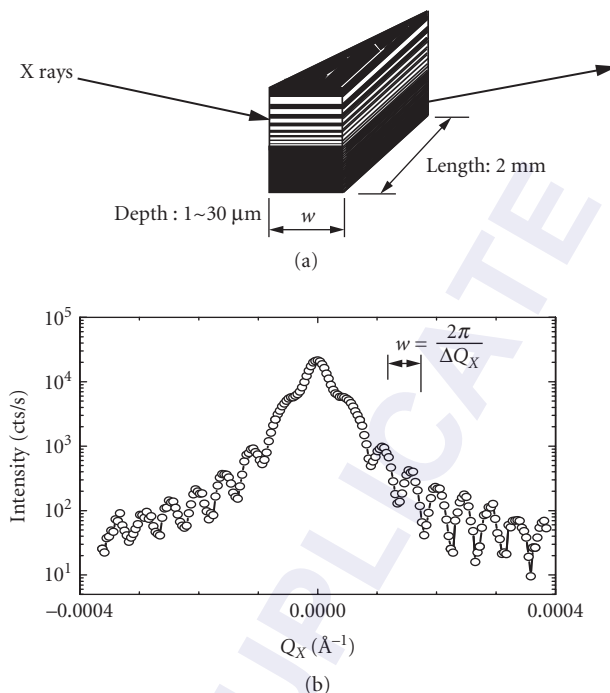


FIGURE 4 (a) Schematic illustration of an MLL made to have a tapered shape so that the optical depth w varies in a direction transverse to the optical axis. This local tapered width can be measured from fringe spacings observed in a transverse (rocking) scan, as shown in (b). (Reprinted from Ref. 25. Copyright 2007, with permission from the American Institute of Physics.)

The resultant structure must be formed into an actual lens, and this requires processing steps such as dicing, thinning, and polishing. The multilayer must be robust enough to withstand these steps without delamination, and this imposes constraints on the internal strain. That sufficiently thick multilayers could be produced and were robust enough to withstand the needed processing was first demonstrated in 2005 for W/Si as well as for Mo/Si periodic multilayers.¹⁹

Another important observation was made in the preliminary studies on thick periodic multilayers. Just as for Laue-case diffraction from thin crystals, thickness fringes occurred that permitted an accurate determination of the thickness in transmission. Far-field diffraction data showing the fringes is shown in Fig. 4.¹⁹ The diffraction efficiency depends strongly on this thickness, and its calibration is important for comparisons to theory.

WSi₂/Si Bilayers and Metrology Based on Scanning Electron Microscope Data

If thick MLLs are to have well-defined zones throughout, then interface roughness should ideally not build up progressively. Kinetic roughening of a growing surface is a complex phenomenon.²⁰ Furthermore, diffusion across the interface driven by a chemical potential gradient should be obviated as much as possible. Evidence for interdiffusion for W/Si multilayers has been reported by Windt et al.¹⁸ For these reasons, WSi₂/Si bilayers were chosen for MLLs. The choice has proven to be

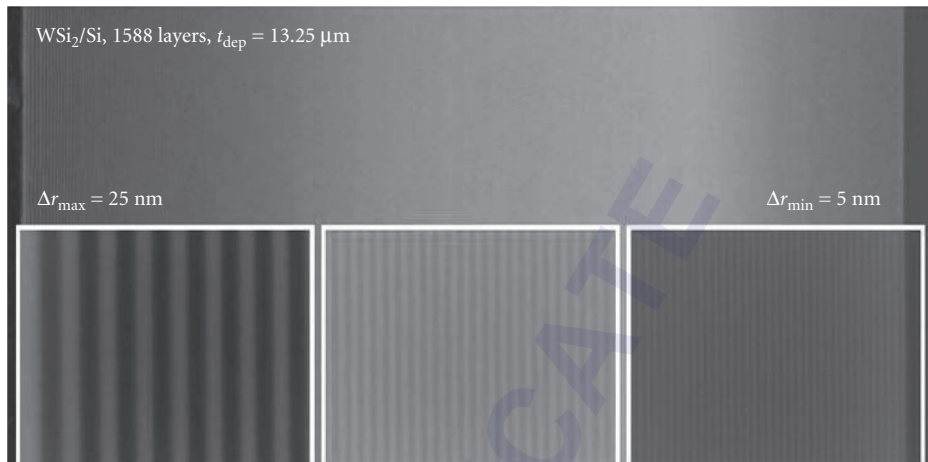


FIGURE 5 A scanning electron microscope image of a partial MLL having a 5-nm-thick outermost zone. The structure was used to obtain a focus of 16 nm for x rays of 0.064-nm (19.5-keV) wavelength. (See also color insert.) (Reprinted from Ref. 14. Copyright 2008, with permission from American Institute of Physics.)

propitious. X-ray reflectivity measurements made in situ, as a multilayer was built up, revealed that the growth of WSi_2 smoothed an underlying Si surface, and that there was not a net build up of roughness, albeit for only five periods.²¹

We note that MoSi_2/Si MLLs are also achievable. A focus of 28.5 nm has recently been reported.²² Such MLLs may be preferable at photon energies below ~ 12 keV due to the L absorption edges of W.

SEM images of cross sections of MLL wafers have proven to be very useful for analyses of the sputtering procedures. Such an SEM image is shown in Fig. 5.¹⁴

Not only must the interface roughness of each interface be well controlled, but also each layer must be grown at the Fresnel zone position with good accuracy. Ignoring the second term in Eq. (2), one finds that the width of the n th zone is given by

$$(\Delta r_n)^{-1} = (r_n - r_{n-1})^{-1} = 2r_n / \lambda f \quad (4)$$

This very useful result has been used to evaluate MLL wafers prior to making lenses and to provide feedback information for the crystal growth. Images taken by scanning electron microscopes on cross-sectioned MLL wafers were used to produce plots of Δr_n versus r_n , such as shown in Fig. 6.^{23,24} Ideally these plots should be linear, and deviations from linearity provided clues on how to adjust the sputtering process.²⁴

Lens Processing

Lenses can be successfully made starting from a deposited multilayer wafer.²⁵ The process is shown illustratively in Fig. 7. The first step is to form a “sandwich structure” by gluing a covering wafer on the multilayer surface with an epoxy. The entire structure is then diced on a high-speed saw, and the resultant bar-shaped sections are blocked with protective end caps. The bars are then bonded to a polishing fixture with the plane of the multilayers perpendicular to the plane of the fixture. The mounted bar is then thinned and polished on one side. The bar is then released, flipped over, and polished on the other side. These steps are needed because sawing damage to the multilayer cannot be avoided and damage on both sides needs to be removed. After mechanical polishing, ion beam milling is needed on each side to remove possible damage introduced by the mechanical polishing. After fine manual polishing, the final MLL is then mounted on a Mo holder and is then ready for beamline use.

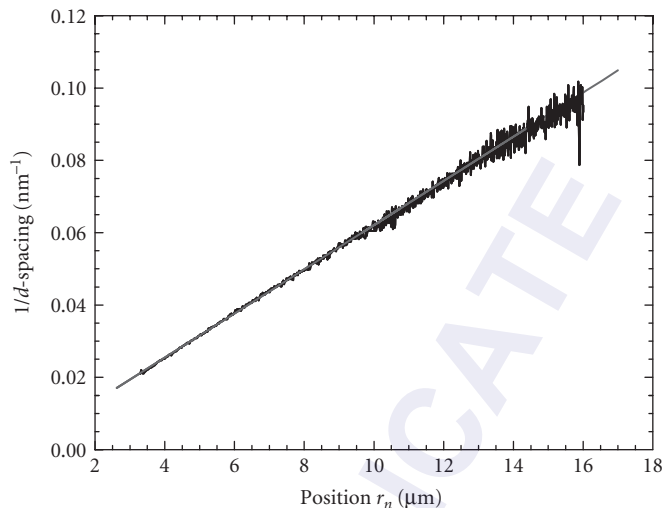


FIGURE 6 Results of analyses of electron microscope (SEM) images for layer thicknesses. The inverse of layer thicknesses when plotted as a function of MLL radius is shown in as the jagged curve. The straight line shows the desired linear functional dependence for Fresnel zones.

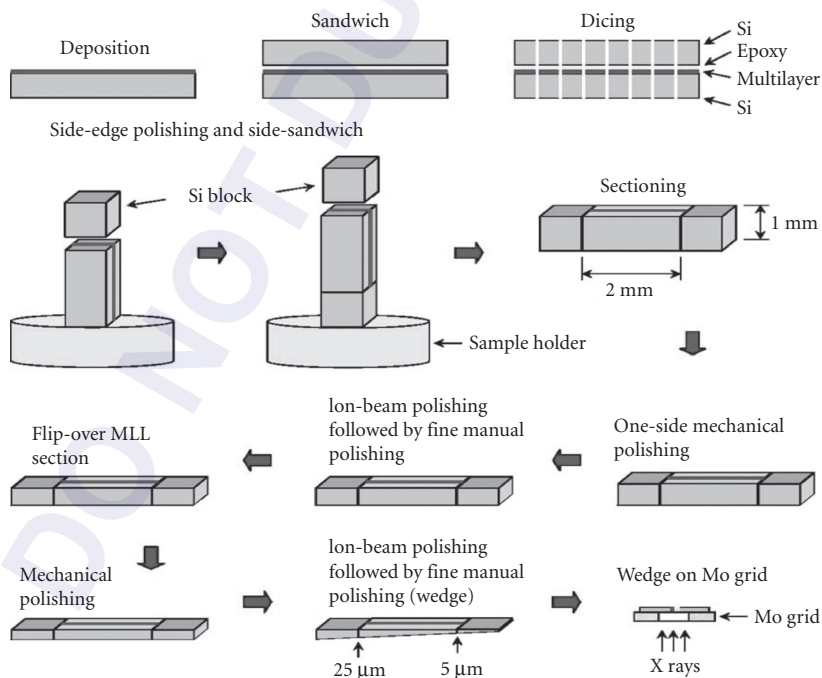


FIGURE 7 Cartoon illustrating the sequence of steps used to process a lens suitable for use with x rays, starting from an as-sputtered MLL wafer. (See also color insert.) (Reprinted from Ref. 25. Copyright 2007, with permission from the American Institute of Physics.)

42.4 INSTRUMENTAL BEAMLINE ARRANGEMENT AND MEASUREMENTS

Measurements of the focus have been made at beamlines 26 ID, 12 BM, and 8 ID at the Advanced Photon Source.^{8,14} The setup used for the measurements is shown in Fig. 8. For all the measurements made to date only a partial MLL was used. That is, an incomplete linear zone plate structure was studied. These were incomplete in two ways. First, not all zones were deposited. That is, sputter depositions were done by first sputtering an outermost zone, but the sputtering was halted before growing the inner most zones. Second, only one side above the center line in Fig. 2c was studied so that only a single tilt angle was employed. For the measurements, a horizontal diffraction plane was used. The MLL interfaces were positioned parallel to a vertical plane, and tilting was done by rotating MLLs about a vertical axis. An upstream beamline slit to aperture the horizontal source size was used to increase the horizontal coherence length, and reduction of this slit aperture proved to be needed to achieve the smallest focuses. X rays were monochromatized with a standard Si(111) double bounce arrangement. The energy was left fixed during all measurements. Data were obtained both at 19.5 and 29.5 keV.

Focus-size measurements were made by driving a specially made “analyzer” horizontally through the focused beam as shown in Fig. 8. The analyzer was a sputtered Pt thin-film processed in a manner similar to that of the MLLs themselves. In this way, a Pt thin-film segment was illuminated edge-on

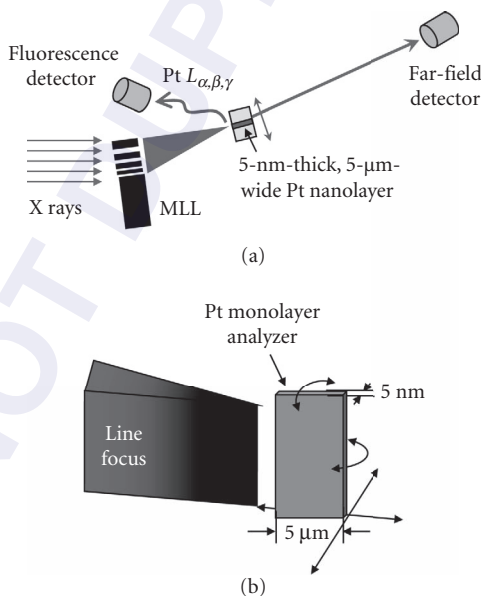


FIGURE 8 (a) Experimental layout of MLL focusing measurements. For the measurements, a specially made analyzer made from a Pt thin film deposited by sputtering was scanned through the focal plane. One detector was positioned to count Pt fluorescence photons and another was positioned to count diffracted photons in the far field. (b) Schematic illustration of the radiation pattern near the focal plane and the Pt analyzer. The two must be well aligned for correct measurement of the width of the line focus.

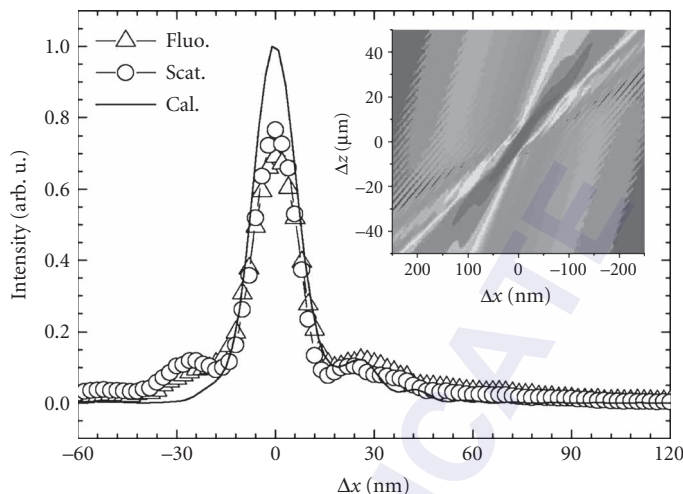


FIGURE 9 Measured and calculated intensity profiles for the focus of the MLL shown in Fig. 5. The FWHM values are 17.6 nm from fluorescence data and 15.6 nm from far-field scattering data. These should be compared to a calculated value of 15.0 nm. The calculated intensities were scaled according to measured and calculated efficiencies of 30 and 32 percent, respectively. The calculated results do not include the effects of vibrations, finite analyzer width, and finite source size. The inset shows the calculated isophote pattern near the focal plane. (See also color insert.) (Reprinted from Ref. 14. Copyright 2008, with permission from the American Institute of Physics.)

(or nearly so) by the line focus of an MLL. Angular alignment of the Pt edge to that of the MLL focal line was critical in obtaining the smallest attainable focus. As shown schematically in Fig. 8, two signals were recorded to measure the profile of a focused beam as the analyzer was scanned. A p-i-n detector was positioned at roughly 90° to record Pt fluorescence while a scintillation counter collected the signal in the far field. The far-field signal recorded scattering from the analyzer including, in particular, reflectivity from the Pt thin film. Recent results for the measurement of the focus at 19.5 keV are shown in Fig. 9.¹⁴ The far-field detector was also used to measure the efficiency.

Efficiencies of MLLs were measured by two methods. First, as shown in Fig. 10, without an analyzer present and with the detector set to observe the direct beam that passes through an MLL, one observes a dip in the transmitted intensity as the MLL is scanned through the incident beam. With the MLL tilted far from a Bragg condition, there is a dip due to absorption. But with the MLL tilted to a Bragg angle, a significantly larger dip was found. The difference between the dips is akin to extinction in crystal diffraction, since intensity is removed from the direct beam by Laue-case diffraction. Unlike for thin Fresnel lenses (or gratings), the Bragg condition depletes diffraction intensity from other orders, and, as a consequence, the extinction measurement is a reasonable measure of the overall efficiency. The second method to obtain the diffraction efficiency is based on the observation that a $\theta - 2\theta$ diffraction scan, where θ is the tilting angle, yields a measure of the local diffraction efficiency. This follows since the zone spacings satisfy the Bragg condition progressively as θ is increased so that the local efficiency is measured. An example far-field scan of the local diffraction efficiency is shown in Fig. 11. Calculated values are also shown in Fig. 11. The overall diffraction efficiency can be obtained by integration of the diffraction scan data and by normalizing to the incident beam intensity. The two methods have been found to agree within a few percent. At 19.5 keV the efficiency for the MLL shown in Fig. 5 was found to be 31 percent.⁸

Results were also obtained at 29.25 keV for the MLL shown in Fig. 5. Fluorescence data for scans through the focus revealed a FWHM of 16 nm, as shown in Fig. 12. The measured efficiency in this case was 17 percent.

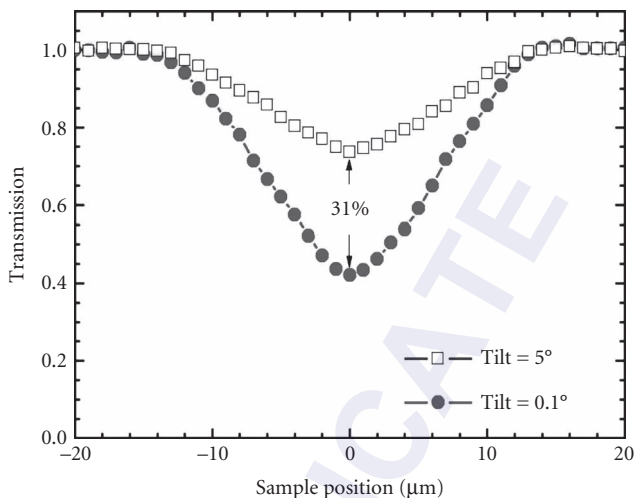


FIGURE 10 The fraction of the photons transmitted through the MLL of Fig. 5 at the Bragg condition (circles) and away from the Bragg condition (squares) as the MLL is scanned through the incident beam. The difference between the two is due to x rays that are directed to the focus only when the Bragg condition is satisfied, a phenomenon known as extinction.

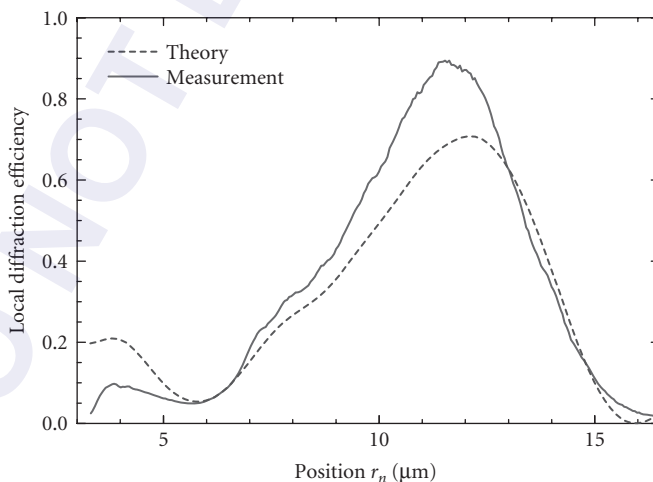


FIGURE 11 Measured and calculated focusing efficiency as a function of radius for an MLL with 5-nm-thick outermost zone width. The integrated efficiencies are 33 percent measured and 30 percent calculated. The calculations were made with coupled wave theory and show the transition from kinematic to dynamic properties at d-spacings of ~ 10 nm. The larger diffraction efficiency is most likely due to measurement uncertainties.

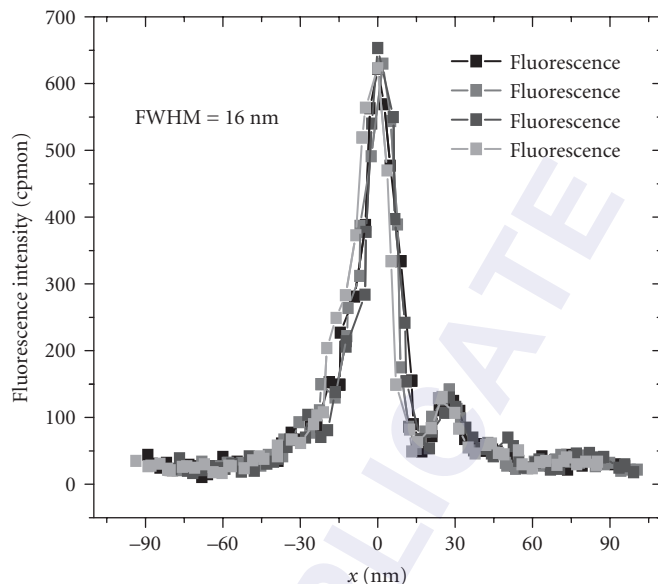


FIGURE 12 Measured line focuses of the MLL in Fig. 5 with 0.042-nm (29.25-keV) x rays. Results from four scans of the fluorescence intensity are shown. The efficiency was measured to be 17 percent. A FWHM value of 16 nm applies. (See also color insert.)

42.5 TAKAGI-TAUPIN CALCULATIONS

There does not appear to be a fundamental reason why x-ray focus sizes should not reach down to atomic dimensions. The coupled-wave-theory (CWT) used by Maser and Schmahl approximates an MLL as consisting of regions that are locally periodic and solves for local grating solutions for the diffracted waves. The orders of diffraction that result from the local gratings are then spliced to solve for the dominant diffracted waves. For values of the NA approaching unity, this approximate procedure breaks down. However, other theoretical approaches have been applied to explore a possible smallest focus on fundamental grounds. A focus size of 0.83 nm was calculated by Schroer for a wedged MLL,²⁶ and Pfeiffer et al.²⁷ report finding no lower limit for a flat MLL in a 1-to-1 imaging condition. In addition, a new diffraction theory has been developed⁹ based on the Takagi-Taupin equations for crystal diffraction,²⁸ which starts from the MLL layer structure, a priori. This new diffraction treatment proceeds by making a coordinate transformation that maps the zone plate spacings to a pseudo-Fourier series for the charge density. The electric field waves are then solved as in the case of a bulk crystal. Possible focuses below 1 nm are also predicted by the Takagi-Taupin approach.⁹

42.6 WEDGED MLLS

The new Takagi-Taupin-based theory has been applied to the case of an MLL in which each interface is tilted to a separate angle so that all zones may satisfy the Bragg condition for a divergent fan of radiation arising from a point source (see Fig. 1). When viewed internally the zones have wedge-like shapes, and this more or less “ideal” MLL structure has been called a wedged MLL (wMLL).²⁹ It was recently demonstrated that such a structure can be deposited by sputtering past a mask. The effect of the mask is to create a lateral gradient in the deposited film thickness, and, since for each

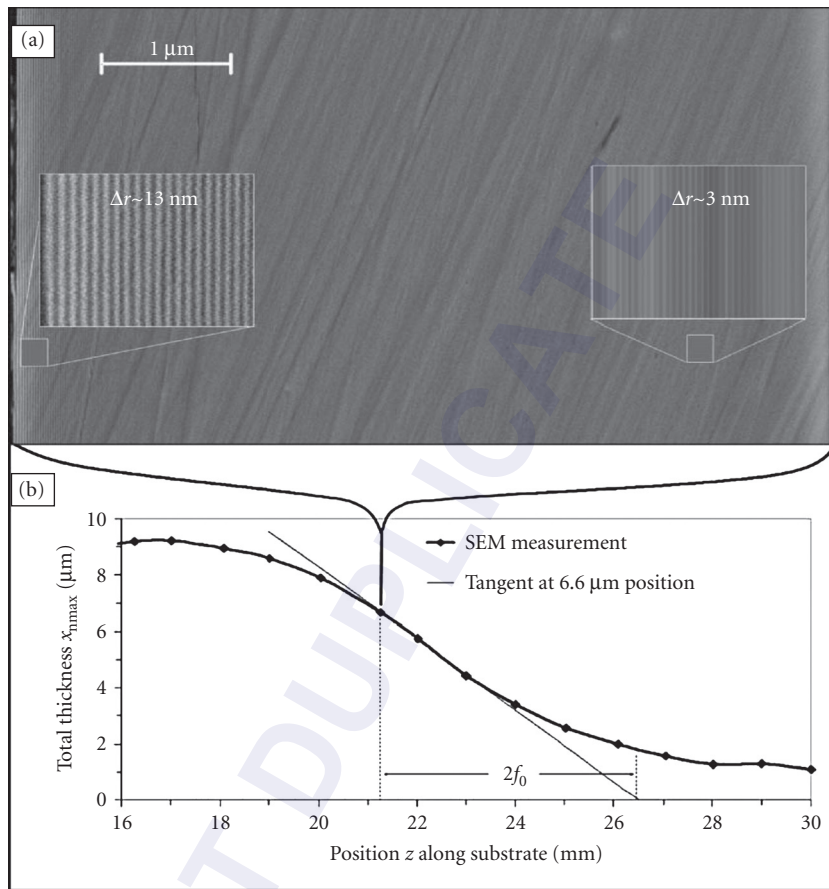


FIGURE 13 Ideal (wedged) MLL. (a) SEM image for one wafer location of a multilayer grown by sputtering past a mask over the surface of the wafer. The linear mask is designed to produce a gradient that creates a set of ideal MLL structures of varying focal lengths. The thickness of the total multilayer structure is also changed accordingly, and, as shown in (b), demonstrates the lateral gradient. (Reprinted from Ref. 29. Copyright 2008, with permission from the American Institute of Physics.)

sputtered layer the ratio of its thickness at the entrance surface to its thickness at the exit surface is the same, an entire wMLL can be built up by sputtering past a single mask. A proof-of-principle result is shown in Fig. 13. SEM data at one location are shown in Fig. 13a, and the total thickness as a function of position is shown in Fig. 13b. These data reveal that the total deposited thickness varies laterally as required. The lateral variation required is different for different focal lengths, and one such sputtered wafer can, in principle, provide MLLs for a series of focal lengths. Results of zone thickness analyses of the SEM micrograph shown in Fig. 13a are shown in Fig. 14a.

The expected results as calculated by the Takagi-Taupin approach for a wMLL having a 3-nm outermost zone width are shown in Fig. 14b. The width in the focal plane is only 2.4 nm. We note that sputtered layer thicknesses as small as 0.75 nm have been demonstrated,³⁰ and this value can be used to estimate the smallest feasible focus presently achievable with sputtering. For x rays with an energy of 19.5 keV, a wMLL having an outermost zone width of 0.75 nm is predicted to achieve a focus of 0.7 nm, a value of atomic dimensions. The efficiency of such an MLL is calculated to be 50 percent.

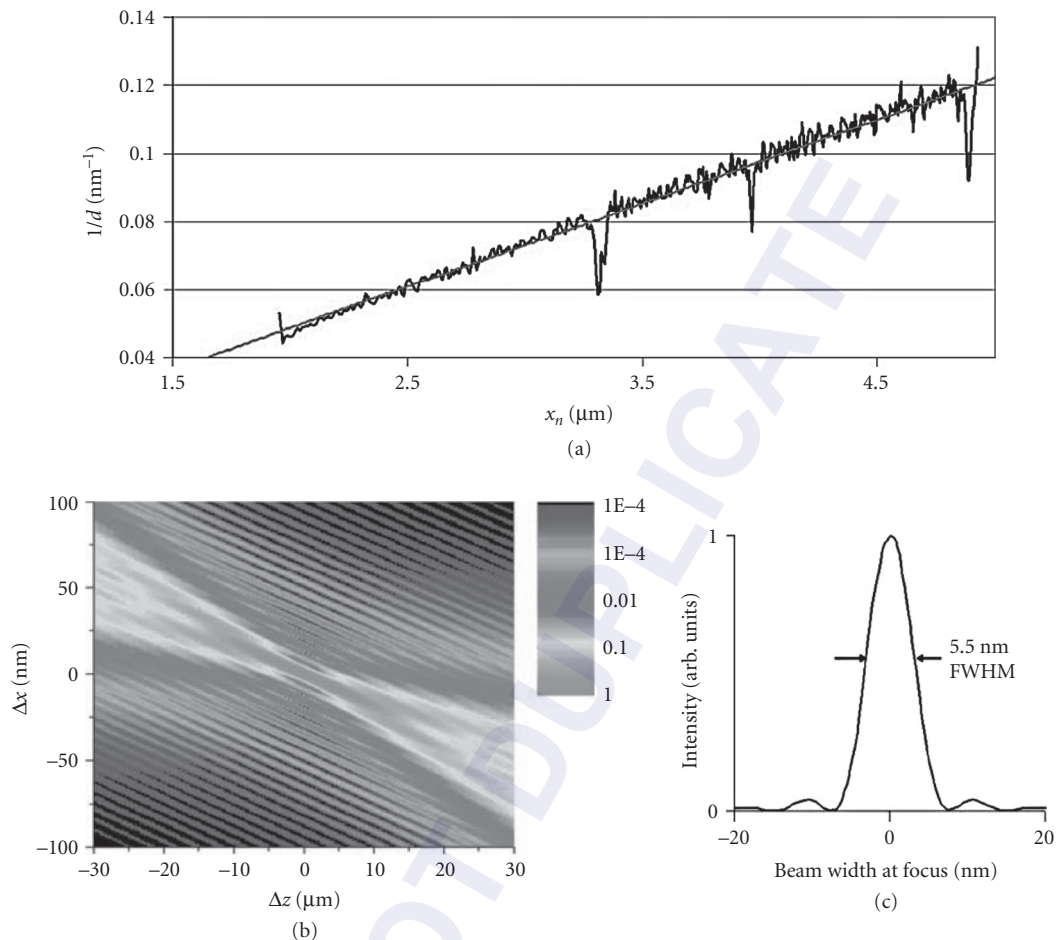


FIGURE 14 Analyses for ideal (wedged) MLL structure of Fig. 13a. The inverse layer thickness vs. radius is shown in (a). The calculated isophote pattern and intensity in the focal plane are shown in (b) and (c), respectively. (See also color insert.) (Reprinted from Ref. 29. Copyright 2008, with permission from the American Institute of Physics.)

42.7 MMLs WITH CURVED INTERFACES

The Takagi-Taupin-based calculations can be extended to curved interfaces. That curved interfaces might be more ideal is suggested by the observation that the Bragg condition cannot be met both at the entrance and at the exit surfaces of a wedged MLL. That is, the angle of a ray from a point source to any given interface at the entrance side is not the same as at the exit side. This simple observation is born out by a rigorous treatment, with the result that the surfaces must, in general, lie on concentric ellipses. For a parallel incident beam, these are concentric parabolas. If thin enough zones can somehow be deposited, the calculations predict that focuses of a few angstroms should be possible. That is, diffraction theory does not rule out true atomic dimensions for the focus of an MLL. The calculated result is shown in Fig. 15. We stress that achieving such a focus will require the actual deposition of such a structure, and at the present time this appears to be beyond the limits of magnetron sputtering.

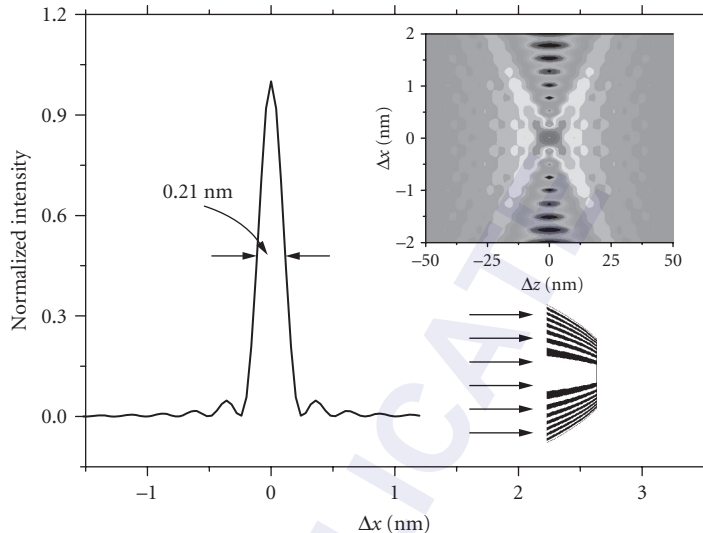


FIGURE 15 Calculated intensity in the focal plane for an MLL having parabolically shaped interfaces and an outermost zone width of 0.25 nm. The lower inset shows a cartoon of the MLL, and the upper inset shows an isophote pattern around the focus. (See also color insert.) (Reprinted from Ref. 9. Copyright 2007, with permission from the American Physical Society. <http://link.aps.org/abstract/PRB/v76/e115438>.)

42.8 MLL PROSPECTS

We consider now several practical limits that apply to the useful application of MLLs.³¹ First, useful lenses should have reasonably large focal lengths to provide space for samples and sample environments. Somewhat arbitrarily, we examine the consequences for limiting the focal length to 1 mm.

By combining the first term in Eq. (2) with Eq. (4), we obtain

$$r_N = \lambda f / (2\Delta r_N) \quad (5)$$

Here N is the outermost zone number given by

$$N = \lambda f / (2\Delta r_N)^2 \quad (6)$$

It is convenient to write Eq. (5) in terms of the total deposited thickness needed for a full linear zone plate, $D_N = 2r_N$, and the x-ray photon energy, $E = hc/\lambda$. Here $hc = 1.24$ keV-nm is the product of Planck's constant and the speed of light. The net result is

$$D_N = f hc / (2\Delta r_N E) \quad (7)$$

This equation is plotted in Fig. 16 for $f = 1$ mm and for $\Delta r_N = 1, 3, 10,$ and 30 nm. As shown in Fig. 16, for a 1-nm focus, a total of $100 \mu\text{m}$ must be deposited for an MLL that functions at 15 keV. A very important observation is that higher photon energies relax this requirement to smaller thicknesses.

The related plot of N as a function of E is shown in Fig. 17. The significance of this figure relates to chromatic blurring which occurs for an incident bandpass, $\Delta E/E$, greater than $1/N$.³² So, for example, to avoid chromatic blurring for a 1-nm focus at 15 keV, a monochromator upstream of the MLL is needed with a specification $E/\Delta E = 2 \times 10^4$.

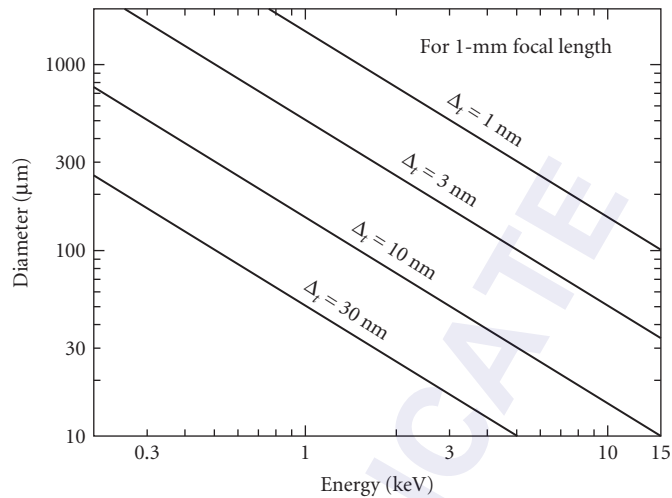


FIGURE 16 Required MLL diameter as a function of x-ray photon energy with a fixed focal length of 1 mm. Diameters are shown for focuses of 1, 3, 10, and 30 nm. For a completed MLL of a certain diameter, the focus size decreases with increasing energy.

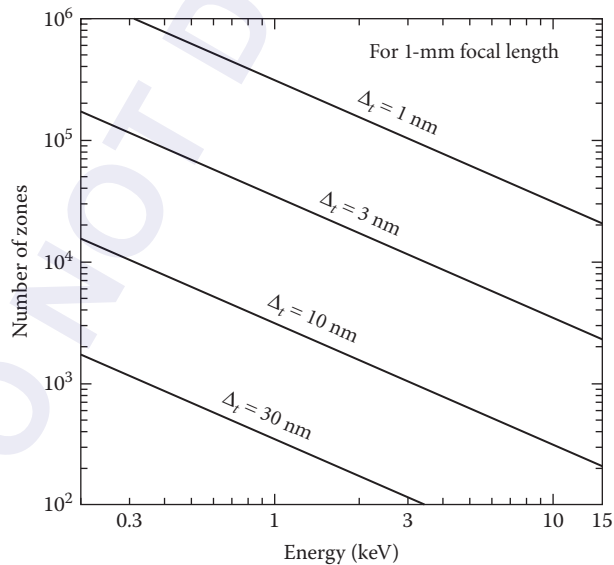


FIGURE 17 Required number for Fresnel zones as a function of x-ray photon energy with a fixed focal length of 1 mm. Diameters are shown for focuses of 1, 3, 10, and 30 nm. For a completed MLL with a certain number of zones, the focus size decreases with increasing energy.

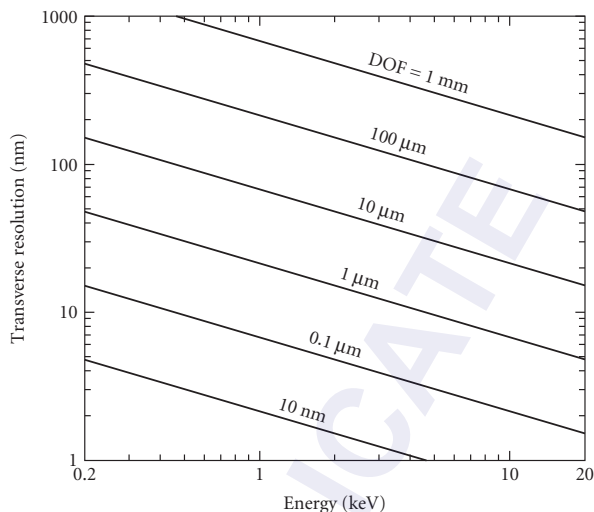


FIGURE 18 The transverse resolution (the width of the outermost zone) plotted as a function of photon energy for several depths of focus (DOF). Here the focal length has been fixed at 1 mm. For a given resolution, the DOF increases with photon energy.

Finally, we consider the depth of focus (DOF). This is also a very important consideration, since it determines the thickness of a sample that can be studied. The DOF is given by

$$\text{DOF} = \lambda/\text{NA}^2 \quad (8)$$

Here the NA is $D_N/2f$. The lateral resolution Δr_N is plotted in Fig. 18 as a function of photon energy for DOF values of 1 mm, 100 μm , 10 μm , 1 μm , 0.1 μm , and 10 nm. For a 1-nm resolution the DOF is 48 nm at 15 keV.

As shown in Figs. 16, 17, and 18, MLLs are increasingly achievable and propitious as one considers ever harder x rays. Smaller total deposited thickness, relaxed monochromaticity requirements, and larger depth of focus all apply for hard x rays compared to soft x rays.

42.9 SUMMARY

Multilayer Laue lenses have the potential to provide hard x-ray beams focused to unprecedented dimensions that approach the atomic scale. A focus of 5 nm is on the horizon. We have reviewed the diffraction theory as well as the experimental results that support this vision. MLLs are especially promising for hard x rays with energies above 20 keV, since fewer layers are required for very small focuses. A primary task yet to be accomplished is to refine the sputtering to increase the number of total layers.

42.10 ACKNOWLEDGMENTS

We are grateful for the support of the Scientific User Facilities/Advanced Photon Source (APS), Materials Science Division (MSD), and Center for Nanoscale Materials at Argonne National Laboratory. We are also grateful to the Electron Microscopy Center of MSD for the use of their facilities. Work at the APS is supported by the U.S. Department of Energy, Office of Science, Office of Basic Energy Sciences, under Contract No. DE-AC02-06 CH 11357.

42.11 REFERENCES

1. A. Stanton, "Wilhelm Conrad Röntgen on a New Kind of Rays: Translation of a Paper Read before the Würzburg Physical and Medical Society, 1895," *Nature* **53**:274–276, doi:10.1038/053274b0.
2. B. L. Henke, E. M. Gullikson, and J. C. Davis, *At. Data Nucl. Data Tables* **54**:181 (1993).
3. A. Snigirev, V. Kohn, I. Snigireva, B. Lengeler, *Nature* **384**:49 (1996).
4. C. G. Schroer and B. Lengeler, *Phys. Rev. Lett.* **94**:054802 (2005).
5. W. H. Zachariasen, *Theory of X-Ray Diffraction in Crystals*, Constable and Co., London, 1945; A. Authier, *Dynamical Theory of X-Ray Diffraction*, Oxford University Press, 2001.
6. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon, New York, 1980.
7. J. Maser, G. B. Stephenson, S. Vogt, W. Yun, A. Macrander, H. C. Kang, C. Liu, R. Conley, *Proc. SPIE* **5539**:185 (2004).
8. H. C. Kang, J. Maser, G. B. Stephenson, C. Liu, R. Conley, A. T. Macrander, and S. Vogt, *Phys. Rev. Lett.* **96**:127401 (2006).
9. H. Yan, J. Maser, A. T. Macrander, Q. Shen, S. Vogt, B. Stephenson, and H. C. Kang, *Phys. Rev. B* **76**:115438 (2007).
10. M. Ritala, K. Kukli, A. Rahtu, P. I. Raisenen, M. Leskela, T. Sajavaara, and J. Keinonen, *Science* **288**:319 (2000).
11. V. Swaminathan and A. T. Macrander, *Materials Aspects of GaAs and InP Based Structures*, Prentice Hall, Englewood Cliffs, N.J., 1991.
12. J. Kirz, *J. Opt. Soc. Am.* **64**:301 (1974).
13. J. Maser, B. Lai, W. Yun, S. D. Shastri, Z. Cai, W. Rodrigues, S. Xu, and E. Trakhtenberg, *Proc. SPIE* **4783**:74 (2002).
14. H. C. Kang, H. Yan, R. P. Winarski, M. V. Holt, J. Maser, C. Liu, R. Conley, S. Vogt, A. T. Macrander, and G. B. Stephenson, *Appl. Phys. Lett.* **92**:221114 (2008).
15. J. Maser and G. Schmahl, *Opt. Commun.* **89**:355 (1992).
16. T. W. Barbee, *Opt. Eng.* **25**:898 (1986).
17. E. Spiller, *Soft X-Ray Optics*, SPIE, Bellingham, W. A., 1994.
18. D. L. Windt, F. E. Christensen, W. W. Craig, C. Hailey, F. A. Harrison, M. Jimenez-Garate, R. Kalyanaraman, and P. H. Mao, *J. Appl. Phys.* **88**:460 (2000).
19. H. C. Kang, G. B. Stephenson, C. Liu, R. Conley, A. T. Macrander, J. Maser, S. Bajt, and H. N. Chapman, *Appl. Phys. Lett.* **86**:151109 (2005).
20. A.-L. Barabási and H. E. Stanley, *Fractal Concepts in Surface Growth*, Cambridge University Press, Cambridge, England, 1995.
21. Y.-P. Wang, H. Zhou, L. Zhou, R. L. Headrick, A. T. Macrander, and A. S. Özcan, *J. Appl. Phys.* **101**:023503 (2007).
22. T. Koyama, S. Ichimaru, T. Tsuji, H. Takano, Y. Kagoshima, T. Ohchi, H. Takenaka, *Appl. Phys. Express* **1**:117003 (2008).
23. C. Liu, R. Conley, A. T. Macrander, J. Maser, H. C. Kang, M. A. Zurbuchen, and G. B. Stephenson, *J. Appl. Phys.* **98**:113519 (2005).
24. C. Liu, R. Conley, A. T. Macrander, J. Maser, H. C. Kang, and G. B. Stephenson, *Thin Solid Films* **515**:654 (2006).
25. H. C. Kang, G. B. Stephenson, C. Liu, R. Conley, R. Khachatryan, M. Wiczorek, A. T. Macrander, H. Yan, J. Maser, J. Hiller, and R. Koratala, *Rev. Sci. Instrum.* **78**:046103 (2007).
26. C. G. Schroer, *Phys. Rev. B* **74**:033405 (2006).
27. F. Pfeiffer, C. David, J. F. van der Veen, C. Bergemann, *Phys. Rev. B* **73**:245331 (2006).
28. S. Takagi and H. H. Wills, *Acta Crystallogr.* **15**:1311 (1962); D. Taupin, *Bull. Soc. Fr. Mineral. Cristallogr.* **87**:469 (1964).
29. R. Conley, C. Liu, J. Qian, C. Kewish, A. T. Macrander, H. Yan, H. C. Kang, J. Maser, and G. B. Stephenson, *Rev. Sci. Instrum.* **79**:053104 (2008).
30. Y. Chu, C. Liu, D. Mancini, F. DeCarlo, A. T. Macrander, B. Lai, and D. Shu, *Rev. Sci. Instrum.* **73**:1485 (2002).
31. C. Jacobsen, private communication.
32. Center for X-Ray Optics and Advanced Light Source, *X-Ray Data Booklet*, <http://xdb.lbl.gov/>.

POLARIZING CRYSTAL OPTICS

Qun Shen

*National Synchrotron Light Source II
Brookhaven National Laboratory
Upton, New York*

43.1 INTRODUCTION

Being able to produce and analyze a general polarization of an electromagnetic wave has long benefited scientists and researchers in the field of visible light optics, as well as those engaged in studying the optical properties of condensed matter.¹⁻³ In the x-ray regime, however, such abilities have been very limited because of the weak interaction of x rays with materials, especially for production and analysis of circularly polarized x-ray beams. The situation has changed significantly in recent years. The growing interest in studying magnetic and anisotropic electronic materials by x-ray scattering and spectroscopic techniques has initiated many new developments in both the production and the analysis of specially polarized x rays. Routinely available, high-brightness synchrotron radiation sources can now provide naturally collimated x rays that can be easily manipulated by special x-ray optics to generate energy-tunable as well as polarization-tunable x-ray beams. The recent developments in x-ray phase retarders and multiple-Bragg-beam interference have allowed complete analyses of general elliptical polarization of x rays from special optics and special insertion devices. In this article, we will review these recent advances, especially in the area of the production and detection of circular polarization.

As for any electromagnetic wave,¹⁻³ a general x-ray beam defined by

$$E(\mathbf{r}, t) = (E_{\sigma}\hat{\sigma} + E_{\pi}e^{i\varepsilon}\hat{\pi})e^{i(\mathbf{k}\cdot\mathbf{r} - \omega t)} \quad (1)$$

can be linearly polarized if $\varepsilon = 0$ or 180° , circularly polarized if $\varepsilon = \pm 90^{\circ}$ and $E_{\sigma} = E_{\pi}$, and elliptically polarized for other values of ε , E_{σ} , and E_{π} . If $\varepsilon =$ random values, then the x-ray beam is unpolarized or

has an unpolarized component. The polarization is in general characterized by three Stokes-Poincaré parameters (P_1, P_2, P_3):

$$\begin{aligned} P_1 &= \frac{I_{0^\circ} - I_{90^\circ}}{I_{0^\circ} + I_{90^\circ}} = \frac{E_\sigma^2 - E_\pi^2}{E_\sigma^2 + E_\pi^2} \\ P_2 &= \frac{I_{45^\circ} - I_{-45^\circ}}{I_{45^\circ} + I_{-45^\circ}} = \frac{2E_\sigma E_\pi \cos \epsilon}{E_\sigma^2 + E_\pi^2} \\ P_3 &= \frac{I_+ - I_-}{I_+ + I_-} = \frac{2E_\sigma E_\pi \sin \epsilon}{E_\sigma^2 + E_\pi^2} \end{aligned} \quad (2)$$

which represent the degrees of the 0° and 90° [σ (perpendicular) and π (parallel)] linear polarization, the $\pm 45^\circ$ -tilted linear polarization, and the left- and right-handed circular polarization, respectively.¹⁻³ The unpolarized portion in the beam is characterized by its fraction P_0 of the total intensity, given by $P_0 = 1 - (P_1^2 + P_2^2 + P_3^2)^{1/2}$. The unpolarized component is generally related to the incoherency in the x-ray beam, where the phase between the σ and the π wavefields is not well-defined, and can exist only in partially coherent radiation.¹

43.2 LINEAR POLARIZERS

Today's high-brightness synchrotron sources provide natural, linearly polarized x rays on the orbital plane of the storage ring. The degree of linear polarization is usually better than from 90 to 95 percent. The actual degree of linear polarization is not important for most x-ray experiments, as long as it is known and unchanged during the course of an experiment. Certain types of synchrotron experiments, such as linear dichroism spectroscopy, magnetic scattering, and atomic and nuclear resonant scattering, do require or prefer a higher degree of linear polarization. One of the most demanding experiments is an optical activity measurement⁴ in which the plane of polarization is defined and its rotation is measured to a high precision.

Using a $2\theta = 90^\circ$ Bragg reflection (Fig. 1a) from a perfect crystal such as silicon, the natural linear polarization from a synchrotron source can be further enhanced. In general, the linear polarization P'_1 after a Bragg reflection in the vertical diffraction plane can be improved from the incoming polarization P_1 by a factor given by

$$P'_1 = \frac{(1 + P_1) - (1 - P_1)p}{(1 + P_1) + (1 - P_1)p} \quad (3)$$

where p is the polarization factor of the reflection: $p = \cos^2 2\theta$ for a kinematic and $p = |\cos 2\theta|$ for a dynamical Bragg reflection. Obviously, $P_1 \leq P'_1 \leq 1$, depending on the actual 2θ angle used. The efficiency of a 90° Bragg reflection polarizer is very high, only limited by the effective reflectivity, which is usually better than 50 percent for a kinematic crystal and at least 90 percent for a perfect dynamical crystal.

Another type of linear polarizer for x rays is based on the principle of Borrmann transmission (Fig. 1b) in dynamical diffraction.⁵ When a strong Bragg reflection [e.g., Si (220)], is excited in the Laue transmission mode, one of the σ -wave field branches has a smaller-than-usual absorption coefficient and is thus preferentially transmitted through the crystal in the forward direction. This anomalously transmitted beam is therefore highly σ -polarized, and the Laue-diffracting crystal (typically with $\mu t \sim 10$) can be used as a linear polarizer. This method has the advantage that the insertion of the polarizer does not alter the beam path, which is a common and desirable feature in visible light optics. However, because of its extremely narrow angular acceptance range (a fraction of the reflection

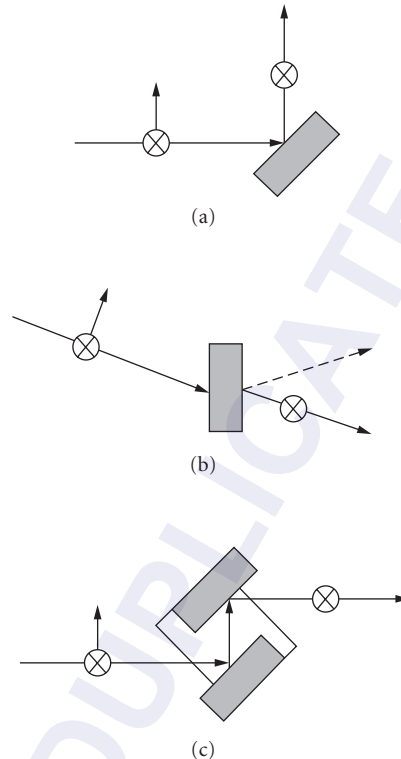


FIGURE 1 Possible linear polarizers for x rays: (a) 90° Bragg reflection; (b) Bormann transmission; and (c) multiple-bounce Bragg reflection with possible detuning.

angular width) and its intensity loss due to small but nonzero absorption, the Bormann transmission method is seldom used in practice as a linear polarizer.

The 90° Bragg reflection polarizer can be used in series to increase the purity of linear polarization even further. This can be easily achieved in practice using a channel-cut single crystal (Fig. 1c), in which the incident x-ray beam is diffracted multiple times between the two parallel crystal slabs. The highest degree of linear polarization (~ 99.99 percent) can be provided by a channel-cut crystal with a weak link between the two slabs.^{6,7} By slightly detuning the two crystal slabs, the intensity ratio of the π -polarized to the σ -polarized x rays can be suppressed by another factor of from 10 to 100. It has another advantage that the requirement of $2\theta = 90^\circ$ can be further relaxed and the polarizer is effectively tunable in a wider energy range (e.g., ± 20 percent) at 9 keV using Si (440).

Synchrotron experiments designed to study linear dichroism in materials sometimes require fast switching between two perpendicular linear polarization states in an incident x-ray beam. Although the switching can be done by a simple 90° rotation of a linear polarizer around the incident beam, a rapid switching can be achieved by a couple of methods. One is to use two simultaneously excited Bragg reflections whose scattering planes are perpendicular (or close to perpendicular) to each other. The other is to insert an x-ray half-wave phase retarder based on the effect of diffractive birefringence in dynamical diffraction from perfect crystals.⁸ In the latter method, if one uses a transmission-type phase retarder, the incident beam direction and position are not affected by the polarization switching, which is a significant advantage.

43.3 LINEAR POLARIZATION ANALYZERS

In general, any linear polarizer can be used as an analyzer for linear polarizations P_1 and P_2 . The simplest form of linear analyzer is to measure the 90° elastic scattering from an amorphous target⁹ or a powder sample.¹⁰ For better precisions, Bragg reflections with $2\theta = 90^\circ$ from single crystals are usually used. Again, one can achieve very high precision by means of a multiply bounced Bragg reflection from a channel-cut perfect crystal. The requirement for $2\theta = 90^\circ$ can be relaxed if one takes into account a scale factor involving the actual 2θ angle used in the measurement, as given in the following.

All these methods are based on measuring the intensity variation I_b as a function of the rotation angle χ of the scattering plane around the incident beam (Fig. 2) with respect to a reference polarization direction (e.g., the σ direction):

$$I_b(\chi) = \frac{1}{2}[(1+p) + (P_1 \cos 2\chi + P_2 \sin 2\chi)(1-p)] \quad (4)$$

where p is again the polarization factor of the scattering process as defined in Eq. (3). For $2\theta = 90^\circ$, $p = 0$.

It is straightforward to show that for an arbitrary $2\theta \neq 0$, the incident linear polarizations P_1 and P_2 can be obtained by measuring two difference-over-sum ratios,

$$P_1(2\theta) = \frac{I_b(0^\circ) - I_b(90^\circ)}{I_b(0^\circ) + I_b(90^\circ)} \quad (5)$$

$$P_2(2\theta) = \frac{I_b(45^\circ) - I_b(-45^\circ)}{I_b(45^\circ) + I_b(-45^\circ)}$$

and using a proper scale factor involving the actual 2θ ,

$$\begin{bmatrix} P_1 \\ P_2 \end{bmatrix} = \frac{1+p}{1-p} \begin{bmatrix} P_1(2\theta) \\ P_2(2\theta) \end{bmatrix} \quad (6)$$

Equations (5) and (6) are useful when a single Bragg reflection is used for linear polarization analyses over a wide energy range.^{11,12}

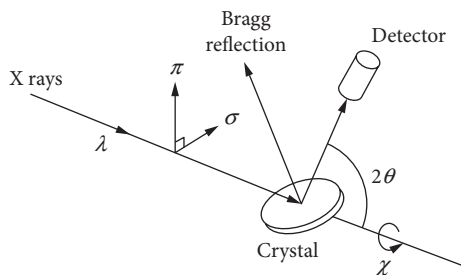


FIGURE 2 Bragg reflection from a crystal serves as a linear polarization analyzer.

43.4 PHASE PLATES FOR CIRCULAR POLARIZATION

In principle, circular polarization can be produced from linearly polarized radiation by one of the following three effects (Fig. 3): circular absorption dichroism, linear birefringence or double refraction, and reflection birefringence analogous to Fresnel rhomb. The last two effects only work with linearly polarized radiation, while the first effect can be used to convert both linear and unpolarized radiation.

Circular absorption dichroism is an effect of differential absorption between the left- and right-handed circular polarizations. For x rays, the natural circular dichroism in materials is a very weak effect. To date the only x-ray circular dichroism that has been studied extensively is the magnetic

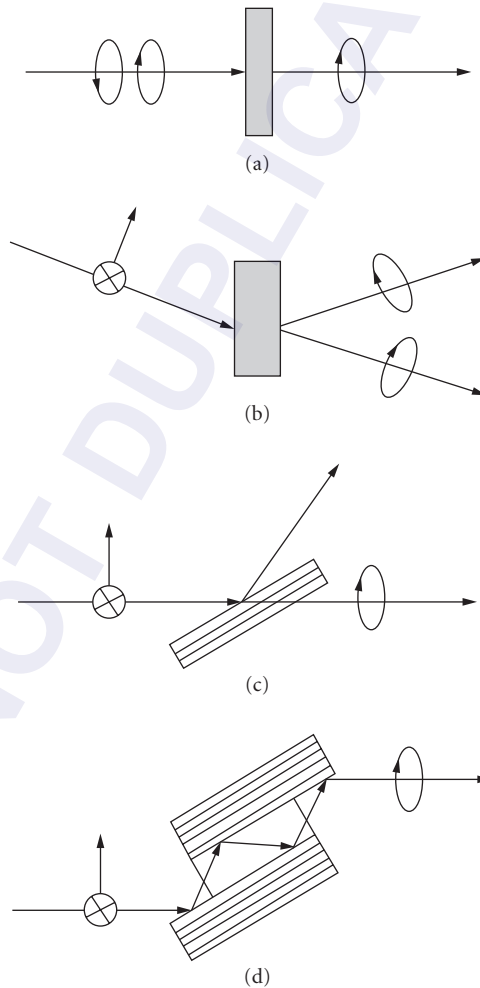


FIGURE 3 Circular phase plates for x rays based on: (a) circular absorption dichroism; (b) linear birefringence in Laue geometry; (c) linear birefringence in Bragg transmission geometry; and (d) Fresnel rhomb in Bragg reflection geometry.

circular dichroism in magnetic materials. However, the largest difference experimentally observed in absorption coefficients between the left- and the right-handed circular polarization is on the order of 1 percent, which is too small for any practical use as an x-ray circular polarizer.¹³

Fresnel rhomb is based on the principle that the reflected waves from the surface of a material can have different phase shifts between two orthogonal linear states, σ and π . The x-ray analog of this occurs at a Bragg reflection from a crystal and arises from the difference in the intrinsic angular widths (Darwin widths) of the σ and the π diffracted waves.^{14,15} The width for the π polarization is smaller than that for the σ polarization. Because of a continuous phase shift of 180° from one side of the reflection width to the other, a difference in the phase shift can be obtained between the σ and the π polarizations if the angular position is selected to be near the edge of the angular width. Experimentally, a multiply bounced Bragg reflection is needed to make a $\pm 90^\circ$ phase shift. The phase shift is independent of the crystal thickness, but this method requires a highly collimated incident beam, about 1/10 of the reflection width, which is usually the limiting factor for its throughput.¹⁶

The linear birefringence or double refraction effect relies on the difference in the magnitudes of the wavevectors of two orthogonal linear polarization states, σ and π , when a plane wave travels through a crystalline material. Because of this difference, a phase shift between the σ and the π traveling waves can be accumulated through the thickness of the birefringent material.¹⁷

$$\Delta = 2\pi(K_\sigma - K_\pi)t \quad (7)$$

where t is the thickness, and K_σ and K_π are the magnitudes of the wavevectors inside the crystal for the σ and π wavefields, respectively. When the phase shift Δ reaches $\pm 90^\circ$, circularly polarized radiation is generated, and such a device is usually termed a *quarter-wave phase plate* or a *quarter-wave phase retarder*.

For x rays, large birefringence effects exist near strong Bragg reflections in relatively perfect crystals. Depending on the diffraction geometry, one can have three types of transmission birefringence phase retarders: Laue transmission, Laue reflection, and Bragg transmission, as illustrated in Fig. 3*b* and *c*. The Laue reflection-type^{9,18} works at full excitation of a Bragg reflection, while the Laue and the Bragg transmission types^{19–22} work at the tails of a Bragg reflection, which has the advantage of a relaxed angular acceptance. In the past few years, it has been demonstrated that the Bragg transmission-type phase retarders are very practical x-ray circular polarizers. With good-quality diamond single crystals, such circular phase-retarders can tolerate a larger angular divergence and their throughputs can be as high as 0.25, with a degree of circular polarization in the range of from 95 to 99 percent. The handedness of the circular polarization can be switched easily by setting the diffracting crystal to either side of the Bragg reflection rocking curve. There have been some excellent review articles^{20–22} in this area and the reader is referred to them for more details.

43.5 CIRCULAR POLARIZATION ANALYZERS

Circular polarization P_3 of an x-ray beam can be qualitatively detected by magnetic Compton scattering¹⁸ and by magnetic circular dichroism. In general, these techniques are not suitable for quantitative polarization determination because these effects are relatively new and because of the uncertainties in the materials themselves and in the theories describing the effects.

Two methods have been developed in the past few years for quantitative measurements of circular polarization in the x-ray regime. One is to use a quarter-wave phase plate to turn the circular polarization into linear polarization which can then be measured using a linear polarization analyzer (Fig. 4*a*), as described in the previous sections. This method is entirely analogous to the similar techniques in visible optics and has been used by a couple of groups to characterize special insertion devices and quarter-wave phase plates.^{20,23}

Multiple-beam Bragg diffraction (MBD) is the other way to measure the degree of circular polarization of x rays (Fig. 4*b*). This technique makes use of the phase shift δ between the σ and the π wave

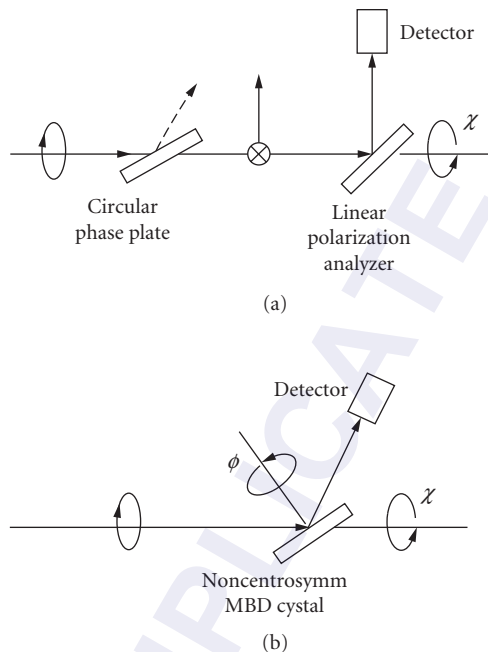


FIGURE 4 Two types of optics for analyzing circularly polarized x rays: (a) circular phase plate plus a linear polarization analyzer and (b) multiple-beam Bragg diffraction from a noncentrosymmetric crystal.

fields that arises from the phase-sensitive interference²⁴ and possible polarization mixing in an MBD process in a single crystal.²⁵ The phase shift δ is insensitive to crystal perfection and thickness as well as x-ray wavelength, since it is strictly determined by the crystal structure factors. Thus the MBD method has a broad applicable energy range and a good tolerance in the angular divergence of an x-ray beam.

Multiple-beam Bragg diffraction is also called simultaneous Bragg reflection, or *Umweganregung* (detour), as first noted by Renninger.²⁶ It occurs when two or more sets of atomic planes satisfy the Bragg's law simultaneously inside a crystal. A convenient way (Fig. 5) of realizing a multiple-beam

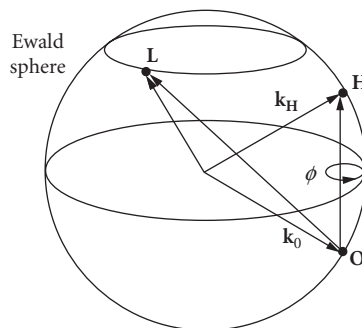


FIGURE 5 Ewald sphere construction showing a multiple-beam diffraction geometry.

reflection condition is exciting first one Bragg reflection, say \mathbf{H} , then rotating the diffracting crystal around \mathbf{H} to bring another reciprocal node, say \mathbf{L} , onto the Ewald sphere.²⁷ The rotation around the \mathbf{H} is represented by the azimuthal angle ϕ . It has been known^{28–30} that the simultaneously excited Bragg waves interfere with each other and give rise to a diffracted intensity that is sensitive to the *structural phase triplet* $\delta = \arg(F_{\mathbf{L}}F_{\mathbf{H-L}}/F_{\mathbf{H}})$, where $F_{\mathbf{H}}$ is the structure factor of the main reflection \mathbf{H} and $F_{\mathbf{L}}$ and $F_{\mathbf{H-L}}$ are the structure factors of the detoured reflections \mathbf{L} and $\mathbf{H-L}$. Another effect that exists in MBD is polarization mixing due to the intrinsic double Thomson scattering process,^{31,32} which causes a part of the π wave field amplitude to be scattered into the σ channel and vice versa. The combination of these two effects makes the MBD intensity sensitive to the phase difference between the σ and the π wave fields and thus to the circular polarization in the incident x-ray beam.

For maximum sensitivity to circular polarizations, a noncentrosymmetric crystal structure can be chosen as the analyzer and an MBD reflection with $\delta = \pm 90^\circ$ can be used. A complete determination of all Stokes-Poincaré parameters can be obtained either by measuring *three* independent MBD profiles^{25,33} and solving for (P_1, P_2, P_3) , or by separate measurements^{1,12,34} of P_1 and P_2 . This latter method provides a very convenient way for complete determination of Stokes-Poincaré parameters. Once the analyzer crystal is oriented and the proper multiple reflection is found, a complete determination of (P_1, P_2, P_3) , requires four rocking curve measurements of the integrated 2-beam intensities, $I_b(0^\circ)$, $I_b(90^\circ)$, $I_b(-45^\circ)$, and $I_b(45^\circ)$, at $\chi = 0^\circ, 90^\circ, -45^\circ$, and 45° , respectively, and two rocking curve measurements of the MBD intensities, $I(+\Delta\phi)$ and $I(-\Delta\phi)$, at each side of the multiple reflection peak. The (P_1, P_2, P_3) are then obtained by taking the three pairs of difference-over-sum ratios,¹² with a scale factor for P_3 involving the multiple-beam structure factors.

We would like to point out that for an MBD reflection to have maximum sensitivity to circular polarization, $\delta = \pm 90^\circ$ is desired. This condition can be conveniently satisfied by choosing a noncentrosymmetric crystal such as GaAs. However, as we have previously shown,³² because of the *dynamical* phase shift *within* the full excitation of a multiple reflection, an MBD in a centrosymmetric crystal such as Si or Ge can also be sensitive to circular polarization. The only drawback is that it requires an extremely high angular collimation. With the arrival of low-emittance undulator-type sources, one should keep in mind that the use of Si or Ge crystals for MBD circular polarimetry may be feasible for well-collimated x-ray beams.

43.6 ACKNOWLEDGMENTS

The author is grateful to B. W. Batterman, K. D. Finkelstein, S. Shastri, and D. Walko for many useful discussions and collaborations. This work is supported by the National Science Foundation through CHESS under Grant No. DMR-97-13424.

43.7 REFERENCES

1. J. D. Jackson, *Classical Electrodynamics*, 2nd ed., John Wiley and Sons, New York, 1975.
2. D. J. Griffiths, *Introduction to Electrodynamics*, Prentice-Hall, Englewood Cliffs, New Jersey, 1975.
3. M. Born and E. Wolf, *Principles of Optics*, 6th ed., Pergamon, New York, 1983.
4. M. Hart and A. R. D. Rodrigues, "Optical Activity and the Faraday Effect at X-Ray Frequencies," *Phil. Mag.* **B43**:321 (1981).
5. B. W. Batterman and H. Cole, "Dynamical Diffraction of X Rays by Perfect Crystals," *Rev. Mod. Phys.* **36**:681 (1964).
6. M. Hart, "X-Ray Polarization Phenomena," *Phil. Mag.* **B38**:41 (1978).
7. D. P. Siddons, M. Hart, Y. Amemiya, and J. B. Hastings, "X-Ray Optical Activity and the Faraday Effect in Cobalt and Its Compounds," *Phys. Rev. Lett.* **64**:1967 (1990).
8. C. Giles, C. Malgange, J. Goulon, C. Vettier, F. de Bergivin, A. Freund, P. Elleaume, E. Dartyge, A. Fontaine, C. Giorgetti, and S. Pizzini, "X-Ray Phase Plate for Energy Dispersive and Monochromatic Experiments," *SPIE Proceedings* **2010**:136 (1993).

9. J. A. Golovchenko, B. M. Kincaid, R. A. Levesque, A. E. Meixner, and D. R. Kaplan, "Polarization Pendellosung and the Generation of Circularly Polarized X Rays with a Quarter-Wave Plate," *Phys. Rev. Lett.* **57**:202 (1986).
10. G. Materlik and P. Suortti, "Measurements of the Polarization of X Rays from a Synchrotron Source," *J. Appl. Cryst.* **17**:7 (1984).
11. Q. Shen and K. D. Finkelstein, "A Complete Characterization of X-Ray Polarization State by Combination of Single and Multiple Bragg Reflections," *Rev. Sci. Instrum.* **64**:3451 (1993).
12. Q. Shen, S. Shastri, and K. D. Finkelstein, "Stokes Polarimetry for X Rays Using Multiple-Beam Diffraction," *Rev. Sci. Instrum.* **66**:1610 (1995).
13. K. D. Finkelstein and Q. Shen, "Feasibility of Using Magnetic Circular Dichroism to Produce Circularly Polarized X Rays," unpublished.
14. O. Brummer, Ch. Eisenschmidt, and H. Hoche, "Polarization Phenomena of X Rays in the Bragg Case," *Acta Crystallogr.* **A40**:394 (1984).
15. B. W. Batterman, "X-Ray Phase Plate," *Phys. Rev.* **B45**:12677 (1992).
16. S. Shastri, K. D. Finkelstein, Q. Shen, B. W. Batterman, and D. Walko, "Undulator Test of a Bragg Reflection Elliptical Polarizer at ~ 7.1 keV," *Rev. Sci. Instrum.* **66**:1581 (1994).
17. V. Belyakov and V. Dmitrienko, "Polarization Phenomena in X-Ray Optics," *Sov. Phys. Uspek.* **32**:697 (1989).
18. D. M. Mills, "Phase-Plate Performance for the Production of Circularly Polarized X Rays," *Nucl. Instrum. Meth.* **A266**:531 (1988).
19. K. Hirano, K. Izumi, T. Ishikawa, S. Annaka, and S. Kikuta, "An X-Ray Phase Plate Using Bragg Case Diffraction," *Jpn. J. Appl. Phys.* **30**:L407 (1991).
20. C. Giles, C. Malgrange, J. Goulon, F. de Bergivin, C. Vettier, A. Fontaine, E. Dartyge, S. Pizzini, F. Baudelet, and A. Freund, "Perfect Crystal and Mosaic Crystal Quarter-Wave Plates for Circular Magnetic X-Ray Dichroism Experiments," *Rev. Sci. Instrum.* **66**:1549 (1995).
21. K. Hirano, T. Ishikawa, and S. Kikuta, "Development and Application of X-Ray Phase Retarders," *Rev. Sci. Instrum.* **66**:1604 (1995).
22. Y. Hasegawa, Y. Ueji, K. Okitsu, J. M. Ablett, D. P. Siddons, and Y. Amemiya, "Tunable X-Ray Polarization Reflector with Perfect Crystals," *Acta Cryst.* **A55**:955–962 (1999).
23. T. Ishikawa, K. Hirano, and S. Kikuta, "Complete Determination of Polarization State in the Hard X-Ray Region," *J. Appl. Cryst.* **24**:982 (1991).
24. Q. Shen and K. D. Finkelstein, "Solving the Phase Problem with Multiple-Beam Diffraction and Elliptically Polarized X Rays," *Phys. Rev. Lett.* **45**:5075 (1990).
25. Q. Shen and K. D. Finkelstein, "Complete Determination of X-Ray Polarization Using Multiple-Beam Bragg Diffraction," *Phys. Rev.* **B45**:5075 (1992).
26. M. Renninger, "Umweganregung, eine bisher unbeachtete Wechselwirkungserscheinung bei Raumgitterinterferenzen," *Z. Phys.* **106**:141 (1937).
27. H. Cole, F. W. Chambers, and H. M. Dunn, "Simultaneous Diffraction: Indexing Umweganregung Peaks in Simple Cases," *Acta Crystallogr.* **15**:138 (1962).
28. M. Hart and A. R. Lang, *Phys. Rev. Lett.* **7**:120–121 (1961).
29. Q. Shen, "A New Approach to Multi-Beam X-Ray Diffraction Using Perturbation Theory of Scattering," *Acta Crystallogr.* **A42**:525 (1986).
30. Q. Shen, "Solving the Phase Problem Using Reference-Beam X-Ray Diffraction," *Phys. Rev. Lett.* **80**:3268–3271 (1998).
31. Q. Shen, "Polarization State Mixing in Multiple-Beam Diffraction and Its Application to Solving the Phase Problem," *SPIE Proceedings* **1550**:27 (1991).
32. Q. Shen, "Effects of a General X-Ray Polarization in Multiple-Beam Bragg Diffraction," *Acta Crystallogr.* **A49**:605 (1993).
33. K. Hirano, T. Mori, A. Iida, R. Colella, S. Sasaki, and Q. Shen, "Determination of the Stokes-Poincaré Parameters for a Synchrotron X-Ray Beam by Multiple Bragg Scattering," *Jpn. J. Appl. Phys.* **35**:5550–5552 (1996).
34. J. C. Lang and G. Srajer, "Bragg Transmission Phase Plates for the Production of Circularly Polarized X Rays," *Rev. Sci. Instrum.* **66**:1540–1542 (1995).

This page intentionally left blank.

DO NOT DUPLICATE

SUBPART

5.3

REFLECTIVE OPTICS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

IMAGE FORMATION WITH GRAZING INCIDENCE OPTICS

James E. Harvey

*CREOL, The College of Optics and Photonics
University of Central Florida
Orlando, Florida*

44.1 GLOSSARY

ACV	autocovariance function
AUDPSF	average unregistered detected point spread function
AXAF	Advanced X-Ray Astrophysical Facility
DPSF	<i>detected</i> point spread function
ESA	European Space Agency
EUV	extreme ultraviolet
EUVE	extreme ultraviolet explorer
GOES	Geostationary Orbiting Environmental Satellite
GPSF	geometrical point spread function
HH	hyperboloid-hyperboloid
HPR	half power radius
MTF	modulation transfer function
OFOV	operational field-of-view
PSF	point spread function
RDPSF	registered detected point spread function
ROSAT	Rontgensatellit
SAX	Italian x-ray astronomy satellite
SSPSF	surface scatter point spread function
SXI	solar x-ray imager
WS	Wolter-Schwarzschild

44.2 INTRODUCTION TO X-RAY MIRRORS

Conventional mirrors have traditionally exhibited useful reflectances at x-ray wavelengths only for grazing incidence angles. The first two-dimensional image produced by deflecting x rays in a controlled manner was obtained by Kirkpatrick and Baez in 1948 with two grazing incidence mirrors as

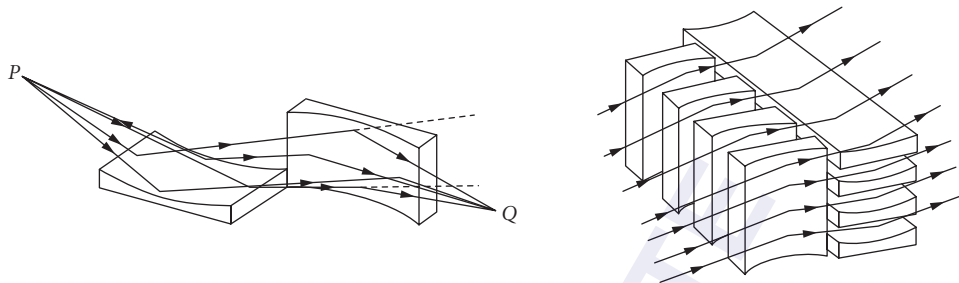


FIGURE 1 (a) The Kirkpatrick-Baez telescope consists of two orthogonal grazing incidence parabolic sheet mirrors and (b) a multiple stack of several mirrors can be used to substantially increase the collecting area.

illustrated in Fig. 1a.¹ The extremely small collecting area of such an imaging system can be alleviated by constructing the multiplate Kirkpatrick-Baez mirror of Fig. 1b.

In 1952 Hans Wolter published a paper in which he discussed several rotationally symmetric grazing incidence x-ray telescope systems.² The Wolter Type I telescope consists of a coaxial paraboloid (primary mirror) and hyperboloid (secondary mirror) as illustrated in Fig. 2a. The focus of the paraboloid is coincident with the rear focus of the hyperboloid and the reflection occurs on the inside of both mirrors. The Wolter Type II telescope also consists of a coaxial paraboloid and hyperboloid. However, the focus of the paraboloid is coincident with the front focus of the hyperboloid and the reflection occurs on the inside of the paraboloid and the outside of the hyperboloid. This system, illustrated in Fig. 2b, is the grazing incidence analog to the classical Cassegrain telescope. The Wolter Type III telescope consisting of a grazing incidence paraboloid and ellipsoid is illustrated in Fig. 2c. The focus of the paraboloid is coincident with the front focus of the ellipsoid, and the reflection occurs on the outside of the paraboloid and the inside of the ellipsoid.

The Wolter Type I telescope typically has a grazing angle of less than a degree and is used for hard x rays (greater than 1 keV). The Wolter Type II telescope typically has a grazing angle of approximately 10 degrees and is used for soft x rays and the extreme ultraviolet (EUV). The Wolter Type III optical design is not practical for astronomical telescopes; however, a finite-conjugate version of it is frequently used for x-ray microscopes.³ (See Chap. 51.)

The grazing incidence optical configurations are all free of spherical aberration; however, they exhibit severe field curvature, coma, and astigmatism. They are also cumbersome and difficult to fabricate and align, and scattering effects from imperfectly polished surfaces severely degrade image quality for these very short wavelengths. Aschenbach has written a nice review of these scattering effects in x-ray telescopes.⁴

Primarily due to much improved optical surface metrology capabilities, the conventional optical fabrication techniques of grinding and polishing glass substrates are resulting in major advances in the resolution of grazing incidence x-ray telescopes.^{5,6} The European ROSAT (Rontgensatellit) telescope⁷ consisting of four nested Wolter Type I grazing incidence telescopes provided substantial improvement in both effective collecting area and resolution over the Einstein Observatory which was launched in 1978,^{8,9} and NASA's Chandra Observatory, called the Advanced X-Ray Astrophysical Facility (AXAF) during the years it was being designed and fabricated, has demonstrated that the technology exists to produce large Wolter Type I grazing incidence x-ray telescopes with subarcsecond resolution. This progress in grazing incidence x-ray optics performance is illustrated in Fig. 3.¹⁰

The very smooth surfaces required of high-resolution x-ray optics have been achieved by tedious and time-consuming optical polishing efforts of skilled opticians. AXAF and ROSAT were thus very expensive to produce.

When such high resolution is not required, other optical materials and fabrication techniques may be applicable. The extreme ultraviolet explorer (EUVE) was a NASA-funded astronomy mission intended to perform an all-sky survey in the 70 to 760 Å spectral region. The deep survey and spectroscopic portion of the mission utilized a Wolter Type II grazing incidence telescope built at the Space Sciences Laboratory

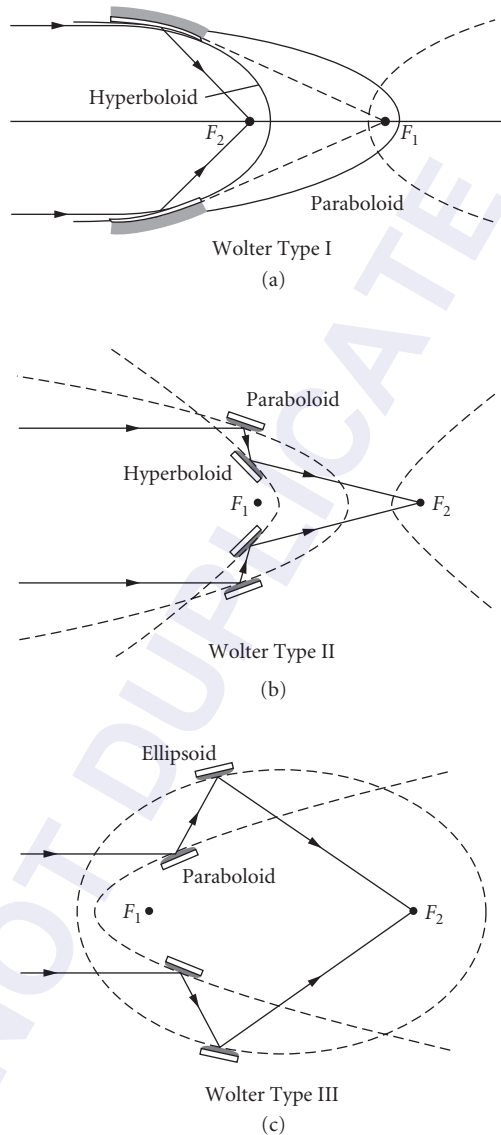


FIGURE 2 Hans Wolter's three classical, confocal, cylindrically symmetric grazing incidence x-ray telescope designs.

at the University of California, Berkeley. Its mirrors were fabricated from aluminum substrates by forging, rough machining, diamond turning to the desired figure, nickel plating, polishing, and coating with gold. An image half-power-width of approximately 1.5 arcsec was achieved.¹¹

Smooth x-ray mirror surfaces have been achieved without any labor-intensive polishing by merely dipping diamond-turned metal substrates in lacquer, then depositing tungsten or gold coatings to yield the desired high reflectance.^{12,13}

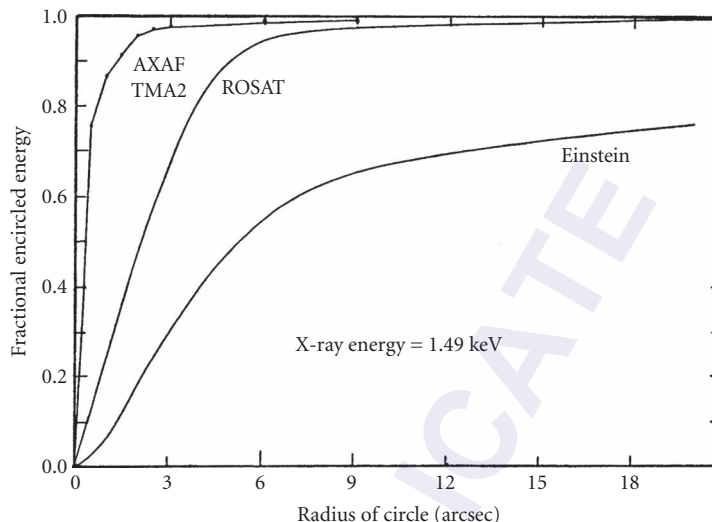


FIGURE 3 Fractional encircled energy plots at 1.49 keV for the Einstein Observatory, the European ROSAT telescope, and the AXAF technology mirror assembly (TMA2) are compared.

The apparent smoothness of the lacquer-coated surfaces and a desire for light weight and a high throughput or filling factor led Petre and Serlemitsos to develop the concept of tightly nested conical foil x-ray telescopes.^{14,15} These conical foil x-ray imaging mirrors are discussed in more detail in Chap. 48.

Still other novel optical fabrication concepts for grazing incidence x-ray optical surfaces include a variety of replication techniques.¹⁶ The Italian x-ray astronomy satellite (SAX) consisted of 30 nested coaxial mirrors electroformed over conical mandrels to a thickness ranging from 0.2 to 0.4 mm.¹⁷ And the European Space Agency (ESA) provided a dramatic increase in collecting area with its high throughput x-ray spectroscopy XMM mission featuring several modules of 58 tightly nested confocal Wolter Type I telescopes fabricated with a metal/epoxy replication technique.¹⁸

Another attempt to avoid the problems of grazing incidence telescope configurations has resulted in the practice of depositing enhanced reflectance multiple high and low atomic number thin films onto more conventional Schwarzschild configurations (or other normal incidence designs). This rapidly emerging technology¹⁹⁻²¹ is currently effective for wavelengths longer than about 40 Å and then only for very narrow bandwidths.

These advances in the fabrication of x-ray optics, along with new technologies for the production of x rays such as synchrotron sources, free electron lasers, and laser generated x-ray plasmas are stimulating renewed efforts in the areas of x-ray/EUV astronomy, soft x-ray microscopy, and x-ray/EUV microlithography.

44.3 OPTICAL DESIGN AND RESIDUAL ABERRATIONS OF GRAZING INCIDENCE TELESCOPES

The two-mirror grazing incidence telescopes described above (Wolter Types I, II, and III) were axially symmetric, confocal, and were designed by following the principles of on-axis stigmatic imaging developed over 300 years ago by Newton, Gregory, and Cassegrain. Wolter did not

attempt to optimize his first designs for off-axis or finite conjugate imaging, but he did write a second paper later that same year formulating a completely aplanatic (corrected for both spherical aberration and coma) version of his designs; that is, the Wolter-Schwarzschild designs.²² Van Speybroeck and Chase used computerized ray tracing algorithms to empirically determine the parametric effects of varying designs on the imaging performance of the Wolter²³ and Wolter-Schwarzschild²⁴ Type-I telescopes. Their findings, published in 1972 and 1973, were extremely useful but lacked the identification and interpretation of conventional aberrations (i.e., defocus, spherical aberration, coma, astigmatism, field curvature, distortion, etc.). In 1977, Werner²⁵ attempted the computational optimization of a Wolter Type I telescope by relaxing the surface shape constraint to that of a generalized axial polynomial. This resulted in almost flat imaging response across the field of view but significantly compromised the on-axis performance. Also in 1977, Winkler and Korsch²⁶ published an apparently decisive and thorough formulation of two-mirror grazing incidence aberration theory. The results showed, however due to their limited precision, that any classical Wolter type telescope was *already* aplanatic. Even Wolter himself would not agree with this as evidenced by his second paper.²² In 1979, a paper by Cash et al.²⁷ concluding that standard, near normal incidence aberration theory could be *exactly* applied to grazing incidence optical elements. However, Korsch²⁸ showed (even with his low precision) that there exists a first order coma term *not* present in normal aberration theory for a single mirror. Furthermore, Nariai²⁹ stated in 1987 that “it is not possible to use ordinary aberration theory because the expansion of the aberration in series of powers on the height of the object and on the radius of the pupil does not converge, . . . etc.” Nariai³⁰ later showed analytically that all aberrations in his own expansion must be integrated over the entire annular pupil, apparently, aberration *coefficients* in grazing incidence systems are themselves a function of pupil coordinates. Saha^{31–34} published several extensive theoretical and numerical analyses of the aberrations of generalized Wolter Type II as well as some Wolter Type I configurations. These papers provided much-needed insight and understanding concerning the aberrations of grazing incidence x-ray optical systems. The aberration theory of Wolter type x-ray telescopes is briefly reviewed in the Chap. 45.

A Wolter Type I grazing incidence x-ray telescope made up of a paraboloid and hyperboloid is illustrated in Fig. 4. The equation for a paraboloid with its vertex at z_p is given by

$$r_p^2 = 2r_p(z - z_p) \quad (1)$$

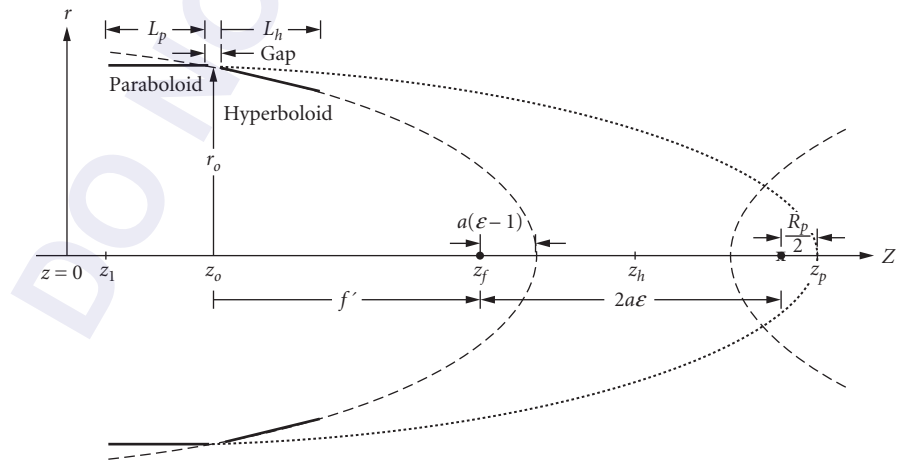


FIGURE 4 Wolter Type I grazing incidence telescope configuration.

where R_p is the paraboloid vertex radius of curvature and r_p is the radius of the paraboloid at the axial position z . The equation for a hyperboloid centered at z_h is given by

$$\frac{(z-z_h)^2}{a^2} - \frac{r_h^2}{b^2} = 1 \quad (2)$$

where a and b are the semimajor and semiminor axes of the hyperboloid. The eccentricity of the hyperboloid is determined by a and b

$$\varepsilon = \sqrt{\frac{a^2}{b^2} + 1} \quad (3)$$

The separation of the two hyperboloid foci is given by $2a\varepsilon$. If we superpose the rear hyperboloid focus with the paraboloid focus, the front hyperboloid focus becomes the system focus and $f' = z_f - z_o$ becomes the nominal focal length of the telescope. This is accomplished by positioning the origin of our coordinate system an arbitrary distance z_1 in front of the front edge of the paraboloid mirror and setting

$$z_p = z_1 + L_p + gap/2 + f' + 2a\varepsilon + R_p/2 \quad \text{and} \quad z_h = z_1 + L_p + gap/2 + f' + a\varepsilon \quad (4)$$

where L_p is the length of the paraboloid mirror and gap is the width of the gap between the paraboloid and the hyperboloid.

The optical prescription of a classical Wolter Type I x-ray telescope can thus be completely defined by the three independent parameters R_p , a , and b (or R_p , a , and ε). An optimized (maximized effective collecting area) Wolter Type I telescope can be obtained if we require the grazing angles of reflection from the paraboloid and the hyperboloid to be equal near their point of intersection. This constraint reduces the number of independent parameters defining the optical prescription to two.

For our purposes it is more convenient to choose the telescope radius at the intersection of the paraboloid and the hyperboloid, r_o , and the nominal focal length of the telescope f' as the parameters defining the optical prescription. The grazing angle at the joint is then given by

$$\alpha = \frac{1}{4} \arctan\left(\frac{r_o}{f'}\right) \quad (5)$$

The actual focal length, as measured from the system principal/nodal point, is slightly larger than the nominal focal length

$$f = f' + \frac{r_o^2}{2f'} \quad (6)$$

and the plate scale is the reciprocal of this focal length, expressed in arcsec per micrometers.

In addition to the telescope radius r_o and the nominal focal length f' , the remaining optical design parameters include the length of the paraboloid mirror L_p , the length of the hyperboloid mirror L_h , and the width of the gap between the two mirror elements. From these input parameters the actual dimensions of the mirror elements can be calculated as well as the obscuration ratio of the collecting aperture which determines both the geometrical collecting area and the diffraction-limited image characteristics.

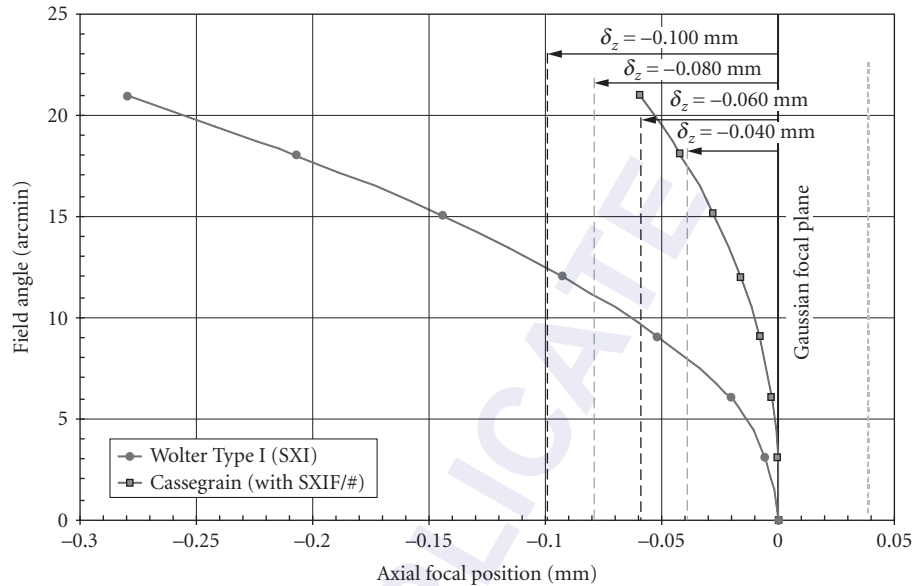


FIGURE 5 Field curvature exhibited by the SXI baseline design compared with that of a normal incidence Cassegrain telescope with the same focal ratio.

The classical Wolter Type I x-ray telescope design produces an ideal on-axis geometrical point image (zero spherical aberration); however, field curvature is a dominant limiting factor determining the off-axis performance of grazing incidence x-ray telescopes if a flat detector or focal plane is used. The severe field curvature of a classical Wolter Type I grazing incidence x-ray telescope is demonstrated in Fig. 5 by using the prescription of the solar x-ray imager (SXI) baseline design. SXI is a complementary, add-on instrument designed primarily for use on the GOES next generation geosynchronous weather satellites. However, its modular design is suitable for installation on many other spacecraft platforms. Its primary mission is to continuously observe the full solar disc, including coronal holes, active regions, flares, and coronal mass ejections.³⁵ The axial focal position was determined by minimizing the geometrical root-mean-square (rms) image size obtained from geometrical ray trace data.

If covering a significant field-of-view, the focal plane of such systems is frequently despaced to improve the off-axis performance; although, this results in a degraded (defocused) on-axis image. Geometrical optical performance from ray trace data is conveniently expressed in terms of rms image radius in arcsec. This quantity is calculated and plotted as a function of field angle for several different axial positions of the focal plane, see Fig. 6. Also shown for comparison is the performance curve that would be achieved with a curved detector conforming to the optimally curved focal surface.

Note that the curve for the best focal surface in Fig. 6 appears to have a linear and a quadratic component. This is consistent with findings of Van Speybroek and Chase.²³ For small field angles, the linear component dominates and will be associated with a comalike aberration.^{36,37} Similarly, the quadratic component of the curve will be associated with an astigmatism-like aberration.^{36,37} The curve corresponding to the Gaussian image plane also appears to consist of a linear and a quadratic component. The linear component (coma) is same as for the best focal surface as evidenced by the slope at small field angles. However, the quadratic component is significantly larger since it contains a contribution from both astigmatism and field curvature. Without attempting to develop any rigorous aberration theory, we will, throughout this chapter, refer to the image degradation (indicated by rms image size) that is linear with field angle as coma, and the degradation that is quadratic with field

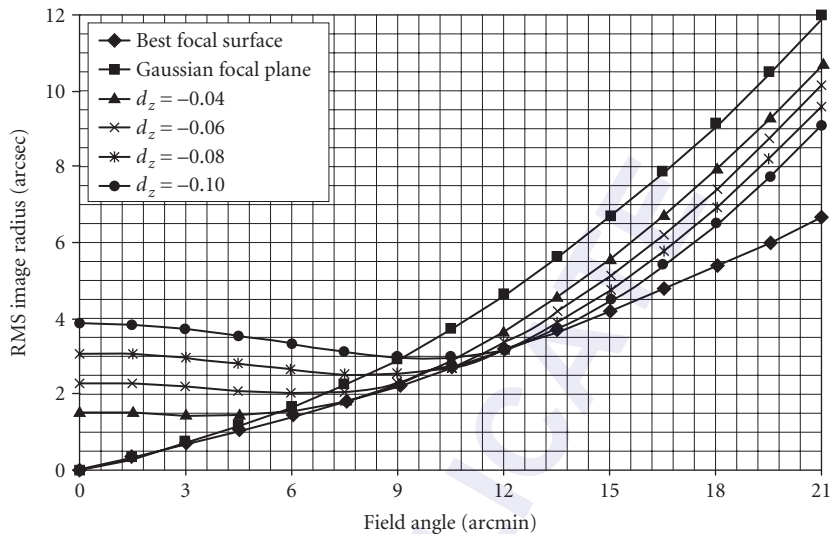


FIGURE 6 Geometrical performance of a classical Wolter Type I telescope for different axial positions of the focal plane.

angle as a combination of astigmatism and field curvature. Similarly, we will consider any image degradation on-axis to be caused by a combination of defocus and spherical aberration. Clearly, we are not trying to distinguish between various orders of aberrations. For example, the linear component of these curves includes all orders of linear coma. Likewise, the quadratic component includes third-order field curvature and astigmatism as well as all higher-order aberration terms that have a quadratic dependence on field angle; this includes the fifth-order aberration usually referred to as oblique spherical aberration.^{25,29–34} There are also, no doubt, cubic and higher-order contributions to the curves in Fig. 6; however, they do not appear to play a significant role for field angles less than 21 arcmin.

Despacing the focal plane of the classical Wolter Type I grazing incidence telescope clearly balances field curvature with defocus, thus improving the wide-field performance at the expense of the small-field performance. There are no additional design variables available to further correct or balance aberrations. This classical Wolter Type I design produces a stigmatic image on-axis in the Gaussian focal plane and has thus been used in virtually every x-ray stellar telescope built in the last 40 years, including the Einstein Observatory,³⁸ ROSAT,³⁹ and AXAF (the Chandra Observatory).⁴⁰

However, for wide-field x-ray imaging systems, stigmatic imaging on-axis is no longer an appropriate image quality requirement. Consider a solar physics application where the imaging system is a staring x-ray telescope pointed at the center of the sun. Sunspots or solar flares can appear anywhere on the solar disc; hence, the *field-weighted-average resolution element* as degraded by all error sources is an appropriate image quality criterion.⁴¹ This led to a departure from the conventional Wolter Type I grazing incidence x-ray telescope design which optimizes the on-axis resolution and is thus used for most stellar x-ray telescopes. Instead, a family of optimal hyperboloid-hyperboloid (HH) grazing incidence x-ray telescope designs was developed, where each member of the family is optimized for a different operational field-of-view (OFOV).⁴²

A schematic diagram of the resulting generalized Wolter Type I grazing incidence x-ray telescope is shown in Fig. 7. Note that the front focus of the primary mirror does *not* coincide with the rear focus of the secondary mirror as is the case with the classical Wolter Type I design. This *confocal delta* is indicated as the quantity Δ_{ps} in Fig. 7. Similarly, the system focal plane does *not* lie at the front focus of the secondary mirror. This displacement is indicated as Δf .

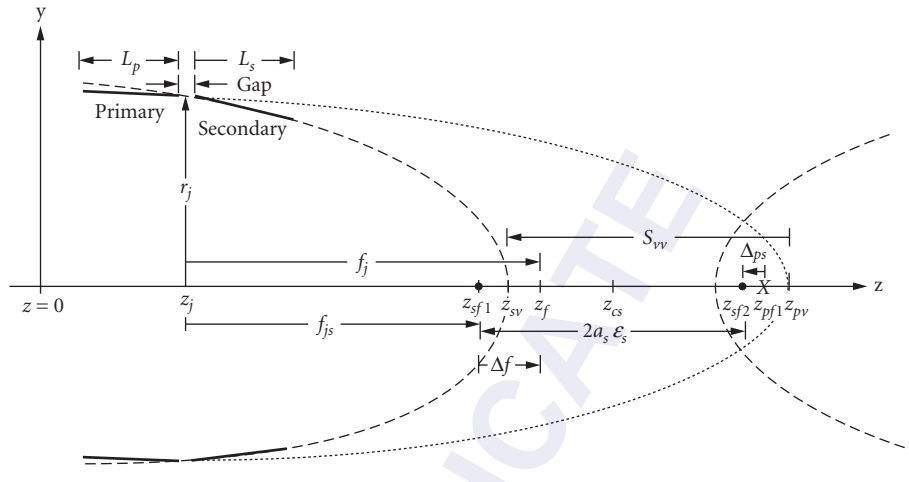


FIGURE 7 Hyperboloid-hyperboloid grazing incidence x-ray telescope design.

The above HH grazing incidence x-ray telescope thus has five design variables that can be varied in the design optimization process; the vertex radii of curvature of the two mirrors, their conic constants, and the vertex-to-vertex separation. For the SXI first-order design parameters a detailed comparison of the geometrical performance of the optimally despaced classical Wolter Type I design, the optimally despaced Wolter-Schwarzschild (WS) design, and the optimum HH design was performed. Figure 8 shows that the optimally despaced WS design and the optimum HH design always significantly

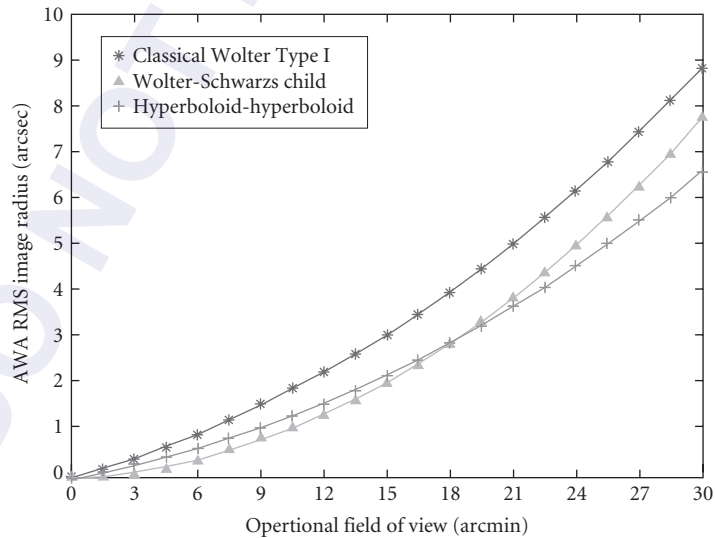


FIGURE 8 Comparison of the field weighted-average rms image radius versus OFOV for three different types of grazing incidence x-ray telescopes for wide-field imaging applications.

outperforms the optimally despaced classical Wolter Type I design. However, the optimally despaced WS design only outperforms the optimum HH design for OFOVs less than approximately 18 arcmin.⁴² Note that for OFOVs greater than approximately 18 arcmin, the *optimum HH design outperforms the optimally despaced WS design*.

44.4 IMAGE ANALYSIS FOR GRAZING INCIDENCE X-RAY OPTICS

A complete systems engineering analysis of x-ray telescope performance requires that we look at the effects of aperture diffraction, geometrical aberrations, surface scattering, and all other potential error sources such as assembly and alignment errors and metrology errors that appear in the mirror manufacturer's error budget tree as shown in Fig. 9. Detector effects must also be accurately modeled, particularly if it utilizes a staring mosaic detector array.

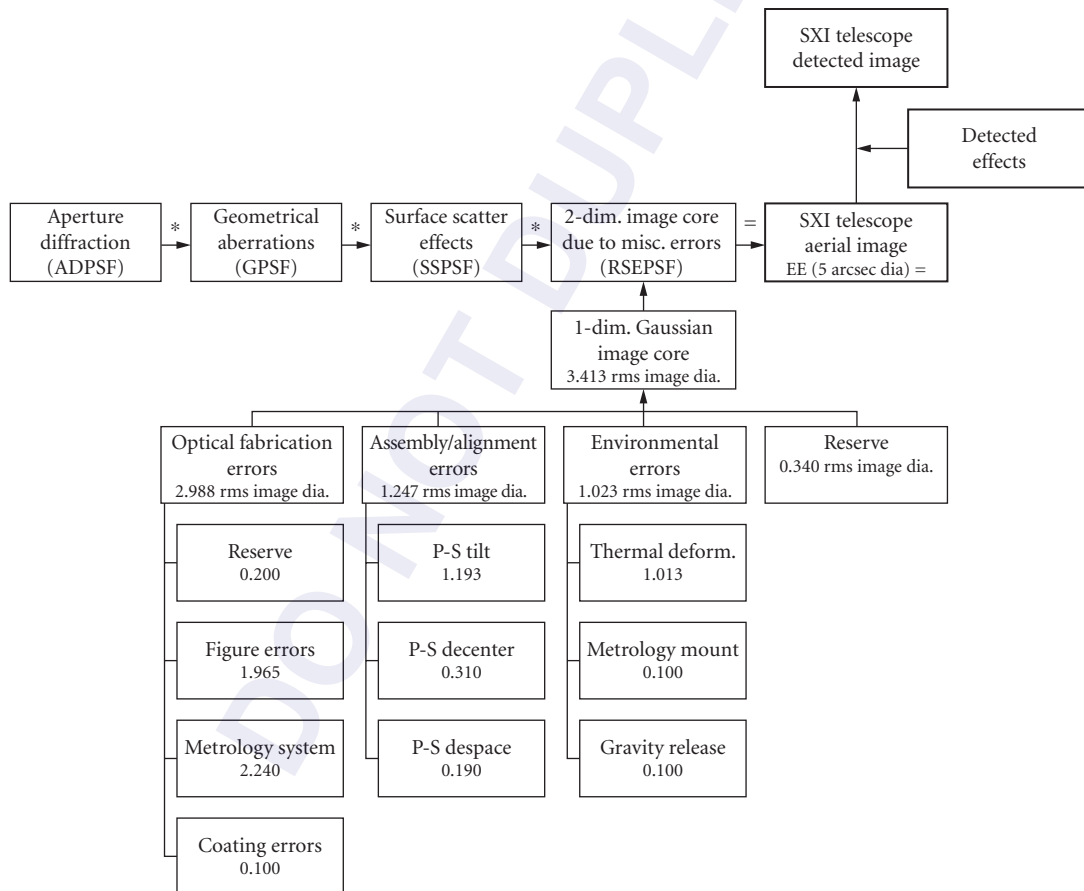


FIGURE 9 Error budget tree for manufacturing the SXI mirrors (* is a symbolic notation for the convolution operation).

Aperture diffraction effects are not necessarily negligible since the effects of the high obscuration ratios inherent to grazing incidence optics off-sets the effects of the very short x-ray wavelengths.^{43,44} Furthermore, most conventional, commercially available optical design and analysis codes do not adequately model the diffraction effects of highly obscured annular apertures. A semiempirical expression that does not require the high sampling density and associated computational problems when performing numerical convolutions can be used.⁴⁴

Likewise, most commercially available optical design and analysis codes do not provide aberration coefficients for highly obscured annular apertures. However, a high density spot diagram for each field angle of interest can be produced with a standard ray trace program (such as Code V or ZEMAX), the ray intercept data can be used to construct a ray *density* function in the focal plane which can be referred to as a geometrical point spread function (GPSF).

When light is reflected from an imperfect optical surface, the reflected radiation consists of a specularly reflected component and a diffusely reflected (or scattered) component. The scattered light behavior can be calculated from the surface characteristics (see Vol. I, Chap. 8, "Surface Scattering," by Eugene L. Church and Peter Z. Takacs). The surface autocovariance (ACV) function or the surface power spectral density (PSD) function completely determines the scattered light behavior. For the SXI program, we used NASA's EEGRAZE code to calculate the scatter profile for the SXI telescope. Either assumed or measured surface metrology data (PSD function) must be supplied as input for the EEGRAZE code. This scatter profile was then used to construct the surface scatter point spread function (SSPSF).

And finally, the effects of all other potential error sources appearing in the SXI mirror manufacturer's error budget tree were modeled as a single contribution to the final rms image core diameter. The aerial image produced by the x-ray telescope was then modeled by the *convolution* of the individual PSFs associated with the respective error sources as illustrated schematically in Fig. 10.⁴⁵

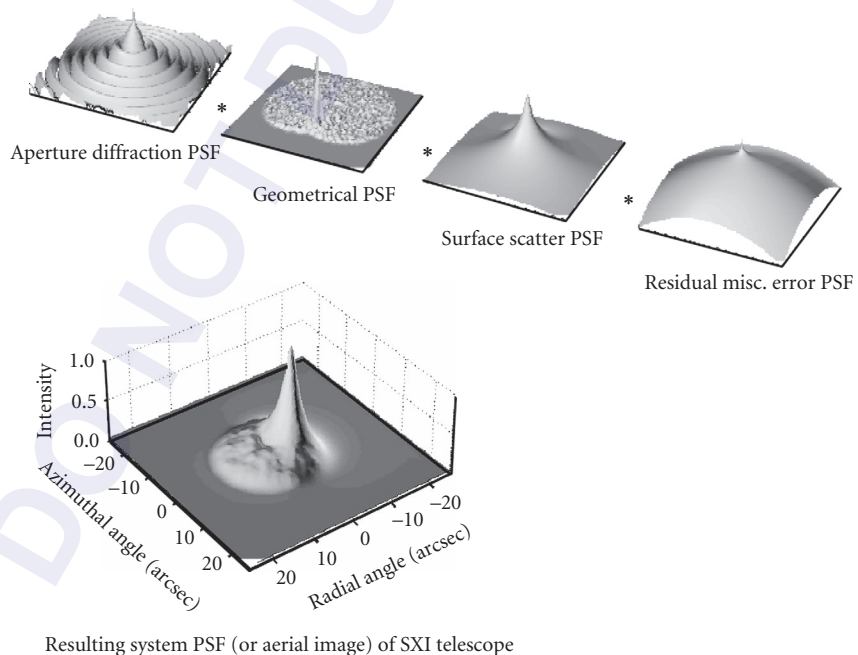


FIGURE 10 Illustration of the PSFs for the individual image degradation mechanisms and their resulting convolution (*aerial image*) of the SXI telescope for a field angle of 15 arcmin and a wavelength of 44.7 Å. (See also color insert.)

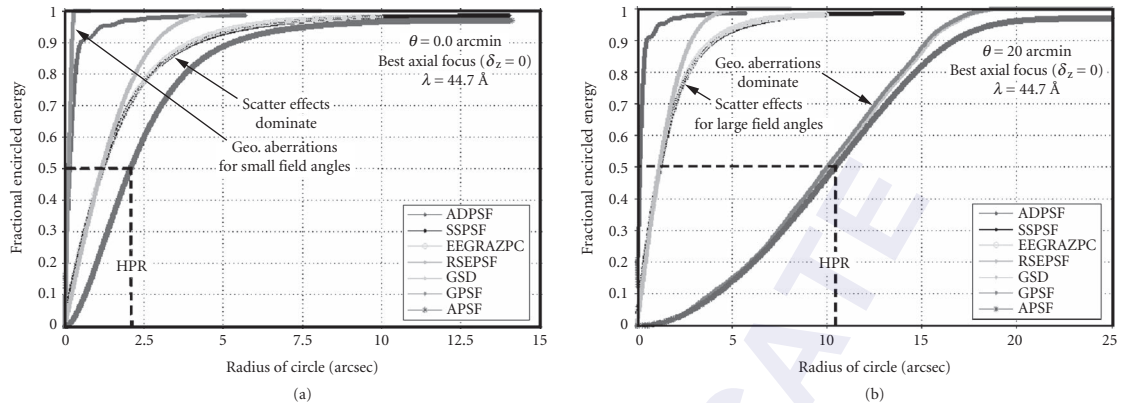


FIGURE 11 The fractional encircled energy of the aerial image and the four functions contributing to it provide insight into the image quality of a grazing incidence x-ray telescope. Note that scatter effects dominate geometrical aberrations for small field angles and geometrical aberrations dominate scatter effects for large field angles. (See also color insert.)

Fractional encircled energy plots of these four functions and the aerial image provide insight concerning the relative effects of the various image degradation mechanisms as shown in Fig. 11a and 11b. Note that for small field angles the scattering effects dominate geometrical aberrations, whereas for large field angles the geometrical aberrations dominate the scattering effects. Note also that the half power radius (HPR) of the on-axis aerial image is 2.0 arcsec and the HPR of the aerial image at a 20 arcmin field angle is 10.5 arcsec.

When a single detector is scanned over an aerial image, the detected image (in the scan direction) can be modeled by the convolution of the aerial image with the detector; or conversely, one can multiply the modulation transfer function (MTF) of the imaging system by the detector MTF. However, mosaic detector arrays employ a discrete sampling interval that causes these systems to exhibit a particular kind of *local shift variance* which causes the appearance of the reconstructed image to vary with the location of the aerial point spread function (PSF) relative to the sampling (i.e., pixel) grid.⁴⁶ The MTF approach to system performance analysis is thus not directly applicable to systems utilizing staring mosaic detector arrays.

Figure 12 illustrates the example of a Gaussian aerial PSF slightly larger than a detector pixel being sampled (averaging over each detector pixel) to produce a *detected* point spread function (DPSF).

An interpolation scheme is then used to reconstruct a smooth DPSF. This is done for three situations: (1) when the aerial PSF is precisely “registered” at the center of a detector pixel, (2) when the aerial PSF is positioned on the boundary between two detector pixels, and (3) when the aerial PSF is positioned at a point where four detector pixels meet. This registration error can result in as much as a 40 percent variation in the calculated HPR of the reconstructed DPSF.

For an application where the telescope is being operated as a staring telescope recording fine detail in an extended image (random location of aerial PSF on pixel), the “average unregistered” detected point spread function (AUDPSF) is given by the convolution of the registered detected point spread function (RDPSF) by the unit cell of the sampling grid.⁴⁷

Figure 13 illustrates the calculation of both the reconstructed registered DPSF and the reconstructed average unregistered DPSF. Since the aerial PSF is represented as a dense numerical array, the averaging over the individual pixels is referred to as a “binning” operation. Care is taken to precisely “register” the sampling detector grid by positioning it so as to maximize the signal produced by a given pixel. We then use a cubic interpolation technique to reconstruct the RDPSF. Finally, we convolve by the unit cell of the sampling grid to produce the AUDPSF.⁴⁷

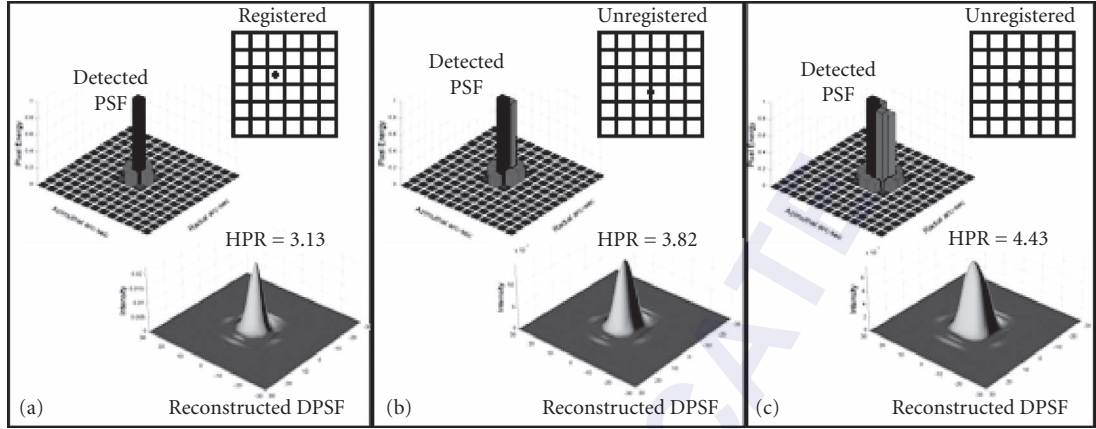


FIGURE 12 The DPSF and the reconstructed DPSF for: (a) precisely “registered” aerial PSF; (b) aerial PSF centered on boundary between two pixels; and (c) aerial PSF positioned where four pixels meet. (See also color insert.)

We can now calculate the field-weighted-average HPR of either the aerial image or the AUDPSF for any specific operational field-of-view (OFOV).

$$\text{HPR}_{fwa} = \frac{1}{A_T} \int_{\theta=0}^{\text{OFOV}} \text{HPR}(\theta) 2\pi\theta d\theta \quad (7)$$

where $A_T = \pi(\text{OFOV})^2$

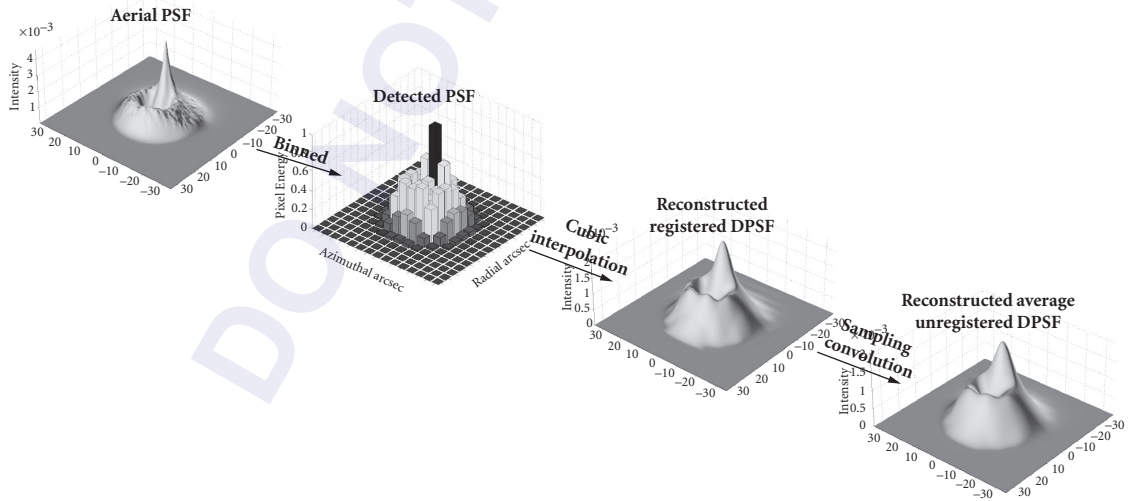


FIGURE 13 A graphical illustration of the numerical computation technique for modeling both the reconstructed “registered” DPSF and the reconstructed “average unregistered” DPSF is indicated. (See also color insert.)

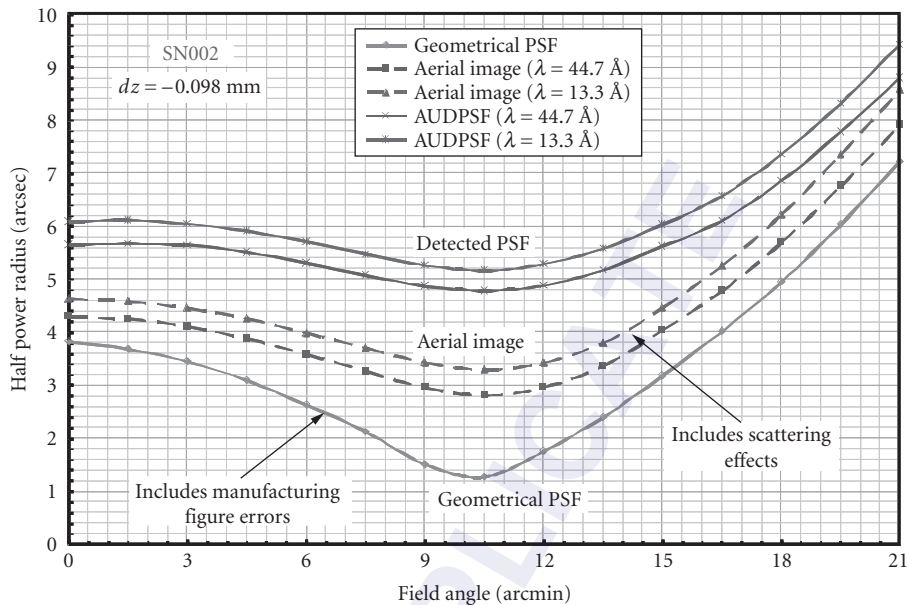


FIGURE 14 Comparison of the predicted HPR versus field angle of the geometrical PSF, the aerial image, and the average unregistered detected PSF for the SXI telescope at two different wavelengths (44.7 and 13.3 Å). (See also color insert.)

Defining a spatial (angular) resolution element as a circle of radius $HPR(\theta)$, we can also calculate the number of spatial resolution elements N in that OFOV

$$N = \text{number of res. ele.} = 2 \int_{\theta=0}^{\text{OFOV}} \frac{\theta}{HPR^2(\theta)} d\theta \quad (8)$$

Minimizing the HPR_{fwa} for a given OFOV maximizes the number of spatial resolution elements N , and thus also maximizes the amount of information contained in the image.

A comparison of the predicted HPR versus field angle of the geometrical PSF (geometrical aberrations only, as determined by ray trace analysis), the aerial image (including all system errors except detector effects), and the averaged unregistered detected PSF for the SN002 SXI telescope that was launched on GOES-13 is illustrated in Fig. 14.^{48,49} Predictions were performed for two different wavelengths, 44.7 and 13.3 Å. These curves provide insight concerning the relative effect of various error sources upon the image quality, and certainly demonstrate the inadequacy of merely performing a geometrical ray trace analysis of these grazing incidence x-ray telescopes. Similar graphs were obtained for each of the other “as-manufactured” SXI telescopes.⁵⁰

44.5 VALIDATION OF IMAGE ANALYSIS FOR GRAZING INCIDENCE X-RAY OPTICS

The SXI was launched on GOES-13 on May 24, 2006, and the “first light” image was recorded with on July 6, 2006. Figure 15 shows one of the first raw on-orbit images, which was described as “exquisite” by at least one solar physicist.

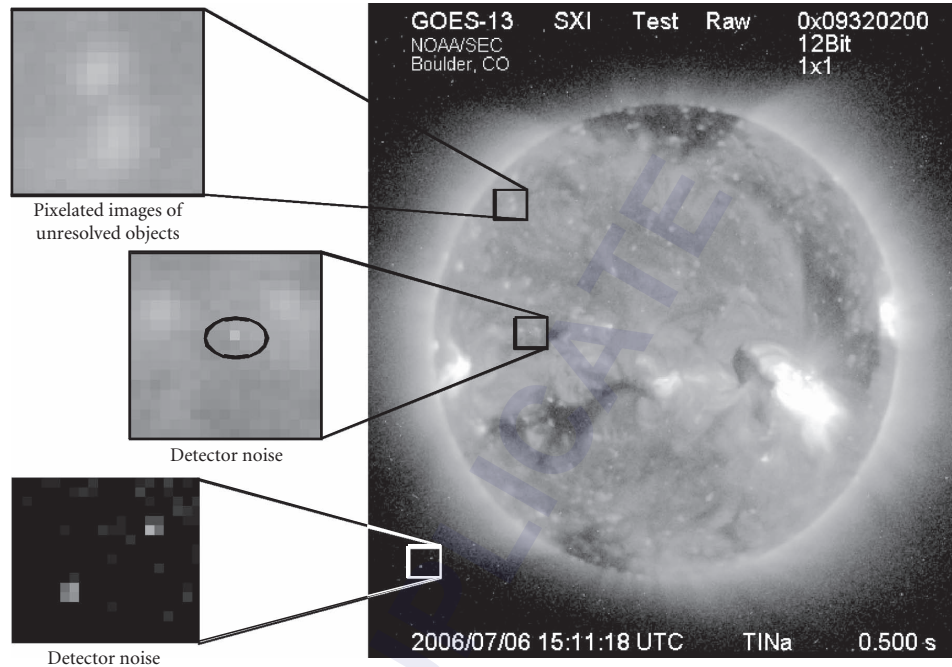


FIGURE 15 On-orbit solar image with three small regions highlighted and magnified for detailed inspection. This allows one to distinguish between images of unresolved bright features on the sun and merely detector noise. (See also color insert.)

Upon close inspection, it is reasonable to assume that the abundance of small features in Fig. 15 consist of the images of *unresolved bright spots* on the sun. This assumption becomes more obvious as we highlight and magnify small areas in the raw on-orbit image. As examples, we illustrate three separate regions of the solar image that we have highlighted and magnified. The first highlighted area shows two closely spaced pixelated images of unresolved features on the sun. Note that they are 6 to 8 pixels in diameter which is consistent with the numerical predictions of the detected point spread function in Fig. 13. The second highlighted area shows a very small bright spot which, when magnified, is obviously merely detector noise (a single detector pixel with unusually high responsivity). Similarly the third highlighted area depicts detector noise which should not be mistaken with actual bright spots in the solar corona as they would have to be at least as large as a point spread function.

In order to extract quantitative data from the raw on-orbit images to compare with our image quality predictions, the solar disc was divided into seven radial zones. The sizes of small prominent features in each radial zone were then measured (and normalized by the known size of the solar disc) and averaged. The apparent (measured) diameter of these features was assumed to be equal to twice their half power radius (HPR).⁵¹

These average-measured HPRs were then compared to the results of the detailed image quality predictions illustrated earlier by superposing these average values for each radial zone on the performance prediction curves presented in Fig. 14. This extremely simple and straightforward comparison of raw on-orbit experimental data with the detailed and exhaustive numerical predictions of our systems engineering analysis of image quality is displayed in Fig. 16. The excellent agreement of this quantitative data with image quality predictions provides an experimental validation of the detailed systems engineering analysis of image quality required for grazing incidence x-ray optics.

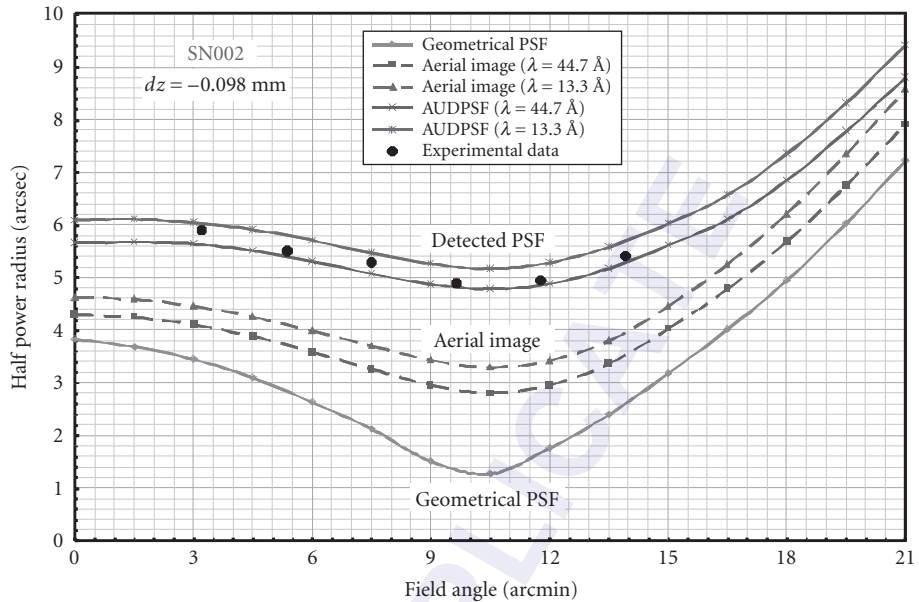


FIGURE 16 Experimental validation of an exhaustive systems engineering analysis of image quality for grazing incidence x-ray telescopes, including the modeling of surface scatter and detector effects. (See also color insert.)

44.6 REFERENCES

1. P. Kirkpatrick and A. V. Baez, "Formation of Optical Images by X-Rays," *J. Opt. Soc. Am.* **38**:776 (1948).
2. H. Wolter, "Mirror Systems with Glancing Incidence on Image Producing Optics for X-Rays," *Ann. Phys.* **10**:94 (1952).
3. J. E. Harvey, K. L. Lewotsky, and A. Kotha, "Performance Predictions of a Schwarzschild Imaging Microscope for Soft X-Ray Applications," *Opt. Eng.* **35**:2423–2436 (Aug. 1996).
4. B. R. Aschenbach, "X-Ray Telescopes," *Rep. Prog. Phys.* **48**:579–629 (1985), The Institute of Physics, Great Britain.
5. L. P. Van Speybroeck, "Grazing Incidence Optics for the U.S. High Resolution X-Ray Astronomy Program," *Opt. Eng.* **27**:1398–1403 (1988).
6. A. Slomba, R. Babish, and P. Glenn, "Mirror Surface Metrology and Polishing for AXAF/TMA," *Proc. SPIE* **597**:40 (1985).
7. B. R. Aschenbach, "Design, Construction, and Performance of the ROSAT High-Resolution Mirror Assembly," *Appl. Opt.* **27**:1404–1413 (1988).
8. R. Giacconi, et al., "The Einstein (HEAO 2) X-Ray Observatory," *Astrophys. J.* **230**:540 (1979).
9. L. P. Van Speybroeck, "Einstein Observatory (HEAO B) Mirror Design and Performance," *Proc. SPIE* **184**:2 (1979).
10. J. E. Harvey, "Recent Progress in X-Ray Imaging," presented at the *AIAA Space Programs and Technologies Conference*, Huntsville, AL, Sept. 1990.
11. S. Bowyer and J. Green, "Fabrication, Evaluation, and Performance of Machined Metal Grazing Incidence Telescopes," *Appl. Opt.* **27**:1414–1422 (1988).
12. R. C. Catura, E. G. Joki, D. T. Roethig, and W. J. Brookover, "Lacquer Coated X-Ray Optics," *Proc. SPIE* **640**:140–144 (1986).

13. J. A. Nousek, et al., "Diamond-Turned Lacquer-Coated Soft X-Ray Telescope Mirrors," *Appl. Opt.* **27**:1430–1432 (1988).
14. R. Petre and P. J. Serlemitsos, "Conical Imaging Mirrors for High-Speed X-Ray Telescopes," *Appl. Opt.* **24**:1833 (1985).
15. R. Petre, P. J. Serlemitsos, F. E. Marshall, K. Jahoda, and H. Kunieda, "In Flight Performance of the Broad-Band X-Ray Telescope," *Proc. SPIE* **1546**:72–81 (1991).
16. Y. Matsui, M. P. Ulmer, and P. Z. Takacs, "X-Ray and Optical Profiler Analysis of Electroformed X-Ray Optics," *Appl. Opt.* **27**:1558–1563 (1988).
17. O. Citterio, et al., "Optics for the X-Ray Imaging Concentrators Aboard the X-Ray Astronomy Satellite SAX," *Proc. SPIE* **830**:139 (1987).
18. W. Egle, H. Bulla, P. Kaufmann, B. Aschenbach, and H. Brauningner, "Production of the First Mirror Shell for ESA's XMM Telescope by Application of a Dedicated Large Area Replication Technique," *Opt. Eng.* **29**:1267 (1990).
19. E. Spiller, "Reflective Multilayer Coatings in the Far UV Region," *Appl. Opt.* **15**:2333 (1976).
20. A. B. C. Walker, Jr., T. W. Barbee, Jr., R. B. Hoover, and J. F. Lindblom, "Soft X-Ray Images of the Solar Corona with a Normal-Incidence Cassegrain Multilayer Telescope," *Science* **241**:1781 (1988).
21. N. M. Ceglio (ed.), Thirty-one papers presented at the 1991 SPIE conference on "Multilayer Optics for Advanced X-Ray Applications," *Proc. SPIE* **1547** (1991).
22. H. Wolter, "Generalized Schwarzschild Mirror Systems with Glancing Incidence as Optics for X-Rays," *Ann. Phys.* **10**:286 (1952).
23. L. P. Van Speybroeck and R. C. Chase, "Design Parameters of Paraboloid-Hyperboloid Telescopes for X-Ray Astronomy," *Appl. Opt.* **11**:440–445 (1972).
24. R. C. Chase and L. P. Van Speybroeck, "Wolter-Schwarzschild Telescopes for X-Ray Astronomy," *Appl. Opt.* **12**:1042–1044 (1973).
25. W. Werner, "Imaging Properties of Wolter I Type X-Ray Telescopes," *Appl. Opt.* **16**:764–773 (1977).
26. C. E. Winkler and D. Korsch, "Primary Aberrations for Grazing Incidence," *Appl. Opt.* **16**:2464–2469 (1977).
27. W. Cash, D. L. Sheeley, and J. H. Underwood, "Space Optics Imaging X-Ray Optics Workshop," *Proc. SPIE* **184**:228 (1979).
28. D. Korsch, *Reflective Optics*, Chapter 11, Academic Press Inc., Boston, MA, 1991, pp. 282–284.
29. K. Nariai, "Geometrical Aberrations of a Generalized Wolter Type I Telescope," *Appl. Opt.* **26**:4428–4432 (1987).
30. K. Nariai, "Geometrical Aberrations of a Generalized Wolter Type I Telescope 2: Analytical Study," *Appl. Opt.* **27**:345–350 (1988).
31. T. T. Saha, "Transverse Ray Aberrations for Paraboloid-Hyperboloid Telescopes," *Appl. Opt.* **24**:1856–1863 (1985).
32. T. T. Saha, "Transverse Ray Aberrations of Wolter Type I Telescopes," *Proc. SPIE* **640**:10–19 (1986).
33. T. T. Saha, "General Surface Equations for Glancing Incidence Telescopes," *Appl. Opt.* **26**:658–663 (1987).
34. T. T. Saha, "Aberrations for Grazing Incidence Telescopes," *Appl. Opt.* **27**:1492–1498 (1988).
35. P. L. Bornman, D. Speich, J. Hirman, V. Pizzo, R. Grubb, Balch, and G. Heckman, "GOES Solar X-Ray Imager: Overview and Operational Goals," in *GOES-8 and Beyond*, E. R. Washwell, ed., *Proc. SPIE* **2812**:309–319 (1996).
36. H. H. Hopkins, *Wave Theory of Aberrations*, Clarendon Press, Oxford, 1950.
37. V. N. Majahan, *Optical Imaging and Aberrations*, SPIE Optical Engineering Press, Bellingham, WA, 1998.
38. L. P. Van Speybroeck, "Einstein Observatory (HEAO B) Mirror Design and Performance," in *Space Optics-Imaging X-Ray Optics Workshop*, M. Weisskopf, ed., *Proc. SPIE* **184**:2–11 (1979).
39. B. Aschenbach, "Design, Construction, and Performance of the ROSAT High-Resolution Mirror Assembly," *Appl. Opt.* **27**:1404–1413 (1988).
40. L. P. Van Speybroeck, "Grazing Incidence Optics for the U.S. High-resolution X-Ray Astronomy Program," *Opt. Eng.* **27**:1398–1403 (1988).
41. P. L. Thompson and J. E. Harvey, "Development of an Image Quality Criterion for Wide-Field Applications of Grazing Incidence X-Ray Telescopes," *Proc. SPIE* **3766-13**:162–172 (1999).
42. J. E. Harvey, A. Krywonos, P. L. Thompson, and T. T. Saha, "Grazing Incidence Hyperboloid-Hyperboloid Designs for Wide-Field X-Ray Imaging Applications," *Appl. Opt.* **40**:136–144 (Jan. 2001).

43. J. E. Harvey, "Diffraction Effects in Grazing Incidence X-Ray Telescopes," *J. X-Ray Sci. Tech.* **3**:68–76 (1991).
44. P. L. Thompson and J. E. Harvey, "A Systems Engineering Analysis of Aplanatic Wolter Type I X-Ray Telescopes," *Opt. Eng.* **39**:1677–1691 (Jun. 2000).
45. J. E. Harvey and A. Krywonos, "A Systems Engineering Analysis of Image Quality," *Proc. SPIE* **4093B-50**: 379–388 (2000).
46. S. K. Park, R. Schowengerdt, and M. Kaczynski, "Modulation-Transfer-Function Analysis for Sampled Image Systems," *Appl. Opt.* **2**:2572–2582 (Aug. 1, 1984).
47. G. D. Boreman, *Modulation Transfer Function in Optical and Electro-Optical Systems*, SPIE Press, *Tutorial Texts in Optical Engineering* **TT52**:41 (2001).
48. J. E. Harvey, M. Atanassova, and A. Krywonos, "Including Detector Effects in the Design of Grazing Incidence X-Ray Telescopes," presented at *SPIE's International Symposium on Astronomical Telescopes and Instrumentation*, Glasgow, Scotland, June 2004; published in *Proc. SPIE* **5497-76**:90–99 (Jun. 2004).
49. M. I. Atanassova, "Optimizing the Performance of As-Manufactured Grazing Incidence X-Ray Telescopes Using Mosaic Detector Arrays," Ph.D. Dissertation, University of Central Florida, Orlando, FL, 2005.
50. J. E. Harvey, M. Atanassova, and A. Krywonos, "Systems Engineering Analysis of Five 'As-Manufactured', SXI Telescopes," *Proc. SPIE* **5867-15**:114–124 (Aug. 2005).
51. J. E. Harvey, A. Krywonos, M. Atanassova, and P. L. Thompson, "The Solar X-Ray Imager (SXI) on GOES-13: Design, Analysis, and On-Orbit Performance," presented at *SPIE's International Symposium on Optics and Photonics*, San Diego, CA, August 2007; published in *Proc. SPIE* **6689-11**:6890I–66890I-9 (Aug. 2007).

ABERRATIONS FOR GRAZING INCIDENCE OPTICS

Timo T. Saha

*NASA/Goddard Space Flight Center
Greenbelt, Maryland*

45.1 GRAZING INCIDENCE TELESCOPES

Large number of grazing incidence telescope configurations have been designed and studied (see Chaps. 33 and 47 in this volume). Wolter¹ telescopes are commonly used in astronomical applications. Wolter telescopes consist of a paraboloidal primary mirror and a hyperboloidal or an ellipsoidal secondary mirror. There are eight possible combinations of Wolter telescopes.² Out of these possible designs only type 1 and type 2 telescopes are widely used. Type 1 telescope is typically used for x-ray applications and type 2 telescopes are used for EUV applications.

Wolter-Schwarzschild (WS) telescopes³ offer improved image quality over a small field of view. The WS designs are stigmatic and free of third-order coma and, therefore, the point spread function (PSF) is significantly better over a small field of view. Typically the image is more symmetric about its centroid. The eight designs of WS telescopes have not been widely used because the surface equations are complex parametric equations complicating the analysis, and typically the resolution requirements are too low to take full advantage of the WS designs.

There are several other design options. Most notable are wide-field x-ray telescope designs. Polynomial designs were originally suggested by Burrows⁴ and hyperboloid-hyperboloid designs for solar physics applications were designed by Harvey.⁵

No general aberration theory exists for grazing incidence telescopes that would cover all the design options. Several authors have studied the aberrations of grazing incidence telescopes.⁶⁻⁹ A comprehensive theory of Wolter type 1 and 2 telescopes has been developed.^{10,11} Later this theory was expanded to include all possible combinations of grazing incidence and also normal incidence paraboloid-hyperboloid and paraboloid-ellipsoid telescopes.¹² In this chapter the aberration theory of Wolter-type telescopes is briefly reviewed.

45.2 SURFACE EQUATIONS

The surface equations of the grazing incidence telescopes can be combined to a general form covering large number of design options. General surface equations have been developed for grazing incidence telescopes.² Unfortunately, these equations are too complex to be a basis for the aberration

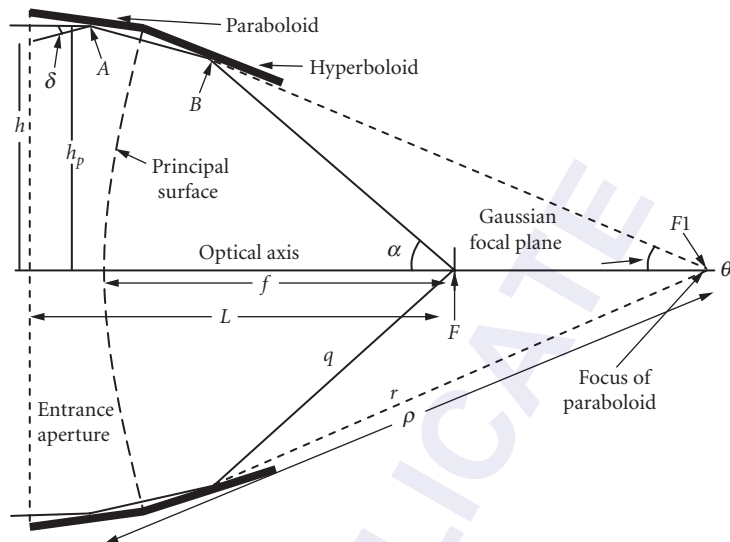


FIGURE 1 Cross-section of Wolter type 1 telescope showing the optical components, ray paths, and the parameters.

theory. The equations of the Wolter telescopes presented in the cylindrical coordinate system are relatively simple. Assuming an incoming ray hits the telescope entrance aperture at radial and azimuthal location (h, β) , it strikes the primary mirror at a location A , the secondary mirror at a location B , and the focal plane at a location F , as shown in Figs. 1 and 2. The extension of the ray which hits the secondary at the location B would intersect the optical axis at $F1$, at the common focus of the paraboloid and hyperboloid.

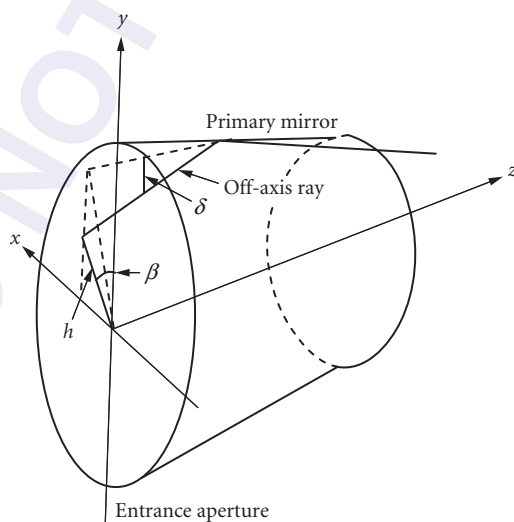


FIGURE 2 Primary mirror of the grazing incidence telescope showing entrance aperture and incoming off-axis ray.

Assuming the primary and the secondary mirror are surfaces of revolution and the primary mirror focus coincides with the secondary mirror focus, then the surface equation of the primary mirror can be written as

$$1/\rho = (\cos \theta - 1)/R_1 \quad (1)$$

where ρ is the distance from the paraboloid-hyperboloid focus F1 to a point on the primary mirror, θ is an angle this ray makes with the optical axis, and R_1 is the vertex radius of curvature of the primary mirror.

The secondary hyperboloid can be expressed either as a function of angle α or θ

$$1/q = (\varepsilon \cos \alpha - 1)/R_2 \quad (2)$$

$$1/r = (1 - \varepsilon \cos \theta)/R_2 \quad (3)$$

where q is the distance from the focus F to the point B on the secondary mirror, α is an angle q makes with the optical axis, R_2 is the vertex radius of curvature of the secondary mirror, ε is the eccentricity of the secondary, and r is the distance from the point B on the secondary to the focus F1 of the primary.

The principal surface^{13,14} of the Wolter telescopes is defined by the intersection points of the extensions of the incoming rays and the extensions of the rays reflected on the secondary, q . The principal surface goes through the intersection points of the primary mirror and secondary mirror. The principal surface of the Wolter telescopes is always a paraboloid.¹³ It is useful to define the focal length f of the telescope to be the distance from the vertex of the principal surface to the telescope focus F . Quite often the focal length is defined as the axial distance from telescope focus to the primary-secondary surface intersection plane.

A useful relation for the paraboloid-hyperboloid or paraboloid-ellipsoid telescopes is¹⁵

$$h_p = 2f \tan(\alpha/2) \quad (4)$$

This equation ties the telescope object side to the image side. Equation (4) shows that the Wolter telescopes do not satisfy the Abbe's sine condition

$$h_p = f \sin \alpha \quad (5)$$

If the angle α is small, as is the case in x-ray optical systems, then the trigonometric terms could be expanded in α . The difference in Eqs. (4) and (5) would be in the third-order term indicating that the third-order coma should be small in the Wolter telescopes and the telescopes nearly satisfy the Abbe's sine condition.

45.3 TRANSVERSE RAY ABERRATION EXPANSIONS

Transverse ray aberration (TRA) expansions are expressed as functions of entrance aperture coordinates (h, β) shown in Fig. 2. The expansions presented here are based on the format introduced by Cox.¹⁶ The derivation is simple but rather lengthy. In the derivation an off-axis ray making an angle δ is mathematically traced from the primary mirror to the secondary mirror and to the focal plane. The surface intersection points on the primary and on the secondary are solved with respect to an on-axis ray and expanded in δ and h . The image plane intersection point (H_x, H_y) of this ray are the final image coordinates. They are then expressed as functions

of entrance aperture coordinates (h, β) and image height $H_0 (=f \tan(\delta))$. The TRA expansions of Wolter telescopes are

$$H_x = H_0(A_2 h^2 \sin(2\beta) + B_2 h^4 \sin(2\beta)) + H_0^2 \sin\beta [A_3 h + B_3 h^3 + C_1 h^5 + \cos^2(\beta) B_4 h^3] + H_0^3 \sin(2\beta) (B_5 h^2 + C_2 h^4) / 2 \quad (6)$$

$$H_y = H_0 [1 + A_2 h^2 (2 + \cos(2\beta)) + B_2 h^4 (1.5 + \cos(2\beta) - \lambda_1 h^4)] + H_0^2 \cos\beta \{A_4 h + h^3 [B_3 + \lambda_2 + B_4 (1 + \cos^2(\beta))]\} + H_0^3 \cos^2(\beta) (B_5 h^2 + C_2 h^4) \quad (7)$$

In the derivation and resulting expansions only significant terms for grazing incidence telescopes are kept. The image height terms (H_0) higher than third order are dropped. The radial height, h , terms higher than fourth order are also dropped. The parameters of the TRA expansions are listed as follows in Eqs. (8) through (19) as functions of the telescope basic parameters f , R_1 , R_2 , ϵ , and L . The parameter L is the length of the telescope from the entrance aperture to the focal plane.

$$A_2 = 1/(4f^2) \quad (8)$$

$$A_3 = [K/(2f) - 1/R_1 - 1/R_2]/f \quad (9)$$

$$A_4 = [3K/(2f) - 1/R_1 - 1/R_2]/f \quad (10)$$

$$B_2 = 1/(8f^4) \quad (11)$$

$$B_3 = [K/f - 1/R_1 - 1/R_2 - 2f/(R_1 R_2)]/(4f^3) \quad (12)$$

$$B_4 = K/(2f^4) \quad (13)$$

$$B_5 = [-K + 1 + 2(R_1/R_2) + (R_1/(2f))]/(f^2 R_1 R_2) \quad (14)$$

$$\lambda_1 = 1/(16f^4) \quad (15)$$

$$\lambda_2 = [-K/(4f) + 1/R_1 + 1/R_2]/f^3 \quad (16)$$

$$C_1 = -1/(4f^4 R_1 R_2) \quad (17)$$

$$C_2 = (-1/R_1 + 1/R_2)/(2f^2 R_1 R_2) \quad (18)$$

$$K = [L - 2\epsilon R_2/(\epsilon^2 - 1)]/f + R_1/R_2 \quad (19)$$

The Wolter telescopes are stigmatic on-axis and, therefore, the designs are free of all orders of spherical aberration terms. The first term in the expansion is the third-order coma. The coma coefficient is proportional to the inverse square of the focal length and the coma term is inversely proportional to the square of the telescopes' f -number. It does not depend on the location of the entrance aperture or the other parameters.

The third-order aberration terms A_3 and A_4 are proportional to second order of the image height (H_0) and first-order aperture height h . Both terms depend on the location of the entrance aperture. Astigmatism and field curvature can be derived from these terms.¹⁰

All fifth-order terms are represented in the expansions relative to Cox's work.¹⁶ As the third-order spherical aberration, the fifth-order spherical aberration term is zero since the Wolter telescopes are stigmatic. Coefficients B_2 and λ_1 represent fifth-order circular coma. The term including the B_3 , B_4 , and λ_2 is the so-called *astrolate aberration*.¹⁶ If coefficient $B_4 = 0$, the term is called *fifth-order oblique spherical aberration*. The terms represented by B_5 coefficient is the fifth-order elliptical coma aberration. Approximations suitable for grazing incidence telescopes were made in the derivation of B_5 coefficient.

Two seventh-order terms C_1 and C_2 proportional to $H_0^2 h^5$ and $H_0^3 h^4$ are approximations. Exact solutions are very complex formulas of the basic parameters.

Typically the seventh-order terms and the fifth-order terms become more important when the grazing angles of the surfaces decrease. Expanding the TRA equations as a function of radial height h is not the best choice in case of grazing incidence systems. For example, expanding the TRA equations as a function of Δh ($h - h_{\text{int}}$), where h_{int} is the radial height at the primary secondary intersection plane, could lead to fewer terms and aberration coefficients more meaningful for the grazing incidence telescopes.

The RMS image size can be represented as a function of the aberrations coefficients.¹⁷ The resulting equation is rather complex function of the aberration coefficients and the field angle.

45.4 CURVATURE OF THE BEST FOCAL SURFACE

All combinations of Wolter telescopes suffer from large curvature of the best focal surface. This limits the field of view of the telescopes. The shape of the best focal surface is parabolic. The radius of curvature R_d of this surface can be estimated from the TRA equation¹¹

$$1/R_d = f \left[A_3 + A_4 + (R_{\text{max}}^2 + R_{\text{min}}^2)(B_3 + B_4 + \lambda_2/2) + \frac{C_1}{3} \frac{R_{\text{max}}^6 - R_{\text{min}}^6}{R_{\text{max}}^2 - R_{\text{min}}^2} \right] \quad (20)$$

where R_{max} and R_{min} are the maximum and minimum radial heights of the entrance aperture, respectively. If only the third-order terms A_3 and A_4 are included, the equation represents third-order field curvature.¹⁰ In case of grazing incidence telescopes the fifth-order term $(B_3 + B_4 + \lambda_2/2)$ is comparable to the third-order term. If the grazing angles are small, even the seventh-order term C_1 cannot be omitted.

Alternative equation for the shape of the best focal surface is given by Shealy.¹⁷ In Shealy's paper the RMS image radius is formally calculated from the TRA equations [Eqs. (6) and (7)].

45.5 ABERRATION BALANCING

In case of Wolter type 1 telescopes the aberration equations suggest that for the optimum design the primary and the secondary should be as close to each other as possible. Separating the primary and the secondary, increases the radial heights on the primary, and therefore the image size.

The best focal surface of the Wolter type 1 telescopes always curves toward the telescope. The largest term in the aberration coefficients is the sum of the inverse of the radii of curvatures $(1/R_1 + 1/R_2)$. For Wolter type 1 telescopes, the radii of curvatures are both negative and these quantities in the terms A_3 , A_4 , and B_3 cannot cancel each other.

In case of Wolter type 2 telescopes, R_1 is negative and R_2 is positive. The radius of curvature can be optimized. It turns out that for all the practical designs the R_2 is always smaller than R_1 and the radius of curvature of the best focal surface is negative and curving toward the telescope. The aberrations are optimized when the primary and secondary are as close to the primary-secondary surface intersection point as possible.

The field of curvature could be improved by moving the entrance aperture away from the telescope. The K parameter would get bigger since the length L of the telescope would increase. Having the entrance aperture far in front of the telescope may not be practical design. For example, the vignetting would increase rapidly as a function of the off-axis angle.

The aberration equations [Eqs. (6) and (7)] presented in this paper are derived in terms of conventional parameters. The equations are shown as functions of the entrance aperture coordinates using the formulation introduced by Cox.¹⁶ The TRA polynomials and the OPD-polynomial have also been derived as functions of the exit pupil coordinates¹² using the traditional formulation shown, for example, in *Handbook of Optics*.¹⁸ The derivation includes all the terms shown in this paper (as a function of entrance aperture coordinates). The expansions are valid for all the combinations of Wolter telescopes and also for all the combinations of normal incidence paraboloid-hyperboloid or paraboloid-ellipsoid telescopes.

45.6 ON-AXIS ABERRATIONS

Rigid-Body Motions

Rigid-body motions and low spatial frequency errors of the primary and secondary are the most important on-axis image aberrations. These errors typically degrade the on-axis resolution of the grazing incidence telescopes and limit the encircled energy performance of the telescopes. Glenn¹⁹ introduced an orthonormal set of Legendre-Fourier (L-F) polynomials for cylindrical mirrors which are used to describe the low-order errors of the primary and secondary mirrors. The L-F polynomials have been implemented in the optical surface analysis code (OSAC) ray trace code.²⁰

The TRA aberration expansions have been derived for the rigid body motions and low order L-F polynomials.²¹ The rigid-body motions of the primary and secondary mirrors are despace, decenter, tilt, and defocus errors. The TRA equations can be derived following the similar principle used in the derivation of TRA expansions of off-axis aberrations shown in Sec. 45.5.

The defocus term (Δz) can be expressed as^{11,21}

$$H_x = -\frac{\Delta z}{f} h \sin(\beta) \quad (21)$$

$$H_y = -\frac{\Delta z}{f} h \cos(\beta) \quad (22)$$

In the equations only the first-order term in radial height h is kept and the higher-order terms are omitted. The defocus terms are proportional to the radial height h . Therefore, this term can be used to optimize the off-axis aberration terms that are also proportional to $h \sin(\beta)$ - $h \cos(\beta)$ pair. This principle was used to find the best focal surface and the radius of curvature of the best focal surface.

Despace surface errors of the primary ($n = 1$) and secondary ($n = 2$) can be expressed as surface radial height errors Δh_n and surface axial errors Δz_n as

$$\Delta h_n = 0 \quad (23)$$

$$\Delta z_n = \text{constant} \quad (24)$$

The despace image terms of the primary mirror can be approximated by

$$H_x = -\Delta z_1 \frac{4fR_1^2}{h^3} \sin(\beta) \quad (25)$$

$$H_y = -\Delta z_1 \frac{4fR_1^2}{h^3} \cos(\beta) \quad (26)$$

The despace image terms of the secondary mirror are

$$H_x = \Delta z_2 \left(\frac{4fR_2^2}{h^3} + \frac{h}{f} \right) \sin(\beta) \quad (27)$$

$$H_y = \Delta z_2 \left(\frac{4fR_2^2}{h^3} + \frac{h}{f} \right) \cos(\beta) \quad (28)$$

The leading term in the expansions is now proportional to inverse of h^3 . The off-axis aberrations shown above do not have terms proportional to inverse of h . This is because the primary mirror and secondary mirror are stigmatic and confocal.

Decentering the mirror leads to radial height error of the primary ($n = 1$) and secondary ($n = 2$) that can be written in terms of radial error as

$$\Delta h_n = e_{01n} \cos(\beta) + f_{01n} \sin(\beta) \quad (29)$$

where e_{01n} and f_{01n} are the amount of decenter error. The approximate TRA aberrations of the decentered primary mirror for the cosine component are

$$H_x = e_{011} \frac{2fR_1}{h^2} \sin(2\beta) \quad (30)$$

$$H_y = e_{011} \frac{2fR_1}{h^2} \cos(2\beta) \quad (31)$$

The equations of decentered secondary are similar

$$H_x = -e_{012} \frac{2fR_1}{h^2} \sin(2\beta) \quad (32)$$

$$H_y = e_{012} \left[1 - \frac{2fR_1}{h^2} \cos(2\beta) \right] \quad (33)$$

The tilt error of the mirrors can also be expressed in terms of radial height error and axial translation of the primary ($n = 1$) and secondary ($n = 2$). The radial and axial errors are

$$\Delta h_n = z_n (E_{11n} \cos(\beta) + F_{11n} \sin(\beta)) \quad (34)$$

$$\Delta z_n = -h_n (E_{11n} \cos(\beta) + F_{11n} \sin(\beta)) \quad (35)$$

where E_{11n} and F_{11n} are the amount of tilt error in radians when the mirror is rotated about x axis and y axis, respectively. The approximate expansions of the primary mirror are

$$H_x = -E_{111} f \left(\frac{h_{10}}{h} \right)^2 \sin(2\beta) \quad (36)$$

$$H_y = -E_{111} f \left[1 + \left(\frac{h_{10}}{h} \right)^2 \cos(2\beta) \right] \quad (37)$$

where h_{10} is the radial height of the primary mirror at the center of the mirror. The equations of the secondary mirror are

$$H_x = E_{112} \frac{f h_{10} h_{20}}{h^2} \sin(2\beta) \quad (38)$$

$$H_y = E_{112} f \frac{h_{10}}{h_{20}} \left[1 + \left(\frac{h_{10}}{h} \right)^2 \cos(2\beta) \right] \quad (39)$$

where h_{20} is the radial height of the secondary mirror at the center of the mirror.

The $\sin(2\beta)$ - $\cos(2\beta)$ relationship of the primary and secondary components shown in the equations for the decenter [Eqs.(30) to (33)] and tilt [Eqs. (36) to (39)] errors is typical of coma. Note that the components are now inversely proportional to h^2 . The off-axis coma term is directly proportional to h^2 .

Axial and Circumferential Slope Errors and TRA Equations

Assuming the primary mirror figure error Δh_1 and the secondary mirror figure error Δh_2 are known as functions of the surface axial coordinate and circumferential coordinate. Then, if the grazing angles of the mirrors are small, the on-axis aberrations can be evaluated easily from the surface slope errors. The TRA equations of the primary mirror of the axial surface errors are

$$H_x = -2f \frac{\partial \Delta h_1}{\partial z_1} \sin(\beta) \quad (40)$$

$$H_y = -2f \frac{\partial \Delta h_1}{\partial z_1} \cos(\beta) \quad (41)$$

where $\partial\Delta h_1/\partial z_1$ is its axial slope error. The secondary mirror TRA equations for the axial slope errors are

$$H_x = 2q \frac{\partial\Delta h_2}{\partial z_2} \sin(\beta) \quad (42)$$

$$H_y = 2q \frac{\partial\Delta h_2}{\partial z_2} \cos(\beta) \quad (43)$$

where q is the distance from the secondary to the telescope focus and $\partial\Delta h_2/\partial z_2$ is the axial slope error of the secondary. The q variable can be approximated by q_0 that is the distance from the center point of the surface to the telescope focus.

The approximated TRA equations for the primary mirror circumferential slope errors are

$$H_x = f \sin(2I_1) \frac{1}{h} \frac{\partial\Delta h_1}{\partial\beta} \cos(\beta) \quad (44)$$

$$H_y = -f \sin(2I_1) \frac{1}{h} \frac{\partial\Delta h_1}{\partial\beta} \sin(\beta) \quad (45)$$

where I_1 is the grazing angle on the primary and $\partial\Delta h_1/(h\partial\beta)$ is the circumferential slope error of the primary mirror. The approximate TRA equations of the secondary for the circumferential slope error are

$$H_x = -2f \sin(I_2) \frac{h_2}{h} \frac{\partial\Delta h_2}{\partial\beta} \cos(\beta) \quad (46)$$

$$H_y = 2f \sin(I_2) \frac{h_2}{h} \frac{\partial\Delta h_2}{\partial\beta} \sin(\beta) \quad (47)$$

where I_2 is the grazing angle of the surface, h_2 is the radial height of the surface, and $\partial\Delta h_2/h_2\partial\beta$ is the circumferential slope of the surface. The variables I_2 and h_2 can be approximated by the grazing angle and radial height at the midpoint on the surface.

The TRA equations of the circumferential errors include the grazing angle I_1 or I_2 . Typically, the grazing angles are small (0.5 to few degrees). Therefore, the circumferential errors have miniscule effect on the image at the focal plane compared to the axial slope errors and, in many cases, can be ignored.

45.7 REFERENCES

1. H. Wolter, "Mirror Systems with Glancing Incidence as Image-Producing Optics for X-Rays," *Ann. Phys.* **10**:94 (1952).
2. T. T. Saha, "General Surface Equations for Glancing Incidence Telescopes," *Appl. Opt.* **26**:658–663 (1987).
3. H. Wolter, "Generalized Schwarzschild Systems of Mirrors with Glancing Reflection as Optical System for X-Rays," *Ann. Phys.* **10**:286 (1952).
4. C. Burrows, R. Burg, and R. Giacconi, "Optimal Grazing Incidence Optics and Its Applications to Wide-Field X-Ray Imaging," *Astrophys. J.* **392**:760–765 (1992).
5. J. E. Harvey, A. Krywonos, P. L. Thompson, and T. T. Saha, "Grazing Incidence Hyperboloid-Hyperboloid Designs for Wide-Field X-Ray Imaging Applications," *Appl. Opt.* **40**:136–144 (2001).
6. L. VanSpeybroeck and R. Chase, "Design Parameters of Paraboloid-Hyperboloid Telescopes for X-Ray Astronomy," *Appl. Opt.* **11**:440 (1972).
7. H. Wolter, "Estimation of Image Aberrations for X-Ray Telescopes," *Opt. Acta* **18**:425 (1971).
8. W. Werner, "Imaging Properties of Wolter Type 1 X-Ray Telescopes," *Appl. Opt.* **16**:764 (1977).
9. C. Winkler and D. Korsch, "Primary Aberrations for Grazing Incidence," *Appl. Opt.* **16**:2464 (1977).
10. T. T. Saha, "Transverse Ray Aberrations of Paraboloid-Hyperboloid Telescopes," *Appl. Opt.* **24**:1856–1863 (1985).

11. T. T. Saha, "Transverse Ray Aberrations of Wolter Type 1 Telescopes," *Proc. SPIE* **640**:10–19 (1986).
12. T. T. Saha, "Aberrations for Grazing Incidence Telescopes," *Appl. Opt.* **27**:1492–1498 (1988).
13. H. P. Brueggemann, *Conic Mirrors*, Focal Press, London, 1968.
14. J. D. Mangus, "Optical Design of Glancing Incidence XUV Telescopes," *Appl. Opt.* **9**:1019 (1970).
15. T. T. Saha, "A Generalized Sine Condition and Performance Comparison of Wolter Type 2 and Wolter-Schwarzschild Telescopes," *Proc. SPIE* **444**:112–117 (1984).
16. A. Cox, *A System of Optical Design*, Focal Press, London, 1964.
17. D. L. Shealy and T. T. Saha, "Formula for the RMS Blur Circle Radius of Wolter Telescopes Based on Aberration Theory," *Appl. Opt.* **29**:2433–2439 (1990).
18. W. G. Driscoll, ed., *Handbook of Optics*, McGraw-Hill, New York, 1978.
19. P. Glenn, "Set of Orthonormal Surface Error Descriptors for Near-Cylindrical Optics," *Opt. Eng.* **23**:384–390 (1984).
20. R. J. Noll, P. Glenn, and J. F. Osantowski, "Optical Surface Analysis Code (OSAC)," *Proc. SPIE* **362**:78–85 (1983).
21. T. T. Saha, "Image Defects from Surface and Alignment Errors in Grazing Incidence Telescopes," *Opt. Eng.* **29**:1296–1305 (1990).

This page intentionally left blank.

DO NOT DUPLICATE

Peter Z. Takacs

*Brookhaven National Laboratory
Upton, New York*

46.1 GLOSSARY

d	sampling distance on surface
DFT	two-sided discrete Fourier transform
f	spatial frequency
$F[*]$	Fourier transform operator
H	transfer function in frequency domain
$K(m)$	bookkeeping factor
M	surface slope function
PSD	power spectral density function
Rq, Mq	root-mean-square (RMS) roughness, and slope error
$S_1(f)$	one-sided height profile power spectral density function
$S_1'(f)$	one-sided slope profile power spectral density function
$W(n)$	window function in spatial domain
w	sensor pixel width projected onto surface
$Z(x)$	surface height profile
$z(x)$	residual surface height profile after detrending

46.2 INTRODUCTION

The earliest use of mirrors for producing images with x rays dates from the work of Kirkpatrick and Baez in 1948 who used an orthogonal pair of shallow cylindrical mirrors to produce a focal spot.¹ Interest in x-ray mirrorfabrication technology grew rapidly in the 1970s and 1980s as a result of several factors: the initiation of a number of space-based x-ray telescope projects, the development of diamond machining of precision optical components, the development of bright laboratory x-ray sources for x-ray microscopy, and the development of dedicated synchrotron radiation source

facilities with many available beam lines. Owing to the nature of grazing incidence x-ray optics, conventional interferometric techniques for surface figure measurement are not easily employed in the testing of these aspheric surfaces. Surface profilometry, either contact or noncontact, has been the method of choice for measuring the figure of Wolter telescopes, x-ray microprobe optics, and synchrotron radiation mirrors. This has necessitated the development of specialized instruments with measurement capabilities far beyond those of the typical coordinate-measuring machine (CMM) found in precision machining shops. Grazing incidence x-ray mirror metrology requires measurement precision and accuracy in the nanometer range, while typical CMM machines are only able to achieve measurements at the micron level.

46.3 SURFACE FINISH METROLOGY

Surface roughness reduction has been a major consideration in the use of grazing incidence optics in x-ray telescope applications^{2,3} and in synchrotron beam line instrumentation.⁴ Optical systems used to reflect x rays at extreme grazing incidence angles are often made from far off-axis segments of cylinders, paraboloids, and toroids. These aspheric surfaces usually have a long tangential radius and a short sagittal radius, which necessitates fabrication processes that depart from conventional optical polishing techniques. Conventional full-contact pitch lap polishing works well for flat and spherical surfaces, but is difficult to achieve for aspheric surfaces. Novel figuring techniques, such as single-point diamond turning and ductile grinding, leave large microroughness levels that need to be reduced by polishing.

Measuring the roughness of aspheric optics at the sub-nanometer level was difficult before the advent of noncontact optical profilers. The collection of papers by Bennett⁵ chronicles the development of various surface roughness–measurement instruments and techniques over the past several decades. Of particular relevance to x-ray optics are the noncontact optical-measurement methods. Early optical profilers were based around interference microscopes with images recorded on film or by eye. Roughness was estimated by observing irregularities in the shape of multiple-beam interference fringes, as in the fringes of equal chromatic order (FECO) technique⁶ or by analyzing differential interference contrast (Nomarski) microscope images.⁷ Contact methods were more quantitative. Stylus profilers, such as the Talystep and Dektak, used a small radius diamond stylus that was dragged across the surface to record the surface topography. Analog data was recorded on a chart recorder. The stylus profilers were difficult to use for x-ray optics because they could only handle small samples and were not suited for measuring the surface of large optics. There was always the risk of damaging the surface of an optic coated with a soft metal-reflecting layer with the diamond stylus. Methods were developed to make replicas of surface roughness with collodion films that could then be analyzed in electron microscopes or by stylus instruments. This hybrid measurement technique was not very repeatable and suffered from accuracy problems. Sommargren developed a noncontact optical scanning profiler with Angstrom-level accuracy that utilized a precision rotary stage to scan a probe beam around a central stationary reference beam.⁸ The first noncontact optical profiling instrument based upon a computer-controlled phase-measuring interferometer (PMI) was developed by Koliopoulos^{9,10} and became available as a commercial instrument in the early 1980s. Similar instruments are now available from several sources.

Digital optical profilers revolutionized surface roughness-measurement technology for x-ray optics. Since they were true noncontact measurements, it was possible to configure the microscope-based PMI instruments to measure the entire surface of full-size mirrors. Manufacturers were able to get rapid feedback about their polishing processes and were able to make significant improvements in producing low-scatter superpolished surfaces for x-ray optics. In parallel with these instrument developments, advances in x-ray scattering theory were made to connect the performance of grazing incidence optics at x-ray wavelengths with roughness measurements made by the various surface-profiling techniques at visible wavelengths. A critical part of connecting the surface-finish measurements to the functional performance at x-ray wavelengths was the understanding of the measurement bandwidth of each profiling technique (frequency footprint) and how it is related to scattering at grazing incidence through the power spectral density function.^{11–17} (See Chap. 44 in this volume.)

46.4 SURFACE FIGURE METROLOGY

Full-Aperture Techniques

Conventional optical interferometry of grazing incidence aspheric optics is difficult to implement. Speer used a Linnik interferometer arrangement to visualize the figure error on toroidal surfaces at EUV wavelengths,¹⁸ but most interferometry techniques at grazing incidence on cylindrical aspheres result in poor spatial sampling of errors along the surface owing to the extreme foreshortening of the surface in the aperture. Grazing angle of incidence testing of long flat and large radius spherical mirrors can be done by using double-pass Fizeau techniques.¹⁹ This geometry allows one to foreshorten the long axis of a mirror to fit within the interferometer aperture that is usually 4 or 6 inches in diameter. Careful system calibration needs to be done with this test geometry to ensure that subtle systematic errors are corrected in order to achieve nanometer accuracy. Full-aperture measurements of x-ray telescope optics at x-ray wavelengths have been made at specialized test facilities, such as the X-Ray Calibration Facility at Marshall Space Flight Center that has a 518-m long source distance, and the PANTER facility at the Max-Planck-Institute near Munich. These are not interferometric tests, rather they measure directly the quality of the image of an x-ray point source.

Height Profilometry

Practical figure metrology methods for aspheric optics are generally based on some type of profilometry. Surface-profilometry instruments for grazing incidence metrology can be classified into two broad categories: height-measuring instruments and slope-measuring instruments, based on the fundamental measurement provided by the machine. Height-measuring instruments measure the distance between the test surface and a reference surface; slope-measuring instruments measure either the differential phase shift between two closely spaced probe beams or the angle of a reflected laser beam directly, and are usually based upon some variant of the optical lever principle.²⁰⁻²² A number of surface profilers with nanometer precision were originally developed for microroughness investigations. Their scan ranges were limited to a few millimeters. Bennett et al.²³ describe a number of commercially available mechanical stylus-based profiling instruments that have been modified to extend the measurement length to the 100- to 200-mm range. Carl Zeiss has developed a 3D coordinate-measuring machine, the M400, which utilizes a laser-based distance interferometer on all three axes with 10-nm positioning resolution.²⁴ This machine is capable of measuring slope errors with a repeatability of 0.1 arc sec (0.5 μ rad) over the 400 \times 600 \times 200 mm measurement volume. Other scanning stylus profilers were developed at the National Physical Laboratory²⁵ and at Cranfield Institute.²⁶ A noncontact fringe-scanning technique was developed at Perkin-Elmer to measure the distance between an aspheric test surface and a toroidal test plate. The fringe scanner measured the position of Fizeau fringes formed in the air gap between the two surfaces and converted the air-gap profile into a surface height profile.^{27,28} This instrument was used to measure the axial profile of the HEAO-2 and AXAF Wolter telescope mirrors.

Slope Profilometry

Slope-profilometry instruments generally operate on the principle of the optical lever, where an angular deviation of a probe beam is magnified and recorded with a high-gain amplifier. Jones chronicles the development of sensitive optical lever instruments capable of measuring nanoradian angular deflections.²² An optical lever device was developed by George Random to measure slope errors on x-ray optics several centimeters in length.^{29,30} The Random Devices Slope Scanner used a clever air-bearing slide mechanism to move the probe beam along a circular arc that matched the curvature of the part under test.³¹ This motion always kept the surface near the focus of the laser beam spot. The deflection of the reflected beam was measured with a resolution of a few tenths of an arc second by the differential signal from a pair of photodiodes. As DeCew notes, the limitation in the measurement accuracy was not in the optical system, but was in the tolerances of the mechanical translation system.

Makosh and coworkers developed a number of surface-profiling instruments at IBM in Germany based on the principle of differential heterodyne interferometry.³²⁻³⁴ Two probe beams generated by a Wollaston prism were focussed onto a test surface. The spot separation could be varied from 2 μm to 1.5 mm by changing the Wollaston prism or the microscope objective. The phase shift between the two spots was a measure of the local surface slope, and the slope was integrated to provide surface height roughness information. They point out that the differential height measurement is insensitive to mechanical vibrations and air turbulence and is suitable for long-distance measurements.

A slope-measuring system was developed by Ennos and Virdee at NPL that was based on the principle of laser autocollimation.^{35,36} An unfocused laser beam was reflected off a cylindrical asphere and an astigmatic line image was focused onto a position-sensitive quad cell detector so that the detector sensed the slope errors in the axial direction of the scan. An innovation in the Ennos and Virdee system was the incorporation of an optical reference beam to monitor changes in the pointing direction of the laser beam during the scan. The measurement precision was on the order of 0.5 μrad , corresponding to a displacement of the line focus by 0.1 μm on the detector, but beam pointing drifts of 1.5 μrad over the 10-min scan period needed to be corrected from the measured profiles. The instrument could reproducibly measure height profiles to within ± 5 nm over 16-mm travel lengths.

Bristow and Arackellian reported on the development of a slope-measuring profiler based upon a modification of a surface microroughness-profiling instrument.³⁷ The original was based on the principle of Nomarski differential interference contrast but was modified to act as a differential polarization interferometer.³⁸ The MP2000 had an intrinsic lateral resolution on the order of 1 μm , owing to the focal properties of the microscope objective used to focus the beam onto the surface. An astigmatic autofocus control system was added to the instrument to extend the depth of focus and allow the profiler to scan over longer distances on curved parts.³⁹ The basic instrument had a maximum scan range of 100 mm.

Glenn described a novel surface-profiling instrument that, unlike the height- and slope-measuring instruments, measured surface curvature directly on aspheric parts.^{40,41} The Bauer Model 100 used a laser pencil beam that was passed through a calcite prism to generate two spatially separated and orthogonally polarized beams. Both beams were directed onto the test surface by a steerable mirror so that normal incidence could be maintained over surface slope changes of up to $\pm 30^\circ$. The reflected beams were directed to separate quad-cell detectors, one of which could be translated to accommodate the average surface curvature. The differential measurement between each separate beam position was a direct measurement of the variation of the surface curvature away from the average. Because the instrument measures the intrinsic curvature of the surface, it is immune to vibration or tilt errors between the optical head and the part during the scan period. This instrument is particularly well-suited for azimuthal circularity measurements on x-ray telescope mirrors and mandrels.

Weingärtner, Schulz, Estler, and coworkers at the PTB in Braunschweig, Germany, have developed a number of precision figure-measuring techniques based on multiple-beam slope profilometry with large lateral shear distances.⁴²⁻⁴⁷ Their measurement techniques and analysis algorithms correct for errors in the translation stages and rigid-body motions of the test piece during the measurement. Schulz has extended this work to combine a scanning microinterferometer that measures distance at 16 discrete positions with an autocollimator to measure pitch error during the scan.^{48,49} A sketch of an autocollimator is shown in Fig. 1. The interferometer is scanned over the surface with a 0.2-mm step size, determined by the interferometer sensor resolution. The resultant height measurements are stitched together to produce a high-resolution surface profile with sub-nanometer uncertainty.

The long trace profiler (LTP) is optimized for measuring the surface figure of synchrotron beam line optics. Its design is based upon the principle of the pencil beam interferometer, originally developed by von Bieren in 1982.^{50,51} Qian designed an improved version of the beam-splitting system for the original LTP I at BNL that produces zero-optical path difference probe beams with a spacing that could be adjusted from complete overlap to several millimeters apart.⁵²⁻⁵⁴ The optical head design is quite simple, with no internal moving parts during the measurement to ensure long-term stability and maintain high accuracy. An advantage of the pencil beam interferometer over microscope-based systems is the effectively infinite depth of field of the probe beam, which relaxes the tolerances on test piece alignment and allows for rapid setup and testing. An LTP slope profile of a single-crystal silicon cylinder mirror is shown in Fig. 2. This mirror exhibits several significant defects and was returned to the manufacturer for rework.

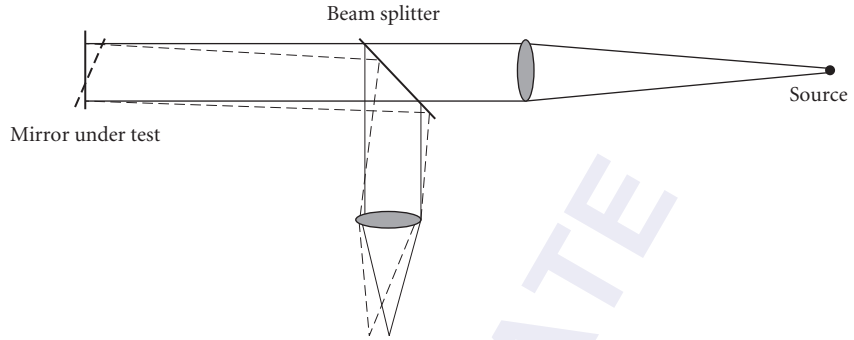


FIGURE 1 Sketch of a visible light autocollimator.

The LTP can be operated in two basic configurations. The standard configuration places the optical head on a linear translation stage and moves the entire optical head over the surface that is being measured. The alternative configuration places the stationary optical head off to the side with the test beam directed horizontally to a penta prism that sends the probe beam down to the test surface. The penta prism is scanned over the surface; the optical head remains fixed. This is the original measurement method proposed by von Bieren.⁵⁰ The nano-optic-measuring (NOM) machine at BESSY II also uses a scanning penta prism with a fixed autocollimator providing the main probe beam.⁵⁵ Each configuration has advantages and disadvantages. The moving optical head requires a high-quality translation stage to minimize pitch angle errors during the scan. A reference beam must be used to correct for residual pitch errors during the scan caused by the sag of the translation stage and irregularities in the motion. The moving penta prism, however, always redirects the incident beam by exactly 90° , independent of

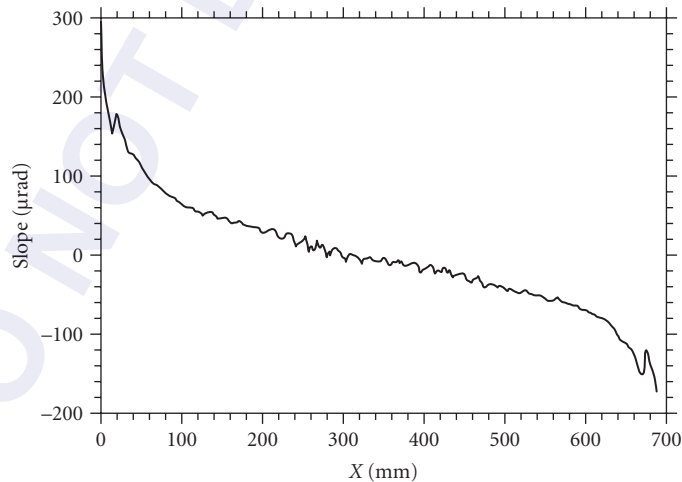


FIGURE 2 Slope profile measurement on a 700-mm long silicon cylinder mirror made with an LTP. Sampling step size is 2 mm. Mean has been subtracted from the data (detrend 0). Profile shows that the surface has an overall convex curvature (tilted profile down to the right) with significant edge roll-off (change in slope at each end). Also, a polishing defect with a 20-mm period is evident in the center of the surface. The slope profile emphasizes high-frequency surface defects. (See also color insert.)

the errors in the translation stage. Use of a reference beam is optional in this case.⁵⁶⁻⁶² Use of the penta prism also greatly relaxes the tolerances on the translation mechanism. The penta prism LTP can be configured to provide measurements in difficult environments, such as inside complete x-ray telescope cylinders^{60,63-66} and in situ at SR beam lines during actual operations.^{61,67}

46.5 PRACTICAL PROFILE ANALYSIS CONSIDERATIONS

Nomenclature

Most profiling instruments today produce a map of surface height or slope topography over a 2D area. Early profiling instruments, optical and contact, were restricted to single 1D line profiles. Some authors refer to surface profilometry over a 2D area as *3D profilometry* and over a line as *2D profilometry*, which introduces some ambiguity into the nomenclature. For purposes of clarity, 1D will refer to a linear scan over one line and 2D will refer to a scan over an area.

Detrending

Most x-ray mirrors are made today with extremely smooth surfaces by processes that produce isotropic roughness. In this case it is sufficient to consider surface statistics in one dimension to be descriptive of the entire surface. Most x-ray mirrors have overall figures that are either plano, spherical, or cylindrical, or are a far off-axis segment of a conic section. The current trend is to produce aspheric surfaces by bending substrates that are prefigured with a long-radius sphere or cylinder. In order to assess the underlying residual roughness, it is usually necessary to subtract the best-fit conic section from the measured data before computing statistical properties. The magnitude of the overall figure is usually many orders of magnitude greater than the residual surface roughness and needs to be removed by mathematical manipulation of the measured profile. The residual profile z can be viewed as the difference between the measured profile, Z_{raw} , and the deterministic figure:

$$z(x_n) = Z_{\text{raw}}(x_n) - [A + Bx_n + Cx_n^2 + Dx_n^3 + \dots] \quad (1)$$

where the profile is sampled at discrete points, x_n , separated by equal intervals, d , and the figure is described by some polynomial function of position. The process of removing the deterministic part of the measured profile is called *detrending* and we denote removing the mean as a “detrend0” or “D0” process, a first order detrend as a “detrend1” or “D1” process, and so forth. Most often, the detrending function is determined by a least-squares fit to the measured data. Any convenient polynomial function can be used to describe this function. In certain cases it may be more useful to do a straight-line detrend between the endpoints of the profile. The simple polynomial shown in Eq. (1) has terms that can be identified with various rigid body alignments: A = piston (DC) offset, B = tilt. The second order curvature term C is an intrinsic surface property, independent of body orientation. The third-order term D is related to the ellipticity of the surface, but its value depends on which part of the ellipse is being viewed. For instruments that measure slope M , directly, the slope space profiles can be detrended in a manner analogous to Eq. (1), but the identification of the polynomial coefficients with rigid-body parameters differs from the distance space identifications.

Simple polynomials are not orthogonal over the measured region, so the values of lower-order terms change when higher-order terms are added to the detrending. Other orthogonal polynomial functions, such as sinusoids, Chebyshev, and Legendre may be better suited for linear traces, while Zernike polynomials may be better suited for 2D area detrending, depending upon the application. The simple polynomial, however, has a simple connection to the radius of curvature of the surface when $R \gg L$, (L = the trace length): $R = 1/2C$. For the other polynomials, some combination of two or more coefficients is required to estimate the radius. Note that this discussion has involved fitting a 1D

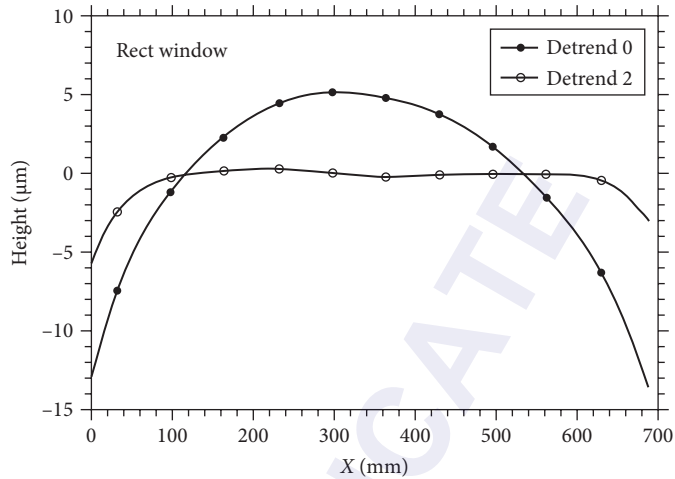


FIGURE 3 Height profile calculated by integrating the slope profile of Fig. 1. Solid circles: mean height subtracted (detrend 0); open circles: second order polynomial subtracted (detrend 2). The radius of curvature extracted from the second order term coefficient is 3.572 km. The residual profile shows that the surface has a “kink” in the center that separates it into two distinct segments with slightly different slopes. This low frequency defect is not evident in the slope profile of Fig. 1. (See also color insert.)

polynomial to a linear profile that is only a function of one coordinate. Equation (1) can be generalized to two variables when fitting a 2D polynomial to the 2D surface. However, when a 2D polynomial is fit to an area, the individual row or column profiles will generally not be fully 1D detrended. In computing 1D statistical quantities from 2D data, one must be aware that the extracted profiles may need further detrending.

A typical detrending example is shown in Fig. 3. The detrend0 height profile has only had the piston (DC) term removed. The radius of curvature can be extracted by further detrending with a second order polynomial function ($D2$). The radius of curvature of the central region of the surface between 100 and 600 mm is derived from the coefficient of the x^2 term, giving $R = 3.572$ km as the best-fit radius. When this fit is subtracted from the $D0$ curve, the edge roll-off at each end of the mirror becomes visible. The region over which the detrending polynomial is applied, and over which the statistical roughness and slope-error parameters are computed, depends upon how the mirror will be used in a synchrotron beamline and on the clear aperture of the illuminated region. Surface errors outside the clear aperture can generally be ignored and should be excluded from the statistics.

Power Spectral Density Function

Various statistical quantities can be computed from profiler data. Standards exist that define the various surface texture parameters that can be derived from surface profile measurement.^{68,69} For high-performance x-ray optics, the most useful descriptor of surface roughness is the power spectral density function (PSD). Church and collaborators have shown that the PSD computed from normal incidence visible light profilometry measurements can be used to predict the performance of grazing incidence optics at x-ray wavelengths in SR beam lines. (See Chap. 8 in Vol. I and Chap. 44 of this volume for a discussion on the connection between PSD and optical scatter.) A detailed description of the definition of the PSD function and related issues involving calculations from sampled data can be found in SEMI standard MF1811-0704⁷⁰ and in the volume on scattered light by Stover.⁷¹

The following paragraphs highlight parameters derived from profile measurements important for x-ray optic characterization.

The “periodogram estimator” for a 1D profile $Z(n)$ is used to define the 1D PSD function^{70–72}

$$S_1(m) = \frac{2d}{N} |\text{DFT}(m)|^2 \cdot K(m) \quad (2)$$

where DFT is the *two-sided discrete Fourier transform* of the N real data points:

$$\text{DFT}(m) = \sum_{n=1}^N e^{i2\pi \frac{(n-1)(m-1)}{N}} \cdot Z(n) \cdot W(n) \quad (3)$$

d = sampling period in one direction, N is the number of sampled points, m is the spatial frequency index, where $f_m = (m-1)/Nd$ is the value of the spatial frequency, $W(n)$ is a window function (see following section), and $K(m)$ is a bookkeeping factor to ensure that Parseval’s theorem is satisfied, that is, that the variance of the distance space profile is equal to the variance of the frequency space spectrum (see the following section).

A number of useful bandwidth-limited statistical parameters and functions can be derived from the PSD function in frequency space. The RMS roughness, Rq , over a given spatial frequency bandwidth is given by

$$Rq^2(\text{low}, \text{hi}) = \frac{1}{Nd} \sum_{m=\text{low}}^{\text{hi}} S_1(m) \quad (4)$$

where low and hi are the indices corresponding to the desired spatial frequency range. When the frequency indices correspond to the full bandwidth, this number is identically equal to the RMS roughness computed in distance space (Parseval’s theorem). Measurements from instruments that have different bandwidths can be compared by restricting the calculation of RMS roughness to the bandwidth that is common to both instruments.

The “Bookkeeping Factor” for the PSD

Most discrete Fourier transform algorithms used in computing libraries today can efficiently calculate the transform for an arbitrary number of real or complex data points. See the section in *Numerical Recipes*⁷³ for a discussion of practical considerations in Fourier transform calculation. The form of the DFT shown in Eq. (3) is known as the “two-sided” DFT, since the m -index runs over all N frequency terms. But since we are starting with N real numbers, the $|\text{DFT}|^2$ will be symmetric about the Nyquist frequency, defined to be $f_{\text{Ny}} = 1/2d$ and will consist of only $N/2$ independent numbers (assuming N is even—see the following discussion). The numbers beyond the Nyquist frequency correspond to the negative frequencies in the input signal and have the same amplitude as the positive frequencies when the input data are real numbers. We can then effectively ignore the numbers beyond the Nyquist frequency and only need consider terms over one half of the N points. When we do this, we restrict m to range over approximately $N/2$ points. However, the missing power in the negative frequencies needs to be added to the positive frequency terms, hence the factor of 2 in the numerator of Eq. (2).

An additional bookkeeping consideration depends on whether N is an odd or even number. In order to ensure that the total power, according to Parseval’s theorem, is satisfied in both distance and frequency space, careful consideration must be given to the terms at the extremes of the frequency interval. In our realization of the Fourier transform, the DC term occurs at frequency index $m = 1$. The fundamental frequency is at $m = 2$ with a value of $1/Nd = 1/L$, where L is the total trace length. The difficulty arises in specifying the index of the Nyquist frequency for a set of N discrete sampled points that can be even or odd. When N is even, the Nyquist frequency term occurs at a single index position, $m = (N/2) + 1$. When N is odd, the Nyquist frequency power is split between adjacent points

$m = 1 + (N - 1)/2$ and $m = 2 + (N - 1)/2$. When the terms in the two-sided spectrum above the Nyquist frequency are discarded and the remaining terms are doubled to generate the one-sided spectrum, the resultant DC term for both N even and odd must be reduced by $1/2$, and the Nyquist term for the $N =$ even case must also be reduced by $1/2$. Hence the bookkeeping factor, $K(m)$:

$$\begin{aligned} \text{for } N \text{ even} \quad K(m) &= \begin{cases} 1/2 & \text{for } m = 1 \text{ or } m = (N/2) + 1 \\ 1 & \text{for all other } m \end{cases} & \text{where } m = 1, 2, \dots, (N/2) + 1 \\ \\ \text{for } N \text{ odd} \quad K(m) &= \begin{cases} 1/2 & \text{for } m = 1 \\ 1 & \text{for all other } m \end{cases} & \text{where } m = 1, 2, \dots, \left(\frac{N-1}{2}\right) + 1 \end{aligned}$$

Windowing

Proper preparation of raw input profile data is necessary to obtain meaningful results from statistical calculations. Detrending is usually required to remove the gross figure terms from the measured data so that the underlying surface roughness can be seen. Even after the gross figure has been removed by detrending, edge discontinuities in the residual profile can introduce spurious power into the spectrum and hide the true nature of the underlying roughness spectrum. Methods have been developed in the signal-processing literature to deal with this “leakage” problem.⁷³⁻⁷⁵ These methods are collectively known as “prewhitening” techniques. A simple method to preprocess the data is known as “windowing.” A window function, $W(n)$, in distance space is applied to the residual profile before the PSD is computed to smooth the spectrum somewhat and to minimize the spectral leakage from edge discontinuities. The window function used should be normalized to unity so as not to introduce any additional scale factors that would distort the magnitude of the spectrum. Most common window functions tend to enhance the lowest two or three spatial frequencies which are generally related to the deterministic figure components, but the power in these frequencies should already have been minimized by the detrending process. We prefer to use the Blackman window for processing smooth-surface residual profile data. The Blackman window, normalized to unit area, is defined as

$$W(n) = \sqrt{\frac{2}{1523}} [21 - 25\cos(2\pi(n-1)/N) + 4\cos(4\pi(n-1)/N)] \quad (6)$$

The Blackman window applied to the detrended profiles of Fig. 3 are shown in Fig. 4. One can see that this window function forces the edge discontinuities to go to zero. This has the beneficial effect of reducing the discontinuity between the derivatives of the profile at each end point. Discontinuities at the end points, or large spurious spikes in the data, produce large high-frequency ringing effects in the DFT coefficients that distort the underlying surface spectrum. Figure 5 shows the results of computing the PSD for the windowed and unwindowed profiles of Figs. 3 and 4. Application of the window function effectively eliminates the contamination of the high frequency content in each profile, even for significant edge discontinuities.

Instrument Transfer Function Effects

The ideal surface-profiling instrument has a unity transfer function response over an infinite spatial period bandwidth. In other words, the measurement does not distort the intrinsic surface properties. In the real world, however, all measuring instruments have a response function that varies over

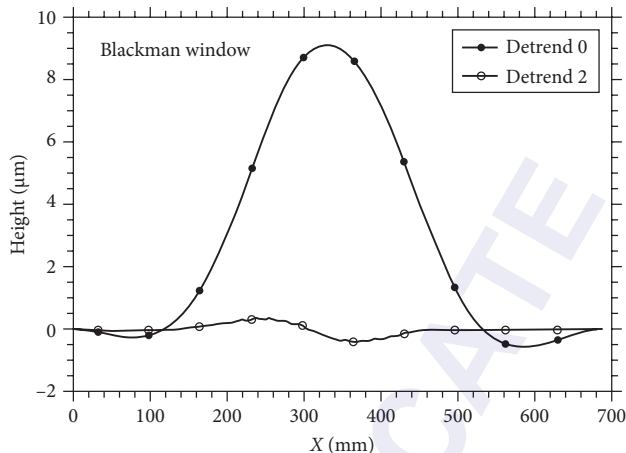


FIGURE 4 The height profiles of Fig. 3 with a Blackman window applied. The edge discontinuities are minimized by this function. Although the shape of the profile is distorted, the average statistical properties of the underlying function are not changed. (See also color insert.)

a limited bandwidth. For optical profiling instruments, the transfer function is limited mainly by the numerical aperture of the objective and the sampling properties of the detector. The transfer function of an unobscured objective for incoherent illumination is given by⁷⁶

$$H_{\text{obj}} = \frac{\pi}{2} (\cos^{-1} \Omega - \Omega \sqrt{1 - \Omega^2}) \quad (7)$$

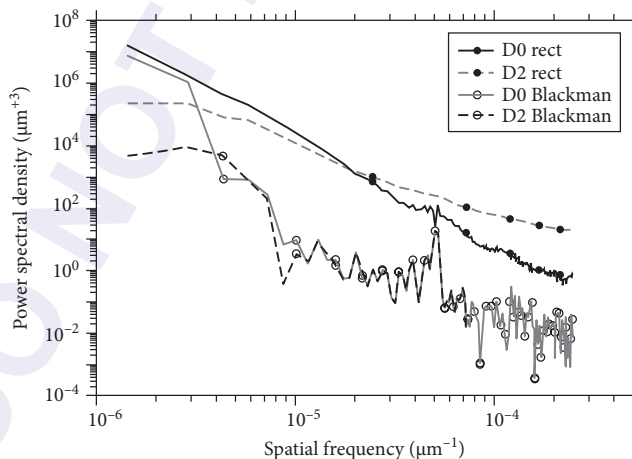


FIGURE 5 PSD curves computed from the four profiles in Figs. 3 and 4. The two upper curves from the unwinded data show severe contamination effects due to the strong edge discontinuities that introduce spurious power into all frequencies. The lower curves show how the Blackman filter eliminates the discontinuity contamination, allowing the underlying surface spectral characteristics to become visible. (See also color insert.)

where $\Omega = \lambda f/2NA$ and NA is the numerical aperture of the objective lens. The cutoff frequency where H_{obj} goes to zero is at $f_{\text{cutoff}} = 2NA/\lambda$; the transfer function is zero for all higher frequencies. This function may need to be modified if the objective lens in the profiling instrument contains a central obscuration, as in the case of a Mirau objective. The transfer function of an ideal 1D linear sensor array is given by

$$H_{\text{arr}} = \frac{\sin(\pi wf)}{\pi wf} \quad (8)$$

where w is the pixel width in one dimension. In most cases the pixel width is equal to the sampling distance, $w = d$, and the attenuation at the Nyquist frequency, $f = 1/2d$, is $\sin(\pi/2)/(\pi/2) = 0.63$. Since the transfer function of the array is still above zero beyond the Nyquist frequency, and the optical cutoff frequency is also generally beyond the Nyquist frequency, the measured spectral density will usually contain aliasing from frequencies beyond the Nyquist. This usually results in a flattening of the measured spectrum at the highest spatial frequencies. This effect is obvious when measurements are made with more than one magnification objective with an optical profiler. The combined optical and pixel sampling transfer functions can be used as an inverse filter to restore the high-frequency content of the measured spectrum and give a better estimate of the intrinsic surface power spectral density:

$$\bar{S} = S_{\text{meas}} \cdot H_{\text{obj}}^{-1} H_{\text{arr}}^{-1} \quad (9)$$

There are practical limitations of this technique, such as the need to avoid singularities in the inverse filter that cause the resultant PSD estimate to blow up. In practice, one must impose a practical cutoff frequency before the Nyquist frequency or before the zero in the sampling function to avoid significant distortion of the spectrum by noise and aliasing.⁷⁷ Other complications to this simple restoration filter approach occur when, unknown to the user, the signals from the sensor are preprocessed inside the measuring instrument, such as when adjacent rows of pixels are averaged to smooth out amplifier gain differences in 2D array sensors. In this case, the H_{arr} filter needs to be modified with correction factors.^{78,79}

Slope Measurement Analysis

Surface slope profiles can be derived from surface height measurements by finite difference calculation:

$$M(x_n) = \frac{1}{d} [Z(x_{n+1}) - Z(x_n)] \quad n = 1, 2, \dots, (N - 1) \quad (10)$$

Note that there is always one less slope point than there are height points in this calculation. Conversely, height profiles can be generated from slope profiles by numerical integration:

$$Z(x_n) = Z_0 + d \sum_{i=1}^n M(x_i) \quad n = 1, 2, \dots, N \quad (11)$$

Note that this latter calculation involves an arbitrary constant, Z_0 , which corresponds to a rigid body piston orientation of the part. Surface slope and height profiles can also be computed by Fourier differentiation and integration. The formal relationships between the slope and height transforms are

$$\begin{aligned} F[M(x)]_m &= i2\pi f_m F[z(x)]_m \\ F[z(x)]_m &= \frac{-i}{2\pi f_m} F[M(x)]_m \end{aligned} \quad (12)$$

where $F[*]$ is the Fourier transform operator. Care must be exercised in implementing these expressions with a DFT to ensure that the frequency terms f_m encompass both the negative and positive frequencies around zero and multiply the corresponding transformed height and slope numbers. Also, the DC term at $f = 0$ must be excluded from the denominator in the slope-to-height transform.

Of particular interest to users of x-ray optics is the slope PSD function, S' , which is related to the height PSD by

$$S'_1(m) = (2\pi f_m)^2 S_1(m) \quad (13)$$

The slope S' spectrum can also be calculated directly from data that is generated by profilers that measure slope by substituting the slope data, $M(n)$, for the $Z(n)$ data in Eq. (3). Bandwidth-limited RMS slope numbers, M_q , can be calculated from the slope S' function in a manner analogous to Eq. (4). Conversely, a measured slope profile PSD curve can be used to predict the height spectrum by the inverse of the above expression:

$$S_1(m) = \frac{1}{(2\pi f_m)^2} S'_1(m) \quad (14)$$

Care must be exercised in this case to exclude the DC term ($f_m = 0$) from the calculation, as it will cause the result to be nonsense.

46.6 REFERENCES

1. P. Kirkpatrick and A. V. Baez, "Formation of Optical Images by X-Rays," *J. Opt. Soc. Am.* **38**:766–774 (1948).
2. B. Aschenbach, "Design, Construction, and Performance of the Rosat High-Resolution X-Ray Mirror Assembly," *Appl. Opt.* **27**(8):1404–1413 (1988).
3. L. Van Speybroeck, *AXAF Mirror Fabrication—Element Fabrication*, http://hea-www.harvard.edu/asc/news_02/subsection3_4_3.html (1994).
4. M. Howells and P. Z. Takacs, "Use of Diamond Turned Mirrors for Synchrotron Radiation," *Nucl. Instrum. Methods* **195**:251–257 (1982).
5. J. M. Bennett, *Surface Finish and Its Measurement*, Washington, D.C.: Optical Society of America, p. 918, 1992.
6. S. Tolansky, *Multiple-Beam Interferometry of Surfaces and Films*, Oxford: Clarendon Press, p. 187, 1948.
7. J. S. Hartman, R. L. Gordon, and D. L. Lessor, "Quantitative Surface Topography Determination by Nomarski Reflection Microscopy. 2: Microscope Modification, Calibration, and Planar Sample Experiments," *Appl. Opt.* **19**(17):2998–3009 (1980).
8. G. E. Sommargren, "Optical Heterodyne Profilometry," *Appl. Opt.* **20**:610 (1981).
9. C. L. Koliopoulos and J. C. Wyant, "Profilometer for Diamond-Turned Optics Using a Phase-Shifting Interferometer," *J. Opt. Soc. Am.* **70**(12):1591–1591 (1980).
10. B. Bhushan, J. C. Wyant, and C. L. Koliopoulos, "Measurement of Surface Topography of Magnetic Tapes by Mirau Interferometry," *Appl. Opt.* **24**:1489–1497 (1985).
11. E. L. Church, "The Precision Measurement and Characterization of Surface Finish," *Precision Surface Metrology, Proc. SPIE* **429**, J. C. Wyant, ed., pp. 86–95 (1983).
12. E. L. Church and P. Z. Takacs, "Use of an Optical-Profilometer Instrument for the Measurement of the Figure and Finish of Optical-Quality Surfaces," *WEAR* **109**:241–257 (1986).
13. E. L. Church and P. Z. Takacs, "Statistical and Signal Processing Concepts in Surface Metrology," *Optical Manufacturing, Testing, and Aspheric Optics, Proc. SPIE* **645**:107–115 (1986).
14. E. L. Church and P. Z. Takacs, "The Interpretation of Glancing-Incidence Scattering Measurements," *Grazing Incidence Optics, Proc. SPIE* **640**:126–133 (1986).
15. E. L. Church and P. Z. Takacs, "Spectral and Parameter Estimation Arising in the Metrology of High-Performance Mirror Surfaces," *Proc. of ICASSP '86*, S. Saito, ed., pp. 185–188 (1986).

16. E. L. Church and P. Z. Takacs, "Effects of the Optical Transfer Function in Surface Profile Measurements," *Surface Characterization and Testing II*, J. E. Grievenkamp and M. Young, eds., *Proc. SPIE* **1164**:46–59 (1989).
17. E. L. Church and P. Z. Takacs, "Prediction of Mirror Performance from Laboratory Measurements," in *X-ray/EUV Optics for Astronomy and Microscopy*, R. Hoover, ed., *Proc. SPIE* **1160**:323–336 (1989).
18. R. J. Speer, M. Chrisp, D. Turner, S. Mrowka, and K. Tregidgo, "Grazing Incidence Interferometry: The Use of the Linnik Interferometer for Testing Image-Forming Reflection Systems," *Appl. Opt.* **18**(12):2003–2012 (1979).
19. P. H. Langenbeck, "New Developments in Interferometry—VI. Multipass Interferometry," *Appl. Opt.* **8**(3):543–552 (1969).
20. W. H. Wilson and T. D. Eps, *Proc. Phys. Soc.* **32**:326 (1920).
21. R. V. Jones, "Some Points in the Design of Optical Levers," *Proc. Phys. Soc. B* **64**:469–482 (1951).
22. R. V. Jones, *Instruments and Experiences: Papers on Measurement and Instrument Design*, Wiley Series in Measurement Science and Technology, P. H. Sydenham, ed., John Wiley & Sons, Ltd., p. 485, 1988.
23. J. M. Bennett, V. Elings, and K. Kjoller, "Recent Developments in Profiling Optical Surfaces," *Appl. Opt.* **32**(19):3442–3447 (1993).
24. K. Becker and E. Heynacher, "M400—A Coordinate Measuring Machine with 10 nm Resolution," *In-Process Optical Metrology for Precision Machining*, *Proc. SPIE* **802**:209–216 (1987).
25. M. Stedman and V. W. Stanley, "Machine for the Rapid and Accurate Measurement of Profile," *Proc. SPIE* **163**:99–102 (1979).
26. W. J. Wills-Moren and P. B. Leadbeater, "Stylus Profilometry of Large Optics," presented at San Diego, Calif., *Advanced Optical Manufacturing and Testing*, *Proc. SPIE* **1333**:183–194 (1990).
27. P. S. Young, "Fabrication of the High-Resolution Mirror Assembly for the HEAO-2 X-Ray Telescope," *Proc. SPIE* **184**:131–138 (1979).
28. A. Sarnik and P. Glenn, "Mirror Figure Characterization and Analysis for the Advanced X-Ray Astrophysics Facility/Technology Mirror Assembly (AXAF/TMA) X-Ray Telescope," *Grazing Incidence Optics for Astronomical and Laboratory Applications*, *Proc. SPIE* **830**, S. Bowyer and J. C. Green, eds., pp. 29–36 (1987).
29. J. K. Silk, *A Grazing Incidence Microscope for X-Ray Imaging Applications*, *Annals of the New York Academy of Sciences* **342**:116–129 (1980).
30. R. H. Price, "X-Ray Microscopy Using Grazing Incidence Reflection Optics," *Low Energy X-Ray Diagnostics*, AIP Conf. Proc. **75**, D. T. Attwood and B. L. Henke, eds., pp. 189–199 (1981).
31. A. E. DeCew, Jr. and R. W. Wagner, "An Optical Lever for the Metrology of Grazing Incidence Optics," *Optical Manufacturing, Testing, and Aspheric Optics*, *Proc. SPIE* **645**:127–132 (1986).
32. G. Makosch and B. Solf, "Surface Profiling by Electro-Optical Phase Measurements," *High Resolution Soft X-Ray Optics*, *Proc. SPIE* **316**, E. Spiller, ed., pp. 42–53 (1981).
33. G. Makosch and B. Drollinger, "Surface Profile Measurement with a Scanning Differential AC Interferometer," *Appl. Opt.* **23**(24):4544–4553 (1984).
34. G. Makosch, "LASSI—a Scanning Differential AC Interferometer for Surface Profile and Roughness Measurement," *Surface Measurement and Characterization*, *Proc. SPIE* **1009**, J. M. Bennett, ed., pp. 244–253 (1988).
35. A. E. Ennos and M. S. Virdee, "High Accuracy Profile Measurement of Quasi-Conical Mirror Surfaces by Laser Autocollimation," *Prec. Eng.* **4**:5–9 (1982).
36. A. E. Ennos and M. V. Virdee, "Precision Measurement of Surface Form by Laser Autocollimation," *Industrial Applications of Laser Technology*, *Proc. SPIE* **398**:252–257 (1983).
37. T. C. Bristow and K. Arackellian, "Surface Roughness Measurements using a Nomarski Type Scanning Instrument," *Metrology: Figure and Finish*, *Proc. SPIE* **749**:114–118 (1987).
38. J. M. Eastman and J. M. Zavislan, "A New Optical Surface Microprofiling Instrument," *Precision Surface Metrology*, *Proc. SPIE* **429**:56–64 (1983).
39. T. C. Bristow, G. Wagner, J. R. Bietry, and R. A. Auriemma, "Surface Profile Measurements on Curved Parts," *Surface Characterization and Testing II*, J. Grievenkamp and M. Young, eds., *Proc. SPIE* **1164**:134–141 (1989).
40. P. Glenn, "Angstrom Level Profilometry for Sub-Millimeter to Meter Scale Surface Errors," *Advanced Optical Manufacturing and Testing*, *Proc. SPIE* **1333**, G.M. Sanger, P. B. Ried, and L. R. Baker, eds., pp. 326–336 (1990).

41. P. Glenn, "Robust, Sub-Angstrom Level Mid-Spatial Frequency Profilometry," *Advanced Optical Manufacturing and Testing, Proc. SPIE* **1333**, G. M. Sanger, P. B. Ried, and L. R. Baker, eds., pp. 175–181 (1990).
42. I. Weingartner, M. Schulz, and C. Elster, "Novel Scanning Technique for Ultra-Precise Measurement of Topography," *Proc. SPIE* **3782**:306–317 (1999).
43. M. Schulz, P. Thomsen-Schmidt, and I. Weingaertner, "Reliable Curvature Sensor for Measuring the Topography of Complex Surfaces," *Optical Devices and Diagnostics in Materials Science, Proc. SPIE* **4098**:84–93 (2000).
44. P. Thomsen-Schmidt, M. Schulz, and I. Weingaertner, "Facility for the Curvature-Based Measurement of the Nanotopography of Complex Surfaces," *Proc. SPIE* **4098**:94–101 (2000).
45. I. Weingartner and C. Elster, "System of Four Distance Sensors for High-Accuracy Measurement of Topography," *Proc. Eng.* **28**(2):164–170 (2004).
46. I. Weingaertner, M. Wurm, R. Geckeler, et al., "Novel Scheme for the Ultra-Precise and Fast Measurement of the Nanotopography of Large Wafers," *Proc. SPIE* **4779**:13–22 (2002).
47. J. Illemann, R. Geckeler, I. Weingaertner, et al., "Topography Measurement of Nanometer Synchrotron-Optics," *Proc. SPIE* **4782**:29–37 (2002).
48. M. Schulz and C. Elster, "Traceable Multiple Sensor System for Measuring Curved Surface Profiles with High Accuracy and High Lateral Resolution," *Opt. Eng.* **45**(6) (2006).
49. M. Schulz, C. Elster, and I. Weingartner, "Coupled Distance Sensor Systems for High-Accuracy Topography Measurement: Accounting for Scanning Stage and Systematic Sensor Errors," *Proc. Eng.* **30**(1):32–38 (2006).
50. K. von Bieren, "Pencil Beam Interferometer for Aspherical Optical Surfaces," *Laser Diagnostics*, S. Holly, ed., *Proc. SPIE* **343**:101–108 (1982).
51. K. von Bieren, "Interferometry of Wavefronts Reflected Off Conical Surfaces," *Appl. Opt.* **22**:2109–2114 (1983).
52. P. Z. Takacs, S. Qian, and J. Colbert, "Design of a Long-Trace Surface Profiler," *Metrology—Figure and Finish, Proc. SPIE* **749**, B. Truax, ed., pp. 59–64 (1987).
53. P. Z. Takacs, S. K. Feng, E. L. Church, S. Qian, and W. Liu, "Long Trace Profile Measurements on Cylindrical Aspheres," *Advances in Fabrication and Metrology for Optics and Large Optics, Proc. SPIE* **966**, J. B. Arnold and R. A. Parks, eds., pp. 354–364 (1988).
54. P. Z. Takacs, K. Furenlid, R. DeBiasse, and E. L. Church, "Surface Topography Measurements over the 1 meter to 10 micrometer Spatial Period Bandwidth," *Surface Characterization and Testing II, Proc. SPIE* **1164**, J. E. Grievenkamp and M. Young, eds., pp. 203–211 (1989).
55. H. Lammert, F. Siewert, and T. Zeschke, *The Nano-Optic-Measuring Machine NOM at BESSY—Further Improvement of the Measuring Accuracy*, in *BESSY Annual Report 2005*. BESSY GmbH: Berlin, Germany. pp. 481–486 (2006).
56. S. Qian, W. Jark, and P. Z. Takacs, "The Penta-Prism LTP: A Long-Trace-Profiler with Stationary Optical Head and Moving Penta-Prism," *Rev. Sci Instrum.* **66**:2187 (1995).
57. S. Qian, W. Jark, P. Z. Takacs, K. J. Randall, and W. Yun, "In-Situ Surface Profiler for High Heat Load Mirror Measurement," *Opt. Eng.* **34**(2):396–402 (1995).
58. S. Qian, W. Jark, and P. Z. Takacs, "The Penta-Prism LTP: A Long-Trace-Profiler with Stationary Optical Head and Moving Penta-Prism," *Rev. Sci Instrum.* **66**:2562–2569 (1995).
59. S. Qian, W. Jark, P. Z. Takacs, K. J. Randall, Z. Xu, and W. Yun, "In-Situ Long Trace Profiler for Measurement of Mirror Profiles at Third Generation Synchrotron Facilities," *Rev. Sci Instrum.* **67**:3369 (1996).
60. S. Qian, H. Li, and P. Z. Takacs, "Penta-Prism Long Trace Profiler (PPLTP) for Measurement of Grazing Incidence Space Optics," *Multilayer and Grazing Incidence X-Ray/EUV Optics III, Proc. SPIE*, **2805**, R. Hoover and A. B. C. Walker, Jr., eds., pp. 108–114 (1996).
61. S. Qian, W. Jark, G. Sostero, A. Gambitta, et al., *Penta-prism LTP Detects First In-Situ Distortion Profile*, *Synchrotron Radiation News* **9**(3), pp. 42–44 (1996).
62. S. Qian, W. Jark, G. Sostero, A. Gambitta, et al., "Advantages of the In-Situ LTP Distortion Profile Test on High-Heat-Load Mirrors and Applications," *Proc. SPIE* **2856**, L. E. Berman and J. Arthur, eds., pp. 172–182 (1996).
63. H. Li, X. Li, M. W. Grindel, and P. Z. Takacs, "Measurement of X-Ray Telescope Mirrors Using A Vertical Scanning Long Trace Profiler," *Opt. Eng.* **35**(2):330–338 (1996).

64. H. Li, P. Z. Takacs, and T. Oversluizen, "Vertical Scanning Long Trace Profiler: a Tool for Metrology of X-Ray Mirrors," *Materials, Manufacturing, and Measurement for Synchrotron Radiation Mirrors, Proc. SPIE* **3152**:180–187 (1997).
65. M. Gubarev, T. Kester, and P. Z. Takacs, "Calibration of a Vertical-Scan Long Trace Profiler at MSFC," *Optical Manufacturing and Testing IV, Proc. SPIE* **4451**, H. P. Stahl, ed., pp. 333–339 (2001).
66. P. Z. Takacs, H. Li, X. Li, and M. W. Grindel, "3-D X-Ray Mirror Metrology with a Vertical Scanning Long Trace Profiler," *Rev. Sci. Instrum.* **67**, G. K. Shenoy and J. L. Dehmer, eds., pp. 3368–3369 (1996).
67. P. Z. Takacs, S. Qian, K. J. Randall, W. Yun, and H. Li, "Mirror Distortion Measurements with an In-Situ LTP," *Advances in Mirror Technology for Synchrotron X-Ray and Laser Applications, Proc. SPIE* **3447**:117–124 (1998).
68. ASME, *ASME B46.1-2002—Surface Texture (Surface Roughness, Waviness, and Lay)*, The American Society of Mechanical Engineers: New York (2002).
69. ISO, *ISO 4287:1997 Geometrical Product Specifications (GPS)—Surface Texture: Profile Method—Terms, Definitions and Surface Texture Parameters*, International Organization for Standardization, 1997.
70. SEMI, *SEMI MF1811-0704 Guide for Estimating the Power Spectral Density Function and Related Finish Parameters from Surface Profile Data*, Semiconductor Equipment and Materials International: San Jose, Calif., 2004.
71. J. C. Stover, *Optical Scattering: Measurement and Analysis*, 2d ed., Bellingham, WA: SPIE Press, p. 340, 1995.
72. E. L. Church and H. C. Berry, "Spectral Analysis of the Finish of Polished Optical Surfaces," *WEAR* **83**:189–201 (1982).
73. W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling, *Numerical Recipes—the Art of Scientific Computing*, Cambridge: Cambridge University Press. p. 818, 1986.
74. D. E. Newland, *An Introduction to Random Vibrations, Spectral and Wavelet Analysis*, 3d ed., New York: Wiley, p. 477, 1993.
75. C. K. Yuen and D. Fraser, *Digital Spectral Analysis*, London: Pitman Publishing Ltd., p. 156, 1979.
76. M. Born and E. Wolf, *Principles of Optics*. 6th ed., New York: Pergamon Press, p. 808, 1980.
77. E. L. Church, T. V. Vorburger, and J. C. Wyant, "Direct Comparison of Mechanical and Optical Measurements of the Finish of Precision Machined and Optical Surfaces," *Opt. Eng.* **24**(3):388–395 (1985).
78. V. V. Yashchuk, et al., "Cross-Check of Different Techniques for Two-Dimensional Power Spectral Density Measurements of X-Ray Optics," *Advances in Metrology for X-Ray and EUV Optics, Proc. SPIE* **5921**, L. Assoufid, P. Z. Takacs, and J. S. Taylor, eds., p. 59210G (2005).
79. V. V. Yashchuk, A. D. Franck, S. C. Irick, M. R. Howells, A. A. MacDowell, and W. R. McKinney, "Two-Dimensional Power Spectral Density Measurements of X-Ray Optics with the Micromap Interferometric Microscope," *Nano- and Micro-Metrology, Proc. SPIE* **5858**, H. Ottevaere, P. DeWolf, and D. S. P. O. Wiersma, eds., p. 58580A (2005).

This page intentionally left blank.

DO NOT DUPLICATE

ASTRONOMICAL X-RAY OPTICS

Marshall K. Joy and Brian D. Ramsey

*National Aeronautics and Space Administration
Marshall Space Flight Center
Huntsville, Alabama*

Over the past three decades, grazing incidence optics has transformed observational x-ray astronomy into a major scientific discipline at the cutting edge of research in astrophysics and cosmology. This chapter summarizes the design principles of grazing incidence optics for astronomical applications, describes the capabilities of the current generation of x-ray telescopes and the techniques used in their fabrication, and explores avenues of future development.

47.1 INTRODUCTION

The first detection of a cosmic x-ray source outside of our solar system was made during a brief rocket flight less than 50 years ago.¹ This flight, which discovered a bright source of x rays in the Scorpius Constellation, was quickly followed by other suborbital experiments and in 1970, the first dedicated x-ray astronomy satellite, UHURU, was launched into an equatorial orbit from Kenya.² UHURU, which used mechanically collimated gas-filled detectors, operated for just over 2 years and produced a catalog of 339 cosmic x-ray sources.

While UHURU significantly advanced the discipline, the real revolution in x-ray astronomy came about with the introduction of grazing-incidence optics aboard the Einstein observatory in 1978.³ Focusing optics provide an enormous increase in signal to noise ratio by concentrating source photons into a tiny region of the detector, thereby reducing the detector-area-dependent background to a very small value. Despite a modest collecting area, less than that of the UHURU, the Einstein observatory had two-to-three orders of magnitude more sensitivity, enabling emission from a wide range of sources to be detected and changing our view of the x-ray sky.

Since that time, payloads have increased in capability and sophistication. The current “flagship” x-ray astronomy missions are the U.S.-led Chandra observatory and the European-led XMM-Newton observatory. Chandra represents the state of the art in astronomical x-ray optics with sub-arcsecond on-axis angular resolution and about 0.1 m² of effective collecting area.⁴ Its sensitivity is over five orders of magnitude greater than that of UHURU, despite having only slightly greater collecting area; in deep fields Chandra resolves more than 1000 sources per square degree.

The XMM-Newton observatory, designed for high throughput, has nearly 0.5 m^2 of effective collecting area, and 15-arcsecond-level angular resolution.⁵

In addition to providing a considerable increase in sensitivity and enabling fine imaging, x-ray optics also permit the use of small-format, high-performance focal-plane detectors. Both Chandra and XMM feature fine-pixel silicon imagers with energy resolutions an order of magnitude greater than the earlier gas-filled detectors. Currently planned missions will utilize imaging x-ray calorimeters that will offer one to two orders of magnitude further spectroscopic improvement.

This chapter describes the optics used in, or with potential for use in, x-ray astronomy. The missions described above use mirror geometries based on a design first articulated by Wolter⁶ and these, and their fabrication techniques, are described in Sec. 47.2. An alternate mirror configuration, termed Kirkpatrick-Baez, which has not seen use yet in astronomy but offers potential future benefit, is described in Sec. 47.3. Payloads designed to extend the range of x-ray focusing optics into the hard-x-ray region are detailed in Sec. 47.4. Finally, new developments offering the promise of ultra-high angular resolution for future missions are discussed in Sec. 47.5.

47.2 WOLTER X-RAY OPTICS

Optical Design and Angular Resolution

Wolter optics are formed by grazing-incidence reflections off two concentric conic sections (a paraboloid and hyperboloid, or a paraboloid and ellipsoid—see Chap. 44). The most common case (Wolter type I) is conceptually similar to the familiar Cassegrain optical telescope: the incoming parallel beam of x rays first reflects from the parabolic section and then from the hyperbolic section, forming an image at the focus (Fig. 1). To increase the collecting area, reflecting shells of different diameters are nested, with a common focal plane.

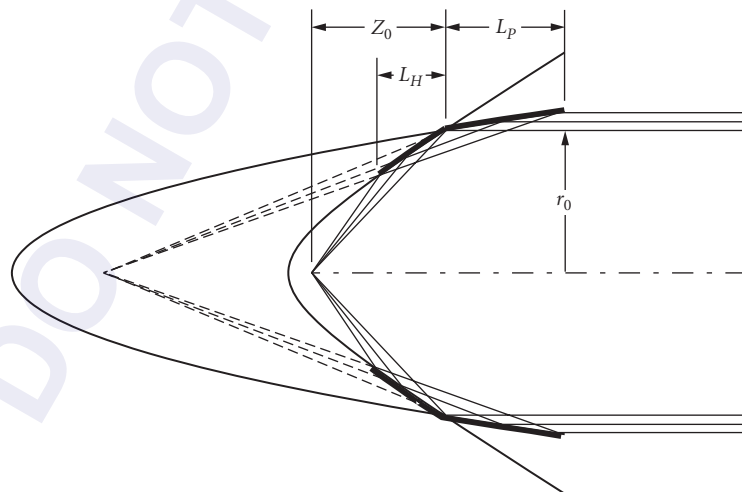


FIGURE 1 Geometry of a Wolter type I x-ray optic. Parallel light incident from the right is reflected at grazing incidence on the interior surfaces of the parabolic and hyperbolic sections; the image plane is at the focus of the hyperboloid.

A Wolter I optic can be described by four quantities:

1. The focal length, Z_0 , which is defined as the distance from the intersection of the paraboloid and hyperboloid to the focal point
2. The mean grazing angle α , which is defined in terms of the radius of the optic at the intersection plane r_0 :

$$\alpha \equiv \frac{1}{4} \arctan \frac{r_0}{Z_0} \quad (1)$$

3. The ratio of the grazing angles of the paraboloid and the hyperboloid ξ , measured at the intersection point
4. The length of the paraboloid L_p

Wolter optics produce a curved image plane, and have aberrations which can cause the angular resolution of the optic to be significantly worsened for x-ray sources that are displaced from the optical axis of the telescope (see Chap. 45); for these reasons, designs are usually optimized using detailed ray-trace simulations (see Chap. 35). However, a good approximation to the optimum design can be readily obtained using the results of Van Speybroeck and Chase.⁷ The highest x-ray energy that the optic must transmit largely determines the mean grazing angle (see “X-Ray Reflectivity” section), which in turn constrains the focal ratio of the optic [Eq. (1)]. The grazing angles on the parabolic and hyperbolic sections are usually comparable, so $\xi \approx 1$. With the diameter and length of the optic as free parameters, the curves in Van Speybroeck and Chase can be used to estimate the angular resolution and the collecting area for different designs. Very high resolution and good off-axis performance is possible with Wolter I optics. Figure 2 presents sub-arcsecond angular resolution images from the Chandra X-Ray Observatory and wide field images from the XMM-Newton Observatory.

Mirror Figure and Surface Roughness

Irregularities in the mirror surface will cause light to be scattered out of the core of the x-ray image, degrading the angular resolution of the telescope. If an incoming x ray strikes an area of the mirror surface that is displaced from the ideal height by an amount σ , the resulting optical path difference is given by

$$\text{OPD} = 2\sigma \sin \alpha \quad (2)$$

and the corresponding phase difference is

$$\Delta = \frac{4\pi\sigma \sin \alpha}{\lambda} \quad (3)$$

where α is the grazing angle and λ is the x-ray wavelength. For a uniformly rough mirror surface with a gaussian height distribution, RMS values of σ and Δ can be used to calculate the scattered intensity relative to the total intensity:^{8,9}

$$\frac{I_s}{I_0} = 1 - e^{-\Delta_{\text{RMS}}^2} \quad (4)$$

This result implies that high-quality x-ray reflectors must have exceptionally smooth surfaces: in order for the scattered intensity I_s to be small, Δ_{RMS} must be $\ll 1$. For a grazing angle of 0.5° and a wavelength of 1 \AA , a surface with an RMS microroughness of 4 \AA will scatter approximately 20 percent of the incident intensity. A surface roughness of 8 \AA will scatter more than 50 percent of the incident intensity at this x-ray wavelength and grazing angle.

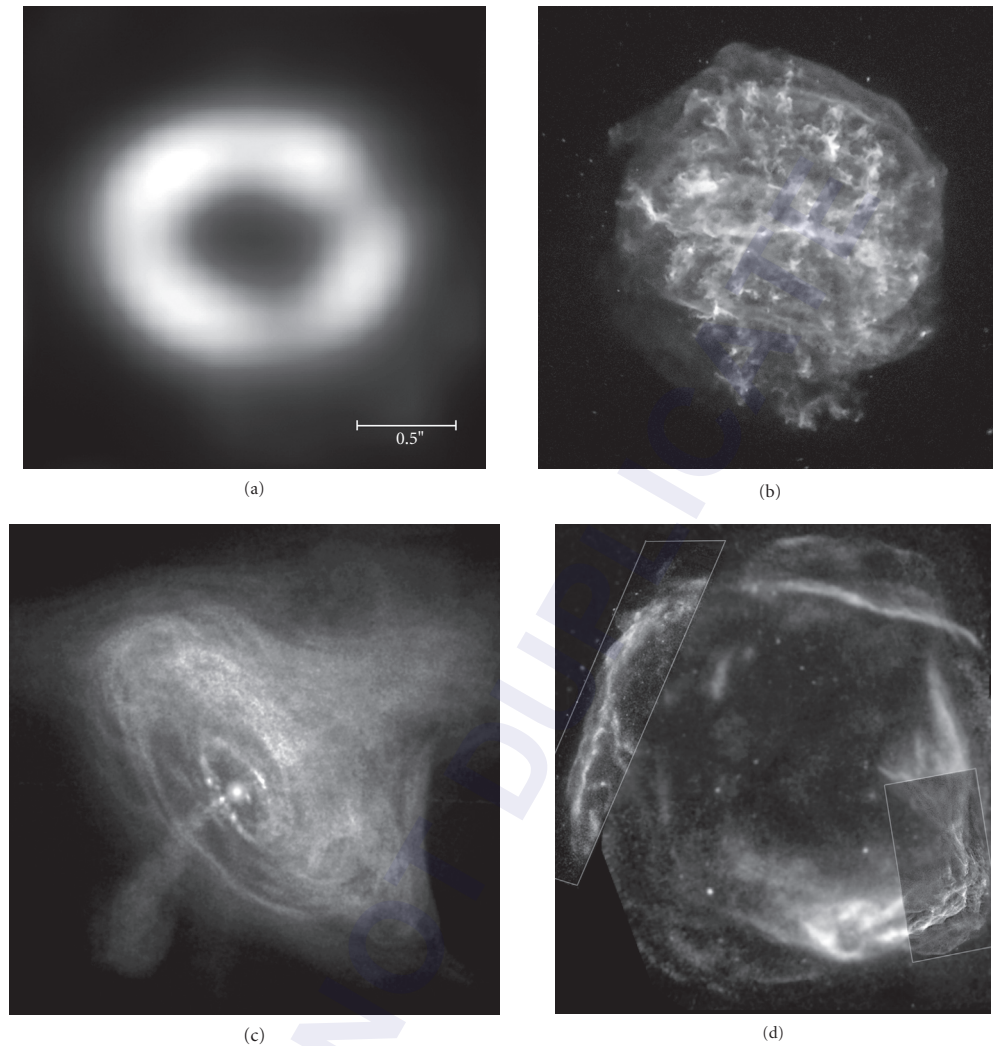


FIGURE 2 (a) A Chandra x-ray image of the fiery ring surrounding the 1987 supernova explosion in the Dorado constellation. Subarcsecond angular resolution is required to resolve the structure surrounding the supernova remnant (<http://chandra.harvard.edu/photo/2005/sn87a/>). (NASA/CXC/PSU/S.Park & D. Burrows.) (b) A Chandra x-ray image of the supernova remnant G292.0+1.8. The colors in the image encode the x-ray energies emitted by the supernova remnant; the center of G292.0+1.8 contains a region of high energy x-ray emission from the magnetized bubble of high-energy particles that surround the pulsar, a rapidly rotating neutron star that remained behind after the original, massive star exploded (<http://chandra.harvard.edu/photo/2007/g292/>). (NASA/CXC/Penn State/S.Park et al.) (c) Chandra x-ray image of the Crab Nebula—the remains of a nearby supernova explosion which was seen on Earth in 1054 AD. At the center of the bright nebula is a rapidly spinning neutron star, or pulsar, that emits pulses of radiation 30 times a second (<http://chandra.harvard.edu/photo/2002/0052>). (NASA/CXC/ASU/J. Hester et al.) (d) XMM-Newton and Chandra x-ray images RCW 86, an expanding ring of debris that was created after a massive star in the Milky Way collapsed onto itself and exploded. Both the XMM-Newton and Chandra images show low-energy x-rays in red, medium energies in green and high energies in blue. The Chandra observations focused on the northeast (left-hand) and southwest (lower right) side of RCW 86, and show that x-ray radiation is produced both by high-energy electrons accelerated in a magnetic field (blue) as well as heat from the blast itself (red). These images demonstrate the large field of view and moderate angular resolution of XMM-Newton, compared to the smaller field of view and high angular resolution provided by Chandra (<http://chandra.harvard.edu/photo/2007/2snr/>). (NASA/CXC/ESA/Univ. of Utrecht/J. Vink et al.) (See also color insert.)

To calculate the effect of surface imperfections on the x-ray point spread function in detail, it is necessary to measure the distribution of the deviations on all spatial scales of the optic, from the microroughness on submicron scales to the overall slope error of the full mirror. The surface deviations are characterized by the power spectral density, and the resulting x-ray scattering can be calculated using the methods described by Church,¹⁰ Aschenbach,⁸ O'Dell et al.,¹¹ and Hughes et al.¹²

X-Ray Reflectivity

For most astronomical applications, x-ray telescopes are required to have high reflection efficiency across the operational energy band, and also to function at the highest possible energy for a given focal length and diameter. The overall reflection efficiency and the high energy response can both be increased by applying a high density coating to the polished reflecting surfaces of the optics, such as nickel, gold, platinum, or iridium (the latter element having the highest density and the best reflectivity).¹³ As discussed earlier, the uniformity and smoothness of the reflective coating is critical to the imaging performance of the optic.

The reflection efficiency of an x-ray optic is strongly dependent on energy, due to the presence of atomic absorption edges and the rapid decrease in efficiency at high energies. The reflectivity can be readily calculated using the Fresnel reflection formulae expressed in terms of the complex dielectric constant of the reflecting surface, κ .¹³ Dielectric constants can be easily derived from experimentally determined atomic scattering factors, f , using the relation

$$\kappa = 1 - \frac{r_e \rho N_A}{\pi A_w} \lambda^2 f \quad (5)$$

where r_e is the classical electron radius, N_A is Avogadro's number, and ρ and A_w are the density and molecular weight of the reflecting material. An extensive compilation of atomic scattering factors for elements 1 to 92 over the energy range 50 eV to 30 keV is given by Henke, Gullikson, and Davis¹⁴ (also see Chap. 36).

Effective Collecting Area

The useful collecting area of a grazing incidence optic depends on several factors: the size and number of mirror shells, the projection of the grazing incidence mirrors onto the sky, blockage of the aperture by support structures, x-ray reflection efficiency, and vignetting of off-axis sources. The effective collecting area A_e is the convolution of all of these terms; it is highly energy dependent due to the mirror reflectivity, and falls to zero at the high energy cutoff of the optic.

Technologies for Fabricating Wolter-Type X-Ray Optics

There are currently three primary methods for fabricating nested Wolter x-ray optics for astronomical applications. The approach that has produced the highest angular resolution involves the fabrication of Zerodur glass shells, a few centimeters thick, which are ground, figured, highly polished, and coated on the interior reflecting surface. Optics of this type have been built and flown on the Einstein observatory,³ ROSAT,¹⁵ and the Chandra X-Ray Observatory.¹⁶ Einstein and ROSAT each had approximately 5 arcsec angular resolution, while Chandra has sub-arcsecond on-axis resolution. These optics represent what might be called the traditional approach to mirror fabrication which involves meticulous figuring and polishing of relatively heavy substrates. The resulting mirrors deliver superb performance, but in an effort to avoid the high cost necessarily incurred by this labor-intensive work and the significant weight of the mirrors and their support structures, other fabrication techniques have been developed that trade angular resolution for light weight, higher throughput, and lower cost.

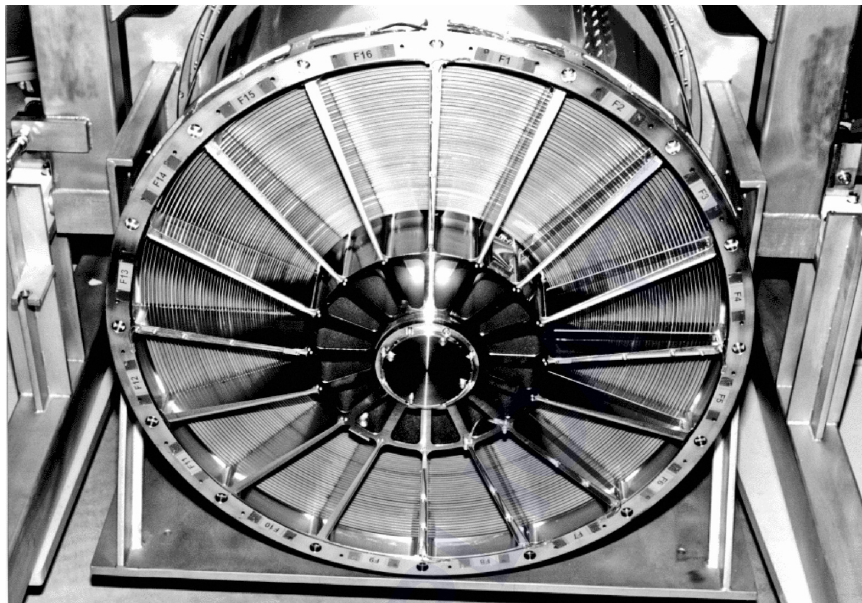


FIGURE 3 An x-ray optics module for the XMM observatory. Fifty-eight electroformed nickel Wolter I optics are nested to increase the effective x-ray collecting area. (See also color insert.) (Photo courtesy of ESA.)

A mirror-fabrication process used widely in x-ray astronomy is electroformed nickel replication (ENR). This technique makes use of a master or “mandrel” which is used to replicate identical mirror shells. The required parabolic and hyperbolic surfaces are machined into the outer surface of the mandrel, which is then highly polished and coated or chemically treated. A thin nickel shell is electroplated onto the mandrel; the shell and mandrel (which have different coefficients of thermal expansion) are then separated by cooling. The attraction of the process is that the resulting shells are full cylinders that contain both parabolic and hyperbolic segments of the Wolter geometry. This makes them inherently stable and results in reasonable angular resolution despite their thin (1 mm or less) walls. A good example of the use of the ENR process is in the mirrors for the XMM-Newton Observatory which has three optics modules, each containing 58 replicated Wolter-I shells (Fig. 3), which provide high effective collecting area and 15-arcsecond angular resolution.¹⁷

A third approach to mirror fabrication is to utilize reflector segments mounted in a housing to give the desired mirror geometry. This technique opens the possibility of compact, very-light-weight, densely-packed optics modules. To date, many segmented optics have flown on satellite missions and all have utilized thin aluminum reflectors coated (typically via an epoxy replication off a smooth mandrel) to improve their surface roughness. The current SUZAKU¹⁸ x-ray astronomy mission has several foil segment mirror modules (Fig. 4). Typically these have around 175 nested shells, each consisting of $8 \times 150 \mu\text{m}$ thick reflectors. The on-orbit measured angular resolution is around 2 arcminutes.

Recently there has been a significant optics development effort utilizing slumped glass segmented reflectors rather than aluminum. The attraction of the glass is that it is available with very smooth surfaces due to the fabrication process used. Certain glass types, notably float borosilicate glass, is obtainable with sub-5-Å smoothness making it an ideal x-ray reflector even for demanding multilayer coatings. The NuSTAR mission,¹⁹ recently selected as a NASA small explorer, utilizes float glass reflectors “slumped” by heating the glass to approximately 600°C on a suitably figured pyrex mandrel. The reflectors are coated with multilayers for hard-x-ray responsivity (see Sec. 47.4), and

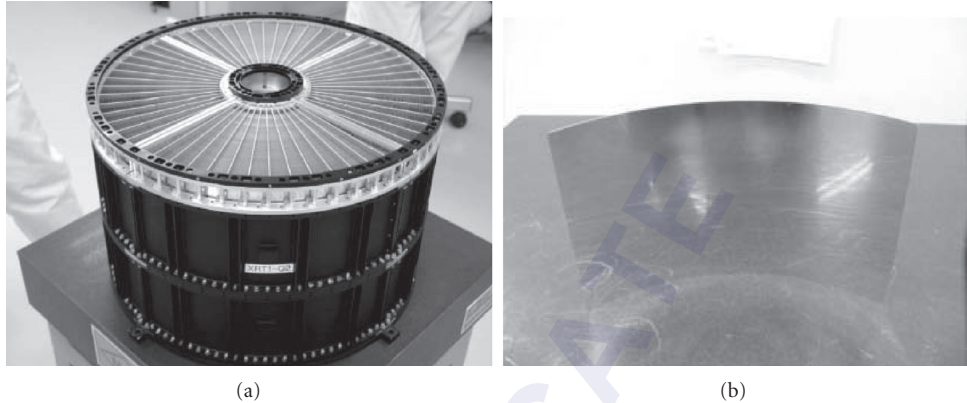


FIGURE 4 The segmented foil mirrors aboard the SUZAKU spacecraft. (a) A complete mirror module and (b) a single aluminum foil reflector coated with gold. (See also color insert.)

are assembled into a module through the use of machined precision graphite spacers inserted and glued between successive shells. The spacers are machined in situ, before each successive shell is glued, to prevent error buildup. The expected angular resolution for the NuSTAR optics is approximately 50 arcseconds.

47.3 KIRKPATRICK-BAEZ OPTICS

The first demonstration of grazing incidence x-ray imaging was performed by Kirkpatrick and Baez²⁰ using mirrors with one-dimensional parabolic curvature (Fig. 5a). An optic of this type focuses in one plane only, so two mirrors are mounted orthogonal to each other to achieve 2D imaging (Fig. 5b); they can also be stacked to increase the effective collecting area (Fig. 5c). Design studies of multielement Kirkpatrick-Baez systems have been carried out by Van Speybroeck et al.,²¹ Weisskopf,²² and Kast.²³

Kirkpatrick-Baez optics has not been widely used in astronomical applications because the required grazing angles are a factor of two larger than those of a Wolter design, for the same optic diameter and focal length. (This can be seen from Fig. 7: An x-ray reflecting in the horizontal plane receives the full angular displacement in a single reflection, while a Wolter optic achieves the same angular displacement via two half-angle reflections.) The larger grazing angles in the Kirkpatrick-Baez design significantly reduce the high energy efficiency of the optic. However, optimal packing schemes and ease of manufacture may eventually permit Kirkpatrick-Baez designs to overcome this disadvantage; for example, mass-produced flat plates can be substituted for the curved optics in cases where the curvature of the mirrors is small (Fig. 6). In the small angle approximation the projected width of the flat plate is

$$\delta = \frac{Lr_0}{2F} \quad (6)$$

where L is the length of the plate, r_0 is the distance of the plate from the optical axis, and F is the focal length. The limiting angular resolution due to the flat plate geometry, $\Delta\theta$, is

$$\Delta\theta(\text{arcsec}) \sim 10^5 \frac{Lr_0}{F^2} \quad (7)$$

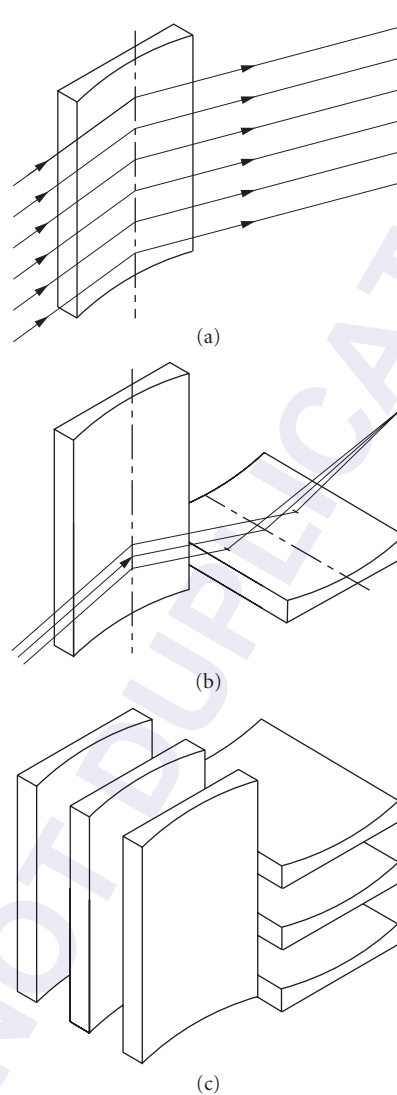


FIGURE 5 Kirkpatrick-Baez optics. (a) Mirrors with one-dimensional parabolic curvature focus in one plane only, producing a line image; (b) two mirrors mounted orthogonal to each other can produce 2D imaging; and (c) parabolic Kirkpatrick-Baez optics can be nested to increase the effective collecting area.

Equation (7) indicates that the error introduced by the flat plate approximation can be made reasonably small: $\Delta\theta = 10$ arcsec for a mirror 0.1 m in length and $r_0 = 0.1$ m away from the axis of a 10-m focal length optic. Analytic expressions for the number of flat plates, their spacing, and the resulting packing fraction are given by Weisskopf.²²

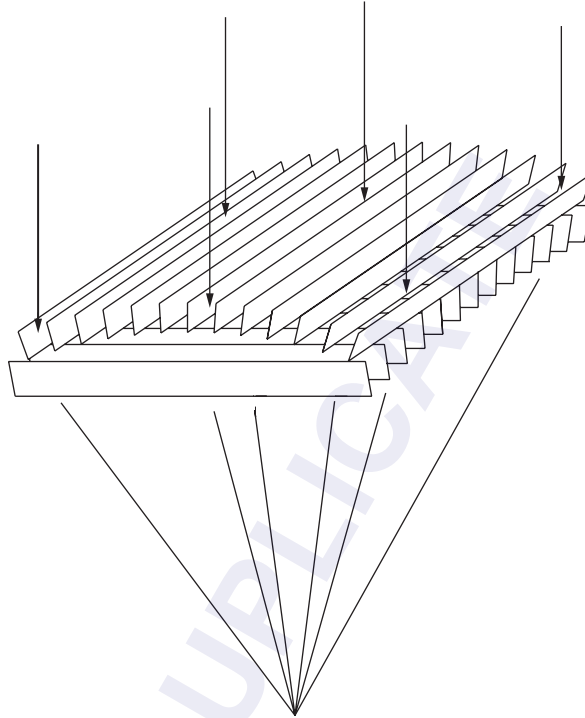


FIGURE 6 Flat plates can be substituted for the parabolic Kirkpatrick-Baez optics in cases where the curvature of the mirrors is small, and 2D imaging can be achieved by stacking two orthogonal sets of flat reflectors.

47.4 HARD X-RAY OPTICS

As noted in Sec. 47.1, x-ray optics has revolutionized the field of x-ray astronomy, with current state-of-the-art telescopes having more than five orders of magnitude greater sensitivity than the first satellite-borne instruments. This statement is only true, however, below approximately 10 keV where x-ray optics have been used extensively. The hard-x-ray region, above 10 keV, is very important for study as in this energy regime sources transit from thermal to nonthermal emission mechanisms. In addition, nuclear lines appear and obscured objects become visible. Despite this significance, this energy region remains relatively unexplored at high sensitivities and fine angular scales.

The difficulty inherent in obtaining high-energy response is that the critical angle, below which x rays can be reflected, is given by

$$\phi_c \sim 0.93 \cdot \lambda \cdot \sqrt{\rho}$$

where λ is the wavelength of the incident x ray (nanometers) and ρ is the density of the reflecting medium (grams per cubic centimeter). Thus at 6 keV, for example, a nickel surface will reflect x rays up to an angle of approximately 1/2 degree, but at 60 keV, this will be only 3 arcminutes and thus the projected collecting area of hard-x-ray mirror shells becomes very small.

There have recently been several balloon payloads developed with hard-x-ray telescope systems using slightly different approaches to obtaining high-energy response. The HERO balloon payload^{24,25} uses an array of shallow-graze-angle, iridium-coated optics fabricated using the electroformed nickel replication process, while the HEFT²⁶ and the InFOCUS/SUMIT²⁷ payloads use multilayers (see Chap. 41) to extend the graze angles over which useful reflectivity can be obtained. Both InFOCUS/SUMIT and HEFT utilize a segmented optics approach, the former with aluminum reflectors and the latter with slumped glass segments as described in “Technologies for Fabricating Wolter-Type X-Ray Optics” section.

There are also various satellite missions with hard-x-ray focusing optics that have been approved and are in the early stages of development. Among these are the US NuSTAR small explorer mission,^{19,28} the Japanese NeXT (ASTRO-h) mission.^{27,29}

NuSTAR is a small explorer mission that has recently been approved for flight with a planned launch in 2011. It is an outgrowth of the HEFT balloon program, described previously, but with an increased focal length and glass-mirror-module size to give an effective area around 200 cm² at 50 keV and useful response to above 70 keV. This, coupled with the long integration times available on orbit, results in a two-order-of-magnitude increase in sensitivity over currently orbiting (non-focusing) hard-x-ray instruments. Launched in an equatorial orbit for low background, NuSTAR will have a 2-year mission lifetime.

The Japanese NeXT mission builds on the InFOCUS/SUMIT experience, and will utilize multilayer-coated aluminum foil optics for extended response up to approximately 60 keV. It is projected to fly in the 2013/2014 timeframe.

47.5 TOWARD HIGHER ANGULAR RESOLUTION

The next generation of x-ray missions requires ever-higher angular resolutions and collecting areas and it is difficult to construct grazing incidence x-ray optics that combine both of these. The Chandra X-ray Observatory currently delivers sub-arcsecond angular resolution, a superlative achievement for an x-ray telescope, but still an order of magnitude below the resolution achieved by optical observatories, and more than three orders of magnitude below the diffraction limit. Improving conventional Wolter x-ray optics to produce milli-arcsecond (or better) angular resolution is a formidable technical challenge. One alternative is to exploit an approach that is widely used in the radio and optical spectral regions: interferometry.

Several fundamental questions about the feasibility of x-ray interferometry are immediately apparent. First, interferometry requires that the phase errors be controlled to a fraction of a wavelength; is it possible to figure and control optics for x-ray wavelengths on the scale of Angstroms? Second, cosmic x-ray sources are notoriously faint (especially the cosmologically interesting sources at high redshift), and the effective collecting areas of x-ray telescopes are low; is it possible to build an x-ray interferometer that would collect a sufficient number of photons to construct a useful image in a finite amount of time?

The scientific case for very-high-angular resolution x-ray imaging proves to be strong. Many of the most interesting objects in the universe (including active galactic nuclei, accretion disks, and the black holes that power them) are very small (< 1 milli-arcsec), very hot ($T > 10^8$ K), and radiate a large amount of thermal energy ($B_\nu \propto T^4$) primarily in the x-ray spectral region. Obtaining the highest scientific returns, however, will require microarcsecond x-ray angular resolution.³⁰

The answers to some of the technical feasibility questions are interesting and rather unexpected. Figure 7 shows a schematic drawing of an x-ray interferometer, in which the beams from two x-ray telescopes are mixed to form fringes. Grazing incidence mirrors are used for both the telescope optics and the beam-combining optics.

How accurately must the grazing incidence optics be figured and controlled in order to achieve an acceptable phase error at x-ray wavelengths? For a $1/2$ degree grazing incidence optic, $\lambda = 2$ Å, and a $\lambda/10$ path difference, Eq. (2) indicates that the tolerance on surface deviations is large

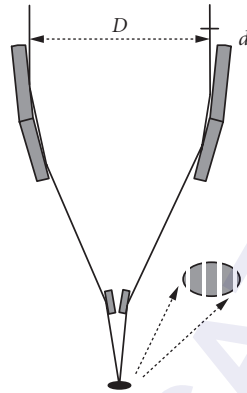


FIGURE 7 Schematic drawing of a grazing-incidence x-ray interferometer.

compared to the wavelength: $\sigma \approx 10 \text{ \AA}$. Thus, for typical grazing angles, the factor of $\sin\alpha$ in Eq. (2) relaxes the tolerances on the mirror surface by approximately two orders of magnitude. Significantly, this brings the required x-ray optics and metrology solidly into a regime that has been demonstrated in the laboratory, suggesting the future possibility of an ultra-high-angular resolution x-ray astronomy mission.

47.6 REFERENCES

1. R. Giacconi, H. Gursky, F. Paolini, and B. Rossi, "Evidence for X-Rays from Sources Outside the Solar System," *Phys. Rev. Lett.* **9**:435 (1962).
2. R. Giacconi, E. Kellogg, P. Gorenstein, H. Gursky, and H. Tananbaum, "An X-Ray Scan of the Galactic Plane from UHURU," *Ap. J.* **165**: L27 (1971).
3. R. Giacconi, G. Branduardi, U. Briel, A. Epstein, D. Fabricant, E. Feigelson, W. Forman, et al., "The Einstein/HEAO 2/X-Ray Observatory," *Ap. J.* **230**:540 (1979).
4. M. C. Weisskopf and S. L. O'Dell, "Calibration of the AXAF Observatory: Overview," *Proc. SPIE* **3113**:2 (1997).
5. P. Gondoin, B. R. Aschenbach, M. W. Beijersbergen, R. Egger, F. A. Jansen, Y. Stockman, and Tock, J. -P., "Calibration of the First XMM Flight Mirror Module: I. Image Quality," *Proc. SPIE* **3444**:278 (1998).
6. H. Wolter, "Grazing Incidence Mirror Systems as Imaging Optics for X-Rays," *Ann. Phys.* **10** (Ser. 6):94 (1952).
7. L. P. Van Speybroeck, and R. C. Chase, "Design Parameters of Paraboloid-Hyperboloid Telescopes for X-Ray Astronomy," *Appl. Opt.* **11**:440 (1972).
8. B. Aschenbach, "X-Ray Telescopes," *Rep. Prog. Phys.* **48**:579–629 (1985).
9. P. Beckmann, and A. Spizzichino, *The Scattering of Electromagnetic Waves from Rough Surfaces*, Oxford: Pergamon (1963).
10. E. L. Church, "Role of Surface Topography in X-Ray Scattering," *Proc. SPIE* **184**: 196 (1979).
11. S. L. O'Dell, R. F. Elsner, J. J. Kolodziejczak, M. Weisskopf, J. P. Hughes, and L. P. van Speybroeck, "X-Ray Evidence for Particulate Contamination on the AXAF VETA-I Mirrors," *Proc. SPIE* **1742**:171 (1993).
12. J. P. Hughes, D. Schwartz, A. Szentgyorgyi, L. Van Speybroeck, and P. Zhao, "Surface Finish Quality of the Outer AXAF Mirror Pair Based on X-Ray Measurements of the VETA-I," *Proc. SPIE* **1742**:152 (1993).
13. R. F. Elsner, S. L. O'Dell, and M. C. Weisskopf, "Effective Area of the AXAF X-Ray Telescope: Dependence upon Dielectric Constants of Coating Materials," *J. X-Ray Sci. Technol.* **3**, 35–44 (1991).

14. B. L. Henke, E. M. Gullikson, and J. C. Davis, "X-Ray Interactions: Photoabsorption, Scattering, Transmission, and Reflection at $E = 50\text{--}30000$ eV, $Z = 1\text{--}92$," *Atomic Data and Nuclear Data Tables* **54**(2): 181–342 (1993). This data, along with useful computational tools, is also available on the internet at http://www-cxro.lbl.gov/optical_constants.
15. B. Aschenbach, "First Results from the X-Ray Astronomy Mission ROSAT," *Rev. Mod. Astron.* **4**:173 (1991).
16. M.C. Weisskopf, "Five Years of Operation of the Chandra X-Ray Observatory," *Proc. SPIE* **5488**:25 (2004).
17. Y. Stockman, P. Barzin, H. Hansen, E. Mazy, J. P. Tock, D. de Chambure, R. Laine, D. Kampf, et al., "XMM Flight Model Mirror Module Testing," *Proc. SPIE* **3766**:51 (1999).
18. K. Mitsuda, M. Bautz, H. Inoue, R. L. Kelley, K. Koyama, H. Kunieda, K. Makishima, et al., "The X-Ray Observatory Suzaku," *Publications of the ASJ* **59**:1–7 (2007).
19. F. A. Harrison, F. E. Christensen, W. Craig, C. Hailey, W. Baumgartner, C. M. H. Chen, J. Chonko, et al., "Development of the HEFT and NuSTAR Focusing Telescopes," *Exp. Astr.* **20**:131 (2005).
20. P. Kirkpatrick, and A. V. Baez, "Formation of Optical Images by X-Rays," *J. Opt. Soc. Am.* **38**:776 (1948).
21. L. P. Van Speybroeck, R. C. Chase, and T. F. Zehnpfennig, "Orthogonal Mirror Telescopes for X-Ray Astronomy," *Appl. Opt.* **10**:945 (1971).
22. M. C. Weisskopf, "Design of Grazing-Incidence X-Ray Telescopes," *Appl. Opt.* **12**:1436 (1973).
23. J. W. Kast, "Scanning Kirkpatrick-Baez X-Ray Telescope to Maximize Effective Area and Eliminate Spurious Images," *Appl. Opt.* **14**:537 (1975).
24. B. D. Ramsey, C. D. Alexander, J. A. Apple, C. M. Benson, K. L. Dietz, R. F. Elsner, D. Engelhaupt, et al., "First Images from HERO: A Hard-X-Ray Focusing Telescope," *J. Astrophys.* **568**:432–435 (2002).
25. B. D. Ramsey, R. Elsner, D. Engelhaupt, M. Gubarev, J. J. Dell Kolodziejczak, C. O. S. L. O Speegle, and M. C. Weisskopf, "The Development of Hard-X-Ray Optics at NASA/MSFC," *Proc. SPIE* **5168**:129–135 (2004).
26. J. E. Koglin, F. E. Christensen, J. Chonko, W. W. Craig, T. R. Decker, M. A. Jimenez-Garate, K. S. Gunderson, et al., "Development and Production of Hard X-Ray Multilayer Optics for HEFT," *Proc. SPIE* **4851**: 607–618 (2003).
27. Y. Ogasaka, K. Tamura, T. Miyazawa, Y. Fukaya, T. Iwahara, N. Sasaki, A. Furuzawa, et al., "Thin-Foil Multilayer-Supermirror Hard X-Ray Telescope for InFOC μ S/SUMIT Balloon Experiments and NeXT Satellite Program," *Proc. SPIE* **6688**:03 (2007).
28. J. E. Koglin, F. E. Christensen, W. Craig, T. R. Decker, C. J. Hailey, F. A. Harrison, C. Hawthorn, et al., "NuSTAR Hard X-Ray Optics," *Proc. SPIE* **5900**(31) (2005).
29. Y. Ogasaka, K. Tamura, R. Shibata, A. Furuzawa, T. Okajima, K. Yamashita, Y. Tawara, H. Kuneida, "NeXT Hard X-Ray Telescope," *Proc. SPIE* **5488**:148–155 (2004).
30. K. Gendreau, N. White, S. Owens, W. Cash, A. Shipley, and M. Joy, "The MAXIM X-Ray Interferometry Mission Concept Study," *Liege International Astrophysical Colloquia* **36**:11 (2001).

Ladislav Pina

*Faculty of Nuclear Sciences and Physical Engineering
Czech Technical University
Prague, V Holesovickach 2*

48.1 INTRODUCTION

Improved source development and the growing number of types of x-ray sources have generated an increased interest in detection, spectroscopy, and x-ray imaging. Synchrotron, free electron laser (FEL), laser plasma, tokamak plasma, capillary discharge, and Z-pinch x-ray sources need various diagnostics and offer many applications in science and technology. X-ray imaging and x-ray optical systems play a major role. The study of refraction, reflection, and diffraction in the x-ray region have led to the development of new optical elements and systems. They allow extended use of x rays in microscopy, computed tomography, nondestructive testing (NDT), and material science, as well as in the semiconductor industry, arts, security, and astronomy. Demands for extremely wide field of view (FOV) and cost-effective optical systems have led to the concept of multifoil optics (MFO). The basic principles and designs of grazing incidence multifoil optics are reviewed in this chapter. Recent developments and typical applications for multifoil optics in astronomy, extreme ultraviolet (EUV) lithography, and microscopy are described. The requirements for these applications are discussed.

48.2 GRAZING INCIDENCE OPTICS

Atomic scattering functions for x rays had been studied by Henke¹ and are tabulated in Chap. 36. The real part of complex index of refraction in the x-ray region approaches unity so the focusing effect is small in comparison with classic visible radiation optics. The imaginary part of the index causes absorption that is not negligible for most materials. Therefore, refraction optics had played no significant role until recently (see Chap. 37). The theory of reflection of x rays from a mirror surface shows the possibilities and limitations for optics in the x-ray spectral region (see Chaps. 26 and 44). Reflection optics have been more suitable for most applications. However, x rays can reflect only at grazing incidence angles and, thus, only a relatively small solid angle of collection can be used. Different solutions can be applied to enlarge the collection angle such as nested Wolter-type optics described in Chaps. 44 to 46, or multifoil optics (MFO). The microroughness of the reflecting surface plays a critical role in the x-ray spectral region. Microroughness on no larger than an atomic scale is necessary for acceptable x-ray reflection² and optic figure errors no larger than micro-

nanometer scales are necessary for high-resolution imaging. Grazing incidence x-ray optics can be used to collect the radiation in a wide wavelength range. While grazing incidence optics are commonly used in space x-ray telescopes, they can be also successfully used for laboratory imaging, as well as for collecting x rays from the laboratory sources. The critical angle is relatively large in case of the EUV radiation. It can be up to 15° for gold-coated mirrors with surface microroughness below 1 nm and radiation wavelength around 10 nm, or better for some other materials. Beside flat mirrors, the most commonly used optics are toroidal mirrors, ellipsoidal mirrors, or parabolic mirrors. Kirkpatrick-Baez,³ Wolter,⁴ Schwarzschild, and polycapillary systems are also commonly used and based on reflection. Multilayer structures can be deposited on the surfaces of mirrors to achieve the desired spectral properties or higher solid angles of collection.

48.3 MULTIFOIL LOBSTER-EYE OPTICS

The first representative of the MFO is lobster-eye optics (LE). The name of these optics arises because lobsters have eyes which operate in a similar manner.⁵ See also Chap. 49.

Lobster-Eye Geometry

X-ray optics offer an excellent opportunity to achieve very wide fields of view. A one-dimensional lobster-eye geometry was originally suggested by Schmidt,⁶ based upon flat reflectors. The device consists of a set of flat reflecting surfaces. The plane reflectors are arranged in a uniform radial pattern around the perimeter of a cylinder of radius R (Fig. 1a). X rays from a given direction are focused to a line on the surface of a cylinder of radius $R/2$. The azimuthal angle is determined directly from the centroid of the focused image. Used at a glancing angle for x rays of wavelength 1 nm and longer, this device can be used to focus a sizable portion of an intercepted beam of parallel x rays. Focusing is not perfect and the image size is finite. On the other hand, this type of focusing device

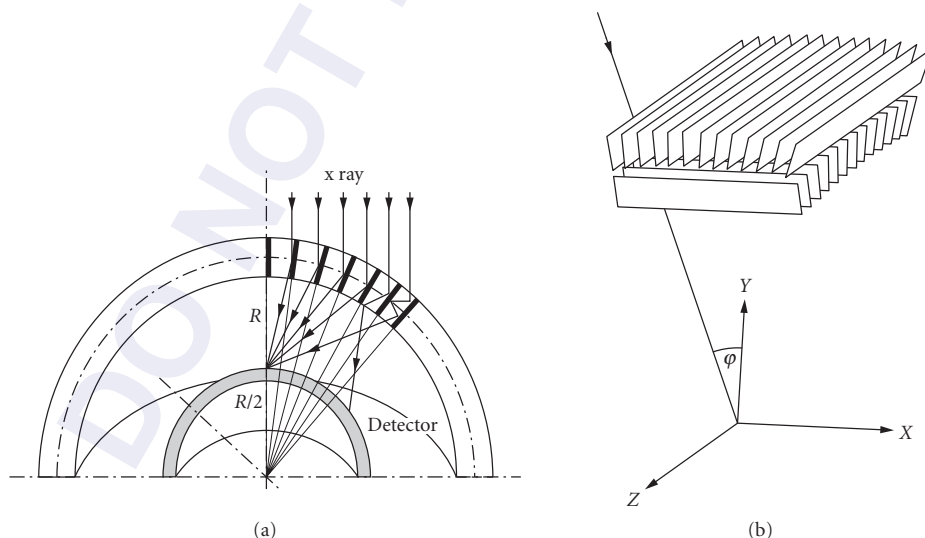


FIGURE 1 (a) Geometry of a lobster-eye x-ray optic. Parallel light incident from the top is reflected at grazing incidence on the set of flat mirrors. The image is formed on a sphere with $R/2$ radius. (b) The Schmidt lobster-eye objective in the double-focusing arrangement. (Courtesy of A. Inneman.)

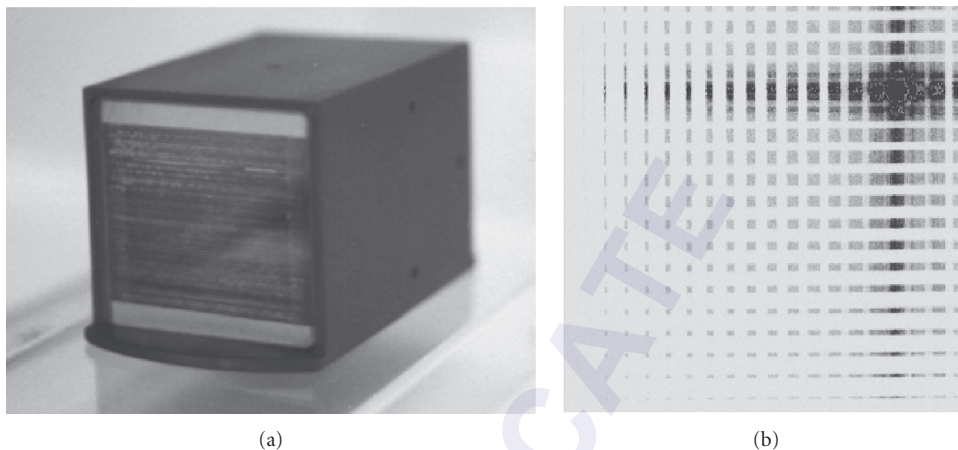


FIGURE 2 (a) The mini (24×24 mm, 0.1-mm-thick foils spaced by 0.3 mm) Schmidt LE module. (b) The mini Schmidt LE module focal spot image (8-keV x rays from microfocus x-ray tube, image area 12.3×12.3 mm). Maximum beam intensity in the focal spot is $680\times$ higher than the intensity of unfocused beam. (Courtesy of Reflex.)

offers a wide field of view, up to a maximum of 2π with coded apertures. It appears practically possible to achieve an angular resolution of the order of one-tenth of a degree or better. Two such systems in sequence, with orthogonal stacks of reflectors, form a double-focusing device (Figs. 1b and 2). Such device offers a field of view of up to 1000 square degrees at a moderate angular resolution.

48.4 MULTIFOIL KIRKPATRICK-BAEZ OPTICS

Schmidt LE systems with flat mirrors offer unsurpassed FOV, but have serious imaging and focusing limitations. Ideally, a good LE system should have mirrors as short as possible, very thin and densely packed. In laboratory applications, where the FOV is often not of the key importance, better focusing systems can be used. Kirkpatrick and Baez³ proposed combination of two mirrors in orthogonal

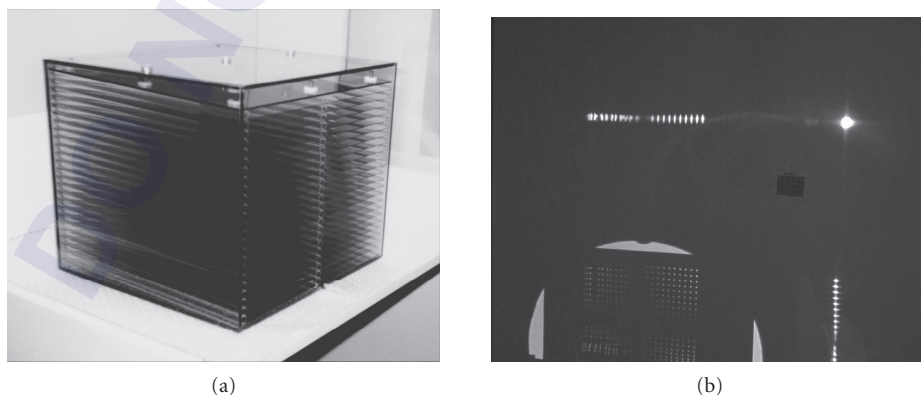


FIGURE 3 Kirkpatrick-Baez test Au-coated glass foils system (a) and VIS focal image (see also color insert) (b) studied for XEUS project. (Courtesy of Reflex.)

configuration (KB system) in order to achieve two-dimensional imaging. With this mirror combination they were the first to demonstrate experimentally grazing incidence x-ray imaging. Multifoil versions of KB systems (Fig. 3—example of MFO KB test module) have been studied for astronomical purposes by Van Speybroeck et al.,⁷ Weisskopf,⁸ Kast,⁹ and Gorenstein.¹⁰ A laboratory MFO system with elliptical mirrors was recently designed and tested in EUV lithography.

48.5 SUMMARY

Wide-field lobster-eye x-ray telescopes are expected to play an important role in future x-ray astrophysics missions and projects. These devices allow the study of novel science, including such important fields such as gamma ray bursts (GRB). Use of wide field x-ray optics will allow the signal-to-noise ratio to be increased compared to measurements with nonfocusing devices. The expected limiting sensitivity of LE telescopes is roughly 10^{-12} ergcm⁻²s⁻¹ for daily observation in the soft x-ray range. The scientific applications are expected to be very broad, covering numerous types and categories of variable and transient x-ray sources including x-ray binaries, AGN, blazars, Supernovae, x-ray counterparts of gamma ray bursts (including orphan afterglows), x-ray flashes, cataclysmic variables, and the like. Laboratory multifoil optical condenser allows extended use of x rays in microscopy and in EUV lithography.

48.6 REFERENCES

1. B. L. Henke, E. M. Gullikson, and J. C. Davis, "X-Ray Interactions: Photoabsorption, Scattering, Transmission, and Reflection at $E = 50\text{--}30000$ eV, $Z = 1\text{--}92$," *Atomic Data and Nuclear Data Tables*, **54**(2):181–342 (1993). This data, along with useful computational tools, is also available on the internet at http://www-cxro.lbl.gov/optical_constants.
2. P. Beckmann and A. Spizzichino. *The Scattering of Electromagnetic Waves from Rough Surfaces*, Oxford: Pergamon (1963).
3. P. Kirkpatrick and A. V. Baez, *J. Opt. Soc. Am.* **38**:776 (1948).
4. H. Wolter, *Ann. Phys.* **10**:94 (1952).
5. J. R. P. Angel, *Astrophys. J.* **233**:364–373 (1979).
6. W. H. K. Schmidt, *Nucl. Instr. Meth.* **127**:285 (1975).
7. L. P. Van Speybroeck, R. C. Chase, and T. F. Zehnpennif, *Appl. Opt.* **10**:945 (1971).
8. M. C. Weisskopf, *Appl. Opt.* **12**:1436 (1973).
9. J. W. Kast, "Scanning Kirkpatrick-Baez X-Ray Telescope to Maximize Effective Area and Eliminate Spurious Images," *Appl. Opt.* **14**:537 (1975).
10. P. Gorenstein, *SPIE* **3444**:382 (1998).

Marco W. Beijersbergen

Cosine Research B.V./Cosine Science & Computing B.V.
Leiden University
Leiden, Netherlands

49.1 INTRODUCTION

For x rays, the ratio of effective collecting area to the total reflecting surface area is small. This quantity, defined as the surface utility, for a Wolter-I optic, equals

$$s = R^2(\alpha, E) \sin(\alpha)/2 \quad (1)$$

where α is the grazing incidence angle on both mirrors in the Wolter-I optic, and $R(\alpha, E)$ the reflectivity at that grazing angle and the energy of the incident radiation. This quantity is only determined by the reflective properties of the surface, and is generally small. A typical value is 0.0015 for Au at 1 keV, which means that 1 cm² mirror only results in 0.0015 cm² effective collecting area. The challenge is therefore to create large surface area with high accuracy. For a given surface area the mass is determined by the density and thickness of the mirrors. However, thinner mirrors have lower stiffness and it is therefore more difficult to achieve a good figure and therefore angular resolution of the optic. An x-ray optic is therefore a trade-off between mass and resolution. The resolution can be improved by increasing the stiffness for a given mirror thickness. One possible solution is the use of a spacer that supports the mirrors. A more drastic implementation are pore optics, where the spacing of the spacers is of the same order as the spacing between the mirrors. There are currently three technologies that employ pore optics to implement x-ray pore optics: microchannel plate, silicon stacks, and micromachined silicon.¹⁻¹⁵

49.2 GLASS MICROPORE OPTICS

Introduction

Glass micropore optics are made from square glass fibers that are arranged in an appropriate configuration, after which the cores of the fibers are etched away such that thin mirrors remain. The resulting mirror thickness can be as low as 1 μm . With a glass density of about 2.3, this results in several orders of magnitude reduction in mass compared to 1-mm nickel shells or 0.1-mm-thick

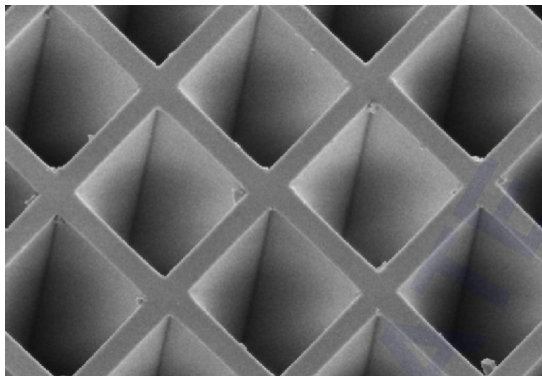


FIGURE 1 SEM image of a square-pore square-pack microchannel plate that acts as micropore optics. The pores are $20 \times 20 \mu\text{m}$, the wall thickness is $2 \mu\text{m}$.

aluminium foils. These optics can be produced with the same technology used for the manufacturing of microchannel plates, but using square instead of round glass blocks from which the fibers are drawn, as shown in Fig. 1. The process results in multifibers, as shown in Fig. 2, a square array of typically 50×50 fibers. These can be arranged in the desired geometry in a block. Plates are cut from the block and the core of the fibers are etched away, leaving square holes in a regular grid, as shown in Fig. 3.¹⁶

Glass micropore optics can be manufactured from different types of glass. In the etching process heavier elements are etched away preferentially and the resulting surface is mostly composed of SiO_x . It has not yet been possible to coat the inside of the pores uniformly with a material with a higher atomic number.

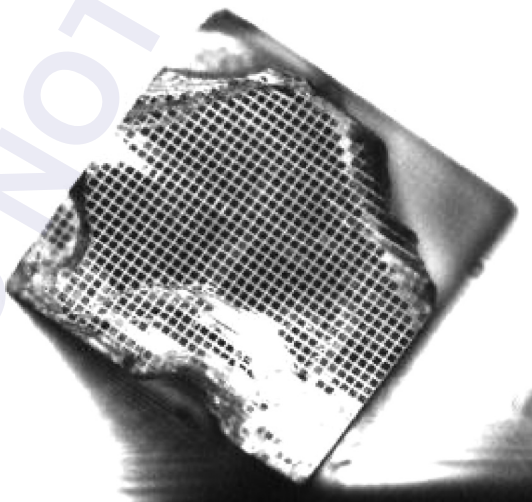


FIGURE 2 The front side of a multifiber with a size of $0.7 \times 0.7 \text{ mm}^2$. The front side of the unprocessed fiber is irregular and therefore not the entire area is properly focussed.

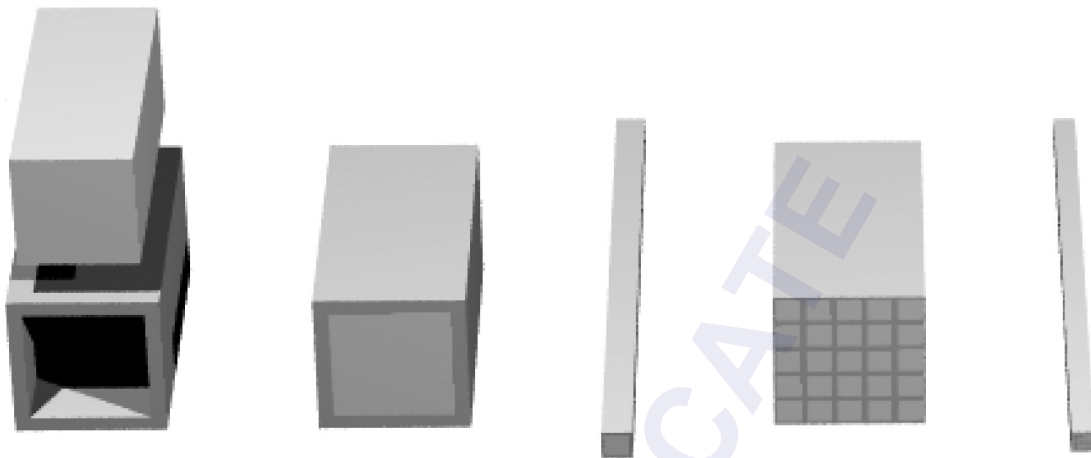


FIGURE 3 Manufacturing of a microchannel plate optics. A square core is inserted into a square cladding, and fibers are drawn. The fibers are stacked in a block, and multifibers are drawn.

In the drawing process the surface of the core glass is elongated by several orders of magnitude, thereby reducing the surface roughness significantly to about 1 to 2 nm. This makes them suitable for x-ray reflection. The typical pore size varies between 10 and 100 μm , and the wall thickness between 2 and 20 μm . This results in an open area ratio of typically 65 percent.

Lobster-Eye Optics

Multifibers can be arranged in a regular rectangular arrays, resulting in the so-called square-pore square-pack plates. If these are slumped over a spherical surface, a lobster-eye optic results,¹⁷ as shown in Fig. 4. Radiation that reflects twice inside the pore from orthogonal walls will be focused into a point. However, a large fraction of the incoming radiation will be reflected from one wall only

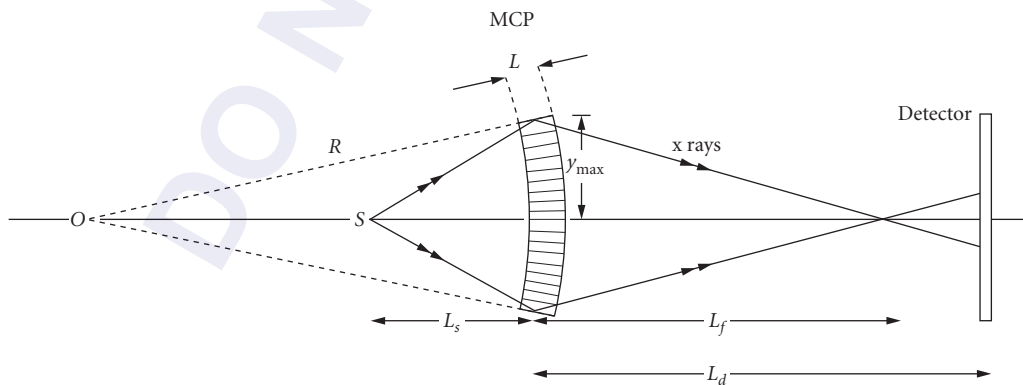


FIGURE 4 Creating a collimated beam with a slumped multichannel plate. (From Ref. 16.)



FIGURE 5 The crucifix image produced by a microchannel plate optic with square packing in a confocal arrangement. The half-energy width of the spot is about 6 arcmin. Data taken at the beam line of the Space Research Centre of Leicester University at 1 keV. (From Ref. 18.)

or go straight through the pores. The resulting image from a point source is called a crucifix image, as shown in Fig. 5.¹⁸ The crucifix point spread function (PSF) of the lobster optic results in a typical image resolution of a few arcminutes half-energy width. With a sufficiently large signal-to-noise ratio it is possible to deconvolve the image using the known PSF.

The lobster-eye optic only focuses light from an annulus on the optic. The collecting area will therefore be smaller than the total open fraction of the optic.

The field of view of a large curved lobster-eye optic is large. This optic is therefore well suited for sky survey instruments. A large-FOV instrument based on a lobster-eye optic is being developed for the International Space Station.¹⁹

The lobster-eye optic can be used for sources at finite distance, in which case a thin lens formula applies, as shown in Fig. 4,

$$\frac{1}{L_s} - \frac{1}{L_f} = \frac{2}{R_{\text{slump}}} \quad (2)$$

Glass micropore optics have been used to image x-ray fluorescence of biological and geological samples.^{20,21}

Wolter Optics

If the square fibers are stacked in concentric rings or sectors, as shown in Fig. 6, the radiation that is reflected from the azimuthal walls will result in a focus, as shown in Fig. 7. A second plate can be used to eliminate the extreme coma from a single plate, as shown in Fig. 8, resulting in proper imaging.^{22,23} In the case of a source at infinity the required slump radii are

$$f = \frac{1}{4}R_1 = \frac{3}{4}R_2 \quad (3)$$

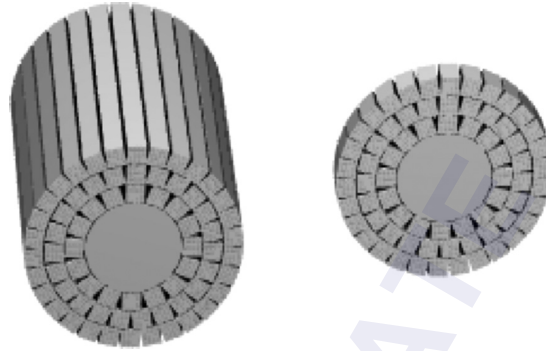


FIGURE 6 Radial microchannel plate optic are produced by stacking multifibers radially, cutting plates and etching the cores.

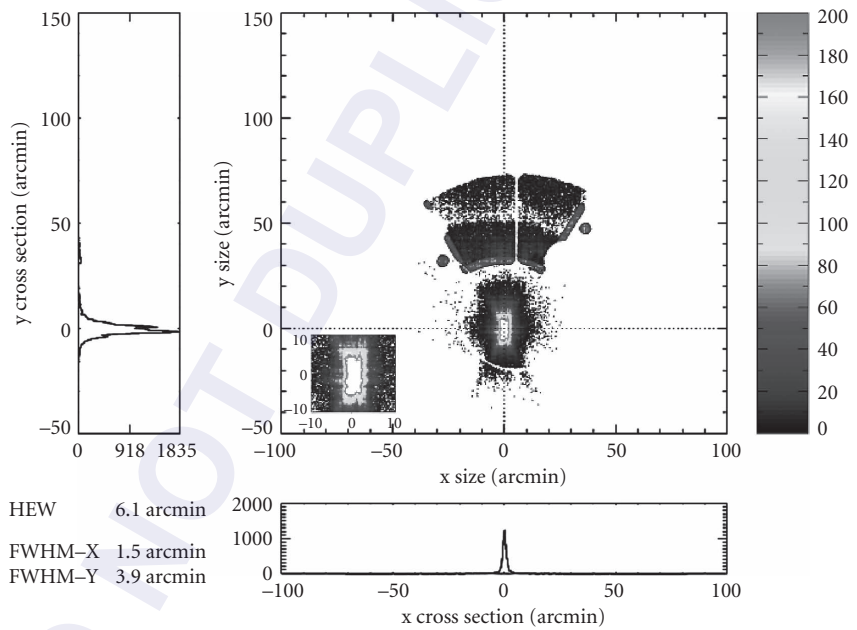


FIGURE 7 The focus of a single-plate radially packed microchannel plate optic. (See also color insert.)

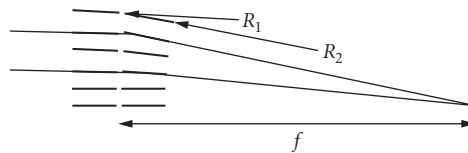


FIGURE 8 Two radially stacked plates that are slumped to different radii mimic a Wolter-I optical geometry and provide true focusing.

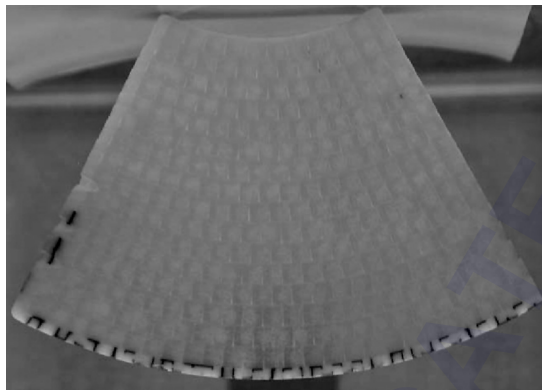


FIGURE 9 A segment of a radially stacked micropore optic. Two such plates behind each other act as a conical approximation to a Wolter-I optic. The width of the segment is 35 mm. (See also color insert.)

When the ratio of the pore size and thickness of the plate is properly chosen, the entire optic contributes to the focus of an on-axis source. The geometrical collecting area is therefore the open area ratio times the aperture.

The imaging performance of a glass micropore optic will be determined by the optical geometry, manufacturing errors that are imminent in the production of glass surfaces, as well as diffraction due to the pore structure. The focal plane is in the near field of the diffraction of the pores at typical x-ray wavelengths and Fresnel diffraction treatment is required to calculate the diffraction pattern. The actual geometry is basically a conical approximation to a Wolter-I. This will result in an image that has the size of a pore. However, in the case of concentric rings of multifibers, the multifibers will not be appropriately slumped in the azimuthal direction, and the focal spot will have the size of a multifiber. For a 1-mm multifiber size and a focal length of 1 m this corresponds to 3 arcmin. Diffraction at x-ray energies can be comparable depending on energy and pore size.²⁴

The optic can be created in off-axis segments, so that a larger collecting area in a single focus can be obtained, as shown in Fig. 9. Such an optic is being developed for an x-ray imaging spectrometer for BepiColombo, a European mission to Mercury^{25,26} and could be used for medium-resolution timing experiments.²⁷

49.3 SILICON PORE OPTICS

Silicon pore optics are produced by stacking silicon plates with ribs on the backside onto a conical mandrel, as shown in Fig. 10. Silicon wafers produced by the semiconductor industry are used. The chemomechanical polishing process that was developed for high-quality wafers results in a surface that has very good quality as x-ray mirror. Typical surface roughnesses are 0.1 nm and the flatness is better than a few nanometers over 25 mm. The plates are stacked such that the ribs make contact with the next plate and form a strong bond. This way stacks of a large number of concentric mirrors can be made. The porous structure, shown in Fig. 11, is almost as stiff as a solid silicon block. The x rays will reflect from the front side of the wafer toward the focus. A second stack is required to eliminate the severe coma caused by a single reflection. The current imaging resolution is 17 arcsec HEW over the full area and below 5 arcsec over selected regions.²⁸

Metallic and multilayer coatings can be applied to the surface inside the pore before stacking, provided that the coating is laid down in stripes, leaving uncoated silicon for the bonding.

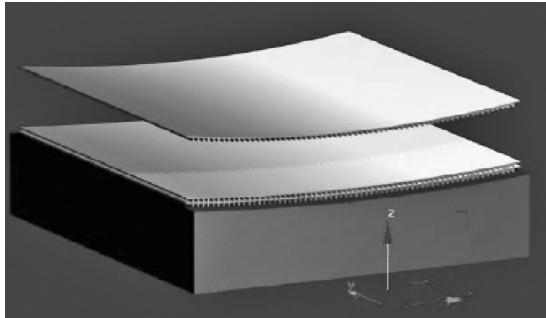


FIGURE 10 Silicon pore optics are produced by stacking silicon plates with ribs on the backside onto a concave mandrel.



FIGURE 11 A silicon pore optics produced by stacking ribbed silicon plates. The pore size is approximately $1 \times 1 \text{ mm}^2$, the width of the stack is 68 mm.

The resulting optic has a typical pore size of $1 \times 1 \text{ mm}^2$ and a typical mirror thickness of $170 \mu\text{m}$. The mass of these optics is about 1/10 of nickel replicated optics and is comparable to foil optics. The resolution is potentially better. This technology is well suited for conical approximations to a Wolter telescope. In that geometry the resolution will be inversely proportional to the focal length for a given pore size, and therefore this optic is best suited for large focal lengths. A closer approximation to a Wolter-I optic would require bending the plates in two directions, leading to additional distortion which has to be balanced against the benefit of the improved design.

The resulting stacks can be mounted as smaller modules into a larger structure. Note that the combination of two stacks forms a lens and therefore the alignment of a so-called tandem is much less critical than the alignment of an individual stack.

This optic is the baseline technology for the XEUS/IXO mission concept and would allow a collecting area of more than 3 m^2 in space.

49.4 MICROMACHINED SILICON

Careful etching of silicon can result in surfaces that are sufficiently smooth for x-ray reflection. With this method a structure can be etched into silicon that forms a geometry that can be used for imaging.²⁹ To get good surface quality requires etching along silicon crystal planes, which limits the geometries that can be realized.

49.5 REFERENCES

1. S. W. Wilkins, A. W. Stevenson, K. A. Nugent, H. Chapman, and S. Steenstrup, "On the Concentration, Focusing and Collimation of X-Rays and Neutrons Using Microchannel Plates and Configurations of Holes," *Rev. Sci. Instrum.* **60**:1026–1036 (1989).
2. H. N. Chapman, K. A. Nugent, and S. W. Wilkins, "X-Ray Focusing Using Square Channel-Capillary Arrays," *Rev. Sci. Instrum.* **62**:1542–1561 (1991).
3. G. W. Fraser, J. E. Lees, J. F. Pearson, M. R. Sims, and K. Roxburgh, "X-Ray Focusing Using Microchannel Plates," *Proc. SPIE* **1546**: 41–52 (1991).

4. P. Kaaret and P. Geissbühler, "Lobster Eye X-Ray Optics Using Microchannels Plates," *Proc. SPIE* **1546**:82–90 (1992); P. Kaaret, P. Geissbühler, A. Chen, and E. Glavinas, "X-Ray Focusing Using Microchannel Plates," *Appl. Opt.* **31**:7339–7343 (1991).
5. G. W. Fraser, A. N. Brunton, J. E. Lees, J. F. Pearson, and W. B. Feller, "X-Ray Focusing Using Square-Pore Microchannel Plates First Observation of Cruciform Image Structure," *NIM A* **324**:404–407 (1993).
6. G. W. Fraser, A. N. Brunton, J. E. Lees, and D. Lemberson, "Production of Quasi-Parallel X-Ray Beams Using Microchannel Plate X-Ray Lenses," *NIM A* **334**:579–588 (1993).
7. A. N. Brunton, G. W. Fraser, J. E. Lees, W. B. Feller, and P. L. White, "X-Ray Focusing with 11 μm Square Pore Microchannel Plates," in *X-Ray and Ultraviolet Sensors and Applications*, *Proc. SPIE* **2519** (1995).
8. I. C. E. Turcu, A. N. Brunton, G. W. Fraser, and J. E. Lees, "Microchannel Plate (MCP) Focusing Optics for a Repetitive Laser-Plasma Source," in *Applications of Laser Plasma Radiation II*, *Proc. SPIE* **2523** (1995).
9. A. G. Peele, K. A. Nugent, A. V. Rode, K. Gabel, M. C. Richardson, R. Strack and W. Siegmund, "X-Ray Focusing with Lobster-Eye Optics: A Comparison of Theory with Experiment," *Appl. Opt.* **35**:4420–4425 (1996).
10. A. N. Brunton, J. E. Lees, G. W. Fraser, and A. S. Tremsin, "MCP-Based X-Ray Collimators for Lithography of Semiconductor Devices," *Proc. SPIE* **2805**:212 (1996).
11. A. G. Peel and W. Zhang, "Lobster-Eye All-Sky Monitors: Comparison of One- and Two-Dimensional Designs," *Rev. Sci. Instrum.* **69**:2785–2793 (1998).
12. A. N. Brunton, A. P. Martin, G. W. Fraser, and W. B. Feller, "A Study of 8.5 μm Microchannel Plate X-Ray Optics," *NIM A* **431**:356–365 (1999).
13. M. W. Beijersbergen, M. Bavdaz, E. J. Buis, and D. H. Lumb, "Micro-Pore X-Ray Optics Developments and Application to an X-Ray Timing Mission," *Proc. SPIE Int. Soc. Opt. Eng.* **5488**:468 (2004).
14. M. Beijersbergen, S. Kraft, R. Gunther, et al., "Silicon Pore Optics: Novel Lightweight High-Resolution X-Ray Optics Developed for XEUS," *Proc. SPIE* **5488** (2004).
15. M. J. Collon, M. W. Beijersbergen, K. Wallace, M. Bavdaz, R. Fairbend, J. Séguy, E. Schyns, M. Krummy, and M. Freyberg, "X-Ray Imaging Glass Micro-Pore Optics," *Proc. SPIE* **6688**:668812 (2007).
16. G. W. Fraser, A. N. Brunton, J. E. Lees, J. F. Pearson, R. Willingale, D. L. Emberson, W. B. Feller, M. Stedman, and J. Haycocks, "Development of Microchannel Plate MCP X-Ray Optics," *Proc. SPIE* **2011**:215–226 (1993).
17. J. R. P. Angel, "Lobster Eyes as X-Ray Telescopes," *Astrophys. J.* **233**:364–373 (1979).
18. T. J. Norton, P. F. Morrissey, J. P. Haas, L. J. Payne, J. Carbone, and R. A. Kimble, "Photon-Counting Intensified Random-Access Charge Injection Device," in S. Fineschi, B. E. Woodgate, and R. A. Kimble (eds.), *Ultraviolet and X-Ray Detection, Spectroscopy, and Polarimetry III*, *Proc. SPIE* **3765**:452 (1999).
19. G. W. Fraser, A. N. Brunton, N. P. Bannister, J. F. Pearson, M. Ward, T. J. Stevenson, D. J. Watson, et al., "LOBSTER-ISS: An Imaging X-Ray All-Sky Monitor for the International Space Station," *Proc. SPIE* **4497**:115 (2002).
20. A. P. Martin, A. N. Brunton, G. W. Fraser, A. D. Holland, A. Keay, J. Hill, N. Nelms, et al., "Imaging X-Ray Fluorescence Spectroscopy Using Microchannel Plate X-Ray Optics," *X-Ray Spectrometry* **28**:64–70 (1999).
21. G. J. Price, G. W. Fraser, J. F. Pearson, J. P. Nussey, I. B. Hutchinson, A. D. Holland, K. Turner, and D. Pullan, "Prototype Imaging X-Ray Fluorescence Spectrometer Based on Microchannel Plate Optics," *Rev. Sci. Instrum.* **75**:2314 (2004).
22. M. W. Beijersbergen, M. Bavdaz, A. J. Peacock, E. Tomaselli, G. W. Fraser, and A. N. Brunton, "Novel Micropore X-Ray Optics Produced with Microchannel Plate Technology," *Proc. SPIE Int. Soc. Opt. Eng.* **4012**:218 (2000).
23. G. J. Price, A. N. Brunton, M. W. Beijersbergen, G. W. Fraser, M. Bavdaz, and J. -P. Boutot, "X-Ray Focusing with Wolter Microchannel Plate Optics," *NIM A* **490**:276 (2002).
24. A. L. Mieremet and M. W. Beijersbergen, "Fundamental Spatial Resolution of an X-Ray Pore Optic," *Appl. Opt.* **44**:7098–7105 (2005).
25. A. Owens, M. Bavdaz, M. W. Beijersbergen, A. N. Brunton, G. W. Fraser, D. Martin, P. Nieminen, A. J. Peacock, and M. G. Pia, "HERMES: An Imaging X-Ray Fluorescence Spectrometer for the BepiColombo Mission to Mercury," *Proc. SPIE* **4506**:136–145 (2001).

26. R. Willingale, G. W. Fraser, and J. F. Pearson, "Optimization of Square Pore Optics for the X-Ray Spectrometer on Bepi-Columbo," *Proc. SPIE* **5900**:590012 (2005).
27. M. Bavdaz, D. H. Lumb, A. Peacock, and M. Beijersbergen, "MCP-Optics for X-Ray Timing, X-Ray Timing 2003: Rossi and Beyond," *AIP Conf. Proc.* **714**:443-446, held 3-5 November, 2003 in Cambridge, Ma. P. Kaaret, F. K. Lamb, and J. H. Swank, (ed.). Melville, NY: American Institute of Physics (2004).
28. M. J. Collon, R. Günther, M. Ackermann, E. J. Buis, G. Vacanti, M. W. Beijersbergen, M. Bavdaz, K. Wallace, M. Freyberg, and M. Krumrey, "Performance of Silicon Pore Optics," *Proc. SPIE* **7011**:70111E (2008).
29. Y. Ezoe, M. Koshiishi, M. Mita, K. Mitsuda, A. Hoshino, Y. Ishisaki, Z. Yang, T. Takano, and R. Maeda, "Micropore X-Ray Optics Using Anisotropic Wet Etching of (110) Silicon Wafers," *Appl. Opt.* **45**:8932-8938 (2006).

This page intentionally left blank.

DO NOT DUPLICATE

ADAPTIVE X-RAY OPTICS

Ali Khounsary

Argonne National Laboratory
Argonne, Illinois

50.1 INTRODUCTION

The field of adaptive optics (AO) has its origin in the second half of the twentieth century^{1,2} in connection with improving the image quality of celestial objects taken by Earth-based telescopes. Prior to this, and for some 300 years, telescopes had evolved significantly, from primitive constructs to sophisticated systems that incorporated the latest in design, optical material, polishing, and assembly.³ Further improvement in image quality required collecting more light and, more importantly, dealing with the blurring caused by atmospheric turbulence. These two requirements have led to the developments of active and adaptive optics, respectively.

To collect more light, telescopes with large primary mirrors several meters in diameter were proposed. To control optical aberrations due to mechanical and thermal effects in such large systems, *active optics*⁴ techniques were developed and implemented that allowed automatic adjustments via built-in corrective optical elements operating at fairly low temporal frequency ($\ll 1$ Hz).

To address the degrading effects of atmospheric turbulence and to obtain sharper images, *adaptive optics* was developed and implemented.⁵ Using a bright star or a beacon,⁶ wavefront distortions due to atmospheric turbulence are measured in real time, and this information is fed into a system of many high-speed (kHz range) actuators operating on a deformable mirror that adjusts its shape, as shown in Fig. 1, to correct for the atmospheric distortion to produce sharper images. A deformable mirror of this type is typically about 10 to 20 cm in diameter, located behind the telescope focus, and its axial range of deformation is typically on the order of a few microns.

Developments in adaptive optics in the visible or IR region of the spectrum have continued with military, communication, astronomy, energy, and medical applications (see Chap. 5). These efforts are now being extended to the x-ray region, notably in x-ray astronomy and synchrotron-based x-ray systems.

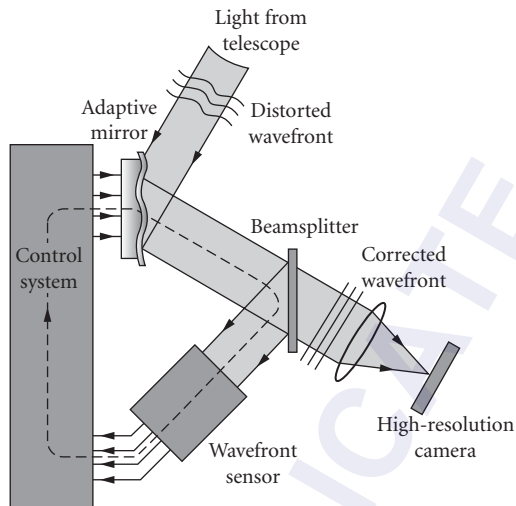


FIGURE 1 Light from the celestial object of interest and a nearby “guide star” passes through the telescope optics and propagates within the adaptive system shown. The light from the star is continuously measured by a high-speed wavefront sensor and analyzed to determine atmospheric distortion. This information is used to change the shape of the deformable mirror in order to cancel out the distortion due to atmosphere, which allows the high-resolution camera to obtain sharper images of the stars and celestial objects. (Courtesy of Claire Max, Center for Adaptive Optics, University of California at Santa Cruz, California.)

50.2 ADAPTIVE OPTICS IN X-RAY ASTRONOMY

Unlike optical astronomy, space-based hard x-ray astronomy relies on grazing incident optics to collect an image. This severely restricts the amount of light that can be collected because the available optical surface is limited. To address this problem, light-weight optics composed of a large number of nested confocal mirrors with very thin walls have been developed (see Chap. 47). Mirrors are individually produced by electroforming on a suitably shaped mandrel and then releasing them.^{7,8} For high-angular-resolution hard x-ray astronomy, telescopes with rigid optics are presently preferred, as they maintain their figure during the rigor of launch and beyond. The light collection area is severely restricted, however, due to weight considerations. Research is continuing in the use of adaptive optics techniques for in situ control and figure correction.^{9,10} For soft x-ray telescopes operating near normal-incident angles, on the other hand, adaptive optics is being developed to correct figure errors of the primary mirror in an effort to approach diffraction limit performance.¹¹

50.3 ACTIVE AND ADAPTIVE OPTICS FOR SYNCHROTRON- AND LAB-BASED X-RAY SOURCES

The terms active optics and adaptive optics have distinct meanings in the *optical telescope community*, denoting low-speed and high-speed correction systems, respectively. Elsewhere, this distinction is not

always observed. For example, within the *x-ray optics community*, the terms active, adaptive, bendable, and deformable are used interchangeably, yet only a few systems are truly adaptive in the sense of high-speed automatic sensing, feedback, and correction. In what follows, the term adaptive is used to denote optics that during operation undergo complex deformation or rigid body motions that require (manual or automatic) monitoring and active feedback, or optics that are of the self-correcting type, which have built-in self-adjustment mechanisms by design. Other systems are termed active optics.

In recent years, several groups of researchers have begun exploring adaptive optics technologies for application to synchrotron and lab-based x-ray optics. This trend is expected to accelerate because of its impact on beamline performance and throughput. More significantly, because of the relative ease with which adaptive technologies can be implemented at the so-called mezzzo (sub-mm) level (using MEMS, micro-sensors, and related processing and technologies), an entirely new generation of x-ray optical systems could result.

The ultimate goal would be to develop a new class of enabling adaptive optical systems that can, in real time, sense and compensate for undesirable changes (a) in the optics itself, (e.g., due to beam heat load), (b) in the laboratory environment (e.g., thermal, mechanical), and (c) possibly in the x-ray source (temporal and spatial). As an example, it would be possible to focus and keep a nanofocused x-ray beam on a small sample for an extended period of time, a goal that might best be achieved not by expansive control of a long beamline but by active/adaptive control of the optics.

Major x-ray optical systems currently in use on synchrotron x-ray beamlines or in lab-based x-ray systems include monochromators, mirrors, multilayer-coated mirrors, zone plates, polycapillaries, and compound refractive lenses. Many are active optics, as they have provisions for angular and spatial movements and adjustments. They generally lack automatic sensing, feedback, and control because, unlike optical astronomy, which requires an automatic and relatively high-speed feedback system on the order of kHz, many laboratory and synchrotron x-ray optics rely on painstakingly slow operator feedback and adjustments. In many instances “sensing” is difficult, labor intensive, or complex, and not amenable to effective automation. Thus, of the three components of an adaptive optical system (sensor, optics, and control), the sensor remains most problematic. Currently, for example, when bending a mirror to focus an x-ray beam, intensive operator intervention, in the form of repeated knife-edge scans and bending adjustments, is necessary to arrive at the optimal mirror shape for finest focus. With an adaptive x-ray system, the mirror surface profile could be automatically measured (optically or directly by x rays) and the bending adjusted to arrive at the best surface profile.

Thus, incorporating adaptive techniques in synchrotron- and lab-based x-ray optics optical systems is highly desirable.

Progress in three cases is described in this chapter:

- Bendable optics: Controlled elastic bending of optical substrates into versatile shapes necessary to focus, collimate, or collect x rays.
- Figure correction: Controlled correction of long spatial wavelength figure errors (due to manufacturing, thermal load, mounting, gravity, etc.).
- Novel optics: Development of new types of optical systems, such as x-ray telescopes for both tabletop and synchrotron x-ray sources, where adaptive techniques form the basis of the design rather than an enhancing feature.

Bendable Optics

Because of the widespread need to collect and focus x rays on one hand, and the availability of higher-quality flat substrates on the other, a number of bendable optical systems (composed of bent mirrors, monochromators, or both) have been developed.^{12–17} An additional impetus is the difficulty and expense in manufacturing high-quality mirrors with elliptical or higher-order polynomial profiles.^{18–21} A variety of bending tools and techniques using electrical, mechanical, and thermal means have been developed, although most are mechanically bent by applying moments at or near the ends of the mirrors, as shown in Fig. 2.

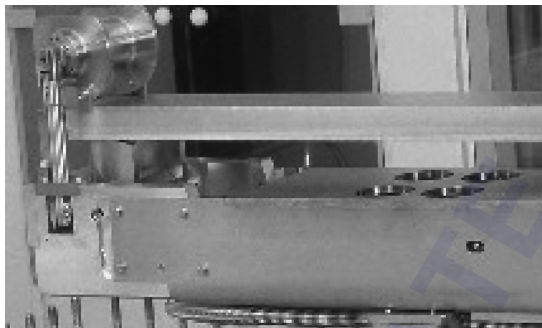


FIGURE 2 A mirror bender designed by IRELEC for ESRF. The system applies two controlled bending moments, by means of two electrical actuators, at each end of a mirror. The moments are controlled independently to be able to bend the optics into elliptical or cylindrical profiles. (Courtesy of IRELEC, France.)

The design principle underlying bendable optics is simple elastic beam theory. It is well known that a flat beam of uniform cross section can be bent to the arc of a circle of radius R given by $(1/R) = (M/EI)$, where M is the applied bending moment at its ends, E is the modulus of elasticity, and I is the beam's moment of inertia.²² This one-dimensional treatment is applicable to long mirror substrates where the length is much larger than other dimensions. Application of two equal moments at the ends produces a circular arc; two unequal moments produce an ellipse. More generally, this equation can be written as a function of location x along the length of the mirror and measured from an arbitrary origin as

$$\frac{1}{R(x)} = \frac{M(x)}{EI(x)}$$

where $R(x)$ is the *local* radius of curvature. Thus, a homogeneous elastic substrate may be bent into a range of longitudinal profiles by the appropriate application of moments along its length.^{14,18,23} Alternatively, because of its dependence on moment of inertia, even with the application of a constant moment M , the curvature $R(x)$ can be made to vary along the mirror length by using a substrate with a tailored nonuniform cross section (width or thickness). Thus, within limits, a flat mirror can be bent into one with an arbitrarily varying longitudinal surface profile, and indeed many such optics have been made and are in operation.^{20, 24–26}

With the few exceptions, these bendable systems do not employ adaptive technology, although most of the necessary apparatus, except for an automatic sensing system, are in place. An automatic sensing and feedback system could provide real-time enhanced dynamic bending options and automatic adjustments in response to environmental changes, and the capacity to correct long wavelength tangential slope errors in systems with multiple actuators.

Although bendable optical systems have successfully produced focused x-ray beams down to submicron size, it should be noted that they have some practical limitations.^{18,27,28} First, they are sensitive to the bender adjustment, especially when adaptive techniques are not used. Second, they are bulky and harder to cool in comparison to the prefigured alternative. Lastly, in practice, bent mirrors have not been able to match the precise figures possible with prefigured rigid optics to produce nanometer-size focal spots.^{29,30}

Adaptive X-Ray Optics

For the reasons described earlier, only a few truly adaptive x-ray optical systems that use actuators and robust sensing techniques, such as Shack-Hartmann wavefront sensing, have been developed to date. A few of the main systems are described next.



FIGURE 3 A bimorph mirror, 600 mm in length, composed of four bimorph plates each with four electrodes, installed at the GM/CA-CAT beamline at the Advanced Photon Source in the United States.³⁶ (Courtesy of SESO, France.)

Bimorph Mirrors One of the promising adaptive optical systems developed for x-ray applications is the bimorph mirror system. Piezoelectric bimorph mirrors³¹ were first proposed for adaptive control of laser³² and astronomical mirrors,³³ and prototype systems were built.^{34,35} A bimorph mirror consists of two faceplates with a sheet of active elements sandwiched symmetrically between them. The active elements are made of two piezoelectric plates glued together with their polarization vectors in the same direction perpendicular to the plate face. Each plate is coated with a thin metallic electrode at the glued surface. When a voltage is applied at the interface electrodes, one plate contracts and the other expands, bending the mirror assembly spherically. To change the radius, the applied voltage is changed. To change the mirror shape locally, a number of smaller embedded piezoelectric plates can be used to provide versatile shape control, as shown in Fig. 3.

In the early 1990s, the European Synchrotron Radiation Facility (ESRF) began exploring active and adaptive technologies in connection with thermally induced deformation in cooled x-ray mirrors for high-heat-load x-ray beamlines.^{37,38} A system composed of multiple discrete piezoelectric actuators and a Shack-Hartmann wavefront sensor³⁹ was built, as shown in Fig. 4.⁴⁰ It was recognized, however, that thermally induced deformations are rather smooth and slow varying along the mirror length, and their correction requires only a gentle bending for which a bimorph mirror with one or a small number of electrodes is most appropriate.⁴¹ For focusing x rays that require an elliptically bent mirror, a bimorph composed of a few electrodes is sufficient;⁴² short bimorph mirrors were fabricated and bent to an elliptical shape in a Kirkpatrick-Baez geometry⁴³ (see Chap. 44) to focus x-ray beams in two directions.⁴⁴ The design has been evolving over time and improvements are gradually being made in fabrication, control, feedback, and performance.^{45–48} Presently about 150 of these bimorph mirrors are in use on synchrotron x-ray beams worldwide, a number that is likely to grow as they are recognized as an economical approach to adaptive optics for synchrotrons and x-ray astronomy.⁴⁹

It is important to note that one of the severe limitations on the use of bimorph mirrors at synchrotron facilities is that the mirror temperature must remain below 60°C, and thus this technology is only suitable for monochromatic beams or other low thermal load beams.

Adaptive Optics for Thermal and Environmental Control and Correction As noted, one potential use of adaptive optics in synchrotron and perhaps lab-based x-ray systems is for the control and correction of surface deformations due to beam thermal load, gravity, or environmental factors. Thermally induced deformations are particularly important, and an adaptive optics system that allows cooling to take place while providing dynamic correction is highly desirable. For reference,

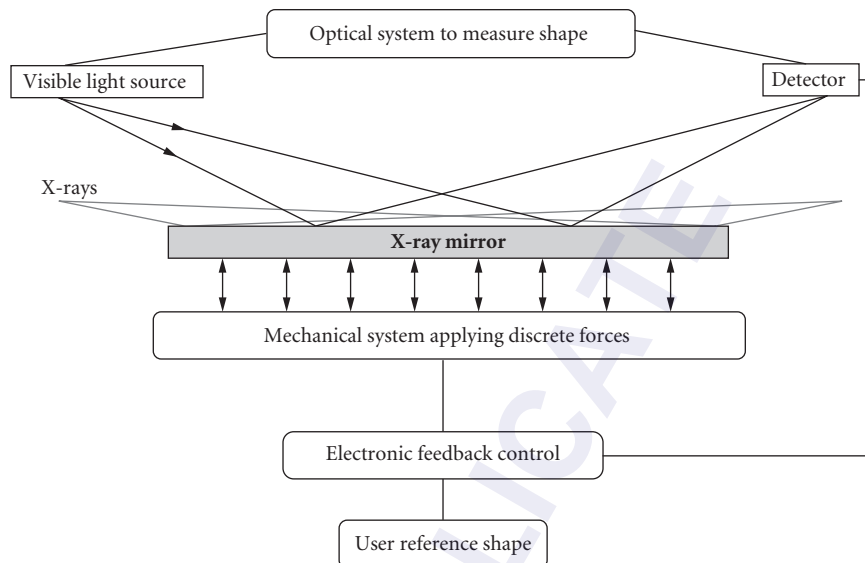


FIGURE 4 Schematic of a cooled adaptive mirror system for synchrotron x-ray applications. A Shack-Hartmann wavefront sensor is used to measure the thermally induced deformation of the optical surface and provide feedback to 22 (in two rows of 11) piezoelectric actuators that flex the 1-m-long mirror. Each actuator has a maximum stroke of 40 micrometers.⁴⁰ (Courtesy of Jean Susini, ESRF Grenoble, France.)

it should be noted that the third-generation synchrotron x-ray facilities that have come online since the early 1990s generate very small but powerful x-ray beams with low divergence⁵⁰ but with thermal loads from a few hundred to thousands of watts. Optical substrates subjected to high-heat-load x-ray beams typically deform into a convex shape. For reflective optics, this results in increased beam divergence; for refractive optics, it results in loss of x-ray beam intensity. These optical elements, therefore, must be cooled to remove the heat, and their adverse thermal deformations corrected, internally or externally. While most corrections require external means, under some circumstances it is possible to build into the design of the system some untraditional adaptive features such that the sensing, feedback, and correction are performed by the optics itself.⁵¹ In one design, two properly devised cooling blocks are optimally positioned on the long sides of a mirror close to the reflecting surface. When a high-heat-load x-ray beam strikes the mirror at grazing angles, the heat is deposited along a narrow longitudinal footprint on the optical surface. The mirror immediately deforms into a convex shape, but within several minutes a reverse thermal moment is automatically generated within the mirror body that reverses the initial curvature and largely flattens the optical surface.^{51,52}

In general, however, thermal distortions cannot be mitigated by internal reversal, and traditional adaptive techniques must be implemented to correct the deformations. One such system was developed at the ESRF.⁴⁰ In this system, shown in Fig. 4, a Shack-Hartmann sensor monitors the mirror surface to determine thermally induced deformations and provides feedback to a system of 22 piezoelectric actuators to correct the figure of the thermally deformed mirror.

X-Ray Monochromators As noted, x-ray monochromators (see Chaps. 30 and 39) installed on high-heat-load x-ray synchrotron beamlines can undergo thermal deformation that effectively reduces the intensity of the diffracted beam. Unlike mirrors, however, monochromators intercept the beam typically at large angles and absorb most of the white incident beam, and thus they are subjected to a substantially higher heat flux. They must therefore be efficiently cooled and their thermal

deformation managed. This thermal management issue is likely to become more acute as a result of the upgrades planned at synchrotron facilities. Presently most high-heat-load silicon monochromators, by necessity, employ cryogenic cooling around 125 K to exploit silicon's high thermal conductivity and nearly zero thermal expansion coefficients around this temperature. Thermal deformation is insignificant, requiring little or no corrections. However, the cooling systems are somewhat complex and water cooling is preferred, if possible, especially for lower-heat-load beamlines. A number of "adaptive" techniques to overcome undesirable thermal distortions were explored earlier. For example, in one design, a silicon monochromator with a thin diffracting surface was water-cooled with water jets on its back while pressurized helium was used on the front side of the diffracting surface to flatten its convex shape and increase diffracted intensity.⁵³ In another design, a silicon monochromator, consisting of a thin crystal faceplate with small cooling channels in it, was bonded to and supported by beryllium through which some 20 actuators were used to flatten the faceplate.⁵⁴

While many investigators have examined adaptive techniques for the control of thermal effects, some have examined active control of an optical system thermally. For example, to increase the throughput or resolution in a monochromator subjected to a divergent incident beam, heat can be applied to generate a thermal gradient across the face of the diffracting crystal such that the d-spacing variation matches incident angles of the divergent beam, i.e., $\lambda = 2d \sin \theta_B$ is a constant. This is useful when the divergences of a beam appreciably exceed the acceptance angles of the monochromator used. Gains of up to two orders of magnitude in intensity, relative to standard flat isothermal crystal systems, are expected.⁵⁵ A number of other schemes to use thermal gradients to actively control optics have also been suggested or implemented with varying degrees of success.^{56,57}

Novel Optics

Microchannel Plates (MCPs) Microchannel plates, also known as "micro pore optics" (MPOs), (see Chap. 49) were originally developed to provide high-resolution images in low-light conditions (e.g., night vision goggles) using image intensifiers. More recently, they have been used in x-ray detectors and have also been investigated for x-ray focusing⁵⁸ in a manner similar to the lobster eye vision proposed earlier.⁵⁹ A typical MCP is a planar glass, on the order of 1 mm thick, composed of thousands of small channels. To focus x rays using MCPs, these channels must be curved. This can be accomplished by spherically slumped or adaptively controlled bending of the plate so that incident rays on the channel walls of the MCP are reflected onto a screen producing an image of the object. Developments in this field continue through various design improvements, smoother walls, and more accurate channel alignment. Lobster-ISS (International Space Station) is one such system that is expected to be deployed in 2010 aboard the ISS.^{60,61} While lobster x-ray optics was originally intended for astronomical applications to view a broad swath of the sky, these compact x-ray focusing devices are being explored for lab-based sources and possible synchrotron use.⁶²

Adaptive Microstructured Optical Arrays (MOAs) Microstructured optical arrays are a new class of x-ray focusing optical systems under development that is conceptually similar to polycapillary and microchannel plane optics but differs in that x-rays are guided by a single (rather than multiple) grazing angle reflection from each MOA.^{63,64} In a sense, MOAs are a discretized version of a polycapillary optics, allowing considerable flexibility in the design of the channels and layout of the arrays. Figure 5 shows a sketch of a one-dimensional MOA system composed of two arrays focusing x rays from S to S' . The second array shown is bent. Bending one or both arrays provides a variable focal length. Adaptive control, for example, of the piezoelectric material coated on the support structure of the arrays allows bending, changing the dimensions of the channels, or correcting optical aberrations. Preliminary MOAs are being made by the deep silicon etching process.⁶⁵ MOAs hold considerable potential for x-ray optics in general and adaptive optics in particular. For example, it is estimated that an optimally designed MOA could provide two orders of magnitude more focused flux than the corresponding zone-plate, primarily due to the larger aperture afforded by an MOA.⁶⁵

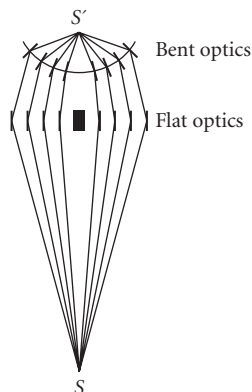


FIGURE 5 Schematic of a microstructured optical array system used to focus x rays from a source S to focus S' . This one-dimensional system is composed of two arrays, the second of which (the top one) is bent. (Courtesy of Alan Michette, Kings College, U.K.)

50.4 CONCLUSIONS

Adaptive x-ray optics, as a field, is in its early stages of development. Most x-ray optical systems in use today are passive because stability is valued more than versatility.

However, development of robust adaptive systems is likely to provide the means to achieve enhanced throughput, resolution, stability, and versatility. Concerted effort is needed to develop cost-effective and reliable systems that provide real improvement over the existing optical systems. Actuation mechanisms (thermal, piezoelectric, magnetic, electrostatic, motor driven) and sensing techniques suitable in a variety of x-ray environments must be developed. Development of novel x-optics based on adaptive techniques would enhance throughput and versatility of both synchrotron- and lab-based systems.

50.5 REFERENCES

1. H. W. Babcock, "The Possibility of Compensating Astronomical Seeing," *Pub. Astr. Soc. Pac.* **65**:229–236 (1953).
2. V. P. Linnik, "On the Possibility of Reducing the Influence of Atmospheric Seeing on the Image Quality of Stars," (1957) article translated and reprinted in F. Merkle, ed., *ICO-16 Satellite Conference on Active and Adaptive Optics*, Vol. 48 of *ESO Conference and Workshop Proc.*, 535–537, ESO, Garching, 1993.
3. R. N. Wilson, *Reflecting Telescope Optics I*, 2nd ed., Springer, Berlin, 2007.
4. R. K. Tyson, (ed.), *Adaptive Optics Engineering Handbook*, Marcel Dekker, Inc., New York, 2000.
5. J. M. Beckers, "Adaptive Optics for Astronomy: Principles, Performance, and Applications," *Ann. Rev. Astron. Astrophys.* **31**:13–62 (1993).
6. N. Hubin and L. Noethe, "Active Optics, Adaptive Optics, and Laser Guide Stars," *Science* **262**:1390–1394 (1993).
7. O. Citterio, G. Bonelli, G. Conti, E. Mattaini, and B. Sacco, "High Throughput Optics for X-Ray Astronomy," *IL Nuovo Cimento* **13**(2):375–389 (1990).

8. R. Shibata, Y. Ogasaka, K. Tamura, A. Furuzawa, Y. Haba, Y. Tawara, H. Kunieda, et al., "Hard X-Ray Mirrors by Multiplayer Replication: Development and Applications," *Proc. 8th Int. Conf. X-Ray Microscopy*, pp. 159–161, 2006.
9. M. Ulmer, Private Communication, 2008.
10. O. Cugat, S. Basrour, C. Divoux, P. Mounaix, and G. Reyne, "Deformable Magnetic Mirror for Adaptive Optics," *Sensors and Actuators A*. **89**:1–9 (2001).
11. S. Kitamoto, H. Takano, H. Saitoh, N. Yamamoto, T. Kohmura, K. Suga, H. Sekiguchi, Y. Ohkawa, J. Kanai, and S. Chiba, "Development of an Ultra-High Precision X-Ray Telescope, Stellar-Mass, Intermediate-Mass, and Supermassive Black Holes," S. Mineshige and K. Makishima (eds.) *Prog. of Theor. Phys. Suppl.* **155**:363–364 (2004).
12. W. Ehrenberg, "X-Ray Optics," *Nature* **160**:4062 (1947).
13. W. Ehrenberg, "X-Ray Optics: The Production of Converging Beams by Total Reflection," *J. Opt. Soc. Am.* **39**:741–746 (1949).
14. J. A. Howell and P. Horowitz, "Ellipsoidal and Bent Cylindrical Condensing Mirrors for Synchrotron Radiation," *Nucl. Instrum. Meth.* **125**:225–230 (1975).
15. *Proc. Workshop on X-Ray Instrumentation for Synchrotron Radiation*, SSRL Report No. 78/04, VII/36–37, 1978.
16. J. B. Leigh and G. Rosenbaum, "Synchrotron X-Ray Sources: A New Tool in Biological Structural and Kinetic Analysis," *Annu. Rev. Biophys. Bioeng.* **5**:239–270 (1976).
17. G. E. Ice and C. J. Sparks, "Conical Geometry for Sagittal Focusing as Applied to X-Rays from Synchrotrons," *Oak Ridge National Laboratory Report*, ORNL/TM-12327, 1993.
18. S. J. Chen, C. K. Kuan, S. Y. Perng, D. J. Wang, H. C. Ho, T. C. Tseng, Y. C. Lo, and C. T. Chen, "New Focusing Mirror System for Synchrotron Radiation Infrared Beamlines," *Opt. Eng.* **43**:3077–3082 (2004).
19. H. A. Padmore, M. R. Howells, S. Irick, T. Renner, R. Sandler, and Y.-M. Koo, "New Schemes for Producing High-Accuracy Elliptical X-Ray Mirrors by Elastic Bending," *Proc. SPIE* **2856**:145–156 (1996).
20. P. Eng, M. Newville, M. L. Rivers, and S. R. Sutton, "Dynamically Figured Kirkpatrick Baez X-Ray Microfocusing Optics," *Proc. SPIE* **3449**:145–156 (1998).
21. B. W. Adams and K. Attenkofer, "An Active-Optic X-Ray Fluorescence Analyzer with High Energy Resolution, Large Solid Angle Coverage, and a Large Tuning Range," *Rev. Sci. Instrum.* **79**:023–102 (2008).
22. W. C. Young, *Roark's Formula for Stress and Strain*, 6th ed., McGraw-Hill, New York, 1989.
23. M. R. Howells and D. Lunt, "Design for an Adjustable-Curvature, High-Power, X-Ray Mirror Based on Elastic Bending," *Opt. Eng.* **32**:1981–1989 (1993).
24. D. Turner and J. M. Bennett, "An Elliptical Reflector Formed by Bending a Cantilever," Imperial College (1971).
25. J. H. Underwood, "Generation of a Parallel X-Ray Beam and Its Use in Testing Collimators," *Space Sci. Instrum.* **3**:259–270 (1977).
26. G. E. Ice and C. J. Sparks, "A Simple Cantilevered Mirror for Focusing Synchrotron Radiation," *Nucl. Instrum. Meth. A* **266**:394–398 (1988).
27. D. Lunt, J. Bender, D. W. Ewing, and W. R. McKinney, "XUV Synchrotron Optical Components for the Advanced Light Source," *Proc. SPIE* **1740**:161–172 (1993).
28. M. R. Howells, D. Camble, R. M. Duarte, S. Irick, A. A. MacDowell, H. A. Padmore, T. R. Renner, et al., "Theory and Practice of Elliptically Bent X-Ray Mirrors," *Opt. Eng.* **39**:2748–2762 (2000).
29. W. Liu, G. E. Ice, J. Z. Tischler, A. Khounsary, C. Liu, L. Assoufid, and A. T. Macrander, "Short Focal Length Kirkpatrick-Baez Mirrors for a Hard X-Ray Nanoprobe," *Rev. Sci. Instrum.* **76**:113701 (2005).
30. H. Mimura, H. Yumoto, S. Matsuyama, Y. Sano, K. Yamamura, Y. Mori, and M. Yabashi, "Efficient Focusing of Hard X Rays to 25 nm by a Total Reflection Mirror," *Appl. Phys. Lett.* **90**:051903 (2007).
31. J. H. McElroy, P. E. Thompson, H. E. Walker, E. H. Johnson, D. J. Radecki, and R. S. Reynolds, "Laser Tuners Using Circular Piezoelectric Benders," *Appl. Opt.* **14**:1297–1302 (1972).
32. P. V. Nikolaev and A. V. Smirnov, "Model Wavefront Correctors," *Sov. J. Opt. Technol.* **54**:693–700 (1987).
33. F. Roddier, "A New Concept in Adaptive Optics: Curvature Sensing and Compensation," *Appl. Opt.* **27**:1223–1225 (1998).
34. F. Roddier and F. Forbes, "Curvature Sensing and Compensation," in *Adaptive Optics in Solar Observation*, O. Envold (ed.), p. 176, 1987.
35. A. V. Ikramov, S. V. Ramanov, I. M. Roshchupkin, A. G. Safronov, and A. O. Sulimov, "Bimorph Adaptive Mirror," *Sov. J. Quantum Electron.* **22**:163–166 (1992).

36. R. F. Fischetti, D. W. Yoder, S. Xu, S. Stepanov, O. Makarov, R. Benn, S. Corcoran, et al., "Optical Performance of the GM/CA-CAT Canted Undulator Beamlines for Protein Crystallography," *Synch. Rad. Instrum* **879**:754–757 (2007).
37. J. Susini, G. Marot, L. Zhang, R. Ravelet, and P. Jagourel, "Conceptual Design of an Adaptive X-Ray Mirror Prototype for the ESRF," *Rev. Sci. Instrum.* **63**:489–492 (2002).
38. J. P. Gaffard, R. Ravelet, and C. Boyer, "X-Ray Adaptive Mirror: Principle and State of the Art," *Proc. SPIE* **1739**:474–488 (1992).
39. J. Schwiegerling and D. R. Neal, "Historical Development of the Shack-Hartmann Wavefront Sensor," in *Robert Shannon and Roland Shack*, J. E. Harvey and R. B. Hooker (eds.), *SPIE* 132–139, 2005.
40. J. Susini, R. Baker, M. Krumrey, W. Schwegle, and A. Kwick, "Adaptive X-Ray Mirrors Prototype: First Results," *Rev. Sci. Instrum.* **66**:2048–2052 (1995).
41. J. P. Gaffard and P. Jagourel, "An Active Bimorph Structure for X-Ray Gratings and Mirrors," *Proc. SPIE* **2856**:197–206 (1997).
42. J. Susini, D. Laberge, and L. Zhang, "Compact Active/Adaptive X-Ray Mirror: Bimorph Piezoelectric Flexible Mirror," *Rev. Sci. Instrum.* **66**:2229–2231 (1995).
43. P. Kirkpatrick and A. V. Baez, "Formation of Optical Images by X-Rays," *J. Opt. Soc. Am.* **38**:766–774 (1948).
44. J. Susini, D. Laberge, and O. Hignette, "R & D Program on Bimorph Mirrors at the ESRF," *Proc. SPIE* **2856**:130–144 (1997).
45. R. Signorato, "R & D Program on Multi-Segmented Piezoelectric Bimorph Mirrors at the ESRF: Status Report," *Proc. SPIE* **3447**:20–31 (1998).
46. R. Signorato and T. Ishikawa, "R & D on Third Generation Multi-Segmented Piezoelectric Bimorph Mirror Substrates at SPring8," *Nucl. Instrum. Meth. A* **467–468**:271–274 (2001).
47. T. C. Tseng, S. -J. Chen, Z. C. Yen, S. Y. Perng, D. J. Wang, C. K. Kuan, J. R. Chen, and C. T. Chen, "Design and Fabrication of Aspherical Bimorph PZT Optics," *Nucl. Instrum. Meth. A* **467–468**:294–297 (2001).
48. R. Signorato, D. Haesermann, M. Somayazulu, and J. F. Carre, "Performance of an Adaptive μ -Focusing Kirkpatrick-Baez System for High Pressure Studies at the Advanced Photon Source," *Proc. SPIE* **5193**:112–123 (2004).
49. P. Doel, C. Atkins, S. Thompson, D. Brooks, J. Yao, C. Feldman, R. Willingale, T. Button, D. Zhang, and A. James, "Development of Piezoelectric Actuators for Active X-Ray Optics," *Proc. SPIE* **6705**:7050M (2007).
50. E. E. Koch, *Handbook on Synchrotron Radiation*, Vol. 1a, Elsevier Science Publishers, B. V. North Hollands, The Netherlands, 1983.
51. A. Khounsary and W. Yun, "On Optimal Contact Cooling of High-Heat-Load X-Ray Mirrors," *Rev. Sci. Instrum.* **67**:3354 (1996).
52. Y. Li, A. M. Khounsary, and S. Nair, "How and Why Side Cooling of High-Heat-Load Optics Works," *Proc. SPIE* **5533**:157 (2004).
53. L. E. Berman and M. Hart, "Adaptive Crystal for High Power Synchrotron Sources," *Nucl. Instrum. Meth. A* **302**:558–562 (1991).
54. D. Dezoret, R. Marmoret, A. K. Freund, A. Kwick, and R. Ravelet, "Design of an Adaptive Cooled First Crystal for an X-Ray Monochromator," *Proc. SPIE* **2279**:544–549 (1994).
55. G. S. Knapp and R. K. Smither, "High Resolution Monochromator Systems using Thermal Gradient Induced Variable Bragg Spacing," *Nucl. Instrum. Meth. A* **246**(1–3):365–367 (1986).
56. R. K. Smither, "Variable Focus Crystal Diffraction Lens," *Rev. Sci. Instrum.* **60**:2044–2047 (1990).
57. M. Popovici and W. B. Yelon, "On the Optical Design of Heat-Loaded Double Crystal Monochromators for Synchrotron Radiation," *J. Appl. Cryst.* **25**:471–476 (1992).
58. G. J. Price, A. N. Brunton, M. W. Beijersbergen, G. W. Fraser, M. Bavdaz, J. -P. Boutot, R. Fairbend, S. -O. Flyck, A. Peacock, and E. Tomaselli, "X-Ray Focusing with Wolter Microchannel Plate Optics," *Nucl. Instrum. Meth. A* **490**:276–289 (2002).
59. J. R. P. Angel, "Lobster Eyes as X-Ray Telescopes," *Astron. J.* **233**:364 (1979).
60. J. F. Pearson, N. P. Bannister, and G. W. Fraser, "Lobster-ISS: All-Sky X-Ray Imaging from the International Space Station," *Astron. Nachr./AN* **324**(1–2):168 (2003).
61. T. Roberts, "Future X-Ray Astronomy Missions," UK-XRA 2005, Leicester University, 2005.

62. J. Mutz, O. Bonnet, R. Fairbend, E. Schyns, and J. Seguy, "Micro-Pore Optics: From Planetary X-Rays to Industrial Market," *Proc. SPIE* **6479**:64790F.1–8 (2007).
63. P. D. Prewett and A. G. Michette, "MOXI: A Novel Microfabricated Zoom Lens for X-Ray Imaging," *Proc. SPIE* **4145**:180–187 (2001).
64. M. Y. Al Aioubi, P. D. Prewett, S. E. Huq, V. Djakov, and A. G. Michette, "A Novel MOEMS Based Adaptive Optics for X-Ray Focusing," *Microelec. Eng.* **83**:1321–1325 (2006).
65. A. Michette, T. Button, C. Dunare, C. Feldman, M. Forkard, D. Hart, C. Mcfaul, et al., "Active Micro-Structured Arrays for X-Ray Optics," *Proc. SPIE* **6705**:670502.1–11 (2007).

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

THE SCHWARZSCHILD OBJECTIVE

Franco Cerrina

*Department of Electrical and Computer Engineering
University of Wisconsin
Madison, Wisconsin*

51.1 INTRODUCTION

The Schwarzschild objective is based on the use of two almost concentric spherical surfaces: a small convex mirror and a larger concave facing each other as shown in Fig. 1.

While this design is widely attributed to Schwarzschild, previous descriptions of the optical system had already been analyzed and published by Paul and Chretien. The design is simple and elegant, and well suited for optical systems with small field of view and high resolution. The Schwarzschild objective is an evolution of the Cassegrain telescope, where the primary and secondary are both nonspherical elements, providing good aberration correction and large field of view. Aspherical optics are, however, difficult and expensive to manufacture, and simpler designs are desirable. The Schwarzschild objective replaces the aspheres with spherical elements, as shown in Fig. 1.

The Schwarzschild objective has found its primary use in microscopy and astronomy at wavelengths where glass lenses are not suitable or where a truly achromatic optical is needed, like in microspectroscopy systems. The applications are thus mainly in the Infrared (FTIR microscopes) and the UV; recently, the Schwarzschild objective has been extended to the Extreme UV region (≈ 13 nm) for microscopy^{1,2-5} and for the development of advanced lithography.⁶ Some of the first X-ray images of the sky were also acquired using Schwarzschild objectives.²

In the Schwarzschild objective each spherical surface forms an aberrated image, and it is a simple exercise in third-order expansion to show that the aberrations can be made to compensate each other. This is because the surfaces have curvatures of different sign, and hence each aberration term is of the opposite sign. All the third-order Seidel aberrations are exactly zero on axis, and only some fifth-order coma is present. However, the focal plane of the Schwarzschild objective is not flat; because of the spherical symmetry, the focal plane forms a spherical surface. Thus, the Schwarzschild objective has excellent on-axis imaging properties, and these extend to a reasonably large field. It is ideally suited for scanning systems, where the field of view is very narrow.

The excellent imaging properties of the Schwarzschild objective are demonstrated by an analysis of the objective's modulation transfer function (MTF). In Fig. 2 one can see the excellent imaging property of the Schwarzschild objective, essentially diffraction limited, as well as the effect of central obscuration (increase of side wings). As shown in Fig. 3 the curve extends to the diffraction limit

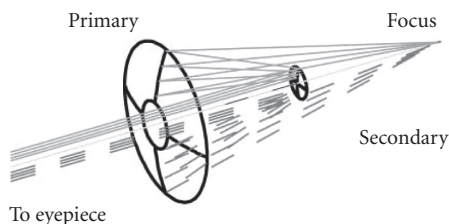


FIGURE 1 Layout of a Schwarzschild objective used as a microscope objective.

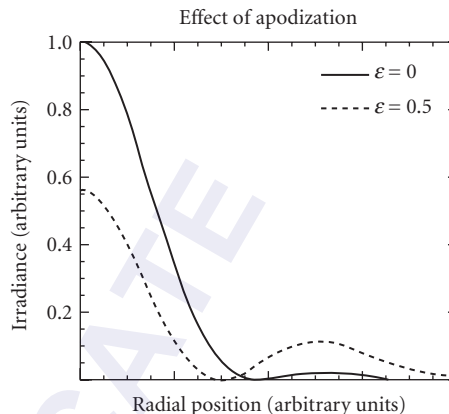


FIGURE 2 Point image of a Schwarzschild objective used as a microscope objective at a wavelength of 13 nm. ϵ refers to the amount of obstruction

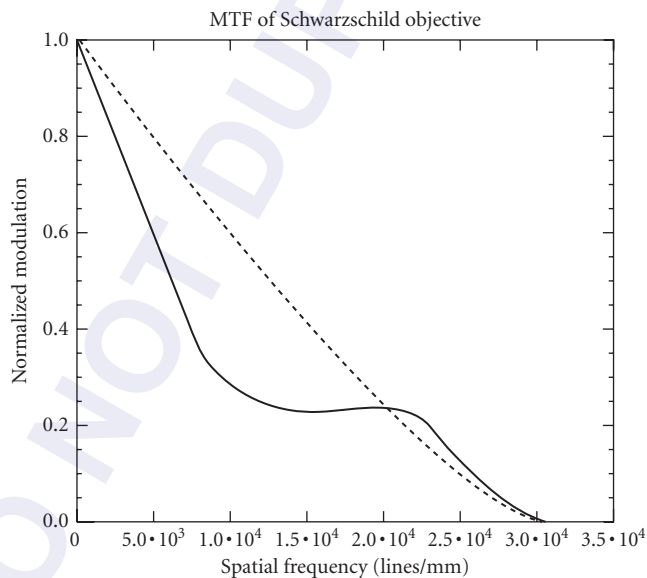


FIGURE 3 Modulation transfer function (MTF) of the previous Schwarzschild objective used as a microscope objective. Notice that in this figure the dotted line corresponds to $\epsilon = 0$ (no obstruction).

and remains very close to the ideal (diffraction limited) MTF curve even at wavelengths as short as 13 nm. However, in the midspatial frequency region, the MTF deviates considerably from the ideal. This behavior is typical of optical systems with central obstruction, and in general does not affect the ultimate resolution of the objective.

51.2 APPLICATIONS TO X-RAY DOMAIN

The high resolution of the microscope is achieved through the use of near normal-incidence optics; off-axis systems of comparable speed have unacceptable aberrations. This poses a big problem in the soft x-ray (EUV) region, where normal-incidence reflectivity is essentially negligible. The development of interference multilayer coatings for the soft x-ray region [now often called *Extreme UV* (EUV)] has in part solved this problem and made possible sophisticated optics.² (Also, see Chap. 41.) The combination of Schwarzschild objective and multilayer coatings, in its simplicity, is a very appealing design for the EUV region. Because the difference in optical path is so small, it is possible to achieve diffraction-limited performance even at wavelengths of 13 nm with good overall transmission. At wavelengths longer than about $2d = 40 \text{ \AA}$ (300 eV) it is possible to use interference filters to increase the reflectivity of surfaces in the soft x rays. These filters, often called “multilayers,” are, in effect, synthetic Bragg crystals formed by alternating layers with large optical contrast.² From another point of view, they are a $\lambda/4$ stack in the very short x-ray region. The possibility of using near normal optical surfaces clearly simplifies the design and allows greater freedom to the optical designer. Furthermore, the relatively high reflectivity of the multilayers (up to 60 percent) allows the design of multiple surface optics.

Several microscopes have been built at synchrotron facilities, and are operated successfully.³ Nonsynchrotron sources do not have enough brightness to deliver the flux required for practical experiments. The best resolution achieved with a Schwarzschild objective is of the order of 90 nm at the MAXIMUM microscope⁷ at the advanced light source. For a numerical aperture of 0.2 at the sample, the diffraction-limited resolution is approximately given by the Rayleigh criterion; at a wavelength of 13 nm, we have $\delta = \lambda/2 \text{ NA} = 32 \text{ nm}$. However, one is limited by the flux (i.e., by the finite brightness of the source) and by mounting and surface errors.³ Diffraction limit operation has not yet been demonstrated.

High resolution is particularly easy to achieve if only a small imaging field is required as in the case of a scanning microscope. The requirement of UHV conditions for the surface physics experiments forced the design of in situ alignment systems, by including piezodriven actuators in the mirror holders. During the alignment procedure, the use of an at-wavelength knife edge test made it possible to reach the ultimate resolution of approximately 900 \AA .⁴ Excellent reflectivities were achieved by using Mo-Si multilayers for the region below 100 eV, and Ru-B₄C for the region around 135 eV.

An example of a Schwarzschild objective is shown in Fig. 4. The two mirrors are coated with a Mo-Si multilayer for 92 eV operation. Notice the groove on the large concave mirror for stress-free clamping. These mirrors are then mounted in the casing shown in Fig. 5, which includes piezoelectric adjustments suitable for ultrahigh vacuum operation. Finally, Fig. 6 shows some images acquired using this objective in a scanning x-ray microscope.^{5,7}

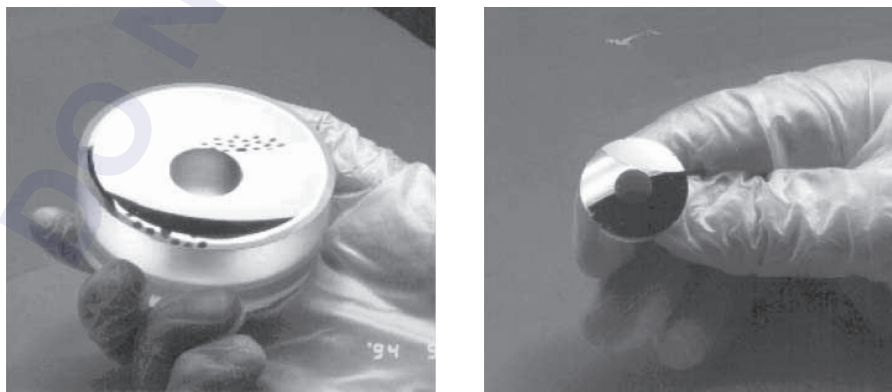


FIGURE 4 Mirrors forming a Schwarzschild objective during assembly.

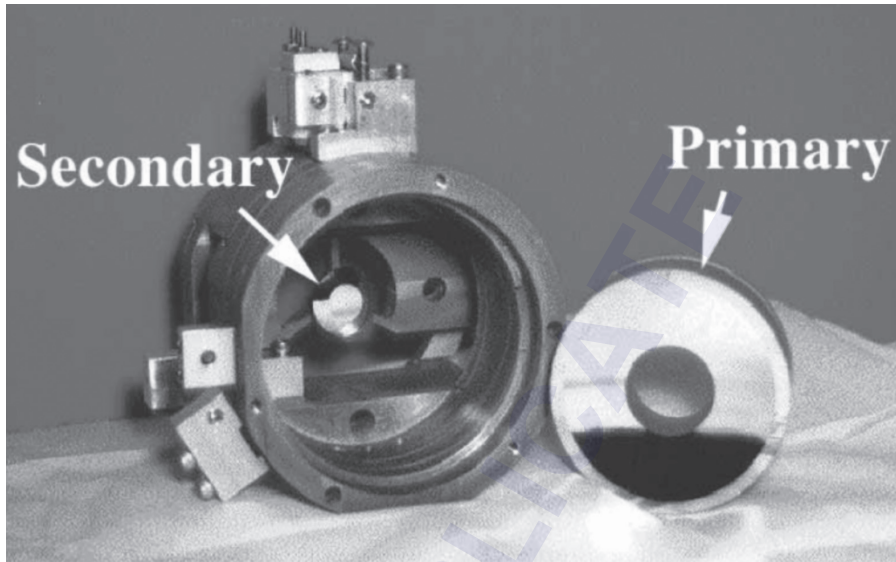


FIGURE 5 Schwarzschild objective in its mount with in situ adjustments.

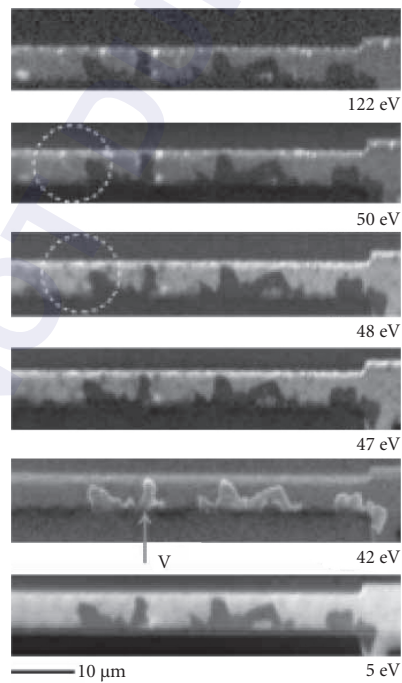


FIGURE 6 Images acquired from an Al-Cu interconnect wire in an integrated circuit after failure using the Schwarzschild objective as part of a scanning photoemission microscope.⁵ The images correspond to different electron kinetic energy and show chemical contrast.

51.3 REFERENCES

1. I. Lovas, W. Santy, E. Spiller, R. Tibbetts, and J. Wilczynski, *SPIE High Resolution X-Ray Optics* **316**:90–97 (1981).
2. D. Attwood, *Phys. Today* **45**:24 (1992).
3. F. Cerrina, *J. Electr. Spectr. and Rel. Phenom.* **76**:9–19 (1995).
4. A. K. Ray-Chaudhuri, W. Ng, S. Liang, and F. Cerrina, *Nucl. Instr. and Methods in Physics* **A347**:364–371 (1992).
5. H. H. Solak, G. F. Lorusso, S. Singh-Gasson, and F. Cerrina, *Appl. Phys. Lett.* **74**(1):22 (Jan. 1999).
6. C. W. Gwyn, R. Stulen, D. Sweeney, and D. Attwood, *Journ. Vac. Sci. and Techn.* **16**:3142–3149 (1998).
7. W. Ng, A. K. Ray-Chaudhuri, S. Liang, S. Singh, H. Solak, F. Cerrina, G. Margaritondo, et. al. *Synchr. Rad. News* **7**:25–29 (Mar./Apr. 1994).

This page intentionally left blank.

DO NOT DUPLICATE

SINGLE CAPILLARIES

Donald H. Bilderback and Sterling W. Cornaby

*Cornell High Energy Synchrotron Source
School of Applied and Engineering Physics
Cornell University
Ithaca, New York*

52.1 BACKGROUND

Monocapillary x-ray optics can be used to increase the x-ray flux per square micrometer onto a small sample while also controlling the divergence of the x-ray beam. These optics efficiently collect and transport x rays of all energies up to a cutoff energy that is dependent on the capillary material and shape. The past decade has seen the rapid development of elliptically figured monocapillary optics that are designed to condense an x-ray beam and produce a highly demagnified image of the x-ray source (tube or synchrotron) at the sample position. Monocapillary optics are being used in a wide variety of applications such as x-ray diffraction, x-ray fluorescence, small angle x-ray scattering, confocal x-ray microscopy, and so on.

52.2 DESIGN PARAMETERS

Monocapillary optics relies on total external reflection of the x rays from the internal surface of the glass tube to transport x rays. To keep the x rays from being absorbed in the wall of the capillary, the angle of incidence must be kept below the critical angle, which is typically less than 4 milliradians (0.23°). Glass materials that have been used to fabricate capillary optics are borosilicate (Pyrex), lead-based, and silica glasses. The composition of these typical glasses is shown in Table 1.

The critical angle, θ_c , is equal to $(2\delta)^{1/2}$ where δ is the refractive index decrement of the material at the energy of the x-ray photon. For borosilicate glass (Corning 7740) with a density of 2.23 g/cm^3 , θ_c is approximately $(3.8 \times 10^{-2})/E$ radians where E is the x-ray energy in keV.

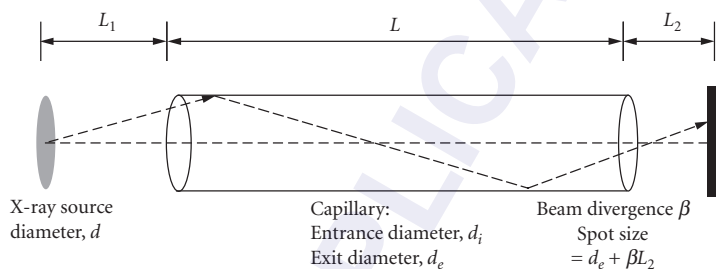
The simplest form of a monocapillary optic is a hollow straight glass tube, as shown in Fig. 1.

Straight capillary tubes were first used in the 1920s when Jentswch and Nahring demonstrated that x rays¹ can be guided down the length of the tube by multiple reflections of x rays at constant angle from the inner glass surface. The intensity of the beam is proportional to the solid angle subtended by the capillary entrance² or to $(d_i/L_1)^2$. The ideal intensity gain over a pinhole at the exit is proportional to $[(L + L_1)/L_1]^2$ provided that the glancing angle at entrance of the capillary is less

TABLE 1 Composition of Typical Glass Starting Materials and Their Basic Physical Properties

Glass Code	Type	Density (g/cm ³)	Softening Point (°C)	Weight Percent (%)							
				SiO ₂	Na ₂ O	K ₂ O	CaO	MgO	PbO	B ₂ O ₃	Al ₂ O ₃
0080	Soda lime	2.47	696	73.6	16	0.6	5.2	3.6	0	0	1.0
7050	Borosilicate	2.25	703	67.3	4.6	1.0	0	0.2	0	24.6	1.7
7740	Borosilicate	2.23	820	80.5	3.8	0.4	0	0	0	12.9	2.2
7900	96% Silica	2.18	1500	96.3	0.2	0.2	0	0	0	2.9	0.4
7910	99% Silica	2.18	1500	99.5	0	0	0	0	0	0	0
8870	High lead	4.28	580	35	0	7.2	0	0	58	0	0

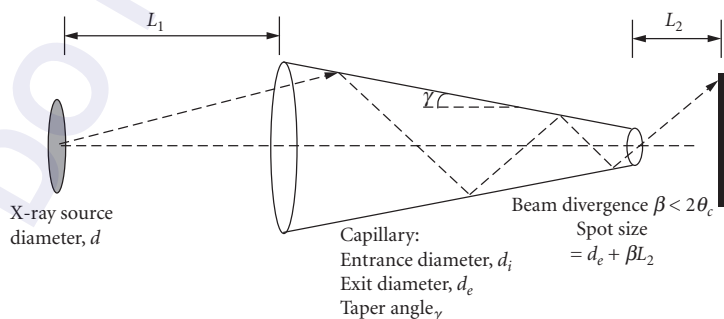
Source: Corning Glass Works, Corning, NY.

**FIGURE 1** Design parameters of a straight capillary tube geometry (diameter of capillary entrance = diameter of exit) and its divergence, β .

than or equal to θ_c . These properties have been used by several investigators to create small beams of x rays that ranged from 10 to 200 μm in diameter.³⁻¹⁰

Condensing capillary optics are figured so that the x-ray beam is compressed with each reflection, producing a gain in intensity. These optics produce a smaller, more intense beam at the expense of divergence. Stern et al.¹¹ were the first to quantitatively describe how x rays would propagate down a linearly tapered or optic. The design parameters for this optic are shown in Fig. 2 and are primarily driven by the following parameters:

1. Divergence of the exit beam, and spot size
2. Distance from the source, source dimensions, and source beam divergence

**FIGURE 2** Design parameters for a condensing capillary optic including its entrance capillary diameter, d_i , its exit diameter, d_e , and its tapering angle, γ .

The acceptance angle, at a given x-ray energy, of the capillary is $2\theta_c$. Upon reflection, the angle of incidence increases by 2γ and only photons for which the angle of incidence is less than the critical angle will emerge from the capillary. The emerging beam has a divergence which is generally less than $2\theta_c$ and the sample must be placed close to the exit end of the capillary to preserve the small spot size. The beam is smallest at the tip and the sample should be positioned no further away than 20 to 100 times the exit diameter.¹² The design of condensing capillary optics relies on the use of ray tracing computer models that allow the capillary shape, figure errors, surface roughness, materials, and source parameters to be included.^{13–15} Advanced condensing capillary designs are being investigated that include using elliptical and hybrid combinations of capillary shapes.^{16,17}

The present state-of-the-art in single channel capillaries are single-bounce capillary optics, which produce a well-defined focus. A single hollow capillary tube is shaped into a parabolic or elliptical curve, (a conic section of eccentricity less than 1) rotated about an axis of symmetry. With the well-controlled shape, the x rays need only one bounce from the inner surface to be directed to the capillary optic's focus. For synchrotron applications, the source is typically located many meters away from the optic and the incident radiation has a low divergence.

The key parameters in the design of these optics are:

1. The working distance (distance F in the diagram). This sets the demagnification that determines the spot size. The spot size is approximately $d^*(L_2/2 + F)/L_1$.
2. The divergence of the x-ray beam incident upon the sample. The full divergence of the optic is $\theta_D = (\text{tip inner diameter})/F$. For total external reflection, $\theta_D < 4\theta_c$.

The elliptically shaped optics were first developed at the University of Melbourne, Australia, and have demonstrated the ability to produce a focused spot of 40 to 50 μm with monochromatic light, and a gain factor of 700 (the increase in flux through a small $5 \times 5 \mu\text{m}$ aperture with the optic) with a divergence angle of approximately 6 milliradians.^{18,19}

Single-bounce monocapillary optics have been subsequently used for a number of years at Cornell High Energy Synchrotron Source (CHESS). They have been made in an array of sizes, with focal lengths ranging from 20 to 150 mm, divergences from 2 to 10 mrad, and have produced spot sizes between 5 and 50 μm , with gains in intensity ranging from 10 to 1000. The single-bounce optics available at CHESS generally have inner base diameter sizes ranging from 80 μm to 1.2 mm, inner tip diameter sizes ranging from 40 μm to 1 mm, and tube lengths ranging from 40 to 150 mm. The real spot size and divergence of the x rays in the image plane depends on the size of the source, the shape of the glass tube, and the slope errors introduced during the manufacturing (Fig. 3).^{20,21}

Elliptically shaped single-bounce capillary x-ray optics are focusing optics and not imaging optics. The optic's focus is an extremely distorted image of the source since the rays are effectively smeared over 720° of rotation.^{20,21} Additionally, the optic's image is smeared due to the optic's range of magnifications; the rays closest to exit tip have the strongest focusing and those furthest from the

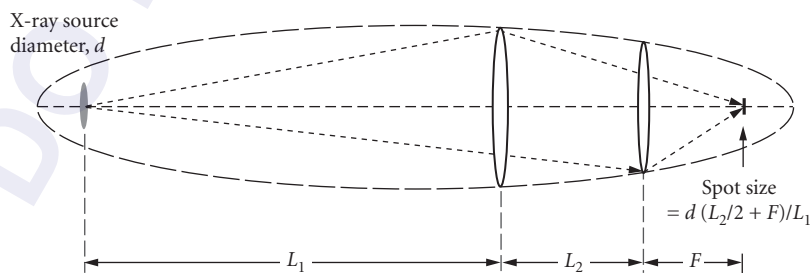


FIGURE 3 Focusing ellipse for single-bounce capillary. L_1 is the distance from the x-ray source to the capillary entrance. L_2 is the length of the capillary and F is distance from the tip of the capillary to the focus.

tip have the weakest focusing (or smallest angular deflection). Grazing incident optics, however, can produce images with a double-bounce Wolter mirror design, with one bounce from an ellipsoid surface and a second bounce from a hyperboloid surface.^{22,23}

52.3 FABRICATION

Capillary optics are typically fabricated by heating the glass tube and then pulling the glass out of the furnace in a controlled way. The shape is varied by changing the rate at which the glass is fed into and pulled out of the furnace and by the effects of surface tension. A schematic of a typical puller is shown in Fig. 4.

This puller, located at the Cornell High Energy Synchrotron Source (CHESS), controls the furnace motion, glass tension, and the rate at which the glass enters and leaves the furnace.²¹ This puller has been used to fabricate capillary optics which are 4 to 20 cm in length starting with outer tubing diameter up to 3 mm. The furnace limits the peak temperature to about 900°C, so that low temperature glasses, such as borosilicate glass, soda lime glass, etc., can be pulled with the present equipment.

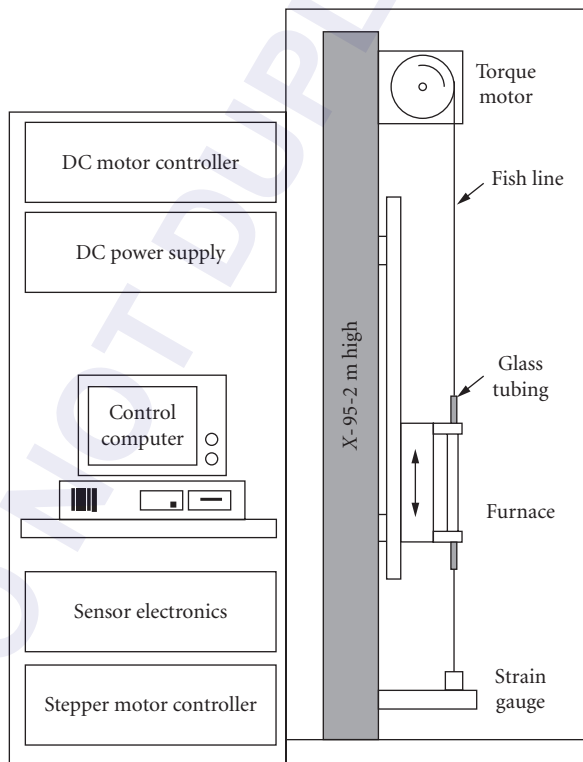


FIGURE 4 A glass tube is suspended in an electric furnace from a piece of fish line that is attached to a strain gauge at the bottom. The torque motor keeps a constant tension as the glass yields during drawing. The furnace is programmed to move based on the amount of glass yielding to make the desired elliptical, parabolic, etc., shape. (See also color insert.)

The fabrication of high efficiency capillary optics requires three factors: (1) a diameter versus length profile that approaches an ideal shape, (2) a concentric capillary bore, and (3) glass that is smooth on an atomic scale. Other glass pullers have been fabricated for capillary drawing.^{21,24,25}

A second method has been developed to fabricate capillary optics from metals. In this technique, the capillary optic is replicated from an ultrasmooth mandrel through sputtering, vacuum evaporation, or other coating techniques, followed by electroforming.²⁶ The mandrel is expendable and is formed by the precision etching of a glass or metal wire into the shape of the desired bore of the capillary optic. This technique has been successfully used to fabricate gold and copper paraboloidal imaging optics that produce a focused spot size of less than 10 μm in diameter with a collimated synchrotron radiation source. Advantages of this technique include wide latitude in the selection of materials compromising the optics, good control of the figure and straightness of the part, and high thermal conductivity.

52.4 APPLICATIONS OF SINGLE-BOUNCE CAPILLARY OPTICS

Applications that have used the single-bounce monocapillary optics are very diverse. They include high pressure powder diffraction, high resolution microdiffraction (μXRD), micro-x-ray fluorescence (μXRF), confocal x-ray fluorescence, microprotein crystallography, Laue protein crystallography, microsmall angle x-ray scattering (μSAXS), x-ray absorption fine-structure (μXAFS), and x-ray absorption near edge structure (μXANES).^{19,21,27-32} The small micrometer-sized x-ray beams have been used on a very large array of samples, including proteins, synthetic and natural fibers, polymers, foils, paintings from antiquity, solid-state devices, biological tissues, and so on. At CHESS, the single-bounce capillary optics have been used to focus beams with bandwidths ranging from 0.01 to 30 percent (they are achromatic), produced spot sizes as small as 5 μm , and total flux throughput ranging from 10^9 to 10^{13} photons/second, as seen in Fig. 5.

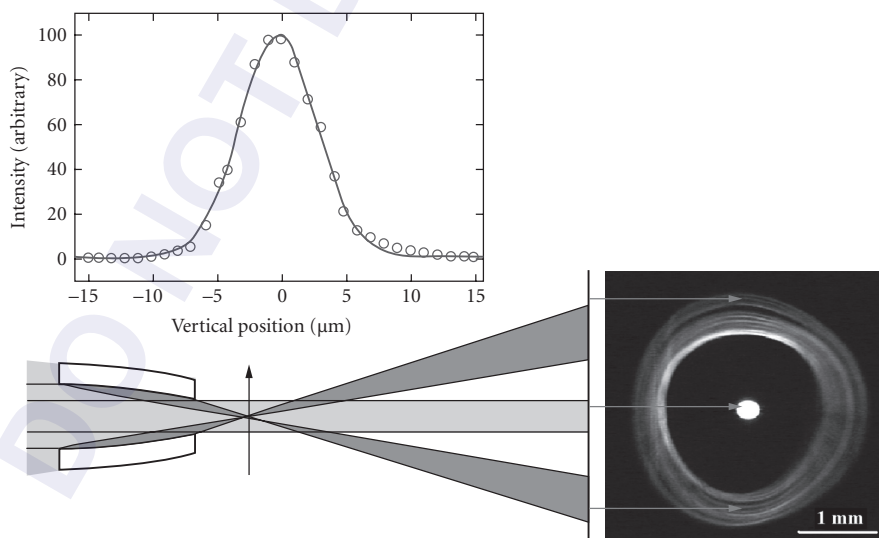


FIGURE 5 Upper panel: Profile of intensity versus 5 μm vertical pinhole position at focus of a 9 milliradians capillary producing a spot size of 5 μm FWHM at a distance of 20 mm beyond the tip of the capillary.²¹ Lower panel: The far-field image shows the direct beam (center dot) passing through the capillary and the once-reflected beam forming the outer ring of intensity. The structure in the ring is due to slope-error imperfections arising from the pulling process. (See also color insert.)

The highest flux density achieved with these optics is 3×10^{12} photons/second in a $10 \mu\text{m}$ spot size ($\sim 4 \times 10^{10}$ photons/second/ μm^2) at the advance photon source (APS).³² The uses of the small x-ray beams produced with these optics, in conjunction with synchrotron sources, are still in infancy. Also, they are just beginning to be used in conjunction with microfocusing x-ray tubes and synchrotron radiation sources for functions such as condensers in x-ray microscopes.³³

52.5 APPLICATIONS OF CONDENSING CAPILLARY OPTICS

Condensing capillary optics are used to produce small spots of x rays whose diameter ranges from 1 to about $25 \mu\text{m}$ with x-ray tubes and from 0.1 to $10 \mu\text{m}$ with synchrotron sources. A major application area is in determining the composition and structure of heterogeneous samples through x-ray fluorescence and diffraction. In these techniques, capillaries are used to increase the flux density of the x-ray beam onto the sample. By scanning a sample through the beam and recording the resulting x-ray fluorescence spectra, one can obtain the spatial distribution of major, minor, and trace elements contained within the sample. The diffracted radiation can be used to identify the crystallographic structure of particular regions of the sample. Diffraction data has been obtained with x-ray beams from 0.8 to $0.05 \mu\text{m}$ in diameter.^{17,34–36} Yamamoto developed an integrated spectrometer that he used for the simultaneous measurement of local strain and trace metal contaminants in integrated circuits.^{17,35,37–39} York has developed and implemented multiple capillary-based microdiffractometers for defect analysis in magnetic disk read heads.³⁴ Capillary-based microbeams have been used to study the impact of crystallographic structure on the sensitivity of IC interconnects to the effects of electromigration.^{35,37–39} Other laboratory based microfluorescence instruments have been developed to produce x-ray microbeams that range from $200 \mu\text{m}$ to approximately $5 \mu\text{m}$ in diameter.^{7–10,16,38–41}

52.6 CONCLUSIONS

Single-bounce monocapillary optics are still in the developmental stage. They are highly useable for making micrometer-sized x-ray beams, but there is still further room for improvement. The most serious limitations with single-bounce optics arises from slope errors and lack of centerline straightness during the time of manufacture. Presently, slope errors are at the $50 \mu\text{rad}$ level and centerline errors are generally $1 \mu\text{m}$ or less for the Cornell puller. Slope errors on the order $10 \mu\text{rad}$ or less will be required to take full advantage of third generation synchrotron sources. One advantage these grazing incidence devices have over many other competing microbeam optics (e.g., zone plates, refractive lenses, etc.) is that they are achromatic; they work over a wide energy range without changing the position or the spot size and they are very easy to add or align on a beam line when x-ray microbeam capability is required.

52.7 ACKNOWLEDGMENTS

CHESS is supported by the National Science Foundation and NIH-NIGMS via NSF grant DMR-0225180.

52.8 REFERENCES

1. F. Jentzch and E. Nahrung, "Die Fortleitung von Licht—und Rontgenstrahlen durch Rohren," *Zeitschr. F. Techn. Phys.* **12**:185 (1931).
2. W. T. Vetterling and R. V. Pound, "Measurements on an X-Ray Light Pipe at 5.9 and 14.4 keV," *J. Opt. Soc. Am.* **66**(10):1048 (1976).
3. D. Mosher and S. Staphanakis, "X-Ray Light Pipes," *Appl. Phys. Lett.* **29**:105 (1976).

4. P. S. Chung and R. H. Pantell, "Properties of X-Ray Guide Tubes," *Electron. Lett.* **13**:527 (1977).
5. H. Nakazawa, "X-Ray Guide Tube for Diffraction Experiments," *J. Appl. Crystallog.* **16**:239 (1983).
6. A. Rindby, "Application of Fiber Technique in the X-Ray Region," *Nucl. Instrum. Methods A* **249**:536 (1986).
7. P. Engström, S. Larsson, A. Rindby, and B. Stocklassa, "A 200 μm X-Ray Microbeam Spectrometer," *Nucl. Inst. Meth. Phys. Res.* **B36**:222 (1989).
8. D. A. Carpenter, "An Improved Laboratory X-Ray Source for Microfluorescence Analysis," *X-Ray Spectrometry* **18**:253 (1989).
9. H. Fukumoto, Y. Kobayashi, M. Kurahashi, and A. Kawase, "Development of a High Spatial Resolution X-Ray Fluorescence Spectrometer and its Application to Quantitative Analysis of Biological Systems," *Advances in X-Ray Analysis* **35**:1285 (1992).
10. S. Shimomura and H. Nakazawa, "Scanning X-Ray Analytical Microscope Using X-Ray Guide Tube," *Advances in X-Ray Analysis* **35**:1289 (1992).
11. E. A. Stern, Z. Kalman, A. Lewis, and K. Lieberman, "Simple Method for Focussing X Rays Using Tapered Capillaries," *Appl. Opt.* **27**(24):5135 (1988).
12. D. H. Bilderback, D. J. Theil, R. Pahl, and K. E. Brister, "X-Ray Applications with Glass-Capillary Optics," *J. Synchrotron Rad.* **1**:37 (1994).
13. L. Vincze, K. Janssens, A. Rindby, and F. Adams, "Detailed Ray-Tracing Code for Capillary Optics," *X-Ray Spectrometry* **24**:27 (1995).
14. D. X. Balaic and K. A. Nugent, "The X-Ray Optics of Tapered Capillaries," *Appl. Opt.* **34**:7263 (1995).
15. D. J. Thiel, "Ray-Tracing Analysis of Capillary Concentrators for Macromolecular Crystallography," *J. Synchrotron Rad.* **5**:820 (1998).
16. A. Attaelmanan, P. Voglis, A. Rindby, S. Larsson, and P. Engström, "Improved Capillary Optics Applied to Microbeam X-Ray Fluorescence: Resolution and Sensitivity," *Rev. Sci. Instrum.* **66**(1):24 (1995).
17. N. Yamamoto, "A Micro-Fluorescent/Diffracted X-Ray Spectrometer with a Micro-X-Ray Beam Formed by a Fine Glass Capillary," *Rev. Sci. Instrum.* **69**(9):3051 (1996).
18. D. X. Balaic, K. A. Nugent, Z. Barnea, R. Garrett, and S. W. Wilkins, "Focussing of X-Rays by Total Reflection from a Paraboloidally-Tapered Glass Capillary," *J. Synchrotron Rad.* **2**:296 (1995).
19. D. X. Balaic, Z. Barnea, K. A. Nugent, R. Garrett, J. N. Varghese, and S. W. Wilkins, "Protein Crystal Diffraction Patterns Using Capillary-Focused Synchrotron X-Ray Beam," *J. Synchrotron Rad.* **3**:289 (1996).
20. R. Huang and D. H. Bilderback, "Single-Bounce Monocapillaries for Focusing Synchrotron Radiation: Modeling, Measurements and Theoretical Limits," *J. Synchrotron Rad.* **13**:74–84 (2006).
21. S. Cornaby, "The Handbook of X-Ray Single-Bounce Monocapillary Optics, Including Design of the Optics and Synchrotron Applications," Dissertation, Cornell University, May 2008.
22. H. Wolter, "Spiegelsysteme Streifenden Einfalls Als Abbildende Optiken für Röntgenstrahlen," *Annalen der Physik* **445**(1–2):94–114 (1952).
23. H. Takano, S. Aoki, M. Kumegawa, N. Watanabe, T. Ohhigashi, T. Aota, K. Yamamoto, et al., "X-Ray Scattering Microscope with a Wolter Mirror," *Rev. Sci. Instrum.* **73**(7):2629–2633 (2002).
24. A. Rindby, "Production of Capillaries," *Proc. First International Developers Workshop on Glass Capillary Optics for X-Ray Microbeam Applications*, Don Bilderback and Per Engstrom, eds., Cornell University, 157 (October 18–19, 1996).
25. M. Haller, J. Heck, S. Kneip, A. Knöchel, F. Lechtenberg, and M. Radtke, "Stretching Tapered Capillaries from Bubbles," *Proc. First International Developers Workshop on Glass Capillary Optics for X-Ray Microbeam Applications*, Don Bilderback and Per Engstrom, eds., Cornell University, 165 (October 18–20, 1996).
26. G. Hirsch, "Tapered Capillary Optics," US Patent—5,772,903 (1998).
27. A. A. Sirenko, A. Kazimirov, S. Cornaby, D. H. Bilderback, B. Neubert, P. Bruckner, F. Scholz, V. Shneidman, and A. Ougazzaden, "Microbeam High Angular Resolution X-Ray Diffraction in InGaN/GaN Selective-Area-Grown Ridge Structures," *Appl. Phys. Lett.* **89**:181926-1 to 181926-3 (2006).
28. K. Limburg, R. Huang, and D. Bilderback, "Fish Otolith Trace Element Maps: New Approaches with Synchrotron Microbeam X-Ray Fluorescence," *X-Ray Spectrometry* **36**:336–342 (2007).
29. A. R. Woll, J. Mass, C. Bisulca, R. Huang, D. H. Bilderback, S. Gruner, and N. Gao, "Development of Confocal X-Ray Fluorescence (XRF) Microscopy at the Cornell High Energy Synchrotron Source," *Appl. Phys. a-Mater* **83**(2):235–238 (2006).

30. J. S. Lamb, S. Cornaby, K. Andresen, L. Kwok, H. Y. Park, X. Y. Qiu, D. M. Smilgies, D. H. Bilderback, and L. Pollack, "Focusing Capillary Optics for use in Solution Small-angle X-Ray Scattering," *J. Appl. Crystallog.* **40**:193–195 (2007).
31. S. Cornaby, T. Szebenyi, R. Huang, and D. H. Bilderback, "Design of Single-Bounce Monocapillary X-Ray Optics," *Advances in X-Ray Analysis* **50**:194–200 (2006).
32. R. A Barrea, R. Huang, S. Cornaby, D. Bilderback, and T. Irving, "High Flux Hard X-Ray Microprobe Using a Double Focused Undulator Beam and a Single-Bounced Monocapillary," *J. Synchrotron Rad.* **16**:76–82 (2009).
33. X. Zeng, F. Duerwer, M. Feser, C. Huang, A. Lyon, A. Tkachuk, and W. Yun, "Ellipsoidal and Parabolic Glass Capillaries as Condensers for X-Ray Microscopes," *App. Opt.* **47**(13):2376–2381 (2008).
34. B. R. York, "Chemical Analysis Using X-Ray Microbeam Diffraction (XRMD) with Tapered Capillary Optics," *Microbeam Analysis*, Proc. 28th Annual MAS Meeting, 11 (1994).
35. N. Yamamoto, Y. Homma, S. Sakata, and Y. Hosokawa, "X-Ray Spectrometer with a Submicron X-Ray Beam for ULSI Microanalysis," *Mat. Res. Soc. Symp. Proc.* **338**:209 (1994).
36. D. H. Bilderback, S. A. Hoffman, and D. J. Theil, "Nanometer Spatial Resolution Achieved in Hard X-Ray Imaging and Laue Diffraction Experiments," *Science* **263**:201–203 (1994).
37. N. Yamamoto and Y. Hosokawa, "Development of an Innovative 5 $\mu\text{m}\phi$ Focussed X-Ray Beam Energy Dispersive Spectrometer and its Applications," *Jpn. J. Appl. Phys.* **27**(11):L2203 (1988).
38. N. Yamamoto and S. Sakata, "Applications of a Micro X-Ray Spectrometer and a Computer Simulator for Stress Analyses in Al Interconnects," *Jpn. J. Appl. Phys.* **28**(11):L2065 (1989).
39. N. Yamamoto and S. Sakata, "Strain Analysis in Fine Al Interconnections by X-Ray Diffraction Spectrometry Using Micro X-Ray Beam," *Jpn. J. Appl. Phys.* **34**:L664 (1995).
40. P. -C. Wang, G. S. Cargill, I. C. Noyan, and C. -K. Hu, "Electromigration-Induced Stress in Aluminum Conductor Lines measured by X-Ray Microdiffraction," *Appl. Phys. Lett.* **72**:1296 (1998).
41. D. A. Carpenter, A. Gorin, and J. T. Shor, "Analysis of Heterogeneous Materials with X-Ray Microfluorescence and Microdiffraction," *Advances in X-Ray Analysis* **38**:557 (1995).

POLYCAPILLARY X-RAY OPTICS

Carolyn MacDonald

*University at Albany
Albany, New York*

Walter Gibson*

*X-Ray Optical Systems, Inc.
East Greenbush, New York*

53.1 INTRODUCTION

Polycapillary optics are arrays of a large number of small hollow glass tubes.^{1,2} X rays are guided down these curved and tapered tubes by multiple reflections in a manner analogous to the way fiber optics guide light. Like micropore and multifoil optics (see Chaps. 48 and 49), they differ from single bore capillaries in that the focusing or collecting effects come from the overlap of the beams from hundreds of thousands of channels, rather than from the action within a single tube. A cross section of a polycapillary fiber is shown in Fig. 1. As for single bore capillaries, x rays can be transmitted down a curved hollow tube as long as the tube is small enough, and bent gently enough, to keep the angles of incidence less than the critical angle for total reflection θ_c . The critical angle for borosilicate glass is approximately

$$\theta_c \approx \frac{30 \text{ keV}}{E} \text{ mrad} \quad (1)$$

which is approximately 30 mrad (1.7°) for 1-keV-photons and 1.5 mrad (0.086°) for 20-keV photons. The angles are somewhat larger for leaded glass. As shown in Fig. 2, the angle of incidence for a ray near one edge (toward the center of curvature) increases with tube diameter. The requirement that the incident angles remain less than the critical angle necessitates the use of tiny tubes. However, mechanical limitations prohibit the manufacture of capillary fibers with outer diameters smaller than about 300 μm . For this reason, polycapillary fibers, which have tube diameters that are much smaller than the fiber diameter, are employed. Typical channel sizes are between 2 and 50 μm , although some research has been performed with channel sizes down to submicron sizes.³

While large area optics can be produced by stringing thousands of such fibers through lithographically produced metal grids to produce a multifiber lens, as shown in Fig. 3, most commercial optics are one-piece, monolithic optics as sketched in Fig. 4.

Systems involving the use of a large number of capillary channels were first suggested by M. A. Kumahov and his collaborators in 1986.⁴ The development and study of polycapillary optics and applications^{2,5-13} have been pursued since 1990. Polycapillary optics are well-suited for broadband or divergent radiation. They have been used as focusing collectors for x-ray astronomy, to produce large area

*In memoriam Distinguished Professor Emeritus Walter M. Gibson.

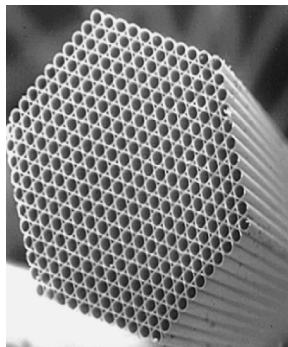


FIGURE 1 Cross-sectional scanning electron micrograph of a polycapillary fiber with 0.55-mm outer diameter and 50- μm -diameter channels.

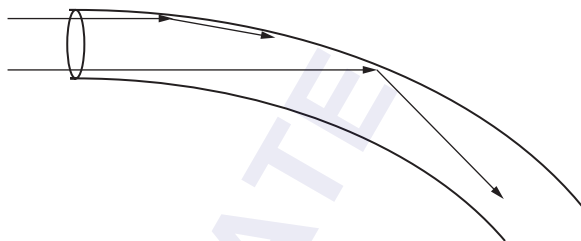


FIGURE 2 X rays traveling in a bent capillary tube. The ray entering at the bottom (closest to the center of curvature) strikes at a large angle. (Adapted from Ref. 8.)



FIGURE 3 Multifiber collimating lens constructed from over a thousand individual polycapillary fibers strung through a metal grid. The lens is 10 cm long with an output of 20×20 mm. The fibers are parallel at the output end (shown) and at the input end point to a common focal spot.



FIGURE 4 Sketch of the interior channels of a monolithic polycapillary optic. Monolithic optics can be focusing, as shown, or collimating, as in Fig. 3. (Sketch from Ref. 14.)

collimated beams for wafer analysis, and to provide small focused beams for protein crystallography with low power x-ray sources. They are also being developed for a number of medical applications, including the removal of Compton scattering with the resultant improvement in contrast and resolution in mammography,^{9,11} the production of monochromatic parallel beams for high-contrast imaging in a clinical setting,¹⁵ and the detection and localization of radioactive tracers in microscintigraphy.¹⁶⁻¹⁹

53.2 SIMULATIONS AND DEFECT ANALYSIS

The realization of numerous applications has been advanced by the development of simulation analyses which allow for increasingly accurate assessment of optics defects. These computer codes, like Shadow²⁰ (see Chap. 35), are generally based on Monte Carlo simulations of geometrical optics trajectories and provide essential information on performance, design, and potential applications of capillary optics.²¹ Some simulations also allow for the roughness^{22,23} and waviness^{24,25} of the capillary walls, as well as channel blockage and profile error to be taken into account. Optics performance over a range of energy from 10 to 80 keV can often be matched with one or two fitting parameters. For submicron channels, wave effects become significant.

Detailed measurements of polycapillary fibers, including transmission, absorption, and exit divergence, have been performed as a function of length, bend radius, x-ray source position, x-ray source geometry, and x-ray energy.^{9,12,26-28} Computer-automated systems for fiber measurement have repeatability within 1 percent.²⁹ Very good agreement is found between simulation and experimental results for a wide range of geometries.

Profile Error

Bending the channels increases the x-ray incidence angles, as shown in Fig. 2. Because the critical angle θ_c is inversely proportional to the x-ray photon energy, bending decreases the x-ray transmission down the channels most significantly at higher-photon energies. Lenses can be designed to have curvatures which deliberately discriminate against high-photon energies, for example, to reject higher-order harmonics from monochromators. Experimental data taken on a nominally straight fiber is shown in Fig. 5. The transmission drops off above 45 keV. The model cannot fit the entire spectra with bending alone. A bending radius smaller than 100 m would be necessary to fit the mid-range energies, but underestimates the high-energy transmission.

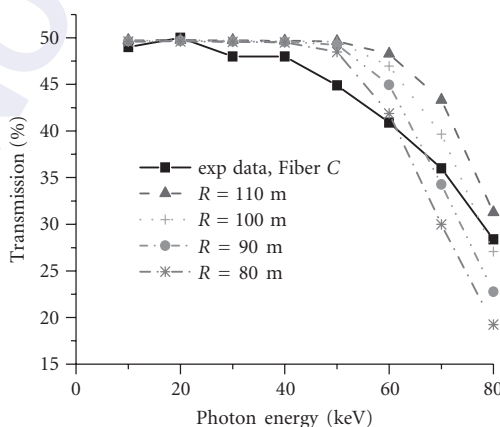


FIGURE 5 Transmission spectra of a nearly straight fiber simulated with different bending curvatures alone and compared with experimental data. (From Ref. 25.)

Waviness

Midrange spatial frequency slope errors, i.e., surface oscillations with wavelengths shorter than the capillary length and longer than the wavelength of the roughness, are often called waviness. The detailed shape of the channel walls is unknown, but waviness is modeled as a random tilt of the glass wall. The tilt angles are assumed to have a Gaussian distribution with width σ .²⁵ For high quality glass and photon energies less than 200 keV, σ is much smaller than the critical angle θ_c . Consideration is taken in the simulation of the fact that the surface tilt angle will affect the probability of x-ray impact on that surface.²⁴ The effect of waviness on fiber transmission is shown in Fig. 6. Waviness is primarily responsible for reduction of transmission at midrange energies, and additionally, for an increased reduction in transmission as the source is moved away from the fiber axis. A simulation fit including waviness and bending for a single 0.5-mm diameter fiber with 10 μm channels is shown in Fig. 7.³⁰ Most borosilicate and lead glass optics have simulation fitting parameters which give a Gaussian width for the waviness of 0.12 to 0.15 mrad. This is in agreement with the slope-error data of the Cornell group.³¹

Roughness

Roughness is treated as reducing the reflectivity of the glass (see Chap. 46). Roughness only slightly decreases the specular reflectivity at low angles, and so has almost no impact on the transmission spectra, but becomes increasingly important under circumstances in which the angle and number of reflections increase. Surface roughness must be considered to model the effects of moving the source away from the focal point.

Blockage

Another defect that is seen occasionally in borosilicate glass optics, and more prevalently in lead glass fibers,^{27,32} is a drop in transmission at low energies, as shown in Fig. 8. Reasonable agreement is obtained over the whole range of photon energies by assuming that a glass layer of fixed thickness

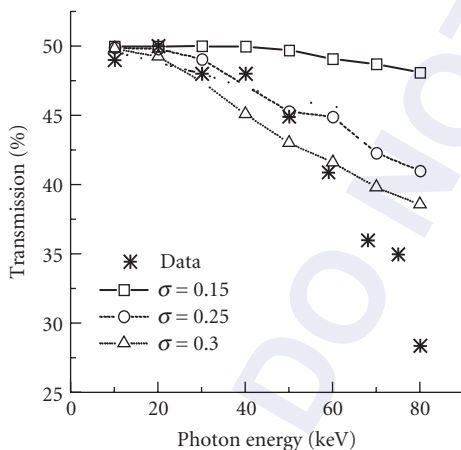


FIGURE 6 Simulations of transmission spectra for a fiber with waviness values from 0.15 to 0.3 mrad compared with the experimental data. The simulations do not include the roughness or bending. (From Ref. 24.)

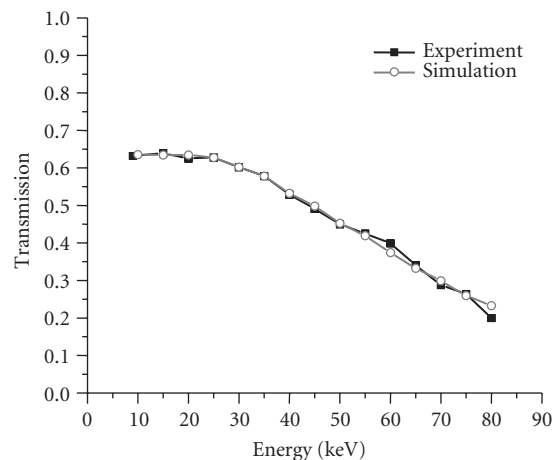


FIGURE 7 Transmission of a single fiber of the type used for a multifiber lens, compared to a simulation with fitting parameters waviness = 0.15 mrad and unintentional central axis bending radius $R = 120$ m. (From Ref. 30.)

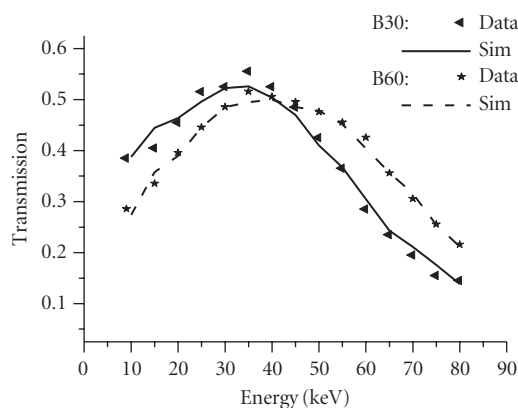


FIGURE 8 Transmission of two similar lead glass fibers, 30 and 60 mm in length. The simulation fits include 17 and 33 μm of glass layer, respectively, or 0.55 μm of blockage per mm of length. (Figure from Ref. 27.)

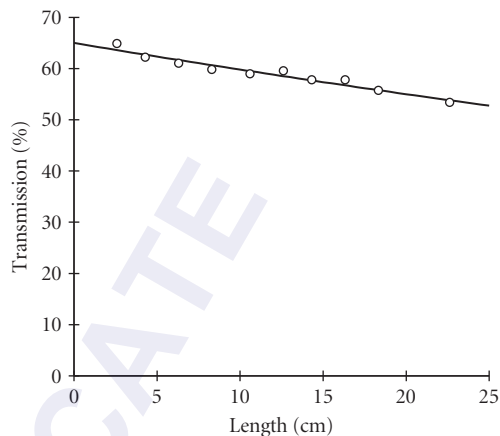


FIGURE 9 Measured transmission at 8 keV, as a function of length, for borosilicate polycapillary fibers, with 17 μm channels. Solid line is an exponential fit with decay length = 120 cm. (Adapted from Ref. 9.)

blocks the channels. The increase in required layer thickness with fiber length is consistent with a stochastic random model of glass inclusions. This random probability of glass inclusions would cause the transmission to drop exponentially with optic length, as shown in Fig. 9.

53.3 RADIATION RESISTANCE

Because the x-ray optical properties of materials depend on total electron density, the optical constants are insensitive to changes in electronic state. Color centers that form rapidly in glass during exposure to intense radiation are not indicative of a change in the x-ray transmission of a polycapillary optic. Thin fibers exposed to intense beams undergo reversible deformation. However, rigid optics, if annealed in situ at 100°C, were shown to withstand in excess of 2 MJ/cm² of white-beam bending-magnet radiation without measurable change in performance at 8 keV.³³

53.4 ALIGNMENT AND MEASUREMENT

Standard techniques have been developed for aligning and characterizing polycapillary optics. A typical setup is shown in Fig. 10.³⁴ Depending on the source geometry and desired beam placement, either the source location or the optic position is translated perpendicular to the optic axis in small steps, producing a measurement of intensity versus relative source position, as shown in Fig. 11.³⁵ The plot is symmetric and Gaussian, which indicates good alignment of the source, optic and detector. In order to determine the focal distance of the optic, source scans are performed at different distances from the source. At the focal distance, the ratio of the width of the scan curve to the optic-to-source distance, called the source scan angle, should be the smallest, as shown in Fig. 12.⁹

Transmission is the ratio of the number of photons passing along the channels to the number incident on the front face of the optic. Transmission with respect to the source-optic distance is also shown in Fig. 12. The highest transmission and the lowest source angle occurred at the focal distance.

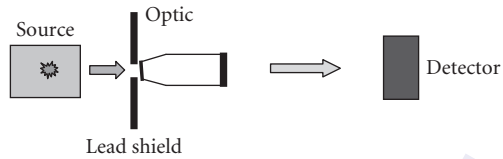


FIGURE 10 Set up for optics measurement.

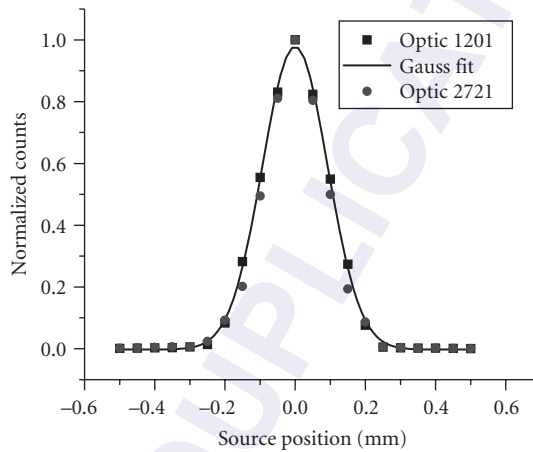


FIGURE 11 Source scan plot at 17.5 keV for two different optics with input focal lengths ranging from 48 to 56 mm. (From Ref. 34.)

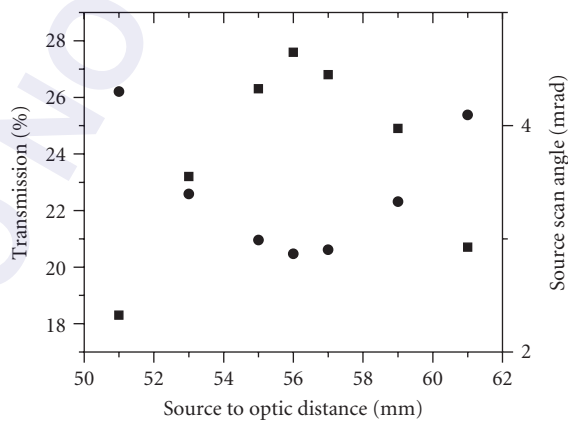


FIGURE 12 Transmission (■) and source angle (●) with respect to the source optic. The transmission is the highest at the designed focal distance of 56 mm. (From Ref. 34.)

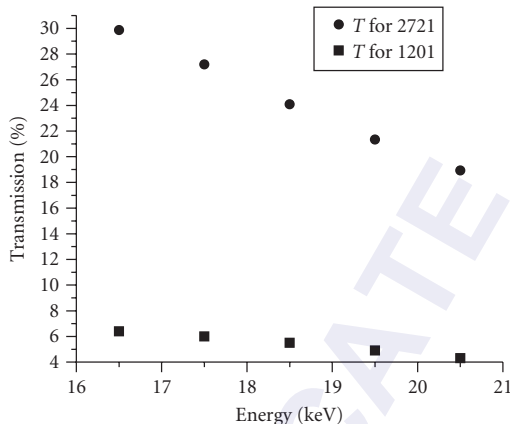


FIGURE 13 Transmission with respect to energy for optic 1201 (■) and 2721 (●).

Transmission with respect to energy can then be measured, as shown in Fig. 13.³⁴ At 17.5 keV, the transmission is 6 percent for optic 1201 and 27 percent for optic 2721 with the source at the focal distance. The transmission of optic 2721 was about 4.5 times larger than that of optic 1201 because the optics were designed with different radii of curvature R of the outer channel and different channel size c , resulting in different maximum incident angle.² The maximum incident angle for a parallel input beam is about $\theta_i = \sqrt{2c/R}$.² For optic 1201, $c = 10 \mu\text{m}$ and $R \approx 1.5 \text{ m}$, giving a maximum incident angle of 3.9 mrad, equal to the critical angle for 8 keV, so the transmission is poor for the measured energies. Instead, for optic 2721 $c = 7 \mu\text{m}$ and $R \approx 2.3 \text{ m}$, giving a maximum incident angle of 2.3 mrad, equal to the critical angle for 13-keV photons. A channel halfway from the center to the outer edge had a maximum incident angle of 1.6 mrad, equal to the critical angle for 18 keV. Thus, the overall optics performance should drop off starting in this range of 13–18 keV. Measured and calculated intensity after the optic is shown in Fig. 14 for a low power source. The measured values were

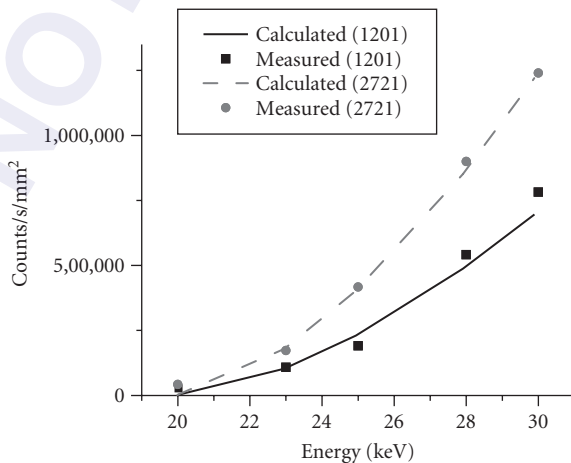


FIGURE 14 Comparison of calculated and measured intensity after optic 1201 (■) and 2721 (●) for a source current of 0.1 mA. (From Ref. 34.)

in good agreement with calculated values. The counts per unit time and unit area after optic 2721 are about 1.7 times larger than that after optic 1201. This factor is less than $T_{2721}/T_{1201} = 4.5$. That is because the two optics have the different focal distances and different areas.

53.5 COLLIMATION

As shown in Fig. 15, the output from a multifiber polycapillary collimating optic has both global divergence α and local divergence β . Even if the fibers are parallel ($\alpha = 0$), the output divergence β is not zero, but is determined by the critical angle and therefore the x-ray energy.

The exit divergence from capillary optics is measured by rotating a high-quality crystal in the beam and measuring the angular width of a Bragg peak. Since the Darwin width and mosaicity of the crystal are typically much smaller than the exit divergence from the optic, the measurement yields the divergence directly. The result at 8 keV, for a 5-mm diameter monolithic collimating optic, is a full width at half maximum of 3 mrad, which is approximately given by the critical angle for total reflection.³⁶

The divergence of the beam, and therefore the angular resolution of a diffraction measurement using polycapillary optics, does not depend on the source size, unlike the case for pinhole collimation. Thus larger, higher-power sources may be used without adversely affecting the resolution of the measurement. The maximum useful x-ray source spot size is limited by the acceptance area of the collimating lens, which is about 1 mm wide and 15 mm long for a multifiber lens, well matched for a typical rotating anode. Conversely, the divergence does not decrease for smaller sources because waviness increases the angle of reflection for x-ray photons and thus the average angle at which they exit the fiber. For example, for the case of a small source with a local divergence of 2.4 mrad, a simulation at 8 keV with no channel wall defects produces a divergence less than the critical angle, but for a simulation including a typical waviness of 0.15 mrad, the divergence grows to 3.9 mrad, which matches the measured value.²⁴

A comparison of measured and simulated transmission for a collimating lens is shown in Fig. 16. The output varies by less than 3 percent across the face of the optic, as shown in Fig. 17. The transmission of a 3-cm square multifiber collimating lens developed for astrophysical applications is shown as a function of photon energy in Fig. 18. The “gain” obtained from any optic ultimately should be a figure of merit defined by a particular application. For a diffraction measurement which requires 2D collimation, e.g., strain measurements, the intensity ratio delivered within a given output divergence angle is a reasonable figure.

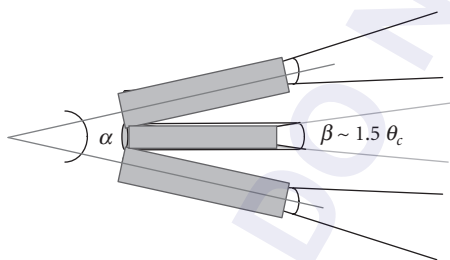


FIGURE 15 The output divergence from a multifiber collimating optic is characterized by global and local divergence. (From Ref. 2.)

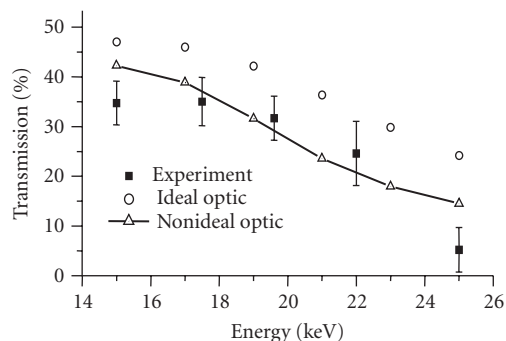


FIGURE 16 Comparison of experimental and simulated transmission as a function of energy for a multifiber lens with a 250 mm input focal length. The nonideal simulation includes a waviness of 0.15 mrad and an unintentional bend with a radius of 125 m for the center fibers. (From Ref. 15.)

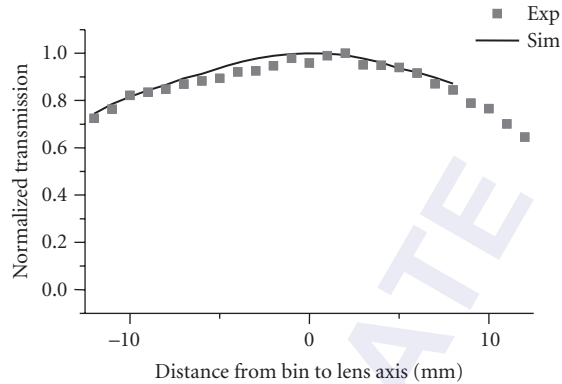


FIGURE 17 Measured and simulated local transmission of the lens in Fig. 16, at 8 keV. The uniformity scan was carried out by scanning a 5×6 mm lead aperture across the output beam. (From Ref. 30.)

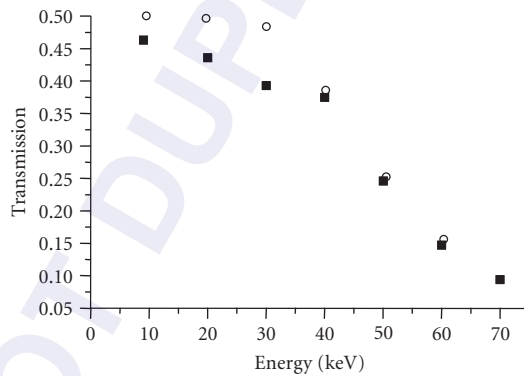


FIGURE 18 Simulated⁶ (circles) and measured⁷ (squares) transmission of a 3-cm square multifiber collimating lens with 2-m focal length.

53.6 FOCUSING

Polycapillary lenses can be used to collect broadband divergent radiation and redirect it toward a focal spot. The first Kumakhov capillary lens was built to demonstrate focusing. It had channel sizes of $300 \mu\text{m}$ and was about 1 m long.³⁷ Smaller channel sizes allow for tighter bending and shorter lenses.

Focusing the beam increases the intensity on a small sample, compared to pinhole collimation. The intensity gain depends on the spot size produced by the optic, which is determined by the capillary channel size c , output focal length f , and x-ray critical angle θ_c .

$$d_{\text{spot}} \approx \sqrt{c^2 + (1.3 \cdot f_{\text{out}} \cdot \theta_c)^2} \quad (2)$$

The factor 1.3 is an experimentally determined parameter that arises from the fact that most of the beam has a divergence less than the maximum possible divergence of $2\theta_c$ produced by reflection.

TABLE 1 Results for Monolithic Focusing Optic in a Synchrotron Beam³⁹

X-Ray Energy (keV)	Spot Size (mm)	Transmission (%)	Measured Gain 350- μm Pinhole	Calculated Gain 350- μm Pinhole	Calculated Gain 90- μm Pinhole	Calculated Gain 10- μm Pinhole
6	0.09	36	78	81	645	911
8	0.08	49	96	110	933	1359
10	0.09	39	83	87	624	842
12	0.09	39	74	87	654	903
white	0.17	42	11	89	243	266

The critical angle θ_c at 20 keV is 1.5 mrad. A lens with $c = 3.4 \mu\text{m}$ and $f_{\text{out}} = 9 \text{ mm}$ has a predicted spot size of 18 μm . An intensity distribution measurement, made by scanning a small pinhole, gave a FWHM of 21 μm .³⁸

For the case of a very small sample with diameter σ , the gain, relative to pinhole collimation, is given by

$$\text{Gain} = \frac{P_{\text{lens}}}{P_{\text{pinhole}}} = \left(\frac{d_{\text{spot}}^2}{f_{\text{in}}^2} \right) T_{\text{Lens}} \left(\frac{L_{\text{source-sample}}}{\sigma} \right)^2 \quad (3)$$

For the slightly convergent lens used to produce lysozyme diffraction data,³⁶ the computed gain for a 0.5-mm pinhole is 124. The measured intensity gain is 110. The intensity gain for a more highly convergent lens with an output focal spot of 21 μm is 2400 at a source to sample distance of 100 mm. Because of the divergence from each channel, lenses with smaller focal lengths have smaller spot sizes, as do measurements at higher photon energies.

Polycapillary optics can also be used to focus parallel beam radiation for astrophysical or synchrotron applications. A 5-mm diameter lens with a 17-mm focal length was measured using knife-edge scans in a synchrotron beam.³⁹ The measured spot size, transmission, and gain are shown in Table 1 along with calculated gains.

53.7 APPLICATIONS

Energy Filtering

The dependence of the critical angle for reflection on photon energy results in an energy-dependent transmission, which can be varied with lens parameters, as shown in Figs. 16 and 18. Thus, capillary optics can be used as a low pass filter to remove high-energy bremsstrahlung photons above a given energy threshold. With this low-pass filter, high-anode voltages can be used to increase the intensity of the characteristic lines without increasing the high-energy background. An example of the effect of a monolithic optic designed to pass 8-keV Cu K_{α} radiation is shown in Fig. 19.³⁶ The lens reduces the Cu K_{β} 9 keV peak and suppresses the high-energy bremsstrahlung. For low-resolution diffraction applications, the energy filtration provided by the optic allowed the monochromator to be replaced with a simple filter. Albertini has proposed using polycapillary optics as a low-pass-filter for energy-dispersive x-ray diffractometry and reflectometry.⁴⁰

X-Ray Fluorescence

Polycapillary optics are widely used in x-ray fluorescence (XRF). Details of x-ray fluorescence analysis are described in Chap. 29. Focusing the beam from a divergent source can result in a large intensity increase compared to pin hole collimation. This allows for flexible system development. Formica designed a system for *in situ* thin-film deposition.⁴¹ Nikitina,⁴² Buzanich,⁴³ and Vittiglio⁴⁴

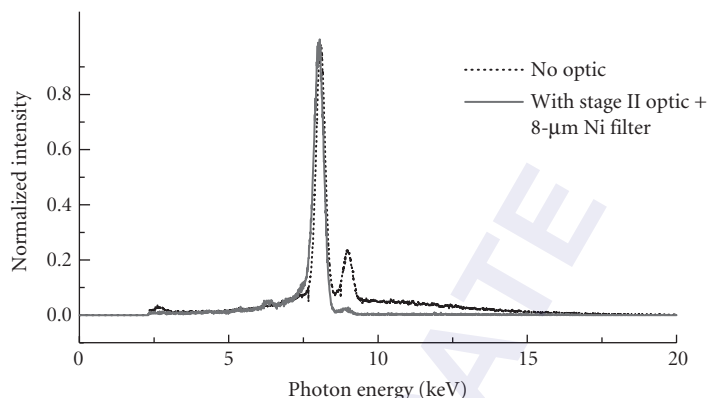


FIGURE 19 Spectrum of a copper tube source with and without a slightly focusing optic. The optic suppresses the high-energy bremsstrahlung. The nickel filter suppresses the $K\beta$ peak. (From Ref. 36.)

have demonstrated portable systems for materials and archaeometric analysis. Luo has developed a system for combinatorial materials studies.⁴⁵ Feldkamp has developed a system consistent with microtomography.⁴⁶ Kanngiesser provided detailed analysis of the beam shape and spectrum from focusing polycapillary optics to allow for accurate quantitative analysis.⁴⁷ Langer showed that polycapillary focussing makes possible the analysis of Kossel patterns from the fluorescent excitation using a tube source.⁴⁸ Instead of using the focussing optic on the excitation side, Smit uses a focusing optic to collect the fluorescence radiation in a synchrotron beam system,⁴⁹ Zitnik in a proton beam system,⁵⁰ and Alberi for particle induced x-ray emission (PIXE).⁵¹ With a conventional x-ray source, two optics can be used to both excite and collect the fluorescent radiation, as shown in Fig. 20.⁵² This provides the double benefit of enhanced signal intensity and three-dimensional spatial resolution. The three-dimensional resolution arises from the overlap of the cone of irradiation of the first lens and the cone of collection of the second. Sun uses a similar arrangement for energy-dispersive micro-x-ray diffraction.⁵³ Alternatively, Yang uses a collimating optic to increase the excitation intensity for total reflection x-ray fluorescence,⁵⁴ and Sun uses a reversed collimating optic in combination with a torroidal mirror for synchrotron radiation EXAFS.⁵⁵

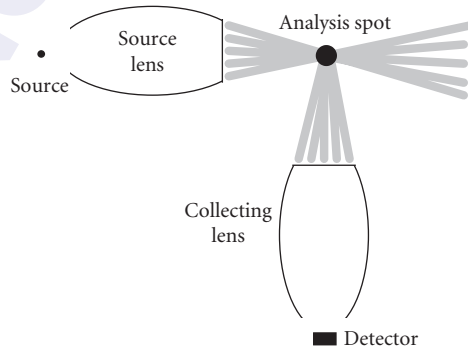


FIGURE 20 Sketch of microfluorescence experiment, showing that overlap of irradiation and collection volumes yields three-dimensional spatial resolution. (From Ref. 2.)

Single Crystal Diffraction

Collimating Significant reduction in data collection times for single crystal diffraction can be achieved with collimating polycapillary optics. The divergence from the optic, for example, 0.19° at 8 keV, is less than the ω crystal oscillations typically employed to increase the density of reflections captured in a single image in protein crystallography.⁵⁶ Using a low-power x-ray source which allows close access to the beam spot, the x-ray intensity obtained with a 20-W x-ray source was comparable to that achievable with a 3.5-kW rotating anode source with pinhole collimation. Chicken egg white lysozyme data taken with a 20-W source and a collimating lens in 20 minutes per frame produced an R -factor (variance between the measured and model structure factors) of 5.1 percent and resolution of 1.6 Å, as good as equivalent rotating anode data taken with the same or longer exposure time. The data shown in Fig. 21 was taken in 10 minutes with the collimating optic and 20-W source. Gubarev reported high-quality data extending to 1.7 Å with a 47-W system with greater flux than for a rotating anode using graded multilayer mirrors.⁵⁷

Focused For focused beam diffraction, the volume of reciprocal space that is accessed in a single measurement is greatly increased compared to parallel beam geometries. Additionally, the small irradiation spot on the sample reduces angular broadening due to sample size effects and allows the detector to be placed close enough to the sample that a small imaging detector will intercept diffracted beams at high 2θ angles. The reduction in angular broadening due to the small source size is significant because the diffraction spot is not isotropically broadened by the convergence of the focused beam.⁵⁸

Figure 22 displays a sketch of the diffraction condition for a single crystal with a monochromatic convergent beam. Diffraction conditions are satisfied for the two incident beam directions, \mathbf{k}_0 and \mathbf{k}_1 , when they make the same angle with the reciprocal lattice vector, \mathbf{G} . Thus, changing from \mathbf{k}_0 to \mathbf{k}_1 rotates the diffraction triangle of \mathbf{k}_0 , \mathbf{G} , and \mathbf{k}_f about the vector \mathbf{G} by an angle ϕ . This results in the diffracted beam \mathbf{k}_f moving to trace out a tangential line on the detector. The resultant streak, shown in Fig. 23 has length

$$x = 2z \tan\left(\frac{\phi}{2}\right) = 2 \frac{D}{\cos(2\theta)} \tan\left(\frac{\phi}{2}\right) \quad (1)$$

The maximum value of ϕ is the convergence angle. There is no broadening in the radial direction. The “straight through” beam shown as the vector D will vary within the range of incoming beam directions for different reciprocal lattice vectors.

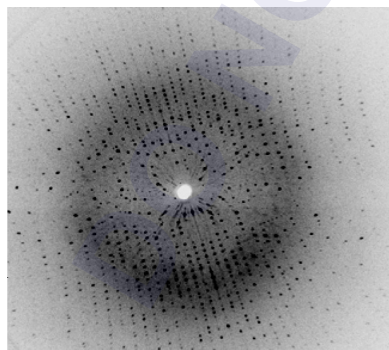


FIGURE 21 Lysozyme diffraction image taken with a 20-W microfocus source and collimating polycapillary optic in 10 minutes. (From Ref. 56.)

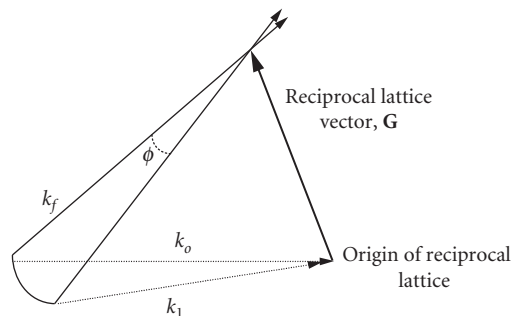


FIGURE 22 Ewald sphere description of focused beam diffraction on a single crystal.

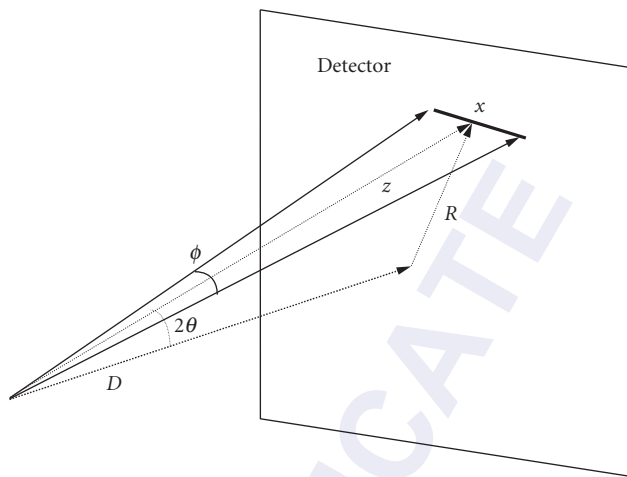


FIGURE 23 Diffraction streak due to beam focusing.

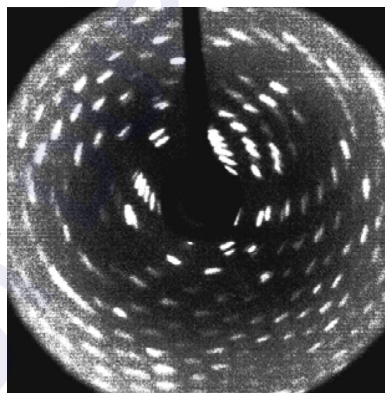


FIGURE 24 Lysozyme diffraction image taken in 5 min with a 1-mA tube source and a focusing polycapillary optic. (From Ref. 14.)

The effects of the one-dimensional streaking are shown in Fig. 24 for a single crystal egg-white lysozyme diffraction pattern taken with a 2.1° focusing angle.¹⁴ Serious overlap problems were not encountered except in low index directions, which are of less interest for structure determination. However, so long as the streaked diffraction spots are narrow compared to their separation, they can, in principle, be analyzed. For protein crystals with unit cell dimensions $< 200 \text{ \AA}$, such as lysozyme, this does not preclude structural determinations. However, for cell dimensions $> 200 \text{ \AA}$, the diffraction spots are not completely separated. Patterns with smaller convergence angles can be analyzed with conventional software and give good results.^{56,59} Patterns from crystals with smaller, less complex, unit cells will have lower-diffraction spot density and therefore less potential overlap than for protein crystals.

A direct comparison was made of data quality and collection time on a rotating anode system using a single egg-white lysozyme crystal, with and without polycapillary optics.⁵⁶ The direct beam intensity gain and the diffracted beam signal gain were both a factor of 20 for a lens with a 0.3° convergence. There was no degradation in data quality. Li reported intensity and resolution improvement with the use of a slightly focusing polycapillary optic for protein crystallography compared to double focusing mirrors.⁶⁰

Powder Diffraction

Collimating Reductions in data collection time can also be obtained for powder diffraction with collimating polycapillary optics. In addition, the parallel beam geometry provides insensitivity to sample preparation, shape, position, and transparency. Clapp reported significant reduction in sensitivity to sample preparation.⁶¹ Chen reported reduction in sensitivity to sample shifts and acceptability for whole pattern fitting.⁶² Secondly, the symmetric beam profile and enhanced flux gives much improved particle statistics and measurement statistics.

The nearly Gaussian peaks produced by the polycapillary collimating optics give more precise peak localization. For example, a pinhole with a divergence of 0.2 mrad produces peak localization errors approximately equal to the divergence, while the system with the polycapillary optic on the same sample produced peak localization that was detector limited to 0.04° .⁶³ Misture reported a significant intensity gain, reduction in peak shape and peak localization errors (also see Chap. 28).⁶⁴ The constant peak width and resolution throughout the diffraction space facilitates very high precision residual stress and texture analysis and reciprocal space mapping. The peak shape is ideally suited to phase identification and full pattern analysis of phase content using wavelet transforms.

In addition, polycapillary optics can be used as high-resolution soller slits. A schematic for an experiment comparing conventional and polycapillary soller slits is shown in Fig. 25. Replacing the soller slit with a polycapillary straight angular filter provided approximately a factor of two gain in signal-to-noise ratio for the identification of the ζ phase on a galvannaed steel sheet.⁶⁵ Leoni performed detailed analysis of the beam shape and correction algorithms for a system similar to that of Fig. 25.⁶⁶

Focusing For focusing polycapillary optics, the peak resolution for powder diffraction continues to be much smaller than the beam divergence, even for highly convergent beams, and agrees well with a simple geometrical model.^{63,67}

Scatter Rejection in Imaging

The earliest announcements of the discovery of x rays included a radiograph of a human hand.⁶⁸ One hundred years after its first use, simple radiography is still the most common medical imaging modality. The relative simplicity of the required apparatus and the speed with which the image is

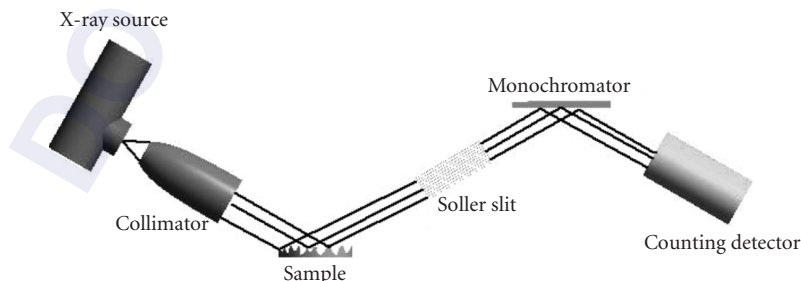


FIGURE 25 Schematic of diffraction setup showing polycapillary collimator and normal soller slit configuration. (From Ref. 65.)

obtained (with the resultant low cost) make x-ray imaging an extremely important diagnostic tool. However, Compton scattering can result in substantial image degradation. In a conventional medical imaging system, scatter is partially removed by inserting a grid with lead ribbons parallel to the incoming beam. Alternatively, scatter can be removed by inserting a polycapillary optic between the patient and the detector. Because capillary optics have an angular acceptance that is limited by the very small critical angle, scattered photons are not transported down optics channels, but are largely absorbed by the glass walls of the capillary optic. The scatter transmission of a polycapillary optic, measured by moving the source off axis to an angle much larger than the critical angle, was less than 1 percent at 20 keV,⁹ resulting in a contrast enhancement for a Lucite phantom of a factor of 1.4 compared to a conventional antiscatter grid, with a maximum enhancement of more than a factor of 2 in regions of low contrast or high scatter.¹¹ Measurements at 40 keV showed a contrast enhancement of more than a factor of 4.⁶⁹ Lead glass optics performed well at higher energies.⁷⁰

While not true imaging optics, polycapillary fibers can transmit an image in the same manner as a coherent fiber bundle. Polycapillary optics can be used to magnify and demagnify images. Magnification can improve image resolution, particularly if detector limited. Magnification is conventionally performed by increasing the air gap between the patient and the detector, but this results in a loss of resolution due to the geometrical blurring from the finite source size, as shown in Fig. 26.

However, magnification using polycapillary optics, even with a large focal spot, would not be subject to a loss in resolution, as shown in Fig. 26. Measured modulation transfer functions (MTF) for computed radiography image plate detectors using a magnifying polycapillary optic are compared to air gap magnification in Fig. 27.¹¹ The limiting MTF was increased by the polycapillary magnification factor of 1.8. The resolution was not degraded by the capillary structure, which was on a smaller scale (20- μm channel size) than the desired resolution. The optic yields an MTF increase at all spatial frequencies, which may be diagnostically more significant than the increase in limiting MTF. Magnifying capillary optics provide simultaneous contrast enhancement and resolution increase. Efforts have been made to increase the size of the optics. A transmission image of a multioptic jig is shown in Fig. 28.⁷¹

A technical difficulty with many direct x-ray detectors is that they are difficult to manufacture with large enough areas for medical imaging. With the use of x-ray optics, demagnification could be used to match the image size to a CCD, CID,⁷² or other digital detectors. A common digital imaging modality is

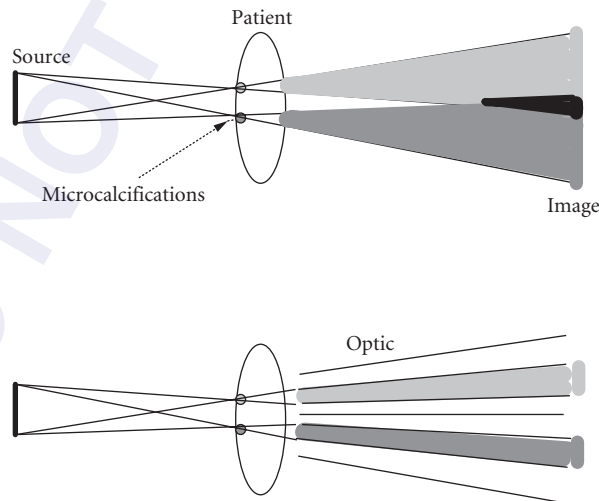


FIGURE 26 Air gap magnification (top) showing degradation of image due to finite source size. Magnification with a long tapered polycapillary optic (bottom) showing no increased blurring after the exit plane of the patient.

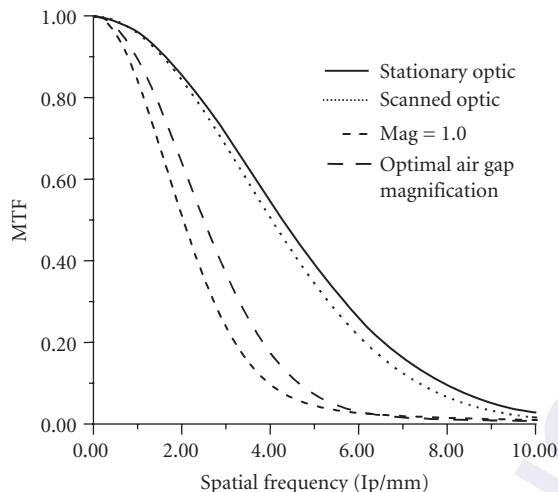


FIGURE 27 Modulation transfer function of a computed radiography plate with no optic, with airgap magnification, and with a polycapillary optic. (From Ref. 11.)

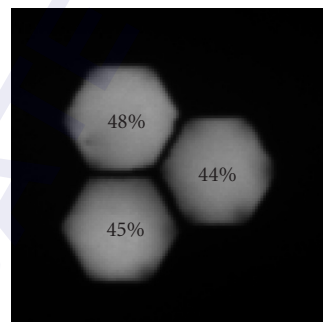


FIGURE 28 Transmission image of triad of monolithic optics. (From Ref. 71.)

the use of a fused fiber optic following a phosphor screen to demagnify the resulting visible light image for recording with CCD technology. The use of capillary x-ray optics to demagnify the x-ray image in connection with direct x-ray sensitive detectors may be more efficient and would avoid the loss of resolution in the optical conversion process.

Monochromatic Imaging

In conventional radiography, subject contrast arises from relatively low differences in absorption coefficients between different tissue types. The already low subject contrast is further reduced in a conventional system by averaging over relatively large-energy bandwidths. Synchrotron measurements using monochromatic beams have demonstrated higher contrast, but synchrotrons are not clinically available. Using monochromator crystals with a conventional source without an optic is not practical because the low intensity of the diffracted beam will not allow imaging *in vivo* before motion blur occurs. Polycapillary collimating optics can allow sufficient diffracted beam intensity to make clinical monochromatic imaging possible without a synchrotron.¹⁵ Measurements at 17.5 keV showed subject contrast enhancement of a factor of 2, in agreement with theoretical calculations. This contrast enhancement is in addition to that expected from the reduction of scattered radiation.

The output divergence of the collimating optic affects the resolution and is an important parameter, especially for low-energy and high-resolution modalities. Good angular resolution was achieved even with a large spot source. Resolution was measured by recording a knife-edge shadow with a restimable phosphor-computed-radiography image plate. For a silicon crystal, the width results primarily from the energy spread of the incident radiation. The image in Fig. 29 shows the resolved Mo $K\alpha$ energy doublet.³⁴ The divergence causes the fiber structure to blur and the field to become more uniform as the beam propagates away from the optic, as shown in Fig. 30.¹⁵

In addition to the monochromatization, the use of a uniform parallel beam would eliminate the variations in resolution and magnification of objects on the entry and exit side of the patient. A parallel, tunable, monochromatic source would facilitate the use of a number of innovative imaging techniques such as refraction contrast,⁷³ K-edge tuning,⁷⁴ dual energy,⁷⁵ and phase imaging. Refraction contrast is observed for highly monochromatic beams when a second crystal is placed after the specimen. Gradients

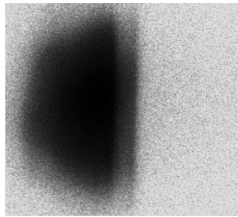


FIGURE 29 Knife-edge image. The extra shadow is due to the $K\alpha$ doublet. The image plate is 550 mm from the knife-edge.

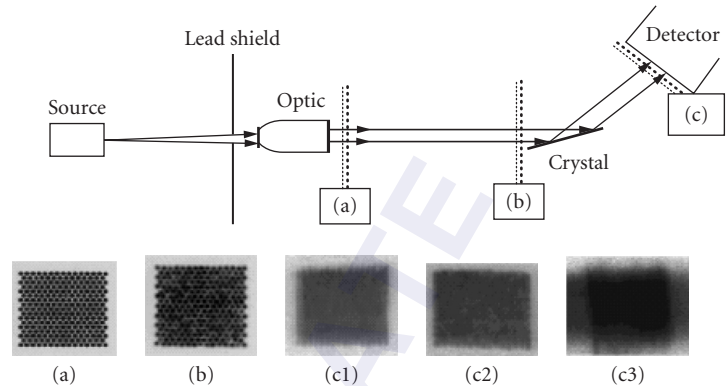


FIGURE 30 X-ray images taken at the three indicated locations along the beam. The fiber structure is clearly visible in (a) and blurs as the beam propagates. Image (c1) is for silicon, (c2) for mica, and (c3) for graphite.

in the index of refraction deflect the beam at the edges of features. The deflected beams are not diffracted by the second crystal, resulting in greatly enhanced contrast in the vicinity of feature edges. K-edge tuning is the selection of a photon energy just above the absorption edge for an element that is preferentially concentrated in, e.g., cancerous tissues. Dual-energy imaging uses the tunability of the source to make a subtraction image taken at energies just above and below the critical edge. Monochromatization with the Bragg angle at 45° produces a polarized beam, as shown in Fig. 31.⁷⁶

Scintigraphy

Unlike external beam radiography, collimators are required for nuclear imaging (see Chap. 32). In addition, detectors with energy discrimination are required because only a small fraction of the photons from the radioactive source are transmitted through the collimator. The majority of the radiation is scattered by the surrounding tissue, to create a broad source with approximately the same total

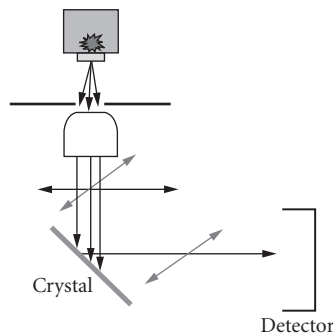


FIGURE 31 Set up for polarization by diffraction at 90° . After diffraction by the crystal, the beam is polarized in the out-of-plane direction.

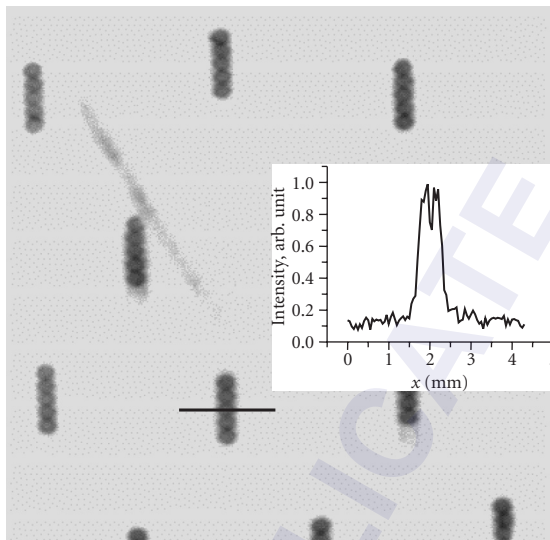


FIGURE 32 Images taken with set of brachytherapy seeds show five beads per seed and a ring-like structure. The source to optic distance was $z = 18$ mm.

intensity. Since Compton scattering causes the photons to lose energy, energy discrimination can be used to distinguish between rays which have traveled directly to the detector from those which have been scattered by the tissue. However, energy-sensitive detectors tend to have relatively poor spatial resolution compared to radiographic detectors. Testing has been performed of a different mechanism to discriminate against scatter and provide high resolution. In this study, a polycapillary optic was used as a high resolution collimator, paired with a radiographic detector with a large number of pixels.^{16–19}

Images taken with a parallel 20-mm-thick, 30-mm-wide lead glass polycapillary parallel hole collimator and a computed radiography plate demonstrated the ability to image features of individual ^{125}I ion exchange beads within brachytherapy seeds even in the presence of 30 mm of tissue-equivalent scatter material, as shown in Fig. 32. Calculations using simple geometrical models were in good agreement with measured signal-to-background ratios and count rates. The calculations and measurements showed that for highly heterogeneous radiation distributions high signal-to-background ratios could be achieved by using the high resolution to “dilute” the scatter background without energy discrimination. In addition, the high resolution collimator/radiographic detector system is significantly more compact than typical γ camera systems. Because energy sensitivity is no longer required, a wide range of radiographic detectors with high-resolution could be employed. The polycapillary collimator could be either parallel or converging.

Therapy

Conventional x-ray radiation therapy is currently performed with high energy x-ray or gamma radiation. The patient is exposed to multiple parallel beams, produced with slit collimators, with an intersection at the tumor. High-energy photons are chosen to minimize the absorbed skin dose relative to the dose at the tumor, although energies as low as 100 keV are employed in orthovoltage modalities. The choice of high energies to reduce skin dose is necessary because, in an unfocused beam, the intensity is necessarily higher near the point of entry than at the tumor site. Use of focusing optics at energies less than 100 keV might substantially increase the tumor dose relative to the skin dose.

Neutron beams can also be focussed by polycapillary optics (see Chaps. 63 and 64).^{77–79} A focused neutron beam could be used in boron neutron capture therapy (BNCT), based on the selective delivery of a boronated pharmaceutical to cancerous tissue followed by irradiation with thermal neutrons.⁸⁰ Perrung⁸¹ has described a BNCT system using capillary optics. Mayer has performed an experimental study to determine the radiation dose introduced by focusing a cold neutron beam through a tissue equivalent mass.⁸² With the use of radiochromic images, it was found that neutron scattering within the material did not significantly alter the focal image. As a result, the deposition of dose peaks quickly at the depth of the focal point. This procedure could be useful in treating near-surface regions such as ocular melanomas.

53.8 SUMMARY

Polycapillary x-ray optics are a powerful control technology for x-ray beams. Using polycapillary optics to collimate the output from a point source provides much higher intensity than pinhole collimation, particularly if two dimensional collimation is required. Focusing the beam yields even higher intensity gains.

Polycapillary optics can also be used to perform low-pass spectral filtering, which allows the use of increased source voltage. Further, the optics also remove the connection between source size and resolution, which allows the use of increased source current. Pairing a polycapillary optic with a diffracting crystal results in efficient production of monochromatic beams for imaging and analysis.

Tapered or straight polycapillary optics can also be used as angular filters, to replace soller slits in diffraction systems, to reduce scatter fraction in imaging thick objects, and to provide resolution in microscintigraphy.

53.9 ACKNOWLEDGMENTS

The authors are grateful for useful discussions and data from a large number of collaborators, including Carmen Abreu, David Aloisi, Simon Bates, Ayhan Bingolbali, David Bittel, Cari, Dan Carter, Heather Chen, Patrick Conlon, Greg Downing, Ning Gao, David Gibson, Mikhail Gubarev, Joseph Ho, Frank Hoffman, Huapeng Huang, Huimin Hu, Abrar Hussein, Chris Jezewski, Kardiawarman, John Kimball, Ira Klotzko, David Kruger, Danhong Li, Dip Mahato, Kevin Matney, David Mildner, Johanna Mitchell, Charles Mistretta, Robin Moresi, Noor Mail, Scott Owens Rohrbach, Wally Pepler, Sushil Padiyar, Igor Ponomarev, Bimal Rath, Christine Russell, Robert Schmitz, Francisca Sugiro, Suparmi, Christi Trufus-Feinberg, Johannes Ullrich, Hui Wang, Lei Wang, Russel Youngman, Brian York, Qi Fan Xiao, and Wei Zhou and for grant support from the Department of Commerce, NASA, NIH, and the Breast Cancer Research Program.

53.10 REFERENCES

1. V. A. Arkadev, M. A. Kumakhov, et al., *Sov. Physics. Usp.* **32**:3, March 1989.
2. C. A. MacDonald, "Applications and Measurements of Polycapillary X-Ray Optics," *J. X-Ray Science and Tech.* **6**:32–47, 1996.
3. A. Bjeoumikhov, S. Bjeoumikhova, H. Rieseemeier, M. Radtke, and R. Wedell, "Propagation of Synchrotron Radiation through Nanocapillary Structures," *Phys. Lett. A* **366**(4–5):283–288, July 2, 2007.
4. V. A. Arkd'ev, A. I. Kolomitsev, M. A. Kumakhov, I.Yu. Ponorave, I. A. Khodeev, Yu. P. Chertov, and I. M. Shakparonov, "Wide-Band X-Ray Optics with a Large Angular Aperture," *Sov. Phys. Usp.* **32**(3):271, March 1989.

5. I. L. Klotzko, Q. F. Xiao, D. M. Gibson, R. G. Downing, W. M. Gibson, Karnaukhov A., and C. J. Jezewski, "Investigation of Glass Polycapillary Collimator for Use in Proximity Based X-Ray Lithography," *Proc. SPIE* **2523**:175–182, 1995.
6. C. H. Russell, W. M. Gibson, M. V. Gubarev, F. A. Hofmann, M. K. Joy, C. A. MacDonald, L. Wang, Qi-Fan Xiao, and R. Youngman, "Application of Polycapillary Optics for Hard X-Ray Astronomy," in *Grazing Incidence and Multilayer X-Ray Optical Systems*, R. B. Hoover and A. B. C. Walker II, eds., *Proc. SPIE* **3113**:369–377, 1997.
7. C. H. Russell, M. Gubarev, J. Kolodziejczak, M. Joy, C. A. MacDonald, and W. M. Gibson, "Polycapillary X-Ray Optics for X-Ray Astronomy," in *Advances in X-Ray Analysis*, 43, *Proc. 48th Denver X-Ray Conference*, 1999.
8. C. C. Abreu and C. A. MacDonald, "Beam Collimation, Focusing, Filtering, and Imaging with Polycapillary X-Ray and Neutron Optics," *Phys. Med.* **XIII**(3):79–89, 1997.
9. C. C. Abreu, D. G. Kruger, C. A. MacDonald, C. A. Mistretta, W. W. Peppler, and Q. F. Xiao, "Measurements of Capillary X-Ray Optics with Potential for Use in Mammographic Imaging," *Med. Phys.* **22**(11)(Pt. 1):1793–1801, November 1995.
10. W. M. Gibson, C. A. MacDonald, and M. S. Kumakhov, "The Kumakhov Lens: A New X-Ray and Neutron Optics with Potential for Medical Applications," in *Technology Requirements for Biomedical Imaging*, S. K. Mun, ed., I.E.E.E. Press, New York, #2580, 1991.
11. D. G. Kruger, C. C. Abreu, E. G. Hendee, A. Kocharian, W. W. Peppler, C. A. Mistretta, and C. A. MacDonald, "Imaging Characteristics of X-Ray Capillary Optics in Mammography," *Med. Phys.* **23**(2):187–196, February 1996.
12. C. A. MacDonald, C. C. Abreu, S. Budkov, et al., "Quantitative Measurements of the Performance of Capillary X-Ray Optics," in *Multilayer and Grazing Incidence X-Ray/EUV Optics II*, R. B. Hoover and A. Walker, eds. *Proc. SPIE*, SPIE, Bellingham, Wash., 2011, 1993.
13. W. M. Gibson and C. A. MacDonald, "Polycapillary Kumakhov Optics: A Status Report," in *X-Ray and UV Detectors*, *Proc. SPIE*, Bellingham, Wash., 2278, 1994.
14. S. M. Owens, F. A. Hoffman, C. A. MacDonald, and W. M. Gibson, "Microdiffraction Using Collimating and Convergent Beam Polycapillary Optics," *Advances in X-Ray Analysis*, vol. 41, *Proc. 46th Annual Denver X-Ray Conference*, Steamboat Springs, Colorado, August 4–8, 1997.
15. F. R. Sugiro, D. Li, and C. A. MacDonald, "Beam Collimation with Polycapillary X-Ray Optics for High Contrast High Resolution Monochromatic Imaging," *Med. Phys.* **31**:3288, 2004.
16. C. A. MacDonald, N. Mail, W. M. Gibson, S. M. Jorgensen, and E. L. Ritman, "Micro Gamma Camera Optics with High Sensitivity and Resolution," in M. J. Flynn, ed., *Physics of Medical Imaging*, *Proc. SPIE* **5745**:1–6, 2005.
17. S. M. Jorgensen, M. S. Chmelik, D. R. Eaker, C. A. MacDonald, and E. L. A. Ritman, "Polycapillary X-Ray Optics-Based Integrated Micro-SPECT/CT Scanner," in *Developments in X-Ray Tomography IV*, U. Bonse, ed., *Proc. SPIE* **5535**:36–42, 2004.
18. W. M. Gibson, C. A. MacDonald, and N. Mail, "Potential for Radioscintigraphy with Polycapillary Optics," in A. M. Khounary, C. A. MacDonald, eds., *Advances in Laboratory-Based X-Ray Sources and Optics III*, *Proc. SPIE* **4781**:104–111, 2002.
19. N. Mail, C. MacDonald, and W. M. Gibson, "Microscintigraphy with High-Resolution Collimators and Radiographic Imaging Detectors," *Med. Phys.* **39**(2):645–655, 2009.
20. A. Liu, "The X-Ray Distribution after a Focussing Polycapillary—A Shadow Simulation," *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms* **243**(1):223–226, January, 2006.
21. D. Hampai, G. Cappuccio, G. Cibin, S. B. Dabagov, and V. Sessa, "Modeling of X-Ray Transport through Polycapillary Optics," *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, *Proc. 10th International Symposium on Radiation Physics—ISRP 10*, **580**(1):85–89, September 21, 2007.
22. D. Bittel and J. Kimball, *J. Appl. Phys.* **74**(2):877–883, 1993.
23. J. Harvey, chap. 11, In *Handbook of Optics*, Volume II, M. Bass, ed., McGraw-Hill, New York, 1996.
24. H. Wang, W. Lei, W. M. Gibson, and C. A. MacDonald, "Simulation Study of Polycapillary X-Ray Optics," in *X-Ray Optics, Instruments, and Missions*, R. B. Hoover and A. B. C. Walker II, ed., *Proc. SPIE* **3444**:643–651, 1998.
25. L. Wang, B. K. Rath, W. M. Gibson, J. C. Kimball, and C. A. MacDonald, *J. Appl. Phys.* **80**(7):3628–3638, 1996.

26. J. B. Ullrich, V. Kovantsev, and C. A. MacDonald, *J. Appl. Phys.* **74**(10):5933–5939, 1993.
27. Suparmi, Cari, W. Lei, H. Wang, W. M. Gibson, and C. A. MacDonald, *J. Appl. Phys.* **90**:5363–5368, 2001.
28. C. A. MacDonald and W. M. Gibson, “Applications and Advances in Polycapillary Optics,” *X-Ray Spectrometry* **32** (3):258–268, 2003.
29. B. Rath, R. Youngman, and C. A. MacDonald, *Rev. Sci. Instrum.* **65**:3393–3398, 1994.
30. F. R. Sugiro, S. D. Padiyar, and C. A. MacDonald, “Characterization of Pre- and Post- Patient X-Ray Polycapillary Optics for Mammographic Imaging” in C. A. MacDonald and A. M. Khounsary, eds., *Advances in Laboratory-Based X-Ray Sources and Optics*, Proc. SPIE **4144**:204, 215, 2000.
31. D. Bilderbeck and E. Fontes, *AIP Conference Proceedings* **417**:147–155, 1997.
32. Cari, C. A. MacDonald, W. M. Gibson, C. D. Alexander, M. K. Joy, C. H. Russell, and Z. W. Chen, “Characterization of a Long Focal Length Polycapillary Optic for High Energy X Rays,” in *Advances in Laboratory-Based X-Ray Sources and Optics*, C.A. MacDonald and Ali M. Khounsary, eds., Proc. SPIE **4144**:183–192, 2000.
33. B. K. Rath, W. M. Gibson, L. Wang, B. E. Homan, and C. A. MacDonald, “Measurement and Analysis of Radiation Effects in Polycapillary X-Ray Optics,” *J. Appl. Phys.* **83**(12):7424–7435, June 15, 1998.
34. D. Li, F. R. Sugiro, and C. A. MacDonald, “Source-Optic Optimization for Compact Monochromatic Imaging,” in *X-Ray Sources and Optics*, C. A. MacDonald, A. T. Macrander, T. Ishikawa, C. Morawe, J. L. Wood, eds., Proc. SPIE **5537**:105–114, 2004.
35. D. Li, N. Mail and C. A. MacDonald, “A Comparison of Doubly Curved Crystal and Polycapillary Optics for Monochromatic Beam Production from a Clinical Source,” in M. J. Flynn, ed., *Physics of Medical Imaging*, Proc. SPIE **5745**:754–763, 2005.
36. S. M. Owens, J. B. Ullrich, I. Y. Ponomarev, D. C. Carter, R. C. Sisk, J. X. Ho, and W. M. Gibson, “Polycapillary X-Ray Optics for Macromolecular Crystallography,” in R. B. Hoover and F. P. Doty, eds., *Hard X-Ray/Gamma-Ray and Neutron Optics, Sensors, and Applications*, Proc. SPIE, 2859, 1996.
37. W. M. Gibson and M. A. Kumakhov, *Yearbook of Science & Technology*. McGraw-Hill, New York, 488–490, 1993.
38. F. A. Hoffman, N. Gao, S. M. Owens, W. M. Gibson, C. A. MacDonald, and S. M. Lee, “Polycapillary Optics for In-Situ Process Diagnostics,” in *In Situ Process Diagnostics and Intelligent Materials Processing*, *Materials Research Society Proc.*, P. A. Rosenthal, W. M. Duncan, J. A. Woollam, eds., **502**:133–138, 1998.
39. F. A. Hofmann, C. A. Freinberg-Trufas, S. M. Owens, S. D. Padiyar, and C. A. MacDonald, “Focusing of Synchrotron Radiation with Polycapillary Optics,” *Beam Interactions with Materials and Atoms: Nuclear Instruments and Methods B*, **133**:145–150, 1997.
40. V. Rossi Albertini, B. Paci, A. Generosi, S. B. Dabagov, O. Mikhin, and M. A. Kumakhov, “On the Use of Polycapillary Structures to Improve Laboratory Energy-Dispersive X-Ray Diffractometry and Reflectometry,” *Spectrochimica Acta Part B: Atomic Spectroscopy* **62**(11):1203–1207, November 2007.
41. S. P. Formica and S. M. Lee, “X-Ray Fluorescence System for Thin Film Composition Analysis during Deposition,” *Thin Solid Films* **491**(1–2):71–77, November 22, 2005.
42. S. V. Nikitina, A. S. Shcherbakov, and N. S. Ibraimov, *Rev. Sci. Instrum.* **70**:2950 1999.
43. G. Buzanich, P. Wobrauschek, C. Strelti, A. Markowicz, D. Wegrzynek, E. Chinea-Cano, and S. Bamford, “A Portable Micro-X-Ray Fluorescence Spectrometer with Polycapillary Optics and Vacuum Chamber for Archaeometric and Other Applications,” *Spectrochimica Acta Part B: Atomic Spectroscopy* **62**(11):1252–1256, November 2007.
44. G. Vittiglio, S. Bichlmeier, P. Klinger, et al., “A Compact [mu]-XRF Spectrometer for (in situ) Analyses of Cultural Heritage and Forensic Materials,” *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms* **213**, 5th Topical Meeting on Industrial Radiation and Radioisotope Measurement Applications **213**:693–698, January 2004.
45. Z. Luo, B. Geng, J. Bao, C. Liu, W. Liu, C. Gao, Z. Liu, and X. Ding, “High-Throughput X-Ray Characterization System for Combinatorial Materials Studies,” *Rev. Sci. Instrum.* **76**:095105, 2005.
46. J. M. Feldkamp, C. G. Schroer, J. Patommel, B. Lengeler, T. F. Gunzler, M. Schweitzer, C. Stenzel, M. Dieckmann, and W. H. Schroeder, “Compact X-Ray Microtomography System for Element Mapping and Absorption Imaging,” *Rev. Sci. Instrum.* **78**:073702, 2007.
47. B. Kanngie[ss]er, N. Kemf, and W. Malzer, “Spectral and Lateral Resolved Characterisation of X-Ray Microbeams,” *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms* **198**(3–4):230–237, December 2002.

48. E. Langer, S. Dabritz, W. Hauffe, and M. Haschke, "Advances in X-Ray Excitation of Kossel Patterns by a Focusing Polycapillary Lens," *Applied Surface Science, 13th Applied Surface Analysis Workshop—AOFA 13* 252(1):240–244, September 30, 2005.
49. Z. Smit, K. Janssens, K. Proost, and I. Langus, "Confocal μ -XRF Depth Analysis of Paint Layers," *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms, Proc. Sixteenth International Conference on Ion Beam Analysis*, 219,220:35–40, June 2004.
50. M. Zitnik, P. Pelicon, N. Grlj, A. G. Karydas, D. Sokaras, R. Schutz, and B. Kanngiesser, *Appl. Phys. Lett.* 93:094104 (2008).
51. R. Alberti, A. Bjeoumikhov, N. Grassi, C. Guazzoni, T. Klatka, A. Longoni, and A. Quattrone, "Use of Silicon Drift Detectors for the Detection of Medium-Light Elements in PIXE," *Nuclear Instruments and Methods in Physics Research Section B: Beam Interactions with Materials and Atoms, Accelerators in Applied Research and Technology - Proceedings of the 9th European Conference on Accelerators in Applied Research and Technology* 266(10):2296–2300, May 2008.
52. W. M. Gibson and M. A. Kumakhov, in *X-Ray Detector Physics and Applications, Proc. SPIE* 1736:172, 1992.
53. T. Sun, M. Zhang, X. Ding, Z. Liu, X. Lin, and H. Liu, "Characterization of a Polycapillary X-Ray Lens for Application in Confocal Three-Dimensional Energy-Dispersive Micro X-Ray Diffraction Experiments," *J. Appl. Cryst.* 40:1169–1173, 2007.
54. J. Yang, D. Zhao, Q. Xu, and X. Ding, "Development, Application of Glancing Incident X-Ray Fluorescence Spectrometry Using Parallel Polycapillary X-Ray Lens," *Appl. Surface Science*, In Press, Accepted Manuscript, Available online October 14, 2008.
55. T. Sun, Y. Xie, Z. Liu, T. Liu, T. Hu, and X. Ding, "Application of a Combined System of Polycapillary X-Ray Lens and Toroidal Mirror in Micro-X-Ray-Absorption Fine-Structure Facility," *Appl. Phys.* 99:094907, 2006.
56. F. A. Hofmann, W. M. Gibson, C. A. MacDonald, D. A. Carter, J. X. Ho, and J. R. Ruble, "Polycapillary Optic—Source Combinations for Protein Crystallography," *J. Appl. Crystallogr.* 34:330–335, 2001.
57. M. Gubarev, E. Ciszak, I. Ponomarev, W. Gibson, and M. Joy, "A Compact X-Ray System for Macromolecular Crystallography," *Rev. Sci. Instrum.* 71:3900, 2000.
58. C. A. MacDonald, S. M. Owens, and W. M. Gibson, "Polycapillary X-Ray Optics for Microdiffraction," *J. Appl. Crystallogr.* 32:160–167, 1999.
59. J. X. Ho, E. H. Snell, C. R. Sisk, J. R. Ruble, D. C. Carter, S. M. Owens, and W. M. Gibson, "Stationary Crystal Diffraction with a Monochromatic Convergent X-Ray Beam Source and Application for Macromolecular Crystal Data Collection," *Acta Cryst.* D54:200–214, 1998.
60. P.-W. Li and R.-C. Bi, "Applications of Polycapillary X-Ray Optics in Protein Crystallography," *J. Appl. Cryst.* 31:806–811, 1998.
61. R. A. Clapp and M. Haller, "Parallel Beam Methods in Powder Diffraction and Texture in the Laboratory," *Adv. X-Ray Anal.* 43:135–140 (2000).
62. X. Chen, S. Bates, and K. R. Morris, "Quantifying Amorphous Content of Lactose Using Parallel Beam X-Ray Powder Diffraction and Whole Pattern Fitting," *J. Pharma. Biomed. Ana.* 26(1):63–72, August 2001.
63. W. Zhou, D. N. Mahato, C. A. MacDonald, "Analysis of Powder X-Ray Diffraction Resolution Using Collimating and Focusing Polycapillary Optics," *Thin Solid Films*, accepted.
64. S. T. Mixture and M. Haller, "Application of Polycapillary Optics for Parallel Beam Powder Diffraction," *Adv. X-Ray Anal.* 43:248–253, 2000.
65. W. M. Gibson, H. Huang, J. Nicolich, P. Klein, and C. A. MacDonald, "Polycapillary Optics for Angular Filtering of X Rays in Two Dimensions," in *Proc. 50th 2001 Denver X-Ray Conference, Advances in X-Ray Analysis* 45:F-58, 2002.
66. M. Leoni, U. Welzel, and P. Scardi, "Polycapillary Optics for Materials Science Studies: Instrumental Effects and Their Correction," *J. Res. NIST* 109(1):27–48, 2004.
67. A. Bingölbali, W. Zhou, D. N. Mahato, and C. A. MacDonald, "Focused Beam Powder Diffraction with Polycapillary and Curved Crystal Optics," in *Advances in X-Ray/EUV Optics and Components III*, Proc. SPIE, A. M. Khounsary, C. Morawe, S. Goto, eds., SPIE, Bellingham Wash., 7077, 2008.
68. W. C. Röntgen, "On a New Form of Radiation," *Nature*, 53, 274–276, January 23, 1896, English Translation from *Sitzungsberichte der Würzburger Physik-med. Gesellschaft*, 1895.
69. Cari, Suparmi, W. M. Gibson, and C. A. MacDonald, "Contrast Enhancement Measurements Using Polycapillary X-Ray Optics at 20–40 keV," in L.E. Antonuk, M.J. Yaffe, eds., *Medical Imaging 2001:Physics of Medical Imaging, Proc. SPIE* 4320:163–170, 2001.

70. Suparmi, Cari, Lei Wang, Hui Wang, W. M. Gibson, and C. A. MacDonald, "Measurement and Analysis of Leaded glass Capillary Optic Performance for Hard X Ray Applications," *J. Appl. Phys.* **90**(10):5363–5368, 2001.
71. R. E. Ross, C. D. Bradford, and W. W. Pepler, "Optimization of X-Ray Capillary Optics for Mammography," *Med. Imaging*, 2002.
72. R. Wentink, J. Carbone, D. Aloisi, W. M. Gibson, C. A. MacDonald, Q. E. Hanley, R. E. Fields, and M. B. Denton, "Charge Injection Device (CID) Technology: An Imaging Solution for Photon and Particle Imaging Applications," *Proc. SPIE* **2279**, 1994.
73. V. A. Somenkov, A. K. Tklich, and S. S. Shilstein, "X-Ray Refraction Radiography of Biological Objects," *Zhurnal Tekhnicheskoy Fiziki* **11**:197, 1991.
74. F. E. Carroll, "Generation of Soft X-Rays by Using the Free Electron Laser as a Proposed Means of Diagnosing and Treating Breast Cancer," *Lasers Surg. Med.* **11**:72–78, 1991.
75. P. C. Johns, D. J. Drost, M. J. Yaffe, and A. Fenster, "Dual-energy Mammography: Initial Experimental Results," *Med. Phys.* **12**:297–304, May/June 1985.
76. R. Schmitz, A. Bingölbali, A. Hussain, and C. A. MacDonald, "Development of Polarized and Monochromatic X-Ray Beams from Tube Sources in Advances," in A. M. Khounsary, C. Morawe, S. Goto, eds., *X-Ray/EUV Optics and Components III*, *Proc. SPIE* **7077**, 2008.
77. H. Chen, R. G. Downing, D. F. R. Mildner, W. M. Gibson, M. A. Kumakhov, I. Yu Ponomarev, and M. V. Gubarev, "Guiding and Focusing Neutron Beams using Capillary Optics," *Nature* **357**:391, 1992.
78. Q. F. Xiao, H. Chen, V. A. Sharov, D. F. R. Mildner, R. G. Downing, N. Gao, and D. M. Gibson, "Neutron Focusing Optic for Submillimeter Materials Analysis," *Rev. Sci. Instrum.* **65**:3399–3402, 1994.
79. D. F. R. Mildner, H. H. Chen-Mayer, and R. G. Downing, "Characteristic of a Polycapillary Neutron Focusing Lens," *Proceed. of Intl. Symp. on Adv. in Neutron Optics and Related Research Facilities*, March 19–21 (1996), *J. Phys. Soc. Jpn.*
80. R. F. Barth, A. H. Soloway, and R. G. Fairchild, "Boron Neutron Capture Therapy for Cancer," *Sci. Am.* **263**(4):100–103, 1990.
81. A. J. Peurrung, "Capillary Optics for Neutron Capture Therapy," *Med. Phys. Apr.* **23**(4):487–494, 1996.
82. R. Mayer, J. Welsh, and H. Chen-Mayer, "Focused Neutron Beam Dose Deposition Profiles in Tissue Equivalent Materials: A Pilot Study of BNCT," presented at the *5th International Conference on Neutron Techniques*, Crete, Greece, June 9–15, 1996.

This page intentionally left blank.

DO NOT DUPLICATE

SUBPART

5.4

X-RAY SOURCES

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

X-RAY TUBE SOURCES

Susanne M. Lee

*GE Global Research
Nikayuna, New York*

Carolyn MacDonald

*University at Albany
Albany, New York*

54.1 INTRODUCTION

An x-ray measurement system includes a source, optics (at least pinhole optics or apertures), and a detector. To optimize the system, the requirements of the application and the properties of the available sources and detectors must be considered. This chapter discusses the properties of the most common type of source, x-ray tubes, and how they affect optics optimization. Even though tube sources have been around for more than 100 years, recent advances have been made in tube design and performance. These include the method by which the electrons are generated, submicron focal spot sizes, sources with multiple targets, and the ability to direct the x-ray beam to a desired location.

The most common x-ray tubes generate x rays by accelerating electrons toward a target anode. The electron beam impact on the target generates x rays, as shown in Fig. 1. The mechanisms for x-ray generation and the resulting spectra are discussed in Sec. 54.2. Electron generation in the cathode is discussed in Sec. 54.3. The design of the anode, considerations in the choice of anode material, and limitations in the electron spot size and shape on the anode are discussed in Sec. 54.4. General optimization considerations are discussed in Sec. 54.5.

X-ray tube sources can be unipolar or bipolar, depending on how the electron-accelerating potentials are distributed between the source anode and cathode. In unipolar sources, such as the one shown in Fig. 1, the anode can be grounded, and the accelerating potential created by keeping the cathode at a large negative voltage. Alternatively, the electron source can be grounded with the anode at positive voltage. One advantage of anode grounding is the elimination of high-voltage standoffs, allowing much smaller distances between the x-ray generation point on the anode and the x-ray source exit window. This is beneficial for many x-ray optics, since short input optic distances allow the optic to collect from large source solid angles. In bipolar circuits, the anode is usually held at half the full potential between the cathode and anode; the cathode is then held at the same negative potential. By splitting the potential, smaller high-voltage standoffs and lighter electrical cables can be used, which increase the portability of such sources. Since portability is important in many nondestructive testing applications, industrial x-ray tubes are frequently bipolar.

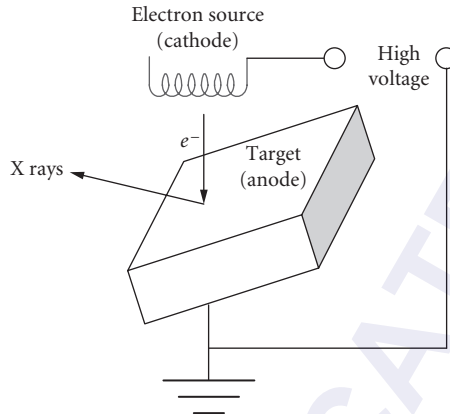


FIGURE 1 Schematic of a common unipolar anode-grounded electron impact source. Electrons emitted from the cathode, usually inside a vacuum vessel, are accelerated toward a metal anode by an accelerating voltage, which typically ranges from 10,000 to 500,000 V.

54.2 SPECTRA

Two basic processes convert the kinetic energy of the electrons into x rays, bremsstrahlung radiation, and characteristic emission.

Bremsstrahlung Continuous Radiation

The high-energy electrons emitted from the cathode impact the target and are decelerated by their interaction with the nuclei and electrons of the target material, converting the kinetic energy of the impacting electrons into electromagnetic energy. The x rays generated by this mechanism, known as braking radiation or bremsstrahlung, have a broad energy spectrum, with the high-energy cut-off determined by the maximum kinetic energy of the incident electrons, $K_{\max} = eV_{\text{tube}}$, where e is the electron charge and V_{tube} is the accelerating voltage. The efficiency of converting incident electron kinetic energy into continuum x-ray radiation energy is approximately¹

$$\eta = \frac{P_{\text{brems}}}{P_{\text{electrical}}} = \frac{1}{2} CZ_{\text{target}} (K_{\max})^{\alpha} \cong \frac{C}{2} Z_{\text{target}} (eV_{\text{tube}}) \quad (1)$$

where P_{brems} is the total bremsstrahlung radiation power, $P_{\text{electrical}}$ is the electrical power of the tube, $P_{\text{electrical}} = I_{\text{tube}} V_{\text{tube}}$, $\alpha \cong 1$, C is Kramer's constant, $C \cong 2.2 \times 10^{-9}/eV$, K_{\max} is the maximum kinetic energy of the electrons, $K_{\max} \sim eV_{\text{tube}}$, and Z_{target} is the atomic number of the target material. The efficiency of bremsstrahlung production is low, even for a high Z material at high voltage. For example, for a tungsten target with an electron accelerating potential of 100 kV, only 0.8% of the impacting electron beam energy is converted into x rays. Most of the energy in the incident electron beam is converted to heat.

The total bremsstrahlung power from Eq. (1) is

$$P_{\text{brem}} = \eta P_{\text{electrical}} \cong \frac{C}{2e} Z_{\text{target}} (eV_{\text{tube}})^2 I_{\text{tube}} \quad (2)$$

where I_{tube} is the tube current. The emitted power depends roughly on the square of the electron-accelerating voltage.

For thin targets, the bremsstrahlung energy emission spectrum is usually regarded as flat, independent of x-ray energy, up to the maximum energy, where the intensity vanishes.² The differential bremsstrahlung energy intensity, $E d\sigma$, between E and $E + dE$ is then a constant, i.e.,

$$E d\sigma \sim (\text{const}) Z_{\text{target}} dE \quad (3)$$

where $d\sigma$ is the differential bremsstrahlung cross section. For thicker target materials, the impacting electrons undergo multiple scattering. The bremsstrahlung spectrum then for thicker targets becomes a superposition of thin-target bremsstrahlung spectra from thin layers at increasing depths, with successively lower cutoff energies due to the electron energy losses as the electrons penetrate into the material. At moderate tube voltages, up to approximately 100 kV, the shape of this spectrum, shown in Fig. 2, is roughly linear in x-ray energy.^{3,4}

While the x-ray spectrum produced by the electrons is roughly triangular for thick targets, if the target is thick enough and the electron penetration depth deep enough, the x rays can be absorbed by the target as they exit it, which results in a preferential decrease in the low energy x-ray intensity actually emitted by the tube. A similar loss in low energy photons is produced by intentional filtration and by absorption due to the source and detector windows and any air path the x rays traverse. As a result, while the Kramer approximation reproduces the general shape of the Bremsstrahlung from a real x-ray tube, it does not model the detailed spectrum well.

Accurate bremsstrahlung cross sections are difficult to calculate from first principles. A number of cross sections have been proposed with varying degrees of success (see Nakamori et al. for an excellent review of the literature).⁵ Birch and Marshall,⁶ proposed a more precise semi-empirical model where

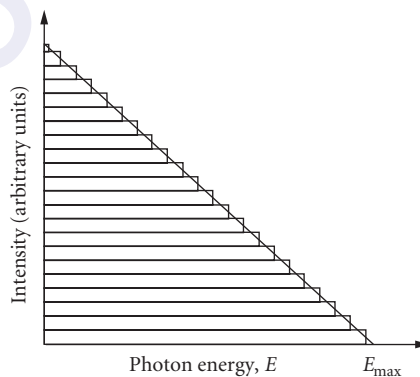


FIGURE 2 Triangular shape of bremsstrahlung intensity from a thick target, showing the superposition of flat spectra from thin-target bremsstrahlung with decreasing kinetic energy as the electrons penetrate into the material.

the bremsstrahlung cross section $d\sigma$ was modified to include a polynomial energy dependence and a relativistic correction factor,

$$\begin{aligned}
 Ed\sigma &= \rho N_A \left(1 + \frac{K_e}{m_0 c^2}\right) Q(E) dE \\
 &= \rho N_A \left(1 + \frac{K_e}{m_0 c^2}\right) \left(\frac{Z_{\text{target}}^2}{K_e}\right) \left[0.0503 - 0.94597 \left(\frac{E}{K_e}\right) + 0.1553 \left(\frac{E}{K_e}\right)^2\right. \\
 &\quad \left.+ 1.1632 \left(\frac{E}{K_e}\right)^3 - 0.6818 \left(\frac{E}{K_e}\right)^4\right] dE
 \end{aligned} \tag{4}$$

where N_A is Avagadro's number, ρ is the target material density, and $m_0 c^2$ is the electron rest mass energy. The polynomial coefficients were determined by fitting their calculated bremsstrahlung intensity to measured spectra.

Birch and Marshall also included x-ray attenuation in the target. In their model, the number of bremsstrahlung x rays having energy between E and $E + dE$ is given by⁶

$$N(E)dE = \frac{1}{A} \int_E^{eV_{\text{tube}}} d\sigma \left(\frac{1}{\rho} \frac{dK_e}{dx}\right)^{-1} \left[e^{-\frac{\mu(E)(eV_{\text{tube}}^2 - K_e^2) \cos \gamma}{\rho C_{\text{TW}} \sin \alpha}} \right] dK_e \tag{5}$$

where A is an empirical correction factor, $d\sigma$ is the differential bremsstrahlung cross section for x-ray emission with x-ray energy between E and $E + dE$, K_e is the electron kinetic energy, x is the electron depth into the target, $\mu(E)$ is the linear x-ray absorption coefficient in the target material at the x-ray energy E , γ is the angle the incident electron makes with the target, C_{TW} is the empirically determined Thomson-Whiddington constant, and α is the angle the exiting x ray makes with the target surface. The expression in brackets $[1/\rho(dk_e/dx)^{-1}]$ is the inverse of the electron stopping power of the target and tends to be a slowly varying function of the electron energy, such as $(a + be^{-cK_e})$, where a , b , and c are estimated from experimental data, usually with a least squares fit. The exponential factor in Eq. (5) accounts for the absorption of generated x rays in the target material. Both the electron stopping power and x-ray absorption in the target have been tabulated as a function of electron energy.⁷ Note that Eq. (5) gives the number of photons per energy interval and must therefore be multiplied by the photon energy E and integrated to obtain the total bremsstrahlung power.

Characteristic Radiation

If the kinetic energy of the electron impacting the target exceeds the electron binding energy of the target material, bound target electrons can be ejected. Then either a free electron, or one bound in an upper shell of the target atom, can fill the vacant lower energy state, in the process emitting a photon. If the vacant lower energy state is an inner shell, the emitted photons will be in the x-ray regime for most common target materials. Because these emitted x rays have an energy that depends on the core electronic levels of the target material, they are known as characteristic x rays. The lines are enumerated according to the vacancy core level, with K levels having principal quantum number $n = 1$, L levels $n = 2$, and M levels $n = 3$. The α transitions have $\Delta n = 1$; the β transitions $\Delta n = 2$. The energy of the core atomic level with quantum number n is given by

$$E_n = -E_0 \frac{(Z - \zeta)^2}{n^2} \tag{6}$$

where E_0 is the Rydberg energy,

$$E_0 = \frac{me^4}{8\epsilon_0^2 h^2} = 13.6 \text{ eV} \quad (7)$$

Z is the atomic number of the anode material, and ζ is a screening constant that describes the reduction in effective nuclear charge experienced by outer shell electrons. The screening constant increases with n ; for K-level electrons, $\zeta \sim 3$; for L levels, $\zeta \sim 12$ for elements with moderate Z . For copper, with $Z = 29$, the energy of a K-shell electron is -9 keV. The $1/n^2$ dependence in Eq. (6) results in bunching of the upper levels so that the L energy is just -1 keV. When an electron from the L shell fills a vacancy in the K shell, a copper $K\alpha$ photon is emitted with energy 8.0 keV. The $K\beta$ line energy, which is due to filling the K shell vacancy with the slightly higher-energy M-shell electron, is 8.9 keV. The levels, and hence emission lines, are further split according to the total angular momentum quantum number J , resulting in doublet and multiplet lines.

Characteristic emission requires overcoming the inner-shell electron binding energy of the target atoms. (For low Z targets the binding energy of the electron in the solid is nearly equal to the ionization energy of a gaseous atom, although relaxation effects can be significant for tungsten.) No emission will be seen unless the tube voltage imparts to the impacting electron a kinetic energy greater than the electron binding energy in the target. For example, to produce the K-line emission from a copper target, the tube voltage must be larger than 9 kV. Above this minimum voltage V_{\min} the characteristic photon emission rate increases with tube current I_{tube} and tube voltage V_{tube} roughly as

$$F = A \frac{I_{\text{tube}}}{e} \left(\frac{V_{\text{tube}}}{V_{\min}} - 1 \right)^p \quad (8)$$

where A is a constant determined by the target material, eV_{tube} is the maximum electron impact energy, and the exponent p ranges from 1.2 at high tube voltages to 1.7 at low tube voltages. At high tube voltages the electron penetration into the target is greater, leading to greater x-ray reabsorption in the target. The coefficient A gives the ratio of the number of characteristic photons emitted per incident electron when the tube voltage is twice the minimum voltage needed to excite the characteristic line, and is typically around 1 photon per 1000 to 6000 electrons. An example of a measurement and fit to a single value of A is shown in Fig. 3.⁸ This measurement was performed with a small pinhole to avoid flooding the detector. Measurements can also be made with filters of known

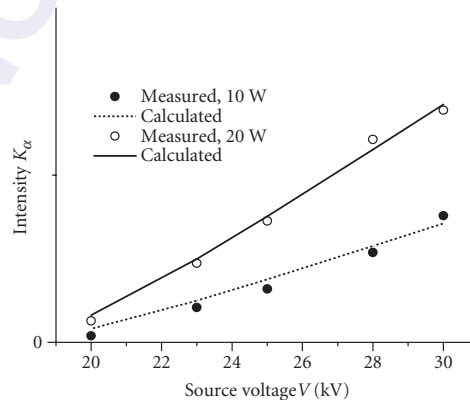


FIGURE 3 Intensity of Mo K_α doublet versus anode voltage at 20 W (\circ) and at 10 W (\bullet). The solid and dashed lines are the calculated intensities.

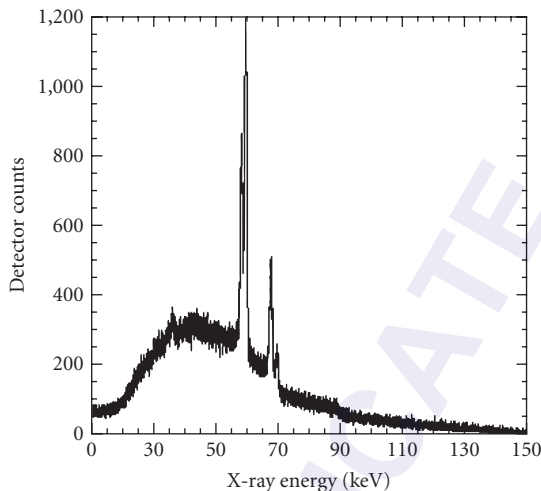


FIGURE 4 Source spectrum from a tungsten anode operated at 150 kV. The characteristic lines are superimposed on the bremsstrahlung, which has a maximum energy of 150 keV, according to Eq. (2). At low energies the spectrum differs significantly from the triangular approximation due to x-ray absorption in the target, air, windows of the x-ray tube, and the detector.

transmission. In either case, care must be taken so that errors in pinhole alignment, pinhole size, or filter thickness and composition do not affect the result, as the measured count rate must be orders of magnitude smaller than the source emission rate for most detectors.

The width of the characteristic emission lines is quite small; for Cu $K\alpha$ it is 4 eV and for Mo $K\alpha$ it is 8 eV.⁹ As a result, while the characteristic emission rate is typically similar in magnitude to the bremsstrahlung rate, the characteristic radiation is much brighter, and the diffracted intensity from, e.g., a narrow bandwidth monochromator, will be very much higher for the characteristic radiation.

Chapter 36 contains tables of the characteristic lines and corresponding energies for the elements commonly used in x-ray analysis. *The International Tables for X-Ray Crystallography*⁴ contain data on the efficiency with which various characteristic x-ray energies are produced as a function of incident electron beam energies. A typical spectrum from a tungsten anode, where both the characteristic lines and bremsstrahlung continuum can be seen, is shown in Fig. 4.

Spectral Selection

The Effect of Anode Material on Spectrum A consequence of the Z dependence in Eq. (8) is that increasing the atomic number of the anode slowly increases the emission energy of a particular emission line; e.g., the $K\alpha$ energy for W ($Z = 74$) is 59.3 keV compared to 8.0 keV for Cu ($Z = 29$). For that reason, for applications where monochromatic beams are desired, the anode material is chosen for the required x-ray energy. For example in mammography (described in Chap. 31), the anode material is chosen to provide a more intense beam in the 20 to 30 keV range, requiring selection of Mo, Rh, or Ag targets. In x-ray fluorescence (XRF) applications (described in Chap. 29), the anode material is chosen to be slightly higher in atomic number than the element to be sampled, since this will maximize the fluorescence emission intensity. In x-ray diffraction, when the sample is weakly diffracting, a Cu rather than Mo or W tube is frequently chosen because, at the lower energies, the cross section for diffraction is higher. In monochromatic x-ray diffraction, where an intense characteristic line is desired with a low intensity continuum, increasing the atomic number will have

a deleterious effect on the ratio of characteristic to continuum emission, since, according to Eq. (1), the continuum intensity will increase. Thus, the highest resolution x-ray diffraction measurements performed with laboratory sources use sources with lower atomic number targets.

In other applications where it is more important that a high total x-ray power is generated, such as white-beam diffraction or some imaging applications, the increase in bremsstrahlung continuum intensity with atomic number is a positive effect, and higher Z sources, like W, are common.

The Effect of Tube Voltage on Spectrum Many tubes use pulsed rather than continuous operation, and even in continuous operation the tube voltage may not be stable, so the tube voltage is usually specified in terms of peak voltage, kVp.

The tube voltage is also used to shape the spectrum. According to Eqs. (2) and (8), the total unfiltered bremsstrahlung power increases approximately as the square of the tube voltage, while the total characteristic intensity increases somewhat slower, especially for high accelerating voltages. The resultant ratio of the *power* in the characteristic line to the *total* bremsstrahlung *power* is highest at a tube voltage corresponding to

$$eV_{\text{tube}} = \frac{\alpha + 1}{\alpha + 1 - p} E_{\text{char}} \cong \frac{2}{2 - p} E_{\text{char}} \quad (9)$$

For low voltages ($p \sim 1.7$), the electron energy eV_{tube} should be roughly six to seven times the characteristic line energy, or about 45 to 55 kVp for Cu $K\alpha$ radiation. For higher voltages ($p \sim 1.2$), the tube voltage need only be approximately 2.5 times the electron binding energy, or about 150 kVp for W $K\alpha$ radiation.

However, the ratio of the characteristic photon emission rate to the bremsstrahlung photon emission increases monotonically with tube voltage. The relative importance of the characteristic versus bremsstrahlung intensities depends on the experiment—in many imaging applications, the image reconstruction process is simplified considerably if the tube voltage does not excite the characteristic energies. In monochromatic x-ray diffraction applications, the characteristic lines are usually preferred; the monochromator, x-ray optics, and type of x-ray detector (photon counting or energy-integrating detector) can impact the characteristic to bremsstrahlung ratio. One consideration with x-ray optics is whether they will pass undesirable high-energy components, e.g., integral multiples of the characteristic energy, which may affect the choice of tube voltage. In addition, at a fixed current, increasing the tube voltage increases the power that must be dissipated in the anode, and thus generally the source size.

Filtering For some applications, such as low-resolution diffraction, a sufficiently monochromatic spectrum can be obtained by employing a simple absorption filter with a K-edge between the $K\alpha$ and $K\beta$ emission lines of the anode material. For example, Ni, with $Z = 28$ is used to absorb the Cu ($Z = 29$) $K\beta$ emission in many laboratory diffractometers. The filter also absorbs the low-energy bremsstrahlung, but the tube voltage should be set only slightly above the characteristic line voltage to avoid creating too much high-energy bremsstrahlung. Aluminum filters are also commonly used in medical imaging systems to remove the low-energy photons that would otherwise be absorbed in the patient as unnecessary dose.

Because the filter preferentially transmits high-energy photons, the bremsstrahlung power output rises more rapidly with tube voltage than the V^2 factor for thick targets of Eq. (2), to more than V^3 under heavy filtration.¹⁰

Multi-Energy Imaging Some materials have similar x-ray absorption properties for a limited range of energies, yet different properties for other energies. By imaging these objects with multiple spectral ranges, the materials can be differentiated from each other. This is known as multi-energy imaging and can be accomplished in several different ways. One way is to exploit the tube-voltage dependence of Eqs. (2) and (5) to produce bremsstrahlung spectra with different x-ray energy distributions. Operating the x-ray source at two different tube voltages will produce x-ray beams with bremsstrahlung spectra having different maximum intensity energies and different high-energy

cutoffs. This technique is known as kVp-switching. Unfortunately, the differences in tube voltage needed to provide clear image differences can be on the order of 100 kV, which is a large voltage difference to switch quickly and stably.

An alternative approach to kVp-switching is to have multiple target materials present in one x-ray source housing. The electron beam is directed toward each target material sequentially, creating images made with different energy spectra. This technique is beginning to see more use in imaging applications such as mammography, where it has been difficult to distinguish between tumors and the surrounding soft tissue due to similar absorption properties. The different energy images together with contrast agents allow other features, such as the vascularization of the tumors, to be seen more readily. One issue with using different target materials is the physical displacement of the targets from each other, which causes the x rays from each target to intersect the object to be imaged at different angles; the different viewing angles then have to be compensated for in the image processing algorithms unless the target material is moved between images.

54.3 CATHODE DESIGN AND GEOMETRY

Hot Filament Sources

Typically, a helical wire filament is heated in vacuum to slightly below the melting point of the filament material ($\sim 1000^\circ\text{C}$ or higher), allowing electrons with high kinetic energy to escape. A grid at a slightly positive potential with respect to the filament is usually placed near the filament to focus and accelerate the electrons in the direction of the target. Because of the high operating temperature of this filament type, substantial warm-up times, often on the order of a half hour, are required before stable filament current is achieved. In addition, the high temperature results in substantial evaporation from the filament, eroding and reducing the filament lifetime. As the filament wire diameter decreases, so does the electron density that can be extracted from the material. This has the negative effect of causing the electron beam intensity to be nonuniform throughout its lifetime. Higher electron extraction density requires hotter temperatures, which results in shorter lifetimes. However, wire filaments are the most physically robust of all electron emitters and can withstand the high-voltage arcs that sometimes occur inside the tube without significant filament degradation.

An alternative type of hot filament electron emitter is the dispenser cathode, which consists of a high electron-emissivity thin film on top of an integral resistor. Such integration reduces the power consumption of the device to around 1/2 W, while providing excellent thermal stability and fast warm-up times, as short as a few seconds, even though the operating temperatures are still on the order of 1000°C . The current densities obtained from dispenser cathodes, $\sim 8\text{ A/cm}^2$, are about four times higher than traditional hot filaments can provide, with similar lifetimes. In addition, since the architecture of dispenser cathodes is different than traditional hot filament cathodes, more uniform current densities are obtained over the cathode lifetime. However, the thin-film nature of these emitters makes them more susceptible to arc damage.

Cold-Cathode Field Emission Sources

Field emission sources, consisting of carbon nanotube (CNT) or small diameter nanorods of high-electron-emissivity materials, are just becoming available commercially. Currently, CNT sources are the more common. Since the diameter of the CNTs used in x-ray sources is so small, a few to tens of nanometers in diameter, very high electric field enhancements occur around the tips of the CNTs, reducing the need to apply high voltages to the whole device. In theory, voltages as low as a few tens of volts applied to a CNT device should be able to extract a sufficient electron beam density for x-ray source operation. In practice, hundreds to thousands of volts are used for the extraction process

on unheated CNT devices, although it is expected that as CNT device–manufacturing processes mature, the operating voltages will decrease. Since the operating temperature is so low, essentially no evaporation should occur, greatly increasing the “filament” lifetime. However, CNT adherence issues arise under arcing in the x-ray tube.

Inverse Geometry Sources

In conventional imaging, the source is well separated from the object and the detector, which is usually about the same size as the illuminated part of the object and in close proximity to it. In inverse geometry imaging applications, the object to be imaged is placed very close to the x-ray source and the detector is far away from both object and source. The small source-object distance requires either a prohibitively large x-ray beam or one that is scanned over the object. The latter approach is the most common for inverse-geometry sources; the electron beam is rastered across a target anode, which is approximately as large as the object being imaged. Rastering is usually accomplished by magnetically and/or electrically steering the electron beam across the target. The detector, which is typically placed about a meter away, is very small, roughly the size of a single pixel detector. Since the x-ray beam scans a large space, this type of source typically does not lend itself well to x-ray optic applications.

Electron Beam Steering and Multispot (Array) Sources

Even in conventional imaging, with the object and detector close to each other and the source farther away, being able to steer the electron beam—usually magnetically—to different positions on the target has the advantage that images from different viewing angles can be acquired simply and quickly. These images can then be combined, for example, to reconstruct a stereoscopic image of the object. Alternatively, multiple electron sources can be present inside a single x-ray source housing to accomplish the same type of imaging. Cold-cathode field emission devices particularly lend themselves to such multispot arrangements since they can be deposited in arrays with a huge number of addressable emitters.

54.4 EFFECT OF ANODE MATERIAL, GEOMETRY, AND SOURCE SIZE ON INTENSITY AND BRIGHTNESS

For electron impact sources, the two fundamental limiting factors on source brightness are the diameter of the electron beam on the target and the maximum power density the target can tolerate without melting, since most of the electron beam power is dissipated in the target as heat. This power density, determined by the electrical power, $I_{\text{tube}} V_{\text{tube}}$, limits the total flux and hence brightness that can be produced by the source, according to Eqs. (10) and (8). The brightness of the beam transmitted by an x-ray optic is limited, as stated in Liouville’s theorem,⁶ to this maximum source brightness.

Transmission Sources

In transmission sources, the electron beam is perpendicular to the target, which is thick enough to stop the electrons, but thin enough for the x rays to traverse the target thickness and exit the tube with minimal intensity losses. In reflection sources, the terms anode and target are used interchangeably, but in transmission targets they are not, since the anode tends to be positioned close to the cathode to provide the greatest electron acceleration, while the target (at ground potential)

is farther away, with electron focusing elements in between. Frequently the target material also acts as the source x-ray window, enabling focal-spot-to-optic distances of a few tens of microns, or the thickness of the transmission target. Another advantage of transmission sources, especially for many types of x-ray optics, is their small focal spot sizes, which can be as small as a few tenths of a micron with appropriate electron redirection and focusing optics. However, transmission source powers tend to be low due to limits on heat dissipation in the targets.

Target Material

For higher power sources metals with good thermal conductivity and high melting temperatures, such as tungsten and copper, are generally chosen. In addition, according to Eq. (10), the total source x-ray intensity increases with atomic number, so the most powerful sources have the largest atomic number targets, such as W or Re.

Stationary Anodes

To attain the maximum source brightness in a reflection source, where the electron beam makes an oblique angle with the anode target material and the target is stationary, direct anode cooling is necessary to prevent target melting. In this case, the target material has to be thick enough to separate the cooling medium from the vacuum required for the electron beam operation. Jets of water are frequently sprayed onto the backside of these thick targets to directly cool them. No matter how the anode is cooled, the greater the x-ray source power, the larger the x-ray beam diameter required to prevent melting. Müller⁷ has shown that the diameter (not, as one might expect, the area) increases linearly with power. For sealed-tube copper x-ray sources, the maximum power is typically 5 kW/mm of x-ray beam diameter. The x-ray beam focal spot size in stationary anode tubes can range from microns to as large as 6 mm or more, for very high voltage tubes (above 500 kVp). Sources with submillimeter diameter focal spots are referred to as micro-focus sources and contain electron-focusing elements that aid in attaining the small spot diameter. Additionally, many microfocus sources also contain electron redirection elements to keep the electron beam precisely on the same target spot. A shift in the x-ray source position of a few hundred microns is not significant when the focal spot is several millimeters in diameter; however, a similar shift of a micron diameter beam can eliminate all x-ray intensity through an optic with a small acceptance angle. Source manufacturers address this issue by using magnetic fields to correct the electron beam direction. These magnetic fields can be generated with simple external permanent magnets or with internal electromagnets connected into a feedback loop that monitors the beam position and adjusts the electromagnet current to re-center the electron beam.

Rotating Anodes

Alternatively, in reflection sources, the heat load can be reduced by rotating the target so rapidly that the electron beam dwell time on the target is too short to thermally overload the target. With advances in bearing lubricants, very fast rotations are possible. Maximum powers of roughly 150 kW/mm of x-ray beam diameter can be achieved without any direct anode cooling. The electrons that are not absorbed by the rotating target are called *secondary electrons*. These electrons, while lower in energy than the incident beam, still carry considerable kinetic energy that must be removed from the system with an electron collector, usually a sizeable piece of copper with a slit in it to let the x rays pass. The high power density of rotating anode sources usually results in larger x-ray focal spot sizes (roughly 0.1 mm to about 3 mm in diameter) than in stationary target sources. Although the x-ray power of rotating anode sources is a couple of orders of magnitude higher than for stationary targets, with accordingly more intense x-ray beams, rotating anode systems typically are larger and more costly to maintain. They also generally require larger source to window distances, limiting the available acceptance angle for small x-ray optics.

Liquid-Metal Anodes

All commercial x-ray sources have solid phase target materials. However, over the years liquid metal targets have been repeatedly proposed. The basic concept consists of a liquid metal, e.g., GaInSn or hot Sn (300°C), flowing turbulently past a thin barrier, e.g., molybdenum foil or glass, that is reasonably electron transparent; the barrier separates the liquid metal from the vacuum in which the cathode has to operate. The electrons from the cathode pass through the vacuum/liquid barrier and interact with the liquid target, producing x rays in much the same way they would from a solid target, with the liquid metal removing the heat generated during electron impact as the metal flows away from the impact region. The bremsstrahlung energy distribution from a liquid target is slightly different than that from a solid target, since the reduced density of atoms in a liquid will skew the energy distribution toward lower energies compared to those from a solid target. Liquid metal anode sources have the potential for much improved continuous loading capability and significantly higher maximum loading (a few hundred kW/mm x-ray beam diameter).

Source Shape

In many reflection sources, the electron beam intersects the target in a line rather than a circular spot. The x rays are then viewed at a small takeoff angle with respect to the target, resulting in what appears to be a nearly square or circular source spot, as shown in Fig. 5. Because the apparent size of the beam is reduced by the sine of the take-off angle, the brightness in the resultant x-ray beam is increased to

$$\text{Brightness} \propto (\text{power})/(\text{apparent area}) = (\text{power})/(A \sin \alpha) \quad (10)$$

where A is the actual electron beam spot area on the target and α is the takeoff angle (see Fig. 5). This can considerably improve the brightness for some applications. However, the x rays are not produced at the surface of the target, but a few microns into the target, depending on the electron stopping power of the target material and the electron voltage. Absorption of the x rays traveling out of the target material will eventually reduce the source brightness at small grazing angles. The radiated power is reduced exponentially with the path length,

$$\text{Brightness} \propto \frac{\text{power out}}{A \sin \alpha} = \frac{\text{power}}{A} \frac{e^{-(d/\sin \alpha)/D}}{\sin \alpha} \quad (11)$$

where d is the depth in which x rays are produced and D the absorption length for the x rays. The factor on the right is plotted in Fig. 6 for $d/D = 10$, which is approximately the case for a typical Cu K_α system. The brightness is usually maximized for a takeoff angle from 6° to 12° . Excessively small takeoff angles are more affected by surface roughness, which may increase with target use.

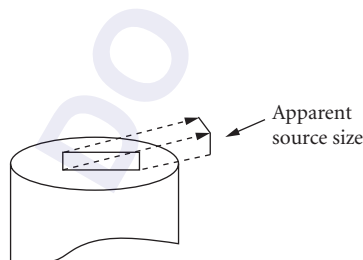


FIGURE 5 Apparent x-ray beam size from a line focus source when viewed at an angle.

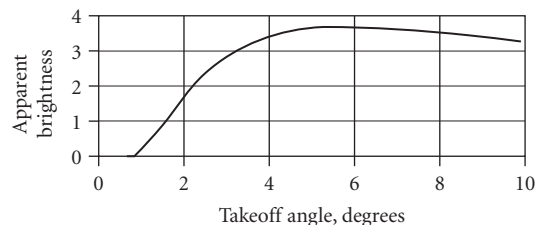


FIGURE 6 Source brightness as a function of takeoff angle for an x-ray production depth to absorption length ratio of 10.

Optics with sufficient depth of field to collect from the entire length of the line source can take advantage of the increased power emitted by line-focus sources compared to a smaller square source. It should be noted, however, that optics designed for very short source-to-optic distances often do not have large depths of field.

Source Depth and Size Measurement

The solid angle from which an optic can collect x rays depends on how close the optic can be placed to the x-ray generation point inside a source. For optics placed external to a source, this is limited by the distance from the vacuum window to the electron beam spot on the target. One method to measure this distance in practice is to move a pinhole in a series of steps transverse to the beam and record the intensity in multiple exposures, as shown in Fig. 7.⁸ The source position X relative to the pinhole is related to the image position Y and the source spot to pinhole distance L by

$$L = \frac{X}{Y} D \quad (12)$$

where D is the pinhole to image plate distance.

Most optics can only collect from relatively small source spots. Thus, source size measurement is very important for optics and system design, and also for optic characterization. Source sizes can be deduced from pinhole images. Because small pinholes may be thick compared to their diameter, it is important to ensure that the measured size is not influenced by pinhole shape, and so multiple images are required for accurate assessment.

Several images should be taken with the detector at different distances D from the pinhole, as shown in the sketch of Fig. 8. The slope of a plot of the image size W versus pinhole-to-detector distance D yields the source diameter S .

$$W = \frac{(S+2r)D}{L} + 2r \quad (13)$$

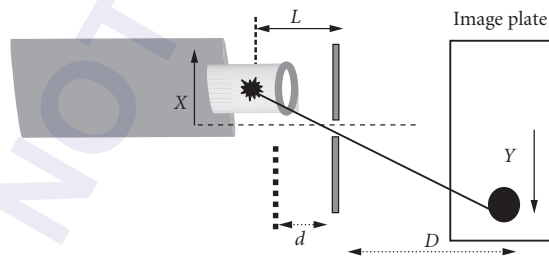


FIGURE 7 Schematic diagram of the setup for a source depth measurement.

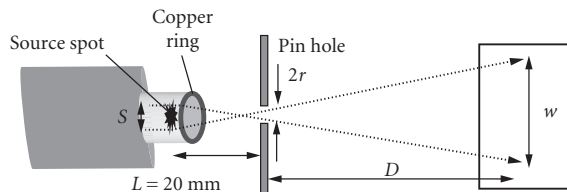


FIGURE 8 Schematic diagram of the setup for source spot measurement.

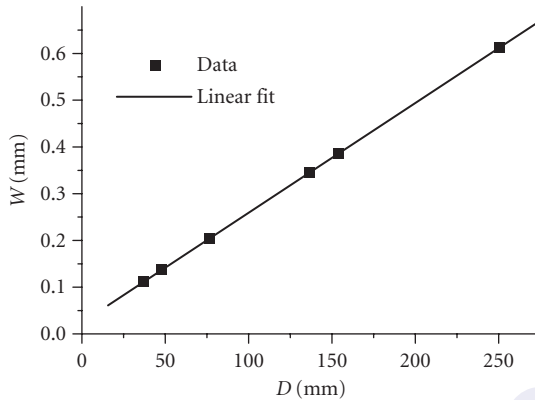


FIGURE 9 Image size W versus pinhole to image plate distance D taken with a 25- μm pinhole. The slope is 0.0024, yielding a source size of 25 μm , and intercept 0.023 mm, in good agreement with the nominal pinhole size.

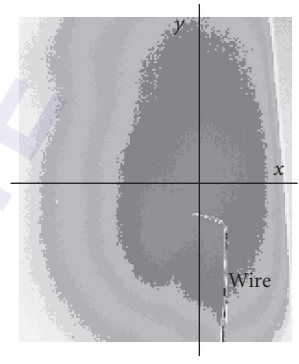


FIGURE 10 Source image at 10 cm from the source spot.

There should be good agreement between the intercept value and the pinhole radius, as shown in Fig. 9.

The beam emitted from a standard tube is often restricted in one direction due to shadowing by the anode, as shown in Fig. 10, which has a cone angle of approximately 15° in the horizontal direction.

Effective Source Size for X-Ray Optics

When applying Liouville's theorem to x-ray optics, the limiting factor is not the source focal spot size, but the area over which the optic can collect photons. Even if the optic can be placed close to the source, the optic collection area may be smaller than the actual focal spot, wasting source photons. However, large source sizes do not detract from beam brightness, if the power density is maintained and the optic only collects from a small area suitable for a specific application. For sources with nonuniform brightness, typically high peak brightness in a central region of the beam, the optic has to be carefully optimized to collect x rays from the smallest possible region that will provide sufficient power for the selected application.

54.5 GENERAL OPTIMIZATION

Trade-Offs between Low- and High-Power Sources

Low-power sources typically have shorter distances between the x-ray generation point on the target and the source exit window, where x-ray optics are usually positioned. Such small source-to-optic distances allow the optics to collect significantly larger solid angles than similar optics applied to larger sources. For applications, e.g., diffraction, where x-ray pencil beams are employed, or microanalysis, which requires small x-ray beam spots, x-ray optics coupled to low-power sources can provide similar intensities to higher-power sources, with greater convenience, portability, ease of maintenance, and lower cost. However, applications such as many medical and nondestructive imaging modalities, that require irradiating a large area in a short time use higher-power sources. These allow shorter exposure times, but are not conducive to efficient x-ray optic coupling.

Parallel Beam

As a result of Liouville's theorem, the x-ray source has to be smaller than the desired parallel beam diameter in order to achieve an intensity gain. The gain is produced by collecting over a solid angle from the source that is larger than the output divergence of the collimated beam. Therefore, it is particularly challenging to provide intense, small-cross-section, nearly parallel beams, as required by some microdiffraction applications.

Focused Beam

An optic can achieve large gains if it collects a diverging beam from a point source and refocuses it into a converging beam. The diameter of the optic is not limited by the application, but by manufacturing and integration constraints. For any optic, the minimum focal spot size will be smaller for a smaller optic-to-focus distance, as the output divergence of the optic is never zero. The best spatial definition of the beam will occur for smaller focal distances.

Effect of Different Optics Types on Brightness

Imaging optics such as pinholes and mirrors are able to conserve brightness. Pinholes are the most basic optic and can be used to define unity gain. They are also useful to further condition the beam and eliminate scatter at the sample location. A straight, untapered, single capillary tube preserves the beam condition across its length, if transport losses are ignored. Therefore, it is identical in performance to placing a pinhole at the input of the capillary. Straight single capillaries can be useful as improved pinholes or to provide more convenient geometries.⁸ (Capillaries are discussed in Chap. 52.) Mirrors are imaging optics and can preserve brightness. An ideal symmetric mirror produces a focal spot equal in size to the source spot. In order to demagnify the source focal spot size, the output focal distance has to be reduced from the symmetric case. (Mirrors are discussed in Chaps. 44 to 47.) The same considerations for mirrors are valid for diffractive optics, with the additional restriction that the Bragg diffraction condition must be met for the characteristic energy of the source target material, to obtain the most intense x-ray beam. (Diffractive optics are discussed in Chaps. 39 to 43.) Tapered single capillary tubes can provide a significantly enhanced flux over simple pinholes. If shaped appropriately, the tapered capillaries can serve as imaging micro-mirrors and fall under the mirror category. Polycapillary optics use many consecutive reflections to achieve a large overall deflection. The acceptance from the multiple capillary tubes can be overlapped so that the optic collects from a large solid angle with correspondingly large total power. However, polycapillary optics are not imaging optics and cannot conserve brightness. (Polycapillary optics are discussed in Chap. 53.)

Choosing a Source/Optic Combination

Choosing an appropriate source/optic combination depends on analyzing the requirements of the application for photon energy, energy bandwidth, beam size, beam divergence, and total power. If total x-ray power is the only consideration, a large focal spot source with high Z target material, high voltage, high current, small source-to-sample distance, and no optic is the most appropriate choice. If an optic is to be employed, it must be carefully designed to take into account the source geometry. Large focal spot sources with large source-to-output window distances are not well matched with optics designed to produce small bright beams. Detector properties, discussed in Chaps. 60 to 62, also affect the system optimization.

54.6 REFERENCES

1. H. Compton and S. K. Allison, *X Rays in Theory and Experiment*, 2nd ed., D. Van Nostrand, New York, 1935.
2. H. A. Kramers, "On the Theory of X-Ray Absorption and of the Continuous X-Ray Spectrum," *Phil. Mag.* **46**:836, 1923.
3. W. Potts, "Electron Impact X-Ray Sources," in *X-Ray Science and Technology*, A. G. Michette and C. J. Buckley (eds.), IOP Publishing, London, 1993.
4. J. L. Goldstein, D. E. Newbury, P. Echlin, D. C. Joy, C. Fiori, and E. Lifshin, *Scanning Electron Microscopy and X-Ray Microanalysis*, Plenum Press, New York, 1981, p. 97.
5. N. Nakamori, K. Yamano, M. Yamada, and H. Kanamori, "Effect of Electron Energy Distribution on Bremsstrahlung Spectrum," *Jpn. J. Appl. Phys.* **32**:4019–4025, 1993.
6. Birch and M. Marshal, "Computation of Bremsstrahlung X-Ray Spectra and Comparison with Spectra Measured with a Ge(Li) Detector," *Phys. Med. Biol.* **24**:505–517, 1979.
7. M. Tucker, G. T. Barnes, and D. P. Chakraborty, "Semiempirical Model for Generating Tungsten Target X-Ray Spectra," *Med. Phys.* **18**:211–218, 1991.
8. N. Mail, W. M. Gibson, and C. A. MacDonald, "Molybdenum Microfocus Source Coupling to Polycapillary Optics for Powder Diffraction," in *Advances in Laboratory-Based X-Ray Sources and Optics III*, A. M. Khounsary and C. A. MacDonald, (eds.), *SPIE* **4781**:87–95, 2002.
9. Thompson, D. Atword, E. Gullikson, et al., "X-Ray Data Booklet" CXRO: Berkeley, 2002. Web site: <http://www-cxro.lbl.gov>, accessed 31 May 2009.
10. F. O'Foghlu and G. A. Johnson, "Voltage Waveform Effects on Output and Penetration of W- and Mo-Anode Mammographic Tubes," *Phys. Med. Biol.* **26**:291–303, 1981.

This page intentionally left blank.

DO NOT DUPLICATE

SYNCHROTRON SOURCES

Steven L. Hulbert

*National Synchrotron Light Source
Brookhaven National Laboratory
Upton, New York*

Gwyn P. Williams

*Free Electron Laser
Thomas Jefferson National Accelerator Facility
Newport News, Virginia*

55.1 INTRODUCTION

Synchrotron radiation is a bright, broadband, polarized, pulsed source of electromagnetic radiation extending from the far-infrared to the hard x-ray region. *Brightness*, defined as flux per unit area per unit solid angle, is normally a more important quantity than flux or intensity, particularly in throughput-limited applications which utilize only a small fraction of the transverse phase space of the emitted radiation or a small energy bandwidth, or both.

It is well known from classical theory of electricity and magnetism that accelerating charges emit electromagnetic radiation. In the case of synchrotron radiation, relativistic electrons are accelerated in a circular orbit and emit electromagnetic radiation in a broad spectral range. The visible portion of this spectrum was first observed on April 24, 1947 at General Electric's Schenectady facility by Floyd Haber, a machinist working with the synchrotron team, although the first theoretical predictions were by Liénard¹ in the latter part of the 1800s. An excellent early history with references is presented by Blewett² and a history covering the development of the utilization of synchrotron radiation is presented by Hartman.³

Synchrotron radiation covers the entire electromagnetic spectrum from the far-infrared (or THz) and infrared, through the visible, ultraviolet, and x-ray regions and into the very hard x-ray range up to energies of 100 kilovolts and above. If the charged particles are of low mass, such as electrons, and if they are traveling relativistically, the emitted radiation is very intense and highly collimated, with opening angles, which depend inversely on the energy of the particle, on the order of 1 milliradian. In electron storage rings there are two distinct types of sources of synchrotron radiation: dipole (bending) magnets and insertion devices. Insertion devices are further classified as either wigglers, which act like a sequence of bending magnets with alternating polarities, or undulators, which are also multiperiod alternating magnet systems but in which the beam deflections are small resulting in coherent interference of the emitted light.

In typical storage rings used as synchrotron radiation sources, several bunches of up to $\sim 10^{12}$ electrons circulate in vacuum, guided by magnetic fields. The bunches are typically several 10s of centimeters long, so that the light is pulsed, being on for a few 10s to a few 100s of picoseconds, and off for several 10s to a few 100s of nanoseconds depending on the particular machine and the radio-frequency cavity which restores the energy lost to synchrotron radiation. The revolution time for a ring of circumference 30 m is 100 ns, so that each bunch of $\sim 10^{12}$ electrons is seen 10^7 times per second, giving a

current of ~ 1 A (assuming single bunch filling). In linacs that are used as light sources the electrons are also in bunches, but these are usually much shorter resulting in pulses of < 1 ps.

The most important characteristic of accelerators built specifically as synchrotron radiation sources is that they have a magnetic focusing system which is designed to concentrate the electrons into bunches of very small transverse cross section and to keep the electron transverse velocities small. The combination of high intensity with small opening angles and small source dimensions results in the very high brightness.

The first synchrotron radiation sources to be used were operated parasitically on existing high-energy physics or accelerator development programs. These were not optimized for brightness, and were usually accelerators rather than storage rings, meaning that the electron beams were constantly being injected, accelerated, and extracted. Owing to the successful use of these sources for scientific programs, a second generation of dedicated storage rings was built starting in the early 1980s. In the mid 1990s, a third generation of sources was built, this time based largely on insertion devices, especially undulators of various types. A fourth generation of accelerator-based photon sources is now coming on line, based on what is called multiparticle coherent emission, in which coherence along the path of the electrons, or longitudinal coherence, plays the major role. This is achieved by microbunching the electrons on a length scale comparable to or smaller than the scale of the wavelengths emitted. The emission is then proportional to the square of the number of electrons N which, if N is 10^{12} , can be a very large enhancement. These sources can approach the theoretical diffraction limit of source emittance (the product of solid angle and area).

55.2 THEORY OF SYNCHROTRON RADIATION EMISSION

General

The theory describing synchrotron radiation emission is based on classical electrodynamics and can be found in the works of Tomboulion and Hartman,⁴ Schwinger,⁵ Jackson,⁶ Winick,⁷ Hofmann,⁸ Krinsky, Perlman, and Watson,⁹ and Kim¹⁰. A quantum description, presented by Sokolov and Ternov,¹¹ is quantitatively equivalent.

Here we present a phenomenological description in order to highlight the general concepts involved. Electrons in circular motion radiate in a dipole pattern as shown schematically in Fig. 1a.

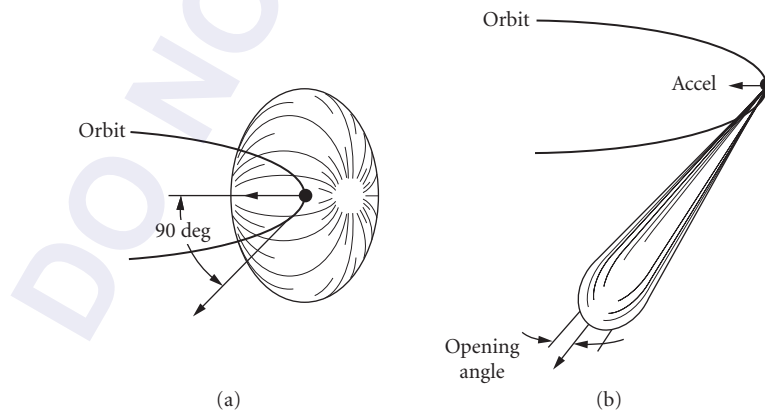


FIGURE 1 Conceptual representation of the radiation pattern from a charged particle undergoing circular acceleration at (a) nonrelativistic and (b) relativistic velocities.

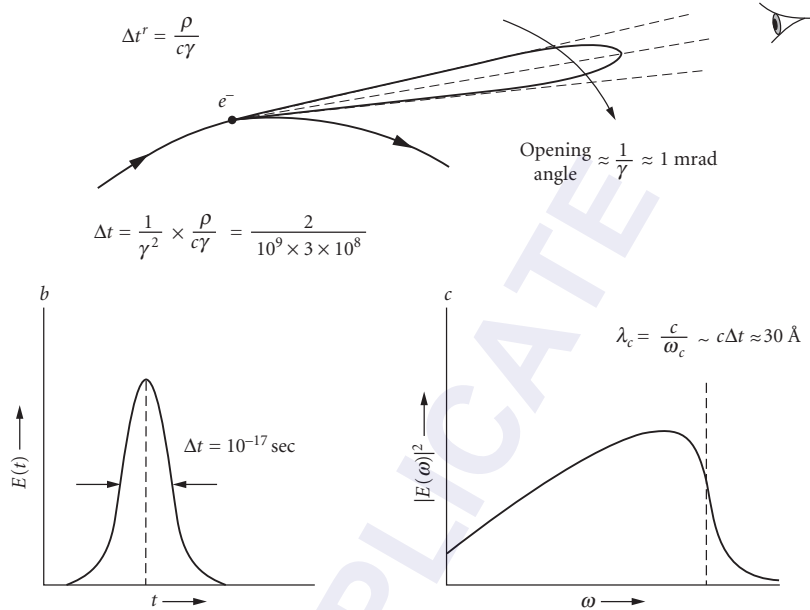


FIGURE 2 Illustration of the derivation of the spectrum emitted by a charged particle in a storage ring.

As the electron energies increase and the particles start traveling at relativistic velocities, this dipole pattern appears different to an observer in the rest frame of the laboratory. Special relativity tells us that angles θ_r in a transmitting object are related to those in the receiving frame θ_t by

$$\tan \theta_r = \frac{\sin \theta_t}{\gamma (\cos \theta_t - \beta)} \quad (1)$$

with γ , the ratio of the mass of the electron to its rest mass, being given by E/m_0c^2 , E being the electron energy, m_0 the electron rest mass, and c the velocity of light. β is the ratio of electron velocity v to the velocity of light c . Since $\beta \approx 1$ for electrons travelling at relativistic energies, the peak of the dipole emission pattern in the particle frame, $\theta_t = 90^\circ$, transforms to $\theta_r \approx \tan \theta_t \approx \gamma^{-1}$ in the laboratory frame as shown in Fig. 1b. Thus γ^{-1} is a typical opening angle of the radiation in the laboratory frame. For an electron viewed in passing by an observer, as shown in Fig. 2, the duration of the pulse produced by a particle under circular motion of radius ρ will be $\rho/\gamma c$ in the particle frame, or $\rho/\gamma c \times 1/\gamma^2$ in the laboratory frame owing to time dilation. The Fourier transform of this function will contain frequency components up to the reciprocal of this time interval. For a storage ring with a radius of 2 meters and $\gamma = 1000$, corresponding to a stored electron beam energy of ~ 500 MeV, the time interval is 10^{-17} seconds, which corresponds to light of wavelength 30 Å.

Bending Magnet Radiation

For an electron storage ring, the relationship between the electron beam energy E , bending radius ρ , and field B is

$$\rho = \frac{E}{ecB} = \frac{E[\text{GeV}]}{0.300B[\text{T}]} \quad (2)$$

where γ , the ratio of the mass of the electron to its rest mass is given by $\gamma = E/m_0c^2 = E/0.511 \text{ MeV} = 1957E$ [GeV] and λ_c , which is defined as the wavelength for which half the power is emitted above and half below, is

$$\lambda_c = \frac{4\pi\rho}{3\gamma^3} = 5.59 \text{ \AA} \frac{\rho[m]}{E^3[\text{GeV}^3]} = \frac{18.6 \text{ \AA}}{B[T]E^2[\text{GeV}^2]} \quad (3)$$

The critical frequency and photon energy are

$$\omega_c = \frac{2\pi c}{\lambda_c} = \frac{3c\gamma^3}{2\rho} \quad \varepsilon_c[eV] = \hbar\omega_c(eV) = 665.5E^2[\text{GeV}^2]B[T] \quad (4)$$

The angular distribution of synchrotron radiation flux emitted by electrons moving through a bending magnet with a circular trajectory in the horizontal plane is given¹⁰ by

$$\begin{aligned} \frac{d^2F_{bm}(\omega)}{d\theta d\psi} &= \frac{3\alpha}{4\pi^2} \gamma^2 \frac{\Delta\omega}{\omega} \frac{I}{e} \left(\frac{\omega}{\omega_c}\right)^2 (1 + \gamma^2\psi^2)^2 \left[K_{2/3}^2(\xi) + \frac{\gamma^2\psi^2}{1 + \gamma^2\psi^2} K_{1/3}^2(\xi) \right] \\ &= 1.326 \times 10^{13} \text{ photons/sec/mrad}^2/0.1\% \text{ bandwidth} \\ &\quad \times E^2[\text{GeV}^2] I[A] (1 + \gamma^2\psi^2)^2 \left(\frac{\omega}{\omega_c}\right)^2 \left[K_{2/3}^2(\xi) + \frac{\gamma^2\psi^2}{1 + \gamma^2\psi^2} K_{1/3}^2(\xi) \right] \end{aligned} \quad (5)$$

where θ is the observation angle in the horizontal plane, ψ the observation angle in the vertical plane, α the fine structure constant (1/137), ω the light frequency, I the beam current, and $\xi = (\omega/2\omega_c)(1 + \gamma^2\psi^2)^{3/2}$. The subscripted K s are modified Bessel functions of the second kind. The $K_{2/3}$ term represents light linearly polarized parallel to the electron orbit plane, while the $K_{1/3}$ term represents light linearly polarized perpendicular to the orbit plane.

If one integrates over all vertical angles, then the total intensity is

$$\begin{aligned} \frac{dF_{bm}(\omega)}{d\theta} &= \frac{\sqrt{3}}{2\pi} \alpha\gamma \frac{\Delta\omega}{\omega} \frac{I}{e} \frac{\omega}{\omega_c} \int_{\omega/\omega_c}^{\infty} K_{5/3}(y) dy \\ &= 2.457 \times 10^{13} \text{ photons/sec/mrad}/0.1\% \text{ bandwidth} \times E[\text{GeV}] I[A] \frac{\omega}{\omega_c} \int_{\omega/\omega_c}^{\infty} K_{5/3}(y) dy \end{aligned} \quad (6)$$

The Bessel functions can be computed easily using algorithms of Kostroun:¹²

$$K_\nu(x) = h \left\{ \frac{e^{-x}}{2} + \sum_{r=1}^{\infty} e^{-x \cosh(rh)} \cosh(vrh) \right\} \quad (7)$$

and

$$\int_x^{\infty} K_\nu(\eta) d\eta = h \left\{ \frac{e^{-x}}{2} + \sum_{r=1}^{\infty} e^{-x \cosh(rh)} \frac{\cos h(vrh)}{\cos h(rh)} \right\} \quad (8)$$

for all x and for any fractional order ν , where h is some suitable interval such as 0.5. In evaluating the series, the sum is terminated when the r th term is small, $<10^{-5}$, for example.

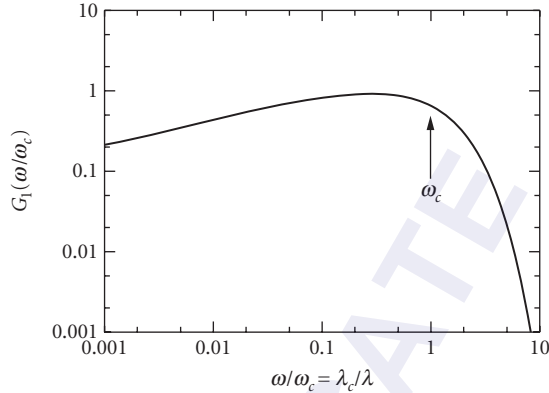


FIGURE 3 Universal synchrotron radiation output curve.

In Fig. 3 we plot the universal function

$$G_1\left(\frac{\omega}{\omega_c}\right) = \frac{\omega}{\omega_c} \int_{\omega/\omega_c}^{\infty} K_{5/3}(y) dy$$

from Eq. (6), so that the photon energy dependence of the flux from a given ring can be calculated readily. It is found that the emission falls off exponentially as $e^{-\lambda_c/\lambda}$ for wavelengths shorter than λ_c , but only as $\lambda^{-1/3}$ at longer wavelengths.

The vertical angular distribution is more complicated. For a given ring and wavelength, there is a characteristic natural opening angle for the emitted light. The opening angle increases with increasing wavelength. If we define ψ as the vertical angle relative to the orbital plane, and if the vertical angular distribution of the emitted flux is assumed to be Gaussian in shape, then the *rms* divergence σ_ψ is defined as $1/\sqrt{2\pi}$ times the ratio of Eqs. (5) and (6) evaluated at $\psi = 0$:

$$\sigma_\psi = \sqrt{\frac{2\pi}{3}} \frac{1}{\gamma} \left(\frac{\omega}{\omega_c}\right)^{-1} \frac{\int_{\omega/\omega_c}^{\infty} K_{5/3}(y) dy}{K_{2/3}^2(\omega/2\omega_c)} \quad (9)$$

In reality, the distribution is not Gaussian, especially in view of the fact that the distribution for the vertically polarized component vanishes in the horizontal plane ($\psi = 0$). However, σ_ψ defined by Eq. (9) is still a simple and useful measure of the angular divergence. The photon energy (ω) dependence of the electron-energy-independent quantity $\gamma\sigma_\psi$ is plotted in Fig. 4. At $\omega = \omega_c$, $\sigma_\psi = 0.647/\gamma$. The asymptotic values of σ_ψ can be obtained from the asymptotic values of the Bessel functions and are

$$\sigma_\psi \approx \frac{1.07}{\gamma} \left(\frac{\omega}{\omega_c}\right)^{-1/3} \quad \omega \ll \omega_c \quad (10)$$

and

$$\sigma_\psi \approx \frac{0.58}{\gamma} \left(\frac{\omega}{\omega_c}\right)^{-1/2} \quad \omega \gg \omega_c \quad (11)$$

In Fig. 5 we show examples of the normalized vertical angular distributions of both parallel and perpendicularly polarized synchrotron radiation for a selection of wavelengths.

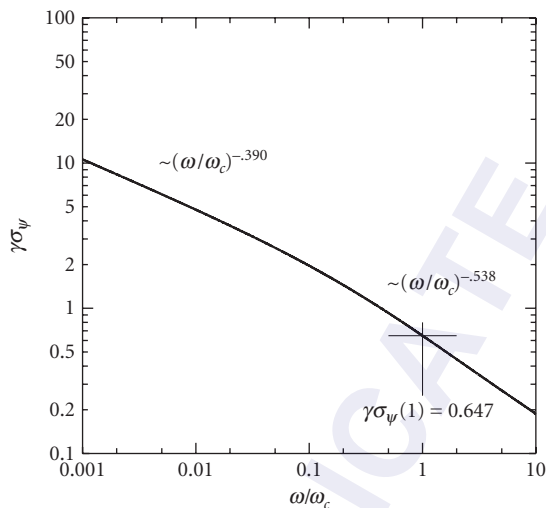


FIGURE 4 Plot of the normalized vertical opening angle $\gamma\sigma_\psi$ for bending magnet radiation.

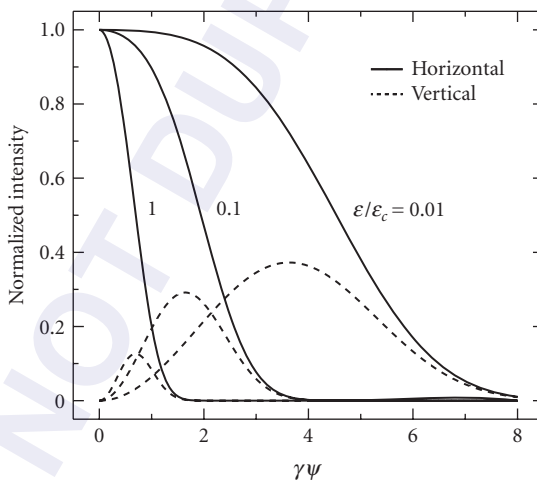


FIGURE 5 Normalized intensities of horizontal and vertical polarization components, as functions of the vertical observation angle for different photon energies.

Circular Polarization and Aperturing for Magnetic Circular Dichroism

Circularly polarized radiation is a valuable tool for the study of electronic, magnetic, and geometric structure of a wide variety of materials. The dichroic response in the soft x-ray spectral region (100 to 1500eV) is especially important because in this energy range almost every element has a strong dipole transition from a sharp core level to its lowest unoccupied state.¹³

The production of bright sources of circularly polarized *soft x-rays* is therefore a topic of keen interest, and is a problem which has seen a multitude of solutions, from special insertion devices

(crossed undulators, helical undulators, elliptically polarized undulators/wigglers) to optical devices (multiple-bounce reflectors/multilayers and quarter-wave plates). However, standard bending magnet synchrotron radiation sources are good sources of elliptically polarized soft x-rays when viewed from either above or below the orbital plane.

As discussed by Chen,¹³ a practical solution involves acceptance of a finite vertical angular range, $\psi_{\text{off}} - \Delta\psi/2 < \psi < \psi_{\text{off}} + \Delta\psi/2$ centered about any vertical offset angle $\psi = \psi_{\text{off}}$ or, equivalently, about $\psi = -\psi_{\text{off}}$. This slice of bending magnet radiation exhibits a circular polarization¹⁴

$$P_c = -\frac{2A_h A_v}{(A_h^2 + A_v^2)} \quad (12)$$

where $A_h = K_{2/3}(\xi)$ and $A_v = \gamma\psi/(1 + \gamma^2\psi^2)^{1/2} K_{1/3}(\xi)$ are proportional to the square-roots of the horizontally and vertically polarized components of bending magnet flux [Eq. (5)], i.e., A_h and A_v are proportional to the horizontal and vertical components of the electric field, respectively. P_c depends on the vertical angle ψ , electron energy γ and, through ξ , the emitted photon energy ω/ω_c . In Fig. 6 we plot values of P_c vs $\gamma\psi$ and ω/ω_c for $\gamma = 1565$ ($E = 0.8\text{GeV}$) and $\rho = 1.91\text{m}$ ($h\nu_{\text{crit}} = 594\text{eV}$).

Magnetic circular dichroism (MCD) measures the normalized difference of the absorption of right circular and left circular light. Assuming no systematic error, the signal to noise ratio in such a measurement defines a figure of merit.

$$\text{MCD figure of merit} = (\text{average circular polarization}) \times (\text{flux fraction})^{1/2} \quad (13)$$

where

$$\text{average circular polarization} = \frac{\int_{\psi_{\text{off}} - \frac{\Delta\psi}{2}}^{\psi_{\text{off}} + \frac{\Delta\psi}{2}} P_c(\psi) \frac{dF}{d\psi} d\psi}{\int_{\psi_{\text{off}} - \frac{\Delta\psi}{2}}^{\psi_{\text{off}} + \frac{\Delta\psi}{2}} \frac{dF}{d\psi} d\psi} \quad (14)$$

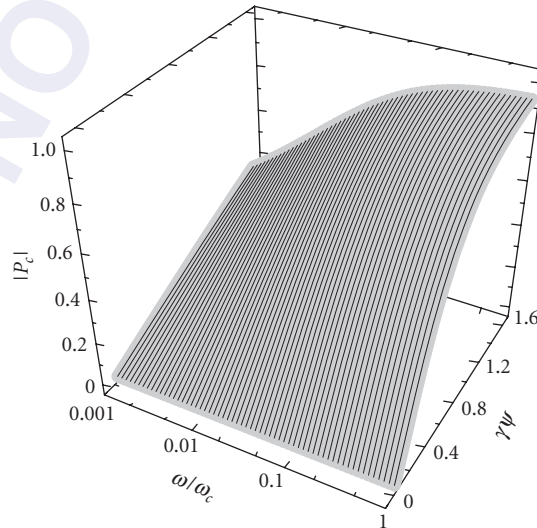


FIGURE 6 P_c versus $\gamma\psi$ versus ω/ω_c .

and the fraction of the total (vertically-integrated) flux emitted into the vertical slice $\psi = \psi_{\text{off}} \pm \Delta\psi/2$ is

$$\text{flux fraction} = \frac{1}{dF_{bm}(\omega)/d\theta} \int_{\psi_{\text{off}} - \Delta\psi/2}^{\psi_{\text{off}} + \Delta\psi/2} \frac{d^2 F_{bm}(\omega)}{d\theta d\psi} d\psi \quad (15)$$

Here $d^2 F_{bm}(\omega)/d\theta d\psi$ is the angular dependence of the bending magnetic flux from Eq. (5) and $dF_{bm}(\omega)/d\theta$ is the vertically integrated flux from Eq. (6). For a 0.8 GeV storage ring (e.g., the VUV ring at the National Synchrotron Light Source (NSLS), Upton, NY USA), the choices of ψ and $\Delta\psi$ that maximize the MCD figure of merit are 0.5 mrad and 0.66 mrad, respectively. This yields a flux fraction ~ 0.3 , a circular polarization ~ 0.65 and a figure of merit ~ 0.35 .

Bending Magnet Power

Integration of $\hbar\omega(d^2 F_{bm}(\omega)/d\theta d\psi)$ from Eq. (5) over all frequencies ω yields the angular distribution of power radiated by a bending magnet:

$$\begin{aligned} \frac{d^2 P_{bm}}{d\theta d\psi} &= \int_0^\infty \hbar\omega \frac{d^2 F_{bm}(\omega)}{d\theta d\psi} d\omega = \frac{I}{e} \frac{\alpha \hbar c \gamma^5}{2\pi\rho} \frac{7}{16} F(\gamma\psi) \\ &= 18.082 \text{ W/mrad}^2 \times \frac{E^5 [\text{GeV}^5] I [\text{A}]}{\rho [\text{m}]} F(\gamma\psi) \end{aligned} \quad (16)$$

which is independent of the horizontal angle θ as required by symmetry, and the vertical angular dependence is contained in the factor

$$F(\gamma\psi) = \frac{1}{(1 + \gamma^2\psi^2)^{5/2}} \left[1 + \frac{5}{7} \frac{\gamma^2\psi^2}{(1 + \gamma^2\psi^2)} \right] \quad (17)$$

The first term in $F(\gamma\psi)$ represents the component of the bending magnet radiation parallel to the orbital plane, while the second represents the perpendicular polarization component. $F(\gamma\psi)$ and its polarization components are plotted versus $\gamma\psi$ in Fig. 7. Note that the area under the F_{parallel} curve is approximately seven times greater than that for $F_{\text{perpendicular}}$.

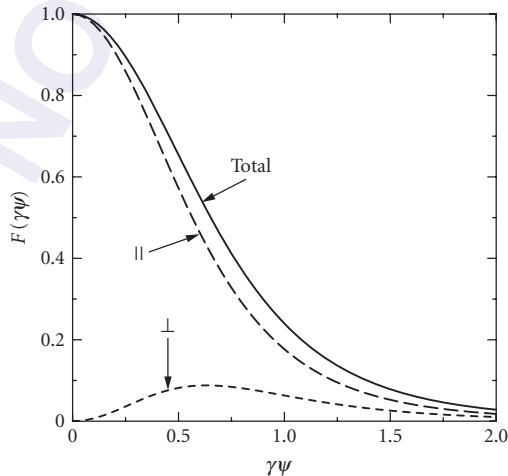


FIGURE 7 Vertical angle dependence of bending magnet power, $F(\gamma\psi)$, versus $\gamma\psi$.

Integrating Eq. (17) over the out-of-orbital-plane (vertical) angle ψ yields the total power radiated per unit in-orbital-plane (horizontal) angle θ :

$$\frac{dP_{bm}}{d\theta} = \frac{I}{e} \frac{hc\alpha\gamma^4}{3\pi\rho} = 14.08 \text{ W/mrad} \times \frac{E^4[\text{GeV}^4]I[\text{A}]}{\rho[\text{m}]} \quad (18)$$

For example, a 1.0-GeV storage ring with 2 m radius bends generates 7.04 W/mrad/Amp of stored current. By contrast, a 2.5 GeV machine with 7-m radius bends generates 78.6 W/mrad/A and a 7 GeV machine with 39 m radius bends generates 867 W/mrad/A.

Bending Magnet Brightness

Thus far we have calculated the emitted flux in photons per second per milliradian² of solid angle. In order to calculate the brightness we need to include the source size. In these calculations we calculate the central (or maximum) brightness, for which we use the natural opening angle to define both the horizontal and vertical angles. Using vertical angles larger than this will not increase the flux as there is no emission. Using larger horizontal angles will increase the flux proportionately as all horizontal angles are filled with light, but owing to the curvature of the electron trajectory, the *average* brightness will actually be less. The brightness expression^{15,16} is

$$B_{bm} = \frac{\frac{d^2 F_{bm}}{(d\theta d\psi)} \Big|_{\psi=0}}{2\pi \sum_x \sum_y} \quad (19)$$

where

$$\sum_x = [\varepsilon_x \beta_x + \eta_x^2 \sigma_E^2 + \sigma_r^2]^{1/2} \quad \text{and} \quad \sum_y = \left[\varepsilon_y \beta_y + \sigma_r^2 + \frac{\varepsilon_y^2 + \varepsilon_y \gamma_y \sigma_r^2}{\sigma_\psi^2} \right]^{1/2} \quad (20)$$

ε_x and ε_y are the electron beam emittances in the horizontal and vertical directions respectively, β_x and β_y are the electron beam beta functions in the horizontal and vertical planes, η_x is the dispersion function in the horizontal plane, and σ_E is the rms value of the relative energy spread. All the electron beam parameters are properties of a particular storage ring. The diffraction-limited source size is $\sigma_r = \lambda/4\pi\sigma_\psi$. The effective source sizes (\sum_x and \sum_y) are photon energy dependent via the natural opening angle σ_ψ and the diffraction-limited source size σ_r .

55.3 INSERTION DEVICES (UNDULATORS AND WIGGLERS)

General

Insertion devices are periodic magnetic structures installed in straight sections of storage rings, as illustrated in Fig. 8, in which the vertical magnetic field varies approximately sinusoidally along the axis of the undulator. The resulting motion of the electrons is also approximately sinusoidal, but in the horizontal plane. One can understand the nature of the spectra emitted from these devices by again studying the electric field as a function of time, and this is shown in Fig. 9. This shows that the electric field and hence its Fourier transform, the spectrum, depend critically on the magnitude of the beam deflection in the device. At one extreme, when the magnetic fields are high, as in Fig. 9a, the deflection is large and the electric field is a series of pulses similar to those obtained from a

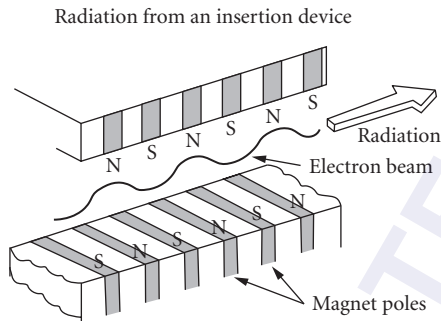


FIGURE 8 Schematic of an insertion device.

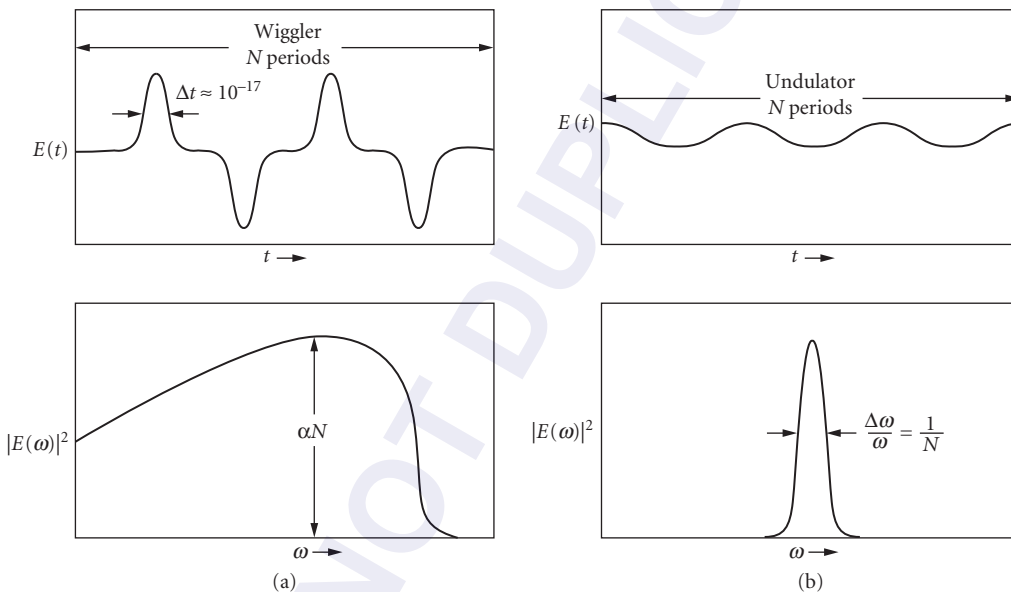


FIGURE 9 Conceptual representation of the electric fields emitted as a function of time by an electron in (a) a wiggler and (b) an undulator, with the corresponding spectra.

dipole. Such a device is termed a “wiggler.” The Fourier transform for the wiggler is N times that of a single dipole. At the other extreme, as in Fig. 9b, the deflection of the electron beam is such that the electric field as a function of time is sinusoidal, and the Fourier transform is then a single peak with a width proportional to the inverse of the length of the wavetrain L^* according to $\lambda^2/\Delta\lambda = L^*$. L^* is obtained by dividing the real length of the device L by γ^2 because of relativistic effects. Thus for a meter long device emitting at a wavelength $\lambda = 10$ nm in a machine of energy 0.5 GeV ($\gamma \sim 1000$), corresponding to, for example, 1 cm period length and magnetic field strength $B = 1.5$ T, we get $\lambda^2/\Delta\lambda = 10^{-6}$ meters, and $\lambda/\Delta\lambda = 100$. Note that $\Delta\lambda/\lambda \sim 1/N$ as expected. Interference occurs in an undulator since the electric field from one part of the electron path is added coherently to that from adjacent parts.

Formal Treatment

We assume that the motion of an electron in an insertion device is sinusoidal, and that we have a magnetic field in the vertical (y) direction varying periodically along the z direction, with

$$B_y = -B_0 \sin(2\pi z/\lambda_u) \quad 0 \leq z \leq N\lambda_u \quad (21)$$

where B_0 is the peak magnetic field, λ_u is the period length, and N the number of periods. By integrating the equation of motion, the electron transverse velocity $c\beta_x$ is found to be

$$\beta_x = \frac{K}{\gamma} \cos(2\pi z/\lambda_u) \quad (22)$$

where

$$K = \frac{eB_0\lambda_u}{2\pi mc} = 0.934\lambda_u[\text{cm}]B_0[\text{T}] \quad (23)$$

is a dimensionless parameter which is proportional to the deflection of the electron beam. The maximum slope of the electron trajectory is $\delta = (K/\gamma)$. In terms of δ , we define an undulator as a device in which $\delta \leq \gamma^{-1}$, which corresponds to $K \leq 1$. When K is large, the device is called a wiggler. In most insertion devices the field can be changed either electromagnetically or mechanically, and in some cases K can vary between the two extremes of undulator and wiggler operation.

Wigglers

For the wiggler, the flux distribution is given by $2N$ (where N is the number of magnetic periods) times the appropriate bending magnet formulae in Eqs. (5) and (6). However, ρ or B must be taken at the point in the path of the electron which is tangent to the direction of observation. For a horizontal angle θ ,

$$\varepsilon_c(\theta) = \varepsilon_{c\text{max}} \sqrt{1 - (\theta/\delta)^2} \quad (24)$$

where

$$\varepsilon_{c\text{max}}[\text{keV}] = 0.665E^2[\text{GeV}^2]B_0[\text{T}] \quad (25)$$

from Eq. (4). Integration over θ , which is usually performed numerically, gives the wiggler flux.

The calculation of the brightness of wigglers needs to take into account the depth-of-field effects, i.e., the contribution to the apparent source size from different poles. The expression for the brightness of wigglers is

$$B_w = \frac{d^2 F_w}{d\theta d\psi} \sum_{\pm} \sum_{n=-\frac{1}{2}N}^{\frac{1}{2}N} \frac{1}{2\pi} \times \frac{\exp\left[-\frac{1}{2}\left(\frac{x_o^2}{\sigma_x^2 + z_{n\pm}^2 \sigma_x'^2}\right)\right]}{\left[(\sigma_x^2 + z_{n\pm}^2 \sigma_x'^2)\left(\frac{\varepsilon_y^2}{\sigma_\psi^2} + \sigma_y^2 + z_{n\pm}^2 \sigma_y'^2\right)\right]^{1/2}} \quad (26)$$

where $z_{n\pm} = \lambda_w[n \pm (1/4)]$, λ_w is the wiggler period, and σ_ψ is identical to Eq. (9), but evaluated, in the wiggler case, as the instantaneous radius at the tangent to the straight-ahead ($\theta = \psi = 0$) direction (i.e., minimum ρ , maximum ε_c), $\sigma_x = \sqrt{\varepsilon_x \beta_x}$ and $\sigma_y = \sqrt{\varepsilon_y \beta_y}$ are the rms transverse beam sizes,

while $\sigma'_x = \sqrt{\varepsilon_x/\beta_x}$ and $\sigma'_y = \sqrt{\varepsilon_y/\beta_y}$ are the angular divergences of the electron beam in the horizontal and vertical directions respectively. The exponential factor in Eq. (26) arises because wigglers have two source points separated by $2x_o$, where

$$x_o = \frac{K \lambda_w}{\gamma 2\pi} \quad (27)$$

The summations in Eq. (26) must be performed for each photon energy because σ_ψ is photon-energy dependent.

Undulators

The interference which occurs in an undulator, i.e., when K is moderate ($K \leq 1$), produces sharp peaks in the forward direction at a fundamental ($n = 1$) and all odd harmonics ($n = 3, 5, 7, \dots$) as shown for a zero emittance ($\varepsilon = 0$) electron beam in Fig. 10a (dotted line). In the $\varepsilon = 0$ case, the even harmonics ($n = 2, 4, 6, \dots$) peak off-axis and do not appear in the forward direction. For real ($\varepsilon \neq 0$) electron beams, the spectral shape, angular distribution, and peak brightness are strongly dependent on the emittance and energy spread of the electron beam as well as the period and magnitude of the insertion device field.

In general, the effect of electron beam emittance is to cause all harmonics to appear in the forward direction (solid line in Fig. 10a). The effect of angle integration on the spectrum in Fig. 10a is shown in Fig. 10b, a spectrum which is independent of electron beam emittance except for the presence of “noise” in the zero emittance case. The effect of electron beam emittance on the angular distribution of the fundamental, second, and third harmonics of this device is shown in Fig. 10c, which also nicely demonstrates the dependence on harmonic number.

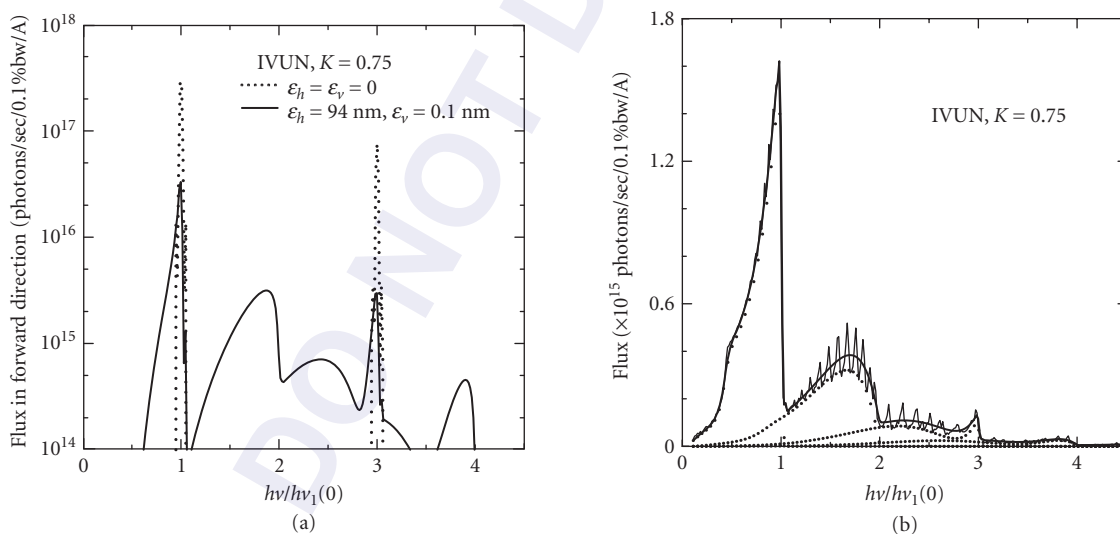


FIGURE 10 Spectral output and angular distribution of the emission from the early-2000s vintage NSLS In-Vacuum Undulator (IVUN) for $K = 0.75$. (a) Spectral output in the forward direction, with (solid line) and without (dotted line) the effect of electron beam emittance; (b) angle-integrated spectral output with (solid line) and without (faint solid line) the effect of electron beam emittance, and the decomposition into harmonics ($n = 1, 2, 3, 4$) (dotted lines); and (c) angular distribution of the first three harmonics ($n = 1, 2, 3$), with and without the effect of electron beam emittance. Emittance values used: 94-nm horizontal, 0.1-nm vertical.

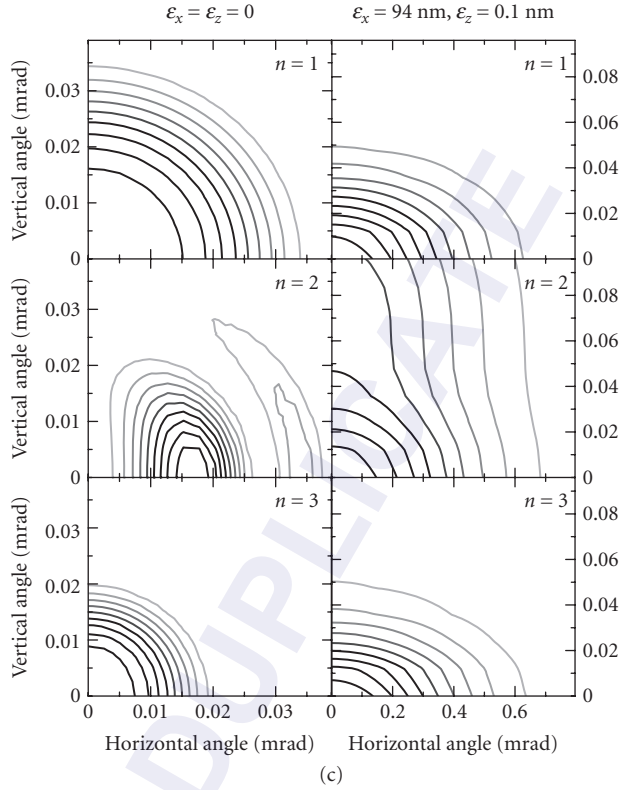


FIGURE 10 (Continued)

The peak wavelengths of the emitted radiation λ_n are given by

$$\lambda_n = \frac{\lambda_u}{2n\gamma^2} \left(1 + \frac{K^2}{2} + \gamma^2 \theta^2 \right) \quad n=1, 3, 5, \dots \quad (28)$$

where λ_u is the undulator period length. They soften as the square of the deviation angle θ away from the forward direction.

Of main interest is the intense central cone of radiation. An approximate formula for flux integrated over the central cone (for the odd harmonics) is

$$\begin{aligned} F_u(K, \omega) &= \pi \alpha N \frac{\Delta \omega}{\omega} \frac{I}{e} Q_n(K) \quad n=1, 3, 5, \dots \\ &= 1.431 \times 10^{14} \text{ photons/sec/0.1\% bandwidth} \times I[A] N Q_n(K) \end{aligned} \quad (29)$$

where

$$Q_n(K) = \left(1 + \frac{K^2}{2} \right) \frac{F_n(K)}{n} \quad n=1, 3, 5, \dots \quad (30)$$

and

$$F_n(K) = \frac{K^2 n^2}{(1 + K^2/2)^2} \left\{ J_{(n-1)/2} \left[\frac{nK^2}{4 \left(1 + \frac{1}{2}K^2\right)} \right] - J_{(n+1)/2} \left[\frac{nK^2}{4 \left(1 + \frac{1}{2}K^2\right)} \right] \right\}^2 \quad (31)$$

Here $J_{(n \pm 1)/2}$, $n = 1, 3, 5, \dots$ are the integer Bessel functions of the first kind: J_0, J_1, J_2, \dots

To calculate the undulator flux angular distribution and spectral output into arbitrary solid angle, one can use freely available codes such as Urgent¹⁷ (R. P. Walker and B. Diviacco). To include magnetic field errors (e.g., measured values), use Ur¹⁸ (R. J. Dejus and A. Luccio), SRW¹⁹ (O. Chubar and P. Elleaume), or Spectra²⁰ (T. Tanaka and H. Kitamura).

The brightness of an undulator B_u is approximated by dividing the central cone flux by the effective angular divergence, $\Sigma'_x (\Sigma'_y)$, and by the effective source size, $\Sigma_x (\Sigma_y)$, in the horizontal (vertical) directions. These are given by convolution of the Gaussian distributions of the electron beam and the diffraction-limited photon beam, in both angle and space:

$$\Sigma_{x'} = \sqrt{\sigma_x^2 + \sigma_r^2} \quad \Sigma_{y'} = \sqrt{\sigma_y^2 + \sigma_r^2} \quad (32)$$

$$\Sigma_x = \sqrt{\sigma_x^2 + \sigma_r^2} \quad \Sigma_y = \sqrt{\sigma_y^2 + \sigma_r^2} \quad (33)$$

Thus, B_u is given by

$$B_u = \frac{F_u}{(2\pi)^2 \Sigma_x \Sigma_y \Sigma_{x'} \Sigma_{y'}} \quad (34)$$

The diffraction-limited emittance of a photon beam is the minimum value in the inequality

$$\varepsilon = \sigma_r \sigma_{r'} \geq \frac{\hbar}{2} = \frac{\lambda}{4\pi} \quad (35)$$

where ε is the photon emittance and λ is the wavelength, in direct analogy to the Heisenberg uncertainty principle in nonrelativistic quantum mechanics. The space versus angle separation of this minimum emittance is energy and harmonic dependent.²¹ For the exact harmonic frequency in the forward direction, given by Eq. (28) with $\theta = 0$, there appears to be consensus that σ_r and $\sigma_{r'}$ are given by

$$\sigma_{r'} = \sqrt{\frac{\lambda}{2L}} \quad \text{and} \quad \sigma_r = \frac{\sqrt{2\lambda L}}{4\pi} \quad (36)$$

On the other hand, at the peak of the angle-integrated undulator spectrum, which lies a factor of $[1 - (1/nN)]$ below the exact harmonic energy, σ_r and $\sigma_{r'}$ are given by

$$\sigma_{r'} = \sqrt{\frac{\lambda}{L}} \quad \text{and} \quad \sigma_r = \frac{\sqrt{\lambda L}}{4\pi} \quad (37)$$

It is clear from Eqs. (32) and (33) that the choice of expression for σ_r and $\sigma_{r'}$ can have a non-negligible effect on the undulator brightness value especially for small beam size and opening angle. Lacking a specific functional form for σ_r and $\sigma_{r'}$ as a function of photon energy, we generally use Eq. (36) in evaluating the expression for undulator spectral brightness from Eq. (34).

Insertion Device Power

The Schwinger⁵ formula for the distribution of radiated power from an electron in a sinusoidal trajectory, which applies with reasonable approximation to undulators and, to a lesser extent, wigglers, reduces²² to

$$\frac{d^2P}{d\theta d\psi} = P_{\text{total}} \frac{21\gamma^2}{16\pi K} G(K) f_K(\gamma\theta, \gamma\psi) \quad (38)$$

where the total (angle-integrated) radiated power is

$$P_{\text{total}} = \frac{N}{6} \frac{Z_0 I 2\pi e c}{\lambda_u} \gamma^2 K^2 = 633.0 \text{ W} \times E^2 [\text{GeV}^2] B_0^2 [T^2] L [m] I [A] \quad (39)$$

where N is the number of undulator or wiggler periods, Z_0 is the vacuum impedance (377Ω), I is the storage ring current, e is the electronic charge, c is the speed of light, $L = N\lambda_u$ is the length of the insertion device,

$$G(K) = \frac{K}{(1 + K^2)^{7/2}} \left(K^6 + \frac{24}{7} K^4 + 4K^2 + \frac{16}{7} \right) \quad (40)$$

and

$$f_K(\gamma\theta, \gamma\psi) = \frac{16}{7\pi} \frac{K}{G(K)} \int_{-\pi}^{\pi} d\alpha \left(\frac{1}{D^3} - \frac{4(\gamma\theta - K \cos \alpha)^2}{D^5} \right) \sin^2 \alpha \quad (41)$$

where

$$D = 1 + \gamma^2 \psi^2 + (\gamma\theta - K \cos \alpha)^2 \quad (42)$$

The integral in the expression for f_K is best evaluated numerically.

For $K > 1$, which includes all wigglers and much of the useful range of undulators, an approximate formula for the angle dependence of the radiated power is

$$f_K(\gamma\theta, \gamma\psi) = \sqrt{1 - \left(\frac{\gamma\theta}{K} \right)^2} F(\gamma\psi) \quad (43)$$

where $F(\gamma\psi)$ is the bending magnet formula from Eq. (17). This form clearly indicates the strong weakening of insertion device power as θ increases, vanishing at $\theta = \pm K/\gamma$.

Since $f_K(0,0)$ is normalized to unity, the radiated power density in the forward direction (i.e., along the undulator axis) is

$$\frac{d^2P}{d\theta d\psi} (\theta = 0, \psi = 0) = P_{\text{total}} \frac{21\gamma^2}{16\pi K} G(K) = 10.84 \text{ W/mrad}^2 \times B_0 [T] E^4 [\text{GeV}^4] I [A] N G(K) \quad (44)$$

Polarization of Undulators and Wigglers

The polarization properties of the light emitted by wigglers are similar to those of dipoles. For both sources the radiation is elliptically polarized when observed at some angle away from the orbital plane as given by Eq. (5). For radiation from planar undulators, however, the polarization is always linear. The polarization direction, which is in the horizontal plane when observed from that plane,

rotates in a complicated way at other directions of observation. A comprehensive analysis of the polarization from undulators has been carried out by Kitamura.²³ The linear polarization of the undulator radiation is due to the symmetry of the electron trajectory within each period. The polarization can in fact be controlled by a deliberate breaking of this symmetry. Circularly polarized radiation can be produced by a helical undulator, in which the series of dipole magnets is arranged such that each period is rotated by a fixed angle with respect to the previous one. To generate variable polarization, one can use a pair of planar undulators oriented at right angles to each other. The amplitude of the radiation from these so-called crossed undulators is a linear superposition of two parts, one linearly polarized along the x direction and another linearly polarized along the y direction, x and y being orthogonal to the electron beam direction. By varying the relative phase of the two amplitudes by means of a variable-field magnet between the undulators, it is possible to modulate the polarization in an arbitrary way. The polarization can be linear and switched between two mutually perpendicular directions, or it can be switched between left and right circular polarization. For this device to work, it is necessary to use a monochromator with a sufficiently small band-pass, so that the wave trains from the two undulators are stretched and overlap. Also the angular divergence of the electron beam should be sufficiently small or the fluctuation in relative phase will limit the achievable degree of polarization. A planar undulator whose pole boundaries are tilted away from a right angle with respect to the axial direction can be used as a helical undulator if the electron trajectory lies a certain distance above or below the midplane of the device.

Transverse Spatial Coherence

As shown by Kim²⁴ and utilized in the brightness formulae given earlier, in wave optics the phase-space area of a radiation beam is given by the ratio of flux (F_0) to brightness (B_0). A diffraction-limited photon beam (no electron size or angular divergence contribution) occupies the minimum possible phase-space area. From Eqs. (32) to (37) this area is

$$(2\pi\sigma_r\sigma_r')^2 = (2\pi\varepsilon)^2 = \left(\frac{\lambda}{2}\right)^2 \quad (45)$$

Thus, the phase space occupied by a single Gaussian mode radiation beam is $(\lambda/2)^2$, and such a beam is referred to as completely transversely coherent. It then follows that the transversely coherent flux of a radiation beam is

$$F_{\text{coherent}} = \left(\frac{\lambda}{2}\right)^2 B_0 \quad (46)$$

and the degree of transverse spatial coherence is

$$\frac{F_{\text{coherent}}}{F_0} = \left(\frac{\lambda}{2}\right)^2 \frac{B_0}{F_0} \quad (47)$$

Conversely, the number of Gaussian modes occupied by a beam is

$$\frac{F_0}{F_{\text{coherent}}} = \frac{F_0}{B_0(\lambda/2)^2} \quad (48)$$

Transverse spatial coherence is the quantity which determines the throughput of phase-sensitive devices such as Fresnel zone plates used for x-ray microscopy.

55.4 COHERENCE OF SYNCHROTRON RADIATION EMISSION IN THE LONG WAVELENGTH LIMIT

We now discuss coherent effects which depend upon the phase coherence between electric fields emitted by different electrons within a bunch. These effects naturally become stronger as the emission wavelength increases to values greater than the bunch length, the so-called long wavelength limit, which is generally in the far-infrared range or beyond for standard storage ring parameters. In the evolution from 3rd to 4th generation sources, the bunch length decreases, thereby pushing the coherent emission to shorter wavelengths, into the UV, VUV, and beyond. One way to generate shorter bunches is with single-pass devices such as linacs, wherein the electrons can be closer packed since relaxation processes which occur in storage rings may not have had time to develop.

In this section, we confine our discussion to coherent synchrotron radiation emission in the long wavelength limit. We describe two types of coherence: (i) longitudinal, or temporal (sometimes called multiparticle or super-radiant) coherence; and (ii) transverse or lateral coherence.

Temporal (Longitudinal/Multiparticle/Super-Radiant) Coherence

Under normal circumstances in a storage ring, the bunch length is considerably longer than the emitted wavelength, and since the electrons are randomly distributed, there is no phase correlation between the emitted electric fields from different electrons. For N_e particles emitting, then, the amplitude of the electric field is generated from the statistical noise, which is proportional to $N_e^{1/2}$. This implies that the power, which goes like to E^2 , is proportional to N_e .

In situations in which the wavelength is much longer than the bunch length, the phase differences between the electric fields of the electrons are small compared to the wavelength, and the field is N_e times that of a single electron. In this regime, the intensity is N_e^2 times that of a single electron, or N_e times greater than the incoherent process. This effect is very large because N_e is generally quite large: the number of electrons per bunch can easily be in the nanoCoulomb range, i.e., containing 10^{10} to 10^{11} electrons.

Long-wavelength coherent emission from bunches of relativistic charged particles was first described theoretically by Nodvick and Saxon²⁵ in 1954. For the situation of multiple electrons, the more general expression for the flux emitted by an electron *bunch* as a function of frequency (ω) and solid angle (Ω) is derived by extending the expressions derived earlier for a single electron [Eq. (5)], to a system of N_e electrons, thus:

$$\frac{d^2 F_{bm}(\omega)}{d\theta d\psi} = 1.326 \times 10^{13} \text{ photons/sec/mrad}^2 / 0.1\% \text{ bandwidth} \times [N_e [1 - f(\omega)] + N_e^2 f(\omega)] \\ \times E^2 [\text{GeV}^2] I [A] (1 + \gamma^2 \psi^2)^2 \left(\frac{\omega}{\omega_c} \right)^2 \left[K_{2/3}^2(\xi) + \frac{\gamma^2 \psi^2}{1 + \gamma^2 \psi^2} K_{1/3}^2(\xi) \right] \quad (49)$$

where the term $N_e^2 f(\omega)$ represents the coherent enhancement and includes the form factor $f(\omega)$, which is the Fourier transform of the normalized longitudinal particle distribution within the bunch, i.e.,

$$f(\omega) = \left| \int_{-\infty}^{\infty} e^{i\omega \tilde{z}/c} S(z) dz \right|^2 \quad (50)$$

where $S(z)$ is the distribution function for particles in the bunch, measured relative to the bunch center.

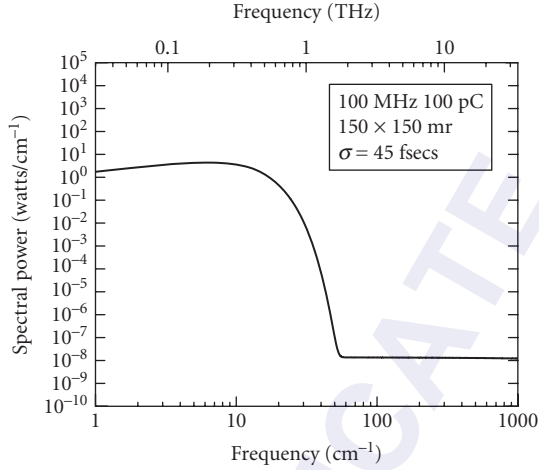


FIGURE 11 Long-wavelength coherent enhancement of synchrotron radiation spectra power for short electron bunches. Electron beam parameters used: 100-pC charge, 100-MHz repetition rate, 45-fs ($1 - \sigma$) bunch length. Emission angle range considered: 150 mrad (horiz.), 150 mrad (vert.).

If we assume that the electron bunch has a longitudinal Gaussian particle distribution of width σ_z , the form factor will be

$$f(\omega) = e^{-\left(\frac{\omega\sigma_z}{c}\right)^2} = e^{-4\pi^2\sigma_z^2\left(\frac{1}{\lambda}\right)^2} \quad (51)$$

where λ is the wavelength of the light at frequency ω .

The long-wavelength coherent enhancement is shown graphically for a chosen set of beam parameters in Fig. 11.

Transverse, or Lateral, Coherence

Transverse coherence was defined and described earlier in this chapter in the context of emission from undulators. In the long wavelength limit, even the emission from bending magnets becomes transversely coherent. In this section, we derive the interesting result that, in the long wavelength limit, the flux emitted into the natural opening angle (transverse coherent limit) is constant for all wavelengths, independent of storage ring parameters such as electron energy and bend radius.

1. Natural opening angle for synchrotron radiation for $\omega < \omega_c$:

$$\theta_{\text{rms}} = 0.8282 \left(\frac{\lambda(m)}{2\pi\rho(m)} \right)^{1/3} \text{ radians} \quad (52)$$

2. For $\omega < \omega_c$, the universal function

$$G_1\left(\frac{\omega}{\omega_c}\right) = \frac{\omega}{\omega_c} \int_{\omega/\omega_c}^{\infty} K_{5/3}(y) dy$$

from Eq. (6) and Fig. 3 can be approximated by^{8,11}

$$G_1\left(\frac{\omega}{\omega_c}\right) = \frac{1.33 \times 8 \times \pi}{9 \times \sqrt{3}} \left(\frac{\omega}{\omega_c}\right)^{1/3} = 2.144 \left(\frac{\omega}{\omega_c}\right)^{1/3}$$

3. Therefore, the formula for angle integrated flux (F) for $\omega < \omega_c$ into horizontal collection angle θ is

$$F = 5.27 \times 10^{16} \text{ photons/sec/0.1\% bandwidth} \times \theta[\text{rad}] E[\text{GeV}] I[\text{A}] \left(\frac{\omega}{\omega_c}\right)^{1/3} \quad (53)$$

4. Critical wavelength λ_c is given by

$$\lambda_c [m] = 5.59 \times 10^{-10} \frac{\rho [m]}{E^3 [\text{GeV}^3]} = \frac{2\pi c [\text{m/sec}]}{E^3 [\text{GeV}^3]} \quad (54)$$

so that

$$\omega_c [\text{sec}^{-1}] = \frac{2\pi c [\text{m/sec}] E^3 [\text{GeV}^3]}{5.59 \times 10^{-10} \rho [m]} \quad (55)$$

5. Substituting Eqs. (52) and (55) into Eq. (53) and using $4 \times \theta_{\text{rms}}$ for the horizontal collection angle θ , we obtain:

$$F = 5.374 \times 10^{16} \text{ photons/sec/0.1\% bandwidth} \times E[\text{GeV}] I[\text{A}] \times \left(\frac{5.59 \times 10^{-10} \omega [\text{sec}^{-1}] \rho [m]}{2\pi c [\text{m/sec}] E^3 [\text{GeV}^3]}\right)^{1/3} \times 4 \times 0.8282 \left(\frac{2\pi c [\text{m/sec}]}{2\pi \rho [m] \omega [\text{sec}^{-1}]}\right)^{1/3} \quad (56)$$

or

$$F = 7.94 \times 10^{13} \text{ photons/sec/0.1\% bandwidth} \times I[\text{A}] \quad (57)$$

Figure 12 shows a flux plot based on an exact calculation for a synchrotron storage ring source with 3-GeV electron energy and 5-meter bend radius. This plot agrees with Eq. (57) from $\sim 2 \times 10^{-6} \omega_c$

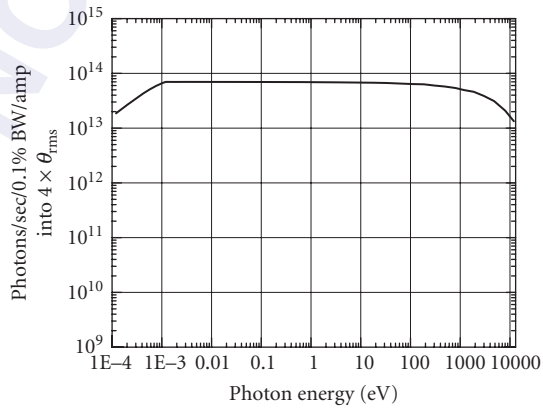


FIGURE 12 Synchrotron radiation spectral flux emitted into $4\theta_{\text{rms}}$ (horiz. \times vert.) for a storage ring with 3-GeV electron energy and 5-meter bend radius.

to $\sim\omega_c$. The curve turns downward at the low-energy end owing to the finite size of the ring vacuum chamber. At the upper energy end it also turns downward owing to a failure of the approximation in Eq. (53) at and above ω_c .

For consistency as well as interest, another route to the result in Eq. (57) is via the undulator flux formula given by Eq. (29), with $K = 1$, $n = 1$, and $N = 1$. In this case, intended to approximate a bending magnet, $F_n(K) = 0.37$ one obtains, from Eq. (29), a flux value of $F = 7.99 \times 10^{13}$ photons/sec/0.1% bandwidth $\times I [A]$, in remarkable agreement with Eq. (57).

55.5 CONCLUSION

We have attempted to compile the formulae needed to calculate the flux, brightness, polarization (linear and circular), and power produced by the three standard storage ring synchrotron radiation sources: bending magnets, wigglers and undulators. Where necessary, these formulae have contained reference to the emittance (ϵ) of the electron beam, as well as to the electron beam size (σ) and its divergence (σ'). For all three types of sources, the source phase space area, i.e., the spatial and angular extent of the effective (real) source, is a convolution of its electron and photon components. For a more detailed description of these properties, see Ref. 26 and references therein.

55.6 REFERENCES

1. A. Liénard, "Champ électrique et magnétique, produit par une charge électrique concentrée en un point et animée d'un mouvement quelconque," *L'Eclairage Electrique* **16**:5 (1898).
2. J. P. Blewett, "Synchrotron Radiation—1873–1947," *Nucl. Instrum. Methods* **A266**:1 (1988).
3. P. C. Hartman, "Introductory Remarks," *Nucl. Instrum. Methods* **195**:1 (1982).
4. D. H. Tomboulion and P. L. Hartman, "Spectral and Angular Distribution of Ultraviolet Synchrotron Radiation from the 300-MeV Cornell Synchrotron," *Phys. Rev.* **102**:1423 (1956).
5. J. Schwinger, "On the Classical Radiation of Accelerated Electrons," *Phys. Rev.* **75**:1912 (1949).
6. J. D. Jackson, *Classical Electrodynamics* (Wiley, New York) 1975.
7. H. Winick, *Synchrotron Radiation Research* (Plenum Press, New York) Chapter 2, 1980.
8. A. Hofmann, *The Physics of Synchrotron Radiation*, (Cambridge University Press, Cambridge, England) 2004, ISBN 0 521 30826 7; A. Hofmann, *Phys. Rep.* **64**:253 (1980).
9. S. Krinsky, M. L. Perlman, and R. E. Watson, *Handbook of Synchrotron Radiation*, E. E. Koch (ed.) (North-Holland, Amsterdam) Chapter 2, 1983.
10. K. J. Kim, "Physics of Particle Accelerators," *AIP Proceedings* **184**:565 (1989).
11. A. A. Sokolov and I. M. Ternov, "Synchrotron Radiation" (Cambridge University Press, Cambridge, England) 1968, ISBN 0 521 30826 7.
12. V. O. Kostroun, "Simple Numerical Evaluation of Modified Bessel Functions $K_\nu(x)$ of Fractional Order and the Integral $\int_{x_0}^{\infty} K_\nu(\eta)\delta\eta$," *Nucl. Instrum. Methods* **172**:371 (1980).
13. C. T. Chen, "Raytracing, Chopper and Guideline for Double-Headed Dragon Monochromators," *Rev. of Scientific Instr.* **63**:1229 (1992).
14. M. Born and E. Wolf, *Principles of Optics*, (Pergamon Press, London) 1964.
15. Lawrence Berkeley Laboratory Publication 643 Rev. 2, "X-Ray Data Booklet," Berkeley, Calif. (1989).
16. S. L. Hulbert and J. M. Weber, "Flux and Brightness Calculations for Various Synchrotron Radiation Sources," *Nucl. Instrum. Methods* **A319**:25 (1992).
17. R. P. Walker and B. Diviacco, "URGENT—A Computer Program for Calculating Undulator Radiation Spectral, Angular, Polarization and Power Density Properties," *Rev. Scientific Instr.* **63**:392 (1992).
18. R. P. Dejus and A. Luccio, "Program UR: General Purpose Code for Synchrotron Radiation Calculations," *Nucl. Instrum. Methods* **A347**:61 (1994).

19. O. Chubar and P. Elleaume, "Accurate and Efficient Computation of Synchrotron Radiation in the Near Field Region," *Proc. of the EPAC98 Conf.* **22–26**:1177–1179 (1998).
20. T. Tanaka and H. Kitamura, "SPECTRA: A Synchrotron Radiation Calculation Code," *J. Synch. Rad.* **8**:1221 (2001).
21. M. R. Howells and B. M. Kincaid, *LBL Report #34751* (1993).
22. K. J. Kim, "Angular Distribution of Undulator Power for an Arbitrary Deflection Parameter K," *Nucl. Instrum. Methods* **A246**:67 (1986).
23. H. Kitamura, "Polarization of Undulator Radiation," *Jpn. J. Appl. Phys.* **19**:L185 (1980).
24. K. J. Kim, "Brightness, Coherence and Propagation Characteristics of Synchrotron Radiation," *Nucl. Instrum. Methods* **A246**:71 (1986). See also earlier work of A. M. Kondratenko and A. N. Skrinsky, "Use of Radiation of Electron Storage Rings in X-Ray Holography of Objects," *Opt. Spectroscopy* **42**:189 (1977) and F. Zernike, "The Concept of Degree of Coherence and Its Application to Optical Problems," *Physica V* **785** (1938).
25. J. S. Nodvick and D. S. Saxon, "Suppression of Coherent Radiation by Electrons in a Synchrotron," *Phys. Rev.* **96**:180 (1954).
26. S. L. Hulbert and G. P. Williams, "Synchrotron Radiation Sources," in *Vacuum Ultraviolet Spectroscopy I*, J. A. R. Samson and D. L. Ederer (eds.), *Experimental Methods in the Physical Sciences* **31**:1 (Academic Press, San Diego, Calif.) 1998.

This page intentionally left blank.

DO NOT DUPLICATE

Alan Michette

*King's College London
United Kingdom*

56.1 INTRODUCTION

Laser plasmas are produced by focusing a pulsed laser beam, e.g., Nd:YAG at 1.064 μm , possibly frequency multiplied, or KrF excimer at 249 nm, onto a target, typically a tape,¹ liquid droplet or jet² or a gas-jet.³ The required irradiance, that is the focused intensity per unit area per unit time, is in the range $\sim 10^{17}$ to 10^{19} Wm^{-2} , which heats the target material to $\sim 10^6$ K and thus ionizes it to produce the plasma. The need for high irradiance means that high beam energies, small focused spot sizes and short pulse lengths must be used. For repetitive systems¹ pulse energies in the range ~ 10 mJ to 1 J are used with focal spot sizes of ~ 10 μm . For single pulse systems, using much higher pulse energies, focal spot sizes can be much larger. Normally pulse lengths are in the range of a few picoseconds to several nanoseconds, but lasers with much shorter pulse lengths, in the tens of femtoseconds range, have also been used.⁴

For picosecond–nanosecond systems the emission is generated from the ionized material. The spectral characteristics depend mainly on the target material, with the proviso that the irradiance must be sufficient to produce the ionic state required to give a particular spectral feature. The use of tape targets, compared to liquids or gases, allows a wider range of materials to be used, but the effects of increased particulate debris emission must be alleviated by a low-pressure buffer gas. For light materials it is common to almost strip the atoms and emission is then from hydrogen- and helium-like ions. The spectrum largely consists of characteristic line emission, with small contributions from bremsstrahlung and recombination radiation. The lines typically have bandwidths $\Delta\lambda/\lambda \sim 10^{-4}$.

For a repetitive source operated at ~ 100 Hz, the brilliance in a particular line, e.g., the H-like carbon Lyman- α line at a wavelength of 3.37 nm, can be comparable to that of a second-generation synchrotron, albeit only at this wavelength. Use of a heavier Z target, such as copper or gold, means that the ions are less fully stripped and the emission from many closely spaced ionic energy levels merges into a quasi-continuum. Such targets also give rise to more bremsstrahlung. Although the overall emission is higher than that from light targets it is considerably lower than that of a synchrotron, and the peak emission can be much less than that from light materials.

56.2 CHARACTERISTIC RADIATION

Characteristic radiation arises as a direct consequence of the quantum nature of the plasma ions. Electrons bound in ionic systems can only occupy discrete energy levels, with only one electron in each quantum state—the Pauli Exclusion Principle. In the ground state all the lower energy levels will be filled. If, by some process, an electron is removed from a filled level, either by promoting it to a higher energy level (excitation) or by completely removing it from the ion (ionization), then the empty level may be filled by an electron from a higher occupied state, which loses energy by emitting radiation. Because of the discreteness of the energy levels the transition between two defined states will always give a photon of the same energy, which is characteristic of the ion involved since the energies of the bound states depend on the nuclear charge and on the screening effects of other electrons.

Energies of Spectral Lines

For an ion with just one electron—a hydrogenic ion—the energies of the electron states are

$$E_n = -\frac{1}{2n^2} \frac{\mu}{\hbar^2} \left(\frac{Ze^2}{4\pi\epsilon_0} \right)^2 \quad (1)$$

where the negative sign indicates a bound state, n is the principal quantum number, Z is the nuclear charge, and μ is the reduced mass of the electron-nucleus system

$$\mu = \frac{m_e M}{m_e + M} \quad (2)$$

M being the mass of the nucleus. The energy of the photon emitted when an electron makes a transition from an initial level n_i to a final level n_f is thus

$$E = \hbar\omega = \frac{1}{2} \frac{\mu}{\hbar^2} \left(\frac{Ze^2}{4\pi\epsilon_0} \right)^2 \left(\frac{1}{n_f^2} - \frac{1}{n_i^2} \right) \quad (3)$$

As is well known, for atomic hydrogen the transitions result in long wavelength radiation, for example, the Lyman series in the ultraviolet, corresponding to $n_f = 1$. For hydrogenic ions, Eq. (3) shows that the spectral line energies are a factor of Z^2 higher than those of atomic hydrogen (and the wavelengths are correspondingly shorter); for hydrogenic carbon, for example, the Lyman- α line is at a wavelength of 3.37nm, compared to 121.6nm in atomic hydrogen.

Spectral Line Intensities

Characteristic radiation as described here corresponds to spontaneous emission, for which the transition rate, in the electric-dipole approximation, is⁵

$$R_{if} = \frac{\omega_{if}^3}{3\pi\epsilon_0 \hbar c^3} |\langle \psi_f | \mathbf{er} | \psi_i \rangle|^2 \quad (4)$$

where $\langle \psi_f | \mathbf{er} | \psi_i \rangle$ is the matrix element describing the transition from the initial state with wave function ψ_i via the electric dipole operator \mathbf{er} to the finale state ψ_f . The intensity of the emission is given by multiplying Eq. (4) by the energy $\hbar\omega_{if}$ of each emitted photon, the statistical weight g_i of the initial state, which is the number of degenerate configurations that can give rise to the transition,

and the probability F_f that an electron has been removed from the level to which the system eventually relaxes. Thus the intensity of a characteristic line is given by

$$I_{if} = \frac{g_i F_f \omega_{if}^4}{3\pi\epsilon_0 c^3} |\langle \psi_f | e\mathbf{r} | \psi_i \rangle|^2 \quad (5)$$

If the matrix element in the electric dipole transition rate is zero, then the transition cannot take place. This is embodied in the familiar electric dipole selection rules,

$$\Delta l = \pm 1 \quad \Delta m = 0, \pm 1 \quad \Delta j = 0, \pm 1 \quad \text{with } j=0 \rightarrow 0 \text{ forbidden} \quad (6)$$

where l is the orbital angular momentum, m is the magnetic quantum number and j is the total angular momentum.

For higher-order transitions the selection rules are more complicated. For electric quadrupole radiation, for example,

$$\Delta l = 0, \pm 2 \quad \Delta j = 0, \pm 1, \pm 2 \quad \text{with } 0 \rightarrow 0 \quad \frac{1}{2} \rightarrow \frac{1}{2} \text{ \& } 0 \leftrightarrow 1 \text{ forbidden} \quad (7)$$

and the parities of the initial and final states must be the same. For magnetic dipole transitions the parity must also remain unchanged, with

$$\Delta l = 0 \quad \Delta j = 0, \pm 1 \quad \text{with } 0 \rightarrow 0 \text{ forbidden} \quad (8)$$

The relative intensities in different spectral lines can, in principle, be calculated from Eq. (5), so long as all states with the same energy difference $E_i - E_f$ are taken into account. The intensities are characterized by the *oscillator strength*

$$\phi_{if} = \frac{2m_e \omega_{if}}{3\hbar} |\langle \psi_f | \mathbf{r} | \psi_i \rangle|^2 \quad (9)$$

where, for hydrogenic systems,

$$\sum_f \phi_{if} = 1 \quad (10)$$

the sum being over all final states that contribute to emission or absorption of a photon of energy $\hbar\omega_{if}$, including the continuum if appropriate. In terms of the oscillator strength, the transition rate, Eq. (4), becomes

$$R_{if} = \frac{2\hbar\alpha}{m_e c^2} \omega_{if}^2 |\phi_{if}| \quad (11)$$

where α is the fine structure constant and the modulus of ϕ_{if} is used since, by convention, the oscillator strength is defined as negative for emission ($E_f < E_i$) and positive for absorption. As hydrogenic wave functions can be calculated exactly, the oscillator strengths and transition rates are straightforward to determine.

For ions with Z electrons the sum rule, Eq. (10), becomes

$$\sum_f \phi_{if} = Z \quad (12)$$

and the oscillator strengths can only be estimated since approximate methods must be used to determine the wave functions. The oscillator strengths depend on the quantum numbers n , l , and m of the initial and final states, particularly on the magnetic quantum number m and hence on the

polarization of the radiation. It is often useful to use an oscillator strength averaged over the initial and final state magnetic quantum numbers, since this is independent of the polarization,

$$\bar{\varphi}_{if}(n_i, l_i, n_f, l_f) = \frac{1}{2l_i + 1} \sum_{m_i = -l_i}^{l_i} \sum_{m_f = -l_f}^{l_f} \varphi_{if}(n_i, l_i, m_i, n_f, l_f, m_f) \quad (13)$$

The sum rule, Eq. (10) or (12), is still obeyed for the average oscillator strength. The corresponding relative line intensities depend on the statistical weights $2(2l + 1)$ of the initial state and the spontaneous emission coefficient as well as on the oscillator strengths.⁶

For line emission from hydrogenic ions, as the atomic number increases relativistic effects become increasingly important.⁷ This causes the average oscillator strengths to decrease; for a low- Z element such as carbon the decreases are only $\sim 0.1\%$, for a medium- Z element such as copper they are about 3%, while for high- Z elements such as gold or uranium they are about 20% and 30%, respectively.

Spectral Line Shapes and Widths

The discussion of line emission so far has implicitly assumed that each spectral line has an exactly defined energy, with corresponding frequency and wavelength. However, in practice there is always a spread of energies in each line as the result of several factors.⁸ Since the spread depends upon the plasma properties, such as temperature and density, measurement of line shapes and widths provides a useful diagnostic of the plasma state.

Natural Line Width Natural line widths arise since atomic and ionic levels do not have precisely determined energies, because interactions with real or virtual photons result in transitions between levels. Hence the levels have finite lifetimes τ and so, as a result of the uncertainty principle, there is an energy spread

$$\Delta E \sim \frac{h}{\tau} \quad (14)$$

with a corresponding frequency spread

$$\Delta \nu \sim \frac{1}{\tau} \quad (15)$$

For transitions from a state i to a state j , the lifetime τ_{ij} is given by

$$\tau_{ij} = \frac{2}{\sum_{k < i} A_{ik} + \sum_{k < j} A_{jk}} \quad (16)$$

where A_{ij} is the spontaneous emission coefficient for the transition $i \rightarrow j$ and the sums are over all spontaneous transitions from the initial and final states, respectively. If the final state is the ground state then the second summation is equal to zero.

The amplitude $A(t)$ of a state with lifetime τ decays exponentially with time

$$A(t) \propto e^{-t/\tau} e^{2\pi i \nu_0 t} \quad (17)$$

where ν_0 is the central transition frequency. The shape of the spectral line is given by the square of the Fourier transform of the amplitude, resulting in a Lorentzian distribution

$$I(\nu) = I(\nu_0) \frac{1}{1 + [2\pi(\nu - \nu_0)\tau]^2} \quad (18)$$

where $I(\nu)$ is the intensity at frequency ν . The full width at half maximum (FWHM) of this distribution is

$$\Delta\nu_L = \frac{1}{\pi\tau} \quad (19)$$

In terms of wavelength, which is often how spectra are presented, the line shape is still (obviously) Lorentzian and the FWHM is

$$\Delta\lambda = \frac{\lambda_0^2}{\pi\tau c} \quad (20)$$

Doppler Broadening If the emitting ions are in thermal motion, there will be a Doppler shift $\Delta\nu_D$ on the frequency of the emitted radiation. In a plasma, the particles are normally nonrelativistic,

$$\Delta\nu_D = |\nu - \nu_0| = \nu_0 \frac{v}{c} \quad (21)$$

where v is the particle speed along the line of sight. Radiation from particles with a velocity component away from an observer will thus be red shifted, while a velocity component toward the observer causes a blue shift. For particles in random motion the spectral lines are therefore broadened and, if the velocity distribution $f(v)$ is Maxwellian with a mean speed v_m , then

$$f(v) \propto e^{-v^2/2v_m^2} \quad (22)$$

The resulting line profile, found simply by making the substitution $v \rightarrow c(\nu/\nu_0 - 1)$ in Eq. (22), is

$$I(\nu) = I(\nu_0) \exp\left[-\frac{(\nu - \nu_0)^2 c^2}{2v_m^2 \nu_0^2}\right] \quad (23)$$

which is a Gaussian distribution with FWHM

$$\Delta\nu_G = \nu_0 \frac{v_m}{c} 2(2\ln 2)^{1/2} \quad (24)$$

Clearly, the broadening is larger if the mean speed is higher, so that Doppler broadening can dominate in hot plasmas and is a measure of the plasma temperature. In terms of the wavelength, Eqs. (23) and (24) become

$$I(\lambda) = I(\lambda_0) \exp\left[-\frac{(\lambda - \lambda_0)^2 c^2}{2v_m^2 \lambda_0^2}\right] \quad (25)$$

and

$$\Delta\lambda_G = \lambda_0 \frac{v_m}{c} 2(2\ln 2)^{1/2} \quad (26)$$

Pressure Broadening In plasmas neighbouring particles affect the emission profiles of spectral lines. There are several contributions to this *pressure broadening*. Collisions, of either ions or electrons, with the emitting ion can change the ionic state, effectively shortening the lifetimes of the energy levels involved in the emission. By a similar argument to that in the discussion of natural broadening this leads to a Lorentzian profile but with a larger width. The width is governed by the

mean time between collisions, and so this process can be dominant in dense plasmas, and is a measure of the plasma density.

Nearby ions and electrons also generate electric and magnetic fields, which change the ionic energy levels of the emitting ion and hence contribute to the line width. The major contributor in plasmas is normally the electric field which, in general, produces a shift in the energy levels proportional to the square of the electric field—the quadratic Stark effect. This leads to asymmetric broadening and hence to shifts in the central frequencies of spectral lines. In hydrogenic ions the effect is linear, except in the ground state where the linear effect disappears, and hence the broadening is symmetric with no shift of the central frequency. The effect on line shapes, apart from in hydrogenic ions, is very complicated and requires each case to be treated individually,⁸ beyond the scope of the present discussion. However, it is possible to use simple arguments to obtain an idea of how line profiles and widths scale with plasma properties.⁹

In the nearest neighbour approximation, an ion in a range of distances dr about r from another ion experiences a high field; the probability of it being affected is proportional to the volume $4\pi r^2 dr$. The frequency shift $|v - v_0|$ caused by this is proportional to the field strength, i.e., to r^{-2} , so that

$$\begin{aligned} |v - v_0| &\propto r^{-2} & r &\propto |v - v_0|^{-1/2} & dr &\propto |v - v_0|^{-3/2} dv \\ r^2 dr &\propto |v - v_0|^{-5/2} dv \end{aligned} \quad (27)$$

and the resulting line profile is

$$I(v) \propto |v - v_0|^{-5/2} \quad (28)$$

This clearly breaks down for frequencies close to the central frequency, since otherwise the intensity $I(v_0)$ would become infinite. The line profile, Eq. (28), was derived using the nearest neighbour approximation which is only valid for the distance below which the electric field of the perturbing ion dominates over that due to other ions. This distance is determined by the volume containing about one ion, i.e., $r \approx [3/(4\pi n_1)]^{1/3}$, where n_1 is the ion number density. Thus the line width, which is proportional to r^{-2} , is

$$\Delta v_s \propto n_1^{2/3} \quad (29)$$

Combinations of Line Broadening Mechanisms The preceding discussions show that, if the line broadening mechanism can be identified, properties of the plasma, such as density and temperature, can be determined from measured line profiles and widths. In general, however, the emission lines will be affected by more than one process, and thus combinations of the various line broadening mechanisms must be considered. In addition, the measured profiles will be modified by the instrument used to obtain the information; this may be dominant and will not, in general, provide a well defined profile.

For two independent effects, with profiles $I_1(v)$ and $I_2(v)$ leading to FWHM line widths Δ_1 and Δ_2 , the combined profile is the convolution

$$I(v) = \int I_1(v - v') I_2(v') dv' \quad (30)$$

If the two profiles are the same the combined profile has the same shape; for Lorentzian distributions the combined width is

$$\Delta_L = \Delta_1 + \Delta_2 \quad (31)$$

while for Gaussian distributions it is

$$\Delta_G = (\Delta_1^2 + \Delta_2^2)^{1/2} \quad (32)$$

If the two profiles are not the same then there is no general simple form for the convolved line shape. The combination of a Lorentzian and a Gaussian profile leads to the Voigt function

$$I_V(q) \propto \int_{-\infty}^{\infty} \frac{e^{-y^2}}{(x-y)^2 + a^2} dy \quad (33)$$

where

$$q = \frac{2(\ln 2)^{1/2}}{\Delta V_G} (v - v_0) \quad \text{and} \quad a = (\ln 2)^{1/2} \frac{\Delta V_L}{\Delta V_G} \quad (34)$$

Equation (33) does not have an analytic solution; the integration has to be done numerically for each frequency involved and so is very expensive computationally. An alternative to the Voigt function, which has been used in many other fields^{10–17} but not widely in the analysis of x-ray emission from plasmas, is the Pearson VII function.^{18–20} This has the form

$$I_P(\lambda) = I(\lambda_0) [1 + P(\lambda)^2]^{-M} \quad (35)$$

where

$$P(\lambda) = \frac{2(\lambda - \lambda_0) \sqrt{2^{1/M} - 1}}{\Delta_p} \quad (36)$$

In Eqs. (35) and (36) M is the Pearson parameter and Δ_p is the FWHM. If $M = 1$ the distribution is clearly Lorentzian, while as M becomes large the distribution is indistinguishable from a Gaussian. For intermediate values of M the Pearson VII distribution is more general than the Lorentzian and Gaussian functions due to the extra free parameter (M), and it makes no assumptions about the line profile. It can also provide a good approximation to a Voigt function, compared to which it is much less computationally intensive since it does not require numerical integration. For $M < 1$ the distribution has broader wings than a Lorentzian function; this corresponds to no known broadening mechanism but even then the Pearson VII distribution can still provide a convenient functional form for line profile fitting.

A disadvantage of the Pearson VII distribution is that there are no direct relationships between the plasma properties and the fit parameters Δ_p and M . The widths of the Lorentzian and Gaussian distributions, which can be extracted from fits to the Voigt function can, however, be related to plasma properties such as density and temperature. But if the spectra are time and spatially integrated the plasma parameters cannot readily be extracted, so that there is no advantage in using the Lorentzian, Gaussian, or Voigt profiles. Unlike the Voigt function, the Pearson VII distribution can be integrated analytically to provide the total line intensity

$$I_P^{\text{tot}} = \frac{\Delta_p I_0 \sqrt{\pi}}{2\sqrt{2^{1/M} - 1}} \frac{\Gamma\left(M - \frac{1}{2}\right)}{\Gamma(M)} \quad (37)$$

where Γ is the gamma function,

$$\Gamma(g) = \int_0^{\infty} y^{g-1} e^{-y} dy \quad (38)$$

A common characteristic of the Lorentzian, Gaussian, Voigt, and Pearson VII profiles is that they are symmetrical about the mean. However, some line-broadening mechanisms, e.g., Stark broadening, can result in asymmetric profiles. In this case a suitable fitting function is the Pearson IV distribution^{18,20}

$$I_{P\text{IV}}(\lambda) = I_{P\text{VII}}(\lambda) \exp[-a \tan^{-1} P(\lambda)] \quad (39)$$

which clearly reduces to the Pearson VII distribution when a , the asymmetry parameter, is zero. The disadvantages of the Pearson IV distribution are the extra parameter required in the fit, and the total

intensity in the line has to be obtained by numerical integration since the distribution is not analytically integrable. Hence the fitting is more computationally expensive than for the Pearson VII distribution, but the Pearson IV function may be useful in analysis of high spectral resolution data where Stark broadening may be observed.

Reabsorption of Spectral Lines The discussion of spectral lines has so far assumed that the emission profiles are not modified by reabsorption within the plasma. This may not always be the case, especially in optically thick plasmas.

The linear absorption coefficient α_λ , which is a function of the wavelength of the radiation, describes the amount of attenuation over a thickness x of a medium

$$I_\lambda(x) = I_\lambda(0)e^{-\alpha_\lambda x} \quad (40)$$

The transmission of radiation can be the same through a tenuous extended plasma or through a dense compact plasma. It may not be necessary, or possible, for a remote observer to distinguish these, in which case the optical thickness (or depth) τ_λ is a useful concept; this is defined as the absorption coefficient (which may depend upon position) integrated between two points x_1 and x_2 ,

$$\tau_\lambda = \int_{x_1}^{x_2} \alpha_\lambda dx \quad (41)$$

An optically thin medium has $\tau_\lambda < 1$ while for an optically thick medium $\tau_\lambda \gg 1$; a plasma may be optically thin for some wavelengths and optically thick for others.

In an optically thick plasma, some of the radiation in a spectral line may be reabsorbed, which can alter the line profile and apparent FWHM. If there is no reabsorption, a beam of radiation passing through the plasma will gain in intensity due to spontaneous and stimulated emission. However, this cannot continue indefinitely since the intensity cannot become larger than that of a blackbody at the effective plasma temperature. Hence, as the blackbody level is approached, absorption must balance the emission and the intensity saturates. Since a spectral line is most intense at its central wavelength the emitted intensity close to the line center is the first to be limited. Thus the line becomes flattened near to its center, which modifies the profile and increases the apparent line width.

In a nonuniform plasma, for example, if the temperature decreases near to the edge—which is often the case—the observed central line intensity can be decreased since the observed emission near to the line center is from the cooler plasma, which has a lower blackbody intensity. This is known as line inversion and can lead to confusion in the interpretation of plasma spectra, as it may appear as two closely spaced lines.

56.3 BREMSSTRAHLUNG

The preceding discussion of line emission refers to the situation in which an electron bound in an ion makes a transition to another bound state of the same ion (bound-bound radiation). In this section the transition of a free electron (or collection of electrons) to another free state is considered—free-free radiation, caused by the acceleration of the electron in the Coulomb field of an ion. According to Maxwell's equations, accelerated charges emit electromagnetic radiation and so they lose energy; the emitted radiation is thus called "braking radiation" or *bremstrahlung*. An important difference between bound-bound and free-free radiation is that free electrons do not have quantised energies, and hence the emission spectrum is a continuum.

Although in practice many electrons will encounter many ions, it is instructive to consider, in the first instance, the interaction of a single electron with a single ion. In this case the radiated energy in the frequency range $d\omega$ is given by

$$\left. \frac{dW}{d\omega} \right|_{\omega=\omega_0} \approx \frac{Z^2 e^6}{(4\pi\epsilon_0)^3} \frac{8}{3\pi m_e^2 c^3} \frac{1}{v^2 b^2} \quad (42)$$

where v is the electron speed, b is the impact parameter (the perpendicular distance between the ion and the initial electron direction) and ω_0 is the central frequency of the emitted radiation. If the single electron collides with a random assembly of ions, with a number density $n_i \text{ m}^{-3}$, the power spectrum is given by multiplying Eq. (42) by $n_i v$ and integrating over the impact parameter

$$\frac{dP}{d\omega} = n_i v \int_{b_{\min}}^{b_{\max}} \frac{dW}{d\omega} 2\pi b db \approx \frac{Z^2 e^6}{(4\pi\epsilon_0)^3} \frac{16n_i}{3\pi m_e^2 c^3 v} \int_{b_{\min}}^{b_{\max}} \frac{1}{b} db = \frac{Z^2 e^6}{(4\pi\epsilon_0)^3} \frac{16n_i}{3\pi m_e^2 c^3 v} \ln \frac{b_{\max}}{b_{\min}} \quad (43)$$

Naïvely, Eq. (43) might suggest that a single electron encountering a group of ions would radiate infinite power, with a minimum impact parameter $b_{\min} = 0$ (head-on collision) and a maximum of $b_{\max} = \infty$. This is clearly incorrect; the explanation is that b_{\max} is determined by shielding—since the electron is only accelerated if the nuclear charge is not shielded by surrounding particles—and b_{\min} is governed by the Uncertainty Principle, since the electron is only close to an ion for a finite time and so the impact parameter cannot be determined precisely.

In the nonquantum, nonrelativistic limit the ratio b_{\max}/b_{\min} may be obtained in terms of the impact parameter b_{90} for which the electron is scattered through an angle of 90° . For the case $b \gg b_{90}$ this leads to

$$\frac{b_{\max}}{b_{\min}} = \frac{2\gamma v}{\omega b_{90}} \quad (44)$$

where γ is Euler's constant

$$\gamma = \lim_{n \rightarrow \infty} \left(\sum_{k=1}^n \frac{1}{k} - \ln n \right) = 0.5772157\dots \quad (45)$$

Equation (43) may generally be written as

$$\frac{dP}{d\omega} = \frac{Z^2 e^6}{(4\pi\epsilon_0)^3} \frac{16n_i}{3\pi m_e^2 c^3 v} \frac{\pi}{\sqrt{3}} G \quad (46)$$

where, for $b \gg b_{90}$,

$$G = \frac{\sqrt{3}}{\pi} \ln \frac{2\gamma v}{\omega b_{90}} \quad (47)$$

Equation (46) can be used for other cases, including when quantum and relativistic effects are important (neither of which are normally appropriate for laser-generated plasmas used as x-ray sources), so long as an appropriate choice is made for the Gaunt factor, G , which never differs substantially from unity; in many cases using $G = 1$ is a sufficient approximation. Several authors have calculated Gaunt factors for bremsstrahlung under a range of conditions, including relativistic, and extensive tabulations or graphical representations are available.^{21,22}

Bremsstrahlung from a Collection of Electrons In a plasma the electron density n_e is usually larger than the ion density n_i , and so bremsstrahlung from a collection of electrons must be considered. The radiated power must then be integrated over the electron velocity distribution $f(\mathbf{v})$ and the spectral power emitted into the whole solid angle per unit angular frequency is then

$$4\pi w(\omega) = \frac{Z^2 e^6}{(4\pi\epsilon_0)^3} \frac{16\pi n_i}{3\sqrt{3}\pi m_e^2 c^3} \int \frac{G(\omega, \mathbf{v})}{v} f(\mathbf{v}) d^3\mathbf{v} \quad (48)$$

For an isotropic Maxwellian velocity distribution, as usually assumed,

$$4\pi w(\omega) = n_e n_i \frac{Z^2 e^6}{(4\pi\epsilon_0)^3} \frac{16\pi}{3\sqrt{3}m_e^2 c^3} \left(\frac{2m_e}{\pi T_{eV}}\right)^{1/2} e^{-\hbar\omega/T_{eV}} \bar{G}(\omega, T_{eV}) \quad (49)$$

where T_{eV} is the plasma temperature in electronvolts and $\hbar\omega$ is the emitted photon energy. Values of the Maxwell-averaged Gaunt factor, $\bar{G}(\omega, T_{eV})$ are typically ~ 1 , and this is normally sufficient to give an accurate enough approximation of the bremsstrahlung intensity.

56.4 RECOMBINATION RADIATION

In some circumstances in plasmas recombination radiation, emitted when an electron is captured by an ion, can be a dominant part of the spectrum, competing with and modifying the bremsstrahlung component. In recombination radiation the emitted energy is clearly greater than the free-electron kinetic energy, and the final-state electron is bound with a discrete energy spectrum E_n . For hydrogenic ions of charge Ze , the emitted spectral power is given by

$$4\pi w(\omega) = n_e n_i \frac{Z^2 e^6}{(4\pi\epsilon_0)^3} \frac{16\pi}{3\sqrt{3}m_e^2 c^3} \left(\frac{2m_e}{\pi T_{eV}}\right)^{1/2} e^{-\hbar\omega/T_{eV}} F_n \quad (50)$$

where

$$F_n = \frac{Z^2 e^4 m_e}{(4\pi\epsilon_0 \hbar)^2 n^3} \frac{1}{T_{eV}} G_n e^{-E_n/T_{eV}} \quad (51)$$

and G_n is the Gaunt factor for capture of the free electron to the energy level with principal quantum number n and energy E_n . Equation (50) is the same as that for bremsstrahlung from a collection of electrons, Eq. (49), except that F_n has replaced the Maxwell-averaged Gaunt factor $\bar{G}(\omega, T_{eV})$. The ratio of powers emitted in recombination radiation and bremsstrahlung is thus

$$\frac{F_n}{\bar{g}(\omega, T_{eV})} = \frac{Z^2 e^4 m_e}{(4\pi\epsilon_0 \hbar)^2 n^3} \frac{1}{T_{eV}} \frac{G_n}{\bar{g}(\omega, T_{eV})} e^{-E_n/T_{eV}} \quad (52)$$

It must be remembered that this is only for hydrogenic ions, but Eq. (52) gives an idea of the order of magnitude of recombination radiation compared to bremsstrahlung. Putting in the values of the constants, and recalling that the ratio of the Gaunt factors will be of order unity,

$$\frac{F_n}{\bar{G}(\omega, T_{eV})} \approx \frac{27.2 Z^2}{n^3 T_{eV}} e^{13.6 Z^2/n^2 T_{eV}} \quad (53)$$

which shows that recombination radiation can be many orders of magnitude more intense than bremsstrahlung for low temperature plasmas and for recombination to inner ionic levels, i.e., small n .

56.5 REFERENCES

1. I. C. E. Turcu, I. N. Ross, P. Tenda, C. W. Wharton, R. A. Meldrum, H. Daido, M. S. Schulz, et al., "Picosecond Excimer Laser-Plasma X-Ray Source for Microscopy, Biochemistry, and Lithography," *Proc. SPIE* **2015**:243–260 (1993).
2. H. M. Hertz, L. Rymell, M. Berglund, and L. Malmqvist, "Debris-Free Liquid-Target Laser-Plasma Soft X-Ray Source for Microscopy and Lithography," in *X-Ray Microscopy and Spectromicroscopy*, Berlin: Springer, 1998, pp. V-3–V-13.

3. H. Fiedorowicz, A. Bartnik, H. Szczurek, H. Daido, N. Sakaya, V. Kmetik, Y. Kato, et al., "Investigation of Soft X-Ray Emission from a Gas Puff Target Irradiated with a Nd:YAG Laser," *Opt. Comm.* **163**:103–114 (1999).
4. C. Reich, P. Gibbon, I. Uschmann, and E. Förster, "Yield Optimization and Time Structure of Femtosecond Laser Plasma K_{α} Sources," *Phys. Rev. Lett.* **84**:4846–4849 (2000).
5. B. H. Bransden and C. J. Joachain, *Physics of Atoms and Molecules* 2nd ed., London: Longmans, 2002, pp. 195–197.
6. H. A. Bethe and E. E. Salpeter, *Quantum Mechanics of One- and Two-Electron Atoms*, Berlin: Springer, 1957, p. 265.
7. Y. G. Pal'chikov, "Relativistic Transition Probabilities and Oscillator Strengths in Hydrogen-Like Atoms," *Physica Scripta* **57**:581–593 (1998).
8. R. G. Breene, *The Shift and Shape of Spectral Lines*, Oxford: Pergamon Press, 1962.
9. I. H. Hutchinson, *Principles of Plasma Diagnostics*, Cambridge: Cambridge University Press, 1987, pp. 218–221.
10. H. Toraya, M. Yoshimura, and S. Somiya, "A Computer Program for the Deconvolution of X-Ray Diffraction Profiles with the Composite of Pearson Type VII Functions," *J. Appl. Cryst.* **16**:653–657 (1983).
11. C. P. Lafrance, J. Debigare, and R. E. Prudhomme, "Study of Crystalline Orientation in Drawn Ultra-High Molecular Weight Polyethylene Films," *J. Polym. Sci. Pol. Phys.* **31**:255–264 (1993).
12. M. Oetzel and G. Heger, "Laboratory X-Ray Powder Diffraction: A Comparison of Different Geometries with Special Attention to the Usage of the $Cu K_{\alpha}$ Doublet," *J. Appl. Cryst.* **32**:799–807 (1999).
13. A. Santoro, R. J. Cava, D. W. Murphy, and R. S. Roth, "Use of the Pearson Type VII Distribution in the Neutron Profile Refinement of the Structures of $LiReO_3$ and Li_2ReO_3 ," *Proc. AIP* **89**:162–165 (1982).
14. D. R. Noakes, "Magnetic Field Distributions at Interstitial Sites in Nondilute Alloys," *Phys. Rev. B* **44**:5064–5072 (1991).
15. V. Villani, R. Pucciariello, and G. Ajroldi, "Calorimetric Study of the Room-Temperature Transitions of Polytetrafluoroethylene—The Influence of Thermal History," *J. Polym. Sci. Pol. Phys.* **29**:1255–1259 (1991).
16. R. Oven, D. G. Ashworth, and M. D. J. Bowyer, "Formulas for the Distribution of Ions Under an Ideal Mask," *J. Phys. D* **25**:1235–1237 (1992).
17. W. Wulfhekel and J. M. Cadogan, "Mössbauer Line-Sharpening—Application to Magnetically Split Spectra," *Hyperfine Interact.* **92**:1195–1202 (1994).
18. W. P. Elderton and N. L. Johnson, *Systems of Frequency Curves*, London: Cambridge University Press, 1969.
19. A. G. Michette and S. J. Pfauntsch, "Laser Plasma X-Ray Line Spectra Fitted Using the Pearson VII Function," *J. Phys. D Appl. Phys.* **33**:1186–1190 (2000).
20. S. J. Pfauntsch, *Developments in Soft X-Ray Laboratory Systems for Microscopy and Cellular Probing*, Ph.D. Thesis (London University), 2001.
21. P. J. Brussard and H. C. van der Hulst, "Approximation Formulas for Nonrelativistic Bremsstrahlung and Average Gaunt Factors for a Maxwellian Electron Gas," *Rev. Mod. Phys.* **34**:507–520 (1962).
22. T. R. Carson, "Coulomb Free-Free Gaunt Factors," *Astron. Astrophys.* **189**:319–324 (1988).

This page intentionally left blank.

DO NOT DUPLICATE

Victor Kantsyrev

*Physics Department
University of Nevada
Reno, Nevada*

57.1 INTRODUCTION

In pinch systems, the source of x-ray radiation is the plasma that is produced by an electrical current discharge and compressed by the associated magnetic field. Numerous forms of pinch plasma are used as radiation sources (Z-pinch,¹ θ -pinch,² etc.). The most common type of pinch plasma source is the Z-pinch. In the Z-pinch system, plasma is created by applying a high-voltage pulse from a fast generator between an anode and cathode; its name refers to the direction of the earliest experimental facility, where the current flowed down the vertical discharged tube, the z axis on a normal mathematical diagram. Both electrodes have cylindrical geometry and are placed in a vacuum chamber, where vacuum is typically held at 10^{-4} to 10^{-5} torr. The anode-cathode gap is initially filled with gas from a fast valve gas-puff system or connected by an array of tiny (micrometers in diameter) wires, placed around the z axis. The typical Z-pinch plasma implosion process is illustrated in Fig. 1, where the current starts to flow through outer layers of plasma due to a skin effect, and the Z-pinch plasma implodes as a shell. After the discharge current I starts to flow axially through anode-cathode gap, it ionizes the wires or gas, produces the azimuthal magnetic field \mathbf{B} , and implodes the plasma by magnetic field forces. During compression of plasma into a hot dense column and its stagnation, the gas dynamic pressure equalizes the magnetic forces (Phase III in Fig. 1),³ and the kinetic energy is converted to thermal energy and to a radiation pulse. The implosion typically lasts from 50 ns to several hundred microseconds. For a plasma that consists of higher atomic number elements, the radiative collapse can happen when strong line or continuum radiation leads to plasma cooling, which decreases the gas-dynamic pressure, and allows more significant plasma compression, to a diameter of 10 to 100 μm . Under such conditions, Z-pinch loads reveal small, localized sources of x-ray radiation (“hot” or “bright” spots, the size of which may be as small as 1 to 5 μm) formed at random points along the Z-pinch axis.

In Z-pinch plasmas, EUV and x-ray radiation occurs by ionization, excitation, and recombination processes which involve multicharged ions,⁴ as shown in Fig. 2.⁵ Bremsstrahlung emission is negligible. For high-current Z-pinch systems (more than 3 to 10 MA) high-atomic-number plasmas (Pt, W, etc.), are typically optically thick, and plasma radiation can be approximated by a black-body spectrum with some presence of spectral lines of multicharged ions.¹ At university-scale Z-pinch generators, with currents of about 0.5 to 1 MA, the total EUV/x-ray radiation yield can be as high as 10 to 30% of the electromagnetic energy delivered to plasma.⁶

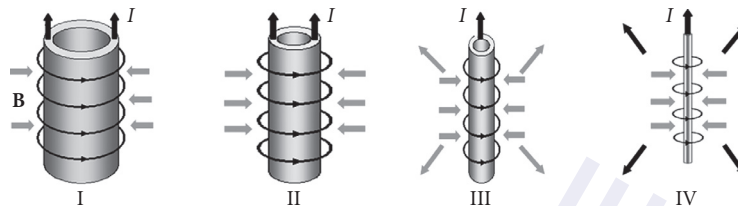


FIGURE 1 Z-pinch plasma implosion, I indicating current, and B , the magnetic field. The grey arrows pointing to the central axis show magnetic forces that compress plasma; the grey arrows outside the plasma show optical and extreme ultraviolet (EUV) emission, and black arrows showing x-ray emission. Phase I: acceleration of plasma shell from diameter 2 to 12 cm; Phase II: run-in, start of optical and EUV emission; Phase III: stagnation at diameter 1 to 4 mm; generation of softer x-rays and keV bursts; plasma reaches the thermalization stage with high concentration of multi-charged ions; Phase IV: for plasma that consists of higher atomic number elements, radiative collapse can happen when strong line radiation occurrence leads to plasma radiation cooling, a decrease in gas-dynamic pressure and additional plasma implosion to diameter 10 to 100 μm .

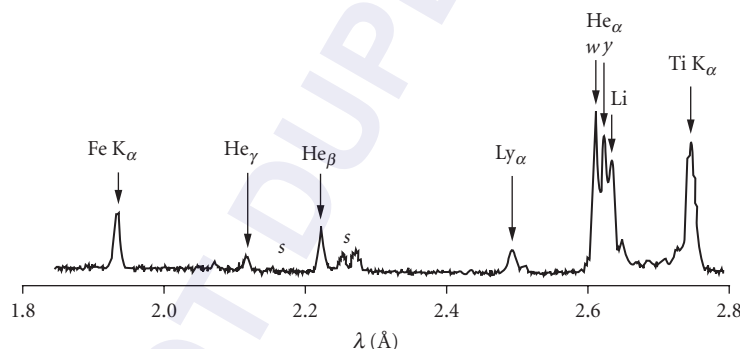


FIGURE 2 Typical spectrum of Ti X-pinch obtained at the Z-pinch Zebra generator. Peak current is 1 MA. The most intense and important spectral lines for diagnostics include the He-like lines, resonance line $\text{He}_{\alpha}(w)$, intercombination line y , resonance lines He_{β} and He_{γ} , and H-like resonance line Ly_{α} . They also include the satellite lines Li and S of Li-like ions to above mentioned resonance lines and “cold” $K\alpha$ lines of Ti and Fe. The “cold” lines were generated by electron beam existing in X-pinch plasma; Ti line was from material of wires and Fe line was from stainless steel anode.

57.2 TYPES OF Z-PINCH RADIATION SOURCES

Z-pinch plasma radiation sources have higher conversion coefficients of input electrical power to EUV/x-ray radiation compared to laser plasma or conventional x-ray sources. Z-pinch systems can use a wide range of materials for creation of plasma, including metals, dielectrics, and gases. Z-pinch radiation sources can be divided into multishot and single-shot systems. The first type includes plasma focus, gas-puff, vacuum spark, and capillary discharge devices. They can operate even in a repetition pulse regime (up to 0.1 to 0.5 Hz), but not for a lengthy period of time (from tens to several thousand shots). Their parameters are shown in Table 1.

TABLE 1 Characteristics of Pinch Sources

Source ^a	Mode ^b / Current ^c (MA)	Average Source Size ^d (mm)	X-Ray Burst Duration ^e (ns)	Total Yield ^f (kJ)
Plasma focus	multi 0.1–2	1–50	20–100	0.1–50
Gas-puff	multi 0.1–15	0.5–50	5–100	0.1–100
Vacuum spark	multi 0.05–0.4	0.01–5	5–20	0.05–0.25
Wire-array	single 0.3–20	0.5–20	5–50	0.1–1,800
X-pinch	single 0.1–1	0.002–5	0.1–10	0.01–10

^aThe minimum values given for sources parameters are for table-top devices.

^bThe mode refers to multishot or single-shot operational mode.

^cCurrent is peak current.

^dAverage source size is the size of the emitting region observed in single shot.

^eX-ray burst duration is the FWHM.

^fThe total yield is the total EUV/x-ray radiation yield.

The plasma focus is similar to a plasma gun device, where plasma sheath ejects between two cylindrical coaxial electrodes (positioned in chamber fill with H₂, D₂, He, or another gas under pressure of several torr) and implodes on central axis near electrodes' edge.⁷ Plasma focus facilities with current from 100 kA to 2 MA are mainly applied in fusion research, x-ray lithography and micromachining, surface material modification, and x-ray sources for medical purposes.

For a gas-puff source,⁸ the anode-cathode gap in the vacuum chamber fills first with a supersonic hollow gas stream from a fast valve (50 to 100- μ sec-duration pulse), and then high voltage is applied to electrodes and the plasma implodes, generating an x-ray burst 10 to 100 ns in duration. The gas-puff devices cover the range from compact table-top machines with current near 100 kA to generators with 5 to 15 MA of current. The large facilities are mainly pulse sources of K-shell radiation (Ne, Ar, and Kr) with a yield more than 30 kJ, applied in fusion research and x-ray laser studies. Table-top gas-puff devices are used in fusion studies, x-ray lithography and micromachining, x-ray laser studies, and x-ray spectroscopy.

A vacuum spark⁹ is a simple compact device that consists of two electrodes positioned in a vacuum chamber and separated by a 5 to 10 mm gap. A needle-shaped anode electrode is connected to a fast capacitor with an initial stored energy of 1 to 3 kJ. The cathode electrode includes a small plasma gun. After the plasma gun ejects a cold plasma cloud into the inter-electrode gap, the current begins to flow, evaporating and ionizing anode material. For designs without a plasma gun, the cold plasma is created by a pulsed laser beam that focuses on flat cathode. The vacuum spark current peak is typically smaller than 200 to 400 kA. The plasma column implodes, with the formation of hot spots that generate x-ray bursts. Vacuum spark generators are used in x-ray spectroscopy, x-ray lithography and micromachining, and as a testing facility for new diagnostics. The vacuum spark's main shortcoming is instabilities of the x-ray yield from shot to shot.

In capillary discharge devices,¹⁰ the plasma is formed on the inner surface of a ceramic capillary (with inner diameter several millimeters, and length 100 to 150 mm) under a high-voltage pulse, and implodes in a narrow column on the capillary central axis. Such devices are table-top and generate mainly EUV bursts. They became widely applied in EUV/X-ray lasers¹⁰ in the last decade.

Single-shot x-ray sources are wire array and X-pinch. In these sources, the wires, anode, and cathode electrodes should be changed after each shot. These sources generate maximum radiation yields and powers in EUV and x-ray spectral regions. X-ray yields of 1.8 MJ and powers of 250 TW were obtained for facilities such as Sandia National Laboratories Z-generator with wire array loads.¹¹ In a Z-pinch source with a conventional wire array load, as shown in Fig. 3a,¹² the anode and cathode are initially connected by a cylindrical array (or double cylindrical array, also called nested array) of wires several micrometers in diameter that are placed around a central z axis. The wire number varies from 4 to 60 (on university-scale generators) to 200 to 300 (on multi-MA machines); array diameter is from 8 to 10 mm to 30 to 40 mm. The length of wires connecting anode and cathode is 10 to 20 mm. The mechanisms of the implosion and x-ray burst generation are generally

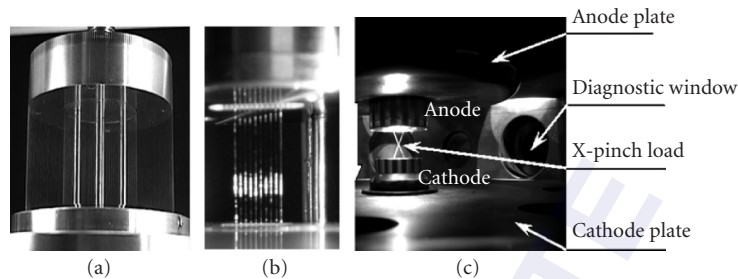


FIGURE 3 (a) A nested cylindrical wire array.¹² The central rods that are supporting the frame should be removed before the shot. (b) Single planar wire array.¹³ The gap between the wires is 1 mm. The supporting rod at the right should be removed before the shot. (c) X-pinch load in the vacuum chamber.¹⁴

the same, as described in Fig. 1. Source parameters are shown in Table 1. For the development of more effective plasma sources than conventional cylindrical arrays, new load geometries would help to achieve plasma of higher temperature and density, and reduce the source size and possibly shape of radiation pulse. A single planar array, as shown in Fig. 3b, consists of tiny wires placed in one-plane row parallel to each other between anode and cathode plates.¹³ Multiplanar arrays consist of 2 to 3 such rows placed parallel to each other between the anode and cathode plates.⁶ Compact single planar wire array and multiplanar wire arrays (with the width of just several mm) as well as compact cylindrical wire arrays (diameter smaller than 4 to 6 mm) show better data than any other tested loads at 1-MA generators. The maximum electron temperature T_e and density n_e are 1 to 1.4 keV and 10^{21} to 10^{22} cm^{-3} , respectively.^{6,13} The maximum radiation power in short nanosecond-scale rise-time x-ray bursts is around 1 TW and the yield is 23 to 25 kJ from 50 to 60 kJ of electrical energy delivered by the generator to the Z-pinch load.^{6,13}

An X-pinch load is formed by crossing two or more wires at one point, as shown in Fig. 3c. X-pinch plasmas have been actively investigated for currents from 0.1 to 1 MA.^{14,15} The unique property of X-pinch as “a point source” is that the emitting region is localized near the wire crossing point and the source size varies from 1 to 2 μm up to hundreds of micrometers. X-pinch plasmas can radiate just one or two bursts with emission times ranging from a few nanoseconds to less than 100 ps,¹⁵ as described in Table 1. At 1-MA university-scale generators, the total radiation yield can be 10 kJ, which is half of the output of planar wire array and compact cylindrical array sources, but, the burst power is comparable due to the shorter X-pinch x-ray burst duration.¹⁴

57.3 CHOICE OF OPTICS FOR Z-PINCH SOURCES

All Z-pinch plasma sources are characterized by the presence of not only powerful x-ray/EUV bursts, but also intense neutral and ion beams that can damage the filters (micrometer-thick plastic or metal films) typically used to protect x-ray/EUV optics as zone plates, transmission gratings, and multilayer and grazing incidence systems (see Chaps. 38, 40, 41, 44, 48, and 49). One of the solutions is the application of strong permanent magnets positioned between plasma and optics. The magnetic field deflects ion beams from the filter and significantly increases the filter lifetime. Table-top plasma focus, gas-puff, vacuum spark, and capillary discharge devices can be used with zone plates in applications which require monochromatic radiation, such as x-ray microscopy. Those sources also can be applied in x-ray lithography, which needs high x-ray flux, with multilayer or grazing incidence optics (high-bandpass x-ray optical systems).

For more powerful pinch sources, like 1-MA University-scale generators, the situation is more difficult, because the sensitive x-ray optical systems, even with magnetic protection, sustain damage even at a distance of 1 to 1.5 m from plasma,¹⁴ and the efficiency of the optics decreases dramatically.

Relatively inexpensive high-bandpass EUV/x-ray glass capillary optics (see Chaps. 52 and 53) can be used with such powerful x-ray Z-pinch sources.¹⁴ These can be placed at a distance of several centimeters from a z-pinch source and changed after each shot. Glass capillary optics can be adapted for any source configuration, such as a plasma column in wire array source or a point-type source in an x-pinch system. Estimates show¹⁴ that for a 1-MA wire array (linear source) or x-pinch (point source) with glass capillary optics with a 2 to 3 mm diameter focusing spot, it is possible to reach energy densities of 2 to 10 J/cm² and flux densities of (1 to 5) × 10⁹ W/cm². Because a typical surface melting threshold is about 1 to 2 J/cm², such Z-pinch sources with glass capillary optics can be used for surface modification research.

57.4 REFERENCES

1. C. Deeney, C. A. Coverdale, M. R. Douglas, et al., "Titanium K-shell X-Ray Radiation from High Velocity Wire Array Implosions on the 20-MA Z Accelerator," *Phys. Plasmas* **6**:2081 (1999).
2. M. H. Elghazaly, A. M. Abd Elbaky, A. H. Bassyouni, et al., "Spectroscopic Studies of Plasma Temperature in a Low-Pressure Hydrogen in Plasma," *J. Quant. Spectrosc. Rad. Transf.* **61**:503 (1999).
3. A. Velikovich, J. Davis, Y. K. Chong, et al., "Implosion Instabilities and Mitigation," *Minicourse on the Physics of Z-Pinches*, ICOPS, Monterey, CA, 2005.
4. N. R. Pereira and J. Davis, "X-Rays from Z-Pinches on Relativistic Electron-Beam Generators," *J. Appl. Physics* **64**:R1 (1988).
5. A. S. Shlyaptseva, S. B. Hansen, V. L. Kantsyrev, et al., "X-Ray Spectropolarimetry of High-Temperature Plasmas," *Rev. Scientific Instr.* **72**:1241 (2001).
6. V. L. Kantsyrev, L. I. Rudakov, A. S. Safronova, et al., "Double Planar Wire Array as a Compact Plasma Radiation Source," *Phys. Plasmas* **15**:030704 (2008).
7. N. V. Filipov, T. I. Fillipova, and V. P. Vinogradov, "Dense High-Temperature Plasma in a Noncylindrical Z-Pinch Compression," *Nucl. Fusion Suppl.* **2**:577 (1962).
8. F. C. Young, R. J. Comisso, D. P. Murphy, et al., "Measurement and Analysis of Continuum Radiation from a Large-Diameter Long Implosion Time Argon Gas Puff Z-Pinch at 6 MA," *IEEE Trans. Plasma Sci.* **34**(5):2312 (2006).
9. H. Flemberg, "X-Ray Spectra of High-Temperature Plasma," *Ark. Mat. Astron. Fys. A* **28**(18):1 (1942).
10. J. J. Rocca, D. P. Clark, J. L. A. Chilla, et al., "Energy Extraction and Achievement of the Saturation Limit in a Discharge-Pumped Table-Top Soft X-Ray Amplifier," *Phys. Rev. Lett.* **77**:1476 (1996).
11. M. K. Matzen, M. A. Sweeney, R. G. Adams, et al., "Pulsed-Power-Driven High Energy Density Physics and Inertial Confinement Fusion Research," *Phys. Plasmas* **12**:055503 (2005).
12. C. Deeney, "History of Z-Pinches," *Minicourse on the Physics of Z-Pinches*, ICOPS, Monterey, CA, 2005.
13. V. L. Kantsyrev, L. I. Rudakov, A. S. Safronova, et al., "Properties of a Planar Wire Arrays Z-Pinch Source and Comparisons with Cylindrical Arrays," *High Energ. Dens. Physics* **3**:136 (2007).
14. V. L. Kantsyrev, D. A. Fedin, A. S. Shlyaptseva, et al., "High-Z 0.9–1.0 MA X-Pinch as a Possible Backlighter in 50–100 keV and Sub-keV-10 Kev Spectral Regions and a Powerful Soft X-Ray Source for Surface Modification Research," *Rev. Scientific Instr.* **74**:1935 (2003).
15. S. A. Pikuz, D. B. Sinars, T. A. Shelkovenko, et al., "High Energy Density Z-Pinch Plasma Conditions with Picosecond Time Resolution," *Phys. Rev. Lett.* **89**:035003 (2002).

This page intentionally left blank.

DO NOT DUPLICATE

Greg Tallents

*University of York
York, United Kingdom*

58.1 FREE-ELECTRON LASERS

The original concept for a free-electron laser was proposed by Madey in 1971¹ with free-electron laser output observed by Elias in 1976². In a free-electron laser, stimulated emission occurs from relativistic electrons passing through a spatially periodic transverse magnetic field (known as an undulator). In free-electron lasers operating at wavelength λ , the electrons interact with the laser electric field and arrange themselves into “microbunches” separated by one wavelength λ so that the electric field arising from the stimulated emission by each electron is coherently additive at each spatial point. As the intensity of laser output is proportional to the square of the laser electric field, if N electrons are present, the free-electron laser intensity is proportional to N^2 . The simplest form of free-electron lasing is where spontaneously emitted photons are amplified, so-called self-amplified spontaneous emission (SASE) output. The frequency bandwidth of output in SASE mode is usually large ($\nu/\Delta\nu \approx 10$), potentially enabling ultrashort laser pulses ($<10^{-15}$ s) to be produced. However, large frequency bandwidth implies a short coherence length that can be limiting for some applications. To reduce the free-electron laser bandwidth, increase the coherence length and increase the output brilliance (photons $\text{s}^{-1}\text{mm}^{-2}\text{mrad}^{-2}$ 0.1%BW); it will be possible to “seed” the amplification process with photons produced by infrared laser harmonic radiation. The production of harmonic radiation is discussed in this chapter.

In 2008, free-electron laser output of ≈ 70 μJ per pulse had been produced down to a fundamental wavelength $\lambda \approx 6$ nm together with significantly shorter wavelengths, but less intense harmonic components at the FLASH free-electron laser facility at Deutsches Elektronen-Synchrotron (DESY) in Hamburg.³ Free-electron laser operation will be extended to an unprecedented short wavelength of ≈ 0.1 nm with the European X-Ray Free-Electron Laser (XFEL) to be also constructed at DESY with the Linac Coherent Light Source (LCLS) at the Stanford Linear Accelerator Center (SLAC) due to come on-line in 2009 and with the Spring-8 Compact SASE Source (SCSS) under development in Japan.

Free-electron lasers (FELs) are intense (peak brilliance up to 10^{34} photons $\text{s}^{-1}\text{mm}^{-2}$ 0.1% bandwidth is expected for XFEL),³ have high transverse coherence, and may be able to be ultimately focused to focal diameters $\sim \lambda$. The high-focused irradiance (initially up to 10^{18} Wcm^{-2} , ultimately possibly up to 10^{30} Wcm^{-2} assuming focusing to λ dimensions) and high photon energy (ultimately exceeding 10 keV) will enable many novel experiments in high-energy density physics, biological imaging, and other fields.

Future possibilities for basic physics research include the production of Unruh radiation⁴ and electron-positron pair production. However, FELs are big, multiuser facilities of large capital and running cost. A review in 2008 led to the decision in the United Kingdom not to proceed with the 4GLS free-electron laser source producing extreme ultraviolet (EUV) radiation at $\lambda > 10$ nm.⁵ The review suggested that FELs are uniquely useful for harder x-ray experiments ($\lambda < 5$ nm), but that they have strong competitors from other laser sources in the EUV.

58.2 HIGH HARMONIC PRODUCTION

Experiments with 10 to 100-fs-duration visible and near-infrared laser pulses incident into inert gas targets such as Ne, Ar, or Xe targets have shown that harmonics are produced at wavelengths >4 nm with an efficiency $\sim 10^{-6}$ into each harmonic.^{6,7} Electrons are tunnel-ionised from gas atoms by the laser and then accelerated in the laser electric field. Provided linear polarisation is employed, the electron is accelerated back to the atom when the laser electric field reverses direction. Returning back to the atom, the electron recombines, emitting radiation up to the energy of the accelerated electron (\approx atom ionisation energy + $3.17 \times$ electron ponderomotive energy).⁸ Due to this harmonic production being closely associated with each oscillation of the driving laser electric field, the harmonics have to a good approximation the temporal duration and beam qualities of the optical laser creating them. The harmonics are each produced with approximately equal efficiency out to a maximum harmonic number (in the so-called “plateau” regime). The maximum harmonic photon energy $h\nu_{\max}$ in the plateau is

$$h\nu_{\max} \cong I_p + \left(3.17 \times 10^{-13} \frac{eV}{(W/cm^2)\mu m^2} \right) I_{\text{laser}} \lambda_{\mu m}^2 \quad (1)$$

where I_p is the atom ionization energy (in eV), I_{laser} is the laser irradiance (in Wcm^{-2}) and $\lambda_{\mu m}$ is the driving laser wavelength (in microns). For example, a focused laser irradiance of $3 \times 10^{14} Wcm^{-2}$ can be expected to produce strong harmonics up to $h\nu_{\max} \approx 100$ eV, equivalent to a wavelength of 12 nm.

Harmonics from the laser irradiation of solid targets have been recently demonstrated down to wavelengths ≈ 0.33 nm with 10^{-6} efficiency from 500-fs laser pulses.⁹ Such a development is exciting as the harmonic wavelength is in the x-ray regime. To achieve x-ray harmonics, a large laser system capable of focusing to $\sim 10^{21} Wcm^{-2}$ with high contrast ($>10^{10}$) is required. Much less stringent requirements are needed for the production of EUV harmonics. An intensity of $10^{19} Wcm^{-2}$ and contrast $> 10^{-6}$ is sufficient.¹⁰ Harmonics from solid surfaces are emitted in the specular reflection direction for the incident laser in a narrow cone much smaller than the incident laser cone size. As they arise from the relativistic oscillation by the incoming laser electric field of the vacuum-solid interface, the beam quality of the harmonic radiation depends strongly on the uniformity of the solid target irradiation. However, comprehensive measurements of solid target high harmonic beam quality have not yet been made.

58.3 PLASMA-BASED EUV LASERS

Until the late 1990s, plasma-based EUV lasers had only been pumped with high-energy laser pulses. For example, the record shortest wavelength EUV laser with saturated output was produced at 5.9 nm, but required ~ 60 to 120 J of energy per pulse.¹¹ It was not feasible to propose that such lasers could be used to produce EUV lasing for applications where high average power was required as it is not possible to have a high-repetition-rate laser of such energy per pulse. Developments in the last 10 years have completely changed this picture and have enabled demonstrations showing repetitively pulsed EUV laser action (10 to 50 nm) at high average power. Discharge pumped EUV laser action

at 46.9 nm was demonstrated with a capillary discharge,¹² but significantly shorter wavelength (<20 nm), high repetition rate lasing has been achieved with infrared laser pumping in plasmas.¹³

Repetitively pulsed (>10 Hz), short duration (<ps) infrared lasers capable of pumping EUV lasing from initially solid targets were developed in the late 1990s. Resulting EUV gain durations are short with <ps pumping pulses, so “traveling wave” pumping is needed whereby the pumping laser is incident along the ~3 to 10 mm target length coincident with the speed of propagation ($\approx c$) of the amplifying EUV beam. The efficiency of generating EUV lasing from solid targets works best with electron collisional excitation via laser pumping using two-pulse irradiation. A prepulse laser generates an expanding plasma. A main pulse heats, ionizes, and excites this plasma to generate gain via electron collisional excitation at some later time. The prepulse irradiation ensures that the volume of the gain region is large, that the main laser pulse is strongly absorbed, and that density gradients at the time of arrival of the main pulse are small so that refraction of the EUV lasing beam is reduced. A recent significant development, grazing incidence pumping,¹⁴ improves the pumping efficiency even more by offering a controlled way of achieving pumping at the optimum electron density in the plasma, as shown in Fig. 1. There is an optimum density for gain above which collisional de-excitation destroys the population inversion and below which the gain coefficient drops due to the decrease in density. In grazing incidence pumping, laser light penetrates to a turning point determined by the angle of incidence tuned to correspond to the optimum plasma density. In grazing-incidence pumping, there is also an inherent (close to velocity c) traveling wave excitation of gain. EUV lasing has been produced with laser pumping energies <1 J using the grazing incidence pumping technique.¹⁵

The development of plasma-based extreme ultraviolet (EUV) lasers has been a very successful investigation that initially started in the 1970s, but has matured to a state where saturated output has been achieved down to the record short wavelength of 5.9 nm and the atomic, plasma, and propagation physics are known to good accuracy (see Ref. 16). Plasma-based EUV lasers (10 to 20 nm) offer a route to the achievement of both high peak brightness (typically 10^{24} photons s^{-1} mm^{-2} $mrad^{-2}$ has been achieved) and high average power (e.g., 10^{14} photons s^{-1} mm^{-2} $mrad^{-2}$ under operation at 10 Hz). This compares to, for example, third generation synchrotron light sources that produce smaller to comparable peak brightness and comparable average brightness for a spectral bandwidth $\Delta\nu/\nu$ of 0.1% over similar spectral ranges. Plasma-based EUV laser output is extremely narrowband ($\Delta\nu/\nu < 10^{-4}$), so coherence lengths are typically 100 μm to 1 mm. The pulse energy in a plasma-based EUV laser is

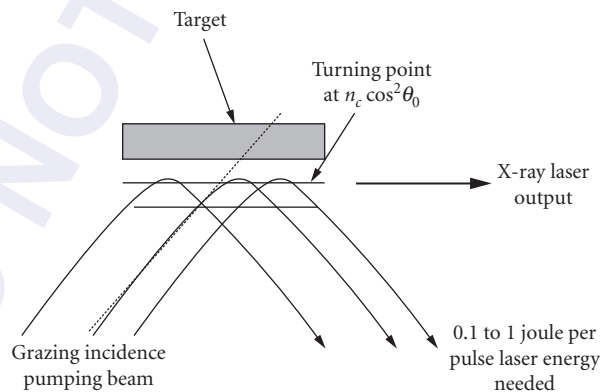


FIGURE 1 Schematic diagram illustrating grazing-incidence pumping into a preformed plasma. Heating occurs preferentially at the turning point of the beam (at electron density $n_c \cos^2 \theta_0$, where n_c is the electron critical density for the pump laser wavelength). The angle of incidence θ_0 is tuned so that the pump laser beam turning point corresponds to the optimum density for EUV laser gain (0.1 to 1 joule per pulse laser energy needed).

typically 1 μJ to 1 mJ, giving 0.01- to 10-mW average power (at rep rates of 10 Hz) and 1- to 1000-MW peak power.

An outstanding issue with plasma-based EUV lasers is the beam spatial coherence. The phase profile of these lasers needs to be improved considerably while maintaining their power output to enable their use for the high quality optics testing needed, for example, by the EUV lithography industry. At present plasma-based EUV lasers have high Fresnel number ($\sim 10^4$) with EUV lasers effectively composed of a large number (\approx Fresnel number) of “beamlets” that interfere to create output that is a complex speckle pattern.¹⁷ With short pulse (\sim ps) pumping, each beamlet results from a single spontaneous emission and the phase of each beamlet is approximately uniform transversely, but random with respect to other beamlets. The temporal duration of each beamlet is typically 3 ps and is close to the Fourier transform limit.¹⁸

58.4 REFERENCES

1. J. M. Madey, “Stimulated Emission of Bremsstrahlung in a Periodic Magnetic Field,” *J. Appl. Phys.* **42**:1906 (1971).
2. L. R. Elias, “Observation of Stimulated Emission Radiation by Relativistic Electrons in Spatially Periodic Transverse Magnetic Field,” *Phys. Rev. Lett.* **36**:717 (1976).
3. V. Ayvazyan, N. Baboi, J. Bähr, V. Balandin, B. Beutner, A. Brandt, I. Bohnet, et al., “First Operation of a Free-Electron Laser Generating GW Power Radiation at 32 nm Wavelength,” *Eur. J. Phys.* **D37**:207 (2006); see also <http://xfelinfo.desy.de/en/start/2/index.html>, accessed 6 May 2009.
4. P. Chen and T. Tajima, “Testing Unruh Radiation with Ultraintense Lasers,” *Phys. Rev. Lett.* **83**:256 (1999).
5. B. W. J. McNeil, N. R. Thompson, D. J. Dunning, J. G. Karssenber, P. J. M. van der Slot, and K-J Boller, “A Design for the Generation of Temporally Coherent Radiation Pulses in the VUV and beyond by a Self-Seeding High Gain Free Electron Laser Amplifier,” *New J. Phys.* **9**:82 (2007).
6. C. Altucci, R. Bruzese, C. de Lisio, M. Nisoli, G. Cerullo, S. Stagira, S. De Silvestri, et al., “Features of High-Order Harmonic Generation in the 30 fs and the Sub-10 fs Regimes,” *J. Opt.* **A2**:289 (2000).
7. T. Ditmire, J. K. Crane, H. Nguyen, L. B. DaSilva, and M. D. Perry, “Energy-yield and Conversion Efficiency Measurements of High Order Harmonic Radiation,” *Phys. Rev.* **A51**:R902 (1995).
8. P. Corkum, “Plasma Perspectives on Strong Field Multiphoton Ionization,” *Phys. Rev. Lett.* **71**:1994 (1993).
9. B. Dromey, S. Kar, C. Bellei, D. C. Carroll, R. J. Clarke, J. S. Green, S. Kneip, et al., “Bright Multi-keV Harmonic Generation from Relativistically Oscillating Plasma Surfaces,” *Phys. Rev. Lett.* **99**:085001 (2007).
10. P. A. Norreys, M. Zepf, S. Moustazis, A. P. Fewes, J. Zhang, P. Lee, M. Bakarezos, et al., “Efficient Extreme UV Harmonics Generated from Picosecond Laser Pulse Interactions with Solid Targets,” *Phys. Rev. Lett.* **76**:1832 (1996).
11. R. Smith, G. J. Tallents, J. Zhang, G. Eder, S. McCabe, G. J. Pert, and E. Wolfrum, “Saturation Behavior of Two X-ray Lasing Transitions in Ni-Like Dy,” *Phys. Rev.* **A59**:R47 (1999).
12. B. R. Benware, C. D. Macchietto, C. H. Moreno, and J. J. Rocca, “Demonstration of a High Average Power Tabletop Soft X-Ray Laser,” *Phys. Rev. Lett.* **81**:5805 (1998).
13. Y. Wang, M. A. Larotonda, B. M. Luther, D. Alessi, M. Berrill, V. N. Shlyaptsev, and J. J. Rocca, “Demonstration of High Repetition Rate Tabletop Soft X-Ray Lasers with Saturated Output at Wavelengths Down to 13.9 nm and Gain Down to 10.9 nm,” *Phys. Rev.* **A72**:053807 (2005).
14. R. Keenan, J. Dunn, P. K. Patel, D. F. Price, R. F. Smith, and V. N. Shlyaptsev, “High Repetition Rate Grazing Incidence Pumped X-Ray Laser Operating at 18.9 nm,” *Phys. Rev. Lett.* **94**:103901 (2005).
15. J. Tummler, K. A. Janulewicz, G. Priebe, and P. V. Nickles, “10-Hz Grazing-Incidence Pumped Ni-Like Mo X-Ray Laser,” *Phys. Rev.* **E72**:037401 (2005).
16. G. J. Tallents, “The Physics of Soft X-Ray Lasers Pumped by Electron Collisions in Laser Plasmas,” *J. Phys.* **D36**:R259 (2003).
17. O. Guilbaud, A. Klisnick, K. Cassou, S. Kazamias, D. Ros, G. Jamelot, D. Joyeux, and D. Phalippou, “Origin of Microstructures in Picosecond X-Ray Laser Beams,” *Europhys. Lett.* **74**:823 (2006).
18. P. Mistry, M. H. Edwards, and G. J. Tallents, “X-Ray Laser Pulses at the Fourier Transform Limit,” *Phys. Rev.* **A75**:013818 (2006).

INVERSE COMPTON X-RAY SOURCES

Frank Carroll

*MXISystems
Nashville, Tennessee*

59.1 INTRODUCTION

Beams of pulsed, tunable, monochromatic hard x rays have been a long sought after goal. Such “dream beams,” that are tunable from energies as low as 8 keV to well into the *hundreds of keV*, have now become available in compact units that can rival synchrotrons in many ways, delivering fluxes of photons of varying bandwidths that approach or exceed the output of those much larger facilities at some energies. Additionally, these monochromatic beams of x rays can be produced in cone beam geometries that are very useful for covering larger areas, even as large as humans.

All of this has come about because we have learned to harness a process called inverse Compton scattering (Thomson backscattering). In normal Compton scattering, a photon incoherently scattering from a nearly stationary electron exchanges some of its energy and momentum with the electron, resulting in a lower energy photon. In inverse Compton scattering, the electron is traveling at high velocity, and the photon gains energy. In one embodiment of this process, an electron beam is accelerated to near-relativistic energies typically between 12 and 50 MeV and then focused down to a focal spot that is anywhere from 3 to 100 micron in size in an area designated the interaction zone (IZ). In like fashion, an intense (terawatt) laser beam is also focused down to a similar sized focal spot and counterpropagated against the packet of electrons in a head-to-head collision (180° geometry), as shown in Fig. 1.

The Rayleigh ranges of the beams are aligned so that they completely overlap in the IZ and such that the packets of light and electrons both reach the IZ at the same instant. In that collision, the laser photons are Doppler shifted by the inverse Compton process to x-ray energies. Hence a light photon goes in and an x-ray photon comes out. X rays are generated in a somewhat slowly diverging cone beam along the axis and in the direction traveled by the electron beam, as shown in Fig. 1. These x rays typically exit the machine through a beryllium vacuum window for use in various applications. Since the accelerator can be tuned, the x rays emanating from the machine are tunable. Since all of the electrons are not at exactly the same energy, some reduction in monochromaticity of the x rays is seen, but since the laser light is nearly monochromatic, the x rays, on the whole, are nearly monochromatic.¹⁻⁵

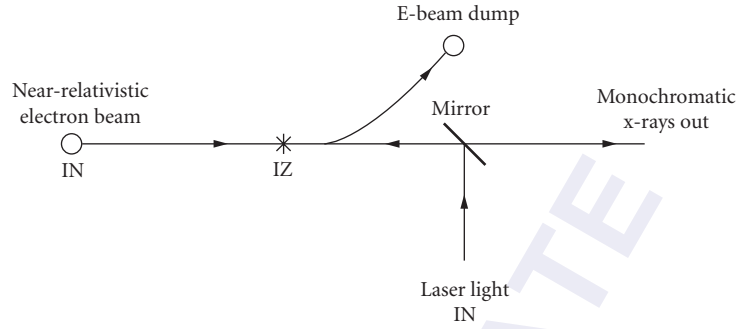


FIGURE 1 Inverse Compton process in practice.

59.2 INVERSE COMPTON CALCULATIONS

The theoretical x-ray yields from such sources can be computed directly from the Thomson scattering cross section and some basic beam-geometry and flux considerations. The Thomson cross-section σ_T is

$$\sigma_T = \frac{8\pi}{3} r_e^2 = 6.652 \times 10^{-29} \text{ m}^2 \quad (1)$$

where r_e is the classical electron radius. If one uses a diffraction-limited optical beam and an electron beam that is assumed to be smaller than the optical spot size, then this gives an x-ray yield of

$$N_x(\Delta V_x) = \frac{QeU_L\sigma_T}{hcqZ_r\Delta V_x} \quad (2)$$

where Qe is the charge in an electron bunch, U_L is the energy (in joules) in a laser pulse, h is Planck's constant, c is the speed of light, q is the electron charge, Z_r is the Rayleigh range of the laser focal spot (which is matched to the electron beam focal parameters for a given desired ΔV_x), and ΔV_x is the fractional energy spread of the x-ray output which is useful to the application. For a system with 10 J of laser light, 1 nC of electron charge, and an $f/10$ final focus of the laser, the output is about 5×10^9 x rays into a 10% bandwidth. The brilliance, is $7 \times 10^{31} [\text{m}^2 \cdot \text{sr} \cdot \text{s} \cdot (10\% \text{ bandwidth})]^{-1}$ assuming a 10-ps output pulse.

Modifications to such machines can include using light of a shorter wavelength, thereby keeping gradients within the accelerator smaller. In addition, lasers used for these machines can be run at 20 Hz instead of 10 Hz, yielding almost double the light and hence greater x-ray flux.

An additional and extremely important enhancement comes about by using smaller electron beam focal spots at the interaction zone, which, when using a higher current of electrons, increases the effective cross section for interaction with the light photons, creating significantly more x rays as well.

When all of these factors are used together, one gains anywhere from one to three orders of magnitude in the output flux of the x rays.⁶⁻⁸

59.3 PRACTICAL DEVICES

Synchrotrons are, of course, broadband white radiation sources (see Chap. 55), which use monochromators to deliver beams of narrow bandwidth (see Chap. 39) to multiple user beamlines. Inverse Compton sources typically deliver bandwidths of 0.1 to 10%, and can be further monochromatized to narrower bandwidth using various x-ray optics as well,⁹⁻¹² but this diminishes their brightness relative to the typical synchrotron source. However, Compton sources are much more compact,

TABLE 1 Comparison of Synchrotrons to Available Compton Sources

	(Generation 2) Machine*	(Generation 3) Compton Source†	Typical Synchrotrons‡
E-beam:	Linac running in “single pulse” mode Up to 75 MeV 0.5 nanocoulombs/pulse Emittance $\cong 3.00\pi$ mm-mrad Copper photocathode	Same Same 2–3 nanocoulombs/pulse Emittance $\cong 0.7\pi$ mm-mrad Copper-Mg photocathode	Synchrotron ring 0.1–7.0 GeV 50–5200 mA current Variable N/A
Laser:	Tabletop terawatt 1054-nm IR light 20 J/5 ps pulse Once every 5 min	Tabletop terawatt 532-nm green light 1.5 J/5 ps pulse 10–20 Hz rep rate	N/A N/A N/A N/A
X-ray beam:	$10^{10}/5$ ps shot 3×10^7 photons/s Tunable from 12–50 keV Change energy 30 min to 4 h 0.1–10% bandwidth Conebeam geometry (large areas covered in a single shot) Effective focal spot 100–300 μ	10^9 – 10^{10} photons/5 ps shot $2.23 \times 10^{11-12}$ photons/s Tunable from 10–100 keV Change energy in 1/10th of a second Same Same but focusable to 100 μ (round beam) or 300 μ in a square beam, if desired Effective focal spot 3–6 μ	Quasi-CW Typically $1.4 \times 10^{13-16}$ photons/s Variable from 1–57 keV Variable 1×10^{-4} $100 \times 200 \mu$ Effectively collimated
Shielding vault:	None	None	Complex

*Inverse Compton source operational at the Vanderbilt University Medical Free-Electron Laser Facility, Nashville, TN since April 2001.^{4,5}

†MXISystems, Inc., Inverse Compton source, 2008.¹⁹

‡Advanced Design Consulting USA, Inc., Synchrotron Primer—World’s Light Sources. Section 4-2.¹³

affordable, and available. A comparison of typical synchrotron beamlines to two published compact Compton sources is shown in Table 1.

While inverse Compton x-ray sources have been proposed by many over the years, there have only been a few devices successfully built and operated. (Ting 500 eV,¹⁴ Carroll 14 to 18 keV,^{1,2} Ruth 10 to 35 keV,^{10,15} Carroll 12 to 53 keV,^{3,4} Gibson/Pelaides 70 keV,¹⁶ and more recently Barty/T-Rex¹⁷ at 776 keV.) Some of the more compact devices are now being commercialized for both medical and nonmedical purposes in the ranges from 8 to 120 keV (MXISystems, Inc.)^{18,19} and 10 to 35 keV (Lyncean Technologies, Inc.)^{20,21}

59.4 APPLICATIONS

With the advent of this enabling technology come broad-ranging applications. In medicine these uses encompass both the diagnostic and therapeutic arenas.

Due to inherent contrast differences between tissues in the human body, monochromatic x-ray beams can make cancerous tissues stand out from normal tissues without the use of contrast materials, while at the same time lowering radiation dose significantly to the patient.^{22,23} This is particularly true in studies such as mammography, where 3-D compressionless studies become feasible, allowing the examination to be performed more accurately and with far less discomfort to the patient while requiring anywhere from 5 to 60 times less radiation dose depending on how the study is performed (see Chap. 31).

The extremely small effective focal spot size at the interaction zone results in sufficient lateral coherence to allow one to take advantage of the disparate refractive indices of various tissues and use the diffractive effects produced at edges within tissues to perform phase contrast imaging (see Chap. 27). This type of imaging may be capable of delivering 100 to 1000 times the information than that achieved using absorption imaging (what we have done for the past 100 years).^{24–30}

By detecting and back projecting low angle scatter emanating from an irradiated tissue or organ, one can conceivably perform noninvasive biopsies of tissues deep within the body of the patient, discerning cancerous from noncancerous tissues in a painless fashion.³¹

Therapeutically, monochromatic x rays can be tuned to the binding energies of the innermost electron shells of any atom in the periodic chart, thereby knocking K- or L-shell electrons from their orbits, causing cascades of energy (the Auger cascade) concentrating all of this energy within nanometers of the atom. By attaching a drug containing a heavy metal target (such as Pt, Gd, or I) to the DNA of cancerous cells, one can stimulate Auger cascades within the DNA causing double-stranded breaks that for the most part will not heal. This translates into treatment of cancers with 3 to 5 times less radiation than is now used by the best techniques.^{32,33}

59.5 INDUSTRIAL/MILITARY/CRYSTALLOGRAPHIC USES

These narrow bandwidth, pulsed x rays (with pulse lengths of 5 to 10 ps) are produced through widely tunable energy ranges. These are extremely useful in evaluation of defects in newer carbon composite airplane parts,^{34,35} spacecraft parts and subsystems, examination of corrosion and cracks in our aging aluminum aircraft fleet, internal analysis of turbines at full power, performance of protein crystallography,³⁶ and a host of other nondestructive testing applications.

59.6 REFERENCES

1. F. E. Carroll, J. W. Waters, R. R. Price, C. A. Brau, C. F. Roos, N. H. Tolk, D. R. Pickens, and W. H. Stephens, "Near-Monochromatic X-Ray Beams Produced by the Free Electron Laser and Compton Backscatter," *Invest. Radiol.* **25**:465–471, 1990.
2. C. Pellegrini, "On Some Methods of X-Ray Production from Relativistic Electron Beams," *J. X-Ray Sci. Technol.* **4**:275–289, 1994.
3. F. E. Carroll, J. W. Waters, R. H. Traeger, M. H. Mendenhall, W. W. Clark, and D. A. Brau, "Production of Tunable, Monochromatic X-Rays by the Vanderbilt Free-Electron Laser," *Proceedings SPIE. Free Electron Laser Challenges II* **3614**:139–146, 1999.
4. F. E. Carroll, "Perspectives: Tunable, Monochromatic X-Rays: A New Paradigm in Medicine," *AJR* **179**: 583–590, 2002.
5. F. E. Carroll, M. H. Mendenhall, R. H. Traeger, C. Brau, and J. W. Waters, "Pulsed, Tunable, Monochromatic X-Rays from a Compact Source: New Opportunities," *AJR* **181**:1197–1202, 2003.
6. W. J. Brown, S. G. Anderson, C. P. J. Barty, S. M. Betts, R. Booth, J. K. Crane, R. R. Cross, et al., "Experimental Characterization of an Ultrafast Thomson Scattering X-Ray Source with Three-Dimensional Time and Frequency-Domain Analysis," *Phys. Rev. ST Accel. Beams.* **7**:060702–12, 2004.
7. F. V. Hartemann, W. J. Brown, D. J. Gibson, S. G. Anderson, A. M. Tremaine, P. T. Springer, A. J. Wooton, E. P. Hartouni, and C. P. J. Barty, "High Energy Scaling of Compton Light Sources," *Phys. Rev. ST Accel. Beams* **8**:100702–17, 2005.
8. C. A. Brau, "Oscillations in the Spectrum of Nonlinear Thomson-Backscattered Radiation," *Phys. Rev. ST Accel. Beams* **7**:020701–9, 2004.
9. M. A. Kumakhov and F. F. Komarov, "Multiple Reflection from Surface X-Ray Optics," *Phys. Rep.* **191**:289, 1990.
10. P. A. Tompkins, C. C. Abreau, F. E. Carroll, Q. F. Xiao, and C. A. MacDonald, "Use of Capillary Optics as a Beam Intensifier for a Compton X-Ray Source," *Med. Phys.* **21**:1777–1784, 1994.

11. www.rigako.com/optics/index.htm, accessed July 16, 2009.
12. F. R. Sugiuro, D. Li, and C. A. MacDonald, "Beam Collimation with Polycapillary X-Ray Optics for High Contrast High Resolution Monochromatic Imaging," *Med. Phys.* **31**(12):3288–3297, 2004.
13. www.adc9001.com/index.php?src=primers_synchrotron, accessed Sept. 16, 2008.
14. A. Ting, R. Fischer, A. Fisher, C. I. Moore, B. Hafizi, R. Elton, K. Krushelnick, et al., "Demonstration Experiment of a Laser Synchrotron Source for Tunable, Monochromatic X-Rays at 500 eV," *Nucl. Instrum. Meth. Phys. Res. A* **375**:ABS68–70, 1996.
15. Z. Huang and R. D. Ruth, "Laser-Electron Storage Ring," *Phys. Rev. Lett.* **80**(5):976–9, 1998.
16. D. J. Gibson, S. G. Anderson, C. P. J. Barty, S. M. Betts, R. Booth, W. J. Brown, J. K. Crane, et al., "PLEIADES: A Picosecond Compton Scattering X-Ray Source for Advanced Backlighting and Time-Resolved Material Studies," *Phys. Plasmas* **11**:2857 (2004).
17. Internal LLNL communication. https://publicaffairs.llnl.gov/news/news_releases/2008/NR-08-04-03.html, accessed July 16, 2009.
18. US Patents 6,332,017 B1, 2001 and 6,687,333 B2, 2004. System and Method for Producing Pulsed Monochromatic X-Rays.
19. MXISystems, Inc. www.mxisystems.com, accessed July 16, 2009.
20. International Patent WO 101926 A2, 2005 and US Patents 5,825,847, 1998 and 6,035,015, 2000.
21. Lyncean. Inc. www.lynceantech.com, accessed July 16, 2009.
22. P. C. Johns and M. J. Yaffe, "X-Ray Characterization of Normal and Neoplastic Breast Tissues," *Phys. Med. Biol.* **32**:675–695, 1987.
23. F. E. Carroll, J. W. Waters, W. W. Andrews, R. R. Price, D. R. Pickens, R. Willcott, P. Tompkins, et al., "Attenuation of Monochromatic X-Rays by Normal and Abnormal Breast Tissues," *Invest. Radiol.* **29**:266–272, 1994.
24. T. Takeda, A. Momose, Y. Itai, J. Wu, and K. Hirano, "Phase-Contrast Imaging with Synchrotron X-Rays for Detecting Cancer Lesions: Preliminary Investigation," *Acad. Radiol.* **2**:799–803, 1995.
25. D. Chapman, W. Thomlinson, R. E. Johnston, et al., "Diffraction Enhanced X-Ray Imaging," *Phys. Med. Bio.* **42**:2015–2025, 1997.
26. U. Kleuker, P. Suortti, W. Weyrich, and P. Spanne, "Feasibility Study of X-Ray Diffraction Computed Tomography for Medical Imaging," *Phys. Med. Biol.* **43**:2911–2923, 1998.
27. E. D. Pisano, E. R. Johnston, D. Chapman, G. Iacocca, C. A. Livasy, D. B. Washburn, D. E. Sayers, et al., "Human Breast Cancer Specimens: Diffraction-Enhanced Imaging with Histologic Correlation-Improved Conspicuity of Lesion Detail Compared with Digital Radiography," *Radiology* **214**:895–901, 2000.
28. E. F. Donnelly and R. R. Price, "Quantification of the Effect of KVP on Edge-Enhancement Index in Phase-Contrast Radiography," *Med. Phys.* **29**(6):999–1002, 2002.
29. E. F. Donnelly, R. R. Price, R. David, and D. R. Pickens, "Dual Focal-Spot Imaging for Phase Extraction in Phase-Contrast Radiography," *Med. Phys.* **30**(9), September 2003.
30. A. Momose, "Phase-Sensitive Imaging and Phase Tomography Using X-Ray Interferometers," *Opt. Exp.* **11**(19):2303–2314, 2003.
31. P. C. Johns, R. J. Leclair, and M. P. Wismayer, "Medical X-Ray Imaging with Scattered Photons," *SPIE Reg. Meet. Optoelectron. Photon. Imaging SPIE* **TD01**:355–357, 2002.
32. J. P. Pignol, E. Rakovitch, D. Beachey, and C. L. Sech, "Clinical Significance of Atomic Inner Shell Ionization (ISI) and Auger Cascade for Radiosensitization Using IUdR, BUdR, Platinum Salts, or Gadolinium Porphyrin Compounds," *Int. J. Radiat. Oncol. Biol. Phys.* **55**(4):1082–1091, 2003.
33. F. E. Carroll, "Tunable Monochromatic X-Rays, an Enabling Technology for Molecular/Cellular Imaging and Therapy," *J. Cell. Biochem.* **90**:502–508, 2003.
34. J. Martin-Herrero and Ch. Germain, "Microstructure Reconstruction of Fibrous C/C Composites from X-Ray Microtomography," *Carbon* **45**:1242–1253, 2007.
35. P. M. Pawar and R. Ganguli, "On the Effect of Progressive Damage on Composite Helicopter Rotor System Behavior," *Comp. Struc.* **78**:410–423, 2007.
36. F. V. Hartemann, H. A. Baldis, A. K. Kerman, A. Le Fol, N. C. Luhmann, and B. Rupp, "Three-Dimensional Theory of Emittance in Compton Scattering and X-Ray Protein Crystallography," *Phys. Rev. E* **64**:016501, 2001.

This page intentionally left blank.

DO NOT DUPLICATE

SUBPART

5.5

X-RAY DETECTORS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

INTRODUCTION TO X-RAY DETECTORS

Walter Gibson*

*X-Ray Optical Systems, Inc.
Rensselaer, New York*

Peter Siddons

*Brookhaven National Laboratory
Upton, New York*

60.1 INTRODUCTION

X rays can interact with gases, liquids, or solids to produce electronic, atomic, molecular, or collective excitations. Measurement of these excitations can be used to determine the x-ray intensity, energy, and position. A large variety of phenomena have been used for such detection. These include ionization of atoms and molecules, production of electrons and holes in semiconductors, excitation of both short-lived and metastable electronic states in liquids and solids, and excitation of phonons and Cooper pairs in superconductors. Selection of the detector or detectors to be used for a given application depends on the particular requirements and constraints inherent in that application, and will typically involve such factors as intensity, energy resolution, position resolution, size, convenience, and cost. To help with consideration of the kinds of trade-offs involved in selection of a particular detector, and, in the context of this *Handbook*, the factors that influence the optimization of a particular optic-detector system, this chapter summarizes the detector choices in terms of their type and function. There are no comprehensive reviews specifically devoted to x-ray detectors. Most books cover detectors for several types of ionizing radiation.¹⁻⁴ A useful summary specifically devoted to x-ray detectors is given by Buckley,⁵ and a summary of x-ray detectors for astronomy is given by Fraser.⁶

60.2 DETECTOR TYPE

Ionization

Ionization Chamber Ionization chambers contain a gas in a region of high electric field, usually produced with flat plates. X rays incident on a gas can excite bound electrons, resulting in their separation from the atom with energy equal to the difference between the absorbed x-ray energy and the binding energy of the electron. The emitted electron can, if it has sufficient energy, result in further ionization.

*This volume is dedicated in memory of Walter Gibson.

The average x-ray (actually secondary electron, because these produce most of the ionization) energy loss to produce an electron-ion pair depends on the gas, but it is typically about 30 eV (electronvolts). This is considerably in excess of the energy needed to ionize the atom, with the difference going into dissociation and excitation of molecules. The electrons and positive ions can then be separated by an applied electric field. The resulting charge current in the external circuit is a measure of the number of electron-ion pairs produced. Because the electrons move much faster than the heavy positive ions, the short-term charge pulse is primarily a measure of the electron motion. Ionization chambers have frequently been used as spectrometers for energetic ions, such as alpha particles or fission fragments, because they have shorter range than x rays, permitting an electrode configuration that allows only electron motion to be measured.

The longer range of x rays precludes easy separation of electron and ion motion in a simple two-plate chamber. The addition of screening grids can allow such a separation, but the typical applications of these devices do not usually justify the additional complexity. Also, the statistical fluctuations are larger than in some other spectrometers, so ion chambers are rarely used to measure x-ray energy distributions. However, ion chambers are the basis of frequently used radiation-monitoring devices, such as the pencil-type electrometers that are used as personal monitors. Another widespread application is to measure the x-ray beam intensity in synchrotron-based experiments. In this application, the ion chamber can be relatively thin (absorbing only a very small but proportional fraction of the incident x rays), and the resulting current flow in the external circuit is used to measure the x-ray intensity. For example, an ion chamber before a sample and one after can be used to measure the x-ray absorption in the sample as in x-ray absorption fine structure (EXAFS)^{7,8} and x-ray absorption near-edge structure (XANES)⁹ experiments (see Chap. 30). The gas in an ion chamber can just be air at atmospheric pressure as in simple monitoring devices. Frequently, when a more consistent measurement is desirable, a rare gas such as helium or argon is used to suppress electron-ion recombination, often with added methane or other gas to increase the electron drift velocity. In the EXAFS applications mentioned here, it is common practice to mix the measurement rare gas with helium in varying fractions in order to tailor the absorbed radiation fraction for optimal measurement precision. Clearly, one does not want to have all of the radiation absorbed by the first of the two chambers. A general discussion of ionization chambers can be found in reviews of radiation detectors.¹⁻⁴

Proportional Counter If the electric field gradient in an ionization chamber is increased high enough, even low-energy electrons released during the x-ray energy loss process can be accelerated enough to produce further ionization. Although, for charged particles or electrons, a parallel plate geometry is frequently used, for x rays, thin wires are often used as the positive electrode so that the acceleration takes place in a relatively confined space, close to the anode. Over a relatively wide applied voltage range, the charge multiplication increases with applied field and is constant at a given applied field. This multiplication results in charge amplification that is fast, low noise, and because it takes place close to the anode, provides effective separation from the ion current. The multiplication factor can be large, typically in the range between 10^4 and 10^6 . The multiplication gain can be stabilized and made less sensitive to the bias voltage as well as the particle energy or counting rate by addition of a small amount of a polyatomic gas to the rare gas. A typical gas mixture is 10 percent methane and 90 percent argon, known as "P-10 gas." Carbon dioxide is also a common quench gas, and has the added benefit that it is much less prone to polymerization than methane. Polymerized methane can collect on the fine wires and cause severe localized gain reduction. This is particularly undesirable in a position-sensitive detector since it can distort the position linearity. Proportional counters have been extensively used as x-ray spectrometers. They are particularly useful for low-energy x rays because the intrinsic detector amplification reduces the effects of external amplifier noise, and because they can be used in a windowless mode. Measurements have been made for energies less than 1 keV (kiloelectronvolt).¹⁰ At high x-ray energies, the sensitivity is limited by the reduced absorption of x rays in the gas, although the energy discrimination for even high-energy x rays that are absorbed remains good so long as the secondary electrons are stopped in the detector. The effective energy range is dependent on the gas pressure and type, the detector size and shape, and on the presence of applied magnetic or electric fields, which can cause the secondary electrons

to adopt spiral paths. X rays with energy as high as 100 keV have been measured with special detector designs.¹¹

One of the most useful applications of proportional counters for x rays is position determination.¹² If the charge current in an anode wire is collected from both ends of the wire, resistive separation in the wire can be used to determine the position along the wire at which the charge is collected. Position resolution of fractions of millimeters can be achieved this way. Furthermore, if a number of parallel anode wires are used, the relative charge collection on adjacent wires can be used to determine the position of the x-ray absorption between the wires. In this way, two-dimensional position determinations can be achieved. At the same time, energy discrimination can be obtained by combining the charge signal for a given event from both ends of the anode wire and from adjacent wires. This has been used for a number of applications, notably measurement of x-ray diffraction distributions and imaging of astrophysical x-ray sources.¹³ Another popular method of encoding the position is the delay-line technique. In this method, the anode avalanche induces charge in a segmented cathode which is connected to the nodes of a lumped-element L-C delay line. Measurement of the arrival time of the induced charge pulse at either end of the delay line provides a measurement of the position of the event along the wire. This method avoids the need to fabricate the rather delicate resistive wire needed in the charge-division method, and such detectors are therefore more robust. This same method can also be used in the second dimension, encoding the multiple anode wires described above. A third technique uses the same interpolation of partial charges induced on the segmented cathode as described earlier for the multi-anode 2-D detector. In either plane, this technique involves significantly more complexity than the delay-line method.

Geiger Counter If the electric field in an x-ray proportional counter is increased beyond the proportional region (often by making the diameter of the anode wire smaller), the electron multiplication increases catastrophically, resulting in large-current pulses that are no longer proportional to the x-ray energy. This is called the *Geiger region*. In fact, the multiplication is large enough that it is usually necessary to include components in the gas mixture that will quench the resulting electrical discharge. One benefit of this mode of detection is the large amplification that makes it possible to detect x rays with simplified electronics. This is why the most common use of such detectors is in portable or remote radiation monitors.

Semiconductor Detector In gases, x-ray absorption can produce electrons and positive ions that can be separated and measured by an applied electric field. In solids, the situation is somewhat different. In a conductor, the electrons that are excited by absorption of x rays are quickly dissipated in energy and lost among the abundant conduction electrons in the material, and the positive charge is almost immediately neutralized by conduction electrons. Furthermore, an electric field cannot be maintained in the material. In an insulator, it is easy to establish an electric field to sweep out excited electrons. However, the positive charge remains behind, trapped and immobile. The electric field due to this trapped charge quickly cancels the applied field. As a result, initial charge pulses from incident x rays decrease in size and intensity as a result of such polarization effects.¹⁴ However, for semiconductors, the situation is much more favorable. By doping with an appropriate impurity or formation of a surface barrier, a rectifying *pn* junction can be formed that allows an electric field to be established in the solid. The thickness of the field region, called the *depletion layer* because it is free of mobile charge carriers, depends on the applied voltage and on the resistivity of the bulk semiconductor.

X rays absorbed within the depletion layer ionize bound electrons with energy equal to the difference between the absorbed x ray and the binding energy of the electron. These energetic electrons are free to move in the solid and produce further ionization by interaction with bound electrons. In this sense, the situation is similar to x-ray absorption in ionization chambers. Indeed, semiconductor detectors are sometimes referred to as *solid-state ion chambers*. However, there is an important difference. Instead of the ionization producing a heavy, slow moving positive ion as in a gas, or an immobile positive charge as in an insulator, it produces a positively charged "hole." The hole is a vacancy in the electronic valence band of the solid, with low effective mass and high mobility, closer to that of free or conduction electrons than to positive ions. The negatively charged electrons and

positively charged holes are swept apart by the electric field in the depletion layer, producing a current pulse in the external circuit. Contrary to the ion chamber, both the electron and hole motion contribute to the short-term charge pulse. The amplitude of the charge pulse is proportional to the energy of the absorbed x ray. The energy resolution depends on the statistics of the ionization process and on the electronic noise in the associated electronics, which depends in part on the capacitance of the semiconductor *pn* junction. Detailed reviews of semiconductor detectors are given in Refs. 1, 15, and 16.

To increase the depletion layer thickness and therefore the sensitive volume of the detector, very high-resistivity semiconductor material is required. For silicon, this is achieved by drifting lithium ions through the diode at an elevated temperature in the presence of an electric field.^{17–19} The lithium compensates bound impurities in *p*-type silicon. Such detectors are called *lithium-drifted silicon*, or Si(Li), detectors. An alternative is to use very high-purity germanium (HpGe). Both Si(Li) and HpGe detectors are cooled (typically by liquid nitrogen) to reduce electronic noise due to thermally generated electrons and holes [and for Si(Li) to prevent out diffusion of the compensating lithium ions]. Typically, the energy resolution for such detectors is between 120 and 140 eV.

The availability of high-resistivity silicon, grown using the floating zone technique, has made it possible to fabricate detectors using planar fabrication techniques taken from the microchip industry, which perform almost as well as the Si(Li) devices, but without requiring cryogenics.²⁰ Good performance is achieved using Peltier coolers to bring the detector to around -30°C .

This same technology has allowed a new generation of devices based on manipulating the photogenerated charges within the silicon bulk to bring the charge produced over a large area to a low-capacitance collection point several millimeters away from the point of generation. This provides a large-area detector with a small readout capacitance, and hence good collection solid angle while maintaining low electronic noise. These devices are called silicon drift detectors (SDDs).²¹ Single units are commercially available and compete strongly with Si(Li) for low to moderate energies. Since they are low-capacitance devices they can handle large count rates and maintain good resolution. Arrays of such devices are becoming available and should soon become ubiquitous in certain applications.

High electric fields in solid state devices can lead to avalanche charge multiplication in a similar manner to the gas proportional counter. In silicon, these devices are called avalanche photodiodes (APDs).²² Gains of around 100 are possible, and the high fields lead to fast pulses, of order 1 nanosecond. Consequently, they have a high count rate capability, several 10s of megahertz. They have some energy resolution, around 20 percent.

Wide bandgap semiconductors, such as CdZnTe or CdTe, have been used because of their stronger absorption and the possibility of room temperature operation. Such materials frequently exhibit trapping of minority carriers (typically holes) and therefore can show reduced resolution and sometimes polarization effects. They are of particular interest for use in arrays (typically by segmentation of one of the electrodes) for medical or astrophysical applications.^{23,24}

Channel Electron Multipliers X rays and visible photons can be detected by photocathodes and electron multipliers. Photoelectrons that are released from thin films or surfaces, accelerated, and electrostatically focused on another surface produce increasing numbers of secondary electrons as they move from surface to surface, until they are collected on an anode and measured as a charge pulse in an electronic circuit. The first such multipliers, similar to photomultiplier tubes, were introduced in the early 1960s. These were largely replaced with simpler and more stable semiconducting glass or ceramic tubes with a voltage applied between the input and output of the tube. These can be thought of as continuous-dynode electron multipliers. Modern channel electron multipliers (CEM) are made from lead glasses, which contain PbO, SiO₂, and several weight percent of alkali metal oxides. These are light, compact, and quite rugged and are used in x-ray pulse-counting detectors in applications ranging from electron microscopy to x-ray astronomy. The electron gain depends on the secondary electron yield of the surfaces employed, the applied voltage, and the channel diameter and length. If gas is present in the channel, ionization can take place, resulting in further electron multiplication and higher apparent gain. However, ionization

can lead to feedback resulting from positive ions accelerating back up the channel, leading to long-term “ringing” of the signal due to electrons released by accelerated ions. Ion feedback can be suppressed if the channel is operated at very high vacuum. A more convenient solution is to curve the channel so that positive ions do not accelerate to high-enough velocity to produce secondary electrons before they strike the channel wall. Gain saturation results when the supply of charge to the channel wall lags the depletion rate by secondary electron emission. In practice, operational gains of from 10^4 to 10^6 are possible.⁴

Very high counting rates, and especially imaging of x rays, are possible by the use of microchannel plates (MCP), which are arrays of very fine bore CEMs, typically with channel diameters and channel separations of a few 10s of micrometers.²⁵ Curving the channels in MCPs to suppress ion feedback is difficult and expensive, so two (or more) tilted channel MCPs are put together to form “chevron” multipliers in modern detectors. Energy resolution in CEM detectors is poor, and their efficiency decreases with increasing x-ray energy. Their simplicity, low cost, and, for MCPs, high spatial resolution make them very useful for many applications. Microchannel plates are also used as electron multipliers for low-noise amplification when used with scintillation detectors.

Scintillation

Incident x rays can excite electrons in atoms or molecules to excited electronic states, which then deexcite by returning to their ground state with accompanying photon emission. The excitation can be caused directly by the x ray or, more usually, by secondary energetic electrons excited by the x ray. The emitted photons (usually in the optical region of the spectrum) can be measured by film and also by photomultiplier tubes or by other electron multipliers such as MCPs. Various plastics, liquids, organic crystals, inorganic crystals, and gases can act as scintillators.²⁶ Polycrystalline solid light-emitting materials are sometimes called *phosphors*. The deexcitation (and therefore, light emission) can be almost instantaneous or delayed. Both types are used for x-ray detectors.

Activated Phosphor The most used scintillator detectors that also give the photon energy spectrum involve scintillating crystals with rapid response. These are often single crystals of sodium or cesium iodide with thallium added to activate the photo response and shift emission into the visible range. Because alkali halides are frequently hygroscopic, they must be sealed. Scintillators are almost never used for measurement of x rays with energy less than about 1 keV. The scintillating crystal is commonly mounted directly onto the face of the vacuum photomultiplier tube, which provides photoelectron conversion, and low-noise amplification of the electrons. The conversion efficiency, η , which is the fraction of incident energy that appears as scintillation, is typically about 5 percent for CsI(Tl) and 12 percent for NaI(Tl).²⁷ The quantum detection efficiency (QDE) is defined as the number of photoelectrons produced per incident photon and depends on absorption in the cladding, the thickness and absorption coefficient of the scintillator, and the nature of the photocathode of the photomultiplier tube.⁵ The energy resolution of a well-matched NaI(Tl) system is $\Delta E/E \sim 1.7E^{-1/2}$. Activated alkali halide scintillation detectors have working ranges of a few kiloelectronvolts to a few megaelectronvolts. They have good dynamic range with upper detection rates limited by the scintillator decay time, which is approximately 1 μ s for NaI(Tl). Organic scintillators have faster (nanosecond) decay times but have lower efficiency and energy resolution. Scintillation detectors are convenient to use, can be large (several square centimeters), are readily available commercially, and are used for a large number of x-ray applications.

Scintillators are often used for x-ray imaging applications. In such cases, energy discrimination is often not important, and other phosphors, such as rare-earth oxides activated with terbium [e.g., $Gd_2O_2S(Tb)$, $La_2O_2(Tb)$, $Y_2O_2(Tb)$] or organic phosphors [e.g., tetraphenyl butadiene (TPB)] are used. The phosphor is deposited as a thin layer (1 to 10 μ m) directly on a glass plate, which is viewed by a sensitive television camera, or directly on the face of a fused-glass fiber optic, which is tapered to allow the intensity distribution to match the size and shape of a charge-coupled (CCD) or charge-injected (CID) imaging camera. This type of scintillator-fiber-optic-CCD imaging detector has become the standard for synchrotron-based protein crystallography beamlines, since it

produces a digital image which is amenable to computer analysis. At the time of writing, however, it is becoming clear that large area direct-detection semiconductor detectors will soon overtake these CCD-based detectors for this application, since they provide better point-spread function and faster readout.²⁸ Fast readout becomes important when the exposure time is much less than the detector readout time.

Restimulable Phosphor In the last decade, a different type of scintillation detector, commonly referred to as an imaging or computed radiography plate, has become widely used for x-ray imaging.^{29,30} This makes use of phosphor materials in which the electron excitation is to a metastable state that can have long (minutes to hours) decay time. After exposure, the plate is scanned by a laser beam that stimulates the metastable excited states to deexcite, producing light that is detected by a photomultiplier tube. The time of the light emission gives the position of the x-ray exposure and the intensity of the light, the x-ray intensity. Such detectors have a number of useful features:

1. They produce a digital record that leads conveniently to computer processing, transmission, and storage.
2. They have a linear response (as opposed to film) with a large dynamic range.
3. The thickness and nature of the phosphor used can be varied to accommodate the need (e.g., low-energy x rays, high-energy x rays, neutrons, etc.), and the size of the plate can be adapted to existing systems. For example, such plates are frequently used to replace photographic film.
4. The photostimulable plates can be reused by “erasing” the stored image by exposure to intense white light for a few seconds. At the present time, a restimulable scanning plate system (including exposure plates, scanner, and eraser) can be less expensive and more flexible than other commercially available position-sensitive detectors, although some with automatic readout features are comparable or more expensive. Such systems, although faster than film-based systems, do not have the rapid readout needed for some applications. In particular, they are much slower than the CCD-based detectors discussed above. In addition, the spatial resolution of these systems is limited to around 50 μm , which makes them inferior to film for high-resolution applications.

Film

The oldest, and still the most widely used x-ray detector, is silver halide-based photographic film. Indeed, the mysterious exposure of light-protected film led to discovery of x rays by Roentgen in 1895, and to demonstration of their use for medical imaging, still the most widely used and important application. Metastable excitation of grains of silver halide crystals is induced by x rays or secondary electrons. The resulting latent image is developed by chemical reduction of the excited grains to produce the familiar photographic image. Sometimes a phosphor coating or an adjacent phosphor plate is used to increase the sensitivity to the penetrating x radiation, with the halide excitation being produced by the secondary light emission from the phosphor. An extensive discussion of the use of film for optical imaging is given in Vol. I of this *Handbook*. Much of that discussion also applies to x-ray excitation. Detailed discussion of the use of film for x-ray detection and imaging is given in Refs. 1 and 5.

Cryogenic

Recent development of superconducting cryogenic x-ray detectors provides opportunity for important new applications, particularly in materials analysis and astrophysics. These make use of techniques to reach and maintain very low temperature (30 to 70 mK) for long periods of time.³¹ They use the measurement of the temperature rise (microcalorimeter),^{32–35} or tunnel junction current³⁶ in superconducting materials.

Microcalorimeter Measurement of the temperature rise induced by absorption of an x-ray photon has been demonstrated by the use of semiconductor (typically Si or Ge) thermistors whose electrical conductivity is temperature dependent, and with bimetallic transition-edge sensors (TES). Of these, the TES detectors are the most studied and appear to have the highest potential for x-ray spectrometry. They operate by holding the temperature of the metallic absorber at the transition edge between the superconducting and normal state. At the transition edge, the conductivity is extremely temperature sensitive. The temperature rise that is induced by an absorbed x ray is detected by measuring the conductivity with a superconducting quantum (SQUID) detector. The transition-edge temperature is typically about 50 mK and is stabilized by thermoelectric feedback. Energy resolution as low as 2 eV and counting rates as high as 500 counts per second (cps) have been reported.^{37,38} Both the energy resolution and the thermal recovery time (therefore the counting rate) are affected by the thermal capacitance of the absorber, which must be as low as possible. Absorber sizes of $300 \times 300 \times 50 \mu\text{m}$ are typical for a very high-resolution detector. Detector arrays are under development to increase the effective area, the count rate capability, and to provide imaging for astrophysical applications.³⁰ Alternatively, polycapillary focusing optics³⁹ have been used to increase the effective area (to $>7 \text{ mm}^2$).

Superconducting Tunneling Junction (STJ) Another type of cryogenic detector involves x-ray-induced breakup of superconducting Cooper pairs, which leads to decreased tunneling current across a superconducting Josephson junction. Although such detectors must also be small and maintained at low temperature to keep the thermally induced current across the junction low, the dependence of the counting rate on the thermal capacitance is relaxed. Such detectors have demonstrated energy resolution $<20 \text{ eV}$ and counting rate $>20,000 \text{ cps}$.³³ Again, arrays or optics have been used to increase the effective area.

60.3 SUMMARY

A list of properties of a number of single-pixel x-ray detectors is provided in Table 1 for comparative purposes. Spatial resolution for typical position-sensitive or imaging detectors is given in Table 2.

TABLE 1 Typical Parameters That Affect Detector Choice, Including Counting Rate Limitations, Availability of Current Mode Operation, and Energy Resolution

Detector	Intensity measurements		
	Pulse Counting (cps)	Current Mode	Energy Resolution
Ionization			
Ionization chamber	1–10 k	Yes	1–10 keV
Proportional counter	10 k		2–20 keV
Geiger counter	50 k		
Semiconductor	0–100 k		
Si(Li)			120 eV
HpGe			140 eV
CdZnTe			250 eV
Pin diode		Yes	
Scintillation detector	100–500 k		20–50%
Cryogenic			
Microcalorimeter detector	100–500 k		2–10 eV
STJ detector	2–20 k		20–50 eV

TABLE 2 Position-Sensitive X-Ray Detectors

Detector	Spatial Resolution (μm)
Film	20–200
Ionization	
Multiwire proportional counter	100–300
Semiconductor	
Arrays	
Segmented electrode	100–500
Charge-coupled detector (CCD)	20–50
Charge injection detector (CID)	20–50
Active pixel	
CMOS	20–50
Amorphous silicon	50–100
Pulse height dispersive	100–300
Scintillation	
Segmented	500–2000
Video	50–100
Multichannel plate readout	20–100
Restimulable phosphor	50–200
Cryogenic	
Array	
Microcalorimeter	200–500
Tunnel junction	300–500
Dispersive	500–1000

60.4 REFERENCES

1. W. J. Price, *Nuclear Radiation Detection*, 2nd ed., McGraw-Hill, New York, 1964.
2. G. G. Knoll, *Radiation Detection and Measurement*, Wiley, New York, 1979.
3. W. R. Leo, *Techniques for Nuclear and Particle Physics Experiments*, Springer-Verlag, Heidelberg, 1987.
4. C. F. G. Delaney and E. C. Finch, *Radiation Detectors*, Oxford University Press, Oxford, 1992.
5. C. J. Buckley, in *X-Ray Science and Technology*, A. G. Michette and C. J. Buckley, (eds.), Institute of Physics Publishing, Bristol, 1993, pp. 207–252.
6. G. Fraser, *X-Ray Detectors in Astronomy*, Cambridge University Press, Cambridge, 1989.
7. S. J. Gurman, “EXAFS Studies in Materials Science,” *J. Mat. Sci.* **17**(6):1541–1570 (1982).
8. E. A. Stearn and S. M. Heald, in “Investigation of Soft Radiation by Proportional Counters,” *Handbook on Synchrotron Radiation*, Vol. 1b, E. E. Koch, (ed.), North Holland, Amsterdam, 1983, p. 955.
9. D. G. Stearns and M. B. Stearns, “Microscopic Methods in Metals,” in *Topics in Current Physics*, U. Gonser, (ed.), 49, Springer, Berlin, 1983, p. 153.
10. S. C. A. Curran, A. L. Cockroft, and J. Angus, *Phil. Mag.* **40**:929 (1949).
11. B. D. Ramsey, J. J. Kolodziejczak, M. A. Fulton, J. A. Mir, and M. C. Weisskopf, “Development of Microstrip Proportional Counters for X-Ray Astronomy,” *Proc. SPIE* **2280**:110–18 (1994).
12. G. Charpak, R. Bouclier, T. Bressani, J. Favier, and C. Zupancic, *Nucl. Instr. and Meth. in Phys. Res.* **62**: 262–268 (1968).
13. B. D. Ramsey, R. A. Austin, and R. Decher, *Sp. Sci. Rev.* **69**:139–204 (1994).
14. R. Hofstader, *Nucleonics* **4**:2 (1949); *Proc. IRE* **38**:721 (1950).
15. G. L. Miller, W. M. Gibson, and P. F. Donovan, *Ann. Rev. Nucl. Sci.* **12**:189 (1962).
16. W. M. Gibson, G. L. Miller, and P. F. Donovan, in *Alpha, Beta, and Gamma Spectroscopy*, K. Siegbahn (ed.), North Holland, Amsterdam, 1964.

17. E. M. Pell, *J. Appl. Phys.* **31**:291 (1960).
18. J. H. Elliott, *Nucl. Instr. and Meth.* **12**:60 (1961).
19. J. L. Blankenship and C. J. Borkowski, *IRE Trans. on Nucl. Sci.* **NS-9**:213 (1963).
20. J. Kemmer, *Nucl. Instr. and Meth. in Phys. Res. A* **226**(1):89–93 (1984).
21. E. Gatti and P. Rehak, *Nucl. Instr. and Meth.* **225**:608 (1984).
22. A. Q. R. Baron and S. L. Ruby, *Nucl. Instr. and Meth. A* **343** (2-3):517–526 (1994).
23. J. E. Gindley, T. A. Prince, N. Gehrels, J. Tueller, C. J. Hailey, B. D. Ramsey, M. C. Weisskopf, P. Ubertini, and G. K. Skinner, *SPIE Proc.* **2518**:202–210 (1995).
24. F. P. Doty, H. B. Barber, F. L. Augustine, J. F. Butler, B. A. Apotovsky, E. T. Young, and W. Hamilton, *Nucl. Instr. and Meth. In Phys. Res.* **A253**:356–360 (1994).
25. S. Dhawan, *IEEE Trans. Nucl. Sci.* **NS-28**:672 (1981).
26. R. B. Murray, in *Nuclear Instruments and Their Uses*, A. H. Snell (ed.), Wiley, New York, 1962.
27. R. L. Heath, R. Hofstader, and E. B. Hughes, *Nucl. Instr. and Meth.* **162**:431 (1979).
28. E. F. Eikenberry, et al., *Nucl. Instr. and Meth. A* **501**(1):21 pp. 260–266 (2003).
29. H. Nanto, K. Murayama, T. Usuda, F. Endo, Y. Hirai, S. Tahiguchi, and N. Takeguchi, *J. Appl. Phys.* **74**:1445–1447 (1993).
30. H. Nanto, Y. Hirai, F. Endo, and M. Ikeda, *SPIE Proc.* **2278**:108–117 (1994).
31. M. Altman, G. Angloher, M. Buhler, T. Hertrich, J. Hohne, M. Huber, J. Jochum, et. al., *Proc. of 8th Int. Workshop on Low Temperature Detectors* (1999).
32. D. A. Wollman, K. D. Irwin, G. C. Hilton, L. L. Dulcie, D. E. Newbury, and J. M. Martinis, *J. Microscopy* 196–223 (1997).
33. D. McCammon, W. Cui, M. Juda, P. Plucinsky, J. Zhang, R. L. Kelley, S. S. Holt, G. M. Madejski, S. H. Moseley, and A. E. Symkowiak, *Nucl. Phys.* **A527**:821 (1991).
34. D. Lesyna, D. Di Marzio, S. Gottesman, and M. Kesselman, *J. Low Temp. Phys.* **93**:779 (1993).
35. E. Silver, M. LeGros, N. Madden, J. Beeman, and E. Haller, *X-Ray Spectrom.* **25**:115–122 (1996).
36. M. Frank, L. J. Hiller, J. B. Le Grand, C. A. Mears, S. E. Labov, M. A. Lindeman, H. Netel, and D. Chow, *Rev. Sci. Instrum.* **69**:25 (1998).
37. K. D. Irwin, G. C. Hilton, J. M. Martinis, S. Deiker, N. Bergren, S. W. Nam, D. A. Rudman and D. A. Wollman, “A Mo–Cu Superconducting Transition-Edge Microcalorimeter with 4.5 eV Energy Resolution at 6 keV,” *Nucl. Instrum. Meth. A* **444**(1–2):184–187 (2000).
38. D. A. Wollman, G. C. Hilton, K. D. Irwin, N. F. Bergren, D. A. Rudman, D. E. Newbury, and J. M. Martinis, 1999 NCSL Workshop and Symposium.
39. D. A. Wollman, C. Jezewski, G. C. Hilton, Q.-F. Xiao, L. L. Dulcie, and J. M. Martinis, *Microsc. Microanal.* **4**:172–173 (1998).

This page intentionally left blank.

DO NOT DUPLICATE

ADVANCES IN IMAGING DETECTORS

Aaron Couture

*GE Global Research Center
Niskayuna, New York*

Recent advances have allowed digital x-ray detectors to become widely available commercially, replacing many traditional film based systems. Digital x-ray detectors have also enabled a wide variety of new applications. This chapter summarizes the current state of the art for digital x-ray imaging detectors.

61.1 INTRODUCTION

X-ray detectors have a history of over one hundred years of providing rapid, simple, and low cost images of internal structure for a wide range of applications. Medical and dental x-ray imaging are used heavily in screening (e.g., mammography and dental check-up imaging), diagnostic (e.g., cardiac and angiographic imaging), and surgical (e.g., cardiac catheterization) procedures (see Chap. 31). Over a million medical x-ray procedures per day are performed worldwide. Industrial x-ray imaging is used to penetrate dense materials and provide high-resolution information about the internal composition and structure of manufactured parts. Security applications include airport screening of baggage and personal items, imaging of cargo containers and rail cars, as well as detection of land mines. In the past decade, digital x-ray detectors have begun to replace film screen systems due to decreased cost, improved performance, and workflow. Digital x-ray detectors continue to enable new applications not previously possible: tele-radiology, computer aided detection (CAD), digital subtraction contrast enhancement, and computed tomography (CT) are prominent examples. The cumulative growth of digital systems in the market has also accelerated, with the number of medical systems in the field from less than one thousand in the year 2000 to more than 30,000 in 2007.

Each imaging application brings new and different performance requirements for x-ray imaging detectors. Some of these competing detector requirements include: size, resolution, readout rate, noise, image quality, cost, weight, portability, and power. The following section will summarize a number of key technologies that represent the current state of the art in x-ray imaging detectors. A wide variety of concepts for digital x-ray detectors have been explored both academically and commercially. The technologies detailed in this chapter represent recently developed digital x-ray detectors. Photo-stimulated phosphors (CR plate) detectors, which currently compose a large fraction of the digital x-ray market, are not included here, but are introduced in Chap. 60.

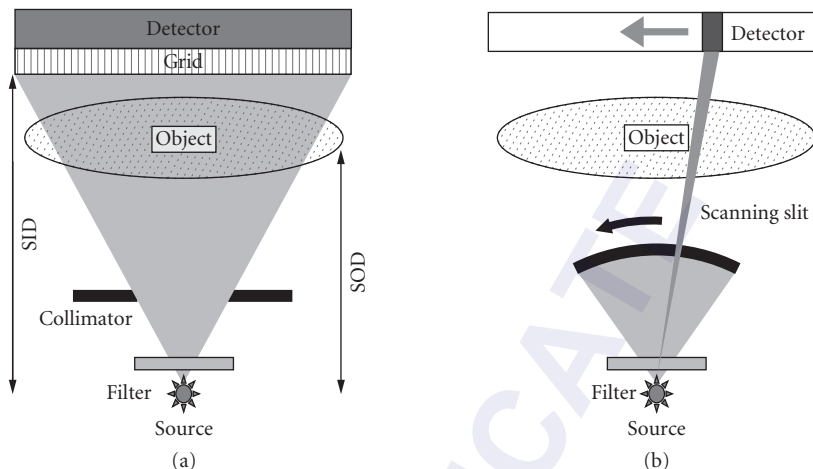


FIGURE 1 (a) Area type x-ray imaging system using a collimator to define the x-ray beam as well as a grid for scatter rejection. Source to object (SOD) and source to image distance (SID) are shown. (b) Slot type x-ray imaging system with fan beam scanning illumination and linear detector.

X-ray imaging is typically performed in one of two geometries (Fig. 1): area detection or slot detection. In area detection, the x-ray source is typically a few millimeters in size and is collimated to expose only the active area of the detector. The magnification of the system is the ratio of the source-to-image distance (SID) to the source-to-object distance (SOD). Especially in thick specimens, scatter can significantly degrade x-ray images. For area detection, a grid is typically inserted directly on the detector surface. Additionally, increasing the object to image distance decreases the scattered radiation hitting the detector; however, the fraction of the object imaged also decreases. In radiographic applications the maximum dimension to be imaged in a single shot is the chest, which sets the typical size for imaging at roughly 40 cm. In slot detection, a collimated fan beam and single slit linear detector are simultaneously scanned across the object during x-ray illumination. Slot detectors have the capability to reject scatter, and also reduce the size and cost of the detector. However, the x-ray tube is subjected to higher power and image acquisition is slower than area detection.

Another classification of x-ray detectors is the mode used for readout, including integrating and continuous. Integrating detectors collect signal from a gated x-ray source prior to triggering the readout of the detector. Detector readout is performed following illumination by the x-ray source. Data conversion electronics can be multiplexed so that a single converter channel is dedicated to an entire column of pixels. Continuous readout detectors output the instantaneous x-ray signal intensity during constant illumination. Depending on the gain, continuous readout detectors may also provide photon counting and energy information (see Chap. 62). Continuous readout requires a single conversion channel dedicated to each pixel. To maintain low power and cost, x-ray imaging detectors typically use integrating readout mode. Finally, x-ray imaging systems can be designed for single shot applications, termed radiographic, as well as for sequenced “video-like” imaging, termed fluoroscopic. For fluoroscopic applications, the fraction of signal from prior frames that contribute to later frames is quantified by the detector lag (0 to 100%).

For virtually every imaging application it is important to minimize x-ray dose. For medical applications x-ray dose can be harmful to both patients and doctors. For industrial and security applications, tube power output is limited, and x-ray dose determines the amount of time required to produce an image. The detective quantum efficiency (DQE) is a metric used widely to quantify the image quality (IQ) of x-ray imaging detectors. Note that this is different from the quantum detector efficiency used for single pixel detectors. The detective quantum efficiency quantifies the ability of a detector to accurately

transfer an x-ray input image into output electronic or digital signal, normalized to input x-ray dose (X). A detector with high DQE can deliver equivalent IQ for lower x-ray dose. It is computed as a function of spatial frequency f ,

$$\text{DQE}(f) = \frac{[S \cdot \text{MTF}(f)]^2}{\text{NPS}(f) \cdot X \cdot C} \quad (1)$$

where C is the x-ray fluence, S is the x-ray conversion efficiency, and NPS is the noise power spectrum, the spatial or temporal noise added to the image. MTF is the modulation transfer function, the ratio of the amplitude of the output image at a spatial frequency f to the input amplitude. MTF is a measure of the spatial resolution.

The spatial resolution of the detector is influenced not only by the pixel pitch, but also by spreading of signal to adjacent pixels. The conversion efficiency of the detector is a function of the efficiency of the absorbing layer to generate electronic signal as well as the fraction of detector area that is sensitive to x rays (fill-factor). In detectors that have electronic noise, the DQE of a detector degrades with lower x-ray dose. Electronic noise is suppressed in detectors that have inherent gain, such as image intensifiers. Flat panel detectors with no active gain have inherent electronic noise related to the transfer of small amounts of charge from the pixels in the array to electronics bonded to the panel. Some x-ray imaging detectors add electronic amplification in order to improve performance at low x-ray doses.

X-ray imaging detectors must also be insensitive to gain hysteresis, sometimes referred to as ghosting. An x-ray detector exposed to both high and low dose conditions must continue to have a uniform response. If the detector is hysteretic, the gain for regions of high exposure can be modified, leading to contrast appearing in subsequent images. Ghosting can lead to image artifacts building up over time. For acceptable quality medical images, the gain hysteresis must be limited to a few percent.

X-ray scatter can degrade the image quality in x-ray images and is especially a problem with thick specimens. The imaging geometry as well as the use of scatter-rejection grids can greatly improve the contrast-to-noise ratio (CNR) under such conditions. To improve the rejection ratio, grids can be manufactured with septa focused on the x-ray source and made from materials with high stopping power. The imperfect x-ray transmission of the grid can degrade the DQE, so high transmission fraction is important in maintaining image quality. Recently, the use of high aspect ratio microlithography (HARM) has been applied to manufacturing focused metal grids.¹ The high aspect ratio as well as dense material could potentially offer enhanced scatter rejection while maintaining good x-ray transmission.

61.2 FLAT-PANEL DETECTORS

Introduction

The most common x-ray imaging detectors available include: flat-panel indirect-conversion detectors, flat-panel direct-conversion detectors, and charge-couple-device (CCD) based detectors. Figure 2 shows a diagram of these three types of detectors, including an electrical schematic of a typical pixel. Each detector includes an absorbing layer responsible for stopping a fraction of incident x rays, a conversion layer resulting in electric charge, and a switching matrix responsible for storing and reading out charge. Flat panel detectors are typically based on amorphous or polysilicon on glass substrates, while CCD detectors are fabricated on single crystal silicon.

Both indirect and direct-conversion flat panel x-ray detectors utilize an array of amorphous silicon (a-Si) thin film transistors (TFT) to form an active array of switches for signal readout. Amorphous silicon has a number of advantages for use in x-ray detectors. First, the liquid crystal display (LCD) industry is based on very similar active matrix a-Si TFT panels. Market pressures have driven displays toward large area, low cost, and low defects. Photolithography, thin film deposition, and etching equipment can be

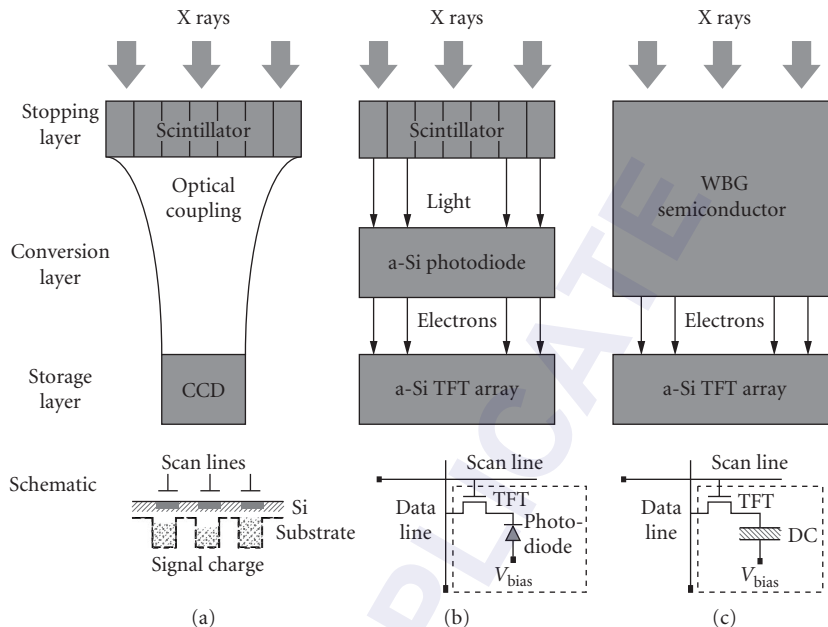


FIGURE 2 (a) Scintillator-based CCD detector with optical coupling; (b) indirect-conversion flat panel detector with scintillator layer optically coupled to photodiode layer, integrated onto a-Si TFT panel; and (c) direct-conversion detector with wide band gap (WBG) semiconductor absorption layer deposited directly onto a-Si TFT panel.

directly incorporated into manufacturing of x-ray detectors. Flat panel x-ray detectors with active area of more than 40 cm and pixel pitches less than 100 μm are available. Second, the a-Si TFT has very low leakage, allowing storage of signal charge for many seconds during x-ray exposure. Last, amorphous silicon is naturally radiation hard.

Figure 2 shows a schematic for pixels in indirect and direct-conversion flat panel detectors. The TFT switching matrix and readout method for flat panel x-ray detectors is the same. Flat panel detectors utilize scan lines that are individually energized to address rows of TFT pixels, as well as data lines that transfer signal charge to readout electronics at the border of the panel. The readout is multiplexed so that all pixels on a single line are converted simultaneously, utilizing a single converter for each data line. Multiplexed readout coupled with a panel design optimized for fast switching has enabled the use of flat panel detectors in fluoroscopic medical applications with 30 Hz frame rates. Readout electronics are designed to satisfy competing requirements for rapid readout, low noise, or low power.

Flat panel detectors also have advantages for compact design and weight. Portable radiography detectors are currently available with thicknesses of a few centimeters and weights under 15 lb. Wireless digital x-ray detectors have recently been demonstrated for radiography.

Flat Panel: Indirect Conversion

Indirect-conversion flat panel x-ray detectors utilize a layer of scintillator material to absorb x rays and emit visible light photons. The scintillator material must be high density in order to efficiently stop x rays, must efficiently convert absorbed x-ray energy into visible photons, and efficiently transfer visible light photons to the photosensing layer. The scintillator material must be a uniform layer across the imaging area, typically hundreds of micrometer thick and up to 40 cm in size. The scintillator

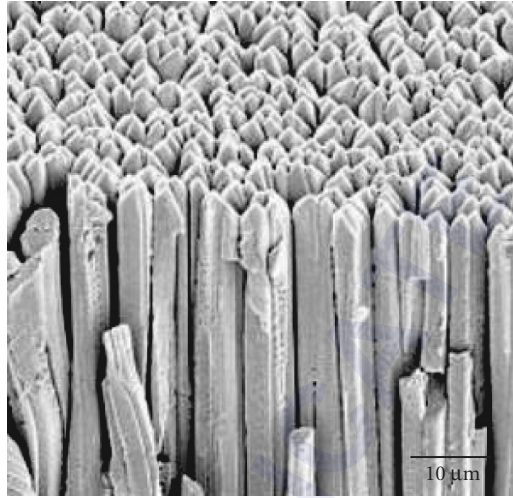


FIGURE 3 Thin film deposited CsI(Tl) scintillator, typically deposited in thicknesses from 100 μm to millimeters. Needle microstructure is responsible for enhancing scintillator spatial resolution.

material is typically a deposited thin film, a plastic sheet, or a solid plate. Thallium doped cesium iodide, CsI(Tl), is widely used as a scintillating material due to its high conversion efficiency, limited after-glow, and capability for thin film deposition, as well as an emission spectrum that matches well with a-Si photodiode absorption. A single x ray can generate thousands of optical photons. The optical photons are emitted isotropically, leading to spreading of the light signal before it is absorbed in the photosensing layer, degrading the spatial resolution and hence the DQE of the detector. Thin film deposited CsI(Tl) can be grown with a structured needle morphology that tends to limit lateral spreading of scintillation light, as shown in Fig. 3. Use of a structured scintillator allows thick layers (500 μm and greater) to be deposited while maintaining good spatial resolution (MTF) of the detector. Thallium doping boosts the photoelectron yield and also shifts the emission peak from 315 to 550 nm; both effects improve the conversion efficiency for a-Si diode photo detectors. In addition, lithographically patterned structuring of thin film deposited scintillator has also been investigated.² The visible light generated by the scintillator is captured by an amorphous silicon photodiode layer and converted to electric charge that can be stored in the pixel using the a-Si TFT array. In addition to requirements related to conversion efficiency and light spreading, which directly affect the DQE of the detector, the scintillator also has requirements related to lag, after-glow, and hysteresis. Typically lag related to the CsI(Tl) can be less than a few percent, and the gain hysteresis is also limited to less than a few percent.

The typical pixel design of an indirect-conversion flat panel x-ray detector is shown in Fig. 2. An amorphous silicon photodiode layer is deposited onto a TFT active matrix, and a common bias is provided to bias the array. Leakage through the photodiode degrades the dynamic range of the detector and can lead to spatial artifacts and temperature sensitivity. However, amorphous silicon diodes have been reported³ with leakage values of <100 pA/cm². Fluoroscopic lag in direct-conversion flat panels can be less than 10 percent.⁴ An etched mesa structure is the most typical photodiode geometry. However, in order to maximize the active area covered by the photodiode array and increase conversion efficiency of the detector, continuous photodiode layers have also been investigated.⁵

The DQE of indirect-conversion flat panel detectors is improved by the x-ray absorption in dense, high brightness scintillators. DQE for mammography detectors at 8.5 mR has been reported in excess of 80 percent at 0 frequency, with greater than 30 percent at 5 lp/mm (5 line pairs per mm is 100 μm

resolution). The DQE for these detectors is degraded mainly by light spreading in the scintillator as well as electronic noise in the TFT readout. Electronic noise sources originating from pixel switching as well as transferring the charge on the data line are significant for indirect-conversion flat panel detectors. Active areas of research include TFT pixel configurations with gain stage at each pixel as well as optimization of panel and converter electronics in order to reduce the electronic noise of the panel.⁶ Additionally, some designs incorporate an additional storage capacitor device to increase the charge that can be stored in each pixel.⁷ Diode switching can also be used in place of the TFT array.⁸ Lightweight and rugged substrates have been investigated, as well as the use of ink-jet printing of organic electronics for the TFT and photodiode.⁹

Flat Panel: Direct Conversion

Direct-conversion detectors, as shown in Fig. 2, combine the x-ray absorbing layer and conversion layer into one material. One advantage of direct-conversion detectors is the reduced manufacturing cost, due to a reduced number of layers on the TFT flat panel. Additionally, since the signal charge created in the conversion layer is subjected to a strong electric field, there is no degradation of the spatial resolution due to spreading of the charge in the direct-conversion material. Last, the fraction of area that is sensitive to x rays is large due to shaping of the electric field in the conversion layer. Direct-conversion detectors are commercially available for medical as well as industrial imaging, typically for single shot imaging systems. Fluoroscopic mode detectors also have been demonstrated.¹⁰

Direct-conversion materials must be optimized for multiple requirements simultaneously: high x-ray absorption (high Z, high density), efficient charge collection, low dark current, good uniformity, low lag, and long term reliability, and stability. The properties of amorphous selenium, mercuric iodide, and lead iodide, three materials that have been investigated for direct-conversion x-ray imaging, are summarized in Table 1. Amorphous selenium is currently the most widely used material in commercially available direct-conversion detectors. The x-ray absorption efficiency scales with the atomic number as Z^4 . Formation energy relates to the energy required to create a charge pair, and is inversely proportional to the number of charges liberated from each absorbed x ray. The mobility lifetime product relates to the ability to remove charge from the conversion layer. The electric field relates to the ability to operate the detector electronics at lower voltages.

Figure 4 shows the thin film layer structure for direct-conversion materials. Lead and mercuric iodide are deposited directly onto the TFT array, using physical vapor deposition (PVD) or screen-printing of semiconductor powder incorporated in a polymer binder (particle in binder, PIB). The bias electrode is deposited in a separate step, and a polymer layer is used to encapsulate the full structure for environmental isolation. Amorphous Se is typically deposited in a layered p-type/intrinsic/n-type (PIN) structure with PVD using lightly doped p-Se and n-Se as blocking contacts that reduce the leakage current, as shown in Fig. 4b. With bias applied, x rays that are absorbed in the direct-conversion material generate charge pairs that are swept to the pixel and bias contacts. Charge is stored in the pixel until readout of the detector, which occurs in the same way for indirect-conversion detectors, as described in the beginning of Sec. 61.2. Leakage current in the direct-conversion material can degrade

TABLE 1 Properties of Typical Direct-Conversion Materials Used in X-Ray Imaging Detectors^{11,12,13}

	HgI ₂	PbI ₂	a-Se
Atomic number (Z)	80, 53	82, 53	34
Band gap (Ev)	2.1	2.3	2.2
Charge pair formation energy (eV)	5	5.5	42
Mobility-lifetime product ($\mu\tau$, cm/V ²)	10 ⁻⁵ (h) 10 ⁻⁵ (e ⁻)	10 ⁻⁶ (h) 10 ⁻⁷ (e ⁻)	10 ⁻⁶ (h) 10 ⁻⁶ (e ⁻)
Electric field (V/ μ m)	0.2–1	0.2–1	10

Mobility-lifetime products for holes (h) and electrons (e) can also be dependent on temperature.

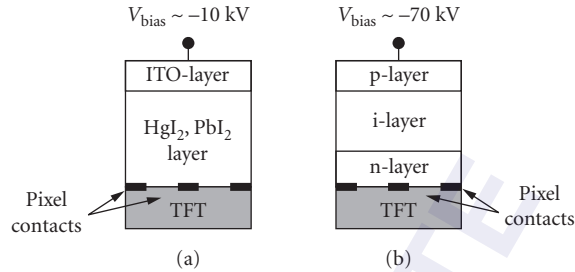


FIGURE 4 Typical configuration and biasing for direct-conversion thin films for (a) mercuric and lead iodide and (b) amorphous selenium. Indium tin oxide (ITO) and p-type/intrinsic/n-type (PIN) top contacts are important in blocking injection of carriers into the semiconductor as well as for removing signal charge.

the dynamic range of the detector and can also generate artifacts if the leakage is not stable with time or temperature. Leakage currents have been reported as low as $700\text{pA}/\text{cm}^2$ for mercuric iodide.¹⁴ Gain hysteresis and ghosting have been well studied during the development of direct-conversion detectors. The zero frequency DQE performance for direct-conversion flat panel detectors is reported to be in excess of 70 percent for amorphous selenium mammography detectors.¹⁵

61.3 CCD DETECTORS

A third class of digital x-ray imaging detectors utilizes a scintillation layer coupled to a high efficiency optical system which transfers the image onto a smaller (2 to 3 cm^2) CCD photo sensor. Applications for CCD detectors vary from medical imaging, mammography, and video rate industrial imaging, to protein crystallography.¹⁶ CCD systems have advantages in cost and low electronic noise related to the photo sensor; however, some degradation in DQE occurs due to inefficiency in the optical system that couples the scintillator to the CCD. A higher reduction ratio used to couple signal to the CCD results in lower efficiency, due to losses in the optical system that are fundamentally limited by solid angle considerations.

Small area prototype flat-panel CCD detectors have been demonstrated with a scintillator deposited directly on the photo sensor.¹⁷ An additional stage of optical amplification can be added in order to compensate for the inefficiency of the optical demagnification.¹⁶ The main elements of CCD x-ray imaging detectors are shown in Fig. 2. They include a scintillating layer for absorbing x rays which outputs visible light that is coupled to an optical system that focuses the image onto the CCD photosensor. The CCD is composed of rows of metal gates patterned on silicon with implanted regions that act as columns. The gates are biased to define the separation of pixels, and photocharges that are generated in the depletion region of the silicon are stored in the charge wells. To read out charge, the biasing voltage on the gates can be changed to shift charge to neighboring rows. Charge transport in CCDs can be highly efficient, resulting in readout noise levels of less than 10 electrons.¹⁶

Figure 5 shows three examples of common configurations for CCD detectors including both area detectors and slot scanners. Figure 5a and b show two tiling methods used to increase the size of the x-ray image, while limiting the demagnification ratio for the detector. Fiber optic tapers used for optical coupling are fused bundles of glass fiber light guides. The bundles are heated and shaped into tapers for demagnification, with near theoretical efficiencies.¹⁶ Tiles must not appear in the final x-ray image, so there must be no change in imaging performance across the boundary of adjacent units. Figure 5c shows the configuration for a CCD array used in a slot scanning geometry. Tiled units with scintillator and fiber optics are constructed into a linear detector that is translated across the imaging plane.

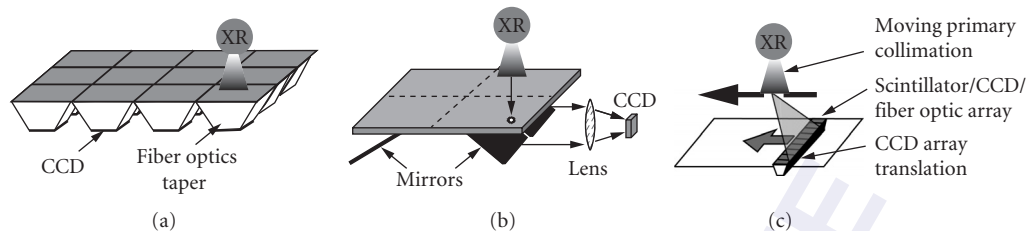


FIGURE 5 Three configurations of scintillator/CCD digital x-ray imaging detectors: (a) tiled detector elements with scintillator coupled to a fiber optic taper and CCD; (b) tiled detector with single scintillator plate plus multiple mirrors optically coupled to CCD; and (c) slot scanning geometry, linear tiled detector elements with scintillator coupled to a fiber optic taper and CCD.

Scintillators for CCD systems can be either optically coupled solid plate or plastic film, or thin film deposited directly on optical elements. The scintillating material has similar requirements and performance to that for indirect-conversion flat panel detectors; however the spectral response of the scintillator must be tailored to the absorption of CCD detectors that peak at red visible wavelengths.

61.4 CONCLUSION

Numerous digital x-ray imaging technologies have created a rapidly expanding medical and industrial x-ray imaging market. This chapter has focused on the most mature technologies; however, rapid market growth is currently driving the development of higher performance and lower cost systems. A number of new technologies are currently under development that have the potential to further improve image quality, increase detector readout rates, and reduce system cost.

61.5 REFERENCES

1. K. Fischer, B. Chadhuri, H. Guckel, and C. Tang, "Fabrication of Two-Dimensional X-Ray Anti-Scatter Grids for Mammography," *Advances in X-ray Optics, Proceedings of SPIE* **4145**:227–234 (August 2000).
2. V. Nagarkar, S. Tipnis, V. Gaysinskiy, S. Miller, A. Karellas, and S. Vedantham, "New Design of a Structured CsI(Tl) Screen for Digital Mammography," *Proceedings of SPIE* **5030**:541–546 (June 2003).
3. S. Tchakarov, P. Cabarrocas, U. Dutta, P. Chatterjee, and B. Equer, "Experimental Study and Modeling of Reverse-Bias Dark Currents in PIN Structures Using Amorphous and Polymorphous Silicon," *Journal of Applied Physics* **94**:7317–7327 (2003).
4. P. Granfors, R. Aufrichtig, G. Possin, B. Giambattista, Z. Huang, J. Liu, and B. Ma, "Performance of a $41 \times 41 \text{ cm}^2$ Amorphous Silicon Flat Panel X-Ray Detector Designed for Angiographic and R&F Imaging Applications," *Medical Physics* **30**:2715–2726 (2003).
5. R. Weisfield, W. Yao, T. Speaker, K. Zhou, R. Colbeth, and C. Proano, "Performance Analysis of a 127-Micron Pixel Large-Area TFT/Photodiode Array with Boosted Fill Factor," *Proceedings of SPIE* **5368**:338–348 (May 2004).
6. L. Antonuk, Y. Li, H. Du, Y. El-Mohri, Q. Zhao, J. Yamamoto, A. Sawant, Y. Wang, and Z. Su, "Investigation of Strategies to Achieve Optimal DQE Performance from Indirect Detection, Active Matrix Flat-Panel Imagers (AMFPIs) through Novel Pixel Amplification Architectures," *Proceedings of SPIE* **5745**:18–31 (February 2005).
7. D. Albagli, S. Han, A. Couture, H. Hudspeth, C. Collazo, and P. Granfors, "Performance of Optimized Amorphous Silicon, Cesium-Iodide Based Large Field-of-View Detector for Mammography," *Proceedings of SPIE* **5745**:1078–1086 (February 2005).

8. C. van Berkel, M. Powell, and S. Deane, "Physics of a-Si:H Switching Diodes," *Journal of Non-Crystalline Solids* **164–166**(2): 653–658 (1993).
9. R. Street, W. Wong, S. Ready, R. Lujan, A. Arias, M. Chabynyc, A. Salleo, R. Apte, and L. Antonuk, "Printed Active-Matrix TFT Arrays for X-Ray Imaging," *Proceedings of SPIE* **5745**:7–17 (February 2005).
10. D. Hunt, O. Tousignant, and J. Rowlands, "Evaluation of the Imaging Properties of an Amorphous Selenium-based Flat Panel Detector for Digital Fluoroscopy," *Medical Physics* **33**:1166–1175 (2004).
11. D. Alexiev, N. Dytlewski, M. I. Reinhard, and L. Mo, "Characterisation of Single-Crystal Mercuric Iodide," *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* **517**:226–229 (2004).
12. R. Street, M. Mulato, S. Ready, R. Lau, J. Ho, K. VanSchuylenbergh, M. Schieber, et al., "Comparative Study of PbI₂ and HgI₂ as Direct Detector Materials for High Resolution X-ray Image Sensors," *Proceedings of SPIE* **4320**:1–12 (June 2001).
13. G. Belev, and S.O. Kasap, "Amorphous Selenium as an X-Ray Photoconductor," *Journal of Non-Crystalline Solids* **345–346**:484–488 (2004).
14. G. Zentai, L. Partain, R. Pavlyuchkova, C. Proano, G. Virshup, L. Melekhov, A. Zuck, et al., "Mercuric Iodide and Lead Iodide X-Ray Detectors for Radiographic and Fluoroscopic Medical Imaging," *Proceedings of SPIE* **5030**:77–91 (February 2003).
15. J. Jesneck, R. Saunders, E. Samei, J. Zia, and J. Lo, "Detector Evaluation of a Prototype Amorphous Selenium Based Full Field Digital Mammography System," *Proceedings of SPIE* **5745**:478–485 (February 2005).
16. S. Gruner, M. Tate, and E. Eikenberry, "Charge-Coupled Device Area X-Ray Detectors," *Review of Scientific Instruments* **73**:2815–2842 (2002).
17. D. Scheffer, "A Wafer Scale Active Pixel CMOS Image Sensor for Generic X-Ray Radiology," *Proceedings of SPIE* **6510**:651000–651001 (February 2007).

This page intentionally left blank.

DO NOT DUPLICATE

X-RAY SPECTRAL DETECTION AND IMAGING

Eric Lifshin

*College of Nanoscale Science and Engineering
University at Albany
Albany, New York*

The concept of an x-ray image is usually associated with a radiograph (also see Chap. 31).¹ X rays are transmitted through an object and an image is formed on a detector placed on the side opposite the source. Contrast is based on point-to-point variation in absorption and the resolution is determined by scattering effects and by the pixel size of the detector. Radiographs are collected in parallel, that is, the entire image is formed at one time either traditionally by film or more recently by solid state imagers of various types. The transmission mode is of great value because the penetrating power of x rays makes it possible to see structure within objects be they humans or metallurgical castings. Images containing surface detail of the type created with a reflection optical microscope are difficult, if not impossible, to form because lens systems similar to those used in optical microscopes are generally not available with the exception of zone plates and reflective multilayer optics (see Chaps. 40 and 41). The index of refraction for most materials for x rays is so close to one that glass or other lens materials used for focusing the visible, infrared, or ultraviolet parts of the electromagnetic spectrum will not work for x rays. Thus, images are formed by transmitting x rays through the object of interest onto an image plane. In the most conventional forms of radiography the magnification on the detector plane is unity and any magnification in the image is the result of optically or electronically magnifying the detected image.

Spectral imaging, on the other hand, often involves creating a map of an object based on the spatial distribution of x rays of one or more particular energies. This is usually not done in transmission by the formation of a kind of mono-energetic radiograph, but rather by examining the emission of x rays excited from an area of a material selectively bombarded with some form of ionizing radiation. Furthermore, parallel imaging is not used as in the case for radiography. Instead the ionizing radiation is focused to a point on the object surface and then scanned over that surface in the form of a two-dimensional raster analogous to the electron beam scan on a cathode ray tube (CRT) display such as a conventional television tube. The image is formed in a serial or sequential mode, but if the scanning is fast enough and the persistence of the CRT phosphor long enough, it will look like an image formed in parallel. The focused (or highly collimated) ionizing beam used to cause x-ray emission can consist of electrons, ions, or x rays. In addition to the mode of scanning the beam over the surface of the sample, the beam can also be stationary and the sample physically scanned under a static beam. Either mode of operation forms the basis for a variety of analytical microscopy techniques including scanning electron microscopy (SEM), the analytical electron microscopy (AEM), proton induced x-ray emission (PIXIE), x-ray microfluorescence (XRF), and focused ion beam microscopy (FIB).

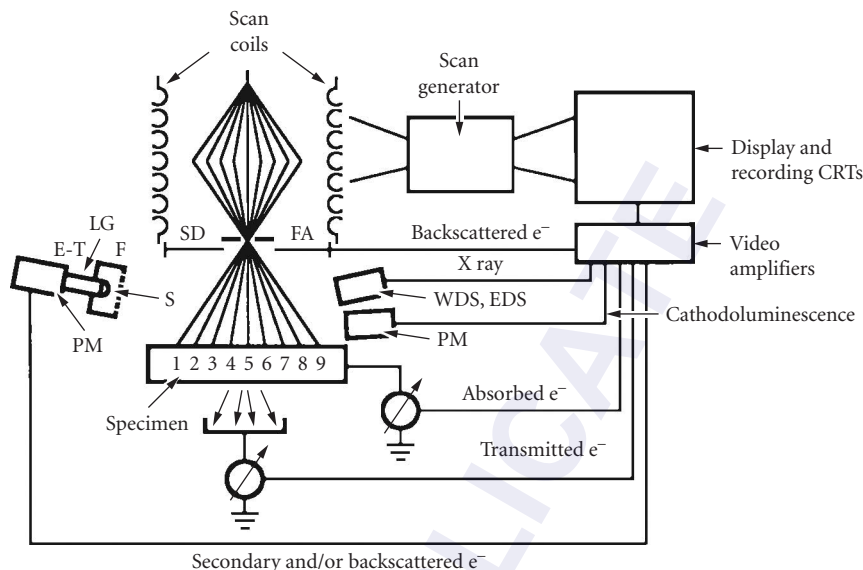


FIGURE 1 Basic SEM operation. An electron beam is scanned over the surface of sample covering the surface plane. Points 1 to 9 shown here are just in one dimension. The scan generator also controls the scan on a CRT display. The brightness of the CRT is determined by the intensity of a selected signal which as shown here could be the secondary electron signal detected with an Everhart-Thornley detector (ET) or various other detectors including EDS and WDS x-ray detectors.

As an example, the basic operation of an SEM, the most popular embodiment of this approach, is shown in Fig. 1. A focused beam of electrons is scanned over the surface of a specimen causing the emission of secondary electrons at each point it strikes. The emitted secondary electron signal is detected and its intensity used to modulate the brightness of a synchronously scanning CRT or other display controlled by the same scan generator. The magnification of the image is simply the ratio of the distance scanned on the display to that scanned on the sample. It is simply varied by changing the scan area on the sample. Variations in the intensity of the secondary electrons emitted from each point forms the basis of the observed image contrast. This point-to-point variation signal can be due to differences in the secondary electron yield arising from differences in surface topography or the material properties of the sample. SEMs typically operate in the magnification range from about $20\times$ up over $1,000,000\times$ with their ultimate resolution determined by how finally the electron beam can be focused, the volume from which the detected secondary electrons originate, and signal to noise ratio associated with the small electron beam currents used. Figure 1 also shows the use of x-ray detectors used to create images based on the point-to-point variation of the x-ray emission of a single energy.

Before discussing x-ray imaging it is first useful to review how x-ray analysis is done from a single point, as for example, one based on a detail observed in an SEM image. Currently the most common spectrometers/detectors used are the energy dispersive detector (EDS) and the wavelength dispersive (WDS) detector. A detailed description of the operation of an EDS detector can be found in Ref. 2 and an introduction is given in Chap. 60. The basic concept is that the x-ray energy of each photon detected is used to create electron-hole pairs in a solid state device. The number of electron-hole pairs produced is proportional to the energy of the x-ray photon detected. The initial electron pulse is converted into a voltage pulse whose size remains proportional to the x-ray energy. The pulses are then sorted by voltage and displayed as a histogram of pulse intensity versus voltage. The use of reference samples makes it easy to adjust the gain such that the x scale is calibrated to give energy rather than voltage.

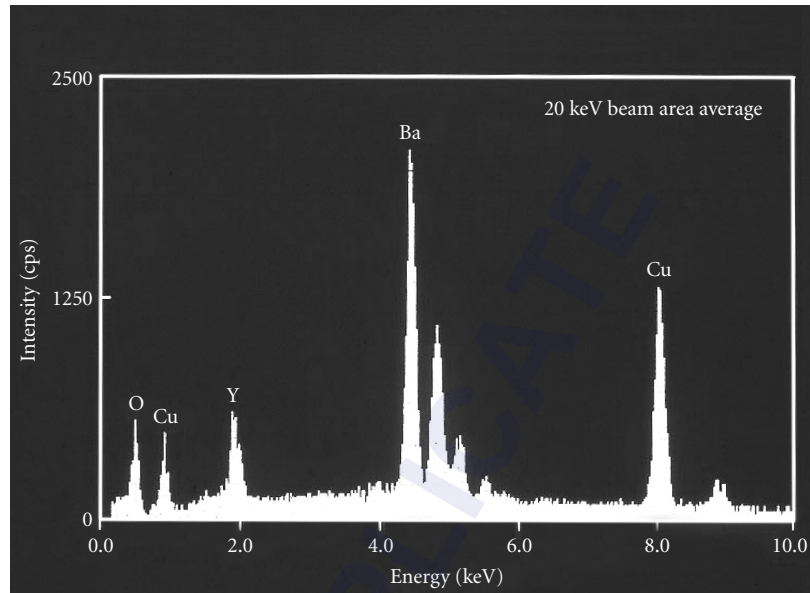


FIGURE 2 EDS spectrum of high temperature superconductor.

Figure 2 gives a typical x-ray spectrum, in this case of a high temperature superconductor, YBCO. Peaks can usually be identified easily as a result of Moseley's law, which relates the energy of a given spectral series to that of the atomic number of the element from which it originates. The characteristic lines observed are the result of the primary electron beam ejecting core shell electrons and the subsequent emission of x-ray photons with an energy equal to the difference between that of the core shell from which the electron was ejected and that of the shell from which the electron drops down to fill the vacancy (see Chap. 29). The range of x-ray energies used in SEM-EDS analysis is typically from about 100 eV up to about 20 KeV, which covers all elements from Be on up in atomic number. The lines used are either K, L, or M lines depending on whether the core shell electron was ejected from the first, second, or third energy level as defined by the principal quantum number of that level being 1, 2, or 3. To ensure an adequate signal intensity the electron beam energy used is typically 2 or 3 times the excitation energy of the core shell electrons associated with a particular element and shell. Since SEMs usually operate at 30 keV or less the lines used in an analysis are generally the K lines for the first third of the periodic table, the L lines for the middle third, and the M lines for the highest atomic number elements. Exceptions occur in cases of serious peak overlaps in the observed spectra.

EDS detectors have been used with SEMs since the late 1960s and since that time performance has improved significantly. In particular the energy resolution as measured at Mn $K\alpha$ has improved from about 500 eV to better than 130 eV. The development of thin windows capable of withstanding an atmosphere difference in pressure between the detector and the specimen environment combined with very low noise electronics has made Be $K\alpha$ (0.109 keV) analysis possible. Very sophisticated user interfaces have been developed to provide for the rapid identification of peaks, background removal (mostly from the x-ray continuum), peak area determination (needed for quantitative analysis), and a range of data display options including x-ray mapping. The past few years have seen another major advance that will revolutionize EDS analysis. It is the development of silicon drift detectors that will eventually replace lithium drifted detectors of the type that has dominated the field for well over 30 years. These detectors³ have all of the features mentioned above, but can increase throughput by 10 to 100 times resulting in count rate capabilities as high as a million counts per second when those signal levels are available. Until recently it was nearly impossible to get count rates greater than about 5000 counts per second without

paying a serious penalty in detector resolution. This was a serious limitation in x-ray mapping for two reasons in particular. First, when doing a point-by-point map the dwell time per pixel may be minimal (on the order of microseconds) depending on the scan rate used and the number of pixels used. Thus the number of counts obtained even from an element present at a high concentration level could be so small that counting statistics would make it very difficult to determine that the signal is above the background level. Second, in EDS the count rate of importance is the total count rate entering the detector. If a minor constituent is to be mapped then its proportion of the total count rate may be so low that it will be totally obscured by the high count rate of the major peaks. The only way to overcome the above difficulties is to take very long scans, but that may lead to problems arising from instrument drift and/or sample contamination. An example of how these difficulties can be overcome with the use of the new silicon drift technology is shown in Fig. 3 where x-ray maps were obtained in 12 s.

In an SEM, characteristic x-ray photon production is a very inefficient process compared to the generation of secondary electrons where the yield and collection of secondary electrons can even be greater than one per incident electron. In fact, the number of x-ray photons collected can be less than one per billion incident electrons. The number of incident electrons is closely related to the electron probe size on the sample and thus for example a 100 nm probe that could contain 10^{10} electrons striking the sample per second may yield only a relatively small number of measured x rays of a given elemental line even from a pure element sample. At 1.0 nm, which is about the smallest size of probes now in use, the yield is considerably less. The spatial resolution of x-ray maps collected in this way is ultimately determined by the x-ray excitation volume, which depends strongly on the electron beam energy, the sample composition, and the spectral line selected. Figure 4 shows a Monte Carlo simulation of electron scattering in copper for different beam energies.³ It clearly shows that the volume of electron scattering can be considerably larger than the electron beam size. Although electron energies are tracked until they are essentially zero, as long as the energy exceeds the ionization energy for a given spectral line, x rays will be produced. In this example that volume will be larger for the lower excitation energy Cu $L\alpha$ line than for the Cu $K\alpha$ line. While reducing the beam energy may help it must be pointed out that the intensity of the emitted line is a strong function of the ratio of the beam energy to

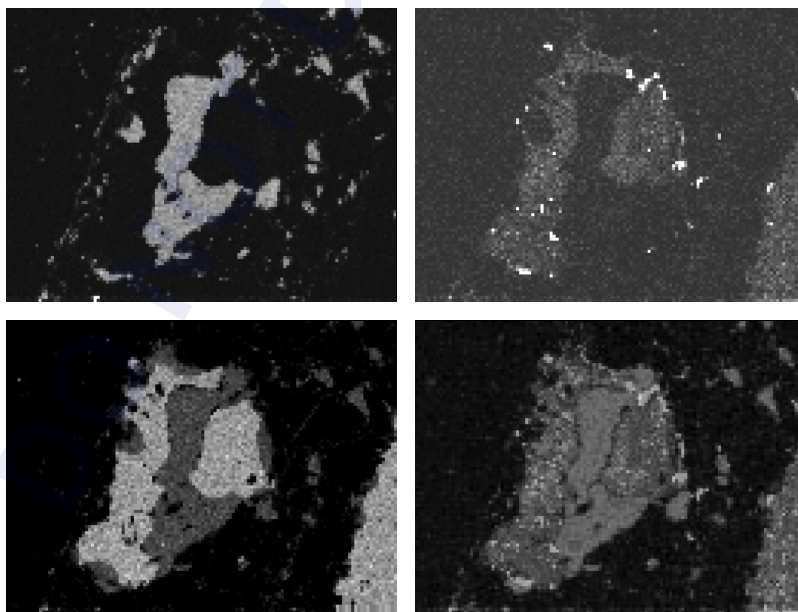


FIGURE 3 A 12-s x-ray spectrum image (128×96 pixels; 1 ms) obtained with a silicon drift detector. (See also color insert.) (Courtesy of Dale Newbury NIST.)

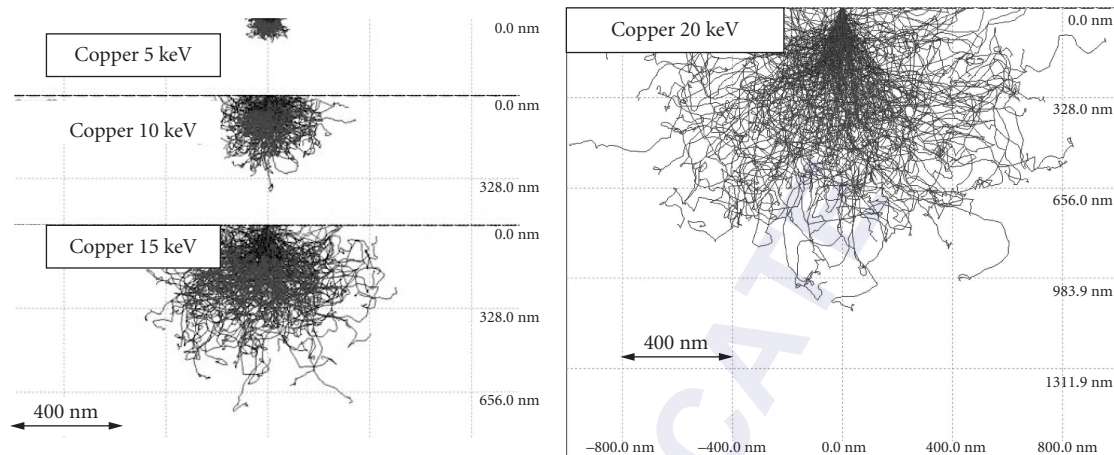


FIGURE 4 Monte Carlo simulations of electron trajectories in copper at various beam energies. Note the significant increase in the scattering volume with increasing energy. (See also color insert.)

the excitation energy, so the resolution approaches its best value as the generated x-ray intensity goes to zero. Traditionally x-ray microanalysis has had a resolution limit of about $1.0\ \mu\text{m}$ for the reasons cited when thick samples are studied. A combination of careful line selection and beam energy can drop this value to about $100\ \text{nm}$ at best. To do better, thin samples must be prepared to limit electron scattering and higher beam energies used. This is the approach of the analytical transmission electron microscope where the limit of spatial resolution for chemical analysis can be around $10\ \text{nm}$.

If larger areas are to be examined and very high spatial resolution is not required, x-ray mapping by x-ray fluorescence can be very effective. It also provides better detection limits since x-ray excited x-ray spectra have much less background due to less of the continuum in the spectrum (the continuum background arises only from scattering of the primary x-ray source since x rays do not generate continuum spectra from the sample). Figure 5 is an example of an x-ray map generated using x-ray excitation and a scanning specimen stage. While the spatial resolution of this type of image can be extended down to about $15\ \mu\text{m}$, it is expected that, as more advanced synchrotrons are implemented, x-ray probes based on collimated beams could be $50\ \text{nm}$ or less. Focused high energy ion beams, particularly protons, can also be used to create x-ray maps, but once again the yield can be low and the resolution will not be much better than $1000\ \text{nm}$. Medium energy focused ion beams, such as the $30\ \text{keV}$ used in an FIB, can be focused to less than $10\ \text{nm}$; however, the x-ray generation cross-sections are low, as are the beam currents, so that no measurable x rays are observed.

As stated previously, most x-ray spectroscopic imaging involves point-by-point (serial) data collection, by either scanning the probing beam or the specimen stage. There has also been some recent work with controlled-drift detectors to create two-dimensional x-ray detectors that provide spectral information from an array of points without scanning. For example, Castoldi et al. use polycapillary optics (see Chap. 53) to excite an array of points on a specimen and then another capillary optic to image the excited x rays on the surface of a two-dimensional detector with a $36\ \text{mm}^2$ detector containing 180 by $180\ \mu\text{m}$ pixels.⁴ In this manner they were able to image the $\text{Cu K}\alpha$ line on a fine pitch printed circuit board used as a sample. While this resolution is considerably less than that of the scanning methods used in the SEM, the method may still be useful for dynamic studies of samples where the sampling time may be very short since the parallel detection process can be quite fast. Another approach to parallel detection has been used in conjunction with CdZnTe and CdTe detectors based on pixelated CMOS signal processing in which both the numbers of x-ray photons and their average energy can be determined. This work has been recently described by Kruger et al.,⁵ but the spatial resolution is also limited due to the pixel sizes currently possible.

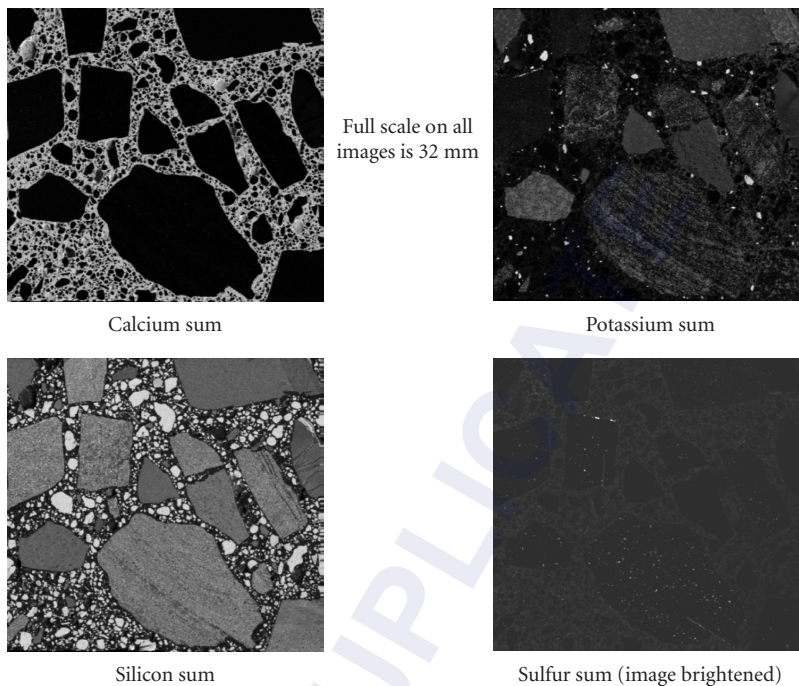


FIGURE 5 Micro-XRF maps. (Courtesy of J. M. Davis, NIST.)

62.1 REFERENCES

1. E. Krestel, *Imaging Systems for Medical Diagnostics*, Siemens Aktiengesellschaft, Berlin, pp. 221–249, (1990).
2. J. I. Goldstein, D. E. Newbury, D. C. Joy, C. E. Lyman, P. Echlin, E. Lifshin, L. Sawyer, and J. Michael, *Scanning Electron Microscopy and X-Ray Microanalysis*, 2nd ed., Kluwer Academic/Plenum Press, New York, NY, pp. 297–323 (2003).
3. D. Drouin, A. R. Couture, D. Joly, X. Tastet, V. Aimez, and R. Gauvin, “CASINO V2.42—a Fast and Easy-To-Use Modeling Tool for Scanning Electron Microscopy and Microanalysis Users” *Scanning* **29**(3):92–101 (2007).
4. A. Castoldi, C. Guazzoni, R. Hartman, and L. Strueder, “Application of Controlled-Drift Detectors to Spectroscopic X-Ray Imaging,” *IEEE Nuclear Science Symposium Conference Record* 1003–1008 (2007).
5. H. Krueger, J. Fink, E. Kraft, N. Wermes, P. Fischer, I. Peric, C. Herrmann, M. Overdick, and W. Ruetten, “CIX—A Detector for Spectral Enhanced X-Ray Imaging by Simultaneous Counting and Integrating,” arXiv.org, e-Print Archive, *Physics*, 1–12, arXiv:0802.2247v1 [physics.ins-det] (2008.)

SUBPART

5.6

NEUTRON OPTICS AND APPLICATIONS

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

David Mildner

*NIST Center for Neutron Research
National Institute of Standards and Technology
Gaithersburg, Maryland*

63.1 NEUTRON PHYSICS

The neutron is a subatomic particle with zero charge, a mass of $m = 1.00897$ atomic mass units, a spin of $1/2$, and a magnetic moment of $\mu_n = -1.9132$ nuclear magnetons. These four basic properties make thermal and cold neutrons not only such a rich and useful scientific tool for the investigation of condensed matter, but also a basis for observing many beautiful optical phenomena with remarkable properties,¹ as well as for constructing numerous optical devices.² “Thermal” neutrons are those in thermal equilibrium with their surroundings near room temperature. When a beam of energy, E , given by

$$E = \frac{2\pi^2\hbar^2}{m\lambda^2} \quad (1)$$

is selected from a thermal distribution, the de Broglie wavelength, λ , is comparable to interatomic distances in condensed matter, where $h = 2\pi\hbar$ is Planck's constant. Consequently neutron *diffraction* in condensed matter is analogous to x-ray diffraction. However, a typical thermal neutron has an energy of 0.025 eV or 4×10^{-21} J, much lower than electromagnetic radiation of comparable wavelength. Consequently, the mass and energy are such that the frequency of the radiation is comparable with the vibrational frequencies found in materials, which makes the measurement of the *inelastic scattering* a useful probe of these vibrations. Hence, neutrons are ideally suited for the study of the atomic structure and dynamics in condensed matter. The magnetic moment of the neutron interacts with those of unpaired electrons in magnetic materials, giving rise to *magnetic diffraction* and *inelastic scattering*. Again, the wavelength and energy of thermal neutrons are ideal for the study of the magnetic structure and dynamics of spin systems. As a result of zero charge, the neutron has only a short-ranged interaction with the nucleus. This means that the amplitude of the interaction is small, so that neutrons penetrate into the bulk of most materials. Moreover, the interaction probability of neutrons varies irregularly with the nuclear isotope, unlike x rays whose amplitudes increase monotonically with the atomic number. Finally, the neutron spin of $1/2$ enables a neutron beam to exist in one of two polarized states. When the neutron scatters from a nucleus of nonzero spin, the strength of the interaction depends on the relative orientation of the neutron and the nuclear spin.

The propagation of neutron de Broglie waves in a potential field is analogous to the propagation of light waves in a medium with a continuously variable refractive index. The potential can be

gravitational, magnetic, or nuclear. For example, slow neutrons follow a parabolic path under the effect of gravity as in classical mechanics. Neutrons in a constant magnetic field experience a torque and undergo precession. In a nonuniform magnetic field they experience a force that depends on the relative orientation of the spin and field vectors. Thus, a nonuniform magnetic field is a birefringent medium for an unpolarized neutron beam, with results analogous to the Stern-Gerlach experiment. Neutrons can be focused by refraction in an inhomogeneous magnetic field provided by a magnetic hexapole. Neutron waves in bulk nonmagnetic materials interact with the atomic nuclei. Generally, the interaction gives rise to an isotropic spherical wave that reradiates the neutron beam incident on the nucleus. The potential of the neutron-nucleus interaction has an imaginary part that represents neutron absorption. A scattering length, b , that depends on a point-like interaction potential, describes the scattering of a beam of neutrons. Generally, b is positive (but not always) and varies with different isotopes, and even with the same isotope, depending on the relative spin orientations of the neutron and isotope.

The similarity of the mathematical descriptions for neutron wave and light propagation gives rise to phenomena that are analogous to those of classical optics. In fact, virtually all the well-known classical optical phenomena that are characteristic of light and x rays have also been demonstrated with neutrons. In geometric optics, there are not only the refraction and reflection of neutrons by materials, but also special properties in magnetic media. In wave optics, there is Bragg diffraction from crystalline materials, and so on, as for x rays, but there are also other phenomena completely analogous to classical optics. Various neutron optics experiments have demonstrated quantum-mechanical phenomena on a macroscopic scale. For example, the perfect crystal neutron interferometer can measure the relative phase between two plane wave states, with the large separation of the beams enabling easy access for material phase-shifting devices and magnetic fields. These provide both quantitative verification of various fundamental quantum-mechanical principles applied to neutrons and accurate measurements of neutron scattering lengths. Finally, neutron optical devices have been developed to transport, collimate, focus, monochromate, filter, polarize, or otherwise manipulate neutron beams for applications in both basic and applied research, such as the study of the microscopic structure and dynamics of materials.

The optical phenomena arise from coherent elastic scattering of neutrons in condensed matter.³⁻⁵ The neutron wave function, $\psi(\mathbf{r})$, can be described by a one-body stationary Schrödinger equation, upon which all neutron optics is based, viz.,

$$[-(\hbar^2/2m)\nabla^2 + V(\mathbf{r})]\psi(\mathbf{r}) = E\psi(\mathbf{r}) \quad (2)$$

in which E is the incident neutron energy. The optical potential $V(\mathbf{r})$ represents the effective interaction of the neutron with the medium. The scattering of a neutron by a single bound nucleus is based on the Born approximation and uses the Fermi pseudopotential, $V(\mathbf{r}) = (2\pi\hbar^2/m)b\delta(\mathbf{r})$, where \mathbf{r} is the neutron position relative to the nucleus, to represent the effective interaction between the neutron and the nucleus. The wave function $\psi(\mathbf{r})$ provides a description of coherent elastic scattering and all neutron optical phenomena. There are other scattering processes (incoherent elastic scattering and inelastic scattering) collectively referred to as *diffuse scattering*. In addition, the incident neutron might be absorbed by the nucleus. Both of these cause attenuation of the coherent wave, $\psi(\mathbf{r})$, in the medium, so that the potential $V(\mathbf{r})$ and the bound scattering length, b , are in general complex.

The interactions of neutrons in bulk nonmagnetic matter are with atomic nuclei, and for neutrons traveling in a bulk medium the potential $V(\mathbf{r})$ may be replaced by a summation of pseudopotentials centered at each nucleus labeled i , given by

$$V(\mathbf{r}) = \sum_i (2\pi\hbar^2/m)b_i\delta(\mathbf{r}-\mathbf{r}_i) \quad (3)$$

This aggregate potential results in a coherent scattered wave that is the sum of the incident plane wave and the superposition of the spherical waves emanating from each nucleus. For a homogeneous system (gas, liquid, or amorphous solid) the optical potential has a constant value $V_0 = (2\pi\hbar^2/m)\rho b$, where ρ is the atom density and b is the average bound coherent scattering length per atom.

The Schrödinger equation is a macroscopic equation that describes coherent elastic scattering and all neutron optical phenomena in terms of the interaction of the neutron wave with a potential barrier.

The general solution in a medium of constant potential, V_0 , can be expressed as a superposition of plane waves where the magnitudes of the wave vectors can be determined by the incident neutron energy,

$$E = \frac{(\hbar\mathbf{k})^2}{2m} = \frac{(\hbar\mathbf{k}')^2}{2m} + V_0 \quad (4)$$

where \mathbf{k} and \mathbf{k}' are the incident and secondary wave vectors. Note that the neutron momentum $\mathbf{p} = \hbar\mathbf{k}$, and the magnitude of $|\mathbf{k}| = 2\pi/\lambda$, where λ is the neutron wavelength. For elastic scattering, the two equations, $|\mathbf{k}| = |\mathbf{k}'|$ (conservation of energy) and $\mathbf{k} = \mathbf{k}' + \mathbf{q}$ (conservation of momentum) combine, where \mathbf{q} is the scattering vector whose magnitude is given by

$$|\mathbf{q}| = (4\pi/\lambda)\sin(\varphi/2) \quad (5)$$

where φ is the scattering angle, that is, the angle between \mathbf{k} and \mathbf{k}' , and $\hbar\mathbf{q}$ is the momentum transferred from the neutron to the scattering system. The directions of the various wave vectors and their corresponding amplitudes are determined by requiring $\psi(\mathbf{r})$ and $\nabla\psi(\mathbf{r})$ to be continuous at the boundary between media. Some neutron optical phenomena are well described by the kinematic theory of diffraction. In geometric optics, neutron trajectories obey the same laws of reflection and refraction as in classical optics, though true mirror reflection only occurs for ultracold neutrons. Other phenomena require the dynamical diffraction theory that takes into account the interchange between the transmitted and reflected waves. Goldberger and Seitz⁶ have shown from the theory of dispersion that, with respect to neutron optics, all materials behave like a continuous macroscopic medium with a refractive index. In general, the propagation of de Broglie neutron waves in a potential field $V(\mathbf{r})$ is analogous to the propagation of light waves in a medium with a continuously variable refractive index that is defined by

$$n(\mathbf{r}) = [1 - V(\mathbf{r})/E]^{1/2} \quad (6)$$

In bulk media, $V(\mathbf{r})$ is replaced by V_0 . Sears¹ has given a more rigorous and comprehensive treatment of dispersion theory.

63.2 SCATTERING LENGTHS AND CROSS SECTIONS

The study of the neutron optics within materials requires understanding of scattering lengths and cross sections. This is also necessary for neutron scattering measurements in the study of the structure and dynamics of condensed matter. The scattering of a neutron by a single bound nucleus is described within the Born approximation by the Fermi pseudopotential because $V_0/E \ll 1$. The scattering length is a measure of the strength of the interaction of the neutron with the nucleus, and the intensity of neutron scattering depends on the cross section of the sample. Fortunately, for most nuclei this involves only s wave scattering, and consequently the scattering lengths and cross sections for thermal neutrons are independent of the neutron wave vector \mathbf{k} (or wavelength $\lambda = 2\pi/|\mathbf{k}|$), whereas the absorption cross sections σ_a are inversely proportional to \mathbf{k} (or inversely with velocity v), and therefore increase linearly with λ . For most nuclides, $V(\mathbf{r})$ is a strongly attractive potential and, therefore, the scattering length is positive. Indeed, there are only a few nuclides, such ^1H , ^7Li , ^{48}Ti , and ^{55}Mn , plus a few others, that have negative scattering lengths.

The general theory of neutron scattering lengths and cross sections has been summarized by Sears,^{7,8} and a compilation of recommended values is given in atomic data and nuclear data tables.⁹ The bound scattering length b of a nucleus is in general complex, given by

$$b = b' - ib'' \quad (7)$$

The scattering cross section is given by $\sigma_s = 4\pi\langle|b|^2\rangle$, where the brackets denote a statistical average over neutron and nuclear spin states, and is divided into coherent and incoherent contributions. The coherent scattering depends on the spatial correlations between scattering nuclei and gives rise to variations in scattered intensity as a function of the scattering vector. This, therefore, gives information regarding the atomic structure of the sample under study. The coherent scattering is responsible for the neutron optical effects that include reflection, refraction, diffraction, and interference. The incoherent scattering on the other hand is independent of the scattering vector and arises from the variance of the scattering lengths at each nuclear site. It depends only on the constituent nuclei and the density of the scattering material. In neutron diffraction measurements, including reflectometry and small-angle scattering, it is the incoherent scattering that gives rise to a featureless background, and its subtraction is important for data analysis, while the coherent scattering gives rise to the interference pattern.

The absorption depends on the imaginary part of the scattering length, and the absorption cross section is given by $\sigma_a = (4\pi/k)b''$. In practice, the imaginary part of the scattering length is of the same order of magnitude as the real part only when the absorption is large ($\sigma_a \approx 10^4\sigma_v$). Values of σ_a are tabulated for a neutron velocity of 2200 m s⁻¹ (equivalent to an energy of 25.3 meV, a wave vector $k \approx 3.49 \text{ \AA}^{-1}$, or a wavelength $\lambda \approx 1.798 \text{ \AA}$). Nuclides, such as ³He, ⁶Li, and ¹⁰B shown in Table 1, with strong neutron absorption have an imaginary part to their scattering length, and are the principal converters from neutrons to charged particles that are found in the majority of neutron detectors (see Sec. 63.8). There are also a few nuclides like ¹¹³Cd or ¹⁵⁷Gd that have low energy (n, γ) resonances, these are used for beam collimation and definition, and can also be used as foil detectors.

Furthermore, the absorption is usually large for only one spin state of the compound nucleus. For nuclides such as ¹¹³Cd, ¹⁵⁵Gd, or ¹⁵⁷Gd, the absorption is large only in the $J = I + 1/2$ state (neutron and nuclear spins are parallel), and $b'' = 0$. For nuclides such as ³He or ¹⁰B, the absorption is large only in the $J = I - 1/2$ state (neutron and nuclear spins are antiparallel), and $b'' = 0$. (This forms the basis of the nuclear-spin polarizing ³He filter; see Sec. 63.7). For a few nuclides such as ⁶Li, the absorption is large for both $J = I \pm 1/2$ states. Finally, the dependence of the scattering length on the particular isotope means that there is a further incoherence for a given element depending on the abundance of the various constituent isotopes.

Consider only elastic scattering and take the potential given in Eq. (3). The differential scattering cross section may be written as the sum of two contributions

$$N \frac{d\sigma}{d\Omega}(\varphi) = \sum_{jj'}^N \langle b \rangle^2 \exp(i\mathbf{q} \cdot \mathbf{r}_{jj'}) + \sum_{j \neq j'}^N (\langle b^2 \rangle - \langle b \rangle^2) \quad (8)$$

where the first term is the differential coherent scattering cross section, dependent on the distance $\mathbf{r}_{jj'}$ between nuclei labeled j and j' , and summed over all nuclei N within the scattering system,

TABLE 1 Scattering Lengths and Cross Sections for ³He, ⁶Li, and ¹⁰B

	$I(\pi)^*$	ϵ^\dagger (%)	b_+^\ddagger	b_-^\ddagger	b_c^\ddagger	b_i^\ddagger	σ_c^\S	σ_i^\S	σ_s^\S	$\sigma_a^{\S\S}$
³ He	1/2(+)	0.00014	4.30	10.07	5.74	-2.5	4.42	1.53	6.0	5333
				-5.93i	-1.48i	+2.57i				
⁶ Li	1(+)	7.5	0.67	4.67	2.00	-1.89	0.51	0.46	0.97	940
			-0.076i	-0.63i	-0.26i	+0.26i				
¹⁰ B	3(+)	20.0	-4.2	5.2	-0.1	-4.7	0.14	3.0	3.1	3835
				-2.49i	-1.07i	+1.23i				

* $I(\pi)$ is the spin and parity of the nuclear ground state.

† ϵ is the natural isotopic abundance.

‡Values of scattering length b are in 10^{-13} cm or fm.

§Values of cross sections σ are in 10^{-24} cm² or barns.

§§ σ_a is the absorption cross section at $\lambda = 1.798 \text{ \AA}$.

Note: b_+ , b_- , σ_c , σ_i and σ_s are defined later in the text, and b_+ and b_- refer to the scattering lengths of the nonzero spin isotopes aligned parallel and antiparallel to the incident neutron spin.

where N is the number of nuclei or atoms, and the second term is the differential incoherent scattering cross section. The coherent scattering length of a mixture of nuclei is the mean scattering length, whereas the incoherent scattering length is the standard deviation of the scattering lengths from that mean. Scattering amplitudes add for coherent scattering, whereas scattering intensities add for incoherent scattering. Incoherent scattering arises from two contributions: one from spin incoherence for those nuclei that have nonzero spin (there are different scattering lengths depending whether the neutron and nuclear spins are parallel or antiparallel), and the other from isotope incoherence. The latter arises because a given element may have different nuclides, each with its own scattering length.

The neutron has a spin $1/2$, and if the nucleus has a nonzero spin I , it may be aligned either parallel or antiparallel to the incident neutron spin s . This gives rise to spin incoherence, and the bound scattering length is spin dependent with

$$b = b_c + \frac{2b_i \mathbf{s} \cdot \mathbf{I}}{\sqrt{I(I+1)}} \quad (9)$$

where b_c and b_i are the bound coherent and incoherent scattering lengths. The coupling of these spins gives rise to two different scattering lengths b_+ and b_- , each with different statistical weights,

$$w_+ = \frac{I+1}{2I+1}$$

and

$$w_- = \frac{I}{2I+1} \quad (10)$$

The bound coherent scattering length for the nucleus with spin I is therefore

$$b_c = \langle b \rangle = w_+ b_+ + w_- b_- \quad (11)$$

and the bound incoherent scattering length is

$$b_i = [\langle b^2 \rangle - \langle b \rangle^2]^{1/2} = \sqrt{(w_+ w_-)(b_+ - b_-)} \quad (12)$$

An important example is ^1H for which the spin dependent scattering lengths are $b_+ \approx 10.85$ fm and $b_- \approx -47.51$ fm. The result is $b_c \approx -3.74$ fm and $b_i \approx 25.27$ fm. The spin incoherent scattering can be observed using polarized neutrons, either by aligning the nuclear spins within the sample or by polarization analysis. Unless otherwise stated, it is assumed that both the incident neutron beam and the nuclear spins are unpolarized.

Each nuclide has a scattering length given by Eq. (9). The total scattering cross section σ_s is the sum of the coherent and incoherent scattering cross sections, σ_c and σ_i , with

$$\sigma_c = 4\pi |b_c|^2$$

and

$$\sigma_i = 4\pi |b_i|^2 \quad (13)$$

unless $I = 0$ in which case $b_i = 0$ and $\sigma_i = 0$. The total scattering cross section is given by

$$\sigma_s = 4\pi \langle |b|^2 \rangle \quad (14)$$

where the brackets $\langle \dots \rangle$ denote a statistical average over the neutron and nuclear spins. The absorption cross section is given by

$$\sigma_a = (4\pi/k) \langle b'' \rangle \quad (15)$$

where k is the magnitude of the incident neutron wave vector. The absorption cross section is determined by the imaginary part of the coherent scattering length. It is only when neutron and nucleus are both polarized that the imaginary part of the incoherent scattering length contributes to σ_a .

Each element may be composed of different isotopes j of natural abundance ϵ_j (such that $\sum_j \epsilon_j = 1$), each with its own coherent and incoherent scattering lengths, b_{cj} and b_{ij} (except that $b_{ij} = 0$ for $I = 0$). The coherent scattering length for the element is given by

$$b_c = \sum_j \epsilon_j b_{cj} \quad (16)$$

the total scattering cross section is

$$\sigma_s = \sum_j \epsilon_j \sigma_{sj} = 4\pi \sum_j \epsilon_j b_{cj}^2 \quad (17)$$

the bound coherent scattering cross section is

$$\sigma_c = 4\pi \left(\sum_j \epsilon_j b_{cj} \right)^2 = 4\pi |b_c|^2 \quad (18)$$

the bound incoherent scattering cross section is

$$\sigma_i = \sigma_s - \sigma_c = 4\pi \sum_j \epsilon_j b_{cj}^2 - 4\pi \left(\sum_j \epsilon_j b_{cj} \right)^2 = 4\pi |b_i|^2 \quad (19)$$

and the absorption cross section is

$$\sigma_a = \sum_j \epsilon_j \sigma_{aj} = (4\pi/k) \sum_j \epsilon_j b_{cj}'' \quad (20)$$

The incoherent scattering, which is independent of scattering angle φ is composed of contributions from the spin and isotope incoherence, where

$$\sigma_i(\text{spin}) = \sum_j \epsilon_j \sigma_{ij} = 4\pi \sum_j \epsilon_j |b_{ij}|^2 = 4\pi \sum_j \epsilon_j w_{+j} w_{-j} |b_{+j} - b_{-j}|^2 \quad (21)$$

and

$$\sigma_i(\text{isotope}) = 4\pi \sum_j \epsilon_j \epsilon_{j'} |b_{cj} - b_{c j'}|^2 \quad (22)$$

The cross sections for a molecule are determined from tables⁹ of scattering lengths and cross sections that give values of b_c , σ_p , and σ_a for each constituent element (where σ_i is the sum of the spin

and isotope contributions) as well as for each isotope. If f_p is the fraction of atoms p within the molecule, the average coherent scattering length per atom is

$$\langle b \rangle = \sum_p f_p b_{cp} \quad (23)$$

so that the coherent scattering cross section per atom is

$$\sigma_c = 4\pi \langle b \rangle^2 = 4\pi \left(\sum_p f_p b_{cp} \right)^2 \quad (24)$$

whereas the incoherent scattering cross section is the addition of the individual cross sections $\sigma_i = \sum_p f_p \sigma_{ip}$, so that the total scattering cross section per atom is

$$\sigma_s = \sigma_c + \sigma_i = 4\pi \left(\sum_p f_p b_{cp} \right)^2 + \sum_p f_p \sigma_{ip} \quad (25)$$

The absorption cross section per atom is the sum of the individual cross sections

$$\sigma_a = \sum_p f_p \sigma_{ap} \quad (26)$$

For a molecule of molecular mass M and macroscopic mass density ρ_M , the molecular number density is given by $N = N_A \rho_M / M$, where N_A is Avogadro's number, $\approx 0.6023 \cdot 10^{24}$ molecules/mol. The macroscopic cross sections therefore become

$$\sum_{\text{coh}} = N \sigma_c = 4\pi N \left(\sum_p f_p b_{cp} \right)^2$$

$$\sum_{\text{inc}} = N \sigma_i$$

and

$$\sum_{\text{abs}} = N \sigma_a \quad (27)$$

where the units of the macroscopic cross section are cm^{-1} , provided that the microscopic cross section σ is in units of 10^{-24} cm^2 (or b), and the scattering length b is in 10^{-12} cm .

Scattering Length Density

While the coherent scattering length is an important quantity determined by the neutron-nucleon strong force interaction, it is the scattering length density of the composition that determines the refraction, reflection, and diffraction properties of the material. If there are s atoms per molecule, the scattering length density of the molecule is given by $sN \sum_p f_p b_{cp} = \rho b$, where $\rho = sN = sN_A \rho_M / M$ is the atom number density, and b is the average scattering length per atom. The product ρb is called the scattering length density of a material. The scattered intensity in small-angle scattering is proportional to the square of the difference in the scattering length densities between two regions of different constituents. For example, this contrast ($\rho_A b_A - \rho_B b_B$) might be from inhomogeneities of molecule A (or pores for which $\rho_A b_A = 0$) in a background or solvent of molecules B.

TABLE 2 Neutron Scattering Lengths and Cross Sections for Hydrogen, Oxygen, and Water

	A^*	$I(\pi)^\dagger$	ϵ^*	b_c^{\ddagger}	b_i^{\ddagger}	σ_c^{\S}	σ_i^{\S}	σ_s^{\S}	$\sigma_a^{\S\ddagger}$
H	1	1/2(+)	99.985	-3.7391		1.7568	80.26	82.02	0.3326
	2	1(+)	0.015	-3.741	25.274	1.7583	80.27	82.03	0.3326
O				5.803		0.4232	0.0	4.232	0.00019
	16	0(+)	99.762	5.803	0	0.4232	0	4.232	0.0001
	17	5/2(+)	0.038	5.78	0.18	0.420	0.004	4.20	0.236
H ₂ O				-0.56		0.118	56.05	56.17	0.222
D ₂ O				6.38		5.11	0.02	5.13	0.0004

* A is the mass number and ϵ is the natural isotopic abundance in per cent.

$^\dagger I(\pi)$ is the spin and parity of the nuclear ground state.

‡ Values of scattering length b are in 10^{-13} cm or fm.

§ Values of cross sections σ are in 10^{-24} cm² or barns.

$^\S\sigma_a$ is the absorption cross section at $\lambda = 1.798$ Å.

An important example is the contrast used in small-angle scattering that can be obtained using deuterated solvents for matching the scattering length density to one part of a given system in order to highlight another part. This is available on account of the greatly different scattering properties of hydrogen and deuterium. The values are given in Table 2. Consider light water, H₂O, with $d \approx 1$ g cm⁻³, and molecular mass 18. The number density of molecules $N \approx 0.03346 \cdot 10^{24}$ cm⁻³. The coherent scattering lengths are $b_H \approx -3.74$ fm and $b_O \approx 5.80$ fm, so that the scattering length density is $(2b_H + b_O)N = \rho b \approx -5.60 \times 10^{-7}$ Å⁻². The average coherent scattering length is $\langle b \rangle = (2/3b_H + 1/3b_O) \approx -0.56$ fm. The coherent scattering cross section for H₂O is $\sigma_c = 3 \times 4\pi \langle b \rangle^2 \approx 0.118$ b per molecule, and $\Sigma_{\text{coh}} \approx 0.0039$ cm⁻¹. This should be compared with the much larger incoherent scattering cross section of light water, 160.52 b/molecule (principally due to the spin incoherence of H, but also the isotope incoherence from the 0.015% of D), and $\Sigma_{\text{inc}} \approx 5.37$ cm⁻¹.

Similarly, consider heavy water, D₂O, which has the number density of molecules $N \approx 0.03329 \cdot 10^{24}$ cm⁻³. The coherent scattering lengths are $b_D \approx 6.67$ fm and $b_O \approx 5.80$ fm, so that the scattering length density is $(2b_D + b_O)N = \rho b \approx 6.37 \times 10^{-6}$ Å⁻². The average coherent scattering length is $\langle b \rangle = (2/3b_D + 1/3b_O) \approx 6.38$ fm. The average coherent scattering cross section for D₂O is $\sigma_c = 3 \times 4\pi \langle b \rangle^2 \approx 15.35$ b/molecule, and $\Sigma_{\text{coh}} \approx 0.511$ cm⁻¹. This compares with the incoherent scattering cross section of heavy water, 4.10 b/molecule, and $\Sigma_{\text{inc}} \approx 0.137$ cm⁻¹.

The scattering length density may be varied over a wide range using a mixture of deuterated and protonated solvents, and in particular light and heavy water. If x is the volume fraction of heavy water in the mixture, the scattering length density is given by $(1-x)(\rho b)_{\text{H}_2\text{O}} + x(\rho b)_{\text{D}_2\text{O}} = (6.96x - 0.56) \times 10^{-6}$ Å⁻¹. As the deuterated fraction increases, the scattering length density increases linearly, and at 8 percent deuteration the coherent scattering is at a minimum, after which it increases quadratically. On the other hand both the incoherent and the absorption cross sections decrease steadily upon increasing deuteration. The matching of the scattering length density (effectively, the index of refraction) of a material with the liquid in which it is immersed gives rise to a sharp minimum in the diffuse scattering.¹⁰ This technique allows the measurement of the refractive index of a wide range of materials by matching the value of the neutron scattering density, ρb , of the material with a mixture of H₂O and D₂O, which can cover a wide range as shown in Fig. 1. This is analogous to the Christiansen filter in classical optics. The refractive index matching has become an important method for determining the macromolecular structure of biological and polymeric systems using neutron scattering by varying the contrast between different parts of the structure and the supporting medium. This is now the standard method for small-angle neutron scattering from biological macromolecules.¹¹ This technique of refractive index matching is also available for nuclear and magnetic components of the refractive index by matching one of the two spin states to the nonmagnetic material in order to highlight the other.

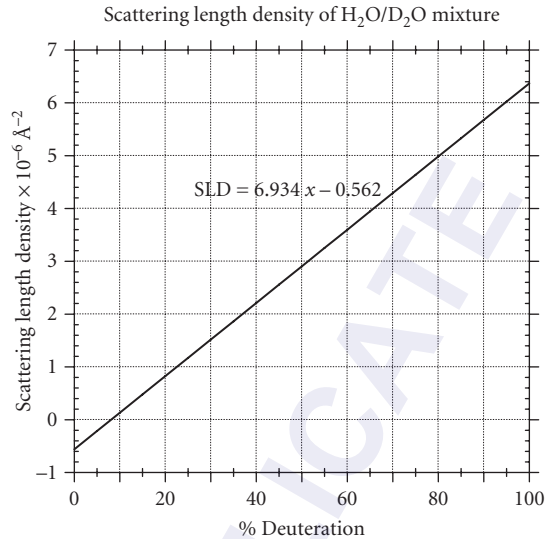


FIGURE 1 The scattering density of water as a function of the fraction x that is deuterated. The scattering lengths of H and D are very different (-0.37×10^{-12} cm and 0.66×10^{-12} cm). The scattering densities of all biological molecules lie between the limits of pure H₂O (-0.56×10^{-12} cm⁻²) and pure D₂O (6.3×10^{-12} cm⁻²) and can therefore be matched by some fraction of deuterated water.

In the case of a mixture of light and heavy water the various scattering properties depend on the fraction x of deuteration. There is H/D exchange between molecules so that there are fractions $(1-x)^2$, $2x(1-x)$, and x^2 of H₂O, HDO and D₂O, respectively. The scattering length density is given by $N[(1-x)b_{\text{H}_2\text{O}} + xb_{\text{D}_2\text{O}}]$. The coherent scattering cross section per molecule is given by $\sigma_{\text{coh}} = 4\pi \times 3[x\langle b_c \rangle_{\text{D}_2\text{O}} + (1-x)\langle b_c \rangle_{\text{H}_2\text{O}}]^2$.

The true incoherent scattering is given by $\sum_p f_p b_{p\text{-inc}}^2 = (1-x)b_{\text{H}_2\text{O}\text{-inc}}^2 + xb_{\text{D}_2\text{O}\text{-inc}}^2$. None of these quantities are changed by the H/D exchange. However, the compositional scattering is given by $\sum_{pp'} f_p f_{p'} (b_{cp} - b_{cp'})^2 = x(1-x)(b_{c\text{H}_2\text{O}} - b_{c\text{D}_2\text{O}})^2/2$. This is smaller by 50 percent of that if there were no H/D exchange.

Neutron Attenuation

The Lambert-Beer law of attenuation $T = \exp(-\mu\ell)$ relating the transmission of radiation through matter also applies to the attenuation of a neutron beam in a given medium of length ℓ . (This equation is exact for absorption and approximate for scattered neutrons, but converges to the exact value as the detector is moved further downstream from the sample.) The attenuation coefficient μ is usually written using the macroscopic cross section $\Sigma = \rho\sigma_T$, where $\sigma_T = \sigma_a + \sigma_s$ is the total cross section. The total transmission through an object that is composed of isotopes i , each with a number density ρ_i , thickness t_i , and total cross section σ_{Ti} , is given by $T = \exp[-\sum(\rho_i\sigma_{Ti}t_i)]$. Because the cross sections vary irregularly between elements in the periodic table, neutron radiography can give good contrast between neighboring elements, unlike x-ray radiography. Certain heavy elements also have

large resonances. Moreover, neutrons have an increased sensitivity relative to x rays for light elements, such as hydrogen, carbon, nitrogen, and oxygen. In particular, the anomalously large scattering cross section for hydrogen makes neutron radiography important for the location of hydrogenous regions within materials, and complementary to x-ray radiography for a number of applications. The technique has been important in studying low temperature hydrogen fuel cells.

Neutron radiography¹² is used in both industrial and fundamental research for nondestructive testing, with tomographic and image reconstruction methods. The spatial resolution of a neutron image is defined as the minimum spatial separation required for resolving two point-like objects that have been blurred by the imaging system. The spatial resolution of the image due to the geometry of the neutron beam arrangement is determined by the aperture diameter D and the aperture-detector separation distance L . A value of $L/D \approx 600$ is nearly optimal when using a standard radiography camera (resolution $\approx 250 \mu\text{m}$). Increasing L/D reduces the neutron intensity and requires longer acquisition times. For higher resolution neutron radiography detectors (resolution $\approx 25 \mu\text{m}$) much higher values of L/D (≥ 1200) are necessary.

63.3 NEUTRON SOURCES

There are several types of sources of thermal neutron beams that have adequate fluxes for useful measurements.^{13,14} (Though note that the resulting neutron beams are characterized by relatively low neutron current densities compared to beams on x-ray synchrotron sources.) In each case these sources produce fast neutrons that must be slowed down or moderated to thermal energies with velocities useful for neutron spectroscopy. The moderator of area $\approx 100 \text{ cm}^2$ becomes the effective source that produces a neutron density with a quasi-Maxwellian-Boltzmann velocity distribution approximately in thermal equilibrium with the moderator. The performance of a neutron source is usually characterized by its brilliance defined by $d^4\Phi/d\Omega d\lambda dA dt$ in units of neutrons $\text{ster}^{-1} \text{ \AA}^{-1} \text{ cm}^{-2} \text{ s}^{-1}$. A fully moderated thermal neutron source produces beams having a Maxwell-Boltzmann velocity distribution, such that the neutron beam has a wavelength spectrum that varies with the neutron wavelength as

$$I(\lambda) \propto \lambda^{-5} \exp[-(\lambda_T/\lambda)^2] \quad (28)$$

where λ_T is the wavelength of a neutron with an energy $k_B T$

$$\lambda_T = h(2mk_B T)^{-1/2} \quad (29)$$

and where k_B is the Boltzmann constant and T is the temperature of the moderator. Fig. 2 shows the fully moderated wavelength spectrum for $T = 293.6 \text{ K}$.

Thermal nuclear reactors produce neutrons using the fission reaction, and are generally steady-state sources, though there are also pulsed reactors. Each fission event releases a huge amount of energy (200 MeV) in the form of kinetic energy of the fission fragments, gamma rays, and several fast neutrons. The fission fragments are the major source of heating within the reactor core, and the heat removal is the limitation to the accessible power at research reactors and therefore to the available neutron beam currents. Furthermore, the gammas and fast neutrons require careful shielding, and to produce a useful thermal neutron beam, the fast neutrons must be moderated by over five orders of magnitude from around 2 MeV to $\approx 25 \text{ meV}$ by some hydrogenous material (usually light or heavy water).

Pulsed sources generally require a particle accelerator (though they may operate in a quasi-continuous mode). An electron linac produces bremsstrahlung photons when the electrons are stopped in a heavy metal target, and these excite resonance interactions to produce evaporation neutrons. (This can also be achieved from the (γ, n) reaction in light nuclei such as D or Be.) Alternatively,

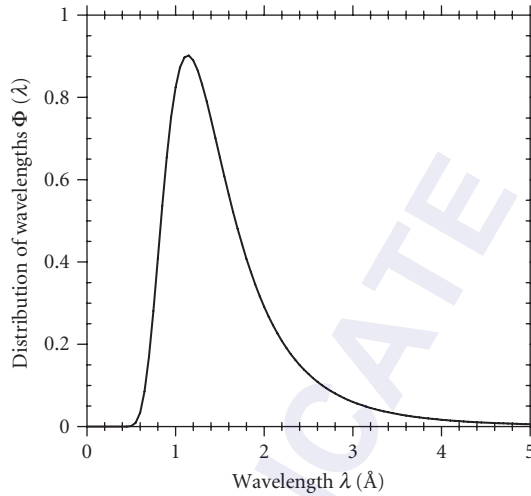


FIGURE 2 The normalized neutron flux distribution of a completely moderated Maxwellian spectrum, $\Phi(\lambda) = 2\Phi_0 (\lambda_T^4 / \lambda^5) \exp[-\lambda_T^2 / \lambda^2]$, corresponding to a temperature T of 293.6 K, or $\lambda_T = 1.798$ Å, or energy 25.3 MeV, or velocity 2200 m s⁻¹. The mean wavelength is $\sqrt{(\pi/2)} \lambda_T = 2.25$ Å, and the most probable wavelength is $\sqrt{(2/5)} \lambda_T = 1.137$ Å (corresponding to an energy of 63.2 MeV). Note that there are always many unwanted fast and epithermal neutrons, whatever the source, that need to be removed from the beam.

beams of high energy (≈ 100 MeV) H⁻ ions produced in a linear accelerator are stripped of their electrons as they are injected into the proton synchrotron ring. After further acceleration to much higher energies (≈ 800 MeV) they smash into the heavy metal target to produce many spallation products including high-energy neutrons. As before, these neutrons must be moderated to thermal energies by some hydrogenous material to provide neutron beams suitable for experimental work.

Whatever source is used, an important consideration is the background and the acceptable signal-to-noise for any neutron measurement on account of high energy gamma rays and fast neutrons ($E > 1$ keV) produced by the source. While thermal neutrons can be stopped by materials with a high absorption cross section (such as cadmium, gadolinium, boron, and lithium), the shielding from fast neutrons requires that they first be slowed down to thermal energies before the capture cross section in one of those materials becomes appreciable. A suitable combination of various materials is required; heavy metals such as iron can enable neutron inelastic collisions, hydrogenous materials such as concrete can moderate the neutron spectrum, and finally absorbing material such as boron can capture the moderated neutrons. Heavy concrete having all three items is often used, but detailed analysis of the shielding¹⁵ requires careful computer simulations for optimum results.

While shielding from the direct gammas from the source (either the reactor core or the spallation target) is often provided by a combination of lead, steel, and concrete, the experimental arrangement must take into account the capture gammas resulting from thermal neutrons absorbed in materials to produce both prompt and long-term activation. Both cadmium and gadolinium emit high-energy gammas upon neutron capture, and so must be used judiciously. Boron also produces a 0.478 MeV gamma, so that lithium is the preferred shielding material in the measurement area. On the other hand gadolinium makes excellent beam definition as a knife-edge aperture, provided other

material shields the capture gammas. Other aperture materials need to be thicker edged and can cause extra diffuse scatter. (Note that both boron and lithium, together with helium, are also neutron detector materials on account of their high absorption cross sections.)

Neutrons produced in a reactor come into thermal equilibrium with the moderator, and have an energy spectrum with a Maxwellian distribution characteristic of the temperature of the moderator. (In the case of a pulsed source thermal equilibrium is only approximate because time-of-flight measurements require a thin moderator to obtain adequate time resolution.) As indicated in Eq. (29), the lower the temperature of the moderator the lower the average neutron energy and the longer the available wavelengths. Consequently, it is best to have a cold moderator (often hydrogen or a mixture of hydrogen and deuterium) to produce large amounts of the long-wavelength neutrons for high-resolution measurements. Indeed, the greater availability of cold neutrons has increased the number of beam lines and experiments that involve neutron optics.

Neutron beams produced at both sources must be transported from the moderator to the experimental hall, either through beam tubes, or better, through neutron guides using total reflection. This allows both thermal and subthermal beams to be brought with little loss to a region of much lower source background. Straight guides look directly at the source, so that either curved guides or neutron filters must be used to remove these source gamma rays and fast neutrons from the beam. The various filters that are used are either Bragg cut-off filters (such as the liquid nitrogen-cooled polycrystalline beryllium filter), or single-crystal filters (such as the sapphire that can be used at room temperature, or magnesium fluoride or bismuth at cryogenic temperatures).

The determination of which source might be more appropriate is dependent on the particular measurement required.¹⁶ Though for some experiments the wide energy spectrum of the polychromatic beam after filtration is sometimes used, generally a monochromatic beam is used for neutron spectroscopy, whereas high-resolution powder diffraction is best performed using the entire spectrum on a pulsed source. However, the most appropriate beam to use depends on the particular measurement envisaged, though most neutron optics research has been performed on cold sources at reactors.

The steady state reactor source can use either a crystal monochromator or a velocity selector to determine the wavelength. The resolution of the measurement can be tailored to the particular problem by suitable choices of collimation and neutron wavelength, and is symmetric and approximately Gaussian. The monochromator crystal selects a narrow wavelength band typically in the range $\Delta\lambda/\lambda \approx 10^{-2}$ from the white beam according to Bragg's law. The velocity selector is usually composed of helical channels formed by absorbing blades that rotate about an axis parallel to the neutron beam. The mean neutron velocity is determined by the rotation speed, and tilting the rotation axis by a small angle enables a change in the coarse wavelength resolution $\Delta\lambda/\lambda \approx 0.1$ to 0.2. Disc choppers composed of absorbing material with slits rotating in phase at high speed on an axis parallel to the beam and placed a distance apart can also monochromate the beam. The rotation speed and the slit widths determine the wavelength resolution, which may be increased by using counter-rotating choppers. Finally a Fermi chopper with narrow curved slits that rotates about an axis perpendicular to the beam can also monochromate the beam.

The pulsed source measurement requires special considerations because it is by necessity performed in time-of-flight mode synchronized to the source frequency. The instrument collects the data by measuring the time interval t taken for the neutron pulse to travel a distance L from the moderator to its detection point. This allows a determination of the neutron energy or wavelength by

$$\lambda = (h/m)(t/L) \quad (30)$$

Since the relative uncertainty in the length of the flight path is small, the time uncertainty $\Delta t/t$ dominates the resolution of the pulsed source instrument. Hence the resolution of the measurement depends on the pulse width Δt produced at the thin (≈ 50 mm) moderator, and is asymmetric and wavelength dependent. Consequently, the neutron moderation time must be kept to a minimum, while maximizing the beam intensities at long wavelengths emanating from the moderator. Heterogeneous poisoning through a layer of neutron absorbing material within the moderator keeps the moderation

properties while reducing the pulse width at low energies. Cooling the moderator to low temperatures shifts the spectrum to long wavelengths while keeping the pulse width narrow. Finally, grooved or fluted moderators can enable greater beam currents by increasing the effective volume for moderation and yet decreasing the effective width for pulse emission at long wavelengths.

63.4 NEUTRON OPTICAL DEVICES

Neutrons stream from the source to the sample through evacuated tubes that might include a number of optical devices. A simple example is the neutron guide that uses grazing incidence reflective optics to transport the neutron beam. Because neutron sources are very weak compared with x-ray synchrotron sources, efforts are made to increase the neutron current density at the sample by various focusing methods. For example, the simple pinhole collimation used in small-angle scattering can be enhanced by multiple confocal pinhole collimations that converge at a point on the two-dimensional area detector. Though the overall current on the sample has increased, the resolution of the measurement is the same as for pinhole collimation. Focusing can be achieved using various optical devices. These elements use the optical characteristics of thermal and cold neutrons, such as diffractive optics (monochromators to select specific wavelengths), reflective optics using grazing incidence (guides transport the beam, mirrors deflect the beam, and various reflection devices can focus the beam), or refractive optics (prisms also refract the beam and concave lenses focus the beam). Many of these optical devices can also be adapted specifically for polarized neutrons. In addition, filters use neutron absorption to limit the range of transmitted wavelengths. However, note that the results are limited by Liouville's theorem that shows that the phase space density of neutrons cannot be increased using energy conserving methods; indeed, it can only decrease, by the use of imperfect devices that absorb, scatter, or reflect fractions of the neutron beam. Consequently, there are limitations on the ability to focus the neutron beam onto the sample, because any increase in neutron current density is attained at the expense of the beam divergence at the sample. If, however, increased beam divergence can be tolerated, even in one dimension, this can be useful. There are other modes of focusing available. Focusing of neutrons in momentum space becomes important for diffraction and scattering applications. There is also time focusing for improving the resolution for time-of-flight measurements¹³ on instruments when the source, sample, analyzer (if any), and detector have extended areas. In this case the orientation of the various components is arranged so that neutrons with different velocities arrive simultaneously at the detector, having had the same momentum change, despite the broad range of velocities within the incident beam.

Neutron Collimation

Neutrons may be collimated by defining the beam path with pinhole or slit apertures in a material of high absorption cross section, such as cadmium or gadolinium, that is essentially black to thermal neutrons. The aperture sizes and the distance between them define the beam direction and the divergence. The transmission may be calculated using phase space acceptance diagrams indicating the position and angle in each of the two orthogonal transverse directions. The maximum angle transmitted relative to the beam centerline is given by the collimation angle. Often collimation is required in only one dimension (that of the scattering plane of the measurement). Soller collimators, composed of a number of long thin blades of neutron absorbing material spaced equally apart, allow large transmission, but with narrow divergence in the one plane. These are often used to define the wavelength resolution of single-crystal monochromator instruments. Mini collimators composed of thin single-crystal wafers coated with gadolinium are also used for the purpose. (Oscillating radial collimators with radial blades have also been used in conjunction with position-sensitive detectors to reduce the background from the sample environment.) The transmission function of an ideal collimator is triangular in shape. The transmission may be increased by the use of

highly polished reflecting surfaces, such that transmission is uniform in position and angle up to θ_c , the critical angle of the surface material, provided the device, now a guide, is sufficiently long. Furthermore, neutron benders have been constructed from stretched mylar blades, coated with copper and separated by spacers.

Neutron Guides

The principle of total external reflection from smooth surfaces at small grazing angles enables neutron guide tubes¹⁷ to be used as channels to transport high-intensity, thermal, or cold neutron beams over relatively long distances (≈ 100 m) with only small losses in intensity. They can supply beams for multiple instruments for neutron scattering experiments at locations of low background with significant improvement in signal-to-noise. Gamma rays and fast neutrons emitted from the source decrease with the square of the distance from the source, whereas the guide maintains the same beam current, neglecting reflectivity losses. These neutron guides are common at research reactors, particularly those with cold sources because their efficiency increases with wavelength. Neutron guides are analogous to fiber optics or light pipes in ordinary optics. A neutron entering the guide tube with an angle of incidence at the wall that is less than the critical angle $\theta_c = \lambda \sqrt{(\rho b/\pi)} = \lambda/\lambda_c$ for its particular wavelength is transported along the hollow tube by multiple total reflections.

The guides are usually made of highly polished boron-glass plates of rectangular cross section that are plated with nickel, which has the highest critical angle per unit wavelength, $\theta_c/\lambda \approx 17.3$ mrad nm⁻¹ or about 1° nm⁻¹, of all elements. Nickel ($b \approx 10.3$ fm) has a scattering length density $\rho b \approx 9.37 \times 10^{-6}$ Å⁻² with $\lambda_c = (\pi/\rho b)^{1/2} \approx 58$ nm and is the reference to which other guides are compared. Isotopically pure ⁵⁸Ni ($b \approx 14.4$ fm) has a greater scattering length density of 1.31×10^{-5} Å⁻² and $\lambda_c \approx 49$ nm, and has a critical angle per unit wavelength, $\theta_c/\lambda \approx 20.4$ mrad nm⁻¹. The use of ⁵⁸Ni therefore results in a 40 percent increase in transmission. Further increases in transmission may be obtained by using glass coated with supermirror layers on top of nickel. The superposition of many diffraction peaks from alternating layers of thin films having a large contrast in scattering length density acts as a two-dimensional crystal to extend the total reflection angle beyond that of nickel. The best contrast is obtained with nickel-titanium layers because Ti is one of the few elements that have a negative scattering length ($b \approx -3.44$ fm). The contrast is given by $\rho b = (\rho b)_{Ni} - (\rho b)_{Ti} \approx (9.37 + 1.95) \times 10^{-6} = 1.13 \times 10^{-5}$ Å⁻². The multilayer consists of a large number of thin layers with a range of d-spacings up to $d = \sqrt{(\pi/4\rho b)}$ such that there is a continuous diffraction using Bragg's law $\theta = \lambda/2d$ at small angles. The resulting supermirror is characterized by the ratio of the effective critical scattering vector $q_c (= 4\pi\theta_c/\lambda)$ to that of natural nickel. Large area coatings with a ratio of 3.6 and with adequate reflectivity are routinely available for use in optical devices and higher values are under development.¹⁸ In practice, there is some loss in reflectivity beyond the nickel total reflection angle, and it drops off considerably at highest values of q_c . This is caused by a number of factors: the roughness of the substrate surface, the interdiffusion between layers, and the stresses induced within the material. Sears¹⁹ has reviewed the theory of multilayer neutron monochromators, and Anderson²⁰ has reviewed the development of supermirrors.

Guides are composed of many sections that might typically be 750 mm long, with cross sections 200 mm high and 60 mm wide. Neglecting any reflection losses, the solid angle of the beam emerging from the exit of the guide is $4\theta_c^2$, providing that the guide entrance is uniformly illuminated in space and angle. The effective collimation achieved by these rectangular neutron guides corresponds to $2\theta_c$ in each orthogonal direction, and therefore depends on wavelength. Straight guides have the maximum transmission (proportional to the square of the wavelength) and require filters to obtain good transmission characteristics without the unwanted radiation. However, many guides are curved to take the beam away from the direct streaming of fast neutrons and gammas, and have good transmission for slow neutrons above a characteristic wavelength. Fast neutrons and gamma rays pass through the glass wall and are eliminated from the beam, so that guides can transport highly collimated beams of thermal neutrons, devoid of gammas and fast neutrons, to regions of low-background far from the source. If the guide sections are

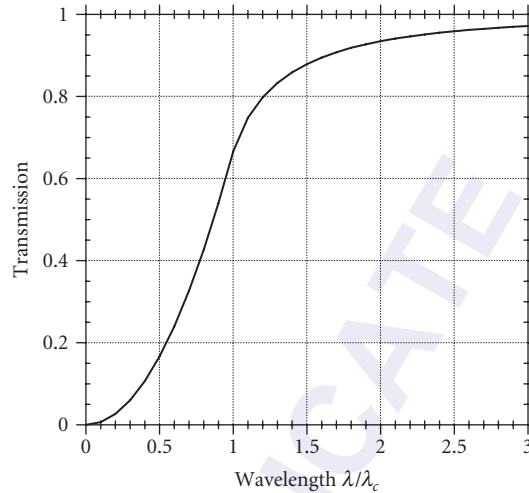


FIGURE 3 The transmission efficiency through a curved guide or neutron bender of length at least $\sqrt{(8wR)}$ as a function of wavelength λ , relative to the straight guide, assuming that the reflectivity of the surfaces of guide coating is 100 percent. For $\lambda < \lambda_c$, $T = (2/3)(\lambda/\lambda_c)^2$, and for $\lambda > \lambda_c$, $T = (2/3)(\lambda/\lambda_c)^2 \{1 - [1 - (\lambda/\lambda_c)^{-2}]^{3/2}\}$.

placed along a polygonal approximation to a curved guide with a gentle curvature, the transmission is defined by a characteristic wavelength $\lambda_c = (\rho b/\pi)^{-1/2} \psi_c$, where the characteristic angle of the curved guide $\psi_c = (2w/R)^{1/2}$, with w and R being the width and the radius of curvature of the guide. For wavelengths less than λ_c , transmission only occurs by garland reflections along the outer (concave) wall. Fig. 3 shows that at λ_c , the transmission is only 2/3 that for the straight guide. For wavelengths greater than λ_c , zig-zag trajectories are also allowed and the transmission increases asymptotically toward that of the straight guide. Consequently, the curved guide acts as a filter provided that the length of the guide is greater than $\sqrt{(8wR)}$ to avoid the direct line-of-sight with the source.

Guides are also found on pulsed source instruments. They maintain the incident beam current density over long distances to enable more instruments to surround the source. However, the pulse frequency, f , limits the wavelength range to $\Delta\lambda = (h/m)(1/fL)$ available on an instrument when the detectors are at a distance L from the source, on account of the frame overlap of the slowest neutrons by the fastest neutrons from the subsequent pulse. Disc choppers and filters often provide the means to overcome this wavelength limitation on a beam line.

The current density of a neutron beam at the sample may be increased by the use of a converging guide. The gain of the converging guide depends on three factors: (1) the convergence of the guide, (2) the ratio of the critical angles on the converging and straight guides, and (3) the neutron wavelength. The greatest gains are achieved when the critical angle of the converging guide is much greater than the divergence of the incoming beam defined by the critical angle of the straight guide. Even with supermirror coatings, intensity gains for the two-dimensional converging guide are limited to perhaps an order of magnitude on account of the increase in absolute angle that a neutron trajectory makes with the device axis upon successive reflections.²¹

More recently, the use of "ballistic" guides reduces losses from nonperfect reflection. Neutrons enter a guide that initially diverges and then converges to its exit. The central part of the guide has a coating of lower q_c (and therefore higher reflectivity) than the entrance and exit of the guide. Neutron trajectories have a smaller number of reflections with higher reflectivity, resulting in a greater neutron

current density than the equivalent straight guide. The details of the transmission of these ballistic guides as a function of wavelength require computer simulation.²²

Neutron Filters

The use of long wavelength neutrons at reactor and pulsed sources often requires a filter to remove epithermal neutrons (energies greater than ≈ 1 eV), and large perfect single crystals can produce a beam relatively free of fast neutron background. They are also used to reduce the higher-order reflections from crystal monochromators. A useful neutron filter material²³ must have wavelength dependent cross sections such that the total cross section is low at thermal energies of interest, but large at epithermal and higher energies. The efficiency of the filter depends on the magnitude of the various cross sections. The absorption cross section is usually linearly dependent on the neutron wavelength and is always independent of temperature. The coherent Bragg scattering cross section depends on the neutron wavelength and the crystal temperature, orientation, and perfection, and can be reduced by suitable orientations and by using highly perfect crystals. The incoherent elastic scattering cross section is usually small and is independent of the wavelength. The inelastic or phonon scattering cross section also varies with wavelength and depends on the crystal temperature T . At low temperatures it varies as $T^{7/2}$ and becomes linear with T at high temperatures. At energies well below $k\Theta_D$, where Θ_D is the Debye temperature of the crystal, the single phonon cross section is also linear dependent on the wavelength. At higher energies ($E > k\Theta_D$), the multiphonon cross section increases with energy and temperature, and rises to the free atom cross section at much higher energies.

Typical single-crystal filters include silicon, quartz (SiO_2), sapphire (Al_2O_3), magnesium oxide, and magnesium fluoride (Table 3). These filters show a minimum in the total cross section around (0.2 to 0.3 nm), with a small linear increase at longer wavelength, and a sharp increase at shorter wavelengths. These materials all have low absorption and incoherent scattering cross sections. MgO and MgF_2 have relatively low Debye temperatures and the cross sections in the thermal region are reduced considerably at cryogenic temperatures. On the other hand there is little to be gained by lowering the temperature of Al_2O_3 from room temperature on account of its high Debye temperature ($\Theta_D = 1032$ K). Polycrystalline beryllium cooled to 77 K is frequently used as a Bragg cutoff filter for neutrons with wavelengths less than 0.4 nm (energies above 5 MeV), since no Bragg scattering is possible for sufficiently long wavelengths such that $\lambda \geq 2d_{\max}$, where d_{\max} is the largest plane spacing. Pyrolytic graphite crystal oriented with the c axis along the incident beam direction allows certain energy windows to exist for which there is no Bragg reflections available to scatter the neutrons. Finally, resonance filters that have high absorption resonance cross sections can eliminate specific wavelengths from the beam to remove $\lambda/2$ contamination from the scattered beam.

TABLE 3 Possible Filter Materials

	M^*	ρ_M^*	N^\dagger	$\sum_i b_i^\ddagger$	ρb^\S	Σ_{fa}^\P
SiO_2	60.08	2.65	0.0266	1.576	4.19	0.253
Al_2O_3	101.96	3.98	0.0234	2.431	5.71	0.328
MgO	40.31	3.58	0.0535	1.118	5.98	0.384
MgF_2	62.32	3.18	0.0307	1.668	5.13	0.334

* M is the molecular mass in g/mol and ρ_M is density in g cm^{-3} .

$^\dagger N$ is the molecule number density in 10^{24} cm^{-3} .

$^\ddagger \sum_i b_i$ is the molecular coherent scattering length in 10^{-12} cm .

$^\S \rho b = N \sum_i b_i$ is the scattering length density in 10^{-6} \AA^{-2} .

$^\P \Sigma_{fa}$ is the free atom macroscopic cross section in cm^{-1} at high energies.

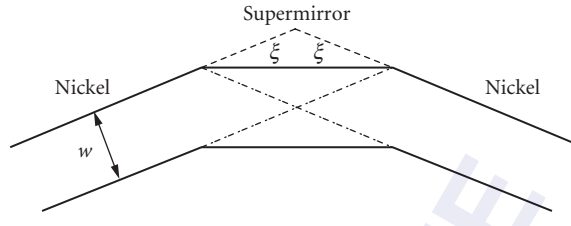


FIGURE 4 A schematic drawing of a neutron optical filter in which the initial long guide of nickel, of width w , is followed at an angle, ξ , by a short length of supermirror guide, and then another nickel guide at a further angle, ξ . If the angle ξ is made equal to the critical θ_c of nickel at the shortest wavelength of interest, then the supermirror must have a critical angle of at least $2\theta_c$.

A different concept is the neutron optical filter,²⁴ in which the beam may be deflected by a small angle with a high transmission above a cutoff wavelength, without the spatial asymmetry that is introduced by the curved guide. This is illustrated in Fig. 4. An initial long nickel guide is followed by a short length of supermirror guide, offset by a small angle, ξ , after which the nickel guide continues at a further small angle, ξ . The length of the intermediate supermirror guide is such that there is no direct line of sight through the entire system. The advantage of such a geometry is that it has the property that a parallel beam is transported unchanged, unlike the case of a curved guide. The supermirror critical angle must be at least twice that of nickel. A modification is to have the central section tapered and with a different critical angle. Consequently, this type of filter can be designed for particular experimental arrangements.

63.5 REFRACTION AND REFLECTION

All isotopes have a scattering length b that characterizes the neutron-nucleus interaction (albeit some have a complex value).⁷ Consequently, all materials have an index of refraction that depends on the scattering length of the isotopes that compose the material. The basic difference is that, as for x rays, the index of refraction for neutrons has a value of $n < 1$, whereas $n > 1$ for light. Moreover, the neutron wavelengths are about three orders of magnitude smaller than for visible light. This means that in general much larger instruments are required to observe the optical phenomena with the same spatial resolution. In addition, the most intense thermal neutron sources are orders of magnitude weaker than conventional light and x-ray sources. This means that the phenomena are generally not as easy to observe, nor is the resolution as good. However, successful neutron optical devices have been developed and are now commonplace in instrumentation for scientific research studies, particularly those that use cold- or long-wavelength neutrons.

The index of refraction of a medium is defined by $n = k'/k = \lambda/\lambda'$, and for slow neutrons is given by

$$n^2 = 1 - \frac{\lambda^2 \rho b}{\pi} \quad (31)$$

where $\lambda = 2\pi/|\mathbf{k}|$ is the incident neutron wavelength, and ρ is the average atom density of the material. Since $V/E \ll 1$, Eq. (31) may be approximated by $n = 1 - \lambda^2 \rho b / 2\pi$. For thermal neutrons with wavelength $\lambda > 0.1$ nm, the refractive index for most materials differs from unity by a few parts in 10^6 . Consequently, the various optical features are similar to those in x-ray optics, and various analogous features have been demonstrated. However, the small deviation of the refractive index from unity means

that focusing through either refraction or reflection is weak. Note that at the boundary between two media, A and B, the relative index of refraction is given by

$$n^2 = 1 - (\lambda^2/\pi)(\rho_A b_A - \rho_B b_B) \quad (32)$$

where $\rho_A b_A$ is the scattering length density of medium A.

A critical wavelength is defined by $\lambda_c = (\pi/\rho b)^{1/2}$. If b is real and positive, the index of refraction is real if $\lambda < \lambda_c$ and imaginary if $\lambda > \lambda_c$. The phenomenon of total reflection for all angles of incidence occurs when n is purely imaginary, that is when $\lambda > \lambda_c$. Typically for materials $\lambda_c \approx 140$ nm, which corresponds to a neutron energy of 4×10^{-8} eV or a temperature of 0.5 mK. It is therefore only for ultracold neutrons²⁵ that true mirror reflection can occur. Mirror reflections for ultracold neutrons can be used to confine neutrons in material cavities or bottles. Smooth, clean surfaces are needed for good storage, though this is limited by β decay of the free neutron (with a half-life of ≈ 887 s). In practice, the storage time is shorter because of losses caused by inelastic scattering from surface impurities. Ultracold neutrons are useful for precision measurements of the free neutron lifetime, the neutron magnetic dipole moment and the search for its electrical dipole moment.

Note that the usual form of the refractive index is given by Eq. (31). This neglects local field effects (the nucleus scatters not only the incident wave but also waves scattered by all other nuclei). It also neglects the attenuation of the amplitude due to diffuse scattering of the coherent wave in the medium. Sears¹ has given a more rigorous and comprehensive treatment of dispersion theory. Fluctuations in the scattering length density may be due to both number density fluctuations as well as to the effects of spin and isotope disorder that produce fluctuations in the scattering length. Coherent scattering gives rise to the neutron optical effects that include reflection, refraction, diffraction, and so on, and depend on the orientation of the system relative to the incident beam, whereas diffuse incoherent scattering is distributed isotropically.

Refraction and Mirror Reflection

The total external reflection of neutrons from material surfaces has been demonstrated, analogous to the total external reflection of x rays and the total internal reflection of light. Because the index of refraction is close to unity, total reflection is possible only at glancing incident angles, θ , on a material surface when $\cos \theta > n$. This is the neutron analogue to Snell's law for light (the continuity of the tangential component of the wavevector). The measurement of the critical angle for various surfaces using monochromatic neutrons can be used for the determination of scattering lengths with a relative uncertainty of ≈ 0.2 percent. The reflectivity, R , is given by Fresnel's law,

$$R = \left| \frac{1 - [1 - (\theta_c/\theta)^2]^{1/2}}{1 + [1 - (\theta_c/\theta)^2]^{1/2}} \right|^2 \quad (33)$$

where $\theta_c = \lambda/\lambda_c = \lambda(\rho b/\pi)^{1/2}$, for nonmagnetic materials. If $b > 0$, $R = 1$ for $\theta < \theta_c$, the critical angle for total reflection. Because $(1 - n^2) \approx 10^{-5}$ for thermal neutrons, $\theta_c \approx 3$ mrad $\approx 0.2^\circ$. Hence, the measurement of the reflectivity can be used to determine both the magnitude and sign of the bound coherent scattering length.

The mean optical potential of neutrons in most materials is comparable in magnitude with the gravitational potential energy corresponding to a height difference of order of 1 m. This is the principle of the gravity refractometer,²⁶ in which neutrons from a well-collimated initial horizontal beam fall under gravity in a long flight path onto the horizontal surface of a liquid. All neutrons achieve a critical vertical velocity at a height $h_0 = (2\pi\hbar^2/m^2g)\rho b$ and penetrate the surface of the liquid; otherwise, they will be totally reflected by the mirror and be detected in the counter. Consequently, h_0 is a measure of the quantity ρb , and hence of the scattering length, b . This experiment, illustrated in Fig. 5 demonstrates the unique particle-wave properties of neutron optics, with the motion obeying classical physics and the refraction obeying quantum physics. The significant point is that the

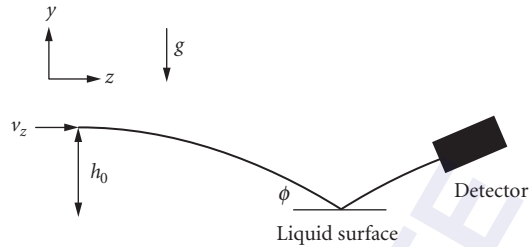


FIGURE 5 A schematic diagram of the neutron gravity spectrometer. A neutron of initial horizontal velocity v_z and vertical velocity $v_y = 0$ falls a vertical distance h and strikes the liquid surface at a grazing angle of $\phi = \sqrt{2gh}/v_z$ after traveling a horizontal distance $v_z \sqrt{2h/g}$. The critical height for total reflection is $h_0 = (2\pi\hbar^2/m^2g)\rho b$, independent of v_x or λ .

measurement is independent of the neutron wavelength. This allows the use of a more intense neutron beam to achieve a measurement with high statistical accuracy, resulting a most accurate method to determine b for a liquid.

Grazing-Angle Reflection

Focusing systems based on reflective optics do not suffer from the chromatic aberrations of refractive optics. For example, the advantage of an ellipsoidal mirror is that, neglecting gravity, all the trajectories emanating from one focus (the source) reach the other focus regardless of the neutron wavelength. However, single bounce mirror systems suffer from high aberrations. The difficulty with grazing incidence optics is that a near parallel beam incident is required and so that the optical element must be placed far enough from the source to ensure small incident beam divergence. Kirkpatrick-Baez neutron mirrors²⁷ using two successive reflections in orthogonal directions can efficiently focus neutron beams into small areas with a maximum divergence that is limited by the mirror critical angle. The size of the focal spot is primarily determined by geometrical demagnification of the source and by figure errors in the mirror shape. Approximately two orders of magnitude in neutron current density increase can be achieved within a spot of diameter $\approx 100 \mu\text{m}$ using crossed mirrors.

Grazing angle reflective optics based on two-bounce Wolter geometries are used extensively in x-ray astronomy because they minimize optical aberrations for off-axis trajectories (see Chap. 64). They can be also designed to focus cold and thermal neutron beams using consecutive reflections from parabolic and hyperbolic surfaces. The tilt angle of the hyperbolic section is 3 times larger than the tilt angle of the parabolic section to preserve the grazing angle of the trajectory. The mirrors are fabricated using an electroformed nickel replication process. They have a cylindrical form with different diameters but with the same focal length, such that they can be nested to increase the system throughput. Nested replicated optics can improve the focused beam intensity by increasing the incident area of the quasi-parallel beam.

Polycapillary Optics

Polycapillary optics (see Chap. 53) provides a better method for increasing the neutron current density at a point by using many narrow guides. This allows the possibility of a greater curvature than for macroguides for a given wavelength because the transmission characteristics of a curved guide

depend on the ratio of the transverse to longitudinal dimensions of the guide. These miniature versions of neutron guides can be used to transport, bend, and even focus thermal neutron beams. Glass polycapillary fibers with thousands of hollow channels recently developed for transporting and focusing x rays can also be used for cold (≈ 0.4 nm) neutron beams. The advantage of the narrow (≈ 10 μm) channels is that the neutron beam can be bent more sharply than for the wide (≈ 50 mm) guides. The transmission properties of the capillaries depend on the internal diameter and bending radius of the channels, the glass composition, and the surface smoothness.

The focusing of neutrons comes at a price. Liouville's theorem requires that an increase in neutron current density is obtained with a necessary increase in beam divergence. This means that real-space focusing has limited applications in condensed matter research using neutron scattering for which angular divergence is important, whereas neutron absorption techniques in analytical and materials research do not require a monochromatic neutron beam or a small angular divergence. Consequently, a high-intensity beam focused by a neutron capillary lens onto a small sample area can be useful for improving both the detection limits of individual elements and the spatial resolution of the measurement. Polycapillary fiber lenses have been produced and tested for analytical applications.²⁸ The gain depends on the divergence and the energy spectrum of the incident beam. Within the typical focal spot size of 0.5 mm, limited by the outer diameter of the fibers, the gain can be nearly two orders of magnitude. For the monolithic lens, a fused version of the polycapillary lens, the focal spot size is less than 200 μm with a gain of about 10.

Prisms and Lenses

The prism deflection of neutrons has been observed.²⁹ Because the refractive power, $(n - 1)$, is negative for most materials, prism deflection is in the opposite direction compared with ordinary light optics. Deflection angles of a few seconds of arc can be measured with high accuracy, allowing the measurements of refractive indices to 1 part in 10^4 . Because $(n - 1)$ is proportional to λ^2 , prism deflection of neutrons is highly dispersive. It has been used in a fused-quartz converging lens (which is concave in shape) and the focusing of 20 nm neutrons has been demonstrated.³⁰ More recently,³¹ the focusing of shorter (9 to 20) nm wavelength neutrons has been demonstrated using a series of many closely spaced symmetric biconcave MgF_2 lenses originally produced for infrared optics, with gains of approximately 15 and focal lengths of (1 to 6) m, which is well matched to small-angle neutron scattering. Magnesium fluoride is preferred to quartz because it has a higher scattering length density ($(\rho b)_{\text{MgF}_2} = 5.13 \times 10^{-6} \text{ \AA}^{-2}$, $(\rho b)_{\text{SiO}_2} = 4.19 \times 10^{-6} \text{ \AA}^{-2}$). This compound refractive lens is analogous to that for x rays (see Chap. 37).

Focusing now plays an increasing role in small-angle neutron scattering instruments in order to attain lower values of scattering vector without reducing the current density on the sample, by decreasing the extent of the beam penumbra at the detector. Pinhole collimation using the optimum equal-flight-path configuration, with the source aperture size twice that of the sample, results in a beam profile that is conical at the detector. The intensity on sample is proportional to the area of both apertures. A focusing optic consisting of a series of biconcave lenses enables the direct beam profile at the detector to be both narrower and uniform. When the detector is placed at the image position of the source with respect to the optic the beam profile at the detector is uniform, and its size is independent of the size of the sample aperture.

The focal length for a set of N biconcave lenses of radius of curvature $|r|$ is given by

$$f = |r|/2N(1 - n) \quad (34)$$

where the index of refraction is given by $n = 1 - \rho b \lambda^2 / 2\pi$. Hence the focal length of the biconcave lens is $f = \pi |r| / (N \rho b \lambda^2)$. Since f varies inversely with the square of the wavelength, refractive lenses are strongly chromatic, with large chromatic aberrations. For perfect focusing with the lens placed immediately before the sample, the usual lens equation is valid with $1/L_1 + 1/L_2 = 1/f$, where L_1 and L_2 are the incident (source to sample) and secondary (sample to detector) flight paths, respectively.

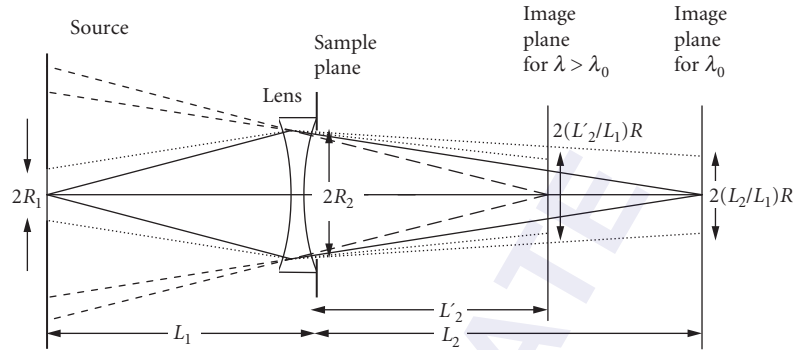


FIGURE 6 A schematic diagram of a focusing lens arrangement with source of radius R_1 at a distance L_1 from the sample with an aperture of radius R_2 . The focusing lens with a focal length f_0 is placed in front of the sample such that the source is imaged (continuous lines) with a radius $(L_2/L_1)R_1$ at a distance L_2 from the lens for a wavelength λ_0 . For another wavelength $\lambda (> \lambda_0)$ the source is imaged (long dashed lines) with a radius $(L'_2/L_1)R_1$ at a distance $L'_2 (< L_2)$ such that $1/L'_2 - 1/L_2 - 1/f_0[(\lambda/\lambda_0)^2 - 1]$. (See also color insert.)

However, the chromatic aberration determines that this is only true at one wavelength. If f_0 is the focal length of the lens at a wavelength λ_0 such that the system is focused at that wavelength, then the focal length of the lens at some other wavelength λ is given by $f = f_0(\lambda_0/\lambda)^2$. This is illustrated in Fig. 6. Since the index of refraction of material lenses is very close to unity, the focusing is weak. The radius of curvature of the lens must be small, and a large number of lenses means that the part of the beam travels through a large thickness of material. If a Fresnel lens is used the thickness can be reduced considerably, decreasing the attenuation. This is particularly useful for cold neutrons.

The use of both long wavelength neutrons and long flight paths results in the transmitted beam falling under the influence of gravity by an amount that is wavelength dependent. That is, the neutron's vertical position changes by an amount $y = -(g/2)(m\lambda/h)^2L$ over a distance L , where g is the acceleration due to gravity. Since the incident beam has some divergence, this smearing results in an oval shape to the neutron beam spot at the detector. A prism can compensate for the chromatic aberration caused by gravity because the refractive index for neutrons has a wavelength dependence ($\sim \lambda^2$) similar to that for gravity. An "antigravity device" has been demonstrated using single-crystal prisms such that the transmitted beam is restored to the instrument centerline defined by the source and sample aperture centers. Others have studied both magnetic prisms and the compound refractive prism made of single crystal elements.

63.6 DIFFRACTION AND INTERFERENCE

The Bragg diffraction by single crystals, powders, amorphous materials, and liquids is analogous to that for electromagnetic waves, with formulas similar to those for x-ray diffraction except the magnitude is determined by the coherent scattering lengths. The diffraction pattern is simply a function of scattering angle and is proportional to the square of the Fourier transform of the structure of the diffracting sample. The small-angle neutron scattering from inhomogeneous materials also has expressions similar to those for x rays. The availability of the (partial) substitution of hydrogen by deuterium allows the ability to change the relative scattering powers for different features within the structure.

Most neutron scattering instruments use a large crystal monochromator to select neutrons of a particular wavelength λ from the polychromatic beam from the source to be incident onto the sample.

The mean wavelength reflected from the crystal depends on the Bragg angle θ and the monochromator crystal plane spacing d according to Bragg's law, $v\lambda = 2d \sin \theta$, where v is the order of diffraction. Usually the first order ($v = 1$) is desired because the reflection is most intense. A neutron filter may be desired to remove the higher-order contamination. Monochromators are characterized by their peak reflectivity and their mosaic width obtained by a rocking curve, rotating the crystal between two perfect crystals. The intrinsic angular reflectivity width for a particular wavelength of perfect crystals, called the Darwin width, is far too narrow for their use as neutron monochromators in diffractometers, and though the peak reflectivity is very high, the integrated reflectivity is low. It is necessary to deform the crystals in order to introduce dislocations to form many slightly misaligned crystals, called a mosaic crystal. Alternatively, the perfect crystal may be bent such that the crystal orientation varies across the face of the crystal, such that the crystal lattice spacing has a gradient across its depth (see Chap. 39 for x-ray monochromators). In both cases, the mosaic of the crystal is described by the width of the Gaussian angular crystal plane distribution obtained in a rocking curve measurement. The reflectivity may be reduced by extinction and absorption, and sometimes by simultaneous parasitic scattering. However, perfect crystals are used for some applications, such as high-resolution backscattering instruments, neutron interferometry, and ultra small-angle scattering.

The wavelength resolution of the diffracted beam depends on the angular collimation both before and after the monochromator and on the mosaic spread of the crystal. For maximum reflected intensity at a given resolution it is usual to make the collimations comparable, and the mosaic spread large. Typical crystal mosaic spreads are between 0.2° and 0.5° . The collimation out of the scattering plane should be as relaxed as possible to increase the reflected intensity. Typical neutron monochromators are graphite, silicon, germanium, copper, and beryllium. Monochromators may be used in reflection (Bragg) geometry or in Laue (transmission) geometry. Neutron monochromator efficiencies are higher in reflection geometry than in transmission geometry.

Exact analytic general solutions^{32,33} of the Darwin equations have been derived that describe the multiple Bragg reflection of x rays or neutrons in a mosaic slab crystal. Both the intensity reflection coefficient R and the transmission coefficient T , where $R + T = 1$, can be expressed in compact form for both the Bragg case (reflection geometry) and the Laue case (transmission geometry). For the Bragg case

$$R = b / \{ [a(a+2b)]^{1/2} \coth [a(a+2b)]^{1/2} + (a+b) \}$$

$$T = \frac{[a(a+2b)]^{1/2}}{[a(a+2b)]^{1/2} \cosh [a(a+2b)]^{1/2} + (a+b) \sinh [a(a+2b)]^{1/2}} \quad (35)$$

And for the Laue case

$$R = 1/2 \exp(-a) [1 - \exp(-2b)]$$

$$T = 1/2 \exp(-a) [1 + \exp(-2b)] \quad (36)$$

The dimensionless quantities are defined by $a = \mu t / \sin(\theta + \alpha)$ and $b = \sigma t / \sin(\theta + \alpha)$, where the angles θ and α , respectively, are the Bragg angle and the angle between the reflecting planes and the crystal surface, and t is the thickness of the monochromator crystal. The linear coefficients σ and μ , respectively, are the Bragg reflection coefficient, and the attenuation coefficient for all other processes other than Bragg reflection, including absorption, incoherent scattering, and coherent inelastic scattering. The quantity $\mu = \rho(\sigma_a + \sigma_i)$, where σ_a and σ_i are the absorption and incoherent scattering cross sections and ρ is the atom number density. The quantity σ depends on the incident neutron wavelength λ and the Bragg angle θ and is nonzero only if Bragg's law is satisfied. It is given by $\sigma = Q_c W(\theta - \theta_{\text{hkl}})$, where $W(\theta - \theta_{\text{hkl}})$ is a normalized rocking curve, and is sharply peaked at $\theta = \theta_{\text{hkl}}$, and its shape is characteristic of the mosaic structure of the crystal with Miller indices hkl . Here $Q_c = \lambda^3 |F_{\text{hkl}}|^2 / v_0^2 \sin 2\theta_{\text{hkl}}$, where F_{hkl} is the structure factor for the unit cell of volume v_0 . This enables the calculation of the reflectivity of an absorbing crystal of finite thickness in situations

where the Bragg planes make an arbitrary angle with the surface of the crystal. The symmetric case is the usual monochromator geometry ($\alpha = 0$) that results in a reflected intensity with an asymmetric intensity distribution across the beam. The asymmetric case is known as a Fankuchen-cut crystal ($\alpha \neq 0$), and for $\alpha > 0$ the beam is compressed and the diffracted beam is more intense.

A considerable gain in intensity for a neutron spectrometer is obtained with an improved monochromator with an anisotropic mosaic structure by focusing in both the horizontal and the vertical directions, to form a double focusing monochromator³⁴ using many crystals assembled together, with little degradation of instrumental resolution. Vertical focusing increases the intensity for a diffraction measurement without affecting the resolution. This is achieved by placing the crystals along an arc of radius of curvature R_V in the vertical plane given by

$$R_V = \frac{2L_1L_2 \sin \theta_B}{L_1 + L_2} \quad (37)$$

where L_1 and L_2 are the distances of the effective source and the sample from the monochromator, or the distances from the sample and the detector if the crystal is used as an analyzer. Horizontal focusing involves Bragg crystal optics (see Chap. 22) and requires the bending of long single crystals in a suitable mechanical device, or placing smaller crystals along an arc of radius of curvature R_H in the horizontal plane given by

$$R_H = \frac{2L_1L_2}{(L_1 + L_2) \sin \theta_B} \quad (38)$$

Diffraction Effects

Various important measurements have illustrated the wave nature of neutrons. Shull³⁵ has demonstrated Fraunhofer slit diffraction using a monochromatic beam of neutrons diffracted from a perfect single crystal of silicon through a slit that was analyzed by a similar crystal. The diffraction broadening of the transmitted beam was shown for several very narrow slit widths to agree with the classical formula $(\sin x/x)^{1/2}$ for diffraction from a single slit. Later experiments³⁶ with greater precision performed for single slit and double slit geometries have showed excellent agreement with wave diffraction theory. The diffraction of thermal neutrons by ruled gratings at grazing incidence has also been demonstrated.³⁷

The Fresnel diffraction of an opaque straight edge³⁸ has been found to be in excellent agreement with theory. The focusing of slow neutrons by zone plates (phase gratings of variable spacings) has been demonstrated.^{30,39} Conventional zone plates, analogous to those for x rays (see Chap. 40), use constructive interference of diffracted waves through open zones, with absorption of radiation in dark zones. Neutron zone plates use phase reversing in alternate zones and can achieve greater focusing efficiency. Zone plates can be produced by photolithography on silicon wafer substrates from a computer-drawn pattern. The resolution is limited by the inherent chromatic aberration of the zone plate. Cylindrical zone plates have been used as a condensing lens.⁴⁰

Interference

There have been various demonstrations of the interference of coherent beams using thermal, cold, and ultracold neutrons.^{41,42} Two beam interferometers allow the determination of the relative phases of two coherent wavefronts and the investigation of different influences, such as force fields and material media, on the phase of a wave. This is again analogous with classical optics. When neutrons are incident on perfect crystals under Bragg conditions, there is a dynamic interchange of flux between the incident beam and the Bragg diffracted beam inside the crystal medium. The radiation density traveling in either the Bragg or the forward direction is a periodic function of depth

in the crystal. This results in extinction effects, with an associated primary extinction length of order (1 to 100) μm . The dynamical theory of diffraction leads to anomalous transmission effects, and there have been numerous experiments using perfect crystals to observe these special effects of neutron wave propagation in periodic media. For example, Shull⁴³ has observed the Pendellösung interference fringe structures in the diffraction from single-crystal silicon.

The superposition of two spatially separated parts of one coherent wavefront produces interference patterns from which the relative phases of the two parts can be deduced. The effective source for the original wavefront is a beam with a spatial coherence produced by passage through a narrow slit from which a partially coherent cylindrical wave emerges. This beam illuminates a pair of closely spaced narrow slits, and the diffracted radiation gives rise to interference fringes in the plane of observation,³⁶ analogous to Young's slits in classical optics.

A similar experiment⁴⁴ has shown the interference of two spatially separated parts of one wavefront using the boundary between two domains of opposite magnetization in a crystal of Fe-3%Si. In traversing the ferromagnetic foil on either side of a domain wall in the sample, the spin of the neutron precesses in opposite directions. The two parts of the wavefunction acquire a relative phase shift that shows up in the interference pattern in the plane of observation. The expected pattern is the superposition of two Fresnel straight-edge patterns that for zero phase shift combine to give a uniform distribution. For π phase shift, destructive interference occurs in the center of the pattern, giving a deep minimum. The relative phase shift depends on the angle of precession of the neutron spin, which is proportional to the thickness of the region of the magnetic field traversed. The results of the experiment demonstrate that a phase shift of π (destructive interference) occurs when the spin is rotated by odd multiples of 2π , verifying the distinctive behavior of spin 1/2 particles under rotation, that the neutron is a spinor.

The coherent splitting of the amplitude of a wave by means of partial reflection from thin films gives rise to interference effects in visible light. For neutrons this partial reflection from the front and back surfaces of various thin films occurs at glancing angles close to critical reflection. The resultant interference phenomena give valuable information about the structure of the interfaces involved.⁴⁵ Various neutron optical devices based on thinfilm interference have been developed, such as bandpass monochromators, supermirrors that are highly efficient neutron polarizers,⁴⁶ reflectivity and polarizing characteristics of Fe-Ge multilayers,⁴⁷ and devices involving multiple-beam interference effects, analogous to classical optics (see Chap. 41 on multilayers for x rays). Multiple-beam interference has been demonstrated⁴⁷ with ultracold neutrons in which thin films of copper and gold are used in the construction of a resonant structure analogous to a Fabry-Perot cavity.

Perfect Crystal Interferometers

Bonse and Hart⁴⁸ have pioneered the x-ray interferometer that uses Bragg reflection by crystal flats to split and recombine two coherent beams. The crystal flats are cut from a monolithic piece of large, highly perfect single-crystal silicon leaving a central backbone. This ensures that all the crystal planes in each crystal flat are perfectly aligned, since each remains part of the original monolithic crystal. This idea has been applied to neutron interferometry, and involves triple Laue transmission geometry optics.⁴⁹ This design is analogous to the Mach-Zehnder interferometer of classical optics. The interferometer illustrated in Fig. 7 comprises three identical, perfect crystal flats cut perpendicular to a set of strongly reflecting planes, with distances between flats usually a few centimeters. There are of order 10^9 oscillations of the neutron wave function in each beam path, so that there are stringent requirements on the microphonic and thermal stability of the device. This requires vibrational isolation of the interferometer from the reactor hall, plus thermal isolation. The three crystal slabs must be aligned within the Darwin width of a few seconds of arc for thermal neutrons in silicon. The precise definition of the wavelength is accomplished by the Bragg diffraction of the interferometer. Placing an object of a given thickness in one of the beams leads to a phase shift or change in the interference contrast. The path difference must be small compared with the coherence length of the beam that depends on the degree of beam monochromatization. There have been many applications of neutron interferometry in fundamental physics, such as gravitationally induced quantum interference and tests of nonstandard quantum mechanics.^{41,42}

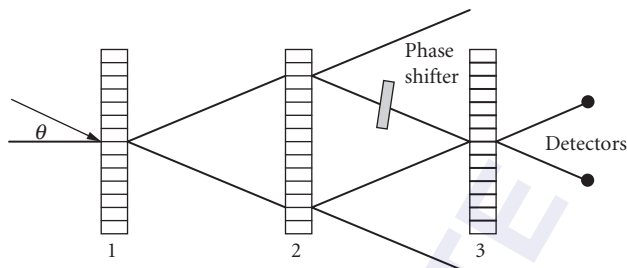


FIGURE 7 A schematic diagram of the triple Laue neutron interferometer cut from a monolithic perfect crystal of silicon in which the first slab splits the incident beam into two coherent parts, with the second slab acting as a mirror to bring the beams together at the third slab. Interference fringes are observed in the detectors by rotating the phase shifter in one of the paths.

The analog of the classical Rayleigh interferometer is a two-crystal monolithic device cut from a large silicon crystal.⁵⁰ The incident beam is restricted by a narrow entrance slit and is Bragg-reflected by the first silicon crystal slab. Inside the crystal, the diffracted neutron beam fills the entire Borrmann triangle, and the beam that leaves the back face of the crystal is broad, having a width that is dependent on the crystal thickness. There is a conjugate focal point on the exit face of the second crystal, from which neutrons along the complementary interfering rays leave the device in the incident and diffracted beam directions. If the central rays are blocked by a cadmium absorber, the parallelogram ray diagram is analogous to the three-crystal interferometer. Another variant uses Bragg reflection geometry rather than Laue transmission geometry. The incident beam reflected from the front surface of the first crystal is brought back to interfere with the beam originating from the back face of this crystal by using a second parallel Bragg reflecting crystal. The advantage is that the spatial definition of the outgoing interfering beam is as sharp as the beam incident on the interferometer.

Perfect crystal neutron interferometry is a convenient and most precise technique for measuring coherent scattering lengths solid, liquids, and gases. It has the advantage of not being limited to liquids as the gravity reflectometer. The period of the intensity oscillations as a function of the orientation of a slab of material that is rotated about an axis perpendicular to the plane of the interferometer enables the determination of the refractive index and hence the scattering length of the slab material. There is also a differential technique using a phase shifter in one of the beams. Neutron interferometry measures the difference in the average neutron-nuclear potential of the sample traversed by the two beams. Intensity oscillations are observed by varying the atom density, with more rapid interference oscillations for larger scattering lengths. The sensitivity of the technique allows a measurement of the hydrogen content in various transition metals at a level of about 0.05 atom fraction.⁵¹

63.7 POLARIZATION TECHNIQUES

Neutron Polarization

It is the spin of $1/2$ of the neutron that results in neutron optics having phenomena quite different from x-ray optics, and enables neutrons beams to be polarized. The neutron spin, \mathbf{s} , interacts with the nuclear spin, \mathbf{I} , and the optical potential, $V(\mathbf{r})$, and the bound scattering length, b , are both spin-dependent for $I \neq 0$. The total spin, $\mathbf{J} = \mathbf{I} + \mathbf{s}$, is a constant of the motion with eigenvalues $J = I \pm 1/2$. The interaction potential of the magnetic dipole moment μ_n of the neutron in a magnetic field \mathbf{B} is

given by $V_{\text{mag}} = -\boldsymbol{\mu}_n \cdot \mathbf{B}$. When a neutron beam traverses a magnetic field, the field defines a quantization axis such that the neutron spin is either “up” or “down” relative to the field direction. That is, the spins of neutrons may be oriented such that their projections are either parallel (+) or antiparallel (–) to the magnetic field. Because there are two possible spin states, a beam with a fraction, ϕ_+ , of its neutrons in the spin “up” state and a fraction, ϕ_- , in the “down” state has a magnitude of its polarization, \mathbf{P} , in the field direction given by

$$|\mathbf{P}| = \phi_+ - \phi_- = 2\phi_+ - 1 = 1 - 2\phi_- \quad (39)$$

The total optical potential in a medium containing a magnetic field is the sum of the nuclear and magnetic potentials. Consequently, the index of refraction may be written

$$n^2 = 1 - \lambda^2 [\rho b / \pi \pm m \mu_n B / (2\pi^2 \hbar^2)] \quad (40)$$

This is often written as $n^2 = 1 - \lambda^2 \rho(b \pm p) / \pi$, where p is considered a magnetic scattering length, the sign of which depends on the neutron spin orientation relative to magnetic field, \mathbf{B} , in the medium.

This polarization property of a neutron beam gives many special capabilities to neutron scattering and neutron optics.⁵² The refractive index depends on the orientation of the spin with respect to the magnetization direction, and therefore there are two critical glancing angles of reflection, with polarizing mirrors having a greater critical angle for the + state than for the – state. This property of reflection of neutrons from magnetized mirrors is widely used for obtaining polarized neutrons. From Eq. (40), if $\rho b < m \mu_n B / (2\pi \hbar^2)$, neutrons with spins antiparallel to the magnetization are not reflected, and only those neutrons with parallel spin orientation are reflected. However, values of p are less than b , so that devices work with an angular range such that reflection occurs for + spins and transmits for the – spins. Supermirror polarizers⁵³ can be used in transmission or reflection geometry. In transmission geometry the required (–) spin state beam is transmitted through the supermirror and its substrate and the unwanted (+) spin state is reflected from the beam and absorbed elsewhere. In reflection geometry the required (+) spin state beam is reflected, and the unwanted (–) spin state is transmitted through the supermirror and absorbed in the substrate.

Polarizing mirrors are generally ferromagnetic films deposited on silicon substrates and $m = 3$ Fe/Si polarizing mirrors are common. Polarized beams are obtained using total reflection from magnetized iron and cobalt mirrors, with the refractive index of the ferromagnetic material in one spin state matched to that of the nonmagnetic material. Polarizing mirrors produce high polarization and good transmission, but are restricted to low incident angles and wavelengths longer than 2 Å. In addition, thin (≈ 100 nm) evaporated films of magnetic alloys have been used. Two-dimensional multilayer structures have alternate layers of a (magnetized) magnetic material whose refractive index for the (–) spin neutrons is matched to the second (nonmagnetic) layer, for example, in Fe-Ge multilayers. Supermirrors can produce a broadband-polarized beam using evaporated multilayer films deposited to provide a lattice spacing gradient. The reflectivity profile may be extended in angle beyond the simple mirror critical glancing angle through thin-film interference.

A ferromagnetic single crystal may simultaneously polarize and monochromate a neutron beam. A magnetic field applied perpendicular to the scattering vector saturates the magnetic moment along the field direction. Diffraction occurs corresponding to plane spacings with an intensity that is dependent on both the nuclear and the magnetic structure factors. The sign of the latter depends on the direction of the neutron spin relative to the magnetic field. Ideally, if the nuclear and magnetic structure factors are equal in magnitude, so that for one spin state the intensity of diffraction is zero and the other spin state has twice the structure for diffraction. The naturally occurring ferromagnets are the 3d elements Fe, Ni, and Co, but none are good single crystal polarizers. The nuclear scattering lengths for iron and nickel are too large for their weak magnetic moments, while cobalt is a strong neutron-absorbing element. Alloys are found to have a better match of the structure factors, and $\text{Co}_{0.92}\text{Fe}_{0.08}$ (220 reflection with $d = 1.6$ Å) and magnetite Fe_3O_4 are used as polarizing monochromators, but Heusler alloy Cu_2MnAl (111 reflection with $d = 3.43$ Å) is often

preferred, with a reflectivity of 25 to 30 percent for the reflected spin state and a neutron polarization of 95 percent.

There are various polarizing filters that absorb one spin state while transmitting the other spin state. The filter thickness is a compromise between the desired polarization P that increases and the transmission T that decreases with filter thickness. High polarization ($P > 0.96$) with reasonable transmission ($T > 0.2$) requires that the spin-independent cross section is small compared to the spin dependent absorption cross section. Resonance absorption polarization filters depend on the spin dependent absorption cross section of polarized nuclei at their nuclear resonant energy. Cooling ^{149}Sm or ^{151}Eu in a magnetic field can provide efficient polarizing filters with reasonable nuclear polarization over a wide energy range. A better filter is obtained with nuclear-spin polarized ^3He that takes advantage of the huge absorption cross section of ^3He that comes from only the $(-)$ spin state, so that it preferably absorbs neutrons with spins antiparallel to the ^3He spin, while transmitting the $(+)$ spin state.

Neutron capture by a ^3He nucleus results in an excited ^4He nuclear state that subsequently decays into a proton and a triton. The cross section for capture of thermal neutrons with spin antiparallel to the ^3He nuclear spin is 10666 b , whereas that for neutrons with spin parallel is essentially zero. Theoretically with a sufficient gas density of 100 percent polarized ^3He , all neutrons with antiparallel spin are absorbed, while nearly all neutrons with parallel spin are transmitted, resulting in 100 percent neutron polarization and 50 percent transmission. For an initially unpolarized neutron beam that passes through a cell filled with ^3He with less than 100 percent polarization, the transmitted beam has a polarization P_n given by

$$P_n(\lambda) = (n_+ - n_-)/(n_+ + n_-) = \tanh[\sigma(\lambda)\rho_{\text{He}}LP_{\text{He}}] \quad (41)$$

and a transmission T_n given by

$$T_n(\lambda) = T_0 \cosh[\sigma(\lambda)\rho_{\text{He}}LP_{\text{He}}] \quad (42)$$

where n_+ and n_- are the numbers of neutrons with parallel and antiparallel spin, $\sigma(\lambda)$ is the absorption cross section for unpolarized neutrons of wavelength λ , ρ_{He} is the number density of ^3He , L is the length of the cell, and P_{He} is the nuclear polarization of the ^3He . For $P_{\text{He}} = 0$, T_0 is the transmission through an evacuated cell, neglecting the transmission loss of the glass windows. ^3He can be polarized either by spin exchange⁵⁴ with optically pumped Rb or metastability-exchange optical pumping methods,⁵⁵ and ^3He nuclear polarizations near 75 percent have been attained. The characteristics of the filter depend on its opacity, which is defined by the product of the cell length, the ^3He gas pressure, and the neutron wavelength (proportional to the absorption cross section). Figure 8 shows the neutron polarization P_n and the transmission T_n for a ^3He nuclear polarization $P_{\text{He}} = 0.75$ as a function of the cell opacity. This indicates that there is a trade-off between neutron polarization and transmission, and a figure-of-merit is usually given by $P_n^2 T_n$. Note that shorter wavelength neutrons require a thicker target to produce the same polarization.

A spin flipper is a device that can reverse the direction of polarization of a beam. A thin flat coil with a relatively strong field inside is placed perpendicular to the neutron beam. Within the coil the beam polarization direction is no longer parallel to the magnetic field and it begins to precess around the field. The particular field strength of the flipper and its direction determine whether the neutrons precess exactly by π or $\pi/2$ by the time the beam has passed through the flipper.

Larmor Precession

When a polarized neutron beam enters a homogeneous magnetic field B that is perpendicular to the neutron magnetic moment, the direction of the neutron spin precesses in B with a Larmor frequency given by $\omega_L = 2|\mu_n B|/\hbar = \gamma B$, where γ is the gyromagnetic ratio of the neutron. This phenomenon is utilized by the spin echo spectrometer,⁵⁶ in which the incoming and outgoing velocities of the neutron

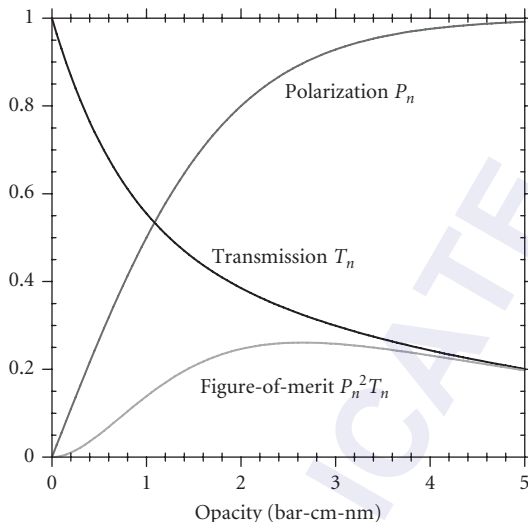


FIGURE 8 Neutron polarization P_n , transmission T_n , and figure-of-merit $P_n^2 T_n$ as a function of opacity, the product of the ^3He pressure in bar, assuming a room temperature cell, the cell length L in cm, and the wavelength λ in nm, for a nuclear polarization $P_{\text{He}} = 0.75$ of the ^3He .

are measured, so that small changes in velocity that may take place on scattering can be inferred from the shift in phase in the Larmor precession angle. This is effectively a time-of-flight method in which polarized neutrons precess in the magnetic field to provide its own clock. Such a technique enables energy resolution in the neV region, orders of magnitude smaller than for conventional instruments.

Cold neutrons from a neutron guide are coarsely monochromatized by a velocity selector and polarized by supermirrors. The spins are made horizontal by a $\pi/2$ spin flipper. The polarized neutron beam enters the first magnetic field B_1 of length L_1 and the spins precess many times. The total precession angle of a neutron spin is $\phi_1 = \gamma B_1 L_1 / v_1$, where v_1 is the neutron velocity. After scattering the spins are reversed in direction by a π flipper, and then precess in the reverse direction. The total precession angle of a neutron spin in the second magnetic field B_2 of length L_2 is $\phi_2 = \gamma B_2 L_2 / v_2$, where v_2 is the final neutron velocity. The spins are flipped again by $\pi/2$, and the polarization of the neutron beam is analyzed. If the scattering is elastic, $v_1 = v_2$, and if $B_1 L_1 = B_2 L_2$, the final beam has the same polarization. If there is a change in energy $\hbar\omega$ upon scattering where ω is small, the total change in the precession angle is given by $\phi_{\text{tot}} = \gamma BL / v^3 (\hbar\omega/m)$.

The analogue to the Stern-Gerlach effect arises from the interaction between the neutron magnetic moment μ_n and an inhomogeneous magnetic field that bends the neutron trajectory by imparting an acceleration

$$d^2 \mathbf{r} / dt^2 = \mp |\mu/m| |\nabla| \mathbf{B} | \quad (43)$$

A superconducting hexapole magnet of typical length ≈ 2 m can simultaneously focus and polarize a neutron beam. If the incident neutron beam is unpolarized, the trajectory of those neutrons with spins oriented parallel to the direction of the magnetic field bend toward the magnet axis, and those antiparallel bend away from the axis. As a result half the neutrons are spin polarized about the local magnetic field and are focused a distance that is inversely proportional to λ . Because the effect index of refraction of the magnetic field lens is near unity, the focusing is weak and depends on the strength of the hexapole magnetic field.

63.8 NEUTRON DETECTION

Neutrons are uncharged and do not directly ionize atoms, and consequently they easily penetrate most materials. They cannot be detected directly, and their detection requires a nuclear reaction within a converter that results in the emission of a charged particle or ionizing radiation that can more readily interact with other materials.^{13,57} (Fast neutrons may be detected by proton recoil or after thermalization.) The efficiency of absorption of neutrons incident normally on an absorber of thickness ℓ is given by $1 - \exp[-\rho\sigma_a\ell]$, where ρ is the atom number density of the absorbing nuclei and σ_a the absorption cross section. The three most common reactions used for the detection of slow neutrons all have cross sections that increase linearly with wavelength (see Table 1). All three reactions are exothermic with Q values much greater than the energy of slow/thermal neutrons (≈ 25 meV) so that there is no energy discrimination. (Note that energy or wavelength determination of a neutron beam requires a crystal spectrometer or a mechanical monochromator.) The choice of the reaction depends on its ability to discriminate against gammas. The reaction Q value determines the energy released, and the greater the energy, the easier it is to discriminate against gammas.

1. The $^{10}\text{B}(n,\alpha)$ reaction has a thermal (velocity of 2200 m s⁻¹, or wavelength of $\lambda \approx 0.179$ nm) absorption cross section of 3835 b, and a Q value of 2.792 MeV to decay directly the ground state (6 % branch ratio), or a Q of 2.310 MeV for decay to an excited state with the subsequent emission of a 478-keV gamma (94%) with a half life of about 10^{-3} s. The reaction products are an alpha particle with an energy of either 1.77 MeV (6%) or 1.47 MeV (94%), and a ^7Li ion with an energy 1.01 MeV (6%) or 0.84 MeV (94%). This reaction is used in BF_3 gas proportional counters or in B-loaded scintillators. The natural isotopic abundance of ^{10}B is 19.8 percent.
2. The $^6\text{Li}(n,\alpha)$ reaction with a thermal absorption cross section of 940 b and a Q value of 4.78 MeV emits no gamma and is commonly used in neutron scintillator detectors. The reaction products are an alpha particle with an energy of 2.05 MeV and a ^3H ion with an energy of 2.73 MeV. While the absorption cross section is lower for this reaction than for ^{10}B , the Q value of the reaction is higher and therefore more ionized charge is produced. The natural isotopic abundance of ^6Li is 7.4 percent.
3. The $^3\text{He}(n,p)$ reaction, with an unpolarized thermal absorption cross section of 5333 b and a Q value of 0.764 MeV, is used in ^3He gas-proportional counters. The reaction products are a proton with an energy of 0.573 MeV and a ^3H ion (triton) with an energy of 0.191 MeV. Though the absorption cross section is much higher than the other two isotopes, it is relatively expensive. The lower Q value can be compensated by high electron multiplication in ^3He proportional counters.

In all three cases, the reactions products are emitted in opposite directions randomly oriented in space. These reactions are also used in the one- and two-dimensional detectors, and often form the basis of other detection methods, such as imaging detectors and neutron imaging plates. (They are also key reactions used in neutron depth profiling.) In general cadmium and gadolinium, or other rare earth elements, that have large thermal absorption neutron cross sections are not considered good detector materials, because they result in the copious production of gammas.

Neutron detectors may be operated in either current mode or in pulse mode. In the former, the detector events are integrated either over a short time to provide a beam monitor, or over a longer time as in a charge coupled device camera for beam alignment, or in a photographic plate or an imaging plate detector used for neutron radiography or sometimes in neutron diffraction. These detectors can operate with high intensity and can have high spatial resolution. Generally they have a large dynamic range, but may suffer from low efficiency, gamma sensitivity, and high noise. In the pulse mode individual neutrons are detected in a gas counter or a scintillator detector, and sometimes a converter foil; these are used in most neutron spectroscopy measurements. These detectors generally have high efficiency, low noise, and can discriminate against gammas, though they may be limited in their count rate and have limited spatial resolution. Discrete detector and multiwire proportional

counters (MWPCs) with coincidence encoding require little signal processing and can operate at relatively high counting rates. Signals from each detection element can be processed individually, allowing high data rates. Resolution can be improved if the charge (or scintillation light) from each neutron detection event is distributed over several detection elements.

Gas Detectors

The gas detector may be operated in either a current or a pulse mode. When the charged particles from the reaction pass through a gas, the primary mode of interaction is ionization and the excitation of gas molecules along the particle track. The average energy lost by the incident particle per ion pair formed depends on the particular gas, the type of radiation and its energy. Typical values are about (30 to 35) eV per ion pair. For instance a 1 MeV particle that is fully stopped within the gas produces slightly greater than 3×10^4 ion pairs. The operation of the ionization chamber is based on the collection of all charge created by direct ionization within the gas by the application of an electric field. Two metallic electrodes surround the fill gas of ^3He or $^{10}\text{BF}_3$. The two emitted particles slow down in the gas and the charges of the electron-ion pairs produced along their tracks are collected on the plates using an electric field strength of ≈ 100 V/cm. The ^3He reaction produces 0.764 MeV kinetic energy in the charged particles and results in about 2.5×10^4 electron-ion pairs. The ionization chamber is usually used in current mode and serves as a low efficiency beam monitors on most neutron beam lines.

By increasing the electric field within the gas, the original ion pairs created within the gas are accelerated toward the electrodes, leading to additional ionization of the gas molecules by the accelerated electrons. This secondary ionization leads to an avalanche of charge, and the density of the electron cloud grows exponentially with distance as the avalanche progresses. It terminates when the central anode wire has collected all the free electrons produced. This is the ubiquitous proportional counter, with the total charge reaching the anode proportional to the number of primary electrons; and the gas multiplication or gain can be as high as $\approx 10^5$. Individual neutron events are easily detected. The charge collection time is about 10 μs , so that the proportional counter is count rate limited to about 3×10^4 s^{-1} . The fill gas must not have appreciable electron attachment, and must quench any remaining charge to prevent the migration of the heavier ions. This is generally performed by inert gases and the best is P10 gas (90% argon and 10% CH_4 by volume). Partial pressure in the ^3He counter is adjusted for the required neutron absorption efficiency, and a pressure of (10 to 20) bar can be used. When the reaction takes place near the edge of the detector housing, some of the charge is lost by the “wall effect” such that the full charge avalanche does not take place. However, a heavy gas with a high stopping power to reduce the charged particle range may ameliorate this. The differential pulse height spectrum is such that a lower level discriminator may be used to eliminate the gamma ray background.

The gas proportional counter is noiseless and has high charge amplification with good gamma discrimination and can stand high radiation fields. They are usually cylindrical, though detectors with elliptical cross sections are also available. The neutron sensitive fill gas is usually ^3He . It may be made position-sensitive by having an amplifier at both ends of the anode wire. In the charge division method, the collected charge is divided in proportional to the resistance to ground at each end along the wire. The positional sensitivity depends on the highly resistive anode wire. The other method is the observation of the rise time from preamplifiers at either end of the resistive anode wire. Typical spatial resolution may be as high as 1.5 mm. The multiwire gas proportional counter MWPC is a two-dimensional detector with a large number of parallel anode wires in a plane mounted between two cathodes composed of strips in orthogonal directions. The multiple electrodes are enclosed within a large single gas chamber. As before the gas multiplication of electrons is directed toward the anode wire, and positive charge is induced on the cathodes. These enable the positional determination in the two directions using coincidence encoding. The positional resolution is greater than 1 mm. The charge collection time is about 1 μs , so that the two-dimensional proportional counter is limited to about 3×10^5 events s^{-1} . At high detection rates, the space charge reduces the electric field around the wire and hence the charge gain. The microstrip detector⁵⁸ uses a photolithography technique, and

replaces the usual multiwire anode plane to reduce space-charge effects and gives higher counting rates and better spatial resolution.

Neutron Scintillators

All gas detectors require some form of gamma discrimination, and only glass scintillators lead to tolerable separation between neutron and gamma signals. The energy from the charged particle raises electrons within the material to excited states that decay back to the ground state by the emission of photons. The scintillator must be transparent to its own emitted radiation. Often the scintillating material must be mixed with a compound that contains the absorbing nuclei, such as a lithium compound dispersed in a matrix of ZnS(Ag). A ZnS(Cu, Ag, Au) scintillator is superior because it gives green light that is better for CCD (charge-coupled device) detection. ${}^6\text{Li}$ -based scintillators are favored over ${}^{10}\text{B}$ -based scintillators, because the greater energy released by the ${}^6\text{Li}$ reaction results in greater light production. Fused $\text{B}_2\text{O}_3/\text{ZnS}$ has less effective gamma discrimination. The crystals often contain small amounts of an activating impurity to provide the necessary electronic states to produce suitable scintillating properties. The light produced in the scintillator may be directly coupled to a light detection device such as a photomultiplier tube, though this may also be a channel-plate amplifier, a CCD camera, or photographic film. Alternatively, the scintillator can be coupled through light guide or optical fibers, or through a system of lenses and mirrors. A lithium-loaded ZnS scintillator with an image intensifier is often used to monitor the profile of the neutron beam with the output of observed on a screen.

The range of the ionizing particles is only a few microns within the scintillator material, enabling much greater spatial resolution. A thin detector limits the gamma sensitivity and the same time limits the detection efficiency. Typical scintillators include ZnS/ ${}^6\text{LiF}$ activated with Ag, ${}^6\text{LiI}$ crystal activated with Eu, and ${}^6\text{Li}$ glass activated with Ce, with the first mostly used, having $\approx 1.6 \times 10^5$ photons produced per absorbed neutron. Only a fraction of these photons are collected in the photomultiplier or other light-sensitive device, so that the quantum efficiency of a typical photocathode may be ≈ 20 percent. Photomultiplier tubes also have to suffer from dark current, which can be decreased by cooling the tube. The spread in pulse amplitudes from a scintillator detector is typically 10 to 20 percent and is the limitation on the spatial resolution of scintillator position-sensitive detectors (PSDs).

Scintillating glass fibers that are loaded with both ${}^6\text{Li}$ and Ce can be operated with a spatial resolution of $100\ \mu\text{m}$. The triton emitted from the reaction can excite a Ce^{3+} ion, and the subsequent deexcitation of the electron to the ground state results in the emission of visible light that is transmitted through the fiber and collected by the photomultiplier tube. These fibers can be arranged to form a large area position-sensitive detector. The principal disadvantage of these glass fibers is the sensitivity to gammas.

Photographic film may have a large dynamic range with high spatial resolution of $\approx 100\ \mu\text{m}$, but as an integrating device is sensitive to gammas. The Anger camera consists of a scintillator that is coupled by optical fibers to a number of photomultipliers. The positional encoding of the neutron event is determined by centroiding the output signal, with a precision much smaller than the size of the individual detector element, with a resolution of several mm.

Solid State Detectors

The deposition of the particle energy in a semiconductor results in the formation of electron-hole pairs. The energy required to form such pairs is $\approx 3\ \text{eV}$. FWHM pulse-height resolutions are generally better than 0.5 percent. Semiconductor diodes have a depletion layer to form the active region of the detector. The depletion layer thickness is adjusted by the bias voltage on the diode, with thicknesses up to several mm in high purity Si or Ge. The electrons and holes are collected with no amplification, analogous to the gas ionization chamber, and B or Li are common additions in semiconductor devices as the neutron absorber. The ranges of the primary charged particles is short, so that the position resolution should be good. The detectors are also sensitive to gammas.

The cross sections for neutron induced fission reactions in ^{233}U , ^{235}U , and ^{239}Pu are relatively large at low neutron energies. The reaction Q values are extremely large, ≈ 200 MeV, with approximately 160 MeV being released as kinetic energy of the fission fragments. These produce a large charge in the ionization chamber.

When the neutron-sensitive material is a solid film, the density of the material is so high that the path lengths in the film are very short, and the probability of the charged particle escaping and reaching the detection medium is reduced. The maximum detection efficiency is obtained with an optimum foil thickness that depends on the range of the charged particle in the foil, and this depends on whether the foil detection is in transmission or backscattering geometry. A problem with foils is that only one charged particle is detected and even this loses varying amounts of its energy before entering the charged particle detector. Hence the pulse height spectrum is broad, and gamma discrimination is difficult.

Gadolinium foil can also be used in conjunction with a microchannel plate (MCP) detector where the MCP acts as an electron multiplier. Greater neutron detection efficiency is obtained when the neutron absorbing material (such as ^6Li or ^{10}B) is imbedded within the glass of the MCP. The reaction products create secondary electrons that are attracted by a positive bias and cause an electron avalanche that travels down the channel. This results in thousands of electrons emerging from the MCP due to a single neutron. The position of the electronic pulses may be determined directed using a two-dimensional wire grid or indirectly on a phosphor screen to produce a proportional light image. Typical MCPs consists of ^{10}B -doped glass with thousands of $5\ \mu\text{m}$ diameter individual channels. Spatial resolution approaching $10\ \mu\text{m}$ is possible with good detection efficiency. Sometimes these MCPs are arranged in a chevron configuration to increase the electron gain. Also microsphere plate detectors operate on principles similar to MCPs, and spatial resolutions better than $250\ \mu\text{m}$ have been achieved. These devices are essentially insensitive to gammas.

Imaging Plates

Photographic film can be used to detect the light produced in a neutron-sensitive scintillator screen or in conjunction with a neutron-sensitive foil such as Gd to detect the charge reaction directly. The path length of the charge particles is quite short in the high-density photographic emulsion, enabling good spatial resolution with photographic film. The best spatial resolution of $\approx 20\ \mu\text{m}$ is achieved with good quality film in contact with a Gd foil enriched with ^{157}Gd . A resolution of 0.5 mm is obtained with a standard neutron polaroid camera using a scintillator screen of $\text{ZnS}(\text{Ag})^6\text{LiF}$. The optical densities of the photographic films can be digitized with optical scanners. These integrating detectors have no data rate limitations, but can provide no time information and therefore cannot be used in time-of-flight applications.

The concept of storage photostimulable phosphors or imaging plates involves radiation trapping in a phosphor-coated detector plate and the release of the trapped energy as light when the plate is stimulated under scanning in a reader with a fine laser beam of typical spot size of $50\ \mu\text{m}$. A photomultiplier tube that detects the light is coordinated with the scanned positions to create a digital image of the radiation field. The light intensity is linear with the received radiation level over a wide dynamic range ($> 10^5$). The information on the detector plate can be erased after scanning so that it can be reused. It is far more sensitive than film autoradiography. The detectors have a spatial resolution ranging from 25 to $200\ \mu\text{m}$ with a single detector plate of size up to 350×420 mm.

Reusable photostimulable phosphors suitable to neutron detection incorporate converter material such as Gd_2O_3 or LiF within the phosphor plate itself. Yet the inherent sensitivity of the imaging plate to gamma radiation produce high background interference. In addition, residual radioactivity from neutron activated materials, particularly from europium, and the depletion of neutron-sensitive materials in the imaging plate after repeated usage may require frequent calibration. An alternative is a two-step transfer process within which a foil is placed in the neutron field and then is transfer later to the imaging plate to record the radioactivity distribution in the converter and reproduce the neutron field information. Dysprosium is the most efficient neutron converter, decays with a half life of 2.35 h almost entirely by beta emission, and the activity is linear for over 5 orders of magnitude.

63.9 REFERENCES

1. V. F. Sears, *Neutron Optics: An Introduction to the Theory of Neutron Optical Phenomena and Their Applications*, Oxford University Press, Oxford, 1989.
2. M. Utsuro (ed.), "Advance in Neutrons Optics and Related Research Facilities," *J. Phys. Soc. Japan* **65** (suppl. A): (1996).
3. W. Marshall and S. W. Lovesey, *Theory of Thermal Neutron Scattering: The Use of Neutrons for the Investigation of Condensed Matter*, Clarendon Press, Oxford, 1971.
4. G. L. Squires, *Thermal Neutron Scattering*, Cambridge University Press, Cambridge, 1978.
5. S. W. Lovesey, *Theory of Neutron Scattering from Condensed Matter*, Oxford University Press, Oxford, 1984.
6. M. L. Goldberger and F. Seitz, "Theory of the Refraction and the Diffraction of Neutrons by Crystals," *Phys. Rev.* **71**:294–310 (1947).
7. V. F. Sears, "Neutron Scattering Lengths and Cross Sections," in *Methods of Experimental Physics: Neutron Scattering*, K. Sköld and D. L. Price (eds.), vol. 23A, Academic Press, San Diego, CA, 1986, pp. 521–550.
8. V. F. Sears, "Neutron Scattering Lengths and Cross Sections," *Neutron News* **3**(3):26–37 (1992).
9. L. Koester, H. Rauch, and E. Seymann, "Neutron Scattering Lengths: A Survey of Experimental Data and Methods," *Atomic Data and Nuclear Data Tables* **49**:65–120 (1991).
10. L. Koester and H. Ungerer, "Coherent Scattering Amplitudes Measured by Small Angle Scattering of Neutrons," *Z. Phys.* **219**:300–310 (1969).
11. H. B. Stuhmann, "Molecular Biology," in *Methods of Experimental Physics: Neutron Scattering*, K. Sköld and D. L. Price (eds.), vol. 23C, Academic Press, San Diego, CA, 1986, pp. 367–403.
12. *Neutron Radiography*, Proceedings of the 8th World Conference on Neutron Radiography, Gaithersburg, MD, M. Arif and R.G. Downing, (eds.), DEStech Publications Inc., Lancaster, PA, 2008.
13. K. M. Beckurts and K. Wirtz, *Neutron Physics*, Springer-Verlag, Berlin, 1964.
14. J. M. Carpenter and W. B. Yelon, "Neutron Sources," in *Methods of Experimental Physics: Neutron Scattering*, K. Sköld and D. L. Price (eds.), vol. 23A, Academic Press, San Diego, CA, 1986, pp. 99–196.
15. D. H. Kim, C. S. Gil, J.-D. Kim, and J. Chang, "Generation and Validation of a Shielding Library Based on ENDF/B-VI.8," *Radiation Protection Dosimetry* **115**:232–237 (2005).
16. C. G. Windsor, "Experimental Techniques," in *Methods of Experimental Physics: Neutron Scattering*, K. Sköld and D. L. Price (eds.), vol. 23A, Academic Press, San Diego, CA, 1986, pp. 197–257.
17. H. Maier-Leibnitz and T. Springer, "The Use of Neutron Optical Devices on Beam-Hole Experiments," *React. Sci. Tech. (J. Nucl. Energy A/B)* **17**:217–225 (1963).
18. P. Böni, "New Concepts for Neutron Instrumentation," *Nucl. Instrum. Meth. A* **586**:1–8 (2008).
19. V. F. Sears, "Theory of Multilayer Neutron Monochromators," *Acta Cryst.* **A39**:601–608 (1983).
20. I. Anderson, "From 1 to m: The Development of Supermirrors," *SPIE Proc.* **1738**:118–124, *Neutron Optical Devices and Applications*, C. F. Majkrzak and J. L. Wood (eds.) (1992).
21. J. R. D. Copley and C. F. Majkrzak, "Calculations and Measurement of the Performance of Converging Neutron Guides," *SPIE Proc.* **983**:93–104, *Thin-Films, Neutron Optical Devices: Mirrors, Supermirrors, Multilayer Monochromators, Polarizers, and Beam Guides* (C. Majkrzak, ed.) (1988).
22. C. Schanzer, P. Böni, U. Filges, and T. Hils, "Advanced Geometries for Ballistic Neutron Guides," *Nucl. Instrum. Meth. A* **529**:63–68 (2001).
23. A. K. Freund, "Cross-sections of Materials used as Neutron Monochromators and Filters," *Nucl. Instrum. Meth.* **213**:495–501 (1983).
24. J. B. Hayter, "Neutron Optics at the Advanced Neutron Source," *SPIE Proc.* **1738**:2–7, *Neutron Optical Devices and Applications*, C. F. Majkrzak and J. L. Wood, (eds.) (1992).
25. A. Steyerl, "Neutron Physics," in *Springer Tracts in Modern Physics*, vol. 80, pp. 57–130, Springer-Verlag, Berlin, 1977.
26. L. Koester, "Absolutmessung der Kohärenten Streulänge von Quecksilber mit den Neutronen-Schwerkraft Refraktometer am FRM," *Z. Phys.* **182**:326–328 (1965).
27. G. E. Ice, C. R. Hubbard, B. C. Larson, J. W. L. Pang, J. D. Budai, S. Spooner, and S. Vogel, "Kirkpatrick-Baez Microfocusing Optics for Thermal Neutrons," *Nucl. Instrum. Meth. A* **539**:312–320 (2005).

28. Q. F. Xiao, H. Chen, V. A. Sharov, D. F. R. Mildner, R. G. Downing, N. Gao, and D. M. Gibson, "Neutron Focusing Optic for Submillimeter Materials Analysis," *Rev. Sci. Instrum.* **65**:3399–3402 (1994).
29. C. S. Schneider and C. G. Schull, "Forward Magnetic Scattering Amplitude of Iron for Thermal Neutrons," *Phys. Rev.* **B3**:830–835 (1971).
30. R. Gähler, J. Kalus, and W. Mampe, "An Optical Instrument for the Search of a Neutron Charge," *J. Phys.* **E13**:546–548 (1980).
31. M. R. Eskildsen, P. L. Gammel, E. D. Isaacs, C. Detlefs, K. Mortensen, and D. J. Bishop, "Compound Refractive Optics for the Imaging and Focusing of Low-Energy Neutrons," *Nature* **391**:563–566 (1998).
32. V.F. Sears, "Bragg Reflection in Mosaic Crystals I. General Solution of the Darwin Equations," *Acta Cryst.* **A53**:35–45, and "Bragg Reflection in Mosaic Crystals II. Neutron Monochromator Properties," *Acta Cryst.* **A53**:46–54 (1997).
33. H.-C. Hu, "Universal Treatment of X-ray and Neutron Diffraction in Crystals I. Theory," *Acta Cryst.* **A53**:484–492, and "Universal Treatment of X-ray and Neutron Diffraction in Crystals II. Extinction," *Acta Cryst.* **A53**:493–504 (1997).
34. C. Broholm, "Proposal for a Doubly Focusing Cold Neutron Spectrometer at NIST," *Nucl. Instrum. Meth. A* **369**:169–179 (1996).
35. C. G. Schull, "Single-Slit Diffraction of Neutrons," *Phys. Rev.* **179**:752–754 (1969).
36. A. Zeilinger, R. Gähler, C. G. Schull, and W. Treimer, "Experimental Status and Recent Results of Neutron Interference Optics," *AIP Conf. Proc.* **89**:93 (1981).
37. H. Kurz and H. Rauch, "Diffraction of Thermal Neutrons by a Ruled Grating," *Z. Phys.* **220**:419–426 (1969).
38. R. Gähler, A. G. Klein, and A. Zeilinger, "Neutron Optical Tests of Nonlinear Wave Mechanics," *Phys. Rev.* **A23**:1611–1617 (1981).
39. P. D. Kearney, A. G. Klein, G. I. Opat, and R. A. Gähler, "Imaging and Focusing of Neutrons by a Zone Plate," *Nature* **287**:313–314 (1980).
40. A. G. Klein, P. D. Kearney, G. I. Opat, and R. A. Gähler, "Focusing of Slow Neutrons with Cylindrical Zone Plates," *Phys. Lett.* **83A**:71–73 (1981).
41. S. A. Werner and A. G. Klein, "Neutron Optics," in *Methods of Experimental Physics: Neutron Scattering*, K. Sköld and D. L. Price (eds.), vol. 23A, Academic Press, San Diego, CA, 1986, pp. 259–337.
42. H. Rauch and S. A. Werner, *Neutron Interferometry: Lessons in Experimental Quantum Mechanics*, Oxford University Press, Oxford, 2000.
43. C. G. Schull, "Observation of Pendellösung Fringe Structure in Neutron Diffraction," *Phys. Rev. Lett.* **21**:1585–1589 (1968).
44. A. G. Klein and G. I. Opat, "Observation of 2π Rotations by Fresnel Diffraction of Neutrons," *Phys. Rev. Lett.* **37**:238–240 (1976).
45. R. R. Highfield, R. K. Thomas, P. G. Cummins, D. P. Gregory, J. Mingins, J. B. Hayter, and O. Schärpf, "Critical Reflection of Neutrons from Langmuir-Blodgett Films on Glass," *Thin Solids Films* **99**:165–172 (1983).
46. F. Mezei, "Novel Polarized Neutron Devices: Supermirror and Spin Component Amplifier," *Commun. Phys.* **1**:81 (1976); F. Mezei and P. A. Dagleish, "Corrigendum and First Experimental Evidence on Neutron Supermirrors," *Commun. Phys.* **4**:41 (1977).
47. K.-A. Steinhauser, A. Steyerl, H. Scheckhofer, and S. S. Malik, "Observation of Quasibound States of the Neutron in Matter," *Phys. Rev. Lett.* **44**:1306–1309 (1981).
48. U. Bonse and M. Hart, "An X-Ray Interferometer," *Appl. Phys. Lett.* **6**:155–156 (1965).
49. H. Rauch, W. Treimer, and U. Bonse, "Test of a Single Crystal Neutron Interferometer," *Phys. Lett.* **47A**:369–371 (1974).
50. A. Zeilinger, C. G. Schull, M. A. Horne, and G. Squires, "Two-Crystal Neutron Interferometry," in *Neutron Interferometry*, U. Bonse and H. Rauch (eds.), Oxford University Press, London, 1979, pp. 48–59.
51. H. Rauch, E. Seidl, A. Zeilinger, W. Bausspiess, and U. Bonse, "Hydrogen Detection in Metals by Neutron Interferometry," *J. Appl. Phys.* **49**:2731–2734 (1978).
52. W. G. Williams, *Polarized Neutrons*, Oxford University Press, Oxford, 1988.
53. P. Böni, D. Clemens, M. S. Kumar, and C. Pappas, "Applications of Remanent Supermirror Polarizers," *Physica B* **267–268**:320–327 (2007).

54. W. C. Chen, G. Armstrong, Y. Chen, B. Collett, R. Erwin, T. R. Gentile, G. L. Jones, J. W. Lynn, S. McKenney, and J. E. Steinberg, "³He Spin Filters for a Thermal Neutron Triple Axis Spectrometer," *Physica B* **397**:168–171 (2007).
55. E. Lelièvre-Berna, "Mid-Term Report on the NM13 Neutron Spin Filter Project," *Physica B* **397**:162–167 (2007).
56. F. Mezei, *Neutron Spin Echo*, Springer-Verlag, Berlin, 1980.
57. G. F. Knoll, *Radiation Detection and Measurement*, 2nd ed., chap. 14, John Wiley & Sons, New York, 1989.
58. A. Oed, "Micro Pattern Structures for Gas Detectors," *Nucl. Instrum. Meth. A* **471**:109–114 (2001).

DO NOT DUPLICATE

This page intentionally left blank.

DO NOT DUPLICATE

GRAZING-INCIDENCE NEUTRON OPTICS

Mikhail Gubarev and Brian Ramsey

*NASA/Marshall Space Flight Center
Huntsville, Alabama*

64.1 INTRODUCTION

Due to their unique sensitivity to light elements, neutrons are an important tool for progress in material and life sciences, medical research, and therapy as well as for fundamental physics. However, the availability of useful neutron beams is limited as neutron sources typically have relatively low brilliance (compared, e.g., to x-ray sources.) Focusing and concentrating of neutrons, through the use of suitable optics, offers a way to significantly increase neutron current density at the samples under investigations. Such a capability becomes particularly important for experiments with very small samples. In addition, neutron focusing can also improve the signal to noise ratio for such applications as small angle neutron scattering (SANS) analysis and neutron crystallography, when the neutron beam can be focused onto the detector.

This chapter is concerned with focusing neutron optics, that is, the optical elements capable of producing an image of the source through a limited but determined number of reflections. Other optical elements, such as capillary optics and neutron guides, utilize a stochastic number of reflections to transport or concentrate neutrons. In these, neutrons undergo multiple reflections from guide or capillary walls to emerge in a new direction. These elements are described in Chaps. 53 and 63.

64.2 TOTAL EXTERNAL REFLECTION

The optical properties of materials with respect to neutrons are characterized by their refractive indices¹ which are a function of the neutron wavelength

$$n = 1 - \delta + i\beta \quad (1)$$

with

$$\delta = \frac{\lambda^2 N_d b}{2\pi}$$

and

$$\beta = \frac{\lambda N_a \sigma_a}{4\pi}$$

where N_a is the atomic number density, b is the coherent scattering length, and σ_a is the absorption cross-section. For most materials (with the exception of those containing Li, B, Cd, Sm, or Gd) the absorption cross-section is almost zero and Eq. (1) can be reduced to

$$n = 1 - \delta$$

The parameter δ is strongly dependent on neutron wavelength and in the case of thermal neutrons, is about 10^{-5} for most elements and their isotopes. So, as is the case for x rays, the refractive index for neutrons is slightly less than unity. Therefore, thermal and cold neutrons can be reflected from smooth surfaces at shallow “grazing-incidence” angles (total external reflection) or be refracted at the boundaries of different materials. Both reflective and refractive optics can be used to focus and concentrate neutrons. The advantage of neutron reflective optics, reviewed here, is that they are achromatic compared to refractive optics.

Total external reflection occurs when neutrons are incident on a surface at angles below the critical angle ϑ_{cr} , given by

$$\vartheta_{cr} = \sqrt{2\delta} \quad (2)$$

so that at 10 Å, for example, the critical angle for reflection off nickel is approximately 1° .

Because the neutron absorption cross section is negligible for most materials, the reflectivity, which can be calculated by the Fresnel equations (see, e.g., Chap. 63), is almost 100 percent below the critical angle. This permits the development of efficient grazing-incidence neutron optics.

64.3 DIFFRACTIVE SCATTERING AND MIRROR SURFACE ROUGHNESS REQUIREMENTS

Mirror surfaces must be quite smooth as diffractive scattering off surface irregularities produces the wings of a mirror’s point-spread function (PSF). The grating equation provides the relationship between the diffractive scattering angle θ for neutrons with wavelength λ incident, at an angle α , onto a surface with roughness at spatial frequency f

$$\theta \sin \alpha = f \lambda \quad (3)$$

The neutron flux excluded from an angular aperture with radius θ , the total integrated scatter (TIS) is given by

$$\text{TIS}(>\theta) = \int_{\theta'} \psi_{sc} 2\pi \vartheta' d\vartheta' = \left(\frac{4\pi \sin \alpha}{\lambda} \right)^2 \int_{\theta \sin \alpha / \lambda} \text{PSD}(f') df' = (4\pi \sin \alpha)^2 \left(\frac{\sigma}{\lambda} \right)^2 \quad (4)$$

where ψ_{sc} is the single reflection point spread function, PSD is the power spectrum density function of the surface, and σ is the surface roughness [the rms value, calculated for the spatial frequencies

$f > \theta(\sin\alpha/\lambda)$]. The formula describes the single reflection case and for double reflection, as for many optical configurations, the single reflection point spread function has to be multiplied by a factor of two.

At a given neutron wavelength and surface roughness, the excluded fraction of neutrons increases with grazing angle, so the worst case occurs at the critical angle of the surface material. For a nickel surface the critical angle is given by

$$\vartheta_{\text{cr}} \text{ (rad)} = 1.73 \times 10^{-3} \times \lambda \text{ (\AA)} \quad (5)$$

So, the worst case TIS for a nickel surface (single reflection) is

$$\text{TIS}_s(>\theta_{\text{cr}}) \approx 5 \times 10^{-4} \sigma_s^2 \quad (6)$$

for spatial wavelengths $1/f < 60/\theta$, where σ_s is the surface roughness (the rms value in Å). In practice, optical surfaces can be polished to the level of a few Å, so neutron scattering due to the surface microroughness can be kept to acceptable levels. Effects calculated above, though, will be doubled for 2-bounce mirror systems.

64.4 IMAGING FOCUSING OPTICS

Optical elements such as elliptical, toroidal, Kirpatrick-Baez (KB) or Wolter mirrors (see Chap. 44), and “tapered” (elliptical or parabolic) reflectors (see Chap. 52) utilize single or double reflections to image neutron sources and to produce fine focusing. The shallow critical angles for neutrons lead to mirror designs with large mirror lengths and small apertures. With low neutron currents, especially for thermal neutron sources, this dictates the use of nested systems to improve throughput. Fabrication of large super-polished high-figure-quality surfaces is a challenging task in general as figure errors and the roughness of the mirror surface determine the optical performance.

The specific properties of neutrons and the typical size of the neutron sources pose additional challenges (over similar x-ray optics) for development of high performance optics for beam control. Neutron sources are typically extended and also less brilliant compared to x-ray sources. Thus on-axis optical schemes, originally developed for point-like x-ray sources, may not deliver optimum optical performance. The major portion of an extended neutron source would be off the mirror’s optical axis, and this results in higher optical aberrations compared to an on-axis point source if classical focusing schemes are used. These aberrations can be minimized to some extent by collimating a neutron source down for applications when the neutron current is not so critical and by optimizing the figure of the mirror to boost the off axis optical performance. In addition, the effects of gravity, which need not be considered for x rays, are significant in neutron applications. The dependence of velocity on wavelength makes a grazing incidence neutron optic gravitationally chromatic.

Optical Elements: Configurations

Depending on the experimental task and resources available, the optimal optical configuration for a specific neutron application may include single, double, or more reflections.

Optical schemes using single reflection for neutron focusing can be divided into two subgroups based on neutron source configuration. An elliptical scheme is applicable for divergent point-like sources. In this case the source is placed at the first ellipsoid focal point and the neutrons are focused at the second focal point of the ellipsoid. A parabolic scheme would be applicable for quasi-parallel neutron beams, where a single-bounce optical element would be placed after straight neutron guides. The schematics of neutron focusing based on these optical configurations are shown in Figs. 1 and 2.

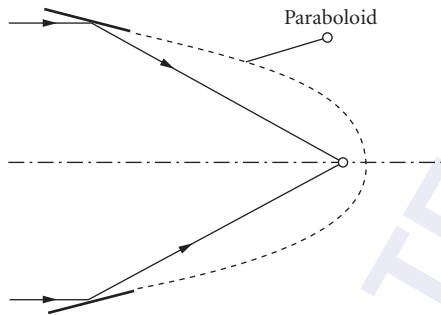


FIGURE 1 A single parabolic reflector used to focus a quasi-parallel input neutron beam.

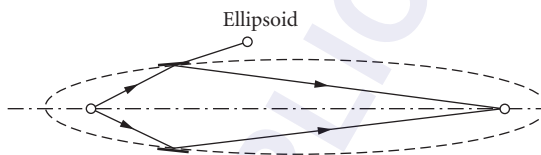


FIGURE 2 A single ellipsoidal reflector used to focus a divergent point source of neutrons.

A significant problem with any grazing-incidence focusing system is image degradation due to astigmatism and spherical and coma aberrations. The use of true ellipsoidal (for finite source) and paraboloidal (for infinite source) mirrors would eliminate the astigmatism and spherical aberration components, but in practice the figure of single-bounce grazing incidence mirrors is commonly approximated to ease mirror fabrication. Astigmatism, the optical aberration due to the difference in the meridional and sagittal focal distances of a mirror, can be eliminated when a toroidal mirror is used for neutron focusing. The mirror is astigmatism-free, if the sagittal and meridional radii of a toroidal mirror are set such that the corresponding focal distances are equal. This requirement is satisfied if the ratio between the sagittal and meridional radii of the mirror is

$$\frac{R_s}{R_m} = \sin^2 \alpha$$

where R_s and R_m are the sagittal and meridional radii of the mirror curvature and α is the grazing angle. Toroidal neutron mirrors have been used to focus neutrons at the ILL Spin-Echo Spectrometer² and at the Jülich SANS instrument and reflectometer.^{3,4} (The Jülich toroidal mirror has now been transferred to the FRM neutron source in Garching.⁵) Elliptical and parabolic tapered neutron guides have also been used as a single-bounce systems.⁶⁻¹⁰

The principal disadvantage of single-bounce systems is off-axis aberrations, especially coma. The off-axis coma is due to the difference in angular magnification for neutrons reflected from different areas of the mirror located on the same mirror meridian and it cannot be corrected using only one reflection. Image deterioration caused by this optical aberration is directly proportional to the mirror length. The aberration might in practice be severe because neutron sources are extended and the quasi-parallel neutron beams formed by straight neutron guides have nonzero divergence, so in both these cases many neutrons are nonaxial. To limit the focal spot blur of a single-bounce optical system due to off-axis coma one can therefore decrease the mirror length and/or collimate the neutron source down. However, both lead to reduced throughput, but this can be compensated in some cases

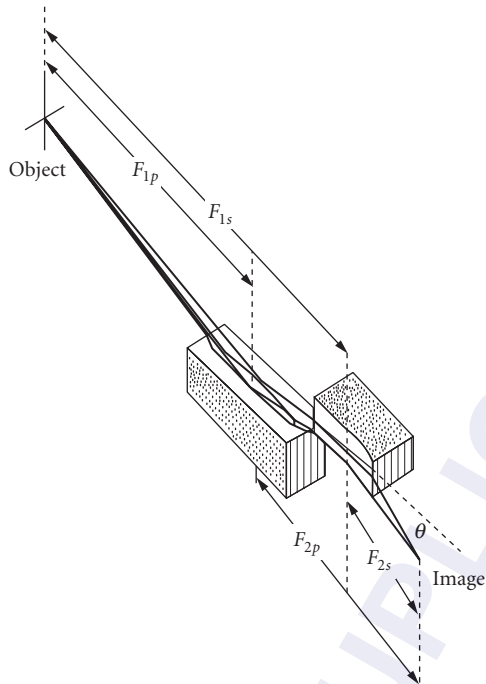


FIGURE 3 The Kirkpatrick-Baez mirror configuration uses a single reflection from each of two orthogonal mirrors. (Courtesy of Gene Ice, Oak Ridge National Laboratory.)

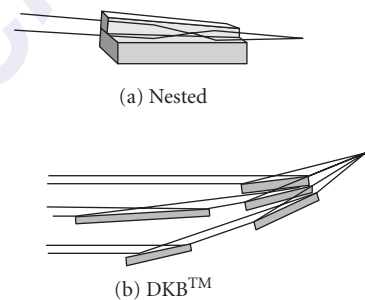


FIGURE 4 Two KB optical configurations for a neutron microscope. (Courtesy of Gene Ice, Oak Ridge National Laboratory.)

by the use of a nested optical system. A distinct advantage of single-bounce systems, however, is the ease of construction and alignment.

Another way to simplify the fabrication of grazing incidence mirrors is to use an optical scheme with two or more consecutive reflections from separate surfaces orthogonal to each other. Such a scheme was first described by Kirkpatrick and Baez (KB) for spherical mirrors¹¹ and is shown schematically in Fig. 3. The KB optical scheme, which has elliptical and parabolic mirrors for focusing beams from finite and infinite sources, respectively, eliminates astigmatism and spherical aberrations. However, because the reflections appear in orthogonal planes the KB scheme suffers from off-axis coma the same way as the single-bounce optical systems. The KB mirrors are relatively easy to fabricate but are relatively difficult to align due to sensitivity to the tilt errors. However, if bendable mirrors are used the misalignment can be often corrected by refocusing-rebending.

The KB systems provide good optical performance in small apertures if short mirrors are used, making the system attractive if neutrons need to be focused onto small samples. To increase the aperture, nested systems are needed and the KB systems lends itself readily to this with the mirrors stacked together in two orthogonal planes. A few nesting schemes have been proposed for KB microscopes and two examples are shown in Fig. 4. In one method, shown in Fig. 4a, mirrors are nested against each other so that the first reflection can be in either the horizontal or vertical plane. Another method, termed deflected KB (DKB), utilizes flat deflection mirrors to “steer” the incoming flux into the elliptical focusing mirrors, enabling a larger collecting area. A one-dimensional schematic illustrating focusing in one plane is depicted in Fig. 4b with the nested elliptical mirrors shown at the right. Inner rays enter directly, but outer rays are first deflected by the plane mirrors on the left. A second system, orthogonal to this, would provide focusing in the other dimension. Greater neutron current gains are

possible with more deflections. However, the increase in effective aperture due to reflection from an additional mirror have to be traded against the increase in image blur due to the neutron scattering from additional mirror surface's microroughness. In cases where the KB optical system is intended for neutron imaging, the difference in magnifications in orthogonal planes, proportional to the distance between the two orthogonal mirrors, needs to be taken into account.

The KB scheme has been successfully used for the neutron focusing.¹² The smallest focal spot achieved to date was using bendable mirrors with the figure adjusted to approximate an ellipse. The focal spot for neutrons with wavelength of 0.1 nm was measured to be $89 \times 90 \mu\text{m}$ (FWHM).¹³

The astigmatism of grazing incidence mirrors can be easily corrected by the use of optical geometries with an even number of reflections from the confocal surfaces of revolution. The on-axis coma is eliminated while the off-axis coma can be also be reduced for large apertures. With these aberrations either eliminated or reduced, the image degradation due to the obscured aperture diffraction would need to be taken into account.¹⁴ Since each reflection from a nonideal surface, in a multibounce system, results in additional scattering due to microroughness, the two bounce systems have been developed the most. Such systems, first discussed by Wolter,¹⁵ are used extensively in x-ray astronomy and can also be designed for use with cold and thermal neutron beams. In Wolter geometries two mirrors with surface figures of second order such as a paraboloid, hyperboloid, or ellipsoid are arranged coaxially. A schematic of two of the most used Wolter mirror configurations is shown in Fig. 5.

Wolter geometries with paraboloid, hyperboloid, and ellipsoid mirrors are optimized to provide superior optical performance on-axis and are capable of delivering better optical performance compared to single-bounce and KB systems. It has been previously shown in the literature that the application of a polynomial approximation to the surface of grazing incidence mirror shells^{16,17} as well as defocusing¹⁸ enhances the global performance of the mirror over the entire field-of-view. Nested Wolter-1 geometry optics can greatly improve the focused beam intensity by increasing the incident beam area while keeping the optical aberrations low. For reference, a picture of an x-ray optic module with 12 nested mirrors is shown in Fig. 6.

The focusing capabilities of neutron imaging optics based on Wolter geometries has been successfully demonstrated using a beam of cold neutrons.¹⁹ A test mirror originally designed as a 1/10-scale

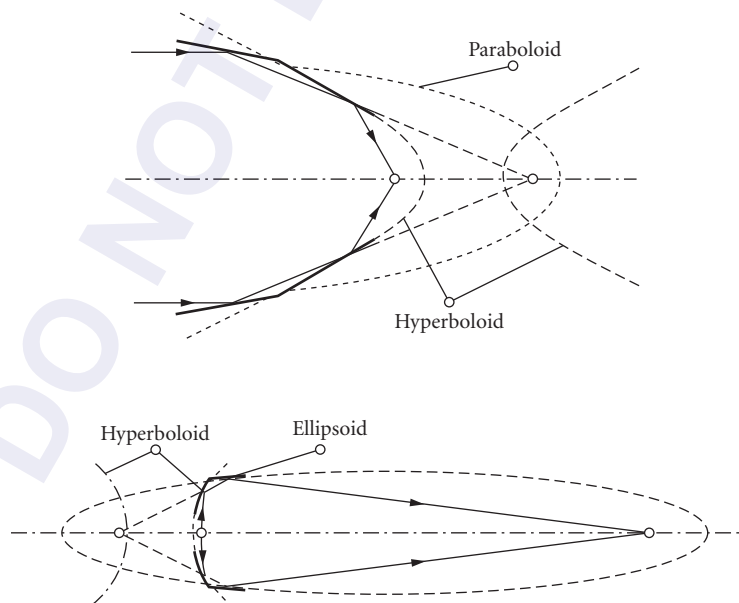


FIGURE 5 Commonly used Wolter-1 mirror configurations. (See also color insert.)



FIGURE 6 An x-ray optic module with 12 Wolter-1 nested mirrors.

version of the innermost mirror of NASA's Chandra X-Ray Observatory²⁰ was fabricated using an electroformed nickel replication process, wherein a thin nickel mirror shell is electroformed off a superpolished and figured mandrel.²¹ In the experiment the area of the optic illuminated by the neutron beam was estimated to be 17.7 mm². A focal-spot size, defined mostly by the divergence of the incoming beam, was measured to be 1.15 mm (FWHM).

Optical Elements: Materials

The largest possible critical angle observed for an isotope is that of Ni⁵⁸, and so this is the material of choice for grazing incidence optics. Typically, optical elements are fabricated from another more traditional material such as zerodur, glass, or silicon, and then the mirror surface is coated with Ni⁵⁸. In the case of the nickel replication technique, natural nickel can be used to produce the mirror shell with a small loss of critical angle. In some applications it is desirable to use nonmagnetic materials to preserve neutron polarization and for these Cu⁶⁵ with an aluminum-layer overcoat to prevent oxidation can be used.³

Multilayer coatings, typically alternating layers of nickel and titanium, can be deposited on the surface of a curved neutron mirror to increase the acceptance angle of the mirror. Such a neutron "supermirror" efficiently reflect neutrons with wavelengths λ when the mirror angle is less than a critical angle,¹³ given by

$$\vartheta_{cr}(\text{rad}) = m \times 1.73 \times 10^{-3} \times \lambda(\text{\AA}) \quad (7)$$

where m is the ratio of the maximum supermirror angle compared to the maximum efficient scattering angle of a nickel surface. Typical practical values of m are 3 to 4.

64.5 REFERENCES

1. V. F. Sears, *Neutron Optics*, Oxford University Press, Oxford, 1989, p. 64.
2. C. Hayes, C. Lartigue, A. Kolmar, J. R. D. Copley, B. Alefeld, F. Mezei, D. Richter, and T. Springer, "The Focusing Mirror at the ILL Spin-Echo Spectrometer IN15: Experimental Results," *Journal of the Physical Society of Japan* 65(suppl. A):312–315 (1996).

3. B. Alefeld, C. Hayes, F. Mezei, D. Richter, and T. Springer, "High-Resolution Focusing SANS with Toroidal Neutron Mirror," *Physica B* **234–236**:1052–1054 (1997).
4. B. Alefeld, L. Dohmen, D. Richter, and Th. Brückel, "X-Ray Space Technology for Focusing Small-Angle Neutron Scattering and Neutron Reflectometry," *Physica B* **283**:330–332 (2000).
5. Institut Fur Festkorperforschung: http://www.fz-juelich.de/iff/pics_pdfs/ism/2007/PNI_Kentzinger.pdf, accessed May 2009.
6. S. Yamada, T. Shinohara, H. Sasao, T. Oku, J. Suzuki, H. Matsue, and H. M. Shimizu, "Development of a Multichannel Parabolic Guide for Thermal Neutron Beam Focusing," *Physica B* **385–386**:1243–1246 (2006).
7. N. Kardjilov, P. Böni, A. Hilger, M. Strobl, and W. Treimer, "Characterization of a Focusing Parabolic Guide Using Neutron Radiography Method," *Nuclear Instruments and Methods in Physics Research A* **542**: 248–252 (2005).
8. T. Hils, P. Boeni, and J. Stahn, "Focusing Parabolic Guide for Very Small Samples," *Physica B* **350**:166–168 (2004).
9. S. Mühlbauer, M. Stadlbauer, P. Böni, C. Schanzer, J. Stahn, and U. Filges, "Performance of an Elliptically Tapered Neutron Guide," *Physica B* **385–386**:1247–1249 (2006).
10. S. Mühlbauer, P. G. Niklowitz, M. Stadlbauer, R. Georgii, P. Link, J. Stahn, and P. Böni, "Elliptic Neutron Guides—Focusing on Tiny Samples," *Nuclear Instruments and Methods in Physics Research A* **586**:77–80 (2008).
11. P. Kirkpatrick and A. Baez, "Formation of Optical Images by X Rays," *Journal of Optical Society of America* **38**: 766 (1948).
12. G. E. Ice, C. R. Hubbard, B. C. Larson, J. W. L. Pang, J. D. Budai, S. Spooner, and S. C. Vogel, "Kirkpatrick-Baez Microfocusing Optics for Thermal Neutrons," *Nuclear Instruments and Methods in Physics Research A* **539**:312–320 (2005).
13. G. E. Ice, C. R. Hubbard, B. C. Larson, J. W. L. Pang, J. D. Budai, S. Spooner, S. C. Vogel, R. B. Rogge, J. H. Fox, and R. L. Donabarger, "High-Performance Kirkpatrick-Baez Supermirrors for Neutron Milli- and Micro-Beams," *Materials Science and Engineering, A* **437**:120–125 (2006).
14. P. L. Thompson and J. E. Harvey, "Aplanatic Wolter Type 1 Telescope Design: Is there a Practical Advantage?," *Proceedings of SPIE* **3444**:526–542 (1998).
15. H. Wolter, "Spiegelsysteme streifend einfallend als abbildende optiken für Röntgenstrahlen," *Annalen der Physik* **6**(10):94 (1952).
16. P. Conconi, G. Pareschi, S. Campana, G. Chincarini, and G. Tagliaferri, "Wide-Field X-Ray Imaging for Future Missions, Including XEUS," *Proceedings of SPIE* **5168**:334–345 (2004).
17. C. J. Burrows, R. Burg, and R. Giacconi, "Optimal Grazing-Incidence Optics and Its Application to Wide-Field X-Ray Imaging," *Astrophysical Journal* **392**:760 (1992).
18. L. P. VanSpeybroeck and R. C. Chase, "Design Parameters of Paraboloid-Hyperboloid Telescopes for X-Ray Astronomy," *Applied Optics* **11**(2):440–445 (1972).
19. M. V. Gubarev, B. D. Ramsey, D. E. Engelhaupt, J. M. Burgess, and D. F. R. Mildner, "An Evaluation of Grazing-Incidence Optics for Neutron Imaging," *Nuclear Instruments and Methods in Physics Research B* **265**(2):626–630 (2007).
20. Chandra X-Ray Center: http://cxc.harvard.edu/cdo/about_chandra/overview_cxo.html, accessed May 2009.
21. B. D. Ramsey, "Replicated Nickel Optics for the Hard-X-Ray Region," *Experimental Astronomy* **20**(1):85–92 (2006).

INDEX

Index note: The *f* after a page number refers to a figure, the *n* to a note, and the *t* to a table.

- Abbe's sine condition, 45.3
- Aberrations:
of gratings and monochromators, 38.7
in grazing incidence optics, 45.1–45.8, 45.2*f*
in grazing incidence telescopes, 44.6–44.12, 44.7*f*, 44.9*f*–44.11*f*
in grazing-incidence neutron optics, 64.4, 64.6
and imaging through atmospheric turbulence, 4.17–4.20, 4.17*f*–4.18*f*, 4.19*t*, 4.20*t*, 4.27–4.30, 4.27*f*
Zernike modes of, 5.11, 5.11*t*, 5.12*f*
- Absolute method (scatterometer calibration), 1.15
- Absorption:
and atmospheric optics, 3.4–3.5, 3.4*f*
of Cr³⁺, 2.19–2.21, 2.19*f*, 2.20*f*
electro-, 13.55–13.56, 13.56*f*, 13.58
excited state, 14.6
measurements of, 2.2–2.13, 2.4*f*, 2.6*f*, 2.8*f*, 2.10*f*, 2.12*f*
molecular, 3.12–3.15, 3.13*f*, 3.22–3.23, 3.22*f*, 3.23*f*
in neutron optics, 63.6
photo-, 36.1
in photonic crystal fibers, 11.20–11.21, 11.20*f*
and reabsorption of spectra, 56.8
and x-ray optics, 26.7
- ac Kerr effect, 7.11
- Achromaticity, of fiber-based couplers, 16.3–16.4
- Acoustically rotated tangential phase matching, 6.26*f*, 6.27
- Acousto-optic devices, 6.3–6.45
and acousto-optic diffraction, 6.4, 6.9
Bragg cells (wideband), 6.30, 6.31*t*
deflectors, 6.22–6.30, 6.24*f*, 6.26*f*, 6.28*f*, 6.29*t*
figures of merit for, 6.16–6.17
materials for, 6.16–6.22
- Acousto-optic devices (*Cont.*):
modulators, 6.31–6.35, 6.34*t*
acousto-optic frequency shifters, 6.35
and Bragg diffraction, 6.4, 6.6, 6.7, 6.14
image (scophony), 6.34–6.35
principle of operation, 6.32, 6.32*f*
optical birefringence in, 6.17
propagation and attenuation in, 6.17
theory of acousto-optic interaction, 6.5–6.16
elasto-optic and roto-optic effects, 6.5–6.6
of finite geometry, 6.14–6.16
frequency characteristics, 6.12–6.14, 6.13*f*, 6.14*f*
phase matching, 6.9–6.12, 6.10*f*
plane wave analysis, 6.6–6.9
tunable filters as [*see* Acousto-optic tunable filters (AOTFs)]
wideband AO Bragg cells, 6.30, 6.31*t*
- Acousto-optic diffraction, 6.4, 6.9
- Acousto-optic figures of merit, 6.16–6.17
- Acousto-optic frequency shifters (AOFS), 6.35
- Acousto-optic interaction, theory of, 6.5–6.16
elasto-optic and roto-optic effects, 6.5–6.6
of finite geometry, 6.14–6.16
frequency characteristics, 6.12–6.14, 6.13*f*, 6.14*f*
phase matching, 6.9–6.12, 6.10*f*
plane wave analysis, 6.6–6.9
- Acousto-optic modulators, 6.23*t*, 6.31–6.35, 6.34*t*
acousto-optic frequency shifters, 6.35
image (scophony), 6.34–6.35
principle of operation, 6.32, 6.32*f*
- Acousto-optic tunable filters (AOTFs), 6.23*t*
collinear beam, 6.43, 6.43*f*, 6.45*t*
long-infrared, 6.42

- Acousto-optic tunable filters (AOTFs) (*Cont.*):
 and longitudinal spatial modulation, 6.12
 mid-infrared, 6.42
 noncritical phase-matching, 6.37, 6.39–6.42,
 6.38f, 6.43t
 principle of operation, 6.36–6.39,
 6.37f, 6.38f
 ultraviolet, 6.42
- Activated phosphor detectors, 60.7–60.8
- Active mode locking, 20.17
- Active optics, 50.1, 50.2
- Adaptive microstructured optical arrays
 (MOAs), 50.7–50.8, 50.8f
- Adaptive optics (AO), 5.1–5.46, 50.1
 concepts and components of, 5.2–5.5, 5.3f
 designing systems of, 5.38–5.46
 requirements, 5.39, 5.40t
 results for 3.5-m telescope systems,
 5.43–5.45, 5.43f–5.45f
 results for 10-m telescope systems,
 5.45–5.46, 5.45f–5.46f
 subaperture size, 5.40–5.43, 5.42f
- as enabling technology, 5.2
- hardware and software implementation for,
 5.21–5.38
 higher-order wavefront sensing techniques,
 5.36–5.37
 laser beacons, 5.27–5.34, 5.28f–5.31f,
 5.33f
 real-time processors, 5.34–5.35, 5.34f,
 5.35f
- Shack-Hartmann technique, 5.23–5.27,
 5.23f, 5.25f, 5.26f
- tracking, 5.21–5.23, 5.22f
- wavefront correctors, 5.37–5.38, 5.38f
- and imaging through atmospheric
 turbulence, 4.35–4.36
- and turbulence, 5.5–5.21
- anisoplanatism, 5.19
- atmospheric tilt and Strehl ratio, 5.14,
 5.14f–5.15f
- Fried's coherence diameter and spatial
 scale, 5.9–5.13, 5.11t, 5.12f, 5.13f
- higher-order phase fluctuations, 5.18–5.19
- on imaging and spectroscopy, 5.19–5.21,
 5.19f, 5.20f
- Kolmogorov model, 5.5–5.6
- tracking requirements, 5.15–5.18, 5.17f
- variation of n and C^2n parameters,
 5.6–5.8, 5.7f, 5.8f
- Adaptive x-ray optics, 50.1–50.8
 hard vs. soft x-ray telescopes, 50.2
 history of, 50.1, 50.2f
 for synchrotron and lab-based sources,
 50.2–50.7, 50.4f, 50.6f, 50.8f
- Add/drop multiplexers, optical, 21.2, 21.8,
 21.8f, 21.9f, 21.12
- ADONIS system, 5.35
- Advanced X-Ray Astrophysical Facility (AXAF),
 44.4, 44.6f, 44.10
- Aerosols, 3.6, 3.7, 3.10–3.11, 3.10f, 3.11f
- Air cladding, for photonic crystal fibers, 25.2,
 25.21
- Air-filled photonic crystal fibers, 25.20–25.21
- Airy disc, 37.6
- Airy distribution, of thin lenses, 40.3
- Akhieser loss, 6.17
- Algorithms, singular-value-decomposition,
 5.25, 5.26
- All solid-core photonic crystal fibers, 25.21
- All-fiber monolithic systems (fiber lasers),
 25.9–25.11, 25.16–25.18, 25.16f
- All-optical switching, in OTDM communication
 networks, 20.22, 20.23f
- All-solid bandgap guiding fibers, 11.15, 11.16,
 11.17f
- Aluminum oxide (Al_2O_3), 2.19, 2.19f, 2.21, 2.22
- Amplification:
 chirped pulse, 25.2, 25.32, 25.33
 lump, 21.44
 by SOAs, 19.1f–19.2f, 19.2, 19.22–19.27
 CWDM systems, 19.27
 DWDM systems, 19.25–19.27, 19.25f,
 19.26f
 single-channel systems, 19.22–19.24,
 19.22f–19.24f
- Amplified spontaneous emission (ASE), 9.13,
 10.11, 14.3, 14.6
- Amplified spontaneous emission (ASE) noise,
 19.3, 19.4f, 19.9, 19.18, 19.24, 19.35
- Amplified spontaneous emission (ASE) power,
 19.20
- Amplifiers:
 in communication systems, 9.13–9.14
 erbium-doped fiber, 14.4–14.7, 14.5f,
 21.38–21.41, 21.38f–21.42f
 erbium/ytterbium-doped fiber, 14.7–14.8
 field-effect transistor, 13.70
 linear optical, 19.14, 19.27
 optical fiber, 14.1–14.11, 14.2t

- Amplifiers (*Cont.*):
 parametric, 11.23, 11.24, 14.10–14.11
 praseodymium-doped fiber, 14.7
 Raman, 14.8–14.9, 14.8f, 14.10f,
 21.42–21.44, 21.42f–21.44f
 rare-earth-doped, 14.2–14.4, 14.3f
 semiconductor optical [see Semiconductor
 optical amplifiers (SOAs)]
 semiconductor vs. fiber, 9.13–9.14
 slab-coupled optical waveguide, 19.21
 in WDM networks, 21.37–21.44, 21.37f
 ytterbium-doped fiber, 14.7
- Amplitude modulation, 7.22–7.24, 7.23f, 7.24f
- Amplitude modulation index, 7.22
- Amplitude zone plates, 40.4–40.5, 40.5f
- Amplitude-modulated (AM) acoustic waves, 6.32
- Amplitude-modulated (AM) systems, optical
 fibers in, 9.16
- Amplitude-shift-keying (ASK), 21.28, 21.29
- Analog modulation, 6.32
- Analog to digital conversion, 20.4, 20.8, 20.8f
- Analog transmission, on optical fibers,
 9.15–9.17, 9.16f
- Anger cameras, 63.33
- Angle of arrival, 4.23–4.26, 4.25f, 4.26f
- Angle of deflection, 6.40
- Angle scans, for scatterometers, 1.14
- Angled stripe designs, of SOAs, 19.8, 19.8f
- Angular apertures, of NPM AOTFs, 6.41
- Angular dispersion, 38.6
- Angular resolution:
 in astronomical x-ray optics, 47.10–47.11, 47.11f
 in Wolter x-ray optics, 47.2–47.3, 47.2f
- Angular spectrum of plane waves (ASW)
 approach, 6.14–6.15
- Anisoplanatism, 5.19
 and laser beacons, 5.27–5.29, 5.28f–5.30f
 and subaperture size, 5.42–5.43
- Anisotropic acoustic beam collimation, 6.25
- Anisotropic diffraction, 6.10
- Anisotropy, 17.2
- Annealed proton exchange (APE), 7.30
- Annular pupils, 4.10–4.16, 4.11f–4.15f, 4.15t
- Anodes, as x-ray tube sources, 54.8–54.9,
 54.12, 54.13
- ANSI (American National Standards Institute),
 23.4, 23.6
- Antiguides, for VCSELs, 13.46
- Antireflection (AR) coatings, 19.8, 19.8f, 19.20
- Antiresonance, 11.12
- Anti-Stokes scattering, 10.5
- Anti-Stokes waves, 10.9
- Apertures:
 angular, 6.41
 Gaussian, 37.6
 numerical, 9.4, 25.2, 25.18, 42.2
 order-selecting, 40.4, 40.5
 pinhole, 32.3
 and subaperture size, 5.40–5.43, 5.42f
- Apodization, 17.6–17.7, 17.7f
- Arabidopsis thaliana*, 37.10–37.11, 37.10f
- Area-detection x-ray imaging, 61.2, 61.2f
- ARROW structures, 11.15
- Arsenic trisulfide (As_2S_3), 12.7
- ASML Alpha Tool, 34.2, 34.3f
- Astigmatism, 13.12, 64.4, 64.6
- ASTM (American Society for Testing
 Materials), 1.4
- Astralate aberrations, 45.4
- Astronomical x-ray optics, 47.1–47.11
 angular resolution of, 47.10–47.11, 47.11f
 hard, 47.9–47.10
 history of, 47.1–47.2
 Kirkpatrick-Baez optics, 47.7–47.8, 47.8f, 47.9f
 Wolter, 47.2–47.7, 47.2f, 47.4f, 47.6f, 47.7f
- Astronomy, x-ray, 33.1–33.4, 33.3t
- ATM (asynchronous transfer mode), fiber optic
 standards and, 23.6
- Atmosphere, standard, 3.6–3.11, 3.7t, 3.8f–3.11f
- Atmospheric Infrared Sounder (AIRS), 3.39
- Atmospheric optics, 3.3–3.45
 and absorption of light, 3.4–3.5, 3.4f
 and atmospheric optical transmission,
 3.22–3.26, 3.22f–3.27f
 and composition of standard atmosphere,
 3.6–3.11, 3.7t, 3.8f–3.11f
 and global climate change, 3.43–3.45, 3.44f
 and meteorological optics, 3.40–3.43,
 3.41f–3.43f
 remote sensing in, 3.36–3.40, 3.37f–3.40f
 and theory of interaction of light with
 atmosphere, 3.11–3.22
 inelastic optical processes, 3.21–3.22
 Mie scattering, 3.16–3.18, 3.17f–3.19f
 molecular absorption, 3.12–3.15, 3.13f
 molecular emission and thermal spectral
 radiance, 3.18, 3.20, 3.20f
 molecular Rayleigh scattering, 3.15–3.16
 surface reflectivity and multiple scattering,
 3.21, 3.21f

- Atmospheric optics (*Cont.*):
 turbulence in, 3.26, 3.28–3.36
 beam spreading, 3.32–3.33, 3.33*f*
 beam wander, 3.31–3.32
 imaging and heterodyne detection, 3.34
 parameters for, 3.28–3.31, 3.29*f*, 3.30*f*
 scintillation, 3.34–3.36, 3.35*f*
- Atmospheric tilt, 5.14, 5.14*f*–5.15*f*
- Atmospheric turbulence, imaging through,
 4.1–4.37
 aberration variance and approximate Strehl
 ratio for, 4.27–4.28, 4.28*f*
 and adaptive optics, 4.35–4.36
 angle of arrival fluctuations, 4.23–4.26, 4.25*f*,
 4.26*f*
 and covariance and variance of expansion
 coefficients, 4.20–4.22, 4.21*t*,
 4.22*f*–4.23*f*
 Kolmogorov turbulence and atmospheric
 coherence length, 4.7–4.10, 4.8*f*, 4.9*f*
 long-exposure images, 4.3–4.7
 modal correction of turbulence, 4.28–4.30,
 4.29*t*, 4.30*f*
 and modal expansion of aberration function,
 4.17–4.20, 4.17*f*–4.18*f*, 4.19*t*, 4.20*t*
 and resolution of telescopes, 4.2–4.3
 short-exposure image, 4.31–4.35,
 4.31*f*–4.34*f*, 4.35*t*
 and systems with annular pupils, 4.10–4.16,
 4.11*f*–4.15*f*, 4.15*t*
- Atomic energy levels, 2.2–2.5, 2.4*f*
- Atomic scattering, x-ray optics and, 36.1, 36.2*f*
- Atomic spectra, 2.13
- AT&T (American Telephone & Telegraph), 21.1
- Attenuated total reflectance (ATR) waveguides,
 12.11
- Attenuation:
 defined, 21.13
 fiber, 21.13–21.14
 for fiber optic communication links, 15.7
 Lambert-Beer law of, 63.11
 linear, 31.1, 31.2
 neutron, 63.11–63.12
 in optical fibers, 9.4, 9.5*f*
 in photonic crystal fibers, 11.19–11.22,
 11.20*f*, 11.21*f*
 x-ray, 31.1–31.4, 31.2*f*, 31.3*f*
- Attenuators:
 for fiber-based couplers, 16.4
 for networking, 18.2, 18.9
 variable optical, 21.12, 21.13*f*
- Auger cascades, 59.4
- Auger energies, 36.3*t*, 36.9*t*
- Auger excitation peaks, 29.3
- Auger recombination, 13.27–13.28
- Autocovariance (ACV) function, 44.13
- Avalanche photodiode (APD) receivers, 9.8,
 9.10–9.11
- Avalanche photodiodes, 13.63, 13.71–13.73
- Average unregistered detected point spread
 function (AUDPSF), 44.14, 44.15*f*
- Axial and circumferential slope errors,
 45.7–45.8
- Babinet-Soleil compensators, 7.22
- Backward Brillouin scattering, 11.25
- Ballistic neutron guides, 63.17–63.18
- Balloon payloads, in astronomical x-ray optics,
 47.10
- Band edges, of photonic crystal fibers,
 11.9–11.10
- Band gap energy, 19.2
- Band-filling modulators, 13.62
- Bandpass response, in acousto-optic
 interaction, 6.15–6.16
- Bandwidth:
 defined, 20.1
 fiber, 21.2–21.3, 21.2*f*
- Batman doping, 25.20
- Bacons, laser, 5.27–5.34
 focus anisoplanatism, 5.27–5.29, 5.28*f*–5.30*f*
 mesospheric sodium laser beams, 5.32–5.34,
 5.33*f*
 Rayleigh, 5.30–5.32, 5.31*f*
- Beam spatial coherence, 58.4
- Beam splitters, 18.6, 18.6*f*
- Beam spreading, 3.32–3.33, 3.33*f*
- Beam steering, 6.27–6.29, 54.11
- Beam wander, 3.31–3.32
- Beamlines, for multilayer Laue lenses,
 42.9–42.10, 42.9*f*–42.12*f*
- Beer-Lambert law, 3.11, 3.20, 3.21
- Bend, of liquid crystals, 8.22, 8.23*f*
- Bend loss, of photonic crystal fibers,
 11.21–11.22, 11.21*f*
- Bendable optics, 50.3–50.4, 50.4*f*
- Bending magnet synchrotron radiation sources:
 brightness of, 55.9
 power of, 55.8–55.9, 55.8*f*
 radiation from, 55.3–55.6, 55.5*f*, 55.6*f*
- Bent crystals, 39.5–39.6, 39.5*f*–39.6*f*
- BepiColombo mission, 49.6

- Bessel functions, 6.6
- Bias, of electroabsorption modulators, 13.59
- Bidirectional reflectance distribution function (BRDF or BDRF), 1.4–1.6, 1.5f, 3.21
- Bidirectional scatter distribution function (BSDF), 1.6–1.7
- Bidirectional scatter distribution function (BSDF) scatterometers, 1.8, 1.8f
- Bidirectional transmission distribution function (BTDF), 1.6
- Bimorph mirrors, 5.37, 5.37f, 50.4, 50.4f, 50.6f
- Binary digits, in OTDM networks, 20.8
- Binary units, for electro-optic modulators, 7.27
- Biphase coding, 20.9, 20.9f
- Birefringence, 6.17, 8.19, 11.17
- Birefringent diffraction bandshapes, 6.13, 6.14, 6.14f
- Birefringent fibers, 16.5–16.6
- Birefringent phased array deflectors, 6.29
- Birefringent tangential phase matching, 6.25
- Bit error rate (BER), for fiber optic communications, 15.1–15.5, 15.3f, 15.4f, 15.8, 15.13, 15.15, 15.17
- Bit error ratio (BER), DQPSK and, 21.34
- Bit rate, for optical fibers, 9.12
- Blackman window function, 46.9, 46.10f
- Bloch waves, 11.22
- Blockage defects, in polycapillary x-ray optics, 53.4, 53.5, 53.5f
- Boltzmann population factor, 3.14, 3.20
- Bormann transmission, 43.3, 43.3f
- Bormann triangle, 63.27
- Born approximation, 27.2, 63.4, 63.5
- Boron neutron capture therapy (BNCT), 53.19
- Boron [$^{10}\text{B}(n, \alpha)$] reaction, 63.31
- Boundary quality, of multilayers, 41.5–41.7, 41.7t
- Bragg angle, 6.9, 30.5, 39.3
- Bragg cells (wideband), in acousto-optic devices, 6.27, 6.30, 6.31t
- Bragg condition, 6.32, 42.3f, 42.4, 42.10, 42.11f
- Bragg diffraction:
 - and acousto-optic modulators, 6.4, 6.6, 6.7, 6.14
 - and brightness of x-ray tube sources, 54.16
 - far and near, 6.8–6.9, 6.12
 - in neutron optics, 63.23
 - order of, 20.14
 - and phase matching equations, 6.11
 - of x-rays, 42.2
- Bragg fibers, 11.4
- Bragg geometry, for crystal monochromators, 39.2–39.4, 39.2f, 39.3f, 39.6
- Bragg gratings:
 - and DBR lasers, 20.14
 - fiber, 17.1–17.9
 - applications, 17.8–17.9, 17.8f
 - chirped, 21.22–21.23, 21.23f, 21.25–21.26, 21.25f
 - fabrication, 17.4–17.8, 17.5f–17.7f
 - and fiber lasers, 25.8, 25.16, 25.18, 25.30, 25.31
 - long-period gratings vs., 24.9, 24.11
 - photosensitivity, 17.2–17.3
 - properties of, 17.3–17.4, 17.4f
 - sensors based on, 24.5–24.8, 24.6f–24.7f
 - volume, 25.29, 25.30
- Bragg limit, 6.8
- Bragg mirrors, 13.28, 13.44
- Bragg planes, of crystal monochromators, 39.1, 39.2, 39.5
- Bragg reflection:
 - of crystals, 40.9
 - in interferometers, 63.26, 63.27
 - and linear polarization, 43.2–43.4, 43.3f, 43.4f, 43.6, 43.8
 - and liquid crystals, 8.10
 - in monochromators, 39.1, 63.24
 - and multilayers, 41.2
 - and x-ray absorption spectroscopy, 30.2, 30.4
- Bragg reflection monochromator, 39.4, 39.5
- Bragg reflector lasers, 13.7, 13.28–13.29, 13.29f
- Bragg reflectors, 13.45
- Bragg scattering cross section, 63.18
- Bragg transmission phase retarders, 43.6
- Bragg wavelength, 17.3, 20.15, 24.7
- Bragg-Brentano powder diffractometer, 28.1, 28.2f, 28.3, 28.3f
- Bragg-Fresnel lenses, 40.9–40.10, 40.9f
- Bragg's law, 26.8, 28.3, 30.1, 39.1, 39.4, 43.7, 63.14, 63.16, 63.24
- Bragg-symmetric crystals, 35.3
- Braking radiation (*see* Bremsstrahlung radiation)
- Bremsstrahlung photons, 63.12
- Bremsstrahlung radiation, 31.3
 - continuous, 54.4–54.6, 54.5f
 - from laser-generated plasmas, 56.1, 56.8–56.10
 - from pinch plasma sources, 57.1
 - and x-ray fluorescence, 29.3, 29.5, 29.6, 29.11

- Brewster angle, in VUV and x-ray region, **41.9**
- Bridge fiber method, of mode matching, **25.17**, **25.17f**
- Brightness:
 - of synchrotron radiation, **55.1**, **55.9**
 - x-ray tube sources, **54.11–54.15**, **54.13f–54.15f**
- Brillouin frequency shift, **10.8**
- Brillouin scattering:
 - backward, **11.25**
 - forward, **11.25**, **11.26**
 - photonic crystal fibers, **11.25–11.26**, **11.25f**
 - stimulated, **10.1**, **10.7–10.9**, **15.8**, **21.20**, **25.6**
- Brillouin zones, **11.9**
- Broadband transmission, **3.23–3.24**, **3.25f–3.26f**
- Broadening:
 - Doppler, **3.14**, **5.32**, **56.5**
 - homogeneous, **20.1**
 - Lorentzian, **3.14**
 - pressure, **56.5–56.6**
 - spectral, **56.6–56.8**
- Bulk electro-optic modulators, **7.3**, **7.16–7.28**
 - amplitude modulation, **7.22–7.24**, **7.23f**, **7.24f**
 - frequency modulation, **7.24–7.25**, **7.25f**
 - phase modulation, **7.18–7.20**
 - polarization modulation (dynamic retardation), **7.20–7.22**, **7.20f**, **7.21f**
 - scanners, **7.26–7.28**, **7.26f–7.28f**
- Bulk lasers, **25.5–25.6**
- Bunches, electron, **55.17**
- Buried heterostructure (BH) lasers, **13.5**
- Cable television (CATV), **9.16–9.17**
- Calamitic liquid crystals (LCs), **8.4**, **8.5**, **8.9**, **8.11f**
- Calar Alto telescope, **5.27**
- Cameras:
 - Anger, **63.33**
 - gamma, **32.2**
- Canada–France–Hawaii Telescope, **5.23**
- Capillary discharge devices, **57.3**, **57.3t**
- Capillary optics, **28.5** (*See also* Monocapillary x-ray optics; Polycapillary optics)
- Carbon dioxide lasers, **12.3t**, **12.9**, **12.13**
- Carbon nanotubes (CNTs), **54.10–54.11**
- Carrier sense multiple access with collision detection (CSMA/CD), **23.7**
- Carrier-suppressed return-to-zero (CSRZ) formats, in WDM networks, **21.30**, **21.31f**, **21.36t**, **21.37f**
- Carrier-to-noise ratio (CNR), for optical fibers, **9.15–9.16**, **9.16f**
- Cassegrain telescopes, **44.4**
- Catadioptric lens systems, **35.1**
- Cathodes, as x-ray tube sources, **54.10–54.11**
- Catoptric lens systems, **35.1**
- Cauchy equations, for liquid crystals, **8.20**, **8.21**
- Čerenkov effects, **11.24**
- Chalcogenides:
 - in fiber lasers, **25.27t**, **25.29**
 - as infrared fibers, **12.2**, **12.2f**, **12.3t**, **12.6**, **12.7**, **12.7f**, **12.13**
- Chandra X-Ray Observatory, **33.2–33.4**, **33.3t**, **44.4**, **44.10**, **47.1**, **47.4f**, **47.5**, **47.10**, **64.7**
- Channel electron multipliers (CEMs), **60.6–60.7**
- Channel power equalization, for EDFAs, **21.41**, **21.42f**
- Characteristic radiation:
 - Bremsstrahlung radiation as, **56.8–56.10**
 - from laser-generated plasmas, **56.2–56.10**
 - recombination radiation as, **56.10**
 - spectral lines as, **56.2–56.8**
 - from x-ray tube sources, **54.6–54.8**, **54.7f**
- Charge-coupled devices (CCDs):
 - in adaptive optics, **5.21**
 - as x-ray detectors, **60.7**, **60.8**, **60.10t**
 - x-ray imaging detectors in, **61.7–61.8**, **61.8f**
- Charge-injected devices (CIDs), **60.7**, **60.10t**
- Chebyshev polynomials, **46.6**
- Chemical vapor deposition (CVD), **25.26**
- Chip screening, of SOAs, **19.17**, **19.17f**
- Chiral liquid crystals (LCs), **8.8–8.13**, **8.11f**
- Chirally coupled core (CCC) fibers, **25.2**, **25.19f**, **25.21–25.22**
- Chirp and chirping, **20.1**
 - of electro-absorption modulators, **13.58**
 - of fiber Bragg gratings, **17.7**
 - frequency, **13.1**, **13.17–13.18**, **13.17f**
 - of lasers, **9.8**
 - and optical fibers, **10.3**, **10.11**
 - of solitons, **22.3**
- Chirped fiber Bragg gratings (FBGs), **21.22–21.23**, **21.23f**, **21.25–21.26**, **21.25f**
- Chirped pulse amplification (CPA), **25.2**, **25.32**, **25.33**

- Cholesteric liquid crystal display (Ch-LCD), 8.32
- Chromatic dispersion:
and fiber optic communication links, 15.9, 15.10
in WDM networks, 21.14–21.16, 21.15f, 21.16f
- Chromium ions (Cr^{3+}):
absorption and photoluminescence of, 2.19–2.21, 2.19f, 2.20f
optical spectroscopy of, 2.9–2.10, 2.10f
polarization spectroscopy of, 2.21–2.22, 2.23f
- Circuits, optical tank, 20.21, 20.21f
- Circuit-switched networks, 21.7–21.10, 21.8f–21.9f
- Circular gratings, 40.1
- Circular polarization:
analyzers for, 43.6–43.8, 43.7f
phase plates for, 43.5f
and synchrotron radiation, 55.6–55.7
- Circulators, for networking, 18.3, 18.3f, 18.10
- Circulators, optical: fiber Bragg gratings, 17.8f, 17.9
- Circumferential coordinates, 45.7
- Circumferential slope errors, 45.8
- Cladding:
defined, 25.2
photonic crystal fibers in, 11.7–11.11, 11.8f–11.10f
- Cleaving, of photonic crystal fibers, 11.26
- Climate change, global, 3.43–3.45, 3.44f
- Clock recovery, in OTDM networks, 20.21–20.22, 20.21f
- Coarse wavelength division multiplexing (CWDM) systems, 19.27
- Coatings:
antireflection, 19.8, 19.8f, 19.20
for infrared optical fibers, 12.4, 12.7, 12.9
reflective (multilayers), 41.1–41.10
and calculation of multilayer properties, 41.3–41.4
for diffractive imaging, 41.9–41.10, 41.10f
fabrication methods and performance of, 41.4–41.9, 41.5f, 41.6t, 41.7t, 41.8f
properties of, 41.1–41.3, 41.2f
- Coblentz sphere, 1.10, 1.10f, 1.11
- Code mark inversion (CMI), 20.9
- Code V (program), 35.1
- Coding, in OTDM networks, 20.9–20.10, 20.9f
- Coherence:
beam spatial, 58.4
in coherent x-ray optics, 27.5
in long wavelength limit, 55.17–55.20, 55.18f, 55.19f
mutual coherence function (MCF), 4.4, 4.7, 4.10
temporal, 55.17–55.18, 55.18f
transverse, 55.18–55.20, 55.19f
transverse spatial, 55.16
- Coherence collapse, from laser diodes, 13.22–13.23
- Coherence diameter, Fried's, 5.2, 5.7, 5.9–5.10
- Coherence length:
atmospheric, 4.7–4.10
Fried's, 4.8
and imaging through turbulence, 4.8–4.10, 4.8f, 4.9f
spatial, 27.2
- Coherent area, of atmosphere, 4.16
- Coherent diffraction microscopy, 27.4–27.5, 27.4f, 27.5f
- Coherent Doppler LIDAR, 3.38, 3.39f
- Coherent scattering, 26.7, 63.7, 63.8
- Coherent x-ray optics, 27.1–27.5, 27.3f–37.5f
- Coiling, of LMA fibers, 25.18
- Cold cathode fluorescent lamps (CCFLs), 8.30, 8.30f
- Cold-cathode field emission, 54.10–54.11
- Collimating polycapillary optics, 53.14, 53.14f
- Collimating single crystal diffraction, 53.12, 53.12f
- Collimation and collimators:
anisotropic acoustic beam, 6.25
neutron, 63.15–63.16
in neutron and x-ray optics, 26.11, 26.11f
in polycapillary x-ray optics, 53.8–53.9, 53.8f–53.9f
with refractive x-ray lenses, 37.8
Soller, 28.2, 28.2f, 28.3
in SPECT imaging, 32.3
- Collinear beam acousto-optic tunable filters (CBAOTFs), 6.43, 6.43f, 6.45t
- Collision length, for solitons, 22.9–22.10, 22.15
- Coma, of grazing incidence telescopes, 44.9, 44.10

- Communication networks and systems:
 fiber-optic standards for, 23.1–23.8
 ATM/SONET, 23.6
 ESCON, 23.1–23.2, 23.2f
 Ethernet, 23.7
 FDDI, 23.2–23.3, 23.3f
 Fibre Channel standard, 23.4, 23.5f, 23.5t
 InfiniBand, 23.8, 23.8t
 optical fibers in, 9.3–9.17
 analog transmission, 9.15–9.17
 bit rate, 9.12
 distance limits, 9.12–9.13
 fiber for, 9.4–9.7, 9.5f, 9.6f
 fiber-optic networks, 9.14–9.15
 optical amplifiers, 9.13–9.14
 photodetectors, 9.8
 receiver sensitivity, 9.8–9.11
 repeater spacing, 9.12–9.13
 technology, 9.4–9.8
 transmitting sources, 9.7–9.8
 optical time-division multiplexed, 20.1–20.25
 analog to digital conversion, 20.8, 20.8f
 binary digits and line coding in, 20.8–20.10, 20.9f
 device technology, 20.12–20.24
 history of, 20.3
 interleaving in, 20.6–20.7, 20.7f
 modulation, 20.17–20.20, 20.17f, 20.19f, 20.20f
 multiplexing and demultiplexing, 20.3–20.12, 20.7f, 20.22, 20.23f
 optical clock recovery, 20.21–20.22, 20.21f
 sampling, 20.4–20.6, 20.5f, 20.6f
 serial vs. parallel, 20.12, 20.13f
 timing recovery, 20.10–20.12, 20.10f, 20.11f
 transmitters, 20.12–20.17, 20.14f–20.16f
 ultrahigh-speed, 20.23–20.24, 20.24f
 solitons in, 22.1–22.17
 classical solitons, 22.2–22.4, 22.2f–22.4f
 design of transmission systems, 22.5–22.7
 dispersion-managed solitons, 22.12–22.15, 22.13f, 22.14f
 frequency-guiding filters, 22.7–22.9
 wavelength division multiplexing, 22.9–22.12
- Communication networks and systems (*Cont.*):
 wavelength-division multiplexed, 21.1–21.44
 carrier-suppressed return-to-zero and duobinary, 21.30–21.33, 21.31f, 21.32f
 chromatic dispersion in, 21.14–21.16, 21.15f, 21.16f
 circuit and packet switching in, 21.7–21.11, 21.8f–21.11f
 dispersion and nonlinearities of, 21.16–21.26, 21.17f–21.21f, 21.23f–21.27f
 DPSK and DQSK, 21.33–21.36, 21.33f–21.35f, 21.36t, 21.37t
 fiber attenuation and optical power loss, 21.13–21.14
 fiber bandwidth, 21.2–21.3, 21.2f
 fiber system impairments, 21.13–21.26
 history of, 21.1–21.2
 network reconfigurability, 21.12–21.13, 21.12f, 21.13f
 optical amplifiers in, 21.37–21.44, 21.37f, 21.38f–21.42f
 optical modulation formats, 21.27–21.36, 21.28f–21.30f
 point-to-point links, 21.4
 in real systems, 21.3–21.4, 21.3f, 21.4f
 star, ring, and mesh topologies, 21.5–21.7, 21.5f–21.7f
 wavelength-routed networks, 21.5, 21.5f
 WDM dispersion managed soliton transmission, 22.15–22.17
- Commutators, 20.1, 20.7f
 Compensators, Babinet-Soleil, 7.22
 Complementary metal oxide semiconductors (CMOS) signal processing, 62.5
 Complex index of refraction, for x-ray optics, 48.1
 Compton radiation sources, 55.2–55.3, 55.3t
 Compton scattering, 59.1
 and circular polarization, 43.6
 inverse, 59.1
 and microfocus x-ray fluorescence, 29.5, 29.8f
 and polycapillary x-ray optics, 53.3, 53.15, 53.18
 and refractive x-ray lenses, 37.5, 37.7
 and x-ray attenuation, 31.2
 and x-ray optics, 26.7, 36.1
 Computed tomography (CT), 31.5–31.7, 31.6f, 31.7f
 Condensing monocapillary x-ray optics, 52.6

- Cone effect, of laser beacons, 5.27
- Cone-beam computed tomography (CT), 31.7, 31.7f
- Configurational relaxation, in solids, 2.14–2.17, 2.15f–2.18f
- Confinement:
 - in double heterostructure laser diodes, 13.4–13.6, 13.4f, 13.6f
 - and photonic crystal fibers, 11.22
 - in semiconductor optical amplifiers, 19.6
- Connector losses, in fiber optic communication, 15.7, 15.8t
- Constellation-X Observatory, 33.4
- Continuous readout x-ray detectors, 61.2
- Continuous-wave (CW) lasers, 25.4, 25.5, 25.5f, 25.7
 - diode lasers, 13.49
 - dye lasers, 5.32, 5.33f, 5.34
- Contrast ratio, for fiber optic modulation, 13.57
- Contrast-to-noise ration (CNR), of x-ray detectors, 61.3
- Controlled-drift x-ray detectors, 62.5
- Converging neutron guides, 63.17
- Coordinate-measuring machines (CMMs), 46.2
- Coordinates, circumferential, 45.7
- Copper distributed data interface (CDDI), 23.3
- Core (term), 25.2
- Core drilling, of fiber lasers, 25.26–25.27
- Core-cladding index difference, of photonic crystal fibers, 11.12–11.17, 11.12f–11.17f
- Cornell High Energy Synchrotron Sources (CHESS), 52.3, 52.5
- Coulomb explosion, 2.5, 2.6
- Coulomb fields, of ions, 56.8
- Coulomb interactions, of 3d electrons, 2.9
- Coulomb repulsion, 2.5
- Coupled quantum wells (CQWs), 13.58
- Coupled-wave-theory (CWT), for multilayer Laue lenses, 42.4, 42.12
- Couplers and coupling:
 - directional, 18.2, 18.3, 18.3f, 18.9, 18.9f
 - evanescent, 25.11
 - fiber-based, 16.1–16.6, 16.2f, 16.5f
 - nonfused fiber, 25.10
 - vertical, 13.60
 - wavelength-selective, 14.2
- Coupling coefficient, of DBR lasers, 13.28
- Coupling loss, of fiber-optic components, 18.1
- Covariance mapping, 2.5
- Cross-gain modulation (XGM), 19.12, 19.13f, 19.27, 19.29–19.30, 19.29f, 19.32, 19.35–19.36
- Crossover frequency, liquid crystals and, 8.16
- Cross-phase modulation (XPM):
 - in optical fibers, 10.3–10.4
 - and SOAs, 19.13, 19.30–19.32, 19.31f, 19.33f, 19.35–19.36
 - and solitons, 22.5, 22.13–22.15
 - in WDM networks, 21.19
- Crucifix images, 49.4, 49.4f
- Cryogenic x-ray detectors, 60.9, 60.9t, 60.10t
- Crystal diffraction:
 - and EDXRF, 29.6–29.7, 29.8f–29.9f
 - single, 53.12–53.14, 53.12f–53.13f
 - and WDXRF, 29.2
- Crystal interferometers, 63.26–63.27, 63.27f
- Crystal monochromators:
 - and bent crystals, 39.1–39.6, 39.2t, 39.3f, 39.5f–39.6f
 - in neutron optics, 63.23–63.25
- Crystal optics:
 - and electro-optic modulators, 7.3–7.4, 7.4f, 7.8t–7.10t
 - and the index ellipsoid, 7.3–7.7, 7.4f–7.6f, 7.8f–7.10f
 - polarizing, 43.1–43.8, 43.3f–43.5f, 43.7f
 - and x-ray diffraction, 28.3–28.4, 28.4f
- Crystalline infrared fibers, 12.2t, 12.3t, 12.7–12.10, 12.8f, 12.10f
- Crystals:
 - bent, 39.5–39.6, 39.5f–39.6f
 - Bragg-symmetric, 35.3
 - Darwin width of, 63.24, 63.26
 - doubly curved, 29.6–29.7, 29.8f
 - in electro-optic modulators, 7.33–7.34
 - Fankuchen-cut, 63.25
 - mosaic, 39.2
 - organic, 7.33–7.34
 - SHADOW code for, 35.2
 - [See also Liquid crystals; Photonic crystal fibers (PCFs)]
- Current:
 - dark, 13.69, 13.73
 - saturation, 13.69
 - trap-assisted thermal generation, 13.69
- Current confinement, in VCSELs, 13.45, 13.45f
- Curvature of best focal surface, 45.5

- Damage, optical, **13.54**
- Damage threshold, of fiber lasers, **25.7**
- Dark current, **13.69, 13.73**
- Darwin widths, **43.6, 63.24, 63.26**
- Data communication systems, **15.1–15.2**
- Data transmission formats, **20.9–20.10, 20.9f**
- dc Kerr effect, **7.11**
- De Broglie neutron waves, **63.3, 63.5**
- Debye-Waller factor, **41.4**
- Decentering errors, for grazing incidence optics, **45.6–45.7**
- Decommutator, **20.1, 20.7f**
- Deep saturation regime, for fiber amplifiers, **14.4**
- Deflectors, **6.22–6.31, 6.23t**
 - anisotropic acoustic beam collimation by, **6.25**
 - birefringent phased array, **6.29**
 - high-resolution, **6.29, 6.29t**
 - isotropic AO diffraction by, **6.23–6.24, 6.24f**
 - phase array beam steering by, **6.27–6.29, 6.28f**
 - resolution of, **6.25**
 - tangential phase matching by, **6.25–6.27, 6.26f**
- Defocus errors, for grazing incidence optics, **45.6**
- Deformable mirrors, **5.4, 5.4f, 5.37–5.38, 5.37f, 5.38f**
- Degree of modulation, for electro-optic modulators, **7.35–7.36**
- Dektak stylus profiler, **46.2**
- Delay-line technique, **60.5**
- Demultiplexers and demultiplexing, **20.1**
 - for networking, **18.4, 18.10–18.11, 18.10f, 18.11f**
 - in OTDM networks, **20.7–20.8, 20.7f, 20.22, 20.23f**
 - terahertz asymmetric optical, **20.22**
- Dense wavelength division multiplexing (DWDM):
 - and AOTFs, **6.43, 6.44**
 - and optical fiber amplifiers, **14.1**
 - and SOAs, **19.25–19.27, 19.25f, 19.26f**
- Density of states (DOS), for photonic crystal fibers, **11.8, 11.9f**
- Depletion layer, of semiconductor detectors, **60.5**
- Depletion region, of pin photodiodes, **13.64**
- Deposition:
 - modified chemical vapor, **25.2, 25.21, 25.26, 25.28**
 - outside vapor, **25.2, 25.26**
 - physical vapor, **61.6**
 - vapor axial, **25.3, 25.26**
- Deposition, chemical vapor, **25.26**
- Depth of focus (DOF), of multilayer Laue lenses, **42.17, 42.17f**
- Depth-of-amplitude modulation, **7.22**
- Depth-of-phase modulation, **7.19**
- Despace errors, for grazing incidence optics, **45.6**
- Destructive ports, of M-Z interferometers, **21.34**
- Detected point spread function (DPSF), **44.14, 44.15f**
- Detection and detectors:
 - activated phosphor, **60.7–60.8**
 - collision, **23.7**
 - direct-conversion flat panel, **61.4f, 61.6–61.7, 61.6t, 61.7f**
 - energy dispersive, **62.2–62.4**
 - for fiber optic systems, **13.2–13.3, 13.63–13.73**
 - avalanche photodiodes, **13.71–13.73**
 - MSM detectors, **13.73**
 - pin diodes, **13.63–13.71, 13.64f, 13.66f, 13.66t, 13.68f**
 - Schottky photodiodes, **13.73**
 - heterodyne, **3.34**
 - light detection and ranging (LIDAR), **3.38–3.39, 3.38f, 25.25**
 - for medical imaging, **31.4, 31.4f, 61.2**
 - MSM photoconductive, **13.63, 13.73**
 - in neutron optics, **63.31–63.34**
 - photodetectors, **9.8**
 - semiconductor, **60.5**
 - solid state, **63.33–63.34**
 - x-ray, **31.4f**
 - controlled-drift, **62.5**
 - cryogenic, **60.8–60.9, 60.9t, 60.10t**
 - film, **60.8, 60.9t, 60.10t**
 - ionization, **60.3–60.7, 60.9t, 60.10t**
 - scintillation, **60.7–60.8, 60.9t, 60.10t**
 - for x-ray imaging, **61.1–61.8**
 - CCD detectors, **61.7–61.8, 61.8f**
 - flat panel detectors, **61.3–61.7, 61.4f, 61.6t, 61.7f**
 - geometries for and classifications of, **61.1–61.3, 61.2f**

- Detective quantum efficiency (DQE), of x-ray detectors, **61.2**, **61.3**, **61.5–61.7**
- Detrending, in x-ray mirror metrology, **46.6–46.7**, **46.7f**
- Detuning, **13.29**, **30.3**
- Dichroism:
 defined, **25.2**
 magnetic circular, **55.7–55.9**, **55.7f**
- Dielectric properties, of liquid crystals, **8.14–8.18**, **8.15f**, **8.18f**
- Differential group delay (DGD), **21.17**
- Differential phase shift keying (DPSK), **19.27**, **21.33–21.34**, **21.36**, **21.36t**, **21.37t**
- Differential quadrature phase-shift-keying (DQPSK), **21.34–21.36**, **21.35f**, **21.36t**, **21.37t**
- Differential quantum efficiency, **13.9**
- Diffraction:
 acousto-optic, **6.4**, **6.9**
 anisotropic, **6.10**
 and birefringent diffraction bandshapes, **6.13**, **6.14**, **6.14f**
 Bragg
 and acousto-optic modulators, **6.4**, **6.6**, **6.7**, **6.14**
 and brightness of x-ray tube sources, **54.16**
 far and near, **6.8–6.9**, **6.12**
 in neutron optics, **63.23**
 order of, **20.14**
 and phase matching equations, **6.11**
 of x-rays, **42.2**
 by crystals, **39.2–39.3**
 and energy-dispersive x-ray fluorescence, **29.6–29.7**, **29.8f–29.9f**
 Fraunhofer, **42.2**, **63.25**
 Fresnel, **27.2**, **27.4**, **40.9**, **63.25**
 by gratings and monochromators, **38.1–38.3**, **38.2f**
 isotropic
 and acousto-optic interactions, **6.9–6.10**, **6.10f**, **6.13**, **6.13f**
 by deflectors, **6.23–6.24**, **6.24f**
 and isotropic diffraction bandshape, **6.13**, **6.13f**
 kinematic theory of, **63.5**
 lysozyme, **53.12–53.14**, **53.12f**, **53.13f**
 near- and far-field, **27.2–27.3**, **27.3f**
 in neutron optics, **63.3**, **63.25**
 order of, **18.6**
- Diffraction (*Cont.*):
 powder, **53.14**, **53.14f**
 single crystal, **53.12–53.14**, **53.12f–53.13f**
 and wavelength-dispersive x-ray fluorescence, **29.2**
 x-ray, **28.5–28.6**, **28.6f**
 and x-ray optics, **26.7**, **26.8**
- Diffraction efficiencies, of zone plates, **40.4–40.8**, **40.5f**, **40.7f**, **40.8f**
- Diffraction geometry, tilted, **42.4**
- Diffractional imaging, **41.9–41.10**, **41.10f**
- Diffractional optics, **26.8**, **26.8f**, **26.9f**
- Diffractional scattering, **64.2–64.3**
- Diffractionometers, **28.1–28.3**, **28.2f**, **28.3f**
- Diffuse scattering, **63.4**
- Digital, analog conversion to, **20.4**, **20.8**, **20.8f**
- Digital displays, in medical imaging, **31.8–31.9**, **31.8f–31.10f**
- Digital modulation, **6.33**
- Digital on-off-keying receivers, **9.9–9.11**
- Digital tomosynthesis, **31.7–31.8**
- Digital Video Broadcasters Measurement Group standards, **15.4**
- Diode pumps, **25.12–25.13**
- Diodes:
 laser, **13.3–13.24**
 double heterostructure, **13.3–13.8**, **13.3f**, **13.4f**, **13.6f**, **13.7f**
 noise characteristics of, **13.18–13.24**, **13.19f–13.21f**
 operating characteristics of, **13.8–13.13**, **13.10f**
 transient response of, **13.13–13.18**, **13.14f**, **13.16f**, **13.17f**
- light-emitting, **13.1**, **13.36–13.42**, **13.38f**
 edge-emitting, **13.40**
 operating characteristics of, **13.40–13.42**, **13.40f**, **13.41–13.42**
 and optical fibers, **9.7**
 surface-emitting, **13.38–13.40**, **13.38f**
 and transmissive TFT LCDs, **8.29–8.31**, **8.30f**
 pin, **13.2**, **13.63–13.71**, **13.66t**
 dark current, **13.69**
 geometry of, **13.64–13.65**, **13.64f**
 noise, **13.70–13.71**
 sensitivity, **13.65–13.66**, **13.66f**
 speed, **13.67–13.68**
 and unitraveling-carrier (UTC) photodiodes, **13.68–13.69**, **13.68f**
 (*See also* Photodiodes)

- Dioptric lens systems, **35.1**
- Direct modulation, **20.17–20.18, 20.17f**
- Direct-conversion flat panel detectors, **61.4f, 61.6–61.7, 61.6t, 61.7f**
- Directional couplers, for networking, **18.2, 18.3, 18.3f, 18.9, 18.9f**
- Directionality (isolation), of fiber-optic components, **18.1**
- Discrete electroabsorption modulators, **13.59**
- Dispenser cathodes, **54.10**
- Dispersion:
- angular, **38.6**
 - chromatic
 - and fiber optic communication links, **15.9, 15.10**
 - in WDM networks, **21.14–21.16, 21.15f, 21.16f**
 - in fiber optic communication links, **15.9–15.11, 15.10f**
 - by gratings and monochromators, **38.6–38.7**
 - group velocity, **11.18, 11.18f, 11.19f**
 - in optical fibers, **9.5–9.7, 9.6f**
 - theory of, **63.5**
 - in WDM networks, **21.20–21.26, 21.21f**
 - chromatic, **21.14–21.16, 21.15f, 21.16f**
 - electronic solutions for, **21.26, 21.27f**
 - fixed dispersion compensation, **21.22–21.23, 21.23f**
 - tunable dispersion compensation, **21.23–21.26, 21.24f–21.26f**
- Dispersion compensating fiber (DCF), **15.10, 21.22**
- Dispersion length, of solitons, **22.3**
- Dispersion optimizing fiber, **15.10**
- Dispersion-managed (DM) solitons, **22.12–22.15, 22.13f, 22.14f**
- Dispersion-shifted fibers, **9.7**
- Displacement vectors, bulk modulators and, **7.18**
- Displays, for medical imaging, **31.8–31.9**
- Distance limits, of optical fibers, **9.12–9.13**
- Distorted object approach, in coherent x-ray optics, **27.2–27.4**
- Distortion(s):
- in fiber optic communication links, **15.5–15.6, 15.6f**
 - in optical fibers, **9.17**
- Distributed Bragg reflector (DBR) lasers, **9.8, 13.7, 13.28–13.29, 13.29f, 20.14**
- Distributed feedback (DFB) lasers:
- in fiber optic systems, **13.7, 13.30–13.32, 13.30f, 13.31f**
 - optical fibers for, **9.8**
 - in OTDM communication networks, **20.14–20.15, 20.15f**
- Distributed feedback (DFB) threshold, of fiber optic systems, **13.30–13.32**
- Distributed Raman amplifiers (DRAs), **21.44**
- Divalent rare-earth ions, **2.11**
- Divided voltage method, for transreflective LCDs, **8.35, 8.35f**
- Dopant profiling, for LMA fibers, **25.18**
- Dopants, for fiber lasers, **25.22–25.26, 25.23t**
- Doppler broadening, **3.14, 5.32, 56.5**
- Doppler effect, **2.3, 6.9**
- Doppler LIDAR systems, **3.38–3.39**
- Doppler profiles, of CW lasers, **5.32, 5.34**
- Doppler shifts, **10.7**
- Doppler-dominated lineshapes, **3.23**
- Double heterostructure laser diodes, **13.3–13.8, 13.3f, 13.4f, 13.6f, 13.7f**
- Double heterostructure waveguides, **19.3f, 19.5**
- Double-bounce Wolter mirrors, **52.4**
- Doubly curved crystals (DCCs), **29.6–29.7, 29.8f**
- Dragon* (monochromator) systems, **38.3**
- Drive circuitry, of LEDs, **13.42**
- Drive power, of NPM AOTFs, **6.41**
- Dual attach FDDI nodes, **23.3**
- Dual frequency effect, on LCDs, **8.18**
- Dual-cell-gap transreflective LCDs, **8.32–8.33, 8.33f**
- Duobinary formats, for WDM networks, **21.32–21.33, 21.32f, 21.36, 21.36t, 21.37t**
- Duplexers, for networking, **18.4, 18.4f, 18.10–18.11**
- Dye-doped polymers, in electro-optic modulators, **7.34**
- Dynamic Jahn-Teller effect, **2.9**
- Dynamic range:
- of fiber amplifiers, **14.3**
 - of fiber optic communication links, **15.5–15.6, 15.6f**
 - of wideband Bragg cells, **6.30**
- Dynamic retardation, of modulators, **7.20–7.22, 7.20f, 7.21f**
- Dynamical theory of diffraction, for crystals, **39.2**

- Edge-coupled pin waveguides, **13.68**
- Edge-defined film-fed growth (EFG) technique, **12.9, 12.10f**
- Edge-emitting lasers, **13.3**
- Edge-emitting light-emitting diodes (E-LEDs), **13.36, 13.37, 13.40**
- EEGRAZE code, **44.13**
- Effective group refractive index, **13.13**
- Efficiency:
- detective quantum, **61.2, 61.3, 61.5–61.7**
 - differential quantum, **13.9**
 - diffraction, **40.4–40.8, 40.5f, 40.7f, 40.8f**
 - external quantum, **13.9**
 - external slope, **13.9**
 - gain, **14.6**
 - internal quantum, **13.9, 13.10**
 - modulation, **7.36**
 - quantum, **13.65, 13.66**
 - quantum detection, **60.7, 61.2–61.3**
 - radiative quantum, **14.7**
- Eigenpolarization, **7.13–7.14, 7.14f**
- Einstein coefficient for spontaneous emission, **2.13**
- Einstein Observatory, **44.4, 44.10, 47.1, 47.5**
- Elastic constants, of liquid crystals, **8.22, 8.23f**
- Elastic scattering, in neutron optics, **63.5**
- Elasto-optic effect, **6.5–6.6**
- Electric dipole selection rules, **56.3**
- Electrical injection:
- in laser diodes, **13.4, 13.5**
 - in VCSELs, **13.45, 13.45f**
- Electrical time domain multiplexed (ETDM) transmission, **20.3, 20.25**
- Electroabsorption, **13.55–13.56, 13.56f**
- applying fields in semiconductors, **13.58**
- Electroabsorption modulators (EAMs):
- in fiber optic systems, **13.55–13.60, 13.56f, 13.57f**
 - in OTDM networks, **20.18, 20.20, 20.20f**
- Electroluminescence, of LEDs, **13.37**
- Electron beam steering, **54.11**
- Electron beam energies, **36.3t–36.6t**
- Electron excitation, WDXRF and, **29.2**
- Electron linacs, **63.12**
- Electron-beam lithography (EBL), **40.8**
- Electrons, **2.9, 54.12, 55.17, 56.2, 58.1**
- Electro-optic effect, **7.6–7.16**
- and eigenpolarization/phase velocity indices of refraction, **7.13–7.16, 7.14f, 7.16f**
 - Jacobi method, **7.11–7.13**
 - linear, **7.7, 7.8t, 7.10**
 - and lithium niobate modulators, **13.49–13.51, 13.51f**
 - quadratic (Kerr), **7.9t–7.10t, 7.11**
- Electro-optic modulators, **7.1–7.39**
- applications for, **7.36–7.39, 7.37f–7.38f**
 - bulk modulators, **7.16–7.28**
 - amplitude modulation, **7.22–7.24, 7.23f, 7.24f**
 - frequency modulation, **7.24–7.25, 7.25f**
 - phase modulation, **7.18–7.20**
 - polarization modulation (dynamic retardation), **7.20–7.22, 7.20f, 7.21f**
 - scanners, **7.26–7.28, 7.26f–7.28f**
 - crystal optics and the index ellipsoid, **7.3–7.7, 7.4f–7.6f, 7.8f–7.10f**
 - and electro-optic effect, **7.6–7.16**
 - eigenpolarization/phase velocity indices of refraction, **7.13–7.16, 7.14f, 7.16f**
 - Jacobi method, **7.11–7.13**
 - linear, **7.7, 7.8t, 7.10**
 - quadratic (Kerr), **7.9t–7.10t, 7.11**
 - and Euler angles, **7.39**
 - in fiber optic systems, **13.61**
 - geometries, **7.16–7.18, 7.17f**
 - light propagation in, **7.3**
 - longitudinal, **7.16, 7.17, 7.17f**
 - materials for, **7.33–7.34**
 - in OTDM networks, **20.18–20.19, 20.19f**
 - performance criteria for, **7.34–7.36**
 - polymer modulators, **13.55**
 - transverse modulators, **7.16, 7.17, 7.17f**
 - traveling wave modulators, **7.28–7.30, 7.29f**
 - waveguide modulators, **7.30–7.32, 7.31f–7.33f**
- Electro-optic sampling, **7.36–7.37, 7.37f–7.38f**
- Electrorefraction, **13.2**
- Electrorefractive modulators, **13.61–13.62**
- Elements (chemical), x-ray properties of, **36.1–36.9, 36.3t–36.9t**
- Elliptical reflectors, **64.3, 64.4, 64.4f**
- Emission:
- amplified spontaneous, **9.13, 10.11, 14.3, 14.6**
 - cold-cathode field, **54.10–54.11**
 - molecular, **3.18, 3.20, 3.20f**
 - particle-induced x-ray, **29.4**
 - self-amplified spontaneous, **58.1**
 - spontaneous, **2.13, 19.3, 20.1**
 - stimulated, **20.2**

- Emission lines, K-shell and L-shell, **36.3t–36.8t**
- Endlessly single-mode photonic crystal fiber (ESM-PCF), **11.12, 11.13, 11.21, 11.21f**
- End-pumped schemes, for fiber lasers, **25.9–25.10, 25.28**
- Energy(-ies):
- atomic, **2.2–2.5, 2.4f**
 - Auger, **36.3t, 36.9t**
 - band gap, **19.2**
 - conservation of, **6.9**
 - electron binding, **36.3t–36.6t**
 - equipartition, **22.6**
 - Fermi level, **30.1**
 - filtering of, **53.10**
 - photon, **36.7t–36.8t**
(*See also specific sources, e.g.: Fiber lasers*)
- Energy dispersive detectors (EDS), **62.2–62.4**
- Energy-dispersive x-ray fluorescence (EDXRF), **29.3–29.11**
- with doubly curved crystal diffraction, **29.6–29.7, 29.8f–29.9f**
 - monocapillary micro-XRF, **29.4**
 - polycapillary micro-XRF, **29.4–29.6, 29.5f, 29.6f**
 - ultrahigh resolution, **29.9–29.11, 29.9f–29.11f**
- Entrance slits, of gratings and monochromators, **38.7**
- Environmental control and correction, for adaptive optics, **50.5, 50.6**
- Epitaxial growth, **13.2**
- Epithermal neutrons, **63.18**
- Equatorial divergence, x-ray diffraction and, **28.1**
- Equipartition energy, **22.6**
- Erbium-doped fiber amplifiers (EDFAs):
- energy levels, **14.4**
 - fast power transients, **21.39–21.41, 21.39f, 21.40f**
 - gain flattening, **14.6–14.7, 21.38–21.39, 21.39f**
 - gain formation, **14.4–14.5, 14.5f**
 - gain peaking, **21.38, 21.38f**
 - noise, **14.6**
 - pump wavelength options, **14.5–14.6**
 - semiconductor amplifiers vs., **9.13, 9.14, 14.1, 14.2t**
 - static gain dynamic and channel power equalization, **21.41, 21.41f–21.42f**
 - in WDM networks, **21.2–21.3, 21.2f**
- Erbium-doped fibers, **25.23t, 25.24, 25.32, 25.33**
- Erbium-doped yttrium aluminum garnet (Er:YAG) lasers, **12.3t, 12.6, 12.13**
- Erbium/ytterbium-doped fiber amplifiers (EYDFAs), **14.2, 14.2t, 14.7–14.8**
- ESCON (Enterprise System Connection) standard, **23.1–23.2, 23.2f**
- ESO telescope, **5.35**
- Etch and regrowth fabrication, of electroabsorption modulators, **13.59–13.60**
- Ethernet standard, **23.7**
- Euler angles, **7.13, 7.39**
- Euler's constant, **56.9**
- European Synchrotron Radiation Facility (ESRF), **37.5, 37.8, 50.5, 50.6**
- European X-Ray Free-Electron Laser (XFEL), **58.1**
- Evanescent coupling, **25.11**
- Evanescent wave spectroscopy (EWS), **12.13**
- Evanescent-wave coupled pin waveguides, **13.68**
- Excitation, of electrons, **56.2**
- Excitation spectroscopy, **2.15f, 2.21**
- Excited state absorption (ESA), **14.6**
- Exit slits, of gratings and monochromators, **38.7**
- Expansion coefficients, in atmospheric optics, **4.20–4.22, 4.21t, 4.22f–4.23f**
- Extended Gordon-Haus effect, **22.11**
- Extended x-ray absorption fine structure (EXAFS), **30.2, 30.4**
- External circuits, noise from, **13.70**
- External mirrors, **13.32–13.33, 13.33f**
- External modulation, in OTDM networks, **20.18–20.20, 20.19f, 20.20f**
- External optical feedback, from laser diodes, **13.21–13.24**
- External quantum efficiency, **13.9**
- External slope efficiency, **13.9**
- Extinction ratio, **7.35, 15.13**
- Extreme ultraviolet explorer (EUV), **44.4, 44.5**
- Extreme ultraviolet (EUV) lasers, **58.2–58.4, 58.3f**
- Extreme ultraviolet lithography (EUV-L), **34.1–34.6**
- and EUV-interferometric lithography, **34.4–34.5, 34.5f**
 - limitations of, **34.5–34.6, 34.5f, 34.6f**
 - and multifoil optics, **48.1**
 - and multilayers, **41.5, 41.7, 41.8**
 - in semiconductor industry, **34.1–34.2, 34.2t**
 - technology for, **34.2–34.5, 34.3f, 34.4f**

- Extreme ultraviolet (EUV) region,
Schwarzschild objective for, 51.3
- Extreme ultraviolet-interferometric lithography
(EUV-IL), 34.4–34.5, 34.5f
- Extrinsic Fabry-Perot interferometric (EFPI)
sensors, 24.2–24.4, 24.2f, 24.3f
- Extrusion, of fiber lasers, 25.27
- Eye degradations, of optical fiber receivers, 9.11
- Eye openings, optical fiber receivers and, 9.11
- Eye safety, 15.1
- Fabry-Perot cavities, 13.65, 19.19, 20.21, 20.21f,
24.2
- Fabry-Perot filters, 13.33, 22.8
- Fabry-Perot interferometers, 3.36, 24.2–24.5,
24.2f–24.4f
- Fabry-Perot lasers, 9.11, 13.12–13.13, 13.29,
15.11, 15.14
- Fabry-Perot semiconductor lasers, 20.13–20.14,
20.14f
- Fankuchen-cut crystals, 63.25
- Far Bragg diffraction acousto-optic interaction
and, 6.8
- Far field, 4.10
- Faraday rotators, 18.7, 18.7f, 18.10
- Far-field diffraction, 27.2–27.3, 27.3f
- FASCODE program, 3.23, 3.24f
- Fast power transients, for EDFAs, 21.39–21.41,
21.45f
- FDDI (fiber distributed data interface) standard,
13.41, 23.2–23.3, 23.3f
- Feedback:
in fiber optic systems, 13.30–13.32
from laser diodes, 13.21–13.24
from SOAs, 19.8, 19.8f
- Fermi choppers, 63.14
- Fermi level, of energy, 30.1
- Fermi pseudopotential, 63.4, 63.5
- Ferroelectric smectic phase, of liquid crystals,
8.12, 8.12f
- Fiber (material), for optical fibers, 9.4–9.7
- Fiber amplifiers, 14.1–14.11
categories and features of, 14.1–14.2, 14.2t
erbium-doped
energy levels, 14.4
fast power transients, 21.39–21.41, 21.39f,
21.40f
gain flattening, 14.6–14.7, 21.38–21.39,
21.39f
gain formation, 14.4–14.5, 14.5f
- Fiber amplifiers, erbium-doped (*Cont.*):
gain peaking, 21.38, 21.38f
noise, 14.6
pump wavelength options, 14.5–14.6
semiconductor amplifiers vs., 9.13, 9.14,
14.1, 14.2t
static gain dynamic and channel power
equalization, 21.41, 21.41f–21.42f
in WDM networks, 21.2–21.3, 21.2f
erbium/ytterbium-doped, 14.2, 14.2t, 14.7–14.8
infrared fibers for, 12.3t
parametric, 14.10–14.11
praseodymium-doped fiber amplifiers
(PDFAs), 14.7
Raman, 14.8–14.9, 14.8f, 14.9, 14.10f
rare-earth-doped, 14.2–14.4, 14.3f
semiconductors vs., 9.13–9.14
ytterbium-doped, 14.7
- Fiber attenuation, optical power loss and,
21.13–21.14
- Fiber bandwidth, of WDM networks,
21.2–21.3, 21.2f
- Fiber Bragg gratings (FBGs), 17.1–17.9
applications, 17.8–17.9, 17.8f
chirped, 21.22–21.23, 21.23f, 21.25–21.26,
21.25f
fabrication, 17.4–17.8, 17.5f–17.7f
and fiber lasers, 25.8, 25.16, 25.18, 25.30, 25.31
long-period gratings vs., 24.9, 24.11
photosensitivity of, 17.2–17.3
properties of, 17.3–17.4, 17.4f
sensors based on, 24.5–24.8, 24.6f–24.7f
- Fiber lasers, 25.1–25.33
architectures, 25.9–25.18, 25.19f
all-fiber monolithic systems, 25.16–25.18,
25.16f
free space, 25.13–25.15, 25.14f, 25.15f
pumping techniques, 25.9–25.13, 25.11f,
25.12f
bulk lasers vs., 25.5–25.6
dopants for, 25.22–25.26, 25.23t
fabrication of, 25.26–25.29, 25.27t
growth of, 25.5, 25.5f
history of, 25.3–25.4, 25.4f
infrared fibers for, 12.3t
limitations of, 25.6–25.7
LMA designs for, 25.18–25.22, 25.19f
operation of, 25.7–25.8
spectral and temporal modalities of,
25.29–25.33

- Fiber length, for rare-earth-doped amplifiers, 14.2–14.3, 14.3f
 - Fiber optic amplifiers, 14.1–14.11
 - categories and features of, 14.1–14.2, 14.2t
 - erbium-doped
 - energy levels, 14.4
 - fast power transients, 21.39–21.41, 21.39f, 21.40f
 - gain flattening, 14.6–14.7, 21.38–21.39, 21.39f
 - gain formation, 14.4–14.5, 14.5f
 - gain peaking, 21.38, 21.38f
 - noise, 14.6
 - pump wavelength options, 14.5–14.6
 - semiconductor amplifiers vs., 9.13, 9.14, 14.1, 14.2t
 - static gain dynamic and channel power equalization, 21.41, 21.41f–21.42f
 - in WDM networks, 21.2–21.3, 21.2f
 - erbium/ytterbium-doped, 14.7–14.8
 - parametric, 14.10–14.11
 - praseodymium-doped, 14.7
 - Raman fiber, 14.8–14.9, 14.8f, 14.10f
 - rare-earth-doped, 14.2–14.4, 14.3f
 - ytterbium-doped, 14.7
 - Fiber optic chemical sensors, 12.3t
 - Fiber optic communication links, 15.1–15.17
 - distortions and dynamics range of, 15.5–15.6, 15.6f
 - figures of merit for, 15.2–15.6, 15.3f, 15.4f
 - link budget analysis for, 15.6–15.17
 - extinction ratio, 15.13
 - installation loss, 15.6–15.7, 15.8t
 - optical power penalties, 15.8–15.17, 15.10f
 - Fiber optic communication standards, 23.1–23.8
 - ATM/SONET, 23.6
 - ESCON, 23.1–23.2, 23.2f
 - Ethernet, 23.7
 - FDDI, 23.2–23.3, 23.3f
 - Fibre Channel standard, 23.4, 23.5f, 23.5t
 - InfiniBand, 23.8, 23.8t
 - Fiber optic networking, micro-optics-based
 - components for, 18.1–18.12
 - attenuators, 18.2, 18.9
 - beam splitters, 18.6, 18.6f
 - circulators, 18.3, 18.3f, 18.10
 - directional couplers, 18.2, 18.3, 18.3f, 18.9, 18.9f
 - Faraday rotators, 18.7, 18.7f
 - filters, 18.6
 - Fiber optic networking, micro-optics-based components for (*Cont.*):
 - gratings, 18.5–18.6, 18.6f
 - GRIN-rod lenses, 18.7, 18.8, 18.8f
 - isolators, 18.3, 18.10, 18.10f
 - mechanical switches, 18.4, 18.5, 18.5f, 18.11, 18.11f, 18.12f
 - MEMS mirrors and switches, 18.8, 18.8f, 18.11, 18.12f
 - multiplexers/demultiplexers/duplexers, 18.4, 18.4f, 18.10–18.11
 - network functions, 18.2–18.5
 - polarizers, 18.7, 18.7f
 - power splitters, 18.2–18.3, 18.2f, 18.9, 18.9f
 - prisms, 18.5, 18.5f
 - Fiber optic networks and systems, 9.14–9.15
 - detectors in, 13.2–13.3, 13.63–13.73
 - avalanche photodiodes, 13.71–13.73
 - MSM detectors, 13.73
 - pin diodes, 13.63–13.71, 13.64f, 13.66f, 13.66t, 13.68f
 - Schottky photodiodes, 13.73
 - modulators in, 13.2, 13.48–13.63
 - electroabsorption, 13.55–13.60, 13.56f, 13.57f
 - electro-optic, 13.61
 - electrorefractive, 13.61–13.62
 - lithium niobate, 13.48–13.55, 13.49f, 13.51f, 13.54f
 - semiconductor interferometric, 13.63
 - sources for, 13.1–13.48
 - distributed Bragg reflector lasers, 13.28–13.29, 13.29f
 - distributed feedback lasers, 13.30–13.32, 13.30f, 13.31f
 - laser diodes, 13.3–13.24, 13.3f, 13.4f, 13.6f, 13.7f, 13.10f, 13.14f, 13.16f, 13.17f, 13.19f–13.21f
 - light-emitting diodes, 13.36–13.42, 13.38f, 13.40f
 - quantum well lasers, 13.24–13.28, 13.25f–13.27f
 - strained layer quantum well lasers, 13.26–13.28, 13.26f, 13.27f
 - tunable lasers, 13.32–13.36, 13.33f–13.36f
 - vertical cavity surface-emitting lasers, 13.42–13.48, 13.43f, 13.45f
- [See also related topics, e.g.: Optical time-division multiplexed (OTDM)]
- communication networks

- Fiber optic sensors, **24.1–24.13**
 extrinsic Fabry-Perot interferometric,
24.2–24.4, 24.2f, 24.3f
 fiber Bragg grating, **24.5–24.8, 24.6f–24.7f**
 intrinsic Fabry-Perot interferometric sensors,
24.4–24.5, 24.4f
 long-period grating sensors, **24.8–24.13,**
24.9f–24.12f, 24.11t, 24.13t
- Fiber pigtail connection, **13.8**
- Fiber pulling, for fiber lasers, **25.26**
- Fiber Raman lasers, **10.7**
- Fiber-based couplers, **16.1–16.6, 16.2f, 16.5f**
- Fibre Channel Arbitrated Loop (FC-AL), **23.4**
- Fibre Channel standard, **23.2, 23.4, 23.5f, 23.5t**
- Field of view (FOV):
 and multifoil optics, **48.1**
 for scatterometers, **1.6, 1.10, 1.12**
- Field-effect transistor (FET) amplifiers, **13.70**
- Field-weighted-average resolution, **44.10**
- Fifth-order oblique spherical aberration, **45.4**
- Figures of merit:
 acousto-optic, **6.16–6.17**
 fiber optic communication links, **15.2–15.6,**
15.3f, 15.4f
- Filling factor, electro-optic effect and, **13.50**
- Film x-ray detectors, **60.8, 60.9t, 60.10t**
- Filtered backprojection (FBP), **32.2**
- Filters and filtering:
 Bragg, **63.14**
 Fabry-Perot, **13.33, 22.8**
 in fiber optic networking, **18.6**
 frequency-guiding, **22.7–22.9**
 guided-mode resonance, **25.2, 25.30**
 Mach-Zehnder, **14.6, 21.39**
 for networking, **18.6**
 neutron, **63.18–63.19, 63.18t, 63.19f**
 of x-ray tube source spectra, **54.9**
 [See also Acousto-optic tunable filters
 (AOTFs)]
- Finite-difference time-domain (FDTD)
 analysis, for PCFs, **11.7**
- Fitting error, **5.41**
- Fixed dispersion compensation, **21.22–21.23,**
21.23f
- Fizeau techniques, for surface figure metrology,
46.3
- Flame brushing, of fiber Bragg gratings, **17.2**
- FLASH free-electron laser facility, **58.1**
- Flat panel detectors, for x-ray imaging,
61.3–61.7, 61.4f, 61.6t, 61.7f
- Fluorescence:
 laser-induced, **3.21**
 x-ray, **28.1, 54.8, 62.5, 62.6f**
 and polycapillary x-ray optics,
53.10–53.11, 53.11f
 and x-ray diffraction, **28.5–28.6, 28.6f**
- Fluorescence line narrowing (FLN), **2.13–2.14,**
2.14f
- Fluorescent lamps, cold cathode, **8.30, 8.30f**
- Fluorides, for fiber lasers, **25.28–25.29**
- Fluoroaluminate glass, **12.4, 12.4t**
- Fluorozirconate glass (ZBLAN), **12.5, 12.5f**
 and fiber lasers, **25.3, 25.24, 25.27t, 25.28**
 fluoroaluminate glass vs., **12.4, 12.4t**
- Flux, of electron bunches, **55.17**
- Focal surfaces, of grazing incidence optics,
45.5
- Focus anisoplanatism, **5.27–5.29, 5.28f–5.30f,**
5.42–5.43
- Focused beams, **54.16**
- Focused single crystal diffraction, **53.12–53.14,**
53.12f–53.13f
- Focusing:
 of gratings and monochromators, **38.3–38.6,**
38.3f, 38.4t–38.5t
 in grazing-incidence neutron optics,
64.3–64.7, 64.4f–64.7f
 Kerr, **7.39**
 in neutron and x-ray optics, **26.9–26.11, 63.22**
 with refractive x-ray lenses, **37.8–37.11,**
37.9f, 37.10f
 (See also Hard x-rays, nanofocusing of)
- Focusing polycapillary x-ray optics, **53.9–53.10,**
53.10t, 53.14
- Forward Brillouin scattering, **11.25, 11.26**
- Forward error correction (FEC) coding, **19.27,**
21.26
- Fourier crosstalk matrix, **32.3**
- Fourier differentiation, **46.11–46.12**
- Fourier intervals, **27.4**
- Fourier transform spectroscopy, **2.5, 2.6f**
- Fourier transforming infrared spectrometer
 (FTIR), **1.14**
- Fourier transforms, **6.15, 16.4, 46.8, 55.9, 55.10**
- Four-wave mixing (FWM):
 in optical fibers, **10.2, 10.9–10.11, 10.11f**
 and SOAs, **19.13, 19.14f, 19.27, 19.33–19.35,**
19.33f, 19.34f
 and solitons, **22.11, 22.15**
 in WDM networks, **21.19–21.20, 21.20f**

- Frames, in OTDM, **20.7**
- Frank elastic constants, **8.22**
- Franz-Keldysh effect, **13.59**
- Fraunhofer diffraction, **42.2, 63.25**
- Fraunhofer regime, in coherent x-ray optics, **27.2**
- Free carriers, in fiber optic systems, **13.4**
- Free electron lasers (FELs), **29.4, 41.9–41.10, 41.10f, 48.1, 58.1–58.2**
- Free space fiber lasers, **25.9, 25.10, 25.13–25.15, 25.14f, 25.15f**
- Freedericksz threshold voltage, of LC cells, **8.27**
- Frequency:
- and acousto-optic interaction, **6.12–6.14, 6.13f, 6.14f, 6.16**
 - crossover, **8.16**
 - of electro-optic modulators, **7.35**
 - Greenwood, **5.19, 5.22, 5.42**
 - of liquid crystals, **8.17–8.18, 8.18f**
 - Nyquist, **27.4**
 - Stokes, **10.8, 21.42**
 - Tyler, **5.17**
- Frequency chirping, **13.1, 13.17–13.18, 13.17f**
- Frequency modulation, **7.24–7.25, 7.25f**
- Frequency shifters, in electro-optic phase shifts, **13.51**
- Frequency-division multiplexing (FDM), **9.16**
- Frequency-guiding filters, **22.7–22.9**
- Frequency-modulated (FM) systems, **9.16**
- Fresnel diffraction, **27.2, 27.4, 40.9, 63.25**
- Fresnel equations, **41.3, 64.2**
- Fresnel integrals, **6.15**
- Fresnel phase zones, **27.3**
- Fresnel reflection, **13.6, 13.53, 17.3, 25.8**
- Fresnel rhomb, **43.5–43.6, 43.5f**
- Fresnel waves, **27.2**
- Fresnel zone plates, **40.2, 42.3, 42.3f, 55.16**
- Fresnel zones, MLLs and, **42.16f**
- Fresnel's law, **63.20**
- Fresnel-Soret zone plates, **40.2**
- Fried's coherence diameter, **5.2, 5.7, 5.9–5.10**
- Fried's coherence length, **4.8**
- Fringes, in sensor signals, **24.3**
- Fringes of equal chromatic order (FECO) technique, **46.2**
- Full width at half-maximum (FWHM), of fiber Bragg gratings, **17.6–17.7, 17.7f, 24.8**
- Full-aperture techniques, in surface figure metrology, **46.3**
- Fusion splicing, **25.18**
- Gabor zone plates, **40.4, 40.5, 40.7**
- Gain:
- of avalanche photodiodes, **13.71–13.72**
 - of EDFAs, **14.4–14.5, 14.5f**
 - polarization-dependent, **14.9, 19.18–19.20, 21.18**
 - Raman, **10.5, 10.6, 21.42f**
 - of SOAs, **19.4–19.6, 19.4f–19.6f**
- Gain clamping, **19.14–19.15, 19.14f**
- Gain dynamics, of SOAs, **19.12–19.13, 19.13f, 19.14f**
- Gain efficiency, of EDFAs, **14.6**
- Gain flattening, **14.6–14.7, 21.38–21.39, 21.39f**
- Gain peaking, **21.38, 21.38f**
- Gain per unit length, of lasers, **13.7, 13.8**
- Gain ripple, **19.8, 19.8f, 19.19, 19.19f**
- Gain saturation, **13.15, 13.17**
- Gain-guided index antiguided fibers (GG IAG), **25.2, 25.19f, 25.22**
- Gain-guided laser diodes, **13.5**
- Gallium phosphide (GaP), **6.16, 6.17t, 6.29t, 6.30, 6.31t, 6.34t**
- Gamma cameras, **32.2**
- Gamma rays, in SPECT imaging, **32.1–32.2**
- Gas detectors, **63.32–63.33**
- Gas-puff sources, **57.3, 57.3t**
- Gaunt factors, **56.9–56.10**
- Gaussian apertures, **37.6**
- Gaussian approximation of sensitivity, for optical fiber receivers, **9.10**
- Gaussian error function, **15.3**
- Gaussian image point, **38.5–38.6**
- Gaussian integral, for noise, **15.2, 15.4**
- Gaussian line profiles, **3.14, 56.5–56.7**
- Gaussian spectra, **15.13**
- Gaussian statistics, for wave propagation, **5.9**
- Gaussian transmission, **37.5**
- Gauss-Seidel iterative method, **3.21**
- Geiger counters, **60.5, 60.9t**
- Geiger region, **60.5**
- Gemini North telescope, **5.20, 5.21f**
- GENLN2 (code), **3.23**
- Geometrical point spread function (GPSF), **44.13**
- Germanate:
- in fiber lasers, **25.27t, 25.28**
 - in optical fibers, **12.3–12.4, 12.6, 12.6f**
- Glancing angle, of crystal monochromators, **39.1, 39.2**

- Glass:
- for fiber lasers, 25.27*t*, 25.28–25.29
 - fluorozirconate (ZBLAN), 12.5, 12.5*f*
 - and fiber lasers, 25.3, 25.24, 25.27*t*, 25.28
 - fluoraluminate glass vs., 12.3*t*, 12.4, 12.4*t*
 - fluoroaluminate, 12.4, 12.4*t*
 - heavy-metal fluoride, 12.1–12.5, 12.2*f*, 12.3*t*, 12.4*t*, 12.5*f*
 - heavy-metal oxide, 12.2*t*, 12.3–12.7, 12.3*t*, 12.4*t*, 12.5*f*–12.7*f*
 - for hollow waveguides, 12.2*f*, 12.11–12.13, 12.12*f*
 - negative core-cladding index difference, 11.14, 11.15
 - Raman bands of, 11.24
 - Rayleigh scattering in, 11.21
 - Zerodur, 47.5
- Glass micro-pore optics, 49.1–49.6, 49.2*f*–49.6*f*
- Global climate change, 3.43–3.45, 3.44*f*
- Goebel mirrors, 26.10
- GOES-13 satellite, 44.16–44.17, 44.17*f*
- Gooch-Tarry first minimum condition, 8.26
- Goos-Hanschen shifts, 13.55
- Gordon-Haus effect, 22.7–22.8, 22.11
- Graded index (GRIN) fibers, 15.17
- Graded index separate confinement heterostructure (GRINSCH), 13.4*f*, 13.5
- Graded index-rod (GRIN-rod) lenses, 18.7, 18.8, 18.8*f*, 18.10, 18.11*f*
- Gradient tilt (G-tilt), 4.3
 - and adaptive optics, 5.14–5.16
 - and angle of arrival, 4.23, 4.25, 4.26
- Gradients, of wavefronts, 5.23
- Graphical user interface (GUI), for SHADOW code, 35.3
- Grasshopper* monochromator, 38.3
- Grating equation, for VUV and soft x-ray region, 38.2–38.3
- Gratings:
- Bragg, 20.14, 25.29, 25.30
 - circular, 40.1
 - fiber Bragg, 17.1–17.9
 - applications, 17.8–17.9, 17.8*f*
 - chirped, 21.22–21.23, 21.23*f*, 21.25–21.26, 21.25*f*
 - fabrication, 17.4–17.8, 17.5*f*–17.7*f*
 - and fiber lasers, 25.8, 25.16, 25.18, 25.30, 25.31
 - long-period gratings vs., 24.9, 24.11
 - photosensitivity, 17.2–17.3
 - Gratings, fiber Bragg (*Cont.*):
 - properties of, 17.3–17.4, 17.4*f*
 - sensors based on, 24.5–24.8, 24.6*f*–24.7*f*
 - Fresnel diffraction of, 40.9
 - Hill, 17.5
 - long-period, 21.39, 21.39*f*, 24.10–24.11, 24.11*t*, 24.13*t*
 - for networking, 18.5–18.6, 18.6*f*
 - sampled, 13.34, 13.34*f*
 - short-period, 24.6
 - superstructure, 13.34, 13.35*f*
 - in VUV and soft x-ray region, 38.1–38.8
 - diffraction properties, 38.1–38.3, 38.2*f*
 - dispersion properties, 38.6–38.7
 - efficiency, 38.8
 - focusing properties, 38.3–38.6, 38.3*f*, 38.4*t*–38.5*t*
 - resolution properties, 38.7
- Grazing incidence optics:
- aberrations of, 44.6–44.12, 44.7*f*, 44.9*f*–44.11*f*, 45.1–45.8, 45.2*f*
 - image formation with, 44.3–44.18
 - and multifoil optics, 48.1–48.2
 - and pumping by EUV lasers, 58.3, 58.4*f*
 - telescopes with, 44.6–44.12, 44.7*f*, 44.9*f*–44.11*f*
 - and x-ray mirrors, 44.3–44.6, 44.4*f*–44.6*f*
- Grazing incidence x-ray optics, 44.12–44.18, 44.12*f*–44.18*f*
- Grazing-angle reflection, 63.21
- Grazing-incidence neutron optics, 64.1–64.7
 - diffractive scattering and mirror surface roughness, 64.2–64.3
 - and imaging focusing optics, 64.3–64.7, 64.4*f*–64.7*f*
 - total external reflection, 64.1–64.2
- Green flashes, 3.43
- Green's function solution, 6.14
- Greenwood frequency, 5.19, 5.22, 5.42
- Group velocity dispersion (GVD), 11.18, 11.18*f*, 11.19*f*, 15.11
- Grüneisen constant, 6.17
- GTWave technology, 25.11
- Guidance, in photonic crystal fibers, 11.11–11.26
 - attenuation mechanisms, 11.19–11.22, 11.20*f*, 11.21*f*
 - birefringence, 11.17
 - Brillouin scattering, 11.25–11.26, 11.25*f*

- Guidance, in photonic crystal fibers (*Cont.*):
 group velocity dispersion, 11.18, 11.18f, 11.19f
 Kerr nonlinearities, 11.22–11.24, 11.24f
 negative core-cladding index difference, 11.14–11.17, 11.14f–11.17f
 positive core-cladding index difference, 11.12–11.14, 11.12f, 11.13f
 Raman scattering, 11.24
 resonance and antiresonance, 11.12
- Guided-mode resonance filters (GMRFs), 25.2, 25.30
- Guides, neutron, 63.15–63.18
- Half-period zones, of zone plates, 40.1
- Half-wave voltage, of electro-optic modulators, 7.19, 7.22
- Hard x-ray beamlines, SHADOW code for, 35.5, 35.5f
- Hard x-ray optics, astronomical, 47.9–47.10
- Hard x-ray telescopes, 50.2
- Hard x-rays, nanofocusing of, 42.1–42.17
 history of, 42.2–42.4, 42.2f, 42.3f
 instrumental beamline arrangement and measurements for, 42.9–42.10, 42.9f–42.12f
 limitations of, 42.15–42.17, 42.16f–42.17f
 with magnetron-sputtered MLLs, 42.5–42.7, 42.6f–42.8f
 on MLLs with curved interfaces, 42.14, 42.15f
 Takagi-Taupin calculations for, 42.12–42.14
 volume diffraction calculations for, 42.4–42.5, 42.5f
 with wedged MLLs, 42.12–42.13, 42.13f, 42.14f
- Hardware implementation, for adaptive optics, 5.21–5.38
 higher-order wavefront sensing techniques, 5.36–5.37
 laser beacons, 5.27–5.34, 5.28f–5.31f, 5.33f
 real-time processors, 5.34–5.35, 5.34f, 5.35f
 Shack-Hartmann technique, 5.23–5.27, 5.23f, 5.25f, 5.26f
 tracking, 5.21–5.23, 5.22f
 wavefront correctors, 5.37–5.38, 5.38f
- Heavy water, scattering in, 63.10
- Heavy-hole (HH) bands, of strained layer quantum well lasers, 13.27
- Heavy-metal fluoride glass (HMFG) fibers, 12.1–12.5, 12.2f, 12.3t, 12.4t, 12.5f
- Heavy-metal oxide glass fibers, 12.2t, 12.3–12.7, 12.3t, 12.4t, 12.5f–12.7f
- HEFT balloon payloads, 47.10
- Height profilometry, 46.3
- Helium atoms, 2.3
- Helium [$^3\text{He}(n, p)$] reaction, 63.31
- Helmholtz equation, 5.8
- Hermetic enclosure, for laser diodes, 13.7, 13.7f
- Hermite-Gaussian functions, 11.7
- HERO balloon payload, 47.10
- Heterodyne detection, 3.34
- Heterostructures, of fiber optic devices, 13.2
- High aspect ratio microlithography (HARM), 61.3
- High harmonic production, of x-ray lasers, 58.2
- High Resolution Doppler Imager (HRDI), 3.36, 3.37f
- Higher-order mode (HOM) fibers, 25.2, 25.19f, 25.22
- High-power spectrally controlled fiber lasers, 25.29
- High-power ultrashort pulse technologies, for fiber lasers, 25.32–25.33
- High-power USP oscillators, 25.32–25.33
- High-reflectivity mirrors, 41.7–41.8, 41.8f
- High-resolution (HR) acousto-optic deflectors, 6.29, 6.29t
- Hill gratings, 17.5
- HITRAN database, 3.14, 3.22–3.23, 3.22f, 3.23f
- HITRAN-PC program, 3.26
- Hobby-Eberly telescopes, 5.2
- Holeburning:
 optical, 2.13, 2.14f
 spatial, 20.2
- Hollow glass waveguides (HGWs), 12.2f, 12.11–12.13, 12.12f
- Hollow waveguides, 12.2f, 12.2t, 12.3t, 12.10–12.13, 12.12f
- Hollow-core photonic crystal fibers:
 attenuation in, 11.20–11.22, 11.20f
 birefringence of, 11.17
 and group velocity dispersion, 11.18, 11.19f
 Kerr effects for, 11.23
 and negative core-cladding index difference, 11.14–11.15, 11.14f
- Holmium-doped fibers, 25.23t, 25.25–25.26
- Homogeneous broadening, 20.1

- Hot bands, in molecular spectroscopy, 2.5
 Hot filament sources, of x-ray tubes, 54.10
 HRCam system, 5.23
 Huang-Rhys factor, 2.15
 Huesler alloy, 63.28–63.29
 Hufnagel model, of atmospheric turbulence, 3.29
 Hufnagel-Valley model, of atmospheric turbulence, 3.30, 5.7, 5.8
 Huygens-Fresnel approximation, 3.31–3.33
 Hybrid mode locking, 20.17
 Hybrid network topologies, for WDM networks, 21.7
 Hydrogen (H^-) ions, negative, 2.3
 Hydrogen loading, 17.2
 Hyperboloid-hyperboloid (HH) grazing incidence x-ray telescopes, 44.10–44.12, 44.11f
- IBM, 23.1, 23.2
 Ice, in standard atmosphere, 3.6, 3.42, 3.42f
 Ideal receivers, 9.9
 Ideality factor, of semiconductor diodes, 13.69
 IEEE (Institute of Electrical and Electronics Engineers) standards, 23.2, 23.7
 Image (scophony) AO modulators, 6.34–6.35
 Image receptors, in medical imaging, 31.4, 31.4f
 Image wander, 4.3 (*See also* Angle of arrival)
 Images:
 long-exposure, 4.3–4.7
 from Nomarski microscope, 46.2
 short-exposure, 4.3, 4.31–4.35, 4.31f–4.34f, 4.35t
 x-ray, 31.1–31.4
 Imaging:
 and atmospheric turbulence, 3.34
 diffractive, 41.9–41.10, 41.10f
 with grazing incidence optics, 44.3–44.18
 with grazing-incidence neutron optics, 64.3–64.7, 64.4f–64.7f
 medical, 31.1–31.10
 applications of, 31.9, 31.10
 digital displays, 31.8–31.9, 31.8f–31.10f
 digital tomosynthesis, 31.7–31.8
 and inverse Compton x-ray sources, 59.3–59.4
 and polycapillary x-ray optics, 53.14–53.16, 53.15f–53.16f
 radiography, 31.1–31.4, 31.2f–31.4f
 tomography, 31.1, 31.5–31.7, 31.5f–31.7f
 x-ray detectors for, 61.2
 Imaging (*Cont.*):
 molecular, 32.1
 monochromatic, 53.16–53.17, 53.17f
 multi-energy, 54.9–54.10
 nuclear, 53.17, 53.18, 53.18f
 with refractive x-ray lenses, 37.6–37.7, 37.6f, 37.7f
 scatter rejection in, 53.14–53.16, 53.15f–53.16f
 SPECT, 32.1–32.3
 and spectroscopy, 5.19–5.21, 5.19f, 5.20f
 thermal, 12.3t
 through atmospheric turbulence, 4.1–4.37
 aberration variance and approximate Strehl ratio for, 4.27–4.28, 4.28f
 and adaptive optics, 4.35–4.36
 angle of arrival fluctuations, 4.23–4.26, 4.25f, 4.26f
 and covariance and variance of expansion coefficients, 4.20–4.22, 4.21t, 4.22f–4.23f
 Kolmogorov turbulence and atmospheric coherence length, 4.7–4.10, 4.8f, 4.9f
 long-exposure images, 4.3–4.7
 modal correction of turbulence, 4.28–4.30, 4.29t, 4.30f
 and modal expansion of aberration function, 4.17–4.20, 4.17f–4.18f, 4.19t, 4.20t
 and resolution of telescopes, 4.2–4.3
 short-exposure image, 4.31–4.35, 4.31f–4.34f, 4.35t
 and systems with annular pupils, 4.10–4.16, 4.11f–4.15f, 4.15t
 Imaging detectors (x-ray), 61.1–61.8
 CCD detectors, 61.7–61.8, 61.8f
 flat panel detectors, 61.3–61.7, 61.4f, 61.6t, 61.7f
 geometries for and classifications of, 61.1–61.3, 61.2f
 Imaging plates, in neutron optics, 63.34
 Impact ionization coefficient, 13.72
 Incidence angle, constant, 38.2
 Incident power measurement, in scatterometers, 1.14–1.15
 Incoherent scattering, 26.7, 31.2, 63.7, 63.8
 Index ellipsoid, of electro-optic modulators, 7.4–7.7, 7.5f–7.6f

- Index of refraction:
 complex, **48.1**
 and fiber Bragg gratings, **17.2**
 and Kolmogorov turbulence, **4.7**
 of liquid crystals, **8.18–8.19**
 in neutron optics, **63.19–63.20**
 phase velocity, **7.15–7.16, 7.16f**
 of photonic crystal fibers, **11.9–11.10**
 structure function of, **5.6–5.7**
- Index-guided laser stripes, **13.6**
- Indirect modulation, in OTDM networks,
20.17–20.18, 20.17f
- Indirect-conversion flat panel detectors,
61.4–61.6, 61.4f
- Inelastic optical processes, **3.21–3.22**
- Inelastic scattering, **21.20, 63.3**
- In-fiber devices, for photonic crystal fibers,
11.27, 11.28
- InfiniBand standard, **23.8, 23.8t**
- InFOCUS/SUMIT balloon payloads, **47.10**
- Infrared (IR) optical fibers, **12.1–12.13**
 applications, **12.13**
 categories and properties of, **12.1–12.3,**
12.2f, 12.2t, 12.3t
 crystalline, **12.2t, 12.3t, 12.7–12.10, 12.8f,**
12.10f
 heavy-metal oxide glass in, **12.2t–12.4t,**
12.3–12.7, 12.5f–12.7f
 in hollow waveguides, **12.2t, 12.3t,**
12.10–12.13, 12.12f
- Injection:
 electrical, **13.4, 13.5, 13.45, 13.45f**
 in LEDs, **13.37**
- Injection seeding, optical clock recovery and,
20.21–20.22
- In-line semiconductor optical amplifiers
 (SOAs), **19.24**
- Inner scale of turbulence, **4.7**
- Inorganic crystals, **7.33**
- In-plane switching (IPS) cells, **8.16, 8.26, 8.28f**
- Input saturation power, for fiber amplifiers,
14.3
- Insertion devices:
 for synchrotron radiation sources,
55.9–55.16, 55.10f, 55.12f, 55.13f
- Insertion loss, **7.36, 13.53, 18.1**
- Installation loss, for fiber optic communication
 links, **15.6–15.7, 15.8t**
- Instrument signature, of scatterometers, **1.6,**
1.11–1.13, 1.11f, 1.13t
- Instrument transfer function effects, in x-ray
 mirror metrology, **46.9–46.11**
- Instrumental line spread function, **38.7**
- Integrated Mach-Zehnder interferometers,
20.18–20.19, 20.19f
- Integrated planar lightwave circuits (iPLCs),
21.12, 21.13f
- Integrated-optic modulators, **7.3, 7.30–7.32,**
7.31f–7.33f
- Integrating x-ray detectors, **61.2**
- Intelligent Physical Protocol Enhanced Physical
 Project, **23.4**
- Intensity, of spectral lines, **56.2–56.4**
- Intensity modulation and modulators, **6.34,**
13.57–13.58
- Interaction zone, for Compton scattering, **59.1**
- Interband processes, SOAs and, **19.12, 19.13**
- Interconnected switchable networks, **23.4**
- Interference:
 multiple-Bragg-beam, **43.1**
 in neutron optics, **63.25–63.27, 63.27f**
 Nomarski differential, **46.4**
 Pendellösung, **63.26**
 in SOAs, **19.23**
 in x-ray optics, **26.8, 26.8f, 26.9f**
- Interferometers:
 Bragg reflection in, **63.26, 63.27**
 crystal, **63.26–63.27, 63.27f**
 Fabry-Perot, **3.36, 24.2–24.5, 24.2f–24.4f**
 Laue transmission in, **63.26**
 Linnik, **46.2**
 Mach-Zehnder, **21.34**
 and Bragg grating sensors, **24.8**
 and DPSK, **21.33–21.34, 21.33f**
 and electro-optic modulators, **7.22, 7.24,**
7.32, 7.38
 and fiber Bragg gratings, **17.8**
 integrated, **20.18–20.19, 20.19f**
 in OTDM networks, **20.22, 20.23f**
 SOAs in, **19.31–19.32, 19.31f, 19.33f, 19.36**
 and supercontinuum generation, **11.23**
 Michelson, **17.8–17.9, 17.8f**
 perfect crystal, **63.26–63.27, 63.27f**
 phase-measuring, **46.2**
 Sagnac, **20.22**
 ultrafast nonlinear, **19.32**
 unbalanced nonlinear, **20.22**
- Interferometric lithography (IL), **34.4**
- Interferometric Mach-Zehnder modulators,
13.51–13.52, 13.54–13.55

- Interferometric method, of FBG fabrication, 24.6, 24.6f
- Interferometric modulators, 13.51–13.52
- Intergovernmental Panel on Climate Change, 3.44
- Interleaving, in OTDM networks, 20.6–20.7, 20.7f
- Intermodal dispersion, in optical fibers, 9.5–9.6
- Intermodulation distortions (IMDs), 15.5–15.6
- Intermodulation (IM) products, of acousto-optic devices, 6.30
- Internal quantum efficiency, 13.9, 13.10
- Internal writing technique, for FBGs, 17.4–17.5
- International Space Station, 49.4
- International Telecommunications Union standards, 15.15
- Inter-switch links (ISLs), 23.4
- Intersymbol interference (ISI), 19.23
- Intraband processes, of SOAs, 19.12
- Intramodal dispersion, 9.5–9.6, 9.6f
- Intrinsic Fabry-Perot interferometric (IFPI) sensors, 24.4–24.5, 24.4f
- Inverse Compton scattering, 59.1
- Inverse Compton x-ray sources, 59.1–59.4, 59.3t
- Inverse Raman effect, 10.5, 10.6, 14.9
- Ionization, of electrons, 56.2
- Ionization chambers, 60.3–60.4, 60.9t
- Ionization x-ray detectors, 60.3–60.7, 60.9t, 60.10t
- Ionizing radiation, fiber optic communication links and, 15.17
- ISO standards, 23.2
- Isolation (directionality), of fiber-optic components, 18.1
- Isolators, for networking, 18.3, 18.10, 18.10f
- Isoplastic angle, 5.19
- Isotropic diffraction:
 - and acousto-optic interactions, 6.9–6.10, 6.10f, 6.13, 6.13f
 - by deflectors, 6.23–6.24, 6.24f
- Iterative phasing technique, for coherent diffraction microscopy, 27.4
- Jacobi method, of electro-optic effect, 7.11–7.13
- Jahn-Teller effect, 2.9
- Jitter:
 - in fiber optic communication links, 15.15–15.16
 - and solitons, 22.6–22.8, 22.11, 22.16
- Jitter transfer function (JTF), 15.16
- Johansson bent/ground focusing monochromator, 39.5
- Johnson noise, of pin diodes, 13.70
- Jülich SANS instrument, 64.4
- Kagomé lattice, 11.5f, 11.11, 11.16
- Karhunen-Loève functions, 4.36
- Keck Observatory, 5.27
- Keck telescopes, 4.36, 5.2, 5.27
- Kerr cells, 7.34
- Kerr effect, 7.11, 14.11, 19.34, 20.1, 20.22, 22.3
- Kerr electro-optic effect, 7.9t–7.10t, 7.11
- Kerr focusing, 7.39
- Kerr interactions, 10.2
- Kerr nonlinearities, 7.38–7.39, 11.22–11.24, 11.24f
- Kerr-lens mode-locking (KLM), 7.11, 7.39
- Kinematic theory of diffraction, in neutron optics, 63.5
- Kinematical theory of diffraction, for crystals, 39.2, 39.3
- Kirkpatrick-Baez (KB) mirrors, 44.4, 44.4f, 63.21, 64.5–64.6, 64.5f
- Kirkpatrick-Baez (KB) optics, 47.7–47.8, 47.8f, 47.9f, 50.5
- Kirkpatrick-Baez (KB) systems, 48.3–48.4, 48.3f
- Kolmogorov model of turbulence, 5.5–5.6, 5.11
- Kolmogorov spatial power spectral density, 5.6
- Kolmogorov spectrum, 3.28, 3.31
- Kolmogorov turbulence, 4.3, 4.7–4.10, 4.8f, 4.9f, 4.27, 4.30, 4.36
- Kossel patterns, 53.11
- Kramer approximation, for Bremsstrahlung radiation, 54.5
- Kramers-Kronig relations, 13.61, 17.3, 20.22
- Kronecker delta, 4.19
- Kronig-Penney model, 2.12
- KRS-5 fiber, 12.7–12.8
- K-shell emission lines, of elements, 36.3t–36.8t
- Kumahov capillary lenses, 53.9
- Lab-based radiation sources, 50.2–50.7, 50.4f, 50.6f, 50.8f
- Lamb dip spectroscopy, 2.5, 2.7f

- Lamb shifts, 2.2, 2.3
- Lambert-Beer law of attenuation, 63.11
- Lamor precession, 63.29–63.30
- Lamps, cold cathode fluorescent, 8.30, 8.30f
- Langmuir-Blodgett techniques, 7.34
- Large Binocular Telescope (LBT), 5.5
- Large flat-mode fibers, in fiber lasers, 25.19f, 25.20
- Large-mode-area (LMA) fibers, 25.2, 25.4, 25.5
 in all-fiber monolithic systems, 25.16
 chirally coupled core fibers, 25.21–25.22
 designs, 25.19f
 equations for, 25.8–25.9
 and photonic crystal fibers, 25.20–25.21
 techniques using, 25.18–25.20
- Laser beacons (laser guide star sensing), 5.21, 5.23, 5.27–5.34
 focus anisoplanatism of, 5.27–5.29, 5.28f–5.30f
 and mesospheric sodium laser beams, 5.32–5.34, 5.33f
 Rayleigh, 5.30–5.32, 5.31f
- Laser beam scanning, by high-resolution deflectors, 6.29
- Laser diodes, 13.3–13.24
 double heterostructure, 13.3–13.8, 13.3f, 13.4f, 13.6f, 13.7f
 noise characteristics of, 13.18–13.24, 13.19f–13.21f
 operating characteristics of, 13.8–13.13, 13.10f
 transient response of, 13.13–13.18, 13.14f, 13.16f, 13.17f
- Laser guide star (LGS) sensing, 5.21, 5.23, 5.27 (See also Laser beacons)
- Laser mode locking, 7.38–7.39
- Laser-generated plasmas, 56.1–56.10
 and Bremsstrahlung radiation, 56.8–56.10
 and recombination radiation, 56.10
 sources of, 56.1
 and spectral line emission, 56.2–56.8
- Laser-heated pedestal growth (LHPG) technique, 12.9–12.10, 12.10f
- Laser-induced fluorescence, 3.21
- Laser-induced-breakdown spectroscopy (LIBS), 3.39
- Lasers:
 Bragg reflector, 13.7
 bulk, 25.5–25.6
 buried heterostructure, 13.5
- Lasers (*Cont.*):
 carbon dioxide, 12.3t, 12.9, 12.13
 chirp of, 9.8
 continuous-wave, 25.4, 25.5, 25.5f, 25.7
 diode lasers, 13.49
 dye lasers, 5.32, 5.33f, 5.34
 distributed Bragg reflector, 9.8, 13.7, 13.28–13.29, 13.29f, 20.14
 distributed feedback
 in fiber optic systems, 13.30–13.32, 13.30f, 13.31f
 and optical fibers, 9.8
 in OTDM networks, 20.14–20.15, 20.15f
 quarter-wavelength-shifted grating, 13.31–13.32, 13.31f
 edge-emitting, 13.3
 and electroabsorption modulators, 13.60
 erbium-doped yttrium aluminum garnet (Er:YAG), 12.3t, 12.6, 12.13
 European X-ray Free-Electron Laser, 58.1
 external cavity, 13.22, 13.33
 extreme ultraviolet, 58.2–58.4, 58.3f
 Fabry-Perot, 13.12–13.13, 13.29
 multilongitudinal mode, 9.7–9.8
 semiconductor, 20.13–20.14, 20.14f
 fiber, 25.1–25.33
 all-fiber monolithic systems, 25.16–25.18, 25.16f
 architectures, 25.9–25.18, 25.12f, 25.19f
 bulk lasers vs., 25.5–25.6
 dopants, 25.22–25.26, 25.23t
 equations for, 25.8–25.9
 fabrication of, 25.26–25.29, 25.27t
 free space, 25.13–25.15, 25.14f, 25.15f
 growth of, 25.5, 25.5f
 history of, 25.3–25.4, 25.4f
 limitations, 25.6–25.7
 LMA fiber designs, 25.18–25.22, 25.19f
 operation of, 25.7–25.8
 pumping techniques, 25.9–25.13, 25.11f, 25.12f
 spectral and temporal modalities, 25.29–25.33
 fiber Raman, 10.7
 free electron, 41.9–41.10, 41.10f, 48.1, 58.1–58.2
 free space fiber, 25.9, 25.10, 25.13–25.15, 25.14f, 25.15f
 mesospheric sodium, 5.32–5.34, 5.33f
 mode-locked, 20.15–20.17, 20.16f

Lasers (*Cont.*):

- multiple quantum well, 13.24–13.25
- and optical fibers, 9.7–9.8
- in OTDM communication networks, 20.13–20.17, 20.14f, 20.16f
- phase noise (finite line width) of, 9.8
- planar buried heterostructure, 13.6
- plasma-based EUV, 58.2–58.4, 58.3f
- pulsed-dye, 5.32
- quantum well, 13.24–13.28, 13.25f–13.27f
- relative intensity noise of, 9.11
- ridge waveguide, 13.6
- for scatterometers, 1.8
- semiconductor, 13.1, 20.13–20.14, 20.14f
- single-longitudinal-mode, 9.8
- sum-frequency, 5.32
- tunable, 13.32–13.36, 13.33f–13.36f
- vertical cavity surface-emitting, 9.7n, 13.42–13.48, 13.43f, 13.45f, 19.14
- and wavelength-division multiplexing, 9.8
- x-ray, 58.1–58.4, 58.3f
- Lattice-matched compositions, SOAs and, 19.11
- Lattice-matched epitaxial layers, of fiber optic devices, 13.2
- Laue crystals, 35.3
- Laue equation, 39.1
- Laue geometry, for crystal monochromators, 39.2f, 39.4–39.6
- Laue lenses, multilayer [*see* Multilayer Laue lenses (MLLs)]
- Laue phase retarders, 43.6
- Laue transmission, 26.8, 26.9f, 63.24, 63.26
- Laue-diffracting crystals, 43.2
- Layer adding method, multiple scattering and, 3.21
- LBLRTM (code), 3.23
- Leakage channel fibers, 25.21
- Legendre polynomials, 46.6
- Legendre-Fourier (L-F) polynomials, 45.6
- Lenses:
 - Airy distribution of, 40.3
 - Bragg-Fresnel, 40.9–40.10, 40.9f
 - catadioptric systems of, 35.1
 - catoptric systems of, 35.1
 - dioptric systems of, 35.1
 - GRIN-rod, 18.7, 18.8, 18.8f, 18.10, 18.11f
 - Kerr, 7.11, 7.39
 - Kumahov capillary, 53.9

Lenses (*Cont.*):

- multilayer Laue, 42.1–42.17
 - with curved interfaces, 42.14, 42.15f
 - history of, 42.2–42.4, 42.2f, 42.3f
 - instrumental beamline arrangement and measurements, 42.9–42.10, 42.9f–42.12f
 - limitations of, 42.15–42.17, 42.16f–42.17f
 - magnetron-sputtered, 42.5–42.7, 42.6f–42.8f
 - Takagi-Taupin calculations, 42.12–42.14
 - volume diffraction calculations, 42.4–42.5, 42.5f
 - wedged, 42.12–42.13, 42.13f, 42.14f
 - and x-ray/neutron optics, 26.10
- in neutron optics, 63.22–63.23, 63.23f
- refractive x-ray, 37.3–37.11
 - applications of, 37.11
 - history of, 37.3
 - nanofocusing, 37.8–37.11, 37.9f, 37.10f
 - parabolic, 37.4–37.8, 37.4f, 37.6f, 37.7f
 - zone plates as, 40.3–40.4
- Leslie viscosity coefficients, 8.24
- Lifetime, photon, 20.1
- Light:
 - absorption of, 3.4–3.5, 3.4f
 - out-of-plane profile of, 13.11
 - propagation of, 7.3
 - retroreflection of guided, 13.6–13.7
 - spatial characteristics of, 13.11–13.12, 13.46
 - spectral characteristics of, 13.12–13.13
 - theory of interaction of atmosphere and, 3.11–3.22
 - inelastic optical processes, 3.21–3.22
 - Mie scattering, 3.16–3.18, 3.17f–3.19f
 - molecular absorption, 3.12–3.15, 3.13f
 - molecular emission and thermal spectral radiance, 3.18, 3.20, 3.20f
 - molecular Rayleigh scattering, 3.15–3.16
 - surface reflectivity and multiple scattering, 3.21, 3.21f
- Light detection and ranging (LIDAR) systems, 3.38–3.39, 3.38f, 25.25
- Light out vs. current in (L-I curve), 13.9–13.11, 13.10f, 13.46–13.47
- Light-emitting diodes (LEDs), 13.1, 13.36–13.42, 13.38f
 - edge-emitting, 13.40
 - operating characteristics of, 13.40–13.42, 13.40f
 - and optical fibers, 9.7

- Light-emitting diodes (LEDs) (*Cont.*):
 surface-emitting, 13.38–13.40, 13.38f
 and transmissive TFT LCDs, 8.29–8.31, 8.30f
- Light-hole (LH) bands, 13.27
- Linac Coherent Light Source (LCLS), 58.1
- Line edge roughness (LER), in extreme ultraviolet lithography, 34.6
- Line width, of lasers, 9.8
- Linear attenuation coefficient, 31.1, 31.2
- Linear dispersion, 38.7
- Linear electro-optic (Pockels) effect:
 and electro-optic modulators, 7.6–7.11, 7.8t, 7.11
 and liquids, 7.34
 in OTDM networks, 20.1, 20.18, 20.19
- Linear optical amplifiers (LOAs), 19.14, 19.27
- Linear polarization analyzers, 43.4, 43.4f
- Linear polarizers, 43.2–43.3, 43.3f
- Linear regime, for rare-earth-doped fiber amplifiers, 14.3
- Linear tomography, 31.5, 31.5f
- Linearity, of scatterometers, 1.15
- Line-by-line transmission programs, 3.23, 3.24f
- Linewidth enhancement factor, 13.17, 13.20
- Link budget analysis, for fiber optic communication:
 installation loss, 15.6–15.7, 15.8t
 optical power penalties, 15.8–15.17, 15.10f, 15.11f
- Linnik interferometers, 46.2
- Liouville's theorem, 54.11, 54.15, 54.16, 63.15, 63.22
- Liquid crystal (LC) cells, 8.25–8.28, 8.26f, 8.27f, 8.28f
- Liquid crystal displays (LCDs), 8.29–8.35, 61.3
 reflective, 8.31–8.32, 8.31f
 transmissive TFT, 8.29–8.31, 8.30f
 transreflective, 8.32–8.35, 8.33f, 8.34f
- Liquid crystals (LCs):
 composition of, 8.2–8.4, 8.3f
 dielectric properties of, 8.14–8.18, 8.15f, 8.18f
 elastic properties of, 8.22–8.23, 8.23f
 limitations of, 8.37
 optical properties, 8.18–8.22, 8.19f, 8.20f, 8.22f
 phase transitions of, 8.13, 8.13f, 8.14f
 phases of, 8.8–8.13, 8.11f–8.12f
 physical properties of, 8.13–8.23
 in polymer/liquid crystal composites, 8.36–8.37, 8.36f–8.37f
- Liquid crystals (LCs) (*Cont.*):
 types of, 8.4–8.8, 8.4f–8.6f, 8.7t, 8.9t
 viscosities of, 8.23–8.25, 8.24f
- Liquid-crystal-on-silicon (LCoS) displays, 8.32
- Liquid-metal anodes, 54.13
- Liquids, in electro-optic modulators, 7.34
- Lithium niobate (LiNbO_3), 6.30, 6.31, 6.31t
- Lithium niobate (LiNbO_3) modulators, 13.2, 13.48–13.55, 13.49f
 electro-optic effect, 13.49–13.51, 13.51f
 electro-optic polymer, 13.55
 high-speed operation of, 13.52–13.53
 insertion loss in, 13.53
 as Mach-Zehnder modulators, 13.51–13.52, 13.54–13.55, 13.54f
 phase modulation by, 13.51
 photorefractivity and optical damage of, 13.54
 polarization independence of, 13.53
 Y-branch interferometric, 13.51–13.52
- Lithium [${}^6\text{Li}(n, \alpha)$] reaction, 63.31
- Lithium-drifted silicon x-ray detectors, 60.6
- Lithography:
 electron-beam, 40.8
 extreme ultraviolet, 34.1–34.6, 34.2t, 34.3f–34.6f
 extreme ultraviolet-interferometric, 34.4–34.5, 34.5f
 high aspect ratio microlithography, 61.3
 interferometric, 34.4
 optical, 34.1
- L-I-V measurements, for SOAs, 19.17, 19.17f
- Lobster-eye (LE) optics, 48.2–48.4, 48.2f, 48.3f, 49.3–49.4, 49.3f–49.4f
- Lobster-ISS system, 50.7
- Local area networks (LANs), 9.14, 21.7
 power splitters and couplers in, 18.2, 18.3f
 standards for, 23.2, 23.6, 23.7
- Local gain per unit length, of lasers, 13.8
- Local oscillators, 9.13
- Local shift variance, in grazing incidence x-ray optics, 44.14
- Log-amplitude structure function, 4.5
- Long trace profiler (LTP), 46.4, 46.5, 46.5f
- Long-exposure images, 4.3–4.7
- Long-exposure MCF, 4.10
- Longitudinal electro-optic modulators, 7.16, 7.17, 7.17f
- Longitudinal spatial modulation (LSM), 6.12, 6.17

- Long-period gratings (LPGs), **21.39**, **21.39f**, **24.8–24.13**, **24.9f–24.12f**, **24.11t**, **24.13t**
- Long-wave infrared (LWIR) AOTFs, **6.42**
- Lorentzian broadening, **3.14**
- Lorentzian lineshapes, of spectra, **2.13**, **56.4–56.7**
- Loss:
- Akhieser, **6.17**
 - bend, **11.21–11.22**, **11.21f**
 - connector, **15.7**, **15.8t**
 - coupling, **18.1**
 - insertion, **7.36**, **13.53**, **18.1**
 - installation, **15.6–15.7**, **15.8t**
 - optical power, **21.13–21.14**
 - polarization-dependent, **19.18**, **21.18**
 - radiation induced, **15.17**
 - splice, **15.7**, **15.8t**
 - transmission, **15.7**
- Loss, of fiber-optic components, **18.1**
- Low cost (LC) FDDI, **23.3**
- Low leakage guidance, **11.16**, **11.17**
- Low-frequency fluctuations (LFF), of lasers, **13.23**
- LOWTRAN program, **3.23–3.24**, **3.25f–3.26f**
- L-shell emission lines, of elements, **36.3t–36.8t**
- Lump amplification, **21.44**
- Lyman series, **56.2**
- Lyotropic liquid crystals (LCs), **8.3**, **8.5f**
- Lysozyme diffraction, **53.12–53.14**, **53.12f**, **53.13f**
- Mach-Zehnder devices, as fiber-based couplers, **16.4–16.5**, **16.5f**
- Mach-Zehnder filters, **14.6**, **21.39**
- Mach-Zehnder interferometers (MZIs):
- and Bragg grating sensors, **24.8**
 - and DPSK, **21.33–21.34**, **21.33f**
 - and electro-optic modulators, **7.22**, **7.24**, **7.32**, **7.38**
 - and fiber Bragg gratings, **17.8**
 - integrated, **20.18–20.19**, **20.19f**
 - in OTDM networks, **20.22**, **20.23f**
 - SOAs in, **19.31–19.32**, **19.31f**, **19.33f**, **19.36**
 - and supercontinuum generation, **11.23**
- Mach-Zehnder modulators:
- interferometric, **13.51–13.52**, **13.54–13.55**, **13.54f**, **13.63**
 - in WDM networks, **21.30**, **21.31f**, **21.32f**
- Magnesium oxide (MgO), **2.20**, **2.20f**, **2.22**
- Magnetic circular dichroism (MCD), **55.7–55.9**, **55.7f**
- Magnetron-sputtered multilayer Laue lenses (MLLs), **42.5–42.7**, **42.6f–42.8f**
- Magnification, by gratings and monochromators, **38.7**
- Maier-Saupe mean-field coupling constant, **8.25**
- Mammography, **54.8**, **59.3–59.4**
- Manchester (biphase) coding, **20.9**, **20.9f**
- Markov random process approximation, for beam wander, **3.31**, **3.32**
- Masks, for extreme ultraviolet lithography, **34.2**, **34.6**
- Master oscillator power amplifier (MOPA) systems, **25.2**
- free space designs, **25.13–25.15**, **25.15f**
 - monolithic designs, **25.16–25.17**, **25.16f**
 - nanosecond designs, **25.31–25.32**
 - narrow linewidths, **25.30**
 - ultrashort systems, **25.33**
- MathCad program, **5.39**
- Mathematica program, **5.39**
- Maximum frequency deviation, of electro-optic modulators, **7.35**
- Maximum likelihood sequence estimations (MLSEs), **21.26**
- Maximum refractive index, for PCFs, **11.9–11.10**
- Maximum tolerable input jitter (MTIJ), **15.16**
- Maximum-likelihood expectation maximization (MLEM), **32.2**
- Maxwell-averaged Gaunt factors, **56.10**
- Maxwell-Boltzmann distribution, for EDFA spectra, **14.4**
- Maxwell's equations:
- and Bremsstrahlung radiation, **56.8**
 - for fiber-based couplers, **16.2**, **16.3**
 - and photonic crystal fibers, **11.3**, **11.6**, **11.20**
 - and wave propagation, **5.8**
- Mean-field theory, of liquid crystals, **8.22–8.23**
- Mechanical switches, for networking, **18.4**, **18.5**, **18.5f**, **18.11**, **18.11f**, **18.12f**
- Media access control (MAC), of FDDI, **23.3**
- Media interface connectors (MICs), **23.3**
- Medical imaging, **31.1–31.10**
- applications of, **31.9**, **31.10**
 - digital displays, **31.8–31.9**, **31.8f–31.10f**
 - digital tomosynthesis, **31.7–31.8**
 - and inverse Compton x-ray sources, **59.3–59.4**
 - and polycapillary x-ray optics, **53.14–53.16**, **53.15f–53.16f**
 - radiography, **31.1–31.4**, **31.2f–31.4f**
 - tomography, **31.1**, **31.5–31.7**, **31.5f–31.7f**
 - x-ray detectors for, **61.2**

- Medicine, nuclear, 32.1–32.4
- Mellin transforms, 5.16
- Mesh topologies, of WDM networks, 21.6
- Mesospheric sodium lasers, 5.32–5.34, 5.33f
- Metallization, for SOAs, 19.16, 19.16f
- Metal-semiconductor-metal (MSM) photoconductive detectors, 13.63, 13.73
- Meteorological optics, 3.40–3.43, 3.41f–3.43f
- Metrology:
- and magnetron-sputtered MLLs, 42.6–42.7, 42.7f, 42.8f
 - scatterometers in, 1.16
 - surface figure, 46.3–46.6, 46.5f
 - surface finish, 46.2
 - x-ray mirror, 46.1–46.2
 - history of, 46.1–46.2
 - profile analysis considerations, 46.6–46.12, 46.7f, 46.10f
 - surface figure metrology, 46.3–46.6, 46.5f
 - surface finish metrology, 46.2
- Metropolitan area networks (MANs), 9.14, 21.7
- Michelson interferometers, 17.8–17.9, 17.8f
- Microbunches, of electrons, 58.1
- Microcalorimeter detectors, 29.9–29.11, 29.11f, 60.9, 60.9t
- Microchannel plate (MCP) detectors, 63.34
- Microchannel plates (MCPs), 49.2, 49.3f, 49.5f, 50.7, 60.7
- Micro-electromechanical systems (MEMS)
- mirrors and switches, 18.8, 18.8f, 18.11, 18.12f
- Microfocus x-ray fluorescence (MXRF):
- with doubly curved crystal diffraction, 29.6–29.7, 29.8f–29.9f
 - monocapillary, 29.4
 - polycapillary, 29.4–29.6, 29.5f, 29.6f
 - ultrahigh resolution, 29.9–29.11, 29.9f–29.11f
- Microfocusing, with refractive x-ray lenses, 37.7–37.8
- Micro-optics-based components, for networking, 18.1–18.12
- attenuators, 18.2, 18.9
 - beam splitters, 18.6, 18.6f
 - circulators, 18.3, 18.3f, 18.10
 - directional couplers, 18.2, 18.3, 18.3f, 18.9, 18.9f
 - Faraday rotators, 18.7, 18.7f
 - filters, 18.6
 - gratings, 18.5–18.6, 18.6f
- Micro-optics-based components, for networking
(*Cont.*):
- GRIN-rod lenses, 18.7, 18.8, 18.8f
 - isolators, 18.3, 18.10, 18.10f
 - mechanical switches, 18.4, 18.5, 18.5f, 18.11, 18.11f, 18.12f
 - MEMS mirrors and switches, 18.8, 18.8f, 18.11, 18.12f
 - multiplexers/demultiplexers/duplexers, 18.4, 18.4f, 18.10–18.11
 - network functions, 18.2–18.5
 - polarizers, 18.7, 18.7f
 - power splitters, 18.2–18.3, 18.2f, 18.9, 18.9f
 - prisms, 18.5, 18.5f
- Micro-pore optics, 49.1–49.6
- Microscopes and microscopy:
- coherent diffraction, 27.4–27.5, 27.4f, 27.5f
 - Nomarski microscope, 46.2
 - scanning electron
 - and magnetron-sputtered MLLs, 42.6–42.7, 42.7f, 42.8f, 42.13, 42.13f
 - and x-ray spectral detection, 62.1–62.3, 62.2f
 - x-ray, 37.6
- Microsource devices, 28.7
- Microstrip detectors, 63.32–63.33
- Microstructured optical arrays (MOAs), adaptive, 50.7–50.8, 50.8f
- Mid-wave infrared (MWIR) AOTFs, 6.42
- Mie scattering, 3.12, 3.16–3.18, 3.17f–3.19f
- Mie theory, 11.7
- Miesowicz viscosity coefficients, 8.24
- Mirages, 3.42–3.43, 3.42f, 3.43f
- Mirror reflectivity, for VCSELs, 13.44
- Mirror surface roughness:
- and grazing-incidence neutron optics, 64.2–64.3
 - and Wolter x-ray optics, 47.3–47.5, 47.4f
- Mirrors:
- bimorph, 50.4, 50.4f, 50.6f
 - Bragg, 13.28, 13.44
 - deformable, 5.4, 5.4f, 5.37–5.38, 5.37f, 5.38f
 - double-bounce Wolter, 52.4
 - external, 13.32–13.33, 13.33f
 - Goebel, 26.10
 - high-reflectivity, 41.7–41.8, 41.8f
 - Kirkpatrick-Baez, 44.4, 44.4f, 63.21, 64.5–64.6, 64.5f
 - in micro-electromechanical systems, 18.8, 18.8f, 18.12f

- Mirrors (*Cont.*):
 for neutron optics, 63.20–63.21, 63.21f
 nonlinear optical loop, 20.22
 point-spread function of, 64.2
 polarizing, 63.28
 semiconductor laser amplifier loop optical, 20.22
 semiconductor saturable absorber, 25.3, 25.32
 and SHADOW code, 35.4, 35.5
 Wolter configurations, 52.4, 64.6, 64.6f
 (*See also* X-ray mirrors)
- Modal approach, to wavefront error correction, 4.35
- Modal dispersion, 15.9
- Modal filtering, 11.12–11.13, 11.13f
- Modal gain per unit length, 13.7, 13.8
- Modal noise, 15.16–15.17
- Mode field adaptors (MFAs), 25.2, 25.17–25.18
- Mode field diameter (MFD), 25.2
- Mode hopping, 13.13
- Mode matching, in mode field adaptors, 25.17, 25.17f
- Mode partition noise, 20.1
 for fiber optic communication links, 15.11–15.13, 15.11f
 for laser diodes, 13.19–13.20, 13.20f
- Mode transformers, photonic crystal fibers and, 11.26–11.27, 11.27f
- Mode-locked fiber lasers, 25.32–25.33
- Mode-locked lasers, 20.15–20.17, 20.16f
- Moderators, for neutron optics, 63.12, 63.14–63.15
- Modified chemical vapor deposition (MCVD), 25.2, 25.21, 25.26, 25.28
- MODTRAN program, 3.24
- Modulated grating (MG) reflectors, 13.36
- Modulation:
 amplitude, 7.22–7.24, 7.23f, 7.24f
 cross-gain, 19.12, 19.13f, 19.27, 19.29–19.30, 19.29f, 19.32, 19.35–19.36
 cross-phase
 in optical fibers, 10.3–10.4
 and SOAs, 19.13, 19.30–19.32, 19.31f, 19.33f, 19.35–19.36
 and solitons, 22.5, 22.13–22.15
 in WDM networks, 21.19
 degree of, for electro-optic modulators, 7.35–7.36
 depth-of-amplitude, 7.22
- Modulation (*Cont.*):
 depth-of-phase, 7.19
 digital, 6.33
 direct, 20.17–20.18, 20.17f
 frequency, 7.24–7.25, 7.25f
 longitudinal spatial, 6.12
 in OTDM communication networks
 direct and indirect, 20.17–20.18, 20.17f
 external, 20.18–20.20, 20.19f, 20.20f
 percent, 7.35
 phase, 7.18–7.20
 bulk electro-optic modulators, 7.18–7.20
 by lithium niobate modulators, 13.51
 polarization, 7.20–7.22, 7.20f, 7.21f
 pulse code, 20.8
 self-phase
 in optical fibers, 10.3–10.4
 and solitons, 22.3–22.4, 22.3f, 22.4f
 in WDM networks, 21.18–21.19, 21.19f
 spatial light, 6.4, 6.9
 transverse spatial, 6.11–6.12, 6.23, 6.30, 6.31
 in WDM networks, 21.27–21.36
 basic concepts, 21.27–21.29, 21.28f–21.30f
 carrier-suppressed return-to-zero and duobinary, 21.30–21.33, 21.31f, 21.32f
 comparisons of, 21.36, 21.36t, 21.37t
 DPSK and DQSK, 21.33–21.36, 21.33f–21.35f, 21.36t, 21.37t
- Modulation bandwidth, electro-optic modulators, 7.34
- Modulation efficiency, of electro-optic modulators, 7.36
- Modulation error ratio (MER), 15.4–15.5
- Modulation formats, for WDM networks, 21.27–21.36
 basic concepts, 21.27–21.29, 21.28f–21.30f
 carrier-suppressed return-to-zero and duobinary, 21.30–21.33, 21.31f, 21.32f
 comparisons of, 21.36, 21.36t, 21.37t
 DPSK and DQSK, 21.33–21.36, 21.33f–21.35f, 21.36t, 21.37t
- Modulation instability, 10.3
- Modulation response, of laser diodes, 13.16–13.17, 13.16f
- Modulation transfer functions (MTFs):
 for acousto-optic modulators, 6.32
 for polycapillary x-ray optics, 53.15, 53.16f
 for Schwarzschild objectives, 51.1–51.2, 51.2f
 and SPECT imaging, 32.3
 for x-ray detectors, 61.3

- Modulators:
- acousto-optic, 6.23*t*, 6.31–6.35, 6.32*f*, 6.34*t*
 - acousto-optic frequency shifters, 6.35
 - and Bragg diffraction, 6.4, 6.6, 6.7, 6.14
 - image (scophony), 6.34–6.35
 - multi-mode interference, 13.2, 13.35
 - principle of operation, 6.32, 6.32*f*
 - band-filling, 13.62
 - electroabsorption
 - in fiber optic systems, 13.55–13.60, 13.56*f*, 13.60
 - in OTDM networks, 20.18, 20.20, 20.20*f*
 - electro-optic, 7.1–7.39, 13.61, 20.18–20.19, 20.19*f*
 - applications for, 7.36–7.39
 - bulk modulators, 7.16–7.28, 7.21*f*, 7.24*f*–7.28*f*
 - crystal optics and the index ellipsoid, 7.3–7.7, 7.4*f*–7.6*f*, 7.8*f*–7.10*f*
 - and electro-optic effect, 7.6–7.16, 7.8*t*–7.10*t*, 7.14*f*, 7.16*f*
 - electro-optic sampling, 7.36–7.37, 7.37*f*–7.38*f*
 - and Euler angles, 7.39
 - in fiber optic systems, 13.61
 - geometries, 7.16–7.18, 7.17*f*
 - laser mode locking, 7.38–7.39
 - light propagation in, 7.3
 - materials, 7.33–7.34
 - in OTDM networks, 20.18–20.19
 - performance criteria, 7.34–7.36
 - sensors, 7.38
 - traveling wave modulators, 7.28–7.30, 7.29*f*
 - waveguide or integrated-optic modulators, 7.30–7.32, 7.31*f*–7.33*f*
 - electrorefractive, 13.61–13.62
 - integrated-optic, 7.3, 7.30–7.32, 7.31*f*–7.33*f*
 - interferometric Mach-Zehnder, 13.51–13.52, 13.54–13.55, 13.63
 - lithium niobate, 13.2, 13.48–13.55, 13.49*f*
 - Mach-Zehnder
 - interferometric, 13.51–13.52, 13.54–13.55, 13.54*f*, 13.63
 - in WDM networks, 21.30, 21.31*f*, 21.32*f*
 - Nipi, 13.62
 - semiconductor interferometric, 13.63
 - separate confinement heterostructure for, 13.4*f*, 13.5
 - waveguide, 7.30–7.32, 7.31*f*–7.33*f*, 13.56, 13.57*f*
 - Molecular absorption, 3.12–3.15, 3.13*f*
 - Molecular absorption line database, 3.22–3.23, 3.22*f*, 3.23*f*
 - Molecular emission, 3.18, 3.20, 3.20*f*
 - Molecular gases, in standard atmosphere, 3.6, 3.7*t*, 3.8*f*–3.9*f*
 - Molecular imaging, 32.1
 - Molecular spectroscopy, 2.5–2.6, 2.6*f*, 2.7*f*
 - Molecular targets, for SPECT imaging, 32.1
 - Monin-Obukhov similarity theory, 3.31
 - Monocapillary x-ray optics, 28.5, 52.1–52.6, 52.2*f*–52.5*f*, 52.2*t*
 - Monochromatic imaging, in polycapillary x-ray optics, 53.16–53.17, 53.17*f*
 - Monochromators:
 - Bragg reflection, 39.4, 39.5
 - Bragg reflections in, 39.1, 63.24
 - crystal
 - and bent crystals, 39.1–39.6, 39.2*t*, 39.3*f*, 39.5*f*–39.6*f*
 - in neutron optics, 63.23–63.25
 - Dragon* systems of, 38.3
 - Grasshopper*, 38.3
 - Johansson bent/ground focusing, 39.5
 - and SHADOW code, 35.4, 35.4*f*
 - spherical-grating, 38.3
 - synchrotron radiation, 39.6
 - toroidal-grating, 38.3
 - in VUV and soft x-ray region, 38.1–38.8
 - diffraction properties, 38.1–38.3, 38.2*f*
 - dispersion properties, 38.6–38.7
 - efficiency of, 38.8
 - focusing properties, 38.3–38.6, 38.3*f*, 38.4*t*–38.5*t*
 - resolution properties, 38.7
 - x-ray, 30.1–30.4, 50.6–50.7
 - and x-ray diffraction, 28.3, 28.4
 - Monolithic fiber laser resonators, 25.16
 - Monolithic tunable lasers, 13.33–13.36, 13.34*f*–13.36*f*
 - Mosaic crystals, 39.2
 - Mt. Pinatubo, 3.10, 3.18, 3.39
 - Mt. Wilson telescope, 5.27
 - Mueller matrices, 1.14, 19.18–19.19
 - Multichannel Bragg cells (MCBC), 6.30–6.31
 - Multicore fibers, 25.2, 25.22
 - Multidomain vertical alignment (MVA) cells, 8.25, 8.27–8.28
 - Multi-energy imaging, 54.9–54.10
 - Multi-fiber push on (MPO) connectors, 23.8

- Multifoil Kirkpatrick-Baez optics, 48.3–48.4, 48.3f
- Multifoil lobster-eye optics, 48.2–48.4, 48.2f, 48.3f
- Multifoil optics (MFO), 48.1–48.4, 48.2f, 48.3f
- Multilayer Laue lenses (MLLs):
 with curved interfaces, 42.14, 42.15f
 and hard x-rays, 42.1–42.17
 history of, 42.2–42.4, 42.2f, 42.3f
 instrumental beamline arrangement
 and measurements for, 42.9–42.10, 42.9f–42.12f
 limitations of, 42.15–42.17, 42.16f–42.17f
 with magnetron-sputtered MLLs, 42.5–42.7, 42.6f–42.8f
 on MLLs with curved interfaces, 42.14, 42.15f
 Takagi-Taupin calculations for, 42.12–42.14
 volume diffraction calculations for, 42.4–42.5, 42.5f
 with wedged MLLs, 42.12–42.13, 42.13f, 42.14f
 history of, 42.2–42.4, 42.2f, 42.3f
 instrumental beamline arrangement and measurements of, 42.9–42.10, 42.9f–42.12f
 limitations of, 42.15–42.17, 42.16f–42.17f
 magnetron-sputtered, 42.5–42.7, 42.6f–42.8f
 Takagi-Taupin calculations for, 42.12–42.14
 volume diffraction calculations for, 42.4–42.5, 42.5f
 wedged, 42.12–42.13, 42.13f, 42.14f
 and x-ray/neutron optics, 26.10
- Multilayers (reflective coatings), 41.1–41.10
 and calculation of multilayer properties, 41.3–41.4
 for diffractive imaging, 41.9–41.10, 41.10f
 fabrication and performance of, 41.4–41.9, 41.5f, 41.6t, 41.7t, 41.8f
 periodic, 42.5–42.6, 42.6f
 properties of, 41.1–41.3, 41.2f
 and x-ray diffraction, 28.5, 28.5f
- Multilongitudinal mode Fabry-Perot laser, 9.7–9.8
- Multimode fibers, for E-LEDs, 13.40
- Multimode interference (MMI) modulators, 13.2, 13.35
- Multimode interferometric Mach-Zehnder modulators, 13.54–13.55, 13.54f
- Multipath interference noise, 15.13–15.14
- Multiple Mirror Telescope (MMT), 5.5
- Multiple quantum well (MQW) lasers, 13.24–13.25
- Multiple quantum wells (MQW), 20.20
- Multiple scattering, 3.21, 3.21f
- Multiple-beam Bragg diffraction (MBD), 43.6–43.8, 43.7f
- Multiple-Bragg-beam interference, 43.1
- Multiplexers and multiplexing:
 for networking, 18.4, 18.10–18.11, 18.11f
 optical add/drop, 21.2, 21.8, 21.8f, 21.9f
 in OTDM networks, 20.1, 20.3–20.12, 20.5f–20.11f, 20.13f
 parallel, 20.12
 serial, 20.12
 time-division, 9.12, 20.3, 21.3
 (See also Wavelength division multiplexing)
- Multiplication, of avalanche photodiodes, 13.71–13.72
- Multiquantum wells (MQWs), 2.11, 19.7, 19.11f, 19.21
- Multiwire proportional counters (MWPCs), 63.31–63.32
- Mutual coherence function (MCF), 4.4, 4.7, 4.10
- Nanofocusing, of hard x-rays (see Hard x-rays, nanofocusing of)
- Nanofocusing lenses (NFLs), 37.8–37.11, 37.9f, 37.10f
- Nano-optic-measuring (NOM) machine, 46.5
- Nanosecond fiber systems, 25.30–25.32
- Narrow linewidth fiber lasers, 25.29–25.30
- Natural guide star (NGS) sensing, 5.21
- Natural line width, of spectral lines, 56.4–56.5
- Near Bragg diffraction, 6.8–6.9, 6.12
- Near field, 4.10
- Near-field diffraction, 27.2–27.3, 27.3f
- Negative core-cladding index difference, of photonic crystal fibers, 11.14–11.17, 11.14f–11.17f
- Negative hydrogen (H^-) ions, 2.3
- Negative orders of radiation, 40.1
- Nematic phase, of liquid crystals, 8.8, 8.11f, 8.11
- Neodymium-doped fibers, 25.23–25.24, 25.23t
- Networking, micro-optics-based components
 for, 18.1–18.12
 attenuators, 18.2, 18.9
 beam splitters, 18.6, 18.6f

- Networking, micro-optics-based components
for (*Cont.*):
circulators, 18.3, 18.3f, 18.10
directional couplers, 18.2, 18.3, 18.3f, 18.9, 18.9f
Faraday rotators, 18.7, 18.7f
filters, 18.6
gratings, 18.5–18.6, 18.6f
GRIN-rod lenses, 18.7, 18.8, 18.8f
isolators, 18.3, 18.10, 18.10f
mechanical switches, 18.4, 18.5, 18.5f, 18.11,
18.11f, 18.12f
MEMS mirrors and switches, 18.8, 18.8f,
18.11, 18.12f
multiplexers/demultiplexers/duplexers, 18.4,
18.4f, 18.10–18.11
network functions, 18.2–18.5
polarizers, 18.7, 18.7f
power splitters, 18.2–18.3, 18.2f, 18.9, 18.9f
prisms, 18.5, 18.5f
(*See also related topics, e.g.: Communication
networks and systems*)
- Neutron attenuation, 63.11–63.12
Neutron collimation, 63.15–63.16
Neutron filters, 63.18–63.19, 63.18t, 63.19f
Neutron gravity spectrometer, 63.21f
Neutron guides, 63.15–63.18, 63.17f
Neutron optics, 63.3–63.34
detection in, 63.31–63.34
devices for, 63.15–63.19, 63.17f, 63.18t, 63.19f
diffraction and interference in, 63.23–63.27,
63.27f
grazing-incidence, 64.1–64.7
diffractive scattering and mirror surface
roughness, 64.2–64.3
imaging focusing optics, 64.3–64.7,
64.4f–64.7f
materials of optical elements, 64.7
total external reflection, 64.1–64.2
and neutron physics, 63.3–63.5
and neutron sources, 63.12–63.15, 63.13f
polarization techniques for, 63.27–63.30,
63.30f
refraction and reflection in, 63.19–63.23,
63.21f, 63.23f
scattering lengths and cross sections,
63.5–63.12, 63.6t, 63.10t
neutron attenuation, 63.11–63.12
scattering length density, 63.9–63.11, 63.11f
and x-ray optics, 26.5–26.11, 26.8f, 26.9f,
26.11f, 36.2f
- Neutron physics, 63.3–63.5
Neutron polarization, 63.27–63.29, 63.30f
Neutron scintillators, 63.33
Neutron zone plates, 63.25
Neutrons:
epithermal, 63.18
MCP detectors for, 63.34
scattering cross sections of, 63.6–63.9, 63.6t,
63.10t
scattering length densities of, 63.9–63.11, 63.11f
scattering lengths of, 63.5–63.9, 63.6t, 63.10t
thermal, 63.3
total integrated scatter of, 64.2–64.3
(*See also Neutron optics*)
- NeXT spacecraft, 47.10
Nipi modulators, 13.62
Noise:
ASE, 19.3, 19.4f, 19.9, 19.18, 19.24, 19.35
avalanche photodiodes, 13.72–13.73
of EDFAs, 14.6
in fiber optic communication links,
15.11–15.14, 15.16–15.17
of laser diodes, 13.18–13.24, 13.20f
modal, 15.16–15.17
mode partition, 13.19–13.20, 15.11–15.13, 15.11f
multipath interference, 15.13–15.14
phase (linewidth), 13.20–13.21, 13.21f
of pin diodes, 13.70–13.71
relative intensity, 13.18–13.19, 13.19f, 15.14
of SOAs, 19.3, 19.4f, 19.9, 19.18, 19.20,
19.24, 19.35
[*See also Signal-to-noise ratio (SNR)*]
Noise equivalent BDSF (NEBDSF), for
scatterometers, 1.6, 1.8, 1.12–1.13
Noise equivalent power (NEP), 1.13, 13.71
Noise figure, of EDFAs, 14.6
Nomarski differential interference, 46.4
Nomarski microscope, images from, 46.2
Noncritical phase-matching acousto-optic
tunable filters (NPM AOTFs), 6.37, 6.38f,
6.39–6.42
angle of deflection, 6.40
angular aperture, 6.41
for long-infrared, 6.42
for mid-infrared, 6.42
optical throughput, 6.41
performance of, 6.42–6.44, 6.43t
resolution, 6.40
sidelobe suppression, 6.41–6.42
transmission and drive power, 6.41
tuning relation, 6.39–6.40
for ultraviolet, 6.42

- Nonfused fiber couplers, 25.10
- Nonlinear acoustic (NA) interaction, in acousto-optic devices, 6.30
- Nonlinear distortion, 9.17
- Nonlinear effects:
- of fiber lasers, 25.6
 - four-wave mixing, 10.2, 10.9–10.11, 10.11*f* in optical fibers, 10.1–10.12
 - self- and cross-phase modulation, 10.3–10.4
 - stimulated Brillouin scattering, 10.1, 10.7–10.9
 - stimulated Raman scattering, 10.1, 10.4–10.7, 10.5*f*
- Nonlinear length, of solitons, 22.3
- Nonlinear optical loop mirrors (NOLM), 20.22
- Nonlinear optics, WDM networks and, 21.18–21.20, 21.19*f*, 21.20*f*
- Nonlinear Schrödinger equation (NSE), 10.4, 22.2
- Non-return-to-zero differential-phase-shift-keying (NRZ-DPSK) format, 21.28
- Non-return-to-zero (NRZ) format:
- in OTDM networks, 20.8, 20.9*f*, 20.12
 - in WDM networks, 21.16, 21.29*f*, 21.32*f*
- Non-return-to-zero on-off keying (NRZ-OOK), 21.34, 21.36*t*, 21.37*t*
- Nonzero dispersion-shifted fiber (NZDSF), 21.21, 21.21*f*, 21.34*f*
- Nuclear imaging, 53.17, 53.18, 53.18*f*
- Nuclear medicine, 32.1–32.4
- Numerical aperture (NA), 9.4, 25.2, 25.18, 42.2
- NuSTAR spacecraft, 47.6, 47.7, 47.10
- Nyquist frequency, 27.4
- Nyquist frequency power, 46.8–46.9, 46.11
- Nyquist noise, 13.70
- Objectives, Schwarzschild, 26.10, 51.1–51.3, 51.2*f*–51.4*f*
- Observatories:
- Chandra, 33.2–33.4, 33.3*t*, 44.4, 44.10, 47.1, 47.4*f*, 47.5, 47.10, 64.7
 - Constellation-X, 33.4
 - Einstein, 44.4, 44.10, 47.1, 47.5
 - ROSAT, 47.5
 - SOHO, 41.3
 - Suzaku, 33.3–33.4, 33.3*t*
 - TRACE, 41.3
 - W. M. Keck, 5.27
 - XMM-Newton, 33.3, 47.2, 47.4*f*, 47.6, 47.6*f*
 - XMM-Newton observatory, 33.3*t*
 - x-ray, 33.1–33.4, 33.3*t*
- On-axis aberrations, 45.6–45.8
- On-axis optics, 64.3
- On-axis tangential phase matching, 6.25–6.26
- On-blaze condition, 38.2
- $1 \times N$ power splitters, 16.1, 16.4
- 1D profilometry, 46.6
- On-off keying (OOK), in WDM networks, 21.29, 21.30, 21.34, 21.36*t*, 21.37*t*
- Open fiber control (OFC), for Fibre Channel standard, 23.4
- Optical absorption, measurements of, 2.2–2.13, 2.4*f*, 2.6*f*–2.8*f*, 2.10*f*, 2.12*f*
- Optical add/drop multiplexers (OADMs), 21.2, 21.8, 21.8*f*, 21.9*f*, 21.12
- Optical amplifiers:
- communications applications for, 9.14
 - semiconductor vs. fiber, 9.13–9.14
 - in WDM networks, 21.37–21.44, 21.37*f*
 - EDFA, 21.38–21.41, 21.38*f*–21.42*f*
 - Raman, 21.42–21.44, 21.42*f*–21.44*f*
 - (*See also* Optical fiber amplifiers)
- Optical burst switching (OBS), 21.11
- Optical circulators, 17.8*f*, 17.9
- Optical clock recovery, 20.21–20.22, 20.21*f*
- Optical crossconnects (OXC), 21.5*f*, 21.6, 21.8, 21.8*f*, 21.10, 21.10*f*
- Optical damage, 13.54, 25.6
- Optical electric field, bulk modulators and, 7.18
- Optical fiber amplifiers, 14.1–14.11
- categories and features of, 14.1–14.2, 14.2*t*
 - erbium-doped
 - energy levels, 14.4
 - fast power transients, 21.39–21.41, 21.39*f*, 21.40*f*
 - gain flattening, 14.6–14.7, 21.38–21.39, 21.39*f*
 - gain formation, 14.4–14.5, 14.5*f*
 - gain peaking, 21.38, 21.38*f*
 - noise, 14.6
 - pump wavelength options, 14.5–14.6
 - semiconductor amplifiers vs., 9.13, 9.14, 14.1, 14.2*t*
 - static gain dynamic and channel power equalization, 21.41, 21.41*f*–21.42*f*
 - in WDM networks, 21.2–21.3, 21.2*f*
- erbium/ytterbium-doped, 14.7–14.8
- parametric, 14.10–14.11
- praseodymium-doped, 14.7
- Raman fiber, 14.8–14.9, 14.8*f*, 14.10*f*
- rare-earth-doped, 14.2–14.4, 14.3*f*
- ytterbium-doped, 14.7

- Optical fiber sensors, 24.1–24.13
 extrinsic Fabry-Perot interferometric, 24.2–24.4, 24.2*f*, 24.3*f*
 fiber Bragg grating, 24.5–24.8, 24.6*f*–24.7*f*
 intrinsic Fabry-Perot interferometric sensors, 24.4–24.5, 24.4*f*
 long-period grating sensors, 24.8–24.13, 24.9*f*–24.12*f*, 24.11*t*, 24.13*t*
- Optical fibers:
 in communication systems, 9.3–9.17
 analog transmission, 9.15–9.17
 bit rate, 9.12
 distance limits, 9.12–9.13
 fiber for, 9.4–9.7, 9.5*f*, 9.6*f*
 fiber-optic networks, 9.14–9.15
 optical amplifiers, 9.13–9.14
 photodetectors, 9.8
 receiver sensitivity, 9.8–9.11
 repeater spacing, 9.12–9.13
 technology, 9.4–9.8
 transmitting sources, 9.7–9.8
- infrared, 12.1–12.13
 applications, 12.13
 categories and properties of, 12.1–12.3, 12.2*f*, 12.2*t*, 12.3*t*
 crystalline, 12.2*t*, 12.3*t*, 12.7–12.10, 12.8*f*, 12.10*f*
 heavy-metal oxide glass in, 12.2*t*–12.4*t*, 12.3–12.7, 12.5*f*–12.7*f*
 hollow waveguides, 12.2*t*, 12.3*t*, 12.10–12.13, 12.12*f*
- nonlinear effects in, 10.1–10.12
 four-wave mixing, 10.2, 10.9–10.11, 10.11*f*
 self- and cross-phase modulation, 10.3–10.4
 stimulated Brillouin scattering, 10.1, 10.7–10.9
 stimulated Raman scattering, 10.1, 10.4–10.7, 10.5*f*
- [See also related topics, e.g.: Photonic crystal fibers (PCFs)]
- Optical holeburning (OHB), 2.13, 2.14*f*
 Optical insertion loss, of electro-optic modulators, 7.36
 Optical Kerr effect, 7.11
 Optical lithography (OL), 34.1
 Optical mode conditioners, 23.7
 Optical power dependence, of electroabsorption modulators, 13.59
- Optical power penalties, for fiber optic communication links, 15.8–15.17, 15.10*f*, 15.11*f*
 Optical signal-to-noise ratio (OSNR):
 of SOAs, 19.24–19.27
 for WDM networks, 21.20, 21.28, 21.34
 Optical spectrum analyzers (OSAs), 19.18
 Optical strength, of turbulence (C_n^2), 5.6–5.8, 5.7*f*, 5.8*f*
 Optical tank circuits, 20.21, 20.21*f*
 Optical throughput, of NPM AOTFs, 6.41
 Optical time-division multiplexed (OTDM) communication networks:
 and all-optical switching for demultiplexing, 20.22, 20.23*f*
 device technology, 20.12–20.24
 direct and indirect modulation in, 20.17–20.18, 20.17*f*
 external modulation in, 20.18–20.20, 20.19*f*, 20.20*f*
 history of, 20.3
 multiplexing in, 20.1, 20.3–20.12, 20.5*f*–20.11*f*, 20.13*f*
 and optical clock recovery, 20.21–20.22, 20.21*f*
 serial vs. parallel, 20.12, 20.13*f*
 transmitters in, 20.12–20.17, 20.14*f*–20.16*f*
 ultrahigh-speed OTDM, 20.23–20.24, 20.24*f*
 and WDM, 21.2
- Optical transfer function (OTF):
 and adaptive optics, 5.19–5.20, 5.20*f*
 and atmospheric turbulence, 4.3, 4.6–4.7
 of systems with annular pupils, 4.10–4.13, 4.12*f*, 4.13*f*
- Optical transmission, atmospheric, 3.22–3.26, 3.22*f*–3.27*f*
 Optical-electrical field overlap parameter, of electro-optic effect, 13.50
 Optically detected magnetic resonance (ODMR), 2.23–2.24, 2.24*f*
 Optically rotated tangential phase matching, 6.26*f*, 6.27
- Order selecting aperture (OSA), of zone plates, 40.4, 40.5
 Ordered dye-doped polymers, 7.34
 Ordered subsets expectation maximum (OSEM), in SPECT imaging, 32.2
 Organic crystals, 7.33–7.34
 Oscillator strength, 56.3–56.4

- Oscillators:
 and fiber lasers, 25.13, 25.14f,
 25.30–25.33
 high-power USP, 25.32–25.33
 local, 9.13
 parametric, 11.23, 11.24
 Q-switched, 25.30–25.31
 voltage-controlled, 20.11, 20.11f
 [See also Master oscillator power amplifier
 (MOPA) systems]
- Outer scale of turbulence, 4.7
- Out-of-plane profile, of emitted light, 13.11
- Outside vapor deposition (OVD), 25.2,
 25.26
- Overlapping integral, of acousto-optic
 interaction, 6.15
- Pacific Northwest National Laboratory
 program, 3.26
- Packaging, of SOAs, 19.17, 19.17f
- Packet-switched networks, 21.7, 21.10–21.11,
 21.10f–21.11f
- Parabolic reflectors, for neutron beams, 64.3,
 64.4, 64.4f
- Parallel beams, and x-ray tube sources, 54.16
- Parallel multiplexing, 20.12
- Parallel-hole collimators, 32.3
- Parametric amplifiers and oscillators, 11.23,
 11.24, 14.2, 14.10–14.11
- Parseval's theorem, 46.8
- Partial coherence length, 4.10
- Particle-induced x-ray emission (PIXE), 29.4
- Particulate matter, in standard atmosphere,
 3.6–3.7, 3.9–3.11, 3.10f, 3.11f
- Passbands, 20.1
- Passive mode locking, 20.17
- Path function, for gratings and monochromators,
 38.3–38.5, 38.3f, 38.4t–38.5t
- Patterned vertical alignment (PVA) cells,
 8.27–8.28, 8.28f
- Pauli exclusion principle, 56.2
- Pearson IV function, 56.7–56.8
- Pearson VII function, 56.7, 56.8
- Pendellösung interference, 63.26
- Perfect crystal interferometers, 63.26–63.27,
 63.27f
- Periodogram estimator, for 1D profiles, 46.8
- Periodic multilayers, of MLLs, 42.5–42.6,
 42.6f
- PETRA III (synchrotron source), 37.10
- Phase aberration function:
 correction of, 4.28–4.30
 modal expansion of, 4.17–4.20, 4.17f–4.18f,
 4.19t, 4.20t
- Phase array beam steering, 6.27–6.29, 6.28f
- Phase fluctuations, adaptive optics and,
 5.18–5.19
- Phase mask method, of FBG fabrication,
 17.5–17.6, 17.6f, 17.8, 24.7, 24.7f
- Phase matching:
 for acousto-optic devices, 6.9–6.12, 6.10f
 birefringent tangential, 6.25
 on-axis tangential, 6.25–6.26
 optically rotated tangential, 6.26f, 6.27
 tangential, 6.12, 6.13, 6.17, 6.25–6.27,
 6.26f
- Phase modulation, 7.18–7.20, 13.51 (See also
 Cross-phase modulation; Self-phase
 modulation)
- Phase modulation index, 7.19
- Phase noise, 9.8, 13.1, 13.20–13.21, 13.21f
- Phase plates, for circular polarization,
 43.5–43.6, 43.5f
- Phase retarders, 41.9, 43.6
- Phase structure function, 4.5, 5.9, 5.10
- Phase transitions, of liquid crystals, 8.13–8.14,
 8.13f, 8.14f
- Phase velocity indices of refraction, 7.15–7.16,
 7.16f
- Phase zone plates, 40.5–40.7, 40.7f
- Phase-locked loops (PLLs), 20.11, 20.11f
- Phase-measuring interferometers
 (PMIs), 46.2
- Phase-space acceptance, by gratings and
 monochromators, 38.7
- Phosphates, for fiber lasers, 25.27t, 25.28
- Phosphor x-ray detectors, 60.7–60.8, 60.10t
- Photoabsorption, 36.1
- Photodetectors, 9.8
- Photodiodes:
 avalanche, 13.63, 13.71–13.73
 pin, 13.64–13.66
 resonant, 13.65
 Schottky, 13.63, 13.73
 unitraveling-carrier, 13.68–13.69, 13.68f
- Photographic film, neutron detection with,
 63.33, 63.34
- Photomultiplier tubes (PMT), 31.5
- Photonic bandgaps (PBGs), 11.2–11.3, 11.8,
 11.8f, 11.10f, 11.11, 11.11f, 11.14

- Photonic crystal fibers (PCFs), 11.1–11.28, 25.2
 all solid-core, 25.21
 Bragg fibers, 11.4
 in cladding, 11.7–11.11, 11.8*f*–11.11*f*
 cleaving and splicing of, 11.26
 design and fabrication of, 11.4–11.6, 11.5*f*
 endlessly single-mode, 11.12, 11.13, 11.21, 11.21*f*
 in fiber lasers, 25.19*f*, 25.20–25.21, 25.27
 and guidance
 attenuation mechanisms, 11.19–11.22, 11.20*f*, 11.21*f*
 birefringence, 11.17
 core-cladding index difference, 11.12–11.17, 11.12*f*–11.17*f*
 group velocity dispersion, 11.18, 11.18*f*, 11.19*f*
 Kerr nonlinearities, 11.22–11.24, 11.24*f*
 resonance and antiresonance, 11.12
 scattering, 11.24–11.26, 11.25*f*
 history of, 11.2–11.4, 11.3*f*
 in-fiber devices for, 11.27, 11.28
 mode transformers, 11.26–11.27, 11.27*f*
 modeling and analysis of, 11.6–11.7
- Photonic integrated circuits (PICs), 19.1, 19.36
- Photons:
 energies of, 36.7*t*–36.8*t*
 lifetimes of, 20.1
- Photorefractivity, of lithium niobate modulators, 13.54
- Photoresist, for extreme ultraviolet lithography, 34.2, 34.6
- Photosensitivity, of fiber Bragg gratings, 17.2–17.3
- Photostimulable phosphors, 63.34
- Physical layer (PHY) implementation, of FDDI connectors, 23.2–23.3
- Physical vapor deposition (PVD), 61.6
- Picosecond (unit), 20.1
- Pigtail connection, fiber, 13.8
- Pin diodes, 13.2, 13.63–13.71, 13.66*t*
 dark current, 13.69
 geometry of, 13.64–13.65, 13.64*f*
 noise, 13.70–13.71
 sensitivity, 13.65–13.66, 13.66*f*
 speed, 13.67–13.68
 untraveling-carrier (UTC) photodiodes, 13.68–13.69, 13.68*f*
- Pin holes, in neutron and x-ray optics, 26.7
- Pin junctions, in fiber optic systems, 13.59
- Pin photodiodes, 13.64–13.66
- Pin waveguides, 13.68
- Pinch plasma, 57.1–57.5, 57.2*f*, 57.3*t*, 57.4*f*
- Pinhole apertures, for SPECT imaging, 32.3
- Pitch, of liquid crystals, 8.10
- Planar buried heterostructure (PBH) lasers, 13.6
- Planck radiation law, 3.18, 3.20
- Plane wave analysis, for acousto-optic interaction, 6.6–6.9
- Plasma:
 and atomic spectroscopy, 2.4–2.5
 for extreme ultraviolet lithography, 34.5
 laser-generated, 56.1–56.10
 Bremsstrahlung, 56.8–56.10
 and characteristic radiation, 56.2–56.10
 recombination radiation, 56.10
 spectral line broadening, 56.2–56.8
 pinch, 57.1–57.5, 57.2*f*, 57.3*t*, 57.4*f*
- Plasma focus, for z-pinch radiation, 57.3, 57.3*t*
- Plasma-based EUV lasers, 58.2–58.4, 58.3*f*
- P-n junctions, 20.1
- Pockels cells, 7.33
- Pockels (linear electro-optic) effect:
 and electro-optic modulators, 7.6–7.11, 7.8*t*, 7.11
 and liquids, 7.34
 in OTDM networks, 20.1, 20.18, 20.19
- Pockels' theory, elasto-optic effect and, 6.5, 6.7
- Poincaré sphere, 19.18–19.19
- Point-by-point technique, for fiber Bragg gratings, 17.8
- Point-spread function (PSF), 4.1, 4.34–4.36, 4.34*f*, 44.13, 44.14, 44.15*f*, 64.2
- Point-to-point links, in WDM networks, 21.3*f*, 21.4
- Poisson distribution, 9.10
- Poisson-effect cross coupling, 24.3
- Poisson's equation, 5.36
- Polarization:
 circular, 43.5–43.8, 43.7*f*
 phase plates for, 43.5–43.6, 43.5*f*
 and synchrotron radiation, 55.6–55.9, 55.7*f*
 eigen-, 7.13–7.16, 7.14*f*, 7.16*f*
 of insertion devices, 55.15–55.16
 of laser diodes, 13.13
 linear, 43.2–43.4, 43.3*f*, 43.4*f*, 43.6, 43.8
 neutron, 63.27–63.29, 63.30*f*
 Stokes-Poincaré parameters for, 43.2
 transverse electric, 19.7
 transverse magnetic, 19.7
 and VCSELs, 13.48
 of x-rays, 43.1–43.2

- Polarization analyzers, 43.4, 43.4f, 43.6–43.8, 43.7f
- Polarization dependence, 13.58–13.59, 19.7, 19.7f, 19.32
- Polarization independence, 6.44, 13.53
- Polarization modulation (dynamic retardation), 7.20–7.22, 7.20f, 7.21f
- Polarization scrambling, 19.18
- Polarization spectroscopy, 2.21, 2.22, 2.23f
- Polarization-dependent gain (PDG), 14.9, 19.18–19.20, 21.18
- Polarization-dependent loss (PDL), 19.18, 21.18
- Polarization-maintaining (PM) fibers, 25.3, 25.12
- Polarization-mode dispersion (PMD), 21.16–21.18, 21.17f–21.18f
- Polarizers:
- fiber-based couplers as, 16.5–16.6
 - linear, 43.2–43.3
 - multilayers, 41.9
 - for networking, 18.7, 18.7f, 18.10
- Polarizing crystal optics, 43.1–43.8
- circular polarization analyzers, 43.6–43.8, 43.7f
 - linear polarization analyzers, 43.4, 43.4f
 - linear polarizers, 43.2–43.3, 43.3f
 - phase plates for circular polarization, 43.5–43.6, 43.5f
 - and polarization of x-rays, 43.1–43.2
- Polarizing mirrors, 63.28
- Polycapillary optics:
- and brightness of x-ray tube sources, 54.16
 - collimating, 53.14, 53.14f
 - and neutron optics, 63.21–63.22
 - x-ray diffraction, 28.5
- Polycapillary x-ray optics, 53.1–53.19
- alignment and measurement in, 53.5–53.8, 53.6f, 53.7f
 - applications of, 53.10–53.19, 53.11f–53.18f
 - collimation, 53.8–53.9, 53.8f–53.9f
 - focusing, 53.9–53.10, 53.10t
 - history of, 53.1–53.3, 53.2f
 - radiation resistance in, 53.5
 - simulations and defect analysis in, 53.3–53.5, 53.3f–53.5f
- Polycrystalline (PC) fibers, 12.2, 12.3t, 12.8–12.9, 12.8f
- Polymer stabilized cholesteric texture (PSCT), of liquid crystals, 8.37, 8.37f
- Polymer sustained alignment (PSA) technique, for LC cells, 8.28
- Polymer-dispersed liquid crystals (PDLCs), 8.36, 8.36f
- Polymer/liquid crystal composites, 8.36–8.37, 8.36f–8.37f
- Polymer-stabilized liquid crystals (PSLCs), 8.36, 8.36f, 8.37
- Polynomials:
- Chebyshev, 46.6
 - Legendre, 46.6
 - Legendre-Fourier, 45.6
 - Zernike, 4.17–4.20, 4.20t, 5.10, 46.6
- Population inversions, amplification and, 19.2
- Pore optics, 49.1–49.7, 49.2f–49.6f
- Positive core-cladding index difference, for PFCs, 11.12–11.14, 11.12f, 11.13f
- Positive orders of radiation, 40.1
- Positive-intrinsic-negative (PIN) receivers, 9.8–9.10
- Positron emission tomography (PET), 32.1
- Postprocessing, of SOAs, 19.16, 19.17
- Powder diffraction, in polycapillary x-ray optics, 53.14, 53.14f
- Power:
- ASE, 19.20
 - channel, for EDFAs, 21.41, 21.42f
 - and dispersion-managed solitons, 22.12
 - for extreme ultraviolet lithography, 34.5
 - incident, of scatterometers, 1.14–1.15
 - input saturation, 14.3
 - insertion device, 55.15
 - noise equivalent, 1.13, 13.71
 - of NPM AOTFs, 6.41
 - Nyquist frequency, 46.8–46.9, 46.11
 - saturation output, 14.3, 14.4
 - Stokes, 14.9
 - of synchrotron radiation, 55.8–55.9, 55.8f
- Power amplifiers, 14.4
- Power dependence, of electroabsorption modulators, 13.59
- Power loss, 21.13–21.14
- Power penalties, of fiber optic devices, 9.11, 15.8–15.17
- Power per unit bandwidth, 7.34–7.35
- Power spectral density (PSD) function, 41.6, 41.7, 46.7–46.9
- Power splitters, 16.1, 16.4, 18.2–18.3, 18.2f, 18.3f, 18.9, 18.9f
- Poynting vectors, 7.3, 7.5, 42.2

- Praseodymium-doped fiber amplifiers (PDFAs), 14.2, 14.2*t*, 14.7
- Pre-amp semiconductor optical amplifiers (SOAs), 19.23
- Prefocusing, with refractive x-ray lenses, 37.8
- Preform manufacture, of fiber lasers, 25.26–25.27
- Pressure broadening, of spectral lines, 56.5–56.6
- Prewhitening techniques, in x-ray mirror metrology, 46.9
- Prisms, 18.5, 18.5*f*, 63.22
- Processors, real-time, 5.34–5.35, 5.34*f*, 5.35*f*
- Profile analysis, in x-ray mirror metrology, 46.6–46.12, 46.7*f*, 46.10*f*
- Profile errors, in polycapillary x-ray optics, 53.3, 53.3*f*
- Profilometry:
height, 46.3
1D, 2D, and 3D, 46.6
slope, 46.3–46.6, 46.5*f*
- Proportional counters, x-ray detectors and, 60.4–60.5, 60.9*t*, 60.10*t*
- Pulse amplification, chirped, 25.2, 25.32, 25.33
- Pulse code modulation, 20.8
- Pulsed-dye lasers, 5.32
- Pump cladding, 25.3
- Pump wavelength, EDFAs and, 14.5–14.6
- Pumping techniques, for fiber lasers, 25.9–25.13, 25.11*f*, 25.12*f*
- Pumps, for rare-earth-doped amplifiers, 14.2–14.3, 14.3*f*
- Pupils, annular, 4.10–4.16, 4.11*f*–4.15*f*, 4.15*t*
- Q factor, for OTDM networks, 20.11
- Q-switched oscillators, 25.30–25.31
- Quadratic (Kerr) electro-optic effect, 7.6, 7.9*t*–7.10*t*, 7.11
- Quadratic Stark effect, 56.6
- Quantum detection efficiency (QDE), 60.7, 61.2–61.3
- Quantum dots, 19.12, 19.21
- Quantum efficiency, of pin photodiodes, 13.65, 13.66
- Quantum limit, of digital on-off keying receivers, 9.9
- Quantum well (QW) lasers, 13.24–13.28, 13.25*f*, 13.26*f*
- Quantum wells (QWs), 20.1
coupled, 13.58
and EA modulators, 20.20
and electrorefractive modulators, 13.62
in fiber optic systems, 13.4
and SOAs, 19.7, 19.11, 19.12
and VCSELs, 13.43
- Quantum-confined Stark effect (QCSE), 13.2, 13.56, 13.56*t*, 20.1
- Quarter-wave phase plates, 43.6
- Quarter-wavelength-shifted gratings, 13.31–13.32, 13.31*f*
- Quaternary structure, of fiber optic devices, 13.2
- Radiance, thermal spectral, 3.18, 3.20, 3.20*f*
- Radiation:
Bremsstrahlung, 31.3
continuous, 54.4–54.6, 54.5*f*
from laser-generated plasmas, 56.1, 56.8–56.10
from pinch plasma sources, 57.1
and x-ray fluorescence, 29.3, 29.5, 29.6, 29.11
characteristic
Bremsstrahlung radiation as, 56.8–56.10
from laser-generated plasmas, 56.2–56.10
recombination radiation as, 56.10
spectral lines as, 56.2–56.8
from x-ray tube sources, 54.6–54.8, 54.7*f*
Compton sources of, 55.2–55.3, 55.3*t*
ionizing, 15.17
lab-based sources of, 50.2–50.7, 50.4*f*, 50.6*f*, 50.8*f*
from laser-generated plasmas, 56.2–56.10
negative orders of, 40.1
Planck radiation law, 3.18, 3.20
positive orders of, 40.1
recombination, 56.10
synchrotron sources of (*see* Synchrotron radiation sources)
Unruh, 58.2
X-pinch sources of, 57.3, 57.3*t*, 57.4
zero order of, 40.1
- Radiation induced loss, 15.17
- Radiation resistance, of polycapillary x-ray optics, 53.5
- Radiative quantum efficiency, of PDFAs, 14.7
- Radiography, 31.1–31.4, 31.2*f*–31.4*f*, 62.1
- Radiometry, infrared fibers for, 12.3*t*

- Rainbows, 3.41, 3.41*f*
- Raman amplifiers, 19.27, 21.42–21.44, 21.42*f*–21.44*f*
- Raman bands, of glass, 11.24
- Raman effect, inverse, 10.5, 10.6, 14.9
- Raman fiber amplifiers, 14.1, 14.2, 14.2*t*, 14.8–14.9, 14.8*f*, 14.10*f*
- Raman gain, 10.5, 10.6, 21.42*f*
- Raman resonance, 14.8
- Raman scattering:
 - and atmospheric optics, 3.12, 3.21
 - and configurational relaxation of solids, 2.15–2.17, 2.18*f*
 - stimulated, 25.6
 - and fiber optic communication links, 15.8
 - and optical fiber amplifiers, 14.2, 14.8
 - in optical fibers, 10.1, 10.4–10.7, 10.5*f*
 - and photonic crystal fiber guidance, 11.23, 11.24, 11.26
 - and WDM networks, 21.20
- Raman threshold, 10.7
- Raman-Nath diffraction regime, 6.4, 6.6
- Raman-Nath equations, 6.6–6.7
- Random Device Slope Scanner, 46.3
- Rare-earth ions, 2.7–2.8, 2.11
- Rare-earth-doped fiber lasers, 25.22–25.26, 25.23*t*
- Rare-earth-doped optical fiber amplifiers, 14.1, 14.2–14.4, 14.3*f*
- Ray tracing, for x-ray optics, 35.1–35.6, 35.4*f*, 35.5*f*
- Rayleigh backscattering, 10.8, 14.3
- Rayleigh beacons, 5.29–5.31, 5.31*f*
- Rayleigh criterion, 7.26, 34.1, 40.3, 42.2, 51.3
- Rayleigh range, inverse Compton scattering and, 59.1, 59.2
- Rayleigh resolution limit, 11.13
- Rayleigh scattering:
 - in glass, 11.21
 - and green flashes, 3.43
 - for HMFG fibers, 12.4
 - and laser beacons, 5.27, 5.30, 5.32
 - and MXRF, 29.5, 29.8*f*
 - and optical fibers, 9.4, 9.12
 - and Raman fiber amplifiers, 14.9
 - in theory of interaction of light and atmosphere, 3.12, 3.15–3.16
 - and x-ray optics, 26.7
- Reabsorption, of spectral lines, 56.8
- Real-time processors, for adaptive optics, 5.34–5.35, 5.34*f*, 5.35*f*
- Receivers:
 - avalanche photodiode, 9.8, 9.10–9.11
 - digital on-off-keying, 9.9–9.11
 - ideal, 9.9
 - in OTDM networks, 20.7–20.8
 - positive-intrinsic-negative, 9.8–9.10
 - in scatterometers, 1.9–1.10, 1.9*f*
 - sensitivity of fiber optic, 9.8–9.11
- Receptors, image, 31.4, 31.4*f*
- Reciprocal linear dispersion, 38.7
- Recombination radiation, 56.10
- Reconfigurability, of WDM networks, 21.12–21.13, 21.12*f*, 21.13*f*
- Reconfigurable optical add/drop multiplexers (ROADMs), 21.12
- Red-shifts, in peak gain wavelength, 25.24
- Reference method (scatterometer calibration), 1.15
- Reflectance, in interaction of light and atmosphere, 3.21, 3.21*f*
- Reflection(s):
 - Bragg
 - of crystals, 40.9
 - in interferometers, 63.26, 63.27
 - and linear polarization, 43.2–43.4, 43.3*f*, 43.4*f*, 43.6, 43.8
 - and liquid crystals, 8.10
 - in monochromators, 63.24
 - and multilayers, 41.2
 - simultaneous, 43.7 [See also Multiple-beam Bragg diffraction (MBD)]
 - and x-ray absorption spectroscopy, 30.2, 30.4
 - Fresnel, 13.6, 13.53, 17.3, 25.8
 - grazing-angle, 63.21
 - in neutron and x-ray optics, 26.8–26.9, 63.20–63.21, 63.21*f*
 - reference, for Fabry-Perot sensors, 24.2
 - retro-, 13.6–13.7
 - sensing, for Fabry-Perot sensors, 24.2
 - total external, 64.1–64.2
 - total internal, 11.2, 11.3, 11.14, 11.15
- Reflective LCDs, 8.31–8.32, 8.31*f*
- Reflective semiconductor optical amplifiers (SOAs), 19.28
- Reflectivity, of Wolter x-ray optics, 47.5

- Reflectors:
 Bragg, 13.45
 modulated grating, 13.36
 for neutron beams, 64.3, 64.4, 64.4f
- Refraction:
 electro-, 13.2
 in neutron optics, 63.19–63.23, 63.21f, 63.23f
 in x-ray optics, 26.6–26.8, 26.8f
- Refractive index (*see* Index of refraction)
- Refractive x-ray lenses, 37.3–37.11
 applications of, 37.11
 history of, 37.3
 nanofocusing, 37.8–37.11, 37.9f, 37.10f
 parabolic, 37.4–37.8, 37.4f, 37.6f, 37.7f
- Registered detected point spread function (RDPSF), 44.14, 44.15f
- Relative intensity noise (RIN):
 in fiber optic communication links, 15.2, 15.14
 of laser diodes, 13.18–13.19, 13.19f, 13.23
 and optical fibers, 9.11, 9.16
- Relaxation, configurational, 2.14–2.17, 2.15f–2.18f
- Relaxation oscillations, of laser diodes, 13.14–13.16
- Remote sensing, in atmospheric optics, 3.36–3.40, 3.37f–3.40f
- Repeater spacing, for optical fibers, 9.12–9.13
- Resolution:
 angular, 47.2–47.3, 47.2f, 47.10–47.11, 47.11f
 of deflectors, 6.25, 6.29, 6.30
 field-weighted-average, 44.10
 of gratings and monochromators, 38.7
 of NPM AOTFs, 6.40
 spatial, 26.10–26.11
 of telescopes, 4.2–4.3
 in x-ray imaging, 62.3–62.5
- Resolution enhancement techniques (RET), for extreme ultraviolet lithography, 34.1
- Resonance, of PCFs, 11.12
- Resonant cavity light-emitting diodes (RC-LEDs), 13.39, 13.40
- Resonant circuits, electro-optic modulators and, 7.34
- Resonant photodiodes, 13.65
- Resonators, 25.13, 25.16
- Responsivity, of pin diodes, 13.65–13.66, 13.66f
- Restimulable phosphor detectors, 60.8, 60.10t
- Retarders, phase, 41.9, 43.6
- Retroreflection, of guided light, 13.6–13.7
- Return-to-zero differential quadrature phase-shift-keying (RZ-DQPSK) format, 21.35, 21.35f, 21.36t, 21.37t
- Return-to-zero differential-phase-shift-keying (RZ-DPSK) format, 21.28, 21.34f, 21.36t, 21.37t
- Return-to-zero (RZ) format, 20.8, 20.9f, 20.10f, 21.16, 21.29f, 21.31f
- Return-to-zero on-off keying (RZ-OOK), 21.34, 21.36t
- Rician density function, 3.36
- Ridge waveguide (RWG) laser, 13.6
- Rigid-body motions, 45.6–45.7
- Ring topologies, for WDM networks, 21.5–21.6, 21.5f–21.7f
- Robustness, of solitons, 22.4–22.5
- ROSAT observatory, 47.5
- ROSAT (Rontgensatellit) telescope, 44.4, 44.6f, 44.10
- Rose model, for x-ray attenuation, 31.2
- Rotating anodes, as x-ray tube sources, 54.12
- Rotationally parabolic profiles, for refractive x-ray lenses, 37.4–37.8, 37.4f, 37.6f, 37.7f
- Rotators, Faraday, 18.7, 18.7f, 18.10
- Roto-optic effect, 6.6
- Roughness, polycapillary x-ray optics, 53.4
- Rowland circle, 39.5, 39.6
- Rowland spherical grating, 38.6
- Russell Saunders coupling, 2.11
- Rydberg energy, 54.7
- Rytov transformation, 5.9
- SAGE II satellite system, 3.39, 3.40f
- Sagittal-focusing geometry, for monochromators, 39.6
- Sagnac interferometers, 20.22
- Sagnac loops, 17.8
- Sample mounts, for scatterometers, 1.9
- Sampled gratings (SG), 13.34, 13.34f
- Sampling, in OTDM networks, 20.1, 20.4–20.6, 20.5f, 20.6f
- Sapphire, 12.3t, 12.9–12.10, 12.10f
- Saturated output power, of rare-earth-doped amplifiers, 14.4
- Saturation, of SOAs, 19.9–19.10, 19.10f
- Saturation current, of semiconductor diodes, 13.69
- Saturation output power, of rare-earth-doped amplifiers, 14.3

- Saturation regime, for rare-earth-doped amplifiers, 14.3
- Scaling law, or photonic crystal fibers, 11.7
- Scan rate, of high-resolution deflectors, 6.29
- Scanners, 7.26–7.28, 7.26f–7.28f
- Scanning electron microscopy (SEM):
 and magnetron-sputtered MLLs, 42.6–42.7, 42.7f, 42.8f, 42.13, 42.13f
 and x-ray fluorescence, 29.2, 29.3
 and x-ray spectral detection, 62.1–62.3, 62.2f
- Scatter function (cosine-corrected BRDF), 1.6
- Scatter rejection, polycapillary x-ray optics and, 53.14–53.16, 53.15f–53.16f
- Scattering:
 anti-Stokes, 10.5
 Brillouin
 backward, 11.25
 forward, 11.25, 11.26
 photonic crystal fibers, 11.25–11.26, 11.25f
 stimulated, 10.1, 10.7–10.9, 15.8, 21.20, 25.6
 Compton, 59.1
 and circular polarization, 43.6
 inverse, 59.1
 and MXRF, 29.5, 29.8f
 and polycapillary x-ray optics, 53.3, 53.15, 53.18
 and refractive x-ray lenses, 37.5, 37.7
 and x-ray attenuation, 31.2
 and x-ray optics, 26.7, 36.1
 elastic, 63.5
 in heavy water, 63.10
 incoherent, 26.7, 31.2, 63.7, 63.8
 inelastic, 63.3
 and photonic crystal fibers, 11.24–11.26, 11.25f
 Raman, 11.24
 and atmospheric optics, 3.12, 3.21
 and configurational relaxation of solids, 2.15–2.17, 2.18f
 Rayleigh, 11.21
 backscattering, 10.8, 14.3
 in glass, 11.21
 and green flashes, 3.43
 for HMFG fibers, 12.4
 and laser beacons, 5.27, 5.30, 5.32
 and MXRF, 29.5, 29.8f
 and optical fibers, 9.4, 9.12
 and Raman fiber amplifiers, 14.9
 Rayleigh (Cont.):
 in theory of interaction of light and atmosphere, 3.12, 3.15–3.16
 and x-ray optics, 26.7
 by silica-air photonic crystal fibers, 11.25
 small angle neutron, 64.1
 stimulated Raman, 25.6
 and fiber optic communication links, 15.8
 and optical fiber amplifiers, 14.2, 14.8
 in optical fibers, 10.1, 10.4–10.7, 10.5f
 and photonic crystal fiber guidance, 11.23, 11.24, 11.26
 and WDM networks, 21.20
 Thompson, 26.7
 Thomson backscattering, 59.1
 in water, 63.10, 63.10t, 63.11f
- Scattering cross sections, for neutrons, 63.6–63.9, 63.6t, 63.10t
- Scattering length densities, of neutrons, 63.9–63.11, 63.11f
- Scattering lengths, of neutrons, 63.5–63.9, 63.6t, 63.10t
- Scatterometers, 1.3–1.16
 BSDF, 1.8, 1.8f
 calibration of, 1.14–1.15
 configurations and components for, 1.7–1.11, 1.8f–1.10f
 error analysis of, 1.15
 incident power measurement of, 1.14–1.15
 instrument signature and quality, 1.11–1.13, 1.11f, 1.13t
 measurement issues with, 1.13–1.14
 in metrology, 1.16
 specifications of, 1.5–1.7, 1.5f
- Schawlow-Townes linewidth, 13.21
- Schottky barriers, 13.58, 13.69
- Schottky photodiodes, 13.63, 13.73
- Schroder-Van Laar equation, 8.14
- Schrödinger equation, 63.4–63.5
- Schwarzschild objectives, 26.10, 51.1–51.3, 51.2f–51.4f
- Schwarzschild optics, 44.6, 48.2
- Schwinger formula, 55.15
- Scintigraphy, 32.1, 53.17, 53.18, 53.18f
- Scintillation, 3.26, 3.34–3.36, 3.35f
- Scintillation x-ray detectors, 60.7–60.8, 60.9t, 60.10t
- Scintillator-based flat panel detectors, 61.4, 61.4f, 61.7–61.8, 61.8f
- Scintillators, neutron, 63.33

- Secondary electrons, **54.12**
- Segmented deformable mirrors, **5.37, 5.37f**
- Selective area epitaxy, **13.60**
- Self-amplified spontaneous emission (SASE)
mode, of lasers, **58.1**
- Self-healing ring networks, for WDM,
21.6, 21.6f
- Self-phase modulation (SPM):
in optical fibers, **10.3–10.4**
and solitons, **22.3–22.4, 22.3f, 22.4f, 22.13**
in WDM networks, **21.18–21.19, 21.19f**
- Self-saturation effect, of EDFAs, **14.6**
- Self-steepening, **10.4**
- Semiconductor (x-ray) detectors, **29.3, 60.5–60.6, 60.9t, 60.10t**
- Semiconductor Equipment and Materials International (SEMI), **1.4**
- Semiconductor interferometric modulators,
13.2, 13.63
- Semiconductor laser amplifier loop optical mirrors (SLALOM), **20.22**
- Semiconductor laser amplifiers (SLAs),
9.13–9.14
- Semiconductor lasers, **13.1**
- Semiconductor optical amplifiers (SOAs),
19.1–19.36
amplification in, **19.1f–19.2f, 19.2, 19.22–19.27, 19.22f–19.26f**
ASE noise, **19.3, 19.4f**
confinement factor, **19.6**
device characterization, **19.17–19.21, 19.17f, 19.19f–19.21f**
fabrication, **19.15–19.17, 19.15f–19.17f**
gain, **19.4–19.6, 19.4f–19.6f**
gain clamping, **19.14–19.15, 19.14f**
gain dynamics, **19.12–19.13, 19.13f, 19.14f**
gain ripple and feedback reduction,
19.8, 19.8f
history of, **19.1**
material systems, **19.11–19.12, 19.11f–19.12f**
noise figure, **19.9**
nonlinear applications, **19.29–19.36, 19.29f, 19.31f, 19.33f, 19.34f**
and photonic integrated circuits, **19.36**
polarization dependence, **19.7, 19.7f**
saturation, **19.9–19.10, 19.10f**
switching and modulation, **19.28, 19.28f**
- Semiconductor saturable absorber mirror (SESAM), **25.3, 25.32**
- Semiconductors:
and electroabsorption modulators, **13.58**
in electro-optic modulators, **7.34**
and extreme ultraviolet lithography,
34.1–34.2, 34.2t
fiber amplifiers vs., **9.13–9.14**
signal processing in complementary metal oxide, **62.5**
- Sensing reflection, for Fabry-Perot sensors, **24.2**
- Sensitivity (responsivity), of pin diodes,
13.65–13.66, 13.66f
- Sensors:
electro-optic modulators, **7.38**
extrinsic Fabry-Perot interferometric (EFPI),
24.2–24.4, 24.2f, 24.3f
fiber optic chemical, **12.3t**
intrinsic Fabry-Perot interferometric (IFPI),
24.4–24.5, 24.4f
optical fiber, **24.1–24.13**
comparison of, **24.13**
extrinsic Fabry-Perot interferometric,
24.2–24.4, 24.2f, 24.3f
fiber Bragg grating, **24.5–24.8, 24.6f–24.7f**
intrinsic Fabry-Perot interferometric,
24.4–24.5, 24.4f
long-period grating sensors, **24.8–24.13, 24.9f–24.12f, 24.11t, 24.13t**
Shack-Hartmann, **5.21, 5.36, 5.40–5.43**
Shack-Hartmann wavefront, **50.5, 50.6f**
transition-edge, **60.9**
- Separate confinement heterostructure (SCH),
in fiber optic modulators, **13.4f, 13.5**
- Serial byte connection (SBCON) standard, **23.1**
- Serial multiplexing, **20.12**
- Servo lag, subaperture size and, **5.41–5.42**
- Shack-Hartmann sensors, **5.21, 5.36, 5.40–5.43**
- Shack-Hartmann technique, for adaptive optics, **5.23–5.27, 5.23f, 5.25f, 5.26f**
- Shack-Hartmann wavefront sensors, **50.5, 50.6f**
- SHADOW code, **35.2–35.5, 35.4f, 35.5f**
- SHADOWVUI interface, **35.3, 35.5, 35.5f**
- Shane telescope, **5.27, 5.32**
- Shielding, in neutron optics, **63.13–63.14**
- Short-exposure images, **4.3, 4.31–4.35, 4.31f–4.34f, 4.35t**
- Short-period gratings (SPGs), **24.6**
- Shot noise, **13.70, 13.73**
- Sidelobe suppression, in AOTFs, **6.41–40.2, 6.44–6.45**

- Signal processing, by wideband AO Bragg cells, 6.30
- Signal-dependent noises, 9.11
- Signal-to-noise ratio (SNR):
 for electro-optic modulators, 7.16
 for fiber optic communication links, 15.2, 15.4, 15.5
 for laser diodes, 13.19
 in neutron optics, 63.13
 optical
 for SOAs, 19.24–19.27
 for WDM networks, 21.20, 21.28, 21.34
 for optical fibers, 9.9
 for solitons, 22.1, 22.6–22.8
- Silica fibers, in lasers, 25.27*t*, 25.28
- Silica-air photonic crystal fibers, 11.14, 11.14*f*, 11.25
- Silicon pore optics, 49.6–49.7, 49.7*f*
- Silver halide fibers, 12.3*t*, 12.8–12.9, 12.8*f*
- Symbol-X spacecraft, 47.10
- Simultaneous Bragg reflection (multiple-beam Bragg diffraction), 43.6–43.8, 43.7*f*
- Single attach FDDI nodes, 23.3
- Single capillaries (*see* Monocapillary x-ray optics)
- Single crystal diffraction, 53.12–53.14, 53.12*f*–53.13*f*
- Single crystal (SC) fibers, as infrared optical fibers, 12.3*t*, 12.9–12.10, 12.10*f*
- Single mirror resonators, 25.13
- Single photon emission computed tomography (SPECT), 32.1–32.4
- Single-bounce monocapillary x-ray optics, 52.5–52.6, 52.5*f*
- Single-cell-gap transreflective LCDs, 8.33–8.35, 8.34*f*
- Single-channel systems, of SOAs, 19.22–19.24, 19.22*f*–19.24*f*
- Single-longitudinal-mode (SLM) lasers, 9.8
- Single-mode excitation technique (SMET), 25.3, 25.20
- Single-mode fibers (SMFs):
 dispersion in, 21.16, 21.16*f*, 21.21, 21.21*f*
 for E-LEDs, 13.40
 and photonic crystal fibers, 11.2, 11.11, 11.13, 11.23–11.25
- Singular-value-decomposition (SVD) algorithms, 5.25, 5.26
- Slab-coupled optical waveguide amplifiers (SCOWAs), 19.21
- Slater integrals, 2.8
- Slope errors, 38.7, 45.7–45.8
- Slope measurement analysis, in x-ray mirror metrology, 46.11–46.12
- Slope profilometry, 46.3–46.6
- Slot-detection x-ray imaging, 61.2, 61.2*f*
- Slowly varying envelope approximation (SVEA), 10.2, 10.6
- Small angle neutron scattering (SANS), 64.1
- Smectic phase, of liquid crystals, 8.8–8.12, 8.11*f*, 8.12*f*
- Snell's law, 19.8, 26.8
- Soft x-ray region, gratings and monochromators in, 38.1–38.8
 diffraction properties, 38.1–38.3, 38.2*f*
 dispersion properties, 38.6–38.7
 efficiency, 38.8
 focusing properties, 38.3–38.6, 38.3*f*, 38.4*t*–38.5*t*
 resolution properties, 38.7
- Soft x-ray telescopes, 50.2
- Soft x-rays, circularly polarized, 55.6–55.7
- Software implementation, for adaptive optics, 5.21–5.38
 higher-order wavefront sensing techniques, 5.36–5.37
 laser beacons, 5.27–5.34, 5.28*f*–5.31*f*, 5.33*f*
 real-time processors, 5.34–5.35, 5.34*f*, 5.35*f*
 Shack-Hartmann technique, 5.23–5.27, 5.23*f*, 5.25*f*, 5.26*f*
 tracking, 5.21–5.23, 5.22*f*
 wavefront correctors, 5.37–5.38, 5.38*f*
- SOHO observatory, 41.3
- Solar x-ray Imager (SXI), 44.9, 44.12–44.13, 44.16–44.17, 44.16*f*, 44.17*f*
- Solid core photonic crystal fibers:
 attenuation in, 11.20, 11.21
 Brillouin scattering in, 11.25
 group velocity dispersion, 11.18, 11.18*f*, 11.19*f*
 Kerr effects for, 11.23
- Solid state detectors, in neutron optics, 63.33–63.34
- Solids:
 configurational relaxation in, 2.14–2.17, 2.15*f*–2.18*f*
 optical absorption measurements of, 2.7–2.13, 2.8*f*, 2.10*f*, 2.12*f*
 zero-phonon lines in, 2.13–2.14, 2.14*f*
- Solid-state ion chambers, 60.5–60.6

- Solitons:
- classical, 22.2–22.4, 22.2*f*–22.4*f*
 - in communication systems, 22.1–22.17
 - dispersion-managed, 22.12–22.15, 22.13*f*, 22.14*f*
 - effects of, 22.1–22.2
 - errors, 22.6
 - and frequency-guiding filters, 22.7–22.9
 - in optical amplifiers, 9.14
 - and optical fibers, 10.1
 - properties of, 22.4–22.5
 - self-frequency shift cancellation, 11.24
 - transmission systems for, 22.5–22.7
 - and wavelength division multiplexing, 22.2, 22.9–22.12, 22.15–22.17
- Soller collimators, 28.2, 28.2*f*
- Soller slits, 26.7
- SOR telescope, 5.27, 5.32, 5.34, 5.38
- Source debris, in extreme ultraviolet lithography, 34.5
- Source depth measurements, for x-ray tube sources, 54.14, 54.14*f*
- Source spot measurements, for x-ray tube sources, 54.14, 54.14*f*, 54.15, 54.15*f*
- Source-model technique, for photonic crystal fibers, 11.7
- Space-charge region, of pin photodiodes, 13.64
- Spatial coherence, 4.10
- Spatial coherence length, 27.2
- Spatial frequency spectrum, 6.16
- Spatial hole burning, 13.13, 20.2
- Spatial light modulation (SLM), 6.4, 6.9
- Spatial resolution, 26.10–26.11
- Spatial scale, adaptive optics and, 5.10–5.13, 5.11*t*, 5.12*f*, 5.13*f*
- Spectra (code), 55.14
- Spectral and temporal modalities, of fiber lasers, 25.29–25.33
- Spectral hole burning, 13.13
- Spectral inversion, four-wave mixing and, 10.11
- Spectral lines:
 - broadening of
 - Doppler, 3.14, 5.32, 56.5
 - homogeneous, 20.1
 - Lorentzian, 3.14
 - pressure, 56.5–56.6
 - spectral, 56.6–56.8
 - from laser-generated plasmas, 56.2–56.8
- Spectral radiance, thermal, 3.18, 3.20, 3.20*f*
- Spectrometers, 1.14, 63.21*f*
- Spectroscopy:
 - and absorption/photoluminescence of Cr³⁺, 2.19–2.21, 2.19*f*, 2.20*f*
 - in adaptive optics, 5.19–5.21, 5.19*f*, 5.20*f*
 - atomic, 2.4–2.5
 - evanescent wave, 12.13
 - excitation, 2.15*f*, 2.21
 - Fourier transform, 2.5, 2.6*f*
 - and homogeneous lineshapes of spectra, 2.13–2.17, 2.14*f*–2.18*f*
 - Lamb dip, 2.5, 2.7*f*
 - Laser-induced-breakdown, 3.39
 - measurements from, 2.1–2.24
 - molecular, 2.5–2.6, 2.6*f*, 2.7*f*
 - optical, 2.9–2.10, 2.10*f*
 - optical absorption measurements
 - of atomic energy levels, 2.2–2.5, 2.4*f*
 - of solids, 2.7–2.13, 2.8*f*, 2.10*f*, 2.12*f*
 - polarization, 2.21–2.22, 2.23*f*
 - time-domain, 7.37, 7.37*f*–7.38*f*
 - time-of-flight, 2.5
 - x-ray absorption, 30.1–30.5, 30.2*f*–30.5*f*
 - x-ray absorption near edge, 30.2, 30.2*f*, 30.4*f*
 - Zeeman, 2.23–2.24, 2.24*f*
- Spectrum(-a):
 - atomic, 2.13
 - homogeneous lineshapes of, 2.13–2.17, 2.14*f*–2.18*f*
 - Kolmogorov, 3.28, 3.31
 - spatial frequency, 6.16
 - von Kármán, 3.28, 5.6
 - of x-ray tube sources, 54.4–54.10, 54.5*f*, 54.8*f*
 - (See also Spectral lines)
- Speed:
 - of avalanche photodiodes, 13.72
 - of pin diodes, 13.67–13.68
- Spherical aberrations, 45.4
- Spherical-grating monochromators (SGMs), 38.3
- Spin flippers, 63.29
- Splay, of liquid crystals, 8.22, 8.23*f*
- Splice losses, for fiber optic communication links, 15.7, 15.8*t*
- Splicing, of PCFs, 11.26
- Split-off band, of strained layer quantum well lasers, 13.27
- Spontaneous emission, 2.13, 19.3, 20.1
- Spring-8 Compact SASE Source (SCSS), 58.1
- Spurious-free dynamic range (SFDR), 15.6
- SRW (code), 55.14

- Stack and draw method, of fiber laser fabrication, 25.27
- Stacked actuator continuous facesheet deformable mirrors, 5.37*f*, 5.38, 5.38*f*
- Standard atmosphere, composition of, 3.6–3.11, 3.7*t*, 3.8*f*–3.11*f*
- Standard dispersion-shifted fibers (DSFs), 24.10–24.11, 24.11*t*, 24.13*t*
- Star topologies, for WDM networks, 21.5–21.6, 21.5*f*, 21.7*f*
- Starfire Optical Range, 5.23
- Stark effect, 2.21, 20.20, 25.24, 56.6
- Stark levels, of EDFAs, 14.4, 14.5, 14.5*f*
- Static gain dynamic, for EDFAs, 21.41, 21.41*f*
- Static Jahn-Teller effect, 2.9
- Stationary anodes, of x-ray tube sources, 54.12
- Step-and-scan stage, of extreme ultraviolet lithography, 34.2
- Stern-Gerlach experiment, 63.4
- Stimulated Brillouin scattering (SBS), 10.1, 10.7–10.9, 15.8, 21.20, 25.6
- Stimulated emission, 20.2
- Stimulated Raman scattering (SRS), 25.6
and fiber optic communication links, 15.8
and optical fiber amplifiers, 14.2, 14.8
in optical fibers, 10.1, 10.4–10.7, 10.5*f*
and photonic crystal fiber guidance, 11.23, 11.24, 11.26
and WDM networks, 21.20
- Stokes frequency, 10.8, 21.42
- Stokes intensity, 10.6, 10.7
- Stokes power, 14.9
- Stokes shifts, 2.15–2.16
- Stokes waves, 10.4–10.9, 14.2
- Stokes-Poincaré parameters, for polarization, 43.2, 43.8
- Strain, in fiber optic devices, 13.2
- Strained layer quantum well lasers, 13.26–13.28, 13.26*f*, 13.27*f*
- Strained-layer superlattices (SLs), 2.11–2.13, 2.12*f*
- Strehl ratio:
and adaptive optics, 5.14, 5.14*f*–5.15*f*, 5.17, 5.17*f*, 5.20–5.22, 5.22*f*, 5.29, 5.29*f*, 5.35, 5.35*f*, 5.39, 5.40, 5.43, 5.43*f*–5.46*f*, 5.46
and imaging through atmospheric turbulence, 4.13, 4.14*f*, 4.15*t*, 4.27–4.28, 4.28*f*, 4.32–4.33, 4.34*f*, 4.36
- Strip loading, of RWG lasers, 13.6
- Subaperture size, for adaptive optics, 5.40–5.43, 5.42*f*
- Sum-frequency lasers, 5.32
- Superconducting quantum (SQUID) detector, 60.9
- Superconducting tunneling junctions (STJ) detectors, 60.9, 60.9*t*
- Supercontinuum (SC) generation, 11.23, 11.24*f*
- Supermirrors, 41.7, 41.8, 64.7
- Superstructure gratings (SSGs), 13.34, 13.35*f*
- Surface axial coordinates, in grazing incidence optics, 45.7
- Surface figure metrology, 46.3–46.6, 46.5*f*
- Surface finish metrology, 46.2
- Surface-emitting light-emitting diodes (S-LEDs), 13.36–13.40, 13.38*f*
- Suzaku observatory, 33.3–33.4, 33.3*t*
- SUZUKU spacecraft, 47.6, 47.7*f*
- Swift Burst Alert Telescope (BAT), 33.1
- Switches and switching:
for fiber-based couplers, 16.4
for networking, 18.4, 18.5, 18.5*f*, 18.11, 18.11*f*, 18.12*f*
and SOAs, 19.28, 19.28*f*
- Synchronous digital hierarchy (SDH), 9.17
- Synchronous optical networks (SONETs), 9.17, 21.6, 23.6
- Synchrotron beamlines, SHADOW code and, 35.1–35.2, 35.4
- Synchrotron radiation (SR) monochromators, 39.6
- Synchrotron radiation sources, 55.1–55.20
adaptive x-ray optics for, 50.2–50.7, 50.4*f*, 50.6*f*, 50.8*f*
coherence of, 55.17–55.20, 55.18*f*–55.19*f*
Compton sources vs., 55.2–55.3, 55.3*t*
history of, 55.1–55.2
insertion devices, 55.9–55.16, 55.10*f*, 55.12*f*, 55.13*f*
as linear polarizers, 43.2, 43.3
and refractive x-ray lenses, 37.5, 37.11
theory of synchrotron radiation emission, 55.2–55.9, 55.2*f*, 55.3*f*, 55.5*f*–55.8*f*
- Takagi-Taupin calculations, for MLLs, 42.12–42.14
- Talystep stylus profiler, 46.2
- Tanabe-Sugano diagrams, 2.9–2.11, 2.10*f*, 2.19
- Tangential phase matching (TPM), 6.12, 6.13, 6.17, 6.25–6.27, 6.26*f*

- Tapered fiber bundles (TFBs), 25.3, 25.10–25.11, 25.11*f*, 25.16
- Tapered fiber coupler process, 16.1–16.3, 16.2*f*
- Tapered fiber method, for adaptor fabrication, 25.17, 25.17*f*
- Teflon coatings, for infrared optical fibers, 12.4, 12.7, 12.9
- Telecommunication systems, data communication vs., 15.1–15.2
- Telescopes:
- and adaptive optics, 5.2–5.5, 5.3*f*, 5.4*f*, 5.43–5.46, 5.43*f*–5.46*f*
 - Burst Alert, 33.1
 - Calar Alto, 5.27
 - Canada-France-Hawaii, 5.23
 - Cassegrain, 44.4
 - ESO, 5.35
 - Gemini North, 5.20, 5.21*f*
 - with grazing incidence optics, 44.6–44.12, 44.7*f*, 44.9*f*–44.11*f*, 45.1
 - hard vs. soft x-ray, 50.2
 - Hobby-Eberly, 5.2
 - hyperboloid-hyperboloid (HH) grazing incidence x-ray, 44.10–44.12, 44.11*f*
 - Keck, 4.36, 5.2, 5.27
 - Large Binocular, 5.5
 - Mt. Wilson, 5.27
 - Multiple Mirror, 5.5
 - resolution of, 4.2–4.3
 - ROSAT, 44.4, 44.6*f*
 - Shane, 5.27
 - SOR, 5.27
 - Swift Burst Alert, 33.1
 - 10-m, 5.45–5.46, 5.45*f*–5.46*f*
 - 3.5-m, 5.43–5.45, 5.43*f*–5.45*f*
 - wide-field lobster-eye, 48.4
 - Wolter, 44.4, 44.5*f*, 44.6–44.10, 44.7*f*, 44.10*f*, 45.1–45.5, 45.2*f*
 - Wolter-Schwarzschild, 44.7, 44.11, 45.1
- Tellurites, for fiber lasers, 25.27*t*, 25.28
- Tellurium dioxide (TeO_2), 6.17, 6.21*t*, 6.25, 6.29*t*, 6.34*t*, 6.39, 6.42
- Temperature:
- and laser diodes, 13.11
 - and liquid crystals, 8.17, 8.21–8.22, 8.22*f*
 - and long-period grating sensors, 24.13, 24.13*t*
- Temporal coherence, of synchrotron radiation sources, 55.17–55.18, 55.18*f*
- 10-m telescope systems, 5.45–5.46, 5.45*f*–5.46*f*
- Terabit (unit), 20.2
- Terahertz asymmetric optical demultiplexers (TOAD), 20.22
- Ternary layers, in fiber optic devices, 13.2
- Thermal control and correction, for adaptive optics, 50.5, 50.6
- Thermal imaging, 12.3*t*
- Thermal neutrons, 63.3
- Thermal (Johnson) noise, 13.70
- Thermal runaway, of lasers, 13.11
- Thermal spectral radiance, 3.18, 3.20, 3.20*f*
- Thermally expanded cores (TECs), 25.3, 25.17, 25.17*f*
- Thermoelectric coolers (TECs), 19.4*f*, 19.5
- Thin film transistors (TFTs), 61.3, 61.4, 61.6
- Thin lenses, zone plates as, 40.3–40.4
- Thompson scattering, 26.7
- Thomson backscattering, 59.1
- Thomson scattering cross-section, 59.2
- 3D profilometry, 46.6
- 3.5-m telescopes, 5.43–5.45, 5.43*f*–5.45*f*
- Thulium-doped fibers, 25.23*t*, 25.25, 25.32
- Tight binding approximations, for photonic crystal fibers, 11.8
- Tilt:
- atmospheric, 5.14, 5.14*f*–5.15*f*
 - gradient (G-tilt), 4.3
 - and adaptive optics, 5.14–5.16
 - and angle of arrival, 4.23, 4.25, 4.26
 - wavefront (Z-tilt), 4.3, 4.23–4.25, 5.14, 5.15
 - Zernike, 4.23–4.26
- Tilt errors, for grazing incidence optics, 45.7
- Tilt-corrected phase variance, 4.31–4.32, 4.31*f*
- Tilted diffraction geometry, 42.4
- Time shifts, of solitons, 22.15, 22.16
- Time-division multiplexing (TDM), 9.12, 20.3, 21.3
- Time-domain spectroscopy (TDS) systems, 7.37, 7.37*f*–7.38*f*
- Time-of-flight (TOF) spectroscopy, 2.5
- Timing recovery, in OTDM networks, 20.10–20.12, 20.10*f*, 20.11*f*
- Tokens, of FDDI connectors, 23.3
- Tomography, 31.1, 31.5–31.7, 31.5*f*–31.7*f*
- Tomosynthesis, digital, 31.7–31.8
- Toroidal reflectors, for neutron beams, 64.4
- Toroidal-grating monochromators (TGMs), 38.3
- Total external reflection, 64.1–64.2

- Total integrated scatter (TIS):
of neutrons, **64.2–64.3**
and scatterometers, **1.4, 1.6–1.7, 1.10–1.11**
- Total internal reflection (TIR), of photonic crystal fibers, **11.2, 11.3, 11.14, 11.15**
- TRACE observatory, **41.3**
- Tracking, for adaptive optics, **5.15–5.18, 5.17f, 5.21–5.23, 5.22f**
- Transient response:
of light-emitting diodes (LEDs), **13.41–13.42**
- Transient response, of laser diodes, **13.13–13.18, 13.14f, 13.16f, 13.17f**
- Transition-edge sensor (TES) microcalorimeter detectors, **29.9–29.11, 29.9f, 29.11f**
- Transition-edge sensors (TES), **60.9**
- Transition-metal ions, spectra of, **2.8–2.9, 2.8f**
- Transmission:
of amplitude modulators, **7.22**
analog, **9.15–9.17, 9.16f**
atmospheric optical, **3.22–3.26, 3.22f–3.27f**
Bormann, **43.3, 43.3f**
broadband, **3.23–3.24, 3.25f–3.26f**
electrical time domain multiplexed, **20.3, 20.25**
formats for data, **20.9–20.10, 20.9f**
Gaussian, **37.5**
Laue, **26.8, 26.9f, 63.24, 63.26**
line-by-line, **3.23, 3.24f**
of NPM AOTFs, **6.41**
with optical fibers, **9.7–9.8, 9.15–9.17**
in OTDM networks, **20.7–20.8, 20.12–20.17, 20.14f–20.16f**
polycapillary x-ray optics, **53.5–53.8, 53.6f, 53.7f**
in polycapillary x-ray optics, **53.5–53.8, 53.6f, 53.7f**
for solitons, **22.5–22.7**
WDM dispersion managed soliton, **22.15–22.17**
of x-ray tube sources, **54.11–54.12**
- Transmission loss, in fiber optic links, **15.7**
- Transmissive thin-film transistor (TFT) LCDs, **8.29–8.31, 8.30f, 8.31f, 8.32–8.35, 8.34f**
- Transparent fiber systems, **14.1**
- Transreflective LCDs (TR-LCDs), **8.32–8.35, 8.33f, 8.34f**
- Transverse coherence, **55.18–55.20, 55.19f**
- Transverse effective wavelength, for PCFs, **11.10**
- Transverse electric (TE) polarization, **19.7**
- Transverse electric (TE) waveguide mode, of laser diodes, **13.13**
- Transverse electro-optic modulators, **7.16, 7.17, 7.17f**
- Transverse holographic technique, for fiber Bragg gratings, **17.5, 17.5f**
- Transverse magnetic (TM) mode, of laser diodes, **13.13**
- Transverse magnetic (TM) polarization, **19.7**
- Transverse ray aberration (TRA) equations, **45.3–45.8**
- Transverse spatial coherence, **55.16**
- Transverse spatial modulation (TSM), **6.11–6.12, 6.23, 6.30, 6.31**
- Trap-assisted thermal generation current, **13.69**
- Traveling wave electro-optic modulators, **7.28–7.30, 7.29f**
- Traveling wave pumping, by EUV lasers, **58.3**
- Traveling-wave amplification, **19.3, 19.4f**
- Tree topologies, of WDM networks, **21.6, 21.7**
- Trivalent rare-earth ions, **2.11**
- Tube voltage, of x-ray tube sources, **54.9**
- Tunable dispersion compensation, in WDM networks, **21.23–21.26, 21.24f–21.26f**
- Tunable filters, acousto-optic, **6.23t, 6.35–6.45**
collinear beam, **6.43, 6.43f, 6.45t**
and longitudinal spatial modulation, **6.12**
long-infrared, **6.42**
mid-infrared, **6.42**
noncritical phase-matching, **6.37, 6.39–6.42, 6.38f, 6.43t**
principle of operation, **6.36–6.39, 6.37f, 6.38f**
ultraviolet, **6.42**
- Tunable lasers, for fiber optic systems, **13.32–13.36, 13.33f–13.36f**
- Tungsten silicide/silicon (WSi₂/Si) bilayers, in MLLs, **42.6, 42.7**
- Tuning relation, of NPM AOTFs, **6.39–6.40**
- Turbulence:
and adaptive optics, **5.5–5.21**
anisoplanatism, **5.19**
atmospheric tilt and Strehl ratio, **5.14, 5.14f–5.15f**
Fried's coherence diameter and spatial scale, **5.9–5.13, 5.11t, 5.12f, 5.13f**
higher-order phase fluctuations, **5.18–5.19**
imaging, **4.35–4.36, 5.19–5.21, 5.19f, 5.20f**
Kolmogorov model, **5.5–5.6**
tracking requirements, **5.15–5.18, 5.17f**

- Turbulence, and adaptive optics (*Cont.*):
 variation of n and C^2n parameters,
 5.6–5.8, 5.7f, 5.8f
 in atmospheric optics, 3.26, 3.28–3.36
 beam spreading, 3.32–3.33, 3.33f
 beam wander, 3.31–3.32
 imaging and heterodyne detection, 3.34
 parameters for, 3.28–3.31, 3.29f, 3.30f
 scintillation, 3.34–3.36, 3.35f
- Hufnagel model of, 3.29
- Hufnagel-Valley model of, 3.30, 5.7, 5.8
- imaging through atmospheric, 4.1–4.37
 aberration variance and approximate
 Strehl ratio for, 4.27–4.28, 4.28f
 adaptive optics, 4.35–4.36
 angle of arrival fluctuations, 4.23–4.26,
 4.25f, 4.26f
 expansion coefficients, 4.20–4.22, 4.21t,
 4.22f–4.23f
- Kolmogorov turbulence and atmospheric
 coherence length, 4.7–4.10, 4.8f, 4.9f
 long-exposure images, 4.3–4.7
 modal correction of turbulence,
 4.28–4.30, 4.29t, 4.30f
 modal expansion of aberration function,
 4.17–4.20, 4.17f–4.18f, 4.19t, 4.20t
 and resolution of telescopes, 4.2–4.3
 short-exposure image, 4.31–4.35,
 4.31f–4.34f, 4.35t
 systems with annular pupils, 4.10–4.16,
 4.11f–4.15f, 4.15t
- inner scale of, 4.7
- Kolmogorov, 4.3, 4.7–4.10, 4.8f, 4.9f, 4.27,
 4.30, 4.36
- Kolmogorov model of, 5.5–5.6, 5.11
 optical strength of (C^2n), 5.6–5.8, 5.7f, 5.8f
 outer scale of, 4.7
 phase structure function of, 4.5
- Turn-on delay, of laser diodes, 13.13–13.14,
 13.14f
- Twist, of liquid crystals, 8.22, 8.23f
- Twist grain boundary (TGB) phases, of liquid
 crystals, 8.11
- Twisted nematic (TN) cells, 8.16, 8.25–8.26, 8.26f
- 2D profilometry, 46.6
- Two-wave interaction (acousto-optic), 6.8
- Tyler frequency, 5.17
- UHURU satellite, 47.1
- Ultrafast nonlinear interferometers (UNIs), 19.32
- Ultrahigh-speed OTDM, 20.23–20.24, 20.24f
- Ultrashort pulses (USPs), for fiber lasers,
 25.32–25.33
- Ultraviolet (UV) AOTFs, 6.42
- Umweganregung* (multiple-beam Bragg
 diffraction), 43.6–43.8, 43.7f
- Unbalanced nonlinear interferometers (UNI),
 20.22
- Undulators, 55.12–55.14, 55.12f, 55.13f, 58.1
- Unitraveling-carrier (UTC) photodiodes,
 13.68–13.69, 13.68f
- Unity divergence ratio, for AO modulators,
 6.32–6.33
- Unruh radiation, 58.2
- Upper Atmospheric Research Satellite (UARS),
 3.36, 3.37f
- Ur (code), 55.14
- Urgent (code), 55.14
- U.S. Air Force Cambridge Research
 Laboratories, 3.22
- U.S. Department of Defense, 5.27
- U.S. National Institute of Standards and
 Technology (NIST) database, 3.26
- U.S. Standard Atmosphere, 3.5, 3.6, 3.8f,
 3.9f
- Vacuum spark source, of z-pinch radiation,
 57.3, 57.3t
- Vacuum ultraviolet (VUV) region, gratings
 and monochromators in, 38.1–38.8, 38.2f,
 38.3f, 38.4t–38.5t
- Van der Waals theory, 8.23
- Vapor axial deposition (VAD), 25.3, 25.26
- Variable optical attenuators (VOAs), 21.12,
 21.13f
- Verdet constant, 18.7
- Vertical alignment (VA) cells, 8.25, 8.27–8.28,
 8.28
- Vertical cavity surface-emitting lasers
 (VCSELs), 13.1–13.2, 13.43f
 commercial, 13.48
 electrical injection and current confinement
 for, 13.45, 13.45f
 light out vs. current in (L-I curve) of,
 13.46–13.47
 and linear optical amplifiers, 19.14
 and mirror reflectivity, 13.44
 and optical fibers, 9.7n
 polarization of, 13.48

- Vertical cavity surface-emitting lasers (VCSELs)
 (*Cont.*):
 and quantum wells, 13.43
 and spatial characteristics of emitted light,
 13.46
 spectral characteristics of, 13.47–13.48
- Vertical coupling, of electroabsorption
 modulators, 13.60
- Virtual tributaries, of SONET, 23.6
- Viscosities, of liquid crystals, 8.23–8.25, 8.24*f*
- Voigt function, 56.7
- Voigt lineshape profiles, 3.23
- Voltage, of x-ray tube sources, 54.9
- Voltage-controlled oscillators (VCOs), 20.11,
 20.11*f*
- Volume Bragg gratings (VBGs), 25.29, 25.30
- Volume diffraction calculations, for MLLs,
 42.4–42.5, 42.5*f*
- Von Kármán spectrum, 3.28, 5.6
- V-parameter, of fiber lasers, 25.3
- W. M. Keck Observatory, 5.27
- Wadsworth condition, 38.6
- Wall effect, in gas detectors, 63.32
- Water:
 scattering in, 63.10, 63.10*t*, 63.11*f*
 in standard atmosphere, 3.6, 3.10–3.11, 3.11*f*
- Wave propagation, reciprocity of, 4.9–4.10
- Wave structure function, 4.6, 5.9*n*
- Wavefront correctors, for adaptive optics,
 5.37–5.38, 5.38*f*
- Wavefront errors, 4.35, 5.40–5.41
- Wavefront sensing techniques, 5.36–5.37
- Wavefront tilt (Z-tilt), 4.3, 4.23–4.25, 5.14, 5.15
- Wavefronts, gradients of, 5.23
- Waveguide confinement factor, 13.4*f*, 13.5
- Waveguide modulators, 7.30–7.32, 7.31*f*–7.33*f*,
 13.56, 13.57*f*
- Waveguides:
 attenuated total reflectance, 12.11
 double heterostructure, 19.3*f*
 evanescent-wave coupled pin, 13.68
 of SOAs, 19.15–19.16, 19.15*f*–19.16*f*
- Wavelength(s):
 attenuation vs., 15.7
 Bragg, 17.3, 20.15, 24.7
 of liquid crystals, 8.19–8.22, 8.20*f*
 pump, 14.5–14.6
 transverse effective, 11.10
- Wavelength blockers (WB), 21.12, 21.13*f*
- Wavelength dispersion, 9.6–9.7
- Wavelength dispersive detectors (WDS),
 62.2
- Wavelength division multiplexing (WDM):
 dense
 and AOTFs, 6.43, 6.44
 and optical fiber amplifiers, 14.1
 and SOAs, 19.25–19.27, 19.25*f*, 19.26*f*
 and dispersion-managed solitons,
 22.15–22.17
 and ESCON standard, 23.2
 and fiber-based couplers, 16.1, 16.4
 and solitons, 22.2, 22.8–22.12,
 22.15–22.17
 (*See also* Wavelength-division multiplexing
 networks)
- Wavelength scans, for scatterometers, 1.14
- Wavelength selective switches (WSS), 21.12,
 21.13*f*
- Wavelength-dispersive x-ray fluorescence
 (WDXRF), 29.2–29.3, 29.2*f*
- Wavelength-division multiplexing (WDM)
 networks, 18.4, 18.6, 21.1–21.44
 architecture of
 circuit and packet switching, 21.7–21.11,
 21.8*f*–21.11*f*
 network reconfigurability, 21.12–21.13,
 21.12*f*, 21.13*f*
 point-to-point links, 21.4
 star, ring, and mesh topologies, 21.5–21.7,
 21.5*f*–21.7*f*
 wavelength-routed networks, 21.5, 21.5*f*
 fiber bandwidth, 21.2–21.3, 21.2*f*
 fiber system impairments, 21.13–21.26
 chromatic dispersion, 21.14–21.16, 21.15*f*,
 21.16*f*
 dispersion and nonlinearities manage-
 ment, 21.20–21.26, 21.21*f*,
 21.23*f*–21.27*f*
 fiber attenuation and optical power loss,
 21.13–21.14
 fiber nonlinearities, 21.18–21.20, 21.19*f*,
 21.20*f*
 polarization-mode dispersion,
 21.16–21.18, 21.17*f*–21.18*f*
 history of, 21.1–21.2
 optical amplifiers in, 21.37–21.44, 21.37*f*
 EDFA, 21.38–21.41, 21.38*f*–21.42*f*
 Raman, 21.42–21.44, 21.42*f*–21.44*f*
 and optical fibers, 9.8, 9.12, 9.13, 9.15

- Wavelength-division multiplexing (WDM)
 networks (*Cont.*):
 optical modulation formats for, 21.27–21.36
 basic concepts, 21.27–21.29, 21.28f–21.30f
 carrier-suppressed return-to-zero and
 duobinary, 21.30–21.33, 21.31f,
 21.32f
 comparisons of, 21.36, 21.36t, 21.37t
 DPSK and DQSK, 21.33–21.36, 21.33f–
 21.35f, 21.36t, 21.37t
 in real systems, 21.3–21.4, 21.3f, 21.4f
 Wavelength-routed networks, 21.5, 21.5f
 Wavelength-selective couplers (WSCs), 14.2
 Waviness, in polycapillary x-ray optics, 53.4,
 53.4f
 Wedged multilayer Laue lenses (wMLLs),
 42.12–42.13, 42.13f, 42.14f
 Well and barrier intermixing, 13.60
 Wide area networks (WANs), 9.14, 21.7
 Wideband AO Bragg cells, 6.27, 6.30, 6.31t
 Wide-field lobster-eye telescopes, 48.4
 Wigglers, 55.11–55.12
 Window structure, of SOAs, 19.8, 19.8f
 Windowing, in x-ray mirror metrology, 46.9,
 46.10f
 Wire array sources, of z-pinch radiation,
 57.3–57.4, 57.3t, 57.4f
 Wollaston prism, 46.4
 Wolter geometries, 63.21
 Wolter mirror configurations, 52.4, 64.6, 64.6f
 Wolter optics, 26.10, 48.1, 48.2, 49.4–49.6,
 49.5f–49.6f
 Wolter telescopes, 44.4, 44.5f, 44.6–44.10, 44.7f,
 44.10f, 45.1–45.5, 45.2f
 Wolter x-ray optics, 47.2–47.7, 47.2f, 47.4f,
 47.6f, 47.7f
 Wolter-Schwarzschild (WS) telescopes, 44.7,
 44.11, 45.1
 XEUS/IXO mission, 49.7
 XMM mission, 44.6
 XMM-Newton observatory, 33.3, 33.3t, 47.2,
 47.4f, 47.6, 47.6f
 X-pinch sources of radiation, 57.3, 57.3t,
 57.4
 X-ray absorption fine structure (EXAFS), 29.3,
 60.4
 X-ray absorption near edge spectroscopy
 (XANES), 30.2, 30.2f, 30.4f
 X-ray absorption near-edge structure
 (XANES), 29.3, 60.4
 X-ray absorption spectroscopy (XAS),
 30.1–30.5, 30.2f–30.5f
 X-ray astronomy, 33.1–33.4, 33.3t
 X-ray astronomy satellite (SAX), 44.6
 X-ray detectors, 60.3–60.10
 cryogenic, 60.8–60.9, 60.9t, 60.10t
 film, 60.8, 60.9t, 60.10t
 ionization, 60.3–60.7, 60.9t, 60.10t
 scintillation, 60.7–60.8, 60.9t, 60.10t
 X-ray diffraction (XRD), 28.1–28.7,
 28.3f–28.6f
 X-ray fluorescence (XRF), 29.1–29.11, 54.8
 energy-dispersive, 29.3–29.11, 29.5f, 29.6f,
 29.8f–29.11f
 history of, 29.1
 and polycapillary x-ray optics, 53.10–53.11,
 53.11f
 wavelength-dispersive, 29.2–29.3, 29.2f
 and x-ray diffraction, 28.1
 and x-ray imaging, 62.5, 62.6f
 X-ray imaging detectors, 61.1–61.8, 61.2f
 CCD detectors, 61.7–61.8, 61.8f
 flat panel detectors, 61.3–61.7, 61.4f, 61.6t,
 61.7f
 X-ray lasers, 58.1–58.4, 58.3f
 X-ray mapping, 62.4, 62.5, 62.6f
 X-ray microscopes, 37.6
 X-ray mirrors:
 and grazing incidence optics, 44.3–44.6,
 44.3f–44.6f
 metrology of, 46.1–46.12
 history of, 46.1–46.2
 profile analysis considerations, 46.6–46.12,
 46.7f, 46.10f
 surface figure metrology, 46.3–46.6, 46.5f
 surface finish metrology, 46.2
 X-ray monochromators, 30.1–30.4,
 50.6–50.7
 X-ray observatories, 33.1–33.4, 33.3t
 X-ray optics:
 adaptive, 50.1–50.8
 hard vs. soft x-ray telescopes, 50.2
 history of, 50.1, 50.2f
 synchrotron and lab-based sources,
 50.2–50.8, 50.4f, 50.6f, 50.8f
 astronomical, 47.1–47.11
 angular resolution of, 47.10–47.11, 47.11f
 hard, 47.9–47.10

- X-ray optics, astronomical (*Cont.*):
 history of, 47.1–47.2
 Kirkpatrick-Baez optics, 47.7–47.8, 47.8f, 47.9f
 Wolter, 47.2–47.7, 47.2f, 47.4f, 47.6f, 47.7f
 coherent, 27.1–27.5, 27.3f–27.5f
 gratings and monochromators in, 38.1–38.8
 diffraction properties, 38.1–38.3, 38.2f
 dispersion properties, 38.6–38.7
 efficiency, 38.8
 focusing properties, 38.3–38.6, 38.3f, 38.4t–38.5t
 resolution properties, 38.7
 and inverse Compton x-ray sources, 59.1–59.4, 59.3t
 and medical imaging, 31.1–31.4, 31.2f–31.4f
 monocapillary, 52.1–52.6, 52.2f–52.5f, 52.2t
 multifoil, 48.1–48.4, 48.2f, 48.3f
 and neutron optics, 26.5–26.11, 26.8f, 26.9f, 26.11f, 36.2f
 polycapillary, 53.1–53.19, 53.2f
 alignment and measurement, 53.5–53.8, 53.6f, 53.7f
 collimation, 53.8–53.9, 53.8f–53.9f
 energy filtering, 53.10
 focusing, 53.9–53.10, 53.10t
 monochromatic imaging, 53.16–53.17, 53.17f
 powder diffraction, 53.14, 53.14f
 radiation resistance, 53.5
 scatter rejection in imaging, 53.14–53.16, 53.15f–53.16f
 scintigraphy, 53.17, 53.18, 53.18f
 simulations and defect analysis, 53.3–53.5, 53.3f–53.5f
 single crystal diffraction, 53.12–53.14, 53.12f, 53.13f
 therapy, 53.18–53.19
 x-ray fluorescence, 53.10–53.11, 53.11f
 and pore optics, 49.1–49.7
 ray tracing for, 35.1–35.6, 35.4f, 35.5f
 and Schwarzschild objective, 51.3, 51.3f–51.4f
 spectral detection and imaging in, 62.1–62.5, 62.2f–62.6f
 and x-ray properties of materials, 36.1–36.9
 Auger energies, 36.3t, 36.9t
 electron binding energies, 36.3t–36.6t
- X-ray optics, and x-ray properties of materials (*Cont.*):
 photoabsorption and scattering, 36.1
 photon energies, 36.7t–36.8t
 x-ray and neutron optics, 36.2f
- X-ray spectral detection and imaging, 62.1–62.5, 62.2f–62.6f
- X-ray tube sources, 54.3–54.16
 brightness and intensity of, 54.11–54.15, 54.13f–54.15f
 cathode design and geometry, 54.10–54.11
 characteristics of, 54.3, 54.4f
 optimization of, 54.15–54.16
 spectra of, 54.4–54.10, 54.5f, 54.8f
- X-rays:
 circularly polarized soft, 55.6–55.7
 nanofocusing of hard, 42.1–42.17
 history of, 42.2–42.4, 42.2f, 42.3f
 instrumental beamline arrangement and measurements for, 42.9–42.10, 42.9f–42.12f
 limitations of, 42.15–42.17, 42.16f–42.17f
 with magnetron-sputtered MLLs, 42.5–42.7, 42.6f–42.8f
 on MLLs with curved interfaces, 42.14, 42.15f
 Takagi-Taupin calculations for, 42.12–42.14
 volume diffraction calculations for, 42.4–42.5, 42.5f
 with wedged MLLs, 42.12–42.13, 42.13f, 42.14f
 polarization of, 43.1–43.2
- Y-branch interferometric modulators, 13.51–13.52
- Young's modulus, of infrared optical fibers, 12.3t, 12.9
- Ytterbium-doped fiber amplifiers (YDFAs), 14.2, 14.7
- Ytterbium-doped fibers, 25.23t, 25.24–25.25, 25.31, 25.33
- ZBLAN (fluorozirconate glass), 12.5, 12.5f
 and fiber lasers, 25.3, 25.24, 25.27t, 25.28
 fluoraluminate glass vs., 12.3t, 12.4, 12.4t
- Zeeman effect, 2.21
- Zeeman spectroscopy, 2.23–2.24, 2.24f
- Zemax (program), 35.1
- Zernike annular expansion coefficients, 4.25

- Zernike coefficients, variance of, 4.22, 4.22*f*, 4.23*f*, 4.36
- Zernike decomposition elements, 5.37
- Zernike modes of aberrations, 5.11, 5.11*t*, 5.12*f*
- Zernike polynomials, 4.17–4.20, 4.20*t*, 5.10, 46.6
- Zernike tilts, 4.23–4.26
- Zero order, of radiation, 40.1
- Zerodur glass shells, 47.5
- Zero-phonon lines, in solids, 2.13–2.14, 2.14*f*
- Zonal approach, to wavefront error correction, 4.35
- Zone plate law, 42.4
- Zone plates, 40.1–40.10
 - amplitude, 40.4–40.5, 40.5*f*
 - and Bragg-Fresnel lenses, 40.9–40.10, 40.9*f*
 - diffraction efficiencies of, 40.4–40.8, 40.5*f*, 40.7*f*, 40.8*f*
 - Fresnel, 40.2, 42.3, 42.3*f*, 55.16
 - geometry of, 40.1–40.3, 40.2*f*
 - neutron, 63.25
 - phase, 40.5–40.7, 40.7*f*
 - as thin lenses, 40.3–40.4
- Z-pinch plasma, 57.1–57.5, 57.2*f*, 57.3*t*, 57.4*f*
- Z-tilt (wavefront tilt), 4.3, 4.23–4.25, 5.14, 5.15

DO NOT DUPLICATE

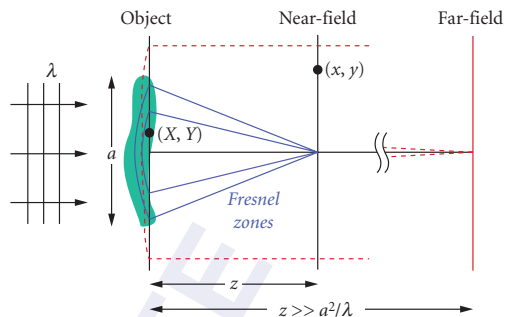


FIGURE 27.1 Schematic illustration of coherent x-ray wave propagation with a distorted object approach both for nearfield Fresnel diffraction, where an object extends into multiple Fresnel zones (solid lines), and for far-field Fraunhofer diffraction, where an object occupies only the center of the first Fresnel zone (dashed lines).

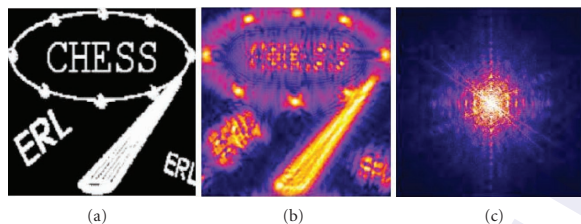


FIGURE 27.2 Simulated diffraction amplitudes $|F(x, y)|$, of an amplitude object (a) of $10 \mu\text{m} \times 10 \mu\text{m}$, with $l = 1 \text{ \AA}$ x rays, at image-to-object distance (b) $z = 2 \text{ mm}$ and (c) $z = \infty$, using the unified distorted object approach Eq. (3) with $Nz = 500$ zones in (b) and $Nz = 0$ in (c). Notice that the diffraction pattern changes from noncentrosymmetric in the near-field (b) to centrosymmetric in the far-field (c).

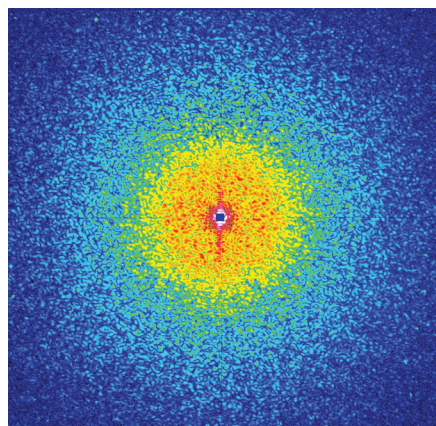


FIGURE 27.3 Example of a coherent x-ray diffraction pattern from a gold nanofoam specimen of $\sim 2 \mu\text{m}$ in size, using 7.35-keV coherent x rays. The corner of the image corresponds to $\sim 8 \text{ nm}$ spatial frequency.

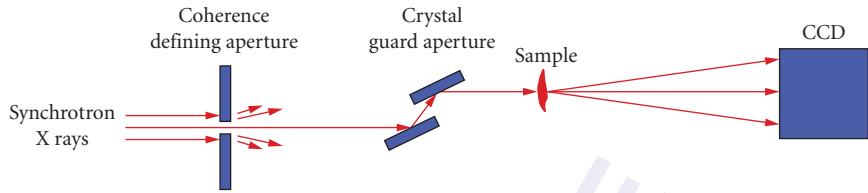


FIGURE 27.4 Concept of a perfect-crystal guard aperture in coherent diffraction imaging experiments for the purpose of eliminating unwanted parasitic scattering background in order to achieve high signal-to-noise in a diffraction pattern.

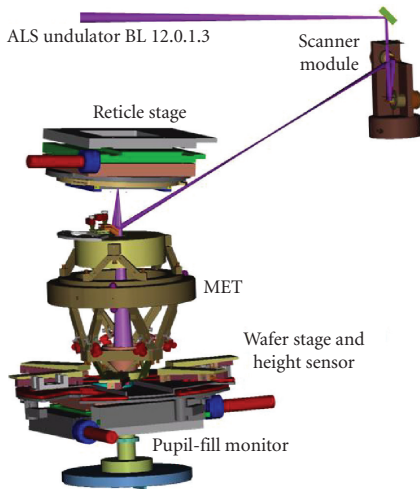


FIGURE 34.1 EUV exposure tool. The design includes 4 mirrors, and the mask (reticle) and wafer location. The whole system is under vacuum, and the condenser optical system is not shown. The whole system is enclosed in a vacuum chamber, not shown. The overall size is of several meters. This tool is installed at the Advanced Light Source, Berkeley.



FIGURE 34.2 ASML alpha demo tool during installation. The whole system is under vacuum, and a plasma source is used to generate the EUV radiation (not shown).

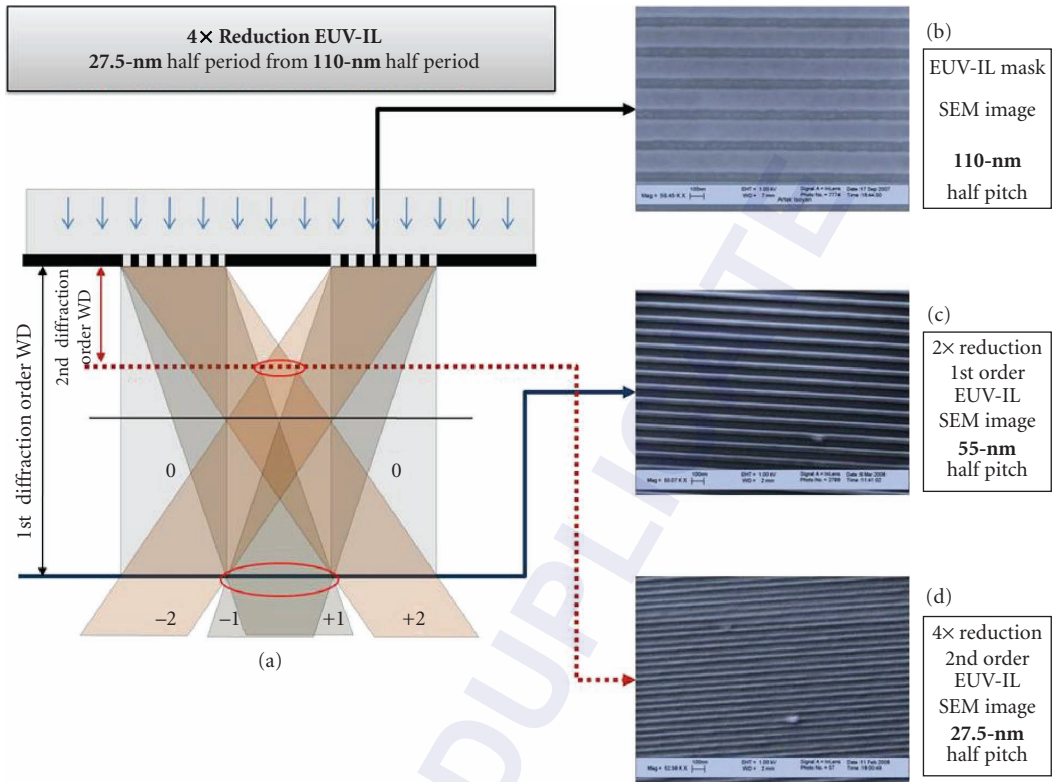


FIGURE 34.3 EUV interferometric lithography. The diffraction gratings are illuminated by a synchrotron, and the diffracted beams interfere as shown. The beams overlap creating 1st and 2nd order interference patterns of excellent visibility. Right, SEM images of the grating, and of the first- and second-order exposures. Notice the relative period of the images—the 1× period is half of that of the diffracting grating, and the 2× is 1/4.

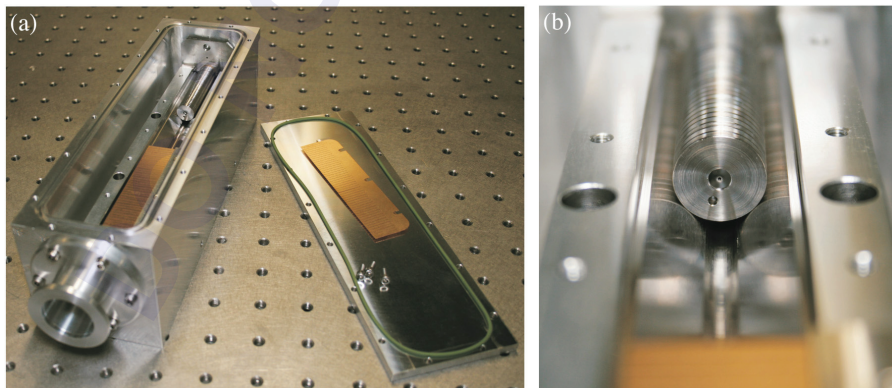


FIGURE 37.2 (a) Housing with partly assembled Be lens, (b) stack of Be lenses. Each individual lens is centered inside of a hard metal coin. They are aligned along the optical axis by stacking the coins in a high precision v-groove.

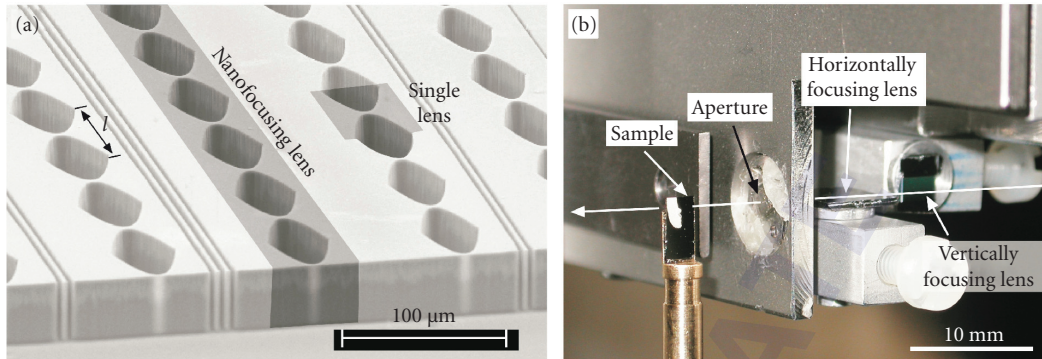


FIGURE 37.5 (a) Array of nanofocusing lenses made of silicon. A large number of single lenses are aligned behind each other to form a nanofocusing lens. Several nanofocusing lenses with different radius of curvature R are placed in parallel onto the same substrate. (b) Scanning microprobe setup with two crossed nanofocusing lenses. An aperture defining pinhole is placed behind the second lens.

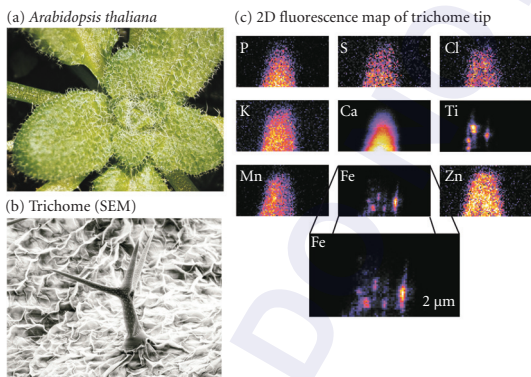


FIGURE 37.7 (a) Photograph of the plant *Arabidopsis thaliana*, (b) secondary electron micrograph of a leaf hair (trichome), (c) two-dimensional fluorescence map of the tip of the trichome at 100-nm spatial resolution. While most elements are homogeneously distributed, iron (Fe) and titanium (Ti) are localized on the level of 100 nm.

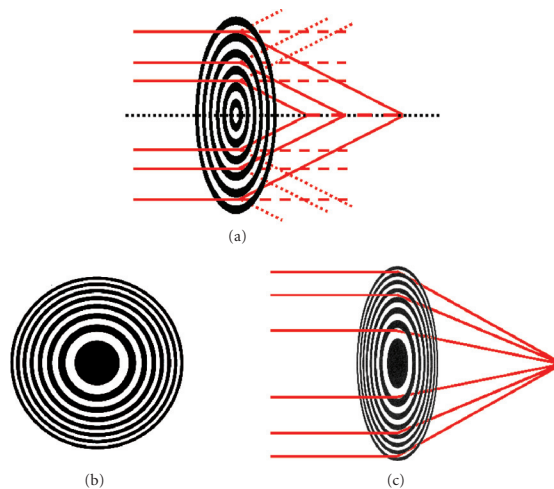


FIGURE 40.1 Diffraction by (a) a circular grating of constant period and (b, c) a zone plate.

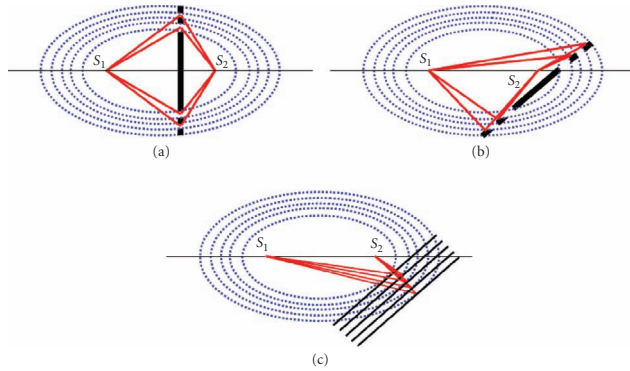


FIGURE 40.6 Structure of Bragg-Fresnel lenses.

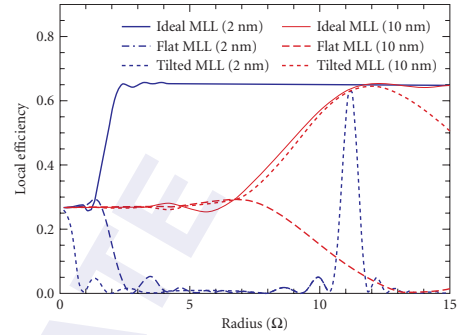


FIGURE 42.3 Calculated diffraction efficiency at 0.064 nm (19.5 keV) for outermost zone widths of both 10 nm (red) and 2 nm (blue) as a function of radius for ideal (wedged), flat, and tilted MLLs. For flat MLLs efficiencies do not exceed 26 percent and only very low-order zones diffract in the 2-nm case. For an outermost zone width of 10 nm, the ideal and tilted cases have almost the same performance, but for a 2-nm outermost zone, the ideal case is far superior. For the 10-nm tilted case, only a sharp Bragg peak is seen at the radius for which a Bragg condition is satisfied. This figure shows that meeting the Bragg condition everywhere becomes increasingly important for outermost zones less than 10 nm in order to ensure a high efficiency throughout most of the MLL.

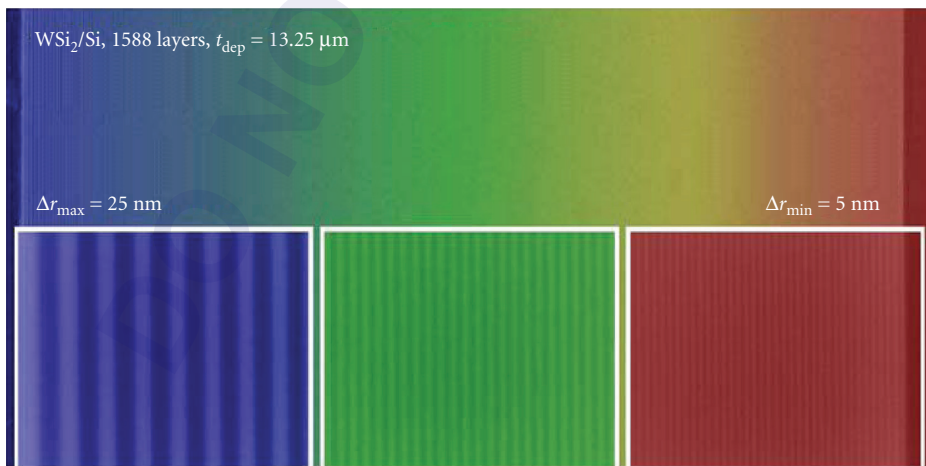


FIGURE 42.5 A scanning electron microscope image of a partial MLL having a 5-nm-thick outermost zone. The structure was used to obtain a focus of 16 nm for x rays of 0.064-nm (19.5-keV) wavelength.

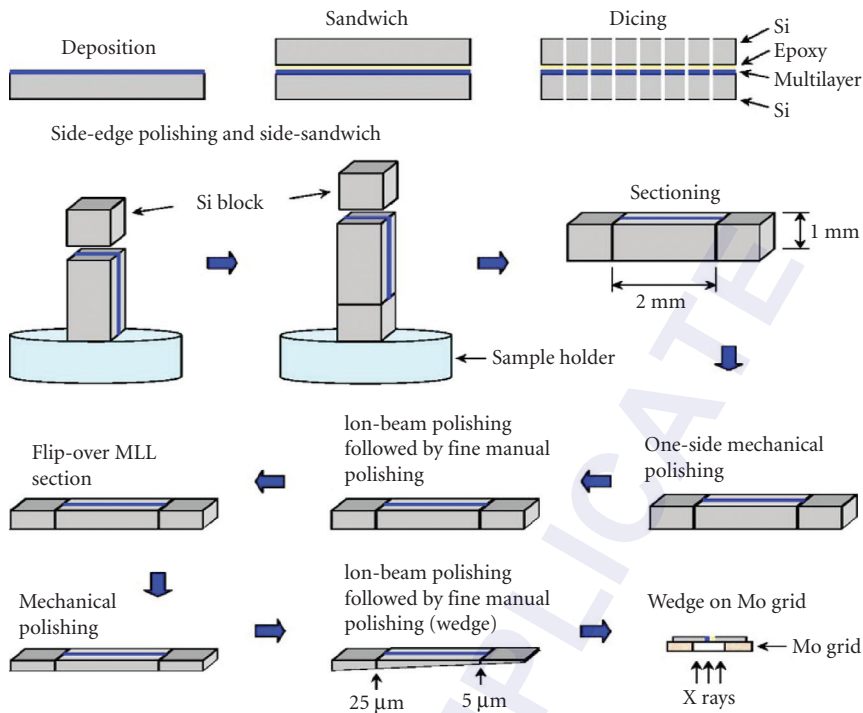


FIGURE 42.7 Cartoon illustrating the sequence of steps used to process a lens suitable for use with x rays, starting from an as-sputtered MLL wafer.

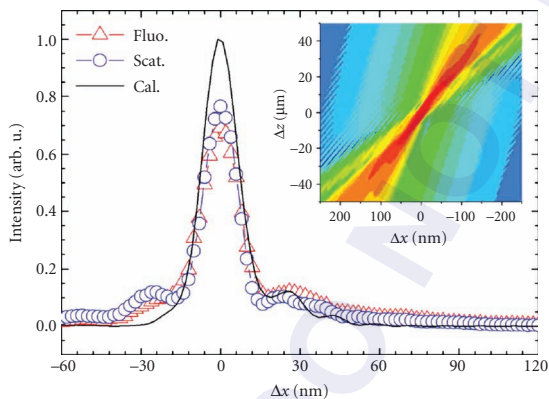


FIGURE 42.9 Measured and calculated intensity profiles for the focus of the MLL shown in Fig. 42.5. The FWHM values are 17.6 nm from fluorescence data and 15.6 nm from far-field scattering data. These should be compared to a calculated value of 15.0 nm. The calculated intensities were scaled according to measured and calculated efficiencies of 30 and 32 percent, respectively. The calculated results do not include the effects of vibrations, finite analyzer width, and finite sources size. The inset shows the calculated isophote pattern near the focal plane.

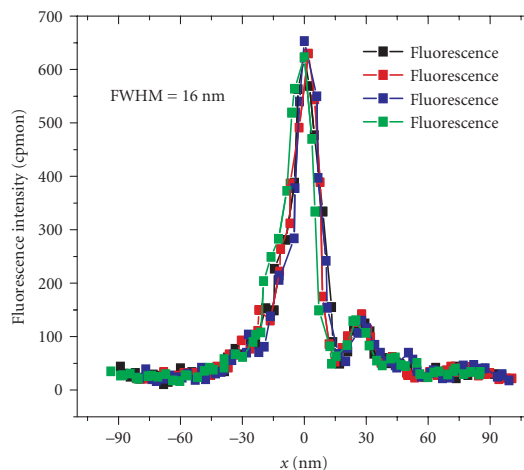


FIGURE 42.12 Measured line focuses of the MLL in Fig. 42.5 with 0.042-nm (19.25-keV) x rays. Results from four scans of the fluorescence intensity are shown. The efficiency was measured to be 17 percent. A FWHM value of 16 nm applies.

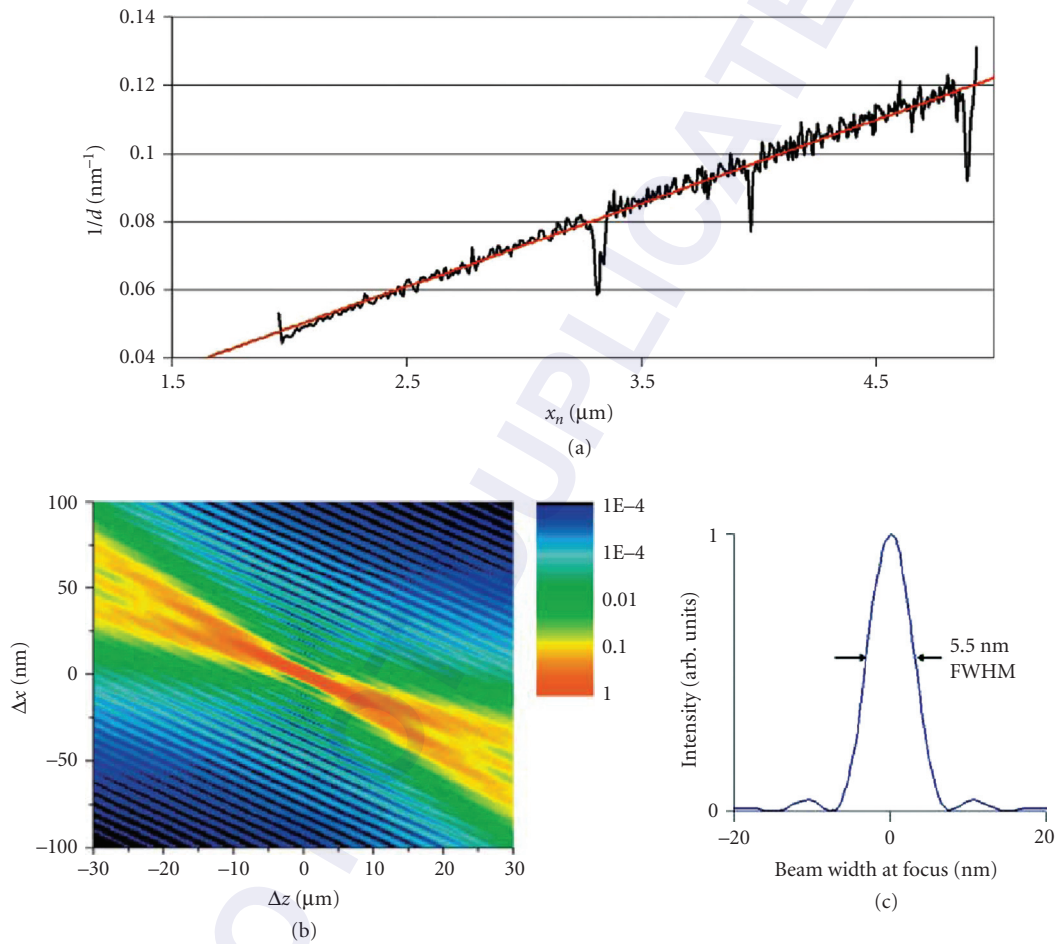


FIGURE 42.14 Analyses for ideal (wedged) MLL structure of Fig. 42.13a. The inverse layer thickness vs. radius is shown in (a). The calculated isophote pattern and intensity in the focal plane are shown in (b) and (c), respectively.

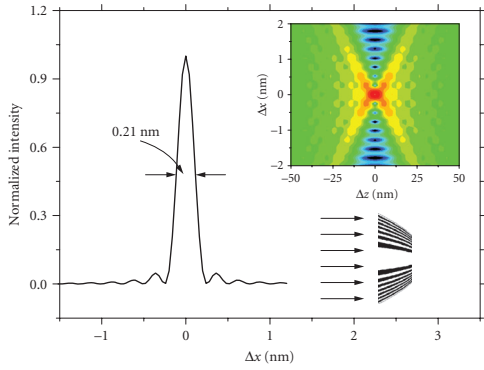


FIGURE 42.15 Calculated intensity in the focal plane for an MLL having parabolically shaped interfaces and an outermost zone width of 0.25 nm. The lower inset shows a cartoon of the MLL, and the upper inset shows an isophote pattern around the focus.

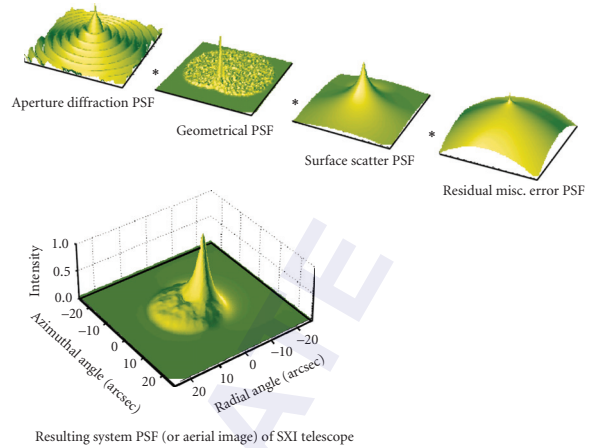


FIGURE 44.10 Illustration of the PSFs for the individual image degradation mechanisms and their resulting convolution (aerial image) of the SXI telescope for a field angle of 15 arcmin and a wavelength of 44.7 Å.

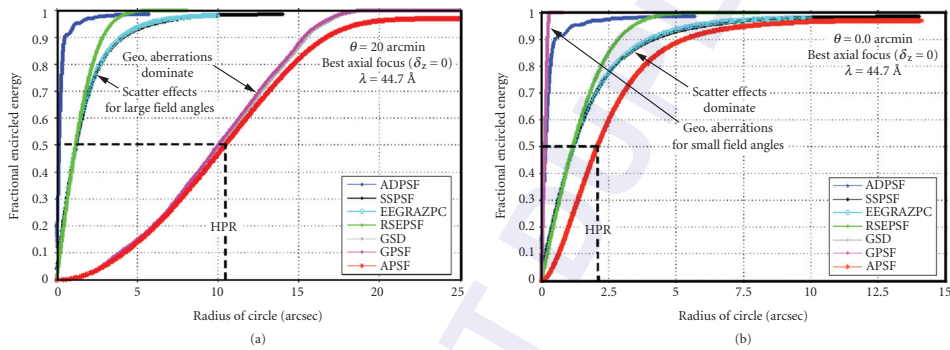


FIGURE 44.11 The fractional encircled energy of the aerial image and the four functions contributing to it provide insight into the image quality of a grazing incidence x-ray telescope. Note that scatter effects dominate geometrical aberrations for small field angles and geometrical aberrations dominate scatter effects for large field angles.

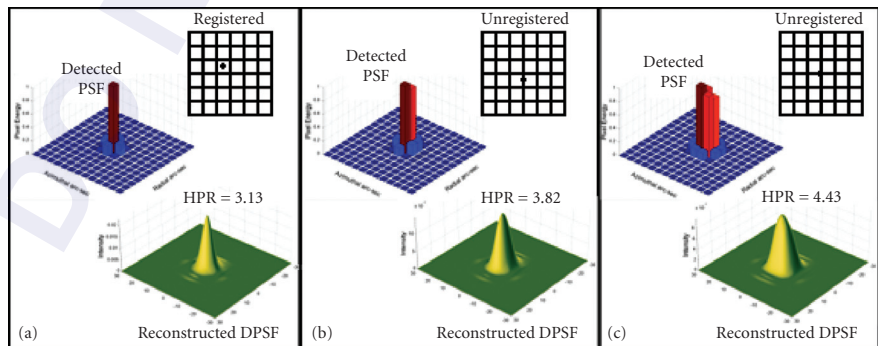


FIGURE 44.12 The DPSF and the reconstructed DPSF for (a) precisely “registered” aerial PSF, (b) aerial PSF centered on boundary between two pixels, and (c) aerial PSF positioned where four pixels meet.

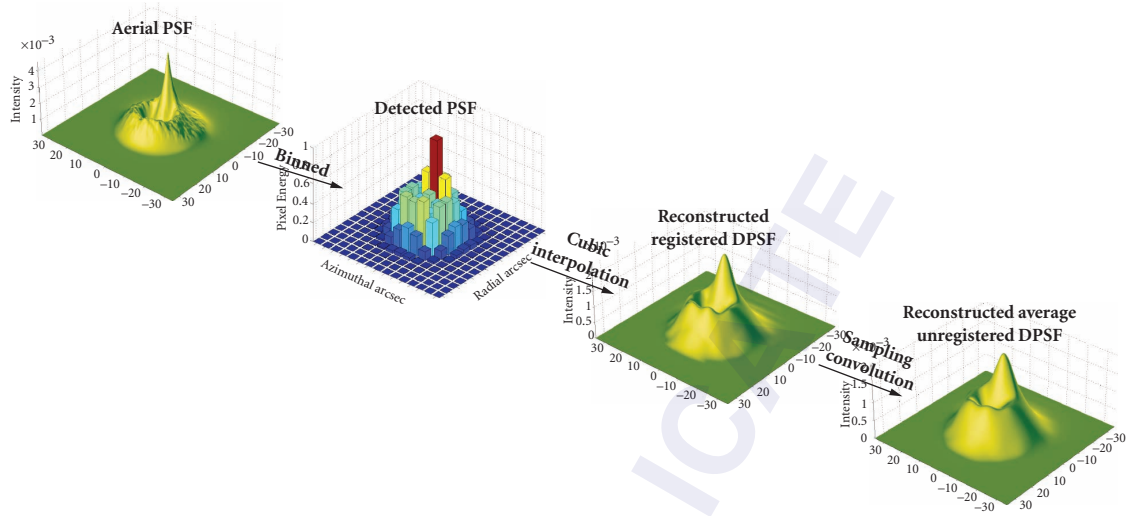


FIGURE 44.13 A graphical illustration of the numerical computation technique for modeling both the reconstructed “registered” DPSF and the reconstructed “average unregistered” DPSF is indicated.

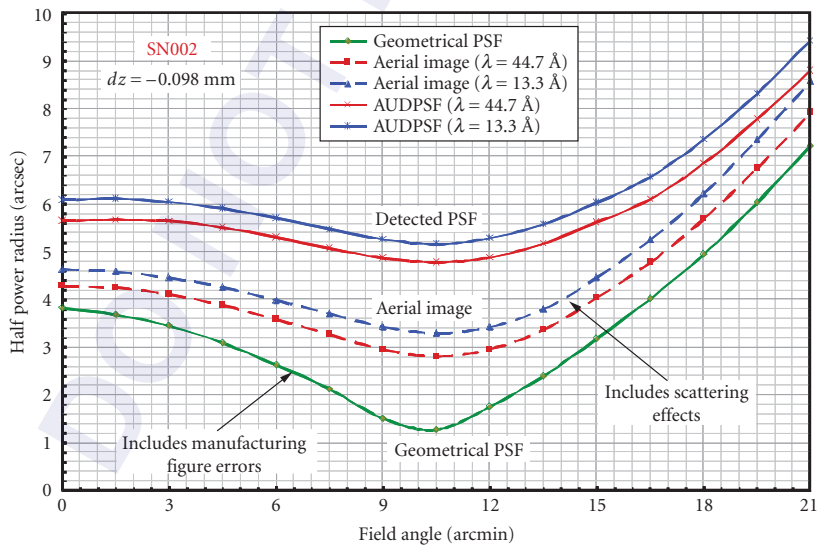


FIGURE 44.14 Comparison of the predicted HPR versus field angle of the geometrical PSF, the aerial image, and the average unregistered detected PSF for the SXI telescope at two different wavelengths (44.7 and 13.3 Å).

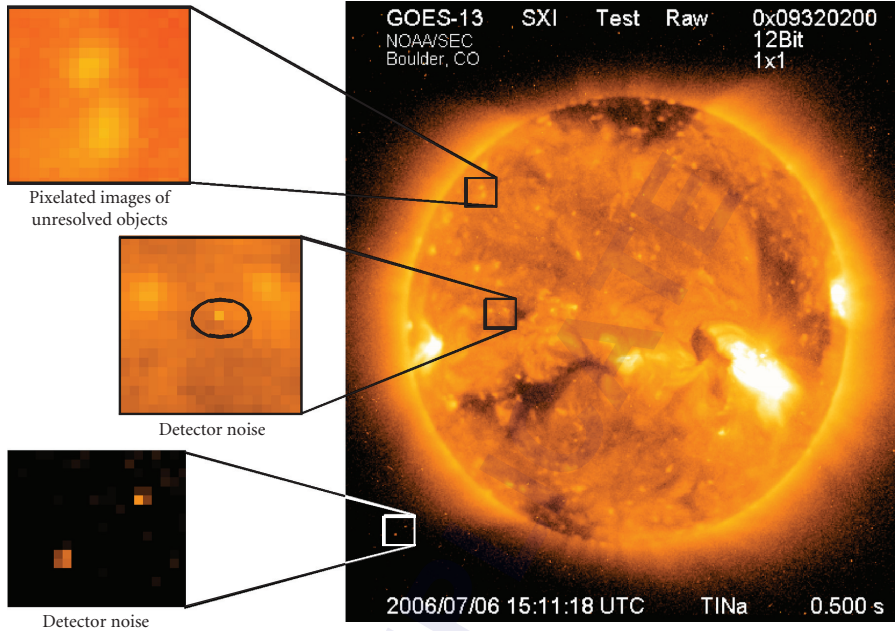


FIGURE 44.15 On-orbit solar image with three small regions highlighted and magnified for detailed inspection. This allows one to distinguish between images of unresolved bright features on the sun and merely detector noise.

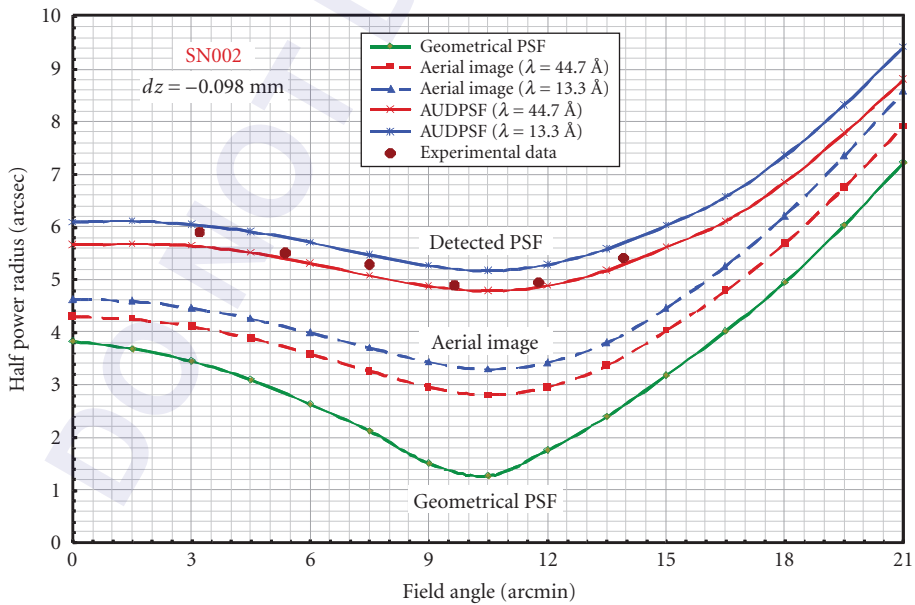


FIGURE 44.16 Experimental validation of an exhaustive systems engineering analysis of image quality for grazing incidence x-ray telescopes, including the modeling of surface scatter and detector effects.

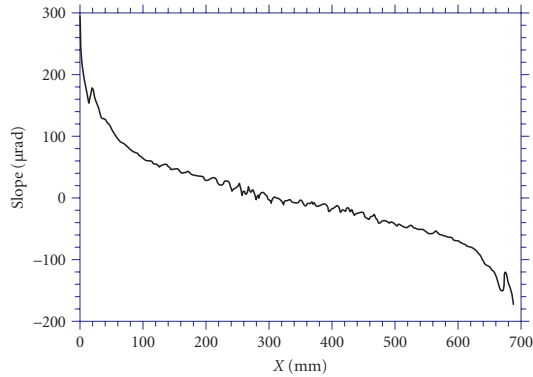


FIGURE 46.2 Slope profile measurement on a 700-mm long silicon cylinder mirror made with an LTP. Sampling step size is 2 mm. Mean has been subtracted from the data (detrrend 0). Profile shows that the surface has an overall convex curvature (tilted profile down to the right) with significant edge roll-off (change in slope at each end). Also, a polishing defect with a 20-mm period is evident in the center of the surface. The slope profile emphasizes high-frequency surface defects.

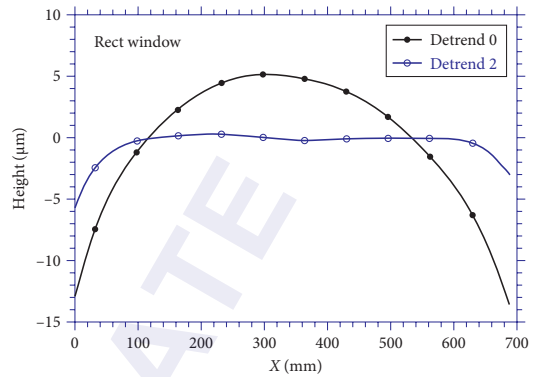


FIGURE 46.3 Height profile calculated by integrating the slope profile of Fig. 46.1. Solid circles: mean height subtracted (detrrend 0); open circles: second order polynomial subtracted (detrrend 2). The radius of curvature extracted from the second order term coefficient is 3.572 km. The residual profile shows that the surface has a “kink” in the center that separates it into two distinct segments with slightly different slopes. This low frequency defect is not evident in the slope profile of Fig. 46.1.

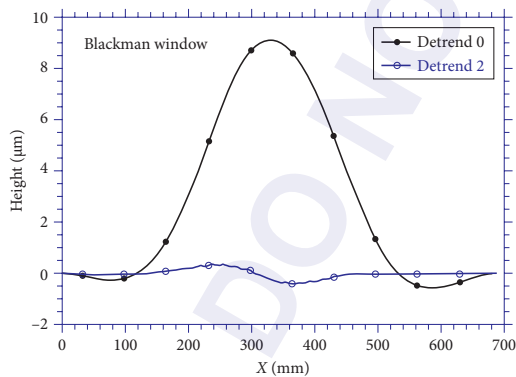


FIGURE 46.4 The height profiles of Fig. 46.3 with a Blackman window applied. The edge discontinuities are minimized by this function. Although the shape of the profile is distorted, the average statistical properties of the underlying function are not changed.

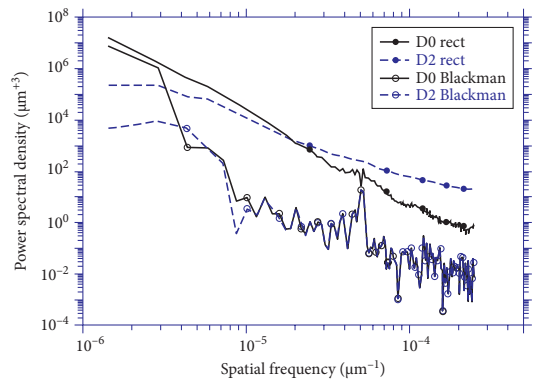
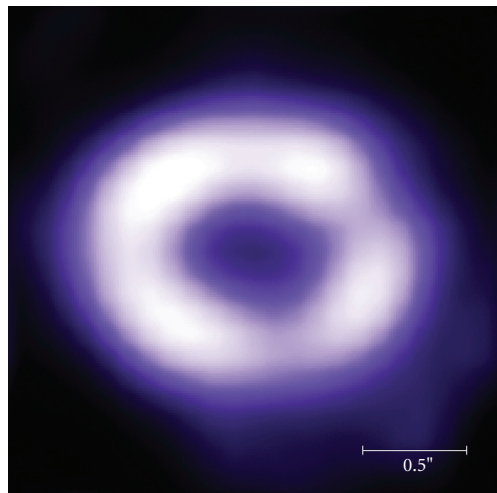
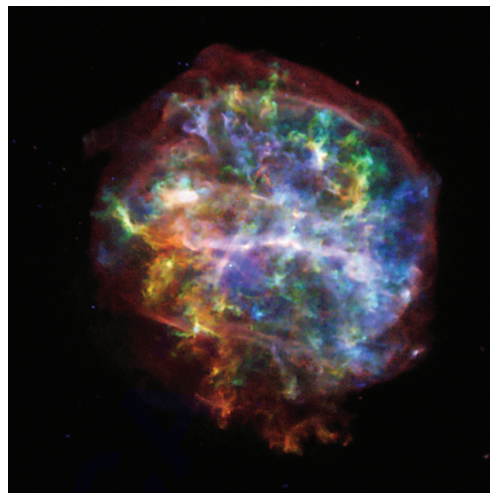


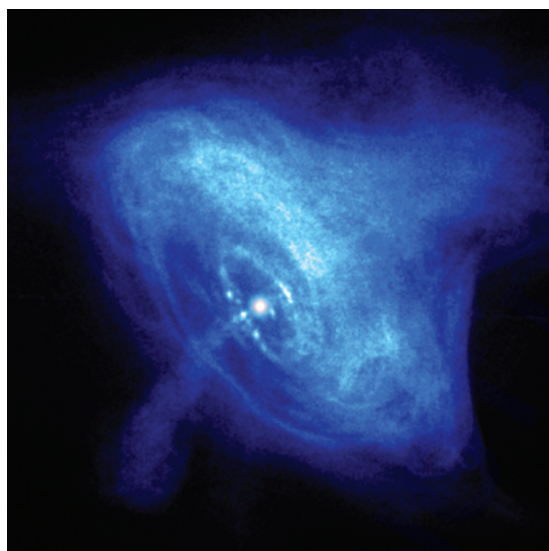
FIGURE 46.5 PSD curves computed from the four profiles in Figs. 46.3 and 46.4. The two upper curves from the unwindowed data show severe contamination effects due to the strong edge discontinuities that introduce spurious power into all frequencies. The lower curves show how the Blackman filter eliminates the discontinuity contamination, allowing the underlying surface spectral characteristics to become visible.



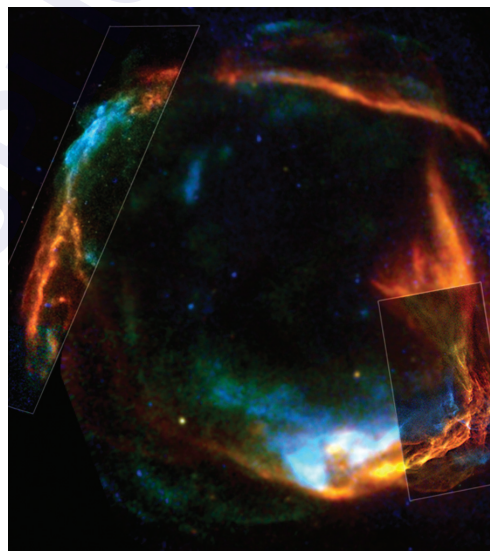
(a)



(b)



(c)



(d)

FIGURE 47.2 (a) A Chandra x-ray image of the fiery ring surrounding the 1987 supernova explosion in the Dorado constellation. Subarcsecond angular resolution is required to resolve the structure surrounding the supernova remnant (<http://chandra.harvard.edu/photo/2005/sn87a/>). (NASA/CXC/PSU/S.Park & D. Burrows.) (b) A Chandra x-ray image of the supernova remnant G292.0+1.8. The colors in the image encode the x-ray energies emitted by the supernova remnant; the center of G292.0+1.8 contains a region of high energy x-ray emission from the magnetized bubble of high-energy particles that surround the pulsar, a rapidly rotating neutron star that remained behind after the original, massive star exploded (<http://chandra.harvard.edu/photo/2007/g292/>). (NASA/CXC/Penn State/S.Park et al.) (c) Chandra x-ray image of the Crab Nebula—the remains of a nearby supernova explosion which was seen on Earth in 1054 AD. At the center of the bright nebula is a rapidly spinning neutron star, or pulsar, that emits pulses of radiation 30 times a second (<http://chandra.harvard.edu/photo/2002/0052>). (NASA/CXC/ASU/J. Hester et al.) (d) XMM-Newton and Chandra x-ray images RCW 86, an expanding ring of debris that was created after a massive star in the Milky Way collapsed onto itself and exploded. Both the XMM-Newton and Chandra images show low-energy x-rays in red, medium energies in green and high energies in blue. The Chandra observations focused on the northeast (left-hand) and southwest (lower right) side of RCW 86, and show that x-ray radiation is produced both by high-energy electrons accelerated in a magnetic field (blue) as well as heat from the blast itself (red). These images demonstrate the large field of view and moderate angular resolution of XMM-Newton, compared to the smaller field of view and high angular resolution provided by Chandra (<http://chandra.harvard.edu/photo/2007/2snr/>). (NASA/CXC/ESA/Univ. of Utrecht/J. Vink et al.)

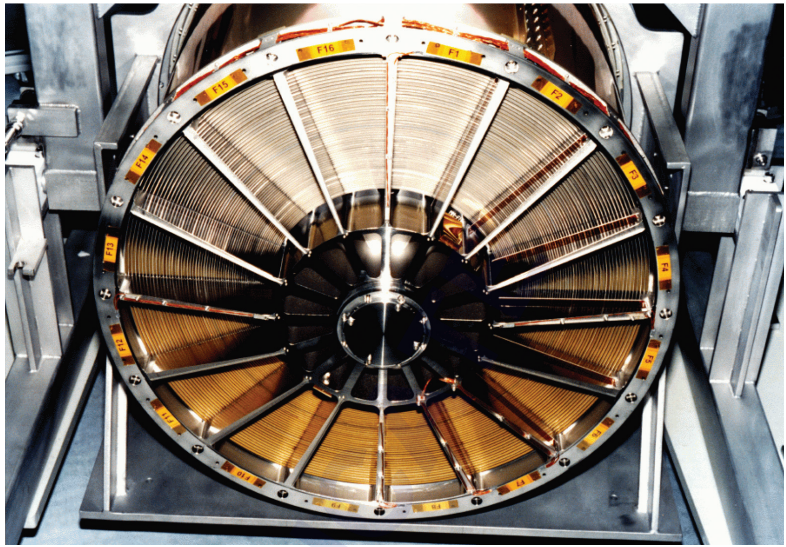
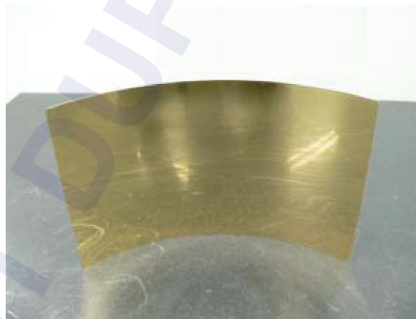


FIGURE 47.3 An x-ray optics module for the XMM observatory. Fifty-eight electroformed nickel Wolter I optics are nested to increase the effective x-ray collecting area.



(a)



(b)

FIGURE 47.4 The segmented foil mirrors aboard the SUZAKU spacecraft. (a) A complete mirror module and (b) a single aluminum foil reflector coated with gold.

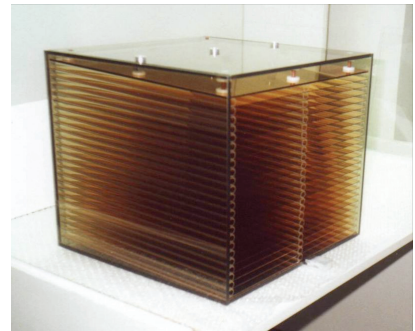


FIGURE 48.3a Kirkpatrick-Baez test Au coated glass foils system studied for XEUS project.

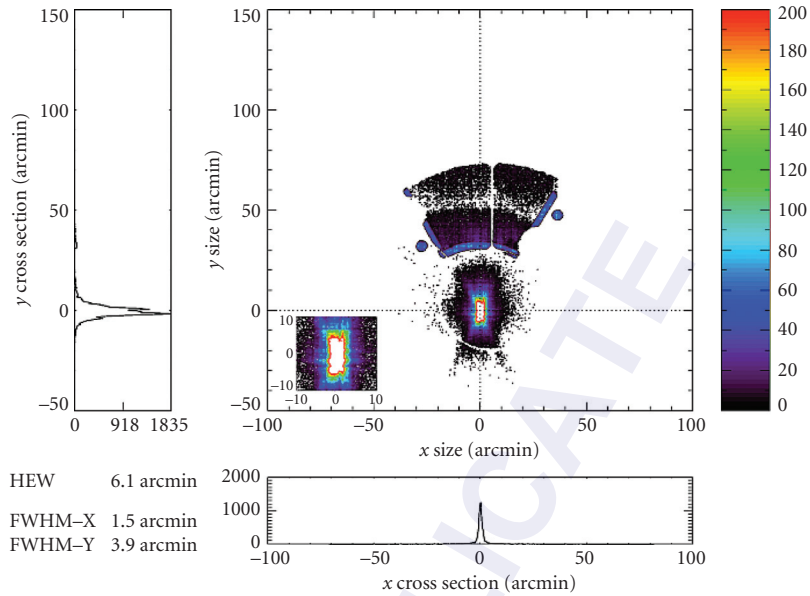


FIGURE 49.7 The focus of a single plate radially packed microchannel plate optic.

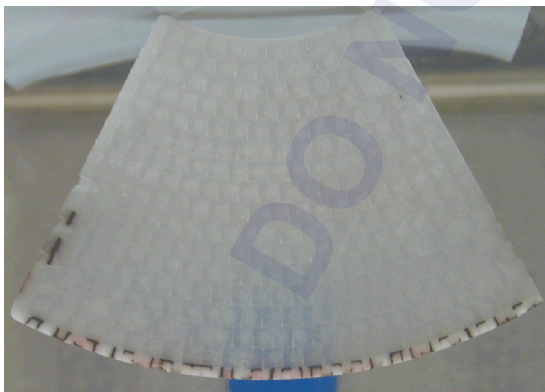


FIGURE 49.9 A segment of a radially stacked micropore optic. Two such plates behind each other act as a conical approximation to a Wolter-I optic. The width of the segment is 35 mm.

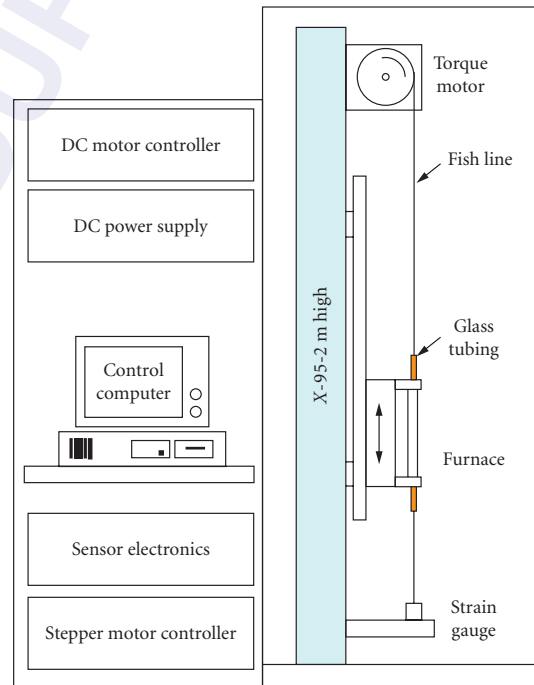


FIGURE 52.4 A glass tube is suspended in an electric furnace from a piece of fish line that is attached to a strain gauge at the bottom. The torque motor keeps a constant tension as the glass yields during drawing. The furnace is programmed to move based on the amount of glass yielding to make the desired elliptical, parabolic, etc., shape.

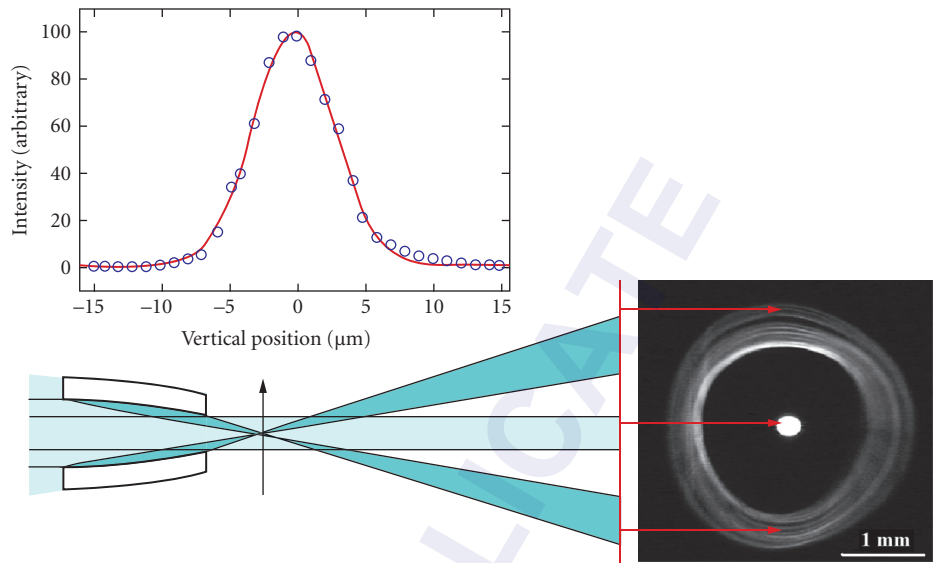


FIGURE 52.5 Upper panel: Profile of intensity versus $5\ \mu\text{m}$ vertical pinhole position at focus of a 9 milliradians capillary producing a spot size of $5\ \mu\text{m}$ FWHM at a distance of 20 mm beyond the tip of the capillary. Lower panel: The far-field image shows the direct beam (center dot) passing through the capillary and the once-reflected beam forming the outer ring of intensity. The structure in the ring is due to slope-error imperfections arising from the pulling process.

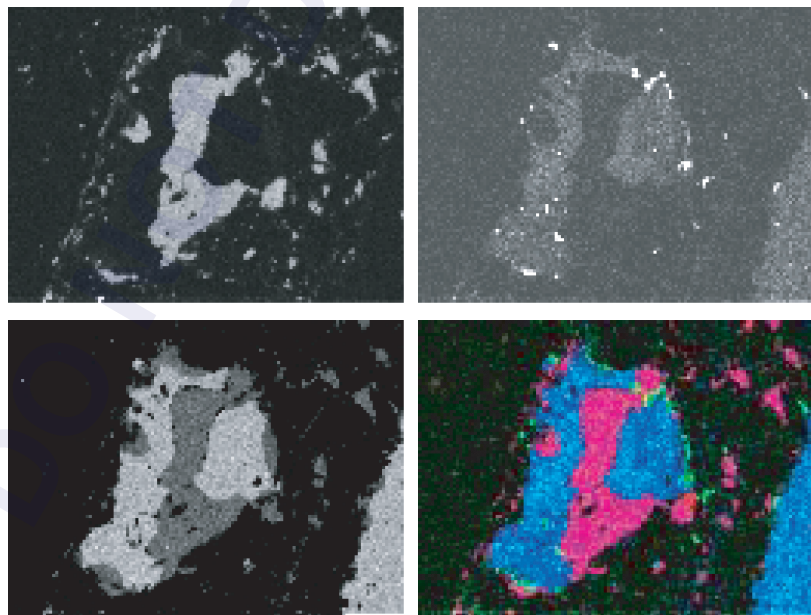


FIGURE 62.3 A 12-s x-ray spectrum image (128×96 pixels; 1 ms) obtained with a silicon drift detector.

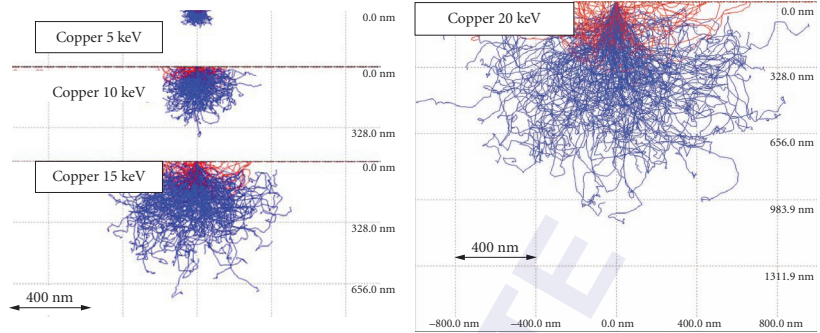


FIGURE 62.4 Monte Carlo simulations of electron trajectories in copper at various beam energies. Note the significant increase in the scattering volume with increasing energy.

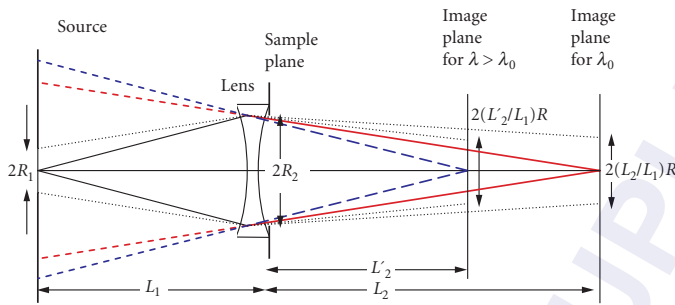


FIGURE 63.6 A schematic diagram of a focusing lens arrangement with source of radius R_1 at a distance L_1 from the sample with an aperture of radius R_2 . The focusing lens with a focal length f_0 is placed in front of the sample such that the source is imaged (continuous lines) with a radius $(L_2/L_1)R_1$ at a distance L_2 from the lens for a wavelength λ_0 . For another wavelength $\lambda (> \lambda_0)$ the source is imaged (long dashed lines) with a radius $(L'_2/L_1)R_1$ at a distance $L'_2 (< L_2)$ such that $1/L'_2 - 1/L_2 - 1/f_0[(\lambda/\lambda_0)^2 - 1]$.

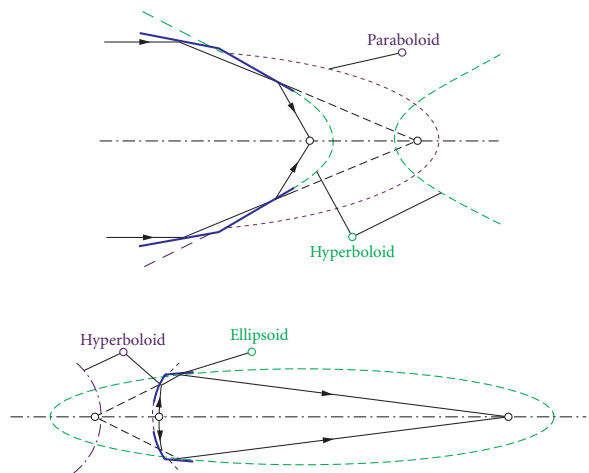


FIGURE 64.5 Commonly used Wolter-1 mirror configurations.

INDEX

Index note: The *f* after a page number refers to a figure, the *n* to a note, and the *t* to a table.

- Abbe illuminated eyepieces, **II**:12.12, 12.12*f*
Abbe illumination system, **II**:39.23, 39.23*f*, 39.34, 39.35*f*
Abbe numbers, **I**:29.36; **IV**:2.23, 2.28, 2.29*f*
and axial chromatic aberrations, **I**:17.22
of binary optics, **I**:23.6, 23.6*f*
in gradient index optics, **I**:24.3, 24.7
of molded microlenses, **I**:22.12, 22.12*t*
of reflective and catadioptric objectives, **I**:29.17
Abbe's sine condition, **I**:29.36; **II**:34.19; **V**:45.3
and reflective and catadioptric objectives, **I**:29.34
of stigmatic conditioning, **I**:1.30–1.31, 17.10
Abbe-Porter experiments, **I**:11.1
Abbe's prisms, **I**:19.3*t*, 19.7, 19.7*f*–19.8*f*
ABC model, of surface finish, **I**:8.14–8.15
Abel transform, **I**:8.13
Abelès method, **I**:12.11–12.12
Aberrated wavefronts, **I**:2.12, 2.13
Aberration coefficients, **II**:3.10–3.11
Aberration control, **IV**:3.8
Aberration curves (in lens design), **II**:2.1–2.6
considerations for, **II**:2.5–2.6
field plots of, **II**:2.4–2.5
transverse ray plots of, **II**:2.2–2.4
Aberrations (in general):
balancing of, **II**:11.30, 11.35–11.36, 11.36*t*
in binary optics, **I**:23.4–23.7, 23.5*f*, 23.6*f*
chromatic, **I**:1.91–1.92
and contact lenses, **III**:20.24–20.25, 20.25*f*
defined, **I**:1.28
defocus as, **I**:1.85–1.86
evaluation of, **II**:3.9–3.11
and forward scattering of light, **III**:1.21
and general aspheres, **I**:29.3
in gradient index optics, **I**:24.3
of gratings and monochromators, **V**:38.7
in grazing incidence optics, **V**:45.1–45.8, 45.2*f*
Aberrations (in general) (*Cont.*):
in grazing incidence telescopes, **V**:44.6–44.12, 44.7*f*, 44.9*f*–44.11*f*
in grazing-incidence neutron optics, **V**:64.4, 64.6
higher-order, **I**:29.37
and imaging through atmospheric turbulence, **V**:4.17*f*–4.18*f*, 4.17–4.20, 4.19*t*, 4.20*t*, 4.27*f*, 4.27–4.30
in instrumental optics, **I**:1.85
in intraocular lenses, **III**:21.9–21.10, 21.10*f*
of point images, **I**:1.85–1.92, 1.88*f*
polarization, **I**:15.35*f*–15.37*f*, 15.35–15.37
pupil, **I**:1.76
ray, **I**:1.87–1.88, 1.88*f*
Seidel, **I**:29.38
spherical, **I**:1.90, 29.7, 29.8, 29.15, 29.21, 29.37, 29.38
and Stiles-Crawford effect, **III**:8.8
and stop position, **I**:1.92
and stop size, **I**:1.92
in systems with rotational symmetry, **I**:1.89–1.90
third-order, **I**:1.90–1.91, 29.38
transverse primary chromatic, **I**:17.22
wavefront, **I**:1.86*f*, 1.86–1.88
Zernike modes of, **V**:5.11, 5.11*t*, 5.12*f*
(*see also specific aberrations, e.g.*: Axial chromatic aberration)
Aberrations (human eye), **III**:1.3, 16.6
absence of, **III**:1.12–1.14, 1.13*f*, 1.14*f*
and AO-controlled light delivery
to generate aberrations, **III**:15.24
longitudinal chromatic aberration, **III**:15.22
transverse chromatic aberration, **III**:15.22–15.23
astigmatism, **III**:15.2

- Aberrations (human eye) (*Cont.*):
 chromatic, **III**:1.19*f*, 1.19–1.20
 and accommodation, **III**:1.34
 age-related, **III**:14.14
 correcting color coordinates for, **III**:10.17
 correction of, **III**:1.25–1.26, 1.26*f*
 longitudinal, **III**:15.22
 and macular pigment, **III**:1.9
 transverse, **III**:15.22–15.23
 with visual instruments, **III**:1.28
 control of, **III**:1.3
 correction of, **III**:1.25–1.26, 1.26*f*, 10.17
 defocus, **III**:15.2
 and depth of focus, **III**:1.28
 higher-order, **III**:16.6, 16.7
 and idiosyncratic peculiarities, **III**:1.6
 image quality
 for aberration-free eye, **III**:1.12–1.14
 calculating, **III**:1.21–1.22
 monochromatic, **III**:1.4, 1.14–1.19
 age-related, **III**:14.12–14.14, 14.13*f*
 correction of, **III**:1.25, 1.26, 1.26*f*
 off-axis, **III**:1.18*f*, 1.18–1.19
 on the visual axis, **III**:1.15–1.18
 and observed optical performance,
III:1.23–1.26
 properties of, **III**:15.4–15.7, 15.5*f*, 15.6*f*
 and pupil diameter, **III**:1.8
 and refractive surgery, **III**:1.15
- Aberrometers, **III**:1.15, 1.23, 12.6
- Ablation rate (refractive correction), **I**:28.54;
III:16.16–16.18, 16.18*f*
- Abney's law, obedience to, **III**:10.44, 11.37
- Above threshold ionization (ATI),
IV:21.14–21.17
 plateau extension of, **IV**:21.19
 quasi-classical, **IV**:21.16*f*, 21.16–21.17
 relativistic electron, **IV**:21.20, 21.21*f*
 strong field interactions with atoms,
IV:21.14–21.17, 21.15*f*, 21.16*f*
- Abrasion resistance, of antireflection coatings,
IV:7.31–7.32, 7.32*f*
- Abrasive forming, polymer, **IV**:3.11–3.12
- Absolute detectors, **II**:34.27–34.30
 electrical substitution radiometers,
II:34.27–34.29
 photoionization devices, **II**:34.29
 predictable quantum efficiency devices,
II:34.29–34.30
- Absolute instruments, **I**:1.29
- Absolute measurements, **II**:34.20–34.37
 accuracy and traceability of, **II**:34.21
 error propagation in, **II**:34.22
 error types in, **II**:34.21–34.23
 relative vs., **II**:34.20–34.21
 and uncertainty estimates, **II**:34.21–34.23
- Absolute method (scatterometer calibration),
V:1.15
- Absolute responsivity units (A/W), **II**:34.31
- Absolute sources (of radiation), **II**:34.23–34.27
 blackbody radiator, **II**:34.23–34.24
 blackbody simulators, **II**:34.24–34.26
 synchrotron radiation, **II**:34.26–34.27
- Absorbance, water, **IV**:1.10
- Absorbers:
 black, **IV**:7.104
 saturable, **IV**:18.5–18.11
 fast, **IV**:18.9–18.10
 self-amplitude modulation, **IV**:18.5–18.7,
 18.6*f*, 18.7*f*
 semiconductor saturable absorber mirrors,
IV:18.3, 18.10–18.11
 slow, **IV**:18.7–18.9, 18.8*f*
- Absorbing compounds (in black surfaces),
IV:6.15
- Absorbing media:
 antireflection coatings for, **IV**:7.26, 7.27
 in multilayer reflectors, **IV**:7.37–7.38, 7.38*f*
 in photodetectors, **II**:26.4*f*, 26.5
 radiant power transfer through, **II**:34.13
- Absorbing substrate (AS) chips, **II**:17.7, 17.7*t*
- Absorptance:
 defined, **II**:35.4
 measurement of, **II**:35.10
 of metals, **IV**:4.6, 4.39, 4.40*f*–4.42*f*, 4.48, 4.49
 and emittance, **IV**:4.49, 4.49*f*, 4.50*t*, 4.51*t*
 and mass attenuation coefficients for
 photons, **IV**:4.48*t*
 of optical coatings, **IV**:7.12–7.13
 in thermal detectors, **II**:28.2
 and transmittance/reflectance, **II**:35.7, 35.8,
 35.8*t*
 of water, **IV**:1.5*t*, 1.9
- Absorption:
 and atmospheric optics, **V**:3.4*f*, 3.4–3.5
 bio-optical models for, **IV**:1.27*f*, 1.27–1.29,
 1.28*t*
 of Cr³⁺, **V**:2.19*f*, 2.19–2.21, 2.20*f*
 cutoff filters based on, **IV**:7.60, 7.60*f*
 defect-related, **IV**:5.37–5.39, 5.38*f*, 5.39*f*

- Absorption (*Cont.*):
 defined, **II**:35.4
 by detritus in water, **IV**:1.26–1.27
 in dielectrics, **IV**:4.4
 direct interband, **IV**:8.27–8.28
 by dissolved organic matter, **IV**:1.22–1.23, 1.23*t*
 electro-, **V**:13.55–13.56, 13.56*f*, 13.58
 excited state, **V**:14.6
 fundamental edge of (*see* Fundamental absorption edge)
 impurity-related, **IV**:5.37–5.39, 5.38*f*, 5.39*f*, 5.51*f*, 5.51–5.52
 interband (*see* Interband absorption)
 lattice, **IV**:5.13–5.20, 5.15*f*, 5.17*t*, 5.18*f*–5.19*f*, 5.20*t*, 5.21*f*
 lenticular, **III**:1.9
 of light in laser cooling, **IV**:20.4
 magnetoabsorption, **IV**:5.51*f*, 5.51–5.52
 measurement of, **IV**:5.64
 measurements of, **V**:2.2–2.13, 2.4*f*, 2.6*f*, 2.8*f*, 2.10*f*, 2.12*f*
 molecular, **V**:3.12–3.15, 3.13*f*, 3.22*f*, 3.22–3.23, 3.23*f*
 in multilayer reflectors, **IV**:7.40–7.41, 7.41*f*
 multiphonon, **IV**:5.16–5.17, 5.17*t*, 5.18*f*
 in neutron optics, **V**:63.6
 nonlinear, **IV**:16.7–16.9
 by organic detritus, **IV**:1.25–1.27, 1.25*t*, 1.26*f*
 in overdense plasmas, **IV**:21.47*f*, 21.47–21.48
 phonon, **IV**:5.13–5.16, 5.15*f*
 photo-, **V**:36.1
 in photonic crystal fibers, **V**:11.20*f*, 11.20–11.21
 by phytoplankton, **IV**:1.23–1.25, 1.24*f*–1.25*f*
 quantum resonance, **II**:22.16, 22.17
 and reabsorption of spectra, **V**:56.8
 resonance, **IV**:21.47*f*, 21.47–21.48
 by sea water, **IV**:1.21, 1.22*t*
 in solids, **IV**:8.27–8.28
 spectral, **IV**:1.26*f*, 1.26–1.27
 stimulated, **II**:16.7–16.8, 16.8*f*
 superlinear, **IV**:5.57
 two-photon, **IV**:5.56
 by water, **IV**:1.20–1.29
 and x-ray optics, **V**:26.7
- Absorption coefficient(s), **II**:32.2–32.4, 32.3*f*, **IV**:1.23*t*
 for optical constants, **IV**:5.9–5.10
 of *pin* photodiodes, **II**:25.8, 25.9*f*
- Absorption coefficient(s) (*Cont.*):
 spectral
 of natural waters, **IV**:1.27*f*
 for phytoplankton, **IV**:1.24, 1.24*f*, 1.25*f*, 1.25*t*
 of visible array detectors, **II**:32.2–32.3, 32.3*f*
 of water, **IV**:1.5*t*, 1.7*f*, 1.10, 1.17, 1.18*t*
 natural water, **IV**:1.20–1.21
 sea water, **IV**:1.21, 1.22*t*
- Absorption coefficient dependence, **IV**:5.20, 5.21*f*
- Absorption cross section, **I**:7.5, 31.3
- Absorption index, **I**:12.6
- Absorption rate, **II**:23.8
- Absorption saturation, **IV**:16.21*f*, 16.21–16.22
- Absorption spectrum, **IV**:5.12, 22.6, 22.6*f*, 22.7*f*
- Absorption transitions:
 direct interband, **IV**:5.22*f*–5.23*f*, 5.22–5.23, 5.25*f*
 indirect, **IV**:5.22–5.24, 5.24*f*–5.25*f*
- ac Kerr effect, **V**:7.11
- Ac lamps, **II**:15.32*f*
- ac Stark effect, **IV**:16.3*t*, 16.7
- Acceleration:
 bubble, **IV**:21.41–21.42, 21.42*f*
 direct laser, **IV**:21.43
 electron, **IV**:21.39–21.42, 21.40*f*
 MeV proton, **IV**:21.54, 21.54*f*
 plasma beat wave, **IV**:21.41
 target normal sheath, **IV**:21.54
- Accent lighting, **II**:40.14, 40.14*f*
- Acceptance (*étendue*), **I**:1.22, 1.81, 13.7
- Access time, for optical disk data, **I**:35.6
- Accessories, for cameras, **I**:25.16–25.17, 25.18*f*
- Accommodating intraocular lenses,
III:14.29–14.30, 21.1, 21.14, 21.18–21.19
- Accommodation, **III**:1.3, 1.29–1.36, 12.3, 21.2
 accuracy of, **III**:1.32–1.34, 1.33*f*, 1.34*f*
 in aging eyes, **III**:1.35*f*, 1.35–1.36, 14.4, 21.3
 application to instrumentation, **III**:1.34–1.35
 and change in lens, **III**:1.5
 with computer work, **III**:23.10–23.11
 defined, **III**:12.1, 13.1, 16.1, 21.1, 23.1
 dynamics of, **III**:1.31*f*, 1.31–1.32
 and extraocular muscle movement, **III**:12.16
 fluctuations in, **III**:1.32
 with head-mounted displays, **III**:13.31, 25.10–25.12
 and ocular aberration, **III**:1.17–1.18
 and presbyopia, **III**:14.8, 14.9*f*, 14.29–14.30
 and refractive errors, **III**:12.15–12.16

- Accommodation (*Cont.*):
 and spherical ametropias, **III**:16.5
 stability of, **III**:1.32
 vergence input, **III**:1.29, 1.30, 1.34
- Accommodative demand, **III**:13.30
 with contact lenses, **III**:20.26–20.30,
 20.27*f*–20.29*f*, 20.29*t*
 defined, **III**:20.1
- Accommodative miosis, **III**:1.30
- Accordion solutions, **IV**:15.28
- Accuracy:
 of absolute measurements, **II**:34.21
 of CGHs, **II**:14.6*f*, 14.6–14.7, 14.7*f*
 as measure of systemic errors, **II**:12.2
- Acellular, **III**:16.1
- Achromatic antireflection coatings, **IV**:7.29
- Achromatic beam splitters, **IV**:7.62*f*–7.65*f*,
 7.62–7.65
- Achromatic detection:
 and chromatic adaptation, **III**:11.47–11.49,
 11.48*f*
 measurements favoring, **III**:11.37
- Achromatic doublets (lenses), **I**:17.22–17.25,
 17.23*f*–17.25*f*, 17.24*t*
- Achromatic lenses, athermalized, **II**:1.16, 1.16*f*
- Achromatic mechanism (color vision),
III:11.1, 11.37
- Achromatic retardation plates, **I**:13.48–13.52,
 13.50*f*, 13.53*t*
- Achromatic signals, multiplexing of chromatic
 signals and (color vision), **III**:11.76–11.79,
 11.78*f*
- Achromaticity, of fiber-based couplers,
V:16.3–16.4
- Achromatism, **II**:1.14–1.15, 1.15*f*
- Acktar black coatings, **IV**:6.55
- Acoustic optical modulators (AOMs), **IV**:14.19,
 20.13
- Acoustic phonons, **IV**:5.14–5.16
- Acoustically rotated tangential phase matching,
V:6.26*f*, 6.27
- Acousto-optic cells, **I**:11.8–11.9, 11.9*f*
- Acousto-optic correlators, **I**:11.10–11.12, 11.11*f*
- Acousto-optic devices, **V**:6.3–6.45
 and acousto-optic diffraction, **V**:6.4, 6.9
 Bragg cells (wideband), **V**:6.30–6.31, 6.31*t*
 deflectors, **V**:6.22–6.29, 6.24*f*, 6.26*f*, 6.28*f*,
 6.29*t*
 figures of merit for, **V**:6.16–6.17
 materials for, **V**:6.16–6.22
- Acousto-optic devices (*Cont.*):
 modulators, **V**:6.31–6.35, 6.34*t*
 acousto-optic frequency shifters, **V**:6.35
 and Bragg diffraction, **V**:6.4, 6.6, 6.7, 6.14
 image (scophony), **V**:6.33, 6.34
 principle of operation, **V**:6.32, 6.32*f*
 optical birefringence in, **V**:6.17
 propagation and attenuation in, **V**:6.17
 theory of acousto-optic interaction,
V:6.5–6.16
 elasto-optic and roto-optic effects,
V:6.5–6.6
 of finite geometry, **V**:6.14–6.16
 frequency characteristics, **V**:6.12–6.14,
 6.13*f*, 6.14*f*
 phase matching, **V**:6.9–6.12, 6.10*f*
 plane wave analysis, **V**:6.6–6.9
 tunable filters as [*see* Acousto-optic tunable
 filters (AOTFs)]
 wideband AO Bragg cells, **V**:6.30, 6.31*t*
- Acousto-optic diffraction, **V**:6.4, 6.9
- Acousto-optic figures of merit, **V**:6.16–6.17
- Acousto-optic frequency shifters (AOFS),
V:6.35
- Acousto-optic interaction, theory of,
V:6.5–6.16
 elasto-optic and roto-optic effects, **V**:6.5–6.6
 of finite geometry, **V**:6.14–6.16
 frequency characteristics, **V**:6.12–6.14,
 6.13*f*, 6.14*f*
 phase matching, **V**:6.9–6.12, 6.10*f*
 plane wave analysis, **V**:6.6–6.9
- Acousto-optic modulators, **V**:6.23*t*, 6.31–6.35,
 6.34*t*
 acousto-optic frequency shifters, **V**:6.35
 image (scophony), **V**:6.34
 principle of operation, **V**:6.32, 6.32*f*
- Acousto-optic scanners, **I**:30.44–30.45
- Acousto-optic tunable filters (AOTFs), **V**:6.23*t*
 collinear beam, **V**:6.43, 6.43*f*, 6.45*t*
 long-infrared, 6.42
 and longitudinal spatial modulation, **V**:6.12
 mid-infrared, **V**:6.42
 noncritical phase-matching, **V**:6.37,
 6.39–6.42, 6.38*f*, 6.43*t*
 principle of operation, **V**:6.36–6.39, 6.37*f*,
 6.38*f*
 ultraviolet, **V**:6.40
- Actinic effects (of radiation), **II**:34.6, 34.7
- Actinic ultraviolet action spectrum, **II**:36.17

- Actinometry:
 conversions between radiometry/photometry
 and, **II**:34.12, 34.12*t*
 defined, **II**:34.7, 34.11
- Action spectra (ocular radiation),
III:7.2, 7.3, 7.4*f*
- Activated phosphor detectors, **V**:60.7–60.8
- Activated-phosphor sources (of radiation),
II:15.49
- Active athermalization, **II**:6.24, 6.24*f*
- Active autofocus systems, for cameras,
I:25.11–25.12, 25.12*f*
- Active devices:
 fabrication of, **I**:21.14
 for integrated optics, **I**:21.25–21.31,
 21.26*f*–21.31*f*
- Active imaging, **II**:31.29–31.30
- Active layer removal, in PIC manufacturing,
I:21.19
- Active mechanical athermalization, **II**:8.11, 8.11*f*
- Active mode locking, **V**:20.17
- Active nonlinear optical phenomena, **IV**:5.54,
 5.54*t*
- Active optical limiting, **IV**:13.1–13.3
- Active optics, **V**:50.1, 50.2
- Active pixel sensors, **I**:26.2, 26.8*f*, 26.8–26.9
- Active scanning, **I**:30.4
- Active-passive transitions, in PICs,
I:21.19–21.20
- Acutance:
 defined, **II**:30.2
 of photographic systems, **II**:29.17–29.19,
 29.18*t*, 29.19*f*
- Adachi dispersion model, **IV**:2.15, 2.22
- Adaptation:
 in aging eyes, **III**:14.15
 in contrast detection, **III**:2.25–2.27
 multiplicative and subtractive, **III**:2.26
 phoria, **III**:13.21–13.22
 in vision, **II**:40.9
 in wearing contact lenses, **III**:12.12
- Adaptive microstructured optical arrays
 (MOAs), **V**:50.7–50.8, 50.8*f*
- Adaptive optics (AO) (in general), **V**:5.1–5.46,
 50.1
 concepts and components of, **V**:5.2–5.5, 5.3*f*
 designing systems of, **V**:5.38–5.46
 requirements, **V**:5.39, 5.40*t*
 results for 3.5-m telescope systems,
V:5.43*f*–5.45*f*, 5.43–5.45
- Adaptive optics (AO) (in general), designing
 systems of (*Cont.*):
 results for 10-m telescope systems,
V:5.45*f*–5.46*f*, 5.45–5.46
 subaperture size, **V**:5.40–5.43, 5.42*f*
 as enabling technology, **V**:5.2
 hardware and software implementation for,
V:5.21–5.38
 higher-order wavefront sensing techniques,
V:5.36–5.37
 laser beacons, **V**:5.27–5.34, 5.28*f*–5.31*f*,
 5.33*f*
 real-time processors, **V**:5.34*f*, 5.34–5.35, 5.35*f*
 Shack-Hartmann technique, **V**:5.23*f*,
 5.23–5.27, 5.25*f*, 5.26*f*
 tracking, **V**:5.21–5.23, 5.22*f*
 wavefront correctors, **V**:5.37–5.38, 5.38*f*
 and imaging through atmospheric
 turbulence, **V**:4.35–4.36
 and turbulence, **V**:5.5–5.21
 anisoplanatism, **V**:5.19
 atmospheric tilt and Strehl ratio, **V**:5.14,
 5.14*f*–5.15*f*
 Fried's coherence diameter and spatial
 scale, **V**:5.9–5.13, 5.11*t*, 5.12*f*, 5.13*f*
 higher-order phase fluctuations,
V:5.18–5.19
 on imaging and spectroscopy, **V**:5.19*f*,
 5.19–5.21, 5.20*f*
 Kolmogorov model, **V**:5.5–5.6
 tracking requirements, **V**:5.15–5.18, 5.17*f*
 variation of n and C_n^2 parameters,
V:5.6–5.8, 5.7*f*, 5.8*f*
- Adaptive optics (AO) (in retinal microscopy
 and vision), **III**:1.25, 15.1–15.24
- AO-controlled light delivery to the retina,
III:15.22–15.24
 alignment, **III**:15.23
 in an AO SLO, **III**:15.23
 conventional AO vision systems, **III**:15.23
 to generate aberrations, **III**:15.24
 longitudinal chromatic aberration,
III:15.22
 measuring activity of individual cones,
III:15.24
 transverse chromatic aberration,
III:15.22–15.23
 uses of, **III**:15.23–15.24
 for correcting monochromatic aberrations,
III:1.25

- Adaptive optics (AO) (in retinal microscopy and vision) (*Cont.*):
 and correction for SCE-1, **III**:9.5
 defined, **III**:15.1, 17.1
 history of, **III**:15.2–15.3
 imaging of the retina, **III**:15.16–15.22
 contrast and resolution, **III**:15.21–15.22
 flood-illuminated AO ophthalmoscope, **III**:15.16*f*, 15.16–15.17, 15.17*f*
 optical coherence tomography, **III**:15.19–15.21, 15.20*f*, 15.21*f*
 scanning laser ophthalmoscope, **III**:15.17–15.19, 15.18*f*, 15.19*f*, 17.9
 implementation of, **III**:15.7*f*, 15.7–15.15
 control system, **III**:15.12–15.15
 wavefront corrector, **III**:15.9–15.12, 15.10*f*, 15.11*f*
 wavefront sensor, **III**:15.8*f*, 15.8–15.9
 properties of ocular aberrations, **III**:15.4–15.7, 15.5*f*, 15.6*f*
 wavefront corrector, **III**:15.9–15.12, 15.10*f*, 15.11*f*
 wavefront sensor, **III**:15.8*f*, 15.8–15.9
- Adaptive optics retina cameras, **III**:15.3, 15.12 (*see also* Ophthalmoscopes)
- Adaptive x-ray optics, **V**:50.1–50.8
 hard vs. soft x-ray telescopes, **V**:50.2
 history of, **V**:50.1, 50.2*f*
 for synchrotron and lab-based sources, **V**:50.2–50.7, 50.4*f*, 50.6*f*, 50.8*f*
- Add/drop multiplexers, optical, **V**:21.2, 21.8, 21.8*f*, 21.9*f*, 21.12
- Addition, as analog operation, **I**:11.2
- Additive damping, **II**:3.18
- Additive pulse modelocking (APM), **IV**:18.3, 18.14
- Additives, polymers as, **IV**:3.5
- Additivity:
 in color matching, **III**:10.8
 of color opponency and color appearance, **III**:11.67
 with equilibrium colors, **III**:11.66
 field, **III**:11.51, 11.52*f*
 obedience to Abney's law, **III**:10.44, 11.37
 Stiles-Crawford effect, **III**:9.10*f*–9.11*f*
- Adiabatic approximation, **II**:23.21 (*see also* Markov approximation)
- Adiabatic ionization stabilization, **IV**:21.20, 21.21, 21.22*f*
- Adiabatic potentials, **IV**:21.23, 21.24*f*
- Adjustment tasks:
 judgment tasks vs., **III**:3.2
 psychophysical measurement of, **III**:3.4–3.5
 magnitude production, **III**:3.5
 matching, **III**:3.4–3.5
 nulling, **III**:3.4
 threshold, **III**:3.4, 3.5*f*
- ADONIS system, **V**:5.35
- Advanced Photo System (APS), **II**:30.21, 30.26, 30.27*t*
- Advanced X-Ray Astrophysical Facility (AXAF), **V**:44.4, 44.6*f*, 44.10
- Aerial cameras, **I**:25.20
- Aerial images, **I**:1.26
- Aerogels, **IV**:6.15
- Aeroglaze Z series, **IV**:6.36*f*, 6.37, 6.37*f*, 6.39, 6.39*f*–6.42*f*
- Aeroglaze Z302, **IV**:6.14
- Aeroglaze Z306, **IV**:6.14, 6.17, 6.29*f*, 6.35
- Aerosol Polarimeter Sensor (APS), **I**:15.38
- Aerosols, **V**:3.6, 3.7, 3.10*f*, 3.10–3.11, 3.11*f*
- Aesthetics and emotional characteristics (electronic imaging), **III**:24.9–24.10
 color, art, and emotion, **III**:24.10
 image quality, **III**:24.9
 virtual reality and presence, **III**:24.9
- Affine transformations, **I**:1.57
- Afocal Cassegrain-Mersenne telescope objective, **I**:29.9
- Afocal Gregorian-Mersenne telescope objective, **I**:29.12
- Afocal lenses, **I**:18.1–18.22
 for binoculars, **I**:18.13–18.14, 18.14*f*
 catadioptric, **I**:18.21–18.22, 18.22*f*
 Galilean, **I**:18.15*f*, 18.15–18.17, 18.16*f*
 Gaussian, **I**:1.45, 1.46*f*, 1.53*f*, 1.53–1.54, 1.54*f*
 Gaussian analysis, **I**:18.4–18.6
 and focusing lenses, **I**:18.2–18.4, 18.3*f*, 18.5*f*
 and optical invariant, **I**:18.7
 subjective aspects of, **I**:18.6, 18.7*f*
 Keplerian, **I**:18.7–18.14
 and eye relief manipulation, **I**:18.8–18.10, 18.9*f*, 18.10*f*
 field-of-view limitations in, **I**:18.11
 finite conjugate afocal relays, **I**:18.11–18.12, 18.12*f*
 thin-lens model of, **I**:18.7–18.8, 18.8*f*
 paraxial matrix methods, **I**:1.70
 for periscopes, **I**:18.19, 18.19*f*

- Afocal lenses (*Cont.*):
 reflecting, **I**:18.19–18.21, 18.20*f*, 18.21*f*
 in relay trains, **I**:18.17*f*, 18.17–18.19, 18.18*f*
 for scanners, **I**:18.13, 18.13*f*
 for telescopes, **I**:18.10–18.11, 18.11*f*
- Afocal magnification, **I**:18.5–18.6
- Afocal objectives, **I**:29.9, 29.12, 29.29–29.30
- Afocal systems:
 as attachments, **II**:1.8, 1.9*f*
 first-order layout for, **II**:1.7, 1.7*f*
- After-effect (color vision):
 and McCollough effect, **III**:11.76
 and orientation of contours, **III**:11.75*f*,
 11.75–11.76
 and orientation-selectivity, **III**:11.80
- Afterimage, **III**:23.1
- Against-the-rule astigmatism, **III**:1.6
- Age-related changes in vision, **III**:1.5, 1.7,
 1.35–1.36, 14.1–14.30
 in accommodation, **III**:1.35*f*, 1.35–1.36,
 14.4, 21.3
 in bovine lenses, **III**:19.10
 changes in the eye causing, **III**:21.3–21.4
 in color vision, **III**:14.15, 14.17
 and demographic changes in world, **III**:14.2*f*,
 14.2–14.4
 in depth and stereovision, **III**:14.22
 economic/social implications of, **III**:14.4
 and fluorescence, **III**:1.21
 and health implications of living longer,
III:14.3–14.4
 index of diffusion for, **III**:1.23
 and level of aberration, **III**:1.18
 minimal, **III**:14.22
 ocular diseases, **III**:14.22–14.27
 age-related macular degeneration,
III:14.24–14.25
 cataract, **III**:14.24
 diabetic retinopathy, **III**:14.25–14.26
 glaucomas, **III**:14.26–14.27
 life-span environmental radiation damage,
III:14.22–14.23
 in optics, **III**:14.4–14.14
 accommodation and presbyopia, **III**:14.8,
 14.9*f*
 anterior and posterior chambers, **III**:14.5,
 14.6
 cornea, **III**:14.5, 14.6*f*
 eye size, **III**:14.11
 lens, **III**:14.7–14.8
- Age-related changes in vision, in optics (*Cont.*):
 pupil, **III**:14.6–14.7, 14.7*f*
 retina, **III**:14.9–14.11
 tears, **III**:14.4–14.5
 transparency/cataract, **III**:14.8
 presbyopic corrections for, **III**:14.27–14.30
 accommodation restoration,
III:14.29–14.30
 contact lenses, **III**:14.27–14.28
 intraocular lenses, **III**:14.28–14.29
 noncataract related refractive surgeries,
III:14.29
 spectacles, **III**:14.27
 in pupil diameter, **III**:1.8, 1.18
 in retinal image quality, **III**:14.11–14.14
 chromatic aberration, **III**:14.14
 intraocular scatter, **III**:14.12
 monochromatic aberrations,
III:14.12–14.14, 14.13*f*
 in retinal reflectance, **III**:1.11
 and RMS wavefront error, **III**:1.15–1.17,
 1.16*t*
 in scattering, **III**:1.20
 in sensitivity, **III**:14.15, 14.16*f*
 in spatial vision, **III**:14.17–14.19, 14.18*f*
 in temporal vision, **III**:14.19–14.21, 14.20*f*
 in transmittance, **III**:1.9
 in visual acuity, **III**:4.13
 in visual field, **III**:14.21, 14.21*f*
- Age-related macular degeneration (AMD),
III:14.24–14.25
 defined, **III**:14.1
 OFDI at 1050 nm in, **III**:18.17, 18.17*f*
 OFDI in diagnosis of, **III**:18.8
- Agfachrome, **II**:29.14
- Agile beam steering, **I**:30.51–30.63
 with decentered-lens and microlens arrays,
I:30.57–30.60, 30.58*f*–30.60*f*, 30.62–30.63
 with digital micromirror devices, **I**:30.60–30.61
 with gimbal-less two-axis scanning
 micromirrors, **I**:30.61–30.62, 30.62*f*
 phased-array, **I**:30.52–30.57, 30.53*f*,
 30.62–30.63
- Ahrens method of spar cutting, **I**:13.12
- Ahrens Nicol prisms, **I**:13.16*f*, 13.17
- Ahrens prisms, **I**:13.14
- AHU pelloid, **II**:30.4
- Air cladding, for photonic crystal fibers,
V:25.2, 25.21

- Aircraft, synthetic aperture radar for, **I**:11.6–11.7, 11.7*f*
- Air-filled photonic crystal fibers, **V**:25.20–25.21
- Air-spaced doublet lens, **II**:6.7
- Air-spaced triplet lens, **II**:6.21, 6.22*f*
- Airway beacon lamps, **II**:15.11
- Airy diffraction pattern, **III**:1.12
- Airy disks, **III**:4.1, 4.4*f*, 4.5; **V**:37.6
in confocal microscopy, **I**:28.50
defined, **I**:3.26
of DIC microscopes, **I**:28.39
of microscopes, **I**:28.17–28.19, 28.18*f*, 28.19*f*
of solid-state cameras, **I**:26.15
and vector diffraction, **I**:3.33
- Airy distribution, of thin lenses, **V**:40.3
- Airy equation, **I**:12.10
- Airy pattern, **I**:17.38
- Airy-Drude formula, **I**:16.5
- Akhieser loss, **V**:6.17
- Akzo (company), **IV**:6.14
- Akzo Nobel paints, **IV**:6.39, 6.42*f*, 6.43*f*
- Algorithms, singular-value-decomposition, **V**:5.25, 5.26
- Alignment:
of eye (*see* Eye alignment)
of photoreceptors, **III**:8.5–8.8, 8.6*f*, 8.7*f*
- A-line, **III**:18.1
- All solid-core photonic crystal fibers, **V**:25.21
- Allan Deviation, **II**:22.2–22.4
- Allan Variance method, **II**:22.2, 22.3
- All-dielectric color selective (dichroic) beam splitters, **IV**:7.65–7.66, 7.67*f*
- All-dielectric reflectors:
broadband, **IV**:7.39, 7.40*f*, 7.45*f*–7.47*f*, 7.45–7.47
interference filters with, **IV**:7.81
- Allelotropia, **III**:13.7
- All-fiber monolithic systems (fiber lasers), **V**:25.9–25.11, 25.16*f*, 25.16–25.18
- All-optical switching, in OTDM communication networks, **V**:20.22, 20.23*f*
- Alloy disordering, **II**:19.24
- All-solid band gap guiding fibers, **V**:11.15, 11.16, 11.17*f*
- Alphanumeric displays, LED, **II**:17.31–17.32
- Altaite (PbTe), **IV**:2.40*t*, 2.44*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.64*t*, 2.69*t*
- Altenhof objectives, **I**:29.32–29.33
- Alternative Paul three-mirror objective, **I**:29.28
- Alumina (Al₂O₃), **IV**:2.38*t*, 2.46*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.70*t*, 2.76*t*
- Aluminized phosphor-screen/window assembly, **II**:31.14, 31.15*f*
- Aluminum:
absorptance of, **IV**:4.40*f*, 4.48*t*, 4.51*t*
anodized, **IV**:6.5*f*, 6.38*f*, 6.58*f*
black velvet cloth on, **IV**:6.31*f*
commando cloth on, **IV**:6.31*f*
diamond turning and, **II**:10.4
flame-sprayed, **IV**:6.57
grooved and blazed, **IV**:6.21, 6.24*f*, 6.25*f*
optical constants for, **IV**:4.11
optical properties of, **IV**:4.12*t*, 4.20*f*, 4.21*f*
penetration depth of, **IV**:4.47*f*
physical properties of, **IV**:4.52*t*, 4.54*t*
reflectance of, **IV**:4.27*t*–4.28*t*, 4.40*f*, 4.44*f*, 4.46*f*
roughened, **IV**:6.21, 6.22*f*, 6.30*f*
sandblasted, **IV**:6.45*f*
thermal properties of
coefficient of linear thermal expansion, **IV**:4.56*t*, 4.57*f*
elastic properties, **IV**:4.69*t*
at room temperature, **IV**:4.55*t*
specific heat, **IV**:4.65*t*, 4.66*f*
strength and fracture properties, **IV**:4.70*t*
thermal conductivity, **IV**:4.58*t*, 4.59*f*–4.60*f*
and UV light, **IV**:6.21
- Aluminum alloys:
reflectance of, **IV**:4.44*f*
thermal conductivity of, **IV**:4.59*f*–4.60*f*
- Aluminum gallium arsenide (AlGaAs) emitters, **II**:17.32
- Aluminum gallium arsenide (AlGaAs) LEDs, **II**:17.17, 17.17*t*, 17.28*f*
- Aluminum gallium arsenide (AlGaAs) quantum well photodetectors, **II**:25.16*f*, 25.16–25.17, 25.17*f*
- Aluminum gallium arsenide (AlGaAs) substrate, **II**:17.22
- Aluminum gallium nitride (AlGaN) alloy photovoltaic detectors, **II**:24.46
- Aluminum gallium nitride (AlGaN) substrate, **II**:17.22
- Aluminum indium gallium nitride (AlInGaN) material systems, **II**:18.1–18.2, 18.2*f*, 18.4

- Aluminum indium gallium phosphide (AlInGaP) LEDs, **II**:17.18, 17.19f
- Aluminum indium gallium phosphide (AlInGaP) material systems, **II**:18.1, 18.5
- Aluminum indium gallium phosphide (AlInGaP) substrate, **II**:17.22
- Aluminum mirrors, **II**:6.20f; **IV**:7.106f–7.108f, 7.106–7.108
- Aluminum oxide (Al_2O_3), **V**:2.19, 2.19f, 2.21, 2.22
- Aluminum oxynitride ($\text{Al}_{23}\text{O}_{27}\text{N}_5$) (ALON), **IV**:2.38t, 2.44t, 2.47t, 2.50t, 2.55t, 2.60t, 2.76t
- Alvarez plates, **I**:22.16
- Alvarez-Humphrey plates, **I**:22.37
- Amacrine cells, **III**:2.10, 2.11
- Ambient lighting, **II**:40.12, 40.13f, 40.15f
- Ambient temperature electrical substitution radiometers, **II**:34.27
- Amblyopia, **III**:2.34, 2.35, 12.1, 12.16
- American Conference of Governmental Industrial Hygienists (ACGIH), **III**:7.9
- American Institute of Physics (AIP), **II**:36.3
- American National Standards Institute (ANSI), **II**:4.11, 36.2; **III**:7.9, 7.11
- American Society for Testing and Materials (ASTM), **II**:37.11; **IV**:6.17
- Ames Perfect Diffuse Reflector, **IV**:6.7f
- Ames Research Center, **IV**:6.34
- Ames 24E, **IV**:6.7f, 6.26, 6.26f, 6.27f, 6.34
- Ames 24E2, **IV**:6.27f, 6.28f, 6.34
- Ametropias, **III**:1.6–1.7, 12.4, 16.4–16.5, 16.5f
- astigmatism, **III**:16.5–16.6
- correcting, **III**:12.16, 16.8f
- defined, **III**:12.1, 13.1
- spherical, **III**:16.5
- uncorrected, **III**:12.16
- AMI (amplified MOS imager) MOS readout, **II**:32.21
- Amici lenses, **I**:17.10, 17.10f
- Amici prisms, **I**:19.3t, 19.11, 19.12f, 20.6f
- Ammonium phosphate ($\text{NH}_4\text{H}_2\text{PO}_4$, ADP), **IV**:2.40t, 2.45t, 2.48t, 2.52t, 2.57t, 2.64t
- Ammosoc-Delone-Krainov (ADK) ionization rate, **IV**:21.13
- Amorphous materials, **IV**:2.5, 12.26t [*see also* Glass(es)]
- Amorphous silicon photoconductors, **II**:32.4f, 32.31, 32.32
- Ampere law, **IV**:2.6
- Amplification, **II**:16.9
- chirped pulse, **V**:25.2, 25.32, 25.33
- in color CRTs, **III**:22.6–22.9, 22.7f
- lump, **V**:21.44
- by SOAs, **V**:19.1f–19.2f, 19.2, 19.22–19.27
- CWDM systems, **V**:19.27
- DWDM systems, **V**:19.25f, 19.25–19.27, 19.26f
- single-channel systems, **V**:19.22f–19.24f, 19.22–19.24
- Amplified spontaneous emission (ASE), **V**:9.13, 10.11, 14.3, 14.6
- Amplified spontaneous emission (ASE) noise, **V**:19.3, 19.4f, 19.9, 19.18, 19.24, 19.35
- Amplified spontaneous emission (ASE) power, **V**:19.20
- Amplifiers, **II**:27.2
- in communication systems, **V**:9.13–9.14
- erbium-doped fiber, **V**:14.4–14.7, 14.5f, 21.38f–21.42f, 21.38–21.41
- erbium/ytterbium-doped fiber, **V**:14.7–14.8
- field-effect transistor, **V**:13.70
- flashlamp pumped Nd:glass, **IV**:21.5
- linear optical, **V**:19.14, 19.27
- lock-in, **IV**:5.64
- optical fiber, **V**:14.1–14.11, 14.2t
- optical parametric, **IV**:23.13–23.14
- parametric, **V**:11.23, 11.24, 14.10–14.11
- praseodymium-doped fiber, **V**:14.7
- properties of, **II**:16.3
- for PZTs, **II**:22.18
- Raman, **IV**:15.4, 15.4f, 22.15; **V**:14.8f, 14.8–14.9, 14.10f, 21.42f–21.44f, 21.42–21.44
- rare-earth-doped, **V**:14.2–14.4, 14.3f
- selection of, **II**:27.10–27.12
- semiconductor optical [*see* Semiconductor optical amplifiers (SOAs)]
- semiconductor vs. fiber, **V**:9.13–9.14
- slab-coupled optical waveguide, **V**:19.21
- Ti:sapphire, **IV**:21.5
- transconductance, **II**:27.11f, 27.11–27.12
- voltage, **II**:27.10–27.11
- in WDM networks, **V**:21.37f, 21.37–21.44
- ytterbium-doped fiber, **V**:14.7
- Amplifying media, antireflection coatings for, **IV**:7.26, 7.27
- Amplitude, of waves, **I**:2.4, 2.5, 12.5

- Amplitude division, interference by, **I**:2.14, 2.19–2.28
 and extended sources, **I**:2.20
 and Fizeau interferometers, **I**:2.24–2.26, 2.25*f*
 and fringes of equal inclination, **I**:2.20–2.22, 2.21*f*, 2.22*f*
 and fringes of equal thickness, **I**:2.22–2.24, 2.23*f*
 and Michelson interferometers, **I**:2.26*f*–2.27*f*, 2.26–2.28
 plane-parallel plate, **I**:2.19, 2.20*f*, 2.30*f*, 2.30–2.33, 2.32*f*, 2.33*f*
 thin films, **I**:2.24
- Amplitude gating, **II**:21.7
- Amplitude modulation (AM), **II**:19.36; **V**:7.22–7.24, 7.23*f*, 7.24*f*
- Amplitude modulation index, **V**:7.22
- Amplitude penetration depth, **I**:12.5
- Amplitude reflection coefficients, Fresnel, **I**:12.7–12.8, 12.10
- Amplitude response, frequency vs., **II**:22.6–22.7
- Amplitude scattering matrix, **I**:7.10, 7.13 (*see also* Jones matrix)
- Amplitude transmission coefficients, Fresnel, **I**:12.8
- Amplitude zone plates, **V**:40.4–40.5, 40.5*f*
- Amplitude-modulated (AM) acoustic waves, **V**:6.32
- Amplitude-modulated (AM) systems, optical fibers in, **V**:9.16
- Amplitude-shift-keyed (ASK) transmission, **I**:21.30
- Amplitude-shift-keying (ASK), **V**:21.28, 21.29
- AMTIR-1 glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.68*t*
- AMTIR-3 glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.68*t*
- Analog modulation, **V**:6.32
- Analog optical signal and image processing, **I**:11.1–11.20
 Fourier transforms in, **I**:11.3*f*, 11.3–11.5, 11.5*f*
 and fundamental analog operations, **I**:11.2–11.3
 incoherent processing, **I**:11.17–11.20, 11.18*f*, 11.19*f*
 and spatial filtering, **I**:11.5–11.6, 11.6*f*
 of synthetic aperture radar data, **I**:11.6–11.8, 11.7*f*–11.8*f*
 of temporal signals, **I**:11.8–11.12, 11.9*f*–11.11*f*
 of two-dimensional images, **I**:11.12–11.17, 11.13*f*
- Analog to digital conversion, **V**:20.4, 20.8, 20.8*f*
- Analog transmission, **I**:21.32–21.34, 21.33*f*, 21.34*f*; **V**:9.15–9.17, 9.16*f*
- Analytical density, **II**:29.14
- Analytical signal representation, in coherence theory, **I**:5.2–5.3
- Analyzed states, of polarizers, **I**:15.19
- Analyzer vectors, **I**:15.11
- Anamorphic afocal attachments (anamorphosers), **I**:18.16*f*, 18.16–18.17
- Anamorphic error control, **I**:30.49, 30.50
- Anastigmatic (term), **I**:29.36
- Anastigmatic objectives, **I**:29.12–29.13
 -ance (suffix), **II**:35.3; **IV**:4.5
- Anger cameras, **V**:63.33
- Angle α , **III**:1.20
- Angle characteristic function, **I**:1.14–1.15, 1.15*f*, 1.17
- Angle measurement, **II**:12.10–12.17
 autocollimeters for, **II**:12.11*f*, 12.11–12.12, 12.12*f*
 interferometric methods of, **II**:12.14
 levels (tools) for, **II**:12.13*f*, 12.13–12.14, 12.14*f*
 mechanical methods of, **II**:12.10–12.11, 12.11*f*
 in prisms, **II**:12.14–12.16, 12.15*f*–12.17*f*
 theodolites for, **II**:12.13
- Angle of arrival, **V**:4.23–4.26, 4.25*f*, 4.26*f*
- Angle of deflection, **V**:6.40
- Angle of incidence, **I**:1.23, 1.39, 13.48
- Angle of incidence, obliquity of, **III**:8.20
- Angle scans, for scatterometers, **V**:1.14
- Angle solves, **II**:3.6
- Angled stripe designs, of SOAs, **V**:19.8, 19.8*f*
- Angle-of-incidence effects, for cutoff filters, **IV**:7.56
- Angle-point characteristic function, **I**:1.16, 1.17
- Angular apertures, of NPM AOTFs, **V**:6.41
- Angular change, of optical beams, **I**:30.5
- Angular correction function, **I**:5.6, 5.7*f*
- Angular dilution, **II**:39.6
- Angular dispersion, **V**:38.6
- Angular distribution, **IV**:1.12, 21.8–21.9, 21.9*f*
- Angular magnification, **I**:1.52, 1.78, 18.4
- Angular resolution:
 in astronomical x-ray optics, **V**:47.10–47.11, 47.11*f*
 in Wolter x-ray optics, **V**:47.2*f*, 47.2–47.3

- Angular resolution, in diffraction-limited eye, **III**:1.12
- Angular scan, **I**:30.28–30.29
- Angular sensitivity, in two-material periodic multilayers theory, **IV**:7.37
- Angular spectrum of plane waves (ASW) approach, **V**:6.14–6.15
- Angular spectrum representation, **I**:5.14*f*, 5.14–5.15
- Angular uniformity, **II**:39.31
- Anharmonic oscillator model, of second-order nonlinear optical susceptibility, **IV**:10.7–10.9, 10.8*f*
- Aniridia, **III**:9.1
- Aniseikonia, **III**:1.41–1.42
 - binocular factors in, **III**:12.16–12.17
 - defined, **III**:9.1, 12.1, 13.1, 25.1
 - distortion from interocular anisomagnification, **III**:13.17–13.18
- Anisometropia, **III**:1.7, 1.41, 1.42
 - binocular factors in, **III**:12.16–12.17
 - defined, **III**:12.1, 23.1
 - interocular blur suppression with, **III**:13.18–13.19
 - as problem with computer work, **III**:23.11
- Anisophoria, optically induced, **III**:13.26
- Anisoplanatism, **V**:5.19
 - and laser beacons, **V**:5.27–5.29, 5.28*f*–5.30*f*
 - and subaperture size, **V**:5.42–5.43
- Anisotropic acoustic beam collimation, **V**:6.25
- Anisotropic crystals, propagation of light and, **IV**:8.8–8.11, 8.9*t*, 8.10*f*
- Anisotropic diffraction, **V**:6.10
- Anisotropic scattering, **IV**:12.7
- Anisotropy, **I**:35.26; **V**:17.2
- Annealed proton exchange (APE), **V**:7.30
- Annealed proton exchange (APE) process:
 - for fiber optic gyroscopes, **I**:21.35, 21.36, 21.36*t*
 - for LiNbO₃ waveguides, **I**:21.16–21.17
- Annealing:
 - of glass, **IV**:2.28
 - of optical surfaces, **IV**:19.3
- Annular flanges, **II**:6.3–6.4, 6.4*f*
- Annular polynomials:
 - for defocus, **II**:11.38*f*, 11.39
 - Zernike, **II**:11.13–11.21, 11.14*f*, 11.16*f*, 11.17*t*–11.21*t*
- Annular pupils, **V**:4.10–4.16, 4.11*f*–4.15*f*, 4.15*t*
- Annunciator assemblies, **II**:17.30
- Anodes:
 - in photomultipliers, **II**:27.6, 27.7, 27.7*f*
 - as x-ray tube sources, **V**:54.8–54.9, 54.12, 54.13
- Anodized aluminum surface, **IV**:6.5*f*, 6.38*f*, 6.58*f*
- Anodized surface treatments, **IV**:6.44–6.49
 - (see also *specific anodized treatments*, e.g.: Martin Black)
- Anodized surfaces, **IV**:6.3*t*
- Anomalous diffraction, **I**:7.5, 7.6*f*
- Anomalous (negative) dispersion, **IV**:4.4, 18.11
- Anomalous reflection colors, **II**:30.17
- Anomalous skin effect, **IV**:21.49
- ANSI (American National Standards Institute), **V**:23.4, 23.6
- Anterior chamber, **III**:14.5, 14.6, 16.3
- Antiblooming, **II**:32.9, 32.10*f*
- Antiguides, for VCSELs, **V**:13.46
- Antihalation undercoat (AHU) layers, **II**:30.4
- Antimony flint glass, **IV**:2.43*t*
- Antinodal points, of lens systems, **I**:17.7
- Antiprincipal points, of lens systems, **I**:17.7
- Antireflection (AR) coatings, **IV**:7.15–7.32; **V**:19.8, 19.8*f*, 19.20
 - of absorbing and amplifying media, **IV**:7.26, 7.27
 - homogeneous-layer, **IV**:7.16–7.23, 7.17*f*–7.19*f*, 7.20*t*–7.21*t*, 7.22*f*–7.23*f*
 - inhomogeneous and structured, **IV**:7.23–7.26, 7.24*f*, 7.26*f*
 - at nonnormal angle of incidence, **IV**:7.28*f*–7.31*f*, 7.28–7.31
 - nonoptical properties of, **IV**:7.31–7.32, 7.32*f*
 - surface reflections and optical performance, **IV**:7.15–7.16, 7.16*f*
 - of surfaces carrying thin films, **IV**:7.27–7.28, 7.28*f*
 - universal, **IV**:7.26, 7.27*f*
- Antiresonance, **V**:11.12
- Antiresonant Fabry-Perot saturable absorber (A-FPSA), **IV**:18.3, 18.11
- Anti-Stokes four-wave mixing, coherent, **IV**:15.2*t*, 15.3*t*, 15.4, 15.4*f*
- Anti-Stokes scattering, **IV**:15.1–15.3, 15.3*t*; **V**:10.5
 - coherent Raman, **IV**:15.4, 15.4*f*, 15.34, 15.42, 15.42*t*, 15.43*f*
 - multiple, **IV**:15.2*f*
 - Raman, **IV**:15.32–15.34, 15.33*f*, 15.35*f*
 - shifted Raman, **IV**:16.15, 16.15*f*
 - stimulated Raman, **IV**:15.2*t*

- Anti-Stokes shift, **IV**:15.2, 15.43
 Anti-Stokes waves, **IV**:15.1, 15.43; **V**:10.9
 AOM transducers, **II**:22.20
 APART (stray light analysis program), **II**:7.11; **IV**:6.19
 Aperture(s), **I**:1.74–1.76, 1.75*f*
 angular, **V**:6.41
 circular
 diffraction of light from, **I**:3.6*f*, 3.6–3.7, 3.7*f*, 3.9–3.11
 Fraunhofer patterns for, **I**:3.25, 3.26, 3.27*f*
 data about, **II**:3.4
 double-slit, **I**:3.26–3.28, 3.27*f*, 3.28*f*
 Gaussian, **V**:37.6
 image space numerical, **I**:1.79
 linear, **I**:30.54
 nonideal, **II**:34.35*f*, 34.35–34.36
 numerical, **I**:1.79, 17.9; **II**:34.20, 39.1; **V**:9.4, 25.2, 25.18, 42.2
 in optical design software, **II**:3.6
 order-selecting, **V**:40.4, 40.5
 pinhole, **V**:32.3
 rectangular, **I**:3.19–3.20
 Fraunhofer patterns for, **I**:3.25, 3.26
 Fresnel diffraction from, **I**:3.19–3.20, 3.20*f*
 and subaperture size, **V**:5.40–5.43, 5.42*f*
 uniformly illuminated, **I**:30.9, 30.10, 30.10*f*
 Aperture flash mode, **II**:39.7
 Aperture placement (in stray light suppression), **II**:7.5–7.10
 aperture stops, **II**:7.6–7.7, 7.7*f*, 7.8*f*
 field stops, **II**:7.7, 7.8*f*, 7.9*f*, 34.18–34.19, 34.19*f*
 Lyot stops, **II**:7.8*f*–7.11*f*, 7.8–7.10
 Aperture stops, **I**:1.74, 1.75*f*, 17.8, 29.5, 29.36; **II**:1.4, 3.4, 7.6–7.7, 7.7*f*, 7.8*f*, 34.18, 34.19*f*
 Aperture-scanning microscopy, **I**:28.53–28.54
 Aphakia, **III**:12.14, 14.28
 correction of, **III**:12.14–12.15
 defined, **III**:9.1, 12.1, 14.1, 21.1
 Aphakic hazard, **II**:36.17
 Aplanatic (term), **I**:29.36
 Aplanatic lenses, **I**:17.5, 17.11*f*, 17.11–17.12
 Aplanatic objectives, **I**:29.11–29.13
 Aplanatic optical systems, **II**:34.19–34.20
 Aplanats, **II**:39.8, 39.9*f*
 Apodization, **II**:39.7; **III**:6.14, 21.17, 21.18; **V**:17.6–17.7, 17.7*f*
 defined, **III**:8.1
 Stiles-Crawford effect, **III**:8.8
 Apostilb (unit), **II**:34.43, 36.7, 36.8*t*
 Apovortex surfaces, **II**:39.11
 Apparent optical properties (AOPs), of water, **IV**:1.4, 1.5*t*–1.6*t*, 1.12–1.13
 Appliqués, black, **IV**:6.15
 Approximate transfer factor (ATF), **II**:4.1, 4.4
 Aqueous humor, **III**:1.3*f*, 14.5, 14.6, 14.26, 16.3
Arabidopsis thaliana, **V**:37.10*f*, 37.10–37.11
 Arbitrary phase profiles, **I**:23.8
 Arbitrary systems, paraxial matrix methods for, **I**:1.67
 Arc radiation sources (*see specific arcs, e.g.:* Argon arcs)
 Area image sensor arrays, **II**:32.24–32.32
 about, **II**:32.2
 CCD
 frame transfer, **II**:32.26–32.28, 32.27*f*, 32.28*f*
 interline transfer, **II**:32.28–32.32, 32.29*f*–32.31*f*
 performance of, **II**:32.32
 image area dimensions for, **II**:32.25*t*
 MOS, **II**:32.25–32.26, 32.26*f*
 Area-detection x-Ray imaging, **V**:61.2, 61.2*f*
 Area-solid-angle-product, **I**:1.22; **II**:39.5 (*see also Étendue*)
 Argon arcs, **II**:15.12, 15.13
 Argon ion lasers, **II**:16.14, 16.15*f*, 16.30
 Arrayed waveguide gratings (AWGs), **I**:21.24
 Array-mode selection, semiconductor, **II**:19.27
 Array-mode stability, semiconductor, **II**:19.27
 ARROW structures, **V**:11.15
 Arsenic antisite (As_{Ga}), **IV**:18.3
 Arsenic triselenide glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*
 Arsenic trisulfide (As_2S_3), **V**:12.7
 Arsenic trisulfide glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*
 Artificial sources (of radiation), **II**:15.3–15.53
 about, **II**:15.3–15.4
 commercial, **II**:15.13–15.53
 activated-phosphor sources, **II**:15.49
 blackbody simulators, **II**:15.14, 15.15*f*, 15.16*f*
 carbon arc sources, **II**:15.21–15.24, 15.23*f*, 15.24*f*, 15.25*t*–15.27*t*, 15.28*f*
 concentrated arc lamps, **II**:15.47–15.48, 15.48*f*, 15.49*f*
 glow modulator tubes, **II**:15.49, 15.50*f*, 15.51*f*, 15.52*t*

- Artificial sources (of radiation), commercial
(*Cont.*):
 high-energy sources, **II**:15.40
 high-pressure enclosed arc, **II**:15.24,
 15.28–15.34, 15.29*f*–15.35*f*
 hydrogen and deuterium arc lamps,
II:15.49, 15.53*f*
 incandescent nongaseous sources,
II:15.15–15.21, 15.17*f*–15.22*f*
 low-pressure enclosed arc, **II**:15.35–15.47,
 15.36*f*, 15.36*t*–15.43*t*, 15.44*f*–15.47*f*,
 15.46*t*, 15.47*t*
 special-purpose sources, **II**:15.53
 luminaire optics for, **II**:40.45*f*, 40.45–40.47,
 40.46*f*
 and radiation law, **II**:15.4–15.7, 15.5*f*, 15.5*t*,
 15.6*f*
 standardized laboratory sources, **II**:15.7–15.13
 baseline standard of radiation,
II:15.9, 15.9*f*, 15.10*f*, 15.12*f*
 blackbody cavity theory, **II**:15.7–15.9, 15.8*f*
 working standards of radiation,
II:15.9–15.13, 15.10*f*, 15.12*f*, 15.13*f*
 ASAP (optical software), **II**:7.25
 ASML Alpha Tool, **V**:34.2, 34.3*f*
 ASP (stray light analysis program), **IV**:6.19
 Aspheric intraocular lenses, **III**:21.1, 21.10–21.12,
 21.11*f*, 21.11*t*, 21.12*t*
 Aspheric lenses, **II**:39.8, 39.9, 39.9*f*; **III**:20.1,
 20.20–20.23, 20.21*f*–20.23*f*
 Aspheric measuring system, **II**:10.12*f*
 Aspheric surfaces, **II**:3.5; **IV**:3.8–3.9
 Aspherical optics fabrication, **II**:9.7*f*, 9.7–9.8
 Aspherical surfaces:
 and axial gradients, **I**:24.3
 and reflective/catadioptric objectives, **I**:29.3
 in systems of revolution, **I**:1.35
 Aspherical wavefront measurement,
II:13.23–13.27
 holographic compensators, **II**:13.25, 13.25*f*,
 13.26*f*
 infrared interferometry, **II**:13.25
 Moiré tests, **II**:13.26–13.27
 refractive or reflective compensators,
II:13.24, 13.24*f*, 13.25
 sub-Nyquist interferometry, **II**:13.27
 two-wavelength interferometry, **II**:13.25,
 13.26
 wavefront stitching, **II**:13.27, 13.27*f*
 Assembly of polymers, mechanical, **IV**:3.14*f*,
 3.14–3.16, 3.15*f*
 Assembly tolerances, **II**:5.8
 Associative memory, optical, **IV**:12.34
 Asthenopia (eyestrain), **III**:23.2
 Astigmatic dial, **III**:12.7*f*
 Astigmatic difference (term), **I**:1.43
 Astigmatism, **I**:1.90, 1.91, 29.34, 29.37; **II**:2.3,
 2.3*f*; **III**:1.6, 1.7, 1.7*f*, 1.18*f*, 12.4, 16.5–16.6;
V:13.12, 64.4, 64.6
 and binocular instrumentation,
III:1.41–1.42
 and contact lenses, **III**:20.15
 correction of, **III**:16.9
 cylindrical correction, **III**:13.16, 13.17
 with hydrogel contact lenses, **III**:12.13
 with spectacle lenses, **III**:12.9, 15.2
 with spherical contact lenses, **III**:12.12
 defined, **III**:12.1, 23.1
 determining axis of, **III**:12.6
 following cataract surgery, **III**:21.13
 irregular, **III**:16.6
 measuring, **III**:12.7
 off-axis, **III**:1.18, 1.18*f*
 ASTM (American Society for Testing
 Materials), **V**:1.4
 Astralate aberrations, **V**:45.4
 Astronomical telescopes, **I**:18.10; **II**:1.7*f*
 Astronomical x-ray optics, **V**:47.1–47.11
 angular resolution of, **V**:47.10–47.11, 47.11*f*
 hard, **V**:47.9–47.10
 history of, **V**:47.1–47.2
 Kirkpatrick-Baez optics, **V**:47.7–47.8, 47.8*f*,
 47.9*f*
 Wolter, **V**:47.2*f*, 47.2–47.7, 47.4*f*, 47.6*f*, 47.7*f*
 Astronomy:
 radio, **I**:5.23
 x-ray, **V**:33.1–33.4, 33.3*t*
 Asymmetric color matching, **III**:11.27, 11.29,
 11.30*f*
 and changes in overall intensity,
III:11.68
 and chromatic adaptation, **III**:11.69
 and color constancy, **III**:11.71
 Athermal glasses, **IV**:2.26
 Athermal laser beam expanders, **II**:8.13–8.14
 Athermalization, **II**:1.15–1.16, 1.16*f*, 6.22–6.24;
IV:3.9
 active, **II**:6.24, 6.24*f*
 intrinsic, **II**:8.7*f*, 8.7–8.8

- Athermalization (*Cont.*):
 mechanical, **II**:8.8–8.12
 active, **II**:8.11, 8.11*f*
 by image processing, **II**:8.12
 part active, part passive, **II**:8.11–8.12, 8.12*f*
 passive, **II**:8.8*f*–8.10*f*, 8.8–8.10
 optical, **II**:8.12–8.15, 8.13*t*
 about, **II**:8.12–8.13
 athermal laser beam expanders, **II**:8.13–8.14
 diffractive optics usage, **II**:8.15
 of separated components, **II**:8.14, 8.15
 three-material solutions, **II**:8.14,
 8.14*t*, 8.15*t*
 passive, **II**:6.22, 6.23*f*, 6.24
 single material design, **II**:6.22, 6.23*f*
- Athermalized achromatic lenses, **II**:1.16, 1.16*f*
- ATM (asynchronous transfer mode), fiber optic standards and, **V**:23.6
- Atmosphere, standard, **V**:3.6–3.11, 3.7*t*, 3.8*f*–3.11*f*
- Atmospheric Infrared Sounder (AIRS), **V**:3.39
- Atmospheric optics, **V**:3.3–3.45
 and absorption of light, **V**:3.4*f*, 3.4–3.5
 and atmospheric optical transmission, **V**:3.22*f*–3.27*f*, 3.22–3.26
 and composition of standard atmosphere, **V**:3.6–3.11, 3.7*t*, 3.8*f*–3.11*f*
 and global climate change, **V**:3.43–3.45, 3.44*f*
 and meteorological optics, **V**:3.40–3.43, 3.41*f*–3.43*f*
 remote sensing in, **V**:3.36–3.40, 3.37*f*–3.40*f*
 and theory of interaction of light with atmosphere, **V**:3.11–3.22
 inelastic optical processes, **V**:3.21–3.22
 Mie scattering, **V**:3.16–3.18, 3.17*f*–3.19*f*
 molecular absorption, **V**:3.12–3.15, 3.13*f*
 molecular emission and thermal spectral radiance, **V**:3.18, 3.20, 3.20*f*
 molecular Rayleigh scattering, **V**:3.15–3.16
 surface reflectivity and multiple scattering, **V**:3.21, 3.21*f*
 turbulence in, **V**:3.26, 3.28–3.36
 beam spreading, **V**:3.32–3.33, 3.33*f*
 beam wander, **V**:3.31–3.32
 imaging and heterodyne detection, **V**:3.34
 parameters for, **V**:3.28–3.31, 3.29*f*, 3.30*f*
 scintillation, **V**:3.34–3.36, 3.35*f*
- Atmospheric particles, scattering by, **I**:7.2
- Atmospheric tilt, **V**:5.14, 5.14*f*–5.15*f*
- Atmospheric turbulence, imaging through, **V**:4.1–4.37
 aberration variance and approximate Strehl ratio for, **V**:4.27–4.28, 4.28*f*
 and adaptive optics, **V**:4.35–4.36
 angle of arrival fluctuations, **V**:4.23–4.26, 4.25*f*, 4.26*f*
 and covariance and variance of expansion coefficients, **V**:4.20–4.22, 4.21*t*, 4.22*f*–4.23*f*
 Kolmogorov turbulence and atmospheric coherence length, **V**:4.7–4.10, 4.8*f*, 4.9*f*
 long-exposure images, **V**:4.3–4.7
 modal correction of turbulence, **V**:4.28–4.30, 4.29*t*, 4.30*f*
 and modal expansion of aberration function, **V**:4.17*f*–4.18*f*, 4.17–4.20, 4.19*t*, 4.20*t*
 and resolution of telescopes, **V**:4.2–4.3
 short-exposure image, **V**:4.31*f*–4.34*f*, 4.31–4.35, 4.35*t*
 and systems with annular pupils, **V**:4.10–4.16, 4.11*f*–4.15*f*, 4.15*t*
- Atom interferometry, **IV**:11.22–11.23, 11.24*f*
- Atomic beams:
 brightening of, **IV**:20.27*f*, 20.27–20.28
 collimation of, **IV**:20.15*f*, 20.15–20.16
 slowing of, **IV**:20.11–20.13, 20.12*f*, 20.12*t*, 20.13*f*
- Atomic clocks, **IV**:20.28
- Atomic coherence, maximal, **IV**:14.28–14.32, 14.29*f*–14.32*f*
- Atomic energy levels, **V**:2.2–2.5, 2.4*f*
- Atomic funnels, **IV**:20.27
- Atomic ionization, **IV**:21.3
- Atomic layer deposition, **IV**:7.11
- Atomic (gain) noise, **II**:23.34–23.35
- Atomic oxygen, black surfaces and, **IV**:6.16–6.17
- Atomic resonance, **IV**:22.2–22.9, 22.3*f*
 about, **IV**:22.2–22.5, 22.3*f*
 double, **IV**:22.5–22.9
 features of, **IV**:22.5–22.6, 22.6*f*
 tunable double resonance, **IV**:22.6–22.9, 22.7*f*, 22.8*f*
- Atomic scattering, x-ray optics and, **V**:36.1, 36.2*f*
- Atomic spectra, **V**:2.13
- Atomic systems, coherence in, **IV**:14.4*f*, 14.4–14.5

- Atom-laser interactions:
- applications of, **IV**:20.26–20.39
 - atomic beam brightening, **IV**:20.27*f*, 20.27–20.28
 - atomic clocks, **IV**:20.28
 - Bose-Einstein condensation, **IV**:20.35–20.37, 20.36*f*
 - dark states, **IV**:20.37–20.39, 20.38*f*
 - optical lattices, **IV**:20.31–20.34, 20.32*f*–20.34*f*
 - ultracold collisions, **IV**:20.28–20.31, 20.30*f*, 20.31*f*
 - cooling atoms with, **IV**:20.3–20.21
 - below Doppler limit, **IV**:20.17–20.21, 20.18*f*–20.20*f*
 - history of, **IV**:20.3–20.4
 - optical molasses, **IV**:20.13–20.17, 20.14*f*–20.16*f*
 - properties of lasers, **IV**:20.4–20.6
 - slowing atomic beams, **IV**:20.11–20.13, 20.12*f*, 20.12*t*, 20.13*f*
 - theoretical description, **IV**:20.6–20.11, 20.9*f*
 - trapping atoms with, **IV**:20.21–20.39
 - magnetic traps, **IV**:20.21–20.23, 20.22*f*
 - magneto-optical traps, **IV**:20.24*f*, 20.24–20.25, 20.26*f*
 - optical traps, **IV**:20.23*f*, 20.23–20.24
- Atoms (generally):
- electronic structure of, **I**:10.12–10.16, 10.13*f*–10.15*f*
 - in motion, **IV**:20.8–20.10, 20.9*f*
 - multielectron, **I**:10.10–10.11
 - one-electron, **I**:10.7–10.9, 10.8*f*, 10.9*f*
 - spectra of, **I**:10.3
 - trapping of neutral (*see* Trapping atoms)
 - two-level, **IV**:20.6–20.8
- Atoms, strong field interactions with,
- IV**:21.10–21.21
 - above threshold ionization, **IV**:21.14–21.17, 21.15*f*, 21.16*f*
 - ionization stabilization, **IV**:21.20–21.21, 21.22*f*
 - Keldysh parameter, **IV**:21.10
 - multiphoton and quasi-classical regimes, **IV**:21.10
 - multiphoton ionization, **IV**:21.10–21.12, 21.11*f*
 - relativistic effects, **IV**:21.19–21.20, 21.21*f*
 - Atoms, strong field interactions with (*Cont.*):
 - rescattering effects, **IV**:21.18*f*, 21.18–21.19, 21.19*f*
 - tunnel ionization, **IV**:21.12*f*, 21.12–21.14, 21.14*f*
 - Atrophic (dry) age-related macular degeneration, **III**:14.1
 - AT&T (American Telephone & Telegraph), **V**:21.1
 - Attention, in human vision, **III**:24.6, 24.7
 - Attenuated total reflectance (ATR) waveguides, **V**:12.11
 - Attenuation:
 - defined, **V**:21.13
 - fiber, **V**:21.13–21.14
 - for fiber optic communication links, **V**:15.7
 - Lambert-Beer law of, **V**:63.11
 - linear, **V**:31.1, 31.2
 - neutron, **V**:63.11–63.12
 - in optical fibers, **V**:9.4, 9.5*f*
 - in photonic crystal fibers, **V**:11.19–11.22, 11.20*f*, 11.21*f*
 - in water
 - beam, **IV**:1.40–1.41, 1.41*f*, 1.42*f*
 - diffuse and Jerlov water types, **IV**:1.42–1.46, 1.43*t*–1.45*t*, 1.44*f*, 1.45*f*
 - x-ray, **V**:31.1–31.4, 31.2*f*, 31.3*f*
 - Attenuation functions, of water, **IV**:1.13
 - Attenuators:
 - for fiber-based couplers, **V**:16.4
 - for networking, **V**:18.2, 18.9
 - neutral, **IV**:7.105, 7.105*f*
 - variable optical, **V**:21.12, 21.13*f*
 - Attosecond optics, **II**:21.1–21.9
 - about, **II**:21.2
 - driving lasers in, **II**:21.4–21.6
 - carrier-envelope offset frequency, **II**:21.5
 - carrier-envelope phase, **II**:21.5*f*, 21.6
 - carrier-envelope phasemeter, **II**:21.6
 - chirped pulse amplification, **II**:21.5
 - chirped pulse amplifiers, **II**:21.6
 - single-shot *f*-to-*2f* interferometer, **II**:21.6
 - high-harmonic generation, **II**:21.2, 21.2*f*
 - phase-matching in, **II**:21.4
 - ponderomotive potential in, **II**:21.3
 - pulse characterization, **II**:21.8*f*, 21.8–21.9
 - attosecond streak camera, **II**:21.9
 - FROG-CRAB, **II**:21.9
 - RABITT, **II**:21.9
 - second-order autocorrelator, **II**:21.9

- Attosecond optics (*Cont.*):
 pulse generation, **II**:21.6–21.8, 21.7*f*
 amplitude gating, **II**:21.7
 attosecond pulse train, **II**:21.6, 21.7
 double optical gating, **II**:21.8
 polarization gating, **II**:21.7–21.8
 two-color gating, **II**:21.7
 quantum trajectories in, **II**:21.3–21.4
 semiclassical model of, **II**:21.3
 single isolated pulses in, **II**:21.4
 strong field approximation in, **II**:21.3
- Attosecond pulse, **II**:21.8*f*, 21.8–21.9
 attosecond streak camera, **II**:21.9
 FROG, **II**:21.9
 FROG-CRAB, **II**:21.9
 generation of, **II**:21.6–21.8, 21.7*f*; **IV**:21.31
 amplitude gating, **II**:21.7
 attosecond pulse train, **II**:21.6, 21.7
 double optical gating, **II**:21.8
 polarization gating, **II**:21.7–21.8
 two-color gating, **II**:21.7
- RABITT, **II**:21.9
- Attosecond pulse train, **II**:21.6, 21.7
 Attosecond streak camera, **II**:21.9
 Auger cascades, **V**:59.4
 Auger energies, **V**:36.3*t*, 36.9*t*
 Auger excitation peaks, **V**:29.3
 Auger recombination, **V**:13.27–13.28
 Augmented resolution, **I**:30.12–30.14, 30.13*f*
 Augur recombination, **II**:19.17, 19.17*f*
- Autocollimeters:
 angle measurement with, **II**:12.11*f*,
 12.11–12.12, 12.12*f*
 curvature measurement with, **II**:12.19–12.20,
 12.20*f*
 defined, **II**:12.11–12.12
- Autocorrelation noise, **III**:18.11–18.12
 Autocovariance (ACV) function, **V**:44.13
 Autofocus, of cameras, **I**:25.11–25.15,
 25.12*f*–25.15*f*
 Autofocus SLRs, **I**:25.12–25.14, 25.13*f*
 Autokeratometers, **III**:12.6
 Automatic brightness control (ABC), **II**:31.18
 Automatic focusing, optical disks and,
I:35.12–35.14, 35.13*f*
 Automatic spherometers, **II**:12.19
 Automatic tracking, on optical disks,
I:35.14*f*–35.16*f*, 35.14–35.17
- Autorefractors, **III**:12.6
 Autoset levels (tools), **II**:12.14, 12.14*f*
- Avalanche multiplication, **II**:25.9
 Avalanche photodetectors (APDs):
 high-speed, **II**:26.17–26.20, 26.18*f*, 26.20*f*,
 26.21*f*
 improvements in, **II**:26.3
 Avalanche photodiode (APD) receivers, **V**:9.8,
 9.10–9.11
 Avalanche photodiodes, **II**:24.62–24.70,
 24.63*f*–24.70*f*, 24.72*f*, 24.72–24.73, 24.73*f*,
 25.8–25.10, 25.9*f*; **V**:13.63, 13.71–13.73
 defined, **II**:24.10
 germanium, **II**:24.70*f*, 24.72*f*, 24.72–24.73,
 24.73*f*
 InGaAs, **II**:24.66*f*–24.69*f*, 24.66–24.70
 silicon, **II**:24.62–24.65, 24.63*f*–24.66*f*
- Average degree of polarization (average DoP),
I:14.32–14.33
 Average unregistered detected point spread
 function (AUDPSF), **V**:44.14, 44.15*f*
- Avogadro's number, **II**:34.11
- Axial ametropia, **III**:20.1
 Axial and circumferential slope errors,
V:45.7–45.8
 Axial aniseikonia, **III**:13.17, 13.18
 Axial astigmatism, **I**:29.34, 29.37
 Axial chromatic aberration, **II**:1.14, 2.2, 2.3*f*
 Axial color, **I**:1.91, 29.9, 29.37
 Axial edge lift, **III**:20.1
 Axial gap prevention, **II**:6.21, 6.22*f*
 Axial gradient lenses, **I**:24.3–24.5, 24.4*f*
 Axial gradients, **III**:19.5
 Axial image point, **I**:1.27
 Axial (longitudinal) magnification, **I**:1.28, 1.52,
 17.5
 Axial rays, **II**:1.4, 1.11*f*, 1.12
 Axial resolution, **III**:18.1
 Axial thickness, of polymers, **IV**:3.10
 Axicons, **I**:11.7, 11.7*f*
 Axis wander, of prisms, **I**:13.15
 AxoScan Mueller matrix polarimeters, **I**:15.33
 Azimuth, **I**:15.10, 16.16
 Azimuth angle, **II**:35.5
 Azomethine dyes, **II**:30.10*f*
 excited state properties of, **II**:30.11–30.12,
 30.12*f*
 formation of, **II**:30.10
 photochemistry of, **II**:30.11
- Babinet compensators, **I**:13.53–13.55, 13.54*f*
 Babinet principle, **I**:3.9–3.11, 3.10*f*, 3.11*f*, 3.13

- Babinet-Soleil compensators, **I**:13.55*f*, 13.55–13.56; **V**:7.22
- Back central optical radius (BCOR), **III**:20.3
- Back focal length (BFL), of camera lenses, **I**:27.2, 27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
- Back light, **II**:40.43, 40.44*f*
- Back vertex power (BVP), **III**:12.1, 20.1
 contact lenses, **III**:20.7–20.8, 20.8*t*
 spectacle lenses, **III**:12.9
- Background temperature, **II**:24.10
- Background-limited performance (BLIP), of infrared detector arrays, **II**:33.24
- Backing, film, **II**:29.4
- Backlighting, **II**:40.1, 40.12, 40.47, 40.47*f*, 40.48*f*
- Backscatter and backscattering, **II**:20.13–20.15, 20.15*f*
 coherent, **I**:9.14*f*, 9.14–9.15
 enhanced, **I**:6.5*f*, 6.5–6.7
- Backward Brillouin scattering, **V**:11.25
- Backward Raman amplifiers, **IV**:15.4, 15.4*f*
- Backward Raman generators, **IV**:15.4, 15.4*f*
- Backward Raman scattering, **IV**:15.41, 21.38*f*, 21.39
- Backward trace, **II**:39.7
- Bacteria, in water, **IV**:1.14
- Baffles:
 cone-shaped secondary, **II**:7.3*f*, 7.3–7.4, 7.4*f*
 design of, **IV**:6.19
 in extreme environments, **IV**:6.18
 with integrating cavities, **II**:39.26
 in lighting design, **II**:40.41, 40.45*f*, 40.46
 selection process for black, **IV**:6.10*f*–6.11*f*, 6.10–6.12, 6.12*t*–6.13*t*
 shields for, **II**:7.9*f*, 7.9–7.10, 7.10*f*
 in stray light suppression, **II**:7.10, 7.11
 two-stage, **II**:7.10
- Baker relays, **I**:18.18, 18.18*f*
- Baker super-Schmidt objective, **I**:29.21
- Baker-Nunn objective, **I**:29.22
- Balanced spherical aberrations, **II**:11.30
- Ball Black, **IV**:6.5*f*, 6.6*f*, 6.50, 6.51*f*, 6.53*f*
- Ballasts:
 in fluorescent lamps, **II**:40.32–40.33
 in HID lamps, **II**:40.36
- Ballistic neutron guides, **V**:63.17–63.18
- Balloon payloads, in astronomical x-ray optics, **V**:47.10
- Balmer α -spectra, **I**:10.7–10.8, 10.8*f*
- Balmer β -transition, **I**:10.9
- Band edges, of photonic crystal fibers, **V**:11.9–11.10
- Band gap energy, **V**:19.2
- Band pass, **II**:38.8
- Band structures:
 defined, **IV**:9.3
 of solids, **IV**:8.24–8.27, 8.26*f*
- Band-filling modulators, **V**:13.62
- Bandpass filters, **IV**:7.73–7.96
 about, **IV**:7.73, 7.77*f*–7.78*f*, 7.77–7.78
 angular properties of, **IV**:7.91–7.94, 7.92*f*, 7.93*f*
 with multiple peaks, **IV**:7.90, 7.91*f*
 narrow- and medium-, **IV**:7.78–7.83, 7.79*f*, 7.80*f*, 7.82*f*–7.88*f*
 nonpolarizing, **IV**:7.66, 7.67*f*
 square-top multicavity, **IV**:7.82*f*–7.88*f*, 7.82–7.83
 stability and temperature dependence of, **IV**:7.94
 very narrow, **IV**:7.83, 7.88–7.89, 7.89*f*
 wedge filters, **IV**:7.90, 7.91, 7.91*f*
 wide, **IV**:7.90, 7.90*f*
 wide-angle, **IV**:7.93*f*, 7.93–7.94
 for XUV and x-ray regions, **IV**:7.94–7.96, 7.95*f*–7.96*f*
- Bandpass response, in acousto-optic interaction, **V**:6.15–6.16
- Bandwidth:
 of amplifiers, **II**:27.10
 defined, **V**:20.1
 fiber, **V**:21.2*f*, 21.2–21.3
 gain-bandwidth, **II**:26.17
 normalization of, **II**:36.14–36.16, 36.15*t*, 36.16*f*
 of photomultipliers, **II**:27.7
 “Bang-bang” zoom, **II**:1.12
- Banker lamps, **II**:40.12, 40.46, 40.46*f*
- Bar spherometers, **II**:12.19, 12.19*f*
- Bar-code reading, **II**:17.34
- Bare source light, **II**:40.43
- Bare states, **IV**:14.3
- Barium beta borate (BBO), **IV**:2.38*t*, 2.46*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.75*t*, 17.1, 18.12
- Barium beta borate (BBO) optical parametric oscillators, **IV**:10.18, 10.19*f*
- Barium crown glass, **IV**:2.41*t*
- Barium dense flint glass, **IV**:2.42*t*
- Barium flint glass, **IV**:2.42*t*

- Barium fluoride (BaF_2), **IV**:2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.69*t*, 2.76*t*
- Barium strontium titanate (BST), **II**:28.11, 28.12
- Barium titanate (BaTiO_3), **IV**:2.38*t*, 2.44*t*, 2.47*t*, 2.55*t*, 2.76*t*, 12.13*t*, 12.14–12.16, 12.16*f*
- Barrel distortion, **I**:1.91
- Barrel length (BRL), of camera lenses, **I**:27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
- Barrier suppression ionization (BSI), **IV**:21.14, 21.14*f*
- Barrier suppression ionization thresholds, **IV**:21.26, 21.27*f*
- Baryta layer, **II**:30.5
- Base curve:
 - of contact lenses, **III**:12.1, 12.12
 - of spectacle lenses, **III**:12.1, 12.9
- Base curve radius (BCR):
 - for contact lenses, **III**:20.3, 20.4*f*, 20.5
 - defined, **III**:20.1
- Baseline, **III**:13.1
- Batman doping, **V**:25.20
- Batwing lenses, **II**:40.12
- Baud rate, **II**:17.33
- Beacon lamps, airway, **II**:15.11
- Beacons, laser, **V**:5.27–5.34
 - focus anisoplanatism, **V**:5.27–5.29, 5.28*f*–5.30*f*
 - mesospheric sodium laser beams, **V**:5.32–5.34, 5.33*f*
 - Rayleigh, **V**:5.30–5.32, 5.31*f*
- Beam attenuation, in water, **IV**:1.13, 1.40–1.41, 1.41*f*, 1.42*f*
- Beam attenuation coefficient, **IV**:1.7*f*, 1.10
- Beam cleanup, **IV**:12.30
- Beam deviation and displacement, **I**:19.2
- Beam diffusers, **I**:23.13
- Beam propagation, split-step, **IV**:12.10
- Beam propagation method (BPM), **I**:21.8
- Beam separators, **III**:5.10
- Beam shapers and shaping:
 - in binary optics, **I**:23.13
 - for optical disks, **I**:35.8, 35.9*f*, 35.10*f*
- Beam spatial coherence, **V**:58.4
- Beam splitters (BSs), **I**:13.41–13.42; **II**:13.7, 34.32; **III**:5.9*f*, 5.9–5.10; **IV**:7.61–7.67, 7.62*f*–7.68*f*, 23.2, 23.2*f*, 23.14; **V**:18.6, 18.6*f*
- Beam spreading, **V**:3.32–3.33, 3.33*f*
- Beam steering, **V**:6.27–6.29, 54.11 (*see* Agile beam steering)
- Beam transformers, **II**:39.18, 39.18*f*
- Beam walkoff time, **IV**:15.30
- Beam wander, **V**:3.31–3.32
- Beam-forming illumination systems, **II**:39.22, 39.39
- Beamlines, for multilayer Laue lenses, **V**:42.9*f*–42.12*f*, 42.9–42.10
- Beam-smearing faceted reflectors, **II**:39.39
- Beam-splitter gratings, **I**:23.11, 23.12
 - Dammann approach, for binary gratings, **I**:23.12
- Beam-splitter prisms, **I**:13.18–13.22
 - Foster, **I**:13.7, 13.18*f*, 13.21–13.22
 - Glan-Thompson, **I**:13.18*f*, 13.22
 - Rochon, **I**:13.7, 13.18*f*, 13.18–13.21, 13.24
 - Sénarmont, **I**:13.7, 13.18, 13.18*f*, 13.21
 - Wollaston, **I**:13.7, 13.18, 13.18*f*, 13.21, 13.24
- Beating (phenomena), **IV**:18.19
- Becke line, **I**:28.27
- Beer-Lambert law (Beer's law), **II**:16.9, 34.35, 38.5; **III**:2.7; **V**:3.11, 3.20, 3.21
- Beer's law, **IV**:8.7, 8.28
- Beilby-layer polarizers, **I**:13.28
- Bell-clamping (edging fabrication step), **II**:9.6
- Bellcore, **II**:19.39, 19.41
- Bend, of liquid crystals, **V**:8.22, 8.23*f*
- Bend loss, of photonic crystal fibers, **V**:11.21*f*, 11.21–11.22
- Bendable optics, **V**:50.3–50.4, 50.4*f*
- Bending magnet synchrotron radiation sources:
 - brightness of, **V**:55.9
 - power of, **V**:55.8*f*, 55.8–55.9
 - radiation from, **V**:55.3–55.6, 55.5*f*, 55.6*f*
- Benes architecture, for switches, **I**:21.34–21.35
- Bent crystals, **V**:39.5*f*–39.6*f*, 39.5–39.6
- BepiColombo mission, **V**:49.6
- Berremann calculus, **I**:15.32
- Bertrand lenses, **I**:28.8, 28.44*f*
- Bertrand-type Feussner prisms, **I**:13.23
- Beryllium:
 - absorptance of, **IV**:4.48*t*, 4.50*t*
 - optical constants for, **IV**:4.11
 - optical properties of, **IV**:4.12*t*, 4.21*f*, 4.26*f*
 - penetration depth of, **IV**:4.47*f*
 - physical properties of, **IV**:4.52*t*, 4.54*t*
 - reflectance of, **IV**:4.28*t*–4.29*t*, 4.45*f*, 4.46*f*
 - surfaces of, **IV**:6.51, 6.52, 6.53*f*, 6.58*f*
 - thermal properties of
 - coefficient of linear thermal expansion, **IV**:4.56*t*, 4.57*f*
 - elastic properties, **IV**:4.69*t*

- Beryllium, thermal properties of (*Cont.*):
 at room temperature, **IV**:4.55*t*
 specific heat, **IV**:4.65*t*, 4.68*f*
 strength and fracture properties, **IV**:4.70*t*
 thermal conductivity, **IV**:4.58*t*, 4.59*f*–4.60*f*
- Bessel functions, **I**:7.12, 7.14–7.15, 17.38, 28.17;
III:8.1, 8.13, 8.22; **V**:6.6
- Best-correction, **III**:1.6
- Beta cloth, **IV**:6.57, 6.58*f*
- Betatron resonance, **IV**:21.42–21.43
- Bevel gauges, **II**:12.10, 12.11, 12.11*f*
- Bevel placement (on vanes), **II**:7.13, 7.14*f*
- Bezold-Brücke hue shift, **III**:11.1, 11.67*f*,
 11.67–11.68
- Bias, of electroabsorption modulators, **V**:13.59
- Bias angle, **II**:31.13
- Biased *pin* photodetectors, **II**:26.6*f*
- Biaxial crystals, **IV**:8.8, 8.9*t*, 8.10, 8.10*f*
- Bichromatic test mixtures, sensitivity to,
III:11.39, 11.40
- Biconical reflectance, **II**:35.5*t*, 35.6*f*, 35.6*t*
- Bidirectional reflectance distribution function (BRDF, BDRF), **I**:8.4, 15.38*f*–15.40*f*,
 15.38–15.39; **II**:7.1, 7.18–7.19, 7.22, 35.5,
 35.5*t*, 35.6*f*, 35.6*t*, 35.13, 37.9; **V**:1.4–1.6,
 1.5*f*, 3.21
- Bidirectional scatter distribution function (BSDF), **II**:7.2, 7.23, 7.24*f*, 7.25*f*, 35.13;
IV:6.1, 6.9–6.10, 6.18–6.19; **V**:1.6–1.7
- Bidirectional scatter distribution function (BSDF) scatterometers, **V**:1.8, 1.8*f*
- Bidirectional transmittance distribution function (BTDF), **II**:35.3, 35.13; **V**:1.6
- Bifocal jump, with monocular magnification,
III:13.15*f*, 13.15–13.16
- Bifocal lenses, **I**:23.12, 23.12*f*; **III**:12.10*f*, 14.27
 contact lenses, **III**:12.13, 14.28
 intraocular lenses, **III**:14.28
 for presbyopia, **III**:12.8, 12.10
 vertical prism introduced by, **III**:13.28–13.29
- Bihemispherical reflectance, **II**:35.5*t*, 35.6*f*, 35.6*t*
- Billet's split lens, **I**:2.16, 2.17*f*
- Bimorph mirrors, **III**:15.1, 15.10, 15.10*f*, 15.11;
V:5.37, 5.37*f*, 50.4, 50.4*f*, 50.6*f*
- Binary digits, in OTDM networks, **V**:20.8
- Binary optics, **I**:23.1–23.17
 fabrication of
 mask layout, **I**:23.14–23.16, 23.15*f*, 23.15*t*
 micromachining techniques, **I**:23.16,
 23.16*f*, 23.17
- Binary optics (*Cont.*):
 and geometrical optics, **I**:23.2–23.9
 aberration correction, **I**:23.4–23.7, 23.5*f*,
 23.6*f*
 analytical models, **I**:23.2–23.4, 23.6
 micro-optics, **I**:23.7*f*–23.8*f*, 23.7–23.8
 optical performance, **I**:23.8–23.9, 23.9*f*,
 23.10*t*
 and scalar diffraction theory, **I**:23.10–23.13,
 23.11*t*, 23.12*f*, 23.13*f*
 and vector diffraction theory, **I**:23.13–23.14,
 23.14*f*
- Binary units, for electro-optic modulators,
V:7.27
- Binning, **II**:38.10
- Binocular cues, **III**:13.3, 13.7
- Binocular disparities, **III**:2.40, 13.1, 13.4–13.5
 induced by prisms or lenses, **III**:13.28
 and perceived direction, **III**:13.7–13.8
- Binocular field, **III**:1.3
 horizontal angular extent of, **III**:1.38*f*
 stereopsis in, **III**:1.38–1.42
- Binocular fusion, **III**:13.1, 13.12–13.13
- Binocular instrumentation:
 and chromostereopsis, **III**:1.20
 differential focusing for, **III**:1.7
 tolerances in, **III**:1.41–1.42
- Binocular parallax, **III**:13.1, 13.22
- Binocular rivalry, **III**:13.1, 13.19
- Binocular rivalry suppression, **III**:13.12, 13.33
- Binocular stereoscopic discrimination,
III:2.40–2.41, 2.41*f*
- Binocular vision factors, **III**:13.1–13.35
 in computer work, **III**:23.11
 coordination of eyes, **III**:13.20–13.25
 binocular parallax, **III**:13.22
 cross-coupling and direction/distance of
 gaze, **III**:13.23*f*, 13.23–13.24
 defocus and efforts to clear vision,
III:13.22–13.23
 intrinsic stimuli to vergence, **III**:13.21–13.22
 perceived distance, **III**:13.24
 zone of clear and single binocular vision,
III:13.24–13.25, 13.25*f*
 distortion by monocular magnification,
III:13.13–13.16
 bifocal jump, **III**:13.15*f*, 13.15–13.16
 from convergence responses to prism,
III:13.19

- Binocular vision factors, distortion by
monocular magnification (*Cont.*):
discrepant views of objects/images,
III:13.16
motion parallax, **III**:13.14–13.15
perspective distortion, **III**:13.13, 13.14*f*
stereopsis, **III**:13.13, 13.14
distortion from interocular aniso-
magnification, **III**:13.16–13.19
aniseikonia, **III**:13.17–13.18
interocular blur suppression with
anisometropia, **III**:13.18–13.19
lenses and prisms, **III**:13.16–13.17, 13.17*f*
eye alignment, **III**:13.20–13.25
magnification induced errors of,
III:13.25–13.27
prism induced errors of, **III**:13.27–13.29
eye movements, **III**:13.19–13.20
focus and responses to distance, **III**:13.30
fusion and suppression, **III**:13.12–13.13
gaze control, **III**:13.29–13.30
and head mounted visual display systems,
III:13.31–13.35
distance conflicts, **III**:13.31–13.32
optical errors, **III**:13.34–13.35
spatial location conflicts, **III**:13.33, 13.34*f*
visual-vestibular conflicts, **III**:13.32–13.33
lens effects on vergence and phoria,
III:13.25–13.27
perceived direction, **III**:13.7–13.10
corresponding retinal points, **III**:13.8
horopter, **III**:13.8–13.9, 13.9*f*, 13.10*f*
vertical horopter, **III**:13.9
perceived space, **III**:13.3–13.7
binocular cues, **III**:13.7
extraretinal information for eye
movements, **III**:13.7
kinetic cues, **III**:13.4–13.7, 13.5*f*, 13.6*f*
monocular cues, **III**:13.3–13.7
prisms
distortion from interocular aniso-
magnification, **III**:13.16–13.17
effects on vergence and phoria,
III:13.25–13.27
errors of alignment with, **III**:13.25–13.29
in refractive errors, **III**:12.15–12.17
aniseikonia, **III**:12.16–12.17
anisometropia, **III**:12.16–12.17
convergence and accommodation,
III:12.15–12.16
Binocular vision factors (*Cont.*):
stereopsis, **III**:13.11*f*, 13.11–13.12
visual field, **III**:13.3
Binocularity, **III**:23.1
Binoculars, **I**:18.13–18.14, 18.14*f*
Binomial vectors, of space curves, **I**:1.19
Biological waveguides, **III**:8.1–8.29
in cochlear hair cells and human hair,
III:8.24–8.26
in fiber-optic plant tissues, **III**:8.26–8.28,
8.27*f*
models of, **III**:8.8–8.15, 8.13*f*, 8.14*f*
physical assumptions in, **III**:8.22
and propagation of light, **III**:8.21, 8.22
retinal layer of rods and cones, **III**:8.8–8.9,
8.9*f*, 8.10*f*, 8.12–8.15
three-segment model, **III**:8.9, 8.11–8.15
and photoreceptors, **III**:8.3–8.5
modal patterns in monkey/human
receptors, **III**:8.19–8.24, 8.22*f*
orientation and alignment, **III**:8.5–8.8,
8.6*f*, 8.7*f*
photoreceptor optics, **III**:8.3
quantitative observations of single receptors,
III:8.15, 8.16*f*, 8.17, 8.18*f*, 8.19*f*
in sponges, **III**:8.28–8.29
and Stiles-Crawford effect of the first kind,
III:8.3
Bio-optical models, of absorption, **IV**:1.27*f*,
1.27–1.29, 1.28*t*
Biotar lenses, **I**:17.28
Biphase coding, **V**:20.9, 20.9*f*
Biplates, **I**:13.56
Bipolar cells (retina), **III**:2.9, 2.10
Bipolar mechanisms (color vision), **III**:11.1,
11.80–11.81
Bipolar transistors, **II**:27.11
Bird-wing mirror, **IV**:12.7, 12.8*f*
Birefringence, **I**:15.6, 15.41; **III**:1.10, 18.1,
18.20, 18.22, 18.25, 18.27; **IV**:8.9, 17.1;
V:6.17, 8.19, 11.17
Birefringent diffraction bandshapes, **V**:6.13,
6.14, 6.14*f*
Birefringent fibers, **V**:16.5–16.6
Birefringent phased array deflectors, **V**:6.29
Birefringent tangential phase matching, **V**:6.25
Bismuth germanium oxide ($\text{Bi}_{12}\text{GeO}_{20}$) (BGO),
IV:2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*
Bismuth triborate (BiB_3O_6) (BIBO), **IV**:2.38*t*,
2.46*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.75*t*, 17.14

- Bistable optical switches, **IV**:16.31
- Bistatic radar cross-section (RCS), in surface scattering, **I**:8.3, 8.4
- Bit error rate (BER):
 DQPSK and, **V**:21.34
 for fiber optic communications, **V**:15.1–15.5, 15.3*f*, 15.4*f*, 15.8, 15.13, 15.15, 15.17
- Bit rate, for optical fibers, **V**:9.12
- Bitoric lenses, **III**:20.1, 20.17–20.20, 20.18*f*–20.20*f*
- Black absorbers, **IV**:7.104, 7.105*f*
- Black chrome, **IV**:6.53, 6.54, 6.54*f*
- Black coatings, **IV**:6.13*t*
- Black cobalt, **IV**:6.53, 6.54*f*
- Black dye, **IV**:6.15
- Black felt contact paper, **IV**:6.30*f*, 6.32*f*
- Black glass, **IV**:6.57
- Black Kapton, **IV**:6.57, 6.57*f*
- Black layer system, **IV**:6.15
- Black nickel, **IV**:6.21, 6.23*f*
- Black paint, **IV**:6.21, 6.24*f*
- Black surfaces, **IV**:6.1–6.59
 creation of, **IV**:6.13–6.16
 environmental degradation of, **IV**:6.16–6.18
 atomic oxygen effects, **IV**:6.16–6.17
 extreme environments, **IV**:6.18
 outgassing, **IV**:6.17
 particle generation, **IV**:6.17–6.18
 for far-infrared applications, **IV**:6.21, 6.26–6.34, 6.28*f*–6.34*f*
 Ames 24E and 24E2, **IV**:6.26*f*, 6.28*f*, 6.34
 Cornell Black, **IV**:6.26*f*, 6.27
 Infrablack, **IV**:6.26*f*, 6.28, 6.28*f*
 multiple-layer approach, **IV**:6.26, 6.26*f*–6.27*f*
 Teflon overcoat, **IV**:6.27
 optical characterization of, **IV**:6.18–6.21, 6.20*t*
 paints and surface treatments, **IV**:6.35–6.58, 6.37*f*, 6.43*f*, 6.53*f*
 Acktar black coatings, **IV**:6.55
 Aeroglaze Z series, **IV**:6.36*f*, 6.37, 6.37*f*, 6.39, 6.39*f*–6.42*f*
 Akzo Nobel paints, **IV**:6.39, 6.42*f*, 6.43*f*
 anodized processes, **IV**:6.44–6.49, 6.47*f*, 6.48*f*, 6.51*f*, 6.53*f*
 black glass, **IV**:6.57
 Black Kapton, **IV**:6.57, 6.57*f*
 carbon nanotubes and nanostructured materials, **IV**:6.55, 6.59*f*
- Black surfaces, paints and surface treatments (*Cont.*):
 Cardinal Black, **IV**:6.36*f*, 6.39, 6.44*f*
 Cat-a-lac Black, **IV**:6.39, 6.42*f*, 6.53*f*
 DeSoto Black, **IV**:6.37*f*, 6.39
 DURACON, **IV**:6.55–6.56
 electrically conductive black paint, **IV**:6.56
 electrodeposited surfaces, **IV**:6.53–6.54, 6.54*f*, 6.55*f*
 etching of electroless nickel, **IV**:6.49–6.50, 6.50*f*, 6.51*f*, 6.53*f*
 flame-sprayed aluminum, **IV**:6.57
 Floquil, **IV**:6.44
 gold blacks, **IV**:6.57
 high-resistivity coatings, **IV**:6.56
 IBM Black (tungsten hexafluoride), **IV**:6.56
 ion beam-sputtered surfaces, **IV**:6.53
 Parson's Black, **IV**:6.44, 6.53*f*
 plasma-sprayed surfaces, **IV**:6.50–6.52, 6.51*f*–6.53*f*
 silicon carbide, **IV**:6.56
 SolarChem, **IV**:6.44, 6.48*f*, 6.53*f*
 specular metallic anodized surfaces, **IV**:6.57, 6.58*f*
 sputtered and CVD surfaces, **IV**:6.56
 3M paints and derivatives, **IV**:6.35–6.37, 6.36*f*, 6.38*f*, 6.53*f*
 ZO-MOD BLACK, **IV**:6.56
 selection process for, **IV**:6.10*f*–6.11*f*, 6.10–6.12, 6.12*t*–6.13*t*
 and substrates, **IV**:6.34–6.35
 types and morphologies of, **IV**:6.1–6.10, 6.2*t*–6.4*t*, 6.5*f*–6.8*f*
 for ultraviolet applications, **IV**:6.21, 6.22*f*–6.25*f*
- Black Tedlar, **IV**:6.57*f*
- Black velvet cloth, **IV**:6.31*f*, 6.33*f*
- Black-and-white (B&W) film, **II**:29.4, 30.24–30.25, 30.25*t*
- Blackbody cavity theory, **II**:15.7–15.9
- Blackbody D star, **II**:24.10
- Blackbody detectivity, **II**:24.10
- Blackbody noise-equivalent power, **II**:24.10
- Blackbody radiation, **II**:15.4–15.6, 15.5*t*
 emittance of, **II**:34.25–34.26
 sources of, **II**:15.14, 15.15*f*, 15.16*f*, 34.23–34.24
 temperature vs., **II**:36.12, 36.12*f*, 36.14, 36.14*f*
 working standards for, **II**:15.14, 15.16*f*

- Blackbody responsivity, **II**:24.10
- Blackbody simulators, **II**:15.14, 15.15*f*, 15.16*f*, 34.24–34.26
- Black-light fluorescent lamps, **II**:15.35, 15.36*t*
- Blackman window function, **V**:46.9, 46.10*f*
- Blazing, gratings and, **IV**:5.60
- Bleaching, in film development, **II**:29.14
- Blindness:
- from cataract, **III**:14.24, 21.3
 - from diabetic retinopathy, **III**:14.25
 - economic impact of, **III**:14.4
 - flash, **II**:40.9
 - from glaucoma, **III**:14.26
 - and increasing life span, **III**:14.3
 - from uncorrected ametropia, **III**:12.4
- Blink response (of eye), **IV**:13.1
- Blip detector (blip condition), **II**:24.10
- Bloch equations, optical, **IV**:11.3–11.6
- Bloch solution, **IV**:8.25
- Bloch sphere, **IV**:11.4
- Bloch waves, **V**:11.22
- Bloch's law, **III**:2.28
- Blockage defects, in polycapillary x-ray optics, **V**:53.4, 53.5, 53.5*f*
- Blocked impurity band (BIB), **II**:33.7
- Blocking contacts, **II**:26.3
- Blocking filters, **II**:38.8
- Blocks (optical disk data), **I**:35.6–35.7
- Blood gas analysis, **II**:17.34
- Blooming:
- antiblooming, **II**:32.9, 32.10*f*
 - in image sensors, **II**:32.6, 32.9
 - thermal, **IV**:16.22
- Blue emitters, in LED technology, **II**:17.18, 17.19
- Blue light, color film and, **II**:29.13, 29.13*f*, 30.3–30.4
- Blue semiconductor lasers, **II**:19.7
- Blue to red (B/R) ratio, **III**:8.28
- Blue-enhanced photodiodes, **II**:24.55*f*, 24.61*f*, 24.61–24.62, 24.62*f*
- Blue-light photochemical injury, **III**:7.4, 7.10
- Blur, **III**:1.28
- in correction for SCE-1, **III**:9.4–9.6, 9.14
 - with monovision, **III**:14.28
 - suppression of, **III**:13.18–13.19, 14.28
 - (see also Astigmatism; Defocus)
- Blur filters, **II**:32.34, 32.34*f*
- “Boat grown” technique, **II**:17.21
- Bode representation, of servo system, **II**:22.5–22.6, 22.6*f*
- Bohr frequency condition, **I**:10.4
- Bohr's theory of hydrogen, **I**:10.3
- Bolometers, **II**:24.5, 28.3–28.5, 28.4*f*
- about, **II**:28.1
 - carbon, **II**:28.5
 - detectivity of perfect, **II**:24.17, 24.18*f*
 - germanium low-temperature, **II**:24.31–24.32, 24.32*f*, 24.33*f*
 - indium antimonide hot-electron, **II**:24.29, 24.30, 24.30*f*, 24.31*f*
 - as infrared detectors, **II**:33.9–33.10
 - metal, **II**:28.4, 28.7*t*
 - properties of, **II**:28.7*t*
 - resistive arrays of, **II**:28.10*f*, 28.10–28.11
 - semiconductor, **II**:28.4–28.5
 - superconducting, **II**:28.5
 - thermistor, **II**:24.24*f*, 24.24–24.25, 24.25*f*, 28.7*t*
- Boltzmann population factor, **V**:3.14, 3.20
- Bombardment, of carbon surface, **IV**:6.8*f*
- Bonded mountings, **II**:6.13–6.15, 6.15*f*, 6.16*f*
- Boresight tolerances, **II**:5.8
- Bormann triangle, **V**:63.27
- Born approximation, **III**:8.1, 8.15; **V**:27.2, 63.4, 63.5
- Born series, **I**:9.4
- Born-Oppenheimer approximation, **I**:10.19, 10.20, 10.22
- Boron, **IV**:5.82*f*, 5.83
- Boron Black, **IV**:6.50, 6.51, 6.52*f*
- Boron carbide, **IV**:6.51
- Boron neutron capture therapy (BNCT), **V**:53.19
- Boron [$^{10}\text{B}(n, \alpha)$] reaction, **V**:63.31
- Boron-doped silicon (Si:B) detectors, **II**:24.95*f*, 24.96
- Borosilicate crown glass, **IV**:2.41*t*
- Bormann transmission, **V**:43.3, 43.3*f*
- Bose-Einstein condensation (BEC), **II**:23.39; **IV**:14.22, 20.26, 20.35–20.37, 20.36*f*
- Bose-Einstein distribution of states, **IV**:2.16
- Bose-Einstein statistics, **II**:23.9
- Bouillotte lamps, **II**:40.12, 40.46, 40.46*f*
- Boulder Damage Symposium, **IV**:19.1–19.2
- Boules, glass, **II**:9.3
- Bound excitons, **IV**:5.26*t*, 5.29, 5.46
- Bound modes, of optical waveguides, **I**:21.3

- Boundary conditions (of optics):
 defined, **II**:3.17
 methods for handling, **II**:3.18–3.19
 specification of, **II**:4.12
- Boundary quality, of multilayers, **V**:41.5–41.7, 41.7*t*
- Bound-electronic optical Kerr effect, **IV**:16.3*t*, 16.12–16.13, 16.13*f*
- Bovine lenses, **III**:19.8–19.11, 19.11*f*
- Bowman's layer, **III**:16.4
- Bowman's membrane (cornea), **III**:14.5
- Boxcar averaging, of modulated signal sources, **II**:27.13, 27.13*f*, 27.15
- Boynnton illusion, **III**:11.72, 11.73*f*, 11.78
- Brace half-shade plates, **I**:13.57
- Bragg angle, **V**:6.9, 30.5, 39.3
- Bragg cell spectrum analyzers, **I**:11.9–11.10, 11.10*f*
- Bragg cells, **I**:11.11–11.12, 11.19; **V**:6.27, 6.31*t*, 6.30
- Bragg condition, **V**:6.32, 42.3*f*, 42.4, 42.10, 42.11*f*
- Bragg diffraction, **I**:11.9, 11.9*f*
 and acousto-optic modulators, **V**:6.4, 6.6, 6.7, 6.14
 and brightness of x-ray tube sources, **V**:54.16
 far and near, **V**:6.8–6.9, 6.12
 in neutron optics, **V**:63.23
 order of, **V**:20.14
 and phase matching equations, **V**:6.11
 of x-rays, **V**:42.2
- Bragg fibers, **V**:11.4
- Bragg geometry, for crystal monochromators, **V**:39.2*f*, 39.2–39.4, 39.3*f*, 39.6
- Bragg gratings, **IV**:22.9–22.11, 22.10*f*
 and DBR lasers, **V**:20.14
 fiber, **V**:17.1–17.9
 applications, **V**:17.8*f*, 17.8–17.9
 chirped, **V**:21.22–21.23, 21.23*f*, 21.25*f*, 21.25–21.26
 fabrication, **V**:17.4–17.8, 17.5*f*–17.7*f*
 and fiber lasers, **V**:25.8, 25.16, 25.18, 25.30, 25.31
 long-period gratings vs., **V**:24.9, 24.11
 photosensitivity, **V**:17.2–17.3
 properties of, **V**:17.3–17.4, 17.4*f*
 sensors based on, **V**:24.5–24.8, 24.6*f*–24.7*f*
 volume, **V**:25.29, 25.30
- Bragg limit, **V**:6.8
- Bragg mirrors, **V**:13.28, 13.44
- Bragg planes, of crystal monochromators, **V**:39.1, 39.2, 39.5
- Bragg reflection, **IV**:20.33
 of crystals, **V**:40.9
 in interferometers, **V**:63.26, 63.27
 and linear polarization, **V**:43.2–43.4, 43.3*f*, 43.4*f*, 43.6, 43.8
 and liquid crystals, **V**:8.10
 in monochromators, **V**:39.1, 63.24
 and multilayers, **V**:41.2
 and x-ray absorption spectroscopy, **V**:30.2, 30.4
- Bragg reflection filters, **I**:21.30
- Bragg reflection monochromator, **V**:39.4, 39.5
- Bragg reflector lasers, **V**:13.7, 13.28–13.29, 13.29*f*
- Bragg reflectors, **I**:21.8, 21.30; **II**:19.41; **V**:13.45
- Bragg regime, **I**:30.39–30.41, 30.41*f*, 30.42*f*
- Bragg scattering cross section, **V**:63.18
- Bragg transmission phase retarders, **V**:43.6
- Bragg wavelength, **V**:17.3, 20.15, 24.7
- Bragg-Brentano powder diffractometer, **V**:28.1, 28.2*f*, 28.3, 28.3*f*
- Bragg-Fresnel lenses, **V**:40.9*f*, 40.9–40.10
- Bragg's law, **V**:26.8, 28.3, 30.1, 39.1, 39.4, 43.7, 63.14, 63.16, 63.24
- Bragg-symmetric crystals, **V**:35.3
- Braking radiation (*see* Bremsstrahlung radiation)
- Brashear-Hastings prisms, **I**:19.3*t*, 19.25, 19.25*f*
- Bravais biplates, **I**:13.56
- Breault Research Organization, **IV**:6.1
- Bremsstrahlung heating, inverse, **IV**:21.37, 21.37*f*
- Bremsstrahlung photons, **V**:63.12
- Bremsstrahlung radiation, **V**:31.3
 continuous, **V**:54.4–54.6, 54.5*f*
 from laser generated plasmas, **V**:56.1, 56.8–56.10
 from pinch plasma sources, **V**:57.1
 and x-ray fluorescence, **V**:29.3, 29.5, 29.6, 29.11
- Brewster angle:
 and Airy equation, **I**:12.10
 defined, **I**:12.12
 and extinction ratio, **I**:12.18, 12.22, 12.22*f*, 12.23*f*, 12.24
 and polarization, **I**:12.15
 in VUV and x-ray region, **V**:41.9
- Brewster angle prisms, **I**:13.13

- Brewster angle reflection polarizers, **I**:12.16*f*,
12.16–12.18, 12.17*f*, 13.34–13.37,
13.34*t*–13.36*t*
- Brewster angle transmission polarizers,
I:12.18–12.24, 12.19*t*–**I**:12.20*t*, 12.21*f*,
13.37–13.39, 13.38*t*–13.39*t*
- Brewster's angle, **IV**:8.12, 8.23
- Bridge fiber method, of mode matching,
V:25.17, 25.17*f*
- Bridge mirror, **IV**:12.7, 12.8*f*
- Bridgeman technique, **II**:17.21
- Bright field microscopy, **I**:28.25, 28.27*f*,
28.27–28.28
- Bright (coupled) state, **IV**:14.4
- Brightening, atomic beam, **IV**:20.27*f*,
20.27–20.28
- Brightness:
of carbon arcs, **II**:15.23*f*
in color CRTs, **III**:22.7, 22.19
defined, **III**:11.1, 11.8, 23.1
luminance vs., **II**:34.40
and luminous efficiency, **III**:11.70, 11.70*f*
in monochrome CRTs, **III**:22.8*f*
perception of, **II**:40.4, 40.4*f*
of scene, **II**:31.4–31.5
of synchrotron radiation, **V**:55.1, 55.9
in visual acuity tests, **III**:4.8
x-ray tube sources, **V**:54.11–54.15,
54.13*f*–54.15*f*
- Brightness matching, **III**:10.43–10.45
- Brillouin frequency shift, **V**:10.8
- Brillouin gain, **IV**:15.48
- Brillouin scattering, **I**:31.30; **IV**:21.38
backward, **V**:11.25
for crystals and glasses, **IV**:2.27
defined, **IV**:2.27
forward, **V**:11.25, 11.26
in measurement, **IV**:5.76, 5.77
in nonlinear optics, **IV**:16.14, 16.18–16.19
photonic crystal fibers, **V**:11.25*f*, 11.25–11.26
Raman vs., **IV**:15.1
in solids, **IV**:8.18
stimulated, **V**:10.1, 10.7–10.9, 15.8, 21.20,
25.6 [see Stimulated Brillouin scattering
(SBS)]
in strong-field physics, **IV**:21.38
- Brillouin spectroscopy, **IV**:5.57–5.58
- Brillouin zone, **IV**:8.25, 8.26*f*, 8.27, 8.29, 9.3,
9.4; **V**:11.9
- Brillouin-enhanced four-wave mixing
(BEFWM), **IV**:15.53, 15.54, 15.54*f*
- British Glare Index (CIBSE), **II**:40.10
- Broad bandwidth solid-state lasers, **II**:16.34*f*,
16.34–16.35
- Broadband light sources, **IV**:5.58–5.59
- Broadband parametric amplification, **IV**:18.12
- Broadband reflectors, all-dielectric,
IV:7.45*f*–7.47*f*, 7.45–7.47
- Broadband SBS slow light, **IV**:22.15
- Broadband transient Raman scattering,
IV:15.28–15.32, 15.29*f*
- Broadband transmission, **V**:3.23–3.24,
3.25*f*–3.26*f*
- Broadening:
Doppler, **V**:3.14, 5.32, 56.5
homogeneous, **V**:20.1
of lineshapes, **I**:10.7
Lorentzian, **V**:3.14
pressure, **V**:56.5–56.6
spectral, **V**:56.6–56.8
in spectral lines, **IV**:14.13, 14.14
- Bromellite (BeO), **IV**:2.38*t*, 2.46*t*, 2.47*t*, 2.50*t*,
2.55*t*, 2.60*t*, 2.70*t*
- Bromine, in light bulbs, **II**:40.30
- Bromyrite (AgBr), **IV**:2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*,
2.55*t*, 2.60*t*, 2.69*t*, 2.76*t*
- Brownian fractals, **I**:8.9, 8.17
- Brownian movement, **I**:28.28
- Bruch's membrane, **III**:8.1, 14.25, 18.1
- Brunning distance-measuring interferometers,
II:12.8*f*, 12.8–12.9
- BS-39B glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- Bubble acceleration, **IV**:21.41–21.42, 21.42*f*
- Buckbee Mears wire-grid polarizers, **I**:13.31
- Buffered direct injection (BDI), **II**:33.19*t*,
33.20*f*, 33.21–33.22
- Build-and-test evaluation (for stray light
suppression), **II**:7.28
- Built-in potential, **II**:25.6
- Bulb blackening, **II**:40.30
- Bulb shield, **II**:40.45*f*, 40.46
- Bulk (term), **IV**:4.3
- Bulk compound semiconductors,
IV:12.20–12.21, 12.20*t*, 12.21*f*
- Bulk electro-optic modulators, **V**:7.3, 7.16–7.28
amplitude modulation, **V**:7.22–7.24, 7.23*f*,
7.24*f*
frequency modulation, **V**:7.24–7.25, 7.25*f*
phase modulation, **V**:7.18–7.20

- Bulk electro-optic modulators (*Cont.*):
 polarization modulation (dynamic retardation), **V**:7.20*f*, 7.20–7.22, 7.21*f*
 scanners, **V**:7.26*f*–7.28*f*, 7.26–7.28
- Bulk lasers, **V**:25.5–25.6
- Bulk material photodetectors, **II**:26.4*f*, 26.5
- Bulk modulus, for metals, **IV**:4.69*t*
- Bulk-grown materials, **II**:17.8
- Bumper function, **III**:11.59, 11.60*f*
- Bunches, electron, **V**:55.17
- Bunsen-Kirchhoff spectrometers, **I**:20.5*f*
- Bunsen-Roscoe law of photochemistry, **III**:7.1–7.3, 7.7
- Buried crescent lasers, **II**:19.24, 19.25*f*
- Buried heterostructure (BH) lasers, **II**:19.8, 19.9*f*, 19.20*t*, 19.24, 19.36*f*; **V**:13.5
- Buried TRS (BTRS) lasers, **II**:19.19, 19.20*t*, 19.21*f*
- Buried V-groove-substrate inner stripe (BVSIS), **II**:19.19, 19.20*t*
- Buried-channel CCDs (BCCDs), **II**:32.14, 33.13
- Buried-channel MOS capacitors, **II**:32.4*f*, 32.7–32.8
- Burnished mounting, **II**:6.3, 6.3*f*
- Butterfly scanners, **I**:30.50*f*, 30.51
- C line, **III**:25.1
- Cable television (CATV), **I**:21.2, 21.32–21.34; **II**:25.11–25.12; **V**:9.16–9.17
- Cadmium germanium diarsenide (CdGeAs₂), **IV**:2.38*t*, 2.45*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.61*t*, 2.74*t*
- Cadmium selenide (CdSe), **IV**:2.38*t*, 2.46*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.61*t*, 2.70*t*, 2.77*t*
- Cadmium selenide (CdSe) photoconductors, **II**:24.49–24.52, 24.52*f*
- Cadmium sulfide (CdS) photoconductors, **II**:24.49–24.52, 24.51*f*–24.53*f*
- Cadmium telluride (CdTe), **IV**:2.39*t*, 2.44*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.69*t*, 2.77*t*
- Cadmium telluride (CdTe) detectors, **II**:24.52, 24.54, 24.54*f*
- Cadmium zinc telluride (CdZnTe) detectors, **II**:24.52
- Caged compounds, in microscopy, **I**:28.55
- Calamitic liquid crystals (LCs), **V**:8.4, 8.5, 8.9, 8.11*f*
- Calar Alto telescope, **V**:5.27
- Calcite (CaCO₃), **IV**:2.38*t*, 2.46*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.61*t*, 2.74*t*, 2.77*t*
 double refraction in, **I**:13.2*f*–13.3*f*, 13.2–13.6, 13.4*t*–13.5*t*
 Feussner prisms of, **I**:13.23
 Rochon prisms of, **I**:13.20
- Calcium molybdate (powelite) (CaMoO₄), **IV**:2.38*t*, 2.45*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.61*t*, 2.72*t*, 2.77*t*
- Calibration:
 artificial sources of radiation for [*see* Artificial sources (of radiation)]
 of color CRTs, **III**:22.20–22.21 (*see also* Characterization, of CRTs)
 errors in, and SCE-1 correction, **III**:9.6–9.13
 legal traceability of, **II**:34.21
 in Maxwellian viewing systems, **III**:5.17*f*, 5.17–5.18
 photometric, **II**:34.42–34.43
 radiometric, **II**:34.31–34.32
 self-calibration, **II**:34.29
 spectroradiometric, **II**:38.11–38.13, 38.11*t*, 38.12*f*
- Calibration transfer devices, **II**:34.31–34.32
- Callier coefficient, **II**:29.7
- Camera formula, for solid-state cameras, **I**:26.13
- Camera lenses:
 classification system for, **I**:27.17, 27.23*t*, 27.24
 design limitations of, **I**:27.1–27.2
 fish-eye, **I**:27.6
 inverted telephoto, **I**:27.2, 27.6, 27.7*f*–27.14*f*
 extreme wide-angle, **I**:27.6, 27.13*f*, 27.14*f*
 highly complex extreme speed, **I**:27.6, 27.9*f*–27.12*f*
 very compact moderate speed, **I**:27.6, 27.7*f*–27.8*f*
 SLR normal lenses, **I**:27.2, 27.3*f*–27.4*f*
 telephoto lenses, **I**:27.6, 27.13, 27.15*f*–27.16*f*
 wide-angle lenses, **I**:27.2, 27.5*f*
 zoom lenses, **I**:27.17, 27.20*f*–27.22*f*
- Cameras, **I**:25.3–25.26
 accessories for, **I**:25.16–25.17, 25.17*f*
 adaptive optics retina cameras, **III**:15.3, 15.12 (*see also* Ophthalmoscopes)
 aerial, **I**:25.20
 Anger, **V**:63.33
 and autoexposure, **I**:25.10–25.11
 and autofocus, **I**:25.11–25.15, 25.12*f*–25.15*f*

- Cameras (*Cont.*):
 characteristics of, **I**:25.3–25.4
 clandestine, **I**:25.21
 color, **III**:10.38–10.40
 critical features of, **I**:25.8, 25.9*f*
 display types, **I**:25.6–25.7
 endoscopic, **I**:25.21, 25.21*f*
 features of, **I**:25.8, 25.18, 25.19*f*
 film for, **I**:25.5, 25.6
 and flash, **I**:25.16, 25.17*f*
 formats for, **I**:25.18
 gamma, **V**:32.2
 high-speed, **I**:25.21–25.22, 25.22*f*
 and image, **I**:25.5
 images from, **I**:25.7
 and instant pictures, **I**:25.8
 large-format, **I**:25.18–25.20
 and red eye, **I**:25.16
 and resolution of fine detail, **I**:25.5–25.6, 25.6*f*
 sewer, **I**:25.22–25.23
 solid-state (*see* Solid-state cameras)
 stereo, **I**:25.23–25.24, 25.23*f*
 streak, **I**:25.24, 25.24*f*
 thermal imaging, **I**:25.25
 and time lag, **I**: 25.8–25.9, 25.9*f*
 underwater, **I**:25.25
 video, **I**:25.7–25.8
 view, **I**:25.19*f*, 25.18–25.20
 for wide-angle photography, **I**:25.26*f*, 25.25
- Canada-France-Hawaii Telescope, **V**:5.23
- Candela (unit), **II**:34.37, 34.39, 36.4, 36.5, 37.3, 37.4
- Candle power, **II**:37.3
- Canon EOS A2E camera, **I**:25.14*f*, 25.15, 25.15*f*, 25.16
- Cantor sets, **I**:8.9
- Capacitive bolometers, **II**:33.10
- Capacitive transimpedance amplifier (CTIA), **II**:33.19*t*, 33.20*f*, 33.22–33.23
- Capacitors, MOS, **II**:32.4*f*, 32.7–32.8
- Capillary, **III**:18.1
- Capillary discharge devices, **V**:57.3, 57.3*t*
- Capillary mercury-arc lamps, **II**:15.30–15.31, 15.31*f*
- Capillary optics, **V**:28.5 (*see also* Monocapillary x-ray optics; Polycapillary optics)
- Capsulorhexis, **III**:21.1
- Carbamide [(NH₄)₂CO], **IV**:2.40*t*
- Carbon arc light sources, **II**:40.40
- Carbon arc sources (of radiation), **II**:15.21–15.24, 15.23*f*, 15.24*f*, 15.25*t*–15.27*t*, 15.28*f*
- Carbon black particles, **IV**:6.15
- Carbon bolometers, **II**:28.5, 28.7*t*
- Carbon deposition, laser-assisted, **IV**:19.5
- Carbon disulfide (CS₂), **IV**:16.13–16.14, 16.14*f*
- Carbon nanotubes (CNTs), **IV**:6.55, 6.59*f*; **V**:54.10–54.11
- Carbon surface bombarded with positive argon ions, **IV**:6.8*f*
- Carbon-black suspensions (CBSs), **IV**:13.10–13.11, 13.11*f*
- Carbon-dioxide lasers, **II**:16.16, 16.16*f*, 16.30; **V**:12.3*t*, 12.9, 12.13
- Cardinal Black, **IV**:6.36*f*, 6.39, 6.44*f*
- Cardinal directions (color vision), **III**:11.1, 11.56
- Cardinal mechanisms (color vision), **III**:11.1, 11.54, 11.56
 evidence for, **III**:11.79
 and sensitivity losses, **III**:11.29
- Cardinal points, of lenses, **I**:1.44, 17.7
- Carey Lea silver (CLS), **II**:29.13, 30.4
- Carl Zeiss prism system, **I**:19.3*t*, 19.16, 19.16*f*
- Carlisle objectives, **I**:29.8
- Carrier confinement, **II**:17.12, 17.12*f*–17.14*f*, 17.13, 17.17
- Carrier density, **II**:19.30–19.33
- Carrier effects, in integrated optics, **I**:21.10–21.12
- Carrier sense multiple access with collision detection (CSMA/CD), **V**:23.7
- Carrier transit time, **II**:26.6*f*, 26.6–26.7
- Carrier trapping, **II**:26.9, 26.9*f*; **IV**:18.21–18.23, 18.22*f*
- Carrier-carrier scattering, **IV**:18.20
- Carrier-envelope offset, **II**:20.4
- Carrier-envelope offset frequency, **II**:21.5
- Carrier-envelope (CE) phase:
 of chirped pulse amplifiers, **II**:21.6
 of lasers, **II**:21.5, 21.5*f*
- Carrier-envelope phasemeters, **II**:21.6
- Carrier-suppressed return-to-zero (CSRZ) formats, in WDM networks, **V**:21.30, 21.31*f*, 21.36*t*, 21.37*f*
- Carrier-to-noise ratio (CNR), **I**:35.24; **V**:9.15–9.16, 9.16*f*
- Cartesian coordinates, **I**:1.20, 1.21
- Cascaded limiters, **IV**:13.6
- Cascaded $x^{(1)}$: $x^{(1)}$ processes, of third-order optical nonlinearities, **IV**:16.20–16.22, 16.21*f*

- Cascaded $x^{(2)}:x^{(2)}$ processes, of third-order optical nonlinearities, **IV**:16.22–16.24, 16.23*f*, 16.24*f*
- Cassegrain design, **II**:7.3*f*, 7.11, 7.14, 7.14*f*, 7.16, 7.16*f*, 7.19, 7.20*f*
- Cassegrain objectives, **I**:29.6, 29.7
 afocal Cassegrain–Mersenne telescope, **I**:29.9
 dual magnification, **I**:29.9–29.10
 with field corrector and spherical secondary, **I**:29.8–29.9
 Houghton–Cassegrain, **I**:29.22–29.23
 Mangin–Cassegrain with correctors, **I**:29.24
 reflective Schmidt–Cassegrain, **I**:29.17
 Schmidt–Cassegrain, **I**:29.16–29.17
 Schmidt–meniscus Cassegrain, **I**:29.21
 with Schwarzschild relay, **I**:29.32
 solid Makutsov–Cassegrain, **I**:29.19
 spherical–primary, with reflective field corrector, **I**:29.9
 three-mirror, **I**:29.30
- Cassegrain telescopes, **I**:18.21; **V**:44.4
- Casting, of polymers, **IV**:3.11
- Cat lenses, **III**:19.9
- Cat mirror, **IV**:12.7, 12.8*f*
- Catadioptric (term), **I**:29.37
- Catadioptric Herschelian objective, **I**:29.27
- Catadioptric lenses, **V**:35.1
 afocal, **I**:18.21–18.22, 18.22*f*
 systems of, **I**:1.9
- Catadioptric objectives (*see* Reflective and catadioptric objectives)
- Cat-a-lac Black, **IV**:6.39, 6.42*f*, 6.53*f*
- Cataract, **III**:1.21, 14.8
 in aging eyes, **III**:14.24, 21.3–21.4
 defined, **III**:14.1, 19.1, 21.1, 23.1
 infrared, **III**:7.7, 7.10
 postsurgical correction of, **III**:12.14
 from radiation, **III**:7.4–7.6, 7.6*f*, 14.23
 treatment of, **III**:12.14
- Cataract lenses, **III**:12.15
- Cataract surgery, **III**:14.4, 14.24
 intraocular lenses in, **III**:21.4, 21.5 (*see also* Intraocular lenses)
 renormalization following, **III**:14.17
- Cathode ray tubes (CRTs), **I**:25.6–25.7, 30.4, 30.25–30.26; **III**:23.1
 color (*see* Color cathode ray tubes) and Computer Vision Syndrome, **III**:23.7, 23.8
- Cathode ray tubes (*Cont.*):
 monochrome, **III**:22.3*f*, 22.3–22.4
 controls for, **III**:22.7, 22.8*f*
 design and operation of, **III**:22.3*f*, 22.3–22.4
 standards for, **III**:22.14
 in optical systems, **III**:5.16
 screen reflections with, **III**:23.5
- Cathodes:
 photo-, **II**:27.6, 27.7*f*
 shielding of, **II**:27.10
 as x-ray tube sources, **V**:54.10–54.11
- Catoptric lens systems, **V**:35.1
- Catoptric systems, **I**:1.9
- Cauchy dispersion formula, **IV**:2.21
- Cauchy equation, **III**:1.19, 1.20; **V**:8.20, 8.21
- Causality, principle of, **IV**:2.8
- Caustics, ray densities and, **I**:1.88
- Cavity(-ies):
 distributed feedback lasers, **II**:16.29
 integrating (*see* Integrating cavities, of nonimaging optics)
 mode-locking, **II**:16.27–16.29, 16.28*f*
 modifying output distribution of, **II**:39.27
 properties of, **II**:16.3
 Q-switching, **II**:16.26–16.27, 16.27*f*
 ring lasers, **II**:16.29
 stability of, **II**:16.23–16.25, 16.24*f*, 16.25*f*
 surface absorption and, **IV**:6.15
 unstable resonators, **II**:16.25–16.26, 16.26*f*
- Cavity dumping, **II**:16.27
- Cavity losses, **II**:23.18
- Cavity resonance, for cw optical parametric oscillators, **IV**:17.2–17.4, 17.3*f*, 17.4*f*
- Cavity-shaped radiometers, **II**:34.28
- Ceilings, illuminated, **II**:40.13*f*
- Cellulose acetate butyrate, **IV**:3.4*t*
- Cellulose acetate film, **II**:29.4
- Cenco Company, **II**:15.47
- Center, of afocal lens, **I**:1.54
- Center for Optics Manufacturing, **II**:9.4
- Center of rotation (eye), **III**:1.42, 13.1
- Center thickness, for contact lenses, **III**:20.6
- Center-of-mass motion of atoms, **II**:23.45; **IV**:20.4
- Central visual processing, **III**:2.12–2.14, 2.13*f*
- Centrally obscured system (*see* Cassegrain design)
- Centration, of spherical lenses, **II**:9.8

- Ceragyrite (AgCl), **IV**:2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.69*t*, 2.76*t*
- Čerenkov effects, **V**:11.24
- Cesium chloride (CsCl), **IV**:2.68*t*
- Cesium iodide (CsI), **IV**:2.39*t*, 2.44*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.62*t*, 2.68*t*, 2.77*t*
- Cesium lithium borate ($\text{CsLiB}_6\text{O}_{10}$) (CLBO), **IV**:2.39*t*, 2.45*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.62*t*
- Chalcogenides:
in fiber lasers, **V**:25.27*t*, 25.29
as infrared fibers, **V**:12.2, 12.2*f*, 12.3*t*, 12.6, 12.7, 12.7*f*, 12.13
- Chalcopyrite, **IV**:2.74*t*
- Chandra X-Ray Observatory, **V**:33.2–33.4, 33.3*t*, 44.4, 44.10, 47.1, 47.4*f*, 47.5, 47.10, 64.7
- Channel electron multipliers (CEMs), **V**:60.6–60.7
- Channel power equalization, for EDFAs, **V**:21.41, 21.42*f*
- Channel stop region, **II**:32.7
- Channeled substrate planar (CSP) lasers, **II**:19.20*t*, 19.21*f*, 19.22, 19.36*f*
- Channels, for wave propagation, **I**:9.16
- Characteristic functions (geometrical optics), **I**:1.13–1.18
angle characteristic function, **I**:1.14–1.15, 1.15*f*, 1.17
angle-point characteristic function, **I**:1.16, 1.17
and expansions about rays, **I**:1.16
and expansions about the axis, **I**:1.16–1.17
ideal, **I**:1.17–1.18
mixed, **I**:1.13
paraxial forms of, **I**:1.17
and paraxial matrices, **I**:1.74
point characteristic function, **I**:1.14
point eikonal, **I**:1.14, 1.17
point-angle characteristic function, **I**:1.15–1.17
- Characteristic radiation:
Bremsstrahlung radiation as, **V**:56.8–56.10
from laser generated plasmas, **V**:56.2–56.10
recombination radiation as, **V**:56.10
spectral lines as, **V**:56.2–56.8
from x-ray tube sources, **V**:54.6–54.8, 54.7*f*
- Characterization:
of CRTs, **III**:22.20–22.34
absolute vs. for interaction, **III**:22.33
calibration vs., **III**:22.20
- Characterization, of CRTs (*Cont.*):
choice of method for, **III**:22.20–22.21
exhaustive, **III**:22.21–22.23
local, **III**:22.24–22.27
model-dependent, **III**:22.27–22.33
of head-mounted displays, **III**:25.7–25.10, 25.8*f*, 25.9*t*–25.10*t*
- Charge injection devices (CIDs), **I**:26.6*f*–26.8*f*, 26.6–26.7; **II**:31.1, 32.20, 33.10–33.11, 33.11*f*, 33.12*f*; **V**:60.7, 60.10*t*
- Charge integration matrix (CIM), **II**:33.10–33.11, 33.11*f*, 33.12*f*
- Charge pumping, **II**:32.9*n*
- Charge spectrographs, **I**:34.8, 34.8*f*
- Charge sweep devices (CSDs), **II**:33.12*f*, 33.13
- Charge-coupled detector area image sensor arrays:
frame transfer, **II**:32.26–32.28, 32.27*f*, 32.28*f*
interline transfer, **II**:32.28–32.32, 32.29*f*–32.31*f*
performance of, **II**:32.32
- Charge-coupled devices (CCDs), **I**:25.7, 26.3–26.5, 26.4*f*–26.6*f*; **II**:25.10, 25.11, 25.11*f*, 31.1, 32.12–32.20, 38.9–38.10, 38.10*t*; **IV**:5.61
in adaptive optics, **V**:5.21
characteristics of, **II**:32.17*f*, 32.17–32.20, 32.18*f*
electronics of, **II**:38.10
image sensing with, **II**:32.8
linear arrays of, **II**:32.21–32.24, 32.22*f*, 32.23*f*
MIS photogate FPAs for, **II**:33.10–33.11, 33.11*f*, 33.12*f*
multilinear arrays of, **II**:32.21, 32.23*f*, 32.24
operation of, **II**:32.12–32.14, 32.13*f*
output of, **II**:32.14–32.15, 32.15*f*
performance of, **II**:32.32
readout from, **II**:32.12–32.21, 32.13*f*
types of, **II**:32.15–32.17, 32.16*f*
as x-ray detectors, **V**:60.7, 60.8, 60.10*t*
x-ray imaging detectors in, **V**:61.7–61.8, 61.8*f*
- Charged area development (CAD), in xerographic systems, **I**:34.4
- Charge-resonance enhanced ionization (CREI), **IV**:21.26
- Chartered Institution of Building Services Engineers (CIBSE), **II**:40.2
- Chebyshev polynomials, **I**:7.15; **V**:46.6
- Chemglaze Z series (*see* Aeroglaze Z series)

- Chemical beam epitaxy (CBE), **I**:21.17, 21.18;
II:19.7
- Chemical vapor deposition (CVD), **IV**:6.56,
7.11; **V**:25.26
- Chemical-assisted ion beam etching (CAIBE),
II:19.39
- Chief rays, **I**:1.75, 17.8, 29.20, 29.37; **III**:4.3
- Chiolite, **I**:13.42
- Chip screening, of SOAs, **V**:19.17, 19.17f
- Chiral liquid crystals (LCs), **V**:8.8–8.13, 8.11f
- Chiral particles, scattering by, **I**:7.2
- Chirally coupled core (CCC) fibers, **V**:25.2,
25.19f, 25.21–25.22
- Chirp and chirping, **V**:20.1
of electro-absorption modulators, **V**:13.58
of fiber Bragg gratings, **V**:17.7
frequency, **V**:13.1, 13.17f, 13.17–13.18
of lasers, **V**:9.8
and optical fibers, **V**:10.3, 10.11
of solitons, **V**:22.3
- Chirped fiber Bragg gratings (FBGs),
V:21.22–21.23, 21.23f, 21.25f,
21.25–21.26
- Chirped multilayers, **IV**:7.47, 7.48
- Chirped pulse amplification (CPA), **II**:21.5,
21.5f, 21.6; **V**:25.2, 25.32, 25.33
- Chirped pulse amplification (CPA) lasers,
IV:21.4f, 21.4–21.5
- Chirped pulse excitation, **IV**:11.25–11.26
- Chirps, in spectroscopy, **IV**:11.3
- Chlorophyll:
and absorption by phytoplankton, **IV**:1.23,
1.24, 1.28
and beam attenuation, **IV**:1.41
and diffuse attenuation, **IV**:1.44
fluorescence by, **IV**:1.49
and remote sensing, **IV**:1.46
- Chloroplasts, **IV**:1.23
- Cholesky decomposition, **I**:14.41, 14.42
- Cholesteric liquid crystal display (Ch-LCD),
V:8.32
- Chopper-stabilized amplifiers, **II**:27.11
- Chopper-stabilized BDI, **II**:33.19t, 33.20f,
33.21–33.22
- Choroid, **III**:18.1
- Chromatic aberration, **I**:1.91–1.92; **II**:2.2–2.4,
2.3f, 2.4f; **III**:1.19f, 1.19–1.20
and accommodation, **III**:1.34
age-related, **III**:14.14
correcting color coordinates for, **III**:10.17
- Chromatic aberration (*Cont.*):
correction of, **I**:23.5–23.6, 23.6f;
III:1.25–1.26, 1.26f
longitudinal, **III**:15.22
and macular pigment, **III**:1.9
transverse, **III**:15.22–15.23
with visual instruments, **III**:1.28
- Chromatic aberration control, **IV**:3.8
- Chromatic adaptation, **III**:11.35f
achromatic detection and, **III**:11.47–11.49,
11.48f
color appearance and, **III**:11.68
and luminous efficiency, **III**:11.37
second-site, **III**:11.17–11.18
and the Sloan notch, **III**:11.49, 11.51
- Chromatic contrast detection, **III**:2.29–2.31,
2.30f
- Chromatic CSFs, **III**:2.29–2.31
- Chromatic detection:
and color appearance, **III**:11.69
on neutral fields, **III**:11.34
- Chromatic discrimination, **III**:11.57–11.62
and chromatic adaptation, **III**:11.69
and color appearance, **III**:11.69
defined, **III**:11.1
and gap effect, **III**:11.72, 11.74
near detection threshold, **III**:11.58–11.59
pedestal experiments, **III**:11.59, 11.60f,
11.61–11.62
- Chromatic dispersion:
and fiber optic communication links, **V**:15.9,
15.10
in WDM networks, **V**:21.14–21.16, 21.15f,
21.16f
- Chromatic mechanisms, **III**:11.1, 11.80–11.81
- Chromatic signals, multiplexing of achromatic
signals and, **III**:11.76–11.79, 11.78f
- Chromatic spatial CSFs, **III**:2.29
- Chromatic temporal CSFs, **III**:2.31
- Chromatic valence data, **III**:11.27, 11.28f
- Chromaticity coordinates, **III**:10.1, 10.4t,
10.20–10.21, 10.22f
- Chromaticity diagrams, **III**:10.20–10.21, 10.22f
- Chromatism, of axial gradients, **I**:24.3–24.6
- Chromium:
absorptance of, **IV**:4.48t, 4.50t
optical properties of, **IV**:4.13t–4.14t, 4.22f
physical properties of, **IV**:4.54t
reflectance in, **IV**:4.30t–4.31t
thermal properties of, **IV**:4.69t

- Chromium ions (Cr^{3+}):
 absorption and photoluminescence of,
 V:2.19*f*, 2.19–2.21, 2.20*f*
 optical spectroscopy of, V:2.9–2.10, 2.10*f*
 polarization spectroscopy of, V:2.21–2.22,
 2.23*f*
- Chromium lasers, II:16.34, 16.35
- Chromogenic film, II:29.14
- Chromophores, III:21.1, 21.20, 21.20*f*
- Chromostereopsis, III:1.20
- CIE 1931 2° color-matching functions,
 III:10.6, 10.12
- CIE 1964 10° color-matching functions,
 III:10.12–10.13, 10.16
- CIE luminous efficiency functions,
 III:10.44–10.45, 11.37, 11.38*f*
- CIELAB, III:10.42–10.43
- Cilia, III:8.1
- Ciliary body, III:1.3*f*
- Ciliary muscle, III:21.2
- Ciliary photoreceptors, III:8.3
- Ciliary ring, III:1.30–1.31
- Circle polynomials, II:11.36*t*, 11.39
 isometric plots/interferograms/PSFs for
 defocus, II:11.38*f*
 and noncircular pupils, II:11.37, 11.39
 radial, II:11.7, 11.9*f*–11.10*f*
 Zernike, I:1.90, 23.3; II:11.4, 11.6–11.12,
 11.8*t*–11.9*t*, 11.9*f*–11.11*f*, 11.12*t*
- Circuits:
 in integrated optics, I:21.21–21.31
 for active devices, I:21.25–21.31,
 21.26*f*–21.31*f*
 for passive devices, I:21.21–21.25,
 21.22*f*–21.25*f*
 optical tank, V:20.21, 20.21*f*
- Circuit-switched networks, V:21.7–21.10,
 21.8*f*–21.9*f*
- Circular analyzers, of polarized light, I:15.18,
 15.19
- Circular apertures:
 diffraction of light from, I:3.6*f*, 3.6–3.7, 3.7*f*,
 3.9–3.11
 Fraunhofer patterns for, I:3.25, 3.26, 3.27*f*
- Circular discs, projected area of, II:36.3*t*
- Circular gratings, V:40.1
- Circular polarization:
 analyzers for, V:43.6–43.8, 43.7*f*
 phase plates for, V:43.5*f*
 and synchrotron radiation, V:55.6–55.7
- Circular polarization (OCT), III:18.21
- Circular polarizers, I:15.17–15.19
- Circular scan, I:30.16, 30.18*f*
- Circulators:
 and fiber Bragg gratings, V:17.8*f*, 17.9
 for networking, V:18.3, 18.3*f*, 18.10
- Circumferential coordinates, V:45.7
- Circumferential slope errors, V:45.8
- Cladding:
 defined, V:25.2
 photonic crystal fibers in, V:11.7–11.11,
 11.8*f*–11.10*f*
- Cladding layers, II:19.4
- Clamshell housings, IV:3.15, 3.15*f*
- Clandestine cameras, I:25.21
- Clarity, perception of visual, II:40.5
- Classical electronic polarization theory,
 IV:2.14
- Classical harmonic oscillator model,
 IV:10.5–10.7, 10.6*f*
- Clausius-Mossotti equation, IV:8.7
- Clausius-Mossotti relationships, IV:2.24
- Clausius-Mossotti theory, I:7.16
- Cleaning, in xerographic systems, I:34.10
- Cleaning and cleanliness, of optical surfaces,
 IV:19.3–19.5
- Clear lens extraction (CLE), III:21.18
- Cleaving, of photonic crystal fibers, V:11.26
- Clebsch-Gordon coefficients, I:31.17
- Climate change, global, V:3.43–3.45, 3.44*f*
- Clip test (of photographic film), II:30.23
- Clipped Lambertian distribution,
 II:39.3–39.4
- Clock generation, II:33.16
- Clock recovery, in OTDM networks, V:20.21*f*,
 20.21–20.22
- Clocks, atomic, IV:20.28
- Closed family, IV:20.38
- Closed HMDs, III:25.3, 25.4
- Closed-loop performance (in servo systems),
 II:22.8
- Closed-loop stability issues (in servo systems),
 II:22.8–22.12, 22.9*f*
 PID controller vs. notch filters, II:22.10*f*,
 22.10–22.11, 22.11*f*
 rule-of-thumb PID design for system with
 transducer resonance, II:22.11–22.12
- Cluster electron heating, IV:21.34, 21.35
- Cluster expansion, IV:21.35, 21.35*f*

- Clusters, strong field interactions with, **IV**:21.31–21.36
- Coulomb explosion, **IV**:21.33–21.34
- intense laser pulse interactions, **IV**:21.35–21.36, 21.36*f*
- ionization mechanisms in, **IV**:21.31–21.33, 21.32*f*
- nanoplasma description, **IV**:21.34–21.35, 21.35*f*
- Coarse wavelength division multiplexing (CWDM) systems, **V**:19.27
- Coarse-grained derivative, **II**:23.21
- Coated spheres, scattering by, **I**:7.14
- Coatings:
- antireflection, **IV**:7.15–7.32; **V**:19.8, 19.8*f*, 19.20
 - of absorbing and amplifying media, **IV**:7.26, 7.27
 - homogeneous-layer, **IV**:7.16–7.23, 7.17*f*–7.19*f*, 7.20*t*–7.21*t*, 7.22*f*–7.23*f*
 - inhomogeneous and structured, **IV**:7.23–7.26, 7.24*f*, 7.26*f*
 - at nonnormal angle of incidence, **IV**:7.28*f*–7.31*f*, 7.28–7.31
 - nonoptical properties of, **IV**:7.31–7.32, 7.32*f*
 - surface reflections and optical performance, **IV**:7.15–7.16, 7.16*f*
 - of surfaces carrying thin films, **IV**:7.27–7.28, 7.28*f*
 - universal, **IV**:7.26, 7.27*f*
 - filters with metallic reflecting, **IV**:7.80–7.81
 - high performance optical multilayer, **IV**:7.96–7.98
 - for infrared optical fibers, **V**:12.4, 12.7, 12.9
 - and interference polarizers, **IV**:7.70*f*–7.72*f*, 7.70–7.72
 - laser-induced damage in, **IV**:19.3–19.4
 - lens specifications for, **II**:4.10
 - measurements on, **IV**:7.12–7.14
 - narrowband reflection, **IV**:7.43, 7.44*f*
 - phase, **IV**:7.101, 7.101*f*–7.104*f*, 7.102
 - of photographic film, **II**:29.4
 - reflection, **IV**:7.106*f*–7.113*f*, 7.106–7.113
 - reflective (multilayers), **V**:41.1–41.10
 - and calculation of multilayer properties, **V**:41.3–41.4
 - for diffractive imaging, **V**:41.9–41.10, 41.10*f*
 - fabrication methods and performance of, **V**:41.4–41.9, 41.5*f*, 41.6*t*, 41.7*t*, 41.8*f*
 - properties of, **V**:41.1–41.3, 41.2*f*
- Coatings (*Cont.*):
- thin-film
 - and antireflection coatings, **IV**:7.27–7.28, 7.28*f*
 - manufacturing of, **IV**:7.10–7.12
 - of metal, **IV**:7.104, 7.104*f*
 - theory and design of, **IV**:7.5–7.10, 7.6*f*, 7.9*f*
 - transmission and reflection of, **IV**:7.3
 - types of, **IV**:3.17–3.18, 6.13*t*
 - for ultrafast optics, **IV**:7.47–7.48, 7.48*f*
 - (*see also specific coatings, e.g.*: Ebanol C coating)
 - Coblentz sphere, **V**:1.10, 1.10*f*, 1.11
 - Coblentz-type thermopiles, **II**:24.23
 - Cochlea, **III**:8.1, 8.26
 - Cochlear hair cells, light guide effect in, **III**:8.24–8.26
 - Coddington's equations, **I**:1.44
 - Code mark inversion (CMI), **V**:20.9
 - CODE V (optical software), **V**:35.1
 - Coding, in OTDM networks, **V**:20.9*f*, 20.9–20.10
 - Coefficient of finesse (interference), **I**:2.31
 - Coefficient of linear thermal expansion, **IV**:4.7, 4.56*t*, 4.57*f*–4.58*f*
 - Coefficient of thermal expansion (CTE), of metals, **IV**:4.6–4.7, 4.10*t*, 4.53, 4.55*t*
 - Coercivity, of optical disk data, **I**:35.17*n*, 35.27, 35.27*f*, 35.28
 - Cogging, of streak cameras, **I**:25.24
 - Coherence, **I**:5.1–5.23, 6.2–6.13
 - analytical signal representation, **I**:5.2–5.3
 - applications of, **I**:5.22–5.23
 - in atomic systems, **IV**:14.4*f*, 14.4–14.5, 14.28–14.32, 14.29*f*–14.32*f*
 - beam spatial, **V**:58.4
 - in binary optics, **I**:23.7, 23.8*f*
 - classical, **I**:5.1–5.2
 - coherence area, **I**:5.3
 - and coherence functions, **I**:5.4–5.9
 - angular correction function, **I**:5.6, 5.7*f*
 - complex degree of coherence, **I**:5.4
 - complex degree of spectral coherence, **I**:5.5
 - cross-spectral density function, **I**:5.5
 - efficient sampling of, **I**:6.10–6.12, 6.11*f*
 - higher-order functions, **I**:5.8–5.9
 - intensity, **I**:5.7
 - mutual coherence function, **I**:5.4
 - radiance, **I**:5.8
 - radiant emittance, **I**:5.7–5.8

- Coherence, and coherence functions (*Cont.*):
 radiant intensity, **I**:5.8
 spectrum and normalized spectrum,
I:5.5–5.6
 coherence time, **I**:5.3
 in coherent x-ray optics, **V**:27.5
 complex degree of, **I**:2.37, 5.4
 and enhanced backscatter, **I**:6.5f, 6.5–6.7
 and general linear systems, **I**:6.3–6.4
 and image formation, **I**:6.9f, 6.9–6.10
 and interference, **I**:2.13, 2.36–2.42
 laser sources, **I**:2.41–2.42, 2.42f
 Michelson stellar interferometers, **I**:2.40f,
 2.40–2.41
 mutual coherence function, **I**:2.36f,
 2.36–2.38
 spatial coherence, **I**:2.38f–2.39f, 2.38–2.40
 temporal coherence, **I**:2.41
 and Koehler-illumination, **I**:6.12f, 6.12–6.13
 and laser modes, **I**:5.23
 and Lau effect, **I**:6.7–6.8, 6.8f
 of light sources, **I**:5.9–5.13
 in long wavelength limit, **V**:55.17–55.20,
 55.18f, 55.19f
 and Lukosz-type super-resolving systems,
I:6.9f, 6.9–6.10
 maximal, **IV**:14.3, 14.28–14.32, 14.29f–14.32f
 measurements of coherence, **I**:5.3–5.4
 mutual coherence function (MCF), **V**:4.4,
 4.7, 4.10
 and noncosmological red shift, **I**:5.23
 and optical image enhancement, **I**:11.14–11.17
 partial, **I**:2.38
 and polarization effects, **I**:5.22
 propagation in, **I**:5.13–5.19, 5.14f–5.16f, 6.4
 and radio astronomy, **I**:5.23
 scalar field amplitude, **I**:5.3
 spatial coherence, **I**:5.3
 and speckle, **I**:5.22
 and spectral representation, **I**:5.22
 and spectrum of light, **I**:5.19–5.22, 5.20f
 and statistical radiometry, **I**:5.22
 temporal, **I**:2.41, 5.3; **V**:55.17–55.18, 55.18f
 time averages in, **I**:6.4–6.5
 transverse, **V**:55.18–55.20, 55.19f
 transverse spatial, **V**:55.16
- Coherence area, **I**:5.3
 Coherence collapse, from laser diodes,
V:13.22–13.23
 Coherence diameter, Fried's, **V**:5.2, 5.7, 5.9–5.10
 Coherence length, **I**:2.19; **III**:18.1; **IV**:10.15
 atmospheric, **V**:4.7–4.10
 Fried's, **V**:4.8
 and imaging through turbulence, **V**:4.8f,
 4.8–4.10, 4.9f
 spatial, **V**:27.2
 Coherence time, **I**:5.3
 Coherence volume, **I**:5.3
 Coherency matrix, **I**:12.29–12.30, 14.41
 Coherent anti-Stokes four-wave mixing,
IV:15.2t, 15.3t, 15.4, 15.4f
 Coherent anti-Stokes Raman scattering
 (CARS), **IV**:15.4, 15.4f, 15.34, 15.42,
 15.42t, 15.43f, 16.3t, 16.4, 16.5f, 16.17f,
 16.17–16.18
 Coherent area, of atmosphere, **V**:4.16
 Coherent arrays, scattering by, **I**:7.2–7.3
 Coherent backscattering, **I**:9.14f, 9.14–9.15
 Coherent control, **IV**:18.21
 Coherent coupling, **III**:1.27–1.28
 Coherent diffraction microscopy, **V**:27.4f,
 27.4–27.5, 27.5f
 Coherent Doppler LIDAR, **V**:3.38, 3.39f
 Coherent excitons, **IV**:18.19–18.20
 Coherent illumination, pupil size and,
III:6.9–6.12, 6.11t
 Coherent image amplification, **IV**:12.29, 12.30
 Coherent mode representation (spectrum of
 light), **I**:5.20–5.21
 Coherent optical image enhancement,
I:11.14–11.17
 Coherent optical transients, **IV**:11.1–11.28
 chirped pulse excitation, **IV**:11.25–11.26
 and cw spectroscopy, **IV**:11.2
 experimental considerations, **IV**:11.26–11.28,
 11.27f
 free polarization decay, **IV**:11.7–11.11, 11.8f,
 11.10f, 11.11f
 Maxwell-Bloch equations, **IV**:11.6–11.7
 optical Bloch equations, **IV**:11.3–11.6
 phase conjugate geometry and optical
 Ramsey fringes, **IV**:11.19–11.22,
 11.20f, 11.21f
 photon echo, **IV**:11.11–11.15, 11.12f, 11.13f,
 11.15f
 stimulated photon echo, **IV**:11.15–11.19,
 11.16f–11.19f
 two-photon transitions and atom
 interferometry, **IV**:11.22–11.23, 11.24f

- Coherent population return (CPR), **IV**:14.1, 14.30, 14.30*f*, 14.31*f*
- Coherent population trapping (CPT), **IV**:14.1, 14.3–14.5, 14.7, 20.37
- Coherent radiation, **I**:30.2, 30.25–30.26
- Coherent Raman scattering, **IV**:15.3
- Coherent scattering, **I**:7.3, 9.2, 9.3, 9.5–9.7, 9.6*f*; **V**:26.7, 63.7, 63.8
- Coherent states, **II**:23.12
- Coherent Stokes Raman scattering (CSRS), **IV**:16.3*t*, 16.17*f*, 16.17–16.18
- Coherent x-ray optics, **V**:27.1–27.5, 27.3*f*–37.5*f*
- Coiling:
of light bulb filament, **II**:40.30
of LMA fibers, **V**:25.18
- Cold cathode fluorescent lamps (CCFLs), **II**:40.32; **V**:8.30, 8.30*f*
- Cold mirrors, **IV**:7.58
- Cold-cathode field emission, **V**:54.10–54.11
- Coleoptile, **III**:8.1, 8.26
- Collagen, **III**:16.1, 16.4
- Collected volatile condensable materials (CVCM), **IV**:6.17
- Collective tunneling, **IV**:21.18
- Collector power (in stray light suppression), **II**:7.2
- Collectors (*see* Concentrators, nonimaging)
- Collet-and-cap housings, **IV**:3.15, 3.15*f*
- Colliding pulse modelocking (CPM), **IV**:18.3
- Collimating polycapillary optics, **V**:53.14, 53.14*f*
- Collimating single crystal diffraction, **V**:53.12, 53.12*f*
- Collimation and collimators, **I**:35.8, 35.9*f*, 35.10*f*
anisotropic acoustic beam, **V**:6.25
atomic beams, **IV**:20.15*f*, 20.15–20.16
autocollimators, **II**:12.12
conic collimators, **II**:39.8, 39.9*f*
neutron, **V**:63.15–63.16
in neutron and x-ray optics, **V**:26.11, 26.11*f*
in polycapillary x-ray optics, **V**:53.8*f*–53.9*f*, 53.8–53.9
with refractive x-ray lenses, **V**:37.8
Soller, **V**:28.2, 28.2*f*, 28.3
in SPECT imaging, **V**:32.3
- Collinear beam acousto-optic tunable filters (CBAOTFs), **V**:6.43, 6.43*f*, 6.45*t*
- Collineation, **I**:1.56–1.63
of conjugate lines, **I**:1.59
of conjugate planes, **I**:1.58–1.59
- Collineation (*Cont.*):
coordinate systems and degrees of freedom for, **I**:1.57
equations of, **I**:1.57–1.58
general properties of, **I**:1.62–1.63
matrix representation of, **I**:1.59–1.60
of rotationally symmetric lenses, **I**:1.60–1.62, 1.62*f*
- Collision length, for solitons, **V**:22.9–22.10, 22.15
- Collisional broadening, emission-line, **II**:16.5
- Collisional broadening, spectral-line, **IV**:14.13
- Collisional heating, **IV**:21.37, 21.37*f*
- Collisional ionization, **IV**:21.31, 21.32
- Collisions, **IV**:20.28–20.31, 20.30*f*, 20.31*f*
excited-state, **IV**:20.29
ground-state, **IV**:20.29
trap loss, **IV**:20.29
- Colloidal silver, **II**:29.13
- Colloids, in water, **IV**:1.14
- Colocalization, single molecule high-resolution, **I**:28.23
- Color(s):
anomalous reflection, **II**:30.17
axial, **I**:1.91, 29.9, 29.37
in electronic imaging, **III**:24.5–24.6, 24.8–24.9
in human visual system, **I**:26.18
lateral, **I**:1.91–1.92
in LEDs, **II**:40.37
and lighting design, **II**:40.7–40.9
mixing of, **II**:40.8
science of, **II**:30.15–30.18, 30.16*f*, 30.17*f*
and visual acuity, **III**:4.10
in xerographic systems, **I**:34.11*f*–34.13*f*, 34.11–34.12
- Color, ocean, **IV**:1.46
- Color aliasing, **II**:32.34
- Color appearance:
and chromatic adaptation, **III**:11.68
and chromatic detection and discrimination, **III**:11.69
and color constancy, **III**:11.71
and color opponency, **III**:11.62–11.66
color-opponent response or valence functions, **III**:11.63, 11.65*f*
hue scaling, **III**:11.63, 11.64*f*
opponent-colors theory, **III**:11.62–11.63

- Color appearance, and color opponency (*Cont.*):
 spectral properties of color-opponent mechanisms, **III**:11.63
 unique hues and equilibrium colors, **III**:11.63–11.66
 defined, **III**:11.2
 and habituation, **III**:11.69–11.70
 phenomenological aspects of, **III**:11.5
 and postreceptoral mechanisms, **III**:11.3, 11.4*f*, 11.5
 and stabilized borders, **III**:11.74–11.75
 trichromacy vs., **III**:11.5
- Color assimilation, **III**:11.2, 11.4*f*
- Color balance and tracking, in color CRTs, **III**:22.19
- Color cameras, **III**:10.38–10.40
- Color cathode ray tubes (color CRTs), **III**:22.1–22.34
 colorimetric calibration/characterization, **III**:22.20–22.34
 absolute vs. characterization for interaction, **III**:22.33
 choice of method for, **III**:22.20–22.21
 exhaustive methods, **III**:22.21–22.23
 local methods, **III**:22.24–22.27
 model-dependent methods, **III**:22.27–22.33
- design and operation of, **III**:22.3–22.13
 electronics and controls, **III**:22.6–22.13, 22.7*f*–22.13*f*
 and monochrome CRTs, **III**:22.3*f*, 22.3–22.4
 shadowmask color CRTs, **III**:22.4–22.6, 22.5*f*, 22.6*f*
- operational characteristics of, **III**:22.13–22.18
 colorimetric standards, **III**:22.14
 spatial characteristics of emitted light, **III**:22.14–22.15
 spatial uniformity, **III**:22.17*f*, 22.17–22.18, 22.18*f*
 stability of output, **III**:22.16*f*, 22.16–22.17, 22.17*f*
 temporal characteristics of emitted light, **III**:22.14–22.16
 timing and synchronization standards, **III**:22.13–22.14
- setup for image display, **III**:22.18–22.19
 brightness, **III**:22.19
 color balance and tracking, **III**:22.19
 contrast, **III**:22.19
 focus, **III**:22.19
 viewing environments, **III**:22.19–22.20
- Color constancy, **III**:10.39–10.40, 11.2
- Color constancy mechanisms, **III**:11.71–11.72
- Color contrast, **III**:11.2, 11.4*f*
- Color convergence, **III**:23.1
- Color coordinate systems, **III**:10.11–10.24
 adjusting cone spectral sensitivities, **III**:10.17–10.18
 colorimetric measurements, **III**:10.23–10.24
 color-matching functions, **III**:10.11–10.13
 cone fundamentals, **III**:10.13–10.14
 limits of color-matching data, **III**:10.15–10.17
 opponent and contrast spaces, **III**:10.18–10.19
 stimulus spaces, **III**:10.11
 visualizing color data, **III**:10.19*f*, 10.19–10.23, 10.22*f*, 10.23*f*
- Color coordinates:
 in color-deficient observers, **III**:10.16
 correcting for chromatic aberrations, **III**:10.17
 of different visual systems, **III**:10.38–10.39
 of surfaces, **III**:10.36
 transformation to CIELAB, **III**:10.42–10.43
- Color data representations, **III**:11.31–11.33, 11.32*f*
- Color density, **II**:29.7–29.8
- Color direction, **III**:11.12
- Color discrimination, **III**:10.40–10.43, 10.41*f*, 14.15, 14.17
- Color filter arrays (CFAs), **I**:26.18; **III**:32.32–32.34, 32.33*f*, 32.34*f*
- Color imaging architectures, **II**:32.32–32.34
 integral filter arrays, **II**:32.32–32.34, 32.33*f*, 32.34*f*
 sequential, **II**:32.32, 32.33*f*
 three-chip, **II**:32.32, 32.33*f*
- Color LCDs, **III**:22.37–22.40
 colorimetry of color pixels, **III**:22.38–22.39
 controls and input standards, **III**:22.39
 geometry of color pixels, **III**:22.37*f*, 22.37–22.38
 spatial variations in output, **III**:22.40
 temporal variations in output, **III**:22.39–22.40
- Color matching, **III**:10.6–10.10
 in color-deficient observers, **III**:10.16
 consistency across observers, **III**:10.9
 critical properties of, **III**:10.8–10.10
 errors in, **III**:10.43, 10.43*f*, 10.44
 Grassmann's laws, **III**:10.8–10.9

- Color matching (*Cont.*):
 maximum saturation method, **III**:10.6*f*,
 10.6–10.7
 Maxwell's method, **III**:10.8, 10.8*f*
 persistence of, **III**:10.9
 trichromatic, **III**:10.7–10.8
 tristimulus values for arbitrary lights,
III:10.9
 uniqueness in, **III**:10.9
- Color negative films, **II**:30.25–30.28, 30.27*t*
- Color opponency:
 and color appearance, **III**:11.62–11.66
 and “forbidden” colors, **III**:11.74
 implied by test measurements, **III**:11.12–11.15,
 11.13*f*, 11.26–11.27, 11.28*f*
 and multiplexing of color and luminance
 signals, **III**:11.76–11.79
 third-level, zero crossings of, **III**:11.83, 11.84
- Color photographic films:
 about, **II**:30.2
 coating of, **II**:29.4
 negative, **II**:30.25–30.28, 30.27*t*
 reversal, **II**:30.22–30.24, 30.23*t*
 structure of, **II**:29.12–29.15, 29.13*f*, 29.14*f*,
 30.3*f*, 30.3–30.5
- Color photographic paper, **II**:30.5
- Color records, **II**:30.4
- Color rendering, **II**:40.8–40.9
- Color rendering index (CRI), **II**:40.8
- Color reversal films, **II**:30.2, 30.22–30.24,
 30.23*t*
- Color sequential systems, **II**:32.32, 32.33*f*
- Color slide films (*see* Color reversal films)
- Color space calculations, **II**:38.4–38.5
- Color space transformation matrix,
III:10.1, 10.31
- Color spaces, **III**:10.11, 10.18–10.19,
 11.31–11.33
 contrast, **III**:10.19
 defined, **III**:11.2
 different directions of, **III**:11.39, 11.41–11.46,
 11.44*f*, 11.46*f*
 detection contours in L,M plane,
III:11.40*f*, 11.41–11.42
 detection in planes other than L,M,
III:11.43–11.45, 11.44*f*
 mechanism interactions, **III**:11.42–11.43
 spatial and temporal CSFs, **III**:11.45–11.46,
 11.46*f*
 DKL, **III**:10.19, 11.32, 11.32*f*, 11.33
- Color spaces (*Cont.*):
 opponent, **III**:10.19
 specifying, **III**:10.11
 stimulus, **III**:10.11
 transformations between, **III**:10.24*t*,
 10.29–10.32
 uniform, **III**:10.40, 10.42
- Color temperature, **II**:34.44, 37.4*t*, 37.6–37.7,
 38.5, 40.8; **III**:23.2
- Color transparency films (*see* Color reversal
 films)
- Color valence, **III**:11.2, 11.63, 11.65*f*, 11.66
- Color vision, age-related changes in,
III:14.15, 14.17
- Color vision mechanisms, **III**:11.1–11.85
 basic model details and limits, **III**:11.31
 chromatic discrimination, **III**:11.57–11.62
 near detection threshold, **III**:11.58–11.59
 pedestal experiments, **III**:11.59, 11.60*f*,
 11.61–11.62
 color and contours, **III**:11.72–11.79,
 11.73*f*–11.75*f*
 color appearance and stabilized borders,
III:11.74–11.75
 contours and after-effects, **III**:11.75*f*,
 11.75–11.76
 gap effect and luminance pedestals,
III:11.72, 11.74
 McCollough effect, **III**:11.76, 11.77*f*
 multiplexing chromatic and achromatic
 signals, **III**:11.76–11.79, 11.78*f*
 color appearance and color opponency,
III:11.62–11.66
 color-opponent response (valence
 functions), **III**:11.63, 11.65*f*
 hue scaling, **III**:11.63, 11.64*f*
 opponent-colors theory, **III**:11.62–11.63
 spectral properties of color-opponent
 mechanisms, **III**:11.63
 unique hues and equilibrium colors,
III:11.63–11.66
 color constancy, **III**:11.71–11.72
 color data representations, **III**:11.31–11.33,
 11.32*f*
 color-appearance mechanisms,
III:11.26–11.31
 color-discrimination mechanisms vs.,
III:11.5–11.8, 11.6*f*, 11.7*f*, 11.81–11.82
 field measurements and first-site
 adaptation, **III**:11.27, 11.29, 11.30*f*

- Color vision mechanisms, color-appearance mechanisms (*Cont.*):
 field measurements and second-site adaptation, **III**:11.29, 11.31
 test measurements and opponency, **III**:11.26–11.27, 11.28*f*
- color-discrimination mechanisms, **III**:11.9–11.26
 color-appearance mechanisms vs., **III**:11.5–11.8, 11.6*f*, 11.7*f*, 11.81–11.82
 field method, **III**:11.11–11.12
 first-site adaptation, **III**:11.15–11.17, 11.16*f*
 opponency implied by test measurements, **III**:11.12–11.15, 11.13*f*
 psychophysical test method, **III**:11.9, 11.11, 11.12
 second-site adaptation, **III**:11.17–11.22, 11.18*f*, 11.19*f*
 sites of limiting noise, **III**:11.20, 11.23*f*, 11.23–11.26, 11.25*f*
- field sensitivities, **III**:11.46–11.57
 achromatic detection and chromatic adaptation, **III**:11.47–11.49, 11.48*f*
 chromatic adaptation and the Sloan notch, **III**:11.49, 11.51
 detection contours and field adaptation, **III**:11.53–11.54, 11.54*f*
 field additivity, **III**:11.51, 11.52*f*
 first- and second-site adaptation, **III**:11.51, 11.52, 11.53*f*
 habituation or contrast adaptation experiments, **III**:11.54–11.56, 11.55*f*
 multiple cone inputs, **III**:11.49, 11.50*f*
 noise-masking experiments, **III**:11.56–11.57
 Stiles' π -mechanisms, **III**:11.46, 11.47, 11.47*f*
- guiding principles of, **III**:11.8–11.9
 linearity of color-opponent mechanisms, **III**:11.66–11.70
 Bezold-Brücke effect and invariant hues, **III**:11.67*f*, 11.67–11.68
 color appearance and chromatic adaptation, **III**:11.68
 color appearance and chromatic detection/discrimination, **III**:11.69
 color appearance and habituation, **III**:11.69–11.70
 luminance and brightness, **III**:11.70, 11.70*f*
 tests of linearity, **III**:11.66–11.67
- Color vision mechanisms (*Cont.*):
 low-level and higher-order mechanisms, **III**:11.79–11.80
 and mechanism concept, **III**:11.9–11.11, 11.10*f*
 nomenclature for, **III**:11.8
 test sensitivities, **III**:11.34–11.46
 to different directions of color space, **III**:11.39, 11.41–11.43, 11.44*f*, 11.45–11.46, 11.46*f*
 luminance, **III**:11.37–11.39, 11.38*f*, 11.40*f*
 to spectral lights, **III**:11.34, 11.35*f*, 11.36*f*, 11.37
 three-stage zone models, **III**:11.82–11.85, 11.85*f*
 unipolar vs. bipolar chromatic mechanisms, **III**:11.80–11.81
- Color-appearance mechanisms, **III**:11.26–11.31
 color-discrimination mechanisms vs., **III**:11.5–11.8, 11.6*f*, 11.7*f*, 11.81–11.82
 defined, **III**:11.2
 field measurements
 and first-site adaptation, **III**:11.27, 11.29, 11.30*f*
 and second-site adaptation, **III**:11.29, 11.31
 test measurements and opponency, **III**:11.26–11.27, 11.28*f*
- Color-center lasers, **II**:16.35; **IV**:18.10
- Color-discrimination mechanisms, **III**:11.9–11.26, 11.10*f*
 color-appearance mechanisms vs., **III**:11.5–11.8, 11.6*f*, 11.7*f*, 11.81–11.82
 defined, **III**:11.2
 field method, **III**:11.11–11.12
 first-site adaptation, **III**:11.15–11.17, 11.16*f*
 low-level and higher-order, **III**:11.79–11.80
 opponency implied by test measurements, **III**:11.12–11.15, 11.13*f*
 psychophysical test method, **III**:11.9, 11.11, 11.12
 second-site adaptation, **III**:11.17–11.22, 11.18*f*, 11.19*f*
 sites of limiting noise, **III**:11.20, 11.23*f*, 11.23–11.26, 11.25*f*
- Colorimeters, **III**:10.23
 Colorimetry, **II**:37.11; **III**:10.1–10.45
 brightness matching and photometry, **III**:10.43–10.45
 color cameras, **III**:10.38–10.40

- Colorimetry (*Cont.*):
- color coordinate systems, **III**:10.11–10.24
 - adjusting cone spectral sensitivities, **III**:10.17–10.18
 - colorimetric measurements, **III**:10.23–10.24
 - color-matching functions, **III**:10.11–10.13
 - cone fundamentals, **III**:10.13–10.14
 - limits of color-matching data, **III**:10.15–10.17
 - opponent and contrast spaces, **III**:10.18–10.19
 - stimulus spaces, **III**:10.11
 - visualizing color data, **III**:10.19*f*, 10.19–10.23, 10.22*f*, 10.23*f*
 - color discrimination, **III**:10.40–10.43, 10.41*f*
 - color matching, **III**:10.6–10.10
 - conventional terms/notation, **III**:10.4, 10.4*t*, 10.5
 - errors in color-matching functions, **III**:10.43, 10.43*f*, 10.44
 - image processing chain, **III**:10.2*f*
 - matrix representations/calculations, **III**:10.24–10.32
 - stimulus representation, **III**:10.24–10.27, 10.25*f*, 10.26*f*
 - transformations between color spaces, **III**:10.24*t*, 10.29–10.32
 - vector representation of data, **III**:10.25*f*, 10.27*f*, 10.27–10.29
 - metamerism, **III**:10.36–10.38, 10.37*f*
 - scope of, **III**:10.3
 - standards for color CRTs, **III**:22.14
 - surfaces and illuminants, **III**:10.32–10.36, 10.33*f*, 10.34*f*
 - trichromacy, **III**:10.4–10.6
 - univariance, **III**:10.4
 - visual systems, **III**:10.38–10.40
- Color-matching functions (CMFs), **III**:10.11–10.13
- defined, **III**:10.1, 10.4*t*
 - limits of data, **III**:10.15–10.17
 - and luminosity function, **III**:10.10
 - and maximum saturation method, **III**:10.7
 - online tabulation of, **III**:10.11
 - specificity of, **III**:10.9
 - standards for, **III**:10.12–10.13
 - tailored to individuals, **III**:10.15
 - transformation of, **III**:10.10, 10.10*f*
- Color-opponent mechanisms:
- linearity of, **III**:11.66–11.70
 - Bezold-Brücke effect and invariant hues, **III**:11.67*f*, 11.67–11.68
 - color appearance and chromatic adaptation, **III**:11.68
 - color appearance and chromatic detection/discrimination, **III**:11.69
 - color appearance and habituation, **III**:11.69–11.70
 - luminance and brightness, **III**:11.70, 11.70*f*
 - tests of linearity, **III**:11.66–11.67
 - and opponent-colors theory, **III**:11.3
 - spectral properties of, **III**:11.63
- Color-selective beam splitters, **IV**:7.65–7.66, 7.66*f*, 7.67*f*
- Coma, **I**:1.31, 24.7, 29.37
- of grazing incidence telescopes, **V**:44.9, 44.10
 - with spherical aberration, **II**:2.4, 2.4*f*
- Combined recombination, **II**:17.3
- Combined servo transducers, **II**:22.19
- Coming's glass molding process, **I**:22.9, 22.9*f*
- Commando cloth, **IV**:6.31*f*, 6.33*f*
- Commercial sources (of radiation), **II**:15.13–15.53
- activated phosphor, **II**:15.49
 - blackbody simulators, **II**:15.14, 15.15*f*
 - carbon arcs, **II**:15.21–15.24, 15.23*f*, 15.24*f*, 15.25*t*–15.27*t*, 15.28*f*
 - concentrated arcs, **II**:15.47–15.49, 15.48*f*, 15.49*f*
 - glow modulator tubes, **II**:15.49, 15.50*f*, 15.51*f*, 15.52*t*
 - high-energy, **II**:15.40
 - high-pressure enclosed arc, **II**:15.24, 15.28–15.34
 - compact-source arcs, **II**:15.31–15.34, 15.32*f*–15.35*f*
 - Lucalox lamps, **II**:15.30, 15.31*f*
 - mercury arcs, **II**:15.30*f*, 15.31*f*, 15.31*t*
 - multivapor arcs, **II**:15.29, 15.31*f*
 - Uviarc, **II**:15.28–15.29, 15.29*f*, 15.30*f*
 - hydrogen and deuterium arcs, **II**:15.49, 15.53*f*
 - incandescent nongaseous, **II**:15.15–15.21
 - comparisons, **II**:15.19, 15.19*f*
 - gas mantle, **II**:15.17, 15.18, 15.19*f*
 - global, **II**:15.17, 15.18*f*

- Commercial sources (of radiation),
 incandescent nongaseous (*Cont.*):
 Nernst glower, **II**:15.14, 15.15, 15.17, 15.17*f*
 quartz-envelope lamps, **II**:15.20, 15.21
 tungsten-filament lamps, **II**:15.19, 15.20,
 15.20*f*–15.22*f*
- low-pressure enclosed arc, **II**:15.35–15.47
 black-light fluorescent lamps, **II**:15.35,
 15.36*t*
- electrodeless discharge lamps, **II**:15.36, 15.44
 germicidal lamps, **II**:15.35
 hollow cathode lamps, **II**:15.35,
 15.37*t*–15.43*t*, 15.44*f*
- Pluecker spectrum tubes, **II**:15.47, 15.47*f*,
 15.47*t*
- spectral lamps, **II**:15.44, 15.45, 15.45*f*,
 15.46*f*, 15.46*t*
- Sterilamps, **II**:15.35, 15.36*f*
- special-purpose, **II**:15.53
- Commission Internationale de l'Éclairage
 (CIE), **II**:40.2; **III**:7.9, 10.1, 10.3, 10.27
 publications from, **II**:37.11
 standard photometric observer, **II**:37.2
- Common glasses, **IV**:2.3
- Common path interferometers, **II**:13.9, 13.11*f*
- Communication networks and systems:
 fiber-optic standards for, **V**:23.1–23.8
 ATM/SONET, **V**:23.6
 ESCON, **V**:23.1–23.2, 23.2*f*
 Ethernet, **V**:23.7
 FDDI, **V**:23.2–23.3, 23.3*f*
 Fibre Channel standard, **V**:23.4, 23.5*f*, 23.5*t*
 InfiniBand, **V**:23.8, 23.8*t*
- optical fibers in, **V**:9.3–9.17
 analog transmission, **V**:9.15–9.17
 bit rate, **V**:9.12
 distance limits, **V**:9.12–9.13
 fiber for, **V**:9.4–9.7, 9.5*f*, 9.6*f*
 fiber-optic networks, **V**:9.14–9.15
 optical amplifiers, **V**:9.13–9.14
 photodetectors, **V**:9.8
 receiver sensitivity, **V**:9.8–9.11
 repeater spacing, **V**:9.12–9.13
 technology, **V**:9.4–9.8
 transmitting sources, **V**:9.7–9.8
- optical time-division multiplexed,
V:20.1–20.25
 analog to digital conversion, **V**:20.8, 20.8*f*
 binary digits and line coding in,
V:20.8–20.10, 20.9*f*
- Communication networks and systems, optical
 time-division multiplexed (*Cont.*):
 device technology, **V**:20.12–20.24
 history of, **V**:20.3
 interleaving in, **V**:20.6–20.7, 20.7*f*
 modulation, **V**:20.17*f*, 20.17–20.20, 20.19*f*,
 20.20*f*
 multiplexing and demultiplexing,
V:20.3–20.12, 20.7*f*, 20.22, 20.23*f*
 optical clock recovery, **V**:20.21*f*,
 20.21–20.22
 sampling, **V**:20.4–20.6, 20.5*f*, 20.6*f*
 serial vs. parallel, **V**:20.12, 20.13*f*
 timing recovery, **V**:20.10*f*, 20.10–20.12,
 20.11*f*
 transmitters, **V**:20.12–20.17, 20.14*f*–20.16*f*
 ultrahigh-speed, **V**:20.23–20.24, 20.24*f*
- solitons in, **V**:22.1–22.17
 classical solitons, **V**:22.2*f*–22.4*f*, 22.2–22.4
 design of transmission systems,
V:22.5–22.7
 dispersion-managed solitons,
V:22.12–22.15, 22.13*f*, 22.14*f*
 frequency-guiding filters, **V**:22.7–22.9
 wavelength division multiplexing,
V:22.9–22.12
- WDM dispersion managed soliton
 transmission, **V**:22.15–22.17
- wavelength-division multiplexed,
V:21.1–21.44
 carrier-suppressed return-to-zero and
 duobinary, **V**:21.30–21.33, 21.31*f*,
 21.32*f*
 chromatic dispersion in, **V**:21.14–21.16,
 21.15*f*, 21.16*f*
 circuit and packet switching in,
V:21.7–21.11, 21.8*f*–21.11*f*
 dispersion and nonlinearities of,
V:21.16–21.26, 21.17*f*–21.21*f*,
 21.23*f*–21.27*f*
 DPSK and DQSK, **V**:21.33*f*–21.35*f*,
 21.33–21.36, 21.36*t*, 21.37*t*
 fiber attenuation and optical power loss,
V:21.13–21.14
 fiber bandwidth, **V**:21.2*f*, 21.2–21.3
 fiber system impairments, **V**:21.13–21.26
 history of, **V**:21.1–21.2
 network reconfigurability, **V**:21.12*f*,
 21.12–21.13, 21.13*f*

- Communication networks and systems,
wavelength-division multiplexed (*Cont.*):
optical amplifiers in, **V**:21.37*f*, 21.37–21.44,
21.38*f*–21.42*f*
optical modulation formats, **V**:21.27–21.36,
21.28*f*–21.30*f*
point-to-point links, **V**:21.4
in real systems, **V**:21.3*f*, 21.3–21.4, 21.4*f*
star, ring, and mesh topologies,
V:21.5*f*–21.7*f*, 21.5–21.7
wavelength-routed networks, **V**:21.5, 21.5*f*
- Communications, out-of-plane coupling for,
IV:9.12
- Commutators, **V**:20.1, 20.7*f*
- Compact disks (CDs), **I**:35.1–35.2, 35.5*n*
- Compact fluorescent lights (CFLs), **II**:40.25*t*,
40.26*t*, 40.28*f*, 40.31
- Compact-source arcs, **II**:15.31–15.34,
15.32*f*–15.35*f*
- Compensators, **I**:13.53–13.56, 13.54*f*, 13.55*f*,
28.38
Babinet-Soleil, **V**:7.22
Dall, **II**:13.24, 13.24*f*
holographic, **II**:13.25
Offner, **II**:13.24, 13.24*f*
reflective, **II**:13.24, 13.24*f*, 13.25
refractive, **II**:13.24, 13.24*f*, 13.25
and tolerances, **II**:3.21, 5.7
- Complementary aperture screens, **I**:3.9–3.11
- Complementary metal oxide semiconductors
(CMOS) signal processing, **V**:62.5
- Complementary metal-oxide semiconductors
(CMOSs), **I**:26.8–26.9
- Complete monolithic FPAs, **II**:33.10
- Complex amplitude, **I**:2.4, 2.5, 11.2
- Complex cells (cortical neurons), **III**:2.14
- Complex degree of coherence, **I**:2.37, 5.4
- Complex degree of spectral coherence, **I**:5.5
- Complex Fresnel relation, **IV**:8.15
- Complex index of refraction, for x-ray optics,
V:48.1
- Complex refractive index, **I**:7.12–7.13, 12.5,
12.6; **IV**:1.16–1.17, 2.8
- Compliance tensors, **IV**:2.30, 2.31*t*
- Composite retardation plates, **I**:13.52, 13.53
- Composites, **IV**:12.26*t*, 12.27*t*
- Compositional modulation, **IV**:5.65, 5.66*t*, 5.67
- Compound elliptical collectors (CECs),
II:39.14, 39.15*f*, 39.27, 39.37
- Compound hyperbolic collectors (CHCs),
II:39.15, 39.15*f*, 39.16*f*, 39.37
- Compound lens, thermal defocus of, **II**:8.4, 8.5*f*
- Compound microscopes, **I**:17.10
- Compound mirror optics configurations,
I:30.15*f*, 30.15–30.16
- Compound parabolic collectors (CPCs),
II:39.13*f*, 39.13–39.14, 39.14*f*, 39.18, 39.19,
39.19*f*
- Compression molding, of polymers, **IV**:3.12
- Compressively-strained QW lasers, **II**:19.16*f*,
19.17
- Compton effect, **II**:23.9
- Compton radiation sources, **V**:55.2–55.3, 55.3*t*
- Compton scattering, **V**:59.1
and circular polarization, **V**:43.6
inverse, **V**:59.1
and microfocus x-ray fluorescence, **V**:29.5,
29.8*f*
and polycapillary x-ray optics, **V**:53.3, 53.15,
53.18
and refractive x-ray lenses, **V**:37.5, 37.7
and x-ray attenuation, **V**:31.2
and x-ray optics, **V**:26.7, 36.1
- Computed tomography (CT), **V**:31.5–31.7,
31.6*f*, 31.7*f*
- Computer graphics software, **II**:40.21–40.23,
40.22*f*–40.24*f*
- Computer numeric control (CNC) systems,
II:9.4
- Computer numerical control (CNC) lathe
turning, **IV**:3.12
- Computer Vision Syndrome (CVS),
III:23.1–23.12
disorders and eye conditions, **III**:23.9–23.12
accommodation, **III**:23.10–23.11
anisometropia, **III**:23.11
binocular vision, **III**:23.11
dry eyes, **III**:23.9
presbyopia, **III**:23.11–23.12
refractive error, **III**:23.10
and work environment, **III**:23.4–23.9
lighting, **III**:23.4–23.5, 23.5*t*
monitor characteristics, **III**:23.6–23.8
screen reflections, **III**:23.5–23.6
work habits, **III**:23.8–23.9
workstation arrangement, **III**:23.8
- Computer-aided design (CAD) software, **II**:40.19
- Computer-Automatic Virtual Environment
(CAVE), **III**:13.32

- Computer-generated holograms (CGHs),
II:14.1–14.9
 about, **II**:14.1–14.3
 accuracy limitations of, **II**:14.6*f*, 14.6–14.7,
 14.7*f*
 discussion of, **II**:14.9
 experimental results from, **II**:14.7*f*,
 14.7–14.9, 14.8*f*
 interferometers using, **II**:14.4*f*, 14.4–14.5,
 14.5*f*
 plotting of, **II**:14.3*f*, 14.3–14.4
 sample, **II**:14.2*f*
- Concave facets, **II**:39.40, 39.40*f*
- Concentrated arc lamps, **II**:15.47–15.49, 15.48*f*,
 15.49*f*
- Concentration:
 of radiation, **II**:39.1, 39.5, 39.6
 of solution, **II**:38.5
- Concentrators, nonimaging, **II**:39.12–39.22
 calculation of, **II**:39.5, 39.6
 compound elliptical collectors, **II**:39.14,
 39.15*f*
 compound hyperbolic collectors, **II**:39.15,
 39.15*f*, 39.16*f*
 compound parabolic collectors, **II**:39.13*f*,
 39.13–39.14, 39.14*f*
 dielectric compound parabolic collectors,
II:39.15, 39.16, 39.16*f*
 edge rays, **II**:39.22
 geometrical vector flux, **II**:39.21–39.22
 inhomogeneous media, **II**:39.22
 integrating cavities with, **II**:39.26, 39.27
 and lenses, **II**:39.16–39.17
 and mirrors, **II**:39.17
 multiple surface concentrators,
II:39.16–39.17, 39.17*f*
 restricted exit angle concentrators with
 lenses, **II**:39.18, 39.18*f*
 RX, **II**:39.17, 39.17*f*
 RXI, **II**:39.17, 39.17*f*
 star, **II**:39.20, 39.21
 tapered lightpipes, **II**:39.12–39.13, 39.13*f*
 θ_1/θ_2 concentrators, **II**:39.18–39.20, 39.19*f*
 2D vs. 3D, **II**:39.20*f*, 39.20–39.21, 39.21*f*
- Concomitant eye motion, **III**:13.1
- Condensers, first-order layout for, **II**:1.10–1.11,
 1.11*f*
- Condensing monocapillary x-ray optics, **V**:52.6
- Condition of detailed balance (term), **II**:23.23
- Condon point, for laser light, **IV**:20.29
- Conductance channels, in speckle patterns,
I:9.16
- Conduction band (CB), **II**:17.4, 17.4*f*, 17.5*f*;
IV:18.3
- Conduction bandgap, **II**:25.3, 25.3*f*
- Conductive magnetic brush (CMB), **I**:34.6
- Conductivity:
 diffraction and, **I**:3.32–3.33
 of metals, **IV**:4.6
 of paints, **IV**:6.12, 6.12*t*
 of polymers, **IV**:3.3–3.4
 of solids, **IV**:8.4
 of water, **IV**:1.16
- Cone contrast spaces, **III**:11.2, 11.32,
 11.41–11.42
- Cone contrasts, **III**:11.2, 11.32
- Cone coordinates, **III**:10.1, 10.11 (*see also*
 Tristimulus values)
- Cone effect, of laser beacons, **V**:5.27
- Cone fundamentals, **III**:10.13–10.14
 adjusting cone spectral sensitivities,
III:10.17–10.18
 defined, **III**:10.1, 10.4*t*
 online tabulation of, **III**:10.11
 primaries yielding, **III**:10.10
 Smith-Pokorny, **III**:10.12
 Stockman and Sharpe, **III**:10.12, 10.16
 tailored to individuals, **III**:10.15
- Cone magno pathway, **III**:2.9, 2.10*f*
- Cone mechanisms, **III**:11.10*f*, 11.11
 defined, **III**:11.2
 and stimulus direction, **III**:11.33
 and test measurements, **III**:11.12–11.14
- Cone parvo pathway, **III**:2.9, 2.10*f*
- Cone pathways, **III**:2.9, 2.10, 2.10*f*
- Cone polymorphism, **III**:10.15
- Cone-beam computed tomography (CT),
V:31.7, 31.7*f*
- Cone-excitation spaces, **III**:11.31–11.32
- Cone-opponent mechanisms, **III**:11.8
 and bichromatic test mixtures, **III**:11.39
 defined, **III**:11.2
 sensitization in, **III**:11.69
 unipolar vs. bipolar, **III**:11.80–11.81
- Cones (cone photoreceptors), **II**:30.15, 30.16*f*,
 34.37, 34.38, 36.8, 36.8*f*, 36.9*f*
 adaptation at, **III**:11.52
 and age-related photopic vision changes,
III:14.15
 alignment of, **III**:8.4

- Cones (cone photoreceptors) (*Cont.*):
 densities and distributions of, **III**:2.7
 in dichromatic observers, **III**:10.14
 directional sensitivity of, **III**:8.5
 ideal “average,” **III**:8.21, 8.21*f*
 inner segments, **III**:2.6*f*
 light-collection area of, **III**:2.8
 linear density of, **III**:2.6*f*
 and maximum saturation color matching,
III:10.7
 measuring activity of, **III**:15.24
 optical waveguide properties of, **III**:14.11
 outer segments of, **III**:2.6*f*
 pedicles of, **III**:8.20
 photocurrent responses of, **III**:2.8*f*
 in retinal layer of rods and cones model of
 biological waveguides, **III**:8.8–8.9, 8.9*f*,
 8.10*f*, 8.12–8.15
 S-cone flicker sensitivity, **III**:11.74–11.75
 spatial distribution of, **III**:2.6
 spectral sensitivities of, **III**:10.3, 10.13, 11.8
 adjusting for individual differences,
III:10.17–10.18
 estimates of, **III**:10.14, 10.14*f*
 loss hypothesis for, **III**:10.14
 in normal and dichromatic observers,
III:10.13
 time constant of photopigment regeneration,
III:2.7
 types and functions of, **III**:2.4, 2.6
 Cone-shaped secondary baffle, **II**:7.3*f*, 7.3–7.4,
 7.4*f*
 Confidence interval (CI), **II**:34.22
 Configuration factor algebra, **II**:34.14
 Configurational coordinate model, of
 lineshapes, **I**:10.22, 10.23*f*
 Configurational relaxation, in solids,
V:2.14–2.17, 2.15*f*–2.18*f*
 Confinement:
 in double heterostructure laser diodes,
V:13.4*f*, 13.4–13.6, 13.6*f*
 and photonic crystal fibers, **V**:11.22
 in semiconductor optical amplifiers, **V**:19.6
 Confocal Cassegrainians (telescopes), **I**:18.21
 Confocal cavity technique, **II**:12.20, 12.20*t*
 Confocal microscopes (confocal microscopy),
I:28.49*f*, 28.49–28.51, 28.51*f*; **III**:17.1–17.10
 clinical, **III**:17.3, 17.6–17.9
 clinical applications using, **III**:17.8–17.9
 defined, **III**:17.1
 Confocal microscopes (confocal microscopy)
 (*Cont.*):
 development of, **III**:17.3, 17.5
 laser scanning, **III**:17.3, 17.7–17.8
 Nipkow disk, **III**:17.4*f*, 17.5, 17.5*f*
 scanning slit, **III**:17.3, 17.6*f*, 17.6–17.9, 17.7*f*
 spatial filtering with, **III**:17.3
 Svishchev, **III**:17.6, 17.6*f*
 theory of confocal microscopy, **III**:17.3, 17.4*f*
 Confocal parabolas, **I**:18.20, 18.20*f*, 18.21
 Confocal parameter, **II**:16.23
 ConfoScan 4 microscope, **III**:17.6, 17.7, 17.9
 Conic collimators, **II**:39.8, 39.9*f*
 Conic constant (term), **I**:1.34, 29.37
 Conic mirrors, **I**:29.3*f*, 29.4*f*
 Conic reflectors, **II**:39.11, 39.11*f*
 Conic surfaces, **II**:3.5
 Conical (term), **II**:35.5
 Conical surfaces, in systems of revolution,
I:1.34–1.35
 Conical-directional reflectance, **II**:35.5*t*, 35.6*f*,
 35.6*t*
 Conical-hemispherical reflectance, **II**:35.5*t*,
 35.6*f*, 35.6*t*
 Conjugate (version) eye movements/position,
III:1.42, 1.43, 13.1, 13.7, 13.20
 Conjugate lines, collineation of, **I**:1.59
 Conjugate matrices, **I**:1.68–1.71, 1.73
 Conjugate mirrors, phase (*see* Phase conjugate
 mirrors)
 Conjugate planes:
 collineation of, **I**:1.58–1.59
 in microscopes, **I**:28.4–28.5
 Connected networks, 3D photonic crystals and,
IV:9.4–9.5
 Connector losses, in fiber optic communication,
V:15.7, 15.8*t*
 Conoscopic imaging, **I**:28.8–28.9
 Conservation, of radiant power transfer,
II:34.13*f*
 “Conservation of complexity,” **II**:3.7
 Conservation of étendue, law of, **I**:1.22
 Constellation-X Observatory, **V**:33.4
 Constraints:
 defined, **II**:3.17
 methods for handling, **II**:3.18–3.19
 Constricted double-heterostructure large optical
 cavity (CDH-LOC), **II**:19.19, 19.20*t*, 19.21*f*
 Constringence, **IV**:2.23 (*see also* Abbe number)
 Constructive interference, **I**:2.7

- Consultative Committee on Photometry and Radiometry (CCPR), **II**:36.2
- Contact lenses, **III**:20.1–20.34
accommodation with, **III**:20.26–20.30, 20.27*f*–20.29*f*, 20.29*t*
anisophoria and, **III**:13.26
for aphakic patients, **III**:12.15
base curve of, **III**:12.1
convergence with, **III**:20.25–20.26
correction with, **III**:12.11–12.14
hydrogel lenses, **III**:12.12–12.13
for presbyopia, **III**:12.13–12.14, 14.27–14.28
rigid lenses, **III**:12.11–12.12
- design considerations, **III**:20.20–20.25
aberrations, **III**:20.24–20.25, 20.25*f*
aspheric lenses, **III**:20.20–20.23, 20.21*f*–20.23*f*
posterior peripheral curve, **III**:20.23–20.24, 20.24*t*
- design parameters, **III**:20.2–20.6, 20.4*t*, 20.5*t*
base curve radius (BCR), **III**:20.3, 20.4*f*, 20.5
center thickness, **III**:20.6
edge thickness, **III**:20.6
optical zone diameter (OZD), **III**:20.5
overall diameter (OAD), **III**:20.5
posterior peripheral curve systems, **III**:20.5, 20.6
- magnification, **III**:20.31–20.33
relative spectacle, **III**:20.32–20.33, 20.33*f*
spectacle, **III**:20.31–20.32
- materials for, **III**:20.3
- power of, **III**:20.6–20.20
back and front vertex power, **III**:20.7–20.8, 20.8*t*
and contact lens as thick lens, **III**:20.6–20.7
effective power, **III**:20.8–20.12, 20.9*t*, 20.10*f*, 20.11*f*
lacrimial lens consideration, **III**:20.12*f*–20.14*f*, 20.12–20.15
residual astigmatism, **III**:20.15
of soft lenses, **III**:20.15–20.16
of toric lenses, **III**:20.16*f*–20.20*f*, 20.16–20.20, 20.18*t*
- prismatic effects with, **III**:20.30–20.31
prism-ballasted contacted lenses, **III**:20.30
unintentional (induced) prism, **III**:20.31
(see also *specific types of lenses*)
- Contact scanners, **II**:32.21, 32.22*f*
- Contact stresses, **II**:6.21
- Contacting, in wafer processing, **II**:17.24
- Contamination, of optical surfaces, **IV**:19.4–19.5
- Contamination control, **IV**:6.16
- Contamination levels (in stray light suppression), **II**:7.18–7.19, 7.18*t*, 7.19*f*–7.21*f*
- Continuous polishers (CPs), **II**:9.7
- Continuous readout x-ray detectors, **V**:61.2
- Continuous ring flanges, **II**:6.4*f*, 6.11
- Continuous wave (cw) power, **II**:19.19, 19.22, 19.22*f*
- Continuous-wave (cw) lasers, **II**:23.18, 34.32; **IV**:7.14, 14.16–14.18, 14.17*f*; **V**:25.4, 25.5, 25.5*f*, 25.7
diode lasers, **V**:13.49
dye lasers, **I**:10.8; **V**:5.32, 5.33*f*, 5.34
- Continuous-wave optical parametric oscillators (cw OPOs), **IV**:17.1–17.31
cavity resonance configurations for, **IV**:17.2–17.4, 17.3*f*, 17.4*f*
for correlated twin beams of light, **IV**:17.28, 17.29*f*, 17.30*f*
for hyperspectral imaging, **IV**:17.27–17.28
limitations of, **IV**:17.30
for metrology and optical frequency synthesis, **IV**:17.28, 17.29
multiple-resonant oscillators, **IV**:17.16–17.21
doubly resonant, **IV**:17.16–17.17
pump enhanced singly resonant, **IV**:17.17–17.20, 17.18*f*–17.20*f*
triply resonant, **IV**:17.20–17.21, 17.21*f*
singly resonant oscillators, **IV**:17.4–17.16
guided-wave nonlinear structures, **IV**:17.15–17.16
MgO:sPPLT in, **IV**:17.14–17.15, 17.15*f*
PPLN crystals in, **IV**:17.4–17.13, 17.6*f*–17.11*f*
QPM nonlinear materials, **IV**:17.13–17.14
in spectroscopy, **IV**:17.21–17.27
high-resolution Doppler-free, **IV**:17.27
photoacoustic, **IV**:17.22*f*–17.24*f*, 17.22–17.24
single-pass absorption, **IV**:17.24–17.27, 17.25*f*, 17.26*f*
- technological advances in, **IV**:17.1–17.2, 17.30–17.31
- Continuous-wave (cw) Q-switched modelocking, **IV**:18.5*f*
- Continuous-wave (cw) spectroscopy, **IV**:11.2
- Continuum excitations, **IV**:18.20

- Continuum pulse generation, **IV**:18.4
- Contours (color vision), **III**:11.72–11.79, 11.73*f*–11.75*f*
 and after-effects, **III**:11.75*f*, 11.75–11.76
 color appearance and stabilized borders, **III**:11.74–11.75
 detection surface/contour, **III**:11.12
 defined, **III**:11.2
 and directions of color spaces, **III**:11.40*f*
 in equiluminant plane, **III**:11.44*f*
 and field adaptation, **III**:11.53–11.54, 11.54*f*
 in the L, M plane, **III**:11.40*f*, 11.41–11.42
 discrimination contours, **III**:11.44*f*
 gap effect and luminance pedestals, **III**:11.72, 11.74
 McCollough effect, **III**:11.76, 11.77*f*
 multiplexing chromatic and achromatic signals, **III**:11.76–11.79, 11.78*f*
 threshold surface/contour, **III**:11.12–11.15, 11.13*f*
 defined, **III**:11.3
 and loss of information, **III**:11.20
 and noise, **III**:11.20, 11.23*f*, 11.23–11.26, 11.25*f*
 and second-site adaptation to steady fields, **III**:11.18*f*
- Contrast:
 in color CRTs, **III**:22.19
 defined, **III**:23.2
 in microscopy
 bright field microscopy, **I**:28.25, 28.27*f*, 28.27–28.28
 dark field microscopy, **I**:28.28
 Hoffman modulation contrast, **I**:28.29
 interference microscopy, **I**:28.33–28.44, 28.35*f*, 28.37*f*, 28.38*f*, 28.40*f*, 28.42*f*, 28.43*f*
 and modulation transfer function, **I**:28.24–28.25, 28.25*f*, 28.26*f*
 phase contrast, **I**:28.28–28.29, 28.29*f*
 SSEE microscopy, **I**:28.29, 28.30, 28.30*f*–28.33*f*, 28.33
 in monochrome CRTs, **III**:22.7, 22.8*f*
 in retinal imaging, **III**:15.21–15.22
 in vision experiments, **III**:3.4
 and visual acuity, **III**:4.12, 4.12*f*
- Contrast coding:
 first-site adaptation, **III**:11.15–11.16
 Weber's law and, **III**:11.15–11.16, 11.16*f*
- Contrast (modulation) color spaces, **III**:10.19
 Contrast constancy, **III**:2.33
 Contrast detection, **III**:2.19–2.31
 adaptation and inhibition, **III**:2.25*f*, 2.25–2.27
 chromatic, **III**:2.29–2.31, 2.30*f*
 eye movements, **III**:2.21
 as function of spatial frequency, **III**:2.20*f*
 optical transfer function, **III**:2.21–2.22, 2.22*f*
 optical/retinal inhomogeneity, **III**:2.24, 2.25*f*
 receptors, **III**:2.22–2.23
 spatial channels, **III**:2.23–2.24
 temporal, **III**:2.27*f*, 2.27–2.29, 2.29*f*
- Contrast discrimination, **III**:2.31*f*, 2.31–2.32
 Contrast estimation, **III**:2.33
 Contrast masking, **III**:2.32*f*, 2.32–2.33
 Contrast ratio, for fiber optic modulation, **V**:13.57
 Contrast sensitivity, **III**:15.1
 Contrast sensitivity functions (CSFs), **III**:2.20*f*, 2.20–2.23, 2.22*f*, 24.3
 and adaptation, **III**:2.25*f*, 2.25–2.27
 age-related changes in, **III**:14.17, 14.19
 chromatic, **III**:2.29–2.31
 in color vision, **III**:11.45–11.46, 11.46*f*
 at different retinal eccentricities, **III**:2.24, 2.25*f*
 of ideal observer, **III**:2.24
 for motion detection, **III**:2.36, 2.37
 spatial, **III**:11.45–11.46, 11.46*f*
 spatio-temporal, **III**:2.28–2.29
 for stereopsis, **III**:2.40, 2.41
 temporal, **III**:2.27–2.29, 2.29*f*, 11.45
- Contrast transfer function (CTF), **I**:4.7–4.8, 4.8*f*, 28.24; **III**:4.5
 Contrast-to-noise ration (CNR), of x-ray detectors, **V**:61.3
 Controlled grinding, **IV**:19.3
 Controlled-drift x-ray detectors, **V**:62.5
 Conventional evaporation, **IV**:7.11
 Convergence, **III**:13.21
 in color CRTs, **III**:22.12, 22.13
 with contact lenses, **III**:20.25–20.26
 defined, **III**:13.1, 23.2
 and egocentric direction, **III**:13.7
 and refractive errors, **III**:12.15–12.16
 Convergence accommodation/convergence interaction (CA/C ratio), **III**:13.24
 Convergence micropsia, **III**:13.19
 Convergence responses to prism, **III**:13.19

- Convergent reflectors, **II**:39.38*f*, 39.38–39.40, 39.39*f*
- Converging neutron guides, **V**:63.17
- Conversion factors:
for English and SI units, **II**:37.7*t*
for photometric and radiometric quantities, **II**:36.11–36.14, 36.12*f*–36.14*f*
- Convex surfaces, testing of, **II**:14.5, 14.5*f*
- Cook objectives, **I**:29.31
- Cooling:
of atoms with atom–laser interactions, **IV**:20.3–20.21
below Doppler limit, **IV**:20.17–20.21, 20.18*f*–20.20*f*
history of, **IV**:20.3–20.4
optical molasses, **IV**:20.13–20.17, 20.14*f*–20.16*f*
properties of lasers, **IV**:20.4–20.6
slowing atomic beams, **IV**:20.11–20.13, 20.12*f*, 20.12*t*, 20.13*f*
theoretical description, **IV**:20.6–20.11, 20.9*f*
- Doppler, **IV**:20.13–20.15, 20.14*f*
laser (*see* Laser cooling)
polarization gradient, **IV**:20.17
Raman, **IV**:20.21
- Cooling rate, for glasses, **IV**:2.5*n*
- Coordinate measurement machines (CMMs), **II**:9.6
- Coordinate measurement method (CMM), **II**:40.53, 40.54
- Coordinate systems:
for aberrations of point images, **I**:1.86
for collineation, **I**:1.57
for Fresnel equations, **I**:12.6–12.7, 12.7*f*
left-handed, **I**:12.6
for Mueller matrices, **I**:14.19–14.20
- Coordinate-measuring machines (CMMs), **V**:46.2
- Coordinates, circumferential, **V**:45.7
- Copepod, **III**:8.1
- Copilia*, **III**:8.4
- Copper, **II**:17.28
absorptance of, **IV**:4.40*f*, 4.48*t*, 4.50*t*
optical properties of, **IV**:4.12*t*–4.13*t*, 4.22*f*
physical properties of, **IV**:4.52*t*–4.54*t*
reflectance of, **IV**:4.29*t*–4.30*t*, 4.40*f*
thermal properties of
coefficient of linear thermal expansion, **IV**:4.56*t*, 4.57*f*
elastic properties, **IV**:4.69*t*
- Copper, thermal properties of (*Cont.*):
at room temperature, **IV**:4.55*t*
specific heat, **IV**:4.65*t*, 4.66*f*
strength and fracture properties, **IV**:4.70*t*
thermal conductivity, **IV**:4.58*t*, 4.60*f*–4.61*f*
- Copper black, **IV**:6.21, 6.23*f*
- Copper distributed data interface (CDDI), **V**:23.3
- Copper gallium sulfide (CuGaS₂), **IV**:2.39*t*, 2.45*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.62*t*, 2.74*t*
- Copper vapor lasers (CVLs), **II**:16.12, 16.13*f*, 16.30
- Copper-doped germanium (Ge:Cu) detectors, **II**:24.84*f*, 24.85*f*, 24.96, 24.97, 24.97*f*–24.99*f*
- Core (term), **V**:25.2
- Core drilling, of fiber lasers, **V**:25.26–25.27
- Core excitons, **IV**:5.26*t*
- Core-cladding index difference, of photonic crystal fibers, **V**:11.12*f*–11.17*f*, 11.12–11.17
- Cornea, **III**:1.3*f*, 1.4–1.6, 16.3–16.4, 21.2
absorption of ultraviolet light at, **III**:1.9
aging-related changes in, **III**:14.5, 14.6*f*
asphericity in, **III**:1.5*f*
and cataract (*see* Cataract)
confocal microscopy of (*see* Confocal microscopes)
as entrance pupil, **III**:1.8
injury to, **III**:7.4
keratometry of, **III**:21.6
laser ablation, **III**:16.11–16.19
ablation profiles, **III**:16.14–16.15
ablation rate, **III**:16.16–16.18, 16.18*f*
corneal photoablation, **III**:16.16, 16.17*f*
Epi-LASIK, **III**:16.12, 16.13, 16.13*f*
LASEK, **III**:16.12, 16.13
LASIK, **III**:16.13–16.14, 16.14*f*
photorefractive keratectomy (PRK), **III**:16.11*f*, 16.11–16.12, 16.12*f*
thermal, photochemical, and photoacoustic effects, **III**:16.18–16.19
optical axis of, **III**:9.5
optical power of, **III**:16.2
and refraction in the eye, **III**:12.3
refractive index of, **III**:14.5
refractive surgery modalities
corneal incisions/implants, **III**:16.9–16.11
laser corneal procedures, **III**:16.11–16.15
and rigid contact lenses, **III**:12.12

- Cornea (*Cont.*):
 spatially modulated excimer laser ablation
 of, **III**:1.25
 spectral sensitivities at, **III**:10.18–10.19
 stray light from, **III**:1.20
 Corneal endothelium, **III**:16.4
 Corneal epithelium, **III**:16.4
 Corneal hydration, **III**:16.18
 Corneal incisions/implants, **III**:16.9–16.11,
 16.10*f*
 Corneal photoablation, **III**:16.16–16.19, 16.17*f*
 ablation rate, **III**:16.16–16.18, 16.18*f*
 thermal, photochemical, and photoacoustic
 effects, **III**:16.18–16.19
 Corneal relaxing incisions, **III**:21.13
 Corneal stroma, **III**:16.4
 Cornell Black, **IV**:6.26*f*, 6.27
 Cornell High Energy Synchrotron Sources
 (CHESS), **V**:52.3, 52.5
 Corner cube prisms, **II**:12.16
 Cornice lighting, **II**:40.13*f*
 Corning ULE, **II**:6.18
 Cornu equation, for refraction index, **IV**:2.22
 Cornu's spiral, **I**:3.16–3.19, 3.18*f*
 Coroneo effect, **III**:7.5, 7.6*f*, 7.7
 Corotron, in xerographic systems, **I**:34.2, 34.3*f*
 Correctors, for reflective and catadioptric
 objectives:
 anaplanatic, anastigmatic Schwarzschild with
 aspheric corrector plate, **I**:29.13
 Cassegrain with spherical secondary and
 field corrector, **I**:29.8–29.9
 Mangin-Cassegrain with correctors, **I**:29.24
 Ritchey-Chretien telescope with two-lens
 corrector, **I**:29.8
 spherical-primary Cassegrain with reflective
 field corrector, **I**:29.9
 three-lens prime focus corrector, **I**:29.10
 Correlated color temperature (CCT), **II**:34.44,
 37.7, 38.5, 40.8
 Correlated double sampling (CDS), **I**:26.11;
II:33.13
 Correlated emission lasers (CELs), **II**:23.42–23.43
 Correlated twin beams of light, **IV**:17.28,
 17.29*f*, 17.30*f*
 Correlators, acousto-optic, **I**:11.10–11.12,
 11.11*f*
 Corresponding retinal points, **III**:13.1, 13.8
 Cortical cataract, **III**:12.14, 14.8
 Cortical neurons, **III**:2.14
 CORTRAN glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*,
 2.67*t*
 Corundum, **IV**:2.70*t*
 Cosine condition, of stigmatic imaging, **I**:1.30,
 1.30*f*
 Cosine law, **II**:37.8, 37.8*f*
 Cosine-to-the-fourth approximation, **I**:1.81;
II:34.16
 Cotton polarizers, **I**:13.21
 Cotyledons, **III**:8.1, 8.26, 8.27*f*
 Couder objective, **I**:29.12
 Coulomb attraction, **IV**:8.31
 Coulomb explosions, **V**:2.5, 2.6
 cluster, **IV**:21.33–21.34
 molecular, **IV**:21.24–21.25
 Coulomb fields, of ions, **V**:56.8
 Coulomb gauge, for solids, **IV**:8.5
 Coulomb interactions, of 3d electrons, **V**:2.9
 Coulomb potentials, **IV**:21.31, 21.32*f*
 Coulomb repulsion, **I**:10.10, 10.12, 10.16; **V**:2.5
 Coupled cavity lasers, **II**:19.37*f*, 19.38
 Coupled plasmon-phonon behavior, **IV**:5.35,
 5.36, 5.36*f*, 5.37*f*
 Coupled quantum wells (CQWs), **V**:13.58
 Coupled resonator structures (CRS),
IV:22.11–22.13, 22.12*f*
 Coupled-dipole method, **I**:7.15
 Coupled-wave-theory (CWT), for multilayer
 Laue lenses, **V**:42.4, 42.12
 Couple-mode theory, in integrated optics,
I:21.8
 Couplers and coupling:
 of circulating pulses, **II**:20.12*f*, 20.12–20.15,
 20.15*f*
 directional, **V**:18.2, 18.3, 18.3*f*, 18.9, 18.9*f*
 étendue and source, **II**:40.41–40.42
 evanescent, **V**:25.11
 fiber-based, **V**:16.1–16.6, 16.2*f*, 16.5*f*
 in film development, **II**:29.14
 gain, **II**:19.29
 in-plane, **IV**:9.10–9.11
 interface, **II**:20.14
 nonfused fiber, **V**:25.10
 out-of-plane, **IV**:9.11–9.12
 output, **II**:16.13
 phase-conjugated, **II**:20.14
 repetition-rate, **II**:20.14–20.15, 20.15*f*
 two-beam, **IV**:12.4*f*, 12.4–12.6, 13.8–13.9
 vertical, **V**:13.60
 wavelength-selective, **V**:14.2

- Coupling coefficient, of DBR lasers, **V**:13.28
- Coupling laser power, **IV**:14.14–14.15
- Coupling (tie) layer, of bandpass filters, **IV**:7.82
- Coupling loss, of fiber-optic components, **V**:18.1
- Coupling noise, resistive, **II**:27.5, 27.6f
- Coupling schemes, **IV**:14.30f (*see also specific coupling schemes, e.g.: Lambda coupling*)
- Covariance mapping, **V**:2.5
- Cove lighting, **II**:40.13f, 40.16f
- Coverslip correction, **I**:28.10–28.11, 28.12f, 28.13
- CR-39 resin (poly-diallylglycol), **IV**:3.11
- Craik-O'Brien-Cornsweet illusion, **III**:11.78
- Craters, surface absorption and, **IV**:6.15
- Creep strength, of metals, **IV**:4.8
- Critical fusion frequency (CFF), **III**:14.19
- Critical illumination, **I**:28.7 (*see Abbe illumination system*)
- Critical objects (in stray light suppression), **II**:7.2
 imaged, **II**:7.4, 7.5f
 real-space, **II**:7.2–7.4, 7.3f, 7.4f
- Cross-coupling, direction/distance of gaze and, **III**:13.23f, 13.23–13.24
- Crossed reflectors, **II**:39.38f
- Crossed string relationship, **II**:39.4, 39.4f
- Cross-gain modulation (XGM), **V**:19.12, 19.13f, 19.27, 19.29f, 19.29–19.30, 19.32, 19.35–19.36
- Cross-Kerr phase shifter, **IV**:23.11
- Crossover frequency, liquid crystals and, **V**:8.16
- Cross-phase modulation (XPM):
 in optical fibers, **V**:10.3–10.4
 and SOAs, **V**:19.13, 19.30–19.32, 19.31f, 19.33f, 19.35–19.36
 and solitons, **V**:22.5, 22.13–22.15
 in WDM networks, **V**:21.19
- Cross-spectral density function, **I**:5.5, 5.9, 5.10, 5.16
- Crown glasses, **IV**:2.28, 2.41t, 2.42t
- Crucifix images, **V**:49.4, 49.4f
- Cryogenic electrical substitution radiometers, **II**:34.28
- Cryogenic x-ray detectors, **V**:60.9, 60.9t, 60.10t
- Crystal diffraction:
 and EDXRF, **V**:29.6–29.7, 29.8f–29.9f
 single, **V**:53.12f–53.13f, 53.12–53.14
 and WDXRF, **V**:29.2
- Crystal interferometers, **V**:63.26–63.27, 63.27f
- Crystal monochromators:
 and bent crystals, **V**:39.1–39.6, 39.2t, 39.3f, 39.5f–39.6f
 in neutron optics, **V**:63.23–63.25
- Crystal optics:
 and electro-optic modulators, **V**:7.3–7.4, 7.4f, 7.8t–7.10t
 and the index ellipsoid, **V**:7.3–7.7, 7.4f–7.6f, 7.8f–7.10f
 polarizing, **V**:43.1–43.8, 43.3f–43.5f, 43.7f
 and x-ray diffraction, **V**:28.3–28.4, 28.4f
- Crystallin proteins, **III**:19.1
- Crystalline infrared fibers, **V**:12.2t, 12.3t, 12.7–12.10, 12.8f, 12.10f
- Crystalline lens, **III**:1.3f, 16.3, 21.2
 and accommodation, **III**:1.30–1.31
 aging-related changes in, **III**:14.7–14.8
 cataract of, **III**:12.14, 21.3 (*see also Cataract; Intraocular lenses*)
 and cone spectral sensitivities, **III**:10.17
 contours of refractive index in, **III**:1.5f
 distribution of refractive index in, **III**:1.5f, 1.5–1.6
 and fluctuations in accommodation, **III**:1.32
 fluorescence in, **III**:1.21
 gradient index structure of, **III**:1.18, 21.2
 growth of, **III**:1.4
 and refraction in the eye, **III**:12.3
 refractive index gradient of, **III**:1.3, 19.12–19.13, 19.13f, 19.14f
 scattered light from, **III**:1.20
 UV absorption in, **III**:1.9
- Crystalline optics, **II**:9.8
- Crystalline-quartz retardation plates, **I**:13.46–13.48
- Crystals:
 anisotropic, **IV**:8.8–8.11, 8.9t, 8.10f
 bent, **V**:39.5f–39.6f, 39.5–39.6
 biaxial, **IV**:8.8, 8.9t, 8.10, 8.10f
 Bragg-symmetric, **V**:35.3
 Darwin width of, **V**:63.24, 63.26
 and dielectric tensor and optical indicatrix, **IV**:2.17–2.19, 2.19f
 and dispersion formulas for refractive index, **IV**:2.21–2.22
 doubly curved, **V**:29.6–29.7, 29.8f
 in electro-optic modulators, **V**:7.33–7.34
 Fankuchen-cut, **V**:63.25
 and glasses, **IV**:2.1–2.77

Crystals (*Cont.*):

- lattice vibration model parameters for, **IV:2.76t–2.77t**
- material properties of, **IV:2.27–2.36**
 - characteristic temperatures, **IV:2.32, 2.33**
 - combinations of, **IV:2.36**
 - correlations of, **IV:2.36**
 - elastic properties, **IV:2.30–2.31, 2.31t**
 - hardness and strength, **IV:2.31–2.32, 2.32f, 2.32t**
 - heat capacity and Debye temperature, **IV:2.33–2.34**
 - material designation and composition, **IV:2.27, 2.29f**
 - naming of, **IV:2.27**
 - thermal conductivity, **IV:2.35f, 2.35–2.36**
 - thermal expansion, **IV:2.34f, 2.34–2.35**
 - unit cell parameters, molecular weight, and density, **IV:2.30**
- mechanical properties of, **IV:2.47t–2.48t**
- mosaic, **V:39.2**
- nonlinear optical, **IV:10.19–10.20, 10.20t–10.22t**
- and nonlinear optical coefficients, **IV:2.26–2.27, 2.27t**
- optical activity of, **IV:8.10–8.11**
- optical applications of, **IV:2.17–2.27**
- as optical materials, **IV:2.4–2.5**
- optical modes of, **IV:2.68t–2.76t**
 - with α -quartz structure, **IV:2.71t**
 - with cesium chloride structure, **IV:2.68t**
 - with chalcopyrite structure, **IV:2.74t**
 - with corundum structure, **IV:2.70t**
 - with cubic perovskite structure, **IV:2.73t**
 - with diamond structure, **IV:2.68t**
 - with fluorite structure, **IV:2.69t**
 - other structures, **IV:2.74t–2.76t**
 - with rutile structure, **IV:2.71t**
 - with scheelite structure, **IV:2.72t**
 - with sodium chloride structure, **IV:2.69t**
 - with spinel structure, **IV:2.73t**
 - with tetragonal perovskite structure, **IV:2.73t**
 - with trigonal selenium structure, **IV:2.70t**
 - with wurtzite structure, **IV:2.70t**
 - with zinblend structure, **IV:2.69t**
- optical properties of, **IV:2.6, 2.8–2.9**
- organic, **V:7.33–7.34**

Crystals (*Cont.*):

- origin and models of, **IV:2.9–2.17, 2.10f**
 - absorption in transparent region, **IV:2.17**
 - electronic transitions, **IV:2.12–2.15, 2.13f**
 - lattice vibrations, **IV:2.11–2.12**
 - multiphoton absorption and refraction, **IV:2.15–2.17, 2.16f, 2.17f**
- and photoelastic coefficients, **IV:2.24**
- physical properties of, **IV:2.37, 2.38t–2.43t**
 - classes and symmetry properties, **IV:2.7t**
 - composition, structure, and density, **IV:2.38t–2.41t**
 - physical constants, **IV:2.8t**
 - symmetry properties, **IV:2.5, 2.6t–2.8t**
- Raman scattering in, **IV:8.19t–8.20t**
- room-temperature dispersion formulas, **IV:2.60t–2.66t**
- room-temperature elastic constants of
 - cubic crystals, **IV:2.44t–2.49t**
 - hexagonal crystals, **IV:2.46t**
 - monoclinic crystals, **IV:2.47t**
 - orthorhombic crystals, **IV:2.46t**
 - tetragonal crystals, **IV:2.44t–IV:2.45t**
- and scatter, **IV:2.27**
- SHADOW code for, **V:35.2**
- thermal properties, **IV:2.50t–2.53t**
- and thermo-optic coefficients, **IV:2.24–2.26**
- and total power law, **IV:2.19–2.20, 2.20f**
- uniaxial, **IV:8.8, 8.9t, 8.10f**
 - [*see also* Liquid crystals; Photonic crystal fibers (PCFs)]
- Cubic crystals, **IV:8.9t, 8.20t**
 - dielectric constants of, **IV:2.18**
 - room-temperature elastic constants of, **IV:2.44t–2.49t**
 - symmetries of, **IV:2.7t**
- Cubic oxides (sillenites), **IV:12.17–12.19, 12.18t, 12.19f**
- Cubic perovskite structure, of crystals and glasses, **IV:2.73t**
- Cubic zirconia ($\text{ZrO}_2\cdot 0.12\text{Y}_2\text{O}_3$), **IV:2.41t, 2.44t, 2.48t, 2.69t**
- Cumulative size distribution, of particles in water, **IV:1.15**
- Curie temperature, **I:35.25**
- Curie temperature, of crystals and glasses, **IV:2.33**
- Current:
 - dark, **V:13.69, 13.73**
 - saturation, **V:13.69**
 - trap-assisted thermal generation, **V:13.69**

- Current confinement, in VCSELs, **V**:13.45, 13.45*f*
- Current density, **II**:19.12*f*, 19.12–19.13, 19.13*f*
- Current-confined constricted double-heterostructure large optical cavity (CC-CDH-LOC), **II**:19.19, 19.20*t*, 19.21*f*
- Curvature:
- of space curves, **I**:1.18–1.19
 - vertex (paraxial), **I**:1.32–1.33
- Curvature measurement, **II**:12.17–12.25
- mechanical methods of, **II**:12.17–12.19, 12.18*f*, 12.19*f*, 12.19*t*
 - optical methods of, **II**:12.19–12.21, 12.20*f*, 12.20*t*, 12.21*f*
- Curvature of best focal surface, **V**:45.5
- Curved surfaces, radial gradients with, **I**:24.7
- Cusp surface, of diamond-turned optics, **II**:10.10, 10.10*f*
- Cutoff, slope of, **IV**:7.56
- Cutoff filters, **IV**:7.53–7.60, 7.54*f*, 7.55*f*, 7.57*f*, 7.59*f*–7.61*f*
- Cutoff wavelength, **II**:24.10
- Cw modelocking, **IV**:18.5*f*
- Cyanine dyes, **II**:30.13, 30.13*f*
- Cyclopean eye, **III**:13.1
- Cyclopean locus, **III**:13.7
- Cyclophoria, **III**:13.35
- Cycloplegia, **III**:1.25, 25.11
- Cyclotron resonance (CR), **IV**:5.12, 5.12*f*, 5.40, 5.47–5.50, 5.48*f*–5.50*f*
- Cyclovergence, **III**:13.1, 13.22, 13.27
- Cylinders, scattering by, **I**:7.14
- Cylindrical lenses, **I**:22.45, 22.46*f*
- Cylindrical wavefronts, **I**:3.13–3.21, 3.14*f*
- and Cornu's spiral, **I**:3.16–3.19, 3.18*f*, 3.19*t*
 - and opaque strip construction, **I**:3.20–3.21
 - from rectangular apertures, **I**:3.19–3.20
 - from straight edges, **I**:3.14–3.16
- Cytochrome oxidase “blobs,” **III**:2.14
- Czerny-Turner monochromators, **I**:31.6
- Czerny-Turner mounts, **I**:20.8*f*, 20.14*t*
- Da Vinci stereopsis, **III**:13.4
- Dall compensators, **II**:13.24, 13.24*f*
- Dall-Kirkham objective, **I**:29.8
- Damage:
- laser-induced [see Laser-induced damage (IID)]
 - optical, **V**:13.54
- Damage threshold, of fiber lasers, **V**:25.7
- Dammann approach, for binary gratings, **I**:23.12
- Damped least-squares (DLS) method, **II**:3.17–3.19
- Damping:
- additive, **II**:3.18
 - of field by reservoir, **II**:23.33–23.34
 - in laser cooling, **IV**:20.19–20.20, 20.20*f*
 - phonon, **IV**:5.14
- Damping factor, **II**:3.18
- Dark counts (of photomultipliers), **II**:27.8
- Dark current, **V**:13.69, 13.73
- absorption coefficient, **II**:25.8, 25.9*f*
 - in CCDs, **II**:32.20
 - correction of, **II**:34.33
 - defined, **II**:24.10
 - diffusion current, **II**:25.7
 - generation-recombination current, **II**:25.7–25.8
 - histogram of, **II**:32.12*n*
 - in photosensing elements, **II**:24.19–24.20, 32.10–32.12, 32.11*f*
 - in *pin* photodiodes, **II**:25.7–25.8
 - quantum efficiency, **II**:25.8
 - responsivity, **II**:25.8
 - tunneling current, **II**:25.8
- Dark decay, of xerographic photoreceptors, **I**:34.3
- Dark field microscopy, **I**:28.28
- Dark focus, **III**:1.33, 1.34, 1.34*f*
- Dark regions, **II**:17.28
- Dark signal correction, **II**:34.33
- Dark (noncoupled) states, **IV**:14.4, 14.6–14.7, 20.37–20.39, 20.38*f*
- Dark-line defects, **II**:17.28
- Darwin widths, **V**:43.6, 63.24, 63.26
- Dashed rays, **II**:1.12, 1.12*f*
- Data communication systems, **V**:15.1–15.2
- Data noise, **I**:35.24
- Data rates, **I**:30.6–30.8
- Data storage, photorefractive holographic, **IV**:12.37
- Data transmission formats, **V**:20.9*f*, 20.9–20.10
- Data-reduction equations, polarimetric, **I**:15.14–15.15
- Daylight:
- as natural light source, **II**:40.40–40.41
 - simulation of, **II**:40.17
 - spectrum of, **II**:40.40*f*

- Daylighting schemes, luminaires for, **II**:40.47–40.50, 40.49*f*–40.51*f*
- Day/night cameras, **II**:31.28–31.29
- Dazzle, **II**:40.9
- dc carbon arcs, **II**:15.25*t*
- dc Kerr effect, **V**:7.11
- dc lamps, **II**:15.32*f*
- De Broglie neutron waves, **V**:63.3, 63.5
- De Valois De Valois zone model, **III**:11.82
- De-Broglie wavelengths, **IV**:8.4
- Debye molar heat capacity, **IV**:2.33–2.34
- Debye temperature, for crystals and glasses, **IV**:2.33–2.34
- Debye-Waller factor, **V**:41.4
- Decay, homogeneous, **IV**:11.10, 11.11*f*
- Decay time, **II**:16.4
- Decentered lens arrays, for agile beam steering, **I**:30.57–30.60, 30.58*f*–30.60*f*, 30.62–30.63
- Decentering errors, for grazing incidence optics, **V**:45.6–45.7
- Decision tasks (in experiments), **III**:3.2
- Decommutator, **V**:20.1, 20.7*f*
- Decomposition:
 - of conjugate matrices, **I**:1.71–1.72
 - of Mueller matrices, **I**:14.33
 - Cholesky decomposition, **I**:14.41, 14.42
 - polar decomposition, **I**:14.39–14.40
 - SVD, **I**:15.25–15.27
 - Weyl's plane-wave, **I**:3.23
- Decorative lighting, **II**:40.14, 40.16*f*
- Deep reactive ion etching (DRIE), **I**:22.7, 22.23, 30.62
- Deep saturation regime, for fiber amplifiers, **V**:14.4
- DEEP SPACE BLACK, **IV**:6.49
- Deep x-ray lithography (DXRL), **I**:22.7
- Deep-diffused stripe (DDS) lasers, **II**:19.23
- Defect modes (in photonic crystals), **IV**:22.11, 22.12*f*
- Defect-related absorption, **IV**:5.37–5.39, 5.38*f*, 5.39*f*
- Deflectors, **V**:6.23*t*, 6.22–6.31
 - anisotropic acoustic beam collimation by, **V**:6.25
 - birefringent phased array, **V**:6.29
 - high-resolution, **V**:6.29*t*, 6.29
 - isotropic AO diffraction by, **V**:6.23–6.24, 6.24*f*
 - phase array beam steering by, **V**:6.27–6.29, 6.28*f*
 - resolution of, **V**:6.25
 - tangential phase matching by, **V**:6.25–6.27, 6.26*f*
- Defocus, **I**:1.82, 1.83, 1.83*f*, 1.85–1.86, 1.90, 1.91
 - aberrations of, **II**:11.30
 - annular polynomials for, **II**:11.38*f*, 11.39
 - and efforts to clear vision, **III**:13.22–13.23
 - with head-mounted displays, **III**:25.9–25.10, 25.9*t*–25.10*t*
 - modulation transfer with, **III**:1.29
 - with monovision, **III**:14.28
 - off-axis, **III**:1.18, 1.19
 - Rayleigh units of, **III**:1.29
 - spectacles for, **III**:15.2
 - thermal, **II**:8.4, 8.5*f*
 - and visual acuity, **III**:4.9, 4.9*f*, 4.10*f*
 - with visual instruments, **III**:1.28
 - (*see also* Blur)
- Defocus errors, for grazing incidence optics, **V**:45.6
- Defocusing:
 - ionization-induced, **IV**:21.43*f*, 21.43–21.44
 - self, **IV**:13.7, 13.8, 19.9–19.11, 19.10*f*
 - thermal, **IV**:13.8
- Deformable mirrors, **III**:15.1, 15.10, 15.10*f*; **V**:5.4, 5.4*f*, 5.37*f*, 5.37–5.38, 5.38*f*
- Degauss, in color CRTs, **III**:22.9
- Degenerate four-wave mixing (DFWM), **IV**:16.27–16.28, 16.28*f*, 18.17
- Degenerate integrated structures, **I**:21.12
- Degree of circular polarization (DoCP), **I**:15.10
- Degree of coherence, **I**:2.37
- Degree of linear polarization (DoLP), **I**:15.10
- Degree of modulation, for electro-optic modulators, **V**:7.35–7.36
- Degree of polarization (DoP), **I**:12.14–12.15, 15.9
- Degree of polarization (DoP) surfaces and maps, **I**:14.31–14.32, 14.32*f*
- Dektak stylus profiler, **V**:46.2
- Delay-line technique, **V**:60.5
- Delta gun CRTs, **III**:22.4, 22.13
- Demultiplexers and demultiplexing, **V**:20.1
 - for networking, **V**:18.4, 18.10*f*, 18.10–18.11, 18.11*f*
 - in OTDM networks, **V**:20.7*f*, 20.7–20.8, 20.22, 20.23*f*
 - terahertz asymmetric optical, **V**:20.22
- Dense crown flint glass, **IV**:2.42*t*
- Dense flint glass, **IV**:2.43*t*
- Dense phosphate crown glass, **IV**:2.41*t*

- Dense wavelength division multiplexing (DWDM):
 and AOTFs, **V**:6.43, 6.44
 and optical fiber amplifiers, **V**:14.1
 and SOAs, **V**:19.25*f*, 19.25–19.27, 19.26*f*
- Density (of photographic films), **II**:29.6–29.8, 29.7*f*
- Density matrix, **IV**:16.21, 16.22
- Density (coherency) matrix, **I**:12.29–12.30, 14.41
- Density of states (DOS), **II**:19.9–19.10, 19.10*f*;
V:11.8, 11.9*f*
- Density-operator approach, to quantum theory of lasers, **II**:23.14–23.33
 derivation of Scully-Lamb master equation, **II**:23.17–23.22
 cavity losses, **II**:23.18
 laser master equation, **II**:23.19*f*, 23.19–23.20
 micromaser master equation, **II**:23.20–23.22
 photon statistics, **II**:23.22–23.27
 laser, **II**:23.22–23.26, 23.23*f*, 23.25*f*
 micromaser, **II**:23.26–23.27, 23.27*f*
 spectrum, **II**:23.28–23.33
 laser field, **II**:23.28–23.31, 23.30*f*
 micromaser field, **II**:23.31–23.33
 time evolution of the field in Jaynes-Cummings model, **II**:23.15*f*, 23.15–23.17
- Dephasing, **IV**:11.15, 14.12–14.14
- Depletion, of charge, **II**:25.6
- Depletion layer, of semiconductor detectors, **V**:60.5
- Depletion layer generation current, **II**:32.10–32.11, 32.11*f*
- Depletion region, of pin photodiodes, **V**:13.64
- Depolarization:
 defined, **I**:15.7
 and diagonal depolarizers, **I**:14.30
 Mueller matrices for, **I**:14.30–14.39
 depolarization index, **I**:14.32
 generators of, **I**:14.33–14.39, 14.36*f*, 14.37*f*
 and volume scattering, **I**:9.16–9.17, 9.17*f*
- Depolarization index, **I**:14.32
- Deposition:
 chemical vapor, **V**:25.26
 modified chemical vapor, **V**:25.2, 25.21, 25.26, 25.28
 outside vapor, **V**:25.2, 25.26
 physical vapor, **V**:61.6
 vapor axial, **V**:25.3, 25.26
- Deposition method, of manufacturing thin-films, **IV**:7.11–7.12
- Depth dependent sensitivity (OCT), **III**:18.13*f*, 18.13–18.14, 18.14*f*
- Depth of field, **I**:1.84, 17.37, 28.22–28.23
- Depth of focus (DOF), **I**:1.84, 17.37, 17.37*f*, 28.22; **II**:4.7–4.8, 4.8*f*
 in human eye, **III**:1.28–1.29
 of multilayer Laue lenses, **V**:42.17, 42.17*f*
 and pupil diameter, **III**:1.8, 1.29, 1.30*f*
- Depth ordering, **III**:13.4, 13.5, 13.10–13.12
- Depth perception, **III**:2.39–2.40
 age-related changes in, **III**:14.22
 stereodepth, **III**:13.11–13.12
- Depth range:
 with OFDI, **III**:18.16
 in SD-OCT, **III**:18.6–18.7
- Depth-of-amplitude modulation, **V**:7.22
- Depth-of-phase modulation, **V**:7.19
- Derivative matrices, **I**:1.73
- Derotation, of polygon scanners, **I**:30.35–30.36
- Derrington Krauskopf Lennie (DKL) space, **III**:11.32, 11.32*f*, 11.33
 defined, **III**:11.2
 equiluminant plane
 hue scaling, **III**:11.63
 multiple mechanisms in, **III**:11.57, 11.58
- Descemet's membrane, **III**:14.5, 16.4
- Design software, optical (*see* Optical design software)
- Designer blacks, **IV**:6.14
- DeSoto Black, **IV**:6.37*f*, 6.39
- Despace errors, for grazing incidence optics, **V**:45.6
- Destructive interference, **I**:2.7
- Destructive ports, of M-Z interferometers, **V**:21.34
- Detailed balance, condition of, **II**:23.23
- Detected point spread function (DPSF), **V**:44.14, 44.15*f*
- Detection and detectors:
 absolute detectors, **II**:34.27–34.30
 activated phosphor, **V**:60.7–60.8
 collision, **V**:23.7
 direct-conversion flat panel, **V**:61.4*f*, 61.6*t*, 61.6–61.7, 61.7*f*
 energy dispersive, **V**:62.2–62.4

- Detection and detectors (*Cont.*):
 for fiber optic systems, **V**:13.2–13.3,
 13.63–13.73
 avalanche photodiodes, **V**:13.71–13.73
 MSM detectors, **V**:13.73
 pin diodes, **V**:13.63–13.71, 13.64*f*, 13.66*f*,
 13.66*t*, 13.68*f*
 Schottky photodiodes, **V**:13.73
 heterodyne, **V**:3.34
 light detection and ranging (LIDAR),
V:3.38*f*, 3.38–3.39, 25.25
 light detectors, **III**:5.21, 5.22*t*
 for medical imaging, **V**:31.4, 31.4*f*, 61.2
 MSM photoconductive, **V**:13.63, 13.73
 in neutron optics, **V**:63.31–63.34
 photodetectors, **V**:9.8
 semiconductor, **V**:60.5
 semiconductor detectors, **IV**:5.61
 solid state, **V**:63.33–63.34
 spectroradiometry, **II**:38.8–38.10, 38.9*t*,
 38.10*t*
 standards for, **II**:38.12–38.13
 x-ray, **V**:31.4*f*
 controlled-drift, **V**:62.5
 cryogenic, **V**:60.8–60.9, 60.9*t*, 60.10*t*
 film, **V**:60.8, 60.9*t*, 60.10*t*
 ionization, **V**:60.3–60.7, 60.9*t*, 60.10*t*
 scintillation, **V**:60.7–60.8, 60.9*t*, 60.10*t*
 for x-ray imaging, **V**:61.1–61.8
 CCD detectors, **V**:61.7–61.8, 61.8*f*
 flat panel detectors, **V**:61.3–61.7, 61.4*f*,
 61.6*t*, 61.7*f*
 geometries for and classifications of,
V:61.1–61.3, 61.2*f*
- Detection surface or contour (color vision),
III:11.12
 defined, **III**:11.2
 and directions of color spaces, **III**:11.40*f*
 in equiluminant plane, **III**:11.44*f*
 and field adaptation, **III**:11.53–11.54, 11.54*f*
 in the L, M plane, **III**:11.40*f*, 11.41–11.42
- Detection tasks, **III**:2.15, 3.2
- Detection threshold:
 color vision
 chromatic discrimination near,
III:11.58–11.59
 two-stage model of, **III**:11.23*f*
 and discrimination thresholds, **III**:4.6–4.7
- Detective quantum efficiency (DQE), **II**:24.10,
 29.1, 29.23; **V**:61.2, 61.3, 61.5–61.7
- Detective time constant, **II**:24.10
- Detectivity:
 of film, **II**:29.23
 of infrared detector arrays, **II**:33.23–33.24
 normalized, **II**:38.9
 of photodetectors, **II**:25.4
 of thermal detectors, **II**:28.3, 28.3*f*
- Detrending, in x-ray mirror metrology,
V:46.6–46.7, 46.7*f*
- Detritus, organic:
 absorption by, **IV**:1.25–1.27, 1.25*t*, 1.26*f*
 in water, **IV**:1.14
- Detuning, **IV**:14.11*f*; **V**:13.29, 30.3
- Deuteranopes, **III**:10.16
- Deuterated triglycine sulfate (DTGS),
II:28.1, 28.7
- Deuterium lamps, **II**:15.13, 15.49, 15.53*f*, 34.31
- Developable film, **II**:29.5
- Development:
 of photographic film, **II**:29.5, 29.12, 29.13*f*,
 30.24
 in xerographic systems, **I**:34.5*f*–34.9*f*,
 34.5–34.10
- Development inhibitor anchimeric releasing
 (DIAR) color correction, **II**:30.18
- Development inhibitor-releasing (DIR)
 chemicals, **II**:30.3
- DeVries-Rose law, **III**:2.26
- Dewar container, **II**:24.10
- Diabetic retinopathy, **III**:14.1, 14.25–14.26
- Diagonal depolarizers, **I**:14.30
- Dialytes (lenses), **I**:17.25, 17.25*f*
- Diamagnification, **I**:1.23
- Diamond (crystal), **IV**:2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*,
 2.55*t*, 2.61*t*, 2.68*t*
- Diamond lattice structure, **IV**:5.6, 5.16, 5.17*t*
 of air spheres, **IV**:9.5
 of crystals and glasses, **IV**:2.68*t*
- Diamond tools, **II**:10.5, 10.8
- Diamond turning, **I**:22.15–22.18, 22.16*t*,
 22.17*f*, 22.18*f*; **II**:6.20, 10.1–10.14;
IV:3.12
 about, **II**:10.1–10.2
 advantages of, **II**:10.2–10.4
 of axicon optical element, **II**:10.3*f*
 machine tools for, **II**:10.6–10.8, 10.7*f*
 materials applicable to, **II**:10.4–10.5, 10.4*t*
 process of, **II**:10.2
 steps in, **II**:10.8–10.9
 traditional optical fabrication vs., **II**:10.6

- Diamond-turned optics:
 cleaning of, **II**:10.9
 metrology of, **II**:10.12*f*, 10.12–10.13, 10.13*f*
 surface finishing of, **II**:10.9*f*–10.11*f*, 10.9–10.11
- Diaphany, **III**:1.20
- Diattenuation, **I**:14.6, 15.7; **III**:18.1
 linear, **I**:14.8, 14.17
 Mueller matrices for, **I**:14.16*f*, 14.16–14.19
- Dichroic beam splitters, **IV**:7.65–7.66, 7.67*f*
- Dichroic polarizers, **I**:13.24–13.33, 13.26*f*, 13.27*f*
 coatings as, **I**:13.28
 measuring polarization of, **I**:13.33
 pyrolytic-graphite polarizers, **I**:13.28–13.29, 13.29*f*
 sheet polarizers, **I**:13.25–13.28
- Dichroism, **I**:15.19, 15.41, 31.17, 31.20, 31.21; **III**:18.1
 defined, **V**:25.2
 magnetic circular, **V**:55.7*f*, 55.7–55.9
- Dichromatic vision, **III**:10.13, 10.14, 10.16
- Dicke narrowing, **IV**:15.9
- Die fab, in wafer processing, **II**:17.24–17.25
- Die-attach materials, **II**:7.28
- Dielectric broadband reflectors, **IV**:7.39, 7.40*f*, 7.45*f*–7.47*f*, 7.45–7.47
- Dielectric color selective (dichroic) beam splitters, **IV**:7.65–7.66, 7.67*f*
- Dielectric compound parabolic collectors (DCPCs), **II**:39.15, 39.16, 39.16*f*
- Dielectric constant, **IV**:2.14, 2.17–2.18
 of crystals and glasses, **IV**:2.6*t*
 dispersion of, **IV**:8.22*f*
 for solids, **IV**:8.15
- Dielectric function, **IV**:4.26*f*, 5.8–5.9, 5.12
- Dielectric impermeability, **I**:21.9
- Dielectric multiple reflection cutoff filters, metal, **IV**:7.59–7.60
- Dielectric multiple reflection filters, metal, **IV**:7.111
- Dielectric potential, **IV**:9.15
- Dielectric properties, of liquid crystals, **V**:8.14–8.18, 8.15*f*, 8.18*f*
- Dielectric reflectors:
 interference filters with, **IV**:7.81
 metal, **IV**:7.81–7.82, 7.108, 7.109, 7.110*f*
- Dielectric solid-state laser gain media, **II**:16.32–16.34
- Dielectric tensor, of crystals and glasses, **IV**:2.17–2.18
- Dielectric total internal reflecting concentrator (DTIRC), **II**:39.17
- Differential amplifiers, **II**:27.11
- Differential detection, of optical disk data, **I**:35.22*f*, 35.22–35.23
- Differential geometry, of rays, **I**:1.19–1.21
- Differential group delay (DGD), **V**:21.17
- Differential measurement technique, **II**:22.13
- Differential-phase-shift-keyed (DPSK), **I**:21.30, 21.32; **V**:19.27, 21.33–21.34, 21.36, 21.36*t*, 21.37*t*
- Differential quadrature phase-shift-keying (DQPSK), **V**:21.34–21.36, 21.35*f*, 21.36*t*, 21.37*t*
- Differential quantum efficiency, **V**:13.9
- Differential reflectivity (DR), **IV**:18.1, 18.5, 18.6, 18.6*f*
- Differential scattering cross-sections (DSCs), **I**:7.8, 8.3, 8.4
- Differential transmission (DT), **IV**:18.1
- Differential transmission (DT) spectroscopy, **IV**:18.18–18.19
- Differential-interference contrast (DIC) microscopy, **I**:28.27, 28.39–28.41, 28.40*f*
- Diffraction, **I**:3.2–3.3, 3.3*f*
 acousto-optic, **V**:6.4, 6.9
 anisotropic, **V**:6.10
 anomalous, **I**:7.5, 7.6*f*
 and birefringent diffraction bandshapes, **V**:6.13, 6.14, 6.14*f*
- Bragg, **I**:11.9, 11.9*f*
 and acousto-optic modulators, **V**:6.4, 6.6, 6.7, 6.14
 and brightness of x-ray tube sources, **V**:54.16
 far and near, **V**:6.8–6.9, 6.12
 in neutron optics, **V**:63.23
 order of, **V**:20.14
 and phase matching equations, **V**:6.11
 of x-rays, **V**:42.2
- from circular apertures, **I**:3.6*f*, 3.6–3.7, 3.7*f*
 and Cornu's spiral, **I**:3.16–3.17
 by crystals, **V**:39.2–39.3
 of cylindrical wavefronts, **I**:3.13–3.21, 3.14*f*
 Cornu's spiral, **I**:3.16–3.19, 3.18*f*, 3.19*t*
 opaque strip construction, **I**:3.20–3.21, 3.21*f*
 from rectangular apertures, **I**:3.19–3.20, 3.20*f*
 from straight edge, **I**:3.14*f*, 3.14–3.16, 3.15*f*

- Diffraction (*Cont.*):
 definition of, **I**:3.6
 from disks, **I**:3.7–3.8
 and energy-dispersive x-ray fluorescence,
V:29.6–29.7, 29.8*f*–29.9*f*
 Fraunhofer, **I**:3.24*f*–3.26*f*, 3.24–3.28; **V**:42.2,
 63.25
 Airy diffraction as, **I**:28.17
 conducting screens for, **I**:3.33*f*
 and gratings, **I**:20.3
 Fresnel, **V**:27.2, 27.4, 40.9, 63.25
 Fresnel-Kirchhoff formula, **I**:3.21, 3.22, 3.32
 by gratings and monochromators,
V:38.1–38.3, 38.2*f*
 Green's function, **I**:3.22–3.23
 Huygens-Fresnel construction, **I**:3.4–3.13
 Babinet principle, **I**:3.9–3.11, 3.10*f*, 3.11*f*
 circular apertures and disks, light from,
I:3.6*f*–3.9*f*, 3.6–3.9
 Fresnel zones, **I**:3.4–3.6
 zone plates, **I**:3.11–3.13
 isotropic
 and acousto-optic interactions, **V**:6.9–6.10,
 6.10*f*, 6.13, 6.13*f*
 by deflectors, **V**:6.23–6.24, 6.24*f*
 and isotropic diffraction bandshape, **V**:6.13,
 6.13*f*
 kinematic theory of, **V**:63.5
 lysozyme, **V**:53.12*f*, 53.12–53.14, 53.13*f*
 mathematical theory of, **I**:3.21–3.29
 diffraction grating, **I**:3.28–3.29, 3.29*f*, 3.30*f*
 Fraunhofer diffraction, **I**:3.24–3.28
 Fresnel and Fraunhofer approximations,
I:3.23–3.24
 in multifocal optics, **III**:21.15–21.18
 near- and far-field, **V**:27.2–27.3, 27.3*f*
 in neutron optics, **V**:63.3, 63.25
 order of, **V**:18.6
 powder, **V**:53.14, 53.14*f*
 Rayleigh-Sommerfeld, **I**:3.9, 3.10, 3.23, 3.29
 and resolution of microscopes, **I**:28.19–28.22,
 28.20*f*, 28.21*f*
 and retinal image quality, **III**:1.12–1.14, 1.21
 scalar diffraction theory, **I**:23.10–23.13,
 23.11*t*, 23.12*f*, 23.13*f*
 single crystal, **V**:53.12*f*–53.13*f*, 53.12–53.14,
 53.13*f*
 by spheres, **I**:7.4
 stationary phase approximation for, **I**:3.29,
 3.31–3.32
- Diffraction (*Cont.*):
 vector, **I**:3.32*f*–3.37*f*, 3.32–3.37, 23.13–23.14,
 23.14*f*
 and wavelength-dispersive x-ray fluorescence,
V:29.2
 x-ray, **V**:28.5–28.6, 28.6*f*
 and x-ray optics, **V**:26.7, 26.8
 Diffraction efficiencies, of zone plates,
V:40.4–40.8, 40.5*f*, 40.7*f*, 40.8*f*
 Diffraction efficiency, **I**:23.9, 23.9*f*, 23.10*t*
 Diffraction focus, **II**:11.35
 Diffraction geometry, tilted, **V**:42.4
 Diffraction grating, **IV**:5.59–5.60
 Diffraction gratings, **I**:3.28–3.29, 3.29*f*, 3.30*f*,
 20.4
 arrayed waveguide, **I**:21.24
 beam-splitter, **I**:23.11, 23.12
 binary, **I**:23.13
 dispersive gratings vs., **I**:20.3–20.4
 interference, **I**:33.14
 multiple beam, **I**:2.28–2.29, 2.29*f*, 2.30*f*
 for PICs, **I**:21.19
 as polarizers, **I**:13.30–13.33
 Diffraction limit, **III**:4.1
 Diffraction limited imaging, **III**:17.1
 Diffraction patterns, three-dimensional,
I:28.19–28.22
 Diffraction theory, **III**:6.1
 Diffraction transfer factor (DTF), **II**:4.1, 4.4
 Diffraction-limited depth of focus, **I**:28.22
 Diffraction-limited lenses, **I**:17.37–17.39,
 17.41*f*–17.42*f*
 Diffraction-limited optics, **II**:4.4, 8.5, 8.6*f*
 Diffraction-related error, **II**:34.33
 Diffraction-type polarizers, **I**:13.30–13.33
 Diffractive contact lenses, **III**:14.28
 Diffractive imaging, **V**:41.9–41.10, 41.10*f*
 Diffractive optics, **I**:33.13, 35.28, 35.28*f*, 35.29;
II:8.15; **V**:26.8, 26.8*f*, 26.9*f*
 Diffractive scattering, **V**:64.2–64.3
 Diffractometers, **V**:28.1–28.3, 28.2*f*, 28.3*f*
 Diffuse attenuation:
 and Jerlov water types, **IV**:1.42–1.46,
 1.43*t*–1.45*t*, 1.44*f*, 1.45*f*
 in water, **IV**:1.13
 Diffuse attenuation coefficient, **IV**:1.12
 Diffuse density, **II**:29.7, 29.7*f*
 Diffuse Infrared Background Explorer, **IV**:6.29*f*
 Diffuse lighting, **II**:40.15
 Diffuse reflectance, **II**:35.4, 35.11, 35.12*f*, 35.13

- Diffuse scattering, **V**:63.4
- Diffuse transmittance, **II**:35.3, 35.9*f*
- Diffused homojunctions, LED, **II**:17.9–17.10, 17.10*f*, 17.11*f*
- Diffusely scattering surfaces (*see* Black surfaces)
- Diffusers, **I**:6.5, 6.5*f*, 6.7, 23.13; **III**:5.10*f*, 5.10–5.11
- with integrating cavities, **II**:39.26
 - perfect reflecting, **II**:37.8
 - perfect transmitting, **II**:37.8
 - plane, **II**:38.14, 38.14*f*, 38.15*f*
 - thin teflon, **II**:38.15*f*
- Diffusion:
- open-tube, **II**:17.24
 - photogenerated charge collection by, **II**:32.5
 - sealed-ampoule, **II**:17.23
 - semisealed-ampoule, **II**:17.24
 - in wafer processing, **II**:17.23–17.24
 - zinc, **II**:17.9–17.10, 17.10*f*
- Diffusion approximation, of radiative transfer, **I**:9.11–9.12
- Diffusion current, **II**:25.7, 26.8, 26.8*f*, 26.9, 32.10, 32.11*f*, 32.11*n*
- Digital, analog conversion to, **V**:20.4, 20.8, 20.8*f*
- Digital displays, in medical imaging, **V**:31.8*f*–31.10*f*, 31.8–31.9
- Digital holographic microscopy (DHM), **I**:28.42, 28.43
- Digital imaging:
- color representations in, **III**:10.35
 - displays for vision research
 - color cathode ray tubes (color CRTs), **III**:22.1–22.34
 - liquid crystal displays (LCDs), **III**:22.34–22.40
- Digital Libraries, **III**:24.7–24.8, 24.10
- Digital light processing (DLP), **I**:30.3, 30.60
- Digital micromirror devices (DMDs), **I**:22.23, 30.60–30.61
- Digital modulation, **V**:6.33
- Digital on-off-keying receivers, **V**:9.9–9.11
- Digital signal processing (DSP), **II**:27.14
- Digital tomosynthesis, **V**:31.7–31.8
- Digital transmission, in integrated optics, **I**:21.31–21.32
- Digital II SSA cameras, **II**:31.24–31.26
- Digital Video Broadcasters Measurement Group standards, **V**:15.4
- Digitized Green's function, **I**:7.15, 7.16
- Diluted magnetic semiconductors (DMSs), **IV**:5.45
- Dilution, in nonimaging optics, **II**:39.6, 39.40
- Diode lasers, external cavity, **II**:22.21–22.23, 22.22*f*
- Diode phototube, **II**:24.6
- Diode pumps, **V**:25.12–25.13
- Diode-pumped solid state lasers, **II**:22.20–22.21
- Diodes:
- laser, **V**:13.3–13.24
 - double heterostructure, **V**:13.3*f*, 13.3–13.8, 13.4*f*, 13.6*f*, 13.7*f*
 - noise characteristics of, **V**:13.18–13.24, 13.19*f*–13.21*f*
 - operating characteristics of, **V**:13.8–13.13, 13.10*f*
 - transient response of, **V**:13.13–13.18, 13.14*f*, 13.16*f*, 13.17*f*
 - light-emitting, **V**:13.1, 13.36–13.42, 13.38*f*
 - edge-emitting, **V**:13.40
 - operating characteristics of, **V**:13.40*f*, 13.40–13.42
 - and optical fibers, **V**:9.7
 - surface-emitting, **V**:13.38*f*, 13.38–13.40
 - and transmissive TFT LCDs, **V**:8.29–8.31, 8.30*f*
 - photo- (*see* Photodiodes)
 - pin, **V**:13.2, 13.63–13.71, 13.66*t*
 - dark current, **V**:13.69
 - geometry of, **V**:13.64*f*, 13.64–13.65
 - noise, **V**:13.70–13.71
 - sensitivity, **V**:13.65–13.66, 13.66*f*
 - speed, **V**:13.67–13.68
 - and unitraveling-carrier (UTC) photodiodes, **V**:13.68*f*, 13.68–13.69 (*see also* Photodiodes)
- Diopter (D), **III**:12.1, 12.4, 13.21, 25.1
- Dioptric errors of focus, **III**:1.13, 1.14, 1.28
- Dioptric lens systems, **V**:35.1
- Dioptric systems, **I**:1.9
- Diplopia (double vision), **III**:23.2
- Dipole active modes, infrared, **IV**:8.16–8.18, 8.17*f*
- Dipole approximation, discrete, **I**:7.16
- Dipole force, **IV**:20.8
- Dipole model of light, **I**:3.33–3.36, 3.37*f*
- Dipper function, **III**:11.59, 11.60*f*
- Dirac delta function, **I**:5.9, 5.12, 6.8
- Dirac equation, for one-electron atom, **I**:10.10
- Dirac series, **I**:30.8

- Direct detector integration (DDI), **II**:33.18, 33.19*t*, 33.20*f*
- Direct excitons, **IV**:5.29
- Direct injection (DI), **II**:33.18–33.21, 33.19*t*, 33.20*f*, 33.21*f*
- Direct interband absorption, in solids, **IV**:8.27–8.28
- Direct (vertical) interband absorption transitions, **IV**:5.22, 5.22*f*–5.23*f*, 5.22–5.23, 5.25*f*
- Direct interferometry, **II**:13.17–13.18
- Direct laser acceleration (DLA), **IV**:21.43
- Direct lighting, **II**:40.14
- Direct modulation, **V**:20.17*f*, 20.17–20.18
- Direct ophthalmoscope, **III**:12.5–12.6
- Direct overwrite (DOW), of optical disks, **I**:35.30, 35.30*f*
- Direct pulse generation, **IV**:18.4
- Direct readouts (DROs), from FPAs, **II**:33.15–33.18
- electronically scanned staring FPAs, **II**:33.16–33.17
- output circuits, **II**:33.18
- time delay integration scanning FPAs, **II**:33.17, 33.17*f*
- X-Y addressing and clock generation, **II**:33.16
- Direct semiconductors, **II**:17.4, 17.5*f*
- Direct-charge-injection silicon (DSI) FPAs, **II**:33.11*f*, 33.13
- Direct-conversion flat panel detectors, **V**:61.4*f*, 61.6*t*, 61.6–61.7, 61.7*f*
- Direct-indirect lighting, **II**:40.15, 40.46*f*, 40.47
- Direction, perception of, **III**:13.7–13.10
- corresponding retinal points, **III**:13.8
- horopter, **III**:13.8–13.9, 13.9*f*, 13.10*f*
- vertical horopter, **III**:13.9
- Direction discrimination, **III**:2.37
- Direction selectivity, **III**:2.14*n*
- Directional (term), **II**:35.5
- Directional conical reflectance, **II**:35.5*t*, 35.6*f*, 35.6*t*
- Directional couplers, for networking, **V**:18.2, 18.3, 18.3*f*, 18.9, 18.9*f*
- Directional emittance, **II**:35.7
- Directional hemispherical emittance, **II**:35.15
- Directional hemispherical reflectance, **II**:35.5*t*, 35.6*f*, 35.6*t*
- Directional spectral absorptance, **II**:35.7, 35.8*t*
- Directional spectral emittance, **II**:35.7
- Directional total absorptance, **II**:35.8*t*
- Directionality (isolation), of fiber-optic components, **V**:18.1
- Direct-reading autocollimators, **II**:12.12
- Direct-vision prisms, **I**:19.2
- Disability glare, **II**:40.9–40.10; **III**:14.4, 14.12, 23.2
- Discharged area development (DAD), **I**:34.4
- Discomfort, visual, **II**:40.9–40.12
- Discomfort glare, **II**:40.9–40.12, 40.11*t*; **III**:14.12, 23.2
- Disconjugate eye movements, **III**:13.1
- Discrete actuator deformable mirrors, **III**:15.1, 15.10, 15.10*f*
- Discrete dipole approximation, **I**:7.16
- Discrete electroabsorption modulators, **V**:13.59
- Discrete energy levels, **II**:16.4
- Discrete signals, incoherent processing of, **I**:11.17–11.20, 11.18*f*, 11.19*f*
- Discrimination contours (color vision), **III**:11.44*f*
- Discrimination experiments (*see* Pedestal experiments)
- Discrimination tasks, **III**:2.15, 3.2
- Discrimination thresholds, **III**:4.6–4.7
- Diseases of the eye:
- age-related, **III**:14.3
- age-related macular degeneration, **III**:14.1, 14.24–14.25
- cataract, **III**:14.24
- diabetic retinopathy, **III**:14.25–14.26
- glaucomas, **III**:14.26–14.27
- life-span environmental radiation damage, **III**:14.22–14.23
- visual impairment secondary to, **III**:14.4 (*see also specific diseases*)
- Disjunctive eye movements (*see* Vergence eye movements)
- Disk PZT transducers, **II**:22.17–22.18
- Disk rotation speed, of optical disks, **I**:35.5–35.6
- Disks, as aperture screens, **I**:3.7–3.11, 3.8*f*
- Dislocation reduction, **II**:18.2, 18.2*f*
- Dispenser cathodes, **V**:54.10
- Dispersion, **I**:20.1; **IV**:18.11
- angular, **V**:38.6
- Cauchy, **IV**:2.21
- chromatic
- and fiber optic communication links, **V**:15.9, 15.10
- in WDM networks, **V**:21.14–21.16, 21.15*f*, 21.16*f*

- Dispersion (*Cont.*):
 for crystals and glasses, **IV**:2.21–2.23,
 2.60*t*–2.66*t*, 2.66*t*–2.68*t*
 Drude, **IV**:2.21–2.22
 in fiber optic communication links,
V:15.9–15.11, 15.10*f*
 by gratings and monochromators,
V:38.6–38.7
 group velocity, **V**:11.18, 11.18*f*, 11.19*f*
 of light, **II**:38.8
 Maxwell-Helmholtz-Drude, **IV**:2.12, 2.21–2.22
 in multilayer reflectors, **IV**:7.40
 normal vs. anomalous, **IV**:4.4
 in optical fibers, **V**:9.5–9.7, 9.6*f*
 principle of, **IV**:2.23
 range of anomalous, **IV**:4.4
 room-temperature, **IV**:2.60*t*–2.68*t*
 in solids, **IV**:8.14–8.16, 8.22*f*
 theory of, **V**:63.5
 in WDM networks, **V**:21.20–21.26, 21.21*f*
 chromatic, **V**:21.14–21.16, 21.15*f*, 21.16*f*
 electronic solutions for, **V**:21.26, 21.27*f*
 fixed dispersion compensation,
V:21.22–21.23, 21.23*f*
 tunable dispersion compensation,
V:21.23–21.26, 21.24*f*–21.26*f*
 Zernike, **IV**:2.22
- Dispersion compensating fiber (DCF), **V**:15.10,
 21.22
- Dispersion length, of solitons, **V**:22.3
- Dispersion optimizing fiber, **V**:15.10
- Dispersion spectrometers, **IV**:5.59–5.60
- Dispersion-managed (DM) solitons,
V:22.12–22.15, 22.13*f*, 22.14*f*
- Dispersion-shifted fibers, **V**:9.7
- Dispersivity, **II**:30.9
- Dispersive prisms and gratings, **I**:20.1–20.15
 configurations of, **I**:20.4–20.15
 diffraction gratings vs., **I**:20.3–20.4
 Eagle configuration, **I**:20.7, 20.11*f*
 Ebert-Fastie configuration, **I**:20.8, 20.12*f*
 Littrow configuration, **I**:20.7*f*, 20.10
 Paschen-Runge configuration, **I**:20.7, 20.11*f*
 Pfund configuration, **I**:20.8*f*, 20.10, 20.13*f*
 in spectrometers, **I**:20.2–20.3, 20.3*f*
 in spectroradiometers, **I**:20.1, 20.2*f*, 20.14*t*
 Wadsworth configuration, **I**:20.5*f*, 20.8,
 20.12*f*
- Displacement current, **II**:26.7
- Displacement vectors, bulk modulators and,
V:7.18
- Displays (in general), **II**:40.1 (*see also specific displays, e.g.: Monolithic LED displays*)
 of cameras, **I**:25.6–25.7
 for medical imaging, **V**:31.8–31.9
 out-of-plane coupling and, **IV**:9.11–9.12
- Displays for vision research:
 color cathode ray tubes (color CRTs),
III:22.1–22.34
 colorimetric calibration of, **III**:22.20–22.34
 design and operation of, **III**:22.3–22.13
 operational characteristics of,
III:22.13–22.18
 setup for image display, **III**:22.18–22.19
 viewing environments, **III**:22.19–22.20
 liquid crystal displays (LCDs), **III**:22.34–22.40
 color LCDs, **III**:22.37–22.40
 monochrome, operational principles of,
III:22.34–22.37
- Disruptive movements (*see Vergence eye movements*)
- Dissipative force, **IV**:20.7
- Dissolved substances, in water, **IV**:1.13
- Distal stimuli (human vision), **III**:4.2
- Distance(s):
 in Gaussian lenses, **I**:1.51–1.53
 hyperfocal, **I**:1.85
 perception of, **III**:13.24
- Distance conflicts, in head mounted display systems, **III**:13.31–13.32
- Distance limits, of optical fibers, **V**:9.12–9.13
- Distance measurement (*see Length measurements*)
- Distorted object approach, in coherent x-ray optics, **V**:27.2–27.4
- Distortion(s), **I**:29.37
 barrel, **I**:1.91
 in fiber optic communication links,
V:15.5–15.6, 15.6*f*
 with head-mounted displays, **III**:25.7
 nonrectilinear, **I**:27.6
 of objectives, **I**:29.6
 in optical fibers, **V**:9.17
 pincushion, **I**:1.91
 pupil, **I**:1.78
 rectilinear, **I**:27.6, 27.13*f*, 27.14*f*
- Distortion plot, **II**:2.4, 2.4*f*
- Distortion tolerances, **II**:5.8
- Distortion-free focusing lenses, **I**:18.4

- Distracting glare, **II**:40.9
- Distributed backscattering, **II**:20.14
- Distributed Bragg reflector (DBR) lasers,
I:21.25, 21.30, 21.32, 21.37, 21.37*f*;
II:19.38, 19.40; **IV**:17.9, 17.10*f*, 17.11*f*;
V:9.8, 13.7, 13.28–13.29, 13.29*f*, 20.14
- Distributed feedback (DFB) lasers, **I**:21.25,
21.29, 21.30, 21.32, 21.38, 21.42; **II**:16.29,
19.36, 19.38; **IV**:9.6
in fiber optic systems, **V**:13.7, 13.30*f*,
13.30–13.32, 13.31*f*
optical fibers for, **V**:9.8
in OTDM communication networks,
V:20.14–20.15, 20.15*f*
- Distributed feedback (DFB) threshold, of fiber
optic systems, **V**:13.30–13.32
- Distributed grating surface-emitting lasers,
II:19.40*f*, 19.40–19.41
- Distributed index of refraction, **I**:24.1 [*see also*
Gradient index (GRIN) optics]
- Distributed Raman amplifiers (DRAs), **V**:21.44
- Distributed-index planar microlenses,
I:22.26–22.31, 22.27*f*–22.30*f*, 22.27*t*, 22.31*t*
- Distribution functions, of water, **IV**:1.6*t*, 1.12
- Distribution temperature, **II**:34.43–34.44, 37.7
- Disturbance of wavefront, **I**:3.5–3.6, 3.14*f*,
3.14–3.15
- Divalent rare-earth ions, **V**:2.11
- Divergence, **III**:13.21
defined, **III**:13.1
in first 6 weeks of life, **III**:13.21
- Divergent reflectors, **II**:39.38*f*, 39.38–39.40,
39.39*f*
- Divided voltage method, for transreflective
LCDs, **V**:8.35*f*, 8.35
- Division-of-amplitude photopolarimeter
(DOAP), **I**:16.15, 16.15*f*, 16.16
- Division-of-amplitude polarimeters, **I**:15.5–15.6
- Division-of-aperture polarimeters, **I**:15.5
- Division-of-wavefront photopolarimeter
(DOWP), **I**:16.14, 16.14*f*
- DKL color space, **III**:10.19
- D-log H curve (for photographic films),
II:29.8*f*, 29.8–29.10
- Documentation, for polymers, **IV**:3.5
- Domes, mounting of, **II**:6.11, 6.12*f*
- Donor-acceptor pair (DAP) transition, **IV**:5.71
- Dopant profiling, for LMA fibers, **V**:25.18
- Dopants, for fiber lasers, **V**:25.22–25.26, 25.23*t*
- Doped extrinsic silicon, **II**:33.7, 33.8*f*
- Doping, substrate, **II**:17.20
- Doppler broadening, **I**:10.2, 10.7, 31.23; **II**:16.5,
16.6, 16.6*f*, 16.9; **V**:3.14, 5.32, 56.5
- Doppler cooling, **IV**:20.13–20.15, 20.14*f*
- Doppler cooling limit, **IV**:20.15
- Doppler effect, **V**:2.3, 6.9
- Doppler LIDAR systems, **V**:3.38–3.39
- Doppler limit, **IV**:20.5
- Doppler linewidth (*see* Full width at half
maximum)
- Doppler OCT, **III**:18.18, 18.18*f*–18.20*f*,
18.18–18.19
- Doppler profiles, of CW lasers, **V**:5.32, 5.34
- Doppler shift, **I**:2.13, 5.23, 11.6, 31.30; **IV**:11.19,
20.4; **V**:10.7
- Doppler temperature, **IV**:20.5, 20.11, 20.15
- Doppler-dominated lineshapes, **V**:3.23
- Dot matrix, **III**:23.2
- Dot pitch, **III**:23.2
- Double atomic resonance, **IV**:22.5–22.9,
22.6*f*–22.8*f*
- Double Dove prisms, **I**:19.3*t*, 19.10, 19.10*f*
- Double heterojunction (DH) LEDs, **II**:17.13,
17.13*f*–17.15*f*
- Double heterostructure laser diodes, **V**:13.3*f*,
13.3–13.8, 13.4*f*, 13.6*f*, 13.7*f*
- Double heterostructure (DH) lasers, **II**:19.4,
19.5*f*, 19.7, 19.12–19.15, 19.18–19.19, 19.19*f*
- Double heterostructure *pin* photodiodes, **II**:26.13
- Double heterostructure waveguides, **V**:19.3*f*,
19.5
- Double ionization, strong field, **IV**:21.18–21.19,
21.19*f*
- Double monochromators, **II**:35.9, 38.15*f*
- Double optical gating, **II**:21.8
- Double phase conjugate mirrors, **IV**:12.7, 12.8*f*
- Double photonic resonance, **IV**:22.11–22.13,
22.12*f*
- Double refraction, in calcite, **I**:13.2–13.6,
13.4*t*–13.5*t*
- Double sampling, correlated, **II**:33.13
- Double-beam spectrophotometers, **I**:31.4–31.5,
31.5*f*; **II**:35.8–35.9
- Double-bounce Wolter mirrors, **V**:52.4
- Double-channel planar buried heterostructure
(DC-PBH) lasers, **II**:19.24, 19.25*f*, 19.34*f*
- Double-Gauss lenses, **I**:17.27–17.28, 17.28*f*, 27.2
- Double-lambda coupling, **IV**:14.24*f*
- Double-pass methods (*retinal image quality*),
III:1.22–1.23

- Double-pass monochromators, **I**:20.9*f*
 Double-pass objective optics, **I**:30.32*f*,
 30.32–30.33
 Double-pass photodetectors, **II**:26.4*f*
 Double-passed two-beam interferometers,
I:32.8
 Double-reflection error control, **I**:30.50*f*,
 30.50–30.51, 30.51*f*
 Double-sided Feynman diagrams,
IV:11.12–11.13, 11.13*f*
 Double-slit apertures, **I**:3.26–3.28, 3.27*f*, 3.28*f*
 Doublet lens, air-spaced, **II**:6.7
 Doublets, achromatic (lenses), **I**:17.22–17.25,
 17.23*f*–17.25*f*, 17.24*t*
 Doubly curved crystals (DCCs), **V**:29.6–29.7,
 29.8*f*
 Doubly resonant optical parametric oscillators
 (DR OPOs), **IV**:10.18
 Doubly resonant oscillators (DROs),
IV:17.2–17.4, 17.3*f*, 17.4*f*, 17.16–17.17
 Dove prisms, **I**:19.3*t*, 19.9, 19.9*f*, 19.10, 19.10*f*
 Downconversion, parametric, **II**:23.14
 Down-conversion, spontaneous parametric,
IV:23.10, 23.13
 Downdwelling average cosine, of water, **IV**:1.6*t*,
 1.7*f*, 1.12
 Downdwelling diffuse attenuation coefficient,
 for sea water, **IV**:1.44*t*–1.45*t*
 Downdwelling irradiance, of water, **IV**:1.5*t*,
 1.7*f*, 1.8
 Downdwelling irradiance diffuse attenuation
 coefficients, for sea water, **IV**:1.43*t*
 Downhill optimizer, **II**:3.17
 Downward scalar irradiance, of water, **IV**:1.5*t*,
 1.7*f*, 1.8
 Draft angle, **II**:39.10
Dragon (monochromator) systems, **V**:38.3
 Drag-wiping (cleaning), **II**:10.9
 Dressed states, **IV**:14.3
 Drift, **III**:1.44
 in CCDs, **II**:32.17
 frequency vs. time, **II**:22.2
 low offset, **II**:27.11
 photogenerated charge collection by, **II**:32.5
 thermocouple junctions as source of, **II**:27.6,
 27.6*f*
 Drift velocity, **IV**:21.7
 Drive circuitry, of LEDs, **V**:13.42
 Drive power, of NPM AOTFs, **V**:6.41
 Driving lasers, in attosecond optics,
II:21.4–21.6, 21.5*f*
 Drop-in assembly, **II**:6.6, 6.6*f*
 Droplet keratopathies, **III**:7.6, 7.7
 Drude approximation, **IV**:5.34
 Drude dispersion formula, for crystals and
 glasses, **IV**:2.21–2.22
 Drude model, **I**:21.10; **IV**:4.4–4.5, 8.15, 8.21,
 8.22*f*, 16.20
 Dry age-related macular degeneration,
III:14.1
 Dry eyes:
 with increasing age, **III**:14.5
 as problem with computer work, **III**:23.9
d-tensor, **IV**:10.11
 Dual attach FDDI nodes, **V**:23.3
 Dual beam detection, **II**:22.13
 Dual frequency effect, on LCDs, **V**:8.18
 Dual in-line octocouplers, **II**:17.32, 17.32*f*
 Dual magnification Cassegrain objective,
I:29.9–29.10
 Dual rotating retarder polarimeters, **I**:15.16,
 15.16*f*
 Dual-cavity PE-SROs, **IV**:17.18–17.20, 17.19*f*,
 17.20*f*
 Dual-cell-gap transreflective LCDs, **V**:8.32–8.33,
 8.33*f*
 Dual-Fock states, **IV**:23.8–23.9
 Ductility, of metals, **IV**:4.8, 4.70
 Dumet (alloy), **II**:40.29
 Duobinary formats, for WDM networks,
V:21.32*f*, 21.32–21.33, 21.36, 21.36*t*,
 21.37*t*
 Duplexers, for networking, **V**:18.4, 18.4*f*,
 18.10–18.11
 DURACON, **IV**:6.55–6.56
 Duty cycle, for input/output scanning, **I**:30.14
 Dye lasers, **II**:16.31–16.32, 16.32*f*, 20.15–20.16
 Dye-doped polymers, in electro-optic
 modulators, **V**:7.34
 Dye-forming reaction, in film development,
II:29.14
 Dyes:
 azomethine, **II**:30.10*f*, 30.11–30.12, 30.12*f*
 cyanine, **II**:30.13, 30.13*f*
 light-absorbing, **II**:30.7
 photographic, **II**:30.10*f*, 30.10–30.13, 30.12*f*
 yellow filter, **II**:30.4
 Dynamic ionization stabilization, **IV**:21.21
 Dynamic Jahn-Teller effect, **V**:2.9

- Dynamic range:
of fiber amplifiers, **V**:14.3
of fiber optic communication links,
V:15.5–15.6, 15.6*f*
of solid-state cameras, **I**:26.11, 26.14
of wideband Bragg cells, **V**:6.30
- Dynamic retardation, of modulators, **V**:7.20*f*,
7.20–7.22, 7.21*f*
- Dynamic scattering, **I**:9.7*f*, 9.7–9.8
- Dynamical theory of diffraction, for crystals,
V:39.2
- Dynodes, in photomultipliers, **II**:27.6–27.9,
27.7*f*
- Dyson interference microscopes, **I**:28.41, 28.42,
28.42*f*
- Dyson lenses, **I**:18.21, 18.22, 18.22*f*
- Dysphotopsia:
defined, **III**:21.1
with intraocular lenses, **III**:21.21
- Eagle configuration, of dispersive prisms,
I:20.7, 20.11*f*
- Ebanol C coating, **IV**:6.56
- Eberhard effects, **II**:30.3
- Ebert-Fastie configuration, **I**:20.8, 20.12*f*,
20.14*t*
- Eccentric field arrangement, of lenses, **I**:18.22
- Eccentric pupil design (*see* Z-system)
- Eccentricity, **I**:1.34, 15.10; **III**:2.3, 2.7, 2.24,
2.25*f*
- ECE (United Nations Economic Commission
for Europe), **II**:40.63–40.64
- Echo signal, **IV**:11.11–11.12, 11.12*f*
- Eclipse burn of the retina, **III**:7.3
- ECP-2200, **IV**:6.27, 6.29*f* (*see also* MH 2200
coating)
- Edge clearance, **III**:20.1
- Edge filters, nonpolarizing, **IV**:7.66, 7.67*f*
- Edge image, in diffraction-limited eye, **III**:1.13
- Edge lit backlight, **II**:40.47, 40.47*f*
- Edge rays, **II**:39.22, 39.38
- Edge response, OTF and, **I**:4.7
- Edge thickness, for contact lenses, **III**:20.6
- Edge-absorbing photodetectors, **II**:26.4*f*
- Edge-coupled pin waveguides, **V**:13.68
- Edge-defined film-fed growth (EFG) technique,
V:12.9, 12.10*f*
- Edge-emitting lasers (ELASERs), **II**:25.15;
V:13.3
- Edge-emitting LEDs (ELEDs), **II**:25.15;
V:13.36, 13.37, 13.40
- Edge-illuminated photodetectors, **II**:26.4*f*, 26.5
- Edging step (of optics fabrication), **II**:9.6
- EEGRAZE code, **V**:44.13
- Effective focal length (EFL):
of camera lenses, **I**:27.3*f*–27.5*f*, 27.7*f*–27.16*f*,
27.18*f*–27.22*f*, 27.25
of Gaussian lenses, **I**:1.48
of microlenses, **I**:22.10
- Effective group refractive index, **V**:13.13
- Effective index method, **I**:12.11, 12.12, 21.4
- Effective mass, **IV**:8.25–8.26, 8.26*f*
- Effective medium theories (EMTs), **I**:16.4, 16.9
- Effective power:
defined, **III**:20.1
for spectacle and contact lenses,
III:20.8–20.12, 20.9*t*, 20.10*f*, 20.11*f*
“Effective” troland, **III**:8.7
- Effective-medium representation, of volume
scattering, **I**:9.8
- Efficiency:
detective quantum, **V**:61.2, 61.3, 61.5–61.7
differential quantum, **V**:13.9
diffraction, **V**:40.4–40.8, 40.5*f*, 40.7*f*, 40.8*f*
external quantum, **V**:13.9
external slope, **V**:13.9
gain, **V**:14.6
internal quantum, **V**:13.9, 13.10
modulation, **V**:7.36
quantum, **V**:13.65, 13.66
quantum detection, **V**:60.7, 61.2–61.3
radiative quantum, **V**:14.7
- Efficiency factors, for scattering by particles,
I:7.5, 7.6*f*
- E-field-dependent electronic polarizability,
IV:19.5
- E-folding distance, **I**:7.13
- “Egg-crate” array, of potential wells, **IV**:20.32
- Ego-center, **III**:13.2, 13.7
- Ego-motion, **III**:13.8, 13.9
- Ehrenfest theorem, **IV**:20.6
- Eigenpolarization, **I**:15.7; **V**:7.13–7.14, 7.14*f*
- Eikonals, **I**:1.12–1.14
- Einstein (unit), **II**:34.11
- Einstein coefficient:
for spontaneous emission, **I**:31.3; **V**:2.13
for stimulated absorption, **I**:10.6
- Einstein Observatory, **V**:44.4, 44.10, 47.1, 47.5
- Einstein’s light quanta, **II**:23.6–23.9, 23.8*f*

- Einstein's particle hypothesis, **II**:23.7
- Einstein-Smoluchowski theory of scattering, **IV**:1.30
- Eisenburg and Pearson two-mirror, three reflection objective, **I**:29.25
- Elastic constants (crystals):
 cubic crystals, **IV**:2.44*t*–2.49*t*
 hexagonal crystals, **IV**:2.46*t*
 liquid crystals, **V**:8.22, 8.23*f*
 monoclinic crystals, **IV**:2.47*t*
 orthorhombic crystals, **IV**:2.46*t*
- Elastic properties, of crystals and glasses, **IV**:2.30–2.31, 2.31*t*
- Elastic scattering, **I**:7.3; **V**:63.5
- Elastic stiffness, of metals, **IV**:4.69, 4.69*t*
- Elastomeric mountings, **II**:6.4, 6.4*f*, 6.5, 6.12
- Elasto-optic coefficients, for crystals and glasses, **IV**:2.21
- Elasto-optic effect, **V**:6.5–6.6
- Electric dipole selection rules, **V**:56.3
- Electric field amplitude, in multilayer systems, **IV**:7.9–7.10
- Electric field reflectance, Fresnel expression for, **IV**:4.5
- Electric fields, **I**:2.3–2.4, 3.2, 3.3
- Electrical contact (light bulb), **II**:40.29*f*
- Electrical injection:
 in laser diodes, **V**:13.4, 13.5
 in VCSELs, **V**:13.45, 13.45*f*
- Electrical parasitics (laser), **II**:19.34*f*, 19.34–19.35
- Electrical substitution radiometers, **II**:34.27–34.29
- Electrical time domain multiplexed (ETDM) transmission, **V**:20.3, 20.25
- Electrical transfer function, with series inductance, **II**:26.13
- Electrically conductive black paint, **IV**:6.56
- Electric-field-modulated reflection spectroscopy, **IV**:5.66*t*, 5.67
- Electroabsorption, **V**:13.55–13.56, 13.56*f*
 applying fields in semiconductors, **V**:13.58
- Electroabsorption modulators (EAMs):
 in fiber optic systems, **V**:13.55–13.60, 13.56*f*, 13.57*f*
 in OTDM networks, **V**:20.18, 20.20, 20.20*f*
- Electrodeless discharge lamps, **II**:15.36, 15.44
- Electrodeless fluorescent lamps, **II**:40.36–40.37
- Electrodeless lamps, **II**:40.25*t*, 40.26*t*, 40.36–40.37
- Electrodeless sulfur lamps (ESLs), **II**:40.36–40.37
- Electrodeposited surfaces, **IV**:6.8*f*, 6.53–6.54, 6.54*f*, 6.55*f*
- Electroless nickel, etched, **IV**:6.5*f*, 6.6*f*, 6.49–6.50, 6.50*f*, 6.51*f*, 6.53*f*
- Electroluminescence, of LEDs, **V**:13.37
- Electroluminescent light sources, **II**:40.37–40.39, 40.38*f*, 40.38*t*, 40.39*f*
- Electromagnetic dipole model of light, **I**:3.33–3.36, 3.37*f*
- Electromagnetic radiation, **III**:23.2
- Electromagnetic spectrum, semiconductor interactions with, **IV**:5.3–5.6, 5.4*f*
- Electromagnetic theory, **I**:12.4
- Electromagnetically induced transparency (EIT), **IV**:14.1–14.36, 22.5–22.9, 22.7*f*, 22.8*f*
 coherence in two- and three-level atomic systems, **IV**:14.4*f*, 14.4–14.5
 and cw lasers, **IV**:14.16–14.18, 14.17*f*
 at few photon level, **IV**:14.32–14.33
 gain and lasing without inversion, **IV**:14.18–14.19
 as interference effect, **IV**:14.2–14.4
 manipulation of optical properties by, **IV**:14.10–14.15, 14.11*f*, 14.13*t*
 coupling laser power, **IV**:14.14–14.15
 dephasing and fluctuations in laser fields, **IV**:14.13
 dephasing in gas phase media, **IV**:14.12–14.13
 dephasing in solids, **IV**:14.13–14.14
 inhomogeneous broadening, **IV**:14.14
 and maximal atomic coherence, **IV**:14.28–14.32, 14.29*f*–14.32*f*
 nonlinear optical frequency conversion, **IV**:14.24*f*, 14.24–14.28, 14.27*f*
 physical concept of, **IV**:14.5–14.10, 14.6*f*, 14.8*f*, 14.9*f*
 pulse propagation effects, **IV**:14.20–14.22
 and pulsed lasers, **IV**:14.15–14.16, 14.16*f*
 and refraction index in dressed atoms, **IV**:14.19–14.20
 in solids, **IV**:14.33–14.36, 14.35*f*, 14.36*f*
 ultraslow light pulses, **IV**:14.22–14.23, 14.23*f*
- Electron(s), **V**:2.9, 54.12, 55.17, 56.2, 58.1
 lifetimes of, **I**:10.6
 in multielectron atoms, **I**:10.10–10.11
 in one-electron atoms, **I**:10.7–10.9, 10.8*f*, 10.9*f*
 and optical spectra, **I**:10.12–10.16, 10.13*f*–10.15*f*

- Electron(s) (*Cont.*):
 relativistic, above threshold ionization,
 IV:21.20, 21.21*f*
 strong field interactions with single,
 IV:21.5–21.10, 21.7*f*, 21.9*f*
- Electron acceleration, IV:21.39–21.42, 21.40*f*
- Electron beam steering, V:54.11
- Electron beams, strong field interactions with
 relativistic, IV:21.9–21.10
- Electron binding energies, V:36.3*t*–36.6*t*
- Electron bombardment (EB), II:31.23
- Electron current, II:26.7
- Electron excitation, WDXRF and, V:29.2
- Electron lenses, II:31.8, 31.8*f*
- Electron linacs, V:63.12
- Electron stochastic heating, IV:21.36
- Electron-beam lithography (EBL), V:40.8
- Electron-bombarded SSAs (EBSSAs),
 II:31.23–31.27
 digital cameras, II:31.24–31.26
 modulation transfer function and limiting
 resolution of, II:31.26–31.27
 proximity-focused, II:31.23*f*, 31.23–31.24,
 31.24*f*
- Electron-hole drops, IV:5.26*t*
- Electron-hole pairs, IV:14.34, 16.31
- Electronic holography, I:33.9–33.14,
 33.11*f*–33.13*f*
- Electronic imaging:
 human vision and (*see* Human vision and
 electronic imaging)
 in ophthalmoscopic methods, III:1.23
- Electronic structure, of atoms, I:10.12–10.16,
 10.13*f*–10.15*f*
- Electronically scanned staring FPAs,
 II:33.16–33.17
- Electronic-speckle pattern interferometry
 (ESPI), I:33.9–33.13, 33.11*f*–33.13*f*
- Electro-optic coefficients, for crystals and
 glasses, IV:2.21
- Electro-optic effect, V:7.6–7.16
 and eigenpolarization/phase velocity indices
 of refraction, V:7.13–7.16, 7.14*f*, 7.16*f*
 Jacobi method, V:7.11–7.13
 linear, V:7.7, 7.8*t*, 7.10
 and lithium niobate modulators,
 V:13.49–13.51, 13.51*f*
 quadratic (Kerr), V:7.9*t*–7.10*t*, 7.11
- Electro-optic holography (EOH), I:33.11–33.13,
 33.12*f*, 33.13*f*
- Electro-optic modulators, V:7.1–7.39
 applications for, V:7.36–7.39, 7.37*f*–7.38*f*
 bulk modulators, V:7.16–7.28
 amplitude modulation, V:7.22–7.24, 7.23*f*,
 7.24*f*
 frequency modulation, V:7.24–7.25, 7.25*f*
 phase modulation, V:7.18–7.20
 polarization modulation (dynamic
 retardation), V:7.20*f*, 7.20–7.22, 7.21*f*
 scanners, V:7.26*f*–7.28*f*, 7.26–7.28
 crystal optics and the index ellipsoid,
 V:7.3–7.7, 7.4*f*–7.6*f*, 7.8*f*–7.10*f*
 and electro-optic effect, V:7.6–7.16
 eigenpolarization/phase velocity indices of
 refraction, V:7.13–7.16, 7.14*f*, 7.16*f*
 Jacobi method, V:7.11–7.13
 linear, V:7.7, 7.8*t*, 7.10
 quadratic (Kerr), V:7.9*t*–7.10*t*, 7.11
 and Euler angles, V:7.39
 in fiber optic systems, V:13.61
 geometries, V:7.16–7.18, 7.17*f*
 light propagation in, V:7.3
 longitudinal, V:7.16, 7.17, 7.17*f*
 materials for, V:7.33–7.34
 in OTDM networks, V:20.18–20.19, 20.19*f*
 performance criteria for, V:7.34–7.36
 polymer modulators, V:13.55
 transverse modulators, V:7.16, 7.17, 7.17*f*
 traveling wave modulators, V:7.28–7.30, 7.29*f*
 waveguide modulators, V:7.30–7.32,
 7.31*f*–7.33*f*
- Electro-optic modulators (EOMs), II:22.14, 22.20
- Electro-optic sampling, V:7.36–7.37, 7.37*f*–7.38*f*
- Electro-optic (gradient) scanners, I:30.45–30.48,
 30.46*f*–30.48*f*
- Electro-optic tensors, I:21.10
- Electro-Optical Industries, Inc. (EOI), II:15.14,
 15.15*f*, 15.16*f*
- Electro-optical modulators, I:15.23
- Electroreflectance, IV:5.66*t*, 5.67
- Electrorefraction, V:13.2
- Electrorefractive modulators, V:13.61–13.62
- Electrorefractive photorefractive (ERPR) effect,
 IV:12.21
- Electrostriction, IV:16.18–16.19, 19.5
- Element wedge, of polymers, IV:3.10
- Elements (chemical), x-ray properties of,
 V:36.1–36.9, 36.3*t*–36.9*t*
- Ellipse blindness, III:7.7
- Ellipsoid, III:8.1

- Ellipsometers and ellipsometry, **I**:16.1–16.21, 16.10*f*–16.12*f*, 16.10–16.18; **IV**:5.5, 5.57, 5.63, 5.66*t*, 5.67–5.69, 5.68*f*, 5.69*f*
 about, **I**:16.2*f*, 16.2–16.3
 applications, **I**:16.21
 for azimuth measurements, **I**:16.16
 conventions, **I**:16.3*f*, 16.3–16.4
 defined, **I**:15.7
 four-detector photopolarimeters, **I**:16.14*f*–16.16*f*, 16.14–16.16
 generalized, **I**:16.19
 interferometric arrangements of, **I**:16.18
 Jones-matrix generalized, **I**:16.19
 modeling and inversion, **I**:16.4–16.9, 16.6*f*–16.8*f*, 16.10*f*
 Mueller-matrix generalized, **I**:16.19–16.21, 16.20*f*, 16.20*t*, 16.21*f*
 multiple-angle-of-incidence, **I**:16.3
 normal-incidence rotating-sample, **I**:16.18
 null, **I**:16.11, 16.12
 perpendicular-incidence, **I**:16.17–16.18, 16.18*f*
 photometric, **I**:16.12–16.14, 16.13*f*, 16.14*f*
 and polarimetry, **I**:15.30–15.32, 15.31*f*, 15.32*f*
 return-path, **I**:16.16–16.17, 16.17*f*
 rotating-analyzer, **I**:16.13, 16.13*f*, 16.14
 rotating-detector, **I**:16.14, 16.14*f*
 spectroscopic, **I**:16.3
 transmission, **I**:16.10
 variable-angle spectroscopic, **I**:16.3
- Ellipsometric angles, **I**:16.3
- Elliptical polarizers, **I**:14.10, 15.17–15.18
- Elliptical polynomials, **II**:11.21, 11.25–11.27, 11.26*t*–11.27*t*, 11.36*t*, 11.38*f*
- Elliptical reflectors, **V**:64.3, 64.4, 64.4*f*
- Elliptical retarders, Mueller matrices for, **I**:14.14
- Ellipticity, of polarization elements, **I**:15.10
- Elongation, of metals, **IV**:4.70, 4.70*t*
- Embedded polarizers, **IV**:7.70–7.71, 7.71*f*, 7.72*f*
Emiliana huxleyi, **IV**:1.15
- Emission:
 amplified spontaneous, **V**:9.13, 10.11, 14.3, 14.6
 cold-cathode field, **V**:54.10–54.11
 molecular, **V**:3.18, 3.20, 3.20*f*
 particle-induced x-ray, **V**:29.4
 self-amplified spontaneous, **V**:58.1
 spontaneous, **V**:2.13, 19.3, 20.1
 stimulated, **V**:20.2 (see Stimulated emission)
- Emission lasers, correlated, **II**:23.42–23.43
- Emission lines, K-shell and L-shell, **V**:36.3*t*–36.8*t*
- Emission linewidth (of radiation), **II**:16.4–16.7, 16.6*f*, 16.7*f*
- Emission-line broadening, **II**:16.4–16.7, 16.6*f*, 16.7*f*
- Emissivity:
 of blackbody cavity, **II**:15.7, 15.8*f*
 tungsten, **II**:40.28*f*
- Emittance, **II**:39.2
 and absorptance, **II**:35.8
 calculating, **II**:34.25–34.26
 for crystals and glasses, **IV**:2.20
 defined, **II**:35.7
 measurement of, **II**:35.14–35.16, 35.15*f*
 of metals, **IV**:4.6, 4.49, 4.49*f*, 4.50*t*, 4.51*t*
 and surface coatings, **IV**:6.19, 6.20*t*
- Emitted photon wavelength, **II**:17.4–17.5
- Emitters:
 AlGaAs, **II**:17.32
 blue, **II**:17.18, 17.19
 GaAsP, **II**:17.32
- Emmetropia, **III**:1.6, 1.32, 12.4, 16.4, 16.5*f*
 age-related, **III**:14.11
 defined, **III**:6.1, 12.1, 13.2
 and focus of collimated light, **III**:12.3*f*
 and Maxwellian viewing, **III**:5.4–5.5
- Emmetropization, **III**:1.7
- Emotional characteristics in electronic imaging (see Aesthetic and emotional characteristics)
- Empirical horopter, **III**:13.8
- Empty field myopia, **III**:1.33
- Empty magnification limit, **I**:28.17
- Emulsions, photographic, **II**:24.100, 24.101*f*, 29.4, 30.7
- Enclosed arcs, **II**:15.24, 15.28–15.47
 high-pressure, **II**:15.24, 15.28–15.34
 capillary mercury-arc lamps, **II**:15.30–15.31, 15.31*f*
 compact-source arcs, **II**:15.31–15.34, 15.32*f*–15.35*f*
 Lucalox lamps, **II**:15.30, 15.31*f*
 mercury arcs, **II**:15.29, 15.30*f*
 multivapor arcs, **II**:15.29, 15.31*f*
 Uviarc, **II**:15.28–15.29, 15.29*f*, 15.30*f*
- low-pressure, **II**:15.35–15.47
 black-light fluorescent lamps, **II**:15.35, 15.36*t*
 electrodeless discharge lamps, **II**:15.36, 15.44

- Enclosed arcs, low-pressure (*Cont.*):
 germicidal lamps, **II**:15.35
 hollow cathode lamps, **II**:15.35,
 15.37–15.43*t*, 15.44*f*
 Pluecker spectrum tubes, **II**:15.47, 15.47*f*,
 15.47*t*
 spectral lamps, **II**:15.44, 15.45, 15.45*f*,
 15.46*f*, 15.46*t*
 Sterilamps, **II**:15.35, 15.36*f*
- End loss, **II**:19.6
- Endlessly single-mode photonic crystal fiber
 (ESM-PCF), **V**:11.12, 11.13, 11.21, 11.21*f*
- Endoscopic cameras, **I**:25.21, 25.21*f*
- Endothelium:
 corneal, **III**:16.4
 defined, **III**:16.1
- End-pumped schemes, for fiber lasers,
V:25.9–25.10, 25.28
- Energy(-ies):
 atomic, **V**:2.2–2.5, 2.4*f*
 Auger, **V**:36.3*t*, 36.9*t*
 band gap, **V**:19.2
 conservation of, **V**:6.9
 electron binding, **V**:36.3*t*–36.6*t*
 equipartition, **V**:22.6
 Fermi level, **V**:30.1
 filtering of, **V**:53.10
 flow of, in solids, **IV**:8.7–8.8
 Landau levels of, **IV**:5.40, 5.42*f*
 levels of, **II**:16.4, 16.7
 luminous, **II**:37.4*t*, 37.6
 measurement of, **II**:34.32
 nomenclature for, **II**:36.4, 36.5
 photon, **V**:36.7*t*–36.8*t*
 radiant, **II**:34.7, 37.4*t*, 37.6
 units of, **II**:34.5–34.6
 (*see also specific sources, e.g.: Fiber lasers*)
- Energy band structure, **II**:17.3*f*–17.5*f*, 17.3–17.6
- Energy bandgap, **II**:25.3, 25.3*f*
- Energy bands:
 magnetic field effects on, **IV**:5.40
 for solids, **IV**:8.25–8.27, 8.26*f*
- Energy dispersive detectors (EDS), **V**:62.2–62.4
- Energy walk-off angle, **IV**:8.9
- Energy-dispersive x-ray fluorescence (EDXRF),
V:29.3–29.11
 with doubly curved crystal diffraction,
V:29.6–29.7, 29.8*f*–29.9*f*
 monocapillary micro-XRF, **V**:29.4
 polycapillary micro-XRF, **V**:29.4–29.6,
 29.5*f*, 29.6*f*
 ultrahigh resolution, **V**:29.9*f*–29.11*f*,
 29.9–29.11
- Energy-time uncertainty principle, **IV**:23.4
- Engineering moduli, for crystals and glasses,
IV:2.37
- English units, and SI units, **II**:37.7, 37.7*t*
- Enhanced backscatter (EBS), **I**:6.5*f*, 6.5–6.7
- Enhanced Martin Black, **IV**:6.46, 6.47
- Enhanced refraction, **IV**:14.20
- Entanglement, quantum (*see* Quantum
 entanglement, in optical interferometry)
- Entrance damage, laser-induced, **IV**:19.3
- Entrance pupil, **I**:1.76, 17.8, 18.4–18.6, 29.37;
II:34.18, 34.19*f*
 and correction for SCE-1, **III**:9.4
 defined, **III**:13.2
 in Maxwellian viewing, **III**:5.7
 in reflectometry, **III**:8.6–8.7
 and retinal illuminance, **III**:1.12
 and retinal irradiance, **III**:2.4
 and stimulus specification, **III**:4.2–4.3, 4.3*f*
- Entrance pupil distance (ENP) (camera lenses),
I:27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*,
 27.25
- Entrance slits, of gratings and monochromators,
V:38.7
- Entrance window, **I**:17.9; **II**:34.19, 34.19*f*
- Environmental control and correction, for
 adaptive optics, **V**:50.5, 50.6
- Environmental degradation, of black surfaces,
IV:6.16–6.18
- Environmental specifications, optical, **II**:4.10
- Environmentally responsible glass, **IV**:2.29–2.30
- Epicotyl, **III**:8.26, 8.27*f*
- Epi-illumination, in microscopes, **I**:28.7*f*,
 28.7–28.9, 28.8*f*
- Epi-LASIK (epithelial laser in situ
 keratomileusis), **III**:16.12, 16.13, 16.13*f*
- Epitaxial growth, **II**:17.8, 17.21; **V**:13.2
- Epitaxial technology (for LEDs), **II**:17.21–17.23
- Epitaxy, **I**:21.17–21.20
- Epithelial layer (cornea), **III**:14.5
- Epithelium:
 corneal, **III**:16.4
 defined, **III**:16.1
- Epithermal neutrons, **V**:63.18
- Eper Laser Black, **IV**:6.56
- Epoxy, in indicator lamps, **II**:17.29

- Equatorial divergence, x-ray diffraction and, **V**:28.1
- Equatorial plane, **III**:19.1
- Equilibrium colors:
 additivity with, **III**:11.66
 and unique hues, **III**:11.63–11.66
 white as, **III**:11.63
- Equilibrium level (accommodation), **III**:1.33
- Equiluminance, difficulty in producing, **III**:11.80
- Equiluminant plane:
 detection and discrimination contours in, **III**:11.44*f*
 of DKL space
 hue scaling, **III**:11.57
 multiple mechanisms in, **III**:11.57, 11.58
- Equipartition energy, **V**:22.6
- Equivalent lenses, **I**:17.20
- Equivalent neutral density (END), **II**:29.15
- Equivalent noise input (ENI), **II**:24.11
- Equivalent particles, in volume scattering, **I**:9.6
- Equivalent reflectance, retinal, **III**:1.11
- “Equivalent” troland, **III**:9.2
- Equivalent veiling luminance (EVL), **II**:40.10; **III**:1.20
- Equivalent-sphere correction, **III**:1.6
- Erasing, in xerographic systems, **I**:34.10
- Erbium-doped fiber amplifiers (EDFAs):
 energy levels, **V**:14.4
 fast power transients, **V**:21.39*f*, 21.39–21.41, 21.40*f*
 gain flattening, **V**:14.6–14.7, 21.38–21.39, 21.39*f*
 gain formation, **V**:14.4–14.5, 14.5*f*
 gain peaking, **V**:21.38, 21.38*f*
 noise, **V**:14.6
 pump wavelength options, **V**:14.5–14.6
 semiconductor amplifiers vs., **V**:9.13, 9.14, 14.1, 14.2*t*
 static gain dynamic and channel power equalization, **V**:21.41, 21.41*f*–21.42*f*
 in WDM networks, **V**:21.2*f*, 21.2–21.3
- Erbium-doped fibers, **V**:25.23*t*, 25.24, 25.32, 25.33
- Erbium-doped yttrium aluminum garnet (Er:YAG) lasers, **V**:12.3*t*, 12.6, 12.13
- Erbium/ytterbium-doped fiber amplifiers (EYDFAs), **V**:14.2, 14.2*t*, 14.7–14.8
- Ergonomics, **III**:23.2
- Error functions, in optical design software, **II**:3.17, 3.19–3.20
- Error types, in absolute measurements, **II**:34.21–34.23
- Erythema, **III**:7.1
- ESCON (Enterprise System Connection) standard, **V**:23.1–23.2, 23.2*f*
- ESO telescope, **V**:5.35
- Estimation tasks, **III**:2.15–2.16
- Etch and regrowth fabrication, of electroabsorption modulators, **V**:13.59–13.60
- Etched electroless nickel surface, **IV**:6.5*f*, 6.6*f*, 6.49–6.50, 6.50*f*, 6.51*f*, 6.53*f*
- Etching:
 for PICs, **I**:21.18–21.19
 surface, **IV**:6.15
- Étendue, **I**:1.22, 1.81, 13.7; **II**:34.15
 defined, **II**:40.42
 geometrical, **II**:38.8
 in nonimaging optics, **II**:39.2, 39.3, 39.4*f*, 39.5
 and source coupling, **II**:40.41–40.42
- Étendue loss, **II**:39.6
- Ethernet standard, **V**:23.7
- Etiolated plant tissues, **III**:8.26–8.28
- Euler angles, **V**:7.13, 7.39
- Euler equations, **I**:1.20
- Euler’s constant, **V**:56.9
- Euphotic zone, of water, **IV**:1.46, 1.46*t*
- European Synchrotron Radiation Facility (ESRF), **V**:37.5, 37.8, 50.5, 50.6
- European X-Ray Free-Electron Laser (XFEL), **V**:58.1
- Evaluation function (of optical software), **II**:3.8–3.16
 of aberrations, **II**:3.9–3.11
 paraxial ray-trace, **II**:3.8–3.9, 3.9*f*
 ray-trace, **II**:3.11–3.13, 3.12*f*
 by spot-diagram analysis, **II**:3.13–3.16
- Evanescence coupling, **V**:25.11
- Evanescence wave spectroscopy (EWS), **V**:12.13
- Evanescence-wave coupled pin waveguides, **V**:13.68
- Evaporated spacers, **IV**:7.83, 7.84*f*, 7.85*f*
- Evaporation method, of manufacturing thin-films, **IV**:7.11
- Event-driven programs (optical software), **II**:3.7
- Exact rays (term), **II**:3.3, 3.11–3.12
- Excimer lasers, **II**:16.30–16.31, 16.31*f*; **III**:16.15–16.16

- Excimers, in fluorescent lamps, **II**:40.31
- Excitance, total integrated, **IV**:2.19
- Excitation(s):
- chirped pulse, **IV**:11.25–11.26
 - continuum, **IV**:18.20
 - electron, **V**:56.2
 - excitonic, **IV**:18.19–18.20
 - photoexcitation, **IV**:5.70*f*
 - single-particle, **IV**:5.81, 5.82*f*
- Excitation spectroscopy, **V**:2.15*f*, 2.21
- Excited state, of azomethine dyes, **II**:30.11–30.12, 30.12*f*
- Excited state absorption (ESA), **IV**:13.5, 16.19; **V**:14.6
- Excited state collisions, **IV**:20.29
- Exciton(s):
- free, luminescence, **IV**:5.72, 5.73*f*
 - in semiconductors, **IV**:5.25–5.29, 5.26*t*, 5.27*f*–5.28*f*, 5.46
 - and solids, **IV**:8.31–8.32
- Exciton gases, **IV**:5.26*t*
- Exciton recombination, **II**:17.6
- Exciton Rydberg, **IV**:8.31
- Excitonic effects, in integrated optics, **I**:21.11
- Excitonic excitations, **IV**:18.19–18.20
- Excitonic magneto-optical effects, **IV**:5.46, 5.47*f*
- Excitonic molecules, **IV**:5.26*t*
- Exhaustive characterization methods (ECM) (color CRTs), **III**:22.21–22.23
- interpolation, **III**:22.22–22.23
 - inverses, **III**:22.23
 - out-of-gamut colors, **III**:22.23
 - sampling, **III**:22.21–22.22
- Exit damage, laser-induced, **IV**:19.3
- Exit pupil, **I**:1.76, 17.8, 18.6, 29.37; **II**:34.18, 34.19*f*; **III**:2.3
- measurement of power at, **III**:5.17
 - in reflectometry, **III**:8.6–8.7, 8.7*f*
 - and retinal illuminance, **III**:1.12
 - and Stiles-Crawford effects, **III**:1.11
- Exit pupil distance (EXP), of camera lenses, **I**:27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
- Exit slits, of gratings and monochromators, **V**:38.7
- Exit window (lens), **I**:17.9
- Exitance, **II**:39.2
- defined, **II**:34.8
 - luminous, **II**:37.4*t*, 37.5, 37.5*f*
 - radiant, **II**:15.4–15.6, 15.5*t*, 15.6*f*, 37.4*t*, 37.5, 37.5*f*
- Exiting beam, of polarimeters, **I**:15.4
- Expansion coefficients, in atmospheric optics, **V**:4.20–4.22, 4.21*t*, 4.22*f*–4.23*f*
- Experimental conditions, **III**:3.2
- Exposure:
- luminous, **II**:37.4*t*, 37.6
 - of photographic films, **II**:29.5–29.6
 - radiant, **II**:37.4*t*, 37.6
 - in xerographic systems, **I**:34.3
- Exposure limits (ELs) (radiation), **III**:7.9–7.11
- exceeding, **III**:7.11
 - for infrared, **III**:7.10–7.11
 - for laser light, **III**:7.12
 - for ultraviolet, **III**:7.9–7.10
 - for visible light, **III**:7.10
- Extended baffle shields, **II**:7.9*f*, 7.9–7.10, 7.10*f*
- Extended boundary condition method (EBCM), **I**:7.15
- Extended Gordon-Haus effect, **V**:22.11
- Extended objects, images of, **I**:1.27
- Extended sources, interference by, **I**:2.20
- Extended wavelength photodetectors, **II**:25.10, 25.10*t*
- Extended x-ray absorption fine structure (EXAFS), **V**:30.2, 30.4
- Exterior lighting, **II**:40.61–40.62, 40.63*t*
- External cavity diode lasers (ECDLs), **II**:22.21–22.23, 22.22*f*
- External circuits, noise from, **V**:13.70
- External limiting membrane (ELM), **III**:8.8, 8.9, 8.10*f*, 8.12
- External mirrors, **V**:13.32–13.33, 13.33*f*
- External modulation, in OTDM networks, **V**:20.18–20.20, 20.19*f*, 20.20*f*
- External optical feedback, from laser diodes, **V**:13.21–13.24
- External pulse compression, **IV**:18.4, 18.11–18.12
- External quantum efficiency, **V**:13.9
- External self-action, **IV**:16.25
- External slope efficiency, **V**:13.9
- Extinction coefficient, of metals, **IV**:4.3, 4.11, 4.12*t*–4.19*t*, 4.20*f*–4.26*f*
- Extinction cross section, **I**:7.5, 7.8
- Extinction paradox, **I**:7.8
- Extinction ratio, **V**:7.35, 15.13
- about, **I**:12.14–12.15, 12.17
 - of polarizers, **I**:12.21–12.24, 12.22*f*, 12.23*f*, 14.6

- Extracapsular cataract extraction (ECCE),
III:12.14, 21.4
- Extraocular muscles, **III**:1.42, 12.8, 12.16
- Extraretinal cues:
 defined, **III**:13.2
 for eye movements, **III**:13.7
 in space perception, **III**:13.3
- Extreme environments, black surface
 degradation in, **IV**:6.18
- Extreme infrared (IR) light, **II**:25.2
- Extreme ultraviolet explorer (EUVE), **IV**:6.21;
V:44.4, 44.5
- Extreme ultraviolet (EUV) lasers, **V**:58.2–58.4,
 58.3*f*
- Extreme ultraviolet (XUV) light:
 bandpass filters for, **IV**:7.94–7.96, 7.95*f*–7.96*f*
 interference polarizers for, **IV**:7.73, 7.76*f*–7.77*f*
 multilayer reflectors for, **IV**:7.42–7.43, 7.53
- Extreme ultraviolet lithography (EUV-L),
V:34.1–34.6
 and EUV-interferometric lithography,
V:34.4–34.5, 34.5*f*
 limitations of, **V**:34.5*f*, 34.5–34.6, 34.6*f*
 and multifoil optics, **V**:48.1
 and multilayers, **V**:41.5, 41.7, 41.8
 in semiconductor industry, **V**:34.1–34.2, 34.2*t*
 technology for, **V**:34.2–34.5, 34.3*f*, 34.4*f*
- Extreme ultraviolet (EUV) region,
 Schwarzschild objective for, **V**:51.3
- Extreme ultraviolet-interferometric lithography
 (EUV-IL), **V**:34.4–34.5, 34.5*f*
- Extreme wide-angle lenses, **I**:27.6, 27.13*f*, 27.14*f*
- Extrinsic Fabry-Perot interferometric (EFPI)
 sensors, **V**:24.2*f*, 24.2–24.4, 24.3*f*
- Extrinsic optical properties:
 of semiconductors, **IV**:5.11
 of solids, **IV**:8.3
- Extrinsic photoconductors, **II**:25.5, 25.5*f*
- Extrinsic photodetectors, **II**:24.7, 24.7*f*
- Extrinsic semiconductor transition, **II**:24.11
- Extrusion, of fiber lasers, **V**:25.27
- Exudative (wet) age-related macular
 degeneration, **III**:14.1, 14.24–14.35
- Eye alignment:
 binocular parallax, **III**:13.22
 cross-coupling and direction/distance of
 gaze, **III**:13.23*f*, 13.23–13.24
 defocus and efforts to clear vision,
III:13.22–13.23
 intrinsic stimuli to vergence, **III**:13.21–13.22
- Eye alignment (*Cont.*):
 magnification induced errors of,
III:13.25–13.27
 in Maxwellian viewing, **III**:5.8
 perceived distance, **III**:13.24
 prism induced errors of, **III**:13.27–13.29
 zone of clear and single binocular vision,
III:13.24–13.25, 13.25*f*
- Eye degradations, of optical fiber receivers,
V:9.11
- Eye loupes, **I**:17.9–17.10
- Eye movements, **III**:1.42–1.45, 13.19–13.20
 analysis of, **III**:24.7
 characteristics of, **III**:1.43–1.44
 in contrast detection, **III**:2.21
 coordination of
 binocular parallax, **III**:13.22
 cross-coupling and direction/distance of
 gaze, **III**:13.23*f*, 13.23–13.24
 defocus and efforts to clear vision,
III:13.22–13.23
 intrinsic stimuli to vergence, **III**:13.21–13.22
 perceived distance, **III**:13.24
 zone of clear and single binocular vision,
III:13.24–13.25, 13.25*f*
 defined, **III**:13.2
 extraretinal information for, **III**:13.7
 gaze control, **III**:13.29–13.30
 and optical flow, **III**:2.39
 stability of fixation, **III**:1.44, 1.45*f*
 types of, **III**:13.19–13.20
- Eye openings, optical fiber receivers and, **V**:9.11
- Eye protectors, for laser hazards, **III**:7.14
- Eye relief, **III**:1.7
- Eye relief (ER), **I**:18.8–18.10, 18.9*f*, 18.10*f*
- Eye safety, **V**:15.1
- Eye space, of afocal lenses, **I**:18.4
- Eye space pupil diameter, **I**:18.6
- Eye tracking, in cameras, **I**:25.14*f*, 25.14–25.15,
 25.15*f*
- Eyepieces, in afocal systems, **I**:18.7
- Eyes:
 bovine lenses, **III**:19.8–19.11, 19.11*f*
 cat lenses, **III**:19.9
 fish lenses, **III**:19.6, 19.6*f*
 gibbon lenses, **III**:19.14
 guinea pig lenses, **III**:19.8
 of invertebrate aquatic worms, **III**:9.14
 model, **III**:21.13
 octopus lenses, **III**:19.7, 19.7*f*

- Eyes (*Cont.*):
 pig lenses, **III**:19.11, 19.12*f*
 porcine lenses, **III**:19.11, 19.12*f*
 primate eye lens, **III**:19.13*f*, 19.14, 19.14*f*
 rabbit lenses, **III**:19.8
 rat lenses, **III**:19.7–19.8
 (*see also* Human eye)
- Eyesight, defined, **III**:23.2
- Eyestrain (asthenopia), **III**:23.2
- F line, **III**:25.1
- Fabrication, optical, **II**:9.3–9.9
 about, **II**:9.3
 aspherical, **II**:9.7*f*, 9.7–9.8
 crystalline, **II**:9.8
 by diamond turning (*see* Diamond turning)
 diamond turning vs. traditional, **II**:10.6
 guide to methods of, **II**:10.3*t*
 material formation for, **II**:9.3–9.4
 methods of, **II**:10.3*t*
 plano, **II**:9.7
 spherical, **II**:9.4–9.6
- Fabry-Perot cavities, **V**:13.65, 19.19, 20.21, 20.21*f*, 24.2
- Fabry-Perot etalon (cavity), **I**:2.33, 2.33*f*
- Fabry-Perot filters, **V**:13.33, 22.8
- Fabry-Perot interference filters, **IV**:7.78–7.82, 7.79*f*, 7.80*f*, 7.92–7.94, 7.93*f*, 7.96
- Fabry-Perot interferometers, **I**:32.4–32.7, 32.7*f*, 32.14; **II**:16.19*f*; **IV**:7.13, 7.39, 7.39*f*, 7.40, 7.89; **V**:3.36, 24.2*f*–24.4*f*, 24.2–24.5
 in dynamic wave meters, **I**:32.17
 in gravitational wave interferometers, **I**:32.21, 32.21*f*
 as heterodyne interferometers, **I**:32.10
 and multiple beam interference, **I**:2.33–2.36, 2.34*f*, 2.35*f*
 and wire-grid polarizers, **I**:13.31
- Fabry-Perot lasers, **V**:9.11, 13.12–13.13, 13.29, 15.11, 15.14
- Fabry-Perot resonators, **IV**:22.11, 22.12*f*
- Fabry-Perot semiconductor lasers, **V**:20.13–20.14, 20.14*f*
- Faceted reflectors, **II**:39.10*f*, 39.39*f*, 39.39–39.41, 39.40*f*
- Failures per 10⁹ hours (FITS), **II**:17.25
- False polarization, **I**:15.38
- False-colored infrared film, **II**:30.22
- Family momentum, **IV**:20.38
- Fankuchen-cut crystals, **V**:63.25
- Fano interferences, **IV**:14.2
- Far Bragg diffraction acousto-optic interaction and, **V**:6.8
- Far field, **V**:4.10
- Far point, **III**:12.1, 12.8, 16.5
- Far ultraviolet radiation, **II**:15.12, 15.13
- Faraday cages, **I**:34.8, 34.8*f*
- Faraday effect, **IV**:5.50–5.51
- Faraday rotation:
 free-carrier, **IV**:5.50–5.51
 interband, **IV**:5.44–5.45, 5.45*f*
- Faraday rotators, **I**:28.46; **V**:18.7, 18.7*f*, 18.10
- Faraday shutters, **I**:25.22
- Faraday's law, **IV**:2.6
- Far-field diffraction, **V**:27.2–27.3, 27.3*f*
- Far-infrared (FIR) lasers, **IV**:5.48
- Far-infrared (FIR) radiation, **II**:24.3, 25.2
 and black surfaces, **IV**:6.21, 6.26–6.34, 6.28*f*–6.34*f*
 Ames 24E and 24E2, **IV**:6.26*f*, 6.28*f*, 6.34
 Cornell Black, **IV**:6.26*f*, 6.27
 Infrablack, **IV**:6.26*f*, 6.28, 6.28*f*
 multiple-layer approach, **IV**:6.26, 6.26*f*–6.27*f*
 Teflon overcoat, **IV**:6.27
 and EIT, **IV**:14.3
- Far-infrared region, multilayer reflectors for, **IV**:7.52, 7.52*f*
- Far-infrared (FIR) telescopes, **IV**:6.48
- Far-off-resonance traps (FORTs), **IV**:20.23
- Farsightedness, **III**:23.2 (*see also* Hyperopia)
- FASCODE program, **V**:3.23, 3.24*f*
- Fast axis, **I**:12.25, 15.7
- Fast ignition, **IV**:21.54–21.55, 21.55*f*
- Fast power transients, for EDFAs, **V**:21.39–21.41, 21.45*f*
- Fast saturable absorbers, **IV**:18.8*f*, 18.9–18.10
- Fatigue, thermal, **II**:17.25
- Fatigue strength, of metals, **IV**:4.8
- Fax machines, **I**:24.6
- FDDI (fiber distributed data interface) standard, **V**:13.41, 23.2–23.3, 23.3*f*
- Federal Motor Vehicle Safety Standards (FMVSS), **II**:40.63
- Feedback:
 in fiber optic systems, **V**:13.30–13.32
 from laser diodes, **V**:13.21–13.24
 optical, **II**:16.2
 resonant optical, **II**:19.38, 19.38*f*
 from SOAs, **V**:19.8, 19.8*f*
 stabilization of, **II**:34.32

- Feedthrough, of optical disk data, **I**:35.14
- Femtosecond x-ray production,
IV:21.52–21.53, 21.53*f*
- Femtoseconds, **II**:20.1; **IV**:5.7
- Fermat's principle, **I**:1.11–1.13, 1.24
- Fermi choppers, **V**:63.14
- Fermi level, **IV**:8.21; **V**:30.1
- Fermi occupation functions, **II**:19.11
- Fermi pseudopotential, **V**:63.4, 63.5
- Fermi's golden rule, **IV**:8.25
- Ferrimagnetism, **I**:35.25, 35.26, 35.26*f*
- Ferroelectric bolometer arrays, **II**:28.11, 28.12,
 28.12*f*
- Ferroelectric detectors, **II**:33.10
- Ferroelectric oxides, **IV**:12.13–12.14, 12.13*t*
- Ferroelectric photorefractive materials,
IV:12.13–12.17, 12.13*t*
 barium titanate, **IV**:12.15–12.16, 12.16*f*
 lithium niobate and lithium tantalate,
IV:12.14
 potassium niobate, **IV**:12.16–12.17
 strontium barium niobate and related
 compounds, **IV**:12.17
 tin hypthiodiphosphate, **IV**:12.17, 12.18*t*
- Ferroelectric smectic phase, of liquid crystals,
V:8.12, 8.12*f*
- Feussner prisms, **I**:13.6, 13.7, 13.22*f*, 13.22–13.23
- FF5 glass (593355), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- Fiber (material), for optical fibers, **V**:9.4–9.7
- Fiber amplifiers, **V**:14.1–14.11
 categories and features of, **V**:14.1–14.2, 14.2*t*
 erbium-doped
 energy levels, **V**:14.4
 fast power transients, **V**:21.39*f*, 21.39–21.41,
 21.40*f*
 gain flattening, **V**:14.6–14.7, 21.38–21.39,
 21.39*f*
 gain formation, **V**:14.4–14.5, 14.5*f*
 gain peaking, **V**:21.38, 21.38*f*
 noise, **V**:14.6
 pump wavelength options, **V**:14.5–14.6
 semiconductor amplifiers vs., **V**:9.13, 9.14,
 14.1, 14.2*t*
 static gain dynamic and channel power
 equalization, **V**:21.41, 21.41*f*–21.42*f*
 in WDM networks, **V**:21.2*f*, 21.2–21.3
 erbium/ytterbium-doped, **V**:14.2, 14.2*t*,
 14.7–14.8
 infrared fibers for, **V**:12.3*t*
 parametric, **V**:14.10–14.11
- Fiber amplifiers (*Cont.*):
 praseodymium-doped fiber amplifiers
 (PDFAs), **V**:14.7
 Raman, **V**:14.8*f*, 14.8–14.9, 14.10*f*
 rare-earth-doped, **V**:14.2–14.4, 14.3*f*
 semiconductors vs., **V**:9.13–9.14
 ytterbium-doped, **V**:14.7
- Fiber attenuation, optical power loss and,
V:21.13–21.14
- Fiber bandwidth, of WDM networks, **V**:21.2*f*,
 21.2–21.3
- Fiber Bragg gratings (FBGs), **V**:17.1–17.9
 applications, **V**:17.8*f*, 17.8–17.9
 chirped, **V**:21.22–21.23, 21.23*f*, 21.25*f*,
 21.25–21.26
 fabrication, **V**:17.4–17.8, 17.5*f*–17.7*f*
 and fiber lasers, **V**:25.8, 25.16, 25.18, 25.30,
 25.31
 long-period gratings vs., **V**:24.9, 24.11
 photosensitivity of, **V**:17.2–17.3
 properties of, **V**:17.3–17.4, 17.4*f*
 sensors based on, **V**:24.5–24.8, 24.6*f*–24.7*f*
- Fiber feedback, **IV**:17.16
- Fiber interferometers, **I**:32.14–32.16, 32.15*f*
- Fiber lasers, **II**:16.34; **V**:25.1–25.33
 architectures, **V**:25.9–25.18, 25.19*f*
 all-fiber monolithic systems, **V**:25.16*f*,
 25.16–25.18
 free space, **V**:25.13–25.15, 25.14*f*, 25.15*f*
 pumping techniques, **V**:25.9–25.13, 25.11*f*,
 25.12*f*
 bulk lasers vs., **V**:25.5–25.6
 dopants for, **V**:25.22–25.26, 25.23*t*
 fabrication of, **V**:25.26–25.29, 25.27*t*
 growth of, **V**:25.5, 25.5*f*
 history of, **V**:25.3–25.4, 25.4*f*
 infrared fibers for, **V**:12.3*t*
 limitations of, **V**:25.6–25.7
 LMA designs for, **V**:25.18–25.22, 25.19*f*
 operation of, **V**:25.7–25.8
 spectral and temporal modalities of,
V:25.29–25.33
- Fiber length, for rare-earth-doped amplifiers,
V:14.2–14.3, 14.3*f*
- Fiber optic amplifiers, **V**:14.1–14.11
 categories and features of, **V**:14.1–14.2, 14.2*t*
 erbium-doped
 energy levels, **V**:14.4
 fast power transients, **V**:21.39*f*,
 21.39–21.41, 21.40*f*

- Fiber optic amplifiers, erbium-doped (*Cont.*):
 gain flattening, **V**:14.6–14.7, 21.38–21.39, 21.39*f*
 gain formation, **V**:14.4–14.5, 14.5*f*
 gain peaking, **V**:21.38, 21.38*f*
 noise, **V**:14.6
 pump wavelength options, **V**:14.5–14.6
 semiconductor amplifiers vs., **V**:9.13, 9.14, 14.1, 14.2*t*
 static gain dynamic and channel power equalization, **V**:21.41, 21.41*f*–21.42*f*
 in WDM networks, **V**:21.2*f*, 21.2–21.3
 erbium/ytterbium-doped, **V**:14.7–14.8
 parametric, **V**:14.10–14.11
 praseodymium-doped, **V**:14.7
 Raman fiber, **V**:14.8*f*, 14.8–14.9, 14.10*f*
 rare-earth-doped, **V**:14.2–14.4, 14.3*f*
 ytterbium-doped, **V**:14.7
- Fiber optic bundles, human receptors as, **III**:8.3–8.5
- Fiber optic chemical sensors, **V**:12.3*t*
- Fiber optic communication links, **V**:15.1–15.17
 distortions and dynamics range of, **V**:15.5–15.6, 15.6*f*
 figures of merit for, **V**:15.2–15.6, 15.3*f*, 15.4*f*
 link budget analysis for, **V**:15.6–15.17
 extinction ratio, **V**:15.13
 installation loss, **V**:15.6–15.7, 15.8*t*
 optical power penalties, **V**:15.8–15.17, 15.10*f*
- Fiber optic communication standards, **V**:23.1–23.8
 ATM/SONET, **V**:23.6
 ESCON, **V**:23.1–23.2, 23.2*f*
 Ethernet, **V**:23.7
 FDDI, **V**:23.2–23.3, 23.3*f*
 Fibre Channel standard, **V**:23.4, 23.5*f*, 23.5*t*
 fibre channel standard, **V**:23.4, 23.5*f*, 23.5*t*
 InfiniBand, **V**:23.8, 23.8*t*
- Fiber optic features:
 of plant tissues, **III**:8.26–8.28, 8.27*f*
 of sponges, **III**:8.28–8.29
- Fiber optic gyroscopes (FOG), **I**:21.2, 21.35–21.37, 21.36*f*, 21.36*t*
- Fiber optic networking, micro-optics-based components for, **V**:18.1–18.12
 attenuators, **V**:18.2, 18.9
 beam splitters, **V**:18.6, 18.6*f*
 circulators, **V**:18.3, 18.3*f*, 18.10
 directional couplers, **V**:18.2, 18.3, 18.3*f*, 18.9, 18.9*f*
- Fiber optic networking, micro-optics-based components for (*Cont.*):
 Faraday rotators, **V**:18.7, 18.7*f*
 filters, **V**:18.6
 gratings, **V**:18.5–18.6, 18.6*f*
 GRIN-rod lenses, **V**:18.7, 18.8, 18.8*f*
 isolators, **V**:18.3, 18.10, 18.10*f*
 mechanical switches, **V**:18.4, 18.5, 18.5*f*, 18.11, 18.11*f*, 18.12*f*
 MEMS mirrors and switches, **V**:18.8, 18.8*f*, 18.11, 18.12*f*
 multiplexers/demultiplexers/duplexers, **V**:18.4, 18.4*f*, 18.10–18.11
 network functions, **V**:18.2–18.5
 polarizers, **V**:18.7, 18.7*f*
 power splitters, **V**:18.2*f*, 18.2–18.3, 18.9, 18.9*f*
 prisms, **V**:18.5, 18.5*f*
- Fiber optic networks and systems, **V**:9.14–9.15
 detectors in, **V**:13.2–13.3, 13.63–13.73
 avalanche photodiodes, **V**:13.71–13.73
 MSM detectors, **V**:13.73
 pin diodes, **V**:13.63–13.71, 13.64*f*, 13.66*f*, 13.66*t*, 13.68*f*
 Schottky photodiodes, **V**:13.73
 modulators in, **V**:13.2, 13.48–13.63
 electroabsorption, **V**:13.55–13.60, 13.56*f*, 13.57*f*
 electro-optic, **V**:13.61
 electrorefractive, **V**:13.61–13.62
 lithium niobate, **V**:13.48–13.55, 13.49*f*, 13.51*f*, 13.54*f*
 semiconductor interferometric, **V**:13.63
 sources for, **V**:13.1–13.48
 distributed Bragg reflector lasers, **V**:13.28–13.29, 13.29*f*
 distributed feedback lasers, **V**:13.30*f*, 13.30–13.32, 13.31*f*
 laser diodes, **V**:13.3*f*, 13.3–13.24, 13.4*f*, 13.6*f*, 13.7*f*, 13.10*f*, 13.14*f*, 13.16*f*, 13.17*f*, 13.19*f*–13.21*f*
 light-emitting diodes, **V**:13.36–13.42, 13.38*f*, 13.40*f*
 quantum well lasers, **V**:13.24–13.28, 13.25*f*–13.27*f*
 strained layer quantum well lasers, **V**:13.26*f*, 13.26–13.28, 13.27*f*
 tunable lasers, **V**:13.32–13.36, 13.33*f*–13.36*f*
 vertical cavity surface-emitting lasers, **V**:13.42–13.48, 13.43*f*, 13.45*f*

- Fiber optic networks and systems (*Cont.*):
 [see also related topics, e.g.: Optical time-division multiplexed (OTDM) communication networks]
- Fiber optic sensors, **V**:24.1–24.13
 extrinsic Fabry-Perot interferometric, **V**:24.2*f*, 24.2–24.4, 24.3*f*
 fiber Bragg grating, **V**:24.5–24.8, 24.6*f*–24.7*f*
 intrinsic Fabry-Perot interferometric sensors, **V**:24.4*f*, 24.4–24.5
 long-period grating sensors, **V**:24.8–24.13, 24.9*f*–24.12*f*, 24.11*t*, 24.13*t*
- Fiber optics:
 fiber alignment for, **II**:17.33
 focal-length measurement with, **II**:12.25
 LED considerations with, **II**:17.33–17.34
 in nonimaging optics, **II**:39.21
- Fiber pigtail connection, **V**:13.8
- Fiber pulling, for fiber lasers, **V**:25.26
- Fiber pump lasers, **IV**:17.7–17.11, 17.9*f*–17.11*f*
- Fiber Raman amplifiers, **IV**:22.15
- Fiber Raman lasers, **V**:10.7
- Fiber squeezers, **I**:15.24
- Fiber-based couplers, **V**:16.1–16.6, 16.2*f*, 16.5*f*
- Fiber-optic octocouplers, **II**:17.33–17.34
- Fiberoptic-coupled (FO) II SSAs, **II**:31.20*f*, 31.20–31.22, 31.21*f*, 31.21*t*
- Fibers:
 photonic bandgap, **IV**:2.23
 slow light propagation in, **IV**:22.13–22.15, 22.14*f*
- Fiber-to-fiber excess loss, **I**:21.13–21.14
- Fibre Channel Arbitrated Loop (FC-AL), **V**:23.4
- Fibre Channel standard, **V**:23.2, 23.4, 23.5*f*, 23.5*t*
- Fibroblast, **III**:16.1
- Fick's law, **I**:9.12
- Field (lens), **I**:1.74
- Fields, of rays, **I**:1.13
- Field additivity (color vision), **III**:11.51, 11.52*f*
- Field angles:
 of apertures, **I**:1.75, 1.75*f*, 1.76
 of Glan-Thompson type prisms, **I**:13.12
- Field curvature, **I**:1.91, 29.7, 29.37 (see also Petzval curvature)
- Field curvature plot, **II**:2.4, 2.4*f*, 2.5
- Field flatness (aberration), **I**:1.91
- Field flattener lenses, **I**:17.28
- Field intensities, of waves, **I**:2.5–2.6
- Field lenses, **I**:1.82, 1.82*f*, 17.10; **II**:1.8, 1.10, 1.10*f*
- Field measurements (color vision):
 and first-site adaptation, **III**:11.27, 11.29, 11.30*f*
 and second-site adaptation, **III**:11.29, 11.31
- Field method (color vision), **III**:11.11–11.12
 defined, **III**:11.2
 evidence for higher-order mechanisms from, **III**:11.79–11.80
- Field of view (FOV), **I**:1.74; **II**:3.4, 7.11, 24.11, 31.1; **III**:8.4
 for head-mounted displays, **III**:25.2, 25.5, 25.7–25.8, 25.9*f*
 in Keplerian afocal lenses, **I**:18.11
 in Maxwellian viewing, **III**:6.6*f*, 6.6–6.7
 and multifoil optics, **V**:48.1
 for reflective and catadioptric objectives, **I**:29.34–29.35, 29.35*f*, 29.36*f*
 for scatterometers, **V**:1.6, 1.10, 1.12
 in telescopes, **I**:18.15*f*, 18.15–18.16
- Field patch trace, **II**:39.7–39.8
- Field plots, of aberration curves, **II**:2.4–2.5
- Field quality, in optical systems, **III**:5.11
- Field sensitivities (color vision), **III**:11.46–11.57
 achromatic detection and chromatic adaptation, **III**:11.47–11.49, 11.48*f*
 chromatic adaptation and the Sloan notch, **III**:11.49, 11.51
 detection contours and field adaptation, **III**:11.53–11.54, 11.54*f*
 field additivity, **III**:11.51, 11.52*f*
 first- and second-site adaptation, **III**:11.51, 11.52, 11.53*f*
 habituation or contrast adaptation experiments, **III**:11.54–11.56, 11.55*f*
 multiple cone inputs, **III**:11.49, 11.50*f*
 noise-masking experiments, **III**:11.56–11.57
 Stiles' π -mechanisms, **III**:11.46, 11.47, 11.47*f*
- "Field" sensitivity method (color vision), **III**:11.11–11.12
- Field size (lens), **I**:28.13
- Field stops, **I**:1.74, 17.9, 29.5, 29.37; **II**:7.7, 7.8*f*, 7.9*f*, 34.18–34.19, 34.19*f*
- Field-effect transistor (FET) amplifiers, **V**:13.70
- Field-effect transistors (FETs), **I**:21.38; **II**:27.10
- Field-enhanced pyroelectric arrays (see Ferroelectric bolometer arrays)
- Field-flattened Schmidt objective, **I**:29.14–29.15

- Fields, in color CRTs, **III**:22.9*f*, 22.9–22.10, 22.10*f*
- Field-weighted-average resolution, **V**:44.10
- Fifth-order oblique spherical aberration, **V**:45.4
- Figures of merit (FOM), **II**:24.13, 33.23–33.28
- acousto-optic, **V**:6.16–6.17
 - fiber optic communication links, **V**:15.2–15.6, 15.3*f*, 15.4*f*
 - for infrared photodetectors, **II**:25.12
 - minimum resolvable temperature (MRT), **II**:33.27*f*, 33.27–33.28
 - NE ΔT , **II**:33.25*f*, 33.26*f*
 - in spectroradiometry, **II**:38.5–38.6
- Filament notching, **II**:40.30
- Filaments:
 - lamp, **II**:15.20*f*
 - light bulb, **II**:40.25, 40.27, 40.29*f*, 40.29–40.30
- Fill factor, of binary optics, **I**:23.8
- Fill gases, light bulb, **II**:40.29*f*, 40.30
- Filling factor, electro-optic effect and, **V**:13.50
- Film(s):
 - and black surfaces, **IV**:6.15
 - camera, **I**:25.5, 25.6
 - photographic (*see* Photographic films)
 - polymer, **IV**:12.23–12.25, 12.26*t*–12.27*t*
 - semiconductor-doped dielectric, **IV**:18.11
 - thin (*see* Thin-film coatings)
- Film x-ray detectors, **V**:60.8, 60.9*t*, 60.10*t*
- Filtered backprojection (FBP), **V**:32.2
- Filters (filtering), **IV**:7.1–7.114
 - antireflection coatings, **IV**:7.15–7.32
 - of absorbing and amplifying media, **IV**:7.26, 7.27
 - homogeneous-layer, **IV**:7.16–7.23, 7.17*f*–7.19*f*, 7.20*t*–7.21*t*, 7.22*f*–7.23*f*
 - inhomogeneous and structured, **IV**:7.23–7.26, 7.24*f*, 7.26*f*
 - at nonnormal angle of incidence, **IV**:7.28*f*–7.31*f*, 7.28–7.31
 - nonoptical properties of, **IV**:7.31–7.32, 7.32*f*
 - surface reflections and optical performance, **IV**:7.15–7.16, 7.16*f*
 - of surfaces carrying thin films, **IV**:7.27–7.28, 7.28*f*
 - universal, **IV**:7.26, 7.27*f*
 - bandpass, **IV**:7.73–7.96
 - about, **IV**:7.73, 7.77*f*–7.78*f*, 7.77–7.78
 - angular properties of, **IV**:7.91–7.94, 7.92*f*, 7.93*f*
 - Filters (filtering), bandpass (*Cont.*):
 - with multiple peaks, **IV**:7.90, 7.91*f*
 - narrow- and medium-, **IV**:7.78–7.83, 7.79*f*, 7.80*f*, 7.82*f*–7.88*f*
 - stability and temperature dependence of, **IV**:7.94
 - very narrow, **IV**:7.83, 7.88–7.89, 7.89*f*
 - wedge filters, **IV**:7.90, 7.91, 7.91*f*
 - wide, **IV**:7.90, 7.90*f*
 - for XUV and x-ray regions, **IV**:7.94–7.96, 7.95*f*–7.96*f*
 - beam splitters, **IV**:7.61–7.67, 7.62*f*–7.68*f*
 - achromatic beam splitters, **IV**:7.62*f*–7.65*f*, 7.62–7.65
 - color-selective beam splitters, **IV**:7.65–7.66, 7.66*f*, 7.67*f*
 - geometrical considerations for, **IV**:7.61–7.62
 - blocking, **II**:38.8
 - blur, **II**:32.34, 32.34*f*
 - Bragg, **V**:63.14
 - with coatings
 - measurements on, **IV**:7.12–7.14
 - transmission and reflection of, **IV**:7.3
 - for coherent optical image enhancement, **I**:11.14–11.17
 - cutoff, heat-control, and solar-cell cover, **IV**:7.53–7.60
 - cutoff filters, **IV**:7.53–7.60, 7.54*f*, 7.55*f*, 7.57*f*, 7.59*f*–7.61*f*
 - heat reflectors, **IV**:7.58
 - solar-cell cover filters, **IV**:7.58
 - Fabry-Perot, **V**:13.33, 22.8
 - in fiber optic networking, **V**:18.6
 - frequency-guiding, **V**:22.7–22.9
 - guided-mode resonance, **V**:25.2, 25.30
 - high performance optical multilayer coatings, **IV**:7.96–7.98, 7.97*f*
 - infrared, **II**:40.12
 - intensity, **III**:5.14, 5.14*t*
 - interference, **II**:34.36
 - interference polarizers and polarizing beam splitters, **IV**:7.69–7.73, 7.70*f*–7.72*f*, 7.76*f*–7.77*f*
 - with low reflection, **IV**:7.104*f*–7.105*f*, 7.104–7.106
 - in II SSA cameras, **II**:31.7
 - Mach-Zehnder, **I**:21.23, 21.24*f*, 21.25; **V**:14.6, 21.39
 - matched, **IV**:12.28–12.29, 12.29*f*, 12.30*f*

- Filters (filtering) (*Cont.*):
- multilayer reflectors, **IV**:7.39–7.53
 - all-dielectric broadband reflectors, **IV**:7.39, 7.40*f*, 7.45*f*–7.47*f*, 7.45–7.47
 - coatings for ultrafast optics, **IV**:7.47–7.48, 7.48*f*
 - for far-infrared region, **IV**:7.52, 7.52*f*
 - graded reflectivity mirrors, **IV**:7.52
 - imperfections in, **IV**:7.40–7.43, 7.41*f*–7.43*f*
 - for interferometers and lasers, **IV**:7.39*f*–7.40*f*, 7.39–7.40
 - narrowband reflection coatings, **IV**:7.43, 7.44*f*
 - rejection filters, **IV**:7.48–7.50, 7.49*f*–7.51*f*
 - resonant reflectors, **IV**:7.43–7.45, 7.44*f*
 - for soft X ray and XUV regions, **IV**:7.53
 - narrowband, **I**:3.3
 - for networking, **V**:18.6
 - neutral density, **II**:40.52
 - neutral filters, **IV**:7.67, 7.67*f*–7.68*f*
 - neutron, **V**:63.18*t*, 63.18–63.19, 63.19*f*
 - notch, **II**:22.10*f*, 22.10–22.11, 22.11*f*
 - novelty, **IV**:12.32, 12.33*f*–12.35*f*
 - for pattern recognition, **I**:11.12–11.14
 - phase coatings, **IV**:7.101, 7.101*f*–7.104*f*, 7.102
 - reflection, **IV**:7.5, 7.5*f*
 - reflection coatings and, **IV**:7.106*f*–7.113*f*, 7.106–7.113
 - spatial, **I**:11.5–11.6, 11.6*f*
 - special purpose coatings, **IV**:7.113–7.114, 7.114*f*
 - theory of, **IV**:7.1*f*, 7.2
 - thin-film coatings
 - and antireflection coatings, **IV**:7.27–7.28, 7.28*f*
 - manufacturing of, **IV**:7.10–7.12
 - of metal, **IV**:7.104, 7.104*f*
 - theory and design of, **IV**:7.5–7.10, 7.6*f*, 7.9*f*
 - transmission, **IV**:7.3–7.5, 7.4*f*
 - for two or three spectral regions, **IV**:7.98*f*–7.101*f*, 7.98–7.100
 - two-material periodic multilayers theory for, **IV**:7.32–7.38, 7.33*f*–7.38*f*
 - UV, **II**:40.12
 - of x-ray tube source spectra, **V**:54.9
 - [*see also* Acousto-optic tunable filters (AOTFs)]
- Filtrate absorption, **IV**:1.21
- Finesse:
 - coefficient of, **I**:2.31
 - of Fabry-Perot etalon, **I**:2.34, 2.35
 - of interferometers, **I**:32.6
- Finish models, for surface scattering, **I**:8.14–8.15
- Finite conjugate afocal relays, **I**:18.11–18.12, 18.12*f*
- Finite rays, **I**:1.35
- Finite-difference time-domain (FDTD) analysis, for PCFs, **V**:11.7
- Finite-difference time-domain (FDTD) solution (to Maxwell's equations), **IV**:9.3
- Finite-difference time-domain (FDTD) technique, **I**:7.16–7.17
- Finitely distant objects, systems with, **II**:1.6, 1.6*f*
- First-order layout techniques, **II**:1.3–1.16
 - achromatism, **II**:1.14–1.15, 1.15*f*
 - afocal attachments, **II**:1.8, 1.9*f*
 - afocal systems, **II**:1.7, 1.7*f*
 - athermalization, **II**:1.15–1.16, 1.16*f*
 - axial/principal rays, **II**:1.12
 - component power minimization, **II**:1.13, 1.13*f*
 - condensers, **II**:1.10–1.11, 1.11*f*
 - defined, **II**:1.4
 - field lenses, **II**:1.8, 1.10, 1.10*f*
 - magnifiers and microscopes, **II**:1.8
 - ray-tracing, **II**:1.4–1.5
 - reasonableness of layout, **II**:1.13–1.14
 - two-component systems, **II**:1.5–1.7
 - zoom or varifocal systems, **II**:1.11–1.12
- First-order optics, **I**:1.29, 1.37
- First-order retardation plates, **I**:13.46
- First-site adaptation (color vision), **III**:11.11, 11.15–11.17, 11.16*f*
 - defined, **III**:11.2
 - field measurements and, **III**:11.27, 11.29, 11.30*f*
 - and field sensitivities, **III**:11.51, 11.52, 11.53*f*
 - incompleteness of, **III**:11.17
 - signals reaching second site, **III**:11.17
 - Weber's law and contrast coding, **III**:11.15–11.16, 11.16*f*
- Fish lenses, **III**:19.6, 19.6*f*
- Fish-eye lenses, **I**:27.6; **III**:19.2–19.3
- Fissures, and surface absorption, **IV**:6.15
- “Fitness for use,” **II**:17.25
- Fitting error, **V**:5.41
- 5 × 7 matrix LED displays, **II**:17.31–17.32
- Five-axis machining, **II**:10.7*f*

- Fixation, **III**:1.42
 defined, **III**:13.2
 stability of, **III**:1.44, 1.45*f*
- Fixed dispersion compensation, **V**:21.22–21.23, 21.23*f*
- Fixed interferogram evaluation, **II**:13.14–13.15
- Fixed-orientation mirrors, **II**:6.17
- Fixed-pattern noise (FPN), **I**:26.11
- Fixer, film, **II**:29.5
- Fizeau interferometers, **I**:2.24–2.26, 2.25*f*, 32.2, 32.2*f*, 32.17; **II**:12.14, 13.8–13.9, 13.9*f*, 13.10*f*, 13.18, 14.4, 14.5*f*
- Fizeau techniques, for surface figure metrology, **V**:46.3
- Flame brushing, of fiber Bragg gratings, **V**:17.2
- Flame hydrolysis (FHD), **I**:21.13–21.14
- Flame-sprayed aluminum, **IV**:6.57
- Flaming arcs, **II**:15.23, 15.24*f*, 15.26*t*–15.27*t*
- Flanges:
 annular, **II**:6.3–6.4, 6.4*f*
 continuous ring, **II**:6.4*f*, 6.11
- Flash (camera), **I**:25.16, 25.17*f*
- Flash blindness, **II**:40.9
- FLASH free-electron laser facility, **V**:58.1
- Flash lamps, **II**:16.16–16.17, 16.17*f*
- Flashlamp pumped Nd:glass amplifiers, **IV**:21.5
- Flat panel detectors, for x-ray imaging, **V**:61.3–61.7, 61.4*f*, 61.6*t*, 61.7*f*
- Flat-field objective optics, **I**:30.30–30.33
- Flat-medial-field objectives, **I**:29.11
- Flex-Pivots (flexures), **II**:6.19
- Flexural strength, of metals, **IV**:4.70, 4.70*t*
- Flexure, **III**:20.16
- Flexure mountings, **II**:6.5, 6.5*f*, 6.15–6.17, 6.16*f*
- Flicker:
 in fluorescent lamps, **II**:40.32–40.33
 impact of, **II**:40.12
 in incandescent lamps, **II**:40.30
- Flicker floor, **II**:22.3
- Flicker noise, **II**:24.11
- Flicker photometrics, **III**:10.34
- Flight control, **II**:19.3
- Flint glasses, **IV**:2.28, 2.41*t*–2.43*t*
- Flip chip packaging, **II**:18.6, 18.6*f*
- Flip-and-fold approach, **II**:39.29*f*
- Floating diffusion, **II**:32.14, 32.15*f*
- Floating gate output amplifiers, **II**:32.15
- Flood-illuminated AO ophthalmoscope, **III**:15.3, 15.16*f*, 15.16–15.17, 15.17*f*
- Floquet theory, **IV**:21.23
- Floquil, **IV**:6.44
- Fluence, **III**:8.26
 corneal photoablation, **III**:16.16, 16.17, 16.19*f*
 defined, **III**:8.1
- Fluor crown glass, **IV**:2.41*t*
- Fluorescence, **II**:34.13, 40.30
 chlorophyll, **IV**:1.49
 laser-induced, **V**:3.21
 lenticular, **III**:1.21
 stray light as result of, **III**:1.21
 x-ray, **V**:28.1, 54.8, 62.5, 62.6*f*
 and polycapillary x-ray optics, **V**:53.10–53.11, 53.11*f*
 and x-ray diffraction, **V**:28.5–28.6, 28.6*f*
- Fluorescence imaging with one nanometer accuracy (FIONA), **I**:28.23
- Fluorescence line narrowing (FLN), **I**:10.18, 31.29, 31.29*f*; **V**:2.13–2.14, 2.14*f*
- Fluorescent lamps, **II**:40.30–40.33
 applications for, **II**:40.26*t*
 characteristics of, **II**:40.25*t*
 cold cathode, **V**:8.30, 8.30*f*
 construction of, **II**:40.33*f*, 40.34*f*
 elements of, **II**:40.31*f*
 emission spectrum of, **II**:40.32*f*, 40.35*f*
 types of, **II**:40.28*f*
- Fluorescent microscopy, **I**:28.48–28.49
- Fluorides, for fiber lasers, **V**:25.28–25.29
- Fluorite (CaF₂), **IV**:2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.61*t*, 2.69*t*, 2.77*t*
- Fluoro flint glass, **IV**:2.42*t*
- Fluoroalumininate glass, **V**:12.4, 12.4*t*
- Fluoro-silicone/acrylate (F-S/A), **III**:20.1
- Fluorozirconate glass (ZBLAN), **V**:12.5, 12.5*f*
 and fiber lasers, **V**:25.3, 25.24, 25.27*t*, 25.28
 fluoroalumininate glass vs., **V**:12.4, 12.4*t*
- Flux:
 of electron bunches, **V**:55.17
 luminous, **II**:37.4, 37.4*t*, 37.6
 of polarization elements, **I**:15.9
 radiant, **II**:37.3, 37.4*t*
 and radiative transfer, **I**:9.12–9.13, 9.13*f*
 total luminous, **II**:37.4*t*, 37.6
 total radiant, **II**:37.4*t*, 37.6
- Flux budget, for fiber optics, **II**:17.33
- Flux density (*see* Irradiance)
- Fly-by-light (FBL), **II**:19.3
- FM spectroscopy, **II**:22.13–22.14
- F-number, **I**:1.79, 17.9; **II**:34.20
- Focal depth, **II**:4.7–4.8, 4.8*f*

- Focal Gaussian lenses, **I**:1.45–1.53, 1.47*f*
 conjugate equations of, **I**:1.49*f*, 1.49–1.50
 magnifications and distances in, **I**:1.50–1.53
 nodal points of, **I**:1.48, 1.48*f*, 1.49
 principal focal points of, **I**:1.47
 principal planes of, **I**:1.47, 1.47*f*, 1.48
 reduced coordinates of, **I**:1.53
- Focal length, **II**:1.5–1.7, 12.21–12.25; **III**:23.2
 effective
 of camera lenses, **I**:27.3*f*–27.5*f*, 27.7*f*–27.16*f*,
 27.18*f*–27.22*f*, 27.25
 of Gaussian lenses, **I**:1.48
 of microlenses, **I**:22.10
- fiber optics, **II**:12.25
- focimeters, **II**:12.22*f*, 12.22–12.23, 12.23*f*
- Fourier transforms, **II**:12.24
- in gradient index optics, **I**:24.5–24.6
- microlenses, **II**:12.24
- Moiré deflectometry, **II**:12.23, 12.24*f*
- nodal slide bench, **II**:12.22, 12.22*f*
- primary, **I**:3.12
- of surfaces, **I**:1.39–1.40
- in systems of revolution, **I**:1.40
- Talbot autoimages, **II**:12.23, 12.24
- Focal lines, **I**:1.58
- Focal plane arrays (FPAs), **II**:33.3
 hybrid, **II**:33.14–33.23
 microbolometer, **II**:33.13–33.14
- MIS photogate, **II**:33.10–33.11, 33.11*f*, 33.12*f*
- monolithic, **II**:33.10–33.14
- Focal planes, **I**:1.57, 1.70
- Focal plane-to-focal plane conjugate matrices,
I:1.69
- Focal points:
 front and rear, **I**:1.40, 1.47
 of lens systems, **I**:17.7
 principal, **I**:1.47, 1.58
- Focal ratio, **I**:29.5, 29.37; **II**:34.20
- Focal surfaces, of grazing incidence optics,
V:45.5
- Focal-plane-to-focal-plane geometry, **I**:11.3*f*,
 11.3–11.4
- Focimeters, **II**:12.22*f*, 12.22–12.23, 12.23*f*; **III**:12.9
- Fock states, **IV**:23.8–23.9, 23.14
- Focus:
 in color CRTs, **III**:22.19
 in Maxwellian viewing, **III**:5.4–5.5, 5.5*f*, 6.5–6.6
 in monochrome CRTs, **III**:22.7
 range of, **I**:1.85
 and responses to distance, **III**:13.30
- Focus anisoplanatism, **V**:5.27–5.29, 5.28*f*–5.30*f*,
 5.42–5.43
- Focus athermalization techniques, **II**:6.22–6.24
 active athermalization, **II**:6.24, 6.24*f*
 passive athermalization, **II**:6.22, 6.23*f*, 6.24
 single material designs, **II**:6.22, 6.23*f*
- Focus distance, **II**:1.5–1.7
- Focus error signal (FES), **I**:35.12–35.14, 35.13*f*
- Focus of expansion, **III**:13.2
- Focus shift, thermal, **II**:8.2–8.4, 8.3*t*, 8.4*t*
- Focused beams, **IV**:15.41; **V**:54.16
- Focused single crystal diffraction,
V:53.12*f*–53.13*f*, 53.12–53.14
- Focusing:
 of gratings and monochromators, **V**:38.3*f*,
 38.3–38.6, 38.4*t*–38.5*t*
 in grazing-incidence neutron optics,
V:64.3–64.7, 64.4*f*–64.7*f*
 Kerr, **V**:7.39
 in neutron and x-ray optics, **V**:26.9–26.11,
 63.22
 of optical disks, **I**:35.9–35.12, 35.11*f*, 35.12*f*
 with refractive x-ray lenses, **V**:37.8–37.11,
 37.9*f*, 37.10*f*
 (*see also* Hard x-rays, nanofocusing of)
- Focusing polycapillary x-ray optics,
V:53.9–53.10, 53.10*t*, 53.14
- Fog film, **II**:29.9
- Fokker-Planck equation, **II**:23.37;
IV:20.10–20.11
- Font, **III**:23.2
- Foot-candle (unit), **II**:34.43, 36.7, 36.7*t*, 37.7*t*;
III:23.2
- Foot-lambert (unit), **II**:34.43, 36.7, 36.8*t*, 37.7,
 37.7*t*
- Forbes method, **II**:3.20
- Forbidden bands, **IV**:9.2
- “Forbidden” colors, **III**:11.74
- Fore-optics, **II**:38.7
- Form birefringence (term), **I**:15.41
- Form dichroism (term), **I**:15.41
- Förster resonance energy transfer (FRET),
I:28.23
- 45° half-wave linear retarders, Mueller matrices
 for, **I**:14.12*t*
- 45° linear polarizers, **I**:14.10*t*
- 45° quarter-wave linear retarders, Mueller
 matrices for, **I**:14.12*t*
- Forward Brillouin scattering, **V**:11.25, 11.26
- Forward drift velocity, **IV**:21.7

- Forward error correction (FEC) coding,
V:19.27, 21.26
- Forward light, II:40.43, 40.44*f*
- Forward Raman amplifiers, IV:15.4, 15.4*f*
- Forward Raman generators, IV:15.4, 15.4*f*
- Forward-looking infrared (FLIR) systems,
I:30.22, 30.23, 30.51; II:33.4; IV:6.54
- Foster prisms, I:13.7, 13.18*f*, 13.21–13.22
- Foucault prisms, I:13.7, 13.17
- Foucault test, II:12.19, 13.2*f*, 13.2–13.3, 13.3*f*
- Four-detector photopolarimeters (FDPs),
I:16.14*f*–16.16*f*, 16.14–16.16
- Fourier analysis:
of interferograms, II:13.16–13.17, 13.17*f*
for radiative transfer, I:9.11
- Fourier approach to optics, III:4.8
- Fourier crosstalk matrix, V:32.3
- Fourier differentiation, V:46.11–46.12
- Fourier domain filters, I:11.13
- Fourier domain OCT (*see* Spectral domain OCT)
- Fourier intervals, V:27.4
- Fourier Theory of Optics:
defined, III:6.1
in Maxwellian viewing, III:6.10–6.12
- Fourier transform lenses, I:18.12
- Fourier transform plane, I:33.16–33.19, 33.18*f*,
33.20*f*
- Fourier transform spectrophotometers,
II:35.9
- Fourier transform spectroscopy, V:2.5, 2.6*f*
- Fourier transforming infrared spectrometer
(FTIR), V:1.14
- Fourier transforms, III:2.23, 18.1; V:6.15, 16.4,
46.8, 55.9, 55.10
in analog optical and image processing,
I:11.3–11.5, 11.5*f*
and coherence theory, I:5.20
and diffraction, I:3.2, 3.24
for focal plane-to-focal plane matrices,
I:1.69
for focal-length determination, II:12.24
and transfer functions, I:4.2
of uniformly illuminated linear aperture,
I:30.54
- Fourier-transform spectrometers, IV:5.60*f*,
5.60–5.61, 5.72, 7.12
- Four-phase CCDs, II:32.13*f*, 32.13–32.14,
32.16*f*
- Four-photon absorption (4PA), IV:19.10
- Four-powered-mirror lenses, I:18.21, 18.21*f*
- Four-wave mixing (FWM), IV:18.3
coherent anti-Stokes, IV:15.2*t*, 15.3*t*, 15.4,
15.4*f*
degenerate, IV:18.17
and EIT, IV:14.24*f*, 14.24–14.26, 14.27*f*
in optical fibers, V:10.2, 10.9–10.11, 10.11*f*
resonant, IV:14.28
and SOAs, V:19.13, 19.14*f*, 19.27, 19.33*f*,
19.33–19.35, 19.34*f*
and solitons, V:22.11, 22.15
and third-order optical nonlinearities,
IV:16.27–16.28, 16.28*f*
transient, IV:18.17*f*, 18.17–18.18
in WDM networks, V:21.19–21.20, 21.20*f*
- Four-wave mixing phase conjugation, IV:12.6*f*,
12.6–12.7
- Fovea, III:1.3, 1.3*f*, 1.9, 1.15, 2.5*f*, 14.9
and cone spectral sensitivities, III:10.17
defined, III:4.1, 18.1
disparities in, III:2.40
function of, III:13.3
receptor size in, III:2.23
- Foveal avascular zone, III:14.9
- Foveal pit, III:14.11
- Foves, III:2.24
- Fractal model of surface finish, I:8.14
- Fractals:
brownian, I:8.9, 8.17
Fresnel-Kirchhoff approximation for,
I:8.8–8.9
- Fracture toughness:
of crystals and glasses, IV:2.32, 2.32*t*
of metals, IV:4.8, 4.70, 4.70*t*
- Frame interline transfer (FIT) CCDs, II:32.27*f*,
32.29
- Frame transfer (FT) CCD image sensors,
II:32.26–32.28, 32.27*f*, 32.28*f*, 32.32
- Frame transfer (FT) CCD TDI FPAs, II:33.11,
33.12*f*
- Frames:
in color CRTs, III:22.9*f*, 22.9–22.10, 22.10*f*
in OTDM, V:20.7
- Frank elastic constants, V:8.22
- Frankford Arsenal prisms, I:19.3*t*, 19.18*f*–19.24*f*,
19.18–19.24
- Frank-Ritter-type prisms, I:13.6, 13.6*f*,
13.13–13.14
- Franz-Keldysh effect, I:21.11, 21.11*f*; IV:12.22;
V:13.59
- Fraunhofer approximation, I:5.14, 5.16

- Fraunhofer diffraction, **I**:3.24*f*–3.26*f*, 3.24–3.28;
III:6.10, 6.12; **V**:42.2, 63.25
 Airy diffraction as, **I**:28.17
 conducting screens for, **I**:3.33*f*
 and gratings, **I**:20.3
 Fraunhofer regime, in coherent x-ray optics,
V:27.2
 Fraunhofer theory, **I**:7.11
 Fredericksz cells, **I**:14.31, 15.21, 15.22
 Fredericksz threshold voltage, of LC cells,
V:8.27
 Fredkin gate, **IV**:23.11
 Free carriers and free-carrier effects,
IV:5.47–5.52
 in crystals, **IV**:2.15
 cyclotron resonance, **IV**:5.47–5.50,
 5.48*f*–5.50*f*
 in fiber optic systems, **V**:13.4
 free-carrier Faraday rotation, **IV**:5.50–5.51
 impurity magnetoabsorption, **IV**:5.51*f*,
 5.51–5.52
 in semiconductors, **IV**:5.33–5.36, 5.35*f*–5.37*f*,
 5.81, 5.82*f*
 Free electron properties, of solids, **IV**:8.21–8.24
 Drude model, **IV**:8.21, 8.22*f*
 interband transitions in metals, **IV**:8.21
 plasmons, **IV**:8.23–8.24
 reflectivity, **IV**:8.23
 Free polarization decay (FPD), **IV**:11.7–11.11,
 11.8*f*, 11.10*f*, 11.11*f*
 Freedericksz threshold voltage, of LC cells,
 8.27
 Free spectral range (FSR):
 of bandpass filters, **IV**:7.77
 of interferometers, **I**:2.34, 2.35, 2.35*f*, 32.5
 Free (newtonian) viewing, **III**:5.2–5.4, 5.3*f*
 limitations of, **III**:5.4
 retinal illuminance, **III**:5.2–5.3, 5.3*f*
 the troland, **III**:5.3–5.4
 Free-carrier Faraday rotation, **IV**:5.50–5.51
 Free-electron lasers (FELs), **II**:16.36–16.37,
 16.37*f*, 23.43–23.45; **V**:29.4, 41.9–41.10,
 41.10*f*, 48.1, 58.1–58.2
 Free-electron-laser (FEL) lamps, **II**:15.11,
 15.12, 15.13*f*
 Free-exciton (FE) luminescence, **IV**:5.72, 5.73*f*
 Free-space fiber lasers, **V**:25.9, 25.10, 25.13–25.15,
 25.14*f*, 25.15*f*
 Frenet equation, **I**:1.19
 Frenkel excitons, **IV**:5.26, 5.26*t*
 Frequency:
 and acousto-optic interaction, **V**:6.12–6.14,
 6.13*f*, 6.14*f*, 6.16
 crossover, **V**:8.16
 and drift, **II**:22.2
 of electro-optic modulators, **V**:7.35
 Greenwood, **V**:5.19, 5.22, 5.42
 of liquid crystals, **V**:8.18*f*, 8.17–8.18
 Nyquist, **V**:27.4
 phase and amplitude responses vs., **II**:22.6*f*,
 22.6–22.7, 22.7*f*
 stability of, **II**:22.2
 Stokes, **V**:10.8, 21.42
 Tyler, **V**:5.17
 Frequency chirping, **V**:13.1, 13.17*f*,
 13.17–13.18
 Frequency comb, **II**:20.1, 20.2, 20.7–20.9
 Frequency conversion, nonlinear optical,
IV:14.24*f*, 14.24–14.28, 14.27*f*
 Frequency discriminators:
 for laser locking, **II**:22.12–22.14
 optical cavity-based, **II**:22.14–22.16, 22.17*f*
 Frequency domain, **I**:7.17
 Frequency mixing, **IV**:16.3*t*
 Frequency modulation (FM), **II**:19.36, 19.36*f*,
 22.4; **V**:7.24–7.25, 7.25*f*
 Frequency shift keying (FSK), **II**:19.36
 Frequency shifters, in electro-optic phase shifts,
V:13.51
 Frequency-controlled lasers, **II**:22.7, 22.7*f*
 Frequency-division multiplexing (FDM),
V:9.16
 Frequency-guiding filters, **V**:22.7–22.9
 Frequency-modulated (FM) systems, **V**:9.16
 Frequency-modulation interferometers,
I:32.9, 32.9*f*, 32.10
 Frequency-resolved optical gating (FROG),
II:21.8, 21.9
 Frequency-resolved optical gating for complete
 reconstruction of attosecond bursts
 (FROG-CRAB), **II**:21.9
 Frequency-selective-feedback lasers,
II:19.37*f*, 19.38
 Fresnel amplitude reflection coefficients,
I:12.7–12.8
 Fresnel amplitude transmission coefficients,
I:12.8
 Fresnel diffraction, **I**:3.14–3.24; **V**:27.2, 27.4,
 40.9, 63.25

- Fresnel equations, **I**:12.6–12.13, 12.15;
V:41.3, 64.2
 for absorbing materials, **I**:12.10–12.13, 12.13*f*
 coordinate system for, **I**:12.6–12.7, 12.7*f*
 for nonabsorbing materials, **I**:12.8–12.10,
 12.9*f*
- Fresnel integrals, **V**:6.15
- Fresnel intensity reflectivities, **I**:8.6
- Fresnel lenses, **I**:22.31*f*–22.33*f*, 22.31–22.37,
 22.35*f*–22.37*f*; **II**:39.9–39.10, 39.10*f*,
 40.45*f*, 40.46
- Fresnel losses, **I**:22.10
- Fresnel number, **I**:3.25
- Fresnel phase zones, **V**:27.3
- Fresnel propagation kernels, **I**:6.6
- Fresnel reflection, **V**:13.6, 13.53, 17.3, 25.8
- Fresnel reflection coefficients, **I**:8.10; **IV**:7.106
- Fresnel relations, **IV**:8.11, 8.15
- Fresnel rhombs, **I**:13.45, 13.50*f*; **V**:43.5*f*,
 43.5–43.6
- Fresnel Risley prisms, **I**:19.27*f*
- Fresnel waves, **V**:27.2
- Fresnel zone plates, **V**:40.2, 42.3, 42.3*f*, 55.16
- Fresnel zones, **I**:3.4*f*, 3.4–3.6
 for cylindrical wavefronts, **I**:3.6*f*, 3.7*f*,
 3.13–3.14, 3.14*f*
 and Fraunhofer diffraction, **I**:3.25
 MLLs and, **V**:42.16*f*
 opaque strip obstruction, **I**:3.20–3.21, 3.21*f*
- Fresnel-Kirchhoff approximation, **I**:8.5–8.9
- Fresnel-Kirchhoff diffraction formula, **I**:3.21,
 3.22, 3.32
- Fresnel's biprism, **I**:2.16, 2.16*f*
- Fresnel's law, **V**:63.20
- Fresnel's mirror, **I**:2.16, 2.16*f*
- Fresnel-Soret zone plates, **V**:40.2
- Fried's coherence diameter, **V**:5.2, 5.7, 5.9–5.10
- Fried's coherence length, **V**:4.8
- Fringe localization, **I**:2.14
- Fringe power, of polymers, **IV**:3.11
- Fringe visibility (contrast), **I**:2.7–2.8
- Fringe washout (OCT), **III**:18.15
- Fringe-counting interferometers, **I**:32.8, 32.8*f*
- Fringes, in sensor signals, **V**:24.3
- Fringes of equal chromatic order (FECO)
 technique, **V**:46.2
- Fringes of equal inclination, **I**:2.20–2.22, 2.21*f*,
 2.22*f*
- Fringes of equal thickness, **I**:2.22–2.24, 2.23*f*
- Frogs legs mirror, **IV**:12.7, 12.8*f*
- Front focal lengths, **I**:1.40
- Front focal points, **I**:1.40, 1.47
- Front principal planes, **I**:1.48
- Front vertex distance (FVD), camera lens
 performance and, **I**:27.3*f*–27.5*f*,
 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
- Front vertex power (FVP):
 contact lenses, **III**:20.7–20.8, 20.8*t*
 defined, **III**:20.1
- Frontoparallel plane, **III**:13.2
- Fuji Photo Film USA, **II**:29.25
- Fujicolor, **II**:29.14
- Full width at half maximum (FWHM),
II:16.5, 16.6, 20.3, 21.2; **IV**:18.3;
V:17.6–17.7, 17.7*f*, 24.8
- Full width at half maximum (FWHM) points,
I:30.9
- Full-aperture techniques, in surface figure
 metrology, **V**:46.3
- Full-frame arrays, of CCDs, **I**:26.3–26.4, 26.4*f*
- Fully functional polymers, **IV**:12.27*t*
- Functional specifications (optical design),
II:4.2
- Fundamental absorption edge, **IV**:5.21
 absorption near, **IV**:5.21*f*, 5.21–5.22
 high-energy transitions above, **IV**:5.29–5.33,
 5.30*f*–5.34*f*
- Fundamental array mode, **II**:19.27
- Furniture-integrated lighting system, **II**:40.13*f*
- Fused glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- Fusing, in xerographic systems, **I**:34.10
- Fusion neutron production, **IV**:21.53
- Fusion splicing, **V**:25.18
- Gabor objective, **I**:29.20
- Gabor zone plates, **V**:40.4, 40.5, 40.7
- Gain:
 of avalanche photodiodes, **V**:13.71–13.72
 defined, **II**:16.9
 of EDFAs, **V**:14.4–14.5, 14.5*f*
 in photomultipliers, **II**:27.7
 polarization-dependent, **V**:14.9, 19.18–19.20,
 21.18
 Raman, **V**:10.5, 10.6, 21.42*f*
 of SOAs, **V**:19.4*f*–19.6*f*, 19.4–19.6
- Gain clamping, **V**:19.14*f*, 19.14–19.15
- Gain coefficient, **II**:16.9–16.10
- Gain control, in CRTs, **III**:22.8
- Gain coupling, **II**:19.29

- Gain dynamics, of SOAs, **V**:19.12–19.13, 19.13*f*, 19.14*f*
- Gain efficiency, of EDFAs, **V**:14.6
- Gain flattening, **V**:14.6–14.7, 21.38–21.39, 21.39*f*
- Gain medium, **II**:16.3
- Gain narrowing, in steady-state Stokes scattering, **IV**:15.21
- Gain peaking, **V**:21.38, 21.38*f*
- Gain per unit length, of lasers, **V**:13.7, 13.8
- Gain ripple, **V**:19.8, 19.8*f*, 19.19, 19.19*f*
- Gain saturation, **II**:16.10; **IV**:18.8, 18.8*f*; **V**:13.15, 13.17
- Gain stability margin, **II**:22.9–22.10
- Gain without inversion, and EIT, **IV**:14.18–14.19
- Gain-bandwidth (GB), **II**:26.17
- Gain-coupled arrays, **II**:19.27
- Gain-guided index antiguided fibers (GG IAG), **V**:25.2, 25.19*f*, 25.22
- Gain-guided laser diodes, **V**:13.5
- Gain-guided phased array, **II**:19.28*f*
- Galilean lenses, **I**:18.7, 18.15*f*, 18.15–18.17, 18.16*f*
- Galilean telescopes, **I**:18.15; **II**:1.7*f*
- Galileo Galilei, **I**:18.2, 28.3
- Gallagher-Pritchard (GP) model (of trap-loss collisions), **IV**:20.29
- Gallium aluminum arsenide (GaAlAs) LEDs, **II**:17.12, 17.12*f*, 17.13*f*
- Gallium arsenide (GaAs):
 composition, structure, and density of, **IV**:2.39*t*
 dispersion formulas for, **IV**:2.62*t*
 elastic constants of, **IV**:2.44*t*
 lattice vibration model parameters for, **IV**:2.77*t*
 linear-chain model calculations for, **IV**:5.19*f*
 local vibrational modes for, **IV**:5.19*f*
 luminescence in, **IV**:5.72, 5.73, 5.74*f*
 mechanical properties of, **IV**:2.47*t*
 multiphonon absorption of vacuum-grown, **IV**:5.18*f*
 optical modes of, with zinblende structure, **IV**:2.69*t*
 optical properties of, **IV**:2.56*t*
 Raman scattering of, **IV**:5.80, 5.81, 5.81*f*
 thermal properties of, **IV**:2.51*t*
- Gallium arsenide (GaAs) lasers, **II**:19.7
- Gallium arsenide (GaAs) LEDs, **II**:17.8, 17.9, 17.9*f*
- Gallium arsenide (GaAs) semiconductor diode lasers, **II**:16.18*f*, 16.18–16.19
- Gallium arsenide phosphide (GaAsP) emitters, **II**:17.32
- Gallium arsenide phosphide (GaAsP) LEDs, **II**:17.9–17.10, 17.10*f*, 17.15–17.17
 energy band diagram for, **II**:17.5*f*
 homojunction in, **II**:17.10, 17.11*f*
 light degradation in, **II**:17.27*f*
 performance summary of chips in, **II**:17.16*t*
- Gallium arsenide phosphide (GaAsP) photodiodes, **II**:24.49, 24.49*f*, 24.50*f*
- Gallium arsenide phosphide (GaAsP) substrate, **II**:17.22
- Gallium arsenide (GaAs) quantum well photodetectors, **II**:25.16*f*, 25.16–25.17, 25.17*f*
- Gallium nitride (GaN), **II**:17.22; **IV**:2.39*t*, 2.46*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.62*t*, 2.70*t*, 2.77*t*
- Gallium nitride (GaN) photovoltaic detectors, **II**:24.42, 24.43, 24.45*f*, 24.46, 24.46*f*, 24.47
- Gallium phosphide (GaP), **II**:17.16, 17.20–17.22; **IV**:2.39*t*, 2.44*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.62*t*, 2.69*t*, 2.77*t*; **V**:6.16, 6.17*t*, 6.29*t*, 6.30, 6.31*t*, 6.34*t*
- Gallium phosphide (GaP) dynodes, **II**:24.42, 24.44*f*
- Gallium phosphide (GaP) photodiodes, **II**:24.47–24.49, 24.48*f*
- Gallium-doped germanium (Ge:Ga) infrared detectors, **II**:24.100
- Gallium-doped silicon (Si:Ga) infrared detectors, **II**:24.95, 24.95*f*, 24.96, 24.96*f*
- Galvanometer scanners, **I**:30.41–30.44, 30.43*f*, 30.44*f*
- Galvo-driven Brewster plates, **II**:22.18–22.19
- Gamma cameras, **V**:32.2
- Gamma rays, in SPECT imaging, **V**:32.1–32.2
- Ganglion cells, **III**:2.10–2.11
 contrast sensitivity functions of, **III**:2.12*f*
 density of, **III**:2.11
 and LGN, **III**:2.12
 linear density of, **III**:2.6*f*
 midget, **III**:2.10–2.11, 2.10*n*
 off-center, **III**:2.10
 on-center, **III**:2.10
 parasol, **III**:2.10*n*
 P-cells, **III**:11.76, 11.79
 physiology of, **III**:2.11
 transfer functions of, **III**:2.11–2.12

- Ganzfeld, **III**:1.12
- Gap effect (color vision), **III**:11.72, 11.74
- Gap modes (GMs), **IV**:5.17
- Gas chromatography-mass spectroscopy (GC-MS), **II**:33.4
- Gas detectors, **V**:63.32–63.33
- Gas lights, **II**:40.40
- Gas mantle, **II**:15.17–15.19, 15.19*f*
- Gas permeable (GP) contact lenses, **III**:12.12, 20.3
- aberrations and, **III**:20.24
- aspheric, **III**:20.22
- base curve radius for, **III**:20.3, 20.5
- bitoric, power of, **III**:20.17–20.20
- center thickness of, **III**:20.6
- edge thickness of, **III**:20.6
- lacrimal or tear lens with, **III**:20.12–20.15
- OAD/OZD of, **III**:20.5
- posterior peripheral curve systems of, **III**:20.5, 20.6
- and residual astigmatism, **III**:20.15
- Gas phase media, dephasing in, **IV**:14.12–14.13
- Gaseous laser gain media, **II**:16.30–16.31, 16.31*f*
- Gases, strong field nonlinear optics in, **IV**:21.27–21.31, 21.28*f*
- Gas-filled lamps, **II**:34.31
- Gas-puff sources, **V**:57.3, 57.3*t*
- Gate modulation, **II**:33.19*t*, 33.20*f*, 33.22
- Gated integration, **II**:27.12–27.13, 27.13*f*, 27.15
- Gating:
- amplitude, **II**:21.7
- double optical, **II**:21.8
- frequency resolved, **II**:21.8, 21.9
- FROG, **II**:21.9
- FROG-CRAB, **II**:21.9
- polarization, **II**:21.7–21.8
- two-color, **II**:21.7
- Gaunt factors, **V**:56.9–56.10
- Gauss illuminated eyepieces, **II**:12.12
- Gauss law, **IV**:2.6
- Gauss points, of lenses, **I**:1.44
- Gaussian analyses, of lenses:
- afocal, **I**:18.4–18.7, 18.7*f*
- focusing, **I**:18.2–18.4, 18.3*f*, 18.5*f*
- Gaussian apertures, **V**:37.6
- Gaussian approximation of sensitivity, for optical fiber receivers, **V**:9.10
- Gaussian error function, **V**:15.3
- Gaussian image point, **V**:38.5–38.6
- Gaussian integral, for noise, **V**:15.2, 15.4
- Gaussian intensity distribution, **II**:39.28, 39.29*f*
- Gaussian lenses, **I**:1.44–1.55
- afocal, **I**:1.45, 1.46*f*, 1.53*f*, 1.53–1.54, 1.54*f*
- focal, **I**:1.45–1.53, 1.47*f*
- conjugate equations of, **I**:1.49*f*, 1.49–1.50
- magnifications and distances in, **I**:1.50–1.53
- nodal points of, **I**:1.48, 1.48*f*, 1.49
- principal focal points of, **I**:1.47
- principal planes of, **I**:1.47, 1.47*f*, 1.48
- reduced coordinates of, **I**:1.53
- notation for, **I**:1.45*t*
- properties of, **I**:1.54
- systems of, **I**:1.54–1.55
- Gaussian line profiles, **V**:3.14, 56.5–56.7
- Gaussian line shape, **II**:16.6*f*
- Gaussian mode, **II**:16.21
- Gaussian noise, **III**:2.17, 2.18
- Gaussian optics, **I**:1.29, 1.44
- Gaussian parameters (optical design), **II**:4.5–4.6, 4.6*t*
- Gaussian spectra, **V**:15.13
- Gaussian statistics, for wave propagation, **V**:5.9
- Gaussian transmission, **V**:37.5
- Gaussian-shaped beam, **II**:16.22
- Gauss-Seidel iterative method, **V**:3.21
- Gaze control, **III**:13.29–13.30
- Gaze eccentricity, **III**:13.8
- Gaze holding, **III**:1.42
- Gaze shifting, **III**:1.42
- GDx Nerve Fiber Analyzer, **I**:15.41
- Gédamine, **I**:13.9, 13.10, 13.11*f*, 13.20
- Geiger counters, **V**:60.5, 60.9*t*
- Geiger region, **V**:60.5
- Gelbstoff, in water, **IV**:1.13 (*see also* Yellow matter)
- Gemini North telescope, **V**:5.20, 5.21*f*
- General Conference on Weights and Measures (CGPM), **II**:36.2
- General deviation prisms, **I**:19.3*t*, 19.28, 19.28*f*–19.29*f*
- General Electric, **II**:15.29, 15.30, 15.48, 19.29*t*
- General Photonics, **I**:15.24
- General system data (optics), **II**:3.3, 3.4
- Generalized ellipsometry (GE), **I**:16.19
- Generalized Lagrange invariant (étendue), **I**:1.22, 1.81, 13.7
- Generalized pupil function, **III**:2.3
- Generalized Rabi frequency, **IV**:11.4
- Generating step (of optics fabrication), **II**:9.4
- Generation noise, **II**:24.11

- Generation-recombination (GR) current,
II:25.7–25.8
- Generation-recombination (GR) noise,
II:24.11
- GENLN2 (code), **V**:3.23
- Geometrical configuration factor (GCF), **II**:7.1,
 7.2, 7.22; **IV**:6.12
- Geometrical étendue, **II**:38.8
- Geometrical optical transfer function (GOTF),
II:3.16
- Geometrical optics:
 for aberrations of point images, **I**:1.86*f*,
 1.86–1.92
 and binary optics, **I**:23.5*f*–23.9*f*, 23.5–23.9,
 23.10*t*
 characteristic functions of, **I**:1.13–1.18, 1.15*f*
 for collineation, **I**:1.56–1.65, 1.65*f*
 and conservation of étendue, **I**:1.22
 defined, **I**:1.8
 for Gaussian lenses, **I**:1.44–1.55
 afocal, **I**:1.45, 1.46*f*, 1.53*f*, 1.53–1.54, 1.54*f*
 focal, **I**:1.45–1.53, 1.47*f*, 1.47–1.55, 1.48*f*,
 1.49*f*
 notation for, **I**:1.45*t*
 properties of, **I**:1.54
 systems of, **I**:1.54–1.55
 and images about known rays, **I**:1.43–1.44,
 1.44*f*
 for imaging, **I**:1.26–1.31, 1.30*f*
 at interfaces of homogeneous media,
I:1.23–1.26
 lens sizes and fields in, **I**:1.74–1.85
 apertures, **I**:1.74–1.77, 1.75*f*
 cosine-to-the-fourth approximation,
I:1.81
 field lenses, **I**:1.82, 1.82*f*
 fields, **I**:1.74, 1.77, 1.84
 F-number, **I**:1.79
 focus and defocus, **I**:1.82–1.85, 1.83*f*
 irradiance, **I**:1.80
 power per pixel, **I**:1.80
 pupils, **I**:1.76*f*, 1.76–1.79, 1.78*f*
 telecentricity, **I**:1.83–1.84
 total lens étendue, **I**:1.81
 vignetting, **I**:1.81, 1.81*f*, 1.82
 paraxial matrix methods, **I**:1.65–1.74
 and rays, **I**:1.8–1.13
 in heterogeneous media, **I**:1.18–1.22
 paths of, **I**:1.10–1.13
 and skew invariant, **I**:1.23
- Geometrical optics (*Cont.*):
 of systems of revolution, **I**:1.32–1.43
 paraxial optics of, **I**:1.37–1.43
 ray tracing in, **I**:1.35–1.37, 1.36*f*
 surfaces, **I**:1.32–1.35
 unfolded reflections, **I**:1.32
- Geometrical path length, **I**:1.11
- Geometrical point spread function (GPSF),
V:44.13
- Geometrical vector flux, **II**:39.21–39.22
- Geometrical wavefronts, **I**:1.12–1.13
- Geometry-controlled lasers, **II**:19.37*f*, 19.38
- Germanate:
 in fiber lasers, **V**:25.27*t*, 25.28
 in optical fibers, **V**:12.3–12.4, 12.6, 12.6*f*
- Germania glass, fused, **IV**:2.43*t*, 2.49*t*, 2.54*t*,
 2.59*t*, 2.67*t*
- Germanium:
 absorptance of, **IV**:4.48*t*
 in crystal form, **IV**:2.39*t*, 2.44*t*, 2.47*t*, 2.51*t*,
 2.56*t*, 2.62*t*, 2.68*t*
 thermal properties of
 elastic stiffness, **IV**:4.69*t*
 moduli and Poisson's ratio, **IV**:4.69*t*
 strength and fracture properties, **IV**:4.70*t*
- Germanium (Ge) avalanche photodiodes,
II:24.70*f*, 24.72*f*, 24.72–24.73, 24.73*f*
- Germanium (Ge) bolometers, **II**:28.5, 28.7*t*
- Germanium (Ge) detectors:
 copper-doped, **II**:24.84*f*, 24.85*f*, 24.96, 24.97,
 24.97*f*–24.99*f*
 gallium-doped infrared, **II**:24.100
 gold-doped, **II**:24.83–24.85, 24.84*f*–24.86*f*
 intrinsic photodetectors, **II**:24.70*f*–24.73*f*,
 24.70–24.73
 mercury-doped, **II**:24.84*f*, 24.92–24.95,
 24.93*f*–24.95*f*
pn and *pin*, **II**:24.70*f*–24.72*f*, 24.70–24.71
 zinc-doped, **II**:24.84*f*, 24.98–24.100, 24.99*f*
- Germanium gallium arsenide (GeGaAs)
 photodiodes, **II**:34.31
- Germanium (Ge) intrinsic photodetectors,
II:24.70*f*–24.73*f*, 24.70–24.73
- Germanium (Ge) low-temperature bolometers,
II:24.31–24.32, 24.32*f*, 24.33*f*
- Germanium photodiodes, **II**:38.9, 38.9*t*
- Germicidal lamps, **II**:15.35
- Ghost imaging, **IV**:23.13
- “Ghosts” (in gratings), **IV**:5.60
- Gibbon lenses, **III**:19.14, 19.14*f*

- Gilvin, **IV**:1.13 (*see also* Yellow matter)
- Gimbal-less two-axis scanning-micromirror devices (GSMs), **I**:30.61–30.62, 30.62*f*
- Glancing angle, of crystal monochromators, **V**:39.1, 39.2
- Glan-Foucault prisms, **I**:13.7, 13.9, 13.11*f*, 13.12–13.14
- Glan-Taylor prisms, **I**:13.7, 13.9*n*, 13.10*f*, 13.10–13.14, 13.11*f*
- Glan-Thompson prisms, **I**:13.6, 13.6*f*, 13.9–13.12, 13.10*f*, 13.18*f*, 13.22
field angle of, **I**:13.12
and optical spectrometers, **I**:31.7
sheet polarizers vs., **I**:13.27
transmission by, **I**:13.9–13.10, 13.11*f*
- Glan-type prisms, **I**:13.6, 13.6*f*, 13.8–13.15
defects and testing of, **I**:13.14–13.15
Frank-Ritter, **I**:13.6, 13.6*f*, 13.13–13.14
Glan-Foucault, **I**:13.7, 13.9, 13.11*f*, 13.12–13.14
Glan-Taylor, **I**:13.7, 13.9*n*, 13.10*f*, 13.10–13.14, 13.11*f*
Glan-Thompson, **I**:13.6*f*, 13.9–13.12, 13.10*f*, 13.11*f*, 13.27
Lippich, **I**:13.6, 13.6*f*, 13.12–13.13
Marple-Hess, **I**:13.12, 13.13
precautions with, **I**:13.14
- Glare:
with computer work, **III**:23.4–23.5, 23.5*t*
defined, **III**:23.2
and exterior lighting, **II**:40.62
in human eye, **III**:1.20
limiting of, **II**:40.10, 40.41
and visual discomfort, **II**:40.9–40.12, 40.11*t*
and windows, **II**:40.41
- Glare stops (*see* Lyot stops)
- Glass(es):
amorphous, **IV**:12.26*t*
antireflection coatings for, **IV**:7.26–7.28, 7.28*f*
common, **IV**:2.3
and crystals, **IV**:2.1–2.77
defined, **IV**:2.33
for fiber lasers, **V**:25.27*t*, 25.28–25.29
fluorozirconate (ZBLAN), **V**:12.5, 12.5*f*
and fiber lasers, **V**:25.3, 25.24, 25.27*t*, 25.28
fluoroaluminate glass vs., **V**:12.3*t*, 12.4, 12.4*t*
fluoroaluminate, **V**:12.4, 12.4*t*
formation of optical, **II**:9.3–9.4
heavy-metal fluoride, **V**:12.1–12.5, 12.2*f*, 12.3*t*, 12.4, 12.5*f*
- Glass(es) (*Cont.*):
heavy-metal oxide, **V**:12.2*t*, 12.3*t*, 12.3–12.7, 12.4*t*, 12.5*f*–12.7*f*
for hollow waveguides, **V**:12.2*f*, 12.11–12.13, 12.12*f*
hybrid organic-inorganic, **IV**:12.27*t*
material properties of, **IV**:2.27–2.36
characteristic temperatures, **IV**:2.32, 2.33
combinations of, **IV**:2.36
correlations of, **IV**:2.36
elastic properties, **IV**:2.30–2.31, 2.31*t*
hardness and strength, **IV**:2.31–2.32, 2.32*f*, 2.32*t*
heat capacity and Debye temperature, **IV**:2.33–2.34
material designation and composition, **IV**:2.28–2.30, 2.29*f*
naming of, **IV**:2.27
thermal conductivity, **IV**:2.35*f*, 2.35–2.36
thermal expansion, **IV**:2.34*f*, 2.34–2.35
unit cell parameters, molecular weight, and density, **IV**:2.30
mechanical properties of, **IV**:2.49*t*
for microlenses, **I**:22.9–22.10, 22.10*t*, 22.11*f*, 22.12*t*, 22.13*t*, 22.14*f*
negative core-cladding index difference, **V**:11.14, 11.15
for objectives, **I**:29.2, 29.2*t*
optical, **II**:5.9
optical applications of, **IV**:2.17–2.27
dielectric tensor and optical indicatrix, **IV**:2.17–2.19, 2.19*f*
dispersion formulas, **IV**:2.21–2.23
nonlinear optical coefficients, **IV**:2.26–2.27, 2.27*t*
scatter, **IV**:2.27
thermo-optic coefficients, **IV**:2.24*f*, 2.24–2.26
total power law, **IV**:2.19–2.20, 2.20*f*
as optical materials, **IV**:2.4–2.5
optical properties of, **IV**:2.6, 2.8–2.9
lattice vibration model parameters, **IV**:2.76*t*–2.77*t*
origin and models, **IV**:2.9–2.17, 2.10*f*, 2.13*f*, 2.16*f*, 2.17*f*
room-temperature dispersion formulas, **IV**:2.66*t*–2.68*t*
summary table, **IV**:2.59*t*
as photographic film emulsion, **II**:29.4

- Glass(es) (*Cont.*):
 physical properties of, **IV**:2.37, 2.38*t*–2.43*t*
 optical glass reference table, **IV**:2.41*t*–2.43*t*
 physical constants, **IV**:2.8*t*
 specialty, and substrate materials, **IV**:2.43*t*
 symmetry properties, **IV**:2.5, 2.6*t*, 2.8*t*
 Raman bands of, **V**:11.24
 Rayleigh scattering in, **V**:11.21
 sol-gel formed, **I**:24.8
 thermal properties of, **IV**:2.54*t*
 tolerances for, **II**:5.9
 Zerodur, **V**:47.5
- Glass envelope, lightbulb, **II**:40.29, 40.29*f*
 Glass lenses, for spectacles, **III**:12.9
 Glass micro-pore optics, **V**:49.1–49.6,
 49.2*f*–49.6*f*
 Glass scatterers, **IV**:13.8
 Glass-based lasers, **IV**:21.5
 Glass-calcite Rochon prisms, **I**:13.20
 Glass-ceramics, **IV**:2.33
 Glaucoma, **III**:14.3, 18.25–18.27
 aging eyes, **III**:14.26–14.27
 cause of, **III**:14.6
 defined, **III**:14.1
 Glazebrook prisms, **I**:13.6, 13.9 (*see also*
 Glan-Thompson prisms)
 Glazing, window, **II**:40.41
 Global climate change, **V**:3.43–3.45, 3.44*f*
 Globar, **II**:15.17, 15.18*f*, 15.19, 15.19*f*
 Glossiness, **II**:40.5
 Glow lamps, **II**:40.39
 Glow modulator tubes, **II**:15.49, 15.50*f*, 15.51*f*,
 15.52*t*
 Gobs, glass, **II**:9.4
 Goddard Space Flight Center, **IV**:6.35, 6.39
 Goebel mirrors, **V**:26.10
 Goerz prism system, **I**:19.3*t*, 19.17, 19.17*f*
 GOES-13 satellite, **V**:44.16–44.17, 44.17*f*
 Goly cell detectors, **II**:28.2, 28.6, 28.7*t*
 Gold:
 absorptance of, **IV**:4.40*f*, 4.48*t*, 4.50*t*, 4.51*t*
 diamond turning and, **II**:10.5
 optical properties of, **IV**:4.14*t*, 4.23*f*
 physical properties of, **IV**:4.52*t*, 4.54*t*
 reflectance of, **IV**:4.31*t*–4.32*t*, 4.40*f*
 thermal properties of
 coefficient of linear thermal expansion,
IV:4.56*t*, 4.57*f*
 elastic stiffness, **IV**:4.69*t*
 moduli and Poisson's ratio, **IV**:4.69*t*
 Gold, thermal properties of (*Cont.*):
 at room temperature, **IV**:4.55*t*
 specific heat, **IV**:4.65*t*, 4.66*f*
 strength and fracture properties, **IV**:4.70*t*
 thermal conductivity, **IV**:4.58*t*, 4.60*f*–4.61*f*
 Gold black surfaces, **IV**:6.57
 Gold iridite, **IV**:6.21, 6.22*f*
 Gold-doped germanium (Ge:Au) detectors,
II:24.83–24.85, 24.84*f*–24.86*f*
 Gold-germanium (Au-Ge) alloys, **II**:17.24
 Goldpoint blackbody, **II**:15.9
 Gold-zinc (Au-Zn) alloys, **II**:17.24
 Goniometers (goniophotometers), **II**:12.10,
 40.52–40.53, 40.53*f*, 40.54*f*
 Gooch-Tarry first minimum condition, **V**:8.26
 Goos-Hanchen shift, **I**:21.4
 Goos-Hanschen shift, **V**:13.55
 Gordon inequality, **IV**:1.21
 Gordon-Haus effect, **V**:22.7–22.8, 22.11
 Gouffé method, **II**:15.7–15.9, 15.8*f*
 Graded index (GRIN) fibers, **V**:15.17
 Graded index profile, of optical waveguides,
I:21.4
 Graded index separate confinement
 heterostructure (GRIN SCH), **II**:19.14,
 19.14*f*; **V**:13.4*f*, 13.5
 Graded index-rod (GRIN-rod) lenses, **V**:18.7,
 18.8, 18.8*f*, 18.10, 18.11*f*
 Graded reflectivity mirrors, **IV**:7.52
 Graded-index (GRIN) films, **I**:16.9
 Gradient dispersion, **I**:24.3
 Gradient force, **IV**:20.8
 Gradient index (GRIN) optics, **I**:24.1–24.8;
III:19.1–19.15
 analytic solutions of, **I**:24.2
 axial gradient lenses, **I**:24.3–24.5, 24.4*f*
 axial gradients, **III**:19.5
 bovine lenses, **III**:19.8–19.11, 19.11*f*
 cat lenses, **III**:19.9
 eye lens, **III**:19.5–19.6
 fish lenses, **III**:19.6, 19.6*f*
 functional considerations, **III**:19.14–19.15
 guinea pig lenses, **III**:19.8
 human/primate lenses, **III**:19.12–19.14,
 19.13*f*, 19.14*f*
 materials, **I**:24.8
 mathematical representations of, **I**:24.2–24.3
 nature of index gradient, **III**:19.21–19.15
 octopus lenses, **III**:19.7, 19.7*f*
 pig lenses, **III**:19.11, 19.12*f*

- Gradient index (GRIN) optics (*Cont.*):
 rabbit lenses, **III**:19.8
 radial gradients, **I**:24.5*f*, 24.5–24.8, 24.7*f*;
III:19.3–19.5, 19.4*f*, 19.5*f*
 rat lenses, **III**:19.7–19.8
 spherical gradients, **III**:19.2–19.3
- Gradient tilt (G-tilt), **V**:4.3
 and adaptive optics, **V**:5.14–5.16
 and angle of arrival, **V**:4.23, 4.25, 4.26
- Gradient-freeze technique, **II**:17.21
- Gradients, of wavefronts, **V**:5.23
- Grain boundaries, of crystals, **IV**:2.4
- Grains and graininess:
 of photographic films, **II**:29.5
 of photographic images, **II**:29.18*t*,
 29.19–29.22, 29.21*f*
 of silver halide crystals, **II**:29.4
- Granularity, photographic film speed and,
II:30.19
- Graphical user interface (GUI), for SHADOW
 code, **V**:35.3
- Grasshopper* monochromator, **V**:38.3
- Grassmann's laws, **III**:10.8–10.9, 10.28
- Grating acuity tasks, **III**:2.34, 2.34*f*, 2.35*f*
- Grating equation, **I**:23.4; **II**:38.7–38.8;
V:38.2–38.3
- Grating interferometers, **I**:32.4, 32.5*f*
- Grating multiplexers, **I**:23.11, 23.12
- Grating polarizers, **I**:13.33
- Grating surface-emitting laser array, **II**:19.40*f*,
 19.40–19.41
- Gratings (*see* Diffraction gratings; Dispersive
 prisms and gratings)
 Bragg, **IV**:22.9–22.11, 22.10*f*; **V**:20.14, 25.29,
 25.30
 circular, **V**:40.1
 diffraction, **IV**:5.59–5.60
 fiber Bragg, **V**:17.1–17.9
 applications, **V**:17.8*f*, 17.8–17.9
 chirped, **V**:21.22–21.23, 21.23*f*, 21.25*f*,
 21.25–21.26
 fabrication, **V**:17.4–17.8, 17.5*f*–17.7*f*
 and fiber lasers, **V**:25.8, 25.16, 25.18, 25.30,
 25.31
 long-period gratings vs., **V**:24.9, 24.11
 photosensitivity, **V**:17.2–17.3
 properties of, **V**:17.3–17.4, 17.4*f*
 sensors based on, **V**:24.5–24.8, 24.6*f*–24.7*f*
 formation of, **IV**:12.1–12.3, 12.2*f*
 Fresnel diffraction of, **V**:40.9
- Gratings (*see* Diffraction gratings; Dispersive
 prisms and gratings) (*Cont.*):
 Hill, **V**:17.5
 long-period, **V**:21.39, 21.39*f*, 24.10–24.11,
 24.11*t*, 24.13*t*
 Moiré, **IV**:22.11, 22.12*f*
 for networking, **V**:18.5–18.6, 18.6*f*
 sampled, **V**:13.34, 13.34*f*
 short-period, **V**:24.6
 superstructure, **V**:13.34, 13.35*f*
 transmission, **IV**:12.7, 12.8*f*
 in VUV and soft x-ray region, **V**:38.1–38.8
 diffraction properties, **V**:38.1–38.3, 38.2*f*
 dispersion properties, **V**:38.6–38.7
 efficiency, **V**:38.8
 focusing properties, **V**:38.3*f*, 38.3–38.6,
 38.4*t*–38.5*t*
 resolution properties, **V**:38.7
- Gravitational-wave interferometers, **I**:32.21,
 32.21*f*
- Gray code, on optical disks, **I**:35.16
- Gray gel, **II**:30.4
- “Gray world” assumption, **III**:10.39
- Graybody, **II**:35.7
- Gray-scale masks, **I**:22.23*f*–22.25*f*, 22.23–22.25
- Grazing incidence optics:
 aberrations of, **V**:44.6–44.12, 44.7*f*,
 44.9*f*–44.11*f*, 45.1–45.8, 45.2*f*
 image formation with, **V**:44.3–44.18
 and multifoil optics, **V**:48.1–48.2
 and pumping by EUV lasers, **V**:58.3, 58.4*f*
 telescopes with, **V**:44.6–44.12, 44.7*f*,
 44.9*f*–44.11*f*
 x-ray, **V**:44.12*f*–44.18*f*, 44.12–44.18
 and x-ray mirrors, **V**:44.3–44.6, 44.4*f*–44.6*f*
- Grazing-angle reflection, **V**:63.21
- Grazing-incidence neutron optics, **V**:64.1–64.7
 diffractive scattering and mirror surface
 roughness, **V**:64.2–64.3
 and imaging focusing optics, **V**:64.3–64.7,
 64.4*f*–64.7*f*
 total external reflection, **V**:64.1–64.2
- Green flashes, **V**:3.43
- Greens function, **I**:3.22*f*, 3.22–3.23
 in approximations of multiple scattering,
I:9.10
 digitized, **I**:7.15, 7.16
 in scattering, **I**:9.3
- Green light, color film and, **II**:29.13, 29.13*f*
- Green-emitting AlInGaP devices, **II**:17.18

- Greenockite (CdS), **IV**:2.38*t*, 2.46*t*, 2.47*t*, 2.51*t*, 2.56*t*, 2.61*t*, 2.70*t*, 2.77*t*
- Green's function solution, **V**:6.14
- Greenwood frequency, **V**:5.19, 5.22, 5.42
- Gregorian objective, **I**:29.10
- Grinding:
 - aspheric, **II**:9.8
 - controlled, **IV**:19.3
- Grooved regions, on optical disks, **I**:35.15, 35.15*f*, 35.16
- Gross-Pitaevski equation, **IV**:20.34
- Ground loop noise, **II**:27.5, 27.6*f*
- Ground state, **II**:16.4
- Ground-based telescopes, **IV**:6.12
- Ground-state collisions, **IV**:20.29
- Group Delay Dispersion (GDD), **IV**:7.47–7.48, 7.48*f*
- Group Delay (GD) phase changes, **IV**:7.47, 7.48
- Group theory, **IV**:5.6
- Group velocity, **IV**:22.3–22.5
- Group velocity dispersion (GVD), **IV**:11.27, 22.4–22.5; **V**:11.18, 11.18*f*, 11.19*f*, 15.11
- Grown homojunctions, LED, **II**:17.8, 17.9, 17.9*f*
- Growth techniques:
 - for epitaxial layers, **II**:17.21–17.23
 - substrate, **II**:17.20, 17.21
- Grüneisen constant, **V**:6.17
- Grüneisen relationship, **IV**:2.34, 2.35
- GTWave technology, **V**:25.11
- Guard ring, **II**:24.11
- GUERAP (stray light analysis program), **IV**:6.19
- Guidance, in photonic crystal fibers, **V**:11.11–11.26
 - attenuation mechanisms, **V**:11.19–11.22, 11.20*f*, 11.21*f*
 - birefringence, **V**:11.17
 - Brillouin scattering, **V**:11.25*f*, 11.25–11.26
 - group velocity dispersion, **V**:11.18, 11.18*f*, 11.19*f*
 - Kerr nonlinearities, **V**:11.22–11.24, 11.24*f*
 - negative core-cladding index difference, **V**:11.14*f*–11.17*f*, 11.14–11.17
 - positive core-cladding index difference, **V**:11.12*f*, 11.12–11.14, 11.13*f*
 - Raman scattering, **V**:11.24
 - resonance and antiresonance, **V**:11.12
- Guide to the Expression of Uncertainty in Measurement* (IS), **II**:38.6
- Guided-mode resonance filters (GMRFs), **V**:25.2, 25.30
- Guided-wave nonlinear structures, **IV**:17.15–17.16
- Guides, neutron, **V**:63.15–63.18
- Guinea pig lenses, **III**:19.8
- Gunn diodes, **I**:31.21
- Gurney-Mott mechanism, **II**:29.5
- Habituation (color vision), **III**:11.54–11.56, 11.55*f*
 - and color appearance, **III**:11.29, 11.69–11.70
 - defined, **III**:11.2
 - second-site, **III**:11.19*f*, 11.19–11.20, 11.21*f*, 11.22*f*
- Hafnium dioxide-yttrium oxide (HfO₂:Y₂O₃), **IV**:2.56*t*, 2.62*t*, 2.77*t*
- H-aggregates, **II**:30.13
- Haidinger fringes, **I**:2.22, 2.22*f*, 2.27
- Haidinger interferometers, **II**:12.14
- Halation, **II**:30.4
- Half-bleaching constant, **III**:2.7
- Half-period zones, **I**:3.5; **V**:40.1
- Half-shade devices, **I**:13.56–13.57
- Halftoning, low-level vision models in, **III**:24.4
- Half-wave linear retarders, Mueller matrices for, **I**:14.12*t*, 14.13
- Half-wave plates, **I**:12.25, 12.27
- Half-wave retarders, Mueller matrices for, **I**:14.14
- Half-wave voltage, of electro-optic modulators, **V**:7.19, 7.22
- Halite (NaCl), **IV**:2.40*t*, 2.44*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.64*t*, 2.69*t*
- Halle prisms, **I**:13.16*f*, 13.17
- Halogen lamps, **II**:15.11, 15.12, 15.13*f*, 40.25*t*, 40.26*t*, 40.30
- Halon, **II**:38.12–38.13
- Halophosphates, **II**:40.31
- Halos, **III**:1.21
 - in aging eyes, **III**:14.13
 - with glaucoma, **III**:14.6
- Hamiltonian optics, **I**:1.13, 1.43
- Hamiltonian rays, **II**:3.12
- Hamilton's equations for rays, **I**:1.21
- Hanbury-Brown-Twiss (HB&T) effect, **II**:23.13, 23.13*f*, 23.14
- Hankel functions, **III**:8.1
 - in modal analysis, **III**:8.22
 - for waveguides, **III**:8.13

- Haploopia, **III**:13.2
- Haptic, **III**:21.1
- Hard mounting, of optics, **II**:6.1–6.4, 6.3*f*, 6.4*f*
- Hard x-ray beamlines, SHADOW code for, **V**:35.5, 35.5*f*
- Hard x-ray optics, astronomical, **V**:47.9–47.10
- Hard x-ray telescopes, **V**:50.2
- Hard x-rays, nanofocusing of, **V**:42.1–42.17
- history of, **V**:42.2*f*, 42.2–42.4, 42.3*f*
- instrumental beamline arrangement and measurements for, **V**:42.9*f*–42.12*f*, 42.9–42.10
- limitations of, **V**:42.15–42.17, 42.16*f*–42.17*f*
- with magnetron-sputtered MLLs, **V**:42.5–42.7, 42.6*f*–42.8*f*
- on MLLs with curved interfaces, **V**:42.14, 42.15*f*
- Takagi-Taupin calculations for, **V**:42.12–42.14
- volume diffraction calculations for, **V**:42.4–42.5, 42.5*f*
- with wedged MLLs, **V**:42.12–42.13, 42.13*f*, 42.14*f*
- Hardness:
- of crystals and glasses, **IV**:2.31, 2.32*f*
- of metals, **IV**:4.8
- of polymers, **IV**:3.2–3.3
- Hardware implementation, for adaptive optics, **V**:5.21–5.38
- higher-order wavefront sensing techniques, **V**:5.36–5.37
- laser beacons, **V**:5.27–5.34, 5.28*f*–5.31*f*, 5.33*f*
- real-time processors, **V**:5.34*f*, 5.34–5.35, 5.35*f*
- Shack-Hartmann technique, **V**:5.23*f*, 5.23–5.27, 5.25*f*, 5.26*f*
- tracking, **V**:5.21–5.23, 5.22*f*
- wavefront correctors, **V**:5.37–5.38, 5.38*f*
- Harmonic generation:
- high, **IV**:21.27–21.30
- harmonic yield and phase matching, **IV**:21.30
- quasi-classical model, **IV**:21.28*f*, 21.29–21.30
- from solid plasmas, **IV**:21.50–21.52, 21.51*f*
- third-order, **IV**:16.2, 16.3*t*
- in crystals, **IV**:16.14
- energy level diagrams for, **IV**:16.5*f*
- and semiconductors, **IV**:5.56
- Harmonic yield, **IV**:21.30
- Harting-Dove (Dove) prisms, **I**:19.3*t*, 19.9, 19.9*f*, 19.10, 19.10*f*
- Hartman testing, **I**:22.26
- Hartmann equation, for refraction index, **IV**:2.22
- Hartmann test, **II**:13.4–13.6, 13.5*f*
- Hartmann-Shack test, **II**:13.6*f*, 13.6–13.7
- Hartnack-Prazmowski prisms, **I**:13.16*f*, 13.17
- Hartree-Fock variational approach, **I**:10.11
- Hastings (Brashear-Hastings) prisms, **I**:19.3*t*, 19.25, 19.25*f*
- Hausdorff-Besicovitch dimension, **I**:8.8
- H&D curve (*see* D-log H curve)
- Header information, on optical disks, **I**:35.6
- Heading judgments, magnification and, **III**:13.14–13.15
- Headlamps:
- design of, **II**:40.21, 40.23, 40.23*f*
- low-beam, **II**:40.64*f*, 40.64–40.67, 40.65*t*, 40.66*f*, 40.66*t*
- Head-mounted displays (HMDs), **III**:25.1–25.12
- accommodation with, **III**:25.10–25.12
- binocular vision factors in design of, **III**:13.31–13.35
- distance conflicts, **III**:13.31–13.32
- optical errors, **III**:13.34–13.35
- spatial location conflicts, **III**:13.33, 13.34*f*
- visual-vestibular conflicts, **III**:13.32–13.33
- characterizing, **III**:25.7–25.10, 25.8*f*, 25.9*t*–25.10*t*
- common design considerations for, **III**:25.2–25.7
- comfort, **III**:25.2–25.5, 25.3*f*
- functionality, **III**:25.5–25.6, 25.6*f*, 25.7*f*
- safety, **III**:25.2
- usability, **III**:25.2
- Health-care facility lighting, **II**:40.58–40.60, 40.60*t*
- Heat capacity (specific heat):
- of crystals and glasses, **IV**:2.6*t*, 2.33
- of metals, **IV**:4.7, 4.10*t*, 4.53, 4.55, 4.55*t*, 4.65*t*, 4.66*f*–4.69*f*
- Heat reflectors, **IV**:7.58
- Heat sinks, **III**:5.15
- Heat-induced lensing effect, **IV**:16.22
- Heating:
- cluster electron, **IV**:21.34, 21.35
- inverse Bremsstrahlung, **IV**:21.37, 21.37*f*
- $\mathbf{j} \times \mathbf{B}$, **IV**:21.49
- self-heating, **IV**:17.12
- vacuum, **IV**:21.47*f*, 21.48–21.49
- Heat-pipe blackbody furnace, **II**:15.9*f*

- Heavy water, scattering in, **V**:63.10
- Heavy-hole (HH) bands, of strained layer quantum well lasers, **V**:13.27
- Heavy-metal fluoride glass (HMFG) fibers, **V**:12.1–12.5, 12.2*f*, 12.3*t*, 12.4*t*, 12.5*f*
- Heavy-metal oxide glass fibers, **V**:12.2*t*, 12.3*t*, 12.3–12.7, 12.4*t*, 12.5*f*–12.7*f*
- HEFT balloon payloads, **V**:47.10
- Heidelberg Retinal Tomograph (HRT), **III**:17.7, 17.9
- Height profilometry, **V**:46.3
- Height solves, **II**:3.6
- Heisenberg limit, **IV**:23.6–23.7
- Heisenberg number-phase uncertainty relation, **IV**:23.4
- Heisenberg uncertainty principle (HUP), **I**:21.11; **IV**:23.6
- Heisenberg-Langevin approach, to quantum theory of lasers, **II**:23.33–23.35
- Heisenberg's matrix mechanics, **I**:10.3
- Helium atoms, **V**:2.3
- Helium [$^3\text{He}(n, p)$] reaction, **V**:63.31
- Helium-cadmium (He-Cd) lasers, **II**:16.6, 16.15, 16.15*f*, 16.30
- Helium-neon (He-Ne) lasers, **II**:16.15, 16.15*f*, 16.30
- Helmholtz equation, **I**:3.2, 5.9, 5.10, 33.3; **V**:5.8
- Helmholtz invariant, **I**:1.77 (*see also* Two-ray paraxial invariant)
- Helmholtz's reciprocity theorem, **III**:8.1, 8.25, 8.27
- Hemispherical emittance, **II**:35.15
- Hemispherical total absorptance, **II**:35.8*t*
- Hemispherical-conical reflectance, **II**:35.5*t*, 35.6*f*, 35.6*t*
- Hemispherical-directional reflectance, **II**:35.5*t*, 35.6*f*, 35.6*t*
- Hemispherical-spectral absorptance, **II**:35.8*t*
- Henle fibers, **III**:8.2, 8.20
- Hering's law, **III**:13.20, 13.26
- Hermetic enclosure, for laser diodes, **V**:13.7, 13.7*f*
- Hermite-Gaussian functions, **V**:11.7
- Hermitian matrices, **I**:14.19, 14.41
- HERO balloon payload, **V**:47.10
- Herschel condition, **I**:1.31, 29.34, 29.37
- Herschelian objectives, **I**:29.6, 29.27
- Hertz (HZ), **III**:23.2
- Herzberg equation, **IV**:2.22
- Heterochromatic flicker photometry (HFP), **III**:11.37, 11.49
- Heterochromatic modulation photometry (HMP), **III**:11.37
- Heterodyne detection, **V**:3.34
- Heterodyne interferometers, **I**:32.10, 32.10*f*, 32.20, 32.20*f*; **II**:13.22
- Heterogeneous media, rays in, **I**:1.9, 1.18–1.22
- Heterojunction lasers, **II**:19.4
- Heterojunctions, **II**:17.12, 17.12*f*–17.15*f*, 17.13, 17.17, 26.9
- Heterophoria, **III**:13.2, 13.26
- Heterostructures, of fiber optic devices, **V**:13.2
- Hewlett-Packard double-frequency distance-measuring interferometer, **II**:12.9*f*, 12.9–12.10
- Hexagonal crystals:
 - anisotropic, **IV**:8.9*t*
 - dielectric constants of, **IV**:2.18
 - room-temperature elastic constants, **IV**:2.46*t*
 - symmetries of, **IV**:2.7*t*, 8.20*t*
- Hexagonal polynomials, **II**:11.21, 11.22*t*–11.25*t*, 11.36*t*, 11.38*f*, 11.39
- High aspect ratio microlithography (HARM), **V**:61.3
- High emittance-low absorptance coatings, **IV**:6.16
- High order harmonic generation (HHG), **IV**:21.27–21.30
 - harmonic yield and phase matching, **IV**:21.30
 - quasi-classical model, **IV**:21.28*f*, 21.29–21.30
 - from solid plasmas, **IV**:21.50–21.52, 21.51*f*
- High harmonic production, of x-ray lasers, **V**:58.2
- High magnetic field production, **IV**:21.53
- High performance optical multilayer coatings, **IV**:7.96–7.98, 7.97*f*
- High Resolution Doppler Imager (HRDI), **V**:3.36, 3.37*f*
- High-accuracy spectrophotometers, **II**:35.9
- High-brightness visible LEDs (HB-LEDs), **II**:18.1–18.6
 - about, **II**:18.1
 - epitaxial growth of, **II**:18.3
 - packaging of, **II**:18.5*f*, 18.5–18.6, 18.6*f*
 - processing of, **II**:18.3*f*, 18.3–18.4
 - semiconductor material systems for, **II**:18.1–18.2
 - solid-state lighting with, **II**:18.4*f*, 18.4–18.5
 - structure for modern InGaN, **II**:18.2*f*
 - substrates for, **II**:18.2*f*, 18.2–18.3

- High-dry objectives, **I**:28.11, 28.12*f*
 High-dye-yield yellow couplers, **II**:30.6
 High-energy radiation, **II**:15.40, 30.19–30.20
 High-energy transitions above fundamental edge, **IV**:5.29–5.33, 5.30*f*–5.34*f*
 Higher-order aberrations, **I**:29.37; **III**:16.6, 16.7
 correction of, **III**:16.9
 defined, **III**:15.1, 20.2
 Higher-order mechanisms (color vision), **III**:11.80
 Higher-order mode (HOM) fibers, **V**:25.2, 25.19*f*, 25.22
 High-gain oscillators, **II**:20.10–20.12, 20.11*f*
 High-intensity carbon arc lamps, **II**:15.21–15.23, 15.24*f*
 High-intensity discharge (HID) lamps, **II**:40.33–40.36
 applications for, **II**:40.26*t*
 characteristics of, **II**:40.25*t*
 CMH, **II**:40.25*t*, 40.26*t*, 40.33, 40.36
 construction of, **II**:40.34*f*
 emission spectrum of, **II**:40.35*f*
 Hg, **II**:40.25*t*, 40.26*t*, 40.33, 40.36
 HPS, **II**:40.25*t*, 40.26*t*, 40.33
 MH, **II**:40.25*t*, 40.26*t*, 40.33, 40.35, 40.35*f*, 40.36
 High-intensity reciprocity failure, of photographic films, **II**:29.12
 Highlight color, in xerographic systems, **I**:34.12, 34.13*f*
 High-order harmonic generation, **II**:21.2, 21.2*f*
 High-performance miniature systems, **I**:22.5–22.8
 High-power diode lasers, **II**:19.19–19.23, 19.20*t*, 19.21*f*, 19.22*f*
 High-power laser arrays, **II**:19.26–19.30, 19.28*f*, 19.28*t*, 19.29*t*, 19.30*f*
 High-power lasers, **II**:19.24, 19.25*f*, 19.25*t*
 High-power semiconductor lasers, **II**:19.18–19.30
 arrays in, **II**:19.26–19.29, 19.28*f*, 19.28*t*, 19.29–19.30, 19.29*t*, 19.30*f*
 commercial diode, **II**:19.19–19.23, 19.20*t*, 19.21*f*, 19.22*f*
 future directions for, **II**:19.23–19.26, 19.25*f*, 19.25*t*, 19.26*t*, 19.27*f*
 mode-stabilized lasers with reduced facet intensity, **II**:19.18–19.19, 19.19*f*
 High-power spectrally controlled fiber lasers, **V**:25.29
 High-power strained QW lasers, **II**:19.26, 19.26*t*
 High-power ultrashort pulse technologies, for fiber lasers, **V**:25.32–25.33
 High-power USP oscillators, **V**:25.32–25.33
 High-pressure, short-arc xenon lamps, **II**:15.35*f*
 High-pressure enclosed arcs, **II**:15.24, 15.28–15.34
 capillary mercury-arc lamps, **II**:15.30–15.31, 15.31*f*
 compact-source arcs, **II**:15.31–15.34, 15.32*f*–15.35*f*
 Lucalox lamps, **II**:15.30, 15.31*f*
 mercury arcs, **II**:15.29, 15.30*f*
 multivapor arcs, **II**:15.29, 15.31*f*
 Uviarc, **II**:15.28–15.29, 15.29*f*, 15.30*f*
 High-pressure mercury-arc lamps, **II**:15.29*f*, 15.30*f*
 High-Q cavities, **IV**:9.12
 High-reflectance zones, of multilayers, **IV**:7.36–7.37, 7.37*f*
 High-reflectivity mirrors, **V**:41.7–41.8, 41.8*f*
 High-repetition short-pulse lasers, **IV**:11.26–11.27, 11.27*f*
 High-resistivity coatings, **IV**:6.56
 High-resolution (HR) acousto-optic deflectors, **V**:6.29*t*, 6.29
 High-resolution Doppler-free spectroscopy, **IV**:17.27
 High-speed cameras, **I**:25.21, 25.22*f*
 High-speed modulation, of semiconductor lasers, **II**:19.30–19.36, 19.31*f*, **II**:19.36*f*
 High-speed optical recording systems, **II**:19.3
 High-speed photoconductors, **II**:26.20–26.23, 26.21*f*–26.23*f*
 High-speed photodetectors, **II**:26.1–26.24
 about, **II**:26.3
 avalanche photodetectors, **II**:26.17–26.20, 26.18*f*, 26.20*f*, 26.21*f*
 photoconductors, **II**:26.20–26.23, 26.21*f*–26.23*f*
 pin photodiodes, **II**:26.10, 26.12–26.15
 resonant, **II**:26.15, 26.15*f*
 vertically illuminated, **II**:26.3, 26.4*f*, 26.5, 26.10, 26.12*f*, 26.12–26.13
 waveguide, **II**:26.13–26.14, 26.14*f*
 Schottky photodiodes, **II**:26.16, 26.16*f*, 26.17*f*
 speed limitations on, **II**:26.5–26.10
 carrier transit time, **II**:26.6*f*, 26.6–26.7
 carrier trapping, **II**:26.9, 26.9*f*
 diffusion current, **II**:26.8, 26.8*f*, 26.9
 packaging, **II**:26.9–26.10, 26.10*f*, 26.11*f*
 RC time constant, **II**:26.7*f*, 26.7–26.8
 structures of, **II**:26.3–26.5, 26.4*f*

- High-speed photographic films, **II**:30.18–30.20
- High-voltage power supply (HVPS), **II**:31.1, 31.9, 31.10*f*
- Hilbert space formulation, **I**:5.2
- Hilbert transforms, **IV**:2.8–2.9
- Hill gratings, **V**:17.5
- Hindle mounts, **II**:6.19, 6.19*f*
- HITRAN database, **V**:3.14, 3.22*f*, 3.22–3.23, 3.23*f*
- HITRAN-PC program, **V**:3.26
- Hobby-Eberly telescopes, **V**:5.2
- Hoffman modulation contrast microscopy, **I**:28.29
- Hole boring, **IV**:21.50, 21.50*f*
- Hole current, **II**:26.7
- Hole-accumulated photodiodes (HADs), **II**:32.4*f*, 32.8
- Holeburning:
 - optical, **V**:2.13, 2.14*f*
 - spatial, **V**:20.2
- Hollow cathode lamps, **II**:15.35, 15.37*t*–15.43*t*, 15.44*f*
- Hollow glass waveguides (HGWs), **V**:12.2*f*, 12.11–12.13, 12.12*f*
- Hollow lightpipes, **II**:39.30–39.31
- Hollow waveguides, **V**:12.2*f*, 12.2*t*, 12.3*t*, 12.10–12.13, 12.12*f*
- Hollow-core photonic crystal fibers:
 - attenuation in, **V**:11.20*f*, 11.20–11.22
 - birefringence of, **V**:11.17
 - and group velocity dispersion, **V**:11.18, 11.19*f*
 - Kerr effects for, **V**:11.23
 - and negative core-cladding index difference, **V**:11.14*f*, 11.14–11.15
- Holmium-doped fibers, **V**:25.23*t*, 25.25–25.26
- Holograms, computer-generated (*see* Computer-generated holograms)
- Holographic compensators, **II**:13.25
- Holographic inspection, **I**:33.16–33.19, 33.17*f*–33.18*f*, 33.20*f*–33.22*f*
- Holographic memory, **I**:33.24–33.25
- Holographic microscopes, **I**:28.42, 28.43, 28.43*f*
- Holographic optical elements (HOEs), **I**:33.13–33.16, 33.15*f*, 35.28–35.29, 35.29*f*; **IV**:12.31
- Holographic optics, **I**:33.13
- Holographic scanners, **I**:30.38–30.41, 30.40*f*–30.42*f*
- Holographic storage:
 - data, **IV**:12.37
 - photorefractive, **IV**:12.36–12.37
- Holography, **I**:33.1–33.25
 - electronic, **I**:33.9–33.14, 33.11*f*–33.13*f*
 - and holographic inspection, **I**:33.16–33.19, 33.17*f*–33.18*f*, 33.20*f*–33.22*f*
 - and holographic memory, **I**:33.24–33.25
 - and interferometry, **I**:33.4–33.9, 33.8*f*
 - and lithography, **I**:33.22*f*–33.24*f*, 33.22–33.24
 - optical elements for, **I**:33.13–33.16, 33.15*f*
 - principles of, **I**:33.2–33.4, 33.3*f*
 - real-time, **IV**:12.28–12.29, 12.29*f*, 12.30*f*
- Holtronic Technologies holographic system, **I**:33.23
- Homogeneity, in polymeric optics, **IV**:3.7
- Homogeneous broadening, **V**:20.1
 - emission-line, **II**:16.5, 16.6, 16.6*f*, 16.9
 - of lineshapes, **I**:10.7
 - spectral-line, **IV**:14.13
- Homogeneous coordinates, of collineation matrix, **I**:1.59
- Homogeneous decay, **IV**:11.10, 11.11*f*
- Homogenous index, **III**:19.1
- Homogeneous media, **I**:1.9, 1.23–1.26
- Homogeneous polarization elements:
 - defined, **I**:15.7
 - in Mueller matrices, **I**:14.25–14.26
- Homogeneous reflectors, **II**:39.39
- Homogeneous sources of light, **I**:5.11–5.12, 5.19
- Homogeneous temperature change, **II**:8.2–8.6, 8.3*t*, 8.4*t*, 8.5*f*, 8.6*f*
- Homogeneous-layer antireflection coatings, **IV**:7.16–7.23, 7.17*f*–7.19*f*, 7.20*t*–7.21*t*, 7.22*f*, 7.23*f*
- Homogeneously broadened systems, **IV**:11.18–11.19, 11.19*f*
- Homojunction lasers, **II**:19.4
- Homojunctions, LED, **II**:17.8–17.10, 17.9*f*–17.11*f*
- Hong-Ou-Mandel effect, **IV**:23.10, 23.14
- Hooke's law, **IV**:8.14, 8.21
- Hopkin's formula, **I**:5.13
- Horizontal cells (retina), **III**:2.9
- Horizontal half-wave linear retarders, Mueller matrices for, **I**:14.12*t*
- Horizontal illuminance, **II**:40.7, 40.18*f*
- Horizontal linear polarizers, **I**:14.9, 14.10*t*
- Horizontal quarter-wave linear retarders, Mueller matrices for, **I**:14.12*t*
- Horizontal scanning, in color CRTs, **III**:22.10–22.11, 22.11*f*
- Horizontally illuminated photodetectors, **II**:26.4*f*, 26.5

- Horizontal/vertical size and position, in color CRTs, **III**:22.12
- Horopter, **III**:2.40, 13.2, 13.8–13.9, 13.9*f*, 13.10*f*
- Hot bands, in molecular spectroscopy, **V**:2.5
- Hot cathode fluorescent lamps, **II**:40.32
- Hot filament sources, of x-ray tubes, **V**:54.10
- Hot mirrors, **IV**:7.58
- Hot-electron bolometers, **II**:24.29, 24.30, 24.30*f*, 24.31*f*
- Hottel strings, **II**:39.4, 39.14
- Houghton objectives, **I**:29.22–29.23
- Housings, lens, **IV**:3.15*f*, 3.15–3.16
- HRCam system, **V**:5.23
- HTF-1 glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.68*t*
- Huang-Rhys factor, **V**:2.15
- Huang-Rhys parameter, **I**:10.22
- Hub mounting, **II**:6.17, 6.18*f*
- Hubble telescope, **I**:29.8; **II**:11.4, 13.24
- Hue, **II**:40.5
- Hue cancellation, **III**:11.26–11.27
- Hue scaling, **III**:11.63, 11.64*f*
- Huesler alloy, **V**:63.28–63.29
- Hufnagel model, of atmospheric turbulence, **V**:3.29
- Hufnagel-Valley model, of atmospheric turbulence, **V**:3.30, 5.7, 5.8
- Hughes Airborne Optical Adjunct Coating, **IV**:6.49
- Human eye, **II**:30.15–30.16, 30.16*f*, 34.6; **III**:1.1–1.45, 1.3*f*
- aberrations in, **III**:1.3 [*see also* Aberrations (in human eye)]
 - accommodation response, **III**:1.29–1.36
 - accuracy of, **III**:1.32–1.34
 - age-dependent changes in, **III**:1.35–1.36
 - application to instrumentation, **III**:1.34–1.35
 - dynamics of, **III**:1.31–1.32
 - stability of, **III**:1.32
 - vergence input, **III**:1.34
 - active optical limiting by, **IV**:13.1
 - age-related changes in, **III**:1.5, 1.7, 14.4–14.14 (*see also* Adaptive optics)
 - accommodation, **III**:1.35*f*, 1.35–1.36
 - accommodation and presbyopia, **III**:14.8, 14.9*f*
 - anterior and posterior chambers, **III**:14.5, 14.6
 - cornea, **III**:14.5, 14.6*f*
 - eye size, **III**:14.11
- Human eye, age-related changes in (*Cont.*):
- fluorescence, **III**:1.21
 - index of diffusion for, **III**:1.23
 - lens, **III**:14.7–14.8
 - and level of aberration, **III**:1.18
 - pupil, **III**:14.6–14.7, 14.7*f*
 - pupil diameter, **III**:1.8, 1.18
 - retina, **III**:14.9–14.11
 - retinal reflectance, **III**:1.11
 - RMS wavefront error, **III**:1.15–1.17, 1.16*t*
 - scattering, **III**:1.20
 - tears, **III**:14.4–14.5
 - transmittance, **III**:1.9
 - transparency/cataract, **III**:14.8
 - ametropia, **III**:1.6–1.7
 - cones in, **II**:36.8, 36.8*f*, 36.9*f*
 - cornea, **III**:16.3, 16.4, 16.4*f*
 - damage from light exposure, **III**:5.18*f*, 5.18–5.19
 - damage thresholds for, **IV**:13.3, 13.3*f*
 - depth-of-focus, **III**:1.8, 1.28–1.29, 1.30*f*
 - models of, **III**:1.36–1.38
 - paraxial, **III**:1.36–1.37, 1.37*f*
 - wide-angle, **III**:1.38
 - monochromatic ocular aberrations, **III**:1.4, 1.14–1.19
 - off-axis, **III**:1.18*f*, 1.18–1.19
 - on the visual axis, **III**:1.15–1.18
 - movements of, **III**:1.42–1.45
 - characteristics of, **III**:1.43–1.44
 - stability of fixation, **III**:1.44, 1.45*f*
 - ocular parameters, **III**:1.4*f*, 1.4–1.6, 1.5*f*
 - ocular radiometry, **III**:1.11–1.12
 - optical components of, **III**:2.2
 - optics of, **III**:16.2–16.3, 16.3*f*
 - pupil diameter, **III**:1.8–1.9
 - refractive elements in, **III**:12.3
 - retinal illuminance, **III**:1.11–1.12
 - retinal image quality, **III**:1.12–1.28
 - in aberration-free eye, **III**:1.12–1.14
 - calculation from aberration data, **III**:1.21–1.22
 - chromatic aberration, **III**:1.19*f*, 1.19–1.20
 - comparison between methods, **III**:1.23
 - effects of aberration correction, **III**:1.25–1.26
 - intraocular scattered light, **III**:1.20–1.21
 - lenticular fluorescence, **III**:1.21
 - monochromatic ocular aberrations, **III**:1.14–1.19

- Human eye, retinal image quality (*Cont.*):
 observed optical performance,
III:1.23–1.25
 ophthalmoscopic (double-pass) methods,
III:1.22–1.23
 psychophysical comparison method,
III:1.22
 and pupil diameter, **III**:1.8
 on the visual axis, **III**:1.21–1.28
 with visual instruments, **III**:1.27–1.28
 retinal reflectance, **III**:1.11
 rods in, **II**:36.8*f*, 36.8–36.10, 36.9*f*
 size of, **III**:14.11
 stereopsis, **III**:1.38–1.42
 aniseikonia, **III**:1.41–1.42
 stereoscopic and related instruments,
III:1.41
 tolerances in binocular instrumentation,
III:1.41–1.42
 Stiles-Crawford effect, **III**:1.10–1.11
 structure of, **III**:21.2–21.3
 transmittance, **III**:1.9–1.11
 visual vs. optical axis orientation in, **III**:1.3
 wavelengths detectable by the, **II**:36.8*f*,
 36.8–36.10, 36.9*f*
 (*see also specific topics*)
 Human hair, light guide effect in, **III**:8.24–8.26
 Human vision and electronic imaging,
III:24.1–24.11
 analysis of image features, **III**:24.6–24.9
 attention and region of interest, **III**:24.7
 Digital Libraries applications, **III**:24.7–24.8
 user interface design, **III**:24.8–24.9
 visualization, **III**:24.8
 information sources for, **III**:24.11
 perception of imaging artifacts, **III**:24.2–24.6
 early color vision, **III**:24.5
 embedding digital watermarks, **III**:24.5
 image quality and compression,
III:24.3–24.4
 limitations of early vision models,
III:24.5–24.6
 rendering and halftoning, **III**:24.4
 target detection in medical images, **III**:24.5
 representation of aesthetic and emotional
 characteristics, **III**:24.9–24.10
 color, art, and emotion, **III**:24.10
 image quality, **III**:24.9
 virtual reality and presence, **III**:24.9
 Human visual system (HVS), color in, **I**:26.18
 Humidity specifications, for lenses, **II**:4.10
 Hund's rule, **I**:10.12
 Hurst dimension, **I**:8.8
 Huygens wavelets, **I**:5.16
 Huygens-Fresnel approximation, **V**:3.31–3.33
 Huygens-Fresnel construction (diffraction),
I:3.4–3.13
 Babinet principle, **I**:3.9–3.11
 for double refraction in calcite, **I**:13.5
 Fresnel zones in, **I**:3.4–3.6
 and light from circular apertures and disks,
I:3.6–3.9
 zone plates in, **I**:3.11–3.13, 3.12*f*
 Hybrid arrays, pyroelectric, **II**:28.11*f*,
 28.11–28.12, 28.12*f*
 Hybrid FPAs:
 direct readout architectures of, **II**:33.15–33.18
 electronically scanned staring FPAs,
II:33.16–33.17
 output circuits, **II**:33.18
 TDI scanning FPAs, **II**:33.17, 33.17*f*
 X-Y addressing and clock generation,
II:33.16
 input circuits of, **II**:33.18–33.23, 33.19*t*, 33.20*f*
 buffered direct injection, **II**:33.19*t*,
 33.20*f*, 33.21
 capacitive transimpedance amplifier,
II:33.19*t*, 33.20*f*, 33.22–33.23
 chopper-stabilized BDI, **II**:33.19*t*, 33.20*f*,
 33.21–33.22
 direct detector integration, **II**:33.18,
 33.19*t*, 33.20*f*
 direct injection, **II**:33.18–33.21, 33.19*t*,
 33.20*f*, 33.21*f*
 gate modulation, **II**:33.19*t*, 33.20*f*, 33.22
 thermal expansion match in, **II**:33.14
 Hybrid mode locking, **V**:20.17
 Hybrid network topologies, for WDM networks,
V:21.7
 Hybrid organic-inorganic composites, glasses,
 and sol-gels, **IV**:12.27*t*
 Hybrid reflectors (*see* Faceted reflectors)
 Hyde maxim, **II**:3.22
 Hydrogel contact lenses (*see* Soft contact lenses)
 Hydrogen, Bohr's theory of, **I**:10.3
 Hydrogen arc lamps, **II**:15.49, 15.53*f*
 Hydrogen (H⁻) ions, negative, **V**:2.3
 Hydrogen loading, **V**:17.2
 Hydrologic optics, **IV**:1.3 (*see also* Water)
 Hydrophilic, defined, **III**:16.1

- Hydrostatic pressure, **IV**:5.66*t*
- Hyperacuity, **III**:2.34, 2.35, 4.13–4.16, 4.14*f*
 across ages, **III**:14.22
 defined, **III**:4.1
 and superresolution, **III**:4.15
- Hyperbolic cumulative size distribution,
IV:1.15
- Hyperboloid-hyperboloid (HH) grazing
 incidence x-ray telescopes, **V**:44.10–44.12,
 44.11*f*
- Hyperfocal distance, **I**:1.85
- Hypermetropia (hyperopia), **III**:1.6, 16.5,
 16.5*f*, 16.8
- Hyperopia, **III**:1.18, 12.4
 defined, **III**:13.2
 and focus of collimated light, **III**:12.3*f*
- Hyper-Raman scattering, **IV**:15.2*t*, 15.3, 15.3*t*
- Hyperspectral imaging, **IV**:17.27–17.28
- Hypo (film fixer), **II**:29.5
- Hypocotyl, **III**:8.2, 8.26, 8.27*f*
- Hysteresis instability, of metals, **IV**:4.10
- Hysteresis loops, of optical disks, **I**:35.27,
 35.27*f*, 35.28
- IBM, **V**:23.1, 23.2
- IBM Black, **IV**:6.56
- Ice, in standard atmosphere, **V**:3.6, 3.42, 3.42*f*
- Ideal imaging (term), **I**:1.28, 1.38
- Ideal mode-locked lasers, **II**:20.7
- Ideal observer, in judgment experiments,
III:3.6
- Ideal receivers, **V**:9.9
- Ideal thermal detectors, **II**:28.2–28.3, 28.3*f*
- Ideality factor, of semiconductor diodes,
V:13.69
- Ideal-observer models (color vision), **III**:11.26
- Ideal-observer theory, **III**:2.16–2.19, 2.19*f*, 2.24
- Identification tasks, **III**:2.15, 2.16
- Identity matrix, **I**:14.8
- IEEE (Institute of Electrical and Electronics
 Engineers) standards, **V**:23.2, 23.7
- Ilford Photo Corporation, **II**:29.25
- Illinois Institute of Technology, **IV**:6.35
- Illuminance, **II**:37.4*t*, 37.5, 37.5*f*, 39.2*t*
 defined, **II**:34.11, 34.40, 40.1; **III**:23.2
 guidelines on levels of, **II**:40.7*t*
 and lighting design, **II**:40.7
 and luminance, **II**:37.9, 37.9*f*
 measurement of, **III**:5.17, 5.18
 retinal, **II**:34.40–34.42; **III**:1.11–1.12
- Illuminance (*Cont.*):
 SI units of, **III**:5.3
 uniformity of, **II**:40.7
 unit conversions for, **II**:36.7*t*, 36.8*t*
 units of, **II**:34.43
- Illuminance meters, **II**:34.42, 40.51, 40.52*f*
- Illuminants:
 metamerism for, **III**:10.38
 reflected from surfaces, **III**:10.32–10.36
- Illuminated ceilings, **II**:40.13*f*
- Illuminated eyepieces, **II**:12.12, 12.12*f*
- Illuminated objects (in stray light suppression),
II:7.5, 7.5*f*, 7.6*f*
- Illuminating Engineering Society (IES),
II:40.19
- Illumination:
 critical, **I**:28.7
 epi-, **I**:28.7*f*, 28.7–28.9, 28.8*f*
 guidelines on, **II**:40.7*t*
 for input/output scanning, **I**:30.14
 in nonimaging objects, **II**:39.1
 trans-, **I**:28.5–28.7, 28.6*f*
 (*see also* Uniform illumination, of
 nonimaging optics)
- Illumination area, surface scattering and,
I:8.12
- Illumination Engineering Society of North
 America (IESNA), **II**:36.2, 36.3, 37.11,
 40.2, 40.7*t*
- Illumination subsystem, of nonimaging optics,
II:39.22
- Image(s), **I**:1.26
 about known rays, **I**:1.43–1.44, 1.44*f*
 aerial, **I**:1.26
 from cameras, **I**:25.5, 25.7
 of extended objects, **I**:1.27
 in focal Gaussian lenses, **I**:1.51
 formation of, **I**:6.9*f*, 6.9–6.10, 17.5*f*,
 17.5–17.8, 17.6*f*
 latent, **I**:34.1–34.4, 34.2*f*–34.3*f*, 34.2–34.4
 long-exposure, **V**:4.3–4.7
 medial, **I**:29.11, 29.37
 from Nomarski microscope, **V**:46.2
 point, **I**:1.85–1.92, 1.88*f*
 of points, **I**:1.27
 received, **I**:1.26
 recorded, **I**:1.26
 short-exposure, **V**:4.3, 4.31*f*–4.34*f*,
 4.31–4.35, 4.35*t*
 in systems of revolution, **I**:1.42

- Image(s) (*Cont.*):
 two-dimensional, **I**:11.12–11.17, 11.13*f*
 virtual, **I**:29.38
 x-ray, **V**:31.1–31.4
- Image (scophony) AO modulators, **V**:6.34–6.35
- Image artifacts:
 and perceived fidelity, **III**:24.5
 perception of (*see* Perception, of imaging artifacts)
- Image compression, early techniques for, **III**:24.4
- Image derotation, **I**:30.36
- Image dissectors, **II**:39.21, 39.21*f*
- Image distance, **I**:18.3
- Image distribution, for cameras, **I**:25.7
- Image enhancement, coherent optical, **I**:11.14–11.17
- Image erectors, **I**:18.10
- Image features in vision, **III**:24.6–24.9
 attention and region of interest, **III**:24.7
 Digital Libraries applications, **III**:24.7–24.8
 user interface design, **III**:24.8–24.9
 visualization, **III**:24.8
- Image formation, **III**:2.2–2.4
- Image height, **I**:1.27, 18.4; **II**:1.4
- Image intensifiers (IIs), **II**:31.7–31.18,
 31.8*f*–31.10*f*
 defined, **II**:31.8
 input window/photocathode assemblies for,
II:31.10–31.12, 31.11*f*, 31.12*f*
 MCP IIs, **II**:31.7, 31.9, 31.9*f*, 31.10*f*
 and microchannel plates, **II**:31.12*f*,
 31.12–31.14, 31.13*f*, 31.13*t*
 phosphor screen assemblies for,
II:31.14–31.16, 31.14*t*, 31.15*f*
 proximity-focused MCP IIs, **II**:31.16–31.18,
 31.17*t*, 31.18*f*, 31.19*f*
- Image inversion, **I**:19.2
- Image irradiance, **II**:4.7
- Image lag, **II**:32.6
- Image overlap, **III**:13.4
- Image plane, **I**:1.27
- Image point, **I**:1.27
- Image processing, **II**:8.12 (*see also* Analog optical signal and image processing)
- Image quality, **II**:4.6–4.7
 for head-mounted displays, **III**:25.2, 2
 5.5–25.7, 25.7*f*
 in human eye (*see* Retinal image quality)
 subjective approach to, **III**:24.9
- Image receptors:
 in medical imaging, **V**:31.4, 31.4*f*
 in xerographic systems, **I**:34.1
- Image reversion, **I**:19.2
- Image rotation, **I**:30.35–30.36
- Image sampling, by photoreceptors, **III**:2.4–2.9
- Image sensors, **II**:32.2–32.12, 32.3*f*, 32.21–32.34
 antiblooming in, **II**:32.9, 32.10*f*
 area arrays of, **II**:32.24–32.32, 32.25*t*
 CCD performance, **II**:32.32
 frame transfer CCDs, **II**:32.26–32.28,
 32.27*f*, 32.28*f*
 interline transfer CCDs, **II**:32.28–32.32,
 32.29*f*–32.31*f*
 MOS, **II**:32.25–32.26, 32.26*f*
 color imaging with, **II**:32.32–32.34, 32.33*f*,
 32.34*f*
 dark current in, **II**:32.10–32.12, 32.11*f*
 junction photodiodes, **II**:32.3–32.6, 32.4*f*, 32.6*f*
 linear arrays of, **II**:32.21–32.24, 32.22*f*, 32.23*f*
 MOS capacitors, **II**:32.7–32.8
 photoconductors, **II**:32.8–32.9
 pinned photodiodes, **II**:32.8
- Image size, **II**:1.4
- Image space, **I**:1.26, 1.83*f*, 1.83–1.84
- Image space numerical aperture, **I**:1.79
- Image specifications, for lenses, **II**:4.3, 4.6–4.8,
 4.8*f*
- Image structure, of photographic systems, **II**:29.17
- Image wander, **V**:4.3 (*see also* Angle of arrival)
- Imaged critical objects, **II**:7.4, 7.5*f*
- Image-forming cone (bundle), **I**:1.74
- Image-forming rays, **I**:1.74
- Image-intensified (II) electronic imaging,
II:31.1–31.30
 about, **II**:31.2–31.3, 31.3*f*
 applications for, **II**:31.27–31.30
 active imaging, **II**:31.29–31.30
 day/night cameras, **II**:31.28–31.29
 mosaic II SSA cameras, **II**:31.29
 optical multichannel analyzers,
II:31.27–31.28
 range gating and LADAR, **II**:31.28
 image-intensifier modules of, **II**:31.7–31.18,
 31.8*f*–31.10*f*
 input window/photocathode assemblies,
II:31.10–31.12, 31.11*f*, 31.12*f*
 microchannel plates, **II**:31.12*f*,
 31.12–31.14, 31.13*f*, 31.13*t*

- Image-intensified (II) electronic imaging,
 image-intensifier modules of (*Cont.*):
 phosphor screens, **II**:31.14–31.16, 31.14*t*,
 31.15*f*
 proximity-focused MCP IIs, **II**:31.16–31.18,
 31.17*t*, 31.18*f*, 31.19*f*
 optical interface of, **II**:31.3–31.7
 considerations, **II**:31.6*f*, 31.6–31.7
 photometry and camera lens,
II:31.5–31.6
 quantum limited imaging conditions,
II:31.3–31.4
 radiometry, **II**:31.4–31.5
 self-scanned arrays, **II**:31.19*f*, 31.19–31.27
 electron-bombarded, **II**:31.23*f*, 31.23–31.27,
 31.24*f*
 fiberoptic-coupled, **II**:31.20*f*, 31.20–31.22,
 31.21*f*, 31.21*t*
 lens-coupled, **II**:31.22*f*, 31.22–31.23
 Image-intensified self-scanned arrays (II SSAs),
II:31.19*f*, 31.19–31.27
 for active imaging, **II**:31.29–31.30
 camera for, **II**:31.5–31.6
 electron-bombarded, **II**:31.23*f*, 31.23–31.27,
 31.24*f*
 fiberoptic-coupled, **II**:31.20*f*, 31.20–31.22,
 31.21*f*, 31.21*t*
 lens-coupled, **II**:31.22*f*, 31.22–31.23
 Imaging:
 and atmospheric turbulence, **V**:3.34
 conscopic vs. orthoscopic, **I**:28.8–28.9
 diffractive, **V**:41.9–41.10, 41.10*f*
 fluorescence, **I**:28.23
 in geometrical optics, **I**:1.26–1.31
 ghost, **IV**:23.13
 with grazing incidence optics, **V**:44.3–44.18
 with grazing-incidence neutron optics,
V:64.3–64.7, 64.4*f*–64.7*f*
 hyperspectral, **IV**:17.27–17.28
 ideal, **I**:1.17, 1.28–1.29, 1.38
 medical, **V**:31.1–31.10
 applications of, **V**:31.9, 31.10
 digital displays, **V**:31.8*f*–31.10*f*,
 31.8–31.9
 digital tomosynthesis, **V**:31.7–31.8
 and inverse Compton x-ray sources,
V:59.3–59.4
 and polycapillary x-ray optics,
V:53.14–53.16, 53.15*f*–53.16*f*
 radiography, **V**:31.1–31.4, 31.2*f*–31.4*f*
 Imaging, medical (*Cont.*):
 tomography, **V**:31.1, 31.5, 31.5*f*–31.7*f*,
 31.5–31.7
 x-ray detectors for, **V**:61.2
 in microscopes, **I**:28.44–28.54
 aperture-scanning microscopy,
I:28.53–28.54
 confocal microscopy, **I**:28.49*f*,
 28.49–28.51, 28.51*f*
 fluorescent microscopy, **I**:28.48–28.49
 light field microscopy, **I**:28.53
 polarizing microscopes, **I**:28.44*f*,
 28.44–28.48, 28.46*f*, 28.47*f*
 structured illumination, **I**:28.52
 molecular, **V**:32.1
 monochromatic, **V**:53.16–53.17, 53.17*f*
 multi-energy, **V**:54.9–54.10
 Newtonian equation for, **I**:17.8
 nuclear, **V**:53.17, 53.18, 53.18*f*
 in polarimeters, **I**:15.6
 pupil, **I**:1.76
 quantum, **IV**:23.13–23.14
 with refractive x-ray lenses, **V**:37.6*f*,
 37.6–37.7, 37.7*f*
 scatter rejection in, **V**:53.14–53.16,
 53.15*f*–53.16*f*
 SPECT, **V**:32.1–32.3
 and spectroscopy, **V**:5.19*f*, 5.19–5.21, 5.20*f*
 stigmatic, **I**:1.29–1.31, 1.30*f*
 thermal, **V**:12.3*t*
 through atmospheric turbulence,
V:4.1–4.37
 aberration variance and approximate
 Strehl ratio for, **V**:4.27–4.28, 4.28*f*
 and adaptive optics, **V**:4.35–4.36
 angle of arrival fluctuations, **V**:4.23–4.26,
 4.25*f*, 4.26*f*
 and covariance and variance of expan-
 sion coefficients, **V**:4.20–4.22, 4.21*t*,
 4.22*f*–4.23*f*
 Kolmogorov turbulence and atmospheric
 coherence length, **V**:4.7–4.10,
 4.8*f*, 4.9*f*
 long-exposure images, **V**:4.3–4.7
 modal correction of turbulence,
V:4.28–4.30, 4.29*t*, 4.30*f*
 and modal expansion of aberration
 function, **V**:4.17*f*–4.18*f*, 4.17–4.20,
 4.19*t*, 4.20*t*
 and resolution of telescopes, **V**:4.2–4.3

- Imaging, through atmospheric turbulence (*Cont.*):
 short-exposure image, **V**:4.31*f*–4.34*f*,
 4.31–4.35, 4.35*t*
 and systems with annular pupils,
V:4.10–4.16, 4.11*f*–4.15*f*, 4.15*t*
 time-gated, **IV**:15.42, 15.43, 15.44*f*
- Image-space scanners, **I**:30.18–30.23
- Imaging detectors (x-ray), **V**:61.1–61.8
 CCD detectors, **V**:61.7–61.8, 61.8*f*
 flat panel detectors, **V**:61.3–61.7, 61.4*f*, 61.6*t*,
 61.7*f*
 geometries for and classifications of,
V:61.1–61.3, 61.2*f*
- Imaging plates, in neutron optics, **V**:63.34
- Imaging science, **III**:24.1
- Immersion lenses, **I**:17.11
- Impact ionization coefficient, **V**:13.72
- Impedance:
 in amplifiers, **II**:27.10–27.11
 in photodetectors, **II**:24.19–24.20
- Imperfect polarizers, **I**:13.33
- Impermeability, dielectric, **I**:21.9
- Impulse response function (IRF),
III:14.19–14.21, 14.20*f*
- Impurity band conduction (IBC), **II**:33.7
- Impurity magnetoabsorption, **IV**:5.51*f*,
 5.51–5.52
- Impurity-related absorption, **IV**:5.37–5.39,
 5.38*f*, 5.39*f*
- Impurity-related vibrational optic effects,
IV:5.17–5.20, 5.18*f*–5.19*f*, 5.20*t*, 5.21*f*
- Incandescence, **II**:40.25
- Incandescent sources (of radiation), **II**:40.25,
 40.27–40.30
 calibration of, **II**:34.31
 characteristics of, **II**:40.25*t*
 elements of, **II**:40.29*f*
 nongaseous, **II**:15.15–15.21
 comparisons of, **II**:15.19, 15.19*f*
 gas mantle, **II**:15.17, 15.18, 15.19*f*
 globar, **II**:15.17, 15.18*f*
 Nernst glower, **II**:15.14, 15.15, 15.17, 15.17*f*
 quartz-envelope lamps, **II**:15.20, 15.21
 tungsten-filament lamps, **II**:15.19, 15.20,
 15.20*f*–15.22*f*
 tungsten emissivity in, **II**:40.28*f*
- Incidence angle, constant, **V**:38.2
- Incident power measurement, in scatterometers,
V:1.14–1.15
- Inclination factor, **I**:3.5
- Incoherent arrays, scattering by, **I**:7.2–7.3
- Incoherent light sources, **I**:5.12, 5.18–5.19;
III:6.1
- Incoherent processing, of discrete signals,
I:11.17–11.20, 11.18*f*, 11.19*f*
- Incoherent radiation, **I**:30.2, 30.26, 30.27
- Incoherent scattering, **I**:9.2–9.5; **V**:26.7, 31.2,
 63.7, 63.8
- Incoherent targets, pupil size and, **III**:6.7–6.9,
 6.8*f*, 6.9*t*
- Incomplete polarimeters, **I**:15.4
- Incomplete sample-measuring polarimeters,
I:15.16–15.17
- Increment contrast, **III**:2.31, 2.32
- Increment thresholds (color vision), **III**:11.16*f*
- Incremental cone-excitation space, **III**:11.2,
 11.31, 11.32
- Index contrast, 3D photonic crystals and,
IV:9.4
- Index ellipsoid:
 of crystals and glasses, **IV**:2.18–2.19, 2.19*f*
 of electro-optic modulators, **V**:7.4–7.7,
 7.5*f*–7.6*f*
- Index ellipsoid equation, **I**:21.9–21.10
- Index grating, **IV**:12.1–12.3, 12.2*f*
- Index of absorption, **IV**:2.8
- Index of refraction (*see* Refractive index)
- Index-guided laser stripes, **V**:13.6
- Index-guided lasers, **II**:19.8, 19.27, 19.28*f*
- Indicator lamps, LED, **II**:17.27*f*, 17.29*f*,
 17.29–17.30
- Indirect (nonvertical) absorption transitions,
IV:5.22–5.24, 5.24*f*–5.25*f*
- Indirect excitons, **IV**:5.29
- Indirect glare, **II**:40.9
- Indirect interband transitions, **IV**:8.29*f*,
 8.29–8.30
- Indirect lighting, **II**:40.14, 40.15, 40.15*f*, 40.16*f*,
 40.46*f*
- Indirect modulation, in OTDM networks,
V:20.17*f*, 20.17–20.18
- Indirect semiconductors, **II**:17.4, 17.4*f*, 17.5*f*,
 17.6
- Indirect-conversion flat panel detectors, **V**:61.4,
 61.4*f*, 61.4–61.6
- Indium antimonide (InSb) hot-electron
 bolometers, **II**:24.29, 24.30, 24.30*f*, 24.31*f*
- Indium antimonide (InSb) intrinsic
 photovoltaic detectors, **II**:24.80–24.83,
 24.82*f*, 24.83*f*

- Indium arsenide (InAs) photovoltaic detectors, **II**:24.75, 24.77*f*–24.79*f*, 24.77–24.78
- Indium gallium arsenic phosphide (InGaAsP) laser material system, **II**:19.7
- Indium gallium arsenide (InGaAs) detectors, **II**:24.65–24.70, 24.66*f*–24.69*f*
- Indium gallium arsenide (InGaAs) photodetectors, **II**:25.10, 25.10*t*
- Indium gallium arsenide (InGaAs) photodiodes, **II**:24.66*f*–24.69*f*, 24.66–24.70, 34.31
- Indium gallium nitride (InGaN) HB-LEDs, **II**:18.2*f*
- Indium phosphide (InP), **IV**:12.21
- Indium phosphide (InP) laser material system, **II**:19.7
- Induced transparency, **IV**:21.52
- Induced-transmission filters, **IV**:7.83, 7.88*f*
- Induction lamps (ILs), **II**:40.36–40.37
- Inductive pickup noise, **II**:27.5, 27.6*f*
- Inductively coupled plasma (ICP), **II**:18.3
- Industrial lighting, **II**:40.60–40.61, 40.61*f*
- Inelastic optical processes, **V**:3.21–3.22
- Inelastic scattering, **V**:21.20, 63.3
of light, **IV**:5.76*f*, 5.76–5.83, 5.78*f*–5.82*f*
and polarization, **IV**:1.47–1.49, 1.48*f*, 1.49*f*
- Inertial confinement fusion (ICF), **IV**:21.54, 21.55*f*
- “Infant mortality period,” **II**:17.26, 17.26*f*
- Infants:
corneal endothelium in, **III**:14.5
lenses in, **III**:14.8
size of eye globe in, **III**:14.11
- Infectious film developers, **II**:29.5
- In-fiber devices, for photonic crystal fibers, **V**:11.27, 11.28
- InfiniBand standard, **V**:23.8, 23.8*t*
- Infinitely distant objects, systems with, **II**:1.5–1.6, 1.6*f*
- Influence function, **III**:15.1
- InFOCUS/SUMIT balloon payloads, **V**:47.10
- Information capacity, of photographic systems, **II**:29.24
- Information theory, visual resolution and, **III**:4.15–4.16
- Information-processing model, **III**:2.15*f*, 2.15–2.16
- Infrared, **IV**:6.26, 6.26*f*, 6.28, 6.28*f*, 6.48, 6.48*f*
- Infrared cataract, **III**:7.7, 7.10
- Infrared detector arrays, **II**:33.1–33.31
about, **II**:33.3–33.4
applications for, **II**:33.4
current status of, **II**:33.28*f*, 33.28–33.30, 33.29*f*, 33.29*t*
future trends and technology directions of, **II**:33.30*f*, 33.30–33.31, 33.31*f*
hybrid FPAs, **II**:33.14–33.23
detector interface input circuit, **II**:33.18–33.23, 33.19*t*, 33.20*f*, 33.21*f*
hybrid readout, **II**:33.15–33.23
readout, **II**:33.17*f*
thermal expansion match in, **II**:33.14
monolithic FPAs, **II**:33.10–33.14
direct-charge-injection silicon FPAs, **II**:33.11*f*, 33.13
microbolometer FPAs, **II**:33.13–33.14
MIS photogate FPAs, **II**:33.10–33.11, 33.11*f*, 33.12*f*
scanning and staring, **II**:33.14
silicon FPAs, **II**:33.11*f*, 33.11–33.13, 33.12*f*
operating principles of, **II**:33.7–33.10, 33.8*f*, 33.9*f*
performance of, **II**:33.23–33.28
detectivity, **II**:33.23–33.24
minimum resolvable temperature, **II**:33.27*f*, 33.27–33.28
NE ΔT , **II**:33.24–33.27, 33.25*f*, 33.26*f*
percentage of BLIP, **II**:33.24
scanning and staring, **II**:33.6*f*, 33.6–33.7
spectral bands for, **II**:33.4–33.5, 33.6*f*
- Infrared detectors:
gallium-doped germanium, **II**:24.100
gallium-doped silicon, **II**:24.95, 24.95*f*, 24.96, 24.96*f*
- Infrared (IR) dipole active modes, **IV**:8.16–8.18, 8.17*f*
- Infrared emitting diodes (IREDs), **I**:25.14, 25.14*f*
- Infrared film, **II**:30.22
- Infrared filters, **II**:40.12
- Infrared interferometry, **II**:13.25
- Infrared LED chips, **II**:17.8, 17.9, 17.9*f*
- Infrared (IR) optical fibers, **V**:12.1–12.13
applications, **V**:12.13
categories and properties of, **V**:12.1–12.3, 12.2*f*, 12.2*t*, 12.3*t*
crystalline, **V**:12.2*t*, 12.3*t*, 12.7–12.10, 12.8*f*, 12.10*f*

- Infrared (IR) optical fibers (*Cont.*):
 heavy-metal oxide glass in, **V**:12.2*t*–12.4*t*,
 12.3–12.7, 12.5*f*–12.7*f*
 in hollow waveguides, **V**:12.2*t*, 12.3*t*,
 12.10–12.13, 12.12*f*
- Infrared photodetectors, **II**:25.12, 25.15
- Infrared (IR) radiation, **II**:34.6, 40.41
 damage from, **III**:7.2
 exposure limits for, **III**:7.10–7.11
 extreme, **II**:25.2
 far, **II**:24.3, 25.2
 forward looking, **II**:33.4
 long-wavelength, **II**:24.3, 33.3–33.5, 33.6*f*
 medium-wavelength, **II**:24.3, 25.2, 33.3,
 33.5, 33.6*f*
 near, **II**:24.3, 25.2
 short-wavelength, **II**:24.3, 33.3, 33.5
 single-order plates and, **I**:13.47
 very long-wavelength, **II**:24.3
- Infrared radiometry, standards for,
II:15.11–15.12, 15.12*f*
- Infrared (IR) region:
 absorption in, **IV**:5.19*f*
 all-dielectric reflectors for, **IV**:7.39, 7.40*f*
 multilayer reflectors for far-, **IV**:7.52, 7.52*f*
- Infrared (IR) suppressing filters, **IV**:7.58, 7.58*f*
- Inherent optical properties (IOPs), of water,
IV:1.4, 1.5*t*, 1.9–1.12, 1.10*f*
- Inhibition, **III**:2.25*f*, 2.25–2.27
- Inhomogeneous antireflection coatings,
IV:7.23–7.26, 7.24*f*–7.26*f*
- Inhomogeneous broadening:
 emission-line, **II**:16.5, 16.6*f*
 spectral-line, **IV**:14.14
- Inhomogeneous media, **II**:39.22
- Inhomogeneous (heterogeneous) media, rays
 in, **I**:1.9, 1.18–1.22
- Inhomogeneous optics, **I**:24.1 [*see also*
 Gradient index (GRIN) optics]
- Inhomogeneous polarization elements,
I:14.25–14.26, 14.26*f*, 15.7, 15.20
- Inhomogeneous reflectors, **II**:39.39
- Injection:
 electrical, **V**:13.4, 13.5, 13.45, 13.45*f*
 in LEDs, **V**:13.37
- Injection molding, of polymers, **IV**:3.2,
 3.12–3.13
- Injection seeding, optical clock recovery and,
V:20.21–20.22
- Injection-locked lasers, **II**:19.37*f*, 19.38
- Injuries, radiation, **III**:7.3–7.7, 7.5*f*
 cataract, **III**:7.5–7.6, 7.6*f*
 droplet keratopathies, **III**:7.6, 7.7
 infrared cataract, **III**:7.7
 mechanisms of, **III**:7.2–7.3
 action spectra, **III**:7.2, 7.4*f*
 exposure duration and reciprocity,
III:7.2–7.3
 photokeratitis, **III**:7.4
 photoretinitis, **III**:7.7
 pterygium, **III**:7.6, 7.7
- In-line gun CRTs, **III**:22.5
- In-line semiconductor optical amplifiers
 (SOAs), **V**:19.24
- Inner ionization, of cluster, **IV**:21.31–21.32,
 21.32*f*
- Inner scale of turbulence, **V**:4.7
- Inorganic crystals, **V**:7.33
- Inorganic particles, in water, **IV**:1.14–1.15
- In-plane coupling, **IV**:9.10–9.11
- In-plane switching (IPS) cells, **V**:8.16, 8.26, 8.28*f*
- Input circuits, of hybrid FPAs, **II**:33.18–33.23,
 33.19*t*, 33.20*f*
 buffered direct injection, **II**:33.19*t*, 33.20*f*,
 33.21
 capacitive transimpedance amplifier,
II:33.19*t*, 33.20*f*, 33.22–33.23
 chopper-stabilized BDI, **II**:33.19*t*, 33.20*f*,
 33.21–33.22
 direct detector integration, **II**:33.18, 33.19*t*,
 33.20*f*
 direct injection, **II**:33.18–33.21, 33.19*t*,
 33.20*f*, 33.21*f*
 gate modulation, **II**:33.19*t*, 33.20*f*, 33.22
- Input optics, **II**:38.7
- Input planes, translations of, **I**:1.68
- Input saturation power, for fiber amplifiers,
V:14.3
- Input standards, for color LCDs, **III**:22.39
- Input windows, of image intensifiers, **II**:31.9*f*,
 31.9–31.12, 31.11*f*, 31.12*f*
- Input/output scanning, **I**:30.2, 30.4–30.6,
 30.25–30.34
 objective, preobjective, and postobjective,
I:30.28–30.29, 30.29*f*, 30.30*f*
 objective optics, **I**:30.30–30.33, 30.32*f*–30.33*f*
 power density and power transfer of,
I:30.25–30.28, 30.27*f*, 30.28*f*
 resolution of, **I**:30.8–30.14, 30.10*f*, 30.10*t*,
 30.11*t*, 30.12*f*–30.13*f*

- Insertion devices:
 for synchrotron radiation sources,
 V:55.9–55.16, 55.10*f*, 55.12*f*, 55.13*f*
- Insertion loss, V:7.36, 13.53, 18.1
- Inspection, holographic, I:33.16–33.19,
 33.17*f*–33.18*f*, 33.20*f*–33.22*f*
- Installation loss, for fiber optic communication
 links, V:15.6–15.7, 15.8*t*
- Instant pictures, I:25.8
- Instantaneous coefficient of linear thermal
 expansion, IV:4.7
- Institute for Electrical and Electronic
 Engineering (IEEE), II:36.3
- Instrument myopia, III:1.34–1.35
- Instrument signature, of scatterometers, V:1.6,
 1.11*f*, 1.11–1.13, 1.13*t*
- Instrument transfer function effects, in x-ray
 mirror metrology, V:46.9–46.11
- Instrumental line spread function,
 V:38.7
- Instrumental polarization, I:12.15
- Instrumentation, spectroscopic,
 IV:5.58–5.61, 5.59*f*
 detectors, IV:5.61
 and light sources, IV:5.58–5.59
 broadband, IV:5.58–5.59
 laser, IV:5.59
 spectrometers and monochromators,
 IV:5.59–5.61
 dispersion spectrometers, IV:5.59–5.60
 Fourier-transform spectrometers, IV:5.60*f*,
 5.60–5.61, 5.72
- Insulating magnetic brush (IMB),
 I:34.6
- Insulators, lightbulb, II:40.29*f*
- Integral color filter arrays (CFAs),
 II:32.32–32.34, 32.33*f*, 32.34*f*
- Integral density, II:29.14
- Integrated circuits:
 photonic
 of III-V materials, I:21.17–21.20
 in integrated optics, I:21.2
 in WDM systems, I:21.37, 21.38
 polarizers for, I:13.57
- Integrated lasers, with 45° mirror, II:19.39*f*,
 19.39–19.40
- Integrated Mach-Zehnder interferometers,
 V:20.18–20.19, 20.19*f*
- Integrated optic circuits (IOCs),
 I:21.2
- Integrated optics (IO), I:21.1–21.41
 applications of, I:21.31–21.39
 analog transmission, I:21.32–21.34, 21.33*f*,
 21.34*f*
 digital transmission, I:21.31–21.32
 fiber optic gyroscopes, I:21.35–21.37,
 21.36*f*, 21.36*t*
 silicon photonics transmission, I:21.38,
 21.39
 switching, I:21.34*f*, 21.34–21.35
 WDM systems, I:21.37*f*–21.39*f*, 21.37–21.38
 circuit elements of, I:21.21–21.31
 active devices, I:21.25–21.31, 21.26*f*–21.31*f*
 passive devices, I:21.21–21.25, 21.22*f*–21.25*f*
 future trends in, I:21.39–21.41
 advanced integration, I:21.40–21.41
 shift from R&D to manufacturing,
 I:21.39–21.40
 materials and fabrication techniques for,
 I:21.13–21.21
 ion-exchanged glass waveguides, I:21.13
 LiNbO₃ and LiTaO₃, I:21.16–21.17, 21.28,
 21.33
 PICs of III-V materials, I:21.17–21.20
 silicon photonics, I:21.14–21.16, 21.15*f*
 thin film oxides, I:21.13–21.14
 physics of, I:21.3–21.12
 carrier effects, I:21.10–21.12
 index of refraction, I:21.8–21.9
 linear electro-optical effect, I:21.9–21.10
 nonlinear effects, I:21.12
 optical waveguides, I:21.3*f*–21.5*f*,
 21.3–21.8, 21.7*f*
 thermal effects, I:21.12
- Integrated planar lightwave circuits (iPLCs),
 V:21.12, 21.13*f*
- Integrated transmittance, II:35.3
- Integrated-optic modulators, V:7.3, 7.30–7.32,
 7.31*f*–7.33*f*
- Integrating cavities, of nonimaging optics,
 II:39.24–39.27
 efficiency vs. luminance, II:39.26, 39.27*f*
 modifying cavity output distribution,
 II:39.27
 with nonimaging concentrator/collectors,
 II:39.26, 39.27
 nonuniformities with spherical, II:39.24*f*,
 39.24–39.26, 39.25*f*
- Integrating spheres, III:5.10, 5.11

- Integrating spheres (devices), **II**:35.9, 35.11*f*,
35.11–35.13, 37.9–37.10, 37.10*f*
- Integrating x-ray detectors, **V**:61.2
- Integrating-bucket phase shifting, **II**:13.21, 13.21*f*
- Intelligent Physical Protocol Enhanced Physical
Project, **V**:23.4
- “Intelligent” state, **IV**:23.8
- Intense laser pulses:
cluster interactions with, **IV**:21.35–21.36, 21.36*f*
plasma instabilities driven by, **IV**:21.38*f*,
21.38–21.39
- Intensified CCD (ICCD), **I**:26.3
- Intensity, **I**:5.7, 11.2
color vision
defined, **III**:11.2
dependence of hue on, **III**:11.67*f*,
11.67–11.68
of test stimulus, **III**:11.12
defined, **II**:34.9, 40.1
luminous, **II**:39.2*t*
nomenclature for, **II**:36.4
radiant, **II**:39.2*t*
of spectral lines, **V**:56.2–56.4
- Intensity interferometers, **I**:32.19
- Intensity modulation and modulators, **V**:6.34,
13.57–13.58
- Intensity of light, controlling, **III**:5.13–5.15,
5.14*t*
- Intensity reflection coefficients, **I**:12.8–12.12
- Intensity scaling (light), **III**:10.24, 10.25, 10.25*f*
- Intensity transmission coefficients, **I**:12.9–12.10
- Interacting beams, propagation of, **IV**:16.3
- Interaction zone, for Compton scattering,
V:59.1
- Interband absorption, of semiconductors,
IV:5.21–5.33
absorption near fundamental edge, **IV**:5.21*f*,
5.21–5.22
direct transitions, **IV**:5.22*f*–5.23*f*, 5.22–5.23
excitons, **IV**:5.25–5.29, 5.26*t*, 5.27*f*–5.28*f*
high-energy transitions above fundamental
edge, **IV**:5.29–5.33, 5.30*f*–5.34*f*
indirect transitions, **IV**:5.23–5.24, 5.24*f*–5.25*f*
near fundamental edge, **IV**:5.21*f*, 5.21–5.22
polaritons, **IV**:5.29
- Interband magneto-optical effects,
IV:5.42–5.46, 5.43*f*
excitonic, **IV**:5.46, 5.47*f*
Faraday rotation, **IV**:5.44–5.45, 5.45*f*
magnetoreflexion, **IV**:5.43, 5.44, 5.44*f*
- Interband processes, SOAs and, **V**:19.12, 19.13
- Interband transitions, of solids, **IV**:8.27–8.32
direct interband absorption, **IV**:8.27–8.28
excitons, **IV**:8.31–8.32
indirect transitions, **IV**:8.29*f*, 8.29–8.30
joint density of states, **IV**:8.28, 8.29*f*
in metals, **IV**:8.21
multiphoton absorption, **IV**:8.30–8.31
selection rules and forbidden transitions,
IV:8.28, 8.29
- Interconnected switchable networks, **V**:23.4
- Interface coupling, **II**:20.14
- Interference, **I**:2.3–2.42
by amplitude division, **I**:2.19–2.28
extended source, **I**:2.20
Fizeau interferometers, **I**:2.24–2.26, 2.25*f*
fringes of equal inclination, **I**:2.20–2.22,
2.21*f*, 2.22*f*
fringes of equal thickness, **I**:2.22–2.24, 2.23*f*
Michelson interferometer, **I**:2.26*f*–2.27*f*,
2.26–2.28
plane-parallel plate, **I**:2.19, 2.20*f*
thin films, **I**:2.24
applications of, **I**:2.42
and coherence, **I**:2.36*f*, 2.36–2.42,
2.38*f*–2.40*f*, 2.42*f*
constructive and destructive, **I**:2.7
effects of, **I**:2.5–2.14
aberrated wavefronts, **I**:2.12, 2.13
coherence, **I**:2.13
interference fringes, **I**:2.6*f*, 2.6–2.8, 2.7*t*
plane wave and spherical wave, **I**:2.9–2.11,
2.10*f*
temporal beats, **I**:2.13
two plane waves, **I**:2.8–2.9, 2.9*f*
two spherical waves, **I**:2.11–2.12, 2.12*f*, 2.13*f*
- multiple beam, **I**:2.28–2.36
diffraction gratings, **I**:2.28–2.29, 2.29*f*, 2.30*f*
Fabry-Perot interferometer, **I**:2.33–2.36,
2.34*f*, 2.35*f*
plane-parallel plates, **I**:2.30*f*, 2.30–2.33,
2.32*f*, 2.33*f*
- multiple-Bragg-beam, **V**:43.1
in neutron optics, **V**:63.25–63.27, 63.27*f*
Nomarski differential, **V**:46.4
order of, **I**:2.7
Pendellösung, **V**:63.26
in SOAs, **V**:19.23
by wavefront division, **I**:2.14–2.19,
2.15*f*–2.18*f*

- Interference (*Cont.*):
 and wavefronts, **I**:2.4–2.5, 2.5*f*
 and waves, **I**:2.3–2.4, 2.5*f*
 in x-ray optics, **V**:26.8, 26.8*f*, 26.9*f*
- Interference effect, EIT as, **IV**:14.2–14.4
- Interference filters, **II**:34.36
 Fabry-Perot, **IV**:7.78–7.82, 7.79*f*, 7.80*f*,
 7.92–7.94, 7.93*f*, 7.96
 Mach-Zehnder, **I**:21.23, 21.24*f*, 21.25
- Interference fringes, **I**:2.6*f*, 2.6–2.8, 2.7*t*
- Interference gratings, **I**:33.14
- Interference microscopy, **I**:28.33–28.44
 differential-interference contrast microscopes,
I:28.39–28.41, 28.40*f*
 Dyson microscopes, **I**:28.41, 28.42, 28.42*f*
 holographic, **I**:28.42, 28.43, 28.43*f*
 Jamin-Lebedev microscopes, **I**:28.38*f*,
 28.38–28.39
 Linnik microscopes, **I**:28.36–28.38, 28.37*f*
 Mach-Zehnder microscopes, **I**:28.36, 28.37*f*
 Mirau microscopes, **I**:28.41, 28.42, 28.42*f*
 and optical coherence tomography,
I:28.43–28.44
 optical path difference (OPD) in,
I:28.33–28.34, 28.35*f*, 28.36
- Interference polarizers, **I**:13.39–13.41,
 13.40*f*; **IV**:7.69–7.73, 7.70*f*–7.72*f*,
 7.76*f*–7.77*f*
- Interfering transition pathways, **IV**:14.2–14.3
- Interferograms, **II**:13.14–13.18
 from direct interferometry, **II**:13.17–13.18
 fixed, **II**:13.14–13.15
 Fourier analysis of, **II**:13.16–13.17,
 13.17*f*
 interpolation of, **II**:13.15–13.16
- Interferometers and interferometry,
I:32.1–32.21; **II**:13.7–13.12; **III**:5.24
 of angles, **II**:12.14
 atom, **IV**:11.22–11.23, 11.24*f*
 Bragg reflection in, **V**:63.26, 63.27
 Brunning distance-measuring, **II**:12.8*f*,
 12.8–12.9
 common-path, **II**:13.9, 13.11*f*
 computer-generated holograms for, **II**:14.4*f*,
 14.4–14.5, 14.5*f*
 crystal, **V**:63.26–63.27, 63.27*f*
 direct, **II**:13.17–13.18
 distance-measuring, **II**:12.7*f*–12.9*f*,
 12.7–12.10
 double-passed two-beam, **I**:32.8
- Interferometers and interferometry (*Cont.*):
 Fabry-Perot, **I**:32.4–32.7, 32.7*f*, 32.14;
II:16.19*f*; **IV**:7.13, 7.39, 7.39*f*, 7.40, 7.89;
V:3.36, 24.2*f*–24.4*f*, 24.2–24.5
 in dynamic wave meters, **I**:32.17
 in gravitational wave interferometers,
I:32.21, 32.21*f*
 as heterodyne interferometers, **I**:32.10
 and multiple beam interference,
I:2.33–2.36, 2.34*f*, 2.35*f*
 and wire-grid polarizers, **I**:13.31
 fiber, **I**:32.14–32.16, 32.15*f*
 finesse of, **I**:32.6
 Fizeau, **I**:2.24–2.26, 2.25*f*, 32.2, 32.2*f*, 32.17;
II:12.14, 13.8–13.9, 13.9*f*, 13.10*f*, 13.18,
 14.4, 14.5*f*
 free spectral range of, **I**:2.34, 2.35, 2.35*f*, 32.5
 frequency-modulation, **I**:32.9, 32.9*f*, 32.10
 fringe-counting, **I**:32.8, 32.8*f*
 grating, **I**:32.4, 32.5*f*
 gravitational-wave, **I**:32.21, 32.21*f*
 Haidinger, **II**:12.14
 heterodyne, **I**:32.10, 32.10*f*, 32.20, 32.20*f*;
II:13.22
 and holography, **I**:33.4–33.9, 33.8*f*
 infrared, **II**:13.25
 intensity, **I**:32.19
 and interferometric optical switches,
I:32.19
 and interferometric wave meters, **I**:32.16*f*,
 32.16–32.17, 32.17*f*
 laser-Doppler, **I**:32.12–32.13, 32.13*f*
 laser-feedback, **I**:32.13–32.14, 32.14*f*
 lateral-shearing, **II**:12.14, 13.9–13.12, 13.11*f*,
 13.12*f*
 Laue transmission in, **V**:63.26
 Linnik, **V**:46.2
 Mach-Zehnder, **I**:21.10, 21.12, 21.14, 21.16,
 21.32, 21.40, 32.3, 32.3*f*, 33.5; **IV**:23.2*f*,
 23.2–23.4; **V**:21.34
 and Bragg grating sensors, **V**:24.8
 and DPSK, **V**:21.33*f*, 21.33–21.34
 and electro-optic modulators, **V**:7.22, 7.24,
 7.32, 7.38
 and fiber Bragg gratings, **V**:17.8
 integrated, **V**:20.18–20.19, 20.19*f*
 in OTDM networks, **V**:20.22, 20.23*f*
 SOAs in, **V**:19.31*f*, 19.31–19.32, 19.33*f*,
 19.36
 and supercontinuum generation, **V**:11.23

- Interferometers and interferometry (*Cont.*):
of medium distances, **II**:12.6–12.10,
12.7f–12.9f
Michelson, **I**:2.26f–2.27f, 2.26–2.28, 32.2,
32.3f, 32.21, 33.4; **II**:12.5, 12.6, 12.14;
IV:7.42, 7.104f; **V**:17.8f, 17.8–17.9
Michelson stellar, **I**:2.40f, 2.40–2.41, 32.19,
32.19f
micro-interferometers, **II**:10.13, 10.13f
multilayer reflectors for, **IV**:7.39, 7.39f
multiple-pass, **II**:13.13
multiple-reflection, **II**:13.13
Newton, **I**:2.25
Nomarski, **I**:32.4, 32.5f
nonreacting, **II**:12.7
nulling, **I**:32.20–32.21
perfect crystal, **V**:63.26–63.27, 63.27f
phase conjugate, **IV**:12.32, 12.33f, 12.34f
phase-conjugate, **I**:32.17, 32.18, 32.18f
phase-locked, **I**:32.11–32.12, 32.12f
phase-measuring, **V**:46.2
phase-shifting, **I**:32.10–32.11, 32.11f;
II:13.18f–13.20f, 13.18–13.23
heterodyne interferometer, **II**:13.22
integrating bucket method, **II**:13.21,
13.21f
phase errors, **II**:13.22
phase stepping method of, **II**:13.20, 13.20f
phase-lock method, **II**:13.23, 13.23f
simultaneous measurement, **II**:13.22
two steps plus one method, **II**:13.21, 13.22
point diffraction, **II**:13.11f
polarization, **I**:32.4, 32.5f
pulse-train, **II**:20.12, 20.12f
quantum entanglement in, **IV**:23.1–23.15
concepts and equations for, **IV**:23.1–23.4,
23.2f, 23.4f
digital approaches, **IV**:23.7–23.9
Heisenberg limit, **IV**:23.6–23.7
N00N state, **IV**:23.9–23.12, 23.10f, 23.11f
and quantum imaging, **IV**:23.13–23.14
and remote sensing, **IV**:23.14–23.15
shot-noise limit, **IV**:23.4–23.6, 23.5f
radial-shearing, **II**:13.12, 13.13f
reversing-shearing, **II**:13.12, 13.13f
rotational-shearing, **II**:13.12, 13.13f
Sagnac, **I**:21.35, 21.36, 21.36f, 32.3–32.4,
32.4f; **V**:20.22
second-harmonic, **I**:32.17–32.18, 32.18f
sensitivity of, **II**:13.13–13.14, 13.14f
- Interferometers and interferometry (*Cont.*):
shearing, **I**:32.4, 32.6f
single-shot *f*-to-2*f*, **II**:21.6
of small distances, **II**:12.5, 12.6
stellar, **I**:32.19f, 32.19–32.21, 32.20f
sub-Nyquist, **II**:13.27
and third-order optical nonlinearities,
IV:16.28–16.29
three-beam, **I**:32.7f, 32.7–32.8
time-domain atom, **IV**:11.22–11.24, 11.24f
two-wavelength, **II**:13.25, 13.26
and two-wavelength interferometry, **I**:32.9
Twyman-Green, **I**:2.28, 32.2, 32.9, 33.5;
II:13.7f, 13.7–13.8, 13.8f, 13.18
ultrafast nonlinear, **V**:19.32
unbalanced nonlinear, **V**:20.22
Young's two pinhole, **I**:6.3
Zernike phase-contrast method applied to,
II:13.13–13.14, 13.14f
(*see also specific interferometers, e.g.: Lateral-
shearing interferometers*)
- Interferometric arrays, **I**:32.20–32.21
Interferometric ellipsometry, **I**:16.18
Interferometric lithography (IL), **V**:34.4
Interferometric Mach-Zehnder modulators,
V:13.51–13.52, 13.54–13.55
Interferometric method, of FBG fabrication,
V:24.6, 24.6f
Interferometric modulators, **V**:13.51–13.52
Interferometric optical switches, **I**:32.19
Interferometric plots, for orthonormal
aberrations, **II**:11.36–11.37, 11.37f, 11.38f
Interferometric wave meters, **I**:32.16f,
32.16–32.17, 32.17f
Intergovernmental Panel on Climate Change,
V:3.44
Interior lighting, **II**:40.55–40.61
for health-care facilities, **II**:40.58–40.60,
40.60t
for industry, **II**:40.60–40.61, 40.61f
for offices, **II**:40.55, 40.56t
for residences, **II**:40.57, 40.58, 40.59t
for retail, **II**:40.55–40.57, 40.56t–40.58t
Interlaced monitors, **III**:23.2
Interlayer interimage effects (IIEs), **II**:30.19
Interleaving, in OTDM networks, **V**:20.6–20.7,
20.7f
Interline transfer, of CCDs, **I**:26.4–26.5, 26.5f,
26.6f

- Interline transfer (IT) CCD image sensors, **II**:32.28–32.32, 32.29*f*–32.31*f*
- Interline-transfer (IT) CCD FPAs, **II**:33.11–33.13, 33.12*f*
- Intermodal dispersion, in optical fibers, **V**:9.5–9.6
- Intermodulation distortions (IMDs), **V**:15.5–15.6
- Intermodulation (IM) products, of acousto-optic devices, **V**:6.30
- Internal quantum efficiency, **V**:13.9, 13.10
- Internal self-action, **IV**:13.7, 16.25
- Internal writing technique, for FBGs, **V**:17.4–17.5
- Internally processed (IP) photocathodes, **II**:31.24
- International Association for Physical Sciences of the Ocean (IAPSO), **IV**:1.4, 1.5*t*–1.6*t*
- International Astronomical Union (IAU), **II**:36.3
- International Bureau of Weights and Measures (BIPM), **II**:36.2
- International candle (unit), **II**:37.3
- International Commission on Illumination (CIE), **II**:36.2 (*see also* Commission Internationale de l'Éclairage)
- International Commission on Non-Ionizing Radiation Protection (ICNIRP), **III**:7.9, 7.12
- International Committee for Weights and Measures (CIPM), **II**:36.2, 38.6
- International Electrotechnical Commission (IEC), **III**:7.9, 7.12
- International Graphics Exchange Specification (IGES), **II**:40.19
- International Space Station, **V**:49.4
- International Standards Organization (ISO), **II**:4.10, 4.11, 36.2, 40.19; **III**:7.9
- International System of Units (*see* SI units)
- International Telecommunications Union standards, **V**:15.15
- International Union of Pure and Applied Physics (IUPAP), **II**:36.2
- Interocular aniso-magnification, distortion from, **III**:13.16–13.19
- aniseikonia, **III**:13.17–13.18
- interocular blur suppression with anisometropia, **III**:13.18–13.19
- lenses and prisms, **III**:13.16–13.17, 13.17*f*
- Interphotoreceptor matrix (IPM), **III**:8.8, 8.9
- Interpupillary distance (IPD), **III**:1.39–1.41, 1.40*f*
- Intersection points, ray tracing for, **I**:1.36, 1.36*f*
- Interstitial matrix, **III**:8.2
- Inter-switch links (ISLs), **V**:23.4
- Intersymbol interference (ISI), **V**:19.23
- Intervallence band absorption (IVBA), **II**:19.17
- Intervalley scattering, **IV**:18.20
- Interwoven pulse trains, **II**:20.13
- Intraband magneto-optical effects, **IV**:5.47–5.52
- cyclotron resonance, **IV**:5.47–5.50, 5.48*f*–5.50*f*
- free-carrier Faraday rotation, **IV**:5.50–5.51
- impurity magnetoabsorption, **IV**:5.51*f*, 5.51–5.52
- Intraband processes, of SOAs, **V**:19.12
- Intracapsular cataract extraction (ICCE), **III**:12.14
- Intracavity singly-resonant oscillators (IC-SROs), **IV**:17.3*f*
- Intracorneal ring segments, **III**:16.10*f*, 16.10–16.11
- Intramodal dispersion, **V**:9.5–9.6, 9.6*f*
- Intraocular lenses (IOLs), **III**:12.15, 16.9, 21.1–21.22
- accommodating, **III**:14.29–14.30, 21.18–21.19
- and aging of the eye, **III**:21.3–21.4
- aspheric, **III**:21.10–21.12, 21.11*f*, 21.12*t*
- for cataract correction, **III**:14.8
- and cataract surgery, **III**:21.4
- defined, **III**:14.1–14.2, 21.1
- design of, **III**:21.5–21.20
- accommodating lenses, **III**:21.18–21.19
- aspheric lenses, **III**:21.10–21.12, 21.11*f*, 21.12*t*
- chromophores, **III**:21.20, 21.20*f*
- IOL aberrations, **III**:21.9–21.10, 21.10*f*
- IOL power, **III**:21.5–21.9, 21.7*t*, 21.8*f*
- multifocal lenses, **III**:21.14–21.18, 21.15*f*, 21.16*f*, 21.16*t*, 21.18*f*
- phakic lenses, **III**:21.19
- testing IOLs, **III**:21.13–21.14
- toric lenses, **III**:21.13, 21.12*f*
- effective lens position (ELP), **III**:21.7–21.9
- multifocal, **III**:21.14–21.18, 21.15*f*, 21.16*f*, 21.16*t*, 21.18*f*
- phakic, **III**:21.19
- for presbyopic correction, **III**:14.28–14.29

- Intraocular lenses (IOLs) (*Cont.*):
 side effects of, **III**:21.20–21.22
 dysphotopsia, **III**:21.21
 posterior capsule opacification,
III:21.21–21.22
 and structure of the eye, **III**:21.2–21.3
 toric, **III**:21.13, 21.12*f*
- Intraocular pressure (IOP), **III**:14.2, 14.26,
 14.27, 16.3
- Intraocular scatter, age-related, **III**:14.12
- Intrinsic athermalization, **II**:8.7*f*, 8.7–8.8
- Intrinsic Fabry-Perot interferometric (IFPI)
 sensors, **V**:24.4*f*, 24.4–24.5
- Intrinsic infrared detectors, **II**:33.7, 33.8*f*
- Intrinsic optical properties:
 of semiconductors, **IV**:5.11
 of solids, **IV**:8.3
- Intrinsic photoconductors, **II**:25.5, 25.5*f*
- Intrinsic photodetectors, **II**:24.7, 24.7*f*
 germanium, **II**:24.70*f*–24.73*f*, 24.70–24.73
 indium antimonide photovoltaic,
II:24.80–24.83, 24.82*f*, 24.83*f*
- Intrinsic semiconductor transition, **II**:24.11
- Invar 36, **IV**:4.10*t*, 4.52*t*, 4.55*t*, 4.69*t*, 4.70*t*
- Invariance properties, of rays, **I**:1.10
- Invariant hues:
 and Bezold-Brücke effect, **III**:11.67*f*,
 11.67–11.68
 defined, **III**:11.2
- Invariants, **I**:18.7; **II**:1.11
- Inverse Bremsstrahlung heating, **IV**:21.37,
 21.37*f*
- Inverse Compton scattering, **V**:59.1
- Inverse Compton x-ray sources, **V**:59.1–59.4,
 59.3*t*
- Inverse dielectric tensor, **IV**:2.6*t*, 2.19
- Inverse filters, for coherent optical image
 enhancement, **I**:11.14–11.17
- Inverse Galilean telescopes, **I**:18.15, 18.16,
 18.16*f*
- Inverse Raman effect, **V**:10.5, 10.6, 14.9
- Inverse square law, **II**:34.14, 37.8
- Inverse systems, conjugate matrices for, **I**:1.71
- Inversion layer, in MOS transistors, **II**:25.11,
 25.11*f*
- Invertebrate aquatic worm, eye structures in,
III:9.14
- Inverted channel substrate planar (ICSP) lasers,
II:19.20*t*, 19.23
- Inverted telephoto lenses, **I**:27.2, 27.6,
 27.7*f*–27.14*f*
 highly complex extreme speed, **I**:27.6,
 27.9*f*–27.12*f*
 with nonrectilinear distortion, **I**:27.6
 with rectilinear distortion correction, **I**:27.6,
 27.13*f*, 27.14*f*
 very compact moderate speed, **I**:27.6,
 27.7*f*–27.8*f*
- Involute reflectors, **II**:39.11–39.12, 39.12*f*
- Ion beam-sputtered surfaces, **IV**:6.53
- Ion bombardment strip lasers, **II**:19.8, 19.9*f*
- Ion current measurement devices, **II**:34.29
- Ion exchange process, **I**:24.8
- Ion-assisted deposition, **IV**:7.11
- Ion-beam sputtering, **IV**:7.11, 7.14
- Ion-exchanged glass waveguides, **I**:21.13
- Ion-implanted semiconductors, **IV**:18.21
- Ionization:
 above-threshold (*see* Above threshold
 ionization)
 atomic, **IV**:21.3
 barrier suppression, **IV**:21.14, 21.14*f*
 in clusters, **IV**:21.31–21.33, 21.32*f*
 collisional, **IV**:21.31, 21.32
 double, **IV**:21.18–21.19, 21.19*f*
 of electrons, **V**:56.2
 inner, **IV**:21.31–21.32, 21.32*f*
 molecular tunnel (*see* Molecular tunnel
 ionization)
 multiphoton, **IV**:21.10–21.12, 21.11*f*
 outer, **IV**:21.32–21.33
 stabilization of, **IV**:21.20–21.21, 21.22*f*
 threshold (*see* Threshold ionization)
 tunnel (*see* Tunnel ionization)
- Ionization chambers, **V**:60.3–60.4, 60.9*t*
- Ionization distance, **IV**:21.25*f*, 21.25–21.26,
 21.27*f*
- Ionization ignition, **IV**:21.31
- Ionization rate, ADK, **IV**:21.13
- Ionization stabilization, **IV**:21.20–21.21, 21.22*f*
- Ionization x-ray detectors, **V**:60.3–60.7, 60.9*t*,
 60.10*t*
- Ionization-induced defocusing, **IV**:21.43*f*,
 21.43–21.44
- Ionized arsenic antisite (As_{Ga}^+), **IV**:18.3
- Ionizing radiation, fiber optic communication
 links and, **V**:15.17
- Ion-plating process, **IV**:7.11

- Ions, tri-positive rare earth, **I**:10.16–10.18, 10.16*t*, 10.17*f*
- IRG 2 glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- IRG 9 glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- IRG 11 glass, **IV**:2.43*t*, 2.49*t*, 2.59*t*, 2.68*t*
- IRG 100 glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.68*t*
- Iris (eye), **III**:1.3*f*, 16.3, 21.2; **IV**:13.1
- Iris (lens), **I**:17.8
- Iron:
- absorptance of, **IV**:4.40*f*, 4.48*t*, 4.50*t*
 - optical properties of, **IV**:4.15*t*, 4.23*f*
 - physical properties of, **IV**:4.54*t*
 - reflectance of, **IV**:4.32*t*–4.33*t*, 4.40*f*
 - thermal properties of
 - coefficient of linear thermal expansion, **IV**:4.56*t*, 4.57*f*
 - elastic stiffness, **IV**:4.69*t*
 - moduli and Poisson's ratio, **IV**:4.69*t*
 - at room temperature, **IV**:4.55*t*
 - specific heat, **IV**:4.65*t*, 4.67*f*
 - thermal conductivity, **IV**:4.58*t*, 4.62*f*–4.63*f*
- Irradiance, **I**:3.3
- of circular patterns and disks, **I**:3.7–3.8
 - of complementary aperture screen patterns, **I**:3.10–3.11
 - defined, **II**:34.8, 36.5, 37.4*t*, 37.5, 39.2*t*
 - and diffraction gratings, **I**:3.28, 3.29
 - of double-slit patterns, **I**:3.27, 3.27*f*
 - excitance and emittance vs., **II**:39.2
 - image specifications for, **II**:4.7
 - of lambertian objects, **I**:1.80
 - spectral, **II**:36.14, 38.1–38.2, 38.11*t*, 38.13*f*–38.16*f*, 38.13–38.16
 - of straight-edge patterns, **I**:3.15–3.17, 3.18*f*
 - as vector, **I**:3.33, 3.37*f*
 - of water, **IV**:1.5*t*, 1.8–1.9
 - and zone plates, **I**:3.11, 3.12*f*
- Irradiance reflectance (irradiance ratio), of water, **IV**:1.6*t*, 1.7*f*, 1.12, 1.46–1.47, 1.47*f*
- Irradiance response units, **II**:34.31
- Irradiated plasma, **IV**:21.46*f*, 21.46–21.47
- Irradiation, microbeam, **I**:28.54
- Irregular astigmatism, **III**:16.6
- ISO standards, **V**:2.2
- Isoaccommodation circle, **III**:13.23, 13.23*f*
- Isochromatic CFSs, **III**:2.30, 2.30*f*
- Isodiscrimination contours, **III**:10.40
- Isoelectronic dopants, **II**:17.16
- Isoelectronic trap, **II**:17.6, 17.6*f*
- Isosindical contours, **III**:19.1
- Isolated pulses, in attosecond optics, **II**:21.4
- Isolating direction (color vision), **III**:11.2, 11.33
- Isolation (directionality), of fiber-optic components, **V**:18.1
- Isolators, for networking, **V**:18.3, 18.10, 18.10*f*
- Isoluminant CFSs, **III**:2.30, 2.30*f*
- Isoluminant patterns, **III**:2.30
- Isometric plots, for orthonormal aberrations, **II**:11.36–11.37, 11.37*f*, 11.38*f*
- Isoplanatic patch, **III**:15.5–15.7, 15.6*f*
- Isoplanatic angle, **V**:5.19
- Isotope broadening, **II**:16.6
- Isotropic (term), **II**:36.4, 39.3
- Isotropic crystals:
- dielectric constants of, **IV**:2.18
 - symmetries of, **IV**:2.7*t*
- Isotropic diffraction:
- and acousto-optic interactions, **V**:6.9–6.10, 6.10*f*, 6.13, 6.13*f*
 - by deflectors, **V**:6.23–6.24, 6.24*f*
- Isotropic homogenous spheres, scattering by, **I**:7.11–7.14
- Isotropic point source, **II**:36.4
- Isotropic solids, **IV**:8.8, 8.9*t*
- Isovergence circle, **III**:13.23, 13.23*f*
- Iterated rays, **II**:3.12
- Iterative phasing technique, for coherent diffraction microscopy, **V**:27.4
- Iturriaga R., **IV**:1.25
- ivity (suffix), **II**:35.3; **IV**:4.5
- J** × **B** heating, **IV**:21.49
- Jackson cross-cylinder check test, **III**:12.7
- Jacobi method, of electro-optic effect, **V**:7.11–7.13
- J-aggregates, **II**:30.13, 30.14
- Jahn-Teller effect, **V**:2.9
- Jamin-Lebedev microscopes, **I**:28.38*f*, 28.38–28.39
- Javal's rule, **III**:12.5
- Jaynes-Cummings model, **II**:23.15*f*, 23.15–23.17
- Jellett-Cornu prisms, **I**:13.56
- Jerlov water types, **IV**:1.42–1.46
- Jet polishing, **II**:9.6
- Jitter:
- in fiber optic communication links, **V**:15.15–15.16
 - and solitons, **V**:22.6–22.8, 22.11, 22.16
- Jitter transfer function (JTF), **V**:15.16

- J-K model, of surface finish, **I**:8.13, 8.15–8.17
- Johansson bent/ground focusing monochromator, **V**:39.5
- Johnson noise, of pin diodes, **V**:13.70 (*see also* Thermal noise)
- Johnson noise power density, **II**:28.3
- Joint density of states, **IV**:8.28, 8.29*f*
- Jones (unit), **II**:24.11, 24.13
- Jones calculus, **I**:12.29–12.30
- Jones matrix, **I**:7.10
 - in ellipsometry, **I**:15.30, 16.19
 - and Mueller matrices, **I**:14.3, 14.22–14.24, 14.27, 14.33
 - tensor product of, **I**:14.23
- Jones matrix formalism, **III**:18.1, 18.20, 18.22–18.24
- Joule (unit), **II**:34.5–34.6, 37.6
- Judd's three-stage Müller zone theory, **III**:11.6, 11.7*f*, 11.8
- Judd-Vos modified 2° color-matching functions, **III**:10.12, 10.45
- Judd-Vos modified function, **II**:36.10
- Judgment tasks:
 - adjustment tasks vs., **III**:3.2
 - psychophysical measurement of, **III**:3.6–3.8
 - ideal observer, **III**:3.6
 - rating scale, **III**:3.6–3.8, 3.7*f*
 - response time, **III**:3.8
 - two-alternative forced choice (2afc), **III**:3.8
 - yes-no, **III**:3.6
- Jülich SANS instrument, **V**:64.4
- Jumps, electric-field, **IV**:9.6
- Junction photodiodes, **II**:32.3–32.6, 32.4*f*, 32.6*f*
- Junge (hyperbolic) cumulative size distribution, **IV**:1.15
- Just noticeable distortion level (JND), **III**:24.4
- JV model (of trap-loss collisions), **IV**:20.30
- “K,” **III**:20.2
- K (optical constant of water), **IV**:1.17, 1.17*f*
- Kagomé lattice, **V**:11.5*f*, 11.11, 11.16
- Kaleidoscope effect, **II**:39.28
- Kane momentum matrix elements, **IV**:8.26
- Karhunen-Loève functions, **V**:4.36
- Keck Observatory, **V**:5.27
- Keck telescope, **II**:11.4; **V**:4.36, 5.2, 5.27
- Keldysh parameter, **IV**:21.10
- Keldysh-Faisal-Reiss (KFR) theories, **IV**:21.12
- Keplerian afocal lenses, **I**:18.7–18.14
 - in binoculars, **I**:18.13–18.14, 18.14*f*
 - and eye relief manipulation, **I**:18.8–18.10, 18.9*f*, 18.10*f*
 - field-of-view limitations in, **I**:18.11
 - finite conjugate afocal relays, **I**:18.11–18.12, 18.12*f*
 - for scanners, **I**:18.13, 18.13*f*
 - in terrestrial telescopes, **I**:18.10–18.11, 18.11*f*
 - thin-lens model, **I**:18.7–18.8, 18.8*f*
- Keplerian telescopes, **I**:18.10
- Keratoconus, **III**:9.2, 20.24
- Keratocytes, **III**:16.1, 16.4
- Keratometer mire, **III**:12.5
- Keratometry, **III**:12.5, 21.6
- Kerr cell shutters, **I**:25.21
- Kerr cells, **I**:31.9; **V**:7.34
- Kerr effect, **I**:15.23, 35.21, 35.23, 35.25; **IV**:18.11–18.15; **V**:7.11, 14.11, 19.34, 20.1, 20.22, 22.3
 - longitudinal, **IV**:18.11–18.15, 18.12*f*
 - optical, **IV**:16.11–16.14
 - Raman-induced, **IV**:16.3*t*, 16.12, 16.17
 - transverse, **IV**:18.11, 18.14–18.15
- Kerr electro-optic effect, **V**:7.9*t*–7.10*t*, 7.11
- Kerr focusing, **V**:7.39
- Kerr interactions, **V**:10.2
- Kerr lens modelocking (KLM), **IV**:16.25, 18.3, 18.14–18.15; **V**:7.11, 7.39
- Kerr nonlinearities, **V**:7.38–7.39, 11.22–11.24, 11.24*f*
- Kerr-lensing (*see* Self-focusing)
- Kerr-type nonlinearity, **IV**:14.33
- Kick operator, **II**:23.17
- Kinematic lens mounts, **I**:22.7, 22.8*f*
- Kinematic theory of diffraction, in neutron optics, **V**:63.5
- Kinematical theory of diffraction, for crystals, **V**:39.2, 39.3
- Kinetic cues, for perceived space, **III**:13.4–13.7, 13.5*f*, 13.6*f*
- Kirchhoff's law, **II**:34.25, 35.7, 35.8*t*; **IV**:4.6
- Kirkpatrick-Baez (KB) mirrors, **V**:44.4, 44.4*f*, 63.21, 64.5*f*, 64.5–64.6
- Kirkpatrick-Baez (KB) optics, **V**:47.7–47.8, 47.8*f*, 47.9*f*, 50.5
- Kirkpatrick-Baez (KB) systems, **V**:48.3*f*, 48.3–48.4
- Kitaev phase estimation algorithm, **IV**:23.12
- Klein bottle, **I**:14.15

- Kleinman **d**-tensor, **IV**:10.11
 Knapp's law, **III**:13.17, 13.26
 Knife-edge test (*see* Foucault test)
 Knill-LaFlamme-Milburn (KLM) scheme, **IV**:23.12
 Knoop test, **IV**:2.31, 2.32*f*
 Kodachrome film, **II**:29.14, 30.23
 Kodacolor, **II**:29.14
 Kodak Cameo Motordrive camera, **I**:25.16, 25.16*f*
 Kodak Cobra flash, **I**:25.16, 25.16*f*
 Kodak DCS 200 camera, **I**:25.7
 Kodak Fun Saver Portait 35, **I**:25.16, 25.16*f*
 Kodak glass molding process, **I**:22.9, 22.10
 Kodak Gold film, **II**:30.25
 Kodak Royal Gold film, **II**:30.25
 Kodak Technical Pan Film, **II**:29.19*t*
 Koehler (Köhler) illumination, **II**:1.11, 1.11*f*,
 39.23*f*, 39.23–39.24, 39.34, 39.35*f*; **III**:5.6
 and coherence theory, **I**:6.12*f*, 6.12–6.13
 in microscopes, **I**:28.4–28.5, 28.5*f*, 28.7, 28.9
 Kolmogorov model of turbulence, **V**:5.5–5.6, 5.11
 Kolmogorov spatial power spectral density, **V**:5.6
 Kolmogorov spectrum, **V**:3.28, 3.31
 Kolmogorov turbulence, **V**:4.3, 4.7–4.10, 4.8*f*,
 4.9*f*, 4.27, 4.30, 4.36
 Kopelevich model of absorption, **IV**:1.28
 Korsch objectives, **I**:29.30–29.32, 29.34
 Kossel patterns, **V**:53.11
 Kramer approximation, for Bremsstrahlung
 radiation, **V**:54.5
 Kramers states, **I**:31.19, 31.19*f*, 31.20
 Kramers-Kronig (K-K) relations, **I**:21.11;
V:13.61, 17.3, 20.22
 dielectric-constant, **IV**:2.9, 2.12, 2.22
 dispersion, **IV**:5.10–5.11, 16.1, 16.9–16.11
 for solids, **IV**:8.15
 Kronecker delta, **V**:4.19
 Kronecker delta function, **I**:5.20
 Kronig-Penney model, **V**:2.12
 KRS-5 crystals, **IV**:2.40*t*, 2.44*t*, 2.48*t*, 2.53*t*,
 2.58*t*, 2.65*t*
 KRS-5 fiber, **V**:12.7–12.8
 KRS-6 crystals, **IV**:2.40*t*, 2.44*t*, 2.48*t*, 2.53*t*,
 2.58*t*, 2.65*t*
 K-shell emission lines, of elements, **V**:36.3*t*–36.8*t*
 Kubelka-Munk theory, **I**:9.13
 Kumahov capillary lenses, **V**:53.9

 La Jolla colorimeters, **III**:5.10
 Lab-based radiation sources, **V**:50.2–50.7, 50.4*f*,
 50.6*f*, 50.8*f*
 Laboratory sources (of radiation), **II**:15.7–15.13
 baseline standard for, **II**:15.9, 15.9*f*, 15.10*f*,
 15.12*f*
 blackbody cavity theory, **II**:15.7–15.9, 15.8*f*
 working standards for, **II**:15.9–15.13, 15.10*f*,
 15.12*f*, 15.13*f*
 Labsphere, **II**:38.12
 Lacrimal lens:
 and contact lens power, **III**:20.12*f*–20.14*f*,
 20.12–20.15
 defined, **III**:20.2
 Ladder coupling, **IV**:14.1, 14.6*f*, 14.24*f*
 Lag, **III**:23.2
 Lagrange invariant, **I**:1.22, 1.41, 1.77, 1.81, 13.7,
 17.5, 30.11
 Lagrangian rays, **II**:3.12
 Lamb dip spectroscopy, **V**:2.5, 2.7*f*
 Lamb shift(s), **I**:10.4; **II**:23.13; **V**:2.2, 2.3
 Lambda coupling, **IV**:14.1, 14.6*f*, 14.8*f*, 14.9*f*,
 14.26, 14.27
 Lambda Research Corporation, **II**:7.27
 Lambert (unit), **II**:34.43, 36.7, 36.8*t*
 Lambert-Beer law of attenuation, **V**:63.11
 Lambertian approximation (of radiant flux
 transfer), **II**:34.14–34.18
 and lambertian sources, **II**:34.14–34.17,
 34.15*f*, 34.16*f*
 radiant flux transfer through lambertian
 reflecting sphere, **II**:34.17–34.18
 Lambertian black surface, **IV**:6.14
 Lambertian objects, irradiance for, **I**:1.80
 Lambertian surface, **II**:37.8
 Lambert's cosine law, **II**:37.8, 37.8*f*
 Lambert's law, **I**:5.12
 Lamellas, **III**:16.1, 16.4
 Lamipol structures, **I**:13.57
 Lamor precession, **V**:63.29–63.30
 Lamp safety standards, **III**:7.14–7.15
 Lamps:
 cold cathode fluorescent, **V**:8.30, 8.30*f*
 configurations of, **II**:15.20*f*
 modeling of, **II**:40.17
 standards for, **II**:15.11
 (*see also specific types of lamps, e.g.*: Airway
 beacon lamps)
Lamps for Scientific Purposes (G. M. B. H.
 Osram), **II**:15.20
 Land (term), **II**:34.35, 34.35*f*
 Landau levels (of energy), **IV**:5.40, 5.42*f*
 Landé interval formula, **I**:10.17

- Landolt C test, **III**:4.8
- Landolt fringe, **I**:13.14, 13.17, 13.18
- Landolt ring target, **III**:2.34
- Lands, of optical disks, **I**:35.3, 35.5*f*
- Landscape lenses, **I**:17.17*f*–17.20*f*, 17.17–17.20, 17.19*t*
- Langmuir waves, **IV**:21.38
- Langmuir-Blodgett techniques, **V**:7.34
- Lanthanum glass, **IV**:2.42*t*, 2.43*t*
- Laporte selection rule, **I**:10.10
- Lapping step (of optics fabrication), **II**:9.5
- Large Binocular Telescope (LBT), **V**:5.5
- Large field (color matching):
 defined, **III**:10.9
 standards for, **III**:10.11, 10.12–10.13
- Large flat-mode fibers, in fiber lasers, **V**:25.19*f*, 25.20
- Large-area detectors, **II**:25.12
- Large-format cameras, **I**:25.18–25.20
- Large-format film, **I**:25.6
- Large-mode-area (LMA) fibers, **V**:25.2, 25.4, 25.5
 in all-fiber monolithic systems, **V**:25.16
 chirally coupled core fibers, **V**:25.21–25.22
 designs, **V**:25.19*f*
 equations for, **V**:25.8–25.9
 and photonic crystal fibers, **V**:25.20–25.21
 techniques using, **V**:25.18–25.20
- LASEK (laser subepithelial keratomileusis), **III**:16.12, 16.13
- Laser(s), **II**:12.7, 16.1–16.37
 about, **II**:16.2–16.3
 Bragg reflector, **V**:13.7
 bulk, **V**:25.5–25.6
 buried heterostructure, **V**:13.5
 carbon dioxide, **V**:12.3*t*, 12.9, 12.13
 chirp of, **V**:9.8
 continuous wave (CW) dye, **I**:10.8
 continuous-wave, **IV**:7.14, 14.16–14.18, 14.17*f*; **V**:25.4, 25.5, 25.5*f*, 25.7
 diode lasers, **V**:13.49
 dye lasers, **V**:5.32, 5.33*f*, 5.34
 diagram of, **II**:16.2*f*
 distributed Bragg reflector, **I**:21.25, 21.30, 21.32, 21.37, 21.37*f*; **V**:9.8, 13.7, 13.28–13.29, 13.29*f*, 20.14
 distributed feedback, **I**:21.25, 21.29, 21.30, 21.32, 21.38, 21.40
 in fiber optic systems, **V**:13.30*f*, 13.30–13.32, 13.31*f*
 and optical fibers, **V**:9.8
- Laser(s), distributed feedback (*Cont.*):
 in OTDM networks, **V**:20.14–20.15, 20.15*f*
 quarter-wavelength-shifted grating, **V**:13.31*f*, 13.31–13.32
- edge-emitting, **V**:13.3
 and electroabsorption modulators, **V**:13.60
 electromagnetic spectrum involving, **II**:16.2, 16.3, 16.3*f*
- erbium-doped yttrium aluminum garnet (Er:YAG), **V**:12.3*t*, 12.6, 12.13
- European X-ray Free-Electron Laser, **V**:58.1
- external cavity, **V**:13.22, 13.33
- extreme ultraviolet, **V**:58.2–58.4, 58.3*f*
- eye injury from, **III**:7.3, 7.4
- Fabry-Perot, **V**:13.12–13.13, 13.29
 multilongitudinal mode, **V**:9.7–9.8
 semiconductor, **V**:20.13–20.14, 20.14*f*
- fiber, **V**:25.1–25.33
 all-fiber monolithic systems, **V**:25.16*f*, 25.16–25.18
 architectures, **V**:25.9–25.18, 25.12*f*, 25.19*f*
 bulk lasers vs., **V**:25.5–25.6
 dopants, **V**:25.22–25.26, 25.23*t*, 25.23–25.26
 equations for, **V**:25.8–25.9
 fabrication of, **V**:25.26–25.29, 25.27*t*
 free space, **V**:25.13–25.15, 25.14*f*, 25.15*f*
 growth of, **V**:25.5, 25.5*f*
 history of, **V**:25.3–25.4, 25.4*f*
 limitations, **V**:25.6–25.7
 LMA fiber designs, **V**:25.18–25.22, 25.19*f*
 operation of, **V**:25.7–25.8
 pumping techniques, **V**:25.9–25.13, 25.11*f*, 25.12*f*
 spectral and temporal modalities, **V**:25.29–25.33
- fiber Raman, **V**:10.7
- free electron, **V**:41.9–41.10, 41.10*f*, 48.1, 58.1–58.2
- free space fiber, **V**:25.9, 25.10, 25.13–25.15, 25.14*f*, 25.15*f*
- hazards related to, **III**:7.11–7.14, 7.12*f*, 7.13*f*
 accidents, **III**:7.12–7.14
 and eye protectors, **III**:7.14
 safety standards, **III**:7.12
- for interferometry, **I**:2.41–2.42
- and laser gain medium, **II**:16.4–16.19
 emission linewidth and line broadening of radiating species, **II**:16.4–16.7, 16.6*f*, 16.7*f*
 energy levels and radiation, **II**:16.4

- Laser(s), and laser gain medium (*Cont.*):
 gain saturation, **II**:16.10
 optimization of output coupling from
 laser cavity, **II**:16.13, 16.14
 population inversions, **II**:16.8–16.10,
 16.12–16.13, 16.13*f*, 16.14*f*
 pumping techniques to produce inversions,
II:16.14–16.19
 stimulated absorption and emission,
II:16.7–16.8, 16.8*f*
 threshold conditions, **II**:16.10–16.12, 16.11*f*
 as light sources, **II**:40.39
 in measurement, **II**:12.2, 12.6
 mesospheric sodium, **V**:5.32–5.34, 5.33*f*
 mode-locked, **V**:20.15–20.17, 20.16*f*
 multilayer reflectors for, **IV**:7.39*f*, 7.39–7.40
 multiple quantum well, **V**:13.24–13.25
 and optical cavities or resonators, **II**:16.19*f*,
 16.19–16.25
 configurations and cavity stability,
II:16.23–16.25, 16.24*f*, 16.25*f*
 longitudinal laser modes, **II**:16.20, 16.20*f*
 transverse laser modes, **II**:16.21*f*–16.23*f*,
 16.21–16.23
 and optical fibers, **V**:9.7–9.8
 in OTDM communication networks,
V:20.13–20.17, 20.14*f*, 20.16*f*
 phase noise (finite line width) of, **V**:9.8
 planar buried heterostructure, **V**:13.6
 plasma-based EUV, **V**:58.2–58.4, 58.3*f*
 probability flow diagram for, **II**:23.23*f*
 pulsed-dye, **V**:5.32
 quantum theory of (*see* Quantum theory of
 lasers)
 quantum well, **V**:13.24–13.28, 13.25*f*–13.27*f*
 as radiometric characterization tool, **II**:34.32
 relative intensity noise of, **V**:9.11
 ridge waveguide, **V**:13.6
 for scatterometers, **V**:1.8
 semiconductor, **V**:13.1, 20.13–20.14, 20.14*f*
 semiconductor arrays of, **II**:19.26–19.29,
 19.28*f*, 19.28*t*
 single-longitudinal-mode, **V**:9.8
 special laser cavities, **II**:16.25–16.29
 distributed feedback lasers, **II**:16.29
 mode-locking, **II**:16.27–16.29, 16.28*f*
 Q-switching, **II**:16.26–16.27, 16.27*f*
 ring lasers, **II**:16.29
 unstable resonators, **II**:16.25–16.26, 16.26*f*
 and speckle fields, **III**:5.21
- Laser(s) (*Cont.*):
 sum-frequency, **V**:5.32
 theoretical description of, **IV**:20.6–20.11
 atoms in motion, **IV**:20.8–20.10, 20.9*f*
 Fokker-Planck equation, **IV**:20.10–20.11
 force on two-level atom, **IV**:20.6–20.7
 tunable, **V**:13.32–13.36, 13.33*f*–13.36*f*
 two-dimensional high-power arrays of,
II:19.29–19.30, 19.29*t*, 19.30*f*
 as two-level system, **II**:20.23–20.24, 20.25*t*
 types of, **II**:16.29–16.37, 16.31*f*, 16.32*f*,
 16.34*f*, 16.35*f*, 16.37*f*; **III**:5.20*t*
 vertical cavity surface-emitting, **V**:9.7*n*,
 13.42–13.48, 13.43*f*, 13.45*f*, 19.14
 and wavelength-division multiplexing, **V**:9.8
 x-ray, **V**:58.1–58.4, 58.3*f*
 (*see also related topics, e.g.:* Laser
 stabilization)
- Laser ablation, **III**:16.11–16.19
 ablation profiles, **III**:16.14–16.15
 ablation rate, **III**:16.16–16.18, 16.18*f*
 corneal photoablation, **III**:16.16, 16.17*f*
 Epi-LASIK, **III**:16.12, 16.13, 16.13*f*
 LASEK, **III**:16.12, 16.13
 LASIK, **III**:16.13–16.14, 16.14*f*
 photorefractive keratectomy (PRK),
III:16.11*f*, 16.11–16.12, 16.12*f*
 thermal, photochemical, and photoacoustic
 effects, **III**:16.18–16.19
- Laser beacons (laser guide star sensing), **V**:5.21,
 5.23, 5.27–5.34
 focus anisoplanatism of, **V**:5.27–5.29, 5.28*f*–5.30*f*
 and mesospheric sodium laser beams,
V:5.32–5.34, 5.33*f*
 Rayleigh, **V**:5.30–5.32, 5.31*f*
- Laser beam expanders, athermal, **II**:8.13–8.14
- Laser beam scanning, by high-resolution
 deflectors, **V**:6.29
- Laser Black, **IV**:6.56
- Laser cavities, **II**:16.25–16.29
 distributed feedback lasers, **II**:16.29
 mode-locking, **II**:16.27–16.29, 16.28*f*
 Q-switching, **II**:16.26–16.27, 16.27*f*
 ring lasers, **II**:16.29
 unstable resonators, **II**:16.25–16.26, 16.26*f*
- Laser conditioning, **IV**:19.4
- Laser cooling, **IV**:20.3–20.21, 20.26–20.39
 about, **IV**:20.3–20.4
 in atomic beam brightening, **IV**:20.27*f*,
 20.27–20.28

- Laser cooling (*Cont.*):
 in atomic clocks, **IV**:20.28
 below Doppler limit, **IV**:20.17–20.21,
 20.18*f*–20.20*f*
 in Bose-Einstein condensation,
IV:20.35–20.37, 20.36*f*
 in dark states, **IV**:20.37–20.39, 20.38*f*
 defined, **IV**:20.3
 history of, **IV**:20.3–20.4
 in optical lattices, **IV**:20.31–20.34,
 20.32*f*–20.34*f*
 by optical molasses, **IV**:20.13–20.17,
 20.14*f*–20.16*f*
 properties of, **IV**:20.4–20.6
 Sisyphus, **IV**:20.19
 by slowing of atomic beams, **IV**:20.11–20.13,
 20.12*f*, 20.12*t*, 20.13*f*
 in ultracold collisions, **IV**:20.28–20.31,
 20.30*f*, 20.31*f*
- Laser damage threshold (LDT), of coatings,
IV:7.13–7.14, 7.18
- Laser detection and ranging (LADAR), **II**:31.28,
 31.30
- Laser diodes, **V**:13.3–13.24
 double heterostructure, **V**:13.3*f*, 13.3–13.8,
 13.4*f*, 13.6*f*, 13.7*f*
 noise characteristics of, **V**:13.18–13.24,
 13.19*f*–13.21*f*
 operating characteristics of, **V**:13.8–13.13,
 13.10*f*
 transient response of, **V**:13.13–13.18, 13.14*f*,
 13.16*f*, 13.17*f*
- Laser direct write (LDW) fabrication,
I:22.19–22.23, 22.20*f*–22.22*f*
- Laser field, spectral properties of the,
II:23.28–23.31
- Laser fields, **IV**:14.13
- Laser gain media, **II**:16.4–16.19
 dielectric solid-state, **II**:16.32–16.34
 gaseous, **II**:16.30–16.31, 16.31*f*
 liquid, **II**:16.31–16.32, 16.32*f*
 properties associated with, **II**:16.4–16.19
 emission linewidth and line broadening
 of radiating species, **II**:16.4–16.7,
 16.6*f*, 16.7*f*
 energy levels and radiation, **II**:16.4
 gain saturation, **II**:16.10
 optimization of output coupling from
 laser cavity, **II**:16.13, 16.14
- Laser gain media, properties associated with
(*Cont.*):
 population inversions, **II**:16.8–16.10,
 16.12–16.13, 16.13*f*, 16.14*f*
 pumping techniques to produce
 inversions, **II**:16.14–16.19
 stimulated absorption and emission,
II:16.7–16.8, 16.8*f*
 threshold conditions with mirrors,
II:16.10–16.12, 16.11*f*
 in vacuum, **II**:16.36–16.37, 16.37*f*
- Laser-generated plasmas, **V**:56.1–56.10
 and Bremsstrahlung radiation, **V**:56.8–56.10
 and recombination radiation, **V**:56.10
 sources of, **V**:56.1
 and spectral line emission, **V**:56.2–56.8
- Laser guide star (LGS) sensing, **V**:5.21, 5.23,
 5.27 (*see also* Laser beacons)
- Laser Interferometer Gravitational Wave
 Observatory (LIGO), **IV**:23.1, 23.7
- Laser Light Detection and Ranging (LIDAR)
 systems, **IV**:23.1
- Laser light sources, **IV**:5.59
- Laser linewidth, **II**:23.34–23.35
- Laser locking, frequency discriminators for,
II:22.12–22.14
- Laser master equation, **II**:23.19*f*, 23.19–23.20
- Laser mode locking, **V**:7.38–7.39
- Laser modes, coherence theory and, **I**:5.23
- Laser noise, **I**:35.12, 35.24
- Laser phase-transition analogy, **II**:23.35–23.40,
 23.37*t*, 23.38*f*, 23.39*f*
- Laser photocoagulation, **III**:14.26
- Laser photon statistics, **II**:23.22–23.26, 23.23*f*,
 23.25*f*
- Laser power combining, **IV**:12.30, 12.31
- Laser power measurement, **II**:34.32
- Laser power modulation (LPM), **I**:35.17–35.19
- Laser pulses, **IV**:21.35–21.36, 21.36*f*, 21.38*f*,
 21.38–21.39
- Laser radar (LIDAR) systems, **I**:30.51
- Laser ray tracing, **III**:16.7
- Laser resonators, **II**:16.20*f*, 16.23*f*–16.26*f*,
 16.23–16.26
- Laser scanning, **II**:40.54
- Laser scanning confocal microscopes, **III**:17.1,
 17.3, 17.7–17.8
- Laser scanning ophthalmoscopes, **III**:17.3
- Laser scribing, **II**:17.24–17.25
- Laser speckle, **I**:9.14

- Laser stabilization, **II**:22.1–22.24
 about, **II**:22.1
 Allan Deviation, **II**:22.2–22.3
 and frequency discriminators for laser locking, **II**:22.12–22.14
 frequency vs. time drift, **II**:22.2
 future directions for, **II**:22.23–22.24
 and optical cavity-based frequency discriminators, **II**:22.14–22.16, 22.17*f*
 quantifying frequency stability, **II**:22.2
 and quantum resonance absorption, **II**:22.16, 22.17
 representative/example designs of, **II**:22.20–22.23, 22.22*f*
 and servos, **II**:22.5–22.12
 Bode representation of servos, **II**:22.5–22.6, 22.6*f*
 closed-loop performance, **II**:22.8
 closed-loop stability issues, **II**:22.8–22.12, 22.9*f*–22.1*f*, 22.10–22.12
 measurement noise, **II**:22.7–22.8
 phase and amplitude responses, **II**:22.6*f*, 22.6–22.7, 22.7*f*
 spectral noise density, **II**:22.3–22.5
 and transducers, **II**:22.17–22.20
- Laser Stark spectroscopy, **I**:31.27*f*, 31.27–31.29, 31.28*f*
- Laser stripe structures, **II**:19.8, 19.9*f*
- Laser technology, for strong field interactions, **IV**:21.4*f*, 21.4–21.5
- Laser-assisted chemical etching (LACE), **I**:22.45
- Laser-Doppler interferometers, **I**:32.12–32.13, 32.13*f*
- Laser-feedback interferometers, **I**:32.13–32.14, 32.14*f*
- Laser-heated pedestal growth (LHPG) technique, **V**:12.9–12.10, 12.10*f*
- Laser-induced breakdown (LIB), **IV**:19.6–19.9, 19.8*f*
- Laser-induced damage (LID), **IV**:19.1–19.11
 avoidance of, **IV**:19.5–19.6
 and critical NLO parameters, **IV**:19.9*f*, 19.9–19.11, 19.10*f*
 estimates of, **IV**:19.2
 mechanisms of, **IV**:19.6–19.9, 19.8*f*
 and nonlinear optical effects, **IV**:19.5
 package-induced, **IV**:19.4–19.5
 surface damage, **IV**:19.2–19.4
- Laser-induced fluorescence, **V**:3.21
- Laser-induced-breakdown spectroscopy (LIBS), **V**:3.39
- Lasers and Optronics Buying Guide*, **II**:15.14
- LASIK (laser-assisted keratomileusis) surgery, **III**:1.6, 12.14, 16.13–16.14, 16.14*f*, 23.2
- Lasing, without inversion, **II**:23.40–23.42, 23.41*f*
- Lasing without inversion (LWI), **IV**:14.1, 14.3–14.4, 14.18–14.19
- Latent image (LI), **II**:30.7
- Latent images, in xerographic systems, **I**:34.1–34.4, 34.2*f*–34.3*f*
- Latent-image speck, **II**:29.5
- Lateral aberration curves (*see* Transverse ray plots)
- Lateral aniseikonia, **III**:13.17
- Lateral antiblooming, **II**:32.9, 32.10*f*
- Lateral chromatic aberration (*see* Transverse chromatic aberration)
- Lateral color, **I**:1.91–1.92, 29.14, 29.37; **II**:1.14, 2.5, 2.5*f*
- Lateral geniculate nucleus (LGN), **III**:2.12–2.14
 anatomy of, **III**:2.12–2.13, 2.13*f*
 physiology of, **III**:2.13–2.14
- Lateral magnification, **I**:17.5
- Lateral *pin* photodetectors, **II**:25.15*f*, 25.15–25.16
- Lateral resolution, of microscopes, **I**:28.17–28.19
- Lateral-collection photodetectors, **II**:26.4*f*
- Lateral-shearing interferometers, **II**:12.14, 13.9–13.12, 13.11*f*, 13.12*f*
- Lathe assembly technique, **II**:6.7–6.8, 6.8*f*
- Lattice absorption, semiconductor, **IV**:5.13–5.20
 impurity-related vibrational optic effects, **IV**:5.17–5.20, 5.18*f*–5.19*f*, 5.20*t*, 5.21*f*
 multiphonon absorption, **IV**:5.16–5.17, 5.17*t*
 phonons, **IV**:5.13–5.16, 5.15*f*
- Lattice vibrations:
 of crystals and glasses, **IV**:2.11–2.12, 2.76*t*–2.77*t*
 linear-chain model of, **IV**:5.8
 optical, **IV**:20.31–20.34, 20.32*f*–20.34*f*
 in solids, **IV**:8.16–8.18, 8.17*f*, 8.19*t*–8.20*t*
- Lattice-matched compositions, SOAs and, **V**:19.11
- Lattice-matched epitaxial layers, of fiber optic devices, **V**:13.2
- Lau effect, **I**:6.7–6.8, 6.8*f*

- Laue crystals, **V**:35.3
 Laue equation, **V**:39.1
 Laue geometry, for crystal monochromators, **V**:39.2*f*, 39.4–39.6
 Laue lenses, multilayer [see Multilayer Laue lenses (MLLs)]
 Laue phase retarders, **V**:43.6
 Laue transmission, **V**:26.8, 26.9*f*, 63.24, 63.26
 Laue-diffracting crystals, **V**:43.2
 Laurent half shades, **I**:13.56
 Layer adding method, multiple scattering and, **V**:3.21
 LBLRTM (code), **V**:3.23
 Le Grand eye model, **III**:2.3
 Lead salt lasers, **II**:19.7–19.8
 Lead selenide (PbSe) detectors, **II**:24.76*f*, 24.78, 24.79, 24.79*f*–24.82*f*
 Lead sulfide (PbS) photoconductors, **II**:24.73–24.74, 24.74*f*–24.77*f*, 24.74*t*
 Lead tin telluride (PbSnTe) photovoltaic detectors, **II**:24.92, 24.93*f*
 Lead titanate (PbTiO₃), **IV**:2.40*t*, 2.45*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.64*t*
 Lead vapor, **IV**:14.15, 14.16
 Lead-in wires, lightbulb, **II**:40.29, 40.29*f*
 Leakage channel fibers, **V**:25.21
 Leakage current, **II**:24.19–24.20, 32.10–32.12, 32.11*f*
 Leaky waveguides, **I**:21.3
 Leaky-mode arrays, **II**:19.29
 Least-squares method, **II**:3.17–3.19
 LEDs, as light source, 5.15
 Left circular polarizers, Mueller matrices for, **I**:14.10*t*
 Left half-wave circular retarders, Mueller matrices for, **I**:14.12*t*
 Left-circularly polarized light, **I**:12.27
 Left-handed coordinate systems, **I**:12.6
 Legacy films, **II**:30.23–30.25
 Legal traceability (of calibration), **II**:34.21
 Legendre polynomials, **II**:11.5, 11.30; **V**:46.6
 Legendre transformations, **I**:1.13, 1.15, 1.16
 Legendre-Fourier (L-F) polynomials, **V**:45.6
 Leica Summitar lenses, **I**:17.28, 17.28*f*
 Leman prisms, **I**:19.3*t*, 19.13, 19.13*f*
 Length measurements, **II**:12.2–12.10
 interferometric measurement, **II**:12.5–12.10, 12.7*f*–12.9*f*
 stadia and range finders, **II**:12.2–12.4, 12.3*f*, 12.4*f*
 Length measurements (*Cont.*):
 standards for, **II**:12.2
 time-based and optical radar, **II**:12.4, 12.5, 12.6*f*
 Length-to-aperture (*L/A*) ratio, for prisms, **I**:13.7
 Lens(es), **I**:17.3–17.39; **II**:39.32–39.37, 39.33*f*
 Airy distribution of, **V**:40.3
 assembly adjustment of, **II**:6.8–6.11, 6.10*f*
 axial separations in, **I**:1.52
 bovine, **III**:19.8–19.11, 19.11*f*
 Bragg-Fresnel, **V**:40.9*f*, 40.9–40.10
 cardinal points of, **I**:1.44, 17.7
 cat, **III**:19.9
 catadioptric systems of, **V**:35.1
 catoptric systems of, **V**:35.1
 coating specifications for, **II**:4.10
 component specifications for, **II**:4.2, 4.3
 computer-assisted surfacing technologies for, **III**:12.9
 and concentrators, **II**:39.16–39.18, 39.18*f*
 conjugate matrices for, **I**:1.71
 for correction of refractive error, **III**:12.4
 data entry for, **II**:3.2–3.8
 defined, **I**:1.9
 dioptric systems of, **V**:35.1
 distortion from interocular anisomagnification, **III**:13.16–13.17
 effective focal length of, **I**:1.48, 22.10, 27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.25
 effects on vergence and phoria, **III**:13.25–13.27
 entrance pupil of, **I**:1.76, 17.8, 18.4–18.6, 29.37
 entrance window of, **I**:17.9
 equivalent, **I**:17.20
 exit pupil of, **I**:1.76, 17.8, 18.6, 29.37
 exit window of, **I**:17.9
 field size of, **I**:28.13
 fields of, **I**:1.74
 fish, **III**:19.6, 19.6*f*
 fish-eye, **III**:19.2–19.3
 F-number and numerical aperture of, **I**:17.9
 Gauss points of, **I**:1.44
 Gaussian analyses of, **I**:18.2–18.7, 18.3*f*, 18.5*f*, 18.7*f*
 geometrical optics for, **I**:1.74–1.85
 apertures, **I**:1.74–1.77, 1.75*f*
 cosine-to-the-fourth approximation, **I**:1.81
 field lenses, **I**:1.82, 1.82*f*

Au: There is any requirement of volume no. for “LEDs, as light source, 5.15”

- Lens(es), geometrical optics for (*Cont.*):
 fields, **I**:1.74, 1.77, 1.84
 F-number, **I**:1.79
 focus and defocus, **I**:1.82–1.85, 1.83*f*
 irradiance, **I**:1.80
 power per pixel, **I**:1.80
 pupils, **I**:1.76*f*, 1.76–1.79, 1.78*f*
 telecentricity, **I**:1.83–1.84
 total lens étendue, **I**:1.81
 vignetting, **I**:1.81, 1.81*f*, 1.82
 gibbon, **III**:19.14, 19.14*f*
 GRIN-rod, **V**:18.7, 18.8, 18.8*f*, 18.10, 18.11*f*
 guinea pig, **III**:19.8
 of human eye, **III**:1.3*f*, 19.5–19.6 (*see also*
 Cornea; Crystalline lens; Gradient index
 optics)
 humidity specifications for, **II**:4.10
 for II electronic imaging, **II**:31.5–31.6
 image formation in, **I**:17.5*f*, 17.5–17.8, 17.6*f*
 image specifications for, **II**:4.3, 4.6–4.8, 4.8*f*
 inverses of, **I**:1.71
 Kerr, **V**:7.11, 7.39
 Kumahov capillary, **V**:53.9
 lacrimal, **III**:20.2, 20.12*f*–20.14*f*, 20.12–20.15
 literature on nonimaging, **II**:39.32, 39.33
 for magnifiers, **I**:17.9–17.10
 for Maxwellian viewing systems, **III**:5.4
 for microscopes, **I**:28.9–28.17 (*see also*
 Reflective and catadioptric lenses)
 compound, **I**:17.10
 objective lenses, **I**:28.9–28.15, 28.10*t*,
 28.11*t*, 28.12*f*, 28.13*t*, 28.14*f*–28.16*f*
 oculars, **I**:28.16–28.17
 multicomponent assemblies of (*see*
 Multicomponent lens assemblies)
 multilayer Laue, **V**:42.1–42.17
 with curved interfaces, **V**:42.14, 42.15*f*
 history of, **V**:42.2*f*, 42.2–42.4, 42.3*f*
 instrumental beamline arrangement
 and measurements, **V**:42.9*f*–42.12*f*,
 42.9–42.10
 limitations of, **V**:42.15–42.17, 42.16*f*–42.17*f*
 magnetron-sputtered, **V**:42.5–42.7,
 42.6*f*–42.8*f*
 Takagi-Taupin calculations, **V**:42.12–42.14
 volume diffraction calculations,
V:42.4–42.5, 42.5*f*
 wedged, **V**:42.12–42.13, 42.13*f*, 42.14*f*
 and x-ray/neutron optics, **V**:26.10
 natural stop shift of, **I**:22.3
- Lens(es) (*Cont.*):
 negative, **IV**:3.13
 in neutron optics, **V**:63.22–63.23, 63.23*f*
 nodal points of, **I**:1.48, 1.49
 octopus, **III**:19.7, 19.7*f*
 optical center point of, **I**:17.16, 17.17
 optical parameters for, **II**:4.9
 in optical systems, **III**:5.11
 perfect, **II**:4.4
 performance of, **I**:17.29–17.36, 17.30*f*–17.35*f*
 for periscopes, **I**:18.19, 18.19*f*
 porcine, **III**:19.11, 19.12*f*
 positive-powered, **IV**:3.13
 primate, **III**:19.13*f*, 19.14, 19.14*f*
 pupils of, **I**:17.8–17.9
 rabbit, **III**:19.8
 rat, **III**:19.7–19.8
 rays in, **I**:1.35
 reflector/lens-array combinations,
II:39.34–39.37, 39.36*f*, 39.37*f*
 refractive x-ray, **V**:37.3–37.11
 applications of, **V**:37.11
 history of, **V**:37.3
 nanofocusing, **V**:37.8–37.11, 37.9*f*, 37.10*f*
 parabolic, **V**:37.4*f*, 37.4–37.8, 37.6*f*, 37.7*f*
 for scanners, **I**:18.13, 30.57–30.60,
 30.58*f*–30.60*f*
 shape factor of, **I**:17.12–17.13, 17.13*f*
 single element, **I**:17.12–17.17, 17.13*f*–17.16*f*,
 17.17*t*
 single-lens arrays, **II**:39.33–39.34, 39.34*f*
 spherocylindrical, **III**:12.4, 12.5
 stimulus value of, **III**:13.21
 stops of, **I**:17.8–17.9
 systems of, **I**:17.7, 17.20–17.22, 17.21*f*–17.22*f*
 tandem-lens arrays, **II**:39.34, 39.35*f*–39.37*f*
 thermal defocus of, **II**:8.4, 8.5*f*
 thermal focus shift of, **II**:8.2–8.4, 8.3*t*, 8.4*t*
 for vision correction (*see* Contact lenses;
 Intraocular lenses; Spectacles)
 Wyszecki, **III**:9.5
 zone plates as, **V**:40.3–40.4
 (*see also specific types of lenses*)
 Lens axis, **I**:1.32 (*see also* Optical axis)
 Lens capsule, **III**:21.1, 21.2, 21.21–21.22
 Lens design:
 aberration curves in, **II**:2.1–2.6
 software for, **II**:40.20
 and tolerancing calculations, **II**:5.10
 Lens housings, **IV**:3.15*f*, 3.15–3.16

- Lens law, **I**:17.8
 Lens pigment density:
 and color matches, **III**:10.9
 and variations in color matching, **III**:10.15
 Lens setup routine (in optical software), **II**:3.7
 Lens-coupled II SSAs, **II**:31.22*f*, 31.22–31.23
 Lensing, **IV**:16.22, 19.5
 Lenslets, monolithic, **I**:22.25*f*, 22.25–22.26
 Lensometer, **III**:12.9
 Lenticular absorption, **III**:1.9
 Lenticular fluorescence, **III**:1.21
 Leslie viscosity coefficients, **V**:8.24
 Letter acuity task, **III**:2.34
 Levels (tools), **II**:12.13*f*, 12.13–12.14, 12.14*f*
 Lever mechanism mountings, **II**:6.19, 6.19*f*
 LF5 glass (581409), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
 Li Li polarizing beam splitter, **IV**:7.72–7.73, 7.76*f*
 Life-span environmental radiation damage,
 III:14.22–14.23
 Lifetime, photon, **V**:20.1
 Lifetime classification, of photodetector
 materials, **II**:26.5
 Lifetimes, of electrons, **I**:10.6
 Light:
 absorption of, **V**:3.4*f*, 3.4–3.5
 defined, **III**:4.1, 23.3
 out-of-plane profile of, **V**:13.11
 propagation of, **V**:7.3
 retroreflection of guided, **V**:13.6–13.7
 spatial characteristics of, **V**:13.11–13.12, 13.46
 spectral characteristics of, **V**:13.12–13.13
 spectrum of, **II**:25.2
 speed of, **II**:12.2
 theory of interaction of atmosphere and,
 V:3.11–3.22
 inelastic optical processes, **V**:3.21–3.22
 Mie scattering, **V**:3.16–3.18, 3.17*f*–3.19*f*
 molecular absorption, **V**:3.12–3.15, 3.13*f*
 molecular emission and thermal spectral
 radiance, **V**:3.18, 3.20, 3.20*f*
 molecular Rayleigh scattering, **V**:3.15–3.16
 surface reflectivity and multiple scattering,
 V:3.21, 3.21*f*
 Light amplification by stimulated emission of
 radiation [*see* Laser(s)]
 Light detection and ranging (LIDAR) systems,
 II:25.12; **V**:3.38*f*, 3.38–3.39, 25.25
 Light distribution, **II**:40.5
 Light extraction, **II**:17.6–17.8, 17.7*t*
 Light field microscopy, **I**:28.53
 Light flint glass, **IV**:2.42*t*
 Light grasp (étendue), **I**:1.22, 1.81, 13.7
 Light loss factor (LLF), **II**:40.17
 Light modulation, **IV**:5.66*t*
 Light out vs. current in (L-I curve),
 V:13.9–13.11, 13.10*f*, 13.46–13.47
 Light output:
 in color CRTs
 spatial characteristics of, **III**:22.14–22.15
 spatial uniformity of, **III**:22.17*f*,
 22.17–22.18, 22.18*f*
 stability of, **III**:22.16*f*, 22.16–22.17, 22.17*f*
 in color LCDs
 spatial variations in, **III**:22.40
 temporal variations in, **III**:22.39–22.40
 Light pipe reflectometers, rotating,
 IV:5.62, 5.63*f*
 Light piping, **II**:30.6*f*, 30.6–30.7
 Light pollution, **II**:40.43, 40.62
 Light pressure force, **IV**:20.7
 Light propagation, in solids, **IV**:8.4–8.13
 anisotropic crystals, **IV**:8.8–8.11, 8.9*t*, 8.10*f*
 energy flow, **IV**:8.7–8.8
 interfaces, **IV**:8.11–8.13, 8.12*f*, 8.13*f*
 Maxwell's equations, **IV**:8.4–8.6
 wave equations and optical constants,
 IV:8.6–8.7
 Light quanta, **II**:23.6–23.9, 23.8*f*
 Light scattering, **II**:30.5–30.7, 30.6*f*
 Light shelves, **II**:40.48, 40.50*f*
 Light sources, **I**:5.6, 5.9–5.13; **II**:40.24–40.41,
 40.28*f*; **IV**:5.58–5.59
 applications for, **II**:40.26*t*
 carbon arc sources, **II**:40.40
 characteristics of, **II**:40.25*t*
 coherent, **III**:5.19, 5.21
 control of (*see* Optical generation of visual
 stimulus)
 daylight, **II**:40.40*f*, 40.40–40.41
 electrodeless lamps, **II**:40.36–40.37
 electroluminescent sources, **II**:40.37–40.39,
 40.38*f*, 40.38*t*, 40.39*f*
 fluorescent lamps, **II**:40.30–40.33, 40.31*f*, 40.32*f*
 glow lamps, **II**:40.39
 high-intensity discharge lamps,
 II:40.33*f*–40.35*f*, 40.33–40.36
 incandescent sources, **II**:40.25, 40.27–40.30,
 40.28*f*, 40.29*f*
 low-pressure sodium lamps,
 II:40.33*f*–40.35*f*, 40.33–40.36

- Light sources (*Cont.*):
 in Maxwellian viewing, **III**:6.2
 nuclear sources, **II**:40.39
 pure Xe arc lamps, **II**:40.39
 short arc sources, **II**:40.39
 types of, **II**:40.27*f*
 in vision laboratories, **III**:5.19, 5.19*t*, 5.20*t*
- Light stability, **II**:30.10
- Light stabilization, **II**:30.12–30.13
- Light trespass, **II**:40.43, 40.62
- Light-absorbing dye, **II**:30.7
- Lightbulbs, **II**:40.27–40.30
 base of, **II**:40.29, 40.29*f*
 CFL/fluorescent/miniature, **II**:40.28*f*
 elements of, **II**:40.29*f*
 shapes of, **II**:15.20*f*, 15.30*f*
 sizes of, **II**:15.30*f*
 types of, **II**:40.27*t*
- Light-emitting diodes (LEDs), **II**:17.1–17.34;
V:13.1, 13.36–13.42, 13.38*f*
 conversion of, luminous intensity to radiant
 intensity, **II**:36.13, 36.13*f*
 device structures of, **II**:17.8–17.15
 diffused homojunctions, **II**:17.9–17.10,
 17.10*f*, 17.11*f*
 double heterojunctions, **II**:17.13,
 17.13*f*–17.15*f*
 grown homojunctions, **II**:17.8, 17.9,
 17.9*f*
 single heterojunctions, **II**:17.12, 17.12*f*
 edge-emitting, **V**:13.40
 epitaxial technology for, **II**:17.21–17.23
 in integrated optics, **I**:21.2
 lamps with, **II**:40.37–40.39
 applications for, **II**:40.26*t*
 characteristics of, **II**:40.25*t*
 materials/emitted colors of, **II**:40.38*t*
 photonic crystal, **II**:40.39*f*
 structure of, **II**:40.38*f*
 and LED-based products, **II**:17.29*f*–17.31*f*,
 17.29–17.32
 and light extraction, **II**:17.6–17.8, 17.7*t*
 as light source, **III**:5.15
 and light-generation processes, **II**:17.2–17.6,
 17.3*f*–17.6*f*
 material systems for, **II**:17.15–17.19
 AlInGaP system, **II**:17.18, 17.19*f*
 Al_xGa_{1-x}As system, **II**:17.17, 17.17*t*
 blue LED technology, **II**:17.18, 17.19
 GaAs_{1-x}P_x system, **II**:17.15–17.17, 17.16*t*
- Light-emitting diodes (LEDs) (*Cont.*):
 octocouplers in, **II**:17.32*f*, 17.32–17.34
 operating characteristics of, **V**:13.40*f*,
 13.40–13.42
 and optical fibers, **V**:9.7
 and parallel matrix-vector multipliers,
I:11.18
 production levels for, **II**:17.2
 quality/reliability of, **II**:17.25–17.28
 and serial incoherent matrix-vector
 multipliers, **I**:11.17–11.18
 substrate technology for, **II**:17.20–17.21,
 17.20*t*
 surface-emitting, **V**:13.38*f*, 13.38–13.40
 transmission by, **I**:21.32
 and transmissive TFT LCDs, **V**:8.29–8.31,
 8.30*f*
 wafer processing for, **II**:17.23–17.25
 (*see also specific light-emitting diodes, e.g.:*
 High-brightness visible LEDs)
- Light-gathering power (étendue), **I**:1.22,
 1.81, 13.7
- Light-hole (LH) bands, **V**:13.27
- Lighting, **II**:40.1–40.71
 about, **II**:40.1–40.3
 for aging individuals, **III**:14.4
 for computer work, **III**:23.4–23.5, 23.5*t*
 exterior, **II**:40.61–40.62, 40.63*t*
 functions of, **II**:40.12–40.14, 40.13*f*–40.16*f*
 insufficient, **II**:40.9
 interior, **II**:40.55–40.61
 health-care facility lighting,
II:40.58–40.60, 40.60*t*
 industrial lighting, **II**:40.60–40.61, 40.61*f*
 office lighting, **II**:40.55, 40.56*t*
 residential lighting, **II**:40.57, 40.58, 40.59*t*
 retail lighting, **II**:40.55–40.57, 40.56*t*–40.58*t*
 perception of, **II**:40.4*f*, 40.4–40.6
 for transportation, **II**:40.63–40.71
 roadway lighting, **II**:40.67, 40.69–40.71,
 40.70*t*, 40.71*t*
 vehicular lighting, **II**:40.63–40.67, 40.64*f*,
 40.65*t*, 40.66*f*, 40.66*t*, 40.68*t*, 40.69*f*
 vision biology, **II**:40.3–40.6
 (*see also Luminaires*)
- Lighting design, **II**:40.6–40.23
 and color, **II**:40.7–40.9
 and context, **II**:40.6
 and functions of lighting, **II**:40.12–40.14,
 40.13*f*–40.16*f*

- Lighting design (*Cont.*):
 geometries in, **II**:40.13*f*, 40.14–40.15, 40.15*f*,
 40.16*f*
 goals of, **II**:40.6
 and illuminance, **II**:40.7, 40.7*t*
 and properties of objects and impact,
II:40.16
 system layout and simulation in,
II:40.16–40.23, 40.18*f*, 40.22*f*–40.24*f*
 and visual discomfort, **II**:40.9–40.12, 40.11*t*
- Lighting Design and Application (IESNA),
II:39.8
- Lighting design software, **II**:40.21
- Lighting geometries, **II**:40.13*f*, 40.14–40.15,
 40.15*f*, 40.16*f*
- Lighting Handbook (IESNA), **II**:36.7, 40.17
- Lighting measurement, **II**:40.51–40.54
 goniometers, **II**:40.52–40.53, 40.53*t*, 40.54*f*
 illuminance meters, **II**:40.51, 40.52*f*
 luminance meters, **II**:40.52, 40.52*f*
 reflectometers, **II**:40.52
 surface measurement systems, **II**:40.53, 40.54
- Lighting system layout and simulation,
II:40.16–40.23, 40.18*f*
 computer graphics software for,
II:40.21–40.23, 40.22*f*, 40.23*f*
 IGES standard for, **II**:40.19
 optical analysis and design software for,
II:40.20
 optical design and analysis software for,
II:40.20–40.21
 software tools for, **II**:40.18–40.23,
 40.22*f*–40.24*f*
 source modeling software for, **II**:40.19–40.20
 STEP standard for, **II**:40.19
- Light-measuring polarimeters, **I**:15.3–15.4,
 15.11–15.13
- Lightness (term), **II**:40.4
- Light-output degradation, **II**:17.26–17.28,
 17.28*f*
- Lightpipes, **II**:39.27–39.32
 angular uniformity of, **II**:39.31
 applications for, **II**:39.32
 length of, **II**:39.30
 periodic distributions of, **II**:39.30
 shapes of, **II**:39.27–39.30, 39.28*f*–39.30*f*
 solid vs. hollow, **II**:39.30–39.31
 tapered, **II**:39.12–39.13, 39.13*f*, 39.31*f*,
 39.31–39.32
- LightTools (optical software), **II**:7.26
- Light-trap silicon photodiodes, **II**:34.30
- Limbal relaxing incisions, **III**:21.13
- Limbus, **III**:21.1
- Limiters:
 cascaded, **IV**:13.6
 optical, **IV**:12.32
 self-protecting, **IV**:13.9
 tandem, **IV**:13.6, 13.6*f*
- Limiting, optical, **IV**:13.2 (*see also* Passive
 optical limiting)
- Limiting noise, site of, **III**:11.24
- Limiting resolution, of EBSSAs, **II**:31.27
- Linac Coherent Light Source (LCLS), **V**:58.1
- Line defects, **IV**:9.12–9.13
- Line edge roughness (LER), in extreme
 ultraviolet lithography, **V**:34.6
- Linear attenuation coefficient, **V**:31.1, 31.2
- Linear colorimetry models, **III**:10.26–10.27,
 10.27*f*
- Linear diattenuation and diattenuators,
I:14.8, 14.17
- Linear dispersion, **V**:38.7
- Linear electro-optic (Pockels) effect,
I:21.9–21.10; **IV**:12.2–12.3
 and electro-optic modulators, **V**:7.6–7.11,
 7.8*t*, 7.11
 and liquids, **V**:7.34
 in OTDM networks, **V**:20.1, 20.18, 20.19
- Linear image sensor arrays, **II**:32.2, 32.21–32.24,
 32.22*f*, 32.23*f*
- Linear lasers, **II**:20.18–20.19, 20.19*f*
- Linear magnification, **I**:18.4
- Linear models (matrix algebra), **III**:10.1,
 10.46–10.47
- Linear optical absorption, **IV**:16.18
- Linear optical amplifiers (LOAs), **V**:19.14, 19.27
- Linear optical properties (of semiconductors),
IV:5.11–5.39, 5.12*f*, 5.13*t*
 free carriers, **IV**:5.33–5.36, 5.35*f*–5.37*f*
 impurity and defect absorption, **IV**:5.37–5.39,
 5.38*f*, 5.39*f*
 interband absorption, **IV**:5.21–5.33
 absorption near fundamental edge,
IV:5.21*f*, 5.21–5.22
 direct transitions, **IV**:5.22*f*–5.23*f*, 5.22–5.23
 excitons, **IV**:5.25–5.29, 5.26*t*, 5.27*f*–5.28*f*
 high-energy transitions above fundamental
 edge, **IV**:5.29–5.33, 5.30*f*–5.34*f*
 indirect transitions, **IV**:5.23–5.24, 5.24*f*–5.25*f*
 polaritons, **IV**:5.29

- Linear optical properties (of semiconductors)
(*Cont.*):
 lattice absorption, **IV**:5.13–5.20
 impurity-related vibrational optic effects,
 IV:5.17–5.20, 5.18*f*–5.19*f*, 5.20*t*, 5.21*f*
 multiphonon absorption, **IV**:5.16–5.17,
 5.17*t*
 phonons, **IV**:5.13–5.16, 5.15*f*
 and models of dielectric function, **IV**:5.12
- Linear optics, **III**:17.1
- Linear perspective, **III**:13.3
- Linear polarization (OCT), **III**:18.21
- Linear polarization analyzers, **V**:43.4, 43.4*f*
- Linear polarization sensitivity, **I**:14.17
- Linear polarizers, **I**:14.9, 14.10*t*, 15.7;
V:43.2–43.3, 43.3*f*
- Linear regime, for rare-earth-doped fiber
 amplifiers, **V**:14.3
- Linear systems, coherence theory for, **I**:6.3–6.4
- Linear tomography, **V**:31.5, 31.5*f*
- Linear variable differential transformers
 (LVDTs), **II**:10.12
- Linear visual mechanisms, **III**:11.3
- Linear-chain model, of lattice vibrations, **IV**:5.8
- Linearity (in general):
 in paraxial matrix methods, **I**:1.66
 of photoemissive detectors, **II**:24.41
 of scatterometers, **V**:1.15
 of systems of revolution, **I**:1.41
 and transfer functions, **I**:4.2
- Linearity of color-opponent mechanisms,
III:11.66–11.70
- Bezold-Brücke effect and invariant hues,
III:11.67*f*, 11.67–11.68
- color appearance
 and chromatic adaptation, **III**:11.68
 and chromatic detection and
 discrimination, **III**:11.69
 and habituation, **III**:11.69–11.70
 luminance and brightness, **III**:11.70, 11.70*f*
 tests of linearity, **III**:11.66–11.67
- Linearization, in optical systems, **III**:5.15
- Line-by-line transmission programs, **V**:3.23, 3.24*f*
- Line-scanned imaging systems, **II**:31.29
- Line-spread function (LSF), **III**:1.21
 in diffraction-limited eye, **III**:1.13
 in ophthalmoscopic methods, **III**:1.23
- Linewidth:
 of lasers, **V**:9.8
 spectral density and, **II**:22.4–22.5
- Linewidth enhancement factor, **V**:13.17, 13.20
- Link budget analysis, for fiber optic
 communication:
 installation loss, **V**:15.6–15.7, 15.8*t*
 optical power penalties, **V**:15.8–15.17, 15.10*f*,
 15.11*f*
- Linnik interference microscopes, **I**:28.36–28.38,
 28.37*f*
- Linnik interferometers, **V**:46.2
- Lin-perp-lin polarization gradient cooling,
IV:20.17–20.18, 20.18*f*, 20.19*f*
- Liou Brennan schematic eye, **III**:15.6*f*, 15.7
- Liouville's theorem, **I**:1.22; **V**:54.11, 54.15,
 54.16, 63.15, 63.22
- Lippich-type prisms, **I**:13.6, 13.12–13.13
 Glan-Taylor, **I**:13.7, 13.9*n*, 13.10, 13.10*f*,
 13.11*f*, 13.12–13.14
 half-shade, **I**:13.12*n*, 13.56
 Marple-Hess, **I**:13.12, 13.13
- Lippmann emulsions, **II**:29.4
- Lippmann-Bragg holographic mirrors,
IV:7.50
- Liquid crystal (LC) cells, **I**:15.32*f*, 15.33*f*,
 15.33–15.35, 15.34*t*; **V**:8.25–8.28, 8.26*f*,
 8.27*f*, 8.28*f*
- Liquid crystal displays (LCDs), **III**:22.34–22.40;
V:8.29–8.35, 61.3
 color, **III**:22.37–22.40
 colorimetry of color pixels, **III**:22.38–22.39
 controls and input standards, **III**:22.39
 geometry of color pixels, **III**:22.37*f*,
 22.37–22.38
 spatial variations in output, **III**:22.40
 temporal variations in output,
III:22.39–22.40
 and Computer Vision Syndrome, **III**:23.7–23.8
 defined, **III**:23.2
 monochrome, operational principles of,
III:22.34–22.37, 22.35*f*, 22.36*f*
 in optical systems, **III**:5.16
 reflections from, **III**:23.6
 reflective, **V**:8.31*f*, 8.31–8.32
 transmissive TFT, **V**:8.29–8.31, 8.30*f*
 transreflective, **V**:8.32–8.35, 8.33*f*, 8.34*f*
- Liquid crystal (LC) lenses, **I**:22.40–22.41, 22.42*t*
- Liquid crystal on silicon (LCOS) panels,
I:15.28
- Liquid crystal retarders, **I**:15.21–15.23, 15.22*f*
- Liquid crystal variable retarders (LCVRs),
I:15.21–15.23, 15.22*f*

- Liquid crystals (LCs), **IV**:13.12, 13.12*f*
 composition of, **V**:8.2–8.4, 8.3*f*
 dielectric properties of, **V**:8.14–8.18, 8.15*f*,
 8.18*f*
 elastic properties of, **V**:8.22–8.23, 8.23*f*
 limitations of, **V**:8.37
 optical properties, **V**:8.19*f*, 8.18–8.22,
 8.20*f*, 8.22*f*
 phase transitions of, **V**:8.13, 8.13*f*, 8.14*f*
 phases of, **V**:8.8*f*–8.13*f*, 8.11–8.12
 physical properties of, **V**:8.13–8.23
 in polymer/liquid crystal composites,
V:8.36*f*–8.37*f*, 8.36–8.37
 types of, **V**:8.4*f*–8.6*f*, 8.4–8.8, 8.7*t*–8.9*t*
 viscosities of, **V**:8.24*f*, 8.23–8.25
- Liquid Encapsulated Czochralski (LEC)
 technique, **II**:17.20, 17.21
- Liquid immersion development, in xerographic
 systems, **I**:34.9, 34.10
- Liquid laser gain media, **II**:16.31–16.32, 16.32*f*
- Liquid lenses, **I**:22.37–22.41, 22.38*f*–22.41*f*, 22.42*t*
- Liquid phase epitaxy (LPE), **II**:17.21–17.22,
 19.6, 19.19
- Liquid spatial light modulators (LC-SLMs),
III:15.11
- Liquid-crystal-on-silicon (LCoS) displays, **V**:8.32
- Liquid-metal anodes, **V**:54.13
- Liquid-phase epitaxial (LPE) growth, **I**:21.17
- Liquids, in electro-optic modulators, **V**:7.34
- Lissajous-type figures, from Risley prisms,
I:19.25, 19.26*f*
- Lit-appearance modeling, **II**:40.21–40.23,
 40.23*f*, 40.24*f*
- Lithium fluoride (LiF), **IV**:2.39*t*, 2.44*t*, 2.48*t*,
 2.57*t*, 2.69*t*
- Lithium iodate (α -LiIO₃), **IV**:2.39*t*, 2.46*t*, 2.48*t*,
 2.57*t*, 2.75*t*
- Lithium niobate (LiNbO₃), **I**:21.16–21.17,
 21.28, 21.33, 21.39; **IV**:2.39*t*, 2.46*t*, 2.48*t*,
 2.52*t*, 2.57*t*, 2.63*t*, 2.75*t*, 12.14, 17.1,
 17.4–17.13, 17.6*f*–17.11*f*; **V**:6.30,
 6.31, 6.31*t*
- Lithium niobate (LiNbO₃) modulators, **V**:13.2,
 13.48–13.55, 13.49*f*
 electro-optic effect, **V**:13.49–13.51, 13.51*f*
 electro-optic polymer, **V**:13.55
 high-speed operation of, **V**:13.52–13.53
 insertion loss in, **V**:13.53
 as Mach-Zehnder modulators,
V:13.51–13.52, 13.54*f*, 13.54–13.55
- Lithium niobate (LiNbO₃) modulators (*Cont.*):
 phase modulation by, **V**:13.51
 photorefractivity and optical damage of,
V:13.54
 polarization independence of, **V**:13.53
 Y-branch interferometric, **V**:13.51–13.52
- Lithium [⁶Li(*n*, α)] reaction, **V**:63.31
- Lithium tantalate (LiTaO₃), **I**:21.17; **IV**:12.14,
 17.14–17.15, 17.15*f*
- Lithium triborate (LiB₃O₅) (LBO), **IV**:2.39*t*,
 2.46*t*, 2.51*t*, 2.56*t*, 2.63*t*, 2.75*t*, 17.1
- Lithium-calcium-aluminum fluoride
 (LiCaAlF₆) (LiCAF), **IV**:2.39*t*, 2.46*t*, 2.51*t*,
 2.57*t*, 2.63*t*
- Lithium-drifted silicon x-ray detectors,
V:60.6
- Lithographic etching, **II**:18.3*f*, 18.3–18.4
- Lithographic projection lens, **II**:6.10*f*
- Lithography:
 deep x-ray, **I**:22.7
 electron-beam, **V**:40.8
 extreme ultraviolet, **V**:34.1–34.6, 34.2*t*,
 34.3*f*–34.6*f*
 extreme ultraviolet-interferometric,
V:34.4–34.5, 34.5*f*
 high aspect ratio microlithography, **V**:61.3
 and holography, **I**:33.22*f*–33.24*f*, 33.22–33.24
 interferometric, **V**:34.4
 and miniature and micro-optics,
I:22.18–22.25
 with gray-scale masks, **I**:22.23*f*–22.25*f*,
 22.23–22.25
 laser direct write fabrication, **I**:22.19–22.23,
 22.20*f*–22.22*f*
 optical, **V**:34.1
- Littrow configuration, of dispersive prisms,
I:20.7*f*, 20.10
- Littrow mirrors, **I**:20.4
- L-I-V measurements, for SOAs, **V**:19.17, 19.17*f*
- LLF1 glass (548458), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- LLL sensitivity, **II**:31.19, 31.20
- Lloyd's mirror, **I**:2.16, 2.17*f*, 2.18
- Lobster-eye (LE) optics, **V**:48.2*f*, 48.2–48.4,
 48.3*f*, 49.3*f*–49.4*f*, 49.3–49.4
- Lobster-ISS system, **V**:50.7
- Local area networks, **II**:19.3
- Local area networks (LANs), **V**:9.14, 21.7
 power splitters and couplers in, **V**:18.2, 18.3*f*
 standards for, **V**:23.2, 23.6, 23.7

- Local characterization methods (LCM) (color CRTs), **III**:22.24–22.27
 individual colors, **III**:22.24–22.26
 inverses, **III**:22.25
 local regions of color, **III**:22.24–22.26
 one-dimensional color spaces, **III**:22.25, 22.26
 out-of-gamut colors, **III**:22.25
- Local density approximation (LDA), **IV**:5.5
- Local gain per unit length, of lasers, **V**:13.8
- Local oscillators, **V**:9.13
- Local shift variance, in grazing incidence x-ray optics, **V**:44.14
- Local vibrational modes (LVM), **IV**:5.17, 5.18, 5.19*f*, 5.20, 5.20*f*
- Localization, **III**:4.14, 4.15, 8.4
- Localization, in volume scattering, **I**:9.13–9.17, 9.14*f*
- Localized avalanche breakdown, **II**:17.28
- Localized vibration, **IV**:5.82*f*, 5.83
- Lockheed Martin, **IV**:6.46
- Lock-in amplifiers, **II**:27.14, 27.14*f*, 27.15, 38.10; **IV**:5.64
- Log-amplitude structure function, **V**:4.5
- Long duration exposure facility (LDEF), **IV**:6.17
- Long trace profiler (LTP), **V**:46.4, 46.5, 46.5*f*
- Long wavelength lasers, **II**:19.8
- Long wavelength QW lasers, **II**:19.17*f*, 19.17–19.18
- Long-exposure images, **V**:4.3–4.7
- Long-exposure MCF, **V**:4.10
- Longitudinal aberrations, **I**:1.87
- Longitudinal acoustic (LA) phonons, **IV**:5.24, 5.25*f*
- Longitudinal (axial) chromatic aberration (LCA), **III**:1.19, 1.20, 1.28, 8.8, 15.22
- Longitudinal electro-optic modulators, **V**:7.16, 7.17, 7.17*f*
- Longitudinal horopter, **III**:13.8
- Longitudinal Kerr effect, **IV**:18.11–18.15, 18.12*f* (*see also* Self-phase modulation)
- Longitudinal laser modes, **II**:16.19*f*, 16.20, 16.20*f*
- Longitudinal magnification, **I**:1.28, 1.52, 17.5
- Longitudinal optic (LO) phonons, **IV**:5.24, 5.25*f*, 5.79, 5.79*f*, 5.80
- Longitudinal relaxation rate and time, **IV**:11.5
- Longitudinal spatial modulation (LSM), **V**:6.12, 6.17
- Longitudinal-mode (LO) frequencies, for crystals and glasses, **IV**:2.11, 2.12
- Long-period gratings (LPGs), **V**:21.39, 21.39*f*, 24.8–24.13, 24.9*f*–24.12*f*, 24.11*t*, 24.13*t*
- Long-range motion discrimination mechanism, **III**:2.38
- Long-wave infrared (LWIR) AOTFs, **V**:6.42
- Long-wavelength infrared (LWIR), **II**:24.3, 33.3–33.5, 33.6*f*
- Lorentz model:
 of absorption, **IV**:4.4
 of dispersion, **IV**:8.14, 8.21
- Lorentzian broadening, **V**:3.14
- Lorentzian distribution of frequencies, **I**:10.7
- Lorentzian line shape, **II**:16.6*f*
- Lorentzian lineshapes, of spectra, **V**:2.13, 56.4–56.7
- Loss:
 Akhieser, **V**:6.17
 bend, **V**:11.21*f*, 11.21–11.22
 connector, **V**:15.7, 15.8*t*
 coupling, **V**:18.1
 of fiber-optic components, **V**:18.1
 insertion, **V**:7.36, 13.53, 18.1
 installation, **V**:15.6–15.7, 15.8*t*
 optical power, **V**:21.13–21.14
 polarization-dependent, **V**:19.18, 21.18
 radiation induced, **V**:15.17
 splice, **V**:15.7, 15.8*t*
 transmission, **V**:15.7
- Loss hypothesis, **III**:10.14
- Loupes, eye, **I**:17.9–17.10
- Louvers, **II**:40.41, 40.45*f*, 40.46
- Low cost (LC) FDDI, **V**:23.3
- Low inertia scanners, **I**:30.43
- Low leakage guidance, **V**:11.16, 11.17
- Low temperature (LT), **IV**:18.3
- Low-beam headlamps, **II**:40.64*f*, 40.64–40.67, 40.65*t*, 40.66*f*, 40.66*t*
- Low-frequency fluctuations (LFF), of lasers, **V**:13.23
- Low-intensity carbon arc lamps, **II**:15.21–15.24, 15.24*f*, 15.28*f*
- Low-intensity reciprocity failure, of photographic films, **II**:29.12
- Low-level mechanisms (color vision), **III**:11.80
- Low-level-light (LLL) TV imaging, **II**:31.1–31.4
- Low-light-level television (LLTV) systems, **I**:26.3
- Low-order aberrations, **III**:15.1
- Low-order flux models, of radiative transfer, **I**:9.12–9.13, 9.13*f*
- Low-pressure enclosed arcs, **II**:15.35–15.47
 black-light fluorescent lamps, **II**:15.35, 15.36*t*
 electrodeless discharge lamps, **II**:15.36, 15.44

- Low-pressure enclosed arcs (*Cont.*):
 germicidal lamps, **II**:15.35
 hollow cathode lamps, **II**:15.35, 15.37*t*–15.43*t*,
 15.44*f*
 Pluecker spectrum tubes, **II**:15.47, 15.47*f*,
 15.47*t*
 spectral lamps, **II**:15.44, 15.45, 15.45*f*, 15.46*f*,
 15.46*t*
 Sterilamps, **II**:15.35, 15.36*f*
- Low-pressure lamps, **II**:15.36*f*
- Low-pressure sodium (LPS) lamps, **II**:40.33–40.36
 applications for, **II**:40.26*t*
 characteristics of, **II**:40.25*t*
 construction of, **II**:40.34*f*
 emission spectrum of, **II**:40.35*f*
- Low-Q cavities, **IV**:9.12
- Low-speed photographic films, **II**:30.18–30.20
- Low-temperature bolometers, **II**:24.31–24.32,
 24.32*f*, 24.33*f*, 28.5
- Low-temperature (LT) grown photoconductors,
II:26.23
- Low-temperature (LT) molecular beam epitaxy,
IV:18.21
- LOWTRAN program, **V**:3.23–3.24,
 3.25*f*–3.26*f*
- LOX8 glass, **IV**:6.57
- L-shell emission lines, of elements,
V:36.3*t*–36.8*t*
- Lucalox lamps, **II**:15.30, 15.31*f*
- Lu-Chipman polar decomposition,
I:14.39–14.40
- Lukosz-type super-resolving systems,
I:6.9*f*, 6.9–6.10
- Lumen (unit), **II**:36.6, 37.6, 39.2*t*
- Lumen lighting simulation, **II**:40.17
- Luminaires, **II**:40.24–40.50
 applications for, **II**:40.26*t*
 calculation of needed, **II**:40.16–40.17
 characteristics of, **II**:40.25*t*
 classification system for, **II**:40.43–40.45,
 40.43*t*, 40.44*f*
 defined, **II**:40.1
 design of, **II**:40.41–40.50
 conics shapes and intensity distribution,
II:40.43*t*
 etendue and source coupling,
II:40.41–40.42
 luminaire classification system,
II:40.43–40.45, 40.43*t*, 40.44*f*
 methods, **II**:40.42–40.43, 40.43*t*
- Luminaires (*Cont.*):
 light sources for, **II**:40.24–40.41,
 40.27*f*, 40.28*f*
 carbon arc sources, **II**:40.40
 daylight, **II**:40.40*f*, 40.40–40.41
 electrodeless lamps, **II**:40.36–40.37
 electroluminescent sources, **II**:40.37–40.39,
 40.38*f*, 40.38*t*, 40.39*f*
 fluorescent lamps, **II**:40.30–40.33, 40.31*f*,
 40.32*f*
 glow lamps, **II**:40.39
 high-intensity discharge lamps,
II:40.33*f*–40.35*f*, 40.33–40.36
 incandescent sources, **II**:40.25, 40.27–40.30,
 40.28*f*, 40.29*f*
 low-pressure sodium lamps,
II:40.33*f*–40.35*f*, 40.33–40.36
 nuclear sources, **II**:40.39
 pure Xe arc lamps, **II**:40.39
 short arc sources, **II**:40.39
 optics of, **II**:40.45–40.50
 for artificial sources, **II**:40.45*f*,
 40.45–40.47, 40.46*f*
 backlighting, **II**:40.47, 40.47*f*, 40.48*f*
 for daylight sources, **II**:40.47–40.50,
 40.49*f*–40.51*f*
- Luminance, **II**:37.4*t*, 37.5, 37.5*f*, 39.2*t*;
III:11.37–11.39
 calibration of, **II**:34.42
 cone numerosity, **III**:11.38
 defined, **II**:34.11, 34.40, 36.7, 40.1; **III**:2.29*n*,
 11.3, 23.3
 and illuminance, **II**:37.9
 of integrating cavities, **II**:39.26
 luminous efficiency, **III**:11.37
 luminous efficiency functions, **III**:11.37,
 11.38*f*
 multiple cone inputs, **III**:11.49, 11.50*f*
 multiple luminance signals, **III**:11.38–11.39,
 11.40*f*, 11.70, 11.70*f*
 in nonimaging optics, **II**:39.2*t*, 39.3
 uniformity of, **II**:40.7
 units of, **II**:34.43
 and visual acuity, **III**:4.11, 4.11*f*, 4.12
- Luminance contrast, **II**:40.6, 40.10
- Luminance mechanism, **III**:11.11
- Luminance meters, **II**:40.52, 40.52*f*
- Luminance pedestals, **III**:11.72, 11.74
- Luminance ratio, **II**:40.7
- Luminance ratios, for computer work, **III**:23.5

- Luminescence, **I**:31.17, 31.19*f*, 31.19–31.21
 Luminescence excitation spectrometers, **I**:31.11*f*, 31.11–31.12
 Luminescence spectrometers, **I**:31.5–31.12, 31.8*f*, 31.11*f*
 Luminescence spectroscopy, **IV**:5.69–5.75, 5.70*f*, 5.72*f*–5.75*f*
 Luminosity (étendue), **I**:1.22, 1.81, 13.7
 Luminosity function, **III**:10.10, 10.16
 Luminous efficacy, **II**:18.5
 Luminous efficiency, **II**:17.15; **III**:11.37
 and brightness, **III**:11.70, 11.70*f*
 and chromatic adaptation, **III**:11.48*f*
 variations in, **III**:11.33
 Luminous efficiency functions, **III**:10.13, 10.44–10.45, 11.33, 11.37, 11.38*f*
 Luminous energy, **II**:37.4*t*, 37.6
 Luminous exitance, **II**:37.4*t*, 37.5, 37.5*f*
 Luminous exposure, **II**:37.4*t*, 37.6
 Luminous flux, **II**:15.11, 34.10–34.11, 34.39, 34.42, 37.4, 37.4*t*, 37.6, 38.2; **III**:23.3
 Luminous flux density, **II**:34.11
 Luminous intensity, **II**:15.11, 34.11, 34.40, 37.4, 37.4*t*, 39.2*t*
 defined, **III**:23.3
 in stimulus specification, **III**:4.3–4.4
 Lump amplification, **V**:21.44
 Luneburg lens, **I**:1.21, 22.26; **III**:19.3 (*see also* Distributed-index planar microlenses)
 Lux (unit), **II**:34.43, 36.7, 36.7*t*; **III**:23.3
 Luxmeters, **II**:40.52*f*
 Lux-second, **II**:29.6
 Lyddane-Sachs-Teller (LST) relation, **IV**:2.11, 5.14, 8.17
 Lyman series, **V**:56.2
 Lyot coronagraphs, **I**:29.5
 Lyot stops, **I**:29.5, 29.5*f*, 29.37; **II**:7.8*f*–7.11*f*, 7.8–7.10
 Lyotropic liquid crystals (LCs), **V**:8.3, 8.5*f*
 Lysozyme diffraction, **V**:53.12*f*, 53.12–53.14, 53.13*f*
- MacAdam ellipses, **III**:10.40
 Mach-Zehnder devices, as fiber-based couplers, **V**:16.4–16.5, 16.5*f*
 Mach-Zehnder filters, **V**:14.6, 21.39
 Mach-Zehnder (MZ) interference filters, **I**:21.23, 21.24*f*, 21.25
 Mach-Zehnder (MZ) interference microscopes, **I**:28.36, 28.37*f*, 28.39
 Mach-Zehnder interferometers (MZIs), **I**:21.10, 21.12, 21.14, 21.16, 21.32, 21.40, 32.3, 32.3*f*, 33.5; **III**:5.24; **IV**:23.2*f*, 23.2–23.4
 and Bragg grating sensors, **V**:24.8
 and DPSK, **V**:21.33*f*, 21.33–21.34
 and electro-optic modulators, **V**:7.22, 7.24, 7.32, 7.38
 and fiber Bragg gratings, **V**:17.8
 integrated, **V**:20.18–20.19, 20.19*f*
 in OTDM networks, **V**:20.22, 20.23*f*
 SOAs in, **V**:19.31*f*, 19.31–19.32, 19.33*f*, 19.36
 and supercontinuum generation, **V**:11.23
 Mach-Zehnder (MZ) modulators, **I**:21.26, 21.27*f*, 21.28, 21.32, 21.34
 interferometric, **V**:13.51–13.52, 13.54*f*, 13.54–13.55, 13.63
 in WDM networks, **V**:21.30, 21.31*f*, 21.32*f*
 MacNeille polarizers, **IV**:7.70–7.72, 7.71*f*, 7.72*f*, 7.75*f*
 Macrofocal reflectors, **II**:39.11
 Macula lutea (yellow spot), **III**:7.1, 7.11
 Macular pigment (MP), **III**:1.9, 1.11, 14.9–14.10
 Macular pigment density:
 adjustments for individual differences, **III**:10.17
 and color matches, **III**:10.9
 and variations in color matching, **III**:10.15
 Magnesium, as *p*-type impurity, **II**:17.23
 Magnesium oxide (MgO), **V**:2.20, 2.20*f*, 2.22
 Magnesium-oxide doped stoichiometric lithium tantalite (MgO:sPPLT), **IV**:17.14–17.15, 17.15*f*
 Magnetic brush development, in xerographic systems, **I**:34.5*f*–34.7*f*, 34.5–34.7
 Magnetic circular dichroism (MCD), **I**:31.20, 31.21; **V**:55.7*f*, 55.7–55.9
 Magnetic circular polarization (MCP), **I**:31.21
 Magnetic field modulation (MFM), **I**:35.19, 35.19*f*, 35.20; **IV**:5.66*t*
 Magnetic permeability, of water, **IV**:1.16
 Magnetic resonance, optically detected, **I**:31.21–31.23, 31.22*f*, 31.23*f*
 Magnetic resonance imaging (MRI), **III**:19.1
 Magnetic shielding, **II**:27.10
 Magnetic traps, **IV**:20.21–20.23, 20.22*f*
 Magnetoabsorption, impurity, **IV**:5.51*f*, 5.51–5.52
 Magneto-optical (MO) disks, **I**:35.2
 Magneto-optical (MO) modulators, **I**:15.23

- Magneto-optical (MO) properties (of semiconductors), **IV**:5.39–5.52, 5.41*t*
 effect of magnetic field on energy bands, **IV**:5.40, 5.42, 5.42*f*
 interband effects, **IV**:5.42–5.46, 5.43*f*–5.45*f*, 5.47*f*
 intraband or free-carrier effects, **IV**:5.47–5.52, 5.48*f*–5.51*f*
 semiconductor nanostructures, **IV**:5.52
- Magneto-optical (MO) readout, **I**:35.21–35.24, 35.22*f*, 35.24*f*
- Magneto-optical (MO) recording, **I**:35.25–35.28, 35.26*f*, 35.27*f* (*see also* Optical disk data storage)
- Magneto-optical traps (MOTs), **IV**:14.18, 20.24*f*, 20.24–20.25, 20.26*f*
- Magnetopolarons, **IV**:5.50
- Magnetoreflection, **IV**:5.43, 5.44, 5.44*f*
- Magnetorheological finishing (MRF), **II**:9.5
- Magnetron sputtering, **IV**:7.11
- Magnetron-sputtered multilayer Laue lenses (MLLs), **V**:42.5–42.7, 42.6*f*–42.8*f*
- Magnification, **I**:1.28; **II**:1.4
 afocal, **I**:18.5–18.6
 angular, **I**:1.52, 1.78, 18.4
 axial (longitudinal), **I**:1.28, 1.52, 17.5
 with contact lenses, **III**:20.31–20.33
 relative spectacle, **III**:20.32–20.33, 20.33*f*
 spectacle, **III**:20.31–20.32
 dia-, **I**:1.23
 distortion by monocular magnification, **III**:13.13–13.16
 bifocal jump, **III**:13.15*f*, 13.15–13.16
 from convergence responses to prism, **III**:13.19
 discrepant views of objects/images, **III**:13.16
 motion parallax, **III**:13.14–13.15
 perspective distortion, **III**:13.13, 13.14*f*
 stereopsis, **III**:13.13, 13.14
 distortion from interocular anisomagnification, **III**:13.16–13.19
 aniseikonia, **III**:13.17–13.18
 interocular blur suppression with anisometropia, **III**:13.18–13.19
 lenses and prisms, **III**:13.16–13.17, 13.17*f*
 errors of eye alignment induced by, **III**:13.25–13.27
 by focal Gaussian lenses, **I**:1.52–1.53
 by gratings and monochromators, **V**:38.7
- Magnification (*Cont.*):
 lateral, **I**:17.5
 linear, **I**:18.4
 longitudinal, **I**:1.28
 and magnifiers, **I**:17.9–17.10
 in microscopes, **I**:28.3–28.4
 pupil, **I**:1.76
 pupil angular, **I**:1.78
 scan, **I**:30.5, 30.12–30.14
 secondary, **I**:29.12, 29.38
 in systems of revolution, **I**:1.42
 transverse, **I**:1.28, 1.50–1.51
 visual, **I**:1.28
- Magnifiers, first-order layout for, **II**:1.8
- Magnitude (optical aberrations), **III**:15.4
- Magnitude estimation, in psychophysical measurement, **III**:3.8
- Magnitude production, in adjustment experiments, **III**:3.5
- Magnocellular laminae, **III**:2.12, 2.13, 2.13*f*
- Maier-Saupe mean-field coupling constant, **V**:8.25
- Main event loop, **II**:3.7
- Maksutov objectives, **I**:29.19, 29.20
- Maksutov sphere, **II**:14.7, 14.8*f*
- Maksutov test, **II**:14.7–14.9, 14.8*f*
- Malus-Dupin principle, **I**:1.12
- Mammography, **V**:54.8, 59.3–59.4
- Manchester (biphase) coding, **V**:20.9, 20.9*f*
- Mandel Q_M parameter, **II**:23.26
- Mandelbaum's phenomenon, **III**:25.11–25.12
- Mandell-Moore bitoric lens guide, **III**:20.17–20.19, 20.18*t*, 20.19*f*
- Mangin elements, in objective design, **I**:29.6
- Mangin objectives, **I**:29.7, 29.24
- Mankiewicz Gebr. & Co., **IV**:6.35
- Manly Rowe fraction, **IV**:15.15
- Manufacturing error budget, for polymeric optics, **IV**:3.10
- Mapping, of object and image space, **I**:1.27
- Maréchal criterion, **III**:1.15
- Marginal rays, **I**:1.75, 17.8, 29.37
- “Marine snow,” **IV**:1.14, 1.29
- Markov approximation, **II**:23.21
- Markov random process approximation, for beam wander, **V**:3.31, 3.32
- Marple-Hess prisms, **I**:13.12, 13.13
- Martin Black, **IV**:6.3*t*, 6.12, 6.14, 6.15, 6.26*f*, 6.28*f*, 6.46, 6.47*f*, 6.51*f*, 6.53*f*

- Martin Black coating, **II**:7.14*f*–7.16*f*, 7.14–7.17, 7.23, 7.25*f*
- Martin Marietta, **IV**:6.48, 6.53
- Martin Optical Black, **IV**:6.5*f*
- Masers, micro- (*see* Micromasers)
- Mask layout, for binary optics, **I**:23.14–23.16, 23.15*f*, 23.15*t*
- Maskless lithography tool (MLT), **I**:22.23
- Masks (for extreme ultraviolet lithography), **V**:34.2, 34.6
- Masks (in experiments), **III**:3.2
- Mass attenuation coefficients, for photons and metals, **IV**:4.48, 4.48*t*, 4.49
- Mass density, of metals, **IV**:4.6
- Mass-transport process, for miniature and micro-optics, **I**:22.45, 22.46*f*
- Master groups, of zoom lenses, **I**:27.17
- Master oscillator power amplifier (MOPA) systems, **I**:21.30, 21.30*f*; **II**:19.41; **IV**:17.6; **V**:25.2
- free space designs, **V**:25.13–25.15, 25.15*f*
- monolithic designs, **V**:25.16*f*, 25.16–25.17
- nanosecond designs, **V**:25.31–25.32
- narrow linewidths, **V**:25.30
- ultrashort systems, **V**:25.33
- Matched filtering, **IV**:12.28–12.29, 12.29*f*, 12.30*f*
- Matched filters, for pattern recognition, **I**:11.12–11.13
- Matching:
 - in adjustment experiments, **III**:3.4, 3.5
 - color (*see* Color matching)
- Material designation, of crystals and glasses, **IV**:2.27–2.30, 2.29*f*
- Materials:
 - formation of, for optics, **II**:9.3–9.4
 - specifications for, **II**:4.9
 - tolerancing and properties of, **II**:5.9
- MathCad program, **V**:5.39
- Mathematica program, **V**:5.39
- Matrices:
 - amplitude scattering, **I**:7.10
 - for collineation, **I**:1.59–1.60
 - computing polarization with, **I**:12.27–12.30
 - density (coherency), **I**:12.29–12.30, 14.41
 - identity, **I**:14.8
 - Jones, **I**:7.10
 - in ellipsometry, **I**:15.30, 16.19
 - and Mueller matrices, **I**:14.3, 14.22–14.24, 14.27, 14.33
 - tensor product of, **I**:14.23
- Matrices (*Cont.*):
 - paraxial, and geometrical optics, **I**:1.65–1.74
 - Pauli spin, **I**:14.24, 14.41
 - point spread, **I**:15.36, 15.36*f*
 - power, **I**:1.67
 - for radiative transfer, **I**:9.11
 - for single scattering, **I**:9.16, 9.17
 - Stokes, **I**:14.4
 - T*-matrix method, **I**:7.15 (*see also* Mueller matrices)
- Matrix algebra, **III**:10.45–10.48
 - addition and multiplication, **III**:10.45–10.46
 - colorimetry, **III**:10.24–10.32
 - stimulus representation, **III**:10.24–10.27, 10.25*f*, 10.26*f*
 - transformations between color spaces, **III**:10.29–10.32
 - vector representation of data, **III**:10.25*f*, 10.27*f*, 10.27–10.29
 - glossary of notation in, **III**:10.5*t*
 - linear models, **III**:10.46–10.47
 - matrix transposition, **III**:10.46
 - simultaneous linear equations, **III**:10.47–10.48
 - singular value decomposition, **III**:10.48
 - special matrices and vectors, **III**:10.46
 - vectors and matrices, **III**:10.45
- Matrix generators, **I**:14.27
- Matrix theory for multilayer systems, **IV**:7.6–7.10, 7.9*f*
- Matsushita, **II**:19.22
- Matte, **II**:30.3
- Maximal atomic coherence, **IV**:14.28–14.32, 14.29*f*–14.32*f*
- Maximal coherence, **IV**:14.3
- Maximized D star, **II**:24.11
- Maximum frequency deviation, of electro-optic modulators, **V**:7.35
- Maximum likelihood sequence estimations (MLSEs), **V**:21.26
- Maximum refractive index, for PCFs, **V**:11.9–11.10
- Maximum saturation method (color matching), **III**:10.6*f*, 10.6–10.7
- Maximum spectral luminous efficiency (of radiation), **II**:37.2
- Maximum tolerable input jitter (MTIJ), **V**:15.16
- Maximum usable temperature, of metals, **IV**:4.7, 4.55*t*

- Maximum-likelihood expectation maximization (MLEM), **V**:32.2
- Maxwell equations, **I**:7.3
- Maxwell fisheye (lens), **I**:1.21
- Maxwell wave equation, **IV**:2.12
- Maxwell-averaged Gaunt factors, **V**:56.10
- Maxwell-Bloch equations, **IV**:11.6–11.7
- Maxwell-Boltzmann distribution, **IV**:20.11; **V**:14.4
- Maxwell-Boltzmann velocity distribution, **I**:31.24
- Maxwell-Garnett mixing formula, **I**:9.8
- Maxwell-Helmholtz-Drude dispersion formula, **IV**:2.12, 2.21–2.22
- Maxwellian ideal imaging, **I**:1.17, 1.28, 1.38
- Maxwellian view (viewing), **III**:5.4–5.8, 5.5*f*, 6.1–6.14
- advantages of, **III**:5.7
 - control of focus, **III**:5.4–5.5, 5.5*f*
 - defined, **III**:7.1
 - field of view, **III**:6.6*f*, 6.6–6.7
 - focus, **III**:6.5–6.6
 - interferometers, **III**:5.24
 - partial coherence, **III**:6.12–6.14
 - positioning of pupil in, **III**:5.9
 - pupil size, **III**:6.7–6.12
 - coherent illumination, **III**:6.9–6.12, 6.11*t*
 - incoherent target, **III**:6.7–6.9, 6.8*f*, 6.9*t*
 - retinal conjugate plane, **III**:5.5*f*
 - retinal illuminance, **III**:5.6–5.7, 6.3–6.5
 - size, **III**:5.5–5.6
 - spatial frequency content of stimuli, **III**:5.8
 - two-channel, **III**:5.21, 5.23*f*, 5.23–5.24
- Maxwell's electromagnetic theory, **I**:10.3
- Maxwell's equations:
- and binary gratings, **I**:23.13
 - and Bremsstrahlung radiation, **V**:56.8
 - and coherence theory, **I**:5.2, 5.3
 - and diffraction, **I**:3.1, 3.2, 3.4
 - for electric and magnetic fields, **IV**:1.16, 1.17, 5.52
 - for fiber-based couplers, **V**:16.2, 16.3
 - and laws of reflection and refraction, **I**:1.24
 - methods for solving, **IV**:9.2–9.3
 - for optical fields, **IV**:16.24
 - and optical waveguides, **I**:21.3, 21.5
 - and photonic crystal fibers, **V**:11.3, 11.6, 11.20
 - and refractive index, **IV**:2.6
 - for solids, **IV**:8.4–8.6
 - and surface scattering, **I**:8.4
 - and wave propagation, **V**:5.8
- Maxwell's method, for color matching, **III**:10.8, 10.8*f*
- Maxwell's principle, **II**:30.16
- Masers (microwave amplification by z-motion-induced emission of radiation), **II**:23.45
- McCarthy objective, **I**:29.32
- McCollough effect, **III**:11.76, 11.77*f*
- McKinley relays, **I**:18.18, 18.18*f*, 18.19
- Mean coefficient of linear thermal expansion, **IV**:4.7
- Mean-field approximation, **I**:7.16
- Mean-field theory, of liquid crystals, **V**:8.22–8.23
- Mean-sphere correction, **III**:1.6
- Mean-square-spot size (MSS), **II**:3.15
- Mean-time-to-failure (MTTF), **II**:25.13, 25.14
- Measured horopter, **III**:13.8
- Measurement(s), **II**:35.8–35.16
- of absorptance, **II**:35.10
 - on coatings, **IV**:7.12–7.14
 - of emittance, **II**:35.14–35.16, 35.15*f*
 - of lighting, **II**:40.51–40.54, 40.52*f*, 40.53*t*, 40.54*f*
 - of reflectance, **II**:35.10*f*–35.12*f*, 35.10–35.13, 35.14*t*
 - with surfaces and processes, **IV**:6.9–6.10
 - terminology, **II**:35.2*f*, 35.2–35.3
 - of transmittance, **II**:35.8–35.10, 35.9*f*
- (see also *specific types of measurement*, e.g.: Nonlinearity measurement)
- Measurement noise, **II**:22.7–22.8
- Measurement techniques, for semiconductors, **IV**:5.7, 5.56–5.83
- ellipsometry, **IV**:5.67–5.69, 5.68*f*, 5.69*f*
 - inelastic light scattering, **IV**:5.76*f*, 5.76–5.83, 5.78*f*–5.82*f*
 - instrumentation, **IV**:5.58–5.61, 5.59*f*, 5.60*f*, 5.72
 - luminescence, **IV**:5.69–5.75, 5.70*f*, 5.72*f*–5.75*f*
 - modulation spectroscopy, **IV**:5.64*f*–5.65*f*, 5.64–5.67, 5.66*t*, 5.68*f*
 - reflection and transmission/absorption, **IV**:5.62–5.64, 5.63*f*
 - spectroscopic procedures, **IV**:5.56–5.58
- Measurements Assurance Program (MAP), **II**:35.13
- Mechanical assembly, of polymers, **IV**:3.14*f*, 3.14–3.16, 3.15*f*

- Mechanical athermalization, **II**:8.8–8.12
 active, **II**:8.11, 8.11*f*
 by image processing, **II**:8.12
 part active, part passive, **II**:8.11–8.12, 8.12*f*
 passive, **II**:8.8*f*, 8.8–8.10, 8.9*f*, 8.10*f*
- Mechanical cycling, of metals, **IV**:4.10
- Mechanical distances, in focal Gaussian lenses, **I**:1.53
- Mechanical scribing, **II**:17.24
- Mechanical specifications:
 for lenses, **II**:4.9
 optical vs., **II**:4.2
- Mechanical switches, for networking, **V**:18.4, 18.5, 18.5*f*, 18.11, 18.11*f*, 18.12*f*
- Mechanical tolerances, **II**:5.2
- Mechanical vibrations, **II**:27.5, 27.6*f*
- Mechanically clamped mountings, **II**:6.12, 6.13*f*
- Mechanism direction (color vision), **III**:11.3, 11.33
- Mechanist approach (color vision) (*see* Color vision mechanisms)
- Media access control (MAC), of FDDI, **V**:23.3
- Media interface connectors (MICs), **V**:23.3
- Medial images, **I**:29.11, 29.37
- Medical imaging, **V**:31.1–31.10
 applications of, **V**:31.9, 31.10
 digital displays, **V**:31.8*f*–31.10*f*, 31.8–31.9
 digital tomosynthesis, **V**:31.7–31.8
 and inverse Compton x-ray sources, **V**:59.3–59.4
 and polycapillary x-ray optics, **V**:53.14–53.16, 53.15*f*–53.16*f*
 radiography, **V**:31.1–31.4, 31.2*f*–31.4*f*
 tomography, **V**:31.1, 31.5, 31.5*f*–31.7*f*, 31.5–31.7
 x-ray detectors for, **V**:61.2
- Medicine, nuclear, **V**:32.1–32.4
- Medium-bandpass filters, **IV**:7.78–7.83, 7.79*f*, 7.80*f*, 7.82*f*–7.88*f*
- Medium-format film, **I**:25.6
- Medium-wavelength infrared (MWIR) radiation, **II**:24.3, 25.2, 33.3, 33.5, 33.6*f*
- Meinel-Shack objective, **I**:29.28
- Melanin, **III**:14.9–14.10
- Mellin transforms, **I**:11.14; **V**:5.16
- Melt data sheets, for glass, **IV**:2.29
- Melted-resin arrays, **I**:22.42*f*–22.44*f*, 22.42–22.45
- Membrane mirrors (AO), **III**:15.1, 15.10*f*, 15.10–15.12
- Memory:
 holographic, **I**:33.24–33.25
 optical, **IV**:11.25, 12.34
- Memory colors, **II**:30.21
- Meniscus-shaped elements, **IV**:3.13
- Mercury arc lamps, **II**:15.29, 15.29*f*–15.31*f*, 15.34*f*
- Mercury cadmium telluride (HgCdTe) detectors, **II**:24.86–24.92, 24.87*f*
 infrared, **II**:33.7, 33.8*f*
 photoconductors, **II**:24.86*f*, 24.88–24.90, 24.89*f*–24.91*f*
 photodetectors, **II**:25.10, 25.10*t*
 photovoltaic, **II**:24.86*f*, 24.88*f*, 24.90–24.92, 24.91*f*, 24.92*f*
- Mercury cadmium telluride (Hg_{0.78}Cd_{0.22}Te) narrow-gap alloy, **IV**:5.73, 5.74*f*
- Mercury-doped germanium detectors, **II**:24.84*f*, 24.92–24.95, 24.93*f*–24.95*f*
- Mercury-free fluorescent lamps, **II**:40.31
- Mercury-halide fluorescent lamps, **II**:40.33*f*
- Mercury-vapor fluorescent lamps, **II**:40.30, 40.31, 40.33*f*
- Mercury-xenon lamps, **II**:15.34*f*
- Meridional rays, **I**:1.37
- Meridians (meridional planes), **I**:1.27, 1.32
- Meridional rays, **I**:1.35; **II**:3.3
- Merit function, in optical design software, **II**:3.17
- Mersenne objectives, **I**:29.9, 29.12
- Mersenne telescopes, **I**:18.20
- Mesa etching, **II**:18.3*f*, 18.3–18.4
- Mesa photodiodes, **II**:25.14, 25.15*f*
- Mesh topologies, of WDM networks, **V**:21.6
- Mesocotyl, **III**:8.2
- Mesospheric sodium lasers, **V**:5.32–5.34, 5.33*f*
- Mesopic vision, **II**:34.37, 36.9, 37.2
- Metal insulator semiconductor (MIS) photogate FPAs, **II**:33.10–33.11, 33.11*f*, 33.12*f*
- Metal insulator semiconductors (MISs), **II**:33.4
- Metal-dielectric multiple reflection cutoff filters, **IV**:7.59–7.60
- Metal-dielectric multiple reflection filters, **IV**:7.111
- Metal-dielectric reflectors, **IV**:7.81–7.82, 7.108–7.109, 7.109*f*, 7.110*f*
- Metal-insulator semiconductor (MIS) capacitors, **I**:26.2
- Metallic mirrors, mounting of, **II**:6.19–6.20, 6.20*f*

- Metallic reflecting coatings, **IV**:7.80–7.81
- Metallic reflectors, **IV**:7.106f–7.109f, 7.106–7.108
- Metallization, for SOAs, **V**:19.16, 19.16f
- Metalorganic chemical vapor deposition (MOCVD), **II**:19.6–19.7, 19.20t, 19.23
- Metal-organic chemical vapor deposition (MOCVD), **IV**:18.3
- Metal-organic molecular beam epitaxy (MOMBE), **I**:21.17, 21.18
- Metalorganic vapor phase epitaxy, **II**:17.21, 17.22
- Metal-organic vapor-phase epitaxy (MOVPE), **I**:21.17–21.18
- Metal-oxide semiconductors (MOSs):
 linear arrays of, **II**:32.21–32.24, 32.22f, 32.23f
 readouts from, **II**:32.20–32.21
- Metal-oxide-semiconductor (MOS) area array image sensors, **II**:32.25–32.26, 32.26f
- Metal-oxide-semiconductor (MOS) capacitors, **I**:26.2, 26.6; **II**:32.4f, 32.7–32.8
- Metal-oxide-semiconductor (MOS) detectors, **II**:25.11, 25.11f
- Metals, **IV**:4.1–4.70
- absorptance of, **IV**:4.39, 4.40f–4.42f, 4.48, 4.49
 - and emittance, **IV**:4.49, 4.49f, 4.50t, 4.51t
 - and mass attenuation coefficients for photons, **IV**:4.48t
 - aluminum and aluminum alloys
 - absorptance, **IV**:4.40f, 4.48t, 4.51t
 - optical properties, **IV**:4.12t, 4.20f, 4.21f
 - penetration depth, **IV**:4.47f
 - physical properties, **IV**:4.52t, 4.54t
 - reflectance, **IV**:4.27t–4.28t, 4.40f, 4.44f, 4.46f
 - thermal properties, **IV**:4.55t, 4.56t, 4.57f, 4.58t, 4.59f–4.60f, 4.65t, 4.66f, 4.69t, 4.70t
 - beryllium
 - absorptance, **IV**:4.48t, 4.50t
 - optical properties, **IV**:4.12t, 4.21f, 4.26f
 - penetration depth, **IV**:4.47f
 - physical properties, **IV**:4.52t, 4.54t
 - reflectance, **IV**:4.28t–4.29t, 4.45f, 4.46f
 - thermal properties, **IV**:4.55t, 4.56t, 4.57f, 4.58t, 4.59f–4.60f, 4.65t, 4.68f, 4.69t, 4.70t
 - chromium
 - absorptance, **IV**:4.48t, 4.50t
 - optical properties, **IV**:4.13t–4.14t, 4.22f
 - physical properties, **IV**:4.54t
 - reflectance, **IV**:4.30t–4.31t
 - thermal properties, **IV**:4.69t
- Metals (*Cont.*):
- copper
 - absorptance, **IV**:4.40f, 4.48t, 4.50t
 - optical properties, **IV**:4.12t–4.13t, 4.22f
 - physical properties, **IV**:4.52t–4.54t
 - reflectance, **IV**:4.29t–4.30t, 4.40f
 - thermal properties, **IV**:4.55t, 4.56t, 4.57f, 4.58t, 4.60f–4.61f, 4.65t, 4.66f, 4.69t, 4.70t
 - germanium
 - absorptance, **IV**:4.48t
 - thermal properties, **IV**:4.69t, 4.70t
 - gold
 - absorptance, **IV**:4.40f, 4.48t, 4.50t, 4.51t
 - optical properties, **IV**:4.14t, 4.23f
 - physical properties, **IV**:4.52t, 4.54t
 - reflectance, **IV**:4.31t–4.32t, 4.40f
 - thermal properties, **IV**:4.55t, 4.56t, 4.57f, 4.58t, 4.60f–4.61f, 4.65t, 4.66f, 4.69t, 4.70t
 - interband transitions in, **IV**:8.21
 - Invar 36, **IV**:4.10t, 4.52t, 4.55t, 4.69t, 4.70t
 - iron
 - absorptance, **IV**:4.40f, 4.48t, 4.50t
 - optical properties, **IV**:4.15t, 4.23f
 - physical properties, **IV**:4.54t
 - reflectance, **IV**:4.32t–4.33t, 4.40f
 - thermal properties, **IV**:4.55t, 4.56t, 4.57f, 4.58t, 4.62f–4.63f, 4.65t, 4.67f, 4.69t
 - mechanical properties of, **IV**:4.7–4.8
 - for mirror design, **IV**:4.8–4.10, 4.10t
 - molybdenum
 - absorptance, **IV**:4.41f, 4.48f, 4.48t, 4.50t, 4.51t
 - optical properties, **IV**:4.15t–4.16t, 4.24f
 - physical properties, **IV**:4.52t, 4.54t
 - reflectance, **IV**:4.33t–4.35t, 4.41f
 - thermal properties, **IV**:4.55t, 4.56t, 4.58f, 4.58t, 4.62f–4.63f, 4.65t, 4.67f, 4.69t, 4.70t
 - nickel and nickel alloys
 - absorptance, **IV**:4.41f, 4.48t, 4.50t, 4.51t
 - optical properties, **IV**:4.16t–4.17t, 4.24f
 - penetration depth, **IV**:4.47f
 - physical properties, **IV**:4.52t, 4.54t
 - reflectance, **IV**:4.35t–4.36t, 4.41f, 4.47f
 - thermal properties, **IV**:4.55t, 4.56t, 4.57f, 4.58f, 4.62f–4.63f, 4.65t, 4.67f, 4.69t, 4.70t
 - nomenclature for, **IV**:4.3

Metals (*Cont.*):

- optical properties of, **IV**:4.3–4.6, 4.4*f*, 4.11
 - dielectric function, **IV**:4.26*f*
 - extinction coefficient, **IV**:4.11, 4.12*t*–4.19*t*, 4.20*f*–4.26*f*
 - refraction index, **IV**:4.11, 4.12*t*–4.19*t*, 4.21*f*–4.26*f*
- penetration depth of, **IV**:4.47*f*
- physical properties of, **IV**:4.6, 4.49, 4.52*t*–4.54*t*
- platinum
 - absorptance, **IV**:4.41*f*, 4.48*t*, 4.50*t*, 4.51*t*
 - optical properties, **IV**:4.17*t*, 4.25*f*
 - physical properties, **IV**:4.54*t*
 - reflectance, **IV**:4.36*t*–4.37*t*, 4.41*f*
 - thermal properties, **IV**:4.69*t*, 4.70*t*
- reflectance of, **IV**:4.11, 4.27*t*–4.39*t*, 4.40*f*–4.47*f*
- semiconductors and dielectrics vs., **IV**:8.4
- silicon
 - absorptance, **IV**:4.48*t*
 - physical properties, **IV**:4.52*t*
 - reflectance, **IV**:4.46*f*
 - thermal properties, **IV**:4.55*t*, 4.56*t*, 4.58*f*, 4.58*t*, 4.63*f*–4.64*f*, 4.65*t*, 4.68*f*, 4.69*t*, 4.70*t*
- silicon carbide
 - absorptance, **IV**:4.49*f*, 4.50*f*
 - optical properties, **IV**:4.19*t*, 4.25*f*
 - physical properties, **IV**:4.52*t*
 - reflectance, **IV**:4.41*f*, 4.42*f*, 4.46*f*
 - thermal properties, **IV**:4.55*t*, 4.56*t*, 4.58*f*, 4.58*t*, 4.63*f*–4.64*f*, 4.65*t*, 4.68*f*, 4.69*t*, 4.70*t*
- silver
 - absorptance, **IV**:4.42*f*, 4.48*t*, 4.50*t*, 4.51*t*
 - optical properties, **IV**:4.17*t*–4.18*t*, 4.26*f*
 - physical properties, **IV**:4.52*t*, 4.54*t*
 - reflectance, **IV**:4.37*t*–4.38*t*, 4.42*f*
 - thermal properties, **IV**:4.55*t*, 4.56*t*, 4.57*f*, 4.58*t*, 4.60*f*–4.61*f*, 4.65*t*, 4.66*f*, 4.69*t*, 4.70*t*
- stainless steel
 - physical properties, **IV**:4.52*t*
 - thermal properties, **IV**:4.55*t*, 4.56*t*, 4.57*f*, 4.58*t*, 4.62*f*–4.63*f*, 4.65*t*, 4.67*f*, 4.69*t*, 4.70*t*
- steel, **IV**:4.50*t*, 4.51*t*
- tantalum, **IV**:4.50*t*, 4.69*t*, 4.70*t*
- thermal properties of, **IV**:4.6–4.7, 4.53, 4.55
 - coefficient of linear thermal expansion, **IV**:4.56*t*, 4.57*f*, 4.58*f*

Metals, thermal properties of (*Cont.*):

- elastic properties, **IV**:4.69, 4.69*t*
- at room temperature, **IV**:4.55*t*
- specific heat, **IV**:4.65*t*, 4.66*f*–4.69*f*
- strength and fracture properties, **IV**:4.70, 4.70*t*
- thermal conductivity, **IV**:4.58*t*, 4.59*f*–4.63*f*
- titanium, **IV**:4.48*t*, 4.50*t*, 4.52*t*, 4.55*t*
- tungsten
 - absorptance, **IV**:4.42*f*, 4.48*t*, 4.50*t*, 4.51*t*
 - optical properties, **IV**:4.18*t*–4.19*t*, 4.26*f*
 - physical properties, **IV**:4.54*t*
 - reflectance, **IV**:4.38*t*–4.39*t*, 4.42*f*
 - thermal properties, **IV**:4.69*t*, 4.70*t*
- zinc, **IV**:4.48*t*
- Metal-semiconductor-metal (MSM)
 - photodetectors, **II**:26.3, 26.4*f*; **V**:13.63, 13.73
- Metamerism, **III**:10.7, 10.36–10.38, 10.37*f*
- Metamers, **III**:10.1, 10.7
- Meteorological optics, **V**:3.40–3.43, 3.41*f*–3.43*f*
- Meter (unit), **II**:12.2
- Meter, 1875 Treaty of the, **II**:34.20, 36.2
- Meter candle (unit), **II**:34.43, 36.7
- Method of constant stimuli, **III**:3.9
- Metrology, optical, **II**:12.1–12.25
 - angle measurements, **II**:12.10–12.17
 - autocollimeters, **II**:12.11*f*, 12.11–12.12, 12.12*f*
 - interferometric methods, **II**:12.14
 - levels (tools), **II**:12.13*f*, 12.13–12.14, 12.14*f*
 - mechanical methods, **II**:12.10–12.11, 12.11*f*
 - in prisms, **II**:12.14–12.16, 12.15*f*–12.17*f*
 - theodolites, **II**:12.13
 - curvature measurements, **II**:12.17–12.25
 - mechanical methods, **II**:12.17–12.19, 12.18*f*, 12.19*f*, 12.19*t*
 - optical methods, **II**:12.19–12.21, 12.20*f*, 12.20*t*, 12.21*f*
 - cw OPOs for, **IV**:17.28, 17.29
 - of diamond-turned optics, **II**:10.12*f*, 10.12–10.13, 10.13*f*
 - focal length measurements, **II**:12.21–12.25, 12.22*f*–12.24*f*
 - length measurements, **II**:12.2–12.10
 - interferometers, **II**:12.5–12.10, 12.7*f*–12.9*f*
 - stadia and range finders, **II**:12.2–12.4, 12.3*f*, 12.4*f*
 - time-based and optical radar, **II**:12.4, 12.5, 12.6*f*

- Metrology, optical (*Cont.*):
 and magnetron-sputtered MLLs, **V**:42.6–42.7, 42.7*f*, 42.8*f*
 scatterometers in, **V**:1.16
 straightness measurements, **II**:12.10
 surface figure, **V**:46.3–46.6, 46.5*f*
 surface finish, **V**:46.2
 terminology, **II**:12.2
 x-ray mirror, **V**:46.1–46.12
 history of, **V**:46.1–46.2
 profile analysis considerations, **V**:46.6–46.12, 46.7*f*, 46.10*f*
 surface figure metrology, **V**:46.3–46.6, 46.5*f*
 surface finish metrology, **V**:46.2
- Metropolitan area networks (MANs), **V**:9.14, 21.7
- MeV proton acceleration, **IV**:21.54, 21.54*f*
- MH 2200 coating, **IV**:6.37, 6.38*f*
- MHz (MegaHertz), **III**:23.3
- Mica retardation plates, **I**:13.45–13.46
- Mica spacers, **IV**:7.84*f*, 7.88–7.89
- Michaelis-Menton function, **III**:2.14
- Michel Lévy Color Chart, **I**:28.35*f*
- Michelson contrast, in vision experiments, **III**:3.4
- Michelson formula, **III**:4.8
- Michelson interferometer, **I**:2.26*f*–2.27*f*, 2.26–2.28, 32.2, 32.3*f*, 32.21, 33.4; **II**:12.5, 12.6, 12.14; **III**:18.1; **IV**:5.60, 7.42, 7.104*f*; **V**:17.8*f*, 17.8–17.9
- Michelson stellar interferometers, **I**:2.40*f*, 2.40–2.41, 32.19, 32.19*f*
- Microbeam irradiation, **I**:28.54
- Microbolometer FPAs, **II**:33.13–33.14
- Microbunches, of electrons, **V**:58.1
- Microcalorimeter detectors, **V**:29.9–29.11, 29.11*f*, 60.9, 60.9*t*
- Microcavities, in 3D photonic crystals, **IV**:9.6–9.12, 9.7*f*, 9.8*f*
- Microchannel plate (MCP) detectors, **V**:63.34
- Microchannel plate tubes (MCPTs), **II**:24.32, 24.33*f*, 24.40
- Microchannel plates (MCPs), **II**:31.1, 31.9, 31.9*f*, 31.12*f*, 31.12–31.14, 31.13*f*, 31.13*t*; **V**:49.2, 49.3*f*, 49.5*f*, 50.7, 60.7
- Microchannel-plate image intensifiers (MCP IIs), **II**:31.7
 high-voltage power supply for, **II**:31.9, 31.10*f*
 proximity-focused, **II**:31.9, 31.9*f*, 31.16–31.18, 31.17*t*, 31.18*f*, 31.19*f*
- Microcreep strength, of metals, **IV**:4.8
- Microdensitometers, **II**:29.6, 29.15*f*, 29.15–29.16
- Microelectromechanical systems (MEMS), **I**:30.3; **III**:15.11
- Micro-electromechanical systems (MEMS)
 mirrors and switches, **V**:18.8, 18.8*f*, 18.11, 18.12*f*
- Microfocus x-ray fluorescence (MXRF):
 with doubly curved crystal diffraction, **V**:29.6–29.7, 29.8*f*–29.9*f*
 monocapillary, **V**:29.4
 polycapillary, **V**:29.4–29.6, 29.5*f*, 29.6*f*
 ultrahigh resolution, **V**:29.9*f*–29.11*f*, 29.9–29.11
- Microfocusing, with refractive x-ray lenses, **V**:37.7–37.8
- Micro-Fresnel lenses (MFLs), **I**:22.31*f*–22.33*f*, 22.31–22.37, 22.35*f*–22.37*f*
- Microinterferometers, **II**:10.13, 10.13*f*
- Microkeratomes, **III**:16.13, 16.14
- Microlens arrays, **II**:32.30, 32.31, 32.31*f*
- Microlens Nipkow disk confocal microscope, **III**:17.1–17.2
- Microlenses, **II**:12.24
 distributed-index planar, **I**:22.26–22.31, 22.27*f*–22.30*f*, 22.27*t*, 22.31*t*
 micro-Fresnel lenses (MFLs), **I**:22.31*f*–22.33*f*, 22.31–22.37, 22.35*f*–22.37*f*
 molded glass, **I**:22.9–22.10, 22.10*t*, 22.11*f*, 22.12*t*, 22.13*t*, 22.14*f*
 molded plastic, **I**:22.10, 22.12–22.15
- Micromachined membrane MEMS mirrors, **III**:15.11
- Micromachining techniques, for binary optics, **I**:23.16, 23.16*f*, 23.17
- Micromaser master equation, **II**:23.20–23.22
- Micromasers, **II**:23.26–23.27, 23.27*f*, 23.45
- Microminiature lamps, **II**:15.53
- Micromirrors, **I**:22.23, 30.61–30.62, 30.62*f*
- Micro-optical table (MOT) techniques, **I**:22.6, 22.7*f*
- Micro-optics (*see* Miniature and micro-optics)
- Micro-optics-based components, for networking, **V**:18.1–18.12
 attenuators, **V**:18.2, 18.9
 beam splitters, **V**:18.6, 18.6*f*
 circulators, **V**:18.3, 18.3*f*, 18.10
 directional couplers, **V**:18.2, 18.3, 18.3*f*, 18.9, 18.9*f*
 Faraday rotators, **V**:18.7, 18.7*f*
 filters, **V**:18.6
 gratings, **V**:18.5–18.6, 18.6*f*

- Micro-optics-based components, for networking
(*Cont.*):
GRIN-rod lenses, **V**:18.7, 18.8, 18.8*f*
isolators, **V**:18.3, 18.10, 18.10*f*
mechanical switches, **V**:18.4, 18.5, 18.5*f*,
18.11, 18.11*f*, 18.12*f*
MEMS mirrors and switches, **V**:18.8, 18.8*f*,
18.11, 18.12*f*
multiplexers/demultiplexers/duplexers,
V:18.4, 18.4*f*, 18.10–18.11
network functions, **V**:18.2–18.5
polarizers, **V**:18.7, 18.7*f*
power splitters, **V**:18.2*f*, 18.2–18.3, 18.9, 18.9*f*
prisms, **V**:18.5, 18.5*f*
- Microplasma formation, **II**:17.28
- Micro-pore optics, **V**:49.1–49.6
- Microsaccades, **III**:1.44
- Microscopes, **I**:28.1–28.56
aperture-scanning microscopy, **I**:28.53–28.54
bright field microscopy, **I**:28.25, 28.27*f*,
28.27–28.28
coherent diffraction, **V**:27.4*f*, 27.4–27.5, 27.5*f*
compound, **I**:17.10
confocal microscopy, **I**:28.49*f*, 28.49–28.51,
28.51*f*
contrast in
bright field microscopy, **I**:28.25, 28.27*f*,
28.27–28.28
dark field microscopy, **I**:28.28
Hoffman modulation contrast, **I**:28.29
interference microscopy, **I**:28.33–28.44,
28.35*f*, 28.37*f*, 28.38*f*, 28.40*f*, 28.42*f*,
28.43*f*
modulation transfer function,
I:28.24–28.25, 28.25*f*, 28.26*f*
phase contrast, **I**:28.28–28.29, 28.29*f*
SSEE microscopy, **I**:28.29–28.33,
28.30*f*–28.33*f*
dark field microscopy, **I**:28.28
differential-interference contrast,
I:28.39–28.41, 28.40*f*
Dyson, **I**:28.41, 28.42, 28.42*f*
first-order layout for, **II**:1.8
fluorescent microscopy, **I**:28.48–28.49
history of, **I**:28.1–28.3
holographic, **I**:28.42, 28.43, 28.43*f*
imaging modes of, **I**:28.44*f*, 28.44–28.54,
28.46*f*, 28.47*f*, 28.49*f*, 28.51*f*
interference microscopy, **I**:28.35*f*, 28.37*f*,
28.38*f*, 28.42*f*, 28.43*f*
- Microscopes (*Cont.*):
Jamin-Lebedev, **I**:28.38*f*, 28.38–28.39
lenses in, **I**:28.9–28.17, 28.10*t*, 28.11*t*,
28.12*f*–28.16*f*
light field microscopy, **I**:28.53
Linnik, **I**:28.36–28.38, 28.37*f*
Mach-Zehnder, **I**:28.36, 28.37*f*
Mirau, **I**:28.41, 28.42, 28.42*f*
Nomarski, **II**:10.11, 10.11*f*
Nomarski microscope, **V**:46.2
optical arrangements in, **I**:28.3–28.9,
28.4*f*–28.8*f*
optical path difference (OPD) in,
I:28.33–28.34, 28.35*f*, 28.36
resolution, **I**:28.17–28.24, 28.18*f*–28.21*f*
scanning electron
and magnetron-sputtered MLLs,
V:42.6–42.7, 42.7*f*, 42.8*f*, 42.13, 42.13*f*
and x-ray spectral detection, **V**:62.1–62.3,
62.2*f*
specimen manipulation for, **I**:28.54–28.55
SSEE microscopy, **I**:28.29–28.33, 28.30*f*–28.33*f*
traveling, **II**:12.20, 12.21, 12.21*f*
x-ray, **V**:37.6
- Microsource devices, **V**:28.7
- Microstrain, of metals, **IV**:4.8
- Microstrip detectors, **V**:63.32–63.33
- Microstructured optical arrays (MOAs),
adaptive, **V**:50.7–50.8, 50.8*f*
- Microwave powered lamps (*see* Electrodeless
sulfur lamps)
- Microyield strength, of metals, **IV**:4.8, 4.70, 4.70*t*
- Midget ganglion cells, **III**:2.10–2.11, 2.10*n*
- Midsagittal plane, **III**:13.2
- Mid-wave infrared (MWIR) AOTFs, **V**:6.42
- Mie scattering, **I**:7.11, 7.12, 9.17; **IV**:1.15, 1.32,
2.27; **V**:3.12, 3.16–3.18, 3.17*f*–3.19*f*
- Mie theory, **IV**:1.40; **V**:11.7
- Miesowicz viscosity coefficients, **V**:8.24
- Military Sensing Information Analysis Center
(SENSIAC), **II**:15.6
- Miller algorithm, **I**:7.15
- Miller's rule, **IV**:10.9
- Millilambert (unit), **II**:36.7
- Milliphot (unit), **II**:36.7, 36.7*t*
- Miniature and micro-optics, **I**:22.1–22.46, 22.46*f*
and binary optics, **I**:23.7*f*–23.8*f*, 23.7–23.8
design considerations, **I**:22.2–22.8
diamond turning, **I**:22.15–22.18, 22.16*t*,
22.17*f*, 22.18*f*

- Miniature and micro-optics (*Cont.*):
- distributed-index planar microlenses,
 - I**:22.26–22.31, 22.27*f*–22.30*f*, 22.27*t*, 22.31*t*
 - drawn preform cylindrical lenses, **I**:22.45, 22.46*f*
 - high-performance miniature systems,
 - I**:22.5–22.8
 - laser-assisted chemical etching (LACE),
 - I**:22.45
 - liquid lenses, **I**:22.37–22.41, 22.38*f*–22.41*f*, 22.42*t*
 - and lithography, **I**:22.18–22.25, 22.20*f*–22.25*f*
 - mass-transport process, **I**:22.45, 22.46*f*
 - melted-resin arrays, **I**:22.42*f*–22.44*f*, 22.42–22.45
 - micro-Fresnel lenses, **I**:22.31*f*–22.33*f*, 22.31–22.37, 22.35*f*–22.37*f*
 - molded microlenses, **I**:22.8–22.15
 - molded glass, **I**:22.9–22.10, 22.10*t*, 22.11*f*, 22.12*t*, 22.13*t*, 22.14*f*
 - molded plastic, **I**:22.10, 22.12–22.15
 - monolithic lenslet modules, **I**:22.25*f*, 22.25–22.26
- Miniature eye movements, **III**:1.44
- Miniature lamps, **II**:15.53, 40.26*t*, 40.28*f*
- Minimally distinct border (MDB), **III**:11.37
- Minimum motion (MM), **III**:11.37
- Minimum resolvable temperature (MRT) (of infrared detector arrays), **II**:33.2, 33.27*f*, 33.27–33.28
- Minimum signal, for solid-state cameras, **I**:26.13–26.14
- Minkowitz distance-measuring interferometers, **II**:12.7, 12.8, 12.8*f*
- Minority-carrier recombination, **II**:17.2
- Minox camera, **I**:25.21
- Minus filters, **IV**:7.43, 7.48, 7.49, 7.50*f*
- Mirages, **I**:24.1, 24.2*f*; **V**:3.42*f*, 3.42–3.43, 3.43*f*
- Mirau interference microscopes, **I**:28.41, 28.42, 28.42*f*
- Mircolens arrays, for agile beam steering, **I**:30.57–30.60
- Mirror reflectivity, for VCSELs, **V**:13.44
- Mirror scatter relationships (stray light), **II**:7.18
- Mirror surface roughness:
 - and grazing-incidence neutron optics, **V**:64.2–64.3
 - and Wolter x-ray optics, **V**:47.3–47.5, 47.4*f*
- Mirrored tiling, **II**:39.28, 39.30, 39.30*f*
- Mirror-image effect, **I**:12.6
- Mirrors:
 - aluminum, **IV**:7.106*f*–7.108*f*, 7.106–7.108
 - on amplifiers, **II**:16.3
 - bimorph, **V**:50.4, 50.4*f*, 50.6*f*
 - Bragg, **V**:13.28, 13.44
 - compound, **I**:30.15*f*, 30.15–30.16
 - and concentrators, **II**:39.17
 - conic, **I**:29.3*f*, 29.4*f*
 - deformable, **V**:5.4, 5.4*f*, 5.37*f*, 5.37–5.38, 5.38*f*
 - double-bounce Wolter, **V**:52.4
 - external, **V**:13.32–13.33, 13.33*f*
 - Fresnel's, **I**:2.16, 2.16*f*
 - Goebel, **V**:26.10
 - graded reflectivity, **IV**:7.52
 - high-reflectivity, **V**:41.7–41.8, 41.8*f*
 - hot and cold, **I**:7.58
 - Kirkpatrick-Baez, **V**:44.4, 44.4*f*, 63.21, 64.5*f*, 64.5–64.6
 - Lippmann-Bragg holographic, **IV**:7.50
 - Littrow, **I**:20.4
 - Lloyd's, **I**:2.16, 2.17*f*, 2.18
 - metals for, **IV**:4.8–4.10, 4.10*t*
 - in micro-electromechanical systems, **V**:18.8, 18.8*f*, 18.12*f*
 - micromirrors, **I**:22.23, 30.61–30.62, 30.62*f*
 - mounting of [*see* Mounting (of optical components)]
 - for neutron optics, **V**:63.20–63.21, 63.21*f*
 - nonabsorbing, **II**:19.23–19.24
 - nonlinear optical loop, **V**:20.22
 - perfect, **IV**:7.47
 - phase conjugate, **IV**:12.7, 12.8*f*–12.9*f*, 12.33–12.35
 - plane, **I**:1.25
 - point-spread function of, **V**:64.2
 - polarizing, **V**:63.28
 - in reflecting afocal lenses, **I**:18.20, 18.20*f*, 18.21, 18.21*f*
 - reflection from, **I**:1.25
 - for scanners, **I**:30.14–30.16, 30.15*f*, 30.60–30.62, 30.62*f*
 - self-pumped phase conjugate, **IV**:12.7, 12.8*f*–12.9*f*
 - semiconductor laser amplifier loop optical, **V**:20.22
 - semiconductor saturable absorber, **IV**:18.3, 18.10–18.11; **V**:25.3, 25.32
 - semilinear, **IV**:12.7, 12.8*f*
 - and SHADOW code, **V**:35.4, 35.5

- Mirrors (*Cont.*):
 silver, **IV**:7.106*f*–7.108*f*, 7.107–7.109
 super mirrors, **IV**:7.111
 as thin lenses, **I**:1.55
 threshold conditions with, **II**:16.10–16.12, 16.11*f*
 Wolter configurations, **V**:52.4, 64.6, 64.6*f*
 (*see also* X-ray mirrors)
- Misfocus, **I**:1.82
- Mixed characteristic functions, **I**:1.13
- Mixing rods (*see* Lightpipes)
- M&M states, **IV**:23.14
- Mobius strip, **I**:14.15
- Modal approach, to wavefront error correction, **V**:4.35
- Modal dispersion, **II**:17.34; **V**:15.9
- Modal filtering, **V**:11.12–11.13, 11.13*f*
- Modal gain per unit length, **V**:13.7, 13.8
- Modal noise, **V**:15.16–15.17
- Mode:
 defined, **III**:8.2
 and energy density distribution, **III**:8.15
 in hair cells, **III**:8.25–8.26
 in human photoreceptors, **III**:8.16*f*, 8.18*f*
 in monkey/human retinal receptors, **III**:8.19–8.24, 8.22*f*
- Mode field adaptors (MFAs), **V**:25.2, 25.17–25.18
- Mode field diameter (MFD), **V**:25.2
- Mode hopping, **V**:13.13
- Mode matching, in mode field adaptors, **V**:25.17, 25.17*f*
- Mode partition noise, **V**:20.1
 for fiber optic communication links, **V**:15.11*f*, 15.11–15.13
 for laser diodes, **V**:13.19–13.20, 13.20*f*
- Mode transformers, photonic crystal fibers and, **V**:11.26–11.27, 11.27*f*
- Model eyes, **III**:21.13
- Model-dependent characterization methods (MDCM) (color CRTs), **III**:22.27–22.33
 conditions for use, **III**:22.30
 gun independence, **III**:22.27–22.28
 inverse transformations, **III**:22.32
 measurement of parameters, **III**:22.31
 normalization coefficients, **III**:22.31–22.32
 out-of-gamut colors, **III**:22.32–22.33
 partial models, **III**:22.30
 phosphor constancy, **III**:22.28–22.29, 22.32
 phosphor output models, **III**:22.29–22.30
- Mode-locked fiber lasers, **V**:25.32–25.33
- Mode-locked lasers, **II**:16.27–16.29, 16.28*f*, 20.7; **V**:20.15–20.17, 20.16*f*
- Modelocking:
 additive pulse, **IV**:18.3, 18.14
 colliding pulse, **IV**:18.3
 cw, **IV**:18.5*f*
 cw Q-switched, **IV**:18.5*f*
 Kerr lens, **IV**:18.3, 18.14–18.15
 passive, **IV**:18.8*f*, 18.8–18.9, 18.12–18.15
 Q-switched, **IV**:18.4, 18.5, 18.5*f*
 soliton, **IV**:18.8*f*, 18.12–18.14
- Moderators, for neutron optics, **V**:63.12, 63.14–63.15
- Mode-stabilized lasers, **II**:19.18–19.19, 19.19*f*
- Modified chemical vapor deposition (MCVD), **V**:25.2, 25.21, 25.26, 25.28
- Modified lambda coupling, **IV**:14.24*f*
- MODTRAN program, **V**:3.24
- Modulated grating (MG) reflectors, **V**:13.36
- Modulation:
 amplitude, **V**:7.22–7.24, 7.23*f*, 7.24*f*
 cross-gain, **V**:19.12, 19.13*f*, 19.27, 19.29*f*, 19.29–19.30, 19.32, 19.35–19.36
 cross-phase
 in optical fibers, **V**:10.3–10.4
 and SOAs, **V**:19.13, 19.30–19.32, 19.31*f*, 19.33*f*, 19.35–19.36
 and solitons, **V**:22.5, 22.13–22.15
 in WDM networks, **V**:21.19
 degree of, for electro-optic modulators, **V**:7.35–7.36
 depth-of-amplitude, **V**:7.22
 depth-of-phase, **V**:7.19
 digital, **V**:6.33
 direct, **V**:20.17*f*, 20.17–20.18
 frequency, **V**:7.24–7.25, 7.25*f*
 as laser stabilization technique, **II**:22.13–22.14
 longitudinal spatial, **V**:6.12
 in OTDM communication networks
 direct and indirect, **V**:20.17*f*, 20.17–20.18
 external, **V**:20.18–20.20, 20.19*f*, 20.20*f*
 percent, **V**:7.35
 phase, **V**:7.18–7.20
 bulk electro-optic modulators, **V**:7.18–7.20
 by lithium niobate modulators, **V**:13.51
 polarization, **V**:7.20*f*, 7.20–7.22, 7.21*f*
 pulse code, **V**:20.8

- Modulation (*Cont.*):
- self-phase
 - in optical fibers, **V**:10.3–10.4
 - and solitons, **V**:22.3*f*, 22.3–22.4, 22.4*f*
 - in WDM networks, **V**:21.18–21.19, 21.19*f*
 - spatial light, **V**:6.4, 6.9
 - transverse spatial, **V**:6.11–6.12, 6.30, 6.31, 6.32
 - in WDM networks, **V**:21.27–21.36
 - basic concepts, **V**:21.27–21.29, 21.28*f*–21.30*f*
 - carrier-suppressed return-to-zero and duobinary, **V**:21.30–21.33, 21.31*f*, 21.32*f*
 - comparisons of, **V**:21.36, 21.36*t*, 21.37*t*
 - DPSK and DQSK, **V**:21.33*f*–21.35*f*, 21.33–21.36, 21.36*t*, 21.37*t*
- Modulation bandwidth, electro-optic modulators, **V**:7.34
- Modulation (contrast) color spaces, **III**:10.19
- Modulation efficiency, of electro-optic modulators, **V**:7.36
- Modulation error ratio (MER), **V**:15.4–15.5
- Modulation formats, for WDM networks, **V**:21.27–21.36
 - basic concepts, **V**:21.27–21.29, 21.28*f*–21.30*f*
 - carrier-suppressed return-to-zero and duobinary, **V**:21.30–21.33, 21.31*f*, 21.32*f*
 - comparisons of, **V**:21.36, 21.36*t*, 21.37*t*
 - DPSK and DQSK, **V**:21.33*f*–21.35*f*, 21.33–21.36, 21.36*t*, 21.37*t*
- Modulation instability, **V**:10.3
- Modulation noise, **II**:24.11
- Modulation response, of laser diodes, **V**:13.16*f*, 13.16–13.17
- Modulation spectroscopy, **IV**:5.64*f*–5.65*f*, 5.64–5.67, 5.66*t*, 5.68*f*
- Modulation thresholds, **III**:1.22
- Modulation transfer:
- with defocus, **III**:1.29
 - TCA effect on, **III**:1.20
- Modulation transfer function (MTF), **I**:4.3; **II**:4.1, 4.2, 4.4, 4.5*f*, 4.8*f*, 29.18, 29.18*t*, 29.19*f*, 30.5, 31.26–31.27, 33.2, 33.7
 - aberration-derived, **III**:1.22
 - in aging eyes, **III**:14.13
 - calculations, **I**:4.3–4.6, 4.4*f*, 4.5*f*
 - and camera lens performance, **I**:27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.24, 27.25
 - and characteristics of objective detectors, **I**:28.16
- Modulation transfer function (MTF) (*Cont.*):
- and contrast of microscopes, **I**:28.24–28.25, 28.25*f*, 28.26*f*
 - and development in xerographic systems, **I**:34.7
 - diffraction-limited, **I**:4.4*f*, 4.4–4.5, 4.5*f*
 - in diffraction-limited eye, **III**:1.13, 1.13*f*
 - double-pass, **III**:1.23
 - measurements of, **I**:4.6–4.8, 4.8*f*
 - for microscopes, **I**:28.24, 28.25*f*
 - and observed optical performance, **III**:1.24–1.25
 - and off-axis image quality, **III**:1.26–1.27
 - in ophthalmoscopic methods, **III**:1.23
 - and optical quality of IOLs, **III**:21.14
 - for scanners, **I**:4.6
 - and scattered light, **III**:1.21
 - for solid-state cameras, **I**:26.14
 - at specific wavelengths, **I**:17.38, 17.39
 - of uniformly illuminated apertures, **I**:30.9, 30.10, 30.10*f*
 - of visual instruments, **III**:1.28
 - in young adult eyes, **III**:1.24*f*
- Modulation transfer functions (MTFs):
- for acousto-optic modulators, **V**:6.32–6.33
 - for polycapillary x-ray optics, **V**:53.15, 53.16*f*
 - for Schwarzschild objectives, **V**:51.1–51.2, 51.2*f*
 - and SPECT imaging, **V**:32.3
 - for x-ray detectors, **V**:61.3
- Modulators:
- acousto-optic, **V**:6.23*t*, 6.31–6.35, 6.32*f*, 6.34*t*
 - acousto-optic frequency shifters, **V**:6.35
 - and Bragg diffraction, **V**:6.4, 6.6, 6.7, 6.14
 - image (scophony), **V**:6.34–6.35
 - multi-mode interference, **V**:13.2, 13.35
 - principle of operation, **V**:6.32, 6.32*f*
 - band-filling, **V**:13.62
 - electroabsorption
 - in fiber optic systems, **V**:13.55–13.60, 13.56*f*
 - in OTDM networks, **V**:20.18, 20.20, 20.20*f*
 - electro-optic, **II**:22.14, 22.20; **V**:7.1–7.39, 13.61, 20.18–20.19, 20.19*f*
 - applications for, **V**:7.36–7.39
 - bulk modulators, **V**:7.16–7.28, 7.21*f*, 7.24*f*–7.28*f*
 - crystal optics and the index ellipsoid, **V**:7.3–7.7, 7.4*f*–7.6*f*, 7.8*f*–7.10*f*

- Modulators, electro-optic (*Cont.*):
 and electro-optic effect, **V**:7.6–7.16,
 7.8*t*–7.10*t*, 7.14*f*, 7.16*f*
 electro-optic sampling, **V**:7.36–7.37,
 7.37*f*–7.38*f*
 and Euler angles, **V**:7.39
 in fiber optic systems, **V**:13.61
 geometries, **V**:7.16–7.18, 7.17*f*
 laser mode locking, **V**:7.38–7.39
 light propagation in, **V**:7.3
 materials, **V**:7.33–7.34
 in OTDM networks, **V**:20.18–20.19
 performance criteria, **V**:7.34–7.36
 sensors, **V**:7.38
 traveling wave modulators, **V**:7.28–7.30,
 7.29*f*
 waveguide or integrated-optic
 modulators, **V**:7.30–7.32, 7.31*f*–7.33*f*
 electro-optical modulators, **I**:15.23
 electrorefractive, **V**:13.61–13.62
 integrated-optic, **V**:7.3, 7.30–7.32, 7.31*f*–7.33*f*
 intensity, **III**:5.14
 interferometric Mach-Zehnder,
V:13.51–13.52, 13.54–13.55, 13.63
 lithium niobate, **V**:13.2, 13.48–13.55, 13.49*f*
 Mach-Zehnder, **I**:21.26–21.28, 21.27*f*,
 21.32, 21.34
 interferometric, **V**:13.51–13.52, 13.54*f*,
 13.54–13.55, 13.63
 in WDM networks, **V**:21.30, 21.31*f*, 21.32*f*
 magneto-optical modulators, **I**:15.23
 Nipi, **V**:13.62
 photo-elastic, **I**:15.21, 16.13
 polarization (retardance), **I**:15.20–15.24
 semiconductor interferometric, **V**:13.63
 separate confinement heterostructure for,
V:13.4*f*, 13.5
 traveling wave, **I**:21.26
 waveguide, **V**:7.30–7.32, 7.31*f*–7.33*f*, 13.56,
 13.57*f*
- MOFSET bucket brigades, **II**:33.17
 Mohs scale, **IV**:2.31, 2.32*f*
 Moiré deflectometry, **II**:12.23, 12.24*f*
 Moiré grating, **IV**:22.11, 22.12*f*
 Moiré tests, of spherical aberrations,
II:13.26–13.27
 Molasses, optical [*see* Optical molasses (OM)]
 Molding, of polymers, **IV**:3.2, 3.12–3.13
 Molecular absorption, **V**:3.12–3.15, 3.13*f*
 Molecular absorption line database, **V**:3.22*f*,
 3.22–3.23, 3.23*f*
 Molecular alignment, in strong fields,
IV:21.22–21.23, 21.23*t*, 21.24*f*
 Molecular beam epitaxy (MBE), **I**:21.17,
 21.18; **II**:17.21–17.23, 19.6, 19.7, 25.16;
IV:5.7, 18.3
 Molecular dissociation, **IV**:21.23
 Molecular emission, **V**:3.18, 3.20, 3.20*f*
 Molecular Expressions (website), **I**:28.3
 Molecular gases, in standard atmosphere,
V:3.6, 3.7*t*, 3.8*f*–3.9*f*
 Molecular imaging, **V**:32.1
 Molecular orientational Kerr effect, **IV**:16.3*t*
 Molecular scattering, **I**:7.11
 Molecular spectroscopy, **V**:2.5–2.6, 2.6*f*, 2.7*f*
 Molecular targets, for SPECT imaging,
V:32.1
 Molecular tunnel ionization, **IV**:21.25*f*,
 21.25–21.26, 21.27*f*
 Molecular weight, of crystals and glasses,
IV:2.30
 Molecules, strong field interactions with,
IV:21.22–21.26
 Coulomb explosion, **IV**:21.24–21.25
 nuclear motion and alignment in,
IV:21.22–21.23, 21.23*t*, 21.24*f*
 triatomic and larger, **IV**:21.26
 tunnel ionization and ionization distance,
IV:21.25*f*, 21.25–21.26, 21.27*f*
- Molybdenum:
 absorptance of, **IV**:4.41*f*, 4.48*f*, 4.48*t*, 4.50*t*,
 4.51*t*
 optical properties of, **IV**:4.15*t*–4.16*t*, 4.24*f*
 physical properties of, **IV**:4.52*t*, 4.54*t*
 reflectance of, **IV**:4.33*t*–4.35*t*, 4.41*f*
 thermal properties of
 coefficient of linear thermal expansion,
IV:4.56*t*, 4.58*f*
 elastic stiffness, **IV**:4.69*t*
 moduli and Poisson's ratio, **IV**:4.69*t*
 at room temperature, **IV**:4.55*t*
 specific heat, **IV**:4.65*t*, 4.67*f*
 strength and fracture properties, **IV**:4.70*t*
 thermal conductivity, **IV**:4.58*t*, 4.62*f*–4.63*f*
- Moments normalization (technique),
II:36.15–36.16, 36.15*t*, 36.16*f*
 Momentum, family, **IV**:20.38
 Monin-Obukhov similarity theory,
V:3.31

- Monitors:
 and Computer Vision Syndrome, **III**:23.6–23.8
 interlaced, **III**:23.2
 noninterlaced, **III**:23.3
 in optical systems, **III**:5.16
- Monkey photoreceptors, **III**:8.19–8.24, 8.22*f*
- Monocapillary x-ray optics, **V**:28.5, 52.1–52.6, 52.2*f*–52.5*f*, 52.2*t*
- Monocentric Schmidt-Cassegrain objectives, **I**:29.16
- Monocentric systems, **I**:29.37
- Monochromatic aberration correction, **I**:23.6–23.7
- Monochromatic imaging, in polycapillary x-ray optics, **V**:53.16–53.17, 53.17*f*
- Monochromatic ocular aberrations, **III**:1.4, 1.14–1.19
 age-related, **III**:14.12–14.14, 14.13*f*
 correction of, **III**:1.25, 1.26, 1.26*f*
 off-axis, **III**:1.18*f*, 1.18–1.19
 on the visual axis, **III**:1.15–1.18
- Monochromatic sources of light, **I**:5.11
- Monochromatic vision, **III**:10.16
- Monochromators, **II**:35.9, 38.7–38.8, 38.14*f*–38.16*f*, 38.14–38.16; **IV**:5.59–5.61
 Bragg reflection, **V**:39.4, 39.5
 Bragg reflections in, **V**:39.1, 63.24
 crystal
 and bent crystals, **V**:39.1–39.6, 39.2*t*, 39.3*f*, 39.5*f*–39.6*f*
 in neutron optics, **V**:63.23–63.25
- Czerny-Turner, **I**:31.6
- double-pass, **I**:20.9*f*
- Dragon* systems of, **V**:38.3
- Grasshopper*, **V**:38.3
- Johansson bent/ground focusing, **V**:39.5
- Perkin-Elmer Model 99, **I**:20.9*f*
- and SHADOW code, **V**:35.4, 35.4*f*
- spherical-grating, **V**:38.3
- synchrotron radiation, **V**:39.6
- toroidal-grating, **V**:38.3
- Unicam prism-grating double-, **I**:20.10, 20.13, 20.15*f*
- in VUV and soft x-ray region, **V**:38.1–38.8
 diffraction properties, **V**:38.1–38.3, 38.2*f*
 dispersion properties, **V**:38.6–38.7
 efficiency of, **V**:38.8
 focusing properties, **V**:38.3*f*, 38.3–38.6, 38.4*t*–38.5*t*
 resolution properties, **V**:38.7
- Monochromators (*Cont.*):
 x-ray, **V**:30.1–30.4, 50.6–50.7
 and x-ray diffraction, **V**:28.3, 28.4
- Monochrome CRTs, **III**:22.3*f*, 22.3–22.4
 controls for, **III**:22.7, 22.8*f*
 design and operation of, **III**:22.3*f*, 22.3–22.4
 standards for, **III**:22.14
- Monochrome LCDs, operational principles of, **III**:22.34–22.37
- Monoclinic crystals, **IV**:2.7*t*, 2.18, 2.47*t*, 8.9*t*, 8.19*t*
- Monocular cues, for perceived space, **III**:13.3–13.7
- Monocular field:
 horizontal angular extent of, **III**:1.38*f*
 and stereopsis, **III**:1.38–1.42
- Monocular magnification, distortion by, **III**:13.13–13.16
 bifocal jump, **III**:13.15*f*, 13.15–13.16
 from convergence responses to prism, **III**:13.19
 discrepant views of objects/images, **III**:13.16
 motion parallax, **III**:13.14–13.15
 perspective distortion, **III**:13.13, 13.14*f*
 stereopsis, **III**:13.13, 13.14
- Monogon scanners, **I**:30.34–30.36
- Monolithic built-up mirror substrate, **II**:6.18, 6.18*f*
- Monolithic fiber laser resonators, **V**:25.16
- Monolithic FPAs, **II**:33.10–33.14, 33.11*f*, 33.12*f*
- Monolithic LED displays, **II**:17.30
- Monolithic lenslet modules (MLMs), **I**:22.25*f*, 22.25–22.26
- Monolithic silicon bolometers, **II**:28.10*f*, 28.10–28.11
- Monolithic tunable lasers, **V**:13.33–13.36, 13.34*f*–13.36*f*
- Monolithic two-dimensional (2D) laser arrays, **II**:19.39
- Monovision:
 with contact lenses, **III**:14.28
 defined, **III**:12.13
 with head mounted visual displays, **III**:13.32
- Monte Carlo simulations, **II**:39.7
- Mordants, **II**:30.4
- Morel model of absorption, **IV**:1.24, 1.28
- Morse interatomic potential, **IV**:2.16
- Mosaic crystals, **V**:39.2
- Mosaic II SSA cameras, **II**:31.29
- Motion artifacts (OCT), **III**:18.15

- Motion detection/discrimination, **III**:2.36–2.40
 optic flow fields, **III**:2.39f
 thresholds, **III**:2.37f
- Motion parallax, **III**:13.4, 13.5
 defined, **III**:13.2
 with monocular magnification,
III:13.14–13.15
- Motion perception, **III**:13.6–13.7
- Mott-Wannier excitons (*see* Wannier excitons)
- Mt. Pinatubo, **V**:3.10, 3.18, 3.39
- Mt. Wilson telescope, **V**:5.27
- Mounting (of optical components), **II**:6.1–6.24,
 6.17f, 6.17–6.18, 6.18f
 and contact stresses, **II**:6.21
 of domes, **II**:6.11, 6.12f
 hard, **II**:6.2–6.4
 of individual rotationally symmetric optics,
II:6.2–6.5
 lever-mechanism, **II**:6.19, 6.19f
 of moderate-sized mirrors, **II**:6.17–6.20,
 6.20f
 in multicomponent lens assemblies,
II:6.5–6.11
 drop-in assembly, **II**:6.6, 6.6f
 lathe assembly, **II**:6.7–6.8, 6.8f
 lens adjustments at assembly, **II**:6.8–6.11,
 6.10f
 “poker chip” assembly, **II**:6.8, 6.9f
 tightly-toleranced assembly, **II**:6.7, 6.7f
 of small mirrors/prisms, **II**:6.11–6.17
 bonded mountings, **II**:6.13–6.15, 6.15f
 elastomeric mountings for mirrors,
II:6.12
 flexure mountings for small mirrors/
 prisms, **II**:6.15–6.17, 6.16f
 mechanically clamped mountings, **II**:6.12,
 6.13f
 spring-loaded mountings, **II**:6.13, 6.14f
 soft, **II**:6.4–6.5
 and temperature effects, **II**:6.21–6.24, 6.22f
 of windows, **II**:6.11, 6.11f
- Moving electron beam exposure systems
 (MEBES), **I**:23.15, 23.16, 23.16f
- Moving entry and exit pupils method
 (photoreceptor directionality), **III**:8.6
- Mueller calculus, **I**:12.28–12.30
- Mueller matrices, **I**:7.10, 14.1–14.42; **V**:1.14,
 19.18–19.19
 about, **I**:14.3–14.4
 coordinate system, **I**:14.19–14.20
- Mueller matrices (*Cont.*):
 and depolarization, **I**:14.30–14.31
 depolarization index, **I**:14.32
 generators of, **I**:14.33–14.39, 14.36f, 14.37f
 nondepolarizing matrices, **I**:14.24–14.25,
 14.27–14.30
 diattenuators/diattenuation, **I**:14.16f,
 14.16–14.19
 in ellipsometry, **I**:15.30, 16.19–16.21, 16.20f,
 16.20t, 16.21f
 and Jones matrices, **I**:14.22–14.24
 normalization of, **I**:14.19
 physically realizable, **I**:14.40–14.42
 polar decomposition of matrices,
I:14.39–14.40
 and polarimetry, **I**:15.8–15.9, 15.11
 elements of, **I**:15.13–15.14
 in error analysis, **I**:15.28
 singular value decomposition,
I:15.25–15.27
 polarizance, **I**:14.18
 and polarization, **I**:14.7, 14.8, 14.25–14.27,
 14.33
 average degree of polarization,
I:14.32–14.33
 degree of polarization surfaces and maps,
I:14.31–14.32, 14.32f
 ideal polarizers, **I**:14.8–14.10, 14.10t
 nonpolarizing, **I**:14.8
 for radiative transfer, **I**:9.11
 for refraction and reflection, **I**:14.20–14.22
 retarder, **I**:14.11, 14.12t, 14.13–14.15, 14.15f
 for single scattering, **I**:9.16, 9.17
 and Stokes parameters, **I**:14.4–14.6
 transmittance, **I**:14.16–14.17
- Mueller matrix bidirectional reflectance
 distribution function (MMBRDF),
I:15.39, 15.39f, 15.40f
- Mueller matrix formalism, **III**:18.1, 18.20,
 18.22
- Mueller matrix polarimeters, **I**:15.26–15.27
- Mueller polarimeters, **I**:15.4
- Mueller vectors, **I**:15.15
- Mueller-Jones matrices, **I**:14.24, 14.27–14.29 (*see*
also Nondepolarizing Mueller matrices)
- Müller cells, **III**:8.20–8.21, 14.11
- Muller convention, **I**:12.6
- Müller zone model, **III**:11.65f, 11.84, 11.85f
- Müller zone theories, **III**:11.6, 11.7f
- Multichannel Bragg cells (MCBC), **V**:6.30–6.31

- Multichannel detectors, **II**:38.9
- Multicomponent lens assemblies, **II**:6.5–6.11
 drop-in, **II**:6.6, 6.6f
 lathe, **II**:6.7–6.8, 6.8f
 and lens adjustments, **II**:6.8–6.11, 6.10f
 “poker chip,” **II**:6.8, 6.9f
 tightlytoleranced, **II**:6.7, 6.7f
- Multicomponent polarizers, **IV**:7.69
- Multicore fibers, **V**:25.2, 25.22
- Multidomain vertical alignment (MVA) cells,
V:8.25, 8.27–8.28
- Multielectron atoms, **I**:10.10–10.11
- Multi-energy imaging, **V**:54.9–54.10
- Multi-fiber push on (MPO) connectors, **V**:23.8
- Multifocal lenses, **I**:23.12, 23.12f, 23.13f
 contact lenses, **III**:14.28, 20.2
 intraocular lenses, **III**:21.1, 21.14–21.18,
 21.15f, 21.16f, 21.16t, 21.18f
- Multifoil Kirkpatrick-Baez optics, **V**:48.3f,
 48.3–48.4
- Multifoil lobster-eye optics, **V**:48.2f, 48.2–48.4,
 48.3f
- Multifoil optics (MFO), **V**:48.1–48.4, 48.2f,
 48.3f
- Multilayer (ML) coatings, **IV**:7.96–7.98, 19.4
- Multilayer Laue lenses (MLLs):
 with curved interfaces, **V**:42.14, 42.15f
 and hard x-rays, **V**:42.1–42.17
 history of, **V**:42.2f, 42.2–42.4, 42.3f
 instrumental beamline arrangement and
 measurements for, **V**:42.9f–42.12f,
 42.9–42.10
 limitations of, **V**:42.15–42.17, 42.16f–42.17f
 with magnetron-sputtered MLLs,
V:42.5–42.7, 42.6f–42.8f
 on MLLs with curved interfaces, **V**:42.14,
 42.15f
 Takagi-Taupin calculations for,
V:42.12–42.14
 volume diffraction calculations for,
V:42.4–42.5, 42.5f
 with wedged MLLs, **V**:42.12–42.13, 42.13f,
 42.14f
 history of, **V**:42.2f, 42.2–42.4, 42.3f
 instrumental beamline arrangement and
 measurements of, **V**:42.9f–42.12f,
 42.9–42.10
 limitations of, **V**:42.15–42.17, 42.16f–42.17f
 magnetron-sputtered, **V**:42.5–42.7, 42.6f–42.8f
 Takagi-Taupin calculations for, **V**:42.12–42.14
- Multilayer Laue lenses (MLLs) (*Cont.*):
 volume diffraction calculations for,
V:42.4–42.5, 42.5f
 wedged, **V**:42.12–42.13, 42.13f, 42.14f
 and x-ray/neutron optics, **V**:26.10
- Multilayer reflectors, **IV**:7.39–7.53
 of absorbing materials, **IV**:7.37–7.38, 7.38f
 all-dielectric broadband reflectors, **IV**:7.39,
 7.40f, 7.45f–7.47f, 7.45–7.47
 coatings for ultrafast optics, **IV**:7.47–7.48, 7.48f
 for far-infrared region, **IV**:7.52, 7.52f
 graded reflectivity mirrors, **IV**:7.52
 imperfections in, **IV**:7.40–7.43, 7.41f–7.43f
 for interferometers and lasers, **IV**:7.39f–7.40f,
 7.39–7.40
 narrowband reflection coatings, **IV**:7.43,
 7.44f
 rejection filters, **IV**:7.48–7.50, 7.49f–7.51f
 for soft x-ray and XUV regions, **IV**:7.53
 in two-material periodic multilayers theory,
IV:7.37–7.38, 7.38f
- Multilayers, **V**:41.1–41.10
 and calculation of multilayer properties,
V:41.3–41.4
 for diffractive imaging, **V**:41.9–41.10, 41.10f
 fabrication and performance of, **V**:41.4–41.9,
 41.5f, 41.6t, 41.7t, 41.8f
 matrix theory for, **IV**:7.6–7.10, 7.9f
 periodic, **V**:42.5–42.6, 42.6f
 $[(0.5A)B(0.5A)]^N$, **IV**:7.35, 7.36f
 nonabsorbing $[AB]^N$ and $[AB]NA$,
IV:7.32–7.34, 7.33f–7.35f
 $[xH.(1-x)L]^N.xH$, **IV**:7.37
 properties of, **V**:41.1–41.3, 41.2f
 and x-ray diffraction, **V**:28.5, 28.5f
- Multilongitudinal mode Fabry-Perot laser,
V:9.7–9.8
- Multimode fibers, for E-LEDs, **V**:13.40
- Multimode interference (MMI) modulators,
V:13.2, 13.35
- Multimode interferometric Mach-Zehnder
 modulators, **V**:13.54f, 13.54–13.55
- Multipath interference noise, **V**:15.13–15.14
- Multiphase reflectors, **II**:39.39
- Multiphonon absorption, **IV**:5.16–5.17, 5.17t,
 5.18f, 19.6
- Multiphoton absorption:
 of crystals and glasses, **IV**:2.15–2.17,
 2.16f, 2.17f
 of solids, **IV**:8.30–8.31

- Multiphoton ionization (MPI), **IV**:21.10–21.12, 21.11*f*
- Multiphoton refraction, of crystals and glasses, **IV**:2.15–2.17, 2.16*f*, 2.17*f*
- Multiple anti-Stokes scattering, **IV**:15.2*t*
- Multiple beam interference, **I**:2.28–2.36
 diffraction gratings, **I**:2.28–2.29, 2.29*f*, 2.30*f*
 Fabry-Perot interferometers, **I**:2.33–2.36, 2.34*f*, 2.35*f*
 plane-parallel plates, **I**:2.30*f*, 2.30–2.33, 2.32*f*, 2.33*f*
- Multiple bound excitons, **IV**:5.26*t*
- Multiple cavities, in polymeric optics, **IV**:3.10
- Multiple Mirror Telescope (MMT), **V**:5.5
- Multiple quantum well (MQW) lasers, **V**:13.24–13.25
- Multiple quantum well (MQW) LEDs, **II**:18.1, 18.2*f*
- Multiple quantum wells (MQWs), **IV**:12.22*f*, 12.22–12.23; **V**:20.20
- Multiple Raman scattering, **IV**:15.2–15.3, 15.2*t*
- Multiple Raman Stokes generation, **IV**:15.38–15.40, 15.40*f*
- Multiple scattering, **V**:3.21, 3.21*f*
- Multiple (volume) scattering, **I**:9.2, 9.3, 9.8–9.17
 analytical theory of, **I**:9.9*f*, 9.9–9.10
 depolarization, **I**:9.16–9.17, 9.17*f*
 effective-medium representation, **I**:9.8
 radiative transfer, **I**:9.10–9.13, 9.11*f*, 9.13*f*
 speckle patterns, **I**:9.15*f*, 9.15–9.16
 weak localization, **I**:9.13–9.17, 9.14*f*
- Multiple Stokes scattering, **IV**:15.2–15.3, 15.2*t*, 15.3*t*
- Multiple surface concentrators, **II**:39.16–39.17, 39.17*f*
- Multiple wafer transducers, **II**:22.18–22.19
- Multiple-angle-of-incidence ellipsometry (MAIE), **I**:16.3
- Multiple-beam Bragg diffraction (MBD), **V**:43.6–43.8, 43.7*f*
- Multiple-Bragg-beam interference, **V**:43.1
- Multiple-layer surfaces, **IV**:6.26, 6.26*f*–6.27*f*
- Multiple-order retardation plates, **I**:13.47–13.48
- Multiple-pass interferometers, **II**:13.13
- Multiple-reflection filters, **IV**:7.111*f*–7.113*f*, 7.111–7.113
- Multiple-reflection interferometers, **II**:13.13
- Multiple-resonant oscillators, **IV**:17.16–17.21
 doubly-resonant, **IV**:17.16–17.17
 pump-enhanced singly resonant, **IV**:17.17–17.20, 17.18*f*–17.20*f*
 triply-resonant, **IV**:17.20–17.21, 17.21*f*
- Multiple-segment LED displays, **II**:17.31
- Multiple-track read-write with diode laser arrays, **I**:35.28
- Multiplexed image scanning, **I**:30.23, 30.24*f*
- Multiplexed sensors, for fiber interferometers, **I**:32.16
- Multiplexers and multiplexing:
 for networking, **V**:18.4, 18.10–18.11, 18.11*f*
 optical add/drop, **V**:21.2, 21.8, 21.8*f*, 21.9*f*
 in OTDM networks, **V**:20.1, 20.3–20.12, 20.5*f*–20.11*f*, 20.13*f*
 parallel, **V**:20.12
 serial, **V**:20.12
 spatial multiplexing, **I**:23.8
 time-division, **V**:9.12, 20.3, 21.3
 (*see also* Wavelength division multiplexing)
- Multiplicative adaptation, **III**:2.26
- Multiplier photodiodes, **II**:24.11
- Multiplier phototubes [*see* Photomultiplier tubes (PMTs)]
- Multiplier tubes, **II**:24.6
- Multippliers, serial incoherent matrix-vector, **I**:11.17–11.18, 11.18*f*
- Multipop (film exposure technique), **II**:30.22
- Multiquantum well (MQW) lasers, **II**:19.14, 19.15, 19.15*f*, 19.24, 19.25*t*
- Multiquantum wells (MQWs), **V**:2.11, 19.7, 19.11*f*, 19.21
- Multiquantum-well buried heterostructure (MQW BH) lasers, **II**:19.25*t*
- Multishot laser-induced damage, **IV**:19.4
- Multispeed choppers, **II**:15.14
- Multivapor arcs, **II**:15.29, 15.30*f*, 15.31*f*
- Multewire proportional counters (MWPCs), **V**:63.31–63.32
- Multizone intraocular lenses, **III**:14.28–14.29
- Munnerlyn's equations, **III**:16.15
- Murty's lateral shear interferometer, **II**:12.14, 13.12*f*
- Musculoskeletal, **III**:23.3
- Mutual coherence function (MCF), **I**:2.36*f*, 2.36–2.38, 5.4, 5.10, 6.2–6.3; **V**:4.4, 4.7, 4.10
- Mutual intensity, **I**:5.11, 6.3–6.4, 6.7, 6.8
- Mutually incoherent beam couplers, **IV**:12.7, 12.8*f*

- Mutually pumped phase conjugators (MPPCs), **IV**:12.7
- Mylar film, **II**:29.4
- Myoid, **III**:8.2
- Myopia (nearsightedness), **III**:1.6, 1.18, 12.4, 16.5
 correction of, **III**:13.17–13.18, 16.7–16.8
 defined, **III**:13.2, 19.1, 23.3
 and early presbyopia, **III**:12.8
 empty field, **III**:1.33
 and focus of collimated light, **III**:12.3*f*
 instrument, **III**:1.34–1.35
 night (twilight), **III**:1.34
 refractive surgery for, **III**:12.14
- N* (optical constant of water), **IV**:1.17, 1.17*f*
- N* on 1 annealing, **IV**:19.3, 19.4
- N*00*N* state, **IV**:23.9–23.12, 23.10*f*, 23.11*f*
- Nakamura biplates, **I**:13.56
- Nanofocusing, of hard x-rays (*see* Hard x-rays, nanofocusing of)
- Nanofocusing lenses (NFLs), **V**:37.8–37.11, 37.9*f*, 37.10*f*
- Nano-optic-measuring (NOM) machine, **V**:46.5
- Nanoplasma, in clusters, **IV**:21.34–21.35, 21.35*f*
- Nanosecond fiber systems, **V**:25.30–25.32
- Nanostructure semiconductors, **IV**:5.52
- Nanostructuring, **IV**:6.55, 6.59*f*
- Nanotubes, **IV**:6.55, 6.59*f*
- Naperian absorption coefficient, **III**:8.9
- Narrow linewidth fiber lasers, **V**:25.29–25.30
- Narrowband filters, **I**:3.3
- Narrowband reflection coatings (narrowband rejection filters), **IV**:7.43, 7.44*f*, 7.49
- Narrow-bandpass filters, **IV**:7.78–7.83, 7.79*f*, 7.80*f*, 7.82*f*–7.88*f*, 7.88–7.89, 7.89*f*
- NAS polymer, **IV**:3.4*t*
- National Physical Laboratory (NPL), **II**:15.4
- National Search Engine for Standards, **II**:4.11
- Natural broadening, emission-line, **II**:16.5, 16.6
- Natural broadening, of lineshapes, **I**:10.7
- Natural color (NC) film, **II**:30.27
- Natural guide star (NGS) sensing, **V**:5.21
- Natural line width, of spectral lines, **V**:56.4–56.5
- Natural linewidth (of transition), **II**:16.5
- Natural stop shift, of lenses, **I**:22.3
- Natural waters, **IV**:1.3, 1.13–1.15
- N-BaF10 glass (670472), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- N-BaK4 glass (569560), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- N-BaLF4 glass (580537), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- N-BaSF64 glass (704394), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- NbF1 glass (743492), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- N-BK7 glass (517642), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- NBS Black, **IV**:6.49–6.50, 6.50*f*
- Near Bragg diffraction, **V**:6.8–6.9, 6.12
- Near field, **V**:4.10
- Near infrared (NIR) radiation, **II**:24.3, 25.2
- Near triad, **III**:1.30
- Near-field diffraction, **V**:27.2–27.3, 27.3*f*
- Near-IR crystals, visible, **IV**:10.21*t*, 10.22*t*
- Nearpoint, **III**:23.3
- Nearsightedness (*see* Myopia)
- Near-UV spectrum, semiconductor interactions with, **IV**:5.4*f*, 5.5
- Negative color photographic films, **II**:30.25–30.28, 30.27*t*
- Negative core-cladding index difference, of photonic crystal fibers, **V**:11.14*f*–11.17*f*, 11.14–11.17
- Negative hydrogen (H⁻) ions, **V**:2.3
- Negative lenses, **IV**:3.13
- Negative orders of radiation, **V**:40.1
- Nematic phase, of liquid crystals, **V**:8.8, 8.11, 8.11*f*
- Neodymium (Nd) glass lasers, **II**:16.32–16.33
- Neodymium-doped fibers, **V**:25.23*t*, 25.23–25.24
- Neodymium-yttrium vanadate (Nd:VO₄) lasers, **II**:16.33
- Neodymium-yttrium-aluminum-garnet (Nd:YAG) lasers, **II**:16.32–16.33
- Neodymium-yttrium-lithide-fluoride (Nd:YLF) lasers, **II**:16.33
- Neon signs, **II**:40.39
- Neoprene, **IV**:6.32*f*, 6.34*f*
- Nernst glower, **II**:15.14, 15.15, 15.17, 15.17*f*, 15.19, 15.19*f*
- Nesonian illumination (*see* Abbe illumination system)
- Net complex amplitude, **I**:2.5
- Net irradiance, of water, **IV**:1.5*t*, 1.7*f*, 1.9
- Networking, micro-optics-based components for, **V**:18.1–18.12
 attenuators, **V**:18.2, 18.9
 beam splitters, **V**:18.6, 18.6*f*
 circulators, **V**:18.3, 18.3*f*, 18.10

- Networking, micro-optics-based components
for (*Cont.*):
directional couplers, **V**:18.2, 18.3, 18.3*f*,
18.9, 18.9*f*
Faraday rotators, **V**:18.7, 18.7*f*
filters, **V**:18.6
gratings, **V**:18.5–18.6, 18.6*f*
GRIN-rod lenses, **V**:18.7, 18.8, 18.8*f*
isolators, **V**:18.3, 18.10, 18.10*f*
mechanical switches, **V**:18.4, 18.5, 18.5*f*,
18.11, 18.11*f*, 18.12*f*
MEMS mirrors and switches, **V**:18.8, 18.8*f*,
18.11, 18.12*f*
multiplexers/demultiplexers/duplexers,
V:18.4, 18.4*f*, 18.10–18.11
network functions, **V**:18.2–18.5
polarizers, **V**:18.7, 18.7*f*
power splitters, **V**:18.2*f*, 18.2–18.3, 18.9, 18.9*f*
prisms, **V**:18.5, 18.5*f*
(*see also related topics, e.g.*: Communication
networks and systems)
- Networks, 3D photonic crystals and, **IV**:9.4–9.5
Neumann's principle, **IV**:2.5
Neural networks, **IV**:12.33–12.35
Neuronal receptive fields, resolving capacity of
the eye and, **III**:4.6
Neutral arsenic antisite (As_{Ga}^0), **IV**:18.3
Neutral atoms, trapping of, **IV**:20.21–20.39
in atomic beam brightening, **IV**:20.27*f*,
20.27–20.28
in atomic clocks, **IV**:20.28
in Bose-Einstein condensation,
IV:20.35–20.37, 20.36*f*
in dark states, **IV**:20.37–20.39, 20.38*f*
with magnetic traps, **IV**:20.21–20.23, 20.22*f*
with magneto-optical traps, **IV**:20.24*f*,
20.24–20.25, 20.26*f*
in optical lattices, **IV**:20.31–20.34,
20.32*f*–20.34*f*
with optical traps, **IV**:20.23*f*, 20.23–20.24
in ultracold collisions, **IV**:20.28–20.31,
20.30*f*, 20.31*f*
Neutral attenuators, **IV**:7.105, 7.105*f*
Neutral density filters, **II**:40.52
Neutral filters, **IV**:7.67, 7.67*f*–7.68*f*
Neutron attenuation, **V**:63.11–63.12
Neutron collimation, **V**:63.15–63.16
Neutron filters, **V**:63.18*t*, 63.18–63.19, 63.19*f*
Neutron gravity spectrometer, **V**:63.21*f*
Neutron guides, **V**:63.15–63.18, 63.17*f*
Neutron optics, **V**:63.3–63.34
detection in, **V**:63.31–63.34
devices for, **V**:63.15–63.19, 63.17*f*, 63.18*t*,
63.19*f*
diffraction and interference in,
V:63.23–63.27, 63.27*f*
grazing-incidence, **V**:64.1–64.7
diffractive scattering and mirror surface
roughness, **V**:64.2–64.3
imaging focusing optics, **V**:64.3–64.7,
64.4*f*–64.7*f*
materials of optical elements, **V**:64.7
total external reflection, **V**:64.1–64.2
and neutron physics, **V**:63.3–63.5
and neutron sources, **V**:63.12–63.15, 63.13*f*
polarization techniques for, **V**:63.27–63.30,
63.30*f*
refraction and reflection in, **V**:63.19–63.23,
63.21*f*, 63.23*f*
scattering lengths and cross sections,
V:63.5–63.12, 63.6*t*, 63.10*t*
neutron attenuation, **V**:63.11–63.12
scattering length density, **V**:63.9–63.11,
63.11*f*
and x-ray optics, **V**:26.5–26.11, 26.8*f*, 26.9*f*,
26.11*f*, 36.2*f*
Neutron physics, **V**:63.3–63.5
Neutron polarization, **V**:63.27–63.29, 63.30*f*
Neutron production, fusion, **IV**:21.53
Neutron scattering, **I**:9.6
Neutron scintillators, **V**:63.33
Neutron zone plates, **V**:63.25
Neutrons:
epithermal, **V**:63.18
MCP detectors for, **V**:63.34
scattering cross sections of, **V**:63.6*t*,
63.6–63.9, 63.10*t*
scattering length densities of, **V**:63.9–63.11,
63.11*f*
scattering lengths of, **V**:63.5–63.9, 63.6*t*,
63.10*t*
thermal, **V**:63.3
total integrated scatter of, **V**:64.2–64.3
(*see also* Neutron optics)
- Newton interferometers, **I**:2.25
Newton ring method, **I**:29.22
Newtonian form, of Gaussian equations, **I**:18.4
Newtonian imaging equation, **I**:17.8
Newtonian objectives, **I**:29.6

- Newtonian viewing [*see* Optical generation of visual stimulus, free (newtonian) viewing]
- Newton's equation, for Gaussian focal lenses, **I**:1.49
- Newton's ring pattern, **I**:2.10, 2.19, 2.25–2.26
- NeXT spacecraft, **V**:47.10
- Nextel 2010, **IV**:6.35, 6.37
- Nextel Suede Coating Series 3101-C10, **IV**:6.37, 6.38*f*, 6.53*f*
- N-F2 glass (620364), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- N-FKS glass (487704), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- Nickel:
 - absorptance of, **IV**:4.41*f*, 4.48*t*, 4.50*t*, 4.51*t*
 - optical properties of, **IV**:4.16*t*–4.17*t*, 4.24*f*
 - penetration depth, **IV**:4.47*f*
 - physical properties of, **IV**:4.52*t*, 4.54*t*
 - reflectance of, **IV**:4.35*t*–4.36*t*, 4.41*f*, 4.47*f*
 - thermal properties of
 - coefficient of linear thermal expansion, **IV**:4.56*t*, 4.57*f*
 - elastic stiffness, **IV**:4.69*t*
 - moduli and Poisson's ratio, **IV**:4.69*t*
 - at room temperature, **IV**:4.55*t*
 - specific heat, **IV**:4.65*t*, 4.67*f*
 - strength and fracture properties, **IV**:4.70*t*
 - thermal conductivity, **IV**:4.58*f*, 4.62*f*–4.63*f*
- Nickel alloys, **IV**:4.47*f*
- Nicol curtate prisms, **I**:13.16*f*, 13.17
- Nicol-type polarizers, **I**:13.6, 13.6*f*, 13.10*f*, 13.15–13.18
 - conventional, **I**:13.6*f*, 13.15–13.16
 - Glan-type vs., **I**:13.8–13.9
 - trimmed, **I**:13.16*f*, 13.16–13.18
- Night (twilight) myopia, **III**:1.34
- Nikon N8008s camera, **I**:25.7
- Nikon Plan Apochromat, **I**:28.13, 28.14*f*
- Niobium flint glass, **IV**:2.42*t*
- Nipi modulators, **V**:13.62
- Nipkow disk tandem-scanning confocal microscope, **III**:17.2, 17.4*f*, 17.5, 17.5*f*
- Nipkow disks, **I**:28.50, 28.51*f*
- Nippon Sheet Glass, **I**:24.2, 24.6, 24.7
- Nit (unit), **II**:34.43, 36.7, 36.8*t*
- Nitride LEDs, **II**:17.19, 18.3*f*
- Nitrogen doping, **II**:17.16
- Nitrogen-doped GaAsP, **II**:17.16, 17.21–17.22
- N-K5 glass (522595), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- N-KF9 glass (523515), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- N-KzFS4 glass (613443), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- N-LaF2 glass (744447), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- N-LaF33 glass (754523), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- N-LaK10 glass (720504), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- N-LaSF31A glass (883409), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- Nodal planes, of Gaussian lenses, **I**:1.48–1.49
- Nodal plane-to-nodal plane conjugate matrices, **I**:1.69
- Nodal points, of lenses, **I**:1.48, 1.48*f*, 1.49, 17.7
- Nodal rays, **I**:17.16
- Nodal slide bench, **II**:12.22, 12.22*f*
- Noether's theorem, **I**:1.21
- Noise, **II**:27.3–27.6
 - 1/[function], **II**:27.4
 - ASE, **V**:19.3, 19.4*f*, 19.9, 19.18, 19.24, 19.35
 - atomic, **II**:23.34–23.35
 - avalanche photodiodes, **V**:13.72–13.73
 - in CCDs, **II**:32.20, 32.32
 - data, **I**:35.24
 - of EDFAs, **V**:14.6
 - excess, **II**:24.11
 - in fiber optic communication links, **V**:15.11–15.14, 15.16–15.17
 - fixed-pattern, **I**:26.11
 - generation, **II**:24.11
 - generation-recombination, **II**:24.11
 - ground loop, **II**:27.5, 27.6*f*
 - inductive pickup, **II**:27.5, 27.6*f*
 - laser, **I**:35.12, 35.24
 - of laser diodes, **V**:13.18–13.24, 13.20*f*
 - measurement, **II**:22.7–22.8
 - modal, **V**:15.16–15.17
 - mode partition, **V**:13.19–13.20, 15.11*f*, 15.11–15.13
 - modulation, **II**:24.11
 - multipath interference, **V**:15.13–15.14
 - nonessential, **II**:27.4–27.6, 27.5*f*, 27.6*f*
 - and optical disk data, **I**:35.12, 35.23–35.24, 35.24*f*
 - pattern, **I**:26.11–26.12, 26.12*f*; **II**:32.12
 - phase (linewidth), **V**:13.20–13.21, 13.21*f*
 - in photodetectors, **II**:24.19–24.20, 24.20*f*

- Noise (*Cont.*):
 in photoemissive detectors, **II**:24.39, 24.40, 24.41*f*
 of pin diodes, **V**:13.70–13.71
 relative intensity, **V**:13.18–13.19, 13.19*f*, 15.14
 resistive coupling, **II**:27.5, 27.6*f*
 RMS, **II**:24.12
 shot, **I**:26.11; **II**:24.12, 27.3, 27.3*f*, 32.12; **IV**:23.4–23.6, 23.5*f*
 and signal detection, **II**:27.1
 of SOAs, **V**:19.3, 19.4*f*, 19.9, 19.18, 19.20, 19.24, 19.35
 in solid-state cameras, **I**:26.11–26.12, 26.12*f*
 spatial, **II**:33.26*f*, 33.26–33.27
 and stimulated Raman scattering, **IV**:15.35–15.38, 15.39*f*
 stray capacitance, **II**:27.5, 27.6*f*
 temperature, **II**:24.12
 thermal, **II**:24.13, 27.4, 32.20
 [see also Signal-to-noise ratio (SNR)]
 Noise (color vision), threshold contours for, **III**:11.20, 11.23*f*, 11.23–11.26, 11.25*f*
 Noise equivalent BDSF (NEBDSF), for scatterometers, **V**:1.6, 1.8, 1.12–1.13
 Noise equivalent exposure (NEE), **I**:26.10
 Noise equivalent irradiance (NEI), **II**:24.11
 Noise equivalent power (NEP), **II**:24.10, 24.12, 24.14, 25.12, 28.2, 38.9, 38.10, 38.10*t*; **V**:1.13, 13.71
 Noise equivalent quanta (NEQ), **II**:29.1, 29.23
 Noise equivalent temperature difference (NETD), **II**:28.8–28.9, 33.2, 33.24–33.27, 33.25*f*, 33.26*f*
 Noise figure, of EDFAs, **V**:14.6
 Noise masking (color vision):
 defined, **III**:11.3
 experiments in, **III**:11.56–11.57
 Noise spectral density (NSD), **II**:33.2
 Noise spectrum, **II**:24.11
 Nomarski differential interference, **V**:46.4
 Nomarski interferometers, **I**:32.4, 32.5*f*
 Nomarski microscope, **II**:10.11, 10.11*f*; **V**:46.2
 Nomarski prisms, **I**:28.40, 28.41
 Nomogram, **III**:16.1
 Nonabsorbing $[AB]^N$ and $[AB]NA$ multilayers, **IV**:7.32–7.34, 7.33*f*–7.35*f*
 Nonabsorbing mirrors (NAMs), **II**:19.23–19.24
 Nonafocal lenses, **I**:1.46 (see also Focal lenses)
 Nonblackbody radiation source, **II**:15.10
 Noncalcite prisms, **I**:13.23–13.24
 Nonchromogenic film, **II**:29.14
 Noncircular pupils, **II**:11.4, 11.37, 11.39
 Nonconcomitant movement of eyes, **III**:13.2
 Noncosmological red shift, **I**:5.23
 Noncritical phase matching (NCPM), **IV**:17.1
 Noncritical phase-matching acousto-optic tunable filters (NPM AOTFs), **V**:6.37, 6.38*f*, 6.39–6.42
 angle of deflection, **V**:6.40
 angular aperture, **V**:6.39
 for mid-infrared, **V**:6.42
 optical throughput, **V**:6.41
 performance of, **V**:6.42–6.44, 6.43*t*
 resolution, **V**:6.40
 sidelobe suppression, **V**:6.41–6.42
 transmission and drive power, **V**:6.41
 tuning relation, **V**:6.32
 for ultraviolet, **V**:6.42
 Nondegenerate four-wave mixing (NDFWM), **IV**:14.28
 Nondepolarizing Mueller matrices, **I**:14.24–14.25, 14.27–14.30
 Non-diffraction-limited optics, **II**:8.6
 Nondispersive prisms, **I**:19.1–19.29
 Abbe's, **I**:19.3*t*, 19.7, 19.7*f*–19.8*f*
 Amici (roof), **I**:19.3*t*, 19.11, 19.12*f*
 and beam deviation, **I**:19.2
 and beam displacement, **I**:19.2
 Brashear-Hastings, **I**:19.3*t*, 19.25, 19.25*f*
 Carl Zeiss, **I**:19.3*t*, 19.16, 19.16*f*
 Dove, **I**:19.3*t*, 19.9, 19.9*f*, 19.10, 19.10*f*
 Frankford Arsenal, **I**:19.3*t*, 19.18*f*–19.24*f*, 19.18–19.24
 general deviation, **I**:19.3*t*, 19.28, 19.28*f*–19.29*f*
 Goerz, **I**:19.3*t*, 19.17, 19.17*f*
 and image inversion/reversion, **I**:19.2
 Leman, **I**:19.3*t*, 19.13, 19.13*f*
 Pechan, **I**:19.3*t*, 19.11, 19.11*f*
 penta, **I**:19.3*t*, 19.13, 19.14*f*
 Porro, **I**:19.3, 19.3*t*, 19.5*f*, 19.6, 19.6*f*
 retroreflectors, **I**:19.3*t*, 19.28, 19.28*f*
 reversion, **I**:19.3*t*, 19.14, 19.14*f*
 rhomboidal, **I**:19.3*t*, 19.25, 19.25*f*
 right-angle, **I**:19.3, 19.3*t*, 19.4*f*
 Risley, **I**:19.3*t*, 19.25, 19.25*f*–19.27*f*, 19.27
 Schmidt, **I**:19.3*t*, 19.12, 19.12*f*
 Wollaston, **I**:19.3*t*, 19.15, 19.15*f*
 Nonequilibrium errors, **II**:34.25

- Nonessential noise, **II**:27.4–27.6, 27.5*f*, 27.6*f*
 Nonessential ray tracing, **II**:40.20–40.21
 Nonexudative (dry) age-related macular degeneration, **III**:14.1, 14.24–14.35
 Nonfused fiber couplers, **V**:25.10
 Nonhomogeneous polarization elements, **I**:14.25–14.26, 14.26*f*, 15.7, 15.20
 Nonhomogeneous polarization elements (Mueller matrices), **I**:14.25–14.26, 14.26*f*
 Nonideal aperture, **II**:34.35*f*, 34.35–34.36
 Nonimaging concentrators (*see* Concentrators, nonimaging)
 Nonimaging optics, **II**:39.1–39.41
 about, **II**:39.1–39.2
 aspheric lenses in, **II**:39.8, 39.9, 39.9*f*
 calculations for, **II**:39.2–39.6
 clipped Lambertian distribution, **II**:39.3–39.4
 concentration, **II**:39.5, 39.6
 dilution, **II**:39.6
 etendue, **II**:39.2, 39.3, 39.4*f*
 Hottel strings, **II**:39.4, 39.4*f*
 Lambertian, **II**:39.3
 luminance, **II**:39.3
 projected solid angle, **II**:39.5, 39.5*f*
 solid angle, **II**:39.5, 39.5*f*
 concentration of, **II**:39.12–39.22
 calculation, **II**:39.5, 39.6
 compound elliptical collectors, **II**:39.14, 39.15*f*
 compound hyperbolic collectors, **II**:39.15, 39.15*f*, 39.16*f*
 compound parabolic collectors, **II**:39.13*f*, 39.13–39.14, 39.14*f*
 dielectric compound parabolic collectors, **II**:39.15, 39.16, 39.16*f*
 edge rays, **II**:39.22
 geometrical vector flux, **II**:39.21–39.22
 inhomogeneous media, **II**:39.22
 multiple surface concentrators, **II**:39.16–39.17, 39.17*f*
 restricted exit angle concentrators with lenses, **II**:39.18, 39.18*f*
 tapered lightpipes, **II**:39.12–39.13, 39.13*f*
 θ_1/θ_2 concentrators, **II**:39.18–39.20, 39.19*f*
 2D vs. 3D geometries, **II**:39.20*f*, 39.20–39.21, 39.21*f*
 conic reflectors in, **II**:39.11, 39.11*f*
 Fresnel lenses in, **II**:39.9–39.10, 39.10*f*
 involute reflectors in, **II**:39.11–39.12, 39.12*f*
 Nonimaging optics (*Cont.*):
 macrofocal reflectors in, **II**:39.11
 software modeling of, **II**:39.6–39.8
 spherical lenses in, **II**:39.8, 39.9*f*
 terminology, **II**:39.2, 39.2*t*
 uniform illumination of, **II**:39.22–39.41
 classic projection system uniformity, **II**:39.23*f*, 39.23–39.24
 faceted structures, **II**:39.39*f*, 39.39–39.41, 39.40*f*
 integrating cavities, **II**:39.24*f*, 39.24–39.27, 39.25*f*, 39.27*f*
 lens arrays, **II**:39.32–39.37, 39.33*f*–39.37*f*
 lightpipes, **II**:39.27–39.32, 39.28*f*–39.31*f*
 tailored reflectors, **II**:39.37–39.39, 39.38*f*
 Nonimaging software modeling, **II**:39.6–39.8
 Noninterferometric optical testing, **II**:13.1–13.7
 Foucault test, **II**:13.2*f*, 13.2–13.3, 13.3*f*
 Hartmann test, **II**:13.4–13.6, 13.5*f*
 Hartmann-Shack test, **II**:13.6*f*, 13.6–13.7
 Ronchi test, **II**:13.3*f*, 13.3–13.4, 13.4*f*
 Noninterlaced monitors, **III**:23.3
 Nonlinear absorption (NLA), **IV**:16.29
 limiters of, **IV**:13.6*f*, 13.6–13.7
 mechanisms of, **IV**:16.12–16.13, 16.13*f*
 nondegenerate, **IV**:16.27
 optical limiting by, **IV**:13.4*f*, 13.4–13.7, 13.5*f*
 and third-order optical nonlinearities, **IV**:16.7–16.9
 Nonlinear acoustic (NA) interaction, in acousto-optic devices, **V**:6.30
 Nonlinear atom-field interactions, **IV**:11.10
 Nonlinear distortion, **V**:9.17
 Nonlinear effects:
 of fiber lasers, **V**:25.6
 four-wave mixing, **V**:10.2, 10.9–10.11, 10.11*f*
 in integrated optics, **I**:21.12
 in optical fibers, **V**:10.1–10.12
 self- and cross-phase modulation, **V**:10.3–10.4
 stimulated Brillouin scattering, **V**:10.1, 10.7–10.9
 stimulated Raman scattering, **V**:10.1, 10.4–10.7, 10.5*f*
 Nonlinear length, of solitons, **V**:22.3
 Nonlinear optical coefficients, of crystals and glasses, **IV**:2.26–2.27, 2.27*t*
 Nonlinear optical crystals, **IV**:10.19–10.20, 10.20*t*–10.22*t*

- Nonlinear optical (NLO) effects, **IV**:19.5, 19.9f, 19.9–19.11, 19.10f
- Nonlinear optical frequency conversion, **IV**:14.24f, 14.24–14.28, 14.27f
- Nonlinear optical loop mirrors (NOLM), **V**:20.22
- Nonlinear optical properties (of semiconductors), **IV**:5.52–5.56
- Maxwell's equations and polarization power series expansion, **IV**:5.52–5.53, 5.54t
- second-order, **IV**:5.53–5.55
- third-harmonic generation, **IV**:5.56
- third-order, **IV**:5.55
- two-photon absorption, **IV**:5.56
- Nonlinear optics, **III**:17.2, 17.3; **IV**:10.3–10.23
- about, **IV**:10.4–10.5
- conversion efficiencies, **IV**:10.14–10.16
- crystals for, **IV**:10.19–10.20, 10.20t–10.22t
- equations for, **IV**:10.4–10.5
- microscopic origin of, **IV**:10.5–10.10, 10.6f, 10.8f
- and MKS systems, **IV**:10.21–10.23
- optical parametric process in, **IV**:10.16–10.19, 10.17f–10.19f
- phase-matching condition in second-order processes, **IV**:10.12f, 10.12–10.14, 10.13f
- second-order susceptibility tensor in, **IV**:10.10–10.11
- strong field, **IV**:21.27–21.31, 21.28f
- third-order optical nonlinearities, **IV**:16.1–16.31
- cascaded $x^{(1)}:x^{(1)}$ processes, **IV**:16.20–16.22, 16.21f
- cascaded $x^{(2)}:x^{(2)}$ processes, **IV**:16.22–16.24, 16.23f, 16.24f
- four-wave mixing, **IV**:16.27–16.28, 16.28f
- interferometry, **IV**:16.28–16.29
- Kerr effect, **IV**:16.11–16.14, 16.13f, 16.14f
- Kramers-Kronig dispersion relations, **IV**:16.9–16.11
- nonlinear absorption and nonlinear refraction, **IV**:16.7–16.9
- propagation effects, **IV**:16.24–16.26
- and quantum mechanics, **IV**:16.4–16.7, 16.5f
- stimulated scattering, **IV**:16.14–16.19, 16.15f, 16.17f
- terms for, **IV**:16.1–16.3, 16.3t
- third-harmonic generation, **IV**:16.14
- Nonlinear optics, third-order optical nonlinearities (*Cont.*):
- time-resolved excite-probe techniques, **IV**:16.26f, 16.26–16.27
- two-photon absorption, **IV**:16.19–16.20
- Z-scan, **IV**:16.29–16.30, 16.30f
- ultrashort pulse generation, **IV**:18.1–18.23
- Kerr effect, **IV**:18.11–18.15, 18.12f
- saturable absorbers, **IV**:18.5–18.11, 18.6f–18.8f
- semiconductor ultrafast nonlinearities, **IV**:18.15–18.23, 18.16f, 18.17f, 18.22f
- and ultrafast lasers, **IV**:18.3–18.5, 18.4f, 18.5f
- WDM networks and, **V**:21.18–21.20, 21.19f, 21.20f
- Nonlinear reflectivity, **IV**:18.6–18.7, 18.7f
- Nonlinear refraction (NLR), **IV**:13.4f, 13.7–13.8, 16.7–16.9, 16.29–16.30, 16.30f
- Nonlinear scattering, **IV**:13.4f, 13.8
- Nonlinear Schrödinger equation (NLSE), **IV**:16.25; **V**:10.4, 22.2
- Nonlinear susceptibility, of crystals and glasses, **IV**:2.26
- Nonlinear Thomson scattering, **IV**:21.8–21.9, 21.9f
- Nonlinear visual mechanisms, **III**:11.3
- Nonlinearity correction factor, **II**:34.33
- Nonlinearity measurement, **II**:34.34–34.35
- Nonneutral (color) density, of photographic films, **II**:29.7–29.8
- Nonnormal angle of incidence, antireflection coatings at, **IV**:7.28f–7.31f, 7.28–7.31
- Nonnormal-incidence reflection:
- in Brewster angle reflection polarizers, **I**:13.34–13.37, 13.34t–13.36t
- in pile-of-plates polarizers, **I**:12.15–12.18, 12.16f, 12.17f
- in polarizing beam splitters, **I**:13.41–13.42
- Nonnormal-incidence transmission:
- in Brewster angle transmission polarizers, **I**:13.37–13.39, 13.38t–13.39t
- in interference polarizers, **I**:13.39–13.41, 13.40f
- in pile-of-plates polarizers, **I**:12.18–12.24, 12.19t–12.20t, 12.21f
- in polarizing beam splitters, **I**:13.41–13.42
- Nonpolarizing beam splitters, **IV**:7.63, 7.64f–7.65f, 7.65

- Nonpolarizing edge and bandpass filters, **IV**:7.66, 7.67*f*
- Nonpolarizing elements, **I**:15.7
- Nonpolarizing Mueller matrices, **I**:14.8
- Nonradiating sources of light, **I**:5.18
- Nonradiative recombination, **II**:17.2
- Nonreacting interferometers, **II**:12.7
- Nonreactive evaporation, **IV**:7.11
- Nonrectilinear distortion, **I**:27.6
- Nonresonant degenerate four-wave mixing, **IV**:16.27–16.28, 16.28*f*
- Non-return-to-zero differential-phase-shift-keying (NRZ-DPSK) format, **V**:21.28
- Non-return-to-zero (NRZ) format, **I**:35.17
in OTDM networks, **V**:20.8, 20.9*f*, 20.12
in WDM networks, **V**:21.16, 21.29*f*, 21.32*f*
- Non-return to zero inverted (NRZI) scheme, **I**:35.17
- Non-return-to zero on-off keying (NRZ-OOK), **V**:21.34, 21.36*t*, 21.37*t*
- Nonrotationally symmetric systems, **I**:1.74
- Nonsequential ray tracing, **II**:39.6
- Nonsequential surfaces data, **II**:3.6–3.7
- Nonspherical particles, scattering by, **I**:7.15–7.17
- Nonuniform rational B-splines (NURBS), **II**:39.6, 40.41
- Nonuniformity, of radiation distribution, **II**:34.25, 34.28, 34.35
- Nonzero dispersion-shifted fiber (NZDSF), **V**:21.21, 21.21*f*, 21.34*f*
- No-phonon (NP) photoluminescence, **IV**:5.72, 5.73*f*
- Normal congruence rays, **I**:1.10
- Normal equations, **II**:3.18
- Normal vectors, **I**:1.18, 1.19
- Normal-incidence rotating-sample ellipsometers (NIRSE), **I**:16.18
- Normalization, of Mueller matrices, **I**:14.19
- Normalized detectivity, **II**:38.9
- Normalized detector irradiance (NDI), **II**:7.22
- Normalized detuning, **IV**:14.11*f*
- Normalized spectrum, coherence functions for, **I**:5.5–5.6
- Normalized vector potential, **IV**:21.6
- Normalized water-leaving irradiance, **IV**:1.46–1.48, 1.48*f*
- Notch filters, **II**:22.10*f*, 22.10–22.11, 22.11*f*; **IV**:7.43
- Novelty filters, **IV**:12.32, 12.33*f*–12.35*f*
- np silicon photodiodes, **II**:34.30
- N-PK52A glass (497816), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- N-PSK53A glass (618634), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- N-SF6 (805254), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- N-SK10 glass (623570), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- N-SSK5 glass (658509), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- Nuclear cataract, **III**:7.5, 12.14, 14.8
- Nuclear imaging, **V**:53.17, 53.18, 53.18*f*
- Nuclear light sources, **II**:40.39
- Nuclear medicine, **V**:32.1–32.4
- Nuclear motion, in strong fields, **IV**:21.22–21.23, 21.23*t*, 21.24*f*
- Null correctors, optical, **IV**:3.16
- Null ellipsometers, **I**:16.11, 16.12
- Null optics, polymers and, **IV**:3.16–3.17
- Nulling, in adjustment experiments, **III**:3.4
- Nulling interferometers, **I**:32.20–32.21
- Numeric displays, LED, **II**:17.30*f*, 17.30–17.31, 17.31*f*
- Numerical aperture, **I**:1.78, 1.79, 17.9; **III**:17.2
- Numerical aperture (NA), **II**:34.20, 39.1; **V**:9.4, 25.2, 25.18, 42.2
- Numerically controlled machines (CNCs), **II**:12.11
- NuSTAR spacecraft, **V**:47.6, 47.7, 47.10
- Nutting's law, **II**:29.6
- Nyquist condition, **II**:13.27
- Nyquist frequency, **I**:26.16–26.20, 26.17*f*; **V**:27.4
- Nyquist frequency power, **V**:46.8–46.9, 46.11
- Nyquist limit, **III**:2.6, 2.7, 2.11, 2.22–2.23
- Nyquist noise, **V**:13.70 (*see* Thermal noise)
- N-ZK7 glass (508612), **IV**:2.49*t*, 2.54*t*, 2.59*t*, 2.66*t*
- Oak Ridge, **IV**:6.51
- Obedience to Abney's law, **III**:10.44, 11.37 (*see also* Additivity)
- Object counting, reflexive sensors for, **II**:17.34
- Object relief distance, **I**:18.11–18.12, 18.12*f*
- Object space, **I**:1.26, 1.83
- Object space numerical aperture, **I**:1.78
- Object space pupil diameter, **I**:18.6
- Object transparencies, **I**:11.4–11.5, 11.5*f*
- Objective amplitude of accommodation, **III**:1.32
- Objective optics, **I**:30.30–30.33, 30.32*f*–30.33*f*

- Objective scanning, **I**:30.5, 30.28
- Objective speckle, **I**:33.9
- Objective tasks, **III**:2.15*n*
- Objective tone reproduction, **II**:29.16–29.17, 29.17*f*
- Objectives:
- for microscopes, **I**:28.9–28.15, 28.10*t*, 28.11*t*
 - in afocal systems, **I**:18.7
 - corrections for tube length, **I**:28.13, 28.13*t*
 - coverslip correction, **I**:28.10–28.11, 28.12*f*, 28.13
 - design of, **I**:28.13–28.15, 28.14*f*–28.16*f*
 - field size, **I**:28.13
 - working distance, **I**:28.13
 - reflective and catadioptric (*see* Reflective and catadioptric objectives)
 - Schwarzschild, **V**:26.10, 51.1–51.3, 51.2*f*–**V**:51.4*f*
 - for telescopes
 - afocal Cassegrain-Mersenne, **I**:29.9
 - afocal Gregorian-Mersenne, **I**:29.12
 - Ritchey-Chretien, **I**:29.8
 - three-mirror afocal, **I**:29.29–29.30
- Object-space scanners, **I**:30.18–30.23, 30.19*f*–30.23*f*
- Object-to-image distance, **I**:1.51–1.52
- Oblique effect (pattern discrimination), **III**:2.35
- Oblique spherical aberrations, **I**:1.90, 29.15, 29.21, 29.37
- Obliquity factor, **I**:3.5
- Obscurations, of reflective and catadioptric objectives, **I**:29.4*f*, 29.4–29.5
- Observatories:
- Chandra, **V**:33.2–33.4, 33.3*t*, 44.4, 44.10, 47.1, 47.4*f*, 47.5, 47.10, 64.7
 - Constellation-X, **V**:33.4
 - Einstein, **V**:44.4, 44.10, 47.1, 47.5
 - ROSAT, **V**:47.5
 - SOHO, **V**:41.3
 - Suzaku, **V**:33.3*t*, 33.3–33.4
 - TRACE, **V**:41.3
 - W. M. Keck, **V**:5.27
 - XMM-Newton, **V**:33.3, 47.2, 47.4*f*, 47.6, 47.6*f*
 - XMM-Newton observatory, **V**:33.3*t*
 - x-ray, **V**:33.1–33.4, 33.3*t*
- Ocean color, **IV**:1.46
- Octocouplers, **II**:17.32*f*, 17.32–17.34
- Octopus lenses, **III**:19.7, 19.7*f*
- Ocular motility, **III**:23.3
- Ocular parameters, **III**:1.4–1.6
- Ocular radiation hazards, **III**:7.1–7.15
- examples of, **III**:7.8–7.9
 - exposure limits, **III**:7.9–7.11
 - exceeding, **III**:7.11
 - guidelines for visible light, **III**:7.10
 - IR, **III**:7.10–7.11
 - UV, **III**:7.9–7.10
 - injury mechanisms, **III**:7.2–7.3
 - action spectra, **III**:7.2, 7.4*f*
 - exposure duration and reciprocity, **III**:7.2–7.3
 - lamp safety standards, **III**:7.14–7.15
 - laser hazards, **III**:7.11–7.14, 7.12*f*, 7.13*f*
 - accidents, **III**:7.12–7.14
 - and eye protectors, **III**:7.14
 - safety standards, **III**:7.12
 - retinal irradiance calculations, **III**:7.7–7.8, 7.8*f*
 - types of injury, **III**:7.3–7.7, 7.5*f*
 - cataract, **III**:7.5–7.6, 7.6*f*
 - droplet keratopathies, **III**:7.6, 7.7
 - infrared cataract, **III**:7.7
 - photokeratitis, **III**:7.4
 - photoretinitis, **III**:7.7
 - pterygium, **III**:7.6, 7.7
- Ocular radiometry, **III**:1.11–1.12
- Ocular wavefronts, **III**:16.6–16.7, 16.7*t*
- Oculars, for microscopes, **I**:28.16–28.17
- Off-axis, eccentric-pupil Paul-Gregorian objective, **I**:29.28–29.29
- Off-axis angles, **II**:7.23
- Off-axis chromatic aberrations, **II**:2.2–2.4, 2.3*f*, 2.4*f*
- Off-axis double-pass grating spectrograph, **I**:20.10, 20.13*f*
- Off-axis image quality, **III**:1.18*f*, 1.18–1.19, 1.26–1.27, 1.27*f*
- Off-axis irradiance, **II**:34.16, 34.16*f*
- Off-axis rejection (OAR), **II**:7.23
- Office lighting, **II**:40.55, 40.56*t*
- Offner compensators, **II**:13.24, 13.24*f*
- Offner relay, **I**:29.33
- Offset, thermocouple junctions as source of, **II**:27.6, 27.6*f*
- Offset drift, **II**:27.11
- Offset quantum wells, **I**:21.19
- Offset subtraction error, **II**:34.32–34.33
- Ohmic contact, **II**:17.13
- Olympus water immersion microscope, **III**:17.10

- Omnidirectional reflectors, **IV**:9.2
 On-axis aberrations, **V**:45.6–45.8
 On-axis objective optics, **I**:30.30
 On-axis optics, **V**:64.3
 On-axis tangential phase matching, **V**:6.25–6.26
 On-blaze condition, **V**:38.2
 One mode components analysis, **III**:10.34
 $1 \times N$ power splitters, **V**:16.1, 16.4
 $1/[f]$ noise, **II**:25.12, 27.4
 1D profilometry, **V**:46.6
 One-dimensional optical molasses, **IV**:20.15*f*,
 20.15–20.16
 One-electron atoms, **I**:10.7–10.9, 10.8*f*, 10.9*f*
 One-electron transitions, **IV**:5.7
 135° linear polarizers, Mueller matrices for,
I:14.10*t*
 135° half-wave linear retarders, Mueller
 matrices for, **I**:14.12*t*
 135° quarter-wave linear retarders, Mueller
 matrices for, **I**:14.12*t*
 110 photographic film, **II**:30.21, 30.25
 On-off keying (OOK), in WDM networks,
V:21.29, 21.30, 21.34, 21.36*t*, 21.37*t*
 Open aperture Z-scan, **IV**:16.30
 Open arcs, **II**:15.22
 Open fiber control (OFC), for Fibre Channel
 standard, **V**:23.4
 Open systems, **IV**:11.17*f*, 11.17–11.18, 11.18*f*
 Open-loop gain function, **II**:22.9*f*, 22.9–22.10
 Open-tube diffusion, **II**:17.24
 Ophthalmic polarimetry, **I**:15.39, 15.41
 Ophthalmoheliosis, **III**:7.1, 7.5
 Ophthalmoscopes:
 defined, **III**:15.1
 flood-illuminated AO, **III**:15.3, 15.16–15.17
 scanning laser, **III**:15.2, 15.3, 15.17–15.19,
 15.18*f*, 15.19*f*
 Ophthalmoscopic (double-pass) methods
 (retinal image quality), **III**:1.22–1.23
 Opponent color spaces, **III**:10.18–10.19
 Opponent-colors theory, **III**:11.3, 11.5, 11.6*f*,
 11.62–11.63
 Optic axis, of calcite crystals, **I**:13.2, 13.2*f*,
 13.3*f*, 13.3*n*
 Optic disc, **III**:1.3*f*
 Optic fiber, defined, **III**:19.1
 Optic flow, **III**:13.4
 and bifocal jump, **III**:13.15, 13.16
 defined, **III**:13.2
 Optic flow fields, **III**:2.38–2.39, 2.39*f*
 Optic nerve head, **III**:18.2
 Optical aberrations:
 categories of, **III**:17.2
 in human eyes [*see* Aberrations (in human
 eye)]
 Optical absorption, measurements of,
V:2.2–2.13, 2.4*f*, 2.6*f*–2.8*f*, 2.10*f*, 2.12*f*
 Optical absorption spectrometers, **I**:31.2–31.5,
 31.5*f*
 Optical add/drop multiplexers (OADMs),
V:21.2, 21.8, 21.8*f*, 21.9*f*, 21.12
 Optical amplifiers:
 communications applications for, **V**:9.14
 semiconductor vs. fiber, **V**:9.13–9.14
 in WDM networks, **V**:21.37*f*, 21.37–21.44
 EDFA, **V**:21.38*f*–21.42*f*, 21.38–21.41
 Raman, **V**:21.42*f*–21.44*f*, 21.42–21.44
 (*see also* Optical fiber amplifiers)
 Optical analysis software, **II**:40.20
 Optical athermalization, **II**:8.12–8.15,
 8.13*t*–8.15*t*
 Optical axes, **I**:1.32, 14.8, 18.2, 29.5, 29.37
 Optical axis, **III**:1.3*f*
 Optical Bloch equations (OBEs), **IV**:11.3–11.6,
 20.3, 20.6
 Optical burst switching (OBS), **V**:21.11
 Optical cavities, **II**:19.18, 19.19*f*
 Optical cavity technique, **II**:12.20, 12.20*f*
 Optical cavity-based frequency discriminators,
II:22.14–22.16, 22.17*f*
 Optical center point, of lenses, **I**:17.16, 17.17
 “Optical Characterization in Microelectronics
 Manufacturing” (S. Perkowitz, D. G. Seiler,
 W. M. Duncan), **IV**:5.61–5.62
 Optical choppers, **II**:27.14
 Optical circulator, **III**:18.2
 Optical circulators, **V**:17.8*f*, 17.9
 Optical clock recovery, **V**:20.21*f*, 20.21–20.22
 Optical coherence microscopy (OCM),
I:28.44
 Optical coherence tomography (OCT), **I**:22.2,
 22.39, 22.40*f*, 28.43–28.44; **III**:15.3, 15.19,
 18.1–18.30
 at 1050 nm, **III**:18.15*f*–18.17*f*, 18.15–18.17
 autocorrelation noise, **III**:18.11–18.12
 combined with adaptive optics,
 II:15.19–15.21, 15.20*f*, 15.21*f*
 defined, **III**:15.2, 18.2
 depth dependent sensitivity, **III**:18.13*f*,
 18.13–18.14, 18.14*f*

- Optical coherence tomography (OCT) (*Cont.*):
 Doppler OCT, **III**:18.18, 18.18*f*–18.20*f*,
 18.18–18.19
 fringe washout, **III**:18.15
 motion artifacts, **III**:18.15
 optical frequency domain imaging, **III**:18.7*f*,
 18.7–18.9, 18.8*f*
 polarization sensitive OCT, **III**:18.18,
 18.20–18.27, 18.21*f*, 18.23*f*,
 18.25*f*–18.27*f*
 shot-noise-limited detection, **III**:18.12–18.13
 signal to noise ratio, **III**:18.11
 spectral domain OCT, **III**:18.5*f*, 18.5–18.7,
 18.6*f*, 18.9
 noise analysis of, **III**:18.9–18.10, 18.12*f*
 retinal imaging with, **III**:18.27–18.29,
 18.28*f*, 18.29*f*
 sensitivity advantage of, **III**:18.9
 Stratus OCT 3, **III**:15.20–15.21, 15.21*f*
 time domain OCT, **III**:18.3–18.5, 18.4*f*
- Optical communication systems, **II**:19.3
- Optical components:
 purchasing of, **II**:9.9
 specifications for systems vs., **II**:4.3
- Optical confinement factor, **II**:19.11
- Optical constants, **I**:7.12, 12.4–12.6, 16.5
 for coatings, **IV**:7.13
 and dielectric function, **IV**:5.8–5.9
 of metals, **IV**:4.11, 4.12*t*–4.19*t*,
 4.20*f*–4.26*f*
 of solids, **IV**:8.6–8.7, 8.15
 of water, **IV**:1.17, 1.17*f*
- Optical crossconnects (OXC)s, **V**:21.5*f*, 21.6,
 21.8, 21.8*f*, 21.10, 21.10*f*
- Optical damage, **V**:13.54, 25.6
- Optical density, of photographic films,
II:29.6–29.8, 29.7*f*
- Optical design:
 binocular vision factors in, **III**:13.1–13.35
 and refractive errors, **III**:12.1–12.17
 assessment of, **III**:12.5–12.8
 binocular factors, **III**:12.15–12.17
 correction of, **III**:12.8–12.15
 types of, **III**:12.4–12.5
- Optical design software, **II**:3.1–3.24,
 40.20–40.21
 about, **II**:3.2
 and computing environment, **II**:3.21
 data entry for, **II**:3.2–3.8
 design process flowchart, **II**:3.3*f*
- Optical design software (*Cont.*):
 evaluation function of, **II**:3.8–3.16
 aberrations, **II**:3.9–3.11
 paraxial analysis, **II**:3.8–3.9, 3.9*f*
 ray tracing, **II**:3.11–3.13, 3.12*f*
 spot-diagram analysis, **II**:3.13–3.16
 global optimization with, **II**:3.21
 optimization function of, **II**:3.16–3.21
 programming considerations for, **II**:3.7–3.8
 purchasing of, **II**:3.22–3.24
 setup routine in, **II**:3.7
 simulation with, **II**:3.21
- Optical disk data storage, **I**:35.1–35.30
 alternative storage media, **I**:35.29, 35.30
 automatic focusing, **I**:35.12–35.14, 35.13*f*
 automatic tracking, **I**:35.14*f*–35.16*f*,
 35.14–35.17
 data format and layout for, **I**:35.2–35.7,
 35.3*f*–35.5*f*
 developments in, **I**:35.28–35.30
 diffractive optics, **I**:35.28, 35.28*f*, 35.29
 direct overwrite, **I**:35.30, 35.30*f*
 materials for recording, **I**:35.25–35.28,
 35.26*f*, 35.27*f*
 multiple-track read-write with diode laser
 arrays, **I**:35.28
 and optical path, **I**:35.7*f*, 35.7–35.12,
 35.9*f*–35.11*f*
 readout, **I**:35.21–35.24, 35.22*f*, 35.24*f*
 thermomagnetic recording process,
I:35.17–35.20, 35.18*f*–35.20*f*
- Optical electric field, bulk modulators and,
V:7.18
- Optical errors, in head mounted display systems,
III:13.34–13.35
- Optical extent (étendue), **I**:1.22, 1.81, 13.7
- Optical fiber amplifiers, **V**:14.1–14.11
 categories and features of, **V**:14.1–14.2, 14.2*t*
 erbium-doped
 energy levels, **V**:14.4
 fast power transients, **V**:21.39*f*,
 21.39–21.41, 21.40*f*
 gain flattening, **V**:14.6–14.7, 21.38–21.39,
 21.39*f*
 gain formation, **V**:14.4–14.5, 14.5*f*
 gain peaking, **V**:21.38, 21.38*f*
 noise, **V**:14.6
 pump wavelength options, **V**:14.5–14.6
 semiconductor amplifiers vs., **V**:9.13, 9.14,
 14.1, 14.2*t*

- Optical fiber amplifiers, erbium-doped (*Cont.*):
 static gain dynamic and channel power equalization, **V**:21.41, 21.41*f*–21.42*f*
 in WDM networks, **V**:21.2*f*, 21.2–21.3
 erbium/ytterbium-doped, **V**:14.7–14.8
 parametric, **V**:14.10–14.11
 praseodymium-doped, **V**:14.7
 Raman fiber, **V**:14.8*f*, 14.8–14.9, 14.10*f*
 rare-earth-doped, **V**:14.2–14.4, 14.3*f*
 ytterbium-doped, **V**:14.7
- Optical fiber perform, **III**:19.1
- Optical fiber sensors, **V**:24.1–24.13
 extrinsic Fabry-Perot interferometric, **V**:24.2*f*, 24.2–24.4, 24.3*f*
 fiber Bragg grating, **V**:24.5–24.8, 24.6*f*–24.7*f*
 intrinsic Fabry-Perot interferometric sensors, **V**:24.4*f*, 24.4–24.5
 long-period grating sensors, **V**:24.8–24.13, 24.9*f*–24.12*f*, 24.11*t*, 24.13*t*
- Optical fibers, **III**:5.10, 5.11, 19.4
 in communication systems, **V**:9.3–9.17
 analog transmission, **V**:9.15–9.17
 bit rate, **V**:9.12
 distance limits, **V**:9.12–9.13
 fiber for, **V**:9.4–9.7, 9.5*f*, 9.6*f*
 fiber-optic networks, **V**:9.14–9.15
 optical amplifiers, **V**:9.13–9.14
 photodetectors, **V**:9.8
 receiver sensitivity, **V**:9.8–9.11
 repeater spacing, **V**:9.12–9.13
 technology, **V**:9.4–9.8
 transmitting sources, **V**:9.7–9.8
- infrared, **V**:12.1–12.13
 applications, **V**:12.13
 categories and properties of, **V**:12.1–12.3, 12.2*f*, 12.2*t*, 12.3*t*
 crystalline, **V**:12.2*t*, 12.3*t*, 12.7–12.10, 12.8*f*, 12.10*f*
 heavy-metal oxide glass in, **V**:12.2*t*–12.4*t*, 12.3–12.7, 12.5*f*–12.7*f*
 hollow waveguides, **V**:12.2*t*, 12.3*t*, 12.10–12.13, 12.12*f*
- nonlinear effects in, **V**:10.1–10.12
 four-wave mixing, **V**:10.2, 10.9–10.11, 10.11*f*
 self- and cross-phase modulation, **V**:10.3–10.4
 stimulated Brillouin scattering, **V**:10.1, 10.7–10.9
 stimulated Raman scattering, **V**:10.1, 10.4–10.7, 10.5*f*
- Optical fibers (*Cont.*):
 polarizers for, **I**:13.57
 [*see also related topics, e.g.*: Photonic crystal fibers (PCFs)]
- Optical frequency domain imaging (OFDI), **III**:18.3, 18.7*f*, 18.7–18.9, 18.8*f*
 at 1050 nm, **III**:18.15–18.17
 defined, **III**:18.2
 SD-OCT vs., **III**:18.9
- Optical frequency synthesis, **IV**:17.28, 17.29
- Optical generation of visual stimulus, **III**:5.1–5.25
 building an optical system, **III**:5.8–5.18
 alternating of source and retinal planes, **III**:5.8–5.9
 calibration, **III**:5.17*f*, 5.17–5.19
 combining lights, **III**:5.9*f*, 5.9–5.11, 5.10*f*
 controlling intensity, **III**:5.13–5.15
 controlling wavelength, **III**:5.11–5.13, 5.12*t*
 field quality, **III**:5.11
 generating complex patterns, **III**:5.16, 5.16*f*
 lenses, **III**:5.11
 turning field on/off, **III**:5.13, 5.13*t*
- coherent radiation, **III**:5.19, 5.21
- detectors, **III**:5.21, 5.22*t*
- free (newtonian) viewing, **III**:5.2–5.4, 5.3*f*
 limitations of, **III**:5.4
 retinal illuminance, **III**:5.2–5.3, 5.3*f*
 the troland, **III**:5.3–5.4
- light exposure and optical safety, **III**:5.18*f*, 5.18–5.19
- light sources, **III**:5.19, 5.19*t*, 5.20*t*
- Maxwellian viewing, **III**:5.4–5.8, 5.5*f*
 advantages of, **III**:5.7
 control of focus, **III**:5.4–5.5, 5.5*f*
 interferometers, **III**:5.24
 positioning of pupil in, **III**:5.9
 retinal conjugate plane, **III**:5.5*f*
 retinal illuminance, **III**:5.6–5.7
 size, **III**:5.5–5.6
 spatial frequency content of stimuli, **III**:5.8
 two-channel, **III**:5.21, 5.23*f*, 5.23–5.24
 size of stimulus, **III**:5.2
- Optical holeburning (OHB), **I**:10.18; **V**:2.13, 2.14*f*
- Optical holeburning (OHB) spectroscopy, **I**:31.24*f*–31.26*f*, 31.24–31.26
- Optical image transformers, **II**:39.21

- Optical indicatrix (index ellipsoid), of crystals and glasses, **IV**:2.18–2.19, 2.19*f*
- Optical insertion loss, of electro-optic modulators, **V**:7.36
- Optical interconnects, **IV**:12.31, 12.31*f*
- Optical invariants, **I**:18.7
- Optical Kerr effect (OKE), **IV**:16.11–16.14, 16.13*f*, 16.14*f*; **V**:7.11
- Optical lattices, **IV**:20.31–20.34, 20.32*f*–20.34*f*
- Optical limiters, **IV**:12.32
- Optical lithography (OL), **V**:34.1
- Optical matched filtering, for pattern recognition, **I**:11.12–11.14, 11.13*f*
- Optical metrology, **I**:15.35
- Optical mode conditioners, **V**:23.7
- Optical modes, of crystals and glasses, **IV**:2.68*t*–2.76*t*
 with cesium chloride structure, **IV**:2.68*t*
 with chalcopyrite structure, **IV**:2.74*t*
 with corundum structure, **IV**:2.70*t*
 with cubic perovskite structure, **IV**:2.73*t*
 with diamond structure, **IV**:2.68*t*
 with fluorite structure, **IV**:2.69*t*
 other structures, **IV**:2.74*t*–2.76*t*
 with α -quartz structure, **IV**:2.71*t*
 with rutile structure, **IV**:2.71*t*
 with scheelite structure, **IV**:2.72*t*
 with sodium chloride structure, **IV**:2.69*t*
 with spinel structure, **IV**:2.73*t*
 with tetragonal perovskite structure, **IV**:2.73*t*
 with trigonal selenium structure, **IV**:2.70*t*
 with wurtzite structure, **IV**:2.70*t*
 with zinblende structure, **IV**:2.69*t*
- Optical molasses (OM), **IV**:20.13–20.17
 defined, **IV**:20.3, 20.10
 Doppler cooling, **IV**:20.13–20.15, 20.14*f*
 one-dimensional, **IV**:20.15*f*, 20.15–20.16
 three-dimensional, **IV**:20.16*f*, 20.16–20.17
- Optical monitoring, in thin film manufacturing, **IV**:7.11
- Optical multichannel analyzers (OMAs), **II**:31.27–31.28
- Optical parametric amplifiers (OPAs), **IV**:23.13–23.14
- Optical parametric chirped pulse amplification (OPCPA), **IV**:21.5
- Optical parametric oscillators (OPOs), **II**:20.20*f*, 20.20–20.22, 20.21*f*; **IV**:10.18–10.19, 10.19*f*, 14.15 [*see also* Continuous-wave optical parametric oscillators (cw OPOs)]
- Optical parametric (OP) process, **IV**:10.16–10.19, 10.17*f*–10.19*f*
- Optical path, of optical disks, **I**:35.7*f*, 35.7–35.12, 35.9*f*–35.12*f*
- Optical path difference (OPD), **I**:2.7, 28.33–28.34, 28.35*f*, 28.36; **II**:2.1, 2.6, 3.12*f*, 3.12–3.13, 8.1, 8.7, 13.14–13.15
- Optical path length (OPL), **I**:1.11, 2.5
- Optical phonons, **IV**:5.14
- Optical power (of cornea), **III**:16.2 [*see* Radiant flux (power)]
- Optical power dependence, of electroabsorption modulators, **V**:13.59–13.60
- Optical power penalties, for fiber optic communication links, **V**:15.8–15.17, 15.10*f*, 15.11*f*
- Optical processing systems, for synthetic aperture radar data, **I**:11.7–11.8
- Optical pulse(s), **II**:20.2–20.15
 coupling of circulating, **II**:20.12*f*, 20.12–20.15, 20.15*f*
 in high gain oscillators, **II**:20.10–20.12, 20.11*f*
 in ideal cavity, **II**:20.6–20.7
 and pulse train, **II**:20.2–20.9
 single, **II**:20.2–20.3, 20.3*f*
 toward steady-state, **II**:20.9–20.12
- Optical pumping, **II**:16.16*f*, 16.16–16.19, 16.17*f*, 16.18*f*
- Optical radar, **II**:12.4, 12.5
- Optical Ramsey fringes, **IV**:11.20–11.22, 11.21*f*
- Optical responses, classification of, **IV**:5.12, 5.13*t*
- Optical safety, **III**:5.18*f*, 5.18–5.19
- Optical signal-to-noise ratio (OSNR):
 of SOAs, **V**:19.24–19.27
 for WDM networks, **V**:21.20, 21.28, 21.34
- Optical sine theorem, **I**:17.5
- Optical software (for stray light suppression), **II**:7.24–7.27
 advantages/disadvantages of, **II**:7.29, 7.29*f*
 ASAP, **II**:7.25
 CODE V, **II**:7.26
 FRED, **II**:7.25–7.26
 LightTools, **II**:7.26
 SPEOS, **II**:7.27
 TracePro, **II**:7.27
 ZEMAX, **II**:7.26–7.27

- Optical specifications, **II**:4.1–4.12
 about, **II**:4.1–4.2
 element description, **II**:4.8–4.10
 environmental, **II**:4.10
 image, **II**:4.3, 4.6–4.8, 4.8*f*
 mechanical vs., **II**:4.2
 preparing, **II**:4.5–4.6, 4.6*t*
 presentation of, **II**:4.10–4.11
 problems with writing, **II**:4.11–4.12
 for systems vs. components, **II**:4.3
 wavefront, **II**:4.3–4.5, 4.5*t*
- Optical spectrometers (*see* Spectrometers)
- Optical spectroscopy, **IV**:11.2 (*see* Spectroscopy)
- Optical spectrum analyzers (OSAs), **V**:19.18
- Optical Stark effect, **IV**:16.13
- Optical Stiles-Crawford effect, **III**:9.2
- Optical strength, of turbulence (C_n^2), **V**:5.6–5.8, 5.7*f*, 5.8*f*
- Optical tank circuits, **V**:20.21, 20.21*f*
- Optical theorem, **I**:7.8
- Optical throughput, of NPM AOTFs, **V**:6.41
- Optical time-division multiplexed (OTDM)
 communication networks:
 and all-optical switching for demultiplexing, **V**:20.22, 20.23*f*
 device technology, **V**:20.12–20.24
 direct and indirect modulation in, **V**:20.17*f*, 20.17–20.18
 external modulation in, **V**:20.18–20.20, 20.19*f*, 20.20*f*
 history of, **V**:20.3
 multiplexing in, **V**:20.1, 20.3–20.12, 20.5*f*–20.11*f*, 20.13*f*
 and optical clock recovery, **V**:20.21*f*, 20.21–20.22
 serial vs. parallel, **V**:20.12, 20.13*f*
 transmitters in, **V**:20.12–20.17, 20.14*f*–20.16*f*
 ultrahigh-speed OTDM, **V**:20.23–20.24, 20.24*f*
 and WDM, **V**:21.2
- Optical tolerances (*see* Tolerances)
- Optical train, in microscopes, **I**:28.3–28.5, 28.5*f*
- Optical transfer function (OTF), **II**:29.17;
III:2.21–2.22, 2.22*f*
 and adaptive optics, **V**:5.19–5.20, 5.20*f*
 and atmospheric turbulence, **V**:4.3, 4.6–4.7
 calculations of, **I**:4.3, 4.5
 and camera lens performance, **I**:27.3*f*–27.5*f*, 27.7*f*–27.16*f*, 27.18*f*–27.22*f*, 27.24
- Optical transfer function (OTF) (*Cont.*):
 defined, **III**:4.1
 for human eye, **III**:1.21–1.22
 measurements of, **I**:4.6–4.7
 in ophthalmoscopic methods, **III**:1.23
 and optical quality of IOLs, **III**:21.14
 of systems with annular pupils, **V**:4.10–4.13, 4.12*f*, 4.13*f*
 with visual instruments, **III**:1.27
- Optical transmission, atmospheric,
V:3.22*f*–3.27*f*, 3.22–3.26
- Optical traps, **IV**:20.23*f*, 20.23–20.24
- Optical tube length, **I**:17.10
- Optical tweezers, **I**:28.55; **IV**:20.23
- Optical zone diameter (OZD):
 for contact lenses, **III**:20.5
 defined, **III**:20.2
- Opticaldiagnostics.com, **III**:12.17
- Optical-electrical field overlap parameter, of electro-optic effect, **V**:13.50
- Optically detected magnetic resonance (ODMR), **I**:31.21–31.23, 31.22*f*, 31.23*f*;
V:2.23–2.24, 2.24*f*
- Optically generated plasmas, **IV**:16.20
- Optically polished solid spacers, **IV**:7.88–7.89
- Optically rotated tangential phase matching,
V:6.26*f*, 6.27
- Optically-induced phase charge, **IV**:13.4*f*, 13.9
- Optimization function (of optical software),
II:3.16–3.21
 by damped least-squares method,
II:3.17–3.19
 and error functions, **II**:3.19–3.20
 global, **II**:3.21
 multiconfiguration, **II**:3.20
 by orthonormalization, **II**:3.19
 by simulated annealing, **II**:3.19
 and tolerancing, **II**:3.20–3.21
- OPTIS (simulation software), **II**:7.27
- Optoelectronic integrated circuit (OEIC),
I:21.2
- Optoelectronic integrated circuit (OEIC) chip,
II:25.15
- Optokinetic nystagmus, **III**:1.44
- Optokinetic reflex, **III**:13.20
- Optronic Laboratories, Incorporated, **II**:15.49
- Orange peel, of polymers, **IV**:3.11
- Order selecting aperture (OSA), of zone plates,
V:40.4, 40.5
- Ordered dye-doped polymers, **V**:7.34

- Ordered subsets expectation maximum (OSEM), in SPECT imaging, **V**:32.2
- Ordinary rays, **II**:3.12
- Organ of Corti, **III**:8.2, 8.24–8.26
- Organic black dye, **IV**:6.15
- Organic crystals, **IV**:12.23–12.25, 12.26*t*–12.27*t*; **V**:7.33–7.34
- Organic dye lasers, **II**:16.31–16.32, 16.32*f*
- Organic LEDs (OLEDs), **II**:40.37–40.39
- Organic matter:
 - absorption by, **IV**:1.22–1.23, 1.23*t*, 1.25–1.27, 1.25*t*, 1.26*f*
 - passive limiting in, **IV**:13.10
 - in water, **IV**:1.14
- Organic-inorganic composites, hybrid, **IV**:12.27*t*
- Organometallic vapor phase deposition (OMVPE), **II**:19.6–19.7, 19.20*t*, 19.23
- Orlando Black optical coating, **IV**:6.54, 6.55*f*
- Orlando Black surface, **IV**:6.8*f*
- Orthogonal matrices, **I**:14.11
- Orthokeratology, **III**:12.12
- Orthonormal polynomials, **II**:11.3–11.40
 - and aberration balancing, **II**:11.30, 11.35–11.36, 11.36*t*
 - about, **II**:11.4–11.5
 - circle, for noncircular pupils, **II**:11.37, 11.39
 - defined, **II**:11.5–11.6
 - discussion of, **II**:11.39–11.40
 - elliptical, **II**:11.21, 11.25–11.27, 11.26*t*–11.27*t*
 - hexagonal, **II**:11.21, 11.22*t*–11.25*t*
 - isometric, interferometric, and PSF plots for orthonormal aberrations, **II**:11.36–11.37, 11.37*f*, 11.38*f*
 - rectangular, **II**:11.27–11.28, 11.28*t*, 11.29*t*
 - slit, **II**:11.30, 11.35*t*
 - square, **II**:11.30, 11.31*t*–11.34*t*
 - Zernike annular, **II**:11.13–11.21, 11.14*f*, 11.16*f*, 11.17*t*–11.21*t*
 - Zernike circle, **II**:11.6–11.12, 11.8*t*–11.9*t*, 11.9*f*–11.11*f*, 11.12*t*
- Orthonormalization, **II**:3.19
- Orthorhombic crystals, **IV**:2.7*t*, 2.18, 2.46*t*, 8.9*t*, 8.10, 8.19*t*
- Orthoscopic imaging, **I**:28.8
- Orthotomic systems, **I**:1.10, 1.12
- Oscillation(s):
 - relaxation, **II**:16.12, 19.31*f*, 19.31–19.34
 - wavelength interval between, **I**:12.10
- Oscillator models, of optical nonlinearity, **IV**:10.5–10.9, 10.6*f*, 10.8*f*
- Oscillator strength, **V**:56.3–56.4
- Oscillators, **IV**:12.7–12.9, 12.8*f*–12.9*f* (*see also specific oscillators, e.g.*: Raman oscillators)
 - and fiber lasers, **V**:25.13, 25.14*f*, 25.30–25.33
 - high-gain, **II**:20.10–20.12, 20.11*f*
 - high-power USP, **V**:25.32–25.33
 - local, **V**:9.13
 - optical-parametric, **II**:20.20*f*, 20.20–20.22, 20.21*f*
 - parametric, **V**:11.23, 11.24
 - Q-switched, **V**:25.30–25.31
 - voltage-controlled, **V**:20.11, 20.11*f*
 - [*see also* Master oscillator power amplifier (MOPA) systems]
- Osculating planes, of space curves, **I**:1.18, 1.19
- Outer ionization, of cluster, **IV**:21.32–21.33
- Outer product processors, **I**:11.19
- Outer scale of turbulence, **V**:4.7
- Outgassing:
 - of black surfaces, **IV**:6.17
 - of polymers, **IV**:3.4
- Out-of-band radiation errors, **II**:34.36
- Out-of-plane coupling, **IV**:9.11–9.12
- Out-of-plane profile, of emitted light, **V**:13.11
- Output amplifier noise, in CCDs, **II**:32.20
- Output circuits, direct readout architectures, **II**:33.18
- Output coupling mirror, **II**:16.11
- Output gate (OG), **II**:32.14
- Output planes, conjugate matrices for, **I**:1.68
- Output windows, in proximity-focused MCP IIs, **II**:31.9, 31.9*f*
- Outside vapor deposition (OVD), **V**:25.2, 25.26
- Overall diameter (OAD):
 - for contact lenses, **III**:20.5
 - defined, **III**:20.2
- Overdense plasmas, strong field interactions
 - with, **IV**:21.46–21.52
 - high harmonic generation, **IV**:21.50–21.52, 21.51*f*
 - $\mathbf{j} \times \mathbf{B}$ heating and anomalous skin effect, **IV**:21.49
 - ponderomotive steepening and hole boring, **IV**:21.49–21.50, 21.50*f*
 - relativistic effects and induced transparency, **IV**:21.52
 - resonance absorption, **IV**:21.47*f*, 21.47–21.48

- Overdense plasmas, strong field interactions with (*Cont.*):
 structure of irradiated plasma, **IV**:21.46*f*, 21.46–21.47
 vacuum heating, **IV**:21.47*f*, 21.48–21.49
- Overhead lighting, **II**:40.12, 40.13*f*, 40.14, 40.46*f*
- Overillumination, **I**:30.14
- Overlapping integral, of acousto-optic interaction, **V**:6.15
- Overloosening and overtightening (tolerancing problems), **II**:5.11
- Overscan, in color CRTs, **III**:22.12
- Oxide layer, aluminum reflectance and, **IV**:4.44*f*
- Oxide stripe lasers, **II**:19.8, 19.9*f*
- Oxides, thin film, **I**:21.13–21.14
- Oxygen permeability (Dk), **III**:20.2
- Pacific Northwest National Laboratory program, **V**:3.26
- Package-induced laser-induced damage, **IV**:19.4–19.5
- Packages and packaging:
 HB-LED, **II**:18.5*f*, 18.5–18.6, 18.6*f*
 of photodetectors, **II**:26.9–26.10, 26.10*f*, 26.11*f*
 reliability of LED, **II**:17.25–17.26
 of SOAs, **V**:19.17, 19.17*f*
- Packet-switched networks, **V**:21.7, 21.10*f*–21.11*f*, 21.10–21.11
- Paint modeling, **II**:40.17
- Painted surfaces, **IV**:6.2*t*–6.3*t*
- Paints and surface treatments, **IV**:6.35–6.58, 6.37*f*, 6.43*f*, 6.53*f*
 Actar black coatings, **IV**:6.55
 Aeroglaze Z series, **IV**:6.36*f*, 6.37, 6.37*f*, 6.39, 6.39*f*–6.42*f*
 Akzo Nobel paints, **IV**:6.39, 6.42*f*, 6.43*f*
 anodized processes, **IV**:6.44–6.49, 6.47*f*, 6.48*f*, 6.51*f*, 6.53*f*
 black glass, **IV**:6.57
 Black Kapton, **IV**:6.57, 6.57*f*
 carbon nanotubes and nanostructured materials, **IV**:6.55, 6.59*f*
 Cardinal Black, **IV**:6.36*f*, 6.39, 6.44*f*
 Cat-a-lac Black, **IV**:6.39, 6.42*f*, 6.53*f*
 conductive/nonconductive, **IV**:6.12, 6.12*t*
 DeSoto Black, **IV**:6.37*f*, 6.39
 DURACON, **IV**:6.55–6.56
 electrically conductive black paint, **IV**:6.56
- Paints and surface treatments (*Cont.*):
 electrodeposited surfaces, **IV**:6.53–6.54, 6.54*f*, 6.55*f*
 etching of electroless nickel, **IV**:6.49–6.50, 6.50*f*, 6.51*f*, 6.53*f*
 flame-sprayed aluminum, **IV**:6.57
 Floquil, **IV**:6.44
 gold blacks, **IV**:6.57
 high-resistivity coatings, **IV**:6.56
 IBM Black (tungsten hexafluoride), **IV**:6.56
 ion beam-sputtered surfaces, **IV**:6.53
 Parson's Black, **IV**:6.44, 6.53*f*
 plasma-sprayed surfaces, **IV**:6.50–6.52, 6.51*f*–6.53*f*
 silicon carbide, **IV**:6.56
 SolarChem, **IV**:6.44, 6.48*f*, 6.53*f*
 specular metallic anodized surfaces, **IV**:6.57, 6.58*f*
 sputtered and CVD surfaces, **IV**:6.56
 3M paints and derivatives, **IV**:6.35–6.37, 6.36*f*, 6.38*f*, 6.53*f*
 ZO-MOD BLACK, **IV**:6.56
- Pancharatnam phase, **I**:32.11
- Panoramic cameras, **I**:25.26*f*, 25.25
- Panum's fusional area (PFA), **III**:12.15, 13.12, 25.5
- Parabasal optics, **I**:1.43
- Parabolic louver, **III**:23.3
- Parabolic reflectors, for neutron beams, **V**:64.3, 64.4, 64.4*f*
- Paraboloid objective, **I**:29.6
- Parallax stereogram, **I**:25.24
- Parallel beams, and x-ray tube sources, **V**:54.16
- Parallel matrix-vector multipliers, **I**:11.18*f*, 11.18–11.19
- Parallel multiplexing, **V**:20.12
- Parallel-beam scanners, **I**:30.23–30.25, 30.24*f*–30.25*f*
- Parallel-hole collimators, **V**:32.3
- Paralyzing glare, **II**:40.9
- Parametric amplification, **IV**:10.17–10.18, 16.3*t*
- Parametric amplifiers and oscillators, **V**:11.23, 11.24, 14.2, 14.10–14.11
- Parametric downconversion, **II**:23.14
- Parametric oscillators, **IV**:10.18*f*, 10.18–10.19, 10.19*f*
- Parasol ganglion cells, **III**:2.10*n*
- Paratellurite (TeO₂), **IV**:2.40*t*, 2.45*t*, 2.48*t*, 2.52*t*, 2.58*t*, 2.65*t*, 2.76*t*
- Paraxial chief rays, **I**:1.75

- Paraxial curvature, **I**:1.32–1.33
- Paraxial invariant, **I**:1.41, 1.77
- Paraxial limit, of systems of revolution, **I**:1.38
- Paraxial matrices, for geometrical optics,
I:1.65–1.74
- angle instead of reduced angle, **I**:1.72
 - arbitrary systems, **I**:1.67
 - and characteristic functions, **I**:1.74
 - conjugate matrices, **I**:1.68–1.71, 1.73
 - linearity, **I**:1.66
 - nonrotationally symmetric systems, **I**:1.74
 - operation on two rays, **I**:1.68
 - possible zeros, **I**:1.68
 - power matrix, **I**:1.67
 - skew rays, **I**:1.73
 - transfer matrices, **I**:1.66
 - two-ray specification, **I**:1.72
 - unit determinants, **I**:1.67
- Paraxial models (of human eye), **III**:1.36–1.37, 1.37*f*
- Paraxial optics (generally), **I**:1.29, 1.37
- Paraxial optics (of systems of revolution),
I:1.37–1.43
- angle of incidence at a surface, **I**:1.39
 - axial object and image locations, **I**:1.40
 - image location and magnification, **I**:1.42
 - linearity of, **I**:1.41
 - paraxial limit, **I**:1.38
 - principal focal lengths of surfaces,
I:1.39–1.40
 - ray tracing, **I**:1.40
 - reflection and refraction, **I**:1.38
 - switching axial objects and viewing
positions, **I**:1.43
 - three-ray rule, **I**:1.42
 - transfer, **I**:1.38
 - two-ray paraxial invariant, **I**:1.41
- Paraxial pupils, **I**:1.77
- Paraxial ray tracing, **II**:3.5, 3.8–3.9, 3.9*f*
- Paraxial rays, **I**:1.35; **II**:3.3
- Parseval's theorem, **V**:46.8
- Parson's Black, **IV**:6.44, 6.53*f*
- PART (stray light analysis program), **II**:7.11
- Part active, part passive athermalization,
II:8.11–8.12, 8.12*f*
- Partial coherence, in Maxwellian viewing,
III:6.12–6.14
- Partial coherence interferometry, **III**:21.7
- Partial coherence length, **V**:4.10
- Partially polarized light, **I**:15.7
- Particle hypothesis, Einstein's, **II**:23.7
- Particle pumping, **II**:16.14–16.16, 16.15*f*, 16.16*f*
- Particle-induced x-ray emission (PIXE), **V**:29.4
- Particles:
- scattering by, **I**:7.1–7.17
 - coherent vs. incoherent arrays, **I**:7.2–7.3
 - concepts of, **I**:7.4–7.5, 7.6*f*–7.7*f*, 7.8–7.10, 7.9*f*, 7.10*f*
 - isotropic homogenous spheres, **I**:7.11–7.14
 - Mie, **I**:7.11, 7.12
 - nonspherical particles, **I**:7.15–7.17
 - regular particles, **I**:7.14–7.15
 - single particles, **I**:7.2–7.3
 - theories of, **I**:7.3–7.4
 - surface coatings and generation of,
IV:6.17–6.18
 - in water
 - particle size distributions, **IV**:1.15–1.16, 1.16*f*
 - refraction index of, **IV**:1.20
 - scattering by, **IV**:1.30–1.35, 1.31*t*, 1.32*f*, 1.33*f*, 1.34*t*–1.35*t*
- Particulate matter:
- in standard atmosphere, **V**:3.6–3.7, 3.9–3.11, 3.10*f*, 3.11*f*
 - in water, **IV**:1.14–1.15
- Parts per million (ppm), **II**:17.25
- Parvocellular laminae, **III**:2.12, 2.13, 2.13*f*
- Paschen-Runge configuration, **I**:20.7, 20.11*f*, 20.14*t*
- Passband region, transmission in, **IV**:7.53, 7.54
- Passbands, **V**:20.1
- Passivation, in wafer processing, **II**:17.23
- Passive athermalization, **II**:6.22, 6.23*f*, 6.24, 8.8*f*–8.10*f*, 8.8–8.10
- Passive autofocus systems, for cameras, **I**:25.12
- Passive devices, for integrated optics,
I:21.21–21.25, 21.22*f*–21.25*f*
- Passive mode locking, **V**:20.17
- Passive modelocking, **IV**:18.8*f*, 18.8–18.9, 18.12–18.15
- Passive nonlinear optical phenomena, **IV**:5.54, 5.54*t*
- Passive optical limiting, **IV**:13.1–13.12
- active vs., **IV**:13.1–13.3, 13.2*f*, 13.3*f*
 - in materials, **IV**:13.9–13.12, 13.10*f*–13.12*f*
 - by nonlinear absorption, **IV**:13.4*f*–13.6*f*, 13.4–13.7
 - by nonlinear refraction, **IV**:13.4*f*, 13.7–13.8
 - by nonlinear scattering, **IV**:13.4*f*, 13.8

- Passive optical limiting (*Cont.*):
 optically-induced phase charge, **IV**:13.4f, 13.9
 by photorefraction, **IV**:13.8–13.9
- Path function, for gratings and monochromators, **V**:38.3f, 38.3–38.5, 38.4t–38.5t
- Pattern discrimination, **III**:2.35–2.36
- Pattern effect, of light pulses, **II**:19.32, 19.32f
- Pattern noise, **I**:26.11–26.12, 26.12f; **II**:32.12
- Pattern recognition:
 by matched filtering, **IV**:12.28–12.29, 12.29f, 12.30f
 optical matched filtering for, **I**:11.12–11.14, 11.13f
- Patterned vertical alignment (PVA) cells, **V**:8.27–8.28, 8.28f
- Paul objectives, **I**:29.28–29.29
- Pauli exclusion principle, **V**:56.2
- Pauli spin matrices, **I**:14.24, 14.41
- Pearson IV function, **V**:56.7–56.8
- Pearson VII function, **V**:56.7, 56.8
- Pechan prisms, **I**:19.3t, 19.11, 19.11f
- Pedestal contrast, **III**:2.31
- Pedestal control, in CRTs, **III**:22.8
- Pedestal effects (color vision), **III**:11.60f
 and chromatic discrimination, **III**:11.69
 in crossed conditions, **III**:11.59
 defined, **III**:11.3
 and gap effect, **III**:11.72, 11.74
 in uncrossed conditions, **III**:11.59, 11.61
- Pedestal experiments, **III**:3.2, 11.59–11.62, 11.60f
- Pedicles (cones), **III**:8.2, 8.20
- Pellin-Broca prisms, **I**:20.6f
- Penalty function method, **II**:3.18, 3.19
- Pendellösung interference, **V**:63.26
- Pendular states, **IV**:21.23
- Penetration depth, of metals, **IV**:4.47f
- Pen-Ray, **II**:15.36f
- Penta prisms, **I**:19.3t, 19.13, 19.14f
- Pentaprisms, **II**:6.14f, 6.15f, 12.12, 12.12f
- Percept rivalry suppression, **III**:13.14
- Perception:
 of color, **III**:14.17
 defined, **III**:13.2, 23.3
 of direction, **III**:13.7–13.10
 corresponding retinal points, **III**:13.8
 horopter, **III**:13.8–13.9, 13.9f, 13.10f
 vertical horopter, **III**:13.9
 of distance, **III**:13.24
- Perception (*Cont.*):
 Gestalt principles in, **III**:24.8
 of imaging artifacts in vision, **III**:24.2–24.6
 early color vision, **III**:24.5
 embedding digital watermarks, **III**:24.5
 image quality and compression, **III**:24.3–24.4
 limitations of early vision models, **III**:24.5–24.6
 rendering and halftoning, **III**:24.4
 target detection in medical images, **III**:24.5
 of lit environment, **II**:40.1–40.2, 40.4f, 40.4–40.6
 measuring (*see* Psychophysical measurement)
 of size, **III**:13.19
 of space, **III**:13.3–13.7
 binocular cues, **III**:13.7
 distortion of, **III**:13.16–13.19
 extraretinal information for eye movements, **III**:13.7
 kinetic cues, **III**:13.4–13.7, 13.5f, 13.6f
 monocular cues, **III**:13.3–13.7
 and visual effects created by artists, **III**:24.10
- Percepts:
 3-D, **III**:13.3
 defined, **III**:13.2
- Perceptual constancy, **II**:40.5
- Perceptual Subband Image Coder (PIC), **III**:24.4
- Perceptually based image compression, **III**:24.4
- Percus-Yevick approximation, **I**:9.5
- Perfect crystal interferometers, **V**:63.26–63.27, 63.27f
- Perfect diffuse reflectors (PDRs), **IV**:6.7f, 6.27f
- Perfect lens, **II**:4.4
- Perfect mirrors, **IV**:7.47
- Perfect reflecting diffusers, **II**:37.8
- Perfect transmitting diffusers, **II**:37.8
- Perfectly reflecting (PEC) surfaces, **I**:8.10
- Performance, measuring (*see* Psychophysical measurement)
- Periclase (MgO), **IV**:2.44t, 2.48t, 2.52t, 2.57t, 2.63t, 2.69t
- Perimetry, **III**:14.21
- Periodogram estimator, for 1D profiles, **V**:46.8

- Periodic multilayers:
 $[(0.5A)B(0.5A)]^N$, **IV**:7.35, 7.36*f*
of MLLs, **V**:42.5–42.6, 42.6*f*
nonabsorbing $[AB]^N$ and $[AB]NA$,
IV:7.32–7.34, 7.33*f*–7.35*f*
 $[xH.(1-x)L]^N.xH$, **IV**:7.37
- Periodically poled lithium niobate (LiNbO₃)
(PPLN), **IV**:17.1, 17.4–17.13, 17.6*f*–17.11*f*
- Periodically poled lithium tantalate (LiTaO₃)
(PPLT), **IV**:17.14–17.15, 17.15*f*
- Periodically poled potassium titanyl phosphate
(KTiOPO₄) (PPKTP), **IV**:17.2, 17.28,
17.29*f*, 17.30*f*
- Peripheral field, **III**:1.3
and retinal illuminance, **III**:1.12
TCA in, **III**:1.20
- Peripheral retina, **III**:13.3, 14.10
- Periphery cameras, **I**:25.22
- Periscopes, **II**:1.10, 1.10*f*
- Periscopic lenses, **I**:17.27, 17.27*f*, 18.19, 18.19*f*
- Permittivity, of water, **IV**:1.16
- Perovskite, **IV**:2.73*t*
- Perpendicular magnetic anisotropy, **I**:35.26
- Perpendicular-incidence ellipsometers (PIEs),
I:16.17–16.18, 16.18*f*
- Perspective:
defined, **III**:13.2
distorted, with monocular magnification,
III:13.13, 13.14*f*
- Petawatt lasers, **IV**:21.5
- PETRA III (synchrotron source), **V**:37.10
- Petzold volume scattering functions,
IV:1.33, 1.33*f*
- Petzval (field) curvature:
in gradient index optics, **I**:24.4, 24.6, 24.7
of reflective and catadioptric objectives,
I:29.7, 29.11, 29.14, 29.15, 29.32
as wavefront aberration, **I**:1.91
- Petzval lenses, **I**:17.10, 17.28, 17.28*f*, 17.35*f*
- Petzval sums, **I**:28.15, 29.37; **IV**:3.8
- Petzval surface, **II**:2.4, 2.5
- Pfund configuration, of dispersive prisms,
I:20.8*f*, 20.10, 20.13*f*
- Pfund objectives, **I**:29.6
- Phacoemulsification, **III**:12.14, 21.2,
21.4, 21.5
- Phakia, **III**:14.2
- Phakic lenses:
defined, **III**:21.2
intraocular, **III**:16.9, 21.19
- Phase aberration function:
correction of, **V**:4.28–4.30
modal expansion of, **V**:4.17*f*–4.18*f*,
4.17–4.20, 4.19*t*, 4.20*t*
- Phase array beam steering, **V**:6.27–6.29, 6.28*f*
- Phase charge, optically-induced, **IV**:13.4*f*, 13.9
- Phase coatings, **IV**:7.101, 7.101*f*–7.104*f*, 7.102
- Phase conjugate interferometry, **IV**:12.32,
12.33*f*, 12.34*f*
- Phase conjugate mirrors, **IV**:12.7, 12.8*f*–12.9*f*,
12.33–12.35
- Phase conjugated coupling, **II**:20.14
- Phase conjugation:
Brillouin, **IV**:15.48, 15.52*f*–15.54*f*,
15.52–15.54
photo echo and geometry of, **IV**:11.19–11.22,
11.20*f*, 11.21*f*
- Phase contrast microscopy, **I**:28.28–28.29,
28.29*f*
- Phase diffusion coefficient, **II**:23.29
- Phase discrimination, **III**:2.36
- Phase dispersion filters, **IV**:7.89, 7.89*f*
- Phase distortion, nonlinearly induced,
IV:16.28–16.29
- Phase errors, phase-shifting interferometry
and, **II**:13.22
- Phase fluctuations, adaptive optics and,
V:5.18–5.19
- Phase mask method, of FBG fabrication,
V:17.5–17.6, 17.6*f*, 17.8, 24.7, 24.7*f*
- Phase matching:
for acousto-optic devices, **V**:6.9–6.12, 6.10*f*
in attosecond optics, **II**:21.4
birefringent tangential, **V**:6.25
and harmonic yield, **IV**:21.30
in integrated optics, **I**:21.12
and Maxwell-Bloch equations, **IV**:11.6
noncritical, **IV**:17.1
on-axis tangential, **V**:6.25–6.26
optically rotated tangential, **V**:6.26*f*, 6.27
QPM materials, **IV**:17.1, 17.13–17.14
in second-order processes, **IV**:10.12*f*,
10.12–10.14, 10.13*f*
and stimulated Raman scattering, **IV**:15.7,
15.34, 15.34*f*
tangential, **V**:6.12, 6.13, 6.17, 6.25–6.27, 6.26*f*
- Phase modulation, **V**:7.18–7.20, 13.51 (*see*
also Cross-phase modulation; Self-phase
modulation)
- Phase modulation index, **II**:22.4; **V**:7.19

- Phase noise, **V**:9.8, 13.1, 13.20–13.21, 13.21*f*
Phase plates, **I**:28.28
Phase plates, for circular polarization, **V**:43.5*f*, 43.5–43.6
Phase pulling, of transient Raman scattering, **IV**:15.26–15.27, 15.27*f*
Phase response, frequency vs., **II**:22.6*f*, 22.6–22.7, 22.7*f*
Phase retardation, **I**:12.24; **III**:18.2
Phase retarders, **V**:41.9, 43.6
Phase retarding reflectors, **IV**:7.101–7.102, 7.102*f*, 7.103*f*
Phase shifter, cross-Kerr, **IV**:23.11
Phase space, **II**:39.3
Phase stability, Doppler OCT and, **III**:18.18*f*–18.20*f*, 18.18–18.19
Phase stability margin, **II**:22.9
Phase structure function, **V**:4.5, 5.9, 5.10
Phase transfer function (PTF):
 in diffraction-limited eye, **III**:1.13
 and off-axis image quality, **III**:1.26, 1.27
 in ophthalmoscopic methods, **III**:1.23
Phase transfer functions (PTF), **I**:4.3, 4.7
Phase transitions, of liquid crystals, **V**:8.13–8.14, 8.13*f*, 8.14*f*
Phase velocity indices of refraction, **V**:7.15–7.16, 7.16*f*
Phase zone plates, **V**:40.5–40.7, 40.7*f*
Phase-conjugate interferometers, **I**:32.17, 32.18, 32.18*f*
Phased-arrays, for agile beam steering, **I**:30.52–30.57, 30.53*f*, 30.62–30.63
Phase-interruption broadening, emission-line, **II**:16.5
Phase-lock phase shifting, **II**:13.23
Phase-locked interferometers, **I**:32.11–32.12, 32.12*f*
Phase-locked laser arrays, **II**:19.26, 19.27, 19.28*f*, 19.28*t*
Phase-locked loops (PLLs), **V**:20.11, 20.11*f*
Phase-measuring interferometers (PMIs), **V**:46.2
Phase-sensitive detection (PSD), **I**:31.8*f*, 31.8–31.11
Phase-shifting interferometers, **I**:32.10–32.11, 32.11*f*
Phase-shifting interferometry, **II**:13.18*f*–13.20*f*, 13.18–13.23
 heterodyne interferometer, **II**:13.22
 integrating bucket method, **II**:13.21, 13.21*f*
 phase errors, **II**:13.22
Phase-shifting interferometry (*Cont.*):
 phase-lock method, **II**:13.23, 13.23*f*
 phase-stepping method, **II**:13.20, 13.20*f*
 simultaneous measurement, **II**:13.22
 two-steps-plus-one method, **II**:13.21, 13.22
Phasesonium, **IV**:14.3
Phase-space acceptance, by gratings and monochromators, **V**:38.7
Phase-stepping phase shifting, **II**:13.20, 13.20*f*
Phase-transition temperatures, of crystals and glasses, **IV**:2.32, 2.33
Phasors, **I**:5.2
Phonon broadening, emission-line, **II**:16.5
Phonon coupling, **II**:17.16
Phonons, **IV**:2.11
 acoustic, **IV**:5.14–5.16
 coupled plasmon and, **IV**:5.35, 5.36, 5.36*f*, 5.37*f*
 lattice absorption by, **IV**:5.13–5.16, 5.15*f*
 optical, **IV**:5.14
 and Raman scattering, **IV**:5.79*f*–5.81*f*, 5.79–5.81
 transverse optical (TO), **IV**:8.16–8.18
Phoria, **III**:13.21–13.22
 defined, **III**:13.2
 effect of lenses and prisms on, **III**:13.25–13.27
Phoropter, **III**:12.6
Phosphate crown glass, **IV**:2.41*t*
Phosphates, for fiber lasers, **V**:25.27*t*, 25.28
Phosphor, **III**:23.3
Phosphor salts, **II**:40.31
Phosphor screens:
 of image intensifiers, **II**:31.14–31.16, 31.14*t*, 31.15*f*
 in proximity-focused MCP IIs, **II**:31.9, 31.9*f*
Phosphor x-ray detectors, **V**:60.7–60.8, 60.10*t*
Phosphor-type designation system, **II**:31.14*t*
Phot (unit), **II**:34.43, 36.7, 36.7*t*
Photo cell (*see* Photodiodes)
Photo gain, in photoconductors, **II**:25.5*f*, 25.5–25.6
Photoabsorption, **V**:36.1
Photoacoustic spectroscopy (PAS), **IV**:17.22*f*–17.24*f*, 17.22–17.26, 17.26*f*
Photo-activated localization microscopy (PALM), **I**:28.23
Photo-associative spectroscopy (PAS), **IV**:20.30–20.31, 20.31*f*
Photocapacitors, MOS (*see* Metal-oxide-semiconductor capacitors)

Photodetectors (*Cont.*):

- Si photovoltaic detectors, **II**:24.52*f*,
24.54–24.65, 24.55*f*–24.66*f*
- Si:B detectors, **II**:24.95*f*, 24.96
- SiC UV detectors, **II**:24.47, 24.47*f*
- Si:Ga infrared detectors, **II**:24.95, 24.95*f*,
24.96, 24.96*f*
- and signal detection, **II**:27.2
- spectral response of, **II**:24.18, 24.19*f*
- speed of, **II**:24.20, 24.21
- stability of, **II**:24.21, 24.21*f*
- terminology, **II**:24.10–24.13
- thermal detectors, **II**:24.4*f*–24.6*f*, 24.4–24.6
- thermistor bolometers, **II**:24.24*f*, 2
4.24–24.25, 24.25*f*, 28.7*t*
- thermocouples, **II**:24.22*f*, 24.22–24.23
- thermopiles, **II**:24.23*f*, 24.23–24.24
- TiO₂ UV detectors, **II**:24.47, 24.48*f*
- types of, **II**:25.3, 25.4*f*
- uniformity of, **II**:24.20
- (*see also specific photodetectors, e.g.:*
Avalanche photodetectors)
- Photodiode CCDs (PD-CCDs), **II**:33.11*f*, 33.12
- Photodiode linear arrays (PDAs), **II**:38.9, 38.10,
38.10*t*
- Photodiode MOSs (PD-MOSs),
II:33.11*f*, 33.13
- Photodiodes (PDs):
- avalanche, **V**:13.63, 13.71–13.73
- CCD, **II**:33.11*f*, 33.12
- defined, **II**:24.11
- electronics of, **II**:38.10
- GaAsP, **II**:24.49, 24.49*f*, 24.50*f*
- GaP, **II**:24.47–24.49, 24.48*f*
- Ge avalanche, **II**:24.70*f*, 24.72*f*, 24.72–24.73,
24.73*f*
- GeGaAs, **II**:34.31
- InGaAs, **II**:34.31
- InGaAs avalanche, **II**:24.66*f*–24.69*f*,
24.66–24.70
- junction, **II**:32.3–32.6, 32.4*f*, 32.6*f*
- MOS, **II**:33.11*f*, 33.13
- pin, **V**:13.64–13.66
- pin* (*see pin photodiodes*)
- p⁺np*, **II**:32.8
- resonant, **V**:13.65
- Schottky, **V**:13.63, 13.73
- silicon, **II**:34.29, 34.30
- silicon avalanche, **II**:24.62–24.65,
24.63*f*–24.66*f*

Photodiodes (PDs) (*Cont.*):

- silicon *pn*, **II**:24.52*f*, 24.55*f*–24.59*f*,
24.55–24.58
- unitraveling-carrier, **V**:13.68*f*, 13.68–13.69
- UV-enhanced, **II**:24.55*f*, 24.61*f*, 24.61–24.62,
24.62*f*
- (*see also specific photodiodes, e.g.:* Avalanche
photodiodes)
- Photoelastic coefficients, of crystals and glasses,
IV:2.24
- Photo-elastic modulators (PEMs), **I**:15.21,
16.13
- Photoelectromagnetic (PEM) detectors, **II**:24.9,
24.9*f*
- Photoemissive detectors, **II**:24.6, 24.7*f*
- gallium phosphide dynodes, **II**:24.42, 24.44*f*
- linearity of, **II**:24.41
- manufacturers of, **II**:24.42
- noise from, **II**:24.39, 24.40, 24.41*f*
- operating temperature of, **II**:24.40
- photon counting for, **II**:24.42, 24.43*f*, 24.44*f*
- quantum efficiency of, **II**:24.35*f*–24.38*f*,
24.36–24.38
- recommended circuit for, **II**:24.42, 24.43*f*
- response time of, **II**:24.40
- responsivity of, **II**:24.35*f*, 24.38
- sensitive area of, **II**:24.41
- sensitivity of, **II**:24.34, 24.35*f*–24.39*f*
- sensitivity profile of, **II**:24.41
- short-wavelength considerations for,
II:24.34, 24.40*f*
- specifications for, **II**:24.32–24.42, 24.33*f*
- stability of, **II**:24.41
- Photoexcitation, **II**:25.5; **IV**:5.70*f*
- Photogates, **II**:25.10, 25.11, 25.11*f*
- Photographic detectors, **II**:24.9–24.10
- Photographic dyes, **II**:30.10*f*, 30.10–30.13,
30.12*f*
- Photographic emulsions, **II**:24.100, 24.101*f*,
29.4, 30.7
- Photographic film(s), **II**:29.3–29.16, 30.18–30.28
- about, **II**:30.2–30.3
- black-and-white (B&W) film,
II:30.24–30.25, 30.25*t*
- color, **II**:29.12–29.15, 29.13*f*, 29.14*f*
- color negative film, **II**:30.25–30.28, 30.27*t*
- color reversal film, **II**:30.22–30.24, 30.23*t*
- development effects on, **II**:29.12, 29.13*f*
- D-log H curve for, **II**:29.8*f*, 29.8–29.10, 29.10*f*
- exposure of, **II**:29.5–29.6

- Photographic film(s) (*Cont.*):
 grain element of, **II**:29.5
 granularity of, **II**:30.19
 high-speed vs. low-speed, **II**:30.18–30.20
 and light scattering by silver halide crystals,
II:30.5–30.7, 30.6*f*
 microdensitometers, **II**:29.15*f*, 29.15–29.16
 neutron detection with, **V**:63.33, 63.34
 optical density of, **II**:29.6–29.8, 29.7*f*
 processing of, **II**:29.5
 professional vs. amateur film, **II**:30.20–30.22
 reciprocity failure of, **II**:29.11–29.12
 spectral sensitivity of, **II**:29.11, 29.11*f*
 speed of, **II**:30.18–30.20
 structure of color, **II**:30.3*f*, 30.3–30.5
 structure of silver halide photographic layers
 in, **II**:29.4
- Photographic film speed, **II**:29.9–29.10
 in Advanced Photo System, **II**:30.26
 in color negative film, **II**:30.25
 and granularity, **II**:30.19
 high vs. low, **II**:30.18–30.20
 and sensitivity to high-energy radiation,
II:30.19–30.20
- Photographic materials, **II**:30.1–30.28
 about, **II**:30.1–30.2
 dyes, **II**:30.10–30.13
 about, **II**:30.10, 30.10*f*
 excited state properties, **II**:30.11–30.12,
 30.12*f*
 light stabilization methods, **II**:30.12–30.13
 photochemistry of azomethine dyes,
II:30.11
- films, **II**:30.18–30.28
 black-and-white film, **II**:30.24–30.25,
 30.25*t*
 color negative film, **II**:30.25–30.28, 30.27*t*
 color reversal film, **II**:30.22–30.24, 30.23*t*
 professional vs. amateur film,
II:30.20–30.22
 speed, **II**:30.18–30.20
- optics of, **II**:30.2–30.7
 about, **II**:30.2–30.3
 light scatter by silver halide crystals,
II:30.5–30.7, 30.6*f*
 structure of color films, **II**:30.3*f*, 30.3–30.5
 structure of color papers, **II**:30.5
 and photographic spectral sensitizers,
II:30.13–30.18
 about, **II**:30.13–30.14, 30.14*f*
- Photographic materials, and photographic
 spectral sensitizers (*Cont.*):
 color science, **II**:30.15–30.18, 30.16*f*, 30.17*f*
 photophysics of spectral sensitizers on
 silver halide surfaces, **II**:30.14–30.15,
 30.15*f*
 silver halide light detectors, **II**:30.7–30.9, 30.8*f*
- Photographic papers:
 about, **II**:30.2–30.3
 and light scattering by silver halide crystals,
II:30.5–30.7, 30.6*f*
 structure of color, **II**:30.5
- Photographic plates, **I**:3.3
- Photographic recording, in ophthalmoscopic
 methods, **III**:1.23
- Photographic spectral sensitizers (*see* Spectral
 sensitizers, photographic)
- Photographic systems, **II**:29.16–29.25
 acutance of, **II**:29.17–29.19, 29.18*t*, 29.19*f*
 detective quantum efficiency of, **II**:29.23
 graininess in, **II**:29.19–29.22, 29.21*f*
 image structure of, **II**:29.17
 information capacity of, **II**:29.24
 manufacturers of, **II**:29.25
 performance of, **II**:29.16–29.17, 29.17*f*, 29.18*t*
 resolving power of, **II**:29.24
 sharpness in, **II**:29.22
 signal-to-noise ratio of, **II**:29.22–29.23
- Photography, wide-angle, **I**:25.25
- Photoionization detectors, **II**:24.10
- Photoionization devices, **II**:34.29
- Photoionization yield, **II**:34.29
- Photokeratitis, **III**:1.9, 7.4*f*
 defined, **III**:7.1
 from radiation, **III**:7.4
 UV, **III**:7.3
- Photoluminescence (PL), **IV**:5.70*f*, 5.70–5.75,
 5.72*f*–5.75*f*
- Photoluminescence decay time, **I**:31.12–31.15,
 31.13*f*, 31.14*f*
- Photoluminescence excitation (PLE)
 spectroscopy, **IV**:5.75, 5.75*f*
- Photometric efficiency (PE) factor, **III**:9.3
- Photometric ellipsometers, **I**:16.12–16.14,
 16.13*f*, 16.14*f*
- Photometry, **II**:34.37–34.44, 36.1–36.17, 36.3*f*,
 36.3*t*; **III**:9.1–9.15, 10.10, 10.43–10.45,
 11.37
 about, **II**:36.2–36.4
 approximations, **II**:36.10, 36.10*f*, 36.11*f*

- Photometry (*Cont.*):
 basis of physical, **II**:37.1–37.2, 37.2*f*
 calibrations in, **II**:34.42–34.43
 concepts/terminology of, **II**:34.10–34.11, 34.38*t*, 34.39–34.40, 34.43–34.44, 39.2, 39.2*t*
 conversion between radiometric and photometric quantities, **II**:34.12*t*, 36.11–36.12, 36.12*f*, 36.12–36.14, 36.13*f*, 36.14*f*
 defined, **II**:34.37, 37.1
 and human eye, **II**:36.8*f*, 36.8–36.10, 36.9*f*
 illuminance-luminance relationship, **II**:37.9, 37.9*f*
 integrating sphere device, **II**:37.9–37.10, 37.10*f*
 inverse square law, **II**:37.8
 Lambert's cosine law, **II**:37.8, 37.8*f*
 normalization, **II**:36.14–36.17, 36.15*t*, 36.16*f*
 and photopic/scotopic/mesopic vision, **II**:34.37–34.39, 34.38*t*
 practice in, **II**:37.11
 quantities and units in, **II**:37.4–37.8, 37.4*t*, 37.5*f*, 37.7*t*
 radiometry vs., **II**:34.6
 retinal illuminance, **II**:34.40–34.42
 symbols/nomenclature of, **II**:36.5–36.10, 36.6*t*–36.8*t*
 weighting functions, **II**:36.17
- Photomicrographic lamps, **II**:15.47–15.49, 15.48*f*, 15.49*f*
- Photomultiplier tubes (PMTs), **II**:24.11, 24.32–24.34, 24.33*f*, 24.38*f*, 24.38–24.42, 24.39*f*, 24.43*f*, 38.9, 38.9*t*; **IV**:5.61; **V**:31.5
 applications for, **II**:27.6–27.10, 27.7*f*
 base design of, **II**:27.8–27.10
 electronics of, **II**:38.10
- Photon(s), **II**:23.6–23.14, 25.2, 34.30, 36.4
 Einstein's light quanta, **II**:23.6–23.9, 23.8*f*
 photon-photon correlations, **II**:23.13*f*, 23.13–23.14
 quantum electrodynamics, **II**:23.9–23.13
- Photon absorption (PA), **IV**:19.9, 19.10, 19.10*f*
- Photon correlation spectroscopy (PCS), **I**:9.8
- Photon counting, **II**:27.15
 defined, **II**:24.12
 of modulated signal sources, **II**:27.14
 in photoemissive detectors, **II**:24.42, 24.43*f*, 24.44*f*
- Photon density, **II**:19.30–19.35
- Photon detectors:
 background-limited case of, **II**:24.14–24.17, 24.16*f*, 24.16*t*, 24.17*f*
 strong-signal case of, **II**:24.14
- Photon dose, **II**:34.6, 34.11
- Photon echo, **IV**:11.11–11.19
 about, **IV**:11.11–11.15, 11.12*f*, 11.13*f*, 11.15*f*
 stimulated, **IV**:11.15–11.19, 11.16*f*–11.19*f*
- Photon echo spectroscopy, **IV**:11.26–11.27
- Photon Engineering, **II**:7.25
- Photon flux, **II**:34.5, 34.11
- Photon infrared detectors, **II**:33.7, 33.8*f*
- Photon loss, **IV**:23.14–23.15
- Photon migration approach, to radiative transfer, **I**:9.12
- Photon noise, **III**:2.4
- Photon transfer, **I**:26.12
- Photon-flux irradiance (retina), **III**:2.4, 2.7
- Photonic bandgaps (PBGs), **IV**:9.1–9.17; **V**:11.2–11.3, 11.8, 11.8*f*, 11.10*f*, 11.11, 11.11*f*, 11.14
 fibers with, **IV**:2.23
- Maxwell's equations, **IV**:9.2–9.3
 in 3D photonic crystals, **IV**:9.4–9.12
 criteria for, **IV**:9.4–9.5
 examples of, **IV**:9.5, 9.5*f*
 microcavities in, **IV**:9.6–9.12, 9.7*f*, 9.8*f*
 2D periodicity in microcavities of, **IV**:9.8–9.12, 9.9*f*
 in-plane coupling, **IV**:9.10–9.11
 out-of-plane coupling, **IV**:9.11–9.12
 and waveguides, **IV**:9.12–9.17
 in photonic crystals with 2D periodicity, **IV**:9.13*f*, 9.13–9.14
 waveguide bends, **IV**:9.14–9.16, 9.15*f*, 9.16*f*
 waveguide intersections, **IV**:9.16–9.17, 9.17*f*
- Photonic crystal fibers (PCFs), **V**:11.1–11.28, 25.2
 all solid-core, **V**:25.21
 Bragg fibers, **V**:11.4
 in cladding, **V**:11.7–11.11, 11.8*f*–11.11*f*
 cleaving and splicing of, **V**:11.26
 design and fabrication of, **V**:11.4–11.6, 11.5*f*
 endlessly single-mode, **V**:11.12, 11.13, 11.21, 11.21*f*
 in fiber lasers, **V**:25.19*f*, 25.20–25.21, 25.27
 and guidance
 attenuation mechanisms, **V**:11.19–11.22, 11.20*f*, 11.20–11.22, 11.21*f*
 birefringence, **V**:11.17

- Photonic crystal fibers (PCFs), and guidance
(*Cont.*):
core-cladding index difference,
V:11.12*f*–11.17*f*, 11.12–11.17
group velocity dispersion, V:11.18, 11.18*f*,
11.19*f*
Kerr nonlinearities, V:11.22–11.24,
11.24*f*
resonance and antiresonance, V:11.12
scattering, V:11.24–11.26, 11.25*f*
history of, V:11.2–11.4, 11.3*f*
in-fiber devices for, V:11.27, 11.28
mode transformers, V:11.26–11.27, 11.27*f*
modeling and analysis of, V:11.6–11.7
- Photonic crystals, IV:9.3
- Photonic excitation, II:25.3, 25.3*f*
- Photonic integrated circuits (PICs),
V:19.1, 19.36
of III-V materials, I:21.17–21.20
in integrated optics, I:21.2
in WDM systems, I:21.37, 21.38
- Photonic resonance, IV:22.9–22.13, 22.10*f*,
22.12*f*
- Photonics Directory of Optical Industries*,
II:15.14
- Photon-number eigenstates, IV:23.7–23.8
- Photons:
energies of, V:36.7*t*–36.8*t*
lifetimes of, V:20.1
mass attenuation coefficients for, IV:4.48,
4.48*t*, 4.49
Raman interaction with, IV:15.21–15.22
and water, IV:1.11
- Photophobia, III:23.3
- Photopic luminosity function, III:10.1, 10.4*t*
- Photopic luminous efficiency function,
III:10.44
- Photopic retinal illuminance, III:2.4
- Photopic troland, III:9.2
- Photopic vision, II:34.37–34.39, 34.38*t*,
36.8*f*–36.10*f*, 37.2, 40.3
accommodation response, III:1.32–1.34,
1.33*f*, 1.34*f*
age-related changes in, III:14.15, 14.16*f*
color, III:10.3
Stiles-Crawford effect, III:1.10, 1.10*f*
- Photopigment opsin genes, III:10.14
in color-deficient observers, III:10.16
and variations in color matching,
III:10.15
- Photopigment optical density:
and adjustment of cone spectral sensitivities,
III:10.17
adjustments for individual differences,
III:10.18
and color matches, III:10.9
- Photopigments:
absorption spectra of, III:2.7
adjustments for individual differences,
III:10.18
and color matches, III:10.9
in color-deficient observers, III:10.16
time constant of regeneration, III:2.7
variability in, III:10.15
- Photopolarimeters, I:16.13–16.16, 16.14*f*–16.16*f*
- Photoreactive scattering, IV:12.9
- Photoreceptors, III:10.3, 10.4
biological waveguide models of, III:8.8–8.9,
8.9*f*, 8.10*f*, 8.12–8.15
and color appearance, III:11.62
and contrast detection, III:2.22–2.23
defined, III:15.2
directional sensitivity of, III:8.5
dynamic range of, III:2.9
image sampling by, III:2.4–2.9
inhomogeneous, III:10.16
length of outer segments, III:2.7–2.8
modal patterns in, III:8.16*f*, 8.18*f*, 8.19–8.24,
8.22*f*
for neighboring waveguides, III:8.22
radial distribution of transmitted energy,
III:8.22–8.23, 8.23*f*
Snyder and Pask cone model, III:8.21*f*
transfer function, III:8.24*f*
optical waveguide properties of, III:14.11
optics of, III:8.3, 8.10*f*
orientation and alignment of, III:8.5–8.8,
8.6*f*, 8.7*f*
quantitative observations of, III:8.15, 8.16*f*,
8.17, 8.18*f*, 8.19*f*
and resolving capacity of the eye,
III:4.5, 4.6*f*
schematic diagram of, III:8.9*f*
temporal response properties of, III:2.8–2.9
types and functions of, III:2.4, 2.6
waveguiding in, III:8.3–8.5
in xerographic systems, I:34.1, 34.2*f*,
34.2–34.4, 34.3*f*
(*see also* Cones; Rods)
- Photoreflectance (PR), IV:5.66*t*, 5.67

- Photorefraction:
 optical limiting by, **IV**:13.8–13.9
 of third-order optical nonlinearities,
IV:16.22
- Photorefractive effect, **IV**:12.1–12.38
 beam spatial profiles, **IV**:12.10
 devices using, **IV**:12.28–12.38
 associative memories and neural networks,
IV:12.33–12.35
 gain and two-beam coupling,
IV:12.29–12.32, 12.31*f*
 holographic data storage, **IV**:12.37
 holographic storage, **IV**:12.36–12.37
 loss and two-beam coupling,
IV:12.31–12.32, 12.33*f*–12.35*f*
 phase conjugate interferometry, **IV**:12.32,
 12.33*f*, 12.34*f*
 real-time holography, **IV**:12.28–12.29,
 12.29*f*, 12.30*f*
 solitons, **IV**:12.38
 thresholding, **IV**:12.35–12.36, 12.36*f*
 waveguides, **IV**:12.37
 grating formation, **IV**:12.1–12.3, 12.2*f*
 oscillators and self-pumped mirrors,
IV:12.7–12.9, 12.8*f*–12.9*f*
 standard rate equation model, **IV**:12.3–12.4
 stimulated photoreactive scattering,
IV:12.9
 time-dependent effects, **IV**:12.10
 wave interactions, **IV**:12.4–12.7
 anisotropic scattering, **IV**:12.7
 four-wave mixing, **IV**:12.6*f*, 12.6–12.7
 two-beam coupling, **IV**:12.4*f*, 12.4–12.6
- Photorefractive keratectomy (PRK), **III**:12.14,
 16.11*f*, 16.11–16.12, 16.12*f*
- Photorefractive materials, **IV**:12.10–12.25
 bulk compound semiconductors,
IV:12.20–12.21, 12.20*t*, 12.21*f*
 comparison of, **IV**:12.13, 12.13*t*
 cubic oxides (sillenites), **IV**:12.17–12.19,
 12.18*t*, 12.19*f*
 features of, **IV**:12.10–12.11
 ferroelectric, **IV**:12.13–12.17, 12.13*t*
 barium titanate, **IV**:12.15–12.16, 12.16*f*
 lithium niobate and lithium tantalate,
IV:12.14
 potassium niobate, **IV**:12.16–12.17
 strontium barium niobate and related
 compounds, **IV**:12.17
 tin hypothydiphosphate, **IV**:12.17, 12.18*t*
- Photorefractive materials (*Cont.*):
 figures of merit for, **IV**:12.11–12.13
 steady-state performance, **IV**:12.11–12.12
 transient performance, **IV**:12.12, 12.13
 multiple quantum wells, **IV**:12.22*f*,
 12.22–12.23
 organic crystals and polymer films,
IV:12.23–12.25, 12.26*t*–12.27*t*
 passive limiting in, **IV**:13.11, 13.12*f*
- Photorefractive self-oscillation, **IV**:12.7, 12.9
- Photorefractivity, of lithium niobate
 modulators, **V**:13.54
- Photorelaxation, **IV**:5.70*f*
- Photoresist, for extreme ultraviolet lithography,
V:34.2, 34.6
- Photoresponse nonuniformity (PRNU),
I:26.11–26.12, 26.12*f*
- Photoretinitis, **III**:7.1, 7.7
- Photosensitive compounds, in microscopy,
I:28.55
- Photosensitivity, of fiber Bragg gratings,
V:17.2–17.3
- Photostimulable phosphors,
V:63.34
- Photosynthetically available radiation (PAR),
IV:1.5*t*, 1.7*f*, 1.9
- Phototropism, **III**:8.2
- Photovoltage, **IV**:5.66*t*
- Photovoltaic (PV) arrays, **II**:33.4
- Photovoltaic detectors, **II**:24.8, 24.8*f*, 24.9
 aluminum gallium nitride, **II**:24.46
 gallium nitride, **II**:24.42, 24.43, 24.45*f*, 24.46,
 24.46*f*, 24.47
 indium antimonide, **II**:24.80–24.83, 24.82*f*,
 24.83*f*
 indium arsenide, **II**:24.75, 24.77*f*–24.79*f*,
 24.77–24.78
 lead tin telluride, **II**:24.92, 24.93*f*
 mercury cadmium telluride, **II**:24.86*f*, 24.88*f*,
 24.90–24.92, 24.91*f*, 24.92*f*
 silicon, **II**:24.52*f*, 24.54–24.65,
 24.55*f*–24.66*f*
- Photovoltaic Schottky barrier detectors (SBDs),
II:33.7, 33.8*f*
- Physical constants, for crystals and glasses,
IV:2.8*t*
- Physical layer (PHY) implementation, of FDDI
 connectors, **V**:23.2–23.3
- Physical optics, **II**:3.16
- Physical photometry, **II**:34.37, 36.4, 37.2

- Physical properties:
of crystals and glasses, **IV**:2.37, 2.38*t*–2.43*t*
classes and symmetry properties, **IV**:2.7*t*
composition, structure, and density,
IV:2.38*t*–2.41*t*
physical constants, **IV**:2.8*t*
specialty, and substrate materials, **IV**:2.43*t*
symmetry properties, **IV**:2.5, 2.6*t*–2.8*t*
of metals, **IV**:4.6, 4.49, 4.52*t*–4.54*t*
of polymers, **IV**:3.2–3.5, 3.4*t*
- Physical vapor deposition (PVD), **V**:61.6
- Physically realizable Mueller matrices,
I:14.40–14.42
- Physiological optics, **III**:2.2
- Phytoplankton:
absorption by, **IV**:1.23–1.25, 1.24*f*–1.25*f*, 1.28
in water, **IV**:1.14
- Pickups, **II**:3.6
- Picosecond (unit), **IV**:5.7; **V**:20.1
- Picosecond and sub-picosecond relaxation,
I:31.14, 31.14*f*
- Picosecond solid-state lasers, **IV**:18.10
- Pictures, instant, **I**:25.8
- Piezoelectric transducer (PZT), **I**:15.21, 15.24
amplifier strategies for, **II**:22.18
disk vs. tube, **II**:22.17–22.18
- Piezoelectric-based (PZT-based) systems,
II:22.1
- Piezo-optic coefficients, for crystals and glasses,
IV:2.21
- Pig lenses, **III**:19.11, 19.12*f*
- Pigments:
macular, **III**:1.9, 1.11, 14.9–14.10
of the retina, **III**:14.9–14.11
- Pigtail connection, fiber, **V**:13.8
- Pikhtin-Yas'kov formula, **IV**:2.22
- Pile-of-plates polarizers:
nonnormal-incidence reflection,
I:12.15–12.18, 12.16*f*, 12.17*f*
nonnormal-incidence transmission,
I:12.18–12.24
- Pin diodes, **V**:13.2, 13.63–13.71, 13.66*t*
dark current, **V**:13.69
geometry of, **V**:13.64*f*, 13.64–13.65
noise, **V**:13.70–13.71
sensitivity, **V**:13.65–13.66, 13.66*f*
speed, **V**:13.67–13.68
unitraveling-carrier (UTC) photodiodes,
V:13.68*f*, 13.68–13.69
- Pin holes, in neutron and x-ray optics, **V**:26.7
- pin* junctions, **II**:26.3; **V**:13.59
- pin* photodetectors:
biased, **II**:26.6*f*
high-speed, **II**:26.10, 26.12*f*, 26.12–26.15,
26.14*f*, 26.15*f*
lateral, **II**:25.15*f*, 25.15–25.16
- pin* photodiodes, **II**:24.54, 24.55*f*–24.60*f*,
24.58–24.61, 25.4*f*, 25.6–25.10, 25.7*f*,
32.4*f*; **V**:13.64–13.66
absorption coefficient of, **II**:25.8, 25.9*f*
avalanche photodiodes, **II**:25.8–25.10, 25.9*f*
dark current in, **II**:25.7–25.8
diffusion current of, **II**:25.7
equivalent circuit of, **II**:26.7, 26.7*f*
generation-recombination current of,
II:25.7–25.8
germanium, **II**:24.70*f*–24.72*f*, 24.70–24.71
InGaAs, **II**:24.66, 24.67*f*
operating principles of, **II**:25.6–25.10,
25.7*f*, 25.9*f*
quantum efficiency of, **II**:25.8
resonant, **II**:26.15, 26.15*f*
responsivity of, **II**:25.8
silicon, **II**:24.55*f*–24.57*f*, 24.58–24.61, 24.59*f*,
24.60*f*
tunneling current of, **II**:25.8
vertically-illuminated, **II**:26.3, 26.4*f*, 26.5,
26.10, 26.12*f*, 26.12–26.13
waveguide, **II**:26.13–26.14, 26.14*f*
- Pin waveguides, **V**:13.68
- Pinch plasma, **V**:57.1–57.5, 57.2*f*, 57.3*t*, 57.4*f*
- Pincushion distortion, **I**:1.91
- Pinhole apertures, for SPECT imaging, **V**:32.3
- Piston error, **I**:1.91
- Pitch, of liquid crystals, **V**:8.10
- Pitch-based edging (fabrication step), **II**:9.6
- Pixel summation, **II**:38.10
- Pixels, **I**:30.8; **III**:23.3
- Planar buried heterostructure (PBH) lasers,
II:19.24, 19.25*f*, 19.34*f*; **V**:13.6
- Planar diffused silicon photodiodes, **II**:24.56*f*
- Planar lenses, **I**:17.28
- Planar objects, transmissive, **I**:6.3
- Planar photodetectors, **II**:25.14, 25.15*f*
- Planar secondary source of light, **I**:5.9
- Planckian radiation (*see* Blackbody radiation)
- Planck's law, **II**:23.6, 23.7, 34.23, 34.24,
37.10–37.11; **V**:3.18, 3.20
- Plane diffusers, **II**:38.14, 38.14*f*, 38.15*f*
- Plane mirrors, **I**:1.25

- Plane of incidence, **I**:1.23, 12.6
 Plane of polarization, **I**:12.6*n*
 Plane wave analysis, for acousto-optic interaction, **V**:6.6–6.9
 Plane waves, **I**:2.4, 2.5*f*, 3.3, 3.17
 decomposition of, **I**:3.23
 interference of, **I**:2.8–2.9, 2.9*f*
 and spherical waves, **I**:2.9–2.11, 2.10*f*
 Plane-parallel plates, **I**:2.19, 2.20*f*, 2.30*f*, 2.30–2.33, 2.32*f*, 2.33*f*
 Planes of vibration, **I**:12.6
 Plano optics fabrication, **II**:9.7
 Plant tissues, fiber-optic, **III**:8.26–8.28
 Plasma(s), **IV**:21.3–21.4
 and atomic spectroscopy, **V**:2.4–2.5
 for extreme ultraviolet lithography, **V**:34.5
 in fast ignition, **IV**:21.54*f*, 21.54–21.55, 21.55*f*
 in femtosecond x-ray production, **IV**:21.52–21.53, 21.53*f*
 in fusion neutron production, **IV**:21.53
 in high magnetic field production, **IV**:21.53
 laser-generated, **V**:56.1–56.10
 Bremsstrahlung, **V**:56.8–56.10
 and characteristic radiation, **V**:56.2–56.10
 recombination radiation, **V**:56.10
 spectral line broadening, **V**:56.2–56.8
 in MeV proton acceleration, **IV**:21.54
 nano-, **IV**:21.34–21.35, 21.35*f*
 optically generated, **IV**:16.20
 overdense, **IV**:21.46–21.52
 high harmonic generation, **IV**:21.50–21.52, 21.51*f*
 irradiated plasma, **IV**:21.46*f*, 21.46–21.47
 $\mathbf{j} \times \mathbf{B}$ heating and anomalous skin effect, **IV**:21.49
 ponderomotive steepening and hole boring, **IV**:21.49–21.50, 21.50*f*
 relativistic effects and induced transparency, **IV**:21.52
 resonance absorption, **IV**:21.47*f*, 21.47–21.48
 vacuum heating, **IV**:21.47*f*, 21.48–21.49
 wakefield generation and electron acceleration, **IV**:21.39–21.42, 21.40*f*, 21.42*f*
 pinch, **V**:57.1–57.5, 57.2*f*, 57.3*t*, 57.4*f*
 in Raman amplification, **IV**:21.55
 Plasma(s) (*Cont.*):
 underdense, **IV**:21.36–21.46
 direct laser acceleration and betatron resonance, **IV**:21.42–21.43
 intense laser pulses, **IV**:21.38*f*, 21.38–21.39
 inverse Bremsstrahlung heating, **IV**:21.37, 21.37*f*
 ionization-induced defocusing, **IV**:21.43*f*, 21.43–21.44
 ponderomotive channel formation, **IV**:21.42
 self-channeling and self-phase modulation, **IV**:21.44–21.46, 21.45*f*
 Plasma beat wave acceleration, **IV**:21.41
 Plasma focus, for z-pinch radiation, **V**:57.3, 57.3*t*
 Plasma frequency, **IV**:2.15, 8.21
 Plasma-based EUV lasers, **V**:58.2–58.4, 58.3*f*
 Plasma-ion-assisted deposition, **IV**:7.11
 Plasma-sprayed surfaces, **IV**:6.50–6.52, 6.51*f*–6.52*f*
 Plasmons, **IV**:5.33–5.36, 5.35*f*–5.37*f*, 5.58, 8.23–8.24
 Plastic, as photographic film emulsion, **II**:29.4
 Plastic lenses, for spectacles, **III**:12.9
 Plastic spacers, for bandpass filters, **IV**:7.89
 Plastic-packaged LEDs, **II**:17.25–17.26
 Plate equations, mirror design and, **IV**:4.8–4.9
 Plate polarizers, **IV**:7.69–7.70, 7.70*f*
 Platings, diamond turning of, **II**:10.5
 Platinum:
 absorptance of, **IV**:4.41*f*, 4.48*t*, 4.50*t*, 4.51*t*
 optical properties of, **IV**:4.17*t*, 4.25*f*
 physical properties of, **IV**:4.54*t*
 reflectance of, **IV**:4.36*t*–4.37*t*, 4.41*f*
 thermal properties of, **IV**:4.69*t*, 4.70*t*
 Platinum silicon (PtSi) infrared detectors, **II**:33.7, 33.8*f*, 33.29, 33.29*t*
 Pluecker spectrum tubes, **II**:15.47, 15.47*f*, 15.47*t*
 pn junctions, **II**:24.8, 24.8*f*; **V**:20.1
 pn photodetectors, **II**:24.70*f*–24.72*f*, 24.70–24.71
 pn photodiodes, **II**:24.52*f*, 24.54–24.58, 24.55*f*–24.59*f*, 34.30
 p⁺np photodiodes, **II**:32.8
 Pockels cells, **I**:28.45, 31.9; **V**:7.33
 Pockels effect, **I**:15.23, 21.9 (*see also* Linear electro-optic effect)
 for crystals and glasses, **IV**:2.26
 and electro-optic modulators, **V**:7.6–7.11, 7.8*t*, 7.11

- Pockels effect (*Cont.*):
 and liquids, **V**:7.34
 in OTDM networks, **V**:20.1, 20.18, 20.19
- Pockels' theory, elasto-optic effect and, **V**:6.5, 6.7
- Poincaré spheres, **I**:12.27–12.29, 14.4–14.6, 14.5*f*, 14.26*f*, 28.45; **III**:18.2; **V**:19.18–19.19
- Point characteristic function, **I**:1.11, 1.14
- Point defects, **IV**:9.12–9.13
- Point diffraction interferometers, **II**:13.11*f*
- Point eikonal, **I**:1.14, 1.17
- Point images, aberrations of, **I**:1.85–1.92, 1.86*f*
- Point objects, image planes of, **I**:28.19
- Point source irradiance transmittance (PSIT), **II**:7.23
- Point source normalized irradiance transmittance (PSNIT), **II**:7.22–7.23
- Point source power transmittance (PSPT), **II**:7.23
- Point source transmittance (PST), **II**:7.5, 7.6*f*, 7.22–7.23
- Point spread function (PSF), **II**:3.16, 11.36–11.37, 11.37*f*, 11.38*f*
- Point spread matrix, **I**:15.36, 15.36*f*
- Point-angle characteristic function, **I**:1.15–1.17
- Point-by-point technique, for fiber Bragg gratings, **V**:17.8
- Points, images of, **I**:1.27
- Point-spread function (PSF), **V**:4.1, 4.34*f*, 4.34–4.36, 44.13, 44.14, 44.15*f*, 64.2
- Point-spread function (PSF), in human eye, **III**:1.21
 defined, **III**:4.1, 15.2
 in diffraction-limited eye, **III**:1.12–1.14
 direct measurements of, **III**:2.3–2.4
 and image formation, **III**:2.3
 and resolving capacity of the eye, **III**:4.4, 4.5
- Point-to-point approximation (of radiant flux transfer), **II**:34.14
- Point-to-point links, in WDM networks, **V**:21.3*f*, 21.4
- Poisson distribution, **V**:9.10
- Poisson-effect cross coupling, **V**:24.3
- Poisson's equation, **IV**:8.5; **V**:5.36
- Poisson's ratio:
 for crystals, **IV**:2.31
 for metals, **IV**:4.7, 4.69*t*
- “Poker chip” assembly, **II**:6.8, 6.9*f*
- Polacoat dichroic polarizers, **I**:13.25, 13.26, 13.28
- Polanret system, **I**:28.29
- Polar angles, **II**:35.5
- Polar decomposition, of Mueller matrices, **I**:14.39–14.40
- Polarimeters, **I**:15.3–15.6
 AxoScan Mueller matrix, **I**:15.33
 classes of, **I**:15.5
 complete and incomplete, **I**:15.4
 defined, **I**:15.7
 design metrics for, **I**:15.24–15.25
 division-of-amplitude, **I**:15.5–15.6
 division-of-aperture, **I**:15.5
 dual rotating retarder, **I**:15.16
 dual rotating retarder polarimeters, **I**:15.16, 15.16*f*
 imaging, **I**:15.6
 light-measuring, **I**:15.3–15.4, 15.11–15.13
 Mueller, **I**:15.4
 Mueller matrix, **I**:15.26–15.27
 polarization modulation, **I**:15.5
 sample-measuring, **I**:15.4
 incomplete, **I**:15.16–15.17
 for Mueller matrix elements, **I**:15.13*f*, 15.13–15.14
 spectropolarimeters, **I**:15.6
 Stokes, **I**:15.4, 15.5, 15.25
 time-sequential, **I**:15.5
 (*see also* Photopolarimeters)
- Polarimetric data-reduction equations, **I**:15.14–15.15
- Polarimetric measurement equation, **I**:15.12, 15.14–15.15
- Polarimetry, **I**:15.3–15.41
 applications of, **I**:15.29–15.41
 ellipsometry, **I**:15.30–15.32, 15.31*f*, 15.32
 liquid crystal cell and system testing, **I**:15.32*f*, 15.32–15.35, 15.33*f*, 15.34*t*
 ophthalmic polarimetry, **I**:15.39, 15.41
 polarization aberrations, **I**:15.35*f*–15.37*f*, 15.35–15.37
 polarization light scattering, **I**:15.38*f*–15.40*f*, 15.38–15.39, 15.40*f*
 remote sensing, **I**:15.37–15.38
 error analysis in, **I**:15.27–15.29, 15.29*f*
 instruments for (*see* Polarimeters)
 Mueller matrices in, **I**:15.8–15.9, 15.11
 elements of, **I**:15.13–15.14
 in error analysis, **I**:15.28
 singular value decomposition, **I**:15.25–15.27

- Polarimetry (*Cont.*):
 polarimetric data-reduction equations,
 I:15.14–15.15
 polarimetric measurement equation,
 I:15.14–15.15
 and polarization elements, I:15.17,
 15.19–15.20
 polarization generators and analyzers,
 I:15.4–15.5
 polarization (retardance) modulators,
 I:15.20–15.24, 15.22*f*
 Stokes vectors in, I:15.8–15.10
 terms in, I:15.6–15.7
- Polariscopes, Sénarmont, I:12.30
- Polaritons, IV:5.29
- Polarity, III:23.3
- Polarizance, I:12.14*n*, 14.18
- Polarization, I:12.3–12.30
 average degree of, I:14.32–14.33
 circular, V:43.5–43.8, 43.7*f*
 phase plates for, V:43.5*f*, 43.5–43.6
 and synchrotron radiation, V:55.6–55.9,
 55.7*f*
 and coherence theory, I:5.22
 concepts and conventions, I:12.4–12.6
 defined, I:15.8
 degree of, I:12.14–12.15
 of dichroic polarizers, I:13.33
 eigen-, V:7.13–7.16, 7.14*f*, 7.16*f*
 false, I:15.38
 Fresnel equations for, I:12.6–12.13
 for absorbing materials, I:12.10–12.13,
 12.13*f*
 coordinate system for, I:12.6–12.7, 12.7*f*
 for nonabsorbing materials, I:12.8–12.10,
 12.9*f*
 generators and analyzers of, I:15.4–15.5
 and inelastic scattering, IV:1.49
 of insertion devices, V:55.15–55.16
 instrumental, I:12.15
 at interface of solid, IV:8.12–8.13, 8.13*f*
 of laser diodes, V:13.13
 of light entering eye, III:1.10
 linear, V:43.2–43.4, 43.3*f*, 43.4*f*, 43.6, 43.8
 magnetic circular, I:31.21
 matrix methods for computing,
 I:12.27–12.30
 and Mueller matrices, I:14.7, 14.8,
 14.25–14.27, 14.33
 neutron, V:63.27–63.29, 63.30*f*
- Polarization (*Cont.*):
 pile-of-plates polarizers, I:12.15–12.24
 nonnormal-incidence reflection,
 I:12.15–12.18, 12.16*f*, 12.17*f*
 nonnormal-incidence transmission,
 I:12.18–12.24, 12.19*t*–I:12.20*t*, 12.21*f*
 plane of, I:12.6*n*
 relations for polarizers, I:12.14–12.15
 retardation plates, I:12.24–12.27, 12.25*f*,
 12.26*f*
 and size of electric field, IV:10.4–10.5
 Stokes–Poincaré parameters for, V:43.2
 transverse electric, V:19.7
 transverse magnetic, V:19.7
 and VCSELs, V:13.48
 of x-rays, V:43.1–43.2
 (*see also related topics*)
- Polarization aberration function (PAF),
 I:15.35
- Polarization aberrations, I:15.35*f*–15.37*f*,
 15.35–15.37
- Polarization analyzers, I:15.11; V:43.4, 43.4*f*,
 43.6–43.8, 43.7*f*
- Polarization and Directionality of Earth's
 Reflectances (POLDER) instrument,
 I:15.37
- Polarization artifacts, I:15.38
- Polarization coupling, I:15.8
- Polarization critical region, I:15.28
- Polarization dependence, IV:15.41–15.42,
 15.42*t*; V:13.58–13.59, 19.7, 19.7*f*, 19.32
- Polarization gating, II:21.7–21.8
- Polarization gradient cooling, IV:20.17–20.18,
 20.18*f*, 20.19*f*
- Polarization independence, V:6.44, 13.53
- Polarization instruments, I:12.29
- Polarization interferometers, I:32.4, 32.5*f*
- Polarization light scattering, I:15.38*f*–15.40*f*,
 15.38–15.39
- Polarization modulation (dynamic retardation),
 V:7.20*f*, 7.20–7.22, 7.21*f*
- Polarization modulation polarimeters, I:15.5
- Polarization (retardance) modulators,
 I:15.20–15.24, 15.22*f*
- Polarization power series expansion,
 IV:5.52–5.53, 5.54*t*
- Polarization scrambling, V:19.18
- Polarization sensitive OCT (PS-OCT),
 III:18.18, 18.20–18.27, 18.21*f*, 18.23*f*,
 18.25*f*–18.27*f*

- Polarization spectrometers, **I**:31.15–31.23
 optically detected magnetic resonance,
I:31.21–31.23, 31.22*f*, 31.23*f*
 polarized absorption by, **I**:31.15*f*,
 31.15–31.17, 31.18*f*
 polarized absorption/luminescence
 techniques, **I**:31.17, 31.19*f*, 31.19–31.21
 principles of, **I**:31.15, 31.15*f*
 Polarization spectroscopy, **V**:2.21, 2.22, 2.23*f*
 Polarization splitting, **IV**:7.93, 7.93*f*
 Polarization state detectors (PSDs), **I**:16.10
 Polarization state generators (PSGs), **I**:16.10
 Polarization-dependent gain (PDG), **V**:14.9,
 19.18–19.20, 21.18
 Polarization-dependent loss (PDL), **I**:14.17;
V:19.18, 21.18
 Polarization-dependent systems, **II**:34.33
 Polarization-maintaining (PM) fibers, **V**:25.3,
 25.12
 Polarization-mode dispersion (PMD),
V:21.16–21.18, 21.17*f*–21.18*f*
 Polarized absorption, **I**:31.15*f*, 31.15–31.21,
 31.18*f*, 31.19*f*
 Polarized light, **I**:15.8
 Polarizer-compensator-sample analyzer
 (PCSA) ellipsometer arrangement,
I:16.10–16.14, 16.11*f*
 Polarizers, **I**:13.1–13.57
 beam-splitter prisms as, **I**:13.6, 13.18–13.22
 Foster, **I**:13.7, 13.18*f*, 13.21–13.22
 Glan-Thompson, **I**:13.18*f*, 13.22
 Rochon, **I**:13.7, 13.18*f*, 13.18–13.21,
 13.19*f*, 13.24
 Sénarmont, **I**:13.7, 13.18, 13.18*f*, 13.21
 Wollaston, **I**:13.7, 13.18, 13.18*f*, 13.21,
 13.24
 circular, **I**:15.17–15.19
 compensators, **I**:13.53–13.56, 13.54*f*, 13.55*f*
 Cotton, **I**:13.21
 defined, **I**:15.8
 dichroic and diffraction-type, **I**:13.24–13.33,
 13.26*f*, 13.27*f*
 dichroic polarizing coatings, **I**:13.28
 measuring polarization of, **I**:13.33
 pyrolytic-graphite polarizers,
I:13.28–13.29, 13.29*f*
 sheet polarizers, **I**:13.25–13.28
 wire-grid and grating polarizers,
I:13.30–13.33, 13.31*f*, 13.32*t*
 elliptical, **I**:15.17–15.18
 Polarizers (*Cont.*):
 embedded, **IV**:7.70–7.71, 7.71*f*, 7.72*f*
 Feussner prisms, **I**:13.6, 13.7, 13.22*f*,
 13.22–13.23
 fiber-based couplers as, **V**:16.5–16.6
 Glan-Foucault prisms, **I**:13.7, 13.9, 13.11*f*,
 13.12–13.14
 Glan-type prisms, **I**:13.6, 13.6*f*, 13.8–13.15
 Frank-Ritter-type, **I**:13.6, 13.6*f*, 13.13–13.14
 Glan-Foucault, **I**:13.7, 13.9, 13.11*f*,
 13.12–13.14
 Glan-Thompson type, **I**:13.6*f*, 13.9–13.12,
 13.10*f*, 13.11*f*, 13.27
 Lippich-type, **I**:13.6, 13.6*f*, 13.7, 13.9*n*,
 13.10*f*, 13.10–13.14, 13.11*f*,
 13.12–13.13
 half-shade devices, **I**:13.56–13.57
 ideal, **I**:14.8–14.10, 14.10*t*, 15.7
 imperfect, **I**:13.33
 interference, **IV**:7.69–7.73, 7.70*f*–7.72*f*
 linear, **V**:43.2–43.3
 MacNeille, **IV**:7.70–7.72, 7.71*f*, 7.72*f*, 7.75*f*
 miniature, **I**:13.57
 multicomponent, **IV**:7.69
 multilayers, **V**:41.9
 for networking, **V**:18.7, 18.7*f*, 18.10
 Nicol-type, **I**:13.6, 13.6*f*, 13.10*f*, 13.15–13.18
 conventional, **I**:13.6*f*, 13.15–13.16
 Glan-type vs., **I**:13.8–13.9
 trimmed, **I**:13.6*f*, 13.7, 13.16*f*, 13.16–13.18
 noncalcite prisms as, **I**:13.23–13.24
 nonnormal-incidence reflection by
 Brewster angle reflection polarizers,
I:13.34–13.37, 13.34*t*–13.36*t*
 pile-of-plates polarizers, **I**:12.15–12.18,
 12.16*f*, 12.17*f*
 polarizing beam splitters, **I**:13.41–13.42
 nonnormal-incidence transmission by
 Brewster angle transmission polarizers,
I:13.37–13.39, 13.38*t*–13.39*t*
 interference polarizers, **I**:13.39–13.41, 13.40*f*
 pile-of-plates polarizers, **I**:12.18–12.24,
 12.19*t*–12.20*t*, 12.21*f*
 polarizing beam splitters, **I**:13.41–13.42
 plate, **IV**:7.69–7.70, 7.70*f*
 prism, **I**:13.2*f*–13.3*f*, 13.2–13.8, 13.4*t*–13.5*t*,
 13.6*f*
 relations for, **I**:12.14–12.15
 retardation plates as, **I**:13.43–13.53,
 13.43*t*–13.44*t*, 13.50*f*, 13.53*t*

- Polarizing angle, **I**:12.12, 12.15
- Polarizing beam splitter (PBS) prisms, **I**:13.6, 13.41–13.42
- Polarizing beam splitters (PBSs), **I**:13.41–13.42, 35.22*f*, 35.22–35.23; **IV**:7.70*f*–7.72*f*, 7.70–7.73, 7.74*f*–7.75*f*
- Polarizing coatings, dichroic, **I**:13.28
- Polarizing crystal optics, **V**:43.1–43.8
- circular polarization analyzers, **V**:43.6–43.8, 43.7*f*
 - linear polarization analyzers, **V**:43.4, 43.4*f*
 - linear polarizers, **V**:43.2–43.3, 43.3*f*
 - phase plates for circular polarization, **V**:43.5*f*, 43.5–43.6
 - and polarization of x-rays, **V**:43.1–43.2
- Polarizing mirrors, **V**:63.28
- Polaroid Corporation, **II**:29.12, 29.25
- Polaroid dichroic polarizers, **I**:13.25–13.28, 13.26*f*
- Polaroid Instant Color Film, **II**:29.14
- Polaroid “One Film,” **II**:29.14
- Polarons, in magnetic field, **IV**:5.50
- Pole-mounted luminaires, **II**:40.62, 40.63*t*
- Poling process, **IV**:12.14
- Polishers, continuous, **II**:9.7
- Polishing, of optical surfaces, **IV**:19.3
- Polishing step (of optics fabrication), **II**:9.5–9.6, 9.8
- Polyallyl diglycol carbonate, **IV**:3.4*t*
- Polyamide (Nylon), **IV**:3.4*t*
- Polyarylate, **IV**:3.4*t*
- Polycapillary optics, **V**:53.1–53.19
- alignment and measurement in, **V**:53.5–53.8, 53.6*f*, 53.7*f*
 - applications of, **V**:53.10–53.19, 53.11*f*–53.18*f*
 - and brightness of x-ray tube sources, **V**:54.16
 - collimating, **V**:53.14, 53.14*f*
 - collimation, **V**:53.8*f*–53.9*f*, 53.8–53.9
 - focusing, **V**:53.9–53.10, 53.10*t*
 - history of, **V**:53.1–53.3, 53.2*f*
 - and neutron optics, **V**:63.21–63.22
 - radiation resistance in, **V**:53.5
 - simulations and defect analysis in, **V**:53.3*f*–53.5*f*, 53.3–53.5
 - x-ray diffraction, **V**:28.5
- Polycarbonate, **IV**:3.4*t*, 3.6, 3.6*t*, 3.7*t*, 3.12
- Polycarbonate lenses, for spectacles, **III**:12.9
- Polychloro-trifluoroethylene, **IV**:3.4*t*
- Polychromatic radiation, **II**:34.9–34.10
- Polycrystalline (PC) fibers, **V**:12.2, 12.3*t*, 12.8*f*, 12.8–12.9
- Polycrystalline materials, **IV**:2.3 (*see also* Crystals)
- Polycyclohexyl methacrylate (PCHMA), **IV**:3.4*t*, 3.6, 3.6*t*, 3.7*t*
- Poly-diallylglycol (CR-39) resin, **IV**:3.11
- Polyetherimide (PEI), **IV**:3.4*t*, 3.6, 3.6*t*, 3.7*t*
- Polyethersulfone, **IV**:3.4*t*
- Polyethylene terephthalate film, **II**:29.4
- Polygon scanners, **I**:30.34*f*, 30.34–30.38, 30.35*f*
- Polymer composites, **IV**:12.26*t*
- Polymer films, **IV**:12.24
- Polymer stabilized cholesteric texture (PSCT), of liquid crystals, **V**:8.37, 8.37*f*
- Polymer sustained alignment (PSA) technique, for LC cells, **V**:8.28
- Polymer-dispersed liquid crystals (PDLCs), **V**:8.36, 8.36*f*
- Polymeric optics, **IV**:3.1–3.18
- coatings on, **IV**:3.17–3.18
 - design of, **IV**:3.7–3.11
 - aberration control, **IV**:3.8
 - aspheric surfaces, **IV**:3.8–3.9
 - athermalization, **IV**:3.9
 - dimensional variations, **IV**:3.10
 - manufacturing error budget, **IV**:3.10
 - material selection, **IV**:3.8
 - multiple cavities, **IV**:3.10
 - optical figure variations, **IV**:3.10
 - processing considerations, **IV**:3.9
 - specification, **IV**:3.10–3.11
 - strategy, **IV**:3.7–3.8
 - materials for
 - forms of, **IV**:3.2
 - optical properties, **IV**:3.5–3.7, 3.6*t*, 3.7*f*
 - physical properties, **IV**:3.2–3.5, 3.4*t*
 - selection, **IV**:3.1–3.2
 - processing of, **IV**:3.11–3.17
 - abrasive forming, **IV**:3.11–3.12
 - casting, **IV**:3.11
 - compression molding, **IV**:3.12
 - geometry considerations, **IV**:3.13–3.14
 - injection molding, **IV**:3.12–3.13
 - mechanical assembly, **IV**:3.14*f*, 3.14–3.16, 3.15*f*
 - null optics, **IV**:3.16–3.17
 - shrinkage, **IV**:3.14
 - single-point turning, **IV**:3.12
 - testing and qualification, **IV**:3.16
 - vendor selection, **IV**:3.13

- Polymer/liquid crystal composites,
V:8.36*f*–8.37*f*, 8.36–8.37
- Polymers, fully functional, **IV**:12.27*t*
- Polymer-stabilized liquid crystals (PSLCs),
V:8.35, 8.35*f*, 8.36
- Polymethyl pentene, **IV**:3.4*t*
- Polymethylmethacrylate (PMMA),
III:12.11–12.12, 20.3; **IV**:3.4*t*, 3.6, 3.6*t*,
 3.7*t*, 3.12
- Polymorphs, **IV**:2.27
- Polynomial numbering, **II**:11.39
- Polynomial-ordering number, **II**:11.7
- Polynomials:
 Chebyshev, **V**:46.6
 Legendre, **V**:46.6
 Legendre-Fourier, **V**:45.6
 orthonormal (*see* Orthonormal polynomials)
 Zernike, **V**:4.17–4.20, 4.20*t*, 5.10, 46.6
- Polystyrene, **IV**:3.4*t*, 3.6, 3.6*t*, 3.7*t*
- Polystyrene co-butadiene, **IV**:3.4*t*
- Polysulfone, **IV**:3.4*t*
- Polytetrafluoroethylene (PTFE, Teflon),
II:35.13, 38.12–38.13; **IV**:6.27
- Polyvinylidene fluoride, **IV**:3.4*t*
- Ponderomotive channel formation, **IV**:21.42
- Ponderomotive force, **IV**:21.5–21.6
- Ponderomotive potential, **II**:21.3
- Ponderomotive steepening, **IV**:21.49–21.50,
 21.50*f*
- Population inversions, **II**:16.8–16.10,
 16.12–16.13
 amplification and, **V**:19.2
 described, **II**:16.8–16.10
 mechanism for achieving, **II**:16.12–16.13,
 16.13*f*, 16.14*f*
 optical pumping for, **II**:16.16*f*–16.18*f*,
 16.16–16.19
 particle pumping for, **II**:16.14–16.16, 16.15*f*,
 16.16*f*
 semiconductor diode laser pumping for,
II:16.19
- Population trapping, velocity-selective
 coherent, **IV**:20.37
- Porcine lenses, **III**:19.11, 19.12*f*
- Pore optics, **V**:49.1–49.7, 49.2*f*–49.6*f*
- Porro prisms, **I**:19.3, 19.3*t*, 19.5*f*, 19.6, 19.6*f*
- Positive core-cladding index difference, for
 PFCs, **V**:11.12*f*, 11.12–11.14, 11.13*f*
- Positive image (in photography), **II**:29.9
- Positive orders of radiation, **V**:40.1
- Positive-intrinsic-negative (PIN) receivers,
V:9.8–9.10
- Positive-powered lenses, **IV**:3.13
- Positron emission tomography (PET), **V**:32.1
- Posterior capsule opacification (PCO):
 defined, **III**:21.2
 with intraocular lenses, **III**:21.21–21.22
- Posterior chamber, **III**:14.5, 14.6, 16.3
- Posterior limiting lamina (*see* Descemet's
 membrane)
- Posterior peripheral curve (contact lenses),
III:20.5, 20.6, 20.23–20.24, 20.24*t*
- Posterior subcapsular cataract, **III**:14.8
- Postobjective scanning, **I**:30.5, 30.29, 30.29*f*,
 30.30*f*
- Postprocessing, of SOAs, **V**:19.16, 19.17
- Posttreated Martin Black, **IV**:6.47
- Potassium bromide (KBr), **IV**:2.39*t*, 2.44*t*,
 2.48*t*, 2.51*t*, 2.56*t*, 2.62*t*, 2.69*t*, 2.77*t*
- Potassium dihydrogen phosphate (KH_2PO_4)
 (KDP), **II**:10.2; **IV**:2.39*t*, 2.45*t*, 2.48*t*, 2.51*t*,
 2.56*t*, 2.62*t*, 2.75*t*
- Potassium iodide (KI), **IV**:2.39*t*, 2.44*t*, 2.48*t*,
 2.51*t*, 2.56*t*, 2.62*t*, 2.69*t*, 2.77*t*
- Potassium niobate (KNbO_3), **IV**:2.39*t*, 2.46*t*,
 2.48*t*, 2.51*t*, 2.56*t*, 2.62*t*, 2.75*t*, 12.16–12.17
- Potassium tantalate (KTaO_3), **IV**:2.39*t*, 2.44*t*,
 2.48*t*, 2.51*t*, 2.56*t*, 2.62*t*, 2.73*t*
- Potassium titanyl arsenate (KTiOAsO_4) (KTA),
IV:17.1
- Potassium titanyl phosphate (KTiOPO_4)
 (KTP), **IV**:2.39*t*, 2.46*t*, 2.51*t*, 2.56*t*, 2.63*t*,
 2.75*t*, 17.1, 17.2, 17.28, 17.29*f*, 17.30*f*
- Pound-Drever-Hall (PDH) technique, **IV**:17.18
- Powder cloud development, in xerographic
 systems, **I**:34.9, 34.9*f*
- Powder diffraction, in polycapillary x-ray
 optics, **V**:53.14, 53.14*f*
- Powelite (CaMoO_4), **IV**:2.38*t*, 2.45*t*, 2.47*t*,
 2.50*t*, 2.55*t*, 2.61*t*, 2.72*t*, 2.77*t*
- Power:
 ASE, **V**:19.20
 channel, for EDFAs, **V**:21.41, 21.42*f*
 and dispersion-managed solitons, **V**:22.12
 for extreme ultraviolet lithography, **V**:34.5
 of Gaussian lenses, **I**:1.46–1.47, 1.47*f*
 incident, of scatterometers, **V**:1.14–1.15
 input saturation, **V**:14.3
 insertion device, **V**:55.15
 noise equivalent, **V**:1.13, 13.71

- Power (*Cont.*):
 of NPM AOTFs, **V**:6.41
 Nyquist frequency, **V**:46.8–46.9, 46.11
 saturation output, **V**:14.3, 14.4
 Stokes, **V**:14.9
 of synchrotron radiation, **V**:55.8*f*, 55.8–55.9
- Power amplifiers, **V**:14.4
- Power density, **I**:30.25
- Power dependence, of electroabsorption modulators, **V**:13.59
- Power exponential (PEX) model, of surface finish, **I**:8.15
- Power handling capability, of cw lasers, **IV**:7.14
- Power loss, **V**:21.13–21.14
- Power matrix, **I**:1.67
- Power measurement, for lasers, **II**:34.32
- Power penalties, of fiber optic devices, **V**:9.11, 15.8–15.17
- Power per pixel, **I**:1.80
- Power per unit bandwidth, **V**:7.34–7.35
- Power spectra, for surface scattering, **I**:8.12–8.13, 8.13*f*
- Power spectral density (PSD) function, **V**:41.6, 41.7, 46.7–46.9
- Power spectrum of granularity, **II**:29.21
- Power splitters, **V**:16.1, 16.4, 18.2*f*, 18.2–18.3, 18.3*f*, 18.9, 18.9*f*
- Power supply, high-voltage, **II**:31.9, 31.10*f*
- Power transfer, **I**:30.25–30.28, 30.27*f*, 30.28*f*
- Poynting vector, **IV**:8.7
- Poynting vectors, **I**:1.8, 3.3; **V**:7.3, 7.5, 42.2
- Praseodymium, **IV**:14.34–14.35, 14.35*f*, 14.36*f*
- Praseodymium-doped fiber amplifiers (PDFAs), **V**:14.2, 14.2*t*, 14.7
- Pre-amp semiconductor optical amplifiers (SOAs), **V**:19.23
- Predictable quantum efficiency (PQE) devices, **II**:34.29–34.30
- Preflash, of cameras, **I**:25.16, 25.17
- Prefocusing, with refractive x-ray lenses, **V**:37.8
- Preform manufacture, of fiber lasers, **V**:25.26–25.27
- Preloads, **II**:6.2
- Prentice's rule, **III**:12.16, 13.15, 13.18, 13.28, 20.31
- Preobjective scanning, **I**:30.5, 30.28, 30.29*f*
- Presbyopia (old eye/old sight), **III**:1.7, 12.3, 16.5
 accommodation restoration for, **III**:14.29–14.30
 as age-related, **III**:14.8, 14.9*f*
 assessment of, **III**:12.7–12.8
- Presbyopia (old eye/old sight) (*Cont.*):
 clinical onset of, **III**:12.3, 12.7–12.8
 correction of, **III**:14.27–14.30
 accommodation restoration, **III**:14.29–14.30
 bifocals for, **III**:12.8
 contact lenses for, **III**:12.13–12.14, 14.27–14.28
 hyperchromatic lens design for, **III**:14.14
 intraocular lenses for, **III**:14.28–14.29
 noncataract related refractive surgeries for, **III**:14.29
 spectacles for, **III**:12.10–12.11, 14.27
 defined, **III**:12.1, 14.2, 21.2, 23.3
 as problem with computer work, **III**:23.11–23.12
 and UV radiation exposure, **III**:14.23
- Pressure broadening, of spectral lines, **V**:56.5–56.6
- Pressure force, **IV**:20.7
- Prewhitening techniques, in x-ray mirror metrology, **V**:46.9
- Primaries (colors stimuli), **II**:30.16
- Primary focal length (term), **I**:3.12
- Primary lights (primaries), **III**:10.4
 additive mixing of, **III**:10.4, 10.5*f*
 CMFs for, **III**:10.10, 10.11
 in color matching, **III**:10.6, 10.7
 defined, **III**:10.1, 10.4*t*
 for DKL color space, **III**:10.19
 imaginary, **III**:10.10
 in Maxwell's matching method, **III**:10.8, 10.8*f*
 perceptual vs. physical matches of, **III**:10.7
 for stimulus spaces, **III**:10.11
 in trichromatic color matching, **III**:10.7, 10.8
 and vector representations, **III**:10.27–10.29
- Primary position, **III**:1.42, 13.2
- Primary visual cortex, **III**:2.12, 2.14
- Primate eye lens, **III**:19.13*f*, 19.14, 19.14*f*
- Principal plane-to-principal plane conjugate matrices, **I**:1.69
- Principal rays, **II**:1.4, 1.11*f*, 1.12
- Principal rays (term), **I**:1.75, 17.8
- Principal transmittance (term), **I**:12.14–12.16
- Principle dispersion (term), **IV**:2.23
- Prism(s), **II**:40.45*f*, 40.46; **V**:18.5, 18.5*f*, 63.22
 (*see also* Dispersive prisms and gratings; Nondispersive prisms)
 angle measurement in, **II**:12.14–12.16, 12.15*f*–12.17*f*
 axis wander of, **I**:13.15

Prism(s) (*Cont.*):

beam-splitter, **I**:13.7, 13.18*f*, 13.18–13.22
 Bertrand-type Feussner, **I**:13.23
 Brewster angle, **I**:13.13
 of calcite, **I**:13.20, 13.23
 in contact lenses, **III**:20.2
 distortion from interocular aniso-
 magnification, **III**:13.16–13.17, 13.17*f*
 effects on vergence and phoria,
 III:13.25–13.27
 errors of alignment with, **III**:13.27–13.29
 Feussner, **I**:13.6, 13.7, 13.22*f*, 13.22–13.23
 Foster, **I**:13.7, 13.18*f*, 13.21–13.22
 Foucault, **I**:13.7, 13.17
 Frank-Ritter-type, **I**:13.6, 13.6*f*, 13.13–13.14
 Fresnel's biprism, **I**:2.16, 2.16*f*
 Glan-Foucault, **I**:13.7, 13.9, 13.11*f*, 13.12–13.14
 Glan-Taylor, **I**:13.7, 13.9*n*, 13.10*f*,
 13.10–13.14, 13.11*f*
 Glan-Thompson, **I**:13.6, 13.6*f*, 13.9–13.12,
 13.10*f*, 13.18*f*, 13.22
 field angle of, **I**:13.12
 sheet polarizers vs., **I**:13.27
 transmission by, **I**:13.9–13.10, 13.11*f*
 Glan-type, **I**:13.6, 13.6*f*, 13.8–13.15
 Glazebrook, **I**:13.6, 13.9
 Halle, **I**:13.16*f*, 13.17
 Hartnack-Prazmowski, **I**:13.16*f*, 13.17
 Jellet-Cornu, **I**:13.56
 length-to-aperture (L/A) ratio, **I**:13.7
 Lippich-type, **I**:13.6, 13.9*n*, 13.10*f*, 13.11*f*,
 13.12–13.14, 13.12*n*, 13.56
 Marple-Hess, **I**:13.12, 13.13
 mounting of, **II**:6.11–6.17
 bonded mountings, **II**:6.13–6.15, 6.15*f*,
 6.16*f*
 flexure mountings, **II**:6.15–6.17, 6.16*f*
 mechanically clamped mountings,
 II:6.12, 6.13*f*
 spring-loaded mountings, **II**:6.13, 6.14*f*
 Nicol curtate, **I**:13.16*f*, 13.17
 Nomarski, **I**:28.40, 28.41
 noncalcite, **I**:13.23–13.24
 nonuniform magnification of, **III**:13.10*f*
 penta-, **II**:6.14*f*, 6.15*f*, 12.12, 12.12*f*
 polarizing beam splitter, **I**:13.6, 13.41–13.42
 right-angle, **II**:12.15, 12.16, 12.16*f*, 12.17*f*
 Risley, **II**:12.4, 12.4*f*
 Rochon, **I**:13.7, 13.18*f*, 13.18–13.21, 13.19*f*,
 13.19–13.20, 13.24

Prism(s) (*Cont.*):

rotating glass block, **II**:12.4, 12.4*f*
 semifield angle of, **I**:13.7
 Sénarmont, **I**:13.7, 13.18, 13.18*f*, 13.21
 sliding, **II**:12.4, 12.4*f*
 Steeg and Reuter Nicol, **I**:13.17
 stimulus value of, **III**:13.21
 Wollaston, **I**:13.7, 13.18, 13.18*f*, 13.21, 13.24,
 28.39, 28.40, 32.4
 Zerodur, **II**:6.16, 6.16*f*
 Prism diopter, **III**:12.2, 25.1
 Prism polarizers, **I**:13.2*f*–13.3*f*, 13.2–13.8,
 13.4*t*–13.5*t*, 13.8*f* (*see also specific types*)
 Prism spectrometers, **I**:20.2–20.3, 20.3*f*
 Prismatic effects, **III**:12.8
 and bifocal jump, **III**:13.15
 with contact lenses, **III**:20.30–20.31
 prism-ballasted contacted lenses, **III**:20.30
 unintentional (induced) prism, **III**:20.31
 with corrected ametropia, **III**:12.16
 defined, **III**:12.2
 Prismatic facets, **I**:30.34–30.35
 Prism-ballasted contacted lenses, **III**:20.30
 Prism-based monochromators, **IV**:5.59
 Procedural programs (optical software), **II**:3.7
 Processors, real-time, **V**:5.34*f*, 5.34–5.35, 5.35*f*
 Profile analysis, in x-ray mirror metrology,
 V:46.6–46.12, 46.7*f*, 46.10*f*
 Profile errors, in polycapillary x-ray optics,
 V:53.3, 53.3*f*
 Profilometry:
 height, **V**:46.3
 1D, 2D, and 3D, **V**:46.6
 slope, **V**:46.3–46.6, 46.5*f*
 Progressive addition lenses (PALs),
 III:12.10–12.11, 12.11*f*
 Projected area, in photometry/radiometry,
 II:36.3*f*, 36.3–36.4, 36.3*t*
 Projected solid angle (PSA), **II**:39.5, 39.5*f*
 Projection density, **II**:29.7, 29.7*f*
 Projection lenses, **II**:6.6*f*
 Projection systems, **II**:39.23*f*, 39.23–39.24
 Projective transformation, **I**:1.56 (*see also*
 Collineation)
 Propagation effects, on third-order optical
 nonlinearities, **IV**:16.24–16.26
 Propagation of light, coherence theory and,
 I:5.13–5.19, 5.14*f*–5.16*f*
 Propagation of mutual intensity, **I**:6.4

- Proportional counters, x-ray detectors and, **V**:60.4–60.5, 60.9*t*, 60.10*t*
- Proportional integral derivative (PID)
 controllers, **II**:22.10*f*, 22.10–22.12, 22.11*f*
- Proportional-integral (PI) amplifier circuit, **II**:22.6*f*
- Proportionality (color matching), **III**:10.8
- Protanopes, **III**:10.16
- Proton stripe lasers, **II**:19.35*f*
- Proton-bombardment-defined lasers, **II**:19.41
- Proustite (Ag_3AsS_3), **IV**:2.38*t*, 2.46*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*
- Proximal (psychic) convergence, **III**:1.35
- Proximal stimuli (human vision), **III**:4.2
- Proximity-focus electronic lens, **II**:31.8, 31.23*f*, 31.23–31.24, 31.24*f*
- Proximity-focused MCP IIs, **II**:31.9, 31.9*f*, 31.16–31.18, 31.17*t*, 31.18*f*, 31.19*f*
- Pseudoaccommodation, **III**:14.27, 14.29, 14.30
- Pseudo-Brewster angle, **I**:12.13
- Pseudophakia, **III**:12.15, 14.28
 correction of, **III**:12.14–12.15
 defined, **III**:14.2, 21.2
- Psychic convergence, **III**:1.35
- Psychophysical comparison method (retinal image quality), **III**:1.22
- Psychophysical measurement, **III**:3.1–3.10
 of adjustment tasks, **III**:3.4–3.5
 magnitude production, **III**:3.5
 matching, **III**:3.4–3.5
 nulling, **III**:3.4
 threshold, **III**:3.4, 3.5*f*
 definitions related to, **III**:3.2–3.3
 of judgment tasks, **III**:3.6–3.8
 ideal observer, **III**:3.6
 rating scale, **III**:3.6–3.8, 3.7*f*
 response time for, **III**:3.8
 two-alternative forced choice (2afc), **III**:3.8
 yes-no, **III**:3.6
 magnitude estimation in, **III**:3.8
 professional tips for, **III**:3.10
 stimulus sequencing in, **III**:3.9
 method of constant stimuli, **III**:3.9
 sequential estimation methods, **III**:3.9
 of visual acuity, **III**:4.6–4.7
 visual stimuli, **III**:3.3–3.4
- Psychophysical photometry, **II**:34.37
- Psychophysical test method (color vision), **III**:11.9, 11.11, 11.12
- Psychophysics, **III**:4.1, 15.2
- Pterygium, **III**:7.1, 7.6, 7.7
- P-type impurities, **II**:17.23
- Pulse amplification, chirped, **V**:25.2, 25.32, 25.33
- Pulse area, **IV**:11.8
- Pulse code modulation, **V**:20.8
- Pulse excitation, chirped, **IV**:11.25–11.26
- Pulse generation, **IV**:18.4, 21.31
- Pulse height, of photomultipliers, **II**:27.7–27.8
- Pulse propagation, **IV**:14.20–14.22
- Pulse train interferometry, **II**:20.12, 20.12*f*
- Pulse trains:
 about, **II**:20.3–20.5, 20.4*f*, 20.5*f*
 attosecond, **II**:21.6, 21.7
 and backscattering, **II**:20.13–20.15, 20.15*f*
 soliton solution and steady-state, **II**:20.5–20.9
- Pulse width modulation (PWM), in
 thermomagnetic recording, **I**:35.18–35.19
- Pulsed lasers, **II**:23.18; **IV**:14.15–14.16, 14.16*f*
- Pulsed transient Raman scattering, **IV**:15.22–15.25, 15.24*f*–15.26*f*, 15.24*t*
- Pulsed-dye lasers, **V**:5.32
- Pump (incident light), **IV**:15.1
- Pump cladding, **V**:25.3
- Pump depletion, **IV**:15.9, 15.10, 15.15, 15.20, 15.20*f*
- Pump wavelength, EDFAs and, **V**:14.5–14.6
- Pump-enhanced singly resonant oscillators (PE-SROs), **IV**:17.2–17.4, 17.3*f*, 17.4*f*, 17.17–17.20, 17.18*f*–17.20*f*, 17.27–17.28
- Pumping:
 for cw SROs, **IV**:17.2–17.4
 for population inversions, **II**:16.14–16.19
 optical, **II**:16.16*f*–16.18*f*, 16.16–16.19
 particle, **II**:16.14–16.16, 16.15*f*, 16.16*f*
 semiconductor diode laser, **II**:16.19
- Pumping techniques, for fiber lasers, **V**:25.9–25.13, 25.11*f*, 25.12*f*
- Pump-probe (excite-probe) measurements, **IV**:16.26–16.27, 16.27*f*
- Pump-probe spectroscopy, **IV**:18.18–18.19
- Pump-probe techniques, **IV**:18.6, 18.6*f*
- Pump-resonant singly resonant oscillators, **IV**:17.17–17.20, 17.18*f*–17.20*f*
- Pumps, for rare-earth-doped amplifiers, **V**:14.2–14.3, 14.3*f*
- Punch through (in color films), **II**:30.4
- Punctum remotum*, **III**:12.8

- Pupil(s):
 annular, **V**:4.10–4.16, 4.11*f*–4.15*f*, 4.15*t*
 center of, **III**:1.20
 defined, **III**:6.1
 diameter of (*see* Pupil diameter/size)
 interpupillary distance, **III**:1.39–1.41, 1.40*f*
 of lenses, **I**:1.76*f*, 1.76–1.79, 1.78*f*, 17.8–17.9
 in Maxwellian viewing, **III**:5.9
 noncircular, **II**:11.4
 in ophthalmoscopic methods, **III**:1.23
 and retinal illuminance, **III**:1.12
 true area of, **III**:1.10
 and visual acuity, **III**:4.9
 and visual instruments, **III**:1.27–1.28
 (*see also* Entrance pupil; Exit pupil)
- Pupil aberrations, **I**:1.76
- Pupil angular magnification, **I**:1.78
- Pupil conjugate plane:
 in Maxwellian viewing, **III**:5.8–5.9
 rotating a mirror in, **III**:5.16, 5.16*f*
 shutters at, **III**:5.13, 5.13*t*
- Pupil diameter/size, **III**:1.8*f*, 1.8–1.9
 age-related changes in, **III**:1.18, 14.6–14.7, 14.7*f*, 14.14
 change in, **III**:16.3
 and correction for refractive error, **III**:1.23, 1.24
 and correction for SCE-1, **III**:9.4, 9.13
 and depth-of-focus, **III**:1.29, 1.30*f*
 in Maxwellian viewing, **III**:6.3–6.5, 6.7–6.12
 coherent illumination, **III**:6.9–6.12, 6.11*t*
 incoherent target, **III**:6.7–6.9, 6.8*f*, 6.9*t*
 optimal, **III**:15.22
 and RMS wavefront error, **III**:1.15–1.16
 and Stiles-Crawford effect, **III**:1.10
 variation in parameters as function of, **III**:1.13
- Pupil distortion, **I**:1.78
- Pupil function, OTF and, **III**:1.21
- Pupil imaging, **I**:1.76
- Pupil magnification, **I**:1.76
- Purcell-Pennypacker method, **I**:7.15
- Pure diattenuators, **I**:15.8
- Pure retarders, **I**:15.8
- Purity:
 in color CRTs, **III**:22.9
 of optical materials, **IV**:2.5
- Pursuit movements, **III**:1.43
- Push processing (of film), **II**:30.22
- Pushbroom scan, **I**:30.18
- Pyramidal error, **II**:12.14, 12.15, 12.15*f*
- Pyramidal facets, **I**:30.34*f*, 30.34–30.35, 30.35*f*
- Pyrex glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*
- Pyroelectric detectors, **II**:24.6, 24.6*f*, 24.26*f*–24.29*f*, 24.26–24.29, 28.2, 28.6, 28.6*f*, 28.7, 28.7*t*, 33.10
- Pyroelectric electrical substitution radiometers, **II**:34.27–34.28
- Pyroelectric hybrid arrays, **II**:28.11*f*, 28.11–28.12, 28.12*f*
- Pyrolytic-graphite polarizers, **I**:13.28–13.29, 13.29*f*
- Q factor, for OTDM networks, **V**:20.11
- Q-switched lasers, **II**:16.26–16.27, 16.27*f*
- Q-switched modelocking, **IV**:18.4, 18.5, 18.5*f*
- Q-switched oscillators, **V**:25.30–25.31
- Quadratic (Kerr) electro-optic effect, **V**:7.6, 7.9*t*–7.10*t*, 7.11
- Quadratic Stark effect, **V**:56.6
- Qualification, of polymers, **IV**:3.16
- Quality, image, **II**:4.6–4.7
- Quality Assurance of Ultraviolet Measurements in Europe (QASUME), **II**:38.5
- Quality factor (of microcavities), **IV**:9.7–9.8, 9.8*f*
- Quantitative phase microscopy, **I**:28.27
- Quantized center-of-mass motion, of atoms, **II**:23.45
- Quantum box, **II**:19.18
- Quantum cascade lasers, **II**:16.36
- Quantum coherence theory, **I**:5.2
- Quantum coherence tomography, **IV**:23.13
- Quantum detection efficiency (QDE), **V**:60.7, 61.2–61.3
- Quantum dots, **II**:16.7, 19.18, 26.4*f*, 26.5; **IV**:12.25; **V**:19.12, 19.21
- Quantum efficiency (QE), **II**:25.3, 25.4, 34.29
 defined, **II**:24.12
 detective, **II**:29.23
 of photodetectors, **II**:24.18, 24.19
 of photoemissive detectors, **II**:24.35*f*–24.38*f*, 24.36–24.38
 of photomultipliers, **II**:27.7
 of *pin* photodiodes, **II**:25.8; **V**:13.65, 13.66
- Quantum electrodynamic (QED) shifts, **I**:10.4
- Quantum electrodynamics (QED), **II**:9.6, 23.9–23.13

- Quantum entanglement, in optical interferometry, **IV**:23.1–23.15
 concepts and equations for, **IV**:23.1–23.4, 23.2*f*, 23.4*f*
 digital approaches to, **IV**:23.7–23.9
 Heisenberg limit, **IV**:23.6–23.7
 N00N state, **IV**:23.9–23.12, 23.10*f*, 23.11*f*
 and quantum imaging, **IV**:23.13–23.14
 and remote sensing, **IV**:23.14–23.15
 shot-noise limit, **IV**:23.4–23.6, 23.5*f*
- Quantum fluctuations, **IV**:15.38, 15.39*f*, 23.5, 23.5*f*
- Quantum imaging, **IV**:23.13–23.14
- Quantum interferences, **IV**:14.8
- Quantum limit, of digital on-off keying receivers, **V**:9.9
- Quantum limited imaging (QLI), **II**:31.3–31.4
- Quantum mechanical model, for solids, **IV**:8.4, 8.24–8.25
- Quantum mechanics, third-order optical nonlinearities and, **IV**:16.4–16.7, 16.5*f*
- Quantum photodetectors, **II**:24.6–24.10, 24.7*f*–24.9*f*
- Quantum remote sensing, **IV**:23.14–23.15
- Quantum resonance absorption, **II**:22.16, 22.17
- Quantum sensitivity, **II**:30.9
- Quantum theory of lasers, **II**:23.14–23.35
 about, **II**:23.5–23.6
 density-operator approach to, **II**:23.14–23.33
 derivation of Scully-Lamb master equation, **II**:23.17–23.22, 23.19*f*
 photon statistics, **II**:23.22–23.27, 23.23*f*, 23.25*f*, 23.27*f*
 spectral properties, **II**:23.28–23.33, 23.30*f*
 spectrum, **II**:23.28–23.33
 time evolution of the field in Jaynes-Cummings model, **II**:23.15*f*, 23.15–23.17
 Heisenberg-Langevin approach to, **II**:23.33–23.35
- Quantum theory of nonlinear optical susceptibility, **IV**:10.9–10.10
- Quantum trajectories, in attosecond optics, **II**:21.3–21.4
- Quantum well (QW) detectors, **II**:26.4*f*, 26.5
- Quantum well infrared photodetectors (QWIPs), **II**:25.15–25.17, 25.16*f*, 25.17*f*, 33.9
- Quantum well (QW) lasers, **II**:16.7, 19.9–19.18, 19.20*t*; **V**:13.24–13.28, 13.25*f*, 13.26*f*
 GRIN SCH single, **II**:19.14, 19.14*f*
 long wavelength, **II**:19.17*f*, 19.17–19.18
 schematic of, **II**:19.10*f*
 strained, **II**:19.15–19.17, 19.16*f*
 threshold modal gain, **II**:19.12*f*, 19.12–19.15, 19.13*f*, 19.15*f*
- Quantum well (QW) photodetectors, **II**:25.16*f*, 25.16–25.17, 25.17*f*
- Quantum wells (QWs), **V**:20.1
 coupled, **V**:13.58
 and EA modulators, **V**:20.20
 and electrorefractive modulators, **V**:13.62
 in fiber optic systems, **V**:13.4
 and SOAs, **V**:19.7, 19.11, 19.12
 and VCSELs, **V**:13.43
- Quantum wire, **II**:16.7, 19.18, 26.4*f*, 26.5
- Quantum-confined Stark effect (QCSE), **I**:21.11–21.12, 21.32; **V**:13.2, 13.56, 13.56*t*, 20.1
- Quantum-well intermixing, **I**:21.19, 21.20
- Quart-enhanced photoacoustic spectroscopy (QEPAS), **IV**:17.25–17.26, 17.26*f*
- Quarter pitch length, of the rod, **I**:24.6
- Quarter-wave circular retarders, Mueller matrices for, **I**:14.12*t*
- Quarter-wave linear retarders, Mueller matrices for, **I**:14.11, 14.12*t*
- Quarter-wave phase plates, **V**:43.6
- Quarter-wave plates, **I**:12.25–12.27, 12.26*f*
- Quarter-wavelength-shifted gratings, **V**:13.31*f*, 13.31–13.32
- α -Quartz (SiO₂), **IV**:2.40*t*, 2.46*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.64*t*, 2.71*t*
- Quartz retardation plates, **I**:13.46–13.48
- Quartz-envelope lamps, **II**:15.20, 15.21
- Quasi-homogeneous sources of light, **I**:5.11–5.12, 5.19
- Quasi-monochromatic sources of light, **I**:5.11
- Quasi-phase-matched (QPM) nonlinear materials, **IV**:17.1, 17.13–17.14
- Quasi-plane-wave solution (to transverse profile problem), **IV**:12.10
- Quaternary structure, of fiber optic devices, **V**:13.2
- QUEST, **III**:3.9
- QW ridge (QWR) waveguide lasers, **II**:19.19, 19.20*t*

- Rabbit lenses, **III**:19.8
- Rabi cycles and Rabi cycling, **II**:23.21
 defined, **II**:20.24
 off resonance, **II**:20.27
 on resonance, **II**:20.26*f*, 20.26–20.27, 20.27*f*
- Rabi frequency, **IV**:11.4
- Racah parameters, for ion energy levels, **I**:10.12
- Rack-stack laser arrays, **II**:19.29, 19.30*f*
- Radar, synthetic aperture, **I**:11.6–11.8,
 11.7*f*–11.8*f*
- Radial circle polynomial, **II**:11.7, 11.9*f*–11.10*f*
- Radial edge lift, **III**:20.2
- Radial gradients, **I**:24.5*f*, 24.5–24.8, 24.7*f*;
III:19.3–19.5, 19.4*f*, 19.5*f*
- Radial keratotomy, **III**:12.14, 16.9–16.10, 16.10*f*
- Radial shearing interferometers, **II**:13.12, 13.13*f*
- Radial symmetry, of scanners, **I**:30.5
- Radian (rad), **II**:36.3
- Radiance, **I**:5.8; **II**:34.9, 34.9*f*, 37.4*t*, 37.5, 38.2,
 38.11*t*, 38.13*f*–38.16*f*, 38.13–38.16, 39.2*t*
 thermal spectral, **V**:3.18, 3.20, 3.20*f*
 of water, **IV**:1.5*t*, 1.6, 1.7*f*
- Radiance conservation theorem, **II**:34.12–34.13
- Radiance temperature (unit), **II**:37.4*t*, 37.6
- Radiance units, **II**:34.24
- Radiant emittance, **I**:5.7–5.8
- Radiant energy, **II**:34.7, 37.4*t*, 37.6
- Radiant exitance (emittance), **II**:15.4–15.6,
 15.5*t*, 15.6*f*, 35.3, 37.4*t*, 37.5, 37.5*f*
- Radiant exposure, **II**:37.4*t*, 37.6
- Radiant flux (power), **II**:34.7, 34.11, 34.17–34.18,
 36.4, 36.6*t*, 37.3, 37.4*t*, 37.6, 38.2
- Radiant incidence (*see* Irradiance)
- Radiant intensity, **I**:5.8; **II**:36.4, 37.4, 37.4*t*, 39.2*t*
- Radiant power transfer, **II**:34.12–34.13, 34.13*f*
- Radiant transfer approximations,
II:34.13–34.20
 approximate radiance at an image,
II:34.19–34.20
 lambertian, **II**:34.14–34.18, 34.15*f*, 34.16*f*
 point-to-point, **II**:34.14
 radiometric effect of stops and vignetting,
II:34.18–34.19, 34.19*f*
- Radiation, **II**:34.23–34.27
 actinic effects of, **II**:34.6, 34.7
 artificial sources of [*see* Artificial sources (of
 radiation)]
 baseline standard of, **II**:15.9, 15.9*f*, 15.10*f*, 15.12*f*
 from blackbodies, **II**:34.23–34.24
 from blackbody simulators, **II**:34.24–34.26
- Radiation (*Cont.*):
 Bremsstrahlung, **V**:31.3
 continuous, **V**:54.4–54.6, 54.5*f*
 from laser-generated plasmas, **V**:56.1,
 56.8–56.10
 from pinch plasma sources, **V**:57.1
 and x-ray fluorescence, **V**:29.3, 29.5, 29.6,
 29.11
 characteristic
 Bremsstrahlung radiation as, **V**:56.8–56.10
 from laser-generated plasmas, **V**:56.2–56.10
 recombination radiation as, **V**:56.10
 spectral lines as, **V**:56.2–56.8
 from x-ray tube sources, **V**:54.6–54.8, 54.7*f*
 between circular source and detector,
II:34.15*f*, 34.15–34.16
 coherent, **I**:30.2, 30.25–30.26
 commercial sources of [*see* Commercial
 sources (of radiation)]
 Compton sources of, **V**:55.2–55.3, 55.3*t*
 incandescent sources of [*see* Incandescent
 sources (of radiation)]
 incoherent, **I**:30.2, 30.26, 30.27
 infrared, **I**:13.47 [*see* Infrared (IR) radiation]
 ionizing, **V**:15.17
 lab-based sources of, **V**:50.2–50.7, 50.4*f*,
 50.6*f*, 50.8*f*
 from laser-generated plasmas, **V**:56.2–56.10
 and lasers, **II**:16.4
 negative orders of, **V**:40.1
 photographic film speed and sensitivity to
 high-energy, **II**:30.19–30.20
 Planck radiation law, **V**:3.18, 3.20
 polychromatic, **II**:34.9–34.10
 positive orders of, **V**:40.1
 recombination, **V**:56.10
 from synchrotrons, **II**:34.26–34.27
 synchrotron sources of (*see* Synchrotron
 radiation sources)
 through absorbing media, **II**:34.13
 transfer of, **II**:7.21–7.22
 ultraviolet [*see* Ultraviolet (UV) radiation]
 Unruh, **V**:58.2
 working standards of, **II**:15.9–15.13, 15.10*f*,
 15.12*f*, 15.13*f*
 X-pinch sources of, **V**:57.3, 57.3*t*, 57.4
 zero order of, **V**:40.1
- Radiation damage, life-span, **III**:14.22–14.23
- Radiation fields, coherence theory and, **I**:5.15*f*,
 5.15–5.16, 5.16*f*

- Radiation hazards (*see* Ocular radiation hazards)
- Radiation induced loss, **V**:15.17
- Radiation law, **II**:15.4–15.7, 15.5*f*, 15.5*t*, 15.6*f*
- Radiation modes, of optical waveguides, **I**:21.4
- Radiation pressure force, **IV**:20.7
- Radiation resistance, of polycapillary x-ray optics, **V**:53.5
- Radiation resistance, of polymers, **IV**:3.5
- Radiative escape, **IV**:20.29
- Radiative lifetime, **I**:31.12–31.13, 31.13*f*
- Radiative lifetimes, **II**:16.4, 17.4
- Radiative quantum efficiency, of PDFAs, **V**:14.7
- Radiative recombination, **II**:17.2
- Radiative transfer, in volume scattering, **I**:9.10–9.13, 9.11*f*, 9.13*f*
- Radiative transfer theory, **IV**:1.4
- Radiators, blackbody, **II**:34.23–34.24
- Radio astronomy, **I**:5.23
- Radio frequency (rf) modulation, **II**:22.14
- Radiography, **V**:31.1–31.4, 31.2*f*–31.4*f*, 62.1
- Radiometers and radiometry, **II**:34.3–34.37, 36.1–36.17, 36.3*f*, 36.3*t*
- about, **II**:34.5–34.7, 36.2–36.4
- approximate (*see* Radiant transfer approximations)
- cavity-shaped, **II**:34.28
- for color matching, **III**:10.12
- concepts/terminology of, **II**:34.7, 39.2, 39.2*t*
- conversion between radiometric and photometric quantities, **II**:34.12*t*, 36.11–36.14, 36.12*f*–36.14*f*
- correction for SCE-1 in, **III**:9.1–9.15
- defined, **II**:37.1
- electrical substitution, **II**:34.27–34.29
- geometrical concepts of, **II**:34.8–34.9, 34.9*f*
- of II electronic imaging, **II**:31.4–31.5
- illuminance-luminance relationship, **II**:37.9*f*
- infrared fibers for radiometry, **V**:12.3*t*
- integrating sphere device, **II**:37.9–37.10, 37.10*f*
- laser as characterization tool for, **II**:34.32
- normalization, **II**:36.14–36.17, 36.15*t*, 36.16*f*
- ocular radiometry, **III**:1.11–1.12
- photometry vs., **II**:34.6
- Planck's law, **II**:37.10–37.11
- quantities and units in, **II**:37.3–37.7, 37.4*t*, 37.5*f*, 37.7*t*
- spectral dependence of, **II**:34.9–34.10
- statistical radiometry, **I**:5.22
- Stefan-Boltzmann's law, **II**:37.11
- Radiometers and radiometry (*Cont.*):
- symbols/units/nomenclature of, **II**:36.4–36.5
- thermopile-based, **II**:34.27
- weighting functions, **II**:36.17
- Wien's displacement law, **II**:37.11
- Radiometric quantities, of water, **IV**:1.4–1.9, 1.5*t*, 1.7*f*
- Radius of torsion (space curves), **I**:1.19
- Rainbows, **V**:3.41, 3.41*f*
- Raman amplification, **IV**:21.55
- Raman amplifiers, **IV**:15.4, 15.4*f*, 22.15; **V**:19.27, 21.42*f*–21.44*f*, 21.42–21.44
- Raman bands, of glass, **V**:11.24
- Raman cooling, **IV**:20.21
- Raman cross sections, **IV**:15.5
- Raman effect, inverse, **V**:10.5, 10.6, 14.9
- Raman fiber amplifiers, **V**:14.1, 14.2, 14.2*t*, 14.8*f*, 14.8–14.9, 14.10*f*
- Raman gain, **V**:10.5, 10.6, 21.42*f*
- Raman gain coefficients, **IV**:15.16*t*–15.18*t*
- Raman generators, **IV**:15.4, 15.4*f*
- Raman induced Kerr effect (RIKE), **IV**:16.3*t*, 16.12, 16.17
- Raman linewidths, **IV**:15.9, 15.10*f*, 15.11*t*–15.20*t*
- Raman microspectroscopy, **III**:19.1
- Raman modes, of crystals and glasses, **IV**:2.11
- Raman oscillators, **IV**:15.4, 15.4*f*
- Raman resonance, **V**:14.8
- Raman scattering, **I**:31.30–31.31, 31.31*f*; **IV**:15.1–15.43
- anti-Stokes, **IV**:15.4, 15.4*f*, 15.32–15.34, 15.33*f*, 15.35*f*, 15.42, 15.42*t*, 15.43*f*
- and atmospheric optics, **V**:3.12, 3.21
- backward, **IV**:15.41, 21.38*f*, 21.38–21.39
- Brillouin vs., **IV**:15.1
- coherent, **IV**:15.3, 15.4, 15.4*f*, 15.34, 15.42, 15.42*t*, 15.43*f*
- and configurational relaxation of solids, **V**:2.15–2.17, 2.18*f*
- for crystals and glasses, **IV**:2.27
- forward, **IV**:21.38*f*, 21.38–21.39
- measurement with, **IV**:5.76*f*, 5.76–5.83, 5.78*f*–5.82*f*
- and Raman interactions, **IV**:15.2–15.3, 15.3*t*
- regimes of, **IV**:15.3
- in solids, **IV**:8.16, 8.18, 8.19*t*–8.20*t*
- spontaneous, **IV**:15.3, 15.5

- Raman scattering (*Cont.*):
 stimulated, **V**:25.6 (*see* Stimulated Raman scattering)
 and fiber optic communication links, **V**:15.8
 and optical fiber amplifiers, **V**:14.2, 14.8
 in optical fibers, **V**:10.1, 10.4–10.7, 10.5*f*
 and photonic crystal fiber guidance, **V**:11.23, 11.24, 11.26
 and WDM networks, **V**:21.20
 transient, **IV**:15.22–15.32
 broadband effects, **IV**:15.28–15.32, 15.29*f*
 phase pulling, **IV**:15.26–15.27, 15.27*f*
 pulsed, **IV**:15.22–15.25, 15.24*f*–15.26*f*, 15.24*t*
 solitons, **IV**:15.27–15.28, 15.29*f*
 spectral properties, **IV**:15.32
 by water, **IV**:1.48, 1.49, 1.49*f*
 Raman shifts, **IV**:15.2
 Raman sidebands, generation of, **IV**:14.31, 14.32*f*
 Raman spectroscopy, **IV**:5.57–5.58, 7.48, 7.66, 7.83, 7.96, 7.97*f*, 7.98
 Raman susceptibility, **IV**:15.5–15.6
 Raman threshold, **IV**:15.38; **V**:10.7
 Raman transition frequencies, **IV**:15.11*t*–15.15*t*
 Raman-Nath diffraction, **I**:11.9, 11.9*f*
 Raman-Nath diffraction regime, **V**:6.4, 6.6
 Raman-Nath equations, **V**:6.6–6.7
 Ramsden disk, **I**:28.8, 28.16
 Ramsey fringes, **IV**:11.20–11.22, 11.21*f*
 Random Device Slope Scanner, **V**:46.3
 Random-dot stereograms, **III**:2.40
 Range finders, **II**:12.3*f*, 12.3–12.4, 12.4*f*
 Range gating, **II**:31.28–31.30
 Range of focus, **I**:1.85
 Rapid adiabatic passage (RAP), **IV**:14.1
 Rapid Rectilinear lenses, **I**:17.27
 Rare-earth ions, **I**:10.16–10.18, 10.16*t*, 10.17*f*; **V**:2.7–2.8, 2.11
 Rare-earth-doped fiber lasers, **V**:25.22–25.26, 25.23*t*
 Rare-earth-doped optical fiber amplifiers, **V**:14.1, 14.2–14.4, 14.3*f*
 Raster CRTs, **III**:22.9–22.13, 22.13*f*
 Raster output scanning (ROS) systems, **I**:34.4
 Raster scanning, to despeckle light sources, **III**:5.21
 Rat lenses, **III**:19.7–19.8
 Rating scale judgments, **III**:3.6–3.8, 3.7*f*
 Ray aberrations, **I**:1.87–1.88
 Ray densities, **I**:1.88
 Ray displacement, **II**:3.12
 Ray equation, **I**:1.20
 Ray fans, **I**:1.35
 Ray intercept curves, **II**:2.2–2.4, 3.13
 Ray intercept diagrams, **I**:1.87
 Ray optics, **I**:1.8; **III**:8.15
 Ray paths, **I**:1.10–1.13, 24.2
 Ray sets, **II**:3.20
 Ray tracing, **II**:1.4*f*, 1.4–1.5, 3.11–3.13, 3.12*f*; **III**:19.1
 for binary optics, **I**:23.4, 23.6
 in lighting simulation, **II**:40.20
 nonsequential, **II**:39.6
 in optical design software, **II**:3.11–3.13
 paraxial, **II**:3.5, 3.8–3.9, 3.9*f*
 in systems of revolution, **I**:1.35–1.37, 1.36*f*, 1.40
 for x-ray optics, **V**:35.1–35.6, 35.4*f*, 35.5*f*
 Rayleigh backscattering, **V**:10.8, 14.3
 Rayleigh beacons, **V**:5.29–5.31, 5.31*f*
 Rayleigh criterion, **I**:17.37, 28.6, 28.18; **V**:7.26, 34.1, 40.3, 42.2, 51.3
 Rayleigh criterion of resolving power, **I**:3.26
 Rayleigh index, **I**:8.6
 Rayleigh radius value, **I**:30.10
 Rayleigh range, inverse Compton scattering and, **V**:59.1, 59.2
 Rayleigh range of origin, **I**:5.14, 5.16
 Rayleigh resolution, **I**:30.56, 33.17
 Rayleigh resolution limit, **V**:11.13
 Rayleigh scattering, **I**:7.11, 9.17, 31.30
 in crystals and glasses, **IV**:2.10, 2.27
 in glass, **V**:11.21
 and green flashes, **V**:3.43
 for HMFG fibers, **V**:12.4
 and laser beacons, **V**:5.27, 5.30, 5.32
 and MXRE, **V**:29.5, 29.8*f*
 and optical fibers, **V**:9.4, 9.12
 and Raman fiber amplifiers, **V**:14.9
 and scattering by sea water, **IV**:1.30
 in theory of interaction of light and atmosphere, **V**:3.12, 3.15–3.16
 in third-order optical nonlinearities, **IV**:16.14–16.15
 and x-ray optics, **V**:26.7
 Rayleigh scattering extinction coefficient, **V**:3.16
 Rayleigh units, **III**:1.13, 1.29
 Rayleigh's diffraction integral, **I**:5.13

- Rayleigh-Gans approximation, **I**:7.9, 7.9*f*
 Rayleigh-Rice (RR) approximation, **I**:8.4, 8.9–8.11
 Rayleigh-Sommerfeld diffraction, **I**:3.9, 3.10, 3.23, 3.29
 Rayleigh-wing scattering, **IV**:16.13
 Rays, **I**:1.8–1.13
 axial rays, **II**:1.4, 1.11*f*, 1.12
 chief, **I**:1.75, 17.8, 29.20, 29.37
 for collineation, **I**:1.61
 dashed rays, **II**:1.12, 1.12*f*
 defined, **I**:1.8–1.9
 differential geometry of, **I**:1.19–1.21
 direction of, **I**:1.10
 edge rays, **II**:39.22, 39.38
 exact, **II**:3.3, 3.11–3.12
 expansions about, **I**:1.16
 fields of, **I**:1.13
 finite, **I**:1.35
 groups of, **I**:1.10
 hamiltonian rays, **II**:3.12
 Hamilton's equations for, **I**:1.21
 in heterogeneous media, **I**:1.9, 1.18–1.22
 image-forming, **I**:1.74
 images about known rays, **I**:1.43–1.44, 1.44*f*
 invariance properties of, **I**:1.10
 iterated rays, **II**:3.12
 lagrangian rays, **II**:3.12
 in lenses, **I**:1.35
 marginal, **I**:1.75
 meridional, **I**:1.37
 meridional, **I**:1.35
 meridional rays, **II**:3.3
 nodal, **I**:17.16
 normal congruence, **I**:1.10
 ordinary, **II**:3.12
 paraxial, **I**:1.35, 1.75; **II**:3.3
 paraxial matrices for, **I**:1.68, 1.73
 paths of, **I**:1.10–1.13
 principal, **I**:1.75, 17.8
 principal index of, **I**:13.3
 principal rays, **II**:1.4, 1.11*f*, 1.12
 real and virtual, **I**:1.10, 1.35
 reversibility of, **I**:1.9
 skew, **I**:1.35, 1.73
 variational integral of, **I**:1.19
 RC time constant, **II**:26.7*f*, 26.7–26.8
 Reabsorption, of spectral lines, **V**:56.8
 Reactive evaporation, **IV**:7.11
 Reactive force, **IV**:20.8
 Reactive ion etching (RIE), **I**:21.18–21.19; **II**:18.3, 19.39
 Readout, from optical disk data storage, **I**:35.21–35.24
 Readouts (of visible array detectors), **II**:32.12–32.21
 CCD, **II**:32.12–32.20, 32.13*f*, 32.15*f*–32.18*f*
 MOS, **II**:32.20–32.21
 Real pupils, **I**:1.76
 Real rays, **I**:1.10, 1.35
 Real transitions, **IV**:16.4
 Real-space critical objects, **II**:7.2–7.4, 7.3*f*, 7.4*f*
 Real-time holography, **IV**:12.28–12.29, 12.29*f*, 12.30*f*
 Real-time processors, for adaptive optics, **V**:5.34*f*, 5.34–5.35, 5.35*f*
 Rear focal lengths, **I**:1.40
 Rear focal points, **I**:1.40, 1.47
 Rear principal plane, of Gaussian focal lenses, **I**:1.48
 Reasonableness, of layout, **II**:1.13–1.14
 Received images, **I**:1.26
 Receiver operating characteristic (ROC) curves, **III**:3.7*f*, 3.7–3.8
 Receivers:
 avalanche photodiode, **V**:9.8, 9.10–9.11
 digital on-off-keying, **V**:9.9–9.11
 ideal, **V**:9.9
 in OTDM networks, **V**:20.7–20.8
 positive-intrinsic-negative, **V**:9.8–9.10
 in scatterometers, **V**:1.9*f*, 1.9–1.10
 sensitivity of fiber optic, **V**:9.8–9.11
 Receiving surfaces, in imaging, **I**:1.26
 Receptors, image, **V**:31.4, 31.4*f*
 Recessed lighting, **II**:40.13*f*
 Reciprocal linear dispersion, **V**:38.7
 Reciprocity failure, of photographic films, **II**:29.11–29.12
 Reciprocity theorem, coherence theory and, **I**:5.17–5.18
 Reciprocity theorem of optics (Helmholtz), **III**:8.25, 8.27
 Recoil limit, **IV**:20.5
 Recombination:
 combined, **II**:17.3
 exciton, **II**:17.6
 in GaAs, **II**:17.8, 17.9, 17.9*f*
 minority-carrier, **II**:17.2
 nonradiative, **II**:17.2
 radiative, **II**:17.2

- Recombination radiation, **V**:56.10
- Reconfigurability, of WDM networks, **V**:21.12*f*, 21.12–21.13, 21.13*f*
- Reconfigurable optical add/drop multiplexers (ROADMs), **V**:21.12
- Reconstruction of attosecond beating by interference of two-photon transition (RABITT), **II**:21.9
- Recorded images, **I**:1.26
- Recording, of optical disk data, **I**:35.25–35.28, 35.26*f*, 35.27*f*
- Rectangular apertures, **I**:3.19–3.20, 3.20*f*, 3.25, 3.26
- Rectangular polynomials, **II**:11.27–11.28, 11.28*t*, 11.29*t*, 11.36*t*
- Rectilinear distortion correction, **I**:27.6, 27.13*f*, 27.14*f*
- Red eye, cameras and, **I**:25.6
- Red light, and color film, **II**:29.13, 29.13*f*
- Red shift:
 - noncosmological, **I**:5.23
 - in peak gain wavelength, **V**:25.24
- Red to far-red (R/FR) ratio, **III**:8.28
- Rediagonalization, of index ellipsoid equation, **I**:21.10
- Redistribution force, **IV**:20.8
- Redspot Paint and Varnish, **IV**:6.35
- Reduced eye model, **III**:1.36, 1.37, 2.3, 2.4
- Reduction scanners, in linear sensors, **II**:32.21, 32.22*f*
- Reference method (scatterometer calibration), **V**:1.15
- Reference spheres, for wavefronts, **I**:1.86
- Reflectance, **I**:12.17
 - classification of materials by, **II**:35.4*t*
 - defined, **II**:35.4–35.5
 - geometrical definitions of, **II**:35.6*f*
 - and illuminance/luminance, **II**:37.9
 - in interaction of light and atmosphere, **V**:3.21, 3.21*f*
 - at interface of solid, **IV**:8.12
 - measurement of, **II**:35.10*f*–35.12*f*, 35.10–35.13
 - of metals, **IV**:4.5, 4.11, 4.27*t*–4.39*t*, 4.39, 4.40*f*–4.47*f*
 - nomenclature for, **II**:35.5*t*, 35.6*t*
 - nonabsorbing $[AB]^N$ and $[AB]NA$ multilayers, **IV**:7.32–7.34, 7.33*f*–7.35*f*
 - of optical coatings, **IV**:7.12–7.13
 - retinal, **III**:1.11
 - in solids, **IV**:8.22*f*, 8.23
- Reflectance (*Cont.*):
 - spectral, **II**:38.2, 38.17–38.18, 38.18*f*
 - specular, of black coatings, **IV**:6.26*f*
 - standards of, **II**:35.14*t*
 - and surface coatings, **IV**:6.20*t*
 - and transmittance/absorptance, **II**:35.7, 35.8, 35.8*t*
- Reflecting afocal lenses, **I**:18.19–18.21, 18.20*f*
- Reflecting telescopes, **II**:11.4
- Reflection(s):
 - actual/idealized, **II**:35.2*f*
 - from all-dielectric broadband reflectors, **IV**:7.47
- Bragg
 - of crystals, **V**:40.9
 - in interferometers, **V**:63.26, 63.27
 - and linear polarization, **V**:43.2–43.4, 43.3*f*, 43.4*f*, 43.6, 43.8
 - and liquid crystals, **V**:8.10
 - in monochromators, **V**:63.24
 - and multilayers, **V**:41.2
 - simultaneous, **V**:43.7 [*see also* Multiple-beam Bragg diffraction (MBD)]
 - and x-ray absorption spectroscopy, **V**:30.2, 30.4
- of coatings on substrate, **IV**:7.3
- from computer screens, **III**:23.5–23.6
- defined, **II**:35.4
- enhancement of, for filters and coatings, **IV**:7.108–7.109, 7.109*f*–7.110*f*
- in fiber-optic plant tissues, **III**:8.26, 8.28
- filters with low, **IV**:7.104*f*–7.105*f*, 7.104–7.106
- Fresnel, **V**:13.6, 13.53, 17.3, 25.8
- in Gaussian lens systems, **I**:1.55
- grazing-angle, **V**:63.21
- in homogeneous media, **I**:1.25
- magnetoreflexion, **IV**:5.43, 5.44, 5.44*f*
- magnitude of, **IV**:9.2
- measurement of, **IV**:5.62–5.64, 5.63*f*
- Mueller matrices for, **I**:14.21–14.22
- in neutron and x-ray optics, **V**:26.8–26.9, 63.20–63.21, 63.21*f*
- nonnormal-incidence, **I**:12.15–12.18
- and optical performance, **IV**:7.15–7.16, 7.16*f*
- and phase change of nonabsorbing periodic multilayers, **IV**:7.34
- and phase changes, **I**:12.12–12.13
- photoreflexion, **IV**:5.66*t*, 5.67
- and ray tracing, **I**:1.37
- reference, for Fabry-Perot sensors, **V**:24.2

- Reflection(s) (*Cont.*):
 retro-, **V**:13.6–13.7
 sensing, for Fabry-Perot sensors, **V**:24.2
 of systems of revolution, **I**:1.38, 1.39
 total external, **V**:64.1–64.2
 total internal, **I**:13.20, 21.3; **IV**:8.13; **V**:11.2,
 11.3, 11.14, 11.15
 unfolded, **I**:1.32
 upon interfaces separating different media,
III:8.14*f*, 8.14–8.15
 veiling, **II**:40.12
 volumetric, **I**:7.7*f*
- Reflection coatings, **IV**:7.106*f*–7.113*f*, 7.106–7.113
 at angles close to grazing incidence, **IV**:7.111
 enhancement of reflection, **IV**:7.108–7.109,
 7.109*f*–7.110*f*
 metallic reflectors, **IV**:7.106*f*–7.109*f*,
 7.106–7.108
- Reflection coefficient, for optical constants,
IV:5.9–5.10
- Reflection density, of photographic films,
II:29.8
- Reflection filters, **IV**:7.5, 7.5*f*, 7.111*f*–7.113*f*,
 7.111–7.113
- Reflection model, **III**:10.32
- Reflective and catadioptric objectives, **I**:29.1–29.38
 afocal telescope designs
 Cassegrain-Mersenne, **I**:29.9
 Gregorian-Mersenne, **I**:29.12
 three-mirror, **I**:29.29–29.30
 Altnhof, **I**:29.32–29.33
 anastigmatic designs, **I**:29.12–29.13
 aplanatic designs, **I**:29.11–29.13
 Baker-Nunn, **I**:29.22
 Cassegrain designs, **I**:29.6, 29.7
 afocal Cassegrain-Mersenne telescope, **I**:29.9
 dual magnification, **I**:29.9–29.10
 with field corrector and spherical
 secondary, **I**:29.8–29.9
 Houghton-Cassegrain, **I**:29.22–29.23
 Mangin-Cassegrain with correctors, **I**:29.24
 reflective Schmidt-Cassegrain, **I**:29.17
 Schmidt-Cassegrain, **I**:29.16–29.17
 Schmidt-meniscus Cassegrain, **I**:29.21
 with Schwarzschild relay, **I**:29.32
 solid Makutsov-Cassegrain, **I**:29.19
 spherical-primary, with reflective field
 corrector, **I**:29.9
 three-mirror, **I**:29.30
- Reflective and catadioptric objectives (*Cont.*):
 Cook three-mirror, **I**:29.31
 correctors, in designs
 aplanatic, anastigmatic Schwarzschild
 with aspheric corrector plate, **I**:29.13
 Cassegrain with spherical secondary and
 field corrector, **I**:29.8–29.9
 Mangin-Cassegrain with correctors,
I:29.24
 Ritchey-Chretien telescope with two-lens
 corrector, **I**:29.8
 spherical-primary Cassegrain with
 reflective field corrector, **I**:29.9
 three-lens prime focus corrector, **I**:29.10
 Couder, **I**:29.12
 Dall-Kirkham, **I**:29.8
 Eisenburg and Pearson two-mirror, three
 reflection, **I**:29.25
 features of, **I**:29.2–29.5, 29.3*f*–29.5*f*
 field-of-view plots, **I**:29.34–29.35, 29.35*f*,
 29.36*f*
 flat-medial-field designs, **I**:29.11
 Gabor, **I**:29.20
 glass varieties for, **I**:29.2, 29.2*t*
 Herschelian catadioptric, **I**:29.27
 Houghton designs, **I**:29.22–29.23
 Korsch designs, **I**:29.30–29.32, 29.34
 Maksutov designs, **I**:29.19, 29.20
 Mangin designs, **I**:29.7, 29.24
 Mersenne designs, **I**:29.9, 29.12
 Paul designs, **I**:29.28–29.29
 Ritchey-Chretien with two-lens corrector,
I:29.8
 Schmidt designs, **I**:29.14
 Baker super-Schmidt, **I**:29.21
 field-flattened, **I**:29.14–29.15
 reflective, **I**:29.15
 reflective Schmidt-Cassegrain, **I**:29.17
 Schmidt-Cassegrain, **I**:29.16–29.17
 Schmidt-meniscus Cassegrain, **I**:29.21
 Shafer-relayed-virtual, **I**:29.17–29.18, 29.18*f*
 solid, **I**:29.16
 Schwarzschild designs, **I**:29.12–29.13, 29.32
 SEAL, **I**:29.27–29.28
 Shafer designs
 five mirror unobscured, **I**:29.33–29.34
 four mirror unobscured, **I**:29.33
 Shafer relayed virtual Schmidt,
I:29.17–29.18, 29.18*f*
 two-mirror three reflection, **I**:29.25

- Reflective and catadioptric objectives (*Cont.*):
 spherical primaries in designs, **I**:29.9, 29.11, 29.18, 29.30
 for telescopes
 afocal Cassegrain-Mersenne, **I**:29.9
 afocal Gregorian-Mersenne, **I**:29.12
 Ritchey-Chretien, **I**:29.8
 three-mirror afocal, **I**:29.29–29.30
 terminology, **I**:29.36–29.38
 three-mirror designs, **I**:29.28–29.32, 29.34
 Wetherell and Womble three-mirror, **I**:29.31
 Wright, **I**:29.15
 Yolo, **I**:29.26
 Reflective apertures, **III**:5.10
 Reflective compensators, for spherical aberrations, **II**:13.24, 13.24*f*, 13.25
 Reflective LCDs, **V**:8.31*f*, 8.31–8.32
 Reflective Schmidt objective, **I**:29.15
 Reflective Schmidt-Cassegrain objective, **I**:29.17
 Reflective semiconductor optical amplifiers (SOAs), **V**:19.28
 Reflective systems, **I**:1.9
 Reflective-refractive (RX) concentrators, **II**:39.17, 39.17*f*
 Reflectivity:
 at interface of solid, **IV**:8.12
 nonlinear, **IV**:18.6–18.7, 18.7*f*
 of solids, **IV**:8.23
 of Wolter x-ray optics, **V**:47.5
 Reflectivity amplitude, for solids, **IV**:8.15
 Reflectometers, **II**:35.10, 40.52; **IV**:5.62–5.64, 5.63*f*
 Reflectometric sensors, **III**:19.1
 Reflectors:
 Bragg, **V**:13.45
 conic, **II**:39.11, 39.11*f*
 convergent, **II**:39.38*f*, 39.38–39.40, 39.39*f*
 CPC-type (*see* Compound parabolic collectors)
 divergent, **II**:39.8*f*, 39.9*f*, 39.38–39.40
 faceted, **II**:39.39*f*, 39.39–39.41, 39.40*f*
 headlamp, **II**:40.64
 heat, **IV**:7.58, 7.58*f*
 homogeneous/inhomogeneous, **II**:39.39
 involute, **II**:39.11–39.12, 39.12*f*
 and lens array combinations, **II**:39.34–39.37, 39.36*f*, 39.37*f*
 luminaire, **II**:40.45, 40.45*f*
 macrofocal, **II**:39.11
 Reflectors (*Cont.*):
 metal-dielectric, **IV**:7.81–7.82, 7.108, 7.109, 7.110*f*
 modulated grating, **V**:13.36
 multilayer, **IV**:7.39–7.53
 all-dielectric broadband reflectors, **IV**:7.39, 7.40*f*, 7.45*f*–7.47*f*, 7.45–7.47
 coatings for ultrafast optics, **IV**:7.47–7.48, 7.48*f*
 for far-infrared region, **IV**:7.52, 7.52*f*
 graded reflectivity mirrors, **IV**:7.52
 imperfections in, **IV**:7.40–7.43, 7.41*f*–7.43*f*
 for interferometers and lasers, **IV**:7.39*f*–7.40*f*, 7.39–7.40
 narrowband reflection coatings, **IV**:7.43, 7.44*f*
 rejection filters, **IV**:7.48–7.50, 7.49*f*–7.51*f*
 for soft x-ray and XUV regions, **IV**:7.53
 and two-material periodic multilayers theory, **IV**:7.37–7.38, 7.38*f*
 for neutron beams, **V**:64.3, 64.4, 64.4*f*
 omnidirectional, **IV**:9.2
 resonant, **IV**:7.43–7.45, 7.44*f*
 tailored, **II**:39.37–39.39, 39.38*f*
 very low loss, **IV**:7.41–7.42
 Reflexive sensors, **II**:17.34
 Refraction:
 in calcite, **I**:13.2–13.6, 13.4*t*–13.5*t*
 double, **I**:13.2*f*–13.3*f*, 13.2–13.6, 13.4*t*–13.5*t*
 electro-, **V**:13.2
 enhanced, **IV**:14.20
 in Gaussian lens systems, **I**:1.54
 in homogeneous media, **I**:1.24–1.25
 in human eye, **III**:12.3, 16.1
 Mueller matrices for, **I**:14.20–14.21
 in neutron optics, **V**:63.19–63.23, 63.21*f*, 63.23*f*
 nonlinear, **IV**:16.7–16.9
 ray tracing, **I**:1.37
 and retinal image quality, **III**:1.21
 subjective, **III**:12.6–12.7
 in systems of revolution, **I**:1.38, 1.39
 upon interfaces separating different media, **III**:8.14*f*, 8.14–8.15
 in x-ray optics, **V**:26.6–26.8, 26.8*f*
 Refraction gradients, index of, **I**:24.1 [*see also* Gradient index (GRIN) optics]
 Refraction index, **II**:17.34
 Refractive, defined, **III**:23.3
 Refractive ametropia, **III**:20.2

- Refractive compensators, for spherical aberrations, **II**:13.24, 13.24*f*, 13.25
- Refractive error correction, **III**:12.8–12.15, 16.7–16.19
- ametropias, **III**:16.8*f*
- aphakia and pseudophakia, **III**:12.14–12.15
- with contact lenses, **III**:12.11–12.14
- hydrogel, **III**:12.12–12.13
- for presbyopia, **III**:12.13–12.14
- rigid, **III**:12.11–12.12
- with laser ablation, **III**:16.11–16.19
- ablation rate, **III**:16.16–16.18, 16.18*f*
- corneal photoablation, **III**:16.16, 16.17*f*
- refractive surgery modalities, **III**:16.11*f*–16.14*f*, 16.11–16.15
- thermal, photochemical, and photoacoustic effects, **III**:16.18–16.19
- prescribing, **III**:12.8
- with refractive surgery, **III**:12.14, 16.9–16.15
- corneal incisions/implants, **III**:16.9–16.11, 16.10*f*
- intraocular lenses, **III**:16.9
- laser ablation modalities, **III**:16.11*f*–16.14*f*
- modalities for, **III**:16.11–16.15
- with spectacles, **III**:12.9–12.11, 12.10*f*, 12.11*f*
- Refractive errors, **III**:12.1–12.17, 12.3*f*, 12.4*f*, 16.4–16.19, 16.5*f*
- assessment of, **III**:12.5–12.8
- objective tests, **III**:12.5–12.6
- presbyopia, **III**:12.7–12.8
- subjective techniques, **III**:12.6–12.7
- astigmatism, **III**:16.5–16.6
- binocular factors, **III**:12.15–12.17
- aniseikonia, **III**:12.16–12.17
- anisometropia, **III**:12.16–12.17
- convergence and accommodation, **III**:12.15–12.16
- consequences for optical design, **III**:12.18
- correction of (*see* Refractive error correction)
- ocular wavefronts, **III**:16.6–16.7, 16.7*t*
- as problem with computer work, **III**:23.10
- spherical ametropias, **III**:16.5
- types of, **III**:12.4–12.5
- Refractive index (index of refraction), **I**:1.9; **II**:3.6, 34.13; **III**:8.12
- of anisotropic crystals, **IV**:8.8
- for Brewster angle transmission polarizers, **I**:12.21–12.22
- complex, **I**:7.12–7.13, 12.5, 12.6; **IV**:1.16–1.17; **V**:48.1
- Refractive index (index of refraction) (*Cont.*):
- of cornea, **III**:14.5
- of crystals and glasses, **IV**:2.6, 2.8, 8.10
- defined, **III**:19.1
- dispersion formulas for, **IV**:2.21–2.22
- distributed, **I**:24.1
- in dressed atoms, **IV**:14.19–14.20
- and fiber Bragg gratings, **V**:17.2
- in gradient index optics, **I**:24.2–24.3
- in integrated optics, **I**:21.8–21.9
- and Kolmogorov turbulence, **V**:4.7
- in lens of human eye, **III**:1.5*f*, 1.5–1.6
- of liquid crystals, **V**:8.17, 8.18
- of metals, **IV**:4.3, 4.11, 4.12*t*–4.19*t*, 4.21*f*–4.26*f*
- in neutron optics, **V**:63.19–63.20
- of particles in water, **IV**:1.20
- phase velocity, **V**:7.15–7.16, 7.16*f*
- of photonic crystal fibers, **V**:11.9–11.10
- of polarizers, **I**:12.16, 12.18
- in polymeric optics, **IV**:3.6–3.7, 3.6*t*, 3.7*f*
- for rays in heterogeneous media, **I**:1.21–1.22
- of shallow radical gradients, **I**:24.7–24.8
- in solids, **IV**:8.22*f*, 8.23
- structure function of, **V**:5.6–5.7
- and temperature, **IV**:2.24–2.26
- uniformity of, **IV**:2.5
- of water, **IV**:1.16–1.20, 1.18*f*, 1.19*t*–1.20*t* (*see also* Gradient index optics)
- Refractive index spectrum, **IV**:22.6, 22.6*f*, 22.7*f*
- Refractive lens exchange (RLE), **III**:21.18
- Refractive optics, **I**:23.7, 23.8
- Refractive power, defined, **III**:19.1
- Refractive surgery, **III**:1.15, 1.25, 12.14, 16.9–16.15
- change in cornea with, **III**:16.3
- corneal incisions/implants, **III**:16.9–16.11, 16.10*f*
- intraocular lenses, **III**:16.9
- laser corneal procedures, **III**:16.11–16.15
- ablation profiles, **III**:16.14–16.15
- Epi-LASIK, **III**:16.12, 16.13, 16.13*f*
- LASEK, **III**:16.12, 16.13
- LASIK, **III**:16.13–16.14, 16.14*f*
- photorefractive keratectomy (PRK), **III**:16.11*f*, 16.11–16.12, 16.12*f*
- for presbyopic correction, **III**:14.29
- Refractive systems, **I**:1.9

- Refractive x-ray lenses, **V**:37.3–37.11
 applications of, **V**:37.11
 history of, **V**:37.3
 nanofocusing, **V**:37.8–37.11, 37.9*f*, 37.10*f*
 parabolic, **V**:37.4*f*, 37.4–37.8, 37.6*f*, 37.7*f*
- Refresh rate, **III**:23.3
- Region of interest, in human vision, **III**:24.7
- Region of sag (axial gradients), **I**:24.4*f*
- Registered detected point spread function (RDPSF), **V**:44.14, 44.15*f*
- Regular astigmatism, **III**:1.6
- Rejection filters, **IV**:7.48–7.50, 7.49*f*–7.51*f*
- Rejection ratios, of bandpass filters, **IV**:7.77
- Rejection region, for cutoff filters, **IV**:7.56
- Relative directional sensitivity, **III**:9.6
- Relative intensity noise (RIN):
 defined, **III**:18.2
 in fiber optic communication links, **V**:15.2, 15.14
 of laser diodes, **V**:13.18–13.19, 13.19*f*, 13.23
 and optical fibers, **V**:9.11, 9.16
 with SD-OCT, **III**:18.10
- Relative measurements, absolute vs.,
II:34.20–34.21
- Relative spectacle magnification (RSM),
III:20.2, 20.32–20.33
- Relative visual performance (RVP) model,
II:40.5–40.6
- Relativistic effect(s), in strong field interactions:
 with atoms, **IV**:21.19–21.20
 with free electrons, **IV**:21.6–21.8, 21.7*f*
 with overdense plasmas, **IV**:21.52
 self-focusing and self-channeling as,
IV:21.44–21.45, 21.45*f*
 self-phase modulation as, **IV**:21.45*f*,
 21.45–21.46
- Relativistic electron ATI, **IV**:21.20, 21.21*f*
- Relativistic electron beams, strong field interactions with, **IV**:21.9–21.10
- Relativistic suppression of rescattering,
IV:21.20
- Relativity, of conjugate matrices, **I**:1.70
- Relaxation:
 configurational, **V**:2.14–2.17, 2.15*f*–2.18*f*
 longitudinal and transverse, **IV**:11.5
 photorelaxation, **IV**:5.70*f*
- Relaxation oscillation, **II**:16.12, 19.31*f*,
 19.31–19.34; **V**:13.14–13.16
- Relay lenses, **I**:17.10
- Relay trains, in afocal lenses, **I**:18.17*f*,
 18.17–18.19, 18.18*f*
- Rem jet, **II**:30.4
- Remote sensing:
 in atmospheric optics, **V**:3.36–3.40,
 3.37*f*–3.40*f*
 polarimetry and, **I**:15.37–15.38
 quantum, **IV**:23.14–23.15
 in water, **IV**:1.46–1.47
- Remote sensing scanners, **I**:30.2–30.4,
 30.14–30.25
 circular scan, **I**:30.16, 30.18*f*
 compound mirror optics configurations for,
I:30.15*f*, 30.15–30.16
 multiplexed image scanning by, **I**:30.23,
 30.24*f*
 object- and image-space, **I**:30.18–30.23
 parallel-beam, **I**:30.23–30.25, 30.24*f*–30.25*f*
 pushbroom scan, **I**:30.18
 resolution of, **I**:30.6–30.8, 30.7*f*
 rotating wedge, **I**:30.16, 30.17*f*
 single-mirror, **I**:30.14, 30.15*f*
 two-dimensional, **I**:30.18, 30.19*f*
- Remotely processed (RP) photocathodes,
II:31.10, 31.24
- Rendering, low-level vision models in,
III:24.4
- Reorientational Kerr effect, in solids,
IV:16.13–16.14, 16.14*f*
- Repeater spacing, for optical fibers, **V**:9.12–9.13
- Repetition rate coupling, **II**:20.14–20.15, 20.15*f*
- Rescattering effects:
 relativistic suppression of, **IV**:21.20
 in strong field interactions with atoms,
IV:21.18*f*, 21.18–21.19, 21.19*f*
- Rescattering model, semiclassical, **II**:21.3
- Reset gate (RG), **II**:32.14
- Residential lighting, **II**:40.57, 40.58, 40.59*t*
- Residual amplitude modulation (RAM), **II**:22.14
- Residual astigmatism:
 and contact lens power, **III**:20.15
 defined, **III**:20.2
- Resins, **IV**:3.2, 3.11
- Resistive bolometers, **II**:28.10*f*, 28.10–28.11,
 33.9
- Resistive coupling noise, **II**:27.5, 27.6*f*
- Resistivity:
 coatings with high, **IV**:6.56
 of metals, **IV**:4.54*t*
 of polymers, **IV**:3.5

- Resolution:
- angular, **V**:47.2*f*, 47.2–47.3, 47.10–47.11, 47.11*f*
 - of cameras, **I**:25.6*f*, 25.5–25.6
 - defined, **III**:23.3
 - of deflectors, **V**:6.25, 6.29, 6.30
 - field-weighted-average, **V**:44.10
 - of gratings and monochromators, **V**:38.7
 - and information theory, **III**:4.15–4.16
 - of microscopes, **I**:28.17–28.24
 - Airy disk and lateral resolution, **I**:28.17–28.19, 28.18*f*, 28.19*f*
 - depth of field, **I**:28.22–28.23
 - depth of focus, **I**:28.22
 - three-dimensional diffraction pattern, **I**:28.19–28.22, 28.20*f*, 28.21*f*
 - of NPM AOTFs, **V**:6.40
 - and objective optics, **I**:30.33
 - of optical system, **II**:4.6
 - optics of resolving capacity, **III**:4.4*f*, 4.4–4.5
 - contrast-transfer function, **III**:4.5
 - point-spread function, **III**:4.4, 4.5
 - spatial-frequency coordinates, **III**:4.5
 - Rayleigh, **I**:30.56, 33.17
 - in retinal imaging, **III**:15.21
 - of scanners, **I**:30.6–30.14
 - data rates and remote sensing, **I**:30.6–30.8, 30.7*f*
 - input/output scanning, **I**:30.8–30.14, 30.10*f*, 30.10*t*, 30.11*t*, 30.12*f*–30.13*f*
 - of solid-state cameras, **I**:26.15–26.16, 26.16*f*
 - spatial, **V**:26.10–26.11
 - and superresolution, **III**:4.15
 - of telescopes, **V**:4.2–4.3
 - thresholds of, **III**:4.6–4.7, 4.15
 - in x-ray imaging, **V**:62.3–62.5
- Resolution enhancement techniques (RET), for extreme ultraviolet lithography, **V**:34.1
- Resolution limit, **I**:1.80
- Resolving power, photographic, **II**:29.24
- Resonance:
- cyclotron, **IV**:5.12, 5.12*f*, 5.40, 5.47–5.50, 5.48*f*–5.50*f*
 - of PCFs, **V**:11.12
 - photonic, **IV**:22.9–22.13, 22.10*f*, 22.12*f*
 - slow light propagation and atomic, **IV**:22.2–22.9, 22.3*f*, 22.6*f*–22.8*f*
- Resonance absorption, in overdense plasmas, **IV**:21.47*f*, 21.47–21.48
- Resonant cavity light-emitting diodes (RC-LEDs), **V**:13.39, 13.40
- Resonant circuits, electro-optic modulators and, **V**:7.34
- Resonant degenerate four-wave mixing, **IV**:16.28
- Resonant modes (RMs), **IV**:5.17
- Resonant optical feedback, **II**:19.38, 19.38*f*
- Resonant photodetectors, **II**:26.4*f*
- Resonant photodiodes, **V**:13.65
- Resonant *pin* photodiodes, **II**:26.15, 26.15*f*
- Resonant Raman scattering, **IV**:16.15
- Resonant reflectors, **IV**:7.43–7.45, 7.44*f*
- Resonant scanners, **I**:30.41–30.44, 30.43*f*, 30.44*f*
- Resonators, **V**:25.13, 25.16
- Response time:
- defined, **II**:24.12
 - as measure of performance, **III**:3.8
 - of photodetectors, **II**:25.4
 - of photoemissive detectors, **II**:24.40
- Response to direction (gaze control), **III**:13.29–13.30
- Responsive quantum efficiency (*see* Quantum efficiency)
- Responsivity:
- blackbody, **II**:24.10
 - of photodetectors, **II**:24.18, 24.19, 25.4
 - of photoemissive detectors, **II**:24.35*f*, 24.38
 - of pin diodes, **V**:13.65–13.66, 13.66*f*
 - of *pin* photodiodes, **II**:25.8
 - of solid-state cameras, **I**:26.9–26.10
 - spectral, **II**:24.12, 38.3, 38.18–38.19
 - of spectroradiometers, **II**:38.11–38.12, 38.12*f*
- Restimulable phosphor detectors, **V**:60.8, 60.10*t*
- Resting point of accommodation (RPA), **III**:23.3
- Resting positions (eyes), **III**:13.21
- Resting state (accommodation), **III**:1.33
- Restricted exit angle concentrators, **II**:39.18, 39.18*f*
- Retail lighting, **II**:40.55–40.57, 40.56*t*–40.58*t*
- Retardance, **I**:14.6, 15.8
- Retardance modulators, **I**:15.20
- Retardance space, **I**:14.6
- Retardation plates, **I**:12.24–12.27, 12.25*f*, 12.26*f*, 13.43–13.53, 13.43*t*–13.44*t*
 - achromatic, **I**:13.48–13.52, 13.50*f*, 13.53*t*
 - composite, **I**:13.52, 13.53
 - crystalline-quartz, **I**:13.46–13.48
 - defined, **I**:15.8
 - mica, **I**:13.45–13.46

- Retardation plates (*Cont.*):
 quarter-wave and half-wave, **I**:12.24–12.27,
 12.25*f*, 12.26*f*
 rhomb-type, **I**:13.52, 13.53*t*
 variable, **I**:13.53
- Retarder space, **I**:14.14–14.15, 14.15*f*
- Retarders:
 defined, **I**:15.8
 Mueller matrices for, **I**:14.11–14.15, 14.12*t*,
 14.15*f*
 phase, **V**:41.9, 43.6
- Reticles and reticulation, **II**:12.13, 28.11, 28.12
- Retina, **II**:34.37; **III**:1.3*f*, 21.2
 aging-related changes in, **III**:14.9–14.11
 anatomy of, **III**:2.9–2.11
 AO-controlled light delivery to,
III:15.22–15.24
 alignment, **III**:15.23
 in an AO SLO, **III**:15.23
 conventional AO vision systems, **III**:15.23
 to generate aberrations, **III**:15.24
 longitudinal chromatic aberration,
III:15.22
 measuring activity of individual cones,
III:15.24
 transverse chromatic aberration,
III:15.22–15.23
 uses of, **III**:15.23–15.24
 control of alignment in, **III**:8.7
 direction-corresponding points, **III**:13.8
 fovea, **III**:2.24
 injury to, **III**:7.4
 neural pathways in, **III**:2.9, 2.10*f*
 nonfoveal areas of, **III**:2.24
 OCT image of, **III**:18.4*f*
 optic flow fields, **III**:2.38–2.39, 2.39*f*
 physiology of, **III**:2.11–2.12
 pigments of, **III**:14.9–14.11
 and refraction in the eye, **III**:12.3
 and resolving capacity of the eye,
III:4.5–4.6, 4.6*f*
 schematic diagram of, **III**:8.9*f*, 8.10*f*
 UV light damage to, **III**:14.23
- Retina cameras, AO, **III**:15.3, 15.12 (*see also*
 Ophthalmoscopes)
- Retina pigment epithelium (RPE), **III**:8.12
- Retinal burn, **III**:7.7
- Retinal conjugate plane, in Maxwellian viewing,
III:5.5*f*, 5.8–5.9
- Retinal damage, **II**:40.9; **IV**:13.3
- Retinal disorders, photoreceptor orientation/
 realignment after, **III**:8.7
- Retinal disparity, **III**:13.2
- Retinal eccentricity, visual acuity and,
III:4.10–4.11, 4.11*f*
- Retinal illuminance, **II**:34.40–34.42;
III:1.11–1.12
 in free (newtonian) viewing, **III**:5.2–5.4, 5.3*f*
 maximum permissible, **III**:5.18*f*
 in Maxwellian viewing, **III**:5.6–5.7, 6.3–6.5
 in normal viewing, **III**:6.3
 and pupil diameter, **III**:1.8, 1.9
- Retinal image:
 AO ophthalmic applications, **III**:15.16–15.22
 contrast and resolution, **III**:15.21–15.22
 flood-illuminated AO ophthalmoscope,
III:15.16*f*, 15.16–15.17, 15.17*f*
 optical coherence tomography,
III:15.19–15.21, 15.20*f*, 15.21*f*
 scanning laser ophthalmoscope,
III:15.17–15.19, 15.18*f*, 15.19*f*
 light spread in, **III**:4.4–4.5
 relating actual object and, **III**:4.2
- Retinal image disparity, **III**:13.22
- Retinal image quality, **III**:1.12–1.28
 in aberration-free eye, **III**:1.12–1.14,
 1.13*f*, 1.14*f*
 aging-related changes in, **III**:14.11–14.14
 chromatic aberration, **III**:14.14
 intraocular scatter, **III**:14.12
 monochromatic aberrations,
III:14.12–14.14, 14.13*f*
 chromatic aberration, **III**:1.19–1.20
 computing, **III**:2.3
 intraocular scattered light, **III**:1.20–1.21
 lenticular fluorescence, **III**:1.21
 monochromatic ocular aberrations,
III:1.14–1.19
 off-axis, **III**:1.18*f*, 1.18–1.19, 1.26–1.27,
 1.27*f*
 on the visual axis, **III**:1.15–1.18
 in peripheral field, **III**:1.3
 and pupil diameter, **III**:1.8
 variation with field location, **III**:15.5
 on the visual axis, **III**:1.21–1.28
 calculation from aberration data,
III:1.21–1.22
 comparison between methods, **III**:1.23
 effects of aberration correction,
III:1.25–1.26

- Retinal image quality, on the visual axis (*Cont.*):
observed optical performance, **III**:1.23–1.25
ophthalmoscopic (double-pass) methods,
III:1.22–1.23
psychophysical comparison method,
III:1.22
- Retinal imaging, with spectral domain OCT,
III:18.27–18.29, 18.28*f*, 18.29*f*
- Retinal irradiance, **III**:2.4, 7.7–7.8, 7.8*f*
- Retinal layer of rods and cones model of biological waveguides, **III**:8.8–8.9, 8.9*f*, 8.10*f*, 8.12–8.15
assumptions and approximations for, **III**:8.9, 8.12–8.13
electromagnetic validity of, **III**:8.13
- Retinal microscopy:
adaptive optics in, **III**:15.1–15.24
AO-controlled light delivery to the retina,
III:15.22–15.24
control system, **III**:15.12–15.15
history of, **III**:15.2–15.3
imaging of the retina, **III**:15.16–15.22
implementation of, **III**:15.7–15.15
in ophthalmic applications, **III**:15.15–15.24
properties of ocular aberrations,
III:15.4–15.7, 15.5*f*, 15.6*f*
wavefront corrector, **III**:15.9–15.12, 15.10*f*, 15.11*f*
wavefront sensor, **III**:15.8*f*, 15.8–15.9
defined, **III**:15.2
- Retinal nerve fiber layer (RNFL),
III:18.25–18.27
- Retinal neurons, **III**:2.5*f*
ganglion cells, **III**:2.10–2.11
information transmission by, **III**:2.11
- Retinal pigmented epithelium (RPE),
III:14.25, 18.2
- Retinal processing, **III**:2.9–2.12
- Retinal reflectance, **III**:1.11
layers of occurrence for, **III**:1.23
and ophthalmoscopic methods, **III**:1.23
and stray light, **III**:1.20
- Retinal thermal hazard, **II**:36.17
- Retinex theory (Land), **III**:11.71, 11.72
- Retinopathy:
diabetic, **III**:14.1, 14.25–14.26
solar, **III**:7.1, 7.3
- Retinoscopy, **III**:12.5–12.6
- Retrofocus lenses, **I**:27.2 (*see also* Inverted telephoto camera lenses)
- Retroreflection:
of guided light, **V**:13.6–13.7
measurement of, **II**:35.13
- Retro-reflection testing and correction,
I:15.28–15.29, 15.29*f*
- Retroreflectors, **I**:19.3*t*, 19.28, 19.28*f*
- Return-path ellipsometers (RPEs),
I:16.16–16.17, 16.17*f*
- Return-to-zero differential quadrature phase-shift-keying (RZ-DQPSK) format,
V:21.35, 21.35*f*, 21.36*t*, 21.37*t*
- Return-to-zero differential-phase-shift-keying (RZ-DPSK) format, **V**:21.28, 21.34*f*, 21.36*t*, 21.37*t*
- Return-to-zero (RZ) format, **V**:20.8, 20.9*f*, 20.10*f*, 21.16, 21.29*f*, 21.31*f*
- Return-to zero on-off keying (RZ-OOK),
V:21.34, 21.36*t*
- Reverse bias, **II**:26.3
- Reverse telephoto lenses, **I**:17.29, 17.34*f*
- Reverse-proton-exchanged (RPE) PPLN,
IV:17.16
- Reverse-saturable absorption (RSA), **IV**:13.5, 13.6*f*, 13.7
- Reversibility, of rays, **I**:1.9
- Reversing shear interferometers, **II**:13.12, 13.13*f*
- Reversion prisms, **I**:19.3*t*, 19.14, 19.14*f*
- Revolution, systems of, **I**:1.32–1.43
paraxial optics of, **I**:1.37–1.43
ray tracing in, **I**:1.35–1.37
surfaces, **I**:1.32–1.35
unfolded reflections, **I**:1.32
- Rhabdomic photoreceptors, **III**:8.3
- Rhenium, **II**:40.27
- Rhomboidal prisms, **I**:19.3*t*, 19.25, 19.25*f*
- Rhomb-type retardation plates, **I**:13.52, 13.53*t*
- Ribbon-type tungsten filaments, **II**:15.20*f*
- Riccati-Bessel functions, **I**:7.12
- Rician density function, **V**:3.36
- Ridge waveguide (RWG) lasers, **II**:19.8, 19.9*f*;
V:13.6
- Right circular polarizers, Mueller matrices for,
I:14.10*t*
- Right half-wave circular retarders, Mueller matrices for, **I**:14.12*t*
- Right-angle prisms, **I**:19.3, 19.3*t*, 19.4*f*; **II**:12.15, 12.16, 12.16*f*, 12.17*f*
- Right-circularly polarized light,
I:12.27, 12.28*n*

- Rigid contact lenses:
 correction of refractive errors with,
 III:12.11–12.12
 early types of, **III**:20.2, 20.3
 (*see also* Gas permeable contact lenses)
- Rigid-body motions, **V**:45.6–45.7
- Rigidity, of polymers, **IV**:3.3
- Rimless mounting, **III**:12.2, 12.9
- Ring field lens design, **I**:18.22
- Ring flanges, continuous, **II**:6.4f, 6.11
- Ring lasers, **II**:16.29
 with additional Kerr crystal, **II**:20.17f,
 20.17–20.18
 dye, **II**:20.15–20.16
 Ti:sapphire, with saturable absorber,
 II:20.16f, 20.16–20.17, 20.17f
 in two-level system analogy,
 II:20.24, 20.25f
- Ring resonators, **IV**:12.7, 12.8f, 22.11, 22.12f
- Ring topologies, for WDM networks,
 V:21.5f–21.7f, 21.5–21.6
- Rings, aperture, **II**:3.20
- Risley prisms, **I**:19.3t, 19.25f–19.27f,
 19.25–19.27; **II**:12.4, 12.4f
- Ritchey-Chretien objectives, **I**:29.7–29.8
- Ritchey-Chretien primaries, **I**:29.10
- Ritchey-Chretien two-mirror imaging system,
 II:39.17
- RLM lamp, **II**:40.46, 40.46f, 40.47
- RMS (root mean square) contrast, **III**:3.4
- RMS noise, **II**:24.12
- RMS signal, **II**:24.12
- rms-granularity, **II**:29.19–29.21
- Roadway lighting, **II**:40.67, 40.69–40.71
 and disability glare, **II**:40.10
 and discomfort glare, **II**:40.12
 sign lighting, **II**:40.71
 street lighting, **II**:40.69–40.71, 40.70t,
 40.71t
 tunnel lighting, **II**:40.71
- Robertson's correlated color temperature
 calculation, **II**:38.5
- Robustness, of solitons, **V**:22.4–22.5
- Rochon prisms, **I**:13.7, 13.18f, 13.18–13.21,
 13.19f, 13.24
- Rock-salt lattices, **IV**:5.16
- Rod amacrine cells, **III**:2.5f
- Rod bipolar cells, **III**:2.5f
- Rod pathway, **III**:2.9, 2.10, 2.10f
- Rods, **II**:30.15, 30.16f, 34.37, 36.8f, 36.8–36.10,
 36.9f; **III**:2.5f, 10.3
 and age-related scotopic vision changes,
 III:14.15
 alignment of, **III**:8.4
 amacrine cells, **III**:2.10
 and color matching, **III**:10.17
 in color-deficient observers, **III**:10.16
 diameter of, **III**:2.6f
 directional sensitivity of, **III**:8.5
 function of, **III**:2.4
 linear density of, **III**:2.6f
 and maximum saturation color matching,
 III:10.7
 optical standing waves in, **III**:8.17, 8.19f
 optical waveguide properties of, **III**:14.11
 photocurrent responses of, **III**:2.8f
 in retinal layer of rods and cones model of
 biological waveguides, **III**:8.8–8.9, 8.9f,
 8.10f, 8.12–8.15
 spatial distribution of, **III**:2.6
 spectral sensitivities of, **III**:10.18
 time constant of photopigment regeneration,
 III:2.7
- Rome Air Development Center, **II**:7.19
- Ronchi test, **II**:13.3f, 13.3–13.4, 13.4f
- Roof prisms, **I**:19.11, 19.12f (*see also* Amici
 prisms)
- Roof-mirror-lens arrays, **II**:32.21, 32.22f
- Room temperature vulcanizing (RTV) sealing
 compound, **II**:6.4
- Root-mean-square (rms) wavefront error,
 II:4.1, 4.3, 4.7, 4.8; **III**:1.15–1.17, 1.16f
- ROSAT observatory, **V**:44.2, 44.4f, 44.8, 47.5
- Rose model, for x-ray attenuation, **V**:31.2
- Rostock Cornea Module (RCM) microscope,
 III:17.7, 17.9
- Rotating anodes, as x-ray tube sources, **V**:54.12
- Rotating glass block prisms, **II**:12.4, 12.4f
- Rotating light pipe reflectometers, **IV**:5.62, 5.63f
- Rotating retarders, **I**:15.20
- Rotating wave approximation (RWA), **IV**:20.7
- Rotating wedge scanners, **I**:30.16, 30.17f
- Rotating-analyzer ellipsometer (RAE), **I**:16.13,
 16.13f, 16.14
- Rotating-compensator fixed analyzer (RCFA)
 photopolarimeter, **I**:16.14
- Rotating-detector ellipsometer (RODE),
 I:16.14, 16.14f

- Rotating-element photopolarimeters (REPs),
I:16.13
- Rotation sensors, for fiber interferometers,
I:32.14–32.15, 32.15*f*
- Rotational shear interferometers, II:13.12,
13.13*f*
- Rotational spectra, I:10.20–10.22
- Rotationally parabolic profiles, for refractive
x-ray lenses, V:37.4*f*, 37.4–37.8, 37.6*f*, 37.7*f*
- Rotationally symmetric aspheric lenses, II:9.7,
9.7*f*
- Rotationally symmetric lenses, I:1.27,
1.60–1.62, 1.62*f*
- Rotationally symmetric optics:
hard mounting of, II:6.2–6.4, 6.3*f*, 6.4*f*
soft mounting of, II:6.4*f*, 6.4–6.5, 6.5*f*
- Rotationally symmetric systems, I:1.17,
1.89–1.90
- Roto-optic effect, V:6.6
- Roughened aluminum, IV:6.21, 6.22*f*, 6.30*f*
- Roughness:
polycapillary x-ray optics, V:53.4
surface, IV:6.15
- Routers, waveguide grating, I:21.24
- Rowland circle, I:20.5, 20.7, 20.8, 20.10*f*; V:39.5,
39.6
- Rowland spherical grating, V:38.6
- Rubidium, IV:14.17, 14.17*f*
- Rubidium titanyl arsenate (RbTiOAsO₄) (RTA),
IV:17.1
- Rubidium titanyl phosphate (RbTiOPO₄)
(RTP), IV:2.40*t*, 2.46*t*, 2.48*t*, 2.52*t*, 2.57*t*,
2.64*t*
- Ruby lasers, II:16.12, 16.13*f*, 16.32
- Rugate filters, IV:7.49, 7.50
- Rule-of-thumb PID design, II:22.11–22.12
- “Rule-of-thumb” tolerance, II:6.2
- Runs (trial sequences), III:3.9
- Russell Saunders coupling, V:2.11
- Rutile (TiO₂), IV:2.40*t*, 2.45*t*, 2.48*t*, 2.53*t*, 2.58*t*,
2.65*t*, 2.71*t*
- Rydberg constant, I:10.3
- Rydberg energy, V:54.7
- Rydberg states, II:23.21
- Rytov’s transformation, V:5.9
- Rytov’s series of exponential approximations,
I:9.4
- S on 1 annealing, IV:19.3, 19.4
- Saccades, III:1.42, 1.43, 1.43*f*, 13.20
- Saccadic suppression, III:1.43
- Safety, of head-mounted displays, III:25.2
- Safety standards, for laser hazards, III:7.12
- Sag, of surfaces, I:1.32, 1.33*f*
- SAGE II satellite system, V:3.39, 3.40*f*
- Sagittal fans and foci, I:1.35
- Sagittal plane, III:19.1
- Sagittal-focusing geometry, for monochroma-
tors, V:39.6
- Sagnac interferometers, I:21.35, 21.36, 21.36*f*,
32.3–32.4, 32.4*f*; V:20.22
- Sagnac loops, V:17.8
- Salt (NaCl), IV:2.40*t*, 2.44*t*, 2.48*t*, 2.52*t*, 2.57*t*,
2.64*t*, 2.69*t*
- Sample mounts, for scatterometers, V:1.9
- Sampled gratings (SG), V:13.34, 13.34*f*
- Sampled tracking, on optical disks, I:35.16, 35.16*f*
- Sample-measuring polarimeters, I:15.4, 15.13*f*,
15.13–15.14, 15.16–15.17
- Sampling:
in OTDM networks, V:20.1, 20.4–20.6, 20.5*f*,
20.6*f*
with solid-state cameras, I:26.16–26.19,
26.17*f*–26.19*f*
- Sandblasted aluminum, IV:6.45*f*
- SAOBIC processor, I:11.20
- Sapphire (Al₂O₃), V:12.3*t*, 12.9–12.10, 12.10*f*
dispersion formulas for, IV:2.60*t*
elastic constants of, IV:2.46*t*
infrared spectrum of, IV:2.13, 2.13*f*
lattice vibration model parameters for, IV:2.76*t*
mechanical properties of, IV:2.47*t*
optical modes of, IV:2.70*t*
optical properties of, IV:2.55*t*
properties of, IV:2.38*t*
thermal properties of, IV:2.50*t*
Ti:sapphire amplifiers, IV:21.5
Ti:sapphire lasers, IV:18.3
- Sapphire (Ti:Al₂O₃) lasers, titanium-doped,
II:16.34, 16.34*f*
- Sapphire (Ti:Al₂O₃) ring lasers, titanium-
doped, II:20.16*f*, 20.16–20.17, 20.17*f*
- Sapphire substrate (for HB-LEDs), II:18.2,
18.3, 18.5, 18.6
- Satellite spheres, II:39.26
- Satisloh (company), II:9.4–9.6
- Saturable absorbers, IV:18.5–18.11
fast, IV:18.9–18.10
self-amplitude modulation, IV:18.5–18.7,
18.6*f*, 18.7*f*

- Saturable absorbers (*Cont.*):
 semiconductor saturable absorber mirrors,
IV:18.3, 18.10–18.11
 slow, **IV**:18.7–18.9, 18.8*f*
- Saturable absorption, **IV**:13.5
- Saturable Bragg reflectors (SBRs), **IV**:18.3, 18.11
- Saturated absorption spectroscopy,
I:31.24*f*–31.26*f*, 31.24–31.26
- Saturated colors, **II**:40.7
- Saturated output power, of rare-earth-doped
 amplifiers, **V**:14.4
- Saturation, **II**:40.5, 40.9; **V**:19.9–19.10, 19.10*f*
- Saturation current, of semiconductor diodes,
V:13.69
- Saturation equivalent exposure (SEE),
I:26.10–26.11
- Saturation fluence, **IV**:18.5, 18.6
- Saturation output power, of rare-earth-doped
 amplifiers, **V**:14.3
- Saturation regime, for rare-earth-doped
 amplifiers, **V**:14.3
- Savart plates, **I**:13.56
- Sawing (of LEDs), **II**:17.24, 17.25
- Scalar diffraction theory, for binary optics,
I:23.10–23.13, 23.11*t*, 23.12*f*, 23.13*f*
- Scalar field amplitude, **I**:5.3
- Scaling law:
 or photonic crystal fibers, **V**:11.7
 of spectrum of light, **I**:5.21
- Scan error reduction, **I**:30.48–30.51, 30.49*t*,
 30.50*f*, 30.51*f*
- Scan magnification, **I**:30.5, 30.12–30.14
- Scan rate, of high-resolution deflectors, **V**:6.29
- Scanners, **I**:30.1–30.63; **V**:7.26*f*–7.28*f*, 7.26–7.28
 acousto-optic, **I**:30.44–30.45
 agile beam steering, **I**:30.51–30.63
 decentered lens and microlens arrays,
I:30.57–30.60, 30.58*f*–30.60*f*,
 30.62–30.63
 digital micromirror devices, **I**:30.60–30.61
 gimbal-less two-axis scanning
 micromirrors, **I**:30.61–30.62, 30.62*f*
 phased-array, **I**:30.52–30.57, 30.53*f*,
 30.62–30.63
 electro-optic (gradient), **I**:30.45–30.48,
 30.46*f*–30.48*f*
 error reduction in, **I**:30.48–30.51, 30.49*t*,
 30.50*f*, 30.51*f*
 galvanometer and resonant, **I**:30.41–30.44,
 30.43*f*, 30.44*f*
- Scanners (*Cont.*):
 holographic, **I**:30.38–30.41, 30.40*f*–30.42*f*
 input/output scanning, **I**:30.2, 30.4–30.6,
 30.4*t*, 30.25–30.34
 objective, preobjective, and postobjective,
I:30.28–30.29, 30.29*f*, 30.30*f*
 objective optics, **I**:30.30–30.33,
 30.32*f*–30.33*f*
 power density and power transfer of,
I:30.25–30.28, 30.27*f*, 30.28*f*
 resolution of, **I**:30.8–30.14, 30.10*f*, 30.10*t*,
 30.11*t*, 30.12*f*–30.13*f*
 Keplerian afocal lenses for, **I**:18.13, 18.13*f*
 modulation transfer function (MTF) for, **I**:4.6
 monogon and polygon, **I**:30.34*f*,
 30.34–30.38, 30.35*f*
 remote sensing, **I**:30.2–30.4, 30.14–30.25
 circular scan, **I**:30.16, 30.18*f*
 compound mirror optics configurations,
I:30.15*f*, 30.15–30.16
 multiplexed image scanning, **I**:30.23, 30.24*f*
 object- and image-space, **I**:30.18–30.23,
 30.19*f*–30.23*f*
 parallel-beam, **I**:30.23–30.25, 30.24*f*–30.25*f*
 pushbroom scan, **I**:30.18
 rotating wedge, **I**:30.16, 30.17*f*
 single-mirror, **I**:30.14, 30.15*f*
 two-dimensional, **I**:30.18, 30.19*f*
 resolution of, **I**:30.6–30.14
 data rates and remote sensing, **I**:30.6–30.8
 input/output scanning, **I**:30.8–30.14,
 30.10*f*, 30.10*t*, 30.11*t*, 30.12*f*–30.13*f*
- Scanning, active, **I**:30.4
- Scanning arrays, **II**:33.6, 33.6*f*, 33.14
- Scanning electron microscopy (SEM):
 and magnetron-sputtered MLLs,
V:42.6–42.7, 42.7*f*, 42.8*f*, 42.13, 42.13*f*
 and x-ray fluorescence, **V**:29.2, 29.3
 and x-ray spectral detection, **V**:62.1–62.3,
 62.2*f*
- Scanning FPAs, **II**:33.17, 33.17*f*
- Scanning laser ophthalmoscope (SLO), **III**:15.3,
 15.17, 17.9
 applying adaptive optics to, **III**:15.17–15.19,
 15.18*f*, 15.19*f*, 15.23
 AO-controlled stimulus delivery,
III:15.17–15.19
 visual acuity with, **III**:15.23–15.24
 defined, **III**:15.2

- Scanning slit confocal microscopes, **III**:17.2, 17.3, 17.6*f*, 17.6–17.9, 17.7*f*
- Scanning white light interferometry (SWLI), **II**:10.13
- Scatter function (cosine-corrected BRDF), **V**:1.6
- Scatter rejection, polycapillary x-ray optics and, **V**:53.14–53.16, 53.15*f*–53.16*f*
- Scatterance, of water, **IV**:1.5*t*, 1.10
- Scattered light:
 - with cataracts, **III**:14.24
 - intraocular, **III**:1.20–1.21
 - and OTF calculation, **III**:1.21
 - and psychophysical comparison method, **III**:1.22
- Scattered radiation effect, **II**:34.33
- Scattering:
 - angular distribution of, **IV**:1.12
 - anisotropic, **IV**:12.7
 - anti-Stokes, **V**:10.5
 - and backscattering, **I**:6.5*f*, 6.5–6.7, 9.14*f*, 9.14–9.15
 - Brillouin, **I**:31.30
 - backward, **V**:11.25
 - forward, **V**:11.25, 11.26
 - photonic crystal fibers, **V**:11.25*f*, 11.25–11.26
 - stimulated, **V**:10.1, 10.7–10.9, 15.8, 21.20, 25.6
 - by coated spheres, **I**:7.14
 - coherent and incoherent, **I**:9.2, 9.3
 - Compton, **V**:59.1
 - and circular polarization, **V**:43.6
 - inverse, **V**:59.1
 - and MXRF, **V**:29.5, 29.8*f*
 - and polycapillary x-ray optics, **V**:53.3, 53.15, 53.18
 - and refractive x-ray lenses, **V**:37.5, 37.7
 - and x-ray attenuation, **V**:31.2
 - and x-ray optics, **V**:26.7, 36.1
 - in crystals and glasses, **IV**:2.27
 - by cylinders, **I**:7.14
 - Einstein-Smoluchowski theory of, **IV**:1.30
 - elastic, **V**:63.5
 - in heavy water, **V**:63.10
 - incoherent, **V**:26.7, 31.2, 63.7, 63.8
 - inelastic, **V**:63.3
 - Mie, **I**:7.11, 7.12, 9.17
 - molecular, **I**:7.11
 - neutron, **I**:9.6
- Scattering (*Cont.*):
 - nonlinear, **IV**:13.4*f*, 13.8
 - in optical spectrometers, **I**:31.30–31.31, 31.31*f*
 - by particles, **I**:7.1–7.17
 - in coherent vs. incoherent arrays, **I**:7.2–7.3
 - isotropic homogenous spheres, **I**:7.11–7.14
 - Mie scattering, **I**:7.11, 7.12
 - nonspherical particles, **I**:7.15–7.17
 - regular particles, **I**:7.14–7.15
 - single particles, **I**:7.2–7.3
 - theories of, **I**:7.3–7.4
 - volume scattering vs., **I**:9.2–9.3
 - and photographic film, **II**:29.18
 - and photonic crystal fibers, **V**:11.24–11.26, 11.25*f*
 - and polarization, **I**:15.38*f*–15.40*f*, 15.38–15.39
 - Raman, **I**:31.30–31.31, 31.31*f*; **V**:11.24
 - and atmospheric optics, **V**:3.12, 3.21
 - and configurational relaxation of solids, **V**:2.15–2.17, 2.18*f*
 - Rayleigh, **I**:7.11, 9.17, 31.30; **V**:11.21
 - backscattering, **V**:10.8, 14.3
 - in glass, **V**:11.21
 - and green flashes, **V**:3.43
 - for HMFG fibers, **V**:12.4
 - and laser beacons, **V**:5.27, 5.30, 5.32
 - and MXRF, **V**:29.5, 29.8*f*
 - and optical fibers, **V**:9.4, 9.12
 - and Raman fiber amplifiers, **V**:14.9
 - in theory of interaction of light and atmosphere, **V**:3.12, 3.15–3.16
 - and x-ray optics, **V**:26.7
 - re-, **II**:21.3
 - rescattering, **IV**:21.18*f*, 21.18–21.20, 21.19*f*
 - by silica-air photonic crystal fibers, **V**:11.25
 - by silver halide crystals, **II**:30.5–30.7, 30.6*f*
 - small angle neutron, **V**:64.1
 - spectral, **II**:38.10
 - stimulated, **IV**:16.14–16.19
 - stimulated Raman, **V**:25.6
 - and fiber optic communication links, **V**:15.8
 - and optical fiber amplifiers, **V**:14.2, 14.8
 - in optical fibers, **V**:10.1, 10.4–10.7, 10.5*f*
 - and photonic crystal fiber guidance, **V**:11.23, 11.24, 11.26
 - and WDM networks, **V**:21.20
 - surface, **II**:7.23

- Scattering (*Cont.*):
 from surface, **IV**:6.15
 theory of, **I**:9.3–9.4, 9.4*f*
 Thompson, **V**:26.7
 Thomson, **IV**:21.8–21.9, 21.9*f*
 Thomson backscattering, **V**:59.1
 by water
 inelastic and polarization, **IV**:1.47–1.49,
 1.48*f*, 1.49*f*
 measurement, **IV**:1.29–1.30
 particles, **IV**:1.30–1.35, 1.31*t*, 1.32*f*, 1.33*f*,
 1.34*t*–1.35*t*
 pure water and pure sea water, **IV**:1.30
 wavelength dependence, **IV**:1.35–1.40,
 1.35*t*, 1.36*f*, 1.37*t*, 1.38*t*, 1.39*f*, 1.40*t*
 in water, **V**:63.10, 63.10*t*, 63.11*f*
 x-ray, **I**:9.6
 (*see also specific scattering, e.g.*: Anti-Stokes
 scattering)
 [*see also related topics, e.g.*: Volume (multiple)
 scattering]
- Scattering cross section, **I**:7.4; **V**:63.6*t*, 63.6–63.9,
 63.10*t*
- Scattering force, **IV**:20.7
- Scattering length, **I**:9.6; **V**:63.5–63.9, 63.6*t*,
 63.10*t*
- Scattering length densities, of neutrons,
V:63.9–63.11, 63.11*f*
- Scattering losses, multilayer reflectors and,
IV:7.41
- Scattering matrices, **I**:7.10, 16.8–16.9 (*see also*
 Mueller matrices)
- Scattering planes, **I**:7.9
- Scattering potentials, **I**:9.6
- Scattering rate, **IV**:20.4
- Scattering sensors, **II**:17.34
- Scatterometers, **IV**:7.13; **V**:1.3–1.16
 BSDF, **V**:1.8, 1.8*f*
 calibration of, **V**:1.14–1.15
 configurations and components for,
V:1.7–1.11, 1.8*f*–1.10*f*
 error analysis of, **V**:1.15
 incident power measurement of, **V**:1.14–1.15
 instrument signature and quality, **V**:1.11*f*,
 1.11–1.13, 1.13*t*
 measurement issues with, **V**:1.13–**V**:1.14
 in metrology, **V**:1.16
 specifications of, **V**:1.5*f*, 1.5–1.7
- Schawlow-Townes linewidth, **II**:23.18; **V**:13.21
- Scheelite (CaWO_4), **IV**:2.38*t*, 2.45*t*, 2.47*t*, 2.51*t*,
 2.56*t*, 2.61*t*, 2.72*t*, 2.77*t*
- Scheimpflug condition, **I**:1.61, 17.6*f*, 17.7, 18.4,
 25.18, 25.19*f*
- Scheimpflug rule, **I**:18.4, 18.8
- Schell model sources (of light), **I**:5.11
- Schematic eye model, **III**:1.36, 1.37
- Schiefspiegler objectives, **I**:29.26–29.27
- Schlemm's canal, **III**:14.5
- Schmidt objectives, **I**:29.14
 Baker super-Schmidt, **I**:29.21
 field-flattened, **I**:29.14–29.15
 reflective, **I**:29.15
 reflective Schmidt-Cassegrain, **I**:29.17
 Schmidt-Cassegrain, **I**:29.16–29.17
 Schmidt-meniscus Cassegrain, **I**:29.21
 Shafer-relayed-virtual, **I**:29.17–29.18, 29.18*f*
 solid, **I**:29.16
- Schmidt prisms, **I**:19.3*t*, 19.12, 19.12*f*
- Schmidt-Cassegrain design, **II**:7.20
- Schneider GmbH & Co. KG, **II**:9.4
- Schott Glass Technologies, **IV**:6.57
- Schott glasses, **IV**:2.22, 2.23, 2.26
- Schottky barrier detectors (SBDs), **II**:33.7,
 33.8*f*, 33.12–33.13
- Schottky barriers, **V**:13.58, 13.69
- Schottky contact, **II**:26.3
- Schottky junctions, **II**:26.3
- Schottky photodiodes, **II**:26.16, 26.16*f*, 26.17*f*;
V:13.63, 13.73
- Schroder-Van Laar equation, **V**:8.14
- Schrödinger equation, **I**:10.4; **IV**:2.16, 8.24,
 8.25; **V**:63.4–63.5
- Schrödinger's cat, **IV**:23.1, 23.9
- Schwarzschild arrangement, for McCarthy
 objective, **I**:29.33
- Schwarzschild objectives, **I**:29.12–29.13, 29.32;
V:26.10, 51.1–51.3, 51.2*f*–51.4*f*
- Schwarzschild optics, **V**:44.6, 48.2
- Schwinger formula, **V**:55.15
- Scintigraphy, **V**:32.1, 53.17, 53.18, 53.18*f*
- Scintillation, **V**:3.26, 3.34–3.36, 3.35*f*
- Scintillation x-ray detectors, **V**:60.7–60.8,
 60.9*t*, 60.10*t*
- Scintillator-based flat panel detectors, **V**:61.4,
 61.4*f*, 61.7–61.8, 61.8*f*
- Scintillators, neutron, **V**:63.33
- Sclera, **III**:16.3, 21.2
- Scophony TV projection system, **I**:30.45
- Scorotron, in xerographic systems, **I**:34.2, 34.3*f*

- Scotoma, **III**:7.7
- Scotopic troland, **III**:9.2
- Scotopic vision, **II**:34.37–34.39, 34.38*t*, 36.8*f*, 36.9, 36.9*f*, 36.10, 36.11*f*, 37.2, 40.3; **III**:10.3
 age-related changes in, **III**:14.15, 14.16*f*
 in color-deficient observers, **III**:10.16
 Stiles-Crawford effect, **III**:1.10, 1.10*f*
- Screen reflections, with computer work, **III**:23.5–23.6
- Scribing, **II**:17.24–17.25
- Scully-Lamb master equation, **II**:23.15, 23.17–23.22
 cavity losses, **II**:23.18
 laser master equation, **II**:23.19*f*, 23.19–23.20
 micromaser master equation, **II**:23.20–23.22
- Sea water, **IV**:1.3, 1.18*f*, 1.18–1.21, 1.19*t*–1.20*t*, 1.22*t*, 1.30
- SEAL objective, **I**:29.27–29.28
- Sealed beam lights, **II**:40.26*t*
- Sealed-ampoule diffusion, **II**:17.23
- Second Brewster angle, **I**:12.13
- Secondary electrons, **V**:54.12
- Secondary magnification, **I**:29.12, 29.38
- Secondary sources of light, **I**:5.9–5.10
- Secondary spectrum, **I**:29.7, 29.38
- Secondary spectrum, of radiation, **II**:1.14, 1.15
- Second-harmonic generation (SHG), **III**:17.2, 17.9–17.10
- Second-harmonic interferometers, **I**:32.17–32.18, 32.18*f*
- Second-harmonic processes, conversion efficiencies for, **IV**:10.14–10.16
- Second-order autocorrelator, **II**:21.9
- Second-order nonlinear optics:
 anharmonic oscillator model of susceptibility in, **IV**:10.7–10.9, 10.8*f*
 phase matching in, **IV**:10.12*f*, 10.12–10.14, 10.13*f*
 properties of semiconductors in, **IV**:5.53–5.55
- Second-order susceptibility tensor, **IV**:10.10–10.11
- Second-site adaptation (color vision), **III**:11.17–11.22, 11.18*f*, 11.19*f*
 desensitization by steady fields, **III**:11.17, 11.18, 11.18*f*
 field measurements and, **III**:11.29, 11.31
 and field sensitivities, **III**:11.51, 11.52, 11.53*f*
 habituation, **III**:11.19*f*, 11.19–11.20, 11.21*f*, 11.22*f*
- Second-site desensitization (color vision):
 defined, **III**:11.3
 by steady fields, **III**:11.17, 11.18, 11.18*f*, 11.54*f*
- Sectors (optical disk data), **I**:35.4, 35.6–35.7
- Seek operation, on optical disks, **I**:35.17
- Segmented corrector, **III**:15.2
- Segmented deformable mirrors, **V**:5.37, 5.37*f*
- Segmented piston and piston/tip/tilt wavefront correctors, **III**:15.10, 15.10*f*
- Seidel aberrations, **I**:1.90, 29.38; **III**:12.9
- Seidel (third-order monochromatic) aberrations, **II**:3.9–3.10
- Seidel astigmatism, **II**:11.27, 11.30, 11.35, 11.39, 11.40
- Selective area epitaxy, **I**:21.19–21.20; **V**:13.60
- Selenium, **IV**:2.70*t*
- Self-action effects, **IV**:16.25
- Self-amplified spontaneous emission (SASE) mode, of lasers, **V**:58.1
- Self-amplitude modulation (SAM), **IV**:18.3, 18.5–18.8, 18.6*f*, 18.7*f*
- Self-athermalized, **II**:8.7*f*, 8.7–8.8
- Self-calibration, of silicon photodiodes, **II**:34.29
- Self-centering lens springs, **I**:22.8, 22.8*f*
- Self-channeling, **IV**:21.44–21.45, 21.45*f*
- Self-coherence function, **I**:2.41
- Self-defocusing, **IV**:13.7, 13.8, 19.9–19.11, 19.10*f*
- Self-focusing:
 and laser-induced damage, **IV**:19.5, 19.7
 as relativistic effect in strong field interactions, **IV**:21.44–21.45, 21.45*f*
 thermal, **IV**:13.7–13.8
 of third-order optical nonlinearities, **IV**:16.25
- Self-healing ring networks, for WDM, **V**:21.6, 21.6*f*
- Self-heating, **IV**:17.12
- Self-lensing, **IV**:19.5
- Self-modulated wakefield generation, **IV**:21.41
- Self-motion, **III**:13.8, 13.9
- Selfoc lenses, **I**:24.2, 24.7*f*; **II**:32.21, 32.22*f*
- Self-oscillation, **IV**:12.7–12.9
- Self-phase modulation, **II**:20.6
- Self-phase modulation (SPM), **IV**:21.45*f*, 21.45–21.46 (*see also* Longitudinal Kerr effect)
 in optical fibers, **V**:10.3–10.4
 and solitons, **V**:22.3*f*, 22.3–22.4, 22.4*f*, 22.13
 in WDM networks, **V**:21.18–21.19, 21.19*f*
- Self-protecting limiters, **IV**:13.9

- Self-pumped phase conjugate mirrors (SPPCMs), **IV**:12.7, 12.8*f*–12.9*f*
- Self-saturation effect, of EDFAs, **V**:14.6
- Self-scanned array, **II**:31.2
- Self-steepening, **V**:10.4
- Self-trapped excitons, **IV**:5.26*t*
- Sellaite (MgF_2), **IV**:2.40*t*, 2.45*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.63*t*, 2.71*t*
- Sellenite ($\text{Bi}_{12}\text{SiO}_{20}$) BSO, **IV**:2.38*t*, 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.61*t*, 2.75*t*
- Sellmeier dispersion model, for crystals and glasses, **IV**:2.14, 2.15, 2.21–2.23, 2.25
- Sellmeier formula, for refractive index, **IV**:8.14
- Sellmeier model, for crystals and glasses, **IV**:2.10, 2.11
- Selwyn coefficient, **II**:29.20
- Selwyn's law, **II**:29.20
- Semiconductor(s), **IV**:5.1–5.83
- arrays of, **II**:16.36
 - dielectrics and metals vs., **IV**:8.4
 - direct, **II**:17.4, 17.5*f*
 - and electroabsorption modulators, **V**:13.58
 - electromagnetic spectrum interactions with, **IV**:5.3–5.6, 5.4*f*
 - in electro-optic modulators, **V**:7.34
 - and extreme ultraviolet lithography, **V**:34.1–34.2, 34.2*t*
 - fiber amplifiers vs., **V**:9.13–9.14
 - indirect, **II**:17.4, 17.4*f*, 17.5*f*, 17.6
 - ion-implanted, **IV**:18.21
 - linear optical properties of, **IV**:5.11–5.39, 5.12*f*, 5.13*t*, 5.35*f*–5.37*f*
 - free carriers, **IV**:5.33–5.36
 - impurity and defect absorption, **IV**:5.37–5.39, 5.38*f*, 5.39*f*
 - interband absorption, **IV**:5.21*f*–5.25*f*, 5.21–5.33, 5.26*t*, 5.27*f*–5.28*f*, 5.30*f*–5.34*f*
 - lattice absorption, **IV**:5.13–5.20, 5.15*f*, 5.17*t*, 5.18*f*–5.19*f*, 5.20*t*, 5.21*f*
 - low-temperature, **IV**:18.21
 - magneto-optical properties of, **IV**:5.39–5.52, 5.41*t*
 - effect of magnetic field on energy bands, **IV**:5.40, 5.42, 5.42*f*
 - interband effects, **IV**:5.42–5.46, 5.43*f*–5.45*f*, 5.47*f*
 - intra-band or free-carrier effects, **IV**:5.47–5.52, 5.48*f*–5.51*f*
 - semiconductor nanostructures, **IV**:5.52
- Semiconductor(s) (*Cont.*):
- material systems for, **II**:18.1–18.2
 - materials and applications of, **IV**:5.84*t*–5.86*t*
 - measurement techniques for (*see* Measurement techniques, for semiconductors)
 - nanostructure, **IV**:5.52
 - nonlinear, **IV**:5.52–5.56, 5.54*t*
 - optical/dielectric response in, **IV**:5.8–5.11
 - passive limiting in, **IV**:13.9, 13.10*f*
 - properties of substrates for, **II**:17.20*t*
 - signal processing in complementary metal oxide, **V**:62.5
 - structure of, **IV**:5.6–5.7
 - waveband structure of, **II**:17.3*f*–17.5*f*, 17.3–17.6
- Semiconductor bolometers, **II**:28.4–28.5
- Semiconductor (x-ray) detectors, **V**:29.3, 60.5–60.6, 60.9*t*, 60.10*t*
- Semiconductor Equipment and Materials International (SEMI), **V**:1.4
- Semiconductor interferometric modulators, **V**:13.2, 13.63
- Semiconductor laser amplifier loop optical mirrors (SLALOM), **V**:20.22
- Semiconductor laser amplifiers (SLAs), **V**:9.13–9.14
- Semiconductor laser pumping, **II**:16.19
- Semiconductor lasers, **II**:16.35*f*, 16.35–16.36, 19.1–19.43; **V**:13.1
- applications for, **II**:19.3–19.4
 - arrays of, **II**:19.26–19.29, 19.28*f*, 19.28*t*
 - fabrication and configurations of, **II**:19.6–19.8, 19.9*f*
 - gain mechanism of, **II**:19.4
 - high-power semiconductor lasers, **II**:19.18–19.30
 - arrays, **II**:19.26–19.29, 19.28*f*, 19.28*t*
 - commercial, **II**:19.19–19.23, 19.20*t*, 19.21*f*, 19.22*f*
 - future directions for, **II**:19.23–19.26, 19.25*f*, 19.25*t*, 19.26*t*, 19.27*f*
 - mode-stabilized lasers, **II**:19.18–19.19, 19.19*f*
 - two-dimensional, **II**:19.29–19.30, 19.29*t*, 19.30*f*
 - high-speed modulation of, **II**:19.30–19.36, 19.31*f*–19.36*f*
 - operation of, **II**:19.4*f*–19.6*f*, 19.4–19.6
 - quantum cascade lasers, **II**:16.36

- Semiconductor lasers (*Cont.*):
- quantum well lasers, **II**:19.9–19.18
 - GRIN SCH single, **II**:19.14, 19.14*f*
 - long wavelength, **II**:19.17*f*, 19.17–19.18
 - schematic of, **II**:19.10*f*
 - strained, **II**:19.15–19.17, 19.16*f*
 - threshold modal gain, **II**:19.12*f*, 19.12–19.15, 19.13*f*, 19.15*f*
 - spectral properties of, **II**:19.36–19.39, 19.37*f*, 19.38*f*
 - surface-emitting lasers, **II**:19.39–19.41
 - distributed grating, **II**:19.40*f*, 19.40–19.41
 - integrated laser with 45° mirror, **II**:19.39*f*, 19.39–19.40
 - vertical cavity, **II**:16.36, 19.41, 19.42*f*, 19.43*t*
- Semiconductor optical amplifiers (SOAs), **III**:18.2; **V**:19.1–19.36
- amplification in, **V**:19.1*f*–19.2*f*, 19.2, 19.22*f*–19.26*f*, 19.22–19.27
 - ASE noise, **V**:19.3, 19.4*f*
 - confinement factor, **V**:19.6
 - device characterization, **V**:19.17*f*, 19.17–19.21, 19.19*f*–19.21*f*
 - fabrication, **V**:19.15*f*–19.17*f*, 19.15–19.17
 - gain, **V**:19.4*f*–19.6*f*, 19.4–19.6
 - gain clamping, **V**:19.14*f*, 19.14–19.15
 - gain dynamics, **V**:19.12–19.13, 19.13*f*, 19.14*f*
 - gain ripple and feedback reduction, **V**:19.8, 19.8*f*
 - history of, **V**:19.1
 - material systems, **V**:19.11*f*–19.12*f*, 19.11–19.12
 - noise figure, **V**:19.9
 - nonlinear applications, **V**:19.29*f*, 19.29–19.36, 19.31*f*, 19.33*f*, 19.34*f*
 - and photonic integrated circuits, **V**:19.36
 - polarization dependence, **V**:19.7, 19.7*f*
 - saturation, **V**:19.9–19.10, 19.10*f*
 - switching and modulation, **V**:19.28, 19.28*f*
- Semiconductor photodetectors, **II**:25.2, 38.9, 38.9*t*
- Semiconductor saturable absorber mirror (SESAM), **IV**:18.3, 18.10–18.11; **V**:25.3, 25.32
- Semiconductor saturable absorbers, **IV**:18.4, 18.5, 18.5*f*, 18.10
- Semiconductor ultrafast nonlinearities, **IV**:18.15–18.23
- applications of, **IV**:18.19–18.23
 - and carrier trapping, **IV**:18.21–18.23, 18.22*f*
 - in coherent regime, **IV**:18.19–18.20
- Semiconductor ultrafast nonlinearities (*Cont.*):
- and continuum excitations, **IV**:18.20
 - and excitonic excitations, **IV**:18.19–18.20
 - properties of, **IV**:18.16*f*, 18.16–18.17
 - and pump-probe spectroscopy, **IV**:18.18–18.19
 - in thermalization regime, **IV**:18.20–18.21
 - and transient four-wave mixing, **IV**:18.17*f*, 18.17–18.18
- Semiconductor-doped dielectric films, **IV**:18.11
- Semiconductor-doped glasses, **IV**:18.11
- Semiconductors, complementary metal-oxide, **I**:26.8–26.9
- Semifield angles, of prisms, **I**:13.7
- Semilinear mirrors, **IV**:12.7, 12.8*f*
- Semiquantum four-parameter model, **IV**:2.12
- Semisealed-ampoule diffusion, **II**:17.24
- Sénarmont compensators, **I**:13.53, 28.38
- Sénarmont polariscopes, **I**:12.30
- Sénarmont prisms, **I**:13.7, 13.18, 13.18*f*, 13.21
- Senescent changes in vision (*see* Age-related changes in vision)
- Sensing reflection, for Fabry-Perot sensors, **V**:24.2
- Sensitive area, of photoemissive detectors, **II**:24.41
- Sensitivity:
- defined, **II**:24.12
 - film spectral, **II**:29.11, 29.11*f*
 - film speckle and, to high-energy radiation, **II**:30.19–30.20
 - of interferometers, **II**:13.13–13.14, 13.14*f*
 - of photoemissive detectors, **II**:24.34, 24.35*f*–24.39*f*, 24.41
 - of pin diodes, **V**:13.65–13.66, 13.66*f*
 - quantum, **II**:30.9
- Sensitometry variation, with film processing, **II**:29.10, 29.10*f*
- Sensors:
- active pixel, **I**:26.2, 26.8*f*, 26.8–26.9
 - area arrays of, **II**:32.24–32.32, 32.25*t*
 - about, **II**:32.2
 - frame transfer CCD, **II**:32.26–32.28, 32.27*f*, 32.28*f*
 - image area dimensions for, **II**:32.25*t*
 - interline transfer CCD, **II**:32.28–32.32, 32.29*f*–32.31*f*
 - linear image, **II**:32.2, 32.21–32.24, 32.22*f*, 32.23*f*
 - metal-oxide-semiconductor, **II**:32.25–32.26, 32.26*f*

- Sensors (*Cont.*):
- electro-optic modulators, **V**:7.38
 - extrinsic Fabry-Perot interferometric (EFPI), **V**:24.2*f*, 24.2–24.4, 24.3*f*
 - fiber optic chemical, **V**:12.3*t*
 - generalized, **I**:32.15*f*, 32.15–32.16
 - image, **II**:32.2–32.12, 32.3*f*, 32.21–32.34
 - antiblooming in, **II**:32.9, 32.10*f*
 - color imaging with, **II**:32.32–32.34, 32.33*f*, 32.34*f*
 - dark current in, **II**:32.10–32.12, 32.11*f*
 - junction photodiodes, **II**:32.3–32.6, 32.4*f*, 32.6*f*
 - linear arrays of, **II**:32.21–32.24, 32.22*f*, 32.23*f*
 - MOS capacitors, **II**:32.7–32.8
 - photoconductors, **II**:32.8–32.9
 - pinned photodiodes, **II**:32.8
 - intrinsic Fabry-Perot interferometric (IFPI), **V**:24.4*f*, 24.4–24.5
 - LED detectors in, **II**:17.34
 - multiplexed, **I**:32.16
 - optical fiber, **V**:24.1–24.13
 - comparison of, **V**:24.13
 - extrinsic Fabry-Perot interferometric, **V**:24.2*f*, 24.2–24.4, 24.3*f*
 - fiber Bragg grating, **V**:24.5–24.8, 24.6*f*–24.7*f*
 - intrinsic Fabry-Perot interferometric, **V**:24.4*f*, 24.4–24.5
 - long-period grating sensors, **V**:24.8–24.13, 24.9*f*–24.12*f*, 24.11*t*, 24.13*t*
 - rotation, **I**:32.14–32.15, 32.15*f*
 - Shack-Hartmann, **V**:5.21, 5.36, 5.40–5.43
 - Shack-Hartmann wavefront, **V**:50.5, 50.6*f*
 - staggered linear CCD, **II**:32.23*f*, 32.24
 - time-delay-and-integrate linear, **II**:32.23*f*, 32.24
 - transition-edge, **V**:60.9
- Separate confinement heterostructure (SCH), in fiber optic modulators, **V**:13.4*f*, 13.5
- Separate confinement heterostructure waveguide, **II**:19.24
- Separated absorption, grading, and multiplication layer APDs (SAGM APDs), **II**:26.17, 26.18*f*, 26.20*f*
- Separated absorption and multiplication layer APDs (SAM APDs), **II**:26.3, 26.17
- Sequential stimulus sequencing, **III**:3.9
- Serial byte connection (SBCON) standard, **V**:23.1
- Serial incoherent matrix-vector multipliers, **I**:11.17–11.18, 11.18*f*
- Serial multiplexing, **V**:20.12
- Service temperature, of polymers, **IV**:3.3
- Servo lag, subaperture size and, **V**:5.41–5.42
- Servo stability, **II**:22.8
- Servos, **II**:22.5–22.12
 - Bode representation of, **II**:22.5–22.6, 22.6*f*
 - closed-loop performance, **II**:22.8
 - closed-loop stability issues, **II**:22.8–22.12, 22.9*f*–22.11*f*
 - design with time delay, **II**:22.19–22.20
 - measurement noise not a performance limit, **II**:22.7–22.8
 - phase and amplitude responses vs. frequency, **II**:22.6*f*, 22.6–22.7, 22.7*f*
- Seven-segment LED displays, **II**:17.10, 17.11*f*
- Sewer cameras, **I**:25.22–25.23
- Shack-Hartmann sensors, **V**:5.21, 5.36, 5.40–5.43
- Shack-Hartmann technique, for adaptive optics, **V**:5.23*f*, 5.23–5.27, 5.25*f*, 5.26*f*
- Shack-Hartmann wavefront sensing (SHWS), **III**:15.3, 15.8*f*, 15.8–15.9
- Shack-Hartmann wavefront sensors, **III**:15.2, 16.7; **V**:50.5, 50.6*f*
- SHADOW code, **V**:35.2–35.5, 35.4*f*, 35.5*f*
- Shadowmask color CRTs, **III**:22.4–22.6
 - common problems in, **III**:22.5, 22.6
 - geometry of, **III**:22.4–22.5, 22.5*f*, 22.6*f*
- SHADOWVUI interface, **V**:35.3, 35.5, 35.5*f*
- Shafer objectives:
 - five mirror unobscured, **I**:29.33–29.34
 - four mirror unobscured, **I**:29.33
 - Shafer relayed virtual Schmidt, **I**:29.17–29.18, 29.18*f*
 - two-mirror three reflection, **I**:29.25
- Shake off model (of strong field behavior), **IV**:21.18
- Shallow radial gradient index (SRGRIN), **I**:24.7
- Shane telescope, **V**:5.27, 5.32
- Shape factor:
 - of bandpass filters, **IV**:7.77
 - of lenses, **I**:17.12–17.13, 17.13*f*
- Shapes, projected areas of common, **II**:36.3*f*, 36.3–36.4, 36.3*t*
- Sharpness, of photographic images, **II**:29.18, 29.19, 29.22
- Shear, **III**:13.2
- Shear modulus, for metals, **IV**:4.69*t*
- Shearing interferometers, **I**:32.4, 32.6*f*

- Sheet polarizers, **I**:13.25–13.28
- Shells, mounting of, **II**:6.11, 6.12*f*
- Shenker objective, **I**:29.23
- Shenker objectives, **I**:29.23
- Shielding, in neutron optics, **V**:63.13–63.14
- Shift invariance, **I**:4.2
- Shock specifications, for lenses, **II**:4.10
- Short arc light sources, **II**:15.34, 15.35*f*, 40.39
- Short flint glass, **IV**:2.42*t*
- Short pulse, **IV**:11.7
- Short-exposure images, **V**:4.3, 4.31*f*–4.34*f*, 4.31–4.35, 4.35*t*
- Short-period gratings (SPGs), **V**:24.6
- Short-range motion discrimination mechanism, **III**:2.38
- Short-wavelength infrared (SWIR), **II**:24.3, 33.3, 33.5
- Shot noise, **II**:24.12, 27.3, 27.3*f*, 32.12; **III**:18.9; **V**:13.70, 13.73
 defined, **III**:18.2
 with SD-OCT, **III**:18.10
 of solid-state cameras, **I**:26.11
- Shot-noise limit (SNL), **IV**:23.4–23.6, 23.5*f*
- Shot-noise-limited detection, **III**:18.12–18.13
- Shrinkage, polymer, **IV**:3.14
- SI units, **II**:34.20, 37.7, 37.7*t*
- Side lighting, **II**:40.12
- Sidcar TDI, **II**:33.17, 33.17*f*
- Sidelobe suppression, in AOTFs, **V**:6.41–6.42
- Siemens star, **I**:28.30*f*, 28.47*f*
- Sierpinski Gasket, **I**:8.9
- Sign lighting, **II**:40.71
- Signal, for solid-state cameras, **I**:26.9–26.11, 26.13–26.14
- Signal analysis, **II**:27.12–27.15
 boxcar averaging, **II**:27.13, 27.13*f*
 categories of, **II**:27.3
 gated integration, **II**:27.12–27.13, 27.13*f*
 lock-in amplifiers, **II**:27.13, 27.14, 27.14*f*
 photon counting, **II**:27.14
 selection of technique, **II**:27.14–27.15
 transient photon counting, **II**:27.14
 of unmodulated sources, **II**:27.12
- Signal detection, **II**:27.1–27.12
 and amplifiers, **II**:27.10–27.12, 27.11*f*
 and noise sources, **II**:27.3*f*, 27.3–27.6, 27.5*f*, 27.6*f*
 photomultiplier applications in, **II**:27.6–27.10, 27.7*f*
 technique selection for, **II**:27.2*f*, 27.2–27.3
- Signal detection theory, **III**:2.17, 11.20, 11.23*f*, 11.24
- Signal processing, by wideband AO Bragg cells, **V**:6.30
- Signal-dependent noises, **V**:9.11
- Signal-to-noise ratio (SNR), **II**:22.12, 27.1, 27.3, 29.1, 29.22–29.23, 33.2, 38.10; **III**:18.11
 defined, **III**:18.2
 for electro-optic modulators, **V**:7.16
 for fiber optic communication links, **V**:15.2, 15.4, 15.5
 for in-vivo high speed human retinal imaging, **III**:18.15
 for laser diodes, **V**:13.19
 in neutron optics, **V**:63.13
 optical
 for SOAs, **V**:19.24–19.27
 for WDM networks, **V**:21.20, 21.28, 21.34
 of optical disk data, **I**:35.23
 for optical fibers, **V**:9.9
 and preferred frequency, **III**:2.24
 with SD-OCT, **III**:18.7, 18.10
 of solid-state cameras, **I**:26.13–26.14
 for solitons, **V**:22.1, 22.6–22.8
- Sikkens Aerospace Finishes, **IV**:6.39
- Silica fibers, in lasers, **V**:25.27*t*, 25.28
- Silica glass, fused, **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.67*t*
- Silica-air photonic crystal fibers, **V**:11.14, 11.14*f*, 11.25
- Silicon (Si):
 absorptance of, **IV**:4.48*t*
 crystals of, **IV**:2.40*t*, 2.44*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.64*t*, 2.68*t*
 and diamond turning, **II**:10.5
 doped extrinsic, **II**:33.7, 33.8*f*
 IR absorption due to interstitial oxygen in, **IV**:5.19*f*
 multiphonon absorption of vacuum-grown, **IV**:5.18*f*
 physical properties, **IV**:4.52*t*
 reflectance, **IV**:4.46*f*
 Si:Ga infrared detectors, **II**:24.95, 24.95*f*, 24.96, 24.96*f*
 thermal properties of
 coefficient of linear thermal expansion, **IV**:4.56*t*, 4.58*f*
 elastic stiffness, **IV**:4.69*t*
 moduli and Poisson's ratio, **IV**:4.69*t*

- Silicon (Si), thermal properties of (*Cont.*):
 at room temperature, **IV**:4.55*t*
 specific heat, **IV**:4.65*t*, 4.68*f*
 strength and fracture properties, **IV**:4.70*t*
 thermal conductivity, **IV**:4.58*t*, 4.63*f*–4.64*f*
- Silicon avalanche photodiodes (APDs),
II:24.62–24.65, 24.63*f*–24.66*f*
- Silicon bolometers, **II**:28.7*t*
- Silicon carbide (SiC):
 absorptance of, **IV**:4.49*f*, 4.50*f*
 optical properties of, **IV**:4.19*t*, 4.25*f*
 particles of, **IV**:6.15
 physical properties of, **IV**:4.52*t*
 reflectance of, **IV**:4.41*f*, 4.42*f*, 4.46*f*
 as surface material, **IV**:6.56
 thermal properties of
 coefficient of linear thermal expansion,
IV:4.56*t*, 4.58*f*
 elastic stiffness, **IV**:4.69*t*
 moduli and Poisson's ratio, **IV**:4.69*t*
 at room temperature, **IV**:4.55*t*
 specific heat, **IV**:4.65*t*, 4.68*f*
 strength and fracture properties, **IV**:4.70*t*
 thermal conductivity, **IV**:4.58*t*, 4.63*f*–4.64*f*
- Silicon carbide (SiC) LED devices, **II**:17.18
- Silicon carbide (SiC) substrate (for HB-LEDs),
II:18.2, 18.3
- Silicon carbide (SiC) UV detectors, **II**:24.47,
 24.47*f*
- Silicon CCDs (SCCDs), **II**:33.11*f*, 33.11–33.13,
 33.12*f*, 33.17
- Silicon nitride layer, **II**:17.23
- Silicon oxide (SiO₂) passivation, **II**:18.4
- Silicon oxynitride layer, **II**:17.23
- Silicon (Si) photoconductors, **II**:32.4*f*, 32.31,
 32.32
- Silicon (Si) photodiodes, **II**:38.9, 38.9*t*
 avalanche, **II**:24.62–24.65, 24.63*f*–24.66*f*
 high-quality, **II**:34.30
 light-trap, **II**:34.30
np, **II**:34.30
pin, **II**:24.55*f*–24.57*f*, 24.58–24.61, 24.59*f*,
 24.60*f*
pn, **II**:24.52*f*, 24.55*f*–24.59*f*, 24.55–24.58,
 34.30
 self-calibration of, **II**:34.29
 UV- and blue-enhanced, **II**:24.55*f*, 24.61*f*,
 24.61–24.62, 24.62*f*
- Silicon photonics transmission, **I**:21.14–21.16,
 21.15*f*, 21.38–21.40
- Silicon (Si) photovoltaic detectors,
II:24.54–24.65, 24.55*f*, 24.56*f*
 avalanche photodiodes, **II**:24.62–24.65,
 24.63*f*–24.66*f*
pin photodiodes, **II**:24.55*f*–24.57*f*,
 24.58–24.61, 24.59*f*, 24.60*f*
pn photodiodes, **II**:24.52*f*, 24.55*f*–24.59*f*,
 24.55–24.58
 UV- and blue-enhanced photodiodes,
II:24.55*f*, 24.61*f*, 24.61–24.62, 24.62*f*
- Silicon pore optics, **V**:49.6–49.7, 49.7*f*
- Silicone hydrogel contact lenses, **III**:20.2, 20.3
- Silicon-intensifier-target (SIT) vidicons, **II**:31.8
- Silicon-on-insulator planar waveguide,
IV:22.15
- Silicon-on-insulator (SOI) technology,
I:21.14–21.15, 21.15*f*
- Sillenites (cubic oxides), **IV**:12.17–12.19,
 12.18*t*, 12.19*f*
- Silver:
 absorptance of, **IV**:4.42*f*, 4.48*t*, 4.50*t*, 4.51*t*
 colloidal, **II**:29.13
 optical properties of, **IV**:4.17*t*–4.18*t*, 4.26*f*
 physical properties of, **IV**:4.52*t*, 4.54*t*
 reflectance of, **IV**:4.37*t*–4.38*t*, 4.42*f*
 resistivity of, **IV**:4.54*t*
 thermal properties of
 coefficient of linear thermal expansion,
IV:4.56*t*, 4.57*f*
 elastic stiffness, **IV**:4.69*t*
 moduli and Poisson's ratio, **IV**:4.69*t*
 at room temperature, **IV**:4.55*t*
 specific heat, **IV**:4.65*t*, 4.66*f*
 strength and fracture properties, **IV**:4.70*t*
 thermal conductivity, **IV**:4.58*t*, 4.60*f*–4.61*f*
- Silver gallium sulfide (AgGaS₂) (AGS), **IV**:2.38*t*,
 2.44*t*, 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.74*t*, 2.76*t*
- Silver halide crystals, **II**:30.1, 30.5–30.7, 30.6*f*
- Silver halide fibers, **V**:12.3*t*, 12.8*f*, 12.8–12.9
- Silver halide light detectors, **II**:30.7–30.9, 30.8*f*
- Silver halide surfaces, **II**:29.4, 30.14–30.15, 30.15*f*
- Silver mirrors, **IV**:7.106*f*–7.108*f*, 7.107–7.109
- Silver selenogallate (AgGaSe₂), **IV**:2.38*t*, 2.44*t*,
 2.47*t*, 2.50*t*, 2.55*t*, 2.60*t*, 2.74*t*, 2.76*t*
- Simbol-X spacecraft, **V**:47.10
- Simple cells (cortical neurons), **III**:2.14
- Simple lens, thermal focus shift of, **II**:8.2–8.4,
 8.3*t*, 8.4*t*
- Simplified schematic eye model, **III**:1.36, 1.37
- Simulated annealing, **II**:3.19

- Simultaneous Bragg reflection (multiple-beam Bragg diffraction), **V**:43.6–43.8, 43.7*f*
- Simultaneous measurement, in phase-shifting interferometry, **II**:13.22
- Simultaneous multiple surfaces (SMSs), **II**:39.17, 39.17*f*
- Simultaneous vision contact lenses, **III**:12.13–12.14, 14.28, 20.22
- Sine condition, for stigmatic imaging, **I**:1.30–1.31
- Sine plate, **II**:12.10, 12.11*f*
- Single active electron approximation, **IV**:21.10
- Single attach FDDI nodes, **V**:23.3
- Single capillaries (*see* Monocapillary x-ray optics)
- Single crystal diffraction, **V**:53.12, 53.12*f*–53.13*f*, 53.12–53.14
- Single crystal (SC) fibers, as infrared optical fibers, **V**:12.3*t*, 12.9–12.10, 12.10*f*
- Single element lenses, **I**:17.12–17.17, 17.13*f*–17.16*f*, 17.17*t*
- Single heterojunction LEDs, **II**:17.12, 17.12*f*
- Single isolated pulses, **II**:21.4
- Single lens reflex (SLR) cameras:
 autofocus, **I**:25.12–25.14, 25.13*f*
 features of, **I**:25.9*f*
 formats of, **I**:25.18
 lenses for, **I**:27.1–27.2, 27.3*f*–27.4*f*
 normal lenses for, **I**:27.2, 27.3*f*–27.4*f*
 and time lag, **I**:25.8–25.9
- Single material designs, **II**:6.22, 6.23*f*
- Single mirror resonators, **V**:25.13
- Single molecule high-resolution colocalization (SHREC), **I**:28.23
- Single monochromators, **II**:38.15*f*, 38.16*f*
- Single optical pulse, **II**:20.2–20.3, 20.3*f*, 20.6–20.7
- Single photon emission computed tomography (SPECT), **V**:32.1–32.4
- Single point diamond turning (SPDT), **II**:6.1, 6.20
- Single quantum well (SQW) LEDs, **II**:18.1
- Single scattering, **I**:7.2–7.3
 coherent and incoherent, **I**:9.4–9.7, 9.6*f*
 dynamic, **I**:9.7*f*, 9.7–9.8
 and volume scattering, **I**:9.2–9.8
- Single subband edge enhancement (SSEE) microscopy, **I**:28.29–28.33, 28.30*f*–28.33*f*
- Single speckle, **I**:8.17
- Single-bounce monocapillary x-ray optics, **V**:52.5*f*, 52.5–52.6
- Single-cell-gap transreflective LCDs, **V**:8.34*f*, 8.33–8.35
- Single-channel systems, of SOAs, **V**:19.22*f*–19.24*f*, 19.22–19.24
- Single-component development, in xerographic systems, **I**:34.9, 34.9*f*
- Single-frequency lasers, **II**:19.37*f*
- Single-lens arrays, **II**:39.33–39.34, 39.34*f*
- Single-longitudinal-mode lasers, **II**:19.38
- Single-longitudinal-mode (SLM) lasers, **V**:9.8
- Single-mirror scanners, **I**:30.14, 30.15*f*
- Single-mode excitation technique (SMET), **V**:25.3, 25.20
- Single-mode fibers (SMFs):
 dispersion in, **V**:21.16, 21.16*f*, 21.21, 21.21*f*
 for E-LEDs, **V**:13.40
 and photonic crystal fibers, **V**:11.2, 11.11, 11.13, 11.23–11.25
- Single-mode waveguides, **I**:21.4
- Single-order plates, **I**:13.47
- Single-particle excitation, **IV**:5.81, 5.82*f*
- Single-pass absorption spectroscopy, **IV**:17.24–17.27, 17.25*f*, 17.26*f*
- Single-pass photodetectors, **II**:26.4*f*
- Single-point turning, of polymers, **IV**:3.12
- Single-scattering albedo, **IV**:1.5*t*, 1.7*f*
- Single-shot [*f*]-to-2[*f*] interferometers, **II**:21.6
- Singlet lenses, **I**:17.37–17.38
- Single-use cameras, **II**:30.26
- Singly resonant optical parametric (SOP) oscillators, **IV**:10.18, 10.18*f*
- Singly resonant oscillators (SROs), **IV**:17.2–17.16, 17.3*f*, 17.4*f*
 guided-wave nonlinear structures, **IV**:17.15–17.16
 MgO:sPPLT in, **IV**:17.14–17.15, 17.15*f*
 PPLN crystals in, **IV**:17.4–17.13, 17.6*f*–17.11*f*
 QPM nonlinear materials, **IV**:17.13–17.14
- Singular value decomposition (SVD), **I**:15.25–15.27
- Singular-value-decomposition (SVD) algorithms, **V**:5.25, 5.26
- Sink, in molding, **IV**:3.14
- Sinusoidal ray paths, **I**:24.2
- Sisyphus laser cooling, **IV**:20.19
- Site of limiting noise, **III**:11.24
- Size, perceived, **III**:13.19
- Size-of-source effect, **II**:34.33

- Skew invariant, **I**:1.21, 1.23
 Skew movement, **III**:13.2
 Skew ray limits, **II**:39.20, 39.20*f*
 Skew rays, **I**:1.35, 1.73
 Skewness, **I**:1.23; **II**:40.42
 Skin depth, **I**:7.13
 Skin depth (term), **IV**:4.5
 Skin effect, anomalous, **IV**:21.49
 Skot (unit), **II**:36.7
 Skytubes, **II**:40.50*f*
 Skywells, **II**:40.50*f*
 Slab-coupled optical waveguide amplifiers (SCOWAs), **V**:19.21
 Slater integrals, **V**:2.8
 Slater parameters, **I**:10.12
 Sliding prisms, **II**:12.4, 12.4*f*
 Slit polynomials, **II**:11.30, 11.35*t*, 11.36*t*
 Slit-lamp microscopes, **III**:17.2
 Slit-scanning arrangements, in ophthalmoscopic methods, **III**:1.23
 Sloan notch, **III**:11.34, 11.35*f*, 11.37
 and chromatic adaptation, **III**:11.49, 11.51
 and chromatic discrimination, **III**:11.57, 11.58
 Slope errors, **V**:38.7, 45.7–45.8
 Slope measurement analysis, in x-ray mirror metrology, **V**:46.11–46.12
 Slope of cutoff, **IV**:7.56
 Slope profilometry, **V**:46.3–46.6
 Slot interrupters, **II**:17.34
 Slot-detection x-ray imaging, **V**:61.2, 61.2*f*
 Slow axis, **I**:12.25, 15.8
 Slow light propagation, **IV**:22.1–22.15
 and atomic resonance, **IV**:22.2–22.9, 22.3*f*, 22.6*f*–22.8*f*
 in optical fibers, **IV**:22.13–22.15, 22.14*f*
 and photonic resonance, **IV**:22.9–22.13, 22.10*f*, 22.12*f*
 Slow saturable absorbers, **IV**:18.7–18.9, 18.8*f*
 Slowing, of atomic beams, **IV**:20.11–20.13, 20.12*f*, 20.12*t*, 20.13*f*
 Slowly varying envelope approximation (SVEA), **V**:10.2, 10.6
 Small angle neutron scattering (SANS), **V**:64.1
 Small field (color matching):
 defined, **III**:10.9
 standards for, **III**:10.11, 10.12
 Small-perturbation approximation, for surface scattering, **I**:8.9–8.12
 Small-signal gain coefficient, **II**:16.10
 SMARTCUT technique, **I**:21.14, 21.15
 Smectic phase, of liquid crystals, **V**:8.11*f*, 8.8–8.12, 8.12*f*
 Smith invariant, **I**:1.77 (*see also* Two-ray paraxial invariant)
 Smith reflectors, **I**:28.44
 Smoke detectors, **II**:17.34
 Snakes, effect of laser exposures in, **III**:9.13
 Snellen letters, **III**:4.1, 4.8
 Snell's law, **I**:1.24, 1.38, 12.16, 13.3, 13.5, 13.19; **IV**:7.7, 8.11, 8.23; **V**:19.8, 26.8
 Snell's law, biological waveguides and, **III**:8.14, 8.15
 Snow blindness, **III**:1.9, 7.3, 7.4
 Snubbing, **II**:27.9–27.10
 Society of Automotive Engineers (SAE), **II**:40.2, 40.63–40.64
 Society of Photooptical and Instrumentation Engineers (SPIE), **II**:25.2, 39.12
 Sodium chloride (NaCl), **IV**:2.69*t*
 Soft (hydrogel) contact lenses:
 aberrations and, **III**:20.24
 base curve radius for, **III**:20.3, 20.5
 center thickness of, **III**:20.6
 correction of refractive errors with, **III**:12.12–12.13
 edge thickness of, **III**:20.6
 OAD/OZD of, **III**:20.5
 posterior peripheral curve systems of, **III**:20.6
 power of, **III**:20.15–20.16
 toric lenses
 defined, **III**:20.2
 power of, **III**:20.16–20.17
 Soft mounting, **II**:6.1, 6.4*f*, 6.4–6.5
 Soft x-ray region:
 bandpass filters for, **IV**:7.94–7.96, 7.95*f*–7.96*f*
 gratings and monochromators in, **V**:38.1–38.8
 diffraction properties, **V**:38.1–38.3, 38.2*f*
 dispersion properties, **V**:38.6–38.7
 efficiency, **V**:38.8
 focusing properties, **V**:38.3*f*, 38.3–38.6, 38.4*t*–38.5*t*
 resolution properties, **V**:38.7
 interference polarizers for, **IV**:7.73, 7.76*f*–7.77*f*
 multilayer reflectors for, **IV**:7.42, 7.53
 Soft x-ray telescopes, **V**:50.2
 Soft x-rays, circularly polarized, **V**:55.6–55.7

- Software:
- for lighting simulation, **II**:40.18–40.23
 - for nonimaging modeling, **II**:39.6–39.8
 - for optical design (*see* Optical design software)
 - for stray light suppression, **II**:7.24–7.27
- Software implementation, for adaptive optics, **V**:5.21–5.38
- higher-order wavefront sensing techniques, **V**:5.36–5.37
 - laser beacons, **V**:5.27–5.34, 5.28*f*–5.31*f*, 5.33*f*
 - real-time processors, **V**:5.34*f*, 5.34–5.35, 5.35*f*
 - Shack-Hartmann technique, **V**:5.23*f*, 5.23–5.27, 5.25*f*, 5.26*f*
 - tracking, **V**:5.21–5.23, 5.22*f*
 - wavefront correctors, **V**:5.37–5.38, 5.38*f*
- SOHO observatory, **V**:41.3
- Solano objectives, **I**:29.26
- Solar absorbance, **IV**:6.19, 6.21
- Solar collection, **II**:39.1
- Solar keratitis, **III**:14.23
- Solar light pipes (SLPs), **II**:40.49, 40.51*f*
- Solar retinitis, **III**:7.7
- Solar retinopathy, **III**:7.1, 7.3
- Solar x-ray Imager (SXI), **V**:44.9, 44.12–44.13, 44.16*f*, 44.16–44.17, 44.17*f*
- Solar-cell cover filters, **IV**:7.58, 7.59*f*
- SolarChem, **IV**:6.44, 6.45*f*, 6.48*f*, 6.53*f*
- Soleil compensators, **I**:13.55*f*, 13.55–13.56
- Soleil-Babinet compensators, **I**:35.21*n*
- Sol-gel formed glass, **I**:24.8
- Sol-gels, **IV**:12.27*t*
- Solid angles, **II**:34.9, 34.9*f*, 37.4, 39.5, 39.5*f*
- Solid core photonic crystal fibers:
- attenuation in, **V**:11.20, 11.21
 - Brillouin scattering in, **V**:11.25
 - group velocity dispersion, **V**:11.18, 11.18*f*, 11.19*f*
 - Kerr effects for, **V**:11.23
- Solid lightpipes, **II**:39.30–39.31
- Solid Makutsov-Cassegrain objective, **I**:29.19
- Solid Schmidt objective, **I**:29.16
- Solid spacers, for bandpass filters, **IV**:7.83, 7.88
- Solid state spectroscopy, **I**:10.22–10.26, 10.23*f*–10.27*f*
- Solids:
- band structures and interband transitions of, **IV**:8.24–8.32
 - direct interband absorption, **IV**:8.27–8.28
 - energy band structures, **IV**:8.25–8.27, 8.26*f*
 - excitons, **IV**:8.31–8.32
- Solids, band structures and interband transitions of (*Cont.*):
- indirect transitions, **IV**:8.29*f*, 8.29–8.30
 - joint density of states, **IV**:8.28, 8.29*f*
 - multiphoton absorption, **IV**:8.30–8.31
 - quantum mechanical model, **IV**:8.24–8.25
 - selection rules and forbidden transitions, **IV**:8.28, 8.29
- bound electronic optical Kerr effect in, **IV**:16.12–16.13, 16.13*f*
- configurational relaxation in, **V**:2.14–2.17, 2.15*f*–2.18*f*
- dephasing in, **IV**:14.13–14.14
- dispersion relations in, **IV**:8.14–8.16
- EIT in, **IV**:14.33–14.36, 14.35*f*, 14.36*f*
- extrinsic optical properties of, **IV**:8.3
- free-electron properties of, **IV**:8.21–8.24, 8.22*f*
- intrinsic optical properties of, **IV**:8.3
- lattice interactions in, **IV**:8.16–8.18, 8.17*f*, 8.19*t*–8.20*t*
- optical absorption measurements of, **V**:2.7–2.13, 2.8*f*, 2.10*f*, 2.12*f*
- optical properties of, **IV**:8.1–8.32
- propagation of light in, **IV**:8.4–8.13
- anisotropic crystals, **IV**:8.8–8.11, 8.9*t*, 8.10*f*
 - energy flow, **IV**:8.7–8.8
 - interfaces, **IV**:8.11–8.13, 8.12*f*, 8.13*f*
 - Maxwell's equations, **IV**:8.4–8.6
 - wave equations and optical constants, **IV**:8.6–8.7
- reorientational Kerr effect in, **IV**:16.13–16.14, 16.14*f*
- zero-phonon lines in, **V**:2.13–2.14, 2.14*f*
- Solid-state cameras, **I**:26.1–26.20
- applications, **I**:26.3
 - array performance in, **I**:26.9–26.12, 26.12*f*
 - and charge injection devices, **I**:26.6*f*–26.8*f*, 26.6–26.7
 - and charge-coupled devices, **I**:26.3–26.5, 26.4*f*–26.6*f*
 - complementary metal-oxide semiconductor (CMOS), **I**:26.8–26.9
 - modulation transfer function (MTF) for, **I**:26.14
 - performance metrics for, **I**:26.12–26.16, 26.16*f*
 - sampling with, **I**:26.16–26.19, 26.17*f*–26.19*f*
- Solid-state detectors, **IV**:5.61; **V**:63.33–63.34
- Solid-state ion chambers, **V**:60.5–60.6
- Solid-state lasers, **II**:16.12, 16.13, 16.17–16.18, 16.18*f*, 22.20–22.21; **IV**:18.3

- Solid-state lighting, **II**:18.4*f*, 18.4–18.5
- Solid-state photomultipliers (SSPM), **II**:33.9
- Soliton lasers, **IV**:18.14
- Soliton modelocking, **IV**:18.8*f*, 18.12–18.14
- Soliton solution, **II**:20.5–20.9
- Solitons:
- classical, **V**:22.2*f*–22.4*f*, 22.2–22.4
 - in communication systems, **V**:22.1–22.17
 - dispersion-managed, **V**:22.12–22.15, 22.13*f*, 22.14*f*
 - effects of, **V**:22.1–22.2
 - errors, **V**:22.6
 - and frequency-guiding filters, **V**:22.7–22.9
 - in optical amplifiers, **V**:9.14
 - and optical fibers, **V**:10.1
 - photorefractive, **IV**:12.38
 - properties of, **V**:22.4–22.5
 - self-frequency shift cancellation, **V**:11.24
 - and third-order optical nonlinearities, **IV**:16.25–16.26
 - and transient Raman scattering, **IV**:15.27–15.28, 15.29*f*
 - transmission systems for, **V**:22.5–22.7
 - and wavelength division multiplexing, **V**:22.2, 22.9–22.12, 22.15–22.17
- Soller collimators, **V**:28.2, 28.2*f*
- Soller slits, **V**:26.7
- Solves (term), **II**:3.5
- SOR telescope, **V**:5.27, 5.32, 5.34, 5.38
- Source coupling, **II**:40.41–40.42
- Source debris, in extreme ultraviolet lithography, **V**:34.5
- Source depth measurements, for x-ray tube sources, **V**:54.14, 54.14*f*
- Source diameter, for fiber optics, **II**:17.33
- Source modeling, **II**:39.7
- Source modeling software, **II**:40.19–40.20
- Source modulation, **II**:27.3
- Source spot measurements, for x-ray tube sources, **V**:54.14, 54.14*f*, 54.15, 54.15*f*
- Source-model technique, for photonic crystal fibers, **V**:11.7
- Space, perception of, **III**:13.3–13.7
- binocular cues, **III**:13.7
 - distortion of, **III**:13.16–13.19
 - extraretinal information for eye movements, **III**:13.7
 - kinetic cues, **III**:13.4–13.7, 13.5*f*, 13.6*f*
 - monocular cues, **III**:13.3–13.7
- Space curves, differential geometry of, **I**:1.18–1.19
- Space groups, of crystals, **IV**:2.27–2.28
- Space-bandwidth product, **I**:6.9
- Space-based surfaces, **IV**:6.17–6.18, 6.20*t*, 6.46
- Space-charge region, of pin photodiodes, **V**:13.64
- Spacecraft, **IV**:6.12, 6.16–6.17
- Space-integrating correlator, **I**:11.11
- Spacers, for bandpass filters, **IV**:7.83, 7.88–7.89
- Spaciousness, perception of, **II**:40.5
- Spar cutting, Ahrens method of, **I**:13.12
- Sparrow criterion, **I**:28.18, 28.19
- Spatial channels, in contrast detection, **III**:2.23–2.24
- Spatial coherence, **I**:2.38*f*–2.39*f*, 2.38–2.40, 5.3, 5.5; **V**:4.10
- Spatial coherence length, **V**:27.2
- Spatial dilution, **II**:39.6
- Spatial equichromatic CSFs, **III**:11.45–11.46, 11.46*f*
- Spatial equiluminant CSFs, **III**:11.45–11.46, 11.46*f*
- Spatial fidelity (optical aberrations), **III**:15.4
- Spatial filtering, **I**:11.5–11.6, 11.6*f*
- with confocal microscopes, **III**:17.3
 - to despeckle light sources, **III**:5.21
- Spatial frequency, **III**:4.1
- Spatial frequency channels (in vision), **III**:3.2
- Spatial frequency spectrum, **V**:6.16
- Spatial hole burning, **V**:13.13, 20.2
- Spatial impulse response, of photoreceptors, **III**:8.23
- Spatial light modulation (SLM), **V**:6.4, 6.9
- Spatial location conflicts, in head mounted display systems, **III**:13.33, 13.34*f*
- Spatial multiplexing, **I**:23.8
- Spatial noise, **II**:33.26*f*, 33.26–33.27
- Spatial resolution, **V**:26.10–26.11
- Spatial scale, adaptive optics and, **V**:5.10–5.13, 5.11*t*, 5.12*f*, 5.13*f*
- Spatial sinusoids, **III**:2.20
- Spatial solitons, **IV**:16.26
- Spatial vision, age-related changes in, **III**:14.7, 14.17–14.19, 14.18*f*
- Spatial-frequency content, **I**:4.8
- Spatial-frequency coordinates, resolving capacity of the eye and, **III**:4.5
- Spatially-modulated excimer laser ablation of cornea, **III**:1.25

- Spatio-temporal CSFs, **III**:2.29
- Special-purpose sources (of radiation), **II**:15.53
- Specific heat (*see* Heat capacity)
- Specification, in polymeric optics, **IV**:3.10–3.11
- Specifications, optical (*see* Optical specifications)
- Specimen, for microscopes, **I**:28.22, 28.54–28.55
- Speckle:
 - and coherence theory, **I**:5.22
 - laser, **I**:9.14
 - objective and subjective, **I**:33.9
 - single, **I**:8.17
 - in volume scattering, **I**:9.15*f*, 9.15–9.16
- Speckle effects, **II**:39.32
- Speckle fields, **III**:5.19, 5.21
- Spectacle blur, **III**:12.12
- Spectacle lenses:
 - anisophoria and, **III**:13.26
 - base curve of, **III**:12.1, 12.9
 - effective power of, **III**:20.8, 20.9*t*
 - and gaze control, **III**:13.29–13.30
 - materials for, **III**:12.9–12.10
 - optical power of, **III**:12.4
 - power of, **III**:20.6
- Spectacle magnification (SM), **III**:20.2, 20.31–20.32
- Spectacles, **III**:15.2
 - correction by (*see also* Spectacle lenses)
 - for anisometropia, **III**:12.16, 13.17–13.18
 - aphakia, **III**:12.15
 - for presbyopia, **III**:14.27
 - for refractive errors, **III**:12.9–12.11, 12.10*f*, 12.11*f*
 - defined, **III**:15.2
 - and response to distance, **III**:13.30
- Spectra (code), **V**:55.14
- Spectra Diode Labs, **II**:19.29, 19.29*t*, 19.41
- Spectral (term), **II**:35.2, 35.3
- Spectral absorbance, of water, **IV**:1.9
- Spectral absorption, by detritus in water, **IV**:1.26*f*, 1.26–1.27
- Spectral absorption coefficient:
 - for case 1 waters, **IV**:1.27–1.28
 - for phytoplankton, **IV**:1.24, 1.24*f*, 1.25*f*, 1.25*t*
 - upper bound, in sea water, **IV**:1.21, 1.22*t*
 - for waters, **IV**:1.10, 1.17, 1.20–1.21, 1.27*f*
- Spectral and temporal modalities, of fiber lasers, **V**:25.29–25.33
- Spectral beam attenuation coefficient, for water, **IV**:1.10
- Spectral coherence, **I**:5.5
- Spectral color sensitivities:
 - test method for, **III**:11.34, 11.35*f*, 11.36*f*, 11.37
 - in three-stage Müller zone model, **III**:11.7*f*
 - two-color measurements, **III**:11.34
- Spectral D-double star, **II**:24.12
- Spectral density, **II**:22.3–22.4
- Spectral dependence (of radiometric quantities), **II**:34.9–34.10
- Spectral detectivity, **II**:24.12
- Spectral diffuse attenuation coefficient, for water, **IV**:1.12
- Spectral diffusion, **IV**:11.15
- Spectral (Fourier) domain OCT, **III**:18.5*f*, 18.5–18.7, 18.6*f*, 18.9
 - defined, **III**:18.2
 - noise analysis of, **III**:18.9–18.10, 18.12*f*
 - OFDI vs., **III**:18.9
 - retinal imaging with, **III**:18.27–18.29, 18.28*f*, 18.29*f*
 - sensitivity advantage of, **III**:18.9
- Spectral downwelling average cosine, of water, **IV**:1.12
- Spectral downwelling irradiance, of water, **IV**:1.8
- Spectral downward scalar irradiance, of water, **IV**:1.8
- Spectral D-star, **II**:24.12
- Spectral emittance, **II**:35.7, 35.15
- Spectral emittance, of metals, **IV**:4.6
- Spectral errors, **II**:34.36
- Spectral gain narrowing, in steady-state Stokes scattering, **IV**:15.21
- Spectral hole, **IV**:18.20
- Spectral hole burning, **IV**:22.6, 22.7*f*; **V**:13.13
- Spectral inversion, four-wave mixing and, **V**:10.11
- Spectral irradiance, **II**:38.1–38.2, 38.11*t*, 38.13*f*–38.16*f*, 38.13–38.16
- Spectral irradiance calibration transfer devices, **II**:34.31
- Spectral irradiance lamps, **II**:15.11, 15.12, 15.13*f*
- Spectral irradiance reflectance, of water, **IV**:1.12, 1.46, 1.47*f*
- Spectral lambertian source, **II**:34.17
- Spectral lamps, **II**:15.44, 15.45, 15.45*f*, 15.46*f*, 15.46*t*
- Spectral lines:
 - broadening of
 - Doppler, **V**:3.14, 5.32, 56.5
 - homogeneous, **V**:20.1

- Spectral lines, broadening of (*Cont.*):
 Lorentzian, **V**:3.14
 pressure, **V**:56.5–56.6
 spectral, **V**:56.6–56.8
 from laser generated plasmas, **V**:56.2–56.8
- Spectral luminous efficiency, for photopic vision, **II**:36.8, 36.8*f*, 36.9*f*, 36.16*f*, 37.2
- Spectral net irradiance, of water, **IV**:1.9
- Spectral noise density, **II**:22.3–22.5
- Spectral noise equivalent power, **II**:24.12
- Spectral power density, vector representation of, **III**:10.24, 10.25*f*
- Spectral power distribution, **III**:10.3
 in colorimetric measurements, **III**:10.23
 and tristimulus values, **III**:10.11, 10.36–10.37
- Spectral properties:
 of laser field, **II**:23.28–23.31, 23.30*f*
 of micromaser field, **II**:23.31–23.33
 of semiconductor lasers, **II**:19.36–19.39, 19.37*f*, 19.38*f*
- Spectral radiance, **II**:38.2, 38.11*t*, 38.13*f*–38.16*f*, 38.13–38.16
 thermal, **V**:3.18, 3.20, 3.20*f*
 of water, **IV**:1.6
- Spectral radiance calibration transfer devices, **II**:34.31
- Spectral radiance ribbon filament lamps, **II**:15.11, 15.12*f*
- Spectral radiance units, **II**:34.23–34.24
- Spectral reflectance, **II**:35.4, 35.5, 38.2, 38.17–38.18, 38.18*f*
- Spectral response, of photodetectors, **II**:24.18, 24.19*f*
- Spectral responsivity, **II**:24.12, 38.3, 38.18–38.19
- Spectral scatterance, of water, **IV**:1.10
- Spectral scattering, **II**:38.10
- Spectral scattering coefficient, for water, **IV**:1.10
- Spectral sensitivity, of photographic films, **II**:29.11, 29.11*f*
- Spectral sensitizers, photographic, **II**:30.13–30.18
 about, **II**:30.13–30.14, 30.14*f*
 color science of, **II**:30.15–30.18, 30.16*f*, 30.17*f*
 photophysics of, on silver halide surfaces, **II**:30.14–30.15, 30.15*f*
- Spectral transitions, **I**:10.6–10.7
- Spectral transmission, in polymeric optics, **IV**:3.6
- Spectral transmittance, **II**:35.10, 38.2, 38.17, 38.17*f*; **IV**:7.3–7.4
- Spectral upward plane irradiance, of water, **IV**:1.8
- Spectral upward scalar irradiance, of water, **IV**:1.8
- Spectral upwelling irradiance, of water, **IV**:1.8
- Spectral volume scattering function (VSFs), **IV**:1.37–1.38, 1.38*t*
- Spectrally stray radiation errors, **II**:34.36
- Spectralon, **II**:38.12, 38.13
- Spectrographs, charge, **I**:34.8
- Spectrometers, **I**:31.1–31.31; **IV**:5.59–5.61, 5.60*f*, 5.72, 7.12; **V**:1.14, 63.21*f*
 Bunsen-Kirchhoff, **I**:20.5*f*
 dispersive prisms and gratings for, **I**:20.2–20.3, 20.3*f*
 high-resolution techniques of, **I**:31.23–31.29
 fluorescence line narrowing, **I**:31.29, 31.29*f*
 laser Stark spectroscopy of molecules, **I**:31.27*f*, 31.27–31.29, 31.28*f*
 polarized absorption spectrometers, **I**:31.26
 saturated absorption, **I**:31.24*f*–31.26*f*, 31.24–31.26
 and light scattering, **I**:31.30–31.31, 31.31*f*
 luminescence, **I**:31.5–31.12, 31.8*f*, 31.11*f*
 optical absorption, **I**:31.2–31.5, 31.5*f*
 and photoluminescence decay time, **I**:31.12–31.15, 31.13*f*, 31.14*f*
 polarization, **I**:31.15–31.23
 and optically detected magnetic resonance, **I**:31.21–31.23, 31.22*f*, 31.23*f*
 polarized absorption by, **I**:31.15*f*, 31.15–31.17, 31.18*f*
 polarized absorption/luminescence techniques for, **I**:31.17, 31.19*f*, 31.19–31.21
 prism, **I**:20.2–20.3, 20.3*f*
 Unicam double-monochromator, **I**:20.10, 20.13, 20.15*f*
- Spectrophotometers and spectrophotometry, **II**:34.6, 35.8–35.9, 38.17, 38.17*f*
- Spectrophotometric measurements, **I**:12.15
- Spectropolarimetry and spectropolarimeters, **I**:15.6, 15.8
- Spectroradiometers, **I**:20.1, 20.2*f*, 20.14*t*

- Spectroradiometry, **II**:38.1–38.19
 about, **II**:38.1
 calculations for, **II**:38.3–38.5
 calibration of, **II**:38.11–38.13, 38.11*t*, 38.12*f*
 computer software for, **II**:38.11
 detectors in, **II**:38.8–38.10, 38.9*t*, 38.10*t*
 electronics of, **II**:38.10
 errors in, **II**:38.5–38.6
 figures of merit in, **II**:38.5–38.6
 input (fore-) optics in, **II**:38.7
 monochromators in, **II**:38.7–38.8
 quantities used in, **II**:38.1–38.2
 spectroradiometers, **II**:38.18, 38.18*f*
 system designs in, **II**:38.13–38.19
 spectral irradiance/radiance,
 II:38.13*f*–38.16*f*, 38.13–38.16
 spectral reflectance, **II**:38.17–38.18,
 38.18*f*
 spectral responsivity, **II**:38.18–38.19
 spectral transmittance, **II**:38.17, 38.17*f*
- Spectroscopic ellipsometry (SE), **I**:16.3;
IV:5.66*t*
- Spectroscopic lineshapes, **I**:10.6–10.7,
 10.22–10.27
 in solid state spectroscopy, **I**:10.22–10.26,
 10.23*f*–10.27*f*
 of spectral transitions, **I**:10.6–10.7
- Spectroscopic measurement, **IV**:5.56–5.58
- Spectroscopic transition, rates of, **I**:10.4–10.6
- Spectroscopy, **I**:3.29, 10.1–10.22;
IV:17.21–17.27
 and absorption/photoluminescence of Cr³⁺,
V:2.19*f*, 2.19–2.21, 2.20*f*
 in adaptive optics, **V**:5.19*f*, 5.19–5.21, 5.20*f*
 atomic, **V**:2.4–2.5
 continuous-wave, **IV**:11.2
 defined, **IV**:11.2
 differential transmission, **IV**:18.18–18.19
 evanescent wave, **V**:12.13
 excitation, **V**:2.15*f*, 2.21
 Fourier transform, **V**:2.5, 2.6*f*
 high-resolution Doppler-free, **IV**:17.27
 and homogeneous lineshapes of spectra,
V:2.13–2.17, 2.14*f*–2.18*f*
 knowledge derived from, **IV**:5.6
 Lamb dip, **V**:2.5, 2.7*f*
 Laser-Induced-Breakdown, **V**:3.39
 measurements from, **V**:2.1–2.24
 molecular, **V**:2.5–2.6, 2.6*f*, 2.7*f*
 of multielectron atoms, **I**:10.10–10.11
- Spectroscopy (*Cont.*):
 of one-electron atoms, **I**:10.7–10.9, 10.8*f*,
 10.9*f*
 optical, **V**:2.9–2.10, 2.10*f*
 optical absorption measurements
 of atomic energy levels, **V**:2.2–2.5, 2.4*f*
 of solids, **V**:2.7–2.13, 2.8*f*, 2.10*f*, 2.12*f*
 and outer electronic structure, **I**:10.12–10.16,
 10.13*f*–10.15*f*
 photoacoustic, **IV**:17.22*f*–17.24*f*, 17.22–17.26,
 17.26*f*
 photo-associative, **IV**:20.30–20.31, 20.31*f*
 photon correlation, **I**:9.8
 photon echo, **IV**:11.26–11.27
 polarization, **V**:2.21–2.22, 2.23*f*
 pump-probe, **IV**:18.18–18.19
 Raman, **IV**:7.48, 7.66, 7.83, 7.96, 7.97*f*, 7.98
 rates of spectroscopic transition, **I**:10.4–10.6
 single-pass absorption, **IV**:17.24–17.27,
 17.25*f*, 17.26*f*
 solid state, **I**:10.22–10.26
 theoretical basis, **I**:10.3–10.4
 time-dependent, **IV**:11.2
 time-domain, **V**:7.37, 7.37*f*–7.38*f*
 time-of-flight, **V**:2.5
 of tri-positive rare earth ions, **I**:10.16–10.18,
 10.16*t*, 10.17*f*
 and vibrational and rotational spectra,
I:10.18–10.22, 10.19*f*, 10.21*f*
 x-ray absorption, **V**:30.1–30.5, 30.2*f*–30.5*f*
 x-ray absorption near edge, **V**:30.2, 30.2*f*,
 30.4*f*
 Zeeman, **V**:2.23–2.24, 2.24*f*
- Spectrum(a):
 atomic, **V**:2.13
 homogeneous lineshapes of, **V**:2.13–2.17,
 2.14*f*–2.18*f*
 Kolmogorov, **V**:3.28, 3.31
 of light, **I**:5.19–5.22
 coherence functions for, **I**:5.5–5.6
 coherent mode representation of,
 I:5.20–5.21
 limitations, **I**:5.19, 5.20*f*
 for primary sources, **I**:5.6
 scaling law, **I**:5.21
 Wolf shift, **I**:5.21
 power, **I**:8.12–8.13
 of primary light source, **I**:5.6
 rotational, **I**:10.20–10.22
 secondary, **I**:29.7, 29.38

- Spectrum(a) (*Cont.*):
 spatial frequency, **V**:6.16
 vibrational, **I**:10.18–10.20
 von Kármán, **V**:3.28, 5.6
 and wavefront division, **I**:2.17–2.18, 2.18*f*
 of x-ray tube sources, **V**:54.4–54.10, 54.5*f*,
 54.8*f*
 (*see also* Spectral lines)
- Specular baffles, **IV**:6.14
 Specular black surfaces, **IV**:6.14
 Specular microscopes, **III**:17.2
 Specular reflectance, **II**:35.10, 35.13; **IV**:6.26*f*
 Specular reflection, **III**:23.3
 Specular transmittance, **II**:35.3, 35.9*f*
 Specular vanes, **II**:7.17, 7.17*f*
- Speed:
 of avalanche photodiodes, **V**:13.72
 of LEDs, **II**:17.33
 of photodetectors, **II**:24.20, 24.21
 of pin diodes, **V**:13.67–13.68
 (*see also* Photographic film speed)
- Speed discrimination, **III**:2.37
 SPEOS (optical software), **II**:7.27
- Sphere(s):
 aberrations in, **II**:11.30
 integrating (*see* Integrating spheres)
 nonuniformities with integrating, **II**:39.24*f*,
 39.24–39.26, 39.25*f*
 projected area of, **II**:36.3*f*, 36.3–36.4, 36.3*t*
 scattering by, **I**:7.11–7.14
- Spherelike Brillouin zone, **IV**:9.4
 Spherical aberrations, **I**:1.90, 29.7, 29.38; **V**:45.4
 oblique, **I**:29.15, 29.21, 29.37
 zonal, **I**:29.8, 29.38
- Spherical ametropias, **III**:1.6–1.7, 16.5
 hypermetropia, **III**:16.5
 myopia, **III**:16.5
- Spherical gradients, **III**:19.2–19.3
 Spherical lambertian source, **II**:34.17
 Spherical lenses, **I**:6.3; **II**:39.8, 39.9*f*
 Spherical optics fabrication, **II**:9.4–9.6
 Spherical primaries, in objective designs, **I**:29.9,
 29.11, 29.18, 29.30
- Spherical surfaces, in systems of revolution,
I:1.34
- Spherical waves, **I**:2.4, 2.5*f*, 3.2–3.3
 interference from, **I**:2.11–2.12, 2.12*f*, 2.13*f*
 and plane waves, **I**:2.9–2.11, 2.10*f*
- Spherical-grating monochromators (SGMs),
V:38.3
- Spherochromatism, **I**:24.3–24.6, 29.14; **II**:2.2
 Spherocylindrical lenses, **III**:12.4, 12.5
 Spherometers, **II**:12.18*f*, 12.18–12.19, 12.19*f*,
 12.19*t*
- Spherule, **III**:8.2
 Spin flippers, **V**:63.29
 Spinel (MgAl_2O_4), **IV**:2.39*t*, 2.44*t*, 2.48*t*, 2.52*t*,
 2.57*t*, 2.63*t*, 2.73*t*
- Spire Corporation, **IV**:6.51
 Splay, of liquid crystals, **V**:8.22, 8.23*f*
 Splice losses, for fiber optic communication
 links, **V**:15.7, 15.8*t*
- Splicing, of PCFs, **V**:11.26
 Spline surfaces, **II**:39.6
 Split-aperture scanners, **I**:30.15*f*, 30.15–30.16
 Split-off band, of strained layer quantum well
 lasers, **V**:13.27
- Split-step beam propagation, **IV**:12.10
 Spokes, aperture, **II**:3.20
 Sponges, **III**:8.28–8.29
 Spontaneous decay rate, **I**:10.6
 Spontaneous emission, **IV**:9.8, 20.5; **V**:2.13,
 19.3, 20.1
- Spontaneous emission lasers (*see* Correlated
 emission lasers)
- Spontaneous emission rate, **II**:23.8
 Spontaneous parametric down-conversion
 (SPDC), **IV**:23.10, 23.13
- Spontaneous parametric process,
IV:10.16–10.17, 10.17*f*
- Spontaneous Raman scattering, **IV**:15.3, 15.5
 Spot-diagram analysis, **II**:3.13–3.16
 Spring-8 Compact SASE Source (SCSS), **V**:58.1
 Spring-loaded mountings, **II**:6.13, 6.14*f*
 Spurious-free dynamic range (SFDR), **V**:15.6
 Sputtered beryllium surface, **IV**:6.6*f*
 Sputtered surfaces, **IV**:6.56
 Sputtering method, of manufacturing thin
 films, **IV**:7.11, 7.14
- Square polynomials, **II**:11.30, 11.31*t*–11.34*t*,
 11.36*t*
- Square root law, **III**:2.26
 Square-ended Nicol prisms, **I**:13.16*f*, 13.17
 Square-top multicavity bandpass filters,
IV:7.82*f*–7.88*f*, 7.82–7.83
- Squeezing, in quantum optical interferometry,
IV:23.7
- Squirm, **I**:13.14
 SRI DSRC (company), **II**:19.41
 SRW (code), **V**:55.14

- Stability:
of bandpass filters, **IV**:7.94
light, **II**:30.10
of photodetectors, **II**:24.21, 24.21*f*
of photoemissive detectors, **II**:24.41
- Stability of fixation, **III**:1.44, 1.45*f*
- Stabilization, light, **II**:30.12–30.13 (*see also*
Laser stabilization)
- Stable resonators, **II**:16.23, 16.23*f*
- Stack and draw method, of fiber laser
fabrication, **V**:25.27
- Stacked actuator continuous facesheet
deformable mirrors, **V**:5.37*f*, 5.38, 5.38*f*
- Stadia, **II**:12.2–12.3, 12.3*f*
- Staggered linear CCD image sensor, **II**:32.23*f*,
32.24
- Stagnation, **II**:3.17
- Stainless steel:
physical properties of, **IV**:4.52*t*
thermal properties of, **IV**:4.55*t*–4.58*t*, 4.57*f*
4.62*f*–4.63*f*, 4.65*t*, 4.67*f*, 4.69*t*, 4.70*t*
- Staircase APDs, **II**:26.3
- Standard atmosphere, composition of,
V:3.6–3.11, 3.7*t*, 3.8*f*–3.11*f*
- Standard dispersion-shifted fibers (DSFs),
V:24.10–24.11, 24.11*t*, 24.13*t*
- Standard for the Exchange of Product model
data (STEP), **II**:40.19
- Standard Notation (lenses), **III**:12.5
- Standard observer:
CIE 1931, **III**:10.12
CIE 1964, **III**:10.13
defined, **III**:10.1, 10.4*t*
limitations of, **III**:10.15
standardized CMFs as, **III**:10.9
- Standard rate equation model, **IV**:12.3–12.4
- Standard Test Method (of ASTM), **IV**:6.17
- Standards:
baseline, of radiation sources, **II**:15.9*f*,
15.10*f*, 15.12*f*
for detectors, **II**:38.12–38.13
for infrared radiometry, **II**:15.11–15.12, 15.12*f*
international, **II**:4.11
for length measurements, **II**:12.2
for lighting, **II**:40.19
for lighting system layout and simulation,
II:40.19
for optical image quality, **II**:4.6
published, **II**:4.10
of reflectance, **II**:35.14*t*
- Standards (*Cont.*):
search engine for, **II**:4.11
of spectral transmittance, **II**:35.10
for vehicular lighting, **II**:40.63–40.64, 40.66*f*,
40.66*t*
working, of radiation sources, **II**:15.9–15.13,
15.10*f*, 15.12*f*, 15.13*f*
- Star concentrators, **II**:39.20, 39.21
- Star Instruments, **IV**:6.57
- “Star” patterns (human vision), **III**:1.21
- Star topologies, for WDM networks, **V**:21.5*f*,
21.5–21.6, 21.7*f*
- Starfire Optical Range, **V**:5.23
- Staring arrays, **II**:33.6*f*, 33.6–33.7, 33.14
- Staring FPAs, **II**:33.16–33.17, 33.29, 33.29*t*
- Stark chirped rapid adiabatic passage (SCRAP),
IV:14.1
- Stark effect, **I**:21.11–21.12; **V**:2.21, 20.20, 25.24,
56.6
- Stark levels, of EDFAs, **V**:14.4, 14.5, 14.5*f*
- Stark spectroscopy of molecules, laser,
I:31.27–31.29
- Static gain dynamic, for EDFAs, **V**:21.41, 21.41*f*
- Static Jahn-Teller effect, **V**:2.9
- Stationary anodes, of x-ray tube sources,
V:54.12
- Stationary phase approximation, in diffraction,
I:3.29, 3.31–3.32
- Stationary spectroscopy, **IV**:11.2
- Stationary surfaces, Fresnel-Kirchhoff
approximation for, **I**:8.6–8.8
- Statistical radiometry, **I**:5.22
- Steady-state pulse train, **II**:20.5–20.9
- Steady-state Stokes scattering, **IV**:15.7–15.22
gain coefficients in, **IV**:15.16*t*–15.18*t*
gain narrowing in, **IV**:15.21
photon description in, **IV**:15.21–15.22
pump depletion in, **IV**:15.9, 15.10, 15.15,
15.20, 15.20*f*
Raman linewidths in, **IV**:15.9, 15.10*f*,
15.11*t*–15.20*t*
spectral gain narrowing in, **IV**:15.21
steady-state gain in, **IV**:15.8–15.9
- Steege and Reuter Nicol prisms, **I**:13.17
- Steel, **IV**:4.50*t*, 4.51*t*
- Stefan-Boltzmann law, **II**:34.24, 37.11
- Stellar interferometers, **I**:32.19*f*, 32.19–32.21,
32.20*f*
- Step-and-scan stage, of extreme ultraviolet
lithography, **V**:34.2

- Steradian, **III**:6.3
 Steradian (sr), **II**:36.3, 37.4, 37.4*f*
 Stereo cameras, **I**:25.23–25.24, 25.23*f*
 Stereoacuity, **III**:2.41
 Stereopsis, **III**:1.38–1.42, 2.40
 acuity, **III**:1.39–1.40
 age-related changes in, **III**:14.22
 aniseikonia, **III**:1.41–1.42
 Da Vinci, **III**:13.4
 defined, **III**:13.2
 factors affecting, **III**:2.40, 2.41
 with head mounted visual displays, **III**:13.31, 13.32
 with monocular magnification, **III**:13.13, 13.14
 with monovision, **III**:14.28
 in perception of depth, **III**:13.11*f*, 13.11–13.12
 stereoscopic and related instruments, **III**:1.41
 tolerances in binocular instrumentation, **III**:1.41–1.42
 Stereoscopic instruments, **III**:1.41
 Sterilamps, **II**:15.35, 15.36*f*
 Stern-Gerlach experiment, **V**:63.4
 Stiffness, in crystals and glass, **IV**:2.30, 2.31*t*
 Stigmatic imaging, **I**:1.29–1.31, 1.30*f*
 Stilb (unit), **II**:34.43, 36.7, 36.8*t*
 Stiles and Burch (1955) 2° color-matching functions, **III**:10.12
 Stiles and Burch (1959) 10° color-matching functions, **III**:10.12, 10.13
 Stiles' π -mechanisms, **III**:11.16, 11.16*f*, 11.46, 11.47, 11.47*f*
 and contrast coding, **III**:11.16, 11.16*f*
 and field additivity, **III**:11.51
 limitation of, **III**:11.34
 Stiles-Crawford effect, **III**:1.29, 2.8
 of the first kind, **III**:1.21, 8.5 (*see also* Stiles-Crawford effect of the first kind)
 correction for, **III**:9.1–9.15
 just prior to phototransduction, **III**:6.14
 and Maxwellian viewing, **III**:6.2
 optical, **III**:9.2
 and pupil size, **III**:6.7
 of the second kind, **III**:8.3, 8.5, 9.2
 and human eye, **III**:1.11
 and waveguide models, **III**:8.21
 Stiles-Crawford Effect of the first kind (SCE I/SCE-1), **III**:1.21, 8.5
 and biological waveguides, **III**:8.3, 8.6*f*
 Stiles-Crawford Effect of the first kind (SCE I/SCE-1) (*Cont.*):
 correction for, **III**:9.1–9.15
 adaptive optics techniques for, **III**:9.5
 application of approach, **III**:9.13–9.14
 confounds in, **III**:9.4–9.5
 history of, **III**:9.5
 with non-monochromatic stimulus to vision, **III**:9.5
 Photometric Efficiency factor, **III**:9.3
 sample point-by-point SCE-1 estimates, **III**:9.6–9.13, 9.7*f*, 9.9*f*–9.13*f*
 teleological and developmental factors in, **III**:9.14
 and trolands, **III**:9.2, 9.3
 without adaptive optics, **III**:9.5–9.6
 defined, **III**:9.2
 and human eye, **III**:1.10, 1.10*f*, 1.11
 and photoreceptor directionality, **III**:8.6–8.8
 and waveguide analysis, **III**:8.21
 Stimulated absorption, **II**:16.7–16.8, 16.8*f*
 Stimulated Brillouin scattering (SBS), **IV**:15.43–15.54; **V**:10.1, 10.7–10.9, 15.8, 21.20, 25.6
 equations for, **IV**:15.44–15.48, 15.45*f*
 phase conjugation, **IV**:15.48, 15.52*f*–15.54*f*, 15.52–15.54
 Raman vs., **IV**:15.1
 scattering parameters of materials, **IV**:15.48, 15.49*t*–15.51*t*
 and slow light, **IV**:22.14–22.15
 and stimulated photorefractive scattering, **IV**:12.9
 third-order, **IV**:16.18–16.19
 Stimulated emission, **II**:16.2, 16.7–16.8, 16.8*f*, 16.8–16.9, 23.8; **V**:20.2
 Stimulated emission depletion (STED), **I**:28.24
 Stimulated parametric process, **IV**:10.17–10.18
 Stimulated photon echo, **IV**:11.15–11.19, 11.16*f*–11.19*f*
 Stimulated photoreactive scattering (SPS), **IV**:12.9
 Stimulated Raman adiabatic passage (STIRAP), **IV**:14.1, 14.4, 14.5, 14.7
 Stimulated Raman anti-Stokes scattering, **IV**:15.2*t*
 Stimulated Raman scattering (SRS), **IV**:15.3–15.43; **V**:25.6
 anti-Stokes, **IV**:15.2*t*, 15.32–15.34, 15.33*f*, 15.35*f*
 backward, **IV**:15.41

- Stimulated Raman scattering (SRS) (*Cont.*):
 coherent spectroscopy, **IV**:15.42, 15.42*t*, 15.43*f*
 equations for, **IV**:15.6–15.7
 and fiber optic communication links, **V**:15.8
 focused beams, **IV**:15.41
 formulas for, **IV**:15.20*t*
 geometries for, **IV**:15.4, 15.4*f*
 multiple Stokes generation, **IV**:15.38–15.40, 15.40*f*
 and noise, **IV**:15.35–15.38, 15.39*f*
 and optical fiber amplifiers, **V**:14.2, 14.8
 in optical fibers, **V**:10.1, 10.4–10.7, 10.5*f*
 and photonic crystal fiber guidance, **V**:11.23, 11.24, 11.26
 in plasma, **IV**:21.38*f*, 21.38–21.39
 polarization dependence, **IV**:15.41–15.42, 15.42*t*
 Raman susceptibility, **IV**:15.5–15.6
 and slow light, **IV**:22.15
 steady-state Stokes, **IV**:15.7–15.22
 gain coefficients, **IV**:15.16*t*–15.18*t*
 gain narrowing, **IV**:15.21
 photon description, **IV**:15.21–15.22
 pump depletion, **IV**:15.9, 15.10, 15.15, 15.20, 15.20*f*
 Raman linewidths, **IV**:15.9, 15.10*f*, 15.11*t*–15.20*t*
 spectral gain narrowing, **IV**:15.21
 steady-state gain, **IV**:15.8–15.9
 Stokes, **IV**:15.2*t*
 terminology, **IV**:16.3*t*
 third-order, **IV**:16.15*f*, 16.15–16.18, 16.17*f*
 time-gated imaging, **IV**:15.42, 15.43, 15.44*f*
 transient effects of, **IV**:15.22–15.32
 broadband effects, **IV**:15.28–15.32, 15.29*f*
 phase pulling, **IV**:15.26–15.27, 15.27*f*
 pulsed, **IV**:15.22–15.25, 15.24*f*–15.26*f*, 15.24*t*
 solitons, **IV**:15.27–15.28, 15.29*f*
 spectral properties, **IV**:15.32
 and WDM networks, **V**:21.20
 Stimulated Rayleigh-Wing scattering, **IV**:16.3*t*
 Stimulated scattering, **IV**:16.14–16.19 [*see also* Stimulated Brillouin scattering (SBS); Stimulated Raman scattering (SRS)]
 Stimulus color spaces (colorimetry), **III**:10.11
 Stimulus sequencing, in psychophysical measurement, **III**:3.9
 Stimulus specification (visual acuity), **III**:4.2–4.4, 4.3*f*
 Stochastic optical reconstruction microscopy (STORM), **I**:28.23
 Stokes frequency, **V**:10.8, 21.42
 Stokes intensity, **V**:10.6, 10.7
 Stokes matrix, **I**:14.4 (*see also* Mueller matrix)
 Stokes parameters, **III**:18.20
 and Mueller matrices, **I**:14.4–14.6
 and Poincaré sphere, **I**:14.4–14.6, 14.5*f*, 14.26*f*, 14.33
 Stokes polarimeters, **I**:15.4, 15.5, 15.25
 Stokes power, **V**:14.9
 Stokes scattering, **IV**:15.2–15.3, 15.2*t*, 15.3*t* (*see also* Anti-Stokes scattering; Steady-state Stokes scattering)
 Stokes shift, **IV**:15.1, 15.43; **V**:2.15–2.16
 Stokes shifted Raman scattering, **IV**:16.15, 16.15*f*
 Stokes vectors, **I**:12.14*n*, 12.28; **III**:18.2
 and Mueller matrices, **I**:14.15–14.17, 14.19–14.21
 for nonhomogeneous polarization elements, **I**:14.26*f*
 in polarimetry, **I**:15.3, 15.8–15.13
 for radiative transfer, **I**:9.11
 for speckle patterns, **I**:9.17
 Stokes waves, **IV**:15.1, 15.43; **V**:10.4–10.9, 14.2
 Stokes-Poincaré parameters, for polarization, **V**:43.2, 43.8
 Stop lamps, **II**:40.64*f*, 40.67, 40.68*t*, 40.69*f*
 Stop shift, **I**:1.92
 Stop shifting, **II**:2.5, 2.6*f*
 Stops:
 aperture, **I**:1.74, 1.75*f*, 17.8, 29.5, 29.36; **II**:34.18, 34.19*f*
 field, **I**:1.74, 17.9, 29.5, 29.37; **II**:34.18–34.19, 34.19*f*
 of lenses, **I**:17.8–17.9
 Strabismus, **III**:13.2
 S-trace formula, **I**:1.44
 Straddling springs, **II**:6.13, 6.14*f*
 Straight edges, cylindrical wavefronts and, **I**:3.14*f*, 3.14–3.16, 3.15*f*
 Straightness measurement, **II**:12.10
 Strain:
 in crystals and glasses, **IV**:2.6*t*, 2.30, 2.31*t*, 2.36
 in fiber optic devices, **V**:13.2
 Strained layer quantum well lasers, **V**:13.26*f*, 13.26–13.28, 13.27*f*
 Strained QW lasers, **II**:19.15–19.17, 19.16*f*

- Strained-layer superlattices (SLs), **V**:2.11–2.13, 2.12*f*
- Stratified-medium model (SMM), in ellipsometry, **I**:16.4
- Stray capacitance noise, **II**:27.5, 27.6*f*
- Stray light, **II**:29.15*f*, 29.15–29.16; **IV**:6.10, 6.19
- Stray light suppression, **I**:29.5, 29.5*f*; **II**:7.1–7.32
 about, **II**:7.1–7.2
 aperture placement in, **II**:7.5–7.10
 aperture stops, **II**:7.6–7.7, 7.7*f*, 7.8*f*
 field stops, **II**:7.7, 7.8*f*, 7.9*f*
 Lyot stops, **II**:7.8*f*–7.11*f*, 7.8–7.10
 baffles in, **II**:7.10, 7.11
 and BRDF characteristics, **II**:7.23, 7.24*f*
 Cassegrain design with aperture stop at primary (example), **II**:7.3*f*
 contamination levels in, **II**:7.18–7.19, 7.18*t*, 7.19*f*–7.21*f*
 evaluation methods for, **II**:7.27–7.29, 7.29*f*
 illuminated objects in, **II**:7.5, 7.5*f*, 7.6*f*
 imaged critical objects in, **II**:7.4, 7.5*f*
 information sources on, **II**:7.31–7.32
 issues with, **II**:7.30–7.31
 and point source transmittance definitions, **II**:7.22–7.23
 radiation transfer equation for, **II**:7.21–7.22
 real-space critical objects in, **II**:7.2–7.4, 7.3*f*, 7.4*f*
 software for, **II**:7.24–7.27
 and stray radiation paths, **II**:7.22
 strut design in, **II**:7.20, 7.21, 7.21*f*
 and surface scattering characteristics, **II**:7.23
 vane spacing and depth in, **II**:7.13–7.17
 angle considerations, **II**:7.13–7.16, 7.14*f*, 7.15*f*
 bevel placement, **II**:7.13, 7.14*f*
 depth considerations, **II**:7.16, 7.16*f*, 7.17*f*
 specular vanes, **II**:7.17, 7.17*f*
 vanes in, **II**:7.11–7.12, 7.12*f*, 7.13*f*
- Stray radiation paths, **II**:7.9, 7.22
- Streak cameras, **I**:25.23–25.24, 25.24*f*
- Street lighting, **II**:40.69–40.71, 40.70*t*, 40.71*t*
- Strehl index, **I**:8.6
- Strehl ratio, **III**:1.23–1.24, 1.24*f*, 1.28, 15.2
 and adaptive optics, **V**:5.14, 5.14*f*–5.15*f*, 5.17, 5.17*f*, 5.20–5.22, 5.22*f*, 5.29, 5.29*f*, 5.35, 5.35*f*, 5.39, 5.40, 5.43, 5.43*f*–5.46*f*, 5.46
 and imaging through atmospheric turbulence, **V**:4.13, 4.14*f*, 4.15*t*, 4.27–4.28, 4.28*f*, 4.32–4.33, 4.34*f*, 4.36
- Strength:
 of crystals and glasses, **IV**:2.31–2.32, 2.32*f*, 2.32*t*
 of metals, **IV**:4.8, 4.70, 4.70*t*
 of scattering in water, **IV**:1.12
- Stress, **II**:6.2, 6.21
 in crystals and glasses, **IV**:2.6*t*, 2.30, 2.31*t*
 uniaxial, **IV**:5.66*t*
- Stress tolerance, **II**:6.3
- Stretched segment displays, **II**:17.30*f*, 17.30–17.31, 17.31*f*
- Strip loading, of RWG lasers, **V**:13.6
- Strip mirror integrator (SMI), **II**:39.40
- Stroma (cornea), **III**:14.5
- Strong field approximation, **II**:21.3
- Strong field approximation (SFA), **IV**:21.12
- Strong field double ionization, **IV**:21.18–21.19
- Strong field interactions, **IV**:21.1–21.55
 about, **IV**:21.2–21.3
 with atoms, **IV**:21.10–21.21
 above threshold ionization, **IV**:21.14–21.17, 21.15*f*, 21.16*f*
 ionization stabilization, **IV**:21.20–21.21, 21.22*f*
 Keldysh parameter, **IV**:21.10
 multiphoton and quasi-classical regimes, **IV**:21.10
 multiphoton ionization, **IV**:21.10–21.12, 21.11*f*
 relativistic effects, **IV**:21.19–21.20, 21.21*f*
 rescattering effects, **IV**:21.18*f*, 21.18–21.20, 21.19*f*
 tunnel ionization, **IV**:21.12*f*, 21.12–21.14, 21.14*f*
 with clusters, **IV**:21.31–21.36
 Coulomb explosion, **IV**:21.33–21.34
 intense laser pulse interactions, **IV**:21.35–21.36, 21.36*f*
 ionization, **IV**:21.31–21.33, 21.32*f*
 nanoplasma description, **IV**:21.34–21.35, 21.35*f*
 history of, **IV**:21.3–21.4
 laser technology for, **IV**:21.4*f*, 21.4–21.5
 with molecules, **IV**:21.22–21.26
 Coulomb explosion, **IV**:21.24–21.25
 nuclear motion and alignment, **IV**:21.22–21.23, 21.23*t*, 21.24*f*
 triatomic and larger molecules, **IV**:21.26
 tunnel ionization and ionization distance, **IV**:21.25*f*, 21.25–21.26, 21.27*f*

- Strong field interactions (*Cont.*):
 nonlinear optics in gases, **IV**:21.27–21.31
 attosecond pulse generation, **IV**:21.31
 high order harmonic generation,
IV:21.27–21.30, 21.28*f*
 with single electrons, **IV**:21.5–21.10
 interactions with relativistic electron
 beams, **IV**:21.9–21.10
 nonlinear Thomson scattering,
IV:21.8–21.9, 21.9*f*
 ponderomotive force, **IV**:21.5–21.6
 relativistic effects, **IV**:21.6–21.8, 21.7*f*
 with underdense plasmas, **IV**:21.36–21.52
 applications of, **IV**:21.52–21.55,
 21.53*f*–21.55*f*
 direct laser acceleration and betatron
 resonance, **IV**:21.42–21.43
 high harmonic generation,
IV:21.50–21.52, 21.51*f*
 intense laser pulses, **IV**:21.38*f*, 21.38–21.39
 inverse Bremsstrahlung heating, **IV**:21.37,
 21.37*f*
 ionization-induced defocusing, **IV**:21.43*f*,
 21.43–21.44
 $\mathbf{j} \times \mathbf{B}$ heating and anomalous skin effect,
IV:21.49
 ponderomotive channel formation,
IV:21.42
 ponderomotive steepening and hole
 boring, **IV**:21.49–21.50, 21.50*f*
 relativistic effects and induced
 transparency, **IV**:21.52
 resonance absorption, **IV**:21.47*f*,
 21.47–21.48
 self-channeling and self-phase modulation,
IV:21.44–21.46, 21.45*f*
 structure of irradiated plasma, **IV**:21.46*f*,
 21.46–21.47
 vacuum heating, **IV**:21.47*f*, 21.48–21.49
 wakefield generation and electron
 acceleration, **IV**:21.39–21.42, 21.40*f*,
 21.42*f*
 Strong VW reflectometer, **II**:35.10*f*
 Strontium, **IV**:14.15, 14.16*f*
 Strontium barium niobate (SBN), **IV**:12.17
 Strontium fluoride (SrF₂), **IV**:2.40*t*, 2.44*t*, 2.48*t*,
 2.52*t*, 2.57*t*, 2.64*t*, 2.69*t*
 Strontium molybdate (SrMoO₄), **IV**:2.40*t*,
 2.45*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.65*t*, 2.72*t*
 Strontium titanate (SrTiO₃), **IV**:2.40*t*, 2.44*t*,
 2.48*t*, 2.52*t*, 2.57*t*, 2.65*t*, 2.73*t*
 Structured antireflection coatings, **IV**:7.25*f*,
 7.25–7.26, 7.26*f*
 Strut design (in stray light suppression),
II:7.20, 7.21, 7.21*f*
 Styrene acrylonitrile (SAN), **IV**:3.4*t*, 3.6, 3.6*t*,
 3.7*t*
 Subadditivity, **III**:11.51, 11.52*f*
 Subaperture size, for adaptive optics,
V:5.40–5.43, 5.42*f*
 Sub-Doppler absorption spectroscopy,
I:31.26–31.27
 Subjective amplitude of accommodation,
III:1.32
 Subjective refraction, **III**:12.6–12.7
 Subjective speckle, **I**:33.9
 Subjective tasks, **III**:2.15*n*
 Subjective tone reproduction, **II**:29.16
 Submillimeter (SubMM) radiation, **II**:24.3
 Subminiature lamps, **II**:15.53
 Sub-Nyquist interferometry, **II**:13.27
 Substrate(s):
 absorbing, **II**:17.7, 17.7*t*
 for HB-LEDs, **II**:18.2*f*, 18.2–18.3
 LED, **II**:17.20–17.21, 17.20*t*
 mirror, **II**:6.17*f*, 6.17–6.18, 6.18*f*
 transparent, **II**:17.7, 17.7*t*
 Subtractive adaptation, **III**:2.26
 Sum frequency generation (SFG), **IV**:16.2, 16.3*t*
 Sum rules:
 for crystals and glasses, **IV**:2.9
 for dispersion in solids, **IV**:8.14–8.15
 for semiconductors, **IV**:5.11
 Sum-frequency lasers, **V**:5.32
 Suncatchers, **II**:40.48, 40.50*f*
 Super resolution, **IV**:23.10, 23.12
 Superadditivity, **III**:11.51, 11.52*f*
 Super-Beer's law, **IV**:23.15
 Superconducting bolometers, **II**:28.5, 28.7*t*
 Superconducting quantum (SQUID) detector,
V:60.9
 Superconducting tunneling junctions (STJ)
 detectors, **V**:60.9, 60.9*t*
 Supercontinuum (SC) generation, **V**:11.23,
 11.24*f*
 Superlinear absorption of light, **IV**:5.57
 Supermirrors, **IV**:7.111; **V**:41.7, 41.8, 64.7
 Superposition, **II**:39.2, 39.32, 39.33*f*; **III**:10.25,
 10.26*f*

- Superposition-of-sources nonlinearity measurement, **II**:34.33
- Superresolution (visual acuity), **III**:4.1, 4.15
- Supersensitizers, **II**:30.14, 30.15, 30.15f
- Superstructure gratings (SSGs), **V**:13.34, 13.35f
- Superzone construction, for micro-Fresnel lenses, **I**:22.36
- Support wires, lightbulb, **II**:40.29f, 40.30
- Suppression (binocular vision), **III**:13.12–13.13
of blur, **III**:13.18–13.19
defined, **III**:23.3
with head mounted visual displays, **III**:13.33
- Suprathreshold models (color vision), **III**:11.41
- Surface acoustic wave (SAW) devices, **I**:21.16
- Surface axial coordinates, in grazing incidence optics, **V**:45.7
- Surface coatings, **IV**:6.4t
- Surface damage, laser-induced, **IV**:19.2–19.4
- Surface emitting lasers (SELs) (SLASERs), **II**:19.39f, 19.39–19.41, 19.40f, 19.42f, 19.43t, 25.15
- Surface figure metrology, **V**:46.3–46.6, 46.5f
- Surface finish metrology, **V**:46.2
- Surface finishing, of diamond-turned optics, **II**:10.9f–10.11f, 10.9–10.11
- Surface generation current, **II**:32.10, 32.11, 32.11f
- Surface measurement systems, **II**:40.53, 40.54
- Surface micromachined MEMS devices, **III**:15.11
- Surface mount device (SMD) package, **II**:18.5f
- Surface mount LEDs (SMDs), **II**:40.37
- Surface profilometers, **II**:9.6
- Surface scattering, **I**:8.1–8.17; **II**:7.23
finish models for, **I**:8.14–8.15
and finite illumination area, **I**:8.11–8.12
Fresnel-Kirchhoff approximation for, **I**:8.5–8.9
fractal surfaces, **I**:8.8–8.9
statistically stationary surfaces, **I**:8.6–8.8
notation for, **I**:8.2–8.4
and power spectra, **I**:8.12–8.13, 8.13f
Rayleigh-Rice approximation, **I**:8.9–8.12
second-order statistical functions for, **I**:8.12–8.13
statistics for, **I**:8.12–8.15
and surface finish specifications, **I**:8.15–8.17
- Surface-channel CCDs, **II**:32.14
- Surface-channel MOS capacitors, **II**:32.4f, 32.7
- Surface-emitting light-emitting diodes (S-LEDs), **II**:25.15; **V**:13.36–13.40, 13.38f
- Surfaces, **I**:1.32
aspherical, **I**:1.35
conical, **I**:1.34–1.35
illuminants reflected from, **III**:10.32–10.36
metamerism for, **III**:10.38
modeling of, **II**:40.17
spherical, **I**:1.34
- Surface-tension effects, from polymer molding process, **IV**:3.14
- Surround, visual acuity and, **III**:4.12, 4.13f
- Susceptibility:
nonlinear
anharmonic oscillator model of
second-order, **IV**:10.7–10.9, 10.8f
of crystals and glasses, **IV**:2.26
quantum theory of, **IV**:10.9–10.10
Raman, **IV**:15.5–15.6
- Suspended luminaires, **II**:40.13f
- Suspension systems, of resonant scanners, **I**:30.43, 30.43f, 30.44
- Suzaku observatory, **V**:33.3t, 33.3–33.4
- SUZUKU spacecraft, **V**:47.6, 47.7f
- SVGA (super video graphics array), **III**:23.3
- Svishchev confocal microscope, **III**:17.6, 17.6f
- Sweatt lenses, **I**:23.4
- Sweatt model, **I**:23.4, 23.6
- Swept Source OCT (SS-OCT), **III**:18.2 (*see also* Optical frequency domain imaging)
- Swept-carrier time-domain optical memory, **IV**:11.25
- Swift Burst Alert Telescope (BAT), **V**:33.1
- Switches and switching:
for fiber-based couplers, **V**:16.4
in integrated optics, **I**:21.34f, 21.34–21.35
interferometric optical switches, **I**:32.19
for networking, **V**:18.4, 18.5, 18.5f, 18.11, 18.11f, 18.12f
and SOAs, **V**:19.28, 19.28f
and third-order optical nonlinearities, **IV**:16.30–16.31
- Sylvania, **II**:15.49
- Symmetrical lenses, **I**:1.71, 17.26–17.27, 17.27f
- Symmetry (color matching), **III**:10.8
- Symmetry properties, of crystals and glasses, **IV**:2.5, 2.6t–2.8t
- Synchronization standards, for color CRTs, **III**:22.13–22.14
- Synchrotron radiation, **II**:34.26–34.27

- Synchronous digital hierarchy (SDH), **V**:9.17
- Synchronous optical networks (SONETs),
V:9.17, 21.6, 23.6
- Synchronous pumping, **II**:16.29
- Synchrotron beamlines, SHADOW code and,
V:35.1–35.2, 35.4
- Synchrotron radiation (SR) monochromators,
V:39.6
- Synchrotron radiation sources, **V**:55.1–55.20
- adaptive x-ray optics for, **V**:50.2–50.7, 50.4*f*,
50.6*f*, 50.8*f*
 - coherence of, **V**:55.17–55.20, 55.18*f*–55.19*f*
 - Compton sources vs., **V**:55.2–55.3, 55.3*t*
 - history of, **V**:55.1–55.2
 - insertion devices, **V**:55.9–55.16, 55.10*f*,
55.12*f*, 55.13*f*
 - as linear polarizers, **V**:43.2, 43.3
 - and refractive x-ray lenses, **V**:37.5, 37.11
 - theory of synchrotron radiation emission,
V:55.2*f*, 55.2–55.9, 55.3*f*, 55.5*f*–55.8*f*
- Synthetic aperture radar data, **I**:11.6–11.8,
11.7*f*–11.8*f*
- System specifications, for lenses, **II**:4.3
- Système International (SI), **II**:12.2, 36.2, 37.3
(*see also* SI units)
- System-response (SR) function, **I**:8.11
- TABO, **III**:12.5, 12.7
- Taillights, **II**:40.21, 40.22*f*, 40.64*f*, 40.67, 40.68*t*,
40.69*f*
- Tailored (T) reflectors, **II**:39.37–39.39, 39.38*f*
- Tailoring (of uniformity), **II**:39.2
- Takagi-Taupin calculations, for MLLs,
V:42.12–42.14
- Talbot autoimages, **II**:12.23, 12.24
- Talbot's law, **II**:34.33–34.34
- Talystep stylus profiler, **V**:46.2
- Tanabe-Sugano diagrams, **V**:2.9–2.11, 2.10*f*,
2.19
- Tandem limiters, **IV**:13.6, 13.6*f*
- Tandem-lens arrays, **II**:39.34, 39.35*f*–39.37*f*
- Tangent vectors, of space curves, **I**:1.18, 1.19
- Tangential fans and foci, **I**:1.35
- Tangential phase matching (TPM), **V**:6.12,
6.13, 6.17, 6.25–6.27, 6.26*f*
- Tantalum, **IV**:4.50*t*, 4.69*t*, 4.70*t*
- Tantalum crown glass, **IV**:2.42*t*
- Tapered fiber bundles (TFBs), **V**:25.3,
25.10–25.11, 25.11*f*, 25.16
- Tapered fiber coupler process, **V**:16.1–16.3,
16.2*f*
- Tapered fiber method, for adaptor fabrication,
V:25.17, 25.17*f*
- Tapered lightpipes, **II**:39.12–39.13, 39.13*f*,
39.31*f*, 39.31–39.32
- Target normal sheath acceleration, **IV**:21.54
- “Target” sensitivity method (color vision) [*see*
Test sensitivity method (color vision)]
- Task lighting, **II**:40.12, 40.14
- Taylor-Hobson Form TalySurf, **II**:9.6
- Tear lens, **III**:12.12
- Tears, aging-related changes in, **III**:14.4–14.5
- Technical specifications, **II**:4.2
- Teflon coatings, for infrared optical fibers,
V:12.4, 12.7, 12.9
- Teflon overcoat, **IV**:6.27
- Teflon Wet Lubricant, **IV**:6.27
- Tehis method, **II**:40.53, 40.54
- Telecentric distribution, **II**:39.18, 39.18*f*
- Telecentric lenses, **I**:18.12
- Telecentric principle, **III**:6.6–6.7
- Telecentric stop, **I**:17.9
- Telecentricity, **I**:1.83–1.84, 30.31,
30.31*f*, 30.32
- Telecommunication systems, data
communication vs., **V**:15.1–15.2
- Telephoto lenses, **I**:17.29, 27.2, 27.6,
27.7*f*–27.16*f*, 27.13
- Telescope(s):
- and adaptive optics, **V**:5.2–5.5, 5.3*f*, 5.4*f*,
5.43*f*–5.46*f*, 5.43–5.46
 - astronomical, **I**:18.10; **II**:1.7*f*
 - and black surfaces, **IV**:6.21
 - Burst Alert, **V**:33.1
 - Calar Alto, **V**:5.27
 - Canada-France-Hawaii, **V**:5.23
 - Cassegrain, **V**:44.4
 - Cassegrainian, **I**:18.21
 - ESO, **V**:5.35
 - far-infrared, **IV**:6.48
 - field of view in, **I**:18.15*f*, 18.15–18.16
 - Galilean, **I**:18.15; **II**:1.7*f*
 - Gemini North, **V**:5.20, 5.21*f*
 - with grazing incidence optics, **V**:44.6–44.12,
44.7*f*, 44.9*f*–44.11*f*, 45.1
 - ground-based, **IV**:6.12
 - hard vs. soft x-ray, **V**:50.2
 - Hobby-Eberly, **V**:5.2
 - Hubble, **II**:11.4, 13.24

- Telescope(s) (*Cont.*):
 hyperboloid-hyperboloid (HH) grazing incidence x-ray, **V**:44.10–44.12, 44.11*f*
 Keck, **II**:11.4; **V**:4.36, 5.2, 5.27
 Keplerian, **I**:18.10
 Large Binocular, **V**:5.5
 Mersenne, **I**:18.19
 Mt. Wilson, **V**:5.27
 Multiple Mirror, **V**:5.5
 objectives for, **I**:29.8, 29.9, 29.12, 29.29–29.30
 reflecting, **II**:11.4
 resolution of, **V**:4.2–4.3
 ROSAT, **V**:44.4, 44.6*f*
 Shane, **V**:5.27
 SOR, **V**:5.27
 Swift Burst Alert, **V**:33.1
 10-m, **V**:5.45*f*–5.46*f*, 5.45–5.46
 terrestrial, **I**:18.10–18.11, 18.11*f*
 3.5m, **V**:5.43*f*–5.45*f*, 5.43–5.45
 unit magnification Galilean, **II**:12.4, 12.4*f*
 wide-field lobster-eye, **V**:48.4
 Wolter, **V**:44.4, 44.5*f*, 44.6–44.10, 44.7*f*, 44.10*f*, 45.1–45.5, 45.2*f*
 Wolter-Schwarzschild, **V**:44.7, 44.11, 45.1
 Telescopic lenses, **I**:1.46 (*see also* Afocal lenses)
 Telescopic transformations, **I**:1.57
 Telescoping input optics, **II**:38.7
 Telescopy:
 integrated optics and cable, **I**:21.2, 21.32–21.34
 Scophony TV projection system, **I**:30.45
 Tellurites, for fiber lasers, **V**:25.27*t*, 25.28
 Tellurium dioxide (TeO₂), **V**:6.17, 6.21*t*, 6.25, 6.29*t*, 6.34*t*, 6.39, 6.42
 Temperature(s):
 color, **II**:37.4*t*, 37.6–37.7
 correlated color, **II**:37.7, 38.5
 and crystalline-quartz retardation plates, **I**:13.48
 of crystals and glasses, **IV**:2.32, 2.33
 Curie, **I**:35.25
 distribution, **II**:37.7
 and integrated optics, **I**:21.12
 in laser cooling, **IV**:20.5
 and laser diodes, **V**:13.11
 and liquid crystals, **V**:8.17, 8.21–8.22, 8.22*f*
 and long-period grating sensors, **V**:24.13, 24.13*t*
 of metals, **IV**:4.7
 Temperature(s) (*Cont.*):
 and mounted optics, **II**:6.21–6.24, 6.22*f*–6.24*f*
 radiance, **II**:37.4*t*, 37.6
 and refractive index of glasses, **IV**:2.24–2.26
 Temperature control, of PZT transducers, **II**:22.19
 Temperature dependence:
 of bandpass filters, **IV**:7.94
 of line broadening parameters, **IV**:15.19*t*
 of line shift parameters, **IV**:15.19*t*
 Temperature noise, **II**:24.12
 Temperature specifications, for lenses, **II**:4.10
 Temperature-dependence effects, **II**:34.36–34.37
 Templates (for curvature measurement), **II**:12.17
 Temporal beats, interference and, **I**:2.13
 Temporal coherence, **I**:2.41, 5.3; **III**:18.2
 Temporal coherence, of synchrotron radiation sources, **V**:55.17–55.18, 55.18*f*
 Temporal contrast detection, **III**:2.27*f*, 2.27–2.29, 2.29*f*
 Temporal despeckling, **III**:5.21
 Temporal instability, of metals, **IV**:4.9
 Temporal signals, analog processing of, **I**:11.8–11.12, 11.9*f*–11.11*f*
 Temporal vision, age-related changes in, **III**:14.19–14.21, 14.20*f*
 10-m telescope systems, **V**:5.45*f*–5.46*f*, 5.45–5.46
 Tensile-strained QW lasers, **II**:19.16, 19.16*f*, 19.17
 Tensor, electro-optic, **I**:21.10
 Tensor product, of Jones matrices, **I**:14.23
 Tensor properties:
 of crystals and glasses, **IV**:2.5, 2.6*t*
 of third-order optical nonlinearities, **IV**:16.2–16.3
 Tensors:
 compliance, **IV**:2.30
 dielectric, **IV**:2.17–2.18
d-tensor, **IV**:10.11
 inverse dielectric, **IV**:2.6*t*, 2.19
 second-order susceptibility, **IV**:10.10–10.11
 Terabit (unit), **V**:20.2
 Terahertz asymmetric optical demultiplexers (TOAD), **V**:20.22
 Ternary layers, in fiber optic devices, **V**:13.2
 Terrestrial telescopes, **I**:18.10–18.11, 18.11*f*
 Tessar lenses, **I**:17.26, 17.26*f*

- Test plates (for curvature measurement), **II**:12.17
- Test (“target”) sensitivity method (color vision), **III**:11.11, 11.12, 11.34–11.46
 defined, **III**:11.3
 and different directions of color space, **III**:11.39, 11.41–11.43, 11.44f, 11.45–11.46, 11.46f
 luminance, **III**:11.37–11.39, 11.38f, 11.40f
 and spectral lights, **III**:11.34, 11.35f, 11.36f, 11.37
- Testing, **II**:13.1–13.27
 aspherical wavefront measurement, **II**:13.23–13.27
 holographic compensators, **II**:13.25, 13.25f, 13.26f
 infrared interferometry, **II**:13.25
 Moiré tests, **II**:13.26–13.27
 refractive or reflective compensators, **II**:13.24, 13.24f, 13.25
 sub-Nyquist interferometry, **II**:13.27
 two-wavelength interferometry, **II**:13.25, 13.26
 wavefront stitching, **II**:13.27, 13.27f
 computer-generated holograms in (*see* Computer-generated holograms)
 of convex surfaces, **II**:14.5
 interferogram evaluation, **II**:13.14–13.18
 direct interferometry, **II**:13.17–13.18
 fixed interferograms, **II**:13.14–13.15
 Fourier analysis of interferograms, **II**:13.16–13.17, 13.17f
 global and local interpolation of interferograms, **II**:13.15–13.16
 interferometric, **II**:13.7–13.12
 common path interferometer, **II**:13.9, 13.11f
 Fizeau interferometer, **II**:13.8–13.9, 13.9f, 13.10f
 lateral shearing interferometers, **II**:13.9–13.12, 13.11f, 13.12f
 multiple-pass interferometers, **II**:13.13
 multiple-reflection interferometers, **II**:13.13
 radial, rotational, and reversal shearing interferometers, **II**:13.12, 13.13f
 sensitivity of interferometers, **II**:13.13–13.14, 13.14f
 Twyman-Green interferometer, **II**:13.7f, 13.7–13.8, 13.8f
 Zernike phase-contrast method applied to interferometers, **II**:13.13–13.14, 13.14f
- Testing (*Cont.*):
 Knoop, **IV**:2.31, 2.32f
 noninterferometric, **II**:13.1–13.7
 Foucault test, **II**:13.2f, 13.2–13.3, 13.3f
 Hartmann test, **II**:13.4–13.6, 13.5f
 Hartmann-Shack test, **II**:13.6f, 13.6–13.7
 Ronchi test, **II**:13.3f, 13.3–13.4, 13.4f
 phase-shifting interferometry, **II**:13.18f, 13.18–13.23, 13.19f, 13.20f
 heterodyne interferometer, **II**:13.22
 integrating bucket method, **II**:13.21, 13.21f
 phase errors, **II**:13.22
 phase stepping, **II**:13.20, 13.20f
 phase-lock method, **II**:13.23, 13.23f
 simultaneous measurement, **II**:13.22
 two steps plus one method, **II**:13.21, 13.22
 of polymers, **IV**:3.16
 in wafer processing, **II**:17.24
- Tests of visual acuity, **III**:4.7f, 4.7–4.8
 for infants, **III**:4.8
 traditional visual acuity chart, **III**:4.6
- Tetragonal crystals, **IV**:8.9t, 8.19t
 room-temperature elastic constants, **IV**:2.44t–**IV**:2.45t
 symmetries of, **IV**:2.7t
- Tetragonal perovskite, **IV**:2.73t
- Tetrahedral lattice site, **IV**:5.6
- Tewarson, **I**:12.30
- Texas Instruments, **II**:28.12
- Textured graphite surface, **IV**:6.8f
- Thallium bromide (TlBr), **IV**:2.40t, 2.44t, 2.48t, 2.53t, 2.58t, 2.65t, 2.68t
- Thallium chloride (TlCl), **IV**:2.65t, 2.68t
- Thef-number, **II**:38.8
- Theodolites, **II**:12.13
- Theoretical horopter, **III**:13.8, 13.9f
- The Theory of Coherent Atomic Excitation* (B. W. Shore), **IV**:14.3
- Thermal arrays, **II**:28.7–28.12
 about, **II**:28.7–28.8
 noise equivalent temperature difference in, **II**:28.8–28.9
 pyroelectric hybrid, **II**:28.11f, 28.11–28.12, 28.12f
 resistive bolometer, **II**:28.10f, 28.10–28.11
 theoretical limits of, **II**:28.9f, 28.9–28.10
 thermoelectric, **II**:28.12, 28.12f
- Thermal blooming, **IV**:16.22
- Thermal circuit theory, **II**:28.2

- Thermal coefficient of resistance (TCR),
II:33.2, 33.14
- Thermal compensation, **II**:8.1–8.15
 about, **II**:8.2
 and effect of thermal gradients, **II**:8.6–8.7
 and homogeneous thermal effects, **II**:8.2–8.5,
 8.3*t*, 8.4*t*, 8.5*f*
 intrinsic athermalization, **II**:8.7*f*, 8.7–8.8
 mechanical athermalization, **II**:8.8*f*–8.10*f*,
 8.8–8.12, 8.11*f*, 8.12*f*
 optical athermalization, **II**:8.12–8.15,
 8.13*t*–8.15*t*
 tolerable homogeneous temperature change,
II:8.5–8.6, 8.6*f*
- Thermal conductivity:
 of crystals and glasses, **IV**:2.6*t*, 2.35*f*,
 2.35–2.36
 of metals, **IV**:4.7, 4.10*t*, 4.53, 4.55, 4.55*t*,
 4.58*t*, 4.60*f*–4.64*f*
- Thermal control and correction, for adaptive
 optics, **V**:50.5, 50.6
- Thermal cycling, of metals, **IV**:4.10
- Thermal defocus, of compound lens,
II:8.4, 8.5*f*
- Thermal defocusing, **IV**:13.8
- Thermal detector(s), **II**:24.4*f*, 24.4–24.6,
 28.1–28.12, 38.9, 38.9*t*
 arrays of, **II**:28.7–28.12
 about, **II**:28.7–28.8
 noise equivalent temperature difference,
II:28.8–28.9
 pyroelectric hybrid arrays, **II**:28.11*f*,
 28.11–28.12, 28.12*f*
 resistive bolometer arrays, **II**:28.10*f*,
 28.10–28.11
 theoretical limits, **II**:28.9*f*, 28.9–28.10
 thermoelectric arrays, **II**:28.12, 28.12*f*
 bolometer, **II**:24.5*f*, 28.3–28.5, 28.4*f*
 Golay cell, **II**:28.6
 ideal, **II**:28.2–28.3, 28.3*f*
 performance/sensitivity of, **II**:24.17, 24.18*f*
 properties of, **II**:28.7, 28.7*t*
 pyroelectric, **II**:24.6, 24.6*f*, 28.7
 and thermal circuit theory, **II**:28.2
 thermistor, **II**:24.5
 thermocouple, **II**:28.4
 thermopile, **II**:24.5*f*, 28.4–28.5
- Thermal diffusivity, for metals, **IV**:4.10*t*
- Thermal effects, on third-order optical
 nonlinearities, **IV**:16.22
- Thermal expansion, **II**:33.14
 of crystals and glasses, **IV**:2.6*t*, 2.34*f*,
 2.34–2.35
 for metals, **IV**:4.10*t*
 of metals, **IV**:4.6
- Thermal fatigue, **II**:17.25
- Thermal focus shift, **II**:8.2–8.4, 8.3*t*, 8.4*t*
- Thermal gradients, effect of, **II**:8.6–8.7
- Thermal imaging, **V**:12.3*t*
- Thermal imaging cameras, **I**:25.25
- Thermal infrared detectors, **II**:33.7, 33.8*f*
- Thermal instability, of metals, **IV**:4.10
- Thermal neutrons, **V**:63.3
- Thermal noise, **II**:24.13, 27.4, 32.20
- Thermal (Johnson) noise, **V**:13.70
- Thermal optimization, of media, **I**:35.20, 35.20*f*
- Thermal properties:
 of crystals and glasses, **IV**:2.50*t*–2.53*t*, 2.55*t*
 of high-power lasers, **II**:19.26, 19.27*f*
 of metals, **IV**:4.6–4.7, 4.53, 4.55
 coefficient of linear thermal expansion,
IV:4.56*t*, 4.57*f*, 4.58*f*
 elastic properties, **IV**:4.69, 4.69*t*
 at room temperature, **IV**:4.55*t*
 specific heat, **IV**:4.65*t*, 4.66*f*–4.69*f*
 strength and fracture properties, **IV**:4.70,
 4.70*t*
 thermal conductivity, **IV**:4.58*t*, 4.59*f*–4.63*f*
- Thermal runaway, of lasers, **V**:13.11
- Thermal self-focusing, **IV**:13.7–13.8
- Thermal (Lambertian) sources of light,
I:5.12–5.13
- Thermal spectral radiance, **V**:3.18, 3.20, 3.20*f*
- Thermal stability, of plastic packaging
 materials, **II**:17.26
- Thermalization, of free electron and hole
 distributions, **IV**:18.20–18.21
- Thermally expanded cores (TECs), **V**:25.3,
 25.17, 25.17*f*
- Thermistor bolometers, **II**:24.24*f*, 24.24–24.25,
 24.25*f*, 28.7*t*
- Thermistors, **II**:24.5
- Thermocouple junctions, noise from,
II:27.6, 27.6*f*
- Thermocouples, **II**:24.5, 28.7*t*
 about, **II**:28.1
 manufacturers' specifications for, **II**:24.22*f*,
 24.22–24.23
 as thermal detectors, **II**:28.4
- Thermoelectric arrays, **II**:28.12, 28.12*f*

- Thermoelectric coolers (TECs),
V:19.4*f*, 19.5
- Thermomagnetic recording process,
I:35.17–35.20, 35.18*f*–35.20*f*
- Thermomodulation, **IV**:5.66*t*
- Thermo-optic coefficients, of crystals and
 glasses, **IV**:2.21, 2.24*f*, 2.24–2.26
- Thermo-optic effect, **IV**:16.22
- Thermopiles, **II**:24.5, 28.7*t*
 defined, **II**:24.13
 manufacturers' specifications for, **II**:24.23*f*,
 24.23–24.24
 as thermal detectors, **II**:28.4–28.5
- Thermoplastic resins, **IV**:3.2
- Thermoset resins, **IV**:3.2
- θ_1/θ_2 concentrators, **II**:39.18–39.20, 39.19*f*
- Thick lens systems, **I**:1.55
- Thick window chips, **II**:17.7, 17.7*t*
- Thickness errors, for multilayer reflectors,
IV:7.40
- Thin doublet, **II**:1.15–1.16
- Thin film oxides, **I**:21.13–21.14
- Thin film transistors (TFTs), **V**:61.3,
 61.4, 61.6
- Thin lens systems, **I**:1.55
- Thin lenses, **II**:1.5; **V**:40.3–40.4
- Thin teflon diffusers, **II**:38.15*f*
- Thin-disk lasers, **II**:16.18
- Thin-film coatings:
 and antireflection coatings, **IV**:7.27–7.28,
 7.28*f*
 interference and, **I**:2.24
 laser-induced damage to, **IV**:19.4
 manufacturing of, **IV**:7.10–7.12
 of metal, **IV**:7.104, 7.104*f*
 for multiple reflection filters, **IV**:7.111, 7.112,
 7.112*f*, 7.113*f*
 theory and design of, **IV**:7.5–7.10, 7.6*f*, 7.9*f*
- Thin-lens model:
 of Galilean afocal lenses, **I**:18.15, 18.15*f*
 of Keplerian afocal lenses, **I**:18.7–18.8, 18.8*f*
- Third-harmonic generation (THG), **III**:17.2,
 17.9–17.10
- Third-order aberrations, **I**:1.90–1.91, 29.38
- Third-order harmonic generation (THG),
IV:16.2, 16.3*t*
 in crystals, **IV**:16.14
 energy level diagrams for, **IV**:16.5*f*
 and semiconductors, **IV**:5.56
- Third-order optical nonlinearities,
IV:16.1–16.31
 cascaded $x^{(1)}:x^{(1)}$ processes, **IV**:16.20–16.22,
 16.21*f*
 cascaded $x^{(2)}:x^{(2)}$ processes, **IV**:16.22–16.24,
 16.23*f*, 16.24*f*
 and four-wave mixing, **IV**:16.27–16.28, 16.28*f*
 and interferometry, **IV**:16.28–16.29
 Kerr effect, **IV**:16.11–16.14, 16.13*f*, 16.14*f*
 Kramers-Kronig dispersion relations,
IV:16.9–16.11
 nonlinear absorption and nonlinear
 refraction, **IV**:16.7–16.9
 propagation effects, **IV**:16.24–16.26
 and quantum mechanics, **IV**:16.4–16.7, 16.5*f*
 and semiconductors, **IV**:5.55
 stimulated scattering, **IV**:16.14–16.19, 16.15*f*,
 16.17*f*
 terms for, **IV**:16.1–16.3, 16.3*t*
 third-harmonic generation, **IV**:16.14
 and time-resolved excite-probe techniques,
IV:16.26–16.27, 16.27*f*
 two-photon absorption, **IV**:16.19–16.20
 and Z-scan, **IV**:16.29–16.30, 16.30*f*
- 35-mm photographic films, **II**:30.21, 30.25
- Thompson reversed Nicol prisms, **I**:13.16*f*, 13.17
- Thomson backscattering, **V**:59.1
- Thomson CSF, **I**:21.35
- Thomson scattering, **IV**:21.8–21.9, 21.9*f*; **V**:26.7
- Thomson scattering cross-section, **V**:59.2
- Thoria (in incandescent lights), **II**:40.27
- Threaded retaining rings, **II**:6.3, 6.3*f*
- 3.5-m telescopes, **V**:5.43*f*–5.45*f*, 5.43–5.45
- Three-beam interferometers, **I**:32.7*f*, 32.7–32.8
- Three-chip color systems, **II**:32.32, 32.33*f*
- 3D bandgap materials, **IV**:9.2
- Three-dimensional color data, **III**:10.19–10.20,
 10.20*f*
- 3D concentrators, 2D vs., **II**:39.20*f*, 39.20–39.21,
 39.21*f*
- Three-dimensional diffraction patterns,
I:28.19–28.22, 28.20*f*, 28.21*f*
- 3D optical molasses, **IV**:20.16*f*, 20.16–20.17
- 3D photonic crystals, **IV**:9.4–9.8
 criteria for, **IV**:9.4–9.5
 examples of, **IV**:9.5, 9.5*f*
 microcavities in, **IV**:9.7*f*, 9.7–9.8, 9.8*f*
- 3D profilometry, **V**:46.6
- Three-lens prime focus corrector objective,
I:29.10

- Three-level atomic systems, **IV**:14.4*f*, 14.4–14.6, 14.6*f*
- 3M Black Velvet, **IV**:6.14
- 3M Black Velvet 101-C10, **IV**:6.12, 6.35
- 3M Company, **II**:29.25
- 3M Nextel Black Velvet, **IV**:6.35, 6.36*f*
- 3M paints, **IV**:6.35–6.37, 6.38*f*, 6.53*f*
- 3M Nextel Black Velvet, **IV**:6.35, 6.36*f*
- MH 2200, **IV**:6.37
- Nextel 2010, **IV**:6.35, 6.37
- Nextel Suede Coating Series 3101-C10, **IV**:6.37, 6.38*f*, 6.53*f*
- Three-material athermal solutions, **II**:8.14, 8.14*t*, 8.15*t*
- Three-mirror objectives, **I**:29.28–29.32, 29.34
- Three-phase CCDs, **II**:32.15, 32.16*f*
- Three-phase model (ellipsometry), **I**:16.5–16.8, 16.6*f*–16.8*f*
- Three-photon absorption (3PA), **IV**:19.9, 19.10, 19.10*f*
- Three-powered-mirror lenses, **I**:18.20, 18.20*f*, 18.21
- Three-ray rule, **I**:1.42
- Three-segment model of biological waveguides, **III**:8.9, 8.11*f*, 8.11–8.15
- assumptions and approximations for, **III**:8.9, 8.12–8.13
- electromagnetic validity of, **III**:8.13
- Three-stage zone models (color vision), **III**:11.6, 11.82–11.85, 11.85*f*
- De Valois De Valois model, **III**:11.82
- Müller model, **III**:11.6, 11.7*f*, 11.65*f*, 11.84, 11.85*f*
- Three-step rescattering model, **II**:21.3
- Three-wave mixing, **IV**:14.26
- Threshold carrier density, **II**:19.12, 19.12*f*, 19.13, 19.13*f*
- Threshold current, **II**:19.6, 19.6*f*
- Threshold experiments, **III**:11.17
- Threshold ionization:
- absorbance above, **IV**:21.14–21.17, 21.15*f*, 21.16*f*
- defined, **IV**:21.3
- Threshold modal gain, **II**:19.12, 19.12*f*, 19.13, 19.13*f*
- Threshold surface or contour (color vision), **III**:11.12–11.15, 11.13*f*
- defined, **III**:11.3
- and loss of information, **III**:11.20
- Threshold surface or contour (color vision) (*Cont.*):
- and noise, **III**:11.20, 11.23*f*, 11.23–11.26, 11.25*f*
- and second-site adaptation to steady fields, **III**:11.18*f*
- Threshold voltage, **II**:25.11
- Thresholding devices, **IV**:12.35–12.36, 12.36*f*
- Threshold-limit values (TLVs), **III**:7.9
- Thresholds, **III**:3.3*f*
- in adjustment experiments, **III**:3.4
- defined, 3; **III**:3.2
- detection (color vision)
- chromatic discrimination near, **III**:11.58–11.59
- two-stage model of, **III**:11.23*f*
- discrimination, **III**:4.6–4.7
- estimation of, **III**:3.9
- increment (color vision), **III**:11.16*f*
- modulation, **III**:1.22
- for motion detection/discrimination, **III**:2.37*f*
- in psychophysics, **III**:3.2–3.3
- suprathreshold models (color vision), **III**:11.41
- of visual resolution, **III**:4.6–4.7, 4.15
- Throughput, **I**:1.22, 13.7 (*see also* Étendue)
- Thulium-doped fibers, **V**:25.23*t*, 25.25, 25.32
- Tight binding approximations, for photonic crystal fibers, **V**:11.8
- Tightly toleranced assembly, **II**:6.7, 6.7*f*
- Tilt, **III**:13.2
- atmospheric, **V**:5.14, 5.14*f*–5.15*f*
- gradient (G-tilt), **V**:4.3
- and adaptive optics, **V**:5.14–5.16
- and angle of arrival, **V**:4.23, 4.25, 4.26
- wavefront (Z-tilt), **V**:4.3, 4.23–4.25, 5.14, 5.15
- Zernike, **V**:4.23–4.26
- Tilt errors, for grazing incidence optics, **V**:45.7
- Tilt-corrected phase variance, **V**:4.31*f*, 4.31–4.32
- Tilted diffraction geometry, **V**:42.4
- Tilted planes, for collineation, **I**:1.61, 1.62, 1.62*f*
- Tilted-plane processors, **I**:11.8, 11.8*f*
- Time:
- coherence, **I**:5.3
- visual acuity and, **III**:4.12, 4.13*f*
- Time averages, in coherence theory, **I**:6.2*n*, 6.4–6.5
- Time delay integration (TDI), **II**:33.4
- Time delay integration (TDI) linear sensors, **II**:32.23*f*, 32.24

- Time delay integration (TDI) scanning FPAs, **II**:33.17, 33.17*f*
- Time domain, **I**:7.17
- Time domain OCT, **III**:18.2–18.5, 18.4*f*
- Time evolution of the field, **II**:23.15*f*, 23.15–23.17
- Time lag, of cameras, **I**:25.8–25.9, 25.9*f*
- Time shifts, of solitons, **V**:22.15, 22.16
- Time-averaged color mixing, **II**:40.8
- Time-based measurement, **II**:12.2, 12.4, 12.5, 12.6*f*
- Time-dependent error, **II**:34.35
- Time-dependent (transient) spectroscopy, **IV**:11.2
- Time-division multiplexing (TDM), **V**:9.12, 20.3, 21.3
- Time-domain atom interferometers, **IV**:11.22–11.24, 11.24*f*
- Time-domain spectroscopy (TDS) systems, **V**:7.37, 7.37*f*–7.38*f*
- Time-gated imaging, **IV**:15.42, 15.43, 15.44*f*
- Time-integrated intensity, **IV**:11.18
- Time-integrating correlators, **I**:11.11*f*, 11.11–11.12
- Time-of-flight distance measurement, **II**:12.4, 12.5
- Time-of-flight (TOF) measurement of velocity distribution, **IV**:20.13, 20.13*f*
- Time-of-flight (TOF) spectroscopy, **V**:2.5
- Time-resolved excite-probe techniques, **IV**:16.26–16.27, 16.27*f*
- Time-sequential polarimeters, **I**:15.5
- Timing recovery, in OTDM networks, **V**:20.10*f*, 20.10–20.12, 20.11*f*
- Timing standards, for color CRTs, **III**:22.13–22.14
- Tin hypochlorophosphate ($\text{Sn}_2\text{P}_2\text{S}_6$), **IV**:12.17, 12.18*t*
- Tiodize V-E17, **IV**:6.49
- Titanium, **IV**:4.48*t*, 4.50*t*, 4.52*t*, 4.55*t*
- Titanium dioxide, **IV**:6.15
- Titanium oxide (TiO_2) UV detectors, **II**:24.47, 24.48*f*
- Titanium sapphire (Ti:sapphire) amplifiers, **IV**:21.5
- Titanium sapphire (Ti:sapphire) lasers, **IV**:18.3
- Titanium-doped sapphire ($\text{Ti:Al}_2\text{O}_3$) lasers, **II**:16.34, 16.34*f*
- Titanium-doped sapphire ($\text{Ti:Al}_2\text{O}_3$) ring lasers, **II**:20.16*f*, 20.16–20.17, 20.17*f*
- T-matrix method, **I**:7.15
- T-number, **I**:1.79
- Tokens, of FDDI connectors, **V**:23.3
- Tokyo Institute of Technology, **II**:19.41
- Tolerance budgeting, **II**:5.3
- Tolerance verification, **II**:5.3
- Tolerances, **II**:5.2–5.8
 assembly, **II**:5.8
 basis for, **II**:5.2–5.3
 boresight, **II**:5.8
 budgeting of, **II**:5.3
 distortion, **II**:5.8
 optical vs. mechanical, **II**:5.2
 verification of, **II**:5.3
 wavefront, **II**:5.3–5.7, 5.4*f*, 5.5*f*, 5.5*t*, 5.6*t*, 5.7*f*
- Tolerancing, **II**:5.8–5.11
 and aberration balancing, **II**:11.35, 11.36
 about, **II**:5.1–5.2
 and material properties, **II**:5.9
 measurement practices for, **II**:5.8–5.9
 and optimization, **II**:3.20–3.21
 problems in, **II**:5.11
 procedures for, **II**:5.9–5.10
 shop practices for, **II**:5.8
- Tomography, **IV**:23.13; **V**:31.1, 31.5, 31.5*f*–31.7*f*, 31.5–31.7
- Tomosynthesis, digital, **V**:31.7–31.8
- Tone reproduction, **II**:29.16–29.17, 29.17*f*
- Toner, in xerographic systems, **I**:34.1, 34.8*f*, 34.8–34.9
- Tonic level of accommodation, **III**:1.33–1.34
- Topothesy, of fractals, **I**:8.8
- Toric lenses:
 for astigmatism, **III**:12.13
 bitoric, **III**:20.17–20.20, 20.18*f*–20.20*f*, 20.18*t*
 intraocular, **III**:21.12, 21.13, 21.13*f*
 power of, **III**:20.16*f*, 20.16–20.20, 20.17*f*
 prism with, **III**:20.30
- Toroidal reflectors, for neutron beams, **V**:64.4
- Toroidal-grating monochromators (TGMs), **V**:38.3
- Torsion, **I**:1.19; **III**:13.2, 13.8, 13.22
- Torsional vergence, **III**:13.22, 13.27
- Total emittance, of metals, **IV**:4.6
- Total external reflection, **V**:64.1–64.2
- Total flux into a hemisphere, **II**:34.15
- Total hemispherical emittance, **II**:35.15, 35.15*f*
- Total integrated excitation, for crystals and glasses, **IV**:2.19

- Total integrated scatter (TIS):
of neutrons, **V**:64.2–64.3
and scatterometers, **V**:1.4, 1.6–1.7, 1.10–1.11
- Total internal reflection (TIR), **II**:39.12, 39.17, 40.41; **IV**:8.13; **V**:11.2, 11.3, 11.14, 11.15
of optical waveguides, **I**:21.3
of Rochon prisms, **I**:13.20
- Total internal reflection (TIR) Fresnel lenses, **II**:39.10
- Total luminous flux, **II**:37.4*t*, 37.6
- Total mass loss (TML), **IV**:6.17
- Total power law, **IV**:2.19–2.20, 2.20*f*
- Total radiant flux, **II**:37.4*t*, 37.6
- Total strain, of crystals and glasses, **IV**:2.36
- Total transmittance, **II**:35.3, 35.9*f*
- TRACE observatory, **V**:41.3
- Traceability:
of absolute measurements, **II**:34.21
errors in, **II**:34.28
- TracePro (optical software), **II**:7.27
- Track-error signal (TES), **I**:35.14–35.15, 35.15*f*, 35.17
- Tracking, for adaptive optics, **V**:5.15–5.18, 5.17*f*, 5.21–5.23, 5.22*f*
- Tracks, on optical disks, **I**:35.2–35.5, 35.3*f*–35.5*f*
- Transconductance amplifiers, **II**:27.11*f*, 27.11–27.12
- Transducer resonance, **II**:22.8, 22.11–22.12
- Transducers, **II**:22.17–22.20
- Transfer, in xerographic systems, **I**:34.10
- Transfer function (human vision), **III**:2.3
defined, **III**:6.1
of photoreceptors, **III**:8.23–8.24, 8.24*f*
temporal and spatial components of, **III**:2.11, 2.12*f*
- Transfer functions [*see specific functions, e.g.:*
Modulation transfer function (MTF)]
- Transfer matrices, **I**:1.66
- Transfer matrix solution (to Maxwell's equations), **IV**:9.3
- Transformers, in voltage amplifiers, **II**:27.11
- Transient four-wave mixing (TFW), **IV**:18.17*f*, 18.17–18.18
- Transient photon counting, **II**:27.14
- Transient Raman scattering, **IV**:15.22–15.32
broadband effects, **IV**:15.28–15.32, 15.29*f*
phase pulling, **IV**:15.26–15.27, 15.27*f*
pulsed, **IV**:15.22–15.25, 15.24*f*–15.26*f*, 15.24*t*
solitons, **IV**:15.27–15.28, 15.29*f*
spectral properties, **IV**:15.32
- Transient response:
of laser diodes, **V**:13.13–13.18, 13.14*f*, 13.16*f*, 13.17*f*
of light-emitting diodes (LEDs), **V**:13.41–13.42
- Transillumination, in microscopes, **I**:28.5–28.7, 28.6*f*
- Transition-edge sensor (TES) microcalorimeter detectors, **V**:29.9*f*, 29.9–29.11, 29.11*f*
- Transition-edge sensors (TES), **V**:60.9
- Transition-metal ions, spectra of, **V**:2.8*f*, 2.8–2.9
- Transitions, **IV**:16.4 (*see also specific transitions, e.g.:* One-electron transitions)
- Transitivity (color matching), **III**:10.8
- Translation, by scanners, **I**:30.5
- Transmission, **II**:4.7
actual/idealized, **II**:35.2*f*
of amplitude modulators, **V**:7.22
amplitude-shift-keyed, **I**:21.30
analog, **I**:21.32–21.34, 21.33*f*, 21.34*f*; **V**:9.15–9.17, 9.16*f*
atmospheric optical, **V**:3.22*f*–3.27*f*, 3.22–3.26
Bormann, **V**:43.3, 43.3*f*
broadband, **V**:3.23–3.24, 3.25*f*–3.26*f*
of coatings on substrate, **IV**:7.3
for cutoff filters, **IV**:7.54–7.55, 7.55*f*
defined, **II**:35.3
differential-phase-shift-keyed, **I**:21.30, 21.32
digital, **I**:21.31–21.32
electrical time domain multiplexed, **V**:20.3, 20.25
formats for data, **V**:20.9*f*, 20.9–20.10
Gaussian, **V**:37.5
by Glan-Thompson-type prisms, **I**:13.9–13.10, 13.10*f*, 13.11*f*
Laue, **V**:26.8, 26.9*f*, 63.24, 63.26
by LEDs, **I**:21.32
line-by-line, **V**:3.23, 3.24*f*
measurement of, **IV**:5.64
in multilayer systems, **I**:16.8–16.9
nonnormal-incidence, **I**:12.18–12.24
of NPM AOTFs, **V**:6.41
with optical fibers, **V**:9.7–9.8, 9.15–9.17
in OTDM networks, **V**:20.7–20.8, 20.12–20.17, 20.14*f*–20.16*f*
in passband region, **IV**:7.53, 7.54
in polycapillary x-ray optics, **V**:53.5–53.8, 53.6*f*, 53.7*f*
silicon photonics, **I**:21.14–21.16, 21.15*f*

- Transmission (*Cont.*):
 by silicon photonics, **I**:21.38, 21.39
 for solitons, **V**:22.5–22.7
 WDM dispersion managed soliton,
V:22.15–22.17
 of x-ray tube sources, **V**:54.11–54.12
- Transmission coefficient, for optical constants,
IV:5.9–5.10
- Transmission density, of photographic films,
II:29.6–29.7, 29.7*f*
- Transmission ellipsometry, **I**:16.10
- Transmission filters, **IV**:7.3–7.5, 7.83, 7.88*f*
- Transmission grating, **IV**:12.7, 12.8*f*
- Transmission loss, in fiber optic links, **V**:15.7
- Transmissive planar objects, **I**:6.3
- Transmissive sensors, **II**:17.34
- Transmissive thin-film transistor (TFT) LCDs,
V:8.29–8.31, 8.30*f*, 8.31*f*, 8.34*f*, 8.32–8.35
- Transmittance, **II**:35.3
 of Brewster angle transmission polarizers,
I:12.22
 in human eye, **III**:1.9*f*, 1.9–1.11
 at interface of solid, **IV**:8.12
 measurement of, **II**:35.8–35.10, 35.9*f*
 of metals, **IV**:4.6
 of Mueller matrices, **I**:14.16–14.17
 of optical coatings, **IV**:7.12–7.13
 of pile-of-plates polarizers, **I**:12.15–12.17
 principal, **I**:12.14–12.16
 and reflectance/absorptance, **II**:35.7, 35.8,
 35.8*t*
 spectral, **II**:38.2, 38.17, 38.17*f*; **IV**:7.3–7.4
 of spherical lenses, **I**:6.3
 of water, **IV**:1.5*t*
- Transmitted state, of polarizers, **I**:15.19
- Transmitter speed, for fiber optics,
II:17.33
- Transparency, **II**:39.23, 40.5
 and absorption, **IV**:2.17
 aging-related changes in, **III**:14.8 (*see also*
 Cataract)
 EIT [*see* Electromagnetically induced
 transparency (EIT)]
 induced, **IV**:21.52
- Transparency point, **II**:19.5
- Transparent fiber systems, **V**:14.1
- Transparent HMDs, **III**:25.4–25.5
- Transparent prisms, **IV**:5.59
- Transparent substrate (TS) chips,
II:17.7, 17.7*t*
- Transportation lighting, **II**:40.63–40.71
 roadway lighting, **II**:40.67, 40.69–40.71,
 40.70*t*, 40.71*t*
 vehicular lighting, **II**:40.63–40.67, 40.64*f*,
 40.65*t*, 40.66*f*, 40.66*t*, 40.68*t*, 40.69*f*
- Transreflective LCDs (TR-LCDs), **V**:8.32–8.35,
 8.33*f*, 8.34*f*
- Transverse acoustic (TA) phonons, **IV**:5.24,
 5.25*f*
- Transverse (lateral) chromatic aberration
 (TCA), **III**:1.19, 1.20, 15.22–15.23
 measurement of, **III**:8.8
 with visual instruments, **III**:1.28
- Transverse coherence, **V**:55.18–55.20, 55.19*f*
- Transverse effective wavelength, for PCFs,
V:11.10
- Transverse electric (TE) modes, of optical
 waveguides, **I**:21.6, 21.7
- Transverse electric (TE) polarization, **V**:19.7
- Transverse electric (TE) waveguide mode, of
 laser diodes, **V**:13.13
- Transverse electromagnetic mode (TEM),
II:16.21–16.23, 16.22*f*
- Transverse electro-optic modulators, **V**:7.16,
 7.17, 7.17*f*
- Transverse holographic technique, for fiber
 Bragg gratings, **V**:17.5, 17.5*f*
- Transverse junction stripe (TJS) lasers, **II**:19.8,
 19.9*f*, 19.23–19.24, 19.36*f*
- Transverse Kerr effect, **IV**:18.11, 18.12*f*,
 18.14–18.15
- Transverse laser modes, **II**:16.21*f*–16.23*f*,
 16.21–16.23
- Transverse magnetic (TM) modes:
 of laser diodes, **V**:13.13
 of optical waveguides, **I**:21.6, 21.7
- Transverse magnetic (TM) polarization, **V**:19.7
- Transverse magnification, **I**:1.28, 1.50–1.51
- Transverse optical (TO) frequencies, for
 crystals and glasses, **IV**:2.11, 2.12
- Transverse optical (TO) phonons, **IV**:5.24,
 5.25*f*, 5.80, 5.80*f*, 8.16–8.18
- Transverse primary chromatic aberration
 (TPAC), **I**:17.22
- Transverse ray aberration (TRA) equations,
V:45.3–45.8
- Transverse ray aberrations, **I**:1.87
- Transverse ray plots, **II**:2.2–2.4, 3.13
- Transverse relaxation time, **IV**:11.5
- Transverse spatial coherence, **V**:55.16

- Transverse spatial modulation (TSM),
V:6.11–6.12, 6.23, 6.30, 6.31
- Transverse translational scan, **I**:30.28
- Trap loss collisions, **IV**:20.29
- Trap-assisted thermal generation current,
V:13.69
- Trapping atoms, **IV**:20.21–20.39
 applications of, **IV**:20.26–20.39
 and atomic beam brightening, **IV**:20.27*f*,
 20.27–20.28
 and atomic clocks, **IV**:20.28
 and Bose–Einstein condensation,
IV:20.35–20.37, 20.36*f*
 and dark states, **IV**:20.37–20.39, 20.38*f*
 magnetic traps, **IV**:20.21–20.23, 20.22*f*
 magneto-optical traps, **IV**:20.24*f*,
 20.24–20.25, 20.26*f*
 and optical lattices, **IV**:20.31–20.34,
 20.32*f*–20.34*f*
 optical traps, **IV**:20.23*f*, 20.23–20.24
 and ultra-cold collisions, **IV**:20.28–20.31,
 20.30*f*, 20.31*f*
- Traveling microscopes, **II**:12.20, 12.21, 12.21*f*
- Traveling wave electro-optic modulators,
V:7.28–7.30, 7.29*f*
- Traveling wave modulators, **I**:21.26
- Traveling wave photodetectors, **II**:26.4*f*, 26.5,
 26.14*f*
- Traveling wave pumping, by EUV lasers, **V**:58.3
- Traveling-wave amplification, **V**:19.3, 19.4*f*
- Treaty of the Meter of 1875, **II**:34.20, 36.2
- Tree topologies, of WDM networks, **V**:21.6,
 21.7
- Tremor, **III**:1.44
- Trials, experimental, **III**:3.2
- Triatomic molecules, in strong fields, **IV**:21.26
- Trichromacy (trichromatic vision), **III**:10.4–10.6
 color appearance vs., **III**:11.5
 in color matching, **III**:10.7–10.9
 in color-deficient observers, **III**:10.16
- Triclinic crystals, **IV**:2.7*t*, 2.18, 8.9*t*, 8.10
- Trifocal lenses, **III**:12.10, 14.27
- Trigger detection, with ODMR, **I**:31.21
- Trigonal crystals, **IV**:8.9*t*, 8.19*t*
- Trigonal selenium, **IV**:2.70*t*
- Tri-Level highlight color process, **I**:34.12, 34.13*f*
- Trim retarders, for LC panels, **I**:15.33
- Trimmed Nicol-type polarizers, **I**:13.16*f*,
 13.16–13.18
- Trinitron CRTs, **III**:22.5
- Triphosphors, **II**:40.31, 40.32*f*
- Triplet lenses, **I**:17.26, 17.26*f*; **II**:6.21, 6.22*f*
- Triply resonant oscillators (TROs), **IV**:17.2–17.4,
 17.3*f*, 17.4*f*, 17.20–17.21, 17.21*f*
- Tri-positive rare earth ions, **I**:10.16–10.18,
 10.16*t*, 10.17*f*
- Trischiefspiegler objective, **I**:29.27
- Tristimulus values, **II**:38.3–38.4; **III**:10.4
 for arbitrary lights, **III**:10.9
 in colorimetric measurements, **III**:10.23
 and consistency across observers, **III**:10.9
 defined, **III**:10.1, 10.4*t*
 in maximum stimulus method, **III**:10.7
 negative, **III**:10.7
 and spectral power distribution, **III**:10.11
 spectral power distributions from,
III:10.36–10.37
 tailored to individuals, **III**:10.15
 uniqueness of, **III**:10.9
 in vector representations, **III**:10.27–10.29
- Tritanopes, **III**:10.16
- Trivalent rare-earth ions, **V**:2.11
- Trivex lenses, **III**:12.9
- Troffers, fluorescent luminaire, **II**:40.47
- Troland (unit), **II**:34.41–34.42, 37.7, 37.8;
III:5.3–5.4, 6.3
 defined, **III**:2.7*n*
 “effective,” **III**:8.7
 “equivalent,” **III**:9.2
 limitations of, **III**:9.3
 as unit of retinal illuminance, **III**:9.2
- Trough reflectors, **II**:40.46*f*, 40.47
- TRU-Color Diffuse Black, **IV**:6.49
- Trumpet (term), **II**:39.15, 39.16*f*, 39.17
- TRW, **II**:19.27
- T-trace formula, **I**:1.44
- Tube length (objective lenses), **I**:28.13, 28.13*t*
- Tube voltage, of x-ray tube sources, **V**:54.9
- Tubular PZT transducers, **II**:22.17–22.18
- Tunable dispersion compensation, in WDM
 networks, **V**:21.23–21.26, 21.24*f*–21.26*f*
- Tunable double resonance (electromagnetically
 induced transparency), **IV**:22.6–22.9,
 22.7*f*, 22.8*f*
- Tunable filters, acousto-optic, **V**:6.23*t*,
 6.35–6.45
 collinear beam, **V**:6.43, 6.43*f*, 6.45*t*
 long-infrared, **V**:6.42
 and longitudinal spatial modulation, **V**:6.12
 mid-infrared, **V**:6.42

- Tunable filters, acousto-optic (*Cont.*):
 noncritical phase-matching, **V**:6.37, 6.39–6.42, 6.38*f*, 6.43*t*
 principle of operation, **V**:6.36–6.39, 6.37*f*, 6.38*f*
 ultraviolet, **V**:6.42
- Tunable lasers, for fiber optic systems,
V:13.32–13.36, 13.33*f*–13.36*f*
- Tunable phase-dispersion filters, **IV**:7.89, 7.89*f*
- Tungsten:
 absorptance of, **IV**:4.42*f*, 4.48*t*, 4.50*t*, 4.51*t*
 elastic properties of, **IV**:4.69*t*
 extinction coefficient for, **IV**:4.18*t*–4.19*t*, 4.26*f*
 in HID lamps, **II**:40.35
 in incandescent lights, **II**:40.25, 40.27, 40.29
 reflectance of, **IV**:4.38*t*–4.39*t*, 4.42*f*
 refraction index for, **IV**:4.18*t*–4.19*t*, 4.26*f*
 resistivity of, **IV**:4.54*t*
 strength and fracture properties of, **IV**:4.70*t*
- Tungsten hexafluoride, **IV**:6.56
- Tungsten lamps, **II**:15.13, 40.26*t*, 40.28*f*
- Tungsten silicide/silicon (WSi_2/Si) bilayers, in
 MLLs, **V**:42.6, 42.7
- Tungsten-arc lamps, **II**:15.47–15.48, 15.48*f*, 15.49*f*
- Tungsten-filament lamps, **II**:15.11, 15.12, 15.13*f*, 15.19, 15.20, 15.20*f*–15.22*f*, 34.31
- Tungsten-halogen lamps, **II**:15.11, 15.12, 15.13*f*, 40.25*t*, 40.26*t*, 40.30
- Tuning relation, of NPM AOTFs, **V**:6.39–6.40
- Tunnel diagram (*see* Williamson construction)
- Tunnel ionization:
 atomic, **IV**:21.12*f*, 21.12–21.14, 21.14*f*
 molecular, **IV**:21.25*f*, 21.25–21.26, 21.27*f*
 relativistic, **IV**:21.20
- Tunnel lighting, **II**:40.71
- Tunneling, collective, **IV**:21.18
- Tunneling current, **II**:25.8
- Turbulence:
 and adaptive optics, **V**:5.5–5.21
 anisoplanatism, **V**:5.19
 atmospheric tilt and Strehl ratio, **V**:5.14, 5.14*f*–5.15*f*
 Fried's coherence diameter and spatial scale, **V**:5.9–5.13, 5.11*t*, 5.12*f*, 5.13*f*
 higher-order phase fluctuations, **V**:5.18–5.19
 imaging, **V**:4.35–4.36, 5.19*f*, 5.19–5.21, 5.20*f*
 Kolmogorov model, **V**:5.5–5.6
 tracking requirements, **V**:5.15–5.18, 5.17*f*
 variation of n and C_n^2 parameters, **V**:5.6–5.8, 5.7*f*, 5.8*f*
- Turbulence (*Cont.*):
 in atmospheric optics, **V**:3.26, 3.28–3.36
 beam spreading, **V**:3.32–3.33, 3.33*f*
 beam wander, **V**:3.31–3.32
 imaging and heterodyne detection, **V**:3.34
 parameters for, **V**:3.28–3.31, 3.29*f*, 3.30*f*
 scintillation, **V**:3.34–3.36, 3.35*f*
- Hufnagel model of, **V**:3.29
- Hufnagel-Valley model of, **V**:3.30, 5.7, 5.8
 imaging through atmospheric, **V**:4.1–4.37
 aberration variance and approximate Strehl ratio for, **V**:4.27–4.28, 4.28*f*
 adaptive optics, **V**:4.35–4.36
 angle of arrival fluctuations, **V**:4.23–4.26, 4.25*f*, 4.26*f*
 expansion coefficients, **V**:4.20–4.22, 4.21*t*, 4.22*f*–4.23*f*
- Kolmogorov turbulence and atmospheric coherence length, **V**:4.7–4.10, 4.8*f*, 4.9*f*
- long-exposure images, **V**:4.3–4.7
 modal correction of turbulence, **V**:4.28–4.30, 4.29*t*, 4.30*f*
 modal expansion of aberration function, **V**:4.17*f*–4.18*f*, 4.17–4.20, 4.19*t*, 4.20*t*
 and resolution of telescopes, **V**:4.2–4.3
 short-exposure image, **V**:4.31*f*–4.34*f*, 4.31–4.35, 4.35*t*
 systems with annular pupils, **V**:4.10–4.16, 4.11*f*–4.15*f*, 4.15*t*
- inner scale of, **V**:4.7
- Kolmogorov, **V**:4.3, 4.7–4.10, 4.8*f*, 4.9*f*, 4.27, 4.30, 4.36
- Kolmogorov model of, **V**:5.5–5.6, 5.11
 optical strength of (C_n^2), **V**:5.6–5.8, 5.7*f*, 5.8*f*
 outer scale of, **V**:4.7
 phase structure function of, **V**:4.5
- Turning, single-point, **IV**:3.12
- Turn-on delay, of laser diodes, **V**:13.13–13.14, 13.14*f*
- Tutton's test, **I**:13.45
- Tweezers, optical, **IV**:20.23
- Twilight myopia, **III**:1.34
- Twin beams of light, correlated, **IV**:17.28, 17.29*f*, 17.30*f*
- Twin-channel lasers (TCLs), **II**:19.27
- Twin-channel substrate mesa (TCSM) lasers, **II**:19.20*t*, 19.21*f*, 19.23
- Twin-ridge structure (TRS) lasers, **II**:19.19, 19.20*t*, 19.21*f*, 19.22–19.23

- Twist, of liquid crystals, **V**:8.22, 8.23*f*
- Twist grain boundary (TGB) phases, of liquid crystals, **V**:8.11
- Twisted nematic (TN) cells, **V**:8.16, 8.25–8.26, 8.26*f*
- Two-alternative forced choice (2afc), **III**:3.8
- Two-beam coupling:
 optical limiting by, **IV**:13.8–13.9
 photorefractive gain in, **IV**:12.29–12.32, 12.31*f*
 photorefractive loss in, **IV**:12.31–12.32, 12.33*f*–12.35*f*
 and wave interactions, **IV**:12.4*f*, 12.4–12.6
- Two-channel Maxwellian viewing, **III**:5.21, 5.23*f*, 5.23–5.24
- Two-color gating, **II**:21.7
- Two-color threshold experiments, **III**:11.34, 11.37
- Two-component magnetic brush development, in xerographic systems, **I**:34.5*f*–34.7*f*, 34.5–34.7
- Two-component systems, first-order layout for, **II**:1.5–1.7
- 2D (term), **II**:39.4
- 2D concentrators, 3D vs., **II**:39.20*f*, 39.20–39.21, 39.21*f*
- 2D high-power laser arrays, **II**:19.29–19.30, 19.29*t*, 19.30*f*
- Two-dimensional images, analog processing for, **I**:11.12–11.17, 11.13*f*
- 2D photonic crystals, microcavities of, **IV**:9.8–9.12, 9.9*f*
 in-plane coupling, **IV**:9.10–9.11
 out-of-plane coupling, **IV**:9.11–9.12
 waveguides in, **IV**:9.13*f*, 9.13–9.14
- 2D profilometry, **V**:46.6
- Two-dimensional scanners, **I**:30.18, 30.19*f*
- Two-interference pattern distance-measuring interferometer, **II**:12.7, 12.7*f*
- Two-lens systems, **I**:17.20–17.22, 17.21*f*–17.22*f*
- Two-level atoms:
 coherence in, **IV**:14.4*f*, 14.4–14.5
 force on, **IV**:20.6–20.7
 at rest, **IV**:20.7–20.8
- Two-level coupling, **IV**:14.30, 14.30*f*
- Two-material periodic multilayers theory, **IV**:7.32–7.38
 $[(0.5A)B(0.5A)]^N$ multilayers, **IV**:7.35, 7.36*f*
 angular sensitivity, **IV**:7.37
 multilayer reflectors of absorbing materials, **IV**:7.37–7.38, 7.38*f*
- Two-material periodic multilayers theory (Cont.):
 nonabsorbing $[AB]^N$ and $[AB]NA$ multilayers, **IV**:7.32–7.34, 7.33*f*–7.35*f*
 width of high-reflectance zone, **IV**:7.36–7.37, 7.37*f*
 $[xH.(1-x)L]^N.xH$ multilayers, **IV**:7.37
- Two-mirror, three reflection objective, **I**:29.26
- Two-mirror imaging system, **II**:39.17
- Two-phase CCDs, **II**:32.15–32.16, 32.16*f*
- Two-phase model (ellipsometry), **I**:16.5, 16.6*f*
- Two-photon absorption (2PA):
 energy level diagrams for, **IV**:16.5*f*
 and laser-induced damage, **IV**:19.9, 19.10, 19.10*f*
 and optical limiting, **IV**:13.4, 13.5, 13.6*f*
 of semiconductors, **IV**:5.56
 symbols, **IV**:16.8
 in third-order optical nonlinearities, **IV**:16.19–16.20
- Two-photon transitions, **IV**:11.22–11.23, 11.24*f*
- Two-powered-mirror lenses, **I**:18.20, 18.20*f*
- Two-ray paraxial invariant, **I**:1.41
- Two-stage baffle, **II**:7.10
- Two-step rescattering model, **II**:21.3
- Two-steps-plus-one phase shifting, **II**:13.21, 13.22
- Two-wave interaction (acousto-optic), **V**:6.8
- Two-wavelength interferometry, **I**:32.9;
II:13.25, 13.26
- Twyman-Green interferograms, **II**:13.10*f*, 13.18*f*
- Twyman-Green interferometers, **I**:2.28, 32.2, 32.9, 33.5; **II**:13.7*f*, 13.7–13.8, 13.8*f*, 13.18
- Tyler frequency, **V**:5.17
- Type A errors (in absolute measurement), **II**:34.21–34.23
- Type B errors and error sources (in absolute measurement), **II**:34.32–34.37
 defined, **II**:34.21–34.23
 nonideal aperture, **II**:34.35*f*, 34.35–34.36
 nonlinearity of detector, **II**:34.34–34.35
 nonuniformity, **II**:34.35
 offset subtraction, **II**:34.32–34.33
 polarization effects, **II**:34.33
 scattered radiation effect, **II**:34.33
 size-of-source effect, **II**:34.33
 spectral errors, **II**:34.36
 temperature-dependence effects, **II**:34.36–34.37
 time-dependent error, **II**:34.35

- UHURU satellite, **V**:47.1
- ULE glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*
- Ultimate strength, of metals, **IV**:4.8
- Ultra fast nonlinear interferometers (UNIs), **V**:19.32
- Ultra-cold collisions, **IV**:20.26, 20.28–20.31, 20.30*f*, 20.31*f*
- Ultrafast depletion, of semiconductor band states, **IV**:18.21
- Ultrafast lasers, **IV**:11.26, 18.3–18.5, 18.4*f*, 18.5*f*
- Ultrafast optics, coatings for, **IV**:7.47–7.48, 7.48*f*
- Ultrahigh-speed OTDM, **V**:20.23–20.24, 20.24*f*
- Ultrashort cavity microlasers, **II**:19.39
- Ultrashort optics, **II**:20.1–20.28
- about, **II**:20.1–20.2
- cavities with two circulating pulses, **II**:20.15–20.22
- linear lasers, **II**:20.18–20.19, 20.19*f*
- optical parametric oscillators, **II**:20.20*f*, 20.20–20.22, 20.21*f*
- ring dye lasers, **II**:20.15–20.16
- ring lasers, **II**:20.17*f*, 20.17–20.18
- Ti:sapphire ring lasers, **II**:20.16*f*, 20.16–20.17, 20.17*f*
- coupling of circulating pulses, **II**:20.12*f*, 20.12–20.15, 20.15*f*
- optical pulses and pulse trains, **II**:20.2–20.9
- single optical pulse, **II**:20.2–20.3, 20.3*f*
- soliton solution and steady-state pulse train, **II**:20.5–20.9
- train of pulses, **II**:20.3–20.5, 20.4*f*, 20.5*f*
- and quantum mechanical two-level system, **II**:20.22–20.28
- coherent interaction, **II**:20.22–20.23
- experimental demonstration, **II**:20.24–20.27, 20.25*f*–20.27*f*
- impact of analogy, **II**:20.27–20.28
- laser as two-level system, **II**:20.23–20.24, 20.25*t*
- Rabi cycling, **II**:20.26*f*, 20.26–20.27, 20.27*f*
- steady-state pulse, **II**:20.9–20.12, 20.11*f*
- Ultrashort pulse generation, **IV**:18.1–18.23
- Kerr effect, **IV**:18.11–18.15
- longitudinal, **IV**:18.11–18.15, 18.12*f*
- transverse, **IV**:18.14–18.15
- saturable absorbers, **IV**:18.5–18.11
- fast, **IV**:18.9–18.10
- self-amplitude modulation, **IV**:18.5–18.7, 18.6*f*, 18.7*f*
- Ultrashort pulse generation, saturable absorbers (*Cont.*):
- semiconductor saturable absorber mirrors, **IV**:18.3, 18.10–18.11
- slow, **IV**:18.7–18.9, 18.8*f*
- semiconductor ultrafast nonlinearities, **IV**:18.15–18.23
- carrier trapping, **IV**:18.21–18.23, 18.22*f*
- in coherent regime, **IV**:18.19–18.20
- continuum excitations, **IV**:18.20
- excitonic excitations, **IV**:18.19–18.20
- experimental techniques, **IV**:18.17*f*, 18.17–18.19
- properties, **IV**:18.16*f*, 18.16–18.17
- in thermalization regime, **IV**:18.20–18.21
- and ultrafast lasers, **IV**:18.3–18.5, 18.4*f*, 18.5*f*
- Ultrashort pulses (USPs), for fiber lasers, **V**:25.32–25.33
- Ultralow light pulses, **IV**:14.22–14.23, 14.23*f*
- Ultrasonic-assisted machining, **II**:10.5
- Ultraviolet (UV) AOTFs, **V**:6.42
- Ultraviolet (UV) crystals, **IV**:10.22*t*
- Ultraviolet (UV) detectors:
- silicon carbide, **II**:24.47, 24.47*f*
- TiO₂, **II**:24.47, 24.48*f*
- Ultraviolet (UV) enhanced photodiodes, **II**:24.55*f*, 24.61*f*, 24.61–24.62, 24.62*f*
- Ultraviolet (UV) filters, **II**:40.12
- Ultraviolet (UV) light:
- and black surfaces, **IV**:6.21, 6.22*f*–6.25*f*
- metal-dielectric reflectors for, **IV**:7.108–7.109, 7.109*f*
- semiconductor interactions with, **IV**:5.4*f*, 5.5 [see also Extreme ultraviolet (XUV) light]
- Ultraviolet (UV) radiation, **II**:34.6; **III**:14.23
- cataract from, **III**:7.5–7.6
- and color film, **II**:30.3
- damage from, **III**:7.2
- exposure limits for, **III**:7.9–7.10
- far, **II**:15.12, 15.13
- pterygium and droplet keratopathies from, **III**:7.6, 7.7
- spectrum of, **II**:25.2
- vacuum, **II**:24.3
- Umweganregung* (multiple-beam Bragg diffraction), **V**:43.6–43.8, 43.7*f*
- Unaccommodated eye, **III**:16.5
- Unbalanced nonlinear interferometers (UNI), **V**:20.22
- Unblazed gratings, **I**:20.3–20.4

- Uncertainty principle, **IV**:23.4, 23.6
- Uncertainty state, **IV**:23.6
- Uncrossed reflectors, **II**:39.38, 39.38*f*
- Underdense plasmas, strong field interactions with, **IV**:21.36–21.46
- direct laser acceleration and betatron resonance, **IV**:21.42–21.43
 - intense laser pulses, **IV**:21.38*f*, 21.38–21.39
 - inverse Bremsstrahlung heating, **IV**:21.37, 21.37*f*
 - ionization-induced defocusing, **IV**:21.43*f*, 21.43–21.44
 - ponderomotive channel formation, **IV**:21.42
 - self-channeling and self-phase modulation, **IV**:21.44–21.46, 21.45*f*
 - wakefield generation and electron acceleration, **IV**:21.39–21.42, 21.40*f*, 21.42*f*
- Underillumination, for input/output scanning, **I**:30.14
- Underscan, in color CRTs, **III**:22.12, 22.13*f*
- Underwater cameras, **I**:25.24–25.25
- Undulators, **V**:55.12*f*, 55.12–55.14, 55.13*f*, 58.1
- Unfolded reflections, in systems of revolution, **I**:1.32
- Uniaxial crystals, **IV**:8.8, 8.9*t*, 8.10*f*
- Uniaxial stress, **IV**:5.66*t*
- Unicam double-monochromator spectrometer, **I**:20.10, 20.13, 20.15*f*
- Unified Glare Rating (UGR), **II**:40.10–40.11, 40.11*t*
- Uniform color spaces, **III**:10.40, 10.42
- Uniform illumination, of nonimaging optics, **II**:39.22–39.41
- with classic projection systems, **II**:39.23*f*, 39.23–39.24
 - faceted structures in, **II**:39.39*f*, 39.39–39.41, 39.40*f*
 - integrating cavities in, **II**:39.24*f*, 39.24–39.27, 39.25*f*, 39.27*f*
 - lens arrays in, **II**:39.32–39.37, 39.33*f*–39.37*f*
 - lightpipes in, **II**:39.13*f*, 39.27–39.32, 39.28*f*–39.30*f*
 - tailored reflectors, **II**:39.37–39.39, 39.38*f*
- Uniformity:
- angular, **II**:39.31
 - of binary optics, **I**:23.7
 - control of, **II**:39.1–39.2
 - of luminance/illuminance, **II**:40.7, 40.13*f*
 - of photodetectors, **II**:24.20
 - and visual discomfort, **II**:40.9
- Unintentional (induced) prism, with contact lenses, **III**:20.31
- Unipolar chromatic mechanisms, **III**:11.80–11.81
- Unipolar mechanism (color vision), **III**:11.3
- Unique hues, **III**:11.27
- in color discrimination vs. color appearance tests, **III**:11.81–11.82
 - defined, **III**:11.3
 - and equilibrium colors, **III**:11.63–11.66
 - zero crossings, **III**:11.62–11.63
- Unit cell, crystal, **IV**:2.30
- Unit conversions:
- for English and SI units, **II**:37.7, 37.7*t*
 - for illuminance, **II**:36.7*t*, 36.8*t*
 - for photometric and radiometric quantities, **II**:36.11–36.14, 36.12*f*–36.14*f*
- Unit determinants, paraxial matrix methods for, **I**:1.67
- Unit magnification Galilean telescope, **II**:12.4, 12.4*f*
- Unitary matrix, **III**:18.2
- Unitraveling-carrier (UTC) photodiodes, **V**:13.68*f*, 13.68–13.69
- Unity divergence ratio, for AO modulators, **V**:6.32–6.33
- Univariance, **III**:10.4
- Univariant mechanism (color vision):
- defined, **III**:11.3
 - flicker, **III**:11.45
 - suprathreshold, **III**:11.58
- Universal antireflection coatings, **IV**:7.26, 7.27*f*
- Unlit-appearance modeling, **II**:40.21
- Unmodulated signal sources, **II**:27.12
- Unruh radiation, **V**:58.2
- Unstable resonators, **II**:16.25–16.26, 16.26*f*
- Unstrained QW lasers, **II**:19.15–19.16, 19.16*f*
- Uplight, **II**:40.43, 40.44*f*, 40.45
- Upper Atmospheric Research Satellite (UARS), **V**:3.36, 3.37*f*
- Upward plane irradiance, **IV**:1.5*t*, 1.7*f*, 1.8
- Upward scalar irradiance, **IV**:1.5*t*, 1.7*f*, 1.8
- Upwelling average cosine, **IV**:1.6*t*, 1.7*f*
- Upwelling irradiance, **IV**:1.7*f*, 1.8
- Ur (code), **V**:55.14
- Urbach tail model, **IV**:2.14–2.15
- Urbach's rule, **IV**:5.23
- Urea [(NH₄)₂CO], **IV**:2.40*t*, 2.45*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.64*t*
- Urgent (code), **V**:55.14

- U.S. Air Force Cambridge Research Laboratories, **V**:3.22
- U.S. Air Force three-bar target, **II**:4.6
- U.S. Department of Defense, **V**:5.27
- U.S. National Institute of Standards and Technology (NIST) database, **V**:3.26
- U.S. Standard Atmosphere, **V**:3.5, 3.6, 3.8*f*, 3.9*f*
- Useful life period, of LEDs, **II**:17.26, 17.26*f*
- User interfaces, human image features and, **III**:24.8–24.9
- Uviarc, **II**:15.28–15.29, 15.29*f*, 15.30*f*
- Vacuum, laser gain media in, **II**:16.36–16.37, 16.37*f*
- Vacuum heating, **IV**:21.47*f*, 21.48–21.49
- Vacuum lamps, **II**:34.31
- Vacuum spark source, of z-pinch radiation, **V**:57.3, 57.3*t*
- Vacuum ultraviolet (VUV) radiation, **II**:24.3; **IV**:14.3
- Vacuum ultraviolet (VUV) region, gratings and monochromators in, **V**:38.1–38.8, 38.2*f*, 38.3*f*, 38.4*t*–38.5*t*
- Vacuum ultraviolet (VUV) spectrum, **IV**:5.4*f*, 5.5
- Vacuum-metal interfaces, **IV**:4.43*f*
- Valence band (VB), **II**:17.3, 17.4, 17.4*f*; **IV**:18.3
- Valence lighting, **II**:40.13*f*
- Van Cittert-Zernike theorem, **I**:2.38, 2.39, 5.13–5.14, 5.17–5.19, 6.4, 6.8; **III**:6.12
- Van der Waals theory, **V**:8.23
- van Hove singularities, **IV**:8.28
- Vander Lugt filters, **I**:11.13–11.14
- Vanes (in stray light suppression), **II**:7.11–7.17
 defined, **II**:7.11–7.12, 7.12*f*, 7.13*f*
 placement design for, **II**:7.12*f*
 and scatter path, **II**:7.13*f*
 spacing and depth of, **II**:7.13–7.17, 7.14*f*–7.17*f*
- Vanishing point, **III**:13.4, 13.4*f*
- Vantage point, **III**:13.2
- Vapor axial deposition (VAD), **V**:25.3, 25.26
- Vapor exposure, in LED packaging, **II**:17.26
- Vapor phase epitaxy (VPE), **I**:21.17, 21.18; **II**:17.21, 17.22
- Variable optical attenuators (VOAs), **V**:21.12, 21.13*f*
- Variable retardation plates, **I**:13.53
- Variable temperature blackbody, **II**:15.10*f*
- Variable-angle spectroscopic ellipsometry (VASE), **I**:16.3
- Variable-orientation mirrors, **II**:6.17
- Varian (company), **II**:15.34
- Variational integral, of rays, **I**:1.19
- Varifocal systems, first-order layout for, **II**:1.11–1.12
- Vector aberration theory, **I**:29.5
- Vector diffraction, **I**:3.32*f*–3.37*f*, 3.32–3.37, 23.13–23.14, 23.14*f*
- Vector flux, **II**:39.21–39.22
- Vector Huygens secondary source (unit), **I**:3.33–3.36, 3.37*f*
- Vectors:
 electric field, **I**:2.3–2.4
 of space curves, **I**:1.18, 1.19
- Vector-scattering amplitude, **I**:7.8
- Vee (V) coupling, **IV**:14.1, 14.6*f*
- Vehicular lighting, **II**:40.63–40.67, 40.64*f*, 40.65*t*, 40.66*f*, 40.66*t*, 40.68*t*, 40.69*f*
- Veiling reflections, **II**:40.12, 40.14
- Velocity distribution measurement, **IV**:20.13, 20.13*f*
- Velocity-changing collisions, **IV**:11.15
- Velocity-selective coherent population trapping (VSCPT), **IV**:14.5, 20.37
- Verdet constant, **V**:18.7
- Vergence (disjunctive) eye movements, **III**:1.42, 1.44, 13.20, 13.21
- Vergence system, **III**:1.29, 1.30, 1.32, 1.34, 1.43–1.44
 binocular parallax, **III**:13.22
 defocus and effort to clear vision, **III**:13.22–13.23
 effect of lenses and prisms on, **III**:13.25–13.27
 with head mounted visual displays, **III**:13.31–13.32
 intrinsic stimuli to, **III**:13.21–13.22
 and variations of cross-coupling, **III**:13.23–13.24
 zone of clear and single binocular vision, **III**:13.24–13.25, 13.25*f*
- Verification (of tolerance), **II**:5.3
- VeriMask, **I**:33.19
- Vernier acuity, **III**:4.1
- Vernier acuity task, **III**:2.34, 2.35*f*
- Version eye movements (*see* Conjugate eye movements/position)
- Vertex, for figure of revolution, **I**:1.32
- Vertex curvature, **I**:1.32–1.33
- Vertex distance, **III**:12.9, 16.8

- Vertical alignment (VA) cells, **V**:8.25, 8.28f, 8.27–8.28
- Vertical antiblooming, **II**:32.9, 32.10f
- Vertical Bridgeman technique, **II**:17.21
- Vertical cavity lasers, **II**:19.41, 19.42f, 19.43t
- Vertical cavity semiconductor lasers, **II**:16.36
- Vertical cavity surface-emitting lasers (VCSELs), **II**:16.36; **V**:13.1–13.2, 13.43f
- commercial, **V**:13.48
- electrical injection and current confinement for, **V**:13.45, 13.45f
- light out vs. current in (L-I curve) of, **V**:13.46–13.47
- and linear optical amplifiers, **V**:19.14
- and mirror reflectivity, **V**:13.44
- and optical fibers, **V**:9.7n
- polarization of, **V**:13.48
- and quantum wells, **V**:13.43
- and spatial characteristics of emitted light, **V**:13.46
- spectral characteristics of, **V**:13.47–13.48
- Vertical coupling, of electroabsorption modulators, **V**:13.60
- Vertical half-wave linear retarders, Mueller matrices for, **I**:14.12t
- Vertical horopter, **III**:13.9
- Vertical illuminance, **II**:40.7
- Vertical linear polarizers, **I**:14.10t
- Vertical quarter-wave linear retarders, Mueller matrices for, **I**:14.12t
- Vertical scanning, in color CRTs, **III**:22.11–22.12, 22.12f
- Vertically integrated photodiode (VIP) FPAs, **II**:33.10
- Vertically illuminated *pin* photodiodes, **II**:26.3, 26.4f, 26.5, 26.10, 26.12f, 26.12–26.13
- Very dense crown glass, **IV**:2.42t
- Very light flint glass, **IV**:2.41t
- Very low loss reflectors, **IV**:7.41–7.42
- Very-long-wavelength infrared (VLWIR) radiation, **II**:24.3
- Very-long-wavelength semiconductor lasers, **II**:19.7–19.8
- Vestibulo-ocular reflex (VOR), **III**:13.20
- Vestibulo-ocular responses, **III**:1.42
- Vibration(s):
- lattice, **IV**:2.11–2.12, 2.76t–2.77t
- local, **IV**:5.17, 5.18, 5.19f, 5.20, 5.20f, 5.82f, 5.83
- phonon, **IV**:5.14–5.16
- planes of, **I**:12.6
- Vibration specifications, for lenses, **II**:4.10
- Vibrational optic effects, **IV**:5.17–5.20, 5.18f–5.19f, 5.20t, 5.21f
- Vibrational relaxation, **I**:31.14, 31.14f
- Vibrational sidebands, **I**:10.24
- Vibrational spectra, **I**:10.18–10.20
- Vibration-resistant optical reference cavity, **II**:22.16, 22.17f
- Video cameras, **I**:25.7–25.8
- Video display terminal (VDT), **III**:23.3
- Video head sets (*see* Head-mounted displays)
- Video monitors, in optical systems, **III**:5.16
- Video Quality Experts Group (VQEG), **III**:24.4
- Videocassette recorders (VCRs), **I**:35.1
- Video-enhanced differential-interference contrast (VE-DIC), **I**:28.41
- Vieth-Müller circle, **III**:13.8, 13.9, 13.9f
- Vieth-Müller horopter, **III**:2.40
- View cameras, **I**:25.19f, 25.18–25.20
- Viewing environments:
- for color CRTs, **III**:22.19–22.20
- for computer work, **III**:23.4–23.9
- lighting, **III**:23.4–23.5, 23.5t
- monitor characteristics, **III**:23.6–23.8
- screen reflections, **III**:23.5–23.6
- work habits, **III**:23.8–23.9
- workstation arrangement, **III**:23.8
- Vignetting, **I**:1.81, 1.81f, 1.82, 17.8, 29.6, 29.38; **II**:3.4, 34.19
- Virtual images, **I**:29.38
- Virtual phase CCDs, **II**:32.16f, 32.16–32.17
- Virtual pupils, **I**:1.76
- Virtual rays, **I**:1.10
- Virtual reality:
- head-mounted displays for, **III**:25.1–25.12
- characterizing, **III**:25.7–25.10, 25.8f, 25.9t–25.10t
- design considerations, **III**:25.2–25.7
- research in development of, **III**:24.9
- Virtual transitions, **IV**:16.4
- Virtual tributaries, of SONET, **V**:23.6
- Viruses, in water, **IV**:1.14
- Viscosities, of liquid crystals, **V**:8.24f, 8.23–8.25
- Visible array detectors, **II**:32.1–32.34
- about, **II**:32.2
- image sensing elements of, **II**:32.2–32.12, 32.3f
- antiblooming, **II**:32.9, 32.10f
- dark current, **II**:32.10–32.12, 32.11f
- junction photodiode, **II**:32.3–32.6, 32.4f, 32.6f

- Visible array detectors, image sensing elements
of (*Cont.*):
MOS capacitor, **II**:32.7–32.8
photoconductor, **II**:32.8–32.9
pinned photodiode, **II**:32.8
readout elements of, **II**:32.12–32.21
 CCD, **II**:32.12–32.20, 32.13*f*, 32.15*f*–32.18*f*
 MOS, **II**:32.20–32.21
sensor architectures of, **II**:32.21–32.34
 area image sensor arrays, **II**:32.24–32.32,
 32.25*t*, 32.26*f*–32.31*f*
 color imaging, **II**:32.32–32.34, 32.33*f*, 32.34*f*
 linear image sensor arrays, **II**:32.21–32.24,
 32.22*f*, 32.23*f*
- Visible light, hazards from, **III**:7.10
- Visible light photon counters (VLPCs), **II**:33.9
- Visible near-IR nonlinear optical crystals,
IV:10.21*t*, 10.22*t*
- Visible (VIS) radiation, **II**:24.3, 25.2
- Vision, **II**:40.3–40.6, 40.9
 biology of, **II**:40.3–40.4
 defined, **III**:23.4
 and perception, **II**:40.4–40.5
 photopic/scotopic/mesopic, **II**:34.37–34.39,
 37.2
 (*see also* Human eye)
- Vision experiments/research:
 displays for, **III**:22.1–22.40
 color cathode ray tubes (color CRTs),
 III:22.1–22.34
 liquid crystal displays (LCDs),
 III:22.34–22.40
 using Maxwellian view in (*see* Maxwellian
 view)
 (*see also* Psychophysical measurement)
- Vision laboratories, control of light sources in
 (*see* Optical generation of visual stimulus)
- Vision therapy, **III**:23.4
- Visual acuity, **III**:2.33–2.35, 2.34*f*, 2.35*f*,
 4.1–4.13, 12.6, 12.7
 age-related changes in, **III**:14.17, 14.19
 and color, **III**:4.10
 and contrast, **III**:4.12, 4.12*f*
 defined, **III**:4.1, 12.2, 15.2, 23.1
 and defocus, **III**:4.9, 4.9*f*, 4.10*f*
 factors affecting, **III**:4.9–4.13
 hyperacuity, **III**:4.13–4.16, 4.14*f*
 defined, **III**:4.1
 and superresolution, **III**:4.15
 and luminance, **III**:4.11, 4.11*f*, 4.12
- Visual acuity (*Cont.*):
 practice effect on, **III**:4.12
 and pupil, **III**:4.9
 resolution
 and information, **III**:4.15–4.16
 and superresolution, **III**:4.15
 resolving capacity of the eye, **III**:4.4*f*, 4.4–4.5
 contrast-transfer function, **III**:4.5
 point-spread function, **III**:4.4, 4.5
 spatial-frequency coordinates, **III**:4.5
 and retinal eccentricity, **III**:4.10–4.11, 4.11*f*
 retinal limitations, **III**:4.5–4.6, 4.6*f*
 and stage of development/aging, **III**:4.13
 stimulus specification, **III**:4.2–4.4, 4.3*f*
 and surround, **III**:4.12, 4.13*f*
 tests of visual acuity, **III**:4.7*f*, 4.7–4.8
 and time, **III**:4.12, 4.13*f*
 visual resolution threshold determination,
 III:4.6–4.7
- Visual angles:
 angular extent, **III**:5.2
 defined, **III**:2.3, 4.1, 10.1, 10.4*t*
 degrees of, **III**:10.9
 and image formation, **III**:2.3
 steradian of, **III**:6.3
- Visual axis, **III**:2.3
- Visual clarity, perception of, **II**:40.5
- Visual differences predictor, **III**:24.3
- Visual discomfort, **II**:40.9–40.12, 40.11*t*
- Visual discomfort probability (VCP), **II**:40.10
- Visual displays, measuring performance of
 (*see* Psychophysical measurement)
- Visual fields:
 age-related changes in, **III**:14.21, 14.21*f*
 binocular overlap of, **III**:13.7
 defined, **III**:13.2
 perception of, **III**:13.3
- Visual instruments:
 and accommodation response in eye,
 III:1.34–1.35
 and chromostereopsis, **III**:1.20
 exit pupil of, **III**:1.27–1.28
 retinal image quality with, **III**:1.27–1.28
 stereoscopic, **III**:1.41
- Visual magnification, **I**:1.28
- Visual performance, **III**:2.1–2.41
 aging-related changes in, **III**:14.14–14.22
 color vision, **III**:14.15, 14.17
 depth and stereovision, **III**:14.22

- Visual performance, aging-related changes in (*Cont.*):
 minimal, **III**:14.22
 sensitivity under scotopic and photopic conditions, **III**:14.15, 14.16*f*
 spatial vision, **III**:14.17–14.19, 14.18*f*
 temporal vision, **III**:14.19–14.21, 14.20*f*
 visual field, **III**:14.21, 14.21*f*
- binocular stereoscopic discrimination, **III**:2.40–2.41, 2.41*f*
- central visual processing, **III**:2.12–2.14
 anatomy of LGN, **III**:2.12–2.13, 2.13*f*
 physiology of LGN, **III**:2.13–2.14
- contrast detection, **III**:2.19–2.31
 adaptation and inhibition, **III**:2.25–2.27
 chromatic, **III**:2.29–2.31, 2.30*f*
 eye movements, **III**:2.21
 as function of spatial frequency, **III**:2.20*f*
- optical transfer function, **III**:2.21–2.22, 2.22*f*
- optical/retinal inhomogeneity, **III**:2.24, 2.25*f*
- receptors, **III**:2.22–2.23
- spatial channels, **III**:2.23–2.24
- temporal, **III**:2.27*f*, 2.27–2.29, 2.29*f*
- contrast discrimination, **III**:2.31*f*, 2.31–2.32
- contrast estimation, **III**:2.33
- contrast masking, **III**:2.32*f*, 2.32–2.33
- ideal-observer theory, **III**:2.16–2.19, 2.19*f*
- image formation, **III**:2.2–2.4
- image sampling by photoreceptors, **III**:2.4–2.9
- information-processing model for, **III**:2.15*f*, 2.15–2.16
- motion detection/discrimination, **III**:2.36–2.40
 optic flow fields, **III**:2.39*f*
 thresholds, **III**:2.37*f*
- pattern discrimination, **III**:2.35–2.36
- retinal processing, **III**:2.9–2.12
 anatomy of retina, **III**:2.9–2.11
 physiology of retina, **III**:2.11–2.12
- visual acuity, **III**:2.33–2.35, 2.34*f*, 2.35*f*
- Visual photometry, **II**:36.4
- Visual plane, **III**:13.2
- Visual resolution threshold, determination of, **III**:4.6–4.7
- Visual science, **II**:34.37
- Visual stimuli:
 in decision tasks, **III**:3.2
 optical generation of (*see* Optical generation of visual stimulus)
 in psychophysical experiments, **III**:3.3–3.4
 spatial frequency content of, **III**:5.8
- Visual stress, **III**:23.4
- Visual systems (colorimetry), **III**:10.38–10.40
- Visual tasks, **III**:2.15
- Visualization, in electronic imaging, **III**:24.8
- Visual-vestibular conflicts, in head mounted display systems, **III**:13.32–13.33
- Vitreous humor/gel, **III**:1.3*f*, 14.9, 16.3, 18.2
- Vivid color (VC) film, **II**:30.27
- VLSI-CMOS technology, silicon photonics for, **I**:21.15, 21.16
- Voigt function, **V**:56.7
- Voigt lineshape profiles, **V**:3.23
- Volkov state, **IV**:21.12
- Voltage, of x-ray tube sources, **V**:54.9
- Voltage amplifiers, **II**:27.10–27.11
- Voltage-controlled oscillators (VCOs), **V**:20.11, 20.11*f*
- Volume Bragg gratings (VBGs), **V**:25.29, 25.30
- Volume diffraction calculations, for MLLs, **V**:42.4–42.5, 42.5*f*
- Volume imaging ideal, **I**:1.29
- Volume scattering, **I**:9.1–9.17
 multiple scattering, **I**:9.8–9.17
 analytical theory of, **I**:9.9*f*, 9.9–9.10
 depolarization, **I**:9.16–9.17, 9.17*f*
 effective-medium representation, **I**:9.8
 radiative transfer, **I**:9.10–9.13, 9.11*f*, 9.13*f*
 speckle patterns, **I**:9.15*f*, 9.15–9.16
 weak localization, **I**:9.13–9.17, 9.14*f*
 single particle scattering vs., **I**:9.2–9.3
 and single scattering, **I**:9.4–9.8, 9.6*f*, 9.7*f*
 theory of, **I**:9.3–9.4, 9.4*f*
- Volume scattering function (VSFs):
 for sea water and ocean water, **IV**:1.34*t*–1.35*t*
 spectral, **IV**:1.37–1.38, 1.38*t*
 of water, **IV**:1.5*t*, 1.7*f*, 1.31, 1.31*t*, 1.32*f*, 1.33
 wavelength dependence of, **IV**:1.36*f*
- Volumetric reflection, **I**:7.7*f*
- Volumetric scattering cross section, **I**:7.7*f*
- Von Bezold spreading, **III**:11.3, 11.4*f*
- Von Kármán spectrum, **V**:3.28, 5.6
- Von Kries adaptation, **III**:11.3, 11.49, 11.51, 11.68
- V-parameter, of fiber lasers, **V**:25.3

- W. M. Keck Observatory, **V**:5.27
- W point, **IV**:9.4
- Wadsworth condition, **V**:38.6
- Wadsworth configuration, **I**:20.5*f*, 20.8, 20.12*f*, 20.14*t*
- Wafer processing, **II**:17.23–17.25
- Wakefield generation, **IV**:21.39–21.42, 21.40*f*
- “Walking” backward, by solitons, **IV**:15.28
- Wall effect, in gas detectors, **V**:63.32
- Wall slot lighting, **II**:40.13*f*
- Wall-grazing illumination, **II**:40.13*f*
- Wall-washing illumination, **II**:40.13*f*
- WALRUS objective, **I**:29.28
- Wannier excitons, **IV**:5.26–5.27, 5.26*t*, 8.31
- Water, **IV**:1.3–1.50
- absorption, **IV**:1.20–1.29
 - bio-optical models for, **IV**:1.27*f*, 1.27–1.29, 1.28*t*
 - by dissolved organic matter, **IV**:1.22–1.23, 1.23*t*
 - by organic detritus, **IV**:1.25–1.27, 1.25*t*, 1.26*f*
 - by phytoplankton, **IV**:1.23–1.25, 1.24*f*–1.25*f*
 - by polymers, **IV**:3.4
 - by sea water, **IV**:1.21, 1.22*t*
 - apparent optical properties, **IV**:1.12–1.13
 - attenuation
 - beam, **IV**:1.40–1.41, 1.41*f*, 1.42*f*
 - diffuse and Jerlov water types, **IV**:1.42–1.46, 1.43*t*–1.45*t*, 1.44*f*, 1.45*f*
 - constituents of natural waters, **IV**:1.13–1.15
 - dissolved substances, **IV**:1.13
 - particulate substances, **IV**:1.14–1.15
 - electromagnetic properties of, **IV**:1.16–1.17, 1.18*t*
 - inherent optical properties, **IV**:1.9–1.12, 1.10*f*
 - irradiance reflectance and remote sensing, **IV**:1.46–1.47, 1.47*f*
 - particle size distributions, **IV**:1.15–1.16, 1.16*f*
 - pure, **IV**:1.3
 - radiometric quantities, **IV**:1.4–1.9, 1.5*t*–1.6*t*, 1.7*f*
 - refraction index, **IV**:1.18–1.20
 - particles, **IV**:1.20
 - sea water, **IV**:1.18*f*, 1.18–1.20, 1.19*t*–**IV**:1.20*t*
 - scattering, **IV**:1.30
 - inelastic and polarization, **IV**:1.47–1.49, 1.48*f*, 1.49*f*
 - measurement of, **IV**:1.29–1.30
- Water, scattering (*Cont.*):
- by particles, **IV**:1.30–1.35, 1.31*t*, 1.32*f*, 1.33*f*, 1.34*t*–1.35*t*
 - by pure water and pure sea water, **IV**:1.30
 - wavelength dependence of, **IV**:1.35–1.40, 1.35*t*, 1.36*f*, 1.37*t*, 1.38*t*, 1.39*f*, 1.40*t*
 - scattering in, **V**:63.10, 63.10*t*, 63.11*f*
 - in standard atmosphere, **V**:3.6, 3.10–3.11, 3.11*f*
 - terminology and notation, **IV**:1.3–1.4
- Water vapor regained (WVR), **IV**:6.17
- Watercolor illusion/effect, **III**:11.72, 11.73*f*
- Watt (unit), **II**:39.2*t*
- Wave aberration, **III**:15.2
- Wave equation, **I**:12.4–12.6
- Wave equations, for light propagation in solids, **IV**:8.6–8.7
- Wave interactions, photorefractive effect and, **IV**:12.4*f*, 12.4–12.7, 12.6*f*
- Wave modulation distance meter, **II**:12.5, 12.6*f*
- Wave normals, **I**:13.5
- Wave propagation, reciprocity of, **V**:4.9–4.10
- Wave structure function, **V**:4.6, 5.9*n*
- Waveband materials, **II**:8.3*t*, 8.4*t*
- Waveband structure of semiconductors, **II**:17.3*f*–17.5*f*, 17.3–17.6
- Wavefront aberration, **I**:1.86–1.88; **III**:1.15–1.19, 1.16*f*, 1.16*t*, 1.17*f*
- and defocus of visual instruments, **III**:1.28
 - and final retinal image quality, **III**:1.21–1.22
- Wavefront aberration coefficients, **I**:1.90
- Wavefront correctors (AO systems for the eye), **III**:15.3, 15.7*f*, 15.9–15.12, 15.10*f*, 15.11*f*
- Wavefront correctors, for adaptive optics, **V**:5.37–5.38, 5.38*f*
- Wavefront division, **I**:2.14
- Wavefront division, interference by, **I**:2.14–2.19, 2.15*f*–2.18*f*
- Wavefront error (*W*), **II**:4.1, 4.3, 4.7, 4.8; **V**:4.35, 5.40–5.41
- Wavefront measurement, aspherical (*see* Aspherical wavefront measurement)
- Wavefront multiplexers, **I**:23.11, 23.11*t*, 23.12
- Wavefront quality, of binary optics, **I**:23.8
- Wavefront sensing techniques, **V**:5.36–5.37
- Wavefront sensor (AO systems for the eye), **III**:15.3, 15.7*f*, 15.8*f*, 15.8–15.9
- Wavefront stitching, **II**:13.27, 13.27*f*
- Wavefront tilt (*Z*-tilt), **V**:4.3, 4.23–4.25, 5.14, 5.15

- Wavefront tolerancing, **II**:5.3–5.7, 5.4*f*, 5.5*f*, 5.5*t*, 5.6*t*, 5.7*f*
- Wavefronts, **I**:3.4
- aberrated, **I**:2.12, 2.13
 - cylindrical, **I**:3.13–3.21, 3.14*f*
 - Cornu's spiral, **I**:3.16–3.19
 - opaque strip construction, **I**:3.20–3.21
 - from rectangular apertures, **I**:3.19–3.20
 - from straight edge, **I**:3.14–3.16
 - disturbance of, **I**:3.5–3.6
 - for cylindrical wavefronts, **I**:3.13–3.14, 3.14*f*
 - and straight edges, **I**:3.14–3.15
 - geometrical, **I**:1.12–1.13
 - gradients of, **V**:5.23
 - and interference, **I**:2.4–2.5
 - from lenses, **II**:4.3–4.5, 4.5*t*
- Waveguide bends, **IV**:9.14–9.16, 9.15*f*, 9.16*f*
- Waveguide confinement factor, **V**:13.4*f*, 13.5
- Waveguide grating routers (WGRs), **I**:21.24
- Waveguide intersections, **IV**:9.16–9.17, 9.17*f*
- Waveguide modulators, **V**:7.30–7.32, 7.31*f*–7.33*f*, 13.56, 13.57*f*
- Waveguide parameter (*V*-parameter), **III**:8.2, 8.17, 8.20–8.21
- Waveguide photodetectors, **II**:26.4*f*, 26.5
- Waveguide *pin* photodiodes, **II**:26.13–26.14, 26.14*f*
- Waveguides:
- attenuated total reflectance, **V**:12.11
 - biological (*see* Biological waveguides)
 - double heterostructure, **V**:19.3*f*
 - evanescent-wave coupled pin, **V**:13.68
 - integrated optics, **I**:21.3*f*–21.5*f*, 21.3–21.8, 21.7*f*
 - leaky, **I**:21.3
 - and photonic bandgaps, **IV**:9.12–9.17
 - in photonic crystals with 2D periodicity, **IV**:9.13*f*, 9.13–9.14
 - waveguide bends, **IV**:9.14–9.16, 9.15*f*, 9.16*f*
 - waveguide intersections, **IV**:9.16–9.17, 9.17*f*
 - photorefractive, **IV**:12.37
 - silicon-on-insulator planar, **IV**:22.15
 - of SOAs, **V**:19.15*f*–19.16*f*, 19.15–19.16
- Wavelength(s):
- attenuation vs., **V**:15.7
 - Bragg, **V**:17.3, 20.15, 24.7
 - color data as functions of, **III**:10.21, 10.23, 10.23*f*
 - in fiber optics, **II**:17.33–17.34
- Wavelength(s) (*Cont.*):
- of liquid crystals, **V**:8.19–8.22, 8.20*f*
 - and modulation transfer function, **I**:17.38, 17.39
 - in optical systems, **III**:5.11–5.13, 5.12*t*
 - of plane waves, **I**:2.4
 - pump, **V**:14.5–14.6
 - transverse effective, **V**:11.10
- Wavelength blockers (WB), **V**:21.12, 21.13*f*
- Wavelength dependence, of scattering, **IV**:1.35–1.40, 1.35*t*, 1.36*f*, 1.37*t*, 1.38*t*, 1.39*f*, 1.40*t*
- Wavelength dispersion, **V**:9.6–9.7
- Wavelength dispersive detectors (WDS), **V**:62.2
- Wavelength division multiplexing (WDM):
- dense
 - and AOTFs, **V**:6.43, 6.44
 - and optical fiber amplifiers, **V**:14.1
 - and SOAs, **V**:19.25*f*, 19.25–19.27, 19.26*f*
 - and dispersion-managed solitons, **V**:22.15–22.17
 - and ESCON standard, **V**:23.2
 - and fiber-based couplers, **V**:16.1, 16.4
 - and solitons, **V**:22.2, 22.8–22.12, 22.15–22.17 (*see also* Wavelength-division multiplexing networks)
- Wavelength division multiplexing (WDM) systems:
- fabrication of, **I**:21.14
 - filters for, **I**:21.23
 - in integrated optics, **I**:21.37*f*–21.39*f*, 21.37–21.38
- Wavelength errors, **II**:34.36
- Wavelength interval (between oscillations), **I**:12.10
- Wavelength modulation, **IV**:5.66*t*
- Wavelength scans, for scatterometers, **V**:1.14
- Wavelength-dispersive x-ray fluorescence (WDXRF), **V**:29.2*f*, 29.2–29.3
- Wavelength-division multiplexing (WDM) networks, **V**:18.4, 18.6, 21.1–21.44
- architecture of
- circuit and packet switching, **V**:21.7–21.11, 21.8*f*–21.11*f*
 - network reconfigurability, **V**:21.12*f*, 21.12–21.13, 21.13*f*
 - point-to-point links, **V**:21.4
 - star, ring, and mesh topologies, **V**:21.5*f*–21.7*f*, 21.5–21.7
 - wavelength-routed networks, **V**:21.5, 21.5*f*

- Wavelength-division multiplexing (WDM) networks (*Cont.*):
 fiber bandwidth, **V**:21.2*f*, 21.2–21.3
 fiber system impairments, **V**:21.13–21.26
 chromatic dispersion, **V**:21.14–21.16, 21.15*f*, 21.16*f*
 dispersion and nonlinearities
 management, **V**:21.20–21.26, 21.21*f*, 21.23*f*–21.27*f*
 fiber attenuation and optical power loss, **V**:21.13–21.14
 fiber nonlinearities, **V**:21.18–21.20, 21.19*f*, 21.20*f*
 polarization-mode dispersion, **V**:21.16–21.18, 21.17*f*–21.18*f*
 history of, **V**:21.1–21.2
 optical amplifiers in, **V**:21.37*f*, 21.37–21.44
 EDFA, **V**:21.38*f*–21.42*f*, 21.38–21.41
 Raman, **V**:21.42*f*–21.44*f*, 21.42–21.44
 and optical fibers, **V**:9.8, 9.12, 9.13, 9.15
 optical modulation formats for, **V**:21.27–21.36
 basic concepts, **V**:21.27–21.29, 21.28*f*–21.30*f*
 carrier-suppressed return-to-zero and duobinary, **V**:21.30–21.33, 21.31*f*, 21.32*f*
 comparisons of, **V**:21.36, 21.36*t*, 21.37*t*
 DPSK and DQSK, **V**:21.33*f*–21.35*f*, 21.33–21.36, 21.36*t*, 21.37*t*
 in real systems, **V**:21.3*f*, 21.3–21.4, 21.4*f*
 Wavelength-routed networks, **V**:21.5, 21.5*f*
 Wavelengths, de-Broglie, **IV**:8.4
 Wavelength-selective couplers (WSCs), **V**:14.2
 Wavelength-selective switches (WSS), **V**:21.12, 21.13*f*
 Waveplates, **I**:15.8
 Waves, **I**:2.3–2.6
 amplitudes of, **I**:2.4, 2.5, 12.5
 diffraction of, **I**:3.2–3.3, 3.3*f*
 interference of, **I**:2.5–2.6
 plane, **I**:2.4, 2.5*f*, 3.3, 3.17
 decomposition of, **I**:3.23
 interference of, **I**:2.8–2.9, 2.9*f*
 and spherical waves, **I**:2.9–2.11, 2.10*f*
 spherical, **I**:2.4, 2.5*f*, 3.2–3.3
 interference of, **I**:2.11–2.12, 2.12*f*, 2.13*f*
 and plane waves, **I**:2.9–2.11, 2.10*f*
 Waviness, in polycapillary x-ray optics, **V**:53.4, 53.4*f*
 Weak polarization elements (Mueller matrices), **I**:14.26–14.27
 Wearout period, **II**:17.26, 17.26*f*
 Weber contrast, in vision experiments, **III**:3.4
 Weber's law, **III**:11.16*f*
 and cone contrast spaces, **III**:11.32
 and contrast coding, **III**:11.15–11.16, 11.16*f*
 defined, **III**:11.3
 first-site adaptation, **III**:11.15, 11.16*f*
 second-site adaptation, **III**:11.18
 Webvision, **III**:11.5
 Wedge filters, **IV**:7.90, 7.91, 7.91*f*
 Wedged multilayer Laue lenses (wMLLs), **V**:42.12–42.13, 42.13*f*, 42.14*f*
 Weighting functions, **II**:36.17
 Welder's flash, **III**:1.9, 7.4
 Weld-line, in molding, **IV**:3.14
 Well and barrier intermixing, **V**:13.60
 Well capacity, **II**:25.11
 Welsbach mantle, **II**:15.17, 15.18
 Wernicke prisms, **I**:20.6*f*
 Wet age-related macular degeneration, **III**:14.1, 14.24
 Wetherell and Womble objectives, **I**:29.31
 Weyl's integral, **I**:5.17
 Weyl's plane-wave decomposition, **I**:3.23
 Whiffletrees (lever mechanisms), **II**:6.19
 White balance, in color CRTs, **III**:22.9
 White light, **II**:18.4*f*, 18.4–18.5, 40.7, 40.8, 40.24
 White noise, **III**:18.2
 White surfaces, reflectivity of, **II**:17.31
 White-light LEDs, **II**:40.37, 40.38
 WI 9 lamps, **II**:15.21*f*
 WI 14 lamps, **II**:15.21*f*
 WI 16/G lamps, **II**:15.21*f*, 15.22*f*
 WI 17/G lamps, **II**:15.22*f*
 WI 40/G lamps, **II**:15.22*f*
 WI 41/G lamps, **II**:15.22*f*
 Wide area networks (WANs), **V**:9.14, 21.7
 Wide bandpass filters, **IV**:7.90, 7.90*f*
 Wide-angle bandpass filters, **IV**:7.93*f*, 7.93–7.94
 Wide-angle models (of human eye), **III**:1.38
 Wide-angle photography:
 cameras for, **I**:25.26*f*, 25.25
 lenses for, **I**:27.2, 27.5*f*
 with nonrectilinear distortion, **I**:27.6
 with rectilinear distortion correction, **I**:27.6, 27.13*f*, 27.14*f*
 Wideband AO Bragg cells, **V**:6.27, 6.30, 6.31*t*
 Wide-field lobster-eye telescopes, **V**:48.4

- Wide-field objective with Maksutov correction, **I**:29.20
- Wieman-Hansch experiment, **I**:31.26
- Wiener filter, **I**:11.15–11.17
- Wiener spectrum, **II**:29.21
- Wiener-Khinchine theorem, **I**:5.5, 5.21
- Wien's displacement law, **II**:15.7, 34.23, 34.24, 37.11
- Wien's wavelength displacement law, **I**:31.4
- Wigglers, **V**:55.11–55.12
- Wigner distribution function, **I**:5.8
- Williamson construction, **II**:39.12–39.13, 39.13*f*, 39.28, 39.29*f*, 39.31, 39.32
- Wilson H. R., **III**:2.33
- Winchester heads, **I**:35.6
- Window structure, of SOAs, **V**:19.8, 19.8*f*
- Windowing, in x-ray mirror metrology, **V**:46.9, 46.10*f*
- Window/photocathode assemblies, of image intensifiers, **II**:31.10–31.12, 31.11*f*, 31.12*f*
- Windows:
 - and daylight sources, **II**:40.41, 40.47, 40.48, 40.49*f*, 40.50*f*
 - mounting of optical, **II**:6.11, 6.11*f*, 6.12*f*
- Wire array sources, of z-pinch radiation, **V**:57.3*t*, 57.3–57.4, 57.4*f*
- Wire grids, **I**:13.30*n*
- Wire-grid polarizers, **I**:13.30–13.33, 13.31*f*, 13.32*t*
- Wire-wound thermopile arrays, **II**:24.23
- With-the-rule astigmatism, **III**:1.6
- Wolf shift, **I**:5.21
- Wollaston prisms, **I**:13.7, 13.18, 13.18*f*, 13.21, 13.24; **V**:46.4
 - for DIC microscopes, **I**:28.39, 28.40
 - in Nomarski interferometers, **I**:32.4
 - and nondispersive prisms, **I**:19.3*t*, 19.15, 19.15*f*
- Wolter geometries, **V**:63.21
- Wolter mirror configurations, **V**:52.4, 64.6, 64.6*f*
- Wolter optics, **V**:26.10, 48.1, 48.2, 49.4–49.6, 49.5*f*–49.6*f*
- Wolter telescopes, **V**:44.4, 44.5*f*, 44.6–44.10, 44.7*f*, 44.10*f*, 45.1–45.5, 45.2*f*
- Wolter x-ray optics, **V**:47.2*f*, 47.2–47.7, 47.4*f*, 47.6*f*, 47.7*f*
- Wolter-Schwarzschild (WS) telescopes, **V**:44.7, 44.11, 45.1
- Wood lens, **III**:19.3–19.4, 19.4*f*
- Work function (of photons), **II**:25.2
- Working distance:
 - of objective lenses, **I**:28.13
 - retinoscopy, **III**:12.2, 12.6
- Workstation arrangement, for computer work, **III**:23.8
- World Health Organization (WHO), **III**:7.9
- Wright objectives, **I**:29.15
- Write-once-read-many (WORM) technology, **I**:35.2
- Wulfenite (PbMoO₄), **IV**:2.40*t*, 2.45*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.64*t*, 2.72*t*
- Wurtzite (α -ZnS):
 - in crystals and glasses, **IV**:2.70*t*
 - lattices of, **IV**:5.6
 - properties of, **IV**:2.41*t*, 2.44*t*, 2.46*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.66*t*, 2.69*t*, 2.70*t*
- Wysocki lens, **III**:9.5
- Xenon lamps, **II**:15.34*f*, 15.35*f*, 40.31, 40.35*f*
- Xerographic systems, **I**:34.1–34.13, 34.2*f*
 - cleaning and erasing in, **I**:34.10
 - color in, **I**:34.11*f*–34.13*f*, 34.11–34.12
 - control of, **I**:34.11
 - development in, **I**:34.5*f*–34.9*f*, 34.5–34.10
 - fusing in, **I**:34.10
 - and latent image, **I**:34.1–34.4, 34.2*f*–34.3*f*
 - transfer in, **I**:34.10
- XEUS/IXO mission, **V**:49.7
- XMM mission, **V**:44.6
- XMM-Newton observatory, **V**:33.3, 33.3*t*, 47.2, 47.4*f*, 47.6, 47.6*f*
- X-pinch sources of radiation, **V**:57.3, 57.3*t*, 57.4
- X-ray absorption fine structure (EXAFS), **V**:29.3, 60.4
- X-ray absorption near edge spectroscopy (XANES), **V**:30.2, 30.2*f*, 30.4*f*
- X-ray absorption near-edge structure (XANES), **V**:29.3, 60.4
- X-ray absorption spectroscopy (XAS), **V**:30.1–30.5, 30.2*f*–30.5*f*
- X-ray astronomy, **V**:33.1–33.4, 33.3*t*
- X-ray astronomy satellite (SAX), **V**:44.6
- X-ray detectors, **V**:60.3–60.10
 - cryogenic, **V**:60.8–60.9, 60.9*t*, 60.10*t*
 - film, **V**:60.8, 60.9*t*, 60.10*t*
 - ionization, **V**:60.3–60.7, 60.9*t*, 60.10*t*
 - scintillation, **V**:60.7–60.8, 60.9*t*, 60.10*t*
- X-ray diffraction (XRD), **V**:28.1–28.7, 28.3*f*–28.6*f*

- X-ray fluorescence (XRF), V:29.1–29.11, 54.8
 energy-dispersive, V:29.3–29.11, 29.5*f*, 29.6*f*,
 29.8*f*–29.11*f*
 history of, V:29.1
 and polycapillary x-ray optics,
 V:53.10–53.11, 53.11*f*
 wavelength-dispersive, V:29.2*f*, 29.2–29.3
 and x-ray diffraction, V:28.1
 and x-ray imaging, V:62.5, 62.6*f*
- X-ray imaging detectors, V:61.1–61.8, 61.2*f*
 CCD detectors, V:61.7–61.8, 61.8*f*
 flat panel detectors, V:61.3–61.7, 61.4*f*,
 61.6*t*, 61.7*f*
- X-ray lasers, II:16.31; V:58.1–58.4, 58.3*f*
- X-ray mapping, V:62.4, 62.5, 62.6*f*
- X-ray microscopes, V:37.6
- X-ray mirrors:
 and grazing incidence optics, V:44.3–44.6,
 44.3*f*–44.6*f*
 metrology of, V:46.1–46.12
 history of, V:46.1–46.2
 profile analysis considerations,
 V:46.6–46.12, 46.7*f*, 46.10*f*
 surface figure metrology, V:46.3–46.6, 46.5*f*
 surface finish metrology, V:46.2
- X-ray monochromators, V:30.1–30.4, 50.6–50.7
- X-ray observatories, V:33.1–33.4, 33.3*t*
- X-ray optics:
 adaptive, V:50.1–50.8
 hard vs. soft x-ray telescopes, V:50.2
 history of, V:50.1, 50.2*f*
 synchrotron and lab-based sources,
 V:50.2–50.8, 50.4*f*, 50.6*f*, 50.8*f*
 astronomical, V:47.1–47.11
 angular resolution of, V:47.10–47.11, 47.11*f*
 hard, V:47.9–47.10
 history of, V:47.1–47.2
 Kirkpatrick-Baez optics, V:47.7–47.8,
 47.8*f*, 47.9*f*
 Wolter, V:47.2*f*, 47.2–47.7, 47.4*f*,
 47.6*f*, 47.7*f*
 coherent, V:27.1–27.5, 27.3*f*–27.5*f*
 gratings and monochromators in,
 V:38.1–38.8
 diffraction properties, V:38.1–38.3, 38.2*f*
 dispersion properties, V:38.6–38.7
 efficiency, V:38.8
 focusing properties, V:38.3*f*, 38.3–38.6,
 38.4*t*–38.5*t*
 resolution properties, V:38.7
- X-ray optics (*Cont.*):
 and inverse Compton x-ray sources,
 V:59.1–59.4, 59.3*t*
 and medical imaging, V:31.1–31.4,
 31.2*f*–31.4*f*
 monocapillary, V:52.1–52.6, 52.2*f*–52.5*f*,
 52.2*t*
 multifoil, V:48.1–48.4, 48.2*f*, 48.3*f*
 and neutron optics, V:26.5–26.11, 26.8*f*,
 26.9*f*, 26.11*f*, 36.2*f*
 polycapillary, V:53.1–53.19, 53.2*f*
 alignment and measurement, V:53.5–53.8,
 53.6*f*, 53.7*f*
 collimation, V:53.8*f*–53.9*f*, 53.8–53.9
 energy filtering, V:53.10
 focusing, V:53.9–53.10, 53.10*t*
 monochromatic imaging, V:53.16–53.17,
 53.17*f*
 powder diffraction, V:53.14, 53.14*f*
 radiation resistance, V:53.5
 scatter rejection in imaging,
 V:53.14–53.16, 53.15*f*–53.16*f*
 scintigraphy, V:53.17, 53.18, 53.18*f*
 simulations and defect analysis,
 V:53.3*f*–53.5*f*, 53.3–53.5
 single crystal diffraction, V:53.12*f*,
 53.12–53.14, 53.13*f*
 therapy, V:53.18–53.19
 x-ray fluorescence, V:53.10–53.11, 53.11*f*
 and pore optics, V:49.1–49.7
 ray tracing for, V:35.1–35.6, 35.4*f*, 35.5*f*
 and Schwarzschild objective, V:51.3,
 51.3*f*–51.4*f*
 spectral detection and imaging in,
 V:62.1–62.5, 62.2*f*–62.6*f*
 and x-ray properties of materials, V:36.1–36.9
 Auger energies, V:36.3*t*, 36.9*t*
 electron binding energies, V:36.3*t*–36.6*t*
 photoabsorption and scattering, V:36.1
 photon energies, V:36.7*t*–36.8*t*
 x-ray and neutron optics, V:36.2*f*
- X-ray region:
 beam splitters for, IV:7.63
 soft
 bandpass filters for, IV:7.94–7.96,
 7.95*f*–7.96*f*
 interference polarizers for, IV:7.73,
 7.76*f*–7.77*f*
 multilayer reflectors for, IV:7.42, 7.53
- X-ray scattering, I:9.6

- X-ray spectral detection and imaging, **V**:62.1–62.5, 62.2*f*–62.6*f*
- X-ray tube sources, **V**:54.3–54.16
 brightness and intensity of, **V**:54.11–54.15, 54.13*f*–54.15*f*
 cathode design and geometry, **V**:54.10–54.11
 characteristics of, **V**:54.3, 54.4*f*
 optimization of, **V**:54.15–54.16
 spectra of, **V**:54.4–54.10, 54.5*f*, 54.8*f*
- X-rays:
 circularly polarized soft, **V**:55.6–55.7
 nanofocusing of hard, **V**:42.1–42.17
 history of, **V**:42.2*f*, 42.2–42.4, 42.3*f*
 instrumental beamline arrangement and measurements for, **V**:42.9*f*–42.12*f*, 42.9–42.10
 limitations of, **V**:42.15–42.17, 42.16*f*–42.17*f*
 with magnetron-sputtered MLLs, **V**:42.5–42.7, 42.6*f*–42.8*f*
 on MLLs with curved interfaces, **V**:42.14, 42.15*f*
 Takagi-Taupin calculations for, **V**:42.12–42.14
 volume diffraction calculations for, **V**:42.4–42.5, 42.5*f*
 with wedged MLLs, **V**:42.12–42.13, 42.13*f*, 42.14*f*
 polarization of, **V**:43.1–43.2
- X-Y addressing, **II**:33.16
- Y-branch interferometric modulators, **V**:13.51–13.52
- Y-coupled junctions, **II**:19.27, 19.29
- Yellow filter dyes, **II**:30.4
- Yellow light, **II**:29.13, 29.13*f*
- Yellow matter, **IV**:1.13, 1.21–1.23, 1.28
- Yes-no judgments, **III**:3.6
- Yield strength, of metals, **IV**:4.8, 4.70, 4.70*t*
- Yoked eye movements, **III**:13.2, 13.20
- Yolo objectives, **I**:29.26
- Young-Helmholtz trichromatic theory, **III**:11.5
- Young's astigmatic formulae, **I**:1.44
- Young's double slit experiment, **I**:2.14–2.15, 2.15*f*
- Young's fringes, **I**:13.45; **III**:1.22
- Young's modulus:
 for crystals, **IV**:2.30, 2.31
 of infrared optical fibers, **V**:12.3*t*, 12.9
 for metals, **IV**:4.7, 4.10*t*, 4.69*t*
 for molded microlenses, **I**:22.12*t*
 for polymers, **IV**:3.3
- Young's two pinhole interferometer, **I**:6.3
- Young-Thollon half prisms, **I**:20.7*f*
- Ytterbium-doped fiber amplifiers (YDFAs), **V**:14.2, 14.7
- Ytterbium-doped fibers, **V**:25.23*t*, 25.24–25.25, 25.31, 25.33
- Yttria (Y₂O₃), **IV**:2.40*t*, 2.44*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.65*t*, 2.76*t*
- Yttrium aluminum garnet (Y₃Al₅O₁₂) (YAG), **IV**:2.44*t*, 2.53*t*, 2.58*t*, 2.65*t*, 2.76*t*
- Yttrium aluminum garnet (YAG) phosphor, **II**:18.4
- Yttrium lithium fluoride (LiYF₄) (YLF), **IV**:2.45*t*, 2.48*t*, 2.52*t*, 2.57*t*, 2.63*t*, 2.72*t*
- Yttrium vanadate (YVO₄), **IV**:2.40*t*, 2.45*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.65*t*, 2.76*t*
- Yurke state, **IV**:23.8
- ZBLAN (fluorozirconate glass), **IV**:2.43*t*, 2.49*t*, 2.54*t*, 2.59*t*, 2.68*t*; **V**:12.5, 12.5*f*
 and fiber lasers, **V**:25.3, 25.24, 25.27*t*, 25.28
 fluoraluminate glass vs., **V**:12.3*t*, 12.4, 12.4*t*
- Z-dependent oscillary term, **III**:8.22, 8.23
- Zeeko (company), **II**:9.5–9.6
- Zeeman effect, **I**:31.17, 31.18*f*, 31.19–31.21; **V**:2.21
- Zeeman spectroscopy, **V**:2.23–2.24, 2.24*f*
- Zeiss Infinity Color-Corrected Systems Plan Apo, **I**:28.15, 28.15*f*
- Zeiss prism system, **I**:19.3*t*, 19.16, 19.16*f*
- Zeiss sheet polarizers, **I**:13.26
- ZEMAX (optical software), **II**:7.26–7.27; **V**:35.1
- Zenger prisms, **I**:20.6*f*
- Zernike annular expansion coefficients, **V**:4.25
- Zernike coefficient(s):
 for RMS wavefront error, **III**:1.15–1.19, 1.16*t*, 1.17*f*
 variance of, **V**:4.22, 4.22*f*, 4.23*f*, 4.36
- Zernike decomposition elements, **V**:5.37
- Zernike dispersion formula, **IV**:2.22
- Zernike modes of aberrations, **V**:5.11, 5.11*t*, 5.12*f*
- Zernike phase-contrast test, **II**:13.13–13.14, 13.14*f*
- Zernike polynomials, **I**:1.90, 23.3; **II**:5.9; **V**:4.17–4.20, 4.20*t*, 5.10, 46.6
 annular, **II**:11.13–11.21, 11.14*f*, 11.16*t*–11.21*t*
 circle, **II**:11.4, 11.6–11.12, 11.8*t*–11.9*t*, 11.9*f*–11.11*f*, 11.12*t*, 11.39
 defined, **III**:15.2
 for representing ocular aberrations, **III**:15.4, 15.5*f*, 16.6–16.7, 16.7*t*

- Zernike tilts, **V**:4.23–4.26
- Zero dispersion point, for glasses, **IV**:2.23
- Zero order, of radiation, **V**:40.1
- Zerodur glass, **IV**:2.43*t*, 2.49*t*, 2.54*t*
- Zerodur glass shells, **V**:47.5
- Zerodur prisms, **II**:6.16, 6.16*f*
- Zero-phonon lines, in solids, **V**:2.13–2.14, 2.14*f*
- Zero-phonon transitions, **I**:10.24*f*, 10.24–10.26
- Zeros, in paraxial matrices, **I**:1.68
- ZEST, **III**:3.9
- Zinc, **II**:17.23; **IV**:4.48*t*
- Zinc crown glass, **IV**:2.41*t*
- Zinc diffusion, **II**:17.9–17.10, 17.10*f*
- Zinc doping, **II**:17.20
- Zinc oxide (ZnO) doped GaP, **II**:17.16, 17.21–17.22
- Zinc selenide (ZnSe), **II**:17.19; **IV**:2.41*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.66*t*, 2.69*t*
- Zinc telluride (ZnTe), **IV**:2.41*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.66*t*, 2.69*t*
- Zinblende (β -ZnS):
 in crystals and glasses, **IV**:2.69*t*
 lattices of, **IV**:5.6, 5.16, 5.17*t*
 properties of, **IV**:2.41*t*, 2.44*t*, 2.46*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.65*t*, 2.69*t*, 2.70*t*
- Zinc-doped germanium (Ge:Zn) detectors, **II**:24.84*f*, 24.98–24.100, 24.99*f*
- Zinc-germanium diphosphide (ZnGeP₂), **IV**:2.41*t*, 2.45*t*, 2.48*t*, 2.53*t*, 2.58*t*, 2.65*t*
- Zirconia, cubic (ZrO₂:0.12Y₂O₃), **IV**:2.41*t*, 2.44*t*, 2.48*t*, 2.69*t*
- Zirconium arc lamps, **II**:15.47, 15.48*f*
- ZO-MOD BLACK, **IV**:6.56
- Zonal approach, to wavefront error correction, **V**:4.35
- Zonal cavity lighting simulation, **II**:40.17
- Zonal refractive lenses (IOLs), **III**:21.15
- Zonal spherical aberrations, **I**:29.8, 29.38
- Zone models (color vision), **III**:11.5
 three-stage, **III**:11.6, 11.7*f*, 11.8, 11.82–11.85, 11.85*f*
 two-stage, **III**:11.5, 11.6*f*
- Zone plate law, **V**:42.4
- Zone plates, **I**:3.11–3.13, 3.12*f*; **V**:40.1–40.10
 amplitude, **V**:40.4–40.5, 40.5*f*
 and Bragg-Fresnel lenses, **V**:40.9*f*, 40.9–40.10
 diffraction efficiencies of, **V**:40.4–40.8, 40.5*f*, 40.7*f*, 40.8*f*
 Fresnel, **V**:40.2, 42.3, 42.3*f*, 55.16
 geometry of, **V**:40.1–40.3, 40.2*f*
 neutron, **V**:63.25
 phase, **V**:40.5–40.7, 40.7*f*
 as thin lenses, **V**:40.3–40.4
- Zonules, **III**:1.3*f*, 21.2
- Zoom lenses, **I**:27.17, 27.20*f*–27.22*f*
- Zoom systems, **II**:1.11–1.12, 3.20
- Zooplankton, in water, **IV**:1.14
- Z-pinch plasma, **V**:57.1–57.5, 57.2*f*, 57.3*t*, 57.4*f*
- Z-scan:
 for nonlinear optical parameters, **IV**:19.9, 19.9*f*, 19.10*f*
 and third-order optical nonlinearities, **IV**:16.29–16.30, 16.30*f*
- Z-system (eccentric pupil design), **II**:7.11, 7.12*f*, 7.15*f*, 7.15–7.17, 7.17*f*, 7.19, 7.21*f*
- Z-tilt (wavefront tilt), **V**:4.3, 4.23–4.25, 5.14, 5.15